Data Warehouse Service 9.1.0.210

Developer Guide

Issue 01

Date 2025-07-22





Copyright © Huawei Cloud Computing Technologies Co., Ltd. 2025. All rights reserved.

No part of this document may be reproduced or transmitted in any form or by any means without prior written consent of Huawei Cloud Computing Technologies Co., Ltd.

Trademarks and Permissions

HUAWEI and other Huawei trademarks are the property of Huawei Technologies Co., Ltd. All other trademarks and trade names mentioned in this document are the property of their respective holders.

Notice

The purchased products, services and features are stipulated by the contract made between Huawei Cloud and the customer. All or part of the products, services and features described in this document may not be within the purchase scope or the usage scope. Unless otherwise specified in the contract, all statements, information, and recommendations in this document are provided "AS IS" without warranties, quarantees or representations of any kind, either express or implied.

The information in this document is subject to change without notice. Every effort has been made in the preparation of this document to ensure accuracy of the contents, but all statements, information, and recommendations in this document do not constitute a warranty of any kind, express or implied.

Huawei Cloud Computing Technologies Co., Ltd.

Address: Huawei Cloud Data Center Jiaoxinggong Road

Qianzhong Avenue Gui'an New District Gui Zhou 550029

People's Republic of China

Website: https://www.huaweicloud.com/intl/en-us/

i

Contents

1 Before You Start	1
2 GaussDB(DWS) Development Design Proposal	5
2.1 Overview	
2.2 GaussDB(DWS) Connection Management Specifications	9
2.3 GaussDB(DWS) Object Design Specifications	11
2.3.1 DATABASE Object Design	11
2.3.2 USER Object Design	12
2.3.3 Schema Object Design	13
2.3.4 TABLESPACE Object Design	13
2.3.5 TABLE Object Design (Prioritized)	13
2.3.6 INDEX Object Design (Prioritized)	18
2.3.7 VIEW Object Design	19
2.4 GaussDB(DWS) SQL Statement Development Specifications	19
2.4.1 DDL Operations	19
2.4.2 INSERT Operation	20
2.4.3 UPDATE and DELETE Operations	21
2.4.4 SELECT Operation	21
2.5 GaussDB(DWS) Foreign Table Function Development Specifications	25
2.6 GaussDB(DWS) Stored Procedure Development Specifications	26
2.7 Detailed Design Rules for GaussDB(DWS) Objects	27
2.7.1 GaussDB(DWS) Database Object Naming Rules	27
2.7.2 GaussDB(DWS) Database Object Design Rules	28
2.7.2.1 GaussDB(DWS) Database and Schema Design Rules	28
2.7.2.2 GaussDB(DWS) Table Design Rules	29
2.7.2.3 GaussDB(DWS) Column Design Rules	32
2.7.2.4 GaussDB(DWS) Constraint Design Rules	33
2.7.2.5 Design Rules for GaussDB(DWS) Views and Associated Tables	34
2.7.3 GaussDB(DWS) JDBC Configuration Rules	34
2.7.4 GaussDB(DWS) SQL Writing Rules	36
2.7.5 Rules for Using Custom GaussDB(DWS) External Functions (pgSQL/Java)	39
2.7.6 Rules for Using GaussDB(DWS) PL/pgSQL	40
3 Creating and Managing GaussDB(DWS) Database Objects	44

3.1 Creating and Managing GaussDB(DWS) Databases	44
3.2 Creating and Managing GaussDB(DWS) Schemas	
3.3 Creating and Managing GaussDB(DWS) Tables	48
3.4 Selecting a GaussDB(DWS) Table Storage Model	54
3.5 Creating and Managing GaussDB(DWS) Partitioned Tables	58
3.6 Creating and Managing GaussDB(DWS) Indexes	61
3.7 Creating and Using GaussDB(DWS) Sequences	64
3.8 Creating and Managing GaussDB(DWS) Views	66
3.9 Creating and Managing GaussDB(DWS) Scheduled Tasks	67
3.10 Viewing GaussDB(DWS) System Catalogs	70
4 Syntax Compatibility Differences Among Oracle, Teradata, and MySQL	73
5 GaussDB(DWS) Database Security Management	
5.1 GaussDB(DWS) User and Permissions Management	
5.1.1 GaussDB(DWS) Database User Types	
5.1.2 GaussDB(DWS) Database User Management	
5.1.3 Creating a Custom Password Policy for GaussDB(DWS)	
5.1.4 GaussDB(DWS) Database Permissions Management	
5.1.5 Separation of Duties in GaussDB(DWS)	
5.2 GaussDB(DWS) Sensitive Data Management	
5.2.1 GaussDB(DWS) Row-Level Access Control	
5.2.2 GaussDB(DWS) Data Masking	
5.2.3 Encrypting and Decrypting GaussDB(DWS) Strings	
6 GaussDB(DWS) Data Query	
6.1 GaussDB(DWS) Single-Table Query	
6.2 GaussDB(DWS) Multi-Table Join Query	
6.3 GaussDB(DWS) Subquery Expressions	
6.4 GaussDB(DWS) WITH Expressions	
6.5 Usage of GaussDB(DWS) UNION	
6.6 Data Reading/Writing Across Logical Clusters	
6.7 SQL on Hudi	
6.7.1 Introduction to Hudi	
6.7.2 Preparations Before Using Hudi	
6.7.4 Creating a Hudi Data Description (Foreign Table)	
6.7.5 Synchronizing Hudi Tasks	
6.7.6 Querying a Hudi Foreign Table	
6.7.7 Accessing Hudi Tables on MRS	
7 GaussDB(DWS) Sorting Rules	
8 GaussDB(DWS) User-Defined Functions	

8.1 GaussDB(DWS) PL/Java Functions	168
8.2 GaussDB(DWS) PL/pgSQL Functions	179
9 GaussDB(DWS) Stored Procedure	181
9.1 Overview	
9.2 Converting Data Types in GaussDB(DWS) Stored Procedures	181
9.3 GaussDB(DWS) Stored Procedure Array and Record	183
9.3.1 Arrays	183
9.3.2 record	189
9.4 GaussDB(DWS) Stored Procedure Declaration Syntax	191
9.5 Basic Statements of GaussDB(DWS) Stored Procedures	193
9.6 Dynamic Statements of GaussDB(DWS) Stored Procedures	197
9.6.1 Executing Dynamic Query Statements	197
9.6.2 Executing Dynamic Non-query Statements	199
9.6.3 Dynamically Calling Stored Procedures	200
9.6.4 Dynamically Calling Anonymous Blocks	202
9.7 GaussDB(DWS) Stored Procedure Control Statements	203
9.7.1 RETURN Statements	204
9.7.2 Conditional Statements	206
9.7.3 Loop Statements	208
9.7.4 Branch Statements	211
9.7.5 NULL Statements	212
9.7.6 Error Trapping Statements	212
9.7.7 GOTO Statements	214
9.8 Other Statements in a GaussDB(DWS) Stored Procedure	216
9.9 GaussDB(DWS) Stored Procedure Cursor	217
9.9.1 Overview	217
9.9.2 Explicit Cursor	217
9.9.3 Implicit Cursor	221
9.9.4 Cursor Loop	222
9.10 GaussDB(DWS) Stored Procedure Advanced Package	224
9.10.1 DBMS_LOB	224
9.10.2 DBMS_RANDOM	233
9.10.3 DBMS_OUTPUT	234
9.10.4 UTL_RAW	235
9.10.5 DBMS_JOB	237
9.10.6 DBMS_SQL	247
9.11 GaussDB(DWS) Stored Procedure Debugging	258
10 Using PostGIS Extension	261
10.1 PostGIS	261
10.2 Using PostGIS	262
10.3 PostGIS Support and Constraints	263
10.4 OPEN SOURCE SOFTWARE NOTICE (For PostGIS)	275

Tri Using JDBC or ODBC for GaussDB(DWS) secondary Development	323
11.1 Prerequisites	323
11.2 JDBC-Based Development	323
11.2.1 JDBC Development Process	323
11.2.2 JDBC Package and Driver Class	325
11.2.3 Loading a Driver	325
11.2.4 Connecting to a Database	325
11.2.5 Executing SQL Statements	329
11.2.6 Processing Data in a Result Set	332
11.2.7 Common JDBC Development Examples	335
11.2.8 Processing RoaringBitmap Result Sets and Importing It to GaussDB (DWS)	345
11.2.9 JDBC Interfaces	347
11.3 ODBC-Based Development	361
11.3.1 ODBC Package and Its Dependent Libraries and Header Files	363
11.3.2 Configuring a Data Source in the Linux OS	363
11.3.3 Configuring a Data Source in the Windows OS	371
11.3.4 ODBC Development Example	376
11.3.5 ODBC APIs	381
12 GaussDB(DWS) Resource Monitoring	401
12.1 User Resource Monitoring	
12.2 Resource Pool Monitoring	
12.3 Monitoring Memory Resources	
12.4 Instance Resource Monitoring	
12.5 Real-time Top SQL	
12.6 Historical Top SQL	
12.7 Example for Querying for Top SQLs	
13 GaussDB(DWS) Performance Tuning	422
13.1 Overview	
13.2 Performance Diagnosis	
13.2.1 Cluster Performance Analysis	
13.2.2 Slow SQL Analysis	
13.2.2.1 Querying SQL Statements That Affect Performance Most	
13.2.2.2 Checking Blocked Statements	
13.2.3 SQL Diagnosis	
13.2.4 Table Diagnosis	
13.3 System Optimization	
13.3.1 Tuning Database Parameters	
13.3.2 SMP Parallel Execution	
13.3.3 Configuring LLVM	
13.4 SQL Tuning	
13.4.1 SQL Query Execution Process	
13.4.2 SQL Execution Plan	

13.4.3 Execution Plan Operator	457
13.4.4 SQL Tuning Process	462
13.4.5 Updating Statistics	463
13.4.6 Reviewing and Modifying a Table Definition	472
13.4.7 Advanced SQL Tuning	473
13.4.7.1 SQL Self-Diagnosis	473
13.4.7.2 Optimizing Statement Pushdown	477
13.4.7.3 Optimizing Subqueries	484
13.4.7.4 Optimizing Statistics	492
13.4.7.5 Tuning Operators	498
13.4.7.6 Optimizing Data Skew	499
13.4.7.7 Proactive Preheating and Tuning of Disk Cache	505
13.4.7.8 SQL Statement Rewriting Rules	506
13.4.8 Configuring Optimizer Parameters	508
13.4.9 Hint-based Tuning	509
13.4.9.1 Plan Hint Optimization	510
13.4.9.2 Join Order Hints	512
13.4.9.3 Join Operation Hints	516
13.4.9.4 Rows Hints	517
13.4.9.5 Stream Operation Hints	518
13.4.9.6 Scan Operation Hints	521
13.4.9.7 Sublink Name Hints	522
13.4.9.8 Skew Hints	523
13.4.9.9 Hint That Disables Subquery Pull-up	528
13.4.9.10 Drive Hints	529
13.4.9.11 Dictionary Code Hint	531
13.4.9.12 Configuration Parameter Hints	532
13.4.9.13 Hint Errors, Conflicts, and Other Warnings	535
13.4.9.14 Plan Hint Cases	537
13.4.10 Routinely Maintaining Tables	542
13.4.11 Routinely Recreating an Index	544
13.4.12 Automatic Retry upon SQL Statement Execution Errors	545
13.4.13 Query Band Load Identification	548
13.5 SQL Tuning Examples	553
13.5.1 Case: Selecting an Appropriate Distribution Column	553
13.5.2 Case: Creating an Appropriate Index	554
13.5.3 Case: Adding NOT NULL for JOIN Columns	555
13.5.4 Case: Pushing Down Sort Operations to DNs	557
13.5.5 Case: Configuring cost_param for Better Query Performance	558
13.5.6 Case: Adjusting the Partial Clustering Key	562
13.5.7 Case: Adjusting the Table Storage Mode in a Medium Table	564
13.5.8 Case: Reconstructing Partition Tables	565

13.5.9 Case: Adjusting the GUC Parameter best_agg_plan	566
13.5.10 Case: Rewriting SQL Statements and Eliminating Prune Interference	568
13.5.11 Case: Rewriting SQL Statements and Deleting in-clause	570
13.5.12 Case: Setting Partial Cluster Keys	571
13.5.13 Case: Converting from NOT IN to NOT EXISTS	574
14 GaussDB(DWS) System Catalogs and Views	576
14.1 Overview of System Catalogs and System Views	576
14.2 System Catalogs	577
14.2.1 GS_BLOCKLIST_QUERY	577
14.2.2 GS_BLOCKLIST_SQL	578
14.2.3 GS_OBSSCANINFO	579
14.2.4 GS_RESPOOL_RESOURCE_HISTORY	579
14.2.5 GS_WLM_INSTANCE_HISTORY	583
14.2.6 GS_WLM_OPERATOR_INFO	584
14.2.7 GS_WLM_SESSION_INFO	586
14.2.8 GS_WLM_USER_RESOURCE_HISTORY	594
14.2.9 PG_AGGREGATE	596
14.2.10 PG_AM	597
14.2.11 PG_AMOP	599
14.2.12 PG_AMPROC	600
14.2.13 PG_ATTRDEF	600
14.2.14 PG_ATTRIBUTE	601
14.2.15 PG_AUTHID	603
14.2.16 PG_AUTH_HISTORY	605
14.2.17 PG_AUTH_MEMBERS	605
14.2.18 PG_BLOCKLISTS	606
14.2.19 PG_CAST	607
14.2.20 PG_CLASS	608
14.2.21 PG_COLLATION	613
14.2.22 PG_CONSTRAINT	613
14.2.23 PG_CONVERSION	616
14.2.24 PG_DATABASE	616
14.2.25 PG_DB_ROLE_SETTING	617
14.2.26 PG_DEFAULT_ACL	618
14.2.27 PG_DEPEND	619
14.2.28 PG_DESCRIPTION	620
14.2.29 PG_ENUM	621
14.2.30 PG_EXCEPT_RULE	622
14.2.31 PG_EXTERNAL_NAMESPACE	622
14.2.32 PG_EXTENSION	623
14.2.33 PG_EXTENSION_DATA_SOURCE	623
14.2.34 PG FINE DR INFO	

14.2.35 PG_FOREIGN_DATA_WRAPPER	625
14.2.36 PG_FOREIGN_SERVER	625
14.2.37 PG_FOREIGN_TABLE	626
14.2.38 PG_INDEX	626
14.2.39 PG_INHERITS	629
14.2.40 PG_JOB_INFO	629
14.2.41 PG_JOBS	629
14.2.42 PG_LANGUAGE	631
14.2.43 PG_LARGEOBJECT	632
14.2.44 PG_LARGEOBJECT_METADATA	632
14.2.45 PG_MATVIEW	633
14.2.46 PG_NAMESPACE	634
14.2.47 PG_OBJECT	634
14.2.48 PG_OBSSCANINFO	635
14.2.49 PG_OPCLASS	636
14.2.50 PG_OPERATOR	637
14.2.51 PG_OPFAMILY	638
14.2.52 PG_PARTITION	638
14.2.53 PG_PLTEMPLATE	641
14.2.54 PG_PROC	642
14.2.55 PG_PUBLICATION	645
14.2.56 PG_PUBLICATION_NAMESPACE	646
14.2.57 PG_PUBLICATION_REL	647
14.2.58 PG_RANGE	647
14.2.59 PG_REDACTION_COLUMN	648
14.2.60 PG_REDACTION_POLICY	649
14.2.61 PG_RELFILENODE_SIZE	650
14.2.62 PG_RLSPOLICY	651
14.2.63 PG_RESOURCE_POOL	651
14.2.64 PG_REWRITE	653
14.2.65 PG_SECLABEL	653
14.2.66 PG_SHDEPEND	654
14.2.67 PG_SHDESCRIPTION	655
14.2.68 PG_SHSECLABEL	656
14.2.69 PG_STATISTIC	656
14.2.70 PG_STATISTIC_EXT	658
14.2.71 PG_STAT_OBJECT	659
14.2.72 PG_SUBSCRIPTION	663
14.2.73 PG_SYNONYM	664
14.2.74 PG_TABLESPACE	665
14.2.75 PG_TRIGGER	665
14.2.76 PG_TS_CONFIG	666

14.2.77 PG_TS_CONFIG_MAP	667
14.2.78 PG_TS_DICT	667
14.2.79 PG_TS_PARSER	668
14.2.80 PG_TS_TEMPLATE	669
14.2.81 PG_TYPE	669
14.2.82 PG_USER_MAPPING	673
14.2.83 PG_USER_STATUS	674
14.2.84 PG_WORKLOAD_ACTION	674
14.2.85 PGXC_CLASS	675
14.2.86 PGXC_GROUP	675
14.2.87 PGXC_NODE	677
14.2.88 PLAN_TABLE_DATA	679
14.2.89 SNAPSHOT	680
14.2.90 TABLES_SNAP_TIMESTAMP	680
14.2.91 System Catalogs for Performance View Snapshot	681
14.3 System Views	682
14.3.1 ALL_ALL_TABLES	682
14.3.2 ALL_CONSTRAINTS	682
14.3.3 ALL_CONS_COLUMNS	683
14.3.4 ALL_COL_COMMENTS	683
14.3.5 ALL_DEPENDENCIES	683
14.3.6 ALL_IND_COLUMNS	684
14.3.7 ALL_IND_EXPRESSIONS	685
14.3.8 ALL_INDEXES	685
14.3.9 ALL_OBJECTS	686
14.3.10 ALL_PROCEDURES	686
14.3.11 ALL_SEQUENCES	686
14.3.12 ALL_SOURCE	687
14.3.13 ALL_SYNONYMS	687
14.3.14 ALL_TAB_COLUMNS	688
14.3.15 ALL_TAB_COMMENTS	689
14.3.16 ALL_TABLES	689
14.3.17 ALL_USERS	690
14.3.18 ALL_VIEWS	690
14.3.19 DBA_DATA_FILES	690
14.3.20 DBA_USERS	691
14.3.21 DBA_COL_COMMENTS	691
14.3.22 DBA_CONSTRAINTS	691
14.3.23 DBA_CONS_COLUMNS	692
14.3.24 DBA_IND_COLUMNS	692
14.3.25 DBA_IND_EXPRESSIONS	693
14.3.26 DBA IND PARTITIONS	693

14.3.27 DBA_INDEXES	694
14.3.28 DBA_OBJECTS	695
14.3.29 DBA_PART_INDEXES	695
14.3.30 DBA_PART_TABLES	696
14.3.31 DBA_PROCEDURES	697
14.3.32 DBA_SEQUENCES	697
14.3.33 DBA_SOURCE	697
14.3.34 DBA_SYNONYMS	698
14.3.35 DBA_TAB_COLUMNS	698
14.3.36 DBA_TAB_COMMENTS	699
14.3.37 DBA_TAB_PARTITIONS	699
14.3.38 DBA_TABLES	701
14.3.39 DBA_TABLESPACES	701
14.3.40 DBA_TRIGGERS	701
14.3.41 DBA_VIEWS	702
14.3.42 DUAL	702
14.3.43 GET_ALL_TSC_INFO	702
14.3.44 GET_TSC_INFO	703
14.3.45 GLOBAL_COLUMN_TABLE_IO_STAT	703
14.3.46 GLOBAL_REDO_STAT	704
14.3.47 GLOBAL_REL_IOSTAT	705
14.3.48 GLOBAL_ROW_TABLE_IO_STAT	705
14.3.49 GLOBAL_STAT_DATABASE	706
14.3.50 GLOBAL_TABLE_CHANGE_STAT	708
14.3.51 GLOBAL_TABLE_STAT	709
14.3.52 GLOBAL_WORKLOAD_SQL_COUNT	710
14.3.53 GLOBAL_WORKLOAD_SQL_ELAPSE_TIME	711
14.3.54 GLOBAL_WORKLOAD_TRANSACTION	712
14.3.55 GS_ALL_CONTROL_GROUP_INFO	713
14.3.56 GS_BLOCKLIST_QUERY	713
14.3.57 GS_BLOCKLIST_SQL	714
14.3.58 GS_CLUSTER_RESOURCE_INFO	715
14.3.59 GS_COLUMN_TABLE_IO_STAT	715
14.3.60 GS_OBS_READ_TRAFFIC	716
14.3.61 GS_OBS_WRITE_TRAFFIC	716
14.3.62 GS_INSTR_UNIQUE_SQL	717
14.3.63 GS_NODE_STAT_RESET_TIME	722
14.3.64 GS_OBS_LATENCY	722
14.3.65 GS_QUERY_MONITOR	723
14.3.66 GS_QUERY_RESOURCE_INFO	725
14.3.67 GS_REL_IOSTAT	726
14.3.68 GS RESPOOL RUNTIME INFO	726

14.3.69 GS_RESPOOL_RESOURCE_INFO	727
14.3.70 GS_RESPOOL_MONITOR	730
14.3.71 GS_ROW_TABLE_IO_STAT	732
14.3.72 GS_SESSION_CPU_STATISTICS	733
14.3.73 GS_SESSION_MEMORY_STATISTICS	734
14.3.74 GS_SQL_COUNT	
14.3.75 GS_STAT_DB_CU	736
14.3.76 GS_STAT_SESSION_CU	737
14.3.77 GS_TABLE_CHANGE_STAT	737
14.3.78 GS_TABLE_STAT	738
14.3.79 GS_TOTAL_NODEGROUP_MEMORY_DETAIL	739
14.3.80 GS_USER_MONITOR	740
14.3.81 GS_USER_TRANSACTION	
14.3.82 GS_VIEW_DEPENDENCY	
14.3.83 GS_VIEW_DEPENDENCY_PATH	
14.3.84 GS_VIEW_INVALID	
14.3.85 GS_WAIT_EVENTS	
14.3.86 GS_WLM_OPERATOR_INFO	745
14.3.87 GS_WLM_OPERATOR_HISTORY	
14.3.88 GS_WLM_OPERATOR_STATISTICS	
14.3.89 GS_WLM_SESSION_INFO	
14.3.90 GS_WLM_SESSION_HISTORY	
14.3.91 GS_WLM_SESSION_STATISTICS	
14.3.92 GS_WLM_SQL_ALLOW	
14.3.93 GS_WORKLOAD_SQL_COUNT	
14.3.94 GS_WORKLOAD_SQL_ELAPSE_TIME	
14.3.95 GS_WORKLOAD_TRANSACTION	
14.3.96 MPP_TABLES	
14.3.97 PG_AVAILABLE_EXTENSION_VERSIONS	
14.3.98 PG_AVAILABLE_EXTENSIONS	
14.3.99 PG_BULKLOAD_STATISTICS	
14.3.100 PG_COMM_CLIENT_INFO	
14.3.101 PG_COMM_DELAY	
14.3.102 PG_COMM_STATUS	
14.3.103 PG_COMM_RECV_STREAM	
14.3.104 PG_COMM_SEND_STREAM	
14.3.105 PG_COMM_QUERY_SPEED	
14.3.106 PG_CONTROL_GROUP_CONFIG	
14.3.107 PG_CURSORS	
14.3.108 PG_EXT_STATS	
14.3.109 PG_GET_INVALID_BACKENDS	
14.3.110 PG GET SENDERS CATCHUP TIME	

14.3.111 PG_GLOBAL_TEMP_ATTACHED_PIDS	787
14.3.112 PG_GROUP	787
14.3.113 PG_INDEXES	787
14.3.114 PG_JOB	788
14.3.115 PG_JOB_PROC	790
14.3.116 PG_JOB_SINGLE	790
14.3.117 PG_LIFECYCLE_DATA_DISTRIBUTE	792
14.3.118 PG_LOCKS	792
14.3.119 PG_LWLOCKS	794
14.3.120 PG_NODE_ENV	795
14.3.121 PG_OS_THREADS	796
14.3.122 PG_POOLER_STATUS	796
14.3.123 PG_PREPARED_STATEMENTS	797
14.3.124 PG_PREPARED_XACTS	798
14.3.125 PG_PUBLICATION_TABLES	798
14.3.126 PG_QUERYBAND_ACTION	799
14.3.127 PG_REPLICATION_SLOTS	799
14.3.128 PG_ROLES	800
14.3.129 PG_RULES	801
14.3.130 PG_RUNNING_XACTS	802
14.3.131 PG_SECLABELS	802
14.3.132 PG_SEQUENCES	803
14.3.133 PG_SESSION_WLMSTAT	804
14.3.134 PG_SESSION_IOSTAT	806
14.3.135 PG_SETTINGS	807
14.3.136 PG_SHADOW	808
14.3.137 PG_SHARED_MEMORY_DETAIL	809
14.3.138 PG_STATS	809
14.3.139 PG_STAT_ACTIVITY	811
14.3.140 PG_STAT_ALL_INDEXES	814
14.3.141 PG_STAT_ALL_TABLES	815
14.3.142 PG_STAT_BAD_BLOCK	817
14.3.143 PG_STAT_BGWRITER	817
14.3.144 PG_STAT_DATABASE	818
14.3.145 PG_STAT_DATABASE_CONFLICTS	819
14.3.146 PG_STAT_GET_MEM_MBYTES_RESERVED	820
14.3.147 PG_STAT_USER_FUNCTIONS	821
14.3.148 PG_STAT_USER_INDEXES	821
14.3.149 PG_STAT_USER_TABLES	822
14.3.150 PG_STAT_REPLICATION	823
14.3.151 PG_STAT_SYS_INDEXES	824
14.3.152 PG_STAT_SYS_TABLES	824

14.3.153 PG_STAT_XACT_ALL_TABLES	825
14.3.154 PG_STAT_XACT_SYS_TABLES	826
14.3.155 PG_STAT_XACT_USER_FUNCTIONS	827
14.3.156 PG_STAT_XACT_USER_TABLES	827
14.3.157 PG_STATIO_ALL_INDEXES	828
14.3.158 PG_STATIO_ALL_SEQUENCES	828
14.3.159 PG_STATIO_ALL_TABLES	829
14.3.160 PG_STATIO_SYS_INDEXES	829
14.3.161 PG_STATIO_SYS_SEQUENCES	830
14.3.162 PG_STATIO_SYS_TABLES	830
14.3.163 PG_STATIO_USER_INDEXES	831
14.3.164 PG_STATIO_USER_SEQUENCES	831
14.3.165 PG_STATIO_USER_TABLES	832
14.3.166 PG_THREAD_WAIT_STATUS	833
14.3.167 PG_TABLES	845
14.3.168 PG_TDE_INFO	846
14.3.169 PG_TIMEZONE_ABBREVS	847
14.3.170 PG_TIMEZONE_NAMES	847
14.3.171 PG_TOTAL_MEMORY_DETAIL	847
14.3.172 PG_TOTAL_SCHEMA_INFO	
14.3.173 PG_TOTAL_USER_RESOURCE_INFO	850
14.3.174 PG_USER	852
14.3.175 PG_USER_MAPPINGS	853
14.3.176 PG_VIEWS	854
14.3.177 PG_WLM_STATISTICS	854
14.3.178 PGXC_AIO_RESOURCE_POOL_STATS	855
14.3.179 PGXC_BULKLOAD_PROGRESS	857
14.3.180 PGXC_BULKLOAD_INFO	858
14.3.181 PGXC_BULKLOAD_STATISTICS	861
14.3.182 PGXC_COLUMN_TABLE_IO_STAT	862
14.3.183 PGXC_COMM_CLIENT_INFO	863
14.3.184 PGXC_COMM_DELAY	863
14.3.185 PGXC_COMM_RECV_STREAM	864
14.3.186 PGXC_COMM_SEND_STREAM	865
14.3.187 PGXC_COMM_STATUS	867
14.3.188 PGXC_COMM_QUERY_SPEED	867
14.3.189 PGXC_DEADLOCK	868
14.3.190 PGXC_DISK_CACHE_STATS	870
14.3.191 PGXC_DISK_CACHE_ALL_STATS	870
14.3.192 PGXC_DISK_CACHE_PATH_INFO	872
14.3.193 PGXC_GET_STAT_ALL_TABLES	872
14.3.194 PGXC_GET_STAT_ALL_PARTITIONS	873

14.3.195 PGXC_GET_TABLE_SKEWNESS	875
14.3.196 PGXC_GLOBAL_TEMP_ATTACHED_PIDS	875
14.3.197 PGXC_GTM_SNAPSHOT_STATUS	876
14.3.198 PGXC_INSTANCE_TIME	876
14.3.199 PGXC_LOCKWAIT_DETAIL	877
14.3.200 PGXC_INSTR_UNIQUE_SQL	879
14.3.201 PGXC_LOCK_CONFLICTS	882
14.3.202 PGXC_LWLOCKS	883
14.3.203 PGXC_MEMORY_DEBUG_INFO	884
14.3.204 PGXC_NODE_ENV	886
14.3.205 PGXC_NODE_STAT_RESET_TIME	886
14.3.206 PGXC_OBS_IO_SCHEDULER_STATS	887
14.3.207 PGXC_OBS_IO_SCHEDULER_PERIODIC_STATS	888
14.3.208 PGXC_OS_RUN_INFO	890
14.3.209 PGXC_OS_THREADS	891
14.3.210 PGXC_POOLER_STATUS	891
14.3.211 PGXC_PREPARED_XACTS	892
14.3.212 PGXC_REDO_STAT	892
14.3.213 PGXC_REL_IOSTAT	893
14.3.214 PGXC_REPLICATION_SLOTS	893
14.3.215 PGXC_RESPOOL_RUNTIME_INFO	894
14.3.216 PGXC_RESPOOL_RESOURCE_INFO	895
14.3.217 PGXC_RESPOOL_RESOURCE_HISTORY	898
14.3.218 PGXC_ROW_TABLE_IO_STAT	901
14.3.219 PGXC_RUNNING_XACTS	902
14.3.220 PGXC_SETTINGS	902
14.3.221 PGXC_SESSION_WLMSTAT	904
14.3.222 PGXC_STAT_ACTIVITY	906
14.3.223 PGXC_STAT_BAD_BLOCK	910
14.3.224 PGXC_STAT_BGWRITER	910
14.3.225 PGXC_STAT_DATABASE	911
14.3.226 PGXC_STAT_OBJECT	913
14.3.227 PGXC_STAT_REPLICATION	917
14.3.228 PGXC_STAT_TABLE_DIRTY	
14.3.229 PGXC_STAT_WAL	
14.3.230 PGXC_SQL_COUNT	923
14.3.231 PGXC_TABLE_CHANGE_STAT	923
14.3.232 PGXC_TABLE_STAT	924
14.3.233 PGXC_THREAD_WAIT_STATUS	925
14.3.234 PGXC_TOTAL_MEMORY_DETAIL	
14.3.235 PGXC_TOTAL_SCHEMA_INFO	929
14.3.236 PGXC_TOTAL_SCHEMA_INFO_ANALYZE	929

14.3.237 PGXC_TOTAL_USER_RESOURCE_INFO	930
14.3.238 PGXC_USER_TRANSACTION	933
14.3.239 PGXC_VARIABLE_INFO	934
14.3.240 PGXC_WAIT_DETAIL	935
14.3.241 PGXC_WAIT_EVENTS	937
14.3.242 PGXC_WLM_OPERATOR_HISTORY	938
14.3.243 PGXC_WLM_OPERATOR_INFO	939
14.3.244 PGXC_WLM_OPERATOR_STATISTICS	941
14.3.245 PGXC_WLM_SESSION_INFO	944
14.3.246 PGXC_WLM_SESSION_HISTORY	950
14.3.247 PGXC_WLM_SESSION_STATISTICS	958
14.3.248 PGXC_WLM_TABLE_DISTRIBUTION_SKEWNESS	963
14.3.249 PGXC_WLM_USER_RESOURCE_HISTORY	965
14.3.250 PGXC_WLM_WORKLOAD_RECORDS	968
14.3.251 PGXC_WORKLOAD_SQL_COUNT	969
14.3.252 PGXC_WORKLOAD_SQL_ELAPSE_TIME	970
14.3.253 PGXC_WORKLOAD_TRANSACTION	971
14.3.254 PLAN_TABLE	972
14.3.255 PV_FILE_STAT	973
14.3.256 PV_INSTANCE_TIME	974
14.3.257 PV_MATVIEW_DETAIL	974
14.3.258 PV_OS_RUN_INFO	975
14.3.259 PV_SESSION_MEMORY	976
14.3.260 PV_SESSION_MEMORY_DETAIL	976
14.3.261 PV_SESSION_STAT	978
14.3.262 PV_SESSION_TIME	978
14.3.263 PV_TOTAL_MEMORY_DETAIL	978
14.3.264 PV_REDO_STAT	
14.3.265 PV_RUNTIME_ATTSTATS	980
14.3.266 PV_RUNTIME_RELSTATS	982
14.3.267 REDACTION_COLUMNS	983
14.3.268 REDACTION_POLICIES	984
14.3.269 REMOTE_TABLE_STAT	985
14.3.270 SHOW_TSC_INFO	986
14.3.271 SHOW_ALL_TSC_INFO	987
14.3.272 USER_COL_COMMENTS	987
14.3.273 USER_CONSTRAINTS	987
14.3.274 USER_CONS_COLUMNS	988
14.3.275 USER_INDEXES	989
14.3.276 USER_IND_COLUMNS	989
14.3.277 USER_IND_EXPRESSIONS	989
14.3.278 USER IND PARTITIONS	990

14.3.279 USER_JOBS	991	
14.3.280 USER_OBJECTS	992	
14.3.281 USER_PART_INDEXES	993	
14.3.282 USER_PART_TABLES	993	
14.3.283 USER_PROCEDURES	994	
14.3.284 USER_SEQUENCES	994	
14.3.285 USER_SOURCE	995	
14.3.286 USER_SYNONYMS	995	
14.3.287 USER_TAB_COLUMNS	995	
14.3.288 USER_TAB_COMMENTS	996	
14.3.289 USER_TAB_PARTITIONS	997	
14.3.290 USER_TABLES	997	
14.3.291 USER_TRIGGERS	998	
14.3.292 USER_VIEWS	998	
14.3.293 V\$SESSION	999	
14.3.294 V\$SESSION_LONGOPS	999	
15 GUC Parameters of the GaussDB(DWS) Database	1000	
15.1 Viewing GUC Parameters	1000	
15.2 Configuring GUC Parameters	1001	
15.3 GUC Parameter Usage	1003	
15.4 Connection and Authentication	1003	
15.4.1 Connection Settings	1003	
15.4.2 Security and Authentication (postgresql.conf)	1005	
15.4.3 Communication Library Parameters	1006	
15.5 Resource Consumption	1013	
15.5.1 Memory	1013	
15.5.2 Statement Disk Space Control	1023	
15.5.3 Kernel Resources	1024	
15.5.4 Cost-based Vacuum Delay	1024	
15.5.5 Asynchronous I/O Operations	1026	
15.5.6 Disk Caching	1029	
15.6 Parallel Data Import	1029	
15.7 Write Ahead Logs	1031	
15.7.1 Settings	1031	
15.7.2 Checkpoints	1034	
15.8 HA Replication	1035	
15.8.1 Primary Server		
15.9 Query Planning		
15.9.1 Optimizer Method Configuration		
15.9.2 Optimizer Cost Constants		
15.9.3 Genetic Query Optimizer		
5 9 4 Other Ontimizer Ontions		

15.10 Error Reporting and Logging	1080
15.10.1 Logging Destination	
15.10.2 Logging Time	
15.10.3 Logging Content	
15.11 Runtime Statistics	
15.11.1 Query and Index Statistics Collector	1089
15.11.2 Performance Statistics	1094
15.12 Resource Management	1094
15.13 Automatic Cleanup	1107
15.14 Default Settings of Client Connection	1110
15.14.1 Statement Behavior	1110
15.14.2 Zone and Formatting	1117
15.14.3 Other Default Parameters	1121
15.15 Lock Management	1121
15.16 Version and Platform Compatibility	1125
15.16.1 Compatibility with Earlier Versions	1125
15.16.2 Platform and Client Compatibility	1129
15.17 Fault Tolerance	1172
15.18 Connection Pool Parameters	1173
15.19 Cluster Transaction Parameters	1175
15.20 Developer Operations	1178
15.21 Auditing	1199
15.21.1 Audit Switch	1199
15.21.2 Operation Audit	1201
15.22 Transaction Monitoring	1201
15.23 GTM Parameters	1202
15.24 Miscellaneous Parameters	1203
16 GaussDB(DWS) Developer Terms	1213
17 Hybrid Data Warehouse	1231
17.1 Introduction to Hybrid Data Warehouse	1231
17.2 Support and Constraints	1236
17.3 Hybrid Data Warehouse Syntax	1237
17.3.1 CREATE TABLE	1238
17.3.2 INSERT	1246
17.3.3 DELETE	1248
17.3.4 UPDATE	1249
17.3.5 UPSERT	1251
17.3.6 MERGE INTO	1253
17.3.7 SELECT	1254
17.3.8 ALTER TABLE	1256
17.4 Hybrid Data Warehouse Functions	1258
17.5 Hybrid Data Warehouse Binlog	

Data Warehouse Service	2
Developer Guide	

Contents

17.5.1 Subscribing to Hybrid Data Warehouse Binlog	1265
17.5.2 Real-Time Binlog Consumption by Flink	1270

Before You Start

Target Readers

This document is intended for database designers, application developers, and database administrators, and provides information required for designing, building, querying and maintaining data warehouses.

As a database administrator or application developer, you need to be familiar with:

- Knowledge about OSs, which is the basis for everything.
- SQL syntax, which is the necessary skill for database operation.

Prerequisites

Complete the following tasks before you perform operations described in this document:

- Create a GaussDB(DWS) cluster.
- Install a SQL client.
- Connect the SQL client to the default database of the cluster.

For details about these tasks, see **Getting Started with GaussDB(DWS)**.

Reading Guide

If you are a new GaussDB(DWS) user, you are advised to read the following contents first:

- Sections describing the features, functions, and application scenarios of GaussDB(DWS).
- "Getting Started": guides you through creating a data warehouse cluster, creating a database table, uploading data, and testing queries.

If you intend to or are migrating applications from other data warehouses to GaussDB(DWS), you might want to know how GaussDB(DWS) differs from them.

You can find useful information from the following table for GaussDB(DWS) database application development.

Operation	Query Suggestion
Quickly getting started with	Deploy a cluster, connect to the database, and perform some queries by referring to Getting Started .
GaussDB(DWS)	When you are ready to construct a database, load data to tables and compile the query content to operate the data in the data warehouse. Then, you can return to the <i>Data Warehouse Service Database Developer Guide</i> .
Understand the internal architecture of a GaussDB(DWS) data warehouse.	To know more about GaussDB(DWS), go to the GaussDB(DWS) homepage.
Learn how to design tables to achieve the excellent performance.	GaussDB(DWS) Development Design Proposal introduces the design specifications that should be complied with during the development of database applications. Modeling compliant with these specifications fits the distributed processing architecture of GaussDB(DWS) and provides efficient SQL code.
	To accelerate service execution through optimization, refer to GaussDB(DWS) Performance Tuning. Database administrators' experience and judgment play a more significant role in achieving successful performance optimization than instructions and explanations. However, GaussDB(DWS) Performance Tuning still attempts to illustrate the performance optimization methods that can be referred to by application development personnel and new GaussDB(DWS) database administrators.
Loading data	Importing Data describes how to import data to GaussDB(DWS).
	Excellent Practices for Data Import provides key points for quick data import.
Managing users, groups, and database security	GaussDB(DWS) Database Security Management covers database security topics.
Monitoring and optimizing system	GaussDB(DWS) System Catalogs and Views describes the system catalogs where you can query the database status and monitor the query content and process.
performance	You should also refer to Management Guide to learn how to use the GaussDB(DWS) console to check the system running status and monitoring metrics.

SQL Syntax Text Conventions

To better understand how to use the syntax, you can refer to the following description of SQL syntax text conventions.

Format	Description
Uppercase characters	Keywords must be in uppercase.
Lowercase characters	Parameters must be in lowercase.
[]	Items in brackets [] are optional.
	Preceding elements can appear repeatedly.
[x y]	One item is selected from two or more options or no item is selected.
{ x y }	One item is selected from two or more options.
[x y] []	You can choose either multiple parameters or no parameters. If you choose multiple parameters, simply separate them with spaces.
[x y][,]	You can choose either multiple parameters or no parameters. If you choose multiple parameters, simply separate them with commas (,).
{ x y } []	You must select at least one parameter. If you select multiple parameters, separate them with spaces.
{ x y } [,]	You must select at least one parameter. If you select multiple parameters, separate them with commas (,).

Statement

When writing documents, the writers of GaussDB(DWS) try their best to provide guidance from the perspective of commercial use, application scenarios, and task completion. Even so, references to PostgreSQL content may still exist in the document. For this type of content, the following PostgreSQL Copyright is applicable:

Postgres-XC is Copyright © 1996-2013 by the PostgreSQL Global Development Group.

PostgreSQL is Copyright @ 1996-2013 by the PostgreSQL Global Development Group.

Postgres95 is Copyright © 1994-5 by the Regents of the University of California.

IN NO EVENT SHALL THE UNIVERSITY OF CALIFORNIA BE LIABLE TO ANY PARTY FOR DIRECT, INDIRECT, SPECIAL, INCIDENTAL, OR CONSEQUENTIAL DAMAGES, INCLUDING LOST PROFITS, ARISING OUT OF THE USE OF THIS SOFTWARE AND ITS DOCUMENTATION, EVEN IF THE UNIVERSITY OF CALIFORNIA HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

THE UNIVERSITY OF CALIFORNIA SPECIFICALLY DISCLAIMS ANY WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE. THE SOFTWARE

PROVIDED HEREUNDER IS ON AN "AS-IS" BASIS, AND THE UNIVERSITY OF CALIFORNIA HAS NO OBLIGATIONS TO PROVIDE MAINTENANCE, SUPPORT, UPDATES, ENHANCEMENTS, OR MODIFICATIONS.

2 GaussDB(DWS) Development Design Proposal

2.1 Overview

Objective

This document outlines the rules for design and development that need to be followed when developing the GaussDB(DWS) database. The objective is to enhance development efficiency and ensure the continuity and stability of the service.

Application Scope

These specifications apply to all GaussDB(DWS) self-development scenarios, including designing and developing applications and database services.

Terms

Rule: a mandatory requirement that must be followed during database design and development.

Suggestion: an option that you need to consider for the design and development process.

Description: a detailed explanation of a rule or suggestion.

Overall Development and Design Specifications

The table below provides a list of development and design specifications that must be followed during GaussDB(DWS) development. You can click the links to access the corresponding rules for more details.

Table 2-1 GaussDB(DWS) development and design specifications

N o.	Category		Rule/Suggestion
1	Conn ectio	-	Rule 1.1: Configuring Load Balancing for GaussDB(DWS) Clusters
2	n man agem ent		Rule 1.2: Ending the Database Connection After Necessary Operations (Except in Connection Pool Scenarios)
3	regul ation s		Rule 1.3: Ensuring a Started Transaction Is Committed or Rolled Back
4			Rule 1.4: Ensuring the Idle Timeout Duration Is Shorter Than SESSION_TIMEOUT Value When Connection Pool Is Used for Applications
5			Rule 1.5: Restoring Parameters to Default Values in Connections Before Returning Them to the Pool
6			Rule 1.6: Manually Clearing Temporary Tables Created with a Connection Before Returning it to the Pool
7	Obje ct	DATABAS E object	Rule 2.1: Avoiding Direct Usage of Built-in Databases Such As postgres and gaussdb
8	desig n specif	design	Rule 2.2: Selecting the Suitable Database Code During Database Creation
9	icatio ns		Rule 2.3: Choosing the Right Database Type for Compatibility with the Database to Be Created
10		USER object design	Suggestion 2.4: Creating the Objects with Associated Calculations in the Same Database
11			Rule 2.5: Following the Least Privilege Principle and Avoiding Running Services Using Users with Special Permissions
12			Rule 2.6: Avoiding the Use of a Single Database Account for All Services
13		SCHEMA object design	Suggestion 2.7: Avoiding the Creation of Objects Under Other Users' Private Schemas
14		TABLESPA CE object design	Rule 2.8 Avoiding Tablespace Customization

N o.	11.55		Rule/Suggestion
15		TABLE object	Rule 2.9: Selecting the Optimal Distribution Method and Columns During Table Creation
16		design (prioritize d)	Rule 2.10 Selecting an Optimal Storage Type During Table Creation
17			Rule 2.11 Selecting an Optimal Partitioning Policy During Table Creation
18			Suggestion 2.12: Designing Table Columns for Fast and Accurate Queries
19			Suggestion 2.13: Avoiding the Usage of Auto- increment Columns or Data Types
20		INDEX object	Rule 2.14: Creating Necessary Indexes and Selecting Optimal Columns and Sequences for Them
21		design (prioritize d)	Suggestion 2.15: Optimizing Performance by Choosing the Right Index Type and Avoiding Indexes for Column-Store Tables
22		VIEW object design	Suggestion 2.16: Limiting View Nesting to Three Layers
23	SQL devel opme	DDL operation specificati	Suggestion 3.1: Avoiding Performing DDL Operations (Except CREATE) During Peak Hours or in Long Transactions
24	nt specif icatio	ons	Rule 3.2: Specifying the Scope of Objects to Be Deleted When Using DROP
25	ns	INSERT operation	Rule 3.3: Replacing INSERT with COPY for Efficient Multi-Value Batch Insertion
26		specificati ons	Suggestion 3.4: Avoiding Performing Real-time INSERT Operations on Common Column-store Tables
27		UPDATE/ DELETE	Suggestion 3.5: Preventing Simultaneous Updates or Deletions of the Same Row in a Row-store Table
28		operation specificati ons	Suggestion 3.6: Avoiding Frequent or Simultaneous UPDATE and DELETE Operations on Column-store Tables

N o.	3 - 3		Rule/Suggestion
29		Specificati ons for	Rule 3.7: Avoiding Executing SQL Statements That Do Not Support Pushdown
30		the SELECT operation	Rule 3.8: Specifying Association Conditions when Multiple Tables Are Associated
31		(including the query part in all	Rule 3.9: Ensuring Consistency of Data Types in Associated Fields across Multiple Tables
32		syntaxes)	Suggestion 3.10: Avoiding Function Calculation on Association and Filter Condition Fields
33			Suggestion 3.11: Performing Pressure Tests and Concurrency Control for Resource-intensive SQL Statements
34			Rule 3.12: Avoiding Excessive COUNT Operations on Large Row-store Tables
35			Suggestion 3.13: Avoid Getting Large Result Sets (Except for Data Exports)
36			Suggestion 3.14: Avoiding the Usage of SELECT * for Queries
37			Suggestion 3.15: Using WITH RECURSIVE with Defined Termination Condition for Recursion
38			Suggestion 3.16: Setting Schema Prefix for Table and Function Access
39			Suggestion 3.17: Identifying an SQL Statement with a Unique SQL Comment
40			Recommendation 3.18: Restricting SQL Statements to 64 KB in Length
41	Forei gn table	GDS foreign table	Rule 4.1 Deploying GDS on an Independent Server Outside the GaussDB(DWS) Cluster
42	funct ion devel opme nt specif icatio ns	Foreign table for collaborat ive analysis	Rule 4.2 Avoiding Concurrent Access to Multiple Collaborative Analysis Foreign Tables Across Clusters

N o.	Catego	ory	Rule/Suggestion
43	Store d	-	Suggestion 5.1: Simplifying Stored Procedures and Avoiding Nesting
44	dure devel opme nt specificatio ns		Rule 5.2: Avoiding Non-CREATE DDL Operations in Stored Procedures

2.2 GaussDB(DWS) Connection Management Specifications

Rule 1.1: Configuring Load Balancing for GaussDB(DWS) Clusters

■ NOTE

Impact of rule violation:

- Load imbalance causes performance problems and even service interruption.
- When a CN is faulty, services cannot be automatically recovered or the recovery may take a long time.

Solution:

- Configure ELB load balancing and connect the application to the load balancing IP address
- For how to use JDBC for load balancing, see Configuring JDBC to Connect to a Cluster (Load Balancing Mode).

Rule 1.2: Ending the Database Connection After Necessary Operations (Except in Connection Pool Scenarios)

■ NOTE

Impact of rule violation:

- The number of idle connections exceeds the maximum limit, causing connection creation failure.
- Resource overload occurs because there are too many idle connections.

Solution:

- After the connection between the application and the database is established and used, manually end the connection.
- Set the **session_timeout** parameter on the service side to set the idle timeout duration. The connection will be automatically ended when the idle timeout duration expires.

Rule 1.3: Ensuring a Started Transaction Is Committed or Rolled Back

Ⅲ NOTE

Impact of rule violation:

- If a transaction remains uncommitted for an extended period, it blocks operations such as ALTER, thereby affecting all services.
- The number of idle connections exceeds the maximum limit, causing connection creation failure.

Solution:

- **autocommit** is enabled by default, so there is no need to manually commit any transaction unless you modify the default setting.
- If a transaction is explicitly started, it must be explicitly ended (either by committing or rolling back) once the relevant operations are finished.

Rule 1.4: Ensuring the Idle Timeout Duration Is Shorter Than SESSION_TIMEOUT Value When Connection Pool Is Used for Applications

Impact of rule violation:

• The idle timeout mechanism on the service side clears connections in the connection pool, which negatively impacts connection reuse.

Solution:

To ensure everything works correctly, make sure the idle timeout duration of the
connection pool is shorter than the SESSION_TIMEOUT value in GaussDB(DWS). It is
advised to adjust the idle timeout duration rather than modifying the
SESSION_TIMEOUT value.

Rule 1.5: Restoring Parameters to Default Values in Connections Before Returning Them to the Pool

■ NOTE

Impact of rule violation:

 When a connection is reused by another service, the parameters set by the service may also be reused. This can result in performance issues or service errors.

Solution:

 Before returning the connection to the connection pool, use SET SESSION AUTHORIZATION DEFAULT; RESET ALL; to reset parameters.

Notes:

When connection pool is used for applications, if you set the global GUC parameter using **GS_GUC RELOAD** in GaussDB(DWS), restart the connection pool for the changes to be applied. This is because the modification only affects new connections in the connection pool.

Rule 1.6: Manually Clearing Temporary Tables Created with a Connection Before Returning it to the Pool

☐ NOTE

Impact of rule violation:

 When a connection is reused by other services, an error may be reported when a temporary table is created.

Solution:

• Before returning a connection to the connection pool, use **DROP** to clear the temporary table created by the current session.

2.3 GaussDB(DWS) Object Design Specifications

2.3.1 DATABASE Object Design

Rule 2.1: Avoiding Direct Usage of Built-in Databases Such As postgres and gaussdb

■ NOTE

Impact of rule violation:

- If the code or the compatibility setting of the built-in databases does not meet service requirements, you may need to migrate your data again.
- The time for changes to be applied may be prolonged if all services use built-in databases.

Solution:

• To meet the specific requirements of each service, it is recommended to create a dedicated database and allocate it accordingly.

Rule 2.2: Selecting the Suitable Database Code During Database Creation

□ NOTE

Impact of rule violation:

• Selecting the wrong database code may result in displaying garbled characters, and it is not possible to directly change the database code. In such cases, you will need to create a database and import the data again.

Solution:

• It is advisable to set the **ENCODING** to the UTF-8 format during database creation, unless there are specific requirements for a different encoding.

Rule 2.3: Choosing the Right Database Type for Compatibility with the Database to Be Created

∩ NOTE

Impact of rule violation:

Selecting the wrong type can lead to behavior inconsistencies after migrating the
database from a different vendor to GaussDB(DWS). Unfortunately, it is not possible to
directly change the compatibility setting. The only solution is to create a database and
import the data again.

Solution:

GaussDB(DWS) supports compatibility with databases like Teradata, Oracle, and MySQL.
 You can specify DBCOMPATIBILITY to set the compatible database type when creating a database.

Suggestion 2.4: Creating the Objects with Associated Calculations in the Same Database

□ NOTE

Impact of rule violation:

 Cross-database access tends to have poorer performance compared to performing operations within the same database.

Solution:

 If multiple databases are created, it is advisable to create the objects requiring associated calculations in the same database.

2.3.2 USER Object Design

Rule 2.5: Following the Least Privilege Principle and Avoiding Running Services Using Users with Special Permissions

□ NOTE

Impact of rule violation:

 Administrators have full access to a lot of things in the system and using these users to run services can pose security and control risks.

Solution:

 It is advised to use common users for service running, reserving users with special permissions for management operations.

Rule 2.6: Avoiding the Use of a Single Database Account for All Services

◯ NOTE

Impact of rule violation:

 Using a single database user for all services hinders effective service management and control. In abnormal situations, it becomes impossible to isolate specific users for emergency purposes.

Solution:

- Create administrators, service operation users, and O&M users for different purposes.
- Use different users to run different services for improved management and allocation of services and resources.

2.3.3 Schema Object Design

Suggestion 2.7: Avoiding the Creation of Objects Under Other Users' Private Schemas

NOTE

A private schema refers to a schema with the same name as the user when the user is created. This schema is only accessible to the user.

Impact of rule violation:

• When you create an object under someone else's private schema, the permissions for that object are determined by the schema owner.

Solution:

 Create objects under your own private schema to have full control over the object permissions.

2.3.4 TABLESPACE Object Design

Rule 2.8 Avoiding Tablespace Customization

Ⅲ NOTE

Impact of rule violation:

 In a distributed scenario, using a custom tablespace to create a table can result in the table data not being stored in a distributed manner by DN, leading to storage skew.

Solution

• Use the built-in default tablespace when creating a table object.

2.3.5 TABLE Object Design (Prioritized)

Rule 2.9: Selecting the Optimal Distribution Method and Columns During Table Creation

□ NOTE

Impact of rule violation:

• Incorrect distribution method and column selection can cause storage skew, deteriorate access performance, and even overload storage and computing resources.

Solution:

 When creating a table, it is important to use the **DISTRIBUTE BY** clause to explicitly specify the distribution method and distribution columns. The table below provides principles for selecting the distribution columns.

Table 2-2 Distribution column selection

Distribut ion Method	Description	Scenario
Hash	 Table data is distributed to each DN based on the mapping between hash values generated by distribution columns and DNs. Advantage: Each DN contains only part of data, which is space-saving. Disadvantage: The even distribution of data depends heavily on the selection of distribution columns. If the join condition does not include the distribution columns of each node, data communication between nodes will be required. 	Large tables and fact tables
RoundRo	 Table data is distributed to DNs in polling mode. Advantage: Each DN only contains a portion of the data, taking up a small amount of space. Data is evenly distributed in polling mode and does not rely on distribution columns, eliminating data skews. Disadvantage: Using distribution column conditions cannot eliminate or reduce inter-node communication. In this scenario, the performance is inferior to that of HASH. 	Large tables, fact tables, and tables without proper distribution columns

Distribut ion Method	Description	Scenario
Replicati on	Full data in a table is copied to each DN in the cluster. • Advantage: Each DN holds the complete data of the table. The JOIN operation avoids data communication between nodes, reducing network overhead and the overhead of starting and stopping the STREAM thread. • Disadvantage: Each DN retains complete table data, which is redundant and occupies more storage space.	Small tables and dimension tables

Rule 2.10 Selecting an Optimal Storage Type During Table Creation

□ NOTE

Impact of rule violation:

- Row-store tables are not properly used. As a result, the query performance is poor and resources are overloaded.
- Improper use of column-store tables causes CU expansion, poor performance, and resource overload.

Solution:

• When creating a table, use the **orientation** parameter to explicitly specify the storage type. The following table describes the rules for selecting a storage type.

Table 2-3 Storage type selection

Storag e Type	Applicable Scenario	Inapplicable Scenario
Row storage	 DML addition, deletion, and modification: scenarios with many UPDATE and DELETE operations DML query: point query (simple index-based query that returns only a few records) 	DML query: statistical analysis query (with mass data involved in GROUP and JOIN processes) CAUTION When creating a row-store table (orientation is set to row), do not specify the compress attribute or use a row-store compressed table.

Storag e Type	Applicable Scenario	Inapplicable Scenario
Colum n storage	 DML addition, deletion, and modification: INSERT batch import scenario (The number of data records imported to a single partition at a time is approximately 60,000 times the number of DNs or greater.) DML query: statistical analysis query (with mass data involved in GROUP and JOIN processes) 	 DML addition, deletion, and modification: scenarios with many UPDATE/DELETE operations or a small number of INSERT operations DML query: high-concurrency point query

Rule 2.11 Selecting an Optimal Partitioning Policy During Table Creation

◯ NOTE

Impact of rule violation:

Without partitioning, query performance and data governance efficiency will deteriorate. The larger the data volume, the greater the deterioration. The advantages of partitioning include:

- High query performance: The system queries only the relevant partitions rather than the whole table, improving the query efficiency.
- Improved data governance efficiency: In the data lifecycle management scenario, performing **TRUNCATE** or **DROP PARTITION** on historical partitions is much more efficient and effective than using **DELETE**.

Solution:

• Design partitions for tables that contain fields of the time type.

Table 2-4 Partitioning policy selection

Partitioning Policy	Description	Scenario
Range partitioning	Data is stored in different partitions based on the range of partition key values. The partition key ranges are consecutive but not overlapped.	 The date or time field is used as the partition key. Most queries contain partition keys as filter criteria. Periodically delete data based on the partition key.
List partitioning	Partitioning is performed based on a unique list of partition key values.	 A specific number of enumerated values are used as partition key values. Most queries contain partition keys as filter criteria.

Suggestion 2.12: Designing Table Columns for Fast and Accurate Queries

Ⅲ NOTE

Impact of rule violation:

• The system may have limited storage space and low query efficiency.

Solution:

- 1. Design the table columns well for fast queries.
 - If you can select an integer, do not select the floating point or character type.
 - When using variable-length character type, specify the maximum length based on data features.
- 2. Design the table columns well for accurate queries.
 - Use the time type instead of the character type to store time data.
 - Use the minimum numeric type that meets the requirements. Avoid using bigint if int or smallint is sufficient to save space.

3. Correctly use the constraints.

- Add NOT NULL constraints to columns that never have NULL values. The optimizer automatically optimizes the columns in certain scenarios.
- Do not use the **DEFAULT** constraint for fields that can be supplemented at the service layer. Otherwise, unexpected results may be generated during data loading.

4. Avoid unnecessary data type conversion.

- In tables that are logically related, columns having the same meaning should use the same data type.
- Different types of comparison operations cause data type conversion, which may cause index and partition pruning failures and affect query performance.

Suggestion 2.13: Avoiding the Usage of Auto-increment Columns or Data Types

Ⅲ NOTE

Impact of rule violation:

• When auto-increment sequences or data types are heavily used, the GTM may become overloaded and slow down sequence generation.

Solution:

- Set a UUID to obtain a unique ID.
- If the auto-increment sequence must be used and there is no strict requirement for increasing order, you can set the cache, for example, 1000, to reduce the pressure on GTM.

2.3.6 INDEX Object Design (Prioritized)

Rule 2.14: Creating Necessary Indexes and Selecting Optimal Columns and Sequences for Them

□ NOTE

Impact of rule violation:

- Redundant indexes consume unnecessary space and can impact data import efficiency.
- The column sequence in the composite index is incorrect, affecting the query efficiency.

Best practices:

The following conditions must be met when indexes are used:

- The index column should be a column commonly used for filtering or joining conditions.
- The index column should have more distinct values.
- When creating a multi-column combination index, prioritize columns with more distinct values.
- The number of indexes in a single table should be limited to less than five. You can control the number of indexes by combining them.
- In scenarios where data is added, deleted, or modified in batches, delete the index first
 and then add it back after the batch operation is complete to improve performance
 (real-time access may be affected).

Suggestion 2.15: Optimizing Performance by Choosing the Right Index Type and Avoiding Indexes for Column-Store Tables

■ NOTE

Impact of rule violation:

 Incorrect indexes do not improve column-store access and can negatively affect query performance.

Solution:

- Specify the appropriate index type when creating indexes, avoiding the default PSORT index.
- 2. In point queries where small amounts of data need to be retrieved from mass datasets, consider using a B-tree index.
- 3. For high range query performance, create a partial cluster key (PCK) to quickly filter and scan fact tables using the min/max sparse index. Comply with the following rules to create a PCK:
 - [Notice] Only one PCK can be created in a table. A PCK can contain multiple columns, preferably no more than two columns.
 - [Suggestion] Create a PCK for the filter condition column of the expression (e.g., col op const, where op is the operator =, >, >=, <=, and <, and const is a constant value).

2.3.7 VIEW Object Design

Suggestion 2.16: Limiting View Nesting to Three Layers

Impact of rule violation:

- Too many nested views can lead to unstable execution plans and unpredictable time consumption.
- The risk of rebuilding objects on which views depend is high and the probability of lock conflicts increases.

Solution:

Create views based on physical tables.

2.4 GaussDB(DWS) SQL Statement Development Specifications

2.4.1 DDL Operations

Suggestion 3.1: Avoiding Performing DDL Operations (Except CREATE) During Peak Hours or in Long Transactions

Ⅲ NOTE

Impact of rule violation:

DDL operations like **ALTER**, **DROP**, **TRUNCATE**, **REINDEX**, and **VACUUM FULL** have high lock levels and can block services during execution.

- During peak hours, these DDL operations with high lock levels should be avoided to prevent service blockage.
- Long transactions involving DDL operations with held or waited locks can also block services.

Solution:

 Choose off-peak hours or maintenance windows for DDL operations based on service periods. Specify the DDL execution environment and time consumption to avoid service blockage due to long lock waiting duration.

Rule 3.2: Specifying the Scope of Objects to Be Deleted When Using DROP

A DANGER

Impact of rule violation:

Be cautious when using **DROP OBJECT** (e.g., **DATABASE**, **USER/ROLE**, **SCHEMA**, **TABLE**, **VIEW**) as it may cause data loss, especially with **CASCADE** deletions.

- DROP DATABASE: deletes all objects in the database.
- DROP USER: deletes the USER object and its schemas and table objects.
- DROP SCHEMA: deletes all objects in the schema.
- **DROP TABLE**: deletes the **TABLE** object and the indexes and views that depend on it.

Solution:

 Exercise caution when performing the DROP operation and back up data in advance.

2.4.2 INSERT Operation

Rule 3.3: Replacing INSERT with COPY for Efficient Multi-Value Batch Insertion

◯ NOTE

Impact of rule violation:

 Parsing multiple values is time-consuming and resource-intensive, leading to low efficiency when importing data into the database.

Solution:

 Instead of using INSERT VALUES, the frontend should use APIs like CopyManager of JDBC.

Suggestion 3.4: Avoiding Performing Real-time INSERT Operations on Common Column-store Tables

NOTE

Impact of rule violation:

 Importing a small batch of data in real-time to a common column-store table can significantly expand the small CU, occupying a lot of storage space and impacting the query performance.

Solution:

- In real-time **INSERT** scenarios, evaluate the amount of data to be imported at once and the total amount of data. If the total amount of data is small, use row-store tables.
- In the real-time INSERT scenario, import around 60,000 data records to a single table, partition, or DN at a time. The minimum import batch is 5,000 data records.
- In the real-time INSERT scenario, use H-Store column-store tables (for clusters of version 8.3.0 or later).

2.4.3 UPDATE and DELETE Operations

Suggestion 3.5: Preventing Simultaneous Updates or Deletions of the Same Row in a Row-store Table

□ NOTE

Impact of rule violation:

 Concurrent UPDATE and DELETE operations on row-store tables may cause row lock blockage and distributed deadlocks, which can lead to service errors and performance degradation.

Solution:

• Group **UPDATE** and **DELETE** operations by primary key or distribution column. Perform parallel operations between groups while keeping operations within a group serial.

Suggestion 3.6: Avoiding Frequent or Simultaneous UPDATE and DELETE Operations on Column-store Tables

Impact of rule violation:

- Frequent UPDATE and DELETE operations on column-store tables can result in CU bloat, leading to large space occupation and decreased access performance.
- Concurrent UPDATE and DELETE operations on row-store tables may cause row lock blockage and distributed deadlocks, which can lead to service errors and performance degradation.

Solution:

- Design tables with frequent **UPDATE** and **DELETE** operations as row-store tables.
- Group **UPDATE** and **DELETE** operations by primary key or distribution column. Perform parallel operations between groups while keeping operations within a group serial.

2.4.4 SELECT Operation

Rule 3.7: Avoiding Executing SQL Statements That Do Not Support Pushdown

MOTE

GaussDB(DWS) uses a distributed architecture, and to achieve optimal performance, SQL statements need to be pushed down to utilize distributed computing resources.

Impact of rule violation:

• SQL statements that are not pushed down may experience poor execution performance and, in severe cases, can lead to CN resource bottlenecks, impacting overall services.

Solution:

• Do not use syntax or functions that cannot be executed near the data source. For details, see **Optimizing Statement Pushdown**.

Rule 3.8: Specifying Association Conditions when Multiple Tables Are Associated

■ NOTE

Impact of rule violation:

 If no association condition is specified when linking multiple tables, it will result in a Cartesian product calculation. This can lead to an expanded result set, posing risks of performance issues and resource overload.

Solution:

• Specify filter and association conditions for each table during the association process.

Rule 3.9: Ensuring Consistency of Data Types in Associated Fields across Multiple Tables

Ⅲ NOTE

Impact of rule violation:

 Ensure consistent data types for associated fields to avoid unnecessary type conversions, data redistribution issues, and hindered generation of optimal plans.

Solution:

Use the same data type for associated fields when tables are associated.

Suggestion 3.10: Avoiding Function Calculation on Association and Filter Condition Fields

◯ NOTE

Impact of rule violation:

In cases where function calculations are involved in association and filter conditions, the
optimizer may fail to obtain accurate field statistics, impacting execution performance.

Solution:

- When comparing association condition fields, preprocess the data before importing it into the database, especially when calculations are required for comparison.
- When filter criteria are compared with constants, perform function calculation only on constant columns. The following is an example:

SELECT id, from_image_id, from_person_id, from_video_id FROM face_data
WHERE SS.DEL_FLAG = 'N'
AND NVL(SS.DELETE_FLAG, 'N') = 'N'
The modification is as follows:
SELECT id, from_image_id, from_person_id, from_video_id FROM face_data
where SS.DEL_FLAG = 'N'
AND (SS.DELETE_FLAG = 'N' or SS.DELETE_FLAG is null)

Suggestion 3.11: Performing Pressure Tests and Concurrency Control for Resource-intensive SQL Statements

□ NOTE

Impact of rule violation:

 Storage and computing resources are overloaded, and the overall running performance deteriorates.

Solution:

A resource-intensive SQL statement contains:

- A large number of UNION ALL.
- A large number of AGGs (such as COUNT DISTINCT and MAX).
- A lot of JOIN operations for a large number of tables.
- A large number of **STREAM** operators (plan dimension).

Before rolling out, conduct pressure tests and implement concurrency control for these SQL statements. If the resource capacity is exceeded, optimizing the service should be prioritized before reassessing the rollout plan.

Rule 3.12: Avoiding Excessive COUNT Operations on Large Row-store Tables

■ NOTE

If SSDs or other high-performance disk types are used, it may not be necessary to adhere strictly to this rule, but it is still crucial to monitor the I/O consumption.

Impact of rule violation:

 Performing frequent COUNT operations on large row-store tables can consume a significant amount of I/O resources, potentially leading to performance issues if an I/O bottleneck occurs.

Solution:

• Reduce the frequency of **COUNT** operations, use result caching, and collect statistics by partition to minimize I/O consumption.

Suggestion 3.13: Avoid Getting Large Result Sets (Except for Data Exports)

◯ NOTE

Impact of rule violation:

• If you do not need to view all the results, querying ultra-large result sets becomes inefficient and wasteful in terms of resources.

Solution:

- Use the LIMIT clause to retrieve only the necessary result segments.
- Use a cursor to obtain the result sets by segment and set an appropriate value for **FETCH SIZE** if you need to query a large number of result sets.

Suggestion 3.14: Avoiding the Usage of SELECT * for Queries

□ NOTE

Impact of rule violation:

 Querying unnecessary columns increases the computing load and wastes computing resources.

Solution:

• Clearly list the fields required for the query in the **SELECT** statement to improve the query performance.

Suggestion 3.15: Using WITH RECURSIVE with Defined Termination Condition for Recursion

Impact of rule violation:

- In cases where there is no specific termination condition, recursive operations can enter an infinite loop.
- Recursive operations generate duplicate data and occupy excessive resources.

Solution:

 Design proper termination conditions based on the volume and characteristics of the data in the service table.

Suggestion 3.16: Setting Schema Prefix for Table and Function Access

◯ NOTE

Impact of rule violation:

• If the schema name prefix is not specified, the search will be performed sequentially across all tablespaces based on the tablespace list in the current **search_path**. This can lead to accessing unexpected tables due to schema switchover.

Solution:

• To enhance readability, stability, and portability, explicitly specify the schema prefix as **SCHEMA.** when accessing tables and function objects.

Suggestion 3.17: Identifying an SQL Statement with a Unique SQL Comment

Ⅲ NOTE

Impact of rule violation:

• The service's source tracing capability is limited. You can only verify it with R&D engineers using the database, user name, and client IP address.

Solution:

- You are advised to use **query_band**. The following is an example: SET query_band='JobName=abc;AppName=test;UserName=user';
- Add a unique comment for each SQL statement to facilitate troubleshooting and application performance analysis. The following is an example of such comment.
 /* Module name_Tool name_Job name_Step */, for example, /* mca_python_xxxxxx_step1 */ insert into xxx select ... from

Recommendation 3.18: Restricting SQL Statements to 64 KB in Length

◯ NOTE

Impact of rule violation:

 SQL parsing is time-consuming and difficult to maintain. Frequent execution of SQL statements leads to severe log expansion.

Solution

• Set a 64 KB limit on SQL statements.

2.5 GaussDB(DWS) Foreign Table Function Development Specifications

Rule 4.1 Deploying GDS on an Independent Server Outside the GaussDB(DWS) Cluster

Impact of rule violation:

If GDS is deployed in a GaussDB(DWS) cluster, it contends for resources with CNs or DNs in the cluster, which leads to a decline in the performance of both GDS and CNs or DNs.

Solution

- Deploy GDS on an independent server outside the GaussDB(DWS) cluster.
- Ensure that the disk capacity of the GDS server and the network bandwidth between the GDS server and the GaussDB(DWS) cluster are planned according to the requirements.

Rule 4.2 Avoiding Concurrent Access to Multiple Collaborative Analysis Foreign Tables Across Clusters

□ NOTE

Principle description: When cluster A accesses data in cluster B through collaborative analysis, all DNs in cluster A establish connections and active sessions with the CNs in cluster B.

Impact of rule violation:

The CNs in cluster B are overloaded. As a result, the number of connections and active sessions exceeds the threshold, and the access is abnormal.

Solution:

Use foreign tables to access a single table instead of performing concurrent queries on multiple foreign tables. If it is not possible to avoid concurrent queries, calculate and limit the number of concurrent queries based on the number of DNs in cluster A and the normal service volume of cluster B. Additionally, increasing the values of <code>max_active_statements</code> and <code>max_connections</code> can help solve the problem.

2.6 GaussDB(DWS) Stored Procedure Development Specifications

Suggestion 5.1: Simplifying Stored Procedures and Avoiding Nesting

□ NOTE

Impact of rule violation:

• The maintenance cost for complex and nested stored procedures is high, making fault locating and recovery time-consuming.

Solution:

- Avoid using stored procedures altogether or limit their usage to a single layer. Nested stored procedures should be avoided.
- Create a corresponding log table for the stored procedure design and record information before and after key steps in the log table. Follow the steps below to implement this.

Saving and viewing logs

Step 1 Create a log table.

```
CREATE TABLE func_exec_log
(
id varchar2(32) default lower(sys_guid()),
pro_name varchar2(60),
exec_times int,
log_date date,
deal_date date,
log_mesage text
);
```

Step 2 Create a table and import data.

```
CREATE TABLE demo_table(data_id int, data_number int);
INSERT INTO demo_table values(generate_series(1,1000),generate_series(1,1000));
```

Step 3 Create a service stored procedure.

```
CREATE OR REPLACE FUNCTION demo_table_process(out exe_info text)
LANGUAGE plpgsgl
AS $$
declare v count int;
pro_result text;
fun_name text;
exec_times int;
beain
fun_name := 'demo_table_process';
select nvl(max(exec_times), '0') + 1 into exec_times from func_exec_log where pro_name = fun_name;
-- Insert data into the service table.
insert into demo_table values (dbms_random.value(1, 1000)::int,generate_series(1,
dbms_random.value(10000, 20000)::int));
get diagnostics v_count = ROW_COUNT;
exe_info = sysdate || '# step1:insert count:' || v_count || ' rows;';
-- Delete specified data from a service table.
delete from demo_table where data_id = dbms_random.value(1, 1000)::int;
get diagnostics v_count = ROW_COUNT;
exe_info = exe_info || sysdate || '# step2:delete count:' || v_count || ' rows;';
-- Update service table data.
update demo_table set data_number = dbms_random.value(1, 100)::int where data_id =
dbms_random.value(1, 1000)::int;
exe_info = exe_info || sysdate || '# step3:update count:' || sql%rowcount || ' rows';
-- Record logs either before the entire program ends or after each step completes. You can also create a
function specifically for logging purposes.
insert into func_exec_log(pro_name, exec_times, log_date, deal_date, log_mesage) values
```

```
(fun_name,exec_times,sysdate,split_part(regexp_split_to_table(exe_info, ';'), '#',
1),split_part(regexp_split_to_table(exe_info, ';'), '#', 2));
-- EXCEPTION is used to ensure that logs can be properly recorded when the insertion, update, or deletion exits abnormally.

EXCEPTION
WHEN OTHERS THEN
pro_result := exe_info || sysdate || '# exception error message is: ' || sqlerrm; insert into func_exec_log(pro_name, exec_times, log_date, deal_date, log_mesage)
values(fun_name,exec_times,sysdate,split_part(regexp_split_to_table(pro_result, ';'), '#',
1),split_part(regexp_split_to_table(pro_result, ';'), '#', 2));
END; $$;
```

Step 4 Invoke the stored procedure (normal execution).

SELECT demo table process();

Step 5 View the created log table to check the service running status.

SELECT * FROM func_exec_log ORDER BY log_date desc,deal_date,log_mesage;

demodb=# select * from func_exec_log order by log_date id pro_name	desc,deal dat exec_times		deal_date	log_mesage
637343d9f2f10ec605c7687ff700fffe demo_table_process 637343d9fe850e3105c8687ff700fffe demo_table_process 637343d9068a0fd805c9687ff700fffe demo_table_process (3 rows)	† 1 1 1	2022-11-15 15:46:34	2022-11-15 15:46:33	

Step 6 Invoke the stored procedure again to construct an execution exception.

SELECT demo_table_process(); -- Delete the data_number column of demo_table to construct an exception, and then call the stored procedure again.

Step 7 View the log to check the service running status.

----End

Rule 5.2: Avoiding Non-CREATE DDL Operations in Stored Procedures

☐ NOTE

Impact of rule violation:

A stored procedure is a large transaction. If a non-CREATE DDL operation, especially one
with a high lock level, is executed, it can block external access to related tables during
the stored procedure's execution window.

Solution:

Avoid using non-CREATE DDL operations within stored procedures whenever possible. If
there is a necessity to use such operations, carefully assess the duration of the stored
procedures and the potential impact of the DDL operations. It is advised to schedule
non-CREATE DDL operations during off-peak hours when external access services are
less active.

2.7 Detailed Design Rules for GaussDB(DWS) Objects

2.7.1 GaussDB(DWS) Database Object Naming Rules

The name of a database object must contain 1 to 63 characters, start with a letter or underscore (_), and can contain letters, digits, underscores (_), and dollar signs (\$). If the database uses GBK, UTF8, or SQL_ASCII, object names can include Chinese characters. With UTF8 or SQL_ASCII, each Chinese character counts as three, with a limit of 21 characters. With GBK, each Chinese character counts as two, with a limit of 31 characters. The Latin1 character set does not support Chinese characters. The character set format is specified during database creation. For details, see CREATE DATABASE.

• [Proposal] Do not use reserved or non-reserved keywords to name database objects.

You can use **SELECT * FROM pg_get_keywords ()** to obtain GaussDB(DWS) keywords. For other ways to obtain the keywords, see **Keywords** in the *SQL Syntax Reference*.

- [Proposal] Do not use strings enclosed in double quotation marks to define database object names. In GaussDB(DWS), double quotation marks are used to specify that the enclosed database object names are case sensitive. Case sensitivity of database object names makes problem location difficult.
- [Proposal] Use the same naming format for database objects.
 - In a system undergoing incremental development or service migration, you are advised to comply with its historical naming conventions.
 - A database object name consists of letters, digits, and underscores (_);
 and cannot start with a digit. You are advised to use multiple words separated with hyphens (-).
 - You are advised to use intelligible names and common acronyms or abbreviations for database objects. Acronyms or abbreviations that are generally understood are recommended. For example, you can use English words indicating actual business terms. The naming format should be consistent within a cluster.
 - A variable name must be descriptive and meaningful. It must have a prefix indicating its type.
- [Proposal] The name of a table object should indicate its main characteristics, for example, whether it is an ordinary, temporary, or unlogged table.
 - An ordinary table name should indicate the business relevant to a data set.
 - Temporary tables are named in the format of **tmp**_Suffix.
 - Unlogged tables are named in the format of **ul** Suffix.
 - Foreign tables are named in the format of f_Suffix.

2.7.2 GaussDB(DWS) Database Object Design Rules

2.7.2.1 GaussDB(DWS) Database and Schema Design Rules

In GaussDB(DWS), services can be isolated by databases and schemas. Databases share little resources and cannot directly access each other. Connections to and permissions on them are also isolated. Schemas share more resources than databases do. User permissions on schemas and subordinate objects can be controlled using the **GRANT** and **REVOKE** syntax.

- You are advised to use schemas to isolate services for convenience and resource sharing.
- It is recommended that system administrators create schemas and databases and then assign required permissions to users.

Database Design Suggestions

 Create databases as required. Do not use the default gaussdb database of a cluster.

- Create a maximum of three user-defined databases in a cluster.
- To make your database encoding compatible with most characters, you are advised to use the UTF-8 encoding when creating a database.
- Exercise caution when you set ENCODING and DBCOMPATIBILITY
 configuration items during database creation. In GaussDB(DWS),
 DBCOMPATIBILITY can be set to TD, Oracle, or MySQL to be compatible
 with Teradata, Oracle, or MySQL syntax, respectively. Syntax behavior may
 vary with the three modes. For details, see Syntax Compatibility Differences
 Among Oracle, Teradata, and MySQL.
- By default, a database owner has all permissions for all objects in the database, including the deletion permission. Exercise caution when using the deletion permission.

Schema Design Suggestions

- To let a user access an object in a schema, grant the **usage** permission and the permissions for the object to the user, unless the user has the **sysadmin** permission or is the schema owner.
- To let a user create an object in the schema, grant the **CREATE** permission for the schema to the user.
- By default, a schema owner has all permissions for all objects in the schema, including the deletion permission. Exercise caution when using the deletion permission.

2.7.2.2 GaussDB(DWS) Table Design Rules

GaussDB(DWS) uses a distributed architecture. Data is distributed on DNs. Comply with the following principles to properly design a table:

- [Notice] Evenly distribute data on each DN to prevent data skew. If most data is stored on several DNs, the effective capacity of a cluster decreases. Select a proper distribution column to avoid data skew.
- [Notice] Evenly scan each DN when querying tables. Otherwise, DNs most frequently scanned will become the performance bottleneck. For example, when you use equivalent filter conditions on a fact table, the nodes are not evenly scanned.
- [Notice] Reduce the amount of data to be scanned. You can use the pruning mechanism of a partitioned table.
- [Notice] Minimize random I/O. By clustering or local clustering, you can sequentially store hot data, converting random I/O to sequential I/O to reduce the cost of I/O scanning.
- [Notice] Try to avoid data shuffling. To shuffle data is to physically transfer it from one node to another. This unnecessarily occupies many network resources. To reduce network pressure, locally process data, and to improve cluster performance and concurrency, you can minimize data shuffling by using proper association and grouping conditions.

Selecting a Storage Mode

[Proposal] Selecting a storage mode is the first step in defining a table. The storage mode mainly depends on the user's service type. For details, see **Table 2-5**.

Table 2-5 Table storage modes and scenarios

Storage Mode	Application Scenarios
Row storage	Point queries (simple index-based queries that only return a few records)
	Scenarios requiring frequent addition, deletion, and modification
Column storage	Statistical analysis queries (requiring a large number of association and grouping operations)
	Ad hoc queries (using uncertain query conditions and unable to utilize indexes to scan row-store tables)

When creating a table for analysis, make sure to set the **ORIENTATION** to column storage explicitly.

```
CREATE TABLE public.t1
(
id integer not null,
data integer,
age integer
)
WITH (ORIENTATION =COLUMN);
```

Selecting a Distribution Mode

[Proposal] Comply with the following rules to distribute table data.

Table 2-6 Table distribution modes and scenarios

Distribution Mode	Description	Application Scenarios
Hash	Table data is distributed on all DNs in a cluster by hash.	Fact tables containing a large amount of data
Replication	Full data in a table is stored on every DN in a cluster.	Dimension tables and fact tables containing a small amount of data
Round-robin	Each row of the table is sent to each DN in turn. Therefore, data is evenly distributed on each DN.	Fact tables that contain a large amount of data and cannot find a proper distribution column in hash mode

Selecting a Partitioning Mode

Comply with the following rules to partition a table containing a large amount of data:

- [Proposal] Create partitions on columns that indicate certain ranges, such as dates and regions.
- [Proposal] A partition name should show the data characteristics of a partition. For example, its format can be Keyword+Range characteristics.
- [Proposal] Set the upper limit of a partition to **MAXVALUE** to prevent data overflow.

The example of a partitioned table definition is as follows:

```
CREATE TABLE staffS_p1
 staff_ID
           NUMBER(6) not null,
 FIRST_NAME VARCHAR2(20),
LAST_NAME VARCHAR2(25),
           VARCHAR2(25),
 EMAIL
 PHONE_NUMBER VARCHAR2(20),
 HIRE_DATE DATE,
 employment_ID VARCHAR2(10),
             NUMBER(8,2),
 SALARY
 COMMISSION_PCT NUMBER(4,2),
 MANAGER ID NUMBER(6),
 section_ID NUMBER(4)
PARTITION BY RANGE (HIRE_DATE)
 PARTITION HIRE_19950501 VALUES LESS THAN ('1995-05-01 00:00:00'),
 PARTITION HIRE_19950502 VALUES LESS THAN ('1995-05-02 00:00:00'),
 PARTITION HIRE_maxvalue VALUES LESS THAN (MAXVALUE)
);
```

Selecting a Distribution Key

Selecting a distribution key is important for a hash table. An improper distribution key may cause data skew. As a result, the I/O load is heavy on several DNs, affecting the overall query performance. After you select a distribution policy for a hash table, check for data skew to ensure that data is evenly distributed. Comply with the following rules to select a distribution key:

- [Proposal] Select a column containing discrete data as the distribution key, so that data can be evenly distributed on each DN. If a single column is not discrete enough, consider using multiple columns as distribution keys. You can select the primary key of a table as the distribution key. For example, in an employee information table, select the certificate number column as the distribution key.
- [Proposal] If the first rule is met, do not select a column having constant filter conditions as the distribution key. For example, in a query on the dwcjk table, if the zqdh column contains the constant filter condition zqdh='000001', avoid selecting the zqdh column as the distribution key.
- [Proposal] If the first and second rules are met, select the join conditions in a query as distribution keys. If a join condition is used as a distribution key, the data involved in a join task is locally distributed on DNs, which greatly reduces the data flow cost among DNs.

2.7.2.3 GaussDB(DWS) Column Design Rules

Selecting a Data Type

Comply with the following rules to improve query efficiency when you design columns:

- [Proposal] Use the most efficient data types allowed.
 - If all of the following number types provide the required service precision, they are recommended in descending order of priority: integer, floating point, and numeric.
- [Proposal] In tables that are logically related, columns having the same meaning should use the same data type.
- [Proposal] For string data, you are advised to use variable-length strings and specify the maximum length. To avoid truncation, ensure that the specified maximum length is greater than the maximum number of characters to be stored. You are not advised to use CHAR(n), BPCHAR(n), NCHAR(n), or CHARACTER(n), unless you know that the string length is fixed.

For details about string types, see Common String Types.

Common String Types

Every column requires a data type suitable for its data characteristics. The following table lists common string types in GaussDB(DWS).

Table 2-7 Common string types

Parameter	Description	Max. Storage Capacity
CHAR(n)	Fixed-length string, where <i>n</i> indicates the stored bytes. If the length of an input string is smaller than <i>n</i> , the string is automatically padded to <i>n</i> bytes using NULL characters.	10 MB
CHARACTER(n)	Fixed-length string, where <i>n</i> indicates the stored bytes. If the length of an input string is smaller than <i>n</i> , the string is automatically padded to <i>n</i> bytes using NULL characters.	10 MB
NCHAR(n)	Fixed-length string, where <i>n</i> indicates the stored bytes. If the length of an input string is smaller than <i>n</i> , the string is automatically padded to <i>n</i> bytes using NULL characters.	10 MB

Parameter	Description	Max. Storage Capacity
BPCHAR(n)	Fixed-length string, where <i>n</i> indicates the stored bytes. If the length of an input string is smaller than <i>n</i> , the string is automatically padded to <i>n</i> bytes using NULL characters.	10 MB
VARCHAR(n)	Variable-length string, where <i>n</i> indicates the maximum number of bytes that can be stored.	10 MB
CHARACTER VARYING(n)	Variable-length string, where <i>n</i> indicates the maximum number of bytes that can be stored. This data type and VARCHAR(n) are different representations of the same data type.	10 MB
VARCHAR2(n)	Variable-length string, where <i>n</i> indicates the maximum number of bytes that can be stored. This data type is added to be compatible with the Oracle database, and its behavior is the same as that of VARCHAR(n).	10 MB
NVARCHAR2(n)	Variable-length string, where <i>n</i> indicates the maximum number of bytes that can be stored.	10 MB
TEXT	Variable-length string. Its maximum length is 8203 bytes less than 1 GB.	8203 bytes less than 1 GB

2.7.2.4 GaussDB(DWS) Constraint Design Rules

DEFAULT and NULL Constraints

- [Proposal] If all the column values can be obtained from services, you are not advised to use the **DEFAULT** constraint, because doing so will generate unexpected results during data loading.
- [Proposal] Add **NOT NULL** constraints to columns that never have NULL values. The optimizer automatically optimizes the columns in certain scenarios.
- [Proposal] Explicitly name all constraints excluding NOT NULL and DEFAULT.

Partial Cluster Key

A partial cluster key (PCK) is a local clustering technology used for column-store tables. After creating a PCK, you can quickly filter and scan fact tables using min

or max sparse indexes in GaussDB(DWS). Comply with the following rules to create a PCK:

- [Notice] Only one PCK can be created in a table. A PCK can contain multiple columns, preferably no more than two columns.
- [Proposal] Create a PCK on simple expression filter conditions in a query. Such filter conditions are usually in the form of **col op const**, where **col** specifies a column name, **op** specifies an operator (such as =, >, >=, <=, and <), and **const** specifies a constant.
- [Proposal] If the preceding conditions are met, create a PCK on the column having the least distinct values.

Unique Constraint

- [Notice] Both row-store and column-store tables support unique constraints.
- [Proposal] The constraint name should indicate that it is a unique constraint, for example, **UNI** *Included columns*.

Primary Key Constraint

- [Notice] Both row-store and column-store tables support the primary key constraint.
- [Proposal] The constraint name should indicate that it is a primary key constraint, for example, **PK***Included columns*.

Check Constraint

- [Notice] Check constraints can be used in row-store tables but not in columnstore tables.
- [Proposal] The constraint name should indicate that it is a check constraint, for example, **CK***Included columns*.

2.7.2.5 Design Rules for GaussDB(DWS) Views and Associated Tables

View Design

- [Proposal] Do not nest views unless they have strong dependency on each other.
- [Proposal] Try to avoid sort operations in a view definition.

Joined Table Design

- [Proposal] Minimize joined columns across tables.
- [Proposal] Joined columns should use the same data type.
- [Proposal] The names of associated fields should show the associations. For example, they can use the same name.

2.7.3 GaussDB(DWS) JDBC Configuration Rules

Currently, third-party tools are connected to GaussDB(DWS) trough JDBC. This section describes the precautions for configuring the tools.

Connection Parameters

 [Notice] When a third-party tool connects to GaussDB(DWS) through JDBC, JDBC sends a connection request to GaussDB(DWS). By default, the following parameters are added. For details, see the implementation of the ConnectionFactoryImpl JDBC code.

These parameters may cause the JDBC and gsql clients to display inconsistent data, for example, date data display mode, floating point precision representation, and timezone.

If the result is not as expected, you are advised to explicitly set these parameters in the Java connection setting.

- [Proposal] When connecting to the database through JDBC, ensure that the following two time zones are the same:
 - Time zone of the host where the JDBC client is located
 - Time zone of the host where the GaussDB(DWS) server is located

fetchsize

[Notice] To use **fetchsize** in applications, disable the **autocommit** switch. Enabling the **autocommit** switch makes the **fetchsize** configuration invalid.

autocommit

[Proposal] It is recommended that you enable the **autocommit** switch in the code for connecting to GaussDB(DWS) by the JDBC. If **autocommit** needs to be disabled to improve performance or for other purposes, applications need to ensure their transactions are committed. For example, explicitly commit translations after specifying service SQL statements. Particularly, ensure that all transactions are committed before the client exits.

Connection Releasing

[Proposal] You are advised to use connection pools to limit the number of connections from applications. Do not connect to a database every time you run an SQL statement.

[Proposal] After an application completes its tasks, disconnect its connection to GaussDB(DWS) to release occupied resources. You are advised to set the session timeout interval in the task.

[Proposal] Reset the session environment before releasing connections to the JDBC connection tool. Otherwise, historical session information may cause object conflicts.

• If GUC parameters are set in the connection, before you return the connection to the connection pool, run **SET SESSION AUTHORIZATION DEFAULT;RESET ALL;** to clear the connection status.

• If a temporary table is used, delete it before you return the connection to the connection pool.

CopyManager

[Proposal] In the scenario where the ETL tool is not used and real-time data import is required, it is recommended that you use the CopyManager interface driven by the GaussDB(DWS) JDBC to import data in batches during application development.

For how to use CopyManager, see CopyManager.

2.7.4 GaussDB(DWS) SQL Writing Rules

DDL

- [Proposal] In GaussDB(DWS), you are advised to execute DDL operations, such as creating table or making comments, separately from batch processing jobs to avoid performance deterioration caused by many concurrent transactions.
- [Proposal] Execute data truncation after unlogged tables are used because GaussDB(DWS) cannot ensure the security of unlogged tables in abnormal scenarios.
- [Proposal] Suggestions on the storage mode of temporary and unlogged tables are the same as those on base tables. Create temporary tables in the same storage mode as the base tables to avoid high computing costs caused by hybrid row and column correlation.
- [Proposal] The total length of an index column cannot exceed 50 bytes. Otherwise, the index size will increase greatly, resulting in large storage cost and low index performance.
- [Proposal] Do not delete objects using **DROP...CASCADE**, unless the dependency between objects is specified. Otherwise, the objects may be deleted by mistake.

Data Loading and Uninstalling

- [Proposal] Provide the inserted column list in the insert statement. Example: INSERT INTO task(name,id,comment) VALUES ('task1','100','100th task');
- [Proposal] After data is imported to the database in batches or the data increment reaches the threshold, you are advised to analyze tables to prevent the execution plan from being degraded due to inaccurate statistics.
- [Proposal] To clear all data in a table, you are advised to use **TRUNCATE TABLE** instead of **DELETE TABLE**. **DELETE TABLE** is not efficient and cannot release disk space occupied by the deleted data.

Type conversion

- [Proposal] Perform type coercion to convert data types. If you perform implicit conversion, the result may differ from expected.
- [Proposal] During data query, explicitly specify the data type for constants, and do not attempt to perform any implicit data type conversion.

• [Notice] In Oracle compatibility mode, null strings will be automatically converted to NULL during data import. If a null string needs to be reserved, you need to create a database that is compatible with Teradata.

Query Operation

- [Proposal] Do not return a large number of result sets to a client except the ETL program. If a large result set is returned, consider modifying your service design.
- [Proposal] Perform DDL and DML operations encapsulated in transactions. Operations like table truncation, update, deletion, and dropping, cannot be rolled back once committed. You are advised to encapsulate such operations in transactions so that you can roll back the operations if necessary.
- [Proposal] During query compilation, you are advised to list all columns to be queried and avoid using *. Doing so reduces output lines, improves query performance, and avoids the impact of adding or deleting columns on frontend service compatibility.
- [Proposal] During table object access, add the schema prefix to the table object to avoid accessing an unexpected table due to schema switchover.
- [Proposal] The cost of joining more than eight tables or views, especially full
 joins, is difficult to be estimated. You are advised to use the WITH TABLE AS
 statement or other methods to create interim tables to improve the
 readability of SQL statements.
- [Proposal] Do not use Cartesian products or full joins. Cartesian products and full joins will result in a sharp expansion of result sets and poor performance.
- [Notice] Only IS NULL and IS NOT NULL can be used to determine NULL value comparison results. If any other method is used, NULL is returned. For example, NULL instead of expected Boolean values is returned for NULL<>NULL, NULL=NULL, and NULL<>1.
- [Notice] Do not use count(col) instead of count(*) to count the total number of records in a table. count(*) counts the NULL value (actual rows) while count (col) does not.
- [Notice] While executing count(col), the number of NULL record rows is counted as 0. While executing sum(col), NULL is returned if all records are NULL. If not all the records are NULL, the number of NULL record rows is counted as 0.
- [Notice] To count multiple columns using count(), column names must be
 enclosed with parentheses. For example, count ((col1, col2, col3)). Note:
 When multiple columns are used to count the number of NULL record rows, a
 row is counted even if all the selected columns are NULL. The result is the
 same as that when count(*) is executed.
- [Notice] Null records are not counted when count(distinct col) is used to calculate the number of non-null columns that are not repeated.
- [Notice] If all statistical columns are NULL when count(distinct (col1,col2,...)) is used to count the number of unique values in multiple columns, Null records are also counted, and the records are considered the same.
- [Notice] When constants are used to filter data, the system searches for functions used for calculating these two data types based on the data types of the constants and matched columns. If no function is found, the system

converts the data type implicitly. Then, the system searches for a function used for calculating the converted data type.

SELECT * FROM test WHERE timestamp_col = 20000101;

In the preceding example, if **timestamp_col** is the timestamp type, the system first searches for the function that supports the "equal" operation of the timestamp and int types (constant numbers are considered as the int type). If no such function is found, the **timestamp_col** data and constant numbers are implicitly converted into the text type for calculation.

- [Proposal] Do not use scalar subquery statements. A scalar subquery appears
 in the output list of a SELECT statement. In the following example, the part
 enclosed in parentheses is a scalar subquery statement:
 SELECT id, (SELECT COUNT(*) FROM films f WHERE f.did = s.id) FROM staffs_p1 s;
 - Scalar subqueries often result in query performance deterioration. During application development, scalar subqueries need to be converted into equivalent table associations based on the service logic.
- [Proposal] In **WHERE** clauses, the filtering conditions should be sorted. The condition that few records are selected for reading (the number of filtered records is small) is listed at the beginning.
- [Proposal] Filtering conditions in WHERE clauses should comply with unilateral rules. That is, when the column name is placed on one side of a comparison operator, the optimizer automatically performs pruning optimization in some scenarios. Filtering conditions in a WHERE clause will be displayed in col op expression format, where col indicates a table column, op indicates a comparison operator, such as = and >, and expression indicates an expression that does not contain a column name. For example:
 SELECT id, from_image_id, from_person_id, from_video_id FROM face_data WHERE current_timestamp(6) time < '1 days'::interval;</p>

The modification is as follows:

SELECT id, from_image_id, from_person_id, from_video_id FROM face_data where time > current_timestamp(6) - '1 days'::interval;

- [Proposal] Do not perform unnecessary sorting operations. Sorting requires a
 large amount of memory and CPU. If service logic permits, ORDER BY and
 LIMIT can be combined to reduce resource overhead. By default, data in
 GaussDB(DWS) is sorted by ASC & NULL LAST.
- [Proposal] When the ORDER BY clause is used for sorting, specify sorting modes (ASC or DESC), and use NULL FIRST or NULL LAST for NULL record sorting.
- [proposal] Do not rely on only the LIMIT clause to return the result set displayed in a specific sequence. Combine ORDER BY and LIMIT clauses for some specific result sets and use offset to skip specific results if necessary.
- [Proposal] If the service logic is accurate, you are advised to use UNION ALL instead of UNION.
- [Proposal] If a filtering condition contains only an OR expression, convert the
 OR expression to UNION ALL to improve performance. SQL statements that
 use OR expressions cannot be optimized, resulting in slow execution. Example:
 SELECT * FROM scdc.pub_menu

WHERE (cdp= 300 AND inline=301) OR (cdp= 301 AND inline=302) OR (cdp= 302 AND inline=301);

Convert the statement to the following:

SELECT * FROM scdc.pub_menu WHERE (cdp= 300 AND inline=301) union all SELECT * FROM scdc.pub menu WHERE (cdp= 301 AND inline=302) union all SELECT * FROM scdc.pub_menu WHERE (cdp= 302 AND inline=301);

- [Proposal] If an in(val1, val2, va...) expression contains a large number of columns, you are advised to replace it with the in (values (va1), (val2), (val3...) statement. The optimizer will automatically convert the IN constraint into a non-correlated subquery to improve the query performance.
- [Proposal] Replace (not) in with (not) exist when associated columns do not contain NULL values. For example, in the following query statement, if the T1.C1 column does not contain any NULL value, add the NOT NULL constraint to the T1.C1 column, and then rewrite the statements.
 SELECT * FROM T1 WHERE T1.C1 NOT IN (SELECT T2.C2 FROM T2);

Rewrite the statement as follows:

SELECT * FROM T1 WHERE NOT EXISTS (SELECT * FROM T1,T2 WHERE T1.C1=T2.C2);

- If the value of the T1.C1 column will possibly be NULL, the preceding rewriting cannot be performed.
- If T1.C1 is the output of a subquery, check whether the output is NOT NULL based on the service logic.
- [Proposal] Use cursors instead of the LIMIT OFFSET syntax to perform
 pagination queries to avoid resource overheads caused by multiple executions.
 A cursor must be used in a transaction, and you must disable it and commit
 transaction once the query is finished.

2.7.5 Rules for Using Custom GaussDB(DWS) External Functions (pgSQL/Java)

- [Notice] Java UDFs can perform some Java logic calculation. Do not encapsulate services in Java UDFs.
- [Notice] Do not connect to a database in any way (for example, by using JDBC) in Java functions.
- [Notice] Only the data types listed in the following table can be used. Userdefined types and complex data types (Java Array and derived classes) are not supported.
- [Notice] User-defined aggregation functions (UDAFs) and user-defined tablegenerating functions (UDTFs) are not supported.

Table 2-8 PL/Java mapping for default data types

GaussDB(DWS)	Java
BOOLEAN	boolean
"char"	byte
bytea	byte[]
SMALLINT	short
INTEGER	int

GaussDB(DWS)	Java	
BIGINT	long	
FLOAT4	float	
FLOAT8	double	
CHAR	java.lang.String	
VARCHAR	java.lang.String	
TEXT	java.lang.String	
name	java.lang.String	
DATE	java.sql.Timestamp	
TIME	java.sql.Time (stored value treated as local time)	
TIMETZ	java.sql.Time	
TIMESTAMP	java.sql.Timestamp	
TIMESTAMPTZ	java.sql.Timestamp	

2.7.6 Rules for Using GaussDB(DWS) PL/pgSQL

General Principles

- 1. Development shall strictly comply with design documents.
- 2. Program modules shall be highly cohesive and loosely coupled.
- 3. Proper, comprehensive troubleshooting measures shall be developed.
- 4. Code shall be reasonable and clear.
- 5. Program names shall comply with a unified naming rule.
- 6. Fully consider the program efficiency, including the program execution efficiency and database query and storage efficiency. Use efficient and effective processing methods.
- 7. Program comments shall be detailed, correct, and standard.
- 8. The commit or rollback operation shall be performed at the end of a stored procedure, unless otherwise required by applications.
- 9. Programs shall support 24/7 processing. In the case of an interruption, the applications shall provide secure, easy-to-use resuming features.
- 10. Application output shall be standard and simple. The output shall show the progress, error description, and execution results for application maintenance personnel, and provide clear and intuitive reports and documents for business personnel.

Programming Principles

- 1. Use bound variables in SQL statements in the PL/pgSQL.
- 2. **RETURNING** is recommended for SQL statements in PL/pgSQL.
- 3. Principles for using stored procedures:
 - a. Do not use more than 50 output parameters of the Varchar or Varchar2 type in a stored procedure.
 - b. Do not use the LONG type for input or output parameters.
 - c. Use the CLOB type for output strings that exceed 10 MB.
- 4. Variable declaration principles:
 - a. Use **%TYPE** to declare a variable that has the same meaning as that of a column or variable in an application table.
 - b. Use **%ROWTYPE** to declare a record that has the same meaning as that of a row in an application table.
 - c. Each line of a variable declaration shall contain only one statement.
 - d. Do not declare variables of the LONG type.
- 5. Principles for using cursors:
 - a. Explicit cursors shall be closed after being used.
 - b. A cursor variable shall be closed after being used. If the cursor variable needs to transfer data to an invoked application, the cursor shall be closed in the application. If the cursor variable is used only in a stored procedure, the cursor shall be closed explicitly.
 - c. Before using **DBMS_SQL.CLOSE_CURSOR** to close a cursor, use **DBMS_SQL.IS_OPEN** to check whether the cursor is open.
- 6. Principles for collections:
 - a. You are advised to use the **FOR ALL** statement instead of the **FOR** loop statement to reference elements in a collection.
- 7. Principles for using dynamic statements:
 - a. Dynamic SQL shall not be used in the transaction programs of online systems.
 - b. Dynamic SQL statements can be used to implement DDL statements and system control commands in PL/pgSQL.
 - c. Variable binding is recommended.
- 8. Principles for assembling SQL statements:
 - a. You are advised to use bound variables to assemble SQL statements.
 - b. If the conditions for assembling SQL statements contain external input sources, the characters in the input conditions shall be checked to prevent attacks.
 - c. In a PL/pgSQL script, the length of a single line of code cannot exceed 2499 characters.
- 9. Principles for using triggers:
 - a. Triggers can be used to implement availability design in scenarios where differential data logs are irrelevant to service processing.

b. Do not use triggers to implement service processing functions.

Exception Handling Principles

Any error that occurs in a PL/pgSQL function aborts the execution of the function and related transactions. You can use a **BEGIN** block with an **EXCEPTION** clause to catch and fix errors.

- In a PL/pgSQL block, if an SQL statement cannot return a definite result, you
 are advised to handle exceptions (if any) in EXCEPTION. Otherwise,
 unhandled errors may be transferred to the external block and cause program
 logic errors.
- 2. You can directly use the exceptions that have been defined in the system. GaussDB(DWS) does not support custom exceptions.
- 3. A block containing an **EXCEPTION** clause is more expensive to enter and exit than a block without one. Therefore, do not use **EXCEPTION** without need.

Writing Standard

- 1. Variable naming rules:
 - a. The input parameter format of a procedure or function is IN_Parameter_name. The parameter name shall be in uppercase.
 - b. The output parameter format of a procedure or function is **OUT_***Parameter_name*. The parameter name shall be in uppercase.
 - c. The format for input and output parameters in a procedure or function is **IO_**Parameter name, with the parameter name written in uppercase.
 - d. Variables used in procedures and functions shall be composed of **v**_*Variable_name*. The variable name shall be in lower case.
 - e. In query concatenation, the concatenation variable name of the **WHERE** statement shall be **v_where**, and the concatenation variable name of the **SELECT** statement shall be **v select**.
 - f. The record type (TYPE) name shall consist of **T** and a variable name. The name shall be in uppercase.
 - g. A cursor name shall consist of **CUR** and a variable name. The name shall be in uppercase.
 - h. The name of a reference cursor (REF CURSOR) shall consist of **REF** and a variable name. The name shall be in uppercase.
- 2. Rules for defining variable types:
 - a. Use **%TYPE** to declare the type of a variable that has the same meaning as that of a column in an application table.
 - b. Use **%ROWTYPE** to declare the type of a record that has the same meaning as that of a row in an application table.
- 3. Rules for writing comments:
 - a. Comments shall be meaningful and shall not just repeat the code content
 - b. Comments shall be concise and easy to understand.
 - c. Comments shall be provided at the beginning of each stored procedure or function. The comments shall contain a brief function description, author,

compilation date, program version number, and program change history. The format of the comments at the beginning of stored procedures shall be the same.

- d. Comments shall be provided next to the input and output parameters to describe the meaning of variables.
- e. Comments shall be provided at the beginning of each block or large branch to briefly describe the function of the block. If an algorithm is used, comments shall be provided to describe the purpose and result of the algorithm.
- 4. Variable declaration format:

Each line shall contain only one statement. To assign initial values, write them in the same line.

5. Letter case:

Use uppercase letters except for variable names.

6. Indentation:

In the statements used for creating a stored procedure, the keywords **CREATE**, **AS/IS**, **BEGIN**, and **END** at the same level shall have the same indent.

- 7. Statement rules:
 - a. For statements that define variables, Each line shall contain only one statement.
 - b. The keywords **IF**, **ELSE IF**, **ELSE**, and **END** at the same level shall have the same indent.
 - c. The keywords **CASE** and **END** shall have the same indent. The keywords **WHEN** and **ELSE** shall be indented.
 - d. The keywords **LOOP** and **END LOOP** at the same level shall have the same indent. Nested statements or statements at lower levels shall have more indent.

Creating and Managing GaussDB(DWS) Database Objects

3.1 Creating and Managing GaussDB(DWS) Databases

A database is a collection of objects such as tables, indexes, views, stored procedures, and operators. GaussDB (DWS) supports the creation of multiple databases. However, a client program can connect to and access only one database at a time, and cross-database query is not supported.

Template and Default Databases

- GaussDB (DWS) provides two template databases template0 and template1 and a default database gaussdb.
- By default, each newly created database is based on a template database. The GaussDB(DWS) database uses template1 as the template by default. The encoding format is SQL_ASCII, and user-defined character encoding is not allowed. If you need to specify the character encoding when creating a database, use template0 to create the database.
- Do not use a client or any other tools to connect to or to perform operations on both the two template databases.

□ NOTE

You can run the **show server_encoding** command to view the current database encoding.

Creating a Database.

Run the **CREATE DATABASE** statement to create a database.

CREATE DATABASE mydatabase;

□ NOTE

- When you create a database, if the length of the database name exceeds 63 bytes, the
 server truncates the database name and retains the first 63 bytes. Therefore, you are
 advised to set the length of the database name to a value less than or equal to 63
 bytes. Do not use multi-byte characters as object names. If an object whose name is
 truncated mistakenly cannot be deleted, delete the object using the name before the
 truncation, or manually delete it from the corresponding system catalog on each node.
- Database names must comply with the naming convention of SQL identifiers. The current user automatically becomes the owner of this new database.
- If a database system is used to support independent users and projects, store them in different databases.
- If the projects or users are associated with each other and share resources, store them in different schemas in the same database.
- You must have the permission to create a database or the permission that the system administrator owns.

Viewing Databases

To view databases, perform the following steps:

- Run the \label{lambda} meta-command to view the database list of the database system.
- Querying the database list using the pg_database system catalog SELECT datname FROM pg_database;

Modifying a Database

You can use the **ALTER DATABASE** statement modify database configuration such as the database owner, name, and default settings.

- Run the following command to set the default search path for the database: ALTER DATABASE mydatabase SET search_path TO pa_catalog,public,
- Rename the database.
 ALTER DATABASE mydatabase RENAME TO newdatabase;

Deleting a Database

You can run **DROP DATABASE** statement to delete a database. This statement deletes the system catalog of the database and the database directory on the disk. Only the database owner or system administrator can delete a database. A database being accessed by users cannot be deleted, You need to connect to another database before deleting this database.

Run the **DROP DATABASE** statement to delete a database: **DROP DATABASE** *newdatabase*;

3.2 Creating and Managing GaussDB(DWS) Schemas

A schema is the logical organization of objects and data in a database. Schema management allows multiple users to use the same database without interfering with each other. Third-party applications can be added to corresponding schemas to avoid conflicts.

The same database object name can be used in different schemas in a database without causing conflicts. For example, both **a_schema** and **b_schema** can contain

a table named **mytable**. Users with required permissions can access objects across multiple schemas in a database.

If a user is created, a schema named after the user will also be created in the current database.

Public mode

Each database has a schema named **public**. All users have the ability to use the public schema in the database, but only certain roles have the authority to create objects within it.

Creating a Schema

Run the CREATE SCHEMA command to create a schema.
 CREATE SCHEMA myschema;

To create or access an object in the schema, the object name in the command should be composed of the schema name and the object name, which are separated by a dot (.), for example, **myschema.table**.

 Users can create a schema owned by others. For example, run the following command to create a schema named myschema and set the owner of the schema to user jack:

CREATE SCHEMA myschema AUTHORIZATION jack,

If **authorization username** is not specified, the schema owner is the user who runs the command.

Modifying a Schema

- Run the ALTER SCHEMA command to change the schema name. Only the schema owner can change the schema name.
 ALTER SCHEMA schema_name RENAME TO new_name;
- Run the ALTER SCHEMA command to change the schema owner.
 ALTER SCHEMA schema_name OWNER TO new_owner;

Setting the Schema Search Path

The GUC parameter **search_path** specifies the schema search sequence. The parameter value is a series of schema names separated by commas (,). If no schema is specified during object creation, the object will be added to the first schema displayed in the search path. If there are objects with the same name in different schemas and no schema is specified for an object query, the object will be returned from the first schema containing the object in the search path.

• Run the **SHOW** command to view the current search path.

```
SHOW SEARCH_PATH;
search_path
-----
"$user",public
(1 row)
```

The default value of **search_path** is **"\$user",public**. **\$\$user** indicates the name of the schema with the same name as the current session user. If the schema does not exist, **\$\$user** will be ignored. By default, after a user connects to a database that has schemas with the same name, objects will be added to all the schemas. If there are no such schemas, objects will be added to only to the **public** schema.

Run the SET command to modify the default schema of the current session.
 For example, if the search path is set to "myschema, public", myschema is searched first.

SET SEARCH_PATH TO myschema, public,

You can also run the **ALTER ROLE** command to set search_path for a role (user). For example:

ALTER ROLE jack SET search_path TO myschema, public;

Using a Schema

If you want to create or access an object in a specified schema, the object name must contain the schema name. To be specific, the name consists of a schema name and an object name, which are separated by a dot (.).

• Create a table **mytable** in **myschema**. Create a table in **schema_name.table_name** format.

CREATE TABLE myschema.mytable(id int, name varchar(20));

Query all data in the table mytable in myschema.

SELECT * FROM myschema.mytable; id | name ----+-----(0 rows)

Viewing a Schema

Use the current_schema() function to view the current schema.
 SELECT current_schema();

current_schema
----myschema
(1 row)

- To view the owner of a schema, perform the following join query on the system catalogs PG_NAMESPACE and PG_USER. Replace schema_name in the statement with the name of the schema to be queried.
 SELECT s.nspname,u.usename AS nspowner FROM PG_NAMESPACE s, PG_USER u WHERE nspname='schema_name' AND s.nspowner = u.usesysid;
- To view a list of all schemas, query the system catalog PG_NAMESPACE.
 SELECT * FROM PG NAMESPACE;
- Use the PGXC_TOTAL_SCHEMA_INFO view to query the space usage of schemas in the cluster.
 SELECT * FROM PGXC_TOTAL SCHEMA INFO;
- To view a list of tables in a schema, query the system catalog PG_TABLES. For example, the following query will return a table list from PG_CATALOG in the schema.

SELECT distinct(tablename), schemaname FROM PG_TABLES where schemaname = 'pg_catalog';

Schema Permission Control

By default, a user can only access database objects in its own schema. To access objects in other schemas, the user must have the **usage** permission of the corresponding schema.

By granting the **CREATE** permission for a schema to a user, the user can create objects in this schema.

Grant the usage permission of myschema to user jack.
 GRANT USAGE ON schema myschema TO jack;

 Run the following command to revoke the USAGE permission for myschema from iack:

REVOKE USAGE ON schema myschema FROM jack;

Drop Schema

• Run the **DROP SCHEMA** command to delete an empty schema (no database objects in the schema).

DROP SCHEMA IF EXISTS myschema;

 By default, a schema must be empty before being deleted. To delete a schema and all its objects (such as tables, data, and functions), use the CASCADE keyword.

DROP SCHEMA myschema CASCADE;

System Schema

- Each database has a pg_catalog schema, which contains system catalogs and all built-in data types, functions, and operators. pg_catalog is a part of the search path and has the second highest search priority. It is searched after the schema of temporary tables and before other schemas specified in search_path. This search order ensures that database built-in objects can be found. To use a custom object that has the same name as a built-in object, you can specify the schema of the custom object.
- The **information_schema** consists of a collection of views that contain object information in a database. These views obtain system information from the system catalogs in a standardized way.

3.3 Creating and Managing GaussDB(DWS) Tables

Creating a Table

You can run the **CREATE TABLE** command to create a table. When creating a table, you can define the following information:

- Columns and data type of the table.
- Table or column constraints that restrict a column or the data contained in a table. For details, see **Definition of Table Constraints**.
- Distribution policy of a table, which determines how the GaussDB (DWS)
 database divides data between segments. For details, see Definition of Table
 Distribution.
- Table storage format. For details, see Selecting a GaussDB(DWS) Table Storage Model.
- Partition table information. For details, see Creating and Managing GaussDB(DWS) Partitioned Tables.

Example: Use **CREATE TABLE** to create a table **web_returns_p1**, use **wr_item_sk** as the distribution key, and sets the range distribution function through **wr_returned_date_sk**.

```
CREATE TABLE web_returns_p1
(
    wr_returned_date_sk integer,
    wr_returned_time_sk integer,
```

```
wr_item_sk integer NOT NULL,
wr_refunded_customer_sk integer
)
WITH (orientation = column)
DISTRIBUTE BY HASH (wr_item_sk)
PARTITION BY RANGE(wr_returned_date_sk)
(
PARTITION p2019 START(20191231) END(20221231) EVERY(10000),
PARTITION p0 END(maxvalue)
).
```

Definition of Table Constraints

You can define constraints on columns and tables to restrict data in a table. However, there are the following restrictions:

- The primary key constraint and unique constraint in the table must contain a distribution column.
- Column-store tables support the PARTIAL CLUSTER KEY and table-level primary key and unique constraints, but do not support table-level foreign key constraints.
- Only the **NULL**, **NOT NULL**, and **DEFAULT** constant values can be used as column-store table column constraints.

Table 3-1 Table constraints

Constrain t	Description	Example
Check constraint	A CHECK constraint allows you to specify that values in a specific column must satisfy a Boolean (true) expression.	Create the products table. The price column must be positive. CREATE TABLE products (product_no integer, name text, price numeric CHECK (price > 0));
NOT NULL constraint	A NOT NULL constraint specifies that a column cannot have null values. A non-null constraint is always written as a column constraint.	Create the products table. The values of product_no and name cannot be null. CREATE TABLE products (product_no integer NOT NULL, name text NOT NULL, price numeric);

Constrain t	Description	Example
UNIQUE constraint	A UNIQUE constraint specifies that the values in a column or a group of columns are all unique. If DISTRIBUTE BY REPLICATION is not specified, the column table that contains only unique values must contain distribution columns.	Create the products table. The values of product_no must be unique. CREATE TABLE products (product_no integer UNIQUE, name text, price numeric)DISTRIBUTE BY HASH(product_no);
Primary key constraint	A primary key constraint is the combination of a UNIQUE constraint and a NOT NULL constraint. If DISTRIBUTE BY REPLICATION is not specified, the column set with a primary key constraint must contain distributed columns. If a table has a primary key, the column (or group of columns) of the primary key is selected as the distribution keys of the table by default.	Create the products table. The primary key constraint is product_no . CREATE TABLE products (product_no integer PRIMARY KEY, name text, price numeric)DISTRIBUTE BY HASH(product_no);
Partial cluster key	Partial cluster key can minimize or maximize sparse indexes to quickly filter base tables. Partial cluster key can specify multiple columns, but you are advised to specify no more than two columns.	Create the products table with PCK set to product_no : CREATE TABLE products (product_no integer, name text, price numeric, PARTIAL CLUSTER KEY(product_no)) WITH (ORIENTATION = COLUMN);

Definition of Table Distribution

GaussDB(DWS) supports the following distribution modes: replication, hash, and roundrobin.

■ NOTE

The roundrobin distribution mode is supported only by cluster version 8.1.2 or later.

Policy	Description	Scenario	Advantages/Disadvantages
Replicatio n	Full data in a table is stored on each DN in the cluster.	Small tables and dimension tables	 The advantage of replication is that each DN has full data of the table. During the join operation, data does not need to be redistributed, reducing network overheads and reducing plan segments (each plan segment starts a corresponding thread). The disadvantage of replication is that each DN retains the complete data of the table, resulting in data redundancy. Generally, replication is only used for small dimension tables.
Hash	Table data is distributed on all DNs in the cluster.	Fact tables containing a large amount of data	 The I/O resources of each node can be used during data read/write, greatly improving the read/write speed of a table. Generally, a large table (containing over 1 million records) is defined as a hash table.
Polling (Round- robin)	Each row in the table is sent to each DN in turn. Data can be evenly distributed on each DN.	Fact tables that contain a large amount of data and cannot find a proper distribution column in hash mode	 Round-robin can avoid data skew, improving the space utilization of the cluster. Round-robin does not support local DN optimization like a hash table does, and the query performance of Round-robin is usually lower than that of a hash table. If a proper distribution column can be found for a large table, use the hash distribution mode with better performance. Otherwise, define the table as a round-robin table.

Selecting a Distribution Key

If the hash distribution mode is used, a distribution key must be specified for the user table. When a record is inserted, the system hashes it based on the distribution key and then stores it on the corresponding DN.

Select a hash distribution key based on the following principles:

- 1. The values of the distribution key should be discrete so that data can be evenly distributed on each DN. You can select the primary key of the table as the distribution key. For example, for a person information table, choose the ID number column as the distribution key.
- 2. **Do not select the column that has a constant filter.** For example, if a constant constraint (for example, zqdh= '000001') exists on the **zqdh** column in some queries on the **dwcjk** table, you are not advised to use **zqdh** as the distribution key.
- 3. With the above principles met, you can select join conditions as distribution keys, so that join tasks can be pushed down to DNs for execution, reducing the amount of data transferred between the DNs.

For a hash table, an inappropriate distribution key may cause data skew or poor I/O performance on certain DNs. Therefore, you need to check the table to ensure that data is evenly distributed on each DN. You can run the following SQL statements to check for data skew:

select
xc_node_id, count(1)
from tablename
group by xc_node_id
order by xc node id desc;

xc_node_id corresponds to a DN. Generally, over 5% difference between the amount of data on different DNs is regarded as data skew. If the difference is over 10%, choose another distribution key.

4. You are not advised to add a column as a distribution key, especially add a new column and use the SEQUENCE value to fill the column. (Sequences may cause performance bottlenecks and unnecessary maintenance costs.)

View the data in the table.

- Run the following command to query information about all tables in a database in the system catalog pg_tables:
 SELECT * FROM pg_tables;
- Run the \d+ command of the gsql tool to query table attributes: \d+ customer t1;
- Run the following command to query the data volume of table customer_t1:
 SELECT count(*) FROM customer_t1;
- Run the following command to query all data in table **customer_t1**: **SELECT** * **FROM** *customer t1*;
- Run the following command to query data in column c_customer_sk:
 SELECT c_customer_sk FROM customer_t1;
- Run the following command to filter repeated data in column **c_customer_sk**: **SELECT DISTINCT**(*c_customer_sk*) **FROM** *customer_t1*;
- Run the following command to query all data whose column **c_customer_sk** is **3869**:
 - **SELECT * FROM** *customer_t1* **WHERE** *c_customer_sk = 3869*,
- Run the following command to sort data based on column c_customer_sk.
 SELECT * FROM customer_t1 ORDER BY c_customer_sk;

Deleting Data in a Table

⚠ CAUTION

Exercise caution when running the **DROP TABLE** and **TRUNCATE TABLE** statements. After a table is deleted, data cannot be restored.

- Delete the **customer_t1** table from the database. **DROP TABLE** *customer_t1*;
- You can use **DELETE** or **TRUNCATE** to clear rows in a table without removing the definition of the table.

Delete all rows from the **customer_t1** table.

TRUNCATE TABLE customer t1;

Delete all rows from the customer_t1 table.

DELETE FROM *customer_t1*;

Delete all records whose **c_customer_sk** is **3869** from the **customer_t1** table. **DELETE FROM** *customer_t1* **WHERE** *c_customer_sk* = *3869*,

Managing UNLOGGED Tables

UNLOGGED indicates an unlogged table. Unlogged tables are faster than regular tables because data written to them is not written to the WALs. However, an unlogged table is automatically cleared after a crash or unclean shutdown, incurring data loss risks. The contents of an unlogged table are also not replicated to standby servers. Any indexes created on an unlogged table are not automatically logged as well.

Usage scenario: Unlogged tables do not ensure safe data. Users can back up data before using unlogged tables; for example, users should back up the data before a system upgrade. When creating an unlogged table, disable cnretry (that is, set the GUC parameter **max_query_retry_times** to **0**).

Troubleshooting: If data is missing in the indexes of unlogged tables due to some unexpected operations such as an unclean shutdown, users should re-create the indexes with errors.

- Starting from version 9.1.0, UNLOGGED tables are automatically saved in the pg_unlogged tablespace and cannot be moved or assigned to other tablespaces.
- After an earlier version is upgraded to 9.1.0, the UNLOGGED table created in the earlier version is still stored in the original tablespace.

Version 9.1.0 has a script called **switch_unlogged_tablespace.py** that can move unlogged tables to optimize the recovery time objective (RTO). This script works together with the GUC parameter **enable_unlogged_tablespace_compat**.

1. The script is stored in the **\$GPHOME/script** directory. You can use the **-?** command to obtain help information.

- Migrate all unlogged tables (recommended). python3 switch_unlogged_tablespace.py -t switch
- 3. After the migration, the GUC parameter **enable_unlogged_tablespace_compat** is automatically set to **off**.

NOTICE

After the upgrade to 9.1.0, you are advised to perform the following two steps to improve the instance restart RTO:

- 1. Use the **switch_unlogged_tablespace.py** script to migrate all unlogged tables to the **pg_unlogged** tablespace.
- 2. If the old version does not use any unlogged table, you are advised to set the GUC parameter **enable_unlogged_tablespace_compat** to **OFF**.

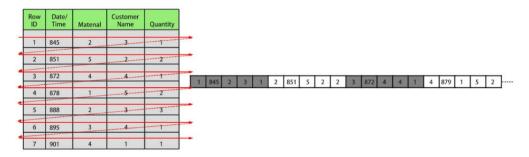
3.4 Selecting a GaussDB(DWS) Table Storage Model

GaussDB(DWS) supports hybrid row and column storage. When creating a table, you can set the table storage mode to row storage or column storage.

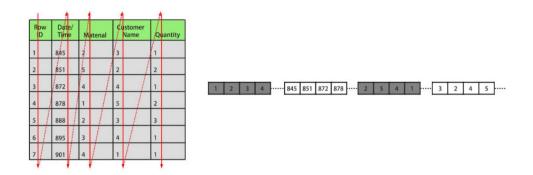
Row storage stores tables to disk partitions by row, and column storage stores tables to disk partitions by column. By default, a table is created in row storage mode. For details about differences between row storage and column storage, see Figure 3-1.

Figure 3-1 Differences between row storage and column storage

Row-based store



Column-based store



In the preceding figure, the upper left part is a row-store table, and the upper right part shows how the row-store table is stored on a disk; the lower left part is a column-store table, and the lower right part shows how the column-store table is stored on a disk.

The row/column storage of a table is specified by the **orientation** attribute in the table definition. The value **row** indicates a row-store table and **column** indicates a column-store table. The default value is **row**. Each storage mode applies to specific scenarios. Select an appropriate mode when creating a table.

Storage Mode	Benefit	Drawback	Application Scenarios
Row storage	Data is stored by row. When you query a row of data, you can quickly locate the	All data in the queried row is read while only a few columns are needed.	1. The number of columns in the table is small, and most fields in the table are queried.
	target row.		2. Point queries (simple index-based query that returns only a few records) are performed.
			3. Add, Delete, Modify, and Query operations on entire rows are frequently performed.
Column storage	Only necessary columns in a query are read. The	It is not suitable for INSERT or UPDATE operations on a	Query a few columns in a table that contains a large number of columns.
homogeneity of data within a column facilitates efficient	small amount of data.	2. Statistical analysis queries (requiring a large number of association and grouping operations)	
	compression.		3. Ad hoc queries (using uncertain query conditions and unable to utilize indexes to scan row-store tables)

Table 3-2 Table storage modes and scenarios

Creating a Row-store Table

For example, to create a row-store table named **customer_t1**, run the following command:

```
CREATE TABLE customer_t1
(
    state_ID CHAR(2),
    state_NAME VARCHAR2(40),
    area_ID NUMBER
);
```

Creating a column-store table.

For example, to create a column-store table named **customer_t2**, run the following command:

```
CREATE TABLE customer_t2
(
state_ID CHAR(2),
```

```
state_NAME VARCHAR2(40),
area_ID NUMBER
)
WITH (ORIENTATION = COLUMN);
```

Table Compression

Table compression can be enabled when a table is created. Table compression enables data in the table to be stored in compressed format to reduce memory usage.

In scenarios where I/O is large (much data is read and written) and CPU is sufficient (little data is computed), select a high compression ratio. In scenarios where I/O is small and CPU is insufficient, select a low compression ratio. Based on this principle, you are advised to select different compression ratios and test and compare the results to select the optimal compression ratio as required. Specify a compressions ratio using the **COMPRESSION** parameter. The supported values are as follows:

- The valid value of column-store tables is **YES**, **NO**, **LOW**, **MIDDLE**, or **HIGH**, and the default value is **LOW**.
- The valid values of row-store tables are YES and NO, and the default is NO.
 (The row-store table compression function is not put into commercial use. To use this function, contact technical support.)

The service scenarios applicable to each compression level are described in the following table.

Compression Level	Application Scenario
LOW	The system CPU usage is high and the disk storage space is sufficient.
MIDDLE	The system CPU usage is moderate and the disk storage space is insufficient.
HIGH	The system CPU usage is low and the disk storage space is insufficient.

For example, to create a compressed column-store table named **customer_t3**, run the following command:

```
CREATE TABLE customer_t3
(

state_ID CHAR(2),
state_NAME VARCHAR2(40),
area_ID NUMBER
)
WITH (ORIENTATION = COLUMN,COMPRESSION=middle);
```

3.5 Creating and Managing GaussDB(DWS) Partitioned Tables

Partitioning refers to splitting what is logically one large table into smaller physical pieces based on specific schemes. The table based on the logic is called a partition cable, and a physical piece is called a partition. Data is stored on these smaller physical pieces, namely, partitions, instead of the larger logical partitioned table. During conditional query, the system scans only the partitions that meet the conditions rather than scanning the entire table improving query performance.

Advantages of partitioned tables:

- Improved query performance. You can search in specific partitions, improving the search efficiency.
- Enhanced availability. If a partition is faulty, data in other partitions is still available.
- Improved maintainability. For expired historical data that needs to be periodically deleted, you can quickly delete it by dropping or truncate partitions.

Supported Table Partition Types

- Range partitioning: partitions are created based on a numeric range, for example, by date or price range.
- List partitioning: partitions are created based on a list of values, such as sales scope or product attribute. Only clusters of 8.1.3 and later versions support this function.

Choosing to Partition a Table

You can choose to partition a table when the table has the following characteristics:

- There are obvious ranges among the fields of the table.
 A table is partitioned based on obvious rangeable fields. Generally, columns such as date, area, and value are used for partitioning. The time column is most commonly used.
- Queries to the table have obvious range characteristics.

 If the queried data fall into specific ranges, its better tables are partitioned so that through partition pruning, only the queried partition needs to be scanned, improving data scanning efficiency and reducing the I/O overhead of data scanning.
- The table contains a large amount of data. Scanning small tables does not take much time, therefore the performance benefits of partitioning are not significant. Therefore, you are advised to partition only large tables. In column-store table, each column is an independent file storage unit, and the minimum storage unit CU can store 60,000 rows of data. Therefore, for column-store partitioned tables, it is recommended that the data volume in each partition be greater than or equal to the number of DNs multiplied by 60,000.

Creating a Range Partitioned Table

Example: Create a table **web_returns_p1** partitioned by the range **wr returned date sk**.

```
CREATE TABLE web_returns_p1
  wr_returned_date_sk
                         integer,
  wr_returned_time_sk
                        integer,
                     integer NOT NULL,
  wr_item_sk
  wr_refunded_customer_sk integer
WITH (orientation = column)
DISTRIBUTE BY HASH (wr_item_sk)
PARTITION BY RANGE (wr_returned_date_sk)
  PARTITION p2016 VALUES LESS THAN (20161231),
  PARTITION p2017 VALUES LESS THAN (20171231),
  PARTITION p2018 VALUES LESS THAN(20181231),
  PARTITION p2019 VALUES LESS THAN(20191231),
  PARTITION pxxxx VALUES LESS THAN(maxvalue)
```

Create partitions in batches, with fixed partition ranges. The following example can be used:

Partition the table **web_returns_p2** by date and time, using time as the partition kev.

```
CRÉATE TABLE web_returns_p2
(
    id integer,
    idle numeric,
    IO numeric,
    scope text,
    IP text,
    time timestamp
)
WITH (TTL='7 days',PERIOD='1 day')
PARTITION BY RANGE(time)
(
    PARTITION P1 VALUES LESS THAN('2022-01-05 16:32:45'),
    PARTITION P2 VALUES LESS THAN('2022-01-06 16:56:12')
);
```

Creating a List Partitioned Table

A list partitioned table can use any column that allows value comparison as the partition key column. When creating a list partitioned table, you must declare the value partition for each partition.

Example: Create a list partitioned table sales_info.

```
CREATE TABLE sales_info
(
```

```
sale_time timestamptz,
period int,
city text,
price numeric(10,2),
remark varchar2(100)
)
DISTRIBUTE BY HASH(sale_time)
PARTITION BY LIST (period, city)
(
PARTITION province1_202201 VALUES (('202201', 'city1'), ('202201', 'city2')),
PARTITION province2_202201 VALUES (('202201', 'city3'), ('202201', 'city4'), ('202201', 'city5')),
PARTITION rest VALUES (DEFAULT)
);
```

Partitioning an Existing Table

A table can be partitioned only when it is created. If you want to partition a table, you must create a partitioned table, load the data in the original table to the partitioned table, delete the original table, and rename the partitioned table as the name of the original table. You must also re-grant permissions on the table to users. For example:

```
CREATE TABLE web_returns_p2
  wr_returned_date_sk
                         integer,
                         integer,
  wr_returned_time_sk
                     integer NOT NULL,
  wr_item_sk
  wr_refunded_customer_sk integer
WITH (orientation = column)
DISTRIBUTE BY HASH (wr_item_sk)
PARTITION BY RANGE(wr_returned_date_sk)
   PARTITION p2016 START(20161231) END(20191231) EVERY(10000),
  PARTITION p0 END(maxvalue)
INSERT INTO web_returns_p2 SELECT * FROM web_returns_p1;
DROP TABLE web_returns_p1;
ALTER TABLE web_returns_p2 RENAME TO web_returns_p1;
GRANT ALL PRIVILEGES ON web_returns_p1 TO dbadmin;
GRANT SELECT ON web_returns_p1 TO jack;
```

Adding a Partition

Run the **ALTER TABLE** statement to add a partition to a partitioned table. For example, to add partition **P2020** to the **web_returns_p1** table, run the following command:

ALTER TABLE web_returns_p1 ADD PARTITION P2020 VALUES LESS THAN (20201231);

Splitting a Partition

The syntax for splitting a partition varies between a range partitioned table and a list partitioned table.

- Run the ALTER TABLE statement to split a partition in a range partitioned table. For example, the partition pxxxx of the table web_returns_p1 is split into two partitions p2020 and p20xx at the splitting point 20201231.
 ALTER TABLE web_returns_p1 SPLIT PARTITION pxxxx AT(20201231) INTO (PARTITION p2020,PARTITION p20xx);
- Run the **ALTER TABLE** statement to split a partition in a list partitioned table. For example, split the partition **province2_202201** of table **sales_inf** into two partitions **province3_202201** and **province4_202201**.

ALTER TABLE sales_info SPLIT PARTITION province2_202201 VALUES(('202201', 'city5')) INTO (PARTITION province3_202201,PARTITION province4_202201);

Merging Partitions

Run the **ALTER TABLE** statement to merge two partitions in a partitioned table. For example, merge partitions **p2016** and **p2017** of table **web_returns_p1** into one partition **p20162017**.

ALTER TABLE web_returns_p1 MERGE PARTITIONS p2016,p2017 INTO PARTITION p20162017;

Deleting a Partition

Run the **ALTER TABLE** statement to delete a partition from a partitioned table. For example, run the following command to delete partition **P2020** from the **web returns p1** table:

ALTER TABLE web_returns_p1 DROP PARTITION P2020,

Querying a Partition

- Query partition p2019.
 SELECT * FROM web_returns_p1 PARTITION (p2019);
 SELECT * FROM web_returns_p1 PARTITION FOR (20201231);
- View partitioned tables using the system catalog dba_tab_partitions.
 SELECT * FROM dba_tab_partitions where table_name='web_returns_p1';

Deleting a Partitioned Table

Run the **DROP TABLE** statement to delete a partitioned table.

DROP TABLE web returns p1;

3.6 Creating and Managing GaussDB(DWS) Indexes

Indexes accelerate the data access speed but also add the processing time of the insert, update, and delete operations. Therefore, before creating an index, consider whether it is necessary and determine the columns where indexes will be created. You can determine whether to add an index for a table by analyzing the service processing and data use of applications, as well as columns that are frequently used as search criteria or need to be sorted.

Index type

- **B-tree**: The B-tree index uses a structure that is similar to the B+ tree structure to store data key values, facilitating index search. B-tree indexes support comparison and range queries.
- **GIN**: GIN indexes are reverse indexes and can process values that contain multiple keys (for example, arrays).
- **GiST**: GiST indexes are suitable for the set data type and multidimensional data types, such as geometric and geographic data types.
- PSORT: PSORT indexes are used to perform partial sort on column-store tables.

Row-based tables support B-tree (default), GIN, and GiST indexes. Column-based tables support PSORT (default), B-Tree, and GIN indexes.

□ NOTE

Create a B-tree index for point queries.

Index Selection Principles

Indexes are created based on columns in database tables. When creating indexes, you need to determine the columns, which can be:

- Columns that are frequently searched: The search efficiency can be improved.
- The uniqueness of the columns and the data sequence structures is ensured.
- Columns that usually function as foreign keys and are used for connections. Then the connection efficiency is improved.
- Columns that are usually searched for by a specified scope. These indexes have already been arranged in a sequence, and the specified scope is contiguous.
- Columns that need to be arranged in a sequence. These indexes have already been arranged in a sequence, so the sequence query time is accelerated.
- Columns that usually use the WHERE clause. Then the condition decision efficiency is increased.
- Fields that are frequently used after keywords, such as ORDER BY, GROUP BY, and DISTINCT.

■ NOTE

- After an index is created, the system automatically determines when to reference
 it. If the system determines that indexing is faster than sequenced scanning, the
 index will be used.
- After an index is successfully created, it must be synchronized with the associated table to ensure new data can be accurately located. Therefore, data operations increase. Therefore, delete unnecessary indexes periodically.

Creating an Index

GaussDB(DWS) supports four methods for creating indexes. For details, see **Table 3-3**.

- After an index is created, the system automatically determines when to reference it. If the system determines that indexing is faster than sequenced scanning, the index will be used
- After an index is successfully created, it must be synchronized with the associated table to ensure new data can be accurately located. Therefore, data operations increase. Therefore, delete unnecessary indexes periodically.

Table 3-3 Indexing Method

Indexing Method	Description
Unique index	Refers to an index that constrains the uniqueness of an index attribute or an attribute group. If a table declares unique constraints or primary keys, GaussDB(DWS) automatically creates unique indexes (or composite indexes) for columns that form the primary keys or unique constraints. Currently, only B-tree can create a unique index in GaussDB(DWS).
Composite index	Refers to an index that can be defined for multiple attributes of a table. Currently, composite indexes can be created only for B-tree in GaussDB(DWS) and a maximum of 32 columns can share a composite index.
Partial index	Refers to an index that can be created for subsets of a table. This indexing method contains only tuples that meet condition expressions.
Expression index	Refers to an index that is built on a function or an expression calculated based on one or more attributes of a table. An expression index works only when the queried expression is the same as the created expression.

- Run the following command to create an ordinary table: CREATE TABLE tpcds.customer_address_bak AS TABLE tpcds.customer_address,
- Create a common index.

You need to query the following information in the **tpcds.customer_address_bak** table:

SELECT ca_address_sk FROM tpcds.customer_address_bak WHERE ca_address_sk=14888,

Generally, the database system needs to scan the

tpcds.customer_address_bak table row by row to find all matched tuples. If the size of the **tpcds.customer_address_bak** table is large but only a few (possibly zero or one) of the WHERE conditions are met, the performance of this sequential scan is low. If the database system maintains an index on the **ca_address_sk** attribute, it can quickly locate the matching tuple by searching only a few layers in the search tree, significantly enhancing data query performance. Furthermore, indexes can improve the update and delete operation performance in the database.

Run the following command to create an index:

CREATE INDEX index_wr_returned_date_sk ON tpcds.customer_address_bak (ca_address_sk);

• Create a unique index.

If a table declares a unique constraint or primary key, GaussDB(DWS) automatically creates a unique index (possibly a multi-column index) on the columns that form the primary key or unique constraint. If no unique constraint or primary key is specified during table creation, you can run the CREATE INDEX statement to create an index.

CREATE UNIQUE INDEX unique_index ON tpcds.customer_address_bak(ca_address_sk);

Create a multi-column index.

Assume you need to frequently query records with **ca_address_sk** being **5050** and **ca_street_number** smaller than **1000** in the

tpcds.customer_address_bak table. Run the following command:

SELECT ca_address_sk,ca_address_id FROM tpcds.customer_address_bak WHERE ca_address_sk = 5050 AND ca_street_number < 1000;

Run the following command to define a multiple-column index on ca address sk and ca street number columns:

CREATE INDEX more_column_index ON

tpcds.customer_address_bak(ca_address_sk ,ca_street_number);

• Create a partition index.

If you only want to find records whose **ca_address_sk** is **5050**, you can create a partial index to facilitate your query.

CREATE INDEX part_index **ON** tpcds.customer_address_bak(ca_address_sk) **WHERE** ca_address_sk = 5050,

Create an expression index.

Assume you need to frequently query records with **ca_street_number** smaller than **1000**, run the following command:

SELECT * FROM tpcds.customer_address_bak WHERE trunc(ca_street_number) < 1000;

The following expression index can be created for this query task: CREATE INDEX para_index ON tpcds.customer_address_bak (trunc(ca_street_number));

Querying an Index

 Run the following command to query all indexes defined by the system and users:

SELECT RELNAME FROM PG_CLASS WHERE RELKIND='i';

• Run the following command to query information about a specified index: \di+ index_wr_returned_date_sk

Recreating an Index

- Recreate the index index_wr_returned_date_sk.
 REINDEX INDEX index_wr_returned_date_sk;
- Recreate all indexes of a table.
 REINDEX TABLE tpcds.customer_address_bak;

Deleting an Index

You can use the **DROP INDEX** statement to delete indexes. **DROP INDEX** *index_wr_returned_date_sk*;

3.7 Creating and Using GaussDB(DWS) Sequences

A sequence is a database object that generates unique integers according to a certain rule and is usually used to generate primary key values.

You can create a sequence for a column in either of the following methods:

- Set the data type of a column to sequence integer. A sequence will be automatically created by the database for this column.
- Use CREATE SEQUENCE to create a new sequenc. Use the nextval('sequence_name') function to increment the sequence and return a

new value. Specify the default value of the column as the sequence value returned by the **nextval(**'sequence_name') function. In this way, this column can be used as a unique identifier.

Creating a Sequence.

Method 1: Set the data type of a column to a sequence integer. For example:

CREATE TABLE 71

(
 id serial,
 name text
):

Method 2: Create a sequence and set the initial value of the **nextval**('sequence_name') function to the default value of a column. You can cache a specific number of sequence values to reduce the requests to the GTM, improving the performance.

 Create a sequence. CREATE SEQUENCE seq1 cache 100;

2. Set the initial value of the **nextval**('sequence_name') function to the default value of a column.

```
CREATE TABLE 72
(
id int not null default nextval('seq1'),
name text
);
```

◯ NOTE

Methods 1 and 2 are similar except that method 2 specifies cache for the sequence. A sequence using cache has holes (non-consecutive values, for example, 1, 4, 5) and cannot keep the order of the values. After a sequence is deleted, its sub-sequences will be deleted automatically. A sequence shared by multiple columns is not forbidden in a database, but you are not advised to do that.

Currently, the preceding two methods cannot be used for existing tables.

Modifying a Sequence

The **ALTER SEQUENCE** statement changes the attributes of an existing sequence, including the owner, owning column, and maximum value.

Associate the sequence with a column.

The sequence will be deleted when you delete the column or the table where the column resides.

ALTER SEQUENCE seq1 OWNED BY T2.id,

Modify the maximum value of serial to 300.
 ALTER SEQUENCE seq1 MAXVALUE 300;

Deleting a Sequence

Run the **DROP SEQUENCE** command to delete a sequence. For example, to delete the sequence named **seq1**, run the following command:

```
DROP SEQUENCE seq1;
```

Precautions

Sequence values are generated by the GTM. By default, each request for a sequence value is sent to the GTM. The GTM calculates the result of the current value plus the step and then returns the result. As GTM is a globally unique node, generating default sequence numbers can cause performance issues. For operations that need frequent sequence number generation, such as bulkload data import, this is not recommended. For example, the **INSERT FROM SELECT** statement has poor performance in the following scenario:

```
CREATE SEQUENCE newSeq1;
CREATE TABLE newT1

(
    id int not null default nextval('newSeq1'),
    name text
);
INSERT INTO newT1(name) SELECT name from T1;
```

To improve the performance, run the following statements (assume that data of 10,000 rows will be imported from *T1* to *newT1*):

```
INSERT INTO newT1(id, name) SELECT id,name from T1;
SELECT SETVAL('newSeq1',10000);
```


Rollback is not supported by sequence functions, including **nextval()** and **setval()**. The value of the setval function immediately takes effect on nextval in the current session in any cases and take effect in other sessions only when no cache is specified for them. If cache is specified for a session, it takes effect only after all the cached values have been used. To avoid duplicate values, use setval only when necessary. Do not set it to an existing sequence value or a cached sequence value.

If BulkLoad is used, set sufficient cache for <code>newSeq1</code> and do not set <code>Maxvalue</code> or <code>Minvalue</code>. To improve the performance, database may push down the invocation of <code>nextval('sequence_name')</code> to DNs. Currently, the concurrent connection requests that can be processed by the GTM are limited. If there are too many DNs, a large number of concurrent connection requests will be sent to the GTM. In this case, you need to limit the concurrent connection of BulkLoad to save the GTM connection resources. If the target table is a replication table (<code>DISTRIBUTE BY REPLICATION</code>), pushdown cannot be performed. If the data volume is large, this will be a disaster for the database. In addition, the database space may be exhausted. After the import is complete, do <code>VACUUM FULL</code>. Therefore, you are not advised to use sequences when <code>BulkLoad</code> is used.

After a sequence is created, a single-row table is maintained on each node to store the sequence definition and value, which is obtained from the last interaction with the GTM rather than updated in real time. The single-row table on a node does not update when other nodes request a new value from the GTM or when the sequence is modified using **setval**.

3.8 Creating and Managing GaussDB(DWS) Views

Views allow users to save queries. Views are not physically stored on disks. Queries to a view run as subqueries. A database only stores the definition of a view and does not store its data. The data is still stored in the original base table. If data in the base table changes, the data in the view changes accordingly. In this sense, a

view is like a window through which users can know their interested data and data changes in the database. A view is triggered every time it is referenced.

Creating a view

Run the **CREATE VIEW** command to create a view. **CREATE OR REPLACE VIEW** MyView **AS SELECT * FROM** tpcds.customer WHERE c_customer_sk < 150;

□ NOTE

The **OR REPLACE** parameter in this command is optional. It indicates that if the view exists, the new view will replace the existing view.

View Details

- View the MyView view. Real-time data will be returned.
 SELECT * FROM myview;
- Run the following command to query the views in the current user:
 SELECT * FROM user_views;
- Run the following command to query all views:
 SELECT * FROM dba views;
- View details about a specified view.

Run the following command to view details about the dba_users view:

SELECT PG_AUTHID.ROLNAME::CHARACTER VARYING(64) AS USERNAME FROM PG_AUTHID;

Rebuilding a View

Run the **ALTER VIEW** command to rebuild a view without entering query statements.

ALTER VIEW myview REBUILD,

Deleting a View

Run the **DROP VIEW** command to delete a view. **DROP VIEW** *myview*;

DROP VIEW ... The **CASCADE** command can be used to delete objects that depend on the view. For example, view A depends on view B. If view B is deleted, view A will also be deleted. Without the CASCADE option, the **DROP VIEW** command will fail.

3.9 Creating and Managing GaussDB(DWS) Scheduled Tasks

GaussDB(DWS) allows users to create scheduled tasks, which are automatically executed at specified time points, reducing O&M workload.

The database follows Oracle's scheduled task feature, allowing you to use the DBMS.JOB advanced package interfaces to create, run, delete, and modify scheduled tasks (including the task ID, task enabling and disabling, task triggering time, triggering interval, and task content).

- The hybrid data warehouse (standalone) does not support scheduled tasks.
- The execution statements of scheduled tasks are not recorded in the Real-time Top SQL logs. The statements can be recorded only in versions later than 8.2.1.
- By default, GaussDB(DWS) uses the UTC time. The execution time of the scheduled task needs to be converted to the time zone of the user.

Periodic Task Management

Step 1 Creates a test table.

CREATE TABLE test(id int, time date);

If the following information is displayed, the table has been created.

CREATE TABLE

Step 2 Create the customized storage procedure.

```
CREATE OR REPLACE PROCEDURE PRC_JOB_1()
AS
N_NUM integer :=1;
BEGIN
FOR I IN 1..1000 LOOP
INSERT INTO test VALUES(I,SYSDATE);
END LOOP;
END;
/
```

If the following information is displayed, the procedure has been created.

CREATE PROCEDURE

Step 3 Create a task.

 Create a task with unspecified job_id and execute the PRC_JOB_1 storage procedure every two minutes.

```
call dbms_job.submit('call public.prc_job_1(); ', sysdate, 'interval "1 minute"', :a);
job
-----
1
(1 row)
```

• Create task with specified job_id.

```
call dbms_job.isubmit(2,'call public.prc_job_1(); ', sysdate, 'interval "1 minute"'); isubmit
-------
(1 row)
```

Step 4 View the created task information about the current user in the **USER_JOBS** view.

Only the system administrator can access this system view. For details about the fields, see **Table 14-341**.

```
postgresselect job,dbname,start_date,last_date,this_date,next_date,broken,status,interval,failures,what from user_jobs;
job | dbname | start_date | last_date | this_date | next_date | broken | status | interval | failures | what
```

Step 5 Stop a task.

```
call dbms_job.broken(1,true);
broken
------
(1 row)
```

Step 6 Start a task.

```
call dbms_job.broken(1,false);
broken
------
(1 row)
```

Step 7 Modify attributes of a task.

 Modify the Next_date parameter information about a task. For example, change the value of Next_date of Job1 to 1 hour.

```
call dbms_job.next_date(1, sysdate+1.0/24);
next_date
-----------
(1 row)
```

 Modify the Interval parameter information of a task. For example, change the value of Interval of Job1 to 1 hour.

```
call dbms_job.interval(1,'sysdate + 1.0/24');
interval
-------
(1 row)
```

 Modify the What parameter information of a JOB. For example, change What of Job1 to insert into public.test values(333, sysdate+5).

```
call dbms_job.what(1,'insert into public.test values(333, sysdate+5);');
what
------
(1 row)
```

Modify Next_date, Interval, and What parameter information of JOB.

```
call dbms_job.change(1, 'call public.prc_job_1();', sysdate, 'interval "1 minute""); change
------
(1 row)
```

Step 8 Delete a job.

```
call dbms_job.remove(1);
remove
------
(1 row)
```

Step 9 Set job permissions.

- During the creation of a job, the job is bound to the user and database that created the job. Accordingly, the user and database are added to **dbname** and **log_user** columns in the **pg_job** system view, respectively.
- If the current user is a DBA user, system administrator, or the user who created the job (**log_user** in **pg_job**), the user has the permissions to delete or modify parameter settings of the job using the remove, change, next_data, what, or interval interface. Otherwise, the system displays a message indicating that the current user has no permission to perform operations on the JOB.
- If the current database is the one that created a job, (that is, **dbname** in **pg_job**), you can delete or modify parameter settings of the job using the remove, change, next_data, what, or interval interface.
- When deleting the database that created a job, (that is, dbname in pg_job), the system associatively deletes the job records of the database.
- When deleting the user who created a job, (that is, **log_user** in **pg_job**), the system associatively deletes the job records of the user.

----End

3.10 Viewing GaussDB(DWS) System Catalogs

In addition to the created tables, a database contains many system catalogs These system catalogs contain cluster installation information and information about various queries and processes in GaussDB(DWS). You can collect information about the database by querying the system catalog.

Querying Database Tables

For example, query the **PG_TABLES** system catalog for all tables in the **public** schema.

SELECT distinct(tablename) FROM pg_tables WHERE SCHEMANAME = 'public';

Information similar to the following is displayed:

```
tablename
------
err_hr_staffs
test
err_hr_staffs_ft3
web_returns_p1
mig_seq_table
films4
(6 rows)
```

Viewing Database Users

You can run the **PG_USER** command to view the list of all users in the database, and view the user ID (**USESYSID**) and permissions.

Ruby 10 t	t t	t *******	default_pool 0
dbadmin 16393 f	 f f	f	default_pool 0
	 f f	f	default_pool 0
	 f f	f	default_pool 0
 (4 rows)			

GaussDB(DWS) uses Ruby to perform routine management and maintenance. You can add **WHERE usesysid > 10** to the **SELECT** statement to filter queries so that only specified user names are displayed.

Viewing and Stopping the Running Query Statements

You can view the running query statements in the **PG_STAT_ALL_INDEXES** view. Do as follows:

Step 1 Set the parameter **track_activities** to **on**.

SET track activities = on;

The database collects the running information about active queries only if the parameter is set to **on**.

Step 2 View the running query statements. Run the following command to view the database names, users, query statuses, and PIDs of the running query statements: SELECT datname, usename, state,pid FROM pg_stat_activity;

If the **state** column is **idle**, the connection is idle and requires a user to enter a command.

To identify only active query statements, run the following command:

SELECT datname, usename, state FROM pg_stat_activity WHERE state != 'idle';

Step 3 To cancel queries that have been running for a long time, use the **PG_TERMINATE_BACKEND** function to end sessions based on the thread ID. SELECT PG_TERMINATE_BACKEND(139834759993104);

If information similar to the following is displayed, the session is successfully terminated:

```
PG_TERMINATE_BACKEND
------
t
(1 row)
```

If information similar to the following is displayed, a user has terminated the current session.

FATAL: terminating connection due to administrator command FATAL: terminating connection due to administrator command

■ NOTE

If the **PG_TERMINATE_BACKEND** function is used to terminate the backend threads of the current session, the gsql client will be reconnected automatically rather than be logged out. The message "The connection to the server was lost." is returned. Attempting reset: Succeeded."

FATAL: terminating connection due to administrator command FATAL: terminating connection due to administrator command The connection to the server was lost. Attempting reset: Succeeded.

----End

4 Syntax Compatibility Differences Among Oracle, Teradata, and MySQL

GaussDB(DWS) supports Oracle, Teradata, and MySQL compatibility types, each matching their respective syntax. Syntax behavior changes based on the selected type.

The database compatibility type can be specified during database creation (using the **DBCOMPATIBILITY** parameter). The following is an example of the syntax. For details, see the **CREATE DATABASE** syntax.

CREATE DATABASE td_compatible_db DBCOMPATIBILITY 'TD'; --Create a Teradata-compatible database. CREATE DATABASE ora_compatible_db DBCOMPATIBILITY 'ORA'; --Create an Oracle-compatible database. CREATE DATABASE mysql_compatible_db DBCOMPATIBILITY 'MYSQL'; --Create a MySQL-compatible database.

Query results:

SELECT datname,datcompatibility FROM PG_DATABASE WHERE datname LIKE '%compatible_db';
datname | datcompatibility
-----td_compatible_db | TD
ora_compatible_db | ORA
mysql_compatible_db | MYSQL
(3 rows)

Table 4-1 Compatibility differences

Compatibility Item	Oracle	Teradata	MySQL
date data type	Converts the date data type to the timestamp data type which stores year, month, day, hour, minute, and second values.	Stores year, month, and date values.	Stores year, month, and date values.

Compatibility Item	Oracle	Teradata	MySQL
Empty string	Only null is available.	Distinguishes empty strings from null values.	Distinguishes empty strings from null values.
Conversion of an empty string to a number	Converts to null .	Converts to 0.	Converts to 0.
Automatic truncation of overlong characters	Not supported.	Supported (set GUC parameter td_compatible_trun cation to ON).	Not supported.
VARCHAR + INT calculation	Converts to BIGINT + INT calculation.	Converts to NUMERIC + NUMERIC calculation.	Converts to BIGINT + INT calculation.
null concatenation	Returns a non- null object after combining a non-null object with null . For example, 'abc' null returns 'abc'.	The strict_text_concat_t d option is added to the GUC parameter behavior_compat_o ptions to be compatible with the Teradata behavior. After the null type is concatenated, null is returned. For example, 'abc' null returns null.	Compatible with MySQL behavior. After the null type is concatenated, null is returned. For example, 'abc' null returns null.

Compatibility Item	Oracle	Teradata	MySQL
Concatenatio n of the char(n) type	Removes spaces and placeholders on the right when the char(n) type is concatenated. For example, cast('a' as char(3)) 'b' returns 'ab'.	After the bpchar_text_withou t_rtrim option is added to the GUC parameter behavior_compat_o ptions, when the char(n) type is concatenated, spaces are reserved and supplemented to the specified length n. Currently, ignoring spaces at the end of a string for comparison is not supported. If the concatenated string contains spaces at the end, the comparison is space-sensitive. For example, cast('a' as char(3)) 'b' returns 'a b'.	Removes spaces and placeholders on the right.
concat(str1,str 2)	Returns the concatenation of all non-null strings.	Returns the concatenation of all non-null strings.	If an input parameter is null , null is returned.
left and right processing of negative values	Returns all characters except the first and last n characters.	Returns all characters except the first and last n characters.	Returns an empty string.

Compatibility Item	Oracle	Teradata	MySQL
lpad(string text, length int [, fill text]) rpad(string text, length int [, fill text])	Fills up the string to the specified length by appending the fill characters (a space by default). If the string is already longer than length then it is truncated (on the right). If fill is an empty string or length is a negative number, null is returned.	If fill is an empty string and the string length is less than the specified length , the original string is returned. If length is a negative number, an empty string is returned.	If fill is an empty string and the string length is less than the specified length, an empty string is returned. If length is a negative number, null is returned.
substr(str, s[, n])	If s is set to 0, the first n characters are returned.	If s is set to 0, the first n characters are returned.	If s is set to 0, an empty string is returned.
substring(str, s[, n]) substring(str [from s] [for n])	If s is set to 0, the first n - 1 characters are returned. If s is < 0, the first s + n - 1 characters are returned. If n is < 0, an error is reported.	If s is set to 0, the first n - 1 characters are returned. If s is < 0, the first s + n - 1 characters are returned. If n is < 0, an error is reported.	If s is set to 0, an empty string is returned. If s is < 0, n characters starting from the last s character are truncated. If n is < 0, an empty string is returned.
trim, ltrim, rtrim, btrim(string[, characters])	Removes the longest string that contains only the characters (a space by default) in the <i>characters</i> from a specified position of the <i>string</i> .	Removes the longest string that contains only the characters (a space by default) in the <i>characters</i> from a specified position of the <i>string</i> .	Removes the string that is equivalent to characters (a space by default) from a specified position of the <i>string</i> .
log(x)	Returns the logarithm with 10 as the base.	Returns the logarithm with 10 as the base.	Returns the natural logarithm.

Compatibility Item	Oracle	Teradata	MySQL
mod(x, 0)	Returns x if the divisor is 0 .	Returns x if the divisor is 0 .	Reports an error if the divisor is 0 .
to_char(date)	The maximum value of the input parameter can only be the maximum value of the timestamp type. The maximum value of the date type is not supported. The return value is of the timestamp type.	The maximum value of the input parameter can only be the maximum value of the timestamp type. The maximum value of the date type is not supported. The return value is of the date type in YYYY/MM/DD format. (The GUC parameter convert_empty_str_to_null_td is enabled.)	Only the timestamp type and the date type support the maximum input value. The return value is of the date type.
to_date, to_timestamp, and to_number processing of empty strings	Returns null .	Returns null. (The convert_empty_str_to_null_td parameter is enabled.)	to_date and to_timestamp returns null. If the parameter passed to to_number is an empty string, 0 is returned.
Return value types of last_day and next_day	Returns values of the timestamp type.	Returns values of the timestamp type.	Returns values of the date type.
Return value type of add_months	Returns values of the timestamp type.	Returns values of the timestamp type.	If the input parameter is of the date type, the return value is of the date type. If the input parameter is of the timestamp type, the return value is of the timestamp type. If the input parameter is of the timestamp type, the return value is of the timestamptz type, the return value is of the timestamptz type, the timestamptz type.

Compatibility Item	Oracle	Teradata	MySQL
CURRENT_TI ME CURRENT_TI ME(p)	Obtains the time of the current transaction. The return value is of the timetz type.	Obtains the time of the current transaction. The return value is of the timetz type.	Obtains the execution time of the current statement. The return value is of the time type.
CURRENT_TI MESTAMP CURRENT_TI MESTAMP(p)	Obtains the execution time of the current statement. The return value is of the timestamptz type.	Obtains the execution time of the current statement. The return value is of the timestamptz type.	Obtains the execution time of the current statement. The return value is of the timestamp type.
CURDATE	Not supported.	Not supported.	Obtains the execution date of the current statement. The return value is of the date type.
CURTIME(p)	Not supported.	Not supported.	Obtains the execution time of the current statement. The return value is of the time type.
LOCALTIME LOCALTIME(p)	Obtains the time of the current transaction. The return value is of the time type.	Obtains the time of the current transaction. The return value is of the time type.	Obtains the execution time of the current statement. The return value is of the timestamp type.
LOCALTIMEST AMP LOCALTIMEST AMP(p)	Obtains the time of the current transaction. The return value is of the timestamp type.	Obtains the time of the current transaction. The return value is of the timestamp type.	Obtains the execution time of the current statement. The return value is of the timestamp type.

Compatibility Item	Oracle	Teradata	MySQL
SYSDATE SYSDATE(p)	Obtains the execution time of the current statement. The return value is of the timestamp(0) type.	Obtains the execution time of the current statement. The return value is of the timestamp(0) type.	Obtains the current system time. The return value is of the timestamp(0) type. This function cannot be pushed down. You are advised to use current_date instead.
now()	Obtains the time of the current transaction. The return value is of the timestamptz type.	Obtains the time of the current transaction. The return value is of the timestamptz type.	Obtains the statement execution time. The return value is of the timestamptz type.
Operator ^	Performs exponentiation.	Performs exponentiation.	Performs the XOR operation.
Expressions GREATEST and LEAST	Returns the comparison results of all non-null input parameters.	Returns the comparison results of all non-null input parameters.	If an input parameter is null , null is returned.
Different input parameter types of CASE, COALESCE, IF, and IFNULL expressions	Reports an error.	Is compatible with behavior of Teradata and supports type conversion between digits and strings. For example, if input parameters for COALESCE are of INT and VARCHAR types, the parameters are resolved as VARCHAR type.	Is compatible with behavior of MySQL and supports type conversion between strings and other types. For example, if input parameters for COALESCE are of DATE, INT, and VARCHAR types, the parameters are resolved as VARCHAR type.
Backquote (`)	Not supported.	Not supported.	Distinguishes MySQL reserved words from common characters.

Running SQL Statements

The following explains how to run SQL statements in a Teradata-compatible database. To see how Oracle and MySQL compatibility types behave differently, switch to either **ora_compatible_db** or **mysql_compatible_db**. Run these SQL

statements, replacing table names with **ora_table** or **mysql_table**, to observe the variations described in the preceding table.

CREATE TABLE td_table(a INT,b VARCHAR(5),c date); INSERT INTO td_table VALUES(1,null,CURRENT_DATE); INSERT INTO td_table VALUES(2,",CURRENT_DATE);

Distinctions Between Empty Strings, NULL Values, and Date Displays

In both Teradata and MySQL compatibility types, empty strings and NULL values are distinct. However, in the Oracle compatibility type, they are considered the same. Additionally, dates convert to timestamps, displaying the year, month, day, hour, minute, and second.

Running SQL Statements in the Teradata Compatibility Type

SELECT a, b, b IS NULL AS null, c FROM td_table;



SELECT CURRENT_DATE;

date 2025-05-07

Running SQL Statements in the Oracle Compatibility Type

SELECT a, b, b IS NULL AS null, c FROM ora_table;



SELECT CURRENT_DATE;

date 2025-05-07

Running SQL Statements in the MySQL Compatibility Type

SELECT a, b, b IS NULL AS null, c FROM mysql_table;



SELECT CURRENT_DATE;

current_date 2025-05-07

Processing Empty Strings

Different from the Oracle database, which processes an empty string as a NULL value, Teradata database converts an empty string to **0** by default. Therefore, when an empty string is queried, value **0** is found.

Similarly, in the Teradata compatibility type, the empty string is converted to **0** of the corresponding numeric type by default. In addition, '-', '+', and ' ' are converted to **0** by default in the Teradata compatibility type, but an error is reported for a decimal point string.

Processing Empty Strings in the Teradata Compatibility Type

SELECT b::int FROM td_table WHERE b = "; 0

Processing Empty Strings in the Oracle Compatibility Type

SELECT b::int FROM ora_table WHERE b = ";

Processing Empty Strings in the MySQL Compatibility Type

SELECT b::int FROM mysql_table WHERE b = ";

0

Automatic Truncation of Overlong Characters

- In the Teradata compatibility type, if td_compatible_truncation is set to on, a long character string will be automatically truncated. If later INSERT statements (not involving foreign tables) insert long strings to columns of char- and varchar-typed columns in the target table, the system will truncate the long strings to ensure no strings exceed the maximum length defined in the target table.
- In the Oracle and MySQL compatibility types, an error is reported when an overlong string is inserted.

Automatic Character Truncation in the Teradata Compatibility Type

SHOW td_compatible_truncation;
SET td_compatible_truncation = ON;
INSERT INTO td_table VALUES(3,'12345678',CURRENT_DATE);
SELECT * FROM td_table WHERE a = 3;

a	b	С
3	12345	2025-05-07

Automatic Character Truncation in the Oracle Compatibility Type

The **td_compatible_truncation** parameter is invalid. If an overlong character is inserted, an error is reported.

```
SHOW td_compatible_truncation;
SET td_compatible_truncation = ON;
INSERT INTO ora_table VALUES(3,'12345678',CURRENT_DATE);

error_msg: data error. STATE: 22001, message: ERROR:
dn 6003 6004: value too long for type character varying(5)
```

Automatic Character Truncation in the MySQL Compatibility Type

Where: referenced column: b

The **td_compatible_truncation** parameter is invalid. If an overlong character is inserted, an error is reported.

```
SHOW td_compatible_truncation;
SET td_compatible_truncation = ON;
INSERT INTO mysql_table VALUES(3,'12345678',CURRENT_DATE);

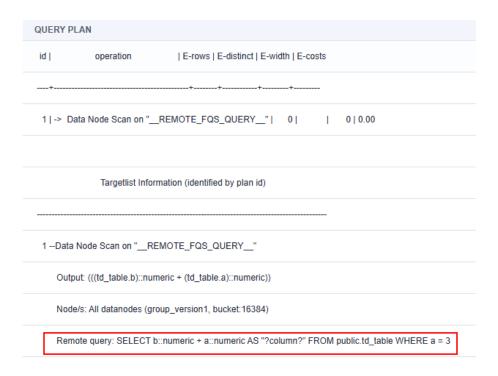
error_msg: data error. STATE: 22001, message: ERROR:
dn_6003_6004: value too long for type character varying(5)
Where: referenced column: b
```

Converting VARCHAR + INT Calculation to NUMERIC + NUMERIC Calculation

- In the Teradata compatibility type, the VARCHAR + INT calculation is converted to the NUMERIC + NUMERIC calculation.
- In the Oracle and MySQL compatibility types, BIGINT + INT calculation is used.

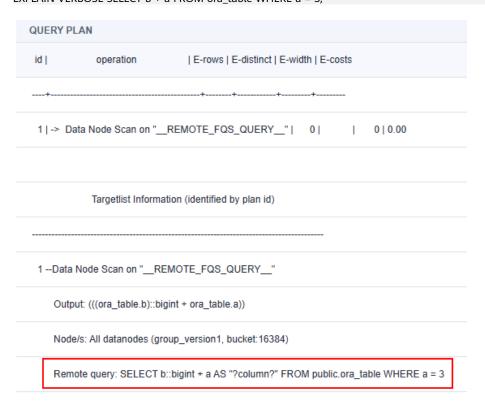
Calculation Conversion in the Teradata Compatibility Type

EXPLAIN VERBOSE SELECT b + a FROM td_table WHERE a = 3;



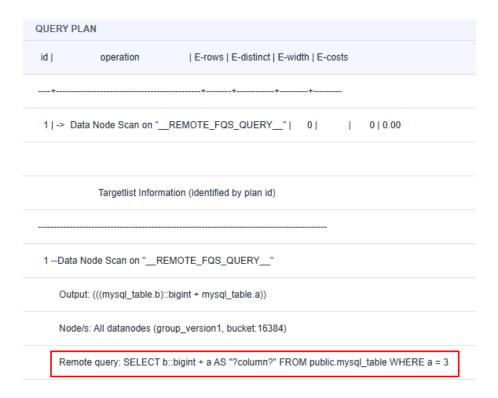
Calculation Conversion in the Oracle Compatibility Type

INSERT INTO ora_table VALUES(3,'12345',CURRENT_DATE); EXPLAIN VERBOSE SELECT b + a FROM ora_table WHERE a = 3;



Calculation Conversion in the MySQL Compatibility Type

INSERT INTO mysql_table VALUES(3,'12345',CURRENT_DATE); EXPLAIN VERBOSE SELECT b + a FROM mysql_table WHERE a = 3;



Concatenating a Null Value

- Teradata compatibility type: The strict_text_concat_td option is added for GUC parameter behavior_compat_options. Concatenating the date, time, number, string, and null will return null.
- Oracle compatibility type: Concatenating a non-null object with NULL values returns a non-null object. For example, 'abc'||null returns 'abc'.
- MySQL compatibility type: Concatenating null types returns null, matching MySQL's behavior. For example, 'abc'||null returns null.

Concatenating a Null Value in the Teradata Compatibility Type

SET behavior_compat_options = 'strict_text_concat_td';
SELECT '2024-02-07 12:12:12'::TIMESTAMP || NULL;
SELECT '12:12:12'::TIME || NULL;
SELECT '12'::TINYINT || NULL;
SELECT 'abc'::CHAR(10) || NULL;



Concatenating a Null Value in the Oracle Compatibility Type

SELECT 'abc'::CHAR(10) || NULL;



Concatenating a Null Value in the MySQL Compatibility Type

SELECT 'abc'::CHAR(10) || NULL;

?column?

abc

Concatenating the char(n) Type

- In the Teradata compatibility type, the **bpchar_text_without_rtrim** option is added to GUC parameter **behavior_compat_options**. When the char(n) type is concatenated, spaces are reserved and padded to the specified length *n*.
- In the Oracle and MySQL compatibility types, spaces are not reserved.

Concatenating the char(n) Type in the Teradata Compatibility Type

If **a** has three spaces and converts to char(10), it fills up to 10 characters since its length is under 10. This uses 10 bytes of storage, with each byte holding 8 bits, totaling 80 bits.

SET behavior_compat_options = 'bpchar_text_without_rtrim';
SELECT bit_length('a '::char(10));



Concatenating the char(n) Type in the Oracle and MySQL Compatibility Types

In the Oracle and MySQL compatibility types, no extra space is reserved. The letter **a** uses just one character and takes up one byte of storage. As a result, the output is **8**.

SELECT bit_length('a '::char(10));

bit_length
8

Using the CONCAT Function

In the Teradata and Oracle compatibility types, **concat(str1,str2)** returns all non-null character strings. In the MySQL compatibility type, it returns **null** if an input parameter contains **null**.

Using the CONCAT Function in the Teradata and Oracle Compatibility Types

SELECT concat(null, 'World!');
concat

World!

Using the CONCAT Function in the MySQL Compatibility Type

SELECT conca	t(null, 'World!');			
concat				
(Null)				

Negative Value Processing in the left and right Functions

In Teradata and Oracle compatibility types, the left function removes the last |n| characters, while the right function removes the first |n| characters. In the MySQL compatibility type, an empty string is returned.

Negative Value Processing in the Teradata and Oracle Compatibility Types



Negative Value Processing in the MySQL Compatibility Type

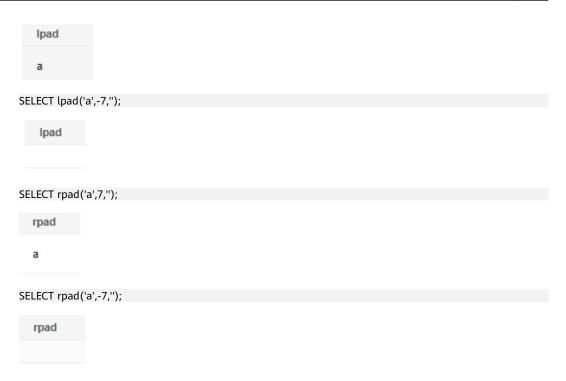
SELECT left('	'abcde', -2);			
left				
SELECT right	t('abcde', -2);			
right				

Empty String Processing in the lpad and rpad Functions

- In the Teradata compatibility type, when **fill** is empty and the input string's length is shorter than **length**, both lpad and rpad functions output the original string unchanged. If the value of **length** is negative, an empty string is returned.
- In the Oracle compatibility type, these functions also return **null**.
- In the MySQL compatibility type, if **fill** is an empty string and the string length is less than the specified **length**, an empty string is returned. If **length** is a negative number, **null** is returned.

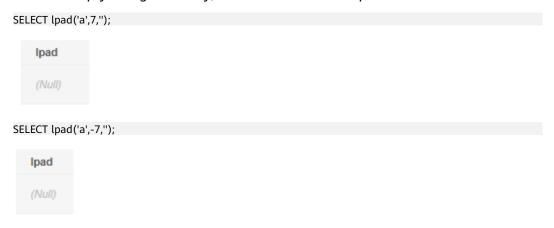
Empty String Processing in the Teradata Compatibility Type (lpad and rpad)

SELECT lpad('a',7,");

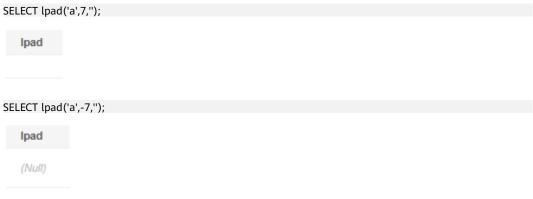


Empty String Processing in the Oracle Compatibility Type (lpad and rpad)

In the Oracle compatibility type, the lpad function returns **null**. The rpad function handles empty strings similarly; no further details are provided here.



Empty String Processing in the MySQL Compatibility Type (lpad and rpad)



Using the SUBSTR Function

In the Teradata and Oracle compatibility types, the substr(str, s[, n]) function gives the first n characters if s is o. In the MySQL compatibility type, it returns an empty string.

Using the SUBSTR Function in the Teradata and Oracle Compatibility Types

SELECT SUBSTR('Hello, World!', 0, 5);



Using the SUBSTR Function in the MySQL Compatibility Type

SELECT SUBSTR('Hello, World!', 0, 5);



Using the SUBSTRING Function

The **substring(str [from s] [for n])** function works differently based on its inputs in the Teradata-compatible and Oracle-compatible modes:

- If **s** is set to **0**, the first **n 1** characters are returned.
- If **s** is less than **0**, the first **s** + **n 1** characters are returned.
- If **n** is less than **0**, an error is reported.

In the MySQL compatibility type:

- If **s** is **0**, an empty string is returned.
- If **s** is less than **0**, characters starting from the last |**s**| character are truncated.
- If **n** is less than **0**, an empty string is returned.

Using the SUBSTRING Function in the Teradata and Oracle Compatibility Types

Example 1: If **s** is **0** and **n** is **5**, the first four characters of **Hello, World!** are returned.

SELECT SUBSTRING('Hello, World!' FROM 0 FOR 5);



Example 2: If **s** is **-1** and **n** is **4**, the first two characters of "Hello, World!" are returned.

SELECT SUBSTRING('Hello, World!' FROM -1 FOR 4);



Example 3: If **n** is **-1**, an error is reported.

SELECT SUBSTRING('Hello, World!' FROM -1 FOR -1);

error_code: DWS.S0010001
error_msg: sql error. STATE: 22011, message: ERROR: negative substring length not allowed
Where: referenced column: substring

Using the SUBSTRING Function in the MySQL Compatibility Type

Example 1: An empty string is returned when **s** is **0** and **n** is **5**.

SELECT SUBSTRING('Hello, World!' FROM 0 FOR 5);



Example 2: If **s** is **-1** and **n** is **4**, the last character is returned.

SELECT SUBSTRING('Hello, World!' FROM -1 FOR 4);



Example 3: If **n** is **-1**, an empty string is returned.

SELECT SUBSTRING('Hello, World!' FROM -1 FOR -1);

substring

Using the BTRIM function

- In the Teradata and Oracle compatibility types, btrim(string[,characters]) removes the longest string that contains only the characters (a space by default) in the characters from a specified position of the string.
- In the MySQL compatibility type, it removes the string that is equivalent to characters (a space by default) from a specified position of the *string*.

Using the BTRIM function in the Teradata and Oracle Compatibility Types

SELECT BTRIM('xxHello Worldxx', 'xxz');



Using the BTRIM function in the MySQL Compatibility Type

SELECT BTRIM('xxHello Worldxx', 'xxz');

btrim

xxHello Worldxx

Using the Log Function

- In the Teradata and Oracle compatibility types, log(x) calculates the base-10 logarithm. For instance, log(100) equals log(10,100).
- In the MySQL compatibility type, it calculates the natural logarithm using base e ($e \approx 2.71828$). For example, log(100) roughly matches log(2.71828,100).

Using the Log Function in the Teradata and Oracle Compatibility Types

SELECT log(100);
log
2

Using the Log Function in the MySQL Compatibility Type

SELECT log(100);
log
4.605170185988092

Using the mod(x,0) Function

In the Teradata and Oracle compatibility types, \mathbf{x} is returned if the divisor of $\mathbf{mod}(\mathbf{x})$ is $\mathbf{0}$. However, in the MySQL compatibility type, this results in an error.

Using the mod(x,0) Function in the Teradata and Oracle Compatibility Types

mod 3

SELECT mod(3,0);

SELECT mod(3,0);

Using the mod(x,0) Function in the MySQL Compatibility Type

error_code: DWS.S0010001
error_msg: sql error. STATE: 22012, message: ERROR: division by zero
Where: referenced column: mod

Using the TO_CHAR Function

- In the Teradata compatibility type, the input parameter's maximum value must be a timestamp, not a date. The function returns a date formatted as YYYY/MM/DD when GUC parameter behavior_compat_options is set to convert_empty_str_to_null_td.
- In the Oracle compatibility type, the input parameter's maximum value must also be a timestamp, not a date. The returned value is a timestamp.
- In the MySQL compatibility type, the maximum value of the input parameter can be a timestamp or date. The return type is date.

Using the TO_CHAR Function in the Teradata Compatibility Type

Using the TO_CHAR Function in the Oracle Compatibility Type

SELECT TO_CHAR(DATE '294276-12-31');

to_char

294276-12-31 00:00:00

SELECT TO_CHAR(DATE '5874897-12-31');

error_code: DWS.S0010011
error_msg: data error. STATE: 22008, message: ERROR: timestamp out of range: "5874897-12-31"
Position: 21
Where: referenced column: to_char

Using the TO_CHAR Function in the MySQL Compatibility Type

SELECT TO_CHAR(DATE '294276-12-31');

to_char

294276-12-31

SELECT TO_CHAR(DATE '5874897-12-31');

to_char

5874897-12-31

to_date, to_timestamp, and to_number Processing of Empty Strings

- In the Teradata compatibility type, if GUC parameter behavior_compat_options is set to convert_empty_str_to_null_td, these functions return null.
- In the Oracle compatibility type, these functions also return **null**.
- In the MySQL compatibility type, **to_date** and **to_timestamp** returns **null**. If the parameter passed to **to_number** is an empty string, **0** is returned.

Empty String Processing in the Teradata and Oracle Compatibility Types (to date, to timestamp, and to number)

```
SET behavior_compat_options = 'convert_empty_str_to_null_td'; --This is only valid for the Teradata-compatible databases.

SELECT TO_DATE(");

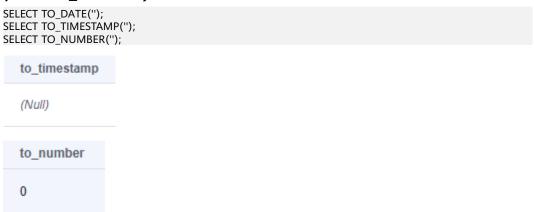
SELECT TO_TIMESTAMP(");

SELECT TO_NUMBER(");

to_number

(Null)
```

Empty String Processing in the MySQL Compatibility Type (to_date, to_timestamp, and to_number)



Using the LAST_DAY and NEXT_DAY Functions

- In the Teradata and Oracle compatibility types, these functions return timestamps.
- In the MySQL compatibility type, these functions return dates.

Using the LAST_DAY and NEXT_DAY Functions in the Teradata and Oracle Compatibility Types

SELECT last_day(to_date('2024-02-07', 'YYYY-MM-DD')) AS cal_result;

cal_result

2024-02-29 00:00:00

SELECT next_day(TIMESTAMP '2024-02-07 00:00:00', 'Sunday')AS cal_result;

cal_result 2024-02-11 00:00:00

Using the LAST_DAY and NEXT_DAY Functions in the MySQL Compatibility Type

SELECT last_day(to_date('2024-02-07', 'YYYY-MM-DD')) AS cal_result;

cal_result 2024-02-29

SELECT next_day(TIMESTAMP '2024-02-07 00:00:00', 'Sunday')AS cal_result;

cal_result 2024-02-11

Using the ADD_MONTHS Function

In the Teradata and Oracle compatibility types, the ADD_MONTHS function returns the timestamp plus integer months.

The ADD_MONTHS function works differently based on the input parameter in the MySQL compatibility type:

- If the input parameter is a date, the function returns a date.
- If the input parameter is a timestamp, the function returns a timestamp.
- If the input parameter is a timestamptz value, the function returns a timestamptz value.

Using the ADD_MONTHS Function in the Teradata and Oracle Compatibility Types

SELECT add_months('2024-02-07'::date,3);

add_months

2024-05-07 00:00:00

Using the ADD_MONTHS Function in the MySQL Compatibility Type

SELECT add_months('2024-02-07'::date,3);

add_months 2024-05-07

SELECT add_months('2024-02-07 00:00:00',3);

add_months 2024-05-07 00:00:00+08

Operator ^

- In the Teradata and Oracle compatibility types, it indicates the exponentiation operation.
- In the MySQL compatibility type, it indicates XOR.

Operator ^ in the Teradata and Oracle Compatibility Types



Operator ^ in the MySQL Compatibility Type

?column?

GREATEST and LEAST Expressions

- In the Teradata and Oracle compatibility types, **GREATEST** and **LEAST** return the comparison results of all non-null input parameters.
- In the MySQL compatibility type, it returns **null** if an input parameter contains **null**.

Using GREATEST and LEAST in the Teradata and Oracle Compatibility Types

SELECT greatest(1,2,3),least(1,2,3),greatest(1,null,3),least(1,null,3);

greatest	least	greatest(1)	least(1)
3	1	3	1

Using GREATEST and LEAST in the MySQL Compatibility Type

SELECT greatest(1,2,3),least(1,2,3),greatest(1,null,3),least(1,null,3);

greatest	least	greatest(1)	least(1)
3	1	(Null)	(Null)

CASE and COALESCE Expressions

- The Teradata compatibility type is compatible with Teradata's behavior and allows converting numbers to strings or vice versa. When using COALESCE with INT and VARCHAR inputs, both values become VARCHAR.
- In the Oracle compatibility type, an error is reported if these expressions are used.
- The MySQL compatibility type is compatible with MySQL's behavior and allows converting data types like strings into others. When using COALESCE with DATE, INT, or VARCHAR inputs, they all convert to VARCHAR.

Using CASE and COALESCE Expressions in the Teradata Compatibility Type

Type resolution for CASE and COALESCE in the Teradata Compatibility Type

- If all inputs are of the same type, and it is not **unknown**, resolve as that type.
- If all inputs are of type **unknown**, resolve as type **text**.
- If inputs are of string type (including unknown which is resolved as type text) and digit type, resolve as the string type. If the inputs are not of the two types, fail.
- If the non-unknown inputs are all of the same type category, choose the input type which is a preferred type in that category, if there is one.
- Convert all inputs to the selected type. Fail if there is not an implicit conversion from a given input to the selected type.

Example 1: Use type resolution with underspecified types in a union as the first example. Here, the unknown-type literal 'b' will be resolved to type **text**.

SELECT text 'a' AS "text" UNION SELECT 'b';

text

a

b

Example 2: Use type resolution in a simple union as the second example. The literal **1.2** is of type **numeric**, and the **integer** value **1** can be cast implicitly to **numeric**, so that type is used.

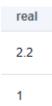
numeric

1.2

1

Example 3: Use type resolution in a transposed union as the third example. Here, since type **real** cannot be implicitly cast to **integer**, but **integer** can be implicitly cast to **real**, the union result type is resolved as **real**.

SELECT 1 AS "real" UNION SELECT CAST('2.2' AS REAL);



Example 4: In the Teradata compatibility type, if input parameters for **COALESCE** are of **int** and **varchar** types, resolve as type **varchar**. In the Oracle compatibility type, an error is reported. Show the execution plan of a statement for querying the types **int** and **varchar** of input parameters for **COALESCE**.

Using CASE and COALESCE Expressions in the Oracle Compatibility Type

In the Teradata compatibility type, if input parameters for **COALESCE** are of **int** and **varchar** types, resolve as type **varchar**. In the Oracle compatibility type, an error is reported. Show the execution plan of a statement for querying the types **int** and **varchar** of input parameters for **COALESCE**.

EXPLAIN VERBOSE select coalesce(a, b) FROM ora_table;

```
error_msg: data type error. STATE: 42804, message: ERROR: COALESCE types integer and character varying cannot be matched Position: 36 Where: referenced column: coalesce
```

Using CASE and COALESCE Expressions in the MySQL Compatibility Type

The MySQL compatibility type is compatible with behavior of MySQL and allows converting data between string and other types. When using **COALESCE** with inputs of DATE, INT, or VARCHAR types, all values are treated as VARCHAR.

5 GaussDB(DWS) Database Security Management

5.1 GaussDB(DWS) User and Permissions Management

5.1.1 GaussDB(DWS) Database User Types

Without separation of permissions, GaussDB(DWS) supports two types of database accounts: administrator and common user. For details about user types and permissions under separation of permissions, see **Separation of Duties in GaussDB(DWS)**.

- The administrator can manage all common users and databases.
- Common users can connect to and access the database, and perform specific database operations and execute SQL statements after being authorized.

Users are authenticated when they log in to the GaussDB(DWS) database. A user can own databases and database objects (such as tables), and grant permissions of these objects to other users and roles. In addition to system administrators, users with the **CREATEDB** attribute can create databases and grant permissions to these databases.

Database User Types

Table 5-1 Database user types

User Type	Description	Allowed Operations	How to Create
Admi nistra tor dbad min	An administrator, also called a system administrator, is an account with the SYSADMIN attribute.	If separation of permissions is not enabled, this account has the highest permission in the system and can perform all operations. The system administrator has the same permissions as the object owner.	 User dbadmin created during cluster creation on the GaussDB(DWS) management console is a system administrator. Use the CREATE USER or ALTER USER syntax to create an administrator. CREATE USER sysadmin WITH SYSADMIN password '{Password}; ALTER USER u1 SYSADMIN;
Com mon user	Use a tool to connect to the database. Have the attributes of specific database system operations, sure as CREATEDB, CREATEROLE, and SYSADMIN. Access database objects. Run SQL statements.		Run the CREATE USER syntax to create a common user. CREATE USER <i>u1</i> PASSWORD '{Password}';
	Private user	A user created with the INDEPENDENT attribute in non-separation-of-permissions mode. Database administrators can manage (DROP, ALTER, and TRUNCATE) objects of private users but cannot access (INSERT, DELETE, SELECT, UPDATE, COPY, GRANT, REVOKE, and ALTER OWNER) the objects before being authorized.	Use the CREATE USER syntax to create a private user. CREATE USER user_independent WITH INDEPENDENT IDENTIFIED BY '{Password}';

5.1.2 GaussDB(DWS) Database User Management

You can use **CREATE USER** and **ALTER USER** to create and manage database users.

- In the non-separation-of-permission mode, a GaussDB(DWS) user account can be created and deleted only by a system administrator or a security administrator with the **CREATEROLE** attribute.
- In separation-of-permission mode, a user account can be created only by a security administrator.

Creating a User

The **CREATE USER** statement is used to create a GaussDB (DWS) user. After creating a user, you can use the user to connect to the database.

- Create common user **u1** and assign the **CREATEDB** attribute to the user. CREATE USER *u1* WITH *CREATEDB* PASSWORD '{Password}';
- To create the system administrator **mydbadmin**, you need to specify the **SYSADMIN** parameter.

CREATE USER mydbadmin sysadmin PASSWORD '{Password}';

- View the created user in the PG_USER view.
 SELECT * FROM pg_user;
- To view user attributes, query the system catalog PG_AUTHID.
 SELECT * FROM pg_authid;

Altering User Attributes

The **ALTER USER** statement is used to alter user attributes, such as changing user passwords or permissions.

Example:

- Rename user u1 to u2.
 ALTER USER u1 RENAME TO u2;
- Grant the CREATEROLE permission to user u1: ALTER USER u1 CREATEROLE;
- For details about how to change the user password, see Setting and Changing a Password.

Locking a User

The **ACCOUNT LOCK** | **ACCOUNT UNLOCK** parameter in the statement is used to lock or unlock a user. A locked user cannot log in to the system. If an account is stolen or illegally accessed, the administrator can manually lock the account. After the account is secured, the administrator can manually unlock the account.

Example:

- To lock user **u1**, run the following command: ALTER USER *u1* ACCOUNT LOCK;
- To unlock user u1, run the following command:
 ALTER USER u1 ACCOUNT UNLOCK;

Deleting a User

The **DROP USER** statement is used to delete one or more GaussDB(DWS) users. An administrator can delete an account that is no longer used. Deleted users cannot be restored.

- If multiple users are deleted at the same time, separate them with commas (,).
- After a user is deleted successfully, all the permissions of the user are also deleted.
- When an account to be deleted is in the active state, it is deleted after the session is disconnected.
- When CASCADE is specified in the DROP USER statement, objects such as tables that depend on the user will be deleted. That is, the objects whose owner is the user are deleted, and the authorizations of other objects to the user are also deleted.

Example:

- -- Delete user u1.
 DROP USER u1;
- Delete account u2 in a cascading manner. DROP USER u2 CASCADE;

5.1.3 Creating a Custom Password Policy for GaussDB(DWS)

When creating or modifying a user, you need to specify a password. GaussDB(DWS) has default password complexity requirements. You can also define database account password policies.

Default GaussDB(DWS) Password Policy

By default, GaussDB(DWS) verifies the password complexity (that is, the GUC parameter **password_policy** is set to **1** by default). The default password policy requires that the password:

- Contain 8 to 32 characters.
- Contain at least three types of the following characters: uppercase letters, lowercase letters, digits, and special characters.
- Cannot be the same as the user name or the user name in reverse order, case insensitive.
- Cannot be the current password or the current password in reverse order.

User-defined Password Policy

The password policy includes the password complexity requirements, password validity period, password reuse settings, password encryption mode, and password retry and lock policies. Different policy items are controlled by the corresponding GUC parameters. For details, see **Security and Authentication (postgresql.conf)**.

Table 5-2 User-defined password policies and corresponding GUC parameters

Password Policy	Parameter	Description	Value Range	Defa ult Value in Gaus sDB(DWS)
Password complexity check	password_p olicy	Specifies whether to check the password complexity when a GaussDB(DW S) account is created or modified.	 O indicates that no password complexity policy is used. Setting this parameter to O leads to security risks. You are advised not to set this parameter to O. I indicates that the default password complexity policy is used. 	1
Password complexity requirement	password_ min_length	Specifies the minimum password length.	An integer ranging from 6 to 999	8
	password_ max_length	Specifies the maximum password length.	An integer ranging from 6 to 999	32
	password_ min_upperc ase	Minimum number of uppercase letters (A-Z)	An integer ranging from 0 to 999 • 0 means no requirements. • 1-999 indicates the minimum number of uppercase letters in the password.	0
	password_ min_lowerc ase	Minimum number of lowercase letters (a-z)	 An integer ranging from 0 to 999 O means no requirements. 1-999 indicates the minimum number of lower letters in the password. 	0

Password Policy	Parameter	Description	Value Range	Defa ult Value in Gaus sDB(DWS)
	password_ min_digital	Minimum number of digits (0-9)	An integer ranging from 0 to 999 • 0 means no requirements. • 1-999 indicates the minimum number of digits in the password.	0
	password_ min_special	Minimum number of special characters (Table 5-3 lists the special characters.)	An integer ranging from 0 to 999 • 0 means no requirements. • 1-999 indicates the minimum number of special characters in the password.	0
Password validity	password_ef fect_time	Password validity period When the number of days in advance a user is notified that the password is about to expire reaches the value of password_no tify_time, the system prompts the user to change the password when the user logs in to the database.	The value is a floating point number ranging from 0 to 999. The unit is day. • 0 indicates the validity period is disabled. • A floating point number from 1 to 999 indicates the validity period of the password. When the password is about to expire or has expired, the system prompts the user to change the password.	90

Password Policy	Parameter	Description	Value Range	Defa ult Value in Gaus sDB(DWS)
	password_n otify_time	Specifies for how many days you are reminded of the password expiry.	 The value is an integer ranging from 0 to 999. The unit is day. O indicates the reminder is disabled. A value ranging from 1 to 999 indicates the number of days prior to password expiration that a user will receive a notification. 	7
Password reuse settings	password_r euse_time	Specifies the number of days after which the password cannot be reused.	 A Floating point number ranging from 0 to 3650. The unit is day. O indicates that the password reuse days are not checked. A positive number indicates that the new password cannot be chosen from passwords in history that are newer than the specified number of days. 	60
	password_r euse_max	Specifies the number of the most recent passwords that the new password cannot be chosen from.	 An integer ranging from 0 to 1000 O indicates that the password reuse times are not checked. A positive number indicates that the new password cannot be chosen from the specified number of the most recent passwords. 	0

Password Policy	Parameter	Description	Value Range	Defa ult Value in Gaus sDB(DWS
Encryption mode	password_e ncryption_t ype	Specifies the password storage encryption mode.	 0 indicates that passwords are encrypted in MD5 mode. The password is encrypted using MD5. This mode is not recommended for users. 1 indicates that passwords are encrypted with SHA-256, which is compatible with the MD5 user authentication method of the PostgreSQL client. The password is stored in ciphertext encrypted by MD5 and SHA256. 2 indicates that password using SHA-256. The password is encrypted using SHA-256. 	1

Password Policy	Parameter	Description	Value Range	Defa ult Value in Gaus sDB(DWS)
Retry and lock	password_lo ck_time	Specifies the duration for a locked account to be automatically unlocked.	a ranging from 0 to 365. The unit is day. e • 0 indicates that the	1
	failed_login _attempts	If the number of incorrect password attempts reaches the value of failed_login_a ttempts, the account is locked and will be automatically unlocked in X (which indicates the value of password_lock_time) seconds.	 An integer ranging from 0 to 1000 O indicates that the automatic locking function does not take effect. A positive number indicates that an account is locked when the number of incorrect password attempts reaches the value of failed_login_attempts. 	10

No. Chara No. Charac No. Charac No. Charact cter ter ter er 1 ~ 9 17 25 < 2 ļ 10 (18 [26 3 11) 19 { 27 **@** > 4 # } 12 20 28 \$] ? 5 13 21 29 6 % 14 = 22 _ 7 Λ 15 23 + : 8 24 & 16 \

Table 5-3 Special characters

Example of User-defined Password Policies

Example 1: Configure the password complexity parameter password_policy.

- 1. Log in to the **GaussDB(DWS) console**.
- 2. Choose **Dedicated Clusters** > **Clusters** in the navigation tree on the left.
- 3. In the cluster list, find the target cluster and click the cluster name. The **Cluster Information** page is displayed.
- 4. Click the **Parameters** tab, change the value of **password_policy**, and click **Save**. The **password_policy** parameter takes effect immediately after being modified. You do not need to restart the cluster.

Figure 5-1 password_policy



Example 2: Configure password_effect_time for password validity period.

- 1. Log in to the GaussDB(DWS) management console.
- 2. In the navigation pane on the left, choose **Clusters**.
- 3. In the cluster list, find the target cluster and click the cluster name. The **Cluster Information** page is displayed.
- 4. Click the **Parameters** tab, change the value of **password_effect_time**, and click **Save**. The modification of **password_effect_time** takes effect immediately. You do not need to restart the cluster.

Modify Records Save Cancel Synchronized ? Parameter Name

Enter a parameter name. Parameter Name ↓≡ partition_mem_batch 256 256 1~65.535 No To optimize the inserting of column-store partitioned tables in batches, data is ca. ssword effect time 90 Day 0 ~ 999 0~2 Specifies the encryption type of user passwords 0 indicates that passwords are e password encryption type password lock time 0 ~ 365 Specifies the duration before an account is automatically unlocked. 0 indicates th.

Figure 5-2 password_effect_time

Setting and Changing a Password

• Both system administrators and common users need to periodically change their passwords to prevent the accounts from being stolen.

For example, to change the password of the user **user1**, connect to the database as the administrator and run the following command:

ALTER USER user1 IDENTIFIED BY 'newpassword' REPLACE 'oldpassword;

■ NOTE

The password must meet input requirements, or the execution will fail.

• An administrator can change its own password and other accounts' passwords. With the permission for changing other accounts' passwords, the administrator can resolve a login failure when a user forgets its password.

To change the password of the user **joe**, run the following command:

ALTER USER joe IDENTIFIED BY 'password;

◯ NOTE

- System administrators are not allowed to change passwords for each other.
- When a system administrator changes the password of a common user, the original password is not required.
- However, when a system administrator changes its own password, the original password is required.
- Password verification

Password verification is required when you set the user or role in the current session. If the entered password is inconsistent with the stored password of the user, an error is reported.

To set the password of the user **joe**, run the following command:

SET ROLE *joe* **PASSWORD** '*password*;

If the following information is displayed, the role setting has been modified: SET ROLE

5.1.4 GaussDB(DWS) Database Permissions Management

Permission Overview

Permissions are used to control whether a user is allowed to access a database object (including schemas, tables, functions, and sequences) to perform operations such as adding, deleting, modifying, querying, and creating a database object.

Permission management in GaussDB(DWS) falls into three categories:

• System permissions

System permissions are also called user attributes, including **SYSADMIN**, **CREATEDB**, **CREATEROLE**, **AUDITADMIN**, and **LOGIN**.

They can be specified only by the **CREATE ROLE** or **ALTER ROLE** syntax. The **SYSADMIN** permission can be granted and revoked using **GRANT ALL PRIVILEGE** and **REVOKE ALL PRIVILEGE**, respectively. System permissions cannot be inherited by a user from a role, and cannot be granted using **PUBLIC**.

Object permissions

Permissions on a database object (table, view, column, database, function, schema, or tablespace) can be granted to a role or user. The **GRANT** command can be used to grant permissions to a user or role. These permissions granted are added to the existing ones.

Permissions

Grant a role's or user's permissions to one or more roles or users. In this case, every role or user can be regarded as a set of one or more database permissions.

If **WITH ADMIN OPTION** is specified, the member can in turn grant permissions in the role to others, and revoke permissions in the role as well. If a role or user granted with certain permissions is changed or revoked, the permissions inherited from the role or user also change.

A database administrator can grant permissions to and revoke them from any role or user. Roles having **CREATEROLE** permission can grant or revoke membership in any role that is not an administrator.

Hierarchical Permission Management

GaussDB(DWS) implements a hierarchical permission management on databases, schemas, and data objects.

- Databases cannot communicate with each other and share very few resources. Their connections and permissions can be isolated. The database cluster has one or more named databases. Users and roles are shared within the entire cluster, but their data is not shared. That is, a user can connect to any database, but after the connection is successful, any user can access only the database declared in the connection request.
- Schemas share more resources than databases do. User permissions on schemas and subordinate objects can be flexibly configured using the **GRANT** and **REVOKE** syntax. Each database has one or more schemas. Each schema contains various types of objects, such as tables, views, and functions. To

- access an object contained in a specified schema, a user must have the **USAGE** permission on the schema.
- After an object is created, by default, only the object owner or system
 administrator can query, modify, and delete the object. To access a specific
 database object, for example, table1, other users must be granted the
 CONNECT permission of database, the USAGE permission of schema, and the
 SELECT permission of table1. To access an object at the bottom layer, a user
 must be granted the permission on the object at the upper layer. To create or
 delete a schema, you must have the CREATE permission on its database.

Schema ______table1 table2 view

Figure 5-3 Hierarchical Permission Management

Roles

The permission management model of GaussDB(DWS) is a typical implementation of the role-based permission control (RBAC). It manages users, roles, and permissions through this model.

A role is a set of permissions.

- The concept of "user" is equivalent to that of "role". The only difference is that "user" has the **login** permission while "role" has the **nologin** permission.
- Roles are assigned with different permissions based on their responsibilities in the database system. A role is a set of database permissions and represents the behavior constraints of a database user or a group of data users.
- Roles and users can be converted. You can use **ALTER** to assign the **login** permission to a role.
- After a role is granted to a user through GRANT, the user will have all the
 permissions of the role. It is recommended that roles be used to efficiently
 grant permissions. For example, you can create different roles of design,
 development, and maintenance personnel, grant the roles to users, and then
 grant specific data permissions required by different users. When permissions
 are granted or revoked at the role level, these permission changes take effect
 for all the members of the role.
- In non-separation-of-duty scenarios, a role can be created, modified, and deleted only by a system administrator or a user with the CREATEROLE attribute. In separation-of-duty scenarios, a role can be created, modified, and deleted only by a user with the CREATEROLE attribute.

To view all roles, query the system catalog **PG_ROLES**.

SELECT * FROM PG_ROLES;

For how to create, modify, and delete a role, see "CREATE ROLE/ALTER ROLE/DROP ROLE" in *SQL Syntax Reference*.

Preset Roles

GaussDB(DWS) provides a group of preset roles. Their names start with **gs_role_**. These roles allow access to operations that require high permissions. You can grant these roles to other users or roles in the database for them to access or use specific information and functions. Exercise caution and ensure security when using preset roles.

The following table describes the permissions of preset roles.

Table 5-4 Permissions of preset roles

Role	Permission
gs_role_signal_bac kend	Invokes functions such as pg_cancel_backend, pg_terminate_backend, pg_terminate_query, pg_cancel_query, pgxc_terminate_query, and pgxc_cancel_query to cancel or terminate sessions, excluding those of the initial users.
gs_role_read_all_s tats	Reads the system status view and uses various extension-related statistics, including information that is usually visible only to system administrators. For example: Resource management views: pgxc_wlm_operator_history pgxc_wlm_operator_statistics pgxc_wlm_session_info pgxc_wlm_session_statistics pgxc_wlm_workload_records pgxc_workload_sql_count pgxc_workload_sql_elapse_time pgxc_workload_transaction Status information views: pgxc_stat_activity pgxc_get_table_skewness table_distribution pgxc_total_memory_detail pgxc_os_run_info pg_nodes_memory pgxc_instance_time pgxc_redo_stat

Role	Permission
gs_role_analyze_a ny	A user with the system-level ANALYZE permission can skip the schema permission check and perform ANALYZE on all tables.
gs_role_vacuum_a ny	A user with the system-level VACUUM permission can skip the schema permission check and perform ANALYZE on all tables.
gs_redaction_polic y	A user with the permission to create, modify, and delete data masking policies and can execute CREATE ALTER DROP REDACTION POLICY on all tables. Clusters of 9.1.0 and later versions support this function.

Restrictions on using preset roles:

- **gs_role_** is the name field dedicated to preset roles in the database. Do not create users or roles starting with **gs_role_** or rename existing users or roles starting with **gs_role_**.
- Do not perform **ALTER** or **DROP** operations on preset roles.
- By default, a preset role does not have the **LOGIN** permission, so there is no preset login password for the role.
- The gsql meta-commands \du and \dg do not display information about preset roles. However, if **PATTERN** is specified, information about preset roles will be displayed.
- If the separation of permissions is disabled, the system administrator and users with the **ADMIN OPTION** permission of preset roles are allowed to perform GRANT and REVOKE operations on preset roles. If the separation of permissions is enabled, the security administrator (with the **CREATEROLE** attribute) and users with the **ADMIN OPTION** permission of preset roles are allowed to perform GRANT and REVOKE operations on preset roles. Example: GRANT gs_role_signal_backend TO user1; REVOKE gs_role_signal_backend FROM user1;

Granting or Revoking Permissions

A user who creates an object is the owner of this object. By default, **Separation of Duties in GaussDB(DWS)** is disabled after cluster installation. A database system administrator has the same permissions as object owners.

After an object is created, only the object owner or system administrator can query, modify, and delete the object, and grant permissions for the object to other users through **GRANT** by default. To enable a user to use an object, the object owner or administrator can run the **GRANT** or **REVOKE** command to grant permissions to or revoke permissions from the user or role.

• Run the **GRANT** statement to grant permissions.

For example, grant the permission of schema **myschema** to role **u1**, and grant the **SELECT** permission of table **myschema.t1** to role **u1**.

GRANT USAGE ON SCHEMA *myschema* TO *u1*;

GRANT *SELECT* ON TABLE *myschema.t1* to *u1*;

Run the REVOKE command to revoke a permission that has been granted.
 For example, revoke all permissions of user u1 on the myschema.t1 table.
 REVOKE ALL PRIVILEGES ON myschema.t1 FROM u1;

5.1.5 Separation of Duties in GaussDB(DWS)

By default, the system administrator with the **SYSADMIN** attribute has the highest permission in the system. To avoid risks caused by centralized permissions, you can enable the separation of permissions to delegate system administrator permissions to security administrators and audit administrators.

- After the separation of permissions is enabled, a system administrator does
 not have the CREATEROLE attribute (security administrator) and
 AUDITADMIN attribute (audit administrator). That is, you do not have the
 permissions for creating roles and users and the permissions for viewing and
 maintaining database audit logs. For details about the CREATEROLE and
 AUDITADMIN attributes, see CREATE ROLE.
- After the separation of permissions is enabled, system administrators have the permissions only for the objects owned by them.

For how to configure permission separation, see Configuring Separation of Duties for the GaussDB(DWS) Cluster

For details about permission changes before and after enabling the separation of permissions, see **Table 5-5** and **Table 5-6**.

Table 5-5 Default user permissions

Object	System Administrator	Security Administrator	Audit Administrato r	Common User
Tables pace	Can create, modify, delete, access, and allocate tablespaces.	Cannot create, modify, delete, or allocate tablespaces, with authorization required for accessing tablespaces.		
Table	Has permissions for all tables.	Has permissions for its own tables, but does not have permissions for other users' tables.		
Index	Can create indexes on all tables.	Can create indexes on their own tables.		
Schem a	Has permissions for all schemas.	Has all permissions for its own schemas, but does not have permissions for other users' schemas.		
Functio n	Has permissions for all functions.	Has permissions for its own functions, has the call permission for other users' functions in the public schema, but does not have permissions for other users' functions in other schemas.		

Object	System Administrator	Security Administrator	Audit Administrato r	Common User
Custo mized view	Has permissions for all views.		s for its own view as for other users	
System catalog and system view	Has permissions for querying all system catalogs and views.	catalogs and vie	s for querying on ews. For details, s System Catalog	ee

Table 5-6 Changes in permissions after the separation of permissions

Objec t	System Administrator	Securi ty Admi nistra tor	Audit Admi nistra tor	Common User
Tables pace	No change	No change		
Table	Permissions reduced Has all permissions for its own tables, but does not have permissions for other users' tables in their schemas.	No change		
Index	Permissions reduced Can create indexes on its own tables.	No change		
Sche ma	Permissions reduced Has all permissions for its own schemas, but does not have permissions for other users' schemas.	No char	nge	
Functi on	Permissions reduced Has all permissions for its own functions, but does not have permissions for other users' functions in their schemas.	No char	nge	
Custo mized view	Permissions reduced Has all permissions for its own views and other users' views in the public schema, but does not have permissions for other users' views in their schemas.	No change		

Objec t	System Administrator	Securi ty Admi nistra tor	Audit Admi nistra tor	Common User
Syste m catalo g and syste m view	No change	No chang e	No chang e	Has no permissio n for viewing any system catalogs or views.

5.2 GaussDB(DWS) Sensitive Data Management

5.2.1 GaussDB(DWS) Row-Level Access Control

The row-level access control feature enables database access control to be accurate to each row of data tables. In this way, the same SQL query may return different results for different users.

You can create a row-level access control policy for a data table. The policy defines an expression that takes effect only for specific database users and SQL operations. When a database user accesses the data table, if a SQL statement meets the specified row-level access control policies of the data table, the expressions that meet the specified condition will be combined by using **AND** or **OR** based on the attribute type (**PERMISSIVE** | **RESTRICTIVE**) and applied to the execution plan in the query optimization phase.

Row-level access control is used to control the visibility of row-level data in tables. By predefining filters for data tables, the expressions that meet the specified condition can be applied to execution plans in the query optimization phase, which will affect the final execution result. Currently, the SQL statements that can be affected include **SELECT**, **UPDATE**, and **DELETE**.

Scenario 1: A table summarizes the data of different users. Users can view only their own data.

```
-- Create users alice, bob, and peter.

CREATE ROLE alice PASSWORD 'password;

CREATE ROLE bob PASSWORD 'password;

CREATE ROLE peter PASSWORD 'password;

-- Create the public.all_data table that contains user information.

CREATE TABLE public.all_data(id int, role varchar(100), data varchar(100));

-- Insert data into the data table.

INSERT INTO all_data VALUES(1, 'alice', 'alice data');

INSERT INTO all_data VALUES(2, 'bob', 'bob data');

INSERT INTO all_data VALUES(3, 'peter', 'peter data');

-- Grant the read permission for the all_data table to users alice, bob, and peter.

GRANT SELECT ON all data TO alice, bob, peter;
```

```
    Enable row-level access control.

ALTER TABLE all_data ENABLE ROW LEVEL SECURITY;
-- Create a row-level access control policy to specify that the current user can view only their own data.
CREATE ROW LEVEL SECURITY POLICY all data rls ON all data USING(role = CURRENT USER);
-- View table details.
\d+ all_data
                   Table "public.all_data"
               Type | Modifiers | Storage | Stats target | Description
Column I
                              | plain |
id | integer |
role | character varying(100) | extended |
data | character varying(100) |
                                     | extended |
Row Level Security Policies:
  POLICY "all_data_rls"
   USING (((role)::name = "current_user"()))
Has OIDs: no
Distribute By: HASH(id)
Location Nodes: ALL DATANODES
Options: orientation=row, compression=no, enable_rowsecurity=true
-- Switch to user alice and run SELECT * FROM all_data.
SET ROLE alice PASSWORD 'password;
SELECT * FROM all_data;
id | role | data
 1 | alice | alice data
(1 row)
EXPLAIN(COSTS OFF) SELECT * FROM all_data;
                 QUERY PLAN
Streaming (type: GATHER)
 Node/s: All datanodes
 -> Seq Scan on all_data
     Filter: ((role)::name = 'alice'::name)
Notice: This query is influenced by row level security feature
(5 rows)
-- Switch to user peter and run SELECT * FROM .all_data.
SET ROLE peter PASSWORD 'password;
SELECT * FROM all_data;
id | role | data
3 | peter | peter data
(1 row)
EXPLAIN(COSTS OFF) SELECT * FROM all_data;
                OUERY PLAN
Streaming (type: GATHER)
 Node/s: All datanodes
  -> Seq Scan on all_data
     Filter: ((role)::name = 'peter'::name)
Notice: This query is influenced by row level security feature
(5 rows)
```

5.2.2 GaussDB(DWS) Data Masking

GaussDB(DWS) provides the column-level dynamic data masking (DDM) function. For sensitive data (such as the ID card number, mobile number, and bank card number), the DDM function is used to redact the original data to protect data security and user privacy.

• Creating a data masking policy for a table

GaussDB(DWS) uses the **CREATE REDACTION POLICY** syntax to create a data masking policy on a table (Do not perform masking), **MASK_FULL** (Mask data into a fixed value), and **MASK_PARTIAL** (Perform partial masking based on the character type, numeric type, or time type.) to specify the application scope of the masking policy.

- Modifying the data masking policy of a table
 - The **ALTER REDACTION POLICY** syntax is used to modify the expression for enabling a masking policy, rename a masking policy, and add, modify, or delete masked columns.
- Deleting the masking policy of a table
 - The **DROP REDACTION POLICY** syntax is used to delete the masking function information of a masking policy on all columns of a table.
- Viewing the masking policy and masked columns

Masking policy information is stored in the system catalog PG_REDACTION_POLICY, and masked column information is stored in the system catalog PG_REDACTION_COLUMN. You can view information about the masking policy and masked columns in the system views REDACTION_POLICIES and REDACTION_COLUMNS.

Precautions

- The data masking feature is controlled by feature_support_options. If error information "REDACTION POLICY is not yet supported, please add enable_data_redaction into feature_support_options." or "Cannot use the feature of data redaction, please drop existed redaction policy or add enable_data_redaction into feature_support_options." is displayed, the data masking feature is not enabled. In this case, you need to contact technical support to set feature support options.
- Generally, you can run the SELECT statement to view the data masking result. If a statement has the following features, sensitive data may be deliberately obtained. In this case, an error will be reported during statement execution.
 - The GROUP BY clause references the Target Entry containing masked columns as the target column.
 - DISTINCT works on the output masked columns.
 - The statement contains CTE.
 - Operations on sets are involved.
 - The target columns of a subquery are not masked columns of the base table, but the expressions or function calls for masked columns of the base table.
- You can use COPY TO or GDS to export the masked data. Due to the irreversibility of the data masking, secondary masking of the data is meaningless.
- Do not set target columns of UPDATE, MERGE INTO, and DELETE statements to masked columns.
- The UPSERT statement allows you to insert update data through EXCLUDED. If data in the base table is updated by referencing masked columns, the data may be modified by mistake. As a result, an error will be reported during the execution.

- In the 8.2.1 cluster version, multiple masking policies can be created for the same table to implement diversified sensitive data classification. The principles for selecting and applying masking policies are as follows:
 - Select the policy with the largest policy_order among multiple candidate policies that meet the requirements of the current session. A larger policy_order indicates a later creation.
 - During data masking, the DML statement inherits only the policy with the largest policy_order.

Examples

The following uses the employee table **emp**, table owner **alice**, and roles **matu** and **july** as an example to illustrate the data masking process. The **emp** table contains private data such as the employee name, mobile number, email address, bank card number, and salary.

Step 1 After connecting to the database as the administrator, create roles **alice**, **matu**, and **july**.

```
CREATE ROLE alice PASSWORD 'password;
CREATE ROLE matu PASSWORD 'password;
CREATE ROLE july PASSWORD 'password;
```

- **Step 2** Grant schema permissions on the current database to **alice**, **matu**, and **july**. GRANT ALL PRIVILEGES on schema *public* to alice, matu, july;
- **Step 3** Switch to role **alice**, create the **emp** table, and insert three pieces of employee information.

```
SET ROLE alice PASSWORD 'password;
```

CREATE TABLE emp(id int, name varchar(20), phone_no varchar(11), card_no number, card_string varchar(19), email text, salary numeric(100, 4), birthday date);

INSERT INTO emp VALUES(1, 'anny', '13420002340', 1234123412341234, '1234-1234-1234-1234', 'smithWu@163.com', 10000.00, '1999-10-02'); INSERT INTO emp VALUES(2, 'bob', '18299023211', 3456345634563, '3456-3456-3456-3456', '66allen_mm@qq.com', 9999.99, '1989-12-12'); INSERT INTO emp VALUES(3, 'cici', '15512231233', NULL, NULL, 'jonesishere@sina.com', NULL, '1992-11-06');

- **Step 4 alice** grants the read permission on the **emp** table to **matu** and **july**.

 GRANT SELECT ON emp TO matu, july;
- Step 5 Create the masking policy mask_emp: Only user alice can view all employee information. User matu and july cannot view employee bank card numbers and salary data. The card_no column is of the numeric type and all of its data is masked into 0 by the MASK_FULL function. The card_string column is of the character type and part of its data is masked by the MASK_PARTIAL function based on the specified input and output formats. The salary column is of the numeric type and the MASK_PARTIAL function is used to mask all digits before the penultimate digit using the number 9.

```
CREATE REDACTION POLICY mask_emp ON emp WHEN (current_user IN ('matu', 'july'))
ADD COLUMN card_no WITH mask_full(card_no),
ADD COLUMN card_string WITH mask_partial(card_string, 'VVVVFVVVVFVVVVFVVVV','WVVV-VVVV-VVVV-VVVV-VVVV,'#',1,12),
ADD COLUMN salary WITH mask_partial(salary, '9', 1, length(salary) - 2);
```

Step 6 Switch to **matu** and **july** and view the employee table **emp**.

```
SET ROLE matu PASSWORD 'password;
SELECT * FROM emp;
```

```
id | name | phone_no | card_no | card_string
                                                              | salary |
                                                                           birthday
                                                   email
 1 | anny | 13420002340 |
                           0 | ####-###-1234 | smithWu@163.com
                                                                          199999,99901
1999-10-02 00:00:00
 2 | bob | 18299023211 |
                          0 | ####-###-###-3456 | 66allen_mm@qq.com | 9999.9990 |
1989-12-12 00:00:00
3 | cici | 15512231233 |
                                        | jonesishere@sina.com |
                                                                    1 1992-11-06 00:00:00
(3 rows)
SET ROLE july PASSWORD 'password';
SELECT * FROM emp:
id | name | phone_no | card_no | card_string |
                                                                           birthday
                                                              | salarv |
1 | anny | 13420002340 |
                           0 | ####-###-1234 | smithWu@163.com
                                                                          | 99999.9990 |
1999-10-02 00:00:00
2 | bob | 18299023211 |
                          0 | ####-###-###-3456 | 66allen_mm@qq.com | 9999.9990 |
1989-12-12 00:00:00
3 | cici | 15512231233 |
                                                                    | 1992-11-06 00:00:00
                                        | jonesishere@sina.com |
(3 rows)
```

Step 7 If you want **matu** to have the permission to view all employee information, but do not want **july** to have. In this case, you only need to modify the effective scope of the policy.

```
SET ROLE alice PASSWORD 'password;
ALTER REDACTION POLICY mask_emp ON emp WHEN(current_user = 'july');
```

Step 8 Switch to users **matu** and **july** and view the **emp** table again, respectively.

```
SET ROLE matu PASSWORD 'password;
SELECT * FROM emp;
                                                                                birthday
id | name | phone_no | card_no | card_string
                                                         email
                                                                   | salary |
 1 | anny | 13420002340 | 1234123412341234 | 1234-1234-1234 | smithWu@163.com
10000.0000 | 1999-10-02 00:00:00
 2 | bob | 18299023211 | 3456345634563456 | 3456-3456-3456 | 66allen_mm@qq.com |
9999.9900 | 1989-12-12 00:00:00
3 | cici | 15512231233 |
                                             | jonesishere@sina.com |
                                                                          | 1992-11-06 00:00:00
(3 rows)
SET ROLE july PASSWORD 'password;
SELECT * FROM emp;
id | name | phone_no | card_no | card_string
                                                   email
                                                              | salary |
                                                                           birthday
1 | anny | 13420002340 |
                         0 | ####-###-###-1234 | smithWu@163.com
1999-10-02 00:00:00
2 | bob | 18299023211 |
                          0 | ####-###-###-3456 | 66allen_mm@qq.com | 9999.9990 |
1989-12-12 00:00:00
3 | cici | 15512231233 |
                                       | jonesishere@sina.com |
                                                                    | 1992-11-06 00:00:00
(3 rows)
```

Step 9 The information in the **phone_no**, **email**, and **birthday** columns is private data. Update masking policy **mask_emp** and add three masked columns.

```
SET ROLE alice PASSWORD 'password;
ALTER REDACTION POLICY mask_emp ON emp ADD COLUMN phone_no WITH mask_partial(phone_no, '*', 4);
ALTER REDACTION POLICY mask_emp ON emp ADD COLUMN email WITH mask_partial(email, '*', 1, position('@' in email));
ALTER REDACTION POLICY mask_emp ON emp ADD COLUMN birthday WITH mask_full(birthday);
```

Step 10 Switch to **july** and view data in the **emp** table.

Step 11 Query **redaction_policies** and **redaction_columns** to view details about the current redaction policy **mask emp**.

```
SELECT * FROM redaction policies;
object_schema | object_owner | object_name | policy_name |
                                                           expression
                                                                           | enable |
policy_description | inherited
public
       | alice | emp | mask_emp | ("current_user"() = 'july'::name) | t
(1 row)
SELECT object_name, column_name, function_info FROM redaction_columns;
object_name | column_name |
                                                    function info
         | card_no | mask_full(card_no)
emp
         | card_string | mask_partial(card_string, 'VVVVFVVVVFVVVV'::text, 'VVVV-VVVV-VVVV-
emp
VVVV'::text, '#'::text, 1, 12)
                   | mask_partial(email, '*'::text, 1, "position"(email, '@'::text))
emp
         | email
                 | mask_partial(salary, '9'::text, 1, (length((salary)::text) - 2))
emp
         salary
         | birthday | mask_full(birthday)
emp
         | phone_no | mask_partial(phone_no, '*'::text, 4)
emp
(6 rows)
```

Step 12 Add the salary_info column. To replace the salary information in text format with *.*, you can create a user-defined masking function. In this step, you can use the PL/pgSQL to define the masking function mask_regexp_salary. To create a masking column, you simply need to customize the function name and parameter list. For details, see GaussDB(DWS) User-Defined Functions.

Step 13 If there is no need to set a redaction policy for the **emp** table, delete redaction policy **mask_emp**.

```
SET ROLE alice PASSWORD 'password';
DROP REDACTION POLICY mask_emp ON emp;
```

----End

5.2.3 Encrypting and Decrypting GaussDB(DWS) Strings

GaussDB(DWS) supports encryption and decryption of strings using the following functions:

gs_encrypt(encryptstr, keystr, cryptotype, cryptomode, hashmethod)

Description: Encrypts an encryptstr string using the keystr key based on the encryption algorithm specified by cryptotype and cryptomode and the HMAC algorithm specified by hashmethod, and returns the encrypted string. cryptotype can be aes128, aes192, and aes256. cryptomode is cbc. hashmethod can be sha256, sha384, sha512, or sm3. Currently, the following types of data can be encrypted: numerals supported in the database; character type; RAW in binary type; and DATE, TIMESTAMP, and SMALLDATETIME in date/time type. The keystr length is related to the encryption algorithm and contains 1 to KeyLen bytes. If cryptotype is aes128, KeyLen is 16; if cryptotype is aes192, KeyLen is 24; if cryptotype is aes256, KeyLen is 32.

Return type: text

Length of the return value: at least $4 \times [(\text{maclen} + 56)/3]$ bytes and no more than $4 \times [(\text{Len} + \text{maclen} + 56)/3]$ bytes, where **Len** indicates the string length (in bytes) before the encryption and **maclen** indicates the length of the HMAC value. If **hashmethod** is **sha256** or **sm3**, **maclen** is **32**; if **hashmethod** is **sha384**, **maclen** is **48**; if **hashmethod** is **sha512**, **maclen** is **64**. That is, if **hashmethod** is **sha256** or **sm3**, the returned string contains 120 to $4 \times [(\text{Len} + 88)/3]$ bytes; if **hashmethod** is **sha384**, the returned string contains 140 to $4 \times [(\text{Len} + 104)/3]$ bytes; if **hashmethod** is **sha512**, the returned string contains 160 to $4 \times [(\text{Len} + 120)/3]$ bytes.

Example:

SELECT gs_encrypt('GaussDB(DWS)', '1234', 'aes128', 'cbc', 'sha256');

gs_encrypt

AAAAAAAAAAACcFjDcCSbop7D87sOa2nxTFrkE9RJQGK34ypgrOPsFJlqggl8tl
+eMDcQYT3po98wPCC7VBfhv7mdBy7lVnzdrp0rdMrD6/zTl8w0v9/s2OA==
(1 row)

∩ NOTE

- A decryption password is required during the execution of this function. For security purposes, the gsql tool does not record this function in the execution history. That is, the execution history of this function cannot be found in **gsql** by paging up and down.
- Do not use the **ge_encrypt** and **gs_encrypt_aes128** functions for the same data table.
- gs_decrypt(decryptstr, keystr, cryptotype, cryptomode, hashmethod)
 Description: Decrypts a decryptstr string using the keystr key based on the encryption algorithm specified by cryptotype and cryptomode and the HMAC algorithm specified by hashmethod, and returns the decrypted string. The keystr used for decryption must be consistent with that used for encryption. keystr cannot be empty.

Return type: text

Example:

SELECT gs_decrypt('AAAAAAAAAACcFjDcCSbop7D87sOa2nxTFrkE9RJQGK34ypgrOPsFJIqggl8tl +eMDcQYT3po98wPCC7VBfhv7mdBy7IVnzdrp0rdMrD6/zTl8w0v9/s2OA==', '1234', 'aes128', 'cbc',

```
'sha256');
gs_decrypt
------
GaussDB(DWS)
(1 row)
```


- A decryption password is required during the execution of this function. For security purposes, the gsql tool does not record this function in the execution history. That is, the execution history of this function cannot be found in **gsql** by paging up and down.
- This function works with the **gs_encrypt** function, and the two functions must use the same encryption algorithm and HMAC algorithm.
- gs_encrypt_aes128(encryptstr,keystr)

Description: Encrypts **encryptstr** strings using **keystr** as the key and returns encrypted strings. The length of **keystr** ranges from 1 to 16 bytes. Currently, the following types of data can be encrypted: numerals supported in the database; character type; RAW in binary type; and DATE, TIMESTAMP, and SMALLDATETIME in date/time type.

Return type: text

Length of the return value: At least 92 bytes and no more than (4*[*Len*/3]+68) bytes, where *Len* indicates the length of the data before encryption (unit: byte).

Example:

□ NOTE

- A decryption password is required during the execution of this function. For security purposes, the gsql tool does not record this function in the execution history. That is, the execution history of this function cannot be found in gsql by paging up and down.
- Do not use the ge_encrypt and gs_encrypt_aes128 functions for the same data table.
- gs_decrypt_aes128(decryptstr,keystr)

Description: Decrypts a **decryptstr** string using the **keystr** key and returns the decrypted string. The **keystr** used for decryption must be consistent with that used for encryption. **keystr** cannot be empty.

Return type: text

Example:

□ NOTE

- A decryption password is required during the execution of this function. For security purposes, the gsql tool does not record this function in the execution history. That is, the execution history of this function cannot be found in **gsql** by paging up and down.
- This function works with the **gs_encrypt_aes128** function.
- qs_hash(hashstr, hashmethod)

Description: Obtains the digest string of a **hashstr** string based on the algorithm specified by **hashmethod**. **hashmethod** can be **sha256**, **sha384**, **sha512**, or **sm3**.

Return type: text

Length of the return value: 64 bytes if **hashmethod** is **sha256** or **sm3**; 96 bytes if **hashmethod** is **sha384**; 128 bytes if **hashmethod** is **sha512**

Example:

md5(string)

Description: Encrypts a string in MD5 mode and returns a value in hexadecimal form.

MD5 is insecure and is not recommended.

Return type: text

Example:

```
SELECT md5('ABC');
md5
-----
902fbdd2b1df0c4f70b4a5d23525e932
(1 row)
```

5.2.4 Using pgcrypto to Encrypt GaussDB(DWS) Data

GaussDB(DWS) 8.2.0 and later provides a built-in cryptographic module pgcrypto. The pgcrypto module allows database users to store certain columns of data after encryption, enhancing sensitive data security. Users without the encryption key cannot read the encrypted data stored in GaussDB(DWS).

The pgcrypto function runs inside database servers, which means that all data and passwords are transmitted in plaintext between pgcrypto and client applications. For security purposes, you are advised to use the SSL connection between the client and the GaussDB(DWS) server.

The functions in the pgcrypto module are as follows.

General Hash Functions

digest()

The digest() function can generate binary hash values by using a specified algorithm. The syntax is as follows:

digest(data text, type text) returns bytea digest(data bytea, type text) returns bytea

data indicates the original data, and type indicates the encryption algorithm (md5, sha1, sha224, sha256, sha384, sha512, or sm3). The return value of the function is a binary string.

Example:

Use the digest() function to encrypt the GaussDB(DWS) string using SHA256 for storage.

hmac()

The hmac() function can calculate the MAC value for data with a key by using a specified algorithm. The syntax is as follows:

hmac(data text, key text, type text) returns bytea hmac(data bytea, key bytea, type text) returns bytea

data indicates the original data, key indicates the encryption key, and type indicates the encryption algorithm (md5, sha1, sha224, sha256, sha384, sha512, or sm3). The return value of the function is a binary string.

Example:

Use **key123** and the SHA256 algorithm to calculate the MAC value for the string **GaussDB(DWS)**.

If both the original data and its encryption result are modified, the digest() function cannot identify the changes. The hmac() function can identify the changes as long as the key is not disclosed.

If the key is longer than the hash block, it will be hashed first, and the hash result will be used as the key.

Cryptographic Hash Functions

The crypt() and gen_salt() functions are used for password hashing. crypt() executes hashes to encrypt data, and gen_salt() generates salted hashes.

The algorithms in crypt() differ from the common MD5 and SHA1 hash algorithms in the following aspects:

- The algorithms used in crypt() are slow. This is the only way to make it difficult for brute-force attackers to crack passwords, which only contain a small amount of data.
- A random value (called salt) is used for encryption, so that users will get different ciphertexts even if they use the same passwords. This can protect passwords for cracking algorithms.
- The encryption results include algorithm types. Passwords can be encrypted using different algorithms for different users.

• Some of the algorithms are self-adaptive. They can slow down computing if it is too fast, and do not cause incompatibility issues with existing passwords.

The following table lists the algorithms supported by the crypt() function.

Table 5-7 Algorithms supported by crypt()

Algorith m	Maximu m Password Length	Adaptabi lity	Salt Bits	Standard Output Length	Description
bf	72	√	128	60	Blowfish-based 2a variation
md5	unlimited	×	48	34	MD5-based algorithm
xdes	8	√	24	20	Extended DES
des	8	×	12	13	Native UNIX algorithm

crypt()

The syntax of crypt() is as follows: crypt(password text, salt text) returns text

This function returns a hash value of the password string in crypt(3) format. The salt parameter is generated by the gen_salt() function.

For the same password, the crypt() function returns a different result each time, because the gen_salt() function generates a different salt each time. During password verification, the previously generated hash result can be used as the salt.

For example, to set a new password, run the following command:

UPDATE ... SET pswhash = crypt('new password', gen_salt('bf',10));

The hash values of the entered password and the stored password are compared.

SELECT (pswhash = crypt('entered password', pswhash)) AS pswmatch FROM ...;

If the entered password is correct, **true** is returned.

Example:

```
create table userpwd(userid int8, pwd text);
CREATE TABLE

insert into userpwd values (1, crypt('this is a pwd', gen_salt('bf',10)));
INSERT 0 1

select crypt('this is a pwd', pwd)=pwd as result from userpwd where userid =1;
result
------
t
(1 row)

select crypt('this is a wrong pwd', pwd)=pwd as result from userpwd where userid =1;
result
-------
```

```
f
(1 row)
```

gen_salt()

The gen_salt() function is used to generate random parameters for **crypt**. The syntax is as follows:

```
gen_salt(type text [, iter_count integer ]) returns text
```

This function generates a random salt string each time. The string determines the algorithm used by the crypt() function. The **type** parameter specifies a hash algorithm (**des**, **xdes**, **md5**, or **bf**) for generating a string. For the xdes and bf algorithms, **iter_count** indicates the number of iterations. A large value indicates a long encryption or cracking time.

The salt generated by an algorithm has a fixed format. For example, in \$2a \$06\$ in the bf algorithm result, 2a indicates the 2a variation of Blowfish, and 06 indicates the number of iterations.

If **iter_count** is ignored, the default number of iterations will be used. The valid **iter_count** values depend on the algorithm used, as shown in the table below. For the xdes algorithm, the number of iterations must be an odd number.

Table 5-8 Iteration co	ounts of	crvpt()
------------------------	----------	---------

Algorithm	Default Value	Min.	Max.
xdes	725	1	16777215
bf	6	4	31

PGP Encryption Functions

The PGP encryption function of GaussDB(DWS) complies with the OpenPGP (RFC 4880) standard, which includes requirements for symmetric key (private key) encryption and asymmetric key (public key) encryption.

An encrypted PGP message consists of the following parts:

- Session key (encrypted symmetric key or public key) of the message
- Data encrypted using the session key

For symmetric key (password) encryption:

- The key is encrypted using the String2Key (S2K) algorithm, which is like a slowed down crypt() algorithm with a random salt. A full-length binary key will be generated.
- 2. If a separate session key is required, a random key will be generated. If it is not required, the S2K key will be used as the session key.
- 3. If the S2K key is directly used for a session, this key will be put in the session key packet. Otherwise, the S2K key will be used to encrypt the session key, and the encryption result will be put in the session key packet.

For public key encryption:

- 1. A random session key is generated.
- 2. This random key is encrypted using the public key and then put in the session key packet.

In either case, the data encryption process is as follows:

- 1. (Optional) Compress data, convert data to UTF-8, or convert newline characters.
- 2. A block consisting of random bytes is added before the data, serving as a random initial value (IV).
- 3. A random prefix and the SHA1 hash value suffix are added to the data.
- 4. The entire content is encrypted using the session key and then placed in the data packet.

Supported PGP encryption functions

pgp_sym_encrypt()

Description: Encrypts a symmetric key.

Syntax:

pgp_sym_encrypt(data text, psw text [, options text]) returns bytea pgp_sym_encrypt_bytea(data bytea, psw text [, options text]) returns bytea

data indicates the data to be encrypted, **psw** indicates the PGP symmetric key, and **options** is used to set options. For details, see **Table 5-9**.

pqp_sym_decrypt()

Description: Decrypts a message encrypted using a PGP symmetric key.

Syntax:

```
pgp_sym_decrypt(msg bytea, psw text [, options text ]) returns text pgp_sym_decrypt_bytea(msg bytea, psw text [, options text ]) returns bytea
```

msg indicates the data to be decrypted, **psw** indicates the PGP symmetric key, and **options** is used to set options. For details, see **Table 5-9**. To avoid generating invalid characters, you are not allowed to use the pgp_sym_decrypt function to decrypt bytea data. You can use the pgp_sym_decrypt_bytea function instead.

pqp pub encrypt()

Description: Encrypts a public key.

Syntax:

```
pgp_pub_encrypt(data text, key bytea [, options text ]) returns bytea pgp_pub_encrypt_bytea(data bytea, key bytea [, options text ]) returns bytea
```

data indicates the data to be encrypted. **key** indicates the PGP public key. If a private key is used as input, an error will be returned. **options** is used to set options. For details, see **Table 5-9**.

pgp_pub_decrypt()

Description: Decrypts a message encrypted using a PGP public key.

Syntax

```
pgp_pub_decrypt(msg bytea, key bytea [, psw text [, options text ]]) returns text pgp_pub_decrypt_bytea(msg bytea, key bytea [, psw text [, options text ]]) returns bytea
```

You can decrypt a message encrypted using a public key. The **key** must be the private key corresponding to the public key used for encryption. If the private key is password protected, specify the password in **psw**. If you have not

specified any password but want to specify this option now, provide an empty password.

To avoid generating invalid characters, you are not allowed to use the pgp_pub_decrypt function to decrypt bytea data. You can use pgp_pub_decrypt_bytea function instead.

The **key** must be the private key corresponding to the public key used for encryption. If the private key is password protected, specify the password in **psw**. If you have not specified any password but want to specify this option now, provide an empty password. The options **parameter** is used to set options. For details, see **Table 5-9**.

pgp_key_id()

Description: Extracts the key ID of the PGP public or private key. If an encrypted message is used as the input, the ID of the key used to encrypt the message will be returned.

Syntax:

pgp_key_id(bytea) returns text

This function can return two special key IDs:

- SYMKEY, indicating that a message is encrypted using a symmetric key.
- ANYKEY, indicating that a message is encrypted using the public key, but the key ID has been deleted. To decrypt the message in this case, you need to try all the keys until you find the correct private key. pgcrypto does not produce such encrypted messages.

Different keys may have the same ID. This situation rarely occurs. In this case, the client application needs to try different keys for decryption, in the same way it deals with **ANYKEY**.

armor()

Description: Converts binary data into PGP ASCII-armor format by the CRC calculation and formatting of a Base64 string.

Syntax:

armor(data bytea [, keys text[], values text[]]) returns text

dearmor()

Description: Performs the reverse conversion.

Syntax:

dearmor(data text) returns bytea

Converts the encrypted data bytea to the PGP ASCII-armor format, or the other way around.

data indicates the data to be converted. If multiple pairs of keys and values are specified, an armor header will be generated for each key-value pair and added to the output. The two arrays are both one-dimensional arrays with the same length, and cannot contain non-ASCII characters.

pgp_armor_headers()

Description: Returns the armor header in the data. pgp_armor_headers(data text, key out text, value out text) returns setof record

The return result is a data row set consisting of key and value columns. Any non-ASCII characters contained in the set are regarded as UTF-8 characters.

Using GnuPG to generate PGP keys

To generate a key, run the following command:

gpg --gen-key

DSA and Elgamal keys are recommended.

To use an RSA key, you must create a DSA or RSA key as the master key used only for signature, and then specify **gpg** --edit-key to add an RSA encryption subkey.

To list keys, run the following command:

gpg --list-secret-keys

To export a public key in ASCII-protected format, run the following command: gpg -a --export KEYID > public.key

To export a private key in ASCII-protected format, run the following command:

gpg -a --export-secret-keys KEYID > secret.key

Before using these keys as the input to the PGP function, run dearmor() on them. Alternatively, if you can process binary data, remove -a from the command.

NOTICE

The PGP encryption function has the following restrictions:

- Signatures are not supported. This function does not check whether the encryption subkey belongs to the master key.
- The encryption key cannot be used as the master key. This constraint does not impose much impact, because it is rarely violated.
- Only one subkey is allowed. This may be a problem, because multiple subkeys are often required. General GPG and PGP keys cannot be used as pgcrypto encryption keys. Their usage is totally different.

PGP function parameters

The option names in the pgcrypto function are similar to those in the GnuPG function. Option values are set using equal signs (=), and the options are separated by commas (,). Example:

pgp_sym_encrypt(data, psw, 'compress-algo=1, cipher-algo=aes256')

Options other than **convert-crlf** can be used only for encryption functions. The decryption function obtains parameters from PGP data.

The most common options are **compress-algo** and **unicode-mode**. You can retain the default values for other options.

Table 5-9 pgcrypto encryption options

Option	Description	Defa ult Valu e	Value	Function
cipher- algo	Cryptographic algorithm	aes12 8	bf, aes128, aes192, aes256, 3des, cast5	pgp_sym_enc rypt, pgp_pub_enc rypt
compre ss-algo	Compression algorithm	0	 0: not compressed 1: ZIP compression 2: ZLIB compression (ZIP + Metadata + CRC) 	pgp_sym_enc rypt, pgp_pub_enc rypt
compre ss-level	Compression level. A high level indicates the compression will be slow, but the data size after compression will be small. 0 disables compression.	6	0, 1-9	pgp_sym_enc rypt, pgp_pub_enc rypt
convert -crlf	Indicates whether to convert \n to \r\n during encryption, and whether to convert \r\n to \n during decryption. RFC4880 requires that \r\n must be used as the newline character in text data storage.	0	0, 1	pgp_sym_enc rypt, pgp_pub_enc rypt, pgp_sym_dec rypt, pgp_pub_dec rypt
disable- mdc	SHA-1 is not used to protect data. It is used only for compatibility with old PGP products.	0	0, 1	pgp_sym_enc rypt, pgp_pub_enc rypt

Option	Description	Defa ult Valu e	Value	Function
sess- key	A separate session key is used. Public key encryption always uses a separate session key. This option is used for symmetric key encryption, which directly uses the S2K key by default.	0	0, 1	pgp_sym_enc rypt
s2k- mode	S2K algorithm	3	 0: Salt is not used. This setting is not recommended. 1: Salt is used, but the number of iterations is fixed. 3: Salt is used, and the number of iterations can be changed. 	pgp_sym_enc rypt
s2k- count	Number of iterations of the S2K algorithm	A rand om value betw een 65,53 6 and 253,9 52.	1024 ≤ Value ≤ 65,011,712	pgp_sym_en crypt and s2k-mode=3
s2k- digest- algo	Digest algorithm used during S2K calculation	sha1	md5, sha1	pgp_sym_enc rypt
s2k- cipher- algo	Password used to encrypt a separate session key	ciphe r- algo algori thm	bf, aes, aes128, aes192, aes256	pgp_sym_enc rypt

Option	Description	Defa ult Valu e	Value	Function
unicode -mode	Whether to convert text data between database internal encoding and UTF-8. If the database already uses UTF-8 encoding, no conversion will be performed, but the message will be marked as UTF-8. If this parameter is not specified, the message will not be marked.	0	0, 1	pgp_sym_enc rypt, pgp_pub_enc rypt

Raw Encryption Functions

Raw encryption functions only run a cipher over data. They do not support any advanced functions of PGP encryption. Therefore, the following problems exist:

- They use user key directly as cipher key.
- No integrity check is performed to check whether the encrypted data was modified.
- You need to associate all encryption parameters yourself, including IV.
- Text data cannot be processed.

With the introduction of PGP encryption, these raw encryption functions are not recommended.

```
encrypt(data bytea, key bytea, type text) returns bytea
decrypt(data bytea, key bytea, type text) returns bytea
encrypt_iv(data bytea, key bytea, iv bytea, type text) returns bytea
decrypt_iv(data bytea, key bytea, iv bytea, type text) returns bytea
```

data indicates the data to be encrypted, and **type** indicates the encryption/decryption method. The syntax of the **type** parameter is as follows:

```
algorithm [ - mode ] [ /pad: padding ]
```

The options of **algorithm** are as follows:

- **bf**: Blowfish algorithm. Synonyms: **BF**, **BF-CBC**; **BLOWFISH**, **BF-CBC**; **BLOWFISH-CBC**, **BF-CBC**; **BLOWFISH-CFB**, **BF-CFB**
- aes: AES algorithm (Rijndael-128, -192, or -256). Synonyms: AES, AES-CBC, RIJNDAEL, AES-CBC, RIJNDAEL, AES-CBC, RIJNDAEL-CBC, AES-CBC, RIJNDAEL-ECB, AES-ECB

- DES algorithm. Synonyms: DES, DES-CBC; 3DES, DES3-CBC, 3DES-ECB, DES3-ECB; 3DES-CBC, DES3-CBC
- CAST5 algorithm. Synonym: CAST5-CBC

The options of **mode** are as follows:

- **cbc**: The next block depends on the previous block. (This is the default value.)
- ecb: Each block is encrypted separately. (This value is used only for tests.)

The options of **padding** are as follows:

- **pkcs**: The data can be of any length. (This is the default value.)
- **none**: The data must be a multiple of cipher block size.

For example, the encryption results of the following functions are the same:

```
encrypt(data, 'fooz', 'bf')
encrypt(data, 'fooz', 'bf-cbc/pad:pkcs')
```

For the **encrypt_iv** and **decrypt_iv** functions, the **iv** parameter indicates the initial value for the CBC mode. This parameter is ignored for ECB. It is truncated or padded with zeroes if not exactly block size. It defaults to all zeroes in the functions without this parameter.

Random Data Functions

• The gen_random_bytes() function is used to generate cryptographically strong random bytes.

gen_random_bytes(count integer) returns bytea

count indicates the number of returned bytes. The value range is 1 to 1024.

Example:

The gen_random_uuid() function is used to return a random UUID of version
 4.

```
SELECT gen_random_uuid();
gen_random_uuid
------
2bd664a2-b760-4859-8af6-8d09ccc5b830
```

6 GaussDB(DWS) Data Query

6.1 GaussDB(DWS) Single-Table Query

Example table:

```
CREATE TABLE newproducts
(
product_id INTEGER NOT NULL,
product_name VARCHAR2(60),
category VARCHAR2(60),
quantity INTEGER
)
WITH (ORIENTATION = COLUMN) DISTRIBUTE BY HASH(product_id);

INSERT INTO newproducts VALUES (1502, 'earphones', 'electronics',150);
INSERT INTO newproducts VALUES (1601, 'telescope', 'toys',80);
INSERT INTO newproducts VALUES (1666, 'Frisbee', 'toys',244);
INSERT INTO newproducts VALUES (1700, 'interface', 'books',100);
INSERT INTO newproducts VALUES (2344, 'milklotion', 'skin care',320);
INSERT INTO newproducts VALUES (3577, 'dumbbell', 'sports',550);
INSERT INTO newproducts VALUES (1210, 'necklace', 'jewels', 200);
```

Simple Queries

Run the **SELECT... FROM...** statement to obtain the result from the database.

```
SELECT category FROM newproducts;
category
------
electr
sports
jewels
toys
books
skin care
toys
(7 rows)
```

Filtering Test Results

Run the WHERE statement to filter the query result and find the queried part.

```
SELECT * FROM newproducts WHERE category='toys';
product_id | product_name | category | quantity
```

Sorting Results

Use the **ORDER BY** statement to sort query results.

```
SELECT product_id,product_name,category,quantity FROM newproducts ORDER BY quantity DESC;
product_id | product_name | category | quantity

3577 | dumbbell | sports | 550
2344 | milklotion | skin care | 320
1666 | Frisbee | toys | 244
1210 | necklace | jewels | 200
1502 | earphones | electronics | 150
1700 | interface | books | 100
1601 | telescope | toys | 80

(7 rows)
```

Limiting the Number of Query Results

If you want the query to return only part of the result, you can use the **LIMIT** statement to limit the number of records returned in the query result.

Aggregated Query

If you want query data comprehensively, you can use the **GROUP BY** statement and aggregate functions to construct an aggregated query.

```
SELECT category, string_agg(quantity,',') FROM newproducts group by category;
category | string_agg
-------
toys | 80,244
books | 100
sports | 550
jewels | 200
skin care | 320
electronics | 150
```

6.2 GaussDB(DWS) Multi-Table Join Query

Join Types

Multiple joins are necessary for accomplishing complex queries. Joins are classified into inner joins and outer joins. Each type of joins have their subtypes.

- Inner join: inner join, cross join, and natural join.
- Outer join: left outer join, right outer join, and full join.

To better illustrate the differences between these joins, the following provides some examples.

Create the sample tables **student** and **math_score** and insert data into them. Set **enable_fast_query_shipping** to **off** (**on** by default), that is, the query optimizer uses the distributed framework. Set **explain_perf_mode** to **pretty** (default value) to specify the **EXPLAIN** display format.

```
CREATE TABLE student(
id INTEGER,
name varchar(50)
);

CREATE TABLE math_score(
id INTEGER,
score INTEGER
);

INSERT INTO student VALUES(1, 'Tom');
INSERT INTO student VALUES(2, 'Lily');
INSERT INTO student VALUES(3, 'Tina');
INSERT INTO student VALUES(4, 'Perry');

INSERT INTO math_score VALUES(1, 80);
INSERT INTO math_score VALUES(2, 75);
INSERT INTO math_score VALUES(4, 95);
INSERT INTO math_score VALUES(6, NULL);

SET enable_fast_query_shipping = off;
SET explain_perf_mode = pretty;
```

Inner Join

Inner join

Syntax:

left_table [INNER] JOIN right_table [ON join_condition | USING (join_column)]

Description: Rows that meet **join_condition** in both the left and right tables are joined and output. Tuples that do not meet **join_condition** are not output.

Example 1: Query students' math scores.

```
SELECT s.id, s.name, ms.score FROM student s JOIN math_score ms on s.id = ms.id;
id | name | score
2 | Lily | 75
 1 | Tom | 80
 4 | Perry | 95
(3 rows)
EXPLAIN SELECT s.id, s.name, ms.score FROM student s JOIN math_score ms on s.id = ms.id;
                       QUERY PLAN
id | operation | E-rows | E-memory | E-width | E-costs
               -----+-----
 1 | -> Streaming (type: GATHER) | 4 | | 13 | 19.47
2 | -> Hash Join (3,4) | 4 | 1MB | 13 | 11.47
3 | -> Seq Scan on math_score ms | 30 | 1MB | 8 | 10.10
      -> Hash | 12 | 16MB | 9 | 1.28
 4 |
          -> Streaming(type: BROADCAST) | 12 | 2MB | 9 | 1.28

-> Seq Scan on student s | 4 | 1MB | 9 | 1.01
 5 İ
 6 |
Predicate Information (identified by plan id)
 2 -- Hash Join (3,4)
      Hash Cond: (ms.id = s.id)
```

```
===== Query Summary =====

System available mem: 1761280KB

Query Max mem: 1761280KB

Query estimated mem: 4400KB

(19 rows)
```

Cross join

Syntax:

left_table CROSS JOIN right_table

Description: Each row in the left table is joined with each row in the right table. The number of final rows is the product of the number of rows on both sides. The product is also called Cartesian product.

Example 2: Cross join of student tables and math score tables.

```
SELECT s.id, s.name, ms.score FROM student s CROSS JOIN math_score ms;
id | name | score
3 | Tina | 80
 2 | Lily | 80
 1 | Tom | 80
 4 | Perry | 80
 3 | Tina |
 2 | Lily |
 1 | Tom |
 4 | Perry |
 3 | Tina | 95
 2 | Lily | 95
1 | Tom | 95
4 | Perry | 95
 2 | Lily | 75
 3 | Tina | 75
 1 | Tom | 75
4 | Perry | 75
(16 rows)
EXPLAIN SELECT s.id, s.name, ms.score FROM student s CROSS JOIN math_score ms;
                   QUERY PLAN
id | operation | E-rows | E-memory | E-width | E-costs
          ------+-----+------
 1 | -> Streaming (type: GATHER) | 120 | | 13 | 19.89
2 | -> Nested Loop (3,4) | 120 | 1MB | 13 | 11.89
3 | -> Seq Scan on math_score ms | 30 | 1MB | 4 | 10.10
4 | -> Materialize | 12 | 16MB | 9 | 1.30
       -> Streaming(type: BROADCAST) | 12 | 2MB | 9 | 1.28

-> Seq Scan on student s | 4 | 1MB | 9 | 1.01
 5 İ
 6 |
 ===== Query Summary =====
System available mem: 1761280KB
Query Max mem: 1761280KB
Query estimated mem: 4144KB
(14 rows)
```

Natural join

Syntax:

left_table NATURAL JOIN right_table

Description: Columns with the same name in left table and right table are joined by equi-join, and the columns with the same name are merged into one column

Example 3: Natural join between the **student** table and the **math_score** table. The columns with the same name in the two tables are the **id** columns, therefore equivalent join is performed based on the **id** columns.

```
SELECT * FROM student s NATURAL JOIN math_score ms;
id | name | score
1 | Tom | 80
4 | Perry | 95
 2 | Lily | 75
(3 rows)
EXPLAIN SELECT * FROM student s NATURAL JOIN math_score ms;
                      QUERY PLAN
id | operation | E-rows | E-memory | E-width | E-costs
 1 | -> Streaming (type: GATHER) | 4 | | 13 | 19.47
2 | -> Hash Join (3,4) | 4 | 1MB | 13 | 11.47
3 | -> Seq Scan on math_score ms | 30 | 1MB | 8 | 10.10
4 | -> Hash | 12 | 16MB | 9 | 1.28
        -> Streaming(type: BROADCAST) | 12 | 2MB | 9 | 1.28

-> Seq Scan on student s | 4 | 1MB | 9 | 1.01
 5 |
  6 |
Predicate Information (identified by plan id)
 2 -- Hash Join (3.4)
      Hash Cond: (ms.id = s.id)
 ===== Query Summary =====
System available mem: 1761280KB
Query Max mem: 1761280KB
Query estimated mem: 4400KB
(19 rows)
```

Outer Join

• Left Join

Syntax:

left_table LEFT [OUTER] JOIN right_table [ON join_condition | USING (join_column)]

Description: The result set of a left outer join includes all rows of left table, not only the joined rows. If a row in the left table does not match any row in right table, the row will be **NULL** in the result set.

Example 4: Perform left join on the **student** table and **math_score** table. The right table data corresponding to the row where ID is 3 in the **student** table is filled with **NULL** in the result set.

• Right join

Syntax:

left_table RIGHT [OUTER] JOIN right_table [ON join_condition | USING (join_column)]

Description: Contrary to the left join, the result set of a right join includes all rows of the right table, not just the joined rows. If a row in the right table does not match any row in right table, the row will be **NULL** in the result set.

Example 5: Perform right join on the **student** table and **math_score** table. The right table data corresponding to the row where ID is 6 in the **math_score** table is filled with **NULL** in the result set.

```
SELECT ms.id, s.name, ms.score FROM student s RIGHT JOIN math_score ms on (s.id = ms.id);
id | name | score
 1 | Tom | 80
 61
 4 | Perry | 95
 2 | Lily | 75
EXPLAIN SELECT ms.id, s.name, ms.score FROM student s RIGHT JOIN math_score ms on (s.id = ms.id);
                        QUERY PLAN
            operation | E-rows | E-memory | E-width | E-costs
id |
 1 | -> Streaming (type: GATHER) | 30 | | 13 | 19.47
2 | -> Hash Left Join (3, 4) | 30 | 1MB | 13 | 11.47
3 | -> Seq Scan on math_score ms | 30 | 1MB | 8 | 10.10
                                | 12 | 16MB | 9 | 1.28
 4
       -> Hash
          -> Streaming(type: BROADCAST) | 12 | 2MB | 9 | 1.28
 5 |
            -> Seq Scan on student s | 4 | 1MB | 9 | 1.01
Predicate Information (identified by plan id)
 2 -- Hash Left Join (3, 4)
     Hash Cond: (ms.id = s.id)
 ===== Query Summary =====
System available mem: 1761280KB
Query Max mem: 1761280KB
Query estimated mem: 5424KB
```

In a right join, **Left** is displayed in the join operator. This is because a right join is actually the process replacing the left table with the right table then performing left join.

Full join

Syntax:

left_table FULL [OUTER] JOIN right_table [ON join_condition | USING (join_column)]

Description: A full join is a combination of a left outer join and a right outer join. The result set of a full outer join includes all rows of the left table and the right table, not just the joined rows. If a row in the left table does not

match any row in the right table, the row will be **NULL** in the result set. If a row in the right table does not match any row in right table, the row will be **NULL** in the result set.

Example 6: Perform full outer join on the **student** table and **math_score** table. The right table data corresponding to the row where ID is 3 is filled with **NULL** in the result set. The left table data corresponding to the row where ID is 6 is filled with **NULL** in the result set.

```
SELECT s.id, s.name, ms.id, ms.score FROM student s FULL JOIN math_score ms ON (s.id = ms.id);
id | name | id | score
2 | Lily | 2 | 75
 4 | Perry | 4 | 95
 1 | Tom | 1 | 80
 3 | Tina | |
  | |6|
(5 rows)
EXPLAIN SELECT s.id, s.name, ms.id, ms.score FROM student s FULL JOIN math_score ms ON (s.id =
                           QUERY PLAN
               operation | E-rows | E-memory | E-width | E-costs
id I
 1 | -> Streaming (type: GATHER) | 30 | 17 | 20.24
2 | -> Hash Full Join (3, 5) | 30 | 1MB | 17 | 12.24
                                                            | 17 | 20.24
        -> Streaming(type: REDISTRIBUTE) | 30 | 2MB | 8 | 11.06

-> Seq Scan on math_score ms | 30 | 1MB | 8 | 10.10
          -> Seq Scan on math_score ms |
  5 |
                                     | 4 | 16MB | 9 | 1.11
        -> Streaming(type: REDISTRIBUTE) | 4 | 2MB | 9 | 1.11
-> Seq Scan on student s | 4 | 1MB | 9 | 1.01
 6 j
 7 |
Predicate Information (identified by plan id)
 2 -- Hash Full Join (3.5)
      Hash Cond: (ms.id = s.id)
 ===== Query Summary =====
System available mem: 1761280KB
Query Max mem: 1761280KB
Query estimated mem: 6496KB
(20 rows)
```

Differences Between the ON Condition and the WHERE Condition in Multi-Table Query

According to the preceding join syntax, except natural join and cross join, the **ON** condition (**USING** is converted to the **ON** condition during query parsing) is used on the join result of both the two tables. Generally, the **WHERE** condition is used in the query statement to restrict the query result. The **ON** join condition and **WHERE** filter condition do not contain conditions that can be pushed down to tables. The differences between **ON** and **WHERE** are as follows:

- The **ON** condition is used for joining two tables.
- WHERE is used to filter the result set.

To sum up, the **ON** condition is used when two tables are joined. After the join result set of two tables is generated, the **WHERE** condition is used.

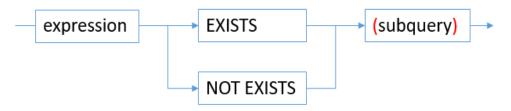
6.3 GaussDB(DWS) Subquery Expressions

A subquery allows you to nest one query within another, enabling more complex data query and analysis.

Subquery Expressions

EXISTS/NOT EXISTS

Before the main query runs, the subquery runs and its result determines if the main query continues. EXISTS returns **true** if the subquery returns at least one row. **NOT EXISTS** returns **true** if the subquery returns no rows.

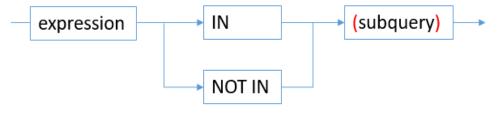


Syntax:

WHERE column_name EXISTS/NOT EXISTS (subquery)

IN/NOT IN

IN and NOT IN are operators that check if a value is in a set of values. **IN** returns **true** when the outer query row matches a subquery row. **NOT IN** returns **true** when the outer query row does not match any subquery row.



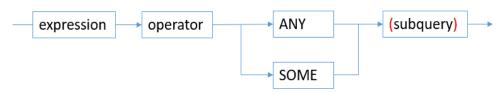
Syntax:

WHERE column_name IN/NOT IN (subquery)

ANY/SOME

ANY indicates that any value in a subquery can match a value in an outer query. **SOME** is the same as **ANY**, but the syntax is different.

The subquery can return only one column. The expression on the left uses operators (=, <>, <, <=, >, >=) to compare the value with each subquery row. The result must be a Boolean value. The result of **ANY** is **true** if any true result is obtained. The result is **false** if no true result is found (including the case where the subquery returns no rows).



Syntax:

WHERE column_name operator ANY/SOME (subquery)

ALL

The subquery on the right must return only one field. The expression on the left uses operators (=, <>, <, <=, >, >=) to compare the value with each subquery row. The result must be a Boolean value. The result of **ALL** is **true** if all rows yield true (including the case where the subquery returns no rows). The result is **false** if any false result is found.



Syntax:

WHERE column_name operator ALL (subquery)

Table 6-1 ALL conditions

Condition	Description
column_name > ALL()	The column_name value must be greater than the maximum value of a set to be true.
column_name >= ALL()	The column_name value must be greater than or equal to the maximum value of a set to be true.
column_name < ALL()	The column_name value must be smaller than the minimum value of a set to be true.
column_name <= ALL()	The column_name value must be smaller than or equal to the minimum value of a set to be true.
column_name <> ALL()	The column_name value cannot be equal to any value in a set to be true.
column_name = ALL()	The column_name value must be equal to any value in a set to be true.

Example

Create the **course** table and insert data into the table.

CREATE TABLE course(cid VARCHAR(10) COMMENT 'No.course',cname VARCHAR(10) COMMENT 'course name',teid VARCHAR(10) COMMENT 'No.teacher');

```
INSERT INTO course VALUES('01' , 'course1' , '02');
INSERT INTO course VALUES('02' , 'course2' , '01');
INSERT INTO course VALUES('03' , 'course3' , '03');
```

Create the **teacher** table and insert data into the table.

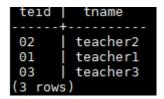
CREATE TABLE teacher(teid VARCHAR(10) COMMENT 'Teacher ID',tname VARCHAR(10) COMMENT'Teacher name');

```
INSERT INTO teacher VALUES('01', 'teacher1');
INSERT INTO teacher VALUES('02', 'teacher2');
INSERT INTO teacher VALUES('03', 'teacher3');
INSERT INTO teacher VALUES('04', 'teacher4');
```

EXISTS/NOT EXISTS example

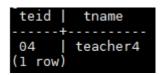
Query the teacher records in the course table.

SELECT * FROM teacher WHERE EXISTS (SELECT * FROM course WHERE course.teid = teacher.teid);



Query the teacher records that are not in the **course** table.

SELECT * FROM teacher WHERE NOT EXISTS (SELECT * FROM course WHERE course.teid = teacher.teid);



• IN/NOT IN example

Query the course table for teacher information based on the teacher ID.

SELECT * FROM course WHERE teid IN (SELECT teid FROM teacher);

Query the information about teachers who are not in the course table.

SELECT * FROM teacher WHERE teid NOT IN (SELECT teid FROM course);

```
teid | tname
-----+-----
04 | teacher4
(1 row)
```

ANY/SOME example

Compare the main query fields on the left with the subquery fields on the right to obtain the required result set.

SELECT * FROM course WHERE teid < ANY (SELECT teid FROM teacher where teid<>'04');

or

SELECT * FROM course WHERE teid < some (SELECT teid FROM teacher where teid<>'04');

ALL example

The value in the **teid** column must be smaller than the minimum value in the set to be true.

SELECT * FROM course WHERE teid < ALL(SELECT teid FROM teacher WHERE teid<>'01');



Important Notes

- Duplicate subquery statements are not allowed in an SQL statement.
- Avoid scalar sub-queries whenever possible. A scalar subquery is a subquery whose result is one value and whose condition expression uses an equal operator.
- Do not use subqueries in the SELECT target columns. Otherwise, the plan cannot be pushed down, affecting the execution performance.
- It is recommended that the nested subqueries cannot exceed two layers. Subqueries cause temporary table overhead. Therefore, complex queries must be optimized based on service logic.

A subquery can be nested in the SELECT statement to implement a more complex query. A subquery can also use the results of other queries in the WHERE clause to better filter data. However, subqueries may cause query performance problems and make code difficult to read and understand. Therefore, when using SQL subqueries in databases such as GaussDB, use them based on the site requirements.

6.4 GaussDB(DWS) WITH Expressions

The WITH expression is used to define auxiliary statements used in large queries. These auxiliary statements are usually called common table expressions (CTE), which can be understood as a named subquery. The subquery can be referenced multiple times by its name in the quey.

An auxiliary statement may use **SELECT**, **INSERT**, **UPDATE**, or **DELETE**. The **WITH** clause can be attached to a main statement, which can be a **SELECT**, **INSERT**, or **DELETE** statement.

SELECT in WITH

This section describes the usage of **SELECT** in a **WITH** clause.

Syntax

```
[WITH [RECURSIVE] with_query [, ...] ] SELECT ...
```

The syntax of with_query is as follows:

```
with_query_name [ ( column_name [, ...] ) ]
AS [ [ NOT ] MATERIALIZED ] ( {select | values | insert | update | delete} )
```

CAUTION

- If you use **MATERIALIZED**, the subquery runs once and its result set is saved. If you use **NOT MATERIALIZED**, the subquery is replaced with its reference in the main query.
- The SQL statement specified by the AS statement of a CTE must be a statement that can return query results. It can be a common SELECT query statement or other data modification statements such as INSERT, UPDATE, DELETE, and VALUES. When using a data modification statement, you need to use the RETURNING clause to return tuples. Example:
 WITH s AS (INSERT INTO t VALUES(1) RETURNING a) SELECT * FROM s;
- A WITH expression indicates the CTE definition in a SQL statement block.
 Multiple CTEs can be defined at the same time. You can specify column names for each CTE or use the aliases of the columns in the query output. Example: WITH s1(a, b) AS (SELECT x, y FROM t1), s2 AS (SELECT x, y FROM t2) SELECT * FROM s1 JOIN s2 ON s1 a=s2 x:

This statement defines two CTEs: **s1** and **s2**. **s1** specifies the column names **a** and **b**, and **s2** does not specify the column names. Therefore, the column names are the output column names **x** and **y**.

- Each CTE can be referenced zero, one, or more times in the main query.
- CTEs with the same name cannot exist in the same statement block. If CTEs with the same name exist in different statement blocks, the CTE in the nearest statement block is referenced.
- An SQL statement may contain multiple SQL statement blocks. Each statement block can contain a WITH expression. The CTE in each WITH expression can be referenced in the current statement block, subsequent CTEs of the current statement block, and sub-layer statement blocks, however, it cannot be referenced in the parent statement block. The definition of each CTE is also a statement block. Therefore, a WITH expression can also be defined in the statement block.

The purpose of SELECT in WITH is to break down complex queries into simple parts. Example:

```
WITH regional_sales AS (
    SELECT region, SUM(amount) AS total_sales
    FROM orders
    GROUP BY region
), top_regions AS (
    SELECT region
    FROM regional_sales
    WHERE total_sales > (SELECT SUM(total_sales)/10 FROM regional_sales)
)

SELECT region,
    product,
    SUM(quantity) AS product_units,
    SUM(amount) AS product_sales
FROM orders
WHERE region IN (SELECT region FROM top_regions)
GROUP BY region, product;
```

The WITH clause defines two auxiliary statements: regional_sales and top_regions. The output of regional_sales is used in top_regions, and the output of top_regions is used in the main SELECT query. This example can be written without WITH. In that case, it must be written with a two-layer nested sub-SELECT statement, making the query longer and difficult to maintain.

Recursive WITH Query

By declaring the keyword **RECURSIVE**, a WITH query can reference its own output.

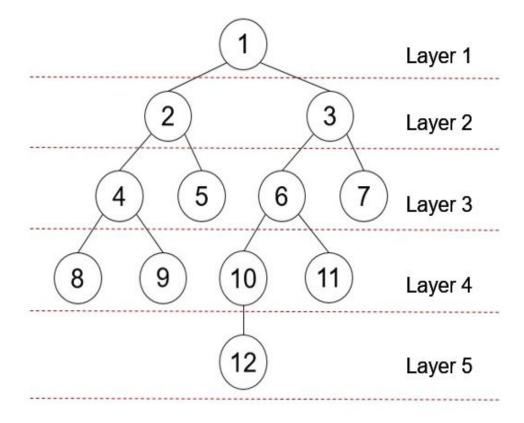
The common form of a recursive WITH query is as follows:

non_recursive_term UNION [ALL] recursive_term

UNION performs deduplication when combining sets. **UNION ALL** directly combines result sets without deduplication. Only recursive items can contain references to the query output.

When using recursive WITH, ensure that the recursive item of the query does not return a tuple. Otherwise, the query will loop infinitely.

The table **tree** is used to store information about all nodes in the following figure.



The table definition statement is as follows:

CREATE TABLE tree(id INT, parentid INT);

The data in the table is as follows:

INSERT INTO tree VALUES(1,0),(2,1),(3,1),(4,2),(5,2),(6,3),(7,3),(8,4),(9,4),(10,6),(11,6),(12,10);

SELECT * FROM tree; id | parentid

+	
1	0
2	1
3	1
4	2
5	2
6	3

```
7 | 3
8 | 4
9 | 4
10 | 6
11 | 6
12 | 10
(12 rows)
```

You can run the following **WITH RECURSIVE** statement to return the nodes and hierarchy information of the entire tree starting from node 1 at the top layer:

```
WITH RECURSIVE nodeset AS

(
-- recursive initializing query

SELECT id, parentid, 1 AS level FROM tree

WHERE id = 1

UNION ALL
-- recursive join query

SELECT tree.id, tree.parentid, level + 1 FROM tree, nodeset

WHERE tree.parentid = nodeset.id
)

SELECT * FROM nodeset ORDER BY id;
```

In the preceding query, a typical **WITH RECURSIVE** expression contains the CTE of at least one recursive query. The CTE is defined as a **UNION ALL** set operation. The first branch is the recursive start query, and the second branch is the recursive join query, the first part is referenced for continuous recursive join. When this statement is executed, the recursive start query is executed once, and the join query is executed several times. The results are added to the start query result set until the results of some join queries are empty.

The command output is as follows:

```
id | parentid | level
               1
 2 |
         1 |
               2
              2
 3 |
         1 |
 4
              3
         2 |
              3
 5
         2 |
 6
         3 |
              3
             3
 7
         3 I
 8 |
         4 |
              4
 91
         4 |
101
         6 I
              4
         6
               4
11
         10 |
12 I
(12 rows)
```

According to the returned result, the start query result contains the result set whose level is 1. The join query is executed for five times. The result sets whose levels are 2, 3, 4, and 5 are output for the first four times. During the fifth execution, there is no record whose parentid is the same as the output result set ID, that is, there is no redundant child node. Therefore, the query ends.

□ NOTE

GaussDB(DWS) supports distributed execution of **WITH RECURSIVE** expressions. **WITH RECURSIVE** involves cyclic calculation. Therefore, GaussDB(DWS) introduces the **max_recursive_times** parameter to control the maximum number of cycles of WITH RECURSIVE. The default value is **200**. If the number of cycles exceeds **200**, an error is reported.

Data Modification Statements in WITH

Use the **INSERT**, **UPDATE**, and **DELETE** commands in the WITH clause. This allows the user to perform multiple different operations in the same query. The following is an example:

```
WITH moved_tree AS (
DELETE FROM tree
WHERE parentid = 4
RETURNING * )
INSERT INTO tree_log
SELECT * FROM moved_tree;
```

The preceding query example actually moves rows from **tree** to **tree_log**. The **DELETE** command in the **WITH** clause deletes the specified rows from **tree**, returns their contents through the **RETURNING** clause, and then the main query reads the output and inserts it into **tree log**.

To retrieve the modified content instead of the target table, the data modification statement in the **WITH** clause should include the **RETURNING** clause. This clause creates a temporary table that can be accessed by the rest of the query. If a data modification statement in the **WITH** statement lacks a **RETURNING** clause, it cannot form a temporary table and cannot be referenced in the remaining queries.

If the **RECURSIVE** keyword is specified, recursive self-reference is not allowed in data modification statements. In some cases, you can bypass this restriction by referencing the output of recursive the **WITH** statement. For example:

```
WITH RECURSIVE included_parts(sub_part, part) AS (
    SELECT sub_part, part FROM parts WHERE part = 'our_product'

UNION ALL

SELECT p.sub_part, p.part

FROM included_parts pr, parts p

WHERE p.part = pr.sub_part
)

DELETE FROM parts

WHERE part IN (SELECT part FROM included_parts);
```

This query will remove all direct or indirect subparts of a product.

The substatements in the **WITH** clause are executed at the same time as the main query. Therefore, when using the data modification statement in a WITH statement, the actual update order is in an unpredictable manner. All statements are executed in the same snapshot, and the effect of the statements is invisible on the target table. This mitigates the unpredictability of the actual order of row updates and means that **RETURNING** data is the only way to convey changes between different **WITH** substatements and the main query.

In this example, the outer layer **SELECT** can return the data before the update.

```
WITH t AS (
    UPDATE tree SET id = id + 1
    RETURNING *)
SELECT * FROM tree;
```

In this example, the external SELECT returns the updated data.

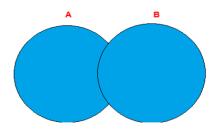
```
WITH t AS (
UPDATE tree SET id = id + 1
RETURNING *)
SELECT * FROM t;
```

The same row cannot be updated twice in a single statement. Otherwise, the update effect will be unpredictable. If only one update takes effect, it is difficult (and sometimes impossible) to predict which one takes effect.

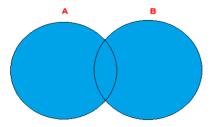
6.5 Usage of GaussDB(DWS) UNION

UNION is a powerful SQL operator that combines the result sets of two or more SELECT statements into one. During combination, the number of columns and data types in the two tables must be the same and correspond to each other. Use the **UNION** or **UNION** ALL keyword between SELECT statements.

UNION removes duplicate rows, while UNION ALL keeps them. Deduplication is time-consuming, so UNION ALL can be faster than UNION if the data sets are already distinct by logic.



The UNION operator combines the results of two queries and removes any duplicates.



The UNION ALL operator combines the results of two queries and keep all the duplicates.

Syntax

SELECT column,... FROM table1 UNION [ALL]SELECT column,... FROM table2

Example

Step 1 Create the student information table **student** (ID, name, gender, and school).

SET current_schema=public; DROP TABLE IF EXISTS student; CREATE table student(sId VARCHAR(10) NOT NULL, sname VARCHAR(10) NOT NULL, sgender VARCHAR(10) NOT NULL, sschool VARCHAR(10) NOT NULL);

Step 2 Insert data into the **student** table.

```
INSERT INTO student VALUES('s01', 'ZhaoLei', 'male', 'NENU');
INSERT INTO student VALUES('s02', 'QianDian', 'male', 'SJTU');
INSERT INTO student VALUES('s03', 'SunFenng', 'male', 'Tongji');
INSERT INTO student VALUES('s04', 'LIYun', 'male', 'CCOM');
INSERT INTO student VALUES('s05', 'ZhouMei', 'female', 'FuDan');
INSERT INTO student VALUES('s06', 'WuLan', 'female', 'WHU');
INSERT INTO student VALUES('s07', 'ZhengZhu', 'female', 'NWAFU');
INSERT INTO student VALUES('s08', 'ZhangShan', 'female', 'Tongji');
```

Step 3 View the student table.

SELECT * FROM student;

Information similar to the following is displayed.

```
sid
         sname
                   sgender
                              sschool
s01
                    male
       ZhaoLei
                              NENU
       LIYun
s04
                    male
                              CCOM
                    female
s07
       ZhengZhu
                              NWAFU
s02
       QianDian
                    male
                              SJTU
s 0 5
       ZhouMei
                    female
                              FuDan
                    female
s08
       ZhangShan
                              Tongji
s03
       SunFenng
                    male
                              Tongji
                              WHU
s06
                    female
       WuLan
(8 rows)
```

Step 4 Create the teacher information table **teacher** (ID, name, gender, and school).

```
DROP TABLE IF EXISTS teacher;
CREATE table teacher(
tid VARCHAR(10) NOT NULL,
tname VARCHAR(10) NOT NULL,
tgender VARCHAR(10) NOT NULL,
tschool VARCHAR(10) NOT NULL);
```

Step 5 Insert data to the **teacher** table.

```
INSERT INTO teacher VALUES('t01' , 'ZhangLei', 'male', 'FuDan');
INSERT INTO teacher VALUES('t02' , 'LiLiang', 'male', 'WHU');
INSERT INTO teacher VALUES('t03' , 'WangGang', 'male', 'Tongji');
```

Step 6 Query the **teacher** table.

SELECT * FROM teacher;

```
tname
                 tgender
                          tschool
                  male
t03
      WangGang
                          Tongji
                  male
      LiLiang
t02
                          WHU
      ZhangLei
                  male
t01 |
                          FuDan
 rows)
```

Step 7 Use **UNION** (combine and deduplicate) to obtain the schools of students and teachers and sort the schools in ascending order by initial letter of the school name.

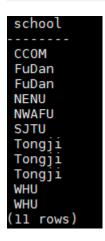
```
SELECT t.school FROM (
SELECT sschool AS school
FROM student
UNION
SELECT tschool AS school
FROM teacher
) t
ORDER BY t.school ASC;
```

Information similar to the following is displayed.

```
school
CCOM
FUDAN
NENU
NWAFU
SJTU
Tongji
WHU
(7 rows)
```

Step 8 Use **UNION ALL** (combine without deduplication) to obtain the schools of all students and teachers and sort the schools by initial letter of the school name in ascending order.

```
SELECT t.school FROM (
SELECT sschool AS school
FROM student
UNION ALL
SELECT tschool AS school
FROM teacher
) t
ORDER BY t.school ASC;
```



Step 9 Use **UNION ALL** (combine the result sets of SQL statements with **WHERE** clause) to get all information about students and teachers from "Tongji' and sort by student and teacher number in ascending order.

```
SELECT t.* FROM (
SELECT Sid AS id,Sname AS name,Sgender AS gender,Sschool AS school
FROM student
WHERE Sschool='Tongji'
UNION ALL
SELECT Tid AS id,Tname AS name,Tgender AS gender,Tschool AS school
FROM teacher
WHERE Tschool='Tongji'
) t
ORDER BY t.id ASC;
```

----End

Summary

In actual service scenarios, pay attention to the following points when using **UNION** and **UNION** ALL:

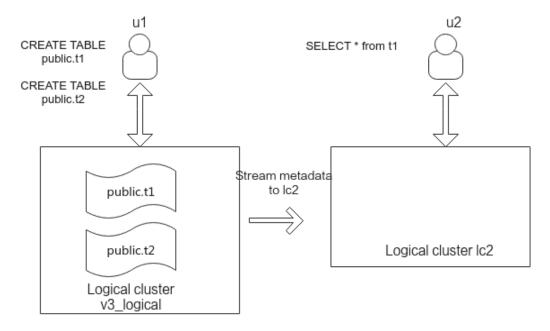
- The number of SQL fields and field types on the left and right sides must be the same.
- Check whether data deduplication (deduplication before combination or during combination) is needed based on service requirements.
- Based on the data volume, valuate the SQL execution efficiency and determine whether to use temporary tables.
- Select UNION or UNION ALL wisely and consider the complexity when writing SQL statements.

6.6 Data Reading/Writing Across Logical Clusters

Scenario

After an associated logical cluster user is created, the query or modification (including Insert, Delete, and Update) submitted by the user is calculated and executed in the associated logical cluster. If the user submits a query or modification request to a table in a different logical cluster, the optimizer generates a cross-logical cluster query or modification plan to enable the user to query or modify the table.

Figure 6-1 Querying data across logical clusters



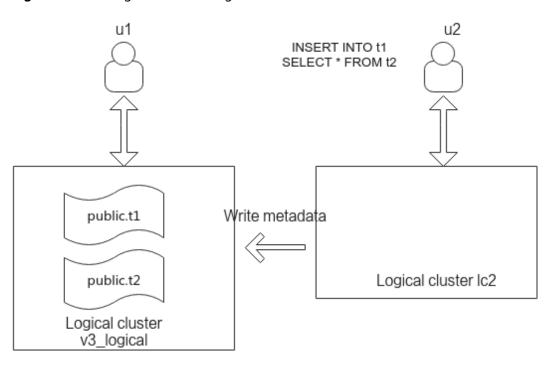


Figure 6-2 Writing data across logical clusters

Procedure

- Step 1 Create a cluster by referring to Creating a DWS 3.0 Cluster with Decoupled Storage and Compute. After a cluster is created, it is converted to a logical cluster v3 logical by default.
- **Step 2** Add three nodes to the elastic cluster, and then add the logical cluster **lc2**.
- Step 3 Create user u1 and associate it with logical cluster v3_logical.

 CREATE USER u1 with SYSADMIN NODE GROUP "v3_logical" password "Password@123";
- Step 4 Create user u2 and associate it with logical cluster lc2.

 CREATE USER u2 with SYSADMIN NODE GROUP "lc2" password "Password@123";
- **Step 5** Log in to the database as user **u1**, create tables **t1** and **t2**, and insert test data into the tables.

```
CREATE TABLE public.t1
(
id integer not null,
data integer,
age integer
)
WITH (ORIENTATION = COLUMN, COLVERSION = 3.0)
DISTRIBUTE BY ROUNDROBIN;

CREATE TABLE public.t2
(
id integer not null,
data integer,
age integer
)
WITH (ORIENTATION = COLUMN, COLVERSION = 3.0)
DISTRIBUTE BY ROUNDROBIN;

INSERT INTO public.t1 VALUES (1,2,10),(2,3,11);
INSERT INTO public.t2 VALUES (1,2,10),(2,3,11);
```

Step 6 Log in to the database as user **u2** and run the commands below to query **t1** and write data.

According to the result, user **u2** can query and write data across logical clusters. SELECT * FROM t1; INSERT INTO t1 SELECT * FROM t2;

----End

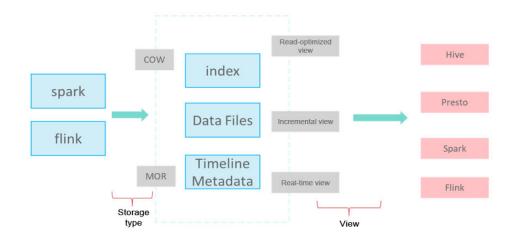
6.7 SQL on Hudi

This feature is supported only by 8.2.1.100 and later versions.

6.7.1 Introduction to Hudi

Apache Hudi indicates Hadoop Upserts Deletes and Incrementals. It is used to manage large analytical datasets stored in the distributed file system (DFS) of the Hadoop big data system.

Hudi is not just a data format. It is also a set of data access methods (similar to the access layer of GaussDB(DWS) storage). In Apache Hudi 0.9, big data components such as Spark and Flink have their own clients. The following figure shows the logical storage of Hudi.



Write Mode

COW: copy-on-write, applicable to scenarios with few updates.

MOR: replication on read. For UPDATE & DELETE, delta log files are written incrementally. During analysis, base and delta log files are compacted asynchronously.

Storage Format

index: index of the primary key. The default value is bloomfilter at the file group level.

data files: base file + delta log file (for updating and deleting base files) **timeline metadata**: manages version logs.

Views

Read-optimized view: reads the base file generated after compaction. The reading of data that is not compacted has some latency (efficient read).

Real-time view: reads the latest data. The base file and delta file are combined during the read (frequent updates).

Incremental view: reads the incremental data written to Hudi, similar to CDC (stream and batch integration).

6.7.2 Preparations Before Using Hudi

Prerequisites

You have created an OBS agency and OBS data source. For details, see **Managing OBS Data Sources**.

Authorizing the Use of OBS Data Sources

Run the **GRANT** command to grant a user the permission to use OBS data sources.

GRANT USAGE ON FOREIGN SERVER server_name TO role_name;

Example:

Run the following command to grant user **sbi_fnd** the permission to access data source **obs hudi**:

GRANT USAGE ON FOREIGN SERVER obs_hudi TO sbi_fnd;

Granting Permissions for Using Foreign Tables

Run the following command to grant a user the permission to use foreign tables:

ALTER USER role_name USEFT;

Example:

Run the following command to grant the foreign table access permission to user **sbi_fnd**:

ALTER USER sbi_fnd USEFT;

6.7.3 Hudi User Interfaces

Querying Real-Time Views and Incremental Views

GaussDB(DWS) provides table-level parameters similar to spark-sql to support real-time and incremental views.

The parameters are described as follows. Replace **SCHEMA.FOREIGN_TABLE** with the actual schema name and foreign table name.

Parameter Value Description hoodie.SCHEMA.FOREIGN_TABLE **SNAPSHOT** Queries the real-time view. .consume.mode **INCREMENTAL** Oueries the incremental view. hoodie.SCHEMA.FOREIGN_TABLE hudi Specifies the start commit .consume. start.timestamp timestamp of incremental synchronization. hoodie.SCHEMA.FOREIGN TABLE hudi Specifies the end commit .consume. ending.timestamp timestamp of incremental synchronization. If this parameter is not specified, the latest commit is used.

Table 6-2 Parameters for querying real-time views and incremental views

Ⅲ NOTE

- The preceding parameters can be set by running the set command and are valid only in the current session. You can run the reset command to restore the default values.
- You can use the system function **pg_catalog.pg_show_custom_settings()** to query the parameter setting details.
- When querying the incremental view of the MOR table, you need to use the WHERE condition to filter the _hoodie_commit_time column to prevent the log file data that is not compacted from being read. This operation is not required for the COW table.

Querying Hudi Foreign Table and Automatically Synchronizing Tasks

GaussDB(DWS) provides a series of system functions to obtain Hudi foreign table information and create Hudi automatic synchronization tasks. The automatic Hudi synchronization task periodically synchronizes data from Hudi foreign tables to GaussDB(DWS) internal tables.

Table 6-3 Hudi system functions

No.	Function	Туре	Functionality
1	pg_show_custom_settings()	Built-in function s	Queries details about the parameter settings of an HUDI foreign table.
2	hudi_get_options(regclass)	Built-in function s	Queries the attributes of an HUDI foreign table (hoodie.properties).
3	hudi_get_max_commit(regcla ss)	Built-in function s	Obtains the latest commit timestamp of the current HUDI foreign table.

No.	Function	Туре	Functionality
4	hudi_sync_task_submit(regcla ss, regclass)	Built-in function	Submits the HUDI automatic synchronization
	hudi_sync_task_submit(regcla ss, regclass, text, text)	S	task.
5	hudi_show_sync_state()	Built-in function s	Obtains the synchronization status of the HUDI automatic synchronization task.
6	hudi_sync(regclass, regclass)	Stored procedur e	Specifies the entry for invoking the HUDI automatic synchronization task.
7	hudi_sync_custom(regclass, regclass, text)	Stored procedur e	Specifies the entry for invoking the HUDI automatic synchronization task. Users can define the mapping between fields in the target table and data source table.
8	hudi_set_sync_commit(regclas s, regclass, text)	Built-in function s	Sets the start timestamp of the first synchronization of the HUDI automatic synchronization task to prevent resynchronization.
	hudi_set_sync_commit(text, text)		Sets the start timestamp of the next synchronization of a HUDI automatic synchronization task. You can use it to sync historical data again or to skip some data.

6.7.4 Creating a Hudi Data Description (Foreign Table)

A foreign table maps data on OBS. GaussDB(DWS) accesses Hudi data on OBS through foreign tables. For details, see section **CREATE FOREIGN TABLE (SQL on OBS or Hadoop)**.

Compared with OBS foreign tables, you only need to set **format** to **hudi** for Hudi foreign tables. For Hudi bucket tables, you need to set **distribute by** to **hash(bk_col1,bk_col2...)**. Only 9.1.0.100 and later versions support Hudi bucket tables.

Obtaining the Definitions of Tables on MRS.

Hudi foreign tables on GaussDB(DWS) are read-only. Before creating a foreign table, you need to specify the number of fields defined in the target data and the type of each field. A Hudi foreign table supports a maximum of 5000 columns.

For example, for a Hudi table on MRS, you can use spark-sql to query the original table definitions:

SHOW create table rtd_mfdt_int_currency_t;

Compiling GaussDB(DWS) Table Definitions

Non-bucket table

Copy the definitions of all columns in the MRS table, perform proper type conversion to adapt to the GaussDB(DWS) syntax, and create an OBS foreign table

```
CREATE FOREIGN TABLE rtd_mfdt_int_currency_ft(
_hoodie_commit_time text,
_hoodie_commit_seqno text,
_hoodie_record_key text,
_hoodie_partition_path text,
_hoodie_file_name text,
...
)SERVER obs_server OPTIONS (
foldername '/erpgc-obs-test-01/s000/sbi_fnd/rtd_mfdt_int_currency_t/',
format 'hudi',
encoding 'utf-8'
)distribute by roundrobin;
```

foldername indicates the storage path of the Hudi data on OBS, which corresponds to **LOCATION** in the Spark-sql table definitions of MRS. The path must end with a slash (/).

Bucket table

Copy the definitions of all columns in the MRS table, perform proper type conversion to adapt to the GaussDB(DWS) syntax, create an OBS foreign table, and specify the hash distribution mode.

```
CREATE FOREIGN TABLE rtd_mfdt_int_currency_ft(
_hoodie_commit_time text,
_hoodie_commit_seqno text,
_hoodie_record_key text,
_hoodie_partition_path text,
_hoodie_file_name text,
...
)SERVER obs_server OPTIONS (
foldername '/erpgc-obs-test-01/s000/sbi_fnd/rtd_mfdt_int_currency_t/',
format 'hudi',
encoding 'utf-8'
)distribute by hash(bk_col1,bk_col2...);
```

foldername indicates the storage path of the Hudi data on OBS, which corresponds to **LOCATION** in the Spark-sql table definitions of MRS. The path must end with a slash (/).

distribute by indicates the distribution column of the bucket table. The value must be the same as that of **hoodie.bucket.index.hash.field** in the **foldername/.hoodie/hoodie.index.properties** file.

6.7.5 Synchronizing Hudi Tasks

Creating a Hudi Task

Migration

If data has been imported to the GaussDB(DWS) table using CDL, use SQL on Hudi to migrate data. Alternatively, use CDM to perform full initialization and then use SQL on Hudi to synchronize incremental data.

Step 1 To create the **hudi.hudi_sync_state** synchronization status table, you must have the administrator permission.

SELECT pg_catalog.create_hudi_sync_table();

Generally, hudi.hudi_sync_state is created only once in each database.

Step 2 To set the CDL synchronization progress, you must have the INSERT and UPDATE permissions on the target table and the SELECT permission on the HUDI foreign table. Otherwise, the synchronization progress cannot be set.

SELECT hudi set sync commit('SCHEMA.TABLE', 'SCHEMA.FOREIGN TABLE', 'LATEST COMMIT');

Where:

- **SCHEMA.TABLE** indicates the name and schema of the target table for data synchronization.
- SCHEMA.FOREIGN_TABLE indicates the name and schema of the OBS foreign table.
- LATEST COMMIT indicates the end time of the Hudi synchronization.

Example: Data has been synchronized to the target table **public.in_rel** from hudi by **20220913152131**. Use SQL on Hudi to continue to export data from the OBS foreign table **hudi_read1**.

SELECT hudi_set_sync_commit('public.in_rel', 'public.hudi_read1', '20220913152131');

Step 3 Submit the Hudi synchronization task.

SELECT hudi_sync_task_submit('SCHEMA.TABLE', 'SCHEMA.FOREIGN_TABLE');

Example: Use SQL on Hudi to continue to export data from the OBS foreign table **hudi_read1** to the target table **public.in_rel**.

SELECT hudi_sync_task_submit('public.in_rel', 'public.hudi_read1');

----End

Creation

If the GaussDB(DWS) table is empty and data is synchronized from Hudi for the first time, run the following command to create a task:

SELECT hudi_sync_task_submit('SCHEMA.TABLE', 'SCHEMA.FOREIGN_TABLE');

Querying Hudi Synchronization Tasks

Query a Hudi synchronization task. In the query result, **task_id uniquely** identifies a Hudi synchronization task.

SELECT * FROM pg_task_show('SQLonHudi');

Suspending Hudi Synchronization Tasks

Query the Hudi task and obtain task_id to suspend the Hudi task.

SELECT pg_task_pause('task_id');

Example:

Suspend the synchronization task whose **task_id** is 64479410-a04c-0700-d150-3037d700fffe.

SELECT pg_task_pause('64479410-a04c-0700-d150-3037d700fffe');

Resuming Hudi Synchronization Tasks

Query the Hudi task, obtain the value of task_id, and resume the Hudi task.

SELECT pg_task_resume('task_id');

Example:

Resume the synchronization task whose **task_id** is **64479410-a04c-0700-d150-3037d700fffe**.

SELECT pg_task_resume('64479410-a04c-0700-d150-3037d700fffe');

Deleting a Hudi Synchronization Task

Query the Hudi task, obtain task id, and delete the Hudi synchronization task.

SELECT pg_task_remove('task_id');

Example:

Delete the synchronization task whose **task_id** is **64479410-a04c-0700-d150-3037d700fffe**.

SELECT pg_task_remove('64479410-a04c-0700-d150-3037d700fffe');

Querying Past Synchronization Information

Use the **hudi_sync_state_history_view** view to query information about past Hudi synchronization tasks. This view is supported only by clusters of version 9.1.0 and later.

SELECT * FROM pg_catalog.hudi_sync_state_history_view;

Table 6-4 hudi_sync_state_history_view columns

Column	Туре	Description
task_id	TEXT	Task ID
target_tbl	TEXT	Name of the synchronization target table
source_ftbl	TEXT	Name of the synchronization source table (foreign table)

Column	Туре	Description
latest_commit	TEXT	Timestamp of the latest successful synchronization
latest_sync_count	BIGINT	Number of rows that are successfully synchronized last time
latest_sync_start	TIMESTAMP WITH TIME ZONE	Start time of the latest synchronization task
latest_sync_end	TIMESTAMP WITH TIME ZONE	Time when the latest synchronization task ends
hudi_flushdisk_tim e	TEXT	Time when the hudi file is flushed to disks

Querying the Status of a Synchronization Task

Use the **hudi_show_sync_state()** function to query the status of a Hudi synchronization task.

SELECT * FROM hudi_show_sync_state();

Resetting a Hudi Synchronization Task with Consecutive Failures

Use the **pg_task_resume()** function to reset a Hudi synchronization task that fails consecutively.

If the number of consecutive failures is greater than or equal to 10, the task is automatically suspended. You need to manually call the **pg_task_resume()** function to reset the task. This function is supported only by clusters of version 9.1.0 and later.

Input parameter: task_id of the Hudi task that fails consecutively

SELECT pg_task_resume('task_id');

6.7.6 Querying a Hudi Foreign Table

You can query data in a Hudi foreign table. By default, it gives you a real-time view. You can set parameters to query the incremental data.

Querying Incremental Data

You can set incremental query parameters to implement incremental query.

SET hoodie.SCHEMA.FOREIGN_TABLE.consume.mode=incremental; SET hoodie.SCHEMA.FOREIGN_TABLE.consume.start.timestamp=start_timestamp, SET hoodie.SCHEMA.FOREIGN_TABLE.consume.ending.timestamp=end_timestamp, SELECT * FROM SCHEMA.FOREIGN_TABLE;

Example:

Query the incremental data of the MOR hudi foreign table public.rtd_mfdt_int_currency_ft from 20221207164617 to 20221207170234.
Where:

SET hoodie.public.rtd_mfdt_int_currency_ft.consume.mode=incremental;
SET hoodie.public.rtd_mfdt_int_currency_ft.consume.start.timestamp=20221207164617;
SET hoodie.public.rtd_mfdt_int_currency_ft.consume.ending.timestamp=20221207170234;
SELECT * FROM public.rtd_mfdt_int_currency_ft where _hoodie_commit_time>20221207164617 and _hoodie_commit_time<=20221207170234;

Querying the Configured Incremental Parameters

You can use the following function to check the incremental parameter configuration.

SELECT * FROM pg_show_custom_settings();

Querying the Properties of a Hudi Foreign Table (hoodie.properties)

Run the following command to query the **hoodie.properties** of the Hudi data on OBS:

SELECT * FROM hudi_get_options('SCHEMA.FOREIGN_TABLE');

Example: Query the hudi properties of the OBS foreign table **rtd_mfdt_int_unit_ft** in the current schema.

SELECT * FROM hudi_get_options('rtd_mfdt_int_unit_ft');

Querying the Maximum Timeline of a Hudi Foreign Table

Run the following command to query the maximum timeline of the hudi data on OBS, that is, the latest submitted data:

SELECT * FROM hudi_get_max_commit('SCHEMA.FOREIGN_TABLE');

Example: Query the maximum timeline of the OBS foreign table **rtd_mfdt_int_unit_ft** in the current schema.

SELECT * FROM hudi_get_max_commit('rtd_mfdt_int_unit_ft');

6.7.7 Accessing Hudi Tables on MRS

SQL on Hudi supports access to hudi tables stored on MRS. This function is supported only by clusters of version 9.1.0 or later.

Prerequisites

You have created an MRS data source. For details, see MRS Data Sources.

SQL on Hudi can read hudi tables stored on MRS. The only difference in usage compared to OBS is when creating data sources.

Accessing Multiple MRS Clusters Concurrently

Due to JDK restrictions, one JVM can store only one Kerberos configuration file at a time. As a result, one GaussDB(DWS) cluster cannot concurrently access hudi tables in multiple MRS clusters through SQL on Hudi. To avoid this issue, do as follows:

- **Step 1** Obtain the krb5.conf file of each MRS cluster from the downloaded client.
- **Step 2** Use the krb5.conf file of any MRS cluster as the file to be combined (cluster A for short).
- **Step 3** Add the KDC domain information of cluster B to **realms** in cluster A's configuration file.

Example:

```
[realms]
CLUSTER.A.COM = {
  admin_server = ClusterA_SERVER_IP:PORT
  kdc = ClusterA_KDC_IP:PORT
  kdc = ClusterA_KDC_IP:PORT
}
CLUSTER.B.COM = {
  admin_server = ClusterB_SERVER_IP:PORT
  kdc = ClusterB_KDC_IP:PORT
  kdc = ClusterB_KDC_IP:PORT
}
```

Step 4 Add the domain information of cluster B to **domain_realm** in cluster A's configuration file.

Example:

```
[domain_realm]
.cluster.a.com = CLUSTER.A.COM
.cluster.b.com = CLUSTER.B.COM
```

Step 5 Replace the original **krb5.conf** file with the combined one in the original path of each node in each cluster.

----End

A CAUTION

The preceding example is for reference only. During actual operations, you need to combine the actual KDC domain information in the cluster' **realms** or **domain realm**.

GaussDB(DWS) Sorting Rules

The collation feature allows specifying the data sorting order and data classification rules in a character set. This alleviates the restriction that the **LC_COLLATE** and **LC_CTYPE** settings of a database cannot be changed after its creation.

Overview

Every expression of a collatable data type has a collation. (The built-in collatable data types are text, varchar, and char. User-defined base types can also be marked collatable, and of course a domain over a collatable data type is collatable.) If the expression is a column reference, the collation of the expression is the defined collation of the column. If the expression is a constant, the collation is the default collation of the data type of the constant. The collation of a more complex expression is derived from the collations of its inputs.

Collation Combination Principles

- The collation of an expression can be the default collation, which means the locale settings defined for the database. It is also possible for an expression's collation to be indeterminate. In such cases, ordering operations and other operations that need to know the collation will fail.
- For a function or operator call, the collation that is derived by examining the
 argument collations is used at run time for performing the specified
 operation. If the result of the function or operator call is of a collatable data
 type, the collation is also used as the defined collation of the function or
 operator expression, in case there is a surrounding expression that requires
 knowledge of its collation.
- The collation derivation of an expression can be implicit or explicit. This distinction affects how collations are combined when multiple different collations appear in an expression. An explicit collation derivation occurs when a **COLLATE** clause is used; all other collation derivations are implicit. When multiple collations need to be combined, the following rules are used:
 - If any input expression has an explicit collation derivation, then all
 explicitly derived collations among the input expressions must be the
 same, otherwise an error is raised. If any explicitly derived collation is
 present, that is the result of the collation combination.

- Otherwise, all input expressions must have the same implicit collation derivation or the default collation. If any non-default collation is present, that is the result of the collation combination. Otherwise, the result is the default collation.
- If there are conflicting non-default implicit collations among the input expressions, then the combination is deemed to have indeterminate collation. This is not an error condition unless the particular function being invoked requires knowledge of the collation it should apply. If it does, an error will be raised at run-time.
- In a CASE expression, the comparison rule is subject to the COLLATE setting in the WHEN clause.
- Explicit COLLATE derivation takes effect only in the current query (CTE or SUBQUERY). Outside the query, implicit derivation takes effect.

Collation Tips

- Do not use multiple collations in the same query statement. Otherwise, exceptional result sets may be generated.
- Do not use multiple COLLATE clauses to specify a collation.

Case-insensitive Collation Support

Since cluster 8.1.3, GaussDB(DWS) has added the built-in case_insensitive collation, which is case-insensitive to character types in some actions (such as sorting, comparison, and hash).

Constraints:

- Supported character types: char, character, nchar, and varchar/character varying/varchar2/nvarchar2/clob/text.
- The character types **char** and **name** are not supported.
- The following encoding formats are not supported: PG_EUC_JIS_2004, PG_MULE_INTERNAL, PG_LATIN10 and PG_WIN874.
- It cannot be specified to LC_COLLATE when CREATE DATABASE is executed.
- Regular expressions are not supported.
- Record comparison of the character type (for example, **record_eq**) is not supported.
- Time series tables are not supported.
- Skew optimization is not supported.
- RoughCheck optimization is not supported.

Examples

The COLLATE clause is specified in the statement.

Set the column attribute to **case_insensitive** when creating a table.

This parameter is specified during table creation and does not need to be specified during query.

CASE expression, which is subject to the COLLATE setting in the WHEN clause.

Implicit derivation across subqueries.

```
SELECT * FROM (SELECT a collate "C" from t1) WHERE a in ('a','b');

a
---
a
b
(2 rows)

SELECT * FROM t1,(SELECT a collate "C" from t1) t2 WHERE t1.a=t2.a;

ERROR: could not determine which collation to use for string hashing

HINT: Use the COLLATE clause to set the collation explicitly.
```

CAUTION

- **collate case_insensitive** is an insensitive sorting, and the result set is uncertain. If sensitive sorting is used after **collate case_insensitive** sorting, the result set may be unstable. Therefore, do not use sensitive sorting and insensitive sorting together in statements.
- If **collate case_insensitive** is used to specify character behaviors as case-insensitive, the performance will be affected. If you require high performance, exercise caution when configuring this parameter.

8 GaussDB(DWS) User-Defined Functions

■ NOTE

- The hybrid data warehouse (deployed in standalone mode) does not support userdefined functions.
- The hybrid data warehouse (standalone) does 8.2.0.100 and later versions support OBS import and export.

8.1 GaussDB(DWS) PL/Java Functions

GaussDB(DWS) supports Java user-defined functions (UDFs) to extend data processing capabilities. Developers can write UDFs in Java and use them in SQL queries or other data processing tasks.

Constraints

Java UDF can be used for some Java logical computing. You are not advised to encapsulate services in Java UDF.

- You are not advised to connect to a database in any way (for example, JDBC) in Java functions.
- Currently, only data types listed in Table 8-1 are supported. Other data types, such as user-defined data types and complex data types (for example, Java array and its derived types) are not supported.
- Currently, UDAF and UDTF are not supported.
- GaussDB(DWS) PL/Java is based on open source PL/Java 1.5.5 and uses JRE 1.8.0.432.

Using Java UDFs

With PL/Java, you can compile Java methods using a Java IDE, deploy the JAR files to the GaussDB(DWS) database, and create functions as a database administrator. For compatibility, use JRE 1.8.0.432 for compiling.

Step 1 Compile a JAR package.

Java method implementation and JAR package archiving can be achieved in an integrated development environment (IDE). The following is a simple example of

compilation and archiving through command lines. You can create a JAR package that contains a single method in the similar way.

First, prepare an **Example.java** file that contains a method for converting substrings to uppercase. In the following example, **Example** is the class name and **upperString** is the method name:

```
public class Example
{
    public static String upperString (String text, int beginIndex, int endIndex)
    {
        return text.substring(beginIndex, endIndex).toUpperCase();
    }
}
```

Then, create a manifest.txt file containing the following content:

```
Manifest-Version: 1.0
Main-Class: Example
Specification-Title: "Example"
Specification-Version: "1.0"
Created-By: 1.6.0_35-b10-428-11M3811
Build-Date: 08/14/2018 10:09 AM
```

Manifest-Version specifies the version of the manifest file. Main-Class specifies the main class used by the .jar file. Specification-Title and Specification-Version are the extended attributes of the package. Specification-Title specifies the title of the extended specification and Specification-Version specifies the version of the extended specification. Created-By specifies the person who created the file. Build-Date specifies the date when the file was created.

Finally, archive the .java file and package it into javaudf-example.jar.

```
javac Example.java
jar cfm javaudf-example.jar manifest.txt Example.class
```

NOTICE

JAR package names must comply with JDK rules. If a name contains invalid characters, an error occurs when a function is deployed or used.

Step 2 Deploy the JAR package.

Place the JAR package on the OBS server using the method described in For details, see "Uploading a File" in *Object Storage Service Console Operation Guide*.. Then, create the AK/SK. For details about how to obtain the AK/SK, see section **Creating Access Keys (AK and SK)**. Log in to the database and run the **gs_extend_library** function to import the file to GaussDB(DWS).

```
SELECT gs_extend_library('addjar', 'obs://bucket/path/javaudf-example.jar accesskey=access_key_value_to_be_replaced secretkey=secret_access_key_value_to_be_replaced region=region_name libraryname=example');
```

For details about how to use the **gs_extend_library** function, see **Manage JAR packages and files**. Change the values of AK and SK as needed. Replace *region_name* with an actual region name.

Step 3 Use a PL/Java function.

Log in to the database as a user who has the **sysadmin** permission (for example, dbadmin) and create the **java_upperstring** function:

CREATE FUNCTION java_upperstring(VARCHAR, INTEGER, INTEGER)
RETURNS VARCHAR
AS 'Example.upperString'
LANGUAGE JAVA;

- The data type defined in the java_upperstring function should be a type in GaussDB(DWS) and match the data type defined in Step 1 in the upperString method in Java. For details about the mapping between GaussDB(DWS) and Java data types, see Table 8-1.
- The AS clause specifies the class name and static method name of the Java method invoked by the function. The format is *Class name.Method name*. The class name and method name must match the Java class and method defined in **Step 1**.
- To use PL/Java functions, set LANGUAGE to JAVA.
- For details about CREATE FUNCTION, see Create functions.

Execute the java_upperstring function.

SELECT java_upperstring('test', 0, 1);

The expected result is as follows:

```
java_upperstring
------
T
(1 row)
```

Step 4 Authorize a common user to use the PL/Java function.

Create a common user named udf user.

CREATE USER udf_user PASSWORD 'password;

This command grants user **udf_user** the permission for the java_upperstring function. Note that the user can use this function only if it also has the permission for using the schema of the function.

```
GRANT ALL PRIVILEGES ON SCHEMA public TO udf_user;
GRANT ALL PRIVILEGES ON FUNCTION java_upperstring(VARCHAR, INTEGER, INTEGER) TO udf_user;
```

Log in to the database as user **udf_user**.

SET SESSION SESSION AUTHORIZATION udf_user PASSWORD 'password;

Execute the java_upperstring function.

SELECT public.java_upperstring('test', 0, 1);

The expected result is as follows:

```
java_upperstring
-----
T
(1 row)
```

Step 5 Delete the function.

If you no longer need this function, delete it. DROP FUNCTION java_upperstring;

Step 6 Uninstall the JAR package.

Use the gs_extend_library function to uninstall the JAR package.

SELECT gs_extend_library('rmjar', 'libraryname=example');

----End

Mapping for Basic Data Types

Table 8-1 PL/Java mapping for default data types

GaussDB(DWS)	Java	
BOOLEAN	boolean	
"char"	byte	
bytea	byte[]	
SMALLINT	short	
INTEGER	int	
BIGINT	long	
FLOAT4	float	
FLOAT8	double	
CHAR	java.lang.String	
VARCHAR	java.lang.String	
TEXT	java.lang.String	
name	java.lang.String	
DATE	java.sql.Timestamp	
TIME	java.sql.Time (stored value treated as local time)	
TIMETZ	java.sql.Time	
TIMESTAMP	java.sql.Timestamp	
TIMESTAMPTZ	java.sql.Timestamp	

UDF Example: SQL Definition and Usage

Manage JAR packages and files.

A database user having the **sysadmin** permission can use the gs_extend_library function to deploy, view, and delete JAR packages in the database. The syntax of the function is as follows:

SELECT gs_extend_library('[action]', '[operation]');

□ NOTE

- action: operation action. The options are as follows:
 - **ls**: Displays JAR packages in the database and checks the MD5 value consistency of files on each node.
 - addjar: deploys a JAR package on the OBS server in the database.
 - rmjar: Deletes JAR packages from the database.
- **operation**: operation string. The format can be either of the following: obs://[bucket]/[source_filepath] accesskey=[accesskey] secretkey=[secretkey] region=[region] libraryname=[libraryname]
 - **bucket**: name of the bucket to which the OBS file belongs. It is mandatory.
 - **source filepath**: file path on the OBS server. Only .jar files are supported.
 - accesskey: key obtained for accessing the OBS service. It is mandatory.
 - secret_key: secret key obtained for the OBS service. It is mandatory.
 - **region**: region where the OBS bucket stored in the JAR package of a user-defined function belongs to. This parameter is mandatory.
 - **libraryname**: user-defined library name, which is used to invoke JAR files in GaussDB(DWS). If **action** is set to **addjar** or **rmjar**, **libraryname** must be specified. If **action** is set to **ls**, **libraryname** is optional. Note that a user-defined library name cannot contain the following characters: /|;&\$<>\'{}"() []~*?!

• Create functions.

PL/Java functions can be created using the **CREATE FUNCTION** syntax and are defined as **LANGUAGE JAVA**, including the **RETURNS** and **AS** clauses.

- To use CREATE FUNCTION, specify the name and parameter type for the function to be created.
- The **RETURNS** clause specifies the return type for the function.
- The AS clause specifies the class name and static method name of the
 Java method to be invoked. If the NULL value needs to be transferred to
 the Java method as an input parameter, specify the name of the Java
 encapsulation class corresponding to the parameter type. For details, see
 UDF Example: Processing NULL Values.
- For details about the syntax, see CREATE FUNCTION.

```
CREATE [ OR REPLACE ] FUNCTION function name
( [ { argname [ argmode ] argtype [ { DEFAULT | := | = } expression ]} [, ...] ])
[ RETURNS rettype [ DETERMINISTIC ] ]
LANGUAGE JAVA
  { IMMUTABLE | STABLE | VOLATILE }
  | [ NOT ] LEAKPROOF
  WINDOW
  { CALLED ON NULL INPUT | RETURNS NULL ON NULL INPUT | STRICT }
  AUTHID CURRENT_USER}
  | { FENCED }
  COST execution_cost
  ROWS result_rows
  | SET configuration_parameter { {TO |=} value | FROM CURRENT}
] [...]
  AS 'class_name.method_name' ( { argtype } [, ...] )
```

Use functions.

During execution, PL/Java searches for the Java class specified by a function among all the deployed JAR packages, which are ranked by name in

alphabetical order, invokes the Java method in the first found class, and returns results.

Delete functions.

PL/Java functions can be deleted by using the **DROP FUNCTION** syntax. For details about the syntax, see DROP FUNCTION.

```
DROP FUNCTION [ IF EXISTS ] function_name [ ( [ {[ argmode ] [ argname ] argtype} [, ...] ] ) [ CASCADE | RESTRICT ] ];
```

To delete an overloaded function (for details, see **UDF Example: Overloaded Functions**), specify **argtype** in the function. To delete other functions, simply specify **function_name**.

Authorize permissions for functions.

Only user **sysadmin** can create PL/Java functions. It can also grant other users the permission to use the PL/Java functions. For details about the syntax, see GRANT.

```
GRANT { EXECUTE | ALL [ PRIVILEGES ] }
ON { FUNCTION {function_name ( [ {[ argmode ] [ arg_name ] arg_type} [, ...] ] )} [, ...]
| ALL FUNCTIONS IN SCHEMA schema_name [, ...] }
TO { [ GROUP ] role_name | PUBLIC } [, ...]
[ WITH GRANT OPTION ];
```

UDF Example: Processing Array Types

GaussDB(DWS) can convert basic array types. You only need to append a pair of square brackets ([]) to the data type when creating a function.

```
CREATE FUNCTION java_arrayLength(INTEGER[])
RETURNS INTEGER
AS 'Example.getArrayLength'
LANGUAGE JAVA;
```

Java code is similar to the following:

```
public class Example
{
   public static int getArrayLength(Integer[] intArray)
   {
      return intArray.length;
   }
}
```

Invoke the following statement:

```
SELECT java_arrayLength(ARRAY[1, 2, 3]);
```

The expected result is as follows:

```
java_arrayLength
------3
(1 row)
```

UDF Example: Processing NULL Values

NULL values cannot be handled for GaussDB(DWS) data types that are mapped and can be converted to simple Java types by default. If you use a Java function to obtain and process the **NULL** value transferred from GaussDB(DWS), specify the Java encapsulation class in the **AS** clause as follows:

```
CREATE FUNCTION java_countnulls(INTEGER[])
RETURNS INTEGER
```

```
AS 'Example.countNulls(java.lang.Integer[])'
LANGUAGE JAVA;
```

Java code is similar to the following:

```
public class Example
{
   public static int countNulls(Integer[] intArray)
   {
      int nullCount = 0;
      for (int idx = 0; idx < intArray.length; ++idx)
      {
        if (intArray[idx] == null)
           nullCount++;
      }
      return nullCount;
   }
}</pre>
```

Invoke the following statement:

```
SELECT java_countNulls(ARRAY[null, 1, null, 2, null]);
```

The expected result is as follows:

UDF Example: Overloaded Functions

PL/Java supports overloaded functions. You can create functions with the same name or invoke overloaded functions from Java code. The procedure is as follows:

Step 1 Create overloaded functions.

For example, create two Java methods with the same name, and specify the methods dummy(int) and dummy(String) with different parameter types.

```
public class Example
{
    public static int dummy(int value)
    {
        return value*2;
    }
    public static String dummy(String value)
    {
        return value;
    }
}
```

In addition, create two functions with the same names as the above two functions in GaussDB(DWS).

```
CREATE FUNCTION java_dummy(INTEGER)
RETURNS INTEGER
AS 'Example.dummy'
LANGUAGE JAVA;

CREATE FUNCTION java_dummy(VARCHAR)
RETURNS VARCHAR
AS 'Example.dummy'
LANGUAGE JAVA;
```

Step 2 Invoke the overloaded functions.

GaussDB(DWS) invokes the functions that match the specified parameter type. The results of invoking the above two functions are as follows:

```
SELECT java_dummy(5);
java_dummy
------
10
(1 row)

SELECT java_dummy('5');
java_dummy
------
5
(1 row)
```

Note that GaussDB(DWS) may implicitly convert data types. Therefore, you are advised to specify the parameter type when invoking an overloaded function.

```
SELECT java_dummy(5::varchar);
java_dummy
------
5
(1 row)
```

In this case, the specified parameter type is preferentially used for matching. If there is no Java method matching the specified parameter type, the system implicitly converts the parameter and searches for Java methods based on the conversion result.

```
SELECT java_dummy(5::INTEGER);
java_dummy
------
10
(1 row)

DROP FUNCTION java_dummy(INTEGER);
SELECT java_dummy(5::INTEGER);
java_dummy
------
5
(1 row)
```

NOTICE

Data types supporting implicit conversion are as follows:

- SMALLINT: It can be converted to the INTEGER type by default.
- **SMALLINT** and **INTEGER**: They can be converted to the **BIGINT** type by default.
- TINYINT, SMALLINT, INTEGER, and BIGINT: They can be converted to the BOOL type by default.
- The following data types can be converted to TEXT by default: CHAR, NAME, BIGINT, INTEGER, SMALLINT, TINYINT, RAW, FLOAT4, FLOAT8, BPCHAR, VARCHAR, NVARCHAR2, DATE, TIMESTAMP, TIMESTAMPTZ, NUMERIC, and SMALLDATETIME.
- The following data types can be converted to VARCHAR by default: TEXT, CHAR, BIGINT, INTEGER, SMALLINT, TINYINT, RAW, FLOAT4, FLOAT8, BPCHAR, DATE, NVARCHAR2, TIMESTAMP, NUMERIC, and SMALLDATETIME.

Step 3 Delete the overloaded functions.

To delete an overloaded function, specify the parameter type for the function. Otherwise, the function cannot be deleted.

DROP FUNCTION java_dummy(INTEGER);

----End

UDF-related GUC Parameters

• udf_memory_limit

A system-level GUC parameter. It is used to limit the physical memory used by each CN or DN for executing UDFs. The default value is **0.05** * max_process_memory. You can use the postgresql.conf file to modify the parameter setting. The modification takes effect only after the database is restarted.

NOTICE

- udf_memory_limit is a part of max_process_memory. When a CN or DN is started, memory calculated by udf_memory_limit minus 200 MB will be reserved for UDF Worker processes. CN and DN processes are different from the UDF Worker process, and the CN and DN processes will save memory for the UDF Worker process.
 - For example, if max_process_memory is set to 10GB on a DN and udf_memory_limit is set to 4GB, the DN can use a maximum of 6.2 GB memory, that is, 10 GB (4 GB 200 MB). This case applies even if no UDF is executed. By default, the value of udf_memory_limit is 0.05 * max_process_memory. Querying the pv_total_memory_detail view will prove that the value of process_used_memory would never exceed the calculation result of max_process_memory (udf_memory_limit 200 MB).
- If the UDF process is disconnected, an error message will be displayed. Example: "memory in UDF Work Process is limited by cgroup: [usage: xxx, max_usage_history: xxx, limit: xxx]." You can learn the current memory usage from this message. In the error information, usage indicates the total physical memory used by the rest of the UDF process after a UDF process is killed. max_usage_history indicates the highest memory usage of the UDF process after the UDF instance is started. limit indicates the maximum memory used by the UDF process. If the value of max_usage_history is close to the value of limit, the memory usage of the current cluster may exceed the limit. In this case, optimize workloads or adjust the value of udf_memory_limit as needed.
- Executing a simplest Java UDF on a CN consumes about 50 MB physical memory. You can set this parameter based on the memory usage and concurrency of Java functions to be used. After this parameter is added, you are not advised to set UDFWorkerMemHardLimit and FencedUDFMemoryLimit.
- If the parallelism of the UDF process is excessively high and the memory usage exceeds the udf_memory_limit value, unexpected situations such as process exit may occur. In this scenario, the execution result may be unreliable. You are advised to set this parameter to reserve sufficient memory based on the site requirements. If the system has the /var/log/messages log, check the log to see whether the memory is insufficient because the cgroup memory limit has been reached. If the memory is severely insufficient, the UDF master process may exit. You can view the UDF log for analysis. The default UDF log path is \$GAUSSLOG/cm/cm_agent/pg_log. For example, if the log below is displayed, the memory resources are insufficient and the UDF master process exits. In this case, you need to check the udf_memory_limit parameter.
 - 0 [BACKEND] FATAL: poll() failed: Bad address, please check the parameter:udf_memory_limit to make sure there is enough memory.

• FencedUDFMemoryLimit

A session-level GUC parameter. It is used to specify the maximum virtual memory used by a single Fenced UDF Worker process initiated by a session. SET FencedUDFMemoryLimit='512MB';

The value range of this parameter is (150 MB, 1G). If the value is greater than 1G, an error will be reported immediately. If the value is less than or equal to 150 MB, an error will be reported during function invoking.

NOTICE

- If **FencedUDFMemoryLimit** is set to **0**, the virtual memory for a Fenced UDF Worker process will not be limited.
- You are advised to use udf_memory_limit to control the physical memory used by Fenced UDF Worker processes. You are not advised to use
 FencedUDFMemoryLimit, especially when Java UDFs are used. If you are clear about the impact of this parameter, set it based on the following information:
 - After a C Fenced UDF Worker process is started, it will occupy about 200 MB virtual memory, and about 16 MB physical memory.
 - After a Java Fenced UDF Worker process is started, it will occupy about 2.5 GB virtual memory, and about 50 MB physical memory.

Exception Handling

If there is an exception in a JVM, PL/Java will export JVM stack information during the exception to a client.

Logging

PL/Java uses the standard Java Logger. Therefore, you can record logs as follows:

Logger.getAnonymousLogger().config("Time is " + new Date(System.currentTimeMillis()));

An initialized Java Logger class is set to the **CONFIG** level by default, corresponding to the **LOG** level in GaussDB(DWS). In this case, log messages generated by Java Logger are all redirected to the GaussDB(DWS) backend. Then, the log messages are written into server logs or displayed on the user interface. MPPDB server logs record information at the **LOG**, **WARNING**, and **ERROR** levels. The SQL user interface displays logs at the **WARNING** and **ERROR** levels. The following table lists mapping between Java Logger levels and GaussDB(DWS) log levels.

Table 8-2 PL/Java log levels

java.util.logging.Level	GaussDB(DWS) Log Level
SERVER	ERROR
WARNING	WARNING
CONFIG	LOG
INFO	INFO
FINE	DEBUG1

java.util.logging.Level	GaussDB(DWS) Log Level
FINER	DEBUG2
FINEST	DEBUG3

You can change Java Logger levels. For example, if the Java Logger level is changed to **SEVERE** by the following Java code, log messages (**msg**) will not be recorded in GaussDB(DWS) logs during **WARNING** logging.

Logger log = Logger.getAnonymousLogger();
Log.setLevel(Level.SEVERE);
log.log(Level.WARNING, msg);

Security Issues

In GaussDB(DWS), PL/Java is an untrusted language. Only user **sysadmin** can create PL/Java functions. The user can grant other users the permission for using the PL/Java functions. For details, see **Authorize permissions for functions**.

In addition, PL/Java controls user access to file systems, forbidding users from reading most system files, or writing, deleting, or executing any system files in Java methods.

8.2 GaussDB(DWS) PL/pgSQL Functions

PL/pgSQL is similar to PL/SQL of Oracle. It is a loadable procedural language.

The functions created using PL/pgSQL can be used in any place where you can use built-in functions. For example, you can create calculation functions with complex conditions and use them to define operators or use them for index expressions.

SQL is used by most databases as a query language. It is portable and easy to learn. Each SQL statement must be executed independently by a database server.

In this case, when a client application sends a query to the server, it must wait for it to be processed, receive and process the results, and then perform some calculation before sending more queries to the server. If the client and server are not on the same machine, all these operations will cause inter-process communication and increase network loads.

PL/pgSQL enables a whole computing part and a series of queries to be grouped inside a database server. This makes procedural language available and SQL easier to use. In addition, the client/server communication cost is reduced.

- Extra round-trip communication between clients and servers is eliminated.
- Intermediate results that are not required by clients do not need to be sorted or transmitted between the clients and servers.
- Parsing can be skipped in multiple rounds of queries.

PL/pgSQL can use all data types, operators, and functions in SQL.

For details about the PL/pgSQL syntax for creating functions, see **CREATE FUNCTION**. As mentioned earlier, PL/pgSQL is similar to PL/SQL of Oracle and is a

loadable procedural language. Its application method is similar to that of **GaussDB(DWS) Stored Procedure**. There is only one difference. Stored procedures have no return values but the functions have.

9 GaussDB(DWS) Stored Procedure

9.1 Overview

What Is a GaussDB(DWS) Stored Procedure?

In GaussDB(DWS), business rules and logics are saved as stored procedures.

A stored procedure is a combination of SQL, PL/SQL, and Java statements. Stored procedures can move the code that executes business rules from applications to databases. In this way, code can be used by multiple programs at a time.

For details about how to create and call a stored procedure, see **CREATE PROCEDURE**.

The functions created using the PL/pgSQL language mentioned in **GaussDB(DWS) PL/pgSQL Functions** are similar to the application methods of stored procedures. Unless otherwise specified, the following sections apply to stored procedures and PL/pgSQL functions.

GaussDB(DWS) Stored Procedure Data Types

A data type refers to a value set and an operation set defined on the value set. A GaussDB(DWS) database consists of tables, each of which is defined by its own columns. Each column corresponds to a data type. GaussDB(DWS) uses corresponding functions to perform operations on data based on data types. For example, GaussDB(DWS) can perform addition, subtraction, multiplication, and division operations on data of numeric values.

9.2 Converting Data Types in GaussDB(DWS) Stored Procedures

Certain data types in the database support implicit data type conversions, such as assignments and parameters invoked by functions. For other data types, you can use the type conversion functions provided by GaussDB(DWS), such as the CAST function, to forcibly convert them.

Table 9-1 lists common implicit data type conversions in GaussDB(DWS).

NOTICE

The valid value range of DATE supported by GaussDB(DWS) is from 4713 B.C. to 294276 A.D.

Table 9-1 Implicit data type conversions

Raw Data Type	Target Data Type	Remarks
CHAR	VARCHAR2	N/A
CHAR	NUMBER	Raw data must consist of digits.
CHAR	DATE	Raw data cannot exceed the valid date range.
CHAR	RAW	N/A
CHAR	CLOB	N/A
VARCHAR2	CHAR	N/A
VARCHAR2	NUMBER	Raw data must consist of digits.
VARCHAR2	DATE	Raw data cannot exceed the valid date range.
VARCHAR2	CLOB	N/A
NUMBER	CHAR	N/A
NUMBER	VARCHAR2	N/A
DATE	CHAR	N/A
DATE	VARCHAR2	N/A
RAW	CHAR	N/A
RAW	VARCHAR2	N/A
CLOB	CHAR	N/A
CLOB	VARCHAR2	N/A
CLOB	NUMBER	Raw data must consist of digits.
INT4	CHAR	N/A

9.3 GaussDB(DWS) Stored Procedure Array and Record

9.3.1 Arrays

Use of Array Types

Before the use of arrays, an array type needs to be defined:

Define an array type immediately after the **AS** keyword in a stored procedure. Run the following statement:

TYPE array_type IS VARRAY(size) OF data_type [NOT NULL];

Its parameters are as follows:

- array_type: indicates the name of the array type to be defined.
- **VARRAY**: indicates the array type to be defined.
- **size**: indicates the maximum number of members in the array type to be defined. The value is a positive integer.
- **data_type**: indicates the types of members in the array type to be created.
- **NOT NULL**: an optional constraint. It can be used to ensure that none of the elements in the array is **NULL**.

- In GaussDB(DWS), an array automatically increases. If an access violation occurs, a null value will be returned, and no error message will be reported. If out-of-bounds write occurs in an array, the message **Subscript outside of limit** is displayed.
- The scope of an array type defined in a stored procedure takes effect only in this storage process.
- It is recommended that you use one of the preceding methods to define an array type. If both methods are used to define the same array type, GaussDB(DWS) prefers the array type defined in a stored procedure to declare array variables.

In GaussDB(DWS) 8.1.0 and earlier versions, the system does not verify the length of array elements and out-of-bounds write because the array can automatically increase. This version adds related constraints to be compatible with Oracle databases. If out-of-bounds write exists, you can configure **varray_verification** in the parameter **behavior_compat_options** to be compatible with previously unverified operations.

Example:

```
-- Declare an array in a stored procedure.

CREATE OR REPLACE PROCEDURE array_proc

AS

TYPE ARRAY_INTEGER IS VARRAY(1024) OF INTEGER;--Define the array type.

TYPE ARRAY_INTEGER_NOT_NULL IS VARRAY(1024) OF INTEGER NOT NULL;-- Defines non-null array types.

ARRINT ARRAY_INTEGER: = ARRAY_INTEGER(); --Declare the variable of the array type.

BEGIN

ARRINT.extend(10);
FOR I IN 1..10 LOOP

ARRINT(I) := I;
END LOOP;
DBMS_OUTPUT.PUT_LINE(ARRINT.COUNT);
DBMS_OUTPUT.PUT_LINE(ARRINT.COUNT);
```

```
DBMS_OUTPUT.PUT_LINE(ARRINT(10));
DBMS_OUTPUT.PUT_LINE(ARRINT(ARRINT.Iast));
DBMS_OUTPUT.PUT_LINE(ARRINT(ARRINT.last));
END;

-- Invoke the stored procedure.
CALL array_proc();
10
1
10
-- Delete the stored procedure.
DROP PROCEDURE array_proc;
```

Declaration and Use of Rowtype Arrays

In addition to the declaration and use of common arrays and non-null arrays in the preceding example, the array also supports the declaration and use of rowtype arrays.

Example:

```
-- Use the COUNT function on an array in a stored procedure.
CREATE TABLE tbl (a int, b int);
INSERT INTO tbl VALUES(1, 2),(2, 3),(3, 4);
CREATE OR REPLACE PROCEDURE array_proc
  CURSOR all tbl IS SELECT * FROM tbl ORDER BY a:
  TYPE tbl_array_type IS varray(50) OF tbl%rowtype; -- Defines the array of the rowtype type. tbl indicates
any table.
  tbl_array tbl_array_type;
  tbl_item tbl%rowtype;
  inx1 int:
BEGIN
  tbl_array := tbl_array_type();
  inx1 := 0;
  FOR tbl_item IN all_tbl LOOP
     inx1 := inx1 + 1;
     tbl_array(inx1) := tbl_item;
  END LOOP;
  WHILE inx1 IS NOT NULL LOOP
     DBMS_OUTPUT.PUT_LINE('tbl_array(inx1).a=' || tbl_array(inx1).a || ' tbl_array(inx1).b=' ||
tbl_array(inx1).b);
     inx1 := tbl_array.PRIOR(inx1);
  END LOOP;
END;
```

The execution output is as follows:

```
call array_proc();
tbl_array(inx1).a=3 tbl_array(inx1).b=4
tbl_array(inx1).a=2 tbl_array(inx1).b=3
tbl_array(inx1).a=1 tbl_array(inx1).b=2
```

Array Related Functions

GaussDB(DWS) supports Oracle-related array functions. You can use the following functions to obtain array attributes or perform operations on the array content.

COUNT

Returns the number of elements in the current array. Only the initialized elements or the elements extended by the EXTEND function are counted.

Use:

varray.COUNT or varray.COUNT()

Example:

```
-- Use the COUNT function on an array in a stored procedure.

CREATE OR REPLACE PROCEDURE test_varray

AS

TYPE varray_type IS VARRAY(20) OF INT;
v_varray varray_type;

BEGIN

v_varray := varray_type(1, 2, 3);

DBMS_OUTPUT.PUT_LINE('v_varray.count=' || v_varray.count);
v_varray.extend;

DBMS_OUTPUT.PUT_LINE('v_varray.count=' || v_varray.count);

END;
```

The execution output is as follows:

```
call test_varray();
v_varray.count=3
v_varray.count=4
```

FIRST and LAST

The FIRST function can return the subscript of the first element. The LAST function can return the subscript of the last element.

Use:

varray.FIRST or varray.FIRST()

varray.LAST or varray.LAST()

Example:

```
-- Use the FIRST and LAST functions on an array in a stored procedure.

CREATE OR REPLACE PROCEDURE test_varray

AS

TYPE varray_type IS VARRAY(20) OF INT;
v_varray varray_type;

BEGIN

v_varray := varray_type(1, 2, 3);

DBMS_OUTPUT.PUT_LINE('v_varray.first=' || v_varray.first);

DBMS_OUTPUT.PUT_LINE('v_varray.last=' || v_varray.last);

END;
```

The execution output is as follows:

```
call test_varray();
v_varray.first=1
v_varray.last=3
```

EXTEND

□ NOTE

The EXTEND function is used to be compatible with two Oracle database operations. In GaussDB(DWS), an array automatically grows, and the EXTEND function is not necessary. For a newly written stored procedure, you do not need to use the EXTEND function.

The EXTEND function can extend arrays. The EXTEND function can be invoked in either of the following ways:

Method 1:

EXTEND contains an integer input parameter, indicating that the array size is extended by the specified length. When the EXTEND function is executed, the values returned by the COUNT and LAST functions will be updated accordingly.

Use:

varray.EXTEND(size)

By default, one bit is added to the end of *varray*.**EXTEND**, which is equivalent to *varray*.**EXTEND(1)**.

Method 2:

EXTEND contains two integer input parameters. The first parameter indicates the length of the extended size. The second parameter indicates that the value of the extended array element is the same as that of the element with the **index** subscript.

Use:

varray.EXTEND(size, index)

Example:

```
--- Use the EXTEND function on an array in a stored procedure.

CREATE OR REPLACE PROCEDURE test_varray

AS

TYPE varray_type IS VARRAY(20) OF INT;
v_varray varray_type;

BEGIN

v_varray := varray_type(1, 2, 3);
v_varray.extend(3);

DBMS_OUTPUT.PUT_LINE('v_varray.count=' || v_varray.count);
v_varray.extend(2,3);

DBMS_OUTPUT.PUT_LINE('v_varray.count=' || v_varray.count);

DBMS_OUTPUT.PUT_LINE('v_varray(7)=' || v_varray(7));

DBMS_OUTPUT.PUT_LINE('v_varray(8)=' || v_varray(7));

END;
```

The execution output is as follows:

```
call test_varray();
v_varray.count=6
v_varray.count=8
v_varray(7)=3
v_varray(8)=3
```

NEXT and PRIOR

The NEXT and PRIOR functions are used for cyclic array traversal. The NEXT function returns the subscript of the next array element based on the input parameter **index**. If the subscript reaches the maximum value, **NULL** is returned. The PRIOR function returns the subscript of the previous array element based on the input parameter **index**. If the minimum value of the array subscript is reached, **NULL** is returned.

Use:

varray.NEXT(index)

varray.PRIOR(index)

Example:

```
-- Use the NEXT and PRIOR functions on an array in a stored procedure.
CREATE OR REPLACE PROCEDURE test_varray
AS
  TYPE varray_type IS VARRAY(20) OF INT;
  v_varray varray_type;
  i int;
BEGIN
  v_varray := varray_type(1, 2, 3);
  i := v_varray.COUNT;
  WHILE I IS NOT NULL LOOP
     DBMS_OUTPUT_LINE('test prior v_varray('||i||')=' || v_varray(i));
     i := v_varray.PRIOR(i);
  END LOOP;
  i := 1;
  WHILE I IS NOT NULL LOOP
     DBMS_OUTPUT.PUT_LINE('test next v_varray('||i||')=' || v_varray(i));
     i := v_varray.NEXT(i);
  END LOOP;
END;
```

The execution output is as follows:

```
call test_varray();
test prior v_varray(3)=3
test prior v_varray(2)=2
test prior v_varray(1)=1
test next v_varray(1)=1
test next v_varray(2)=2
test next v_varray(3)=3
```

EXISTS

Determines whether an array subscript exists.

Use:

varray.EXISTS(index)

Example:

```
-- Use the EXISTS function on an array in a stored procedure.

CREATE OR REPLACE PROCEDURE test_varray

AS

TYPE varray_type IS VARRAY(20) OF INT;
v_varray varray_type;

BEGIN

v_varray := varray_type(1, 2, 3);

IF v_varray.EXISTS(1) THEN

DBMS_OUTPUT.PUT_LINE('v_varray.EXISTS(1)');

END IF;

IF NOT v_varray.EXISTS(10) THEN

DBMS_OUTPUT.PUT_LINE('NOT v_varray.EXISTS(10)');

END IF;

END;

/
```

The execution output is as follows:

```
call test_varray();
v_varray.EXISTS(1)
NOT v_varray.EXISTS(10)
```

TRIM

Deletes a specified number of elements from the end of an array.

Use:

varray.TRIM(size)

varray.**TRIM** is equivalent to varray.**TRIM(1)**, because the default input parameter is **1**.

Example:

```
-- Use the TRIM function on an array in a stored procedure.

CREATE OR REPLACE PROCEDURE test_varray

AS

TYPE varray_type IS VARRAY(20) OF INT;
v_varray varray_type;

BEGIN

v_varray:= varray_type(1, 2, 3, 4, 5);
v_varray.trim(3);

DBMS_OUTPUT.PUT_LINE('v_varray.count' || v_varray.count);
v_varray.trim;

DBMS_OUTPUT.PUT_LINE('v_varray.count:' || v_varray.count);

END;
```

The execution output is as follows:

```
call test_varray();
v_varray.count:2
v_varray.count:1
```

DELETE

Deletes all elements from an array.

Use:

varray.DELETE or varray.DELETE()

Example:

```
-- Use the DELETE function on an array in a stored procedure.

CREATE OR REPLACE PROCEDURE test_varray

AS

TYPE varray_type IS VARRAY(20) OF INT;
v_varray varray_type;

BEGIN

v_varray:= varray_type(1, 2, 3, 4, 5);
v_varray.delete;
DBMS_OUTPUT.PUT_LINE('v_varray.count:' || v_varray.count);

END;
```

The execution output is as follows:

```
call test_varray();
v_varray.count:0
```

LIMIT

Returns the allowed maximum length of an array.

Use:

varray.LIMIT or varray.LIMIT()

Example:

```
-- Use the LIMIT function on an array in a stored procedure.

CREATE OR REPLACE PROCEDURE test_varray

AS

TYPE varray_type IS VARRAY(20) OF INT;
v_varray varray_type;

BEGIN

v_varray := varray_type(1, 2, 3, 4, 5);

DBMS_OUTPUT.PUT_LINE('v_varray.limit:' || v_varray.limit);

END;

/
```

The execution output is as follows:

```
call test_varray();
v_varray.limit:20
```

9.3.2 record

record Variables

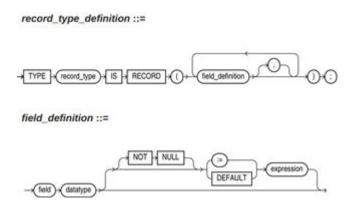
Perform the following operations to create a record variable:

Define a record type and use this type to declare a variable.

Syntax

For the syntax of the record type, see Figure 9-1.

Figure 9-1 Syntax of the record type



The syntax is described as follows:

- record_type: record name
- **field**: record columns
- datatype: record data type
- expression: expression for setting a default value

Ⅲ NOTE

In GaussDB(DWS):

- When assigning values to record variables, you can:
 - Declare a record type and define member variables of this type when you declare a function or stored procedure.
 - Assign the value of a record variable to another record variable.
 - Use SELECT INTO or FETCH to assign values to a record type.
 - Assign the **NULL** value to a record variable.
- The **INSERT** and **UPDATE** statements cannot use a record variable to insert or update data.
- Just like a variable, a record column of the compound type does not have a default value in the declaration.

Examples

```
The table used in the following stored procedure is defined as follows:
CREATE TABLE emp_rec
  empno
                numeric(4,0),
                character varying(10),
  ename
              character varying(9),
  job
  mgr
               numeric(4,0),
  hiredate
               timestamp(0) without time zone,
  sal
              numeric(7,2),
  comm
                numeric(7,2),
  deptno
                numeric(2,0)
with (orientation = column,compression=middle)
distribute by hash (sal);
\d emp_rec
          Table "public.emp_rec"
                                 | Modifiers
 Column |
                   Type
empno | numeric(4,0)
                                   | not null
ename | character varying(10)
     | character varying(9)
        | numeric(4,0)
hiredate | timestamp(0) without time zone |
sal | numeric(7,2)
comm | numeric(7,2)
deptno | numeric(2,0)
-- Perform array operations in the stored procedure.
CREATE OR REPLACE FUNCTION regress_record(p_w VARCHAR2)
RETURNS
VARCHAR2 AS $$
DECLARE
 -- Declare a record type.
 type rec_type is record (name varchar2(100), epno int);
 employer rec_type;
  -- Use %type to declare the record type.
 type rec_type1 is record (name emp_rec.ename%type, epno int not null :=10);
 employer1 rec_type1;
  -- Declare a record type with a default value.
 type rec_type2 is record (
      name varchar2 not null := 'SCOTT',
      epno int not null :=10);
  employer2 rec type2;
  CURSOR C1 IS select ename, empno from emp_rec order by 1 limit 1;
```

```
BEGIN
    -- Assign a value to a member record variable.
   employer.name := 'WARD';
   employer.epno = 18;
   raise info 'employer name: %, epno:%', employer.name, employer.epno;
   -- Assign the value of a record variable to another variable.
   employer1 := employer;
   raise info 'employer1 name: %, epno: %',employer1.name, employer1.epno;
   -- Assign the NULL value to a record variable.
   employer1 := NULL;
   raise info 'employer1 name: %, epno: %',employer1.name, employer1.epno;
    -- Obtain the default value of a record variable.
   raise info 'employer2 name: % ,epno: %', employer2.name, employer2.epno;
    -- Use a record variable in the FOR loop.
   for employer in select ename, empno from emp_rec order by 1 limit 1
         raise info 'employer name: % , epno: %', employer.name, employer.epno;
      end loop;
    -- Use a record variable in the SELECT INTO statement.
   select ename,empno into employer2 from emp_rec order by 1 limit 1;
   raise info 'employer name: %, epno: %', employer2.name, employer2.epno;
    -- Use a record variable in a cursor.
   OPEN C1;
   FETCH C1 INTO employer2;
   raise info 'employer name: %, epno: %', employer2.name, employer2.epno;
   CLOSE C1;
   RETURN employer.name;
END;
LANGUAGE plpgsql;
-- Invoke the stored procedure.
CALL regress_record('abc');
INFO: employer name: WARD, epno:18
INFO: employer1 name: WARD, epno: 18
INFO: employer1 name: <NULL> , epno: <NULL>
INFO: employer2 name: SCOTT ,epno: 10
-- Delete the stored procedure.
DROP PROCEDURE regress_record;
```

9.4 GaussDB(DWS) Stored Procedure Declaration Syntax

Basic Structure

A PL/SQL block can contain a sub-block which can be placed in any section. The following describes the architecture of a PL/SQL block:

 DECLARE: declares variables, types, cursors, and regional stored procedures and functions used in the PL/SQL block.

□ NOTE

This part is optional if no variable needs to be declared.

- An anonymous block may omit the DECLARE keyword if no variable needs to be declared.
- For a stored procedure, **AS** is used, which is equivalent to **DECLARE**. The **AS** keyword must be reserved even if there is no variable declaration part.
- **EXECUTION**: specifies procedure and SQL statements. It is the main part of a program. It is mandatory.

 BEGIN
- **EXCEPTION**: processes errors. It is optional.

EXCEPTION

• END:

NOTICE

You are not allowed to use consecutive tabs in the PL/SQL block, because they may result in an exception when the parameter **-r** is executed using the **gsql** tool.

PL/SQL Block Classification

PL/SQL blocks are classified into the following types:

- Anonymous block: a dynamic block that can be executed only for once. For details about the syntax, see **Anonymous Block**.
- Subprogram: a stored procedure, function, operator, or packages stored in a database. A subprogram created in a database can be called by other programs.

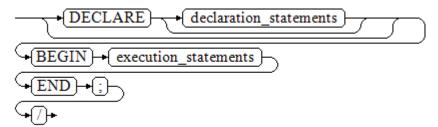
Anonymous Block

An anonymous block applies to a script infrequently executed or a one-off activity. An anonymous block is executed in a session and is not stored.

Syntax

Figure 9-2 shows the syntax diagrams for an anonymous block.

Figure 9-2 anonymous_block::=



Details about the syntax diagram are as follows:

• The execute part of an anonymous block starts with a **BEGIN** statement, has a break with an **END** statement, and ends with a semicolon (;). Type a slash (/) and press **Enter** to execute the statement.

NOTICE

The terminator "/" must be written in an independent row.

- The declaration section includes the variable definition, type, and cursor definition.
- A simplest anonymous block does not execute any commands. At least one statement, even a null statement, must be presented in any implementation blocks.

Examples

The following lists basic anonymous block programs:

```
-- Null statement block:

BEGIN
NULL;
END;
/-- Print information to the console:

BEGIN
dbms_output.put_line('hello world!');
END;
/-- Print variable contents to the console:

DECLARE
my_var VARCHAR2(30);

BEGIN
my_var :='world';
dbms_output.put_line('hello'||my_var);
END;
/-- END;
```

Subprogram

A subprogram stores stored procedures, functions, operators, and advanced packages. A subprogram created in a database can be called by other programs.

9.5 Basic Statements of GaussDB(DWS) Stored Procedures

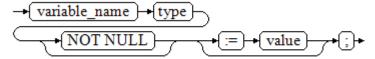
Variable Definition Statement

This section describes the declaration of variables in the PL/SQL and the scope of this variable in codes.

Variable declaration

For details about the variable declaration syntax, see Figure 9-3.

Figure 9-3 declare_variable::=



The syntax is described as follows:

- variable name indicates the name of a variable.
- **type** indicates the type of a variable.
- **value** indicates the initial value of the variable. (If the initial value is not given, **NULL** is taken as the initial value.) **value** can also be an expression.

Examples

```
DECLARE
emp_id INTEGER := 7788; -- Define a variable and assign a value to it.

BEGIN
emp_id := 5*7784; -- Assign a value to the variable.

END;
/
```

In addition to the declaration of basic variable types, **%TYPE** and **%ROWTYPE** can be used to declare variables related to table columns or table structures.

%TYPE attribute

%TYPE declares a variable to be of the same data type as a previously declared variable (for example, a column in a table). For example, if you want to define a **my_name** variable that has the same data type as the **firstname** column in the **employee** table, you can define the variable as follows:

my_name employee.firstname%TYPE

In this way, you can declare **my_name** even if you do not know the data type of **firstname** in **employee**, and the data type of **my_name** can be automatically updated when the data type of **firstname** changes.

%ROWTYPE attribute

%ROWTYPE declares data types of a set of data. It stores a row of table data or results fetched from a cursor. For example, if you want to define a set of data with the same column names and column data types as the **employee** table, you can define the data as follows:

my_employee employee%ROWTYPE

NOTICE

If multiple CNs are used, the **%ROWTYPE** and **%TYPE** attributes of temporary tables cannot be declared in a stored procedure, because a temporary table is valid only in the current session and is invisible to other CNs in the compilation phase. In this case, a message is displayed indicating that the temporary table does not exist.

Variable scope

The scope of a variable indicates the accessibility and availability of a variable in code block. In other words, a variable takes effect only within its scope.

- To define a function scope, a variable must declare and create a BEGIN-END block in the declaration section. The necessity of such declaration is also determined by block structure, which requires that a variable has different scopes and lifetime during a process.
- A variable can be defined multiple times in different scopes, and inner definition can cover outer one.
- A variable defined in an outer block can also be used in a nested block. However, the outer block cannot access variables in the nested block.

Examples

```
DECLARE
emp_id INTEGER :=7788; -- Define a variable and assign a value to it.
outer_var INTEGER :=6688; -- Define a variable and assign a value to it.

BEGIN
DECLARE
emp_id INTEGER :=7799; -- Define a variable and assign a value to it.
inner_var INTEGER :=6688; -- Define a variable and assign a value to it.

BEGIN
dbms_output.put_line('inner emp_id ='||emp_id); -- Display the value as 7799.
dbms_output.put_line('outer_var ='||outer_var); -- Cite variables of an outer block.
END;
dbms_output.put_line('outer emp_id ='||emp_id); -- Display the value as 7788.
END;
//
```

Assignment Statement

Syntax

Figure 9-4 shows the syntax diagram for assigning a value to a variable.

Figure 9-4 assignment_value::=

```
→ (variable_name) → (:=) → (value) → (;) →
```

The syntax is described as follows:

- variable_name indicates the name of a variable.
- **value** can be a value or an expression. The type of **value** must be compatible with the type of **variable name**.

Examples

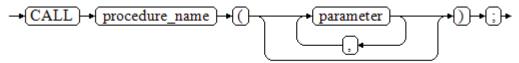
```
DECLARE
emp_id INTEGER := 7788; --Assignment
BEGIN
emp_id := 5; --Assignment
emp_id := 5*7784;
END;
/
```

Call Statement

Syntax

Figure 9-5 shows the syntax diagram for calling a clause.

Figure 9-5 call_clause::=



The syntax is described as follows:

- **procedure_name** specifies the name of a stored procedure.
- **parameter** specifies the parameters for the stored procedure. You can set no parameter or multiple parameters.

Examples

```
-- Create the stored procedure proc_staffs:
CREATE OR REPLACE PROCEDURE proc_staffs
section NUMBER(6),
salary_sum out NUMBER(8,2),
staffs_count out INTEGER
IS
BEGIN
SELECT sum(salary), count(*) INTO salary_sum, staffs_count FROM staffs where section_id = section;
-- Create the stored procedure proc_return:
CREATE OR REPLACE PROCEDURE proc_return
v_num NUMBER(8,2);
v_sum INTEGER;
BEGIN
proc_staffs(30, v_sum, v_num); --Invoke a statement:
dbms_output.put_line(v_sum||'#'||v_num);
RETURN; -- Return a statement.
END;
-- Invoke a stored procedure proc_return:
CALL proc_return();
-- Delete a stored procedure:
DROP PROCEDURE proc_staffs;
DROP PROCEDURE proc_return;
-- Create the function func_return.
CREATE OR REPLACE FUNCTION func_return returns void
language plpgsql
AS $$
DECLARE
v_num INTEGER := 1;
BEGIN
dbms_output.put_line(v_num);
RETURN; -- Return a statement.
END $$;
```

```
-- Invoke the function func_return.

CALL func_return();

1

-- Delete the function:

DROP FUNCTION func_return;
```

9.6 Dynamic Statements of GaussDB(DWS) Stored Procedures

9.6.1 Executing Dynamic Query Statements

You can perform dynamic queries using **EXECUTE IMMEDIATE** or **OPEN FOR** in GaussDB(DWS). **EXECUTE IMMEDIATE** dynamically executes **SELECT** statements and **OPEN FOR** combines use of cursors. If you need to store query results in a data set, use **OPEN FOR**.

EXECUTE IMMEDIATE

Figure 9-6 shows the syntax diagram.

Figure 9-6 EXECUTE IMMEDIATE dynamic_select_clause::=

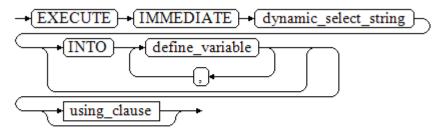
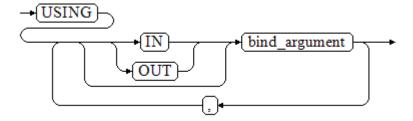


Figure 9-7 shows the syntax diagram for **using_clause**.

Figure 9-7 using_clause-1



The above syntax diagram is explained as follows:

define_variable: specifies variables to store single-line query results.

- **USING IN bind_argument**: specifies where the variable passed to the dynamic SQL value is stored, that is, in the dynamic placeholder of **dynamic select string**.
- **USING OUT bind_argument**: specifies where the dynamic SQL returns the value of the variable.

NOTICE

- In query statements, INTO and OUT cannot coexist.
- A placeholder name starts with a colon (:) followed by digits, characters, or strings, corresponding to bind_argument in the USING clause.
- bind_argument can only be a value, variable, or expression. It cannot be a
 database object such as a table name, column name, and data type. That
 is, bind_argument cannot be used to transfer schema objects for dynamic
 SQL statements. If a stored procedure needs to transfer database objects
 through bind_argument to construct dynamic SQL statements (generally,
 DDL statements), you are advised to use double vertical bars (||) to
 concatenate dynamic_select_clause with a database object.
- A dynamic PL/SQL block allows duplicate placeholders. That is, a
 placeholder can correspond to only one bind_argument in the USING
 clause.

Example

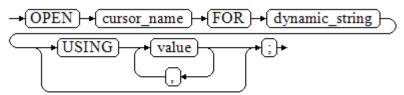
```
--Retrieve values from dynamic statements (INTO clause).
DECLARE
 staff_count VARCHAR2(20);
BEGIN
 EXECUTE IMMEDIATE 'select count(*) from staffs'
   INTO staff_count;
 dbms output.put line(staff count);
END:
-- Pass and retrieve values (the INTO clause is used before the USING clause).
CREATE OR REPLACE PROCEDURE dynamic_proc
 staff id NUMBER(6) := 200:
 first_name VARCHAR2(20);
 salary
           NUMBER(8,2);
BEGIN
 EXECUTE IMMEDIATE 'select first_name, salary from staffs where staff_id = :1'
    INTO first_name, salary
    USING IN staff_id;
 dbms_output.put_line(first_name || ' ' || salary);
END:
-- Invoke the stored procedure.
CALL dynamic_proc();
-- Delete the stored procedure.
DROP PROCEDURE dynamic_proc;
```

OPEN FOR

Dynamic query statements can be executed by using **OPEN FOR** to open dynamic cursors.

For details about the syntax, see Figure 9-8.

Figure 9-8 open_for::=



Parameter description:

- **cursor_name**: specifies the name of the cursor to be opened.
- **dynamic_string**: specifies the dynamic query statement.
- **USING** *value*: applies when a placeholder exists in dynamic_string.

For use of cursors, see GaussDB(DWS) Stored Procedure Cursor.

Example

```
DECLARE
              VARCHAR2(20);
  name
  phone_number VARCHAR2(20);
             NUMBER(8,2);
  salary
  sqlstr
            VARCHAR2(1024);
  TYPE app_ref_cur_type IS REF CURSOR; -- Define the cursor type.
  my_cur app_ref_cur_type; -- Define the cursor variable.
BEGIN
  sqlstr := 'select first_name,phone_number,salary from staffs
     where section_id = :1';
  OPEN my_cur FOR sqlstr USING '30'; -- Open the cursor. using is optional.
  FETCH my_cur INTO name, phone_number, salary; -- Retrieve the data.
  WHILE my cur%FOUND LOOP
      dbms_output.put_line(name||'#'||phone_number||'#'||salary);
      FETCH my_cur INTO name, phone_number, salary;
  END LOOP:
  CLOSE my_cur; -- Close the cursor.
END;
```

9.6.2 Executing Dynamic Non-query Statements

Syntax

Figure 9-9 shows the syntax diagram.

Figure 9-9 noselect::=

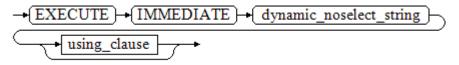
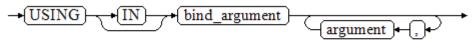


Figure 9-10 shows the syntax diagram for using_clause.

Figure 9-10 using_clause-2



The above syntax diagram is explained as follows:

USING IN bind_argument is used to specify the variable that transfers values to dynamic SQL statements. It is used when a placeholder exists in **dynamic_noselect_string**. That is, a placeholder is replaced by the corresponding *bind_argument* when a dynamic SQL statement is executed. Note that *bind_argument* can only be a value, variable, or expression, and cannot be a database object such as a table name, column name, and data type. If a stored procedure needs to transfer database objects through *bind_argument* to construct dynamic SQL statements (generally, DDL statements), you are advised to use double vertical bars (||) to concatenate *dynamic_select_clause* with a database object. In addition, a dynamic PL/SQL block allows duplicate placeholders. That is, a placeholder can correspond to only one *bind_argument*.

Examples

```
-- Create a table:
CREATE TABLE sections_t1
            NUMBER(4)
 section
 section_name VARCHAR2(30),
 manager_id NUMBER(6),
 place_id
            NUMBER(4)
DISTRIBUTE BY hash(manager_id);
--Declare a variable:
DECLARE
            NUMBER(4) := 280:
 section
 section_name VARCHAR2(30) := 'Info support';
 manager_id NUMBER(6) := 103;
           NUMBER(4) := 1400;
 place_id
 BEGIN

    Execute the auery:

  EXECUTE IMMEDIATE 'insert into sections_t1 values(:1, :2, :3, :4)'
    USING section, section_name, manager_id,place_id;
-- Execute the guery (duplicate placeholders):
  EXECUTE IMMEDIATE 'insert into sections_t1 values(:1, :2, :3, :1)'
    USING section, section_name, manager_id;
-- Run the ALTER statement. (You are advised to use double vertical bars (||) to concatenate the dynamic
DDL statement with a database object.)
  EXECUTE IMMEDIATE 'alter table sections_t1 rename section_name to ' || new_colname;
END;
-- Query data:
SELECT * FROM sections t1;
--Delete the table.
DROP TABLE sections_t1;
```

9.6.3 Dynamically Calling Stored Procedures

This section describes how to dynamically call store procedures. You must use anonymous statement blocks to package stored procedures or statement blocks

and append **IN** and **OUT** behind the **EXECUTE IMMEDIATE...USING** statement to input and output parameters.

Syntax

Figure 9-11 shows the syntax diagram.

Figure 9-11 call_procedure::=

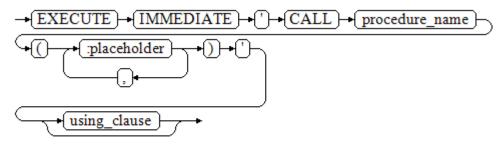
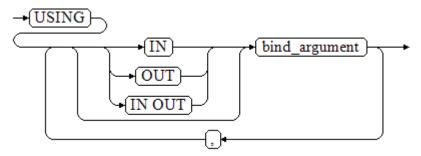


Figure 9-12 shows the syntax diagram for using_clause.

Figure 9-12 using_clause-3



The above syntax diagram is explained as follows:

- **CALL procedure_name**: calls the stored procedure.
- [:placeholder1,:placeholder2, ...]: specifies the placeholder list of the stored procedure parameters. The number of placeholders is the same as the number of parameters. A placeholder name starts with a colon (:) or dollar sign (\$). The colon (:) can be followed by digits, characters, or character strings (excluding digits, characters, or character strings with quotation marks). The dollar sign (\$) can be followed only by digits. A placeholder can correspond to only one bind_argument in the USING clause.
- **USING** [IN|OUT|IN OUT]bind_argument: specifies where the variable passed to the stored procedure parameter value is stored. The modifiers in front of bind_argument and of the corresponding parameter are the same.

```
--Create the stored procedure proc_add.

CREATE OR REPLACE PROCEDURE proc_add

(
```

```
param1 in INTEGER,
  param2 out INTEGER,
  param3 in INTEGER
AS
BEGIN
 param2:= param1 + param3;
END;
DECLARE
  input1 INTEGER:=1;
  input2 INTEGER:=2;
  statement VARCHAR2(200);
  param2
           INTEGER;
BEGIN
  --Declare the call statement.
  statement := 'call proc_add(:col_1, :col_2, :col_3)';(or statement := 'call proc_add($1, $2, $3)';)
  -- Execute the statement.
  EXECUTE IMMEDIATE statement
    USING IN input1, OUT param2, IN input2;
  dbms_output.put_line('result is: '||to_char(param2));
END;
--Delete the stored procedure.
DROP PROCEDURE proc_add;
```

9.6.4 Dynamically Calling Anonymous Blocks

This section describes how to execute anonymous blocks in dynamic statements. Append **IN** and **OUT** behind the **EXECUTE IMMEDIATE...USING** statement to input and output parameters.

Syntax

Figure 9-13 shows the syntax diagram.

Figure 9-13 call_anonymous_block::=

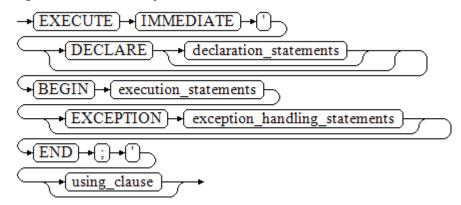
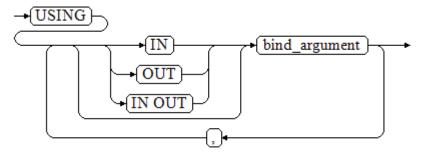


Figure 9-14 shows the syntax diagram for using_clause.

Figure 9-14 using_clause-4



The above syntax diagram is explained as follows:

- The execute part of an anonymous block starts with a **BEGIN** statement, has a break with an **END** statement, and ends with a semicolon (;).
- **USING [IN|OUT|IN OUT]bind_argument**: specifies where the variable passed to the stored procedure parameter value is stored. The modifiers in front of **bind_argument** and of the corresponding parameter are the same.
- The input and output parameters in the middle of an anonymous block are designated by placeholders. The numbers of the placeholders and the parameters are the same. The sequences of the parameters corresponding to the placeholders and the USING parameters are the same.
- Currently in GaussDB(DWS), when dynamic statements call anonymous blocks, placeholders cannot be used to pass input and output parameters in an EXCEPTION statement.

Example

```
-- Create the stored procedure dynamic_proc.
CREATE OR REPLACE PROCEDURE dynamic proc
 staff_id NUMBER(6) := 200;
 first_name VARCHAR2(20);
 salary
            NUMBER(8,2);
BEGIN
--Execute the anonymous block.
  EXECUTE IMMEDIATE 'begin select first_name, salary into :first_name, :salary from staffs where
staff_id= :dno; end;'
    USING OUT first_name, OUT salary, IN staff_id;
 dbms\_output\_line(first\_name || \ ' \ ' \ || \ salary);
END;
-- Invoke the stored procedure.
CALL dynamic_proc();
-- Delete the stored procedure.
DROP PROCEDURE dynamic_proc;
```

9.7 GaussDB(DWS) Stored Procedure Control Statements

9.7.1 RETURN Statements

GaussDB(DWS) provides two methods for returning data: **RETURN** (or **RETURN NEXT**) and **RETURN QUERY**. **RETURN NEXT** and **RETURN QUERY** are used only for functions and cannot be used for stored procedures.

RETURN

Syntax

Figure 9-15 shows the syntax of a return statement.

Figure 9-15 return_clause::=



The syntax is explained as follows:

This statement returns control from a stored procedure or function to a caller.

```
-- Create the stored procedure proc staffs:
CREATE OR REPLACE PROCEDURE proc_staffs
section NUMBER(6),
salary_sum out NUMBER(8,2),
staffs_count out INTEGER
İS
BEGIN
SELECT sum(salary), count(*) INTO salary_sum, staffs_count FROM staffs where section_id = section;
END:
-- Create the stored procedure proc_return:
CREATE OR REPLACE PROCEDURE proc_return
AS
v_num NUMBER(8,2);
v sum INTEGER;
BEGIN
proc_staffs(30, v_sum, v_num); --Call a statement.
dbms\_output.put\_line(v\_sum||'\#'||v\_num);
RETURN; -- Return a statement.
END;
-- Invoke a stored procedure proc_return:
CALL proc_return();
-- Delete a stored procedure:
DROP PROCEDURE proc_staffs;
DROP PROCEDURE proc_return;
--Create the function func_return.
CREATE OR REPLACE FUNCTION func return returns void
language plpgsql
AS $$
DECLARE
v_num INTEGER := 1;
BEGIN
dbms_output.put_line(v_num);
RETURN; -- Return a statement.
```

```
END $$;

-- Invoke the function func_return.

CALL func_return();

1

-- Delete the function:

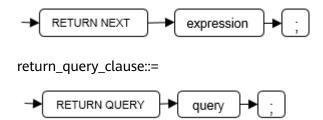
DROP FUNCTION func_return;
```

RETURN NEXT and RETURN QUERY

Syntax

When creating a function, specify **SETOF datatype** for the return values.

return_next_clause::=



The syntax is explained as follows:

If a function needs to return a result set, use **RETURN NEXT** or **RETURN QUERY** to add results to the result set, and then continue to execute the next statement of the function. As the **RETURN NEXT** or **RETURN QUERY** statement is executed repeatedly, more and more results will be added to the result set. After the function is executed, all results are returned.

RETURN NEXT can be used for scalar and compound data types.

RETURN QUERY has a variant **RETURN QUERY EXECUTE**. You can add dynamic queries and add parameters to the queries by using **USING**.

```
CREATE TABLE t1(a int);
INSERT INTO t1 VALUES(1),(10);
-- RETURN NEXT
CREATE OR REPLACE FUNCTION fun_for_return_next() RETURNS SETOF t1 AS $$
DECLARE
 r t1%ROWTYPE;
BEGIN
 FOR r IN select * from t1
 LOOP
   RETURN NEXT r;
 END LOOP;
 RETURN;
END;
$$ LANGUAGE PLPGSQL;
call fun_for_return_next();
а
1
10
(2 rows)
-- RETURN QUERY
```

```
CREATE OR REPLACE FUNCTION fun_for_return_query() RETURNS SETOF t1 AS $

DECLARE
    r t1%ROWTYPE;

BEGIN
    RETURN QUERY select * from t1;

END;

$$

language plpgsql;

call fun_for_return_next();

a
---

1
10
(2 rows)
```

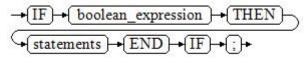
9.7.2 Conditional Statements

Conditional statements are used to decide whether given conditions are met. Operations are executed based on the decisions made.

GaussDB(DWS) supports five usages of IF:

IF_THEN

Figure 9-16 IF_THEN::=



IF_THEN is the simplest form of **IF**. If the condition is true, statements are executed. If it is false, they are skipped.

Examples

```
IF v_user_id <> 0 THEN

UPDATE users SET email = v_email WHERE user_id = v_user_id;

END IF;
```

IF THEN ELSE

Figure 9-17 IF_THEN_ELSE::=

```
→(IF)→(boolean_expression)→(THEN)

→(statements)→(ELSE)→(statements)→(END)→(IF)→(;)→
```

IF-THEN-ELSE statements add **ELSE** branches and can be executed if the condition is **false**.

Examples

```
IF parentid IS NULL OR parentid = "
THEN
RETURN;
ELSE
hp_true_filename(parentid); -- Call the stored procedure.
END IF;
```

IF THEN ELSE IF

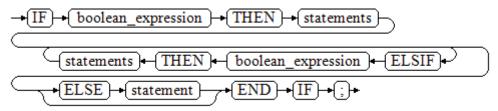
IF statements can be nested in the following way:

```
IF gender = 'm' THEN
    pretty_gender := 'man';
ELSE
    IF gender = 'f' THEN
        pretty_gender := 'woman';
    END IF;
END IF;
```

Actually, this is a way of an **IF** statement nesting in the **ELSE** part of another **IF** statement. Therefore, an **END IF** statement is required for each nesting **IF** statement and another **END IF** statement is required to end the parent **IF**-**ELSE** statement. To set multiple options, use the following form:

• IF_THEN_ELSIF_ELSE

Figure 9-18 IF_THEN_ELSIF_ELSE::=



Examples

```
IF number_tmp = 0 THEN
    result := 'zero';
ELSIF number_tmp > 0 THEN
    result := 'positive';
ELSIF number_tmp < 0 THEN
    result := 'negative';
ELSE
    result := 'NULL';
END IF;</pre>
```

IF_THEN_ELSEIF_ELSE

ELSEIF is an alias of ELSIF.

```
CREATE OR REPLACE PROCEDURE proc_control_structure(i in integer)

AS

BEGIN

IF i > 0 THEN

raise info 'i:% is greater than 0. ',i;

ELSIF i < 0 THEN

raise info 'i:% is smaller than 0. ',i;

ELSE

raise info 'i:% is equal to 0. ',i;

END IF;

RETURN;

END;

/

CALL proc_control_structure(3);

-- Delete the stored procedure.

DROP PROCEDURE proc_control_structure;
```

9.7.3 Loop Statements

Simple LOOP Statements

The syntax diagram is as follows.

Figure 9-19 loop::=

```
\rightarrow LOOP \rightarrow statements \rightarrow END \rightarrow LOOP \rightarrow (;) \rightarrow
```

Examples

```
CREATE OR REPLACE PROCEDURE proc_loop(i in integer, count out integer)

AS

BEGIN

count:=0;
LOOP

IF count > i THEN

raise info 'count is %. ', count;
EXIT;
ELSE

count:=count+1;
END IF;
END LOOP;
END;

/

CALL proc_loop(10,5);
```

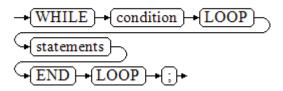
NOTICE

The loop must be exploited together with **EXIT**; otherwise, a dead loop occurs.

WHILE-LOOP Statements

The syntax diagram is as follows.

Figure 9-20 while_loop::=



If the conditional expression is true, a series of statements in the **WHILE** statement are repeatedly executed and the condition is decided each time the loop body is executed.

Examples

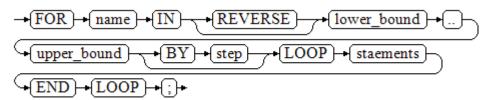
```
CREATE TABLE integertable(c1 integer) DISTRIBUTE BY hash(c1);
CREATE OR REPLACE PROCEDURE proc_while_loop(maxval in integer)
AS

DECLARE
i int :=1;
BEGIN
WHILE i < maxval LOOP
INSERT INTO integertable VALUES(i);
i:=i+1;
END LOOP;
END;
/
-- Invoke a function:
CALL proc_while_loop(10);
-- Delete the stored procedure and table:
DROP PROCEDURE proc_while_loop;
DROP TABLE integertable;
```

FOR_LOOP (Integer variable) Statement

The syntax diagram is as follows.

Figure 9-21 for_loop::=



□ NOTE

- The variable **name** is automatically defined as the **integer** type and exists only in this loop. The variable name falls between lower_bound and upper_bound.
- When the keyword **REVERSE** is used, the lower bound must be greater than or equal to the upper bound; otherwise, the loop body is not executed.

```
-- Loop from 0 to 5:

CREATE OR REPLACE PROCEDURE proc_for_loop()

AS

BEGIN

FOR I IN 0..5 LOOP

DBMS_OUTPUT.PUT_LINE('It is '||to_char(I) || ' time;');

END LOOP;

END;

-- Invoke a function:

CALL proc_for_loop();

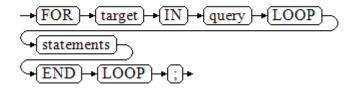
-- Delete the stored procedure:

DROP PROCEDURE proc_for_loop;
```

FOR_LOOP Query Statements

The syntax diagram is as follows.

Figure 9-22 for_loop_query::=



The variable **target** is automatically defined, its type is the same as that in the **query** result, and it is valid only in this loop. The target value is the query result.

Examples

```
-- Display the query result from the loop:

CREATE OR REPLACE PROCEDURE proc_for_loop_query()

AS
    record VARCHAR2(50);

BEGIN
    FOR record IN SELECT spcname FROM pg_tablespace LOOP
    dbms_output.put_line(record);
    END LOOP;

END;
/

-- Invoke a function.

CALL proc_for_loop_query();

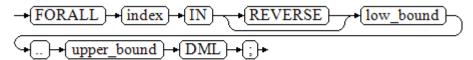
-- Delete the stored procedure.

DROP PROCEDURE proc_for_loop_query;
```

FORALL Batch Query Statements

The syntax diagram is as follows.

Figure 9-23 forall::=



□ NOTE

The variable **index** is automatically defined as the **integer** type and exists only in this loop. The index value falls between low_bound and upper_bound.

```
CREATE TABLE hdfs_t1 (
title NUMBER(6),
```

```
did VARCHAR2(20),
 data_period VARCHAR2(25),
 kind VARCHAR2(25),
 interval VARCHAR2(20),
 time DATE,
 isModified VARCHAR2(10)
DISTRIBUTE BY hash(did);
INSERT INTO hdfs_t1 VALUES( 8, 'Donald', 'OConnell', 'DOCONNEL', '650.507.9833', to_date('21-06-1999',
'dd-mm-yyyy'), 'SH_CLERK' );
CREATE OR REPLACE PROCEDURE proc_forall()
AS
BEGIN
  FORALL i IN 100..120
     insert into hdfs_t1(title) values(i);
END;
-- Invoke a function:
CALL proc_forall();
-- Query the invocation result of the stored procedure:
SELECT * FROM hdfs_t1 WHERE title BETWEEN 100 AND 120;
-- Delete the stored procedure and table:
DROP PROCEDURE proc_forall;
DROP TABLE hdfs_t1;
```

9.7.4 Branch Statements

Syntax

Figure 9-24 shows the syntax diagram.



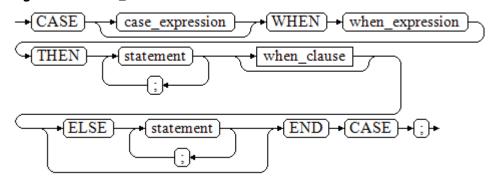
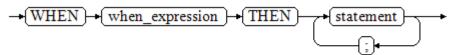


Figure 9-25 shows the syntax diagram for when_clause.

Figure 9-25 when_clause::=



Parameter description:

- case_expression: specifies the variable or expression.
- when_expression: specifies the constant or conditional expression.
- **statement**: specifies the statement to execute.

Examples

```
CREATE OR REPLACE PROCEDURE proc_case_branch(pi_result in integer, pi_return out integer)
  BEGIN
     CASE pi_result
       WHEN 1 THEN
          pi_return := 111;
       WHEN 2 THEN
          pi_return := 222;
       WHEN 3 THEN
          pi_return := 333;
       WHEN 6 THEN
          pi return := 444;
       WHEN 7 THEN
          pi_return := 555;
       WHEN 8 THEN
          pi_return := 666;
       WHEN 9 THEN
          pi_return := 777;
       WHEN 10 THEN
          pi_return := 888;
       ELSE
          pi_return := 999;
     END CASE;
     raise info 'pi_return : %',pi_return ;
END;
CALL proc_case_branch(3,0);
--Delete the stored procedure.
DROP PROCEDURE proc_case_branch;
```

9.7.5 NULL Statements

In PL/SQL programs, a **NULL** statement can be used to indicate "do nothing", which is also known as an empty statement.

A NULL statement acts as a placeholder and can give meaning to certain statements, improving the readability of the program.

Syntax

The following shows example use of NULL statements.

```
DECLARE
...
BEGIN
...
IF v_num IS NULL THEN
NULL; --No data needs to be processed.
END IF;
END;
```

9.7.6 Error Trapping Statements

By default, any error occurring in a PL/SQL function aborts execution of the function, and indeed of the surrounding transaction as well. You can trap errors

and restore from them by using a **BEGIN** block with an **EXCEPTION** clause. The syntax is an extension of the normal syntax for a **BEGIN** block:

```
[<<label>>]
[DECLARE
declarations]
BEGIN
statements
EXCEPTION
WHEN condition [OR condition ...] THEN
handler_statements
[WHEN condition [OR condition ...] THEN
handler_statements
...]
END;
```

If no error occurs, this form of block simply executes all the statements, and then control passes to the next statement after **END**. But if an error occurs inside the executed statement, the statement rolls back and goes to the EXCEPTION list to find the first condition that matches the error. If a match is found, the corresponding **handler_statements** are executed, and then control passes to the next statement after **END**. If no match is found, the error propagates out as though the **EXCEPTION** clause were not there at all:

The error can be caught by an enclosing block with **EXCEPTION**, or if there is none it aborts processing of the function.

The *condition* can be any of those shown in SQL standard error codes. The special condition name **OTHERS** matches every error type except **QUERY CANCELED**.

If a new error occurs within the selected **handler_statements**, it cannot be caught by this **EXCEPTION** clause, but is propagated out. A surrounding **EXCEPTION** clause could catch it.

When an error is caught by an **EXCEPTION** clause, the local variables of the PL/SQL function remain as they were when the error occurred, but all changes to persistent database state within the block are rolled back.

Example:

```
CREATE TABLE mytab (id INT, firstname VARCHAR(20), lastname VARCHAR(20)) DISTRIBUTE BY hash(id);
INSERT INTO mytab(firstname, lastname) VALUES('Tom', 'Jones');
CREATE FUNCTION fun_exp() RETURNS INT
AS $$
DECLARE
  x INT :=0;
  y INT;
BEGIN
  UPDATE mytab SET firstname = 'Joe' WHERE lastname = 'Jones';
  x := x + 1:
  y := x / 0;
EXCEPTION
  WHEN division_by_zero THEN
     RAISE NOTICE 'caught division_by_zero';
    RETURN x:
END;$$
LANGUAGE plpqsql;
CALL fun_exp();
NOTICE: caught division_by_zero
fun_exp
    1
(1 row)
```

When control reaches the assignment to **y**, it will fail with a **division_by_zero** error. This will be caught by the **EXCEPTION** clause. The value returned in the **RETURN** statement will be the incremented value of **x**.

□ NOTE

A block containing an **EXCEPTION** clause is more expensive to enter and exit than a block without one. Therefore, do not use **EXCEPTION** without need.

In the following scenario, an exception cannot be caught, and the entire transaction rolls back. The threads of the nodes participating the stored procedure exit abnormally due to node failure and network fault, or the source data is inconsistent with that of the table structure of the target table during the COPY FROM operation.

Example: Exceptions with **UPDATE/INSERT**

This example uses exception handling to perform either **UPDATE** or **INSERT**, as appropriate:

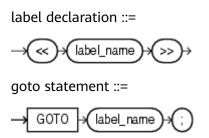
```
CREATE TABLE db (a INT, b TEXT);
CREATE FUNCTION merge_db(key INT, data TEXT) RETURNS VOID AS
BEGIN
  LOOP
-- Try updating the key:
     UPDATE db SET b = data WHERE a = key;
    IF found THEN
       RETURN;
    END IF;
-- Not there, so try to insert the key. If someone else inserts the same key concurrently, we could get a
unique-key failure.
    BEGIN
       INSERT INTO db(a,b) VALUES (key, data);
       RETURN:
    EXCEPTION WHEN unique_violation THEN
     -- Loop to try the UPDATE again:
    FND:
   END LOOP;
END:
LANGUAGE plpqsql;
SELECT merge_db(1, 'david');
SELECT merge_db(1, 'dennis');
-- Delete FUNCTION and TABLE:
DROP FUNCTION merge_db;
DROP TABLE db;
```

9.7.7 GOTO Statements

The **GOTO** statement unconditionally transfers the control from the current statement to a labeled statement. The **GOTO** statement changes the execution logic. Therefore, use this statement only when necessary. Alternatively, you can use

the **EXCEPTION** statement to handle issues in special scenarios. To run the **GOTO** statement, the labeled statement must be unique.

Syntax



Examples

```
CREATE OR REPLACE PROCEDURE GOTO_test()
AS
DECLARE
  v1 int;
BEGIN
  v1 := 0;
    LOOP
    EXIT WHEN v1 > 100;
         v1 := v1 + 2;
          if v1 > 25 THEN
              GOTO pos1;
          END IF;
    END LOOP;
<<pos1>>
v1 := v1 + 10;
raise info 'v1 is %. ', v1;
END;
call GOTO_test();
DROP PROCEDURE GOTO_test();
```

Constraints

The **GOTO** statement has the following constraints:

• The **GOTO** statement does not allow multiple labeled statements even if they are in different blocks.

```
BEGIN
GOTO pos1;
<<pos1>>
SELECT * FROM ...
<<pos1>>
UPDATE t1 SET ...
END;
```

 The GOTO statement cannot transfer control to the IF, CASE, or LOOP statement.

```
BEGIN

GOTO pos1;

If valid THEN

<<pos1>>

SELECT * FROM ...

END IF;

END;
```

• The **GOTO** statement cannot transfer control from one **IF** clause to another, or from one **WHEN** clause in the **CASE** statement to another.

```
BEGIN

IF valid THEN

GOTO pos1;

SELECT * FROM ...

ELSE

<<pos1>>

UPDATE t1 SET ...

END IF;

END;
```

 The GOTO statement cannot transfer control from an outer block to an inner BEGIN-END block.

```
BEGIN
GOTO pos1;
BEGIN
<<pos1>>
UPDATE t1 SET ...
END;
END;
```

 The GOTO statement cannot transfer control from an EXCEPTION block to the current BEGIN-END block but can transfer to an outer BEGIN-END block.

```
BEGIN

<pos1>>

UPDATE t1 SET ...

EXCEPTION

WHEN condition THEN

GOTO pos1;

END;
```

• If the labeled statement in the **GOTO** statement does not exist, you need to add the **NULL** statement.

```
DECLARE
done BOOLEAN;
BEGIN

FOR i IN 1..50 LOOP

IF done THEN

GOTO end_loop;

END IF;

<<end_loop>> -- not allowed unless an executable statement follows

NULL; -- add NULL statement to avoid error

END LOOP; -- raises an error without the previous NULL

END;

/
```

9.8 Other Statements in a GaussDB(DWS) Stored Procedure

Lock Operations

GaussDB(DWS) provides multiple lock modes to control concurrent accesses to table data. These modes are used when Multi-Version Concurrency Control (MVCC) cannot give expected behaviors. Alike, most GaussDB(DWS) commands automatically apply appropriate locks to ensure that called tables are not deleted or modified in an incompatible manner during command execution. For example, when concurrent operations exist, **ALTER TABLE** cannot be executed on the same table.

Cursor Operations

GaussDB(DWS) provides cursors as a data buffer for users to store execution results of SQL statements. Each cursor region has a name. Users can use SQL

statements to obtain records one by one from cursors and grant them to master variables, then being processed further by host languages.

Cursor operations include cursor definition, open, fetch, and close operations.

For the complete example of cursor operations, see **Explicit Cursor**.

9.9 GaussDB(DWS) Stored Procedure Cursor

9.9.1 Overview

To process SQL statements, the stored procedure process assigns a memory segment to store context association. Cursors are handles or pointers to context areas. With cursors, stored procedures can control alterations in context areas.

NOTICE

If JDBC is used to call a stored procedure whose returned value is a cursor, the returned cursor is not available.

Cursors are classified into explicit cursors and implicit cursors. **Table 9-2** shows the usage conditions of explicit and implicit cursors for different SQL statements.

Table 9-2 Cursor usage conditions

SQL Statement	Cursor
Non-query statements	Implicit
Query statements with single-line results	Implicit or explicit
Query statements with multi-line results	Explicit

9.9.2 Explicit Cursor

An explicit cursor is used to process query statements, particularly when the query results contain multiple records.

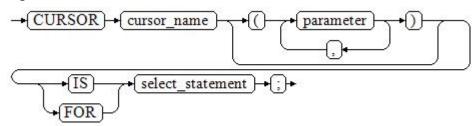
Procedure

An explicit cursor performs the following six PL/SQL steps to process query statements:

Step 1 Define a static cursor: Define a cursor name and its corresponding **SELECT** statement.

Figure 9-26 shows the syntax diagram for defining a static cursor.

Figure 9-26 static_cursor_define::=



Parameter description:

- cursor_name: defines a cursor name.
- **parameter**: specifies cursor parameters. Only input parameters are allowed in the following format:

 parameter_name datatype
- **select_statement**: specifies a query statement.

□ NOTE

The system automatically determines whether the cursor can be used for backward fetches based on the execution plan.

Define a dynamic cursor: Define a **ref** cursor, which means that the cursor can be opened dynamically by a set of static SQL statements. Define the type of the **ref** cursor first and then the cursor variable of this cursor type. Dynamically bind a **SELECT** statement through **OPEN FOR** when the cursor is opened.

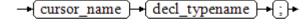
Figure 9-27 and **Figure 9-28** show the syntax diagrams for defining a dynamic cursor.

Figure 9-27 cursor_typename::=

$$\rightarrow$$
 (TYPE) \rightarrow (decl_typename) \rightarrow (IS) \rightarrow (REF) \rightarrow (CURSOR) \rightarrow (;) \rightarrow

GaussDB(DWS) supports the dynamic cursor type **sys_refcursor**. A function or stored procedure can use the **sys_refcursor** parameter to pass on or pass out the cursor result set. A function can return **sys_refcursor** to return the cursor result set.

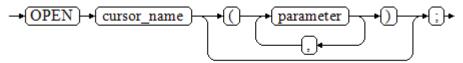
Figure 9-28 dynamic_cursor_define::=



Step 2 Open the static cursor: Execute the **SELECT** statement corresponding to the cursor. The query result is placed in the work area and the pointer directs to the head of the work area to identify the cursor result set. If the cursor query statement contains the **FOR UPDATE** option, the **OPEN** statement locks the data row corresponding to the cursor result set in the database table.

Figure 9-29 shows the syntax diagram for opening a static cursor.

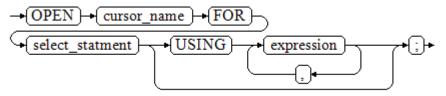
Figure 9-29 open_static_cursor::=



Open the dynamic cursor: Use the **OPEN FOR** statement to open the dynamic cursor and the SQL statement is dynamically bound.

Figure 9-30 shows the syntax diagram for opening a dynamic cursor.

Figure 9-30 open_dynamic_cursor::=

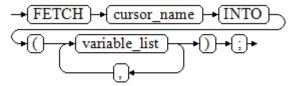


A PL/SQL program cannot use the **OPEN** statement to repeatedly open a cursor.

Step 3 Fetch cursor data: Retrieve data rows in the result set and place them in specified output variables.

Figure 9-31 shows the syntax diagram for fetching cursor data.

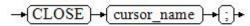
Figure 9-31 fetch_cursor::=



- **Step 4** Process the record.
- **Step 5** Continue to process until the active set has no record.
- **Step 6 Close the cursor**: When fetching and finishing the data in the cursor result set, close the cursor immediately to release system resources used by the cursor and invalidate the work area of the cursor so that the **FETCH** statement cannot be used to fetch data any more. A closed cursor can be reopened using the **OPEN** statement.

Figure 9-32 shows the syntax diagram for closing a cursor.

Figure 9-32 close_cursor::=



----End

Attributes

Cursor attributes are used to control program procedures or learn about program status. When a DML statement is executed, the PL/SQL opens a built-in cursor and processes its result. A cursor is a memory segment for maintaining query results. It is opened when a DML statement is executed and closed when the execution is finished. An explicit cursor has the following attributes:

- **%FOUND**: Boolean attribute, which returns **TRUE** if the last fetch returns a row.
- %NOTFOUND: Boolean attribute, which works opposite to the %FOUND attribute.
- **%ISOPEN**: Boolean attribute, which returns **TRUE** if the cursor has been opened.
- **%ROWCOUNT**: numeric attribute, which returns the number of records fetched from the cursor.

```
-- Specify the method for passing cursor parameters:
CREATE OR REPLACE PROCEDURE cursor_proc1()
DECLARE
  DEPT_NAME VARCHAR(100);
  DEPT_LOC NUMBER(4);
  -- Define a cursor:
  CURSOR C1 IS
    SELECT section_name, place_id FROM sections WHERE section_id <= 50;
  CURSOR C2(sect_id INTEGER) IS
    SELECT section name, place id FROM sections WHERE section id <= sect id;
  TYPE CURSOR_TYPE IS REF CURSOR;
  C3 CURSOR_TYPE;
  SQL_STR VARCHAR(100);
BEGIN
  OPEN C1;-- Open the cursor:
  LOOP
    -- Fetch data from the cursor:
    FETCH C1 INTO DEPT_NAME, DEPT_LOC;
    EXIT WHEN C1%NOTFOUND;
    DBMS_OUTPUT.PUT_LINE(DEPT_NAME||'---'||DEPT_LOC);
  END LOOP;
  CLOSE C1;-- Close the cursor.
  OPEN C2(10);
    FETCH C2 INTO DEPT_NAME, DEPT_LOC;
    EXIT WHEN C2%NOTFOUND:
    DBMS_OUTPUT.PUT_LINE(DEPT_NAME||'---'||DEPT_LOC);
  END LOOP;
  CLOSE C2;
  SQL_STR := 'SELECT section_name, place_id FROM sections WHERE section_id <= :DEPT_NO;';
  OPEN C3 FOR SQL_STR USING 50;
    FETCH C3 INTO DEPT_NAME, DEPT_LOC;
    EXIT WHEN C3%NOTFOUND;
    DBMS_OUTPUT.PUT_LINE(DEPT_NAME||'---'||DEPT_LOC);
  END LOOP;
  CLOSE C3;
END;
CALL cursor_proc1();
```

```
DROP PROCEDURE cursor_proc1;
-- Increase the salary of employees whose salary is lower than CNY3000 by CNY500:
CREATE TABLE staffs_t1 AS TABLE staffs;
CREATE OR REPLACE PROCEDURE cursor_proc2()
DECLARE
 V_EMPNO NUMBER(6);
 V_SAL NUMBER(8,2);
 CURSOR C IS SELECT staff_id, salary FROM staffs_t1;
BEGIN
 OPEN C;
 LOOP
   FETCH C INTO V_EMPNO, V_SAL;
   EXIT WHEN C%NOTFOUND;
   IF V_SAL<=3000 THEN
       UPDATE staffs_t1 SET salary =salary + 500 WHERE staff_id = V_EMPNO;
   END IF;
 END LOOP;
 CLOSE C;
END:
CALL cursor_proc2();
-- Drop the stored procedure:
DROP PROCEDURE cursor_proc2;
DROP TABLE staffs_t1;
-- Use function parameters of the SYS REFCURSOR type:
CREATE OR REPLACE PROCEDURE proc_sys_ref(O OUT SYS_REFCURSOR)
C1 SYS_REFCURSOR;
OPEN C1 FOR SELECT section_ID FROM sections ORDER BY section_ID;
O := C1;
END:
DECLARE
C1 SYS_REFCURSOR;
TEMP NUMBER(4);
BEGIN
proc_sys_ref(C1);
LOOP
 FETCH C1 INTO TEMP;
 DBMS_OUTPUT.PUT_LINE(C1%ROWCOUNT);
 EXIT WHEN C1%NOTFOUND;
END LOOP;
END;
-- Drop the stored procedure:
DROP PROCEDURE proc_sys_ref;
```

9.9.3 Implicit Cursor

The system automatically sets implicit cursors for non-query statements, such as **ALTER** and **DROP**, and creates work areas for these statements. These implicit cursors are named SQL, which is defined by the system.

Overview

Implicit cursor operations, such as definition, opening, value-grant, and closing, are automatically performed by the system. Users can use only the attributes of implicit cursors to complete operations. The data stored in the work area of an

implicit cursor is the latest SQL statement, and is not related to the user-defined explicit cursors.

Format call: SQL%

◯ NOTE

INSERT, **UPDATE**, **DROP**, and **SELECT** statements do not require defined cursors.

Attributes

An implicit cursor has the following attributes:

- SQL%FOUND: Boolean attribute, which returns TRUE if the last fetch returns a row.
- SQL%NOTFOUND: Boolean attribute, which works opposite to the SQL %FOUND attribute.
- **SQL%ROWCOUNT**: numeric attribute, which returns the number of records fetched from the cursor.
- **SQL%ISOPEN**: Boolean attribute, whose value is always **FALSE**. Close implicit cursors immediately after an SQL statement is executed.

Examples

```
-- Delete all employees in a department from the EMP table. If the department has no employees, delete
the department from the DEPT table.
CREATE TABLE staffs_t1 AS TABLE staffs;
CREATE TABLE sections_t1 AS TABLE sections;
CREATE OR REPLACE PROCEDURE proc_cursor3()
  DECLARE
  V_DEPTNO NUMBER(4) := 100;
    DELETE FROM staffs WHERE section_ID = V_DEPTNO;
     -- Proceed based on cursor status:
    IF SQL%NOTFOUND THEN
    DELETE FROM sections_t1 WHERE section_ID = V_DEPTNO;
    END IF;
  END;
CALL proc_cursor3();
-- Drop the stored procedure and the temporary table:
DROP PROCEDURE proc_cursor3;
DROP TABLE staffs_t1;
DROP TABLE sections_t1;
```

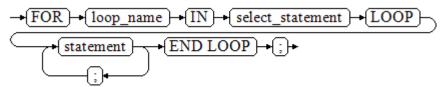
9.9.4 Cursor Loop

The use of cursors in **WHILE** and **LOOP** statements is called a cursor loop. Generally, **OPEN**, **FETCH**, and **CLOSE** statements are needed in cursor loop. The following describes a loop that is applicable to a static cursor loop without executing the four steps of a static cursor.

Syntax

Figure 9-33 shows the syntax diagram for the FOR AS loop.

Figure 9-33 FOR_AS_loop::=



Precautions

- The **UPDATE** operation for the queried table is not allowed in the loop statement.
- The variable *loop_name* is automatically defined and is valid only in this loop. The type and value of *loop_name* are the same as those of the query result of *select statement*.
- The %FOUND, %NOTFOUND, and %ROWCOUNT attributes access the same internal variable in GaussDB(DWS). Transactions and anonymous blocks cannot be accessed by multiple cursors at the same time.

```
BEGIN
FOR ROW_TRANS IN
    SELECT first_name FROM staffs
    DBMS_OUTPUT.PUT_LINE (ROW_TRANS.first_name );
  END LOOP;
END;
-- Create a table:
CREATE TABLE integerTable1( A INTEGER) DISTRIBUTE BY hash(A);
CREATE TABLE integerTable2( B INTEGER) DISTRIBUTE BY hash(B);
INSERT INTO integerTable2 VALUES(2);
-- Multiple cursors share the parameters of cursor attributes:
DECLARE
 CURSOR C1 IS SELECT A FROM integerTable1;--Declare the cursor.
 CURSOR C2 IS SELECT B FROM integerTable2;
 PI A INTEGER:
 PI_B INTEGER;
BFGIN
  OPEN C1;-- Open the cursor.
 OPEN C2;
 FETCH C1 INTO PI_A; ---- The value of C1%FOUND and C2%FOUND is FALSE.
 FETCH C2 INTO PI_B; ---- The value of C1%FOUND and C2%FOUND is TRUE.
-- Determine the cursor status:
 IF C1%FOUND THEN
    IF C2%FOUND THEN
     DBMS_OUTPUT.PUT_LINE('Dual cursor share parameter.');
   END IF;
 END IF;
  CLOSE C1;-- Close the cursor.
 CLOSE C2;
END;
-- Drop the temporary table:
DROP TABLE integerTable1;
DROP TABLE integerTable2;
```

9.10 GaussDB(DWS) Stored Procedure Advanced Package

9.10.1 DBMS_LOB

Related Interfaces

Table 9-3 provides all interfaces supported by the **DBMS_LOB** package.

Table 9-3 DBMS_LOB

API	Description
DBMS_LOB.GETLENGTH	Obtains and returns the specified length of a LOB object.
DBMS_LOB.OPEN	Opens a LOB and returns a LOB descriptor.
DBMS_LOB.READ	Loads a part of LOB contents to BUFFER area according to the specified length and initial position offset.
DBMS_LOB.WRITE	Copies contents in BUFFER area to LOB according to the specified length and initial position offset.
DBMS_LOB.WRITEAPPEN D	Copies contents in BUFFER area to the end part of LOB according to the specified length.
DBMS_LOB.COPY	Copies contents in BLOB to another BLOB according to the specified length and initial position offset.
DBMS_LOB.ERASE	Deletes contents in BLOB according to the specified length and initial position offset.
DBMS_LOB.CLOSE	Closes a LOB descriptor.
DBMS_LOB.INSTR	Returns the position of the Nth occurrence of a character string in LOB.
DBMS_LOB.COMPARE	Compares two LOBs or a certain part of two LOBs.
DBMS_LOB.SUBSTR	Reads the substring of a LOB and returns the number of read bytes or the number of characters.
DBMS_LOB.TRIM	Truncates the LOB of a specified length. After the execution is complete, the length of the LOB is set to the length specified by the newlen parameter.
DBMS_LOB.CREATETEMP ORARY	Creates a temporary BLOB or CLOB.
DBMS_LOB.APPEND	Adds the content of a LOB to another LOB.

DBMS_LOB.GETLENGTH

Specifies the length of a LOB type object obtained and returned by the stored procedure **GETLENGTH**.

The function prototype of **DBMS_LOB.GETLENGTH** is:

```
DBMS_LOB.GETLENGTH (
lob_loc IN BLOB)
RETURN INTEGER;

DBMS_LOB.GETLENGTH (
lob_loc IN CLOB)
RETURN INTEGER;
```

Table 9-4 DBMS_LOB.GETLENGTH interface parameters

Parameter	Description
lob_loc	LOB type object whose length is to be obtained

DBMS_LOB.OPEN

A stored procedure opens a LOB and returns a LOB descriptor. This process is used only for compatibility.

The function prototype of DBMS_LOB.OPEN is:

```
DBMS_LOB.LOB (
lob_loc INOUT BLOB,
open_mode IN BINARY_INTEGER);

DBMS_LOB.LOB (
lob_loc INOUT CLOB,
open_mode IN BINARY_INTEGER);
```

Table 9-5 DBMS_LOB.OPEN interface parameters

Parameter	Description
lob_loc	BLOB or CLOB descriptor that is opened
open_mode IN BINARY_INTEG ER	Open mode (currently, DBMS_LOB.LOB_READWRITE is supported)

DBMS_LOB.READ

The stored procedure **READ** loads a part of LOB contents to BUFFER according to the specified length and initial position offset.

The function prototype of **DBMS_LOB.READ** is:

```
DBMS_LOB.READ (
lob_loc IN BLOB,
amount IN INTEGER,
offset IN INTEGER,
buffer OUT RAW);

DBMS_LOB.READ (
lob_loc IN CLOB,
```

```
amount IN OUT INTEGER,
offset IN INTEGER,
buffer OUT VARCHAR2);
```

Table 9-6 DBMS_LOB.READ interface parameters

Parameter	Description
lob_loc	LOB type object to be loaded
amount	Load data length NOTE If the read length is negative, the error message "ERROR: argument 2 is null, invalid, or out of range." is displayed.
offset	Indicates where to start reading the LOB contents, that is, the offset bytes to initial position of LOB contents.
buffer	Target buffer to store the loaded LOB contents

DBMS_LOB.WRITE

The stored procedure **WRITE** copies contents in BUFFER to LOB variables according to the specified length and initial position offset.

The function prototype of **DBMS_LOB.WRITE** is:

```
DBMS_LOB.WRITE (
lob_loc IN OUT BLOB,
amount IN
              INTEGER,
offset IN
             INTEGER,
buffer IN
             RAW);
DBMS_LOB.WRITE (
               CLOB,
lob_loc IN OUT
              INTEGER,
amount IN
offset IN
             INTEGER,
buffer IN VARCHAR2);
```

Table 9-7 DBMS_LOB.WRITE interface parameters

Parameter	Description
lob_loc	LOB type object to be written
amount	Write data length
	NOTE If the write data is shorter than 1 or longer than the contents to be written, an error is reported.
offset	Indicates where to start writing the LOB contents, that is, the offset bytes to initial position of LOB contents.
	NOTE If the offset is shorter than 1 or longer than the maximum length of LOB type contents, an error is reported.
buffer	Content to be written

DBMS_LOB.WRITEAPPEND

The stored procedure **WRITEAPPEND** copies contents in BUFFER to the end part of LOB according to the specified length.

The function prototype of **DBMS_LOB.WRITEAPPEND** is:

```
DBMS_LOB.WRITEAPPEND (
lob_loc IN OUT BLOB,
amount IN INTEGER,
buffer IN RAW);

DBMS_LOB.WRITEAPPEND (
lob_loc IN OUT CLOB,
amount IN INTEGER,
buffer IN VARCHAR2);
```

Table 9-8 DBMS_LOB.WRITEAPPEND interface parameters

Parameter	Description
lob_loc	LOB type object to be written
amount	Write data length
	NOTE If the write data is shorter than 1 or longer than the contents to be written, an error is reported.
buffer	Content to be written

DBMS_LOB.COPY

The stored procedure **COPY** copies contents in BLOB to another BLOB according to the specified length and initial position offset.

The function prototype of **DBMS_LOB.COPY** is:

```
DBMS_LOB.COPY (
dest_lob IN OUT BLOB,
src_lob IN BLOB,
amount IN INTEGER,
dest_offset IN INTEGER DEFAULT 1,
src_offset IN INTEGER DEFAULT 1);
```

Table 9-9 DBMS_LOB.COPY interface parameters

Parameter	Description
dest_lob	BLOB type object to be pasted
src_lob	BLOB type object to be copied
amount	Replication length. NOTE If the copied data is shorter than 1 or longer than the maximum length of BLOB type contents, an error is reported.
dest_offset	Indicates where to start pasting the BLOB contents, that is, the offset bytes to initial position of BLOB contents. NOTE If the offset is shorter than 1 or longer than the maximum length of BLOB type contents, an error is reported.

Parameter	Description
src_offset	Indicates where to start copying the BLOB contents, that is, the offset bytes to initial position of BLOB contents.
	NOTE If the offset is shorter than 1 or longer than the length of source BLOB, an error is reported.

DBMS_LOB.ERASE

The stored procedure **ERASE** deletes contents in BLOB according to the specified length and initial position offset.

The function prototype of **DBMS_LOB.ERASE** is:

```
DBMS_LOB.ERASE (
lob_loc IN OUT BLOB,
amount IN OUT INTEGER,
offset IN INTEGER DEFAULT 1);
```

Table 9-10 DBMS_LOB.ERASE interface parameters

Parameter	Description
lob_loc	BLOB type object whose contents are to be deleted
amount	Length of contents to be deleted NOTE If the deleted data is shorter than 1 or longer than the maximum length of BLOB type contents, an error is reported.
offset	Indicates where to start deleting the BLOB contents, that is, the offset bytes to initial position of BLOB contents. NOTE If the offset is shorter than 1 or longer than the maximum length of BLOB type contents, an error is reported.

DBMS_LOB.CLOSE

The procedure **CLOSE** disables the enabled contents of LOB according to the specified length and initial position offset.

The function prototype of **DBMS_LOB.CLOSE** is:

```
DBMS_LOB.CLOSE(
src_lob IN BLOB);

DBMS_LOB.CLOSE (
src_lob IN CLOB);
```

Table 9-11 DBMS_LOB.CLOSE interface parameters

Parameter	Description
src_loc	LOB type object to be disabled

• DBMS_LOB.INSTR

This function returns the Nth occurrence position in LOB. If invalid values are entered, **NULL** is returned. The invalid values include offset < 1 or offset > LOBMAXSIZE, nth < 1, and nth > LOBMAXSIZE.

The function prototype of **DBMS_LOB.INSTR** is:

```
DBMS LOB.INSTR (
lob_loc IN
              BLOB,
pattern IN
               RAW.
offset IN
               INTEGER := 1,
         IN
               INTEGER := 1)
nth
RETURN INTEGER;
DBMS_LOB.INSTR (
lob_loc IN
               CLOB,
pattern IN VARCHAR2, offset IN INTEGER := 1, nth IN INTEGER := 1)
RETURN INTEGER;
```

Table 9-12 DBMS_LOB.INSTR interface parameters

Parameter	Description
lob_loc	LOB descriptor to be searched for
pattern	Matched pattern. It is RAW for BLOB and TEXT for CLOB.
offset	For BLOB, the absolute offset is in the unit of byte. For CLOB, the offset is in the unit of character. The matching start position is 1.
nth	Number of pattern matching times. The minimum value is 1.

DBMS_LOB.COMPARE

This function compares two LOBs or a certain part of two LOBs.

- If the two parts are equal, **0** is returned. Otherwise, a non-zero value is returned.
- If the first CLOB is smaller than the second, -1 is returned. If the first CLOB is larger than the second, 1 is returned.
- If any of the amount, offset_1, and offset_2 parameters is invalid, NULL is returned. The valid offset range is 1 to LOBMAXSIZE.

The function prototype of **DBMS_LOB.READ** is:

```
DBMS_LOB.COMPARE (
lob_1 IN BLOB,
lob_2 IN BLOB,
amount IN INTEGER := DBMS_LOB.LOBMAXSIZE,
offset_1 IN INTEGER := 1,
offset_2 IN INTEGER := 1)
RETURN INTEGER;

DBMS_LOB.COMPARE (
lob_1 IN CLOB,
lob_2 IN CLOB,
amount IN INTEGER := DBMS_LOB.LOBMAXSIZE,
offset_1 IN INTEGER := 1,
offset_1 IN INTEGER := 1,
offset_2 IN INTEGER := 1)
RETURN INTEGER;
```

Table 5 15 DBNIS_EGB.CGNITAILE Interface parameters	
Parameter	Description
lob_1	First LOB descriptor to be compared
lob_2	Second LOB descriptor to be compared
amount	Number of characters or bytes to be compared. The maximum value is DBMS_LOB.LOBMAXSIZE.
offset_1	Offset of the first LOB descriptor. The initial position is 1.
offset_2	Offset of the second LOB descriptor. The initial position is 1.

Table 9-13 DBMS LOB.COMPARE interface parameters

• DBMS LOB.SUBSTR

This function reads the substring of a LOB and returns the number of read bytes or the number of characters. If amount > 1, amount < 32767, offset < 1, or offset > LOBMAXSIZE, **NULL** is returned.

The function prototype of **DBMS_LOB.SUBSTR** is:

```
DBMS_LOB.SUBSTR (
lob_loc IN BLOB,
amount IN INTEGER := 32767,
offset IN INTEGER := 1)
RETURN RAW;

DBMS_LOB.SUBSTR (
lob_loc IN CLOB,
amount IN INTEGER := 32767,
offset IN INTEGER := 1)
RETURN VARCHAR2;
```

Table 9-14 DBMS_LOB.SUBSTR interface parameters

Parameter	Description
lob_loc	LOB descriptor of the substring to be read. For BLOB, the return value is the number of read bytes. For CLOB, the return value is the number of characters.
offset	Number of bytes or characters to be read.
buffer	Number of characters or bytes offset from the start position.

DBMS_LOB.TRIM

This stored procedure truncates the LOB of a specified length. After this stored procedure is executed, the length of the LOB is set to the length specified by the **newlen** parameter. If an empty LOB is truncated, no execution result is displayed. If the specified length is longer than the length of LOB, an exception occurs.

The function prototype of **DBMS_LOB.TRIM** is:

```
DBMS_LOB.TRIM (
lob_loc IN OUT BLOB,
newlen IN INTEGER);
```

```
DBMS_LOB.TRIM (
lob_loc IN OUT CLOB,
newlen IN INTEGER);
```

Table 9-15 DBMS_LOB.TRIM interface parameters

Parame ter	Description
lob_loc	BLOB type object to be read
newlen	After truncation, the new LOB length for BLOB is in the unit of byte and that for CLOB is in the unit of character.

DBMS_LOB.CREATETEMPORARY

This stored procedure creates a temporary BLOB or CLOB and is used only for syntax compatibility.

The function prototype of **DBMS_LOB.CREATETEMPORARY** is:

```
DBMS_LOB.CREATETEMPORARY (
lob_loc IN OUT BLOB,
cache IN BOOLEAN,
dur IN INTEGER);

DBMS_LOB.CREATETEMPORARY (
lob_loc IN OUT CLOB,
cache IN BOOLEAN,
dur IN INTEGER);
```

Table 9-16 DBMS_LOB.CREATETEMPORARY interface parameters

Parameter	Description
lob_loc	LOB descriptor
cache	This parameter is used only for syntax compatibility.
dur	This parameter is used only for syntax compatibility.

DBMS_LOB.APPEND

The stored procedure **READ** loads a part of BLOB contents to BUFFER according to the specified length and initial position offset.

The function prototype of **DBMS_LOB.APPEND** is:

```
DBMS_LOB.APPEND (
dest_lob IN OUT BLOB,
src_lob IN BLOB);

DBMS_LOB.APPEND (
dest_lob IN OUT CLOB,
src_lob IN CLOB);
```

Table 9-17 DBMS_LOB.APPEND interface parameters

Parameter	Description
dest_lob	LOB descriptor to be written
src_lob	LOB descriptor to be read

```
-- Obtain the length of the character string.
SELECT DBMS_LOB.GETLENGTH('12345678');
DECLARE
myraw RAW(100);
amount INTEGER :=2:
buffer INTEGER :=1;
DBMS_LOB.READ('123456789012345',amount,buffer,myraw);
dbms_output.put_line(myraw);
end;
CREATE TABLE blob_Table (t1 blob) DISTRIBUTE BY REPLICATION;
CREATE TABLE blob_Table_bak (t2 blob) DISTRIBUTE BY REPLICATION;
INSERT INTO blob_Table VALUES('abcdef');
INSERT INTO blob_Table_bak VALUES('22222');
DECLARE
str varchar2(100) := 'abcdef';
source raw(100);
dest blob;
copyto blob;
amount int;
PSV_SQL varchar2(100);
PSV_SQL1 varchar2(100);
a int :=1;
len int;
BEGIN
source := utl_raw.cast_to_raw(str);
amount := utl_raw.length(source);
PSV_SQL :='select * from blob_Table for update';
PSV_SQL1 := 'select * from blob_Table_bak for update';
EXECUTE IMMEDIATE PSV_SQL into dest;
EXECUTE IMMEDIATE PSV_SQL1 into copyto;
DBMS_LOB.WRITE(dest, amount, 1, source);
DBMS_LOB.WRITEAPPEND(dest, amount, source);
DBMS_LOB.ERASE(dest, a, 1);
DBMS_OUTPUT.PUT_LINE(a);
DBMS_LOB.COPY(copyto, dest, amount, 10, 1);
DBMS_LOB.CLOSE(dest);
RETURN;
END;
--Delete the table.
DROP TABLE blob_Table;
DROP TABLE blob_Table_bak;
```

9.10.2 DBMS RANDOM

Related Interfaces

Table 9-18 provides all interfaces supported by the DBMS_RANDOM package.

Table 9-18 DBMS_RANDOM interface parameters

АРІ	Description
DBMS_RANDO M.SEED	Sets a seed for a random number.
DBMS_RANDO M.VALUE	Generates a random number between a specified low and a specified high.

DBMS_RANDOM.SEED

The stored procedure SEED is used to set a seed for a random number. The DBMS_RANDOM.SEED function prototype is:

DBMS_RANDOM.SEED (seed IN INTEGER);

Table 9-19 DBMS_RANDOM.SEED interface parameters

Parameter	Description
seed	Generates a seed for a random number.

DBMS_RANDOM.VALUE

The stored procedure VALUE generates a random number between a specified low and a specified high. The DBMS_RANDOM.VALUE function prototype is:

DBMS_RANDOM.VALUE(low IN NUMBER, high IN NUMBER) RETURN NUMBER;

Table 9-20 DBMS_RANDOM.VALUE interface parameters

Paramet er	Description
low	Sets the low bound for a random number. The generated random number is greater than or equal to the low.
high	Sets the high bound for a random number. The generated random number is less than the high.

NOTE

The only requirement is that the parameter type is **NUMERIC** regardless of the right and left bound values.

Example

Generate a random number between 0 and 1:

SELECT DBMS_RANDOM.VALUE(0,1);

Generate a random integer ranging from 0 to 100. The random integer is greater than or equal to the specified value of low and less than the specified value of high.

SELECT TRUNC(DBMS_RANDOM.VALUE(0,100));

9.10.3 DBMS_OUTPUT

Related Interfaces

Table 9-21 provides all interfaces supported by the **DBMS_OUTPUT** package.

Table 9-21 DBMS_OUTPUT

API	Description		
DBMS_OUTP UT.PUT_LINE	Outputs the specified text. The text length cannot exceed 32,767 bytes.		
DBMS_OUTP UT.PUT	Outputs the specified text to the front of the specified text without adding a line break. The text length cannot exceed 32,767 bytes.		
DBMS_OUTP UT.ENABLE	Sets the buffer area size. If this interface is not specified, the maximum buffer size is 20,000 bytes and the minimum buffer size is 2,000 bytes. If the specified buffer size is less than 2,000 bytes, the default minimum buffer size is applied.		

DBMS_OUTPUT.PUT_LINE

The **PUT_LINE** procedure writes a row of text carrying a line end symbol in the buffer. The **DBMS_OUTPUT.PUT_LINE** function prototype is:

DBMS_OUTPUT.PUT_LINE (item IN VARCHAR2);

Table 9-22 DBMS_OUTPUT.PUT_LINE interface parameters

Parameter	Description	
item	Specifies the text that was written to the buffer.	

• DBMS_OUTPUT.PUT

The stored procedure **PUT** outputs the specified text to the front of the specified text without adding a linefeed. The **DBMS_OUTPUT.PUT** function prototype is:

DBMS_OUTPUT.PUT (item IN VARCHAR2);

Table 9-23 DBMS_OUTPUT.PUT interface parameters

Parameter	Description	
item	Specifies the text that was written to the specified text.	

DBMS_OUTPUT.ENABLE

The stored procedure **ENABLE** sets the output buffer size. If the size is not specified, it contains a maximum of 20,000 bytes. The **DBMS_OUTPUT.ENABLE** function prototype is:

DBMS_OUTPUT.ENABLE (buf IN INTEGER);

Table 9-24 DBMS_OUTPUT.ENABLE interface parameters

Parameter	Description	
buf	Sets the buffer area size.	

Examples

```
BEGIN

DBMS_OUTPUT.ENABLE(50);

DBMS_OUTPUT.PUT ('hello, ');

DBMS_OUTPUT.PUT_LINE('database!');-- Displaying "hello, database!"

END;

/
```

9.10.4 UTL_RAW

Related Interfaces

Table 9-25 provides all interfaces supported by the UTL_RAW package.

Table 9-25 UTL_RAW

API	Description	
UTL_RAW.CAST_FROM_BI NARY_INTEGER	Converts an INTEGER type value to a binary representation (RAW type).	
UTL_RAW.CAST_TO_BINA RY_INTEGER	Converts a binary representation (RAW type) to an INTEGER type value.	
UTL_RAW.LENGTH	Obtains the length of the RAW type object.	
UTL_RAW.CAST_TO_RAW	Converts a VARCHAR2 type value to a binary expression (RAW type).	

NOTICE

The external representation of the RAW type data is hexadecimal and its internal storage form is binary. For example, the representation of the **RAW** type data **11001011** is 'CB'. The input of the actual type conversion is 'CB'.

UTL_RAW.CAST_FROM_BINARY_INTEGER

The stored procedure **CAST_FROM_BINARY_INTEGER** converts an **INTEGER** type value to a binary representation (**RAW** type).

The UTL_RAW.CAST_FROM_BINARY_INTEGER function prototype is:

UTL_RAW.CAST_FROM_BINARY_INTEGER (
n IN INTEGER,
endianess IN INTEGER)
RETURN RAW;

Table 9-26 UTL_RAW.CAST_FROM_BINARY_INTEGER interface parameters

Paramete r	Description
n	Specifies the INTEGER type value to be converted to the RAW type.
endianess	Specifies the INTEGER type value 1 or 2 of the byte sequence. (1 indicates BIG_ENDIAN and 2 indicates LITTLE-ENDIAN.)

UTL_RAW.CAST_TO_BINARY_INTEGER

The stored procedure CAST_TO_BINARY_INTEGER converts an INTEGER type value in a binary representation (RAW type) to the INTEGER type.

The UTL_RAW.CAST_TO_BINARY_INTEGER function prototype is:

UTL_RAW.CAST_TO_BINARY_INTEGER (
r IN RAW,
endianess IN INTEGER)
RETURN BINARY_INTEGER;

Table 9-27 UTL_RAW.CAST_TO_BINARY_INTEGER interface parameters

Parameter	Description		
r	Specifies an INTEGER type value in a binary representation (RAW type).		
endianess	Specifies the INTEGER type value 1 or 2 of the byte sequence. (1 indicates BIG_ENDIAN and 2 indicates LITTLE-ENDIAN.)		

UTL RAW.LENGTH

The stored procedure LENGTH returns the length of a RAW type object.

The UTL_RAW.LENGTH function prototype is:

UTL_RAW.LENGTH(r IN RAW) RETURN INTEGER;

Table 9-28 UTL_RAW.LENGTH interface parameters

Parameter	Description	
r	Specifies a RAW type object.	

UTL_RAW.CAST_TO_RAW

The stored procedure CAST_TO_RAW converts a VARCHAR2 type object to the RAW type.

The UTL_RAW.CAST_TO_RAW function prototype is:

```
UTL_RAW.CAST_TO_RAW(
c IN VARCHAR2)
RETURN RAW;
```

Table 9-29 UTL_RAW.CAST_TO_RAW interface parameters

Parameter	Description	
С	Specifies a VARCHAR2 type object to be converted.	

Example

Perform operations on RAW data in a stored procedure:

```
CREATE OR REPLACE PROCEDURE proc_raw
AS
str varchar2(100) := 'abcdef';
source raw(100);
amount integer;
BEGIN
source := utl_raw.cast_to_raw(str);--Convert the type.
amount := utl_raw.length(source);--Obtain the length.
dbms_output.put_line(amount);
END;
/
```

Call the stored procedure:

CALL proc_raw();

9.10.5 **DBMS_JOB**

Related Interfaces

Table 9-30 lists all interfaces supported by the DBMS_JOB package.

Table 9-30 DBMS_JOB

Interface	Description
DBMS_JOB.SUBMIT	Submits a job to the job queue. The job number is automatically generated by the system.

Interface	Description		
DBMS_JOB.SUBMIT _NODE	Submits a job to the job queue. The execution node is specified by the user, and the job number is automatically generated by the system.		
DBMS_JOB.ISUBMI T	Submits a job to the job queue. The job number is specified by the user.		
DBMS_JOB.REMOV E	Removes a job from the job queue by job number.		
DBMS_JOB.BROKE N	Disables or enables job execution.		
DBMS_JOB.CHANG E	Modifies user-definable attributes of a job, including the job description, next execution time, and execution interval.		
DBMS_JOB.WHAT	Modifies the job description of a job.		
DBMS_JOB.NEXT_D ATE	Modifies the next execution time of a job.		
DBMS_JOB.INTERV AL	Modifies the execution interval of a job.		
DBMS_JOB.CHANG E_OWNER	Modifies the owner of a job.		
DBMS_JOB.CHANG Modifies the execution node of the scheduled task E_NODE			

• DBMS_JOB.SUBMIT

The stored procedure **SUBMIT** submits a job provided by the system.

A prototype of the DBMS_JOB.SUBMIT function is as follows:

```
DMBS_JOB.SUBMIT(
what IN TEXT,
next_date IN TIMESTAMP DEFAULT sysdate,
job_interval IN TEXT DEFAULT 'null',
job OUT INTEGER);
```

Ⅲ NOTE

When a job is created (using DBMS_JOB), the system binds the current database and the username to the job by default. This function can be invoked by using **call** or **select**. If you invoke this function by using **select**, there is no need to specify output parameters. To invoke this function within a stored procedure, use **perform**.

Parame ter	Typ e	Input/ Output Parame ter	Can Be Empt y	Description
what	text	IN	No	SQL statement to be executed. One or multiple DMLs, anonymous blocks, and SQL statements that invoke stored procedures, or all three combined are supported.
next_dat e	tim esta mp	IN	No	Specifies the next time the job will be executed. The default value is the current system time (sysdate). If the specified time has past, the job is executed at the time it is submitted.
interval	text	IN	Yes	Calculates the next time to execute the job. It can be an interval expression, or sysdate followed by a numeric value, for example, sysdate+1.0/24 . If this parameter is left blank or set to null , the job will be executed only once, and the job status will change to 'd' afterward.
job	inte ger	OUT	No	Specifies the job number. The value ranges from 1 to 32767. When dbms.submit is invoked using select , this parameter can be skipped.

Table 9-31 DBMS_JOB.SUBMIT interface parameters

For example:

select DBMS_JOB.SUBMIT('call pro_xxx();', to_date('20180101','yyyymmdd'),'sysdate+1');

select DBMS_JOB.SUBMIT('call pro_xxx();', to_date('20180101','yyyymmdd'),'sysdate+1.0/24');

CALL DBMS_JOB.SUBMIT('INSERT INTO T_JOB VALUES(1); call pro_1(); call pro_2();', add_months(to_date('201701','yyyymm'),1), 'date_trunc("day",SYSDATE) + 1 +(8*60+30.0)/(24*60)',:jobid);

• DBMS_JOB.SUBMIT_NODE

The stored procedure **SUBMIT** submits a job provided by the system. The execution node is specified by the user. This interface is supported only by clusters of version 8.3.0 or later.

The prototype of the DBMS_JOB.SUBMIT_NODE function is:

```
DMBS_IOB.SUBMIT_NODE(
what IN TEXT,
next_date IN TIMESTAMP DEFAULT sysdate,
job_interval IN TEXT DEFAULT 'null',
job_node IN TEXT DEFAULT NULL,
job OUT INTEGER);
```

Parame ter	Typ e	Input/ Output Parame ter	Can Be Empt y	Description
what	text	IN	No	Specifies the SQL statement to be executed. One or multiple DMLs, anonymous blocks, and SQL statements that invoke stored procedures, or all three combined are supported.
next_dat e	tim esta mp	IN	No	Specifies the next time the job will be executed. The default value is the current system time (sysdate). If the specified time has past, the job is executed at the time it is submitted.
interval	text	IN	Yes	Calculates the next time to execute the job. It can be an interval expression, or sysdate followed by a numeric value, for example, sysdate+1.0/24 . If this parameter is left blank or set to null , the job will be executed only once, and the job status will change to 'd' afterward.
node	text	IN	Yes	Specifies the name of the job execution node.
job	inte ger	OUT	No	Specifies the job number. The value ranges from 1 to 32767. When dbms.submit is invoked using select , this parameter can be skipped.

Table 9-32 DBMS_JOB.SUBMIT_NODE interface parameters

For example:

select DBMS_JOB.SUBMIT_NODE('call pro_xxx();', to_date('20180101','yyyymmdd'),'sysdate +1','coordinator1');

select DBMS_JOB.SUBMIT_NODE('call pro_xxx();', to_date('20180101','yyyymmdd'),'sysdate+1.0/24');

CALL DBMS_JOB.SUBMIT('INSERT INTO T_JOB VALUES(1); call pro_1(); call pro_2();', add_months(to_date('201701','yyyymm'),1), 'date_trunc("day",SYSDATE) + 1 +(8*60+30.0)/(24*60)', 'coordinator1', :jobid);

DBMS_JOB.ISUBMIT

ISUBMIT has the same syntax function as **SUBMIT**, but the first parameter of **ISUBMIT** is an input parameter, that is, a specified job number. In contrast, that last parameter of **SUBMIT** is an output parameter, indicating the job number automatically generated by the system.

For example:

CALL dbms_job.isubmit(101, 'insert_msg_statistic1;', sysdate, 'sysdate+3.0/24');

NOTICE

The pgstats persistence function of GaussDB(DWS) writes the statistics in the memory to the **pg_stat_object** system catalog. If the cluster version is 9.1.0.100 or later, **1** is used as **job_id**. If an earlier cluster version is upgraded to 9.1.0.100 or later and **pg_job** contains tasks, an unoccupied **job_id** is used as the ID of the persistence task. Therefore, when using the **dbms_job.isubmit** interface, ensure that the ID is different from the ID of an existing pgstats persistence task. Otherwise, the task registration fails.

DBMS JOB.REMOVE

The stored procedure **REMOVE** deletes a specified job.

A prototype of the DBMS_JOB.REMOVE function is as follows:

REMOVE(job IN INTEGER);

Table 9-33 DBMS_JOB.REMOVE interface parameters

Para mete r	Туре	Input/ Output Paramet er	Can Be Empty	Description
job	integ er	IN	No	Specifies the job number.

For example:

CALL dbms_job.remove(101);

DBMS_JOB.BROKEN

The stored procedure **BROKEN** sets the broken flag of a job.

A prototype of the DBMS_JOB.BROKEN function is as follows:

DMBS_JOB.BROKEN(
job IN INTEGER,
broken IN BOOLEAN,
next_date IN TIMESTAMP DEFAULT sysdate);

Table 9-34 DBMS JOB.BROKEN interface parameters

Param eter	Туре	Input/ Outpu t Param eter	Ca n Be Em pty	Description
job	integer	IN	No	Specifies the job number.

Param eter	Туре	Input/ Outpu t Param eter	Ca n Be Em pty	Description
broken	boolean	IN	No	Specifies the status flag, true for broken and false for not broken. Setting this parameter to true or false updates the current job. If the parameter is left blank, the job status remains unchanged.
next_da te	timesta mp	IN	Yes	Specifies the next execution time. The default is the current system time. If broken is set to true, next_date is updated to '4000-1-1'. If broken is false and next_date is not empty, next_date is updated for the job. If next_date is empty, it will not be updated. This parameter can be omitted, and its default value will be used in this case.

For example:

CALL dbms_job.broken(101,true);
CALL dbms_job.broken(101,false,sysdate);

• DBMS_JOB.CHANGE

The stored procedure **CHANGE** modifies user-definable attributes of a job, including the job content, next-execution time, and execution interval.

A prototype of the DBMS_JOB.CHANGE function is as follows:

DMBS_JOB.CHANGE(
job IN INTEGER,
what IN TEXT,
next_date IN TIMESTAMP,
interval IN TEXT);

Table 9-35 DBMS_JOB.CHANGE interface parameters

Para met er	Туре	Input/ Output Paramet er	Can Be Empty	Description
job	integ er	IN	No	Specifies the job number.

Para met er	Туре	Input/ Output Paramet er	Can Be Empty	Description
wha t	text	IN	Yes	Specifies the name of the stored procedure or SQL statement block that is executed. If this parameter is left blank, the system does not update the what parameter for the specified job. Otherwise, the system updates the what parameter for the specified job.
next _dat e	time stam p	IN	Yes	Specifies the next execution time. If this parameter is left blank, the system does not update the next_date parameter for the specified job. Otherwise, the system updates the next_date parameter for the specified job.
inter val	text	IN	Yes	Specifies the time expression for calculating the next time the job will be executed. If this parameter is left blank, the system does not update the interval parameter for the specified job. Otherwise, the system updates the interval parameter for the specified job after necessary validity check. If this parameter is set to null , the job will be executed only once, and the job status will change to 'd' afterward.

For example:

CALL dbms_job.change(101, 'call userproc();', sysdate, 'sysdate + 1.0/1440');
CALL dbms_job.change(101, 'insert into tbl_a values(sysdate);', sysdate, 'sysdate + 1.0/1440');

DBMS_JOB.WHAT

The stored procedure **WHAT** modifies the procedures to be executed by a specified job.

A prototype of the DBMS_JOB.WHAT function is as follows:

DMBS_JOB.WHAT(job IN INTEGER, what IN TEXT);

Par am ete r	Туре	Input/ Output Paramet er	Can Be Empty	Description
job	intege r	IN	No	Specifies the job number.
wh at	text	IN	No	Specifies the name of the stored procedure or SQL statement block that is executed.

Table 9-36 DBMS_JOB.WHAT interface parameters

Ⅲ NOTE

- If the value specified by the **what** parameter is one or multiple executable SQL statements, program blocks, or stored procedures, this procedure can be executed successfully; otherwise, it will fail to be executed.
- If the **what** parameter is a simple statement such as insert and update, a schema name must be added in front of the table name.

For example:

CALL dbms_job.what(101, 'call userproc();');
CALL dbms_job.what(101, 'insert into tbl_a values(sysdate);');

DBMS JOB.NEXT DATE

The stored procedure **NEXT_DATE** modifies the next-execution time attribute of a job.

A prototype of the DBMS JOB.NEXT DATE function is as follows:

DMBS_JOB.NEXT_DATE(job IN INTEGER, next_date IN TIMESTAMP);

Table 9-37 DBMS_JOB.NEXT_DATE interface parameters

Parame ter	Туре	Input/ Output Param eter	Can Be Empty	Description
job	integer	IN	No	Specifies the job number.
next_da te	timesta mp	IN	No	Specifies the next execution time.

□ NOTE

If the specified **next_date** value is earlier than the current date, the job is executed once immediately.

For example:

CALL dbms_job.next_date(101,sysdate);

DBMS_JOB.INTERVAL

The stored procedure **INTERVAL** modifies the execution interval attribute of a job.

A prototype of the DBMS_JOB.INTERVAL function is as follows:

DMBS_JOB.INTERVAL(job IN INTEGER, interval IN TEXT);

Table 9-38 DBMS_JOB.INTERVAL interface parameters

Parame ter	Туре	Input / Outp ut Para meter	Can Be Empty	Description
job	intege r	IN	No	Specifies the job number.
interval	text	IN	Yes	Specifies the time expression for calculating the next time the job will be executed. If this parameter is left blank or set to null , the job will be executed only once, and the job status will change to ' d ' afterward. interval must be a valid time or interval type.

For example:

CALL dbms_job.interval(101, 'sysdate + 1.0/1440');

□ NOTE

For a job that is currently running (that is, **job_status** is **'r'**), it is not allowed to use **remove**, **change**, **next_date**, **what**, or **interval** to delete or modify job parameters.

DBMS_JOB.CHANGE_OWNER

The stored procedure ${\it CHANGE_OWNER}$ modifies the owner of a job.

A prototype of the DBMS_JOB.CHANGE_OWNER function is as follows:

DMBS_JOB.CHANGE_OWNER(job IN INTEGER, new_owner IN NAME);

Table 9-39 DBMS_JOB.CHANGE_OWNER interface parameters

Paramet er	Туре	Input/ Output Paramet er	Can Be Empty	Description
job	integer	IN	No	Specifies the job number.

Paramet er	Туре	Input/ Output Paramet er	Can Be Empty	Description
new_own er	name	IN	No	Specifies the new username.

For example:

CALL dbms_job.change_owner(101, 'alice');

DBMS JOB.CHANGE NODE

The stored procedure **CHANGE_NODE** modifies the execution node of the scheduled task. This interface is supported only by clusters of version 8.3.0 or later.

A prototype of the DBMS_JOB.CHANGE_NODE function is:

```
DMBS_JOB.CHANGE_NODE(
job IN INTEGER,
new_node IN text);
```

Table 9-40 DBMS JOB.CHANGE OWNER interface parameters

Paramet er	Туре	Input/ Output Paramet er	Can Be Empty	Description
job	integer	IN	No	Specifies the job number.
new_nod e	text	IN	No	Specifies the new execution node.

For example:

CALL dbms_job.change_node(101, 'coordinator2');

Constraints

- After a new job is created, this job belongs to the current coordinator only, that is, this job can be scheduled and executed only on the current coordinator. Other coordinators will not schedule or execute this job. All coordinators can query, modify, and delete jobs created on other CNs.
- 2. Create, update, and delete jobs only using the procedures provided by the DBMS_JOB package. These procedures synchronize job information between different CNs and associate primary keys between the pg_jobs tables. If you use DML statements to add, delete, or modify records in the pg_jobs table, job information will become inconsistent between CNs and system tables may fail to be associated, compromising internal job management.
- 3. Each user-created task is bound to a CN. If the automatic migration function is not enabled, task statuses cannot be updated in real time when the CN is

- faulty during task execution. When a CN fails, all jobs on this CN cannot be scheduled or executed until the CN is restored manually. Enable the automatic migration function on CNs, so that jobs on the faulty CN will be migrated to other CNs for scheduling.
- 4. For each job, the hosting CN updates the real-time job information (including the job status, last execution start time, last execution end time, next execution start time, the number of execution failures if any) to the **pg_jobs** table, and synchronizes the information to other CNs, ensuring consistent job information between different CNs. In the case of CN failures, job information synchronization is reattempted by the hosting CNs, which increases job execution time. Although job information fails to be synchronized between CNs, job information can still be properly updated in the **pg_jobs** table on the hosting CNs, and jobs can be executed successfully. After a CN recovers, job information such as job execution time and status in its **pg_jobs** table may be incorrect and will be updated only after the jobs are executed again on related CNs.
- 5. For each job, a thread is established to execute it. If multiple jobs are triggered concurrently as scheduled, the system will need some time to start the required threads, resulting in a latency of 0.1 ms in job execution.
- 6. The length of the SQL statement to be executed in a job is limited. The maximum length is 8 KB.

9.10.6 DBMS SQL

Related Interfaces

Table 9-41 lists interfaces supported by the **DBMS_SQL** package.

Table 9-41 DBMS SQL

API	Description
DBMS_SQL.OPEN_CURSOR	Opens a cursor.
DBMS_SQL.CLOSE_CURSOR	Closes an open cursor.
DBMS_SQL.PARSE	Transmits a group of SQL statements to a cursor. Currently, only the SELECT statement is supported.
DBMS_SQL.EXECUTE	Performs a set of dynamically defined operations on the cursor.
DBMS_SQL.FETCHE_ROWS	Reads a row of cursor data.
DBMS_SQL.DEFINE_COLUMN	Dynamically defines a column.
DBMS_SQL.DEFINE_COLUMN_CHAR	Dynamically defines a column of the CHAR type.
DBMS_SQL.DEFINE_COLUMN_INT	Dynamically defines a column of the INT type.

API	Description
DBMS_SQL.DEFINE_COLUMN_LONG	Dynamically defines a column of the LONG type.
DBMS_SQL.DEFINE_COLUMN_RAW	Dynamically defines a column of the RAW type.
DBMS_SQL.DEFINE_COLUMN_TEXT	Dynamically defines a column of the TEXT type.
DBMS_SQL.DEFINE_COLUMN_UNKNOW N	Dynamically defines a column of an unknown type.
DBMS_SQL.COLUMN_VALUE	Reads a dynamically defined column value.
DBMS_SQL.COLUMN_VALUE_CHAR	Reads a dynamically defined column value of the CHAR type.
DBMS_SQL.COLUMN_VALUE_INT	Reads a dynamically defined column value of the INT type.
DBMS_SQL.COLUMN_VALUE_LONG	Reads a dynamically defined column value of the LONG type.
DBMS_SQL.COLUMN_VALUE_RAW	Reads a dynamically defined column value of the RAW type.
DBMS_SQL.COLUMN_VALUE_TEXT	Reads a dynamically defined column value of the TEXT type.
DBMS_SQL.COLUMN_VALUE_UNKNOWN	Reads a dynamically defined column value of an unknown type.
DBMS_SQL.IS_OPEN	Checks whether a cursor is opened.

MOTE

- You are advised to use dbms_sql.define_column and dbms_sql.column_value to define columns.
- If the size of the result set is greater than the value of **work_mem**, the result set will be flushed to disk. The value of **work_mem** must be no greater than 512 MB.
- DBMS_SQL.OPEN_CURSOR

This function opens a cursor and is the prerequisite for the subsequent dbms_sql operations. This function does not transfer any parameter. It automatically generates cursor IDs in an ascending order and returns values to integer variables.

The function prototype of **DBMS_SQL.OPEN_CURSOR** is:

DBMS_SQL.OPEN_CURSOR (
)
RETURN INTEGER;

DBMS_SQL.CLOSE_CURSOR

This function closes a cursor. It is the end of each dbms_sql operation. If this function is not invoked when the stored procedure ends, the memory is still occupied by the cursor. Therefore, remember to close a cursor when you do not need to use it. If an exception occurs, the stored procedure exits but the cursor is not closed. Therefore, you are advised to include this interface in the exception handling of the stored procedure.

The function prototype of DBMS_SQL.CLOSE_CURSOR is:

```
DBMS_SQL.CLOSE_CURSOR (
cursorid IN INTEGER
)
RETURN INTEGER;
```

Table 9-42 DBMS_SQL.CLOSE_CURSOR interface parameters

Parameter Name	Description
cursorid	ID of the cursor to be closed

DBMS_SQL.PARSE

RETURN BOOLEAN;

This function parses the query statement of a given cursor. The input query statement is executed immediately. Currently, only the **SELECT** query statement can be parsed. The statement parameters can be transferred only through the TEXT type. The length cannot exceed 1 GB.

```
The function prototype of DBMS_SQL.PARSE is:

DBMS_SQL.PARSE (
cursorid IN INTEGER,
query_string IN TEXT,
label IN INTEGER
```

Table 9-43 DBMS_SQL.PARSE interface parameters

Parameter Name	Description
cursorid	ID of the cursor whose query statement is parsed
query_string	Query statements to be parsed
language_flag	Version language number. Currently, only 1 is supported.

DBMS_SQL.EXECUTE

This function executes a given cursor. This function receives a cursor ID. The obtained data after is used for subsequent operations. Currently, only the **SELECT** query statement can be executed.

```
The function prototype of DBMS_SQL.EXECUTE is:

DBMS_SQL.EXECUTE(
cursorid IN INTEGER,
)
RETURN INTEGER;
```

Table 9-44 DBMS_SQL.EXECUTE interface parameters

Parameter Name	Description
cursorid	ID of the cursor whose query statement is parsed

DBMS_SQL.FETCHE_ROWS

This function returns the number of data rows that meet query conditions. Each time the interface is executed, the system obtains a set of new rows until all data is read.

```
The function prototype of DBMS_SQL.FETCHE_ROWS is:

DBMS_SQL.FETCHE_ROWS(
cursorid IN INTEGER,
)
RETURN INTEGER;
```

Table 9-45 DBMS_SQL.FETCH_ROWS interface parameters

Parameter Name	Description
curosorid	ID of the cursor to be executed

DBMS_SQL.DEFINE_COLUMN

This function defines columns returned from a given cursor and can be used only for the cursors defined by **SELECT**. The defined columns are identified by the relative positions in the query list. The data type of the input variable determines the column type.

```
The function prototype of DBMS_SQL.DEFINE_COLUMN is:

DBMS_SQL.DEFINE_COLUMN(
cursorid IN INTEGER,
position IN INTEGER,
column_ref IN ANYELEMENT,
column_size IN INTEGER default 1024
)
RETURN INTEGER;
```

Table 9-46 DBMS_SQL.DEFINE_COLUMN interface parameters

Parameter Name	Description
cursorid	ID of the cursor to be executed
position	Position of a dynamically defined column in the query
column_ref	Variable of any type. You can select an appropriate interface to dynamically define columns based on variable types.
column_size	Length of a defined column

DBMS_SQL.DEFINE_COLUMN_CHAR

This function defines columns of the CHAR type returned from a given cursor and can be used only for the cursors defined by **SELECT**. The defined columns are identified by the relative positions in the query list. The data type of the input variable determines the column type.

```
The function prototype of DBMS_SQL.DEFINE_COLUMN_CHAR is:

DBMS_SQL.DEFINE_COLUMN_CHAR(
cursorid IN INTEGER,
position IN INTEGER,
column IN TEXT,
column_size IN INTEGER
)

RETURN INTEGER;
```

Table 9-47 DBMS SQL.DEFINE COLUMN CHAR interface parameters

Parameter Name	Description
cursorid	ID of the cursor to be executed
position	Position of a dynamically defined column in the query
column	Parameter to be defined
column_size	Length of a dynamically defined column

DBMS_SQL.DEFINE_COLUMN_INT

This function defines columns of the INT type returned from a given cursor and can be used only for the cursors defined by **SELECT**. The defined columns are identified by the relative positions in the query list. The data type of the input variable determines the column type.

```
The function prototype of DBMS_SQL.DEFINE_COLUMN_INT is:

DBMS_SQL.DEFINE_COLUMN_INT(
cursorid IN INTEGER,
position IN INTEGER
)

RETURN INTEGER;
```

Table 9-48 DBMS_SQL.DEFINE_COLUMN_INT interface parameters

Parameter Name	Description
cursorid	ID of the cursor to be executed
position	Position of a dynamically defined column in the query

DBMS_SQL.DEFINE_COLUMN_LONG

This function defines columns of a long type (not LONG) returned from a given cursor and can be used only for the cursors defined by **SELECT**. The defined columns are identified by the relative positions in the query list. The data type of the input variable determines the column type. The maximum size of a long column is 1 GB.

The function prototype of **DBMS_SQL.DEFINE_COLUMN_LONG** is:

```
DBMS_SQL.DEFINE_COLUMN_LONG(
cursorid IN INTEGER,
position IN INTEGER
)
RETURN INTEGER;
```

Table 9-49 DBMS_SQL.DEFINE_COLUMN_LONG interface parameters

Parameter Name	Description
cursorid	ID of the cursor to be executed
position	Position of a dynamically defined column in the query

DBMS_SQL.DEFINE_COLUMN_RAW

This function defines columns of the RAW type returned from a given cursor and can be used only for the cursors defined by **SELECT**. The defined columns are identified by the relative positions in the query list. The data type of the input variable determines the column type.

```
The function prototype of DBMS_SQL.DEFINE_COLUMN_RAW is:
```

```
DBMS_SQL.DEFINE_COLUMN_RAW(
cursorid IN INTEGER,
position IN INTEGER,
column IN BYTEA,
column_size IN INTEGER
)
RETURN INTEGER;
```

Table 9-50 DBMS_SQL.DEFINE_COLUMN_RAW interface parameters

Parameter Name	Description
cursorid	ID of the cursor to be executed
position	Position of a dynamically defined column in the query
column	Parameter of the RAW type
column_size	Column length

DBMS_SQL.DEFINE_COLUMN_TEXT

This function defines columns of the TEXT type returned from a given cursor and can be used only for the cursors defined by **SELECT**. The defined columns are identified by the relative positions in the query list. The data type of the input variable determines the column type.

```
The function prototype of DBMS_SQL.DEFINE_COLUMN_TEXT is:
```

```
DBMS_SQL.DEFINE_COLUMN_CHAR(
cursorid IN INTEGER,
position IN INTEGER,
max_size IN INTEGER
)
RETURN INTEGER;
```

Table 9-51 DBMS_SQL.DEFINE_COLUMN_TEXT interface parameters

Parameter Name	Description
cursorid	ID of the cursor to be executed
position	Position of a dynamically defined column in the query
max_size	Maximum length of the defined TEXT type

DBMS_SQL.DEFINE_COLUMN_UNKNOWN

This function processes columns of unknown data types returned from a given cursor and is used only for the system to report an error and exist when the type cannot be identified.

```
The function prototype of DBMS_SQL.DEFINE_COLUMN_UNKNOWN is: DBMS_SQL.DEFINE_COLUMN_CHAR(
```

```
cursorid IN INTEGER, position IN INTEGER, column IN TEXT
)
RETURN INTEGER;
```

Table 9-52 DBMS SQL.DEFINE COLUMN UNKNOWN interface parameters

Parameter Name	Description
cursorid	ID of the cursor to be executed
position	Position of a dynamically defined column in the query
column	Dynamically defined parameter

DBMS_SQL.COLUMN_VALUE

This function returns the cursor element value specified by a cursor and accesses the data obtained by DBMS_SQL.FETCH_ROWS.

```
The function prototype of DBMS_SQL.COLUMN_VALUE is:
```

```
DBMS_SQL.COLUMN_VALUE(
cursorid IN INTEGER,
position IN INTEGER,
column_value INOUT ANYELEMENT
)
RETURN ANYELEMENT;
```

Table 9-53 DBMS_SQL.COLUMN_VALUE interface parameters

Parameter Name	Description
cursorid	ID of the cursor to be executed
position	Position of a dynamically defined column in the query

Parameter Name	Description
column_value	Return value of a defined column

DBMS_SQL.COLUMN_VALUE_CHAR

This function returns the value of the CHAR type in a specified position of a cursor and accesses the data obtained by DBMS_SQL.FETCH_ROWS.

The function prototype of **DBMS_SQL.COLUMN_VALUE_CHAR** is:

DBMS_SQL.COLUMN_VALUE_CHAR(
cursorid IN INTEGER,
position IN INTEGER,
column_value INOUT CHARACTER,
err_num INOUT NUMERIC default 0,
actual_length INOUT INTEGER default 1024
)
RETURN RECORD;

Table 9-54 DBMS_SQL.COLUMN_VALUE_CHAR interface parameters

Parameter Name	Description
cursorid	ID of the cursor to be executed
position	Position of a dynamically defined column in the query
column_value	Return value
err_num	Error No. It is an output parameter and the argument must be a variable. Currently, the output value is -1 regardless of the argument.
actual_length	Length of a return value

DBMS_SQL.COLUMN_VALUE_INT

This function returns the value of the INT type in a specified position of a cursor and accesses the data obtained by DBMS_SQL.FETCH_ROWS. The function prototype of **DBMS_SQL.COLUMN_VALUE_INT** is:

DBMS_SQL.COLUMN_VALUE_INT(
cursorid IN INTEGER,
position IN INTEGER
)
RETURN INTEGER;

Table 9-55 DBMS_SQL.COLUMN_VALUE_INT interface parameters

Parameter Name	Description
cursorid	ID of the cursor to be executed
position	Position of a dynamically defined column in the query

DBMS_SQL.COLUMN_VALUE_LONG

This function returns the value of a long type (not LONG or BIGINT) in a specified position of a cursor and accesses the data obtained by DBMS_SQL.FETCH_ROWS.

The function prototype of **DBMS_SQL.COLUMN_VALUE_LONG** is:

```
DBMS_SQL.COLUMN_VALUE_LONG(
cursorid IN INTEGER,
position IN INTEGER,
length IN INTEGER,
off_set IN INTEGER,
column_value INOUT TEXT,
actual_length INOUT INTEGER default 1024
)
RETURN RECORD;
```

Table 9-56 DBMS_SQL.COLUMN_VALUE_LONG interface parameters

Parameter Name	Description
cursorid	ID of the cursor to be executed
position	Position of a dynamically defined column in the query
length	Length of a return value
off_set	Start position of a return value
column_value	Return value
actual_length	Length of a return value

DBMS_SQL.COLUMN_VALUE_RAW

This function returns the value of the RAW type in a specified position of a cursor and accesses the data obtained by DBMS_SQL.FETCH_ROWS.

```
The function prototype of DBMS SQL.COLUMN VALUE RAW is:
```

```
DBMS_SQL.COLUMN_VALUE_RAW(
cursorid IN INTEGER,
position IN INTEGER,
column_value INOUT BYTEA,
err_num INOUT NUMERIC default 0,
actual_length INOUT INTEGER default 1024
)
RETURN RECORD;
```

Table 9-57 DBMS_SQL.COLUMN_VALUE_RAW interface parameters

Parameter Name	Description
cursorid	ID of the cursor to be executed
position	Position of a dynamically defined column in the query
column_value	Returned column value

Parameter Name	Description
err_num	Error No. It is an output parameter and the argument must be a variable. Currently, the output value is -1 regardless of the argument.
actual_length	Length of a return value. The value longer than this length will be truncated.

DBMS_SQL.COLUMN_VALUE_TEXT

This function returns the value of the TEXT type in a specified position of a cursor and accesses the data obtained by DBMS_SQL.FETCH_ROWS.

The function prototype of **DBMS_SQL.COLUMN_VALUE_TEXT** is:

DBMS_SQL.COLUMN_VALUE_TEXT(
cursorid IN INTEGER,
position IN INTEGER
)

RETURN TEXT;

Table 9-58 DBMS_SQL.COLUMN_VALUE_TEXT interface parameters

Parameter Name	Description
cursorid	ID of the cursor to be executed
position	Position of a dynamically defined column in the query

DBMS_SQL.COLUMN_VALUE_UNKNOWN

This function returns the value of an unknown type in a specified position of a cursor. This is an error handling interface when the type is not unknown.

The function prototype of **DBMS_SQL.COLUMN_VALUE_UNKNOWN** is:

DBMS_SQL.COLUMN_VALUE_UNKNOWN(
cursorid IN INTEGER,
position IN INTEGER,
COLUMN_TYPE IN TEXT

RETURN TEXT;

Table 9-59 DBMS_SQL.COLUMN_VALUE_UNKNOWN interface parameters

Parameter Name	Description
cursorid	ID of the cursor to be executed
position	Position of a dynamically defined column in the query
column_type	Returned parameter type

DBMS_SQL.IS_OPEN

This function returns the status of a cursor: **open**, **parse**, **execute**, or **define**. The value is **TRUE**. If the status is unknown, an error is reported. In other cases, the value is **FALSE**.

```
The function prototype of DBMS_SQL.IS_OPEN is:

DBMS_SQL.IS_OPEN(
cursorid IN INTEGER
)

RETURN BOOLEAN;
```

Table 9-60 DBMS_SQL.IS_OPEN interface parameters

Parameter Name	Description
cursorid	ID of the cursor to be queried

Examples

```
--Perform operations on raw data in a stored procedure.
create or replace procedure pro_dbms_sql_all_02(in_raw raw,v_in int,v_offset int)
cursorid int;
v_id int;
v_info bytea :=1;
query varchar(2000);
execute_ret int;
define_column_ret_raw bytea :='1';
define_column_ret int;
begin
drop table if exists pro dbms sql all tb1 02;
create table pro_dbms_sql_all_tb1_02(a int ,b blob);
insert into pro_dbms_sql_all_tb1_02 values(1,HEXTORAW('DEADBEEE'));
insert into pro_dbms_sql_all_tb1_02 values(2,in_raw);
query := 'select * from pro_dbms_sql_all_tb1_02 order by 1';
--Open a cursor.
cursorid := dbms_sql.open_cursor();
--Compile the cursor.
dbms_sql.parse(cursorid, query, 1);
--Define a column.
define_column_ret:= dbms_sql.define_column(cursorid,1,v_id);
define_column_ret_raw:= dbms_sql.define_column_raw(cursorid,2,v_info,10);
-- Execute the cursor.
execute_ret := dbms_sql.execute(cursorid);
exit when (dbms_sql.fetch_rows(cursorid) <= 0);
--Obtain values.
dbms_sql.column_value(cursorid,1,v_id);
dbms_sql.column_value_raw(cursorid,2,v_info,v_in,v_offset);
--Output the result.
dbms_output.put_line('id:'|| v_id || ' info:' || v_info);
end loop;
--Close the cursor.
dbms_sql.close_cursor(cursorid);
end;
--Invoke the stored procedure.
call pro_dbms_sql_all_02(HEXTORAW('DEADBEEF'),0,1);
--Delete the stored procedure.
DROP PROCEDURE pro_dbms_sql_all_02;
```

9.11 GaussDB(DWS) Stored Procedure Debugging

Syntax

RAISE has the following five syntax formats:

Figure 9-34 raise_format::=

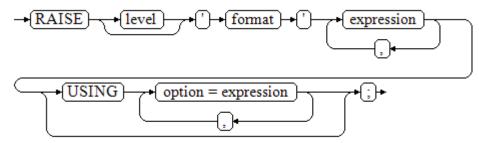


Figure 9-35 raise_condition::=

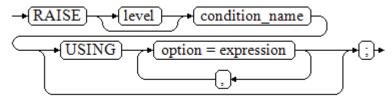


Figure 9-36 raise_sqlstate::=

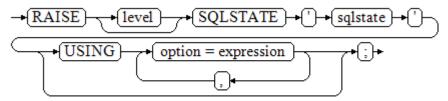


Figure 9-37 raise_option::=

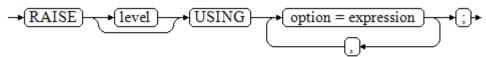


Figure 9-38 raise::=

Parameter description:

- The level option is used to specify the error level, that is, DEBUG, LOG, INFO, NOTICE, WARNING, or EXCEPTION (default). EXCEPTION reports an error that normally terminates the current transaction and the others only generate information at their levels. The log_min_messages and client_min_messages parameters control whether the error messages of specific levels are reported to the client and are written to the server log.
- **format**: specifies the error message text to be reported, a format character string. The format character string can be appended with an expression for insertion to the message text. In a format character string, **%** is replaced by the parameter value attached to format and **%%** is used to print **%**. For example:

--v_job_id replaces % in the character string.
RAISE NOTICE 'Calling cs_create_job(%)',v_job_id;

- option = expression: inserts additional information to an error report. The keyword option can be MESSAGE, DETAIL, HINT, or ERRCODE, and each expression can be any character string.
 - MESSAGE: specifies the error message text. This option cannot be used in a RAISE statement that contains a format character string in front of USING.
 - DETAIL: specifies detailed information of an error.
 - HINT: prints hint information.
 - ERRCODE: designates an error code (SQLSTATE) to a report. A condition name or a five-character SQLSTATE error code can be used.
- condition_name: specifies the condition name corresponding to the error code
- sqlstate: specifies the error code.

If neither a condition name nor an **SQLSTATE** is designated in a **RAISE EXCEPTION** command, the **RAISE EXCEPTION** (**P0001**) is used by default. If no message text is designated, the condition name or SQLSTATE is used as the message text by default.

NOTICE

If the **SQLSTATE** designates an error code, the error code is not limited to a defined error code. It can be any error code containing five digits or ASCII uppercase rather than **00000**. Avoid using error codes that end with three zeros as they are type codes and can be captured by the entire category.

■ NOTE

The syntax described in **Figure 9-38** does not append any parameter. This form is used only for the **EXCEPTION** statement in a **BEGIN** block so that the error can be re-processed.

Examples

Display error and hint information when a transaction terminates:

CREATE OR REPLACE PROCEDURE proc_raise1(user_id in integer) AS REGIN

RAISE EXCEPTION 'Noexistence ID --> %',user_id USING HINT = 'Please check your user ID';

```
END;
/
call proc_raise1(300011);
-- Execution result:
ERROR: Noexistence ID --> 300011
HINT: Please check your user ID

Two methods are available for setting SQLSTATE:
CREATE OR REPLACE PROCEDURE proc_raise2(user_id in integer)
AS
BEGIN
RAISE 'Duplicate user ID: %',user_id USING ERRCODE = 'unique_violation';
END;
/
\set VERBOSITY verbose
call proc_raise2(300011);
```

If the main parameter is a condition name or **SQLSTATE**, the following applies:

RAISE division_by_zero;

RAISE SQLSTATE '22012';

ERROR: Duplicate user ID: 300011

LOCATION: exec_stmt_raise, pl_exec.cpp:3482

For example:

-- Execution result:

SQLSTATE: 23505

```
CREATE OR REPLACE PROCEDURE division(div in integer, dividend in integer)

AS

DECLARE
res int;

BEGIN

IF dividend=0 THEN

RAISE division_by_zero;

RETURN;

ELSE

res := div/dividend;

RAISE INFO 'division result: %', res;

RETURN;

END IF;

END;

/

call division(3,0);

-- Execution result:

ERROR: division_by_zero
```

Alternatively:

RAISE unique_violation USING MESSAGE = 'Duplicate user ID: ' || user_id;

10 Using PostGIS Extension

10.1 PostGIS

□ NOTE

- The third-party software that the PostGIS Extension depends on needs to be installed separately. If you need to use PostGIS, submit a service ticket or contact technical support to submit an application.
- If the error message "ERROR: EXTENSION is not yet supported." is displayed, the PostGIS software package is not installed. Contact technical support.

GaussDB(DWS) provides PostGIS Extension (PostGIS-2.4.2 and PostGIS-3.2.2). PostGIS Extension is a spatial database extender for PostgreSQL. It provides the following spatial information services: spatial objects, spatial indexes, spatial functions, and spatial operators. PostGIS Extension complies with the OpenGIS specifications.

In GaussDB(DWS), PostGIS Extension depends on the listed third-party open-source software.

- PostGIS 2.4.2 depends on the following third-party open-source software:
 - Geos 3.6.2
 - Proj 4.9.2
 - Json 0.12.1
 - Libxml2 2.7.1
 - Gdal 1.11.0
- PostGIS 3.2.2 depends on the following third-party open-source software:
 - Geos-3.11.0
 - Proj-6.0.0
 - Json 0.12.1
 - Libxml2 2.7.1
 - Sqlite3
 - protobuf-c 1.4.1

protobuf 3.6.1

10.2 Using PostGIS

◯ NOTE

- The third-party software that the PostGIS Extension depends on needs to be installed separately. If you need to use PostGIS, submit a service ticket or contact technical support to submit an application.
- If the error message "ERROR: EXTENSION is not yet supported." is displayed, the PostGIS software package is not installed. Contact technical support.

Creating PostGIS Extension

Run the **CREATE EXTENSION** command to create PostGIS Extension.

CREATE EXTENSION postqis;

Using PostGIS Extension

Use the following function to invoke a PostGIS Extension:

SELECT GisFunction (Param1, Param2,.....);

GisFunction is the function, and **Param1** and **Param2** are function parameters. The following SQL statements are a simple illustration for PostGIS use. For details about related functions, see *PostGIS 2.4.2 Manual*.

Example 1: Create a geometry table.

CREATE TABLE cities (id integer, city_name varchar(50)); SELECT AddGeometryColumn('cities', 'position', 4326, 'POINT', 2);

Example 2: Insert geometry data.

INSERT INTO cities (id, position, city_name) VALUES (1,ST_GeomFromText('POINT(-9.5 23)',4326),'CityA'); INSERT INTO cities (id, position, city_name) VALUES (2,ST_GeomFromText('POINT(-10.6 40.3)',4326),'CityB'); INSERT INTO cities (id, position, city_name) VALUES (3,ST_GeomFromText('POINT(20.8 30.3)',4326), 'CityC');

Example 3: Calculate the distance between any two cities among three cities.

SELECT p1.city_name,p2.city_name,ST_Distance(p1.position,p2.position) FROM cities AS p1, cities AS p2 WHERE p1.id > p2.id;

Deleting PostGIS Extension

Run the following command to delete PostGIS Extension from GaussDB(DWS):

DROP EXTENSION postgis [CASCADE];

If PostGIS Extension is the dependee of other objects (for example, geometry tables), you need to add the **CASCADE** keyword to delete all these objects.

10.3 PostGIS Support and Constraints

Supported Data Types

In GaussDB(DWS), PostGIS Extension support the following data types:

- box2d
- box3d
- geometry_dump
- geometry
- geography
- raster

Ⅲ NOTE

If PostGIS is used by a user other than the creator of the PostGIS, set the following GUC parameters:

SET behavior_compat_options = 'bind_procedure_searchpath';

Supported Operators and Functions

□ NOTE

The **ST_Intersects** function in PostGIS uses a caching strategy that enables a high cache hit ratio for the spatial data structures of foreign tables. When there is a significant disparity in the width between the inner and foreign tables, caching the wide table's data avoid the repeated loading of large objects, leading to significant performance enhancements. Practically, leveraging **Join Order Hints** to designate a wide table as the foreign table ensures that the execution plan is optimized for such scenarios.

Table 10-1 Operators and functions supported by PostGIS2.4.2

Category	Function
Management functions	AddGeometryColumn, DropGeometryColumn, DropGeometryTable, PostGIS_Full_Version, PostGIS_GEOS_Version, PostGIS_Liblwgeom_Version, PostGIS_Lib_Build_Date, PostGIS_Lib_Version, PostGIS_PROJ_Version, PostGIS_Scripts_Build_Date, PostGIS_Scripts_Installed, PostGIS_Version, PostGIS_LibXML_Version, PostGIS_Scripts_Released, Populate_Geometry_Columns, UpdateGeometrySRID

Category	Function
Geometry constructors	ST_BdPolyFromText, ST_BdMPolyFromText, ST_Box2dFromGeoHash, ST_GeogFromText, ST_GeographyFromText, ST_GeogFromWKB, ST_GeomCollFromText, ST_GeomFromEWKB, ST_GeomFromEWKT, ST_GeometryFromText, ST_GeomFromGeoHash, ST_GeomFromGML, ST_GeomFromGeoJSON, ST_GeomFromKML, ST_GMLToSQL, ST_GeomFromText, ST_GeomFromWKB, ST_LineFromMultiPoint, ST_LineFromText, ST_LineFromWKB, ST_LinestringFromWKB, ST_MakeBox2D, ST_3DMakeBox, ST_MakeEnvelope, ST_MakePolygon, ST_MakePoint, ST_MakePointM, ST_MLineFromText, ST_MPointFromText, ST_MPolyFromText, ST_Point, ST_PointFromGeoHash, ST_PointFromText, ST_PointFromWKB, ST_Polygon, ST_PolygonFromText, ST_WKBToSQL, ST_WKTToSQL
Geometry accessors	GeometryType, ST_Boundary, ST_CoordDim, ST_Dimension, ST_EndPoint, ST_Envelope, ST_ExteriorRing, ST_GeometryN, ST_GeometryType, ST_InteriorRingN, ST_IsClosed, ST_IsCollection, ST_IsEmpty, ST_IsRing, ST_IsSimple, ST_IsValid, ST_IsValidReason, ST_IsValidDetail, ST_M, ST_NDims, ST_NPoints, ST_NRings, ST_NumGeometries, ST_NumInteriorRings, ST_NumPatches, ST_NumPoints, ST_PatchN, ST_PointN, ST_SRID, ST_StartPoint, ST_Summary, ST_X, ST_XMax, ST_XMin, ST_Y, ST_YMax, ST_YMin, ST_Z, ST_ZMax, ST_Zmflag, ST_ZMin
Geometry editors	ST_AddPoint, ST_Affine, ST_Force2D, ST_Force3D, ST_Force3DZ, ST_Force3DM, ST_Force4D, ST_ForceCollection, ST_ForceSFS, ST_ForceRHR, ST_LineMerge, ST_CollectionExtract, ST_CollectionHomogenize, ST_Multi, ST_RemovePoint, ST_Reverse, ST_Rotate, ST_RotateX, ST_RotateY, ST_RotateZ, ST_Scale, ST_Segmentize, ST_SetPoint, ST_SetSRID, ST_SnapToGrid, ST_Snap, ST_Transform, ST_Translate, ST_TransScale
Geometry outputs	ST_AsBinary, ST_AsEWKB, ST_AsEWKT, ST_AsGeoJSON, ST_AsGML, ST_AsHEXEWKB, ST_AsKML, ST_AsLatLonText, ST_AsSVG, ST_AsText, ST_AsX3D, ST_GeoHash
Operators	&&, &&&, &<, &< , &>, <<, << , =, >>, @, &>, >>, ~, ~=, <->, <#>

Category	Function
Spatial relationships and measurements	ST_3DClosestPoint, ST_3DDistance, ST_3DDWithin, ST_3DDFullyWithin, ST_3DIntersects, ST_3DLongestLine, ST_3DMaxDistance, ST_3DShortestLine, ST_Area, ST_Azimuth, ST_Centroid, ST_ClosestPoint, ST_Contains, ST_ContainsProperly, ST_Covers, ST_CoveredBy, ST_Crosses, ST_LineCrossingDirection, ST_Disjoint, ST_Distance, ST_HausdorffDistance, ST_MaxDistance, ST_DistanceSphere, ST_DistanceSpheroid, ST_DFullyWithin, ST_DWithin, ST_Equals, ST_HasArc, ST_Intersects, ST_Length, ST_Length2D, ST_3DLength, ST_Length_Spheroid, ST_Length2D_Spheroid, ST_3DLength_Spheroid, ST_LongestLine, ST_OrderingEquals, ST_Overlaps, ST_Perimeter, ST_Perimeter, ST_Perimeter, ST_RelateMatch, ST_ShortestLine, ST_Touches, ST_Within
Geometry processing	ST_Buffer, ST_BuildArea, ST_Collect, ST_ConcaveHull, ST_ConvexHull, ST_CurveToLine, ST_DelaunayTriangles, ST_Difference, ST_Dump, ST_DumpPoints, ST_DumpRings, ST_FlipCoordinates, ST_Intersection, ST_LineToCurve, ST_MakeValid, ST_MemUnion, ST_MinimumBoundingCircle, ST_Polygonize, ST_Node, ST_OffsetCurve, ST_RemoveRepeatedPoints, ST_SharedPaths, ST_Shift_Longitude, ST_Simplify, ST_SimplifyPreserveTopology, ST_Split, ST_SymDifference, ST_Union, ST_UnaryUnion
Linear referencing	ST_LineInterpolatePoint, ST_LineLocatePoint, ST_LineSubstring, ST_LocateAlong, ST_LocateBetween, ST_LocateBetweenElevations, ST_InterpolatePoint, ST_AddMeasure
Miscellaneous functions	ST_Accum, Box2D, Box3D, ST_Expand, ST_Extent, ST_3Dextent, Find_SRID, ST_MemSize
Exceptional functions	PostGIS_AddBBox, PostGIS_DropBBox, PostGIS_HasBBox
Raster Management Functions	AddRasterConstraints, DropRasterConstraints, AddOverviewConstraints, DropOverviewConstraints, PostGIS_GDAL_Version, PostGIS_Raster_Lib_Build_Date, PostGIS_Raster_Lib_Version, and ST_GDALDrivers, and UpdateRasterSRID
Raster Constructors	ST_AddBand, ST_AsRaster, ST_Band, ST_MakeEmptyRaster, ST_Tile, and ST_FromGDALRaster

Category	Function
Raster Accessors	ST_GeoReference, ST_Height, ST_IsEmpty, ST_MetaData, ST_NumBands, ST_PixelHeight, ST_PixelWidth, ST_ScaleX, ST_ScaleY, ST_RasterToWorldCoord, ST_RasterToWorldCoordX, ST_RasterToWorldCoordY, ST_Rotation, ST_SkewX, ST_SkewY, ST_SRID, ST_Summary, ST_UpperLeftX, ST_UpperLeftY, ST_Width, ST_WorldToRasterCoord, ST_WorldToRasterCoordX, ST_WorldToRasterCoordY
Raster Band Accessors	ST_BandMetaData, ST_BandNoDataValue, ST_BandIsNoData, ST_BandPath, ST_BandPixelType, and ST_HasNoBand
Raster Pixel Accessors and Setters	ST_PixelAsPolygon, ST_PixelAsPolygons, ST_PixelAsPoint, ST_PixelAsPoints, ST_PixelAsCentroid, ST_PixelAsCentroids, ST_Value, ST_NearestValue, ST_Neighborhood, ST_SetValue, ST_SetValues, ST_DumpValues, and ST_PixelOfValue
Raster Editors	ST_SetGeoReference, ST_SetRotation, ST_SetScale, ST_SetSkew, ST_SetSRID, ST_SetUpperLeft, ST_Resample, ST_Rescale, ST_Reskew, and ST_SnapToGrid, ST_Resize, and ST_Transform
Raster Band Editors	ST_SetBandNoDataValue and ST_SetBandIsNoData
Raster Band Statistics and Analytics	ST_Count, ST_CountAgg, ST_Histogram, ST_Quantile, ST_SummaryStats, ST_SummaryStatsAgg, and ST_ValueCount
Raster Outputs	ST_AsBinary, ST_AsGDALRaster, ST_AsJPEG, ST_AsPNG, and ST_AsTIFF
Raster Processing	ST_Clip, ST_ColorMap, ST_Intersection, ST_MapAlgebra, ST_Reclass, and ST_Union ST_Distinct4ma, ST_InvDistWeight4ma, ST_Max4ma, ST_Mean4ma, ST_Min4ma, ST_MinDist4ma, ST_Range4ma, ST_StdDev4ma, and ST_Sum4ma, ST_Aspect, ST_HillShade, ST_Roughness, ST_Slope, ST_TPI, ST_TRI, Box3D, ST_ConvexHull, ST_DumpAsPolygons, and ST_ Envelope, ST_MinConvexHull, ST_Polygon, ST_Contains, ST_ContainsProperly, ST_Covers, ST_CoveredBy, ST_Disjoint, ST_Intersects, and ST_Overlaps, ST_Touches, ST_SameAlignment, ST_NotSameAlignmentReason, ST_Within, ST_DWithin, and ST_DFullyWithin
Raster Operators	&&, &<, &>, =, @, ~=, and ~

Table 10-2 Operators and functions supported by PostGIS3.2.2

Category	Function
Management functions	AddGeometryColumn, DropGeometryColumn, DropGeometryTable, PostGIS_Full_Version, PostGIS_GEOS_Version, PostGIS_Liblwgeom_Version, PostGIS_Lib_Build_Date, PostGIS_Lib_Version, PostGIS_PROJ_Version, PostGIS_Scripts_Build_Date, PostGIS_Scripts_Installed, PostGIS_Version, PostGIS_LibXML_Version, PostGIS_Scripts_Released, Populate_Geometry_Columns, UpdateGeometrySRID, PostGIS_Libprotobuf_Version, PostGIS_Wagyu_Version
Geometry constructors	ST_BdPolyFromText, ST_BdMPolyFromText, ST_Box2dFromGeoHash, ST_GeneratePoints, ST_GeogFromText, ST_GeographyFromText, ST_GeogFromWKB, ST_GeomCollFromText, ST_GeomFromEWKB, ST_GeomFromEWKT, ST_GeomFromEWKB, ST_GeomFromGeoHash, ST_GeomFromGML, ST_GeomFromGeoJSON, ST_GeomFromKML, ST_GMLToSQL, ST_GeomFromText, ST_GeomFromWKB, ST_LineFromMultiPoint, ST_LineFromText, ST_LineFromWKB, ST_LinestringFromWKB, ST_MakeBox2D, ST_3DMakeBox, ST_MakeEnvelope, ST_MakePolygon, ST_MakePoint, ST_MakePointM, ST_MLineFromText, ST_MPointFromText, ST_MPolyFromText, ST_Point, ST_Points, ST_PointFromGeoHash, ST_PointFromText, ST_PointFromWKB, ST_Polygon, ST_PolygonFromText, ST_WKBToSQL, ST_WKTToSQL, Geography_Distance_Knn, Geometry_Distance_Cpa, Geometry_Hash, ST_3Dlineinterpolate, ST_AsEncodedPolyline
Geometry accessors	GeometryType, ST_Boundary, ST_CoordDim, ST_Dimension, ST_EndPoint, ST_Envelope, ST_ExteriorRing, ST_GeometryN, ST_GeometryType, ST_InteriorRingN, ST_IsClosed, ST_IsCollection, ST_IsEmpty, ST_IsPolygonCCW, ST_IsPolygonCW, ST_IsRing, ST_IsSimple, ST_IsValid, ST_IsValidReason, ST_IsValidDetail, ST_M, ST_NDims, ST_NPoints, ST_NRings, ST_NumGeometries, ST_NumInteriorRings, ST_NumInteriorRing, ST_NumPatches, ST_NumPoints, ST_PatchN, ST_PointN, ST_SRID, ST_StartPoint, ST_Summary, ST_X, ST_XMax, ST_XMin, ST_Y, ST_YMax, ST_YMin, ST_Z, ST_ZMax, ST_Zmflag, ST_ZMin, ST_Wrapx, ST_Asmvt

Category	Function
Geometry editors	ST_AddPoint, ST_Affine, ST_Force2D, ST_Force3D, ST_Force3DZ, ST_Force3DM, ST_Force4D, ST_ForceCollection, ST_ForcePolygonCCW, ST_ForcePolygonCW, ST_ForceSFS, ST_ForceRHR, ST_LineMerge, ST_CollectionExtract, ST_CollectionHomogenize, ST_Multi, ST_Normalize, ST_RemovePoint, ST_Reverse, ST_Rotate, ST_RotateX, ST_RotateY, ST_RotateZ, ST_Scale, ST_Segmentize, ST_SetPoint, ST_SetSRID, ST_SnapToGrid, ST_Snap, ST_Transform, ST_Translate, ST_TransScale, ST_AsmvtGeom, ST_isvalidTrajectory, ST_linefromencodedpolyline, ST_lineinterpolatepoints, ST_MaximuminScribedCircle, ST_OrientedEnvelope, ST_QuantizeCoordinates, ST_ReducePrecision, ST_Scroll, ST_SetEffectiveArea, ST_simplifyvw, ST_square, ST_squaregrid, ST_Swapordinates, ST_Voronoilines, ST_VoronoiPolygons
Geometry outputs	ST_AsBinary, ST_AsEWKB, ST_AsEWKT, ST_AsGeoJSON, ST_AsGML, ST_AsHEXEWKB, ST_AsKML, ST_AsLatLonText, ST_AsSVG, ST_AsText, ST_AsTwkb, ST_AsX3D, ST_GeoHash, Json, Jsonb, ST_GeomfromGeojson
Operators	&& , &&& , &< , &< , &> , << , << , =, >> , @ , &> , >> , ~, ~=, <-> , <#> , <-> , = , <<->>
Spatial relationships and measurements	ST_3DClosestPoint, ST_3DDistance, ST_3DDWithin, ST_3DDFullyWithin, ST_3DIntersects, ST_3DLongestLine, ST_3DMaxDistance, ST_3DShortestLine, ST_Area, ST_Azimuth, ST_Centroid, ST_ClosestPoint, ST_Contains, ST_ContainsProperly, ST_Covers, ST_CoveredBy, ST_Crosses, ST_LineCrossingDirection, ST_Disjoint, ST_Distance, ST_HausdorffDistance, ST_MaxDistance, ST_DistanceSphere, ST_DistanceSpheroid, ST_DFullyWithin, ST_DWithin, ST_Equals, ST_HasArc, ST_Intersects, ST_Length, ST_Length2D, ST_3DLength, ST_LengthSpheroid, ST_Length2DSpheroid, ST_LengthSpheroid, ST_Length2DSpheroid, ST_LongestLine, ST_MinimumBoundingRadius, ST_OrderingEquals, ST_Overlaps, ST_Perimeter, ST_Perimeter2D, ST_3DPerimeter, ST_PointOnSurface, ST_Project, ST_Relate, ST_RelateMatch, ST_ShortestLine, ST_Touches, ST_Within, _ST_DistancerectTree, _ST_DistancerectTreeCached, _ST_SorTableHash

Category	Function
Geometry processing	ST_Buffer, ST_BuildArea, ST_ClipByBox2D, ST_ClusterDBSCAN, ST_ClusterIntersecting, ST_ClusterKMeans, ST_ClusterWithin, ST_Collect, ST_ConcaveHull, ST_ConvexHull, ST_CurveToLine, ST_DelaunayTriangles, ST_Difference, ST_Dump, ST_DumpPoints, ST_DumpRings, ST_FlipCoordinates, ST_Intersection, ST_LineToCurve, ST_MakeValid, ST_MemUnion, ST_MinimumBoundingCircle, ST_Polygonize, ST_Node, ST_OffsetCurve, ST_RemoveRepeatedPoints, ST_SharedPaths, ST_ShiftLongitude, ST_Simplify, ST_SimplifyPreserveTopology, ST_Split, ST_Subdivide, ST_SymDifference, ST_Union, ST_UnaryUnion, ST_BoundingDiagonal, ST_ChaikinsMoothing, ST_ClosestPointofApproach, ST_CollectionExtract, ST_CPAwithin, ST_DistanceCPA, ST_DumpSegments, ST_EstimatedExtent, ST_Filterbym, ST_SetEffectiveArea, ST_Forcecurve
Linear referencing	ST_LineInterpolatePoint, ST_LineLocatePoint, ST_LineSubstring, ST_LocateAlong, ST_LocateBetween, ST_LocateBetweenElevations, ST_InterpolatePoint, ST_AddMeasure
Miscellaneous functions	Array_Agg, Box2D, Box3D, ST_Expand, ST_Extent, ST_3Dextent, Find_SRID, ST_MemSize
Exceptional functions	PostGIS_AddBBox, PostGIS_DropBBox, PostGIS_HasBBox

Spatial Indexes

In GaussDB(DWS), PostGIS Extension supports Generalized Search Tree (GiST) spatial indexes. This index type is inapplicable to partitioned tables. Different from B-tree indexes, GiST indexes are suitable for any type of unconventional data structure and can effectively improve the retrieval efficiency of geometric and geographic data.

Run the following command to create a GiST index:

CREATE INDEX indexname ON tablename USING gist (geometryfield);

Extension Constraints

- Only row-store tables are supported. Column-store indexes are not supported.
- Only Oracle-compatible databases are supported.
- The topology object management module, Topology, is not supported.
- BRIN indexes are not supported.
- The **spatial_ref_sys** table can only be queried during scale-out.

Plug-in Upgrade Compatibility

When upgrading from PostGIS 2.4.2 to 3.2.2, note that certain functions may become incompatible or not fully forward compatible. This can lead to inconsistencies in the functionality before and after the upgrade. Therefore, it is necessary to assess the impact of upgrade incompatibility on your services.

The following table provides compatibility details for related functions.

Category	PostGIS 2.4.2	PostGIS 3.2.2
Added functions in 3.2.2	N/A	ST_IsPolygonCW(geometry)
	N/A	ST_IsPolygonCCW(geometry)
	N/A	ST_PointInsideCircle(geometry,floa t8,float8,float8)
	N/A	ST_ForcePolygonCW(geometry)
	N/A	ST_ForcePolygonCCW(geometry)
	N/A	ST_Normalize(geom geometry)
	N/A	ST_AsTWKB(geom geometry, prec int4 default 0, prec_z int4 default 0, prec_m int4 default 0, with_sizes boolean default false, with_boxes boolean default false)
	N/A	ST_AsTWKB(geom geometry[], ids bigint[], prec int4 default 0, prec_z int4 default 0, prec_m int4 default 0, with_sizes boolean default false, with_boxes boolean default false)
	N/A	ST_MakeLine (geometry[])
	N/A	ST_TileEnvelope(zoom integer, x integer, y integer, bounds geometry DEFAULT 'SRID=3857;LINESTRING(-2003750 8.342789244 -20037508.342789244, 20037508.342789244 20037508.342789244)'::geometry, margin float8 DEFAULT 0.0)
	N/A	ST_ClusterIntersect- ing(geometry[])
	N/A	ST_ClusterWithin(geometry[], float8)
	N/A	ST_ClusterDBSCAN (geometry, eps float8, minpoints int)

Category	PostGIS 2.4.2	PostGIS 3.2.2
	N/A	ST_Scale(geometry,geometry,origi n geometry)
	N/A	ST_GeneratePoints(area geometry, npoints integer, seed integer)
	N/A	ST_FrechetDistance(geom1 geometry, geom2 geometry, float8 default -1)
	N/A	ST_Points(geometry)
	N/A	ST_ClipByBox2d(geom geometry, box box2d)
	N/A	ST_Subdivide(geom geometry, maxvertices integer DEFAULT 256, gridSize float8 DEFAULT -1.0)
	N/A	ST_ClusterIntersecting (geometry)
	N/A	ST_ClusterWithin (geometry, float8)
	N/A	ST_ClusterKMeans(geom geometry, k integer, max_radius float8 default null)
	N/A	ST_AsText(geometry, int4)
	N/A	ST_AsEWKT(geography, int4)
	N/A	_ST_CoveredBy(geog1 geography, geog2 geography)
	N/A	ST_Point(float8, float8, srid integer)
Functions no longer	ST_3DLength_spheroid(ge ometry, spheroid)	N/A
supported in 3.2.2	ST_length2d_spheroid(geo metry, spheroid)	N/A
	ST_locate_between_measu res(geometry, float8, float8)	N/A
	ST_locate_along_measure(geometry, float8)	N/A
	ST_Buffer(geometry,float8, text)	N/A
	ST_GeneratePoints(area geometry, npoints integer)	N/A

Category	PostGIS 2.4.2	PostGIS 3.2.2
	ST_Combine_BBox(box3d, geometry)	N/A
	ST_Combine_BBox(box2d, geometry)	N/A
	pgis_abs_in(cstring)	N/A
	pgis_abs_out(pgis_abs)	N/A
	pgis_abs (internallength = 16, input = pgis_abs_in, output = pgis_abs_out, alignment = double)	N/A
	ST_MemUnion(geometry)	N/A
	pgis_geometry_accum_fin alfn(pgis_abs)	N/A
	ST_MakeLine (geometry)	N/A
	ST_Accum (geometry)	When a single geometry data record is input, its value is output directly. However, if multiple geometry data records are provided, an error message will be displayed, indicating that multiple records are unsupported.
	_ST_AsKML(int4,geometry, int4, text)	N/A
	ST_MemUnion(geometry)	N/A
	_ST_AsGeoJson(int4, geometry, int4, int4)	N/A
	ST_AsGeoJson(gj_version int4, geom geometry, maxdecimaldigits int4 DEFAULT 15, options int4 DEFAULT 0)	N/A
	_ST_DWithin(geography, geography, float8, boolean)	N/A
	ST_point_inside_circle(geo metry,float8,float8,float8)	N/A
	ST_CurveToLine(geometry)	N/A

Category	PostGIS 2.4.2	PostGIS 3.2.2
	ST_Shift_Longitude(geome try)	Use ST_ShiftLongitude instead.
	ST_find_extent(text,text,te xt) and ST_find_extent(text,text)	Use ST_FindExtent instead.
	ST_mem_size(geometry)	Use ST_MemSize instead.
	ST_length_spheroid(geom etry, spheroid)	Use ST_LengthSpheroid instead.
	ST_distance_spheroid(geo m1 geometry, geom2 geometry,spheroid)	Use ST_DistanceSpheroid instead.
	ST_force_2d(geometry)	Use ST_Force2D instead.
	ST_force_3dz(geometry)	Use ST_Force3DZ instead.
	ST_force_3d(geometry)	Use ST_Force3D instead.
	ST_force_3dm(geometry)	Use ST_Force3DM instead.
	ST_force_4d(geometry)	Use ST_Force4D instead.
	ST_force_collection(geome try)	Use ST_ForceCollection instead.
	ST_line_locate_point(geo m1 geometry, geom2 geometry)	Use ST_LineLocatePoint instead.
	ST_line_interpolate_point(geometry, float8)	Use ST_LineInterpolatePoint instead.
	ST_Buffer(geometry,float8)	Use ST_Buffer instead.
Functions with parameter type changed in 3.2.2	pgis_geometry_accum_tra nsfn(pgis_abs, geometry)	pgis_geometry_accum_transfn(inte rnal, geometry)
	pgis_geometry_accum_tra nsfn(pgis_abs, geometry, float8)	pgis_geometry_accum_transfn(inte rnal, geometry, float8)
	pgis_geometry_accum_tra nsfn(pgis_abs, geometry, float8, int)	pgis_geometry_accum_transfn(inte rnal, geometry, float8, int)
	pgis_geometry_union_final fn(pgis_abs)	pgis_geometry_union_finalfn(inter nal)
	pgis_geometry_collect_fin alfn(pgis_abs)	pgis_geometry_collect_finalfn(inter nal)

Category	PostGIS 2.4.2	PostGIS 3.2.2
	pgis_geometry_polygonize _finalfn(pgis_abs)	pgis_geometry_polygonize_finalfn(internal)
	pgis_geometry_clusterinter secting_finalfn(pgis_abs)	pgis_geometry_clusterintersect- ing_finalfn(internal)
	ST_Union (geometry)	ST_Union (geometry)
	ST_Collect (geometry)	ST_Collect (geometry)
	ST_Buffer(geometry,float8, integer)	ST_Buffer(geom geometry, radius float8, quadsegs integer)
Functions with API changed in 3.2.2	ST_AsKML(version int4, geom geometry, maxdecimaldigits int4 DEFAULT 15, nprefix text DEFAULT null)	ST_AsKML(geom geometry, maxdecimaldigits int4 DEFAULT 15, nprefix TEXT default ' ')
	ST_AsKML(geom geometry, maxdecimaldigits int4 DEFAULT 15)	ST_AsKML(geom geometry, maxdecimaldigits int4 DEFAULT 15, nprefix TEXT default ' ')
	ST_AsKML(version int4, geom geometry, maxdecimaldigits int4 DEFAULT 15, nprefix text DEFAULT null)	ST_AsGML(version int4, geog geography, maxdecimaldigits int4 DEFAULT 15, options int4 DEFAULT 0, nprefix text DEFAULT 'gml', id text DEFAULT '')
Functions with default parameters changed in 3.2.2	ST_SymDifference(geom1 geometry, geom2 geometry)	ST_SymDifference(geom1 geometry, geom2 geometry, gridSize float8 DEFAULT -1.0)
	ST_UnaryUnion(geometry)	ST_UnaryUnion(geometry, gridSize float8 DEFAULT -1.0)
	ST_AsGeoJson(geom geometry, maxdecimaldigits int4 DEFAULT 15, options int4 DEFAULT 0)	ST_AsGeoJson(geom geometry, maxdecimaldigits int4 DEFAULT 9, options int4 DEFAULT 8)
	ST_Buffer(geometry,float8, cstring)	ST_Buffer(geom geometry, radius float8, options text DEFAULT ' ')
	_ST_DWithin(geography, geography, float8, boolean)	_ST_DWithin(geog1 geography, geog2 geography, tolerance float8, use_spheroid boolean DEFAULT true)
	ST_IsValidDetail(geometry)	ST_IsValidDetail(geom geometry, flags int4 DEFAULT 0)

Category	PostGIS 2.4.2	PostGIS 3.2.2
	ST_CurveToLine(geom geometry, tol float8, toltype integer, flags integer)	ST_CurveToLine(geom geometry, tol float8 DEFAULT 32, toltype integer DEFAULT 0, flags integer DEFAULT 0)
Functions supported hash join and merge join in 3.2.2	OPERATOR =	OPERATOR =
Function changes in version 3.2.2	ST_ConcaveHull weak verification	Strong verification is added to ST_ConcaveHull to ensure more rigorous verification.
Functions with commutor undefined in 3.2.2	OPERATOR &< ,OPERATOR &< ,OPERATOR &>	OPERATOR &< ,OPERATOR &< ,OPERATOR &>

10.4 OPEN SOURCE SOFTWARE NOTICE (For PostGIS)

This document contains open source software notice for the product. And this document is confidential information of copyright holder. Recipient shall protect it in due care and shall not disseminate it without permission.

Warranty Disclaimer

This document is provided "as is" without any warranty whatsoever, including the accuracy or comprehensiveness. Copyright holder of this document may change the contents of this document at any time without prior notice, and copyright holder disclaims any liability in relation to recipient's use of this document.

Open source software is provided by the author "as is" and any express or implied warranties, including, but not limited to, the implied warranties of merchantability and fitness for a particular purpose are disclaimed. In no event shall the author be liable for any direct, indirect, incidental, special, exemplary, or consequential damages (including, but not limited to, procurement of substitute goods or services; loss of data or profits; or business interruption) however caused and on any theory of liability, whether in contract, strict liability, or tort (including negligence or otherwise) arising in any way out of the use of open source software, even if advised of the possibility of such damage.

Copyright Notice And License Texts

Software: postgis-2.4.2

Copyright notice:

"Copyright (C) 1996-2015 Free Software Foundation, Inc.

Copyright (C) 1989, 1991 Free Software Foundation, Inc.,

51 Franklin Street, Fifth Floor, Boston, MA 02110-1301

Copyright 2008 Kevin Neufeld

Copyright (c) 2009 Walter Bruce Sinclair

Copyright 2006-2013 Stephen Woodbridge.

Copyright (c) 2008 Walter Bruce Sinclair

Copyright (c) 2012 TJ Holowaychuk <tj@vision-media.ca>

Copyright (c) 2008, by Attractive Chaos <attractivechaos@aol.co.uk>

Copyright (c) 2001-2012 Walter Bruce Sinclair

Copyright (c) 2010 Walter Bruce Sinclair

Copyright 2006 Stephen Woodbridge

Copyright 2006-2010 Stephen Woodbridge.

Copyright (c) 2006-2014 Stephen Woodbridge.

Copyright (c) 2017, Even Rouault <even.rouault at spatialys.com>

Copyright (C) 2004-2015 Sandro Santilli <strk@kbt.io>

Copyright (C) 2008 Mark Cave-Ayland <mark.cave-ayland@siriusit.co.uk>

Copyright 2015 Nicklas Avén <nicklas.aven@jordogskog.no>

Copyright 2008 Paul Ramsey

Copyright (C) 2012 Sandro Santilli <strk@kbt.io>

Copyright 2012 Sandro Santilli <strk@kbt.io>

Copyright (C) 2014 Sandro Santilli <strk@kbt.io>

Copyright 2013 Olivier Courtin <olivier.courtin@oslandia.com>

Copyright 2011 Sandro Santilli <strk@kbt.io>

Copyright 2015 Daniel Baston

Copyright 2009 Olivier Courtin <olivier.courtin@oslandia.com>

Copyright 2014 Kashif Rasul kashif.rasul@gmail.com and

Shoaib Burq <saburq@gmail.com>

Copyright 2013 Sandro Santilli <strk@kbt.io>

Copyright 2010 Paul Ramsey cleverelephant.ca>

Copyright (C) 2017 Sandro Santilli <strk@kbt.io>

Copyright (C) 2015 Sandro Santilli <strk@kbt.io>

Copyright (C) 2009 Paul Ramsey cleverelephant.ca>

Copyright (C) 2011 Sandro Santilli <strk@kbt.io>

Copyright 2010 Olivier Courtin <olivier.courtin@oslandia.com>

Copyright 2014 Nicklas Avén

Copyright 2011-2016 Regina Obe

Copyright (C) 2008 Paul Ramsey

Copyright (C) 2011-2015 Sandro Santilli <strk@kbt.io>

Copyright 2010-2012 Olivier Courtin <olivier.courtin@oslandia.com>

Copyright (C) 2015 Daniel Baston dbaston@gmail.com/

Copyright (C) 2013 Nicklas Avén

Copyright (C) 2016 Sandro Santilli <strk@kbt.io>

Copyright 2017 Darafei Praliaskouski <me@komzpa.net>

Copyright (C) 2011-2012 Sandro Santilli <strk@kbt.io>

Copyright (C) 2007-2008 Mark Cave-Ayland

Copyright (C) 2001-2006 Refractions Research Inc.

Copyright 2015 Daniel Baston dbaston@gmail.com/

Copyright 2009 David Skea < David. Skea@gov.bc.ca>

Copyright (C) 2012-2015 Sandro Santilli <strk@kbt.io>

Copyright 2001-2006 Refractions Research Inc.

Copyright (C) 2004 Refractions Research Inc.

Copyright 2011-2014 Sandro Santilli <strk@kbt.io>

Copyright 2009-2010 Sandro Santilli <strk@kbt.io>

Copyright 2015-2016 Daniel Baston dbaston@gmail.com/

Copyright 2011-2015 Sandro Santilli <strk@kbt.io>

Copyright 2007-2008 Mark Cave-Ayland

Copyright 2012-2013 Oslandia <infos@oslandia.com>

Copyright (C) 2015-2017 Sandro Santilli <strk@kbt.io>

Copyright (C) 2001-2003 Refractions Research Inc.

Copyright 2016 Sandro Santilli <strk@kbt.io>

Copyright 2011 Kashif Rasul <kashif.rasul@gmail.com>

Copyright (C) 2014 Nicklas Avén

Copyright (C) 2011 Sandro Santilli <strk@kbt.io>

Copyright (C) 2011-2014 Sandro Santilli <strk@kbt.io>

Copyright (C) 1984, 1989-1990, 2000-2015 Free Software Foundation, Inc.

Copyright (C) 2011 Paul Ramsey

Copyright 2001-2003 Refractions Research Inc.

Copyright 2009-2010 Olivier Courtin <olivier.courtin@oslandia.com>

Copyright 2010-2012 Oslandia

Copyright 2006 Corporacion Autonoma Regional de Santander

Copyright 2013 Nicklas Avén

Copyright 2011-2016 Arrival 3D, Regina Obe

Copyright (C) 2009 David Skea < David. Skea@gov.bc.ca>

Copyright (C) 2017 Sandro Santilli <strk@kbt.io>

Copyright (C) 2010 - Oslandia

Copyright (C) 2006 Mark Leslie <mark.leslie@lisasoft.com>

Copyright (C) 2008-2009 Mark Cave-Ayland <mark.cave-ayland@siriusit.co.uk>

Copyright (C) 2010 Olivier Courtin <olivier.courtin@camptocamp.com>

Copyright 2010 Nicklas Avén

Copyright 2012 Paul Ramsey

Copyright 2011 Nicklas Avén

Copyright 2002 Thamer Alharbash

Copyright 2011 OSGeo

Copyright (C) 2008 Mark Cave-Ayland <mark.cave-ayland@siriusit.co.uk>

Copyright (C) 2004-2007 Refractions Research Inc.

Copyright 2010 LISAsoft Pty Ltd

Copyright 2010 Mark Leslie

Copyright (c) 1999, Frank Warmerdam

Copyright 2009 Mark Cave-Ayland <mark.cave-ayland@siriusit.co.uk>

Copyright (c) 2007, Frank Warmerdam

Copyright 2008 OpenGeo.org

Copyright (C) 2008 OpenGeo.org

Copyright (C) 2009 Mark Cave-Ayland <mark.cave-ayland@siriusit.co.uk>

Copyright 2010 LISAsoft

Copyright (C) 2010 Mark Cave-Ayland <mark.cave-ayland@siriusit.co.uk>

Copyright (c) 1999, 2001, Frank Warmerdam

Copyright (C) 2016-2017 Bj?rn Harrtell <bjorn@wololo.org>

Copyright (C) 2017 Danny G?tte <danny.goette@fem.tu-ilmenau.de>

^copyright^

Copyright 2012 (C) Paul Ramsey cleverelephant.ca>

Copyright (C) 2006 Refractions Research Inc.

Copyright 2001-2009 Refractions Research Inc.

Copyright (C) 2010 Olivier Courtin <olivier.courtin@oslandia.com>

By Nathan Wagner, copyright disclaimed,

this entire file is in the public domain

Copyright 2009-2011 Olivier Courtin <olivier.courtin@oslandia.com>

Copyright (C) 2001-2005 Refractions Research Inc.

Copyright 2001-2011 Refractions Research Inc.

Copyright 2009-2014 Sandro Santilli <strk@kbt.io>

Copyright (C) 2008 Paul Ramsey cleverelephant.ca>

Copyright (C) 2007 Refractions Research Inc.

Copyright (C) 2010 Sandro Santilli <strk@kbt.io>

Copyright 2012 J Smith <dark.panda@gmail.com>

Copyright 2009 - 2010 Oslandia

Copyright 2009 Oslandia

Copyright 2001-2005 Refractions Research Inc.

Copyright 2016 Paul Ramsey cleverelephant.ca>

Copyright 2016 Daniel Baston dbaston@gmail.com/

Copyright (C) 2011 OpenGeo.org

Copyright (c) 2003-2017, Troy D. Hanson http:troydhanson.github.com/uthash/

Copyright (C) 2011 Regents of the University of California

Copyright (C) 2011-2013 Regents of the University of California

Copyright (C) 2010-2011 Jorge Arevalo <jorge.arevalo@deimos-space.com>

Copyright (C) 2010-2011 David Zwarg <dzwarg@azavea.com>

Copyright (C) 2009-2011 Pierre Racine <pierre.racine@sbf.ulaval.ca>

Copyright (C) 2009-2011 Mateusz Loskot <mateusz@loskot.net>

Copyright (C) 2008-2009 Sandro Santilli <strk@kbt.io>

Copyright (C) 2013 Nathaneil Hunter Clay <clay.nathaniel@gmail.com

Copyright (C) 2013 Nathaniel Hunter Clay <clay.nathaniel@gmail.com>

Copyright (C) 2013 Bborie Park <dustymugs@gmail.com>

Copyright (C) 2013 Nathaniel Hunter Clay <clay.nathaniel@gmail.com>

(C) 2009 Mateusz Loskot <mateusz@loskot.net>

Copyright (C) 2009 Mateusz Loskot <mateusz@loskot.net>

Copyright (C) 2009-2010 Mateusz Loskot <mateusz@loskot.net>

Copyright (C) 2009-2010 Jorge Arevalo < jorge.arevalo@deimos-space.com>

Copyright (C) 2012 Regents of the University of California

Copyright (C) 2013 Regents of the University of California

Copyright (C) 2012-2013 Regents of the University of California

Copyright (C) 2009 Sandro Santilli <strk@kbt.io>

"

License: The GPL v2 License.

GNU GENERAL PUBLIC LICENSE

Version 2, June 1991

Copyright (C) 1989, 1991 Free Software Foundation, Inc.

51 Franklin St, Fifth Floor, Boston, MA 02110-1301

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

Preamble

The licenses for most software are designed to take away your freedom to share and change it. By contrast, the GNU General Public License is intended to guarantee your freedom to share and change free software--to make sure the software is free for all its users. This General Public License applies to most of the Free Software Foundation's software and to any other program whose authors

commit to using it. (Some other Free Software Foundation software is covered by the GNU Library General Public License instead.) You can apply it to your programs, too.

When we speak of free software, we are referring to freedom, not price. Our General Public Licenses are designed to make sure that you have the freedom to distribute copies of free software (and charge for this service if you wish), that you receive source code or can get it if you want it, that you can change the software or use pieces of it in new free programs; and that you know you can do these things.

To protect your rights, we need to make restrictions that forbid anyone to deny you these rights or to ask you to surrender the rights. These restrictions translate to certain responsibilities for you if you distribute copies of the software, or if you modify it.

For example, if you distribute copies of such a program, whether gratis or for a fee, you must give the recipients all the rights that you have. You must make sure that they, too, receive or can get the source code. And you must show them these terms so they know their rights.

We protect your rights with two steps: (1) copyright the software, and (2) offer you this license which gives you legal permission to copy, distribute and/or modify the software.

Also, for each author's protection and ours, we want to make certain that everyone understands that there is no warranty for this free software. If the software is modified by someone else and passed on, we want its recipients to know that what they have is not the original, so that any problems introduced by others will not reflect on the original authors' reputations.

Finally, any free program is threatened constantly by software patents. We wish to avoid the danger that redistributors of a free program will individually obtain patent licenses, in effect making the program proprietary. To prevent this, we have made it clear that any patent must be licensed for everyone's free use or not licensed at all.

The precise terms and conditions for copying, distribution and modification follow.?

GNU GENERAL PUBLIC LICENSE

TERMS AND CONDITIONS FOR COPYING, DISTRIBUTION AND MODIFICATION

0. This License applies to any program or other work which contains a notice placed by the copyright holder saying it may be distributed under the terms of this General Public License. The "Program", below, refers to any such program or work, and a "work based on the Program" means either the Program or any derivative work under copyright law: that is to say, a work containing the Program or a portion of it, either verbatim or with modifications and/or translated into another language. (Hereinafter, translation is included without limitation in the term "modification".) Each licensee is addressed as "you".

Activities other than copying, distribution and modification are not covered by this License; they are outside its scope. The act of running the Program is not restricted, and the output from the Program is covered only if its contents constitute a work based on the Program (independent of having been made by running the Program). Whether that is true depends on what the Program does.

1. You may copy and distribute verbatim copies of the Program's source code as you receive it, in any medium, provided that you conspicuously and appropriately publish on each copy an appropriate copyright notice and disclaimer of warranty; keep intact all the notices that refer to this License and to the absence of any warranty; and give any other recipients of the Program a copy of this License along with the Program.

You may charge a fee for the physical act of transferring a copy, and you may at your option offer warranty protection in exchange for a fee.

- 2. You may modify your copy or copies of the Program or any portion of it, thus forming a work based on the Program, and copy and distribute such modifications or work under the terms of Section 1 above, provided that you also meet all of these conditions:
- a) You must cause the modified files to carry prominent notices stating that you changed the files and the date of any change.
- b) You must cause any work that you distribute or publish, that in whole or in part contains or is derived from the Program or any part thereof, to be licensed as a whole at no charge to all third parties under the terms of this License.
- c) If the modified program normally reads commands interactively when run, you must cause it, when started running for such interactive use in the most ordinary way, to print or display an announcement including an appropriate copyright notice and a notice that there is no warranty (or else, saying that you provide a warranty) and that users may redistribute the program under these conditions, and telling the user how to view a copy of this License. (Exception: if the Program itself is interactive but does not normally print such an announcement, your work based on the Program is not required to print an announcement.)

These requirements apply to the modified work as a whole. If identifiable sections of that work are not derived from the Program, and can be reasonably considered independent and separate works in themselves, then this License, and its terms, do not apply to those sections when you distribute them as separate works. But when you distribute the same sections as part of a whole which is a work based on the Program, the distribution of the whole must be on the terms of this License, whose permissions for other licensees extend to the entire whole, and thus to each and every part regardless of who wrote it.

Thus, it is not the intent of this section to claim rights or contest your rights to work written entirely by you; rather, the intent is to exercise the right to control the distribution of derivative or collective works based on the Program.

In addition, mere aggregation of another work not based on the Program with the Program (or with a work based on the Program) on a volume of a storage or

distribution medium does not bring the other work under the scope of this License.

- 3. You may copy and distribute the Program (or a work based on it, under Section 2) in object code or executable form under the terms of Sections 1 and 2 above provided that you also do one of the following:
- a) Accompany it with the complete corresponding machine-readable source code, which must be distributed under the terms of Sections 1 and 2 above on a medium customarily used for software interchange; or,
- b) Accompany it with a written offer, valid for at least three years, to give any third party, for a charge no more than your cost of physically performing source distribution, a complete machine-readable copy of the corresponding source code, to be distributed under the terms of Sections 1 and 2 above on a medium customarily used for software interchange; or,
- c) Accompany it with the information you received as to the offer to distribute corresponding source code. (This alternative is allowed only for noncommercial distribution and only if you received the program in object code or executable form with such an offer, in accord with Subsection b above.)

The source code for a work means the preferred form of the work for making modifications to it. For an executable work, complete source code means all the source code for all modules it contains, plus any associated interface definition files, plus the scripts used to control compilation and installation of the executable. However, as a special exception, the source code distributed need not include anything that is normally distributed (in either source or binary form) with the major components (compiler, kernel, and so on) of the operating system on which the executable runs, unless that component itself accompanies the executable.

If distribution of executable or object code is made by offering access to copy from a designated place, then offering equivalent access to copy the source code from the same place counts as distribution of the source code, even though third parties are not compelled to copy the source along with the object code.

- 4. You may not copy, modify, sublicense, or distribute the Program except as expressly provided under this License. Any attempt otherwise to copy, modify, sublicense or distribute the Program is void, and will automatically terminate your rights under this License. However, parties who have received copies, or rights, from you under this License will not have their licenses terminated so long as such parties remain in full compliance.
- 5. You are not required to accept this License, since you have not signed it. However, nothing else grants you permission to modify or distribute the Program or its derivative works. These actions are prohibited by law if you do not accept this License. Therefore, by modifying or distributing the Program (or any work based on the Program), you indicate your acceptance of this License to do so, and all its terms and conditions for copying, distributing or modifying the Program or works based on it.

- 6. Each time you redistribute the Program (or any work based on the Program), the recipient automatically receives a license from the original licensor to copy, distribute or modify the Program subject to these terms and conditions. You may not impose any further restrictions on the recipients' exercise of the rights granted herein. You are not responsible for enforcing compliance by third parties to this License.
- 7. If, as a consequence of a court judgment or allegation of patent infringement or for any other reason (not limited to patent issues), conditions are imposed on you (whether by court order, agreement or otherwise) that contradict the conditions of this License, they do not excuse you from the conditions of this License. If you cannot distribute so as to satisfy simultaneously your obligations under this License and any other pertinent obligations, then as a consequence you may not distribute the Program at all. For example, if a patent license would not permit royalty-free redistribution of the Program by all those who receive copies directly or indirectly through you, then the only way you could satisfy both it and this License would be to refrain entirely from distribution of the Program.

If any portion of this section is held invalid or unenforceable under any particular circumstance, the balance of the section is intended to apply and the section as a whole is intended to apply in other circumstances.

It is not the purpose of this section to induce you to infringe any patents or other property right claims or to contest validity of any such claims; this section has the sole purpose of protecting the integrity of the free software distribution system, which is implemented by public license practices. Many people have made generous contributions to the wide range of software distributed through that system in reliance on consistent application of that system; it is up to the author/donor to decide if he or she is willing to distribute software through any other system and a licensee cannot impose that choice.

This section is intended to make thoroughly clear what is believed to be a consequence of the rest of this License.

- 8. If the distribution and/or use of the Program is restricted in certain countries either by patents or by copyrighted interfaces, the original copyright holder who places the Program under this License may add an explicit geographical distribution limitation excluding those countries, so that distribution is permitted only in or among countries not thus excluded. In such case, this License incorporates the limitation as if written in the body of this License.
- 9. The Free Software Foundation may publish revised and/or new versions of the General Public License from time to time. Such new versions will be similar in spirit to the present version, but may differ in detail to address new problems or concerns.

Each version is given a distinguishing version number. If the Program specifies a version number of this License which applies to it and "any later version", you have the option of following the terms and conditions either of that version or of any later version published by the Free Software Foundation. If the Program does

not specify a version number of this License, you may choose any version ever published by the Free Software Foundation.

10. If you wish to incorporate parts of the Program into other free programs whose distribution conditions are different, write to the author to ask for permission. For software which is copyrighted by the Free Software Foundation, write to the Free Software Foundation; we sometimes make exceptions for this. Our decision will be guided by the two goals of preserving the free status of all derivatives of our free software and of promoting the sharing and reuse of software generally.

NO WARRANTY

- 11. BECAUSE THE PROGRAM IS LICENSED FREE OF CHARGE, THERE IS NO WARRANTY FOR THE PROGRAM, TO THE EXTENT PERMITTED BY APPLICABLE LAW. EXCEPT WHEN OTHERWISE STATED IN WRITING THE COPYRIGHT HOLDERS AND/OR OTHER PARTIES PROVIDE THE PROGRAM "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESSED OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE. THE ENTIRE RISK AS TO THE QUALITY AND PERFORMANCE OF THE PROGRAM IS WITH YOU. SHOULD THE PROGRAM PROVE DEFECTIVE, YOU ASSUME THE COST OF ALL NECESSARY SERVICING, REPAIR OR CORRECTION.
- 12. IN NO EVENT UNLESS REQUIRED BY APPLICABLE LAW OR AGREED TO IN WRITING WILL ANY COPYRIGHT HOLDER, OR ANY OTHER PARTY WHO MAY MODIFY AND/OR REDISTRIBUTE THE PROGRAM AS PERMITTED ABOVE, BE LIABLE TO YOU FOR DAMAGES, INCLUDING ANY GENERAL, SPECIAL, INCIDENTAL OR CONSEQUENTIAL DAMAGES ARISING OUT OF THE USE OR INABILITY TO USE THE PROGRAM (INCLUDING BUT NOT LIMITED TO LOSS OF DATA OR DATA BEING RENDERED INACCURATE OR LOSSES SUSTAINED BY YOU OR THIRD PARTIES OR A FAILURE OF THE PROGRAM TO OPERATE WITH ANY OTHER PROGRAMS), EVEN IF SUCH HOLDER OR OTHER PARTY HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

END OF TERMS AND CONDITIONS

How to Apply These Terms to Your New Programs

If you develop a new program, and you want it to be of the greatest possible use to the public, the best way to achieve this is to make it free software which everyone can redistribute and change under these terms.

To do so, attach the following notices to the program. It is safest to attach them to the start of each source file to most effectively convey the exclusion of warranty; and each file should have at least the "copyright" line and a pointer to where the full notice is found.

<one line to give the program's name and a brief idea of what it does.>

Copyright (C) <year> <name of author>

This program is free software; you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation; either version 2 of the License, or (at your option) any later version.

This program is distributed in the hope that it will be useful,but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details.

You should have received a copy of the GNU General Public License along with this program; if not, write to the Free Software Foundation, Inc., 51 Franklin St, Fifth Floor, Boston, MA 02110-1301

Also add information on how to contact you by electronic and paper mail.

If the program is interactive, make it output a short notice like this when it starts in an interactive mode:

Gnomovision version 69, Copyright (C) year name of author

Gnomovision comes with ABSOLUTELY NO WARRANTY; for details type 'show w'.

This is free software, and you are welcome to redistribute it under certain conditions; type `show c' for details.

The hypothetical commands `show w' and `show c' should show the appropriate parts of the General Public License. Of course, the commands you use may be called something other than `show w' and `show c'; they could even be mouse-clicks or menu items--whatever suits your program.

You should also get your employer (if you work as a programmer) or your school, if any, to sign a "copyright disclaimer" for the program, if necessary. Here is a sample; alter the names:

Yoyodyne, Inc., hereby disclaims all copyright interest in the program 'Gnomovision' (which makes passes at compilers) written by James Hacker.

<signature of Ty Coon>, 1 April 1989 Ty Coon, President of Vice

This General Public License does not permit incorporating your program into proprietary programs. If your program is a subroutine library, you may consider it more useful to permit linking proprietary applications with the library. If this is what you want to do, use the GNU Library General Public License instead of this License.

Software:Geos

Copyright notice:

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Refractions Research Inc.

Copyright (C) 2013 Sandro Santilli <strk@keybit.net>

- Copyright (C) 2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2009 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2005-2011 Refractions Research Inc.
- Copyright (C) 2009 Ragi Y. Burhum <ragi@burhum.com>
- Copyright (C) 2010 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2005 2006 Refractions Research Inc.
- Copyright (C) 2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2006-2011 Refractions Research Inc.
- Copyright (C) 2011 Sandro Santilli <strk@keybit.net
- Copyright (C) 2009-2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2016 Daniel Baston
- Copyright (C) 2008 Sean Gillies
- Copyright (C) 2009 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2006 Refractions Research Inc.
- Copyright (C) 2012 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2009 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2008-2010 Safe Software Inc.
- Copyright (C) 2006-2007 Refractions Research Inc.
- Copyright (C) 2005-2007 Refractions Research Inc.
- Copyright (C) 2007 Refractions Research Inc.
- Copyright (C) 2014 Mika Heiskanen <mika.heiskanen@fmi.fi>
- Copyright (C) 2009-2010 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2010 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2009 Mateusz Loskot
- Copyright (C) 2005-2009 Refractions Research Inc.
- Copyright (C) 2001-2009 Vivid Solutions Inc.
- Copyright (C) 2012 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2006 Wu Yongwei
- Copyright (C) 2012 Excensus LLC.
- Copyright (C) 1996-2015 Free Software Foundation, Inc.

Copyright (c) 1995 Olivier Devillers <Olivier.Devillers@sophia.inria.fr>

Copyright (C) 2007-2010 Safe Software Inc.

Copyright (C) 2010 Safe Software Inc.

Copyright (C) 2006 Refractions Research

Copyright 2004 Sean Gillies, sgillies@frii.com

Copyright (C) 2011 Mateusz Loskot <mateusz@loskot.net>

Copyright (C) 2015 Nyall Dawson < nyall dot dawson at gmail dot com>

Original code (2.0 and earlier)copyright (c) 2000-2006 Lee Thomason (www.grinninglizard.com)

Original code (2.0 and earlier)copyright (c) 2000-2002 Lee Thomason (www.grinninglizard.com)

License: LGPL V2.1

GNU LESSER GENERAL PUBLIC LICENSE

Version 2.1, February 1999

Copyright (C) 1991, 1999 Free Software Foundation, Inc. 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

Copyright (C) 2005-2011 Refractions Research Inc.

Copyright (C) 2009 Ragi Y. Burhum <ragi@burhum.com>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2005 2006 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006-2011 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net

Copyright (C) 2009-2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2016 Daniel Baston

Copyright (C) 2008 Sean Gillies

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Refractions Research Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2008-2010 Safe Software Inc.

Copyright (C) 2006-2007 Refractions Research Inc.

Copyright (C) 2005-2007 Refractions Research Inc.

Copyright (C) 2007 Refractions Research Inc.

Copyright (C) 2014 Mika Heiskanen <mika.heiskanen@fmi.fi>

Copyright (C) 2009-2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Mateusz Loskot

Copyright (C) 2005-2009 Refractions Research Inc.

Copyright (C) 2001-2009 Vivid Solutions Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Wu Yongwei

Copyright (C) 2012 Excensus LLC.

Copyright (C) 1996-2015 Free Software Foundation, Inc.

Copyright (c) 1995 Olivier Devillers <Olivier.Devillers@sophia.inria.fr>

Copyright (C) 2007-2010 Safe Software Inc.

Copyright (C) 2010 Safe Software Inc.

Copyright (C) 2006 Refractions Research

Copyright 2004 Sean Gillies, sgillies@frii.com

Copyright (C) 2011 Mateusz Loskot <mateusz@loskot.net>

Copyright (C) 2015 Nyall Dawson < nyall dot dawson at gmail dot com>

Original code (2.0 and earlier)copyright (c) 2000-2006 Lee Thomason (www.grinninglizard.com)

Original code (2.0 and earlier)copyright (c) 2000-2002 Lee Thomason (www.grinninglizard.com)

License: LGPL V2.1

GNU LESSER GENERAL PUBLIC LICENSE

Version 2.1, February 1999

Copyright (C) 1991, 1999 Free Software Foundation, Inc. 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

- Copyright (C) 2005-2011 Refractions Research Inc.
- Copyright (C) 2009 Ragi Y. Burhum <ragi@burhum.com>
- Copyright (C) 2010 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2005 2006 Refractions Research Inc.
- Copyright (C) 2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2006-2011 Refractions Research Inc.
- Copyright (C) 2011 Sandro Santilli <strk@keybit.net
- Copyright (C) 2009-2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2016 Daniel Baston
- Copyright (C) 2008 Sean Gillies
- Copyright (C) 2009 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2006 Refractions Research Inc.
- Copyright (C) 2012 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2009 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2008-2010 Safe Software Inc.
- Copyright (C) 2006-2007 Refractions Research Inc.
- Copyright (C) 2005-2007 Refractions Research Inc.
- Copyright (C) 2007 Refractions Research Inc.
- Copyright (C) 2014 Mika Heiskanen <mika.heiskanen@fmi.fi>
- Copyright (C) 2009-2010 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2010 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2009 Mateusz Loskot
- Copyright (C) 2005-2009 Refractions Research Inc.
- Copyright (C) 2001-2009 Vivid Solutions Inc.
- Copyright (C) 2012 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2006 Wu Yongwei
- Copyright (C) 2012 Excensus LLC.
- Copyright (C) 1996-2015 Free Software Foundation, Inc.
- Copyright (c) 1995 Olivier Devillers <Olivier.Devillers@sophia.inria.fr>
- Copyright (C) 2007-2010 Safe Software Inc.
- Copyright (C) 2010 Safe Software Inc.

Copyright (C) 2006 Refractions Research

Copyright 2004 Sean Gillies, sgillies@frii.com

Copyright (C) 2011 Mateusz Loskot <mateusz@loskot.net>

Copyright (C) 2015 Nyall Dawson < nyall dot dawson at gmail dot com>

Original code (2.0 and earlier)copyright (c) 2000-2006 Lee Thomason (www.grinninglizard.com)

Original code (2.0 and earlier)copyright (c) 2000-2002 Lee Thomason (www.grinninglizard.com)

License: LGPL V2.1

GNU LESSER GENERAL PUBLIC LICENSE

Version 2.1, February 1999

Copyright (C) 1991, 1999 Free Software Foundation, Inc. 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

Copyright (C) 2005-2011 Refractions Research Inc.

Copyright (C) 2009 Ragi Y. Burhum <ragi@burhum.com>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2005 2006 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006-2011 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net

Copyright (C) 2009-2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2016 Daniel Baston

Copyright (C) 2008 Sean Gillies

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Refractions Research Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2008-2010 Safe Software Inc.

Copyright (C) 2006-2007 Refractions Research Inc.

Copyright (C) 2005-2007 Refractions Research Inc.

Copyright (C) 2007 Refractions Research Inc.

Copyright (C) 2014 Mika Heiskanen <mika.heiskanen@fmi.fi>

Copyright (C) 2009-2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Mateusz Loskot

Copyright (C) 2005-2009 Refractions Research Inc.

Copyright (C) 2001-2009 Vivid Solutions Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Wu Yongwei

Copyright (C) 2012 Excensus LLC.

Copyright (C) 1996-2015 Free Software Foundation, Inc.

Copyright (c) 1995 Olivier Devillers < Olivier. Devillers@sophia.inria.fr>

Copyright (C) 2007-2010 Safe Software Inc.

Copyright (C) 2010 Safe Software Inc.

Copyright (C) 2006 Refractions Research

Copyright 2004 Sean Gillies, sgillies@frii.com

Copyright (C) 2011 Mateusz Loskot <mateusz@loskot.net>

Copyright (C) 2015 Nyall Dawson < nyall dot dawson at gmail dot com>

Original code (2.0 and earlier)copyright (c) 2000-2006 Lee Thomason (www.grinninglizard.com)

Original code (2.0 and earlier)copyright (c) 2000-2002 Lee Thomason (www.grinninglizard.com)

License: LGPL V2.1

GNU LESSER GENERAL PUBLIC LICENSE

Version 2.1, February 1999

Copyright (C) 1991, 1999 Free Software Foundation, Inc. 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

Copyright (C) 2005-2011 Refractions Research Inc.

Copyright (C) 2009 Ragi Y. Burhum <ragi@burhum.com>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

- Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2005 2006 Refractions Research Inc.
- Copyright (C) 2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2006-2011 Refractions Research Inc.
- Copyright (C) 2011 Sandro Santilli <strk@keybit.net
- Copyright (C) 2009-2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2016 Daniel Baston
- Copyright (C) 2008 Sean Gillies
- Copyright (C) 2009 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2006 Refractions Research Inc.
- Copyright (C) 2012 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2009 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2008-2010 Safe Software Inc.
- Copyright (C) 2006-2007 Refractions Research Inc.
- Copyright (C) 2005-2007 Refractions Research Inc.
- Copyright (C) 2007 Refractions Research Inc.
- Copyright (C) 2014 Mika Heiskanen <mika.heiskanen@fmi.fi>
- Copyright (C) 2009-2010 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2010 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2009 Mateusz Loskot
- Copyright (C) 2005-2009 Refractions Research Inc.
- Copyright (C) 2001-2009 Vivid Solutions Inc.
- Copyright (C) 2012 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2006 Wu Yongwei
- Copyright (C) 2012 Excensus LLC.
- Copyright (C) 1996-2015 Free Software Foundation, Inc.
- Copyright (c) 1995 Olivier Devillers <Olivier.Devillers@sophia.inria.fr>
- Copyright (C) 2007-2010 Safe Software Inc.
- Copyright (C) 2010 Safe Software Inc.
- Copyright (C) 2006 Refractions Research
- Copyright 2004 Sean Gillies, sgillies@frii.com
- Copyright (C) 2011 Mateusz Loskot <mateusz@loskot.net>

Copyright (C) 2015 Nyall Dawson < nyall dot dawson at gmail dot com>

Original code (2.0 and earlier)copyright (c) 2000-2006 Lee Thomason (www.grinninglizard.com)

Original code (2.0 and earlier)copyright (c) 2000-2002 Lee Thomason (www.grinninglizard.com)

License: LGPL V2.1

GNU LESSER GENERAL PUBLIC LICENSE

Version 2.1, February 1999

Copyright (C) 1991, 1999 Free Software Foundation, Inc. 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

Copyright (C) 2005-2011 Refractions Research Inc.

Copyright (C) 2009 Ragi Y. Burhum <ragi@burhum.com>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2005 2006 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006-2011 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net

Copyright (C) 2009-2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2016 Daniel Baston

Copyright (C) 2008 Sean Gillies

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Refractions Research Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2008-2010 Safe Software Inc.

Copyright (C) 2006-2007 Refractions Research Inc.

Copyright (C) 2005-2007 Refractions Research Inc.

Copyright (C) 2007 Refractions Research Inc.

Copyright (C) 2014 Mika Heiskanen <mika.heiskanen@fmi.fi>

Copyright (C) 2009-2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Mateusz Loskot

Copyright (C) 2005-2009 Refractions Research Inc.

Copyright (C) 2001-2009 Vivid Solutions Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Wu Yongwei

Copyright (C) 2012 Excensus LLC.

Copyright (C) 1996-2015 Free Software Foundation, Inc.

Copyright (c) 1995 Olivier Devillers < Olivier. Devillers@sophia.inria.fr>

Copyright (C) 2007-2010 Safe Software Inc.

Copyright (C) 2010 Safe Software Inc.

Copyright (C) 2006 Refractions Research

Copyright 2004 Sean Gillies, sgillies@frii.com

Copyright (C) 2011 Mateusz Loskot <mateusz@loskot.net>

Copyright (C) 2015 Nyall Dawson < nyall dot dawson at gmail dot com>

Original code (2.0 and earlier)copyright (c) 2000-2006 Lee Thomason (www.grinninglizard.com)

Original code (2.0 and earlier)copyright (c) 2000-2002 Lee Thomason (www.grinninglizard.com)

License: LGPL V2.1

GNU LESSER GENERAL PUBLIC LICENSE

Version 2.1, February 1999

Copyright (C) 1991, 1999 Free Software Foundation, Inc. 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

Copyright (C) 2005-2011 Refractions Research Inc.

Copyright (C) 2009 Ragi Y. Burhum <ragi@burhum.com>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2005 2006 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006-2011 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net

Copyright (C) 2009-2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2016 Daniel Baston

Copyright (C) 2008 Sean Gillies

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Refractions Research Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2008-2010 Safe Software Inc.

Copyright (C) 2006-2007 Refractions Research Inc.

Copyright (C) 2005-2007 Refractions Research Inc.

Copyright (C) 2007 Refractions Research Inc.

Copyright (C) 2014 Mika Heiskanen <mika.heiskanen@fmi.fi>

Copyright (C) 2009-2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Mateusz Loskot

Copyright (C) 2005-2009 Refractions Research Inc.

Copyright (C) 2001-2009 Vivid Solutions Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Wu Yongwei

Copyright (C) 2012 Excensus LLC.

Copyright (C) 1996-2015 Free Software Foundation, Inc.

Copyright (c) 1995 Olivier Devillers < Olivier. Devillers@sophia.inria.fr>

Copyright (C) 2007-2010 Safe Software Inc.

Copyright (C) 2010 Safe Software Inc.

Copyright (C) 2006 Refractions Research

Copyright 2004 Sean Gillies, sgillies@frii.com

Copyright (C) 2011 Mateusz Loskot <mateusz@loskot.net>

Copyright (C) 2015 Nyall Dawson < nyall dot dawson at gmail dot com>

Original code (2.0 and earlier)copyright (c) 2000-2006 Lee Thomason (www.grinninglizard.com)

Original code (2.0 and earlier)copyright (c) 2000-2002 Lee Thomason (www.grinninglizard.com)

License: LGPL V2.1

GNU LESSER GENERAL PUBLIC LICENSE

Version 2.1, February 1999

Copyright (C) 1991, 1999 Free Software Foundation, Inc. 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

Copyright (C) 2005-2011 Refractions Research Inc.

Copyright (C) 2009 Ragi Y. Burhum <ragi@burhum.com>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2005 2006 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006-2011 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net

Copyright (C) 2009-2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2016 Daniel Baston

Copyright (C) 2008 Sean Gillies

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Refractions Research Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2008-2010 Safe Software Inc.

Copyright (C) 2006-2007 Refractions Research Inc.

Copyright (C) 2005-2007 Refractions Research Inc.

Copyright (C) 2007 Refractions Research Inc.

Copyright (C) 2014 Mika Heiskanen <mika.heiskanen@fmi.fi>

Copyright (C) 2009-2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Mateusz Loskot

Copyright (C) 2005-2009 Refractions Research Inc.

Copyright (C) 2001-2009 Vivid Solutions Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Wu Yongwei

Copyright (C) 2012 Excensus LLC.

Copyright (C) 1996-2015 Free Software Foundation, Inc.

Copyright (c) 1995 Olivier Devillers <Olivier.Devillers@sophia.inria.fr>

Copyright (C) 2007-2010 Safe Software Inc.

Copyright (C) 2010 Safe Software Inc.

Copyright (C) 2006 Refractions Research

Copyright 2004 Sean Gillies, sgillies@frii.com

Copyright (C) 2011 Mateusz Loskot <mateusz@loskot.net>

Copyright (C) 2015 Nyall Dawson < nyall dot dawson at gmail dot com>

Original code (2.0 and earlier)copyright (c) 2000-2006 Lee Thomason (www.grinninglizard.com)

Original code (2.0 and earlier)copyright (c) 2000-2002 Lee Thomason (www.grinninglizard.com)

License: LGPL V2.1

GNU LESSER GENERAL PUBLIC LICENSE

Version 2.1, February 1999

Copyright (C) 1991, 1999 Free Software Foundation, Inc. 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

Copyright (C) 2005-2011 Refractions Research Inc.

Copyright (C) 2009 Ragi Y. Burhum <ragi@burhum.com>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2005 2006 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006-2011 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net

Copyright (C) 2009-2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2016 Daniel Baston

Copyright (C) 2008 Sean Gillies

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Refractions Research Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2008-2010 Safe Software Inc.

Copyright (C) 2006-2007 Refractions Research Inc.

Copyright (C) 2005-2007 Refractions Research Inc.

Copyright (C) 2007 Refractions Research Inc.

Copyright (C) 2014 Mika Heiskanen <mika.heiskanen@fmi.fi>

Copyright (C) 2009-2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Mateusz Loskot

Copyright (C) 2005-2009 Refractions Research Inc.

Copyright (C) 2001-2009 Vivid Solutions Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Wu Yongwei

Copyright (C) 2012 Excensus LLC.

Copyright (C) 1996-2015 Free Software Foundation, Inc.

Copyright (c) 1995 Olivier Devillers <Olivier.Devillers@sophia.inria.fr>

Copyright (C) 2007-2010 Safe Software Inc.

Copyright (C) 2010 Safe Software Inc.

Copyright (C) 2006 Refractions Research

Copyright 2004 Sean Gillies, sgillies@frii.com

Copyright (C) 2011 Mateusz Loskot <mateusz@loskot.net>

Copyright (C) 2015 Nyall Dawson < nyall dot dawson at gmail dot com>

Original code (2.0 and earlier)copyright (c) 2000-2006 Lee Thomason (www.grinninglizard.com)

Original code (2.0 and earlier)copyright (c) 2000-2002 Lee Thomason (www.grinninglizard.com)

License: LGPL V2.1

GNU LESSER GENERAL PUBLIC LICENSE

Version 2.1, February 1999

Copyright (C) 1991, 1999 Free Software Foundation, Inc. 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

Copyright (C) 2005-2011 Refractions Research Inc.

Copyright (C) 2009 Ragi Y. Burhum <ragi@burhum.com>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2005 2006 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006-2011 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net

Copyright (C) 2009-2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2016 Daniel Baston

Copyright (C) 2008 Sean Gillies

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Refractions Research Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2008-2010 Safe Software Inc.

Copyright (C) 2006-2007 Refractions Research Inc.

Copyright (C) 2005-2007 Refractions Research Inc.

Copyright (C) 2007 Refractions Research Inc.

Copyright (C) 2014 Mika Heiskanen <mika.heiskanen@fmi.fi>

Copyright (C) 2009-2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Mateusz Loskot

Copyright (C) 2005-2009 Refractions Research Inc.

Copyright (C) 2001-2009 Vivid Solutions Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Wu Yongwei

Copyright (C) 2012 Excensus LLC.

Copyright (C) 1996-2015 Free Software Foundation, Inc.

Copyright (c) 1995 Olivier Devillers <Olivier.Devillers@sophia.inria.fr>

Copyright (C) 2007-2010 Safe Software Inc.

Copyright (C) 2010 Safe Software Inc.

Copyright (C) 2006 Refractions Research

Copyright 2004 Sean Gillies, sgillies@frii.com

Copyright (C) 2011 Mateusz Loskot <mateusz@loskot.net>

Copyright (C) 2015 Nyall Dawson < nyall dot dawson at gmail dot com>

Original code (2.0 and earlier)copyright (c) 2000-2006 Lee Thomason (www.grinninglizard.com)

Original code (2.0 and earlier)copyright (c) 2000-2002 Lee Thomason (www.grinninglizard.com)

License: LGPL V2.1

GNU LESSER GENERAL PUBLIC LICENSE

Version 2.1, February 1999

Copyright (C) 1991, 1999 Free Software Foundation, Inc. 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

Copyright (C) 2005-2011 Refractions Research Inc.

Copyright (C) 2009 Ragi Y. Burhum <ragi@burhum.com>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2005 2006 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006-2011 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net

Copyright (C) 2009-2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2016 Daniel Baston

Copyright (C) 2008 Sean Gillies

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Refractions Research Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2008-2010 Safe Software Inc.

Copyright (C) 2006-2007 Refractions Research Inc.

Copyright (C) 2005-2007 Refractions Research Inc.

Copyright (C) 2007 Refractions Research Inc.

Copyright (C) 2014 Mika Heiskanen <mika.heiskanen@fmi.fi>

Copyright (C) 2009-2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Mateusz Loskot

Copyright (C) 2005-2009 Refractions Research Inc.

Copyright (C) 2001-2009 Vivid Solutions Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Wu Yongwei

Copyright (C) 2012 Excensus LLC.

Copyright (C) 1996-2015 Free Software Foundation, Inc.

Copyright (c) 1995 Olivier Devillers <Olivier.Devillers@sophia.inria.fr>

Copyright (C) 2007-2010 Safe Software Inc.

Copyright (C) 2010 Safe Software Inc.

Copyright (C) 2006 Refractions Research

Copyright 2004 Sean Gillies, sgillies@frii.com

Copyright (C) 2011 Mateusz Loskot <mateusz@loskot.net>

Copyright (C) 2015 Nyall Dawson < nyall dot dawson at gmail dot com>

Original code (2.0 and earlier)copyright (c) 2000-2006 Lee Thomason (www.grinninglizard.com)

Original code (2.0 and earlier)copyright (c) 2000-2002 Lee Thomason (www.grinninglizard.com)

License: LGPL V2.1

GNU LESSER GENERAL PUBLIC LICENSE

Version 2.1, February 1999

Copyright (C) 1991, 1999 Free Software Foundation, Inc. 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

Copyright (C) 2005-2011 Refractions Research Inc.

Copyright (C) 2009 Ragi Y. Burhum <ragi@burhum.com>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2005 2006 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006-2011 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net

Copyright (C) 2009-2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2016 Daniel Baston

Copyright (C) 2008 Sean Gillies

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Refractions Research Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2008-2010 Safe Software Inc.

Copyright (C) 2006-2007 Refractions Research Inc.

Copyright (C) 2005-2007 Refractions Research Inc.

Copyright (C) 2007 Refractions Research Inc.

Copyright (C) 2014 Mika Heiskanen <mika.heiskanen@fmi.fi>

Copyright (C) 2009-2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Mateusz Loskot

Copyright (C) 2005-2009 Refractions Research Inc.

Copyright (C) 2001-2009 Vivid Solutions Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Wu Yongwei

Copyright (C) 2012 Excensus LLC.

Copyright (C) 1996-2015 Free Software Foundation, Inc.

Copyright (c) 1995 Olivier Devillers <Olivier.Devillers@sophia.inria.fr>

Copyright (C) 2007-2010 Safe Software Inc.

Copyright (C) 2010 Safe Software Inc.

Copyright (C) 2006 Refractions Research

Copyright 2004 Sean Gillies, sgillies@frii.com

Copyright (C) 2011 Mateusz Loskot <mateusz@loskot.net>

Copyright (C) 2015 Nyall Dawson < nyall dot dawson at gmail dot com>

Original code (2.0 and earlier)copyright (c) 2000-2006 Lee Thomason (www.grinninglizard.com)

Original code (2.0 and earlier)copyright (c) 2000-2002 Lee Thomason (www.grinninglizard.com)

License: LGPL V2.1

GNU LESSER GENERAL PUBLIC LICENSE

Version 2.1, February 1999

Copyright (C) 1991, 1999 Free Software Foundation, Inc. 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

Copyright (C) 2005-2011 Refractions Research Inc.

Copyright (C) 2009 Ragi Y. Burhum <ragi@burhum.com>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2005 2006 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006-2011 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net

Copyright (C) 2009-2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2016 Daniel Baston

Copyright (C) 2008 Sean Gillies

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Refractions Research Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2008-2010 Safe Software Inc.

Copyright (C) 2006-2007 Refractions Research Inc.

Copyright (C) 2005-2007 Refractions Research Inc.

Copyright (C) 2007 Refractions Research Inc.

Copyright (C) 2014 Mika Heiskanen <mika.heiskanen@fmi.fi>

Copyright (C) 2009-2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Mateusz Loskot

Copyright (C) 2005-2009 Refractions Research Inc.

Copyright (C) 2001-2009 Vivid Solutions Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Wu Yongwei

Copyright (C) 2012 Excensus LLC.

Copyright (C) 1996-2015 Free Software Foundation, Inc.

Copyright (c) 1995 Olivier Devillers <Olivier.Devillers@sophia.inria.fr>

Copyright (C) 2007-2010 Safe Software Inc.

Copyright (C) 2010 Safe Software Inc.

Copyright (C) 2006 Refractions Research

Copyright 2004 Sean Gillies, sgillies@frii.com

Copyright (C) 2011 Mateusz Loskot <mateusz@loskot.net>

Copyright (C) 2015 Nyall Dawson < nyall dot dawson at gmail dot com>

Original code (2.0 and earlier)copyright (c) 2000-2006 Lee Thomason (www.grinninglizard.com)

Original code (2.0 and earlier)copyright (c) 2000-2002 Lee Thomason (www.grinninglizard.com)

License: LGPL V2.1

GNU LESSER GENERAL PUBLIC LICENSE

Version 2.1, February 1999

Copyright (C) 1991, 1999 Free Software Foundation, Inc. 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

- Copyright (C) 2005-2011 Refractions Research Inc.
- Copyright (C) 2009 Ragi Y. Burhum <ragi@burhum.com>
- Copyright (C) 2010 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2005 2006 Refractions Research Inc.
- Copyright (C) 2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2006-2011 Refractions Research Inc.
- Copyright (C) 2011 Sandro Santilli <strk@keybit.net
- Copyright (C) 2009-2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2016 Daniel Baston
- Copyright (C) 2008 Sean Gillies
- Copyright (C) 2009 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2006 Refractions Research Inc.
- Copyright (C) 2012 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2009 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2008-2010 Safe Software Inc.
- Copyright (C) 2006-2007 Refractions Research Inc.
- Copyright (C) 2005-2007 Refractions Research Inc.
- Copyright (C) 2007 Refractions Research Inc.
- Copyright (C) 2014 Mika Heiskanen <mika.heiskanen@fmi.fi>
- Copyright (C) 2009-2010 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2010 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2009 Mateusz Loskot
- Copyright (C) 2005-2009 Refractions Research Inc.
- Copyright (C) 2001-2009 Vivid Solutions Inc.
- Copyright (C) 2012 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2006 Wu Yongwei
- Copyright (C) 2012 Excensus LLC.
- Copyright (C) 1996-2015 Free Software Foundation, Inc.
- Copyright (c) 1995 Olivier Devillers <Olivier.Devillers@sophia.inria.fr>
- Copyright (C) 2007-2010 Safe Software Inc.
- Copyright (C) 2010 Safe Software Inc.

Copyright (C) 2006 Refractions Research

Copyright 2004 Sean Gillies, sgillies@frii.com

Copyright (C) 2011 Mateusz Loskot <mateusz@loskot.net>

Copyright (C) 2015 Nyall Dawson < nyall dot dawson at gmail dot com>

Original code (2.0 and earlier)copyright (c) 2000-2006 Lee Thomason (www.grinninglizard.com)

Original code (2.0 and earlier)copyright (c) 2000-2002 Lee Thomason (www.grinninglizard.com)

License: LGPL V2.1

GNU LESSER GENERAL PUBLIC LICENSE

Version 2.1, February 1999

Copyright (C) 1991, 1999 Free Software Foundation, Inc. 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

Copyright (C) 2005-2011 Refractions Research Inc.

Copyright (C) 2009 Ragi Y. Burhum <ragi@burhum.com>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2005 2006 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006-2011 Refractions Research Inc.

Copyright (C) 2011 Sandro Santilli <strk@keybit.net

Copyright (C) 2009-2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2016 Daniel Baston

Copyright (C) 2008 Sean Gillies

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Refractions Research Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Sandro Santilli <strk@keybit.net>

Copyright (C) 2008-2010 Safe Software Inc.

Copyright (C) 2006-2007 Refractions Research Inc.

Copyright (C) 2005-2007 Refractions Research Inc.

Copyright (C) 2007 Refractions Research Inc.

Copyright (C) 2014 Mika Heiskanen <mika.heiskanen@fmi.fi>

Copyright (C) 2009-2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

Copyright (C) 2009 Mateusz Loskot

Copyright (C) 2005-2009 Refractions Research Inc.

Copyright (C) 2001-2009 Vivid Solutions Inc.

Copyright (C) 2012 Sandro Santilli <strk@keybit.net>

Copyright (C) 2006 Wu Yongwei

Copyright (C) 2012 Excensus LLC.

Copyright (C) 1996-2015 Free Software Foundation, Inc.

Copyright (c) 1995 Olivier Devillers <Olivier.Devillers@sophia.inria.fr>

Copyright (C) 2007-2010 Safe Software Inc.

Copyright (C) 2010 Safe Software Inc.

Copyright (C) 2006 Refractions Research

Copyright 2004 Sean Gillies, sgillies@frii.com

Copyright (C) 2011 Mateusz Loskot <mateusz@loskot.net>

Copyright (C) 2015 Nyall Dawson < nyall dot dawson at gmail dot com>

Original code (2.0 and earlier)copyright (c) 2000-2006 Lee Thomason (www.grinninglizard.com)

Original code (2.0 and earlier)copyright (c) 2000-2002 Lee Thomason (www.grinninglizard.com)

License: LGPL V2.1

GNU LESSER GENERAL PUBLIC LICENSE

Version 2.1, February 1999

Copyright (C) 1991, 1999 Free Software Foundation, Inc. 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

Copyright (C) 2005-2011 Refractions Research Inc.

Copyright (C) 2009 Ragi Y. Burhum <ragi@burhum.com>

Copyright (C) 2010 Sandro Santilli <strk@keybit.net>

- Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2005 2006 Refractions Research Inc.
- Copyright (C) 2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2006-2011 Refractions Research Inc.
- Copyright (C) 2011 Sandro Santilli <strk@keybit.net
- Copyright (C) 2009-2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2016 Daniel Baston
- Copyright (C) 2008 Sean Gillies
- Copyright (C) 2009 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2006 Refractions Research Inc.
- Copyright (C) 2012 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2009 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2008-2010 Safe Software Inc.
- Copyright (C) 2006-2007 Refractions Research Inc.
- Copyright (C) 2005-2007 Refractions Research Inc.
- Copyright (C) 2007 Refractions Research Inc.
- Copyright (C) 2014 Mika Heiskanen <mika.heiskanen@fmi.fi>
- Copyright (C) 2009-2010 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2009 2011 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2010 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2009 Mateusz Loskot
- Copyright (C) 2005-2009 Refractions Research Inc.
- Copyright (C) 2001-2009 Vivid Solutions Inc.
- Copyright (C) 2012 Sandro Santilli <strk@keybit.net>
- Copyright (C) 2006 Wu Yongwei
- Copyright (C) 2012 Excensus LLC.
- Copyright (C) 1996-2015 Free Software Foundation, Inc.
- Copyright (c) 1995 Olivier Devillers <Olivier.Devillers@sophia.inria.fr>
- Copyright (C) 2007-2010 Safe Software Inc.
- Copyright (C) 2010 Safe Software Inc.
- Copyright (C) 2006 Refractions Research
- Copyright 2004 Sean Gillies, sgillies@frii.com
- Copyright (C) 2011 Mateusz Loskot <mateusz@loskot.net>

Copyright (C) 2015 Nyall Dawson < nyall dot dawson at gmail dot com>

Original code (2.0 and earlier)copyright (c) 2000-2006 Lee Thomason (www.grinninglizard.com)

Original code (2.0 and earlier)copyright (c) 2000-2002 Lee Thomason (www.grinninglizard.com)

License: LGPL V2.1

GNU LESSER GENERAL PUBLIC LICENSE

Version 2.1, February 1999

Copyright (C) 1991, 1999 Free Software Foundation, Inc. 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

[This is the first released version of the Lesser GPL. It also counts as the successor of the GNU Library Public License, version 2, hence the version number 2.1.]

Preamble

The licenses for most software are designed to take away your freedom to share and change it. By contrast, the GNU General Public

Licenses are intended to guarantee your freedom to share and change free software--to make sure the software is free for all its users.

This license, the Lesser General Public License, applies to some specially designated software packages--typically libraries--of the Free Software Foundation and other authors who decide to use it. You can use it too, but we suggest you first think carefully about whether this license or the ordinary General Public License is the better strategy to use in any particular case, based on the explanations below.

When we speak of free software, we are referring to freedom of use, not price. Our General Public Licenses are designed to make sure that you have the freedom to distribute copies of free software (and charge for this service if you wish); that you receive source code or can get it if you want it; that you can change the software and use pieces of it in new free programs; and that you are informed that you can do these things.

To protect your rights, we need to make restrictions that forbid distributors to deny you these rights or to ask you to surrender these rights. These restrictions translate to certain responsibilities for you if you distribute copies of the library or if you modify it.

For example, if you distribute copies of the library, whether gratis or for a fee, you must give the recipients all the rights that we gave you. You must make sure that they, too, receive or can get the source code. If you link other code with the library,

you must provide complete object files to the recipients, so that they can relink them with the library after making changes to the library and recompiling it. And you must show them these terms so they know their rights.

We protect your rights with a two-step method: (1) we copyright the library, and (2) we offer you this license, which gives you legal permission to copy, distribute and/or modify the library.

To protect each distributor, we want to make it very clear that there is no warranty for the free library. Also, if the library is modified by someone else and passed on, the recipients should know that what they have is not the original version, so that the original author's reputation will not be affected by problems that might be introduced by others.

Finally, software patents pose a constant threat to the existence of any free program. We wish to make sure that a company cannot effectively restrict the users of a free program by obtaining a restrictive license from a patent holder. Therefore, we insist that any patent license obtained for a version of the library must be consistent with the full freedom of use specified in this license.

Most GNU software, including some libraries, is covered by the ordinary GNU General Public License. This license, the GNU Lesser General Public License, applies to certain designated libraries, and

is quite different from the ordinary General Public License. We use this license for certain libraries in order to permit linking those libraries into non-free programs.

When a program is linked with a library, whether statically or using a shared library, the combination of the two is legally speaking a combined work, a derivative of the original library. The ordinary General Public License therefore permits such linking only if the entire combination fits its criteria of freedom. The Lesser General Public License permits more lax criteria for linking other code with the library.

We call this license the "Lesser" General Public License because it does Less to protect the user's freedom than the ordinary General Public License. It also provides other free software developers Less of an advantage over competing non-free programs. These disadvantages are the reason we use the ordinary General Public License for many libraries. However, the Lesser license provides advantages in certain special circumstances.

For example, on rare occasions, there may be a special need to encourage the widest possible use of a certain library, so that it becomes a de-facto standard. To achieve this, non-free programs must be allowed to use the library. A more frequent case is that a free library does the same job as widely used non-free libraries. In this case, there is little to gain by limiting the free library to free software only, so we use the Lesser General Public License.

In other cases, permission to use a particular library in non-free programs enables a greater number of people to use a large body of free software. For example, permission to use the GNU C Library in

non-free programs enables many more people to use the whole GNU operating system, as well as its variant, the GNU/Linux operating system.

Although the Lesser General Public License is Less protective of the users' freedom, it does ensure that the user of a program that is linked with the Library has the freedom and the wherewithal to run that program using a modified version of the Library.

The precise terms and conditions for copying, distribution and modification follow. Pay close attention to the difference between a "work based on the library" and a "work that uses the library". The

former contains code derived from the library, whereas the latter must be combined with the library in order to run.

GNU LESSER GENERAL PUBLIC LICENSE

TERMS AND CONDITIONS FOR COPYING, DISTRIBUTION AND MODIFICATION

0. This License Agreement applies to any software library or other program which contains a notice placed by the copyright holder or other authorized party saying it may be distributed under the terms of this Lesser General Public License (also called "this License"). Each licensee is addressed as "you".

A "library" means a collection of software functions and/or data prepared so as to be conveniently linked with application programs (which use some of those functions and data) to form executables.

The "Library", below, refers to any such software library or work which has been distributed under these terms. A "work based on the Library" means either the Library or any derivative work under

copyright law: that is to say, a work containing the Library or a portion of it, either verbatim or with modifications and/or translated straightforwardly into another language. (Hereinafter, translation is included without limitation in the term "modification".)

"Source code" for a work means the preferred form of the work for making modifications to it. For a library, complete source code means all the source code for all modules it contains, plus any associated interface definition files, plus the scripts used to control compilation and installation of the library.

Activities other than copying, distribution and modification are not covered by this License; they are outside its scope. The act of running a program using the Library is not restricted, and output from such a program is covered only if its contents constitute a work based on the Library (independent of the use of the Library in a tool for writing it). Whether that is true depends on what the Library does and what the program that uses the Library does.

1. You may copy and distribute verbatim copies of the Library's complete source code as you receive it, in any medium, provided that you conspicuously and appropriately publish on each copy an

appropriate copyright notice and disclaimer of warranty; keep intact all the notices that refer to this License and to the absence of any warranty; and distribute a copy of this License along with the

Library.

You may charge a fee for the physical act of transferring a copy, and you may at your option offer warranty protection in exchange for a fee.

2. You may modify your copy or copies of the Library or any portion of it, thus forming a work based on the Library, and copy and distribute such modifications or work under the terms of Section 1

above, provided that you also meet all of these conditions:

- a) The modified work must itself be a software library.
- b) You must cause the files modified to carry prominent notices stating that you changed the files and the date of any change.
- c) You must cause the whole of the work to be licensed at no charge to all third parties under the terms of this License.
- d) If a facility in the modified Library refers to a function or a table of data to be supplied by an application program that uses the facility, other than as an argument passed when the facility is invoked, then you must make a good faith effort to ensure that, in the event an application does not supply such function or table, the facility still operates, and performs whatever part of

its purpose remains meaningful.

(For example, a function in a library to compute square roots has a purpose that is entirely well-defined independent of the application. Therefore, Subsection 2d requires that any application-supplied function or table used by this function must be optional: if the application does not supply it, the square root function must still compute square roots.)

These requirements apply to the modified work as a whole. If identifiable sections of that work are not derived from the Library, and can be reasonably considered independent and separate works in

themselves, then this License, and its terms, do not apply to those sections when you distribute them as separate works. But when you distribute the same sections as part of a whole which is a work based on the Library, the distribution of the whole must be on the terms of this License, whose permissions for other licensees extend to the entire whole, and thus to each and every part regardless of who wrote it.

Thus, it is not the intent of this section to claim rights or contest your rights to work written entirely by you; rather, the intent is to exercise the right to control the distribution of derivative or

collective works based on the Library.

In addition, mere aggregation of another work not based on the Library with the Library (or with a work based on the Library) on a volume of a storage or distribution medium does not bring the other work under the scope of this License.

3. You may opt to apply the terms of the ordinary GNU General Public License instead of this License to a given copy of the Library. To do this, you must alter all the notices that refer to this License, so that they refer to the ordinary GNU General Public License, version 2, instead of to this License. (If a newer version than version 2 of the ordinary GNU General Public License has appeared, then you can specify that version instead if you wish.) Do not make any other change in these notices.

Once this change is made in a given copy, it is irreversible for that copy, so the ordinary GNU General Public License applies to all subsequent copies and derivative works made from that copy.

This option is useful when you wish to copy part of the code of the Library into a program that is not a library.

4. You may copy and distribute the Library (or a portion or derivative of it, under Section 2) in object code or executable form under the terms of Sections 1 and 2 above provided that you accompany

it with the complete corresponding machine-readable source code, which must be distributed under the terms of Sections 1 and 2 above on a medium customarily used for software interchange.

If distribution of object code is made by offering access to copy from a designated place, then offering equivalent access to copy the source code from the same place satisfies the requirement to

distribute the source code, even though third parties are not compelled to copy the source along with the object code.

5. A program that contains no derivative of any portion of the Library, but is designed to work with the Library by being compiled or linked with it, is called a "work that uses the Library". Such a

work, in isolation, is not a derivative work of the Library, and therefore falls outside the scope of this License.

However, linking a "work that uses the Library" with the Library creates an executable that is a derivative of the Library (because it contains portions of the Library), rather than a "work that uses the library". The executable is therefore covered by this License.

Section 6 states terms for distribution of such executables.

When a "work that uses the Library" uses material from a header file that is part of the Library, the object code for the work may be a derivative work of the Library even though the source code is not. Whether this is true is especially

significant if the work can be linked without the Library, or if the work is itself a library. The threshold for this to be true is not precisely defined by law.

If such an object file uses only numerical parameters, data structure layouts and accessors, and small macros and small inline functions (ten lines or less in length), then the use of the object

file is unrestricted, regardless of whether it is legally a derivative work. (Executables containing this object code plus portions of the Library will still fall under Section 6.)

Otherwise, if the work is a derivative of the Library, you may distribute the object code for the work under the terms of Section 6. Any executables containing that work also fall under Section 6,

whether or not they are linked directly with the Library itself.

6. As an exception to the Sections above, you may also combine or link a "work that uses the Library" with the Library to produce a work containing portions of the Library, and distribute that work

under terms of your choice, provided that the terms permit modification of the work for the customer's own use and reverse engineering for debugging such modifications.

You must give prominent notice with each copy of the work that the Library is used in it and that the Library and its use are covered by this License. You must supply a copy of this License. If the work during execution displays copyright notices, you must include the copyright notice for the Library among them, as well as a reference directing the user to the copy of this License. Also, you must do one of these things:

- a) Accompany the work with the complete corresponding machine-readable source code for the Library including whatever changes were used in the work (which must be distributed under Sections 1 and 2 above); and, if the work is an executable linked with the Library, with the complete machine-readable "work that uses the Library", as object code and/or source code, so that the user can modify the Library and then relink to produce a modified executable containing the modified Library. (It is understood that the user who changes the contents of definitions files in the Library will not necessarily be able to recompile the application to use the modified definitions.)
- b) Use a suitable shared library mechanism for linking with the Library. A suitable mechanism is one that (1) uses at run time a copy of the library already present on the user's computer system,

rather than copying library functions into the executable, and (2) will operate properly with a modified version of the library, if the user installs one, as long as the modified version is interface-compatible with the version that the work was made with.

- c) Accompany the work with a written offer, valid for at least three years, to give the same user the materials specified in Subsection 6a, above, for a charge no more than the cost of performing this distribution.
- d) If distribution of the work is made by offering access to copy from a designated place, offer equivalent access to copy the above specified materials from the same place.
- e) Verify that the user has already received a copy of these materials or that you have already sent this user a copy.

For an executable, the required form of the "work that uses the Library" must include any data and utility programs needed for reproducing the executable from it. However, as a special exception,

the materials to be distributed need not include anything that is normally distributed (in either source or binary form) with the major components (compiler, kernel, and so on) of the operating system on

which the executable runs, unless that component itself accompanies the executable.

It may happen that this requirement contradicts the license restrictions of other proprietary libraries that do not normally accompany the operating system. Such a contradiction means you cannot

use both them and the Library together in an executable that you distribute.

- 7. You may place library facilities that are a work based on the Library side-by-side in a single library together with other library facilities not covered by this License, and distribute such a combined library, provided that the separate distribution of the work based on the Library and of the other library facilities is otherwise permitted, and provided that you do these two things:
- a) Accompany the combined library with a copy of the same work based on the Library, uncombined with any other library facilities. This must be distributed under the terms of the Sections above.
- b) Give prominent notice with the combined library of the fact that part of it is a work based on the Library, and explaining where to find the accompanying uncombined form of the same work.
- 8. You may not copy, modify, sublicense, link with, or distribute the Library except as expressly provided under this License. Any attempt otherwise to copy, modify, sublicense, link with, or distribute the Library is void, and will automatically terminate your rights under this License. However, parties who have received copies, or rights, from you under this License will not have their licenses terminated so long as such parties remain in full compliance.
- 9. You are not required to accept this License, since you have not signed it. However, nothing else grants you permission to modify or distribute the Library or its derivative works. These actions are prohibited by law if you do not accept this

License. Therefore, by modifying or distributing the Library (or any work based on the Library), you indicate your acceptance of this License to do so, and all its terms and conditions for copying, distributing or modifying the Library or works based on it.

10. Each time you redistribute the Library (or any work based on the Library), the recipient automatically receives a license from the original licensor to copy, distribute, link with or modify the Library subject to these terms and conditions. You may not impose any further restrictions on the recipients' exercise of the rights granted herein.

You are not responsible for enforcing compliance by third parties with this License.

11. If, as a consequence of a court judgment or allegation of patent infringement or for any other reason (not limited to patent issues), conditions are imposed on you (whether by court order, agreement or otherwise) that contradict the conditions of this License, they do not excuse you from the conditions of this License. If you cannot distribute so as to satisfy simultaneously your obligations under this License and any other pertinent obligations, then as a consequence you may not distribute the Library at all. For example, if a patent license would not permit royalty-free redistribution of the Library by all those who receive copies directly or indirectly through you, then the only way you could satisfy both it and this License would be to refrain entirely from distribution of the Library.

If any portion of this section is held invalid or unenforceable under any particular circumstance, the balance of the section is intended to apply, and the section as a whole is intended to apply in other circumstances.

It is not the purpose of this section to induce you to infringe any patents or other property right claims or to contest validity of any such claims; this section has the sole purpose of protecting the

integrity of the free software distribution system which is implemented by public license practices. Many people have made generous contributions to the wide range of software distributed through that system in reliance on consistent application of that system; it is up to the author/donor to decide if he or she is willing to distribute software through any other system and a licensee cannot

impose that choice.

This section is intended to make thoroughly clear what is believed to be a consequence of the rest of this License.

- 12. If the distribution and/or use of the Library is restricted in certain countries either by patents or by copyrighted interfaces, the original copyright holder who places the Library under this License may add an explicit geographical distribution limitation excluding those countries, so that distribution is permitted only in or among countries not thus excluded. In such case, this License incorporates the limitation as if written in the body of this License.
- 13. The Free Software Foundation may publish revised and/or new versions of the Lesser General Public License from time to time.

Such new versions will be similar in spirit to the present version, but may differ in detail to address new problems or concerns.

Each version is given a distinguishing version number. If the Library specifies a version number of this License which applies to it and "any later version", you have the option of following the terms and conditions either of that version or of any later version published by the Free Software Foundation. If the Library does not specify a license version number, you may choose any version ever published by the Free Software Foundation.

14. If you wish to incorporate parts of the Library into other free programs whose distribution conditions are incompatible with these, write to the author to ask for permission. For software which is

copyrighted by the Free Software Foundation, write to the Free Software Foundation; we sometimes make exceptions for this. Our decision will be guided by the two goals of preserving the free status

of all derivatives of our free software and of promoting the sharing and reuse of software generally.

NO WARRANTY

15. BECAUSE THE LIBRARY IS LICENSED FREE OF CHARGE, THERE IS NO WARRANTY FOR THE LIBRARY, TO THE EXTENT PERMITTED BY APPLICABLE LAW. EXCEPT WHEN OTHERWISE STATED IN WRITING THE COPYRIGHT HOLDERS AND/OR OTHER PARTIES PROVIDE THE LIBRARY "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESSED OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE. THE ENTIRE RISK AS TO THE QUALITY AND PERFORMANCE OF THE LIBRARY IS WITH YOU. SHOULD THE LIBRARY PROVE DEFECTIVE, YOU ASSUME THE COST OF ALL NECESSARY SERVICING, REPAIR OR CORRECTION.

16. IN NO EVENT UNLESS REQUIRED BY APPLICABLE LAW OR AGREED TO IN WRITING WILL ANY COPYRIGHT HOLDER, OR ANY OTHER PARTY WHO MAY MODIFY AND/OR REDISTRIBUTE THE LIBRARY AS PERMITTED ABOVE, BE LIABLE TO YOU FOR DAMAGES, INCLUDING ANY GENERAL, SPECIAL, INCIDENTAL OR CONSEQUENTIAL DAMAGES ARISING OUT OF THE USE OR INABILITY TO USE THE LIBRARY (INCLUDING BUT NOT LIMITED TO LOSS OF DATA OR DATA BEING

RENDERED INACCURATE OR LOSSES SUSTAINED BY YOU OR THIRD PARTIES OR A FAILURE OF THE LIBRARY TO OPERATE WITH ANY OTHER SOFTWARE), EVEN IF SUCH HOLDER OR OTHER PARTY HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

END OF TERMS AND CONDITIONS

How to Apply These Terms to Your New Libraries

If you develop a new library, and you want it to be of the greatest possible use to the public, we recommend making it free software that everyone can redistribute and change. You can do so by permitting redistribution under these terms (or, alternatively, under the terms of the ordinary General Public License).

To apply these terms, attach the following notices to the library. It is safest to attach them to the start of each source file to most effectively convey the exclusion of warranty; and each file should have at least the "copyright" line and a pointer to where the full notice is found.

<one line to give the library's name and a brief idea of what it does.>

Copyright (C) <year> <name of author>

This library is free software; you can redistribute it and/or modify it under the terms of the GNU Lesser General Public License as published by the Free Software Foundation; either version 2.1 of the License, or (at your option) any later version.

This library is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU

Lesser General Public License for more details.

You should have received a copy of the GNU Lesser General Public License along with this library; if not, write to the Free Software Foundation, Inc., 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301

Also add information on how to contact you by electronic and paper mail.

You should also get your employer (if you work as a programmer) or your school, if any, to sign a "copyright disclaimer" for the library, if necessary. Here is a sample; alter the names:

Yoyodyne, Inc., hereby disclaims all copyright interest in the library `Frob' (a library for tweaking knobs) written by James Random Hacker.

<signature of Ty Coon>, 1 April 1990

Ty Coon, President of Vice

That's all there is to it!

Software: JSON-C

Copyright notice:

Copyright (c) 2004, 2005 Metaparadigm Pte. Ltd.

Copyright (c) 2009-2012 Eric Haszlakiewicz

Copyright (c) 2004, 2005 Metaparadigm Pte Ltd

Copyright (c) 2009 Hewlett-Packard Development Company, L.P.

Copyright 2011, John Resig

Copyright 2011, The Dojo Foundation

Copyright (c) 2012 Eric Haszlakiewicz

Copyright (c) 2009-2012 Hewlett-Packard Development Company, L.P.

Copyright (c) 2008-2009 Yahoo! Inc. All rights reserved.

Copyright (C) 1996, 1997, 1998, 1999, 2000, 2001, 2003, 2004, 2005, 2006,

2007, 2008, 2009, 2010, 2011 Free Software Foundation, Inc.

Copyright (c) 2013 Metaparadigm Pte. Ltd.

License: MIT License

Copyright (c) 2009-2012 Eric Haszlakiewicz

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

Copyright (c) 2004, 2005 Metaparadigm Pte Ltd

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

Software: proj

Copyright notice:

"Copyright (C) 2010 Mateusz Loskot <mateusz@loskot.net>

Copyright (C) 2007 Douglas Gregor <doug.gregor@gmail.com>

Copyright (C) 2007 Troy Straszheim

CMake, Copyright (C) 2009-2010 Mateusz Loskot <mateusz@loskot.net>)

Copyright (C) 2011 Nicolas David <nicolas.david@ign.fr>

Copyright (c) 2000, Frank Warmerdam

Copyright (c) 2011, Open Geospatial Consortium, Inc.

Copyright (C) 1996, 1997, 1998, 1999, 2000, 2001, 2003, 2004, 2005, 2006,

2007, 2008, 2009, 2010, 2011 Free Software Foundation, Inc.

Copyright (c) Charles Karney (2012-2015) <charles@karney.com> and licensed

Copyright (c) 2005, Antonello Andrea

Copyright (c) 2010, Frank Warmerdam

Copyright (c) 1995, Gerald Evenden

Copyright (c) 2000, Frank Warmerdam < warmerdam@pobox.com >

Copyright (c) 2010, Frank Warmerdam < warmerdam@pobox.com>

Copyright (c) 2013, Frank Warmerdam

Copyright (c) 2003 Gerald I. Evenden

Copyright (c) 2012, Frank Warmerdam < warmerdam@pobox.com >

Copyright (c) 2002, Frank Warmerdam

Copyright (c) 2004 Gerald I. Evenden

Copyright (c) 2012 Martin Raspaud

Copyright (c) 2001, Thomas Flemming, tf@ttqv.com

Copyright (c) 2002, Frank Warmerdam < warmerdam@pobox.com >

Copyright (c) 2009, Frank Warmerdam

Copyright (c) 2003, 2006 Gerald I. Evenden

Copyright (c) 2011, 2012 Martin Lambers <marlam@marlam.de>

Copyright (c) 2006, Andrey Kiselev

Copyright (c) 2008-2012, Even Rouault <even dot rouault at mines-paris dot org>

Copyright (c) 2001, Frank Warmerdam

Copyright (c) 2001, Frank Warmerdam < warmerdam@pobox.com>

Copyright (c) 2008 Gerald I. Evenden

11

License: MIT License

Please see above

Software: libxml2 Copyright notice:

"See Copyright for the status of this software.

Copyright (C) 1998-2003 Daniel Veillard. All Rights Reserved.

Copyright (C) 2003 Daniel Veillard.

copy: see Copyright for the status of this software.

copy: see Copyright for the status of this software

copy: see Copyright for the status of this software.

Copyright (C) 2000 Bjorn Reese and Daniel Veillard.

Copy: See Copyright for the status of this software.

See COPYRIGHT for the status of this software

Copyright (C) 2000 Gary Pennington and Daniel Veillard.

Copyright (C) 1996, 1997, 1998, 1999, 2000, 2001, 2003, 2004, 2005, 2006,

2007 Free Software Foundation, Inc.

Copyright (C) 1998 Bjorn Reese and Daniel Stenberg.

Copyright (C) 2001 Bjorn Reese <bre> <bre>breese@users.sourceforge.net>

Copyright (C) 2000 Bjorn Reese and Daniel Stenberg.

Copyright (C) 2001 Bjorn Reese and Daniel Stenberg.

See Copyright for the status of this software

"

License: MIT License

Please see above

11 Using JDBC or ODBC for GaussDB(DWS) Secondary Development

11.1 Prerequisites

If the connection pool mechanism is used during application development, comply with the following specifications:

- If GUC parameters are set in the connection, before you return the connection to the connection pool, run SET SESSION AUTHORIZATION DEFAULT; RESET ALL; to clear the connection status.
- If a temporary table is used, delete it before you return the connection to the connection pool.

If you do not do so, the status of connections in the connection pool will remain, which affects subsequent operations using the connection pool.

Downloading Drivers

For details, see **Downloading the JDBC or ODBC Driver**.

11.2 JDBC-Based Development

11.2.1 JDBC Development Process

Java Database Connectivity (JDBC) is a Java API for executing SQL statements. It provides a unified access interface for multiple relational databases, enabling applications to work with data based on it. GaussDB(DWS) supports JDBC 4.0 and requires JDK 1.6 or later for code compiling. It does not support JDBC-ODBC Bridge. The following figure shows the JDBC application development process.

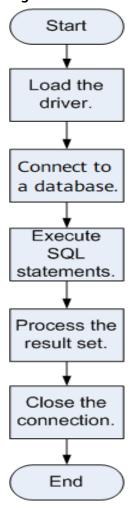


Figure 11-1 JDBC-based application development process

Table 11-1 JDBC development process

Procedure	Description
Load the driver.	Download the JDBC driver and edit and load it in the program.
Connect to a database.	Connect to the database through the JDBC driver.
Execute SQL statements.	Applications operate database data by executing SQL statements.
Process the result set.	Different types of result sets have different application scenarios. Applications need to select the appropriate result set type as needed.
Close the connection.	Make sure to close the database connection after completing the required data operations.

11.2.2 JDBC Package and Driver Class

JDBC Package

Download the **dws_8.x.x_jdbc_driver.zip** software package from the console.

For details, see **Downloading the JDBC or ODBC Driver**.

After the decompression, you will obtain the following JDBC packages in .jar format:

- gsjdbc4.jar: Driver package compatible with PostgreSQL. The class name and class structure in the driver are the same as those in the PostgreSQL driver. All the applications running on PostgreSQL can be smoothly transferred to the current system.
- **gsjdbc200.jar**: This driver package is used when both PostgreSQL and GaussDB(DWS) are accessed in a JVM process. The main class name is **com.huawei.gauss200.jdbc.Driver** and the prefix of the URL for database connection is **jdbc:gaussdb**. Other information of this driver package is the same as that of **gsjdbc4.jar**.

Driver Class

Before creating a database connection, you need to load the database driver class org.postgresql.Driver (decompressed from gsjdbc4.jar) or com.huawei.gauss200.jdbc.Driver (decompressed from gsjdbc200.jar).

□ NOTE

GaussDB(DWS) is compatible with PostgreSQL in the use of JDBC. If two JDBC drivers are used in the same process, class names may conflict.

11.2.3 Loading a Driver

Load the database driver before creating a database connection.

You can load the driver in the following ways:

- Implicitly loading the driver before creating a connection in the code: Class.forName ("org.postgresql.Driver")
- Transferring a parameter during the JVM startup: java -Djdbc.drivers=org.postgresql.Driver jdbctest

- **jdbctest** is the name of a test application.
- If gsjdbc200.jar is used, change the driver class name to "com.huawei.gauss200.jdbc.Driver".

11.2.4 Connecting to a Database

After a database is connected, you can run SQL statements the database to perform operations on data.

Prerequisites

If you use an open source JDBC driver, set **password_encryption_type** to **1**. If it is not **1**, the connection might fail with an error saying "none of the server's SASL authentication mechanisms are supported". To fix this, perform the following steps:

- 1. Check and change the password_encryption_type to 1 if needed with the help of technical support.
- 2. Create a new database user for connection or reset the password of the existing database user.
 - If you use an administrator account, reset the password. For details, see Resetting a Password.
 - If you are a common user, use another client tool (such as Data Studio) to connect to the database and run the ALTER USER statement to change your password.
- Connect to the database.

Here are the reasons why you need to perform these operations:

- MD5 algorithms may by vulnerable to collision attacks and cannot be used for password verification. By default, GaussDB(DWS) disables MD5 password verification, which can cause open source client connection failures. Therefore, you need to adjust the cryptographic algorithm parameter password_encryption_type to enable the MD5 algorithm.
- The database stores the hash digest of passwords instead of password text. During password verification, the system compares the hash digest with the password digest sent from the client (salt operations are involved). If you change your cryptographic algorithm policy, the database cannot generate a new MD5 hash digest for your existing password. For connectivity purposes, you must manually change your password or create a new user. The new password will be encrypted using the hash algorithm and stored for authentication in the next connection.

Function Prototype

JDBC provides the following three database connection methods:

- DriverManager.getConnection(String url);
- DriverManager.getConnection(String url, Properties info);
- DriverManager.getConnection(String url, String user, String password);

Parameter

Table 11-2 Database connection parameters

Parame ter	Description	
url	gsjdbc4.jar database connection descriptor. The descriptor format can be:	
	jdbc:postgresql:database	
	jdbc:postgresql://host/database	
	jdbc:postgresql://host:port/database	
	• jdbc:postgresql://host:port[,host:port][]/database	
	NOTE If gsjdbc200.jar is used, replace jdbc:postgresql with jdbc:gaussdb.	
	database: indicates the name of the database to be connected.	
	 host indicates the name or IP address of the database server. If an ELB is bound to the cluster, set host to the IP address of the ELB. For security purposes, the CN forbids access from other nodes in the cluster without authentication. To access the CN from inside the cluster, deploy the JDBC program on the host where the CN is located and set host to 127.0.0.1. If you do not do so, the error message "FATAL: Forbid remote connection with trust method!" may be displayed. 	
	It is recommended that the service system be deployed outside the cluster. If it is deployed inside, the database performance may be affected.	
	 port: indicates the port number of a database server. By default, the database on port 8000 of the local host is connected. 	
	 Multiple IP addresses and ports can be configured. JDBC balances load by random access and failover, and will automatically ignore unreachable IP addresses. IP addresses are separated using commas. Example: jdbc:postgresql:// 	
	10.10.0.13:8000,10.10.0.14:8000/database	
	If JDBC is used to connect to a cluster, only JDBC connection parameters can be configured in a cluster address. Variables cannot be added.	

Parame ter	Description
info	Database connection properties. Common properties include:
	• user : string type. It indicates the database user establishing a connection.
	• password : string type. It indicates the password of a database user.
	• ssl : Boolean type. It indicates whether the Secure Socket Layer (SSL) is used.
	loggerLevel: string type. It indicates the amount of information that the driver logs and prints to the LogStream or LogWriter specified in the DriverManager. Currently, OFF, DEBUG, and TRACE are supported. DEBUG indicates that only logs of the DEBUG or higher level are printed, generating a few log information. TRACE indicates that logs of the DEBUG and TRACE levels are printed, generating detailed log information. The default value is OFF, indicating that no information will be logged.
	• prepareThreshold: integer type. It indicates the number of PreparedStatement executions required before SQL statements are switched over to servers as prepared statements. The default value is 5.
	batchMode: boolean type. It indicates whether to connect the database in batch mode.
	fetchsize: integer type. It indicates the default fetchsize for statements in the created connection.
	ApplicationName: string type. It indicates an application name. The default value is PostgreSQL JDBC Driver.
	allowReadOnly: boolean type. It indicates whether to enable the read-only mode for connection. The default value is false. If the value is not changed to true, the execution of connection.setReadOnly does not take effect.
	• blobMode : string type. It is used to set the setBinaryStream method to assign values to different data types. The value on indicates that values are assigned to the BLOB data type and off indicates that values are assigned to the bytea data type. The default value is on .
	• connectionExtraInfo : boolean type. It indicates whether the JDBC driver reports the driver deployment path and process owner to the database.
	NOTE The value can be true or false . The default value is true . If connectionExtraInfo is set to true , the JDBC driver reports the driver deployment path and process owner to the database and displays the information in the connection_info parameter (see connection_info). In this case, you can query the information from PG_STAT_ACTIVITY or PGXC_STAT_ACTIVITY .
user	Indicates a database user.

Parame ter	Description
passwor d	Indicates the password of a database user.

Closing the Connection

Make sure to close the database connection after completing the required data operations.

To close the database connection, you can directly invoke the **close** method, for example, **conn.close()**.

Examples

```
//gsjdbc4.jar is used as an example. If gsjdbc200.jar is used, replace the driver class name org.postgresql
with com.huawei.gauss200.jdbc and replace the URL prefix jdbc:postgresql with jdbc:gaussdb.
//The following code encapsulates database connection operations into an interface. The database can then
be connected using an authorized username and password.
public static Connection GetConnection(String username, String passwd) {
     //Set the driver class.
     String driver = "org.postgresql.Driver";
     //Database connection descriptor.
     String sourceURL = "jdbc:postgresql://10.10.0.13:8000/postgres?currentSchema=test";
     Connection conn = null;
     try {
        //Load the driver.
        Class.forName(driver);
     } catch (ClassNotFoundException e ){
        e.printStackTrace();
       return null;
     }
        //Establish a connection.
       conn = DriverManager.getConnection(sourceURL, username, passwd);
        System.out.println("Connection succeed!");
     } catch (SQLException e) {
       e.printStackTrace();
        return null;
     return conn:
```

11.2.5 Executing SQL Statements

Executing an Ordinary SQL Statement

The application performs data (parameter statements do not need to be transferred) in the database by running SQL statements, and you need to perform the following steps:

Step 1 Create a statement object by triggering the createStatement method in Connection.

Statement stmt = con.createStatement();

Step 2 Execute the SQL statement by triggering the executeUpdate method in Statement. int rc = stmt.executeUpdate("CREATE TABLE customer_t1(c_customer_sk INTEGER, c_customer_name VARCHAR(32));");

∩ NOTE

If an execution request (not in a transaction block) received in the database contains multiple statements, the request is packed into a transaction. **VACUUM** is not supported in a transaction block. If one of the statements fails, the entire request will be rolled back.

Step 3 Close the statement object.

stmt.close();

----End

Executing a Prepared SQL Statement

Pre-compiled statements were once complied and optimized and can have additional parameters for different usage. For the statements have been pre-compiled, the execution efficiency is greatly improved. If you want to execute a statement for several times, use a precompiled statement. Perform the following procedure:

Step 1 Create a prepared statement object by calling the prepareStatement method in Connection.

PreparedStatement pstmt = con.prepareStatement("UPDATE customer_t1 SET c_customer_name = ? WHERE c_customer_sk = 1");

- **Step 2** Set parameters by triggering the setShort method in PreparedStatement. pstmt.setShort(1, (short)2);
- **Step 3** Execute the precompiled SQL statement by triggering the executeUpdate method in PreparedStatement.

int rowcount = pstmt.executeUpdate();

Step 4 Close the precompiled statement object by calling the close method in PreparedStatement.

pstmt.close();

----End

Calling a Stored Procedure

Perform the following steps to call existing stored procedures through the JDBC interface in GaussDB(DWS):

- **Step 1** Create a call statement object by calling the prepareCall method in Connection. CallableStatement cstmt = myConn.prepareCall("{? = CALL TESTPROC(?,?,?)}");
- **Step 2** Set parameters by calling the setInt method in CallableStatement.

```
cstmt.setInt(2, 50);
cstmt.setInt(1, 20);
cstmt.setInt(3, 90);
```

Step 3 Register with an output parameter by calling the registerOutParameter method in CallableStatement.

cstmt.registerOutParameter(4, Types.INTEGER); //Register an OUT parameter as an integer.

- **Step 4** Call the stored procedure by calling the execute method in CallableStatement. cstmt.execute();
- **Step 5** Obtain the output parameter by calling the getInt method in CallableStatement. int out = cstmt.getInt(4); //Obtain the OUT parameter.

For example:

```
//The following stored procedure has been created with the OUT parameter:
create or replace procedure testproc
(
    psv_in1 in integer,
    psv_in2 in integer,
    psv_inout in out integer
)
as
begin
    psv_inout := psv_in1 + psv_in2 + psv_inout;
end;
/
```

Step 6 Close the call statement by calling the close method in CallableStatement. cstmt.close();

◯ NOTE

- Many database classes such as Connection, Statement, and ResultSet have a close()
 method. Close these classes after using their objects. Close these actions after using
 their objects. Closing Connection will close all the related Statements, and closing a
 Statement will close its ResultSet.
- Some JDBC drivers support named parameters, which can be used to set parameters by name rather than sequence. If a parameter has a default value, you do not need to specify any parameter value but can use the default value directly. Even though the parameter sequence changes during a stored procedure, the application does not need to be modified. Currently, the GaussDB(DWS) JDBC driver does not support this method.
- GaussDB(DWS) does not support functions containing OUT parameters, or default values of stored procedures and function parameters.

----End

NOTICE

- If JDBC is used to call a stored procedure whose returned value is a cursor, the returned cursor cannot be used.
- A stored procedure and an SQL statement must be executed separately.

Batch Processing

When a prepared statement batch processes multiple pieces of similar data, the database creates only one execution plan. This improves the compilation and optimization efficiency. Perform the following procedure:

Step 1 Create a prepared statement object by calling the prepareStatement method in Connection.

PreparedStatement pstmt = con.prepareStatement("INSERT INTO customer_t1 VALUES (?)");

Step 2 Call the setShort parameter for each piece of data, and call addBatch to confirm that the setting is complete.

pstmt.setShort(1, (short)2);
pstmt.addBatch();

Step 3 Execute batch processing by calling the executeBatch method in PreparedStatement.

int[] rowcount = pstmt.executeBatch();

Step 4 Close the precompiled statement object by calling the close method in PreparedStatement.

pstmt.close();

Do not terminate a batch processing action when it is ongoing; otherwise, the database performance will deteriorate. Therefore, disable the automatic submission function during batch processing, and manually submit every several lines. The statement for disabling automatic submission is **conn.setAutoCommit(false)**.

----End

11.2.6 Processing Data in a Result Set

Setting a Result Set Type

Different types of result sets are applicable to different application scenarios. Applications select proper types of result sets based on requirements. Before executing an SQL statement, you must create a statement object. Some methods of creating statement objects can set the type of a result set. **Table 11-3** lists result set parameters. The connection methods are as follows:

//Create a Statement object. This object will generate a ResultSet object with a specified type and concurrency.

createStatement(int resultSetType, int resultSetConcurrency);

//Create a PreparedStatement object. This object will generate a ResultSet object with a specified type and concurrency.

prepareStatement(String sql, int resultSetType, int resultSetConcurrency);

//Create a CallableStatement object. This object will generate a ResultSet object with a specified type and concurrency.

prepareCall(String sql, int resultSetType, int resultSetConcurrency);

Table 11-3 Result set types

Parameter	Description
resultSetType	Indicates the type of a result set. There are three types of result sets:
	• ResultSet.TYPE_FORWARD_ONLY: The ResultSet object can only be navigated forward. It is the default value.
	ResultSet.TYPE_SCROLL_SENSITIVE: You can view the modified result by scrolling to the modified row.
	 ResultSet.TYPE_SCROLL_INSENSITIVE: The ResultSet object is insensitive to changes in the underlying data source.
	After a result set has obtained data from the database, the result set is insensitive to data changes made by other transactions, even if the result set type is ResultSet.TYPE_SCROLL_SENSITIVE. To obtain up-to-date data of the record pointed by the cursor from the database, call the refreshRow() method in a ResultSet object.
resultSetConcurren- cy	Indicates the concurrency type of a result set. There are two types of concurrency.
	ResultSet.CONCUR_READ_ONLY: The data in a result set cannot be updated except that an updated statement has been created in the result set data.
	ResultSet.CONCUR_UPDATEABLE: changeable result set. The concurrency type for a result set object can be updated if the result set is scrollable.

Positioning a Cursor in a Result Set

ResultSet objects include a cursor pointing to the current data row. The cursor is initially positioned before the first row. The next method moves the cursor to the next row from its current position. When a **ResultSet** object does not have a next row, a call to the next method returns **false**. Therefore, this method is used in the while loop for result set iteration. However, the JDBC driver provides more cursor positioning methods for scrollable result sets, which allows positioning cursor in the specified row. **Table 11-4** lists these methods.

Table 11-4 Methods for positioning a cursor in a result set

Method	Description
next()	Moves cursor to the next row from its current position.
previous()	Moves cursor to the previous row from its current position.

Method	Description
beforeFirst()	Places cursor before the first row.
afterLast()	Places cursor after the last row.
first()	Places cursor to the first row.
last()	Places cursor to the last row.
absolute(int)	Places cursor to a specified row.
relative(int)	Moves cursor forward or backward a specified number of rows.

Obtaining the Cursor Position from a Result Set

This cursor positioning method will be used to change the cursor position for a scrollable result set. JDBC driver provides a method to obtain the cursor position in a result set. Table 11-5 lists the method.

Table 11-5 Method for obtaining the cursor position in a result set

Method	Description
isFirst()	Checks whether the cursor is in the first row.
isLast()	Checks whether the cursor is in the last row.
isBeforeFirst()	Checks whether the cursor is before the first row.
isAfterLast()	Checks whether the cursor is after the last row.
getRow()	Gets the current row number of the cursor.

Obtaining Data in a Result Set

ResultSet objects provide a variety of methods to obtain data from a result set. **Table 11-6** lists the common methods for obtaining data. If you want to know more about other methods, see JDK official documents.

Table 11-6 Common methods for obtaining data from a result set

Method	Description
int getInt(int columnIndex)	Retrieves data of the int type by column.

Method	Description
int getInt(String columnLabel)	Retrieves data of the int type by column name.
String getString(int columnIndex)	Retrieves data of the string type by column.
String getString(String columnLabel)	Retrieves data of the string type by column name.
Date getDate(int columnIndex)	Retrieves data of the date type by column.
Date getDate(String columnLabel)	Retrieves data of the date type by column name.

11.2.7 Common JDBC Development Examples

Example 1

Before completing the following example, you need to create a stored procedure.

```
create or replace procedure testproc
(
    psv_in1 in integer,
    psv_in2 in integer,
    psv_inout in out integer
)
as
begin
    psv_inout := psv_in1 + psv_in2 + psv_inout;
end;
/
```

This example illustrates how to develop applications based on the GaussDB(DWS) JDBC interface.

```
//DBtest.java
//gsjdbc4.jar is used as an example. If gsjdbc200.jar is used, replace the driver class name org.postgresql
with com.huawei.gauss200.jdbc and replace the URL prefix jdbc:postgresql with jdbc:gaussdb.
// This example illustrates the main processes of JDBC-based development, covering database connection
creation, table creation, and data insertion.
import java.sql.Connection;
import java.sql.DriverManager;
import java.sql.PreparedStatement;
import java.sql.SQLException;
import java.sql.Statement;
import java.sql.CallableStatement;
public class DBTest {
 //Establish a connection to the database.
 public static Connection GetConnection(String username, String passwd) {
  String driver = "org.postgresql.Driver";
  String sourceURL = "jdbc:postgresql://localhost:/gaussdb";
  Connection conn = null;
    //Load the database driver.
    Class.forName(driver).newInstance();
```

```
} catch (Exception e) {
   e.printStackTrace();
   return null;
  try {
   //Establish a connection to the database.
   conn = DriverManager.getConnection(sourceURL, username, passwd);
   System.out.println("Connection succeed!");
  } catch (Exception e) {
   e.printStackTrace();
   return null;
  }
  return conn;
 };
 //Run an ordinary SQL statement. Create a customer_t1 table.
 public static void CreateTable(Connection conn) {
  Statement stmt = null;
  try {
   stmt = conn.createStatement();
   //Run an ordinary SQL statement.
   int rc = stmt
      .executeUpdate("CREATE TABLE customer_t1(c_customer_sk INTEGER, c_customer_name
VARCHAR(32));");
   stmt.close();
  } catch (SQLException e) {
   if (stmt != null) {
     try {
      stmt.close();
     } catch (SQLException e1) {
      e1.printStackTrace();
   e.printStackTrace();
 //Run the preprocessing statement to insert data in batches.
 public static void BatchInsertData(Connection conn) {
  PreparedStatement pst = null;
  try {
   //Generate a prepared statement.
   pst = conn.prepareStatement("INSERT INTO customer_t1 VALUES (?,?)");
   for (int i = 0; i < 3; i++) {
     //Add parameters.
     pst.setInt(1, i);
     pst.setString(2, "data " + i);
     pst.addBatch();
   //Run batch processing.
   pst.executeBatch();
   pst.close();
  } catch (SQLException e) {
   if (pst != null) {
     try {
      pst.close();
     } catch (SQLException e1) {
     e1.printStackTrace();
   e.printStackTrace();
```

```
//Run the precompilation statement to update data.
 public static void ExecPreparedSQL(Connection conn) {
  PreparedStatement pstmt = null;
  try {
    pstmt = conn
      .prepareStatement("UPDATE customer_t1 SET c_customer_name = ? WHERE c_customer_sk = 1");
    pstmt.setString(1, "new Data");
    int rowcount = pstmt.executeUpdate();
    pstmt.close();
  } catch (SQLException e) {
    if (pstmt != null) {
     try {
      pstmt.close();
     } catch (SQLException e1) {
      e1.printStackTrace();
    e.printStackTrace();
//Run a stored procedure.
 public static void ExecCallableSQL(Connection conn) {
  CallableStatement cstmt = null;
    cstmt=conn.prepareCall("{? = CALL TESTPROC(?,?,?)}");
    cstmt.setInt(2, 50);
    cstmt.setInt(1, 20);
    cstmt.setInt(3, 90);
    cstmt.registerOutParameter(4, Types.INTEGER); //Register an OUT parameter as an integer.
    cstmt.execute();
    int out = cstmt.getInt(4); //Obtain the out parameter value.
    System.out.println("The CallableStatment TESTPROC returns:"+out);
    cstmt.close();
  } catch (SQLException e) {
    if (cstmt != null) {
     try {
      cstmt.close();
     } catch (SQLException e1) {
      e1.printStackTrace();
    e.printStackTrace();
  * Main process. Call static methods one by one.
 * @param args
 public static void main(String[] args) {
  //Establish a connection to the database.
  Connection conn = GetConnection("tester", "password");
  //Create a table.
  CreateTable(conn);
  //Insert data in batches.
  BatchInsertData(conn);
 //Run the precompilation statement to update data.
  ExecPreparedSQL(conn);
  //Run a stored procedure.
  ExecCallableSQL(conn);
```

```
//Close the connection to the database.

try {
    conn.close();
} catch (SQLException e) {
    e.printStackTrace();
}

}
```

Example 2: High Client Memory Usage

In this example, **setFetchSize** adjusts the memory usage of the client by using the database cursor to obtain server data in batches. It may increase network interaction and damage some performance.

The cursor is valid within a transaction. Therefore, you need to disable the autocommit function.

```
// Disable the autocommit function.
conn.setAutoCommit(false);
Statement st = conn.createStatement();
// Open the cursor and obtain 50 lines of data each time.
st.setFetchSize(50);
ResultSet rs = st.executeQuery("SELECT * FROM mytable");
while (rs.next()){
  System.out.print("a row was returned.");
rs.close();
// Disable the server cursor.
st.setFetchSize(0);
rs = st.executeQuery("SELECT * FROM mytable");
while (rs.next()){
  System.out.print("many rows were returned.");
rs.close():
// Close the statement.
st.close();
```

Retrying SQL Queries for Applications

If the primary DN is faulty and cannot be restored within 40 seconds, its standby is automatically promoted to primary to ensure that the cluster runs properly. Jobs running during the switchover will fail and those started after the switchover will not be affected. To protect upper-layer services from being affected by the failover, refer to the following example to construct a SQL retry mechanism at the service layer.

gsjdbc4.jar is used as an example. If **gsjdbc200.jar** is used, replace the driver class name **org.postgresql** with **com.huawei.gauss200.jdbc** and replace the URL prefix **jdbc:postgresql** with **jdbc:qaussdb**.

```
import java.sql.Connection;
import java.sql.PreparedStatement;
import java.sql.ResultSet;
import java.sql.SQLException;
import java.sql.Statement;
```

```
class ExitHandler extends Thread {
  private Statement cancel_stmt = null;
  public ExitHandler(Statement stmt) {
     super("Exit Handler");
     this.cancel_stmt = stmt;
  public void run() {
     System.out.println("exit handle");
     try {
        this.cancel_stmt.cancel();
     } catch (SQLException e) {
        System.out.println("cancel query failed.");
        e.printStackTrace();
  }
public class SQLRetry {
 //Establish a connection to the database.
  public static Connection GetConnection(String username, String passwd) {
   String driver = "org.postgresql.Driver";
   String sourceURL = "jdbc:postgresql://10.131.72.136:8000/gaussdb";
   Connection conn = null;
   try {
   //Load the database driver.
    Class.forName(driver).newInstance();
   } catch (Exception e) {
    e.printStackTrace();
    return null;
   //Establish a connection to the database.
    conn = DriverManager.getConnection(sourceURL, username, passwd);
    System.out.println("Connection succeed!");
   } catch (Exception e) {
    e.printStackTrace();
    return null;
   return conn;
```

Run an ordinary SQL statement. Create the jdbc_test1 table.

```
} catch (SQLException e1) {
    e1.printStackTrace();
}
e.printStackTrace();
}
```

Run the preprocessing statement to insert data in batches.

```
public static void BatchInsertData(Connection conn) {
   PreparedStatement pst = null;
   try {
    //Generate a prepared statement.
     pst = conn.prepareStatement("INSERT INTO jdbc_test1 VALUES (?,?)");
     for (int i = 0; i < 100; i++) {
     //Add parameters.
      pst.setInt(1, i);
      pst.setString(2, "data " + i);
      pst.addBatch();
    //Run batch processing.
     pst.executeBatch();
     pst.close();
   } catch (SQLException e) {
     if (pst != null) {
      try {
       pst.close();
      } catch (SQLException e1) {
      e1.printStackTrace();
     e.printStackTrace();
```

Run the precompilation statement to update data.

```
private static boolean QueryRedo(Connection conn){
  PreparedStatement pstmt = null;
  boolean retValue = false;
  try {
    pstmt = conn
       .prepareStatement("SELECT col1 FROM jdbc_test1 WHERE col2 = ?");
       pstmt.setString(1, "data 10");
       ResultSet rs = pstmt.executeQuery();
       while (rs.next()) {
         System.out.println("col1 = " + rs.getString("col1"));
      rs.close();
    pstmt.close();
    retValue = true;
   } catch (SQLException e) {
    System.out.println("catch..... retValue " + retValue);
    if (pstmt != null) {
     try {
      pstmt.close();
    } catch (SQLException e1) {
      e1.printStackTrace();
    }
    e.printStackTrace();
   System.out.println("finish.....");
  return retValue;
```

```
}
```

Run a query statement and retry upon a failure. The number of retry times can be configured.

```
public static void ExecPreparedSQL(Connection conn) throws InterruptedException {
     int maxRetryTime = 50;
     int time = 0;
     String result = null;
     do {
       time++;
        try {
System.out.println("time:" + time);
boolean ret = QueryRedo(conn);
if(ret == false){
 System.out.println("retry, time:" + time);
 Thread.sleep(10000);
 QueryRedo(conn);
        } catch (Exception e) {
          e.printStackTrace();
     } while (null == result && time < maxRetryTime);
 }
 * Main process. Call static methods one by one.
  * @param args
* @throws InterruptedException
 public static void main(String[] args) throws InterruptedException {
//Establish a connection to the database.
  Connection conn = GetConnection("testuser", "test@123");
 //Create a table.
  CreateTable(conn);
 //Insert data in batches.
  BatchInsertData(conn);
//Run the precompilation statement to update data.
  ExecPreparedSQL(conn);
 //Close the connection to the database.
  try {
   conn.close();
  } catch (SQLException e) {
    e.printStackTrace();
 }
```

Importing and Exporting Data Through Local Files

When the JAVA language is used for secondary development based on GaussDB(DWS), you can use the CopyManager interface to export data from the database to a local file or import a local file to the database by streaming. The file can be in CSV or TEXT format.

The sample program is as follows. Load the GaussDB(DWS) JDBC driver before running it.

gsjdbc4.jar is used as an example. If **gsjdbc200.jar** is used, replace the driver class name **org.postgresql** with **com.huawei.gauss200.jdbc** and replace the URL prefix **jdbc:postgresql** with **jdbc:qaussdb**.

```
import java.sql.Connection;
import java.sql.DriverManager;
import java.io.IOException;
import java.io.FileInputStream;
import java.io.FileOutputStream;
import java.sql.SQLException;
import org.postgresql.copy.CopyManager;
import org.postgresql.core.BaseConnection;
public class Copy{
   public static void main(String[] args)
    String urls = new String("jdbc:postgresql://10.180.155.74:8000/gaussdb"); //URL of the database
    String username = new String("jack");
String password = new String("*******");
                                                   //Username
                                                  //Password
    String tablename = new String("migration_table"); //Define table information.
    String tablename1 = new String("migration_table_1"); //Define table information.
    String driver = "org.postgresql.Driver";
    Connection conn = null;
    try {
        Class.forName(driver);
        conn = DriverManager.getConnection(urls, username, password);
       } catch (ClassNotFoundException e) {
          e.printStackTrace(System.out);
       } catch (SQLException e) {
          e.printStackTrace(System.out);
```

Import and export data.

```
//Export the query result of migration_table to the local file d:/data.txt.
  copyToFile(conn, "d:/data.txt", "(SELECT * FROM migration_table)");
 } catch (SQLException e) {
 // TODO Auto-generated catch block
 e.printStackTrace();
 } catch (IOException e) {
 // TODO Auto-generated catch block
 e.printStackTrace();
   //Import data from the d:/data.txt file to the migration_table_1 table.
   copyFromFile(conn, "d:/data.txt", migration_table_1);
 } catch (SQLException e) {
 // TODO Auto-generated catch block
     e.printStackTrace();
} catch (IOException e) {
 // TODO Auto-generated catch block
 e.printStackTrace();
   //Export the data from the migration_table_1 table to the d:/data1.txt file.
   copyToFile(conn, "d:/data1.txt", migration_table_1);
 } catch (SQLException e) {
 // TODO Auto-generated catch block
 e.printStackTrace();
 } catch (IOException e) {
 // TODO Auto-generated catch block
e.printStackTrace();
}
  }
public static void copyFromFile(Connection connection, String filePath, String tableName)
```

```
throws SQLException, IOException {
  FileInputStream fileInputStream = null;
     CopyManager copyManager = new CopyManager((BaseConnection);
     fileInputStream = new FileInputStream(filePath);
     copyManager.copyIn("COPY" + tableName + "FROM STDIN", fileInputStream);
  } finally {
     if (fileInputStream != null) {
        try {
          fileInputStream.close();
        } catch (IOException e) {
          e.printStackTrace();
       }
     }
  }
}
public static void copyToFile(Connection connection, String filePath, String tableOrQuery)
      throws SQLException, IOException {
   FileOutputStream fileOutputStream = null;
   try {
      CopyManager copyManager = new CopyManager((BaseConnection)connection);
      fileOutputStream = new FileOutputStream(filePath);
      copyManager.copyOut("COPY" + tableOrQuery + "TO STDOUT", fileOutputStream);
   } finally {
      if (fileOutputStream != null) {
        try {
           fileOutputStream.close();
        } catch (IOException e) {
           e.printStackTrace();
        }
  }
```

Migrating Data from MySQL to GaussDB(DWS)

The following example shows how to use CopyManager to migrate data from MySQL to GaussDB(DWS).

gsjdbc4.jar is used as an example. If **gsjdbc200.jar** is used, replace the driver class name **org.postgresql** with **com.huawei.gauss200.jdbc** and replace the URL prefix **jdbc:postgresql** with **jdbc:gaussdb**.

```
import java.io.StringReader;
import java.sql.Connection;
import java.sql.DriverManager;
import java.sql.ResultSet;
import java.sql.SQLException;
import java.sql.Statement;
import org.postgresql.copy.CopyManager;
import org.postgresql.core.BaseConnection;
public class Migration{
  public static void main(String[] args) {
     String url = new String("jdbc:postgresql://10.180.155.74:8000/gaussdb"); //URL of the database
     String user = new String("jack");
                                             //GaussDB(DWS) username
     String pass = new String("******");
                                                //GaussDB(DWS) password
     String tablename = new String("migration_table"); //Define table information.
     String delimiter = new String("|");
                                               //Define a delimiter.
```

```
String encoding = new String("UTF8");
                                                     //Define a character set.
     String driver = "org.postgresql.Driver";
     StringBuffer buffer = new StringBuffer();
                                                   //Define the buffer to store formatted data.
        //Obtain the query result set of the source database.
        ResultSet rs = getDataSet();
        //Traverse the result set and obtain records row by row.
        //The values of columns in each record are separated by the specified delimiter and end with a
newline character to form strings.
        ////Add the strings to the buffer.
        while (rs.next()) {
           buffer.append(rs.getString(1) + delimiter
                + rs.getString(2) + delimiter
                + rs.getString(3) + delimiter
                + rs.getString(4)
                + "\n");
        rs.close();
        trv {
           //Connect to the target database.
           Class.forName(driver);
           Connection conn = DriverManager.getConnection(url, user, pass);
           BaseConnection baseConn = (BaseConnection) conn;
           baseConn.setAutoCommit(false);
           //Initialize table information.
String sql = "Copy " + tablename + " from STDIN DELIMITER " + "'" + delimiter + "'" + "
ENCODING " + "'" + encoding + "'";
           //Submit data in the buffer.
           CopyManager cp = new CopyManager(baseConn);
           StringReader reader = new StringReader(buffer.toString());
           cp.copyIn(sql, reader);
           baseConn.commit();
           reader.close();
           baseConn.close();
        } catch (ClassNotFoundException e) {
           e.printStackTrace(System.out);
        } catch (SQLException e) {
           e.printStackTrace(System.out);
     } catch (Exception e) {
        e.printStackTrace();
  }
```

Return the query result from the source database.

```
private static ResultSet getDataSet() {
    ResultSet rs = null;
    try {
        Class.forName("com.mysql.jdbc.Driver").newInstance();
        Connection conn = DriverManager.getConnection("jdbc:mysql://10.119.179.227:3306/jack?
useSSL=false&allowPublicKeyRetrieval=true", "jack", "********");
    Statement stmt = conn.createStatement();
    rs = stmt.executeQuery("select * from migration_table");
    } catch (SQLException e) {
        e.printStackTrace();
    } catch (Exception e) {
        e.printStackTrace();
    }
    return rs;
}
```

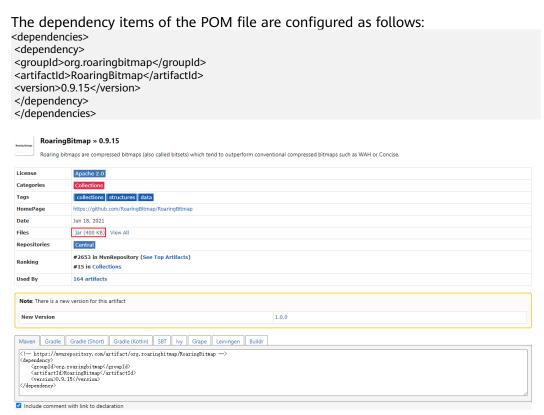
11.2.8 Processing RoaringBitmap Result Sets and Importing It to GaussDB (DWS)

GaussDB(DWS) 8.1.3 and later versions support the RoaringBitmap function. When using the Java language to perform secondary development based on GaussDB(DWS), you can use the CopyManager interface to import a small amount of RoaringBitmap data to GaussDB(DWS).

- Computing capability needs to be improved on the application side if a large amount of data needs to be imported to the database. Otherwise, the import performance will be affected.
- The 64-bit Roaringbitmap generated by Java cannot be imported to the database.

Processing RoaringBitmap Data

Step 1 Visit Maven to download the open-source RoaringBitmap JAR package. Version 0.9.15 is recommended.



Step 2 Invoke the JAR package to convert data to the RoaringBitmap type.

The general process is to declare a Roaring bitmap, call the add() method to convert data of the int type into the Roaringbitmap type, and then serialize the converted data. The sample code is as follows:

```
RoaringBitmap rr2 = new RoaringBitmap ();
for (int i = 1; i < 10000000; i++) {
    rr2.add(i);
}
ByteArrayOutputStream a = new ByteArrayOutputStream();</pre>
```

```
DataOutputStream b = new DataOutputStream(a);
rr2.serialize(b);
```

----End

Data Import

Invoke CopyManager to import data to the database. In this way, a small amount of RoaringBitmap data can be imported to the database without having to be stored locally.

```
//gsjdbc4.jar is used as an example. If gsjdbc200.jar is used, replace the driver class name org.postgresql
with com.huawei.gauss200.jdbc and replace the URL prefix jdbc:postgresql with jdbc:gaussdb.
package rb_demo;
import org.postgresql.copy.CopyManager;
import org.postgresql.core.BaseConnection;
import org.roaringbitmap.RoaringBitmap;
import java.io.ByteArrayInputStream;
import java.io.ByteArrayOutputStream;
import java.io.DataOutputStream;
import java.io.IOException;
import java.io.InputStream;
import java.io.StringReader;
import java.sql.Connection;
import java.sql.DriverManager;
import java.sql.PreparedStatement;
import java.sql.ResultSet;
import java.sql.SQLException;
import java.sql.Statement;
public class rb_demo {
  private static String hexStr = "0123456789ABCDEF";
  public static String bytesToHex(byte[] bytes) {
     StringBuffer sb = new StringBuffer();
     for (int i = 0; i < bytes.length; i++) {
        String hex = Integer.toHexString(bytes[i] & 0xFF);
        if (hex.length() < 2) {
          sb.append(0);
        sb.append(hex);
     return sb.toString();
  }
  public static Connection GetConnection(String username, String passwd) {
     String driver = "org.postgresql.Driver"
String sourceURL = "jdbc:postgresql://10.185.180.161: 8000/gaussdb"; //Database URL
     Connection conn = null;
     try {
    //Load the database driver.
        Class.forName(driver).newInstance();
     } catch (Exception e) {
       e.printStackTrace();
        return null;
     }
     try {
  //Establish a connection to the database.
        conn = DriverManager.getConnection(sourceURL, username, passwd);
        System.out.println("Connection succeed!");
     } catch (Exception e) {
       e.printStackTrace();
```

```
return null;
     }
     return conn;
  public static void main(String[] args) throws IOException {
     RoaringBitmap rr2 = new RoaringBitmap();
     for (int i = 1; i < 10000000; i++) {
       rr2.add(i);
     ByteArrayOutputStream a = new ByteArrayOutputStream();
     DataOutputStream b = new DataOutputStream(a);
     rr2.serialize(b);
Connection conn = GetConnection("test", "Gauss_234"); //User name and password.
     Statement pstmt = null:
     try {
       conn.setAutoCommit(true);
       pstmt = conn.createStatement();
       pstmt.execute("drop table if exists t_rb");
       pstmt.execute("create table t_rb(c1 int, c2 roaringbitmap) distribute by hash (c1);");
       StringReader sr = null;
       CopyManager cm = null;
       cm = new CopyManager((BaseConnection) conn);
       String delimiter = "|";
       StringBuffer tuples = new StringBuffer();
       tuples.append("1" + delimiter + "\\x" + bytesToHex(a.toByteArray()));
       StringBuffer sb = new StringBuffer();
       sb.append(tuples.toString());
       sr = new StringReader(tuples.toString());
       String sql = "copy t_rb from STDIN with (delimiter '|', NOESCAPING)";
long rows = cm.copyIn(sql, sr);//Execute the COPY command to save data to the database.
        pstmt.close();
     } catch (SQLException e) {
       if (pstmt != null) {
          try {
             pstmt.close();
          } catch (SQLException e1) {
             e1.printStackTrace();
       e.printStackTrace();
  }
```

11.2.9 JDBC Interfaces

JDBC interface is a set of API methods for users. This section describes some common interfaces. For other interfaces, see information in JDK1.6 (software package) and JDBC 4.0.

java.sql.Connection

This section describes **java.sql.Connection**, the interface for connecting to a database.

Table 11-7 java.sql.Connection methods

Method	Return Type	Support JDBC 4 or Not
close()	void	Yes
commit()	void	Yes
createStatement()	Statement	Yes
getAutoCommit()	boolean	Yes
getClientInfo()	Properties	Yes
getClientInfo(String name)	String	Yes
getTransactionIsolation()	int	Yes
isClosed()	boolean	Yes
isReadOnly()	boolean	Yes
prepareStatement(String sql)	PreparedStatement	Yes
rollback()	void	Yes
setAutoCommit(boolean autoCommit)	void	Yes
setClientInfo(Properties properties)	void	Yes
setClientInfo(String name,String value)	void	Yes

NOTICE

The interface uses the AutoCommit mode by default, but you can disable it by setting **setAutoCommit** to **false**. This will package all subsequent statements in explicit transactions. Note that you will not be able to execute statements that cannot be executed within transactions.

java.sql.CallableStatement

This section describes **java.sql.CallableStatement**, the stored procedure execution interface.

Table 11-8 java.sql.CallableStatement methods

Method	Return Type	Support JDBC 4 or Not
registerOutParameter(int parameterIndex, int type)	void	Yes
wasNull()	boolean	Yes
getString(int parameterIndex)	String	Yes
getBoolean(int parameterIndex)	boolean	Yes
getByte(int parameterIndex)	byte	Yes
getShort(int parameterIndex)	short	Yes
getInt(int parameterIndex)	int	Yes
getLong(int parameterIndex)	long	Yes
getFloat(int parameterIndex)	float	Yes
getDouble(int parameterIndex)	double	Yes
getBigDecimal(int parameterIndex)	BigDecimal	Yes
getBytes(int parameterIndex)	byte[]	Yes
getDate(int parameterIndex)	Date	Yes
getTime(int parameterIndex)	Time	Yes
getTimestamp(int parameterIndex)	Timestamp	Yes
getObject(int parameterIndex)	Object	Yes

Ⅲ NOTE

- Do not perform batch operations on statements containing **OUT** parameters.
- The following methods are inherited from java.sql.Statement: close, execute, executeQuery, executeUpdate, getConnection, getResultSet, getUpdateCount, isClosed, setMaxRows, and setFetchSize.
- The following methods are inherited from java.sql.PreparedStatement: addBatch, clearParameters, execute, executeQuery, executeUpdate, getMetaData, setBigDecimal, setBoolean, setByte, setBytes, setDate, setDouble, setFloat, setInt, setLong, setNull, setObject, setString, setTime, and setTimestamp.

java.sql.DatabaseMetaData

This section describes **java.sql.DatabaseMetaData**, the interface for defining database objects.

Table 11-9 java.sql.DatabaseMetaData methods

Method	Return Type	Support JDBC 4 or Not
getTables(String catalog, String schemaPattern, String tableNamePattern, String[] types)	ResultSet	Yes
getColumns(String catalog, String schemaPattern, String tableNamePattern, String columnNamePattern)	ResultSet	Yes
getTableTypes()	ResultSet	Yes
getUserName()	String	Yes
isReadOnly()	boolean	Yes
nullsAreSortedHigh()	boolean	Yes
nullsAreSortedLow()	boolean	Yes
nullsAreSortedAtStart()	boolean	Yes
nullsAreSortedAtEnd()	boolean	Yes
getDatabaseProductName()	String	Yes
getDatabaseProductVer- sion()	String	Yes
getDriverName()	String	Yes
getDriverVersion()	String	Yes
getDriverMajorVersion()	int	Yes
getDriverMinorVersion()	int	Yes
usesLocalFiles()	boolean	Yes
usesLocalFilePerTable()	boolean	Yes
supportsMixedCaseIdentifi- ers()	boolean	Yes
storesUpperCaseIdentifiers()	boolean	Yes
storesLowerCaseIdentifiers()	boolean	Yes
supportsMixedCaseQuotedI- dentifiers()	boolean	Yes

Method	Return Type	Support JDBC 4 or Not
storesUpperCaseQuotedI- dentifiers()	boolean	Yes
storesLowerCaseQuotedl- dentifiers()	boolean	Yes
storesMixedCaseQuotedI- dentifiers()	boolean	Yes
supportsAlterTableWithAdd- Column()	boolean	Yes
supportsAlterTableWith- DropColumn()	boolean	Yes
supportsColumnAliasing()	boolean	Yes
nullPlusNonNullIsNull()	boolean	Yes
supportsConvert()	boolean	Yes
supportsConvert(int fromType, int toType)	boolean	Yes
supportsTableCorrelation- Names()	boolean	Yes
supportsDifferentTableCorre- lationNames()	boolean	Yes
supportsExpressionsInOrder- By()	boolean	Yes
supportsOrderByUnrelated()	boolean	Yes
supportsGroupBy()	boolean	Yes
supportsGroupByUnrelated()	boolean	Yes
supportsGroupByBeyondSe- lect()	boolean	Yes
supportsLikeEscapeClause()	boolean	Yes
supportsMultipleResultSets()	boolean	Yes
supportsMultipleTransactions()	boolean	Yes
supportsNonNullableCol- umns()	boolean	Yes
supportsMinimumSQLGram- mar()	boolean	Yes
supportsCoreSQLGrammar()	boolean	Yes

Method	Return Type	Support JDBC 4 or Not
supportsExtendedSQLGram- mar()	boolean	Yes
supportsANSI92EntryLevelS QL()	boolean	Yes
supportsANSI92Intermediate SQL()	boolean	Yes
supportsANSI92FullSQL()	boolean	Yes
supportsIntegrityEnhance- mentFacility()	boolean	Yes
supportsOuterJoins()	boolean	Yes
supportsFullOuterJoins()	boolean	Yes
supportsLimitedOuterJoins()	boolean	Yes
isCatalogAtStart()	boolean	Yes
supportsSchemasInDataMa- nipulation()	boolean	Yes
supportsSavepoints()	boolean	Yes
supportsResultSetHoldabili- ty(int holdability)	boolean	Yes
getResultSetHoldability()	int	Yes
getDatabaseMajorVersion()	int	Yes
getDatabaseMinorVersion()	int	Yes
getJDBCMajorVersion()	int	Yes
getJDBCMinorVersion()	int	Yes

java.sql.Driver

This section describes **java.sql.Driver**, the database driver interface.

Table 11-10 java.sql.Driver methods

Method	Return Type	Support JDBC 4 or Not
acceptsURL(String url)	boolean	Yes
connect(String url, Properties info)	Connection	Yes
jdbcCompliant()	boolean	Yes

Method	Return Type	Support JDBC 4 or Not
getMajorVersion()	int	Yes
getMinorVersion()	int	Yes

java.sql.PreparedStatement

This section describes **java.sql.PreparedStatement**, the interface for preparing statements.

Table 11-11 java.sql.PreparedStatement methods

Method	Return Type	Support JDBC 4 or Not
clearParameters()	void	Yes
execute()	boolean	Yes
executeQuery()	ResultSet	Yes
executeUpdate()	int	Yes
getMetaData()	ResultSetMetaData	Yes
setBoolean(int parameterIndex, boolean x)	void	Yes
setBigDecimal(int parameterIndex, BigDecimal x)	void	Yes
setByte(int parameterIndex, byte x)	void	Yes
setBytes(int parameterIndex, byte[] x)	void	Yes
setDate(int parameterIndex, Date x)	void	Yes
setDouble(int parameterIndex, double x)	void	Yes
setFloat(int parameterIndex, float x)	void	Yes
setInt(int parameterIndex, int x)	void	Yes

Method	Return Type	Support JDBC 4 or Not
setLong(int parameterIndex, long x)	void	Yes
setNString(int parameterIndex, String value)	void	Yes
setShort(int parameterIndex, short x)	void	Yes
setString(int parameterIndex, String x)	void	Yes
addBatch()	void	Yes
executeBatch()	int[]	Yes
clearBatch()	void	Yes

□ NOTE

- addBatch() and execute() can be executed only after clearBatch().
- Calling the **executeBatch()** method does not clear the batch. Clear batch by explicitly calling **clearBatch()**.
- You do not need to use set*() to reuse the values of bounded variables in a batch after they have been added.
- The following methods are inherited from java.sql.Statement: close, execute, executeQuery, executeUpdate, getConnection, getResultSet, getUpdateCount, isClosed, setMaxRows, and setFetchSize.

java.sql.ResultSet

This section describes **java.sql.ResultSet**, the interface for execution result sets.

Table 11-12 java.sql.ResultSet methods

Method	Return Type	Support JDBC 4 or Not
findColumn(String columnLabel)	int	Yes
getBigDecimal(int columnIndex)	BigDecimal	Yes
getBigDecimal(String columnLabel)	BigDecimal	Yes
getBoolean(int columnIndex)	boolean	Yes

Method	Return Type	Support JDBC 4 or Not
getBoolean(String columnLabel)	boolean	Yes
getByte(int columnIndex)	byte	Yes
getBytes(int columnIndex)	byte[]	Yes
getByte(String columnLabel)	byte	Yes
getBytes(String columnLabel)	byte[]	Yes
getDate(int columnIndex)	Date	Yes
getDate(String columnLabel)	Date	Yes
getDouble(int columnIndex)	double	Yes
getDouble(String columnLabel)	double	Yes
getFloat(int columnIndex)	float	Yes
getFloat(String columnLabel)	float	Yes
getInt(int columnIndex)	int	Yes
getInt(String columnLabel)	int	Yes
getLong(int columnIndex)	long	Yes
getLong(String columnLabel)	long	Yes
getShort(int columnIndex)	short	Yes
getShort(String columnLabel)	short	Yes
getString(int columnIndex)	String	Yes
getString(String columnLabel)	String	Yes
getTime(int columnIndex)	Time	Yes
getTime(String columnLabel)	Time	Yes

Method	Return Type	Support JDBC 4 or Not
getTimestamp(int columnIndex)	Timestamp	Yes
getTimestamp(String columnLabel)	Timestamp	Yes
isAfterLast()	boolean	Yes
isBeforeFirst()	boolean	Yes
isFirst()	boolean	Yes
next()	boolean	Yes

□ NOTE

- A statement cannot have multiple open result sets.
- The cursor used to traverse the result set cannot remain in the open state after being committed.

java. sql. Result Set Meta Data

This section describes **java.sql.ResultSetMetaData**, which provides details about ResultSet object information.

Table 11-13 java.sql.ResultSetMetaData methods

Method	Return Type	Support JDBC 4 or Not
getColumnCount()	int	Yes
getColumnName(int column)	String	Yes
getColumnType(int column)	int	Yes
getColumnTypeName(int column)	String	Yes

java.sql.Statement

This section describes **java.sql.Statement**, the interface for executing SQL statements.

Table 11-14 java.sql.Statement methods

Method	Return Type	Support JDBC 4 or Not
close()	void	Yes
execute(String sql)	boolean	Yes
executeQuery(String sql)	ResultSet	Yes
executeUpdate(String sql)	int	Yes
getConnection()	Connection	Yes
getResultSet()	ResultSet	Yes
getQueryTimeout()	int	Yes
getUpdateCount()	int	Yes
isClosed()	boolean	Yes
setQueryTimeout(int seconds)	void	Yes
setFetchSize(int rows)	void	Yes
cancel()	void	Yes

◯ NOTE

setFetchSize can reduce the memory occupied by the result set on the client. Result sets are packaged into cursors and segmented for processing, which will increase the communication traffic between the database and the client, affecting performance.

Database cursors are valid only within their transactions. Therefore, when setting **setFetchSize**, set **setAutoCommit** to **false** and commit transactions on the connection to flush service data to a database.

javax. sql. Connection Pool Data Source

This section describes **javax.sql.ConnectionPoolDataSource**, the interface for data source connection pools.

Table 11-15 javax.sql.ConnectionPoolDataSource methods

Method	Return Type	Support JDBC 4 or Not
getLoginTimeout()	int	Yes
getLogWriter()	PrintWriter	Yes
getPooledConnection()	PooledConnection	Yes

Method	Return Type	Support JDBC 4 or Not
getPooledConnec- tion(String user,String password)	PooledConnection	Yes
setLoginTimeout(int seconds)	void	Yes
setLogWriter(PrintWrit er out)	void	Yes

javax.sql.DataSource

This section describes **javax.sql.DataSource**, the interface for data sources.

Table 11-16 javax.sql.DataSource methods

Method	Return Type	Support JDBC 4 or Not
getConnection()	Connection	Yes
getConnection(String username,String password)	Connection	Yes
getLoginTimeout()	int	Yes
getLogWriter()	PrintWriter	Yes
setLoginTimeout(int seconds)	void	Yes
setLogWriter(PrintWriter out)	void	Yes

javax.sql.PooledConnection

This section describes **javax.sql.PooledConnection**, the connection interface created by a connection pool.

Table 11-17 javax.sql.PooledConnection methods

Method	Return Type	Support JDBC 4 or Not
addConnectionEventListener (ConnectionEventListener listener)	void	Yes
close()	void	Yes
getConnection()	Connection	Yes

Method	Return Type	Support JDBC 4 or Not
removeConnectionEventListener (ConnectionEventListener listener)	void	Yes
addStatementEventListener (StatementEventListener listener)	void	Yes
removeStatementEventListener (StatementEventListener listener)	void	Yes

javax.naming.Context

This section describes **javax.naming.Context**, the context interface for connection configuration.

Table 11-18 javax.naming.Context methods

Method	Return Type	Support JDBC 4 or Not
bind(Name name, Object obj)	void	Yes
bind(String name, Object obj)	void	Yes
lookup(Name name)	Object	Yes
lookup(String name)	Object	Yes
rebind(Name name, Object obj)	void	Yes
rebind(String name, Object obj)	void	Yes
rename(Name oldName, Name newName)	void	Yes
rename(String oldName, String newName)	void	Yes
unbind(Name name)	void	Yes
unbind(String name)	void	Yes

javax.naming.spi.InitialContextFactory

This section describes **javax.naming.spi.InitialContextFactory**, the initial context factory interface.

Table 11-19 javax.naming.spi.InitialContextFactory methods

Method	Return Type	Support JDBC 4 or Not
getInitialContext(Hashtable ,? environment)	Context	Yes

CopyManager

CopyManager is an API interface class provided by the JDBC driver in GaussDB(DWS). It is used to import data to GaussDB(DWS) in batches.

Inheritance relationship of CopyManager

The CopyManager class is in the **org.postgresql.copy** package class and inherits the java.lang.Object class. The declaration of the class is as follows:

public class CopyManager extends Object

Construction method

public CopyManager(BaseConnection connection) throws SQLException

Common methods

Table 11-20 Common methods of CopyManager

Return ed Value	Method	Description	throws
Copyln	copyIn(String sql)	-	SQLException
long	copyIn(String sql, InputStream from)	Uses COPY FROM STDIN to quickly load data to tables in the database from InputStream.	SQLException,IOE xception
long	copyIn(String sql, InputStream from, int bufferSize)	Uses COPY FROM STDIN to quickly load data to tables in the database from InputStream.	SQLException,IOE xception

Return ed Value	Method	Description	throws
long	copyIn(String sql, Reader from)	Uses COPY FROM STDIN to quickly load data to tables in the database from Reader.	SQLException,IOE xception
long	copyIn(String sql, Reader from, int bufferSize)	Uses COPY FROM STDIN to quickly load data to tables in the database from Reader.	SQLException,IOE xception
CopyOu t	copyOut(String sql)	-	SQLException
long	copyOut(String sql, OutputStream to)	Sends the result set of COPY TO STDOUT from the database to the OutputStream class.	SQLException,IOE xception
long	copyOut(String sql, Writer to)	Sends the result set of COPY TO STDOUT from the database to the Writer class.	SQLException,IOE xception

11.3 ODBC-Based Development

Open Database Connectivity (ODBC) is a Microsoft API for accessing databases based on the X/OPEN CLI. The ODBC API alleviates applications from directly operating in databases, and enhances the database portability, extensibility, and maintainability.

Figure 11-2 shows the system structure of ODBC.

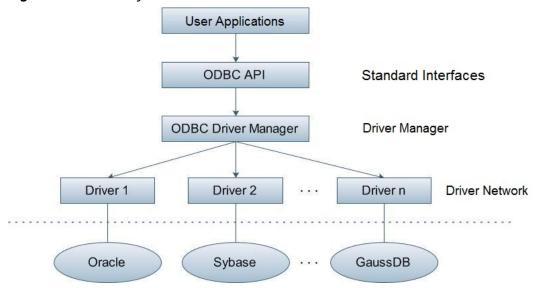


Figure 11-2 ODBC system structure

GaussDB(DWS) supports ODBC 3.5 in the following environments.

Table 11-21 OSs Supported by ODBC

OS	Platform
SUSE Linux Enterprise Server 11 SP1/SP2/SP3/SP4 SUSE Linux Enterprise Server 12 and SP1/SP2/SP3/SP5	x86_64
Red Hat Enterprise Linux 6.4/6.5/6.6/6.7/6.8/6.9/7.0/7.1/7.2/7.3/7.4/7.5	x86_64
Red Hat Enterprise Linux 7.5	ARM64
CentOS 6.4/6.5/6.6/6.7/6.8/6.9/7.0/7.1/7.2/7.3/7.4	x86_64
CentOS 7.6	ARM64
EulerOS 2.0 SP2/SP3	x86_64
EulerOS 2.0 SP8	ARM64
NeoKylin 7.5/7.6	ARM64
Oracle Linux R7U4	x86_64
Windows 7	32-bit
Windows 7	64-bit
Windows Server 2008	32-bit
Windows Server 2008	64-bit

The operating systems listed above refer to the operating systems on which the ODBC program runs. They can be different from the operating systems where databases are deployed.

The ODBC Driver Manager running on UNIX or Linux can be unixODBC or iODBC. Select unixODBC-2.3.0 here as the component for connecting the database.

Windows has a native ODBC Driver Manager. You can locate **Data Sources** (ODBC) by choosing **Control Panel > Administrative Tools**.

■ NOTE

The current database ODBC driver is based on an open source version and may be incompatible with GaussDB(DWS) data types, such as tinyint, smalldatetime, and nvarchar2.

11.3.1 ODBC Package and Its Dependent Libraries and Header Files

Download the ODBC software package from the console.

For details, see **Downloading the JDBC or ODBC Driver**.

ODBC Package for the Linux OS

Obtain the **dws_8.***x.x_***odbc_driver_for_***xxx_xxx_***zip** package from the software package. In the Linux OS, header files (including **sql.h** and **sqlext.h**) and library (**libodbc.so**) are required in application development. These header files and libraries can be obtained from the unixODBC-2.3.0 installation package.

ODBC Package for the Windows OS

Obtain the **dws_8**.*x*.*x*_**odbc**_**driver**_**for**_**windows.zip** package from the software package. In the Windows OS, the required header files and library files are system-resident.

11.3.2 Configuring a Data Source in the Linux OS

The ODBC DRIVER (psqlodbcw.so) provided by GaussDB(DWS) can be used after it has been configured in the data source. To set up a data source, configure the **odbc.ini** and **odbcinst.ini** files on the server. The two files are generated during the unixODBC compilation and installation, and are saved in the **etc** directory by default.

Procedure

Step 1 Obtain the source code package of unixODBC at:

https://sourceforge.net/projects/unixodbc/files/unixODBC/2.3.0/unixODBC-2.3.0.tar.gz/download

Step 2 Currently, unixODBC-2.2.1 is not supported. Assume you are to install unixODBC-2.3.0. Run the following commands. When installing, you can use -- prefix=[your_path] to specify the installation directory. The data source file will be created in [your_path]/etc, and the library file will be generated in [your_path]/lib.

tar zxvf unixODBC-2.3.0.tar.gz cd unixODBC-2.3.0

Open the configure file. If it does not exist, open the configure.ac file. Find LIB_VERSION.

Change the value of LIB_VERSION to 1:0:0 to compile a *.so.1 dynamic library with the same dependency on psqlodbcw.so.

vim configure

./configure --enable-gui=no --prefix=[your_path] # To perform the compilation on a TaiShan server, add the **configure** parameter **--build=aarch64-unknown-linux-gnu**. make

make install

Install unixODBC. If another version of unixODBC has been installed, it will be overwritten after installation.

Step 3 Replace the GaussDB(DWS) client driver.

Decompress dws_8.x.x_odbc_driver_for_xxx_xxx.zip to obtain the psqlodbcw.la and psqlodbcw.so files in the /dws_8.x.x_odbc_driver_for_xxx_xxx/odbc/lib directory.

Step 4 Configure the data source.

1. Configure the ODBC driver file.

Add the following content to the [your_path]/etc/odbcinst.ini file:

[GaussMPP]
Driver64=[your_path]/lib/odbc/psqlodbcw.so
setup=[your_path]/lib/odbc/psqlodbcw.so

For descriptions of the parameters in the odbcinst.ini file, see Table 11-22.

Table 11-22 odbcinst.ini configuration parameters

Parameter	Description	Examples
[DriverName]	Driver name, corresponding to Driver in DSN.	[DRIVER_N]
Driver64	Path of the dynamic driver library	Driver64=/xxx/odbc/lib/ odbc/psqlodbcw.so
setup	Driver installation path, which is the same as the dynamic library path in Driver64.	setup=/xxx/odbc/lib/ odbc/psqlodbcw.so

2. Configure the data source file.

Add the following content to the [your_path]/etc/odbc.ini file:

[MPPODBC]
Driver=GaussMPP
Servername=10.10.0.13 (database server IP address)
Database=gaussdb (database name)
Username=dbadmin (database username)
Password= (database user password)
Port=8000 (database listening port)
Sslmode=allow

For descriptions of the parameters in the odbc.ini file, see Table 11-23.

Table 11-23 odbc.ini configuration parameters

Parameter	Description	Examples
[DSN]	Data source name	[MPPODBC]
Driver	Driver name, corresponding to DriverName in odbcinst.ini Driver=DRIVER_N	
Servername	IP address of the server	Servername=10.145.130. 26
Database	Name of the database to connect to	Database=gaussdb
Username	Name of the database user	Username=dbadmin
Password	Password of the database user	Password= NOTE After a user established a connection, the ODBC driver automatically clears their password stored in memory. However, if this parameter is configured, UnixODBC will cache data source files, which may cause the password to be stored in the memory for a long time. When you connect to an application, you are advised to send your password through an API instead of writing it in a data source configuration file. After the connection has been established, immediately clear the memory segment where your password is stored.
Port	Port ID of the server	Port=8000
Sslmode	Whether to enable the SSL	Sslmode=allow
UseServerSidePre- pare	Whether to enable the extended query protocol for the database. The value can be 0 or 1 . The default value is 1 , indicating	UseServerSidePrepare=1
	that the extended query protocol is enabled.	

Parameter	Description	Examples
UseBatchProtocol	Whether to enable the batch query protocol. If it is enabled, the DML performance can be improved. The value can be 0 or 1. The default value is 1. If this parameter is set to 0, the batch query protocol is disabled (mainly for communication with earlier	UseBatchProtocol=1
	database versions). If this parameter is set to 1 and the support_batch_bind parameter is set to on, the batch query protocol is enabled.	
ConnectionExtral	Whether to display the driver deployment path and process owner in the connection_info parameter mentioned in connection_info	ConnectionExtraInfo=1 NOTE The default value is 1. If this parameter is set to 0, the ODBC driver reports the name and version of the current driver to the database. If this parameter is set to 1, the ODBC driver reports the name, deployment path, and process owner of the current driver to the database and records them in the connection_info parameter (see connection_info). You can query this parameter in PG_STAT_ACTIVITY and PGXC_STAT_ACTIVITY.

Parameter	Description	Examples	
ForExtensionCon- nector	ETL tool performance optimization parameter. It can be used to optimize the memory and reduce the memory usage by the peer CN, to avoid system instability caused by excessive CN memory usage.	ForExtensionConnector=1	
	The value can be 0 or 1 . The default value is 0 , indicating that the optimization item is disabled.		
	Do not set this parameter for other services outside the database system. Otherwise, the service correctness may be affected.		
KeepDisallowPre- mature	Specifies whether the cursor in the SQL statement has the with hold attribute when the following conditions are met: UseDeclareFetch is set to 1, and the application invokes SQLNumResultCols, SQLDescribeCol, or SQLColAttribute after invoking SQLPrepare to obtain the column information of the result set.	KeepDisallowPremature=1 NOTE When UseServerSidePrepare is set to 1, the KeepDisallowPremature parameter does not take effect. To use this parameter, set UseServerSidePrepare to 0. For example, set UseDeclareFetch to 1. KeepDisallowPremature=1 UseServerSidePrepare=0	
	The value can be 0 or 1 . 0 indicates that the with hold attribute is supported, and 1 indicates that the with hold attribute is not supported. The default value is 0 .		

The valid values of **sslmode** are as follows.

Table 11-24 sslmode options

sslmode	Whether SSL Encryption Is Enabled	Description
disable	No	The SSL secure connection is not used.
allow	Probably	The SSL secure encrypted connection is used if required by the database server, but does not check the authenticity of the server.
prefer	Probably	The SSL secure encrypted connection is used as a preferred mode if supported by the database, but does not check the authenticity of the server.
require	Yes	The SSL secure connection must be used, but it only encrypts data and does not check the authenticity of the server.
verify-ca	Yes	The SSL secure connection must be used, and it checks whether the database has certificates issued by a trusted CA.
verify- full	Yes	The SSL secure connection must be used. In addition to the check scope specified by verify-ca , it checks whether the name of the host where the database resides is the same as that on the certificate. This mode is not supported.

Step 5 Enable the SSL mode.

To use SSL certificates for connection, decompress the certificate package contained in the GaussDB(DWS) installation package, and run **source sslcert_env.sh** in a shell environment to deploy certificates in the default location of the current session.

Or manually declare the following environment variables and ensure that the permission for the client.key* series files is set to 600.

export PGSSLCERT= "/YOUR/PATH/OF/client.crt" # Change the path to the absolute path of client.crt. export PGSSLKEY= "/YOUR/PATH/OF/client.key" # Change the path to the absolute path of client.key.

In addition, change the value of **Sslmode** in the data source to **verify-ca**.

- **Step 6** Add the IP address segment of the host where the client is located to the security group rules of GaussDB(DWS) to ensure that the host can communicate with GaussDB(DWS).
- **Step 7** Configure environment variables.

vim ~/.bashrc

Add the following content to the end of the configuration file:

export LD_LIBRARY_PATH=[your_path]/lib/:\$LD_LIBRARY_PATH export ODBCSYSINI=[your_path]/etc export ODBCINI=[your_path]/etc/odbc.ini

It is not recommended to add **LD_LIBRARY_PATH** in the Kylin OS, as it may cause conflicts with the **libssl.so** dynamic library. In the latest version of cluster 9.1.0, the rpath mechanism has been added, so the dependency can be located without **LD_LIBRARY_PATH**.

Step 8 Run the following commands to validate the settings:

source ~/.bashrc

----End

Testing Data Source Configuration

Run the **isql**-v GaussODBC command (GaussODBC is the data source name).

• If the following information is displayed, the configuration is correct and the connection succeeds.



• If error information is displayed, the configuration is incorrect. Check the configuration.

Troubleshooting

• [UnixODBC][Driver Manager]Can't open lib 'xxx/xxx/psqlodbcw.so' : file not found.

Possible causes:

- The path configured in the odbcinst.ini file is incorrect.
 Run ls to check the path in the error information, ensuring that the psqlodbcw.so file exists and you have execution permissions on it.
- The dependent library of **psqlodbcw.so** does not exist or is not in system environment variables.

Run **ldd** to check the path in the error information. If **libodbc.so.1** or other UnixODBC libraries are lacking, configure UnixODBC again following the procedure provided in this section, and add the **lib** directory under its installation directory to **LD_LIBRARY_PATH**. If other libraries are lacking, add the **lib** directory under the ODBC driver package to **LD_LIBRARY_PATH**. Alternatively, you can place the dependency library of **psqlodbcw.so** in the path corresponding to rpath of **psqlodbcw.so**. To view rpath, you can use the **readelf** -**d** command.

- [UnixODBC]connect to server failed: no such file or directory
 Possible causes:
 - An incorrect or unreachable database IP address or port was configured.
 Check the Servername and Port configuration items in data sources.
 - Server monitoring is improper.

If **Servername** and **Port** are correctly configured, ensure the proper network adapter and port are monitored based on database server configurations in the procedure in this section.

- Firewall and network gatekeeper settings are improper.

Check firewall settings, ensuring that the database communication port is trusted.

Check to ensure network gatekeeper settings are proper (if any).

• [unixODBC]The password-stored method is not supported.

Possible causes:

The **sslmode** configuration item is not configured in the data sources.

Solution:

Set it to allow or a higher level. For more details, see Table 11-24.

Server common name "xxxx" does not match host name "xxxxx"

Possible causes:

When **verify-full** is used for SSL encryption, the driver checks whether the host name in certificates is the same as the actual one.

Solution:

To solve this problem, use **verify-ca** to stop checking host names, or generate a set of CA certificates containing the actual host names.

Driver's SQLAllocHandle on SQL_HANDLE_DBC failed

Possible causes:

The executable file (such as the **isql** tool of unixODBC) and the database driver (**psqlodbcw.so**) depend on different library versions of ODBC, such as **libodbc.so.1** and **libodbc.so.2**. You can verify this problem by using the following method:

ldd `which isql` | grep odbc ldd psqlodbcw.so | grep odbc

If the suffix digits of the outputs **libodbc.so** are different or indicate different physical disk files, this problem exists. Both **isql** and **psqlodbcw.so** load **libodbc.so**. If different physical files are loaded, different ODBC libraries with the same function list conflict with each other in a visible domain. As a result, the database driver cannot be loaded.

Solution:

Uninstall the unnecessary unixODBC, such as libodbc.so.2, and create a soft link with the same name and the .so.2 suffix for the remaining libodbc.so.1 library.

• FATAL: Forbid remote connection with trust method!

For security purposes, the CN forbids access from other nodes in the cluster without authentication.

To access the CN from inside the cluster, deploy the ODBC program on the machine where the CN is located and use 127.0.0.1 as the server address. It is recommended that the service system be deployed outside the cluster. If it is deployed inside, the database performance may be affected.

• [unixODBC][Driver Manager]Invalid attribute value

This problem occurs when you use SQL on other GaussDB. The possible cause is that the unixODBC version is not the recommended one. You are advised to run the **odbcinst** --version command to check the unixODBC version.

• authentication method 10 not supported.

If this error occurs on an open source client, the cause may be: The database stores only the SHA-256 hash of the password, but the open source client supports only MD5 hashes.

- The database stores the hashes of user passwords instead of actual passwords.
- In versions earlier than V100R002C80SPC300, the database stores only SHA-256 hashes and no MD5 hashes. Therefore, MD5 cannot be used for user password authentication.
- In V100R002C80SPC300 and later, if a password is updated or a user is created, both types of hashes will be stored, compatible with open-source authentication protocols.
- An MD5 hash can only be generated using the original password, but the password cannot be obtained by reversing its SHA-256 hash. If your database is upgraded from a version earlier than V100R002C80SPC300, passwords in the old version will only have SHA-256 hashes and not support MD5 authentication.

To solve this problem, you can update the user password. Alternatively, create a user, assign the same permissions to the user, and use the new user to connect to the database.

• unsupported frontend protocol 3.51: server supports 1.0 to 3.0

The database version is too early or the database is an open-source database.

Use the driver of the required version to connect to the database.

11.3.3 Configuring a Data Source in the Windows OS

Configure the ODBC data source using the ODBC data source manager preinstalled in the Windows OS.

Procedure

Step 1 Replace the GaussDB(DWS) client driver.

Decompress **GaussDB-9.1.0-Windows-Odbc.tar.gz**. Double-click install **psqlodbc.msi** (for 32-bit OS) or **psqlodbc_x64.msi** (for 64-bit OS).

Step 2 Open Driver Manager.

Use the Driver Manager suitable for your OS to configure the data source. (Assume the Windows system drive is drive C.)

• If you develop 32-bit programs in the 64-bit Windows OS, open the 32-bit Driver Manager at C:\Windows\SysWOW64\odbcad32.exe after you install the 32-bit driver.

Do not open Driver Manager by choosing **Control Panel**, clicking **Administrative Tools**, and clicking **Data Sources (ODBC)**.

Ⅲ NOTE

WoW64 is the acronym for "Windows 32-bit on Windows 64-bit". **C:\Windows \SysWOW64** stores the 32-bit environment on a 64-bit system.

• If you develop 64-bit programs in the 64-bit Windows OS, open the 64-bit Driver Manager at **C:\Windows\System32\odbcad32.exe** after you install the 64-bit driver.

Do not open **Driver Manager** by choosing **Control Panel**, clicking **Administrative Tools**, and clicking **Data Sources (ODBC)**.

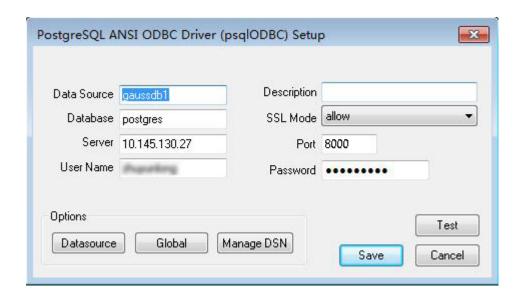
□ NOTE

C:\Windows\System32 stores the environment consistent with the current OS. For technical details, see Windows technical documents.

In a 32-bit Windows OS, open C:\Windows\System32\odbcad32.exe.
 In the Windows OS, click Computer, and choose Control Panel. Click Administrative Tools and click Data Sources (ODBC).

Step 3 Configure the data source.

On the **User DSN** tab, click **Add**, and choose **PostgreSQL Unicode** for setup. (An identifier will be displayed for the 64-bit OS.)



NOTICE

The entered username and password will be recorded in the Windows registry and you do not need to enter them again when connecting to the database next time. For security purposes, you are advised to delete sensitive information before clicking **Save** and enter the required username and password again when using ODBC APIs to connect to the database.

Step 4 Enable the SSL mode.

To use SSL certificates for connection, decompress the certificate package contained in the GaussDB(DWS) installation package, and double-click the **sslcert env.bat** file to deploy certificates in the default location.

NOTICE

The **sslcert_env.bat** file ensures the purity of the certificate environment. When the **%APPDATA%\postgresql** directory exists, a message will be prompted asking you whether you want to remove related directories. If you want to remove related directories, back up files in the directory.

Alternatively, you can copy the client.crt, client.key, client.key.cipher, and client.key.rand files in the certificate file folder to the manually created %APPDATA%\postgresql directory. Change client in the file names to postgres, for example, change client.key to postgres.key. Copy the cacert.pem file to the %APPDATA%\postgresql directory and change its name to root.crt.

Change the value of **SSL Mode** in step 2 to **verify-ca**.

Table 11-25 sslmode options

sslmode	Whether SSL Encryption Is Enabled	Description
disable	No	The SSL secure connection is not used.
allow	Probably	The SSL secure encrypted connection is used if required by the database server, but does not check the authenticity of the server.
prefer	Probably	The SSL secure encrypted connection is used as a preferred mode if supported by the database, but does not check the authenticity of the server.
require	Yes	The SSL secure connection must be used, but it only encrypts data and does not check the authenticity of the server.
verify-ca	Yes	The SSL secure connection must be used, and it checks whether the database has certificates issued by a trusted CA.
verify-full	Yes	The SSL secure connection must be used. In addition to the check scope specified by verify-ca , it checks whether the name of the host where the database resides is the same as that on the certificate. NOTE This mode cannot be used.

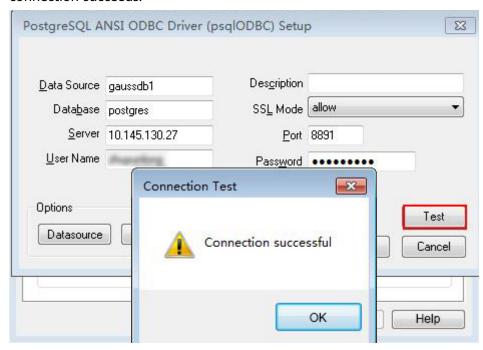
Step 5 Add the IP address segment of the host where the client is located to the security group rules of GaussDB(DWS) to ensure that the host can communicate with GaussDB(DWS).

----End

Testing Data Source Configuration

Click Test.

• If the following information is displayed, the configuration is correct and the connection succeeds.



• If error information is displayed, the configuration is incorrect. Check the configuration.

Troubleshooting

- Server common name "xxxx" does not match host name "xxxxx"

 This problem occurs because when **verify-full** is used for SSL encryption, the driver checks whether the host name in certificates is the same as the actual one. To solve this problem, use **verify-ca** to stop checking host names, or generate a set of CA certificates containing the actual host names.
- connect to server failed: no such file or directory
 Possible causes:
 - An incorrect or unreachable database IP address or port was configured.
 Check the **Servername** and **Port** configuration items in data sources.
 - Server monitoring is improper.
 If Servername and Port are correctly configured, ensure the proper network adapter and port are monitored based on database server configurations in the procedure in this section.
 - Firewall and network gatekeeper settings are improper.
 Check firewall settings, ensuring that the database communication port is trusted.
 - Check to ensure network gatekeeper settings are proper (if any).
- In the specified DSN, the system structures of the drive do not match those of the application.

Possible cause: The bit versions of the drive and program are different.

C:\Windows\SysWOW64\odbcad32.exe is a 32-bit ODBC Drive Manager.

C:\Windows\System32\odbcad32.exe is a 64-bit ODBC Drive Manager.

The password-stored method is not supported.

Possible causes:

sslmode is not configured for the data source. Set this configuration item to **allow** or a higher level to enable SSL connections. For details about **sslmode**, see **Table 11-25**.

authentication method 10 not supported.

If this error occurs on an open source client, the cause may be:

The database stores only the SHA-256 hash of the password, but the open source client supports only MD5 hashes.

□ NOTE

- The database stores the hashes of user passwords instead of actual passwords.
- In versions earlier than V100R002C80SPC300, the database stores only SHA-256 hashes and no MD5 hashes. Therefore, MD5 cannot be used for user password authentication.
- In V100R002C80SPC300 and later, if a password is updated or a user is created, both types of hashes will be stored, compatible with open-source authentication protocols.
- An MD5 hash can only be generated using the original password, but the password cannot be obtained by reversing its SHA-256 hash. If your database is upgraded from a version earlier than V100R002C80SPC300, passwords in the old version will only have SHA-256 hashes and not support MD5 authentication.

To solve this problem, perform the following operations:

- a. Check and change the password_encryption_type to 1 if needed with the help of technical support.
- b. Create a new database user for connection or reset the password of the existing database user.
 - If you use an administrator account, reset the password. For details, see Resetting a Password.
 - If you are a common user, use another client tool (such as Data Studio) to connect to the database and run the ALTER USER statement to change your password.
- c. Connect to the database.
- unsupported frontend protocol 3.51: server supports 1.0 to 3.0
 - The database version is too early or the database is an open-source database. Use the driver of the required version to connect to the database.
- FATAL: GSS authentication method is not allowed because XXXX user password is not disabled.

In some cases, the error is: GSSAPI authentication not supported.

In **pg_hba.conf** of the target CN, the authentication mode is set to **gss** for authenticating the IP address of the current client. However, this authentication algorithm cannot authenticate clients. Change the authentication algorithm to **sha256** and try again.

Note that cross-node connection to the database in the cluster is not supported. If the error is caused by cross-node connection to the CN in the cluster, connect the service program to the database from a node outside the cluster and try again.

11.3.4 ODBC Development Example

Code for Common Functions

The following example shows how to obtain data from GaussDB(DWS) through ODBC.

```
// DBtest.c (compile with: libodbc.so)
#include <stdlib.h>
#include <stdio.h>
#include <sqlext.h>
#ifdef WIN32
#include <windows.h>
#endif
SQLHENV
             V_OD_Env;
                             // Handle ODBC environment
SQLHSTMT
              V OD hstmt;
                              // Handle statement
SQLHDBC
             V_OD_hdbc;
                             // Handle connection
char
          typename[100];
SQLINTEGER value = 100;
SQLINTEGER V_OD_erg,V_OD_buffer,V_OD_err,V_OD_id;
            V_StrLen_or_IndPtr;
SOLLEN
int main(int argc,char *argv[])
    // 1. Apply for an environment handle.
   V_OD_erg = SQLAllocHandle(SQL_HANDLE_ENV,SQL_NULL_HANDLE,&V_OD_Env);
   if ((V_OD_erg != SQL_SUCCESS) && (V_OD_erg != SQL_SUCCESS_WITH_INFO))
       printf("Error AllocHandle\n");
       exit(0);
   // 2. Set environment attributes (version information)
   SQLSetEnvAttr(V_OD_Env, SQL_ATTR_ODBC_VERSION, (void*)SQL_OV_ODBC3, 0);
   // 3. Apply for a connection handle.
   V_OD_erg = SQLAllocHandle(SQL_HANDLE_DBC, V_OD_Env, &V_OD_hdbc);
   if ((V_OD_erg != SQL_SUCCESS) && (V_OD_erg != SQL_SUCCESS_WITH_INFO))
       SQLFreeHandle(SQL_HANDLE_ENV, V_OD_Env);
       exit(0);
   // 4. Set connection attributes.
   SQLSetConnectAttr(V_OD_hdbc, SQL_ATTR_AUTOCOMMIT, SQL_AUTOCOMMIT_ON, 0);
// 5. Connect to the data source. userName and password indicate the username and password for
connecting to the database. Set them as needed.
// If the username and password have been set in the odbc.ini file, you do not need to set userName or
password here, retaining "" for them. However, you are not advised to do so because the username and
password will be disclosed if the permission for odbc.ini is abused.
   V_OD_erg = SQLConnect(V_OD_hdbc, (SQLCHAR*) "gaussdb", SQL_NTS, (SQLCHAR*) "userName", SQL_NTS, (SQLCHAR*) "password", SQL_NTS);
   if ((V_OD_erg != SQL_SUCCESS) && (V_OD_erg != SQL_SUCCESS_WITH_INFO))
      printf("Error SQLConnect %d\n",V_OD_erg);
      SQLFreeHandle(SQL_HANDLE_ENV, V_OD_Env);
      exit(0);
    printf("Connected !\n");
   // 6. Set statement attributes
   SQLSetStmtAttr(V_OD_hstmt,SQL_ATTR_QUERY_TIMEOUT,(SQLPOINTER *)3,0);
    // 7. Apply for a statement handle
   SQLAllocHandle(SQL_HANDLE_STMT, V_OD_hdbc, &V_OD_hstmt);
    // 8. Executes an SQL statement directly
   SQLExecDirect(V_OD_hstmt,"drop table IF EXISTS customer_t1",SQL_NTS);
```

```
SQLExecDirect(V OD hstmt,"CREATE TABLE customer t1(c customer sk INTEGER, c customer name
VARCHAR(32));",SQL_NTS);
   SQLExecDirect(V_OD_hstmt,"insert into customer_t1 values(25,'li')",SQL_NTS);
   // 9. Prepare for execution
   SQLPrepare(V_OD_hstmt,"insert into customer_t1 values(?)",SQL_NTS);
   // 10. Bind parameters
   &value,0,NULL);
   // 11. Execute the ready statement
   SQLExecute(V_OD_hstmt);
   SQLExecDirect(V_OD_hstmt,"select id from testtable",SQL_NTS);
   // 12. Obtain the attributes of a certain column in the result set
SQLColAttribute(V_OD_hstmt,1,SQL_DESC_TYPE_NAME,typename,sizeof(typename),NULL,NULL);
   printf("SQLColAtrribute %s\n",typename);
   // 13. Bind the result set
   SQLBindCol(V_OD_hstmt,1,SQL_C_SLONG, (SQLPOINTER)&V_OD_buffer,150,
         (SQLLEN *)&V_StrLen_or_IndPtr);
   // 14. Collect data using SQLFetch
   V_OD_erg=SQLFetch(V_OD_hstmt);
   // 15. Obtain and return data using SQLGetData
   while(V_OD_erg != SQL_NO_DATA)
      SQLGetData(V\_OD\_hstmt,1,SQL\_C\_SLONG,(SQLPOINTER)\&V\_OD\_id,0,NULL);
      printf("SQLGetData ----ID = %d\n",V_OD_id);
      V OD_erg=SQLFetch(V_OD_hstmt);
   printf("Done !\n");
   // 16. Disconnect from the data source and release handles
   SQLFreeHandle(SQL_HANDLE_STMT,V_OD_hstmt);
   SQLDisconnect(V_OD_hdbc);
   SQLFreeHandle(SQL_HANDLE_DBC,V_OD_hdbc);
   SQLFreeHandle(SQL_HANDLE_ENV, V_OD_Env);
   return(0);
```

Code for Batch Processing

- Enable UseBatchProtocol in the data source and set support_batch_bind to on.
- Use **CHECK_ERROR** to check and print error information.
- This example is used to interactively obtain the DSN, data volume to be processed, and volume of ignored data from users, and insert required data into the test_odbc_batch_insert table.

```
#include <stdio.h>
#include <stdlib.h>
#include <sql.h>
#include <sqlext.h>
#include <string.h>
#include "util.c"
void Exec(SQLHDBC hdbc, SQLCHAR* sql)
  SQLRETURN retcode;
                                 // Return status
  SQLHSTMT hstmt = SQL_NULL_HSTMT; // Statement handle
  SQLCHAR loginfo[2048];
  // Allocate Statement Handle
  retcode = SQLAllocHandle(SQL_HANDLE_STMT, hdbc, &hstmt);
  CHECK_ERROR(retcode, "SQLAllocHandle(SQL_HANDLE_STMT)",
          hstmt, SQL_HANDLE_STMT);
  // Prepare Statement
  retcode = SQLPrepare(hstmt, (SQLCHAR*) sql, SQL_NTS);
```

```
sprintf((char*)loginfo, "SQLPrepare log: %s", (char*)sql);
  CHECK_ERROR(retcode, loginfo, hstmt, SQL_HANDLE_STMT);
  retcode = SQLExecute(hstmt);
  sprintf((char*)loginfo, "SQLExecute stmt log: %s", (char*)sql);
  CHECK_ERROR(retcode, loginfo, hstmt, SQL_HANDLE_STMT);
  retcode = SQLFreeHandle(SQL_HANDLE_STMT, hstmt);
  sprintf((char*)loginfo, "SQLFreeHandle stmt log: %s", (char*)sql);
  CHECK_ERROR(retcode, loginfo, hstmt, SQL_HANDLE_STMT);
int main ()
  SQLHENV henv = SQL_NULL_HENV;
  SQLHDBC hdbc = SQL_NULL_HDBC;
       batchCount = 1000;
  SQLLEN rowsCount = 0;
       ignoreCount = 0;
  int
  SQLRETURN retcode;
  SQLCHAR
             dsn[1024] = {'\0'};
  SQLCHAR
              loginfo[2048];
// Interactively obtain data source names.
  getStr("Please input your DSN", (char*)dsn, sizeof(dsn), 'N');
  Interactively obtain the amount of data to be batch processed.
  getInt("batchCount", &batchCount, 'N', 1);
  do
// Interactively obtain the amount of batch processing data that is not inserted into the database.
     getInt("ignoreCount", &ignoreCount, 'N', 1);
     if (ignoreCount > batchCount)
     {
       printf("ignoreCount(%d) should be less than batchCount(%d)\n", ignoreCount, batchCount);
  }while(ignoreCount > batchCount);
  retcode = SQLAllocHandle(SQL_HANDLE_ENV, SQL_NULL_HANDLE, &henv);
  CHECK_ERROR(retcode, "SQLAllocHandle(SQL_HANDLE_ENV)",
          henv, SQL_HANDLE_ENV);
  // Set ODBC Version
  retcode = SQLSetEnvAttr(henv, SQL_ATTR_ODBC_VERSION,
  (SQLPOINTER*)SQL_OV_ODBC3, 0);
CHECK_ERROR(retcode, "SQLSetEnvAttr(SQL_ATTR_ODBC_VERSION)",
          henv, SQL_HANDLE_ENV);
  // Allocate Connection
  retcode = SQLAllocHandle(SQL_HANDLE_DBC, henv, &hdbc);
  CHECK_ERROR(retcode, "SQLAllocHandle(SQL_HANDLE_DBC)",
          henv, SQL_HANDLE_DBC);
  // Set Login Timeout
  retcode = SQLSetConnectAttr(hdbc, SQL_LOGIN_TIMEOUT, (SQLPOINTER)5, 0);
  CHECK_ERROR(retcode, "SQLSetConnectAttr(SQL_LOGIN_TIMEOUT)",
          hdbc, SQL_HANDLE_DBC);
  // Set Auto Commit
  retcode = SQLSetConnectAttr(hdbc, SQL_ATTR_AUTOCOMMIT,
  (SQLPOINTER)(1), 0);
CHECK_ERROR(retcode, "SQLSetConnectAttr(SQL_ATTR_AUTOCOMMIT)",
          hdbc, SQL_HANDLE_DBC);
  // Connect to DSN
  sprintf(loginfo, "SQLConnect(DSN:%s)", dsn);
  retcode = SQLConnect(hdbc, (SQLCHAR*) dsn, SQL_NTS,
                    (SQLCHAR*) NULL, 0, NULL, 0);
  CHECK_ERROR(retcode, loginfo, hdbc, SQL_HANDLE_DBC);
```

```
// init table info.
  Exec(hdbc, "drop table if exists test_odbc_batch_insert");
  Exec(hdbc, "create table test_odbc_batch_insert(id int primary key, col varchar2(50))");
// The following code constructs the data to be inserted based on the data volume entered by users:
     SQLRETURN retcode;
     SQLHSTMT hstmtinesrt = SQL_NULL_HSTMT;
     SOLCHAR
                  *sql = NULL;
     SQLINTEGER *ids = NULL;
                  *cols = NULL;
     SQLCHAR
                 *bufLenIds = NULL;
     SQLLEN
     SQLLEN
                 *bufLenCols = NULL;
     SQLUSMALLINT *operptr = NULL;
     SQLUSMALLINT *statusptr = NULL;
     SQLULEN
                  process = 0;
// Data is constructed by column. Each column is stored continuously.
     ids = (SQLINTEGER*)malloc(sizeof(ids[0]) * batchCount);
     cols = (SQLCHAR*)malloc(sizeof(cols[0]) * batchCount * 50);
// Data size in each row for a column
     bufLenIds = (SQLLEN*)malloc(sizeof(bufLenIds[0]) * batchCount);
     bufLenCols = (SQLLEN*)malloc(sizeof(bufLenCols[0]) * batchCount);
// Whether this row needs to be processed. The value is SQL_PARAM_IGNORE or SQL_PARAM_PROCEED.
     operptr = (SQLUSMALLINT*)malloc(sizeof(operptr[0]) * batchCount);
     memset(operptr, 0, sizeof(operptr[0]) * batchCount);
// Processing result of the row
// Note: In the database, a statement belongs to one transaction. Therefore, data is processed as a unit.
That is, either all data is inserted successfully or all data fails to be inserted.
     statusptr = (SQLUSMALLINT*)malloc(sizeof(statusptr[0]) * batchCount);
     memset(statusptr, 88, sizeof(statusptr[0]) * batchCount);
     if (NULL == ids \parallel NULL == cols \parallel NULL == bufLenCols \parallel NULL == bufLenIds)
     {
       fprintf(stderr, "FAILED:\tmalloc data memory failed\n");
       goto exit;
     for (int i = 0; i < batchCount; i++)
       ids[i] = i;
       sprintf(cols + 50 * i, "column test value %d", i);
       bufLenIds[i] = sizeof(ids[i]);
       bufLenCols[i] = strlen(cols + 50 * i);
       operptr[i] = (i < ignoreCount) ? SQL\_PARAM\_IGNORE : SQL\_PARAM\_PROCEED;
     }
     // Allocate Statement Handle
     retcode = SQLAllocHandle(SQL_HANDLE_STMT, hdbc, &hstmtinesrt);
     CHECK_ERROR(retcode, "SQLAllocHandle(SQL_HANDLE_STMT)",
             hstmtinesrt, SQL_HANDLE_STMT);
     // Prepare Statement
     sql = (SQLCHAR*)"insert into test_odbc_batch_insert values(?, ?)";
     retcode = SQLPrepare(hstmtinesrt, (SQLCHAR*) sql, SQL_NTS);
     sprintf((char*)loginfo, "SQLPrepare log: %s", (char*)sql);
     CHECK_ERROR(retcode, loginfo, hstmtinesrt, SQL_HANDLE_STMT);
     retcode = SQLSetStmtAttr(hstmtinesrt, SQL_ATTR_PARAMSET_SIZE, (SQLPOINTER)batchCount,
sizeof(batchCount));
     CHECK_ERROR(retcode, "SQLSetStmtAttr", hstmtinesrt, SQL_HANDLE_STMT);
     retcode = SQLBindParameter(hstmtinesrt, 1, SQL_PARAM_INPUT, SQL_C_SLONG, SQL_INTEGER,
sizeof(ids[0]), 0,&(ids[0]), 0, bufLenIds);
     CHECK_ERROR(retcode, "SQLBindParameter for id", hstmtinesrt, SQL_HANDLE_STMT);
     retcode = SQLBindParameter(hstmtinesrt, 2, SQL_PARAM_INPUT, SQL_C_CHAR, SQL_CHAR, 50, 50,
```

```
cols, 50, bufLenCols);
     CHECK_ERROR(retcode, "SQLBindParameter for cols", hstmtinesrt, SQL_HANDLE_STMT);
    retcode = SQLSetStmtAttr(hstmtinesrt, SQL_ATTR_PARAMS_PROCESSED_PTR, (SQLPOINTER)&process,
sizeof(process));
    CHECK ERROR(retcode, "SQLSetStmtAttr for SQL ATTR PARAMS PROCESSED PTR", hstmtinesrt,
SQL_HANDLE_STMT);
     retcode = SQLSetStmtAttr(hstmtinesrt, SQL_ATTR_PARAM_STATUS_PTR, (SQLPOINTER)statusptr,
sizeof(statusptr[0]) * batchCount);
     CHECK ERROR(retcode, "SQLSetStmtAttr for SQL ATTR PARAM STATUS PTR", hstmtinesrt,
SQL HANDLE STMT);
    retcode = SQLSetStmtAttr(hstmtinesrt, SQL_ATTR_PARAM_OPERATION_PTR, (SQLPOINTER)operptr,
sizeof(operptr[0]) * batchCount);
    CHECK_ERROR(retcode, "SQLSetStmtAttr for SQL_ATTR_PARAM_OPERATION_PTR", hstmtinesrt,
SQL_HANDLE_STMT);
     retcode = SQLExecute(hstmtinesrt);
    sprintf((char*)loginfo, "SQLExecute stmt log: %s", (char*)sql);
     CHECK_ERROR(retcode, loginfo, hstmtinesrt, SQL_HANDLE_STMT);
     retcode = SQLRowCount(hstmtinesrt, &rowsCount);
     CHECK_ERROR(retcode, "SQLRowCount execution", hstmtinesrt, SQL_HANDLE_STMT);
    if (rowsCount != (batchCount - ignoreCount))
       sprintf(loginfo, "(batchCount - ignoreCount)(%d) != rowsCount(%d)", (batchCount - ignoreCount),
rowsCount);
       CHECK_ERROR(SQL_ERROR, loginfo, NULL, SQL_HANDLE_STMT);
    else
       sprintf(loginfo, "(batchCount - ignoreCount)(%d) == rowsCount(%d)", (batchCount - ignoreCount),
rowsCount):
       CHECK_ERROR(SQL_SUCCESS, loginfo, NULL, SQL_HANDLE_STMT);
    }
    if (rowsCount != process)
       sprintf(loginfo, "process(%d) != rowsCount(%d)", process, rowsCount);
       CHECK_ERROR(SQL_ERROR, loginfo, NULL, SQL_HANDLE_STMT);
    }
    else
       sprintf(loginfo, "process(%d) == rowsCount(%d)", process, rowsCount);
       CHECK_ERROR(SQL_SUCCESS, loginfo, NULL, SQL_HANDLE_STMT);
    }
     for (int i = 0; i < batchCount; i++)
    {
       if (i < ignoreCount)
       {
          if (statusptr[i] != SQL_PARAM_UNUSED)
            sprintf(loginfo, "statusptr[%d](%d) != SQL_PARAM_UNUSED", i, statusptr[i]);
            CHECK_ERROR(SQL_ERROR, loginfo, NULL, SQL_HANDLE_STMT);
       }
       else if (statusptr[i] != SQL_PARAM_SUCCESS)
          sprintf(loginfo, "statusptr[%d](%d) != SQL_PARAM_SUCCESS", i, statusptr[i]);
          CHECK_ERROR(SQL_ERROR, loginfo, NULL, SQL_HANDLE_STMT);
       }
    }
     retcode = SQLFreeHandle(SQL_HANDLE_STMT, hstmtinesrt);
     sprintf((char*)loginfo, "SQLFreeHandle hstmtinesrt");
     CHECK_ERROR(retcode, loginfo, hstmtinesrt, SQL_HANDLE_STMT);
```

```
exit:
    printf ("\nComplete.\n");

// Connection
    if (hdbc != SQL_NULL_HDBC) {
        SQLDisconnect(hdbc);
        SQLFreeHandle(SQL_HANDLE_DBC, hdbc);
}

// Environment
    if (henv != SQL_NULL_HENV)
        SQLFreeHandle(SQL_HANDLE_ENV, henv);

    return 0;
}
```

11.3.5 ODBC APIs

ODBC APIs are provided for users. This section covers common APIs. For more details on other APIs, see the ODBC Programmer's Reference on **the MSDN** website.

SQLAllocEnv

In ODBC 3.x, **SQLAllocEnv** (a function in ODBC 2.x) was deprecated and replaced by **SQLAllocHandle**. For details, see **SQLAllocHandle**.

SQLAllocConnect

In ODBC 3.x, **SQLAllocConnect** (a function in ODBC 2.x) was deprecated and replaced by **SQLAllocHandle**. For details, see **SQLAllocHandle**.

SQLAllocHandle

Function

SQLAllocHandle allocates environment, connection, or statement handles. It replaces the ODBC 2.*x* functions **SQLAllocEnv**, **SQLAllocConnect**, and **SQLAllocStmt**.

Prototype

```
SQLRETURN SQLAllocHandle(SQLSMALLINT HandleType,
SQLHANDLE InputHandle,
SQLHANDLE *OutputHandlePtr);
```

Table 11-26 SQLAllocHandle parameters

Parameter	Description
HandleType	Handle type allocated by SQLAllocHandle . The value must be one of the following:
	SQL_HANDLE_ENV (environment handle)
	SQL_HANDLE_DBC (connection handle)
	SQL_HANDLE_STMT (statement handle)
	SQL_HANDLE_DESC (description handle)
	The handle application sequence is: SQL_HANDLE_ENV > SQL_HANDLE_DBC > SQL_HANDLE_STMT. The handle applied later depends on the handle applied prior to it.
InputHandle	Type of the new handle to be allocated.
	 If HandleType is SQL_HANDLE_ENV, the value is SQL_NULL_HANDLE.
	If HandleType is SQL_HANDLE_DBC, this must be an environment handle.
	If HandleType is SQL_HANDLE_STMT or SQL_HANDLE_DESC, it must be a connection handle.
OutputHandlePt r	Output parameter: Pointer to a buffer in which the handle returned for the newly allocated data structure is stored.

- SQL_SUCCESS indicates that the call is successful.
- **SQL_SUCCESS_WITH_INFO** indicates warning information.
- SQL_ERROR indicates major errors, such as memory allocation and connection setup failures.
- SQL_INVALID_HANDLE indicates that invalid handles were called. Values
 returned by other APIs are similar to the values returned by the API you have
 used.

Precautions

If **SQLAllocHandle** returns **SQL_ERROR** when it is used to allocate a non-environment handle, it sets **OutputHandlePtr** to **SQL_NULL_HDBC**, **SQL_NULL_HSTMT**, or **SQL_NULL_HDESC**. The application can then call **SQLGetDiagRec**, set **HandleType** and **Handle** to **InputHandle**, and obtain the **SQLSTATE** value. This value can be used to get more information about the function call.

Examples

See ODBC Development Example.

SQLAllocStmt

In ODBC 3.x, **SQLAllocStmt** (a function in ODBC 2.x) was deprecated and replaced by **SQLAllocHandle**. For details, see **SQLAllocHandle**.

SQLBindCol

Function

SQLBindCol is used to associate (bind) columns in a result set to an application data buffer.

Prototype

```
SQLRETURN SQLBindCol(SQLHSTMT StatementHandle,
SQLUSMALLINT ColumnNumber,
SQLSMALLINT TargetType,
SQLPOINTER TargetValuePtr,
SQLLEN BufferLength,
SQLLEN *StrLen_or_IndPtr);
```

Parameters

Table 11-27 SQLBindCol parameters

Parameter	Description
StatementHandl e	Statement handle.
ColumnNumber	Number of the column to be bound. Column numbering begins at 0 and increases in ascending order. Column 0 functions as the bookmark. If no bookmark column is set, column numbering begins at 1 instead.
TargetType	The C data type in the buffer.
TargetValuePtr	Output parameter: pointer to the buffer bound with the column. The SQLFetch function returns data in the buffer. If TargetValuePtr is null, StrLen_or_IndPtr is a valid value.
BufferLength	Length of the buffer to which TargetValuePtr points, in bytes.
StrLen_or_IndPtr	Output parameter: pointer to the length or indicator of the buffer. If StrLen_or_IndPtr is null, no length or indicator is used.

Return values

- **SQL_SUCCESS** indicates that the call is successful.
- **SQL_SUCCESS_WITH_INFO** indicates warning information.
- **SQL_ERROR** indicates major errors, such as memory allocation and connection setup failures.

SQL_INVALID_HANDLE indicates that invalid handles were called. Values
returned by other APIs are similar to the values returned by the API you have
used.

Note

If **SQLBindCol** returns **SQL_ERROR** or **SQL_SUCCESS_WITH_INFO**, the application can then call **SQLGetDiagRec**, set **HandleType** and **Handle** to **SQL_HANDLE_STMT** and **StatementHandle**, and obtain the **SQLSTATE** value. This value can be used to get more information about the function call.

Examples

See **ODBC Development Example**.

SQLBindParameter

Function

SQLBindParameter binds a parameter flag in an SQL statement to a buffer.

Prototype

Table 11-28 SQLBindParameter

Keyword	Description
StatementHandle	Statement handle.
ParameterNumbe r	Parameter marker number, starting at 1 and increasing in an ascending order.
InputOutputType	Input and output parameter types.
ValueType	C data type of the parameter.
ParameterType	SQL data type of the parameter.
ColumnSize	Column size or the expression of the corresponding parameter marker.
DecimalDigits	Decimal number of the column or the expression of the corresponding parameter marker.
ParameterValuePt r	Pointer to the buffer for storing parameter data.

Keyword	Description
BufferLength	Length of the buffer to which the ParameterValuePtr points, in bytes.
StrLen_or_IndPtr	Pointer to the length or indicator of the buffer. If StrLen_or_IndPtr is null, no length or indicator is used.

- SQL SUCCESS indicates that the call is successful.
- **SQL_SUCCESS_WITH_INFO** indicates warning information.
- **SQL_ERROR** indicates major errors, such as memory allocation and connection setup failures.
- SQL_INVALID_HANDLE indicates that invalid handles were called. Values
 returned by other APIs are similar to the values returned by the API you have
 used.

Precautions

If **SQLBindCol** returns **SQL_ERROR** or **SQL_SUCCESS_WITH_INFO**, the application can then call **SQLGetDiagRec**, set **HandleType** and **Handle** to **SQL_HANDLE_STMT** and **StatementHandle**, and obtain the **SQLSTATE** value. This value can be used to get more information about the function call.

Examples

See ODBC Development Example.

SQLColAttribute

Function

SQLColAttribute returns the descriptor information about a column in the result set.

Prototype

SQLRETURN SQLColAttribute(SQLF	HSTMT StatementHandle,
SQLUSMALLINT	ColumnNumber,
SQLUSMALLINT	FieldIdentifier,
SQLPOINTER	CharacterAtrriburePtr,
SQLSMALLINT	BufferLength,
SQLSMALLINT	*StringLengthPtr,
SOI POINTER	NumericAttributePtr):

Table 11-29 SQLColAttribute parameters

Parameter	Description
StatementHandle	Statement handle.

Parameter	Description
ColumnNumber	Column number of the field to be queried, starting at 1 and increasing in an ascending order.
FieldIdentifier	Field identifier of ColumnNumber in IRD.
CharacterAttribu- tePtr	Output parameter: pointer to the buffer that returns FieldIdentifier field value.
BufferLength	FieldIdentifier indicates the buffer length when it refers to an ODBC-defined field and CharacterAttributePtr points to a string or binary buffer. Ignore this parameter if FieldIdentifier is an ODBC.
	 Ignore this parameter if FieldIdentifier is an ODBC- defined field and CharacterAttributePtr points to an integer.
StringLengthPtr	Output parameter: pointer to a buffer in which the total number of valid bytes (for string data) is stored in *CharacterAttributePtr. Ignore the value of BufferLength if the data is not a string.
NumericAttributePt r	Output parameter: pointer to an integer buffer in which the value of the FieldIdentifier field in the ColumnNumber row of the IRD is returned.

- **SQL SUCCESS** indicates that the call is successful.
- **SQL_SUCCESS_WITH_INFO** indicates warning information.
- **SQL_ERROR** indicates major errors, such as memory allocation and connection setup failures.
- SQL_INVALID_HANDLE indicates that invalid handles were called. Values
 returned by other APIs are similar to the values returned by the API you have
 used.

Precautions

If **SQLColAttribute** returns **SQL_ERROR** or **SQL_SUCCESS_WITH_INFO**, the application can then call **SQLGetDiagRec**, set **HandleType** and **Handle** to **SQL_HANDLE_STMT** and **StatementHandle**, and obtain the **SQLSTATE** value. This value can be used to get more information about the function call.

Examples

See ODBC Development Example.

SQLConnect

Function

SQLConnect establishes a connection between a driver and a data source. Using the connection handle, you can obtain crucial information like the program's

status, transaction processing status, and error messages after establishing a connection to the data source.

Prototype

Parameters

Table 11-30 SQLConnect parameters

Parameter	Description
ConnectionHandl e	Connection handle, obtained from SQLAllocHandle .
ServerName	Name of the data source to connect to.
NameLength1	Length of ServerName .
UserName	Database username in the data source.
NameLength2	Length of UserName .
Authentication	Password of the database user in the data source.
NameLength3	Length of Authentication .

Return values

- SQL_SUCCESS indicates that the call is successful.
- SQL_SUCCESS_WITH_INFO indicates warning information.
- **SQL_ERROR** indicates major errors, such as memory allocation and connection setup failures.
- **SQL_INVALID_HANDLE** indicates that invalid handles were called. Values returned by other APIs are similar to the values returned by the API you have used.
- **SQL_STILL_EXECUTING** indicates that the statement is being executed.

Precautions

If **SQLConnect** returns **SQL_ERROR** or **SQL_SUCCESS_WITH_INFO**, the application can then call **SQLGetDiagRec**, set **HandleType** and **Handle** to **SQL_HANDLE_DBC** and **ConnectionHandle**, and obtain the **SQLSTATE** value. This value can be used to get more information about the function call.

Examples

See ODBC Development Example.

SQLDisconnect

Function

SQLDisconnect closes the connection associated with the database connection handle.

Prototype

SQLRETURN SQLDisconnect(SQLHDBC ConnectionHandle);

Parameters

Table 11-31 SQLDisconnect parameters

Parameter	Description
ConnectionHandl e	Connection handle, obtained from SQLAllocHandle.

Return values

- SQL SUCCESS indicates that the call is successful.
- **SQL_SUCCESS_WITH_INFO** indicates warning information.
- **SQL_ERROR** indicates major errors, such as memory allocation and connection setup failures.
- SQL_INVALID_HANDLE indicates that invalid handles were called. Values
 returned by other APIs are similar to the values returned by the API you have
 used.

Precautions

If SQLDisconnect returns SQL_ERROR or SQL_SUCCESS_WITH_INFO, the application can then call SQLGetDiagRec, set HandleType and Handle to SQL_HANDLE_DBC and ConnectionHandle, and obtain the SQLSTATE value. This value can be used to get more information about the function call.

Examples

See **ODBC Development Example**.

SQLExecDirect

Function

SQLExecDirect executes a prepared SQL statement specified in this parameter. This is the fastest execution method for executing only one SQL statement at a time.

Prototype

SQLRETURN SQLExecDirect(SQLHSTMT StatementHandle,
SQLCHAR *StatementText,
SQLINTEGER TextLength);

Table 11-32 SQLExecDirect parameters

Parameter	Description
StatementHandl e	Statement handle, obtained from SQLAllocHandle .
StatementText	SQL statement to be executed. One SQL statement can be executed at a time.
TextLength	Length of StatementText .

- SQL SUCCESS indicates that the call is successful.
- **SQL_SUCCESS_WITH_INFO** indicates warning information.
- **SQL_NEED_DATA** indicates that there are not enough parameters provided to execute the SQL statement.
- **SQL_ERROR** indicates major errors, such as memory allocation and connection setup failures.
- SQL_INVALID_HANDLE indicates that invalid handles were called. Values
 returned by other APIs are similar to the values returned by the API you have
 used.
- **SQL_STILL_EXECUTING** indicates that the statement is being executed.
- **SQL_NO_DATA** indicates that no result set is returned for the SQL statement.

Precautions

If **SQLExecDirect** returns **SQL_ERROR** or **SQL_SUCCESS_WITH_INFO**, the application can then call **SQLGetDiagRec**, set **HandleType** and **Handle** to **SQL_HANDLE_STMT** and **StatementHandle**, and obtain the **SQLSTATE** value. This value can be used to get more information about the function call.

Examples

See **ODBC Development Example**.

SQLExecute

Function

When a statement includes a parameter marker, the **SQLExecute** function executes a prepared SQL statement using the current value of the marker.

Prototype

SQLRETURN SQLExecute(SQLHSTMT StatementHandle);

Table 11-33 SQLExecute parameters

Parameter	Description
StatementHandl e	Statement handle to be executed.

- SQL_SUCCESS indicates that the call is successful.
- **SQL_SUCCESS_WITH_INFO** indicates warning information.
- **SQL_NEED_DATA** indicates that there are not enough parameters provided to execute the SQL statement.
- **SQL_ERROR** indicates major errors, such as memory allocation and connection setup failures.
- **SQL NO DATA** indicates that no result set is returned for the SQL statement.
- **SQL_INVALID_HANDLE** indicates that invalid handles were called. Values returned by other APIs are similar to the values returned by the API you have used.
- **SQL_STILL_EXECUTING** indicates that the statement is being executed.

Precautions

If **SQLExecute** returns **SQL_ERROR** or **SQL_SUCCESS_WITH_INFO**, the application can then call **SQLGetDiagRec**, set **HandleType** and **Handle** to **SQL_HANDLE_STMT** and **StatementHandle**, and obtain the **SQLSTATE** value. This value can be used to get more information about the function call.

Examples

See ODBC Development Example.

SQLFetch

Function

SQLFetch advances the cursor to the next row of the result set and retrieves any bound columns.

Prototype

SQLRETURN SQLFetch(SQLHSTMT StatementHandle);

Table 11-34 SQLFetch parameters

Parameter	Description
StatementHandl e	Statement handle, obtained from SQLAllocHandle .

- SQL_SUCCESS indicates that the call is successful.
- **SQL_SUCCESS_WITH_INFO** indicates warning information.
- **SQL_ERROR** indicates major errors, such as memory allocation and connection setup failures.
- SQL NO DATA indicates that no result set is returned for the SQL statement.
- **SQL_INVALID_HANDLE** indicates that invalid handles were called. Values returned by other APIs are similar to the values returned by the API you have used.
- **SQL_STILL_EXECUTING** indicates that the statement is being executed.

Precautions

If **SQLFetch** returns **SQL_ERROR** or **SQL_SUCCESS_WITH_INFO**, the application can then call **SQLGetDiagRec**, set **HandleType** and **Handle** to **SQL_HANDLE_STMT** and **StatementHandle**, and obtain the **SQLSTATE** value. This value can be used to get more information about the function call.

Examples

See ODBC Development Example.

SQLFreeStmt

In ODBC 3.x, **SQLFreeStmt** (a function in ODBC 2.x) was deprecated and replaced with **SQLFreeHandle**. For details, see **SQLFreeHandle**.

SQLFreeConnect

In ODBC 3.x, **SQLFreeConnect** (a function in ODBC 2.x) was deprecated and replaced with **SQLFreeHandle**. For details, see **SQLFreeHandle**.

SQLFreeHandle

Function

SQLFreeHandle releases resources associated with a specific environment, connection, or statement handle. It replaces the ODBC 2.x functions: **SQLFreeEnv**, **SQLFreeConnect**, and **SQLFreeStmt**.

Prototype

SQLRETURN SQLFreeHandle(SQLSMALLINT HandleType, SQLHANDLE Handle);

Table 11-35 SQLFreeHandle parameters

Parameter	Description		
HandleType	Type of handle to be freed by SQLFreeHandle . The value must be one of the following:		
	SQL_HANDLE_ENV		
	SQL_HANDLE_DBC		
	SQL_HANDLE_STMT		
	SQL_HANDLE_DESC		
	If HandleType is not one of these values, SQLFreeHandle returns SQL_INVALID_HANDLE.		
Handle	Handle to be released.		

- SQL_SUCCESS indicates that the call is successful.
- **SQL_SUCCESS_WITH_INFO** indicates warning information.
- SQL_ERROR indicates major errors, such as memory allocation and connection setup failures.
- SQL_INVALID_HANDLE indicates that invalid handles were called. Values
 returned by other APIs are similar to the values returned by the API you have
 used.

Precautions

If **SQLFreeHandle** returns **SQL_ERROR**, the handle is still valid.

Examples

See ODBC Development Example.

SQLFreeEnv

In ODBC 3.x, **SQLFreeEnv** (a function in ODBC 2.x) was deprecated and replaced with **SQLFreeHandle**. For details, see **SQLFreeHandle**.

SQLPrepare

Function

SQLPrepare prepares an SQL statement to be executed.

Prototype

SQLRETURN SQLPrepare(SQLHSTMT StatementHandle,
SQLCHAR *StatementText,
SQLINTEGER TextLength);

Table 11-36 SQLPrepare parameters

Parameter	Description
StatementHandl e	Statement handle.
StatementText	SQL text string.
TextLength	Length of StatementText .

- SQL_SUCCESS indicates that the call is successful.
- **SQL_SUCCESS_WITH_INFO** indicates warning information.
- **SQL_ERROR** indicates major errors, such as memory allocation and connection setup failures.
- SQL_INVALID_HANDLE indicates that invalid handles were called. Values
 returned by other APIs are similar to the values returned by the API you have
 used.
- **SQL_STILL_EXECUTING** indicates that the statement is being executed.

Precautions

If **SQLPrepare** returns **SQL_ERROR** or **SQL_SUCCESS_WITH_INFO**, the application can then call **SQLGetDiagRec**, set **HandleType** and **Handle** to **SQL_HANDLE_STMT** and **StatementHandle**, and obtain the **SQLSTATE** value. This value can be used to get more information about the function call.

Examples

See ODBC Development Example.

SQLGetData

Function

SQLGetData retrieves data for a single column in the current row of the result set. It can be called multiple times to retrieve data of variable lengths.

Prototype

Table 11-37 SQLGetData parameters

Parameter	Description			
StatementHandle	Statement handle, obtained from SQLAllocHandle .			
Col_or_Param_Nu m	Column number of the data to be returned. The column in the result set are numbered from 1 in ascending orde The number of the bookmark column is 0.			
TargetType	Type identifier of the C data type in the TargetValuePtr buffer. If TargetType is SQL_ARD_TYPE , the driver uses the data type of the SQL_DESC_CONCISE_TYPE field in ARD. If TargetType is SQL_C_DEFAULT , the driver selects a default data type according to the source SQL data type.			
TargetValuePtr	Output parameter: pointer to the pointer that points to the buffer where the data is located.			
BufferLength	Size of the buffer pointed to by TargetValuePtr.			
StrLen_or_IndPtr	Output parameter : pointer to the buffer where the length or identifier value is returned.			

- SQL SUCCESS indicates that the call is successful.
- SQL SUCCESS WITH INFO indicates warning information.
- **SQL_ERROR** indicates major errors, such as memory allocation and connection setup failures.
- **SQL_NO_DATA** indicates that no result set is returned for the SQL statement.
- SQL_INVALID_HANDLE indicates that invalid handles were called. Values
 returned by other APIs are similar to the values returned by the API you have
 used.
- **SQL_STILL_EXECUTING** indicates that the statement is being executed.

Precautions

If **SQLFetch** returns **SQL_ERROR** or **SQL_SUCCESS_WITH_INFO**, the application can then call **SQLGetDiagRec**, set **HandleType** and **Handle** to **SQL_HANDLE_STMT** and **StatementHandle**, and obtain the **SQLSTATE** value. This value can be used to get more information about the function call.

Examples

See ODBC Development Example.

SQLGetDiagRec

Function

SQLGetDiagRec returns the current values of multiple fields of a diagnostic record that contains error, warning, and status information.

Prototype



Parameters

Table 11-38 SQLGetDiagRec parameters

Parameter	Description			
HandleType	Handle type identifier that describes the handle type required for diagnosis. The value must be one of the following:			
	SQL_HANDLE_ENV			
	SQL_HANDLE_DBC			
	SQL_HANDLE_STMT			
	SQL_HANDLE_DESC			
Handle	Handle of the diagnosis data structure. Its type is indicated by HandleType. If HandleType is SQL_HANDLE_ENV , Handle may be shared or non-shared environment handle.			
RecNumber	Status record from which the application seeks information. Status records are numbered from 1.			
SQLState	Output parameter: pointer to a buffer that saves the 5-character SQLSTATE code pertaining to RecNumber.			
NativeErrorPt r	Output parameter: pointer to a buffer that saves the native error code.			
MessageText	Pointer to a buffer that saves text strings of diagnostic information.			
BufferLength	Length of MessageText .			
TextLengthPt r	Output parameter: pointer to the buffer, the total number of bytes in the returned MessageText. If the number of bytes available to return is greater than BufferLength, then the diagnostics information text in MessageText is truncated to BufferLength minus the length of the null termination character.			

Return values

- **SQL_SUCCESS** indicates that the call is successful.
- **SQL_SUCCESS_WITH_INFO** indicates warning information.
- **SQL_ERROR** indicates major errors, such as memory allocation and connection setup failures.

SQL_INVALID_HANDLE indicates that invalid handles were called. Values
returned by other APIs are similar to the values returned by the API you have
used.

Precautions

SQLGetDiagRec does not release diagnostic records for itself. It uses the following returned values to report execution results:

- **SQL_SUCCESS**: The function successfully returns diagnostic information.
- **SQL_SUCCESS_WITH_INFO**: The *MessageText buffer is too small to hold the requested diagnostic message and no diagnostic records are generated.
- SQL_INVALID_HANDLE: The handle specified by **HandType** and **Handle** is invalid.
- **SQL_ERROR**: **RecNumber** is smaller than or equal to zero, or **BufferLength** is smaller than zero.

If an ODBC function returns **SQL_ERROR** or **SQL_SUCCESS_WITH_INFO**, the application can then call **SQLGetDiagRec** and obtain the **SQLSTATE** value. The possible **SQLSTATE** values are listed as follows:

Table 11-39 SQLSTATE values

SQLSTATE	Error	Description		
HY000	General error	An error occurred for which there is no specific SQLSTATE .		
HY001	Memory allocation error	The driver is unable to allocate memory required to support execution or completion of the function.		
HY008	Operation canceled	SQLCancel is called to terminate the statement execution, but the StatementHandle function is still called.		
HY010	Function sequence error	The function is called prior to sending data to data parameters or columns being executed.		
HY013	Memory management error	The function fails to be called. The error may be caused by low memory conditions.		
HYT01	Connection timeout	The connection times out before the data source responds to the request.		
IM001	Function not supported by the driver	A function that is not supported by the StatementHandle driver is called.		

Examples

See ODBC Development Example.

SQLSetConnectAttr

Function

SQLSetConnectAttr sets connection attributes.

Prototype

```
SQLRETURN SQLSetConnectAttr(SQLHDBC ConnectionHandle
SQLINTEGER Attribute,
SQLPOINTER ValuePtr,
SQLINTEGER StringLength);
```

Parameters

Table 11-40 SQLSetConnectAttr parameters

Parameter	Description		
ConnectionHand le	Connection handle.		
Attribute	Attribute to set.		
ValuePtr	Pointer to the value of Attribute . ValuePtr depends on the value of Attribute and can be a 32-bit unsigned integer value or a null-terminated string. If ValuePtr parameter is driver-specific value, it may be signed integer.		
StringLength	If ValuePtr points to a string or a binary buffer, this parameter should be the length of *ValuePtr. If ValuePtr points to an integer, StringLength is ignored.		

Return values

- SQL_SUCCESS indicates that the call is successful.
- **SQL_SUCCESS_WITH_INFO** indicates warning information.
- SQL_ERROR indicates major errors, such as memory allocation and connection setup failures.
- SQL_INVALID_HANDLE indicates that invalid handles were called. Values
 returned by other APIs are similar to the values returned by the API you have
 used.

Precautions

If SQLSetConnectAttr returns SQL_ERROR or SQL_SUCCESS_WITH_INFO, the application can then call SQLGetDiagRec, set HandleType and Handle to SQL_HANDLE_DBC and ConnectionHandle, and obtain the SQLSTATE value. This value can be used to get more information about the function call.

Examples

See ODBC Development Example.

SQLSetEnvAttr

Function

SQLSetEnvAttr sets environment attributes.

Prototype

```
SQLRETURN SQLSetEnvAttr(SQLHENV EnvironmentHandle
SQLINTEGER Attribute,
SQLPOINTER ValuePtr,
SQLINTEGER StringLength);
```

Parameters

Table 11-41 SQLSetEnvAttr parameters

Parameter	Description		
EnvironmentHan dle	Environment handle.		
Attribute	 Environment attribute to be set. Its value must be one of the following: SQL_ATTR_ODBC_VERSION: ODBC version SQL_CONNECTION_POOLING: connection pool attribute SQL_OUTPUT_NTS: string type returned by the driver 		
ValuePtr	Pointer to the value of Attribute . ValuePtr depends on the value of Attribute and can be a 32-bit integer value or a null-terminated string.		
StringLength	If ValuePtr points to a string or a binary buffer, this parameter should be the length of *ValuePtr. If ValuePtr points to an integer, StringLength is ignored.		

Return values

- SQL_SUCCESS indicates that the call is successful.
- **SQL_SUCCESS_WITH_INFO** indicates warning information.
- **SQL_ERROR** indicates major errors, such as memory allocation and connection setup failures.
- SQL_INVALID_HANDLE indicates that invalid handles were called. Values
 returned by other APIs are similar to the values returned by the API you have
 used.

Precautions

If SQLSetEnvAttr returns SQL_ERROR or SQL_SUCCESS_WITH_INFO, the application can then call SQLGetDiagRec, set HandleType and Handle to SQL_HANDLE_ENV and EnvironmentHandle, and obtain the SQLSTATE value. This value can be used to get more information about the function call.

Examples

See ODBC Development Example.

SQLSetStmtAttr

Function

SQLSetStmtAttr sets attributes related to a statement.

Prototype

```
SQLRETURN SQLSetStmtAttr(SQLHSTMT StatementHandle
SQLINTEGER Attribute,
SQLPOINTER ValuePtr,
SQLINTEGER StringLength);
```

Parameters

Table 11-42 SQLSetStmtAttr parameters

Parameter	Description	
StatementHandl e	Statement handle.	
Attribute	Attribute to set.	
ValuePtr	Pointer to the value of Attribute . ValuePtr depends on the value of Attribute and can be a 32-bit unsigned integer value or a pointer to a null-terminated string, a binary buffer, and a driver-specified value. If ValuePtr parameter is driver-specific value, it may be signed integer.	
StringLength	If ValuePtr points to a string or a binary buffer, this parameter should be the length of *ValuePtr. If ValuePtr points to an integer, StringLength is ignored.	

Return values

- SQL_SUCCESS indicates that the call is successful.
- **SQL_SUCCESS_WITH_INFO** indicates warning information.
- **SQL_ERROR** indicates major errors, such as memory allocation and connection setup failures.
- **SQL_INVALID_HANDLE** indicates that invalid handles were called. Values returned by other APIs are similar to the values returned by the API you have used.

Precautions

If SQLSetStmtAttr returns SQL_ERROR or SQL_SUCCESS_WITH_INFO, the application can then call SQLGetDiagRec, set HandleType and Handle to SQL_HANDLE_STMT and StatementHandle, and obtain the SQLSTATE value. This value can be used to get more information about the function call.

Examples

See ODBC Development Example.

12 GaussDB(DWS) Resource Monitoring

GaussDB(DWS) provides multiple dimensional resource monitoring views to show the real-time and historical resource usage of tasks.

12.1 User Resource Monitoring

In the multi-tenant management framework, you can query the real-time usage of all user resources (including the memory, number of CPU cores, storage space, temporary space, operator spilling space, and I/Os) in real time through the system views PG_TOTAL_USER_RESOURCE_INFO and PGXC_TOTAL_USER_RESOURCE_INFO and the function GS_WLM_USER_RESOURCE_INFO. You can also query the system catalog GS_WLM_USER_RESOURCE_HISTORY and system view PGXC_WLM_USER_RESOURCE_HISTORY for the historical usage of user resources.

Precautions

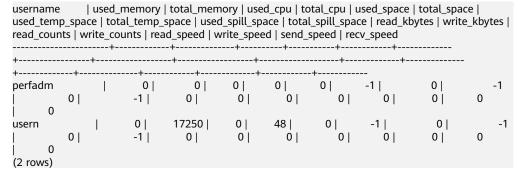
- The CPU, I/O, and memory usage of all jobs on fast and slow lanes (simple jobs on fast lanes and complex jobs on slow lanes) can be monitored.
- Currently, the memory and CPU usage of fast track jobs are not controlled.
 When the fast lane jobs occupy a large number of resources, the used resources may exceed the resource limit.
- In the DN monitoring view, I/O, memory, and CPU display the resource usage and limits of resource pools.
- In the CN monitoring view, I/O, memory, and CPU display the total resource usage and limit of all DN resource pools in the cluster.
- The DN monitoring information is updated every 5 seconds. CNs collect monitoring information from DNs every 5 seconds. Because each instance updates or collects user monitoring information independently, the monitoring information update time on each instance may be different.
- The auxiliary thread automatically invokes the persistence function every 30 seconds to persist user monitoring data. In normal cases, you do not need to do this.

- When there are a large number of users and a large cluster, querying such real-time views will cause network latency due to the real-time communication overhead between CNs and DNs.
- Resources are not monitored for an initial administrator.

Procedure

Query all users' resource quotas and real-time resource usage.
 SELECT * FROM PG_TOTAL_USER_RESOURCE_INFO;

The result view is as follows:



The I/O resource monitoring fields (read_kbytes, write_kbytes, read_counts, write_counts, read_speed, and write_speed) can be available only when the GUC parameter described in enable user metric persistent is enabled.

For details about each column, see **PG_TOTAL_USER_RESOURCE_INFO**.

Query a user's resource quota and real-time resource usage.
 SELECT * FROM GS_WLM_USER_RESOURCE_INFO('username');

The query result is as follows:

• Query all users' resource quotas and historical resource usage. SELECT * FROM GS_WLM_USER_RESOURCE_HISTORY;

The query result is as follows:

For the system catalog GS_WLM_USER_RESOURCE_HISTORY, data in the PG_TOTAL_USER_RESOURCE_INFO view is periodically saved to historical tables only when the GUC parameter enable_user_metric_persistent is enabled.

For details about each column, see GS WLM USER RESOURCE HISTORY.

12.2 Resource Pool Monitoring

Overview

In the multi-tenant management framework, if queries are associated with resource pools, the resources occupied by the queries are summarized to the associated resource pools. You can query the real-time resource usage of all resource pools in the resource pool monitoring view and query the historical resource usage of resource pools in the resource pool monitoring history table.

The resource pool monitoring data is updated every 5s. However, due to the time difference between CNs and DNs, the actual monitoring data update time may be longer than 5s. Generally, the time does not exceed 10s. The resource pool monitoring data is persisted every 30 seconds. The resource pool monitoring logic is basically the same as that of the user resource monitoring. Therefore, the <code>enable_user_metric_persistent</code> and <code>user_metric_retention_time</code> parameters are used to control the persistence and aging of resource pool monitoring data, respectively.

Resources monitored by a resource pool include the running and queuing information of fast and slow lane jobs, and CPU, memory, and logical I/O resource monitoring information. The monitoring views and history tables are as follows:

- Real-time monitoring view of resource pools (single CN):
 GS RESPOOL RUNTIME INFO
- Real-time monitoring view of resource pools (all CNs):
 PGXC_RESPOOL_RUNTIME_INFO
- Real-time monitoring view of resource pool resources (single CN):
 GS_RESPOOL_RESOURCE_INFO
- Real-time monitoring view of resource pool resources (all CNs):
 PGXC_RESPOOL_RESOURCE_INFO
- Historical resource monitoring table of the resource pool (single CN):
 GS_RESPOOL_RESOURCE_HISTORY
- Monitoring view of historical resource pool resources (all CNs): PGXC_RESPOOL_RESOURCE_HISTORY

□ NOTE

- Resource pool monitoring monitors the CPU, I/O, and memory usage of all jobs on the fast and slow lanes.
- Currently, the memory and CPU usage of fast track jobs are not controlled. When the
 fast lane jobs occupy a large number of resources, the used resources may exceed the
 resource limit.
- In the monitoring view of DN resource pools, I/O, memory, and CPU display the resource usage and limits of resource pools.
- In the monitoring view of CN resource pools, I/O, memory, and CPU display the total resource usage and limit of all DN resource pools in the cluster.
- Resource pool monitoring information on DNs is updated every 5 seconds. CNs collect resource pool monitoring information from DNs every 5 seconds. Because each instance updates or collects resource pool monitoring information independently, the monitoring information update time on each instance may be different.
- The auxiliary thread automatically invokes the persistence function every 30 seconds to persist the resource pool monitoring data. In normal cases, you do not need to do this.

Procedure

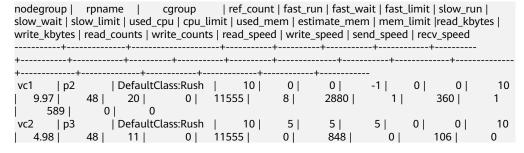
Querying the real-time running status of jobs in a resource pool.
 SELECT * FROM GS_RESPOOL_RUNTIME_INFO;

The result view is as follows:

Where,

- ref_count indicates the number of jobs that reference the current resource pool information. Its value will be retained until the management ends.
- fast_run and slow_run are load management accounting information.
 Their values are valid only when fast_limit and slow_limit are larger than 0.
- c. This view is valid only on CNs. The persistence information is stored in **GS RESPOOL RESOURCE HISTORY**.
- d. For details about each field, see GS RESPOOL RUNTIME INFO.
- Querying the resource quota and real-time resource usage of a resource pool.
 SELECT * FROM GS_RESPOOL_RESOURCE_INFO;

The result view is as follows:



	173		0										
vc2	p4	[PefaultCla	ass:Rush		0	0	0	-1	()	0	10
	0	48	0	0	11555	0)	0	0		0	0	
	0	0	0										
vc1	def	ault_pool	Defaul	tClass:M	edium	0		0	0	-1	0	0	
	-1	0	48	0	0	11555	5	0	0		0	0	
	0	0	0	0									
vc2	def	ault_pool	Defaul	tClass:M	edium	0		0	0	-1	0	0	
	-1	0	48	0	0	11555	5	0	0		0	0	
ĺ	0	0		0									
vc1	p1	[)efaultCla	ass:Rush	2	20	5	5	5	:	3	7	3
1 7	7.98	48	16	768	11555	5	8	265	6	1		32	1
ĺ	543	0	0										
(6 rc	ws)												

- a. This view is valid on both CNs and DNs. The CPU, memory, and I/O usage on a DN indicates the resource consumption of the DN. The CPU, memory, and I/O usage on a CN is the total resource consumption of all DNs in the cluster.
- b. **estimate_mem** is valid only on CNs under dynamic load management. It displays the estimated memory accounting of the resource pool.
- c. I/O monitoring information is recorded only when **enable_logical_io_statistics** is enabled.
- d. For details about each field, see GS_RESPOOL_RESOURCE_INFO.
- Querying the resource quota and historical resource usage of a resource pool.
 SELECT * FROM GS_RESPOOL_RESOURCE_HISTORY ORDER BY timestamp DESC;

The result view is as follows:

```
| nodegroup | rpname |
                                                     | ref_count | fast_run | fast_wait |
                                           cgroup
fast_limit | slow_run | slow_wait | slow_limit | used_cpu | cpu_limit | used_mem | estimate_mem |
mem_limit | read_kbytes | write_kbytes | read_counts | write_counts | read_speed | write_speed |
send_speed | recv_speed
2022-03-04 09:41:57.53739+08 | vc1
                                              | DefaultClass:Rush | 10 |
                                    | p2
                                9.97 |
                                          48 | 20 | 0 | 11555 |
                                                                           0 |
    -1 | 0 | 0 | 10 |
                                 474 |
2320 |
           0 |
                   290 |
                           0 |
                                              0 |
                                              | DefaultClass:Rush |
2022-03-04 09:41:57.53739+08 | vc1
                                                                     20 |
                                    | p1
    5 | 3 | 7 | 3 | 7.98 | 48 | 16 | 768 | 11555 |
| 0 | 237 | 0 | 387 | 0 | 0
                                                                           0 |
                                  387 |
1896 l
                                    | default_pool | DefaultClass:Medium |
2022-03-04 09:41:57.53739+08 | vc2
     0 | -1 | 0 | 0 |
                                  -1 | 0 | 48 | 0 | 0 | 11555 |
                       0 |
      0 1
              0|
                                0 |
                                        0 |
                                                0 |
                                                       0
2022-03-04 09:41:57.53739+08 | vc1
                                     | default_pool | DefaultClass:Medium |
     0 | -1 | 0 | 0 | 0 |
                                  -1 | 0 | 48 | 0 | 0 | 11555 |
                                                0 |
                                0 |
                                        0 |
                                                        n
2022-03-04 09:41:57.53739+08 | vc2
                                              | DefaultClass:Rush |
                                     | p4
                                  0 |
     -1 | 0 | 0 | 10 |
0 | 0 | 0 |
                                        48 | 0 |
                                                   0 | 11555 |
                                                                         0 |
                                                                                  n
                               0 |
                                      0 |
                                               0
2022-03-04 09:41:57.53739+08 | vc2
                                     | p3
                                              | DefaultClass:Rush |
                                                                     10 |
                                                                             5 I
                                                                                    5
     5 | 0 | 0 | 10 | 4.99 | 48 | 11 | 0 | 11555 | 0 | 110 | 0 | 180 | 0 | 0
                                                                           0 |
                                                                                  880
                               180 |
                                         0 |
2022-03-04 09:41:27.335234+08 | vc2
                                     | p3
                                               | DefaultClass:Rush | 10 |
    5 | 0 | 0 | 10 | 4.98 | 48 | 11 |
0 | 107 | 0 | 175 | 0 | 0
                                         48 | 11 |
                                                    0 | 11555 |
```

- a. The monitoring information comes from the resource pool monitoring history table. When **enable_user_metric_persistent** is enabled, the monitoring information is recorded every 30 seconds.
- b. The storage duration of the table data is specified by the **user_metric_retention_time** parameter.

For details about each field, see GS_RESPOOL_RESOURCE_HISTORY.

12.3 Monitoring Memory Resources

Monitoring the Memory

GaussDB(DWS) provides a view for monitoring the memory usage of the entire cluster.

Query the pgxc_total_memory_detail view as a user with sysadmin permissions. SELECT * FROM pgxc_total_memory_detail;

If the following error message is returned during the query, enable the memory management function.

SELECT * FROM pgxc_total_memory_detail; ERROR: unsupported view for memory protection feature is disabled. CONTEXT: PL/pgSQL function pgxc_total_memory_detail() line 12 at FOR over EXECUTE statement

You can set **enable_memory_limit** and **max_process_memory** on the GaussDB(DWS) console to enable memory management. The procedure is as follows:

- 1. Log in to the GaussDB(DWS) console.
- 2. In the navigation pane on the left, click Clusters.
- 3. In the cluster list, find the target cluster and click its name. The **Basic Information** page is displayed.
- Click the Parameter Modification tab, change the value of enable_memory_limit to on, and click Save to save the file.
- 5. Change the value of **max_process_memory** to a proper one. For details about the modification suggestions, see **max_process_memory**. After it is done, click **Save**.
- 6. In the **Modification Preview** dialog box, confirm the modifications and click **Save**. After the modification, restart the cluster for the modification to take effect.

Monitoring the Shared Memory

You can query the context information about the shared memory on the pg_shared_memory_detail view.

SELECT * FROM pg_shared_memory_detail;					
contextname	level	parent	totalsize freesize usedsize		
	++-		++++		
ProcessMemory	0		24576 9840 14736		
Workload manager memo	ry context	1 ProcessMe	emory 2105400 7304 20980)96	
wlm collector hash table	2	Workload manag	ger memory context 8192 3736 445	6	
Resource pool hash table				808	
wlm cgroup hash table	j 2 j	Workload manad	ger memory context 24576 15968 86	808	
(5 rows)		•			

This view lists the context name of the memory, level, the upper-layer memory context, and the total size of the shared memory.

In the database, GUC parameter **memory_tracking_mode** is used to configure the memory statistics collecting mode, including the following options:

- **none:** The memory statistics collecting function is not enabled.
- **normal:** Only memory statistics is collected in real time and no file is generated.
- **executor:** The statistics file is generated, containing the context information about all allocated memory used on the execution layer.

When the parameter is set to **executor**, csv files are generated under the **pg_log** directory of the DN process. The file names are in the format of **memory_track_**<*DN name>_query_*<*queryid>.csv*. The information about the operators executed by the postgres thread of the executor and all stream threads are input in this file during task execution.

The following is an example of the file content:

```
0, 0, ExecutorState, 0, PortalHeapMemory, 0, 40K, 602K, 23
1, 3, CStoreScan_29360131_25, 0, ExecutorState, 1, 265K, 554K, 23
2, 128, cstore scan per scan memory context, 1, CStoreScan_29360131_25, 2, 24K, 24K, 23
3, 127, cstore scan memory context, 1, CStoreScan_29360131_25, 2, 264K, 264K, 23
4, 7, InitPartitionMapTmpMemoryContext, 1, CStoreScan_29360131_25, 2, 31K, 31K, 23
5, 2, VecPartIterator_29360131_24, 0, ExecutorState, 1, 16K, 16K, 23
0, 0, ExecutorState, 0, PortalHeapMemory, 0, 24K, 1163K, 20
1, 3, CStoreScan_29360131_22, 0, ExecutorState, 1, 390K, 1122K, 20
2, 20, cstore scan per scan memory context, 1, CStoreScan_29360131_22, 2, 476K, 476K, 20
3, 19, cstore scan memory context, 1, CStoreScan_29360131_22, 2, 264K, 264K, 20
4, 7, InitPartitionMapTmpMemoryContext, 1, CStoreScan_29360131_22, 2, 23K, 23K, 20
5, 2, VecPartIterator_29360131_21, 0, ExecutorState, 1, 16K, 16K, 20
```

The fields include the output SN, SN of the memory allocation context within the thread, name of the current memory context, output SN of the parent memory context, name of the parent memory context, tree layer No. of the memory context, peak memory used by the current memory context, peak memory used by the current memory context and all its child memory contexts, and plan node ID of the query where the thread is executed.

In this example, the record "1, 3, CStoreScan_29360131_22, 0, ExecutorState, 1, 390K, 1122K, 20" represents the following information about Explain Analyze:

- **CstoreScan_29360131_22** indicates the CstoreScan operator.
- 1122K indicates the peak memory used by the CstoreScan operator.
- **fullexec:** The generated file includes the information about all memory contexts requested by the execution layer.

If the parameter is set to **fullexec**, the output information will be similar to that for **executor**, except that some memory context allocation information may be returned because the information about all memory applications (no matter succeeded or not) is printed. As only the memory application information is recorded, the peak memory used by the memory context is recorded as **0**.

12.4 Instance Resource Monitoring

GaussDB(DWS) provides system catalogs for monitoring the resource usage of CNs and DNs (including memory, CPU usage, disk I/O, process physical I/O, and process logical I/O), and system catalogs for monitoring the resource usage of the entire cluster.

For details about the system catalog **GS_WLM_INSTANCE_HISTORY**, see **GS_WLM_INSTANCE_HISTORY**.

□ NOTE

Data in the system catalog**GS_WLM_INSTANCE_HISTORY** is distributed in corresponding instances. CN monitoring data is stored in the CN instance, and DN monitoring data is stored in the DN instance. The DN has a standby node. When the primary DN is abnormal, the monitoring data of the DN can be restored from the standby node. However, a CN has no standby node. When a CN is abnormal and then restored, the monitoring data of the CN will be lost.

Procedure

Query the latest resource usage of the current instance.
 SELECT * FROM GS WLM INSTANCE HISTORY ORDER BY TIMESTAMP DESC;

The guery result is as follows:

Query the resource usage of the current instance during a specified period.
 SELECT * FROM GS_WLM_INSTANCE_HISTORY WHERE TIMESTAMP > '2022-01-10' AND TIMESTAMP < '2020-01-11' ORDER BY TIMESTAMP DESC;

The query result is as follows:

 To query the latest resource usage of a cluster, you can invoke the pgxc_get_wlm_current_instance_info stored procedure on the CN. SELECT * FROM pgxc_get_wlm_current_instance_info('ALL');

The guery result is as follows:

```
instancename |
                timestamp | used_cpu | free_mem | used_mem | io_await | io_util |
disk_read | disk_write | process_read | process_write | logical_read | logical_write | read_counts |
write counts
coordinator2 | 2020-01-14 21:58:29.290894+08 | 0 | 12010 | 278 | 16.0445 | 7.19561 |
184.431 | 27959.3 | 0 | 10 | coordinator3 | 2020-01-14 21:58:27.567655+08 |
                                         0 | 0 | 0 | 0 |
                                                          0 0
                                                          288 | .964557 | 3.40659 |
                                                0 |
332.468 | 3375.02 | 26 | 13 |
                                          0 |
                                                          0 | 0
                                         0 | 11899 |
                                                         389 | 1.17296 | 3.25 |
datanode1 | 2020-01-14 21:58:23.900321+08 |
329.6 | 2870.4 | 28 | 8 | 13 |
                                                 3 |
                                                          18 |
                                                               6
datanode2 | 2020-01-14 21:58:32.832989+08 | 0 | 11904 | 384 | 17.948 | 8.52148 |
                                               3 |
214.186 | 25894.1 | 28 | 10 |
                                          13 |
                                                          18 | 6
datanode3 | 2020-01-14 21:58:24.826694+08 | 0 | 11894 | 394 | 1.16088 | 3.15 |
                                                                                 328
```

```
2868.8 |
               25 | 10 |
                                                     18 |
coordinator1 | 2020-01-14 21:58:33.367649+08 |
                                                        300 | 9.53286 | 10.05 |
                                            0 | 11988 |
                           0 |
                                                        0 |
43.2 | 55232 |
                  0 |
                                                0 |
coordinator1 | 2020-01-14 21:58:23.216645+08 |
                                            0 | 11988 |
                                                          300 | 1.17085 | 3.21182 |
324.729 | 2831.13 |
                       8 |
                                13 |
                                           0 |
                                                    0 |
                                                            0 |
                                                                     0
(7 rows)
```

• To query the historical resource usage of a cluster, you can call the **pgxc_get_wlm_history_instance_info** stored procedure function on the CN. SELECT * FROM pgxc_get_wlm_history_instance_info('ALL', '2020-01-14 21:00:00', '2020-01-14 22:00:00', 3);

The query result is as follows:

```
| used_cpu | free_mem | used_mem | io_await | io_util |
instancename |
                   timestamp
disk_read | disk_write | process_read | process_write | logical_read | logical_write | read_counts |
write_counts
coordinator2 | 2020-01-14 21:50:49.778902+08 |
                                              0 | 12020 |
                                                            268 | .127371 | .789211 |
                                                     0 |
15.984 | 3994.41 | 0 | 0 |
                                           0 |
                                                             0 |
                                                                      0
coordinator2 | 2020-01-14 21:53:49.043646+08 |
                                                             270 | 30.2902 | 8.65404 |
                                              0 |
                                                   12018 |
276.77 | 16741.8 |
                   3 |
                                1 |
                                           0 |
                                                     0 |
                                                             0 |
                                                                      n
                                                  12018 |
coordinator2 | 2020-01-14 21:57:09.202654+08 |
                                              0 |
                                                            270 | .16051 | .979021 |
                                           0 |
59.9401 l
                      0.1
                                  0 |
           5596 l
                                                     0 |
                                                             0 |
                                                                      0
coordinator3 | 2020-01-14 21:38:48.948646+08 |
                                              0 |
                                                   12012 |
                                                             276 | .0769231 | .00999001
     0 | 35.1648 |
                        0 |
                                  1 |
                                            0 |
                                                     0 |
                                                             0 1
                                                                       0
coordinator3 | 2020-01-14 21:40:29.061178+08 |
                                              0 |
                                                  12012 |
                                                             276 | .118421 | .0199601
     0 | 970.858 |
                       0.1
                                  0.1
                                            0 |
                                                     0 |
                                                              0.1
                                                                    0
                                                  12010 |
coordinator3 | 2020-01-14 21:50:19.612777+08 |
                                                            278 |
                                              0 |
                                                                  24.411 | 11.7665 |
8.78244 | 44641.1 |
                        0 1
                                  0 [
                                                     0 [
                                                              0 [
                                                                       0
datanode1 | 2020-01-14 21:49:42.758649+08 |
                                                  11909 |
                                                            379 | .798776 |
                                              0 |
                                                                             8.02 |
51.2 | 20924.8 |
                   0 |
                              0 |
                                                  0 |
datanode1 | 2020-01-14 21:49:52.760188+08 |
                                              0 |
                                                  11909 |
                                                            379 | 23.8972 |
                                                                              14.1 I
0 | 74760 | 0 | 0 |
                                               0 | 0 |
                                                              0
datanode1 | 2020-01-14 21:50:22.769226+08 |
                                                  11909 |
                                                             379 | 39.5868 |
                                                   0 |
| 19760.8 |
                         0 |
               0 |
                                              0 |
                                  0 1
                                                                0
datanode2
          | 2020-01-14 21:58:02.826185+08 |
                                              0 |
                                                  11905 |
                                                             383 | .351648 |
                 0 |
                                                  0 |
20.8 | 504.8 |
                              0.1
                                                                   0
datanode2 | 2020-01-14 21:56:42.80793+08 |
                                              0 1
                                                  11906 |
                                                            382 | .559748 |
                                                                              .04
   326.4 |
               0 |
                          0 |
                                   0.1
                                              0
                                                     0 |
                                                               0
datanode2 | 2020-01-14 21:45:21.632407+08 |
                                              0 | 11901 |
                                                             387 | 12.1313 | 4.55544 |
3.1968 | 45177.2 |
                     0 |
                                                   0 1
                                                             0 |
datanode3 | 2020-01-14 21:58:14.823317+08 |
                                              0 | 11898 |
                                                            390 | .378205 |
                                                                              .99 [
                                                         0 |
48 | 23353.6 |
                                                 0 |
                0 |
                             0 |
                                                                   0
datanode3 | 2020-01-14 21:47:50.665028+08 |
                                              0 | 11901 |
                                                            387 | 1.07494 |
                                                                             1.19
0 | 15506.4 |
               0 |
                             0 |
                                                0 | 0 |
                                                                  Ω
datanode3 | 2020-01-14 21:51:21.720117+08 |
                                              0 |
                                                  11903 |
                                                            385 | 10.2795 |
0 | 11031.2 | 0 |
                            0 |
                                                0 | 0 |
                                                                  0
coordinator1 | 2020-01-14 21:42:59.121945+08 |
                                              0 | 12020 |
                                                             268 | .0857143 | .0699301
     0 | 6579.02 |
                      0 |
                                  0.1
                                                    0 |
                                                              0 i
                                                                       0
coordinator1 | 2020-01-14 21:41:49.042646+08 |
                                              0 | 12020 |
                                                             268 | 20.9039 |
                                                                            11.3786 |
6042.76 | 57903.7 | 0 | 0 |
                                                   0 |
                                                             0 1
                                                                      0
                                              0 | 12020 |
coordinator1 | 2020-01-14 21:41:09.007652+08 |
                                                             268 | .0446429 |
                                      0 |
0 | 1109.29 | 0 |
                         0 |
                                                0 | 0 |
(18 rows)
```

12.5 Real-time Top SQL

You can query real-time Top SQL in real-time resource monitoring views at different levels. The real-time resource monitoring view records the resource usage (including memory, data spilled to disks, and CPU time) and performance alarm information during job running.

The following table describes the external interfaces of the real-time views.

Table 12-1 Real-time resource monitoring views

Level	Monitored Node	View
Query level/perf	Current CN	GS_WLM_SESSION_STATISTICS
level	All CNs	PGXC_WLM_SESSION_STATISTICS
operator level Current CN		GS_WLM_OPERATOR_STATISTICS
	All CNs	PGXC_WLM_OPERATOR_STATISTICS

□ NOTE

- The view level is determined by the resource monitoring level, that is, the resource_track_level configuration.
- The perf and operator levels affect the values of the query_plan and warning columns in GS_WLM_SESSION_STATISTICS or PGXC_WLM_SESSION_INFO. For details, see SQL Self-Diagnosis.
- Prefixes gs and pgxc indicate views showing single CN information and those showing cluster information, respectively. Common users can log in to a CN in the cluster to query only views with the gs prefix.
- When you query this type of views, there will be network latency, because the views obtain resource usage in real time.
- If an instance fault occurs, some Top SQL statement information may fail to be recorded in real-time resource monitoring views.
- Top SQL statements are recorded in real-time resource monitoring views as follows:
 - Special DDL statements, such as SET, RESET, SHOW, ALTER SESSION SET, and SET CONSTRAINTS, are not recorded.
 - DDL statements, such as CREATE, ALTER, DROP, GRANT, REVOKE, and VACUUM, are recorded.
 - DML statements are recorded, including:
 - the execution of SELECT, INSERT, UPDATE, and DELETE
 - the execution of EXPLAIN ANALYZE and EXPLAIN PERFORMANCE
 - the use of the query-level or perf-level views
 - The entry statements for invoking functions and stored procedures are recorded.
 When the GUC parameter enable_track_record_subsql is enabled, some internal statements (except the DECLARE definition statement) of a stored procedure can be recorded. Only the internal statements delivered to DNs for execution are recorded, and the remaining internal statements are filtered out.
 - The anonymous block statement is recorded. When the GUC parameter enable_track_record_subsql is enabled, some internal statements of an anonymous block can be recorded. Only the internal statements delivered to DNs for execution are recorded, and the remaining internal statements are filtered out.
 - The cursor statements are recorded. If a cursor does not read data from the cache but triggers the condition for delivering the statement to a DN for execution, the cursor statement is recorded and the statement and execution plan are enhanced. However, if the cursor reads data from the cache, the cursor statement is not recorded. When a cursor statement is used in an anonymous block or function and the cursor reads a large amount of data from a DN but is not fully used, the monitoring information about the cursor on the DN cannot be recorded due to the current architecture limitation. The **With Hold** cursor syntax has a special execution logic. It executes queries during transaction committing. If a statement execution error is reported during this period of time, the **aborted** status of the job cannot be recorded in the TopSQL history table.
 - Jobs in a redistribution process are not monitored.
 - The parameters of a statement with placeholders executed by JDBC are generally specified. However, if the length of the parameter and the original statement exceeds 64 KB, the parameter is not recorded. If the statement is a lightweight statement, it is directly delivered to the DN for execution and the parameter is not recorded.
 - In cluster 8.1.3 and later versions, the TopSQL monitoring at the query and perf levels does not affect the query performance. The default value of the GUC parameter resource_track_cost for resource monitoring of statements has been changed to 0. When you query the TopSQL real-time monitoring view, by default, all statements that are being executed are displayed.
 - In 8.1.3 and later versions, if the GUC parameter enable_track_record_subsql for querying the TopSQL monitoring view is enabled, regardless of whether the

- substatement monitoring function is enabled in the service statements, you can view the substatement running information in the TopSQL monitoring view.
- You are advised not to fully enable substatement monitoring in stored procedures, that is, enable_track_record_subsql, in the 8.1.3 cluster version. Because the substatements cannot be filtered by time, fully enabling substatement monitoring may record too many substatements. As a result, archived monitoring tables occupy a large amount of disk space. In the 8.1.3 cluster version, you are advised to enable only the parameters in the corresponding session when querying real-time monitoring information or locating and analyzing some stored procedures. Starting from cluster versions 8.2.1 and later, a new customizable GUC parameter resource_track_subsql_duration is added. By default, it is set to 180 seconds. This parameter allows you to filter and archive substatements based on their execution time.
- Due to specification restrictions, the records of the main statements that are not written to disks in the TopSQL history table are delayed. The records are displayed in the TopSQL history table only when the job is delivered next time.
- In the 8.2.1.200 cluster version, operator_realtime-level top SQL runtime monitoring is added to provide operator-level real-time monitoring. After operator_realtime-level monitoring is enabled, you can query the execution plan and detailed execution information of statements. When you query the operator-level real-time monitoring view of top SQL statements, by default, all statements that are being executed are displayed. However, in stored procedure and cursor scenarios, operator-level real-time monitoring information cannot be displayed. Querying information about all statements imposes great pressure on the CN memory. To ensure job performance, the
 - **pg_stat_get_wlm_realtime_operator_info(queryid)** function is provided for querying a single statement. You can use this function to query the operator execution information of a specified statement. This version does not support the query of historical operator information.
- operator_realtime-level TopSQL runtime monitoring is not supported for lightweight CN statements and stored procedures. In addition, due to the high execution speed of operators, the display of the operator information may lag behind.
- The spill_size field at the query level (job monitoring) and operator level (operator monitoring) varies due to the statistical dimension. The spill size at the query level is the statement files spilled to disks, and the spill size at the operator level is the read and write I/O volume of a specific operator at the logical layer.
- When the GUC parameter enable_stream_operator is set to off, the displayed operator execution information may be inaccurate.

Prerequisites

- The GUC parameter enable_resource_track is set to on. The default value is on.
- The GUC parameter resource_track_level is set to query, perf, operator_realtime, or operator. The default value is query.
- Job monitoring rules are as follows:
 - Jobs whose execution cost estimated by the optimizer is greater than or equal to resource_track_cost.
- If the Cgroups function is properly loaded, you can run the **gs_cgroup -P** command to view information about Cgroups.
- The GUC parameter **enable_track_record_subsql** specifies whether to record internal statements of a stored procedure or anonymous block.

In the preceding prerequisites, enable_resource_track is a system-level parameter that specifies whether to enable resource monitoring, resource_track_level is a

session-level parameter. You can set the resource monitoring level of a session as needed. The following table describes the values of the two parameters.

Table 12-2 Setting the resource monitoring level to collect statistics

enable_resource_ track	resource_track_le vel	Query-Level Information	Operator-Level Information
on(default)	none	Not collected	Not collected
	query(default)	Collected	Not collected
	perf	Collected	Not collected
	operator	Collected	Collected
on(default)	operator_realtime	Collected	Real-time operator monitoring
off	none/query/ operator	Not collected	Not collected

Procedure

- **Step 1** Query for the real-time CPU information in the **gs_session_cpu_statistics** view. **SELECT * FROM** gs_session_cpu_statistics;
- **Step 2** Query for the real-time memory information in the **gs_session_memory_statistics** view.

SELECT * FROM gs_session_memory_statistics;

Step 3 Query for the real-time resource information about the current CN in the gs_wlm_session_statistics view.

SELECT * FROM gs_wlm_session_statistics;

Step 4 Query for the real-time resource information about all CNs in the **pgxc_wlm_session_statistics** view.

SELECT * FROM pgxc_wlm_session_statistics;

Step 5 Query for the real-time resource information about job operators on the current CN in the **gs_wlm_operator_statistics** view.

SELECT * FROM gs_wlm_operator_statistics;

Step 6 Query for the real-time resource information about job operators on all CNs in the **pgxc_wlm_operator_statistics** view.

SELECT * FROM pgxc_wlm_operator_statistics;

Step 7 Query for the load management information about the jobs executed by the current user in the **PG_SESSION_WLMSTAT** view.

SELECT * FROM pg_session_wlmstat;

Step 8 Query the job execution status of the current user on each CN in the **pgxc_wlm_workload_records** view (this view is available when the dynamic load function is enabled, that is, **enable_dynamic_workload** is set to **on**).

SELECT * FROM pgxc_wlm_workload_records;

----End

12.6 Historical Top SQL

You can query historical Top SQL in historical resource monitoring views. The historical resource monitoring view records the resource usage (including memory, data spilled to disks, and CPU time), running status (including errors, termination, and exceptions), and performance alarm information when a job is complete. For queries that abnormally terminate due to FATAL or PANIC errors, their status is displayed as **aborted** and no detailed information is recorded. Status information about query parsing in the optimization phase cannot be monitored.

The following table describes the external interfaces of the historical views.

Level	Monitore d Node	View	
Query level/perf level (recomm ended)	Current CN	History (Internal dump interface. Only statements that have ended in the last three minutes are displayed.)	GS_WLM_SESSION_HISTO RY
		History (all statements)	GS_WLM_SESSION_INFO
	All CNs	History (Internal dump interface. Only statements that have ended in the last three minutes are displayed.)	PGXC_WLM_SESSION_HIS TORY
		History (all statements)	PGXC_WLM_SESSION_INF O
Operator level	Current CN	History (Only statements that have ended in the last three minutes are displayed.)	GS_WLM_OPERATOR_HIS TORY
		History (internal dump interface, all statements)	GS_WLM_OPERATOR_INF O
	All CNs	History (Only statements that have ended in the last three minutes are displayed.)	PGXC_WLM_OPERATOR_ HISTORY
		History (internal dump interface, all statements)	PGXC_WLM_OPERATOR_I NFO

□ NOTE

- The view level is determined by the resource monitoring level, that is, the resource track level configuration.
- The perf and operator levels affect the values of the query_plan and warning columns in GS_WLM_SESSION_STATISTICS or PGXC_WLM_SESSION_INFO. For details, see SQL Self-Diagnosis.
- Prefixes gs and pgxc indicate views showing single CN information and those showing cluster information, respectively. Common users can log in to a CN in the cluster to query only views with the gs prefix.
- If instance fault occurs, some SQL statement information may fail to be recorded in historical resource monitoring views.
- In some abnormal cases, the status information column in the historical Top SQL may be displayed as **unknown**. The recorded monitoring information may be inaccurate.
- The SQL statements that can be recorded in historical resource monitoring views are the same as those recorded in real-time resource monitoring views. For details, see SQL statements recorded in real-time resource monitoring views.
- Historical top SQL statements are recorded only when the GUC parameter enable resource record is enabled.
- You can query historical Top SQL queries and operator-level data only through the PostgreSQL database.
- Historical Top SQL focuses on locating and demarcating query performance problems. It is not used for auditing or recording syntax analysis error statements.
- In 8.2.1 and later cluster versions, the **resource_track_subsql_duration** parameter (default value: 180s) is added to filter out substatements in the stored procedure whose execution time is less than the value of this parameter and archive only substatements whose execution time is greater than the value of this parameter. In 8.2.1 and later versions, the default value of **enable_track_record_subsql** is changed from **off** to **on**, which means substatements in stored procedures are recorded by default. If a substatement is recorded, it must meet the following conditions:
 - In the session where the statement is, the enable_track_record_subsql parameter is enabled.
 - The substatement must be pushed down to DNs for execution. (To prevent TopSQL from recording too many substatements, substatements that are not pushed down to DNs will be filtered out.)
 - The execution time of the substatement exceeds the value of resource_track_subsql_duration in the session.
- By default, the History view queries statements that end in the last 3 minutes. It does this by querying tables. It is actually a temporary view for performance considerations. Since the 8.1.3 cluster version, the real-time monitoring and archiving functions of the TopSQL monitoring have been greatly improved are no performance considerations are needed. Therefore, you are not advised to use the History view.
- In 8.1.3 and later versions, the TopSQL real-time monitoring has no impact on statement performance. You can set the GUC **parameter resource_track_cost** to **0** to monitor the running information of all statements. The statement archiving in the TopSQL history monitoring also has no impact on statement performance. However, when the TPS is high, the following factors need to be considered:
 - Record the disk overhead of all statements. You can estimate the disk space required for archiving a statement as 8 KB, calculate the space usage based on the peak TPS, and adjust the values of resource_track_duration and resource track subsql duration.
 - For memory overhead for caching all statements, you can estimate the memory size required for archiving a statement as 16 KB, and the interval for archiving statements in batches as 5 seconds, then calculate the required peak memory size based on the peak service TPS. The calculation method is as follows: 5 seconds x TPS x 16 KB. The value of session_history_memory GUC (default value: 100 MB)

must be greater than the calculation result to ensure that all statements can be recorded.

Prerequisites

- The GUC parameter enable_resource_track is set to on. The default value is on.
- The GUC parameter **resource_track_level** is set to **query**, **perf**, or **operator**. The default value is **query**. For details, see **Table 12-2**.
- The GUC parameter enable_resource_record is set to on. The default value is on.
- The GUC parameter **resource_track_duration** is less than the sum of the job execution time and queuing time (**60s** by default).
- The GUC parameter enable_track_record_subsql specifies whether to record internal statements of a stored procedure or anonymous block. The default value is on.
- The value of **resource_track_subsql_duration** is less than the execution time of the internal statement in the stored procedure (180s by default).
- Jobs whose sum of the job execution time and queuing time recorded in the real-time resource monitoring view (see Table 12-1) is no less than the value of resource_track_duration are monitored.
- If the Cgroups function is properly loaded, you can run the **gs_cgroup -P** command to view information about Cgroups.

Procedure

Step 1 Query the load records of the current CN after its latest job is complete in the **gs_wlm_session_history** view.

SELECT * FROM gs_wlm_session_history;

Step 2 Query the load records of all the CNs after their latest job are complete in the **pgxc_wlm_session_history** view.

SELECT * FROM pgxc_wlm_session_history;

Step 3 Query the load records of the current CN through the **gs_wlm_session_info** table after the task is complete. To query the historical records successfully, set **enable resource record** to **on**.

SELECT * FROM gs_wlm_session_info;

 Show the 10 queries that consume the most memory (You can specify a query period.):

SELECT * FROM *gs_wlm_session_info* **order by** *max_peak_memory* **desc limit** *10;* **SELECT * FROM** *gs_wlm_session_info* WHERE start_time >= '2022-05-15 21:00:00' and finish_time <='2022-05-15 23:30:00' **order by** *max_peak_memory* **desc limit** *10;*

Show the 10 queries consuming the most CPU resources:

SELECT * FROM gs_wlm_session_info order by total_cpu_time desc limit 10; SELECT * FROM gs_wlm_session_info WHERE start_time >= '2022-05-15 21:00:00' and finish_time <='2022-05-15 23:30:00' order by total_cpu_time desc limit 10;

Step 4 Query for the load records of all the CNs after their jobs are complete in the **pgxc_wlm_session_info** view. To query the historical records successfully, set **enable_resource_record** to **on**.

SELECT * FROM *pgxc_wlm_session_info*;

• Showing the 10 queries on which the CN spends the most time:

SELECT * FROM paxc_wlm_session_info order by duration desc limit 10;

• Query the execution information about a query statement that has been executed. For example, query the execution information about the statement whose **queryid** is **76561193695026478**.

SELECT * FROM *pgxc_wlm_session_info* where queryid = '76561193695026478';

Step 5 Use the pgxc_get_wlm_session_info_bytime function to filter and query the pgxc_wlm_session_info view. To query the historical records successfully, set enable_resource_record to on. You are advised to use this function if the view contains a large number of records.

□ NOTE

A GaussDB(DWS) cluster uses the UTC time by default, which has an 8-hour time difference with the system time. Before queries, ensure that the database time is the same as the system time.

Return the queries started between 2019-09-10 15:30:00 and 2019-09-10
 15:35:00 on all CNs. For each CN, a maximum of 10 queries will be returned.

SELECT * FROM pgxc_get_wlm_session_info_bytime('start_time', '2019-09-10 15:30:00', '2019-09-10 15:35:00', 10);

• Return the queries ended between **2019-09-10 15:30:00** and **2019-09-10 15:35:00** on all CNs. For each CN, a maximum of 10 queries will be returned.

SELECT * FROM pgxc_get_wlm_session_info_bytime('finish_time', '2019-09-10 15:30:00', '2019-09-10 15:35:00', 10);

Step 6 Query the recent resource information of the job operators on the current CN in the **gs_wlm_operator_history** view. Ensure that **resource_track_level** is set to **operator**.

SELECT * FROM gs_wlm_operator_history;

Step 7 Query the recent resource information of the job operators on all the CNs in the pgxc_wlm_operator_history view. Ensure that resource_track_level is set to operator.

SELECT * FROM pgxc_wlm_operator_history;

Step 8 Query the recent resource information of the job operators on the current CN in the **gs_wlm_operator_info** view. Ensure that **resource_track_level** is set to **operator** and **enable_resource_record** to **on**.

SELECT * FROM gs_wlm_operator_info;

Step 9 Query for the historical resource information of job operators on all the CNs in the pgxc_wlm_operator_info view. Ensure that resource_track_level is set to operator and enable_resource_record to on.

SELECT * FROM pgxc_wlm_operator_info;

□ NOTE

- The number of data records that can be retained in the memory is limited due to the preset memory limit. After the real-time query is complete, the data records are imported to historical views. For a query-level view, when the number of queries to be recorded exceeds the upper limit allowed by the memory, the current query cannot be recorded and the next query is performed based on a new rule. On each CN, the memory usage of the query-level historical view is recorded (100 MB by default). You can query the data in the PG TOTAL MEMORY DETAIL view.
- For operator-level views, whether a record can be stored depends on the upper limit allowed by the memory at that time point. If the number of plan nodes plus the number of records in the memory exceeds the upper limit, the record cannot be stored. On each CN, the maximum numbers of real-time and historical operator-level records that can be stored in the memory are max_oper_realt_num (set to 56987 by default) and max_oper_hist_num (set to 113975 by default), respectively. The average number of plan nodes of a query is num_plan_node. Maximum number of concurrent tasks allowed by real-time views on each CN is: num_realt_active = max_oper_realt_num/num_plan_node. Maximum number of concurrent tasks allowed by historical views on each CN is: num_hist_active = max_oper_hist_num/(180/run_time)/num_plan_node.
- In high concurrency, ensure that the number of queries to be recorded does not exceed
 the maximum values set for query- and operator-level views. You can modify the
 memory of the historical query view by configuring the session_history_memory
 parameter. The memory size increases in direct proportion to the maximum number of
 queries that can be recorded.

12.7 Example for Querying for Top SQLs

In this section, TPC-DS sample data is used as an example to describe how to query **Real-time Top SQL** and **Historical Top SQL**.

Configuring Cluster Parameters

To query for historical or archived resource monitoring information about jobs of top SQLs, you need to set related GUC parameters first. The procedure is as follows:

- 1. Log in to the GaussDB(DWS) console.
- 2. On the **Cluster Management** page, locate the required cluster and click the cluster name. The cluster details page is displayed.
- 3. Click the **Parameter Modifications** tab to view the values of cluster parameters.
- Set an appropriate value for parameter resource_track_duration and click Save.

\cap	\cap	П	N	O	т	F
_	_		I V	${}^{\sim}$		_

If **enable_resource_record** is set to **on**, storage space expansion may occur and thereby slightly affects the performance. Therefore, set is to **off** if record archiving is unnecessary.

5. Go back to the **Cluster Management** page, click the refresh button in the upper right corner, and wait until the cluster parameter settings are applied.

Example for Querying for Top SQLs

The TPC-DS sample data is used as an example.

- **Step 1** Open the SQL client tool and connect to your database.
- **Step 2** Run the **EXPLAIN** statement to query for the estimated cost of the SQL statement to be executed to determine whether resources of the SQL statement will be monitored.

By default, only resources of a query whose execution cost is greater than the value of **resource track cost** are monitored and can be queried by users.

For example, run the following statements to query for the estimated execution cost of the SQL statement:

```
SET CURRENT_SCHEMA = tpcds;
EXPLAIN WITH customer_total_return AS
( SELECT sr_customer_sk as ctr_customer_sk,
sr_store_sk as ctr_store_sk,
sum(SR_FEE) as ctr_total_return
FROM store_returns, date_dim
WHERE sr_returned_date_sk = d_date_sk AND d_year =2000
GROUP BY sr_customer_sk, sr_store_sk)
SELECT c_customer_id
FROM customer_total_return ctr1, store, customer
WHERE ctr1.ctr_total_return > (select avg(ctr_total_return)*1.2
FROM customer_total_return ctr2
WHERE ctr1.ctr_store_sk = ctr2.ctr_store_sk)
AND s_store_sk = ctr1.ctr_store_sk
AND s_state = 'TN'
AND ctr1.ctr_customer_sk = c_customer_sk
ORDER BY c_customer_id
limit 100;
```

In the following query result, the value in the first row of the **E-costs** column is the estimated cost of the SQL statement.

Figure 12-1 EXPLAIN result



In this example, to demonstrate the resource monitoring function of top SQLs, set the value of **resource_track_cost** to **100**, which should be lower than the estimated cost in the **EXPLAIN** query result. For more details, see **resource_track_cost**.

■ NOTE

After completing this example, you still need to reset **resource_track_cost** to its default value **100000** or a proper value. An overly small parameter value will compromise the database performance.

Step 3 Run SQL statements.

SET CURRENT SCHEMA = tpcds: WITH customer_total_return AS (SELECT sr_customer_sk as ctr_customer_sk, sr_store_sk as ctr_store_sk, sum(SR_FEE) as ctr_total_return FROM store_returns,date_dim WHERE sr_returned_date_sk = d_date_sk AND d_year =2000 GROUP BY sr_customer_sk ,sr_store_sk) SELECT c_customer_id FROM customer_total_return ctr1, store, customer WHERE ctr1.ctr_total_return > (select avg(ctr_total_return)*1.2 FROM customer_total_return ctr2 WHERE ctr1.ctr_store_sk = ctr2.ctr_store_sk) AND s_store_sk = ctr1.ctr_store_sk AND s_state = 'TN' AND ctr1.ctr_customer_sk = c_customer_sk ORDER BY c_customer_id limit 100;

Step 4 During statement execution, query for the real-time memory peak information about the SOL statement on the current CN.

SELECT query,max_peak_memory,average_peak_memory,memory_skew_percent FROM gs_wlm_session_statistics ORDER BY start_time DESC;

The preceding command queries for the real-time peak information at the query-level. The peak information includes the maximum memory peak among all DNs per second, average memory peak among all DNs per second, and memory usage skew across DNs.

For more examples of querying for the real-time resource monitoring information of top SQLs, see **Real-time Top SQL**.

Step 5 Wait until the SQL statement execution in **Step 3** is complete, and then query for the historical resource monitoring information of the statement.

SELECT query,start_time,finish_time,duration,status FROM gs_wlm_session_history ORDER BY start_time desc;

The preceding command queries for the historical information at the query-level. The peak information includes the execution start time, end time, actual execution time, and execution status. The time unit is ms.

For more examples of querying for the historical resource monitoring information of top SQLs, see **Historical Top SQL**.

Step 6 Wait for 3 minutes after the execution of the SQL statement in **Step 3** is complete, query for the historical resource monitoring information of the statement in the **info** view.

If enable_resource_record is set to on and the execution time of the SQL statement in Step 3 is no less than the value of resource_track_duration, historical information about the SQL statement will be archived to the gs_wlm_session_info view 3 minutes after the execution of the SQL statement is complete.

The **info** view can be queried only when the **postgres** database is connected. Therefore, switch to the **postgres** database before running the following statement:

SELECT query,start_time,finish_time,duration,status FROM gs_wlm_session_info ORDER BY start_time desc;

13 GaussDB(DWS) Performance Tuning

13.1 Overview

Database performance tuning is the process of optimizing database system configuration and SQL queries to improve database performance and efficiency. The purpose includes eliminating performance bottlenecks, reducing response times, increasing throughput and resource utilization, cutting costs, and improving system stability.

This section provides comprehensive guidance for DBAs on performance diagnosis, system tuning, and SQL tuning, as well as practical examples of SQL tuning.

Precautions

- Database performance tuning is a complex and intricate process. To achieve the optimal performance and efficiency, performance tuning must take into consideration multiple factors, such as hardware, software, queries, configuration, and data structures. Engineers performing the performance tuning must be familiar with how database systems work in great detail, including a deep understanding of the system software architecture, software and hardware configurations, database configuration parameters, concurrency control, query handling, and database applications.
- Performance tuning sometimes requires a cluster restart, which may interrupt services. To avoid that, you are advised to schedule performance tuning tasks that require a cluster restart to occur during off-peak hours.

Performance Tuning Process

Figure 13-1 illustrates the performance tuning process.

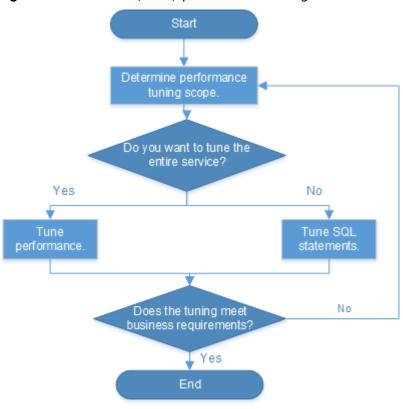


Figure 13-1 GaussDB(DWS) performance tuning

Table 13-1 gives a brief introduction to each phase of the performance tuning process.

Table 13-1 Phase-by-phase introduction to GaussDB(DWS) performance tuning

Phase	Description
Performance Diagnosis	Obtain the CPU, memory, I/O, and network resource usage of each node to check whether these resources are fully utilized and whether any performance bottlenecks exist.
System Optimization	Perform OS and database system-level performance tuning to achieve better utilization of existing CPU, memory, I/O, and network resources, prevent resource conflicts, and improve query throughput.

Phase	Description	
SQL Tuning	Analyze the SQL statements used and determine whether any optimization can be performed. Analysis of SQL statements comprises:	
	 Generating table statistics using ANALYZE: The ANALYZE statement collects statistics about the database table content. Statistical results are stored in the system catalog PG_STATISTIC. The execution plan generator uses these statistics to determine which one is the most effective execution plan. 	
	 Analyzing the execution plan: The EXPLAIN statement displays the execution plan of SQL statements, and the EXPLAIN PERFORMANCE statement displays the execution time of each operator in SQL statements. 	
	 Identifying the root causes of issues: Identify possible causes by analyzing the execution plan and perform specific optimization by modifying database-level SQL optimization parameters. 	
	 Compiling better SQL statements: Compile better SQL statements in the scenarios, such as cache of intermediate and temporary data for complex queries, result set cache, and result set combination. 	

13.2 Performance Diagnosis

13.2.1 Cluster Performance Analysis

The node specifications of different GaussDB(DWS) clusters may vary in terms of the number of CPU cores, memory capacity, and node storage capacity. Different specifications lead to different service handling capacity and performance. Before creating a cluster, you need to select the appropriate cluster specifications based on the actual workloads and application scenario.

If the workloads increase, more resources (such as CPU, memory, and network bandwidth) will be needed in order to maintain the same level of database performance. Insufficient cluster resources will lead to performance issues.

GaussDB(DWS) provides abundant monitoring metrics that you can use to monitor cluster performance and status, including CPU usage, memory usage, disk usage, disk I/O, and network I/O. For any abnormality, you can check the metrics to locate the root cause.

If your service requires additional compute or storage resources, expand the capacity of an existing cluster by adding more nodes to it or changing node specifications through the management console.

13.2.2 Slow SQL Analysis

13.2.2.1 Querying SQL Statements That Affect Performance Most

This section describes how to query SQL statements whose execution takes a long time, leading to poor system performance.

Procedure

Step 1 Query the statements that are run for a long time in the database.

SELECT current_timestamp - query_start AS runtime, datname, usename, query FROM pg_stat_activity where state != 'idle' ORDER BY 1 desc;

After the query, query statements are returned as a list, ranked by execution time in descending order. The first result is the query statement that has the longest execution time in the system. The returned result contains the SQL statement invoked by the system and the SQL statement run by users. Find the statements that were run by users and took a long time.

Alternatively, you can set **current_timestamp - query_start** to be greater than a threshold to identify query statements that are executed for a duration longer than this threshold.

SELECT query FROM pg_stat_activity WHERE current_timestamp - query_start > interval '1 days';

Step 2 Set the parameter **track_activities** to **on**.

```
SET track_activities = on;
```

The database collects the running information about active queries only if the parameter is set to **on**.

Step 3 View the running guery statements.

Viewing **pg stat activity** is used as an example here.

If the **state** column is idle, the connection is idle and requires a user to enter a command.

To identify only active query statements, run the following command:

SELECT datname, usename, state FROM pg_stat_activity WHERE state != 'idle';

- **Step 4** Analyze the status of the query statements that were run for a long time.
 - If the guery statement is normal, wait until the execution is complete.
 - If a query statement is blocked, run the following command to view this query statement:

SELECT datname, usename, state, query FROM pg_stat_activity WHERE waiting = true;

The command output lists a query statement in the block state. The lock resource requested by this query statement is occupied by another session, so this query statement is waiting for the session to release the lock resource.

□ NOTE

Only when the query is blocked by internal lock resources, the **waiting** field is **true**. In most cases, block happens when query statements are waiting for lock resources to be released. However, query statements may be blocked because they are waiting to write in files or for timers. Such blocked queries are not displayed in the **pg_stat_activity** view.

----End

13.2.2.2 Checking Blocked Statements

During database running, query statements are blocked in some service scenarios and run for an excessively long time. In this case, you can forcibly terminate the faulty session.

Procedure

Step 1 View blocked query statements and information about the tables and schemas that block the query statements.

```
SELECT w.query as waiting_query,
w.pid as w_pid,
w.usename as w_user,
l.query as locking_query,
l.pid as l_pid,
l.usename as l_user,
t.schemaname || '.' || t.relname as tablename
from pg_stat_activity w join pg_locks l1 on w.pid = l1.pid
and not l1.granted join pg_locks l2 on l1.relation = l2.relation
and l2.granted join pg_stat_activity l on l2.pid = l.pid join pg_stat_user_tables t on l1.relation = t.relid
where w.waiting;
```

The thread ID, user information, query status, as well as information about the tables and schemas that block the query statements are returned.

Step 2 Run the following command to terminate the required session, where **139834762094352** is the thread ID:

SELECT PG_TERMINATE_BACKEND(139834762094352);

If information similar to the following is displayed, the session is successfully terminated:

```
PG_TERMINATE_BACKEND
------
t
(1 row)
```

If a command output similar to the following is displayed, a user is attempting to terminate the session, and the session will be reconnected rather than being terminated.

FATAL: terminating connection due to administrator command
FATAL: terminating connection due to administrator command
The connection to the server was lost. Attempting reset: Succeeded.

◯ NOTE

If the **PG_TERMINATE_BACKEND** function is used to terminate the background threads of the session, the gsql client will be reconnected rather than be logged out.

13.2.3 SQL Diagnosis

GaussDB(DWS) clusters support SQL diagnosis, which shows the complete execution plans of specific SQL queries. You can search for SQL queries (such as slow queries) using a combination of multiple filter criteria.

To use SQL diagnosis, perform the following steps:

- Step 1 Log in to the GaussDB(DWS) console.
- **Step 2** Choose **Dedicated Clusters** > **Clusters** and locate the cluster to be monitored.
- **Step 3** In the **Operation** column of the target cluster, click **Monitoring Panel**.
- **Step 4** In the navigation pane on the left, choose **Utilities** > **SQL Diagnosis**. The metrics include:
 - Query ID
 - Database
 - Schema Name
 - User Name
 - Client
 - Client IP Address
 - Running Time (ms)
 - CPU Time (ms)
 - Scale-Out Started
 - Completed
 - Details
- **Step 5** On the **SQL Diagnosis** page, you can view the SQL diagnosis information. In the **Details** column of a specified query ID, click **View** to view the detailed SQL diagnosis result, including:
 - Alarm Information
 - SQL Statement
 - Execution Plan



13.2.4 Table Diagnosis

GaussDB(DWS) provides statistics and diagnostic tools for you to learn table status, including:

- Skew Rate: monitors and analyzes uneven data distribution in a cluster, and displays information about the 50 largest tables whose skew rate is higher than 5%.
- Dirty Page Rate: monitors and analyzes dirty pages in a cluster, and displays information about the 50 largest tables whose dirty page rate is higher than 50%.

Skew Rate

Improper distribution columns can cause severe skew during operator computing or data spill to disk. The workloads will be unevenly distributed on DNs, resulting in high disk usage on individual DNs and affecting their performance. After identifying tables with a high skew rate and a relatively large size, you can reselect distribution columns for these tables to have their data redistributed. For details, see How Do I Change Distribution Columns?

Procedure

- **Step 1** Log in to the **GaussDB(DWS) console**.
- **Step 2** Choose **Dedicated Clusters** > **Clusters** and locate the cluster to be monitored.
- **Step 3** In the **Operation** column of the target cluster, click **Monitoring Panel**.
- **Step 4** In the navigation tree on the left, choose **Utilities** > **Table Diagnosis** and click the **Skew Rate** tab. The tables that meet the statistics collection conditions in the cluster are displayed.

----End

Dirty Page Rate

DML operations on tables may generate dirty data, which unnecessarily occupies cluster storage. You can identify tables with a high dirty page rate and a relatively large size, and handle them accordingly. For more information, see **Solution to High Disk Usage and Cluster Read-Only**.

Procedure

- Step 1 Log in to the GaussDB(DWS) console.
- **Step 2** Choose **Dedicated Clusters** > **Clusters** and locate the cluster to be monitored.
- **Step 3** In the **Operation** column of the target cluster, click **Monitoring Panel**.
- **Step 4** In the navigation tree on the left, choose **Utilities** > **Table Diagnosis** and click the **Dirty Page Rate** tab. The tables that meet the statistics collection conditions in the cluster are displayed.

13.3 System Optimization

13.3.1 Tuning Database Parameters

To ensure high performance of the database, you are advised to configure GUC parameters based on available resources and the actual workloads. This section describes some of the common parameters and the recommended configurations for them. For more details, see **Configuring GUC Parameters**.

Parameters Related to Database Memory

Table 13-2 Parameters related to database memory

GUC Parameter	Description	Suggestion
max_process_ memory	Specifies the maximum physical memory available to a single CN/DN.	 On DNs, the value of this parameter is determined based on the server's physical memory and the number of DNs deployed on a single node. Parameter value = (Physical memory – vm.min_free_kbytes) x 0.8/(n + Number of primary DNs). This parameter aims to ensure system reliability, preventing node OOM caused by increasing memory usage. vm.min_free_kbytes indicates OS memory reserved for kernels to receive and send data. Its value is at least 5% of the total memory. That is, max_process_memory = Physical memory x 0.8/ (n + Number of primary DNs). If the cluster scale (number of nodes in the cluster) is smaller than 256, n=1; if the cluster scale is larger than 256 and smaller than 512, n=2; if the cluster scale is larger than 512, n=3. Set this parameter on CNs to the same value as that on DNs. RAM is the maximum memory allocated to the cluster.

GUC Parameter	Description	Suggestion
shared_buffer s	Specifies the size of the shared memory used by GaussDB(DWS). If the value of this parameter is increased, GaussDB(DWS) requires more System V shared memory than the default system setting.	It is recommended that shared_buffers be set to a value less than 40% of the memory. Set it to a large value for row-store tables and a small value for column-store tables. Set this parameter to a large value for row storage and a small value for column storage. For column-store tables: shared_buffers = (Memory of a single server/Number of DNs on the single server) x 0.4 x 0.25 If you want to increase the value of shared_buffers, you also need to increase the value of checkpoint_segments, because a longer period of time is required to write a large amount of new or changed data.
cstore_buffers	Specifies the size of the shared buffer used by column-store tables and column-store tables (ORC, Parquet, and CarbonData) of OBS and HDFS foreign tables.	Column-store tables use the shared buffer specified by cstore_buffers instead of that specified by shared_buffers. When column-store tables are mainly used, reduce the value of shared_buffers and increase that of cstore_buffers. Use cstore_buffers to specify the cache of ORC, Parquet, or CarbonData metadata and data for OBS or HDFS foreign tables. The metadata cache size should be 1/4 of cstore_buffers and not exceed 2 GB. The remaining cache is shared by column-store data and foreign table column-store data.

GUC Parameter	Description	Suggestion
work_mem	Specifies the size of the memory used by internal sequential operations and the Hash table before data is written into temporary disk files.	The default value is 512 MB for small-scale memory (max_process_memory is less than 30 GB) and 2 GB for large-scale memory (max_process_memory is greater than or equal to 30 GB).
		When the specified physical memory is insufficient, work_mem determines whether to write additional operator calculation data into temporary tables based on query characteristics and concurrency. This reduces performance by five to ten times and increases query response times from seconds to minutes.
		 In complex serial query scenarios, each query requires five to ten associated operations. Set work_mem using the following formula: work_mem = 50% of the memory/10.
		 In simple serial query scenarios, each query requires two to five associated operations. Set work_mem using the following formula: work_mem = 50% of the memory/5.
		 For concurrent queries, use the formula: work_mem = work_mem in serialized scenario/Number of concurrent SQL statements.

GUC Parameter	Description	Suggestion
maintenance_ work_mem	Specifies the maximum size of memory used for maintenance operations, involving VACUUM, CREATE INDEX, and ALTER TABLE ADD FOREIGN KEY.	If you set this parameter to the value of work_mem, database dump files can be cleaned up and restored more efficiently. In a database session, only one maintenance operation can be performed at a time. Maintenance is usually performed when there are not much sessions. When the automatic cleanup process is running, up to autovacuum_max_workers times of the memory will be allocated. In this case, set maintenance_work_mem to a value greater than or equal to that of work_mem.

Parameters Related to Queue Concurrency in Databases

GUC Parameter	Description	Suggestion
max_active_st atements (global concurrent queue)	Controls the maximum number of concurrent jobs on a single CN.	All common users' jobs are subject to this threshold, regardless of their complexity. When the number of concurrent jobs reaches the specified threshold, the excess jobs have to wait in a queue. Administrator's jobs are exempt from this limit.
		Set the value of this parameter based on system resources, such as CPU, I/O, and memory resources, to ensure that the system resources can be fully utilized and the system will not be crashed due to excessive concurrent jobs.
parctl_min_co st (local concurrent queue)	Controls the maximum number of concurrent jobs within the same resource pool on a single CN.	The number of concurrent complex jobs are controlled based on their cost.

■ NOTE

When tuning the **max_active_statements** parameter (global concurrent queue), pay attention to the following:

- If max_active_statements is set to -1, which indicates that global concurrency is not limited, users may be disconnected in a high concurrency scenario.
- In a point query scenario, set max_active_statements to 100.
- In an analytical query scenario, set max_active_statements to the number of CPU cores divided by the number of DNs. Generally, its value ranges from 5 to 8.

Database Communication Parameters

By default, nodes in a database cluster communicate using the TCP proxy communication library.

Table 13-3 Database communication parameters

GUC Parameter	Description	Suggestion
comm_quota_ size	comm_quota_size controls the size of data transmitted every time in each flow channel. Its default value is 1M.	In a high concurrency scenario, you can increase its value to improve communication performance, but doing so consumes more memory. Optimize this parameter as needed. If you query the pg_total_memory_detail view of a DN and find that the memory used by the communication layer has reached the threshold of comm_usable_memory, set comm_quota_size to a small value, such as 512K.
comm_usable _memory	comm_usable_memory controls the memory on a DN that can be used for database communication.	The value of this parameter is only used for memory flow control. The default flow control value is 1 MB. If the memory usage exceeds half of the parameter value, the flow control value will be automatically changed to 0.5 MB. If only 20% of the memory specified by the parameter is available, the flow control value will be changed to the allowed minimum, 8 KB.

Database Connection Parameters

Table 13-4 Database connection parameters

GUC Parameter	Description	Suggestion
max_connecti ons	Specifies the maximum number of concurrent connections to the database. This parameter affects the concurrent processing capability of the cluster.	Retain the default value of this parameter on CNs. Set this parameter on DNs to a value calculated using this formula: Number of CNs x Value of this parameter on a CN. If the value of this parameter is increased, GaussDB(DWS) may require more System V shared memory or semaphore, which may exceed the default maximum value of the OS. In this case, modify the value as needed.
max_prepared _transactions	Specifies the maximum number of transactions that can stay in the prepared state simultaneously. If the value of this parameter is increased, GaussDB(DWS) requires more System V shared memory than the default system setting.	The value of max_connections is related to max_prepared_transactions. Before configuring max_connections, ensure that the value of max_prepared_transactions is greater than or equal to that of max_connections. In this way, each session has a prepared transaction in the waiting state.
session_timeo ut	Specifies the maximum duration a database connection can stay idle before it is automatically disconnected.	The value can be an integer in the range 0 to 86400. The minimum unit is second (s). The value 0 disables this timeout mechanism. Generally, you are advised not to set this parameter to 0 .

Other Performance-related Parameters

Table 13-5 Other performance-related parameters

GUC Parameter	Description	Suggestion
enable_dyna mic_workload	Specifies whether to enable dynamic load management. Dynamic load management refers to the automatic queue control of complex queries based on user loads in a database. This fine-tunes system parameters without manual adjustment.	This parameter is enabled by default. Notes: Simple query jobs (which are estimated to require less than 32 MB memory) and non-DML statements (statements other than INSERT, UPDATE, DELETE, and SELECT) have no adaptive load restrictions. Control the upper memory limits for them on a single CN using max_active_statements. In adaptive load scenarios, the value cannot be increased. If you increase it, memory cannot be controlled for certain statements, such as statements that have not been analyzed. Reduce concurrency in the following scenarios, because high concurrency may lead to uncontrollable memory usage. A single tuple occupies excessive memory, for example, a base table contains a column more than 1 MB wide. A query is fully pushed down. A statement occupies a large amount of memory on the CN, for example, a statement that cannot be pushed down or a cursor withholding statement. An execution plan creates a hash table based on the hash join operator, and the table has many duplicate values and occupies a large amount of memory. UDFs are used, which occupy a large amount of memory. When configuring this parameter, you can set query_dop to 0 (adaptive). In this case, the system dynamically selects the optimal degree of parallelism (DOP) for

GUC Parameter	Description	Suggestion
		each query based on resource usage and the execution plan. The enable_dynamic_workload parameter supports the dynamic memory allocation.
bulk_write_rin g_size	Specifies the size of a ring buffer used for parallel data import.	This parameter affects the database import performance. You are advised to increase the value of this parameter on DNs when a large amount of data is to be imported. The default value is 2GB .
data_replicate _buffer_size	Specifies the memory used by queues when the sender sends data pages to the receiver.	The value of this parameter affects the buffer size for data replication between the primary and standby servers. The default value is 16 MB for a CN and 128 MB for a DN. If the server memory is 256 GB, you can increase the value to 512 MB.

13.3.2 SMP Parallel Execution

Complex queries may take a long time. In a system with low concurrency support, this can be a problem. SMP is used to implement operator-level parallel execution, which can effectively speed up queries, improving query performance and resource utilization.

The SMP feature improves performance through operator parallelism but may drive more resource usage, including CPU, memory, network, and I/O. In essence, SMP is a method that trades resources for time, meaning it accelerates queries at the cost of additional resources. It improves system performance in appropriate scenarios and when resources are sufficient, but may also deteriorate performance if used inappropriately. Furthermore, compared with serial processing, SMP generates more candidate plans, which is more time-consuming and may hurt performance.

The SMP feature of GaussDB(DWS) is controlled by the GUC parameter **query_dop**. Users use this parameter to specify an appropriate degree of query parallelism.

Application Scenarios and Constraints for SMP

Applicable Scenarios

Operators supporting parallel processing are used.
 The execution plan contains the following operators:

- a. Scan: Row Storage common table and a line memory partition table sequential scanning, column-oriented storage ordinary table and columnoriented storage partition table sequential scanning, HDFS internal and external table sequence scanning. Surface scanning GDS data can be imported at the same time. All of the above does not support replication tables.
- b. Join: HashJoin, NestLoop
- c. Agg: HashAgg, SortAgg, PlainAgg, and WindowAgg, which supports only **partition by**, and does not support **order by**.
- d. Stream: Redistribute, Broadcast
- e. Other: Result, Subqueryscan, Unique, Material, Setop, Append, VectoRow, RowToVec
- SMP-unique operators are used.

To execute queries in parallel, Stream operators are added for data exchange for the SMP feature. These new operators can be considered as the subtypes of Stream operators.

- a. Local Gather aggregates data of parallel threads within a DN
- b. Local Redistribute redistributes data based on the distributed key across threads within a DN
- c. Local Broadcast broadcasts data to each thread within a DN.
- d. Local RoundRobin distributes data in polling mode across threads within a DN
- e. Split Redistribute redistributes data across parallel threads on different DNs.
- f. Split Broadcast broadcasts data to all parallel DN threads in the cluster.

Among these operators, Local operators exchange data between parallel threads within a DN, and non-Local operators exchange data across DNs.

Example

The TPCH Q1 parallel plan is used as an example.

```
operation
 1 | -> Row Adapter
 2 |
      -> Vector Streaming (type: GATHER)
          -> Vector Sort
             -> Vector Streaming(type: LOCAL GATHER dop: 1/4)
 5 |
                -> Vector Hash Aggregate
                   -> Vector Streaming(type: SPLIT REDISTRIBUTE dop: 4/4)
 6 |
                      -> Vector Hash Aggregate
                         -> Vector Append(9, 10)
 9 1
                            -> Dfs Scan on lineitem
10 |
                            -> Vector Adapter
11 |
                               -> Seq Scan on pg_delta_1423863972 lineitem |
(11 rows)
```

In this plan, implement the Hdfs Scan and HashAgg operator parallel, and adds the Local Gather and Split Redistribute data exchange operator.

In this example, the sixth operator is Split Redistribute, and **dop: 4/4** next to the operator indicates that the degree of parallelism of the sender and receiver is 4. 4 No operator is Local Gather, marked dop: 1/4 above, this operator sender thread parallel degree is 4, while the receiving end thread parallelism degree to 1, that is, lower-layer 5 number Hash Aggregate

operators according to the 4 parallel degree, while the working mode of the port on the upper-layer 1 to 3 number operator according to the executed one by one, 4 number operator is used to achieve intra-DN concurrent threads data aggregation.

You can view the parallelism situation of each operator in the dop information.

Non-Applicable Scenarios

- 1. Small queries are performed, where plan generation may account for a significant portion of the total query time.
- 2. Operators are processed on CNs.
- 3. Statements that cannot be pushed down are executed.
- 4. The **subplan** of a guery and operators containing a subguery are executed.

Impact of Resource Availability on SMP Performance

The SMP architecture accelerates queries at the cost of additional resources. After the plan parallelism is executed, more resources are consumed, including the CPU, memory, I/O, and network bandwidth. As the DOP grows, the resource consumption also increases. If these resources become a bottleneck, SMP cannot improve performance. On the contrary, it may do exactly the opposite. Adaptive SMP is provided to dynamically select the optimal parallel degree for each query based on the resource usage and query requirements. The following information describes the situations that the SMP affects theses resources:

• CPU resources

In a general customer scenario, the system CPU usage rate is not high. Using the SMP parallelism architecture will fully use the CPU resource to improve the system performance. If the number of CPU kernels of the database server is too small and the CPU usage is already high, enabling the SMP parallelism may deteriorate the system performance due to resource compete between multiple threads.

Memory resources

The query parallel causes memory usage growth, but the memory upper limit used by each operator is still restricted by **work_mem**. Assume that **work_mem** is 4 GB, and the degree of parallelism is 2, then the memory upper limit of each concurrent thread is 2 GB. When **work_mem** is small or the system memory is sufficient, running SMP parallelism may push data down to disks. As a result, the query performance deteriorates.

Network bandwidth resources

To execute queries in parallel, data exchange operators are added. Local stream operators exchange data between threads within a DN. Data is exchanged in memory, so it does not impact network performance. Non-local operators exchange data over the network and increase the network load. If the capacity of a network resource has already become a bottleneck, parallelism may hurt performance.

I/O resources

A parallel scan increases I/O resource consumption. It can improve performance only when I/O resources are sufficient.

Other Factors Impacting SMP Performance

Besides the resource factor, other factors may also impact SMP performance, such as uneven data distribution across tables and the degree of system parallelism.

• Impact of data skew on SMP performance

Serious data skew deteriorates parallel execution performance. For example, if the data volume of a value in the join column is much more than that of other values, the data volume of a parallel thread will be much more than that of others after Hash-based data redistribution, resulting in the long-tail issue and poor parallelism performance.

• Impact of system parallelism degree on SMP performance

The SMP feature uses more resources, and unused resources are decreasing in a high concurrency scenario. Therefore, enabling the SMP parallelism will result in serious resource compete among queries. Once resource competes occur, no matter the CPU, I/O, memory, or network resources, all of them will result in entire performance deterioration. In the high concurrency scenario, enabling the SMP will not improve the performance effect and even may cause performance deterioration.

Suggestions for SMP Parameter Settings

To enable the SMP adaptation function, set **query_dop** to **0** and adjust the following parameters to obtain an optimal DOP selection:

comm_usable_memory

If the system memory is large, the value of **max_process_memory** is large. In this case, you are advised to set the value of this parameter to 5% of **max_process_memory**, that is, 4 GB by default.

comm max stream

The recommended value for this parameter is calculated as follows: comm_max_stream = Min(dop_limit x dop_limit x 20 x 2, max_process_memory (bytes) x 0.025/Number of DNs/260). The value must be within the value range of **comm_max_stream**.

max_connections

The recommended value for this parameter is calculated as follows: $max_connections = dop_limit \times 20 \times 6 + 24$. The value must be within the value range of $max_connections$.



In the preceding formulas, **dop_limit** indicates the number of CPUs corresponding to each DN in the cluster. It is calculated as follows: **dop_limit** = Number of logical CPU cores of a single server/Number of DNs of a single server.

SMP Configuration Procedure

NOTICE

The CPU, memory, I/O, and network bandwidth resources are sufficient. In essence, SMP is a method that trades resources for time. After the plan parallelism is executed, resource consumption increases. When these resources become a bottleneck, SMP may deteriorate, rather than improve performance. In addition, it takes a longer time to generate SMP plans than serial plans. Therefore, in TP services that mainly involve short queries or in case resources are insufficient, you are advised to disable SMP by setting **query_dop** to **1**.

Procedure:

- 1. Observe the current system load situation. If the resource is sufficient (the resource usage ratio is smaller than 50%), perform step 2. Otherwise, exit this system.
- Set query_dop to 1 (default value). Use explain to generate an execution plan and check whether the plan can be used in scenarios described in Application Scenarios and Constraints for SMP. If the plan can be used, go to the next step.
- Set query_dop=-value. The value range of the parallelism degree is [1, value].
- 4. Set query_dop=value. The parallelism degree is 1 or value.
- 5. Before the query statement is executed, set **query_dop** to an appropriate value. After the statement is executed, set **query_dop** to **off**. For example: SET query_dop = 0, SELECT COUNT(*) FROM t1 GROUP BY a;

□ NOTE

SET query_dop = 1;

- If resources are sufficient, the higher the degree of parallelism, the better the performance improvement result.
- The SMP parallelism degree supports a session level setting and you are advised to enable SMP before executing queries that meet the requirements. After the execution is complete, disable SMP. Otherwise, SMP may affect services during peak hours.
- SMP adaptability (**query_dop** ≤ 0) depends on resource management. If resource management is disabled, only plans with parallelism degree of only 1 or 2 will be generated.

13.3.3 Configuring LLVM

LLVM dynamic compilation can be used to generate customized machine code for each query to replace original common functions. The query performance is improved by reducing redundant judgment condition and virtual function invocation, and make local data more accurate during actual queries.

LLVM needs to consume extra time to pre-generate intermediate representation (IR) and compile it into code. Therefore, if the data volume is small or if a query itself consumes little time, LLVM actually does more harm than good.

LLVM Application Scenarios and Constraints

Applicable Scenarios

- Expressions supporting LLVM. The query statements that contain the following expressions support LLVM optimization:
 - a. CASE...WHEN...
 - b. IN
 - c. Bool (AND/OR/NOT)
 - d. BooleanTest (IS_NOT_KNOWN/IS_UNKNOWN/IS_TRUE/IS_NOT_TRUE/ IS FALSE/IS NOT FALSE)
 - e. NullTest (IS_NOT_NULL/IS_NULL)
 - f. Operators
 - g. Functions (lpad, substring, btrim, rtrim, and length)
 - h. Nullif

The following data types are supported for expression calculation: bool, tinyint, smallint, int, bigint, float4, float8, numeric, date, time, timetz, timestamp, timestamptz, interval, bpchar, varchar, text, and oid.

Consider using LLVM dynamic compilation and optimization only when expressions are used in the following scenarios:

- **filter** on the **Scan** node in the case of a vectorized executor.
- complicate hash condition, hash join filter, and hash join target in the Hash Join node.
- filter and join filter in the Nested Loop node.
- merge join filter and merge join target in the Merge Join node.
- **filter** in the Group node.
- Operators that can use LLVM:
 - a. Join: HashJoin
 - b. Agg: HashAgg
 - c. Sort

Among them:

- HashJoin supports only Hash Inner Join, and the corresponding hash cond supports comparisons between int4, bigint, and bpchar.
- HashAgg supports sum and avg operations of bigint and numeric data types. Group By statements support int4, bigint, bpchar, text, varchar, timestamp, and the count(*) aggregation operation.
- Sort supports only comparisons between int4, bigint, numeric, bpchar, text, and varchar data types.

With the exception of the operations above, LLVM dynamic compilation and optimization cannot be used. To further confirm, use the explain performance tool to check.

Non-Applicable Scenarios

• LLVM dynamic compilation and optimization are not supported on CNs.

- Tables that have small amounts of data cannot be dynamically compiled using LLVM.
- Query jobs with a non-vectorized execution path cannot be generated.

Other Factors Impacting LLVM Performance

The result of LLVM optimization depends not only on operations and computation in the database, but also on the hardware environment.

• Number of C- functions invoked by query statements

CodeGen cannot be used in all expressions in an entire expression, that is, some expressions use CodeGen while others invoke original C codes for computation. In an entire expression, if more expressions invoke original C codes, LLVM dynamic compilation and optimization may reduce the computational performance. By setting <code>log_min_messages</code> to <code>DEBUG1</code>, you can check expressions that directly invoke C codes.

Memory resources

One of the key LLVM features is to ensure the locality of data, that is, data should be stored in registers whenever possible. Data loading should be reduced at the same time. Therefore, when using LLVM optimization, the value of work_mem must be set as large as required to ensure that the code is processed in the memory using LLVM. Otherwise, performance may deteriorate.

Optimizer cost estimation

The LLVM feature realizes a simple cost estimation model. You can determine whether to use LLVM dynamic compilation and optimization for the current node based on the sizes of tables involved in node computation. If the optimizer understates the actual number of rows involved, the expected performance gains may not be realized. An overestimation will have the same effect.

Recommended Usage of LLVM

LLVM is enabled in the database kernel by default, and users can configure it based on the analysis above. The overall suggestions are as follows:

- Set an appropriate value for work_mem and set it as large as possible. If much data is flushed to disks, you are advised to disable LLVM dynamic compilation and optimization by setting enable_codegen to off.
- Set an appropriate value for codegen_cost_threshold (The default value is 10,000). Ensure that LLVM dynamic compilation and optimization is not used when the data volume is small. After the value is set, if the database performance deteriorates due to the use of LLVM dynamic compilation and optimization, increase the value.
- 3. If a large number of C- functions are invoked, you are advised to disable LLVM dynamic compilation and optimization.
- 4. The constants following the **In** expression cannot exceed 10. Otherwise, LLVM compilation and optimization cannot be performed.

◯ NOTE

If resources are sufficient, the database performance will improve as the data volume increases.

13.4 SQL Tuning

13.4.1 SQL Query Execution Process

The process from receiving SQL statements to the statement execution by the SQL engine is shown in **Figure 13-2** and **Table 13-6**. The texts in red are steps where database administrators can optimize queries.

Figure 13-2 Execution process of query-related SQL statements by the SQL engine

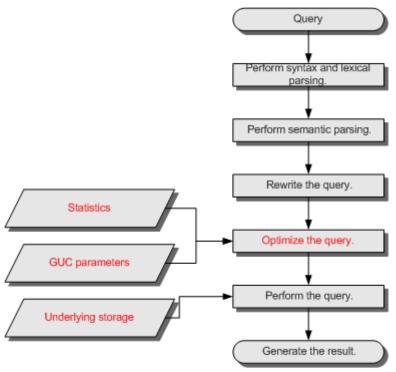


Table 13-6 Execution process of query-related SQL statements by the SQL engine

Step	Description
1. Perform syntax and lexical parsing.	Converts the input SQL statements from the string data type to the formatted structure stmt based on the specified SQL statement rules.
2. Perform semantic parsing.	Converts the formatted structure obtained from the previous step into objects that can be recognized by the database.
3. Rewrite the query statements.	Converts the output of the last step into the structure that optimizes the query execution.

Step	Description
4. Optimize the query.	Determines the execution mode of SQL statements (the execution plan) based on the result obtained from the last step and the internal database statistics. For details about the impact of statistics and GUC parameters on query optimization (execution plan), see Optimizing Queries Using Statistics and Optimizing Queries Using GUC parameters.
5. Perform the query.	Executes the SQL statements based on the execution path specified in the last step. Selecting a proper underlying storage mode improves the query execution efficiency. For details, see Optimizing Queries Using the Underlying Storage.

Optimizing Queries Using Statistics

The GaussDB(DWS) optimizer is a typical Cost-based Optimization (CBO). The database uses the CBO to calculate the number of tuples and execution cost for each execution step in every execution plan. This calculation is based on factors such as the number of table tuples, column width, NULL record ratio, and characteristic values (such as distinct, MCV, and HB values) using specific cost calculation methods. The database then selects the execution plan with the lowest cost for overall execution or for returning the first tuple. These characteristic values are the statistics, which is the core for optimizing a query. Accurate statistics helps the optimizer select the most appropriate query plan. Generally, you can collect statistics of a table or that of some columns in a table using **ANALYZE**. You are advised to periodically execute **ANALYZE** or execute it immediately after you modified most contents in a table.

Optimizing Queries Using GUC parameters

Optimizing queries aims to select an efficient execution mode.

Take the following statement as an example:

```
SELECT count(1)
FROM customer inner join store_sales on (ss_customer_sk = c_customer_sk);
```

During execution of **customer inner join store_sales**, GaussDB(DWS) supports nested loop, merge join, and hash join. The optimizer estimates the result set value and the execution cost under each join mode based on the statistics of the **customer** and **store_sales** tables and selects the execution plan that takes the lowest execution cost.

As described in the preceding content, the execution cost is calculated based on certain methods and statistics. If the actual execution cost cannot be accurately estimated, you need to optimize the execution plan by setting the GUC parameters.

Optimizing Queries Using the Underlying Storage

GaussDB(DWS) supports both row-store and column-store tables. The choice of storage mode ultimately depends on your business needs. Column-store tables are

suitable for computing services that mainly involve associations and aggregations. Row-store tables are better suited for point queries and large-scale updates or deletions.

Optimization methods of each storage mode will be described in details in the performance optimization chapter.

Optimizing Queries by Rewriting SQL Statements

Besides the preceding methods that improve the performance of the execution plan generated by the SQL engine, database administrators can also enhance SQL statement performance by rewriting SQL statements while retaining the original service logic based on the execution mechanism of the database and abundant practical experience.

This requires that the system administrators know the customer business well and have professional knowledge of SQL statements.

13.4.2 SQL Execution Plan

An SQL execution plan is a node tree that displays the detailed steps performed when the GaussDB(DWS) executes an SQL statement.

You can run the **EXPLAIN** command to view the execution plan generated for each query by an optimizer. **EXPLAIN** outputs a row of information for each execution node, showing the basic node type and the expense estimate that the optimizer makes for executing the node.

Execution Plan Information

In addition to setting different display formats for an execution plan, you can use different **EXPLAIN** syntax to display execution plan information in detail. The common usages are as follows. For more usages, see **EXPLAIN Syntax**.

- EXPLAIN *statement*: only generates an execution plan and does not execute. The *statement* indicates SQL statements.
- EXPLAIN ANALYZE statement: generates and executes an execution plan, and displays the execution summary. Then actual execution time statistics are added to the display, including the total elapsed time expended within each plan node (in milliseconds) and the total number of rows it actually returned.
- EXPLAIN PERFORMANCE *statement*: generates and executes the execution plan, and displays all execution information.

To measure the run time cost of each node in the execution plan, the current execution of **EXPLAIN ANALYZE** or **EXPLAIN PERFORMANCE** adds profiling overhead to query execution. Running **EXPLAIN ANALYZE** or **EXPLAIN PERFORMANCE** on a query sometimes takes more time than a normal query. The amount of overhead depends on the nature of the query, as well as the platform being used.

Therefore, if an SQL statement is not finished after being running for a long time, run the **EXPLAIN** statement to view the execution plan and then locate the fault. If the SQL statement has been properly executed, run the **EXPLAIN ANALYZE** or **EXPLAIN PERFORMANCE** statement to check the execution plan and information to locate the fault.

Description of common execution plan keywords:

1. Table access modes

Seq Scan/CStore Scan

Scans all rows of the table in sequence. These are basic scan operators, which are used to scan row-store and column-store tables in sequence.

Index Scan/CStore Index Scan

Scans indexes of row-store and column-store tables. There are indexes in row-store or column-store tables, and the condition column is the index column

The optimizer uses a two-step plan: the child plan node visits an index to find the locations of rows matching the index condition, and then the upper plan node actually fetches those rows from the table itself. Fetching rows separately is much more expensive than reading them sequentially, but because not all pages of the table have to be visited, this is still cheaper than a sequential scan. The upper-layer planning node first sort the location of index identifier rows based on physical locations before reading them. This minimizes the independent capturing overhead.

If there are separate indexes on multiple columns referenced in **WHERE**, the optimizer might choose to use an **AND** or **OR** combination of the indexes. However, this requires the visiting of both indexes, so it is not necessarily a win compared to using just one index and treating the other condition as a filter.

The following Index scans featured with different sorting mechanisms are involved:

Bitmap Index Scan

To use a bitmap index to capture a data page, you need to scan the index to obtain the bitmap and then scan the base table.

Index Scan using index name

Fetches table rows in index order, which makes them even more expensive to read. However, there are so few rows that the extra cost of sorting the row locations is unnecessary. This plan type is used mainly for queries fetching just a single row and queries having an **ORDER BY** condition that matches the index order, because no extra sorting step is needed to satisfy **ORDER BY**.

2. Table connection modes

Nested Loop

Nested-loop is used for queries that have a smaller data set connected. In a Nested-loop join, the foreign table drives the internal table and each row returned from the foreign table should have a matching row in the internal table. The returned result set of all queries should not exceed 10,000. The table that returns a smaller subset will work as a foreign table, and indexes are recommended for connection fields of the internal table.

- (Sonic) Hash Join

A Hash join is used for large tables. The optimizer uses a hash join, in which rows of one table are entered into an in-memory hash table, after which the other table is scanned and the hash table is probed for

matches to each row. Sonic and non-Sonic hash joins differ in their hash table structures, which do not affect the execution result set.

Merge Join

In a merge join, data in the two joined tables is sorted by join columns. Then, data is extracted from the two tables to a sorted table for matching.

Merge join requires more resources for sorting and its performance is lower than that of hash join. If the source data has been sorted, it does not need to be sorted again when merge join is performed. In this case, the performance of merge join is better than that of hash join.

3. Operators

sort

Sorts the result set.

filter

The **EXPLAIN** output shows the **WHERE** clause being applied as a **Filter** condition attached to the **Seq Scan** plan node. This means that the plan node checks the condition for each row it scans, and returns only the ones that meet the condition. The estimated number of output rows has been reduced because of the **WHERE** clause. However, the scan will still have to visit all 10000 rows. As a result, the cost is not decreased. It increases a bit (by 10000 x **cpu_operator_cost**) to reflect the extra CPU time spent on checking the **WHERE** condition.

LIMIT

LIMIT limits the number of output execution results. If a **LIMIT** condition is added, not all rows are retrieved.

Execution Plan Display Format

GaussDB(DWS) provides four display formats: **normal**, **pretty**, **summary**, and **run**. You can change the display format of execution plans by setting **explain_perf_mode**.

 normal indicates that the default printing format is used. Figure 13-3 shows the display format.

Figure 13-3 Example of an execution plan in normal format

pretty indicates that the optimized display mode of GaussDB(DWS) is used. A
new format contains a plan node ID, directly and effectively analyzing
performance. Figure 13-4 is an example.

Figure 13-4 Example of an execution plan using the pretty format

- summary indicates that the analysis result based on such information is printed in addition to the printed information in the format specified by pretty.
- **run** indicates that in addition to the printed information specified by **summary**, the database exports the information as a CSV file.

Common Types of Plans

GaussDB(DWS) has three types of distributed plans:

- Fast Query Shipping (FQS) plan
 - The CN directly delivers statements to DNs. Each DN executes the statements independently and summarizes the execution results on the CN.
- Stream plan
 - The CN generates a plan for the statements to be executed and delivers the plan to DNs for execution. During the execution, DNs use the Stream operator to exchange data.
- Remote-Query plan

After generating a plan, the CN delivers some statements to DNs. Each DN executes the statements independently and sends the execution result to the CN. The CN executes the remaining statements in the plan.

The existing tables **tt01** and **tt02** are defined as follows:

```
CREATE TABLE tt01(c1 int, c2 int) DISTRIBUTE BY hash(c1);
CREATE TABLE tt02(c1 int, c2 int) DISTRIBUTE BY hash(c2);
```

Type 1: FQS plan, all statements pushed down

Two tables are joined, and the join condition is the distribution column of each table. If the stream operator is disabled, the CN directly sends statements to each DN for execution. The result is summarized on the CN.

```
SET enable_stream_operator=off;
SET explain_perf_mode=normal;

EXPLAIN (VERBOSE on,COSTS off) SELECT * FROM tt01,tt02 WHERE tt01.c1=tt02.c2;

QUERY PLAN

Data Node Scan on "__REMOTE_FQS_QUERY__"
Output: tt01.c1, tt01.c2, tt02.c1, tt02.c2
Node/s: All datanodes
Remote query: SELECT tt01.c1, tt01.c2, tt02.c1, tt02.c2 FROM dbadmin.tt01, dbadmin.tt02 WHERE tt01.c1

= tt02.c2
(4 rows)
```

Type 2: Non-FQS plan, some statements pushed down

Two tables are joined and the join condition contains non-distribution columns. If the stream operator is disabled, the CN delivers the base table scanning statements to each DN. Then, the JOIN operation is performed on the CN.

```
SET enable_stream_operator=off;
SET explain_perf_mode=normal;
EXPLAIN (VERBOSE on, COSTS off) SELECT * FROM tt01, tt02 WHERE tt01.c1=tt02.c1;
                    QUERY PLAN
Hash Join
 Output: tt01.c1, tt01.c2, tt02.c1, tt02.c2
 Hash Cond: (tt01.c1 = tt02.c1)
 -> Data Node Scan on tt01 "_REMOTE_TABLE_QUERY_"
     Output: tt01.c1, tt01.c2
     Node/s: All datanodes
     Remote query: SELECT c1, c2 FROM ONLY dbadmin.tt01 WHERE true
     Output: tt02.c1, tt02.c2
     -> Data Node Scan on tt02 "_REMOTE_TABLE_QUERY_"
         Output: tt02.c1, tt02.c2
         Node/s: All datanodes
         Remote query: SELECT c1, c2 FROM ONLY dbadmin.tt02 WHERE true
(13 rows)
```

Type 3: Stream plan, no data exchange between DNs

Two tables are joined, and the join condition is the distribution column of each table. DNs do not need to exchange data. After generating a stream plan, the CN delivers the plan except the Gather Stream part to DNs for execution. The CN scans the base table on each DN, performs hash join, and sends the result to the CN.

```
SET enable fast query shipping=off;
SET enable_stream_operator=on;
EXPLAIN (VERBOSE on, COSTS off) SELECT * FROM tt01, tt02 WHERE tt01.c1=tt02.c2;
              OUERY PLAN
Streaming (type: GATHER)
  Output: tt01.c1, tt01.c2, tt02.c1, tt02.c2
  Node/s: All datanodes
  -> Hash Join
      Output: tt01.c1, tt01.c2, tt02.c1, tt02.c2
      Hash Cond: (tt01.c1 = tt02.c2)
      -> Seg Scan on dbadmin.tt01
         Output: tt01.c1, tt01.c2
         Distribute Key: tt01.c1
      -> Hash
          Output: tt02.c1, tt02.c2
          -> Seq Scan on dbadmin.tt02
              Output: tt02.c1, tt02.c2
              Distribute Key: tt02.c2
(14 rows)
```

Type 4: Stream plan, with data exchange between DNs

When two tables are joined and the join condition contains non-distribution columns, and the stream operator is enabled (SET enable_stream_operator=on), a stream plan is generated, which allows data exchange between DNs. For table **tt02**, the base table is scanned on each DN. After the scanning, the **Redistribute Stream** operator performs hash calculation based on **tt02.c1** in the **JOIN** condition, sends the hash calculation result to each DN, and then performs JOIN on each DN, finally, the data is summarized to the CN.

Type 5: Remote-Query plan

unship_func cannot be pushed down and does not meet partial pushdown requirements (subquery pushdown). Therefore, you can only send base table scanning statements to DNs and collect base table data to the CN for calculation.

```
postgres=> CREATE FUNCTION unship_func(integer,integer) returns integer
postgres-> AS 'select $1 + $2;'
postgres-> LANGUAGE SQL volatile
postgres-> returns null on null input;
CREATE FUNCTION
```

```
stgres=> SET explain_perf_mode=pretty;
er
ostgres=> EXPLAIN VERBOSE SELECT unship_func(tt01.c1,tt01.c2) FROM tt01 J0IN tt02 on tt01.c1=tt02.c1;
QUERY PLAN
                                                                                               | E-rows | E-distinct | E-width | E-costs
 id |
                                           operation
         -> Hash Join (2,3)
-> Data Node Scan on tt01 "_REMOTE_TABLE_QUERY_"
-> Hash
                                                                                                         30
30
30
30
                 -> Data Node Scan on tt02 "_REMOTE_TABLE_QUERY_"
                SQL Diagnostic Information
SQL is not plan-shipping reason: Function unship_func() can not be shipped
Predicate Information (identified by plan id)
  1 --Hash Join (2,3)
Hash Cond: (tt01.c1 = tt02.c1)
               Targetlist Information (identified by plan id)
  1 --Hash Join (2,3)
Output: (tt01.c1 + tt01.c2)
2 --Data Node Scan on tt01 "_REMOTE_TABLE_QUERY_"
Output: tt01.c1, tt01.c2
Node/s: All datanodes
Remote query: SELECT c1, c2 FROM ONLY dbadmin.tt01 WHERE true
  3 --Hash
Output: tt02.c1
4 --Data Node Scan on tt02 "_REMOTE_TABLE_QUERY_"
Output: tt02.c1
Node/s: All datanodes
Remote query: SELECT c1 FROM ONLY dbadmin.tt02 WHERE true
Parser runtime: 0.055 ms
Planner runtime: 0.528 ms
Unique SQL Id: 1780774145
37 rows)
```

EXPLAIN PERFORMANCE Description

You can use **EXPLAIN ANALYZE** or **EXPLAIN PERFORMANCE** to check the SQL statement execution information and compare the actual execution and the optimizer's estimation to find what to optimize. **EXPLAIN PERFORMANCE** provides the execution information on each DN, whereas **EXPLAIN ANALYZE** does not.

Tables are defined as follows:

```
CREATE TABLE tt01(c1 int, c2 int) DISTRIBUTE BY hash(c1);
CREATE TABLE tt02(c1 int, c2 int) DISTRIBUTE BY hash(c2);
```

The following SQL query statement is used as an example:

SELECT * FROM tt01,tt02 WHERE tt01.c1=tt02.c2;

The output of EXPLAIN PERFORMANCE consists of the following parts:

Execution Plan



The plan is displayed as a table, which contains 11 columns: id, operation, Atime, A-rows, E-rows, E-distinct, Peak Memory, E-memory, A-width, E-width, and E-costs. Table 13-7 describes the columns.

Table 13-7 Execution column description

Column	Description	
id	ID of an execution operator.	
operation	Name of an execution operator. The operator of the Vector prefix refers to a vectorized execution engine operator, which exists in a query containing a column-store table. Streaming is a special operator. It implements the core data shuffle function of the distributed architecture. Streaming has three types, which correspond to different data shuffle functions in the distributed architecture:	
	 Streaming (type: GATHER): The CN collects data from DNs. Streaming(type: REDISTRIBUTE): Data is redistributed to all the DNs based on selected columns. Streaming(type: BROADCAST): Data on the current DN is broadcast to all other DNs. 	

Column	Description	
A-time	Execution time of an operator on each DN. Generally, Atime of an operator is two values enclosed by square brackets ([]), indicating the shortest and longest time for completing the operator on all DNs, including the execution time of the lower-layer operators.	
	Note: In the entire plan, the execution time of a leaf node is the execution time of the operator, while the execution time of other operators includes the execution time of its subnodes.	
A-rows	Actual rows output by an operator.	
E-rows	Estimated rows output by each operator.	
E-distinct	Estimated distinct value of the hashjoin operator.	
Peak Memory	Peak memory used when the operator is executed on each DN. The left value in [] is the minimum value, and the right value in [] is the maximum value.	
E-memory	Estimated memory used by each operator on a DN. Only operators executed on DNs are displayed. In certain scenarios, the memory upper limit enclosed in parentheses will be displayed following the estimated memory usage.	
A-width	The actual width of each line of tuple of the current operator. This parameter is valid only for the heavy memory operator is displayed, including: (Vec)HashJoin, (Vec)HashAgg, (Vec) HashSetOp, (Vec)Sort, and (Vec)Materialize operator. The (Vec)HashJoin calculation of width is the width of the right subtree operator, it will be displayed in the right subtree.	
E-width	Estimated width of the output tuple of each operator.	
E-costs	 Estimated execution cost of each operator. E-costs are defined by the optimizer based on cost parameters, habitually grasping disk page as a unit. Other overhead parameters are set by referring to E-costs. The cost of each node (the E-costs value) includes the cost of all of its child nodes. Overhead reflects only what the optimizer is concerned about, but does not consider the time that the result row passed to the client. Although the time may play a very important role in the actual total time, it is ignored by the optimizer, because it cannot be changed by modifying the plan. 	

2. SQL Diagnostic Information

SQL self-diagnosis information. Performance optimization points identified during optimization and execution are displayed. When **EXPLAIN** with the **VERBOSE** attribute (built-in **VERBOSE** of **EXPLAIN PERFORMANCE**) is executed on DML statements, SQL self-diagnosis information is also generated to help locate performance issues.

3. Predicate Information (identified by plan id)

```
Predicate Information (identified by plan id)

2 --Hash Join (3,4)

Hash Cond: (tt01.cl = tt02.c2)

3 --Seq Scan on dbadmin.tt01

Filter: (tt01.cl >= 202007)

5 --Seq Scan on dbadmin.tt02

Filter: (tt02.c2 >= 202007)
```

This part displays the filtering conditions of the corresponding execution operator node, that is, the information that does not change during the entire plan execution, mainly the join conditions and filter information.

8.3.0 and later cluster versions support the display the information of **CU Predicate Filter** and **Pushdown Predicate Filter(will be pruned)** related to dictionary plans.

4. Memory Information (identified by plan id)

```
Memory Information (identified by plan id)
Coordinator Query Peak Memory:
         Query Peak Memory: 2MB
DataNode Query Peak Memory
         dn 6001 6002 Query Peak Memory: OMB
         dn_6003_6004 Query Peak Memory: OMB
  dn_6005_6006 Query Peak Memory: 0MB
1 --Streaming (type: GATHER)
         Peak Memory: 56KB, Estimate Memory: 512MB
  2 -- Hash Join (3,4)
         dn_6001_6002 Peak Memory: 8KB, Estimate Memory: 1024KB
         dn_6003_6004 Peak Memory: 8KB, Estimate Memory: 1024KB
dn_6005_6006 Peak Memory: 8KB, Estimate Memory: 1024KB
  3 --Seq Scan on dbadmin.tt01
         dn_6001_6002 Peak Memory: 32KB, Estimate Memory: 1024KB
         dn_6003_6004 Peak Memory: 32KB, Estimate Memory: 1024KB
dn_6005_6006 Peak Memory: 32KB, Estimate Memory: 1024KB
  4 --Hash
         dn 6001 6002 Buckets: 0 Batches: 0
                                                     Memory Usage: 0kB
         dn 6003 6004 Buckets: 0
                                      Batches: 0
                                                     Memory Usage: 0kB
         dn_6005_6006 Buckets: 0 Batches: 0
                                                     Memory Usage: 0kB
```

The memory usage part displays the memory usage information printed by certain operators (mainly Hash and Sort), including **peak memory**, **estimate memory**, **control memory**, **operator memory**, **width**, **auto spread num**, and **early spilled**; and spill details, including **spill Time(s)**, **inner/outer partition spill num**, **temp file num**, spilled data volume, and **written disk I/O** [*min*, *max*]. The Sort operator does not display the number of files written to disks, and displays disks only when displaying sorting methods.

5. Targetlist Information (identified by plan id)

```
Targetlist Information (identified by plan id)

1 --Streaming (type: GATHER)
        Output: tt01.c1, tt01.c2, tt02.c1, tt02.c2
        Node/s: All datanodes

2 --Hash Join (3,4)
        Output: tt01.c1, tt01.c2, tt02.c1, tt02.c2

3 --Seq Scan on dbadmin.tt01
        Output: tt01.c1, tt01.c2
        Distribute Key: tt01.c1

4 --Hash
        Output: tt02.c1, tt02.c2

5 --Seq Scan on dbadmin.tt02
        Output: tt02.c1, tt02.c2
        Distribute Key: tt02.c2
```

This part displays the output target column information of each operator.

In 8.3.0 and later cluster versions, the dictionary parameters **Dict Optimized** and **Dict Decoded** can be displayed, indicating dictionary columns and dictionary codes, respectively.

6. DataNode Information (identified by plan id)

This part displays the execution time of each operator (including the execution time of filtering and projection, if any), CPU usage, and buffer usage.

Operator execution information

```
4 --Vector Aggregate
dn_6001_6002 (actual time=54076.307..54076.308 rows=1 loops=1) (projection time=65.581)
dn_6003_6004 (actual time=68597.121..68597.122 rows=1 loops=1) (projection time=64.801)
```

The execution information of each operator consists of three parts:

- dn_6001_6002/dn_6003_6004 indicates the information about the execution node. The information in the brackets is the actual execution information.
- actual time indicates the actual execution time. The first number indicates the duration from the time when the operator is executed to the time when the first data record is output. The second number indicates the total execution time of all data records.
- rows indicates the number of output data rows of the operator.

- **loops** indicates the number of execution times of the operator. Note that for a partitioned table, scan on each partition is counted as a scan. Scan on a new partition is counted as a new scan.
- CPU information

```
dn 6001 6002 (CPU: ex c/r=44, ex row=29999729, ex cyc=1346672758, inc cyc=162228870602) dn 6003 6004 (CPU: ex c/r=44, ex row=29999728, ex cyc=1343711688, inc cyc=205791299646)
```

Each operator execution process has CPU information. **cyc** indicates the number of CPU cycles, and **ex cyc** indicates the number of cycles of the current operator, excluding its subnodes. **inc cyc** indicates the number of cycles, including subnodes, **ex row** indicates the number of data rows output by the current operator, and **ex c/r** indicates the mean of **ex cyc** and **ex row**.

Buffer information

```
dn_6001_6002 (Buffers: shared hit=1725 dirtied=117)
dn_6003_6004 (Buffers: shared hit=1723 dirtied=117)
```

Buffers indicates the buffer information, including the read and write operations on shared blocks and temporary blocks.

Shared blocks contain tables and indexes, and temporary blocks are disk blocks used in sorting and materialization. The number of blocks displayed on the upper-layer node contains the number of blocks used by all its subnodes.

 Disk cache information (supported only by V3 tables or foreign tables with decoupled storage and compute and colversion set to 3.0 in 9.1.0.100 or later)

```
dn_6001_6002 (Disk Cache: hit=146611 miss=47 error=47 errorCode=152 scanBytes=18325.02MB remoteReadBytes=959.00MB loadTime=13109.5) (Column 3.0: preloadStep=20 preloadSubmitTime=1230.6 preloadMaitTime=1724.9 preloadMaitCount=28) (0BS 1/0: count=349 averageRTI=238.7 averageLatency=233.0) (Disk Cache: hit=140641 miss=45 error=45 errorCode=152 scanBytes=17508.64MB remoteReadBytes=902.00MB loadTime=13271.7) (Column 3.0: preloadStep=20 preloadSubmitTime=1202.9 preloadMaitTime=3061.5 preloadMaitCount=20) (0BS 1/0: count=327 averageRTI=231.5 averageLatency=263.3 latencyGtls=1) dn_6005_6006 (Disk Cache: hit=130736 miss=47 error=47 errorCode=152 scanBytes=17494.91MB remoteReadBytes=909.00MB loadTime=12752.5) (Column 3.0: preloadStep=20 preloadSubmitTime=1109.4 preloadMaitTime=2723.4 preloadMaitCount=13) (0BS 1/0: count=327 averageRTI=194.3 averageLatency=188.8 latencyGtls=3)
```

Disk Cache indicates the hit information and data read information of the disk cache. (supported by V3 tables or foreign tables with storage and compute decoupled)

miss indicates the number of disk cache misses. **hit** indicates the number of disk cache hits. For details about errorCode, see

disk_cache_error_code in Table 14-158. error indicates the number of times errorCode is generated. scanBytes indicates the amount of data queried by scan, remoteReadBytes indicates the amount of data read on OBS, and loadTime indicates the time for loading data from the disk cache. To improve OBS efficiency, adjacent request blocks are combined. Alternatively, the minimum granularity for writing requests to the disk cache is block (1 MB by default). As a result, the value of scanBytes may be less than that of remoteReadBytes.

Column 3.0: prefetch parameters and prefetch process information of V3 tables with storage and compute decoupled. (This parameter is supported only by V3 tables with storage and compute decoupled and is displayed after prefetch parameters are enabled.)

preloadStep indicates the prefetch step, **preloadSubmitTime** indicates the I/O request submission time in the prefetch process, **preloadWaitTime** indicates the I/O request waiting time in the prefetch

process, and **preloadWaitCount** indicates the number of waiting I/O requests in the prefetch process.

OBS I/O indicates details about an OBS I/O request. (supported by V3 tables or foreign tables with storage and compute decoupled)

count indicates the total number of OBS I/O requests. averageRTT indicates the average round trip time (RTT) of OBS I/O requests. The unit is μ s. averageLatency indicates the average latency of OBS I/O requests. The unit is μ s. latencyGt1s indicates the number of OBS I/O requests whose latency exceeds 1s. latencyGt10s indicates the number of OBS I/O requests whose latency exceeds 10s. retryCount indicates the total number of OBS I/O request retries. rateLimitCount indicates the total number of times that OBS I/O requests are under flow control.

7. User Define Profiling

```
User Define Profiling

Plan Node id: 1 Track name: coordinator get datanode connection cn_5001 (time=9.306 total_calls=1 loops=1)

Plan Node id: 1 Track name: coordinator begin transaction cn_5001 (time=0.002 total_calls=1 loops=1)

Plan Node id: 1 Track name: coordinator send command cn_5001 (time=0.113 total_calls=3 loops=1)

Plan Node id: 1 Track name: coordinator get the first tuple cn_5001 (time=0.091 total_calls=12 loops=1)
```

User-defined information, including the time when CNs and DNs are connected, the time when DNs are connected, and some execution information at the storage layer.

8. Query Summary

```
Datanode executor start time [dn_6005_6006, dn_6001_6002]: [0.360 ms,0.483 ms]
Datanode executor run time [dn_6001_6002, dn_6003_6004]: [0.008 ms,0.009 ms]
Datanode executor end time [dn_6003_6004, dn_6005_6006]: [0.036 ms,0.066 ms]
Remote query poll time: 2.649 ms, Deserialze time: 0.000 ms
System available mem: 1761280KB
Query Max mem: 1761280KB
Query estimated mem: 328KB
Enqueue time: 0.030 ms
Coordinator executor start time: 0.083 ms
Coordinator executor run time: 13.044 ms
Coordinator executor end time: 0.034 ms
Parser runtime: 0.060 ms
Planner runtime: 0.539 ms
Query Id: 218706056932222840
Unique SQL Id: 2641724793
Total runtime: 13.906 ms
```

The total execution time and network traffic, including the maximum and minimum execution time in the initialization and end phases on each DN, initialization, execution, and time in the end phase on each CN, and the system available memory during the current statement execution, and statement estimation memory information.

- DataNode executor start time: start time of the DN executor. The format is [min_node_name, max_node_name]: [min_time, max_time].
- DataNode executor run time: running time of the DN executor. The format is [min node name, max node name]: [min time, max time].
- DataNode executor end time: end time of the DN executor. The format is [min_node_name, max_node_name]: [min_time, max_time].

- Remote query poll time: poll waiting time for receiving results
- **System available mem**: available system memory
- Query Max mem: maximum query memory.
- **Enqueue time**: enqueuing time
- Coordinator executor start time: start time of the CN executor
- Coordinator executor run time: CN executor running time
- **Coordinator executor end time**: end time of the CN executor
- **Parser runtime**: parser running time
- **Planner runtime**: optimizer execution time
- Network traffic, or, the amount of data sent by the stream operator
- **Query Id**: query ID.
- Unique SQL ID: constraint SQL ID
- Total runtime: total execution time

NOTICE

- The difference between A-rows and E-rows shows the deviation between the optimizer estimation and actual execution. Generally, if the deviation is large, the plan generated by the optimizer cannot be trusted, and you need to modify the deviation value.
- If the difference of the A-time values is large, it indicates that the operator computing skew (difference between execution time on DNs) is large and that manual performance tuning is required. Generally, for two adjacent operators, the execution time of the upper-layer operator includes that of the lower-layer operator. However, if the upper-layer operator is a stream operator, its execution time may be less than that of the lower-layer operator, as there is no driving relationship between threads.
- Max Query Peak Memory is often used to estimate the consumed memory of SQL statements, and is also used as an important basis for setting a memory parameter during SQL statement optimization. Generally, the output from EXPLAIN ANALYZE or EXPLAIN PERFORMANCE is provided for the input for further optimization.

13.4.3 Execution Plan Operator

Operator Introduction

In an SQL execution plan, each step indicates a database operator, also called an execution operator. In GaussDB(DWS), operators are the building blocks of data processing. By combining them effectively and optimizing their sequence and execution, you can significantly improve data processing efficiency.

GaussDB(DWS) operators are classified into scan operators, control operators, materialization operators, join operators, and other operators.

Scan Operators

A scan operator scans data in a table, processing one tuple at a time for the upper-layer node. It operates at the leaf node of the query plan tree and can scan tables, result sets, linked lists, and subquery results. The following table lists common scan operators.

Table 13-8 Scan operators

Operator	Description	Scenario
SeqScan	Sequential scanning	It is a basic operator used to scan physical tables in sequence, not an index-assisted scan.
IndexScan	Index scanning	Indexes are created for the attributes involved in selection conditions.
IndexOnlySca n	Obtaining a tuple from an index	The index column completely overwrites the result set column.
BitmapScan(B itmapIndexSc an, BitmapHeapS can)	Obtaining a tuple using a bitmap	BitmapIndexScan uses indexes for attributes to scan data and returns a bitmap. BitmapHeapScan then uses this bitmap to retrieve tuples.
TidScan	Obtaining a tuple by tuple tid	1. WHERE conditions(like CTID = tid or CTID IN (tid1, tid2,)); 2. UPDATE/DELETE WHERE CURRENT OF cursor;
SubqueryScan	Subquery scanning	Another query plan tree (subplan) is used as the scanning object to scan tuples.
FunctionScan	Function scanning	FROM function_name
ValuesScan	Values linked list scanning	It scans the given tuple set in VALUES clauses.
ForeignScan	External table scanning	It queries external tables.
CteScan	CTE table scanning	It scans the subquery defined by the WITH clause in the SELECT query.

Join Operators

In relational algebra, a join operation is equivalent to a join operator. Take a simple example: joining two tables, t1 and t2. There are several types of joins, including inner join, left join, right join, full join, semi join, and anti join. These joins can be implemented using three methods: Nestloop, HashJoin, and MergeJoin.

Table 13-9 Join operators

Operator	Description	Scenario	Implementation Feature
NestLoop	Nested loop join, which is a brute force approach. It scans the inner table for each row.	Inner Join, Left Outer Join, Semi Join, Anti Join	It is used for queries that have a smaller subset connected. In a nested loop, the foreign table drives the internal table. Each row returned by the foreign table is retrieved from the internal table to find the matched row. Therefore, the result set returned by the entire query cannot be greater than 10,000. The table with a smaller subset returned is used as the foreign table. It is recommended that indexes be created for the join fields in the internal table.
MergeJoi n	A merge join on ordered input sorts both the inner and outer tables, identifies the first and last matching rows, and then joins tuples at a time. Equijoin.	Inner Join, Left Outer Join, Right Outer Join, Full Outer Join, Semi Join, Anti Join	In a merge join, data in the two joined tables is sorted by join columns. Then, data is extracted from the two tables to a sorted table for matching. A merge join requires more resources for sorting and its performance is lower than that of a hash join. However, if the source data has been pre-sorted and no more sorting is needed during the merge join, its performance excels.
(Sonic) Hash Join	Hash join: The inner and outer tables use the join column's hash value to create a hash table. Matching values are then stored in the same bucket. The two ends of an equal join must be of the same type and support hash.	Inner Join, Left Outer Join, Right Outer Join, Full Outer Join, Semi Join, Anti Join	A hash Join is used for large tables. The optimizer creates a hash table in memory using the join key and the smaller table. It then scans the larger table and uses the hash table to quickly identify matching rows. While Sonic and non-Sonic hash joins have different internal structures, this does not impact the final result set.

Materialization Operators

Materialization operators are a class of nodes that can cache tuples. During execution, many extended physical operations can be performed only after all tuples are obtained, such as aggregation function operations and sorting without indexes. Materialization operators can cache all the tuples.

Table 13-10 Materialization operators

Operator	Description	Scenario
Material	Materializatio n	Caches the subnode result.
Sort	Sorting	ORDER BY clause, which is used for join, group, and set operations and works with Unique.
Group	Grouping	GROUP BY clause.
Agg	Executes aggregate functions.	 Aggregate functions such as COUNT, SUM, AVG, MAX, and MIN. DISTINCT clause. UNION deduplication.
		4. GROUP BY clause.
WindowAgg	Window functions	WINDOW clause.
Unique	Deduplication (with sorted lower-layer data)	 DISTINCT clause. UNION deduplication.
Hash	HashJoin auxiliary node	Constructs a hash table and use it together with HashJoin.
SetOp	Processing set operations	INTERSECT/INTERSECT ALL, EXCEPT/EXCEPT ALL
LockRows	Processing row-level locks	SELECT FOR SHARE/UPDATE

Control Operators

Control operators are a type of node that handles exceptional scenarios and executes custom workflows.

Table 13-11 Control operators

Operator	Description	Scenario
Result	Performing calculation directly	 Table scanning is not included. The INSERT statement contains only one VALUES clause.
ModifyTable	INSERT/UPDATE/ DELETE upper-layer node	INSERT, UPDATE, and DELETE
Append	Appending	1. UNION(ALL) 2. Table inheritance
MergeAppend	Appending (ordered input)	1. UNION(ALL) 2. Table inheritance
RecursiveUnio n	Processing the UNION subquery defined recursively in the WITH clause	WITH RECURSIVE SELECT statement
BitmapAnd	Bitmap logical AND operation	BitmapScan for multi-dimensional index scanning
BitmapOr	Bitmap logical OR operation	BitmapScan for multi-dimensional index scanning
Limit	Processing the LIMIT clause	OFFSET LIMIT

Other Operators

Other operators include Stream and RemoteQuery. There are three types of Stream operators: Gather stream, Broadcast stream, and Redistribute stream.

- Gather stream: Each source node sends its data to the target node for aggregation.
- Broadcast stream: A source node sends its data to N target nodes for calculation.
- Redistribute stream: Each source node calculates the hash value of its data based on the join condition, distributes the data based on the hash value, and sends the data to the corresponding target node.

Table 13-12 Other Operators

Operator	Description	Scenario
Stream	Multi-node data exchange	When a distributed query plan is executed, data is exchanged between nodes.

Operator	Description	Scenario
Partition Iterator	Partition iterator	Scans partitioned tables and iteratively scans each partition.
RowToVec	Rows-to- column conversion	Hybrid row-column.
DfsScan / DfsIndexScan	HDFS table (index) scanning	HDFS table scanning.

13.4.4 SQL Tuning Process

You can analyze slow SQL statements to optimize them.

Procedure

- Step 1 Collect all table statistics associated with the SQL statements. In a database, statistics indicate the source data of a plan generated by a planner. If statistics are unavailable or out of date, the execution plan may seriously deteriorate, leading to low performance. According to past experience, about 10% performance problem occurred because no statistics are collected. For details, see Updating Statistics.
- **Step 2** Review and modify the table definition.
- Step 3 Generally, some SQL statements can be converted to its equivalent statements in all or certain scenarios by rewriting queries. SQL statements are simpler after they are rewritten. Some execution steps can be simplified to improve the performance. The query rewriting method is universal in all databases. SQL Statement Rewriting Rules describes several optimization methods by rewriting SQL statements.
- **Step 4** View the execution plan to find out the cause. If the SQL statements have been running for a long period of time and not ended, run the **EXPLAIN** command to view the execution plan and then locate the fault. If the SQL statement has been executed, run the **EXPLAIN ANALYZE** or **EXPLAIN PERFORMANCE** command to check the execution plan and actual running situation and then accurately locate the fault. For details about the execution plan, see **SQL Execution Plan**.
- Step 5 For details about EXPLAIN or EXPLAIN PERFORMANCE, the reason why SQL statements are slowly located, and how to solve this problem, see Advanced SQL Tuning.
- **Step 6** Specify a join order; join, stream, or scan operations; number of rows in a result; or redistribution skew information to optimize an execution plan, improving query performance. For details, see **Hint-based Tuning**.
- **Step 7** To maintain high database performance, you are advised to perform **Routinely Maintaining Tables** and **Routinely Recreating an Index**.

Step 8 (Optional) Improve performance by using operators if resources are sufficient in GaussDB(DWS). For details, see **SMP Parallel Execution**.

----End

13.4.5 Updating Statistics

In a database, statistics indicate the source data of a plan generated by a planner. If statistics are unavailable or out of date, the execution plan may seriously deteriorate, leading to low performance.

Scenario

The **ANALYZE** statement collects statistics on database table contents. These statistics will be stored in the **PG_STATISTIC** system catalog. Then, the query optimizer uses the statistics to work out the most efficient execution plan.

After executing batch **INSERT** and **DELETE** operations, you are advised to run the **ANALYZE** statement on the table or the entire database to update statistics. By default, 30,000 rows of statistics are sampled. That is, the default value of the GUC parameter **default_statistics_target** is **100**. If the total number of rows in the table exceeds 1,600,000, you are advised to set **default_statistics_target** to **-2**, indicating that 2% of the statistics are collected.

For an intermediate table generated during the execution of scripts or stored procedures in batch, you also need to run the **ANALYZE** statement.

If there are multiple inter-related columns in a table and the conditions or grouping operations based on these columns are involved in the query, collect statistics about these columns so that the query optimizer can accurately estimate the number of rows and generate an effective execution plan.

Generating Statistics

- Update statistics on a single table.
 ANALYZE tablename.
- Update the statistics of the entire database. ANALYZE;
- Collect statistics from multiple columns.
 - Collect statistics on the column_1 and column_2 columns of the tablename table.

ANALYZE tablename ((column_1, column_2));

 --Add declarations for the column_1 and column_2 columns of the tablename table.

ALTER TABLE tablename ADD STATISTICS ((column_1, column_2));

 Collect the statistics of a single column and statistics of multiple declared columns.

ANALYZE tablename;

 Delete the statistics of column_1 and column_2 in the tablename table or their declarations.

ALTER TABLE tablename DELETE STATISTICS ((column_1, column_2));

NOTICE

- After the statistics are declared for multiple columns by running the ALTER
 TABLE Tablename ADD STATISTICS statement, the system collects the
 statistics about these columns next time ANALYZE is performed on the table or
 the entire database. To collect the statistics, run the ANALYZE statement.
- Use EXPLAIN to show the execution plan of each SQL statement. If rows=10
 (the default value, probably indicating the table has not been analyzed) is
 displayed in the SEQ SCAN output of a table, run the ANALYZE statement for
 this table.

Improving the Quality of Statistics

ANALYZE samples data from a table based on the random sampling algorithm and calculates table data features based on the samples. The number of samples can be specified by the **default_statistics_target** parameter. The value of **default_statistics_target** ranges from **-100** to **10000** and the default value is **100**.

- If the value of default_statistics_target is greater than 0, the number of samples is 300 x default_statistics_target. This means a larger value of default_statistics_target indicates a larger number of samples, larger memory space occupied by samples, and longer time required for calculating statistics.
- If the value of default_statistics_target is smaller than 0, the number of samples is default_statistics_target/100 x Total number of rows in the table. A smaller value of default_statistics_target indicates a larger number of samples. If the value of default_statistics_target is smaller than 0, the sampled data is written to the disk. In this case, the samples do not occupy memory. However, the calculation still takes a long time because the sample size is too large.

When **default_statistics_target** is negative, the number of samples is calculated as **default_statistics_target** divided by 100, multiplied by the total number of rows in the table. This sampling mode is also known as percentage sampling.

Automatic Statistics Collection

When the **autoanalyze** parameter is turned on, the optimizer will automatically collect statistics if it finds that there are no statistics in the table or if the data changes exceed a certain threshold. This ensures that the optimizer has the information it needs to make precise decisions.

In a cost-based optimizer (CBO) model, statistics play a crucial role in determining whether a query plan is generated. Therefore, it is crucial to have timely and effective statistics.

- Table-level statistics are stored in relpages and reltuples of pg_class.
- Column-level statistics, stored in pg_statistics and accessible through the pg_statistics view, provide information on the percentage of NULL values, percentage of distinct values, high-frequency MCV values, and histograms.

Collection condition: If there is a substantial change in data volume (default threshold is **10%**), indicating a shift in data characteristics, the system will initiate the collection of statistics again.

Overall policy: The system enables dynamic sampling to collect statistics promptly and polling sampling to ensure persistent statistics. To ensure fast query performance with response times in seconds, it is recommended to use manual sampling.

Basic Rules

Table 13-13 Typical sampling methods

Functio n	Description	Feature	Constrain t
Auto samplin g	After making significant changes to the data in a job, you need to manually run the ANALYZE command.	 In normal mode, statistics are stored in system catalogs and shared globally. A level-4 lock is applied, preventing concurrent operations on a table. In light mode, statistics are stored in memory and shared globally. A level-1 lock is applied, allowing concurrent operations on a table. In force mode, you can perform forcible sampling even when statistics are locked, in addition to the normal mode functionalities. Syntax: ANALYZE tablename; ANALYZE (light force) tablename; 	N/A
Polling samplin g	Background thread operates according to a threshold. Polling maintenance statistics	Only the normal mode is supported. Statistics are stored in system catalogs and shared. A level-4 lock is applied, preventing concurrent operations on a table. Related GUC parameters: • autovacuum • autovacuum_mode • autovacuum_analyze_threshol d • autovacuum_analyze_scale_factor	Asynchro nous polling triggering

Functio n	Description	Feature	Constrain t
Dynami c samplin g	Depending on the threshold, the query parsing process can take several dozen seconds. Real-time maintenance statistics	 In normal mode, statistics are stored in system catalogs and shared globally. A level-4 lock is applied, preventing concurrent operations on a table. In light mode, statistics are stored in memory and shared globally. A level-1 lock is applied, allowing concurrent operations on a table. Related GUC parameters: autoanalyze autoanalyze_mode 	Real-time triggering upon query In lightweig ht scenarios, persistenc e relies on polling sampling.
Forcible samplin g	Uses SQL hints to forcefully gather statistics for each query.	Used in data feature-sensitive scenarios to ensure real-time and up-to-date query statistics. Usage: select /*+ lightanalyze (t1 1) */ from t1; (1: forcible sampling; 0: sampling disabled)	The SQL statement needs to be modified.
Collecti ng partitio n statistic s	Collects incremental information by partition and combines it globally.	Used in ultra-large partitioned tables to ensure accurate query cost estimation after partition pruning.	This method takes up more storage space but provides greater accuracy.
Collecti ng statistic s from multipl e column s	Gather statistics from multiple columns.	Used to filter multiple columns simultaneously to ensure accurate query cost estimation.	You need to select target columns manually and use temporar y tables.
Collecti ng expressi on statistic s	Collects statistics on a column based on expression functions.	Used in batch expression filtering scenarios to ensure accurate query cost estimation.	Manual identificat ion is required.

Functio n	Description	Feature	Constrain t
Collecti ng expressi on index statistic s	Automatically collects statistics for created expression indexes.	Used in the point query expression filtering scenario to ensure accurate query cost estimation.	Manual identificat ion is required.
Freezin g statistic s	Freezes table-level statistics to prevent changes.	Used in scenarios where data features are extremely stable to prevent sampling and query plan changes. Used in scenarios where data features are highly variable to ensure sampling for each query. Parameter: table-level attribute analyze_mode	N/A
Modifyi ng statistic s	Directly modifies statistics after manual calculation.	Used to maintain a low sampling ratio with manual calibration. Usage: select approx_count_distinct(col_nam e) from table_name; alter table set (n_distinct=xxx)	N/A
Copyin g partitio n informa tion	Copies statistics from old partitions to new ones.	Used for partitioned tables with minimal data feature changes to reduce statistics collection overhead.	N/A
Statistic al informa tion inferenc e	Automatically calculates more accurate statistics based on existing data.	Controlled by the GUC parameter enable_extrapolation_stats.	N/A
Backing up and restorin g statistic s	Backs up statistics to an SQL statement using the EXPLAIN (STAT ON) command.	Used for scenario reproduction or statistics restoration.	Statistics are exported as SQL statement s.

Scenarios and Strategies

The table below outlines typical data processing scenarios and the corresponding strategies for collecting statistics.

Table 13-14 Statistics collection strategies

Scenario	Description	Strategy
Incremental stream processing	Incremental data flow changes with no reasonable time for ANALYZE.	Enable dynamic sampling to automatically collect and share statistics globally.
Online batch processing (Data lake)	Data processing and querying occur concurrently, requiring stable queries.	Enable dynamic sampling or complete data processing and ANALYZE within a transaction. begin; truncate table or partition; copy/merge/insert overwrite ANALYZE (light) tablename; end;
Partition parallel processing	Concurrent data processing in different partitions	Enable dynamic or manual light sampling and collect statistics concurrently for the same table.
Flat-wide table scenario	Wide table with over 100 columns	 Enable automatic predicate management for dynamic sampling. Collect statistics only on the first N columns. Set column-level participation in sampling based on common query predicates.
Large table scenario	Large data volume with changes not reaching the threshold Variable statistics	Lower the threshold for triggering dynamic sampling.
Feature- sensitive scenario	Changeable data features causing unstable query plans, requiring forcible collection.	 Lower the threshold for triggering dynamic sampling. Use the HINT mode in SQL statements for light dynamic sampling. Clear and freeze statistics, recollecting them for each query without sharing.
High- concurrency scenario	Concurrent queries (over 10) are performed on the same table, triggering dynamic sampling and resource usage.	 Disable concurrency, and other queries use outdated statistics. Generate the latest statistics before querying (under development).

Scenario	Description	Strategy
Streaming performance sensitivity	Stream processing with queries responded in seconds or high resource usage	Disable dynamic sampling at the table or SQL level and use background polling sampling.
Batch performance sensitivity	Batch processing with queries responded in seconds or high resource usage	Manually collect statistics during processing.

Resource Consumption

Table 13-15 Resource consumption

Category	Sub-Category	Description
CPU	Predicate column management	Automatically manage predicates and collect statistics only on queried columns.
		Manually mask non-predicate columns.
	Ultra-long column statistics	Data type that can be truncated, counting only the first 1,024 characters.
I/O	30,000 samples are collected by default.	Related to the number of columns, partitions, and small CUs, not table size.
Memory	Buffer usage	At most one slot in the cstore buffer is occupied.
	Memory zero copy	Directly calculate statistics from buffer samples without organizing into tuples.
	Memory adaptation	Configure the system to use temporary tables for sampling when memory is insufficient. Prevent temporary table creation triggered by queries using the analyze_stats_mode parameter.
	Memory size	Control maximum memory usage during ANALYZE with the maintenance_work_mem parameter. Exceeding memory limits results in data being written to disks or reduced samples.

Category	Sub-Category	Description
Lock	Level-4 lock	(Normal mode) Applied in distributed mode, conflicting with DDL , VACUUM , ANALYZE , and REINDEX but not with addition, deletion, or modification.
	Level-1 lock	(Light mode) Only local level-1 lock is supported, conflicting only with DDL statements.

Accuracy and Reliability

Table 13-16 Accuracy/Reliability

Accura cy/ Reliabi lity	Item	Description
Accura cy	Sampli ng size	Configurable to adapt to table size with the default_statistics_target parameter.
	Sampli ng rando mness	 Optimize reservoir and range sampling with the analyze_sample_mode parameter. Enhance randomness of random number calculation with the random_function_version parameter.
	Global sharing	Statistics can be shared across sessions and nodes.
	Modifyi ng	Background thread checks and broadcasts the global modification count in polling mode.
	count broadc ast	The job thread can also directly broadcast the modification count by specifying the tuple_change_sync_threshold parameter.
		Cross-CN modification and query have minimal impact. The modification count is broadcast and synchronized in asynchronous mode.
	Adjusti ng the CU sampli ng ratio	Increase CU sampling ratio if the CU filling rate is low, using the cstore_cu_sample_ratio parameter.
	Stabiliz ing distinct values	Use the n_distinct parameter to stabilize distinct values after random sampling without increasing the sampling ratio.

Accura cy/ Reliabi lity	Item	Description
	Statisti cal inform ation calcula tion	Use the enable_extrapolation_stats parameter to calculate more accurate statistics based on old statistics during distortion estimation.
Reliabil ity	CN fault	Dynamic sampling is unaffected by other CN faults, and statistics are not synchronized. Query quality on the current CN remains unaffected.
	CN restora tion	Forcibly perform dynamic sampling and global synchronization during queries after CN recovery.
	DN fault	Dynamic sampling of the logical cluster is unaffected by faults in other logical clusters.

O&M Monitoring

GaussDB(DWS) offers a comprehensive view of the **ANALYZE** running mode and different execution stages by adding comments after the **ANALYZE** command. This information is primarily presented through the following views:

- query column in the pgxc_stat_activity view
- wait_status column in the pgxc_thread_wait_status view

The format of the **ANALYZE** command is **--Action-RunMode-StatsMode-SyncMode**.

• Values and meanings of **Action**:

{"begin", "finished", "lock FirstCN", "estimate rows", "statistics", "sample rows", "calc stats"};

begin: indicates the start of the process; **finished**: indicates the end of the process; **lock FirstCN**: applies a lock from the FirstCN; **estimate rows**: estimates the number of rows in the first phase; **statistics**: executes **ANALYZE** in the second phase; **sample rows**: collects samples in the second phase; **calc stats**: calculates statistics in the second phase.

• Values and meanings of **RunMode**:

{"manual", "backend", "normal runtime", "light runtime", "light runtime inxact", "light estimate rows", "light manual"};

manual: indicates the manual mode; backend: indicates the background polling mode; normal runtime: indicates the normal dynamic sampling; light runtime: indicates the light dynamic samplin; light runtime inxact: indicates the light dynamic sampling in a transaction; light estimate rows indicates the light estimation function only; light manual: indicates the manual light mode.

 Values and meanings of StatsMode: {"dynamic", "memory", "smptbl"}; **dynamic**: indicates adaptive selection of memory or temporary table placement samples; **memory**: uses only internal storage samples; **smptbl**: uses only temporary table placement samples.

• Values and meanings of **SyncMode**: {"sync", "nosync"};

sync: Statistics are synchronized to all CNs; **nosync**: Statistics are not synchronized.

Example:

Viewing Statistics

- Check the dynamically sampled memory statistics.
 - Retrieve table-level memory statistics.
 SELECT * FROM pv_runtime_relstats;
 - Retrieve column-level memory statistics.
 SELECT * FROM pv_runtime_attstats;
- Check the system catalog statistics.
 - Check the table-level system catalog statistics.
 select relname, relpages, reltuples from pg_class;
 - Check the column-level system catalog statistics.
 SELECT * FROM pg_stats;
- Check the latest time when statistics are collected.

Dynamic sampling stores statistics in memory without modifying the timestamp of the system catalog.

SELECT * FROM pg_object;

13.4.6 Reviewing and Modifying a Table Definition

In a distributed framework, data is distributed on DNs. Data on one or more DNs is stored on a physical storage device. To properly define a table, you must:

- Evenly distribute data on each DN to avoid the available capacity decrease
 of a cluster caused by insufficient storage space of the storage device
 associated with a DN. Specifically, select a proper distribution key to avoid
 data skew.
- 2. **Evenly assign table scanning tasks on each DN** to avoid that a DN is overloaded by the table scanning tasks. Specifically, do not select columns in the equivalent filter of a base table as the distribution key.
- 3. **Reduce the data volume scanned** by using the partition pruning mechanism.
- 4. **Avoid the use of random I/O** by using clustering or partial clustering.
- 5. **Avoid data shuffle** to reduce the network pressure by selecting the **join-condition** column or **group by** column as the distribution column.

The distribution column is the core for defining a table. **Figure 13-5** shows the procedure of defining a table. The table definition is created during the database design and is reviewed and modified during SQL tuning.

Start Service characteristic Data characteristics elect a storage mode Select partial-clustering Select distribution columns columns and a partition key No Sort distribution columns by Small table? Use a hash table Yes Create a table based on the distribution columns with high Use a replication table. priority. Define the table Import sample data No End

Figure 13-5 Defining a table

For details about how to review and modify table definitions, see **Table Optimization Practices**.

13.4.7 Advanced SQL Tuning

13.4.7.1 SQL Self-Diagnosis

Performance issues may occur when you run the INSERT/UPDATE/DELETE/SELECT/MERGE INTO or CREATE TABLE AS statement. The product supports automatic performance diagnosis and saves related diagnosis information to Real-time Top SQL. When enable_resource_track is set to on, the diagnosis information is dumped to Historical Top SQL. You can query the warning column in the GS_WLM_SESSION_STATISTICS, GS_WLM_SESSION_HISTORY, and GS_WLM_SESSION_INFO views to obtain reference information for performance tuning.

- Alarms that can trigger SQL self-diagnosis depend on the settings of resource_track_level.
 - When **resource_track_level** is set to **query**, you can diagnose alarms such as uncollected multi-column/single-column statistics, unpruned partitions, and failure of pushing down SQL statements. When **resource_track_level** is set to **perf** or **operator**, all alarms can be diagnosed.
- Whether a SQL plan will be diagnosed depends on the settings of resource track cost.
 - A SQL plan will be diagnosed only if its execution cost is greater than **resource_track_cost**. You can use the **EXPLAIN** keyword to check the plan execution cost
- When EXPLAIN PERFORMANCE or EXPLAIN VERBOSE is executed, SQL selfdiagnosis information, except the ones without multi-column statistics, will be generated. For details, see SQL Execution Plan.

Alarms Related to SQL Execution Performance

Currently, the following alarms on performance issues will be reported:

1. Statistics of a single column or multiple columns are not collected.

If statistics of a single column or multiple columns are not collected, an alarm is reported. To handle this alarm, you are advised to perform **ANALYZE** on related tables. For details, see **Updating Statistics** and **Optimizing Statistics**.

If no statistics are collected for the OBS foreign table and HDFS foreign table in the query statement, an alarm indicating that statistics are not collected will be reported. Because the **ANALYZE** performance of the OBS foreign table and HDFS foreign table is poor, you are not advised to perform **ANALYZE** on these tables. Instead, you are advised to use the **ALTER FOREIGN TABLE** syntax to modify the **totalrows** attribute of the foreign table to correct the estimated number of rows.

Example alarms:

The statistics about a table are not collected.

```
Statistic Not Collect schema_test.t1
```

The statistics about a single column are not collected.

```
Statistic Not Collect schema_test.t2(c1)
```

The statistics about multiple columns are not collected.

```
Statistic Not Collect schema_test.t3((c1,c2))
```

The statistics about a single column and multiple columns are not collected.

```
Statistic Not Collect
schema_test.t4(c1)
schema_test.t5((c1,c2))
```

2. Partitions are not pruned.

When a partitioned table is queried, the partition is pruned based on the constraints on the partition key to improve the query performance. However, the partition table may not be pruned due to improper constraints, deteriorating the query performance. For details, see Case: Rewriting SQL Statements and Eliminating Prune Interference.

SQL statements are not pushed down.

The cause details are displayed in the alarms. For details, see **Optimizing Statement Pushdown**.

The potential causes for the pushdown failure are as follows:

Caused by functions

The function name is displayed in the diagnosis information. Function pushdown is determined by the **shippable** attribute of the function. For details, see the **CREATE FUNCTION** syntax.

Caused by syntax

The diagnosis information displays the syntax that causes the pushdown failure. For example, if the statement contains the **With Recursive**, **Distinct On**, or **row** expression and the return value is of the record type, an alarm is reported, indicating that the syntax does not support pushdown.

Example alarms:

SQL is not plan-shipping
"enable_stream_operator" is off

SQL is not plan-shipping
"Distinct On" can not be shipped

SQL is not plan-shipping
"v_test_unshipping_log" is VIEW that will be treated as Record type can't be shipped

4. Vectorized plans are not supported.

For SQL statements that cannot use vectorized plans, detailed reasons why vectorized plans cannot be used are reported.

Common reasons are as follows:

- The target column contains functions whose return type is a set.
- The target column or query condition, the distribution key of the Stream operator, and the Limit and Offset clauses contain expressions that cannot be vectorized (such as geospatial types, array expressions, Row expressions, XML expressions, and functions whose parameters or return values contain the refcursor type).
- The **Group By** clause contains an array-equivalent judgment statement.
- GC FDW and LOG FDW do not support vectorization.
- The plan contains operators such as Cte Scan, Recursive Union, Merge Append, and Lock Rows.

Example alarms:

```
SQL is un-vectorized
Function regexp_split_to_table that returns set is un-vectorized

SQL is un-vectorized
Array expression is un-vectorized

SQL is un-vectorized
Function array_agg is un-vectorized

SQL is un-vectorized
RecursiveUnion is un-vectorized
```

5. In a hash join, the larger table is used as the inner table.

An alarm will be reported if the number of rows in the inner table reaches or exceeds 10 times of that in the foreign table, more than 100,000 inner-table rows are processed on each DN in average, and data has been flushed to disks. You can check the **query_plan** column in **GS_WLM_SESSION_HISTORY** to check whether hash joins are used. In this scenario, you need to adjust the sequence of the HashJoin internal and foreign tables. For details, see **Join Order Hints**.

Example alarm:

Execute diagnostic information PlanNode[7] Large Table is INNER in HashJoin "Vector Hash Aggregate"

In the preceding command, 7 indicates the operator whose ID is 7 in the query_plan column.

6. **nestloop** is used in a large-table equivalent join.

An alarm will be reported if nested loop is used in an equivalent join where more than 100,000 larger-table rows are processed on each DN in average. You can check the **query_plan** column of **GS_WLM_SESSION_HISTORY** to see if nested loop is used. In this scenario, you need to adjust the table join mode

and disable the NestLoop join mode between the current internal and foreign tables. For details, see **Join Operation Hints**.

Example alarm:

Execute diagnostic information

PlanNode[5] Large Table with Equal-Condition use Nestloop"Nested Loop"

7. A large table is broadcasted.

An alarm will be reported if more than 100 thousand of rows are broadcasted on each DN in average. In this scenario, the broadcast operation of the Broadcast lower-layer operator needs to be disabled. For details, see **Stream Operation Hints**.

Example alarm:

Execute diagnostic information

PlanNode[5] Large Table in Broadcast "Streaming(type: BROADCAST dop: 1/2)"

Data skew occurs.

An alarm will be reported if the number of rows processed on any DN exceeds 100 thousand, and the number of rows processed on a DN reaches or exceeds 10 times of that processed on another DN. Generally, this alarm is generated due to storage layer skew or computing layer skew. For details, see Optimizing Data Skew.

Example alarm:

Execute diagnostic information

PlanNode[6] DataSkew:"Seq Scan", min_dn_tuples:0, max_dn_tuples:524288

9. The index is improper.

During base table scanning, an alarm is reported if the following conditions are met:

- For row-store tables:
 - When the index scanning is used, the ratio of the number of output lines to the number of scanned lines is greater than 1/1000 and the number of output lines is greater than 10,000.
 - When sequential scanning is used, the number of output lines to the number of scanned lines is less than 1/1000, the number of output lines is less than or equal to 10,000, and the number of scanned lines is greater than 10,000.
- For column-store tables:
 - When the index scanning is used, the ratio of the number of output lines to the number of scanned lines is greater than 1/10000 and the number of output lines is greater than 100.
 - When sequential scanning is used, the number of output lines to the number of scanned lines is less than 1/10,000, the number of output lines is less than or equal to 100, and the number of scanned lines is greater than 10,000.

For details, see **Tuning Operators**. You can also refer to **Case: Creating an Appropriate Index** and **Case: Setting Partial Cluster Keys**.

Example alarms:

Execute diagnostic information

PlanNode[4] Indexscan is not properly used:"Index Only Scan", output:524288, filtered:0, rate:1.00000

PlanNode[5] Indexscan is ought to be used:"Seq Scan", output:1, filtered:524288, rate:0.00000

The diagnosis result is only a suggestion for the current SQL statement. You are advised to create an index only for frequently used filter criteria.

10. Estimation is inaccurate.

An alarm will be reported if the maximum number or the estimated maximum number of rows processed on a DN is over 100,000, and the larger number reaches or exceeds 10 times of the smaller one. In this scenario, you can refer to **Rows Hints** to correct the estimation on the number of rows, so that the optimizer can re-design the execution plan based on the correct number.

Example alarm:

Execute diagnostic information
PlanNode[5] Inaccurate Estimation-Rows: "Hash Join" A-Rows:0, E-Rows:52488

Constraints

- 1. An alarm contains a maximum of 2048 characters. If the length of an alarm exceeds this value (for example, a large number of long table names and column names are displayed in the alarm when their statistics are not collected), a warning instead of an alarm will be reported.

 WARNING, "Planner issue report is truncated, the rest of planner issues will be skipped"
- 2. If a query statement contains the **Limit** operator, alarms of operators lower than **Limit** will not be reported.
- 3. For alarms about data skew and inaccurate estimation, only alarms on the lower-layer nodes in a plan tree will be reported. This is because the same alarms on the upper-level nodes may be triggered by problems on the lower-layer nodes. For example, if data skew occurs on the **Scan** node, data skew may also occur in operators (for example, **Hashagg**) at the upper layer.

13.4.7.2 Optimizing Statement Pushdown

Statement Pushdown

Currently, the GaussDB(DWS) optimizer can use three methods to develop statement execution policies in the distributed framework: generating a statement pushdown plan, a distributed execution plan, or a distributed execution plan for sending statements.

- A statement pushdown plan pushes query statements from a CN down to DNs for execution and returns the execution results to the CN.
- In a distributed execution plan, a CN compiles and optimizes query statements, generates a plan tree, and then sends the plan tree to DNs for execution. After the statements have been executed, execution results will be returned to the CN.
- A distributed execution plan for sending statements pushes queries that can be pushed down (mostly base table scanning statements) to DNs for execution. Then, the plan obtains the intermediate results and sends them to the CN, on which the remaining queries are to be executed.

When sending statements through a distributed execution plan, DNs send numerous intermediate results to CNs. However, certain statements cannot be pushed down and must be executed on CNs, leading to performance bottlenecks in bandwidth, storage, and computing. Therefore, you are not advised to use the query statements that only the third policy is applicable to.

Statements cannot be pushed down to DNs if they have Functions That Do Not Support Pushdown or Syntax That Does Not Support Pushdown. Generally, you can rewrite the execution statements to solve the problem.

Viewing Whether the Execution Plan Has Been Pushed Down to DNs

Perform the following procedure to quickly determine whether the execution plan can be pushed down to DNs:

Step 1 Set the GUC parameter **enable_fast_query_shipping** to **off** to use the distributed framework policy for the query optimizer.

```
SET enable_fast_query_shipping = off,
```

Step 2 View the execution plan.

If the execution plan contains Data Node Scan, the SQL statements cannot be pushed down to DNs. If the execution plan contains Streaming, the SQL statements can be pushed down to DNs.

For example:

```
select
count(ss.ss_sold_date_sk order by ss.ss_sold_date_sk)c1
from store_sales ss, store_returns sr
where
sr.sr_customer_sk = ss.ss_customer_sk;
```

The execution plan is as follows, which indicates that the SQL statement cannot be pushed down.

----End

Syntax That Does Not Support Pushdown

SQL syntax that does not support pushdown is described using the following table definition examples:

DISTRIBUTE BY hash(C_CUSTKEY); CREATE TABLE test_stream(a int, b float);--float does not support redistribution. postgresCREATE TABLE sal_emp (c1 integer[]) DISTRIBUTE BY replication;

The **returning** statement cannot be pushed down.

```
postgresexplain update customer1 set C_NAME = 'a' returning c_name;
                   QUERY PLAN
Update on customer1 (cost=0.00..0.00 rows=30 width=187)
 Node/s: All datanodes
 Node expr: c_custkey
 -> Data Node Scan on customer1 " REMOTE TABLE QUERY " (cost=0.00..0.00 rows=30 width=187)
     Node/s: All datanodes
```

If columns in **count(distinct expr)** do not support redistribution, they do not support pushdown.

```
postgresexplain verbose select count(distinct b) from test_stream;
                          QUERY PLAN
                                          ------ Aggregate (cost=2.50..2.51 rows=1 width=8)
 Output: count(DISTINCT test_stream.b)
 -> Data Node Scan on test_stream "_REMOTE_TABLE_QUERY_" (cost=0.00..0.00 rows=30 width=8)
     Output: test_stream.b
     Node/s: All datanodes
     Remote query: SELECT b FROM ONLY public.test_stream WHERE true
(6 rows)
```

Statements using **distinct on** cannot be pushed down.

```
postgresexplain verbose select distinct on (c_custkey) c_custkey from customer1 order by c_custkey;
                           QUERY PLAN
                                          ------ Unique (cost=49.83..54.83 rows=30 width=8)
 Output: customer1.c_custkey
 -> Sort (cost=49.83..52.33 rows=30 width=8)
     Output: customer1.c_custkey
     Sort Key: customer1.c_custkey
     -> Data Node Scan on customer1 " REMOTE TABLE QUERY " (cost=0.00..0.00 rows=30
width=8)
         Output: customer1.c_custkey
         Node/s: All datanodes
         Remote query: SELECT c_custkey FROM ONLY public.customer1 WHERE true
(9 rows)
```

In a statement using **FULL JOIN**, if the column specified using **JOIN** does not support redistribution, the statement does not support pushdown.

```
postgresexplain select * from test_stream t1 full join test_stream t2 on t1.a=t2.b;
                             QUERY PLAN
                                          ------ Hash Full Join (cost=0.38..0.82 rows=30
width=24)
 Hash Cond: ((t1.a)::double precision = t2.b)
 -> Data Node Scan on test_stream "_REMOTE_TABLE_QUERY_" (cost=0.00..0.00 rows=30 width=12)
     Node/s: All datanodes
 -> Hash (cost=0.00..0.00 rows=30 width=12)
      -> Data Node Scan on test_stream "_REMOTE_TABLE_QUERY_" (cost=0.00..0.00 rows=30
width=12)
         Node/s: All datanodes
(7 rows)
```

A statement containing array expressions cannot be pushed down.

```
postgresexplain verbose select array[c_custkey,1] from customer1 order by c_custkey;
```

```
QUERY PLAN
                                         ------ Sort (cost=49.83..52.33 rows=30 width=8)
 Output: (ARRAY[customer1.c_custkey, 1::bigint]), customer1.c_custkey
 Sort Key: customer1.c_custkey
 -> Data Node Scan on "__REMOTE_SORT_QUERY__" (cost=0.00..0.00 rows=30 width=8)
     Output: (ARRAY[customer1.c_custkey, 1::bigint]), customer1.c_custkey
     Node/s: All datanodes
     Remote query: SELECT ARRAY[c_custkey, 1::bigint], c_custkey FROM ONLY public.customer1
WHERE true ORDER BY 2
(7 rows)
```

• Subplans that are shared among multiple threads and cannot be pushed down.

```
postgres=# explain verbose select c_custkey in (select c_custkey from customer1) b from customer1;

QUERY PLAN

Data Node Scan on customer1 "_REMOTE_TABLE_QUERY_" (cost=2.50..5.00 rows=1000 width=8)
Output: (hashed SubPlan 1)
Node/s: All datanodes
Remote query: SELECT c_custkey FROM ONLY public.customer1 WHERE true
SubPlan 1
-> Data Node Scan on customer "_REMOTE_TABLE_QUERY_" (cost=0.00..0.00 rows=1000 width=8)
Output: public.customer.c_custkey
Node/s: All datanodes
Remote query: SELECT c_custkey FROM ONLY public.customer1 WHERE true
(9 rows)
```

The following table describes the scenarios where a statement containing
 WITH RECURSIVE cannot be pushed down in the current version, as well as the causes.

No.	Scenario	Cause of Not Supporting Pushdown
1	The query contains foreign tables or HDFS tables.	LOG: SQL can't be shipped, reason: RecursiveUnion contains HDFS Table or ForeignScan is not shippable (In this table, LOG describes the cause of not supporting pushdown.)
		In the current version, queries containing foreign tables or HDFS tables do not support pushdown.
2	Multiple Node Groups	LOG: SQL can't be shipped, reason: With-Recursive under multi-nodegroup scenario is not shippable
		In the current version, pushdown is supported only when all base tables are stored and computed in the same Node Group.
3	WITH recursive t_result AS (SELECT dm,sj_dm,name,1 as level FROM test_rec_part WHERE sj_dm > 10 UNION SELECT t2.dm,t2.sj_dm,t2.name ' > '	LOG: SQL can't be shipped, reason: With-Recursive does not contain "ALL" to bind recursive & none-recursive branches
	t1.name,t1.level+1 FROM t_result t1 JOIN test_rec_part t2 ON t2.sj_dm = t1.dm) SELECT * FROM t_result t;	ALL is not used for UNION . In this case, the return result is deduplicated.

No.	Scenario	Cause of Not Supporting Pushdown
4	WITH RECURSIVE x(id) AS (select count(1) from pg_class where oid=1247 UNION ALL SELECT id+1 FROM x WHERE id < 5), y(id) AS (select count(1) from pg_class where oid=1247 UNION ALL SELECT id+1 FROM x WHERE id < 10) SELECT y.*, x.* FROM y LEFT JOIN x USING (id) ORDER BY 1;	LOG: SQL can't be shipped, reason: With-Recursive contains system table is not shippable A base table contains the system catalog.
5	WITH RECURSIVE t(n) AS (VALUES (1) UNION ALL SELECT n+1 FROM t WHERE n < 100) SELECT sum(n) FROM t;	LOG: SQL can't be shipped, reason: With-Recursive contains only values rte is not shippable Only VALUES is used for scanning base tables. In this case, the statement can be executed on the CN, and DNs are unnecessary.
6	select a.ID,a.Name, (with recursive cte as (select ID, PID, NAME from b where b.ID = 1 union all select parent.ID,parent.PID,parent.NAME from cte as child join b as parent on child.pid=parent.id where child.ID = a.ID) select NAME from cte limit 1) cName from (select id, name, count(*) as cnt from a group by id,name) a order by 1,2;	LOG: SQL can't be shipped, reason: With-Recursive recursive term correlated only is not shippable The correlation conditions of correlated subqueries are only in the recursion part, and the non-recursion part has no correlation condition.
7	WITH recursive t_result AS (select * from(SELECT dm,sj_dm,name,1 as level FROM test_rec_part WHERE sj_dm < 10 order by dm limit 6 offset 2) UNION all SELECT t2.dm,t2.sj_dm,t2.name ' > ' t1.name,t1.level+1 FROM t_result t1 JOIN test_rec_part t2 ON t2.sj_dm = t1.dm) SELECT * FROM t_result t;	LOG: SQL can't be shipped, reason: With-Recursive contains conflict distribution in none-recursive(Replicate) recursive(Hash) The replicate plan is used for limit in the non-recursion part but the hash plan is used in the recursion part, resulting in conflicts.

No.	Scenario	Cause of Not Supporting Pushdown
8	with recursive cte as (select * from rec_tb4 where id<4 union all select h.id,h.parentID,h.name from (with recursive cte as (select * from rec_tb4 where id<4 union all select h.id,h.parentID,h.name from rec_tb4 h inner join cte c on h.id=c.parentID) SELECT id ,parentID,name from cte order by parentID) h inner join cte c on h.id=c.parentID) SELECT id ,parentID,name from cte order by parentID,1,2,3;	LOG: SQL can't be shipped, reason: Recursive CTE references recursive CTE "cte" recursive of multiple-layers are nested. That is, a recursive is nested in the recursion part of another recursive.

Functions That Do Not Support Pushdown

This module describes the variability of functions. The function variability in GaussDB(DWS) is as follows:

IMMUTABLE

Indicates that the function always returns the same result if the parameter values are the same.

STABLE

Indicates that the function cannot modify the database, and that within a single table scan it will consistently return the same result for the same parameter values, but that its result varies by SQL statements.

VOLATILE

Indicates that the function value can change even within a single table scan, so no optimizations can be made.

The volatility of a function can be obtained by querying its **provolatile** column in **pg_proc**. The value **i** indicates immutable, **s** indicates stable, and **v** indicates volatile. The valid values of the **proshippable** column in **pg_proc** are **t**, **f**, and **NULL**. This column and the **provolatile** column together describe whether a function is pushed down.

- If the **provolatile** of a function is **i**, the function can be pushed down regardless of the value of **proshippable**.
- If the **provolatile** of a function is **s** or **v**, the function can be pushed only if the value of **proshippable** is **t**.
- CTEs containing random are not pushed down, because pushdown may lead to incorrect results.

For a UDF, you can specify the values of **provolatile** and **proshippable** during its creation. For details, see CREATE FUNCTION.

In scenarios where a function does not support pushdown, perform one of the following as required:

- If it is a system function, replace it with a functionally equivalent one.
- If it is a UDF function, check whether its **provolatile** and **proshippable** are correctly defined.

Example: UDF

Define a user-defined function that generates fixed output for a certain input as the **immutable** type.

Use the TPCDS sales information as an example. You need to define a function to obtain the discount information.

```
CREATE FUNCTION func_percent_2 (NUMERIC, NUMERIC) RETURNS NUMERIC
AS 'SELECT $1 / $2 WHERE $2 > 0.01'
LANGUAGE SQL
VOLATILE;
```

Run the following statement:

```
SELECT func_percent_2(ss_sales_price, ss_list_price)
FROM store sales;
```

The execution plan is as follows:

```
Data Node Scan on store_sales "_REMOTE_TABLE_QUERY_"
  Output: func_percent_2(store_sales.ss_sales_price, store_sales.ss_list_price)
  Remote query: SELECT ss_sales_price, ss_list_price FROM ONLY store_sales WHERE true
(3 rows)
```

func_percent_2 is not pushed down, and **ss_sales_price** and **ss_list_price** are executed on a CN. In this case, a large amount of resources on the CN is consumed, and the performance deteriorates as a result.

In this example, the function returns certain output when certain input is entered. Therefore, we can modify the function to the following one:

```
CREATE FUNCTION func_percent_1 (NUMERIC, NUMERIC) RETURNS NUMERIC
AS 'SELECT $1 / $2 WHERE $2 > 0.01'
LANGUAGE SQL
IMMUTABLE;
```

Run the following statement:

```
SELECT func_percent_1(ss_sales_price, ss_list_price)
FROM store sales;
```

The execution plan is as follows:

```
Data Node Scan on "_REMOTE_FQS_QUERY_" (cost=0.00..0.00 rows=0 width=0)
Output: (func_percent_1(store_sales.ss_sales_price, store_sales.ss_list_price))
Node/s: All datamodes
Remote query: SELECT public.func_percent_1(ss_sales_price, ss_list_price) AS func_percent_1 FROM public.store_sales
(4 rows)
```

func_percent_1 is pushed down to DNs for quicker execution. (In TPCDS 1000X, where three CNs and 18 DNs are used, the query efficiency is improved by over 100 times).

Example 2: Pushing Down the Sorting Operation

Learn more information in Case: Pushing Down Sort Operations to DNs.

13.4.7.3 Optimizing Subqueries

What Is a Subquery

When an application runs a SQL statement to operate the database, a large number of subqueries are used because they are more clear than table join. Especially in complicated query statements, subqueries have more complete and independent semantics, which makes SQL statements clearer and easy to understand. Therefore, subqueries are widely used.

In GaussDB(DWS), subqueries can also be called sublinks based on the location of subqueries in SQL statements.

- Subquery: corresponds to a scope table (RangeTblEntry) in the query parse tree. That is, a subquery is a SELECT statement following immediately after the FROM keyword.
- Sublink: corresponds to an expression in the query parsing tree. That is, a sublink is a statement in the WHERE or ON clause or in the target list.
 In conclusion, a subquery is a scope table and a sublink is an expression in the query parsing tree. A sublink can be found in constraint conditions and expressions. In GaussDB(DWS), sublinks can be classified into the following types:
 - exist sublink: corresponding to the **EXIST** and **NOT EXIST** statements.
 - any_sublink: corresponding to the OP ANY(SELECT...) statement. OP can be the IN, <, >, or = operator.
 - all_sublink: corresponding to the OP ALL(SELECT...) statement. OP can be the IN, <, >, or = operator.
 - rowcompare_sublink: corresponding to the RECORD OP (SELECT...) statement.
 - expr_sublink: corresponding to the (SELECT with a single target list item) statement.
 - array_sublink: corresponding to the **ARRAY(SELECT...)** statement.
 - cte_sublink: corresponding to the WITH(...) statement.

NOTICE

When a subquery in a query statement is in a different execution cluster than the parent query and includes filter criteria related to the parent query, the filter criteria cannot be applied across clusters in logical cluster mode. This leads to additional operator overhead, potentially resulting in poorer performance compared to a single-cluster setup.

The sublinks commonly used in OLAP and HTAP are exist_sublink and any_sublink. The sublinks are pulled up by the optimization engine of GaussDB(DWS). Because of the flexible use of subqueries in SQL statements, complex subqueries may affect query performance. Subqueries are classified into non-correlated subqueries and correlated subqueries.

Non-correlated subquery

The execution of a subquery is independent from any attribute of outer queries. In this way, a subquery can be executed before outer queries.

Example:

```
select t1.c1,t1.c2
from t1
where t1.c1 in (
  select c2
  from t2
  where t2.c2 IN (2,3,4)
);
                 QUERY PLAN
Streaming (type: GATHER)
  Node/s: All datanodes
  -> Hash Right Semi Join
      Hash Cond: (t2.c2 = t1.c1)
      -> Streaming(type: REDISTRIBUTE)
          Spawn on: All datanodes
          -> Seq Scan on t2
              Filter: (c2 = ANY ('{2,3,4}'::integer[]))
      -> Hash
          -> Seq Scan on t1
(10 rows)
```

Correlated subquery

The execution of a subquery depends on some attributes of outer queries which are used as **AND** conditions of the subquery. In the following example, **t1.c1** in the **t2.c1** = **t1.c1** condition is a dependent attribute. Such a subquery depends on outer queries and needs to be executed once for each outer query.

Example:

```
select t1.c1,t1.c2
from t1
where t1.c1 in (
  select c2
  from t2
  where t2.c1 = t1.c1 AND t2.c2 in (2,3,4)
);
                      QUERY PLAN
Streaming (type: GATHER)
  Node/s: All datanodes
  -> Seq Scan on t1
      Filter: (SubPlan 1)
      SubPlan 1
       -> Result
           Filter: (t2.c1 = t1.c1)
           -> Materialize
                 -> Streaming(type: BROADCAST)
                   Spawn on: All datanodes
          -> Seq Scan on t2
                       Filter: (c2 = ANY ('{2,3,4}'::integer[]))
(12 rows)
```

GaussDB(DWS) SubLink Optimization

A subquery is pulled up to join with tables in outer queries, preventing the subquery from being converted into the combination of a subplan and broadcast. You can run the **EXPLAIN** statement to check whether a subquery is converted into the combination of a subplan and broadcast.

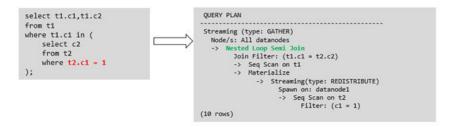
Example:

```
select t1.c1,t1.c2
from t1
where t1.c1 in (
    select c2
    from t2
    where t2.c1 = t1.c1
);

Streaming (type: GATHER)
    Node/s: All datanodes
-> Seq Scan on t1
    Filter: (SubPlan 1)
    SubPlan 1
-> Result
    Filter: (t2.c1 = t1.c1)
-> Materialize
-> Streaming(type: BROADCAST)
    Spawn on: All datanodes
-> Seq Scan on t2

(11 rows)
```

- Sublink-release supported by GaussDB(DWS)
 - Pulling up the IN sublink
 - The subquery cannot contain columns in the outer query (columns in more outer queries are allowed).
 - The subquery cannot contain volatile functions.



- Pulling up the **EXISTS** sublink

The **WHERE** clause must contain a column in the outer query. Other parts of the subquery cannot contain the column. Other restrictions are as follows:

- The subguery must contain the **FROM** clause.
- The subquery cannot contain the **WITH** clause.
- The subquery cannot contain aggregate functions.
- The subquery cannot contain a SET, SORT, LIMIT, WindowAgg, or HAVING operation.
- The subquery cannot contain volatile functions.

```
select t1.c1,t1.c2
from t1
where exists (
    select c2
from t2
    where t2.c1 = t1.c1
);

QUERY PLAN

Streaming (type: GATHER)
Node/s: All datanodes
-> Hash Semi Join
Hash Cond: (t1.c1 = t2.c1)
-> Seq Scan on t1
-> Hash
-> Seq Scan on t2
(7 rows)
```

Pulling up an equivalent query containing aggregation functions
 The WHERE condition of the subquery must contain a column from the outer query. Equivalence comparison must be performed between this column and related columns in tables of the subquery. These conditions must be connected using AND. Other parts of the subquery cannot contain the column. Other restrictions are as follows:

- The expression in the WHERE condition of the subquery must be table columns.
- After the SELECT keyword of the subquery, there must be only one output column. The output column must be an aggregate function (for example, MAX), and the parameter (for example, t2.c2) of the aggregate function cannot be columns of a table (for example, t1) in outer queries. The aggregate function cannot be COUNT.

```
For example, the following subquery can be pulled up:

select * from t1 where c1 > (

    select max(t2.c1) from t2 where t2.c1=t1.c1

);
```

The following subquery cannot be pulled up because the subquery has no aggregation function.

```
select * from t1 where c1 >(
select t2.c1 from t2 where t2.c1=t1.c1
);
```

The following subquery cannot be pulled up because the subquery has two output columns:

```
select * from t1 where (c1,c2) >(
     select max(t2.c1),min(t2.c2) from t2 where t2.c1=t1.c1
);
```

- The subquery must be a FROM clause.
- The subquery cannot contain a GROUP BY, HAVING, or SET operation.
- The subquery can only be inner join.

```
For example, the following subquery can be pulled up:
select * from t1 where c1 >(
select max(t2.c1) from t2 full join t3 on (t2.c2=t3.c2) where t2.c1=t1.c1
);
```

- The target list of the subquery cannot contain the function that returns a set.
- The WHERE condition of the subquery must contain a column from the outer query. Equivalence comparison must be performed between this column and related columns in tables of the subquery. These conditions must be connected using AND. Other parts of the subquery cannot contain the column. For example, the following subquery can be pulled up:

```
select * from t3 where t3.c1=(
    select t1.c1
    from t1 where c1 >(
        select max(t2.c1) from t2 where t2.c1=t1.c1
));
```

If another condition is added to the subquery in the previous example, the subquery cannot be pulled up because the subquery references to the column in the outer query. Example:

```
select * from t3 where t3.c1=(
    select t1.c1
    from t1 where c1 >(
        select max(t2.c1) from t2 where t2.c1=t1.c1 and t3.c1>t2.c2
));
```

Pulling up a sublink in the **OR** clause

If the **WHERE** condition contains a **EXIST**-related sublink connected by **OR**,

for example,

```
select a, c from t1
where t1.a = (select avg(a) from t3 where t1.b = t3.b) or
exists (select * from t4 where t1.c = t4.c);
```

The procedure for promoting the OR clause of an EXIST-related subquery in an OR-ed join is as follows:

- i. Extract opExpr from the OR clause in the WHERE condition. The value is t1.a = (select avg(a) from t3 where t1.b = t3.b).
- ii. The opExpr contains a subquery. If the subquery can be pulled up, the subquery is rewritten as elect avg(a), t3.b from t3 group by t3.b, generating the NOT NULL condition t3.b is not null. The opExpr is replaced with this NOT NULL condition. In this case, the SQL statement changes to:

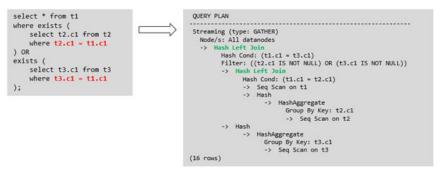
```
select a, c from t1 left join (select avg(a) avg, t3.b from t3 group by t3.b) as t3 on (t1.a = avg and t1.b = t3.b) where t3.b is not null or exists (select * from t4 where t1.c = t4.c);
```

iii. Extract the EXISTS sublink exists (select * from t4 where t1.c = t4.c) from the OR clause to check whether the sublink can be pulled up. If it can be pulled up, it is converted into select t4.c from t4 group by t4.c, generating the NOT NULL condition t4.c is not null. In this case, the SQL statement changes to:

in this case, the SQL statement changes to:

from t1 left join (select avg(a) avg, t3.b from t3 group by t3.b) as t3 on (t1.a = avg and t1.b = t3.b)

left join (select t4.c from t4 group by t4.c) where t3.b is not null or t4.c is not null;



Sublink-release not supported by GaussDB(DWS)

Except the sublinks described above, all the other sublinks cannot be pulled up. In this case, a join subquery is planned as the combination of a subplan and broadcast. As a result, if tables in the subquery have a large amount of data, query performance may be poor.

If a correlated subquery joins with two tables in outer queries, the subquery cannot be pulled up. You need to change the parent query into a **WITH** clause and then perform the join.

Example:

```
select distinct t1.a, t2.a from t1 left join t2 on t1.a=t2.a and not exists (select a,b from test1 where test1.a=t1.a and test1.b=t2.a);
```

The parent query is changed into:

```
with temp as
(
select * from (select t1.a as a, t2.a as b from t1 left join t2 on t1.a=t2.a)
)
select distinct a,b
from temp
where not exists (select a,b from test1 where temp.a=test1.a and temp.b=test1.b);
```

- The subquery (without **COUNT**) in the target list cannot be pulled up.

Example:

```
explain (costs off)
select (select c2 from t2 where t1.c1 = t2.c1) ssq, t1.c2
from t1
where t1.c2 > 10;
```

The execution plan is as follows:

```
explain (costs off)
select (select c2 from t2 where t1.c1 = t2.c1) ssq, t1.c2
from t1
where t1.c2 > 10:
              QUERY PLAN
Streaming (type: GATHER)
  Node/s: All datanodes
  -> Seq Scan on t1
     Filter: (c2 > 10)
     SubPlan 1
       -> Result
           Filter: (t1.c1 = t2.c1)
           -> Materialize
               -> Streaming(type: BROADCAST)
                   Spawn on: All datanodes
                   -> Seq Scan on t2
(11 rows)
```

The correlated subquery is displayed in the target list (query return list). Values need to be returned even if the condition **t1.c1=t2.c1** is not met. Therefore, use a left outer join to join **t1** and **t2** so that the SSQ can return padding values when the condition **t1.c1=t2.c1** is not met.

◯ NOTE

ScalarSubQuery (SSQ) and Correlated-ScalarSubQuery (CSSQ) are described as follows:

- SSQ: a sublink that returns only a single row and column scalar value
- CSSQ: an SSQ containing conditions

The preceding SQL statement can be changed into:

```
with ssq as
(
    select t2.c1, t2.c2 from t2
)
select ssq.c2, t1.c2
from t1 left join ssq on t1.c1 = ssq.c1
where t1.c2 > 10;
```

The execution plan after the change is as follows:

```
QUERY PLAN

Streaming (type: GATHER)
Node/s: All datanodes
-> Hash Right Join
Hash Cond: (t2.c1 = t1.c1)
-> Seq Scan on t2
-> Hash
-> Seq Scan on t1
```

```
Filter: (c2 > 10)
(8 rows)
```

In the preceding example, the SSQ is pulled up to right join, preventing poor performance caused by the combination of a subplan and broadcast when the table (**T2**) in the subquery is too large.

The subquery (with COUNT) in the target list cannot be pulled up.

Example:

```
select (select count(*) from t2 where t2.c1=t1.c1) cnt, t1.c1, t3.c1 from t1,t3 where t1.c1=t3.c1 order by cnt, t1.c1;
```

The execution plan is as follows:

```
QUERY PLAN
Streaming (type: GATHER)
 Node/s: All datanodes
 -> Sort
     Sort Key: ((SubPlan 1)), t1.c1
     -> Hash Join
         Hash Cond: (t1.c1 = t3.c1)
         -> Seg Scan on t1
         -> Hash
             -> Seq Scan on t3
         SubPlan 1
          -> Aggregate
              -> Result
                  Filter: (t2.c1 = t1.c1)
                  -> Materialize
                     -> Streaming(type: BROADCAST)
                         Spawn on: All datanodes
                          -> Seq Scan on t2
(17 rows)
```

The correlated subquery is displayed in the target list (query return list). Values need to be returned even if the condition **t1.c1=t2.c1** is not met. Therefore, use a left outer join to join **t1** and **t2** so that the SSQ can return padding values when the condition **t1.c1=t2.c1** is not met. However, **COUNT** is used, which requires that **0** is returned when the condition is not met. **case-when NULL then 0 else count(*)** can be used.

The preceding SQL statement can be changed into:

The execution plan after the change is as follows:

```
-> Hash
-> HashAggregate
Group By Key: t2.c1
-> Seq Scan on t2
-> Hash
-> Seq Scan on t3
(15 rows)
```

Pulling up nonequivalent subqueries

Example:

```
select t1.c1, t1.c2
from t1
where t1.c1 = (select agg() from t2.c2 > t1.c2);
```

Nonequivalent subqueries cannot be pulled up. You can perform join twice (one CorrelationKey and one rownum self-join) to rewrite the statement.

You can rewrite the statement in either of the following ways:

Subquery rewriting

CTE rewriting

```
WITH dt as
(
    select t1.rowid, agg() aggref
    from t1,t2
    where t1.c2 > t2.c2 group by t1.rowid
)
select t1.c1, t1.c2
from t1, dt
where t1.rowid = dt.rowid AND
t1.c1 = dt.aggref;
```

NOTICE

- Currently, GaussDB(DWS) does not have an effective way to provide globally unique row IDs for tables and intermediate result sets. Therefore, the rewriting is difficult. It is recommended that this issue is avoided at the service layer or by using t1.xc_node_id + t1.ctid to associate row IDs. However, the high repetition rate of xc_node_id leads to low association efficiency, and xc_node_id+ctid cannot be used as the join condition of hash join.
- If the AGG type is **COUNT(*)**, **0** is used for data padding if **CASE-WHEN** is not matched. If the type is not **COUNT(*)**, **NULL** is used.
- CTE rewriting works better by using share scan.

More Optimization Examples

1. Change the base table to a replication table and create an index on the filter column.

```
create table master_table (a int);
create table sub_table(a int, b int);
select a from master_table group by a having a in (select a from sub_table);
```

In this example, a correlated subquery is contained. To improve the query performance, you can change **sub_table** to a replication table and create an index on the **a** column.

2. Modify the **SELECT** statement, change the subquery to a **JOIN** relationship between the primary table and the parent query, or modify the subquery to improve the query performance. Ensure that the subquery to be used is semantically correct.

```
explain (costs off)select * from master_table as t1 where t1.a in (select t2.a from sub_table as t2 where t1.a = t2.b);

QUERY PLAN

Streaming (type: GATHER)
Node/s: All datanodes

-> Seq Scan on master_table t1
Filter: (SubPlan 1)
SubPlan 1
-> Result
Filter: (t1.a = t2.b)
-> Materialize
-> Streaming(type: BROADCAST)
Spawn on: All datanodes
-> Seq Scan on sub_table t2

(11 rows)
```

In the preceding example, a subplan exists in the plan. To delete the subplan, modify the statement as follows:

In this way, the subplan is replaced by the semi-join between the two tables, greatly improving the execution efficiency.

13.4.7.4 Optimizing Statistics

What Is Statistic Optimization

GaussDB(DWS) generates optimal execution plans based on the cost estimation. Optimizers need to estimate the number of data rows and the cost based on statistics collected using ANALYZE. Therefore, the statistics is vital for the estimation of the number of rows and cost. Global statistics are collected using ANALYZE: relpages and reltuples in the pg_class table; stadistinct, stanullfrac, stanumbersN, stavaluesN, and histogram_bounds in the pg_statistic table.

Example 1: Poor Query Performance Due to the Lack of Statistics

The query performance is often significantly impacted by the absence of statistics for tables or columns involved in the query.

The structure of the example table is as follows:

```
CREATE TABLE LINEITEM
                          NOT NULL
L ORDERKEY
               BIGINT
, L_PARTKEY
              BIGINT
                         NOT NULL
, L SUPPKEY
              BIGINT
                         NOT NULL
, L_LINENUMBER BIGINT
                          NOT NULL
, L_QUANTITY DECIMAL(15,2) NOT NULL
, L_EXTENDEDPRICE DECIMAL(15,2) NOT NULL
, L_DISCOUNT DECIMAL(15,2) NOT NULL
, L_TAX
          DECIMAL(15,2) NOT NULL
, L_RETURNFLAG CHAR(1)
, L_LINESTATUS CHAR(1)
                           NOT NULL
                           NOT NULL
, L_SHIPDATE
              DATE
                         NOT NULL
, L_COMMITDATE DATE
                           NOT NULL
, L_RECEIPTDATE DATE
                         NOT NULL
, L_SHIPINSTRUCT CHAR(25) NOT NULL
, L_SHIPMODE CHAR(10)
                          NOT NULL
L_COMMENT VARCHAR(44) NOT NULL
) with (orientation = column, COMPRESSION = MIDDLE) distribute by hash(L_ORDERKEY);
CREATE TABLE ORDERS
O ORDERKEY
               BIGINT
                          NOT NULL
, O_CUSTKEY
             BIGINT
                         NOT NULL
, O_ORDERSTATUS CHAR(1) NOT NULL
, O_TOTALPRICE DECIMAL(15,2) NOT NULL
, O_ORDERDATE DATE NOT NULL
, O_ORDERPRIORITY CHAR(15) NOT NULL
                        NOT NULL
, O_CLERK
             CHAR(15)
, O_SHIPPRIORITY BIGINT
                         NOT NULL
, O_COMMENT VARCHAR(79) NOT NULL
)with (orientation = column, COMPRESSION = MIDDLE) distribute by hash(O_ORDERKEY);
```

The query statements are as follows:

```
explain verbose select
count(*) as numwait
from
lineitem 11,
orders
where
o_orderkey = l1.l_orderkey
and o_orderstatus = 'F'
and l1.l_receiptdate > l1.l_commitdate
and not exists (
select
from
lineitem 13
where
l3.l_orderkey = l1.l_orderkey
and l3.l_suppkey <> l1.l_suppkey
and l3.l_receiptdate > l3.l_commitdate
order by
numwait desc;
```

You can perform the following operations to check whether **ANALYZE** has been executed on the tables or columns involved in the query to collect statistics.

1. Execute **EXPLAIN VERBOSE** to analyze the execution plan and check the warning information.

WARNING:Statistics in some tables or columns(public.lineitem(l_receiptdate,l_commitdate,l_orderkey, l_suppkey), public.orders(o_orderstatus,o_orderkey)) are not collected. HINT:Do analyze for them in order to generate optimized plan.

2. To determine if poor query performance was caused by a lack of statistics in certain tables or columns, check if the following information exists in the log file located in the **pg_log** directory.

2017-06-14 17:28:30.336 CST 140644024579856 20971684 [BACKEND] LOG:Statistics in some tables or columns(public.lineitem(l_receiptdate, l_commitdate,l_orderkey, .l_suppkey), public.orders(o_orderstatus,o_orderkey)) are not collected. 2017-06-14 17:28:30.336 CST 140644024579856 20971684 [BACKEND] HINT:Do analyze for them in order to generate optimized plan.

After confirming that **ANALYZE** has not been executed on the relevant tables or columns, you can execute **ANALYZE** on the tables or columns reported in the WARNING or logs to resolve the issue of slow query performance due to a lack of statistics

Example 2: Setting cost_param to Optimize Query Performance

For details, see Case: Configuring cost_param for Better Query Performance.

Example 3: Optimization is Not Accurate When Intermediate Results Exist in the Query Where JOIN Is Used for Multiple Tables

Symptom: Query the personnel who have checked in an Internet cafe within 15 minutes before and after the check-in of a specified person.

```
SELECT
C.WBM,
C.DZQH,
C.DZ,
B.ZJHM,
B.SWKSSJ,
B.XWSJ
FROM
b_zyk_wbswxx A,
b_zyk_wbswxx B,
b_zyk_wbcs C
WHERE
A.ZJHM = '522522*****3824'
AND A.WBDM = B.WBDM
AND A.WBDM = C.WBDM
AND abs(to_date(A.SWKSSJ,'yyyymmddHH24MISS') - to_date(B.SWKSSJ,'yyyymmddHH24MISS')) <
INTERVAL '15 MINUTES'
ORDER BY
B.SWKSSJ,
B.ZJHM
limit 10 offset 0
```

Figure 13-6 shows the execution plan. This query takes about 12s.

Figure 13-6 Using an unlogged table (1)

Optimization analysis:

- 1. In the execution plan, index scan is used for node scanning, the **Join Filter** calculation in the external **NEST LOOP IN** statement consumes most of the query time, and the calculation uses the string addition and subtraction, and unequal-value comparison.
- 2. Use an unlogged table to record the Internet access time of the specified person. The start time and end time are processed during data insertion, and this reduces subsequent addition and subtraction operations.

```
//Create a temporary unlogged table.
CREATE UNLOGGED TABLE temp_tsw
ZJHM
          NVARCHAR2(18),
WBDM
           NVARCHAR2(14),
SWKSSJ_START NVARCHAR2(14),
SWKSSJ_END NVARCHAR2(14),
       NVARCHAR2(70),
WRM
DZQH
          NVARCHAR2(6),
D7
         NVARCHAR2(70).
IPDZ
       NVARCHAR2(39)
)
//Insert the Internet access record of the specified person, and process the start time and end time.
INSERT INTO
temp_tsw
SELECT
A.ZJHM,
A.WBDM,
to_char((to_date(A.SWKSSJ,'yyyymmddHH24MISS') - INTERVAL '15
MINUTES'), 'yyyymmddHH24MISS'),
to_char((to_date(A.SWKSSJ,'yyyymmddHH24MISS') + INTERVAL '15
MINUTES'), 'yyyymmddHH24MISS'),
B.WBM,B.DZQH,B.DZ,B.IPDZ
FROM
b zyk wbswxx A,
b_zyk_wbcs B
WHERE
A.ZJHM='522522*****3824' AND A.WBDM = B.WBDM
//Query the personnel who have check in an Internet cafe before and after 15 minutes of the check-in
of the specified person. Convert their ID card number format to int8 in comparison.
```

```
SELECT
A.WBM.
A.DZQH,
A.DZ.
A.IPDZ,
B.ZJHM,
B.XM,
to_date(B.SWKSSJ,'yyyymmddHH24MISS') as SWKSSJ,
to_date(B.XWSJ,'yyyymmddHH24MISS') as XWSJ,
FROM temp_tsw A,
b_zyk_wbswxx B
WHERE
A.ZJHM <> B.ZJHM
AND A.WBDM = B.WBDM
AND (B.SWKSSJ)::int8 > (A.swkssj_start)::int8
AND (B.SWKSSJ)::int8 < (A.swkssj_end)::int8
order by
B.SWKSSJ,
B.ZJHM
limit 10 offset 0
```

The query takes about 7s. Figure 13-7 shows the execution plan.

Figure 13-7 Using an unlogged table (2)

- In the previous plan, Hash Join has been executed, and a Hash table has been created for the large table b_zyk_wbswxx. The table contains large amounts of data, so the creation takes long time.
 - **temp_tsw** contains only hundreds of records, and an equal-value connection is created between **temp_tsw** and **b_zyk_wbswxx** using wbdm (the Internet cafe code). Therefore, if **JOIN** is changed to **NEST LOOP JOIN**, index scan can be used for node scanning, and the performance will be boosted.
- 4. Execute the following statement to change **JOIN** to **NEST LOOP JOIN**. SET enable_hashjoin = off;

Figure 13-8 shows the execution plan. The query takes about 3s.

Figure 13-8 Using an unlogged table (3)

```
Limit (cost=240002336196.14..240002336196.17 rows=10 width=190)

> Sort (cost=240002336196.14..240002336196.17 rows=240 width=190)

Sort Kery: b.swkspj, b.sjhm

>> Stremming (type: GATHER) (cost=240002336190.35..240002336190.95 rows=240 width=190)

Node/s: All datanodes

>> Limit (cost=1000097341.26..1000097341.29 rows=10 width=190)

| Sort Key: b.swkspj, b.sjhm
| >> Nexted Loop (cost=100000970940.26..10000097341.26.width=190)

> Sort Key: b.swkspj, b.sjhm
| >> Nexted Loop (cost=10000000000.00..10000097282.36 rows=2726 width=190)

>> Stremming(type: BROADCAST) (cost=0.00..122.00 rows=240 width=256)

Spawn on: All datanodes
| -> Seg Scan on temp_tsw a (cost=0.00..120.00 rows=240 width=256)

>> Partition therefor (cost=0.00..9648.34 rows=273 width=77)

Iterations: 25

| Sep Scan on temp_tsw a (cost=0.00..101 rows=10 width=256)

-> Partition therefor (cost=0.00..9648.34 rows=273 width=77)

| Totation therefor (cost=0.00..9648.34 rows=273 width=77)

| Totation time ((wbdm)::text (a.wbdm):text) AND ((swkspj)::bigint) (a.swkspj_start)::bigint)
| AND ((swkspj)::bigint (a.swkspj_end)::bigint))
| Selected Partitions: 1..25
```

5. Save the query result set in the unlogged table for paging display.

If paging display needs to be achieved on the upper-layer application page, change the **offset** value to determine the result set on the target page. In this way, the previous query statement will be executed every time after a page turning operation, which causes long response latency.

To resolve this problem, you are advised to use the unlogged table to save the result set.

```
//Create an unlogged table to save the result set.
CREATE UNLOGGED TABLE temp_result
WBM
        NVARCHAR2(70),
DZQH NVARCHAR2(6),
DΖ
    NVARCHAR2(70),
IPDZ NVARCHAR2(39),
ZJHM NVARCHAR2(18),
XM
    NVARCHAR2(30),
SWKSSJ date,
XWSJ date,
SWZDH NVARCHAR2(32)
//Insert the result set to the unlogged table. The insertion takes about 3s.
temp_result
SELECT
A.WBM.
A.DZQH,
A.DZ,
A.IPDZ,
B.ZJHM,
to_date(B.SWKSSJ,'yyyymmddHH24MISS') as SWKSSJ,
to_date(B.XWSJ,'yyyymmddHH24MISS') as XWSJ,
B.SWZDH
FROM temp_tsw A,
b_zyk_wbswxx B
WHERE
A.ZJHM <> B.ZJHM
AND A.WBDM = B.WBDM
AND (B.SWKSSJ)::int8 > (A.swkssj_start)::int8
AND (B.SWKSSJ)::int8 < (A.swkssj_end)::int8
//Perform paging query on the result set. The paging query takes about 10 ms.
SELECT
FROM
temp_result
ORDER BY
SWKSSJ,
```

ZJHM LIMIT 10 OFFSET 0;



Collecting global statistics using ANALYZE improves query performance. If a performance problem occurs, you can use plan hint to adjust the query plan to the previous one. For details, see **Hint-based Tuning**.

13.4.7.5 Tuning Operators

Background

A query goes through many steps to produce its final result. Often, the whole query slows down because one step takes too long. This slow step is called a bottleneck. To fix this, you can run the **EXPLAIN ANALYZE** or **PERFORMANCE** command to find the bottleneck.

In the example below, the **Hashagg** operator takes up 66% of the total time: $(51016 - 13535)/56476 \approx 66\%$. So, the **Hashagg** operator is the bottleneck. Start by optimizing this operator to improve performance.



Operator Tuning Example

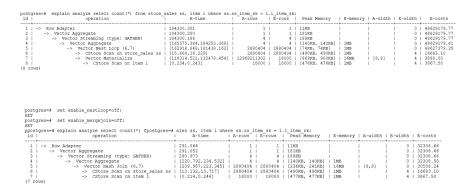
Example 1: When scanning a base table with SeqScan, filtering large amounts of data, like in point queries or range scans, can be slow. Create an index on the condition column and use IndexScan to speed up the process.

```
explain (analyze on, costs off) select * from store sales where ss sold_date sk = 2450944;
          operation | A-time | A-rows | Peak Memory | A-width
id |
                          ---+----+----
 1 | -> Streaming (type: GATHER) | 3666.020
                                           | 3360 | 195KB
 2 | -> Seq Scan on store_sales | [3594.611,3594.611] | 3360 | [34KB, 34KB] |
Predicate Information (identified by plan id)
 2 -- Seq Scan on store_sales
     Filter: (ss_sold_date_sk = 2450944)
     Rows Removed by Filter: 4968936
create index idx on store_sales_row(ss_sold_date_sk);
CREATE INDEX
explain (analyze on, costs off) select * from store_sales_row where ss_sold_date_sk = 2450944;
id |
            operation
                         | A-time | A-rows | Peak Memory | A-width
1 | -> Streaming (type: GATHER) | 81.524 | 3360 | 195KB |
2 | -> Index Scan using idx on store sales row | [13.352,13.352] | 3360 | [34KB, 34KB] |
```

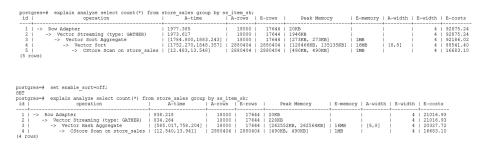
For instance, a full table scan returned 3,360 records but took 3.6 seconds. After indexing the **ss sold date sk** column, the scan time dropped to 13 milliseconds.

Example 2: If **NestLoop** is chosen to join two tables and the row count is high, the join can be slow. In one case, **NestLoop** took 181 seconds. By setting

enable_mergejoin and **enable_nestloop** to **off** and letting the optimizer choose **HashJoin**, the join time was reduced to over 200 milliseconds.



Example 3: **HashAgg** usually performs better. For large result sets, if **Sort** and **GroupAgg** are used, set **enable_sort** to off. **HashAgg** takes much longer time than **Sort** and **GroupAgg** used together.



13.4.7.6 Optimizing Data Skew

Data skew breaks the balance among nodes in the distributed MPP architecture. If the amount of data stored or processed by a node is much greater than that by other nodes, the following problems may occur:

- Storage skew severely limits the system capacity. The skew on a single node hinders system storage utilization.
- Computing skew severely affects performance. The data to be processed on the skew node is much more than that on other nodes, deteriorating overall system performance.
- Data skew severely affects the scalability of the MPP architecture. During storage or computing, data with the same values is often placed on the same node. Therefore, even if you add nodes after a data skew occurs, the skew data (data with the same values) is still placed on the node and affects the system capacity or performance bottleneck.

GaussDB(DWS) provides a complete solution for data skew, including storage and computing skew.

Data Skew in the Storage Layer

In the GaussDB(DWS) database, data is distributed and stored on each DN. You can improve the query efficiency by using distributed execution. However, if data skew occurs, bottlenecks exist on some DNs during distribution execution,

affecting the query performance. This is because the distribution column is not properly selected. This can be solved by adjusting the distribution column.

For example:

```
explain performance select count(*) from inventory;
5 -- CStore Scan on Imz.inventory
     dn 6001 6002 (actual time=0.444..83.127 rows=42000000 loops=1)
     dn_6003_6004 (actual time=0.512..63.554 rows=27000000 loops=1)
     dn_6005_6006 (actual time=0.722..99.033 rows=45000000 loops=1)
     dn_6007_6008 (actual time=0.529..100.379 rows=51000000 loops=1)
     dn_6009_6010 (actual time=0.382..71.341 rows=36000000 loops=1)
     dn_6011_6012 (actual time=0.547..100.274 rows=51000000 loops=1)
     dn_6013_6014 (actual time=0.596..118.289 rows=60000000 loops=1)
     dn_6015_6016 (actual time=1.057..132.346 rows=63000000 loops=1)
     dn_6017_6018 (actual time=0.940..110.310 rows=54000000 loops=1)
     dn_6019_6020 (actual time=0.231..41.198 rows=21000000 loops=1)
     dn_6021_6022 (actual time=0.927..114.538 rows=54000000 loops=1)
     dn_6023_6024 (actual time=0.637..118.385 rows=60000000 loops=1)
     dn_6025_6026 (actual time=0.288..32.240 rows=15000000 loops=1)
     dn_6027_6028 (actual time=0.566..118.096 rows=60000000 loops=1)
     dn 6029 6030 (actual time=0.423..82.913 rows=42000000 loops=1)
     dn_6031_6032 (actual time=0.395..78.103 rows=39000000 loops=1)
     dn 6033 6034 (actual time=0.376..51.052 rows=24000000 loops=1)
     dn_6035_6036 (actual time=0.569..79.463 rows=39000000 loops=1)
```

In the performance information, you can view the number of scan rows of each DN in the inventory table. The number of rows of each DN differs a lot, the biggest is 63000000 and the smallest value is 15000000. This value difference on the performance of data scan is acceptable, but if the join operator exists in the upper-layer, the impact on the performance cannot be ignored.

Generally, the data table is hash distributed on each DN; therefore, it is important to choose a proper distribution column. Use the **table_skewness ()** function to check the data skew of the **inventory** table on each DN. The query result is as follows:

```
select table_skewness('inventory');
        table_skewness
("dn_6015_6016
                    ",63000000,8.046%)
("dn 6013 6014
                    ",60000000,7.663%)
                    ",60000000,7.663%)
("dn_6023_6024
("dn 6027 6028
                    ",60000000,7.663%)
("dn_6017_6018
                    ",54000000,6.897%)
("dn_6021_6022
                    ".54000000.6.897%)
                    ",51000000,6.513%)
("dn 6007 6008
                    ",51000000,6.513%)
("dn_6011_6012
("dn 6005 6006
                    ",45000000,5.747%)
                    ",42000000,5.364%)
("dn_6001_6002
                    ",42000000,5.364%)
("dn_6029_6030
("dn_6031_6032
                    ",39000000,4.981%)
("dn 6035 6036
                    ",39000000,4.981%)
                    ",36000000,4.598%)
("dn 6009 6010
                    ",27000000,3.448%)
("dn 6003 6004
("dn_6033_6034
                     ,24000000,3.065%)
                    ",21000000,2.682%)
("dn 6019 6020
("dn_6025_6026
                    ",15000000,1.916%)
(18 rows)
```

The table definition indicates that the table uses the <code>inv_date_sk</code> column as the distribution column, which causes a data skew. Based on the data distribution of each column, change the distribution column to <code>inv_item_sk</code>. The skew status is as follows:

```
select table_skewness('inventory');
table_skewness
```

```
",43934200,5.611%)
("dn_6001_6002
("dn_6007_6008
                   ",43829420,5.598%)
("dn_6003_6004
                   ".43781960,5.592%)
                   ",43773880,5.591%)
("dn_6031_6032
                   ",43763280,5.589%)
("dn_6033_6034
                   ",43683600,5.579%)
("dn_6011_6012
("dn_6013_6014
                    ",43551660,5.562%)
                   ",43546340,5.561%)
("dn_6027_6028
                   ",43508700,5.557%)
("dn_6009_6010
("dn_6023_6024
                   ".43484540.5.554%)
                    ",43466800,5.551%)
("dn 6019 6020
                   ",43458500,5.550%)
("dn 6021 6022
                   ",43448040,5.549%)
("dn_6017_6018
("dn_6015_6016
                    ",43247700,5.523%)
                   ",43200240,5.517%)
("dn 6005 6006
                   ",43181360,5.515%)
("dn_6029_6030
("dn_6025_6026
                   ",43179700,5.515%)
                    ",42960080,5.487%)
("dn_6035_6036
(18 rows)
```

Data skew is solved.

In addition to the **table_skewness()** function, you can use the **table_distribution** function and the **PGXC_GET_TABLE_SKEWNESS** view to efficiently query the data skew status of each table.

Data Skew in the Computing Layer

Even if data is balanced across nodes after you change the distribution key of a table, data skew may still occur during a query. If data skew occurs in the result set of an operator on a DN, skew will also occur during the computing that involves the operator. Generally, this is caused by data redistribution during the execution.

During a query, JOIN keys and GROUP BY keys are not used as distribution columns. Data is redistributed among DNs based on the hash values of data on the keys. The redistribution is implemented using the Redistribute operator in an execution plan. Data skew in redistribution columns can lead to data skew during system operation. After the redistribution, some nodes will have much more data, process more data, and will have much lower performance than others.

In the following example, the **s** and **t** tables are joined, and **s.x** and **t.x** columns in the join condition are not their distribution keys. Table data is redistributed using the **REDISTRIBUTE** operator. Data skew occurs in the **s.x** column and not in the **t.x** column. The result set of the **Streaming** operator (**id** being **6**) on datanode2 has data three times that of other DNs and causes a skew.

```
select * from skew s,test t where s.x = t.x order by s.a limit 1;
id l
                operation
                                      | A-time
1 | -> Limit
                                       | 52622.382
2 | -> Streaming (type: GATHER)
                                               | 52622.374
                                       | [30138.494,52598.994]
       -> Limit
                                       [30138.486,52598.986]
 4 |
         -> Sort
 5 |
           -> Hash Join (6,8)
                                          | [30127.013,41483.275]
             -> Streaming(type: REDISTRIBUTE) | [11365.110,22024.845]
 6
              -> Seq Scan on public.skew s | [2019.168,2175.369]
 7 I
 8 I
             -> Hash
                                        [2460.108,2499.850]
               -> Streaming(type: REDISTRIBUTE) | [1056.214,1121.887]
 91
                  -> Seq Scan on public.test t | [310.848,325.569]
10 I
6 -- Streaming(type: REDISTRIBUTE)
     datanode1 (rows=5050368)
```

```
datanode2 (rows=15276032)
datanode3 (rows=5174272)
datanode4 (rows=5219328)
```

Computing skew is more difficult to detect than storage skew. To solve computing skew, GaussDB provides the Runtime Load Balance Technology (RLBT) solution, controlled by the **skew_option** parameter. The RLBT solution addresses how to detect and solve data skew.

Detect data skew.

The solution first checks whether skew data exists in redistribution columns used for computing. RLBT can detect data skew based on statistics, specified hints, or rules.

Detection based on statistics

Run the **ANALYZE** statement to collect statistics on tables. The optimizer will automatically identify skew data on redistribution keys based on the statistics and generate optimization plans for queries having potential skew. When the redistribution key has multiple columns, statistics information can be used for identification only when all columns belong to the same base table.

The statistics information can only provide the skew of the base table. If a column in the base table is skewed, or other columns have filtering conditions, or after the join of other tables, we cannot determine whether the skewed data still exists on the skewed column. If **skew_option** is **normal**, it indicates that the skew data still exists, and the base tables will be optimized to solve skew. If **skew_option** is **lazy**, it indicates that no more skew data exists and the optimization will stop.

Detection based on specified hints

The intermediate results of complex queries are difficult to estimate based on statistics. In this case, you can specify hints to provide the skew information, based on which the optimizer optimizes queries. For details about the syntax of hints, see **Skew Hints**.

Detection based on rules

In a business intelligence (BI) system, a large number of SQL statements having outer joins (including left joins, right joins, and full joins) are generated, and many NULL values will be generated in empty columns that have no match for outer joins. If JOIN or GROUP BY operations are performed on the columns, data skew will occur. RLBT can automatically identify this scenario and generate an optimization plan for NULL value skew.

Solve computing skew.

Join and **Aggregate** operators are optimized to solve skew.

- Join optimization

Skew and non-skew data is separately processed. Details are as follows:

a. When redistribution is required on both sides of a join:

Use **PART_REDISTRIBUTE_PART_ROUNDROBIN** on the side with skew. Specifically, perform round-robin on skew data and redistribution on non-skew data.

Use **PART_REDISTRIBUTE_PART_BROADCAST** on the side with no skew. Specifically, perform broadcast on skew data and redistribution on non-skew data.

b. When redistribution is required on only one side of a join:

Use **PART_REDISTRIBUTE_PART_ROUNDROBIN** on the side where redistribution is required.

Use **PART_LOCAL_PART_BROADCAST** on the side where redistribution is not required. Specifically, perform broadcast on skew data and retain other data locally.

c. When a table has **NULL** values padded:

Use **PART_REDISTRIBUTE_PART_LOCAL** on the table. Specifically, retain the **NULL** values locally and perform redistribution on other data.

In the example query, the **s.x** column contains skewed data and its value is **0**. The optimizer identifies the skew data in statistics and generates the following optimization plan:

```
id |
                       operation
                                                          A-time
1 | -> Limit
                                                    | 23642.049
                                                            23642.041
2 |
   -> Streaming (type: GATHER)
3 |
       -> Limit
                                                    | [23310.768,23618.021]
                                                    | [23310.761,23618.012]
         -> Sort
4 |
           -> Hash Join (6,8)
                                                       | [20898.341,21115.272]
5 |
             -> Streaming(type: PART REDISTRIBUTE PART ROUNDROBIN) |
6 I
[7125.834,7472.111]
7 |
             -> Seg Scan on public.skew s
                                                          | [1837.079,1911.025]
                                                    | [2612.484,2640.572]
8 |
             -> Hash
9
               -> Streaming(type: PART REDISTRIBUTE PART BROADCAST) | [1193.548,1297.894]
                                                         | [314.343,328.707]
10 I
                 -> Seq Scan on public.test t
 5 -- Vector Hash Join (6,8)
     Hash Cond: s.x = t.x
     Skew Join Optimized by Statistic
 6 -- Streaming (type: PART REDISTRIBUTE PART ROUNDROBIN)
     datanode1 (rows=7635968)
     datanode2 (rows=7517184)
     datanode3 (rows=7748608)
     datanode4 (rows=7818240)
```

In the preceding execution plan, **Skew Join Optimized by Statistic** indicates that this is an optimized plan used for handling data skew. The **Statistic** keyword indicates that the plan optimization is based on statistics; **Hint** indicates that the optimization is based on hints; **Rule** indicates that the optimization is based on rules. In this plan, skew and non-skew data is separately processed. Non-skew data in the **s** table is redistributed based on its hash values, and skew data (whose value is **0**) is evenly distributed on all nodes in round-robin mode. In this way, data skew is solved.

To ensure result correctness, the **t** table also needs to be processed. In the **t** table, the data whose value is **0** (skew value in the **s.x** table) is broadcast and other data is redistributed based on its hash values.

In this way, data skew in JOIN operations is solved. The above result shows that the output of the **Streaming** operator (**id** being **6**) is balanced and the end-to-end performance of the query is doubled.

If the stream operator type in the execution plan is **HYBRID**, the stream mode varies depending on the skew data. The following plan is an example:

EXPLAIN (nodes OFF, costs OFF) SELECT COUNT(*) FROM skew_scol s, skew_scol1 s1 WHERE s.b = s1.c:

```
QUERY PLAN
id |
                                                  operation
1 | -> Aggregate
2 | -> Streaming (type: GATHER)
     -> Aggregate
     -> Hash Join (5,7)
4 |
       -> Streaming(type: HYBRID)
5 |
6
          -> Seq Scan on skew_scol s
         -> Hash
7 |
8 |
           -> Streaming(type: HYBRID)
             -> Seq Scan on skew_scol1 s1
9 |
Predicate Information (identified by plan id)
4 -- Hash Join (5,7)
Hash Cond: (s.b = s1.c)
Skew Join Optimized by Statistic
5 --Streaming(type: HYBRID)
Skew Filter: (b = 1)
Skew Filter: (b = 0)
8 -- Streaming (type: HYBRID)
Skew Filter: (c = 0)
Skew Filter: (c = 1)
```

Data 1 has skew in the **skew_scol** table. Perform **ROUNDROBIN** on skew data and **REDISTRIBUTE** on non-skew data.

Data 0 is the side with no skew in the **skew_scol** table. Perform **BROADCAST** on skew data and **REDISTRIBUTE** on non-skew data.

As shown in the preceding figure, the two stream types are **PART REDISTRIBUTE PART ROUNDROBIN** and **PART REDISTRIBUTE PART BROADCAST**. In this example, the stream type is **HYBRID**.

Aggregate optimization

For aggregation, data on each DN is deduplicated based on the **GROUP BY** key and then redistributed. After the deduplication on DNs, the global occurrences of each value will not be greater than the number of DNs. Therefore, no serious data skew will occur. Take the following query as an example:

select c1, c2, c3, c4, c5, c6, c7, c8, c9, count(*) from t group by c1, c2, c3, c4, c5, c6, c7, c8, c9 limit 10;

The command output is as follows:

A large amount of skew data exists. As a result, after data is redistributed based on its **GROUP BY** key, the data volume of datanode1 is hundreds of thousands of times that of others. After optimization, a GROUP BY operation

id | operation A-time 1 | -> Streaming (type: GATHER) | 10961.337 | [10953.014,10953.705] 2 | -> HashAggregate -> HashAggregate [10952.957,10953.632] 3 | 4 -> Streaming(type: REDISTRIBUTE) | [10952.859,10953.502] -> HashAggregate | [10084.280,10947.139] 5 I -> Seq Scan on public.t | [4757.031,5201.168] 6 | Predicate Information (identified by plan id) 3 -- HashAggregate Skew Agg Optimized by Statistic 4 -- Streaming(type: REDISTRIBUTE) datanode1 (rows=17)

is performed on the DN to deduplicate data. After redistribution, no data skew occurs.

Applicable scope

Join operator

datanode2 (rows=8) datanode3 (rows=8) datanode4 (rows=14)

- nest loop, merge join, and hash join can be optimized.
- If skew data is on the left to the join, inner join, left join, semi join, and anti join are supported. If skew data is on the right to the join, inner join, right join, right semi join, and right anti join are supported.
- For an optimization plan generated based on statistics, the optimizer checks whether it is optimal by estimating its cost. Optimization plans based on hints or rules are forcibly generated.
- Aggregate operator
 - array_agg, string_agg, and subplan in agg qual cannot be optimized.
 - A plan generated based on statistics is affected by its cost, the plan_mode_seed parameter, and the best_agg_plan parameter. A plan generated based on hints or rules are not affected by them.

13.4.7.7 Proactive Preheating and Tuning of Disk Cache

This function is supported only in 9.1.0.200 or later.

Overview

In the storage-compute decoupling architecture, user data is stored in OBS to reduce storage costs. Each query generates network I/Os to retrieve data from OBS. To improve query speed, reduce storage costs, and minimize performance loss, the architecture provides disk cache capability. Pre-queried data is cached locally, enhancing performance.

The LRU2Q algorithm manages the disk cache with three queues: A1in, A1out, and Am. Data first enters A1in. If A1in is full (adjustable via a GUC parameter,

default is 0.25 times the total queue size), data moves to A1out. Data enters Am only when hit in A1out. Am queue holds the hottest data.

For common queries, LRU2Q is sufficient. However, frequent joins of large and small tables can degrade performance if small tables are frequently evicted from A1in to A1out.

Tuning Syntax

A new tuning policy allows directly adding small table data to Am, ensuring it remains hot and reducing network I/Os during joins. The syntax formats are as follows:

1. Perform the actual query and add data directly to Am. explain warmup hot select ...;

2. Query data in the sequence A1in > A1out > Am.

explain warmup select ...;

3. No actual query operation is performed, and the **explain** logic remains unchanged.

explain select ...;

13.4.7.8 SQL Statement Rewriting Rules

Based on the database SQL execution mechanism and a large number of practices, summarize finds that: using rules of a certain SQL statement, on the basis of the so that the correct test result, which can improve the SQL execution efficiency. You can comply with these rules to greatly improve service query efficiency.

Replacing UNION with UNION ALL

UNION eliminates duplicate rows while merging two result sets but **UNION ALL** merges the two result sets without deduplication. Therefore, replace **UNION** with **UNION ALL** if you are sure that the two result sets do not contain duplicate rows based on the service logic.

Adding NOT NULL to the join column

If there are many **NULL** values in the **JOIN** columns, you can add the filter criterion **IS NOT NULL** to filter data in advance to improve the **JOIN** efficiency.

• Converting **NOT IN** to **NOT EXISTS**

nestloop anti join must be used to implement **NOT IN**, and **Hash anti join** is required for **NOT EXISTS**. If no **NULL** value exists in the **JOIN** column, **NOT IN** is equivalent to **NOT EXISTS**. Therefore, if you are sure that no **NULL** value exists, you can convert **NOT IN** to **NOT EXISTS** to generate **hash joins** and to improve the query performance.

As shown in the following figure, the **t2.d2** column does not contain null values (it is set to **NOT NULL**) and **NOT EXISTS** is used for the query.

SELECT * FROM t1 WHERE NOT EXISTS (SELECT * FROM t2 WHERE t1.c1=t2.d2);

The generated execution plan is as follows:

Figure 13-9 NOT EXISTS execution plan

```
id I
                    operation
  1 | -> Streaming (type: GATHER)
  2
      -> Hash Right Anti Join (3, 5)
  3
            -> Streaming(type: REDISTRIBUTE)
  4
             -> Seq Scan on t2
  5 |
            -> Hash
               -> Seq Scan on t1
Predicate Information (identified by plan id)
  2 -- Hash Right Anti Join (3, 5)
       Hash Cond: (t2.d2 = t1.c1)
(13 rows)
```

Use hashagg.

If a plan involving groupAgg and SORT operations generated by the **GROUP BY** statement is poor in performance, you can set **work_mem** to a larger value to generate a **hashagg** plan, which does not require sorting and improves the performance.

• Replace functions with **CASE** statements

The GaussDB(DWS) performance greatly deteriorates if a large number of functions are called. In this case, you can modify the pushdown functions to **CASE** statements.

Do not use functions or expressions for indexes.

Using functions or expressions for indexes stops indexing. Instead, it enables scanning on the full table.

 Do not use != or <> operators, NULL, OR, or implicit parameter conversion in WHERE clauses.

Split complex SQL statements.

You can split an SQL statement into several ones and save the execution result to a temporary table if the SQL statement is too complex to be tuned using the solutions above, including but not limited to the following scenarios:

- The same subquery is involved in multiple SQL statements of a task and the subquery contains large amounts of data.
- Incorrect Plan cost causes a small hash bucket of subquery. For example, the actual number of rows is 10 million, but only 1000 rows are in hash bucket.
- Functions such as substr and to_number cause incorrect measures for subqueries containing large amounts of data.
- BROADCAST subqueries are performed on large tables in multi-DN environment.

13.4.8 Configuring Optimizer Parameters

This section introduces key CN parameters that affect optimization of SQL statements in GaussDB(DWS). For details about the parameter configuration method, see **Configuring GUC Parameters**.

Table 13-17 CN parameters

Parameter/ Reference Value	Description			
enable_nestloop=o n	Specifies how the optimizer uses Nest Loop Join . If this parameter is set to on , the optimizer preferentially uses Nest Loop Join . If it is set to off , the optimizer preferentially uses other methods, if any. NOTE To temporarily change the value of this parameter in the current database connection (that is, the current session), run the following SQL statement: SET enable_nestloop to off;			
	By default, this parameter is set to on . Change the value as required. Generally, nested loop join has the poorest performance among the three JOIN methods (nested loop join, merge join, and hash join). You are advised to set this parameter to off .			
enable_bitmapscan =on	Specifies whether the optimizer uses bitmap scanning. I the value is on , bitmap scanning is used. If the value is off , it is not used. NOTE If you only want to temporarily change the value of this parameter during the current database connection (that is, the current session), run the following SQL statements: SET enable_bitmapscan to off;			
	The bitmap scanning applies only in the query condition where a > 1 and b > 1 and indexes are created on columns a and b . During performance tuning, if the query performance is poor and bitmapscan operators are in the execution plan, set this parameter to off and check whether the performance is improved.			
enable_fast_query_ shipping=on	Specifies whether the optimizer uses a distribution framework. If the value is on , the execution plan is generated on both CNs and DNs. If the value is off , the distribution framework is used, that is, the execution plan is generated on the CNs and then sent to DNs for execution.			
	NOTE To temporarily change the value of this parameter in the current database connection (that is, the current session), run the following SQL statement: SET enable_fast_query_shipping to off;			

Parameter/ Reference Value	Description
enable_hashagg=o n	Specifies whether to enable the optimizer's use of Hashaggregation plan types.
enable_hashjoin=o n	Specifies whether to enable the optimizer's use of Hash- join plan types.
enable_mergejoin= on	Specifies whether to enable the optimizer's use of Hashmerge plan types.
enable_indexscan= on	Specifies whether to enable the optimizer's use of index- scan plan types.
enable_indexonlysc an=on	Specifies whether to enable the optimizer's use of index- only-scan plan types.
enable_seqscan=on	Specifies whether the optimizer uses bitmap scanning. It is impossible to suppress sequential scans entirely, but setting this variable to off allows the optimizer to preferentially choose other methods if available.
enable_sort=on	Specifies the optimizer sorts. It is impossible to fully suppress explicit sorts, but setting this variable to off allows the optimizer to preferentially choose other methods if available.
enable_broadcast= on	Specifies whether enable the optimizer's use of data broadcast. In data broadcast, a large amount of data is transferred on the network. When the number of transmission nodes (stream) is large and the estimation is inaccurate, set this parameter to off and check whether the performance is improved.
enable_redistribute =on This parameter is supported only by clusters of version 8.2.1.300 or later.	Controls the query optimizer's use of data transmission in local redistribute and split redistribute redistribution modes. This parameter corresponds to enable_broadcast. The optimizer may overestimate the cost of local broadcast and split broadcast. As a result, the optimizer selects the local redistribute or split redistribute redistribution plan. This may cause performance deterioration. Therefore, when the actual data volume of the network transmission node (stream) is small, you can set this parameter to off so that the optimizer preferentially selects the broadcast mode. Then you can check whether the performance is improved.
rewrite_rule	Specifies whether the optimizer enables a specific rewriting rule.

13.4.9 Hint-based Tuning

13.4.9.1 Plan Hint Optimization

In plan hints, you can specify a join order, join, stream, and scan operations, the number of rows in a result, and redistribution skew information to tune an execution plan, improving query performance.

Function

Plan hints can be specified using the keywords such as **SELECT**, **INSERT**, **UPDATE**, **MERGE**, and **DELETE**, in the following format:

/*+ <plan hint> */

You can specify multiple hints for a query plan and separate them by spaces. A hint specified for a query plan does not apply to its subquery plans. To specify a hint for a subquery, add the hint following the keyword of this subquery.

For example:

select /*+ <plan_hint1> <plan_hint2> */ * from t1, (select /*+ <plan_hint3> */ from t2) where 1=1;

In the preceding command, <plan_hint1> and <plan_hint2> are the hints of a query, and <plan_hint3> is the hint of its subquery.

NOTICE

If a hint is specified in the **CREATE VIEW** statement, the hint will be applied each time this view is used.

If the random plan function is enabled (**plan_mode_seed** is set to a value other than 0), the specified hint will not be used.

Supported Hints

Currently, the following hints are supported:

- Join order hints (leading)
- Join operation hints, excluding the semi join, anti join, and unique plan hints
- Rows hints
- Stream operation hints
- Scan operation hints, supporting only tablescan, indexscan, and indexonlyscan
- Sublink name hints
- Skew hints, supporting only the skew in the redistribution involving Join or HashAgg
- Hint used for **Agg** distribution columns Only clusters of 8.1.3.100 and later versions support this function.
- Hint that disables subquery pull-up. Only clusters of 8.2.0 and later versions support this function.
- Configuration parameter hints. For details about supported parameters, see
 Configuration Parameter Hints.

Precautions

- Sort, Setop, and Subplan hints are not supported.
- Hints do not support SMP or Node Groups.
- Hints cannot be used for the target table of the **INSERT** statement.

Examples

The following is the original plan and is used for comparing with the optimized ones:

```
explain
select i_product_name product_name
,i_item_sk item_sk
,s_store_name store_name
,s_zip store_zip
,ad2.ca_street_number c_street_number
,ad2.ca_street_name c_street_name
,ad2.ca_city c_city
,ad2.ca_zip c_zip
,count(*) cnt
,sum(ss_wholesale_cost) s1
,sum(ss_list_price) s2
,sum(ss_coupon_amt) s3
FROM store_sales
,store_returns
,store
,customer
,promotion
,customer_address ad2
,item
WHERE ss_store_sk = s_store_sk AND
ss_customer_sk = c_customer_sk AND
ss_item_sk = i_item_sk and
ss_item_sk = sr_item_sk and
ss_ticket_number = sr_ticket_number and
c_current_addr_sk = ad2.ca_address_sk and
ss_promo_sk = p_promo_sk and
i_color in ('maroon','burnished','dim','steel','navajo','chocolate') and
i_current_price between 35 and 35 + 10 and
i_current_price between 35 + 1 and 35 + 15
group by i_product_name
,i_item_sk
,s_store_name
,s_zip
,ad2.ca_street_number
,ad2.ca_street_name
,ad2.ca_city
,ad2.ca_zip
```

id	operation	I	E-rows	ı	E-memory	E-width	1	E-costs
1	·	ı	6	ī		273	1	3401632.49
2	-> Vector Streaming (type: GATHER)	1	6	Ī		273	1	3401632.49
3	-> Vector Hash Aggregate	1	6	Ī	16MB	273	1	3401630.82
4	-> Vector Streaming(type: REDISTRIBUTE)	1	6	Ī	1MB	169	1	3401630.78
5	-> Vector Hash Join (6,21)	1	6	Ī	16MB	169	1	3401630.42
6	-> Vector Hash Join (7,20)	1	7	I	43MB	173	1	3400343.15
7	-> Vector Streaming(type: REDISTRIBUTE)	1	7	I	1MB	123	1	3395775.64
8	-> Vector Hash Join (9,19)	1	7	I	27MB	123	1	3395775.48
9	-> Vector Streaming(type: REDISTRIBUTE)	1	7	I	1MB	123	1	3386294.72
10	-> Vector Hash Join (11,18)	1	7	I	16MB	123	1	3386294.56
11	-> Vector Hash Join (12,14)	1	7	I	19MB	112	1	3384018.02
12	-> Vector Partition Iterator	1	287999764	I	1MB	12	1	227383.99
13	-> Partitioned CStore Scan on store_returns	1	287999764	I	1MB	12	1	227383.99
14	-> Vector Hash Join (15,17)	1	1516824	I	16MB	124	1	3065686.08
15	-> Vector Partition Iterator	1	2879987999	I	1MB	66	1	2756066.50
16	-> Partitioned CStore Scan on store_sales	1	2879987999	I	1MB	66	1	2756066.50
17	-> CStore Scan on item	1	158	I	1MB	58	1	4051.25
18	-> CStore Scan on store	1	24048	I	1MB	19	1	2264.00
19	-> CStore Scan on customer	1	12000000	I	1MB	1 8	1	12923.00
20	-> CStore Scan on customer_address ad2	1	6000000	I	1MB	58	1	5770.00
21	-> CStore Scan on promotion	1	36000	I	1MB	1 4	1	1268.50
(21	rows)							

13.4.9.2 Join Order Hints

Function

Theses hints specify the join order and outer/inner tables.

Syntax

• Single-layer parentheses (), specifying only the join order. The order of internal and foreign tables is not specified.

```
leading(join_table_list)
leading(@block_name join_table_list)
```

• Double parentheses (()), specifying the join order and outer/inner tables. The outer/inner tables are specified by the outermost parentheses.

```
leading((join_table_list))
leading(@block_name (join_table_list))
```

• Single-layer square brackets [], specifying the join order of [] and the sequence of internal and foreign tables.

```
leading[join_table_list]
leading[@block_name join_table_list]
```

Combination of single-layer parentheses () and single-layer square brackets
[], specifying the join order and the sequence of internal and foreign tables at
any layer. The parentheses () specify only the join order, but not the sequence
of internal and foreign tables. The square brackets [] specify both the join
order and the sequence of internal and foreign tables.

```
leading(join_table_list1 [join_table_list2])
leading[join_table_list1 [join_table_list2]]
leading[join_table_list1 (join_table_list2)]
leading(@block_name join_table_list1 [join_table_list2])
leading(@block_name join_table_list1 [join_table_list2]]
leading(@block_name join_table_list1 (join_table_list2)]
```

NOTICE

Single-layer square brackets [] can be used together with single-layer parentheses () to specify the sequence of internal and foreign tables of any layer. Single-layer [] and double-layer () cannot be used together.

Parameter Description

join_table_list

Specifies the tables to be joined. The values can be table names or table aliases. If a subquery is pulled up, the value can also be the subquery alias. Separate the values with spaces. You can add parentheses to specify the join priorities of tables.

To prevent semantic errors, tables in the list must meet the following requirements:

- The tables must exist in the query or its subquery to be pulled up.
- The table names must be unique in the query or subquery to be pulled up. If they are not, their aliases must be unique.
- A table appears only once in the list.
- An alias (if any) is used to represent a table.

 The syntax format of the table is as follows: [schema.]table[@block_name]

The table name can contain the schema name or block name before the subquery statement block is promoted. If the subquery statement block is optimized and rewritten by the optimizer, the value of **block_name** is different from that of **block_name** in leading.

• If a table has an alias, the alias is preferentially used to represent the table.

block_name

Specifies the block name of the statement block. It indicates that the hint takes effect in the subquery statement block corresponding to the block name.

NOTICE

- By default, a block name is generated for a statement.
- CN lightweight statements do not generate block names.
- Block names can be generated for the CREATE TABLE AS SELECT, SELECT INTO, SELECT, INSERT, UPDATE, DELETE, and MERGE statements.
- The naming rule of a block name is as follows:
 - A block name is automatically generated for the SELECT, INSERT, UPDATE, DELETE, and MERGE statements. The naming format of a block name for these statements is sel\$n, ins\$n, upd\$n, del\$n, and mer \$n, respectively, where n starts from 1. The number of statements of different types is not accumulated, but the number of statements of the same type is accumulated.

Example:

INSERT INTO t SELECT * FROM t1 WHERE a1 IN (select * from t2);	
col#2	
sel\$2	
sel\$1	
36.41	
ins\$1	

 Recursively assigns a block name to each statement block before the optimizer is used.

First, assign block names to the existing statements block based on the statement type, then traverse the statement blocks in the following sequence, and assign block name to the statement blocks in the statement blocks:

- 1. Traverse the target column.
- 2. Traverse the target column in the source table of the MERGE statement.
- 3. Traverse actions (update or insert) in the MERGE statement.
- 4. Traverse the returning clause.
- 5. Traverse the Join and Where conditions in From. (The join condition takes precedence over the Where condition.)
- 6. For a set operation, traverse each branch of the set (UNION, INTERSECT, and EXCEPT).
- 7. Traverse the HAVING clause.
- 8. Traverse the LIMIT OFFSET clause.
- 9. Traverse the LIMIT COUNT clause.
- 10.Traverse CTE
- 11.Traverse the table after From.
- 12. Traverse the UPSERT clause.
- In the rewriting phase of the optimizer, rewriting optimization is performed due to FUL LJOIN, cte inline, materialized view rewriting, INLIST2JOIN, OR conversion, multi count(distinct), Magic Set, lazyagg, and subquery/sublink promotion, a new subquery is constructed. In this case, the recursive processing during block name assignment is also applied to the newly constructed subquery. The number of block names is accumulated.

• In the optimizer rewriting phase, when a subquery is promoted, the table in the inner subquery is promoted to the outer query, and the inner subquery is eliminated. In this case, the promoted table may have the same name as the table in the outer queries. Therefore, the block name to which the promoted table belongs is recorded in the table to distinguish two tables with the same name but are from different query blocks.

For example:

leading(t1 t2 t3 t4 t5): **t1**, **t2**, **t3**, **t4**, and **t5** are joined. The join order and outer/inner tables are not specified.

leading(t1 t2 t3 t4 t5): **t1**, **t2**, **t3**, **t4**, and **t5** are joined in sequence. The table on the right is used as the inner table in each join.

leading(t1 (t2 t3 t4) t5): First, **t2**, **t3**, and **t4** are joined and the outer/inner tables are not specified. Then, the result is joined with **t1** and **t5**, and the outer/inner tables are not specified.

leading(t1 (t2 t3 t4) t5): First, **t2**, **t3**, and **t4** are joined and the outer/inner tables are not specified. Then, the result is joined with **t1**, and **(t2 t3 t4)** is used as the inner table. Finally, the result is joined with **t5**, and **t5** is used as the inner table.

leading((t1 (t2 t3) t4 t5)) leading((t3 t2)): First, **t2** and **t3** are joined and **t2** is used as the inner table. Then, the result is joined with **t1**, and **(t2 t3)** is used as the inner table. Finally, the result is joined with **t4** and then **t5**, and the table on the right in each join is used as the inner table.

leading[t1 [t2 t3]] is equivalent to leading((t1 (t2 t3))) leading((t2 t3)).

leading(t1 [t2 t3]) is equivalent to leading(t1 t2 t3) leading((t2 t3)).

leading[@sel\$1 t1@sel\$1 [t2@sel\$2 t3@sel\$2]] indicates that t2 and t3 are located in the subquery. After the subquery is promoted, t2 and t3 are joined, and then the join table is joined to t1. Where t2 is a foreign table, t3 is an internal table, t1 is a foreign table. The join table of t2 and t3 is an internal table.

Examples

Hint the guery plan in **Examples** as follows:

explain

 $select / *+ leading(((((store_sales store) promotion) item) customer) ad 2) store_returns) leading((store_store_sales)) */ i_product_name product_name ...$

First, **store_sales** and **store** are joined and **store_sales** is the inner table. Then, The result is joined with **promotion**, **item**, **customer**, **ad2**, and **store_returns** in sequence. The optimized plan is as follows:

WARNIN id	G: Duplicated or conflict hint: Leading(store_sales store), will be discarded. operation	E-rows	-	E-width	E-costs
1	->- Row Adapter	l 6	•	273	16308094.34
2	-> Vector Streaming (type: GATHER)	1 6	I	1 273	16308094.34
3	-> Vector Hash Aggregate	1 6	16MB	273	16308092.67
4	-> Vector Hash Join (5,20)	l 6	585MB	169	16308092.63
5	-> Vector Streaming(type: REDISTRIBUTE)	1320811	1MB	181	16069870.93
6	-> Vector Hash Join (7,19)	1320811	43MB	181	16061891.00
7	-> Vector Streaming(type: REDISTRIBUTE)	1320811	1MB	131	16056566.78
8	-> Vector Hash Join (9,18)	1320811	27MB	131	16048586.85
9	-> Vector Streaming(type: REDISTRIBUTE)	1383248	1MB	131	16038321.62
10	-> Vector Hash Join (11,17)	1383248	16MB	131	16029664.50
11	-> Vector Hash Join (12,16)	2626366951	16MB	1 73	15751384.88
12	-> Vector Hash Join (13,14)	2750085660	2156MB	1 77	14226077.19
13	-> CStore Scan on store	24048	1MB	l 19	2264.00
14	-> Vector Partition Iterator	2879987999	1MB	l 66	2756066.50
15	-> Partitioned CStore Scan on store_sales	2879987999	1MB	l 66	2756066.50
16	-> CStore Scan on promotion	36000	I 1MB	1 4	1268.50
17	-> CStore Scan on item	158	1MB	1 58	4051.25
18	-> CStore Scan on customer	12000000	1MB	1 8	12923.00
19	-> CStore Scan on customer_address ad2	6000000	1MB	1 58	5770.00
20	-> Vector Partition Iterator	287999764	1MB	1 12	227383.99
21	-> Partitioned CStore Scan on store returns	287999764	1MB	1 12	227383.99
(21 ro	wa)				

For details about the warning at the top of the plan, see **Hint Errors, Conflicts, and Other Warnings**.

13.4.9.3 Join Operation Hints

Function

Specifies the join method. It can be nested loop join, hash join, or merge join.

Syntax

[no] nestloop|hashjoin|mergejoin([@block_name] table_list)

Parameter Description

- no indicates that the specified hint will not be used for a join.
- block_name indicates the block name of the statement block. For details, see block_name.
- *table_list* specifies the tables to be joined. The values are the same as those of **join_table_list** but contain no parentheses.

For example:

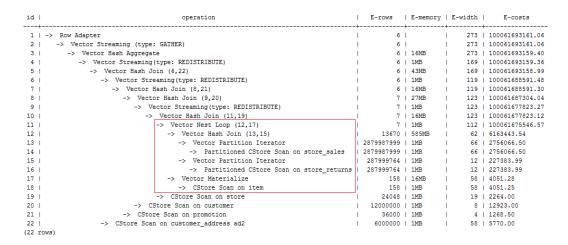
no nestloop(t1 t2 t3): nestloop is not used for joining t1, t2, and t3. The three tables may be joined in either of the two ways: Join t2 and t3, and then t1; join t1 and t2, and then t3. This hint takes effect only for the last join. If necessary, you can hint other joins. For example, you can add no nestloop(t2 t3) to join t2 and t3 first and to forbid the use of nestloop.

Examples

Hint the query plan in **Examples** as follows:

```
explain
select /*+ nestloop(store_sales store_returns item) */ i_product_name product_name ...
```

nestloop is used for the last join between **store_sales**, **store_returns**, and **item**. The optimized plan is as follows:



13.4.9.4 Rows Hints

Function

These hints specify the number of rows in an intermediate result set. Both absolute values and relative values are supported.

Syntax

rows([@block_name] table_list #|+|-|* const)

Parameter Description

- block_name indicates the block name of the statement block. For details, see block name.
- #,+,-, and * are operators used for hinting the estimation. # indicates that the original estimation is used without any calculation. +,-, and * indicate that the original estimation is calculated using these operators. The minimum calculation result is 1. table_list indicates the tables to be joined. The values are the same as those of table_list in Join Operation Hints.
- *const* can be any non-negative number and supports scientific notation.

For example:

rows(t1 #5): The result set of t1 is five rows.

rows(t1 t2 t3 *1000): Multiply the result set of joined t1, t2, and t3 by 1000.

Suggestion

- The hint using * for two tables is recommended, because this hint will take effect for a join as long as the two tables appear on both sides of this join. For example, if the hint is rows(t1 t2 * 3), the join result of (t1 t3 t4) and (t2 t5 t6) will be multiplied by 3 because t1 and t2 appear on both sides of the join.
- **rows** hints can be specified for the result sets of a single table, multiple tables, function tables, and subquery scan tables.

Examples

Hint the query plan in **Examples** as follows:

```
explain select /*+ rows(store_sales store_returns *50) */ i_product_name product_name ...
```

Multiply the result set of joined **store_sales** and **store_returns** by 50. The optimized plan is as follows:

id			_	E-width E-costs
1	-> Row Adapter	312		273 3401656.58
2	-> Vector Streaming (type: GATHER)	312		273 3401656.58
3	-> Vector Hash Aggregate	312	16MB	273 3401634.91
4	-> Vector Streaming(type: REDISTRIBUTE)	313	1MB	169 3401634.39
5	-> Vector Hash Join (6,21)	313	43MB	169 3401633.06
6	-> Vector Streaming(type: REDISTRIBUTE)	313	1MB	119 3397065.38
7	-> Vector Hash Join (8,20)	313	27MB	119 3397064.31
8	<pre>-> Vector Streaming(type: REDISTRIBUTE)</pre>	328	1MB	119 3387583.37
9		328	16MB	119 3387582.18
10	-> Vector Hash Join (11,18)	344	16MB	123 3386294.74
11	-> Vector Hash Join (12,14)	360	19MB	112 3384018.02
12	-> Vector Partition Iterator	287999764	1MB	12 227383.99
13	-> Partitioned CStore Scan on store_returns	287999764	1MB	12 227383.99
14	-> Vector Hash Join (15,17)	1516824	16MB	124 3065686.08
15	-> Vector Partition Iterator	2879987999	1MB	66 2756066.50
16	-> Partitioned CStore Scan on store_sales	2879987999	1MB	66 2756066.50
17	-> CStore Scan on item	158	1MB	58 4051.25
18	-> CStore Scan on store	24048	1MB	19 2264.00
19	-> CStore Scan on promotion	36000	1MB	4 1268.50
20	-> CStore Scan on customer	12000000	1MB	8 12923.00
21	-> CStore Scan on customer_address ad2	6000000	1MB	58 5770.00
(21 1	nows)			

The estimation value after the hint in row 11 is **360**, and the original value is rounded off to 7.

13.4.9.5 Stream Operation Hints

Function

Specifies the stream method, which can be broadcast, redistribute, or specifying the distribution key for **Agg** redistribution.

□ NOTE

Specifies the hint for the distribution column during the Agg process. This parameter is supported only by clusters of version 8.1.3.100 or later.

Syntax

[no] broadcast | redistribute([@block_name] table_list) | redistribute ([@block_name] (*) (columns))

Parameter Description

- **no** indicates that the hinted stream method is not used. When the hint is specified for the distribution columns in the **Agg** redistribution, **no** is invalid.
- block_name indicates the block name of the statement block. For details, see block_name.
- *table_list* specifies the tables to be joined. For details, see **Parameter Description**.
- When hints are specified for distribution columns, the asterisk (*) is fixed and the table name cannot be specified.
- **columns** specifies one or more columns in the **GROUP BY** clause. When there are no **GROUP BY** clauses, it can specify the columns in the **DISTINCT** clause.

- The specified distribution column must be specified using the column sequence number or column name in group by or distinct. The columns in count(distinct) can only be specified using column names.
- For a multi-layer query, you can specify the distribution column hint at each layer. The hint takes effect only at the corresponding layer.
- The column specified in count(distinct) takes effect only for two-level hashagg plans. Otherwise, the specified distribution column is invalid.
- If the optimizer finds that redistribution is not required after estimation, the specified distribution column is invalid.

Tips

- Generally, the optimizer selects a group of non-skew distribution keys for data redistribution based on statistics. If the default distribution keys have data skew, you can manually specify the distribution columns to avoid data skew.
- When selecting a distribution key, select a group of columns with high distinct values as the distribution key based on data distribution features. In this way, data can be evenly distributed to each DN after redistribution.
- After writing hints, you can run explain verbose to print the execution plan and check whether the specified distribution key is valid. If the specified distribution key is invalid, a warning is displayed.

Example

 Hint the query plan in Examples as follows: explain

select /*+ no redistribute(store_sales store_returns item store) leading(((store_sales store_returns item store) customer)) */ i_product_name product_name ...

In the original plan, the join result of **store_sales**, **store_returns**, **item**, and **store** is redistributed before it is joined with **customer**. After the hinting, the redistribution is disabled and the join order is retained. The optimized plan is as follows:

```
id |
                                                                                                        | E-memory | E-width |
                                          operation
                                                                                              E-rows
         Row Adapter
                                                                                                                          273 | 5718448.94
         -> Vector Streaming (type: GATHER)
                                                                                                                                5718448.94
            -> Vector Hash Aggregate
                                                                                                      6 | 16MB
                                                                                                                                5718447.27
                   Vector Streaming(type: REDISTRIBUTE)
                                                                                                      6 | 1MB
                                                                                                                           169 | 5718447.23
                   -> Vector Hash Join (6,21)
                                                                                                      6 | 16MB
                                                                                                                          169 | 5718446.86
                         Vector Hash Join (7,20)
                                                                                                                                5717159.60
                         -> Vector Streaming(type: REDISTRIBUTE)
                                                                                                        | 1MB
                                                                                                                          123 | 5712592.09
                               Vector Hash Join (9,18)
-> Vector Hash Join (10,17)
                                                                                                          585MB
                                                                                                                           123 | 5712591.93
                                                                                                                          123 | 3386294.56
                                                                                                        | 16MB
                                  -> Vector Hash Join (11,13)
-> Vector Partition Iterator
                                                                                                          1.9MB
                                                                                                                          112 | 3384018.02
                                                                                             287999764
12
                                        -> Partitioned CStore Scan on store_returns
                                                                                             287999764 | 1MB
                                                                                                                           12 | 227383.99
                                      -> Vector Hash Join (14,16)
                                                                                                1516824 | 16MB
                                                                                                                          124 | 3065686.08
                                        -> Vector Partition Iterator
                                                                                            2879987999 | 1MB
                                                                                                                           66 | 2756066.50
14
15
                                           -> Partitioned CStore Scan on store_sales
                                                                                            2879987999 | 1MB
                                                                                                                           66 | 2756066.50
                                         -> CStore Scan on item
                                                                                                                                4051.25
17
                                   -> CStore Scan on store
                                                                                                  24048 | 1MB
                                                                                                                           19 | 2264.00
                               -> Vector Streaming(type: BROADCAST)
                                                                                                                                2176297.36
                                                                                             288000000 | 1MB
                                                                                              12000000 | 1MB
                                                                                                                             8 | 12923.00
19 |
                                      CStore Scan on customer
                         -> CStore Scan on customer_address ad2
CStore Scan on promotion
                                                                                                6000000 I 1MB
                                                                                                                           58 | 5770.00
(21 rows)
```

Specifies the distribution columns for Agg redistribution.

explain (verbose on, costs off, nodes off) select /*+ redistribute ((*) $(2\ 3))$ */ a1, b1, c1, count(c1) from t1 group by a1, b1, c1 having count(c1) > 10 and sum(d1) > 100

In the following example, the last two columns of the specified **GROUP BY** columns are used as distribution keys.

```
QUERY PLAN
                   operation
  1 | -> Streaming (type: GATHER)
  2 |
       -> HashAggregate
  3 |
           -> Streaming(type: REDISTRIBUTE)
  4 |
              -> Seq Scan on public.tl
     Predicate Information (identified by plan id)
  2 --HashAggregate
       Filter: ((count(t1.c1) > 10) AND (sum(t1.d1) > 100))
Targetlist Information (identified by plan id)
  1 --Streaming (type: GATHER)
       Output: al, bl, cl, (count(cl))
  2 -- HashAggregate
       Output: al, bl, cl, count(cl)
       Group By Key: tl.al, tl.bl, tl.cl
  3 -- Streaming(type: REDISTRIBUTE)
       Output: al, bl, cl, dl
       Distribute Key: bl, cl
  4 -- Seq Scan on public.tl
        Output: al, bl, cl, dl
   ===== Query Summary =====
System available mem: 24862720KB
Query Max mem: 24862720KB
Query estimated mem: 3138KB
(30 rows)
```

• If the statement does not contain the **GROUP BY** clause, specify the distinct column as the distribution columns.

```
explain (verbose on, costs off, nodes off) select /*+ redistribute ((*) (3 1)) */ distinct a1, b1, c1 from t1;
```

```
QUERY PLAN
_____
 id |
                  operation
  1 | -> Streaming (type: GATHER)
  2 | -> HashAggregate
  3 | -> Streaming(type: REDISTRIBUTE)
4 | -> Seg Scan on public t1
              -> Seq Scan on public.tl
Targetlist Information (identified by plan id)
  1 --Streaming (type: GATHER)
       Output: al, bl, cl
  2 --HashAggregate
       Output: al, bl, cl
       Group By Key: tl.al, tl.bl, tl.cl
  3 -- Streaming (type: REDISTRIBUTE)
       Output: al, bl, cl
       Distribute Key: cl, al
  4 --Seq Scan on public.tl
       Output: al, bl, cl
   ===== Query Summary =====
System available mem: 24862720KB
Query Max mem: 24862720KB
Query estimated mem: 3136KB
(25 rows)
```

13.4.9.6 Scan Operation Hints

Function

These hints specify a scan operation, which can be **tablescan**, **indexscan**, or **indexonlyscan**.

Syntax

[no] tablescan|indexscan|indexonlyscan([@block_name] table [index])

Parameter Description

- **no** indicates that the specified hint will not be used for a join.
- block_name indicates the block name of the statement block. For details, see block_name.
- *table* specifies the table to be scanned. You can specify only one table. Use a table alias (if any) instead of a table name.

□ NOTE

 The syntax format of the table is as follows: [schema.]table[@block_name]

The table name can contain the schema name or block name before the subquery statement block is promoted. If the subquery statement block is optimized and rewritten by the optimizer, the value of **block_name** is different from that of **block_name** in leading.

- If a table has an alias, the alias is preferentially used to represent the table.
- *index* indicates the index for **indexscan** or **indexonlyscan**. You can specify only one index.

indexscan and **indexonlyscan** hints can be used only when the specified index belongs to the table.

Scan operation hints can be used for row-store tables, column-store tables, HDFS tables, HDFS foreign tables, OBS tables, and subquery tables. HDFS tables include primary tables and delta tables. The delta tables are invisible to users. Therefore, scan operation hints are used only for primary tables.

If **indexscan** is specified, **indexscan** or **indexonlyscan** takes effect. **indexscan** and **indexonlyscan** can also take effect at the same time. When **indexscan** and **indexonlyscan hints** appear at the same time, **indexonlyscan** takes effect first.

Example

To specify an index-based hint for a scan, create an index named **i** on the **i_item_sk** column of the **item** table.

```
create index i on item(i_item_sk);
```

Hint the query plan in **Examples** as follows:

```
explain select /*+ indexscan(item i) */ i_product_name product_name ...
```

item is scanned based on an index. The optimized plan is as follows:

id	operation	E-rows	E-memory	E-width	E-costs
1	-> Row Adapter	6		273	100061674938.26
2	-> Vector Streaming (type: GATHER)	6		273	100061674938.26
3	-> Vector Hash Aggregate	6	16MB	273	100061674936.59
4	-> Vector Streaming(type: REDISTRIBUTE)	6	1MB	169	100061674936.55
5	-> Vector Hash Join (6,21)	6	43MB	169	100061674936.19
6	-> Vector Streaming(type: REDISTRIBUTE)	6	1MB	119	100061670368.67
7	-> Vector Hash Join (8,20)	6	16MB	119	100061670368.50
8	-> Vector Hash Join (9,19)	7	27MB	123	100061669081.23
9	-> Vector Streaming(type: REDISTRIBUTE)	7	1MB	123	100061659600.47
10	-> Vector Hash Join (11,18)	7	16MB	123	100061659600.31
11	-> Vector Nest Loop (12,17)	7	1MB	112	100061657323.77
12	-> Vector Hash Join (13,15)	13670	585MB	62	6163443.54
13	-> Vector Partition Iterator	2879987999	1MB	66	2756066.50
14	-> Partitioned CStore Scan on store_sales	2879987999	1MB	66	2756066.50
15	-> Vector Partition Iterator	287999764	1MB	12	227383.99
16	-> Partitioned CStore Scan on store_returns	287999764	1MB	12	227383.99
17	-> CStore Index Scan using i on item	1	1MB	58	4.01
18	-> CStore Scan on store	24048	1MB	19	2264.00
19	-> CStore Scan on customer	12000000	1MB	8	12923.00
20	-> CStore Scan on promotion	36000	1MB	4	1268.50
21	-> CStore Scan on customer_address ad2	6000000	1MB	58	5770.00
(21 rd	ws)				

13.4.9.7 Sublink Name Hints

Function

These hints specify the name of a sublink block.

Syntax

blockname ([@block_name] table)

Precautions

- This block name hint is used by an outer query only when a sublink is pulled up. Currently, only the Agg equivalent join, IN, and EXISTS sublinks can be pulled up. This hint is usually used together with the hints described in the previous sections.
- The subquery after the **FROM** keyword is hinted by using the subquery alias. In this case, **block_name hint** becomes invalid.
- If a sublink contains multiple tables, the tables will be joined with the outerquery tables in a random sequence after the sublink is pulled up. In this case, **blockname** also becomes invalid.

Parameter Description

- block_name indicates the block name of the statement block. For details, see block name.
- *table* indicates the name you have specified for a sublink block.

□ NOTE

 The syntax format of the table is as follows: [schema.]table[@block name]

The table name can contain the schema name or block name before the subquery statement block is promoted. If the subquery statement block is optimized and rewritten by the optimizer, the value of **block_name** is different from that of **block_name** in leading.

• If a table has an alias, the alias is preferentially used to represent the table.

Example

explain select /*+nestloop(store_sales tt) */ * from store_sales where ss_item_sk in (select /* +blockname(tt)*/ i_item_sk from item group by 1);

tt indicates the sublink block name. After being pulled up, the sublink is joined with the outer-query table **store_sales** by using **nestloop**. The optimized plan is as follows:

id	operation	1	E-rows	I	E-memory	l	E-width	I	E-costs
						+-			
	-> Row Adapter		1439994000			ı			325105765847.91
2	-> Vector Streaming (type: GATHER)	ı	1439994000	ı		l	216	I	325105765847.91
3	-> Vector Nest Loop Semi Join (4, 6)	1	1439994000	1	1MB	L	216	I	325026664615.00
4	-> Vector Partition Iterator	1	2879987999	1	1MB	L	216	I	2756066.50
5	-> Partitioned CStore Scan on store_sales	1	2879987999	1	1MB	L	216	I	2756066.50
6	-> Vector Materialize	1	300000	1	16MB	L	4	I	4176.25
7	-> Vector Hash Aggregate	1	300000	1	16MB	L	4	I	3988.75
8	-> CStore Scan on item	1	300000	1	1MB	L	4	I	3832.50
(8 ro	ws)								

13.4.9.8 Skew Hints

Function

Theses hints specify redistribution keys containing skew data and skew values, and are used to optimize redistribution involving Join or HashAgg.

Precautions

- Skew hints are used only if redistribution is required and the specified skew information matches the redistribution information.
- Skew hints are controlled by the GUC parameter **skew_option**. If the parameter is disabled, skew hints cannot be used for solving skew.
- Currently, skew hints support only the table relationships of the ordinary table
 and subquery types. Hints can be specified for base tables, subqueries, and
 WITH ... AS clauses. Unlike other hints, a subquery can be used in skew hints
 regardless of whether it is pulled up.
- Use an alias (if any) to specify a table where data skew occurs.
- You can use a name or an alias to specify a skew column as long as it is not ambiguous. The columns in skew hints cannot be expressions. If data skew occurs in the redistribution that uses an expression as a redistribution key, set the redistribution key as a new column and specify the column in skew hints.
- The number of skew values must be an integer multiple of the number of columns. Skew values must be grouped based on the column sequence, with each group containing a maximum of 10 values. You can specify duplicate values to group skew columns having different number of skew values. For example, the c1 and c2 columns of the t1 table contains skew data. The skew value of the c1 column is a1, and the skew values of the c2 column are b1 and b2. In this case, the skew hint is skew(t1 (c1 c2) ((a1 b1)(a1 b2))). (a1 b1) is a value group, where NULL is allowed as a skew value. Each hint can contain a maximum of 10 groups and the number of groups should be an integer multiple of the number of columns.
- In the redistribution optimization of Join, a skew value must be specified for skew hints. The skew value can be left empty for HashAgg.
- If multiple tables, columns, or values are specified, separate items of the same type with spaces.
- The type of skew values cannot be forcibly converted in hints. To specify a string, enclose it with single quotation marks (' ').

Syntax

- Specify single-table skew.
 skew([@block_name] table (column) [(value)])
- Specify intermediate result skew. skew([@block_name] (join_rel) (column) [(value)])

Parameter Description

- block_name indicates the block name of the statement block. For details, see block name.
- table specifies the table where skew occurs.

□ NOTE

 The syntax format of the table is as follows: [schema.]table[@block_name]

The table name can contain the schema name or block name before the subquery statement block is promoted. If the subquery statement block is optimized and rewritten by the optimizer, the value of **block_name** is different from that of **block_name** in leading.

- If a table has an alias, the alias is preferentially used to represent the table.
- **join_rel** specifies two or more joined tables. For example, **(t1 t2)** indicates that the result of joining **t1** and **t2** tables contains skew data.
- **column** specifies one or more columns where skew occurs.
- value specifies one or more skew values.

Example:

Specify single-table skew.

Each skew hint describes the skew information of one table relationship. To describe the skews of multiple table relationships in a query, specify multiple skew hints.

Skew hints have the following formats:

- One skew value in one column: skew(t (c1) (v1))
 Description: The v1 value in the c1 column of the t table relationship causes skew in query execution.
- Multiple skew values in one column: skew(t (c1) (v1 v2 v3 ...))
 Description: Values including v1, v2, and v3 in the c1 column of the t table relationship cause skew in query execution.
- Multiple columns, each having one skew value: skew(t (c1 c2) (v1 v2))
 Description: The v1 value in the c1 column and the v2 value in the c2 column of the t table relationship cause skew in query execution.
- Multiple columns, each having multiple skew values: skew(t (c1 c2) ((v1 v2) (v3 v4) (v5 v6) ...))

Description: Values including v1, v3, and v5 in the c1 column and values including v2, v4, and v6 in the c2 column of the t table relationship cause skew in query execution.

NOTICE

In the last format, parentheses for skew value groups can be omitted, for example, **skew(t (c1 c2) (v1 v2 v3 v4 v5 v6 ...))**. In a skew hint, either use parentheses for all skew value groups or for none of them.

Otherwise, a syntax error will be generated. For example, **skew(t (c1 c2) (v1 v2 v3 v4 (v5 v6) ...))** will generate an error.

• Specify intermediate result skew.

If data skew does not occur in base tables but in an intermediate result during query execution, specify skew hints of the intermediate result to solve the skew. The format is **skew((t1 t2) (c1) (v1))**.

Description: Data skew occurs after the table relationships **t1** and **t2** are joined. The **c1** column of the **t1** table contains skew data and its skew value is **v1**.

c1 can exist only in a table relationship of **join_rel**. If there is another column having the same name, use aliases to avoid ambiguity.

Suggestion

- For a multi-level query, write the hint on the layer where data skew occurs.
- For a listed subquery, you can specify the subquery name in a hint. If you know data skew occurs on which base table, directly specify the table.
- Aliases are preferred when you specify a table or column in a hint.

Examples

Specify single-table skew.

• Specify hints in the original query.

For example, the original query is as follows:

```
explain
with customer_total_return as
(select sr_customer_sk as ctr_customer_sk
,sr_store_sk as ctr_store_sk
,sum(SR_FEE) as ctr_total_return
from store_returns
,date_dim
where sr returned date sk = d date sk
and d_year =2000
group by sr_customer_sk
,sr_store_sk)
select c_customer_id
from customer_total_return ctr1
,customer
where ctr1.ctr_total_return > (select avg(ctr_total_return)*1.2
from customer_total_return ctr2
where ctr1.ctr_store_sk = ctr2.ctr_store_sk)
and s_store_sk = ctr1.ctr_store_sk
and s state = 'NM'
and ctr1.ctr_customer_sk = c_customer_sk
order by c_customer_id
limit 100;
```

```
| E-width | E-costs
Row Adapter
-> Vector Limit
-> Vector Streaming (type: GATHER)
-> Vector Limit
-> Vector Sort
-> Vector Hash Join (7,29)
-> Vector Streaming (typ
                                                                                                                                                                                                                                                                                                                                                                                                                                                                   100 |
100 |
2400 |
2400 |
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  20 | 911254.47
20 | 911254.47
20 | 911325.75
20 | 911247.62
                                                                                                                                                                                                                                                                                                                                                                                                                                                 2400 | 1MB
3684816 | 16MB
3684817 | 41MB(12374MB)
3684817 | 394KB
3684817 | 16MB
11054450 | 16MB
50247501 | 397MB(12671MB)
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           911631.21
                                                                                        ector Hash Join (7,29)

Vector Streaming(type: REDISTRIBUTE)

-> Vector Hash Join (9,19)

-> Vector Hash Join (10,18)

-> Vector Hash Augregate

-> Vector Hash Join (13,15)

-> Vector Hash Join (13,15)

-> Vector Fartition Iterator

-> Partitioned CStore Soan on store_returns

-> Vector Streaming(type: REDISTRIBUTE)

-> Vector Partition Iterator

-> Pertitioned CStore Soan on date_dim

-> CStore Soan on store

-> CStore Soan on store

-> Vector Hash Aggregate
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                20 | 905379.41
4 | 883010.31
4 | 861302.05
44 | 427109.71
54 | 395302.57
22 | 358663.76
22 | 294300.51
6 | 227383.99
4 | 975.56
4 | 910.65
4 | 910.65
                                                                                                                                                                                                                                                                                                                                                                                                                                                 50247501
                                                                                                                                                                                                                                                                                                                                                                                                                                                 50247501
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     16MB
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   1MB
1MB
384KB
1MB
1MB
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  4 | 910.65

4 | 1006.39

68 | 426707.38

36 | 416239.03

54 | 395302.52

22 | 358663.76

22 | 294300.51

26 | 227383.99

4 | 975.56

4 | 910.65

4 | 910.65

24 | 12923.00
                                                                                                                    -> CStore Scan on store

Vector Hash Aggregate | 192

-> Vector Studymery Scan on ctr2 | 50247501

-> Vector Streaming(type: REDISTRIBUTE) | 50247501

-> Vector Streaming(type: REDISTRIBUTE) | 50247501

-> Vector Hash Join (24,26) | 50247501

-> Vector Partition Iterator | 287999764

-> Vector Streaming(type: BROADCAST) | 8712

-> Vector Partitioned CStore Scan on store returns | 363

-> Partitioned CStore Scan on date_dim | 363

-> Partitioned CStore Scan on date_dim | 363

-> Scan on customer | 12000000
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     16MB
1MB
397MB (12671MB)
384KB
16MB
                                                                        -> CStore Scan on customer
```

Specify the hints of HashAgg in the inner **with** clause and of the outer Hash Join. The query containing hints is as follows:

```
explain
with customer_total_return as
(select /*+ skew(store_returns(sr_store_sk sr_customer_sk)) */sr_customer_sk as ctr_customer_sk
,sr_store_sk as ctr_store_sk
,sum(SR_FEE) as ctr_total_return
from store returns
,date_dim
where sr_returned_date_sk = d_date_sk
and d_year =2000
group by sr_customer_sk
,sr_store_sk)
select /*+ skew(ctr1(ctr_customer_sk)(11))*/ c_customer_id
from customer total return ctr1
,store
,customer
where ctr1.ctr_total_return > (select avg(ctr_total_return)*1.2
from customer_total_return ctr2
where ctr1.ctr_store_sk = ctr2.ctr_store_sk)
and s_store_sk = ctr1.ctr_store_sk
and s state = 'NM'
and ctr1.ctr_customer_sk = c_customer_sk
order by c_customer_id
limit 100:
```

The hints indicate that the **group by** in the inner **with** clause contains skew data during redistribution by HashAgg, corresponding to the original Hash Agg operators 10 and 21; and that the **ctr_customer_sk** column in the outer **ctr1** table contains skew data during redistribution by Hash Join, corresponding to operator 6 in the original plan. The optimized plan is as follows:

```
| 1 -> Row Adapter | 100 | 20 | 1061778.14 | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit | 2 | -> Vector Limit |
```

To solve data skew in the redistribution, Hash Agg is changed to double-level Agg operators and the redistribution operators used by Hash Join are changed in the optimized plan.

Modify the query and then specify hints.

For example, the original query and its plan are as follows:

explain select count(*) from store_sales_1 group by round(ss_list_price);

Columns in hints do not support expressions. To specify hints, rewrite the query as several subqueries. The rewritten query and its plan are as follows:

```
explain select count(*)
```

```
from (select round(ss_list_price),ss_hdemo_sk
from store_sales_1)tmp(a,ss_hdemo_sk)
group by a;
```

Ensure that the service logic is not changed during the rewriting.

Specify hints in the rewritten query as follows:

```
explain
select /*+ skew(tmp(a)) */ count(*)
from (select round(ss_list_price),ss_hdemo_sk
from store_sales_1)tmp(a,ss_hdemo_sk)
group by a;
```

The plan shows that after Hash Agg is changed to double-layer Agg operators, redistributed data is greatly reduced and redistribution time shortened.

You can specify hints in columns in a subquery, for example:

```
explain
select /*+ skew(tmp(b)) */ count(*)
from (select round(ss_list_price) b,ss_hdemo_sk
from store_sales_1)tmp(a,ss_hdemo_sk)
group by a;
```

13.4.9.9 Hint That Disables Subquery Pull-up

Function

To optimize query logic, the optimizer usually pulls up subqueries for execution. However, sometimes the pulled up subqueries do not run much faster than others, and may even be slower due to enlarged search scope. In this case, you can specify the **no merge** hint to disable pull-up. This hint is not recommended in most cases.

Syntax

```
no merge[@block_name]
no merge ([@block_name1] subquery_name[@block_name2])
```

Description

- block_name indicates the block name of the statement block. For details, see block_name.
- subquery_name indicates the name of a subquery. It can also be a view or CTE name. The specified subquery will not be unnested during logic optimization. If subquery_name is not specified, the current query will not be unnested.

Example

Create tables t1, t2, and t3.

```
create table t1(a1 int,b1 int,c1 int,d1 int);
create table t2(a2 int,b2 int,c2 int,d2 int);
create table t3(a3 int,b3 int,c3 int,d3 int);
```

The original statement is as follows:

explain select * from t3, (select a1,b2,c1,d2 from t1,t2 where t1.a1=t2.a2) s1 where t3.b3=s1.b2;

id	operation				E-width	
	-> Hash Join (2,6)	ī	44450			754.31
2	-> Hash Join (3,4)	L	8885	Ī	28	182.11
3	-> Seq Scan on t3	L	1776	Ī	16	27.76
4	-> Hash	L	1776	I	12	27.76
5	-> Seq Scan on t2	L	1776	I	12	27.76
6	-> Hash	L	1776	I	8	27.76
7	-> Seq Scan on tl	L	1776	Ī	8	27.76

In this query, you can use the following methods to disable the pull-up of subquery **s1**:

- Method 1: explain select /*+ no merge(s1) */ * from t3, (select a1,b2,c1,d2 from t1,t2 where t1.a1=t2.a2) s1 where t3.b3=s1.b2;
- Method 2: explain select * from t3, (select /*+ no merge */ a1,b2,c1,d2 from t1,t2 where t1.a1=t2.a2) s1 where t3.b3=s1.b2;

Outcome:

id	•			E-width	•
		1	8880	•	+ 443.03
2	-> Hash Join (3,4)	Ĺ	8885	16	182.11
3	-> Seq Scan on tl	Ĺ	1776	J 8	27.76
4	-> Hash	Ĺ	1776	12	27.76
5	-> Seq Scan on t2	Ĺ	1776	12	27.76
6	-> Hash	Ĺ	1776	16	27.76
7	-> Seq Scan on t3	ī	1776	16	27.76

13.4.9.10 Drive Hints

Function

When generating a query plan, the optimizer uses the dynamic planning or genetic algorithm to enumerate possible join paths and selects the optimal path. When the number of join tables increases, the search space may expand greatly. In this case, you can use a drive hint to specify the fact table in the query and use heuristic search to narrow the search range. This hint takes effect only when the GUC parameter join_search_mode is set to heuristic.

Syntax

[no] drive (table)

Description

- no indicates that the table specified by the hint is not a drive table.
- *table* specifies the table specified by the hint. You can specify only one table. Use a table alias (if any) instead of a table name.

Example

Create tables t1, t2, and t3.

```
create table t1(a1 int,b1 int,c1 int,d1 int);
create table t2(a2 int,b2 int,c2 int,d2 int);
create table t3(a3 int,b3 int,c3 int,d3 int);
```

The original statement is as follows:

explain select * from t1,t2,t3 where t1.b1=t2.b2 and t1.c1=t3.c3 and t2.d2=t3.d3;

	QUERY PLAN								
id	operation				E-memory			-	
1	-> Streaming (type: GATHER)	I	4			i			40.03
2	-> Hash Join (3,9)	I .	4	Ī	1MB	I .	48	1	34.03
3	<pre>-> Streaming(type: BROADCAST)</pre>	I	40	1	2MB	I	32	1	23.64
4	-> Hash Join (5,7)	I	20	1	1MB	I .	32	1	22.06
5	-> Streaming(type: BROADCAST)	I	40	1	2MB	I	16	1	11.67
6	-> Seq Scan on tl	I	20	1	1MB	I	16	1	10.10
7	-> Hash	I	20	1	16MB	I	16	1	10.10
8	-> Seq Scan on t2	I	20	1	1MB	I .	16	1	10.10
9	-> Hash	I	20	1	16MB	I	16	1	10.10
10	-> Seq Scan on t3	I	20	1	1MB	I	16	1	10.10

In the preceding query, you can specify table **t3** as the drive table so that **t3** is preferentially joined with other tables. This reduces the search paths generated by the plan and changes the join sequence. :

explain select /*+ drive(t3) */ * from t1, t2, t3 where t1.b1=t2.b2 and t1.c1=t3.c3 and t2.d2=t3.d3;

After hints are used:

	QUERY PLA	N								
id	operation	- 1	E-rows	ı	E-memory	L	E-width	E-costs		
	+	-+		+		+-	+			
1	-> Streaming (type: GATHER)	-	4			L	48	40.03		
2	-> Hash Join (3,9)	-1	4	1	1MB	L	48	34.03		
3	<pre>-> Streaming(type: BROADCAST)</pre>		40		2MB	L	32	23.64		
4	-> Hash Join (5,7)	1	20	ı	1MB	r	32	22.06		
5	-> Streaming(type: BROADCAST)	1	40	ī	2MB	r	16	11.67		
6	-> Seq Scan on t3(DRIVE)	ī	20	ī	1MB	ï	16	10.10		
7	-> Hash	1	20	ī	16MB	ř.	16	10.10		
8	-> Seq Scan on tl	ī	20	ī	1MB	ï	16	10.10		
9	-> Hash	ī	20	ī	16MB	ï	16	10.10		
10	-> Seq Scan on t2	i	20	i	1MB	i	16	10.10		

A CAUTION

If the GUC parameter <code>join_search_mode</code> is set to <code>heuristic</code> and the number of joined tables exceeds the value of <code>from_collapse_limit</code>, the optimizer automatically identifies drive tables and displays them in the plan. If the final drive selected by the optimizer is incorrect, you can use <code>no drive hint</code> to correct the selection.

13.4.9.11 Dictionary Code Hint

Function

Specifies a column to construct a dictionary code. The comparison of character strings in the dictionary code is converted into the comparison of numbers, which accelerates the query speed such as **group by** and **filter**. This hint is supported only by clusters of version 8.3.0 or later.

Precautions

• Currently, only the new version hatore tables are supported (the table-level parameter **enable hatore opt** is set to **on**).

Syntax

/* + (no) dict(table (column)) */

Parameter description

- dict(table (column))
 Column with the dictionary encoding table enabled.
- no dict(table (column))
 Column with the dictionary encoding table disabled.

Example

SELECT /*+ dict (bitmaptbl_high (server_ip)) */ distinct(server_ip) FROM bitmaptbl_high WHERE scope_name='saetataetaeta' ORDER BY server_ip;

The generated plan is as follows. server_ip uses the dictionary encoding:

You can use no **dict to** disable **server_ip** from using dictionary encoding.

SELECT /*+ no dict (bitmaptbl_high (server_ip)) */ distinct(server_ip) FROM bitmaptbl_high WHERE scope_name='saetataetaeta' ORDER BY server_ip;

13.4.9.12 Configuration Parameter Hints

Function

A hint, or a GUC hint, specifies a configuration parameter value when a plan is generated.

Precautions

- If a parameter set by hint takes effect at the statement level, the hint must be written to the top-level query instead of the subquery. For UNION, INTERSECT, EXCEPT, and MINUS statements, you can write the GUC hint at the statement level to any SELECT clause that participates in the set operation. The configuration parameters set by the GUC hint take effect on each SELECT clause that participates in the set operation.
- When a subquery is pulled up, all GUC hints on the subquery are discarded.
- If a parameter is set by both the statement-level GUC hint and the subquery-level GUC hint, the subquery-level GUC hint takes effect in the corresponding subquery, and the statement-level GUC hint takes effect in other subqueries of the statement.

Syntax

set [global]([@block_name] guc_name guc_value)

Parameters

- **global** indicates that the parameter set by hint takes effect at the statement level. If **global** is not specified, the parameter takes effect only in the subquery where the hint is located.
- block_name indicates the block name of the statement block. For details, see block_name.
- **guc_name** indicates the name of the configuration parameter specified by hint.
- **guc_value** indicates the value of a configuration parameter specified by hint.

Currently, GUC hints support only some configuration parameters. Some parameters cannot be configured at the subquery level and can only be configured at the statement level. The following table lists the supported parameters.

Table 13-18 Configuration parameters supported by GUC hints

Parameter	Configured at the Subquery Level (Yes/No)
agg_max_mem	Yes
agg_redistribute_enhancement	Yes

Parameter	Configured at the Subquery Level (Yes/No)
best_agg_plan	Yes
cost_model_version	No
cost_param	No
enable_array_optimization	No
enable_bitmapscan	Yes
enable_broadcast	Yes
enable_csqual_pushdown	No
enable_redistribute	Yes
enable_extrapolation_stats	Yes
enable_fast_query_shipping	No
enable_force_vector_engine	No
enable_hashagg	Yes
enable_hashfilter	No
enable_hashjoin	Yes
enable_index_nestloop	Yes
enable_indexonlyscan	Yes
enable_indexscan	Yes
enable_join_pseudoconst	Yes
enable_mergejoin	Yes
enable_mixedagg	No
enable_nestloop	Yes
enable_nodegroup_debug	No
enable_partition_dynamic_pruning	Yes
enable_seqscan	Yes
enable_sonic_hashagg	No
enable_sonic_hashjoin	Yes
enable_sort	Yes
enable_stream_ctescan	No
enable_tidscan	Yes

Parameter	Configured at the Subquery Level (Yes/No)
enable_value_redistribute	Yes
enable_vector_engine	No
expected_computing_nodegroup	No
force_bitmapand	Yes
from_collapse_limit	Yes
join_collapse_limit	Yes
join_num_distinct	Yes
outer_join_max_rows_multipler	Yes
prefer_hashjoin_path	No
qrw_inlist2join_optmode	Yes
qual_num_distinct	Yes
query_dop	No
query_max_mem	No
query_mem	No
rewrite_rule	No
setop_optmode	Yes
skew_option	Yes
stream_ctescan_max_estimate_mem	No
stream_ctescan_pred_threshold	No
stream_ctescan_refcount_threshold	No
windowagg_pushdown_enhancement	No
index_selectivity_cost	Yes
index_cost_limit	Yes

Examples

Hint the query plan in **Examples** as follows:

```
explain
select /*+ set global(query_dop 0) */ i_product_name product_name
...
```

This hint indicates that the **query_dop** parameter is set to **0** when the plan for a statement is generated, which means the SMP adaptation function is enabled. The generated plan is as follows:

id	operation			E-memory		
1	-> Row Adapter	i	1			19595.89
2	-> Vector Sonic Hash Aggregate	l l	1		230	19595.89
3	-> Vector Streaming (type: GATHER)	l l	3		230	19595.89
4	-> Vector Sonic Hash Aggregate	I	3	16MB	230	19595.66
5	-> Vector Nest Loop (6,28)	l .	3	1MB	126	19595.62
6	-> Vector Nest Loop (7,27)	l .	3	1MB	130	19291.57
7	-> Vector Streaming(type: LOCAL GATHER dop: 1/2)	l .	3	4MB	118	19279.41
8	-> Vector Nest Loop (9,24)	1	3	1MB	118	19279.38
9	-> Vector Streaming(type: SPLIT REDISTRIBUTE dop: 2/2)	I .	3	4MB	82	18117.66
10	-> Vector Nest Loop (11,21)	l .	3	1MB	82	18117.61
11	-> Vector Streaming(type: SPLIT REDISTRIBUTE dop: 2/2)	l .	3	4MB	82	16195.20
12	-> Vector Sonic Hash Join (13,15)	1	3	16MB	82	16195.15
13	-> Vector Partition Iterator	28	7514	1MB	12	1110.42
14	-> Partitioned CStore Scan on store_returns	28	7514	1MB	12	1110.42
15	-> Vector Streaming(type: LOCAL BROADCAST dop: 2/2)	l .	2764	4MB	94	14718.42
16	-> Vector Sonic Hash Join (17,19)	l .	1382	16MB	94	14699.69
17	-> Vector Partition Iterator	288	0404	1MB	39	11541.07
18	-> Partitioned CStore Scan on store_sales	288	0404	1MB	39	11541.07
19	-> Vector Streaming(type: LOCAL BROADCAST dop: 2/2)	l .	16	4MB	55	1947.12
20	-> CStore Scan on item	l .	8	1MB	55	1947.00
21	-> Vector Materialize	10	0000	16MB	8	1797.41
22	-> Vector Streaming(type: LOCAL REDISTRIBUTE dop: 2/2)	10	0000	4MB	8	1714.07
23	-> CStore Scan on customer	10	0000	1MB	8	703.67
24	-> Vector Materialize	5	0000 I	16MB	44	1099.22
25	-> Vector Streaming(type: LOCAL REDISTRIBUTE dop: 2/2)	5	0000	4MB	44	1057.55
26	-> CStore Scan on customer_address ad2	5	0000	1MB	44	552.33
27	-> CStore Scan on store	l .	36	1MB	20	12.01
28	-> CStore Scan on promotion	l .	900 I	1MB	1 4	300.30
(28 :	cows)					

13.4.9.13 Hint Errors, Conflicts, and Other Warnings

Plan hints change an execution plan. You can run **EXPLAIN** to view the changes.

Hints containing errors are invalid and do not affect statement execution. The errors will be displayed in different ways based on statement types. Hint errors in an **EXPLAIN** statement are displayed as a warning on the interface. Hint errors in other statements will be recorded in debug1-level logs containing the **PLANHINT** keyword.

Hint Error Types

Syntax errors.

An error will be reported if the syntax tree fails to be reduced. The No. of the row generating an error is displayed in the error details.

For example, the hint keyword is incorrect, no table or only one table is specified in the **leading** or **join** hint, or no tables are specified in other hints. The parsing of a hint is terminated immediately after a syntax error is detected. Only the hints that have been parsed successfully are valid.

For example:

leading((t1 t2)) nestloop(t1) rows(t1 t2 #10)

The syntax of **nestloop(t1)** is wrong and its parsing is terminated. Only **leading(t1 t2)** that has been successfully parsed before **nestloop(t1)** is valid.

- Semantic errors.
 - An error will be reported if the specified tables do not exist, multiple tables are found based on the hint setting, or a table is used more than once in the **leading** or **join** hint.
 - An error will be reported if the index specified in a scan hint does not exist.
 - If multiple tables with the same name exist after a subquery is pulled up and some of them need to be hinted, add aliases for them to avoid name duplication.
- Duplicated or conflicted hints.

If hint duplication or conflicts occur, only the first hint takes effect. A message will be displayed to describe the situation.

- Hint duplication indicates that a hint is used more than once in the same query, for example, nestloop(t1 t2) nestloop(t1 t2).
- A hint conflict indicates that the functions of two hints with the same table list conflict with each other.

For example, if **nestloop (t1 t2) hashjoin (t1 t2)** is used, **hashjoin (t1 t2)** becomes invalid. **nestloop(t1 t2)** does not conflict with **no mergejoin(t1 t2)**.

NOTICE

The table list in the **leading** hint is disassembled. For example, **leading** (t1 t2 t3) will be disassembled as **leading(t1 t2) leading((t1 t2) t3)**, which will conflict with **leading(t2 t1)** (if any). In this case, the latter **leading(t2 t1)** becomes invalid. If two hints use duplicated table lists and only one of them has the specified outer/inner table, the one without a specified outer/inner table becomes invalid.

- A hint becomes invalid after a sublink is pulled up.
 - In this case, a message will be displayed. Generally, such invalidation occurs if a sublink contains multiple tables to be joined, because the table list in the sublink becomes invalid after the sublink is pulled up.
- Unsupported column types.
 - Skew hints are specified to optimize redistribution. They will be invalid if their corresponding columns do not support redistribution.
- Specified hints are not used.
 - If hashjoin or mergejoin is specified for non-equivalent joins, it will not be used.
 - If indexscan or indexonlyscan is specified for a table that does not have an index, it will not be used.
 - If indexscan hint or indexonlyscan is specified for a full-table scan or for a scan whose filtering conditions are not set on index columns, it will not be used.
 - The specified **indexonlyscan** hint is used only when the output column contains only indexes.
 - In equivalent joins, only the joins containing equivalence conditions are valid. Therefore, the leading, join, and rows hints specified for the joins without an equivalence condition will not be used. For example, t1, t2, and t3 are to be joined, and the join between t1 and t3 does not contain an equivalence condition. In this case, leading(t1 t3) will not be used.
 - To generate a streaming plan, if the distribution key of a table is the same as its join key, redistribute specified for this table will not be used. If the distribution key and join key are different for this table but the same for the other table in the join, redistribute specified for this table will be used but broadcast will not.
 - If a hint for an **Agg** distribution column is not used, the possible causes are as follows:

- The specified distribution key contains data types that do not support redistribution.
- Redistribution is not required in the execution plan.
- Wrong distribution key sequence numbers are executed.
- For AP functions that use the GROUPING SETS and CUBE clauses, hints are not supported for distribution keys in window aggregate functions.

□ NOTE

Specifies the hint for the distribution column druing the Agg process.. This parameter is supported only by clusters of version 8.1.3.100 or later.

- If no sublink is pulled up, the specified blockname hint will not be used.
- For unused skew hints, the possible causes are:
 - The plan does not require redistribution.
 - The columns specified by hints contain distribution keys.
 - Skew information specified in hints is incorrect or incomplete, for example, no value is specified for join optimization.
 - Skew optimization is disabled by GUC parameters.
- For unused guc hints, the possible causes are:
 - The configuration parameter does not exist.
 - The configuration parameter is not supported by GUC hints.
 - The configuration parameter value is invalid.
 - The statement-level GUC hint is not written in the top-level guery.
 - The configuration parameter set by the GUC hint at the subquery level cannot be set at the subquery level.
 - The subquery where the GUC hint is located is pulled up.

13.4.9.14 Plan Hint Cases

This section takes the statements in TPC-DS (Q24) as an example to describe how to optimize an execution plan by using hints in 1000X+24DN environments. For example:

```
select avg(netpaid) from
(select c_last_name
,c_first_name
,s_store_name
,ca_state
,s_state
,i_color
,i_current_price
,i_manager_id
,i_units
,i_size
```

```
,sum(ss_sales_price) netpaid
from store_sales
,store_returns
.store
,item
,customer
,customer_address
where ss_ticket_number = sr_ticket_number
and ss_item_sk = sr_item_sk
and ss_customer_sk = c_customer_sk
and ss_item_sk = i_item_sk
and ss_store_sk = s_store_sk
and c_birth_country = upper(ca_country)
and s_zip = ca_zip
and s_market_id=7
group by c_last_name
,c_first_name
,s_store_name
,ca_state
.s state
i_color,
i current price,
,i_manager_id
i_units,
,i_size);
```

 The original plan of this statement is as follows and the statement execution takes 110s:

Figure 13-10 Statement initial plan



In this plan, the performance of the layer-10 **broadcast** is poor because the number of rows estimated by the layer-11 operator is 2,140, which is much lower than the actual number of rows. The inaccurate estimation is mainly caused by the underestimated number of rows in layer-13 hash join. In this layer, **store_sales** and **store_returns** are joined (based on the **ss_ticket_number** and **ss_item_sk** columns in **store_sales** and the **sr_ticket_number** and **sr_item_sk** columns in **store_returns**) but the multicolumn correlation is not considered.

2. After the **rows** hint is used for optimization, the plan is as follows and the statement execution takes 318s:

```
select avg(netpaid) from (select /*+rows(store_sales store_returns * 11270)*/ c_last_name ...
```

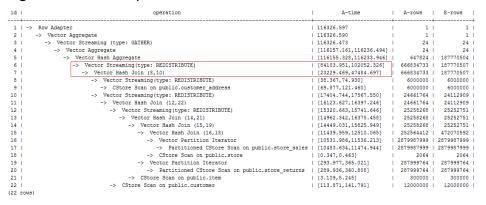
Figure 13-11 Using rows hints for optimization

The execution takes a longer time because layer-9 **redistribute** is slow. Considering that data skew does not occur at layer-9 **redistribute**, the slow redistribution is caused by the slow layer-8 **hashjoin** due to data skew at layer-18 **redistribute**.

3. Data skew occurs at layer-18 redistribute because customer_address has a few different values in its two join keys. Therefore, plan customer_address as the last one to be joined. After the hint is used for optimization, the plan is as follows and the statement execution takes 116s:

```
select avg(netpaid) from
(select /*+rows(store_sales store_returns *11270)
leading((store_sales store_returns store item customer) customer_address)*/
c_last_name ...
```

Figure 13-12 Hint optimization



Most of the time is spent on layer-6 **redistribute**. The plan needs to be further optimized.

4. The last layer redistribute contains skew. Therefore, it takes a long time. To avoid the data skew, plan the **item** table as the last one to be joined because the number of rows is not reduced after **item** is joined. After the hint is used for optimization, the plan is as follows and the statement execution takes 120s:

```
select avg(netpaid) from
(select /*+rows(store_sales store_returns *11270)
leading((customer_address (store_sales store_returns store customer) item))
c_last_name ...
```

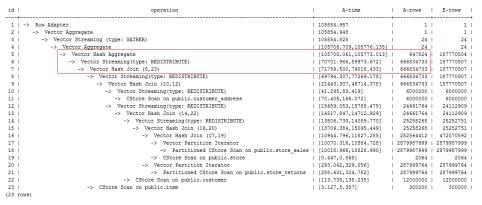
operation A-rows E-rows 1 | -> Row Adapter 120377.258 1 | 1 Vector Aggregate 120377.245 Vector Streaming (type: GATHER) 120377.091 -> Vector Aggregate [120184.884,120301.704] 187770504 Vector Hash Aggregate [120183.119,120297.845] 647824 Vector Streaming(type: REDISTRIBUTE)
-> Vector Hash Join (8,22) [87775.682,106070.878] [22323.764,49878.523] 666834733 187770507 187770507 666834733 -> Vector Hash Join (9,11)
-> Vector Streaming(type: REDISTRIBUTE) [21129.236,45208.255] [37.859,75.412] 666834733 187770507 6000000 6000000 10 | -> CStore Scan on public.customer_address Vector Streaming(type: REDISTRIBUTE) [74.798,114.449] 6000000 [15714.458,15824.928] 24661764 24112909 -> Vector Hash Join (13,21) [14637.516,14955.464] [13898.593,14333.200] 24661764 24112909 Vector Streaming(type: REDISTRIBUTE) 25258268 25252751 -> Vector Hash Join (15,19)
-> Vector Hash Join (16,18) [14166.917,15378.244] [11272.239,12052.532] 25258268 25252751 252564412 472070592 16 17 -> Vector Partition Iterator
-> Partitioned CStore Scan on public.store sales [10409.566,11127.981] 2879987999 2879987999 [10365.838.11077.601] 2879987999 -> CStore Scan on public.store Vector Partition Iterator [0.431,0.609] 2064 287999764 287999764 [343,780,408,254] Partitioned CStore Scan on public.store returns [339.844,403.923] 287999764 287999764 CStore Scan on public.customer [117,234,163,598] 12000000 Vector Streaming(type: BROADCAST) [44.571,130.129] 7200000 CStore Scan on public.item [4.169, 6.347]

Figure 13-13 Modifying hints and executing statements

Data skew occurs after the join of **item** and **customer_address** because **item** is broadcasted at layer-22. As a result, layer-6 **redistribute** is still slow.

5. Add a hint to disable **broadcast** for **item** or add a **redistribute** hint for the join result of **item** and **customer_address**. After the hint is used for optimization, the plan is as follows and the statement execution takes 105s: select avg(netpaid) from (select /*+rows(store_sales store_returns *11270) leading((customer_address (store_sales store_returns store customer) item)) no broadcast(item)*/ c_last_name ...

Figure 13-14 Execution plan



6. The last layer uses single-layer **Agg** and the number of rows is greatly reduced. Set **best_agg_plan** to **3** and change the single-layer **Agg** to a double-layer **Agg**. The plan is as follows and the statement execution takes 94s. The optimization ends.

Figure 13-15 Final optimization plan

id	operation	1	A-time	1	A-rows	ļ.	E-rows	ļ
1			94004.670	1	1		1	٠.
2	-> Vector Aggregate	1	94004.655	-1	1	1	1	í
3	-> Vector Streaming (type: GATHER)	1	94004.504	-1	24	1	24	Ĺ
4	-> Vector Aggregate	1	[93833.832,93928.052]	-1	24	1	24	ı
5	-> Vector Hash Aggregate	1	[93832.460,93926.412]	-1	647824	1	187770507	Ĺ
6	-> Vector Streaming(type: REDISTRIBUTE)	1	[93640.866,93787.939]	-1	647824	1	183912384	ĺ
7	-> Vector Hash Aggregate	1	[93687.544,93791.242]	-1	647824	1	183912384	ı
8	-> Vector Hash Join (9,24)	1	[70025.469,72773.161]	-1	666834733	1	187770507	ı
9	-> Vector Streaming(type: REDISTRIBUTE)	1	[68242.223,71275.972]	-1	666834733	1	187770507	ı
10	-> Vector Hash Join (11,13)	1	[21421.136,44830.306]	-1	666834733	l -	187770507	ı
11	-> Vector Streaming(type: REDISTRIBUTE)	1	[35.444,71.328]	-1	6000000	1	6000000	ı
12	-> CStore Scan on public.customer_address	1	[67.246,119.224]	-1	6000000	1	6000000	ı
13	-> Vector Streaming(type: REDISTRIBUTE)	1	[16089.853,16212.570]	-1	24661764	1	24112909	ı
14	-> Vector Hash Join (15,23)	1	[14822.972,15188.942]	-1	24661764	1	24112909	ı
15	-> Vector Streaming(type: REDISTRIBUTE)	1	[14061.867,14604.162]	-1	25258268	1	25252751	ı
16	-> Vector Hash Join (17,21)	1	[13949.756,15492.311]	-1	25258268	1	25252751	ı
17	-> Vector Hash Join (18,20)	1	[10935.742,12160.719]	-1	252564412	1	472070592	I
18	-> Vector Partition Iterator	1	[10052.958,11194.962]	- 1	2879987999	1 2	2879987999	ı
19	-> Partitioned CStore Scan on public.store_sales	1	[10008.415,11143.984]	-1	2879987999	1 2	2879987999	I
20	-> CStore Scan on public.store		[0.452,0.839]	- 1	2064		2064	
21	-> Vector Partition Iterator	1	[298.235,332.736]	- 1	287999764	1	287999764	I
22	-> Partitioned CStore Scan on public.store_returns	1	[294.067,327.629]	-1	287999764	1	287999764	I
23	-> CStore Scan on public.customer		[114.377,145.156]	-1	12000000		12000000	
24	-> CStore Scan on public.item	1	[3.150,3.530]	-1	300000	1	300000	I
(24:	cows)							

If the query performance deteriorates due to statistics changes, you can use hints to optimize the query plan. Take TPCH-Q17 as an example. The query performance deteriorates after the value of **default_statistics_target** is changed from the default one to **-2** for statistics collection.

 If default_statistics_target is set to the default value 100, the plan is as follows:

Figure 13-16 Default statistics

2. If default_statistics_target is set to -2, the plan is as follows.

Figure 13-17 Changes in statistics

 After the analysis, the cause is that the stream type is changed from BroadCast to Redistribute during the join of the lineitem and part tables.
 You can use a hint to change the stream type back to BroadCast. The figure below shows an example.

Figure 13-18 Statements

```
select /*+ no redistribute(part lineitem) */
    sum(l extendedprice) / 7.0 as avg yearly
from
   lineitem,
   part
where
   p partkey = 1 partkey
    and p brand = 'Brand#23'
    and p_container = 'MED BOX'
    and l_quantity < (
       select
            0.2 * avg(l quantity)
       from
            lineitem
        where
            l partkey = p partkey
    );
```

13.4.10 Routinely Maintaining Tables

To ensure proper database running, after INSERT and DELETE operations, you need to routinely do **VACUUM FULL** and **ANALYZE** as appropriate for customer scenarios and update statistics to obtain better performance.

Related Concepts

You need to routinely run **VACUUM**, **VACUUM FULL**, and **ANALYZE** to maintain tables, because:

- VACUUM FULL reclaims disk space occupied by updated or deleted data and combines small-size data files.
- VACUUM maintains a visualized mapping to track pages that contain arrays
 visible to other active transactions. A common index scan uses the mapping to
 obtain the corresponding array and check whether pages are visible to the
 current transaction. If the array cannot be obtained, the visibility is checked by
 fetching stack arrays. Therefore, updating the visible mapping of a table can
 accelerate unique index scans.
- **VACUUM** can avoid old data loss caused by duplicate transaction IDs when the number of executed transactions exceeds the database threshold.
- **ANALYZE** collects statistics on tables in databases. The statistics are stored in the PG_STATISTIC system catalog. Then, the query optimizer uses the statistics to work out the most efficient execution plan.

Procedure

Step 1 Run the **VACUUM** or **VACUUM FULL** command to reclaim disk space.

VACUUM:

Do VACUUM to the table:

VACUUM customer,

VACUUM

This command can be concurrently executed with database operation commands, including **SELECT**, **INSERT**, **UPDATE**, and **DELETE**; excluding **ALTER TABLE**.

Do **VACUUM** to the partitioned table:

VACUUM customer_par PARTITION (P1);

VACUUM

VACUUM FULL:

VACUUM FULL customer;

VACUUM

VACUUM FULL needs to add exclusive locks on tables it operates on and requires that all other database operations be suspended.

When reclaiming disk space, you can query for the session corresponding to the earliest transactions in the cluster, and then end the earliest long transactions as needed to make full use of the disk space.

- a. Run the following command to query for oldestxmin on the GTM: select * from pgxc_gtm_snapshot_status();
- b. Run the following command to query for the PID of the corresponding session on the CN. *xmin* is the oldestxmin obtained in the previous step. select * from pgxc_running_xacts() where xmin=1400202010;

Step 2 Do **ANALYZE** to update statistical information.

ANALYZE customer,

ANALYZE

Do **ANALYZE VERBOSE** to update statistics and display table information.

ANALYZE VERBOSE customer;

ANALYZE

You can use **VACUUM ANALYZE** at the same time to optimize the query.

VACUUM ANALYZE customer,

VACUUM

VACUUM and **ANALYZE** cause a substantial increase in I/O traffic, which may cause poor performance of other active sessions. Therefore, you are advised to set by specifying the **vacuum_cost_delay** parameter.

Step 3 Delete a table

DROP TABLE customer, DROP TABLE customer_par,

DROP TABLE part,

If the following output is displayed, the index has been deleted.

DROP TABLE

----End

Maintenance Suggestion

- Routinely do VACUUM FULL to large tables. If the database performance deteriorates, do VACUUM FULL to the entire database. If the database performance is stable, you are advised to monthly do VACUUM FULL.
- Routinely do VACUUM FULL to system catalogs, mainly PG_ATTRIBUTE.
- The automatic vacuum process (AUTOVACUUM) in the system automatically runs the VACUUM and ANALYZE statements to reclaim the record space marked as the deleted state and to update statistics related to the table.

13.4.11 Routinely Recreating an Index

Context

When data deletion is repeatedly performed in the database, index keys will be deleted from the index page, resulting in index distention. Recreating an index routinely improves query efficiency.

The database supports B-tree, GIN, and psort indexes.

- Recreating a B-tree index helps improve guery efficiency.
 - If massive data is deleted, index keys on the index page will be deleted.
 As a result, the number of index pages reduces and index bloat occurs.
 Recreating an index helps reclaim wasted space.
 - In the created index, pages adjacent in its logical structure are adjacent in its physical structure. Therefore, a created index achieves higher access speed than an index that has been updated for multiple times.
- You are advised not to recreate a non-B-tree index.

Rebuilding an Index

Use either of the following two methods to recreate an index:

• Run the **DROP INDEX** statement to delete an index and run the **CREATE INDEX** statement to create an index.

When you delete an index, a temporary exclusive lock is added in the parent table to block related read/write operations. When you create an index, the write operation is locked but the read operation is not. The data is read and scanned by order.

- Run the **REINDEX** statement to recreate an index:
 - When you run the REINDEX TABLE statement to recreate an index, an exclusive lock is added to block related read/write operations.
 - When you run the REINDEX INTERNAL TABLE statement to recreate an index for a desc table (), an exclusive lock is added to block read/write operations on the table.

Procedure

Assume the ordinary index areaS_idx exists in the **area_id** column of the imported table **areaS**. Use either of the following two methods to recreate an index:

 Run the DROP INDEX statement to delete the index and run the CREATE INDEX statement to create an index.

- Delete an index.DROP INDEX areaS_idx;DROP INDEX
- Create an index.
 CREATE INDEX areaS_idx ON areaS (area_id);
 CREATE INDEX
- Run the **REINDEX** statement to recreate an index.
 - Run the REINDEX TABLE statement to recreate an index. REINDEX TABLE areaS; REINDEX
 - Run the REINDEX INTERNAL TABLE statement to recreate an index for a desc table ().
 REINDEX INTERNAL TABLE areaS;
 REINDEX

13.4.12 Automatic Retry upon SQL Statement Execution Errors

With automatic retry (referred to as CN retry), GaussDB(DWS) retries an SQL statement when the execution of a statement fails. If an SQL statement sent from the **gsql** client, JDBC driver, or ODBC driver fails to be executed, the CN can automatically identify the error reported during execution and re-deliver the task to retry.

The restrictions of this function are as follows:

- Functionality restrictions:
 - CN retry increases execution success rate but does not guarantee success.
 - CN retry is enabled by default. In this case, the system records logs about temporary tables. If it is disabled, the system will not record the logs.
 Therefore, do not repeatedly enable and disable CN retry when temporary tables are used. Otherwise, data inconsistency may occur after a CN retry following a primary/standby switchover.
 - CN retry is enabled by default. In this case, the unlogged keyword is ignored in the statement for creating unlogged tables and thereby ordinary tables will be created by using this statement. If CN retry is disabled, the system records logs about unlogged tables. Therefore, do not repeatedly enable and disable CN retry when unlogged tables are used. Otherwise, data inconsistency may occur after a CN retry following a primary/standby switchover.
 - When GDS is used to export data, CN retry is supported. The existing mechanism checks for duplicate files and deletes duplicate files during data export. Therefore, you are advised not to repeatedly export data for the same foreign table unless you are sure that files with the same name in the data directory need to be deleted.
- Error type restrictions:
 - Only the error types in **Table 13-19** are supported.
- Statement type restrictions:
 - Support single-statement CN retry, stored procedures, functions, and anonymous blocks. Statements in transaction blocks are not supported.
- Statement restrictions of a stored procedure:
 - If an error occurs during the execution of a stored procedure containing EXCEPTION (including statement block execution and statement

execution in EXCEPTION), the stored procedure can be retried. If an internal error occurs, the stored procedure will retry first, but if the error is captured by **EXCEPTION**, the stored procedure cannot be retried.

- Packages that use global variables are not supported.
- **DBMS JOB** is not supported.
- **UTL_FILE** is not supported.
- If the stored procedure has printed information (such as dbms_output.put_line or raise info), the printed information will be output repeatedly when retry occurs, and "Notice: Retry triggered, some message may be duplicated." will be output before the repeated information.
- Cluster status restrictions:
 - Only DNs or GTMs are faulty.
 - The cluster can be recovered before the number of CN retries reaches the allowed maximum (controlled by max_query_retry_times). Otherwise, CN retry may fail.
 - CN retry is not supported during scale-out.
- Data import restrictions:
 - The COPY FROM STDIN statement is not supported.
 - The **gsql \copy from** metacommand is not supported.
 - JDBC CopyManager copyIn is not supported.

Table 13-19 lists the error types supported by CN retry and the corresponding error codes. You can use the GUC parameter **retry_ecode_list** to set the list of error types supported by CN retry. You are not advised to modify this parameter. To modify it, contact the technical support.

Table 13-19 Error types supported by CN retry

Error Type	Error Code	Remarks
CONNECTION_RESET_BY_PEER	YY00 1	TCP communication errors: Connection reset by peer (communication between the CN and DNs)
STREAM_CONNECTION_RESET_BY _PEER	YY00 2	TCP communication errors: Stream connection reset by peer (communication between DNs)
LOCK_WAIT_TIMEOUT	YY00 3	Lock wait timeout
CONNECTION_TIMED_OUT	YY00 4	TCP communication errors: Connection timed out
SET_QUERY_ERROR	YY00 5	Failed to deliver the SET command: Set query

Error Type	Error Code	Remarks	
OUT_OF_LOGICAL_MEMORY	YY00 6	Failed to apply for memory: Out of logical memory	
SCTP_MEMORY_ALLOC	YY00 7	SCTP communication errors: Memory allocate error	
SCTP_NO_DATA_IN_BUFFER	YY00 8	SCTP communication errors: SCTP no data in buffer	
SCTP_RELEASE_MEMORY_CLOSE	YY00 9	SCTP communication errors: Release memory close	
SCTP_TCP_DISCONNECT	YY01 0	SCTP communication errors: TCP disconnect	
SCTP_DISCONNECT	YY01 1	SCTP communication errors: SCTP disconnect	
SCTP_REMOTE_CLOSE	YY01 2	SCTP communication errors: Stream closed by remote	
SCTP_WAIT_POLL_UNKNOWN	YY01 3	Waiting for an unknown poll: SCTP wait poll unknown	
SNAPSHOT_INVALID	YY01 4	Snapshot invalid	
ERRCODE_CONNECTION_RECEIVE _WRONG	YY01 5	Connection receive wrong	
OUT_OF_MEMORY	5320 0	Out of memory	
CONNECTION_FAILURE	0800 6	GTM errors: Connection failure	
CONNECTION_EXCEPTION	0800 0	Failed to communicate with DNs due to connection errors: Connection exception	
ADMIN_SHUTDOWN	57P0 1	System shutdown by administrators: Admin shutdown	
STREAM_REMOTE_CLOSE_SOCKET	XX00 3	Remote socket disabled: Stream remote close socket	
ERRCODE_STREAM_DUPLICATE_Q UERY_ID	XX00 9	Duplicate query id	
ERRCODE_STREAM_CONCURRENT _UPDATE	YY01 6	Stream concurrent update	
ERRCODE_LLVM_BAD_ALLOC_ERR OR	CG00 3	Memory allocation error: Allocate error	

Error Type	Error Code	Remarks
ERRCODE_LLVM_FATAL_ERROR	CG00 4	Fatal error
HashJoin temporary file reading error (ERRCODE_HASHJOIN_TEMP_FILE _ERROR).	F001 1	File error
Buffer file reading error (ERRCODE_BUFFER_FILE_ERROR)	F001 2	File reading error
Partition number error (ERRCODE_PARTITION_NUM_CHANGED).	4500	During scanning on a list partition table, it is found that the number of partitions is different from that in the optimization phase. This problem usually occurs when the queries and ADD/DROP partitions are concurrently executed. (This error is supported only by clusters of version 8.1.3 or later.)
Unmatched schema name (ERRCODE_UNMATCH_OBJECT_SC HEMA)	42P3 0	Unmatched schema name

To enable CN retry, set the following GUC parameters:

 Mandatory GUC parameters (required by both CNs and DNs) max_query_retry_times

CAUTION

If CN retry is enabled, temporary table data is logged. For data consistency, do not switch the enabled/disabled status for CN retry when the temporary tables are being used by sessions.

 Optional GUC parameters cn_send_buffer_size max_cn_temp_file_size

13.4.13 Query Band Load Identification

Overview

GaussDB(DWS) implements load identification and intra-queue priority control based on query_band. It provides more flexible load identification methods and identifies load queues based on job types, application names, and script names. Users can flexibly configure query_band identification queues based on service

scenarios. In addition, priority control of job delivery in the queue is implemented. In the future, priority control of resources in the queue will be gradually implemented.

Administrators can configure the queue associated with query_band and estimate the memory limit based on service scenarios and job types to implement more flexible load control and resource management and control. If query_band is not configured for the service or the user does not associate query_band with an action, the queue associated with the user and the priority in the queue is used by default.

Load Behaviors Supported by query_band

query_band is a session-level GUC parameter. It is a job identifier of the character data type. Its value can be any string. However, for easier differentiation and configuration, query_band only identifies key-value pairs. For example:

SET query_band='JobName=abc;AppName=test;UserName=user';

JobName=abc, **AppName=test**, and **UserName=user** are independent key-value pairs. Specifications of the query_band key-value pairs:

- query_band is set in key-value pair mode, that is, 'key=value'. Multiple query_band key-value pairs can be set in a session. Multiple key-value pairs are separated by semicolons (;). The maximum length of both the **query_band** key-value pair and parameter value is 1024 characters.
- The query_band key-value pair supports the following valid characters: digits 0 to 9, uppercase letters A to Z, lowercase letters a to z, '.', '-', '_', and '#'.

query_band is configured, and identifies load behaviors, using key-value pairs. The supported load behaviors are described in **Table 13-20**.

Table 13-20 L	oad be	haviors su	upported	by QUERY_	BAND

Туре	Behavior	Behavior Description
Workload management (workload)	Resource pool (respool)	query_band associated with a resource pool
Workload management (workload)	Priority	Priority in the queue
Order	Queue (respool) Currently, this field is invalid and is used for future extension.	query_band query order

The "Type" is used to classify load behaviors. Different load behaviors may belong to a same type. For example, both "Resource pool" and a "Priority" belong to

"Workload management". The "Behavior" indicates a load behavior associated with a query_band key-value pair. The "Behavior description" describes a specific load behavior. The "Order" in the "Type" is used to indicate the priority of the query_band load behavior identification. When a session has multiple query_band key-value pairs, the query_band key-value pair with a smaller order value is preferentially used to identify a load behavior. Each query_band key-value pair can have multiple associated load behaviors, while one load behavior can only have one associated key-value pair. The query_band load behavior is described as follows:

- Resource pool: query_band can be associated with resource pools. During job
 execution, if a resource pool is associated with query_band, the resource pool
 is used in preference. Otherwise, the resource pool associated with the user is
 used.
 - When query_band is associated with a resource pool, an error is reported if the resource pool does not exist, and the association fails.
 - When query_band is associated with a resource pool, the dependency between query_band and the resource pool is recorded.
 - When a resource pool associated with query_band is deleted, a message is displayed indicating that the resource pool fails to be deleted because of the dependency between query_band and the resource pool.
- Intra-queue priority: query_band can be associated with job priorities, including high, medium, and low. Rush is provided as a special priority (green channel). The default priority is medium. In practice, most jobs use the medium priority, low-priority jobs use the low priority, and privileged jobs use the high priority. It is not recommended that a large number of jobs use the high priority. The rush priority is used only in special scenarios and is not recommended in normal cases.

The intra-queue priority is used to implement the queuing priority.

- In the static load management scenario, when the CN concurrency is insufficient, CN global queuing is triggered. The CN global queue is a priority queue.
- In the dynamic load management scenario, if the DN memory is insufficient, CCN global queuing is triggered. The CCN global queue is a priority queue.
- When the resource pool concurrency or memory is insufficient, resource pool queuing is triggered. The resource pool queue is a priority queue.

The preceding priority queues comply with the following scheduling rules:

- Jobs with a higher priority are scheduled first.
- After all jobs with a high priority are scheduled, jobs with a low priority are scheduled.
- In dynamic load management scenarios, the CN global queue does not support the query_band priority.
- Order: The identification order of query_bands can be configured. The default order value is -1. Except the default order value, there are no two query_bands with the same order value. The query_band order is verified when being configured. If there are query_bands with the same order value, the order values are recursively increased by 1 until there are no query_bands with the same order value.

- If a session has multiple query_band key-value pairs, the query_band key-value pair with a smaller order value is used for load identification.
- 0 is the smallest order value, and the default order value -1 is the largest order value.
- If the query_bands are all of the same order value, the anterior query_band is used for load identification.
- For example, if in set query_band='b=1;a=3;c=1'; b=1, the order value of b=1 is -1, a=3 is 4, c=1 is 1, c=1 is used as the query_band for load identification. This design enables load administrators to adjust load scheduling.

Application and Configuration of query_band

- The pg_workload_action cross-database system catalog is used to store the query band action and order. For details, see PG WORKLOAD ACTION.
- The default action and order are not stored in the **pg_workload_action** system catalog. If a non-default action is set for query_band, the default action is also displayed when actions are queried. The message <query_band information not found> is displayed when the action and order to be queried are the default query band action.
- The gs_wlm_set_queryband_action function sets the query_band sequence.
 The maximum length of the first parameter, that is, the query_band key value pair, is 63 characters. For the second parameter, it is case insensitive and multiple actions are separated by semicolons (;). order is the default parameter and its default value is -1. For details, see gs_wlm_set_queryband_action.
- The gs_wlm_set_queryband_order function sets the query_band sequence. The maximum length of the first parameter, that is, a query_band key value pair, is 63 characters. The value of query_band must be greater than or equal to -1. Except the default value -1, the value of query_band order must be unique. When setting the query_band order, if there are query_bands with the same order value, the original order value is increased by 1. For details, see gs wlm set queryband order.
- You can use the **gs_wlm_get_queryband_action** function to query the **query_band** action. For details, see **gs_wlm_set_queryband_action**.
- **pg_queryband_action** provides the system view for querying all query_band actions. For details, see **PG_OUERYBAND_ACTION**.
- The query_band priority is displayed as an integer in the load management view (PG_SESSION_WLMSTAT). The mapping between numbers and priorities is as follows:
 - 0: not controlled by load management
 - 1: low
 - 2: medium
 - 4: high
 - 8: rush
- Permission control: Except initial users, other users have the permission to set and query query_band only when they are authorized.

Ⅲ NOTE

• When all running jobs are canceled in batches or the maximum number of concurrent jobs in a queue is 1 and only one queue is running jobs, the CN may be triggered to automatically wake up jobs. As a result, jobs are not delivered by priority.

Examples

Step 1 Set the associated resource pool to **p1**, priority to **rush**, and order to **1** for query_band **JobName to abc**.

```
SELECT * FROM gs_wlm_set_queryband_action('JobName=abc','respool=p1;priority=rush',1);
gs_wlm_set_queryband_action
------
t
(1 row)
```

Step 2 Change the associated resource pool to **p2** for query_band **JobName=abc**.

```
SELECT * FROM gs_wlm_set_queryband_action('JobName=abc','respool=p2');
gs_wlm_set_queryband_action
-----
t
(1 row)
```

Step 3 Change the priority to **high** for query_band **JobName=abc**.

```
SELECT * FROM gs_wlm_set_queryband_action('JobName=abc','priority=high');
gs_wlm_set_queryband_action
------
t
(1 row)
```

Step 4 Change the order to **3** for query_band **JobName=abc**.

```
SELECT * FROM gs_wlm_set_queryband_order('JobName=abc',3);
gs_wlm_set_queryband_order
------
t
(1 row)
```

Step 5 Query the load behaviors associated with query_band.

Step 6 In **query_band**, set the priority of **AppName=test** to **Low**, associate the user with the resource pool, and use the default sequence.

```
SELECT * FROM gs_wlm_set_queryband_action('AppName=test','priority=low');
gs_wlm_set_queryband_action
------
t
(1 row)
```

Step 7 Query the load behaviors associated with query_band.

Step 8 In **query_band**, cancel all the workload behaviors associated with **JobName=abc** and set them to default behaviors.

```
SELECT * FROM gs_wlm_set_queryband_action('JobName=abc','respool=null;priority=medium',-1); NOTICE: The respool of query_band(JobName=abc) will be removed.
```

```
NOTICE: The priority of query_band(JobName=abc) will be removed.
gs_wlm_set_queryband_action
------
t
(1 row)
```

Step 9 Query the load behaviors associated with query_band.

----End

13.5 SQL Tuning Examples

13.5.1 Case: Selecting an Appropriate Distribution Column

Distribution columns are used to distribute data to different nodes. A proper distribution key can avoid data skew.

When performing join query, you are advised to select the join condition in the query as the distribution key. When a join condition is used as a distribution key, related data is distributed locally on DNs, reducing the cost of data flow between DNs and improving the query speed.

Before optimization

Use a as the distribution column of t1 and t2. The table definition is as follows:

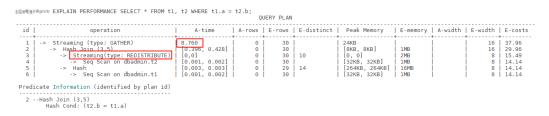
```
CREATE TABLE t1 (a int, b int) DISTRIBUTE BY HASH (a);
CREATE TABLE t2 (a int, b int) DISTRIBUTE BY HASH (a);
```

The following guery is executed:

```
SELECT * FROM t1, t2 WHERE t1.a = t2.b;
```

In this case, the execution plan contains **Streaming(type: REDISTRIBUTE)**, that is, the DN redistributes data to all DNs based on the selected column. This will cause a large amount of data to be transmitted between DNs, as shown in **Figure 13-19**.

Figure 13-19 Selecting an appropriate distribution column (1)



After optimization

Use the join condition in the query as the distribution key and run the following statement to changethe distribution key of **t2** as **b**:

ALTER TABLE t2 DISTRIBUTE BY HASH (b);

After the distribution column of table **t2** is changed to column **b**, the execution plan does not contain **Streaming(type: REDISTRIBUTE)**. This reduces the amount of communication data between DNs and reduces the execution time from 8.7 ms to 2.7 ms, improving query performance, as shown in **Figure 13-20**.

Figure 13-20 Selecting an appropriate distribution column (2)

```
| EXPLAIN PERFORMANCE SELECT * FROM t1, t2 WHERE t1.a = t2.b;
| OUERY PLAN | | OUERY PLAN | | OUERY PLAN | | OUERY PLAN | | OUERY PLAN | | OUERY PLAN | | OUERY PLAN | | OUERY PLAN | | OUERY PLAN | | OUERY PLAN | | OUERY PLAN | | OUERY PLAN | | OUERY PLAN | OUERY PLAN | | OUERY PLAN | OUERY PLAN | | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUERY PLAN | OUER
```

13.5.2 Case: Creating an Appropriate Index

Creating a proper index can accelerate the retrieval of data rows in a table. Indexes occupy disk space and reduce the speed of adding, deleting, and updating rows. If data needs to be updated very frequently or disk space is limited, you need to limit the number of indexes. Create indexes for large tables. Because the more data in the table, the more effective the index is. You are advised to create indexes on:

- Columns that need to be queried frequently
- Joined columns. For a query on joined columns, you are advised to create a composite index on the joined columns. For example, if the join condition is select * from t1 join t2 on t1.a=t2.a and t1.b=t2.b. You can create a composite index on the a and b columns of table t1.
- Columns having filter criteria (especially scope criteria) of a where clause
- Columns that appear after order by, group by, and distinct

Before optimization

The column-store partitioned table orders is defined as follows:

Run the SQL statement to query the execution plan when no index is created. It is found that the execution time is 48 milliseconds.

EXPLAIN PERFORMANCE SELECT * FROM orders WHERE o_custkey = '1106459';



After optimization

The filtering condition column of the **where** clause is **o_custkey**. Add an index to the **o_custkey** column.

CREATE INDEX idx_o_custkey ON orders (o_custkey) LOCAL;

Run the SQL statement to query the execution plan after the index is created. It is found that the execution time is 18 milliseconds.



13.5.3 Case: Adding NOT NULL for JOIN Columns

If there are many **NULL** values in the **JOIN** columns, you can add the filter criterion **IS NOT NULL** to filter data in advance to improve the **JOIN** efficiency.

Before optimization

```
SELECT
FROM
( ( SELECT
 STARTTIME STTIME,
 SUM(NVL(PAGE DELAY MSEL,0)) PAGE DELAY MSEL,
 SUM(NVL(PAGE_SUCCEED_TIMES,0)) PAGE_SUCCEED_TIMES,
 SUM(NVL(FST_PAGE_REQ_NUM,0)) FST_PAGE_REQ_NUM,
 SUM(NVL(PAGE_AVG_SIZE,0)) PAGE_AVG_SIZE,
 SUM(NVL(FST_PAGE_ACK_NUM,0)) FST_PAGE_ACK_NUM,
 SUM(NVL(DATATRANS_DW_DURATION,0)) DATATRANS_DW_DURATION,
 SUM(NVL(PAGE_SR_DELAY_MSEL,0)) PAGE_SR_DELAY_MSEL
FROM
 PS.SDR_WEB_BSCRNC_1DAY SDR
 INNER JOIN (SELECT
   BSCRNC ID,
   BSCRNC_NAME,
   ACCESS_TYPE,
   ACCESS_TYPE_ID
  FROM
   nethouse.DIM_LOC_BSCRNC
   GROUP BY
   BSCRNC_ID,
   BSCRNC_NAME,
   ACCESS_TYPE,
   ACCESS_TYPE_ID) DIM
 ON SDR.BSCRNC_ID = DIM.BSCRNC_ID
 AND DIM.ACCESS_TYPE_ID IN (0,1,2)
 INNER JOIN nethouse.DIM_RAT_MAPPING RAT
 ON (RAT.RAT = SDR.RAT)
WHERE
( (STARTTIME >= 1461340800
```

```
AND STARTTIME < 1461427200) )
AND RAT.ACCESS_TYPE_ID IN (0,1,2)
GROUP BY STTIME ) );
```

Figure 13-21 shows the execution plan.

Figure 13-21 Adding NOT NULL for JOIN columns (1)



After optimization

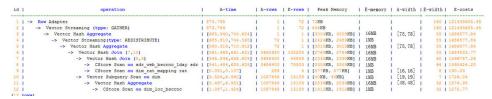
- 1. As shown in Figure 13-21, the sequential scan phase is time consuming.
- 2. The JOIN performance is poor because a large number of null values exist in the JOIN column **BSCRNC_ID** of the PS.SDR_WEB_BSCRNC_1DAY table.

Therefore, you are advised to manually add **NOT NULL** for **JOIN** columns in the statement, as shown below:

```
SELECT
FROM
((SELECT
 STARTTIME STTIME,
 SUM(NVL(PAGE_DELAY_MSEL,0)) PAGE_DELAY_MSEL,
 SUM(NVL(PAGE_SUCCEED_TIMES,0)) PAGE_SUCCEED_TIMES,
 SUM(NVL(FST_PAGE_REQ_NUM,0)) FST_PAGE_REQ_NUM,
 SUM(NVL(PAGE_AVG_SIZE,0)) PAGE_AVG_SIZE,
 SUM(NVL(FST_PAGE_ACK_NUM,0)) FST_PAGE_ACK_NUM,
 SUM(NVL(DATATRANS_DW_DURATION,0)) DATATRANS_DW_DURATION,
 SUM(NVL(PAGE_SR_DELAY_MSEL,0)) PAGE_SR_DELAY_MSEL
 PS.SDR_WEB_BSCRNC_1DAY SDR
 INNER JOIN (SELECT
   BSCRNC_ID,
   BSCRNC_NAME,
   ACCESS_TYPE,
   ACCESS_TYPE_ID
  FROM
   nethouse.DIM_LOC_BSCRNC
  GROUP BY
   BSCRNC ID,
   BSCRNC_NAME,
   ACCESS TYPE,
   ACCESS_TYPE_ID) DIM
 ON SDR.BSCRNC_ID = DIM.BSCRNC_ID
 AND DIM.ACCESS_TYPE_ID IN (0,1,2)
 INNER JOIN nethouse.DIM_RAT_MAPPING RAT
 ON (RAT.RAT = SDR.RAT)
WHERE
 ( (STARTTIME >= 1461340800
 AND STARTTIME < 1461427200) )
 AND RAT.ACCESS_TYPE_ID IN (0,1,2)
 and SDR.BSCRNC_ID is not null
GROUP BY
 STTIME ) ) A;
```

Figure 13-22 shows the execution plan.

Figure 13-22 Adding NOT NULL for JOIN columns (2)



13.5.4 Case: Pushing Down Sort Operations to DNs

In an execution plan, more than 95% of the execution time is spent on **window agg** performed on the CN. In this case, **sum** is performed for the two columns separately, and then another **sum** is performed for the separate sum results of the two columns. After this, trunc and sorting are performed in sequence. You can try to rewrite the statement into a subquery to push down the sorting operations.

Before optimization

The table structure is as follows:

CREATE TABLE public.test(imsi int,L4_DW_THROUGHPUT int,L4_UL_THROUGHPUT int) with (orientation = column) DISTRIBUTE BY hash(imsi);

The guery statements are as follows:

```
SELECT COUNT(1) over() AS DATACNT,
IMSI AS IMSI_IMSI,
CAST(TRUNC(((SUM(L4_UL_THROUGHPUT) + SUM(L4_DW_THROUGHPUT))), 0) AS
DECIMAL(20)) AS TOTAL_VOLOME_KPIID
FROM public.test AS test
GROUP BY IMSI
ORDER BY TOTAL VOLOME KPIID DESC LIMIT 10;
```

The execution plan is as follows:

```
QUERY PLAN
     operation
                                       A-time | A-rows | E-rows | E-distinct | Peak Memory | E-
memory | A-width | E-width | E-costs
 1 | -> Row Adapter
                                      2862.008
                                                     | 10 | 10 |
                                                                         | 31KB
     | 28 | 48360.42
                                                      10 |
 2 |
     -> Vector Limit
                                     | 2861.969
                                                              10 |
                                                                        18KB
        28 | 48360.42
 3 |
                                     | 2861.946
                                                       10 | 1000000 |
                                                                          | 479KB
        -> Vector Sort
        | 28 | 50860.39
         -> Vector WindowAgg
                                                       | 1000000 | 1000000 |
                                                                                 I 69987KB
                                        2166.759
            | 28 | 26750.75
 5 |
           -> Vector Streaming (type: GATHER) | 136.813
                                                          | 1000000 | 1000000 |
208KB
                 | 28 | 15500.75
             -> Vector Sonic Hash Aggregate | [71.374, 73.640] | 1000000 | 1000000 |
 6 |
                                                                                    | [14MB,
14MB] | 96MB(2919MB) | [31,31] | 28 | 15032.00
              -> CStore Scan on public.test | [2.957, 2.994] | 1000000 | 1000000 |
                                                                                 | [1MB,
1MB] | 1MB | | 12 | 1282.00
```

As we can see, both **window agg** and **sort** are performed on the CN, which is time consuming.

After optimization

Modify the statement to a subquery statement, as shown below:

```
SELECT COUNT(1) over() AS DATACNT, IMSI_IMSI, TOTAL_VOLOME_KPIID
FROM (SELECT IMSI AS IMSI_IMSI,
CAST(TRUNC(((SUM(L4_UL_THROUGHPUT) + SUM(L4_DW_THROUGHPUT))),
0) AS DECIMAL(20)) AS TOTAL_VOLOME_KPIID
FROM public.test AS test
GROUP BY IMSI
ORDER BY TOTAL_VOLOME_KPIID DESC LIMIT 10);
```

Perform **sum** on the **trunc** results of the two columns, take it as a subquery, and then perform **window agg** for the subquery to push down the sorting operation to DNs, as shown below:

```
QUERY PLAN
                                 | A-time | A-rows | E-rows | E-distinct | Peak
Memory | E-memory | A-width | E-width | E-costs
 NOW Adapter | 955.277 | 10 | 5 | | 31KB
| 1572KB
4 | -> Vector Limit | [0.018, 0.018] | 10 | 10 | 1MB | 28 | 25836.97
                                                                    | [8KB, 8KB] |
         -> Vector Streaming(type: BROADCAST)
                                          | [0.014, 0.014] | 20 |
[719KB, 719KB] | 28 | 25837.12
 6 | -> Vector Limit
                                  | [927.730, 934.283] | 20 | 20 |
                                                                     | [8KB, 8KB]
         | 28 | 25836.85
-> Vector Sort
| 1MB
                                  | [927.720, 934.269] | 20 | 1000000 |
                                                                       | [463KB.
463KB] | 16MB | [32,32] | 28 | 27086.82 | -> Vector Sonic Hash Aggregate | [456.841, 461.077] | 1000000 | 1000000 |
[15MB, 15MB] | 96MB(2916MB) | [31,31] | 28 | 15032.00
9 | -> CStore Scan on public.test | [2.959, 3.014] | 1000000 | 1000000 |
1MB] | 1MB | 12 | 1282.00
                                                                           Ι [1MB.
```

The optimized SQL statement greatly improves the performance by reducing the execution time from 2.862s to 0.955s. Note that the optimization result in this example is for reference only. Due to the uncertainty of **WindowAgg**, the optimized result set is related to the actual service.

13.5.5 Case: Configuring cost_param for Better Query Performance

The cost_param parameter is used to control use of different estimation methods in specific customer scenarios, allowing estimated values to be close to onsite values. This parameter can control various methods simultaneously by performing AND (&) operations on the bit for each method. A method is selected if its value is not **0**.

Scenario 1: Before Optimization

If **bit0** of **cost_param** is set to **1**, an improved mechanism is used for estimating the selection rate of non-equi-joins. This method is more accurate for estimating the selection rate of joins between two identical tables. The following example describes the optimization scenario when **bit0** of **cost_param** is set to **1**. In V300R002C00 and later, **cost_param & 1=0** is not used. That is, an optimized formula is selected for calculation.

The selection rate indicates the percentage for which the number of rows meeting the join conditions account of the **JOIN** results when the **JOIN** relationship is established between two tables.

The table structure is as follows:

```
CREATE TABLE LINEITEM
L ORDERKEY BIGINT NOT NULL
, L_PARTKEY BIGINT NOT NULL
, L_SUPPKEY BIGINT NOT NULL
, L_LINENUMBER BIGINT NOT NULL
, L_QUANTITY DECIMAL(15,2) NOT NULL
, L_EXTENDEDPRICE DECIMAL(15,2) NOT NULL
, L_DISCOUNT DECIMAL(15,2) NOT NULL
, L_TAX DECIMAL(15,2) NOT NULL
, L_RETURNFLAG CHAR(1) NOT NULL
, L_LINESTATUS CHAR(1) NOT NULL
, L_SHIPDATE DATE NOT NULL
, L_COMMITDATE DATE NOT NULL
, L_RECEIPTDATE DATE NOT NULL
, L_SHIPINSTRUCT CHAR(25) NOT NULL
, L_SHIPMODE CHAR(10) NOT NULL
, L_COMMENT VARCHAR(44) NOT NULL
) with (orientation = column, COMPRESSION = MIDDLE) distribute by hash(L_ORDERKEY);
CREATE TABLE ORDERS
O_ORDERKEY BIGINT NOT NULL
, O_CUSTKEY BIGINT NOT NULL
, O_ORDERSTATUS CHAR(1) NOT NULL
, O_TOTALPRICE DECIMAL(15,2) NOT NULL
, O_ORDERDATE DATE NOT NULL
, O_ORDERPRIORITY CHAR(15) NOT NULL
, O_CLERK CHAR(15) NOT NULL
, O_SHIPPRIORITY BIGINT NOT NULL
, O_COMMENT VARCHAR(79) NOT NULL
)with (orientation = column, COMPRESSION = MIDDLE) distribute by hash(O_ORDERKEY);
```

The query statements are as follows:

```
explain verbose select
count(*) as numwait
from
lineitem l1,
orders
where
o_orderkey = l1.l_orderkey
and o orderstatus = 'F'
and l1.l_receiptdate > l1.l_commitdate
and not exists (
select
from
lineitem l3
where
l3.l_orderkey = l1.l_orderkey
and l3.l_suppkey <> l1.l_suppkey
and l3.l_receiptdate > l3.l_commitdate
order by
numwait desc;
```

The following figure shows the execution plan. (When **verbose** is used, **distinct** is added for column selection which is controlled by **cost off/on**. The hash join rows show the estimated number of distinct values and the other rows do not.)

id					E-width E-costs
1	-> Row Adapter	1	1		8 39.36
2	-> Vector Sort	1	1		8 39.36
3	-> Vector Aggregate	1	1		8 39.34
4	-> Vector Streaming (type: GATHER)	1	2		8 39.34
5	-> Vector Aggregate	1	2		8 39.25
6	-> Vector Hash Anti Join (7, 10)	1	2	4, 5	0 39.24
7	-> Vector Hash Join (8,9)	1	2	200, 1	16 26.12
8	-> CStore Scan on public.lineitem 11	1	7		16 13.05
9	-> CStore Scan on public.orders	1	1		8 13.05
10	-> CStore Scan on public.lineitem 13	1	7		16 13.05

Scenario 1: After Optimization

These queries are from Anti Join connected in the **lineitem** table. When **cost_param & bit0** is **0**, the estimated number of Anti Join rows greatly differs from that of the actual number of rows, compromising the query performance. You can estimate the number of Anti Join rows more accurately by setting **cost_param & bit0** to **1** to improve the query performance. The optimized execution plan is as follows:

id	operation	E-rows	E-memory	E-width	E-costs
1	-> Row Adapter	1		0	9104892.37 9
2	-> Vector Sort	1	1	0	9104892.379
3	-> Vector Aggregate	1	1	0	9104892.358
4	-> Vector Streaming (type: GATHER)	48	1	0	9104892.358
5	-> Vector Aggregate	48	1MB	0	9104890.825
6	-> Vector Hash Join (7.12)	2526630903	929MB	0	8973295.454
7	-> Vector Hash Anti Join (8. 10)	1999996587	3178MB	8	7198231.14
8	-> Vector Partition Iterator	1999996587	1MB	16	3000158.25
9	-> Partitioned CStore Scan on public.lineitem 11	1999996587	1MB	16	3000158.25 1
10	-> Vector Partition Iterator	1999996587	1MB	16	3000158.25
11	-> Partitioned CStore Scan on public.lineitem 13	1999996587	1MB	16	3000158.25
12	-> Vector Partition Iterator	730839014	1MB	8	589611.00
13	-> Partitioned CStore Scan on public.orders	730839014	1MB	8	589611.00

Scenario 2: Before Optimization

If **bit1** is set to **1** (**set cost_param=2**), the selection rate is estimated based on multiple filter criteria. The lowest selection rate among all filter criteria, but not the product of the selection rates for two tables under a specific filter criterion, is used as the total selection rate. This method is more accurate when a close correlation exists between the columns to be filtered. The following example describes the optimization scenario when **bit1** of **cost_param** is set to **1**.

The table structure is as follows:

```
CREATE TABLE NATION
(
N_NATIONKEYINT NOT NULL
, N_NAMECHAR(25) NOT NULL
, N_REGIONKEYINT NOT NULL
, N_COMMENTVARCHAR(152)
) distribute by replication;
CREATE TABLE SUPPLIER
(
S_SUPPKEYBIGINT NOT NULL
, S_NAMECHAR(25) NOT NULL
, S_ADDRESSVARCHAR(40) NOT NULL
, S_NATIONKEYINT NOT NULL
, S_PHONECHAR(15) NOT NULL
, S_ACCTBALDECIMAL(15,2) NOT NULL
, S_COMMENTVARCHAR(101) NOT NULL
```

```
) distribute by hash(S_SUPPKEY);
CREATE TABLE PARTSUPP
(
PS_PARTKEYBIGINT NOT NULL
, PS_SUPPKEYBIGINT NOT NULL
, PS_AVAILQTYBIGINT NOT NULL
, PS_SUPPLYCOSTDECIMAL(15,2)NOT NULL
, PS_COMMENTVARCHAR(199) NOT NULL
) distribute by hash(PS_PARTKEY);
```

The query statements are as follows:

```
set cost_param=2;
explain verbose select
nation.
sum(amount) as sum_profit
from
select
n_name as nation,
l_extendedprice * (1 - l_discount) - ps_supplycost * l_quantity as amount
from
supplier,
lineitem.
partsupp,
nation
where
s_suppkey = l_suppkey
and ps_suppkey = l_suppkey
and ps_partkey = l_partkey
and s_nationkey = n_nationkey
) as profit
group by nation
order by nation;
```

When **bit1** of **cost param** is **0**, the execution plan is shown as follows:

id	operation				E-disting				
1	-> Sort	1	1						61.52
2	-> HashAggregate	1	1	1		- 1	208	I	61.51
3	-> Streaming (type: GATHER)	1	2	1		- 1	208	I	61.51
4	-> HashAggregate	1	2	1		- 1	208		61.36
5	-> Hash Join (6,7)	1	2	1	20, 15	- 1	176		61.33
6	-> Seg Scan on public.nation	1	40	I		- 1	108		20.20
7	-> Hash		2	ı		- 1	76	ı	41.04
8	-> Hash Join (9,16)	1	2	I	10, 13	- 1	76		41.04
9	-> Streaming(type: REDISTRIBUTE)		2	ı			88	I	27.73
10	-> Hash Join (11,14)	1	2	I	10, 13	- 1	88		27.62
11	-> Streaming(type: REDISTRIBUTE)	1	20	I		- 1	70		14.19
12	-> Row Adapter		21	ı			70		13.01
13	-> CStore Scan on public.lineitem	1	20	I		- 1	70		13.01
14	-> Hash	1	21	I		- 1	34		13.13
15	-> Seg Scan on public.partsupp		20	ı		- 1	34	ı	13.13
16	-> Hash	1	21	I		- 1	12		13.13
17	-> Seg Scan on public.supplier	1	20	1		- 1	12	I	13.13

Scenario 2: After Optimization

In the preceding queries, the hash join criteria of the supplier, lineitem, and partsupp tables are setting lineitem.l_suppkey to supplier.s_suppkey and lineitem.l_partkey to partsupp.ps_partkey. Two filter criteria exist in the hash join conditions. lineitem.l_suppkey in the first filter criteria and lineitem.l_partkey in the second filter criteria are two columns with strong relationship of the lineitem table. In this situation, when you estimate the rate of the hash join conditions, if cost_param & bit1 is 0, the selection rate is estimated based on multiple filter criteria. The lowest selection rate among all filter criteria, but not the product of the selection rates for two tables under a specific filter criterion, is used as the total selection rate. This method is more accurate when a

close correlation exists between the columns to be filtered. The plan after optimization is shown as follows:

```
id |
                                                                      | E-rows | E-distinct | E-width | E-costs
       -> HashAggregate
                                                                           10 I
                                                                                                 208 | 64.23
3 1
          -> Streaming (type: GATHER)
                                                                           20 I
                                                                                                208 | 64.23
            -> HashAggregate
                                                                                                208 | 62.71
                                                                           20 | 20, 10
                -> Hash Join (6,7)
                                                                                                176 | 62.46
                  -> Seg Scan on public.nation
                   -> Hash
                                                                                                 76 | 41.97
                     -> Hash Join (9,16)
                                                                           20 | 10, 13
                                                                                                 76 | 41.97
                        -> Streaming(type: REDISTRIBUTE)
                                                                           20 I
                                                                                                 82 | 28.54
                           -> Hash Join (11,14)
                                                                           20 | 10, 13
                                                                                                 82 | 27.63
10 |
                              -> Streaming(type: REDISTRIBUTE)
                                -> Row Adapter
                                                                                                 70 | 13.01
                                    -> CStore Scan on public.lineitem |
13 I
                                                                           20 I
                                                                                                 70 | 13.01
                              -> Hash
                                                                                                 12 | 13.13
14 I
                                                                           21 I
                                -> Seg Scan on public.supplier
                                                                                                 12 | 13.13
                                                                           20 |
                            -> Seg Scan on public.partsupp
                                                                                                 34 | 13.13
```

13.5.6 Case: Adjusting the Partial Clustering Key

Partial Cluster Key (PCK) is an index technology that uses min/max indexes to quickly scan base tables in column storage. Partial cluster key can specify multiple columns, but you are advised to specify no more than two columns. It can be used to accelerated gueries on large column-store tables.

Before Optimization

Create a column-store table **orders_no_pck** without partial clustering (PCK). The table is defined as follows:

Run the following SQL statement to query the execution plan of a point query:

```
EXPLAIN PERFORMANCE
SELECT * FROM orders_no_pck
WHERE o_orderkey = '13095143'
ORDER BY o_orderdate;
```

As shown in the following figure, the execution time is 48 ms. Check **Datanode Information**. It is found that the filter time is 19 ms and the CUNone ratio is 0.

```
| Gaussdb-> EXPLAIN PERFORMANCE | Gaussdb-> SELECT * FROM orders no pck | Gaussdb-> SELECT * FROM orders no pck | Gaussdb-> SELECT * FROM orders no pck | Gaussdb-> Children orders with the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of the content of t
```

After Optimization

The created column-store table **orders_pck** is defined as follows:

Use **ALTER TABLE** to set the **o_orderkey** field to **PCK**:

```
postgres=> ALTER TABLE orders_pck ADD PARTIAL CLUSTER KEY(o_orderkey);
ALTER TABLE
```

Run the following SQL statement to query the execution plan of the same point query SQL statement again:

```
EXPLAIN PERFORMANCE
SELECT * FROM orders_pck
WHERE o_orderkey = '13095143'
ORDER BY o_orderdate;
```

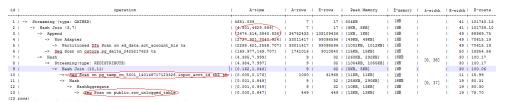
As shown in the following figure, the execution time is 5 ms. Check **Datanode Information**. It is found that the filter time is 0.5 ms and the CUNone ratio is 82. The higher the CUNone ratio, the higher performance that the PCK will bring.

13.5.7 Case: Adjusting the Table Storage Mode in a Medium Table

In GaussDB(DWS), row-store tables use the row execution engine, and column-store tables use the column execution engine. If both row-store table and column-store tables exist in a SQL statement, the system will automatically select the row execution engine. The performance of a column execution engine (except for the indexscan related operators) is much better than that of a row execution engine. Therefore, a column-store table is recommended. This is important for some medium result set dumping tables, and you need to select a proper table storage type.

Before Optimization

During the test at a site, if the following execution plan is performed, the customer expects that the performance can be improved and the result can be returned within 3s.



After Optimization

It is found that the row engine is used after analysis, because both the temporary plan table input_acct_id_tbl and the medium result dumping table row_unlogged_table use a row-store table.

After the two tables are changed into column-store tables, the system performance is improved and the result is returned by 1.6s.

id !	operation	A-time	A-rows	E-rows	Peak Memory	E-memory	A-width	E-width	E-costs
1	-> Row Adapter	1 1567.367	1 7	1 17	38KB	1	1	1 41	101758.52
2 1	-> Vector Streaming (type: GATHER)	1 1567.349	1 7	17	393KB	1	1	1 41	1 101758.52
3 1	-> Vector Hash Join (4,8)	[8.130,1529.101]	1 7	1 17	[2362KB, 2446KB]	16MB	1	1 41	1 101757.48
4 1	-> Vector Append	1 [542.823,1452.479]	1 5681770	108109436	[1KB, 1KB]	I 1MB	1	1 49	88969.75
5 1	-> Partitioned Dfs Scan on sd_data.act_account_his ta	[295.796,1195.830]	1 3940754	99098596	[861KB, 1012KB]	I 1MB	1	1 49	70615.19
61	-> Vector Adapter	1 [236.065,260.284]	1 1741016	9010840	[129KB, 129KB]	1 1MB	i i	1 50	1 18354.56
7 1	-> Seg Scan on catore.pg_delta_2425217623 ta	[152.595,168.048]	1 1741016	9010840	[15KB, 15KB]	1 1MB	1	1 50	18354.56
8 1	-> Vector Streaming(type: REDISTRIBUTE)	[7.727,12.981]	1 9	32	[1092KB, 1141KB]	I 1MB	[0, 40]	1 30	1 118.56
9 1	-> Vector Hash Join (10,11)	[0.132,4.955]	1 9	32	[2217KB, 2217KB]	1 16MB	1	1 30	1 118.45
10	-> CStore Scan on pg temp_cn 5001 140148155066112.input_acct_id_tbl_tb	1 [4.372,4.372]	1 999	31968	[207KB, 207KB]	1 136B	1	1 11	81.00
11	-> Vector Hash Aggregate	[[0.062, 0.209]	1 9	32	[2225KB, 2225KB]	1 16MB	[0, 35]	1 19	1 33.67
12	-> CStore Scan on public col unlogged table	1 [0.011,0.107]	1 449	449	[541KB, 598KB]	I 1MB	1	1 19	1 32.08

13.5.8 Case: Reconstructing Partition Tables

Partitioning refers to splitting what is logically one large table into smaller physical pieces based on specific schemes. The table based on the logic is called a partitioned table, and a physical piece is called a partition. Generally, partitioning is applied to tables that have obvious ranges. Partitions on such tables allow scanning on a small part of data, improving the query performance.

During query, partition pruning is used to minimize bottom-layer data scanning to narrow down the overall scope of scanning in a table. Partition pruning means that the optimizer can automatically extract partitions to be scanned based on the partition key specified in the **FROM** and **WHERE** statements. This avoids full table scanning, reduces the number of data blocks to be scanned, and improves performance.

Before Optimization

Create a non-partition table **orders_no_part**. The table definition is as follows:

Run the following SQL statement to query the execution plan of the non-partition table:

```
EXPLAIN PERFORMANCE
SELECT count(*) FROM orders_no_part WHERE
o_orderdate >= '1996-01-01 00:00:00'::timestamp(0);
```

As shown in the following figure, the execution time is 73 milliseconds, and the full table scanning time is 44 to 45 milliseconds.

After Optimization

Create a partitioned table **orders**. The table is defined as follows:

Run the SQL statement again to query the execution plan of the partitioned table. The execution time is 40 ms, in which the table scanning time is only 13 ms. The smaller the value of **Iterations**, the better the partition pruning effect.

```
EXPLAIN PERFORMANCE
SELECT count(*) FROM orders_no_part WHERE
o_orderdate >= '1996-01-01 00:00:00'::timestamp(0);
```

As shown in the following figure, the execution time is 40 milliseconds, and the table scanning time is only 13 milliseconds. A smaller **Iterations** value indicates a better partition pruning effect.

```
Saussch- EPU-AIM PERFORMANCE

Gaussch- StellClount() Selbo cders WHEE

10 | Operation | A-time | A-rows | E-rows | E-distinct | Peak Memory | E-memory | A-width | E-width | E-costs

1 | -> Row Adapter | 149.925 | 1 | 1 | 1868 | | 8 | 22382.64

2 | -> Vector Aggregate | 46.915 | 1 | 1 | 177KB | | 8 | 22382.64

3 | -> Vector Streaming (type: GATHER) | 48.873 | 3 | 80KB | | 8 | 22382.64

4 | --> Vector Aggregate | 46.967 | 3 | 3 | 80KB | | 8 | 22382.64

5 | -> Vector Aggregate | (40.967, 13.5.999) | 599865 | 584855 | (1386, 13868) | 1MB | 8 | 22382.64

5 | -> Vector Aggregate | (13.995, 13.5.999) | 599865 | 584855 | (1276, 1766) | 1MB | 8 | 277881.69

Prodicate Information (identified by plan id)

5 --Vector Agrithmed Store Scan on public.orders | 6-Partitioned Citore Scan on public.orders | Filter: (orders.o orderdate >= 1996-8018 88:80:xitimestamp(0) without time zone)

Partitions Selected by Static Province S.-7
```

13.5.9 Case: Adjusting the GUC Parameter best_agg_plan

Symptom

The t1 table is defined as follows:

```
create table t1(a int, b int, c int) distribute by hash(a);
```

Assume that the distribution column of the result set provided by the agg lower-layer operator is setA, and the group by column of the agg operation is setB, the agg operations can be performed in two scenarios in the stream framework.

Scenario 1: setA is a subset of setB.

In this scenario, the aggregation result of the lower-layer result set is the correct result, which can be directly used by the upper-layer operator. For details, see the following figure:

Scenario 2: setA is not a subset of setB.

In this scenario, the Stream execution framework is classified into the following three plans:

hashagg+gather(redistribute)+hashagg

redistribute+hashagg(+gather)

hashagg+redistribute+hashagg(+gather)

GaussDB(DWS) provides the guc parameter **best_agg_plan** to intervene the execution plan, and forces the plan to generate the corresponding execution plan. This parameter can be set to **0**, **1**, **2**, and **3**.

- When the value is set to 1, the first plan is forcibly generated.
- When the value is set to **2** and if the **group by** column can be redistributed, the second plan is forcibly generated. Otherwise, the first plan is generated.
- When the value is set to **3** and if the **group by** column can be redistributed, the third plan is generated. Otherwise, the first plan is generated.
- When the value is set to **0**, the query optimizer chooses the most optimal plan by the three preceding plans' evaluation cost.

Possible impacts are as follows:

```
set best_agg_plan to 1;
SET
explain select b,count(1) from t1 group by b;
         operation | E-rows | E-width | E-costs
id l
----+------+-----+-----+-----
1 | -> HashAggregate | 8 | 4 | 15.83
 2 | -> Streaming (type: GATHER) | 25 | 4 | 15.83
3 | -> HashAggregate | 25 | 4 | 14.33
4 | -> Seq Scan on t1 | 30 | 4 | 14.14
(4 rows)
set best_agg_plan to 2;
explain select b,count(1) from t1 group by b;
id | operation | E-rows | E-width | E-costs
1 | -> Streaming (type: GATHER) | 30 | 4 | 15.85
2 | -> HashAggregate | 30 | 4 | 14.60
3 | -> Streaming(type: REDISTRIBUTE) | 30 | 4 | 14.45
4 | -> Seq Scan on t1 | 30 | 4 | 14.14
(4 rows)
set best_agg_plan to 3;
explain select b,count(1) from t1 group by b;
id | operation | E-rows | E-width | E-costs
                              -----+----
1 | -> Streaming (type: GATHER) | 30 | 4 | 15.84
2 | -> HashAggregate | 30 | 4 | 14.59
    -> Streaming(type: REDISTRIBUTE) | 25 | 4 | 14.59

-> HashAggregate | 25 | 4 | 14.33

-> Seq Scan on t1 | 30 | 4 | 14.14
 3 |
 4 |
 5 I
(5 rows)
```

Summary

Generally, the optimizer chooses an optimal execution plan, but the cost estimation, especially that of the intermediate result set, has large deviations, which may result in large deviations in agg calculation. In this case, you need to use best_agg_plan to adjust the agg calculation model.

When the aggregation convergence ratio is very small, that is, the number of result sets does not become small obviously after the agg operation (5 times is a critical point), you can select the redistribute+hashagg or hashagg+redistribute+hashagg execution mode.

13.5.10 Case: Rewriting SQL Statements and Eliminating Prune Interference

A filter criterion that contains the expression of partition key cannot be used for pruning. As a result, the query statement scans almost all data in the partitioned table.

Before Optimization

t_ddw_f10_op_cust_asset_mon indicates the partitioned table. **year_mth** indicates the partition key. This field is an integer consisting of the **year** and **mth** values.

The following figure shows the tested SQL statements.

```
SELECT
count(1)

FROM t_ddw_f10_op_cust_asset_mon b1

WHERE b1.year_mth < substr('20200722',1,6)

AND b1.year_mth + 1 >= substr('20200722',1,6);
```

The test result shows that the table scan of the SQL statement takes 10 seconds. The execution plan of the SQL statement is as follows.

```
EXPLAIN (ANALYZE ON, VERBOSE ON)
SELECT
 count(1)
FROM t_ddw_f10_op_cust_asset_mon b1
WHERE b1.year_mth < substr('20200722',1 ,6 )
AND b1.year_mth + 1 >= cast(substr('20200722',1,6)) AS int);
                                              OUERY PLAN
                                       | A-time | A-rows | E-rows | E-
           operation
distinct | Peak Memory | E-memory | A-width | E-width | E-costs
| 10662.260 | 1 | 1 |
| | 8 | 593656.42

2 | -> Streaming (type: GATHER)

| 136KB | | 8 | 593656.42

3 | -> Aggregate

| [24KB, 24KB] | 1MB | 8 | 593646.42

4 | -> Partition Iterator
                                            | 10662.172 | 4 |
                                             | [9692.785, 10656.068] | 4 | 4
                                             | [8787.198, 9629.138] | 16384000 |
16384000 | 32752850 | | [32KB, 32KB] | 1MB | | 0 | 573175.88
               SQL Diagnostic Information
Partitioned table unprunable Qual
    table public.t ddw f10 op cust asset mon b1:
    left side of expression "((year_mth + 1) > 202008)" invokes function-call/type-conversion
         Predicate Information (identified by plan id)
 4 -- Partition Iterator
    Iterations: 6
 5 -- Partitioned Seq Scan on public.t_ddw_f10_op_cust_asset_mon b1
```

```
Filter: ((b1.year_mth < 202007::bigint) AND ((b1.year_mth + 1) >= 202007))
Rows Removed by Filter: 81920000
Partitions Selected by Static Prune: 1..6
```

After Optimization

After analyzing the execution plan of the statement and checking the SQL selfdiagnosis information in the execution plan, the following diagnosis information is found:

The filter criterion contains the expression (year_mth + 1) > 202008. A filter criterion that contains the expression of partition key cannot be used for pruning. As a result, the query statement scans almost all data in the partitioned table.

Compared with the original SQL statement, the expression (year_mth + 1) > 202008 is derived from the expression b1.year_mth + 1 > substr('20200822',1,6). Based on the diagnosis information, the SQL statement is modified as follows.

```
SELECT
count(1)
FROM t_ddw_f10_op_cust_asset_mon b1
WHERE b1.year_mth <= substr('20200822',1 ,6 )
AND b1.year_mth > cast(substr('20200822',1 ,6 ) AS int) - 1;
```

After the modification, the SQL statement execution information is as follows. The alarm indicating that the pruning is not performed is cleared. After the pruning, the score of the partition to be scanned is 1, and the execution time is shortened from 10 seconds to 3 seconds.

```
EXPLAIN (analyze ON, verbose ON)
SELECT
  count(1)
FROM t_ddw_f10_op_cust_asset_mon b1
WHERE b1.year_mth < substr('20200722',1,6)
AND b1.year_mth >= cast(substr('20200722',1,6) AS int) - 1;
                                                     QUERY PLAN
_____
id I
                     operation
                                                      A-time | A-rows | E-rows | E-
distinct | Peak Memory | E-memory | A-width | E-width | E-costs
                                               3009.796
 1 | -> Aggregate
                                                                | 1| 1|
                  8 | 501541.70
32KB
 2 | -> Streaming (type: GATHER)
                                                       3009.718
                                                                            4 |
      |136KB | | |
                                8 | 501541.70
 3 |
     -> Aggregate
                                                   | [2675.509, 3003.298] |
                                                                          4 |
      | [24KB, 24KB] | 1MB | 8 | 501531.70
                                                   | [1820.725, 2053.836] | 16384000 |
 4 |
        -> Partition Iterator
16380697 | | [16KB, 16KB] | 1MB | | 0 | 491293.75
          -> Partitioned Seq Scan on public.t_ddw_f10_op_cust_asset_mon b1 | [1420.972, 1590.083] |
 5 |
16384000 | 16380697 |
                        | [16KB, 16KB] | 1MB |
                                                     0 | 491293.75
        Predicate Information (identified by plan id)
 4 -- Partition Iterator
    Iterations: 1
 5 -- Partitioned Seq Scan on public.t_ddw_f10_op_cust_asset_mon b1
```

```
Filter: ((b1.year_mth < 202007::bigint) AND (b1.year_mth >= 202006)) Partitions Selected by Static Prune: 6
```

13.5.11 Case: Rewriting SQL Statements and Deleting inclause

Before Optimization

in-clause/any-clause is a common SQL statement constraint. Sometimes, the clause following **in** or **any** is a constant. For example:

```
select
count(1)
from calc_empfyc_c1_result_tmp_t1
where ls_pid_cusr1 in ('20120405', '20130405');

Or

select
count(1)
from calc_empfyc_c1_result_tmp_t1
where ls_pid_cusr1 in any('20120405', '20130405');
```

Some special usages are as follows:

```
SELECT

ls_pid_cusr1,COALESCE(max(round((current_date-bthdate)/365)),0)

FROM calc_empfyc_c1_result_tmp_t1 t1,p10_md_tmp_t2 t2

WHERE t1.ls_pid_cusr1 = any(values(id),(id15))

GROUP BY ls_pid_cusr1;
```

Where **id** and **id15** are columns of p10_md_tmp_t2. ls_pid_cusr1 = any(values(id), (id15)) equals t1. ls_pid_cusr1 = id or t1. ls_pid_cusr1 = id15.

Therefore, join-condition is essentially an inequality, and nestloop must be used for this join operation. The execution plan is as follows:

```
Streaming (type: GATHER) (cost=1641432384.14..1641432583.98 rows=3840 width=49)
Node/s: All datamodes

> Insert on channel.calc_empfyc_cl_result_age_tmp (cost=164142380.14..1641423283.88 rows=3840 width=49)

-> HashAggregate (cost=164123280.14..16414323283.38 rows=3840 width=25)

Output: ti.ls_pid_cust; (COALESCETERS(Imagk(cound((('2017-03-25 00:00:00':timestamp without time zone - t2.bphdats) / 365::double precision))::numeric)

Group by May: ti.ls_pid_cust; (magk(cound((('2017-03-25 00:00:00':timestamp without time zone - t2.bphdats) / 365::double precision))::numeric, 0)))

Distribute May: ti.ls_pid_cust; (COALESCETERS, (magk(cound((('2017-03-25 00:00:00':timestamp without time zone - t2.bphdats) / 365::double precision))::numeric, 0)))

Span on: All datamodes

-> HashAggregate (cost=020714660.07..22071460.65 rows=3968 width=25)

Output: ti.ls_pid_cust, magk(cound((('2017-03-25 00:00:00':timestamp without time zone - t2.bphdats) / 365::double precision))::numeric, 0))

Group By May: ti.ls_pid_cust; (ask(cound((('2017-03-25 00:00:00'):timestamp without time zone - t2.bphdats) / 365::double precision))::numeric, 0))

Group By May: ti.ls_pid_cust; (ask(cound)):timestamp without time zone - t2.bphdats) / 365::double precision))::numeric, 0))

Group By May: ti.ls_pid_cust; (ask(cound)):timestamp without time zone - t2.bphdats) / 365::double precision))::numeric, 0))

Group By May: ti.ls_pid_cust; (ask(cound)):timestamp without time zone - t2.bphdats) / 365::double precision))::numeric, 0))

Group By May: ti.ls_pid_cust; (ask(cound)):timestamp without time zone - t2.bphdats) / 365::double precision))::numeric, 0))

Group By May: ti.ls_pid_cust; (ask(cound)):timestamp without time zone - t2.bphdats) / 365::double precision))::numeric, 0))

Group By May: ti.ls_pid_cust; (ask(cound)):timestamp without time zone - t2.bphdats) / 365::double precision))::numeric, 0))

Group By May: ti.ls_pid_cust; (ask(cust)):timestamp without time zone - t2.bphdats / 250:double precision))::numeric, 0))

Group By May: ti.ls_pid_cust; (
```

After Optimization

The test result shows that both result sets are too large. As a result, nestloop is time-consuming with more than one hour to return results. Therefore, the key to performance optimization is to eliminate nestloop, using more efficient hashjoin. From the perspective of semantic equivalence, the SQL statements can be written as follows:

```
select
ls_pid_cusr1,COALESCE(max(round(ym/365)),0)
from
```

Note: Use **UNION ALL** instead of **UNION** if possible. **UNION** eliminates duplicate rows while merging two result sets but **UNION ALL** merges the two result sets without deduplication. Therefore, replace **UNION** with **UNION ALL** if you are sure that the two result sets do not contain duplicate rows based on the service logic.

The optimized SQL queries consist of two equivalent join subqueries, and each subquery can be used for hashjoin in this scenario. The optimized execution plan is as follows:

id operation	A-time	A-rows	E-rows	Peak Memory	E-memory	A-width
1 -> Streaming (type: GATHER)	6737.281	0	192	292KB		1
2 Insert on channel.calc_empfyc_cl_result_age_tmp	[4665.024,4990.666]	0	192	[1108KB, 1108KB]	1MB	1
3 -> HashAggregate	[4664.996,4990.641]	0	192	[[12KB, 12KB]	16MB	1
4 -> Streaming(type: REDISTRIBUTE)	[4664.991,4990.637]	0	3392	[[2090KB, 2090KB]	1MB	1
5	[3416.939,4958.348]		3392	[[14KB, 14KB]	16MB	1
6 -> Append	[3416.936,4958.340]	0	4011	[[1KB, 1KB]	1MB	T
7 -> Hash Join (8,9)	[2011.226,3080.697]	0	3947	[[6KB, 6KB]	1MB	1
8 -> Seq Scan on channel.pl0_md_tmp_t2 t2	[803.782,1238.984]	443525717	443523360	[12KB, 12KB]	1MB	1
9 -> Hesh	[4.357,328.979]	252608	252608	[482KB, 482KB]	16MB	[35, 39]
1D > Streaming(type: BROADCAST)	[2.345,326.320]	252608	252608	[[2090KB, 2090KB]	1MB	1
11 -> Seq Scan on channel.calc empfyc cl result tmp tl tl	[0.011,0.030]	3947	3947	[11KB, 11KB]	1MB	1
12 -> Hash Join (13,14)	[1376.258,2066.110]	0	64	[[5KB, 5KB]	1MB	1
13 -> Seq Scan on channel.p10 md tmp t2 t2	[777.552,1388.499]	443525717	443523360	[12KB, 12KB]	1MB	1
14 -> Hash	[2.812,4.217]	252608	252608	[482KB, 482KB]	16MB	[23, 27]
15 -> Streaming(type: BROADCAST)	[1.276, 1.868]	252608	252608	[2090KB, 2090KB] [1MB	i .
16	[0.010,0.033]	3947	3947	[11KB, 11KB]	1MB	1

Before the optimization, no result is returned for more than 1 hour. After the optimization, the result is returned within 7s.

13.5.12 Case: Setting Partial Cluster Keys

You can add **PARTIAL CLUSTER KEY**(*column_name*[,...]) to the definition of a column-store table to set one or more columns of this table as partial cluster keys. In this way, each 70 CUs (4.2 million rows) will be sorted based on the cluster keys by default during data import and the value range is narrowed down for each of the new 70 CUs. If the **where** condition in the query statement contains these columns, the filtering performance will be improved.

Before Optimization

```
The partial cluster key is not used. The table is defined as follows:
```

```
CREATE TABLE lineitem
(
L_ORDERKEY BIGINT NOT NULL
, L_PARTKEY BIGINT NOT NULL
, L_SUPPKEY BIGINT NOT NULL
, L_LINENUMBER BIGINT NOT NULL
, L_QUANTITY DECIMAL(15,2) NOT NULL
, L_EXTENDEDPRICE DECIMAL(15,2) NOT NULL
, L_DISCOUNT DECIMAL(15,2) NOT NULL
, L_TAX DECIMAL(15,2) NOT NULL
, L_RETURNFLAG CHAR(1) NOT NULL
, L_LINESTATUS CHAR(1) NOT NULL
```

```
, L_SHIPDATE DATE NOT NULL
, L_COMMITDATE DATE NOT NULL
, L_RECEIPTDATE DATE NOT NULL
, L_SHIPINSTRUCT CHAR(25) NOT NULL
, L_SHIPMODE CHAR(10) NOT NULL
, L_COMMENT
               VARCHAR(44) NOT NULL
with (orientation = column)
distribute by hash(L_ORDERKEY);
sum(l_extendedprice * l_discount) as revenue
from
lineitem
where
l_shipdate >= '1994-01-01'::date
and l_shipdate < '1994-01-01'::date + interval '1 year'
and l_discount between 0.06 - 0.01 and 0.06 + 0.01
and l_quantity < 24;
```

After the data is imported, perform the query and check the execution time.

Figure 13-23 Partial cluster keys not used

id	operation	A-time	A-	rows	E-row	5	Peak Mo	emory	A-width	E-widt	th	E-costs	
1 t D	*******	1052 150				+	1000			+		205002.00	
	Adapter Vector Aggregate	1653.156 1653.146		1		<u> </u>	12KB 184KB					205803.90 205803.90	
	> Vector Aggregate > Vector Streaming (type: GATHER)	1653.070		1			174KB					205803.90	
4 1	-> Vector Streaming (type, GATHER)	[1481.497,1481.497]		1			[225KB,	225KB1				205803.90	
5 1	-> CStore Scan on public.lineitem			14160	11148							205246.40	
(5 rows)	- cocoro ocan on pascicicinateon	[11031001,11031001]	1 -	11100	11110	- 1	(/ JEND)	70210				200210110	

Figure 13-24 CU loading without partial cluster keys

```
5 --CStore Scan on public.lineitem
datanodel (actual time=40.623..1405.004 rows=114160 loops=1)
datanodel (RoughCheck CU: CUNone: 0, CUSome: 101)
datanodel (LLVM Optimized)
datanodel (Buffers: shared hit=18385 read=23)
datanodel (CPU: ex c/r=31917, ex cyc=3643646206, inc cyc=3643646206)
```

After Optimization

In the **where** condition, both the **l_shipdate** and **l_quantity** columns have a few distinct values, and their values can be used for min/max filtering. Therefore, modify the table definition as follows:

```
CREATE TABLE lineitem
L_ORDERKEY BIGINT NOT NULL
, L_PARTKEY
            BIGINT NOT NULL
            BIGINT NOT NULL
, L_SUPPKEY
, L_LINENUMBER BIGINT NOT NULL
, L_QUANTITY DECIMAL(15,2) NOT NULL
, L_EXTENDEDPRICE DECIMAL(15,2) NOT NULL
, L_DISCOUNT DECIMAL(15,2) NOT NULL
, L_TAX
          DECIMAL(15,2) NOT NULL
, L_RETURNFLAG CHAR(1) NOT NULL
, L_LINESTATUS CHAR(1) NOT NULL
, L_SHIPDATE DATE NOT NULL
, L_COMMITDATE DATE NOT NULL
, L RECEIPTDATE DATE NOT NULL
, L_SHIPINSTRUCT CHAR(25) NOT NULL
, L_SHIPMODE
              CHAR(10) NOT NULL
, L_COMMENT
               VARCHAR(44) NOT NULL
, partial cluster key(l_shipdate, l_quantity)
with (orientation = column)
distribute by hash(L_ORDERKEY);
```

Import the data again, perform the query, and check the execution time.

Figure 13-25 Partial cluster keys used

id	operation	ļ ļ	A-time	ļ	A-rows	E-ro	ows	ı	Peak Me	mory	A-width	E-width	E-costs
1 -	> Row Adapter	459	9.539	Ī	1		1	12	2KB			44	205693.85
2	-> Vector Aggregate	459	9.528		1		1	18	84KB		ĺ	44	205693.85
3 j	-> Vector Streaming (type: GATHER)	459	9.452		1		1	1.	74KB		i	44	205693.85
4	-> Vector Aggregate	j [28	85.177,285.177]		1		1 j	[2	225KB,	225KB]	i	44	205693.79
5 j	-> CStore Scan on public.lineitem	j [24	49.757,249.757]		114160	894	475 j	U	792KB,	792KB]	i	12	205246.40
(5 rows													

Figure 13-26 CU loading with partial cluster keys

```
5 --CStore Scan on public.lineitem
datanodel (actual time=23.017..249.757 rows=114160 loops=1)
datanodel (RoughCheck CU: CUNone: 84, CUSome: 17)
datanodel (LLVM Optimized)
datanodel (Buffers: shared hit=2853 read=23)
datanodel (CPU: ex c/r=5673, ex cyc=647656146, inc cyc=647656146)
```

After partial cluster keys are used, the execution time of **5-- CStore Scan on public.lineitem** decreases by 1.2s because 84 CUs are filtered out.

Optimization

- Select partial cluster keys.
 - The following data types support cluster keys: character varying(n), varchar(n), character(n), char(n), text, nvarchar2, timestamp with time zone, timestamp without time zone, date, time without time zone, and time with time zone.
 - Smaller number of distinct values in a partial cluster key generates higher filtering performance.
 - Columns that can filter out larger amount of data is preferentially selected as partial cluster keys.
 - If multiple columns are selected as partial cluster keys, the columns are used in sequence to sort data. You are advised to select a maximum of three columns.
- Modify parameters to reduce the impact of partial cluster keys on the import performance.

After partial cluster keys are used, data will be sorted when they are imported, affecting the import performance. If all the data can be sorted in the memory, the keys have little impact on import. If some data cannot be sorted in the memory and is written into a temporary file for sorting, the import performance will be greatly affected.

The memory used for sorting is specified by the **psort_work_mem** parameter. You can set it to a larger value so that the sorting has less impact on the import performance.

The volume of data to be sorted is specified by the PARTIAL_CLUSTER_ROWS parameter of the table. Decreasing the value of this parameter reduces the amount of data to be sorted at a time. PARTIAL_CLUSTER_ROWS is usually used along with the MAX_BATCHROW parameter. The value of PARTIAL_CLUSTER_ROWS must be an integer multiple of the MAX_BATCHROW value. MAX_BATCHROW specifies the maximum number of rows in a CU.

13.5.13 Case: Converting from NOT IN to NOT EXISTS

nestloop anti join must be used to implement **NOT IN**, while you can use **Hash anti join** to implement **NOT EXISTS**. If no **NULL** value exists in the **JOIN** column, **NOT IN** is equivalent to **NOT EXISTS**. Therefore, if you are sure that no **NULL** value exists, you can convert **NOT IN** to **NOT EXISTS** to generate **hash joins** and to improve the query performance.

Before Optimization

Create two base tables t1 and t2.

```
CREATE TABLE t1(a int, b int, c int not null) WITH(orientation=row); CREATE TABLE t2(a int, b int, c int not null) WITH(orientation=row);
```

Run the following SQL statement to query the **NOT IN** execution plan:

EXPLAIN VERBOSE SELECT * FROM t1 WHERE t1.c NOT IN (SELECT t2.c FROM t2);

The following figure shows the statement output.

```
id |
                                             | E-rows | E-distinct | E-width | E-costs
                    operation
     -> Streaming (type: GATHER)
                                                                           12
                                                                                78.98
        -> Nested Loop Anti Join (3, 4)
                                                    6
                                                                           12
              Seq Scan on public.tl
                                                   60
                                                                           12
                                                                                18.18
              Materialize
                                                  360
              -> Streaming(type: BROADCAST)
                                                  360
                                                                                30.45
                 -> Seq Scan on public.t2
                                                   60
            Predicate Information (identified by plan id)
 2 --Nested Loop Anti Join (3, 4)
       Join Filter: ((tl.c = t2.c) OR (tl.c IS NULL) OR (t2.c IS NULL))
```

According to the returned result, nest loops are used. As the OR operation result of NULL and any value is NULL,

t1.c NOT IN (SELECT t2.c FROM t2)

the preceding condition expression is equivalent to:

t1.c <> ANY(t2.c) AND t1.c IS NOT NULL AND ANY(t2.c) IS NOT NULL

After Optimization

The query can be modified as follows:

```
SELECT * FROM t1 WHERE NOT EXISTS (SELECT * FROM t2 WHERE t2.c = t1.c);
```

Run the following statement to query the execution plan of **NOT EXISTS**:

EXPLAIN VERBOSE SELECT * FROM t1 WHERE NOT EXISTS (SELECT 1 FROM t2 WHERE t2.c = t1.c);

```
QUERY PLAN
 id |
                                                                          | E-rows | E-distinct | E-width | E-costs
                                 operation
                                                                                                                          | 54.56
| 40.56
| 20.12
| 18.18
| 20.12
| 20.12
| 18.18
                                                                                                                     12 |
12 |
12 |
12 |
12 |
4 |
                                                                                    6
         -> Streaming (type: GATHER)
             -> Hash Anti Join (3, 5)
-> Streaming(type: REDISTRIBUTE)
-> Seq Scan on public.tl
  3 | 4 | 5 |
                                                                                  60
60
59
                                                                                          10
                  -> Hash
                                                                                          10
                       -> Streaming(type: REDISTRIBUTE)
-> Seq Scan on public.t2
                                                                                   60
                                                                                   60
Predicate Information (identified by plan id)
   2 -- Hash Anti Join (3, 5)
           Hash Cond: (tl.c = t2.c)
```

14 GaussDB(DWS) System Catalogs and Views

14.1 Overview of System Catalogs and System Views

System catalogs are used by GaussDB(DWS) to store structure metadata. They are a core component the GaussDB(DWS) database system and provide control information for the database system. These system catalogs contain cluster installation information and information about various queries and processes in GaussDB(DWS). You can collect information about the database by querying the system catalog.

System views provide ways to query system catalogs and internal database status. If some columns in one or more tables in a database are frequently searched for, an administrator can define a view for these columns, and then users can directly access these columns in the view without entering search criteria. A view is different from a basic table. It is only a virtual object rather than a physical one. A database only stores the definition of a view and does not store its data. The data is still stored in the original base table. If data in the base table changes, the data in the view changes accordingly. In this sense, a view is like a window through which users can know their interested data and data changes in the database. A view is triggered every time it is referenced.

In separation of duty, non-administrators have no permission to view system catalogs and views. In other scenarios, system catalogs and views are either visible only to administrators or visible to all users. Some of the following system catalogs and views have marked the need of administrator permissions. They are accessible only to administrators.

NOTICE

- Do not add, delete, or modify system catalogs or system views. Manual modification or damage to system catalogs or system views may cause system information inconsistency, system control exceptions, or even cluster unavailability.
- System catalogs do not support toast and cannot be stored across pages. The size of a page is 8 KB, and the length of each field in the system catalog must be less than 8 KB.

14.2 System Catalogs

14.2.1 GS_BLOCKLIST_QUERY

GS_BLOCKLIST_QUERY records job blocklist and exception information. This table uses unique_sql_id as the unique index to collect statistics on job exception information and record blocklist information. You can associate
GS_BLOCKLIST_QUERY with GS_WLM_SESSION_INFO to obtain the query column and execution information of a job.

NOTICE

- The schema of the **GS_BLOCKLIST_QUERY** system catalog is **dbms_om**.
- The **GS_BLOCKLIST_QUERY** system catalog can be queried only in the **postgres** database. If it is queried in other databases, an error is reported.
- The **GS_BLOCKLIST_QUERY** system catalog contains unique indexes, which are distributed on DNs in hash mode. The distributed column is **unique sql id**.
- Generally, constant values are ignored during Unique SQL ID calculation in DML statements. However, constant values cannot be ignored in DDL, DCL, and parameter setting statements. A unique_sql_id may correspond to one or more queries.

GaussDB(DWS) also provides the **GS_BLOCKLIST_QUERY** view for querying job blocklist and exception information. This view can directly display the **query** column. This view depends on **GS_WLM_SESSION_INFO**. If the **GS_WLM_SESSION_INFO** table is large, the query may take a long time.

Table 14-1 GS BLOCKLIST QUERY columns

Name	Туре	Referenc e	Description
unique_sql_id	Bigint	N/A	Query ID generated based on the query parsing tree.
block_list	Boolean	N/A	Check whether a job is in the blocklist.

Name	Туре	Referenc e	Description
except_num	Integer	N/A	Query the number of job exceptions.
except_time	Timestamp	N/A	Query the time when the last job exception occurred.

14.2.2 GS BLOCKLIST SQL

GS_BLOCKLIST_SQL records job blocklist and exception information. This table uses sql_hash as the unique index to collect statistics on job exception information and record blocklist information. You can associate GS_BLOCKLIST_SQL with GS_WLM_SESSION_INFO to obtain the query column and execution information of a job.

GaussDB(DWS) also provides the **GS_BLOCKLIST_SQL** view for querying job blocklist and exception information. This view can directly display the **query** column. This view depends on **GS_WLM_SESSION_INFO**. If the **GS_WLM_SESSION_INFO** table is large, the query may take a long time.

This system catalog is supported only by clusters of version 9.1.0.200 or later.

NOTICE

- The schema of the **GS_BLOCKLIST_SQL** system catalog is **dbms_om**.
- The **GS_BLOCKLIST_SQL** system catalog can be queried only in the **postgres** database. If it is queried in other databases, an error is reported.
- The **GS_BLOCKLIST_SQL** system catalog contains unique indexes, which are distributed on DNs in hash mode. The distributed column is **sql_hash**.
- Generally, constant values are ignored during sql_hash calculation in DML statements. However, constant values cannot be ignored in DDL, DCL, and parameter setting statements. A sql_hash may correspond to one or more queries.

Table 14-2 GS_BLOCKLIST_SQL columns

Name	Туре	Referenc e	Description
sql_hash	Text	N/A	sql_hash generated based on the query parsing tree.
block_list	Boolean	N/A	Check whether a job is in the blocklist.
except_num	Integer	N/A	Query the number of job exceptions.

Name	Туре	Referenc e	Description
except_time	Timestamp	N/A	Query the time when the last job exception occurred.

14.2.3 GS_OBSSCANINFO

GS_OBSSCANINFO defines the OBS runtime information scanned in cluster acceleration scenarios. Each record corresponds to a piece of runtime information of a foreign table on OBS in a query.

Table 14-3 GS_OBSSCANINFO columns

Name	Туре	Reference	Description
query_id	Bigint	-	Specifies a query ID.
user_id	Text	-	Specifies a database user who performs queries.
table_name	Text	-	Specifies the name of a foreign table on OBS.
file_type	Text	-	Specifies the format of files storing the underlying data.
time_stamp	time_st am	-	Specifies the scanning start time.
actual_time	double	-	Specifies the scanning execution time in seconds.
file_scanned	Bigint	-	Specifies the number of files scanned.
data_size	double	-	Specifies the size of data scanned in bytes.
billing_info	Text	-	Specifies the reserved fields.

14.2.4 GS_RESPOOL_RESOURCE_HISTORY

The **GS_RESPOOL_RESOURCE_HISTORY** table records the historical monitoring information about a resource pool on both CNs and DNs.

Table 14-4 GS_RESPOOL_RESOURCE_HISTORY columns

Name	Туре	Description
Timestamp	Timestamp	Time when resource pool monitoring information is persistently stored
nodegroup	Name	Name of the logical cluster of the resource pool. The default value is installation .
rpname	Name	Resource pool name
cgroup	Name	Name of the Cgroup associated with the resource pool
ref_count	Int	Number of jobs referenced by the resource pool. The number is counted regardless of whether the jobs are controlled by the resource pool. This parameter is valid only on CNs.
fast_run	Int	Number of running jobs in the fast lane of the resource pool. This parameter is valid only on CNs.
fast_wait	Int	Number of jobs queued in the fast lane of the resource pool. This parameter is valid only on CNs.
fast_limit	Int	Limit on the number of concurrent jobs in the fast lane in a resource pool. This parameter is valid only on CNs.
slow_run	Int	Number of running jobs in the slow lane of the resource pool. This parameter is valid only on CNs.
slow_wait	Int	Number of jobs queued in the slow lane of the resource pool. This parameter is valid only on CNs.
slow_limit	Int	Limit on the number of concurrent jobs in the slow lane in a resource pool. This parameter is valid only on CNs.
used_cpu	Double	Average number of CPUs used by the resource pool in a 5s monitoring period. The value is accurate to two decimal places.
		 On a DN, it indicates the number of CPUs used by the resource pool on the current DN.
		On a CN, it indicates the total CPU usage of resource pools on all DNs.

Name	Туре	Description
cpu_limit	Int	It indicates the upper limit of available CPUs for resource pools. If the CPU share is limited, this parameter indicates the available CPUs for GaussDB(DWS). If the CPU limit is specified, this parameter indicates the available CPUs for associated Cgroups.
		On a DN, it indicates the upper limit of available CPUs for the resource pool on the current DN.
		On a CN, it indicates the total upper limit of available CPUs for resource pools on all DNs.
used_mem	Int	Memory used by the resource pool, in MB.
		On a DN, it indicates the memory usage of the resource pool on the current DN.
		On a CN, it indicates the total memory usage of resource pools on all DNs.
estimate_me m	Int	Estimated memory used by the jobs running in the resource pools on the current CN. This parameter is valid only on CNs.
mem_limit	Int	Upper limit of available memory for the resource pool (unit: MB).
		On a DN, it indicates the upper limit of available memory for the resource pool on the current DN.
		On a CN, it indicates the total upper limit of available memory for resource pools on all DNs.
read_kbytes	Bigint	Number of logical read bytes in the resource pool within a 5s monitoring period (unit: KB).
		 On a DN, it indicates the number of logical read bytes in the resource pool on the current DN.
		On a CN, it indicates the total logical read bytes of resource pools on all DNs.
write_kbytes	Bigint	Number of logical write bytes in the resource pool within a 5s monitoring period (unit: KB).
		On a DN, it indicates the number of logical write bytes in the resource pool on the current DN.
		On a CN, it indicates the total logical write bytes of resource pools on all DNs.

Name	Туре	Description
read_counts	Bigint	Number of logical reads in the resource pool within a 5s monitoring period.
		On a DN, it indicates the number of logical reads in the resource pool on the current DN.
		 On a CN, it indicates the total number of logical reads in resource pools on all DNs.
write_counts	Bigint	Number of logical writes in the resource pool within a 5s monitoring period.
		 On a DN, it indicates the number of logical writes in the resource pool on the current DN.
		On a CN, it indicates the total number of logical writes in resource pools on all DNs.
read_speed	Double	Average rate of logical reads of the resource pool in a 5s monitoring period, in KB/s.
		On a DN, it indicates the logical read rate of the resource pool on the current DN.
		On a CN, it indicates the overall logical read rate of resource pools on all DNs.
write_speed	Double	Average rate of logical writes of resource pools in a 5s monitoring period, in KB/s.
		 On a DN, it indicates the logical write rate of the resource pool on the current DN.
		On a CN, it indicates the overall logical write rate of resource pools on all DNs.
send_speed	Double	Average network sending rate of the resource pool in a 5-second monitoring period, in KB/s.
		 On a DN, it indicates the network sending rate of the resource pool on the current DN.
		On a CN, it indicates that the cumulative sum of the network sending rates of the resource pool on all DNs.
recv_speed	Double	Average network receiving rate of the resource pool in a 5-second monitoring period, in KB/s.
		On a DN, it indicates the network receiving rate of the resource pool on the current DN.
		On a CN, it indicates that the cumulative sum of the network receiving rates of the resource pool on all DNs.

14.2.5 GS_WLM_INSTANCE_HISTORY

The **GS_WLM_INSTANCE_HISTORY** system catalog stores information about resource usage related to CN or DN instances. Each record in the system table indicates the resource usage of an instance at a specific time point, including the memory, number of CPU cores, disk I/O, physical I/O of the process, and logical I/O of the process.

Table 14-5 GS_WLM_INSTANCE_HISTORY column

Column	Туре	Description
instancena me	Text	Instance name
Timestamp	Timestamp with time zone	Timestamp
used_cpu	Int	CPU usage of an instance
free_mem	Int	Unused memory of an instance (unit: MB)
used_mem	Int	Used memory of an instance (unit: MB)
io_await	Real	Specifies the io_wait value (average value within 10 seconds) of the disk used by an instance.
io_util	Real	Specifies the io_util value (average value within 10 seconds) of the disk used by an instance.
disk_read	Real	Specifies the disk read rate (average value within 10 seconds) of an instance (unit: KB/s).
disk_write	Real	The disk write rate (average value within 10 seconds) of an instance (unit: KB/s).
process_rea d	Bigint	Specifies the read rate (excluding the number of bytes read from the disk pagecache) of the corresponding instance process that reads data from a disk. (Unit: KB/s)
process_wri te	Bigint	Specifies the write rate (excluding the number of bytes written to the disk pagecache) of the corresponding instance process that writes data to a disk within 10 seconds. (Unit: KB/s)
logical_read	Bigint	CN instance: N/A DN instance: Specifies the logical read byte rate of the instance in the statistical interval (10 seconds). (Unit: KB/s)

Column	Туре	Description
logical_writ e	Bigint	CN instance: N/A DN instance: Specifies the logical write byte rate of the instance within the statistical interval (10 seconds). (Unit: KB/s)
read_counts	Bigint	CN instance: N/A DN instance: Specifies the total number of logical read operations of the instance in the statistical interval (10 seconds).
write_count s	Bigint	CN instance: N/A DN instance: Specifies the total number of logical write operations of the instance in the statistical interval (10 seconds).

14.2.6 GS_WLM_OPERATOR_INFO

GS_WLM_OPERATOR_INFO records operators of completed jobs. The data is dumped from the kernel to a system catalog. If the GUC parameter **enable_resource_record** is set to **on**, the system periodically imports records in **GS_WLM_OPERATOR_HISTORY** to this system catalog. You are not advised to enable this function because it occupies storage space and affects performance. You are advised to disable it after performance locating and monitoring tasks are complete.

NOTICE

- The schema of the **GS_WLM_OPERATOR_INFO** system table is **dbms_om**.
- The **GS_WLM_OPERATOR_INFO** system catalog can be queried only in the **postgres** database. If it is queried in other databases, an error is reported.

Table 14-6 GS_WLM_OPERATOR_INFO columns

Column	Туре	Description
nodename	Text	Name of the CN where the statement is executed
queryid	Bigint	Internal query_id used for statement execution
pid	Bigint	Backend thread ID
plan_node_id	Integer	plan_node_id of the execution plan of a query
plan_node_nam e	Text	Name of the operator corresponding to plan_node_id

Column	Туре	Description
start_time	Timestamp with time zone	Time when an operator starts to process the first data record
duration	Bigint	Total execution time of an operator. The unit is ms.
query_dop	Integer	Degree of parallelism (DOP) of the current operator
estimated_rows	Bigint	Number of rows estimated by the optimizer
tuple_processed	Bigint	Number of elements returned by the current operator
min_peak_mem ory	Integer	Minimum peak memory used by the current operator on all DNs. The unit is MB.
max_peak_me mory	Integer	Maximum peak memory used by the current operator on all DNs. The unit is MB.
average_peak_ memory	Integer	Average peak memory used by the current operator on all DNs. The unit is MB.
memory_skew_ percent	Integer	Memory usage skew of the current operator among DNs
min_spill_size	Integer	Minimum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
max_spill_size	Integer	Maximum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
average_spill_si ze	Integer	Average spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
spill_skew_perc ent	Integer	DN spill skew when a spill occurs
min_cpu_time	Bigint	Minimum execution time of the operator on all DNs. The unit is ms.
max_cpu_time	Bigint	Maximum execution time of the operator on all DNs. The unit is ms.
total_cpu_time	Bigint	Total execution time of the operator on all DNs. The unit is ms.
cpu_skew_perce nt	Integer	Skew of the execution time among DNs.

Column	Туре	Description
warning	Text	Warning. The following warnings are displayed:
		Sort/SetOp/HashAgg/HashJoin spill
		2. Spill file size large than 256MB
		3. Broadcast size large than 100MB
		4. Early spill
		5. Spill times is greater than 3
		6. Spill on memory adaptive
		7. Hash table conflict

14.2.7 GS_WLM_SESSION_INFO

GS_WLM_SESSION_INFO records load management information about a completed job executed on all CNs. The data is dumped from the kernel to a system catalog. If the GUC parameter **enable_resource_record** is set to **on**, the system periodically imports records in **GS_WLM_SESSION_HISTORY** to this system catalog. You are not advised to enable this function because it occupies storage space and affects performance. You are advised to disable it after performance locating and monitoring tasks are complete.

NOTICE

- The schema of the **GS_WLM_SESSION_INFO** system table is **dbms_om**.
- The **GS_WLM_SESSION_INFO** system catalog can be queried only in the **postgres** database. If it is queried in other databases, an error is reported.

Table 14-158 lists the columns in the GS_WLM_SESSION_INFO system catalog.

Table 14-7 GS_WLM_SESSION_HISTORY columns

Column	Туре	Description
datid	OID	OID of the database this backend is connected to
dbname	Text	Name of the database the backend is connected to
schemaname	Text	Schema name
nodename	Text	Name of the CN where the statement is run
username	Text	User name used for connecting to the backend
application_na me	Text	Name of the application that is connected to the backend

Column	Туре	Description
client_addr	inet	IP address of the client connected to this backend. If this column is null, it indicates either that the client is connected via a Unix socket on the server machine or that this is an internal process such as autovacuum.
client_hostnam e	Text	Host name of the connected client, as reported by a reverse DNS lookup of client_addr . This column will only be non-null for IP connections, and only when log_hostname is enabled.
client_port	Integer	TCP port number that the client uses for communication with this backend, or -1 if a Unix socket is used
query_band	Text	Job type, which can be set using the GUC parameter query_band and is a null string by default.
block_time	Bigint	Duration that a statement is blocked before being executed, including the statement parsing and optimization duration. The unit is ms.
start_time	Timestamp with time zone	Time when the statement starts to be run
finish_time	Timestamp with time zone	Time when the statement execution ends
duration	Bigint	Execution time of a statement. The unit is ms.
estimate_total_ time	Bigint	Estimated execution time of a statement. The unit is ms.
status	Text	Final statement execution status. Its value can be finished (normal) or aborted (abnormal). The statement status here is the execution status of the database server. If the statement is successfully executed on the database server but an error is reported in the result set, the statement status is finished .
abort_info	Text	Exception information displayed if the final statement execution status is aborted .
resource_pool	Text	Resource pool used by the user
control_group	Text	Cgroup used by the statement

Column	Туре	Description
estimate_mem ory	Integer	Estimated memory used by a statement on a single instance. The unit is MB.
min_peak_mem ory	Integer	Minimum memory peak of a statement across all DNs. The unit is MB.
max_peak_me mory	Integer	Maximum memory peak of a statement across all DNs. The unit is MB.
average_peak_ memory	Integer	Average memory usage during statement execution. The unit is MB.
memory_skew_ percent	Integer	Memory usage skew of a statement among DNs.
spill_info	Text	Statement spill information on all DNs. None indicates that the statement has not been spilled to disks on any DNs.
		All: The statement has been spilled to disks on all DNs. [a:b]: The statement has been spilled to disks on a of b DNs.
min_spill_size	Integer	Minimum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
max_spill_size	Integer	Maximum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
average_spill_si ze	Integer	Average spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
spill_skew_perc ent	Integer	DN spill skew when a spill occurs
min_dn_time	Bigint	Minimum execution time of a statement across all DNs. The unit is ms.
max_dn_time	Bigint	Maximum execution time of a statement across all DNs. The unit is ms.
average_dn_tim e	Bigint	Average execution time of a statement across all DNs. The unit is ms.
dntime_skew_p ercent	Integer	Execution time skew of a statement among DNs.
min_cpu_time	Bigint	Minimum CPU time of a statement across all DNs. The unit is ms.

Column	Туре	Description
max_cpu_time	Bigint	Maximum CPU time of a statement across all DNs. The unit is ms.
total_cpu_time	Bigint	Total CPU time of a statement across all DNs. The unit is ms.
cpu_skew_perce nt	Integer	CPU time skew of a statement among DNs.
min_peak_iops	Integer	Minimum IOPS peak of a statement across all DNs. It is counted by ones in a column-store table and by ten thousands in a row-store table.
max_peak_iops	Integer	Maximum IOPS peak of a statement across all DNs. It is counted by ones in a column-store table and by ten thousands in a row-store table.
average_peak_i ops	Integer	Average IOPS peak of a statement across all DNs. It is counted by ones in a column-store table and by ten thousands in a row-store table.
iops_skew_perc ent	Integer	I/O skew across DNs.
warning	Text	Warning. The following warnings and warnings related to SQL self-diagnosis tuning are displayed:
		1. Spill file size large than 256MB
		2. Broadcast size large than 100MB
		3. Early spill
		4. Spill times is greater than 3
		5. Spill on memory adaptive
		6. Hash table conflict
queryid	Bigint	Internal query ID used for statement execution
query	Text	Statement to be executed. A maximum of 64 KB of strings can be retained.

Column	Туре	Description
query_plan	Text	Execution plan of a statement.
		Specification restrictions:
		Execution plans are displayed only for DML statements.
		2. In 8.2.1.100 and later versions, the number of data binding times is added to the execution plans of Parse Bind Execute (PBE) statements to facilitate statement analysis. The number of data binding times is displayed in the format of PBE bind times : <i>Times</i> .
node_group	Text	Logical cluster of the user running the statement
pid	Bigint	PID of the backend thread of the statement
lane	Text	Fast/Slow lane where the statement is executed
unique_sql_id	Bigint	ID of the normalized unique SQL.
session_id	Text	Unique identifier of a session in the database system. Its format is session_start_time.tid.node_name.
min_read_bytes	Bigint	Minimum I/O read bytes of a statement across all DNs. The unit is byte.
max_read_byte s	Bigint	Maximum I/O read bytes of a statement across all DNs. The unit is byte.
average_read_b ytes	Bigint	Average I/O read bytes of a statement across all DNs.
min_write_byte s	Bigint	Minimum I/O write bytes of a statement across all DNs.
max_write_byte s	Bigint	Maximum I/O write bytes of a statement across all DNs.
average_write_ bytes	Bigint	Average I/O write bytes of a statement across all DNs.
recv_pkg	Bigint	Total number of communication packages received by a statement across all DNs.
send_pkg	Bigint	Total number of communication packages sent by a statement across all DNs.
recv_bytes	Bigint	Total received data of the statement stream, in byte.

Column	Туре	Description
send_bytes	Bigint	Total sent data of the statement stream, in byte.
stmt_type	Text	Query type corresponding to the statement.
except_info	Text	Information about the exception rule triggered by the statement.
unique_plan_id	Bigint	ID of the normalized unique plan.
sql_hash	Text	Normalized SQL hash.
plan_hash	Text	Normalized plan hash.
use_plan_baseli ne	Text	Indicates whether the bound plan is used for executing the current statement. If it is used, the name of the plan_baseline column in pg_plan_baseline is displayed.
outline_name	Text	Name of the outline used for the statement plan.
loader_status	Text	 The JSON string for storing import and export service information is as follows. address: indicates the IP address of the peer cluster. The port number is displayed for the source cluster. direction: indicates the import and export service type. The value can be gds to file, gds from file, gds to pipe, gds from pipe, copy from or copy to. min/max/total_lines/bytes: indicates the minimum value, maximum value, total lines, and bytes of the import and export
parse_time	Bigint	statements on all DNs. Total parsing time before the statement is queued (including lexical and syntax parsing, optimization rewriting, and plan generation time), in milliseconds. This column is available only in clusters of version 8.3.0.100 or later.
disk_cache_hit_ ratio	numeric(5,2)	Disk cache hit rate. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
disk_cache_disk _read_size	Bigint	Total size of data read from disk cache, in MB. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.

Column	Туре	Description
disk_cache_disk _write_size	Bigint	Total size of data written to disk cache, in MB. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
disk_cache_rem ote_read_size	Bigint	Total size of data read remotely from OBS due to disk cache read failure, in MB. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
disk_cache_rem ote_read_time	Bigint	Total number of times data is read remotely from OBS due to disk cache read failure. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
vfs_scan_bytes	Bigint	Total number of bytes scanned by the OBS virtual file system in response to upper-layer requests, in bytes. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
vfs_remote_rea d_bytes	Bigint	Total number of bytes actually read from OBS by the OBS virtual file system, in bytes. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
preload_submit _time	Bigint	Total time for submitting I/O requests in the prefetching process, in microseconds. This column only applies to OBS 3.0 tables with storage and compute decoupled.
preload_wait_ti me	Bigint	Total time for waiting for I/O requests in the prefetching process, in microseconds. This column only applies to OBS 3.0 tables with storage and compute decoupled.
preload_wait_c ount	Bigint	Total number of times that the prefetching process waits for I/O requests. This column only applies to OBS 3.0 tables with storage and compute decoupled.
disk_cache_loa d_time	Bigint	Total time for reading from disk cache, in microseconds. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
disk_cache_conf lict_count	Bigint	Number of times a block in the disk cache produces a hash conflict. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.

Column	Туре	Description
disk_cache_erro r_count	Bigint	Number of disk cache read failures. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
disk_cache_erro r_code	Bigint	 Error code for disk cache read failures. Multiple error codes may be generated. If the disk cache fails to be read, OBS remote read is initiated and cache blocks are rewritten. The error code types are as follows: This column only applies to OBS 3.0 tables and foreign tables. 1: A hash conflict occurs in the disk cache block. 2: The generation time of the disk cache block is later than that of the OldestXmin transaction. 4: Invoking the pread system when reading cache files from the disk cache failed. 8: The data version of the disk cache block does not match. 16: The version of the data written to the write cache does not match the latest version. 32: Opening the cache file corresponding to the cache block failed. 64: The size of the data read from the disk cache does not match.
		128: The CSN recorded in the disk cache block does not match.
obs_io_req_avg _rtt	Bigint	Average Round Trip Time (RTT) for OBS I/O requests, in microseconds. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
obs_io_req_avg _latency	Bigint	Average delay for OBS I/O requests, in microseconds. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
obs_io_req_late ncy_gt_1s	Bigint	Number of OBS I/O requests with a latency exceeding 1 second. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
obs_io_req_late ncy_gt_10s	Bigint	Number of OBS I/O requests with a latency exceeding 10 seconds. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.

Column	Туре	Description
obs_io_req_cou nt	Bigint	Total number of OBS I/O requests. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
obs_io_req_retr y_count	Bigint	Total number of retries for OBS I/O requests. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
obs_io_req_rate _limit_count	Bigint	Total number of times OBS I/O requests are flow-controlled. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.

14.2.8 GS_WLM_USER_RESOURCE_HISTORY

The **GS_WLM_USER_RESOURCE_HISTORY** system catalog stores information about resources used by users. The data of this table is stored on both CNs and DNs. Each record in the system table indicates the resource usage of a user at a time point, including the memory, number of CPU cores, storage space, temporary space, operator spill space, logical I/O traffic, number of logical I/O times, and logical I/O rate. The memory, CPU, and I/O monitoring items record only the resource usage of complex jobs.

Data in the **GS_WLM_USER_RESOURCE_HISTORY** system table comes from the **PG_TOTAL_USER_RESOURCE_INFO** view.

Table 14-8 GS_WLM_USER_RESOURCE_HISTORY column

Column	Туре	Description
username	Text	Username
Timestam p	Timestamp with time zone	Timestamp
used_me mory	Int	 Memory size used by a user, in MB. DN: The memory used by users on the current DN is displayed. CN: The total memory usage of users on all DNs is displayed.

Column	Туре	Description
total_me mory	Int	Memory used by the resource pool, in MB. 0 indicates that the available memory is not limited and depends on the maximum memory available in the database (max_dynamic_memory). A calculation formula is as follows:
		total_memory = max_dynamic_memory * parent_percent * user_percent
		CN: The sum of maximum available memory on all DNs is displayed.
used_cpu	Real	Number of CPU cores in use
total_cpu	Int	Total number of CPU cores of the Cgroup associated with a user on the node
used_spac e	Bigint	Used storage space (unit: KB)
total_spac e	Bigint	Available storage space (unit: KB)1 indicates that the storage space is not limited.
used_tem p_space	Bigint	Used temporary storage space (unit: KB)
total_tem p_space	Bigint	Available temporary storage space (unit: KB)1 indicates that the maximum temporary storage space is not limited.
used_spill _space	Bigint	Space occupied by operators spilled to disk (unit: KB)
total_spill _space	Bigint	Available storage space for operator spill to disk (unit: KB). The value -1 indicates that the space is not limited.
read_kbyt es	Bigint	Byte traffic of read operations in a monitoring period (unit: KB)
write_kby tes	Bigint	Byte traffic of write operations in a monitoring period (unit: KB)
read_cou nts	Bigint	Number of read operations in a monitoring period.
write_cou nts	Bigint	Number of write operations in a monitoring period.
read_spee d	Real	Byte rate of read operations in a monitoring period (unit: KB)
write_spe ed	Real	Byte rate of write operations in a monitoring period (unit: KB)

Column	Туре	Description
send_spee d	Double	Network sending rate in a monitoring period, in KB/s.
recv_spee d	Double	Network receiving rate in a monitoring period, in KB/s.

14.2.9 PG_AGGREGATE

pg_aggregate records information about aggregation functions. Each entry in **pg_aggregate** is an extension of an entry in **pg_proc**. The **pg_proc** entry carries the aggregate's name, input and output data types, and other information that is similar to ordinary functions.

Table 14-9 PG_AGGREGATE columns

Column	Туре	Reference	Description
aggfnoid	regproc	PG_PROC.oid	PG_PROC OID of the aggregate function.
aggtransfn	regproc	PG_PROC.oid	Transition function.
aggcollectfn	regproc	PG_PROC.oid	Aggregate function.
aggfinalfn	regproc	PG_PROC.oid	Final function (zero if none).
aggsortop	OID	PG_OPERATOR.oid	Associated sort operator (zero if none).
aggtranstype	OID	PG_TYPE.oid	Data type of the aggregate function's internal transition (state) data.
agginitval	Text	-	Initial value of the transition state. This is a text column containing the initial value in its external string representation. If this column is null, the transition state value starts out null.
agginitcollect	Text	-	Initial value of the collection state. This is a text column containing the initial value in its external string representation. If this column is null, the collection state value starts out null.

14.2.10 PG_AM

PG_AM records information about index access methods. There is one row for each index access method supported by the system.

Table 14-10 PG_AM columns

Column	Туре	Reference	Description
OID	OID	-	Row identifier (hidden attribute; displayed only when explicitly selected).
amname	Name	-	Name of the access method.
amstrategies	Smallint	-	Number of operator strategies for this access method, or zero if access method does not have a fixed set of operator strategies.
amsupport	Smallint	-	Number of support routines for this access method.
amcanorder	boolean	-	Whether the access method supports ordered scans sorted by the indexed column's value.
amcanorderbyo p	boolean	-	Whether the access method supports ordered scans sorted by the result of an operator on the indexed column.
amcanbackward	boolean	-	Whether the access method supports backward scanning.
amcanunique	boolean	-	Whether the access method supports unique indexes.
amcanmulticol	boolean	-	Whether the access method supports multi-column indexes.
amoptionalkey	boolean	-	Whether the access method supports a scan without any constraint for the first index column.
amsearcharray	boolean	-	Whether the access method supports ScalarArrayOpExpr searches.
amsearchnulls	boolean	-	Whether the access method supports IS NULL/NOT NULL searches.

Column	Туре	Reference	Description
amstorage	boolean	-	Whether an index storage data type can differ from a column data type.
amclusterable	boolean	-	Whether an index of this type can be clustered on.
ampredlocks	boolean	-	Whether an index of this type manages fine-grained predicate locks.
amkeytype	OID	PG_TYPE.oid	Type of data stored in index, or zero if not a fixed type.
aminsert	regproc	PG_PROC.oid	"Insert this tuple" function.
ambeginscan	regproc	PG_PROC.oid	"Prepare for index scan" function.
amgettuple	regproc	PG_PROC.oid	"Next valid tuple" function, or zero if none.
amgetbitmap	regproc	PG_PROC.oid	"Fetch all valid tuples" function, or zero if none.
amrescan	regproc	PG_PROC.oid	"(Re)start index scan" function.
amendscan	regproc	PG_PROC.oid	"Clean up after index scan" function.
ammarkpos	regproc	PG_PROC.oid	"Mark current scan position" function.
amrestrpos	regproc	PG_PROC.oid	"Restore marked scan position" function.
ammerge	regproc	PG_PROC.oid	"Merge multiple indexes" function.
ambuild	regproc	PG_PROC.oid	"Build new index" function.
ambuildempty	regproc	PG_PROC.oid	"Build empty index" function.
ambulkdelete	regproc	PG_PROC.oid	Bulk-delete function.
amvacuumclean up	regproc	PG_PROC.oid	Post- VACUUM cleanup function.
amcanreturn	regproc	PG_PROC.oid	Function to check whether index supports index-only scans, or zero if none
amcostestimate	regproc	PG_PROC.oid	Function to estimate cost of an index scan.
amoptions	regproc	PG_PROC.oid	Function to parse and validate reloptions for an index.

14.2.11 PG AMOP

PG_AMOP records information about operators associated with access method operator families. There is one row for each operator that is a member of an operator family. A family member can be either a search operator or an ordering operator. An operator can appear in more than one family, but cannot appear in more than one search position nor more than one ordering position within a family.

Table 14-11 PG_AMOP columns

Name	Туре	Reference	Description
OID	OID	-	Row identifier (hidden attribute; must be explicitly selected)
amopfamily	OID	PG_OPFAMILY.oid	Operator family this entry is for
amoplefttype	OID	PG_TYPE.oid	Left-hand input data type of operator
amoprighttype	OID	PG_TYPE.oid	Right-hand input data type of operator
amopstrategy	Smallint	-	Number of operator strategies
amoppurpose	Char	-	Operator purpose, either s for search or o for ordering
amopopr	OID	PG_OPERATOR.oid	OID of the operator
amopmethod	OID	PG_AM.oid	Index access method the operator family is for
amopsortfamily	OID	PG_OPFAMILY.oid	If it is a sort operator, the item is sorted according to the B-Tree operator family. If it is a search operator, the value is 0.

A "search" operator entry indicates that an index of this operator family can be searched to find all rows satisfying **WHERE indexed_column operator constant**. Obviously, such an operator must return a Boolean value, and its left-hand input type must match the index's column data type.

An "ordering" operator entry indicates that an index of this operator family can be scanned to return rows in the order represented by **ORDER BY indexed_column operator constant**. Such an operator could return any sortable data type, though

again its left-hand input type must match the index's column data type. The exact semantics of **ORDER BY** are specified by the **amopsortfamily** column, which must reference a B-tree operator family for the operator's result type.

14.2.12 PG AMPROC

PG_AMPROC records information about the support procedures associated with the access method operator families. There is one row for each support procedure belonging to an operator family.

Table 14-12 PG AMPROC columns

Name	Туре	Reference	Description
OID	OID	N/A	Row identifier (hidden attribute; displayed only when explicitly selected)
amprocfamily	OID	PG_OPFAMILY.oid	Operator family this entry is for
amproclefttype	OID	PG_TYPE.oid	Left-hand input data type of associated operator
amprocrightty pe	OID	PG_TYPE.oid	Right-hand input data type of associated operator
amprocnum	Smallin t	N/A	Support procedure number
amproc	regproc	PG_PROC.oid	OID of the procedure

The usual interpretation of the **amproclefttype** and **amprocrighttype** columns is that they identify the left and right input types of the operator(s) that a particular support procedure supports. For some access methods these match the input data type(s) of the support procedure itself, for others not. There is a notion of "default" support procedures for an index, which are those with **amproclefttype** and **amprocrighttype** both equal to the index opclass's **opcintype**.

14.2.13 PG ATTRDEF

PG_ATTRDEF stores default values of columns.

Table 14-13 PG_ATTRDEF columns

Column	Туре	Description
adrelid	OID	Table to which the column belongs
adnum	Smallint	Column No.
adbin	pg_node_tree	Internal representation of the column's default value

Column	Туре	Description
adsrc	Text	Internal representation of the human- readable default value
adbin_on_updat e	pg_node_tree	Internal representation of the value of on_update_expr
adsrc_on_updat e	Text	Internal representation of the human- readable value of on_update_expr

14.2.14 PG_ATTRIBUTE

PG_ATTRIBUTE records information about table columns.

Table 14-14 PG_ATTRIBUTE columns

Column	Туре	Description
attrelid	OID	Table to which the column belongs
attname	Name	Column name
atttypid	OID	Column type
attstattarget	Integer	Controls the level of details of statistics collected for this column by ANALYZE .
		 A zero value indicates that no statistics should be collected.
		A negative value says to use the system default statistics target.
		The exact meaning of positive values is data type-dependent.
		For scalar data types, attstattarget is both the target number of "most common values" to collect, and the target number of histogram bins to create.
attlen	Smallint	Copy of pg_type.typlen of the column's type
attnum	Smallint	Number of a column.
attndims	Integer	Number of dimensions if the column is an array; otherwise, the value is 0.
attcacheoff	Integer	This column is always -1 on disk. When it is loaded into a row descriptor in the memory, it may be updated to cache the offset of the columns in the row.

Column	Туре	Description
atttypmod	Integer	Type-specific data supplied at table creation time (for example, the maximum length of a varchar column). This column is used as the third parameter when passing to type-specific input functions and length coercion functions. The value will generally be -1 for types that do not need ATTTYPMOD.
attbyval	boolean	Copy of pg_type.typbyval of the column's type
attstorage	Char	Copy of pg_type.typstorage of this column's type
attalign	Char	Copy of pg_type.typalign of the column's type
attnotnull	boolean	A not-null constraint. It is possible to change this column to enable or disable the constraint.
atthasdef	boolean	Indicates that this column has a default value, in which case there will be a corresponding entry in the pg_attrdef table that actually defines the value.
attisdropped	boolean	Whether the column has been dropped and is no longer valid. A dropped column is still physically present in the table but is ignored by the analyzer, so it cannot be accessed through SQL.
attislocal	boolean	Whether the column is defined locally in the relation. Note that a column can be locally defined and inherited simultaneously.
attcmprmode	tinyint	Compressed modes for a specific column The compressed mode includes: • ATT_CMPR_NOCOMPRESS • ATT_CMPR_DELTA • ATT_CMPR_DICTIONARY • ATT_CMPR_PREFIX • ATT_CMPR_NUMSTR
attinhcount	Integer	Number of direct ancestors this column has. A column with an ancestor cannot be dropped nor renamed.
attcollation	OID	Defined collation of a column
attacl	aclitem[]	Permissions for column-level access
attoptions	text[]	Property-level options

Column	Туре	Description
attfdwoptions	text[]	Property-level external data options
attinitdefval	bytea	attinitdefval stores the default value expression. ADD COLUMN in a row-store table must use this column.
attkvtype	tinyint	 kv_type attribute of a column. Values: 0 indicates the default value, which is used for non-time series tables. 1 indicates TSTAG, a dimension attribute, which is used only for time series tables. 2 indicates TSFIELD, a metric attribute, which is used only for time series tables. 3 indicates TSTIME, a time attribute, which is used only for time series tables.

Example

Query the field names and field IDs of a specified table. Replace **t1** and **public** with the actual table name and schema name, respectively.

SELECT attname,attnum FROM pg_attribute WHERE attrelid=(SELECT pg_class.oid FROM pg_class JOIN pg_namespace ON relnamespace=pg_namespace.oid WHERE relname='t1' and nspname='public') and attnum>0;

attname	attnum
product_id product_name product_quanti (3 rows)	1 2 ty 3

14.2.15 PG_AUTHID

PG_AUTHID records information about the database authentication identifiers (roles). The concept of users is contained in that of roles. A user is actually a role whose rolcanlogin has been set. Any role, whether the rolcanlogin is set or not, can use other roles as members.

For a cluster, only one **pg_authid** exists which is not available for every database. It is accessible only to users with system administrator rights.

Table 14-15 PG_AUTHID columns

Column	Туре	Description
OID	OID	Row identifier (hidden attribute; must be explicitly selected)
rolname	Name	Role name
rolsuper	boolean	Whether the role is the initial system administrator with the highest permission

Column	Туре	Description
rolinherit	boolean	Whether the role automatically inherits permissions of roles it is a member of
rolcreaterole	boolean	Whether the role can create more roles
rolcreatedb	boolean	Whether the role can create databases
rolcatupdate	boolean	Whether the role can directly update system catalogs. Only the initial system administrator whose usesysid is 10 has this permission. It is not available for other users.
rolcanlogin	boolean	Whether a role can log in, that is, whether a role can be given as the initial session authorization identifier.
rolreplication	boolean	Indicates that the role is a replicated one (an adaptation syntax and no actual meaning).
rolauditadmin	boolean	Indicates that the role is an audit user.
rolsystemadmin	boolean	Indicates that the role is an administrator.
rolconnlimit	Integer	Limits the maximum number of concurrent connections of a user on a CN1 means no limit.
rolpassword	Text	Password (possibly encrypted); NULL if no password.
rolvalidbegin	Timestam p with time zone	Account validity start time; NULL if no start time
rolvaliduntil	Timestam p with time zone	Password expiry time; NULL if no expiration
rolrespool	Name	Resource pool that a user can use
roluseft	boolean	Whether the role can perform operations on foreign tables
rolparentid	OID	OID of a group user to which the user belongs
roltabspace	Text	Storage space of the user permanent table
rolkind	Char	Special type of user, including private users, logical cluster administrators, and common users.
rolnodegroup	OID	OID of a node group associated with a user. The node group must be a logical cluster.
roltempspace	Text	Storage space of the user temporary table

Column	Туре	Description
rolspillspace	Text	Operator disk spill space of the user
rolexcpdata	Text	Reserved column
rolauthinfo	Text	Additional information when LDAP authentication is used. If other authentication modes are used, the value is NULL .
rolpwdexpire	Integer	Password expiration time. Users can change their password before it expires. After the password expires, only the administrator can change the password. The value -1 indicates that the password never expires.
rolpwdtime	Timestam p with time zone	Time when a password is created
roluuid	Bigint	Role identifier. This column is available only in clusters of version 9.1.0 or later.

14.2.16 PG_AUTH_HISTORY

PG_AUTH_HISTORY records the authentication history of the role. It is accessible only to users with system administrator rights.

Table 14-16 PG_AUTH_HISTORY columns

Column	Туре	Description
roloid	OID	Role identifier
passwordtime	Timestamp with time zone	Time of password creation and change
rolpassword	Text	Role password that is encrypted using MD5 or SHA256, or that is not encrypted

14.2.17 PG_AUTH_MEMBERS

PG_AUTH_MEMBERS records the membership relations between roles.

Table 14-17 PG_AUTH_MEMBERS columns

Column	Туре	Description
roleid	OID	ID of a role that has a member

Column	Туре	Description
member	OID	ID of a role that is a member of ROLEID
grantor	OID	ID of a role that grants this membership
admin_option	boolean	Whether a member can grant membership in ROLEID to others

14.2.18 PG_BLOCKLISTS

PG_BLOCKLISTS records query filtering rules. This system catalog is supported only by clusters of version 9.1.0.100 or later.

Table 14-18 PG_BLOCKLISTS columns

Column	Туре	Description
block_name	Name	Name of a query filtering rule
role	OID	User OID bound to the query filtering rule
client_addr	inet	IP address of the client bound to the query filtering rule
application_na me	Name	Name of the client bound to the query filtering rule
unique_sql_id	INT8	unique_sql_id that matches the query filtering rule
sql_hash	Name	sql_hash that matches the query filtering rule
block_type	INT4	Type of the statement bound to the query filtering rule. The type can be SELECT, UPDATE, INSERT, DELETE, or MERGE.
partition_num	INT4	Estimated maximum number of partitions to be scanned
table_num	INT4	Estimated maximum number of tables to be scanned
estimate_row	INT4	Estimated maximum number of rows to be scanned
query_band	Text	Type of the job that is actively identified
sql	Text	SQL statement that matches the query filtering rule
created_time	Timestamp with time zone	Timestamp when a query filtering rule is created or modified

Column	Туре	Description
resource_pool	Name	Name of the resource pool to which the statement intercepted by the query filtering rule is switched. This column is available only in clusters of version 9.1.0.200 or later.
max_active_nu m	Integer	Maximum number of concurrent statements intercepted by the query filtering rule. If the value is lower than the specified limit, execution proceeds normally. However, if the value is equal to or exceeds the limit, an error is reported and the statements are intercepted. This column is available only in clusters of version 9.1.0.200 or later.
is_warning	Boolean	Whether an error or alarm is reported when a statement is intercepted by the query filtering rule. • false indicates that an error is reported when a statement is intercepted. The default value is false. • true indicates that an alarm is
		generated when a statement is intercepted.
		This column is available only in clusters of version 9.1.0.200 or later.

14.2.19 PG_CAST

PG_CAST records conversion relationships between data types.

Table 14-19 PG_CAST columns

Column	Туре	Description
castsource	OID	OID of the source data type
casttarget	OID	OID of the target data type
castfunc	OID	OID of the conversion function. If the value is 0 , no conversion function is required.

Column	Туре	Description	
castcontext	Char	Conversion mode between the source and target data types	
		• e indicates that only explicit conversion can be performed (using the CAST or :: syntax).	
		• i indicates that only implicit conversion can be performed.	
		a indicates that both explicit and implicit conversion can be performed between data types.	
castmethod	Char	Conversion method	
		• f indicates that conversion is performed using the specified function in the castfunc column.	
		b indicates that binary forcible conversion rather than the specified function in the castfunc column is performed between data types.	

14.2.20 PG_CLASS

PG_CLASS records database objects and their relations.

Table 14-20 PG_CLASS columns

Column	Туре	Description	
OID	OID	Row identifier (hidden attribute; must be explicitly selected)	
relname	Name	Name of an object, such as a table, index, or view	
relnamespace	OID	OID of the namespace that contains the relationship	
reltype	OID	Data type that corresponds to this table's row type (the index is 0 because the index does not have pg_type record)	
reloftype	OID	OID is of composite type. 0 indicates other types.	
relowner	OID	Owner of the relationship	
relam	OID	Specifies the access method used, such as B-tree and hash, if this is an index	
relfilenode	OID	Name of the on-disk file of this relationship. If such file does not exist, the value is 0 .	

Column	Туре	Description
reltablespace	OID	Tablespace in which this relationship is stored. If its value is 0 , the default tablespace in this database is used. This column is meaningless if the relationship has no on-disk file.
relpages	Double precisio n	Size of the on-disk representation of this table in pages (of size BLCKSZ). This is only an estimate used by the optimizer.
reltuples	Double precisio n	Number of rows in the table. This is only an estimate used by the optimizer.
relallvisible	Integer	Number of pages marked as all visible in the table. This column is used by the optimizer for optimizing SQL execution. It is updated by VACUUM, ANALYZE, and a few DDL statements such as CREATE INDEX.
reltoastrelid	OID	OID of the TOAST table associated with this table. The OID is 0 if no TOAST table exists.
		The TOAST table stores large columns "offline" in a secondary table.
reltoastidxid	OID	OID of the index for a TOAST table. The OID is 0 for a table other than a TOAST table.
reldeltarelid	OID	OID of a Delta table Delta tables belong to column-store tables. They store long tail data generated during data import.
reldeltaidx	OID	OID of the index for a Delta table
relcudescrelid	OID	OID of a CU description table CU description tables (Desc tables) belong to column-store tables. They control whether storage data in the HDFS table directory is visible.
relcudescidx	OID	OID of the index for a CU description table
relhasindex	boolean	Its value is true if this column is a table and has (or recently had) at least one index. It is set by CREATE INDEX but is not immediately cleared by DROP INDEX . If the VACUUM process detects that a table has no index, it clears the relhasindex column and sets the value to false .
relisshared	boolean	Its value is true if the table is shared across all databases in the cluster. Only certain system catalogs (such as pg_database) are shared.

Column	Туре	Description	
relpersistence	Char	 p indicates a permanent table. u indicates a non-log table. t indicates a local temporary table. g indicates a global temporary table. 	
relkind	Char	 r indicates an ordinary table. i indicates an index. S indicates a sequence. v indicates a view. c indicates the composite type. t indicates a TOAST table. f indicates a foreign table. m indicates a materialized view. 	
relnatts	Smallin t	Number of user columns in the relationship (excluding system columns) pg_attribute has the same number of rows corresponding to the user columns.	
relchecks	Smallin t	Number of constraints on a table. For details, see PG_CONSTRAINT.	
relhasoids	boolean	Its value is true if an OID is generated for each row of the relationship.	
relhaspkey	boolean	Its value is true if the table has (or once had) a primary key.	
relhasrules	boolean	Its value is true if the table has rules. See table PG_REWRITE to check whether it has rules.	
relhastriggers	boolean	Its value is true if the table has (or once had) triggers. For details, see PG_TRIGGER .	
relhassubclass	boolean	Its value is true if the table has (or once had) any inheritance child table.	
relcmprs	tinyint	Whether the compression feature is enabled for the table. Note that only batch insertion triggers compression so ordinary CRUD does not trigger compression.	
		• 0 indicates other tables that do not support compression (primarily system tables, on which the compression attribute cannot be modified).	
		1 indicates that the compression feature of the table data is NOCOMPRESS or has no specified keyword.	
		• 2 indicates that the compression feature of the table data is COMPRESS.	

Column	Туре	Description	
relhasclusterkey	boolean	Whether the local cluster storage is used	
relrowmoveme nt	boolean	 Whether the row migration is allowed when the partitioned table is updated true indicates that the row migration is allowed. false indicates that the row migration is not allowed. 	
parttype	Char	 Whether the table or index has the property of a partitioned table p indicates that the table or index has the property of a partitioned table. n indicates that the table or index does not have the property of a partitioned table. v indicates that the table is the value partitioned table in the HDFS. 	
relfrozenxid	xid32	All transaction IDs before this one have been replaced with a permanent ("frozen") transaction ID in this table. This column is used to track whether the table needs to be vacuumed in order to prevent transaction ID wraparound (or to allow pg_clog to be shrunk). The value is 0 (InvalidTransactionId) if the relationship is not a table. To ensure forward compatibility, this column is reserved. The relfrozenxid64 column is added to record the information.	
relacl	aclite m[]	Access permissions The command output of the query is as follows: rolename=xxxx/yyyyAssigning privileges to a role =xxxx/yyyyAssigning the permission to public xxxx indicates the assigned privileges, and yyyy indicates the roles that are assigned to the privileges. For details about permission descriptions, see Table 14-21.	
reloptions	text[]	Access-method-specific options, as "keyword=value" strings	
relfrozenxid64	Xid	All transaction IDs before this one have been replaced with a permanent ("frozen") transaction ID in this table. This column is used to track whether the table needs to be vacuumed in order to prevent transaction ID wraparound (or to allow pg_clog to be shrunk). The value is 0 (InvalidTransactionId) if the relationship is not a table.	

Table 14-21 Description of privileges

Parameter	Description
r	SELECT (read)
w	UPDATE (write)
a	INSERT (insert)
d	DELETE
D	TRUNCATE
х	REFERENCES
t	TRIGGER
X	EXECUTE
U	USAGE
С	CREATE
С	CONNECT
Т	TEMPORARY
А	ANALYZE ANALYSE
L	ALTER
Р	DROP
v	VACUUM
arwdDxtA, vLP	ALL PRIVILEGES (used for tables)
*	Authorization options for preceding permissions

Examples

View the OID and relfilenode of a table.

SELECT oid,relname,relfilenode FROM pg_class WHERE relname = 'table_name';

Count row-store tables.

SELECT 'row count:'||count(1) as point FROM pg_class WHERE relkind = 'r' and oid > 16384 and reloptions::text not like '%column%' and reloptions::text not like '%internal_mask%';

Count column-store tables.

SELECT 'column count:'||count(1) as point FROM pg_class WHERE relkind = 'r' and oid > 16384 and reloptions::text like '%column%';

Query the comments of all tables in the database:

SELECT relname as tablename,obj_description(relfilenode,'pg_class') as comment FROM pg_class;

14.2.21 PG_COLLATION

PG_COLLATION records the available collations, which are essentially mappings from an SQL name to operating system locale categories.

Table 14-22 PG_COLLATION columns

Name	Туре	Reference	Description
OID	OID	N/A	Row identifier (hidden attribute; displayed only when explicitly selected)
collname	Name	N/A	Collation name (unique per namespace and encoding)
collnamespace	OID	PG_NAMESPACE.oi	OID of the namespace that contains this collation
collowner	OID	PG_AUTHID.oid	Owner of the collation
collencoding	Integer	N/A	Encoding in which the collation is applicable, or -1 if it works for any encoding
collcollate	Name	N/A	LC_COLLATE for this collation object
collctype	Name	N/A	LC_CTYPE for this collation object

14.2.22 PG_CONSTRAINT

PG_CONSTRAINT records check, primary key, unique, and foreign key constraints on the tables.

Table 14-23 PG_CONSTRAINT columns

Column	Туре	Description
conname	Name	Constraint name (not necessarily unique)
connamespace	OID	OID of the namespace that contains the constraint
contype	Char	 c indicates check constraints. f indicates foreign key constraints. p indicates primary key constraints. u indicates unique constraints. t indicates trigger constraints.
condeferrable	boolean	Whether the constraint can be deferrable

Column	Туре	Description
condeferred	boolean	Whether the constraint can be deferrable by default
convalidated	boolean	Whether the constraint is valid Currently, only foreign key and check constraints can be set to false.
conrelid	OID	Table containing this constraint. The value is 0 if it is not a table constraint.
contypid	OID	Domain containing this constraint. The value is 0 if it is not a domain constraint.
conindid	OID	ID of the index associated with the constraint
confrelid	OID	Referenced table if this constraint is a foreign key; otherwise, the value is 0 .
confupdtype	Char	 Foreign key update action code a indicates no action. r indicates restriction. c indicates cascading. n indicates that the parameter is set to null. d indicates that the default value is used.
confdeltype	Char	 Foreign key deletion action code a indicates no action. r indicates restriction. c indicates cascading. n indicates that the parameter is set to null. d indicates that the default value is used.
confmatchtype	Char	 Foreign key match type f indicates full match. p indicates partial match. u indicates simple match (not specified).
conislocal	boolean	Whether the local constraint is defined for the relationship
coninhcount	Integer	Number of direct inheritance parent tables this constraint has. When the number is not 0 , the constraint cannot be deleted or renamed.
connoinherit	boolean	Whether the constraint can be inherited

Column	Туре	Description
consoft	boolean	Whether the column indicates an informational constraint.
conopt	boolean	Whether you can use Informational Constraint to optimize the execution plan.
conkey	smallint[]	Column list of the constrained control if this column is a table constraint
confkey	smallint[]	List of referenced columns if this column is a foreign key
conpfeqop	oid[]	ID list of the equality operators for PK = FK comparisons if this column is a foreign key
conppeqop	oid[]	ID list of the equality operators for PK = PK comparisons if this column is a foreign key
conffeqop	oid[]	ID list of the equality operators for FK = FK comparisons if this column is a foreign key
conexclop	oid[]	ID list of the per-column exclusion operators if this column is an exclusion constraint
conbin	pg_node_tr ee	Internal representation of the expression if this column is a check constraint
consrc	Text	Human-readable representation of the expression if this column is a check constraint

NOTICE

- **consrc** is not updated when referenced objects change; for example, it will not track renaming of columns. You are advised to use **pg_get_constraintdef()** to extract the definition of a check constraint instead of depending on this column.
- **pg_class.relchecks** must be consistent with the number of check-constraint entries in this table for each relationship.

Example

Query whether a specified table has a primary key.

```
CREATE TABLE t1
(

C_CUSTKEY BIGINT ,
C_NAME VARCHAR(25) ,
C_ADDRESS VARCHAR(40) ,
C_NATIONKEY INT ,
C_PHONE CHAR(15) ,
C_ACCTBAL DECIMAL(15,2),
CONSTRAINT C_CUSTKEY_KEY PRIMARY KEY(C_CUSTKEY,C_NAME)
```

```
)
DISTRIBUTE BY HASH(C_CUSTKEY,C_NAME);

SELECT conname FROM pg_constraint WHERE conrelid = 't1'::regclass AND contype = 'p';
conname
------
c_custkey_key
(1 row)
```

14.2.23 PG_CONVERSION

PG_CONVERSION records encoding conversion information.

Table 14-24 PG_CONVERSION columns

Column	Туре	Reference	Description
OID	OID	N/A	Row identifier (hidden attribute; displayed only when explicitly selected)
conname	Name	N/A	Conversion name (unique in a namespace)
connamespace	OID	PG_NAMESPACE. oid	OID of the namespace that contains this conversion
conowner	OID	PG_AUTHID.oid	Owner of the conversion
conforencoding	Integer	N/A	Source encoding ID
contoencoding	Integer	N/A	Destination encoding ID
conproc	regproc	PG_PROC.oid	Conversion procedure
condefault	boolean	N/A	Whether the default conversion is used

14.2.24 PG_DATABASE

PG_DATABASE records information about the available databases.

Table 14-25 PG_DATABASE columns

Column	Туре	Description
datname	Name	Database name
datdba	OID	Owner of the database, usually the user who created it
encoding	Integer	Character encoding for this database
		You can use pg_encoding_to_char() to convert this number to the encoding name.

Column	Туре	Description
datcollate	Name	Sequence used by the database
datctype	Name	Character type used by the database
datistemplate	boolean	Whether this column can serve as a template database
datallowconn	boolean	If false then no one can connect to this database. This column is used to protect the template0 database from being altered.
datconnlimit	Integer	Maximum number of concurrent connections allowed on this database1 indicates no limit.
datlastsysoid	OID	Last system OID in the database
datfrozenxid	xid32	Tracks whether the database needs to be vacuumed in order to prevent transaction ID wraparound. To ensure forward compatibility, this column is reserved. The datfrozenxid64 column is
		added to record the information.
dattablespace	OID	Default tablespace of the database
datcompatibility	Name	 ORA: compatibility mode ORA: compatible with the Oracle database TD: compatible with the Teradata database MySQL: compatible with the MySQL
		database
datacl	aclitem[]	Permission to access the database.
datfrozenxid64	Xid	Tracks whether the database needs to be vacuumed in order to prevent transaction ID wraparound.

14.2.25 PG_DB_ROLE_SETTING

PG_DB_ROLE_SETTING records the default values of configuration items bonded to each role and database when the database is running.

Table 11 20 1 G_DB_KOLE_SETTING COLUMNS			
Column	Туре	Description	
setdatabase	OID	Database corresponding to the configuration items; the value is 0 if the database is not specified.	
setrole	OID	Role corresponding to the configuration items; the value is 0 if the role is not specified.	
setconfig	text[]	Default value of configuration items when the database is running.	

Table 14-26 PG DB ROLE SETTING columns

14.2.26 PG_DEFAULT_ACL

PG_DEFAULT_ACL records the initial privileges assigned to the newly created objects.

Table 14-27 PG_DEFAULT_ACL columns

Column	Туре	Description
defactrole	OID	ID of the role associated with the permission
defaclnamespace	OID	Namespace associated with the permission; the value is 0 if no ID
defaclobjtype	Char	Object type of the permission: • r indicates a table or view. • S indicates a sequence. • f indicates a function. • T indicates a type.
defaclacl	aclitem[]	Access permissions that this type of object should have on creation

Examples

Run the following command to view the initial permissions of the new user role1:

You can also run the following statement to convert the format:

SELECT pg_catalog.pg_get_userbyid(d.defaclrole) AS "Granter", n.nspname AS "Schema", CASE d.defaclobjtype WHEN 'r' THEN 'table' WHEN 'S' THEN 'sequence' WHEN 'f' THEN 'function' WHEN 'T' THEN 'type' END AS "Type", pg_catalog.array_to_string(d.defaclacl, E', ') AS "Access privileges" FROM pg_catalog.pg_default_acl d LEFT JOIN pg_catalog.pg_namespace n ON n.oid = d.defaclnamespace ORDER BY 1, 2, 3;

If the following information is displayed, **user1** grants **role1** the read permission on schema **user1**.

```
Granter | Schema | Type | Access privileges
-------
user1 | user1 | table | role1=r/user1
(1 row)
```

14.2.27 PG_DEPEND

PG_DEPEND records the dependency relationships between database objects. This information allows **DROP** commands to find which other objects must be dropped by **DROP CASCADE** or prevent dropping in the **DROP RESTRICT** case.

See also **PG_SHDEPEND**, which provides similar functionality for recording dependencies between objects that are shared between database clusters.

Table	14-28	PG	DEPEND	columns
-------	-------	----	--------	---------

Column	Туре	Reference	Description
classid	OID	PG_CLASS.oid	OID of the system catalog the dependent object is in
objid	OID	Any OID column	OID of the specific dependent object
objsubid	Integer	N/A	For a table column, this is the column number (the objid and classid refer to the table itself). For all other object types, this column is 0 .
refclassid	OID	PG_CLASS.oid	OID of the system catalog the referenced object is in
refobjid	OID	Any OID column	OID of the specific referenced object
refobjsubid	Integer	N/A	For a table column, this is the column number (the refobjid and refclassid refer to the table itself). For all other object types, this column is 0 .
deptype	Char	N/A	A code defining the specific semantics of this dependency relationship

In all cases, a **pg_depend** entry indicates that the referenced object cannot be dropped without also dropping the dependent object. However, there are several subflavors defined by **deptype**:

DEPENDENCY_NORMAL (n): A normal relationship between separately-created objects. The dependent object can be dropped without affecting the referenced object. The referenced object can only be dropped by specifying CASCADE, in which case the dependent object is dropped, too. Example: a table column has a normal dependency on its data type.

- DEPENDENCY_AUTO (a): The dependent object can be dropped separately from the referenced object, and should be automatically dropped (regardless of RESTRICT or CASCADE mode) if the referenced object is dropped. Example: a named constraint on a table is made autodependent on the table, so that it will go away if the table is dropped.
- DEPENDENCY_INTERNAL (i): The dependent object was created as part of creation of the referenced object, and is only a part of its internal implementation. A DROP of the dependent object will be disallowed outright (We'll tell the user to issue a DROP against the referenced object, instead). A DROP of the referenced object will be propagated through to drop the dependent object whether CASCADE is specified or not. For example, a trigger used to enforce a foreign key constraint is set to an item internally dependent on its constraint in PG_CONSTRAINT.
- DEPENDENCY_EXTENSION (e): The dependent object is a member of the
 extension that is the referenced object. (For details, see PG_EXTENSION). The
 dependent object can be dropped via DROP EXTENSION on the referenced
 object. Functionally this dependency type acts the same as an internal
 dependency, but it is kept separate for clarity and to simplify gs_dump.
- DEPENDENCY_PIN (p): There is no dependent object. This indicates that the system itself depends on the referenced object, and therefore the object cannot be deleted. Entries of this type are created only by **initdb**. The columns with dependent object are all zeroes.

Examples

Query the table that depends on the database object sequence **serial1**:

Query the OID of the sequence serial1 in the system catalog PG_CLASS.
 SELECT oid FROM pg_class WHERE relname ='serial1';
 OID
 17815
 (1 row)

2. Use the system catalog **PG_DEPEND** and the OID of **serial1** to obtain the objects that depend on **serial1**.

3. Obtain the OID of the table that depends on the serial1 sequence based on the refobjid field and query the table name. The result indicates that the table **customer address** depends on **serial1**.

```
SELECT relname FROM pg_class where oid='17812';
relname
------
customer_address
(1 row)
```

14.2.28 PG_DESCRIPTION

PG_DESCRIPTION records optional descriptions (comments) for each database object. Descriptions of many built-in system objects are provided in the initial contents of **PG_DESCRIPTION**.

See also **PG_SHDESCRIPTION**, which performs a similar function for descriptions involving objects that are shared across a database cluster.

Table 14-29 PG DESCRIPTION columns

Column	Туре	Reference	Description
objoid	OID	Any OID column	OID of the object this description pertains to
classoid	OID	PG_CLASS.oid	OID of the system catalog this object appears in
objsubid	Integer	-	For a comment on a table column, this is the column number (the objoid and classoid refer to the table itself). For all other object types, this column is 0 .
description	Text	-	Arbitrary text that serves as the description of this object

14.2.29 PG_ENUM

PG_ENUM records entries showing the values and labels for each enum type. The internal representation of a given enum value is actually the OID of its associated row in **pg_enum**.

Table 14-30 PG_ENUM columns

Column	Туре	Reference	Description
OID	OID	N/A	Row identifier (hidden attribute; displayed only when explicitly selected)
enumtypid	OID	PG_TYPE.oid	OID of pg_type that contains this enum value
enumsortorde r	Real	N/A	Sort position of this enum value within its enum type
enumlabel	Name	N/A	Textual label for this enum value

The OIDs for **PG_ENUM** rows follow a special rule: even-numbered OIDs are guaranteed to be ordered in the same way as the sort ordering of their enum type. That is, if two even OIDs belong to the same enum type, the smaller OID must have the smaller **enumsortorder** value. Odd-numbered OID values need bear no relationship to the sort order. This rule allows the enum comparison routines to avoid catalog lookups in many common cases. The routines that create and alter enum types attempt to assign even OIDs to enum values whenever possible.

When an enum type is created, its members are assigned sort-order positions from 1 to *n*. But members added later might be given negative or fractional values of

enumsortorder. The only requirement on these values is that they be correctly ordered and unique within each enum type.

14.2.30 PG_EXCEPT_RULE

The **PG_EXCEPT_RULE** system catalog stores information about exception rules. An exception rule set consists of multiple exception rules with the same name.

Table 14-31 PG_EXCEPT_RULE

Column	Туре	Description
Name	Name	Name of an exception rule set.
rule	Name	Type of a rule in the exception rule set, or action taken when the current exception rule set is triggered. (For example, it can be blocktime, elapsedtime, spillsize, or an action taken after an exception rule is triggered.)
value	Name	Rule threshold corresponding to the exception rule. If it specifies the action taken after an exception rule is triggered, its value is abort .

14.2.31 PG_EXTERNAL_NAMESPACE

Stores EXTERNAL SCHEMA information. This system catalog is supported only in 8.3.0 and later versions.

Table 14-32 PG_EXTERNAL_NAMESPACE columns

Name	Туре	Description
nspid	OID	External schema OID
srvname	Text	Name of the foreign server
source	Text	Metadata service type
address	Text	Metadata service address
database	Text	Metadata server database
confpath	Text	Path of the configuration file of the metadata server
ensoptions	Text[]	Reserved column, which is left empty currently.
catalog	Text	Metadata server catalog

Example

Query the created EXTERNAL SCHEMA ex1:

SELECT * FROM pg_external_namespace WHERE nspid = (SELECT oid FROM pg_namespace WHERE nspname = 'ex1');

14.2.32 PG EXTENSION

PG_EXTENSION records information about the installed extensions. By default, GaussDB(DWS) has 34 extensions: aio_scheduler, btree_gin, cudesckv, dimsearch, dist_fdw, functional_clog, functional_extension, functional_file, functional_hudi, functional_job, functional_largeobject, functional_memory, functional_other, functional_signal, functional_vacuum, gc_fdw, hdfs_fdw, hstore, log_fdw, operational_backup, operational_cgroup, operational_cudesc, operational_other, operational_replication, operational_restoration, operational_stats, operational_xlog, packages, pgcrypto, pldbgapi, plpgsql, roach_api, tsdb, and uuidossp.

Table 14-33 PG_EXTENSION

Column	Туре	Description
extname	Name	Extension name
extowner	OID	Owner of the extension
extnamespace	OID	Namespace containing the extension's exported objects
extrelocatable	boolean	Whether the extension can be relocated to another schema
extversion	Text	Version number of the extension
extconfig	oid[]	Configuration information about the extension
extcondition	Text[]	Filter conditions for the extension's configuration information

14.2.33 PG_EXTENSION_DATA_SOURCE

PG_EXTENSION_DATA_SOURCE records information about external data source. An external data source contains information about an external database, such as its password encoding. It is mainly used with Extension Connector.

Table 14-34 PG_EXTENSION_DATA_SOURCE columns

Name	Туре	Referenc e	Description
OID	OID	-	Row identifier (hidden attribute; must be explicitly selected)

Name	Туре	Referenc e	Description
srcname	Name	-	Name of an external data source
srcowner	OID	PG_AUTH ID.oid	Owner of an external data source
srctype	Text	-	Type of an external data source. It is NULL by default.
srcversion	Text	-	Type of an external data source. It is NULL by default.
srcacl	aclitem[]	-	Access permissions
srcoptions	Text[]	-	Option used for foreign data sources. It is a keyword=value string.

14.2.34 PG_FINE_DR_INFO

The **PG_FINE_DR_INFO** system catalog records the replay status of the fine-grained DR standby table. This system catalog is supported only by clusters of version 8.2.0.100 or later.

Table 14-35 PG_FINE_DR_INFO columns

Name	Туре	Description
OID	OID	Row identifier (hidden attribute; displayed only when explicitly selected)
relid	OID	OID of the standby fine-grained DR table
lastcsn	Xid	CSN of the last successful playback
lastxmin	Xid	xmin of the last successful playback
lastxmax	Xid	xmax of the last successful playback
laststarttime	Timestamp with time zone	Start time of the last successful playback
lastendtime	Timestamp with time zone	End time of the last successful playback

Examples

Check the playback status of the standby table in the DR cluster.

SELECT * FROM pg_fine_dr_info;				
relid lastcsn lastxmin lastxmax	laststarttime	- [lastendtime	
++++		+		

21132 | 1251610 | 1251609 | 1251611 | 2023-01-04 20:51:58.375136+08 | 2023-01-04 20:51:58.393986+08 (1 row)

14.2.35 PG_FOREIGN_DATA_WRAPPER

PG_FOREIGN_DATA_WRAPPER records foreign-data wrapper definitions. A foreign-data wrapper is the mechanism by which external data, residing on foreign servers, is accessed.

Table 14-36 PG_FOREIGN_DATA_WRAPPER columns

Name	Туре	Reference	Description
OID	OID	N/A	Row identifier (hidden attribute; displayed only when explicitly selected)
fdwname	Name	N/A	Name of the foreign-data wrapper
fdwowner	OID	PG_AUTHID.oid	Owner of the foreign-data wrapper
fdwhandler	OID	PG_PROC.oid	References a handler function that is responsible for supplying execution routines for the foreign-data wrapper. Its value is 0 if no handler is provided.
fdwvalidat or	OID	PG_PROC.oid	References a validator function that is responsible for checking the validity of the options given to the foreign-data wrapper, as well as options for foreign servers and user mappings using the foreign-data wrapper. Its value is 0 if no validator is provided.
fdwacl	aclite m[]	N/A	Access permissions
fdwoptions	Text[]	N/A	Option used for foreign data wrappers. It is a keyword=value string.

14.2.36 PG_FOREIGN_SERVER

PG_FOREIGN_SERVER records the foreign server definitions. A foreign server describes a source of external data, such as a remote server. Foreign servers are accessed via foreign-data wrappers.

Name Reference Description Type OID OID Row identifier (hidden N/A attribute; displayed only when explicitly selected) Name N/A Name of the foreign server srvname srvowner OID PG_AUTHID.oid Owner of the foreign server OID srvfdw PG_FOREIGN_DATA_ OID of the foreign-data wrapper of this foreign server WRAPPER.oid Text Type of the server (optional) N/A srvtype srvversion Text N/A Version of the server (optional) srvacl aclitem[] N/A Access permissions srvoptions Text[] N/A Option used for foreign servers. It is a keyword=value string.

Table 14-37 PG_FOREIGN_SERVER columns

14.2.37 PG_FOREIGN_TABLE

PG_FOREIGN_TABLE records auxiliary information about foreign tables.

Table 14-38 PG FOREIGN TABLE columns

Column	Туре	Description
ftrelid	OID	OID of the foreign table
ftserver	OID	OID of the server where the foreign table is located
ftwriteonly	boolean	Whether data can be written in the foreign table
ftoptions	Text[]	Foreign table options

14.2.38 PG_INDEX

PG_INDEX records part of the information about indexes. The rest is mostly in **PG_CLASS**.

Table 14-39 PG_INDEX columns

Column	Туре	Description
indexrelid	OID	OID of the pg_class entry for this index
indrelid	OID	OID of the pg_class entry for the table this index is for
indnatts	Smallint	Number of columns in an index
indisunique	boolean	This index is a unique index if the value is true .
indisprimary	boolean	This index represents the primary key of the table if the value is true . If this value is true , the value of indisunique is true.
indisexclusion	boolean	This index supports exclusion constraints if the value is true .
indimmediate	boolean	A uniqueness check is performed upon data insertion if the value is true .
indisclustered	boolean	The table was last clustered on this index if the value is true .
indisusable	boolean	This index supports insert/select if the value is true .
indisvalid	boolean	This index is valid for queries if the value is true . If this column is false , this index is possibly incomplete and must still be modified by INSERT/UPDATE operations, but it cannot safely be used for queries. If it is a unique index, the uniqueness property is also not true.
indcheckxmin	boolean	If the value is true , queries must not use the index until the xmin of this row in pg_index is below their TransactionXmin event horizon, because the table may contain broken HOT chains with incompatible rows that they can see.
indisready	boolean	If the value is true , this index is ready for inserts. If the value is false , this index is ignored when data is inserted or modified.
indkey	int2vector	This is an array of indnatts values that indicate which table columns this index creates. For example, a value of 1 3 means that the first and the third columns make up the index key. 0 in this array indicates that the corresponding index attribute is an expression over the table columns, rather than a simple column reference.

Column	Туре	Description
indcollation	oidvector	ID of each column used by the index
indclass	oidvector	For each column in the index key, this column contains the OID of the operator class to use. For details, see PG_OPCLASS.
indoption	int2vector	Array of values that store per-column flag bits. The meaning of the bits is defined by the index's access method.
indexprs	pg_node_tr ee	Expression trees (in nodeToString() representation) for index attributes that are not simple column references. It is a list with one element for each zero entry in INDKEY . NULL if all index attributes are simple references.
indpred	pg_node_tr ee	Expression tree (in nodeToString() representation) for partial index predicate. If the index is not a partial index, the value is null.
indnullstreatment	tinyint	Processing mode of the NULL value in the unique index. This field is valid only if indisunique is set to true .
		Options:
		• 0: NULLS DISTINCT. NULL values are not equivalent and can be inserted repeatedly.
		1: NULLS NOT DISTINCT. NULL values are equivalent and cannot be inserted repeatedly.
		2: NULLS IGNORE. NULL columns are ignored during equivalent comparison. If all index columns are NULL, NULL values can be inserted repeatedly. If part of the index columns are NULL, data can be inserted only if non-null values are different.
		Default value: 0
		NOTE
		 If the current cluster was upgraded from an earlier version to 8.2.0.100, the value of this field is NULL for existing indexes. For newly created indexes, the value of this field is determined by the [NULLS [NOT] DISTINCT NULLS IGNORE] field. The default value is 0.
		 If the current cluster is newly installed and its version is 8.2.0.100, for newly created indexes, the value of this field is determined by the [NULLS [NOT] DISTINCT NULLS IGNORE] field. The default value is 0.

14.2.39 PG_INHERITS

PG_INHERITS records information about table inheritance hierarchies. There is one entry for each direct child table in the database. Indirect inheritance can be determined by following chains of entries.

Table 14-40 PG INHERITS columns

Column	Туре	Reference	Description
inhrelid	OID	PG_CLASS.oid	OID of the child table.
inhparent	OID	PG_CLASS.oid	OID of the parent table.
inhseqno	Integer	-	If there is more than one direct parent for a child table (multiple inheritances), this number tells the order in which the inherited columns are to be arranged. The count starts at 1.

14.2.40 PG_JOB_INFO

PG_JOB_INFO records the execution results of scheduled tasks. The schema of the system catalog is **dbms_om**.

Table 14-41 dbms_om.pg_job_info columns

Column	Туре	Description
job_id	Integer	Job ID
job_db	OID	OID of the database where the task is
start_time	Timestamp with zone	Task execution start time
status	character(8)	Task execution status
end_time	Timestamp with zone	Task execution end time
err_msg	Text	Task execution error information

14.2.41 PG_JOBS

PG_JOBS records detailed information about jobs created by users. Dedicated threads poll the **pg_jobs** table and trigger jobs based on scheduled job execution

time. This table belongs to the Shared Relation category. All job records are visible to all databases.

Table 14-42 PG_JOBS columns

Name	Туре	Description
job_id	Integer Job ID, primary key, unique (with a unique index)	
what	Text	Job content
log_user	OID	Username of the job creator
priv_user	OID	User ID of the job executor
job_db	OID	OID of the database where the job is executed
job_nsp	OID	OID of the namespace where a job is running
job_node	OID	CN node on which the job will be created and executed
is_broken	Boolean	Whether the current job is invalid
start_date	Timestamp without time zone	Start time of the first job execution, accurate to millisecond
next_run_date	Timestamp Scheduled time of the next job execution, accurate to millisecond zone	
failure_count	Smallint Number of consecutive failures	
interval	Text	Job execution interval
last_start_date	Timestamp without time zone	Start time of the last job execution, accurate to millisecond
last_end_date	Timestamp without time zone	End time of the last job execution, accurate to millisecond
last_suc_date	Timestamp without time zone	Start time of the last successful job execution, accurate to millisecond
this_run_date	Timestamp without time zone	Start time of the ongoing job execution, accurate to millisecond

14.2.42 PG_LANGUAGE

PG_LANGUAGE records languages that can be used to write functions or stored procedures.

Table 14-43 PG_LANGUAGE columns

Name	Туре	Reference	Description
OID	OID	N/A	Row identifier (hidden attribute; must be explicitly selected)
lanname	Name	N/A	Name of the language
lanowner	OID	PG_AUTHID .oi	Owner of the language
lanispl	boolean	N/A	The value is false for internal languages (such as SQL) and true for user-defined languages. Currently, gs_dump still uses this to determine which languages need to be dumped, but this might be replaced by a different mechanism in the future.
lanpltrusted	boolean	N/A	Its value is true if this is a trusted language, which means that it is believed not to grant access to anything outside the normal SQL execution environment. Only the initial user can create functions in untrusted languages.
lanplcallfoid	OID	PG_AUTHID.oi d	For external languages, this references the language handler, which is a special function that is responsible for executing all functions that are written in the particular language.
laninline	OID	PG_AUTHID.oi	This references a function that is responsible for executing "inline" anonymous code blocks (DO blocks). The value is 0 if inline blocks are not supported.
lanvalidator	OID	PG_AUTHID.oi d	This references a language validator function that is responsible for checking the syntax and validity of new functions when they are created. The value is 0 if no validator is provided.
lanacl	aclitem[]	N/A	Access permissions

14.2.43 PG_LARGEOBJECT

PG_LARGEOBJECT records the data making up large objects A large object is identified by an OID assigned when it is created. Each large object is broken into segments or "pages" small enough to be conveniently stored as rows in **pg_largeobject**. The amount of data per page is defined to be LOBLKSIZE (which is currently BLCKSZ/4, or typically 2 kB).

It is accessible only to users with system administrator rights.

Table 14-44 PG_LARGEOBJECT columns

Column	Туре	Reference	Description
loid	OID	PG_LARGEOBJECT_ME TADATA.oid	Identifier of the large object that includes this page.
pageno	Integer	-	Page number of this page within its large object (counting from zero).
data	bytea	-	Actual data stored in the large object. This will never be more than LOBLKSIZE bytes and might be less.

Each row of pg_largeobject holds data for one page of a large object, beginning at byte offset (pageno * LOBLKSIZE) within the object. The implementation allows sparse storage: pages might be missing, and might be shorter than LOBLKSIZE bytes even if they are not the last page of the object. Missing regions within a large object are read as 0.

14.2.44 PG_LARGEOBJECT_METADATA

PG_LARGEOBJECT_METADATA records metadata associated with large objects. The actual large object data is stored in **PG_LARGEOBJECT**.

Table 14-45 PG_LARGEOBJECT_METADATA columns

Name	Туре	Reference	Description
OID	OID	N/A	Row identifier (hidden attribute; displayed only when explicitly selected)
lomowner	OID	PG_AUTHID.oid	Owner of the large object
lomacl	aclitem[]	N/A	Access permissions

14.2.45 PG_MATVIEW

PG_MATVIEW records materialized view information about the current node.

Table 14-46 PG_MATVIEW columns

Column	Туре	Description
mvid	OID	OID of the materialized view.
build_mode	Char	Build mode of the materialized view. Id': indicates "deferred", which means that data is contained in the materialized view only when the view is refreshed for the first time. Ii': indicates "immediate", which means that the latest data is included when the materialized view is created.
refresh_method	Char	Refresh method of the materialized view. 'c': indicates complete refresh.
refresh_mode	Char	Refresh mode of the materialized view. 'd': stands for demand, indicating on- demand update.
rewrite_enable	Boolean	Indicates whether query rewriting of the materialized view is supported.
active	Boolean	Indicates whether the materialized view needs to be refreshed.
relnum	Int	Number of materialized view base tables.
start_time	timestamptz	Time when the materialized view is refreshed for the first time. If this parameter is left blank, the first refresh time is the current time plus the interval.
interval	Interval	Interval for refreshing the materialized view.
refresh_time	timestamptz	Last refresh time of the materialized view.
refresh_finish_tim e	timestamptz	End time of the last refresh of a materialized view.

14.2.46 PG_NAMESPACE

PG_NAMESPACE records the namespaces, that is, schema-related information.

Table 14-47 PG_NAMESPACE columns

Column	Туре	Description	
nspname	Name	Name of the namespace	
nspowner	OID	Owner of the namespace	
nsptimeline	Bigint	Timeline when the namespace is created on the DN This column is for internal use and valid only on the DN.	
nspacl	aclitem[]	Access permissions For details, see GRANT and REVOKE.	
permspace	Bigint	Quota of a schema's permanent tablespace	
usedspace	Bigint	Used size of a schema's permanent tablespace	
nsptype	Char	Distinguishes external schemas from common schemas. This parameter is supported only in 8.3.0 and later versions to adapt to LakeFormation features.	
		NOTE The nsptype field is added to distinguish external schemas from common schemas.	
		e indicates an external schema	
		• i indicates a common schema.	

14.2.47 PG_OBJECT

PG_OBJECT records the user creation, creation time, last modification time, and last analyzing time of objects of specified types (types existing in **object_type**).

Table 14-48 PG_OBJECT columns

Column	Туре	Description
object_oid	OID	Object identifier.

Column	Туре	Description
object_type	Char	Object type: • r indicates a table, which can be an ordinary
		table or a temporary table.
		• i indicates an index.
		s indicates a sequence.v indicates a view.
		 p indicates a stored procedure and function.
		• f indicates a foreign table.
creator	OID	ID of the creator.
ctime	Timestamp with time zone	Object creation time.
mtime	Timestamp with time zone	Time when the object was last modified. By default, the ALTER, COMMENT, GRANT/REVOKE, and TRUNCATE operations are recorded.
		object_mtime_record_mode can be used to control whether ALTER, COMMENT, GRANT/ REVOKE, and TRUNCATE operations are recorded.
last_analyze_t ime	Timestamp with time zone	Time when an object is analyzed for the last time.

NOTICE

- Only normal user operations are recorded. Operations before the object upgrade and during the **initdb** process cannot be recorded.
- **ctime** and **mtime** are the start time of the transaction.
- The time of object modification due to capacity expansion is also recorded.

14.2.48 PG_OBSSCANINFO

PG_OBSSCANINFO defines the OBS runtime information scanned in cluster acceleration scenarios. Each record corresponds to a piece of runtime information of a foreign table on OBS in a query.

Name	Туре	Referen ce	Description
query_id	Bigint	N/A	Query ID
user_id	Text	N/A	Database user who performs queries
table_name	Text	N/A	Name of a foreign table on OBS
file_type	Text	N/A	Format of files storing the underlying data
time_stamp	time_stam	N/A	Scanning start time
actual_time	Double	N/A	Scanning execution time, in seconds
file_scanned	Bigint	N/A	Number of files scanned
data_size	Double	N/A	Size of data scanned, in bytes
billing_info	Text	N/A	Reserved column

Table 14-49 PG_OBSSCANINFO columns

14.2.49 PG OPCLASS

PG_OPCLASS defines index access method operator classes.

Each operator class defines semantics for index columns of a particular data type and a particular index access method. An operator class essentially specifies that a particular operator family is applicable to a particular indexable column data type. The set of operators from the family that are actually usable with the indexed column are whichever ones accept the column's data type as their lefthand input.

Table	14-50	PG	OPCLASS	columns
-------	-------	----	----------------	---------

Name	Туре	Reference	Description
OID	OID	-	Row identifier (hidden attribute; must be explicitly selected)
opcmethod	OID	PG_AM.oid	Index access method the operator class is for
opcname	Name	-	Name of the operator class
opcnamespa ce	OID	PG_NAMESPACE.oid	Namespace to which the operator class belongs
opcowner	OID	PG_AUTHID.oid	Owner of the operator class
opcfamily	OID	PG_OPFAMILY.oid	Operator family containing the operator class
opcintype	OID	PG_TYPE.oid	Data type that the operator class indexes

Name	Туре	Reference	Description
opcdefault	boolea n	-	Whether the operator class is the default for opcintype . If it is, its value is true .
opckeytype	OID	PG_TYPE.oid	Type of data stored in index, or zero if same as opcintype

An operator class's **opcmethod** must match the **opfmethod** of its containing operator family. Also, there must be no more than one **pg_opclass** row having **opcdefault** true for any given combination of **opcmethod** and **opcintype**.

14.2.50 PG_OPERATOR

PG_OPERATOR records information about operators.

Table 14-51 PG_OPERATOR columns

Name	Туре	Reference	Description
OID	OID	N/A	Row identifier (hidden attribute; displayed only when explicitly selected)
oprname	Name	N/A	Name of the operator
oprnamespace	OID	PG_NAMESPACE.oid	OID of the namespace that contains this operator
oprowner	OID	PG_AUTHID.oid	Owner of the operator
oprkind	Char	N/A	 b: infix ("both") l: prefix ("left") r: postfix ("right")
oprcanmerge	boolean	N/A	Whether the operator supports merge joins
oprcanhash	boolean	N/A	Whether the operator supports hash joins
oprleft	OID	PG_TYPE.oid	Type of the left operand
oprright	OID	PG_TYPE.oid	Type of the right operand
oprresult	OID	PG_TYPE.oid	Type of the result
oprcom	OID	PG_OPERATOR.oid	Commutator of this operator, if any
oprnegate	OID	PG_OPERATOR.oid	Negator of this operator, if any

Name	Туре	Reference	Description
oprcode	regproc	PG_PROC.oid	Function that implements this operator
oprrest	regproc	PG_PROC.oid	Restriction selectivity estimation function for this operator
oprjoin	regproc	PG_PROC.oid	Join selectivity estimation function for this operator

14.2.51 PG OPFAMILY

PG_OPFAMILY defines operator families.

Each operator family is a collection of operators and associated support routines that implement the semantics specified for a particular index access method. Furthermore, the operators in a family are all "compatible", in a way that is specified by the access method. The operator family concept allows cross-data-type operators to be used with indexes and to be reasoned about using knowledge of access method semantics.

Table 14-52 PG OPFAMILY columns

Name	Туре	Reference	Description
OID	OID	N/A	Row identifier (hidden attribute; displayed only when explicitly selected)
opfmethod	OID	PG_AM.oid	Index method used by the operator family
opfname	Name	N/A	Name of the operator family
opfnamespac e	OID	PG_NAMESPACE.oid	Namespace of the operator family
opfowner	OID	PG_AUTHID.oid	Owner of the operator family

The majority of the information defining an operator family is not in **PG_OPFAMILY**, but in the associated **PG_AMOP**, **PG_AMPROC**, and **PG_OPCLASS**.

14.2.52 PG PARTITION

PG_PARTITION records all partitioned tables, table partitions, toast tables on table partitions, and index partitions in the database. Partitioned index information is not stored in the **PG_PARTITION** system catalog.

Table 14-53 PG_PARTITION columns

Column	Туре	Description
relname	Name	Names of the partitioned tables, table partitions, TOAST tables on table partitions, and index partitions
parttype	Char	 object type r indicates a partitioned table. p indicates a table partition. x indicates an index partition. t indicates a TOAST table.
parentid	OID	OID of the partitioned table in PG_CLASS when the object is a partitioned table or table partition OID of the partitioned index when the object is an index partition
rangenum	Integer	Reserved field.
intervalnum	Integer	Reserved field.
partstrategy	Char	Partition policy of the partitioned table. Only the following policies are supported: r indicates the range partition. v indicates the numeric partition. l: indicates the list partition.
relfilenode	OID	Physical storage locations of the table partition, index partition, and TOAST table on the table partition.
reltablespace	OID	OID of the tablespace containing the table partition, index partition, TOAST table on the table partition
relpages	Double precision	Statistics: numbers of data pages of the table partition and index partition
reltuples	Double precision	Statistics: numbers of tuples of the table partition and index partition
relallvisible	Integer	Statistics: number of visible data pages of the table partition and index partition
reltoastrelid	OID	OID of the TOAST table corresponding to the table partition
reltoastidxid	OID	OID of the TOAST table index corresponding to the table partition
indextblid	OID	OID of the table partition corresponding to the index partition

Column	Туре	Description
indisusable	boolean	Whether the index partition is available
reldeltarelid	OID	OID of a Delta table
reldeltaidx	OID	OID of the index for a Delta table
relcudescrelid	OID	OID of a CU description table
relcudescidx	OID	OID of the index for a CU description table
relfrozenxid	xid32	Frozen transaction ID To ensure forward compatibility, this column is reserved. The relfrozenxid64 column is added to record the information.
intspnum	Integer	Number of tablespaces that the interval partition belongs to
partkey	int2vector	Column number of the partition key
intervaltablespace	oidvector	Tablespace that the interval partition belongs to. Interval partitions fall in the tablespaces in the round-robin manner.
interval	Text[]	Interval value of the interval partition
boundaries	Text[]	Upper boundary of the range partition and interval partition
transit	Text[]	Transit of the interval partition
reloptions	Text[]	Storage property of a partition used for collecting online scale-out information. Same as pg_class.reloptions , it is a keyword=value string.
relfrozenxid64	Xid	Frozen transaction ID
boundexprs	pg_node_t	Partition boundary expression.
	ree	For range partitioning, it is the upper boundary expression of a partition.
		For list partitioning, it is a collection of partition boundary enumeration values.
		The pg_node_tree data is not readable. You can use the expression pg_get_expr to translate the current column into readable information. SELECT pg_get_expr(boundexprs, 0) FROM pg_partition WHERE relname = 'country_202201'; pg_get_expr
		ROW(202201, 'city1'::text), ROW(202201, 'city2'::text) (1 row)

Column	Туре	Description
relmetaversion	Xid	Metadata version information. This column is supported only by clusters of version 9.1.0 or later.

Example

Query the partition information of the partitioned table web_returns_p2.

```
CREATE TABLE web_returns_p2
  wr_returned_date_sk
  wr_returned_time_sk
                            integer,
  wr_item_sk
                       integer NOT NULL,
  wr_refunded_customer_sk integer
WITH (orientation = column)
DISTRIBUTE BY HASH (wr_item_sk)
PARTITION BY RANGE(wr_returned_date_sk)
  PARTITION p2016 START (20161231) END (20191231) EVERY (10000),
  PARTITION p0 END(maxvalue)
SELECT oid FROM pg_class WHERE relname ='web_returns_p2';
97628
SELECT relname, partitype, parentid, boundaries FROM pg_partition WHERE parentid = '97628';
  relname | parttype | parentid | boundaries
                      | 97628 |
web_returns_p2 | r
p2016_0 | p
                     | 97628 | {20161231}
| 97628 | {20171231}
                        97628 | {20161231}
p2016_1
             | p
p2016_2 | p | 97628 | {20171231}
p2016_2 | p | 97628 | {20181231}
p2016_3 | p | 97628 | {20191231}
                | 97628 | {NULL}
           | p
(6 rows)
```

14.2.53 PG_PLTEMPLATE

PG_PLTEMPLATE records template information for procedural languages.

Table 14-54 PG_PLTEMPLATE columns

Column	Туре	Description
tmplname	Name	Name of the language for which this template is used.
tmpltrusted	boolean	The value is true if the language is considered trusted.
tmpldbacreate	boolean	The value is true if the language is created by the owner of the database.
tmplhandler	Text	Name of the call handler function.

Column	Туре	Description
tmplinline	Text	Name of the anonymous block handler. If no name of the block handler exists, the value is null.
tmplvalidator	Text	Name of the verification function. If no verification function is available, the value is null.
tmpllibrary	Text	Path of the shared library that implements languages.
tmplacl	aclitem[]	Access permissions for template (not yet used).

14.2.54 PG_PROC

PG_PROC records information about functions or procedures.

Table 14-55 PG_PROC columns

Column	Туре	Description
proname	Name	Name of the function
pronamespace	OID	OID of the namespace that contains the function
proowner	OID	Owner of the function
prolang	OID	Implementation language or call interface of the function
procost	Real	Estimated execution cost
prorows	Real	Estimate number of result rows
provariadic	OID	Data type of parameter element
protransform	regproc	Simplified call method for this function
proisagg	boolean	Whether this function is an aggregate function
proiswindow	boolean	Whether this function is a window function
prosecdef	boolean	Whether this function is a security definer (such as a "setuid" function)
proleakproof	boolean	Whether this function has side effects. If no leakproof treatment is provided for parameters, the function throws errors.

Column	Туре	Description
proisstrict	boolean	The function returns null if any call parameter is null. In that case the function does not actually be called at all. Functions that are not "strict" must be prepared to process null inputs.
proretset	boolean	The function returns a set, that is, multiple values of the specified data type.
provolatile	Char	Whether the function's result depends only on its input parameters, or is affected by outside factors
		It is i for "immutable" functions, which always deliver the same result for the same inputs.
		 It is s for "stable" functions, whose results (for fixed inputs) do not change within a scan.
		 It is v for "volatile" functions, whose results may change at any time.
pronargs	Smallint	Number of parameters
pronargdefaults	Smallint	Number of parameters that have default values
prorettype	OID	OID of the returned parameter type
proargtypes	oidvecto r	Array with the data types of the function parameters. This array includes only input parameters (including INOUT parameters) and thus represents the call signature of the function.
proallargtypes	oid[]	Array with the data types of the function parameters. This array includes all parameter types (including OUT and INOUT parameters); however, if all the parameters are IN parameters, this column is null. Note that array subscripting is 1-based, whereas for historical reasons, and proargtypes is subscripted from 0.

Column	Туре	Description
proargmodes	"char"[]	Array with the modes of the function parameters.
		• i indicates IN parameters.
		o indicates OUT parameters.
		• b indicates INOUT parameters.
		• v indicates VARIADIC parameters.
		t indicates table-valued parameters.
		If all the parameters are IN parameters, this column is null. Note that subscripts of this array correspond to positions of proallargtypes not proargtypes .
proargnames	Text[]	Array that stores the names of the function parameters. Parameters without a name are set to empty strings in the array. If none of the parameters have a name, this column is null. Note that subscripts correspond to positions of proallargtypes not proargtypes .
proargdefaults	pg_node _tree	Expression tree of the default value. This is the list of PRONARGDEFAULTS elements.
prosrc	Text	A definition that describes a function or stored procedure. In an interpreting language, it is the function source code, a link symbol, a file name, or any body content specified when a function or stored procedure is created, depending on how a language or calling is used.
probin	Text	Additional information about how to call the function. Again, the interpretation is language-specific.
proconfig	Text[]	Function's local settings for run-time configuration variables.
proacl	aclitem[]	Access permissions For details, see GRANT and REVOKE.
prodefaultargpos	int2vect or	Locations of the function default values. Not only the last few parameters have default values.
fencedmode	boolean	Execution mode of a function, indicating whether a function is executed in fence or not fence mode. If the execution mode is fence, the function is executed in the fork process that is reworked. The default value is fence .

Column	Туре	Description
proshippable	boolean	Whether a function can be pushed down to DNs. The default value is false .
		 Functions of the IMMUTABLE type can always be pushed down to the DNs.
		 Functions of the STABLE or VOLATILE type can be pushed down to DNs only if their attribute is SHIPPABLE.
propackage	boolean	Indicates whether the function supports overloading, which is mainly used for the Oracle style function. The default value is false .

Examples

Query the OID of a specified function. For example, obtain the OID **1295** of the **justify_days** function.

```
SELECT oid FROM pg_proc where proname ='justify_days';
OID
-----
1295
(1 row)
```

Query whether a function is an aggregate function. For example, the **justify_days** function is a non-aggregate function.

```
SELECT proisagg FROM pg_proc where proname ='justify_days';
proisagg
------
f
(1 row)
```

14.2.55 PG_PUBLICATION

PG_PUBLICATION records all the publications created in the current database. This system catalog is supported only by clusters of version 8.2.0.100 or later.

Table 14-56 PG_PUBLICATION columns

Column	Туре	Reference	Description
OID	OID	-	Row identifier (hidden attribute; displayed only when explicitly selected)
pubname	Nam e	-	Publication name
pubowner	OID	PG_AUTHID.oid	Publication owner

Column	Туре	Reference	Description
puballtable s	Boole an	-	If its value is true , the publication includes all the tables in the database, including any tables that will be created in the future.
pubinsert	Boole an	-	If its value is true , the INSERT operation is copied for the tables in the publication.
pubupdate	Boole an	-	If its value is true , the UPDATE operation is copied for the tables in the publication.
pubdelete	Boole an	-	If its value is true , the DELETE operation is copied for the tables in the publication.
pubtruncat e	Boole an	-	If its value is true , the TRUNCATE operation is copied for the tables in the publication.

Examples

View all releases.

14.2.56 PG_PUBLICATION_NAMESPACE

PG_PUBLICATION_NAMESPACE records the mapping between publications and schemas in the current database, which is a many-to-many mapping. This system catalog is supported only by clusters of version 8.2.0.100 or later.

Table 14-57 PG_PUBLICATION_NAMESPACE columns

Name	Туре	Reference	Description
OID	OID	-	Row identifier (hidden attribute; displayed only when explicitly selected)
prpubid	OID	PG_PUBLICATION.oid	Publication OID in the mapping
pnnspid	OID	PG_NAMESPACE.oid	Schema OID in the mapping

Examples

View all mappings between publications and schemas.

14.2.57 PG_PUBLICATION_REL

PG_PUBLICATION_REL records the mapping between publications and tables in the current database, which is a many-to-many mapping. This system catalog is supported only by clusters of version 8.2.0.100 or later.

To check detailed information, you are advised to use the PG_PUBLICATION_TABLES view.

Table 14-58 PG_PUBLICATION_REL columns

Name	Туре	Reference	Description
OID	OID	-	Row identifier (hidden attribute; displayed only when explicitly selected)
prpubid	OID	PG_PUBLICATION.oid	Publication OID in the mapping
prrelid	OID	PG_CLASS.oid	OID of the mapped table

Examples

View all mappings between publications and tables.

14.2.58 PG_RANGE

PG_RANGE records information about range types.

This is in addition to the types' entries in PG_TYPE.

Table 14-59 PG_RANGE columns

Column	Туре	Reference	Description
rngtypid	OID	PG_TYPE.oid	OID of the range type

Column	Туре	Reference	Description
rngsubtype	OID	PG_TYPE.oid	OID of the element type (subtype) of this range type
rngcollation	OID	PG_COLLATION.oid	OID of the collation used for range comparisons, or 0 if none
rngsubopc	OID	PG_OPCLASS.oid	OID of the subtype's operator class used for range comparisons
rngcanonica l	regproc	PG_PROC.oid	OID of the function to convert a range value into canonical form, or 0 if none
rngsubdiff	regproc	PG_PROC.oid	OID of the function to return the difference between two element values as double precision , or 0 if none

rngsubopc (plus **rngcollation**, if the element type is collatable) determines the sort ordering used by the range type. **rngcanonical** is used when the element type is discrete.

14.2.59 PG_REDACTION_COLUMN

PG_REDACTION_COLUMN records the information about the redacted columns.

Table 14-60 PG_REDACTION_COLUMN columns

Column	Туре	Description
object_oid	OID	OID of the object to be redacted.
column_attrno	Smallint	attrno of the redacted column.
function_type	Integer	Redaction type. NOTE This column is reserved. It is used only for forward compatibility of redacted column information in earlier versions. The value can be 0 (NONE) or 1 (FULL).
function_parameters	Text	Parameters used when the redaction type is partial (reserved).

Column	Туре	Description
regexp_pattern	Text	Pattern string when the redaction type is regexp (reserved).
regexp_replace_string	Text	Replacement string when the redaction type is regexp (reserved).
regexp_position	Integer	Start and end replacement positions when the redaction type is regexp (reserved).
regexp_occurrence	Integer	Replacement times when the redaction type is regexp (reserved).
regexp_match_parameter	Text	Regular control parameter used when the redaction type is regexp (reserved).
column_description	Text	Description of the redacted column.
function_expr	pg_node_tree	Internal representation of the redaction function.
inherited	bool	Whether a redacted column is inherited from another redacted column.
policy_oid	OID	OID of the masking policy. Supported by clusters of 8.2.1.100 and later versions. It is used to search for masked column information from the metadata in the system catalog.

14.2.60 PG_REDACTION_POLICY

PG_REDACTION_POLICY records information about the object to be redacted.

Table 14-61 PG_REDACTION_POLICY columns

Column	Туре	Description
object_oid	OID	OID of the object to be redacted.
policy_name	Name	Name of the redaction policy.
enable	Boolean	Policy status (enabled or disabled). NOTE The value can be: true: enabled false: disabled
expression	pg_node_tree	Policy effective expression (for users).
policy_description	Text	Description of a policy.
inherited	Bool	Whether a redaction policy is inherited from another redaction policy.
policy_order	float4	Masking policy sequence. This field is supported by 8.2.1.100 and later cluster versions.

14.2.61 PG_RELFILENODE_SIZE

The **PG_RELFILENODE_SIZE** system catalog provides file-level space statistics. Each record in the table corresponds to a physical file on the disk and the size of the file.

Table 14-62 PG_RELFILENODE_SIZE columns

Column	Туре	Description
databasei d	OID	OID of the database that the physical file belongs to If a system catalog is shared across databases, its value is 0 .
tablespac eid	OID	Tablespace OID of the physical file
relfilenod e	OID	Serial number of the physical file
backendi d	Integer	ID of the background thread that creates the physical file. Generally, the value is -1 .

Column	Туре	Description	
type	Integer	Type of the physical file.	
		• The value 0 indicates a data file.	
		• The value 1 indicates an FSM file.	
		The value 2 indicates a VM file.	
		• The value 3 indicates a BCM file.	
		• If the value greater than 4 indicates the total size of the data file and BCM file of the column in a column- store table.	
filesize	Bigint	Size of the physical file, in bytes.	

14.2.62 PG_RLSPOLICY

PG_RLSPOLICY displays the information about row-level access control policies.

Table 14-63 PG_RLSPOLICY columns

Column	Туре	Description
polname	Name	Name of a row-level access control policy
polrelid	OID	Table OID of a row-level access control policy
polcmd	Char	SQL operations affected by a row-level access control policy. The options are *(ALL), r(SELECT), w(UPDATE), and d(DELETE).
polpermi ssive	Boolean	Type of a row-level access control policy NOTE Values of polpermissive: • true: The row-level access control policy is a permissive policy. • false: The row-level access control policy is a restrictive policy.
polroles	oid[]	OID of database user affected by a row-level access control policy
polqual	pg_node _tree	SQL condition expression of a row-level access control policy

14.2.63 PG_RESOURCE_POOL

PG_RESOURCE_POOL records the information about database resource pool.

Table 14-64 PG_RESOURCE_POOL columns

Column	Туре	Description	
respool_name	Name	Name of the resource pool	
mem_percent	Integer	Configured memory percentage. 0 indicates that the memory of the resource pool is not controlled.	
cpu_affinity	Bigint	Reserved column without an actual meaning	
control_group	Name	Name of the Cgroup where the resource pool is located	
active_statements	Integer	Maximum number of concurrent statements in the resource pool	
max_dop	Integer	Maximum number of concurrent simple jobs allowed by the resource pool. -1 and 0 indicate that there are no limitations.	
memory_limit	Name	Estimated memory upper limit for a query.	
parentid	OID	OID of the parent resource pool	
io_limits	Integer	Reserved column without an actual meaning	
io_priority	Text	Reserved column without an actual meaning	
nodegroup	Name	Name of the logical cluster associated with the resource pool. The value is installation for a non-logical cluster.	
is_foreign	Boolean	Whether the resource pool can be used for users outside the logical cluster. If it is set to true , the resource pool controls the resources of common users who do not belong to the current resource pool.	
short_acc	Boolean	Whether to enable short query acceleration for a resource pool. This function is enabled by default.	
		If short query acceleration is enabled, simple queries are controlled on the fast lane.	
		If short query acceleration is disabled, and simple queries are controlled on the slow lane.	
except_rule	Text	Exception rule associated with a resource pool. There can be multiple associated rules, which are separated by commas (,).	
weight	Integer	Resource scheduling weight. Currently, this parameter is used only for network scheduling.	

14.2.64 PG_REWRITE

PG_REWRITE records rewrite rules defined for tables and views.

Table 14-65 PG_REWRITE columns

Column	Туре	Description
rulename	Name	Name of the rule
ev_class	OID	Name of the table that uses the rule
ev_attr	Smallint	Column this rule is for (always 0 to indicate the entire table)
ev_type	Char Event type for this rule: • 1 = SELECT • 2 = UPDATE • 3 = INSERT • 4 = DELETE	
ev_enabled	Char	 Controls in which mode the rule fires O: The rule fires in "origin" and "local" modes. D: The rule is disabled. R: The rule fires in "replica" mode. A: The rule always fires.
is_instead	boolean Its value is true if the rule is an INSTEAD rule.	
ev_qual	pg_node_tr ee Expression tree (in the form of a nodeToString() representation) for the ru qualifying condition	
ev_action	pg_node_tr ee	Query tree (in the form of a nodeToString () representation) for the rule's action
state_change	Timestamp with time zone	Time when the ev_enabled field is updated. This column is available only in clusters of version 9.1.0.200 or later.

14.2.65 PG_SECLABEL

PG_SECLABEL records security labels on database objects.

See also **PG_SHSECLABEL**, which performs a similar function for security labels of database objects that are shared across a database cluster.

	_		
Name	Туре	Reference	Description
objoid	OID	Any OID column	OID of the object this security label pertains to
classoid	OID	PG_CLASS.oid	OID of the system catalog that contains the object
objsubid	Integer	N/A	For a security label on a table column, this is the column number.
provider	Text	N/A	Label provider associated with this label
label	Text	N/A	Security label applied to this object

Table 14-66 PG_SECLABEL columns

14.2.66 PG_SHDEPEND

PG_SHDEPEND records the dependency relationships between database objects and shared objects, such as roles. This information allows GaussDB(DWS) to ensure that those objects are unreferenced before attempting to delete them.

See also **PG_DEPEND**, which performs a similar function for dependencies involving objects within a single database.

Unlike most system catalogs, **PG_SHDEPEND** is shared across all databases of a cluster: there is only one copy of **PG_SHDEPEND** per cluster, not one per database.

Table 14-67 PG_SHDEPEND columns

Name	Туре	Reference	Description
dbid	OID	PG_DATABASE.oid	OID of the database the dependent object is in. The value is 0 for a shared object.
classid	OID	PG_CLASS.oid	OID of the system catalog the dependent object is in.
objid	OID	Any OID column	OID of the specific dependent object
objsubid	Integer	-	For a table column, this is the column number (the objid and classid refer to the table itself). For all other object types, this column is 0 .
refclassid	OID	PG_CLASS.oid	OID of the system catalog the referenced object is in (must be a shared catalog)

Name	Туре	Reference	Description
refobjid	OID	Any OID column	OID of the specific referenced object
deptype	Char	-	Code segment defining the specific semantics of this dependency relationship. See the following text for details.
objfile	Text	-	Path of the user-defined C function library file.

In all cases, a **pg_shdepend** entry indicates that the referenced object cannot be dropped without also dropping the dependent object. However, there are several subflavors defined by **deptype**:

- SHARED DEPENDENCY OWNER (o)
 - The referenced object (which must be a role) is the owner of the dependent object.
- SHARED_DEPENDENCY_ACL (a)
 - The referenced object (which must be a role) is mentioned in the ACL (access control list, i.e., privileges list) of the dependent object. (A **SHARED_DEPENDENCY_ACL** entry is not made for the owner of the object, since the owner will have a **SHARED DEPENDENCY OWNER** entry anyway.)
- SHARED DEPENDENCY PIN (p)

There is no dependent object. This type of entry is a signal that the system itself depends on the referenced object, and so that object must never be deleted. Entries of this type are created only by **initdb**. The columns for the dependent object contain zeroes.

14.2.67 PG_SHDESCRIPTION

PG_SHDESCRIPTION records optional comments for shared database objects. Descriptions can be manipulated with the **COMMENT** command and viewed with gsql's \d commands.

See also **PG_DESCRIPTION**, which performs a similar function for descriptions involving objects within a single database.

Unlike most system catalogs, **PG_SHDESCRIPTION** is shared across all databases of a cluster. There is only one copy of **PG_SHDESCRIPTION** per cluster, not one per database.

Table 14-68 PG_SHDESCRIPTION columns

Name	Туре	Reference	Description
objoid	OID	Any OID column	OID of the object this description pertains to

Name	Туре	Reference	Description
classoid	OID	PG_CLASS.oid	OID of the system catalog where the object resides
description	Text	N/A	Arbitrary text that serves as the description of this object

14.2.68 PG_SHSECLABEL

PG_SHSECLABEL records security labels on shared database objects. Security labels can be manipulated with the **SECURITY LABEL** command.

For an easier way to view security labels, see PG_SECLABELS.

See also **PG_SECLABEL**, which performs a similar function for security labels involving objects within a single database.

Unlike most system catalogs, **PG_SHSECLABEL** is shared across all databases of a cluster. There is only one copy of **PG_SHSECLABEL** per cluster, not one per database.

Table 14-69 PG_SHSECLABEL columns

Name	Туре	Reference	Description
objoid	OID	Any OID column	OID of the object this security label pertains to
classoid	OID	PG_CLASS.oid	OID of the system catalog where the object resides
provider	Text	N/A	Label provider associated with this label
label	Text	N/A	Security label applied to this object

14.2.69 PG_STATISTIC

PG_STATISTIC records statistics about tables and index columns in a database. It is accessible only to users with system administrator rights.

Table 14-70 PG_STATISTIC columns

Column	Туре	Description
starelid	OID	Table or index which the described column belongs to.
starelkind	Char	Type of an object.

Column	Туре	Description
staattnum	Smallint	Number of the described column in the table, starting from 1.
stainherit	boolean	Whether to collect statistics for objects that have inheritance relationship.
stanullfrac	Real	Percentage of column entries that are null.
stawidth	Integer	Average stored width, in bytes, of non-null entries.
stadistinct	Real	Number of distinct, not-null data values in the column for all DNs.
		A value greater than zero is the actual number of distinct values.
		 A value less than zero is the negative of a multiplier for the number of rows in the table. (For example, stadistinct=-0.5 indicates that values in a column appear twice on average.)
		o indicates that the number of distinct values is unknown.
stakindN	Smallint	Code number stating that the type of statistics is stored in Slot N of the pg_statistic row.
		Value range: 1 to 5
staopN	OID	Operator used to generate the statistics stored in Slot N. For example, a histogram slot shows the < operator that defines the sort order of the data. Value range: 1 to 5
stanumbers N	real[]	Numerical statistics of the appropriate type for Slot N. The value is null if the slot kind does not involve numerical values.
		Value range: 1 to 5
stavaluesN	anyarray	Column data values of the appropriate type for Slot N. The value is null if the slot type does not store any data values. Each array's element values are actually of the specific column's data type so there is no way to define these columns' type more specifically than anyarray. Value range: 1 to 5
		Value range: 1 to 5

Column	Туре	Description
stadndistinct	Real	Number of unique non-null data values in the dn1 column.
		 A value greater than zero is the actual number of distinct values.
		 A value less than zero is the negative of a multiplier for the number of rows in the table. (For example, stadistinct=-0.5 indicates that values in a column appear twice on average.)
		• 0 indicates that the number of distinct values is unknown.
staextinfo	Text	Information about extension statistics (reserved)

14.2.70 PG_STATISTIC_EXT

PG_STATISTIC_EXT records extended statistics about tables in a database. The range of extended statistics to be collected is specified by users. Only system administrators can access this system catalog.

Table 14-71 PG_STATISTIC_EXT columns

Column	Туре	Description	
starelid	OID	Table or index which the described column belongs to	
starelkind	Char	Type of an object	
stainherit	boolean	Whether to collect statistics for objects that have inheritance relationship	
stanullfrac	Real	Percentage of column entries that are null	
stawidth	Integer	Average stored width, in bytes, of non-null entries	
stadistinct	Real	Number of distinct, not-null data values in the column for all DNs	
		A value greater than zero is the actual number of distinct values.	
		 A value less than zero is the negative of a multiplier for the number of rows in the table. (For example, stadistinct=-0.5 indicates that values in a column appear twice on average.) 	
		o indicates that the number of distinct values is unknown.	

Column	Туре	Description
stadndistinct	Real	Number of unique non-null data values in the dn1 column
		 A value greater than zero is the actual number of distinct values.
		 A value less than zero is the negative of a multiplier for the number of rows in the table. (For example, stadistinct=-0.5 indicates that values in a column appear twice on average.)
		o indicates that the number of distinct values is unknown.
stakindN	Smallint	Code number stating that the type of statistics is stored in Slot N of the pg_statistic row. Value range: 1 to 5
staopN	OID	Operator used to generate the statistics stored in Slot N. For example, a histogram slot shows the < operator that defines the sort order of the data. Value range: 1 to 5
stakey	int2vector	Array of a column ID
stanumbers N	real[]	Numerical statistics of the appropriate type for Slot N. The value is null if the slot kind does not involve numerical values. Value range: 1 to 5
stavaluesN	anyarray	Column data values of the appropriate type for Slot N. The value is null if the slot type does not store any data values. Each array's element values are actually of the specific column's data type so there is no way to define these columns' type more specifically than anyarray. Value range: 1 to 5
staexprs	pg_node_ tree	Expression corresponding to the extended statistics information.

14.2.71 PG_STAT_OBJECT

Records table statistics and autovacuum efficiency information of the current DB instance, and creates indexes for the **databaseid**, **relid**, and **partid** columns. Update of this system catalog is controlled by the **enable_pg_stat_object** parameter. This system catalog is supported only by clusters of version 8.2.1 or later.

Table 14-72 PG_STAT_OBJECT columns

Column	Туре	Reference	Description
databaseid	OID	PG_DATABAS E.oid	Database OID.
relid	OID	PG_CLASS.oi d	Table OID. It is the OID of the primary table for a partitioned table.
partid	OID	PG_PARTITIO N .oid	Partition OID. If the table is not partitioned, the value is 0 .
numscans	Bigint	N/A	Number of times that sequential scans are started.
tuples_returne d	Bigint	N/A	Number of visible tuples fetched by sequential scans.
tuples_fetche d	Bigint	N/A	Number of visible tuples fetched.
tuples_inserte d	Bigint	N/A	Number of inserted records.
tuples_update d	Bigint	N/A	Number of updated records.
tuples_delete d	Bigint	N/A	Number of deleted records.
tuples_hot_up dated	Bigint	N/A	Number of HOT updates.
n_live_tuples	Bigint	N/A	Number of visible tuples.
last_autovacu um_begin_n_ dead_tuple	Bigint	N/A	Number of tuples deleted before Autovacuum is executed.
n_dead_tuples	Bigint	N/A	Number of tuples deleted after Autovacuum is successful.
changes_since _analyze	Bigint	N/A	Last data modification time after Analyze.
blocks_fetche d	Bigint	N/A	Number of selected pages.
blocks_hit	Bigint	N/A	Number of scanned pages.
cu_mem_hit	Bigint	N/A	Number of CU memory hits.
cu_hdd_sync	Bigint	N/A	Times that CUs are synchronously read from disks.

Column	Туре	Reference	Description
cu_hdd_asyn	Bigint	N/A	Times that CUs are asynchronously read from disks.
data_changed _timestamp	Timestam p with time zone	N/A	Last data modification time.
data_access_ti mestamp	Timestam p with time zone	N/A	Last access time of a table.
analyze_times tamp	Timestam p with time zone	N/A	Last Analyze time.
analyze_count	Bigint	N/A	Total number of Analyze times.
autovac_analy ze_timestamp	Timestam p with time zone	N/A	Last Autoanalyze time.
autovac_analy ze_count	Bigint	N/A	Total number of Autoanalyze times.
vacuum_times tamp	Timestam p with time zone	N/A	Time of the latest Vacuum.
vacuum_coun t	Bigint	N/A	Total number of Vacuum times.
autovac_vacu um_timestam p	Timestam p with time zone	N/A	Last Autovacuum time.
autovac_vacu um_count	Bigint	N/A	Total number of Autovacuum times.
autovacuum_s uccess_count	Bigint	N/A	Total number of successful Autovacuum operations.
last_autovacu um_time_cost	Bigint	N/A	Time spent on the latest successful Autovacuum, in microseconds.
avg_autovacu um_time_cost	Bigint	N/A	Average execution time of successful Autovacuum operations. Unit: μs.
last_autovacu um_failed_co unt	Bigint	N/A	Total number of autovacuum failures since the last successful Autovacuum.

Column	Туре	Reference	Description
last_autovacu um_trigger	Smallint	N/A	Triggering mode of the latest autovacuum, which helps maintenance personnel determine the Vacuum status.
last_autovacu um_oldestxmi n	Bigint	N/A	oldestxmin after the latest successful Autovacuum execution. If the table-level oldestxmin feature is enabled, this field records the value of oldestxmin used by the latest (AUTO)VACUUM of the table.
last_autovacu um_scan_pag es	Bigint	N/A	Number of pages last scanned by autovacuum (only for row-store tables).
last_autovacu um_dirty_pag es	Bigint	N/A	Number of pages last modified by Autovacuum (only for row-store tables).
last_autovacu um_clear_dea dtuples	Bigint	N/A	Number of dead tuples last cleared by Autovacuum (only for row-store tables)
sum_autovacu um_scan_pag es	Bigint	N/A	Total number of pages scanned by Autovacuum since database initialization (only for row-store tables).
sum_autovacu um_dirty_pag es	Bigint	N/A	Number of pages modified by Autovacuum since database initialization (only for row-store tables).
sum_autovacu um_clear_dea dtuples	Bigint	N/A	Total number of dead tuples cleared by Autovacuum since database initialization (only for row-store tables).
last_autovacu um_begin_cu_ size	Bigint	N/A	Size of the CU file before the latest Autovacuum operation (only for column-store tables)
last_autovacu um_cu_size	Bigint	N/A	Size of the CU file after the latest Autovacuum (only for column- store tables)
last_autovacu um_rewrite_si ze	Bigint	N/A	Size of the column-store file last rewritten by autovacuum (only for column-store tables).

Column	Туре	Reference	Description
last_autovacu um_clear_size	Bigint	N/A	Size of the column-store file last cleared by Autovacuum (only for column-store tables).
last_autovacu um_clear_cbtr ee_tuples	Bigint	N/A	Number of cbtree tuples last cleared by Autovacuum (only for column-store tables)
sum_autovacu um_rewrite_si ze	Bigint	N/A	Total size of column-store files rewritten by Autovacuum since database initialization (only for column-store tables).
sum_autovacu um_clear_size	Bigint	N/A	Total size of column-store files cleared by Autovacuum since database initialization (only for column-store tables).
sum_autovacu um_clear_cbtr ee_tuples	Bigint	N/A	Total number of cbtree tuples cleared by Autovacuum since database initialization (only for column-store tables).
last_autovacu um_csn	Bigint	N/A	If the table-level oldestxmin feature is enabled, this field records the CSN value corresponding to the latest oldestxmin value used by the table (AUTO)VACUUM .
last_automerg e_timestamp	Timestam p with time zone	N/A	Last automerge time (only for HStore_opt tables). This column is supported only by 9.1.0.100 and later versions.
last_automerg e_time_cost	Bigint	N/A	Time consumed by the last automerge (only for HStore_opt tables). This column is supported only by 9.1.0.100 and later versions.
last_automerg e_count	Bigint	N/A	Number of records in the last automerge (only for HStore_opt tables). This column is supported only by 9.1.0.100 and later versions.
extra1	Bigint	N/A	Reserved column 1.

14.2.72 PG_SUBSCRIPTION

PG_SUBSCRIPTION records all existing subscriptions.

publisher server.

Column	Туре	Reference	Description
OID	OID	-	Row identifier (hidden attribute; displayed only when explicitly selected)
subdbid	OID	PG_DATABASE.oid	OID of the database that the subscription belongs to
subname	Nam e	-	Name of a subscription
subowner	OID	PG_AUTHID.oid	Owner of a subscription
subenabled	Boole an	-	If it is true , the subscription is enabled and should be replicated.
subconninf o	Text	-	Information about the connection to the database at the publisher end
subslotnam e	Text	-	Name of the replication slot in the publisher database If this parameter is left blank, the value is NONE .
subpublicati ons	Text[]	-	Array of subscribed publication names. These are the references to the publications on the

Table 14-73 PG_SUBSCRIPTION columns

Examples

View all subscriptions.

14.2.73 PG_SYNONYM

PG_SYNONYM records the mapping between synonym object names and other database object names.

Table 14-74 PG_SYNONYM columns

Column	Туре	Description	
synname	Name	Synonym name.	
synnamespace	OID	OID of the namespace where the synonym is located.	
synowner	OID	Owner of a synonym, usually the OID of the user who created it.	
synobjschema	Name	Schema name specified by the associated object.	
synobjname	Name	Name of the associated object.	

14.2.74 PG_TABLESPACE

PG_TABLESPACE records tablespace information.

Table 14-75 PG_TABLESPACE columns

Column	Туре	Description
spcname	Name	Name of the tablespace.
spcowner	OID	Owner of the tablespace, usually the user who created it.
spcacl	aclitem[]	Access permissions. For details, see GRANT and REVOKE .
spcoptions	Text[]	Options of the tablespace.
spcmaxsize	Text	Maximum size of the available disk space, in bytes.

14.2.75 PG_TRIGGER

PG_TRIGGER records the trigger information.

Column	Туре	Description
tgrelid	OID	OID of the table where the trigger is located.
tgname	Name	Trigger name.
tgfoid	OID	Trigger OID.
tgtype	Smallint	Trigger type.

Column	Туре	Description
tgenabled	Char	O: The trigger fires in "origin" or "local" mode.
		D : The trigger is disabled.
		R : The trigger fires in "replica" mode.
		A : The trigger always fires.
tgisinternal	Boolean	Internal trigger ID. If the value is true, it indicates an internal trigger.
tgconstrrelid	OID	Table referenced by the integrity constraint.
tgconstrindid	OID	Index of the integrity constraint.
tgconstraint	OID	OID of the constraint trigger in pg_constraint .
tgdeferrable	Boolean	The constraint trigger is of the DEFERRABLE type.
tginitdeferred	Boolean	whether the trigger is of the INITIALLY DEFERRED type.
tgnargs	Smallint	Input parameters number of the trigger function.
tgattr	int2vector	Column ID specified by the trigger. If no column is specified, an empty array is used.
tgargs	bytea	Parameter transferred to the trigger.
tgqual	pg_node_tree	Indicates the WHEN condition of the trigger. If the WHEN condition does not exist, the value is null.

14.2.76 PG_TS_CONFIG

PG_TS_CONFIG records entries representing text search configurations. A configuration specifies a particular text search parser and a list of dictionaries to use for each of the parser's output token types.

The parser is shown in the **PG_TS_CONFIG** entry, but the token-to-dictionary mapping is defined by subsidiary entries in **PG_TS_CONFIG_MAP**.

Table 14-76 PG_TS_CONFIG columns

Name	Туре	Reference	Description
OID	OID	N/A	Row identifier (hidden attribute; displayed only when explicitly selected)
cfgname	Name	N/A	Text search configuration name

Name	Туре	Reference	Description
cfgnames pace	OID	PG_NAMESPACE.oid	OID of the namespace where the configuration resides
cfgowner	OID	PG_AUTHID.oid	Owner of the configuration
cfgparser	OID	PG_TS_PARSER.oid	OID of the text search parser for this configuration
cfoptions	Text[]	N/A	Configuration options

14.2.77 PG_TS_CONFIG_MAP

PG_TS_CONFIG_MAP records entries showing which text search dictionaries should be consulted, and in what order, for each output token type of each text search configuration's parser.

Table 14-77 PG_TS_CONFIG_MAP columns

Name	Туре	Reference	Description
mapcfg	OID	PG_TS_CONFIG.oi	OID of the PG_TS_CONFIG entry owning this map entry
maptokentype	Intege r	N/A	A token type emitted by the configuration's parser
mapseqno	Intege r	N/A	Order in which to consult this entry
mapdict	OID	PG_TS_DICT.oid	OID of the text search dictionary to consult

14.2.78 PG_TS_DICT

PG_TS_DICT records entries that define text search dictionaries. A dictionary depends on a text search template, which specifies all the implementation functions needed. The dictionary itself provides values for the user-settable parameters supported by the template.

This division of labor allows dictionaries to be created by unprivileged users. The parameters are specified by a text string **dictinitoption**, whose format and meaning vary depending on the template.

Table 14-78 PG_TS_DICT columns

Name	Туре	Reference	Description
OID	OID	N/A	Row identifier (hidden attribute; displayed only when explicitly selected)
dictname	Nam e	N/A	Text search dictionary name
dictnamespace	OID	PG_NAMESPACE.oid	OID of the namespace that contains the dictionary
dictowner	OID	PG_AUTHID.oid	Owner of the dictionary
dicttemplate	OID	PG_TS_TEMPLATE.oid	OID of the text search template for this dictionary
dictinitoption	Text	N/A	Initialization option string for the template

14.2.79 PG_TS_PARSER

PG_TS_PARSER records entries defining text search parsers. A parser splits input text into lexemes and assigns a token type to each lexeme. Since a parser must be implemented by C functions, parsers can be created only by database administrators.

Table 14-79 PG_TS_PARSER columns

Name	Туре	Reference	Description
OID	OID	N/A	Row identifier (hidden attribute; displayed only when explicitly selected)
prsname	Name	N/A	Text search parser name
prsnamespac e	OID	PG_NAMESPACE.oi	OID of the namespace that contains the parser
prsstart	regpro c	PG_PROC.oid	OID of the parser's startup function
prstoken	regpro c	PG_PROC.oid	OID of the parser's next-token function
prsend	regpro c	PG_PROC.oid	OID of the parser's shutdown function

Name	Туре	Reference	Description
prsheadline	regpro c	PG_PROC.oid	OID of the parser's headline function
prslextype	regpro c	PG_PROC.oid	OID of the parser's lextype function

14.2.80 PG_TS_TEMPLATE

PG_TS_TEMPLATE records entries defining text search templates. A template provides a framework for text search dictionaries. Since a template must be implemented by C functions, templates can be created only by database administrators.

Table 14-80 PG_TS_TEMPLATE columns

Name	Туре	Reference	Description
OID	OID	-	Row identifier (hidden attribute; must be explicitly selected)
tmplname	Name	-	Text search template name
tmplnamespac e	OID	PG_NAMESPACE.oid	OID of the namespace that contains the template
tmplinit	regpro c	PG_PROC.oid	OID of the template's initialization function
tmpllexize	regpro c	PG_PROC.oid	OID of the template's lexize function

14.2.81 PG_TYPE

PG_TYPE records the information about data types.

Table 14-81 PG_TYPE columns

Column	Туре	Description
typname	Name	Data type name
typnamesp ace	OID	OID of the namespace that contains this type
typowner	OID	Owner of this type

Column	Туре	Description
typlen	Smallint	Number of bytes in the internal representation of the type for a fixed-size type. But for a variable-length type, typlen is negative. • -1 indicates a "varlena" type (one that has a length word). • -2 indicates a null-terminated C string.
typbyval	boolean	Whether the value of this type is passed by parameter or reference of this column. If the value of TYPLEN is not 1, 2, 4, or 8, you are advised to set TYPBYVAL to false. because values of this type are always passed by reference of this column. TYPBYVAL can be false even if TYPLEN allows passing values.
typtype	Char	 b indicates a basic type. c indicates a composite type, for example, a table's row type. e indicates an enumeration type. p indicates a pseudo type. For details, see typrelid and typbasetype.
typcategory	Char	typcategory is an arbitrary classification of data types that is used by the parser to determine which implicit casts should be "preferred".
typispreferr ed	boolean	Whether data is converted. It is true if conversion is performed when data meets the conversion rules specified by TYPCATEGORY .
typisdefined	boolean	The value is true if the type is defined. The value is false if this is a placeholder entry for a not-yet-defined type. When it is false , type name, namespace, and OID are the only dependable objects.
typdelim	Char	Character that separates two values of this type when parsing array input. Note that the delimiter is associated with the array element data type, not the array data type.
typrelid	OID	If this is a composite type (see typtype), then this column points to the pg_class entry that defines the corresponding table. For a free-standing composite type, the pg_class entry does not represent a table, but it is required for the type's pg_attribute entries to link to. The value is 0 for non-composite types.

Column	Туре	Description
typelem	OID	If typelem is not 0 then it identifies another row in pg_type . The current type can be subscripted like an array yielding values of type typelem . The current type can then be subscripted like an array yielding values of type typelem . A "true" array type is variable length (typlen = -1), but some fixed-length (typlen > 0) types also have nonzero typelem , for example name and point . If a fixed-length type has a typelem , its internal representation must be some number of values of the typelem data type with no other data. Variable-length array types have a header defined by the array subroutines.
typarray	OID	Indicates that the corresponding type record is available in pg_type if the value is not 0 .
typinput	regproc	Input conversion function (text format)
typoutput	regproc	Output conversion function (text format)
typreceive	regproc	Input conversion function (binary format). If no input conversion function, the value is 0 .
typsend	regproc	output conversion function (binary format). If no output conversion function, the value is 0 .
typmodin	regproc	Type modifier input function. The value is 0 if the type does not support modifiers.
typmodout	regproc	Type modifier output function. The value is 0 if the type does not support modifiers.
typanalyze	regproc	Custom ANALYZE function. The value is 0 if the standard function is used.

Column	Туре	Description
typalign	Char	Alignment required when storing a value of this type. It applies to storage on disk as well as most representations of the value inside PostgreSQL. When multiple values are stored consecutively, such as in the representation of a complete row on disk, padding is inserted before a data of this type so that it begins on the specified boundary. The alignment reference is the beginning of the first datum in the sequence. Possible values are:
		• c : char alignment, that is, no alignment needed
		• s : short alignment (2 bytes on most machines)
		• i: int alignment (4 bytes on most machines).
		• d : double alignment (8 bytes on many machines, but by no means all)
		NOTICE For types used in system tables, the size and alignment defined in pg_type must agree with the way that the compiler lays out the column in a structure representing a table row.
typstorage	Char	typstorage tells for varlena types (those with typlen = -1) if the type is prepared for toasting and what the default strategy for attributes of this type should be. Possible values are:
		• p indicates that values are always stored plain.
		 e: Value can be stored in a "secondary" relationship (if the relation has one, see pg_class.reltoastrelid).
		• m : Values can be stored compressed inline.
		 x: Values can be stored compressed inline or stored in secondary storage.
		NOTICE m domains can also be moved out to secondary storage, but only as a last resort (e and x domains are moved first).
typenotnull	boolean	Represents a NOTNULL constraint on a type. Currently, it is used for domains only.
typbasetype	OID	If this is a domain (see typtype), then typbasetype identifies the type that this one is based on. The value is 0 if this type is not a derived type.
typtypmod	Integer	Records the typtypmod to be applied to domains' base types by domains (the value is -1 if the base type does not use typmod). The value is -1 if this type is not a domain.

Column	Туре	Description	
typndims	Integer	Number of array dimensions for a domain that is an array (that is, typbasetype is an array type; the domain's typelem matches the base type's typelem). The value is 0 for types other than domains over array types.	
typcollation	OID	Sequence rule for specified types. Sequencing is not supported if the value is 0.	
typdefaultbi n	pg_node_tr ee	nodeToString() representation of a default expression for the type if the value is non-null. Currently, this column is only used for domains.	
typdefault	Text	The value is null if a type has no associated default value. If typdefaultbin is not null, typdefault must contain a human-readable version of the default expression represented by typdefaultbin . If typdefaultbin is null and typdefault is not, then typdefault is the external representation of the type's default value, which can be fed to the type's input converter to produce a constant.	
typacl	aclitem[]	Access permissions	

14.2.82 PG_USER_MAPPING

PG_USER_MAPPING records the mappings from local users to remote.

It is accessible only to users with system administrator rights. You can use view **PG_USER_MAPPINGS** to query common users.

Table 14-82 PG_USER_MAPPING columns

Column	Туре	Reference	Description
OID	OID	-	Row identifier (hidden attribute; must be explicitly selected)
umuser	OID	PG_AUTHID.oid	OID of the local role being mapped, 0 if the user mapping is public
umserver	OID	PG_FOREIGN_SERVER.	OID of the foreign server that contains this mapping
umoptions	Text[]	-	Option used for user mapping. It is a keyword=value string.

14.2.83 PG_USER_STATUS

PG_USER_STATUS records the states of users that access to the database. It is accessible only to users with system administrator rights.

Table 14-83 PG_USER_STATUS columns

Column	Туре	Description
roloid	OID	ID of the role
failcount	Integer	Specifies the number of failed attempts.
locktime	Timestamp with time zone	Time at which the role is locked
rolstatus	Smallint	 Role state 0: normal 1 indicates that the role is locked for some time because the failed login attempts exceed the threshold 2 indicates that the role is locked by the administrator.
permspac e	Bigint	Size of the permanent table storage space used by a role in the current instance.
tempspac e	Bigint	Size of the temporary table storage space used by a role in the current instance.

14.2.84 PG_WORKLOAD_ACTION

PG_WORKLOAD_ACTION records information about **query_band**.

Table 14-84 PG_WORKLOAD_ACTION columns

Column	Туре	Description
qband	Name	query_band key-value pairs
class	Name	Class of the object associated with query_band
object	Name	Object associated with query_band
action	Name	Action of the object associated with query_band

14.2.85 PGXC_CLASS

PGXC_CLASS records the replicated or distributed information for each table.

Table 14-85 PGXC_CLASS columns

Column	Туре	Description
pcrelid	OID	Table OID
pclocatortype	Char	Locator type
		H: hash
		• M : Modulo
		N: Round Robin
		R: Replicate
pchashalgorithm	Smallint	Distributed tuple using the hash algorithm
pchashbuckets	Smallint	Value of a harsh container
pgroup	Name	Node group name
redistributed	Char	Whether a table has been redistributed
redis_order	Integer	Redistribution sequence
pcattnum	int2vector	Column number used as a distribution key
nodeoids	oidvector_ex tend	List of distributed table node OIDs
options	Text	Extension status information, which is a reserved column in the system

14.2.86 PGXC_GROUP

PGXC_GROUP records node group information. In storage-compute decoupling 3.0 version, each node group in a logical cluster is called a Virtual Warehouse (VW). At the storage KV layer, each VW corresponds to a vgroup.

Table 14-86 PGXC_GROUP columns

Column	Туре	Description
group_name	Name	Node group name

Column	Туре	Description	
in_redistribution	Char	 Whether redistribution is required n indicates that the NodeGroup is not redistributed. y indicates the source NodeGroup in redistribution. t indicates the destination NodeGroup in redistribution. s indicates that the NodeGroup will skip redistribution. 	
group_members	oidvector_ex tend	Node OID list of the node group	
group_buckets	Text	Distributed data bucket group	
is_installation	boolean	Whether to install a sub-cluster	
group_acl	aclitem[]	Access permissions	
group_kind	Char	 i indicates the installation node group, which contains all DNs. n indicates a common non-logical cluster node group. v indicates a logical cluster node group. e indicates the elastic cluster node group. r indicates a replication table node group, which can only be used to create replication tables and can contain one or more logical cluster node groups. 	
group_ckpt_csn	Xid	CSN of the last incremental extraction performed on a node group	
vgroup_id	Xid	ID of the vgroup corresponding to the node group	
vgroup_bucket_count	OID	Number of buckets in the vgroup corresponding to the node group	
group_ckpt_time	Timestamp with time zone	Physical time when the last incremental extraction is performed on a node group	
apply_kv_duration	Integer	Duration of incremental scanning in the last incremental extraction of a node group, in seconds	

Column	Туре	Description
ckpt_duration	Integer	Checkpoint duration in the last incremental extraction of a node group, in seconds

14.2.87 PGXC_NODE

PGXC_NODE records information about cluster nodes.

Table 14-87 PGXC_NODE columns

Column	Туре	Description
node_name	Name	Node name
node_type	Char	Node type C: CN D: DN
node_port	Integer	Port ID of the node
node_host	Name	Host name or IP address of a node. (If a virtual IP address is configured, its value is a virtual IP address.)
node_port1	Integer	Port number of a replication node
node_host1	Name	Host name or IP address of a replication node. (If a virtual IP address is configured, its value is a virtual IP address.)
hostis_primary	boolean	Whether a switchover occurs between the primary and the standby server on the current node
nodeis_primary	boolean	Whether the current node is preferred to execute non-query operations in the replication table
nodeis_preferre d	boolean	Whether the current node is preferred to execute queries in the replication table
node_id	Integer	Node identifier
sctp_port	Integer	Specifies the port used by the TCP proxy communication library or SCTP communication library of the primary node to listen to the data channel.

Column	Туре	Description
control_port	Integer	Specifies the port used by the TCP proxy communication library or SCTP communication library of the primary node to listen to the control channel.
sctp_port1	Integer	Specifies the port used by the TCP proxy communication library or SCTP communication library of the standby node to listen to the data channel.
control_port1	Integer	Specifies the port used by the TCP proxy communication library or SCTP communication library of the standby node to listen to the control channel.
nodeis_central	boolean	Indicates that the current node is the central node.

Examples

Query the CN and DN information of the cluster.

SELECT * FROM pgxc_nod	SELECT * FROM pgxc_node;					
node_name node_type	node_name node_type node_port node_host node_port1 node_host1 hostis_primary					
	nodeis_primary nodeis_preferred node_id					
		oort1 nodeis_central rea				
+		+	+			
+						
-+						
888802358	·	55504 localhost t	†	†		
55505 55507	0 0 f	f				
datanode2 D -905831925	55508 localhost	55508 localhost t	f	f	I	
55509 55511	0 0 f	f				
coordinator1 C	55500 localhost	55500 localhost t	f	f	1	
1938253334	·			·	·	
0 0	0 0 t	f				
datanode3 D	55542 localhost	55542 localhost t	f	f		
-1894792127						
57552 55544	0 0 f	t				
		55546 localhost t	f	f		
-1307323892						
57808 55548	0 0 f	t				
datanode5 D	55550 localhost	55550 localhost t	f	f		
1797586929						
58064 55552	0 0 f	t				
datanode6 D	55554 localhost	55554 localhost t	f	f		
587455710						
58320 55556	0 0 f	t				
datanode7 D	55558 localhost	55558 localhost t	f	f		
-1685037427						
58576 55560	0 0 f	t				
datanode8 D	55562 localhost	55562 localhost t	f	f		
-993847320						
58832 55564	0 0 f	t				
(9 rows)						

14.2.88 PLAN TABLE DATA

PLAN_TABLE_DATA stores the plan information collected by **EXPLAIN PLAN**. Different from the **PLAN_TABLE** view, the system catalog **PLAN_TABLE_DATA** stores the plan information collected by all sessions and users.

Table 14-88 PLAN_TABLE columns

Column	Туре	Description
session_id	Text	Session that inserts the data. Its value consists of a service thread start timestamp and a service thread ID. Values are constrained by NOT NULL .
user_id	OID	User who inserts the data. Values are constrained by NOT NULL .
statement_id	varchar2(30)	Query tag specified by a user
plan_id	Bigint	ID of a plan to be queried
id	Int	Node ID in a plan
operation	varchar2(30)	Operation description
options	varchar2(255)	Operation parameters
object_name	Name	Name of an operated object. It is defined by users.
object_type	varchar2(30)	Object type
object_owner	Name	User-defined schema to which an object belongs
projection	varchar2(400 0)	Returned column information

◯ NOTE

- PLAN_TABLE_DATA records data of all users and sessions on the current node. Only
 administrators can access all the data. Common users can view only their own data in
 the PLAN_TABLE view.
- Data of inactive (exited) sessions is cleaned from PLAN_TABLE_DATA by gs_clean after being stored in this system catalog for a certain period of time (5 minutes by default).
 You can also manually run gs_clean -C to delete inactive session data from the table..
- Data is automatically inserted into PLAN_TABLE_DATA after EXPLAIN PLAN is executed. Therefore, do not manually insert data into or update data in PLAN_TABLE_DATA. Otherwise, data in PLAN_TABLE_DATA may be disordered. To delete data from PLAN_TABLE_DATA, you are advised to use the PLAN_TABLE view.
- Information in the **statement_id**, **object_name**, **object_owner**, and **projection** columns is stored in letter cases specified by users and information in other columns is stored in uppercase.

14.2.89 SNAPSHOT

SNAPSHOT records the start and end time of each performance view snapshot creation. After **enable_wdr_snapshot** is set to **on**, this catalog is created and maintained by the background snapshot thread. It is accessible only to users with system administrator rights.

NOTICE

- This system catalog's schema is **dbms_om**.
- Do not modify or delete this catalog externally. Otherwise, functions related to view snapshots may not work properly.

Table 14-89 dbms_om.snapshot columns

Column	Туре	Description
snapshot_id	Name	Snapshot ID. This column is the primary key and distribution key.
start_ts	Timestamp with time zone	Snapshot start time.
end_ts	Timestamp with time zone	Snapshot end time.

14.2.90 TABLES_SNAP_TIMESTAMP

TABLES_SNAP_TIMESTAMP records the start and end time of the snapshots created for each performance view. After **enable_wdr_snapshot** is set to **on**, this catalog is created and maintained by the background snapshot thread. It is accessible only to users with system administrator rights.

Table 14-90 dbms_om.tables_snap_timestamp columns

Column	Туре	Description
snapshot_id	Name	Snapshot ID. This column is the primary key and distribution key.
db_name	Text	Name of the database to which the view belongs.
tablename	Text	View name.
start_ts	Timestamp with time zone	Snapshot start time.
end_ts	Timestamp with time zone	Snapshot end time.

NOTICE

- This system catalog's schema is **dbms_om**.
- Do not modify or delete this catalog externally. Otherwise, functions related to view snapshots may not work properly.

14.2.91 System Catalogs for Performance View Snapshot

After **enable_wdr_snapshot** is set to **on**, the background snapshot thread creates and maintains a system catalog named in the format of **SNAP_***View name* to record the snapshot result of each performance view. The following system catalogs are accessible only to users with system administrator rights:

- SNAP_PGXC_OS_RUN_INFO
- SNAP_PGXC_WAIT_EVENTS
- SNAP_PGXC_INSTR_UNIQUE_SQL
- SNAP_PGXC_STAT_BAD_BLOCK
- SNAP_PGXC_STAT_BGWRITER
- SNAP_PGXC_STAT_REPLICATION
- SNAP_PGXC_REPLICATION_SLOTS
- SNAP_PGXC_SETTINGS
- SNAP_PGXC_INSTANCE_TIME
- SNAP_GLOBAL_WORKLOAD_TRANSACTION
- SNAP_PGXC_WORKLOAD_SQL_COUNT
- SNAP_PGXC_STAT_DATABASE
- SNAP_GLOBAL_STAT_DATABASE
- SNAP_PGXC_REDO_STAT
- SNAP_GLOBAL_REDO_STAT
- SNAP_PGXC_REL_IOSTAT
- SNAP_GLOBAL_REL_IOSTAT
- SNAP_PGXC_TOTAL_MEMORY_DETAIL
- SNAP_PGXC_NODE_STAT_RESET_TIME
- SNAP_PGXC_SQL_COUNT
- SNAP_GLOBAL_TABLE_STAT
- SNAP_GLOBAL_TABLE_CHANGE_STAT
- SNAP_GLOBAL_COLUMN_TABLE_IO_STAT
- SNAP_GLOBAL_ROW_TABLE_IO_STAT

Except the new **snapshot_id** column (of the bigint type), the definitions of the other columns in these system catalogs are the same as those of the corresponding views, and the distribution key of each system catalog is **snapshot_id**.

For example, SNAP_PGXC_OS_RUN_INFO is used to record snapshots of the PGXC_OS_RUN_INFO view. The snapshot_id column is new, and other columns are the same as those of the PGXC_OS_RUN_INFO view.

NOTICE

- The schema of all above system catalogs is **dbms_om**.
- Do not modify or delete these catalogs externally. Otherwise, functions related to view snapshots may not work properly.

14.3 System Views

14.3.1 ALL_ALL_TABLES

ALL_ALL_TABLES displays the tables or views accessible to the current user.

Table 14-91 ALL_ALL_TABLES columns

Column	Туре	Description
owner	Name	Owner of the table or view
table_name	Name	Name of the table or view
tablespace_name	Name	Tablespace where the table or view is located

14.3.2 ALL_CONSTRAINTS

ALL_CONSTRAINTS displays information about constraints accessible to the current user.

Table 14-92 ALL_CONSTRAINTS columns

Column	Туре	Description
constraint_name	vcharacter varying(64)	Constraint name
constraint_type	Text	 C: Check constraint F: Foreign key constraint P: Primary key constraint U: Unique constraint.
table_name	character varying (64)	Name of constraint-related table
index_owner	character varying(64)	Owner of constraint-related index (only for the unique constraint and primary key constraint)

Column	Туре	Description
index_name	character varying(64)	Name of constraint-related index (only for the unique constraint and primary key constraint)

14.3.3 ALL_CONS_COLUMNS

ALL_CONS_COLUMNS displays information about constraint columns accessible to the current user.

Table 14-93 ALL_CONS_COLUMNS columns

Column	Туре	Description
table_name	character varying(64)	Name of constraint-related table
column_name	character varying(64)	Name of constraint-related column
constraint_name	character varying(64)	Constraint name
position	Smallint	Position of the column in the table

14.3.4 ALL_COL_COMMENTS

ALL_COL_COMMENTS displays column comments of tables and views that the current user can access.

Table 14-94 ALL_COL_COMMENTS columns

Column	Туре	Description
column_name	character varying(64)	Column name
table_name	character varying(64)	Table or view name
owner	character varying(64)	Owner of the table or view
comments	Text	Comments

14.3.5 ALL_DEPENDENCIES

ALL_DEPENDENCIES displays dependencies between functions and advanced packages accessible to the current user.

NOTICE

Currently in GaussDB(DWS), this table is empty without any record due to information constraints.

Table 14-95 ALL_DEPENDENCIES columns

Column	Туре	Description
owner	character varying(30)	Owner of the object
Name	character varying(30)	Object name
type	character varying(17)	Object type
referenced_owner	character varying(30)	Owner of the referenced object
referenced_name	character varying(64)	Name of the referenced object
referenced_type	character varying(17)	Type of the referenced object
referenced_link_name	character varying(128)	Name of the link to the referenced object
schemaid	Numeric	ID of the current schema
dependency_type	character varying(4)	Dependency type (REF or HARD)

14.3.6 ALL_IND_COLUMNS

ALL_IND_COLUMNS displays all index columns accessible to the current user.

Table 14-96 ALL_IND_COLUMNS columns

Column	Туре	Description
index_owner	character varying(64)	Index owner
index_name	character varying(64)	Index name
table_owner	character varying(64)	Table owner
table_name	character varying(64)	Table name
column_name	Name	Column name
column_position	Smallint	Position of a column in the index

14.3.7 ALL_IND_EXPRESSIONS

ALL_IND_EXPRESSIONS displays information about the expression indexes accessible to the current user.

Table 14-97 ALL_IND_EXPRESSIONS columns

Column	Туре	Description
index_owner	character varying(64)	Index owner
index_name	character varying(64)	Index name
table_owner	character varying(64)	Table owner
table_name	character varying(64)	Table name
column_expression	Text	Function-based index expression of a specified column
column_position	Smallint	Position of a column in the index

14.3.8 ALL_INDEXES

ALL_INDEXES displays information about indexes accessible to the current user.

Table 14-98 ALL_INDEXES columns

Column	Туре	Description
owner	character varying(64)	Index owner
index_name	character varying(64)	Index name
table_name	character varying(64)	Name of the table corresponding to the index
uniqueness	Text	Whether the index is unique
generated	character varying(1)	Whether the index name is generated by the system
partitioned	character(3)	Whether the index has the property of the partition table

14.3.9 ALL_OBJECTS

ALL_OBJECTS displays all database objects accessible to the current user.

Table 14-99 ALL_OBJECTS columns

Column	Туре	Description
owner	Name	Owner of the object
object_name	Name	Object name
object_id	OID	OID of the object
object_type	Name	Type of the object
namespace	OID	Namespace containing the object
created	Timestamp with time zone	Object creation time
last_ddl_time	Timestamp with time zone	Last time when the object was modified

NOTICE

For details about the value ranges of **last_ddl_time** and **last_ddl_time**, see **PG_OBJECT**.

14.3.10 ALL_PROCEDURES

ALL_PROCEDURES displays information about all stored procedures or functions accessible to the current user.

Table 14-100 ALL_PROCEDURES columns

Column	Туре	Description
owner	Name	Owner of the object
object_name	Name	Object name

14.3.11 ALL_SEQUENCES

ALL_SEQUENCES displays all sequences accessible to the current user.

• N: It is not a cycle sequence.

Column **Type** Description sequence_owner Name Owner of the sequence Name sequence_name Name of the sequence min_value **Bigint** Minimum value of the sequence max_value **Bigint** Maximum value of the sequence increment_by Bigint Value by which the sequence is incremented cycle_flag character(1) Whether the sequence is a cycle sequence. The value can be **Y** or **N**. • Y: It is a cycle sequence.

Table 14-101 ALL_SEQUENCES columns

14.3.12 ALL SOURCE

ALL_SOURCE displays information about stored procedures or functions accessible to the current user, and provides the columns defined by the stored procedures and functions.

Table 14-102 ALL_SOURCE columns

Column	Туре	Description
owner	Name	Owner of the object
Name	Name	Name of the object
type	Name	Type of the object
Text	Text	Definition of the object

14.3.13 ALL SYNONYMS

ALL_SYNONYMS displays all synonyms accessible to the current user.

Table 14-103 ALL_SYNONYMS columns

Column	Туре	Description
owner	Text	Owner of a synonym
schema_name	Text	Name of the schema to which the synonym belongs

Column	Туре	Description
synonym_name	Text	Synonym name
table_owner	Text	Owner of the associated object
table_schema_nam e	Text	Name of the schema the associated object belongs to
table_name	Text	Name of the associated object

14.3.14 ALL_TAB_COLUMNS

ALL_TAB_COLUMNS displays description of columns of the tables and views that the current user can access.

Table 14-104 ALL_TAB_COLUMNS columns

Column	Туре	Description
owner	character varying(64)	Owner of a table/view
table_name	character varying(64)	Table/View name
column_name	character varying(64)	Column name
data_type	character varying(128)	Data type of a column
column_id	Integer	Column ID generated when an object is created or a column is added
data_length	Integer	Length of the column, in bytes
avg_col_len	Numeric	Average length of a column, in bytes
nullable	bpchar	Whether the column can be empty. For the primary key constraint and non-null constraint, the value is n.
data_precision	Integer	Precision of the data type. This parameter is valid for the numeric data type and NULL for other types.
data_scale	Integer	Number of decimal places. This parameter is valid for the numeric data type and 0 for other data types.
char_length	Numeric	Length of a column, in characters. This parameter is valid only for the varchar, nvarchar2, bpchar, and char types.

Column	Туре	Description
schema	character varying(64)	Namespace that contains the table or view.
kind	Text	Type of the current record. If the column belongs to a table, the value of this column is table . If the column belongs to a view, the value of this column is view .

14.3.15 ALL_TAB_COMMENTS

ALL_TAB_COMMENTS displays comments about all tables and views accessible to the current user.

Table 14-105 ALL_TAB_COMMENTS columns

Column	Туре	Description
owner	character varying(64)	Owner of the table or view
table_name	character varying(64)	Name of the table or view
comments	Text	Comments

14.3.16 ALL_TABLES

ALL_TABLES displays all the tables accessible to the current user.

Table 14-106 ALL_TABLES columns

Column	Туре	Description
owner	character varying(64)	Owner of the table
table_name	character varying(64)	Name of the table
tablespace_name	character varying(64)	Name of the tablespace that contains the table
status	character varying(8)	Whether the current record is valid
temporary	character(1)	Whether the table is a temporary table • Y indicates that it is a temporary table.
		N indicates that it is not a temporary table.

Column	Туре	Description
dropped	character varying	Whether the current record is deleted
		YES indicates that it is deleted.
		NO indicates that it is not deleted.
num_rows	Numeric	Estimated number of rows in the table

14.3.17 ALL_USERS

ALL_USERS displays all users of the database visible to the current user, however, it does not describe the users.

Table 14-107 ALL_USERS columns

Column	Туре	Description
username	Name	Username
user_id	OID	OID of the user

14.3.18 ALL_VIEWS

ALL_VIEWS displays the description about all views accessible to the current user.

Table 14-108 ALL_VIEWS columns

Column	Туре	Description
owner	Name	Owner of the view
view_name	Name	View name
text_length	Integer	Text length of the view
Text	Text	Text in the view

14.3.19 DBA DATA FILES

DBA_DATA_FILES displays the description of database files. It is accessible only to users with system administrator rights.

Table 14-109 DBA_DATA_FILES columns

Column	Туре	Description
tablespace_name	Name	Name of the tablespace to which the file belongs
bytes	Double precision	Length of the file in bytes

14.3.20 DBA_USERS

DBA_USERS displays all user names in the database. It is accessible only to users with system administrator rights.

Table 14-110 DBA_USERS columns

Column	Туре	Description
username	character varying(64)	Username

14.3.21 DBA_COL_COMMENTS

DBA_COL_COMMENTS displays column comments in the tables and views of a database. Only users with system administrator permissions can access this view.

Column	Туре	Description
column_name	character varying(64)	Column name
table_name	character varying(64)	Table or view name
owner	character varying(64)	Owner of the table or view
comments	Text	Comments

14.3.22 DBA_CONSTRAINTS

DBA_CONSTRAINTS displays information about table constraints in database. It is accessible only to users with system administrator rights.

Column	Туре	Description
constraint_name	vcharacter varying(64)	Constraint name

Column	Туре	Description
constraint_type	Text	Constraint type
	C: Check constraint	
		• F : Foreign key constraint
		P: Primary key constraint
		U: Unique constraint.
table_name	character varying(64)	Name of constraint-related table
index_owner	character varying(64)	Owner of constraint-related index (only for the unique constraint and primary key constraint)
index_name	character varying(64)	Name of constraint-related index (only for the unique constraint and primary key constraint)

14.3.23 DBA_CONS_COLUMNS

DBA_CONS_COLUMNS displays information about constraint columns in database tables. It is accessible only to users with system administrator rights.

Column	Туре	Description
table_name	character varying(64)	Name of constraint-related table
column_name	character varying(64)	Name of constraint-related column
constraint_name	character varying(64)	Constraint name
position	Smallint	Position of the column in the table

14.3.24 DBA_IND_COLUMNS

DBA_IND_COLUMNS displays column information about all indexes in the database. It is accessible only to users with system administrator rights.

Column	Туре	Description
index_owner	character varying(64)	Index owner
index_name	character varying(64)	Index name
table_owner	character varying(64)	Table owner

Column	Туре	Description
table_name	character varying(64)	Table name
column_name	Name	Column name
column_position	Smallint	Position of a column in the index

14.3.25 DBA_IND_EXPRESSIONS

DBA_IND_EXPRESSIONS displays the information about expression indexes in the database. It is accessible only to users with system administrator rights.

Column	Туре	Description
index_owner	character varying(64)	Index owner
index_name	character varying(64)	Index name
table_owner	character varying(64)	Table owner
table_name	character varying(64)	Table name
column_expression	Text	Function-based index expression of a specified column
column_position	Smallint	Position of a column in the index

14.3.26 DBA_IND_PARTITIONS

DBA_IND_PARTITIONS displays information about all index partitions in the database. Each index partition of a partitioned table in the database, if present, has a row of records in **DBA_IND_PARTITIONS**. This view is accessible only to users with system administrator rights.

Column	Туре	Description
index_owner	character varying(64)	Name of the owner of the partitioned table index to which the index partition belongs
schema	character varying(64)	Schema of the partitioned index to which the index partition belongs
index_name	character varying(64)	Index name of the partitioned table to which the index partition belongs
partition_nam e	character varying(64)	Name of the index partition

Column	Туре	Description
index_partitio n_usable	boolean	Whether the index partition is available
high_value	Text	Boundary of the table partition corresponding to the index partition. For a range partition, the boundary is the upper boundary. For a list partition, the boundary is the boundary value set.
		Reserved field for forward compatibility. The parameter pretty_high_value is added in version 8.1.3 to record the information.
pretty_high_v alue	Text	Boundary of the table partition corresponding to the index partition. For a range partition, the boundary is the upper boundary. For a list partition, the boundary is the boundary value set.
		The query result is the instant decompilation output of the partition boundary expression. The output of this column is more detailed than that of high_value . The output information can be collation and column data type.
def_tablespac e_name	Name	Tablespace name of the index partition

14.3.27 DBA_INDEXES

DBA_INDEXES displays all indexes in the database. This view is accessible only to users with system administrator rights.

Column	Туре	Description
owner	character varying(64)	Index owner
index_name	character varying(64)	Index name
table_name	character varying(64)	Name of the table corresponding to the index
uniqueness	Text	Whether the index is unique
generated	character varying(1)	Whether the index name is generated by the system

Column	Туре	Description
partitioned	character(3)	Whether the index has the property of the partition table

14.3.28 DBA_OBJECTS

DBA_OBJECTS displays all database objects in the database. This view is accessible only to users with system administrator rights.

Column	Туре	Description
owner	Name	Owner of the object
object_name	Name	Object name
object_id	OID	OID of the object
object_type	Name	Type of the object
namespace	OID	Namespace containing the object
created	Timestamp with time zone	Object creation time
last_ddl_time	Timestamp with time zone	Last time when the object was modified

NOTICE

For details about the value ranges of **last_ddl_time** and **last_ddl_time**, see **PG_OBJECT**.

14.3.29 DBA_PART_INDEXES

DBA_PART_INDEXES displays information about all partitioned table indexes in the database. It is accessible only to users with system administrator rights.

Column	Туре	Description
index_owner	character varying(64)	Name of the owner of the partitioned table index
schema	character varying(64)	Schema of the partitioned table index
index_name	character varying(64)	Name of the partitioned table index

Column	Туре	Description
table_name	character varying(64)	Name of the partitioned table to which the partitioned table index belongs
partitioning_type	Text	Partition policy of the partitioned table NOTE Currently, only range partitioning and list partitioning are supported.
partition_count	Bigint	Number of index partitions of the partitioned table index
def_tablespace_name	Name	Tablespace name of the partitioned table index
partitioning_key_coun t	Integer	Number of partition keys of the partitioned table

14.3.30 DBA_PART_TABLES

DBA_PART_TABLES displays information about all partitioned tables in the database. It is accessible only to users with system administrator rights.

Column	Туре	Description
table_owner	character varying(64)	Name of the owner of the partitioned table
schema	character varying(64)	Schema of the partitioned table
table_name	character varying(64)	Name of the partitioned table
partitioning_type	Text	Partition policy of the partitioned table
		NOTE Currently, only range partitioning and list partitioning are supported.
partition_count	Bigint	Number of partitions of the partitioned table
def_tablespace_name	Name	Tablespace name of the partitioned table
partitioning_key_count	Integer	Number of partition keys of the partitioned table

14.3.31 DBA_PROCEDURES

DBA_PROCEDURES displays information about all stored procedures and functions in the database. This view is accessible only to users with system administrator rights.

Column	Туре	Description
owner	character varying(64)	Owner of the stored procedure or the function
object_name	character varying(64)	Name of the stored procedure or the function
argument_number	Smallint	Number of the input parameters in the stored procedure

14.3.32 DBA_SEQUENCES

DBA_SEQUENCES displays information about all sequences in the database. This view is accessible only to users with system administrator rights.

Column	Туре	Description
sequence_owner	character varying(64)	Owner of the sequence
sequence_name	character varying(64)	Name of the sequence

14.3.33 DBA_SOURCE

DBA_SOURCE displays all stored procedures or functions in the database, and it provides the columns defined by the stored procedures or functions. It is accessible only to users with system administrator rights.

Column	Туре	Description
owner	character varying(64)	Owner of the stored procedure or the function
Name	character varying(64)	Name of the stored procedure or the function
Text	Text	Definition of the stored procedure or the function

14.3.34 DBA_SYNONYMS

DBA_SYNONYMS displays all synonyms in the database. It is accessible only to users with system administrator rights.

Table 14-111 DBA_SYNONYMS columns

Column	Туре	Description
owner	Text	Owner of a synonym
schema_name	Text	Name of the schema to which the synonym belongs
synonym_name	Text	Synonym name
table_owner	Text	Owner of the associated object
table_schema_nam e	Text	Name of the schema the associated object belongs to
table_name	Text	Name of the associated object

14.3.35 DBA_TAB_COLUMNS

DBA_TAB_COLUMNS stores the columns of tables and views. Each column of a table in the database has a row in **DBA_TAB_COLUMNS**. Only users with system administrator permissions can access this view.

Column	Туре	Description
owner	character varying(64)	Owner of a table/view
table_name	character varying(64)	Table/View name
column_name	character varying(64)	Column name
data_type	character varying(128)	Data type of the column
column_id	Integer	Sequence number of the column when a table/view is created
data_length	Integer	Length of the column, in bytes
comments	Text	Comments
avg_col_len	Numeric	Average length of a column, in bytes
nullable	bpchar	Whether the column can be empty. For the primary key constraint and non-null constraint, the value is n.

Column	Туре	Description
data_precision	Integer	Precision of the data type. This parameter is valid for the numeric data type and NULL for other data types.
data_scale	Integer	Number of decimal places. This parameter is valid for the numeric data type and 0 for other data types.
char_length	Numeric	Length of a column, in characters. This parameter is valid only for the varchar, nvarchar2, bpchar, and char types.
schema	character varying(64)	Namespace that contains the table or view.
kind	Text	Type of the current record. If the column belongs to a table, the value of this column is table . If the column belongs to a view, the value of this column is view .

14.3.36 DBA_TAB_COMMENTS

DBA_TAB_COMMENTS displays comments about all tables and views in the database. It is accessible only to users with system administrator rights.

Column	Туре	Description
owner	character varying(64)	Owner of the table or view
table_name	character varying(64)	Name of the table or view
comments	Text	Comments

14.3.37 DBA_TAB_PARTITIONS

DBA_TAB_PARTITIONS displays information about all partitions in the database.

Column	Туре	Description
table_owner	character varying(64)	Owner of the table that contains the partition
schema	character varying(64)	Schema of the partitioned table
table_name	character varying(64)	Table name
partition_name	character varying(64)	Name of the partition

Column	Туре	Description
high_value	Text	Upper boundary of a range partition or boundary value set of a list partition
		Reserved field for forward compatibility. The parameter pretty_high_value is added in version 8.1.3 to record the information.
pretty_high_valu e	Text	Upper boundary of a range partition or boundary value set of a list partition
		The query result is the instant decompilation output of the partition boundary expression. The output of this column is more detailed than that of high_value. The output information can be collation and column data type.
tablespace_name	Name	Name of the tablespace that contains the partition

Example

View the partition information of a partitioned table:

```
CREATE TABLE web_returns_p1
   wr_returned_date_sk
                                        integer,
   wr_returned_time_sk
                                       integer,
   wr_item_sk integer NOT NULL,
   wr_refunded_customer_sk integer
WITH (orientation = column)
DISTRIBUTE BY HASH (wr_item_sk)
PARTITION BY RANGE (wr_returned_date_sk)
   PARTITION p2016 VALUES LESS THAN(20161231),
   PARTITION p2017 VALUES LESS THAN(20171231),
   PARTITION p2018 VALUES LESS THAN (20181231),
   PARTITION p2019 VALUES LESS THAN(20191231),
   PARTITION p2020 VALUES LESS THAN (maxvalue)
SELECT * FROM dba_tab_partitions where table_name='web_returns_p1';
table_owner | schema | table_name | partition_name | high_value | pretty_high_value | tablespace_name

        dbadmin
        | public | web_returns_p1 | p2016
        | 20161231 | 20161231
        | DEFAULT TABLESPACE

        dbadmin
        | public | web_returns_p1 | p2017
        | 20171231 | 20171231
        | DEFAULT TABLESPACE

        dbadmin
        | public | web_returns_p1 | p2018
        | 20181231 | 20181231
        | DEFAULT TABLESPACE

        dbadmin
        | public | web_returns_p1 | p2019
        | 20191231 | 20191231
        | DEFAULT TABLESPACE

        dbadmin
        | public | web_returns_p1 | p2020
        | MAXVALUE | MAXVALUE
        | DEFAULT

dbadmin | public | web_returns_p1 | p2020
                                                                           MAXVALUE | MAXVALUE
                                                                                                                           | DEFAULT
TABLESPACE
(5 rows)
```

14.3.38 DBA_TABLES

DBA_TABLES displays all tables in the database. This view is accessible only to users with system administrator rights.

Column	Туре	Description
owner	character varying(64)	Table owner
table_name	character varying(64)	Table name
tablespace_name	character varying(64)	Name of the tablespace that contains the table
status	character varying(8)	Whether the current record is valid
temporary	character(1)	 Whether the table is a temporary table Y indicates that it is a temporary table. N indicates that it is not a temporary table.
dropped	character varying	 Whether the current record is deleted YES indicates that it is deleted. NO indicates that it is not deleted.
num_rows	Numeric	Estimated number of rows in the table

14.3.39 DBA_TABLESPACES

DBA_TABLESPACES displays information about available tablespaces. It is accessible only to users with system administrator rights.

Table 14-112 DBA_TABLESPACES columns

Column	Туре	Description
tablespace_name	character varying(64)	Name of the tablespace

14.3.40 DBA_TRIGGERS

DBA_TRIGGERS displays information about triggers in the database. This view is accessible only to users with system administrator rights.

Column	Туре	Description
trigger_name	character varying(64)	Trigger name
table_name	character varying(64)	Name of the table that defines the trigger
table_owner	character varying(64)	Owner of the table that defines the trigger

14.3.41 DBA_VIEWS

DBA_VIEWS displays views in the database. This view is accessible only to users with system administrator rights.

Column	Туре	Description	
owner	character varying(64)	Owner of the view	
view_name	character varying(64)	View name	

14.3.42 DUAL

DUAL is automatically created by the database based on the data dictionary. It has only one text column in only one row for storing expression calculation results. It is accessible to all users.

Table 14-113 DUAL columns

Column	Туре	Description
dummy	Text	Expression calculation result

14.3.43 GET_ALL_TSC_INFO

Obtains the TSC information of all nodes again. This view is supported only by clusters of version 8.2.1 or later.

Table 14-114 show_tsc_info() return columns

Column	Туре	Description
node_name	text	Node name
tsc_mult	bigint	TSC conversion multiplier
tsc_shift	bigint	TSC conversion shifts

Column	Туре	Description
tsc_frequency	float8	TSC frequency
tsc_use_freque ncy	boolean	Indicates whether to use the TSC frequency for time conversion.
tsc_ready	boolean	Indicates whether the TSC frequency can be used for time conversion
tsc_scalar_erro r_info	text	Error information about obtaining TSC conversion information
tsc_freq_error_ info	text	Error information about obtaining TSC frequency information

14.3.44 GET_TSC_INFO

Obtains the TSC information of the current node again. This view is supported only by clusters of version 8.2.1 or later.

Table 14-115 show_tsc_info() return columns

Column	Туре	Description
node_name	text	Node name
tsc_mult	bigint	TSC conversion multiplier
tsc_shift	bigint	TSC conversion shifts
tsc_frequency	float8	TSC frequency
tsc_use_freque ncy	boolean	Indicates whether to use the TSC frequency for time conversion.
tsc_ready	boolean	Indicates whether the TSC frequency can be used for time conversion
tsc_scalar_erro r_info	text	Error information about obtaining TSC conversion information
tsc_freq_error_ info	text	Error information about obtaining TSC frequency information

14.3.45 GLOBAL_COLUMN_TABLE_IO_STAT

GLOBAL_COLUMN_TABLE_IO_STAT provides I/O statistics of all column-store tables in the current database. The names, types, and sequences of the columns in the view are the same as those in the **GS_COLUMN_TABLE_IO_STAT** view. For details about the columns, see **Table 14-116**. The value of each statistical column is the sum of the values of the corresponding columns of all nodes.

Column	Туре	Description
schemaname	Name	Namespace of a table
relname	Name	Table name
heap_read	Bigint	Number of blocks logically read in the heap
heap_hit	Bigint	Number of block hits in the heap
idx_read	Bigint	Number of blocks logically read in the index
idx_hit	Bigint	Number of block hits in the index
cu_read	Bigint	Number of logical reads in the Compression Unit
cu_hit	Bigint	Number of hits in the Compression Unit
cidx_read	Bigint	Number of indexes logically read in the Compression Unit
cidx_hit	Bigint	Number of index hits in the Compression Unit

Table 14-116 GS_COLUMN_TABLE_IO_STAT columns

14.3.46 GLOBAL_REDO_STAT

GLOBAL_REDO_STAT displays the total statistics of XLOG redo operations on all nodes in a cluster. Except the **avgiotim** column (indicating the average redo write time of all nodes), the names of the other columns in this view are the same as those in the **PV_REDO_STAT** view. The respective meanings of the other columns are the sum of the values of the same columns in the **PV_REDO_STAT** view on each node.

Table 14-117 GLOBAL_REDO_STAT columns

Column	Туре	Description
phywrts	Bigint	Total number of physical writes on all nodes
phyblkwrt	Bigint	Total number of physical write blocks on all nodes
writetim	Bigint	Total physical write time of all nodes
avgiotim	Bigint	Average redo write time of all nodes
lstiotim	Bigint	Sum of the last write time of all nodes
miniotim	Bigint	Sum of the minimum write time of all nodes
maxiowtm	Bigint	Sum of the maximum write time of all nodes

Ⅲ NOTE

This view is accessible only to users with system administrator rights.

14.3.47 GLOBAL_REL_IOSTAT

GLOBAL_REL_IOSTAT displays the total disk I/O statistics of all nodes in a cluster. The name of each column in this view is the same as that in the **GS_REL_IOSTAT** view, but the column meaning is the sum of the value of the same column in the **GS_REL_IOSTAT** view on each node.

Table 14-118 GLOBAL_REL_IOSTAT columns

Column	Туре	Description
phyrds	Bigint	Total number of disk read times of all nodes
phywrts	Bigint	Total number of disk write times of all nodes
phyblkrd	Bigint	Total number of disk pages read by all nodes
phyblkwrt	Bigint	Total number of disk pages written by all nodes

◯ NOTE

This view is accessible only to users with system administrator rights.

14.3.48 GLOBAL ROW TABLE IO STAT

GLOBAL_ROW_TABLE_IO_STAT provides I/O statistics of all row-store tables in the current database. The names, types, and sequences of the columns in the view are the same as those in the **GS_ROW_TABLE_IO_STAT** view. For details about the columns, see **Table 14-119**. The value of each statistical column is the sum of the values of the corresponding columns of all nodes.

Table 14-119 GS_ROW_TABLE_IO_STAT columns

Column	Туре	Description
schemaname	Name	Namespace of a table
relname	Name	Name of a table
heap_read	Bigint	Number of blocks logically read in the heap
heap_hit	Bigint	Number of block hits in the heap
idx_read	Bigint	Number of blocks logically read in the index
idx_hit	Bigint	Number of block hits in the index
toast_read	Bigint	Number of blocks logically read in the TOAST table

Column	Туре	Description
toast_hit	Bigint	Number of block hits in the TOAST table
tidx_read	Bigint	Number of indexes logically read in the TOAST table
tidx_hit	Bigint	Number of index hits in the TOAST table

14.3.49 GLOBAL_STAT_DATABASE

GLOBAL_STAT_DATABASE displays the status and statistics of databases on all nodes in a cluster.

- When you query the GLOBAL_STAT_DATABASE view on a CN, the respective values of all columns returned, except stats_reset (indicating the status reset time on the current CN), are the sum of values on related nodes in the cluster. Note that the sum range varies depending on the logical meaning of each column in the GLOBAL_STAT_DATABASE view.
- When you query the **GLOBAL_STAT_DATABASE** view on a DN, the query result is the same as that in **Table 14-120**.

Table 14-120 GLOBAL_STAT_DATABASE columns

Column	Туре	Description	Sum Range
datid	OID	Database OID	-
datname	Name	Database name	-
numbackends	Integer	Number of backends currently connected to this database on the current node. This is the only column in this view that reflects the current state value. All columns return the accumulated value since the last reset.	CN
xact_commit	Bigint	Number of transactions in this database that have been committed on the current node	CN
xact_rollback	Bigint	Number of transactions in this database that have been rolled back on the current node	CN
blks_read	Bigint	Number of disk blocks read in this database on the current node	DN

Column	Туре	Description	Sum Range
blks_hit	Bigint	Number of disk blocks found in the buffer cache on the current node, that is, the number of blocks hit in the cache. (This only includes hits in the GaussDB(DWS) buffer cache, not in the file system cache.)	DN
tup_returned	Bigint	Number of rows returned by queries in this database on the current node	DN
tup_fetched	Bigint	Number of rows fetched by queries in this database on the current node	DN
tup_inserted	Bigint	Number of rows inserted in this database on the current node	DN
tup_updated	Bigint	Number of rows updated in this database on the current node	DN
tup_deleted	Bigint	Number of rows deleted from this database on the current node	DN
conflicts	Bigint	Number of queries canceled due to database recovery conflicts on the current node (conflicts occurring only on the standby server). For details, see PG_STAT_DATABASE_CONFLICTS.	CN and DN
temp_files	Bigint	Number of temporary files created by this database on the current node. All temporary files are counted, regardless of why the temporary file was created (for example, sorting or hashing), and regardless of the log_temp_files setting.	DN
temp_bytes	Bigint	Size of temporary files written to this database on the current node. All temporary files are counted, regardless of why the temporary file was created, and regardless of the log_temp_files setting.	DN
deadlocks	Bigint	Number of deadlocks in this database on the current node	CN and DN

Column	Туре	Description	Sum Range
blk_read_time	Double precision	Time spent reading data file blocks by backends in this database on the current node, in milliseconds	DN
blk_write_tim e	Double precision	Time spent writing into data file blocks by backends in this database on the current node, in milliseconds	DN
stats_reset	Timestamp with time zone	Time when the database statistics are reset on the current node	-

14.3.50 GLOBAL_TABLE_CHANGE_STAT

GLOBAL_TABLE_CHANGE_STAT displays the changes of all tables (excluding foreign tables) in the current database. The value of each column that indicates the number of times is the accumulated value since the instance was started.

Table 14-121 GLOBAL_TABLE_CHANGE_STAT columns

Column	Туре	Description
schemaname	Name	Namespace of a table
relname	Name	Table name
last_vacuum	Timestamp with time zone	Time when the last VACUUM operation is performed manually
vacuum_count	Bigint	Number of times of manually performing the VACUUM operation. The value is the sum of the number of times on each CN.
last_autovacuum	Timestamp with time zone	Time when the last VACUUM operation is performed automatically
autovacuum_cou nt	Bigint	Number of times of automatically performing the VACUUM operation. The value is the sum of the number of times on each CN.
last_analyze	Timestamp with time zone	Time when the ANALYZE operation is performed (both manually and automatically)

Column	Туре	Description
analyze_count	Bigint	Number of times of performing the ANALYZE operation (both manually and automatically). The ANALYZE operation is performed on all CNs at the same time. Therefore, the value of this column is the maximum value on all CNs.
last_autoanalyze	Timestamp with time zone	Time when the last ANALYZE operation is performed automatically
autoanalyze_cou nt	Bigint	Number of times of automatically performing the ANALYZE operation. The value is the sum of the number of times on each CN.
last_change	Bigint	Time when the last modification (INSERT, UPDATE, or DELETE) is performed

14.3.51 GLOBAL_TABLE_STAT

GLOBAL_TABLE_STAT displays statistics about all tables (excluding foreign tables) in the current database. The values of **live_tuples** and **dead_tuples** are real-time values, and the values of other statistical columns are accumulated values since the instance was started.

Table 14-122 GLOBAL_TABLE_STAT columns

Column	Туре	Description
schemaname	Name	Namespace of a table
relname	Name	Table name
distribute_mode	Char	Distribution mode of a table. The meaning of this column is the same as that of the pclocatortype column in the pgxc_class system catalog.
seq_scan	Bigint	Number of sequential scans. For a partitioned table, the sum of the number of scans of each partition is displayed.
seq_tuple_read	Bigint	Number of rows scanned in sequence.
index_scan	Bigint	Number of index scans.
index_tuple_read	Bigint	Number of rows scanned by the index.

Column	Туре	Description
tuple_inserted	Bigint	Number of rows inserted. For a replication table, the maximum value of each node is displayed. For a distribution table, the sum of all nodes is displayed.
tuple_updated	Bigint	Number of rows updated. For a replication table, the maximum value of each node is displayed. For a distribution table, the sum of all nodes is displayed.
tuple_deleted	Bigint	Number of rows deleted. For a replication table, the maximum value of each node is displayed. For a distribution table, the sum of all nodes is displayed.
tuple_hot_update d	Bigint	Number of rows with HOT updates. For a replication table, the maximum value of each node is displayed. For a distribution table, the sum of all nodes is displayed.
live_tuples	Bigint	Number of live tuples. The maximum value of each node is displayed. For a distribution table, the sum of all nodes is displayed. This indicator applies only to row store tables.
		This indicator applies only to row-store tables.
dead_tuples	Bigint	Number of dead tuples. The maximum value of each node is displayed. For a distribution table, the sum of all nodes is displayed.
		This indicator applies only to row-store tables.

14.3.52 GLOBAL_WORKLOAD_SQL_COUNT

GLOBAL_WORKLOAD_SQL_COUNT displays statistics on the number of SQL statements executed in all workload Cgroups in a cluster, including the number of **SELECT, UPDATE, INSERT**, and **DELETE** statements and the number of DDL, DML, and DCL statements.

Table 14-123 GLOBAL_WORKLOAD_SQL_COUNT columns

Column	Туре	Description
workload	Name	Workload Cgroup name
select_count	Bigint	Number of SELECT statements
update_count	Bigint	Number of UPDATE statements

Column	Туре	Description
insert_count	Bigint	Number of INSERT statements
delete_count	Bigint	Number of DELETE statements
ddl_count	Bigint	Number of DDL statements
dml_count	Bigint	Number of DML statements
dcl_count	Bigint	Number of DCL statements

14.3.53 GLOBAL_WORKLOAD_SQL_ELAPSE_TIME

GLOBAL_WORKLOAD_SQL_ELAPSE_TIME displays statistics on the response time of SQL statements in all workload Cgroups in a cluster, including the maximum, minimum, average, and total response time of **SELECT**, **UPDATE**, **INSERT**, and **DELETE** statements. The unit is microsecond.

Table 14-124 GLOBAL_WORKLOAD_SQL_ELAPSE_TIME columns

Column	Туре	Description
workload	Name	Workload Cgroup name
total_select_elapse	Bigint	Total response time of SELECT statements
max_select_elapse	Bigint	Maximum response time of SELECT statements
min_select_elapse	Bigint	Minimum response time of SELECT statements
avg_select_elapse	Bigint	Average response time of SELECT statements
total_update_elapse	Bigint	Total response time of UPDATE statements
max_update_elapse	Bigint	Maximum response time of UPDATE statements
min_update_elapse	Bigint	Minimum response time of UPDATE statements
avg_update_elapse	Bigint	Average response time of UPDATE statements

Column	Туре	Description
total_insert_elapse	Bigint	Total response time of INSERT statements
max_insert_elapse	Bigint	Maximum response time of INSERT statements
min_insert_elapse	Bigint	Minimum response time of INSERT statements
avg_insert_elapse	Bigint	Average response time of INSERT statements
total_delete_elapse	Bigint	Total response time of DELETE statements
max_delete_elapse	Bigint	Maximum response time of DELETE statements
min_delete_elapse	Bigint	Minimum response time of DELETE statements
avg_delete_elapse	Bigint	Average response time of DELETE statements

14.3.54 GLOBAL_WORKLOAD_TRANSACTION

GLOBAL_WORKLOAD_TRANSACTION provides the total transaction information about workload Cgroups on all CNs in the cluster. This view is accessible only to users with system administrator rights. It is valid only when the real-time resource monitoring function is enabled, that is, **enable_resource_track** is **on**.

Table 14-125 GLOBAL_WORKLOAD_TRANSACTION columns

Column	Туре	Description
workload	Name	Workload Cgroup name
commit_counter	Bigint	Total number of submission times on each CN
rollback_counter	Bigint	Total number of rollback times on each CN
resp_min	Bigint	Minimum response time of the cluster
resp_max	Bigint	Maximum response time of the cluster
resp_avg	Bigint	Average response time on each CN
resp_total	Bigint	Total response time on each CN

14.3.55 GS ALL CONTROL GROUP INFO

GS_ALL_CONTROL_GROUP_INFO displays all Cgroup information in a database.

Table 14-126 GS_ALL_CONTROL_GROUP_INFO columns

Column	Туре	Description
Name	Text	Name of the Cgroup
type	Text	Type of the Cgroup
gid	Bigint	Cgroup ID
classgid	Bigint	ID of the Class Cgroup to which a Workload belongs
class	Text	Class Cgroup
workload	Text	Workload Cgroup
shares	Bigint	CPU quota allocated to a Cgroup
limits	Bigint	Limit of CPUs allocated to a Cgroup
wdlevel	Bigint	Workload Cgroup level
cpucores	Text	Usage of CPU cores in a Cgroup

14.3.56 GS_BLOCKLIST_QUERY

GS_BLOCKLIST_QUERY is used to query job blocklist and exception information. This view is obtained by associating system catalogs **GS_BLOCKLIST_QUERY** and **GS_WLM_SESSION_INFO**, and deduplicating query results. If the **GS_WLM_SESSION_INFO** table is large, the query may take a long time.

NOTICE

- The schema of the GS BLOCKLIST QUERY view is pg catalog.
- The **GS_BLOCKLIST_QUERY** view can be queried only in the **postgres** database. If it is gueried in other databases, an error is reported.
- Generally, constant values are ignored during unique SQL ID calculation in DML statements. However, constant values cannot be ignored in DDL, DCL, and parameter setting statements. A unique_sql_id may correspond to one or more queries.

Column Type Referenc Description unique_sql_id Bigint N/A Query ID generated based on the query parsing tree. block list Boolean N/A Check whether a job is in the blocklist. Integer N/A Query the number of job except_num exceptions. Timestamp N/A Query the time when the last job except_time exception occurred. Text N/A Statement to be executed. query

Table 14-127 GS_BLOCKLIST_QUERY columns

14.3.57 GS_BLOCKLIST_SQL

GS_BLOCKLIST_SQL is used to query job blocklist and exception information. This view is obtained by associating system catalogs **GS_BLOCKLIST_SQL** and **GS_WLM_SESSION_INFO**, and deduplicating query results. If the **GS_WLM_SESSION_INFO** table is large, the query may take a long time.

This view is supported only by 9.1.0.200 and later cluster versions.

NOTICE

- The schema stored in the GS_BLOCKLIST_SQL view is pg_catalog.
- The **GS_BLOCKLIST_SQL** view can be queried only in the **postgres** database. If it is gueried in other databases, an error is reported.
- Generally, constant values are ignored during sql_hash calculation in DML statements. However, constant values cannot be ignored in DDL, DCL, and parameter setting statements. A sql_hash may correspond to one or more queries.

Table 14-128 GS BLOCKLIST QUERY columns

Column	Туре	Referenc e	Description
sql_hash	Text	N/A	sql_hash generated based on the query parsing tree.
block_list	Boolean	N/A	Whether a job is in the blocklist.
except_num	Integer	N/A	Number of job exceptions.

Column	Туре	Referenc e	Description
except_time	Timestamp	N/A	Time when the last job exception occurred.
query	Text	N/A	Statement to be executed.

14.3.58 GS_CLUSTER_RESOURCE_INFO

GS_CLUSTER_RESOURCE_INFO displays a DN resource summary.

Table 14-129 GS_CLUSTER_RESOURCE_INFO columns

Column	Туре	Description
min_mem_util	Integer	Minimum memory usage of a DN
max_mem_util	Integer	Maximum memory usage of a DN
min_cpu_util	Integer	Minimum CPU usage of a DN
max_cpu_util	Integer	Maximum CPU usage of a DN
min_io_util	Integer	Minimum I/O usage of a DN
max_io_util	Integer	Maximum I/O usage of a DN
used_mem_rate	Integer	Maximum physical memory usage

14.3.59 GS_COLUMN_TABLE_IO_STAT

GS_COLUMN_TABLE_IO_STAT displays the I/O of all column-store tables of the database on the current node. The value of each statistical column is the accumulated value since the instance was started.

Table 14-130 GS_COLUMN_TABLE_IO_STAT columns

Column	Туре	Description
schemaname	Name	Namespace of a table
relname	Name	Table name
heap_read	Bigint	Number of blocks logically read in the heap
heap_hit	Bigint	Number of block hits in the heap
idx_read	Bigint	Number of blocks logically read in the index
idx_hit	Bigint	Number of block hits in the index

Column	Туре	Description
cu_read	Bigint	Number of logical reads in the Compression Unit
cu_hit	Bigint	Number of hits in the Compression Unit
cidx_read	Bigint	Number of indexes logically read in the Compression Unit
cidx_hit	Bigint	Number of index hits in the Compression Unit

14.3.60 GS_OBS_READ_TRAFFIC

Collects statistics on the OBS read traffic and average read bandwidth. The statistical results are aggregated every 10 minutes. This view is supported only by clusters of version 8.2.0 or later.

Column	Туре	Description
nodename	TEXT	Cluster node
hostname	TEXT	Server node
traffic_mb	float8	OBS read traffic statistics during the 10 minutes before logtime
bandwidth_m b_per_s	float8	Average bandwidth, in MB/s
reqcount	Bigint	Number of OBS reads during the 10 minutes before logtime
logtime	Timestamp with time zone	Time when statistics are recorded

Examples

Query statistics on the OBS read traffic and average read bandwidth. The statistical results are aggregated every 10 minutes.

14.3.61 GS_OBS_WRITE_TRAFFIC

Collects statistics on the OBS write traffic and average write bandwidth. The statistical results are aggregated every 10 minutes. This view is supported only by clusters of version 8.2.0 or later.

Column	Туре	Description
nodename	TEXT	Cluster node
hostname	TEXT	Server node
traffic_mb	float8	OBS write traffic statistics during the 10 minutes before logtime
bandwidth_m b_per_s	float8	Average bandwidth, in MB/s
reqcount	Bigint	Number of OBS writes during the 10 minutes before logtime
logtime	Timestamp with time zone	Time when statistics are recorded

Examples

Query statistics on the OBS write traffic and average write bandwidth. The statistical results are aggregated every 10 minutes.

nodenan		traffic_mb bar	ndwidth_mb_per_s reqcou	1 3
			3 .000289970820362525	
16:10:00+0	08			
dn_1	rhel_10_90_45_56	.000354766845703125	5 .000386063466694153	7 2022-10-24
18:50:00+0	08			
dn_1	rhel_10_90_45_56 9	9.34600830078125e-05	5 .000143659648687162	2 2022-11-07
09:20:00+0	08			
dn_1	rhel_10_90_45_56 4	4.10079956054688e-05	5 .000186667253592502	1 2022-11-07
09:30:00+0	08			
dn_1	rhel_10_90_45_56	2048.17834663391	27.2766632219637	2 2022-11-22
16:10:00+0	08			
dn_1	rhel_10_90_45_56	3747.23722648621	28.0842938534546	4 2022-11-22
16:20:00+0	08			
(6 row)				

14.3.62 GS_INSTR_UNIQUE_SQL

Unique SQL Definition

The database parses each received SQL text string and generates an internal parsing tree. The database traverses the parsing tree and ignores constant values in the parsing tree. In this case, an integer value is calculated using a certain algorithm. This integer is used as the Unique SQL ID to uniquely identify this type of SQL. SQL statements with the same Unique SQL ID are called Unique SQL statements.

Examples

Assume that the user enters the following SQL statements in sequence:

```
select * from t1 where id = 1;
select * from t1 where id = 2;
```

The statistics of the two SQL statements are aggregated to the same Unique SQL statement.

select * from t1 where id = ?;

GS_INSTR_UNIQUE_SQL View

The **GS_INSTR_UNIQUE_SQL** view displays the execution information about the Unique SQL statements collected by the current node, including:

- Unique SQL ID and normalized SQL text string. The normalized SQL text is
 described in Examples. Generally, constant values are ignored during Unique
 SQL ID calculation in DML statements. However, constant values cannot be
 ignored in DDL, DCL, and parameter setting statements.
- Number of execution times (number of successful execution times) and response time (SQL execution time in the database, including the maximum, minimum, and total time)
- Cache or I/O information, including the number of physical reads and logical reads of blocks. Only information about successfully executed SQL statements on each DN is collected. The statistical value is related to factors such as the amount of data processed during query execution, used memory, whether the query is executed for multiple times, memory management policy, and whether there are other concurrent queries. The statistical value reflects the number of physical reads and logical reads of the buffer block in the entire query execution process. The statistical value may vary according to the execution time.
- Row activities, such as the number of returned rows, updated rows, inserted rows, deleted rows, sequentially scanned rows, and randomly scanned rows in the result set of the **SELECT** statement. Except that the number of rows returned by the result set is the same as the number of rows in the result set of the **SELECT** statement and is recorded only on the CN, the activity information of other rows is recorded on the DN. The statistical value reflects the row activities during the entire query execution process, including scanning and modifying related system tables, metadata tables, and data tables. The value of this parameter is related to the data volume and related parameter settings. That is, the statistical value is greater than or equal to the scanning and modification times of actual data tables.
- Time distribution, including DB_TIME/CPU_TIME/EXECUTION_TIME/PARSE_TIME/PLAN_TIME/REWRITE_TIME/PL_EXECUTION_TIME/PL_COMPILATION_TIME/NET_SEND_TIME/DATA_IO_TIME. For details, see Table 14-131. The information is collected on both CNs and DNs and is displayed during view query.
- Number of soft and hard parsing times, such as the number of soft parsing times (cache plan) and hard parsing times (generation plan). If the cache plan is executed this time, the number of soft parsing times increases by 1. If the generation plan is regenerated this time, the number of hard parsing times increases by 1. This number is counted on both CNs and DNs and is displayed during view query.

The Unique SQL statistics function has the following restrictions:

• Detailed statistics are displayed only for successfully executed SQL statements. Otherwise, only query, node, and user information are recorded.

- If the Unique SQL statistics collection function is enabled, the CN collects statistics on all received queries, including tool and user queries.
- If an SQL statement contains multiple SQL statements or similar stored procedures, a Unique SQL statement is generated for the outermost SQL statement. The statistics of all sub-SQL statements are summarized to the Unique SQL record.
- The response time statistics of Unique SQL does not include the time of the NET_SEND_TIME phase. Therefore, there is no comparison between EXECUTION_TIME and elapse_time.
- parse_time of clauses cannot be calculated for begin;...;commit and similar transaction blocks.

When a common user accesses the **GS_INSTR_UNIQUE_SQL** view, only the Unique SQL information about the user is displayed. When an administrator accesses the **GS_INSTR_UNIQUE_SQL** view, all Unique SQL information about the current node is displayed. The **GS_INSTR_UNIQUE_SQL** view can be queried on both CNs and DNs. The DN displays the Unique SQL statistics of the local node, and the CN displays the complete Unique SQL statistics of the local node. That is, the CN collects the Unique SQL execution information of the CN from other CNs and DNs and displays the information. You can query the **GS_INSTR_UNIQUE_SQL** view to locate the Top SQL statements that consume different resources, providing a basis for cluster tuning and maintenance.

The background thread checks all Unique SQL statements every hour and deletes the Unique SQL statements whose **last_time** is **instr_unique_sql_timeout** hours ago.

Table 14-131 GS_INSTR_UNIQUE_SQL columns

Column	Туре	Description
node_name	Name	Name of the CN that receives SQL statements
node_id	Integer Node ID, which same as the winde_id in the pgxc_node tall.	
user_name	Name	Username
user_id	OID	User ID
unique_sql_id	Bigint	Normalized Unique SQL ID
query	Text	Normalized SQL text
n_calls	Bigint	Number of successful execution times
min_elapse_time	Bigint	Minimum running time of the SQL statement in the database (unit: µs)

Column	Туре	Description
max_elapse_time	Bigint	Maximum running time of SQL statements in the database (unit: µs)
total_elapse_time	Bigint	Total running time of SQL statements in the database (unit: µs)
n_returned_rows	Bigint	Row activity - Number of rows in the result set returned by the SELECT statement
n_tuples_fetched	Bigint	Row activity - Randomly scan rows (column-store tables/foreign tables are not counted.)
n_tuples_returned	Bigint	Row activity - Sequential scan rows (Column-store tables/foreign tables are not counted.)
n_tuples_inserted	Bigint	Row activity - Inserted rows
n_tuples_updated	Bigint	Row activity - Updated rows
n_tuples_deleted	Bigint	Row activity - Deleted rows
n_blocks_fetched	Bigint	Block access times of the buffer, that is, physical read/I/O
n_blocks_hit	Bigint	Block hits of the buffer, that is, logical read/ cache
n_soft_parse	Bigint	Number of soft parsing times (cache plan)
n_hard_parse	Bigint	Number of hard parsing times (generation plan)

Column	Туре	Description
db_time	Bigint	Valid DB execution time, including the waiting time and network sending time. If multiple threads are involved in query execution, the value of DB_TIME is the sum of DB_TIME of multiple threads (unit: µs).
cpu_time	Bigint	CPU execution time, excluding the sleep time (unit: µs)
execution_time	Bigint	SQL execution time in the query executor, DDL statements, and statements (such as Copy statements) that are not executed by the executor are not counted (unit: µs).
parse_time	Bigint	SQL parsing time (unit: μs)
plan_time	Bigint	SQL generation plan time (unit: μs)
rewrite_time	Bigint	SQL rewriting time (unit: μs)
pl_execution_time	Bigint	Execution time of the plpgsql procedural language function (unit: µs)
pl_compilation_time	Bigint	Compilation time of the plpgsql procedural language function (unit: µs)
net_send_time	Bigint	Network time, including the time spent by the CN in sending data to the client and the time spent by the DN in sending data to the CN (unit: µs)
data_io_time	Bigint	File I/O time (unit: μs)

Column	Туре	Description
first_time	Timestamp with time zone	Time of the first SQL statement execution
last_time	Timestamp with time zone	Time of the last SQL statement execution

14.3.63 GS_NODE_STAT_RESET_TIME

GS_NODE_STAT_RESET_TIME provides the statistics reset time of the current node and returns a timestamp with the time zone.

For details, see the **get_node_stat_reset_time()** function.

When an instance is running, its statistics keep rising. In the following cases, the statistical values in the memory will be reset to $\mathbf{0}$:

- The instance is restarted or a cluster switchover occurs.
- The database is dropped.
- A reset operation is performed. For example, the statistics counter in the database is reset using the pgstat_recv_resetcounter function or the Unique SQL statements are cleared using the reset instr unique sql function.

If any of the preceding events occurs, GaussDB(DWS) will record the time when the statistics are reset. You can query the time using the **get_node_stat_reset_time** function.

14.3.64 GS_OBS_LATENCY

GS_OBS_LATENCY records the average latency of OBS during the 10 minutes before **logtime**. The latency is estimated based on OBS operations. This view is supported only by clusters of version 8.2.0 or later.

Table 14-132 GS_OBS_LATENCY columns

Column	Туре	Description
nodename	Text	Node
hostname	Text	Server node.
latency_ms	Double precision	Average delay of OBS during the 10 minutes before logtime . The unit is ms.
reqcount	Bigint	Number of OBS requests during the 10 minutes before logtime .
logtime	Timestamp with time zone	Time when the delay information is recorded.

14.3.65 GS_QUERY_MONITOR

Displays the running/queuing information and resource usage of ongoing queries. Only queuing and running jobs are displayed. This view can be queried only on CNs and displays only the monitoring information about the main statement. This view is supported only by clusters of 8.2.1.100 and later versions.

Table 14-133 GS_QUERY_MONITOR columns

Column	Туре	Description
usename	Name	Name of the user who performs the query.
nodename	Name	Name of the CN that executes the query.
nodegroup	Name	Name of the cluster where the query is performed. The default cluster name is installation.
rpname	Name	Name of the resource pool associated with the query.
priority	Name	Priority of the query, which can be Rush , High , Medium , and Low .
xact_start	Timestamp	Start time of the transaction to which the query belongs.
query_start	Timestamp	Start time of query execution.
block_time	Bigint	Accumulated queuing time of jobs. Stored procedures and multi-statement task may be queued for multiple times. Unit: second.
duration	Bigint	Running time of a job, excluding the queuing time. Unit: second.
query_band	Text	Job ID, which can be set using the GUC parameter query_band . By default, this parameter is left blank.
attribute	Text	Job attributes:
		Simple: simple job.
		Complicated: complex job.
		This column is invalid before a job is under resource pool management and control. This column is valid only when the job is under or has been under resource pool management and control.

Column	Туре	Description
lane	Text	Resource pool lane where a job is queued or executed:
		fast: fast lane.
		• slow: slow lane.
		This column is invalid before a job is under resource pool management and control. This column is valid only when the job is under or has been under resource pool management and control.
status	Text	Current status of a job. The value can be pending or running .
queue	Text	Job queuing information:
		None: The job is running.
		Global: The job is queued in the global queue of the CN.
		Respool: The job is queued in the resource pool.
		CCN: The job is queued in the CCN.
used_mem	Integer	Maximum peak memory usage of a job across all DNs. The unit is MB.
estimate_me m	Integer	Estimated memory of a job. The unit is MB.
used_cpu	Double precision	Average number of CPU cores occupied by a job since the job starts to run.
read_speed	Integer	Average logical I/O read rate of a job on all DNs. The unit is KB/s.
write_speed	Integer	Average logical I/O write rate of a job on all DNs. The unit is KB/s.
send_speed	Integer	Average transmit rate on all DNs since a job starts to run. The unit is KB/s.
recv_speed	Integer	Average receive rate on all DNs since a job starts to run. The unit is KB/s.
dn_count	Bigint	Number of DNs that execute the job.
stream_count	Bigint	Total number of stream threads of a job on all DNs.
pid	Bigint	ID of the backend thread
lwtid	Integer	Lightweight thread ID of a background thread.
query_id	Bigint	Query ID.

Column	Туре	Description
unique_sql_id	Bigint	ID of the normalized unique SQL.
query	Text	Query that is being executed.

14.3.66 GS_QUERY_RESOURCE_INFO

The **GS_QUERY_RESOURCE_INFO** view displays the resource information about all running jobs on the current DN. This parameter is supported only by clusters of version 9.1.0 or later.

This view can be queried only on DNs. It is used only for O&M operations to locate faults. You are advised not to use this function.

Table 14-134 GS_QUERY_RESOURCE_INFO

Column	Туре	Description
node_name	Text	Instance name, which contains only DNs
user_id	OID	User ID.
queryid	Bigint	Internal query ID used for statement execution.
used_mem	Int	Memory used by the statement on the current DN. The unit is MB.
cpu_time	Bigint	CPU time of a statement on the current DN. The unit is ms.
used_cpu	Double	Number of CPUs used by the statement on the current DN.
spill_size	Bigint	Amount of data spilled to disks on the current DN. The default value is 0 . The unit is MB.
read_bytes	Bigint	Number of logical read bytes used by the statement on the current DN. The unit is KB.
write_bytes	Bigint	Number of logical write bytes used by the statement on the current DN. The unit is KB.
read_count	Bigint	Number of logical reads used by the statement on the current DN.
write_count	Bigint	Number of logical writes used by the statement on the current DN.
read_speed	Int	Logical read rate used by the statement on the current DN. The unit is KB/s.

Column	Туре	Description
write_speed	Int	Logical write rate used by the statement on the current DN. The unit is KB/s.
curr_iops	Int	I/O operations per second of the statement on the current DN. It is recorded as a count in a column-store table and as a count of 10,000 in a row-store table.
send_pkg	Bigint	Total number of communication packages sent by a statement across all DNs.
recv_pkg	Bigint	Total number of communication packages received by a statement across all DNs.
send_bytes	Bigint	Total sent data of the statement stream, in byte.
recv_bytes	Bigint	Total received data of the statement stream, in byte.
send_speed	Int	Network sending rate of the statement on the current DN. The unit is KB/s.
recv_speed	Int	Network receiving rate of the statement on the current DN. The unit is KB/s.

14.3.67 GS_REL_IOSTAT

GS_REL_IOSTAT displays disk I/O statistics on the current node. In the current version, only one page is read or written in each read or write operation. Therefore, the number of read/write times is the same as the number of pages.

Table 14-135 GS_REL_IOSTAT columns

Column	Туре	Description
phyrds	Bigint	Number of disk reads
phywrts	Bigint	Number of disk writes
phyblkrd	Bigint	Number of read pages
phyblkwrt	Bigint	Number of written pages

14.3.68 GS_RESPOOL_RUNTIME_INFO

GS_RESPOOL_RUNTIME_INFO displays information about the running of jobs in all resource pools on the current CN.

Column	Туре	Description
nodegroup	Name	Name of the logical cluster the resource pool belongs to. The default cluster is installation .
rpname	Name	Resource pool name.
ref_count	Int	Number of jobs that reference the resource pool. This count includes both controlled and uncontrolled jobs.
fast_run	Int	Number of jobs currently running in the resource pool's fast lane.
fast_wait	Int	Number of jobs currently queued in the resource pool's fast lane.
slow_run	Int	Number of jobs currently running in the resource pool's slow lane.
slow_wait	Int	Number of jobs currently queued in the resource pool's slow lane.

Table 14-136 GS_RESPOOL_RUNTIME_INFO columns

14.3.69 GS_RESPOOL_RESOURCE_INFO

GS_RESPOOL_RESOURCE_INFO displays job running information about all resource pools on a CN and the information about resource pool usage of an instance (CN/DN).

On a DN, it only displays the monitoring information of the logical cluster that the DN belongs to.

Table 14-137 GS_RESPOOL_RESOURCE_INFO columns

Column	Туре	Description
nodegroup	Name	Name of the logical cluster of the resource pool. The default value is installation .
rpname	Name	Resource pool name
cgroup	Name	Name of the Cgroup associated with the resource pool
ref_count	Int	Number of jobs referenced by the resource pool. The number is counted regardless of whether the job is controlled by the resource pool. This parameter is valid only on CNs.

Column	Туре	Description
fast_run	Int	Number of running jobs in the fast lane of the resource pool. This parameter is valid only on CNs.
fast_wait	Int	Number of jobs queued in the fast lane of the resource pool. This parameter is valid only on CNs.
fast_limit	Int	Limit on the number of concurrent jobs in the fast lane in a resource pool. This parameter is valid only on CNs.
slow_run	Int	Number of running jobs in the slow lane of the resource pool. This parameter is valid only on CNs.
slow_wait	Int	Number of jobs queued in the slow lane of the resource pool. This parameter is valid only on CNs.
slow_limit	Int	Limit on the number of concurrent jobs in the slow lane in a resource pool. This parameter is valid only on CNs.
used_cpu	Double	Average number of CPUs used by the resource pool in a 5s monitoring period. The value is accurate to two decimal places.
		On a DN, it indicates the number of CPUs used by the resource pool on the current DN.
		On a CN, it indicates the total CPU usage of resource pools on all DNs.
cpu_limit	Int	It indicates the upper limit of available CPUs for resource pools. If the CPU share is limited, this parameter indicates the available CPUs for GaussDB(DWS). If the CPU limit is specified, this parameter indicates the available CPUs for associated Cgroups.
		On a DN, it indicates the upper limit of available CPUs for the resource pool on the current DN.
		On a CN, it indicates the total upper limit of available CPUs for resource pools on all DNs.

Column	Туре	Description
used_mem	Int	 Memory size used by the resource pool (unit: MB) On a DN, it indicates the memory usage of the resource pool on the current DN. On a CN, it indicates the total memory usage of resource pools on all DNs.
estimate_me m	Int	Estimated memory used by the jobs running in the resource pools on the current CN. This parameter is valid only on CNs.
mem_limit	Int	 Upper limit of available memory for the resource pool (unit: MB). On a DN, it indicates the upper limit of available memory for the resource pool on the current DN. On a CN, it indicates the total upper limit of available memory for resource pools on all DNs.
read_kbytes	Bigint	 Number of logical read bytes in the resource pool within a 5s monitoring period (unit: KB). On a DN, it indicates the number of logical read bytes in the resource pool on the current DN. On a CN, it indicates the total logical read bytes of resource pools on all DNs.
write_kbytes	Bigint	 Number of logical write bytes in the resource pool within a 5s monitoring period (unit: KB). On a DN, it indicates the number of logical write bytes in the resource pool on the current DN. On a CN, it indicates the total logical write bytes of resource pools on all DNs.
read_counts	Bigint	 Number of logical reads in the resource pool within a 5s monitoring period. On a DN, it indicates the number of logical reads in the resource pool on the current DN. On a CN, it indicates the total number of logical reads in resource pools on all DNs.

Column	Туре	Description
write_counts	Bigint	Number of logical writes in the resource pool within a 5s monitoring period.
		On a DN, it indicates the number of logical writes in the resource pool on the current DN.
		On a CN, it indicates the total number of logical writes in resource pools on all DNs.
read_speed	Double	Average rate of logical reads of the resource pool in a 5s monitoring period.
		On a DN, it indicates the logical read rate of the resource pool on the current DN.
		On a CN, it indicates the overall logical read rate of resource pools on all DNs.
write_speed	Double	Average rate of logical writes of resource pools in a 5s monitoring period, in KB/s.
		On a DN, it indicates the logical write rate of the resource pool on the current DN.
		On a CN, it indicates the overall logical write rate of resource pools on all DNs.
send_speed	Double	Average network sending rate of a resource pool in a 5-second monitoring period. The unit is KB/s.
		On a DN, it indicates the network sending rate of the resource pool on the current DN.
		On a CN, it indicates that the cumulative sum of the network sending rates of the resource pool on all DNs.
recv_speed	Double	Average network receiving rate of a resource pool in a 5-second monitoring period. The unit is KB/s.
		On a DN, it indicates the network receiving rate of the resource pool on the current DN.
		On a CN, it indicates that the cumulative sum of the network receiving rates of the resource pool on all DNs.

14.3.70 GS_RESPOOL_MONITOR

Displays the job running information and resource usage information of all resource pools. This view can be queried only on CNs. This view is supported only by clusters of 8.2.1.100 and later versions.

Table 14-138 GS_RESPOOL_MONITOR columns

Column	Туре	Description	
rpname	Name	Resource pool name.	
nodegroup	Name	Name of the logical cluster the resource pool belongs to. The default value is installation .	
cn_count	Bigint	Number of CNs in the cluster. This parameter is used to determine whether the management and control result of a single CN is proper in a multi-CN environment.	
short_acc	Boolea n	Whether to enable short query acceleration for a resource pool.	
session_count	Bigint	Number of sessions associated with the resource pool, that is, the number of sessions initiated by users associated with the resource pool, including idle and active sessions.	
active_count	Bigint	Number of active sessions associated with the resource pool, that is, the number of sessions that are performing queries.	
global_wait	Bigint	Number of jobs associated with the resource pool that are queued because the number of concurrent jobs on a single CN exceeds the value of max_active_statements.	
fast_run	Bigint	Number of jobs associated with the resource pool that are running on the fast lane.	
fast_wait	Bigint	Number of jobs associated with the resource pool that are queued on the fast lane.	
fast_limit	Bigint	Maximum number of concurrent jobs on the fast lane in a resource pool.	
slow_run	Bigint	Number of jobs associated with the resource pool that are running on the slow lane.	
slow_wait	Bigint	Number of jobs associated with the resource pool that are queued on the slow lane.	
slow_limit	Bigint	Maximum number of concurrent jobs on the slow lane in a resource pool.	
used_mem	Text	Average memory usage of the resource pool on all DNs. The result has been formatted using pg_size_pretty.	
estimate_me m	Text	Total estimated memory of jobs running in the resource pool. The result has been formatted using pg_size_pretty.	

Column	Туре	Description	
mem_limit	Text	Upper limit of the available memory in the resource pool. The result has been formatted using pg_size_pretty.	
query_mem_li mit	Name	Maximum memory that can be used by a single query in a resource pool. This parameter is used to limit the estimated query memory to prevent abnormal queuing caused by overestimation. The estimated memory is used to limit the actually used query memory. The displayed result has been formatted using pg_size_pretty.	
used_cpu	Doubl e precisi on	Average number of CPU cores occupied by a resource pool on all DNs. CPU isolation is performed by node and resource pool. If a single node contains multiple DNs, the number of CPU cores occupied by a resource pool on a single node must be multiplied by the number of DNs.	
cpu_limit	Doubl e precisi on	Average upper limit of available CPUs for a resource pool on all nodes. If CPU Time Limit is enabled, the value is the total number of available CPU cores of GaussDB(DWS). If CPU Usage Limit is enabled, the value is the number of available CPU cores of the associated Cgroup.	
read_speed	Text	Average logical I/O read rate of the resource pool on all DNs. The result has been formatted using pg_size_pretty.	
write_speed	Text	Average logical I/O write rate of the resource pool on all DNs. The result has been formatted using pg_size_pretty.	
send_speed	Text	Average network sending rate of the resource pool on all DNs. The result has been formatted using pg_size_pretty.	
recv_speed	Text	Average receiving rate of the resource pool on all DNs. The result has been formatted using pg_size_pretty .	

14.3.71 GS_ROW_TABLE_IO_STAT

GS_ROW_TABLE_IO_STAT displays the I/O of all row-store tables of the database on the current node. The value of each statistical column is the accumulated value since the instance was started.

Table 14-139 GS_ROW_TABLE_IO_STAT columns

Column	Туре	Description
schemaname	Name	Namespace of a table
relname	Name	Name of a table
heap_read	Bigint	Number of blocks logically read in the heap
heap_hit	Bigint	Number of block hits in the heap
idx_read	Bigint	Number of blocks logically read in the index
idx_hit	Bigint	Number of block hits in the index
toast_read	Bigint	Number of blocks logically read in the TOAST table
toast_hit	Bigint	Number of block hits in the TOAST table
tidx_read	Bigint	Number of indexes logically read in the TOAST table
tidx_hit	Bigint	Number of index hits in the TOAST table

14.3.72 GS_SESSION_CPU_STATISTICS

GS_SESSION_CPU_STATISTICS displays load management information about CPU usage of ongoing complex jobs executed by the current user.

Table 14-140 GS_SESSION_CPU_STATISTICS columns

Column	Туре	Description
datid	OID	OID of the database the backend is connected to.
usename	Name	Username logged in to the backend.
pid	Bigint	Backend thread ID.
start_time	Timestamp with time zone	Start time of statement execution.
min_cpu_time	Bigint	Minimum CPU time of a statement across all DNs. The unit is ms.
max_cpu_time	Bigint	Maximum CPU time of a statement across all DNs. The unit is ms.
total_cpu_tim e	Bigint	Total CPU time of a statement across all DNs. The unit is ms.

Column	Туре	Description
query	Text	Statement currently being executed.
node_group	Text	Logical cluster of the user running the statement.

14.3.73 GS_SESSION_MEMORY_STATISTICS

GS_SESSION_MEMORY_STATISTICS displays load management information about memory usage of ongoing complex jobs executed by the current user.

Table 14-141 GS_SESSION_MEMORY_STATISTICS columns

Column	Туре	Description
datid	OID	OID of the database the backend is connected to.
usename	Name	Username logged in to the backend.
pid	Bigint	Backend thread ID.
start_time	Timestamp with time zone	Start time of statement execution.
min_peak_me mory	Integer	Minimum memory peak of a statement across all DNs, in MB
max_peak_me mory	Integer	Maximum memory peak of a statement across all DNs, in MB
spill_info	Text	Spill information for the statement on all DNs. The options are:
		None : The statement has not been spilled to disks on any DNs.
		All : The statement has been spilled to disks on all DNs.
		[a:b]: The statement has been spilled to disks on a of b DNs.
query	Text	Statement currently being executed.
node_group	Text	Logical cluster of the user running the statement.

14.3.74 GS_SQL_COUNT

GS_SQL_COUNT displays statistics about the five types of statements (SELECT, INSERT, UPDATE, DELETE, and MERGE INTO) executed on the current node of

the database, including the number of execution times, response time (the maximum, minimum, average, and total response time of the other four types of statements except the **MERGE INTO** statement, in microseconds), and the number of execution times of **DDL**, **DML**, and **DCL statements**.

The classification of **DDL**, **DML**, and **DCL** statements in the **GS_SQL_COUNT** view is slightly different from that of the SQL syntax. The details are as follows:

- User-related statements, such as CREATE/ALTER/DROP USER and CREATE/ ALTER/DROP ROLE, are of the DCL type.
- Transaction-related statements such as BEGIN/COMMIT/SET CONSTRAINTS/ ROLLBACK/SAVEPOINT/START are of the DCL type.
- ALTER SYSTEM KILL SESSION is equivalent to the SELECT pg_terminate_backend() statement and is of the DML type.

The classification of other statements is similar to the definition in the SQL syntax.

When a common user queries the **GS_SQL_COUNT** view, only the statistics of this user in the current node can be viewed. When a user with the administrator permissions queries the **GS_SQL_COUNT** view, the statistics of all users in the current node can be viewed. When the cluster or the node is restarted, the statistics are cleared and the counting restarts. The counting is based on the number of queries received by the node, including the queries performed inside the cluster. Statistics about the **GS_SQL_COUNT** view are collected only on CNs, and SQL statements sent from other CNs are not collected. No result is returned when you query the view on a DN.

Table 14-142 GS_SQL_COUNT columns

Column	Туре	Description
node_name	Name	Node name
user_name	Name	Username
select_count	Bigint	Number of SELECT statements.
update_count	Bigint	Number of UPDATE statements
insert_count	Bigint	Number of INSERT statements
delete_count	Bigint	Number of DELETE statements
mergeinto_count	Bigint	Number of MERGE INTO statements
ddl_count	Bigint	Number of DDL statements
dml_count	Bigint	Number of DML statements
dcl_count	Bigint	Number of DCL statements
total_select_elaps e	Bigint	Total response time of SELECT statements
avg_select_elapse	Bigint	Average response time of SELECT statements

Column	Туре	Description
max_select_elaps e	Bigint	Maximum response time of SELECT statements
min_select_elaps e	Bigint	Minimum response time of SELECT statements
total_update_ela pse	Bigint	Total response time of UPDATE statements
avg_update_elap se	Bigint	Average response time of UPDATE statements
max_update_elap se	Bigint	Maximum response time of UPDATE statements
min_update_elap se	Bigint	Minimum response time of UPDATE statements
total_delete_elap se	Bigint	Total response time of DELETE statements
avg_delete_elaps e	Bigint	Average response time of DELETE statements
max_delete_elaps e	Bigint	Maximum response time of DELETE statements
min_delete_elaps e	Bigint	Minimum response time of DELETE statements
total_insert_elaps e	Bigint	Total response time of INSERT statements
avg_insert_elapse	Bigint	Average response time of INSERT statements
max_insert_elaps e	Bigint	Maximum response time of INSERT statements
min_insert_elaps e	Bigint	Minimum response time of INSERT statements

14.3.75 GS_STAT_DB_CU

GS_STAT_DB_CU displays CU hits of each database in each node of a cluster. You can clear it using **gs_stat_reset()**.

Table 14-143 GS_STAT_DB_CU columns

Column	Туре	Description
node_name1	Text	Node name
db_name	Text	Database name

Column	Туре	Description
mem_hit	Bigint	Number of memory hits
hdd_sync_rea d	Bigint	Number of disk synchronous reads
hdd_asyn_rea d	Bigint	Number of disk asynchronous reads

14.3.76 GS_STAT_SESSION_CU

GS_STAT_SESSION_CU displays the CU hit rate of running sessions on each node in a cluster. This data about a session is cleared when you exit this session or restart the cluster.

Table 14-144 GS_STAT_SESSION_CU columns

Column	Туре	Description
node_name1	Text	Node name
mem_hit	Integer	Number of memory hits
hdd_sync_rea d	Integer	Number of disk synchronous reads
hdd_asyn_rea d	Integer	Number of disk asynchronous reads

14.3.77 GS_TABLE_CHANGE_STAT

GS_TABLE_CHANGE_STAT displays the changes of all tables (excluding foreign tables) of the database on the current node. The value of each column that indicates the number of times is the accumulated value since the instance was started.

Table 14-145 GS_TABLE_CHANGE_STAT columns

Column	Туре	Description
schemaname	Name	Namespace of a table
relname	Name	Table name
last_vacuum	Timestamp with time zone	Time when the last VACUUM operation is performed manually
vacuum_count	Bigint	Number of times of manually performing the VACUUM operation

Column	Туре	Description
last_autovacuum	Timestamp with time zone	Time when the last VACUUM operation is performed automatically
autovacuum_cou nt	Bigint	Number of times of automatically performing the VACUUM operation
last_analyze	Timestamp with time zone	Time when the ANALYZE operation is performed (both manually and automatically)
analyze_count	Bigint	Number of times of performing the ANALYZE operation (both manually and automatically)
last_autoanalyze	Timestamp with time zone	Time when the last ANALYZE operation is performed automatically
autoanalyze_cou nt	Bigint	Number of times of automatically performing the ANALYZE operation
last_change	Bigint	Time when the last modification (INSERT, UPDATE, or DELETE) is performed

14.3.78 GS_TABLE_STAT

GS_TABLE_STAT displays statistics about all tables (excluding foreign tables) of the database on the current node. The values of **live_tuples** and **dead_tuples** are real-time values, and the values of other statistical columns are accumulated values since the instance was started.

Table 14-146 GS_TABLE_STAT columns

Column	Туре	Description
schemaname	Name	Namespace of a table
relname	Name	Table name
seq_scan	Bigint	Number of sequential scans. For a partitioned table, the sum of the number of scans of each partition is displayed.
seq_tuple_read	Bigint	Number of rows scanned in sequence.
index_scan	Bigint	Number of index scans.
index_tuple_read	Bigint	Number of rows scanned by the index.
tuple_inserted	Bigint	Number of rows inserted.

Column	Туре	Description
tuple_updated	Bigint	Number of rows updated.
tuple_deleted	Bigint	Number of rows deleted.
tuple_hot_update d	Bigint	Number of rows with HOT updates.
live_tuples	Bigint	Number of live tuples. Query the view on the CN. If ANALYZE is executed, the total number of live tuples in the table is displayed. Otherwise, 0 is displayed. This indicator applies only to row-store tables.
dead_tuples	Bigint	Number of dead tuples. Query the view on the CN. If ANALYZE is executed, the total number of dead tuples in the table is displayed. Otherwise, 0 is displayed. This indicator applies only to rowstore tables.

14.3.79 GS_TOTAL_NODEGROUP_MEMORY_DETAIL

GS_TOTAL_NODEGROUP_MEMORY_DETAIL displays statistics about memory usage of the logical cluster that the current database belongs to in the unit of MB.

Table 14-147 GS_TOTAL_NODEGROUP_MEMORY_DETAIL columns

Column	Туре	Description
ngname	Text	Logical cluster name

Column	Туре	Description	
memorytype	Text	Memory type. The value can be: • ng_total_memory: total memory of the logical	
		cluster	
		 ng_used_memory: memory usage of the logical cluster 	
		• ng_estimate_memory: estimated memory usage of the logical cluster	
		 ng_foreignrp_memsize: total memory of the external resource pool of the logical cluster 	
		 ng_foreignrp_usedsize: memory usage of the external resource pool of the logical instance ng_foreignrp_peaksize: peak memory usage of the external resource pool of the logical cluster ng_foreignrp_mempct: percentage of the external resource pool of the logical cluster to the total memory of the logical cluster 	
		ng_foreignrp_estmsize: estimated memory usage of the external resource pool of the logical cluster	
memorymbytes	Integer	Size of allocated memory-typed memory	

14.3.80 GS_USER_MONITOR

GS_USER_MONITOR displays all users' job running and resource usage information. This view can be queried only on CNs. This view is supported only by clusters of 8.2.1.100 and later versions.

Table 14-148 GS_USER_MONITOR columns

Column	Туре	Description
usename	Name	Username
rpname	Name	Name of the resource pool associated with the user
nodegroup	Name	Name of the logical cluster the resource pool belongs to. The default value is installation .
session_count	Bigint	Number of sessions initiated by the user, including idle and active sessions
active_count	Bigint	Number of active sessions initiated by the user, that is, the number of sessions that are performing queries.
global_wait	Bigint	Number of jobs that are queued because the number of concurrent jobs on a single CN exceeds the value of max_active_statements.

Column	Туре	Description
fast_run	Bigint	Number of jobs that are running on the fast lane of the resource pool among all jobs executed by the user.
fast_wait	Bigint	Number of jobs queued in the fast lane of the resource pool among all jobs executed by the user.
slow_run	Bigint	Number of jobs that are running on the slow lane of the resource pool among all jobs executed by the user.
slow_wait	Bigint	Number of jobs queued in the slow lane of the resource pool among all jobs executed by the user.
used_mem	Bigint	Average memory used by a user on all DNs, in MB.
estimate_me m	Bigint	Total estimated memory used by running jobs, in MB.
used_cpu	Doubl e precisi on	Average number of CPU cores used by a user on all DNs. If a single node contains multiple DNs, the number of CPU cores used by a user on the node must be multiplied by the number of DNs.
read_speed	Bigint	Average logical I/O read rate of a user on all DNs, in KB/s.
write_speed	Bigint	Average logical I/O write rate of a user on all DNs, in KB/s.
send_speed	Bigint	Average data sending rate of a user on all DNs, in KB/s.
recv_speed	Bigint	Average data receiving rate of a user on all DNs, in KB/s.
used_space	Bigint	Used space of user permanent tables, in KB.
space_limit	Bigint	Maximum space that can be used by user permanent tables, in KB. The value -1 indicates that the space size is not limited.
used_temp_sp ace	Bigint	Used space of user temporary tables, in KB.
temp_space_li mit	Bigint	Maximum space that can be used by user temporary tables, in KB. The value -1 indicates that the space size is not limited.
used_spill_spa ce	Bigint	Used space for flushing intermediate result sets, in KB.
spill_space_li mit	Bigint	Maximum space that can be used for flushing intermediate result sets, in KB. The value -1 indicates that the space size is not limited.

14.3.81 GS_USER_TRANSACTION

GS_USER_TRANSACTION provides transaction information about users on a single CN. The database records the number of times that each user commits and rolls back transactions and the response time of transaction commitment and rollback, in microseconds.

Table 14-149 GS_USER_TRANSACTION columns

Column	Туре	Description
usename	Name	Username
commit_counter	Bigint	Number of the commits
rollback_counter	Bigint	Number of rollbacks
resp_min	Bigint	Minimum response time
resp_max	Bigint	Maximum response time
resp_avg	Bigint	Average response time
resp_total	Bigint	Total response time

14.3.82 GS_VIEW_DEPENDENCY

GS_VIEW_DEPENDENCY allows you to query the direct dependencies of all views visible to the current user.

Table 14-150 GS_VIEW_DEPENDENCY columns

Column	Туре	Description
objschema	Name	View space name
objname	Name	View name
refobjschema	Name	Name of the space where the dependent object resides
refobjname	Name	Name of a dependent object
relobjkind	Char	Type of a dependent object • r: table • v: view

14.3.83 GS VIEW DEPENDENCY PATH

GS_VIEW_DEPENDENCY_PATH allows you to query the direct dependencies of all views visible to the current user. If the base table on which the view depends exists and the dependency between views at different levels is normal, you can use this view to query the dependency between views at different levels starting from the base table.

Table 14-151 GS_VIEW_DEPENDENCY_PATH columns

Column	Туре	Description
objschema	Name	View space name
objname	Name	View name
refobjschema	Name	Name of the space where the dependent object resides
refobjname	Name	Name of a dependent object
path	Text	Dependency path

14.3.84 GS_VIEW_INVALID

GS_VIEW_INVALID gueries all unavailable views visible to the current user.

If the basic table, function, or synonym on which the view depends is abnormal, the **validtype** column of the view is displayed as **invalid**. If the system object on which the view depends changes during an upgrade, the **validtype** column of the view is displayed as **invalidInUpgrade**.

Table 14-152 GS_VIEW_INVALID columns

Column	Туре	Description
OID	OID	OID of the view
schemaname	Name	View space name
viewname	Name	Name of the view
viewowner	Name	Owner of the view
definition	Text	Definition of the view
validtype	Text	View validity flag

Column	Туре	Description
last_invalid_time	Timestamp with time zone	Time when a view is invalid. This column is available only in clusters of version 9.1.0.200 or later.

14.3.85 GS_WAIT_EVENTS

GS_WAIT_EVENTS displays statistics about waiting status and events on the current node.

The values of statistical columns in this view are accumulated only when the **enable_track_wait_event** GUC parameter is set to **on**. If **enable_track_wait_event** is set to **off** during statistics measurement, the statistics will no longer be accumulated, but the existing values are not affected. If **enable_track_wait_event** is **off**, 0 row is returned when this view is queried.

Table 14-153 GS_WAIT_EVENTS columns

Column	Туре	Description
nodename	Name	Node name
type	Text	Event type, which can be STATUS, LOCK_EVENT, LWLOCK_EVENT, or IO_EVENT
event	Text	Event name. For details, see PG_THREAD_WAIT_STATUS.
wait	Bigint	Number of times an event occurs. This column and all the columns below are values accumulated during process running.
failed_wait	Bigint	Number of waiting failures. In the current version, this column is used only for counting timeout errors and waiting failures of locks such as LOCK and LWLOCK.
total_wait_time	Bigint	Total duration of the event
avg_wait_time	Bigint	Average duration of the event
max_wait_time	Bigint	Maximum wait time of the event
min_wait_time	Bigint	Minimum wait time of the event

In the current version, for events whose **type** is **LOCK_EVENT**, **LWLOCK_EVENT**, or **IO_EVENT**, the display scope of **GS_WAIT_EVENTS** is the same as that of the corresponding events in the **PG_THREAD_WAIT_STATUS** view.

For events whose **type** is **STATUS**, **GS_WAIT_EVENTS** displays the following waiting status columns. For details, see the **PG_THREAD_WAIT_STATUS** view.

- acquire lwlock
- acquire lock
- wait io
- wait pooler get conn
- wait pooler abort conn
- wait pooler clean conn
- wait transaction sync
- wait wal sync
- wait data sync
- wait producer ready
- create index
- analyze
- vacuum
- vacuum full
- gtm connect
- gtm begin trans
- gtm commit trans
- gtm rollback trans
- gtm create sequence
- gtm alter sequence
- qtm qet sequence val
- gtm set sequence val
- gtm drop sequence
- gtm rename sequence

14.3.86 GS_WLM_OPERATOR_INFO

This view displays the execution information about operators in the query statements that have been executed on the current CN. The information comes from the system catalog **dbms_om**. **gs_wlm_operator_info**.

NOTICE

- The schema of the GS_WLM_OPERATOR_INFO view is pg_catalog.
- The **GS_WLM_OPERATOR_INFO** view can be queried only in the **postgres** database. If the view is queried in other databases, an error is reported.

Table 14-154 GS_WLM_OPERATOR_INFO columns

Column	Туре	Description
nodename	Text	Name of the CN where the statement is executed
queryid	Bigint	Internal query_id used for statement execution
pid	Bigint	Backend thread ID
plan_node_id	Integer	plan_node_id of the execution plan of a query
plan_node_nam e	Text	Name of the operator corresponding to plan_node_id
start_time	Timestamp with time zone	Time when an operator starts to process the first data record
duration	Bigint	Total execution time of an operator. The unit is ms.
query_dop	Integer	Degree of parallelism (DOP) of the current operator
estimated_rows	Bigint	Number of rows estimated by the optimizer
tuple_processed	Bigint	Number of elements returned by the current operator
min_peak_mem ory	Integer	Minimum peak memory used by the current operator on all DNs. The unit is MB.
max_peak_me mory	Integer	Maximum peak memory used by the current operator on all DNs. The unit is MB.
average_peak_ memory	Integer	Average peak memory used by the current operator on all DNs. The unit is MB.
memory_skew_ percent	Integer	Memory usage skew of the current operator among DNs
min_spill_size	Integer	Minimum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
max_spill_size	Integer	Maximum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
average_spill_si ze	Integer	Average spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
spill_skew_perc ent	Integer	DN spill skew when a spill occurs

Column	Туре	Description
min_cpu_time	Bigint	Minimum execution time of the operator on all DNs. The unit is ms.
max_cpu_time	Bigint	Maximum execution time of the operator on all DNs. The unit is ms.
total_cpu_time	Bigint	Total execution time of the operator on all DNs. The unit is ms.
cpu_skew_perce nt	Integer	Skew of the execution time among DNs.
warning	Text	Warning. The following warnings are displayed:
		Sort/SetOp/HashAgg/HashJoin spill
		2. Spill file size large than 256MB
		3. Broadcast size large than 100MB
		4. Early spill
		5. Spill times is greater than 3
		6. Spill on memory adaptive
		7. Hash table conflict

14.3.87 GS_WLM_OPERATOR_HISTORY

GS_WLM_OPERATOR_HISTORY displays the records of operators in jobs that have been executed by the current user on the current CN.

This view is used to query data from GaussDB(DWS). Data in the database is cleared periodically. If the GUC parameter <code>enable_resource_record</code> is set to <code>on</code>, records in the view will be dumped to the system catalog <code>GS_WLM_OPERATOR_INFO</code> every 3 minutes and deleted from the view. If <code>enable_resource_record</code> is set to <code>off</code>, the records will be deleted from the view after the retention period expires. The recorded data is the same as that described in <code>Table 14-155</code>.

Table 14-155 GS_WLM_OPERATOR_INFO columns

Column	Туре	Description
nodename	Text	Name of the CN where the statement is executed
queryid	Bigint	Internal query_id used for statement execution
pid	Bigint	Backend thread ID
plan_node_id	Integer	plan_node_id of the execution plan of a query

Column	Туре	Description
plan_node_nam e	Text	Name of the operator corresponding to plan_node_id
start_time	Timestamp with time zone	Time when an operator starts to process the first data record
duration	Bigint	Total execution time of an operator. The unit is ms.
query_dop	Integer	Degree of parallelism (DOP) of the current operator
estimated_rows	Bigint	Number of rows estimated by the optimizer
tuple_processed	Bigint	Number of elements returned by the current operator
min_peak_mem ory	Integer	Minimum peak memory used by the current operator on all DNs. The unit is MB.
max_peak_me mory	Integer	Maximum peak memory used by the current operator on all DNs. The unit is MB.
average_peak_ memory	Integer	Average peak memory used by the current operator on all DNs. The unit is MB.
memory_skew_ percent	Integer	Memory usage skew of the current operator among DNs
min_spill_size	Integer	Minimum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
max_spill_size	Integer	Maximum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
average_spill_si ze	Integer	Average spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
spill_skew_perc ent	Integer	DN spill skew when a spill occurs
min_cpu_time	Bigint	Minimum execution time of the operator on all DNs. The unit is ms.
max_cpu_time	Bigint	Maximum execution time of the operator on all DNs. The unit is ms.
total_cpu_time	Bigint	Total execution time of the operator on all DNs. The unit is ms.
cpu_skew_perce nt	Integer	Skew of the execution time among DNs.

Column	Туре	Description
warning	Text	Warning. The following warnings are displayed:
		Sort/SetOp/HashAgg/HashJoin spill
		2. Spill file size large than 256MB
		3. Broadcast size large than 100MB
		4. Early spill
		5. Spill times is greater than 3
		6. Spill on memory adaptive
		7. Hash table conflict

14.3.88 GS_WLM_OPERATOR_STATISTICS

GS_WLM_OPERATOR_STATISTICS displays the operators of the jobs that are being executed by the current user.

Table 14-156 GS_WLM_OPERATOR_STATISTICS columns

Column	Туре	Description
queryid	Bigint	Internal query_id used for statement execution
pid	Bigint	ID of the backend thread
plan_node_id	Integer	plan_node_id of the execution plan of a query
plan_node_na me	Text	Name of the operator corresponding to plan_node_id . The maximum length of the operator name is 127 characters (excluding format characters such as spaces).
start_time	Timestamp with time zone	Time when the operator starts to be executed for the first time.
duration	Bigint	Total execution time of the operator from the start to the end, in milliseconds.
status	Text	Execution status of the current operator. The value can be waiting , running , or finished .
query_dop	Integer	DOP of the current operator
estimated_rows	Bigint	Number of rows estimated by the optimizer. If the number of returned estimated rows exceeds int64_max, int64_max is displayed.

Column	Туре	Description
tuple_processe d	Bigint	Total number of elements returned by the current operator on all DNs. If the estimated number of returned rows exceeds int64_max, int64_max is displayed.
min_peak_mem ory	Integer	Minimum peak memory used by the current operator on all DNs. The unit is MB.
max_peak_me mory	Integer	Maximum peak memory used by the current operator on all DNs. The unit is MB.
average_peak_ memory	Integer	Average peak memory used by the current operator on all DNs. The unit is MB.
memory_skew_ percent	Integer	Memory usage skew of the current operator among DNs
min_spill_size	Integer	Minimum logical spilled data among all DNs when a spill occurs, in MB. The default value is 0 .
max_spill_size	Integer	Maximum logical spilled data among all DNs when a spill occurs, in MB. The default value is 0 .
average_spill_si ze	Integer	Average logical spilled data among all DNs when a spill occurs, in MB. The default value is 0 .
spill_skew_perc ent	Integer	DN spill skew when a spill occurs
min_cpu_time	Bigint	Minimum execution time of the operator on all DNs. The unit is ms.
max_cpu_time	Bigint	Maximum execution time of the operator on all DNs. The unit is ms.
total_cpu_time	Bigint	Total execution time of the operator on all DNs. The unit is ms.
cpu_skew_perc ent	Integer	Skew of the execution time among DNs.

Column	Туре	Description
warning	Text	Warning. The following warnings are displayed: 1. Sort/SetOp/HashAgg/HashJoin spill 2. Spill file size large than 256MB 3. Broadcast size large than 100MB 4. Early spill 5. Spill times is greater than 3 6. Spill on memory adaptive 7. Hash table conflict
parent_id	Integer	Parent node ID of the operator node.
exec_count	Integer	Maximum number of times that the operator node can be executed on all DNs.
progress	Text	Progress information of the operator. For the first operator, it is the overall progress of the job. For other operators, it is the progress of the current operator.
min_net_size	Bigint	Minimum network communication data volume (KB) of the operator on all DNs. It mainly applies to network operators.
max_net_size	Bigint	Maximum network communication data volume (KB) of the operator on all DNs. It mainly applies to network operators.
total_net_size	Bigint	Total network communication data volume (KB) of the operator on all DNs. It mainly applies to network operators.
min_read_bytes	Bigint	Minimum amount of data read by the operator from disks on all DNs. The unit is KB.
max_read_byte s	Bigint	Maximum amount of data read by the operator from disks on all DNs. The unit is KB.
total_read_byte s	Bigint	Total amount of data read by the operator from disks on all DNs, in KB.
min_write_byte s	Bigint	Minimum amount of data written by the operator to disks on all DNs. The unit is KB.
max_write_byte s	Bigint	Maximum amount of data written by the operator to disks on all DNs. The unit is KB.
total_write_byt es	Bigint	Total amount of data written by the operator to disks on all DNs, in KB.

14.3.89 GS_WLM_SESSION_INFO

This view displays the execution information about the query statements that have been executed on the current CN. The information comes from the system catalog **dbms_om.gs_wlm_session_info**.

NOTICE

- The schema of the GS_WLM_SESSION_INFO view is pg_catalog.
- The **GS_WLM_SESSION_INFO** view can be queried only in the **postgres** database. If the view is queried in other databases, an error is reported.

Table 14-158 lists the columns in the GS_WLM_SESSION_INFO view.

Table 14-157 GS_WLM_SESSION_HISTORY columns

Column	Туре	Description
datid	OID	OID of the database this backend is connected to
dbname	Text	Name of the database the backend is connected to
schemaname	Text	Schema name
nodename	Text	Name of the CN where the statement is run
username	Text	User name used for connecting to the backend
application_na me	Text	Name of the application that is connected to the backend
client_addr	inet	IP address of the client connected to this backend. If this column is null, it indicates either that the client is connected via a Unix socket on the server machine or that this is an internal process such as autovacuum.
client_hostnam e	Text	Host name of the connected client, as reported by a reverse DNS lookup of client_addr . This column will only be non-null for IP connections, and only when log_hostname is enabled.
client_port	Integer	TCP port number that the client uses for communication with this backend, or -1 if a Unix socket is used
query_band	Text	Job type, which can be set using the GUC parameter query_band and is a null string by default.

Column	Туре	Description
block_time	Bigint	Duration that a statement is blocked before being executed, including the statement parsing and optimization duration. The unit is ms.
start_time	Timestamp with time zone	Time when the statement starts to be run
finish_time	Timestamp with time zone	Time when the statement execution ends
duration	Bigint	Execution time of a statement. The unit is ms.
estimate_total_ time	Bigint	Estimated execution time of a statement. The unit is ms.
status	Text	Final statement execution status. Its value can be finished (normal) or aborted (abnormal). The statement status here is the execution status of the database server. If the statement is successfully executed on the database server but an error is reported in the result set, the statement status is finished .
abort_info	Text	Exception information displayed if the final statement execution status is aborted .
resource_pool	Text	Resource pool used by the user
control_group	Text	Cgroup used by the statement
estimate_mem ory	Integer	Estimated memory used by a statement on a single instance. The unit is MB.
min_peak_mem ory	Integer	Minimum memory peak of a statement across all DNs. The unit is MB.
max_peak_me mory	Integer	Maximum memory peak of a statement across all DNs. The unit is MB.
average_peak_ memory	Integer	Average memory usage during statement execution. The unit is MB.
memory_skew_ percent	Integer	Memory usage skew of a statement among DNs.

Column	Туре	Description
spill_info	Text	Statement spill information on all DNs.
		None indicates that the statement has not been spilled to disks on any DNs.
		All : The statement has been spilled to disks on all DNs.
		[a:b]: The statement has been spilled to disks on a of b DNs.
min_spill_size	Integer	Minimum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
max_spill_size	Integer	Maximum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
average_spill_si ze	Integer	Average spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
spill_skew_perc ent	Integer	DN spill skew when a spill occurs
min_dn_time	Bigint	Minimum execution time of a statement across all DNs. The unit is ms.
max_dn_time	Bigint	Maximum execution time of a statement across all DNs. The unit is ms.
average_dn_tim e	Bigint	Average execution time of a statement across all DNs. The unit is ms.
dntime_skew_p ercent	Integer	Execution time skew of a statement among DNs.
min_cpu_time	Bigint	Minimum CPU time of a statement across all DNs. The unit is ms.
max_cpu_time	Bigint	Maximum CPU time of a statement across all DNs. The unit is ms.
total_cpu_time	Bigint	Total CPU time of a statement across all DNs. The unit is ms.
cpu_skew_perce nt	Integer	CPU time skew of a statement among DNs.
min_peak_iops	Integer	Minimum IOPS peak of a statement across all DNs. It is counted by ones in a column-store table and by ten thousands in a row-store table.

Column	Туре	Description
max_peak_iops	Integer	Maximum IOPS peak of a statement across all DNs. It is counted by ones in a column-store table and by ten thousands in a row-store table.
average_peak_i ops	Integer	Average IOPS peak of a statement across all DNs. It is counted by ones in a column-store table and by ten thousands in a row-store table.
iops_skew_perc ent	Integer	I/O skew across DNs.
warning	Text	Warning. The following warnings and warnings related to SQL self-diagnosis tuning are displayed: 1. Spill file size large than 256MB 2. Broadcast size large than 100MB 3. Early spill
		4. Spill times is greater than 3
		5. Spill on memory adaptive 6. Hash table conflict
queryid	Bigint	Internal query ID used for statement execution
query	Text	Statement to be executed. A maximum of 64 KB of strings can be retained.
query_plan	Text	Execution plan of a statement.
		Specification restrictions:
		Execution plans are displayed only for DML statements.
		2. In 8.2.1.100 and later versions, the number of data binding times is added to the execution plans of Parse Bind Execute (PBE) statements to facilitate statement analysis. The number of data binding times is displayed in the format of PBE bind times : <i>Times</i> .
node_group	Text	Logical cluster of the user running the statement
pid	Bigint	PID of the backend thread of the statement
lane	Text	Fast/Slow lane where the statement is executed
unique_sql_id	Bigint	ID of the normalized unique SQL.

Column	Туре	Description
session_id	Text	Unique identifier of a session in the database system. Its format is session_start_time.tid.node_name.
min_read_bytes	Bigint	Minimum I/O read bytes of a statement across all DNs. The unit is byte.
max_read_byte	Bigint	Maximum I/O read bytes of a statement across all DNs. The unit is byte.
average_read_b ytes	Bigint	Average I/O read bytes of a statement across all DNs.
min_write_byte s	Bigint	Minimum I/O write bytes of a statement across all DNs.
max_write_byte	Bigint	Maximum I/O write bytes of a statement across all DNs.
average_write_ bytes	Bigint	Average I/O write bytes of a statement across all DNs.
recv_pkg	Bigint	Total number of communication packages received by a statement across all DNs.
send_pkg	Bigint	Total number of communication packages sent by a statement across all DNs.
recv_bytes	Bigint	Total received data of the statement stream, in byte.
send_bytes	Bigint	Total sent data of the statement stream, in byte.
stmt_type	Text	Query type corresponding to the statement.
except_info	Text	Information about the exception rule triggered by the statement.
unique_plan_id	Bigint	ID of the normalized unique plan.
sql_hash	Text	Normalized SQL hash.
plan_hash	Text	Normalized plan hash.
use_plan_baseli ne	Text	Indicates whether the bound plan is used for executing the current statement. If it is used, the name of the plan_baseline column in pg_plan_baseline is displayed.
outline_name	Text	Name of the outline used for the statement plan.

Column	Туре	Description
loader_status	Text	The JSON string for storing import and export service information is as follows.
		 address: indicates the IP address of the peer cluster. The port number is displayed for the source cluster.
		2. direction : indicates the import and export service type. The value can be gds to file , gds from file , gds to pipe , gds from pipe , copy from or copy to .
		3. min/max/total_lines/bytes: indicates the minimum value, maximum value, total lines, and bytes of the import and export statements on all DNs.
parse_time	Bigint	Total parsing time before the statement is queued (including lexical and syntax parsing, optimization rewriting, and plan generation time), in milliseconds. This column is available only in clusters of version 8.3.0.100 or later.
disk_cache_hit_ ratio	numeric(5,2)	Disk cache hit rate. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
disk_cache_disk _read_size	Bigint	Total size of data read from disk cache, in MB. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
disk_cache_disk _write_size	Bigint	Total size of data written to disk cache, in MB. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
disk_cache_rem ote_read_size	Bigint	Total size of data read remotely from OBS due to disk cache read failure, in MB. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
disk_cache_rem ote_read_time	Bigint	Total number of times data is read remotely from OBS due to disk cache read failure. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
vfs_scan_bytes	Bigint	Total number of bytes scanned by the OBS virtual file system in response to upper-layer requests, in bytes. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.

Column	Туре	Description
vfs_remote_rea d_bytes	Bigint	Total number of bytes actually read from OBS by the OBS virtual file system, in bytes. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
preload_submit _time	Bigint	Total time for submitting I/O requests in the prefetching process, in microseconds. This column only applies to OBS 3.0 tables with storage and compute decoupled.
preload_wait_ti me	Bigint	Total time for waiting for I/O requests in the prefetching process, in microseconds. This column only applies to OBS 3.0 tables with storage and compute decoupled.
preload_wait_c ount	Bigint	Total number of times that the prefetching process waits for I/O requests. This column only applies to OBS 3.0 tables with storage and compute decoupled.
disk_cache_loa d_time	Bigint	Total time for reading from disk cache, in microseconds. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
disk_cache_conf lict_count	Bigint	Number of times a block in the disk cache produces a hash conflict. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
disk_cache_erro r_count	Bigint	Number of disk cache read failures. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.

Column	Туре	Description
disk_cache_erro r_code	Bigint	Error code for disk cache read failures. Multiple error codes may be generated. If the disk cache fails to be read, OBS remote read is initiated and cache blocks are rewritten. The error code types are as follows: This column only applies to OBS 3.0 tables and foreign tables.
		1: A hash conflict occurs in the disk cache block.
		2: The generation time of the disk cache block is later than that of the OldestXmin transaction.
		• 4: Invoking the pread system when reading cache files from the disk cache failed.
		8: The data version of the disk cache block does not match.
		16: The version of the data written to the write cache does not match the latest version.
		• 32: Opening the cache file corresponding to the cache block failed.
		64: The size of the data read from the disk cache does not match.
		128: The CSN recorded in the disk cache block does not match.
obs_io_req_avg _rtt	Bigint	Average Round Trip Time (RTT) for OBS I/O requests, in microseconds. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
obs_io_req_avg _latency	Bigint	Average delay for OBS I/O requests, in microseconds. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
obs_io_req_late ncy_gt_1s	Bigint	Number of OBS I/O requests with a latency exceeding 1 second. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
obs_io_req_late ncy_gt_10s	Bigint	Number of OBS I/O requests with a latency exceeding 10 seconds. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
obs_io_req_cou nt	Bigint	Total number of OBS I/O requests. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.

Column	Туре	Description
obs_io_req_retr y_count	Bigint	Total number of retries for OBS I/O requests. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
obs_io_req_rate _limit_count	Bigint	Total number of times OBS I/O requests are flow-controlled. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.

14.3.90 GS_WLM_SESSION_HISTORY

GS_WLM_SESSION_HISTORY displays load management information about a completed job executed by the current user on the current CN. The view is used to query data from GaussDB(DWS). The view returns the data queried from the **GS_WLM_SESSION_INFO** table within 3 minutes only if the GUC parameter **enable_resource_track** is set to **on**.

NOTICE

The **GS_WLM_SESSION_HISTORY** view can be queried only in the **postgres** database. If the view is queried in other databases, an error is reported.

Table 14-158 GS_WLM_SESSION_HISTORY columns

Column	Туре	Description
datid	OID	OID of the database this backend is connected to
dbname	Text	Name of the database the backend is connected to
schemaname	Text	Schema name
nodename	Text	Name of the CN where the statement is run
username	Text	User name used for connecting to the backend
application_na me	Text	Name of the application that is connected to the backend
client_addr	inet	IP address of the client connected to this backend. If this column is null, it indicates either that the client is connected via a Unix socket on the server machine or that this is an internal process such as autovacuum.

Column	Туре	Description
client_hostnam e	Text	Host name of the connected client, as reported by a reverse DNS lookup of client_addr . This column will only be non-null for IP connections, and only when log_hostname is enabled.
client_port	Integer	TCP port number that the client uses for communication with this backend, or -1 if a Unix socket is used
query_band	Text	Job type, which can be set using the GUC parameter query_band and is a null string by default.
block_time	Bigint	Duration that a statement is blocked before being executed, including the statement parsing and optimization duration. The unit is ms.
start_time	Timestamp with time zone	Time when the statement starts to be run
finish_time	Timestamp with time zone	Time when the statement execution ends
duration	Bigint	Execution time of a statement. The unit is ms.
estimate_total_ time	Bigint	Estimated execution time of a statement. The unit is ms.
status	Text	Final statement execution status. Its value can be finished (normal) or aborted (abnormal). The statement status here is the execution status of the database server. If the statement is successfully executed on the database server but an error is reported in the result set, the statement status is finished .
abort_info	Text	Exception information displayed if the final statement execution status is aborted .
resource_pool	Text	Resource pool used by the user
control_group	Text	Cgroup used by the statement
estimate_mem ory	Integer	Estimated memory used by a statement on a single instance. The unit is MB.
min_peak_mem ory	Integer	Minimum memory peak of a statement across all DNs. The unit is MB.
max_peak_me mory	Integer	Maximum memory peak of a statement across all DNs. The unit is MB.

Column	Туре	Description
average_peak_ memory	Integer	Average memory usage during statement execution. The unit is MB.
memory_skew_ percent	Integer	Memory usage skew of a statement among DNs.
spill_info	Text	Statement spill information on all DNs.
		None indicates that the statement has not been spilled to disks on any DNs.
		All : The statement has been spilled to disks on all DNs.
		[a:b]: The statement has been spilled to disks on a of b DNs.
min_spill_size	Integer	Minimum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
max_spill_size	Integer	Maximum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
average_spill_si ze	Integer	Average spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
spill_skew_perc ent	Integer	DN spill skew when a spill occurs
min_dn_time	Bigint	Minimum execution time of a statement across all DNs. The unit is ms.
max_dn_time	Bigint	Maximum execution time of a statement across all DNs. The unit is ms.
average_dn_tim e	Bigint	Average execution time of a statement across all DNs. The unit is ms.
dntime_skew_p ercent	Integer	Execution time skew of a statement among DNs.
min_cpu_time	Bigint	Minimum CPU time of a statement across all DNs. The unit is ms.
max_cpu_time	Bigint	Maximum CPU time of a statement across all DNs. The unit is ms.
total_cpu_time	Bigint	Total CPU time of a statement across all DNs. The unit is ms.
cpu_skew_perce nt	Integer	CPU time skew of a statement among DNs.

Column	Туре	Description
min_peak_iops	Integer	Minimum IOPS peak of a statement across all DNs. It is counted by ones in a column-store table and by ten thousands in a row-store table.
max_peak_iops	Integer	Maximum IOPS peak of a statement across all DNs. It is counted by ones in a column-store table and by ten thousands in a row-store table.
average_peak_i ops	Integer	Average IOPS peak of a statement across all DNs. It is counted by ones in a column-store table and by ten thousands in a row-store table.
iops_skew_perc ent	Integer	I/O skew across DNs.
warning	Text	Warning. The following warnings and warnings related to SQL self-diagnosis tuning are displayed: 1. Spill file size large than 256MB 2. Broadcast size large than 100MB 3. Early spill 4. Spill times is greater than 3 5. Spill on memory adaptive 6. Hash table conflict
queryid	Bigint	Internal query ID used for statement execution
query	Text	Statement to be executed. A maximum of 64 KB of strings can be retained.
query_plan	Text	 Execution plan of a statement. Specification restrictions: 1. Execution plans are displayed only for DML statements. 2. In 8.2.1.100 and later versions, the number of data binding times is added to the execution plans of Parse Bind Execute (PBE) statements to facilitate statement analysis. The number of data binding times is displayed in the format of PBE bind times: <i>Times</i>.
node_group	Text	Logical cluster of the user running the statement
pid	Bigint	PID of the backend thread of the statement

Column	Туре	Description
lane	Text	Fast/Slow lane where the statement is executed
unique_sql_id	Bigint	ID of the normalized unique SQL.
session_id	Text	Unique identifier of a session in the database system. Its format is session_start_time.tid.node_name.
min_read_bytes	Bigint	Minimum I/O read bytes of a statement across all DNs. The unit is byte.
max_read_byte s	Bigint	Maximum I/O read bytes of a statement across all DNs. The unit is byte.
average_read_b ytes	Bigint	Average I/O read bytes of a statement across all DNs.
min_write_byte s	Bigint	Minimum I/O write bytes of a statement across all DNs.
max_write_byte s	Bigint	Maximum I/O write bytes of a statement across all DNs.
average_write_ bytes	Bigint	Average I/O write bytes of a statement across all DNs.
recv_pkg	Bigint	Total number of communication packages received by a statement across all DNs.
send_pkg	Bigint	Total number of communication packages sent by a statement across all DNs.
recv_bytes	Bigint	Total received data of the statement stream, in byte.
send_bytes	Bigint	Total sent data of the statement stream, in byte.
stmt_type	Text	Query type corresponding to the statement.
except_info	Text	Information about the exception rule triggered by the statement.
unique_plan_id	Bigint	ID of the normalized unique plan.
sql_hash	Text	Normalized SQL hash.
plan_hash	Text	Normalized plan hash.
use_plan_baseli ne	Text	Indicates whether the bound plan is used for executing the current statement. If it is used, the name of the plan_baseline column in pg_plan_baseline is displayed.

Column	Туре	Description
outline_name	Text	Name of the outline used for the statement plan.
loader_status	Text	The JSON string for storing import and export service information is as follows.
		1. address : indicates the IP address of the peer cluster. The port number is displayed for the source cluster.
		2. direction : indicates the import and export service type. The value can be gds to file , gds from file , gds to pipe , gds from pipe , copy from or copy to .
		3. min/max/total_lines/bytes: indicates the minimum value, maximum value, total lines, and bytes of the import and export statements on all DNs.
parse_time	Bigint	Total parsing time before the statement is queued (including lexical and syntax parsing, optimization rewriting, and plan generation time), in milliseconds. This column is available only in clusters of version 8.3.0.100 or later.
disk_cache_hit_ ratio	numeric(5,2	Disk cache hit rate. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
disk_cache_disk _read_size	Bigint	Total size of data read from disk cache, in MB. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
disk_cache_disk _write_size	Bigint	Total size of data written to disk cache, in MB. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
disk_cache_rem ote_read_size	Bigint	Total size of data read remotely from OBS due to disk cache read failure, in MB. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
disk_cache_rem ote_read_time	Bigint	Total number of times data is read remotely from OBS due to disk cache read failure. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.

Column	Туре	Description	
vfs_scan_bytes	Bigint	Total number of bytes scanned by the OBS virtual file system in response to upper-layer requests, in bytes. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.	
vfs_remote_rea d_bytes	Bigint	Total number of bytes actually read from OBS by the OBS virtual file system, in bytes. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.	
preload_submit _time	Bigint	Total time for submitting I/O requests in the prefetching process, in microseconds. This column only applies to OBS 3.0 tables with storage and compute decoupled.	
preload_wait_ti me	Bigint	Total time for waiting for I/O requests in the prefetching process, in microseconds. This column only applies to OBS 3.0 tables with storage and compute decoupled.	
preload_wait_c ount	Bigint	Total number of times that the prefetching process waits for I/O requests. This column only applies to OBS 3.0 tables with storage and compute decoupled.	
disk_cache_loa d_time	Bigint	Total time for reading from disk cache, in microseconds. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.	
disk_cache_conf lict_count	Bigint	Number of times a block in the disk cache produces a hash conflict. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.	
disk_cache_erro r_count	Bigint	Number of disk cache read failures. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.	

Column	Туре	Description
disk_cache_erro r_code	Bigint	Error code for disk cache read failures. Multiple error codes may be generated. If the disk cache fails to be read, OBS remote read is initiated and cache blocks are rewritten. The error code types are as follows: This column only applies to OBS 3.0 tables and foreign tables.
		1: A hash conflict occurs in the disk cache block.
		2: The generation time of the disk cache block is later than that of the OldestXmin transaction.
		• 4: Invoking the pread system when reading cache files from the disk cache failed.
		8: The data version of the disk cache block does not match.
		16: The version of the data written to the write cache does not match the latest version.
		• 32: Opening the cache file corresponding to the cache block failed.
		64: The size of the data read from the disk cache does not match.
		128: The CSN recorded in the disk cache block does not match.
obs_io_req_avg _rtt	Bigint	Average Round Trip Time (RTT) for OBS I/O requests, in microseconds. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
obs_io_req_avg _latency	Bigint	Average delay for OBS I/O requests, in microseconds. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
obs_io_req_late ncy_gt_1s	Bigint	Number of OBS I/O requests with a latency exceeding 1 second. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
obs_io_req_late ncy_gt_10s	Bigint	Number of OBS I/O requests with a latency exceeding 10 seconds. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
obs_io_req_cou nt	Bigint	Total number of OBS I/O requests. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.

Column	Туре	Description
obs_io_req_retr y_count	Bigint	Total number of retries for OBS I/O requests. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.
obs_io_req_rate _limit_count	Bigint	Total number of times OBS I/O requests are flow-controlled. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.

14.3.91 GS_WLM_SESSION_STATISTICS

GS_WLM_SESSION_STATISTICS displays load management information about jobs being executed by the current user on the current CN.

Table 14-159 GS_WLM_SESSION_STATISTICS columns

Column	Туре	Description	
datid	OID	OID of the database this backend is connected to	
dbname	Name	Name of the database the backend is connected to	
schemaname	Text	Schema name	
nodename	Text	Name of the CN where the statement is executed	
username	Name	User name used for connecting to the backend	
application_nam e	Text	Name of the application that is connected to the backend	
client_addr	inet	IP address of the client connected to this backend. If this column is null, it indicates either that the client is connected via a Unix socket on the server machine or that this is an internal process such as autovacuum.	
client_hostname	Text	Host name of the connected client, as reported by a reverse DNS lookup of client_addr . This column will only be non-null for IP connections, and only when log_hostname is enabled.	
client_port	Integer	TCP port number that the client uses for communication with this backend, or -1 if a Unix socket is used	

Column	Туре	Description	
query_band	Text	Job type, which can be set using the GUC parameter query_band and is a null string by default.	
pid	Bigint	Process ID of the backend	
block_time	Bigint	Block time before the statement is executed. The unit is ms.	
start_time	Timestamp with time zone	Time when the statement starts to be executed	
duration	Bigint	For how long a statement has been executing. The unit is ms.	
estimate_total_ti me	Bigint	Estimated execution time of a statement. The unit is ms.	
estimate_left_ti me	Bigint	Estimated remaining time of statement execution. The unit is ms.	
enqueue	Text	Workload management resource status	
resource_pool	Name	Resource pool used by the user	
control_group	Text	Cgroup used by the statement	
estimate_memor y	Integer	Estimated memory used by a statement on a single instance. The unit is MB.	
min_peak_mem ory	Integer	Minimum memory peak of a statement across all DNs. The unit is MB.	
max_peak_mem ory	Integer	Maximum memory peak of a statement across all DNs. The unit is MB.	
average_peak_m emory	Integer	Average memory usage during statement execution. The unit is MB.	
memory_skew_p ercent	Integer	Memory usage skew of a statement among DNs.	
spill_info	Text	Statement spill information on all DNs. None: The statement has not been spilled to disks on any DNs. All: The statement has been spilled to disks on all DNs. [a:b]: The statement has been spilled to disks on a of b DNs.	
min_spill_size	Integer	Minimum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .	

Column	Туре	Description	
max_spill_size	Integer	Maximum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .	
average_spill_siz e	Integer	Average spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .	
spill_skew_perce nt	Integer	DN spill skew when a spill occurs	
min_dn_time	Bigint	Minimum execution time of a statement across all DNs. The unit is ms.	
max_dn_time	Bigint	Maximum execution time of a statement across all DNs. The unit is ms.	
average_dn_tim e	Bigint	Average execution time of a statement across all DNs. The unit is ms.	
dntime_skew_pe rcent	Integer	Execution time skew of a statement among DNs.	
min_cpu_time	Bigint	Minimum CPU time of a statement across all DNs. The unit is ms.	
max_cpu_time	Bigint	Maximum CPU time of a statement across all DNs. The unit is ms.	
total_cpu_time	Bigint	Total CPU time of a statement across all DNs. The unit is ms.	
cpu_skew_perce nt	Integer	CPU time skew of a statement among DNs.	
min_peak_iops	Integer	Minimum IOPS peak of a statement across all DNs. It is counted by ones in a column-store table and by ten thousands in a row-store table.	
max_peak_iops	Integer	Maximum IOPS peak of a statement across all DNs. It is counted by ones in a column-store table and by ten thousands in a row-store table.	
average_peak_io ps	Integer	Average IOPS peak of a statement across all DNs. It is counted by ones in a column-store table and by ten thousands in a row-store table.	
iops_skew_perce nt	Integer	I/O skew across DNs.	

Column	Туре	Description	
min_read_speed	Integer	Minimum I/O read rate of a statement across all DNs within a monitoring period (5s). The unit is KB/s.	
max_read_speed	Integer	Maximum I/O read rate of a statement across all DNs within a monitoring period (5s). The unit is KB/s.	
average_read_sp eed	Integer	Average I/O read rate of a statement across all DNs within a monitoring period (5s). The unit is KB/s.	
min_write_speed	Integer	Minimum I/O write rate of a statement across all DNs within a monitoring period (5s). The unit is KB/s.	
max_write_spee d	Integer	Maximum I/O write rate of a statement across all DNs within a monitoring period (5s). The unit is KB/s.	
average_write_s peed	Integer	Average I/O write rate of a statement across all DNs within a monitoring period (5s). The unit is KB/s.	
recv_pkg	Bigint	Total number of communication packages received by a statement across all DNs.	
send_pkg	Bigint	Total number of communication packages sent by a statement across all DNs.	
recv_bytes	Bigint	Total received data of the statement stream, in byte.	
send_bytes	Bigint	Total sent data of the statement stream, in byte.	
warning	Text	Warning. The following warnings and warnings related to SQL self-diagnosis tuning are displayed: 1. Spill file size large than 256MB	
		2. Broadcast size large than 100MB	
		3. Early spill	
		4. Spill times is greater than 3	
		5. Spill on memory adaptive6. Hash table conflict	
unique_sql_id	Bigint	ID of the normalized unique SQL.	
queryid	Bigint	Internal query ID used for statement execution	
query	Text	Statement that is being executed	

Column	Туре	Description	
query_plan	Text	Execution plan of a statement	
		Specification restrictions:	
		Execution plans are displayed only for DML statements.	
		2. In 8.2.1.100 and later versions, the number of data binding times is added to the execution plans of Parse Bind Execute (PBE) statements to facilitate statement analysis. The number of data binding times is displayed in the format of PBE bind times : <i>Times</i> .	
node_group	Text	Logical cluster of the user running the statement	
stmt_type	Text	Query type corresponding to the statement.	
except_info	Text	Information about the exception rule triggered by the statement.	
parse_time	Bigint	Total parsing time before the statement is queued (including lexical and syntax parsing, optimization rewriting, and plan generation time), in milliseconds.	
		This column is only supported in version 8.3.0.100 or later.	
unique_plan_id	Bigint	ID of the normalized unique plan.	
sql_hash	Text	Normalized SQL hash.	
plan_hash	Text	Normalized plan hash.	
disk_cache_hit_r atio	numeric(5, 2)	Disk cache hit rate. This column only applies to OBS 3.0 tables and foreign tables.	
disk_cache_disk_ read_size	Bigint	Total size of data read from disk cache, in MB. This column only applies to OBS 3.0 tables and foreign tables.	
disk_cache_disk_ write_size	Bigint	Total size of data written to disk cache, in MB. This column only applies to OBS 3.0 tables and foreign tables.	
disk_cache_remo te_read_size	Bigint	Total size of data read remotely from OBS due to disk cache read failure, in MB. This column only applies to OBS 3.0 tables and foreign tables.	
disk_cache_remo te_read_time	Bigint	Total number of times data is read remotely from OBS due to disk cache read failure. This column only applies to OBS 3.0 tables and foreign tables.	

Column	Туре	Description
block_name	Text	Name of the interception rule that matches the statement. This column is available only in clusters of version 9.1.0.200 or later.

14.3.92 GS_WLM_SQL_ALLOW

The **GS_WLM_SQL_ALLOW** view displays the configured resource management SQL whitelist.

The whitelist contains:

- Default SQL whitelist of the system.
- SQL whitelist specified by the GUC parameter wlm_sql_allow_list.

14.3.93 GS_WORKLOAD_SQL_COUNT

GS_WORKLOAD_SQL_COUNT displays statistics on the number of SQL statements executed in workload Cgroups on the current node, including the number of **SELECT**, **UPDATE**, **INSERT**, and **DELETE** statements and the number of DDL, DML, and DCL statements.

Table 14-160 GS_WORKLOAD_SQL_COUNT columns

Column	Туре	Description
workload	Name	Workload Cgroup name
select_count	Bigint	Number of SELECT statements
update_count	Bigint	Number of UPDATE statements
insert_count	Bigint	Number of INSERT statements
delete_count	Bigint	Number of DELETE statements
ddl_count	Bigint	Number of DDL statements
dml_count	Bigint	Number of DML statements
dcl_count	Bigint	Number of DCL statements

14.3.94 GS_WORKLOAD_SQL_ELAPSE_TIME

GS_WORKLOAD_SQL_ELAPSE_TIME displays statistics on the response time of SQL statements in workload Cgroups on the current node, including the maximum, minimum, average, and total response time of **SELECT**, **UPDATE**, **INSERT**, and **DELETE** statements. The unit is microsecond.

Table 14-161 GS_WORKLOAD_SQL_ELAPSE_TIME columns

Column	Туре	Description
workload	Name	Workload Cgroup name
total_select_elapse	Bigint	Total response time of SELECT statements
max_select_elapse	Bigint	Maximum response time of SELECT statements
min_select_elapse	Bigint	Minimum response time of SELECT statements
avg_select_elapse	Bigint	Average response time of SELECT statements
total_update_elapse	Bigint	Total response time of UPDATE statements
max_update_elapse	Bigint	Maximum response time of UPDATE statements
min_update_elapse	Bigint	Minimum response time of UPDATE statements
avg_update_elapse	Bigint	Average response time of UPDATE statements
total_insert_elapse	Bigint	Total response time of INSERT statements
max_insert_elapse	Bigint	Maximum response time of INSERT statements
min_insert_elapse	Bigint	Minimum response time of INSERT statements
avg_insert_elapse	Bigint	Average response time of INSERT statements
total_delete_elapse	Bigint	Total response time of DELETE statements
max_delete_elapse	Bigint	Maximum response time of DELETE statements
min_delete_elapse	Bigint	Minimum response time of DELETE statements

Column	Туре	Description
avg_delete_elapse	Bigint	Average response time of DELETE statements

14.3.95 GS_WORKLOAD_TRANSACTION

GS_WORKLOAD_TRANSACTION provides transaction information about workload cgroups on a single CN. The database records the number of times that each workload Cgroup commits and rolls back transactions and the response time of transaction commitment and rollback, in microseconds.

Table 14-162 GS_WORKLOAD_TRANSACTION columns

Column	Туре	Description
workload	Name	Workload Cgroup name
commit_counter	Bigint	Number of the commits
rollback_counter	Bigint	Number of rollbacks
resp_min	Bigint	Minimum response time
resp_max	Bigint	Maximum response time
resp_avg	Bigint	Average response time
resp_total	Bigint	Total response time

14.3.96 MPP_TABLES

MPP_TABLES displays information about tables in PGXC_CLASS.

Table 14-163 MPP_TABLES columns

Name	Туре	Description
schemaname	name	Name of the schema that contains the table
tablename	name	Name of a table
tableowner	name	Owner of the table
tablespace	name	Tablespace where the table is located.
pgroup	name	Name of a node cluster.
nodeoids	oidvector_extend	List of distributed table node OIDs

14.3.97 PG_AVAILABLE_EXTENSION_VERSIONS

PG_AVAILABLE_EXTENSION_VERSIONS displays the extension versions of certain database features.

Table 14-164 PG AVAILABLE EXTENSION VERSIONS columns

Column	Туре	Description
Name	Name	Extension name
version	Text	Version name
installed	boolean	The value is true if the version of this extension is currently installed.
superuser	boolean	The value is true if only system administrators are allowed to install this extension.
relocatable	boolean	The value is true if an extension can be relocated to another schema.
schema	Name	Name of the schema that the extension must be installed into. The value is NULL if the extension is partially or fully relocatable.
requires	name[]	Names of prerequisite extensions. The value is NULL if there are no prerequisite extensions.
comment	Text	Comment string from the extension's control file

14.3.98 PG_AVAILABLE_EXTENSIONS

PG_AVAILABLE_EXTENSIONS displays the extended information about certain database features.

Table 14-165 PG_AVAILABLE_EXTENSIONS columns

Column	Туре	Description
Name	Name	Extension name.
default_version	Text	Name of default version. The value is NULL if none is specified.
installed_version	Text	Currently installed version of the extension. The value is NULL if no version is installed.
comment	Text	Comment string from the extension's control file.

14.3.99 PG_BULKLOAD_STATISTICS

On any normal node in a cluster, **PG_BULKLOAD_STATISTICS** displays the execution status of the import and export services. Each import or export service corresponds to a record. This view is accessible only to users with system administrators rights.

Table 14-166 PG_BULKLOAD_STATISTICS columns

Name	Туре	Description
node_name	Text	Node name
db_name	Text	Database name
query_id	Bigint	Query ID. It is equivalent to debug_query_id.
tid	Bigint	ID of the current thread
lwtid	Integer	Lightweight thread ID
session_id	Bigint	GDS session ID
direction	Text	Service type. The options are gds to file , gds from file , gds to pipe , gds from pipe , copy from , and copy to .
query	Text	Query statement
address	Text	Location of the foreign table used for data import and export
query_start	Timestamp with time zone	Start time of data import or export
total_bytes	Bigint	Total size of data to be processed
		This parameter is specified only when a GDS common file is to be imported and the record in the row comes from a CN. Otherwise, left this parameter unspecified.
phase	Text	Execution phase of the current service import and export. The options are INITIALIZING, TRANSFER_DATA, and RELEASE_RESOURCE.
done_lines	Bigint	Number of lines that have been transferred
done_bytes	Bigint	Number of bytes that have been transferred

14.3.100 PG_COMM_CLIENT_INFO

PG_COMM_CLIENT_INFO stores the client connection information of a single node. (You can query this view on a DN to view the information about the connection between the CN and DN.)

Table 14-167 PG_COMM_CLIENT_INFO columns

Column	Туре	Description
node_name	Text	Current node name.
арр	Text	Client application name
tid	Bigint	Thread ID of the current thread.
lwtid	Integer	Lightweight thread ID of the current thread.
query_id	Bigint	Query ID. It is equivalent to debug_query_id .
socket	Integer	It is displayed if the connection is a physical connection.
remote_ip	Text	Peer node IP address.
remote_port	Text	Peer node port.
logic_id	Integer	If the connection is a logical connection, sid is displayed. If -1 is displayed, the current connection is a physical connection.

14.3.101 PG_COMM_DELAY

PG_COMM_DELAY displays the communication library delay status for a single DN.

Table 14-168 PG_COMM_DELAY columns

Column	Туре	Description
node_name	Text	Node name
remote_name	Text	Name of the node with the maximum latency in connecting to the peer end.
remote_host	Text	IP address of the peer.
stream_num	Integer	Number of logical stream connections used by the current physical connection.
min_delay	Integer	Minimum delay of the current physical connection. The unit is microsecond.
average	Integer	Average delay of the current physical connection. The unit is microsecond.

Column	Туре	Description
max_delay	Integer	Maximum delay of the current physical connection. The unit is microsecond.
		NOTE If its value is -1, the latency detection has timed out. In this case, re-establish the connection between nodes and then perform the query.

14.3.102 PG_COMM_STATUS

PG_COMM_STATUS displays the communication library status for a single DN.

Table 14-169 PG_COMM_STATUS columns

Column	Туре	Description
node_name	Text	Node name.
rxpck/s	Integer	Receiving rate of the communication library on a node. The unit is byte/s.
txpck/s	Integer	Sending rate of the communication library on a node. The unit is byte/s.
rxkB/s	Bigint	Receiving rate of the communication library on a node. The unit is KB/s.
txkB/s	Bigint	Sending rate of the communication library on a node. The unit is KB/s.
buffer	Bigint	Size of the buffer of the Cmailbox.
memKB(libcomm)	Bigint	Communication memory size of the libcomm process, in KB.
memKB(libpq)	Bigint	Communication memory size of the libpq process, in KB.
%USED(PM)	Integer	Real-time usage of the postmaster thread.
%USED (sflow)	Integer	Real-time usage of the gs_sender_flow_controller thread.
%USED (rflow)	Integer	Real-time usage of the gs_receiver_flow_controller thread.
%USED (rloop)	Integer	Highest real-time usage among multiple gs_receivers_loop threads.
stream	Integer	Total number of used logical connections.

14.3.103 PG_COMM_RECV_STREAM

PG_COMM_RECV_STREAM displays the receiving stream status of all the communication libraries for a single DN.

Table 14-170 PG_COMM_RECV_STREAM columns

Column	Туре	Description
node_name	Text	Node name
local_tid	Bigint	ID of the thread using this stream
remote_name	Text	Name of the peer node
remote_tid	Bigint	Peer thread ID
idx	Integer	Peer DN ID in the local DN
sid	Integer	Stream ID in the physical connection
tcp_sock	Integer	TCP socket used in the stream
state	Text	 UNKNOWN: The logical connection is unknown. READY: The logical connection is ready. RUN: The logical connection receives packets normally. HOLD: The logical connection is waiting to receive packets. CLOSED: The logical connection is closed. TO_CLOSED: The logical connection is to be closed. WRITING: Data is being written.
query_id	Bigint	debug_query_id corresponding to the stream
pn_id	Integer	plan_node_id of the query executed by the stream
send_smp	Integer	smpid of the sender of the query executed by the stream
recv_smp	Integer	smpid of the receiver of the query executed by the stream
recv_bytes	Bigint	Total data volume received from the stream. The unit is byte.
time	Bigint	Current life cycle service duration of the stream. The unit is ms.
speed	Bigint	Average receiving rate of the stream. The unit is byte/s.

Column	Туре	Description	
quota	Bigint	Current communication quota value of the stream. The unit is Byte.	
buff_usize	Bigint	Current size of the data cache of the stream. The unit is byte.	

14.3.104 PG_COMM_SEND_STREAM

PG_COMM_SEND_STREAM displays the sending stream status of all the communication libraries for a single DN.

Table 14-171 PG_COMM_SEND_STREAM columns

Column	Туре	Description	
node_name	Text	Node name	
local_tid	Bigint	ID of the thread using this stream	
remote_name	Text	Name of the peer node	
remote_tid	Bigint	Peer thread ID	
idx	Integer	Peer DN ID in the local DN	
sid	Integer	Stream ID in the physical connection	
tcp_sock	Integer	TCP socket used in the stream	
state	Text	 Current status of the stream UNKNOWN: The logical connection is unknown. READY: The logical connection is ready. RUN: The logical connection sends packets normally. HOLD: The logical connection is waiting to send packets. CLOSED: The logical connection is closed. TO_CLOSED: The logical connection is to be closed. WRITING: Data is being written. 	
query_id	Bigint	debug_query_id corresponding to the stream	
pn_id	Integer	<pre>plan_node_id of the query executed by the stream</pre>	
send_smp	Integer	smpid of the sender of the query executed by the stream	

Column	Туре	Description	
recv_smp	Integer	smpid of the receiver of the query executed by the stream	
send_bytes	Bigint	Total data volume sent by the stream. The unit is Byte.	
time	Bigint	Current life cycle service duration of the stream. The unit is ms.	
speed	Bigint	Average sending rate of the stream. The unit is Byte/s.	
quota	Bigint	Current communication quota value of the stream. The unit is Byte.	
wait_quota	Bigint	Extra time generated when the stream waits the quota value. The unit is ms.	

14.3.105 PG_COMM_QUERY_SPEED

PG_COMM_QUERY_SPEED displays traffic information about all queries on a single node.

Table 14-172 PG_COMM_QUERY_SPEED columns

Column	Туре	Description	
node_name	Text	Node name	
query_id	Bigint	debug_query_id corresponding to the stream	
rxkB/s	Bigint	Receiving rate of the query stream (unit: byte/s)	
txkB/s	Bigint	Sending rate of the query stream (unit: byte/s)	
rxkB	Bigint	Total received data of the query stream (unit: byte)	
txkB	Bigint	Total sent data of the query stream (unit: byte)	
rxpck/s	Bigint	Packet receiving rate of the query (unit: packets/s)	
txpck/s	Bigint	Packet sending rate of the query (unit: packets/s)	
rxpck	Bigint	Total number of received packets of the query	

Column	Туре	Description
txpck	Bigint	Total number of sent packets of the query

14.3.106 PG_CONTROL_GROUP_CONFIG

PG_CONTROL_GROUP_CONFIG displays the Cgroup configuration information in the system.

Table 14-173 PG_CONTROL_GROUP_CONFIG columns

Column	Туре	Description
pg_control_group_config	Text	Configuration information of the Cgroup

14.3.107 PG_CURSORS

PG_CURSORS displays available cursors.

Table 14-174 PG_CURSORS columns

Column	Туре	Description	
Name	Text	Cursor name.	
statement	Text	Query statement when the cursor is declared to change.	
is_holdable	boolean	Whether the cursor is holdable (that is, it can be accessed after the transaction that declared the cursor has committed). If it is, its value is true .	
is_binary	boolean	Whether the cursor was declared BINARY. If it was, its value is true .	
is_scrollable	boolean	Whether the cursor is scrollable (that is, it allows rows to be retrieved in a nonsequential manner). If it is, its value is true .	
creation_tim e	Timestamp with time zone	Timestamp of the cursor.	

14.3.108 PG_EXT_STATS

PG_EXT_STATS displays extension statistics stored in the **PG_STATISTIC_EXT** table. The extension statistics means multiple columns of statistics.

Table 14-175 PG_EXT_STATS columns

Column	Туре	Reference	Description
schemaname	Name	PG_NAMESP ACE.nspname	Name of the schema that contains a table
tablename	Name	PG_CLASS.rel name	Name of a table
attname	int2vector	PG_STATISTI C_EXT.stakey	Indicates the columns to be combined for collecting statistics.
inherited	boolean	-	Includes inherited sub-columns if the value is true ; otherwise, indicates the column in a specified table.
null_frac	Real	-	Percentage of column combinations that are null to all records
avg_width	Integer	-	Average width of column combinations. The unit is byte.
n_distinct	Real	-	Estimated number of distinct values in a column combination if the value is greater than 0
			 Negative of the number of distinct values divided by the number of rows if the value is less than 0
			The number of distinct values is unknown if the value is 0.
			NOTE The negated form is used when ANALYZE believes that the number of distinct values is likely to increase as the table grows.
			The positive form is used when the column seems to have a fixed number of possible values. For example, -1 indicates that the number of distinct values is the same as the number of rows for a column combination.

Column	Туре	Reference	Description
n_dndistinct	Real	-	Number of unique not-null data values in the dn1 column combination
			• Exact number of distinct values if the value is greater than 0
			 Negative of the number of distinct values divided by the number of rows if the value is less than 0 For example, if a value in a column combination appears twice in average, n_dndistinct equals -0.5.
			The number of distinct values is unknown if the value is 0.
most_commo n_vals	anyarray	-	List of the most common values in a column combination. If this combination does not have the most common values, most_common_vals_null will be NULL. None of the most common values in most_common_vals is NULL.
most_commo n_freqs	real[]	-	List of the frequencies of the most common values, that is, the number of occurrences of each value divided by the total number of rows. (NULL if most_common_vals is NULL)
most_commo n_vals_null	anyarray	-	List of the most common values in a column combination. If this combination does not have the most common values, most_common_vals_null will be NULL. At least one of the common values in most_common_vals_null is NULL.
most_commo n_freqs_null	real[]	-	List of the frequencies of the most common values, that is, the number of occurrences of each value divided by the total number of rows. (NULL if most_common_vals_null is NULL)

14.3.109 PG_GET_INVALID_BACKENDS

PG_GET_INVALID_BACKENDS displays the information about backend threads on the CN that are connected to the current standby DN.

Table 14-176 PG_GET_INVALID_BACKENDS columns

Column	Туре	Description	
pid	Bigint	Thread ID	
node_name	Text	Node information connected to the backend thread	
dbname	Name	Name of the connected database	
backend_start	Timestamp with time zone	Backend thread startup time	
query	Text	Query statement performed by the backend thread	

14.3.110 PG_GET_SENDERS_CATCHUP_TIME

PG_GET_SENDERS_CATCHUP_TIME displays the catchup information of the currently active primary/standby instance sending thread on a single DN.

Table 14-177 PG_GET_SENDERS_CATCHUP_TIME columns

Column	Туре	Description
pid	Bigint	Current sender thread ID
lwpid	Integer	Current sender lwpid
local_role	Text	Local role
peer_role	Text	Peer role
state	Text	Current sender's replication status
type	Text	Current sender type
catchup_start	Timestamp with time zone	Startup time of a catchup task
catchup_end	Timestamp with time zone	End time of a catchup task
catchup_type	Text	Catchup task type, full or incremental
catchup_bcm_filen ame	Text	BCM file executed by the current catchup task

Column	Туре	Description
catchup_bcm_finis hed	Integer	Number of BCM files completed by a catchup task
catchup_bcm_total	Integer	Total number of BCM files to be operated by a catchup task
catchup_percent	Text	Completion percentage of a catchup task
catchup_remaining _time	Text	Estimated remaining time of a catchup task

14.3.111 PG GLOBAL TEMP ATTACHED PIDS

This view displays information about sessions of resources occupied by global temporary tables on the current node. This view is supported only by clusters of 8.2.1.220 and later versions.

Table 14-178 PG_GLOBAL_TEMP_ATTACHED_PIDS columns

Column	Туре	Description	
schemaname	Name	Schema name	
tablename	Name	Table name	
pid	Bigint	PID of a session	

14.3.112 PG GROUP

PG_GROUP displays the database role authentication and the relationship between roles.

Table 14-179 PG GROUP columns

Column	Туре	Description	
groname	Name	Group name	
grosysid	OID	Group ID	
grolist	oid[]	An array, including all the role IDs in this group	

14.3.113 PG INDEXES

PG_INDEXES displays access to useful information about each index in the database.

Column	Туре	Reference	Description	
schemana me	Name	PG_NAMESP ACE.nspname	Name of the schema that contains tables and indexes	
tablenam e	Name	PG_CLASS.rel name	Name of the table for which the index serves	
indexnam e	Name	PG_CLASS.rel name	Index name	
tablespac e	Name	PG_TABLESPA CE.spcname	Name of the tablespace that contains the index	
indexdef	Text	N/A	Index definition (a reconstructed	

Table 14-180 PG INDEXES columns

Example

Query the index information about a specified table.

Query information about indexes of all tables in a specified schema in the current database.

```
SELECT tablename, indexname, indexdef FROM pg_indexes WHERE schemaname = 'public' ORDER BY
tablename, indexname;
tablename | indexname |
                                                      indexdef
-----+-----
books | books_pkey | CREATE UNIQUE INDEX books_pkey ON books USING btree (id) TABLESPACE
pg_default
books | idx_books_tags_gin | CREATE INDEX idx_books_tags_gin ON books USING gin (tags)
TABLESPACE pg_default
                         | CREATE UNIQUE INDEX c_custkey_key ON customer USING btree
customer | c_custkey_key
(c_custkey, c_name) TABLESPACE pg_default
mytable | idx_mytable_id | CREATE INDEX idx_mytable_id ON mytable USING btree (id) TABLESPACE
pg_default
test1 | idx test id | CREATE INDEX idx test id ON test1 USING btree (id) TABLESPACE pg default
                      CREATE UNIQUE INDEX v0_pkey ON v0 USING btree (c) TABLESPACE pg_default
v0
      | v0_pkey
(6 rows)
```

14.3.114 PG JOB

PG_JOB displays detailed information about scheduled tasks created by users.

The **PG_JOB** view replaces the **PG_JOB** system catalog in earlier versions and provides forward compatibility with earlier versions. The original **PG_JOB** system catalog is changed to the **PG_JOBS** system catalog. For details about **PG_JOBS**, see **PG_JOBS**.

Table 14-181 PG_JOB columns

Column	Туре	Description	
job_id	Bigint	Job ID	
current_postg res_pid	Bigint	If the current job has been executed, the PostgreSQL thread ID of this job is recorded. The default value is -1, indicating that the task is not executed or has been executed.	
log_user	Name	User name of the job creator	
priv_user	Name	User name of the job executor	
dbname	Name	Name of the database where the job is executed	
node_name	Name	CN node on which the job will be created and executed	
job_status	Text	Status of the current job. The value range is r , s , f , d , p , w , or l . The default value is s . The indications are as follows: • r=running	
		s=successfully finished	
		• f=job failed	
		d=disable	
		p=pending	
		w=waiting	
		l=launching	
		NOTE	
		 Note: When you disable a scheduled task (by setting job_queue_processes to 0), the thread monitor the job execution is not started, and the job_status will not be updated. You can ignore the job_status. 	
		 Only when the scheduled task function is enabled (that is, when job_queue_processes is not 0), the system updates the value of job_status based on the real-time job status. 	
start_date	Timestamp without time zone	Start time of the first job execution, precise to millisecond	
next_run_date	Timestamp without time zone	Scheduled time of the next job execution, accurate to millisecond	
failure_count	Smallint	Number of consecutive failures	
interval	Text	Job execution interval	

Column	Туре	Description
last_start_dat e	Timestamp without time zone	Start time of the last job execution, accurate to millisecond
last_end_date	Timestamp without time zone	End time of the last job execution, accurate to millisecond
last_suc_date	Timestamp without time zone	Start time of the last successful job execution, accurate to millisecond
this_run_date	Timestamp without time zone	Start time of the ongoing job execution, accurate to millisecond
nspname	Name	Name of the namespace where a job is running
what	Text	Job content

14.3.115 PG_JOB_PROC

The PG_JOB_PROC view replaces the PG_JOB_PROC system catalog in earlier versions and provides forward compatibility with earlier versions. The original PG_JOB_PROC and PG_JOB system catalogs are merged into the PG_JOBS system catalog in the current version. For details about the PG_JOBS system catalog, see PG_JOBS.

Table 14-182 PG_JOB_PROC columns

Column	Туре	Description
job_id	Bigint	Job ID
what	Text	Job content

14.3.116 PG_JOB_SINGLE

PG_JOB_SINGLE displays job information about the current node.

Table 14-183 PG_JOB_SINGLE columns

Column	Туре	Description
job_id	Bigint	Job ID

Column	Туре	Description	
current_postg res_pid	Bigint	If the current job has been executed, the PostgreSQL thread ID of this job is recorded. The default value is -1, indicating that the task is not executed or has been executed.	
log_user	Name	User name of the job creator	
priv_user	Name	User name of the job executor	
dbname	Name	Name of the database where the job is executed	
node_name	Name	CN node on which the job will be created and executed	
job_status	Text	Status of the current job. The value range is r , s , f , d , p , w , or l . The default value is s . The indications are as follows:	
		• r=running	
		s=successfully finished	
		f=job failed	
		d=disable	
		p=pending	
		w=waiting	
		• l=launching	
		NOTE	
		 Note: When you disable a scheduled task (by setting job_queue_processes to 0), the thread monitor the job execution is not started, and the job_status will not be updated. You can ignore the job_status. 	
		 Only when the scheduled task function is enabled (that is, when job_queue_processes is not 0), the system updates the value of job_status based on the real-time job status. 	
start_date	Timestamp without time zone	Start time of the first job execution, precise to millisecond	
next_run_date	Timestamp without time zone	Scheduled time of the next job execution, accurate to millisecond	
failure_count	Smallint	Number of consecutive failures.	
interval	Text	Job execution interval	
last_start_dat e	Timestamp without time zone	Start time of the last job execution, accurate to millisecond	

Column	Туре	Description	
last_end_date	Timestamp without time zone	End time of the last job execution, accurate to millisecond	
last_suc_date	Timestamp without time zone	Start time of the last successful job execution, accurate to millisecond	
this_run_date	Timestamp without time zone	Start time of the ongoing job execution, accurate to millisecond	
nspname	Name	Name of the namespace where a job is running	
what	Text	Job content	

14.3.117 PG_LIFECYCLE_DATA_DISTRIBUTE

PG_LIFECYCLE_DATA_DISTRIBUTE displays the distribution of cold and hot data in a multi-temperature table of OBS.

Table 14-184 PG_LIFECYCLE_DATA_DISTRIBUTE columns

Column	Туре	Description	
schemaname	Name	Schema name	
tablename	Name	Current table name	
nodename	Name	Node name	
hotpartition	Text	Hot partition on the DN	
coldpartition	Text	Cold partition on the DN	
switchablepar tition	Text	Switchable partition on the DN	
hotdatasize	Text	Data size of the hot partition on the DN	
colddatasize	Text	Data size of the cold partition on the DN	
switchabledat asize	Text	Data size of the switchable partition on the DN	

14.3.118 PG_LOCKS

PG_LOCKS displays information about the locks held by open transactions.

Table 14-185 PG_LOCKS columns

Column	Туре	Reference	Description
locktype	Text	N/A	Type of the locked object: relation, extend, page, tuple, transactionid, virtualxid, object, userlock, and advisory
database	OID	PG_DATABAS E.oid	 OID of the database in which the locked target exists The OID is 0 if the target is a shared object. The OID is NULL if the locked target is a transaction.
relation	OID	PG_CLASS.oid	OID of the relationship targeted by the lock. The value is NULL if the object is neither a relationship nor part of a relationship.
page	Integer	N/A	Page number targeted by the lock within the relationship. If the object is neither a relation page nor row page, the value is NULL .
tuple	Smallint	N/A	Row number targeted by the lock within the page. If the object is not a row, the value is NULL .
virtualxid	Text	N/A	Virtual ID of the transaction targeted by the lock. If the object is not a virtual transaction ID, the value is NULL .
transactionid	Xid	N/A	ID of the transaction targeted by the lock. If the object is not a transaction ID, the value is NULL .
classid	OID	PG_CLASS.oid	OID of the system table that contains the object. If the object is not a general database object, the value is NULL .
objid	OID	N/A	OID of the lock target within its system table. If the target is not a general database object, the value is NULL .
objsubid	Smallint	N/A	Column number for a column in the table. The value is 0 if the target is some other object type. If the object is not a general database object, the value is NULL .

Column	Туре	Reference	Description
virtualtransac tion	Text	N/A	Virtual ID of the transaction holding or awaiting this lock
pid	Bigint	N/A	Logical ID of the server thread holding or awaiting this lock. This is NULL if the lock is held by a prepared transaction.
mode	Text	N/A	Lock mode held or desired by this thread For more information about lock modes, see LOCK.
granted	boolean	N/A	 The value is true if the lock is a held lock. The value is false if the lock is an awaited lock.
fastpath	boolean	N/A	Whether the lock is obtained through fast-path (true) or main lock table (false)
waittime	Timestam p with time zone	N/A	Timestamp when the lock wait starts. This column is available only in clusters of version 9.1.0.200 or later.
holdtime	Timestam p with time zone	N/A	Timestamp when the lock starts to be held. This column is available only in clusters of version 9.1.0.200 or later.

14.3.119 PG_LWLOCKS

PG_LWLOCKS provides information on lightweight locks currently held or being waited for by the current instance. This view is supported only by 9.1.0.200 and later cluster versions.

Table 14-186 PG_LWLOCKS columns

Column	Туре	Description
pid	Bigint	ID of the backend thread.
query_id	Bigint	ID of a query.
lwtid	Integer	Lightweight thread ID of the backend thread.
reqlockid	Integer	ID of the lightweight lock that is being requested by the current thread.

Column	Туре	Description	
reqlock	Text	Name of the lightweight lock corresponding to reqlockid .	
heldlocknums	Integer	Number of lightweight locks obtained by the current thread.	
heldlockid	Integer	Lightweight lock ID obtained by the current thread.	
heldlock	Text	Name of the lightweight lock corresponding to heldlockid .	
heldlockmode	Text	Lightweight lock mode corresponding to heldlockid.	

Example

Use the **PG_LWLOCKS** view to query information about lightweight locks that are being held or waiting for the current instance.

14.3.120 PG_NODE_ENV

PG_NODE_ENVO displays the environmental variable information about the current node.

Table 14-187 PG NODE ENV columns

Column	Туре	Description
node_name	Text	Name of the node
host	Text	Host name of the node
process	Integer	Number of the node process
port	Integer	Port ID of the node
installpath	Text	Installation directory of current node
datapath	Text	Data directory of the node
log_directory	Text	Log directory of the node

14.3.121 PG_OS_THREADS

PG_OS_THREADS displays the status information about all the threads under the current node.

Table 14-188 PG_OS_THREADS columns

Column	Туре	Description
node_name	Text	Name of the node
pid	Bigint	Thread number running under the current node process
lwpid	Integer	Lightweight thread IDs corresponding to the PIDs
thread_name	Text	Thread names corresponding to the PIDs
creation_time	Timestamp with time zone	Creation time of the threads corresponding to the PIDs

14.3.122 PG_POOLER_STATUS

PG_POOLER_STATUS displays the cache connection status in the pooler. **PG_POOLER_STATUS** can only query on the CN, and displays the connection cache information about the pooler module.

Table 14-189 PG_POOLER_STATUS columns

Column	Туре	Description
database	Text	Database name
user_name	Text	Username
tid	Bigint	ID of the thread used for the connection to the CN
node_oid	Bigint	OID of the node connected
node_name	Name	Name of the node connected
in_use	boolean	Whether the connection is in use. The options are: • t (true): The connection is in use. • f (false): The connection is not in use.
fdsock	Bigint	Peer socket
remote_pid	Bigint	Peer thread ID

Column	Туре	Description
session_params	Text	GUC session parameter delivered by the connection

Example

View information about the connection pool **pooler**:

```
select database,user_name,node_name,in_use,count(*) from pg_pooler_status group by 1, 2, 3, 4 order by 5
desc limit 50;
database | user_name | node_name | in_use | count
mydbdemo | user3 | cn_5001 | f |
mydbdemo | user3 | dn_6005_6006 | t
mydbdemo | user3 | dn_6001_6002 | t
mydbdemo | user3
                    | dn_6003_6004 | f
                   | dn_6003_6004 | t
                                             2
mydbdemo | user3
mydbdemo | user3 | dn_6005_6006 | f
gaussdb | user3 | dn_6003_6004 | f
mydbdemo | user3 | cn_5001 | t
                                           1
music | user2 | dn_6003_6004 | f
music | user2 | dn_6005_6006 | f
gaussdb | user1 | dn_6005_6006 | f
(13 rows)
```

14.3.123 PG_PREPARED_STATEMENTS

PG_PREPARED_STATEMENTS displays all prepared statements that are available in the current session.

Table 14-190 PG PREPARED STATEMENTS columns

Column	Туре	Description
Name	Text	Identifier of a prepared statement.
statement	Text	Query string for creating this prepared statement. For prepared statements created through SQL, this is the PREPARE statement submitted by the client. For prepared statements created through the frontend/backend protocol, this is the text of the prepared statement itself.
prepare_time	Timestamp with time zone	Timestamp when the prepared statement is created.
parameter_ty pes	regtype[]	Expected parameter types for the prepared statement in the form of an array of regtype . The OID corresponding to an element of this array can be obtained by casting the regtype value to oid.

Column	Туре	Description
from_sql	boolean	• The value is true if the prepared statement is created using the PREPARE statement.
		The value is false if the instance is created using the frontend/backend protocol.

14.3.124 PG_PREPARED_XACTS

PG_PREPARED_XACTS displays information about transactions that are currently prepared for two-phase commit.

Table 14-191 PG_PREPARED_XACTS columns

Column	Туре	Reference	Description
transaction	Xid	N/A	Numeric transaction identifier of the prepared transaction
gid	Text	N/A	Global transaction identifier that was assigned to the transaction
prepared	Timestamp with time zone	N/A	Time at which the transaction is prepared for commit
owner	Name	PG_AUTHID.rolna me	Name of the user that executes the transaction
database	Name	PG_DATABASE.da tname	Name of the database in which the transaction is executed

14.3.125 PG_PUBLICATION_TABLES

PG_PUBLICATION_TABLES displays the mapping between a publication and its published tables. Unlike the underlying system catalog **PG_PUBLICATION_REL**, this view expands the publications defined as **FOR ALL TABLES** and **FOR ALL TABLES** in **SCHEMA**, in which each publishable table has a row. This view is supported only by clusters of version 8.2.0.100 or later.

Table 14-192 PG_PUBLICATION_TABLES columns

Column	Туре	Description
pubname	Name	Publication name

Column	Туре	Description
schemaname	Name	Name of the schema of a table
tablename	Name	Table name

Examples

Query all published tables.

14.3.126 PG_QUERYBAND_ACTION

PG_QUERYBAND_ACTION displays information about the object associated with **query_band** and the **query_band** query order.

Table 14-193 PG_QUERYBAND_ACTION columns

Column	Туре	Description
qband	Text	query_band key-value pairs
respool_id	OID	OID of the resource pool associated with query_band
respool	Text	Name of the resource pool associated with query_band
priority	Text	Intra-queue priority associated with query_band
qborder	Integer	query_band query order

14.3.127 PG_REPLICATION_SLOTS

PG_REPLICATION_SLOTS displays the replication node information.

Table 14-194 PG_REPLICATION_SLOTS columns

Column	Туре	Description
slot_name	Text	Name of a replication node
plugin	Name	Name of the output plug-in of the logical replication slot
slot_type	Text	Type of a replication node

Column	Туре	Description
datoid	OID	OID of the database on the replication node
database	Name	Name of the database on the replication node
active	boolean	Whether the replication node is active
xmin	Xid	Transaction ID of the replication node
catalog_xmin	Text	ID of the earliest-decoded transaction corresponding to the logical replication slot
restart_lsn	Text	Xlog file information on the replication node
dummy_stand by	boolean	Whether the replication node is the dummy standby node

14.3.128 PG_ROLES

PG_ROLES displays information about database roles.

Table 14-195 PG_ROLES columns

Column	Туре	Reference	Description
rolname	Name	N/A	Role name
rolsuper	boolean	N/A	Whether the role is the initial system administrator with the highest permission
rolinherit	boolean	N/A	Whether the role inherits permissions for this type of roles
rolcreaterole	boolean	N/A	Whether the role can create other roles
rolcreatedb	boolean	N/A	Whether the role can create databases
rolcatupdate	boolean	N/A	Whether the role can update system tables directly. Only the initial system administrator whose usesysid is 10 has this permission. It is not available for other users.
rolcanlogin	boolean	N/A	Whether the role can log in to the database
rolreplication	boolean	N/A	Whether the role can be replicated
rolauditadmi n	boolean	N/A	Whether the role is an audit system administrator

Column	Туре	Reference	Description
rolsystemad min	boolean	N/A	Whether the role is a system administrator
rolconnlimit	Integer	N/A	Limits the maximum number of concurrent connections of a user on a CN1 indicates no limit.
rolpassword	Text	N/A	Not the password (always reads as ********)
rolvalidbegin	Timestamp with time zone	N/A	Account validity start time; null if no start time
rolvaliduntil	Timestamp with time zone	N/A	Password expiry time; null if no expiration
rolrespool	Name	N/A	Resource pool that a user can use
rolparentid	OID	PG_AUTHI D.rolparenti d	OID of a group user to which the user belongs
roltabspace	Text	N/A	The storage space of the user permanent table.
roltempspace	Text	N/A	The storage space of the user temporary table.
rolspillspace	Text	N/A	The operator disk flushing space of the user.
rolconfig	Text[]	N/A	Session defaults for runtime configuration variables
OID	OID	PG_AUTHI D.oid	ID of the role
roluseft	boolean	PG_AUTHI D.roluseft	Whether the role can perform operations on foreign tables
nodegroup	Name	N/A	Name of the logical cluster associated with the role. If no logical cluster is associated, this column is left empty.

14.3.129 PG_RULES

PG_RULES displays information about rewrite rules.

Table 14-196 PG_RULES columns

Column	Туре	Description
schemaname	Name	Name of the schema that contains the table
tablename	Name	Name of the table the rule is for
rulename	Name Rule name	
definition	Text	Rule definition (a reconstructed creation command)

14.3.130 PG_RUNNING_XACTS

PG_RUNNING_XACTS displays information about running transactions on the current node.

Table 14-197 PG_RUNNING_XACTS columns

Column	Туре	Description
handle	Integer	Handle corresponding to the transaction in GTM
gxid	Xid	Transaction ID
state	tinyint	Transaction status (3: prepared or 0: starting)
node	Text	Node name
xmin	Xid	Minimum transaction ID xmin on the node
vacuum	boolean	Whether the current transaction is lazy vacuum
timeline	Bigint	Number of database restarts
prepare_xid	Xid	Transaction ID in the prepared status. If the status is not prepared , the value is 0 .
pid	Bigint	Thread ID corresponding to the transaction
next_xid	Xid	Transaction ID sent from a CN to a DN

14.3.131 PG SECLABELS

PG_SECLABELS displays information about security labels.

Table 14-198 PG_SECLABELS columns

Column	Туре	Reference	Description
objoid	OID	Any OID column	OID of the object this security label pertains to
classoid	OID	PG_CLASS.oid	OID of the system table that contains the object
objsubid	Intege r	N/A	For a security label on a table column, this is the column number (the objoid and classoid refer to the table itself). For all other object types, this column is 0 .
objtype	Text	N/A	Type of the object to which this label applies
objnamespac e	OID	PG_NAMESPACE.oid	OID of the namespace for this object, if applicable; otherwise NULL.
objname	Text	N/A	Name of the object to which the label applies
provider	Text	PG_SECLABEL.provider	Label provider associated with this label
label	Text	PG_SECLABEL.label	Security label applied to this object

14.3.132 PG_SEQUENCES

PG_SEQUENCES displays the sequence attributes on which the current user has permissions. This view is supported only by clusters of version 9.1.0 or later.

Table 14-199 PG_SEQUENCES columns

Column	Туре	Description
schemaname	Name	Name of the namespace.
sequencename	Name	Name of the sequence
sequenceowner	Name	Owner of the sequence
start_value	Bigint	Start value of the sequence.
min_value	Bigint	Minimum value generated by the sequence.
max_value	Bigint	Maximum value generated by the sequence.

Column	Туре	Description
increment_by	Bigint	Amount by which the generated value increases each time in a sequence.
cycle	Boolean	If set to true , the sequence value restarts from the minimum value after reaching the maximum value. If set to false , the sequence value stops generating after reaching the maximum value.
cache_size	Bigint	Size of the sequence cache value.
last_value	Bigint	Most recently generated value of the sequence.

14.3.133 PG_SESSION_WLMSTAT

PG_SESSION_WLMSTAT displays the corresponding load management information about the task currently executed by the user.

Table 14-200 PG_SESSION_WLMSTAT columns

Column	Туре	Description
datid	OID	OID of the database this backend is connected to
datname	Name	Name of the database the backend is connected to
threadid	Bigint	ID of the backend thread
processid	Integer	Thread PID of the backend
usesysid	OID	OID of the user who logged in to the backend
appname	Text	Name of the application that is connected to the backend
usename	Name	Name of the user logged in to the backend
priority	Bigint	Priority of Cgroup where the statement is located
attribute	Text	 Ordinary: default attribute of a statement before it is parsed by the database Simple: simple statements Complicated: complicated statements Internal: internal statement of the database
block_time	Bigint	Pending duration of the statements by now (unit: s)

Column	Туре	Description	
elapsed_time	Bigint	Actual execution duration of the statements by now (unit: s)	
total_cpu_time	Bigint	Total CPU usage duration of the statement on the DN in the last period (unit: s)	
cpu_skew_perce nt	Integer	CPU usage inclination ratio of the statement on the DN in the last period	
statement_mem	Integer	Estimated memory required for statement execution. This column is reserved.	
active_points	Integer	Number of concurrently active points occupied by the statement in the resource pool	
dop_value	Integer	DOP value obtained by the statement from the resource pool	
control_group	Text	Cgroup currently used by the statement	
status	Text	 Status of a statement, including: pending running finished (If enqueue is set to StoredProc or Transaction, this state indicates that only some of the jobs in the statement have been executed. This state persists until the finish of this statement.) aborted: terminated unexpectedly active: normal status except for those above unknown: unknown status 	
enqueue	Text	Current queuing status of the statements, including: • Global: global queuing. • Respool: resource pool queuing. • CentralQueue: queuing on the CCN • Transaction: being in a transaction block • StoredProc: being in a stored procedure • None: not in a queue • Forced None: being forcibly executed (transaction block statement or stored procedure statement are) because the statement waiting time exceeds the specified value	
resource_pool	Name	Current resource pool where the statements are located.	

Column	Туре	Description
query	Text	Text of this backend's most recent query If state is active , this column shows the executing query. In all other states, it shows the last query that was executed.
isplana	Bool	In logical cluster mode, indicates whether a statement occupies the resources of other logical clusters. The default value is f , indicating that resources of other logical clusters are not occupied.
node_group	Text	Logical cluster of the user running the statement
lane	Text	Fast or slow lane for statement queries. • fast: fast lane • slow: slow lane • none: not controlled

14.3.134 PG_SESSION_IOSTAT

PG_SESSION_IOSTAT has been discarded in version 8.1.2 and is reserved for compatibility with earlier versions. This view is invalid in the current version. You can use **PGXC_WLM_SESSION_STATISTICS** to view load management information about jobs being executed on all CNs.

Table 14-201 PG_SESSION_IOSTAT columns

Column	Туре	Description	
query_id	Bigint	Job ID	
mincurriops	Integer	Minimum I/O of the current job across DNs	
maxcurriops	Integer	Maximum I/O of the current job across DNs	
minpeakiops	Integer	Minimum peak I/O of the current job across DNs	
maxpeakiops	Integer	Maximum peak I/O of the current job across DNs	
io_limits	Integer	io_limits set for the job	
io_priority	Text	io_priority set for the job	
query	Text	Job	
node_group	Text	Logical cluster of the user running the job	

14.3.135 PG_SETTINGS

PG_SETTINGS displays information about parameters of the running database.

Table 14-202 PG_SETTINGS columns

Column	Туре	Description
Name	Text	Parameter name
setting	Text	Current value of the parameter
unit	Text	Implicit unit of the parameter
category	Text	Logical group of the parameter
short_desc	Text	Brief description of the parameter
extra_desc	Text	Detailed description of the parameter
context	Text	Context of parameter values including internal, postmaster, sighup, backend, superuser, and user
vartype	Text	Parameter type. It can be bool , enum , integer , real , or string .
source	Text	Method of assigning the parameter value
min_val	Text	Minimum value of the parameter. If the parameter type is not numeric data, the value of this column is null.
max_val	Text	Maximum value of the parameter. If the parameter type is not numeric data, the value of this column is null.
enumvals	Text[]	Valid values of an enum-typed parameter. If the parameter type is not enum, the value of this column is null.
boot_val	Text	Default parameter value used upon the database startup
reset_val	Text	Default parameter value used upon the database reset
sourcefile	Text	Configuration file used to set parameter values. If parameter values are not configured using the configuration file, the value of this column is null.
sourceline	Integer	Row number of the configuration file for setting parameter values. If parameter values are not configured using the configuration file, the value of this column is null.

14.3.136 PG_SHADOW

PG_SHADOW displays properties of all roles that are marked as **rolcanlogin** in **PG_AUTHID**.

This view is not readable to all users because it contains passwords. **PG_USER** is a publicly readable view on **PG_SHADOW** that blanks out the password column.

Table 14-203 PG_SHADOW columns

Column	Туре	Reference	Description
usename	Name	PG_AUTHID.rolnam	User name
usesysid	OID	PG_AUTHID.oid	ID of a user
usecreated b	boolea n	-	Indicates that the user can create databases.
usesuper	boolea n	-	Indicates that the user is an administrator.
usecatupd	boolea n	-	Indicates that the user can update system catalogs. Even the system administrator cannot do this unless this column is true .
userepl	boolea n	-	User can initiate streaming replication and put the system in and out of backup mode.
passwd	Text	-	Password (possibly encrypted); null if none. See PG_AUTHID for details about how encrypted passwords are stored.
valbegin	Timesta mp with time zone	-	Account validity start time; null if no start time
valuntil	Timesta mp with time zone	-	Password expiry time; null if no expiration
respool	Name	-	Resource pool used by the user
parent	OID	-	Parent resource pool
spacelimit	Text	-	The storage space of the permanent table.

Column	Туре	Reference	Description
tempspaceli mit	Text	-	The storage space of the temporary table.
spillspaceli mit	Text	-	The operator disk flushing space.
useconfig	text[]	-	Session defaults for runtime configuration variables

14.3.137 PG_SHARED_MEMORY_DETAIL

PG_SHARED_MEMORY_DETAIL displays usage information about all the shared memory contexts.

Table 14-204 PG_SHARED_MEMORY_DETAIL columns

Column	Туре	Description	
contextname	Text	Name of the memory context.	
level	Smallint	Hierarchy of the memory context.	
parent	Text	Parent memory context.	
totalsize	Bigint	Total size of the shared memory, in bytes.	
freesize	Bigint	Remaining size of the shared memory, in bytes.	
usedsize	Bigint	Used size of the shared memory, in bytes.	

14.3.138 PG_STATS

PG_STATS displays the single-column statistics stored in the **pg_statistic** table.

Table 14-205 PG_STATS columns

Column	Туре	Reference	Description
schemaname	Name	PG_NAMESP ACE.nspname	Name of the schema that contains the table
tablename	Name	PG_CLASS.rel name	Name of the table
attname	Name	PG_ATTRIBU TE.attname	Column name

Column	Туре	Reference	Description
inherited	boolean	-	Includes inherited sub-columns if the value is true ; otherwise, indicates the column in a specified table.
null_frac	Real	-	Percentage of column entries that are null
avg_width	Integer	-	Average width in bytes of column's entries
n_distinct	Real		 Estimated number of distinct values in the column if the value is greater than 0 Negative of the number of distinct values divided by the number of rows if the value is less than 0 The negated form is used when ANALYZE believes that the number of distinct values is likely to increase as the table grows. The positive form is used when the column seems to have a fixed number of possible values. For example, -1 indicates a unique column in which the number of
			distinct values is the same as the number of rows.
n_dndistinct	Real	-	 Number of unique non-null data values in the dn1 column Exact number of distinct values if the value is greater than 0 Negative of the number of distinct values divided by the number of rows if the value is less than 0 (For example, if the value of a column appears twice in average, set n_dndistinct=-0.5.) The number of distinct values is unknown if the value is 0.
most_commo n_vals	anyarray	-	List of the most common values in a column. If this combination does not have the most common values, it will be NULL .

Column	Туре	Reference	Description
most_commo n_freqs	real[]	-	List of the frequencies of the most common values, that is, the number of occurrences of each value divided by the total number of rows. (NULL if most_common_vals is NULL)
histogram_bo unds	anyarray	-	List of values that divide the column's values into groups of equal proportion. The values in most_common_vals, if present, are omitted from this histogram calculation. This field is null if the field data type does not have a < operator or if the most_common_vals list accounts for the entire population.
correlation	Real	-	Statistical correlation between physical row ordering and logical ordering of the column values. It ranges from -1 to +1. When the value is near to -1 or +1, an index scan on the column is estimated to be cheaper than when it is near to zero, due to reduction of random access to the disk. This column is null if the column data type does not have a < operator.
most_commo n_elems	anyarray	-	Specifies a list of non-null element values most often appearing.
most_commo n_elem_freqs	real[]	-	Specifies a list of the frequencies of the most common element values.
elem_count_h istogram	real[]	-	Histogram of the counts of distinct non-null element values.

14.3.139 PG_STAT_ACTIVITY

PG_STAT_ACTIVITY displays information about the current user's queries. If you have the rights of an administrator or the preset role, you can view all information about user queries.

Table 14-206 PG_STAT_ACTIVITY columns

Column	Туре	Description
datid	OID	OID of the database that the user session connects to in the backend
datname	Name	Name of the database that the user session connects to in the backend
pid	Bigint	Backend thread ID
lwtid	Integer	Lightweight thread ID
usesysid	OID	OID of the user logging in to the backend
usename	Name	Username for logging in to the backend.
application_name	Text	Name of the application connected to the backend
client_addr	inet	IP address of the client connected to the backend If this column is null, it indicates either that the client is connected via a Unix socket on the server machine or that this is an internal process such as autovacuum.
client_hostname	Text	Host name of the connected client, as reported by a reverse DNS lookup of client_addr. This column will only be non-null for IP connections, and only when log_hostname is enabled.
client_port	Integer	TCP port number that the client uses for communication with this backend, or -1 if a Unix socket is used
backend_start	Timestamp with time zone	Startup time of the backend process, that is, the time when the client connects to the server.
xact_start	Timestamp with time zone	Time when the current transaction was started, or NULL if no transaction is active. If the current query is the first of its transaction, this column is equal to the query_start column.
query_start	Timestamp with time zone	Time when the currently active query was started, or if state is not active , when the last query was started

Column	Туре	Description
state_change	Timestamp with time zone	Time for the last status change
waiting	boolean	The value is t if the backend is waiting for a lock or node. Otherwise, the value is f .
enqueue	Text	 Queuing status of a statement. Its value can be: waiting in global queue: The statement is queuing in the global concurrent queue. The number of concurrent statements exceeds the value of max_active_statements configured for a single CN. waiting in respool queue: The statement is queuing in the resource pool and the concurrency of simple jobs is limited. The main reason is that the concurrency of simple jobs exceeds the upper limit max_dop of the fast track. waiting in ccn queue: The job is in the CCN queue, which may be global memory queuing, slow lane memory queuing, or concurrent queuing. The scenarios are: The available global memory exceeds the upper limit, the job is queuing in the global memory queue. Concurrent requests on the slow lane in the resource pool exceed the upper limit, which is specified by active_statements. The slow lane memory of the resource pool exceeds the upper limit, that is, the estimated memory of concurrent jobs in the resource pool exceeds the upper limit specified by mem_percent. Empty or no waiting queue: The statement is running.

Column	Туре	Description
state	Text	Current overall state of this backend. Its value can be: • active: The backend is executing queries. • idle: The backend is waiting for new client commands. • idle in transaction: The backend is in a transaction, but there is no statement being executed in the transaction.
		• idle in transaction (aborted): The backend is in a transaction, but there are statements failed in the transaction.
		fastpath function call: The backend is executing a fast-path function.
		• disabled : This state is reported if track_activities is disabled in this backend.
		NOTE Common users can view only their own session status. The state information of other accounts is empty.
resource_pool	Name	Resource pool used by the user
stmt_type	Text	Statement type
query_id	Bigint	ID of a query
query	Text	Text of the most recent query in this backend If state is active , this column shows the running query. In all other states, it shows the last query that was executed.
connection_info	Text	A string in JSON format recording the driver type, driver version, driver deployment path, and process owner of the connected database (for details, see connection_info).

14.3.140 PG_STAT_ALL_INDEXES

PG_STAT_ALL_INDEXES displays statistics about all accesses to a specific index in the current database.

Indexes can be used via either simple index scans or "bitmap" index scans. Bitmap scans can combine the output of multiple indexes using AND or OR rules, but

combining independent row fetching with specific indexes is challenging. Consequently, a bitmap scan increases the index count in pg_stat_all_indexes.idx_tup_read and the table count in pg_stat_all_tables.idx_tup_fetch, while having no effect on pg_stat_all_indexes.idx_tup_fetch.

Table 14-207 PG_STAT_ALL_INDEXES columns

Column	Туре	Description
relid	OID	OID of the table for this index.
indexrelid	OID	OID of this index.
schemaname	Name	Name of the schema this index is in.
relname	Name	Name of the table for this index.
indexrelname	Name	Name of this index.
idx_scan	Bigint	Number of index scans initiated on this index.
idx_tup_read	Bigint	Number of index entries returned by scans on this index.
idx_tup_fetch	Bigint	Number of live table rows fetched by simple index scans using this index.

14.3.141 PG_STAT_ALL_TABLES

PG_STAT_ALL_TABLES displays statistics about accesses to tables in the current database, including TOAST tables.

Table 14-208 PG_STAT_ALL_TABLES columns

Column	Туре	Description
relid	OID	Table OID
schemaname	Name	Schema name of the table
relname	Name	Name of the table
seq_scan	Bigint	Number of sequential scans started on the table
seq_tup_read	Bigint	Number of rows that have live data fetched by sequential scans
idx_scan	Bigint	Number of index scans
idx_tup_fetch	Bigint	Number of rows that have live data fetched by index scans
n_tup_ins	Bigint	Number of rows inserted

Column	Туре	Description
n_tup_upd	Bigint	Number of rows updated
n_tup_del	Bigint	Number of rows deleted
n_tup_hot_up d	Bigint	Number of rows updated by HOT (no separate index update is required)
n_live_tup	Bigint	Estimated number of live rows
n_dead_tup	Bigint	Estimated number of dead rows
last_vacuum	Timestamp with time zone	Last time at which this table was manually vacuumed (excluding VACUUM FULL)
last_autovacu um	Timestamp with time zone	Last time at which this table was automatically vacuumed
last_analyze	Timestamp with time zone	Last time at which this table was analyzed
last_autoanal yze	Timestamp with time zone	Last time at which this table was automatically vacuumed
vacuum_coun t	Bigint	Number of vacuum operations (excluding VACUUM FULL)
autovacuum_ count	Bigint	Number of autovacuum operations
analyze_count	Bigint	Number of ANALYZE operations
autoanalyze_c ount	Bigint	Number of autoanalyze operations
last_data_cha nged	Timestamp with time zone	Last time at which this table was updated (by INSERT/UPDATE/DELETE or EXCHANGE/TRUNCATE/DROP partition). This column is recorded only on the local CN.

Example

Query the last data change time in the **table_test** table:

SELECT last_data_changed FROM PG_STAT_ALL_TABLES WHERE relname ='table_test'; last_data_changed

2024-03-27 10:28:16.277136+08

(1 row)

14.3.142 PG_STAT_BAD_BLOCK

PG_STAT_BAD_BLOCK displays statistics about page or CU verification failures after a node is started.

Table 14-209 PG_STAT_BAD_BLOCK columns

Column	Туре	Description
nodename	Text	Node name.
databaseid	Integer	Database OID.
tablespaceid	Integer	Tablespace OID.
relfilenode	Integer	File object ID.
forknum	Integer	File type.
error_count	Integer	Number of verification failures.
first_time	Timestamp with time zone	Time of the first occurrence.
last_time	Timestamp with time zone	Time of the latest occurrence.

14.3.143 PG_STAT_BGWRITER

PG_STAT_BGWRITER displays statistics about the background writer process's activity.

Table 14-210 PG_STAT_BGWRITER columns

Column	Туре	Description
checkpoints_ti med	Bigint	Number of scheduled checkpoints that have been performed.
checkpoints_r eq	Bigint	Number of requested checkpoints that have been performed.
checkpoint_wr ite_time	Double precision	Time spent on writing files to the disk during checkpoints, in milliseconds.
checkpoint_sy nc_time	Double precision	Time spent on synchronizing data to the disk during checkpoints, in milliseconds.
buffers_check point	Bigint	Number of buffers written during checkpoints.

Column	Туре	Description
buffers_clean	Bigint	Number of buffers written by the background writer.
maxwritten_cl ean	Bigint	Number of times the background writer stopped a cleaning scan because too many buffers were written.
buffers_backe nd	Bigint	Number of buffers written directly by a backend
buffers_backe nd_fsync	Bigint	Number of times that the backend has to execute fsync .
buffers_alloc	Bigint	Number of buffers allocated.
stats_reset	Timestamp with time zone	Time at which these statistics were reset.

14.3.144 PG_STAT_DATABASE

PG_STAT_DATABASE displays the status and statistics of each database on the current node.

Table 14-211 PG_STAT_DATABASE columns

Column	Туре	Description
datid	OID	Database OID
datname	Name	Database name
numbackends	Integer	Number of backends currently connected to this database on the current node. This is the only column in this view that reflects the current state value. All columns return the accumulated value since the last reset.
xact_commit	Bigint	Number of transactions in this database that have been committed on the current node
xact_rollback	Bigint	Number of transactions in this database that have been rolled back on the current node
blks_read	Bigint	Number of disk blocks read in this database on the current node
blks_hit	Bigint	Number of disk blocks found in the buffer cache on the current node, that is, the number of blocks hit in the cache. (This only includes hits in the GaussDB(DWS) buffer cache, not in the file system cache.)

Column	Туре	Description
tup_returned	Bigint	Number of rows returned by queries in this database on the current node
tup_fetched	Bigint	Number of rows fetched by queries in this database on the current node
tup_inserted	Bigint	Number of rows inserted in this database on the current node
tup_updated	Bigint	Number of rows updated in this database on the current node
tup_deleted	Bigint	Number of rows deleted from this database on the current node
conflicts	Bigint	Number of queries canceled due to database recovery conflicts on the current node (conflicts occurring only on the standby server). For details, see PG_STAT_DATABASE_CONFLICTS.
temp_files	Bigint	Number of temporary files created by this database on the current node. All temporary files are counted, regardless of why the temporary file was created (for example, sorting or hashing), and regardless of the log_temp_files setting.
temp_bytes	Bigint	Size of temporary files written to this database on the current node. All temporary files are counted, regardless of why the temporary file was created, and regardless of the log_temp_files setting.
deadlocks	Bigint	Number of deadlocks in this database on the current node
blk_read_time	Double precision	Time spent reading data file blocks by backends in this database on the current node, in milliseconds
blk_write_tim e	Double precision	Time spent writing into data file blocks by backends in this database on the current node, in milliseconds
stats_reset	Timestamp with time zone	Time when the database statistics are reset on the current node

14.3.145 PG_STAT_DATABASE_CONFLICTS

PG_STAT_DATABASE_CONFLICTS displays statistics about database conflicts.

Column Description Type datid OID Database OID. datname Name Database name. confl_tablesp **Bigint** Number of conflicting tablespaces. ace confl_lock Bigint Number of conflicting locks. confl_snapsho Bigint Number of conflicting snapshots. confl_bufferpi Bigint Number of conflicting buffers. confl deadloc Bigint Number of conflicting deadlocks.

Table 14-212 PG_STAT_DATABASE_CONFLICTS columns

14.3.146 PG_STAT_GET_MEM_MBYTES_RESERVED

PG_STAT_GET_MEM_MBYTES_RESERVED displays the current activity information of a thread stored in memory. You need to specify the thread ID (pid in **PG_STAT_ACTIVITY**) for query. If the thread ID is set to **0**, the current thread ID is used. For example:

SELECT pg_stat_get_mem_mbytes_reserved(0);

Table 14-213 PG_STAT_GET_MEM_MBYTES_RESERVED columns

Column	Description
ConnectInfo	Connection information.
ParctlManager	Concurrency management information.
GeneralParams	Basic parameter information.
GeneralParams RPDATA	Basic resource pool information.
ExceptionManager	Exception management information.
CollectInfo	Collection information.
GeneralInfo	Basic information.
ParctlState	Concurrency status information.
CPU INFO	CPU information.
ControlGroup	Cgroup information.
IOSTATE	I/O status information.

14.3.147 PG_STAT_USER_FUNCTIONS

PG_STAT_USER_FUNCTIONS displays user-defined function status information in the namespace. (The language of the function is non-internal language.)

Table 14-214 PG_STAT_USER_FUNCTIONS columns

Column	Туре	Description
funcid	OID	Function OID
schemaname	Name	Schema name
funcname	Name	Name of the function
calls	Bigint	Number of times this function has been called
total_time	Double precision	Total time spent in this function and all other functions called by it
self_time	Double precision	Total time spent in this function itself, excluding other functions called by it

14.3.148 PG_STAT_USER_INDEXES

PG_STAT_USER_INDEXES displays information about the index status of user-defined ordinary tables and TOAST tables.

Table 14-215 PG_STAT_USER_INDEXES columns

Column	Туре	Description
relid	OID	Table OID for the index
indexrelid	OID	OID of this index
schemaname	Name	Name of the schema this index is in
relname	Name	Name of the table for this index
indexrelname	Name	Name of this index
idx_scan	Bigint	Number of index scans
idx_tup_read	Bigint	Number of index entries returned by scans on this index
idx_tup_fetch	Bigint	Number of rows that have live data fetched by index scans

14.3.149 PG_STAT_USER_TABLES

PG_STAT_USER_TABLES displays status information about user-defined ordinary tables and TOAST tables in all namespaces.

Table 14-216 PG_STAT_USER_TABLES columns

Column	Туре	Description
relid	OID	Table OID
schemaname	Name	Schema name of the table
relname	Name	Name of a table
seq_scan	Bigint	Number of sequential scans started on the table
seq_tup_read	Bigint	Number of rows that have live data fetched by sequential scans
idx_scan	Bigint	Number of index scans
idx_tup_fetch	Bigint	Number of rows that have live data fetched by index scans
n_tup_ins	Bigint	Number of rows inserted
n_tup_upd	Bigint	Number of rows updated
n_tup_del	Bigint	Number of rows deleted
n_tup_hot_up d	Bigint	Number of rows updated by HOT (no separate index update is required)
n_live_tup	Bigint	Estimated number of live rows
n_dead_tup	Bigint	Estimated number of dead rows
last_vacuum	Timestamp with time zone	Last time at which this table was manually vacuumed (excluding VACUUM FULL)
last_autovacu um	Timestamp with time zone	Time of the last AUTOVACUUM
last_analyze	Timestamp with time zone	Last time at which this table was analyzed
last_autoanal yze	Timestamp with time zone	Time of the last AUTOANALYZE
vacuum_coun t	Bigint	Number of vacuum operations (excluding VACUUM FULL)

Column	Туре	Description
autovacuum_ count	Bigint	Number of autovacuum operations
analyze_count	Bigint	Number of analyze operations
autoanalyze_c ount	Bigint	Number of autoanalyze operations

14.3.150 PG_STAT_REPLICATION

PG_STAT_REPLICATION displays information about log synchronization status, such as the locations of the sender sending logs and the receiver receiving logs.

Table 14-217 PG_STAT_REPLICATION columns

Column	Туре	Description
pid	Bigint	PID of the thread.
usesysid	OID	User system ID.
usename	Name	Username.
application_n ame	Text	Program name.
client_addr	inet	Client address.
client_hostna me	Text	Client name.
client_port	Integer	Client port number.
backend_start	Timestamp with time zone	Program start time.
state	Text	Log replication state (catch-up or consistent streaming).
sender_sent_l ocation	Text	Location where the sender sends logs.
receiver_write _location	Text	Location where the receiver writes logs.
receiver_flush _location	Text	Location where the receiver flushes logs.
receiver_repla y_location	Text	Location where the receiver replays logs.

Column	Туре	Description
sync_priority	Integer	Priority of synchronous duplication (0 indicates asynchronization).
sync_state	Text	Synchronization state (asynchronous duplication, synchronous duplication, or potential synchronization).

14.3.151 PG_STAT_SYS_INDEXES

PG_STAT_SYS_INDEXES displays the index status information about all the system catalogs in the **pg_catalog** and **information_schema** schemas.

Table 14-218 PG_STAT_SYS_INDEXES columns

Column	Туре	Description
relid	OID	Table OID for the index.
indexrelid	OID	OID of this index.
schemaname	Name	Name of the schema this index is in.
relname	Name	Name of the table for this index.
indexrelname	Name	Name of this index.
idx_scan	Bigint	Number of index scans.
idx_tup_read	Bigint	Number of index entries returned by scans on this index.
idx_tup_fetch	Bigint	Number of rows that have live data fetched by index scans.

14.3.152 PG_STAT_SYS_TABLES

PG_STAT_SYS_TABLES displays the statistics about the system catalogs of all the namespaces in **pg_catalog** and **information_schema** schemas.

Table 14-219 PG_STAT_SYS_TABLES columns

Column	Туре	Description
relid	OID	Table OID
schemaname	Name	Schema name of the table
relname	Name	Name of a table

Column	Туре	Description
seq_scan	Bigint	Number of sequential scans started on the table
seq_tup_read	Bigint	Number of rows that have live data fetched by sequential scans
idx_scan	Bigint	Number of index scans
idx_tup_fetch	Bigint	Number of rows that have live data fetched by index scans
n_tup_ins	Bigint	Number of rows inserted
n_tup_upd	Bigint	Number of rows updated
n_tup_del	Bigint	Number of rows deleted
n_tup_hot_up d	Bigint	Number of rows updated by HOT (no separate index update is required)
n_live_tup	Bigint	Estimated number of live rows
n_dead_tup	Bigint	Estimated number of dead rows
last_vacuum	Timestamp with time zone	Last time at which this table was manually vacuumed (excluding VACUUM FULL)
last_autovacu um	Timestamp with time zone	Last time at which this table was automatically vacuumed
last_analyze	Timestamp with time zone	Last time at which this table was analyzed
last_autoanal yze	Timestamp with time zone	Last time at which this table was automatically analyzed
vacuum_coun t	Bigint	Number of vacuum operations (excluding VACUUM FULL)
autovacuum_ count	Bigint	Number of autovacuum operations
analyze_count	Bigint	Number of analyze operations
autoanalyze_c ount	Bigint	Number of autoanalyze operations

14.3.153 PG_STAT_XACT_ALL_TABLES

PG_STAT_XACT_ALL_TABLES displays the transaction status information about all ordinary tables and TOAST tables in the namespaces.

Column Description Type Table OID relid OID schemaname Name Schema name of the table relname Name Name of a table seq_scan Bigint Number of sequential scans started on the table seq_tup_read Bigint Number of live rows fetched by sequential scans idx_scan Bigint Number of index scans started on the table idx_tup_fetch **Bigint** Number of live rows fetched by index scans n_tup_ins **Bigint** Number of rows inserted n_tup_upd **Bigint** Number of rows updated Number of rows deleted n_tup_del **Bigint**

Number of rows with HOT updates (no

separate index update is required).

Table 14-220 PG_STAT_XACT_ALL_TABLES columns

14.3.154 PG_STAT_XACT_SYS_TABLES

n_tup_hot_up

d

PG_STAT_XACT_SYS_TABLES displays the transaction status information of the system catalog in the namespace.

Table 14-221 PG_STAT_XACT_SYS_TABLES columns

Bigint

Column	Туре	Description
relid	OID	Table OID
schemaname	Name	Schema name of the table
relname	Name	Table name
seq_scan	Bigint	Number of sequential scans started on the table
seq_tup_read	Bigint	Number of live rows fetched by sequential scans
idx_scan	Bigint	Number of index scans started on the table
idx_tup_fetch	Bigint	Number of live rows fetched by index scans
n_tup_ins	Bigint	Number of rows inserted
n_tup_upd	Bigint	Number of rows updated

Column	Туре	Description
n_tup_del	Bigint	Number of rows deleted
n_tup_hot_up d	Bigint	Number of rows with HOT updates (no separate index update is required).

14.3.155 PG_STAT_XACT_USER_FUNCTIONS

PG_STAT_XACT_USER_FUNCTIONS displays statistics about function execution.

Table 14-222 PG_STAT_XACT_USER_FUNCTIONS columns

Column	Туре	Description
funcid	OID	Function OID
schemaname	Name	Schema name
funcname	Name	Name of the function
calls	Bigint	Number of times this function has been called
total_time	Double precision	Total time spent in this function and all other functions called by it
self_time	Double precision	Total time spent in this function itself, excluding other functions called by it

14.3.156 PG_STAT_XACT_USER_TABLES

PG_STAT_XACT_USER_TABLES displays the transaction status information of the user table in the namespace.

Table 14-223 PG_STAT_XACT_USER_TABLES columns

Column	Туре	Description
relid	OID	Table OID
schemaname	Name	Schema name of the table
relname	Name	Name of a table
seq_scan	Bigint	Number of sequential scans started on the table
seq_tup_read	Bigint	Number of live rows fetched by sequential scans
idx_scan	Bigint	Number of index scans started on the table

Column	Туре	Description
idx_tup_fetch	Bigint	Number of live rows fetched by index scans
n_tup_ins	Bigint	Number of rows inserted
n_tup_upd	Bigint	Number of rows updated
n_tup_del	Bigint	Number of rows deleted
n_tup_hot_up d	Bigint	Number of rows with HOT updates (no separate index update is required).

14.3.157 PG_STATIO_ALL_INDEXES

PG_STATIO_ALL_INDEXES displays I/O statistics of all indexes in the current database.

Table 14-224 PG_STATIO_ALL_INDEXES columns

Column	Туре	Description
relid	OID	OID of the index table
indexrelid	OID	OID of this index
schemaname	Name	Name of the schema this index is in
relname	Name	Name of the table for this index
indexrelname	Name	Name of this index
idx_blks_read	Bigint	Number of disk blocks read from the index
idx_blks_hit	Bigint	Number of buffer hits in this index

14.3.158 PG_STATIO_ALL_SEQUENCES

PG_STATIO_ALL_SEQUENCES displays the sequence information in the current database and the I/O statistics of a specified sequence.

Table 14-225 PG_STATIO_ALL_SEQUENCES columns

Column	Туре	Description
relid	OID	OID of this sequence
schemaname	Name	Name of the schema this sequence is in
relname	Name	Name of this sequence
blks_read	Bigint	Number of disk blocks read from the sequence

Column	Туре	Description
blks_hit	Bigint	Number of buffer hits in this sequence

14.3.159 PG_STATIO_ALL_TABLES

PG_STATIO_ALL_TABLES displays I/O statistics about all tables (including TOAST tables) in the current database.

Table 14-226 PG_STATIO_ALL_TABLES columns

Column	Туре	Description
relid	OID	Table OID
schemaname	Name	Schema name of the table
relname	Name	Name of a table
heap_blks_rea d	Bigint	Number of disks read from this table
heap_blks_hit	Bigint	Number of buffer hits in this table
idx_blks_read	Bigint	Number of disk blocks read from the index in this table
idx_blks_hit	Bigint	Number of buffer hits in all indexes on this table
toast_blks_rea d	Bigint	Number of disk blocks read from the TOAST table (if any) in this table
toast_blks_hit	Bigint	Number of buffer hits in the TOAST table (if any) in this table
tidx_blks_read	Bigint	Number of disk blocks read from the TOAST table index (if any) in this table
tidx_blks_hit	Bigint	Number of buffer hits in the TOAST table index (if any) in this table

14.3.160 PG_STATIO_SYS_INDEXES

PG_STATIO_SYS_INDEXES displays the I/O status information about all system catalog indexes in the namespace.

Column Type Description relid OID Table OID for the index. OID of this index. indexrelid OID Name of the schema of the index. schemaname Name Name of the table for this index. relname Name indexrelname Name of this index. Name idx_blks_read **Bigint** Number of disk blocks read from the index. idx blks hit Number of buffer hits in this index. **Bigint**

Table 14-227 PG_STATIO_SYS_INDEXES columns

14.3.161 PG_STATIO_SYS_SEQUENCES

PG_STATIO_SYS_SEQUENCES displays the I/O status information about all the system sequences in the namespace.

Table 14-228 PG_STATIO_SYS_SEQUENCES columns

Column	Туре	Description
relid	OID	OID of this sequence
schemaname	Name	Name of the schema this sequence is in
relname	Name	Name of this sequence
blks_read	Bigint	Number of disk blocks read from the sequence
blks_hit	Bigint	Number of buffer hits in this sequence

14.3.162 PG_STATIO_SYS_TABLES

PG_STATIO_SYS_TABLES displays the I/O status information about all the system catalogs in the namespace.

Table 14-229 PG_STATIO_SYS_TABLES columns

Column	Туре	Description
relid	OID	Table OID
schemaname	Name	Schema name of the table
relname	Name	Name of a table

Column	Туре	Description
heap_blks_read	Bigint	Number of disk blocks read from this table
heap_blks_hit	Bigint	Number of buffer hits in this table
idx_blks_read	Bigint	Number of disk blocks read from the index in this table
idx_blks_hit	Bigint	Number of buffer hits in all indexes on this table
toast_blks_read	Bigint	Number of disk blocks read from the TOAST table (if any) in this table
toast_blks_hit	Bigint	Number of buffer hits in the TOAST table (if any) in this table
tidx_blks_read	Bigint	Number of disk blocks read from the TOAST table index (if any) in this table
tidx_blks_hit	Bigint	Number of buffer hits in the TOAST table index (if any) in this table

14.3.163 PG_STATIO_USER_INDEXES

PG_STATIO_USER_INDEXES displays the I/O status information about all the user relationship table indexes in the namespace.

Table 14-230 PG_STATIO_USER_INDEXES columns

Column	Туре	Description
relid	OID	OID of the table for this index
indexrelid	OID	OID of this index
schemaname	Name	Name of the schema this index is in
relname	Name	Name of the table for this index
indexrelname	Name	Name of this index
idx_blks_read	Bigint	Number of disk blocks read from the index
idx_blks_hit	Bigint	Number of buffer hits in this index

14.3.164 PG_STATIO_USER_SEQUENCES

PG_STATIO_USER_SEQUENCES displays the I/O status information about all the user relation table sequences in the namespace.

 Table 14-231 PG_STATIO_USER_SEQUENCES columns

Column	Туре	Description
relid	OID	OID of this sequence
schemaname	Name	Name of the schema this sequence is in
relname	Name	Name of this sequence
blks_read	Bigint	Number of disk blocks read from the sequence
blks_hit	Bigint	Cache hits in the sequence

14.3.165 PG_STATIO_USER_TABLES

PG_STATIO_USER_TABLES displays the I/O status information about all the user relation tables in the namespace.

Table 14-232 PG_STATIO_USER_TABLES columns

Column	Туре	Description
relid	OID	Table OID
schemaname	Name	Schema name of the table
relname	Name	Name of a table
heap_blks_read	Bigint	Number of disk blocks read from this table
heap_blks_hit	Bigint	Number of buffer hits in this table
idx_blks_read	Bigint	Number of disk blocks read from the index in this table
idx_blks_hit	Bigint	Number of buffer hits in all indexes on this table
toast_blks_read	Bigint	Number of disk blocks read from the TOAST table (if any) in this table
toast_blks_hit	Bigint	Number of buffer hits in the TOAST table (if any) in this table
tidx_blks_read	Bigint	Number of disk blocks read from the TOAST table index (if any) in this table
tidx_blks_hit	Bigint	Number of buffer hits in the TOAST table index (if any) in this table

14.3.166 PG_THREAD_WAIT_STATUS

PG_THREAD_WAIT_STATUS allows you to test the block waiting status about the backend thread and auxiliary thread of the current instance.

Table 14-233 PG_THREAD_WAIT_STATUS columns

Column	Туре	Description
node_name	Text	Current node name
db_name	Text	Database name
thread_name	Text	Thread name
query_id	Bigint	Query ID. It is equivalent to debug_query_id.
tid	Bigint	Thread ID of the current thread
lwtid	Integer	Lightweight thread ID of the current thread
ptid	Integer	Parent thread of the streaming thread
tlevel	Integer	Level of the streaming thread
smpid	Integer	Concurrent thread ID
wait_status	Text	Waiting status of the current thread. For details about the waiting status, see Table 14-234 .
wait_event	Text	If wait_status is acquire lock, acquire lwlock, or wait io, this column describes the lock, lightweight lock, and I/O information, respectively. If wait_status is not any of the three values, this column is empty.

The waiting statuses in the wait_status column are as follows:

Table 14-234 Waiting status list

Value	Description
none	Waiting for no event
acquire lock	Waiting for locking until the locking succeeds or times out
acquire lwlock	Waiting for a lightweight lock
wait io	Waiting for I/O completion
wait cmd	Waiting for network communication packet read to complete
wait pooler get conn	Waiting for pooler to obtain the connection

Value	Description
wait pooler abort conn	Waiting for pooler to terminate the connection
wait pooler clean conn	Waiting for pooler to clear connections
pooler create conn: [nodename], total N	Waiting for the pooler to set up a connection. The connection is being established with the node specified by <i>nodename</i> , and there are <i>N</i> connections waiting to be set up.
get conn	Obtaining the connection to other nodes
set cmd: [nodename]	Waiting for running the SET, RESET, TRANSACTION BLOCK LEVEL PARA SET, or SESSION LEVEL PARA SET statement on the connection. The statement is being executed on the node specified by nodename.
cancel query	Canceling the SQL statement that is being executed through the connection
stop query	Stopping the query that is being executed through the connection
wait node: [nodename](plevel), total N, [phase]	Waiting for receiving the data from a connected node. The thread is waiting for the data from the plevel thread of the node specified by <i>nodename</i> . The data of <i>N</i> connections is waiting to be returned. If <i>phase</i> is included, the possible phases are as follows:
	• begin : The transaction is being started.
	• commit : The transaction is being committed.
	rollback: The transaction is being rolled back.
wait transaction sync: xid	Waiting for synchronizing the transaction specified by <i>xid</i>
wait wal sync	Waiting for the completion of wal log of synchronization from the specified LSN to the standby instance
wait data sync	Waiting for the completion of data page synchronization to the standby instance
wait data sync queue	Waiting for putting the data pages that are in the row storage or the CU in the column storage into the synchronization queue

Value	Description
flush data: [nodename](plevel), [phase]	Waiting for sending data to the plevel thread of the node specified by <i>nodename</i> . If <i>phase</i> is included, the possible phase is wait quota , indicating that the current communication flow is waiting for the quota value.
stream get conn: [nodename], total N	Waiting for connecting to the consumer object of the node specified by <i>nodename</i> when the stream flow is initialized. There are <i>N</i> consumers waiting to be connected.
wait producer ready: [nodename] (plevel), total N	Waiting for each producer to be ready when the stream flow is initialized. The thread is waiting for the procedure of the plevel thread on the <i>nodename</i> node to be ready. There are <i>N</i> producers waiting to be ready.
synchronize quit	Waiting for the threads in the stream thread group to quit when the stream plan ends
nodegroup destroy	Waiting for destroying the stream node group when the stream plan ends
wait active statement	Waiting for job execution under resource and load control.
wait global queue	Waiting for job execution. The job is queuing in the global queue.
wait respool queue	Waiting for job execution. The job is queuing in the resource pool.
wait ccn queue	Waiting for job execution. The job is queuing on the central coordinator node (CCN).
gtm connect	Waiting for connecting to GTM.
gtm get gxid	Wait for obtaining xids from GTM.
gtm get snapshot	Wait for obtaining transaction snapshots from GTM.
gtm begin trans	Waiting for GTM to start a transaction.
gtm commit trans	Waiting for GTM to commit a transaction.
gtm rollback trans	Waiting for GTM to roll back a transaction.
gtm create sequence	Waiting for GTM to create a sequence.
gtm alter sequence	Waiting for GTM to modify a sequence.
gtm get sequence val	Waiting for obtaining the next value of a sequence from GTM.

Value	Description
gtm set sequence val	Waiting for GTM to set a sequence value.
gtm drop sequence	Waiting for GTM to delete a sequence.
gtm rename sequence	Waiting for GTM to rename a sequence.
analyze: [relname], [phase]	The thread is doing ANALYZE to the <i>relname</i> table. If <i>phase</i> is included, the possible phase is autovacuum , indicating that the database automatically enables the AutoVacuum thread to execute ANALYZE .
vacuum: [relname], [phase]	The thread is doing VACUUM to the <i>relname</i> table. If <i>phase</i> is included, the possible phase is autovacuum , indicating that the database automatically enables the AutoVacuum thread to execute VACUUM .
vacuum full: [relname]	The thread is doing VACUUM FULL to the <i>relname</i> table.
create index	An index is being created.
HashJoin - [build hash write file]	The HashJoin operator is being executed. In this phase, you need to pay attention to the execution time-consuming.
	• build hash : The HashJoin operator is creating a hash table.
	write file: The HashJoin operator is writing data to disks.
HashAgg - [build hash write file]	The HashAgg operator is being executed. In this phase, you need to pay attention to the execution time-consuming.
	 build hash: The HashAgg operator is creating a hash table.
	write file: The HashAgg operator is writing data to disks.
HashSetop - [build hash write file]	The HashSetop operator is being executed. In this phase, you need to pay attention to the execution time-consuming.
	• build hash : The HashSetop operator is creating a hash table.
	write file: The HashSetop operator is writing data to disks.
Sort Sort - write file	The Sort operator is being executed. write file indicates that the Sort operator is writing data to disks.

Value	Description
Material Material - write file	The Material operator is being executed. write file indicates that the Material operator is writing data to disks.
wait sync consumer next step	The consumer (receive end) synchronously waits for the next iteration.
wait sync producer next step	The producer (transmit end) synchronously waits for the next iteration.
wait agent release	The current agent is being released (supported by 8.1.2 and later versions).
wait stream task	The stream thread is waiting for being reused (supported by 8.1.2 and later versions).

If wait_status is acquire lwlock, acquire lock, or wait io, there is an event performing I/O operations or waiting for obtaining the corresponding lightweight lock or transaction lock.

The following table describes the corresponding wait events when **wait_status** is **acquire lwlock**. (If **wait_event** is **extension**, the lightweight lock is dynamically allocated and is not monitored.)

Table 14-235 List of wait events corresponding to lightweight locks

wait_event	Description
ShmemIndexLock	Used to protect the primary index table, a hash table, in shared memory
OidGenLock	Used to prevent different threads from generating the same OID
XidGenLock	Used to prevent two transactions from obtaining the same XID
ProcArrayLock	Used to prevent concurrent access to or concurrent modification on the ProcArray shared array
SInvalReadLock	Used to prevent concurrent execution with invalid message deletion
SInvalWriteLock	Used to prevent concurrent execution with invalid message write and deletion
WALInsertLock	Used to prevent concurrent execution with WAL insertion
WALWriteLock	Used to prevent concurrent write from a WAL buffer to a disk

wait_event	Description
ControlFileLock	Used to prevent concurrent read/write or concurrent write/write on the pg_control file
CheckpointLock	Used to prevent multi-checkpoint concurrent execution
CLogControlLock	Used to prevent concurrent access to or concurrent modification on the Clog control data structure
MultiXactGenLock	Used to allocate a unique MultiXact ID in serial mode
MultiXactOffsetControl- Lock	Used to prevent concurrent read/write or concurrent write/write on pg_multixact/offset
MultiXactMemberControl- Lock	Used to prevent concurrent read/write or concurrent write/write on pg_multixact/members
RelCacheInitLock	Used to add a lock before any operations are performed on the init file when messages are invalid
CheckpointerCommLock	Used to send file flush requests to a checkpointer. The request structure needs to be inserted to a request queue in serial mode.
TwoPhaseStateLock	Used to prevent concurrent access to or modification on two-phase information sharing arrays
TablespaceCreateLock	Used to check whether a tablespace already exists
BtreeVacuumLock	Used to prevent VACUUM from clearing pages that are being used by B-tree indexes
AutovacuumLock	Used to access the autovacuum worker array in serial mode
AutovacuumScheduleLock	Used to distribute tables requiring VACUUM in serial mode
SyncScanLock	Used to determine the start position of a relfilenode during heap scanning
NodeTableLock	Used to protect a shared structure that stores CN and DN information
PoolerLock	Used to prevent two threads from simultaneously obtaining the same connection from a connection pool
RelationMappingLock	Used to wait for the mapping file between system catalogs and storage locations to be updated
AsyncCtlLock	Used to prevent concurrent access to or concurrent modification on the sharing notification status

wait_event	Description
AsyncQueueLock	Used to prevent concurrent access to or concurrent modification on the sharing notification queue
SerializableXactHashLock	Used to prevent concurrent read/write or concurrent write/write on a sharing structure for serializable transactions
SerializableFinishedList- Lock	Used to prevent concurrent read/write or concurrent write/write on a shared linked list for completed serial transactions
SerializablePredicateLock- ListLock	Used to protect a linked list of serializable transactions that have locks
OldSerXidLock	Used to protect a structure that records serializable transactions that have conflicts
FileStatLock	Used to protect a data structure that stores statistics file information
SyncRepLock	Used to protect Xlog synchronization information during primary-standby replication
DataSyncRepLock	Used to protect data page synchronization information during primary-standby replication
CStoreColspaceCacheLock	Used to add a lock when CU space is allocated for a column-store table
CStoreCUCacheSweep- Lock	Used to add a lock when CU caches used by a column-store table are cyclically washed out
MetaCacheSweepLock	Used to add a lock when metadata is cyclically washed out
DfsConnectorCacheLock	Used to protect a global hash table where HDFS connection handles are cached
dummyServerInfoCache- Lock	Used to protect a global hash table where the information about computing Node Group connections is cached
ExtensionConnectorLi- bLock	Used to add a lock when a specific dynamic library is loaded or uninstalled in ODBC connection initialization scenarios
SearchServerLibLock	Used to add a lock on the file read operation when a specific dynamic library is initially loaded in GPU-accelerated scenarios
DfsUserLoginLock	Used to protect a global linked table where HDFS user information is stored
DfsSpaceCacheLock	Used to ensure that the IDs of files to be imported to an HDFS table increase monotonically

wait_event	Description
LsnXlogChkFileLock	Used to serially update the Xlog flush points for primary and standby servers recorded in a specific structure
GTMHostInfoLock	Used to prevent concurrent access to or concurrent modification on GTM host information
ReplicationSlotAllocation- Lock	Used to add a lock when a primary server allocates stream replication slots during primary-standby replication
ReplicationSlotControl- Lock	Used to prevent concurrent update of replication slot status during primary-standby replication
ResourcePoolHashLock	Used to prevent concurrent access to or concurrent modification on a resource pool table, a hash table
WorkloadStatHashLock	Used to prevent concurrent access to or concurrent modification on a hash table that contains SQL requests from the CN side
WorkloadIoStatHashLock	Used to prevent concurrent access to or concurrent modification on a hash table that contains the I/O information of the current DN
WorkloadCGroupHash- Lock	Used to prevent concurrent access to or concurrent modification on a hash table that contains Cgroup information
OBSGetPathLock	Used to prevent concurrent read/write or concurrent write/write on an OBS path
WorkloadUserInfoLock	Used to prevent concurrent access to or concurrent modification on a hash table that contains user information about load management
WorkloadRecordLock	Used to prevent concurrent access to or concurrent modification on a hash table that contains requests received by CNs during adaptive memory management
WorkloadIOUtilLock	Used to protect a structure that records iostat and CPU load information
WorkloadNodeGroupLock	Used to prevent concurrent access to or concurrent modification on a hash table that contains Node Group information in memory
JobShmemLock	Used to protect global variables in the shared memory that is periodically read during a scheduled task where MPP is compatible with Oracle
OBSRuntimeLock	Used to obtain environment variables, for example, <i>GAUSSHOME</i> .

wait_event	Description	
LLVMDumpIRLock	Used to export the assembly language for dynamically generating functions	
LLVMParseIRLock	Used to compile and parse a finished IR function from the IR file at the start position of a query	
RPNumberLock	Used by a DN on a computing Node Group to count the number of threads for a task where plans are being executed	
ClusterRPLock	Used to control concurrent access on cluster load data maintained in a CCN of the cluster	
CriticalCacheBuildLock	Used to load caches from a shared or local cache initialization file	
WaitCountHashLock	Used to protect a shared structure in user statement counting scenarios	
BufMappingLock	Used to protect operations on a table mapped to shared buffer	
LockMgrLock	It is used to protect a common lock structure.	
PredicateLockMgrLock	Used to protect a lock structure that has serializable transactions	
OperatorRealTLock	Used to prevent concurrent access to or concurrent modification on a global structure that contains real-time data at the operator level	
OperatorHistLock	Used to prevent concurrent access to or concurrent modification on a global structure that contains historical data at the operator level	
SessionRealTLock	Used to prevent concurrent access to or concurrent modification on a global structure that contains real-time data at the query level	
SessionHistLock	Used to prevent concurrent access to or concurrent modification on a global structure that contains historical data at the query level	
CacheSlotMappingLock	Used to protect global CU cache information	
BarrierLock	Used to ensure that only one thread is creating a barrier at a time	

The following table describes the corresponding wait events when **wait_status** is **wait io**.

Table 14-236 List of wait events corresponding to I/Os

wait_event	Description	
BufFileRead	Reads data from a temporary file to a specified buffer.	
BufFileWrite	Writes the content of a specified buffer to a temporary file.	
ControlFileRead	Reads the pg_control file, mainly during database startup, checkpoint execution, and primary/standby verification.	
ControlFileSync	Flushes the pg_control file to a disk, mainly during database initialization.	
ControlFileSyncUpdate	Flushes the pg_control file to a disk, mainly during database startup, checkpoint execution, and primary/standby verification.	
ControlFileWrite	Writes to the pg_control file, mainly during database initialization.	
ControlFileWriteUpdate	Updates the pg_control file, mainly during database startup, checkpoint execution, and primary/standby verification.	
CopyFileRead	Reads a file during file copying.	
CopyFileWrite	Writes a file during file copying.	
DataFileExtend	Writes a file during file extension.	
DataFileFlush	Flushes a table data file to a disk.	
DataFileImmediateSync	Flushes a table data file to a disk immediately.	
DataFilePrefetch	Reads a table data file asynchronously.	
DataFileRead	Reads a table data file synchronously.	
DataFileSync	Flushes table data file modifications to a disk.	
DataFileTruncate	Truncates a table data file.	
DataFileWrite	Writes a table data file.	
LockFileAddToDataDir- Read	Reads the postmaster.pid file.	
LockFileAddToDataDir- Sync	Flushes the postmaster.pid file to a disk.	
LockFileAddToDataDir- Write	Writes the PID information into the postmaster.pid file.	
LockFileCreateRead	Read the LockFile file %s.lock .	
LockFileCreateSync	Flushes the LockFile file %s.lock to a disk.	

wait_event	Description	
LockFileCreateWRITE	Writes the PID information into the LockFile file %s.lock .	
RelationMapRead	Reads the mapping file between system catalogs and storage locations.	
RelationMapSync	Flushes the mapping file between system catalogs and storage locations to a disk.	
RelationMapWrite	Writes the mapping file between system catalogs and storage locations.	
ReplicationSlotRead	Reads a stream replication slot file during a restart.	
ReplicationSlotRestore- Sync	Flushes a stream replication slot file to a disk during a restart.	
ReplicationSlotSync	Flushes a temporary stream replication slot file to a disk during checkpoint execution.	
ReplicationSlotWrite	Writes a temporary stream replication slot file during checkpoint execution.	
SLRUFlushSync	Flushes the pg_clog , pg_subtrans , and pg_multixact files to a disk, mainly during checkpoint execution and database shutdown.	
SLRURead	Reads the pg_clog , pg_subtrans , and pg_multixact files.	
SLRUSync	Writes dirty pages into the pg_clog, pg_subtrans, and pg_multixact files, and flushes the files to a disk, mainly during checkpoint execution and database shutdown.	
SLRUWrite	Writes the pg_clog , pg_subtrans , and pg_multixact files.	
TimelineHistoryRead	Reads the timeline history file during database startup.	
TimelineHistorySync	Flushes the timeline history file to a disk during database startup.	
TimelineHistoryWrite	Writes to the timeline history file during database startup.	
TwophaseFileRead	Reads the pg_twophase file, mainly during two-phase transaction submission and restoration.	
TwophaseFileSync	Flushes the pg_twophase file to a disk, mainly during two-phase transaction submission and restoration.	
TwophaseFileWrite	Writes the pg_twophase file, mainly during two-phase transaction submission and restoration.	

wait_event	Description	
WALBootstrapSync	Flushes an initialized WAL file to a disk during database initialization.	
WALBootstrapWrite	Writes an initialized WAL file during database initialization.	
WALCopyRead	Read operation generated when an existing WAL file is read for replication after archiving and restoration.	
WALCopySync	Flushes a replicated WAL file to a disk after archiving and restoration.	
WALCopyWrite	Write operation generated when an existing WAL file is read for replication after archiving and restoration.	
WALInitSync	Flushes a newly initialized WAL file to a disk during log reclaiming or writing.	
WALInitWrite	Initializes a newly created WAL file to 0 during log reclaiming or writing.	
WALRead	Reads data from Xlogs during redo operations on two-phase files.	
WALSyncMethodAssign	Flushes all open WAL files to a disk.	
WALWrite	Writes a WAL file.	

The following table describes the corresponding wait events when **wait_status** is **acquire lock**.

Table 14-237 List of wait events corresponding to transaction locks

wait_event	Description	
relation	Adds a lock to a table.	
extend	Adds a lock to a table being scaled out.	
partition	Adds a lock to a partitioned table.	
partition_seq	Adds a lock to a partition of a partitioned table.	
page	Adds a lock to a table page.	
tuple	Adds a lock to a tuple on a page.	
transactionid	Adds a lock to a transaction ID.	
virtualxid	Adds a lock to a virtual transaction ID.	
object	Adds a lock to an object.	

wait_event	Description	
cstore_freespace	Adds a lock to idle column-store space.	
userlock	Adds a lock to a user.	
advisory	Adds an advisory lock.	

14.3.167 PG_TABLES

PG_TABLES displays access to each table in the database.

Table 14-238 PG_TABLES columns

Column	Туре	Reference	Description
schemana me	Name	PG_NAMESPACE.nspname	Name of the schema that contains the table
tablenam e	Name	PG_CLASS.relname	Name of the table
tableown er	Name	pg_get_userbyid(PG_CLAS S.relowner)	Owner of the table
tablespac e	Name	PG_TABLESPACE.spcname	Tablespace that contains the table. The default value is null
hasindexe s	boolean	PG_CLASS.relhasindex	Whether the table has (or recently had) an index. If it does, its value is true . Otherwise, its value is false .
hasrules	boolean	PG_CLASS.relhasrules	Whether the table has rules. If it does, its value is true . Otherwise, its value is false .
hastrigger s	boolean	PG_CLASS.RELHASTRIGGE RS	Whether the table has triggers. If it does, its value is true . Otherwise, its value is false .
tablecreat or	Name	pg_get_userbyid(PG_OBJECT.creator)	Table creator. If the creator has been deleted, no value is returned.
created	Timestam p with time zone	PG_OBJECT.ctime	Time when the table was created.

Column	Туре	Reference	Description
last_ddl_ti me	Timestam p with time zone	PG_OBJECT.mtime	Last time when the cluster was modified.

Example

Query all tables in a specified schema.

```
SELECT tablename FROM PG_TABLES WHERE schemaname = 'myschema';
tablename
------
inventory
product
sales_info
test1
mytable
product_info
customer_info
newproducts
customer_t1
(9 rows)
```

14.3.168 PG_TDE_INFO

PG_TDE_INFO displays the encryption information about the current cluster.

Table 14-239 PG_TDE_INFO columns

Column	Туре	Description
is_encrypt	Text	 Whether the cluster is an encryption cluster f: Non-encryption cluster t: Encryption cluster
g_tde_algo	Text	Encryption algorithm • AES-CTR-128
remain	Text	Reserved columns

Examples

Check whether the current cluster is encrypted, and check the encryption algorithm (if any) used by the current cluster.

```
SELECT * FROM PG_TDE_INFO;
is_encrypt | g_tde_algo | remain
-------
f | AES-CTR-128 | remain
(1 row)
```

14.3.169 PG_TIMEZONE_ABBREVS

PG_TIMEZONE_ABBREVS displays all time zone abbreviations that can be recognized by the input routines.

Table 14-240 PG_TIMEZONE_ABBREVS columns

Column	Туре	Description
abbrev	Text	Time zone abbreviation
utc_offset	interval	Offset from UTC
is_dst	boolean	Whether the abbreviation indicates a daylight saving time (DST) zone. If it does, its value is true . Otherwise, its value is false .

14.3.170 PG_TIMEZONE_NAMES

PG_TIMEZONE_NAMES displays all time zone names that can be recognized by **SET TIMEZONE**, along with their associated abbreviations, UTC offsets, and daylight saving time statuses.

Table 14-241 PG_TIMEZONE_NAMES columns

Column	Туре	Description
Name	Text	Name of the time zone
abbrev	Text	Time zone name abbreviation
utc_offset	interval	Offset from UTC
is_dst	boolean	Whether DST is used. If it is, its value is true . Otherwise, its value is false .

14.3.171 PG_TOTAL_MEMORY_DETAIL

PG_TOTAL_MEMORY_DETAIL displays the memory usage of a certain node in the database.

Table 14-242 PG_TOTAL_MEMORY_DETAIL columns

Column	Туре	Description
nodename	Text	Node name

Column	Туре	Description
memorytype	Text	It can be set to any of the following values:
		• max_process_memory: memory used by a GaussDB(DWS) cluster instance
		 process_used_memory: memory used by a GaussDB(DWS) process
		max_dynamic_memory: maximum dynamic memory
		dynamic_used_memory: used dynamic memory
		dynamic_peak_memory: dynamic peak value of the memory
		• dynamic_used_shrctx : maximum dynamic shared memory context
		dynamic_peak_shrctx: dynamic peak value of the shared memory context
		max_shared_memory: maximum shared memory
		 shared_used_memory: used shared memory
		max_cstore_memory: maximum memory allowed for column store
		 cstore_used_memory: memory used for column store
		max_sctpcomm_memory: maximum memory allowed for the communication library
		 sctpcomm_used_memory: memory used for the communication library
		sctpcomm_peak_memory: memory peak of the communication library
		max_topsql_memory: maximum memory that can be used by top SQL to record historical job monitoring information
		 topsql_used_memory: memory used by top SQL to record historical job monitoring information
		topsql_peak_memory: memory peak of top SQL to record historical job monitoring information
		other_used_memory: other used memory
		gpu_max_dynamic_memory: maximum GPU memory

Column	Туре	Description
		gpu_dynamic_used_memory: sum of the available GPU memory and temporary GPU memory
		 gpu_dynamic_peak_memory: maximum memory used for GPU
		 pooler_conn_memory: memory used for pooler connections
		 pooler_freeconn_memory: memory used for idle pooler connections
		 storage_compress_memory: memory used for column-store compression and decompression
		 udf_reserved_memory: memory reserved for the UDF Worker process
		mmap_used_memory: memory used for mmap
memorymbyte s	Integer	Size of the used memory (MB)

14.3.172 PG_TOTAL_SCHEMA_INFO

PG_TOTAL_SCHEMA_INFO displays the storage usage of all schemas in each database. This view is valid only if use_workload_manager is set to **on**.

Column	Туре	Description
schemaid	OID	Schema OID
schemanam e	Text	Schema name
databaseid	OID	Database OID
databasena me	Name	Database name
usedspace	Bigint	Size of the permanent table storage space used by the schema, in bytes.
permspace	Bigint	Upper limit of the permanent table storage space of the schema, in bytes.

14.3.173 PG_TOTAL_USER_RESOURCE_INFO

PG_TOTAL_USER_RESOURCE_INFO displays the resource usage of all users. Only administrators can query this view. This view is valid only if use_workload_manager is set to **on**.

Table 14-243 PG_TOTAL_USER_RESOURCE_INFO columns

Column	Туре	Description
username	Name	Username
used_memory	Integer	Memory used by a user, in MB.
		On a DN, it indicates the memory used by users on the current DN.
		On a CN, it indicates the total memory used by users on all DNs.
total_memory	Integer	Memory used by the resource pool, in MB. 0 indicates that the maximum available memory is not limited and depends on the maximum available memory of the database (max_dynamic_memory). The calculation formula is as follows:
		total_memory = max_dynamic_memory * parent_percent * user_percent
		On a CN, it indicates the total maximum available memory on all DNs.
used_cpu	Double precision	Number of CPU cores in use. Only the CPU usage of complex jobs in the non-default resource pool is collected, and the value is the CPU usage of the related cgroup.
total_cpu	Integer	Total number of CPU cores of the Cgroup associated with a user on the node
used_space	Bigint	Used permanent table storage space (unit: KB)
total_space	Bigint	Available storage space (unit: KB)1 indicates that the storage space is not limited.
used_temp_sp ace	Bigint	Used temporary table storage space (unit: KB)
total_temp_sp ace	Bigint	Available temporary table storage space (unit: KB)1 indicates that the storage space is not limited.
used_spill_spa ce	Bigint	Size of the used operator flushing space, in KB

Column	Туре	Description
total_spill_spa ce	Bigint	Size of the available operator flushing space, in KB. The value -1 indicates that the operator flushing space is not limited.
read_kbytes	Bigint	On a CN, it indicates the total number of bytes logically read by a user on all DNs in the last 5 seconds, in KB.
		On a DN, it indicates the total number of bytes logically read by a user from the instance startup time to the current time, in KB.
write_kbytes	Bigint	On a CN, it indicates the total number of bytes logically written by a user on all DNs in the last 5 seconds, in KB.
		On a DN, it indicates the total number of bytes logically written by a user from the instance startup time to the current time, in KB.
read_counts	Bigint	On a CN, it indicates the total number of logical reads performed by a user on all DNs in the last 5 seconds.
		On a DN, it indicates the total number of logical reads performed by a user from the instance startup time to the current time.
write_counts	Bigint	On a CN, it indicates the total number of logical writes performed by a user on all DNs in the last 5 seconds.
		On a DN, it indicates the total number of logical writes performed by a user from the instance startup time to the current time.
read_speed	Double precision	On a CN, it indicates the sum of average logical read rates of a user on all DNs in the last 5 seconds, in KB/s.
		On a DN, it indicates the average logical read rate of a user on the DN in the last 5 seconds, in KB/s.
write_speed	Double precision	On a CN, it indicates the sum of average logical write rates of a user on all DNs in the last 5 seconds, in KB/s.
		On a DN, it indicates the average logical write rate of a user on the DN in the last 5 seconds, in KB/s.

Column	Туре	Description
send_speed	Double precision	On a CN, it indicates the sum of the average network sending rates of a user on all DNs in the last 5 seconds, in KB/s.
		On a DN, it indicates the average network sending rate of a user on the DN in the last 5 seconds, in KB/s.
recv_speed	Double precision	On a CN, it indicates the sum of the average network receiving rates of a user on all DNs in the last 5 seconds, in KB/s.
		On a DN, it indicates the average network receiving rate of a user on the DN in the last 5 seconds, in KB/s.

14.3.174 PG_USER

PG_USER displays information about users who can access the database.

Table 14-244 PG_USER columns

Column	Туре	Description
usename	Name	User name
usesysid	OID	ID of this user
usecreatedb	boolean	Whether the user has the permission to create databases
usesuper	boolean	whether the user is the initial system administrator with the highest rights.
usecatupd	boolean	whether the user can directly update system tables. Only the initial system administrator whose usesysid is 10 has this permission. It is not available for other users.
userepl	boolean	Whether the user has the permission to duplicate data streams
passwd	Text	Encrypted user password. The value is displayed as ********.
valbegin	Timestamp with time zone	Account validity start time; null if no start time
valuntil	Timestamp with time zone	Password expiry time; null if no expiration

Column	Туре	Description	
respool	Name	Resource pool where the user is in	
parent	OID	Parent user OID	
spacelimit	Text	The storage space of the permanent table.	
tempspaceli mit	Text	The storage space of the temporary table.	
spillspacelimi t	Text	The operator disk flushing space.	
useconfig	Text[]	Session defaults for run-time configuration variables	
nodegroup	Name	Name of the logical cluster associated with the user. If no logical cluster is associated, this column is left blank.	

Example

Query the current database user list.

```
SELECT usename FROM pg_user;
usename
------
dbadmin
u1
u2
u3
(4 rows)
```

14.3.175 PG_USER_MAPPINGS

PG_USER_MAPPINGS displays information about user mappings.

This is essentially a publicly readable view of **PG_USER_MAPPING** that leaves out the options column if the user has no rights to use it.

Table 14-245 PG_USER_MAPPINGS columns

Column	Туре	Reference	Description
umid	OID	PG_USER_MAPPING.oid	OID of the user mapping
srvid	OID	PG_FOREIGN_SERVER.o id	OID of the foreign server that contains this mapping
srvname	Name	PG_FOREIGN_SERVER.s rvname	Name of the foreign server
umuser	OID	PG_AUTHID.oid	OID of the local role being mapped, 0 if the user mapping is public

Column	Туре	Reference	Description
usename	Name	-	Name of the local user to be mapped
umoption s	text[]	-	User mapping specific options. If the current user is the owner of the foreign server, its value is keyword=value strings. Otherwise, its value is null.

14.3.176 PG VIEWS

PG_VIEWS displays basic information about each view in the database.

Table 14-246 PG_VIEWS columns

Column	Туре	Reference	Description
schemana me	Name	PG_NAMESPACE.nspn ame	Name of the schema that contains the view
viewname	Name	PG_CLASS.relname	View name
viewowne r	Name	PG_AUTHID.rolname	Owner of the view
definition	Text	-	Definition of the view

Example

Query all the views in a specified schema.

14.3.177 PG_WLM_STATISTICS

PG_WLM_STATISTICS displays information about workload management after the task is complete or the exception has been handled. This view has been discarded in 8.1.2. You can use **PGXC_WLM_SESSION_INFO** to view load management records of completed jobs executed on all CNs.

Table 14-247 PG_WLM_STATISTICS columns

Column	Туре	Description
statement	Text	Statement executed for exception handling
block_time	Bigint	Block time before the statement is executed
elapsed_time	Bigint	Elapsed time when the statement is executed
total_cpu_time	Bigint	Total time used by the CPU on the DN when the statement is executed for exception handling
qualification_time	Bigint	Period when the statement checks the inclination ratio
cpu_skew_percent	Integer	CPU usage skew on the DN when the statement is executed for exception handling
control_group	Text	Cgroup used when the statement is executed for exception handling
status	Text	Statement status after it is executed for exception handling • pending: The statement is waiting to be
		executed.
		• running: The statement is being executed.
		• finished : The execution is finished normally.
		abort: The execution is unexpectedly terminated.
action	Text	Actions when statements are executed for exception handling
		abort indicates terminating the operation.
		adjust indicates executing the Cgroup adjustment operations. Currently, you can only perform the demotion operation.
		finish indicates that the operation is normally finished.
queryid	Bigint	Internal query ID used for statement execution
threadid	Bigint	ID of the backend thread

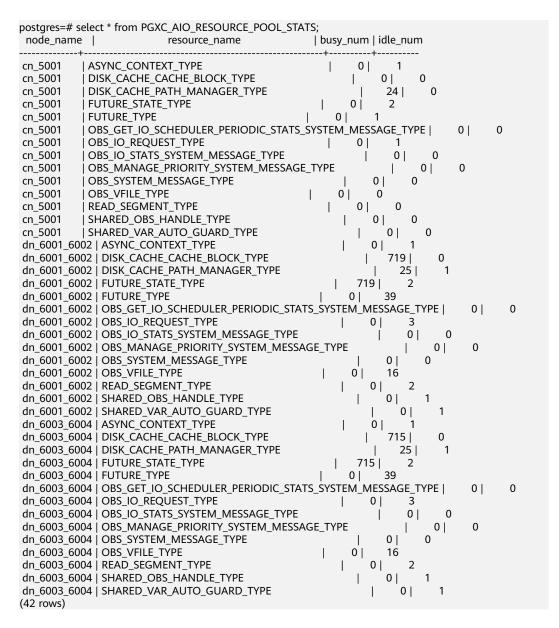
14.3.178 PGXC_AIO_RESOURCE_POOL_STATS

PGXC_AIO_RESOURCE_POOL_STATS queries the status of the asynchronous I/O resource pool usage for all nodes in the cluster. This includes the node name, the name of the asynchronous I/O resource type, the number of asynchronous I/O resources in use, and the number of idle asynchronous I/O resources. This view is supported only by 9.1.0.100 and later cluster versions.

 Table 14-248 PGXC_AIO_RESOURCE_POOL_STATS columns

Column	Туре	Description
node_na me	Text	Node name.
resource_name	Text	Asynchronous I/O resource type. The options are: ASYNC_CONTEXT_TYPE: Asynchronous context resource in the FPT (Future-Promise-Then) framework, at the thread level. DISK_CACHE_CACHE_BLOCK_TYPE: Instance of disk cache granularity block. DISK_CACHE_PATH_MANAGER_TYPE: Cache path manager in the disk cache, at the thread level. FUTURE_STATE_TYPE: Shared state FutureState of Future and Promise in the FPT framework. FUTURE_TYPE: Future in the FPT framework, which provides a non-blocking way to get the result of an asynchronous task. OBS_GET_IO_SCHEDULER_PERIODIC_STATS_SYST_EM_MESSAGE_TYPE: System message of type pgxc_obs_io_scheduler_periodic_stats in the asynchronous scheduling module statistics view. OBS_IO_REQUEST_TYPE: I/O request in the asynchronous scheduling module. OBS_IO_STATS_SYSTEM_MESSAGE_TYPE: System message of the pgxc_obs_io_scheduler_stats view in the asynchronous scheduling statistics module. OBS_MANAGE_PRIORITY_SYSTEM_MESSAGE_TYPE: System message used to adjust priority in the asynchronous scheduling module. OBS_SYSTEM_MESSAGE_TYPE: System message in the asynchronous scheduling module. OBS_SYSTEM_MESSAGE_TYPE: System message in the asynchronous scheduling module. OBS_VFILE_TYPE: Virtual file in OBS, OBS read/ write handle. READ_SEGMENT_TYPE: Entity that merges multiple OBS read requests. SHARED_OBS_HANDLE_TYPE: OBS Handler resource used to connect to the OBS service. SHARED_VAR_AUTO_GUARD_TYPE: Thread-level resource that manages OBS Handler and cached OBS file streams.
busy_num	Bigint	Number of asynchronous I/O resources in use.
idle_num	Bigint	Number of idle asynchronous I/O resources.

Example



14.3.179 PGXC_BULKLOAD_PROGRESS

PGXC_BULKLOAD_PROGRESS displays the progress of the service import. Only GDS common files can be imported. This view is accessible only to users with system administrators rights.

Table 14-249 PGXC_BULKLOAD_PROGRESS columns

Column	Туре	Description
session_id	Bigint	GDS session ID
query_id	Bigint	Query ID. It is equivalent to debug_query_id.
query	Text	Query statement

Column	Туре	Description
progress	Text	Progress percentage

14.3.180 PGXC BULKLOAD INFO

By querying the **PGXC_BULKLOAD_INFO** view on CNs, you can obtain historical statistics information for interconnection, GDS, COPY, and \COPY business executions after they have completed. This view summarizes the historical execution information of import and export business that have already completed on each node of the current cluster (including the interconnection cluster address, import and export business type, maximum, minimum, and total number of rows and bytes written to disk on DNs, etc.), to obtain historical information on import and export business execution and assist in performance troubleshooting.

This view does not record abnormal interruptions of import and export jobs. The data is directly obtained from the system catalog **GS_WLM_SESSION_INFO**, and the **loader_status** field is parsed to obtain import and export service information.

System administrator rights are required to access this view.

Table 14-250 PGXC_BULKLOAD_INFO columns

Column	Туре	Description
datid	OID	OID of the database the backend is connected to.
dbname	Text	Name of the database the backend is connected to.
schemaname	Text	Schema name.
nodename	Text	Name of the CN where the statement is run.
username	Text	Username for connecting to the backend.
application_na me	Text	Name of the application that is connected to the backend.
client_addr	inet	IP address of the client connected to the backend. If this column is null, it indicates that the client is connected via a Unix socket on the server machine or that it is an internal process, such as autovacuum.
client_hostnam e	Text	Host name of the client, which is obtained by reverse DNS lookup of client_addr. This column is only non-null when log_hostname is enabled and IP connection is used.

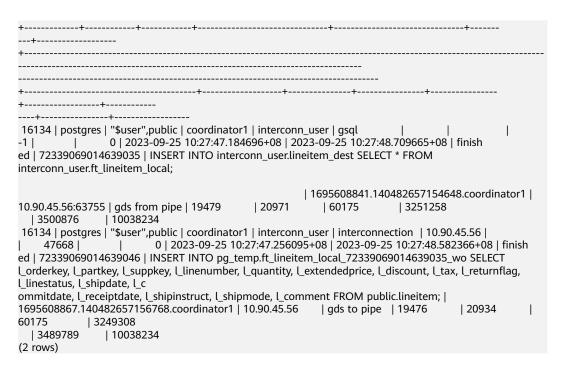
Column	Туре	Description
client_port	Integer	TCP port number used by the client to communicate with the backend. If a Unix socket is used, it is -1.
query_band	Text	Job type, which can be set using the GUC parameter query_band and is a null string by default.
block_time	Bigint	Blocking time before statement execution, including statement parsing and optimization time, in milliseconds.
start_time	Timestamp with time zone	Start time of statement execution.
finish_time	Timestamp with time zone	End time of statement execution.
status	Text	End status of statement execution: finished for normal and aborted for abnormal. The statement status recorded here should be the database server execution status. When the server-side execution is successful and an error occurs when the result set is returned, the statement should be finished.
queryid	Bigint	Internal query ID used for statement execution.
query	Text	Executed statement.
session_id	Text	A session uniquely identified in the database system, in the format of session_start_time.tid.node_name.
address	Text	Server address of the interconnection peer cluster. When not empty, it indicates an interconnection service, and the source cluster will additionally obtain the remote cluster port number.
direction	Text	Type of import and export service, including gds to file, gds from file, gds to pipe, gds from pipe, copy from, and copy to.
min_done_lines	json	Minimum number of rows of a statement across all DNs.
max_done_line s	json	Maximum number of rows of a statement across all DNs.

Column	Туре	Description
total_done_line s	json	Total number of rows of a statement across all DNs.
min_done_byte s	json	Minimum number of bytes of a statement across all DNs.
max_done_byte s	json	Maximum number of bytes of a statement across all DNs.
total_done_byt es	json	Total number of bytes of a statement across all DNs.

- Abnormal interruptions of import and export jobs are not recorded in the view.
- The implementation mechanism of GDS foreign tables and interconnection foreign tables is different. When querying, GDS records the full amount, while interconnection records the actual amount.
- For non-full import and export foreign tables with a limit, due to the special execution plan of limit, the data displayed is collected from one DN, which appears as a maximum value of all and a minimum value of 0.
- If the import and export table is a non-partitioned table:
 - When the GDS partitioned table is small, if one DN has finished collecting data and the other DNs have not started collecting data, they will not collect data. Therefore, when the data volume of GDS from non-partitioned tables is small, the minimum value may be 0, but it is not 0 when the table data volume is large.
 - When exporting non-partitioned tables from the interconnection source cluster, all DNs will be recorded, and only one DN's data will be collected, so the minimum value is 0.
 - When exporting replication tables from the interconnection remote cluster, only one DN will be recorded, so it is equivalent to having only one DN, and the minimum and maximum values are the same.
- Historical monitoring of import and export is implemented by reusing the historical TopSQL function, which follows the precautions, prerequisites, and operation steps of TopSQL. For details, refer to Historical Top SQL.
- Due to the large amount of data recorded by TopSQL, you are advised to query and use it as needed by combining fields such as **start_time** and **finish_time** to improve query performance, or to reduce query frequency.

Example

Use the **PGXC_BULKLOAD_INFO** view to query interconnection import service.



14.3.181 PGXC_BULKLOAD_STATISTICS

PGXC_BULKLOAD_STATISTICS displays real-time statistics about service execution, such as GDS, COPY, and \COPY, on a CN. This view summarizes the real-time execution status of import and export services that are being executed on each node in the current cluster. In this way, you can monitor the real-time progress of import and export services and locate performance problems.

Columns in PGXC_BULKLOAD_STATISTICS are the same as those in PG_BULKLOAD_STATISTICS. This is because PGXC_BULKLOAD_STATISTICS is essentially the summary result of querying PG_BULKLOAD_STATISTICS on each node in the cluster.

This view is accessible only to users with system administrators rights.

Table 14-251 PGXC_BULKLOAD_STATISTICS columns

Name	Туре	Description
node_name	Text	Node name
db_name	Text	Database name
query_id	Bigint	Query ID. It is equivalent to debug_query_id.
tid	Bigint	ID of the current thread
lwtid	Integer	Lightweight thread ID
session_id	Bigint	GDS session ID
direction	Text	Service type. The options are gds to file , gds from file , gds to pipe , gds from pipe , copy from , and copy to .

Name	Туре	Description
query	Text	Query statement
address	Text	Location of the foreign table used for data import and export
query_start	Timestamp with time zone	Start time of data import or export
total_bytes	Bigint	Total size of data to be processed
		This parameter is specified only when a GDS common file is to be imported and the record in the row comes from a CN. Otherwise, left this parameter unspecified.
phase	Text	Current phase. The options are INITIALIZING, TRANSFER_DATA, and RELEASE_RESOURCE.
done_lines	Bigint	Number of lines that have been transferred
done_bytes	Bigint	Number of bytes that have been transferred

14.3.182 PGXC_COLUMN_TABLE_IO_STAT

PGXC_COLUMN_TABLE_IO_STAT provides I/O statistics of all column-store tables of the database on all CNs and DNs in the cluster. Except the **nodename** column of the name type added in front of each row, the names, types, and sequences of other columns are the same as those in the **GS_COLUMN_TABLE_IO_STAT** view. For details about the columns, see **Table 14-252**.

Table 14-252 GS_COLUMN_TABLE_IO_STAT columns

Column	Туре	Description
schemaname	Name	Namespace of a table
relname	Name	Table name
heap_read	Bigint	Number of blocks logically read in the heap
heap_hit	Bigint	Number of block hits in the heap
idx_read	Bigint	Number of blocks logically read in the index
idx_hit	Bigint	Number of block hits in the index
cu_read	Bigint	Number of logical reads in the Compression Unit

Column	Туре	Description
cu_hit	Bigint	Number of hits in the Compression Unit
cidx_read	Bigint	Number of indexes logically read in the Compression Unit
cidx_hit	Bigint	Number of index hits in the Compression Unit

14.3.183 PGXC_COMM_CLIENT_INFO

PGXC_COMM_CLIENT_INFO stores the client connection information of all nodes. (You can query this view on a DN to view the information about the connection between the CN and DN.)

Table 14-253 PGXC_COMM_CLIENT_INFO columns

Column	Туре	Description
node_name	Text	Current node name.
арр	Text	Client application name.
tid	Bigint	Thread ID of the current thread.
lwtid	Integer	Lightweight thread ID of the current thread.
query_id	Bigint	Query ID. It is equivalent to debug_query_id.
socket	Integer	It is displayed if the connection is a physical connection.
remote_ip	Text	Peer node IP address.
remote_port	Text	Peer node port.
logic_id	Integer	If the connection is a logical connection, sid is displayed. If -1 is displayed, the current connection is a physical connection.

14.3.184 PGXC_COMM_DELAY

PGXC_COMM_DELAY displays the communication library delay status for all the DNs.

Table 14-254 PGXC_COMM_DELAY columns

Column	Туре	Description
node_name	Text	Node name

Column	Туре	Description
remote_name	Text	Name of the peer node with the maximum connection latency.
remote_host	Text	IP address of the peer
stream_num	Integer	Number of logical stream connections used by the current physical connection
min_delay	Integer	Minimum delay of the current physical connection. The unit is microsecond.
average	Integer	Average delay of the current physical connection. The unit is microsecond.
max_delay	Integer	Maximum delay of the current physical connection. The unit is microsecond.
		NOTE If its value is -1, the latency detection has timed out. In this case, re-establish the connection between nodes and then perform the query.

14.3.185 PGXC_COMM_RECV_STREAM

PG_COMM_RECV_STREAM displays the receiving stream status of the communication libraries for all the DNs.

Table 14-255 PGXC_COMM_RECV_STREAM columns

Column	Туре	Description
node_name	Text	Node name
local_tid	Bigint	ID of the thread using this stream
remote_name	Text	Name of the peer node
remote_tid	Bigint	Peer thread ID
idx	Integer	Peer DN ID in the local DN
sid	Integer	Stream ID in the physical connection
tcp_sock	Integer	TCP socket used in the stream

Column	Туре	Description
state	Text	Current status of the stream
		UNKNOWN: The logical connection is unknown.
		READY: The logical connection is ready.
		RUN: The logical connection receives packets normally.
		HOLD: The logical connection is waiting to receive packets.
		CLOSED: The logical connection is closed.
		• TO_CLOSED : The logical connection is to be closed.
		WRITING: Data is being written.
query_id	Bigint	debug_query_id corresponding to the stream
pn_id	Integer	<pre>plan_node_id of the query executed by the stream</pre>
send_smp	Integer	smpid of the sender of the query executed by the stream
recv_smp	Integer	smpid of the receiver of the query executed by the stream
recv_bytes	Bigint	Total data volume received from the stream. The unit is byte.
time	Bigint	Current life cycle service duration of the stream. The unit is ms.
speed	Bigint	Average receiving rate of the stream. The unit is byte/s.
quota	Bigint	Current communication quota value of the stream. The unit is Byte.
buff_usize	Bigint	Current size of the data cache of the stream. The unit is byte.

14.3.186 PGXC_COMM_SEND_STREAM

PGXC_COMM_SEND_STREAM displays the sending stream status of the communication libraries for all the DNs.

Table 14-256 PGXC_COMM_SEND_STREAM columns

Column	Туре	Description
node_name	Text	Node name

Column	Туре	Description	
local_tid	Bigint	ID of the thread using this stream	
remote_name	Text	Name of the peer node	
remote_tid	Bigint	Peer thread ID	
idx	Integer	Peer DN ID in the local DN	
sid	Integer	Stream ID in the physical connection	
tcp_sock	Integer	TCP socket used in the stream	
state	Text	 Current status of the stream. UNKNOWN: The logical connection is unknown. READY: The logical connection is ready. RUN: The logical connection sends packets normally. HOLD: The logical connection is waiting to send packets. CLOSED: The logical connection is closed. TO_CLOSED: The logical connection is to be closed. WRITING: Data is being written. 	
query_id	Bigint	debug_query_id corresponding to the stream	
pn_id	Integer	plan_node_id of the query executed by the stream	
send_smp	Integer	smpid of the sender of the query executed by the stream	
recv_smp	Integer	smpid of the receiver of the query executed by the stream	
send_bytes	Bigint	Total data volume sent by the stream. The unit is Byte.	
time	Bigint	Current life cycle service duration of the stream. The unit is ms.	
speed	Bigint	Average sending rate of the stream. The unit is Byte/s.	
quota	Bigint	Current communication quota value of the stream. The unit is Byte.	
wait_quota	Bigint	Extra time generated when the stream waits the quota value. The unit is ms.	

14.3.187 PGXC_COMM_STATUS

PGXC_COMM_STATUS displays the communication library status for all the DNs.

Table 14-257 PGXC_COMM_STATUS columns

Column	Туре	Description
node_name	Text	Node name.
rxpck/s	Integer	Receiving rate of the communication library on a node. The unit is byte/s.
txpck/s	Integer	Sending rate of the communication library on a node. The unit is byte/s.
rxkB/s	Bigint	Receiving rate of the communication library on a node. The unit is KB/s.
txkB/s	Bigint	Sending rate of the communication library on a node. The unit is KB/s.
buffer	Bigint	Size of the buffer of the Cmailbox.
memKB(libcomm)	Bigint	Communication memory size of the libcomm process, in KB.
memKB(libpq)	Bigint	Communication memory size of the libpq process, in KB.
%USED(PM)	Integer	Real-time usage of the postmaster thread.
%USED (sflow)	Integer	Real-time usage of the gs_sender_flow_controller thread.
%USED (rflow)	Integer	Real-time usage of the gs_receiver_flow_controller thread.
%USED (rloop)	Integer	Highest real-time usage among multiple gs_receivers_loop threads.
stream	Integer	Total number of used logical connections.

14.3.188 PGXC_COMM_QUERY_SPEED

PGXC_COMM_QUERY_SPEED displays traffic information about all queries on all nodes.

Table 14-258 PGXC_COMM_QUERY_SPEED columns

Column	Туре	Description
node_name	Text	Node name

Column	Туре	Description
query_id	Bigint	debug_query_id corresponding to the stream
rxkB/s	Bigint	Receiving rate of the query stream (unit: byte/s)
txkB/s	Bigint	Sending rate of the query stream (unit: byte/s)
rxkB	Bigint	Total received data of the query stream (unit: byte)
txkB	Bigint	Total sent data of the query stream (unit: byte)
rxpck/s	Bigint	Packet receiving rate of the query (unit: packets/s)
txpck/s	Bigint	Packet sending rate of the query (Unit: packets/s)
rxpck	Bigint	Total number of received packets of the query
txpck	Bigint	Total number of sent packets of the query

14.3.189 PGXC_DEADLOCK

PGXC_DEADLOCK displays lock wait information generated due to distributed deadlocks.

Currently, **PGXC_DEADLOCK** collects only lock wait information about locks whose **locktype** is **relation**, **partition**, **page**, **tuple**, or **transactionid**.

Table 14-259 PGXC_DEADLOCK columns

Column	Туре	Description
locktype	Text	Type of the locked object
nodename	Name	Name of the node where the locked object resides
dbname	Name	Name of the database where the locked object resides. The value is NULL if the locked object is a transaction.
nspname	Name	Name of the namespace of the locked object
relname	Name	Name of the relation targeted by the lock. The value is NULL if the object is not a relation or part of a relation.

Column	Туре	Description			
partname	Name	Name of the partition targeted by the lock. The value is NULL if the locked object is not a partition.			
page	Integer	Number of the page targeted by the lock. The value is NULL if the locked object is neither a page nor a tuple.			
tuple	Smallint	Number of the tuple targeted by the lock. The value is NULL if the locked object is not a tuple.			
transactioni d	Xid	ID of the transaction targeted by the lock. The value is NULL if the locked object is not a transaction.			
waituserna me	Name	Name of the user who waits for the lock			
waitgxid	Xid	ID of the transaction that waits for the lock			
waitxactstar t	Timestamp with time zone	Start time of the transaction that waits for the lock			
waitqueryid	Bigint	Latest query ID of the thread that waits for the lock			
waitquery	Text	Latest query statement of the thread that waits for the lock			
waitpid	Bigint	ID of the thread that waits for the lock			
waitmode	Text	Mode of the waited lock			
holduserna me	Name	Name of the user who holds the lock			
holdgxid	Xid	ID of the transaction that holds the lock			
holdxactstar t	Timestamp with time zone	Start time of the transaction that holds the lock			
holdqueryid	Bigint	Latest query ID of the thread that holds the lock			
holdquery	Text	Latest query statement of the thread that holds the lock			
holdpid	Bigint	ID of the thread that holds the lock			
holdmode	Text	Mode of the held lock			
waittime	Timestamp with time zone	Timestamp when the lock wait starts. This column is available only in clusters of version 9.1.0.200 or later.			

Column	Туре	Description
holdtime	Timestamp with time zone	Timestamp when the lock starts to be held. This column is available only in clusters of version 9.1.0.200 or later.

14.3.190 PGXC DISK CACHE STATS

PGXC_DISK_CACHE_STATS records the usage of file cache. This system view is supported only by clusters of version 9.1.0 or later.

Table 14-260 PGXC_DISK_CACHE_STATS columns

Column	Туре	Description
node_name	Text	Node name.
total_read	Bigint	Total number of accesses to disk cache.
local_read	Bigint	Total number of times disk cache reads from local disk.
remote_read	Bigint	Total number of times disk cache reads from remote storage.
hit_rate	numeric(5,2)	Hit rate of disk cache.
cache_size	Bigint	Total size of data saved in disk cache, in KB.
fill_rate	numeric(5,2)	Fill rate of disk cache.

Example

Query the hit rate of disk cache on each node.



14.3.191 PGXC_DISK_CACHE_ALL_STATS

PGXC_DISK_CACHE_ALL_STATS records all usage of file cache. This system view is supported only by clusters of version 9.1.0 or later.

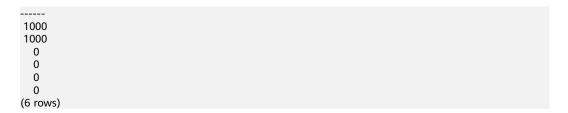
Table 14-261 PGXC_DISK_CACHE_ALL_STATS columns

Column	Туре	Description
node_name	Text	Node name.
total_read	Bigint	Total number of accesses to disk cache.
local_read	Bigint	Total number of times disk cache accesses local disk.
remote_read	Bigint	Total number of times disk cache accesses remote storage.
hit_rate	numeric(5,2)	Hit rate of disk cache.
cache_size	Bigint	Total size of data saved in disk cache, in KB.
fill_rate	numeric(5,2)	Fill rate of disk cache.
temp_file_size	Bigint	Total size of temporary/cold cache files, in KB.
a1in_size	Bigint	Total size of data saved in the alin queue of disk cache, in KB.
a1out_size	Bigint	Total size of data saved in the a1out queue of disk cache, in KB.
am_size	Bigint	Total size of data saved in the am queue of disk cache, in KB.
a1in_fill_rate	numeric(5,2)	Fill rate of the a1in queue in disk cache.
a1out_fill_rate	numeric(5,2)	Fill rate of the a1out queue in disk cache.
am_fill_rate	numeric(5,2)	Fill rate of the am queue in disk cache.
fd	Integer	Number of file descriptors currently in use by disk cache.
pin_block_count	Bigint	Number of pinned blocks in disk cache. This field is supported only by 9.1.0.100 and later cluster versions.

Example

Query the number of file descriptors used by disk cache on each node.

SELECT fd FROM pgxc_disk_cache_all_stats; fd



14.3.192 PGXC DISK CACHE PATH INFO

PGXC_DISK_CACHE_PATH_INFO records information about the hard disk where the file cache is stored. This system view is supported only by clusters of version 9.1.0 or later.

Table 14-262 PGXC_DISK_CACHE_PATH_INFO columns

Column	Туре	Description
path_name	Text	Path name.
node_name	Text	Name of the node the hard disk belongs to.
cache_size	Bigint	Total size of cache files in the hard disk, in bytes.
disk_available	Bigint	Available space in the hard disk, in bytes.
disk_size	Bigint	Total capacity of the hard drive, in bytes.
disk_use_ratio	Double precision	Disk space usage.

Example

Query information about the hard disk used by the file cache.

```
SELECT * FROM pgxc_disk_cache_path_info order by 1;
 path_name | node_name | cache_size | disk_available | disk_size | disk_use_ratio
                                  19619 | 137401716736 | 160982630400 | .146481105479564
dn_6001_6002_0 | dn_6001_6002 |
dn_6001_6002_1 | dn_6001_6002 |
                                  35968 |
                                           137401716736 | 160982630400 | .146481105479564
dn_6003_6004_0 | dn_6003_6004 |
                                  27794 |
                                          121600655360 | 160982630400 | .244634933235629
dn_6003_6004_1 | dn_6003_6004 |
                                  26158 |
                                          121600655360 | 160982630400 | .244634933235629
dn_6005_6006_0 | dn_6005_6006 |
                                  24533
                                           134394839040 | 160982630400 | .165159379579873
dn_6005_6006_1 | dn_6005_6006 | 31065 | 134394839040 | 160982630400 | .165159379579873
```

14.3.193 PGXC_GET_STAT_ALL_TABLES

PGXC_GET_STAT_ALL_TABLES displays information about insertion, update, and deletion operations on tables and the dirty page rate of tables.

Before running **VACUUM FULL** on a system catalog with a high dirty page rate, ensure that no user is performing operations on it.

You are advised to run **VACUUM FULL** to tables (excluding system catalogs) whose dirty page rate exceeds 80% or run it based on service scenarios.

■ NOTE

For clusters of 8.2.0.100 or later, **PGXC_STAT_TABLE_DIRTY** is recommended for querying the dirty page rate.

Table 14-263 PGXC GET STAT ALL TABLES columns

Column	Туре	Description
relid	OID	Table OID
relname	Name	Table name
schemaname	Name	Schema name of the table
n_tup_ins	Numeric	Number of inserted tuples
n_tup_upd	Numeric	Number of updated tuples
n_tup_del	Numeric	Number of deleted tuples
n_live_tup	Numeric	Number of live tuples
n_dead_tup	Numeric	Number of dead tuples
dirty_page_rate	numeric(5,2	Dirty page rate (%) of a table

Examples

Use the view **PGXC_GET_STAT_ALL_TABLES** to query the tables whose dirty page rate is greater than 30%.

SELECT * FROM PGXC_GE* relid relname dirty_page_rate	Γ_STAT_ALL_TABLI schemaname r	n_tup_ins	n_tup_up			/e_tup n	_dead_tup
+							
2840 pg_toast_2619	pg_toast	7415	0	7415	0	291	88.00
9001 pgxc_class	pg_catalog	56331	3	56285	54	143	72.59
53860 reason	dbadmin	9	19	0	9	19	67.86
9025 pg_object	pg_catalog	112858	1179707	11261	9 2	246	429
63.56							
9015 pgxc_node	pg_catalog	15	24	0	15	24	61.54
2606 pg_constraint	pg_catalog	78	0	42	36	42	53.85
1260 pg_authid	pg_catalog	6	6	0	6	6	50.00
(7 rows)							

14.3.194 PGXC_GET_STAT_ALL_PARTITIONS

PGXC_GET_STAT_ALL_PARTITIONS displays information about insertion, update, and deletion operations on partitions of partitioned tables and the dirty page rate of tables.

The statistics of this view depend on the **ANALYZE** operation. To obtain the most accurate information, perform the **ANALYZE** operation on the partitioned table first.

□ NOTE

For clusters of 8.2.0.100 or later, **PGXC_STAT_TABLE_DIRTY** is recommended for querying the dirty page rate.

Table 14-264 PGXC_GET_STAT_ALL_PARTITIONS columns

Column	Туре	Description
relid	OID	Table OID
partid	OID	Partition OID
schemaname	Name	Schema name of the table
relname	Name	Table name
partname	Name	Partition name
n_tup_ins	Numeric	Number of inserted tuples
n_tup_upd	Numeric	Number of updated tuples
n_tup_del	Numeric	Number of deleted tuples
n_live_tup	Numeric	Number of live tuples
n_dead_tup	Numeric	Number of dead tuples
page_dirty_rate	numeric(5,2	Dirty page rate (%) of a table

Example

Query partition tables whose dirty page rate is greater than 30%.

SELECT * FROM PGXC_GET_STAT_ALL_PARTITIONS W relid partid schemaname relname pa n_dead_tup dirty_page_rate	rtname	n_tup_	ins n_tu	ıp_upd			tup
+							
58320 58626 schema_subquery store_hash_par	p1	1	2	0	2	0	2
58430 58706 schema_subquery store_hash_par_i 2 100.00	mor p4	-	1	1	1	0	
58320 58644 schema_subquery store_hash_par	p1		3	0	3	0	3
58430 58770 schema_subquery store_hash_par_i 2 100.00	mor p4	- 1	1	1	1	0	
58320 58643 schema_subquery store_hash_par	p1	1	2	0	2	0	2
58320 58625 schema_subquery store_hash_par 100.00	p1	1	2	0	2	0	2
58320 58579 schema_subquery store_hash_par	p1	1	2	0	2	0	2
58320 58619 schema_subquery store_hash_par	p1	1	3	0	3	0	3
58320 58627 schema_subquery store_hash_par	p1	1	4	0	4	0	4
58320 58657 schema_subquery store_hash_par 100.00	p1	1	3	0	3	0	3
(10 rows)							

14.3.195 PGXC_GET_TABLE_SKEWNESS

PGXC_GET_TABLE_SKEWNESS displays the data skew on tables in the current database. Only the system administrator or the preset role **gs_role_read_all_stats** can access this view.

You are advised to use **PGXC_GET_TABLE_SKEWNESS** to query data skew when there are less than 10,000 tables in a database.

Table 14-265 PGXC_GET_TABLE_SKEWNESS columns

Column	Туре	Description
schemaname	Name	Schema name of a table
tablename	Name	Name of a table
totalsize	Numeric	Total size of a table, in bytes
avgsize	numeric(1000, 0)	Average table size (total table size divided by the number of DNs), which is the ideal size of tables distributed on each DN
maxratio	numeric(10,3)	Ratio of the maximum table size on a single DN to the value of avgsize.
minratio	numeric(10,3)	Ratio of the minimum table size on a single DN to avgsize
skewsize	Bigint	Table skew rate (the maximum table size on a single DN minus the minimum table size on a single DN)
skewratio	numeric(10,3)	Table skew rate (skewsize/avgsize)
skewstddev	numeric(1000, 0)	Standard deviation of table distribution (For two tables of the same size, a larger deviation indicates a more severe skew.)

14.3.196 PGXC_GLOBAL_TEMP_ATTACHED_PIDS

This view displays information about sessions of resources occupied by global temporary tables on CNs. This view is supported only by clusters of 8.2.1.220 and later versions.

ColumnTypeDescriptionnodenameNameNode nameschemanameNameSchema nametablenameNameTable namepidBigintPID of a session

Table 14-266 PG_GLOBAL_TEMP_ATTACHED_PIDS columns

14.3.197 PGXC_GTM_SNAPSHOT_STATUS

PGXC_GTM_SNAPSHOT_STATUS displays transaction information on the current GTM.

Table 14-267 PGXC_GTM_SNAPSHOT_STATUS columns

Column	Туре	Description
xmin	Xid	Minimum ID of the running transactions
xmax	Xid	ID of the transaction next to the executed transaction with the maximum ID
csn	Integer	Sequence number of the transaction to be committed
oldestxmin	Xid	Minimum ID of the executed transactions
xcnt	Integer	Number of the running transactions
running_xids	Text	IDs of the running transactions

14.3.198 PGXC_INSTANCE_TIME

PGXC_INSTANCE_TIME displays the running time of processes on each node in the cluster and the time consumed in each execution phase. Except the node_name column, the other columns are the same as those in the PV_INSTANCE_TIME view. Only the system administrator or the preset role gs_role_read_all_stats can access this view.

Table 14-268 PGXC_INSTANCE_TIME columns

Column	Туре	Description
node_name	Text	Node name
stat_id	Integer	Type ID
stat_name	Text	Name of the runtime type

Column	Туре	Description
value	Bigint	Runtime value

14.3.199 PGXC_LOCKWAIT_DETAIL

PGXC_LOCKWAIT_DETAIL displays detailed information about the lock wait hierarchy on each node in a cluster. If a node has multiple lock wait levels, the entire lock waiting hierarchy is displayed in sequence.

This view is supported only by clusters of version 8.1.3.200 or later.

Table 14-269 PGXC_LOCKWAIT_DETAIL columns

Column	Туре	Description
level	Integer	Level in the lock wait hierarchy. The value starts with 1 and increases by 1 when there is a wait relationship.
node_name	Name	Node name, corresponding to the node_name column in the pgxc_node table.
lock_wait_hi erarchy	Text	Lock wait hierarchy, in the format of <i>node-name</i> : process-ID-> waiting-process-ID-> waiting-
lock_type	Text	Type of the locked object
database	OID	OID of the database where the locked object is.
relation	OID	OID of the relationship of the locked object.
page	Integer	Page index in a relationship
tuple	Smallint	Row number of a page.
virtual_xid	Text	Virtual ID of a transaction.
transaction_i	Xid	Transaction ID.
class_id	OID	OID of the system catalog that contains the object.
obj_id	OID	OID of the object within its system catalog.
obj_subid	Smallint	Column number of a table
virtual_trans action	Text	Virtual ID of the transaction holding or waiting for the lock.
pid	Bigint	ID of the thread holding or awaiting this lock
mode	Text	Lock level

Column	Туре	Description
granted	Boolean	Indicates whether a lock is held.
fastpath	Boolean	Indicates whether to obtain a lock using FASTPATH.
wait_for_pid	Bigint	ID of the thread where a lock conflict occurs.
conflict_mod e	Text	Level of the conflicted lock held by the thread where it is
query_id	Bigint	ID of a query statement.
query	Text	Query statement
application_ name	Text	Name of the application connected to the backend
backend_star t	Timestamp with time zone	Startup time of the backend process, that is, the time when the client connects to the server
xact_start	Timestamp with time zone	Start time of the current transaction
query_start	Timestamp with time zone	Start time of the active query
state	Text	Overall state of the backend
waittime	Timestamp with time zone	Timestamp when the lock wait starts. This column is available only in clusters of version 9.1.0.200 or later.
holdtime	Timestamp with time zone	Timestamp when the lock starts to be obtained. This column is available only in clusters of version 9.1.0.200 or later.

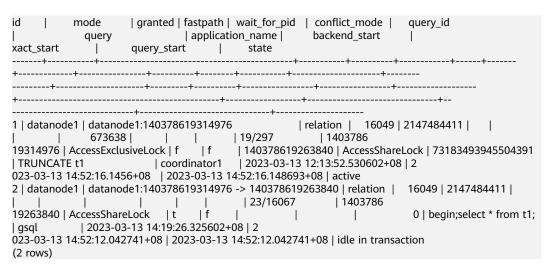
Example

- **Step 1** Connect to the DN, start a transaction, and run the following command: begin;select * from t1;
- **Step 2** Connect to the CN in another window and truncate table **t1**. truncate t1;

In this case, truncation is blocked.

Step 3 Open another window to connect to the CN and run the **select * from pgxc_lockwait_detail;** command.

```
SELECT * FROM PGXC_LOCKWAIT_DETAIL;
level | node_name | lock_wait_hierarchy | lock_type | database | relation | page | tuple |
virtual_xid | transaction_id | class_id | obj_id | obj_subid | virtual_transaction | p
```



----End

14.3.200 PGXC_INSTR_UNIQUE_SQL

PGXC_INSTR_UNIQUE_SQL displays the complete Unique SQL statistics of all CN nodes in the cluster.

Only the system administrator can access this view. The columns in this view are the same as those in the **GS_INSTR_UNIQUE_SQL** view. For details about columns in the view, see **Table 14-270**.

Table 14-270 GS_INSTR_UNIQUE_SQL columns

Column	Туре	Description
node_name	Name	Name of the CN that receives SQL statements
node_id	Integer	Node ID, which is the same as the value of node_id in the pgxc_node table
user_name	Name	Username
user_id	OID	User ID
unique_sql_id	Bigint	Normalized Unique SQL ID
query	Text	Normalized SQL text
n_calls	Bigint	Number of successful execution times
min_elapse_time	Bigint	Minimum running time of the SQL statement in the database (unit: μs)

Column	Туре	Description
max_elapse_time	Bigint	Maximum running time of SQL statements in the database (unit: μs)
total_elapse_time	Bigint	Total running time of SQL statements in the database (unit: µs)
n_returned_rows	Bigint	Row activity - Number of rows in the result set returned by the SELECT statement
n_tuples_fetched	Bigint	Row activity - Randomly scan rows (column-store tables/foreign tables are not counted.)
n_tuples_returned	Bigint	Row activity - Sequential scan rows (Column-store tables/foreign tables are not counted.)
n_tuples_inserted	Bigint	Row activity - Inserted rows
n_tuples_updated	Bigint	Row activity - Updated rows
n_tuples_deleted	Bigint	Row activity - Deleted rows
n_blocks_fetched	Bigint	Block access times of the buffer, that is, physical read/I/O
n_blocks_hit	Bigint	Block hits of the buffer, that is, logical read/ cache
n_soft_parse	Bigint	Number of soft parsing times (cache plan)
n_hard_parse	Bigint	Number of hard parsing times (generation plan)

Column	Туре	Description
db_time	Bigint	Valid DB execution time, including the waiting time and network sending time. If multiple threads are involved in query execution, the value of DB_TIME is the sum of DB_TIME of multiple threads (unit: µs).
cpu_time	Bigint	CPU execution time, excluding the sleep time (unit: µs)
execution_time	Bigint	SQL execution time in the query executor, DDL statements, and statements (such as Copy statements) that are not executed by the executor are not counted (unit: µs).
parse_time	Bigint	SQL parsing time (unit: μs)
plan_time	Bigint	SQL generation plan time (unit: μs)
rewrite_time	Bigint	SQL rewriting time (unit: μs)
pl_execution_time	Bigint	Execution time of the plpgsql procedural language function (unit: µs)
pl_compilation_time	Bigint	Compilation time of the plpgsql procedural language function (unit: µs)
net_send_time	Bigint	Network time, including the time spent by the CN in sending data to the client and the time spent by the DN in sending data to the CN (unit: µs)
data_io_time	Bigint	File I/O time (unit: μs)

Column	Туре	Description
first_time	Timestamp with time zone	Time of the first SQL statement execution
last_time	Timestamp with time zone	Time of the last SQL statement execution

14.3.201 PGXC_LOCK_CONFLICTS

PGXC_LOCK_CONFLICTS displays information about conflicting locks in the cluster.

When a lock is waiting for another lock or another lock is waiting for this one, a lock conflict occurs.

Currently, **PGXC_LOCK_CONFLICTS** collects only information about locks whose **locktype** is **relation**, **partition**, **page**, **tuple**, or **transactionid**.

Table 14-271 PGXC_LOCK_CONFLICTS columns

Column	Туре	Description
locktype	Text	Type of the locked object
nodename	Name	Name of the node where the locked object resides
dbname	Name	Name of the database where the locked object resides. The value is NULL if the locked object is a transaction.
nspname	Name	Name of the namespace of the locked object
relname	Name	Name of the relation targeted by the lock. The value is NULL if the object is not a relation or part of a relation.
partname	Name	Name of the partition targeted by the lock. The value is NULL if the locked object is not a partition.
page	Integer	Number of the page targeted by the lock. The value is NULL if the locked object is neither a page nor a tuple.
tuple	Smallint	Number of the tuple targeted by the lock. The value is NULL if the locked object is not a tuple.
transactionid	Xid	ID of the transaction targeted by the lock. The value is NULL if the locked object is not a transaction.

Column	Туре	Description
username	Name	Name of the user who applies for the lock
gxid	Xid	ID of the transaction that applies for the lock
xactstart	Timestamp with time zone	Start time of the transaction that applies for the lock
queryid	Bigint	Latest query ID of the thread that applies for the lock
query	Text	Latest query statement of the thread that applies for the lock
pid	Bigint	ID of the thread that applies for the lock
mode	Text	Lock mode
granted	Boolean	 TRUE if the lock has been held FALSE if the lock is still waiting for another lock

14.3.202 PGXC_LWLOCKS

PGXC_LWLOCK offers details on lightweight locks that are currently held or being waited for by all instances in the cluster. This view is supported only by 9.1.0.200 and later cluster versions.

Table 14-272 PGXC_LWLOCKS columns

Column	Туре	Description
nodename	Name	Name of the node where the locked object resides
pid	Bigint	ID of the backend thread
query_id	Bigint	ID of a query
lwtid	Integer	Lightweight thread ID of the backend thread
reqlockid	Integer	ID of the lightweight lock that is being requested by the current thread
reglock	Text	Name of the lightweight lock corresponding to reqlockid
heldlocknums	Integer	Number of lightweight locks obtained by the current thread
heldlockid	Integer	Lightweight lock ID obtained by the current thread

Column	Туре	Description
heldlock	Text	Name of the lightweight lock corresponding to heldlockid
heldlockmode	Text	Lightweight lock mode corresponding to heldlockid

Example

Use the **PGXC_LWLOCKS** view to get details on lightweight locks that are currently held or being waited for by all instances in the cluster.

SELECT * FROM pgxc_lwlocks; nodename pid query_id lwtid reqlockid reqlo neldlock heldlockmode			s heldl	ockid
	+		+	
+				
datanode1 139810224193360 78250043525924188 54844		1	1	76390
BUFFER_POOL_LWLOCK Shared	·	·	·	•
datanode1 139810224198200 78250043525924886 54922	1	- 1	1	957438
PGPROC LWLOCK Exclusive		•		
datanode2 140262654050288 0 54832	1	1	7	
WALWriteLock Exclusive	'			
datanode2 140262654052488 78250043525923195 54847	1	1	1	15862
BUFFER POOL LWLOCK Shared	'		* 1	
(4 rows)				
\ · · - · · - /				

14.3.203 PGXC_MEMORY_DEBUG_INFO

PGXC_MEMORY_DEBUG_INFO displays memory error information for each node in the current cluster when executing jobs, making it easy to locate memory error issues. When an error message "memory is temporarily unavailable" is prompted during statement execution, this view can be used to query memory error information for all nodes, which is the same as the memory error information displayed in the log. This view is supported only by clusters of version 8.3.0 or later.

NOTICE

This view only displays the most recent cluster information for errors, and repeated error information will be overwritten. If the same query requests memory multiple times and errors occur, the information will not be updated.

Table 14-273 PGXC_MEMORY_DEBUG_INFO columns

Column	Туре	Description
node_name	Text	Instance name, including CNs and DNs.
query_id	Bigint	ID of the query that is currently requesting memory.

Column	Туре	Description
memory_info	Text	Current instance's memory usage, including: process_used_memory: memory size used by the Gauss DR (DWS) process
		the GaussDB(DWS) process.max_dynamic_memory: maximum dynamic memory.
		dynamic_used_memory: used dynamic memory.
		dynamic_peak_memory: dynamic peak value of memory.
		• dynamic_used_shrctx : maximum dynamic shared memory context.
		 dynamic_peak_shrctx: dynamic peak value of shared memory context.
		• shared_used_memory: used shared memory.
		 cstore_used_memory: memory size used for column store.
		• comm_used_memory : memory size used by the communication library.
		 comm_peak_memory: peak value of memory used by the communication library.
		 other_used_memory: memory size used by other components.
		 topsql_used_memory: memory size used by topsql.
		• large_storage_memory: memory size used for column-store compression and decompression.
		 os_totalmem: total memory size of the operating system.
		os_freeemem: remaining memory size of the operating system.
summary	Text	The total estimated memory consumption and actual memory consumption of jobs on the instance.
abnormal_qu ery	Text	Thread ID and query ID with abnormal memory usage, including two cases:
		Session with the maximum current memory usage.
		Session with the largest difference between estimated memory and actual memory usage.
abnormal_me mory	Text	Memory block with the highest usage, including the maximum shared memory context usage and the maximum common memory context usage.

Column	Туре	Description
top_thread	Text	Information on the top three threads with the highest memory usage:
		context name: memory block currently in use.
		contextlevel: context level.
		sessType: type of the top-level context node.
		totalsize[274,13,260]MB: total memory, released memory, and used memory size of the current memory context, in MB.
create_time	Timestam p with time zone	Time when the memory shortage error occurred.

14.3.204 PGXC_NODE_ENV

PGXC_NODE_ENV displays the environmental variables information about all nodes in a cluster.

Table 14-274 PGXC_NODE_ENV columns

Column	Туре	Description
node_name	Text	Names of all nodes in the cluster.
host	Text	Host names of all nodes in the cluster.
process	Integer	Process IDs of all nodes in the cluster.
port	Integer	Port numbers of all nodes in the cluster.
installpath	Text	Installation directory of all nodes in the cluster.
datapath	Text	Data directories of all nodes in the cluster.
log_directory	Text	Log directories of all nodes in the cluster.

14.3.205 PGXC_NODE_STAT_RESET_TIME

PGXC_NODE_STAT_RESET_TIME displays the time when statistics of each node in the cluster are reset. All columns except **node_name** are the same as those in the **GS_NODE_STAT_RESET_TIME** view. This view is accessible only to users with system administrators rights.

Table 14-275 PGXC_NODE_STAT_RESET_TIME columns

Column	Туре	Description
node_name	Text	Node name
reset_time	Timestamp	Time when statistics on each node are reset

14.3.206 PGXC OBS IO SCHEDULER STATS

PGXC_OBS_IO_SCHEDULER_STATS displays the latest real-time statistics about read/write requests of the OBS I/O Scheduler. This system view is supported only by clusters of version 9.1.0 or later.

Table 14-276 PGXC_OBS_IO_SCHEDULER_STATS columns

Column	Туре	Description
node_name	Text	Node name.
io_type	Char	Type of I/O operation, including:
		• r : read.
		• w: write.
		• s: file operation.
current_bps	INT8	Current bandwidth rate, in KB/s.
best_bps	INT8	Best bandwidth rate achieved recently, in KB/s.
waiting_request_n um	Int	Number of queued requests currently waiting.
mean_request_size	INT8	Average length of requests processed recently, in KB.
total_token_num	Int	Total number of I/O tokens.
available_token_n um	Int	Number of available I/O tokens.
total_worker_num	Int	Total number of working threads.
idle_worker_num	Int	Number of idle working threads.

Example

Step 1 Query statistics about read requests of OBS I/O Scheduler on each node:

SELECT * FROM pgxc_obs_io_scheduler_stats WHERE io_type = 'r' ORDER BY node_name;

node_name | io_type | current_bps | best_bps | waiting_request_num | mean_request_size | total_token_num | available_token_num | total_worker_num | idle_worker_num

+	+	·				
dn_6001_6002 r	26990	26990	0	215	18	16
12	10					
dn_6003_6004 r	21475	21475	10	190	30	30
20	20					
dn_6005_6006 r	12384	12384	36	133	30	27
20	17					

According to the result, this is a snapshot of the statistics at a certain time point when the current I/O scheduler reads I/Os. At this time, the bandwidth is increasing, and **current_bps** is equal to **best_bps**. Take **dn_6003_6004** as an example. You can see that there are queuing requests on the current DN. The value of **total_token_num** is the same as that of **available_token_num**, indicating that the I/O scheduler has not started to process these requests when the view is queried.

Step 2 Wait for a while and initiate the query again.

SELECT * FROM paxc obs io scheduler stats WHERE io type = 'r' ORDER BY node name; node_name | io_type | current_bps | best_bps | waiting_request_num | mean_request_size | total_token_num | available_token_num | total_worker_num | idle_worker_num dn_6001_6002 | r | 13228 | 26990 | 0 | 609 | 18 | 18 12 | 12 dn_6003_6004 | r | 15717 | 21475 | 0 | 30 I 30 622 I 20 | 20 | 18041 | 21767 | 20 dn_6005_6006 | r 0 | 609 | 30 | 20 |

When the queue is empty and the value of available_token_num is equal to that of total_token_num, it indicates that the I/O scheduler has finished processing all requests and there are no new requests in line. The current_bps value is not 0 because it represents the average bandwidth (in bit/s) over a three-second period. Therefore, the displayed value reflects the data from three seconds ago.

Step 3 After a short period of time, the query result is as follows. The value of **current bps** changes to **0**.

SELECT * FROM pgx	_obs_io_	sched	luler_stats WH	IERE io_type = 'r'	ORDER BY n	ode_name;	
node_name io_type current_bps best_bps waiting_request_num mean_request_size total_token_num available_token_num total_worker_num idle_worker_num 							
+t dn 6001 6002 r	· 	 0 l	+ 26990	 0 l	609	18	18
12	12	- 1		- 1			
dn_6003_6004 r		0	21475	0	622	30	30
20 dn 6005 6006 r	20	0	21767	0	609 l	30	30
20	20	١٠	21707	0	009	30	30

----End

14.3.207 PGXC_OBS_IO_SCHEDULER_PERIODIC_STATS

PGXC_OBS_IO_SCHEDULER_PERIODIC_STATS provides statistics on the number of requests and flow control information for different types of OBS I/O Scheduler requests, including read, write, and file operations. This system view is supported only by clusters of version 9.1.0 or later.

The first query result shows the statistics from the cluster startup to the query time, with detailed columns listed in the table below.

Table 14-277 PGXC_OBS_IO_SCHEDULER_PERIODIC_STATS columns

Column	Туре	Description
node_name	Name	Name of a CN or DN, for example, dn_6001_6002.
io_type	Char	Type of I/O operation, including: • R: read • W: write • S: file operation
recent_throttled_req_nu m	Int	Number of times flow control was applied between two query views.
total_throttled_req_num	Int	Total number of times flow control was applied.
last_throttled_dur(s)	INT8	Time interval since the last occurrence of flow control.
waiting_req_num	Int	Number of queued requests currently waiting.
mean_tps	numeric(7,2)	Average TPS (transactions per second) between two query views.
mean_req_size(KB)	INT8	Average length of requests between two query views, in KB.
mean_req_latency(ms)	INT8	Average latency of requests between two query views, in ms.
max_req_latency(ms)	INT8	Maximum latency of requests before two query views, in ms.
mean_bps(KB/s)	INT8	Average read or write speed between two query views, in KB/s.
duration(s)	Int	Time interval between two query views, in seconds.

Example

Run the **SELECT * FROM pgxc_obs_io_scheduler_periodic_stats** statement to query the view content. The following is an example of the query result: SELECT * FROM pgxc_obs_io_scheduler_periodic_stats;

node name | io type | recent throttled reg num | total throttled reg num | last throttled dur(s) | waiting_req_num | mean_tps | mean_req_size(KB) | mean_req_latency(ms) | max_req_latency(ms) | mean_bps(KB/s) | duration(s) -----+--dn_6001_6002 | S 0 | 0 | 0 | 0 1 0.00 0 | 0 | 0 1 0 | 155 dn_6001_6002 | R 0 | 0 | 0 | 0 | 0.00 0 | 0 | 155 0 1 0 | dn_6001_6002 | W 0 | 0 | 0 | 0.00 0 1 0 1 155 0 [cn_5001 0 | | S 0 | 0 | .03 [207 | 0 | 519 | 0 | 155 cn_5001 0 | 0.00 | | R 0 0 | 0 [0 | 0 | 0 | cn_5001 | W 0 | 0 1 0 [.01 0 0 | 288 | 288 | 0 | 155 (6 rows)

To display **0** before the decimal point in the value of **mean_tps**, set the **display_leading_zero** option in the **behavior_compat_options** parameter.

Run the **select * from pgxc_obs_io_scheduler_periodic_stats** statement. The following information is displayed:

SELECT * FROM pgxc_obs_io_scheduler_periodic_stats; node_name | io_type | recent_throttled_req_num | total_throttled_req_num | last_throttled_dur(s) | waiting_req_num | mean_tps | mean_req_size(KB) | mean_req_latency(ms) | max_req_latency(ms) | mean_bps(KB/s) | duration(s) dn_6001_6002 | S 0 1 0 | 0.36 326 | 0 | 0 | 132 | 177 dn_6001_6002 | R 0 | 0 | 0 | 0 | 0.00 0 | 177 0 0 | 0 | dn_6001_6002 | W 0 | 0 1 0 | 0 [0.00 0 | 0 | 0 | 177 cn_5001 | S 0 | 0 | 0 | 0.00 0 | 0 | 0 | 0 | 177 | R 0 | cn_5001 0 | 0 | 0.00 0 | 0 | 0 [0 | 177 cn_5001 0 | 0 | 0.00 | W 0 | 0 | 01 177

14.3.208 PGXC_OS_RUN_INFO

PGXC_OS_RUN_INFO displays the OS running status of each node in the cluster. All columns except **node_name** are the same as those in the **PV_OS_RUN_INFO** view. Only the system administrator or the preset role **gs_role_read_all_stats** can access this view.

Table 14-278 PGXC OS RUN INFO field columns

Column	Туре	Description
node_name	Text	Node name

Column	Туре	Description
id	Integer	ID
Name	Text	Name of the OS running status
value	Numeric	Value of the OS status
comments	Text	Remarks of the OS status
cumulative	boolean	Whether the value of the OS status is cumulative

14.3.209 PGXC_OS_THREADS

PGXC_OS_THREADS displays thread status information under all normal nodes in the current cluster.

Table 14-279 PGXC_OS_THREADS columns

Column	Туре	Description
node_name	Text	Names of all normal nodes currently in the cluster.
pid	Bigint	Thread IDs currently running in the processes of all normal nodes in the cluster.
lwpid	Integer	Lightweight thread IDs corresponding to the PIDs.
thread_name	Text	Thread names corresponding to the PIDs.
creation_time	Timestamp with time zone	Creation time of the threads corresponding to the PIDs.

14.3.210 PGXC_POOLER_STATUS

PGXC_POOLER_STATUS displays the pooler cache connection status of each CN in the current cluster. This view can be queried only on CNs to display the connection cache information of the pooler module on all CNs. The **PGXC_POOLER_STATUS** view is supported only by clusters of version 8.3.0 or later.

Table 14-280 PGXC_POOLER_STATUS columns

Column	Туре	Description
coorname	Text	Name of the CN node.
database	Text	Database name.

Column	Туре	Description
user_name	Text	Username.
tid	Bigint	ID of the thread used for the connection to the CN.
node_oid	Bigint	OID of the node connected to.
node_name	Name	Name of the node connected to.
in_use	boolean	Whether the connection is currently in use. The options are: • t (true): The connection is in use. • f (false): The connection is not in use.
fdsock	Bigint	Peer socket.
remote_pid	Bigint	Peer thread ID.
session_params	Text	GUC session parameters issued by this connection.

14.3.211 PGXC_PREPARED_XACTS

PGXC_PREPARED_XACTS displays the two-phase transactions in the **prepared** phase.

Table 14-281 PGXC_PREPARED_XACTS columns

Column	Туре	Description
pgxc_prepared_xact	Text	Two-phase transactions in the prepared phase.

14.3.212 PGXC_REDO_STAT

PGXC_REDO_STAT displays statistics on redoing Xlogs of each node in the cluster. All columns except **node_name** are the same as those in the **PV_REDO_STAT** view. Only the system administrator or the preset role **gs_role_read_all_stats** can access this view.

Table 14-282 PGXC_REDO_STAT columns

Column	Туре	Description
node_name	Text	Node name
phywrts	Bigint	Number of physical writes

Column	Туре	Description
phyblkwrt	Bigint	Number of physical blocks written
writetim	Bigint	Time taken for physical writes
avgiotim	Bigint	Average time taken per write
lstiotim	Bigint	Time taken for the last write
miniotim	Bigint	Minimum time taken for a write
maxiowtm	Bigint	Maximum time taken for a write

14.3.213 PGXC_REL_IOSTAT

PGXC_REL_IOSTAT displays disk read/write statistics on each node in the cluster. This view is accessible only to users with system administrators rights.

Table 14-283 PGXC_REL_IOSTAT columns

Column	Туре	Description
node_name	Text	Node name
phyrds	Bigint	Number of disk reads
phywrts	Bigint	Number of disk writes
phyblkrd	Bigint	Number of read pages
phyblkwrt	Bigint	Number of written pages

14.3.214 PGXC_REPLICATION_SLOTS

PGXC_REPLICATION_SLOTS displays the replication information of DNs in the cluster. All columns except **node_name** are the same as those in the **PG_REPLICATION_SLOTS** view. This view is accessible only to users with system administrators rights.

Table 14-284 PGXC_REPLICATION_SLOTS columns

Column	Туре	Description
node_name	Text	Node name
slot_name	Text	Name of a replication node
plugin	Name	Name of the output plug-in of the logical replication slot
slot_type	Text	Type of a replication node

Column	Туре	Description
datoid	OID	OID of the database on the replication node
database	Name	Name of the database on the replication node
active	boolean	Whether the replication node is active
xmin	Xid	Transaction ID of the replication node
catalog_xmin	Text	ID of the earliest-decoded transaction corresponding to the logical replication slot
restart_lsn	Text	Xlog file information on the replication node
dummy_stand by	boolean	Whether the replication node is the dummy standby node

14.3.215 PGXC_RESPOOL_RUNTIME_INFO

PGXC_RESPOOL_RUNTIME_INFO displays the running information about all resource pool jobs on all CNs.

Table 14-285 PGXC_RESPOOL_RUNTIME_INFO columns

Column	Туре	Description
nodename	Name	CN name.
nodegroup	Name	Name of the logical cluster the resource pool belongs to. The default cluster is installation .
rpname	Name	Resource pool name.
ref_count	Int	Number of jobs that reference the resource pool. This count includes both controlled and uncontrolled jobs.
fast_run	Int	Number of jobs currently running in the resource pool's fast lane.
fast_wait	Int	Number of jobs currently queued in the resource pool's fast lane.
slow_run	Int	Number of jobs currently running in the resource pool's slow lane.
slow_wait	Int	Number of jobs currently queued in the resource pool's slow lane.

14.3.216 PGXC RESPOOL RESOURCE INFO

PGXC_RESPOOL_RESOURCE_INFO displays the real-time monitoring information about the resource pools on all instances.

- On a DN, it only displays the monitoring information of the logical cluster that the DN belongs to.
- Cluster 8.2.0 and later versions provide the negative memory feedback mechanism. The
 CCN can decrease the estimated memory usage of statements based on their actual
 memory usage on DNs, improving resource utilization by reducing overestimation.
 However, the estimated memory usage on CNs remains unchanged. If the CCN allows
 more jobs to run, the total estimated memory usage in the resource pool monitoring
 view may exceed the memory upper limit of the resource pool.
- Only the operators occupying large memory are under statement memory control. The
 memory, thread initialization costs, and expression costs of the operators with small
 memory usage are not controlled. So the value of used_mem of the resource pool may
 exceed the value of mem_limit to a limited extent.

Table 14-286 PGXC_RESPOOL_RESOURCE_INFO columns

Column	Туре	Description
nodename	Name	Instance name, including CNs and DNs.
nodegroup	Name	Name of the logical cluster of the resource pool. The default value is installation .
rpname	Name	Resource pool name.
cgroup	Name	Name of the Cgroup associated with the resource pool.
ref_count	Int	Number of jobs referenced by the resource pool. The number is counted regardless of whether the jobs are controlled by the resource pool. This parameter is valid only on CNs.
fast_run	Int	Number of running jobs in the fast lane of the resource pool. This parameter is valid only on CNs.
fast_wait	Int	Number of jobs queued in the fast lane of the resource pool. This parameter is valid only on CNs.
fast_limit	Int	Limit on the number of concurrent jobs in the fast lane in a resource pool. This parameter is valid only on CNs.
slow_run	Int	Number of running jobs in the slow lane of the resource pool. This parameter is valid only on CNs.

Column	Туре	Description
slow_wait	Int	Number of jobs queued in the slow lane of the resource pool. This parameter is valid only on CNs.
slow_limit	Int	Limit on the number of concurrent jobs in the slow lane in a resource pool. This parameter is valid only on CNs.
used_cpu	Double	 Average number of CPUs used by the resource pool in a 5s monitoring period. The value is accurate to two decimal places. On a DN, it indicates the number of CPUs used by the resource pool on the current DN. On a CN, it indicates the total CPU usage of resource pools on all DNs.
cpu_limit	Int	It indicates the upper limit of available CPUs for resource pools. If the CPU share is limited, this parameter indicates the available CPUs for GaussDB(DWS). If the CPU limit is specified, this parameter indicates the available CPUs for associated Cgroups. • On a DN, it indicates the upper limit of
		available CPUs for the resource pool on the current DN.
		 On a CN, it indicates the total upper limit of available CPUs for resource pools on all DNs.
used_mem	Int	Memory size used by the resource pool (unit: MB)
		On a DN, it indicates the memory usage of the resource pool on the current DN.
		On a CN, it indicates the total memory usage of resource pools on all DNs.
estimate_me m	Int	Estimated memory used by the jobs running in the resource pools on the current CN. This parameter is valid only on CNs.
mem_limit	Int	Upper limit of available memory for the resource pool (unit: MB).
		On a DN, it indicates the upper limit of available memory for the resource pool on the current DN.
		On a CN, it indicates the total upper limit of available memory for resource pools on all DNs.

Column	Туре	Description
read_kbytes	Bigint	Number of logical read bytes in the resource pool within a 5s monitoring period (unit: KB).
		On a DN, it indicates the number of logical read bytes in the resource pool on the current DN.
		On a CN, it indicates the total logical read bytes of resource pools on all DNs.
write_kbytes	Bigint	Number of logical write bytes in the resource pool within a 5s monitoring period (unit: KB).
		On a DN, it indicates the number of logical write bytes in the resource pool on the current DN.
		On a CN, it indicates the total logical write bytes of resource pools on all DNs.
read_counts	Bigint	Number of logical reads in the resource pool within a 5s monitoring period.
		On a DN, it indicates the number of logical reads in the resource pool on the current DN.
		On a CN, it indicates the total number of logical reads in resource pools on all DNs.
write_counts	Bigint	Number of logical writes in the resource pool within a 5s monitoring period.
		On a DN, it indicates the number of logical writes in the resource pool on the current DN.
		On a CN, it indicates the total number of logical writes in resource pools on all DNs.
read_speed	Double	Average logical read rate of a resource pool in a 5-second monitoring period, in KB/s
		On a DN, it indicates the logical read rate of the resource pool on the current DN.
		On a CN, it indicates the overall logical read rate of resource pools on all DNs.
write_speed	Double	Average logical write rate of a resource pool in a 5-second monitoring period, in KB/s
		On a DN, it indicates the logical write rate of the resource pool on the current DN.
		On a CN, it indicates the overall logical write rate of resource pools on all DNs.

Column	Туре	Description
send_speed	Double	Average network sending rate of a resource pool in a 5-second monitoring period, in KB/s
		On a DN, it indicates the network sending rate of the resource pool on the current DN.
		On a CN, it indicates the sum of the network sending rates of the resource pool on all DNs.
recv_speed	Double	Average network sending rate of a resource pool in a 5-second monitoring period, in KB/s
		On a DN, it indicates the network sending rate of the resource pool on the current DN.
		On a CN, it indicates the sum of the network sending rates of the resource pool on all DNs.

14.3.217 PGXC_RESPOOL_RESOURCE_HISTORY

PGXC_RESPOOL_RESOURCE_HISTORY is used to query historical monitoring information about resource pools on all instances.

Table 14-287 PGXC_RESPOOL_RESOURCE_HISTORY columns

Column	Туре	Description
nodename	Name	Instance name, including CNs and DNs
Timestamp	Timestamp	Time when resource pool monitoring information is persistently stored
nodegroup	Name	Name of the logical cluster the resource pool belongs to. The default cluster is installation .
rpname	Name	Resource pool name
cgroup	Name	Name of the Cgroup associated with the resource pool
ref_count	Int	Number of jobs referenced by the resource pool. The number is counted regardless of whether the jobs are controlled by the resource pool. This parameter is valid only on CNs.
fast_run	Int	Number of running jobs in the fast lane of the resource pool. This parameter is valid only on CNs.
fast_wait	Int	Number of jobs queued in the fast lane of the resource pool. This parameter is valid only on CNs.

Column	Туре	Description
fast_limit	Int	Limit on the number of concurrent jobs in the fast lane in a resource pool. This parameter is valid only on CNs.
slow_run	Int	Number of running jobs in the slow lane of the resource pool. This parameter is valid only on CNs.
slow_wait	Int	Number of jobs queued in the slow lane of the resource pool. This parameter is valid only on CNs.
slow_limit	Int	Limit on the number of concurrent jobs in the slow lane in a resource pool. This parameter is valid only on CNs.
used_cpu	Double	Average number of CPUs used by the resource pool in a 5s monitoring period. The value is accurate to two decimal places.
		On a DN, it indicates the number of CPUs used by the resource pool on the current DN.
		On a CN, it indicates the total CPU usage of resource pools on all DNs.
cpu_limit	Int	It indicates the upper limit of available CPUs for resource pools. If the CPU share is limited, this parameter indicates the available CPUs for GaussDB(DWS). If the CPU limit is specified, this parameter indicates the available CPUs for associated Cgroups.
		On a DN, it indicates the upper limit of available CPUs for the resource pool on the current DN.
		On a CN, it indicates the total upper limit of available CPUs for resource pools on all DNs.
used_mem	Int	Memory used by the resource pool, in MB
		On a DN, it indicates the memory usage of the resource pool on the current DN.
		On a CN, it indicates the total memory usage of resource pools on all DNs.
estimate_me m	Int	Estimated memory used by the jobs running in the resource pools on the current CN. This parameter is valid only on CNs.

Column	Туре	Description
mem_limit	Int	Upper limit of available memory for the resource pool, in MB
		On a DN, it indicates the upper limit of available memory for the resource pool on the current DN.
		On a CN, it indicates the total upper limit of available memory for resource pools on all DNs.
read_kbytes	Bigint	Number of logical read bytes in the resource pool within a 5s monitoring period, in KB
		On a DN, it indicates the number of logical read bytes in the resource pool on the current DN.
		On a CN, it indicates the total logical read bytes of resource pools on all DNs.
write_kbytes	Bigint	Number of logical write bytes in the resource pool within a 5s monitoring period, in KB
		 On a DN, it indicates the number of logical write bytes in the resource pool on the current DN.
		On a CN, it indicates the total logical write bytes of resource pools on all DNs.
read_counts	Bigint	Number of logical reads in the resource pool within a 5s monitoring period
		On a DN, it indicates the number of logical reads in the resource pool on the current DN.
		On a CN, it indicates the total number of logical reads in resource pools on all DNs.
write_counts	Bigint	Number of logical writes in the resource pool within a 5s monitoring period.
		On a DN, it indicates the number of logical writes in the resource pool on the current DN.
		On a CN, it indicates the total number of logical writes in resource pools on all DNs.
read_speed	Double	Average logical read rate of a resource pool in a 5-second monitoring period, in KB/s
		On a DN, it indicates the logical read rate of the resource pool on the current DN.
		On a CN, it indicates the overall logical read rate of resource pools on all DNs.

Column	Туре	Description
write_speed	Double	Average logical write rate of a resource pool in a 5-second monitoring period, in KB/s
		On a DN, it indicates the logical write rate of the resource pool on the current DN.
		On a CN, it indicates the overall logical write rate of resource pools on all DNs.
send_speed	Double	Average network sending rate of a resource pool in a 5-second monitoring period, in KB/s
		On a DN, it indicates the network sending rate of the resource pool on the current DN.
		On a CN, it indicates the sum of the network sending rates of the resource pool on all DNs.
recv_speed	Double	Average network receiving rate of a resource pool in a 5-second monitoring period, in KB/s
		On a DN, it indicates the network receiving rate of the resource pool on the current DN.
		On a CN, it indicates the sum of the network receiving rates of the resource pool on all DNs.

14.3.218 PGXC_ROW_TABLE_IO_STAT

PGXC_ROW_TABLE_IO_STAT provides I/O statistics of all row-store tables of the database on all CNs and DNs in the cluster. Except the **nodename** column of the name type added in front of each row, the names, types, and sequences of other columns are the same as those in the **GS_ROW_TABLE_IO_STAT** view. For details about the columns, see **Table 14-288**.

Table 14-288 GS_ROW_TABLE_IO_STAT columns

Column	Туре	Description
schemaname	Name	Namespace of a table
relname	Name	Name of a table
heap_read	Bigint	Number of blocks logically read in the heap
heap_hit	Bigint	Number of block hits in the heap
idx_read	Bigint	Number of blocks logically read in the index
idx_hit	Bigint	Number of block hits in the index
toast_read	Bigint	Number of blocks logically read in the TOAST table

Column	Туре	Description
toast_hit	Bigint	Number of block hits in the TOAST table
tidx_read	Bigint	Number of indexes logically read in the TOAST table
tidx_hit	Bigint	Number of index hits in the TOAST table

14.3.219 PGXC_RUNNING_XACTS

PGXC_RUNNING_XACTS displays information about running transactions on each node in the cluster. The content is the same as that displayed in **PG_RUNNING_XACTS**.

Table 14-289 PGXC_RUNNING_XACTS columns

Column	Туре	Description
handle	Integer	Handle corresponding to the transaction in GTM
gxid	Xid	Transaction ID
state	tinyint	Transaction status (3: prepared or 0: starting)
node	Text	Node name
xmin	Xid	Minimum transaction ID xmin on the node
vacuum	boolean	Whether the current transaction is lazy vacuum
timeline	Bigint	Number of database restarts
prepare_xid	Xid	Transaction ID in prepared state. If the status is not prepared , the value is 0 .
pid	Bigint	Thread ID corresponding to the transaction
next_xid	Xid	Transaction ID sent from a CN to a DN

14.3.220 PGXC_SETTINGS

PGXC_SETTINGS displays the database running status of each node in the cluster. All columns except **node_name** are the same as those in the **PG_SETTINGS** view. This view is accessible only to users with system administrators rights.

Table 14-290 PGXC_SETTINGS columns

Column	Туре	Description
node_name	Text	Node name
Name	Text	Parameter name
setting	Text	Current value of the parameter
unit	Text	Implicit unit of the parameter
category	Text	Logical group of the parameter
short_desc	Text	Brief description of the parameter
extra_desc	Text	Detailed description of the parameter
context	Text	Context of parameter values including internal, postmaster, sighup, backend, superuser, and user.
vartype	Text	Parameter type. It can be bool , enum , integer , real , or string .
source	Text	Method of assigning the parameter value
min_val	Text	Minimum value of the parameter If the parameter type is not numeric data, the value of this column is null.
max_val	Text	Maximum value of the parameter. If the parameter type is not numeric data, the value of this column is null.
enumvals	Text[]	Valid values of an enum-typed parameter. If the parameter type is not enum, the value of this column is null.
boot_val	Text	Default parameter value used upon the database startup
reset_val	Text	Default parameter value used upon the database reset
sourcefile	Text	Configuration file used to set parameter values. If parameter values are not configured using the configuration file, the value of this column is null.
sourceline	Integer	Row number of the configuration file for setting parameter values. If parameter values are not configured using the configuration file, the value of this column is null.

14.3.221 PGXC_SESSION_WLMSTAT

PGXC_SESSION_WLMSTAT displays load management information about ongoing jobs executed on each CN in the current cluster.

Table 14-291 PGXC_SESSION_WLMSTAT columns

Column	Туре	Description
nodename	Name	Node name.
datid	OID	OID of the database the backend is connected to.
datname	Name	Name of the database the backend is connected to.
threadid	Bigint	ID of the backend thread.
processid	Integer	PID of a backend thread
usesysid	OID	OID of the user who logged in to the backend
appname	Text	Name of the application that is connected to the backend
usename	Name	Name of the user logged in to the backend
priority	Bigint	Priority of Cgroup where the statement is located
attribute	Text	 Ordinary: default attribute of a statement before it is parsed by the database Simple: simple statements Complicated: complicated statements Internal: internal statement of the database
block_time	Bigint	Pending duration of the statements by now (unit: s)
elapsed_time	Bigint	Actual execution duration of the statements by now (unit: s)
total_cpu_time	Bigint	Total CPU usage duration of the statement on the DN in the last period (unit: s)
cpu_skew_perce nt	Integer	CPU usage inclination ratio of the statement on the DN in the last period
statement_mem	Integer	Estimated memory required for statement execution. This column is reserved.
active_points	Integer	Number of concurrently active points occupied by the statement in the resource pool
dop_value	Integer	DOP value obtained by the statement from the resource pool

Column	Туре	Description	
control_group	Text	Cgroup currently used by the statement	
status	Text	Status of a statement, including: • pending • running: The statement is being executed. • finished: The execution is finished normally. (If enqueue is set to StoredProc or Transaction, this state indicates that only	
		 some of the jobs in the statement have been executed. This state persists until the finish of this statement.) aborted: terminated unexpectedly active: normal status except for those above 	
		unknown: unknown status	
enqueue	Text	 Current queuing status of the statements, including: Global: global queuing. Respool: resource pool queuing. CentralQueue: queuing on the CCN Transaction: The statements are in the same transaction block. StoredProc: The statement is in a stored procedure. None: not in a queue Forced None: The transaction block statement or stored procedure statement is being forcibly executed because the statement waiting time exceeds the specified value. 	
resource_pool	Name	Current resource pool where the statements are located.	
query	Text	Text of this backend's most recent query If state is active , this column shows the executing query. In all other states, it shows the last query that was executed.	
isplana	Bool	In logical cluster mode, indicates whether a statement occupies the resources of other logical clusters. The default value is f , indicating that resources of other logical clusters are not occupied.	
node_group	Text	Logical cluster of the user running the statement	

Column	Туре	Description
lane	Text	Fast or slow lane for statement queries.
		fast: fast lane
		• slow: slow lane
		none: not controlled

14.3.222 PGXC_STAT_ACTIVITY

PGXC_STAT_ACTIVITY displays information about the query performed by the current user on all the CNs in the current cluster.

Table 14-292 PGXC_STAT_ACTIVITY columns

Column	Туре	Description
coorname	Text	Name of the CN in the current cluster
datid	OID	OID of the database that the user session connects to in the backend
datname	Name	Name of the database that the user session connects to in the backend
pid	Bigint	ID of the backend thread
lwtid	Integer	Lightweight thread ID of the backend thread
usesysid	OID	OID of the user logging in to the backend
usename	Name	Name of the user logging in to the backend
application_name	Text	Name of the application connected to the backend
client_addr	inet	IP address of the client connected to the backend. If this column is null , it indicates either that the client is connected via a Unix socket on the server machine or that this is an internal process such as autovacuum.
client_hostname	Text	Host name of the connected client, as reported by a reverse DNS lookup of client_addr. This column will only be non-null for IP connections, and only when log_hostname is enabled.

Column	Туре	Description
client_port	Integer	TCP port number that the client uses for communication with this backend, or -1 if a Unix socket is used
backend_start	Timestamp with time zone	Startup time of the backend process, that is, the time when the client connects to the server
xact_start	Timestamp with time zone	Time when the current transaction was started, or NULL if no transaction is active. If the current query is the first of its transaction, this column is equal to the query_start column.
query_start	Timestamp with time zone	Time when the currently active query was started, or time when the last query was started if state is not active
state_change	Timestamp with time zone	Time for the last status change
waiting	boolean	The value is t if the backend is waiting for a lock or node. Otherwise, the value is f .

Column	Туре	Description
enqueue	Text	Queuing status of a statement. Its value can be:
		• waiting in global queue: The statement is in the global concurrent queues.
		waiting in respool queue: The statement is queuing in the resource pool. The scenarios are as follows:
		 When dynamic load balancing is enabled, the number of simple jobs exceeds the upper limit (max_dop) of concurrent jobs on the fast lane.
		 When dynamic load balancing is disabled, the number of simple jobs exceeds the upper limit (max_dop) of concurrent jobs on the fast lane or the number of complex jobs exceeds the upper limit of concurrent jobs on the slow lane.
		waiting in ccn queue: The job is in the CCN queue, which may be global memory queuing, slow lane memory queuing, or concurrent queuing.
		Empty or no waiting queue : The statement is running.

Column	Туре	Description
state	Text	Overall state of the backend. Its value can be:
		active: The backend is executing a query.
		• idle: The backend is waiting for a new client command.
		• idle in transaction: The backend is in a transaction, but there is no statement being executed in the transaction.
		• idle in transaction (aborted): The backend is in a transaction, but there are statements failed in the transaction.
		• fastpath function call: The backend is executing a fast-path function.
		 disabled: This state is reported if track_activities is disabled in this backend.
		NOTE Only system administrators can view the session status of their accounts. The state information of other accounts is empty.
resource_pool	Name	Resource pool used by the user
stmt_type	Text	Type of a user statement
query_id	Bigint	ID of a query
query	Text	Text of this backend's most recent query If the state is active , this column shows the executing query. In all other states, it shows the last query that was executed.
connection_info	Text	A string in JSON format recording the driver type, driver version, driver deployment path, and process owner of the connected database (for details, see connection_info).

Example

Run the following command to view blocked query statements.

SELECT datname,usename,state,query FROM PGXC_STAT_ACTIVITY WHERE waiting = true;

Check the working status of the snapshot thread.

SELECT application_name,backend_start,state_change,state,query FROM PGXC_STAT_ACTIVITY WHERE application_name='WDRSnapshot';

View the running query statements.

View the number of session connections that have been used by postgres. 1 indicates the number of session connections that have been used by **postgres**.

```
SELECT COUNT(*) FROM PGXC_STAT_ACTIVITY WHERE DATNAME='postgres';
count
------
1
(1 row)
```

14.3.223 PGXC_STAT_BAD_BLOCK

PGXC_STAT_BAD_BLOCK displays statistics about page or CU verification failures after all nodes in a cluster are started.

Table 14-293 PGXC_STAT_BAD_BLOCK column

Column	Туре	Description
nodename	Text	Node name.
databaseid	Integer	Database OID.
tablespaceid	Integer	Tablespace OID.
relfilenode	Integer	File object ID.
forknum	Integer	File type.
error_count	Integer	Number of verification failures.
first_time	Timestamp with time zone	Time of the first occurrence.
last_time	Timestamp with time zone	Time of the latest occurrence.

14.3.224 PGXC_STAT_BGWRITER

PGXC_STAT_BGWRITER displays statistics on the background writer of each node in the cluster. All columns except **node_name** are the same as those in the **PG_STAT_BGWRITER** view. This view is accessible only to users with system administrators rights.

Table 14-294 PGXC_STAT_BGWRITER columns

Column	Туре	Description		
node_name	Text	Node name		
checkpoints_ti med	Bigint	Number of scheduled checkpoints that have been performed		
checkpoints_r eq	Bigint	Number of requested checkpoints that have been performed		
checkpoint_wr ite_time	Double precision	Time spent on writing files to the disk during checkpoints, in milliseconds		
checkpoint_sy nc_time	Double precision	Time spent on synchronizing data to the disk during checkpoints, in milliseconds		
buffers_check point	Bigint	Number of buffers written during checkpoints		
buffers_clean	Bigint	Number of buffers written by the background writer		
maxwritten_cl ean	Bigint	Number of times the background writer stopped a cleaning scan because it had written too many buffers		
buffers_backe nd	Bigint	Number of buffers written directly by a backend		
buffers_backe nd_fsync	Bigint	Number of times that a backend has to execute fsync		
buffers_alloc	Bigint	Number of buffers allocated		
stats_reset	Timestamp with time zone	Time at which these statistics were reset		

14.3.225 PGXC_STAT_DATABASE

PGXC_STAT_DATABASE displays the database status and statistics of each node in the cluster. All columns except **node_name** are the same as those in the **PG_STAT_DATABASE** view. This view is accessible only to users with system administrators rights.

Table 14-295 PGXC_STAT_DATABASE columns

Column	Туре	Description	
node_name	Text	Node name	
datid	OID	Database OID	

Column	Туре	Description		
datname	Name	Database name		
numbackends	Integer	Number of backends currently connected to this database on the current node. This is the only column in this view that reflects the current state value. All columns return the accumulated value since the last reset.		
xact_commit	Bigint	Number of transactions in this database that have been committed on the current node		
xact_rollback	Bigint	Number of transactions in this database that have been rolled back on the current node		
blks_read	Bigint	Number of disk blocks read in this database on the current node		
blks_hit	Bigint	Number of disk blocks found in the buffer cache on the current node, that is, the number of blocks hit in the cache. (This only includes hits in the GaussDB(DWS) buffer cache, not in the file system cache.)		
tup_returned	Bigint	Number of rows returned by queries in this database on the current node		
tup_fetched	Bigint	Number of rows fetched by queries in this database on the current node		
tup_inserted	Bigint	Number of rows inserted in this database on the current node		
tup_updated	Bigint	Number of rows updated in this database on the current node		
tup_deleted	Bigint	Number of rows deleted from this database on the current node		
conflicts	Bigint	Number of queries canceled due to database recovery conflicts on the current node (conflicts occurring only on the standby server). For details, see PG_STAT_DATABASE_CONFLICTS.		
temp_files	Bigint	Number of temporary files created by this database on the current node. All temporary files are counted, regardless of why the temporary file was created (for example, sorting or hashing), and regardless of the log_temp_files setting.		

Column	Туре	Description	
temp_bytes	Bigint	Size of temporary files written to this database on the current node. All temporary files are counted, regardless of why the temporary file was created, and regardless of the log_temp_files setting.	
deadlocks	Bigint	Number of deadlocks in this database on the current node	
blk_read_time	Double precision	Time spent reading data file blocks by backends in this database on the current node, in milliseconds	
blk_write_tim e	Double precision	Time spent writing into data file blocks by backends in this database on the current node, in milliseconds	
stats_reset	Timestamp with time zone	Time when the database statistics are reset on the current node	

14.3.226 PGXC_STAT_OBJECT

PGXC_STAT_OBJECT displays statistics and autovacuum efficiency information about tables of all instances in a cluster. This system view is supported only by clusters of version 8.2.1 or later.

Table 14-296 PGXC_STAT_OBJECT columns

Column	Туре	Referenc e	Description
nodename	Name	-	Node name
datname	Name	-	Name of the database where the table is located.
relnamespace	Name	-	Name of the schema where the table is located.
relname	Name	-	Table name.
partname	Name	-	Partition name of the partitioned table
databaseid	OID	PG_DATA BASE.oid	Database OID.
relid	OID	PG_CLAS S.oid	Table OID. It is the OID of the primary table for a partitioned table.

Column	Туре	Referenc e	Description
partid	OID	PG_PARTI TION .oid	Partition OID. If the table is not partitioned, the value is 0 .
numscans	Bigint	-	Number of times that sequential scans are started.
tuples_returne d	Bigint	-	Number of visible tuples fetched by sequential scans.
tuples_fetche d	Bigint	-	Number of visible tuples fetched.
tuples_inserte d	Bigint	-	Number of inserted records.
tuples_update d	Bigint	-	Number of updated records.
tuples_delete d	Bigint	-	Number of deleted records.
tuples_hot_up dated	Bigint	-	Number of HOT updates.
n_live_tuples	Bigint	-	Number of visible tuples.
last_autovacu um_begin_n_ dead_tuple	Bigint	-	Number of tuples deleted before Autovacuum is executed.
n_dead_tuples	Bigint	-	Number of tuples deleted after Autovacuum is successful.
changes_since _analyze	Bigint	-	Last data modification time after Analyze.
blocks_fetche d	Bigint	-	Number of selected pages.
blocks_hit	Bigint	-	Number of scanned pages.
cu_mem_hit	Bigint	-	Number of CU memory hits.
cu_hdd_sync	Bigint	-	Times that CUs are synchronously read from disks.
cu_hdd_asyn	Bigint	-	Times that CUs are asynchronously read from disks.
data_changed _timestamp	Timestamp with time zone	-	Last data modification time.

Column	Туре	Referenc e	Description
data_access_ti mestamp	Timestamp with time zone	-	Last access time of a table.
analyze_times tamp	Timestamp with time zone	-	Last Analyze time.
analyze_count	Bigint	-	Total number of Analyze times.
autovac_analy ze_timestamp	Timestamp with time zone	-	Last Autoanalyze time.
autovac_analy ze_count	Bigint	-	Total number of Autoanalyze times.
vacuum_times tamp	Timestamp with time zone	-	Time of the latest Vacuum.
vacuum_coun t	Bigint	-	Total number of Vacuum times.
autovac_vacu um_timestam p	Timestamp with time zone	-	Last Autovacuum time.
autovac_vacu um_count	Bigint	-	Total number of Autovacuum times.
autovacuum_s uccess_count	Bigint	-	Total number of successful Autovacuum operations.
last_autovacu um_time_cost	Bigint	-	Time spent on the latest successful Autovacuum, in microseconds.
avg_autovacu um_time_cost	Bigint	-	Average execution time of successful Autovacuum operations. Unit: μs.
last_autovacu um_failed_co unt	Bigint	-	Total number of autovacuum failures since the last successful Autovacuum.
last_autovacu um_trigger	Smallint	-	Triggering mode of the latest autovacuum, which helps maintenance personnel determine the Vacuum status.

Column	Туре	Referenc e	Description
last_autovacu um_oldestxmi n	Bigint	-	oldestxmin after the latest successful Autovacuum execution. If the table-level oldestxmin feature is enabled, this field records the value of oldestxmin used by the latest (AUTO)VACUUM of the table.
last_autovacu um_scan_pag es	Bigint	-	Number of pages last scanned by autovacuum (only for row-store tables).
last_autovacu um_dirty_pag es	Bigint	-	Number of pages last modified by Autovacuum (only for row-store tables).
last_autovacu um_clear_dea dtuples	Bigint	-	Number of dead tuples last cleared by Autovacuum (only for row-store tables)
sum_autovacu um_scan_pag es	Bigint	-	Total number of pages scanned by Autovacuum since database initialization (only for row-store tables).
sum_autovacu um_dirty_pag es	Bigint	-	Number of pages modified by Autovacuum since database initialization (only for row-store tables).
sum_autovacu um_clear_dea dtuples	Bigint	-	Total number of dead tuples cleared by Autovacuum since database initialization (only for row-store tables).
last_autovacu um_begin_cu_ size	Bigint	-	Size of the CU file before the latest Autovacuum operation (only for column-store tables).
last_autovacu um_cu_size	Bigint	-	Size of the CU file after the latest Autovacuum (only for column- store tables).
last_autovacu um_rewrite_si ze	Bigint	-	Size of the column-store file last rewritten by autovacuum (only for column-store tables).
last_autovacu um_clear_size	Bigint	-	Size of the column-store file last cleared by Autovacuum (only for column-store tables).

Column	Туре	Referenc e	Description
last_autovacu um_clear_cbtr ee_tuples	Bigint	-	Number of cbtree tuples last cleared by Autovacuum (only for column-store tables).
sum_autovacu um_rewrite_si ze	Bigint	-	Total size of column-store files rewritten by Autovacuum since database initialization (only for column-store tables).
sum_autovacu um_clear_size	Bigint	-	Total size of column-store files cleared by Autovacuum since database initialization (only for column-store tables).
sum_autovacu um_clear_cbtr ee_tuples	Bigint	-	Total number of cbtree tuples cleared by Autovacuum since database initialization (only for column-store tables).
last_autovacu um_csn	Bigint	-	If the table-level oldestxmin feature is enabled, this field records the CSN value corresponding to the latest oldestxmin value used by the table (AUTO)VACUUM.
last_reference _timestamp	Timestamp with time zone	-	Last access time of a table. (This field is supported only by cluster versions 8.3.0 and later.)
			This parameter corresponds to the latest time between data_changed_time_stamp (last modification time) and data_access_timestamp (last access time) in PG_STAT_OBJECT.
extra1	Bigint	-	Reserved field 1.
extra2	Bigint	-	Reserved field 2.
extra3	Bigint	-	Reserved field 3.
extra4	Bigint	-	Reserved field 4.

14.3.227 PGXC_STAT_REPLICATION

PGXC_STAT_REPLICATION displays the log synchronization status of each node in the cluster. All columns except **node_name** are the same as those in the **PG_STAT_REPLICATION** view. This view is accessible only to users with system administrators rights.

Table 14-297 PGXC_STAT_REPLICATION columns

Column	Туре	Description
node_name	Text	Node name
pid	Bigint	PID of the thread
usesysid	OID	User system ID
usename	Name	Username
application_n ame	Text	Program name
client_addr	inet	Client address
client_hostna me	Text	Client name
client_port	Integer	Client port number
backend_start	Timestamp with time zone	Program start time
state	Text	Log replication state (catch-up or consistent streaming)
sender_sent_l ocation	Text	Location where the sender sends logs
receiver_write _location	Text	Location where the receiver writes logs
receiver_flush _location	Text	Location where the receiver flushes logs
receiver_repla y_location	Text	Location where the receiver replays logs
sync_priority	Integer	Priority of synchronous duplication (0 indicates asynchronization)
sync_state	Text	Synchronization state (asynchronous duplication, synchronous duplication, or potential synchronization)

14.3.228 PGXC_STAT_TABLE_DIRTY

PGXC_STAT_TABLE_DIRTY displays statistics about all the tables on all the CNs and DNs in the current cluster, and the dirty page rate of tables on a single CN or DN. This view is supported only by clusters of version 8.1.3 or later.

□ NOTE

The statistics of this view depend on the **ANALYZE** operation. To obtain the most accurate information, perform the **ANALYZE** operation on the table first.

Table 14-298 PGXC_STAT_TABLE_DIRTY columns

Column	Туре	Description
nodename	Text	Node name
schema	Name	Schema name of the table
tablename	Name	Table name
partname	Name	Partition name of the partitioned table
last_vacuum	timestampwith time zone	Time of the last manual VACUUM
last_autovacuum	timestampwith time zone	Time of the last AUTOVACUUM
last_analyze	timestampwith time zone	Time of the last manual ANALYZE
last_autoanalyze	timestampwith time zone	Time of the last AUTOANALYZE
vacuum_count	Bigint	Number of times VACUUM operations
autovacuum_cou nt	Bigint	Number of AUTOVACUUM operations
analyze_count	Bigint	Number of ANALYZE operations
autoanalyze_cou nt	Bigint	Number of AUTOANALYZE_COUNT operations
n_tup_ins	Bigint	Number of rows inserted
n_tup_upd	Bigint	Number of rows updated
n_tup_del	Bigint	Number of rows deleted
n_tup_hot_upd	Bigint	Number of rows with HOT updates
n_tup_change	Bigint	Number of changed rows after ANALYZE
n_live_tup	Bigint	Estimated number of live rows
n_dead_tup	Bigint	Estimated number of dead rows

Column	Туре	Description
dirty_rate	Bigint	Dirty page rate of a single CN or DN
last_data_chang ed	timestampwith time zone	Time when a table was last modified
		This column is left empty in cluster 8.2.0 and later versions. For details, see PG_STAT_ALL_TABLES.

Suggestion

- Before running **VACUUM FULL** on a system catalog with a high dirty page rate, ensure that no user is performing operations on it.
- You are advised to run VACUUM FULL to tables (excluding system catalogs) whose dirty page rate exceeds 80% or run it based on service scenarios.

Scenarios

1. Query the overall dirty page rate of all the user tables in a database.

```
select
  t1.schema,
  t1.tablename,
  t1.total_ins,
  t1.total_upd,
  t1.total del,
  t1. total_tup_hot_upd,
  t1.total_change,
  t1.total_live,
  t1.total_dead,
  t1.total_dirty_rate,
  t1.max_dirty,
  t2.max node,
  t1.min_dirty,
  t2.min_node
from
  (select
     a.schema,
     a.tablename,
     sum(a.n_tup_ins) as total_ins,
     sum(a.n_tup_upd) as total_upd,
     sum(a.n_tup_del) as total_del,
     sum(a.n_tup_hot_upd) as total_tup_hot_upd,
     sum(a.n_tup_change) as total_change,
     sum(a.n_live_tup) as total_live,
     sum(a.n_dead_tup) as total_dead,
     Round((total_dead / (total_dead + total_live + 0.0001) * 100),2) AS total_dirty_rate,
     max(a.dirty_rate) as max_dirty,
     min(a.dirty_rate) as min_dirty
  from pg_catalog.pgxc_stat_table_dirty a where a.partname is null and a.schema not in
('pg_toast','cstore','gs_logical_cluster','sys','dbms_om','information_schema','pg_catalog','dbms_output','
dbms_random','utl_raw','utl_raw dbms_sql','dbms_lob') group by a.tablename, a.schema
  ) t1,
  (select distinct
  tablename, schema,
  first_value(nodename) over(partition by tablename, schema order by dirty_rate) as min_node,
  first_value(nodename) over(partition by tablename, schema order by dirty_rate desc) as max_node
  from (select * from pg_catalog.pgxc_stat_table_dirty)) t2
where t1.tablename = t2.tablename and t1.schema = t2.schema;
```

2. Query the overall dirty page rate of all the tables (user tables and system catalogs) in a database.

```
select
  t1.schema,
  t1.tablename,
  t1.total_ins,
  t1.total_upd,
  t1.total_del,
  t1. total_tup_hot_upd,
  t1.total change,
  t1.total_live,
  t1.total_dead,
  t1.total_dirty_rate,
  t1.max_dirty,
  t2.max node,
  t1.min_dirty,
  t2.min_node
from
  (select
     a.schema,
     a.tablename,
     sum(a.n_tup_ins) as total_ins,
     sum(a.n_tup_upd) as total_upd,
     sum(a.n_tup_del) as total_del,
     sum(a.n_tup_hot_upd) as total_tup_hot_upd,
     sum(a.n_tup_change) as total_change,
     sum(a.n live tup) as total live,
     sum(a.n_dead_tup) as total_dead,
     Round((total_dead / (total_dead + total_live + 0.0001) * 100),2) AS total_dirty_rate,
     max(a.dirty_rate) as max_dirty,
     min(a.dirty_rate) as min_dirty
  from pg_catalog.pgxc_stat_table_dirty a where a.partname is null group by a.tablename, a.schema
  ) t1,
  (select distinct
  tablename, schema,
  first_value(nodename) over(partition by tablename, schema order by dirty_rate) as min_node,
  first_value(nodename) over(partition by tablename, schema order by dirty_rate desc) as max_node
  from (select * from pg_catalog.pgxc_stat_table_dirty)) t2
where t1.tablename = t2.tablename and t1.schema = t2.schema;
```

3. Query all system catalogs in a database.

select * from pgxc_stat_table_dirty where schema in ('pg_toast','cstore','gs_logical_cluster','sys','dbms_om','information_schema','pg_catalog','dbms_output','dbms_random','utl_raw','utl_raw dbms_sql','dbms_lob');

14.3.229 PGXC_STAT_WAL

PGXC_STAT_WAL displays the WAL logs and data page traffic information of the current query. This view is supported only by clusters 8.2.0 and later versions.

Table 14-299 PGXC_STAT_WAL columns

Column	Туре	Description
query_id	Bigint	ID of the current query
query_start	Timesta mp	Start time of the query
global_wal	Bigint	Total number of WAL logs generated by the current query in the cluster, in bytes
global_avg_wal_ speed	Bigint	Average rate of WAL log generation for the current query in the cluster, in byte/s

Column	Туре	Description
global_datapage	Bigint	Total size of data pages generated by the current query in the cluster, in bytes
global_avg_data page_speed	Bigint	Average rate of data page generation for the current query in the cluster, in byte/s
min_wal_node	Text	Name of the instance group that generates the smallest volume of WAL logs in the current query
min_wal	Bigint	Minimum WAL logs generated by a node, in bytes
max_wal_node	Text	Name of the instance group that generates the largest volume of WAL logs in the current query
max_wal	Bigint	Maximum WAL logs generated by a node, in bytes
min_datapage_n ode	Text	Name of the instance group that generates the smallest volume of data pages in the current query
min_data_page	Bigint	Minimum data pages generated by a node, in bytes
max_datapage_n ode	Text	Name of the instance group that generates the largest volume of data pages in the current query
max_data_page	Bigint	Maximum data pages generated by a node, in bytes
avg_wal_per_no de	Bigint	Average WAL logs generated by each node, in bytes
avg_datapage_p er_node	Bigint	Average data pages generated by each node, in bytes
query	Text	Statement that is being executed

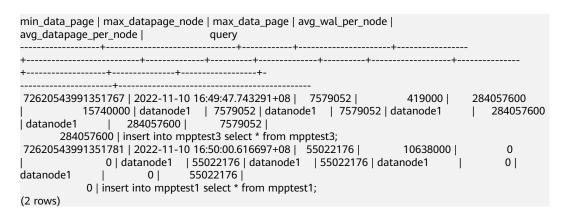
MOTE

When row-store data is imported in batches without indexes, the Xlogs related to logical new pages are generated during data page copy. If the volume of Xlogs is greater than the default value, flow control will be triggered.

Examples

Query the statements that are being executed in the cluster, the total volumes of WAL logs and data pages generated by these statements, their average generation rates, and their distribution on DNs.

```
SELECT * FROM PGXC_STAT_WAL;
query_id | query_start | global_wal | global_avg_wal_speed | global_datapage |
global_avg_datapage_speed | min_wal_node | min_wal_node | max_wal | min_datapage_node |
```



14.3.230 PGXC_SQL_COUNT

PGXC_SQL_COUNT displays the node-level and user-level statistics for the SQL statements of SELECT, INSERT, UPDATE, DELETE, and MERGE INTO and DDL, DML, and DCL statements of each CN in a cluster in real time, identifies query types with heavy load, and measures the capability of a cluster or a node to perform a specific type of query. You can calculate QPS based on the quantities and response time of the preceding types of SQL statements at certain time points. For example, USER1 SELECT is counted as X1 at T1 and as X2 at T2. The SELECT QPS of the user can be calculated as follows: (X2 – X1)/(T2 – T1). In this way, the system can draw cluster-user-level QPS curve graphs and determine cluster throughput, monitoring changes in the service load of each user. If there are drastic changes, the system can locate the specific statement type (such as SELECT, INSERT, UPDATE, DELETE, and MERGE INTO). You can also observe QPS curves to determine the time points when problems occur and then locate the problems using other tools. The curves provide a basis for optimizing cluster performance and locating problems.

Columns in the **PGXC_SQL_COUNT** view are the same as those in the **GS_SQL_COUNT** view. For details, see **Table 14-142**.

□ NOTE

If a **MERGE INTO** statement can be pushed down and a DN receives it, the statement will be counted on the DN and the value of the **mergeinto_count** column will increment by 1. If the pushdown is not allowed, the DN will receive an **UPDATE** or **INSERT** statement. In this case, the **update_count** or **insert_count** column will increment by 1.

14.3.231 PGXC_TABLE_CHANGE_STAT

PGXC_TABLE_CHANGE_STAT displays the changes of all tables of the database on all CNs in the cluster. Except the **nodename** column of the name type added in front of each row, the names, types, and sequences of other columns are the same as those in the **GS_TABLE_CHANGE_STAT** view.

Table 14-300 PGXC_TABLE_CHANGE_STAT columns

Column	Туре	Description
nodename	Name	Node name
schemaname	Name	Table namespace

Column	Туре	Description
relname	Name	Table name
last_vacuum	Timestamp with time zone	Time when the last VACUUM operation is performed manually
vacuum_count	Bigint	Number of times of manually performing the VACUUM operation
last_autovacuum	Timestamp with time zone	Time when the last VACUUM operation is performed automatically
autovacuum_cou nt	Bigint	Number of times of automatically performing the VACUUM operation
last_analyze	Timestamp with time zone	Time when the ANALYZE operation is performed (both manually and automatically)
analyze_count	Bigint	Number of times of performing the ANALYZE operation (both manually and automatically)
last_autoanalyze	Timestamp with time zone	Time when the last ANALYZE operation is performed automatically
autoanalyze_cou nt	Bigint	Number of times of automatically performing the ANALYZE operation
last_change	Bigint	Time when the last modification (INSERT, UPDATE, or DELETE) is performed

14.3.232 PGXC_TABLE_STAT

PGXC_TABLE_STAT provides statistics of all tables of the database on all CNs and DNs in the cluster. Except the **nodename** column of the name type added in front of each row, the names, types, and sequences of other columns are the same as those in the **GS_TABLE_STAT** view.

Table 14-301 PGXC_TABLE_STAT columns

Column	Туре	Description
nodename	Name	Node name
schemaname	Name	Table namespace
relname	Name	Table name

Column	Туре	Description
seq_scan	Bigint	Number of sequential scans. Only row-store tables are counted. For a partitioned table, the sum of the number of scans of each partition is displayed.
seq_tuple_rea d	Bigint	Number of rows scanned in sequence. Only row-store tables are counted.
index_scan	Bigint	Number of index scans. Only row-store tables are counted.
index_tuple_re ad	Bigint	Number of rows scanned by the index. Only row-store tables are counted.
tuple_inserted	Bigint	Number of rows inserted
tuple_updated	Bigint	Number of rows updated
tuple_deleted	Bigint	Number of rows deleted
tuple_hot_upd ated	Bigint	Number of rows with HOT updates.
live_tuples	Bigint	Number of live tuples. Query the view on the CN. If ANALYZE is executed, the total number of live tuples in the table is displayed. Otherwise, 0 is displayed. This indicator applies only to row-store tables.
dead_tuples	Bigint	Number of dead tuples. Query the view on the CN. If ANALYZE is executed, the total number of dead tuples in the table is displayed. Otherwise, 0 is displayed. This indicator applies only to row-store tables.

14.3.233 PGXC_THREAD_WAIT_STATUS

PGXC_THREAD_WAIT_STATUS displays all the call layer hierarchy relationship between threads of the SQL statements on all the nodes in a cluster, and the waiting status of the block for each thread, so that you can easily locate the causes of process response failures and similar phenomena.

The definitions of PGXC_THREAD_WAIT_STATUS view and PG_THREAD_WAIT_STATUS view are the same, because the essence of the PGXC_THREAD_WAIT_STATUS view is the query summary result of the PG_THREAD_WAIT_STATUS view on each node in the cluster.

Table 14-302 PGXC_THREAD_WAIT_STATUS columns

Column	Туре	Description
node_name	Text	Current node name

Column	Туре	Description
db_name	Text	Database name
thread_name	Text	Thread name
query_id	Bigint	Query ID. It is equivalent to debug_query_id .
tid	Bigint	Thread ID of the current thread
lwtid	Integer	Lightweight thread ID of the current thread
ptid	Integer	Parent thread of the streaming thread
tlevel	Integer	Level of the streaming thread
smpid	Integer	Concurrent thread ID
wait_status	Text	Waiting status of the current thread. For details about the waiting status, see Table 14-234.
wait_event	Text	If wait_status is acquire lock, acquire lwlock, or wait io, this column describes the lock, lightweight lock, and I/O information, respectively. If wait_status is not any of the three values, this column is empty.

Example:

Assume you run a statement on coordinator1, and no response is returned after a long period of time. In this case, establish another connection to coordinator1 to check the thread status on it.

Furthermore, you can view the statement working status on each node in the entire cluster. In the following example, no DNs have threads blocked, and there is a huge amount of data to be read, causing slow execution.

```
synchronize quit |
datanode2 | gaussdb | coordinator1 | 20971544 | 140632081299216 | 22975 | 22736 |
synchronize quit |
datanode3 | gaussdb | coordinator1 | 20971544 | 140323627988752 | 22737 | 0 | 0 | wait
node: datanode3 |
datanode3 | gaussdb | coordinator1 | 20971544 | 140323523131152 | 22976 | 22737 | 5 | 0 | net
flush data
datanode3 | gaussdb | coordinator1 | 20971544 | 140323548296976 | 22978 | 22737 |
                                                                                   5 | 1 | net
flush data
datanode4 | gaussdb | coordinator1 | 20971544 | 140103024375568 | 22738 |
node: datanode3
datanode4 | gaussdb | coordinator1 | 20971544 | 140102919517968 | 22979 | 22738 |
synchronize quit
datanode4 | gaussdb | coordinator1 | 20971544 | 140102969849616 | 22980 | 22738 |
                                                                                   5 | 1 |
synchronize quit
coordinator1 | gaussdb | gsql
                               | 20971544 | 140274089064208 | 22579 |
                                                                         | 0|
                                                                                   0 | wait node:
datanode4 |
(13 rows)
```

14.3.234 PGXC_TOTAL_MEMORY_DETAIL

PGXC_TOTAL_MEMORY_DETAIL displays the memory usage in the cluster. Only the system administrator or the preset role **gs_role_read_all_stats** can access this view.

Table 14-303 PGXC_TOTAL_MEMORY_DETAIL columns

Column	Туре	Description
nodename	Text	Node name

Column	Туре	Description
memorytype	Text	Memory name, which can be set to any of the following values:
		 max_process_memory: memory used by a GaussDB(DWS) cluster instance
		 process_used_memory: memory used by a GaussDB(DWS) process
		• max_dynamic_memory: maximum dynamic memory
		• dynamic_used_memory : used dynamic memory
		 dynamic_peak_memory: dynamic peak value of the memory
		 dynamic_used_shrctx: maximum dynamic shared memory context
		 dynamic_peak_shrctx: dynamic peak value of the shared memory context
		 max_shared_memory: maximum shared memory
		shared_used_memory: used shared memory
		 max_cstore_memory: maximum memory allowed for column store
		• cstore_used_memory: memory used for column store
		 max_sctpcomm_memory: maximum memory allowed for the communication library
		 sctpcomm_used_memory: memory used for the communication library
		 sctpcomm_peak_memory: memory peak of the communication library
		 other_used_memory: other used memory
		• gpu_max_dynamic_memory: maximum GPU memory
		 gpu_dynamic_used_memory: sum of the available GPU memory and temporary GPU memory
		 gpu_dynamic_peak_memory: maximum memory used for GPU
		 pooler_conn_memory: memory used for pooler connections
		pooler_freeconn_memory: memory used for idle pooler connections

Column	Туре	Description
		storage_compress_memory: memory used for column-store compression and decompression
		 udf_reserved_memory: memory reserved for the UDF Worker process
		 mmap_used_memory: memory used for mmap
memorymbyte s	Integer	Size of the used memory (MB)

14.3.235 PGXC_TOTAL_SCHEMA_INFO

PGXC_TOTAL_SCHEMA_INFO displays the schema space information of all instances in the cluster, providing visibility into the schema space usage of each instance. This view can be queried only on CNs.

Table 14-304 PGXC_TOTAL_SCHEMA_INFO columns

Column	Туре	Description
schemaname	Text	Schema name.
schemaid	OID	Schema OID.
databasename	Text	Database name.
databaseid	OID	Database OID.
nodename	Text	Instance name.
nodegroup	Text	Node group name.
usedspace	Bigint	Used space size.
permspace	Bigint	Space upper limit.

14.3.236 PGXC_TOTAL_SCHEMA_INFO_ANALYZE

PGXC_TOTAL_SCHEMA_INFO_ANALYZE displays the overall schema space information of the cluster, including the total cluster space, average space of instances, skew ratio, maximum space of a single instance, minimum space of a single instance, and names of the instances with the maximum space and minimum space. It provides visibility into the schema space usage of the entire cluster. This view can be queried only on CNs.

Table 14-305 PGXC_TOTAL_SCHEMA_INFO_ANALYZE columns

Column	Туре	Description
schemaname	Text	Schema name.
databasename	Text	Database name.
nodegroup	Text	Node group name.
total_value	Bigint	Total cluster space in this schema.
avg_value	Bigint	Average space per instance in this schema.
skew_percent	Integer	Skew ratio.
extend_info	Text	The extended information includes the maximum and minimum space values for a single instance, as well as the names of the instances with the maximum and minimum space values.

14.3.237 PGXC_TOTAL_USER_RESOURCE_INFO

The **PGXC_TOTAL_USER_RESOURCE_INFO** view displays real-time resource consumption information of users on all instances. This view is supported only by clusters of version 8.2.0 or later.

Table 14-306 PGXC_TOTAL_USER_RESOURCE_INFO columns

Column	Туре	Description
nodename	Name	Instance name, including CNs and DNs.
username	Name	Username
used_memory	Integer	Used memory (unit: MB)
		On a DN, it indicates a user's memory usage on the current DN.
		On a CN, it indicates a user's total memory usage on all DNs.

Column	Туре	Description
total_memory	Integer	Available memory (unit: MB). 0 indicates that the available memory is not limited and depends on the maximum memory available in the database.
		On a DN, it indicates the memory available to a user on the current DN.
		On a CN, it indicates the total memory available to a user on all DNs.
used_cpu	Double precision	Number of CPU cores in use. Only the CPU usage of complex jobs in the non-default resource pool is collected, and the value is the CPU usage of the related cgroup. On a DN, it indicates a user's CPU core usage
		on the current DN. On a CN, it indicates a user's total CPU core
		usage on all DNs.
total_cpu	Integer	Total number of CPU cores of the Cgroups associated with a user.
		On a DN, it indicates the CPU cores available to a user on the current DN.
		On a CN, it indicates the total CPU cores available to a user on all DNs.
used_space	Bigint	Used permanent table storage space (unit: KB) On a DN, it indicates the size of the permanent table storage space used by a user on the current DN.
		On a CN, it indicates the total size of the permanent table storage space used by a user on all DNs.
total_space	Bigint	Available storage space (unit: KB)1 indicates that the storage space is not limited.
		On a DN, it indicates the size of the permanent table storage space available to a user on the current DN.
		On a CN, it indicates the total size of the permanent table storage space available to a user on all DNs.

Column	Туре	Description
used_temp_sp	Bigint	Used temporary table storage space (unit: KB)
ace		On a DN, it indicates the size of the temporary table storage space used by a user on the current DN.
		On a CN, it indicates the total size of the temporary table storage space used by a user on all DNs.
total_temp_sp ace	Bigint	Available temporary table storage space (unit: KB)1 indicates that the storage space is not limited.
		On a DN, it indicates the size of the temporary table storage space available to a user on the current DN.
		On a CN, it indicates the total size of the temporary table storage space available to a user on all DNs.
used_spill_spa ce	Bigint	Size of space used for operator spill to disk, in KB.
		On a DN, it indicates the space used by a user to spill operators to disk on the current DN.
		On a CN, it indicates the total space used by a user's operators spilled to disk on all DNs.
total_spill_spa ce	Bigint	Size of space available for operator spill to disk, in KB. The value -1 indicates that the space is not limited.
		On a DN, it indicates the space available for a user to spill operators to disk on the current DN.
		On a CN, it indicates the total space available for a user to spill operators to disk on all DNs.
read_kbytes	Bigint	On a CN, it indicates the total number of bytes logically read by a user on all DNs in the last 5 seconds, in KB.
		On a DN, it indicates the total number of bytes logically read by a user from the instance startup time to the current time, in KB.
write_kbytes	Bigint	On a CN, it indicates the total number of bytes logically written by a user on all DNs in the last 5 seconds, in KB.
		On a DN, it indicates the total number of bytes logically written by a user from the instance startup time to the current time, in KB.

Column	Туре	Description
read_counts	Bigint	On a CN, it indicates the total number of logical reads performed by a user on all DNs in the last 5 seconds.
		On a DN, it indicates the total number of logical reads performed by a user from the instance startup time to the current time.
write_counts	Bigint	On a CN, it indicates the total number of logical writes performed by a user on all DNs in the last 5 seconds.
		On a DN, it indicates the total number of logical writes performed by a user from the instance startup time to the current time.
read_speed	Double precision	On a CN, it indicates the average logical read rate of a user on a single DN in the last 5 seconds, in KB/s.
		On a DN, it indicates the average logical read rate of a user on the DN in the last 5 seconds, in KB/s.
write_speed	Double precision	On a CN, it indicates the average logical write rate of a user on a single DN in the last 5 seconds, in KB/s.
		On a DN, it indicates the average logical write rate of a user on the DN in the last 5 seconds, in KB/s.
send_speed	Double precision	On a CN, it indicates the sum of the average network sending rates of a user on all DNs in the last 5 seconds, in KB/s.
		On a DN, it indicates the average network sending rate of a user on the DN in the last 5 seconds, in KB/s.
recv_speed	Double precision	On a CN, it indicates the sum of the average network receiving rates of a user on all DNs in the last 5 seconds, in KB/s.
		On a DN, it indicates the average network receiving rate of a user on the DN in the last 5 seconds, in KB/s.

14.3.238 PGXC_USER_TRANSACTION

PGXC_USER_TRANSACTION provides transaction information about users on all CNs. It is accessible only to users with system administrator rights. This view is valid only if the real-time resource monitoring function is enabled, that is, if **enable_resource_track** is **on**.

Table 14-307 PGXC_USER_TRANSACTION columns

Column	Туре	Description
node_name	Name	Node name.
usename	Name	Username.
commit_counter	Bigint	Number of commits.
rollback_counter	Bigint	Number of rollbacks.
resp_min	Bigint	Minimum response time.
resp_max	Bigint	Maximum response time.
resp_avg	Bigint	Average response time.
resp_total	Bigint	Total response time.

14.3.239 PGXC_VARIABLE_INFO

PGXC_VARIABLE_INFO displays information about XIDs and OIDs of all nodes in a cluster.

Table 14-308 PGXC_VARIABLE_INFO columns

Column	Туре	Description
node_name	Text	Node name.
nextOid	OID	Next OID to be generated under this node.
nextXid	Xid	Next transaction OID to be generated under this node.
oldestXid	Xid	Oldest transaction ID for a node
xidVacLimit	Xid	Critical point for forcing autovacuum.
oldestXidDB	OID	Database OID with the minimum datafrozenxid under this node.
lastExtendCSNL ogpage	Integer	Page number of the last extension of csnlog.
startExtendCSN Logpage	Integer	Starting page number of the csnlog extension.
nextCommitSeq No	Integer	Next CSN to be generated under this node.
latestCompleted Xid	Xid	Latest transaction ID on the node after commit or rollback.

Column	Туре	Description
startupMaxXid	Xid	Last transaction ID before the node shutdown.

14.3.240 PGXC_WAIT_DETAIL

PGXC_WAIT_DETAIL displays detailed information about the SQL waiting hierarchy of all nodes in a cluster. This view is supported only by clusters of version 8.1.3.200 or later.

Table 14-309 PGXC_WAIT_DETAIL columns

Column	Туре	Description
level	Integer	Level in the wait hierarchy. The value starts with 1 and increases by 1 when there is a wait relationship.
lock_wait_hi erarchy	Text	Wait hierarchy, in the format of <i>Node name: Process ID->Node name:Waiting process ID->Node name:Waiting process ID-></i>
node_name	Text	Node name
db_name	Text	Database name
thread_name	Text	Thread name
query_id	Bigint	ID of a query statement
tid	Bigint	Thread ID of the current thread
lwtid	Integer	Lightweight thread ID of the current thread
ptid	Integer	Parent thread of the streaming thread
tlevel	Integer	Level of the streaming thread
smpid	Integer	Concurrent thread ID
wait_status	Text	Waiting status of the current thread
wait_event	Text	Virtual ID of the transaction holding or awaiting this lock
exec_cn	Boolean	SQL execution CN
wait_node	Text	Lock level
query	Text	Query statement
application_ name	Text	Name of the application connected to the backend

Column	Туре	Description
backend_star t	Timestamp with time zone	Startup time of the backend process, that is, the time when the client connects to the server
xact_start	Timestamp with time zone	Start time of the current transaction
query_start	Timestamp with time zone	Start time of the active query
waiting	Boolean	Waiting status
state	Text	Overall state of the backend
waittime	Timestamp with time zone	Timestamp when the lock wait starts. This column is available only in clusters of version 9.1.0.200 or later.
holdtime	Timestamp with time zone	Timestamp when the lock starts to be obtained. This column is available only in clusters of version 9.1.0.200 or later.

Example

- **Step 1** Connect to the CN, start a transaction, and perform the update operation. begin;update td set c2=6 where c1=1;
- **Step 2** Open another window to connect to the CN, start another transaction, and perform the update operation. (Do not update the same record concurrently.) begin;update td set c2=6 where c1=7;

In this case, the update operation is blocked.

- **Step 3** Open another window to connect to the CN node and create an index. create index c2_key on td(c2);
- **Step 4** Run the **select * from pgxc_wait_detail**; command.

```
SELECT * FROM PGXC WAIT DETAIL;
               lock wait hierarchy
                                           | node_name | db_name | thread_name | query_id
level |
          | lwtid | ptid | tlevel | sm
pid | wait_status | wait_event | exec_cn | wait_node |
                                                                          | application_name |
                                                          query
backend_start | xact_st
                              | waiting | state
              query_start
1 | cn_5001:139870843444360
                                          | cn_5001 | postgres | workload | 73183493945299462 |
139870843444360 | 578531 | 0 | 0 | wait node | | t | | WLM fetch collect info from data nodes | workload
2023-03-13 13:56:56.611486+08 | 2023-03-14 11:54
:33.562808+08 | 2023-03-13 13:57:00.262736+08 | t
                                                  active
                                          | cn_5001 | postgres | gsql | 73183493945299204 |
1 | cn 5001:139870843654544
```

```
139870843654544 | 722259 |
                                   0 |
0 | wait node | | t | |
11:52:05.176588+08 | 2023-03-14 11:52
                                   | update td set c2=6 where c1=1;
                                                                         gsql
                                                                                      | 2023-03-14
:19.054727+08 | 2023-03-14 11:53:58.114794+08 | t | active
                                              | cn_5001 | postgres | gsql
                                                                            | 73183493945299218 |
1 | cn_5001:139870843655296
139870843655296 | 722301 |
                                   update td set c2=6 where c1=7;
0 | wait node |
                    | t
                                                                                      | 2023-03-14
11:52:08.084265+08 | 2023-03-14 11:52
:42.978132+08 | 2023-03-14 11:53:59.459575+08 | t
                                                    active
1 | cn_5001:139870843656424
                                             | cn_5001 | postgres | gsql
                                                                            | 73183493945299223 |
139870843656424 | 722344 |
                                      | create index c2_key on td(c2);
0 | acquire lock | relation | t
                                                                          | gsql
                                                                                       | 2023-03-14
11:52:10.967028+08 | 2023-03-14 11:52
:53.463227+08 | 2023-03-14 11:54:00.25203+08 | t
2 | cn_5001:139870843656424 -> cn_5001:139870843655296 | cn_5001 | postgres | gsql
73183493945299218 | 139870843655296 | 722344 |
                 | f
                              | update td set c2=6 where c1=7;
                                                                    | gsql
                                                                                 | 2023-03-14
11:52:08.084265+08 | 2023-03-14 11:52
:42.978132+08 | 2023-03-14 11:53:59.459575+08 | t
(5 rows)
```

----End

14.3.241 PGXC WAIT EVENTS

PGXC_WAIT_EVENTS displays statistics on the waiting status and events of each node in the cluster. The content is the same as that displayed in **GS_WAIT_EVENTS**. This view is accessible only to users with system administrators rights.

Table 14-310 PGXC_WAIT_EVENTS columns

Column	Туре	Description
nodename	Name	Node name.
type	Text	Event type, which can be STATUS, LOCK_EVENT, LWLOCK_EVENT, or IO_EVENT
event	Text	Event name. For details, see PG_THREAD_WAIT_STATUS.
wait	Bigint	Number of times an event occurs. This column and all the columns below are values accumulated during process running.
failed_wait	Bigint	Number of waiting failures. In the current version, this column is used only for counting timeout errors and waiting failures of locks such as LOCK and LWLOCK.
total_wait_time	Bigint	Total duration of the event
avg_wait_time	Bigint	Average duration of the event
max_wait_time	Bigint	Maximum wait time of the event

Column	Туре	Description
min_wait_time	Bigint	Minimum wait time of the event

14.3.242 PGXC_WLM_OPERATOR_HISTORY

PGXC_WLM_OPERATOR_HISTORY displays operator information when a job is finished on all CNs. This view is used to query data from GaussDB(DWS), and the data in the database is cleared periodically every 3 minutes.

Only the system administrator or the preset role **gs_role_read_all_stats** can access this view. For details about columns in the view, see **Table 14-311**.

Table 14-311 GS_WLM_OPERATOR_INFO columns

Column	Туре	Description
nodename	Text	Name of the CN where the statement is executed
queryid	Bigint	Internal query_id used for statement execution
pid	Bigint	Backend thread ID
plan_node_id	Integer	plan_node_id of the execution plan of a query
plan_node_nam e	Text	Name of the operator corresponding to plan_node_id
start_time	Timestamp with time zone	Time when an operator starts to process the first data record
duration	Bigint	Total execution time of an operator. The unit is ms.
query_dop	Integer	Degree of parallelism (DOP) of the current operator
estimated_rows	Bigint	Number of rows estimated by the optimizer
tuple_processed	Bigint	Number of elements returned by the current operator
min_peak_mem ory	Integer	Minimum peak memory used by the current operator on all DNs. The unit is MB.
max_peak_me mory	Integer	Maximum peak memory used by the current operator on all DNs. The unit is MB.
average_peak_ memory	Integer	Average peak memory used by the current operator on all DNs. The unit is MB.

Column	Туре	Description
memory_skew_ percent	Integer	Memory usage skew of the current operator among DNs
min_spill_size	Integer	Minimum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
max_spill_size	Integer	Maximum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
average_spill_si ze	Integer	Average spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
spill_skew_perc ent	Integer	DN spill skew when a spill occurs
min_cpu_time	Bigint	Minimum execution time of the operator on all DNs. The unit is ms.
max_cpu_time	Bigint	Maximum execution time of the operator on all DNs. The unit is ms.
total_cpu_time	Bigint	Total execution time of the operator on all DNs. The unit is ms.
cpu_skew_perce nt	Integer	Skew of the execution time among DNs.
warning	Text	Warning. The following warnings are displayed: 1. Sort/SetOp/HashAgg/HashJoin spill 2. Spill file size large than 256MB 3. Broadcast size large than 100MB 4. Early spill 5. Spill times is greater than 3
		6. Spill on memory adaptive7. Hash table conflict

14.3.243 PGXC_WLM_OPERATOR_INFO

PGXC_WLM_OPERATOR_INFO displays the operator information of completed jobs executed on CNs. The data in this view is obtained from **GS_WLM_OPERATOR_INFO**.

Only the system administrator or the **gs_role_read_all_stats** role can access this view.

NOTICE

The **PGXC_WLM_OPERATOR_INFO** view can be queried only in the **postgres** database. If the view is queried in other databases, an error is reported.

GS_WLM_OPERATOR_INFO lists the columns in the **PGXC_WLM_OPERATOR_INFO** view.

Table 14-312 GS_WLM_OPERATOR_INFO columns

Column	Туре	Description
nodename	Text	Name of the CN where the statement is executed
queryid	Bigint	Internal query_id used for statement execution
pid	Bigint	Backend thread ID
plan_node_id	Integer	plan_node_id of the execution plan of a query
plan_node_nam e	Text	Name of the operator corresponding to plan_node_id
start_time	Timestamp with time zone	Time when an operator starts to process the first data record
duration	Bigint	Total execution time of an operator. The unit is ms.
query_dop	Integer	Degree of parallelism (DOP) of the current operator
estimated_rows	Bigint	Number of rows estimated by the optimizer
tuple_processed	Bigint	Number of elements returned by the current operator
min_peak_mem ory	Integer	Minimum peak memory used by the current operator on all DNs. The unit is MB.
max_peak_me mory	Integer	Maximum peak memory used by the current operator on all DNs. The unit is MB.
average_peak_ memory	Integer	Average peak memory used by the current operator on all DNs. The unit is MB.
memory_skew_ percent	Integer	Memory usage skew of the current operator among DNs
min_spill_size	Integer	Minimum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .

Column	Туре	Description
max_spill_size	Integer	Maximum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
average_spill_si ze	Integer	Average spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
spill_skew_perc ent	Integer	DN spill skew when a spill occurs
min_cpu_time	Bigint	Minimum execution time of the operator on all DNs. The unit is ms.
max_cpu_time	Bigint	Maximum execution time of the operator on all DNs. The unit is ms.
total_cpu_time	Bigint	Total execution time of the operator on all DNs. The unit is ms.
cpu_skew_perce nt	Integer	Skew of the execution time among DNs.
warning	Text	Warning. The following warnings are displayed:
		Sort/SetOp/HashAgg/HashJoin spill
		2. Spill file size large than 256MB
		3. Broadcast size large than 100MB
		4. Early spill
		5. Spill times is greater than 3
		6. Spill on memory adaptive
		7. Hash table conflict

14.3.244 PGXC_WLM_OPERATOR_STATISTICS

PGXC_WLM_OPERATOR_STATISTICS displays the operator information of jobs being executed on CNs. The system administrator can query job operator information of all users in the cluster, while common users can query only their own job operator information.

Table 14-156 lists the columns in the PGXC_WLM_OPERATOR_STATISTICS view.

Table 14-313 GS_WLM_OPERATOR_STATISTICS columns

Column	Туре	Description
queryid	Bigint	Internal query_id used for statement execution
pid	Bigint	ID of the backend thread

Column	Туре	Description
plan_node_id	Integer	plan_node_id of the execution plan of a query
plan_node_na me	Text	Name of the operator corresponding to plan_node_id. The maximum length of the operator name is 127 characters (excluding format characters such as spaces).
start_time	Timestamp with time zone	Time when the operator starts to be executed for the first time.
duration	Bigint	Total execution time of the operator from the start to the end, in milliseconds.
status	Text	Execution status of the current operator. The value can be waiting , running , or finished .
query_dop	Integer	DOP of the current operator
estimated_rows	Bigint	Number of rows estimated by the optimizer. If the number of returned estimated rows exceeds int64_max, int64_max is displayed.
tuple_processe d	Bigint	Total number of elements returned by the current operator on all DNs. If the estimated number of returned rows exceeds int64_max, int64_max is displayed.
min_peak_mem ory	Integer	Minimum peak memory used by the current operator on all DNs. The unit is MB.
max_peak_me mory	Integer	Maximum peak memory used by the current operator on all DNs. The unit is MB.
average_peak_ memory	Integer	Average peak memory used by the current operator on all DNs. The unit is MB.
memory_skew_ percent	Integer	Memory usage skew of the current operator among DNs
min_spill_size	Integer	Minimum logical spilled data among all DNs when a spill occurs, in MB. The default value is 0 .
max_spill_size	Integer	Maximum logical spilled data among all DNs when a spill occurs, in MB. The default value is 0 .
average_spill_si ze	Integer	Average logical spilled data among all DNs when a spill occurs, in MB. The default value is 0 .
spill_skew_perc ent	Integer	DN spill skew when a spill occurs

Column	Туре	Description
min_cpu_time	Bigint	Minimum execution time of the operator on all DNs. The unit is ms.
max_cpu_time	Bigint	Maximum execution time of the operator on all DNs. The unit is ms.
total_cpu_time	Bigint	Total execution time of the operator on all DNs. The unit is ms.
cpu_skew_perc ent	Integer	Skew of the execution time among DNs.
warning	Text	Warning. The following warnings are displayed: 1. Sort/SetOp/HashAgg/HashJoin spill 2. Spill file size large than 256MB 3. Broadcast size large than 100MB 4. Early spill 5. Spill times is greater than 3 6. Spill on memory adaptive 7. Hash table conflict
parent_id	Integer	Parent node ID of the operator node.
exec_count	Integer	Maximum number of times that the operator node can be executed on all DNs.
progress	Text	Progress information of the operator. For the first operator, it is the overall progress of the job. For other operators, it is the progress of the current operator.
min_net_size	Bigint	Minimum network communication data volume (KB) of the operator on all DNs. It mainly applies to network operators.
max_net_size	Bigint	Maximum network communication data volume (KB) of the operator on all DNs. It mainly applies to network operators.
total_net_size	Bigint	Total network communication data volume (KB) of the operator on all DNs. It mainly applies to network operators.
min_read_bytes	Bigint	Minimum amount of data read by the operator from disks on all DNs. The unit is KB.
max_read_byte s	Bigint	Maximum amount of data read by the operator from disks on all DNs. The unit is KB.
total_read_byte s	Bigint	Total amount of data read by the operator from disks on all DNs, in KB.

Column	Туре	Description
min_write_byte s	Bigint	Minimum amount of data written by the operator to disks on all DNs. The unit is KB.
max_write_byte	Bigint	Maximum amount of data written by the operator to disks on all DNs. The unit is KB.
total_write_byt es	Bigint	Total amount of data written by the operator to disks on all DNs, in KB.

14.3.245 PGXC_WLM_SESSION_INFO

PGXC_WLM_SESSION_INFO displays load management information for completed jobs executed on all CNs. The view information comes from the **GS_WLM_SESSION_INFO** system catalog.

NOTICE

The **PGXC_WLM_SESSION_INFO** view can be queried only in the **postgres** database. If the view is queried in other databases, an error is reported.

Table 14-314 PGXC_WLM_SESSION_INFO columns

Column	Туре	Description
datid	OID	OID of the database the backend is connected to
dbname	Text	Name of the database the backend is connected to
schemaname	Text	Schema name
nodename	Text	Name of the CN where the statement is run
username	Text	User name used for connecting to the backend
application_na me	Text	Name of the application that is connected to the backend
client_addr	inet	IP address of the client connected to this backend. If this column is null, it indicates either that the client is connected via a Unix socket on the server machine or that this is an internal process such as autovacuum.

Column	Туре	Description
client_hostnam e	Text	Host name of the connected client, as reported by a reverse DNS lookup of client_addr. This column will only be non-null for IP connections, and only when log_hostname is enabled.
client_port	Integer	TCP port number used by the client to communicate with the backend. If a Unix socket is used, it is -1.
query_band	Text	Job type, which can be set using the GUC parameter query_band and is a null string by default.
block_time	Bigint	Duration that a statement is blocked before being executed, including the statement parsing and optimization duration. The unit is ms.
start_time	Timestamp with time zone	Time when the statement starts to be executed
finish_time	Timestamp with time zone	Time when the statement execution ends
duration	Bigint	Execution time of a statement. The unit is ms.
estimate_total_ time	Bigint	Estimated execution time of a statement. The unit is ms.
status	Text	Final statement execution status. Its value can be finished (normal) or aborted (abnormal). The statement status here is the execution status of the database server. If the statement is successfully executed on the database server but an error is reported in the result set, the statement status is finished .
abort_info	Text	Exception information displayed if the final statement execution status is aborted .
resource_pool	Text	Resource pool used by the user
control_group	Text	Cgroup used by the statement
estimate_mem ory	Integer	Estimated memory used by a statement on a single instance. The unit is MB.
min_peak_mem ory	Integer	Minimum memory peak of a statement across all DNs. The unit is MB.
max_peak_me mory	Integer	Maximum memory peak of a statement across all DNs. The unit is MB.

Column	Туре	Description
average_peak_ memory	Integer	Average memory usage during statement execution. The unit is MB.
memory_skew_ percent	Integer	Memory usage skew of a statement among DNs
spill_info	Text	Spill information for the statement on all DNs. The options are:
		None : The statement has not been spilled to disks on any DNs.
		All : The statement has been spilled to disks on all DNs.
		[a:b]: The statement has been spilled to disks on a of b DNs.
min_spill_size	Integer	Minimum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
max_spill_size	Integer	Maximum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
average_spill_si ze	Integer	Average spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
spill_skew_perc ent	Integer	DN spill skew when a spill occurs
min_dn_time	Bigint	Minimum execution time of a statement across all DNs. The unit is ms.
max_dn_time	Bigint	Maximum execution time of a statement across all DNs. The unit is ms.
average_dn_tim e	Bigint	Average execution time of a statement across all DNs. The unit is ms.
dntime_skew_p ercent	Integer	Execution time skew of a statement among DNs.
min_cpu_time	Bigint	Minimum CPU time of a statement across all DNs. The unit is ms.
max_cpu_time	Bigint	Maximum CPU time of a statement across all DNs. The unit is ms.
total_cpu_time	Bigint	Total CPU time of a statement across all DNs. The unit is ms.
cpu_skew_perc ent	Integer	CPU time skew of a statement among DNs.

Column	Туре	Description
min_peak_iops	Integer	Minimum IOPS peak of a statement across all DNs. It is counted by ones in a column-store table and by ten thousands in a row-store table.
max_peak_iops	Integer	Maximum IOPS peak of a statement across all DNs. It is counted by ones in a column-store table and by ten thousands in a row-store table.
average_peak_i ops	Integer	Average IOPS peak of a statement across all DNs. It is counted by ones in a column-store table and by ten thousands in a row-store table.
iops_skew_perc ent	Integer	I/O skew across DNs
warning	Text	Warning. The following warnings and warnings related to SQL self-diagnosis tuning are displayed: 1. Spill file size large than 256MB 2. Broadcast size large than 100MB 3. Early spill 4. Spill times is greater than 3 5. Spill on memory adaptive 6. Hash table conflict
queryid	Bigint	Internal query ID used for statement execution
query	Text	Statement to be executed. A maximum of 64 KB of strings can be retained.
query_plan	Text	 Execution plan of a statement Specification restrictions: 1. Execution plans are displayed only for DML statements. 2. In 8.2.1.100 and later versions, the number of data binding times is added to the execution plans of Parse Bind Execute (PBE) statements to facilitate statement analysis. The number of data binding times is displayed in the format of PBE bind times: <i>Times</i>.
node_group	Text	Logical cluster of the user running the statement
pid	Bigint	PID of the backend thread for the statement.

Column	Туре	Description
lane	Text	Fast/Slow lane where the statement is executed
unique_sql_id	Bigint	ID of the normalized unique SQL
session_id	Text	Unique identifier of a session in the database system. Its format is session_start_time.tid.node_name.
min_read_bytes	Bigint	Minimum I/O read bytes of a statement across all DNs. The unit is byte.
max_read_byte s	Bigint	Maximum I/O read bytes of a statement across all DNs. The unit is byte.
average_read_b ytes	Bigint	Average I/O read bytes of a statement across all DNs.
min_write_byte s	Bigint	Minimum I/O write bytes of a statement across all DNs.
max_write_byte s	Bigint	Maximum I/O write bytes of a statement across all DNs.
average_write_ bytes	Bigint	Average I/O write bytes of a statement across all DNs.
recv_pkg	Bigint	Total number of communication packages received by a statement across all DNs.
send_pkg	Bigint	Total number of communication packages sent by a statement across all DNs.
recv_bytes	Bigint	Total received data of the statement stream, in byte.
send_bytes	Bigint	Total sent data of the statement stream, in byte.
stmt_type	Text	Query type corresponding to the statement.
except_info	Text	Information about the exception rule triggered by the statement.
parse_time	Bigint	Total parsing time before the statement is queued (including lexical and syntax parsing, optimization rewriting, and plan generation time), in milliseconds. This column is only supported in version 8.3.0.100 or later.
unique_plan_id	Bigint	ID of the normalized unique plan.
sql_hash	Text	Normalized SQL hash.
plan_hash	Text	Normalized plan hash.

Column	Туре	Description
disk_cache_hit_ ratio	numeric(5,2)	Disk cache hit rate. This column only applies to OBS 3.0 tables and foreign tables.
disk_cache_disk _read_size	Bigint	Total size of data read from disk cache, in MB. This column only applies to OBS 3.0 tables and foreign tables.
disk_cache_disk _write_size	Bigint	Total size of data written to disk cache, in MB. This column only applies to OBS 3.0 tables and foreign tables.
disk_cache_rem ote_read_size	Bigint	Total size of data read remotely from OBS due to disk cache read failure, in MB. This column only applies to OBS 3.0 tables and foreign tables.
disk_cache_rem ote_read_time	Bigint	Total number of times data is read remotely from OBS due to disk cache read failure. This column only applies to OBS 3.0 tables and foreign tables.
vfs_scan_bytes	Bigint	Total number of bytes scanned by the OBS virtual file system in response to upper-layer requests, in bytes. This column only applies to OBS 3.0 tables and foreign tables.
vfs_remote_rea d_bytes	Bigint	Total number of bytes actually read from OBS by the OBS virtual file system, in bytes. This column only applies to OBS 3.0 tables and foreign tables.
preload_submit _time	Bigint	Total time for submitting I/O requests in the prefetching process, in microseconds. This column only applies to OBS 3.0 tables.
preload_wait_ti me	Bigint	Total time for waiting for I/O requests in the prefetching process, in microseconds. This column only applies to OBS 3.0 tables.
preload_wait_c ount	Bigint	Total number of times that the prefetching process waits for I/O requests. This column only applies to OBS 3.0 tables.
disk_cache_loa d_time	Bigint	Total time for reading from disk cache, in microseconds. This column only applies to OBS 3.0 tables and foreign tables.
disk_cache_conf lict_count	Bigint	Number of times a block in the disk cache produces a hash conflict. This column only applies to OBS 3.0 tables and foreign tables.
disk_cache_erro r_count	Bigint	Number of disk cache read failures. This column only applies to OBS 3.0 tables and foreign tables.

Column	Туре	Description
disk_cache_erro r_code	Bigint	Error code for disk cache read failures. This column only applies to OBS 3.0 tables and foreign tables.
obs_io_req_avg _rtt	Bigint	Average Round Trip Time (RTT) for OBS I/O requests, in microseconds. This column only applies to OBS 3.0 tables and foreign tables.
obs_io_req_avg _latency	Bigint	Average delay for OBS I/O requests, in microseconds. This column only applies to OBS 3.0 tables and foreign tables.
obs_io_req_late ncy_gt_1s	Bigint	Number of OBS I/O requests with a latency exceeding 1 second. This column only applies to OBS 3.0 tables and foreign tables.
obs_io_req_late ncy_gt_10s	Bigint	Number of OBS I/O requests with a latency exceeding 10 seconds. This column only applies to OBS 3.0 tables and foreign tables.
obs_io_req_cou nt	Bigint	Total number of OBS I/O requests. This column only applies to OBS 3.0 tables and foreign tables.
obs_io_req_retr y_count	Bigint	Total number of retries for OBS I/O requests. This column only applies to OBS 3.0 tables and foreign tables.
obs_io_req_rate _limit_count	Bigint	Total number of times OBS I/O requests are flow-controlled. This column only applies to OBS 3.0 tables and foreign tables.

14.3.246 PGXC_WLM_SESSION_HISTORY

PGXC_WLM_SESSION_HISTORY displays load management records after job execution on all CNs. This view is used to query data from GaussDB(DWS), and the data in the database is cleared periodically every 3 minutes.

NOTICE

The **PGXC_WLM_SESSION_HISTORY** view can be queried only in the **postgres** database. If the view is queried in other databases, an error is reported.

GS_WLM_SESSION_HISTORY lists the columns in the **PGXC_WLM_SESSION_HISTORY** view.

Table 14-315 GS_WLM_SESSION_HISTORY columns

Column	Туре	Description
datid	OID	OID of the database this backend is connected to
dbname	Text	Name of the database the backend is connected to
schemaname	Text	Schema name
nodename	Text	Name of the CN where the statement is run
username	Text	User name used for connecting to the backend
application_na me	Text	Name of the application that is connected to the backend
client_addr	inet	IP address of the client connected to this backend. If this column is null, it indicates either that the client is connected via a Unix socket on the server machine or that this is an internal process such as autovacuum.
client_hostnam e	Text	Host name of the connected client, as reported by a reverse DNS lookup of client_addr . This column will only be non-null for IP connections, and only when log_hostname is enabled.
client_port	Integer	TCP port number that the client uses for communication with this backend, or -1 if a Unix socket is used
query_band	Text	Job type, which can be set using the GUC parameter query_band and is a null string by default.
block_time	Bigint	Duration that a statement is blocked before being executed, including the statement parsing and optimization duration. The unit is ms.
start_time	Timestamp with time zone	Time when the statement starts to be run
finish_time	Timestamp with time zone	Time when the statement execution ends
duration	Bigint	Execution time of a statement. The unit is ms.
estimate_total_ time	Bigint	Estimated execution time of a statement. The unit is ms.

Column	Туре	Description
status	Text	Final statement execution status. Its value can be finished (normal) or aborted (abnormal). The statement status here is the execution status of the database server. If the statement is successfully executed on the database server but an error is reported in the result set, the statement status is finished .
abort_info	Text	Exception information displayed if the final statement execution status is aborted .
resource_pool	Text	Resource pool used by the user
control_group	Text	Cgroup used by the statement
estimate_mem ory	Integer	Estimated memory used by a statement on a single instance. The unit is MB.
min_peak_mem ory	Integer	Minimum memory peak of a statement across all DNs. The unit is MB.
max_peak_me mory	Integer	Maximum memory peak of a statement across all DNs. The unit is MB.
average_peak_ memory	Integer	Average memory usage during statement execution. The unit is MB.
memory_skew_ percent	Integer	Memory usage skew of a statement among DNs.
spill_info	Text	Statement spill information on all DNs.
		None indicates that the statement has not been spilled to disks on any DNs.
		All : The statement has been spilled to disks on all DNs.
		[a:b]: The statement has been spilled to disks on a of b DNs.
min_spill_size	Integer	Minimum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
max_spill_size	Integer	Maximum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
average_spill_si ze	Integer	Average spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .
spill_skew_perc ent	Integer	DN spill skew when a spill occurs

Column	Туре	Description
min_dn_time	Bigint	Minimum execution time of a statement across all DNs. The unit is ms.
max_dn_time	Bigint	Maximum execution time of a statement across all DNs. The unit is ms.
average_dn_tim e	Bigint	Average execution time of a statement across all DNs. The unit is ms.
dntime_skew_p ercent	Integer	Execution time skew of a statement among DNs.
min_cpu_time	Bigint	Minimum CPU time of a statement across all DNs. The unit is ms.
max_cpu_time	Bigint	Maximum CPU time of a statement across all DNs. The unit is ms.
total_cpu_time	Bigint	Total CPU time of a statement across all DNs. The unit is ms.
cpu_skew_perce nt	Integer	CPU time skew of a statement among DNs.
min_peak_iops	Integer	Minimum IOPS peak of a statement across all DNs. It is counted by ones in a column-store table and by ten thousands in a row-store table.
max_peak_iops	Integer	Maximum IOPS peak of a statement across all DNs. It is counted by ones in a column-store table and by ten thousands in a row-store table.
average_peak_i ops	Integer	Average IOPS peak of a statement across all DNs. It is counted by ones in a column-store table and by ten thousands in a row-store table.
iops_skew_perc ent	Integer	I/O skew across DNs.
warning	Text	Warning. The following warnings and warnings related to SQL self-diagnosis tuning are displayed: 1. Spill file size large than 256MB 2. Broadcast size large than 100MB 3. Early spill 4. Spill times is greater than 3 5. Spill on memory adaptive 6. Hash table conflict

Column	Туре	Description		
queryid	Bigint	Internal query ID used for statement execution		
query	Text	Statement to be executed. A maximum of 64 KB of strings can be retained.		
query_plan	Text	 Execution plan of a statement. Specification restrictions: 1. Execution plans are displayed only for DML statements. 2. In 8.2.1.100 and later versions, the number of data binding times is added to the execution plans of Parse Bind Execute (PBE) statements to facilitate statement analysis. The number of data binding times is displayed in the format of PBE bind times: <i>Times</i>. 		
node_group	Text	Logical cluster of the user running the statement		
pid	Bigint	PID of the backend thread of the statement		
lane	Text	Fast/Slow lane where the statement is executed		
unique_sql_id	Bigint	ID of the normalized unique SQL.		
session_id	Text	Unique identifier of a session in the database system. Its format is session_start_time.tid.node_name.		
min_read_bytes	Bigint	Minimum I/O read bytes of a statement across all DNs. The unit is byte.		
max_read_byte s	Bigint	Maximum I/O read bytes of a statement across all DNs. The unit is byte.		
average_read_b ytes	Bigint	Average I/O read bytes of a statement across all DNs.		
min_write_byte s	Bigint	Minimum I/O write bytes of a statement across all DNs.		
max_write_byte s	Bigint	Maximum I/O write bytes of a statement across all DNs.		
average_write_ bytes	Bigint	Average I/O write bytes of a statement across all DNs.		
recv_pkg	Bigint	Total number of communication packages received by a statement across all DNs.		
send_pkg	Bigint	Total number of communication packages sent by a statement across all DNs.		

Column	Туре	Description	
recv_bytes	Bigint	Total received data of the statement stream, in byte.	
send_bytes	Bigint	Total sent data of the statement stream, in byte.	
stmt_type	Text	Query type corresponding to the statement.	
except_info	Text	Information about the exception rule triggered by the statement.	
unique_plan_id	Bigint	ID of the normalized unique plan.	
sql_hash	Text	Normalized SQL hash.	
plan_hash	Text	Normalized plan hash.	
use_plan_baseli ne	Text	Indicates whether the bound plan is used for executing the current statement. If it is used, the name of the plan_baseline column in pg_plan_baseline is displayed.	
outline_name	Text	Name of the outline used for the statement plan.	
loader_status	Text	 The JSON string for storing import and export service information is as follows. address: indicates the IP address of the peer cluster. The port number is displayed for the source cluster. direction: indicates the import and export service type. The value can be gds to file, gds from file, gds to pipe, gds from pipe, copy from or copy to. 	
		3. min/max/total_lines/bytes: indicates the minimum value, maximum value, total lines, and bytes of the import and export statements on all DNs.	
parse_time	Bigint	Total parsing time before the statement is queued (including lexical and syntax parsing, optimization rewriting, and plan generation time), in milliseconds. This column is available only in clusters of version 8.3.0.100 or later.	
disk_cache_hit_ ratio	numeric(5,2	Disk cache hit rate. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.	
disk_cache_disk _read_size	Bigint	Total size of data read from disk cache, in MB. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.	

Column	Туре	Description	
disk_cache_disk _write_size	Bigint	Total size of data written to disk cache, in MB. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.	
disk_cache_rem ote_read_size	Bigint	Total size of data read remotely from OBS due to disk cache read failure, in MB. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.	
disk_cache_rem ote_read_time	Bigint	Total number of times data is read remotely from OBS due to disk cache read failure. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.	
vfs_scan_bytes	Bigint	Total number of bytes scanned by the OBS virtual file system in response to upper-layer requests, in bytes. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.	
vfs_remote_rea d_bytes	Bigint	Total number of bytes actually read from OBS by the OBS virtual file system, in bytes. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.	
preload_submit _time	Bigint	Total time for submitting I/O requests in the prefetching process, in microseconds. This column only applies to OBS 3.0 tables with storage and compute decoupled.	
preload_wait_ti me	Bigint	Total time for waiting for I/O requests in the prefetching process, in microseconds. This column only applies to OBS 3.0 tables with storage and compute decoupled.	
preload_wait_c ount	Bigint	Total number of times that the prefetching process waits for I/O requests. This column only applies to OBS 3.0 tables with storage and compute decoupled.	
disk_cache_loa d_time	Bigint	Total time for reading from disk cache, in microseconds. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.	
disk_cache_conf lict_count	Bigint	Number of times a block in the disk cache produces a hash conflict. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.	

Column	Туре	Description	
disk_cache_erro r_count	Bigint	Number of disk cache read failures. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.	
disk_cache_erro r_code	Bigint	 Error code for disk cache read failures. Multiple error codes may be generated. If the disk cache fails to be read, OBS remote read is initiated and cache blocks are rewritten. The error code types are as follows: This column only applies to OBS 3.0 tables and foreign tables. 1: A hash conflict occurs in the disk cache block. 2: The generation time of the disk cache block is later than that of the OldestXmin transaction. 4: Invoking the pread system when reading cache files from the disk cache failed. 8: The data version of the disk cache block does not match. 16: The version of the data written to the write cache does not match the latest version. 32: Opening the cache file corresponding to the cache block failed. 64: The size of the data read from the disk cache does not match. 	
		128: The CSN recorded in the disk cache block does not match.	
obs_io_req_avg _rtt	Bigint	Average Round Trip Time (RTT) for OBS I/O requests, in microseconds. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.	
obs_io_req_avg _latency	Bigint	Average delay for OBS I/O requests, in microseconds. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.	
obs_io_req_late ncy_gt_1s	Bigint	Number of OBS I/O requests with a latency exceeding 1 second. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.	
obs_io_req_late ncy_gt_10s	Bigint	Number of OBS I/O requests with a latency exceeding 10 seconds. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.	

Column	Туре	Description		
obs_io_req_cou nt	Bigint	Total number of OBS I/O requests. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.		
obs_io_req_retr y_count	Bigint	Total number of retries for OBS I/O requests. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.		
obs_io_req_rate _limit_count	Bigint	Total number of times OBS I/O requests are flow-controlled. This column only applies to OBS 3.0 tables and foreign tables with storage and compute decoupled.		

14.3.247 PGXC_WLM_SESSION_STATISTICS

PGXC_WLM_SESSION_STATISTICS displays load management information about jobs that are being executed on CNs.

Table 14-316 PGXC_WLM_SESSION_STATISTICS columns

Column	Туре	Description	
datid	OID	OID of the database this backend is connected to	
dbname	Name	Name of the database the backend is connected to	
schemaname	Text	Schema name	
nodename	Text	Name of the CN where the statement is executed	
username	Name	User name used for connecting to the backend	
application_nam e	Text	Name of the application that is connected to the backend	
client_addr	inet	IP address of the client connected to this backend. If this column is null, it indicates either that the client is connected via a Unix socket on the server machine or that this is an internal process such as autovacuum.	
client_hostname	Text	Host name of the connected client, as reported by a reverse DNS lookup of client_addr . This column will only be non-null for IP connections, and only when log_hostname is enabled.	

Column	Туре	Description		
client_port	Integer	TCP port number used by the client to communicate with the backend. If a Unix socket is used, it is -1.		
query_band	Text	Job type, which can be set using the GUC parameter query_band and is a null string by default.		
pid	Bigint	ID of the backend thread		
block_time	Bigint	Block time before the statement is executed. The unit is ms.		
start_time	Timestamp with time zone	Time when the statement starts to be executed		
duration	Bigint	For how long a statement has been executing. The unit is ms.		
estimate_total_ti me	Bigint	Estimated execution time of a statement. The unit is ms.		
estimate_left_ti me	Bigint	Estimated remaining time of statement execution. The unit is ms.		
enqueue	Text	Workload management resource status		
resource_pool	Name	Resource pool used by the user		
control_group	Text	Cgroup used by the statement		
estimate_memor y	Integer	Estimated memory used by a statement on a single instance. The unit is MB.		
min_peak_mem ory	Integer	Minimum memory peak of a statement across all DNs. The unit is MB.		
max_peak_mem ory	Integer	Maximum memory peak of a statement across all DNs. The unit is MB.		
average_peak_m emory	Integer	Average memory usage during statement execution. The unit is MB.		
memory_skew_p ercent	Integer	Memory usage skew of a statement among DNs.		

Column	Туре	Description		
spill_info	Text	Spill information for the statement on all DNs. The options are:		
		None : The statement has not been spilled to disks on any DNs.		
		All: The statement has been spilled to disks on all DNs.		
		[a:b]: The statement has been spilled to disks on a of b DNs.		
min_spill_size	Integer	Minimum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .		
max_spill_size	Integer	Maximum spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .		
average_spill_siz e	Integer	Average spilled data among all DNs when a spill occurs. The unit is MB. The default value is 0 .		
spill_skew_perce nt	Integer	DN spill skew when a spill occurs		
min_dn_time	Bigint	Minimum execution time of a statement across all DNs. The unit is ms.		
max_dn_time	Bigint	Maximum execution time of a statement across all DNs. The unit is ms.		
average_dn_tim e	Bigint	Average execution time of a statement across all DNs. The unit is ms.		
dntime_skew_pe rcent	Integer	Execution time skew of a statement among DNs.		
min_cpu_time	Bigint	Minimum CPU time of a statement across all DNs. The unit is ms.		
max_cpu_time	Bigint	Maximum CPU time of a statement across all DNs. The unit is ms.		
total_cpu_time	Bigint	Total CPU time of a statement across all DNs. The unit is ms.		
cpu_skew_perce nt	Integer	CPU time skew of a statement among DNs.		
min_peak_iops	Integer	Minimum IOPS peak of a statement across all DNs. It is counted by ones in a column-store table and by ten thousands in a row-store table.		

Column	Туре	Description		
max_peak_iops	Integer	Maximum IOPS peak of a statement across all DNs. It is counted by ones in a column-store table and by ten thousands in a row-store table.		
average_peak_io ps	Integer	Average IOPS peak of a statement across all DNs. It is counted by ones in a column-store table and by ten thousands in a row-store table.		
iops_skew_perce nt	Integer	I/O skew across DNs.		
min_read_speed	Integer	Minimum I/O read rate of a statement across all DNs within a monitoring period (5s). The unit is KB/s.		
max_read_speed	Integer	Maximum I/O read rate of a statement across all DNs within a monitoring period (5s). The unit is KB/s.		
average_read_sp eed	Integer	Average I/O read rate of a statement across all DNs within a monitoring period (5s). The unit is KB/s.		
min_write_speed	Integer	Minimum I/O write rate of a statement across all DNs within a monitoring period (5s). The unit is KB/s.		
max_write_spee d	Integer	Maximum I/O write rate of a statement across all DNs within a monitoring period (5s). The unit is KB/s.		
average_write_s peed	Integer	Average I/O write rate of a statement across all DNs within a monitoring period (5s). The unit is KB/s.		
recv_pkg	Bigint	Total number of communication packages received by a statement across all DNs.		
send_pkg	Bigint	Total number of communication packages sent by a statement across all DNs.		
recv_bytes	Bigint	Total received data of the statement stream, in byte.		
send_bytes	Bigint	Total sent data of the statement stream, in byte.		

Column	Туре	Description	
warning	Text	Warning. The following warnings and warnings related to SQL self-diagnosis tuning are displayed: 1. Spill file size large than 256MB 2. Broadcast size large than 100MB 3. Early spill	
		4. Spill times is greater than 3	
		5. Spill on memory adaptive	
		6. Hash table conflict	
unique_sql_id	Bigint	ID of the normalized unique SQL.	
queryid	Bigint	Internal query ID used for statement execution	
query	Text	Statement that is being executed	
query_plan	Text	Execution plan of a statement	
		Specification restrictions:	
		Execution plans are displayed only for DML statements.	
		2. In 8.2.1.100 and later versions, the number of data binding times is added to the execution plans of Parse Bind Execute (PBE) statements to facilitate statement analysis. The number of data binding times is displayed in the format of PBE bind times : <i>Times</i> .	
node_group	Text	Logical cluster of the user running the statement	
stmt_type	Text	Query type corresponding to the statement.	
except_info	Text	Information about the exception rule triggered by the statement.	
parse_time	Bigint	Total parsing time before the statement is queued (including lexical and syntax parsing, optimization rewriting, and plan generation time), in milliseconds. This column is only supported in version 8.3.0.100 or later.	
unique_plan_id	Bigint	ID of the normalized unique plan.	
sql_hash	Text	Normalized SQL hash.	
plan_hash	Text	Normalized plan hash.	
disk_cache_hit_r atio	numeric(5, 2)	Disk cache hit rate. This column only applies to OBS 3.0 tables and foreign tables in decoupled storage and compute scenarios.	

Column	Туре	Description		
disk_cache_disk_ read_size	Bigint	Total size of data read from disk cache, in MB. This column only applies to OBS 3.0 tables and foreign tables in decoupled storage and compute scenarios.		
disk_cache_disk_ write_size	Bigint	Total size of data written to disk cache, in MB. This column only applies to OBS 3.0 tables and foreign tables in decoupled storage and compute scenarios.		
disk_cache_remo te_read_size	Bigint	Total size of data read remotely from OBS due to disk cache read failure, in MB. This column only applies to OBS 3.0 tables and foreign tables in decoupled storage and compute scenarios.		
disk_cache_remo te_read_time	Bigint	Total number of times data is read remotely from OBS due to disk cache read failure. This column only applies to OBS 3.0 tables and foreign tables in decoupled storage and compute scenarios.		
block_name	Text	Name of the interception rule that matches the statement.		

14.3.248 PGXC_WLM_TABLE_DISTRIBUTION_SKEWNESS

PGXC_WLM_TABLE_DISTRIBUTION_SKEWNESS displays data skews of tables in the current database. You can quickly query the storage space skew of all tables in the current database on each node. This view is supported only by clusters of version 8.2.1 or later.

Suggestions:

- When you analyze the disk space skew of each table in a database in a large cluster with a large amount of data, the PGXC_WLM_TABLE_DISTRIBUTION_SKEWNESS view delivers better query performance than the gs_table_distribution() function and the PGXC_GET_TABLE_SKEWNESS view. You are advised to use the PGXC_WLM_TABLE_DISTRIBUTION_SKEWNESS view to query the table skew status overview, and then use the gs_table_distribution(schemaname text, tablename text) function to obtain the disk space distribution of a specified table on each node.
- To use this view to query the storage distribution information of a specified table, you must have the **SELECT** permission on the table.
- This function is based on the physical file storage space recorded in the PG_RELFILENODE_SIZE system catalog. Ensure that the GUC parameters use workload manager and enable perm space are enabled.

Column **Type** Description schema name Name Name of the schema where a table is table_name Name Table name total_size Numeric Total storage space of a table on all nodes, in bytes avg_size numeric(1000,0) Average storage space of a table on each node, in bytes Numeric Percentage (%) of the maximum max_percent storage space of a table on each node to the total storage space min_percent Numeric Percentage (%) of the minimum storage space of a table on each node to the total storage space Numeric Skew rate (%) of a table skew percent The formula for calculating the skew rate as follows: Skew rate (skew_percent) = (Maximum value of dnsize - Average value of dnsize) x 100/Maximum value of dnsize.

Table 14-317 PGXC_WLM_TABLE_DISTRIBUTION_SKEWNESS columns

Example

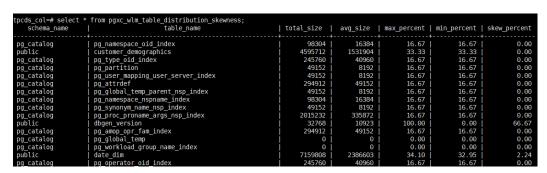
You can use the **PGXC_WLM_TABLE_DISTRIBUTION_SKEWNESS** view to query the table skew status overview, and then use the

gs_table_distribution(schemaname text, tablename text) function to obtain the disk space distribution of a specified table on each node.

Step 1 Use the **PGXC_WLM_TABLE_DISTRIBUTION_SKEWNESS** view to query the table skew status overview.

tpcds_col=# select * from pgxc_wlm_table_distribution_skewness;

The query result is as follows:



The data skew of the **dbgen_version** table is severe.

Step 2 Use the **gs_table_distribution(schemaname text, tablename text)** function to query the disk space distribution of the **dbgen_version** table on each node.

tpcds_col=# select * from gs_table_distribution('public','dbgen_version');

The query result is as follows:

tpcds_col=# select * from schemaname tablename	¯ relkind rel	persistence nodename	dnsize sessionid
public dbgen_versio public dbgen_versio public dbgen_versio (3 rows)	n r p	dn_6001_6000 dn_6005_6000 dn_6003_6004	2 0 5 32768

According to the preceding information, data skew occurs in the disk space occupied by the table on DNs. Most data is stored on **dn_6005_6006**.

----End

14.3.249 PGXC_WLM_USER_RESOURCE_HISTORY

The **PGXC_WLM_USER_RESOURCE_HISTORY** view displays historical information about resource consumption of all users on the corresponding instances. This view is supported only by clusters of version 8.2.0 or later.

Table 14-318 PGXC_WLM_USER_RESOURCE_HISTORY columns

Column	Туре	Description
nodename	Name	Instance name, including CNs and DNs.
username	Text	Username
Timestamp	Timestamp with time zone	Timestamp
used_memory	Integer	Used memory (unit: MB).
		On a DN, it indicates a user's memory usage on the current DN.
		On a CN, it indicates a user's total memory usage on all DNs.
total_memory	Integer	Available memory (unit: MB). 0 indicates that the available memory is not limited and depends on the maximum memory available in the database.
		On a DN, it indicates the memory available to a user on the current DN.
		On a CN, it indicates the total memory available to a user on all DNs.

Column	Туре	Description
used_cpu	Double precision	Number of CPU cores in use. Only the CPU usage of complex jobs in the non-default resource pool is collected, and the value is the CPU usage of the related cgroup.
		On a DN, it indicates a user's CPU core usage on the current DN.
		On a CN, it indicates a user's total CPU core usage on all DNs.
total_cpu	Integer	Total number of CPU cores of the Cgroups associated with a user.
		On a DN, it indicates the CPU cores available to a user on the current DN.
		On a CN, it indicates the total CPU cores available to a user on all DNs.
used_space	Bigint	Used permanent table storage space (unit: KB)
		On a DN, it indicates the size of the permanent table storage space used by a user on the current DN.
		On a CN, it indicates the total size of the permanent table storage space used by a user on all DNs.
total_space	Bigint	Available storage space, in KB. -1 indicates that the storage space is not limited.
		On a DN, it indicates the size of the permanent table storage space available to a user on the current DN.
		On a CN, it indicates the total size of the permanent table storage space available to a user on all DNs.
used_temp_sp	Bigint	Used temporary table storage space (unit: KB)
ace		On a DN, it indicates the size of the temporary table storage space used by a user on the current DN.
		On a CN, it indicates the total size of the temporary table storage space used by a user on all DNs.

Column	Туре	Description
total_temp_sp ace	Bigint	Available temporary table storage space, in KB. -1 indicates that the storage space is not limited.
		On a DN, it indicates the size of the temporary table storage space available to a user on the current DN.
		On a CN, it indicates the total size of the temporary table storage space available to a user on all DNs.
used_spill_spa ce	Bigint	Size of space used for operator spill to disk, in KB.
		On a DN, it indicates displays the size of the operator flushing space used by the user on the current DN.
		On a CN, it indicates the total space used by a user's operators spilled to disk on all DNs.
total_spill_spa ce	Bigint	Size of space available for operator spill to disk, in KB. The value -1 indicates that the space is not limited.
		On a DN, it indicates displays the size of the operator flushing space that can be used by the user on the current DN.
		On a CN, it indicates the total space available for a user to spill operators to disk on all DNs.
read_kbytes	Bigint	On a CN, it indicates total number of bytes read by a user's complex jobs on all DNs in the last 5 seconds. The unit is KB.
		On a DN, it indicates the total number of bytes read by a user's complex jobs from the instance startup time to the current time. The unit is KB.
write_kbytes	Bigint	On a CN, it indicates total number of bytes written by a user's complex jobs on all DNs in the last 5 seconds.
		On a DN, it indicates the total number of bytes written by a user's complex jobs from the instance startup time to the current time. The unit is KB.
read_counts	Bigint	On a CN, it indicates total number of read times of a user's complex jobs on all DNs in the last 5 seconds.
		On a DN, it indicates total number of read times of a user's complex jobs from the instance startup time to the current time.

Column	Туре	Description
write_counts	Bigint	On a CN, it indicates total number of write times of a user's complex jobs on all DNs in the last 5 seconds.
		On a DN, it indicates total number of write times of a user's complex jobs from the instance startup time to the current time.
read_speed	Double precision	On a CN, it indicates the average read rate of a user's complex jobs on a single DN in the last 5 seconds, in KB/s.
		On a DN, it indicates the average read rate of a user's complex jobs on the DN in the last 5 seconds, in KB/s.
write_speed	Double precision	On a CN, it indicates the average write rate of a user's complex jobs on a single DN in the last 5 seconds, in KB/s.
		On a DN, it indicates the average write rate of a user's complex jobs on the DN in the last 5 seconds, in KB/s.
send_speed	Double precision	On a CN, it indicates the sum of the average network sending rates of a user on all DNs in a 5s monitoring period, in KB/s.
		On a DN, it indicates the average network sending rate of a user on the DN in a 5s monitoring period, in KB/s.
recv_speed	Double precision	On a CN, it indicates the sum of the average network receiving rates of a user on all DNs in a 5s monitoring period, in KB/s.
		On a DN, it indicates the average network receiving rate of a user on the DN in a 5s monitoring period, in KB/s.

14.3.250 PGXC_WLM_WORKLOAD_RECORDS

PGXC_WLM_WORKLOAD_RECORDS displays the status of job executed by the current user on CNs. Only the system administrator or the preset role **gs_role_read_all_stats** can access this view. This view is available only when **enable_dynamic_workload** is set to **on**.

Table 14-319 PGXC_WLM_WORKLOAD_RECORDS columns

Column	Туре	Description
node_name	Text	Name of the CN where the job is executed.

Column	Туре	Description
thread_id	Bigint	ID of the backend thread.
processid	Integer	lwpid of the thread.
Timestamp	Bigint	Start time of statement execution.
username	Name	Username logged in to the backend.
memory	Integer	Memory required for the statement.
active_points	Integer	Number of resources consumed by the statement on the resource pool.
max_points	Integer	Maximum number of resources in the resource pool.
priority	Integer	Priority of a job.
resource_pool	Text	Resource pool where a job is.
status	Text	Job execution status. The options are:
		pending
		running
		finished
		aborted
		unknown
control_group	Text	Cgroups used by a job.
enqueue	Text	Queue for the job, including:
		GLOBAL: global queue.
		RESPOOL: resource pool queue.
		ACTIVE: not queued.
query	Text	Statement currently being executed.

14.3.251 PGXC_WORKLOAD_SQL_COUNT

PGXC_WORKLOAD_SQL_COUNT displays statistics on the number of SQL statements executed in workload Cgroups on all CNs in a cluster, including the number of **SELECT**, **UPDATE**, **INSERT**, and **DELETE** statements and the number of DDL, DML, and DCL statements. Only the system administrator or the preset role **gs_role_read_all_stats** can access this view.

Table 14-320 PGXC_WORKLOAD_SQL_COUNT columns

Column	Туре	Description
node_name	Name	Node name.

Column	Туре	Description
workload	Name	Workload Cgroup name.
select_count	Bigint	Number of SELECT statements.
update_count	Bigint	Number of UPDATE statements.
insert_count	Bigint	Number of INSERT statements.
delete_count	Bigint	Number of DELETE statements.
ddl_count	Bigint	Number of DDL statements.
dml_count	Bigint	Number of DML statements.
dcl_count	Bigint	Number of DCL statements.

14.3.252 PGXC_WORKLOAD_SQL_ELAPSE_TIME

PGXC_WORKLOAD_SQL_ELAPSE_TIME displays statistics on the response time of SQL statements in workload Cgroups on all CNs in a cluster, including the maximum, minimum, average, and total response time of **SELECT**, **UPDATE**, **INSERT**, and **DELETE** statements. The unit is microsecond. Only the system administrator or the preset role **gs_role_read_all_stats** can access this view.

Table 14-321 PGXC_WORKLOAD_SQL_ELAPSE_TIME columns

Column	Туре	Description
node_name	Name	Node name.
workload	Name	Workload Cgroup name.
total_select_elapse	Bigint	Total response time of SELECT statements.
max_select_elapse	Bigint	Maximum response time of SELECT statements.
min_select_elapse	Bigint	Minimum response time of SELECT statements.
avg_select_elapse	Bigint	Average response time of SELECT statements.

Column	Туре	Description
total_update_elapse	Bigint	Total response time of UPDATE statements.
max_update_elapse	Bigint	Maximum response time of UPDATE statements.
min_update_elapse	Bigint	Minimum response time of UPDATE statements.
avg_update_elapse	Bigint	Average response time of UPDATE statements.
total_insert_elapse	Bigint	Total response time of INSERT statements.
max_insert_elapse	Bigint	Maximum response time of INSERT statements.
min_insert_elapse	Bigint	Minimum response time of INSERT statements.
avg_insert_elapse	Bigint	Average response time of INSERT statements.
total_delete_elapse	Bigint	Total response time of DELETE statements.
max_delete_elapse	Bigint	Maximum response time of DELETE statements.
min_delete_elapse	Bigint	Minimum response time of DELETE statements.
avg_delete_elapse	Bigint	Average response time of DELETE statements.

14.3.253 PGXC_WORKLOAD_TRANSACTION

PGXC_WORKLOAD_TRANSACTION provides transaction information about workload cgroups on all CNs. Only the system administrator or the preset role <code>gs_role_read_all_stats</code> can access this view. This view is valid only when the real-time resource monitoring function is enabled, that is, when <code>enable_resource_track</code> is <code>on</code>.

Table 14-322 PGXC_WORKLOAD_TRANSACTION columns

Column	Туре	Description
node_name	Name	Node name.
workload	Name	Workload Cgroup name.

Column	Туре	Description
commit_counter	Bigint	Number of the commits.
rollback_counter	Bigint	Number of rollbacks.
resp_min	Bigint	Minimum response time, in microseconds.
resp_max	Bigint	Maximum response time, in microseconds.
resp_avg	Bigint	Average response time, in microseconds.
resp_total	Bigint	Total response time, in microseconds.

14.3.254 PLAN_TABLE

PLAN_TABLE displays the plan information collected by **EXPLAIN PLAN**. Plan information is in a session-level life cycle. After the session exits, the data will be deleted. Data is isolated between sessions and between users.

Table 14-323 PLAN_TABLE columns

Column	Туре	Description
statement_id	varchar2(30)	Query tag specified by a user
plan_id	Bigint	ID of a plan to be queried
id	Int	ID of each operator in a generated plan
operation	varchar2(30)	Operation description of an operator in a plan
options	varchar2(255)	Operation parameters
object_name	Name	Name of an operated object. It is defined by users, not the object alias used in the query.
object_type	varchar2(30)	Object type
object_owner	Name	User-defined schema to which an object belongs
projection	varchar2(400 0)	Returned column information

■ NOTE

- A valid object_type value consists of a relkind type defined in PG_CLASS (TABLE ordinary table, INDEX, SEQUENCE, VIEW, FOREIGN TABLE, COMPOSITE TYPE, or TOASTVALUE TOAST table) and the rtekind type used in the plan (SUBQUERY, JOIN, FUNCTION, VALUES, CTE, or REMOTE QUERY).
- For RangeTableEntry (RTE), **object_owner** is the object description used in the plan. Non-user-defined objects do not have **object owner**.
- Information in the **statement_id**, **object_name**, **object_owner**, and **projection** columns is stored in letter cases specified by users and information in other columns is stored in uppercase.
- PLAN_TABLE supports only SELECT and DELETE and does not support other DML operations.

14.3.255 PV_FILE_STAT

By collecting statistics about the data file I/Os, **PV_FILE_STAT** displays the I/O performance of the data to detect the performance problems, such as abnormal I/O operations.

Table 14-324 PV_FILE_STAT columns

Column	Туре	Description
filenum	OID	File ID.
dbid	OID	Database ID.
spcid	OID	Tablespace ID.
phyrds	Bigint	Number of physical files read.
phywrts	Bigint	Number of physical files written.
phyblkrd	Bigint	Number of physical file blocks read.
phyblkwrt	Bigint	Number of physical file blocks written.
readtim	Bigint	Total duration of file reads, in microseconds.
writetim	Bigint	Total duration of file writes, in microseconds.
avgiotim	Bigint	Average duration of file reads and writes, in microseconds.
lstiotim	Bigint	Duration of the last file read, in microseconds.
miniotim	Bigint	Minimum duration of file reads and writes, in microseconds.
maxiowtm	Bigint	Maximum duration of file reads and writes, in microseconds.

14.3.256 PV INSTANCE TIME

PV_INSTANCE_TIME collects statistics on the running time of processes and the time consumed in each execution phase, in microseconds.

PV_INSTANCE_TIME records time consumption information of the current node. The time consumption information is classified into the following types:

- **DB_TIME**: effective time spent by jobs in multi-core scenarios
- **CPU_TIME**: CPU time spent
- **EXECUTION_TIME**: time spent within executors
- PARSE_TIME: time spent on parsing SQL statements
- PLAN_TIME: time spent on generating plans
- **REWRITE_TIME**: time spent on rewriting SQL statements
- PL_EXECUTION_TIME: execution time of the PL/pgSQL stored procedure
- PL_COMPILATION_TIME: compilation time of the PL/pgSQL stored procedure
- **NET_SEND_TIME**: time spent on the network
- DATA_IO_TIME: I/O time spent

Table 14-325 PV_INSTANCE_TIME columns

Column	Туре	Description
stat_id	Integer	Type ID.
stat_name	Text	Name of the runtime type.
value	Bigint	Runtime value.

14.3.257 PV_MATVIEW_DETAIL

PV_MATVIEW_DETAIL displays detailed information about a materialized view. This view is supported only by clusters of version 9.1.0 200 or later.

Table 14-326 PV_MATVIEW_DETAIL columns

Column	Туре	Description
matview	Text	Materialized view name
baserel	Text	Base table name
partids	oidvector	OID of a specified partition.
contain_entire _rel	Boolean	Whether to create a materialized view based on the entire base table

Column	Туре	Description
build_mode	Text	Build mode of the materialized view.
		• 'd': indicates "deferred", which means that data is contained in the materialized view only when the view is refreshed for the first time.
		'i': indicates "immediate", which means that the latest data is included when the materialized view is created.
refresh_mode	Text	Refresh mode of the materialized view.
		'd': stands for demand, indicating on-demand update.
refresh_meth	Text	Refresh method of the materialized view.
od		'c' indicates a full refresh.
mapping	Text	Mapping between base table partitions and materialized view partitions
active	Boolean	Whether the materialized view needs to be refreshed
refresh_start_t ime	Timestamp with time zone	Start time of the last refresh
refresh_finish_ time	Timestamp with time zone	End time of the last refresh

14.3.258 PV_OS_RUN_INFO

PV_OS_RUN_INFO displays the running status of the current operating system.

Table 14-327 PV_OS_RUN_INFO columns

Column	Туре	Description
id	Integer	ID.
Name	Text	Name of the operating system status.
value	Numeric	Value of the operating system status.
comments	Text	Comments on the operating system status.
cumulative	boolean	Whether the value of the operating system status is cumulative.

14.3.259 PV SESSION MEMORY

PV_SESSION_MEMORY displays statistics about memory usage at the session level in the unit of MB, including all the memory allocated to Postgres and Stream threads on DNs for jobs currently executed by users.

Table 14-328 PV_SESSION_MEMORY columns

Column	Туре	Description	
sessid	Text	Thread start time and ID.	
init_mem	Integer	Memory allocated to the currently executed task before the task enters the executor, in MB.	
used_mem	Integer	Memory allocated to the currently executed task, in MB.	
peak_mem	Integer	Peak memory allocated to the currently executed task, in MB.	

14.3.260 PV_SESSION_MEMORY_DETAIL

PV_SESSION_MEMORY_DETAIL displays statistics about thread memory usage by memory context.

The memory context TempSmallContextGroup collects information about all memory contexts whose value in the **totalsize** column is less than 8192 bytes in the current thread, and the number of the collected memory contexts is recorded in the **usedsize** column. Therefore, the **totalsize** and **freesize** columns for TempSmallContextGroup in the view display the corresponding information about all the memory contexts whose value in the **totalsize** column is less than 8192 bytes in the current thread, and the **usedsize** column displays the number of these memory contexts.

You can run the **SELECT * FROM pv_session_memctx_detail (***threadid,***'');** statement to record information about all memory contexts of a thread into the *threadid_timestamp.log* file in the **/tmp/dumpmem** directory. *threadid* can be obtained from the following table.

Table 14-329 PV_SESSION_MEMORY_DETAIL columns

Column	Туре	Description	
sessid	Text	Thread start time+thread ID (string: timestamp.threadid)	
sesstype	Text	Thread name	
contextname	Text	Text Name of the memory context	
level	Smallint	Hierarchy of the memory context	
parent	Text	Name of the parent memory context	

Column	Туре	Description
totalsize	Bigint	Total size of the memory context, in bytes
freesize	Bigint	Total size of released memory in the memory context, in bytes
usedsize	Bigint	Size of used memory in the memory context, in bytes. For TempSmallContextGroup, this parameter specifies the number of collected memory contexts.

Example

Query the usage of all MemoryContexts on the current node.

Locate the thread in which the MemoryContext is created and used based on **sessid**. Check whether the memory usage meets the expectation based on **totalsize**, **freesize**, and **usedsize** to see whether memory leakage may occur.

totalsize freesize usedsize	esstype	contextname	level	parent
++ ++			+	
	postmaster g		1	
TopMemoryContext	17209904 8081136	9128768		
1667462258.139973631031	040 postgres	SRF multi-call context		5
functionScan_139973631031		68 1722336		
1667461280.139973666686				1
opMemoryContext				
1667450443.139973877479				1
TopMemoryContext	1472544 356088	1116456		
1667462258.139973631031				1
opMemoryContext				
1667461250.139973915236				1
opMemoryContext				
1667450439.139974010144			t	1
opMemoryContext				
1667450439.139974151726	848 WDRSnapsnot	CacheMemoryConte	ext	1
opMemoryContext				
1667450439.139974026925			Xt	1
opMemoryContext				1 11
1667451036.139973746386				1
opMemoryContext				1 11
1667461250.139973950891				1
opMemoryContext 1667450439.139974076212			tovt	1
opMemoryContext			text	
opiwemorycontext 1667450439.139974092994:			tovt	1
opMemoryContext			itext	
1667461254.139973971343				1
opMemoryContext				'
1667461280.139973822945				1
opMemoryContext				1 '1
1667450439.139974202070 [°]				1
opMemoryContext				1 '1
1667450454.139973860697		CacheMemoryContext		1
opMemoryContext				1 '1
0.139975915622720	nostmaster 1	Postmaster	1	I
opMemoryContext	1 1004288 88792		,	1
1667450439.139974218852			ext	1
opMemoryContext				1 1
opinicinion y context	1 3 10230 103400	, 0 1, 00		

1667461250.139973915236096 postgres	TempSmallContextGroup	0
584448 148032	119	
1667462258.139973631031040 postgres	TempSmallContextGroup	0
579712 162128	123	

14.3.261 PV_SESSION_STAT

PV_SESSION_STAT displays session state statistics based on session threads or the **AutoVacuum** thread.

Table 14-330 PV_SESSION_STAT columns

Column	Туре	Description
sessid	Text	Thread ID and thread start time.
statid	Integer	Statistics ID.
statname	Text	Name of the statistics session.
statunit	Text	Unit of the statistics session.
value	Bigint	Value of the statistics session.

14.3.262 PV_SESSION_TIME

PV_SESSION_TIME displays statistics about the running time of session threads and time consumed in each execution phase, in microseconds.

Table 14-331 PV_SESSION_TIME columns

Column	Туре	Description	
sessid	Text	Thread ID and thread start time.	
stat_id	Integer	Statistics ID.	
stat_name	Text	Name of the runtime type.	
value	Bigint	Runtime value.	

14.3.263 PV_TOTAL_MEMORY_DETAIL

PV_TOTAL_MEMORY_DETAIL displays statistics about memory usage of the current database node in the unit of MB.

Table 14-332 PV_TOTAL_MEMORY_DETAIL columns

Name	Туре	Description
nodename	Text	Node name

Name	Туре	Description
memorytype	Text	Memory type. Its value can be:
		max_process_memory: memory used by a GaussDB(DWS) cluster instance
		• process_used_memory: memory used by a GaussDB(DWS) process
		max_dynamic_memory: maximum dynamic memory
		• dynamic_used_memory : used dynamic memory
		dynamic_peak_memory: dynamic peak value of the memory
		dynamic_used_shrctx: maximum dynamic shared memory context
		dynamic_peak_shrctx: dynamic peak value of the shared memory context
		max_shared_memory: maximum shared memory
		• shared_used_memory: used shared memory
		max_cstore_memory: maximum memory allowed for column store
		• cstore_used_memory : memory used for column store
		max_sctpcomm_memory: maximum memory allowed for the communication library
		• sctpcomm_used_memory : memory used for the communication library
		sctpcomm_peak_memory: memory peak of the communication library
		other_used_memory: other used memory
		gpu_max_dynamic_memory: maximum GPU memory
		gpu_dynamic_used_memory: sum of the available GPU memory and temporary GPU memory
		gpu_dynamic_peak_memory: maximum memory used for GPU
		pooler_conn_memory: memory used for pooler connections
		pooler_freeconn_memory: memory used for idle pooler connections
		storage_compress_memory: memory used for column-store compression and decompression
		udf_reserved_memory: memory reserved for the UDF Worker process

Name	Туре	Description	
		• mmap_used_memory: memory used for mmap	
memorymbytes	Integer	Size of allocated memory-typed memory	

14.3.264 PV_REDO_STAT

PV_REDO_STAT displays statistics on redoing Xlogs on the current node.

Table 14-333 PV_REDO_STAT columns

Name	Туре	Description	
phywrts	Bigint	Number of physical writes.	
phyblkwrt	Bigint	Number of physical blocks written.	
writetim	Bigint	Time taken for physical writes.	
avgiotim	Bigint	Average time taken per write.	
lstiotim	Bigint	Time taken for the last write.	
miniotim	Bigint	Minimum time taken for a write.	
maxiowtm	Bigint	Maximum time taken for a write.	

14.3.265 PV_RUNTIME_ATTSTATS

PV_RUNTIME_ATTSTATS displays table-level statistics in the memory generated by autoanalyze. The descriptions of the columns in **PV_RUNTIME_RELSTATS** are the same as those in **PG_STATS**. This view is used only by clusters of version 8.2.0 or later.

Table 14-334 PV_RUNTIME_ATTSTATS columns

Column	Туре	Reference	Description
schemaname	Name	PG_NAMESP ACE.nspname	Name of the schema that contains the table
tablename	Name	PG_CLASS.rel name	Table name
attname	Name	PG_ATTRIBU TE.attname	Column name

Column	Туре	Reference	Description
inherited	boolean	-	If the value is true , the inherited subcolumns are included. If the value is false , only the columns in a specified table are included.
null_frac	Real	-	Percentage of column entries that are null
avg_width	Integer	-	Average width in bytes of column's entries
n_distinct	Real		 If the value is greater than 0, it indicates the estimated number of distinct values in the column. Negative of the number of distinct values divided by the number of rows if the value is less than 0 The negated form is used when ANALYZE believes that the number of distinct values is likely to increase as the table grows. The positive form is used when the column seems to have a fixed number of possible values. For example, -1 indicates a unique column in which the number of distinct values is the same as the number of rows.
n_dndistinct	Real	-	 Number of unique non-null data values in the dn1 column Exact number of distinct values if the value is greater than 0 Negative of the number of distinct values divided by the number of rows if the value is less than 0 (For example, if the value of a column appears twice in average, set n_dndistinct=-0.5.) The number of distinct values is unknown if the value is 0.
most_commo n_vals	anyarray	-	List of the most common values in a column. If this combination does not have the most common values, it will be NULL .

Column	Туре	Reference	Description
most_commo n_freqs	real[]	-	List of the frequencies of the most common values, that is, the number of occurrences of each value divided by the total number of rows. (NULL if most_common_vals is NULL)
histogram_bo unds	anyarray	-	List of values that divide the column's values into groups of equal proportion. The values in most_common_vals, if present, are omitted from this histogram calculation. This field is null if the field data type does not have a < operator or if the most_common_vals list accounts for the entire population.
correlation	Real	-	Statistical correlation between physical row ordering and logical ordering of the column values. It ranges from -1 to +1. When the value is near to -1 or +1, an index scan on the column is estimated to be cheaper than when it is near to zero, due to reduction of random access to the disk. This column is null if the column data type does not have a < operator.
most_commo n_elems	anyarray	-	A list of the most commonly used non-null element values
most_commo n_elem_freqs	real[]	-	A list of the frequencies of the most commonly used element values
elem_count_h istogram	real[]	-	A histogram of the counts of distinct non-null element values

14.3.266 PV_RUNTIME_RELSTATS

PV_RUNTIME_RELSTATS displays table-level statistics in the memory generated by autoanalyze. The descriptions of the columns in **PV_RUNTIME_RELSTATS** are the same as those in **PG_CLASS**. This view is used only by clusters of version 8.2.0 or later.

Table 14-335 PV_RUNTIME_RELSTATS columns

Name	Туре	Description
nspname	Name	Schema name.
relname	Name	Name of an object, such as a table or index.
relpages	Double precision	Size of the on-disk representation of this table in pages (of size BLCKSZ). This is only an estimate used by the optimizer.
reltuples	Double precision	Number of rows in the table. This is only an estimate used by the optimizer.
relallvisible	Integer	Number of pages marked as all visible in the table. This column is used by the optimizer for optimizing SQL execution.
relhasindex	Boolean	Its value is true if this column is a table and has (or recently had) at least one index. It is set by CREATE INDEX but is not immediately cleared by DROP INDEX . If the VACUUM process detects that a table has no index, it clears the relhasindex column and sets the value to false .
changes	Bigint	Total historical modifications in the table by the time the lightweight autoanalyze is triggered.
level	Text	Current phase of the memory statistics generated by the lightweight autoanalyze. It can be local, sendlist, or global.

14.3.267 REDACTION_COLUMNS

REDACTION_COLUMNS displays information about all redaction columns in the current database.

Table 14-336 REDACTION_COLUMNS columns

Name	Туре	Description
object_schema	name	Redacted object schema.
object_owner	Name	Redacted object owner.
object_name	Name	Redacted object name.
column_name	Name	Redacted column name.
function_type	Integer	Redaction type.

Name	Туре	Description
function_parameters	Text	Parameter used when the redaction type is partial (reserved).
regexp_pattern	Text	Pattern string when the redaction type is regexp (reserved).
regexp_replace_string	Text	Replacement string when the redaction type is regexp (reserved).
regexp_position	Integer	Start and end replacement positions when the redaction type is regexp (reserved).
regexp_occurrence	Integer	Replacement times when the redaction type is regexp (reserved).
regexp_match_parameter	Text	Regular control parameter used when the redaction type is regexp (reserved).
function_info	Text	Redaction function information.
column_description	Text	Description of the redacted column.
inherited	Bool	Whether a redacted column is inherited from another redacted column.
policy_name	Name	Name of the data masking policy. This parameter is supported only by clusters of version 8.2.1.100 or later.

14.3.268 REDACTION_POLICIES

REDACTION_POLICIES displays information about all redaction objects in the current database.

Table 14-337 REDACTION_POLICIES columns

Name	Туре	Description
object_schema	name	Redacted object schema.
object_owner	Name	Redacted object owner.
object_name	Name	Redacted object name.
policy_name	Name	Name of the redaction policy.
expression	Text	Policy effective expression (for users).
enable	Boolean	Policy status (enabled or disabled).
policy_description	Text	Policy description.
inherited	Bool	Whether a redacted column is inherited from another redacted column.

14.3.269 REMOTE_TABLE_STAT

REMOTE_TABLE_STAT provides statistics of all tables of the database on all DNs in the cluster. Except the **nodename** column of the name type added in front of each row, the names, types, and sequences of other columns are the same as those in the **GS_TABLE_STAT** view.

Table 14-338 REMOTE_TABLE_STAT columns

Name	Туре	Description
nodename	Name	Node name
schemaname	Name	Table namespace
relname	Name	Table name
seq_scan	Bigint	Number of sequential scans. Only row-store tables are counted. For a partitioned table, the sum of the number of scans of each partition is displayed.
seq_tuple_rea d	Bigint	Number of rows scanned in sequence. Only row-store tables are counted.
index_scan	Bigint	Number of index scans. Only row-store tables are counted.

Name	Туре	Description
index_tuple_re ad	Bigint	Number of rows scanned by the index. Only row-store tables are counted.
tuple_inserted	Bigint	Number of rows inserted
tuple_updated	Bigint	Number of rows updated
tuple_deleted	Bigint	Number of rows deleted
tuple_hot_upd ated	Bigint	Number of rows with HOT updates.
live_tuples	Bigint	Number of live tuples. Query the view on the CN. If ANALYZE is executed, the total number of live tuples in the table is displayed. Otherwise, 0 is displayed. This indicator applies only to row-store tables.
dead_tuples	Bigint	Number of dead tuples. Query the view on the CN. If ANALYZE is executed, the total number of dead tuples in the table is displayed. Otherwise, 0 is displayed. This indicator applies only to row-store tables.

14.3.270 SHOW_TSC_INFO

Queries TSC information about the current node. This view is supported only by clusters of version 8.2.1 or later.

Table 14-339 Parameter

Name	Туре	Description
node_name	text	Node name
tsc_mult	bigint	TSC conversion multiplier
tsc_shift	bigint	TSC conversion shifts
tsc_frequency	float8	TSC frequency.
tsc_use_freque ncy	boolean	Indicates whether to use the TSC frequency for time conversion.
tsc_ready	boolean	Indicates whether the TSC frequency can be used for time conversion
tsc_scalar_erro r_info	text	Error information about obtaining TSC conversion information
tsc_freq_error_ info	text	Error information about obtaining TSC frequency information

14.3.271 SHOW_ALL_TSC_INFO

Queries TSC information about all nodes. This view is supported only by clusters of version 8.2.1 or later.

Table 14-340 Parameter

Name	Туре	Description
node_name	text	Node name
tsc_mult	bigint	TSC conversion multiplier
tsc_shift	bigint	TSC conversion shifts
tsc_frequency	float8	TSC frequency.
tsc_use_freque ncy	boolean	Indicates whether to use the TSC frequency for time conversion.
tsc_ready	boolean	Indicates whether the TSC frequency can be used for time conversion
tsc_scalar_erro r_info	text	Error information about obtaining TSC conversion information
tsc_freq_error_ info	text	Error information about obtaining TSC frequency information

14.3.272 USER_COL_COMMENTS

USER_COL_COMMENTS stores the column comments of the tables and views that the current user can access.

Column	Туре	Description
column_name	character varying(64)	Column name
table_name	character varying(64)	Table or view name
owner	character varying(64)	Owner of the table or view
comments	Text	Comments

14.3.273 USER_CONSTRAINTS

USER_CONSTRAINTS displays the table constraint information accessible to the current user.

Column	Туре	Description	
constraint_name	vcharacter varying(64)	Constraint name	
constraint_type	Text	Constraint type	
		C: Check constraint	
		• F : Foreign key constraint	
		P: Primary key constraint	
		U: Unique constraint.	
table_name	character varying(64)	Name of constraint-related table	
index_owner	character varying(64)	Owner of constraint-related index (only for the unique constraint and primary key constraint)	
index_name	character varying(64)	Name of constraint-related index (only for the unique constraint and primary key constraint)	

Example

Query constraints on a specified table of the current user. Replace **t1** with the actual table name.

14.3.274 USER_CONS_COLUMNS

USER_CONS_COLUMNS displays the information about constraint columns in the tables accessible to the current user.

Column	Туре	Description
table_name	character varying(64)	Name of constraint-related table
column_name	character varying(64)	Name of constraint-related column
constraint_name	character varying(64)	Constraint name
position	Smallint	Position of the column in the table

14.3.275 USER_INDEXES

USER_INDEXES displays index information in the current schema.

Column	Туре	Description
owner	character varying(64)	Index owner
index_name	character varying(64)	Index name
table_name	character varying(64)	Name of the table corresponding to the index
uniqueness	Text	Whether the index is unique
generated	character varying(1)	Whether the index name is generated by the system
partitioned	character(3)	Whether the index has the property of the partition table

14.3.276 USER_IND_COLUMNS

USER_IND_COLUMNS displays column information about all indexes accessible to the current user.

Column	Туре	Description
index_owner	character varying(64)	Index owner
index_name	character varying(64)	Index name
table_owner	character varying(64)	Table owner
table_name	character varying(64)	Table name
column_name	Name	Column name
column_position	Smallint	Position of a column in the index

14.3.277 USER_IND_EXPRESSIONS

USER_IND_EXPRESSIONS displays information about the function-based expression index accessible to the current user.

Column	Туре	Description
index_owner	character varying(64)	Index owner
index_name	character varying(64)	Index name
table_owner	character varying(64)	Table owner
table_name	character varying(64)	Table name
column_expression	Text	Function-based index expression of a specified column
column_position	Smallint	Position of a column in the index

14.3.278 USER_IND_PARTITIONS

USER_IND_PARTITIONS displays information about index partitions accessible to the current user.

Column	Туре	Description
index_owner	character varying(64)	Name of the owner of the partitioned table index to which the index partition belongs
schema	character varying(64)	Schema of the partitioned index to which the index partition belongs
index_name	character varying(64)	Index name of the partitioned table to which the index partition belongs
partition_nam e	character varying(64)	Name of the index partition
index_partitio n_usable	boolean	Whether the index partition is available
high_value	Text	Boundary of the table partition corresponding to the index partition. For a range partition, the boundary is the upper boundary. For a list partition, the boundary is the boundary value set.
		Reserved field for forward compatibility. The parameter pretty_high_value is added in version 8.1.3 to record the information.

Column	Туре	Description
pretty_high_v alue	Text	Boundary of the table partition corresponding to the index partition. For a range partition, the boundary is the upper boundary. For a list partition, the boundary is the boundary value set.
		The query result is the instant decompilation output of the partition boundary expression. The output of this column is more detailed than that of high_value . The output information can be collation and column data type.
def_tablespac e_name	Name	Tablespace name of the index partition

14.3.279 USER_JOBS

USER_JOBS displays all scheduled jobs owned by the current user. This view is accessible only to users with system administrator rights.

Table 14-341 USER_JOBS columns

Column	Туре	Description
job	INT4	Job ID
log_user	name not null	User name of the job creator
priv_user	name not null	User name of the job executor
dbname	name not null	Database in which the job is created
start_date	Timestamp without time zone	Job start time
start_suc	Text	Start time of the successful job execution
last_date	Timestamp without time zone	Start time of the last job execution
last_suc	Text	Start time of the last successful job execution
this_date	Timestamp without time zone	Start time of the ongoing job execution

Column	Туре	Description
this suc	Text	Same as THIS_DATE
next_date	Timestamp without time zone	Schedule time of the next job execution
next suc	Text	Same as next_date
broken	Text	Task status
		Y: the system does not try to execute the task.
		N : the system attempts to execute the task.
status	Char	Status of the current job. The value range is 'r', 's', 'f', 'd'. The default value is 's'. The indications are as follows: • r: running • s: finished
		• f: failed
		• d : aborted
interval	Text	Time expression used to calculate the next execution time. If this parameter is set to null , the job will be executed once only.
failures	Smallint	Number of consecutive failures.
what	Text	Body of the PL/SQL blocks or anonymous clock that the job executes

14.3.280 USER_OBJECTS

USER_OBJECTS displays all database objects accessible to the current user.

Column	Туре	Description
owner	Name	Owner of the object
object_name	Name	Object name
object_id	OID	OID of the object
object_type	Name	Type of the object
namespace	OID	Namespace containing the object
created	Timestamp with time zone	Object creation time
last_ddl_time	Timestamp with time zone	Last time when the object was modified

NOTICE

For details about the value ranges of **last_ddl_time** and **last_ddl_time**, see **PG_OBJECT**.

14.3.281 USER_PART_INDEXES

USER_PART_INDEXES displays information about partitioned table indexes accessible to the current user.

Column	Туре	Description
index_owner	character varying(64)	Name of the owner of the partitioned table index
schema	character varying(64)	Schema of the partitioned table index
index_name	character varying(64)	Name of the partitioned table index
table_name	character varying (64)	Name of the partitioned table to which the partitioned table index belongs
partitioning_type	Text	Partition policy of the partitioned table NOTE Currently, only range partitioning and list partitioning are supported.
partition_count	Bigint	Number of index partitions of the partitioned table index
def_tablespace_name	Name	Tablespace name of the partitioned table index
partitioning_key_coun t	Integer	Number of partition keys of the partitioned table

14.3.282 USER_PART_TABLES

USER_PART_TABLES displays information about partitioned tables accessible to the current user.

Column	Туре	Description
table_owner	character varying(64)	Name of the owner of the partitioned table
schema	character varying(64)	Schema of the partitioned table

Column	Туре	Description
table_name	character varying(64)	Name of the partitioned table
partitioning_type	Text	Partition policy of the partitioned table
		NOTE Currently, only range partitioning and list partitioning are supported.
partition_count	Bigint	Number of partitions of the partitioned table
def_tablespace_name	Name	Tablespace name of the partitioned table
partitioning_key_count	Integer	Number of partition keys of the partitioned table

14.3.283 USER_PROCEDURES

USER_PROCEDURES displays information about all stored procedures and functions in the current schema.

Column	Туре	Description
owner	character varying(64)	Owner of the stored procedure or the function
object_name	character varying(64)	Name of the stored procedure or the function
argument_number	Smallint	Number of the input parameters in the stored procedure

14.3.284 USER_SEQUENCES

USER_SEQUENCES displays sequence information in the current schema.

Column	Туре	Description
sequence_owner	character varying(64)	Owner of the sequence
sequence_name	character varying(64)	Name of the sequence

14.3.285 USER_SOURCE

USER_SOURCE displays information about stored procedures or functions in this mode, and provides the columns defined by the stored procedures or the functions.

Column	Туре	Description
owner	character varying(64)	Owner of the stored procedure or the function
Name	character varying(64)	Name of the stored procedure or the function
Text	Text	Definition of the stored procedure or the function

14.3.286 USER_SYNONYMS

USER_SYNONYMS displays synonyms accessible to the current user.

Table 14-342 USER_SYNONYMS columns

Name	Туре	Description
schema_name	Text	Name of the schema the synonym belongs to.
synonym_name	Text	Synonym name.
table_owner	Text	Owner of the associated object.
table_schema_na me	Text	Name of the schema the associated object belongs to.
table_name	Text	Name of the associated object.

14.3.287 USER_TAB_COLUMNS

USER_TAB_COLUMNS stores information about columns of the tables and views that the current user can access.

Column	Туре	Description
owner	character varying(64)	Owner of a table/view
table_name	character varying(64)	Table/View name

Column	Туре	Description
column_name	character varying(64)	Column name
data_type	character varying(128)	Data type of the column
column_id	Integer	Sequence number of the column when a table/view is created
data_length	Integer	Length of the column, in bytes
comments	Text	Comments
avg_col_len	Numeric	Average length of a column, in bytes
nullable	bpchar	Whether the column can be empty. For the primary key constraint and non-null constraint, the value is n.
data_precision	Integer	Precision of the data type. This parameter is valid for the numeric data type and NULL for other data types.
data_scale	Integer	Number of decimal places. This parameter is valid for the numeric data type and 0 for other data types.
char_length	Numeric	Length of a column, in characters. This parameter is valid only for the varchar, nvarchar2, bpchar, and char types.
schema	character varying(64)	Namespace that contains the table or view.
kind	Text	Type of the current record. If the column belongs to a table, the value of this column is table . If the column belongs to a view, the value of this column is view .

14.3.288 USER_TAB_COMMENTS

USER_TAB_COMMENTS displays comments about all tables and views accessible to the current user.

Column	Туре	Description
owner	character varying(64)	Owner of the table or view
table_name	character varying(64)	Name of the table or view

Column	Туре	Description
comments	Text	Comments

14.3.289 USER_TAB_PARTITIONS

USER_TAB_PARTITIONS displays all table partitions accessible to the current user. Each partition of a partitioned table accessible to the current user has a piece of record in **USER_TAB_PARTITIONS**.

Column	Туре	Description
table_owner	character varying(64)	Owner of the table that contains the partition
schema	character varying(64)	Schema of the partitioned table
table_name	character varying(64)	Table name
partition_name	character varying(64)	Name of the partition
high_value	Text	Upper boundary of a range partition or boundary value set of a list partition
		Reserved field for forward compatibility. The parameter pretty_high_value is added in version 8.1.3 to record the information.
pretty_high_valu e	Text	Upper boundary of a range partition or boundary value set of a list partition
		The query result is the instant decompilation output of the partition boundary expression. The output of this column is more detailed than that of high_value. The output information can be collation and column data type.
tablespace_name	Name	Name of the tablespace that contains the partition

14.3.290 USER_TABLES

USER_TABLES displays table information in the current schema.

Column	Туре	Description
owner	character varying(64)	Table owner
table_name	character varying(64)	Table name
tablespace_name	character varying(64)	Name of the tablespace that contains the table
status	character varying(8)	Whether the current record is valid
temporary	character(1)	 Whether the table is a temporary table Y indicates that it is a temporary table. N indicates that it is not a temporary table.
dropped	character varying	 Whether the current record is deleted YES indicates that it is deleted. NO indicates that it is not deleted.
num_rows	Numeric	Estimated number of rows in the table

14.3.291 USER_TRIGGERS

USER_TRIGGERS displays the information about triggers accessible to the current user.

Column	Туре	Description
trigger_name	character varying(64)	Trigger name
table_name	character varying(64)	Name of the table that defines the trigger
table_owner	character varying(64)	Owner of the table that defines the trigger

14.3.292 USER_VIEWS

USER_VIEWS displays information about all views in the current schema.

Column	Туре	Description
owner	character varying(64)	Owner of the view
view_name	character varying(64)	View name

14.3.293 V\$SESSION

V\$SESSION displays all session information about the current session.

Table 14-343 V\$SESSION columns

Name	Туре	Description
sid	Bigint	OID of the background process of the current activity
serial#	Integer	Sequence number of the active background process, which is 0 in GaussDB(DWS).
user#	OID	OID of the user that has logged in to the background process
username	Name	Name of the user that has logged in to the background process

14.3.294 V\$SESSION_LONGOPS

V\$SESSION_LONGOPS displays the progress of ongoing operations.

Table 14-344 V\$SESSION_LONGOPS columns

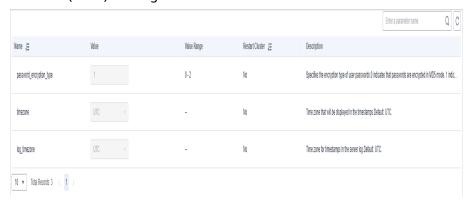
Name	Туре	Description
sid	Bigint	OID of the running background process
serial#	Integer	Sequence number of the running background process, which is 0 in GaussDB(DWS).
sofar	Integer	Completed workload, which is empty in GaussDB(DWS).
totalwork	Integer	Total workload, which is empty in GaussDB(DWS).

15 GUC Parameters of the GaussDB(DWS) Database

15.1 Viewing GUC Parameters

GaussDB(DWS) GUC parameters can control database system behaviors. You can check and adjust the GUC parameters based on your business scenario and data volume.

• After a cluster is installed, you can check database parameters on the GaussDB(DWS) management console.



- You can also connect to a cluster and run SQL commands to check the GUC parameters.
 - Run the SHOW command.
 - □ NOTE

Method 2 can only be used to check the GUC parameter values of CNs, while the GUC parameter values of DNs can be viewed through Method 1: by using the management console.

To view a certain parameter, run the following command:

SHOW server_version;

server_version indicates the database version.

Run the following command to view values of all parameters: **SHOW ALL**:

- Use the **pg_settings** view.

To view a certain parameter, run the following command: SELECT * FROM pg_settings WHERE NAME='server_version';

Run the following command to view values of all parameters: **SELECT * FROM pg_settings**;

15.2 Configuring GUC Parameters

To ensure the optimal performance of GaussDB(DWS), you can adjust the GUC parameters in the database.

Parameter Types and Values

- The GUC parameters of GaussDB(DWS) are classified into the following types:
 - SUSET: database administrator parameters. This type of parameters takes
 effect immediately after they are set. You do not need to restart the
 cluster. If a parameter of this type is set in the current session, the
 parameter takes effect only in the current session.
 - USERSET: common user parameters. This type of parameters takes effect immediately after they are set. You do not need to restart the cluster. If a parameter of this type is set in the current session, the parameter takes effect only in the current session.
 - POSTMASTER: database server parameters. This type of parameters takes
 effect only after the cluster is restarted. After you modify a parameter of
 this type, the system displays a message indicating that the cluster is to
 be restarted. You are advised to manually restart the cluster during offpeak hours for the setting to take effect.
 - SIGHUP: global database parameters. This type of parameters takes effect globally and cannot take effect for single sessions.
 - BACKEND: global database parameters. This type of parameters takes effect globally and cannot take effect for single sessions.
- All parameter names are case insensitive. A parameter value can be an integer, floating point number, string, Boolean value, or enumerated value.
 - The Boolean values can be on/off, true/false, yes/no, or 1/0, and are case-insensitive.
 - The enumerated value range is specified in the enumvals column of the system catalog pg_settings.
- For parameters using units, specify their units during the setting, or default units are used.
 - The default units are specified in the unit column of pg settings.
 - The unit of memory can be KB, MB, or GB.
 - The unit of time can be ms, s, min, h, or d.

Setting GUC Parameters

You can configure GUC parameters in the following ways:

- Method 1: After a cluster is created, log in to the GaussDB(DWS) console and modify the database parameters of the cluster. For details, see Modifying Database Parameters.
- Method 2: Connect to a cluster and run SQL commands to configure the parameters of the SUSET or USERSET type.

Set parameters at database, user, or session levels.

- Set a database-level parameter.

ALTER DATABASE dbname SET paraname TO value;

The setting takes effect in the next session.

Set a user-level parameter.

ALTER USER username SET paraname TO value;

The setting takes effect in the next session.

Set a session-level parameter.

SET paraname TO value;

Parameter value in the current session is changed. After you exit the session, the setting becomes invalid.

Procedure

The following example shows how to set **explain_perf_mode**.

Step 1 View the value of **explain_perf_mode**.

```
SHOW explain_perf_mode;
explain_perf_mode
-----
normal
(1 row)
```

Step 2 Set explain_perf_mode.

Perform one of the following operations:

• Set a database-level parameter.

ALTER DATABASE gaussdb SET explain_perf_mode TO pretty,

If the following information is displayed, the setting has been modified.

ALTER DATABASE

The setting takes effect in the next session.

Set a user-level parameter.

ALTER USER dbadmin SET explain_perf_mode TO pretty,

If the following information is displayed, the setting has been modified.

ALTER USER

The setting takes effect in the next session.

• Set a session-level parameter.

SET explain_perf_mode TO pretty;

If the following information is displayed, the setting has been modified.

SET

Step 3 Check whether the parameter is correctly set.

```
SHOW explain_perf_mode; explain_perf_mode
```

pretty (1 row)

----End

15.3 GUC Parameter Usage

The database provides many operation parameters. Configuration of these parameters affects the behavior of the database system. Before modifying these parameters, learn the impact of these parameters on the database. Otherwise, unexpected results may occur.

Precautions

- If the value range of a parameter is a string, the string should comply with the naming conventions of the path and file name in the OS running the database.
- If the allowed maximum value of a parameter is **INT_MAX**, it indicates the maximum parameter value varies by OS.
- If the allowed maximum value of a parameter is **DBL_MAX**, it indicates the maximum parameter value varies by OS.

15.4 Connection and Authentication

15.4.1 Connection Settings

This section describes parameters related to the connection mode between the client and server.

max connections

Parameter description: Specifies the maximum number of allowed parallel connections to the database. This parameter influences the concurrent processing capability of the cluster.

Type: POSTMASTER

Value range: an integer. For CNs, the value ranges from 100 to 16384. For DNs, the value ranges from 100 to 262143. Because there are internal connections in the cluster, the maximum value is rarely reached. If **invalid value for parameter** "max_connections" is displayed in the log, you need to decrease the max_connections value for DNs.

Default value: **800** for CNs and **5000** for DNs. If the default value is greater than the maximum value supported by kernel (determined when the **gs_initdb** command is executed), an error message will be displayed.

Setting suggestions:

Retain the default value of this parameter on CNs. On a DN, the value of this parameter is calculated as follows:

dop_limit x 20 x 6 + 24: **dop_limit** indicates the number of CPUs of each DN in the cluster. It is calculated as follows: **dop_limit** = Number of logical CPU cores of a single server/Number of DNs of a single server.

The minimum value is 5000.

If the parameter is set to a large value, GaussDB(DWS) requires more SystemV shared memories or semaphores, which may exceed the maximum default configuration of the OS. In this case, modify the value as needed.

NOTICE

The value of max_connections is related to max_prepared_transactions. Before setting max_connections, ensure that the value of max_prepared_transactions is greater than or equal to that of max_connections. In this way, each session has a prepared transaction in the waiting state.

application_name

Parameter description: Specifies the name of the client program connecting to the database.

Type: USERSET

Value range: a string

Default value: gsql

connection info

Parameter description: Specifies the database connection information, including the driver type, driver version, driver deployment path, and process owner. (This is an O&M parameter. Do not configure it by yourself.)

Type: USERSET

Value range: a string

Default value: an empty string

- An empty string indicates that the driver connected to the database does not support
 automatic setting of the connection_info parameter or the parameter is not set by
 users in applications.
- The following is an example of the concatenated value of **connection_info**: {"driver_name":"ODBC","driver_version": "(GaussDB x.x.x build 39137c2d) compiled at 2022-09-23 15:43:11 commit 3629 last mr 5138 debug","driver_path":"/usr/local/lib/psqlodbcw.so","os_user":"omm"}

By default, driver_name, driver_version, driver_path, and os_user are displayed for ODBC, JDBC, and gsql connections. For other connections, driver_name and driver_version are displayed by default. The display of driver_path and os_user is controlled by users (see Connecting to a Database and Configuring a Data Source in the Linux OS).

15.4.2 Security and Authentication (postgresql.conf)

This section describes parameters about how to securely authenticate the client and server.

session timeout

Parameter description: Specifies the maximum idle time without any operations after a connection to the server is established.

Type: USERSET

Value range: an integer ranging from 0 to 86400. The minimum unit is second (s). **0** means to disable the timeout.

Default value: 10 min

NOTICE

- The gsql client of GaussDB(DWS) has an automatic reconnection mechanism. If the initialized local connection of a user to the server times out, gsql disconnects from and reconnects to the server.
- Connections from the pooler connection pool to other CNs and DNs are not controlled by the session_timeout parameter.

ssl_renegotiation_limit

Parameter description: Specifies the traffic volume over the SSL-encrypted channel before the session key is renegotiated. The renegotiation traffic limitation mechanism reduces the probability that attackers use the password analysis method to crack the key based on a huge amount of data but causes big performance losses. The traffic indicates the sum of sent and received traffic.

Type: USERSET

You are advised to retain the default value, that is, disable the renegotiation mechanism. You are not advised to use the **gs_guc** tool or other methods to set the **ssl_renegotiation_limit** parameter in the **postgresql.conf** file. The setting does not take effect.

Value range: an integer ranging from 0 to **INT_MAX**. The unit is KB. **0** indicates that the renegotiation mechanism is disabled.

Default value: 0

failed_login_attempts

Parameter description: Specifies the maximum number of incorrect password attempts before an account is locked. The account will be automatically unlocked after the time specified in **password_lock_time**. For example, incorrect password attempts during login and password input failures when using the **ALTER USER** command

Type: SIGHUP

Value range: an integer ranging from 0 to 1000

- **0** indicates that the automatic locking function does not take effect.
- A positive number indicates that an account is locked when the number of incorrect password attempts reaches the value of **failed_login_attempts**.

Default value: 10

NOTICE

- The locking and unlocking functions take effect only when the values of **failed_login_attempts** and **password_lock_time** are positive numbers.
- failed_login_attempts works with the SSL connection mode of the client to
 identify the number of incorrect password attempts. If PGSSLMODE is set to
 allow or prefer, two connection requests are generated for a password
 connection request. One request attempts an SSL connection, and the other
 request attempts a non-SSL connection. In this case, the number of incorrect
 password attempts perceived by the user is the value of failed_login_attempts
 divided by 2.

15.4.3 Communication Library Parameters

This section describes parameter settings and value ranges for communication libraries.

tcp_keepalives_idle

Parameter description: Specifies the interval between keepalive signal sending in an OS that supports the **TCP_KEEPIDLE** socket parameter. If no keepalive signal is transmitted, the connection is in idle state.

Type: USERSET

NOTICE

- If the OS does not support the TCP_KEEPIDLE parameter, set this parameter to 0.
- The parameter is ignored on the OS where connections are established using the Unix domain socket.

Value range: an integer ranging from 0 to 3600. The unit is second (s).

Default value: 0

tcp_keepalives_interval

Parameter description: Specifies the response time before retransmission when the OS supports the **TCP_KEEPINTVL** socket parameter.

Type: USERSET

Value range: an integer ranging from 0 to 180. The unit is second (s).

Default value: 0

NOTICE

- If the OS does not support the **TCP_KEEPINTVL** parameter, set this parameter to **0**.
- The parameter is ignored on the OS where connections are established using the Unix domain socket.

tcp_keepalives_count

Parameter description: Specifies the number of keepalived signals that can be waited before the GaussDB(DWS) server is disconnected from the client if the OS supports the **TCP_KEEPCNT** socket parameter.

Type: USERSET

NOTICE

- If the OS does not support the TCP_KEEPCNT parameter, set this parameter to 0.
- The parameter is ignored on the OS where connections are established using the Unix domain socket.

Value range: an integer ranging from 0 to 100. **0** indicates that the connection is immediately broken if GaussDB(DWS) does not receive a keepalived signal from the client.

Default value: 0

comm_max_datanode

Parameter description: Specifies the maximum number of DNs supported by the communication library.

Type: USERSET

Value range: an integer ranging from 1 to 8192

Default value: actual number of DNs

NOTICE

Increasing this parameter value takes effect immediately, while decreasing the value takes effect after the cluster is restarted.

comm max stream

Parameter description: maximum number of logical connection data structures cached in the communication library.

Type: SIGHUP

Value range: an integer ranging from 1 to 65535

Default value: 1024

Ⅲ NOTE

If the value of **comm_max_datanode** is small, the process memory is sufficient. In this case, you can increase the value of **comm_max_stream**.

max_stream_pool

Parameter description: Specifies the maximum number of stream threads that can be contained in a stream thread pool. This feature is supported in 8.1.2 or later.

Type: SUSET

Value range: an integer ranging from -1 to INT_MAX. The values **-1** and **0** indicate that the stream thread pool is disabled.

Default value:

- The formula for a new cluster is max_stream_pool=MIN(max_connections, max_process_memory/16/5MB, 1024).
- The formula for a cluster upgraded from versions earlier than 8.3.0.100 is max_stream_pool = MIN(max_connections, max_process_memory/16/5MB, 1024, value-of-the-old-cluster). During the upgrade, the settings for the new cluster are forcibly used. The old value is used if it is smaller.

□ NOTE

- The number of stream threads in a thread pool can be reduced in real time. If the value of this parameter is increased, the number of stream threads is increased to meet the service requirements.
- Generally, you are advised not to change the value of this parameter because the stream thread pool supports the automatic cleanup function.
- If too many idle stream threads occupy the memory, you can decrease the value of this parameter to save the memory.

enable_stream_sync_quit

Parameter description: whether the stream threads exit synchronously when the stream plan ends. This parameter is supported only by clusters of version 8.3.0 or later.

Type: USERSET

Value range: Boolean

• **on** indicates that threads in the stream thread group exit after the stream plan ends.

• **off** indicates that stream threads exit directly after the stream plan ends without waiting for the threads in the stream thread group to exit.

Default value: off

enable_connect_standby

Parameter description: Sets the connection between a CN and a standby DN. This parameter is supported only by clusters of version 8.3.0 or later.

Type: USERSET

Value range: Boolean

- **on** indicates that the CN connects to the standby server.
- **off** indicates that the CN connects to the primary DN.

Default value: off

<u>A</u> CAUTION

- You are not advised to use this parameter in routine services. This parameter
 applies only to O&M operations. You are not advised to use the gs_guc tool for
 global settings. Otherwise, problems such as data inconsistency and result set
 errors may occur.
- Enabling this parameter for a session with temporary tables will delete the temporary table data on DNs and prevent further actions on those tables.

comm quota size

Parameter description: Specifies the maximum size of packets that can be continuously sent by the communication library. When you use a 1GE NIC, a small value ranging from 20 KB to 40 KB is recommended.

Type: USERSET

Value range: an integer ranging from 0 to 102400. The default unit is KB. The value **0** indicates that the quota mechanism is not used.

Default value: 1 MB

comm_usable_memory

Parameter description: Specifies the maximum memory that can be used by the communication library cache on a single DN.

Type: SIGHUP

Value range: an integer ranging from 1 to 256. The default unit is KB. The minimum size cannot be less than 1 GB for installation.

Default value: max_process_memory/8

NOTICE

This parameter must be specifically set based on environment memory and the deployment method. If it is too large, out-of-memory (OOM) may occur. If it is too small, the performance of the communication library may deteriorate.

comm_client_bind

Parameter description: Specifies whether to bind the client of the communication library to a specified IP address when the client initiates a connection.

Type: USERSET

Value range: Boolean

- on indicates that the client is bound to a specified IP address.
- off indicates that the client is not bound to any IP addresses.

NOTICE

If multiple IP addresses of a node in a cluster are on the same communication network segment, set this parameter to **on**. In this case, the client is bound to the IP address specified by **listen_addresses**. The concurrency performance of a cluster depends on the number of random ports because a port can be used only by one client at a time.

Default value: off

comm_no_delay

Parameter description: Specifies whether to use the **NO_DELAY** attribute of the communication library connection. Restart the cluster for the setting to take effect.

Type: USERSET

Value range: Boolean

Default value: off

NOTICE

If packet loss occurs because a large number of packets are received per second, set this parameter to **off** to reduce the total number of packets.

comm_debug_mode

Parameter description: Specifies the debug mode of the communication library, that is, whether to print logs about the communication layer. The setting is effective at the session layer.

NOTICE

When the switch is set to **on**, the number of printed logs is huge, adding extra overhead and reducing database performance. Therefore, set the switch to **on** only in the debug mode.

Type: USERSET

Value range: Boolean

- **on** indicates the detailed debug log of the communication library is printed.
- **off** indicates the detailed debug log of the communication library is not printed.

Default value: off

comm_ackchk_time

Parameter description: Specifies the duration after which the communication library server automatically triggers ACK when no data package is received.

Type: USERSET

Value range: an integer ranging from 0 to 20000. The unit is millisecond (ms). **0** indicates that automatic ACK triggering is disabled.

Default value: 2000

comm timer mode

Parameter description: Specifies the timer mode of the communication library, that is, whether to print timer logs in each phase of the communication layer. The setting is effective at the session layer.

NOTICE

When the switch is set to **on**, the number of printed logs is huge, adding extra overhead and reducing database performance. Therefore, set the switch to **on** only in the debug mode.

Type: USERSET

Value range: Boolean

- **on** indicates the detailed timer log of the communication library is printed.
- **off** indicates the detailed timer log of the communication library is not printed.

Default value: off

comm_stat_mode

Parameter description: Specifies the stat mode of the communication library, that is, whether to print statistics about the communication layer. The setting is effective at the session layer.

NOTICE

When the switch is set to **on**, the number of printed logs is huge, adding extra overhead and reducing database performance. Therefore, set the switch to **on** only in the debug mode.

Type: USERSET

Value range: Boolean

- **on** indicates the statistics log of the communication library is printed.
- **off** indicates the statistics log of the communication library is not printed.

Default value: off

client connection check interval

Parameter description: Specifies the interval for checking the client connection status. This parameter is supported by clusters of version 8.2.0 or later.

Type: USERSET

Value range: an integer ranging from 0 to INT_MAX. The unit is ms. The value **0** indicates that the client connection status is not checked.

Default value: 10000

NOTICE

During a long query executed in a session where a client (such as gsql, JDBC, or ODBC) directly connects to the CN,

- The CN checks the client connection status at the interval specified by client_connection_check_interval. If it detects that the client has been disconnected from the CN, the server terminates the long query and releases related resources to avoid waste of cluster resources.
- The DN checks its connection to the CN at the interval specified by client_connection_check_interval. If the DN detects that it has been disconnected from the CN, it terminates the long query and releases related resources to avoid waste of cluster resources.

conn_recycle_timeout

Parameter description: the interval for reclaiming idle connections between a CN and other nodes to the connection pool. This parameter is supported only by clusters of version 8.2.1 or later.

Type: USERSET

Value range: an integer ranging from 0 to 3600, in second (s). **0** indicates that the

function of reclaiming idle connections is disabled.

Default value: 30

15.5 Resource Consumption

15.5.1 **Memory**

This section describes memory parameters.

NOTICE

Parameters described in this section take effect only after the database service restarts.

max_process_memory

Parameter description: Specifies the maximum physical memory of a database node.

Type: SIGHUP

Value range: an integer ranging from 2 x 1024 x 1024 to INT_MAX/2. The unit is KB.

Default value: Determined based on non-secondary DNs. If multiple DNs are deployed on a server, the value is (Physical memory size) \times 0.8/(1 + Number of primary DNs). If a single DN is deployed on a server, the value is (Physical memory size) \times 0.6. If the calculation result is less than 2 GB, the value is 2 GB by default. The default size of the secondary DN is 12 GB.

Setting suggestions:

- On DNs, the value of this parameter is determined based on the physical system memory and the number of DNs deployed on a single node. If multiple DNs are deployed on a server, the calculation formula for the max_process_memory value is as follows: (Physical memory size vm.min_free_kbytes) x 0.8/(n + Number of primary DNs). If only one DN is deployed on a server, the calculation formula for the max_process_memory value is (Physical memory size vm.min_free_kbytes) x 0.6. This parameter aims to ensure system reliability, preventing node OOM caused by increasing memory usage. vm.min_free_kbytes indicates OS memory reserved for kernels to receive and send data. Its value is at least 5% of the total memory. That is, max_process_memory = Physical memory x 0.8/ (n + Number of primary DNs). If the cluster scale (number of nodes in the cluster) is smaller than 256, n=1; if the cluster scale is larger than 256 and smaller than 512, n=2; if the cluster scale is larger than 512, n=3.
- You are not advised to set this parameter to the minimum threshold.

- Set this parameter on CNs to the same value as that on DNs.
- RAM is the maximum memory allocated to the cluster.
- In GaussDB(DWS) 8.2.0 and later versions, the initial value of max_process_memory is increased to improve memory resource utilization. However, in an unbalanced cluster where a server has two primary DNs running, using the initial value of max_process_memory may cause OOM. In 8.2.0 and later versions, the max_process_memory parameter is changed to the SIGHUP type and can be manually adjusted. The max_process_memory_auto_adjust parameter is added. If a cluster is unbalanced, its CM will dynamically adjust max_process_memory based on the cluster status. The value of max_process_memory is (Physical memory vm.min_free_kbytes) x 0.8/Number of primary DNs.
- In GaussDB(DWS) 8.2.1 or later, the application scope of dynamically adjusting the value of **max_process_memory** is expanded from clusters where each server has only one DN to all cluster deployment modes.
 - If max_process_memory_auto_adjust is set to on, the value of max_process_memory is dynamically adjusted between the upper limit and the lower limit. The lower limit is calculated as follows: (Physical memory size) x 0.8/(1 + Number of primary DNs). The upper limit is specified by the GUC parameter max_process_memory_balanced. (For details about how to set max_process_memory_balanced, contact technical support.)
 - When the cluster works in load balancing mode, the upper limit of max_process_memory is used to improve the overall memory usage of the node. Compared with earlier versions, the memory usage is improved.
 - When the cluster is not in load balancing mode, the lower limit of max_process_memory is used. The overall memory usage of the node is the same as that in versions earlier than 8.2.1.
 - In upgrade scenarios, to ensure forward compatibility, the system does not set max_process_memory_balanced, and max_process_memory uses the value set before the upgrade by default.

max_process_memory_auto_adjust

Parameter description: Specifies whether to enable automatic adjustment for **max_process_memory** parameter. (This parameter is supported only by cluster versions 8.2.0 and later.) In a cluster where each server only has one DN, if this function is enabled, the CM dynamically adjusts the value of **max_process_memory** on the corresponding DN during an active/standby switchover.

Type: SIGHUP

Value range: Boolean

Default value: on

Suggestion: Set this parameter to **on**. For a cluster where each server only has one DN, the initial value of **max_process_memory** is increased in 8.2.0 and later versions to improve memory resource utilization. However, after a primary/ standby switchover, there will be two primary DNs running on the same server. Using the initial value of **max_process_memory** in this case may cause OOM, and you need to let the CM dynamically adjust the value.

shared_buffers

Parameter description: Specifies the size of shared memory used by GaussDB(DWS). If this parameter is set to a large value, GaussDB(DWS) may require more System V shared memory than the default setting.

Type: POSTMASTER

Value range: an integer ranging from 128 to INT_MAX. The unit is 8 KB.

Changing the value of **BLCKSZ** will result in a change in the minimum value of the **shared buffers**.

Default value: The value of this parameter for CNs is half of that for DNs, which is calculated using the formula **POWER(2,ROUND(LOG(2,max_process_memory/18),0))**. If the maximum value allowed by the OS is smaller than 32 MB, this parameter will be automatically changed to the maximum value allowed by the OS during database initialization.

Setting suggestions:

You are advised to set this parameter for DNs to a value greater than that for CNs, because GaussDB(DWS) pushes its most queries down to DNs.

It is recommended that **shared_buffers** be set to a value less than 40% of the memory. Set it to a large value for row-store tables and a small value for column-store tables. For column-store tables: shared_buffers = (Memory of a single server/ Number of DNs on the single server) $\times 0.4 \times 0.25$

If you want to increase the value of **shared_buffers**, you also need to increase the value of **checkpoint_segments**, because a longer period of time is required to write a large amount of new or changed data.

bulk_write_ring_size

Parameter description: Specifies the size of the ring buffer used for data parallel import.

Type: USERSET

Value range: an integer ranging from 16384 to INT_MAX. The unit is KB.

Default value: 2 GB

Setting suggestions: Increase the value of this parameter on DNs if a huge amount of data is to be imported.

buffer_ring_ratio

Parameter description: ring buffer threshold for parallel data export

Type: USERSET

Value range: integer in the range 1–1000

Default value: 250

□ NOTE

- The default value indicates that the threshold is 250/1000 (a quarter) of shared buffers.
- The minimum value is 1/1000 of the value of **shared_buffers**.
- The maximum value is the value of **shared buffers**.

Setting suggestions: If the cache hit ratio is not as expected during export, you are advised to configure this parameter on DNs.

enable_cstore_ring_buffer

Parameter description: Specifies whether to enable column-store RingBuffer. This parameter is supported only by cluster versions 8.2.0 and later.

Type: USERSET

Value range: Boolean

Default value: off

Suggestion: If workloads have been running for a period of time, a large amount of frequently queried data has been stored in the CStoreBuffer, and you want to query large tables that are rarely accessed, you are advised to enable this function before the query and disable it after the query.

temp_buffers

Parameter description: Specifies the maximum size of local temporary buffers used by each database session.

Type: USERSET

Value range: an integer ranging from 800 to INT_MAX/2. The unit is 8 KB.

Default value: 8 MB

□ NOTE

- This parameter can be modified only before the first use of temporary tables within each session. Subsequent attempts to change the value of this parameter will not take effect on that session.
- Based on the value of temp_buffers, a session allocates temporary buffers as required.
 The cost of setting a large value in sessions that do not require many temporary buffers is only a buffer descriptor. If a buffer is used, 8192 bytes will be consumed for it.

max_prepared_transactions

Parameter description: Specifies the maximum number of transactions that can stay in the **prepared** state simultaneously. If this parameter is set to a large value, GaussDB(DWS) may require more System V shared memory than the default setting.

When GaussDB(DWS) is deployed as an HA system, set this parameter on the standby server to the same value or a value greater than that on the primary server. Otherwise, queries will fail on the standby server.

Type: POSTMASTER

Value range: an integer ranging from 0 to 536870911. The value of CN set to **0** indicates that the prepared transaction feature is disabled.

Default value: 800 for CNs and 5000 for DNs

∩ NOTE

Set this parameter to a value greater than or equal to that of **max_connections** to avoid failures in preparation.

work mem

Parameter description: Specifies the memory capacity to be used by internal sort operations and Hash tables before writing to temporary disk files. Sort operations are used for **ORDER BY**, **DISTINCT**, and merge joins. Hash tables are required for Hash joins as well as Hash-based aggregations and **IN** subqueries.

For a complex query, several sort or Hash operations may be running in parallel; each operation will be allowed to use as much memory as this value specifies. If the memory is insufficient, data is written into temporary files. In addition, several running sessions could be performing such operations concurrently. Therefore, the total memory used may be many times the value of **work_mem**.

Type: USERSET

Value range: an integer ranging from 64 to INT_MAX. The unit is KB.

Default value: 512 MB for small-scale memory and 2 GB for large-scale memory (If max_process_memory is greater than or equal to 30 GB, it is large-scale memory. Otherwise, it is small-scale memory.)

Setting suggestions:

If the physical memory specified by **work_mem** is insufficient, additional operator calculation data will be written into temporary tables based on query characteristics and the degree of parallelism. This reduces performance by five to ten times, and prolongs the query response time from seconds to minutes.

- In complex serial query scenarios, each query requires five to ten associated operations. Set **work_mem** using the following formula: **work_mem** = 50% of the memory/10.
- In simple serial query scenarios, each query requires two to five associated operations. Set **work_mem** using the following formula: **work_mem** = 50% of the memory/5.
- For concurrent queries, use the formula: **work_mem** = **work_mem** in serialized scenario/Number of concurrent SQL statements.

NOTICE

Once memory adaptation is enabled, there is no need to use **work_mem** to optimize operator memory usage after collecting statistics. The system generates a plan for each statement and estimates the memory usage of each operator and the entire statement based on the current workload. The system then schedules the queue based on the workload and the overall memory usage of the statement, which can result in statement queuing in high-concurrency scenarios.

query_mem

Parameter description: Specifies the memory used by query. If the value of **query_mem** is greater than 0, the optimizer adjusts the estimated query memory to this value when generating an execution plan.

Type: USERSET

Value range: **0** or an integer greater than 32 MB. The default unit is KB. If the value is set to a negative value or less than 32 MB, the default value **0** is used. In this case, the optimizer does not adjust the estimated query memory.

Default value: 0

query_max_mem

Parameter description: Specifies the maximum memory that can be used by query. If the value of **query_max_mem** is greater than 0, when generating an execution plan, the optimizer uses this value to set the available memory for operators. If job memory usage exceeds the value of this parameter, an error is reported and the job exits.

Type: USERSET

Value range: **0** or an integer greater than 32 MB. The default unit is KB. If the value is less than 32 MB, the system automatically sets this parameter to the default value **0**. In this case, the optimizer does not limit the memory usage of jobs.

Default value: 0

agg_max_mem

Parameter description: Specifies the maximum memory that can be used by the Agg operator when the number of aggregation columns exceeds 5. This parameter takes effect only if the value of **agg_max_mem** is greater than 0. (This parameter is supported only in 8.1.3.200 and later cluster versions.)

Type: USERSET

Value range: **0** or an integer greater than 32 MB. The default unit is KB. If the value is less than 32 MB, the system automatically sets this parameter to the default value **0**. In this case, the memory usage of the Agg operator is not limited based on the value.

Default value:

- If the current cluster is upgraded from an earlier version to 8.1.3 or later, the value in the earlier version is inherited. The default value is **INT MAX**.
- If the current cluster version is 8.1.3 or later, the default value is **2GB**.

enable_rowagg_memory_control

Parameter description: Specifies the upper limit of the memory used by the rowstore agg operator.

Type: USERSET

Value range: Boolean

- **on** indicates that the memory usage limit of the row-store agg operator is enabled. Setting this parameter to **on** can avoid OOM caused by the row-store agg operator, but may affect the agg performance.
- **off** indicates that the memory usage limit of the row-store agg operator is disabled. If this parameter is set to **off**, the system memory may be unavailable.

Default value: on

maintenance_work_mem

Parameter description: Specifies the maximum size of memory to be used for maintenance operations, such as **VACUUM**, **CREATE INDEX**, and **ALTER TABLE ADD FOREIGN KEY**. This parameter may affect the execution efficiency of VACUUM, VACUUM FULL, CLUSTER, and CREATE INDEX.

Type: USERSET

Value range: an integer ranging from 1024 to INT_MAX. The unit is KB.

Default value: 512 MB for small-scale memory and 2 GB for large-scale memory (If max_process_memory is greater than or equal to 30 GB, it is large-scale memory. Otherwise, it is small-scale memory.)

Setting suggestions:

- You are advised to set this parameter to the same value of work_mem so
 that database dump can be cleared or restored more quickly. In a database
 session, only one maintenance operation can be performed at a time.
 Maintenance is usually performed when there are not much sessions.
- When the Automatic Cleanup process is running, up to autovacuum_max_workers times of this memory may be allocated. Set maintenance_work_mem to a value equal to or larger than the value of work_mem.
- If a large amount of data needs to be processed in the cluster, increase the value of this parameter in sessions.

psort_work_mem

Parameter description: Specifies the memory used for internal sort operations on column-store tables before they are written into temporary disk files. This parameter can be used for inserting tables having a partial cluster key or index, creating a table index, and deleting or updating a table.

Type: USERSET

NOTICE

Multiple running sessions may perform partial sorting on a table at the same time. Therefore, the total memory usage may be several times of the **psort_work_mem** value.

Value range: an integer ranging from 64 to INT_MAX. The unit is KB.

Default value: 512 MB

max_loaded_cudesc

Parameter description: Specifies the number of loaded CuDescs per column when a column-store table is scanned. Increasing the value will improve the query performance and increase the memory usage, particularly when there are many columns in the column tables.

Type: USERSET

Value range: an integer ranging from 100 to INT_MAX/2

Default value: 1024

NOTICE

When the value of **max_loaded_cudesc** is set to a large value, the memory may be insufficient.

max_stack_depth

Parameter description: Specifies the maximum safe depth of GaussDB(DWS) execution stack. The safety margin is required because the stack depth is not checked in every routine in the server, but only in key potentially-recursive routines, such as expression evaluation.

Type: SUSET

Take the following into consideration when setting this parameter:

- The ideal value of this parameter is the maximum stack size enforced by the kernel (value of **ulimit -s**).
- Setting this parameter to a value larger than the actual kernel limit means that a running recursive function may crash an individual backend process. In an OS where GaussDB(DWS) can check the kernel limit, such as the SLES, GaussDB(DWS) will prevent this parameter from being set to a value greater than the kernel limit.
- Since not all the OSs provide this function, you are advised to set a specific value for this parameter.

Value range: an integer ranging from 100 to INT_MAX. The unit is KB.

Default value: 2 MB

■ NOTE

2 MB is a small value and will not incur system breakdown in general, but may lead to execution failures of complex functions.

cstore buffers

Parameter description: Specifies the size of the shared buffer used by ORC, Parquet, or CarbonData data of column-store tables and OBS or HDFS column-store foreign tables.

Type: POSTMASTER

Value range: an integer ranging from 16384 to INT MAX. The unit is KB.

Default value: The value of this parameter for CNs is 32 MB, while that for DNs is calculated using the formula **POWER(2,ROUND(LOG(2,max_process_memory/18),0))**.

Setting suggestions:

Column-store tables use the shared buffer specified by **cstore_buffers** instead of that specified by **shared_buffers**. When column-store tables are mainly used, reduce the value of **shared_buffers** and increase that of **cstore_buffers**.

Use **cstore_buffers** to specify the cache of ORC, Parquet, or CarbonData metadata and data for OBS or HDFS foreign tables. The metadata cache size should be 1/4 of **cstore_buffers** and not exceed 2 GB. The remaining cache is shared by column-store data and foreign table column-store data.

dfs_max_memory

Parameter description: Specifies the maximum memory that can be occupied during ORC export. If the memory is insufficient when a wide table is exported, increase the value of this parameter and try again. This parameter is supported only by clusters of version 8.3.0 or later.

Type: USERSET

Value range: an integer ranging from 131072 to 10485760. The unit is KB.

Default value: 262144 KB (256 MB)

schedule splits threshold

Parameter description: Specifies the maximum number of files that can be stored in memory when you schedule an HDFS foreign table. If the number is exceeded, all files in the list will be spilled to disk for scheduling.

Type: USERSET

Value range: an integer ranging from 1 to INT_MAX

Default value: 60000

bulk_read_ring_size

Parameter description: Specifies the ring buffer size used for data parallel export.

Type: USERSET

Value range: an integer ranging from 256 to INT_MAX. The unit is KB.

Default value: 16 MB

check cu size threshold

Parameter description: When inserting data into a column-store table, if the amount of data already inserted in a CU exceeds the value of this parameter, row-level size verification will be performed to prevent the creation of uncompressed CUs larger than 1 GB.

Type: USERSET

Value range: an integer ranging from 0 to 1048576. The unit is KB.

Default value: 1 GB

NOTICE

If row-level size verification fails multiple times, you are advised to temporarily set the parameter to **0** at the session level.

memory_spread_strategy

Parameter description: Specifies the DN memory expansion policy of a customized resource pool. You are advised to set this parameter for services with sufficient memory. After this parameter is set, the query performance is improved. The maximum memory can be the same as that of the default resource pool. However, there may be errors caused by insufficient memory in some service scenarios if the memory is small. This parameter is supported only by clusters of version 8.1.3 or later.

Type: USERSET

Value range: enumerated values

- **none**: indicates that the memory is not expanded.
- **negative**: indicates that the memory and operators are expanded based on the estimated usage.
- crazy: indicates that the memory is directly expanded, which is equivalent to
 the memory expansion policy of the default resource pool. However, there
 may be errors caused by insufficient memory in some service scenarios if the
 memory is small.

Default value: none

async_io_acc_max_memory

Parameter description: Specifies the maximum memory that can be used for asynchronous read/write acceleration in a single task thread. This parameter is supported only by clusters of version 9.0.0 or later.

Type: USERSET

Value range: an integer ranging from 4096 to INT_MAX/2, in KB.

Default value: 128MB

15.5.2 Statement Disk Space Control

This section describes parameters related to statement disk space control, which are used to limit the disk space usage of statements.

sql_use_spacelimit

Parameter description: Specifies the allowed maximum space for files to be spilled to disks in a single SQL statement on a single DN. This parameter limits the space occupied by ordinary tables, temporary tables, and intermediate result sets spilled to disks. System administrators are also restricted by this parameter.

Type: USERSET

Value range: an integer ranging from -1 to INT_MAX. The unit is KB. **-1** indicates no limit.

Default value: Set **sql_use_spacelimit** to 10% of the total disk space of the instance.

□ NOTE

For example, if **sql_use_spacelimit** is set to **100** in the statement, and the data spilled to disks on a single DN exceeds 100 KB, DWS will stop the query and display a message indicating threshold exceeded.

insert into user1.t1 select * from user2.t1;

ERROR: The space used on DN (104 kB) has exceeded the sql use space limit (100 kB).

Handling suggestion:

- Optimize the statement to reduce the data spilled to disks.
- If the disk space is sufficient, increase the value of this parameter.

temp_file_limit

Parameter description: Specifies the total space for files spilled to disks in a single thread. The temporary file can be the one used by sorting or hash tables, or cursors in a session.

This is a session-level setting.

Type: SUSET

Value range: an integer ranging from -1 to INT_MAX. The unit is KB. **-1** indicates no limit.

Default value: Set **temp_file_limit** to 10% of the total disk space of the instance.

NOTICE

This parameter does not apply to disk space occupied by temporary tablespaces used for executing SQL queries.

bi_page_reuse_factor

Parameter description: Specifies the percentage of idle space of old pages that can be reused when page replication is used for data synchronization between

primary and standby DNs in the scenario where data is inserted into row-store tables in batches.

Type: USERSET

Value range: an integer ranging from 0 to 100. The value is a percentage. Value **0** indicates that the old pages are not reused and new pages are requested.

Default value: 70

NOTICE

- You are not advised to set this parameter to a value less than **50** (except **0**). If the idle space of the reused page is small, too much old page data will be transmitted between the primary and standby DNs. As a result, the batch insertion performance deteriorates.
- You are not advised to set this parameter to a value greater than 90. If this
 parameter is set to a value greater than 90, idle pages will be frequently
 queried, but old pages cannot be reused.

15.5.3 Kernel Resources

This section describes kernel resource parameters. Whether these parameters take effect depends on OS settings.

max_files_per_node

Parameter description: Specifies the maximum number of files that can be opened by a single SQL statement on a single node. Generally, you do not need to set this parameter.

Type: SUSET

Value range: an integer ranging from **-1** to **INT_MAX**. The value **-1** indicates that the maximum number is limited.

Default value: 50000

Ⅲ NOTE

- For a newly installed cluster of 9.1.0 or later, the default value of this parameter is 50000
- In an upgrade scenario, if the original cluster supports the max_files_per_node
 parameter, the default value of this parameter remains the same for forward
 compatibility.
- In the upgrade scenario, if the original cluster does not support the max_files_per_node parameter, the default value of this parameter is -1 after the upgrade.
- If error message "The last file name is [%s] and %d files have already been opened on data node [%s] with a maximum of %d files." is displayed during statement execution, increase the value of max_files_per_node.

15.5.4 Cost-based Vacuum Delay

The purpose of cost-based vacuum delay is to allow administrators to reduce the I/O impact of **VACUUM** and **ANALYZE** statements on concurrently active

databases. For example, when maintenance statements such as **VACUUM** and **ANALYZE** do not need to be executed quickly and do not interfere with other database operations, administrators can use this function to achieve this purpose.

NOTICE

Certain operations hold critical locks and should be complete as quickly as possible. In GaussDB(DWS), cost-based vacuum delays do not take effect during such operations. To avoid uselessly long delays in such cases, the actual delay is calculated as follows and is the maximum value of the following calculation results:

- vacuum_cost_delay*accumulated_balance/vacuum_cost_limit
- vacuum_cost_delay*4

During the execution of the ANALYZE | ANALYSE and VACUUM statements, the system maintains an internal counter that keeps track of the estimated cost of the various I/O operations that are performed. When the accumulated cost reaches a limit (specified by **vacuum_cost_limit**), the process performing the operation will sleep for a short period of time (specified by **vacuum_cost_delay**). Then, the counter resets and the operation continues.

By default, this feature is disabled. To enable this feature, set **vacuum_cost_delay** to a value other than 0.

vacuum_cost_delay

Parameter description: Specifies the length of time that the process will sleep when **vacuum_cost_limit** has been exceeded.

Type: USERSET

Value range: an integer ranging from 0 to 100. The unit is millisecond (ms). A positive number enables cost-based vacuum delay and **0** disables cost-based vacuum delay.

Default value: 0

NOTICE

- On many systems, the effective resolution of sleep length is 10 ms. Therefore, setting this parameter to a value that is not a multiple of 10 has the same effect as setting it to the next higher multiple of 10.
- This parameter is set to a small value, such as 10 or 20 milliseconds. Adjusting vacuum's resource consumption is best done by changing other parameters.

vacuum_cost_page_hit

Parameter description: Specifies the estimated cost for vacuuming a buffer found in the shared buffer. It represents the cost to lock the buffer pool, look up the shared Hash table, and scan the page.

Type: USERSET

Value range: an integer ranging from 0 to 10000. The unit is millisecond (ms).

Default value: 1

vacuum_cost_page_miss

Parameter description: Specifies the estimated cost for vacuuming a buffer read from the disk. It represents the cost to lock the buffer pool, look up the shared Hash table, read the desired block from the disk, and scan the block.

Type: USERSET

Value range: an integer ranging from 0 to 10000. The unit is millisecond (ms).

Default value: 2

vacuum_cost_page_dirty

Parameter description: Specifies the estimated cost charged when vacuum modifies a block that was previously clean. It represents the I/Os required to flush the dirty block out to disk again.

Type: USERSET

Value range: an integer ranging from 0 to 10000. The unit is millisecond (ms).

Default value: 20

vacuum cost limit

Parameter description: Specifies the cost limit. The cleanup process will sleep if this limit is exceeded.

Type: USERSET

Value range: an integer ranging from 1 to 10000. The unit is ms.

Default value: 200

15.5.5 Asynchronous I/O Operations

enable_adio_debug

Parameter description: Specifies whether to enable logging related to ADIO, which helps to locate ADIO-related issues. General users are not advised to set this O&M parameter.

Type: SUSET

Value range: Boolean

- on or true indicates the log switch is enabled.
- **off** or **false** indicates the log switch is disabled.

Default value: off

enable fast allocate

Parameter description: Specifies whether the quick allocation switch of the disk space is enabled. This switch can be enabled only in the XFS file system.

Type: SUSET

Value range: Boolean

- on or true indicates that this function is enabled.
- **off** or **false** indicates that the function is disabled.

Default value: off

prefetch_quantity

Parameter description: Specifies the number of row-store prefetches using the

ADIO.

Type: USERSET

Value range: an integer ranging from 1024 to 1048576. The unit is 8 KB.

Default value: 32 MB

backwrite_quantity

Parameter description: Specifies the number of row-store writes using the ADIO.

Type: USERSET

Value range: an integer ranging from 1024 to 1048576. The unit is 8 KB.

Default value: 8MB

cstore_prefetch_quantity

Parameter description: Specifies the number of column-store prefetches using

the ADIO.

Type: USERSET

Value range: an integer. The value range is from 1024 to 1048576 and the unit is

KB.

Default value: 32 MB

cstore_backwrite_quantity

Parameter description: Specifies the number of column-store writes using the

ADIO.

Type: USERSET

Value range: an integer. The value range is from 1024 to 1048576 and the unit is

ΚR

Default value: 8MB

cstore_backwrite_max_threshold

Parameter description: Specifies the maximum number of column-store writes buffered in the database using the ADIO.

Type: USERSET

Value range: an integer ranging from 4096 to INT_MAX/2, in KB

Default value: 2 GB

fast_extend_file_size

Parameter description: Specifies the disk size that the row-store pre-scales using

the ADIO.

Type: SUSET

Value range: an integer. The value range is from 1024 to 1048576 and the unit is

KB.

Default value: 8MB

effective_io_concurrency

Parameter description: Specifies the number of requests that can be simultaneously processed by the disk subsystem. For the RAID array, the parameter value must be the number of disk drive spindles in the array.

Type: USERSET

Value range: an integer ranging from 0 to 1000

Default value: 1

cu preload_max_distance

Parameter description: Specifies the maximum number of CU groups that can be prefetched. This is supported only by clusters of version 9.1.0.100 or later.

Type: USERSET

Value range: an integer ranging from 0 to 1024. The value 0 indicates that

prefetching is disabled.

Default value: 20

cu_preload_count

Parameter description: Specifies the maximum number of CUs to be prefetched. This parameter is supported only by clusters of version 9.1.0 or later.

Type: USERSET

Value range: an integer. The value ranges from 0 to 10000. The value **0** indicates

that prefetching is disabled.

Default value: 600

15.5.6 Disk Caching

The following parameters are supported only by clusters of version 9.1.0 or later.

enable disk cache

Parameter description: Specifies whether to enable file caching. Setting this parameter to **on** only takes effect when **enable_aio_scheduler** is set to **on** and **obs_worker_pool_size** is greater than or equal to 4.

Type: USERSET

Value range: Boolean

Default value: off

enable_disk_cache_recovery

Parameter description: Specifies whether file caching can be restored when the cluster is restarted.

Type: USERSET

Value range: Boolean

Default value: off

15.6 Parallel Data Import

GaussDB(DWS) provides a parallel data import function that enables a large amount of data to be imported in a fast and efficient manner. This section describes parameters for importing data in parallel in GaussDB(DWS).

raise_errors_if_no_files

Parameter description: Specifies whether distinguish between the problems "the number of imported file records is empty" and "the imported file does not exist". If set to **TRUE**, GaussDB(DWS) reports the error "file does not exist" when the issue "the imported file does not exist" occurs.

Type: SUSET

Value range: Boolean

- on indicates the messages of "the number of imported file records is empty" and "the imported file does not exist" are distinguished when files are imported.
- **off** indicates the messages of "the number of imported file records is empty" and "the imported file does not exist" are not distinguished when files are imported.

Default value: off

partition_max_cache_size

Parameter description: To optimize the inserting of column-store partitioned tables in batches, data is cached during the inserting process and then written to

the disk in batches. You can use **partition_max_cache_size** to specify the size of the data buffer. If the value is too large, much memory will be consumed. If it is too small, the performance of inserting column-store partitioned tables in batches will deteriorate.

Type: USERSET

Value range: 4096 to INT_MAX/2. The minimum unit is KB.

Default value: 2 GB

partition_mem_batch

Parameter description: To optimize the performance of batch insert into column-store partitioned tables, data is cached during the inserting process and then written to the disk in batches. If partition_max_cache_size is configured, partition_mem_batch can be used to specify the number of caches. If this parameter is set to a large value, the available cache of each partition will be small, and the performance of batch insert into column-store partitioned tables will deteriorate. If this parameter is set to a small value, the available cache of each partition will be large, consuming much system memory.

Type: USERSET

Value range: 1 to 65535

Default value: 256

gds debug mod

Parameter description: Specifies whether to enable the debug function of Gauss Data Service (GDS). This parameter is used to better locate and analyze GDS faults. After the debug function is enabled, types of packets received or sent by GDS, peer end of GDS during command interaction, and other interaction information about GDS are written into the logs of corresponding nodes. In this way, state switching on the GaussDB state machine and the current state are recorded. If this function is enabled, additional log I/O resources will be consumed, affecting log performance and validity. You are advised to enable this function only when locating GDS faults.

Type: USERSET

Value range: Boolean

- on indicates that the GDS debug function is enabled.
- **off** indicates that the GDS debug function is disabled.

Default value: off

max_copy_data_display

Parameter description: GUC control added for the length of the **rawrecord** field in the copy error table, in the text type. The maximum value is 1 GB minus 8203 bytes (that is, 1073733621 bytes). This parameter is supported only by clusters of version 8.2.1.100 or later.

When this parameter is set, it indicates the maximum number of characters that can be displayed. If the number of characters exceeds the maximum, an ellipsis (...) is displayed at the end.

Type: USERSET

Value range: 0 to 1073733616

Default value: 1024

enable_parallel_batch_insert

Parameter description: Specifies whether to enable the concurrent import function for row-store and column-store tables. This parameter is available only for clusters of version 8.2.0.102 or later.

Type: USERSET

Value range: Boolean

- **on** indicates that concurrent import is enabled. Enabling concurrent data import can improve the importing speed of insert, upsert, update, and delete statements executed on row-store and column-store tables.
- **off** indicates that concurrent import is **disabled**.

Default value: off

15.7 Write Ahead Logs

15.7.1 Settings

synchronous commit

Parameter description: Specifies the synchronization mode of the current transaction.

Type: USERSET

Value range: enumerated values

- **on** indicates synchronization logs of a standby server are updated to disks.
- off indicates asynchronous commit.
- local indicates local commit.
- remote_write indicates synchronization logs of a standby server are written to disks.
- **remote_receive** indicates synchronization logs of a standby server are required to receive data.

Default value: on

wal skip fpw level

Parameter description: Specifies whether to write the whole page to WALs when a page is modified for the first time after a checkpoint on a DN. This parameter

controls if FPWs of different kinds are logged in WAL to avoid space shortage or slow synchronization between primary and standby nodes due to too many WAL logs. This parameter is supported by cluster versions 8.3.0 and later. For versions earlier than 8.3.0, see **write fpi hint**.

Type: USERSET

Value range: an integer

- 0: enables All FPWs.
- 1: specifies that when setting hintbits for the heap tuples is the first change since checkpoint, FPW logs are not written. Whether this function takes effect is not affected by the enable_crc_check and wal_log_hints parameters. If this parameter is set to 1, it is equivalent to write_fpi_hint=off in an earlier version. For details, see write_fpi_hint.
- 2: FPW logs of indexes and FPW logs after the hintbits on the heap page is modified are not recorded.

Default value: 1

wal_decelerate_policy

Parameter description: Specifies the behavior policy after rate limiting is triggered. This parameter is supported only by clusters of 8.2.0 and later versions.

Type: USERSET

Value range: enumerated values

- warning indicates that an alarm is generated but the rate is not limited.
- **decelerate** indicates that the rate will be limited based on policy settings.

Default value: warning

Setting the parameter to **warning** does not affect performance. Setting it to **decelerate** will limit the rate based on policy settings if the rate exceeds the threshold.

wal_write_speed

Parameter description: Specifies the maximum WAL write speed (byte/s) allowed by each query on a single DN. This parameter is supported only by clusters of 8.2.0 or later.

Type: USERSET

Value range: an integer ranging from 1024 to 10240000, in KB.

Default value: 30MB

□ NOTE

The rate of a large number of jobs with index copy and deletion operations will be limited.

wal_decelerate_trigger_threshold

Parameter description: Specifies the threshold of WAL write rate limiting for each query on a single DN. This parameter is supported only by cluster versions 8.2.0 and later.

Type: USERSET

Value range: an integer ranging from 1024 to 10000000000, in KB.

Default value: 128MB

■ NOTE

This function is triggered only if the number of Xlogs generated by a single query is greater than the value of this parameter. DDL operations or a small number of DML operations are not affected.

commit_delay

Parameter description: Specifies the duration of committed data be stored in the WAL buffer.

Type: USERSET

Value range: an integer, ranging from 0 to 100000 (unit: μs). **0** indicates no delay.

Default value: 0

NOTICE

- When this parameter is set to a value other than 0, the committed transaction is stored in the WAL buffer instead of being written to the WAL immediately. Then, the WalWriter process flushes the buffer out to disks periodically.
- If system load is high, other transactions are probably ready to be committed within the delay. If no transactions are waiting to be submitted, the delay is a waste of time.

commit_siblings

Parameter description: Specifies a limit on the number of ongoing transactions. If the number of ongoing transactions is greater than the limit, a new transaction will wait for the period of time specified by **commit_delay** before it is submitted. If the number of ongoing transactions is less than the limit, the new transaction is immediately written into a WAL.

Type: USERSET

Value range: an integer ranging from 0 to 1000

Default value: 5

wal compression

Parameter description: Specifies whether to compress FPI pages.

Type: USERSET

Value range: Boolean

on: enable the compressionoff: disable the compression

Default value: on

NOTICE

- Only zlib compression algorithm is supported.
- For clusters that are upgraded to the current version from an earlier version, this parameter is set to **off** by default. You can run the **gs_guc** command to enable the FPI compression function if needed.
- If the current version is a newly installed version, this parameter is set to **on** by default.
- If this parameter is manually enabled for a cluster upgraded from an earlier version, the cluster cannot be rolled back.

wal_compression_level

Parameter description: Specifies the compression level of zlib compression algorithm when the **wal_compression** parameter is enabled.

Type: USERSET

Value range: an integer ranging from 0 to 9.

- **0** indicates no compression.
- 1 indicates the lowest compression ratio.
- **9** indicates the highest compression ratio.

Default value: 9

15.7.2 Checkpoints

checkpoint_segments

Parameter description: minimum number of WAL segment files in the period specified by **checkpoint_timeout**. The size of each log file is 16 MB.

Type: SIGHUP

Value range: an integer. The minimum value is 1.

Default value: 64

NOTICE

Increasing the value of this parameter speeds up the import of big data. Set this parameter based on **checkpoint_timeout** and **shared_buffers**. This parameter affects the number of WAL log segment files that can be reused. Generally, the maximum number of reused files in the **pg_xlog** folder is twice the number of checkpoint segments. The reused files are not deleted and are renamed to the WAL log segment files which will be later used.

15.8 HA Replication

15.8.1 Primary Server

enable_data_replicate

Parameter description: Specifies the data synchronization mode between the primary and standby servers when data is imported to row-store tables in a database.

Type: USERSET

Value range: Boolean

- on indicates that data pages are used for the data synchronization between
 the primary and standby servers when data is imported to row-store tables in
 a database. This parameter cannot be set to on if replication_type is set to 1.
- **off** indicates that the primary and standby servers synchronize data using Xlogs while the data is imported to a row-store table.

Default value: on

ha_module_debug

Parameter description: This command is used to view the replication status log of a specific data block during data replication.

Type: USERSET

Value range: Boolean

- **on** indicates that the log records the status of each data block during data replication.
- **off** indicates that the status of each data block is not recorded in logs during data replication.

Default value: off

15.9 Query Planning

15.9.1 Optimizer Method Configuration

These configuration parameters provide a crude method of influencing the query plans chosen by the query optimizer. If the default plan chosen by the optimizer for a particular query is not optimal, a temporary solution is to use one of these configuration parameters to force the optimizer to choose a different plan. Better ways include adjusting the optimizer cost constants, manually running **ANALYZE**, increasing the value of the **default_statistics_target** configuration parameter, and adding the statistics collected in a specific column using **ALTER TABLE SET STATISTICS**.

enable_bitmapscan

Parameter description: Controls whether the query optimizer uses the bitmapscan plan type.

Type: USERSET

Value range: Boolean

- on indicates it is enabled.
- off indicates it is disabled.

Default value: on

enable_hashagg

Parameter description: Controls whether the query optimizer uses the Hash aggregation plan type.

Type: USERSET

Value range: Boolean

- **on** indicates it is enabled.
- off indicates it is disabled.

Default value: on

enable_mixedagg

Parameter description: Controls whether the query optimizer uses the Mixed Agg plan type.

Type: USERSET

Value range: Boolean

- **on** indicates that a Mixed Agg query plan is generated for the Grouping Sets statement (including Rollup or Cube) that meets certain conditions.
- **off** indicates it is disabled.

Default value: on

NOTICE

- The default value of this parameter is **on** in a newly installed cluster of 9.1.0.200 or later. In an upgrade scenario, the default value of this parameter is retained for forward compatibility.
- The Mixed Agg query plan can be used to improve the performance of statements dealing with a large amount of data (the data volume of a single DN table is greater than 100 GB).

Mixed Agg is not supported in the following scenarios:

- The data type of the columns in the **GROUP BY** clause do not support hashing.
- The aggregate function uses **DISTINCT** for deduplication or **ORDER BY** for sorting.
- The **GROUPING SETS** clause does not contain empty groups.

enable_hashjoin

Parameter description: Controls whether the query optimizer uses the Hash-join plan type.

Type: USERSET

Value range: Boolean

on indicates it is enabled.off indicates it is disabled.

Default value: on

enable indexscan

Parameter description: Controls whether the query optimizer uses the index-scan plan type.

Type: USERSET

Value range: Boolean

on indicates it is enabled.off indicates it is disabled.

Default value: on

enable_indexonlyscan

Parameter description: Controls whether the query optimizer uses the indexonly-scan plan type.

Type: USERSET

Value range: Boolean

on indicates it is enabled.

• **off** indicates it is disabled.

Default value: on

enable_material

Parameter description: Controls whether the query optimizer uses materialization. It is impossible to suppress materialization entirely, but setting this parameter to **off** prevents the optimizer from inserting materialized nodes.

Type: USERSET

Value range: Boolean

- on indicates it is enabled.
- off indicates it is disabled.

Default value: on

enable_mergejoin

Parameter description: Controls whether the query optimizer uses the merge-join plan type.

Type: USERSET

Value range: Boolean

- on indicates it is enabled.
- off indicates it is disabled.

Default value: off

enable_nestloop

Parameter description: Controls whether the query optimizer uses the nested-loop join plan type to fully scan internal tables. It is impossible to suppress nested-loop joins entirely, but setting this parameter to **off** allows the optimizer to choose other methods if available.

Type: USERSET

Value range: Boolean

- **on** indicates it is enabled.
- off indicates it is disabled.

Default value: off

enable_index_nestloop

Parameter description: Controls whether the query optimizer uses the nested-loop join plan type to scan the parameterized indexes of internal tables.

Type: USERSET

Value range: Boolean

- **on** indicates the query optimizer uses the nested-loop join plan type.
- off indicates the query optimizer does not use the nested-loop join plan type.

Default value: The default value for a newly installed cluster is **on**. If the cluster is upgraded from R8C10, the forward compatibility is retained. If the version is upgraded from R7C10 or an earlier version, the default value is **off**.

left_join_estimation_enhancement

Parameter description: Specifies whether to use the optimized estimated number of rows for left join. This parameter is supported only by clusters of version 8.3.0.100 or later.

Type: USERSET

Value range: Boolean

- on indicates that the optimized value is used.
- off indicates it is disabled.

Default value: off

enable segscan

Parameter description: Controls whether the query optimizer uses the sequential scan plan type. It is impossible to suppress sequential scans entirely, but setting this variable to **off** allows the optimizer to preferentially choose other methods if available.

Type: USERSET

Value range: Boolean

- on indicates it is enabled.
- off indicates it is disabled.

Default value: on

enable sort

Parameter description: Controls whether the query optimizer uses the sort method. It is impossible to suppress explicit sorts entirely, but setting this variable to **off** allows the optimizer to preferentially choose other methods if available.

Type: USERSET

Value range: Boolean

- on indicates it is enabled.
- off indicates it is disabled.

Default value: on

max_opt_sort_rows

Parameter description: Specifies the maximum number of optimized limit+offset rows in an ORDER BY clause. This parameter is supported only by clusters of version 8.3.0 or later.

Type: USERSET

Value range: an integer ranging from 0 to INT_MAX

- If the value is **0**, the parameter does not take effect.
- If this parameter is set to any other value, the optimization takes effect when
 the number of limit+offset rows in the ORDER BY clause is less than the value
 of this parameter. If the number of limit+offset rows in the order by clause is
 greater than the value of this parameter, the optimization does not take
 effect. After the optimization, the time required is reduced, but the memory
 usage may increase.

Default value: 0

enable tidscan

Parameter description: Controls whether the query optimizer uses the Tuple ID (TID) scan plan type.

Type: USERSET

Value range: Boolean

- on indicates it is enabled.
- off indicates it is disabled.

Default value: on

enable_kill_query

Parameter description: In CASCADE mode, when a user is deleted, all the objects belonging to the user are deleted. This parameter specifies whether the queries of the objects belonging to the user can be unlocked when the user is deleted.

Type: SUSET

Value range: Boolean

- **on** indicates the unlocking is allowed.
- **off** indicates the unlocking is not allowed.

Default value: off

enforce oracle behavior

Parameter description: Controls the rule matching modes of regular expressions.

Type: USERSET

Value range: Boolean

- **on** indicates that the ORACLE matching rule is used.
- **off** indicates that the POSIX matching rule is used.

Default value: on

enable_stream_concurrent_update

Parameter description: Controls the use of **stream** in concurrent updates. This parameter is restricted by the **enable_stream_operator** parameter.

Type: USERSET

Value range: Boolean

- **on** indicates that the optimizer can generate stream plans for the **UPDATE** statement.
- **off** indicates that the optimizer can generate only non-stream plans for the **UPDATE** statement.

Default value: on

enable stream ctescan

Parameter description: Specifies whether a stream plan supports ctescan.

Type: USERSET

Value range: Boolean

- **on** indicates that **ctescan** is supported for the stream plan.
- **off** indicates that **ctescan** is not supported for the stream plan.

Default value: on

In upgrade scenarios, the default value of this parameter is forward compatible, and the original value is retained.

enable_stream_operator

Parameter description: Controls whether the query optimizer uses streams.

Type: USERSET

Value range: Boolean

- on indicates it is enabled.
- **off** indicates it is disabled.

Default value: on

enable stream recursive

Parameter description: Specifies whether to push **WITH RECURSIVE** join queries to DNs for processing.

Type: USERSET

Value range: Boolean

• on: WITH RECURSIVE join queries will be pushed down to DNs.

• off: WITH RECURSIVE join queries will not be pushed down to DNs.

Default value: on

enable value redistribute

Parameter description: Specifies whether to generate value redistribute plans. In 8.2.0 and later cluster versions, this parameter takes effect for **rank**, **dense_rank**, and **row number** without the **PARTITION BY** clause.

Type: USERSET

Value range: Boolean

- on indicates that value redistribute plans are generated.
- off indicates that no value redistribute plans are generated.

Default value: on

max_recursive_times

Parameter description: Specifies the maximum number of **WITH RECURSIVE** iterations.

Type: USERSET

Value range: an integer ranging from 0 to INT_MAX

Default value: 200

enable vector engine

Parameter description: Controls whether the query optimizer uses the vectorized executor.

Type: USERSET

C. OSLINSLI

Value range: Boolean

- **on** indicates it is enabled.
- **off** indicates it is disabled.

Default value: on

enable_broadcast

Parameter description: Controls whether the query optimizer uses the broadcast distribution method when it evaluates the cost of stream.

Type: USERSET

Value range: Boolean

on indicates it is enabled.

off indicates it is disabled.

Default value: on

enable_redistribute

Parameter description: Controls whether the query optimizer uses the local redistribute or split redistribute distribution method when estimating the cost of streams. This parameter is supported only by clusters of version 8.3.0 or later.

Type: USERSET

Value range: Boolean

- on indicates that either of the distribution methods is used.
- **off** indicates that none of the distribution methods is used.

Default value: on

enable_change_hjcost

Parameter description: Specifies whether the optimizer excludes internal table running costs when selecting the Hash Join cost path. If it is set to **on**, tables with a few records and high running costs are more possible to be selected.

Type: USERSET

Value range: Boolean

- **on** indicates it is enabled.
- off indicates it is disabled.

Default value: off

enable_fstream

Parameter description: Controls whether the query optimizer uses streams when it delivers statements. This parameter is only used for external HDFS tables.

This parameter has been discarded. To reserve forward compatibility, set this parameter to **on**, but the setting does not make a difference.

Type: USERSET

Value range: Boolean

- on indicates it is enabled.
- off indicates it is disabled.

Default value: off

enable_hashfilter

Parameter description: Controls whether hashfilters can be generated for plans that contain replication tables (including dual and constant tables). This parameter is supported by clusters of version 8.2.0 or later.

Type: USERSET

Value range: Boolean

- **on** indicates that hashfilters can be generated.
- **off** indicates that no hashfilters can be generated.

Default value: on

best agg plan

Parameter description: The query optimizer generates three plans for the aggregate operation under the stream:

- 1. hashagg+gather(redistribute)+hashagg
- redistribute+hashagg(+gather)
- 3. hashagg+redistribute+hashagg(+gather).

This parameter is used to control the query optimizer to generate which type of hashagg plans.

Type: USERSET

Value range: an integer ranging from 0 to 3.

- When the value is set to **1**, the first plan is forcibly generated.
- When the value is set to **2** and if the **group by** column can be redistributed, the second plan is forcibly generated. Otherwise, the first plan is generated.
- When the value is set to **3** and if the **group by** column can be redistributed, the third plan is generated. Otherwise, the first plan is generated.
- When the value is set to **0**, the query optimizer chooses the most optimal plan based on the estimated costs of the three plans above.

Default value: 0

turbo_engine_version

Parameter description: For tables with the turbo storage format specified during table creation (by setting the **enable_turbo_store** parameter to **on** in the table properties), and when the query does not involve merge join or sort agg operators, the executor can use the turbo execution engine, which can significantly improve performance.

Type: USERSET

Value range: an integer ranging from 0 to 3.

- The value **0** indicates that the turbo execution engine is disabled.
- The value 1 indicates that the turbo execution engine is only used for single-table aggregate queries.
- The value **2** indicates that the turbo execution engine is only used for single-table aggregate or multi-table join queries.
- The value 3 indicates that the turbo execution engine can be used to
 accelerate most commonly used operators, except for operators such as
 merge join and sort agg. When the data volume is large and
 turbo_engine_version is set to 3, the occurrence of merge join and sort agg

operators is relatively rare, so turbo execution engine acceleration can be achieved for almost SQL statements.

Default value: 3



enable bucket stream opt

Parameter description: Specifies whether to use the **bucket agg** and **bucket join** policies for level-2 partitioned tables or 3.0 hash distributed tables. It speeds up SQL statement execution by avoiding local data redistribution or broadcast. This is supported only by clusters of version 9.1.0.200 or later.

Type: USERSET

Value range: Boolean

- true: The optimizer uses the bucket agg and bucket join execution policies
 to generate plans when the conditions for the policy to be applied are met. If
 this optimization policy is used, "Bucket Stream: true" is displayed at the end
 of the EXPLAIN statement.
- **false**: The optimizer does not use the **bucket agg** and **bucket join** execution policies to generate plans.

Default value: true

NOTICE

- The default value of this parameter is **true** in a newly installed cluster of 9.1.0.200 or later. In an upgrade scenario, the default value of this parameter is retained for forward compatibility.
- The bucket agg and bucket join execution policies take effect only when the current query has 16 or fewer available CPUs and meets one of the following conditions:
 - 1. The distribution column for level-2 partitions must match the **secondary_part_column** of these partitions. It is recommended that the number of level-2 partitions be the number of DNs multiplied by 12. Supported multiples include 4, 6, 8, 12, and 16.
 - 2. Tables in version 3.0 must use hash distribution, with the number of buckets or DNs exceeding 10.
- If the local stream cost in the plan is low, the query may not select the **bucket agg** and **bucket join** policies.

spill_compression

Parameter description: Specifies the compression algorithm used when the executor operator runs out of memory and needs to spill data to disk. This is supported only by clusters of version 9.1.0.100 or later.

Type: USERSET

Value range: enumerated values

- 'lz4' indicates that the lz4 compression algorithm is used, which provides better performance for scenarios with smaller spill volumes, but requires more storage space.
- **'zstd'** indicates that the zstd compression algorithm is used, which provides better performance for scenarios with larger spill volumes where I/O is the main bottleneck, and requires approximately 2/3 of the storage space used by Iz4.

Default value: 'lz4'

index_selectivity_cost

Parameter description: Controls the cost calculation of cbtree when scanning column-store table indexes (for selectivity > 0.001). This parameter is only supported by clusters of version 8.2.1.100 or later.

Type: USERSET

Value range: a floating point number, which can be -1 or ranges from 0 to 1000.

- If this parameter is set to 0, the index selection rate is not affected by the threshold 0.001.
- If the value is -1, the value is impacted by disable_cost.
- When it is set to other values, the value is the coefficient for cbtree cost calculation.

Default value: -1

index cost limit

Parameter description: threshold for disabling the cost calculation of cbtree during column-store table index scanning. This parameter is supported only by clusters of version 8.2.1.100 or later.

Type: USERSET

Value range: an integer ranging from 0 to 2147483647

- If the value is **0**, the parameter does not take effect.
- If this parameter is set to other values and the number of rows in a table is less than the value of this parameter, the table is not affected by the index selection rate threshold 0.001.

Default value: 0

volatile_shipping_version

Parameter description: Controls the execution scope of volatile functions to be pushed down.

Type: USERSET

Value range: 0, 1, 2, 3

- When set to 3, it extends the support for pushing down InlineCTE when it is only referenced once, on top of the support provided by a value 2. It also extends the support for pushing down the use of volatile functions in UPSERT operations involving replicated tables.
- When the value is **2**, pushdown can be performed when VOLATILE functions are contained in the target column of the copied CTE result.
- If this parameter is set to 1, the nextval, uuid_generate_v1, sys_guid, and uuid functions can be completely pushed down if they are in the target column of a statement.
- If this parameter is set to 0, random functions can be completely pushed down. The **nextval** and **uuid_generate_v1** functions can be pushed down only if **INSERT** contains simple query statements.

Default value: 3

agg_redistribute_enhancement

Parameter description: When the aggregate operation is performed, which contains multiple group by columns and all of the columns are not in the distribution column, you need to select one **group by** column for redistribution. This parameter controls the policy of selecting a redistribution column.

Type: USERSET

Value range: Boolean

- on indicates the column that can be redistributed and evaluates the most distinct value for redistribution.
- off indicates the first column that can be redistributed for redistribution.

Default value: off

enable_valuepartition_pruning

Parameter description: Specifies whether to perform static or dynamic optimization on the partitioned tables in a distributed file system (DFS).

Type: USERSET

Value range: Boolean

- **on** indicates that the DFS partitioned table is dynamically or statically optimized.
- **off** indicates that the DFS partitioned table is not dynamically or statically optimized.

Default value: on

expected_computing_nodegroup

Parameter description: Specifies a computing Node Group or the way to choose such a group. The Node Group mechanism is now for internal use only. You do not need to set it.

During join or aggregation operations, a Node Group can be selected in four modes. In each mode, the specified candidate computing Node Groups are listed for the optimizer to select an appropriate one for the current operator.

Type: USERSET

Value range: a string

- **optimal**: The list of candidate computing Node Groups consists of the Node Group where the operator's operation objects are located and the DNs in the Node Groups on which the current user has the COMPUTE permission.
- **query**: The list of candidate computing Node Groups consists of the Node Group where the operator's operation objects are located and the DNs in the Node Groups where base tables involved in the query are located.
- bind: If the current session user is a logical cluster user, the candidate computing Node Group is the Node Group of the logical cluster associated with the current user. If the session user is not a logical cluster user, the candidate computing Node Group selection rule is the same as that when this parameter is set to query.
- Node Group name:
 - If enable_nodegroup_debug is set to off, the list of candidate computing Node Groups consists of the Node Group where the operator's operation objects are located and the specified Node Group.
 - If enable_nodegroup_debug is set to on, the specified Node Group is used as the candidate Node Group.

Default value: bind

enable_nodegroup_debug

Parameter description: Specifies whether the optimizer assigns computing workloads to a specific Node Group when multiple Node Groups exist in an environment. The Node Group mechanism is now for internal use only. You do not need to set it.

This parameter takes effect only when **expected_computing_nodegroup** is set to a specific Node Group.

Type: USERSET

Value range: Boolean

- **on** indicates that computing workloads are assigned to the Node Group specified by **expected_computing_nodegroup**.
- **off** indicates no Node Group is specified to compute.

Default value: off

stream_multiple

Parameter description: Specifies the weight used for optimizer to calculate the final cost of stream operators.

The base stream cost is multiplied by this weight to make the final cost.

Type: USERSET

Value range: a floating point number ranging from 0 to 10000

Default value: 1

NOTICE

This parameter is applicable only to Redistribute and Broadcast streams.

qrw_inlist2join_optmode

Parameter description: Specifies whether enable inlist-to-join (inlist2join) query rewriting.

Type: USERSET

Value range: a string

disable: inlist2join disabled

cost_base: cost-based inlist2join query rewriting

• rule_base: forcible rule-based inlist2join query rewriting

• A positive integer: threshold of Inlist2join query rewriting. If the number of elements in the list is greater than the threshold, the rewriting is performed.

Default value: disable

enable_inlist_hashing

Parameter description: Specifies whether to use inlist hash optimization. This parameter is supported only by clusters of version 9.1.0 or later.

Type: USERSET

Value range: Boolean

- **on** indicates that inlist hash optimization is enabled.
- **off** indicates that inlist hash optimization is disabled.

Default value: on

setop_optmode

Parameter description: Specifies whether to perform deduplication on the query branch statements of a set operation (UNION/EXCEPT/INTERSECT) without the ALL option.

Type: USERSET

Value range: enumerated values

- **disable**: The guery branch does not perform deduplication.
- **force**: The query branch forcibly performs deduplication.
- **cost**: The optimizer evaluates the costs of query branches with and without deduplication and selects the execution mode with the lower cost.

Default value: cost

NOTICE

- The default value of this parameter is **cost** in a newly installed cluster of 9.1.0.200 or later. In an upgrade scenario, the default value of this parameter is retained for forward compatibility.
- This parameter takes effect only if the execution plan of a SQL statement meets the following conditions:
 - The **UNION**, **EXCEPT**, and **INTERSECT** operations in the SQL statement do not contain the **ALL** option.
 - Data redistribution has been performed on the query branches where the set operation is to be performed.

skew_option

Parameter description: Specifies whether an optimization policy is used

Type: USERSET

Value range: a string

off: policy disabled

normal: radical policy. All possible skews are optimized.

• **lazy**: conservative policy. Uncertain skews are ignored.

Default value: normal

enable_expr_skew_optimization

Parameter description: Specifies whether to use expression statistics in the skew optimization policy. This is supported only by clusters of version 9.1.0.100 or later.

Type: USERSET

Value range: Boolean

- **on** indicates that expression statistics are used to determine whether data skew occurs in the skew optimization policy.
- **off** indicates that expression statistics are not used to determine whether data skew occurs in the skew optimization policy.

Default value: on

prefer_hashjoin_path

Parameter description: whether to preferentially generate hashjoin paths so that other paths with high costs can be pre-pruned to shorten the overall plan generation time. This parameter is supported only by clusters of version 8.2.1 or later.

Type: USERSET

Value range: Boolean

• **on** indicates that the optimization of generating hash join paths in advance is enabled.

off indicates that the optimization of generating hash join paths in advance is disabled.

Default value: on

enable hashfilter test

Parameter description: whether to add hash filters to columns for base table scan to check whether the results meet expectations. In addition, this parameter determines whether to check the DN accuracy when data is inserted (that is, whether the current data should be inserted into the current DN).

Type: USERSET

Value range: Boolean

- on adds a hash filter for the distribution column to the base table scan and performs accurate DN verification during data insertion.
- off does not add a hash filter for the distribution column to the base table scan and does not perform DN verification during data insertion.

Default value: on

NOTICE

- This parameter is valid only for tables distributed in hash mode.
- If this parameter is set to **on**, DN accuracy is verified during data insertion, affecting data insertion performance.

enable cu align 8k

Parameter description: Specifies whether to set the CUs in V3 tables to 8 KB. This parameter is supported only by clusters of version 9.1.0 or later.

Type: USERSET

Value range: Boolean

- **on** indicates that the CUs in V3 tables are set to 8,192 bytes.
- off indicates that the CUs in V3 tables are set to 512 bytes.

Default value: off

enable_cu_batch_insert

Parameter description: Specifies whether to enable the multi-column CU batch write feature for V2 tables. This parameter is supported only by clusters of version 9.1.0 or later.

Type: USERSET

Value range: Boolean

on indicates that the multi-column CU batch write feature is enabled for V2 tables.

• **off** indicates that the multi-column CU batch write feature is disabled for V2 tables.

Default value: off

enable topk optimization

Parameter description: Specifies whether to enable Top K sorting optimization. This is supported only by clusters of version 9.1.0.200 or later.

Type: USERSET

Value range: Boolean

- on indicates that Top K sorting optimization is enabled.
- off indicates that Top K sorting optimization is disabled.

Default value: on

late_read_strategy

Parameter description: Specifies whether to use the late materialization feature. This is supported only by clusters of version 9.1.0.200 or later.

Type: USERSET

Value range: enumerated values

- **topk**: enables the late materialization optimization method for statements that involve both sorting and limiting.
- none: indicates that the late materialization optimization method is not used.

Default value: topk

insert_dop

Parameter description: Specifies the number of concurrent **INSERT DOP** statements. This parameter is available only for clusters of version 9.1.0.200 or later.

Type: USERSET

Value range: an integer ranging from 0 to 64

- Value **0** indicates that the insert concurrency aligns with the query statement in the subplan.
- Value 1 indicates that **insert dop** is disabled.
- Any value bigger than **1** indicate that the insert concurrency is set to the smaller value between **insert_dop** and the query statements in the subplan.

Default value: 1

- Setting **insert_dop** to **0** activates adaptive mode, where concurrency adjusts based on the subplan's query statements. This approach might not work well in all cases.
- Avoid setting insert_dop to a large value, as it can lead to excessive resource consumption. Adjust insert_dop according to your specific session needs.
- Consider both the insert service concurrency and the number of table columns when
 configuring insert_dop. For high concurrency or tables with many columns, use a
 smaller insert_dop value. Disable insert_dop if a table has over 1,000 columns to avoid
 performance issues.

15.9.2 Optimizer Cost Constants

This section describes the optimizer cost constants. The cost variables described in this section are measured on an arbitrary scale. Only their relative values matter, therefore scaling them all in or out by the same factor will result in no differences in the optimizer's choices. By default, these cost variables are based on the cost of sequential page fetches, that is, **seq_page_cost** is conventionally set to **1.0** and the other cost variables are set with reference to the parameter. However, you can use a different scale, such as actual execution time in milliseconds.

seq_page_cost

Parameter description: Specifies the optimizer's estimated cost of a disk page fetch that is part of a series of sequential fetches.

Type: USERSET

Value range: a floating point number ranging from 0 to DBL_MAX

Default value: 1

random_page_cost

Parameter description: Specifies the optimizer's estimated cost of an out-of-sequence disk page fetch.

Type: USERSET

Value range: a floating point number ranging from 0 to DBL MAX

Default value: 4

□ NOTE

- Although the server allows you to set the value of random_page_cost to less than that
 of seq_page_cost, it is not physically sensitive to do so. However, setting them equal
 makes sense if the database is entirely cached in RAM, because in that case there is no
 penalty for fetching pages out of sequence. Also, in a heavily-cached database you
 should lower both values relative to the CPU parameters, since the cost of fetching a
 page already in RAM is much smaller than it would normally be.
- This value can be overwritten for tables and indexes in a particular tablespace by setting the tablespace parameter of the same name.
- Comparing to seq_page_cost, reducing this value will cause the system to prefer index scans and raising it makes index scans relatively more expensive. You can increase or decrease both values at the same time to change the disk I/O cost relative to CPU cost.

cpu tuple cost

Parameter description: Specifies the optimizer's estimated cost of processing each row during a guery.

Type: USERSET

Value range: a floating point number ranging from 0 to DBL_MAX

Default value: 0.01

cpu_index_tuple_cost

Parameter description: Specifies the optimizer's estimated cost of processing each index entry during an index scan.

Type: USERSET

Value range: a floating point number ranging from 0 to DBL_MAX

Default value: 0.005

cpu_operator_cost

Parameter description: Specifies the optimizer's estimated cost of processing each operator or function during a query.

Type: USERSET

Value range: a floating point number ranging from 0 to DBL_MAX

Default value: 0.0025

effective cache size

Parameter description: Specifies the optimizer's assumption about the effective size of the disk cache that is available to a single query.

When setting this parameter you should consider both GaussDB(DWS)'s shared buffer and the kernel's disk cache. Also, take into account the expected number of concurrent queries on different tables, since they will have to share the available space.

This parameter has no effect on the size of shared memory allocated by GaussDB(DWS). It is used only for estimation purposes and does not reserve kernel disk cache. The value is in the unit of disk page. Usually the size of each page is 8192 bytes.

Type: USERSET

Value range: an integer ranging is from 1 to INT_MAX. The unit is 8 KB.

A value greater than the default one may enable index scanning, and a value less than the default one may enable sequence scanning.

Default value: 128MB

allocate_mem_cost

Parameter description: Specifies the query optimizer's estimated cost of creating a Hash table for memory space using Hash join. This parameter is used for optimization when the Hash join estimation is inaccurate.

Type: USERSET

Value range: a floating point number ranging from 0 to DBL_MAX

Default value: 0

smp_thread_cost

Parameter description: Specifies the optimizer's cost for calculating parallel threads of an operator. This parameter is used for tuning if **query_dop** is not suitable for system load management. (This parameter is supported only by clusters of version 8.2.0 or later.)

Type: USERSET

Value range: a floating point number ranging from 1 to 10000

Default value: 1000

15.9.3 Genetic Query Optimizer

This section describes parameters related to genetic query optimizer. The genetic query optimizer (GEQO) is an algorithm that plans queries by using heuristic searching. This algorithm reduces planning time for complex queries and the cost of producing plans are sometimes inferior to those found by the normal exhaustive-search algorithm.

gego

Parameter description: Controls the use of genetic query optimization.

Type: USERSET

Value range: Boolean

- on indicates GEQO is enabled.
- off indicates GEQO is disabled.

Default value: on

NOTICE

Generally, do not set this parameter to **off**. **geqo_threshold** provides more subtle control of GEQO.

geqo_threshold

Parameter description: Specifies the number of **FROM** items. Genetic query optimization is used to plan queries when the number of statements executed is greater than this value.

Type: USERSET

Value range: an integer ranging from 2 to INT_MAX

Default value: 12

NOTICE

- For simpler queries it is best to use the regular, exhaustive-search planner, but for queries with many tables it is better to use GEQO to manage the queries.
- A **FULL OUTER JOIN** construct counts as only one **FROM** item.

geqo_effort

Parameter description: Controls the trade-off between planning time and query plan quality in GEQO.

Type: USERSET

Value range: an integer ranging from 1 to 10

Default value: 5

NOTICE

- Larger values increase the time spent in query planning, but also increase the probability that an efficient query plan is chosen.
- geqo_effort does not have direct effect. This parameter is only used to compute the default values for the other variables that influence GEQO behavior. You can manually set other parameters as required.

geqo_pool_size

Parameter description: Specifies the pool size used by GEQO, that is, the number of individuals in the genetic population.

Type: USERSET

Value range: an integer ranging from 0 to INT_MAX

NOTICE

The value of this parameter must be at least **2**, and useful values are typically from **100** to **1000**. If this parameter is set to **0**, GaussDB(DWS) selects a proper value based on **geqo_effort** and the number of tables.

Default value: 0

geqo_generations

Parameter description: Specifies the number parameter iterations of the algorithm used by GEQO.

Type: USERSET

Value range: an integer ranging from 0 to INT_MAX

NOTICE

The value of this parameter must be at least 1, and useful values are typically from 100 to 1000. If it is set to 0, a suitable value is chosen based on **gego pool size**.

Default value: 0

geqo_selection_bias

Parameter description: Specifies the selection bias used by GEQO. The selection bias is the selective pressure within the population.

Type: USERSET

Value range: a floating point number ranging from 1.5 to 2.0

Default value: 2

geqo_seed

Parameter description: Specifies the initial value of the random number generator used by GEQO to select random paths through the join order search space.

Type: USERSET

Value range: a floating point number ranging from 0.0 to 1.0

NOTICE

Varying the value changes the setting of join paths explored, and may result in a better or worse path being found.

Default value: 0

15.9.4 Other Optimizer Options

default_statistics_target

Parameter description: Specifies the default statistics target for table columns without a column-specific target set via **ALTER TABLE SET STATISTICS**. If this

parameter is set to a positive number, it indicates the number of samples of statistics information. If this parameter is set to a negative number, percentage is used to set the statistic target. The negative number converts to its corresponding percentage, for example, –5 means 5%. During sampling, a random sample size is determined by multiplying the **default_statistics_target** by 300. For example, if the default value is **100**, then 30,000 pages will be randomly read and 30,000 data records will be randomly selected from them to complete the random sampling.

Type: USERSET

Value range: an integer ranging from -100 to 10000

NOTICE

- A larger positive number than the parameter value increases the time required to do **ANALYZE**, but might improve the quality of the optimizer's estimates.
- Changing settings of this parameter may result in performance deterioration. If query performance deteriorates, you can:
 - 1. Restore to the default statistics.
 - 2. Use hints to optimize the query plan. For details, see Hint-based Tuning.
- If this parameter is set to a negative value, the number of samples is greater than or equal to 2% of the total data volume, and the number of records in user tables is less than 1.6 million, the time taken by running **ANALYZE** will be longer than when this parameter uses its default value.
- **AUTOANALYZE** does not allow you to set a sampling size for temporary table sampling. Its default value will be used for sampling.
- If statistics are forcibly calculated based on memory, the sampling size is limited by the **maintenance_work_mem** parameter.

Default value: 100

random function version

Parameter description: Specifies the random function version selected by ANALYZE during data sampling. This feature is supported only in 8.1.2 or later.

Type: USERSET

Value range: enumerated values

- The value **0** indicates that the random function provided by the C standard library is used.
- The value **1** indicates that the optimized and enhanced random function is used.

Default value:

- If the current cluster is upgraded from an earlier version to 8.2.0.100, the default value is **0** to ensure forward compatibility.
- If the cluster version 8.2.0.100 is newly installed, the default value is 1.

constraint exclusion

Parameter description: Controls the query optimizer's use of table constraints to optimize queries.

Type: USERSET

Value range: enumerated values

- **on** indicates the constraints for all tables are examined.
- **off**: No constraints are examined.
- partition indicates that only constraints for inherited child tables and UNION ALL subqueries are examined.

NOTICE

When **constraint_exclusion** is set to **on**, the optimizer compares query conditions with the table's **CHECK** constraints, and omits scanning tables for which the conditions contradict the constraints.

Default value: partition

■ NOTE

Currently, this parameter is set to **on** by default to partition tables. If this parameter is set to **on**, extra planning is imposed on simple queries, which has no benefits. If you have no partitioned tables, set it to **off**.

cursor_tuple_fraction

Parameter description: Specifies the optimizer's estimated fraction of a cursor's rows that are retrieved.

Type: USERSET

Value range: a floating point number ranging from 0.0 to 1.0

NOTICE

Smaller values than the default value bias the optimizer towards using **fast start** plans for cursors, which will retrieve the first few rows quickly while perhaps taking a long time to fetch all rows. Larger values put more emphasis on the total estimated time. At the maximum setting of **1.0**, cursors are planned exactly like regular queries, considering only the total estimated time and how soon the first rows might be delivered.

Default value: 0.1

from_collapse_limit

Parameter description: Specifies whether the optimizer merges sub-queries into upper queries based on the resulting FROM list. The optimizer merges sub-queries

into upper queries if the resulting FROM list would have no more than this many items.

Type: USERSET

Value range: an integer ranging from 1 to INT_MAX

NOTICE

Smaller values reduce planning time but may lead to inferior execution plans.

Default value: 8

join_collapse_limit

Parameter description: Specifies whether the optimizer rewrites **JOIN** constructs (except **FULL JOIN**) into lists of **FROM** items based on the number of the items in the result list.

Type: USERSET

Value range: an integer ranging from 1 to INT_MAX

NOTICE

- Setting this parameter to 1 prevents join reordering. As a result, the join order specified in the query will be the actual order in which the relations are joined. The query optimizer does not always choose the optimal join order. Therefore, advanced users can temporarily set this variable to 1, and then specify the join order they desire explicitly.
- Smaller values reduce planning time but lead to inferior execution plans.

Default value: 8

join_search_mode

Parameter description: plan path search mode.

Type: USERSET

Value range: enumerated values

- **exhaustive**: Traditional dynamic planning and genetic methods are used to search for planned paths.
- heuristic: The heuristic method is used to search for planned paths. This
 method improves the plan generation performance, but there is a possibility
 that the optimal plan is ignored. This setting only takes effect for scenarios
 where a Drive Hint is specified or the number of joined tables exceeds
 from collapse limit.

Default value: heuristic

enable_from_collapse_hint

Parameter description: Specifies whether to rewrite the **FROM** list to make the hint take effect, and then rewrite it again based on the **from_collapse_limit** and **join_collapse_limit** parameters. This parameter is supported by clusters of version 8.2.0 or later.

Type: USERSET

Value range: Boolean

- on indicates that the FROM list is first rewritten in hint mode.
- **off** indicates that the **FROM** list is rewritten without difference.

NOTICE

- If this parameter is enabled, the optimizer preferentially rewrites the **FROM** list in hint mode. However, you can learn whether a hint takes effect only after the plan is generated.
- If this parameter is disabled, the plan is generated in the same way as that in versions earlier than 8.2.0. That is, the plan is generated regardless of whether the table has hints.

Default value: on

plan_mode_seed

Parameter description: This is a commissioning parameter. Currently, it supports only OPTIMIZE_PLAN and RANDOM_PLAN. OPTIMIZE_PLAN indicates the optimal plan, the cost of which is estimated using the dynamic planning algorithm, and its value is 0. RANDOM_PLAN indicates the plan that is randomly generated. If plan_mode_seed is set to -1, you do not need to specify the value of the seed identifier. Instead, the optimizer generates a random integer ranging from 1 to 2147483647, and then generates a random execution plan based on this random number. If plan_mode_seed is set to an integer ranging from 1 to 2147483647, you need to specify the value of the seed identifier, and the optimizer generates a random execution plan based on the seed value.

Type: USERSET

Value range: an integer ranging from -1 to 2147483647

Default value: 0

NOTICE

- If plan_mode_seed is set to RANDOM_PLAN, the optimizer generates different random execution plans, which may not be the optimal. Therefore, to guarantee the query performance, the default value 0 is recommended during upgrade, scale-out, scale-in, and O&M.
- If this parameter is not set to **0**, the specified hint will not be used.

enable hdfs predicate pushdown

Parameter description: Specifies whether the function of pushing down predicates the native data layer is enabled.

Type: SUSET

Value range: Boolean

- **on** indicates this function is enabled.
- off indicates this function is disabled.

Default value: on

windowagg pushdown enhancement

Parameter description: Specifies whether to enable enhanced predicate pushdown for window functions in aggregation scenarios. (This parameter is supported only by clusters of version 8.2.0 or later.)

Type: SUSET

Value range: Boolean

- **on** indicates that the predicate pushdown enhancement for window functions is enabled in aggregation scenarios.
- **off** indicates that the predicate pushdown enhancement for window functions is disabled in aggregation scenarios.

Default value: on

implied_quality_optmode

Parameter description: Specifies how to pass conditions for the equivalent columns in a statement. (This parameter is supported only by clusters of version 8.2.0 or later.)

Type: SUSET

Value range: enumerated values

- **normal** indicates forward compatibility with 8.1.3 and earlier versions, that is, the implied expression behavior is optimized.
- negative indicates that the implied expression behavior is not optimized.
- **positive** indicates that type conversion expressions are optimized in addition to the operations specified by **normal**.

Default value: normal

enable random datanode

Parameter description: Specifies whether the function that random query about DNs in the replication table is enabled. A complete data table is stored on each DN for random retrieval to release the pressure on nodes.

Type: USERSET

Value range: Boolean

on: This function is enabled.off: This function is disabled.

Default value: on

hashagg table size

Parameter description: Specifies the hash table size during HASH AGG execution.

Type: USERSET

Value range: an integer ranging from 0 to INT_MAX/2

Default value: 0

enable_codegen

Parameter description: Specifies whether code optimization can be enabled. Currently, the code optimization uses the LLVM optimization.

Type: USERSET

Value range: Boolean

- **on** indicates code optimization can be enabled.
- off indicates code optimization cannot be enabled.

Default value: off

⚠ CAUTION

- In clusters of version 9.1.0.218 or later, the default value of this parameter is **off**. In clusters of other versions, the default value of this parameter is **on**.
- Currently, the LLVM optimization only supports the vectorized executor and SQL on Hadoop features. You are advised to set this parameter to off in other cases.

codegen_strategy

Parameter description: Specifies the codegen optimization strategy that is used when an expression is converted to codegen-based.

Type: USERSET

Value range: enumerated values

- partial indicates that you can still call the LLVM dynamic optimization strategy using the codegen framework of an expression even if functions that are not codegen-based exist in the expression.
- **pure** indicates that the LLVM dynamic optimization strategy can be called only when all functions in an expression can be codegen-based.

NOTICE

In the scenario where query performance reduces after the codegen function is enabled, you can set this parameter to **pure**. In other scenarios, do not change the default value **partial** of this parameter.

Default value: partial

enable codegen print

Parameter description: Specifies whether the LLVM IR function can be printed in logs.

Type: USERSET

Value range: Boolean

- on indicates that the LLVM IR function can be printed in logs.
- off indicates that the LLVM IR function cannot be printed in logs.

Default value: off

codegen_cost_threshold

Parameter description: The LLVM compilation takes some time to generate executable machine code. Therefore, LLVM compilation is beneficial only when the actual execution cost is more than the sum of the code required for generating machine code and the optimized execution cost. This parameter specifies a threshold. If the estimated execution cost exceeds the threshold, LLVM optimization is performed.

Type: USERSET

Value range: an integer ranging from 0 to INT_MAX

Default value: 10000

llvm_compile_expr_limit

Parameter description: Specifies the limit for compiling expressions with LLVM. If there are more expressions than the limit, only the first ones are compiled and an alarm is generated. (To allow the alarm to be generated, execute **SET analysis_options="on(LLVM_COMPILE)"** before **explain performance** is executed.)

Type: USERSET

Value range: an integer ranging from -1 to INT_MAX

Default value: 500

llvm_compile_time_limit

Parameter description: If the percentage of the LLVM compilation time to the executor running time exceeds the threshold specified by **llvm_compile_time_limit**, an alarm is generated. (To allow the alarm to be

generated, execute **SET analysis_options="on(LLVM_COMPILE)"** before **explain performance** is executed.) This parameter is supported only by clusters of version 8.3.0 or later.

Type: USERSET

Value range: a floating point number ranging from 0.0 to 1.0

Default value: 0.2

enable_constraint_optimization

Parameter description: Specifies whether the informational constraint optimization execution plan can be used for an HDFS foreign table.

Type: SUSET

Value range: Boolean

- on indicates the plan can be used.
- **off** indicates the plan cannot be used.

Default value: on

enable_bloom_filter

Parameter description: Specifies whether the BloomFilter optimization is used.

Type: USERSET

Value range: Boolean

- **on** indicates the BloomFilter optimization can be used.
- **off** indicates the BloomFilter optimization cannot be used.

Default value: on

NOTICE

Scenario: If in a HASH JOIN, the thread of the foreign table contains HDFS tables or column-store tables, the Bloom filter is triggered.

Constraints:

- 1. Only INNER JOIN, SEMI JOIN, RIGHT JOIN, RIGHT SEMI JOIN, RIGHT ANTI JOIN and RIGHT ANTI FULL JOIN are supported.
- 2. JOIN condition of the internal table: It cannot be an expression for HDFS internal or foreign tables. It can be an expression for column-store tables, but only at the non-join layer.
- 3. The join condition of the foreign table must be simple column join.
- 4. When the join conditions of the internal and foreign tables (HDFS) are both simple column joins, the estimated data that can be removed at the plan layer must be over 1/3.
- 5. Joined columns cannot contain NULL values.
- 6. Data type:
 - HDFS internal and foreign tables support SMALLINT, INTEGER, BIGINT, REAL/FLOAT4, DOUBLE PRECISION/FLOAT8, CHAR(n)/CHARACTER(n)/ NCHAR(n), VARCHAR(n)/CHARACTER VARYING(n), CLOB and TEXT.
 - Column-store tables support SMALLINT, INTEGER, BIGINT, OID, "char", CHAR(n)/CHARACTER(n)/NCHAR(n), VARCHAR(n)/CHARACTER VARYING(n), NVARCHAR2(n), CLOB, TEXT, DATE, TIME, TIMESTAMP and TIMESTAMPTZ. The collation of the character type must be C.

runtime_filter_type

Parameter description: Specifies the type of runtime filter used, and only takes effect when **enable_bloom_filter** is enabled. This is supported only by clusters of version 9.1.0.100 or later.

Type: USERSET

Value range: enumerated values

- All indicates that the runtime filters in all scenarios are used.
- Topn_filter indicates the runtime filters in the ORDER BY scenario with LIMIT are used.
- **Bloom_filter** indicates that only runtime filters in join scenarios are used, and a bloom filter is generated for filtering after meeting certain conditions.
- Min_max indicates that only the runtime filters in join scenarios are used and only a min_max filter is generated for filtering.
- **None** indicates that no runtime filters are used, and only the original bloom filter has filtering effect.

Default value: All

NOTICE

- Application scenario: Plan type of the HASH JOIN foreign table in a columnstore table and the ORDER BY plan type with LIMIT in a column-store table.
- Constraints:
 - The usage restrictions for JOIN scenarios are the same as those for the enable_bloom_filter parameter.
 - In the order by scenario with limit, the order by field types only support SMALLINT, INTEGER, BIGINT, "char", CHAR(n)/CHARACTER(n)/ NCHAR(n), VARCHAR(n)/CHARACTER VARYING(n), NVARCHAR2(n), TEXT, DATE, TIME, TIMESTAMP, and TIMESTAMPTZ, and the sorting rules for character types must be specified as C.

runtime_filter_ratio

Parameter description: Specifies the threshold for using bloom filter for fine-grained row-level filtering in join scenarios in runtime filter, and only takes effect when **runtime_filter_type** is set to a value greater than or equal to **Bloom_filter**. This is supported only by clusters of version 9.1.0.100 or later.

Type: USERSET

Value range: a floating point number ranging from 0.0 to 1.0

Default value: 0.01

NOTICE

- Application scenario: HASH JOIN of column-store tables, where the internal table estimate_join_rows/foreign table estimate_join_rows ≤ runtime_filter_ratio. Fine-grained row-level filtering is only recommended for join scenarios where there is a significant difference in data volume between the internal and foreign tables. Improper runtime_filter_ratio settings may lead to degraded performance in join scenarios.
- Usage restrictions: Fine-grained row-level filtering is only supported for join field types of **SMALLINT**, **INTEGER**, **BIGINT**, and **FLOAT**.

runtime_filter_cost_options

Parameter description: Specifies whether to generate a runtime filter plan based on the cost. This is supported only by clusters of version 9.1.0.200 or later.

Type: USERSET

Value range: a string

- **apply_partial**: The runtime filter path can be generated as long as the **build** end contains a table required by a runtime filter on the **probe** end.
- **apply_all**: The runtime filter path can be generated only when the **build** end contains all tables required by the runtime filter that can be applied to the **probe** end.

Default value: ", indicating that runtime filters are not used during plan generation, regardless of the cost.

NOTICE

If both apply_partial and apply_all are set, the setting of apply_all takes effect.

enable extrapolation stats

Parameter description: Specifies whether to use the extrapolation logic based on historical statistics. Using this logic may increase the accuracy of estimation for tables whose statistics have not been collected. However, there is also a possibility that the estimation is too large due to incorrect inference.

Type: USERSET

Value range: Boolean

- **on** indicates that the extrapolation logic is used for data of DATE type based on historical statistics.
- **off** indicates that the extrapolation logic is not used for data of DATE type based on historical statistics.

Default value:

- If the current cluster is upgraded from an earlier version to 8.2.0.100, the default value is **off** to ensure forward compatibility.
- If the cluster version 8.2.0.100 is newly installed, the default value is on.

autoanalyze

Parameter description: Specifies whether to allow automatic statistics collection for a table that has no statistics or a table whose amount of data modification reaches the threshold for triggering ANALYZE when a plan is generated. In this case, AUTOANALYZE cannot be triggered for foreign tables or temporary tables with the ON COMMIT [DELETE ROWS|DROP] option. To collect statistics, you need to manually perform the ANALYZE operation. If an exception occurs in the database during the execution of autoanalyze on a table, after the database is recovered, the system may still prompt you to collect the statistics of the table when you run the statement again. In this case, manually perform the ANALYZE operation on the table to synchronize statistics.

NOTICE

If the amount of data modification reaches the threshold for triggering **ANALYZE**, the amount of data modification exceeds **autovacuum_analyze_threshold** + **autovacuum_analyze_scale_factor** * *reltuples*. *reltuples* indicates the estimated number of rows in the table recorded in **pg class**.

Type: SUSET

Value range: Boolean

- **on** indicates that the table statistics are automatically collected.
- **off** indicates that the table statistics are not automatically collected.

Default value: on

enable analyze partition

Parameter description: Specifies whether to support collecting statistics for a specific partition of a table. After enabling this parameter, you can collect statistics for a specific partition using **ANALYZE table_name PARTITION** (partition_name), and when querying data on this partition of the table, the optimizer will choose to use partition statistics.

Type: USERSET

Value range: Boolean

- **on** indicates supporting collecting statistics for a specific partition of a table.
- **off** indicates that collecting statistics for a specific partition of a table is not supported.

Default value: off

analyze_use_dn_correlation

Parameter description: Specifies whether CNs use correlation statistics of DNs when executing ANALYZE. This is supported only by clusters of version 9.1.0.100 or later.

Type: USERSET

Value range: Boolean

- **on** indicates that CNs use correlation statistics of DNs.
- off indicates that CNs do not use correlation statistics of DNs.

Default value: on

analyze_predicate_column_threshold

Parameter description: Specifies whether to enable ANALYZE operations for predicate columns and the minimum number of columns supported. This is supported only by clusters of version 9.1.0.100 or later.

Type: SIGHUP

Value range: an integer ranging from 0 to 10000

- The value **0** indicates that ANALYZE operations are disabled for predicate columns and predicate columns are not collected or analyzed.
- A value greater than 0 indicates that predicate column collection is enabled and predicate column analysis is performed only on tables whose number of columns is greater than or equal to the value of this parameter.

Default value: 10

enable_runtime_analyze_concurrent

Parameter description: Specifies whether to support concurrent RUNTIME ANALYZE operations on a table. This is supported only by clusters of version 9.1.0.100 or later.

Type: USERSET

Value range: Boolean

- **on** indicates that concurrent operations are supported.
- **off** indicates that concurrent operations are not supported.

Default value: on

analyze_max_columns_count

Parameter description: Specifies the maximum number of columns supported by ANALYZE. This is supported only by clusters of version 9.1.0.100 or later.

Type: USERSET

Value range: an integer ranging from -1 to 10000

- -1 indicates that the number of columns supported by ANALYZE is not limited.
- A value greater than -1 indicates that only columns up to this value will be collected, and any columns beyond this value will not be collected.

Default value: -1

query_dop

Parameter description: Specifies the user-defined degree of parallelism.

Type: USERSET

Value range: an integer ranging from -64 to 64.

[1, 64]: Fixed SMP is enabled, and the system will use the specified degree.

0: SMP adaptation function is enabled. The system dynamically selects the optimal parallelism degree [1,8] (x86 platforms) or [1,64] (Kunpeng platforms) for each query based on the resource usage and query plans.

[-64, -1]: SMP adaptation is enabled, and the system will dynamically select a degree from the limited range.

□ NOTE

- For TP services that mainly involve short queries, if services cannot be optimized through lightweight CNs or statement delivery, it will take a long time to generate an SMP plan. You are advised to set **query_dop** to **1**. For AP services with complex statements, you are advised to set **query_dop** to **0**.
- After enabling concurrent queries, ensure you have sufficient CPU, memory, network, and I/O resources to achieve the optimal performance.
- To prevent performance deterioration caused by an overly large value of **query_dop**, the system calculates the maximum number of available CPU cores for a DN and uses the number as the upper limit for this parameter. If the value of **query_dop** is greater than 4 and also the upper limit, the system resets **query_dop** to the upper limit.

Default value: 1 (0 for cloud flavors with 64 GB or larger memory)

query_dop_ratio

Parameter description: Specifies the DOP multiple used to adjust the optimal DOP preset in the system when **query_dop** is set to **0**. That is, DOP = Preset DOP x query_dop_ratio (ranging from 1 to 64). If this parameter is set to **1**, the DOP cannot be adjusted.

Type: USERSET

Value range: a floating point number ranging from 0 to 64

Default value: 1

debug_group_dop

Parameter description: Specifies the unified DOP parallelism degree allocated to the groups that use the Stream operator as the vertex in the generated execution plan when the value of query_dop is 0. This parameter is used to manually specify the DOP for specific groups for performance optimization. Its format is G1,D1,G2,D2,...,, where G1 and G2 indicate the group IDs that can be obtained from logs and D1 and D2 indicate the specified DOP values and can be any positive integers.

Type: USERSET

Value range: a string

Default value: empty

NOTICE

This parameter is used only for internal optimization and cannot be set. You are advised to use the default value.

enable_analyze_check

Parameter description: Checks whether statistics were collected about tables whose reltuples and relpages are shown as 0 in pg_class during plan generation. This parameter has been discarded in clusters of version 8.1.3 or later, but is

reserved for compatibility with earlier versions. The setting of this parameter does not take effect.

Type: SUSET

Value range: Boolean

on enables the check.

• **off** disables the check.

Default value: on

enable_sonic_hashagg

Parameter description: Specifies whether to use the Hash Agg operator for column-oriented hash table design when certain constraints are met.

Type: USERSET

Value range: Boolean

- **on** indicates that the Hash Agg operator is used for column-oriented hash table design when certain constraints are met.
- **off** indicates that the Hash Agg operator is not used for column-oriented hash table design.

◯ NOTE

- If enable_sonic_hashagg is enabled and certain constraints are met, the Hash Agg
 operator will be used for column-oriented hash table design, and the memory usage of
 the operator can be reduced. However, in scenarios where the code generation
 technology (enabled by enable_codegen) can significantly improve performance, the
 performance of the operator may deteriorate.
- If enable_sonic_hashagg is set to on, when certain constraints are met, the hash
 aggregation operator designed for column-oriented hash tables is used and its name is
 displayed as Sonic Hash Aggregation in the output of the Explain Analyze/Performance
 operation. When the constraints are not met, the operator name is displayed as Hash
 Aggregation.

Default value: on

enable_sonic_hashjoin

Parameter description: Specifies whether to use the Hash Join operator for column-oriented hash table design when certain constraints are met.

Type: USERSET

Value range: Boolean

- **on** indicates that the Hash Join operator is used for column-oriented hash table design when certain constraints are met.
- **off** indicates that the Hash Join operator is not used for column-oriented hash table design.

□ NOTE

- Currently, the parameter can be used only for Inner Join.
- If **enable_sonic_hashjoin** is enabled, the memory usage of the Hash Inner operator can be reduced. However, in scenarios where the code generation technology can significantly improve performance, the performance of the operator may deteriorate.
- If enable_sonic_hashjoin is set to on, when certain constraints are met, the hash join operator designed for column-oriented hash tables is used and its name is displayed as Sonic Hash Join in the output of the Explain Analyze/Performance operation. When the constraints are not met, the operator name is displayed as Hash Join.

Default value: on

enable_sonic_optspill

Parameter description: Specifies whether to optimize the number of hash join or hash agg files spilled to disks in the sonic scenario. This parameter takes effect only when **enable_sonic_hashjoin** or **enable_sonic_hashagg** is enabled.

Type: USERSET

Value range: Boolean

- **on** indicates that the number of files spilled to disks is optimized.
- off indicates that the number of files spilled to disks is not optimized.

For the hash join or hash agg operator that meets the sonic criteria, if this parameter is set to **off**, one file is spilled to disks for each column. If this parameter is set to **on** and the data types of different columns are similar, only one file (a maximum of five files) will be spilled to disks.

Default value: on

expand_hashtable_ratio

Parameter description: Specifies the expansion ratio used to resize the hash table during the execution of the Hash Agg and Hash Join operators.

Type: USERSET

Value range: a floating point number of 0 or ranging from 0.5 to 10

- Value **0** indicates that the hash table is adaptively expanded based on the current memory size.
- The value ranging from 0.5 to 10 indicates the multiple used to expand the hash table. Generally, a larger hash table delivers better performance but occupies more memory space. If the memory space is insufficient, data may be spilled to disks in advance, causing performance deterioration.

Default value: 0

plan_cache_mode

Parameter description: Specifies the policy for generating an execution plan in the **prepare** statement.

Type: USERSET

Value range: enumerated values

- auto indicates that the custom plan or generic plan is selected by default.
- **force_generic_plan** indicates that the **generic plan** is forcibly used.
- **force_custom_plan** indicates that the **custom plan** is forcibly used.

- This parameter is valid only for the **prepare** statement. It is used when the parameterized field in the **prepare** statement has severe data skew.
- The **custom plan** is a plan generated after you run a prepared statement where parameters in the EXECUTE statement are embedded. The **custom plan** generates a plan based on specific parameters in the execute statement. This scheme generates a preferred plan based on specific parameters each time and has good execution performance. The disadvantage is that the plan needs to be regenerated before each execution, resulting in a large amount of repeated optimizer overhead.
- The **generic plan** is a plan generated after you run a prepared statement. The plan policy binds parameters to the plan when you run the EXECUTE statement and execute the plan. The advantage of this solution is that repeated optimizer overheads can be avoided in each execution. The disadvantage is that the plan may not be optimal when data skew occurs for the bound parameter field. When some bound parameters are used, the plan execution performance is poor.

Default value: auto

wlm_query_accelerate

Parameter description: Specifies whether the query needs to be accelerated when short query acceleration is enabled.

Type: USERSET

Value range: an integer ranging from -1 to 1

- -1: indicates that short queries are controlled by the fast lane, and the long queries are controlled by the slow lane.
- **0**: indicates that queries are not accelerated. Both short and long queries are controlled by the slow lane.
- 1: indicates that queries are accelerated. Both short queries and long queries are controlled by the fast lane.

Default value: -1

show_unshippable_warning

Parameter description: Specifies whether to print the alarm for the statement pushdown failure to the client.

Type: USERSET

Value range: Boolean

- **on**: Records the reason why the statement cannot be pushed down in a WARNING log and prints the log to the client.
- **off**: Logs the reason why the statement cannot be pushed down only.

Default value: off

hashjoin_spill_strategy

Parameter description: specifies the hash join policy for spilling data to disks. This feature is supported in 8.1.2 or later.

Type: USERSET

Value range: The value is an integer ranging from 0 to 6.

- **0**: If an inner table is too large to be fully stored in database memory, the table will be partitioned. If the table cannot be further partitioned and there is not enough memory for storing it, the system will check whether the foreign table can be stored in memory and be used to create a hash table. If the foreign table can be stored in the memory and used to create a hash table. HashJoin will be performed. Otherwise, NestLoop will be performed.
- 1: If an inner table is too large to be fully stored in database memory, the table will be partitioned. If the table cannot be further partitioned and there is still not enough memory for storing it, the system will check whether the foreign table can be stored in memory and be used to create a hash table. If both the inner and outer tables are large, a hash join is forcibly performed.
- 2: If the size of the inner table is large and cannot be partitioned after data is spilled to disks for multiple times, HashJoin will be forcibly performed.
- 3: If the size of the inner table is large and cannot be partitioned after data is spilled to disks for multiple times, the system attempts to place the outer table in the available memory of the database to create a hash table. If both the inner and outer tables are large, an error is reported.
- 4: If the size of the inner table is large and cannot be partitioned after data is spilled to disks for multiple times, an error is reported.
- 5: If the inner table is large and cannot be fully stored in database memory, and the foreign table can be fully stored in memory, the foreign table will be used to create a hash table and perform HashJoin. If the foreign table cannot be fully stored in memory, it will be partitioned until the inner and foreign tables cannot be further partitioned. Then, NestLoop will be performed.
- **6**: If the inner table is large and cannot be fully stored in database memory, and the foreign table can be fully stored in memory, the foreign table will be used to create a hash table and perform HashJoin. If the foreign table cannot be fully stored in memory, it will be partitioned until the inner and foreign tables cannot be further partitioned. Then, HashJoin will be forcibly performed.

□ NOTE

- This parameter is valid only for a vectorized hash join operator.
- If the number of distinct values is small and the data volume is large, data may fail to be flushed to disks. As a result, the memory usage is too high and the memory is out of control. If this parameter is set to **0**, the system attempts to swap the inner and outer tables or perform a nested loop join to prevent this problem. However, a nested loop join may deteriorate performance in some scenarios. In this case, this parameter can be set to **1**, **2**, or **6** to forcibly perform HashJoin.
- The value **0** does not take effect for a vectorized full join, and the behavior is the same as that of the value **1**. The system attempts to create a hash table only for the outer table and does not perform a nested loop join.
- If the inner table is too large to be fully stored in memory, but the foreign table can be stored in memory, you are advised to set this parameter to 5 or 6 rather than 0 or 1, directly performing Hashjoin on the foreign table without multiple rounds of partitioning and spill to disk. If a foreign table contains only a small amount of distinct data, creating a hash table using the foreign table may cause performance deterioration. In this case, you can change the value of this parameter to 0 or 1.

Default value: 0

max_streams_per_query

Parameter description: Controls the number of Stream nodes in a query plan. (This parameter is supported only in 8.1.1 and later cluster versions.)

Type: SUSET

Value range: an integer ranging from -1 to 10000.

- -1 indicates that the number of Stream nodes in the query plan is not limited.
- A value within the range 0 to 10000 indicates that when the number of Stream nodes in the query plan exceeds the specified value, an error is reported and the query plan will not be executed.

□ NOTE

- This parameter controls only the Stream nodes on DNs and does not control the Gather nodes on the CN.
- This parameter does not affect the EXPLAIN query plan, but affects EXPLAIN ANALYZE and EXPLAIN PERFORMANCE.

Default value: -1

enable_agg_limit_opt

Parameter description: Specifies whether to optimize **select distinct col from table limit N**. This parameter is valid only if N is less than 16,384. The parameter **table** indicates a column-store table. This parameter is supported only by clusters of version 8.2.0.101 or later.

Type: USERSET

Value range: Boolean

• **on** indicates that the optimization is enabled. After this function is enabled, query results are from different DNs, and you do not need to create a full hash table on each DN, significantly improving query performance.

• **off** indicates that the optimization is disabled.

Default value: on

stream_ctescan_pred_threshold

Parameter description: minimum number of filter criteria contained in a CTE when **enable_stream_ctescan** is set to **on** and the CTE contains only a single table filtering condition. If the value is greater than or equal to the value of this parameter, the share scan mode is used. If the value is less than the value of this parameter, the inline mode is used. This parameter is supported only by clusters of version 8.2.1 or later.

Type: SUSET

Value range: an integer ranging from 0 to INT_MAX

Default value: 2

stream_ctescan_max_estimate_mem

Parameter description: maximum estimated memory value of the CTE when **enable_stream_ctescan** is set to **on**. This parameter must be used together with **stream_ctescan_refcount_threshold**. If the estimated memory is greater than the value of **stream_ctescan_max_estimate_mem** and the number of references is less than the value of **stream_ctescan_refcount_threshold**, the inline mode is used. Otherwise, the sharescan mode is used. This parameter is supported only by clusters of version 8.2.1 or later.

Type: SUSET

Value range: an integer ranging from 32 x 1024 (32 MB) to INT_MAX, in KB.

Default value: 256 MB

stream_ctescan_refcount_threshold

Parameter description: maximum number of times that the CTE can be referenced when **enable_stream_ctescan** is set to **on**. This parameter must be used together with **stream_ctescan_max_estimate_mem**. If the estimated memory is greater than the value of **stream_ctescan_max_estimate_mem** and the number of references is less than the value of

stream_ctescan_refcount_threshold, the inline mode is used. Otherwise, the sharescan mode is used. This parameter is supported only by clusters of version 8.2.1 or later.

Type: SUSET

Value range: an integer ranging from 0 to INT_MAX

Default value: 4

□ NOTE

This parameter takes effect only when the value is greater than 0. When the value is 0, only **stream_ctescan_max_estimate_mem** is used to control the inline behavior.

inlist rough check threshold

Parameter description: Specifies the maximum number of values in the **IN** condition when **enable_csqual_pushdown** is enabled and the filter criterion is **IN** for rough check pushdown. If the number of values in the **IN** filter condition exceeds the value of this parameter, the maximum and minimum values in the **IN** filter condition are used for pushdown. This parameter is supported only by clusters of version 8.2.0.101 or later.

Type: SUSET

Value range: an integer ranging from 0 to 10000

Default value: 100

■ NOTE

If the **IN** condition is executed on the only distribution column of a table, values can be filtered on DNs. In this case, the maximum number of values in the **IN** condition is **inlist_rough_check_threshold** multiplied by the number of DNs.

enable_array_optimization

Parameter description: whether to split the Array type generated by the IN, ANY, or ALL condition into common expressions for execution. This parameter will support multiple optimizations such as vectorized execution, rough check pruning, and partition pruning. This parameter is supported only by clusters of version 8.2.1 or later.

Type: SUSET

Value range: Boolean

- **on** indicates that expressions of the Array type are split for optimization.
- off indicates that expressions of the Array type are not split for optimization.

Default value: on

max_skew_num

Parameter description: controls the number of skew values allowed by the optimizer for redistribution optimization. This parameter is supported only by clusters of version 8.2.1 or later.

Type: SUSET

Value range: an integer ranging from 0 to INT_MAX

Default value: 10

enable_dict_plan

Parameter description: Specifies whether the optimizer uses dictionary encoding to speed up queries that use operators such as **Group By** and **Filter**. This parameter is supported only by clusters of 8.3.0 or later.

Type: USERSET

Value range: Boolean

• **on**: enables the optimizer dictionary encoding.

• **off**: disables the optimizer dictionary encoding.

Default value: off

dict_plan_distinct_limit

Parameter description: Specifies the distinct value of a column in a table. Dictionary encoding is enabled only when the value is less than or equal to the threshold. This parameter is supported only by clusters of 8.3.0 or later.

Type: USERSET

Value range: 0 to INT_MAX

Default value: 10000

The two parameters **dict_plan_distinct_limit** and **dict_plan_duplicate_ratio** determine if dictionary encoding is applied.

dict_plan_duplicate_ratio

Parameter description: Specifies the repetition rate threshold of a column. Dictionary encoding is enabled only when the repetition rate of the column is greater than or equal to the threshold. Dictionary encoding is suitable for columns with a small number of distinct values and a high repetition rate. This parameter is supported only by clusters of 8.3.0 or later.

Type: USERSET

Value range: 0.0 to 100, in percentage

Default value: 90

∩ NOTE

The two parameters **dict_plan_distinct_limit** and **dict_plan_duplicate_ratio** determine if dictionary encoding is applied.

enable_cu_predicate_pushdown

Parameter description:

- Function overview: Specifies whether simple filter criteria are pushed down to the CU for filtering. Enabling this will enhance query performance, particularly when working with the bitmap_columns column and PCK sorting column. It applies to specific WHERE, IS NULL, and IN conditions. This parameter is supported only by clusters of 8.3.0 or later.
- 2. Supported column types:
 - Integer type: INT2, INT4, and INT8
 - Date and time type: DATE, TIMESTAMP, and TIMESTAMPTZ
 - String types: VARCHAR and TEXT

- Numeral type: NUMERIC (a maximum of 19 characters)
- 3. Query conditions: This function supports multiple **WHERE** expressions, including:
 - **IN** expression: matches multiple values.
 - IS NULL / IS NOT NULL condition: checks whether the column value is null.
 - Comparison expressions: greater than (>), less than (<), equal to (=), and not equal to (<>), which is used for range query and exact match.

Type: USERSET

Value range: Boolean

- **on**: Simple filter criteria are pushed down to the CU for filtering.
- **off**: Simple filter criteria are not pushed down to the CU for filtering.

Default value: on

■ NOTE

The basic filter conditions for the dictionary column include the equality operator (=), the **IN** expression, and the **IS (NOT) NULL** condition. This filter is known as the CU Predicate Filter, as it is pushed down to the storage layer and applied during the VectorBatch filling process.

info_constraint_options

Parameter description: Specifies whether or which kind of informational constraints can be created. This is supported only by clusters of version 9.1.0.100 or later.

Type: USERSET

Value range: enumerated values

- none: indicates that no informational constraint can be created.
- **foreign key**: indicates that foreign key constraints can be created.

Default value: none

15.10 Error Reporting and Logging

15.10.1 Logging Destination

log_statement_length_limit

Parameter description: Specifies the length of SQL statements to be printed. If the length of an SQL statement exceeds the specified value, the SQL statement is recorded in the **statement-***Timestamp.***log** file. This parameter is supported only by 9.1.0.200 and later versions.

Type: SUSET

Value range: an integer ranging from 0 to INT_MAX.

Default value: 1024

15.10.2 Logging Time

client_min_messages

Parameter description: Specifies which level of messages are sent to the client. Each level covers all the levels following it. The lower the level is, the fewer messages are sent.

Type: USERSET

NOTICE

When the values of **client_min_messages** and **log_min_messages** are the same, the levels are different.

Valid values: Enumerated values. Valid values: debug5, debug4, debug3, debug2, debug1, info, log, notice, warning, error For details about the parameters, see Table 15-1.

Default value: notice

log_min_messages

Parameter description: Specifies which level of messages will be written into server logs. Each level covers all the levels following it. The lower the level is, the fewer messages will be written into the log.

Type: SUSET

NOTICE

When the values of **client_min_messages** and **log_min_messages** are the same, the levels are different.

Value range: enumerated type. Valid values: debug5, debug4, debug3, debug2, debug1, info, log, notice, warning, error, fatal, panic For details about the parameters, see Table 15-1.

Default value: warning

log_min_error_statement

Parameter description: Specifies which SQL statements that cause errors condition will be recorded in the server log.

Type: SUSET

Value range: enumerated type. Valid values: debug5, debug4, debug3, debug2, debug1, info, log, notice, warning, error, fatal, panic For details about the parameters, see Table 15-1.

□ NOTE

- The default is **error**, indicating that statements causing errors, log messages, fatal errors, or panics will be logged.
- panic: This feature is disabled.

Default value: error

log_min_duration_statement

Parameter description: Specifies the threshold for logging statement execution durations. The execution duration that is greater than the specified value will be logged.

This parameter helps track query statements that need to be optimized. For clients using extended query protocol, durations of the Parse, Bind, and Execute are logged independently.

Type: SUSET

NOTICE

If this parameter and <code>log_statement</code> are used at the same time, statements recorded based on the value of <code>log_statement</code> will not be logged again after their execution duration exceeds the value of this parameter. If you are not using <code>syslog</code>, it is recommended that you log the process ID (PID) or session ID using <code>log_line_prefix</code> so that you can link the current statement message to the last logged duration.

Value range: an integer ranging from -1 to INT_MAX. The unit is millisecond.

- If this parameter is set to **250**, execution durations of SQL statements that run 250 ms or longer will be logged.
- **0**: Execution durations of all the statements are logged.
- -1: This feature is disabled.

Default value: 30min

backtrace min messages

Parameter description: Prints the function's stack information to the server's log file if the level of information generated is greater than or equal to this parameter level.

Type: SUSET

NOTICE

This parameter is used for locating customer on-site problems. Because frequent stack printing will affect the system's overhead and stability, therefore, when you locate the onsite problems, set the value of this parameter to ranks other than **fatal** and **panic**.

Value range: enumerated values

Valid values: debug5, debug4, debug3, debug2, debug1, info, log, notice, warning, error, fatal, panic For details about the parameters, see Table 15-1.

Default value: panic

Table 15-1 explains the message security levels used in GaussDB(DWS). If logging output is sent to **syslog** or **eventlog**, severity is translated in GaussDB(DWS) as shown in the table.

Table 15-1 Message Severity Levels

Severity	Description	syslog	eventlog
debug[1-5]	Provides detailed debug information.	DEBUG	INFORMATIO N
log	Reports information of interest to administrators, for example, checkpoint activity.	INFO	INFORMATIO N
info	Provides information implicitly requested by the user, for example, output from VACUUM VERBOSE.	INFO	INFORMATIO N
notice	Provides information that might be helpful to users, for example, notice of truncation of long identifiers and index created as part of the primary key.	NOTICE	INFORMATIO N
warning	Provides warnings of likely problems, for example, COMMIT outside a transaction block.	NOTICE	WARNING
error	Reports an error that causes a command to terminate.	WARNING	ERROR
fatal	Reports the reason that causes a session to terminate.	ERR	ERROR
panic	Reports an error that caused all database sessions to terminate.	CRIT	ERROR

plog_merge_age

Parameter description: Specifies the output interval of performance log data.

Type: SUSET

NOTICE

This parameter value is in milliseconds. You are advised to set this parameter to a value that is a multiple of 1000. That is, the value is in seconds. Name extension of the performance log files controlled by this parameter is .prf. These log files are stored in the \$GAUSSLOG/gs_profile/<node_name> directory. node_name is the value of pgxc_node_name in the postgres.conf file. You are advised not to use this parameter externally.

Value range: an integer ranging from 0 to INT_MAX. The unit is millisecond (ms).

- **0** indicates that the current session will not output performance log data.
- A value other than 0 indicates the output interval of performance log data. As the value decreases, more log data is generated, which negatively impacts performance.

Default value: 3s

profile_logging_module

Parameter description: Specifies the type of performance logs. When using this parameter, ensure that the value of **plog_merge_age** is not 0. This parameter is a session-level parameter, and you are not advised to use the **gs_guc** tool to set it. Only clusters of 8.1.3 and later versions support this function.

Type: USERSET

Value range: a string

Default value: OBS, HADOOP and REMOTE_DATANODE are enabled. MD is disabled. You can run the **SHOW profile_logging_module** command to view the value.

Setting method: First, you can run **SHOW profile_logging_module** to view which module is controllable. For example, the query output result is as follows:

Open the MD performance log and view the setting. The ALL identifier is equivalent to a shortcut operation. That is, logs of all modules can be enabled or disabled.

15.10.3 Logging Content

debug_pretty_print

Parameter description: Specifies the logs produced by **debug_print_parse**, **debug_print_rewritten**, and **debug_print_plan**. The output format is more

readable but much longer than the output generated when this parameter is set to **off**.

Type: USERSET

Value range: Boolean

- **on** indicates the indentation is enabled.
- **off** indicates the indentation is disabled.

Default value: on

log_duration

Parameter description: Specifies whether to record the duration of every completed SQL statement. For clients using extended query protocols, the time required for parsing, binding, and executing steps are logged independently.

Type: SUSET

Value range: Boolean

- If this parameter is set to **off**, the difference between setting this parameter and setting **log_min_duration_statement** is that when **log_min_duration_statement** is exceeded, the query text is logged, but this parameter does not log it.
- If this parameter is set to **on** and **log_min_duration_statement** has a positive value, all durations are logged but the query text is included only for statements exceeding the threshold. This behavior can be used for gathering statistics in high-load situation.

Default value: on

log_error_verbosity

Parameter description: Specifies the amount of detail written in the server log for each message that is logged.

Type: SUSET

Value range: enumerated values

- **terse** indicates that the output excludes the logging of DETAIL, HINT, QUERY, and CONTEXT error information.
- verbose indicates that the output includes the SQLSTATE error code, the source code file name, function name, and number of the line in which the error occurs.
- default indicates that the output includes the logging of DETAIL, HINT, QUERY, and CONTEXT error information, and excludes the SQLSTATE error code, the source code file name, function name, and number of the line in which the error occurs.

Default value: default

log_lock_waits

Parameter description: If the time that a session used to wait a lock is longer than the value of **deadlock_timeout**, this parameter specifies whether to record

this message in the database. This is useful in determining if lock waits are causing poor performance.

Type: SUSET

Value range: Boolean

- on indicates the information is recorded.
- **off** indicates the information is not recorded.

Default value: off

log_statement

Parameter description: Specifies whether to record SQL statements. For clients using extended query protocols, logging occurs when an execute message is received, and values of the Bind parameters are included (with any embedded single quotation marks doubled).

Type: SUSET

NOTICE

Statements that contain simple syntax errors are not logged even if **log_statement** is set to **all**, because the log message is emitted only after basic parsing has been completed to determine the statement type. If the extended query protocol is used, this setting also does not log statements before the execution phase (during parse analysis or planning). Set **log min error statement** to ERROR or lower to log such statements.

Value range: enumerated values

- **none** indicates that no statement is recorded.
- **ddl** indicates that all data definition statements, such as CREATE, ALTER, and DROP, are recorded.
- mod indicates that all DDL statements and data modification statements, such as INSERT, UPDATE, DELETE, TRUNCATE, and COPY FROM, are recorded.
- **all** indicates that all statements are recorded. The PREPARE, EXECUTE, and EXPLAIN ANALYZE statements are also recorded.

Default value: none

log_temp_files

Parameter description: Specifies whether to record the delete information of temporary files. Temporary files can be created for sorting, hashing, and temporary querying results. A log entry is generated for each temporary file when it is deleted.

Type: SUSET

Value range: an integer ranging from -1 to INT MAX. The unit is KB.

 A positive value indicates that the delete information of temporary files whose values are larger than that of log_temp_files is recorded.

- If the parameter is set to **0**, all the delete information of temporary files is recorded.
- If the parameter is set to -1, the delete information of no temporary files is recorded.

Default value: -1

logging_module

Parameter description: Specifies whether module logs can be output on the server. This parameter is a session-level parameter, and you are not advised to use the **gs_guc** tool to set it.

Type: USERSET

Value range: a string

Default value: **off**. All the module logs on the server can be viewed by running **show logging_module**.

Setting method: First, you can run **show logging_module** to view which module is controllable. For example, the query output result is as follows:

show logging_module; logging_module	
ALL,on(),off(DFS,GUC,HDFS,ORC,SLRU,MEM_CTL,AUTOVAC,ANALYZE,CACHE,ADIO,SSL,G	DS,TBLSPC,WLM,SP
ACE,OBS,EXECUTOR,VEC_EXECUTOR,STREAM,LLVM,OPT,OPT_REWRITE,OPT_JOIN,OPT_A	GG,OPT_SUBPLAN,
OPT_SETOP,OPT_CARD,OPT_SKEW,SMP,UDF,COOP_ANALYZE,WLMCP,ACCELERATE,PLANI	HINT,PARQUET,CARB
ONDATA, SNAPSHOT, XACT, HANDLE, CLOG, TQUAL, EC, REMOTE, CN_RETRY, PLSQL, TEXTSEA	RCH,SEQ,INSTR,CO
MM_IPC,COMM_PARAM,CSTORE,JOB,STREAMPOOL,STREAM_CTESCAN)	
(1 row)	

Controllable modules are identified by uppercase letters, and the special ID ALL is used for setting all module logs. You can control module logs to be exported by setting the log modules to **on** or **off**. Enable log output for SSL:

set logging_module='on(SSL)'; SET show logging_module;	
l	ogging_module
ALL,on(SSL),off(DFS,GUC,HDFS,ORC,SLRU,MEM_CTL,AUTOVA(ACE,OBS,EXECUTOR,VEC_EXECUTOR,STREAM,LLVM,OPT,OPT_FOPT_SETOP,OPT_CARD,OPT_SKEW,SMP,UDF,COOP_ANALYZE,WCCELERATE,PLANHINT,PARQUET,CARBONDATA,SNAPSHOT,XA(RY,PLSQL,TEXTSEARCH,SEQ,INSTR,COMM_IPC,COMM_PARAM,N) (1 row)	RÉWRITE,OPT_JOIN,OPT_AGG,OPT_SÚBPLÁN, /LMCP,A CT,HANDLE,CLOG,TQUAL,EC,REMOTE,CN_RET

SSL log output is enabled.

The ALL identifier is equivalent to a shortcut operation. That is, logs of all modules can be enabled or disabled.

```
set logging_module='off(ALL)';
SFT
show
logging_module;
                            logging_module
ALL,on(),off(DFS,GUC,HDFS,ORC,SLRU,MEM_CTL,AUTOVAC,ANALYZE,CACHE,ADIO,SSL,GDS,TBLSPC,WLM,SP
ACE, OBS, EXECUTOR, VEC EXECUTOR, STREAM, LLVM, OPT, OPT REWRITE, OPT JOIN, OPT AGG, OPT SUBPLAN,
OPT_SETOP,OPT_CARD,OPT_SKEW,SMP,UDF,COOP_ANALYZE,WLMCP,
ACCELERATE, PLANHINT, PARQUET, CARBONDATA, SNAPSHOT, XACT, HANDLE, CLOG, TQUAL, EC, REMOTE, CN_RE
TRY,PLSQL,TEXTSEARCH,SEQ,INSTR,COMM_IPC,COMM_PARAM,CSTORE,JOB,STREAMPOOL,STREAM_CTESCA
(1 row)
set logging_module='on(ALL)';
SFT
show
logging_module;
                 logging_module
ALL,on(DFS,GUC,HDFS,ORC,SLRU,MEM_CTL,AUTOVAC,ANALYZE,CACHE,ADIO,SSL,GDS,TBLSPC,WLM,SPACE,
OBS,EXECUTOR,VEC_EXECUTOR,STREAM,LLVM,OPT,OPT_REWRITE,OPT_JOIN,OPT_AGG,OPT_SUBPLAN,OPT_
SETOP,OPT_CARD,OPT_SKEW,SMP,UDF,COOP_ANALYZE,WLMCP,ACCELE
RATE, PLANHINT, PARQUET, CARBONDATA, SNAPSHOT, XACT, HANDLE, CLOG, TQUAL, EC, REMOTE, CN_RETRY, PLS
```

COMM_IPC logs must be enabled or disabled explicitly. You can run either of the following command to enable the log function of COMM_IPC:

QL,TEXTSEARCH,SEQ,INSTR,COMM_IPC,COMM_PARAM,CSTORE,JOB,STREAMPOOL,STREAM_CTESCAN),off()

```
set logging_module='on(ALL)';
SET
set logging_module='on(COMM_IPC)';
SET
```

After the setting is performed, the log function of the COMM_IPC module will not be automatically disabled. To disable the log function of the COMM_IPC module, you must run the following commands:

```
set logging_module='off(ALL)';
SET
set logging_module='off(COMM_IPC)';
SET
```

Dependency relationship: The value of this parameter depends on the settings of log_min_messages.

enable_unshipping_log

(1 row)

Parameter description: Specifies whether to log statements that are not pushed down. The logs help locate performance issues that may be caused by statements not pushed down.

Type: SUSET

Value range: Boolean

- **on**: Statements not pushed down will be logged.
- off: Statements not pushed down will not be logged.

Default value: on

log_statement_filter_list

Parameter description: Specifies whether to record SQL statements. It sets a collection of error codes, with multiple error codes separated by commas, for example: **GS_001**, **GS_002**. No SQL statement is recorded in the error code log. This parameter is supported only by 9.1.0.200 and later versions.

Type: SUSET

Value range: a string

Default value: an empty string

15.11 Runtime Statistics

15.11.1 Query and Index Statistics Collector

The query and index statistics collector is used to collect statistics during database running. The statistics include the times of inserting and updating a table and an index, the number of disk blocks and tuples, and the time required for the last cleanup and analysis on each table. The statistics can be viewed by querying system view families pg_stats and pg_statistic. The following parameters are used to set the statistics collection feature in the server scope.

track_activities

Parameter description: Collects statistics about the commands that are being executed in session.

Type: SUSET

Value range: Boolean

- **on** indicates that the statistics collection function is enabled.
- off indicates that the statistics collection function is disabled.

Default value: on

track_counts

Parameter description: Collects statistics about data activities.

Type: SUSET

Value range: Boolean

- on indicates that the statistics collection function is enabled.
- off indicates that the statistics collection function is disabled.

■ NOTE

When the database to be cleaned up is selected from the AutoVacuum automatic cleanup process, the database statistics are required. In this case, the default value is set to **on**.

Default value: on

track_io_timing

Parameter description: Collects statistics about I/O invoking timing in the database. The I/O timing statistics can be queried by using the **pg_stat_database** parameter.

Type: SUSET

Value range: Boolean

- If this parameter is set to **on**, the collection function is enabled. In this case, the collector repeatedly queries the OS at the current time. As a result, large numbers of costs may occur on some platforms. Therefore, the default value is set to **off**.
- **off** indicates that the statistics collection function is disabled.

Default value: off

track functions

Parameter description: Collects statistics about invoking times and duration in a function.

Type: SUSET

NOTICE

When the SQL functions are set to inline functions queried by the invoking, these SQL functions cannot be traced no matter these functions are set or not.

Value range: enumerated values

- **pl** indicates that only procedural language functions are traced.
- all indicates that SQL and C language functions are traced.
- **none** indicates that the function tracing function is disabled.

Default value: none

update_process_title

Parameter description: Collects statistics updated with a process name each time the server receives a new SQL statement.

The process name can be viewed on Windows task manager by running the **ps** command.

Type: SUSET

Value range: Boolean

• **on** indicates that the statistics collection function is enabled.

• **off** indicates that the statistics collection function is disabled.

Default value: off

track thread wait status interval

Parameter description: Specifies the interval of collecting the thread status information periodically.

Type: SUSET

Value range: an integer ranging from 0 to 1440, in minutes.

Default value: 30min

enable_save_datachanged_timestamp

Parameter description: Specifies whether to record the time when **INSERT**, **UPDATE**, **DELETE**, or **EXCHANGE/TRUNCATE/DROP PARTITION** is performed on table data.

Type: USERSET

Value range: Boolean

- **on** indicates that the time when an operation is performed on table data will be recorded.
- **off** indicates that the time when an operation is performed on table data will not be recorded.

Default value: on

enable_save_dataaccess_timestamp

Parameter description: Specifies whether to record the last access time of a table. This parameter is supported only by 8.2.1.210 and later cluster versions.

Type: USERSET

Value range: Boolean

- on indicates that the last access time of the table is recorded.
- **off** indicates that the last access time of the table is not recorded.

Default value: off

instr_unique_sql_count

Parameter description: Specifies whether to collect Unique SQL statements and the maximum number allowed.

Type: SIGHUP

Value range: an integer ranging from 0 to INT_MAX

- If it is set to 0, Unique SQL statistics are not collected.
- If the value is greater than **0**, the number of Unique SQL statements collected on the CN cannot exceed the value of this parameter. When the number of collected Unique SQL statements reaches the upper limit, the collection is stopped. In this case, you can increase the value of **reload** to continue the collection.

Default value: 0



If a new value is smaller than the original value, the Unique SQL statistics collected on the CN will be cleared.

track_sql_count

Parameter description: Specifies whether to collect statistics on the number of the **SELECT**, **INSERT**, **UPDATE**, **DELETE**, and **MERGE INTO** statements that are being executed in each session, the response time of the **SELECT**, **INSERT**, **UPDATE**, and **DELETE** statements, and the number of DDL, DML, and DCL statements.

Type: SUSET

Value range: Boolean

- **on** indicates that the statistics collection function is enabled.
- off indicates that the statistics collection function is disabled.

Default value: on

∩ NOTE

- The track_sql_count parameter is restricted by the track_activities parameter.
 - If track_activities is set to on and track_sql_count is set to off, a warning message indicating that track_sql_count is disabled will be displayed when the view gs_sql_count, pgxc_sql_count, gs_workload_sql_count, pgxc_workload_sql_count, global_workload_sql_count, gs_workload_sql_elapse_time, pgxc_workload_sql_elapse_time, or global_workload_sql_elapse_time are queried.
 - If both track_activities and track_sql_count are set to off, two logs indicating that track_activities is disabled and track_sql_count is disabled will be displayed when the views are queried.
 - If **track_activities** is set to **off** and **track_sql_count** is set to **on**, a log indicating that **track_activities** is disabled will be displayed when the views are queried.
- If this parameter is disabled, querying the view returns **0**.

enable_parallel_analyze

Parameter description: Specifies whether to use parallel sampling for internal and foreign table analysis. This parameter is supported only by clusters of version 9.1.0 or later.

Type: USERSET

Value range: Boolean

- **true** indicates that parallel sampling is used for internal and foreign table analysis.
- **false** indicates that parallel sampling is not used for internal and foreign table analysis.

Default value: true

CAUTION

- When enable_parallel_analyze is set to true and analyzing foreign tables, try
 to avoid adding NOT NULL constraints to the target foreign table columns to
 prevent constraint failure due to data source changes. Currently, parallel
 sampling does not support materialized views. If analyze fails due to such
 reasons, set this parameter to false.
- Currently, parallel sampling only supports analyzing ordinary column-store internal tables. This optimization does not take effect when the internal table uses hstore/hstore_opt or is declared as a replicated table.
- Currently, parallel sampling only supports analyzing foreign tables stored in parquet/orc format. This optimization does not take effect when the foreign table is in another format.

parallel_analyze_workers

Parameter description: Specifies the number of concurrent threads for parallel analyze sampling. This parameter is supported only by clusters of version 9.1.0 or later.

Type: USERSET

Value range: an integer ranging from 0 to 64

Default value: 10

∩ NOTE

The value of this parameter should correspond to the cluster load. When the cluster load is low, you can increase the parameter value appropriately based on the cluster configuration to further improve the efficiency of analyze execution.

analyze_sample_multiplier

Parameter description: Specifies the multiplier for the stripe sampling rate used in analyzing foreign tables. This parameter is supported only by clusters of version 9.1.0 or later.

Type: SUSET

Value range: an integer ranging from 0 to 100. **0** indicates that the stripe sampling rate is 100%.

Default value: 3

15.11.2 Performance Statistics

During database operation, accessing locks, disk I/O operations, and handling invalid messages can all be performance bottlenecks for the database. GaussDB(DWS) provides performance statistics methods that can help easily locate performance issues.

Generating Performance Statistics Logs

Parameter description: For each query, the following four parameters control the performance statistics of corresponding modules recorded in the server log:

- The **og_parser_stats** parameter controls the performance statistics of a parser recorded in the server log.
- The log_planner_stats parameter controls the performance statistics of a query optimizer recorded in the server log.
- The **log_executor_stats** parameter controls the performance statistics of an executor recorded in the server log.
- The **log_statement_stats** parameter controls the performance statistics of the whole statement recorded in the server log.

All these parameters can only provide assistant analysis for administrators, which are similar to the getrusage() of the Linux OS.

Type: SUSET

NOTICE

- **log_statement_stats** records the total statement statistics while other parameters only record statistics about each statement.
- The **log_statement_stats** parameter cannot be enabled together with other parameters recording statistics about each statement.

Value range: Boolean

- **on** indicates the function of recording performance statistics is enabled.
- off indicates the function of recording performance statistics is disabled.

Default value: off

15.12 Resource Management

If database resource usage is not controlled, concurrent tasks easily preempt resources. As a result, the OS will be overloaded and cannot respond to user tasks; or even crash and cannot provide any services to users. The GaussDB(DWS) workload management function balances the database workload based on available resources to avoid database overloading.

space_once_adjust_num

Parameter description: In the space control and space statistics functions, specifies the threshold of the number of files processed each time during slow

building and fine-grained calibration. This parameter is supported only by clusters of version 8.1.3 or later.

Type: SIGHUP

Value range: an integer ranging from 0 to INT_MAX

• The value **0** indicates that the slow build and fine-grained calibration functions are disabled.

Default value: 300

□ NOTE

The file quantity threshold affects database resources. You are advised to set the threshold to a proper value.

default_partition_cache_strategy

Parameter description: Specifies the default policy for controlling partition caching. This parameter is supported only by clusters of version 8.3.0 or later.

Type: USERSET

Value range: enumerated values

- cache_each_partition_as_possible enables maximum data caching. Data may not be written to CUs when being inserted into different partitions.
- **flush_when_switch_partition** indicates that data is written to CUs if the data belongs to different partitions during insertion.

Default value: cache each partition as possible

max active statements

Parameter description: Specifies the maximum global concurrency. This parameter applies to a job on a CN.

The database administrator changes the value of this parameter based on system resources (for example, CPU, I/O, and memory resources) so that the system fully supports the concurrency tasks and avoids too many concurrency tasks resulting in system crash.

Type: SIGHUP

Value range: an integer ranging from -1 to INT_MAX. The values **-1** and **0** indicate that the number of concurrent requests is not limited.

Default value: 60

cgroup_name

Parameter description: Specifies the name of the Cgroup in use. It can be used to change the priorities of jobs in the queue of a Cgroup.

If you set **cgroup_name** and then **session_respool**, the Cgroups associated with **session_respool** take effect. If you reverse the order, Cgroups associated with **cgroup_name** take effect.

If the Workload Cgroup level is specified during the **cgroup_name** change, the database does not check the Cgroup level. The level ranges from 1 to 10.

Type: USERSET

You are not advised to set **cgroup_name** and **session_respool** at the same time.

Value range: a string

Default value: DefaultClass:Medium

■ NOTE

DefaultClass:Medium indicates the **Medium** Cgroup belonging to the **Timeshare** Cgroup under the **DefaultClass** Cgroup.

enable_cgroup_switch

Parameter description: Specifies whether the database automatically switches to the **TopWD** group when executing statements by group type.

Type: USERSET

Value range: Boolean

- **on**: The database automatically switches to the **TopWD** group when executing statements by group type.
- **off**: The database does not automatically switch to the **TopWD** group when executing statements by group type.

Default value: off

memory_tracking_mode

Parameter description: Specifies the memory information recording mode.

Type: USERSET

Value range:

- **none**: Memory statistics is not collected.
- **normal:** Only memory statistics is collected in real time and no file is generated.
- **executor:** The statistics file is generated, containing the context information about all allocated memory used by the execution layer.
- **fullexec**: The generated file includes the information about all memory contexts requested by the execution layer.

Default value: none

memory_detail_tracking

Parameter description: Specifies the sequence number of the memory background information distributed in the needed thread and **plannodeid** of the query where the current thread is located.

Type: USERSET

Value range: a string

Default value: empty

NOTICE

It is recommended that you retain the default value for this parameter.

enable resource track

Parameter description: Specifies whether the real-time resource monitoring function is enabled. This parameter must be applied on both CNs and DNs.

Type: SIGHUP

Value range: Boolean

- on indicates the resource monitoring function is enabled.
- off indicates the resource monitoring function is disabled.

Default value: on

enable resource record

Parameter description: Specifies whether resource monitoring records are archived. When this parameter is enabled, records that have been executed are archived to the corresponding **INFO** views (**GS_WLM_SESSION_INFO** and **GS_WLM_OPERATOR_INFO**). This parameter must be applied on both CNs and DNs.

Type: SIGHUP

Value range: Boolean

- **on** indicates that the resource monitoring records are archived.
- **off** indicates that the resource monitoring records are not archived.

Default value: on

The default value of this parameter is **on** for a new cluster. In upgrade scenarios, the default value of this parameter is the same as that of the source version.

enable_track_record_subsql

Parameter description: Specifies whether to enable the function of recording and archiving sub-statements. When this function is enabled, sub-statements in stored procedures and anonymous blocks are recorded and archived to the corresponding **INFO** table (**GS_WLM_SESSION_INFO**). This parameter is a session-level parameter. It can be configured and take effect in the session connected to the CN and affects only the statements in the session. It can also be configured on both the CN and DN and take effect globally.

Type: USERSET

Value range: Boolean

- **on** indicates that the sub-statement resource monitoring records are archived.
- **off** indicates that the sub-statement resource monitoring records are not archived.

Default value: on

block_rule_cost

Parameter description: Specifies the minimum cost to trigger a query filter. This is supported only by clusters of version 9.1.0.200 or later.

Type: USERSET

Value range: an integer ranging from -1 to INT_MAX

- Value **-1** indicates that the system filters all statements without considering the cost.
- Value **0** indicates that all statements whose cost is greater than 0 are intercepted, but special statements (whose cost is 0) are not intercepted.
- If the value is greater than **0** and the statement cost is less than the value of **block_rule_cost**, the statement is not intercepted; otherwise, it is intercepted.

Default value: -1

enable user metric persistent

Parameter description: Specifies whether the user historical resource monitoring dumping function is enabled. When this function is enabled, data in the PG_TOTAL_USER_RESOURCE_INFO view is periodically sampled and saved to the GS_WLM_USER_RESOURCE_HISTORY system catalog, and data in the GS_RESPOOL_RESOURCE_INFO view is periodically sampled and saved to the GS_RESPOOL_RESOURCE_HISTORY system catalog.

Type: SIGHUP

Value range: Boolean

- **on** indicates that the user historical resource monitoring dumping function is enabled.
- **off** indicates that the user historical resource monitoring dumping function is disabled.

Default value: on

user metric retention time

Parameter description: Specifies the retention time of the user historical resource monitoring data.

Type: SIGHUP

Value range: an integer ranging from 0 to 3650. The unit is day.

• If this parameter is set to **0**, user historical resource monitoring data is permanently stored.

• If the value is greater than **0**, user historical resource monitoring data is stored for the specified number of days.

Default value: 7

resource_track_level

Parameter description: Specifies the resource monitoring level of the current session. This parameter is valid only when **enable resource track** is set to **on**.

Type: USERSET

Value range: enumerated values

- none: Resources are not monitored.
- query: enables query-level resource monitoring. If this function is enabled, the plan information (similar to the output information of EXPLAIN) of SQL statements will be recorded in top SQL statements.
- **perf**: enables the perf-level resource monitoring. If this function is enabled, the plan information (similar to the output information of EXPLAIN ANALYZE) that contains the actual execution time and the number of execution rows will be recorded in the top SQL.
- operator_realtime: enables the operator-level resource monitoring. If this
 function is enabled, the operator information of jobs running in real time is
 recorded in the top SQL statements, but is not persisted to the historical top
 SQL statements.
- **operator**: enables the operator-level resource monitoring. If this function is enabled, not only the information including the actual execution time and number of execution rows is recorded in the top SQL statement, but also the operator-level execution information is updated to the top SQL statement.

Default value: query

fast obs tablesize method

Parameter description: Specifies the method for quickly calculating the size of column-store V3 and V3 HStore Opt tables. This parameter is supported only by clusters of version 9.1.0.100 or later.

Type: USERSET

Value range: enumerated values

- **0**: The table size is calculated by listing OBS files.
- 1: The table size is calculated through WLM background statistics using pg_relfilenode_size.
- **2**: The table size is estimated by calculating the maximum offset of each file in cudesc.

Default value: 2

fast obs dbsize method

Parameter description: Specifies the method for quickly calculating the size of database data on OBS. This parameter is supported only by clusters of version 9.1.0.100 or later.

Type: USERSET

Value range: enumerated values

- **0**: The size of the database is directly estimated based on the OBS bucket.
- 1: The size of the entire database is normally calculated in regular mode.

Default value: 0

resource_track_cost

Parameter description: Specifies the minimum execution cost for resource monitoring on statements in the current session. This parameter is valid only when **enable_resource_track** is set to **on**.

Type: USERSET

Value range: an integer ranging from -1 to INT_MAX

- -1 indicates that resource monitoring is disabled.
- A value greater than or equal to **0** indicates that statements whose execution cost exceeds this value will be monitored.

Default value: 0

The default value of this parameter is **0** for a new cluster. In upgrade scenarios, the default value of this parameter is the same as that of the source version.

resource_track_duration

Parameter description: Specifies the minimum statement execution time that determines whether information about jobs of a statement recorded in the real-time view (see Table 12-1) will be dumped to a historical view after the statement is executed. Job information will be dumped from the real-time view (with the suffix statistics) to a historical view (with the suffix history) if the statement execution time is no less than this value. This parameter is valid only when enable_resource_track is set to on.

Type: USERSET

Value range: an integer ranging from 0 to INT_MAX. The unit is second (s).

- 0 indicates that information about all statements recorded in the real-time resource monitoring view (see Table 12-1) will be archived into historical views.
- If the value is greater than **0**, the system archives historical information if the total execution and queuing time of statements in the real-time resource monitoring view (**Table 12-1**) goes over the parameter value.

Default value: 60s

resource_track_subsql_duration

Parameter description: Filters the minimum execution time of substatements in a stored procedure. This parameter is supported only by clusters of version 8.2.1 or later.

If the execution time of a sub-statement in a stored procedure is greater than the value of this parameter, the job information is dumped to the top SQL archive table. This parameter takes effect only when <code>enable_track_record_subsql</code> is set to <code>on</code>

Type: USERSET

Value range: an integer ranging from 0 to INT_MAX. The unit is second (s).

- If the value is **0**, historical information about all substatements in the stored procedure is archived.
- If the value is greater than **0**, historical information is archived when the execution time of a substatement in a stored procedure exceeds the value of this parameter.

Default value: 180s

query_exception_count_limit

Parameter description: Specifies the maximum number of times that a job triggers an exception rule. If the number of times that a job triggers an exception rule reaches the upper limit, the job will be automatically added to the blocklist and cannot be executed anymore. The job can be resumed only after it is removed from the blocklist.

Type: USERSET

Value range: an integer ranging from -1 to INT_MAX

- If the value is -1, the number of times that a job triggers an exception rule is not limited. That is, the job will not be automatically added to blocklist even if it triggers an exception rule for many times.
- If the value is greater than or equal to **0**, the job will be added to the blocklist immediately when the number of times it triggers an exception rule reaches the threshold. The values **0** and **1** both indicate that a job is added to blocklist once the job triggers an exception rule.

Default value: -1

disable_memory_protect

Parameter description: Stops memory protection. To query system views when system memory is insufficient, set this parameter to **on** to stop memory protection. This parameter is used only to diagnose and debug the system when system memory is insufficient. Set it to **off** in other scenarios.

Type: USERSET

Value range: Boolean

• **on** indicates that memory protection stops.

• **off** indicates that memory is protected.

Default value: off

query_band

Parameter description: Specifies the job type of the current session.

Type: USERSET

Value range: a string

Default value: empty

CAUTION

Pay attention to the following when modifying this parameter:

- 1. When this parameter is disabled, it means that the user does not need CCN control function, and the CCN memory negative feedback mechanism will be invalid
- 2. When a job is running, if the value of GUC is changed from **OFF** to **ON**, the CCN memory negative feedback mechanism takes effect. If the concurrency is high, the memory may be temporarily unavailable. After the running job is done, the dynamic load function recovers.
- 3. When a job is running and most jobs are delivered by users from the default resource pool, you are not advised to change the GUC parameter from **enabled** to **disabled**. It may cause a memory error. When there is no job delivered by users from the default resource pool, then you can change the parameter. You are advised to bind a user resource pool before delivering jobs.

wlm memory feedback adjust

Parameter description: Specifies whether to enable memory negative feedback in dynamic load management. (This parameter is supported only by clusters of version 8.2.0 or later.)

Memory is preempted based on the estimated statement memory usage calculated on the CN. If the estimated memory usage of a statement is too high, it will preempt too much memory, causing subsequent jobs to be queued. With the negative memory feedback mechanism, if the cluster memory usage has been overestimated for a period of time, the CCN node will dynamically release some memory for subsequent jobs, improving resource utilization.

Type: SIGHUP

Value range: a string

- on indicates that memory negative feedback is enabled.
- off indicates that memory negative feedback is disabled.
- **on()** enables the memory negative feedback function and specifies the time and estimated memory percentage parameter required to trigger the negative feedback. For example, **on(60,50)** indicates that to trigger the negative feedback mechanism, the memory must be overestimated for 60 consecutive

seconds, and the preempted memory needs must exceed 50% of the available memory. By default, the wait time before the negative feedback mechanism takes effect is 50 seconds. The minimum estimated total memory usage for triggering the mechanism is over 40% of the available system memory.

Default value: on

bbox_dump_count

Parameter description: Specifies the maximum number of core files that are generated by GaussDB(DWS) and can be stored in the path specified by **bbox_dump_path**. If the number of core files exceeds this value, old core files will be deleted. This parameter is valid only if **enable_bbox_dump** is set to **on**.

Type: USERSET

Value range: an integer ranging from 1 to 20

Default value: 8

□ NOTE

When core files are generated during concurrent SQL statement execution, the number of files may be larger than the value of **bbox_dump_count**.

io_limits

Parameter description: This parameter has been discarded in version 8.1.2 and is reserved for compatibility with earlier versions. This parameter is invalid in the current version.

Type: USERSET

Value range: an integer ranging from 0 to 1073741823

Default value: 0

io priority

Parameter description: This parameter has been discarded in version 8.1.2 and is reserved for compatibility with earlier versions. This parameter is invalid in the current version.

Type: USERSET

Value range: enumerated values

- None
- Low
- Medium
- High

Default value: None

session_respool

Parameter description: Specifies the resource pool associated with the current session.

Type: USERSET

If you set **cgroup_name** and then **session_respool**, the Cgroups associated with **session_respool** take effect. If you reverse the order, Cgroups associated with **cgroup_name** take effect.

If the Workload Cgroup level is specified during the **cgroup_name** change, the database does not check the Cgroup level. The level ranges from 1 to 10.

You are not advised to set **cgroup_name** and **session_respool** at the same time.

Value range: a string. This parameter can be set to the resource pool configured through **create resource pool**.

Default value: invalid_pool

enable_transaction_parctl

Parameter description: whether to control transaction block statements and stored procedure statements.

Type: USERSET

Value range: Boolean

- **on**: Transaction block statements and stored procedure statements are controlled.
- **off**: Transaction block statements and stored procedure statements are not controlled.

Default value: on

session_history_memory

Parameter description: Specifies the memory size of a historical query view.

Type: SIGHUP

Value range: an integer ranging from 10 MB to 50% of max_process_memory

Default value: 100MB

topsql_retention_time

Parameter description: Specifies the retention period of historical Top SQL data in the **gs_wlm_session_info** and **gs_wlm_operator_info** tables.

Type: SIGHUP

Value range: an integer ranging from 0 to 3650. The unit is day.

- If it is set to **0**, the data is stored permanently.
- If the value is greater than **0**, the data is stored for the specified number of days.

Default value: 30

CAUTION

- Before setting this GUC parameter to enable the data retention function, delete data from the **gs_wlm_session_info** and **gs_wlm_operator_info** tables.
- The default value of this parameter is **30** for a new cluster. In upgrade scenarios, the default value of this parameter is the same as that of the source version.

transaction_pending_time

Parameter description: maximum queuing time of transaction block statements and stored procedure statements if **enable_transaction_parctl** is set to **on**.

Type: USERSET

Value range: an integer ranging from -1 to INT_MAX. The unit is second (s).

- -1 or 0: No queuing timeout is specified for transaction block statements and stored procedure statements. The statements can be executed when resources are available.
- Value greater than 0: If transaction block statements and stored procedure statements have been queued for a time longer than the specified value, they are forcibly executed regardless of the current resource situation.

Default value: 0

NOTICE

This parameter is valid only for internal statements of stored procedures and transaction blocks. That is, this parameter takes effect only for the statements whose **enqueue** value (for details, see **PG_SESSION_WLMSTAT**) is **Transaction** or **StoredProc**.

enable_concurrency_scaling

Parameter description: Specifies whether to enable the elastic concurrent expansion function. This parameter is supported only by clusters of version 9.1.0.100 or later.

Type: SIGHUP

Value range: Boolean

- **on** indicates that the elastic concurrent expansion function is enabled.
- **off** indicates that the elastic concurrent expansion function is disabled.

Default value: off

concurrency_scaling_max_idle_time

Parameter description: Specifies the maximum idle time of the elastic logical cluster created for concurrent expansion. If it exceeds the set value of this parameter, it will enter the process of destroying the elastic logical cluster. This parameter is supported only by clusters of version 9.1.0.100 or later.

Type: SIGHUP

Value range: an integer ranging from 0 to 60, in minutes. The value **0** indicates that the elastic logical cluster created during concurrent expansion will not be destroyed. Exercise caution when setting this value.

Default value: 5

concurrency_scaling_limit_per_main_vw

Parameter description: Specifies the maximum number of concurrent elastic logical clusters that can be created for each classic logical cluster. This parameter is supported only by clusters of version 9.1.0.100 or later.

Type: SIGHUP

Value range: an integer ranging from 0 to 32

Default value: 5

concurrency_scaling_max_vw_active_statements

Parameter description: Specifies the maximum number of concurrent tasks that can be executed in an elastic logical cluster. This parameter is supported only by clusters of version 9.1.0.100 or later.

Type: SIGHUP

Value range: an integer ranging from -1 to INT_MAX

- Value **0** indicates that the elastic logical cluster does not execute any jobs. Exercise caution when setting this value.
- -1 indicates that the elastic logical cluster is not under concurrency control. Exercise caution when setting this value. This value is supported only by clusters of version 9.1.0.200 or later.

Default value: 60

concurrency_scaling_max_waiting_statements

Parameter description: Specifies the number of queued elastic jobs in the global queue that trigger the creation process of the elastic logical cluster for concurrent expansion. This parameter is supported only by clusters of version 9.1.0.100 or later.

If the number of queued elastic jobs is greater than or equal to the value set for this parameter, and the number of concurrently created elastic logical clusters does not exceed the value set for **concurrency_scaling_limit_per_main_vw**, the process of creating a concurrent expansion elastic logical cluster will be triggered automatically.

Type: SIGHUP

Value range: an integer ranging from 0 to INT_MAX. The value **0** indicates that an elastic logical cluster is created even if there are no queued jobs. If set to **0**, a large number of resources are consumed. Exercise caution when setting this value.

Default value: 10

15.13 Automatic Cleanup

The automatic cleanup process in the system automatically runs the **VACUUM** and **ANALYZE** statements to reclaim the record space marked as deleted and update statistics in the table.

autovacuum_compaction_rows_limit

Parameter description: Specifies the small CU threshold. A CU whose number of live tuples is less than the value of this parameter is considered as a small CU.

Type: USERSET

Value range: an integer ranging from -1 to 5000

Default value: 2500

<u>A</u> CAUTION

- You are advised not to modify this parameter. If you need to modify this parameter, contact technical support.
- If the version is earlier than 9.1.0.100, do not set this parameter, or there may be duplicate data in primary keys.
- If the version is earlier than 9.1.0.100, value **-1** indicates that the 0 CU switch is disabled.
- In version 9.1.0.100, the default value of this parameter is **0**.
- In 9.1.0.200 and later versions, the default value of this parameter is **2500**.

autoanalyze_mode

Parameter description: Specifies the autoanalyze mode. This parameter is supported by clusters of version 8.2.0 or later.

Type: USERSET

Value range: enumerated values

- normal indicates common autoanalyze.
- light indicates lightweight autoanalyze.

Default value:

• If the current cluster is upgraded from an earlier version to 8.2.0, the default value is **normal** to ensure forward compatibility.

If the cluster version 8.2.0 is newly installed, the default value is light.

analyze_stats_mode

Parameter description: Specifies the mode for **ANALYZE** to calculate statistics.

Type: USERSET

Value range: enumerated values

- **memory** indicates that the memory is forcibly used to calculate statistics. Multi-column statistics are not calculated.
- **sample_table** indicates that temporary sampling tables are forcibly used to calculate statistics. Temporary tables do not support this mode.
- dynamic indicates that the statistics calculation mode is determined based on the size of maintenance_work_mem. If maintenance_work_mem can store samples, the memory mode is used. Otherwise, the temporary sampling table mode is used.

Default value:

- If the current cluster is upgraded from an earlier version to 8.2.0.100, the default value is **memory** to ensure forward compatibility.
- If the cluster version 8.2.0.100 is newly installed, the default value is **dynamic**.

analyze_sample_mode

Parameter description: Specifies the sampling model used by ANALYZE.

Type: USERSET

Value range: an integer ranging from 0 to 2

- **0** indicates the default reservoir sampling.
- 1 indicates the optimized reservoir sampling.
- 2 indicates range sampling.

Default value: 0

autovacuum_max_workers

Parameter description: Specifies the maximum number of automatic cleanup threads running at the same time.

Type: SIGHUP

Value range: an integer ranging from 0 to 128. **0** indicates that **autovacuum** is

disabled.

Default value: 6

□ NOTE

This parameter works with autovacuum. The rules for clearing system catalogs and user tables are as follows:

- When autovacuum_max_workers is set to 0, autovacuum is disabled and no tables are cleared.
- When autovacuum_max_workers is set to a value greater than 0 and autovacuum is set to off, the system only clears the system catalogs and column-store tables with delta tables enabled (such as vacuum delta tables, vacuum cudesc tables, and delta merge).
- When autovacuum_max_workers is set to a value greater than 0 and autovacuum is set to on, all tables will be cleared.

autovacuum_max_workers_hstore

Parameter description: Specifies the maximum number of concurrent automatic cleanup threads used for hstore tables in **autovacuum_max_workers**.

Type: SIGHUP

Value range: an integer ranging from 0 to 128. **0** indicates that the automatic cleanup function of HStore tables is disabled.

Default value: 3

□ NOTE

To use HStore tables, set the following parameters, or the HStore performance will deteriorate severely. The recommended settings are as follows:

autovacuum_max_workers_hstore=3, autovacuum_max_workers=6, autovacuum=true

autovacuum naptime

Parameter description: Specifies the interval between two automatic cleanup operations.

Type: SIGHUP

Value range: an integer ranging from 1 to 2147483. The unit is second.

Default value: 60s

autovacuum_vacuum_cost_delay

Parameter description: Specifies the value of the cost delay used in the **autovacuum** operation.

Type: SIGHUP

Value range: an integer ranging from –1 to 100. The unit is ms. **–1** indicates that the normal vacuum cost delay is used.

Default value: 2ms

check crossvw write

Parameter description: Specifies whether to enable cross-VW write detection. This parameter is supported only by clusters of version 9.1.0.100 or later.

Type: USERSET

Value range: an integer, -1 or 1.

- The value **-1** indicates that it is compatible with the capabilities of version 9.0.3. For the v3 table vacuum, it only clears non-last files for all epochs.
- The value 1 indicates checking whether it is a cross-VW write scenario. For the v3 table vacuum, if it is determined to be a non-cross-VW write scenario, it clears non-last files for all epochs, clears the last file for the current epoch, and clears the last file for epochs that are less than the current epoch. If it is determined to be a cross-VW write scenario, CNs will obtain epoch information from all DNs and package it into an epochList to be sent to the metadata VW. The v3 table vacuum will clear non-last files for all epochs and clear the last file for epochs that are less than max{epochList} and not in epochList.

Default value: 1

enable_pg_stat_object

Parameter description: Specifies whether **AUTO VACUUM** updates the **PG_STAT_OBJECT** system catalog. This parameter is supported only by clusters of version 8.2.1 or later.

Type: USERSET

Value range: Boolean

- on indicates that the PG_STAT_OBJECT system catalog is updated during AUTO VACUUM.
- **off** indicates that the **PG_STAT_OBJECT** system catalog is not updated during **AUTO VACUUM**.

Default value: on

15.14 Default Settings of Client Connection

15.14.1 Statement Behavior

This section describes related default parameters involved in the execution of SQL statements.

search_path

Parameter description: Specifies the order in which schemas are searched when an object is referenced with no schema specified. The value of this parameter consists of one or more schema names. Different schema names are separated by commas (,).

Type: USERSET

If the schema of a temporary table exists in the current session, the scheme can be listed in search_path by using the alias pg_temp, for example,
 'pg_temp,public'. The schema of a temporary table has the highest search

priority and is always searched before all the schemas specified in **pg_catalog** and **search_path**. Therefore, do not explicitly specify **pg_temp** to be searched after other schemas in **search_path**. This setting will not take effect and an error message will be displayed. If the alias **pg_temp** is used, the temporary schema will be only searched for database objects, including tables, views, and data types. Functions or operator names will not be searched for.

- The schema of a system catalog, pg_catalog, has the second highest search priority and is the first to be searched among all the schemas, excluding pg_temp, specified in search_path. Therefore, do not explicitly specify pg_catalog to be searched after other schemas in search_path. This setting will not take effect and an error message will be displayed.
- When an object is created without specifying a particular schema, the object will be placed in the first valid schema listed in **search_path**. An error will be reported if the search path is empty.
- The current effective value of the search path can be examined through the SQL function current_schema. This is different from examining the value of search_path, because the current_schema function displays the first valid schema name in search_path.

Value range: a string

□ NOTE

- When this parameter is set to "**\$user**", **public**, a database can be shared (where no users have private schemas, and all share use of public), and private per-user schemas and combinations of them are supported. Other effects can be obtained by modifying the default search path setting, either globally or per-user.
- When this parameter is set to a null string ("), the system automatically converts it into a pair of double quotation marks ("").
- If the content contains double quotation marks, the system considers them as insecure characters and converts each double quotation mark into a pair of double quotation marks.

Default value: "\$user",public

■ NOTE

\$user indicates the name of the schema with the same name as the current session user. If the schema does not exist, **\$user** will be ignored.

current schema

Parameter description: Specifies the current schema.

Type: USERSET

Value range: a string

Default value: "\$user",public

□ NOTE

\$user indicates the name of the schema with the same name as the current session user. If the schema does not exist, **\$user** will be ignored.

default_tablespace

Parameter description: Specifies the default tablespace of the created objects (tables and indexes) when a **CREATE** command does not explicitly specify a tablespace.

- The value of this parameter is either the name of a tablespace, or an empty string that specifies the use of the default tablespace of the current database. If a non-default tablespace is specified, users must have CREATE privilege for it. Otherwise, creation attempts will fail.
- This parameter is not used for temporary tables. For them, the temp_tablespaces is consulted instead.
- This parameter is not used when users create databases. By default, a new database inherits its tablespace setting from the template database.

Type: USERSET

Value range: a string. An empty string indicates that the default tablespace is

used.

Default value: empty

default storage nodegroup

Parameter description: Specifies the Node Group where a table is created by default. This parameter takes effect only for ordinary tables.

Type: USERSET

Value range: a string

- **installation**: indicates that the table is created in the installed Node Group by default.
- random_node_group: indicates that the table is created in a randomly selected Node Group by default. This feature is supported in 8.1.2 or later and is used only in the test environment.
- **roach_group**: indicates that the table is created in all nodes by default. This value is reserved for the Roach tool and cannot be used in other scenarios.
- A value other than the preceding three options indicates that the table is created in a specified Node Group.

Default value: installation

temp_tablespaces

Parameter description: Specifies tablespaces to which temporary objects will be created (temporary tables and their indexes) when a **CREATE** command does not explicitly specify a tablespace. Temporary files for sorting large data are created in these tablespaces.

The value of this parameter is a list of names of tablespaces. When there is more than one name in the list, GaussDB(DWS) chooses a random tablespace from the list upon the creation of a temporary object each time. Except that within a transaction, successively created temporary objects are placed in successive tablespaces in the list. If the element selected from the list is an empty string,

GaussDB(DWS) will automatically use the default tablespace of the current database instead.

Type: USERSET

Value range: a string An empty string indicates that all temporary objects are created only in the default tablespace of the current database. For details, see **default tablespace**.

Default value: empty

check function bodies

Parameter description: Specifies whether to enable validation of the function body string during the execution of **CREATE FUNCTION**. Verification is occasionally disabled to avoid problems, such as forward references when you restore function definitions from a dump.

Type: USERSET

Value range: Boolean

- **on** indicates that validation of the function body string is enabled during the execution of **CREATE FUNCTION**.
- **off** indicates that validation of the function body string is disabled during the execution of **CREATE FUNCTION**.

Default value: on

default_transaction_isolation

Parameter description: Specifies the default isolation level of each transaction.

Type: USERSET

Value range: enumerated values

- **READ COMMITTED**: Only committed data is read. This is the default.
- READ UNCOMMITTED: GaussDB(DWS) does not support READ UNCOMMITTED. If READ UNCOMMITTED is set, READ COMMITTED is executed instead.
- REPEATABLE READ: Only the data committed before transaction start is read.
 Uncommitted data or data committed in other concurrent transactions cannot be read
- SERIALIZABLE: GaussDB(DWS) does not support SERIALIZABLE. If SERIALIZABLE is set, REPEATABLE READ is executed instead.

Default value: READ COMMITTED

default_transaction_read_only_probe

Parameter description: Specifies whether to terminate the execution of special statements (e.g., statements for flushing data to disks and generating new tables or physical files) when the database is about to become read-only (disk usage reaches 90%). The CM module checks and sets the disk usage threshold. It is not

advised to set this parameter. This is supported only by clusters of version 9.1.0.200 or later.

Type: USERSET

Value range: Boolean

- **on** indicates that the execution of the special statement is terminated.
- **off** indicates that the execution of the special statement is not terminated.

Default value: off

default_transaction_deferrable

Parameter description: Specifies the default delaying state of each new transaction. It currently has no effect on read-only transactions or those running at isolation levels lower than serializable.

GaussDB(DWS) does not support the serializable isolation level of each transaction. The parameter is insignificant.

Type: USERSET

Value range: Boolean

- **on** indicates a transaction is delayed by default.
- off indicates a transaction is not delayed by default.

Default value: off

session_replication_role

Parameter description: Specifies the behavior of replication-related triggers and rules for the current session.

Type: USERSET

NOTICE

Setting this parameter will discard all the cached execution plans.

Value range: enumerated values

- **origin** indicates that the system copies operations such as insert, delete, and update from the current session.
- **replica** indicates that the system copies operations such as insert, delete, and update from other places to the current session.
- **local** indicates that the system will detect the role that has logged in to the database when using the function to copy operations and will perform related operations.

Default value: origin

statement timeout

Parameter description: If the statement execution time (starting when the server receives the command) is longer than the duration specified by the parameter, error information is displayed when you attempt to execute the statement and the statement then exits.

Type: USERSET

Value range: an integer ranging from 0 to 2147483647. The unit is ms.

Default value:

- If the current cluster is upgraded from an earlier version to 8.2.0, the value in the earlier version is inherited. The default value is **0**.
- If the cluster version 8.2.0 is newly installed, the default value is 24h.

vacuum_freeze_min_age

Parameter description: Specifies the minimum cutoff age (in the same transaction), based on which **VACUUM** decides whether to replace transaction IDs with FrozenXID while scanning a table.

Type: USERSET

Value range: an integer from 0 to 576460752303423487.

□ NOTE

Although you can set this parameter to a value ranging from **0** to **1000000000** anytime, **VACUUM** will limit the effective value to half the value of **autovacuum_freeze_max_age** by default.

Default value: 5000000000

vacuum_freeze_table_age

Parameter description: Specifies the time that VACUUM freezes tuples while scanning the whole table. **VACUUM** performs a whole-table scan if the value of the **pg_class.relfrozenxid** column of the table has reached the specified time.

Type: USERSET

Value range: an integer from 0 to 576460752303423487.

Ⅲ NOTE

Although users can set this parameter to a value ranging from **0** to **2000000000** anytime, **VACUUM** will limit the effective value to 95% of **autovacuum_freeze_max_age** by default. Therefore, a periodic manual VACUUM has a chance to run before an anti-wraparound autovacuum is launched for the table.

Default value: 15000000000

bytea_output

Parameter description: Specifies the output format for values of the bytea type.

Type: USERSET

Value range: enumerated values

- hex indicates the binary data is converted to the two-byte hexadecimal digit.
- escape indicates the traditional PostgreSQL format is used. It takes the
 approach of representing a binary string as a sequence of ASCII characters,
 while converting those bytes that cannot be represented as an ASCII character
 into special escape sequences.

Default value: hex

xmlbinary

Parameter description: Specifies how binary values are to be encoded in XML.

Type: USERSET

Value range: enumerated values

base64hex

Default value: base64

xmloption

Parameter description: Specifies whether DOCUMENT or CONTENT is implicit when converting between XML and string values.

Type: USERSET

Value range: enumerated values

- document indicates an HTML document.
- content indicates a common string.

Default value: content

gin_pending_list_limit

Parameter description: Specifies the maximum capacity of the pending list when **FASTUPDATE** is enabled for GIN indexes.

Pending List is a data structure specific to GIN indexes and is used to temporarily store index update operations. When the **FASTUPDATE** parameter is enabled for a GIN index, new index items are not directly written to the main index structure, but are stored in the pending list. When conditions are met, the index items are merged into the main index in batches.

If the pending list grows larger than this maximum size, data in the list will be moved to the GIN index data structure in batches. This setting can be overridden for individual GIN indexes by modifying index storage parameters.

Type: USERSET

Value range: an integer ranging from 64 to INT_MAX. The unit is KB.

Default value: 4 MB

15.14.2 Zone and Formatting

This section describes parameters related to the time format setting.

DateStyle

Parameter description: Specifies the display format for date and time values, as well as the rules for interpreting ambiguous date input values.

This variable contains two independent components: the output format specifications (ISO, Postgres, SQL, or German) and the input/output order of year/month/day (DMY, MDY, or YMD). The two components can be set separately or together. The keywords Euro and European are synonyms for DMY; the keywords US, NonEuro, and NonEuropean are synonyms for MDY.

Type: USERSET

Value range: a string

Default value: ISO, MDY

□ NOTE

gs_initdb will initialize this parameter so that its value is the same as that of lc_time.

Suggestion: The ISO format is recommended. Postgres, SQL, and German use abbreviations for time zones, such as **EST**, **WST**, and **CST**.

IntervalStyle

Parameter description: Specifies the display format for interval values.

Type: USERSET

Value range: enumerated values

- **sql_standard** indicates that output matching SQL standards will be generated.
- **postgres** indicates that output matching PostgreSQL 8.4 will be generated when the **DateStyle** parameter is set to **ISO**.
- **postgres_verbose** indicates that output matching PostgreSQL 8.4 will be generated when the **DateStyle** parameter is set to **non_ISO**.
- **iso_8601** indicates that output matching the time interval "format with designators" defined in ISO 8601 will be generated.
- **oracle** indicates the output result that matches the numtodsinterval function in the Oracle database. For details, see numtodsinterval.

NOTICE

The **IntervalStyle** parameter also affects the interpretation of ambiguous interval input.

Default value: postgres

TimeZone

Parameter description: Specifies the time zone for displaying and interpreting time stamps.

Type: USERSET

Value range: a string. You can obtain it by querying the **pg_timezone_names**

view.

Default value: UTC

◯ NOTE

gs_initdb will set a time zone value that is consistent with the system environment.

timezone_abbreviations

Parameter description: Specifies the time zone abbreviations that will be accepted by the server.

Type: USERSET

Value range: a string. You can obtain it by querying the pg_timezone_names view.

Default value: Default

∩ NOTE

Default indicates an abbreviation that works in most of the world. There are also other abbreviations, such as **Australia** and **India** that can be defined for a particular installation.

extra_float_digits

Parameter description: Specifies the number of digits displayed for floating-point values, including float4, float8, and geometric data types. The parameter value is added to the standard number of digits (FLT DIG or DBL DIG as appropriate).

Type: USERSET

Value range: an integer ranging from -15 to 3

- This parameter can be set to **3** to include partially-significant digits. It is especially useful for dumping float data that needs to be restored exactly.
- This parameter can also be set to a negative value to suppress unwanted digits.

Default value: 0

client_encoding

Parameter description: Specifies the client-side encoding type (character set).

Set this parameter as needed. Try to keep the client code and server code consistent to improve efficiency.

Type: USERSET

Value range: encoding compatible with PostgreSQL. **UTF8** indicates that the database encoding is used.

□ NOTE

- You can run the **locale -a** command to check and set the system-supported zone and the corresponding encoding format.
- By default, **gs_initdb** will initialize the setting of this parameter based on the current system environment. You can also run the **locale** command to check the current configuration environment.
- To use consistent encoding for communication within a cluster, you are advised to retain the default value of **client_encoding**. Modification to this parameter in the **postgresql.conf** file (by using the **gs_guc** tool, for example) does not take effect.

Default value: UTF8

Recommended value: SQL_ASCII or UTF8

lc_messages

Parameter description: Specifies the language in which messages are displayed.

Valid values depend on the current system. On some systems, this zone category does not exist. Setting this variable will still work, but there will be no effect. In addition, translated messages for the desired language may not exist. In this case, you can still see the English messages.

Type: SUSET

Value range: a string

■ NOTE

- You can run the **locale -a** command to check and set the system-supported zone and the corresponding encoding format.
- By default, **gs_initdb** will initialize the setting of this parameter based on the current system environment. You can also run the **locale** command to check the current configuration environment.

Default value: C

lc_monetary

Parameter description: Specifies the display format of monetary values. It affects the output of functions such as to_char. Valid values depend on the current system.

Type: USERSET

Value range: a string

□ NOTE

- You can run the **locale -a** command to check and set the system-supported zone and the corresponding encoding format.
- By default, gs_initdb will initialize the setting of this parameter based on the current system environment. You can also run the locale command to check the current configuration environment.

Default value: C

lc_numeric

Parameter description: Specifies the display format of numbers. It affects the output of functions such as to_char. Valid values depend on the current system.

Type: USERSET

Value range: a string

□ NOTE

- You can run the **locale -a** command to check and set the system-supported zone and the corresponding encoding format.
- By default, **gs_initdb** will initialize the setting of this parameter based on the current system environment. You can also run the **locale** command to check the current configuration environment.

Default value: C

lc_time

Parameter description: Specifies the display format of time and zones. It affects the output of functions such as to_char. Valid values depend on the current system.

Type: USERSET

Value range: a string

□ NOTE

- You can run the **locale -a** command to check and set the system-supported zone and the corresponding encoding format.
- By default, **gs_initdb** will initialize the setting of this parameter based on the current system environment. You can also run the **locale** command to check the current configuration environment.

Default value: C

default_text_search_config

Parameter description: Specifies the text search configuration.

If the specified text search configuration does not exist, an error will be reported. If the specified text search configuration is deleted, set

default_text_search_config again. Otherwise, an error will be reported, indicating incorrect configuration.

- The text search configuration is used by text search functions that do not have an explicit argument specifying the configuration.
- When a configuration file matching the environment is determined, gs_initdb will initialize the configuration file with a setting that corresponds to the environment.

Type: USERSET

Value range: a string

□ NOTE

GaussDB(DWS) supports the following two configurations: pg_catalog.english and pg_catalog.simple.

Default value: pg_catalog.english

15.14.3 Other Default Parameters

This section describes the default database loading parameters of the database system.

dynamic_library_path

Parameter description: Specifies the path for saving the shared database files that are dynamically loaded for data searching. When a dynamically loaded module needs to be opened and the file name specified in the **CREATE FUNCTION** or **LOAD** command does not have a directory component, the system will search this path for the required file.

The value of **dynamic_library_path** must be a list of absolute paths separated by colons (:) or by semi-colons (;) on the Windows OS. The special variable **\$libdir** in the beginning of a path will be replaced with the module installation directory provided by GaussDB(DWS). Example:

dynamic_library_path = '/usr/local/lib/postgresql:/opt/testgs/lib:\$libdir'

Type: SUSET

Value range: a string

□ NOTE

If the value of this parameter is set to an empty character string, the automatic path search is turned off.

Default value: \$libdir

gin_fuzzy_search_limit

Parameter description: Specifies the upper limit of the size of the set returned by GIN indexes.

Type: USERSET

Value range: an integer ranging from 0 to INT_MAX. The value 0 indicates no

limit.

Default value: 0

15.15 Lock Management

In GaussDB(DWS), concurrent transactions may cause single-node deadlocks or distributed deadlocks due to resource competition. This section describes parameters used for managing transaction lock mechanisms.

deadlock timeout

Parameter description: Specifies the time, in milliseconds, to wait on a lock before checking whether there is a deadlock condition. When the applied lock exceeds the preset value, the system will check whether a deadlock occurs.

- The check for deadlock is relatively expensive. Therefore, the server does not check it when waiting for a lock every time. Deadlocks do not frequently occur when the system is running. Therefore, the system just needs to wait on the lock for a while before checking for a deadlock. Increasing this value reduces the time wasted in needless deadlock checks, but slows down reporting of real deadlock errors. On a heavily loaded server, you may need to raise it. The value you have set needs to exceed the transaction time. By doing this, the possibility that a lock will be released before the waiter decides to check for deadlocks will be reduced.
- When log_lock_waits is set, this parameter also determines the duration you need to wait before a log message about the lock wait is issued. If you are trying to investigate locking delays, you need to set this parameter to a value smaller than normal deadlock timeout.

Type: SUSET

Value range: an integer ranging from 1 to 2147483647. The unit is millisecond

(ms).

Default value: 1s

ddl_lock_timeout

Parameter description: Indicates the number of seconds a DDL command should wait for the locks to become available. If the time spent in waiting for a lock exceeds the specified time, an error is reported. This parameter is supported only by clusters of version 8.1.3.200 or later.

Type: SUSET

Value range: an integer ranging from 0 to INT_MAX. The unit is millisecond (ms).

- If the value of this parameter is 0, this parameter does not take effect.
- If the value of this parameter is greater than 0, the lock wait time of DDL statements is the value of this parameter, and the lock wait time of other locks is the value of **lockwait_timeout**.

Default value: 0

□ NOTE

This parameter has a higher priority than **lockwait_timeout** and takes effect only for **AccessExclusiveLock**.

ddl select concurrent mode

Parameter description: Specifies the concurrency mode of DDL and **SELECT** statements. This parameter is supported only by clusters of version 8.1.3.320, 8.2.1, or later.

Type: SUSET

Value range: a string

- **none**: DDL and select statements cannot be executed concurrently. Waiting statements are in the lock wait state.
- truncate: When the TRUNCATE statement is blocked by the SELECT statement, the SELECT statement is interrupted and the TRUNCATE statement is executed first.
- exchange: When the EXCHANGE statement is blocked by the SELECT statement, the SELECT statement is interrupted and the EXCHANGE statement is executed first.
- vacuum_full: When the vacuum_full statement is blocked by the SELECT statement, the SELECT statement is interrupted and the vacuum_full statement is executed first. This is supported only by clusters of version 9.1.0.200 or later.
- insert_overwrite: When the insert_overwrite statement is blocked by the SELECT statement, the SELECT statement is interrupted and the insert_overwrite statement is executed first. This is supported only by clusters of version 9.1.0.200 or later.

Default value: none

∩ NOTE

- To reserve time for the SELECT statement to respond to signals, if the value of ddl lock timeout is less than 1 second in the current version, 1 second is used.
- Concurrency is not supported when there are conflicts with locks of higher levels (more than one level). For example, autoanalyze is triggered by SELECT when autoanalyze_mode is set to normal.
- This parameter allows for SELECT statements in either a single statement or a
 transaction block. However, in other versions, it only supports SELECT statements in a
 single statement. For concurrent SELECT operations in a single statement or transaction
 block, learn more information in the description of parameter
 enable_cancel_select_in_txnblock.
- Values other than none can be used together. For example, if this parameter is set to truncate, exchange, the TRUNCATE and EXCHANGE statements are blocked by the SELECT statement. The SELECT statement is interrupted and executed first.

enable cancel select in txnblock

Parameter description: Specifies whether the **SELECT** statement in a transaction block can be interrupted. This parameter is supported only by clusters of version 8.2.1, 9.1.0.200, or later.

Type: USERSET

Value range: Boolean

- **on** indicates that the select operation in the transaction block can be interrupted.
- **off** indicates that the select operation in the transaction block cannot be interrupted.

Default value: on

□ NOTE

- This parameter controls whether the **SELECT** statement in a transaction block can be interrupted by the DDL operation specified in **ddl_select_concurrent_mode**.
- The ddl_select_concurrent_mode parameter controls DDL statements such as TRUNCATE and EXCHANGE, and the enable_cancel_select_in_txnblock parameter controls SELECT statements.

lockwait timeout

Parameter description: Specifies the longest time to wait before a single lock times out. If the time you wait before acquiring a lock exceeds the specified time, an error is reported.

Type: SUSET

Value range: an integer ranging from 0 to INT_MAX. The unit is millisecond (ms).

Default value: 20 min

update_lockwait_timeout

Parameter description: sets the maximum duration that a lock waits for concurrent updates on a row to complete when the concurrent update feature is enabled. If the time you wait before acquiring a lock exceeds the specified time, an error is reported.

Type: SUSET

Value range: an integer ranging from 0 to INT_MAX. The unit is millisecond (ms).

Default value: 2min

partition_lock_upgrade_timeout

Parameter description: Specifies the time to wait before the attempt of a lock upgrade from ExclusiveLock to AccessExclusiveLock times out on partitions.

- When you do MERGE PARTITION and CLUSTER PARTITION on a partitioned table, temporary tables are used for data rearrangement and file exchange. To concurrently perform as many operations as possible on the partitions, ExclusiveLock is acquired for the partitions during data rearrangement and AccessExclusiveLock is acquired during file exchange.
- Generally, a partition waits until it acquires a lock, or a timeout occurs if the
 partition waits for a period of time longer than specified by the
 lockwait_timeout parameter.
- When doing MERGE PARTITION or CLUSTER PARTITION on a partitioned table, you need to acquire AccessExclusiveLock during file exchange. If the lock fails to be acquired, the acquisition is retried in 50 ms. This parameter specifies the time to wait before the lock acquisition attempt times out.
- If this parameter is set to -1, the lock upgrade never times out. The lock upgrade is continuously retried until it succeeds.

Type: USERSET

Value range: an integer ranging from -1 to 3000, in seconds

Default value: 1800

enable_release_scan_lock

Parameter description: Specifies whether a SELECT statement releases a level-1 lock after the statement execution is complete. This parameter reduces DDL conflicts with SELECT locks within transaction blocks. This parameter is supported only by clusters of version 8.3.0 or later.

Type: USERSET

Value range: Boolean

- on indicates that DDL operations will be blocked to wait for the release of cluster locks. The SELECT statement releases the level-1 lock after it finishes, not when the transaction commits.
- off indicates that DDL operations will not be blocked.

Default value: off

vacuum_full_interruptible

Parameter description: Controls the behavior that the **VACUUM FULL** statement gives a lock to other statements. This is supported only by clusters of version 9.1.0.200 or later.

Type: USERSET

Value range: Boolean

- **on** indicates that DDL operations will be blocked to wait for the release of cluster locks. When **VACUUM FULL** blocks other statements, it interrupts the execution and gives the lock to other statements.
- off indicates that DDL operations will not be blocked. When VACUUM FULL
 blocks other statements, it does not interrupt the execution. Other statements
 can be executed only after VACUUM FULL has completed and released the
 lock.

Default value: off

15.16 Version and Platform Compatibility

15.16.1 Compatibility with Earlier Versions

GaussDB(DWS) provides parameter controls for the downward compatibility and external compatibility features of the database. The backward compatibility of the database system can provide support for old versions of database applications. The parameters introduced in this section mainly control the backward compatibility of the database.

array_nulls

Parameter description: Determines whether the array input parser recognizes unquoted NULL as a null array element.

Type: USERSET

Value range: Boolean

- **on** indicates that null values can be entered in arrays.
- **off** indicates backward compatibility with the old behavior. Arrays containing **NULL** values can still be created when this parameter is set to **off**.

Default value: on

backslash quote

Parameter description: Determines whether a single quotation mark can be represented by \' in a string text.

Type: USERSET

NOTICE

When the string text meets the SQL standards, \ has no other meanings. This parameter only affects the handling of non-standard-conforming string texts, including escape string syntax (E'...').

Value range: enumerated values

- **on** indicates that the use of \' is always allowed.
- **off** indicates that the use of \' is rejected.
- **safe_encoding** indicates that the use of \' is allowed only when client encoding does not allow ASCII \ within a multibyte character.

Default value: safe_encoding

default with oids

Parameter description: Determines whether CREATE TABLE and CREATE TABLE AS include an OID field in newly-created tables if neither WITH OIDS nor WITHOUT OIDS is specified. It also determines whether OIDs will be included in tables created by SELECT INTO.

It is not recommended that OIDs be used in user tables. Therefore, this parameter is set to **off** by default. When OIDs are required for a particular table, **WITH OIDS** needs to be specified during the table creation.

Type: USERSET

Value range: Boolean

- on indicates CREATE TABLE and CREATE TABLE AS can include an OID field in newly-created tables.
- **off** indicates **CREATE TABLE** and **CREATE TABLE AS** cannot include any **OID** field in newly-created tables.

Default value: off

escape_string_warning

Parameter description: Specifies a warning on directly using a backslash (\) as an escape in an ordinary character string.

- Applications that wish to use a backslash (\) as an escape need to be modified to use escape string syntax (E'...'). This is because the default behavior of ordinary character strings is now to treat the backslash as an ordinary character in each SQL standard.
- This variable can be enabled to help locate codes that need to be changed.

Type: USERSET

Value range: Boolean

Default value: on

lo_compat_privileges

Parameter description: Determines whether to enable backward compatibility for the privilege check of large objects.

Type: SUSET

Value range: Boolean

on indicates that the privilege check is disabled when users read or modify large objects. This setting is compatible with versions earlier than PostgreSQL 9.0.

Default value: off

quote_all_identifiers

Parameter description: When the database generates SQL, this parameter forcibly quotes all identifiers even if they are not keywords. This will affect the output of EXPLAIN as well as the results of functions, such as pg_get_viewdef. For details, see the **--quote-all-identifiers** parameter of **gs_dump**.

Type: USERSET

Value range: Boolean

- **on** indicates the forcible quotation function is enabled.
- **off** indicates the forcible quotation function is disabled.

Default value: off

sql_inheritance

Parameter description: Determines whether to inherit semantics.

Type: USERSET

Value range: Boolean

off indicates that child tables cannot be accessed by various commands. That is, an ONLY keyword is used by default. This setting is compatible with versions earlier than PostgreSQL 7.1.

Default value: on

standard_conforming_strings

Parameter description: Determines whether ordinary string texts ('...') treat backslashes as ordinary texts as specified in the SQL standard.

- Applications can check this parameter to determine how string texts will be processed.
- It is recommended that characters be escaped by using the escape string syntax (E'...').

Type: USERSET

Value range: Boolean

- **on** indicates that the function is enabled.
- off indicates that the function is disabled.

Default value: on

synchronize_seqscans

Parameter description: Controls sequential scans of tables to synchronize with each other. Concurrent scans read the same data block about at the same time and share the I/O workload.

Type: USERSET

Value range: Boolean

- **on** indicates that a scan may start in the middle of the table and then "wrap around" the end to cover all rows to synchronize with the activity of scans already in progress. This may result in unpredictable changes in the row ordering returned by queries that have no ORDER BY clause.
- **off** indicates that the scan always starts from the table heading.

Default value: on

enable_beta_features

Parameter description: Controls whether certain limited features, such as GDS table join, are available. These features are not explicitly prohibited in earlier versions, but are not recommended due to their limitations in certain scenarios.

Type: USERSET

Value range: Boolean

- **on** indicates that the features are enabled and forward compatible, but may incur errors in certain scenarios.
- off indicates that the features are disabled.

Default value: off

15.16.2 Platform and Client Compatibility

Database systems are widely used across many platforms, and their external compatibility offers a great deal of convenience.

transform_null_equals

Parameter description: Determines whether expressions of the form expr = NULL (or NULL = expr) are treated as expr IS NULL. They return true if expr evaluates to **NULL**, and false otherwise.

- The correct SQL-standard-compliant behavior of expr = NULL is to always return null (unknown).
- Filtered forms in Microsoft Access generate queries that appear to use expr = NULL to test for null values. If you turn this option on, you can use this interface to access the database.

Type: USERSET

Value range: Boolean

- **on** indicates expressions of the form expr = NULL (or NULL = expr) are treated as expr IS NULL.
- off indicates expr = NULL always returns NULL.

Default value: off

∩ NOTE

New users are always confused about the semantics of expressions involving **NULL** values. Therefore, **off** is used as the default value.

td_compatible_truncation

Parameter description: Determines whether to enable features compatible with a Teradata database. You can set this parameter to **on** when connecting to a database compatible with the Teradata database, so that when you perform the INSERT operation, overlong strings are truncated based on the allowed maximum length before being inserted into char- and varchar-type columns in the target table. This ensures all data is inserted into the target table without errors reported.

□ NOTE

- The string truncation function cannot be used if the **INSERT** statement includes a foreign table.
- If inserting multi-byte character data (such as Chinese characters) to database with the
 character set byte encoding (SQL_ASCII, LATIN1), and the character data crosses the
 truncation position, the string is truncated based on its bytes instead of characters.
 Unexpected result will occur in tail after the truncation. If you want correct truncation
 result, you are advised to adopt encoding set such as UTF8, which has no character data
 crossing the truncation position.

Type: USERSET

Value range: Boolean

- **on** indicates overlong strings are truncated.
- **off** indicates overlong strings are not truncated.

Default value: off

behavior_compat_options

Parameter description: Specifies the database compatibility behavior, which consists of multiple items separated by commas (,). Compatibility configurations are applied according to the database type (such as Oracle/Teradata/MySQL). For details, see **Table 15-2**.

Type: USERSET

Value range: a string

Default value: In upgrade scenarios, the default value of this parameter is the same as that in the cluster before the upgrade. When a new cluster is installed, the default value of this parameter is

check_function_conflicts,check_function_shippable,unsupported_set_function_case to prevent serious issues caused by incorrect function attributes that users define.

□ NOTE

- Currently, only items in **Table 15-2** are supported.
- Multiple items are separated by commas (,), for example, set behavior_compat_options='end_month_calculate,display_leading_zero';.
- **strict_concat_functions** and **strict_text_concat_td** are mutually exclusive.
- You are not advised to set **behavior_compat_options** to **'return_null_string'** in Oracle compatibility mode. If this option is set, do not insert query results into tables.

Table 15-2 Compatibility configuration items

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
display_leadi ng_zero	Specifies how floating point numbers are displayed. • If this item is not specified, decimal numbers between -1 and 0, and between 0 and 1, do not display the leading zero before the decimal point. For example, 0.25 is displayed as .25.	ORA TD
	• If this item is specified, decimal numbers between -1 and 0, and between 0 and 1, display the leading zero before the decimal point. For example, 0.25 is displayed as 0.25 .	

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
end_month_c alculate	Specifies the calculation logic of the add_months function. Assuming that the two parameters of the add_months function are param1 and param2, and the sum of the months of param1 and param2 is result: If this item is not specified, and the Day of param1 indicates the last day of a month shorter than result, the Day in the calculation result will equal that in param1. For example: select add_months('2018-02-28',3) from dual; add_months 2018-05-28 00:00:00 If this item is specified, and the Day of param1 indicates the last day of a month shorter than result, the Day in the calculation result will equal that in result. For example: select add_months('2018-02-28',3) from dual; add_months 2018-05-31 00:00:00 (1 row)	ORA TD
compat_anal yze_sample	Specifies the sampling behavior of the ANALYZE operation. If this item is specified, the sample collected by the ANALYZE operation will be limited to around 30,000 records, controlling CN memory consumption and maintaining the stability of ANALYZE.	ORA TD MyS QL
bind_schema _tablespace	Binds a schema with the tablespace with the same name. If a tablespace name is the same as <i>schema_name</i> , default_tablespace will also be set to <i>schema_name</i> if search_path is set to <i>schema_name</i> .	ORA TD MyS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
bind_procedu re_searchpat h	Specifies the search path of the database object for which no schema name is specified. If no schema name is specified for a stored procedure, the search is performed in the schema the stored procedure belongs to. If the stored procedure is not found, the following operations are performed: If this item is not specified, the system reports an error and exits. If this item is specified, the search continues based on the settings of search_path. If the issue persists, the system reports an error and exits.	ORA TD MyS QL
correct_to_nu mber	Controls the compatibility of the to_number() result. If this item is specified, the result of the to_number() function is the same as that of PG11. Otherwise, the result is the same as that of Oracle.	ORA
unbind_divid e_bound	Controls the range check on the result of integer division. If this item is not specified, the division result is checked. If the result is out of the range, an error is reported. In the following example, an out-of-range error is reported because the value of INT_MIN/(-1) is greater than the value of INT_MAX. SELECT (-2147483648)::int / (-1)::int; ERROR: integer out of range If this item is specified, the range of the division result does not need to be checked. In the following example, INT_MIN/(-1) can be used to obtain the output result INT_MAX+1. SELECT (-2147483648)::int / (-1)::int; ?column?	ORA TD
merge_updat e_multi	Specifies whether to perform an update when MERGE INTO is executed to match multiple rows. If this item is specified, no error is reported when multiple rows are matched. Otherwise, an error is reported (same as Oracle).	ORA TD

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
disable_row_ update_multi	Specifies whether to perform an update when multiple rows of a row-store table are matched. If this item is specified, an error is reported when multiple rows are matched. Otherwise, multiple rows can be matched and updated by default.	ORA TD
return_null_s tring	Specifies how to display the empty result (empty string ") of the lpad(), rpad(), repeat(), regexp_split_to_table(), and split_part() functions. • If this item is not specified, the empty string is displayed as NULL. select length(lpad('123',0,'*')) from dual; length (1 row) • If this item is specified, the empty string is displayed as single quotation marks ("). select length(lpad('123',0,'*')) from dual; length 0 (1 row)	ORA
compat_conc at_variadic	Specifies the compatibility of variadic results of the concat() and concat_ws() functions. If this item is specified and a concat function has a parameter of the variadic type, different result formats in Oracle and Teradata are retained. If this item is not specified and a concat function has a parameter of the variadic type, the result format of Oracle is retained for both Oracle and Teradata.	ORA TD

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
convert_strin g_digit_to_nu meric	 Specifies the type casting priority for binary BOOL operations on the CHAR type and INT type. If this item is not specified, the type casting priority is the same as that of PG9.6. After this item is configured, all binary BOOL operations of the CHAR type and INT type are forcibly. 	ORA TD MyS QL
	 After this item is configured, all binary BOOL operations of the CHAR type and INT type are forcibly converted to the NUMERIC type for computation. After this configuration item is set, the CHAR types that are affected include BPCHAR, VARCHAR, NVARCHAR2, and TEXT, and the INT types that are 	
	affected include INT1, INT2, INT4, and INT8. CAUTION This configuration item is valid only for binary BOOL operation, for example, INT2>TEXT and INT4=BPCHAR. Non-BOOL operation is not affected. This configuration item does not support conversion of UNKNOWN operations such as INT>'1.1'. After this configuration item is enabled, all BOOL operations of the CHAR and INT types are preferentially converted to the NUMERIC type for computation, which affects the computation performance of the database. When the JOIN column is a combination of affected types, the execution plan is affected.	

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
check_functio n_conflicts	Controls the check of the custom plpgsql/SQL function attributes. If this parameter is not specified, the IMMUTABLE/STABLE/VOLATILE attributes of a custom function are not checked. If this parameter is specified, the IMMUTABLE attribute of a custom function is checked. If the function contains a table or the STABLE/VOLATILE function, an error is reported during the function execution. In a custom function, a table or the STABLE/VOLATILE function conflicts with the IMMUTABLE attribute, thus function behaviors are not IMMUTABLE in this case. For example, when this parameter is specified, an error is reported in the following scenarios: CREATE OR replace FUNCTION sql_immutable (INTEGER) RETURNS INTEGER AS 'SELECT a+\$1 from shipping_schema.t4 where a=1;' LANGUAGE SQL IMMUTABLE RETURNS NULL ON NULL INPUT; select sql_immutable(1); ERROR: IMMUTABLE function cannot contain SQL statements with relation or Non-IMMUTABLE function. CONTEXT: SQL function "sql_immutable" during startup referenced column: sql_immutable	ORA TD MyS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
varray_verific ation	Indicates whether to verify the array length and array type length. This parameter is compatible with GaussDB(DWS) of versions earlier than 8.1.0. If this parameter is specified, the array length and array type length are not verified. Scenario 1 CREATE OR REPLACE PROCEDURE varray_verification AS TYPE org_varray_type IS varray(5) OF VARCHAR2(2); v_org_varray org_varray_type; BEGIN V_org_varray(1) := '111';If the value exceeds the limit of VARCHAR2(2), the setting will be consistent with that in the historical version and no verification is performed after configuring this option. END; / Scenario 2 CREATE OR REPLACE PROCEDURE varray_verification_i3_1 AS TYPE org_varray_type IS varray(2) OF NUMBER(2); v_org_varray org_varray_type; BEGIN V_org_varray(3) := 1;If the value exceeds the limit of varray(2) specified for array length, the setting will be consistent with that in the historical version and no verification is performed after configuring this option. END; //	ORA TD

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
strict_concat_ functions	Indicates whether the textanycat() and anytextcat() functions are compatible with the return value if there are null parameters. This parameter and strict_text_concat_td are mutually exclusive. In MySQL-compatible mode, this parameter has no impact. If this configuration item is not specified, the returned values of the textanycat() and anytextcat() functions are the same as those in the Oracle database. When this configuration item is specified, if there are null parameters in the textanycat() and anytextcat() functions, the returned value is also null. Different result formats in Oracle and Teradata are retained. If this configuration item is not specified, the returned values of the textanycat() and anytextcat() functions are the same as those in the Oracle database. SELECT textanycat('gauss', cast(NULL as BOOLEAN)); textanycat	ORA TD

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
strict_text_co ncat_td	In Teradata compatible mode, whether the textcat(), textanycat() and anytextcat() functions are compatible with the return value if there are null parameters. This parameter and strict_concat_functions are mutually exclusive. • If this parameter is not specified, the return values of the textcat(), textanycat(), and anytextcat() functions in Teradata-compatible mode are the same as those in GaussDB(DWS). • When this parameter is specified, if the textcat(), textanycat(), and anytextcat() functions contain any null parameter values, the return value is null in Teradata-compatible mode. If this parameter is not specified, the return values of the textcat(), textanycat(), and anytextcat() functions are the same as those in GaussDB(DWS). td_compatibility_db=# SELECT textcat('abc', NULL); textcat	TD

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
compat_displ ay_ref_table	 Sets the column display format in the view. If this parameter is not specified, the prefix is used by default, in the tab.col format. Specify this parameter to the same original definition. It is displayed only when the original definition contains a prefix. SET behavior_compat_options='compat_display_ref_table'; CREATE OR REPLACE VIEW viewtest2 AS SELECT a.c1, c2, a.c3, 0 AS c4 FROM viewtest_tbl a; SELECT pg_get_viewdef('viewtest2'); pg_get_viewdef SELECT a.c1, c2, a.c3, 0 AS c4 FROM viewtest_tbl a; (1 row) 	ORA TD
para_support _set_func	Whether the input parameters of the COALESCE(), NVL(), GREATEST(), and LEAST() functions in a column-store table support multiple result set expressions. If this item is not specified and the input parameter contains multiple result set expressions, an error is reported, indicating that the function is not supported. SELECT COALESCE(regexp_split_to_table(c3,'#'), regexp_split_to_table(c3,'#')) FROM regexp_ext2_tb1 ORDER BY 1 LIMIT 5; ERROR: set-valued function called in context that cannot accept a set When this configuration item is specified, the function input parameter can contain multiple result set expressions. SELECT COALESCE(regexp_split_to_table(c3,'#'), regexp_split_to_table(c3,'#')) FROM regexp_ext2_tb1 ORDER BY 1 LIMIT 5; coalesce	ORA TD

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
disable_selec t_truncate_p arallel	 Controls the DDL lock level such as TRUNCATE in a partitioned table. If this item is specified, the concurrent execution of TRUNCATE and DML operations (such as SELECT) on different partitions is forbidden, and the fast query shipping (FQS) of the SELECT operation on the partitioned table is allowed. You can set this parameter in the OLTP database, where there are many simple queries on partitioned tables, and there is no requirement for concurrent TRUNCATE and DML operations on different partitions. If this item is not specified, SELECT and TRUNCATE operations can be concurrently performed on different partitions in a partitioned table, and the FQS of the partitioned table is disabled to avoid possible inconsistency. 	ORA TD MyS QL
bpchar_text_ without_rtri m	In Teradata-compatible mode, controls the space to be retained on the right during the character conversion from bpchar to text . If the actual length is less than the length specified by bpchar , spaces are added to the value to be compatible with the Teradata style of the bpchar string. Currently, ignoring spaces at the end of a string for comparison is not supported. If the concatenated string contains spaces at the end, the comparison is space-sensitive. The following is an example: td_compatibility_db=# select length('a'::char(10)::text); length 10 (1 row) td_compatibility_db=# select length('a' 'a'::char(10)); length 11 (1 row)	TD

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
convert_empt y_str_to_null_ td	In Teradata-compatible mode, controls the to_date, to_timestamp, and to_number type conversion functions to return null when they encounter empty strings, and controls the format of the return value when the to_char function encounters an input parameter of the date type. Example: If this parameter is not specified: td_compatibility_db=# select to_number("); to_number	TD

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
	td_compatibility_db=# select to_char(date '2020-11-16'); to_char 2020/11/16 (1 row)	
disable_case_ specific	 Determines whether to ignore case sensitivity during character type match. This parameter is valid only in Teradata-compatible mode. If this item is not specified, characters are casesensitive during character type match. If this item is specified, characters are case-insensitive during character type match. After being specified, this item will affect five character types (CHAR, TEXT, BPCHAR, VARCHAR, and NVARCHAR), 12 operators (<, >, =, >=, <=, !=, <>, !=, like, not like, in, and not in), and expressions case when and decode. CAUTION After this item is enabled, the UPPER function is added before the character type, which affects the estimation logic. Therefore, an enhanced estimation model is required. (Suggested settings: cost_param = 16, cost_model_version = 1, join_num_distinct = -20, and qual_num_distinct = 200) 	TD
enable_interv al_to_text	Controls the implicit conversion from the interval type to the text type. • When this option is enabled, the implicit conversion from the interval type to the text type is supported. SELECT TO_DATE('20200923', 'yyyymmdd') - TO_DATE('20200920', 'yyyymmdd') = '3'::text; ?column?	ORA TD MyS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
case_insensiti ve	In MySQL-compatible mode, configure this parameter to specify the case-insensitive input parameters of the locate, strpos, and instr string functions. Currently, this parameter is not configured by default. That is, the input parameter is case-sensitive. Example: If this parameter is not configured, the input parameter is case-sensitive. mysql_compatibility_db=# SELECT LOCATE('sub', 'Substr'); locate 0 (1 row) If this parameter is configured, the input parameter is case-insensitive. mysql_compatibility_db=# SELECT LOCATE('sub', 'Substr'); locate 1 (1 row)	MyS QL
inherit_not_n ull_strict_fun c	Controls the original strict attribute of a function. A function with one parameter can transfer the NOT NULL attribute. func(x) is used an example. If func() is the strict attribute and x contains the NOT NULL constraint, func(x) also contains the NOT NULL constraint. The compatible configuration item is effective in some optimization scenarios, for example, NOT IN and COUNT(DISTINCT) optimization. However, the optimization results may be incorrect in specific scenarios. Currently, this parameter is not configured by default to ensure that the result is correct. However, the performance may be rolled back. If an error occurs, you can set this parameter to roll back to the historical version.	ORA TD MyS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
disable_comp at_minmax_e xpr_mysql	Specifies the method for processing the input parameter null in the greatest/least expression in MySQL-compatible mode. You can configure this parameter to roll back to a historical version. If this parameter is not configured and the input parameter is null, null is returned. mysql_compatibility_db=# SELECT greatest(1, 2, null), least(1, 2, null); greatest least	e MyS QL
	greatest least 	

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
disable_comp at_substr_my sql	Specifies the behavior of the substr/substring function when the start position pos is ≤ 0 in MySQL-compatible mode. You can configure this parameter to roll back to a historical version.	MyS QL
	 If this parameter is not configured, that is, an empty string is returned when pos = 0. When pos < 0, TRUNCATE starts from the last pos character on. mysql_compatibility_db=# SELECT substr('helloworld',0); substr	
	helloworld (1 row) mysql_compatibility_db=# SELECT substring('helloworld',0),substring('helloworld',-2,4); substring substring+ helloworld h (1 row)	

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
disable_comp at_trim_mysq l	Specifies the method for processing the input parameter in the trim/ltrim/rtrim function in MySQL-compatible mode. You can configure this parameter to roll back to a historical version. If this parameter is not configured, the entire substring is matched. mysql_compatibility_db=# SELECT trim('{}{name}{}',','\{}'),trim('xyznamezyx','xyz'); btrim btrim	MyS QL
light_object_ mtime	Specifies whether the mtime column in the pg_object system catalog records object operations. If this parameter is configured, the GRANT, REVOKE, and TRUNCATE operations are not recorded by mtime, that is, the mtime column is not updated. If this parameter is not configured (by default), the ALTER, COMMENT, GRANT, REVOKE, and TRUNCATE operations are recorded by mtime, that is, the mtime column is updated.	ORA TD MyS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
disable_inclu ding_all_mys ql	In MySQL-compatible mode, this parameter controls whether the CREATE TABLELIKE syntax is INCLUDING_ALL.	MyS QL
	By default, this parameter is not set. That is, in MySQL compatibility mode, CREATE TABLE LIKE syntax is INCLUDING_ALL.	
	You can configure this parameter to roll back to a historical version.	
	If this parameter is not set, in MySQL-compatible mode, the CREATE TABLE LIKE syntax is in INCLUDING_ALL. mysqLcompatibility_db=# CREATE TABLE mysqLlike(id int, name varchar(10), score int) DISTRIBUTE BY hash(id) COMMENT 'mysqLlike'; CREATE TABLE mysqLcompatibility_db=# CREATE INDEX index_like ON mysqLlike(name); CREATE INDEX mysqLcompatibility_db=# \d+ mysqLlike;	

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
	If this parameter is set, in MySQL-compatible mode, the CREATE TABLE LIKE syntax is empty. mysql_compatibility_db=# SET behavior_compat_options = 'disable_including_all_mysql'; SET mysql_compatibility_db=# CREATE TABLE mysql_copy LIKE mysql_like; NOTICE: The 'DISTRIBUTE BY' clause is not specified. Using roundrobin as the distribution mode by default. HINT: Please use 'DISTRIBUTE BY' clause to specify suitable data distribution column. CREATE TABLE mysql_db=# \d+ mysql_copy; Table "public.mysql_copy" Column Type Modifiers Storage Stats target Description	
cte_onetime_ inline	 Indicates whether to execute inline for non-stream plans. When this parameter is set, the CTE that is not in a stream plan and is referenced only once executes inline. If this parameter is not set, the CTE that is not in a stream plan and is referenced only once does not execute inline. 	ORA TD MyS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
skip_first_aft er_mysql	Controls whether to ignore the FIRST/AFTER colname syntax in ALTER TABLE ADD/MODIFY/CHANGE COLUMN in MySQL-compatible mode. If this parameter is set, the FIRST/AFTER colname syntax is ignored, and executing this syntax will not result in any errors. mysql_compatibility_db=# SET behavior_compat_options = 'skip_first_after_mysql'; mysql_compatibility_db=# ALTER TABLE t1 ADD COLUMN b text after a; ALTER TABLE If this parameter is not set, the FIRST/AFTER colname syntax is not supported, and executing this syntax causes errors. mysql_compatibility_db=# SET behavior_compat_options = "; mysql_compatibility_db=# ALTER TABLE t1 ADD COLUMN b text after a; ERROR: FIRST/AFTER is not yet supported.	MyS QL
enable_divisi on_by_zero_ mysql	Specifies whether division or modulo operations will result in an error when the divisor is 0 in MySQL-compatible mode. (This configuration item is supported only by clusters of version 8.1.3.110 or later.) • If this parameter is set, NULL is returned if the divisor is 0 in a division or modulo operation. compatible_mysql_db=# SET behavior_compat_options = 'enable_division_by_zero_mysql'; SET compatible_mysql_db=# SELECT 1/0 AS test; test	MyS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
normal_sessi on_id	Indicates whether to generate a session ID in normal format. If this option is set, a session ID in normal format will be generated, which is compatible with session IDs in clusters of version 8.1.3 or earlier. SET behavior_compat_options='normal_session_id'; SELECT pg_current_sessionid(); pg_current_sessionid 1660268184.140594655524608 If this parameter is not set, a session ID in pretty format will be generated. SET behavior_compat_options="; SELECT pg_current_sessionid(); pg_current_sessionid 1660268184.140594655524608.coordinator1 (1 row)	ORA TD MyS QL
disable_jsonb _exact_match	 Specifies whether to check the jsonb type during fuzzy match for binary operators. If this parameter is specified, operators search for matched items within the entire search scope (including the jsonb type) during fuzzy match. This setting is compatible with the match rules of cluster versions later than 8.1.1. SET behavior_compat_options='disable_jsonb_exact_match'; select '2022' - '2'::text; ERROR: cannot delete from scalar If this parameter is not specified, fuzzy match is performed within the search scope, except for the jsonb type. This setting is compatible with the match rules of clusters of version earlier than 8.1.1. SET behavior_compat_options="; select '2022' - '2'::text; ?column?	ORA TD MyS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
merge_into_ with_trigger	 Controls whether the MERGE INTO operation can be performed on tables with triggers. When this option is set, the MERGE INTO operation can be performed on tables with triggers. When the MERGE INTO operation is performed, the trigger on the table is not activated. If this option is not set, an error is reported when the MERGE INTO operation is performed on a table with triggers. 	ORA TD MyS QL
add_column_ default_v_fun c	 Controls whether expression in alter table add column default expression supports volatile functions. If this option is selected, expression in alter table add column default expression supports volatile functions. If this option is not selected, expression in alter table add column default expression does not support volatile functions. If expression contains volatile functions, an error will be reported during statement execution. 	ORA TD MyS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
disable_full_g roup_by_mys ql	Specifies whether to display non-aggregated function query columns after GROUP BY in a query. If this option is specified, the query does not display any non-aggregated function query columns after GROUP BY. SET behavior_compat_options='disable_full_group_by_mysql'; SELECT a,b FROM t1 GROUP BY a; a b + 1 1 2 2 (2 rows) If this option is not specified, the query must display all non-aggregated function query columns after GROUP BY, or an error will be reported. SET behavior_compat_options="; SELECT a,b FROM t1 GROUP BY a; ERROR: column "t1.b" must appear in the GROUP BY clause or be used in an aggregate function LINE 1: SELECT a,b FROM t1 GROUP BY a; CAUTION This parameter must be used together with full_group_by_mode. After configuring this option, if full_group_by_mode is set to notpadding, non-aggregated query columns that are not part of the GROUP BY clause must have consistent data after grouping. Otherwise, the values in that column will be random.	MyS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
disable_gc_fd w_filter_parti al_pushdown	Controls whether filter criteria are pushed down when querying data from a foreign table (of type gc_fdw) in a collaborative analysis scenario. • When this option is specified, if there are factors in the filter criteria that do not meet the pushdown conditions (such as non-immutable functions), all filter criteria will not be pushed down to ensure the consistency of the result set. This behavior is compatible with clusters of version earlier than 8.2.1. Create a table in the source cluster. CREATE TABLE t1(c1 INT, c2 INT, c3 INT) DISTRIBUTE BY HASH(c1); Create a foreign table with the same structure in the local cluster. CREATE SERVER server_remote FOREIGN DATA WRAPPER qc_fdw options(ADDRESS 'address', DBNAME 'dbname', USERNAME 'username', PASSWORD' password'); CREATE FOREIGN TABLE t1(c1 INT, c2 INT, c3 INT) SERVER server_remote; Enable the parameter and see the pushdown behavior. SET behavior_compat_options = 'disable_gc_fdw_filter_partial_pushdown'; EXPLAIN (verbose on,costs off) SELECt * FROM t1 WHERE c1>3 AND c2 < 100 AND now() - '20230101' < c3; QUERY PLAN	ORA TD MyS QL
	Streaming (type: GATHER)	

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
	Output: c1, c2, c3 Node/s: All datanodes -> Foreign Scan on ca_schema.t1 Output: c1, c2, c3 Filter: ((now() - '2023-01-01 00:00:00-08'::timestamp with time zone) < (t1.c3)::interval) Remote SQL: SELECT c1, c2, c3 FROM ca_schema.t1 WHERE ((c1 > 3)) AND ((c2 < 100)) (7 rows)	
ignore_unshi pped_concurr ent_update	Determines whether to ignore new tuples when the UPDATE or DELETE statement is executed in the current session if the statement is not pushed down and the tuples are updated by other sessions. By default, new tuples are not processed.	ORA TD MyS QL
	• If this parameter is specified, new tuples are ignored when the UPDATE or DELETE statement is executed in the current session. If the UPDATE or DELETE statement is successfully executed, data inconsistency occurs in concurrent update scenarios. This behavior is compatible with the behavior in versions earlier than 8.2.1.	
	If this parameter is not set and the UPDATE or DELETE statement executed in the current session detects that tuples have been updated, the UPDATE or DELETE statement of the current session will be reexecuted to ensure data consistency. The number of statement execution retries is controlled by the max_query_retry_times parameter.	
disable_set_g lobal_var_on _datanode	Controls whether the set_config function can be used to set global variables on DNs.	ORA TD
	When this parameter is set, the set_config function cannot be used to set global variables on DNs. By default, this behavior is compatible with the behavior in versions earlier than 8.2.1.	MyS QL
	If this parameter is not set, the set_config function can set global variables on DNs. As a result, the global variable values on CNs and DNs are inconsistent, and errors may occur when the read_global_var function is pushed down.	

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
variadic_null_ check	Whether variadic can transfer the NULL parameter. This function is disabled by default. This parameter is supported only by clusters of version 8.3.0 or later. • When this parameter is set, passing NULL parameters to variadic is not allowed and will result in an error. SET behavior_compat_options = 'variadic_null_check'; SELECT format ('array', VARIADIC NULL); ERROR: VARIADIC parameter must be an array NOTE To be compatible with MySQL, enabling compat_concat_variadic does not take effect for the concat and concat_ws functions, and the NULL parameter can still be passed in. • If this parameter is not set, NULL parameters can be passed to variadic. SET behavior_compat_options = "; SELECT format ('array', VARIADIC NULL); format	ORA TD MyS QL
enable_use_s yscol_in_repli cate_table	Specifies whether oid, ctid, tableoid, or xc_node_id can be used as filter, join, and having conditions during INSERT, UPDATE, MERGE INTO, and DELETE statements are executed on replication tables. This parameter is not set by default. If this parameter is not set and oid, ctid, tableoid, or xc_node_id is used as filter, join, or having conditions when the INSERT, UPDATE, MERGE INTO, or DELETE statements are executed on replication tables, the following error is reported: ERROR: Can not use system column oid/ctid/tableoid/xc_node_id in Replication Table. When this parameter is set, the INSERT, UPDATE, MERGE INTO, and DELETE statements can be executed on replication tables using the system column id, ctid, tableoid, or xc_node_id. CAUTION If oid, ctid, tableoid, or xc_node_id is used as filter, join, and having conditions when the INSERT, UPDATE, MERGE INTO, or DELETE statements are executed on partition tables, the statement may result in cluster core dumps. In this case, exercise caution when setting this parameter.	ORA TD MYS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
enable_force _add_batch	Determines whether GaussDB(DWS) receives U packets in addbatch mode when support_batch_bind is set to on and enable_fast_query_shipping and enable_light_proxy are both set to off. This parameter is not set by default. If this parameter is not set, support_batch_bind is set to on, and enable_fast_query_shipping and enable_light_proxy are both set to off, GaussDB(DWS) does not receive U packets in addbatch mode. If this parameter is set, support_batch_bind is set to on, and enable_fast_query_shipping and enable_light_proxy are both set to off, GaussDB(DWS) receives U packets in addbatch mode. However, packets are imported to the database slowly, which may cause insufficient memory. So, exercise caution when setting this parameter.	ORA TD MYS QL
disable_merg esort_withou t_material	 Controls whether the current stream segment contains materialized operators. If it is, merge sort is used. If this parameter is set and the current stream segment contains materialized operators (material, sort, agg, and CteScan), merge sort can be used. Otherwise, merge sort cannot be used. If this parameter is unset, there is no need to verify whether the current stream segment contains materialized operators to determine whether to use merge sort. 	ORA TD MYS QL
enable_push down_groupi ngset_subqu ery	 Specifies whether conditions from the outer query that are only related to a subquery can be pushed down to the subquery when the subquery contains a grouping set. If the subquery contains grouping sets and this parameter is set, the conditions in the outer query cannot be pushed down to the subquery. If the subquery contains grouping sets and this parameter is not set, the conditions in the outer query can be pushed down to the subquery. 	ORA TD MYS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
enable_whol e_row_var	This parameter mainly involves two scenarios: 1. controlling whether tables or views are allowed to appear in SQL expressions, including but not limited to the target list of queries, GROUP BY lists, etc.; 2. controlling whether non-table records are allowed to appear in SQL expressions. This parameter is supported only by clusters of version 8.3.0 or later. • When this parameter is set, tables or views are allowed to appear in SQL expressions. SET behavior_compat_options = 'enable_whole_row_var'; SELECT a1 FROM t a1; a1 (0 rows) SELECT t FROM (SELECT 1) as t; t (1) (1 rows) • If this parameter is unset, tables or views are not allowed to appear in SQL expressions. SET behavior_compat_options = "; SELECT a1 FROM t a1; ERROR: Table or view cannot appear in expression. Table/view name: t, alias: a1. Please check targetList, groupClause etc. SELECT t FROM (SELECT 1) as t; ERROR: Non-table records cannot appear in expression. alias: t. Please check targetList, groupClause etc.	ORA TD MYS QL
enable_unkn own_datatyp e	 Specifies whether tables containing unknown columns can be created. This parameter is supported only by clusters of version 8.3.0 or later. When this parameter is set, tables containing unknown columns can be created. SET behavior_compat_options = 'enable_unknown_datatype'; CREATE TABLE t(a unknown); WARNING: column "a" has type "unknown" DETAIL: Proceeding with relation creation anyway. CREATE TABLE If this parameter is unset, tables containing unknown columns cannot be created. If the table creation SQL contains an unknown column, an error will be reported. SET behavior_compat_options = "; create table t(a unknown); ERROR: column "a" has type "unknown" 	ORA TD MYS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
alter_distribu te_key_by_pa rtition	 Specifies whether INSERT INTO is executed by partition when ALTER TABLE is used to modify the distribution column of a partitioned table. If this parameter is set, INSERT INTO is executed by partition. The memory usage decreases but the performance deteriorates. If this parameter is unset, INSERT INTO is performed on the entire partitioned table. The performance is good but the memory usage is high. 	ORA TD MYS QL
disable_upda te_returning_ check	 Specifies whether to prevent multiple joins when a replication table is updated with the returning statement. This parameter is supported only by clusters of version 8.3.0 or later. If the parameter is not set, the following error is reported when updating a replication table with a returning statement and involving multiple joins: ERROR: Unsupported FOR UPDATE replicated table joined with other table. Setting this parameter ensures backward compatibility with earlier versions. However, when updating a replication table with a returning statement and involving multiple joins, there may be inconsistencies in the result set. 	ORA TD MYS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
check_functio n_shippable	Controls the check of the custom plpgsql/SQL function attributes. This parameter is supported only by clusters of version 8.3.0 or later. • If this parameter is not specified, the IMMUTABLE/STABLE/VOLATILE attributes of a user-defined	ORA TD MYS QL
	function are not checked. If this parameter is specified, the IMMUTABLE/ STABLE/VOLATILE attributes of user-defined functions are checked based on the following principles: Whitelist: For the three functions in DBMS_OUTPUT, skip the check_function_shippable check. If a user-defined function contains DML statements and the outer layer is IMMUTABLE or SHIPPABLE, the function is pushed down. As a result, an error is reported. If the outer layer of a user-defined function is SHIPPABLE and the inner layer is IMMUTABLE, the function passes the check. If the outer layer of a user-defined function is SHIPPABLE, the function passes the check. If the outer layer of a user-defined function is SHIPPABLE, the function passes the check. If the outer layer of a user-defined function is SHIPPABLE but the inner layer is none of the above, an error is reported. For example, when this parameter is specified, an error is reported in the following scenarios: CREATE OR replace function func_ship(a int) returns int LANGUAGE plpgsql NOT FENCED SHIPPABLE AS \$function\$ begin perform test_ship(); return a; end \$function\$; select func_ship(a) from tt3; ERCCEPTION WHEN OTHERS THEN return a; end \$functionship(a) from tt3; ERROR: parent function is shippable but child is not immutable or	

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
enable_full_s tring_agg	Specifies how string_agg(a, delimiter) over (partition by b order by c) behaves in different situations, such as using full or incremental aggregation in the window. This parameter is supported only by clusters of version 8.3.0 or later. If this parameter is not set, incremental aggregation is used. If this parameter is set, full aggregation is used. By default, this parameter is not set. CREATE TABLE string_agg_dn_col(c1 int, c2 text) WITH(orientation = column) distribute by hash(c1); INSERT INTO string_agg_dn_col values(1, 'hadian'); INSERT INTO string_agg_dn_col values(1, 'hadian'); INSERT INTO string_agg_dn_col values(1, 'nanjing'); SELECT t.c1 AS c1, string_agg(t.c2, ',') OVER(PARTITION BY t.c1 ORDER BY t.c2) AS c2 FROM string_agg_dn_col t ORDER BY c2; c1 c2	ORA TD MYS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
enable_bank er_round	Specifies how numeric types round their values, using the rounding or banker method. This parameter is supported only by clusters of version 8.3.0 or later. Behaviors controlled by parameters include: • Type conversion working when INSERT INTO and ::xxx specify a type, such as integer types (int1, int2, int4, int8), any precision types (decimal, numeric, number), or money types. • Rounding and conversion functions for the numeric type: round(xxx.xx,s), cast('xxx.xx',numeric), or to_char(xxx.xx,'xxx'). • Mathematical calculation of the numeric type. NOTE The banker's rounding rule is as follows: if the digit to be rounded is greater than 5, round up; if it is less than 5, round down; if it is exactly 5, round to the nearest even number. • If this parameter is set, rounding uses the banker method. SET behavior_compat_options = enable_banker_round; SELECT 1.5::int1,1.5::int2,1.5::int4,1.5::int8,1.5::numeric(10,0),1.115::money; int1 int2 int4 int8 numeric money	ORA TD MYS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
	round round numeric numeric to_char to_char +	
create_partiti on_local_inde x	Controls whether the default index created on a partitioned table is a global index or a local index. This parameter is supported only by clusters of version 8.2.1.210 or later. • If this parameter is disabled, a global index is created on a partitioned table by default. SET behavior_compat_options = "; CREATE INDEX sale_id_idx ON sales(sale_id); ERROR: partitioned table does not support global index HINT: please set behavior_compat_options = 'create_partition_local_index'to create local index by default. • If this parameter is enabled, a local index is created on a partitioned table by default. SET behavior_compat_options = create_partition_local_index; CREATE INDEX sale_id_idx ON sales(sale_id); CREATE INDEX	ORA
enable_int_di vision_by_tru ncate	Controls whether the integer division behavior result set returns integers or floating point numbers and the option is compatible with PG or ORA behaviors. If this parameter is set, the integer division result is an integer, the decimal places are truncated, and this parameter is compatible with PG behaviors. SET behavior_compat_options = 'enable_int_division_by_truncate'; SELECT 8::int8 / 3::int8, 8::int4 / 3::int4, 8::int2 / 3::int2, 8::int1 / 3::int1; ?column? ?column? ?column? ?column? ?column? ?column? ?column? ?column? ?column? ?column? ?column ?column ?column ?column ?column ?column ?column? ?column	ORA TD MYS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
select_into_al low_multi_re sult	Controls whether SELECT INTO in a stored procedure can receive multiple rows of results or no result set at all. If this parameter is unset, the following error is reported when the SELECT INTO statement is used to insert multiple values or empty values in a stored procedure: CREATE TABLE PersonTmpTable (id int primary key, name varchar(64), age int, city varchar(512) default null, update_time timestamp default null); INSERT INTO PersonTmpTable VALUES(1,'zhangsan', 23, 'wuhan', '2022-12-10 15:39:32'); INSERT INTO PersonTmpTable VALUE(2,'lisi', 11, 'beijing', '2022-12-13 15:39:32'); INSERT INTO PersonTmpTable VALUE(3,'wangwu', 46, 'xian', '2022-12-1 15:39:32'); INSERT INTO PersonTmpTable VALUE(4,'zhaoliu', 46, 'wuhan', '2022-12-31 15:39:32'); SET behavior_compat_options = "; CREATE OR REPLACE function func_test1 RETURNS int LANGUAGE plpgsql AS \$\$DECLAREtmp INTEGER;BEGINSELECT id INTO tmp FROM PersonTmpTable;dbms_output.put_line(tmp);return tmp;END;\$\$;select func_test1(); ERROR: query returned no rows when process INTO. ERROR: query returned 2 rows more than one row. If this parameter is set, the first row or empty row is inserted when the SELECT INTO statement is used to insert multiple values or empty values in a stored procedure. SET behavior_compat_options = 'select_into_allow_multi_result'; CREATE OR REPLACE function func_test1 RETURNS int LANGUAGE plpgsql AS \$\$DECLAREtmp INTEGER;BEGINSELECT id INTO tmp FROM PersonTmpTable;dbms_output.put_line(tmp);return tmp;END;\$\$;select func_test1(); CREATE OR REPLACE function func_test1 RETURNS int LANGUAGE plpgsql AS \$\$DECLAREtmp INTEGER;BEGINSELECT id INTO tmp FROM PersonTmpTable;dbms_output.put_line(tmp);return tmp;END;\$\$;select func_test1(); CREATE FUNCTION	ORA TD MYS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
orderby_null_f irst	Controls whether NULL values are treated as the minimum value by default when sorting with ORDER BY . If this parameter is set, NULL values are treated as the minimum value by default when sorting with ORDER BY . SET behavior_compat_options = 'orderby_null_first'; SELECT * FROM test ORDER BY a; a b+ 1	D
unsupported _set_function _case	 Specifies whether multiple result set functions can be returned in a CASE WHEN condition. This is supported only by clusters of version 8.3.0.100 or later. This is enabled by default in newly installed clusters of version 9.1.0 or later. If this parameter is set, column storage does not support multiple result set functions in a CASE WHEN condition. CREATE TABLE t1(id int, c1 text) with(orientation=column); INSERT INTO t1 values(1, 'a#1'); SET behavior_compat_options = 'unsupported_set_function_case'; SELECT CASE split_part(regexp_split_to_table(c1, E''),'#',1) when 'a' then c1 else null end from t1; ERROR: set-valued function called in context that cannot accept a set If this parameter is not set, column storage supports multiple result set functions in a CASE WHEN condition. SET behavior_compat_options = "; SELECT CASE split_part(regexp_split_to_table(c1, E''),'#',1) when 'a' then c1 else null end from t1; case a#1 (1 row) 	ORA TD MYS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
enable_chan ge_search_pa th	Specifies whether the search path can be modified after forming a general plan <code>generic_plan</code> . This is supported only by clusters of version 9.1.0 or later. • When this parameter is not set, if a new search path is set and an <code>EXECUTE</code> statement is executed, the database will still search for the corresponding table under the original schema of the table. CREATE SCHEMA s1 CREATE SCHEMA s2 CREATE TABLE abc(f1 INT); SET search_path = \$1; INSERT INTO \$1.abc VALUES(123);INSERT INTO \$2.abc VALUES(456); SET search_path = \$1; PREPPARE p1 AS SELECT f1 FROM abc; EXECUTE p1; f1 123 (1 row) SET search_path = \$2; SELECT f1 FROM abc; f1 123 (1 row) • When this parameter is set, if a new search path is set and an <code>EXECUTE</code> statement is executed, the database will search for the corresponding table in the newly set search path. SET behavior_compat_options = 'enable_change_search_path'; EXECUTE p1; f1 456 (1 row) SET search_path = \$1; EXECUTE p1; f1 456 (1 row)	TD

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
enable_varch ar_to_nvarch ar2	 Specifies whether varchar fields created or updated through DDL statements are automatically switched to nvarchar2 fields. This is supported only by clusters of version 9.1.0 or later. If this parameter is set, varchar fields created or updated through DDL statements are automatically switched to nvarchar2 fields. If this parameter is unset, varchar fields created or updated through DDL statements are not automatically switched to nvarchar2 fields. 	ORA TD MYS QL
normalize_ne gative_zero	Specifies whether the ceil() and round() functions return -0 when processing specific values of the float type. • When this parameter is set, the ceil() function returns 0 when processing (-1,0), and the round() function returns 0 when processing [-0.5, 0). SET behavior_compat_options='normalize_negative_zero'; SELECT ceil(cast(-0.1 as float)); ceil 0 (1 row) • When this parameter is not set, the ceil() function returns -0 when processing (-1,0), and the round() function returns -0 when processing [-0.5, 0). SET behavior_compat_options = "; SELECT ceil(cast(-0.1 as FLOAT)); ceil0 (1 row) SELECT round(cast(-0.1 as FLOAT));	ORA TD MyS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
disable_client _detection_co mmit	Specifies whether to check there is a connection with the client before each transaction is committed. If the connection does not exist, an error is reported, the transaction is rolled back, and data duplication caused by repeated issuance due to disconnection is prevented. If this parameter is not set, the system checks the existence of the client connection before each transaction is committed. If this parameter is set, the system does not check the existence of the client connection before each transaction is committed.	ORA TD MyS QL
change_illeg al_char	Specifies the display of illegal UTF8 characters when reading with GDS. This parameter is supported only by clusters of version 8.3.0.100 or later. When this parameter is enabled, illegal UTF8 characters that are incompatible with GDS are displayed as "*" instead of "?".	MyS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
row_use_pse udo_name	Specifies whether row-related expressions generate pseudo column names for anonymous columns. This is supported only by 9.1.0.100 and later cluster versions. • When this parameter is not set, if there is a corresponding real column name in the row expression, the real column name is used. If it is an anonymous column, pseudo column names f1, f2fn are generated. SELECT row_to_json(row(1,'foo')); row_to_json	ORA TD MyS QL
enable_trunc _orc_string	Controls the foreign table query behavior when the foreign table field is in ORC format and the data type is varchar(n), but the field type in the ORC file is string and the length of the string exceeds n. If this parameter is not set, an error message is returned, indicating that the field is too long. If this parameter is set, the query is responded to, and the result is truncated by the length defined by varchar(n).	ORA TD MyS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
gds_fill_multi _missing_fiel ds	Controls the behavior when the GDS foreign table fault tolerance parameter fill_missing_fields is set to true or on. When fill_missing_fields is set to true or on in a GDS foreign table, any missing columns at the end of a row in the data source file are automatically set to NULL. Before this, only the last column in a row of the data source file can be missing without an error being reported. • If this option is specified, the GDS foreign table tolerates the missing of multiple last columns in a row of the source data file. • If this option is not specified, only the missing of the last column in a row of the data source file is tolerated in the GDS foreign table. This parameter compatible with historical behavior.	ORA TD MyS QL
correct_conv ert_tz	Controls whether the query result of the convert_tz function is consistent with the MySQL behavior. This is supported only by 9.1.0.210 and later cluster versions. • When this configuration item is enabled, the convert_tz function produces the expected results during time zone conversion, aligning with MySQL's behavior. set behavior_compat_options = 'correct_convert_tz'; select convert_tz'(2020-1-1 12:00', '+00:00', 'PRC'); convert_tz	ORA TD MyS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
set_timezone _mysql_style	Controls whether the specified timezone setting using the set timezone=±' Hour:Minute' command, aligns with MySQL's behavior. For example, set timezone='+08:00'/'-07:00'. This is supported only by 9.1.0.210 and later cluster versions. • When configuration item is enabled, set timezone='+08:00' indicates that the time zone is set to GMT+8, aligning with MySQL's behavior. set behavior_compat_options = 'set_timezone_mysqL_style'; set timezone='UTC'; select * from t; a	ORA TD MyS QL

Configuratio n Item	Behavior	Appl icabl e Com pati bilit y Mod e
parquet_int9 6_is_localtim e	Specifies whether to use the database time zone to correct the original timestamp when the INT96 data (dedicated for timestamps) in a foreign table in Parquet format is read. It is available only for clusters of version 9.1.0.210 or later. • If this configuration item is enabled, the parquet INT96 type stores the local time (not UTC+0), and the time zone is not corrected in the read results. set behavior_compat_options = 'parquet_int96_is_localtime'; set time zone 'Asia/Shanghai'; select * from test_parquet_int96; a	ORA TD MyS QL

enable_matview

Parameter description: Specifies whether to enable the materialized view function.

Type: USERSET

Value range: Boolean

- **on**: The materialized view function is enabled.
- off: The materialized view function is disabled.

Default value: off

hive compat options

Parameter description: This parameter specifies whether to convert an empty string to **NULL** when data is inserted into a numeric field using **INSERT INTO VALUES** or into a character field using **INSERT INTO SELECT** in MySQL mode of GaussDB(DWS). This parameter is available only for clusters of version 9.1.0.210 or later.

Type: USERSET

Value range: enumerated values

- empty_str_to_null: An empty string will be converted to NULL when data is inserted into a numeric field or a character field.
- empty: When data is inserted into a numeric field using INSERT INTO
 VALUES or into a character field using INSERT INTO SELECT, empty strings
 are not processed.

Default value: empty

15.17 Fault Tolerance

This section describes parameters used for controlling the methods that the server processes an error occurring in the database system.

exit on error

Parameter description: Specifies whether to terminate the current session.

Type: SUSET

Value range: Boolean

- on indicates that any error will terminate the current session.
- **off** indicates that only a FATAL error will terminate the current session.

Default value: off

omit_encoding_error

Parameter description: When performing character encoding conversion in the database, if a character encoding error occurs and the target character set encoding is UTF-8, the converted character with the error can be ignored and replaced with "?".

Type: USERSET

Value range: Boolean

- **on** indicates that characters that have conversion errors will be ignored and replaced with question marks (?), and error information will be recorded in logs.
- **off** indicates that characters that have conversion errors cannot be converted and error information will be directly displayed.

Default value: off

max_query_retry_times

Parameter description: Specifies the maximum number of retries for the automatic retry feature when a SQL statement encounters an error. Currently, the supported error types for retry include **Connection reset by peer**, **Lock wait timeout**, and **Connection timed out**. Setting this parameter to **0** will disable the retry feature.

Type: USERSET

Value range: an integer ranging from 0 to 20

Default value: 6

retry_ecode_list

Parameter description: Specifies the list of SQL error types that support

automatic retry.

Type: USERSET

Value range: a string

Default value: YY001 YY002 YY003 YY004 YY005 YY006 YY007 YY008 YY009 YY010 YY011 YY012 YY013 YY014 YY015 53200 08006 08000 57P01 XX003 XX009

YY016 CG003 CG004 F0011 F0012 45003 42P30

15.18 Connection Pool Parameters

When a connection pool is used to access the database, database connections are established and then stored in the memory as objects during system running. When you need to access the database, no new connection is established. Instead, an existing idle connection is selected from the connection pool. After you finish accessing the database, the database does not disable the connection but puts it back into the connection pool. The connection can be used for the next access request.

max_pool_size

Parameter description: Specifies the maximum number of connections between a CN's connection pool and another CN/DN.

Type: POSTMASTER

Value range: an integer ranging from 1 to 65535

Default value: 800 for CNs and 5000 for DNs

persistent_datanode_connections

Parameter description: Specifies whether to release the connection for the current session.

Type: USERSET

Value range: Boolean

- off indicates that the connection for the current session will be released.
- **on** indicates that the connection for the current session will not be released.

NOTICE

After this function is enabled, a session may hold a connection but does not run a query. As a result, other query requests fail to be connected. To fix this problem, the number of sessions must be less than or equal to **max active statements**.

Default value: off

cache connection

Parameter description: Specifies whether to reclaim the connections of a connection pool.

Type: USERSET

Value range: Boolean

- **on** indicates that the connections of a connection pool will be reclaimed.
- **off** indicates that the connections of a connection pool will not be reclaimed.

Default value: on

enable force reuse connections

Parameter description: Specifies whether a session forcibly reuses a new connection.

Type: USERSET

Value range: Boolean

- **on** indicates that the new connection is forcibly used.
- **off** indicates that the current connection is used.

Default value: off

syscache_clean_policy

Parameter description: Specifies the policy for clearing the memory and number of idle DN connections. This is supported only by clusters of version 9.1.0.100 or later.

Type: SIGHUP

Value range: a string

This parameter policy consists of three values:

- 1. The first value ranges from 0 to 1 and represents the percentage of total available memory used by DNs. When the percentage of used memory reaches this value, 1/4 of the stream threads will be cleared, and the second value will be evaluated.
- 2. The second value ranges from 0 to 1 and represents the percentage of total available memory used by syscache on DNs. When the percentage of syscache memory usage reaches this value, the third value will be evaluated.
- 3. The third value ranges from 0 to INT_MAX and is measured in MB. It represents the size of syscache memory used by idle threads. When the syscache memory usage of an idle thread reaches this value, the syscache used by that thread will be cleared.

Default value: 0.8,0.3,64

NOTICE

- Before setting this parameter, evaluate the memory usage using views PV_SESSION_MEMORY_DETAIL and PV_TOTAL_MEMORY_DETAIL.
- When setting this parameter, follow the specified format, ensuring that the three values are separated by commas without spaces.
- If the parameter is not set according to the specified format and the setting fails, a WARNING log will be generated in the log, and the parameter value displayed when using the SHOW command to query the parameter will be the last successfully set value. If the setting fails and the system is restarted, the parameter will be set to the default value.
- During the Readcommand phase, if a thread on CN times out after 30 seconds, it will clear DNs if syscache is greater than 256 MB. There are two operations:
 - 1. If the overall memory usage reaches 80%, an auxiliary thread will monitor the memory usage and clear 1/4 of the stream threads. It will also check if syscache usage exceeds 30% of the total memory usage. If it does, it will clear the syscache of Readcommand phase pg threads greater than 64 MB.
 - 2. If a stream thread is idle for more than 30 seconds and syscache usage is greater than 64 MB, it will clear the syscache.

15.19 Cluster Transaction Parameters

This section describes the settings and value ranges of cluster transaction parameters.

transaction_isolation

Parameter description: Specifies the isolation level of the current transaction.

Type: USERSET Value range:

- **READ COMMITTED**: Only committed data is read. This is the default.
- READ UNCOMMITTED: GaussDB(DWS) does not support READ UNCOMMITTED. If READ UNCOMMITTED is set, READ COMMITTED is executed instead.
- **REPEATABLE READ**: Only the data committed before transaction start is read. Uncommitted data or data committed in other concurrent transactions cannot be read.
- SERIALIZABLE: GaussDB(DWS) does not support SERIALIZABLE. If SERIALIZABLE is set, REPEATABLE READ is executed instead.

Default value: READ COMMITTED

transaction_read_only

Parameter description: Specifies that the current transaction is a read-only transaction.

Type: USERSET

Value range: Boolean

- **on** indicates that the current transaction is a read-only transaction.
- off indicates that the current transaction can be a read/write transaction.

Default value: off for CNs and on for DNs

xc_maintenance_mode

Parameter description: Specifies whether the system is in maintenance mode.

Type: SUSET

Value range: Boolean

- on indicates that maintenance mode is enabled.
- **off** indicates that the maintenance mode is disabled.

Default value: off

NOTICE

Enable the maintenance mode with caution to avoid cluster data inconsistencies.

allow_concurrent_tuple_update

Parameter description: Specifies whether to allow concurrent update.

Type: USERSET

Value range: Boolean

- on indicates it is enabled.
- off indicates it is disabled.

Default value: on

gtm_backup_barrier

Parameter description: Specifies whether to create a restoration point for the GTM starting point.

Type: SUSET

Value range: Boolean

- **on** indicates that a restoration point will be created for the GTM starting point.
- **off** indicates that a restoration point will not be created for the GTM starting point.

Default value: off

transaction_deferrable

Parameter description: Specifies whether to delay the execution of a read-only serial transaction without incurring an execution failure. Assume this parameter is set to **on**. When the server detects that the tuples read by a read-only transaction are being modified by other transactions, it delays the execution of the read-only transaction until the other transactions finish modifying the tuples. Currently, this parameter is not used in GaussDB(DWS). Similar to this parameter, the **default_transaction_deferrable** parameter is used to specify whether to allow delayed execution of a transaction.

Type: USERSET

Value range: Boolean

- **on** indicates that the execution of a read-only serial transaction can be delayed.
- **off** indicates that the execution of a read-only serial transaction cannot be delayed.

Default value: off

enforce_two_phase_commit

Parameter description: This parameter is reserved for compatibility with earlier versions. This parameter is invalid in the current version.

enable_show_any_tuples

Parameter description: This parameter is available only in a read-only transaction and is used for analysis. When this parameter is set to **on/true**, all versions of tuples in the table are displayed.

Type: USERSET

Value range: Boolean

- **on/true** indicates that all versions of tuples in the table are displayed.
- off/false indicates that no versions of tuples in the table are displayed.

Default value: off

idle_in_transaction_timeout

Parameter description: duration during which a transaction is allowed to be in the idle state. When a transaction is in the idle state for a period specified by this parameter, the transaction is terminated. This function takes effect only for client connections that are directly connected to CNs and does not take effect for direct DNs or internal connections. This parameter is supported only by clusters of version 8.2.1.100 or later.

Type: USERSET

Value range: 0 to 86400, in second.

Default value: 3600

15.20 Developer Operations

enable_light_colupdate

Parameter description: Specifies whether to enable the lightweight column-store update.

Type: USERSET

Value range: Boolean

- **on** indicates that the lightweight column-store update is enabled.
- **off** indicates that the lightweight column-store update is disabled.

Default value: off

□ NOTE

There is a low probability that an error is reported when lightweight **UPDATE** and backend column-store **AUTOVACUUM** coexist. You can run **ALTER TABLE** to set the table-level parameter **enable_column_autovacuum_garbage** to **off** to avoid this issue. If the table-level parameter **enable_column_autovacuum_garbage** is set to **off**, the backend column-store **AUTOVACUUM** of the table is disabled.

enable_fast_query_shipping

Parameter description: Specifies whether to use the distributed framework for a query optimizer.

Type: USERSET

Value range: Boolean

- **on** indicates that execution plans are generated on CNs and DNs separately.
- **off** indicates that the distributed framework is used. Execution plans are generated on CNs and then sent to DNs for execution.

Default value: on

enable_trigger_shipping

Parameter description: Specifies whether the trigger can be pushed to DNs for execution.

Type: USERSET

Value range: Boolean

- **on** indicates that the trigger can be pushed to DNs for execution.
- **off** indicates that the trigger cannot be pushed to DNs. It must be executed on the CN.

Default value: on

enable_remotejoin

Parameter description: Specifies whether JOIN operation plans can be delivered to DNs for execution.

Type: USERSET

Value range: Boolean

- **on** indicates that JOIN operation plans can be delivered to DNs for execution.
- off indicates that JOIN operation plans cannot be delivered to DNs for execution.

Default value: on

enable_remotegroup

Parameter description: Specifies whether the execution plans of **GROUP BY** and **AGGREGATE** can be delivered to DNs for execution.

Type: USERSET

Value range: Boolean

- **on** indicates that the execution plans of **GROUP BY** and **AGGREGATE** can be delivered to DNs for execution.
- **off** indicates that the execution plans of **GROUP BY** and **AGGREGATE** cannot be delivered to DNs for execution.

Default value: on

enable_remotelimit

Parameter description: Specifies whether the execution plan specified in the LIMIT clause can be pushed down to DNs for execution.

Type: USERSET

Value range: Boolean

• **on** indicates that the execution plan specified in the LIMIT clause can be pushed down to DNs for execution.

• **off** indicates that the execution plan specified in the LIMIT clause cannot be delivered to DNs for execution.

Default value: on

enable_limit_stop

Parameter description: whether to enable the **early stop** optimization for **LIMIT** statements. For a **LIMIT n** statement, if **early stop** is enabled, the CN requests the DN to end the execution after receiving n pieces of data. This method is applicable to complex queries with **LIMIT**. This parameter is supported only by clusters of version 8.1.3.320 or later.

Type: USERSET

Value range: Boolean

- on indicates that early stop is enabled for LIMIT statements.
- **off** indicates that **early stop** is disabled for LIMIT statements.

Default value: on

enable remotesort

Parameter description: Specifies whether the execution plan of the ORDER BY clause can be delivered to DNs for execution.

Type: USERSET

Value range: Boolean

- **on** indicates that the execution plan of the ORDER BY clause can be delivered to DNs for execution.
- **off** indicates that the execution plan of the ORDER BY clause cannot be delivered to DNs for execution.

Default value: on

enable_join_pseudoconst

Parameter description: Specifies whether joining with the pseudo constant is allowed. A pseudo constant indicates that the variables on both sides of a join are identical to the same constant.

Type: USERSET

Value range: Boolean

- **on** indicates that joining with the pseudo constant is allowed.
- **off** indicates that joining with the pseudo constant is not allowed.

Default value: off

cost model version

Parameter description: Specifies the model used for cost estimation in the application scenario. This parameter affects the distinct estimation of the

expression, HashJoin cost model, estimation of the number of rows, distribution key selection during redistribution, and estimation of the number of aggregate rows.

Type: USERSET

Value range: 0, 1, 2, 3, or 4

- **0** indicates that the original cost estimation model is used.
- 1 indicates that the enhanced distinct estimation of the expression, HashJoin cost estimation model, estimation of the number of rows, distribution key selection during redistribution, and estimation of the number of aggregate rows are used on the basis of **0**.
- 2 indicates that the ANALYZE sampling algorithm with better randomicity is used on the basis of 1 to improve the accuracy of statistics collection.
- 3 indicates that the broadcast cost estimation in large cluster scenarios is optimized based on 2 so that the optimizer can select a better plan. This option is supported only by clusters of version 8.3.0 or later.
- 4 indicates that in addition to the optimizations made to the cost estimation
 of hashjoin parallelization, skew, and column-store index ordering in 3, there
 are also optimized row estimations for coalesce expressions and improved
 recognition of skew optimization for subquery constant output columns
 during joins.

Default value: 4

debug_assertions

Parameter description: Specifies whether to enable various assertion checks. This parameter assists in debugging. If you are experiencing strange problems or crashes, set this parameter to **on** to identify programming defects. To use this parameter, the macro USE_ASSERT_CHECKING must be defined (through the configure option **--enable-cassert**) during the GaussDB(DWS) compilation.

Type: USERSET

Value range: Boolean

- **on** indicates that various assertion checks are enabled.
- **off** indicates that various assertion checks are disabled.

∩ NOTE

This parameter is set to **on** by default if GaussDB(DWS) is compiled with various assertion checks enabled.

Default value: off

distribute_test_param

Parameter description: Specifies whether the embedded test stubs for testing the distribution framework take effect. In most cases, developers embed some test stubs in the code during fault injection tests. Each test stub is identified by a unique name. The value of this parameter is a triplet that includes three values: thread level, test stub name, and error level of the injected fault. The three values are separated by commas (,).

Type: USERSET

Value range: a string indicating the name of any embedded test stub.

Default value: -1, default, default

ignore_checksum_failure

Parameter description: Sets whether to ignore check failures (but still generates an alarm) and continues reading data. This parameter is valid only if **enable_crc_check** is set to **on**. Continuing reading data may result in breakdown, damaged data being transferred or hidden, failure of data recovery from remote nodes, or other serious problems. You are not advised to modify the settings.

Type: SUSET

Value range: Boolean

- **on** indicates that data check errors are ignored.
- **off** indicates that data check errors are reported.

Default value: off

default table behavior

Parameter description: behavior type of the default table. This parameter is supported only by clusters of version 8.2.1 or later.

Type: USERSET

Value range: column_btree_index, column_high_compress, column middle compress, or column low compress

- **column_btree_index** indicates that the default index for creating a columnstore table is **B-Tree**.
- column_high_compress indicates that the default compression level of column-store tables is high.
- column_middle_compress indicates that the default compression level of column-store tables is middle.
- column_low_compress indicates that the default compression level of column-store tables is low.

Default value: an empty string

enable_colstore

Parameter description: Specifies whether to create a table as a column-store table by default when no storage method is specified. The value for each node must be the same. This parameter is used for tests. Users are not allowed to enable it.

Type: SUSET

Value range: Boolean

Default value: off

enable_force_vector_engine

Parameter description: Specifies whether to forcibly generate vectorized execution plans for a vectorized execution operator if the operator's child node is a non-vectorized operator. When this parameter is set to **on**, vectorized execution plans are forcibly generated. When **enable_force_vector_engine** is enabled, no matter it is a row-store table, column-store table, or hybrid row-column store table, if the plantree does not contain scenarios that do not support vectorization, the vectorized executor is forcibly used.

Type: USERSET

Value range: Boolean

Default value: off

enable_csqual_pushdown

Parameter description: Specifies whether to deliver filter criteria for a rough check during query.

Type: USERSET

Value range: Boolean

- **on** indicates that a rough check is performed with filter criteria delivered during query.
- **off** indicates that a rough check is performed without filter criteria delivered during query.

Default value: on

explain_dna_file

Parameter description: Specifies the name of a CSV file exported when **explain_perf_mode** is set to **run**.

Type: USERSET

NOTICE

The value of this parameter must be an absolute path plus a file name with the extension .csv.

Value range: a string
Default value: NULL

explain_perf_mode

Parameter description: Specifies the display format of the **explain** command.

Type: USERSET

Value range: normal, pretty, summary, and run

- normal indicates that the default printing format is used.
- pretty indicates that the optimized display mode of GaussDB(DWS) is used. A
 new format contains a plan node ID, directly and effectively analyzing
 performance.
- **summary** indicates that the analysis result based on such information is printed in addition to the printed information in the format specified by **pretty**.
- **run** indicates that in addition to the printed information specified by **summary**, the database exports the information as a CSV file.

Default value: pretty

join_num_distinct

Parameter description: Controls the default distinct value of the join column or expression in application scenarios.

Type: USERSET

Value range: a double-precision floating point number greater than or equal to **– 100**. Decimals may be truncated when displayed on clients.

- If the value is greater than **0**, the value is used as the default distinct value.
- If the value is greater than or equal to -100 and less than 0, it means the percentage used to estimate the default distinct value.
- If the value is **0**, the default distinct value is **200**.

Default value: -20

outer_join_max_rows_multipler

Parameter description: Specifies the maximum number of estimated rows for outer joins.

Type: USERSET

Value range: **0** or a double-precision floating point number greater than or equal to **1**. Decimals may be truncated when displayed on clients.

- If the value is **0**, the estimated number of rows for outer joins is not limited.
- If the value is greater than or equal to **1**, the estimated number of rows cannot exceed a multiple of the number of rows in the foreign table in the outer join.

Default value: 1.1

qual_num_distinct

Parameter description: Controls the default distinct value of the filter column or expression in application scenarios.

Type: USERSET

Value range: a double-precision floating point number greater than or equal to – **100**. Decimals may be truncated when displayed on clients.

- If the value is greater than **0**, the value is used as the default distinct value.
- If the value is greater than or equal to **-100** and less than **0**, it means the percentage used to estimate the default distinct value.
- If the value is **0**, the default distinct value is **200**.

Default value: 200

trace_notify

Parameter description: Specifies whether to generate a large amount of debugging output for the **LISTEN** and **NOTIFY** commands. **client_min_messages** or **log_min_messages** must be **DEBUG1** or lower so that such output can be recorded in the logs on the client or server separately.

Type: USERSET

Value range: Boolean

- on indicates that the function is enabled.
- off indicates that the function is disabled.

Default value: off

trace sort

Parameter description: Specifies whether to display information about resource usage during sorting operations in logs. This parameter is available only when the macro TRACE_SORT is defined during the GaussDB(DWS) compilation. However, TRACE SORT is currently defined by default.

Type: USERSET

Value range: Boolean

- on indicates that the function is enabled.
- **off** indicates that the function is disabled.

Default value: off

zero_damaged_pages

Parameter description: Specifies whether to detect a damaged page header that causes GaussDB(DWS) to report an error, aborting the current transaction.

Type: SUSET

Value range: Boolean

- **on** indicates that the function is enabled.
- off indicates that the function is disabled.

- Setting this parameter to **on** causes the system to report a warning, pad the damaged page with zeros, and then continue with subsequent processing. This behavior will damage data, that is, all rows on the damaged page. However, it allows you to bypass the error and retrieve rows from any undamaged pages that are present in the table. Therefore, it is useful for restoring data that is damaged due to a hardware or software error. In most cases, you are not advised to set this parameter to **on** unless you do not want to restore data from the damaged pages of a table.
- For a column-store table, the system will skip the entire CU and then continue processing. The supported scenarios include the CRC check failure, magic check failure, and incorrect CU length.

Default value: off

replication_test

Parameter description: Specifies whether to enable internal testing on the data replication function.

Type: USERSET

Value range: Boolean

- **on** indicates that internal testing on the data replication function is enabled.
- **off** indicates that internal testing on the data replication function is disabled.

Default value: off

cost_param

Parameter description: Controls use of different estimation methods in specific customer scenarios, allowing estimated values approximating to onsite values. This parameter can control various methods simultaneously by performing AND (&) operations on the bit for each method. A method is selected if its value is not **0**.

If **cost_param & 1** is not set to **0**, an improvement mechanism is selected for calculating a non-equi join selection rate, which is more accurate in estimation of self-join (join between two same tables). In V300R002C00 and later, **cost_param & 1=0** is not used. That is, an optimized formula is selected for calculation.

When **cost_param & 2** is set to a value other than **0**, the selection rate is estimated based on multiple filter criteria. The lowest selection rate among all filter criteria, but not the product of the selection rates for two tables under a specific filter criterion, is used as the total selection rate. This method is more accurate when a close correlation exists between the columns to be filtered.

When **cost_param & 4** is not **0**, the selected debugging model is not recommended when the stream node is evaluated.

When **cost_param & 16** is not **0**, the model between fully correlated and fully uncorrelated models is used to calculate the comprehensive selection rate of two or more filtering conditions or join conditions. If there are many filtering conditions, the strongly-correlated model is preferred.

Type: USERSET

Value range: an integer ranging from 1 to INT_MAX

Default value: 16

convert_string_to_digit

Parameter description: Specifies the implicit conversion priority, which determines whether to preferentially convert strings into numbers.

Type: USERSET

Value range: Boolean

- **on** indicates that strings are preferentially converted into numbers.
- off indicates that strings are not preferentially converted into numbers.

Default value: on

NOTICE

Modify this parameter only when absolutely necessary because the modification will change the rule for converting internal data types and may cause unexpected results.

nls_timestamp_format

Parameter description: Specifies the default timestamp format.

Type: USERSET

Value range: a string

Default value: DD-Mon-YYYY HH:MI:SS.FF AM

enable_partitionwise

Parameter description: Specifies whether to select an intelligent algorithm for joining partitioned tables.

Type: USERSET

Value range: Boolean

- **on** indicates that an intelligent algorithm is selected.
- off indicates that an intelligent algorithm is not selected.

Default value: off

enable_partition_dynamic_pruning

Parameter description: Specifies whether dynamic pruning is enabled during partition table scanning.

Type: USERSET

Value range: Boolean

on: enableoff: disable

Default value: on

max_user_defined_exception

Parameter description: Specifies the maximum number of exceptions. The default value cannot be changed.

Type: USERSET

Value range: an integer

Default value: 1000

datanode_strong_sync

Parameter description: This parameter no longer takes effect.

Type: USERSET

Value range: Boolean

- **on** indicates that forcible synchronization between stream nodes is enabled.
- **off** indicates that forcible synchronization between stream nodes is disabled.

Default value: off

enable_global_stats

Parameter description: Specifies the current statistics mode. This parameter is used to compare global statistics generation plans and the statistics generation plans for a single DN. This parameter is used for tests. Users are not allowed to enable it.

Type: SUSET

Value range: Boolean

- **on** or **true** indicates the global statistics mode.
- **off** or **false** indicates the single-DN statistics mode.

Default value: on

enable_fast_numeric

Parameter description: Specifies whether to enable optimization for numeric data calculation. Calculation of numeric data is time-consuming. Numeric data is converted into int64- or int128-type data to improve numeric data calculation performance.

Type: USERSET

Value range: Boolean

on/true indicates that optimization for numeric data calculation is enabled.

off/false indicates that optimization for numeric data calculation is disabled.

Default value: on

enable row fast numeric

Parameter description: Specifies the format in which numeric data in a row-store table is spilled to disks.

Type: USERSET

Value range: Boolean

- **on/true** indicates that numeric data in a row-store table is spilled to disks in bigint format.
- **off/false** indicates that numeric data in a row-store table is spilled to disks in the original format.

NOTICE

If this parameter is set to **on**, you are advised to enable **enable_force_vector_engine** to improve the query performance of large data sets. However, compared with the original format, there is a high probability that the bigint format occupies more disk space. For example, the TPC-H test set occupies about 7% more space (reference value, may vary depending on the environment).

Default value: off

rewrite rule

Parameter description: Specifies the rewriting rule for enabled optional queries. Some query rewriting rules are optional. Enabling them cannot always improve query efficiency. In a specific customer scenario, you can set the query rewriting rules through the GUC parameter to achieve optimal query efficiency.

This parameter can control the combination of query rewriting rules, for example, there are multiple rewriting rules: rule1, rule2, rule3, and rule4. To set the parameters, you can perform the following operations:

set rewrite_rule=rule1; --Enable query rewriting rule rule1.
set rewrite_rule=rule2,rule3; --Enable query rewriting rules rule2 and rule3.
set rewrite_rule=none; --Disable all optional query rewriting rules.

Type: USERSET

Value range: a string

- none: No optional query rewrite rules are used.
- **Lazyagg**: The Lazy Agg query rewrite rule is used to eliminate aggregate operations in subqueries.
- **magicset**: The Magic Set query rewrite rule is used to push conditions from the main query down to promoted sublinks.
- **uniquecheck**: Uses the Unique Check rewriting rule. (The scenario where the target column does not contain the expression sublink of the aggregate

function can be improved. The function can be enabled only when the value of the target column is unique after the sublink is aggregated based on the associated column. This function is recommended to be used by optimization engineers.)

- **disablerep**: Uses the function that prohibits pulling up sublinks of the replication table. (Disables sublink pull-up for the replication table.)
- **projection_pushdown**: the Projection Pushdown rewriting rule (Removes columns that are not used by the parent query from the subquery).
- **or_conversion**: the OR conversion rewriting rule (eliminates the association OR conditions that are inefficient to execute).
- **plain_lazyagg**: the **Plain Lazy Agg** query rewriting rule (eliminates aggregation operations in a single subquery). This option is supported only by clusters of version 8.1.3.100 or later.
- **eager_magicset**: Uses the **eager_magicset** query rewriting rule (to push down conditions from the main query to subqueries). This option is supported only by clusters of version 8.2.0 or later.
- casewhen_simplification: This rewrite rule uses the CASE WHEN statement to simplify queries. When enabled, it rewrites (case when xxx then const1 else const2)=const1. This option is supported only by clusters of version 8.3.0 or later.
- **outer_join_quality_imply**: When there is an equi-join condition between a left outer join and a right outer join, this rule pushes the expression condition on the outer table's join column down to the inner table's join column. This option is supported only by clusters of version 8.3.0 or later.
- **inlist_merge**: This query rewrite rule uses the **inlist_or_inlist** method to merge **OR** statements with the same base table column. When enabled, it merges and rewrites (**where a in (list1) or a in (list2))** to support **inlist2join**. This option is supported only by clusters of version 8.3.0 or later.
- **subquery_qual_pull_up**: For subqueries that cannot be promoted, if the subquery has filtering conditions on columns that are also used for joining with other tables, this rule extracts the filtering conditions from the subquery and passes them to the other side of the join condition. Currently, only **var op const** forms without type conversion, such as **a > 2**, are supported. When enabled, it is assumed that **outer_join_quality_imply** is also enabled. This is supported only by clusters of version 9.1.0 or later.
- not_distinct_from_opt: When this function is enabled, the is not distinct
 from expression is temporarily rewritten to an equal sign expression so that
 statements containing this syntax can generate hashjoin or mergejoin plans,
 improving execution performance. This is supported only by clusters of version
 9.1.0.218 or later.

Default value: magicset, or_conversion, projection_pushdown, plain_lazyagg, or subquery_qual_pull_up

mv rewrite rule

Parameter description: whether to enable the rewriting rule for the materialized view.

Type: USERSET

Value range: a string

- **none**: No materialized view rewriting rule is used. This value is available only in clusters of version 8.2.1.100 or later.
- **text**: materialized view rewriting rule that uses text matching. This value is available only in clusters of version 8.2.1.100 or later.
- **general**: indicates whether to enable structure matching. This value is supported only by 9.1.0.200 and later cluster versions.
- predicate: specifies whether to enable the expression framework for structure matching rewriting. This parameter is used for preprocessing and regularization of filter and association conditions. This parameter takes effect only when general is enabled. This value is supported only by 9.1.0.210 and later cluster versions.
- **view_delta**: indicates whether to use the primary and foreign key relationships to eliminate redundant joins in materialized views and then match the materialized views with queries. This value is supported only by 9.1.0.210 and later cluster versions.

Default value: text, general, and predicate

NOTICE

general, **predicate**, and **view_delta** are restricted for commercial use. To use them, contact technical support.

enable_compress_spill

Parameter description: Specifies whether to enable the compression function of writing data to a disk.

Type: USERSET

Value range: Boolean

- **on/true** indicates that optimization for writing data to a disk is enabled.
- off/false indicates that optimization for writing data to a disk is disabled.

Default value: on

analysis_options

Parameter description: Specifies whether to enable corresponding features, such as data validation and performance statistics.

Type: USERSET

Value range: a string

- **LLVM_COMPILE** indicates that the codegen compilation time of each thread is displayed on the explain performance page.
- **HASH_CONFLICT** indicates that the log file in the **pg_log** directory of the DN process displays the hash table statistics, including the hash table size, hash chain length, and hash conflict information.
- **STREAM_DATA_CHECK** indicates that a CRC check is performed on data before and after network data transmission.

- TURBO_DATA_CHECK indicates that the data context of the ScalarVector and VectorBatch operators of Turbo is verified. This parameter is supported only by clusters of version 8.3.0.100 or later.
- **KEEP_SAMPLE_DATA**: This parameter retains the sampling data used in each analyze operation in the form of temporary tables. This parameter is supported only by clusters of version 9.1.0 or later.
- BLOCK_RULE: indicates that the time required for checking the query filter is displayed on the explain performance page. This is supported only by 9.1.0.100 and later cluster versions.

Default value: **off(ALL)**, which indicates that no location function is enabled.

resource_track_log

Parameter description: Specifies the log level of self-diagnosis. Currently, this parameter takes effect only in multi-column statistics.

Type: USERSET

Value range: a string

- **summary**: Brief diagnosis information is displayed.
- **detail**: Detailed diagnosis information is displayed.

Currently, the two parameter values differ only when there is an alarm about multi-column statistics not collected. If the parameter is set to **summary**, such an alarm will not be displayed. If it is set to **detail**, such an alarm will be displayed.

Default value: summary

hll default log2m

Parameter description: Specifies the number of buckets for HLL data. The number of buckets affects the precision of distinct values calculated by HLL. As the number of buckets increases, the deviation becomes smaller. The deviation range is as follows: $[-1.04/2^{\log 2m^*1/2}, +1.04/2^{\log 2m^*1/2}]$

Type: USERSET

Value range: an integer ranging from 10 to 16

Default value: 11

hll_default_regwidth

Parameter description: Specifies the number of bits in each bucket for HLL data. A larger value indicates more memory occupied by HLL. **hll_default_regwidth** and **hll_default_log2m** determine the maximum number of distinct values that can be calculated by HLL. For details, see **Table 15-3**.

Type: USERSET

Value range: an integer ranging from 1 to 5

Default value: 5

log2m	regwidth =	regwidth = 2	regwidth = 3	regwidth = 4	regwidth = 5				
10	7.4e+02	3.0e+03	4.7e+04	1.2e+07	7.9e+11				
11	1.5e+03	5.9e+03	9.5e+04	2.4e+07	1.6e+12				
12	3.0e+03	1.2e+04	1.9e+05	4.8e+07	3.2e+12				
13	5.9e+03	2.4e+04	3.8e+05	9.7e+07	6.3e+12				
14	1.2e+04	4.7e+04	7.6e+05	1.9e+08	1.3e+13				
15	2.4e+04	9.5e+04	1.5e+06	3.9e+08	2.5e+13				

Table 15-3 Maximum number of calculated distinct values determined by hll default log2m and hll default requidth

hll_default_expthresh

Parameter description: Specifies the default threshold for switching from the **explicit** mode to the **sparse** mode.

Type: USERSET

Value range: an integer ranging from -1 to 7 -1 indicates the auto mode; **0** indicates that the **explicit** mode is skipped; a value from 1 to 7 indicates that the mode is switched when the number of distinct values reaches $2^{hll_default_expthresh}$.

Default value: -1

hll_default_sparseon

Parameter description: Specifies whether to enable the sparse mode by default.

Type: USERSET

 $\begin{tabular}{ll} \textbf{Valid value}: \textbf{0} and \textbf{1} \textbf{0} indicates that the \textbf{sparse} mode is disabled by default. \textbf{1} \\ \end{tabular}$

indicates that the **sparse** mode is enabled by default.

Default value: 1

hll_max_sparse

Parameter description: Specifies the size of max_sparse.

Type: USERSET

Value range: an integer ranging from -1 to INT_MAX

Default value: -1

enable_compress_hll

Parameter description: Specifies whether to enable memory optimization for HLL.

Type: USERSET

Value range: Boolean

on or **true** indicates that memory optimization is enabled. off or false indicates that memory optimization is disabled.

Default value: off

approx count distinct precision

Parameter description: Specifies the number of buckets in the HyperLogLog++ (HLL++) algorithm. This parameter can be used to adjust the error rate of the approx_count_distinct aggregate function. The number of buckets affects the precision of estimating the distinct value. Having more buckets increases the accuracy of the estimation. The deviation range is as follows: [-1.04/2log2m*1/2, +1.04/2^{log2m*1/2}]

Type: USERSET

Value range: an integer ranging from 10 to 20.

Default value: 17

udf memory limit

Parameter description: Controls the maximum physical memory that can be used when each CN or DN executes UDFs.

Type: POSTMASTER

Value range: an integer ranging from 200 x 1024 to the value of

max process memory and the unit is KB.

Default value: 0.05 * max_process_memory

FencedUDFMemoryLimit

Parameter description: Controls the virtual memory used by each fenced udf worker process.

Type: USERSET

Suggestion: You are not advised to set this parameter. You can set udf_memory_limit instead.

Value range: an integer. The unit can be KB, MB, or GB. 0 indicates that the

memory is not limited.

Default value: 0

enable pbe optimization

Parameter description: Specifies whether the optimizer optimizes the query plan for statements executed in Parse Bind Execute (PBE) mode.

Type: USERSET

Value range: Boolean

- **on** indicates that the optimizer optimizes the query plan.
- **off** indicates that the optimization does not optimize the query plan.

Default value: on

enable_light_proxy

Parameter description: Specifies whether the optimizer optimizes the execution of simple queries on CNs.

Type: USERSET

Value range: Boolean

- **on** indicates that the optimizer optimizes the execution.
- **off** indicates that the optimization does not optimize the execution.

Default value: on

enable_parallel_ddl

Parameter description: Controls whether multiple CNs can concurrently perform DDL operations on the same database object.

Type: USERSET

Value range: Boolean

- **on** indicates that DDL operations can be performed safely and that no distributed deadlock occurs.
- **off** indicates that DDL operations cannot be performed safely and that distributed deadlocks may occur.

Default value: on

gc_fdw_verify_option

Parameter description: Specifies whether to enable the logic for verifying the number of rows in a result set in the collaborative analysis. This parameter is supported only by clusters of version 8.1.3.310 or later.

Type: USERSET

Value range: Boolean

- on indicates that the logic for verifying the number of rows in the result set is enabled. The SELECT COUNT statement is used to obtain the expected number of rows and compare it with the actual number of rows.
- **off** indicates that the logic for verifying the number of rows in the result set is disabled and only the required result set is obtained.

Default value: on

- If this parameter is enabled, the performance deteriorates slightly. In performancesensitive scenarios, you can disable this parameter to improve the performance.
- If the result set row count check fails, an exception will be reported. To enable cooperative analysis logs, set log_min_messages to debug1 and logging_module to 'on(COOP_ANALYZE)'.

show_acce_estimate_detail

Parameter description: When the GaussDB(DWS) cluster is accelerated (acceleration_with_compute_pool is set to on), specifies whether the EXPLAIN statement displays the evaluation information about execution plan pushdown to computing Node Groups. The evaluation information is generally used by O&M personnel during maintenance, and it may affect the output display of the EXPLAIN statement. Therefore, this parameter is disabled by default. The evaluation information is displayed only if the verbose option of the EXPLAIN statement is enabled.

Type: USERSET

Value range: Boolean

- on indicates that the evaluation information is displayed in the output of the EXPLAIN statement.
- **off** indicates that the evaluation information is not displayed in the output of the **EXPLAIN** statement.

Default value: off

full_group_by_mode

Parameter description: Used in conjunction with disable_full_group_by_mysql in behavior_compat_options to control two different behaviors when disable_full_group_by_mysql syntax is enabled.

Type: USERSET

Value range: a string

- nullpadding indicates that NULL values in non-aggregate columns are filled with the non-NULL values in that column, potentially resulting in different rows in the result set.
- **notpadding** indicates that NULL values in non-aggregate columns are not processed, and the entire row data is used, resulting in a random row for non-aggregate columns in the result set.

Default value: notpadding

NOTICE

This parameter only takes effect when **disable_full_group_by_mysql** is enabled in the MySQL-compatible library and non-aggregate columns are present in the query. The two behaviors of this parameter only apply to non-aggregate columns in the query.

enable_cudesc_streaming

Parameter description: Specifies whether to use the cudesc streaming path for accessing data across logical clusters in the decoupled storage and compute architecture. This parameter is supported only by clusters of version 9.1.0 or later.

Type: SUSET

Value range: enumerated values

- off indicates that cudesc streaming is disabled.
- **on** indicates that cudesc streaming is enabled.
- **only_read_on** indicates that cudesc streaming is supported only during data reading.

Default value: on

force_read_from_rw

Parameter description: Forces data to be read from other logical clusters in the decoupled storage and compute architecture (i.e., read data from the logical cluster where the table resides). This parameter is supported only by clusters of version 9.0.0 or later.

Type: USERSET

Value range: Boolean

Default value: off

kv_sync_up_timeout

Parameter description: Specifies the timeout interval for KV synchronization in the decoupled storage and compute architecture. This parameter is supported only by clusters of version 9.0.0 or later.

Type: USERSET

Value range: an integer ranging from 0 to 2147483647

Default value: 10min

enable_insert_foreign_table_dop

Parameter description: Specifies whether to enable DOP acceleration when data is written into an OBS foreign table. The number of DOP threads on each DN is determined by **query_dop**. By adjusting its value, you can control the level of parallelism for your queries. This parameter is supported only in 9.1.0.200 and later versions.

Type: USERSET

Value range: Boolean

- **on** indicates that foreign table DOP acceleration is enabled.
- **off** indicates that foreign table DOP acceleration is disabled.

Default value: off

enable_insert_foreign_table_dop_opt

Parameter description: Specifies whether to enable partition redistribution optimization after **insert dop** is enabled for a foreign table. If the number of partitions to be exported is greater than 10 times the number of partitions, you are advised to enable this function to reduce small files in a single partition and improve export performance. This parameter is supported only by 9.1.0.200 and later versions.

Type: USERSET

Value range: Boolean

- **on** indicates that the **insert dop** redistribution optimization of the partitioned foreign table is enabled.
- **off** indicates that the insert dop redistribution optimization of the partitioned foreign table is disabled.

Default value: off

enable_recyclebin

Parameter description: Specifies whether the recycle bin is enabled or disabled. This parameter is supported only by clusters of versions later than 9.1.0.200.

Type: SUSET

Value range: enumerated values

- off indicates that the recycle bin is disabled.
- **on** indicates that the recycle bin is enabled.

Default value: on

NOTICE

You are advised to use this parameter in scenarios where the drop and truncate operations are not frequently performed.

recyclebin_retention_time

Parameter description: Specifies the retention period of objects in the recycle bin. The objects will be automatically deleted after the retention period expires. This parameter is supported only by clusters of 9.1.0.200 and later versions.

Type: SUSET

Value range: an integer ranging from 1 to 2147483. The unit is s.

Default value: 3600s

NOTICE

You are advised to use this parameter in scenarios where the drop and truncate operations are not frequently performed.

enable hstore binlog table

Parameter description: Specifies whether binlog tables can be created.

Type: SIGHUP

Value range: Boolean

on indicates that binlog tables can be created.

• off indicates that binlog tables cannot be created.

Default value: off

enable_generate_binlog

Parameter description: Specifies whether binlogs are generated for DML operations on binlog tables in the current session. This parameter is available only for clusters of version 9.1.0.200 or later.

Type: USERSET

Value range: Boolean

• **on** indicates that binlogs are generated.

• off indicates that binlogs are not generated.

Default value: on

binlog_consume_timeout

Parameter description: Specifies the duration for cyclically determining whether all binlog records are consumed during binlog table scaling or **VACUUM FULL** operations. This parameter is available only for clusters of version 8.3.0.100 or later. The unit is second.

Type: SIGHUP

Value range: an integer ranging from 0 to 86400

Default value: 3600

15.21 Auditing

15.21.1 Audit Switch

audit enabled

Parameter description: Specifies whether to enable or disable the audit process. After the audit process is enabled, the auditing information written by the background process can be read from the pipe and written into audit files.

Type: SIGHUP

Value range: Boolean

- **on** indicates that the auditing function is enabled.
- off indicates that the auditing function is disabled.

Default value: on

audit_space_limit

Parameter description: Specifies the total disk space occupied by audit files.

Type: SIGHUP

Value range: an integer ranging from 1024 to 1073741824. The unit is KB.

Default value: 1GB

audit_object_name_format

Parameter description: Specifies the format of the object name displayed in the **object name** field of audit logs.

Type: USERSET

Value range: enumerated values

- **single** indicates that the **object_name** field displays a single object name, which is the name of the target object.
- all indicates that the object_name field displays multiple object names.

Default value: single

■ NOTE

If the default value is set to **all**, multiple object names are displayed for SELECT, DELETE, UPDATE, INSERT, MERGE, CREATE TABLE AS, CREATE VIEW AS, DROP USER... CASCADE, DROP OWNED BY... CASCADE, DROP SCHEMA... CASSCADE, DROP TABLE... CASCADE, DROP FOREIGN TABLE... CASCADE, and DROP VIEW... CASCADE.

audit_object_details

Parameter description: whether to record the **object_details** field in audit logs. This field indicates the table name, column name, and column type in the audit statement. This parameter is supported only by clusters of version 8.2.1.100 or later.

Type: USERSET

Value range: Boolean

- **on** indicates that the **object_details** field is recorded during the audit.
- **off** indicates that the **object_details** field is not recorded during the audit.

Default value: off

□ NOTE

- If this parameter is set to **on**, the table name, column name, and column type in the statement will be audited, which may affect the performance. So, exercise caution when setting this parameter to **on**.
- If this parameter is set to on, the object_details field records the following statements:
 SELECT, DELETE, UPDATE, INSERT, MERGE, CREATE TABLE AS SELECT, GRANT, and DECLARE CURSOR. GRANT statements that fail to be executed are not recorded.

15.21.2 Operation Audit

security_enable_options

Parameter description: This parameter lets you use the **grant_to_public**, **grant_with_grant_option**, and **foreign_table_options** functions. These functions are off by default for security. Set this parameter as needed. It works with clusters version 8.2.0 and later.

Type: SIGHUP

Value range: a string

- grant_to_public indicates that grant to public can be used.
- grant_with_grant_option indicates that with grant option can be used.
- **foreign_table_options** allows users to perform operations on foreign tables without explicitly granting the **useft** permission to users.

Default value: empty

□ NOTE

- In a newly installed cluster, this parameter is left blank by default, indicating that none
 of grant_to_public, grant_with_grant_option, and foreign_table_options can be used.
- In upgrade scenarios, the default value of this parameter is forward compatible. If the
 default values of enable_grant_public and enable_grant_option are ON before the
 upgrade, the default value of security_enable_options is grant_to_public,
 grant_with_grant_option after the upgrade.

15.22 Transaction Monitoring

By setting transaction timeout alerts, you can monitor transactions that are automatically rolled back and identify statement issues, as well as monitor statements that take too long to execute.

transaction_sync_naptime

Parameter description: For data consistency, when the local transaction's status differs from that in the snapshot of the GTM, other transactions will be blocked. You need to wait for a few minutes until the transaction status of the local host is consistent with that of the GTM. The **gs_clean** tool is automatically triggered for cleansing when the waiting period on the CN exceeds that of **transaction_sync_naptime**. The tool will shorten the blocking time after it completes the cleansing.

Type: USERSET

Value range: an integer. The minimum value is **0**. The unit is second.

Default value: 5s

□ NOTE

If the value of this parameter is set to **0**, gs_clean will not be automatically invoked for the cleansing before the blocking arrives the duration. Instead, the gs_clean tool is invoked by gs_clean_timeout. The default value is 5 minutes.

transaction sync timeout

Parameter description: For data consistency, when the local transaction's status differs from that in the snapshot of the GTM, other transactions will be blocked. You need to wait for a few minutes until the transaction status of the local host is consistent with that of the GTM. An exception is reported when the waiting duration on the CN exceeds the value of **transaction_sync_timeout**. Roll back the transaction to avoid system blocking due to long time of process response failures (for example, sync lock).

Type: USERSET

Value range: an integer. The minimum value is **0**. The unit is second.

Default value: 10min

□ NOTE

- If the value is 0, no error is reported when the blocking times out or the transaction is rolled back.
- The value of this parameter must be greater than **gs_clean_timeout**. Otherwise, unnecessary transaction rollback will probably occur due to a block timeout caused by residual transactions that have not been deleted by **gs_clean** on a DN.

15.23 GTM Parameters

log_min_messages

Parameter description: Specifies which level of messages will be written into server logs. Each level covers all the levels following it. The lower the level is, the fewer messages will be written into the log.

NOTICE

If the values of **client_min_messages** and **log_min_messages** are the same, they indicate different levels.

Type: SUSET

Valid values: enumerated values. Valid values are debug, debug5, debug4, debug3, debug1, info, log, notice, warning, error, fatal, and panic. For details about the parameters, see Table 15-1.

Default value: warning

15.24 Miscellaneous Parameters

server version

Parameter description: Specifies the server version number in the string format.

Type: INTERNAL ((Fixed parameter. You are not advised to configure this parameter because incorrect values may cause compatibility issues.)

Value range: a string
Default value: 9.2.4

server_version_num

Parameter description: Specifies the server version number in the integer format.

Type: INTERNAL ((Fixed parameter. You are not advised to configure this parameter because incorrect values may cause compatibility issues.)

Value range: an integer

Default value: 90204

enable cluster resize

Parameter description: Indicates whether the current session is for scaling or redistributing data. It should only be used for these specific sessions and not set for other service sessions.

Type: SUSET

Value range: Boolean

- **on** indicates that the current session is for scaling or redistributing data, and allows the execution of specific SQL statements for redistribution.
- **off** indicates that the current session is not for scaling or redistributing data, and does not allow the execution of specific SQL statements for redistribution.

Default value: off

■ NOTE

This parameter is used for internal O&M. Do not set it to on unless absolutely necessary.

dfs_partition_directory_length

Parameter description: Specifies the largest directory name length for the partition directory of a table partitioned by VALUE in the HDFS.

Type: USERSET

Value range: 92 to 7999

Default value: 512

enable hadoop env

Parameter description: Sets whether local row- and column-store tables can be created in a database while the Hadoop feature is used. In the GaussDB(DWS) cluster, it is set to **off** by default to support local row- and column- based storage and cross-cluster access to Hadoop. You are not advised to change the value of this parameter.

Type: USERSET

Value range: Boolean

- **on** or **true**, indicating that local row- and column-store tables cannot be created in a database while the Hadoop feature is used.
- **off** or **false**, indicating that local row- and column-based tables can be created in a database while the Hadoop feature is used.

Default value: off

enable_upgrade_merge_lock_mode

Parameter description: If this parameter is set to **on**, the delta merge operation internally increases the lock level, and errors can be avoided when update and delete operations are performed at the same time.

Type: USERSET

Value range: Boolean

- If this parameter is set to **on**, the delta merge operation internally increases the lock level. In this way, when any two of the **DELTAMERGE**, **UPDATE**, and **DELETE** operations are concurrently performed, an operation can be performed only after the previous one is complete.
- If this parameter is set to off, and any two of the DELTAMERGE, UPDATE, and DELETE operations are concurrently performed to data in a row in the delta table of the HDFS table, errors will be reported during the later operation, and the operation will stop.

Default value: off

job_queue_processes

Parameter description: Specifies the number of jobs that can be concurrently executed.

Type: POSTMASTER

Value range: 0 to 1000

Functions:

Setting job_queue_processes to 0 indicates that the scheduled task function
is disabled and that no job will be executed. (Enabling scheduled tasks may
affect the system performance. At sites where this function is not required,
you are advised to disable it.)

• Setting **job_queue_processes** to a value that is greater than **0** indicates that the scheduled task function is enabled and this value is the maximum number of tasks that can be concurrently processed.

After the scheduled task function is enabled, the **job_scheduler** thread at a scheduled interval polls the **pg_jobs** system catalog. The scheduled task check is performed every second by default.

Too many concurrent tasks consume many system resources, so you need to set the number of concurrent tasks to be processed. If the current number of concurrent tasks reaches <code>job_queue_processes</code> and some of them expire, these tasks will be postponed to the next polling period. Therefore, you are advised to set the polling interval (the <code>interval</code> parameter of the <code>submit</code> API) based on the execution duration of each task to avoid the problem that tasks in the next polling period cannot be properly processed because overlong task execution time.

Note: If the number of parallel jobs is large and the value is too small, these jobs will wait in queues. However, a large parameter value leads to large resource consumption. You are advised to set this parameter to **100** and change it based on the system resource condition.

Default value: 10

ngram_gram_size

Parameter description: Specifies the length of the ngram parser segmentation.

Type: USERSET

Value range: an integer ranging from 1 to 4

Default value: 2

ngram_grapsymbol_ignore

Parameter description: Specifies whether the ngram parser ignores graphical characters.

Type: USERSET

Value range: Boolean

• **on**: Ignores graphical characters.

• **off**: Does not ignore graphical characters.

Default value: off

ngram_punctuation_ignore

Parameter description: Specifies whether the ngram parser ignores punctuations.

Type: USERSET

Value range: Boolean

on: Ignores punctuations.

• off: Does not ignore punctuations.

Default value: on

zhparser_multi_duality

Parameter description: Specifies whether Zhparser aggregates segments in long words with duality.

Type: USERSET

Value range: Boolean

on: Aggregates segments in long words with duality.

• **off**: Does not aggregate segments in long words with duality.

Default value: off

zhparser_multi_short

Parameter description: Specifies whether Zhparser executes long words composite divide.

Type: USERSET

Value range: Boolean

• **on**: Performs compound segmentation for long words.

• **off**: Does not perform compound segmentation for long words.

Default value: on

zhparser_multi_zall

Parameter description: Specifies whether Zhparser displays all single words individually.

Type: USERSET

Value range: Boolean

on: Displays all single words separately.

• **off**: Does not display all single words separately.

Default value: off

zhparser_multi_zmain

Parameter description: Specifies whether Zhparser displays important single words separately.

Type: USERSET

Value range: Boolean

on: Displays important single words separately.

off: Does not display important single words separately.

Default value: off

zhparser_punctuation_ignore

Parameter description: Specifies whether the Zhparser segmentation result ignores special characters including punctuations (\r and \n will not be ignored).

Type: USERSET

Value range: Boolean

- on: Ignores all the special characters including punctuations.
- **off**: Does not ignore all the special characters including punctuations.

Default value: on

zhparser_seg_with_duality

Parameter description: Specifies whether Zhparser aggregates segments in long words with duality.

Type: USERSET

Value range: Boolean

- on: Aggregates segments in long words with duality.
- off: Does not aggregate segments in long words with duality.

Default value: off

acceleration with compute pool

Parameter description: Specifies whether to use the computing resource pool for acceleration when OBS is queried.

Type: USERSET

Value range: Boolean

- **on** indicates that the query covering OBS is accelerated based on the cost when the computing resource pool is available.
- off indicates that no query is accelerated using the computing resource pool.

Default value: off

redact_compat_options

Parameter description: Specifies the compatibility option for calculation using masked data. This parameter is supported only by clusters of version 8.1.3.310 or later.

Type: USERSET

Value range: a string

- **none** indicates that compatibility options are specified.
- **disable_comparison_operator_mask** indicates that comparison operators that do not expose raw data can bypass the data masking check and generate the actual calculation result.

Default value: none

table_skewness_warning_threshold

Parameter description: Specifies the threshold for triggering a table skew alarm.

Type: SUSET

Value range: a floating point number ranging from 0 to 1

Default value: 1

table_skewness_warning_rows

Parameter description: Specifies the minimum number of rows for triggering a table skew alarm.

Type: SUSET

Value range: an integer ranging from 0 to INT_MAX

Default value: 100000

enable_view_update

Parameter description: Enables the view update function or not.

Type: POSTMASTER

Value range: Boolean

- **on** indicates that the view update function is enabled.
- off indicates that the view update function is disabled.

Default value: off

view_independent

Parameter description: Decouples views from tables, functions, and synonyms or not. After the base table is restored, automatic association and re-creation are supported.

Type: SIGHUP

Value range: Boolean

- on indicates that the view decoupling function is enabled. Tables, functions, synonyms, and other views on which views depend can be deleted separately (except temporary tables and temporary views). Associated views are reserved but unavailable.
- **off** indicates that the view decoupling function is disabled. Tables, functions, synonyms, and other views on which views depend cannot be deleted separately. You can only delete them in the cascade mode.

Default value: off

assign_abort_xid

Parameter description: Determines the transaction to be aborted based on the specified XID in a guery.

Type: USERSET

Value range: a character string with the specified XID



This parameter is used only for quick restoration if a user deletes data by mistake (DELETE operation). Do not use this parameter in other scenarios. Otherwise, visible transaction errors may occur.

default distribution mode

Parameter description: Specifies the default distribution mode of a table. This feature is supported only in 8.1.2 or later.

Type: USERSET

Value range: enumerated values

- **roundrobin**: If the distribution mode is not specified during table creation, the default distribution mode is selected according to the following rules:
 - a. If the primary key or unique constraint is included during table creation, hash distribution is selected. The distribution column is the column corresponding to the primary key or unique constraint.
 - b. If the primary key or unique constraint is not included during table creation, round-robin distribution is selected.
- hash: If the distribution mode is not specified during table creation, the default distribution mode is selected according to the following rules:
 - a. If the primary key or unique constraint is included during table creation, hash distribution is selected. The distribution column is the column corresponding to the primary key or unique constraint.
 - b. If the primary key or unique constraint is not included during table creation but there are columns whose data types can be used as distribution columns, hash distribution is selected. The distribution column is the first column whose data type can be used as a distribution column.
 - c. If the primary key or unique constraint is not included during table creation and no column whose data type can be used as a distribution column exists, round-robin distribution is selected.

Default value: roundrobin

◯ NOTE

The default value of this parameter is **roundrobin** for a new GaussDB(DWS) 8.1.2 cluster and is **hash** for an upgrade to GaussDB(DWS) 8.1.2.

object_mtime_record_mode

Parameter description: Sets the update action of the **mtime** column in the **PG_OBJECT** system catalog.

Type: SIGHUP

Value range: a string

- default: ALTER, COMMENT, GRANT/REVOKE, and TRUNCATE operations update the mtime column by default.
- **disable**: The **mtime** column is not updated.
- **disable_acl**: **GRANT** or **REVOKE** operation does not update the **mtime** column.
- **disable_truncate**: **TRUNCATE** operations do not update the **mtime** column.
- disable_partition: Partition ALTER operations do not update the mtime column.

Default value: default

max_volatile_tables

Parameter description: Specifies the maximum number of volatile tables created for each session, including volatile tables and their auxiliary tables. This parameter is supported by clusters of version 8.2.0 or later.

Type: USERSET

Value range: an integer ranging from 0 to INT_MAX

Default value: 300

query cache refresh time

Parameter description: Specifies the cache refresh interval for queries for which the **enable_accelerate_select** parameter takes effect. This parameter is supported only by clusters of version 8.3.0 or later.

Type: USERSET

Value range: a floating point number ranging from 0 to 10000.0, in seconds

Default value: 60.0

vector_engine_strategy

Parameter description: Specifies the vectorization enhancement policy. This parameter is supported only by clusters of version 8.3.0 or later.

Type: USERSET

Value range: enumerated values

• **force** specifies that the vectorization-enhanced plan is forcibly rolled back to the row storage plan when there are scenarios that do not support vectorization.

• **improve** specifies that vectorization enhancement is enabled even when there are scenarios that do not support vectorization.

Default value: improve

default_temptable_type

Parameter description: Specifies the type of temporary table created when **CREATE TABLE** is used to create a temporary table without specifying the table type before **TEMP** or **TEMPORARY**. This parameter is supported only by clusters of version 9.1.0 or later.

Type: USERSET

Value range: enumerated values

- **local**: creates a local temporary table when the type is not specified.
- **volatile**: creates a volatile temporary table when the type is not specified.

Default value: local

pgxc_node_readonly

Parameter description: Specifies whether a CN or DN is an elastic or classic DN. This parameter is supported only by clusters of version 9.1.0 or later.

Type: SUSET

Value range: Boolean

- **on** indicates that the CN or DN is an elastic node.
- off indicates that the CN or DN is a classic node.

Default value: off

hudi sync max commits

Parameter description: Specifies the maximum number of commits for a single synchronization task in Hudi. This parameter is supported only by clusters of version 9.1.0.100 or later.

Type: SIGHUP

Value range: an integer ranging from -1 to INT_MAX

- -1 indicates no limit.
- 0 indicates no limit.
- Any other value indicates the maximum number of commits.

Default value: -1

foreign_table_default_rw_options

Parameter description: Specifies the default permissions when creating a foreign table without specifying them. This parameter is supported only by clusters of version 9.0.3 or later.

Type: USERSET

Value range: a string

READ_ONLY indicates the read-only permission.WRITE_ONLY indicates the write-only permission.

• **READ_WRITE** indicates the read-write permission.

Default value: READ_ONLY

16 GaussDB(DWS) Developer Terms

A~E

Table 16-1 Terms A to E

Term	Description
ACID	Four essential properties that a transaction should have in a DBMS: Atomicity, Consistency, Isolation, and Durability.
cluster ring	A cluster ring consists of several physical servers. The primary-standby-secondary relationships among its DNs do not involve external DNs. That is, none of the primary, standby, or secondary counterparts of DNs belonging to the ring are deployed in other rings. A ring is the smallest unit used for scaling.
Bgwriter	A background write thread created when the database starts. The thread pushes dirty pages in the database to a permanent device (such as a disk).
bit	The smallest unit of information handled by a computer. One bit is expressed as a 1 or a 0 in a binary numeral, or as a true or a false logical condition. A bit is physically represented by an element such as high or low voltage at one point in a circuit, or a small spot on a disk that is magnetized in one way or the other. A single bit conveys little information a human would consider meaningful. A group of eight bits, however, makes up a byte, which can be used to represent many types of information, such as a letter of the alphabet, a decimal digit, or other character.
Bloom filter	Bloom filter is a space-efficient binary vectorized data structure, conceived by Burton Howard Bloom in 1970, that is used to test whether an element is a member of a set. False positive matches are possible, but false negatives are not, in other words, a query returns either "possibly in set (possible error)" or "definitely not in set". In the cases, Bloom filter sacrificed the accuracy for time and space.

Term	Description
CCN	The Central Coordinator (CCN) is a node responsible for determining, queuing, and scheduling complex operations in each CN to enable the dynamic load management of GaussDB(DWS).
CIDR	Classless Inter-Domain Routing (CIDR). CIDR abandons the traditional class-based (class A: 8; class B: 16; and class C: 24) address allocation mode and allows the use of address prefixes of any length, effectively improving the utilization of address space. A CIDR address is in the format of <i>IP address Number of bits in a network ID</i> . For example, in 192.168.23.35/21, 21 indicates that the first 21 bits are the network prefix and others are the host ID.
Cgroups	A control group (Cgroup), also called a priority group (PG) in GaussDB(DWS). The Cgroup is a kernel feature of SUSE Linux and Red Hat that can limit, account for, and isolate the resource usage of a collection of processes.
CLI	Command-line interface (CLI). Users use the CLI to interact with applications. Its input and output are based on texts. Commands are entered through keyboards or similar devices and are compiled and executed by applications. The results are displayed in text or graphic forms on the terminal interface.
СМ	Cluster Manager (CM) manages and monitors the running status of functional units and physical resources in the distributed system, ensuring stable running of the entire system.
CMS	The Cluster Management Service (CMS) component manages the cluster status.
CN	The Coordinator (CN) stores database metadata, splits query tasks and supports their execution, and aggregates the query results returned from DNs.
СТЕ	A common table expression (CTE) is a named temporary result set that exists within the scope of a single statement and that can be referred to later in SELECT, INSERT, UPDATE, or DELETE statements multiple times. It simplifies complex queries, and enhances code readability and maintainability. CTEs can significantly improve the maintainability of complex SQL statements, and have obvious advantages in recursive or step-by-step data processing.
CU	A compression unit is the smallest storage unit of a column- store table.

Term	Description
core file	A file that is created when memory overwriting, assertion failures, or access to invalid memory occurs in a process, causing it to fail. This file is then used for further analysis.
	A core file contains a memory dump, in an all-binary and port- specific format. The name of a core file consists of the word "core" and the OS process ID.
	The core file is available regardless of the type of platform.
core dump	When a program stops abnormally, the core dump, memory dump, or system dump records the state of the working memory of the program at that point in time. In practice, other key pieces of program state are usually dumped at the same time, including the processor registers, which may include the program counter and stack pointer, memory management information, and other processor and OS flags and information. A core dump is often used to assist diagnosis and computer program debugging.
DBA	A database administrator (DBA) instructs or executes database maintenance operations.
DBLINK	An object defining the path from one database to another. A remote database object can be queried with DBLINK.
DBMS	Database Management System (DBMS) is a piece of system management software that allows users to access information in a database. This is a collection of programs that allows you to access, manage, and query data in a database. A DBMS can be classified as memory DBMS or disk DBMS based on the location of the data.
DCL	Data control language (DCL)
DDL	Data definition language (DDL)
DFS	A Distributed File System (DFS) stores files on multiple physical nodes and allows access to the files via unified access interfaces. It is not a specific technique, but a system type.
	For example, Hadoop Distributed File System (HDFS) is a type of DFS and is designed for big data scenarios. It is suitable for high-throughput and batch processing tasks.
DML	Data manipulation language (DML)
DN	Datanode performs table data storage and query operations.
ETCD	The Editable Text Configuration Daemon (ETCD) is a distributed key-value storage system used for configuration sharing and service discovery (registration and search).
ETL	Extract-Transform-Load (ETL) refers to the process of data transmission from the source to the target database.

Term	Description
Extension Connector	Extension Connector is provided by GaussDB(DWS) to process data across clusters. It can send SQL statements to Spark, and can return execution results to your database.
Backup	A backup, or the process of backing up, refers to the copying and archiving of computer data in case of data loss.
backup and restoration	A collection of concepts, procedures, and strategies to protect data loss caused by invalid media or misoperations.
standby server	A node in the GaussDB(DWS) HA solution. It functions as a backup of the primary server. If the primary server is behaving abnormally, the standby server is promoted to primary, ensuring data service continuity.
crash	A crash (or system crash) is an event in which a computer or a program (such as a software application or an OS) ceases to function properly. Often the program will exit after encountering this type of error. Sometimes the offending program may appear to freeze or hang until a crash reporting service documents details of the crash. If the program is a critical part of the OS kernel, the entire computer may crash (possibly resulting in a fatal system error).
encoding	Encoding is representing data and information using code so that it can be processed and analyzed by a computer. Characters, digits, and other objects can be converted into digital code, or information and data can be converted into the required electrical pulse signals based on predefined rules.
encoding technology	A technology that presents data using a specific set of characters, which can be identified by computer hardware and software.
table	A set of columns and rows. Each column is referred to as a field. The value in each field represents a data type. For example, if a table contains people's names, cities, and states, it has three columns: Name, City, and State. In every row in the table, the Name column contains a name, the City column contains a city, and the State column contains a state.
tablespace	A tablespace is a logical storage structure that contains tables, indexes, large objects, and long data. A tablespace provides an abstract layer between physical data and logical data, and provides storage space for all database objects. When you create a table, you can specify which tablespace it belongs to.
concurrency control	A DBMS service that ensures data integrity when multiple transactions are concurrently executed in a multi-user environment. In a multi-threaded environment, GaussDB(DWS) concurrency control ensures that database operations are safe and all database transactions remain consistent at any given time.

Term	Description
query	A request sent to a database, which may include updates, modifications, retrieval, or deletions of information.
query operator	An iterator or a query tree node, which is a basic unit for the execution of a query. Execution of a query can be split into one or more query operators. Common query operators include scan, join, and aggregation.
query fragment	Each query task can be split into one or more query fragments. Each query fragment consists of one or more query operators and can independently run on a node. Query fragments exchange data through data flow operators.
durability	One of the ACID features of database transactions. Durability indicates that transactions that have been committed will permanently survive and not be rolled back.
stored procedure	A group of SQL statements compiled into a single execution plan and stored in a large database system. Users can specify a name and parameters (if any) for a stored procedure to execute the procedure.
OS	An operating system (OS) is loaded by a bootstrap program to a computer to manage other programs in the computer. Other programs are applications or application programs.
secondary server	To ensure high cluster availability, the primary server synchronizes logs to the secondary server if data synchronization between the primary and standby servers fails. If the primary server suddenly breaks down, the standby server is promoted to primary and synchronizes logs from the secondary server for the duration of the breakdown.
BLOB	Binary large object (BLOB) is a collection of binary data stored in a database, such as videos, audio, and images.
recursive query	Recursive query is an advanced technology in database query. It is mainly used to process hierarchical data or used in scenarios where layer-by-layer iterative calculation is required. It traverses and analyzes tree structures, graph structures, and chain relationships through recursive invocation. In the SELECT syntax, WITH RECURSIVE is often used to declare recursive queries. For the examples of using recursive queries, see SELECT .
dynamic load balancing	In GaussDB(DWS), dynamic load balancing automatically adjusts the number of concurrent jobs based on the usage of CPU, I/O, and memory to avoid service errors and to prevent the system from stop responding due to system overload.
segment	A segment in the database indicates a part containing one or more regions. Region is the smallest range of a database and consists of data blocks. One or more segments comprise a tablespace.

F~J

Table 16-2 Terms F to J

Term	Description
Failover	Automatic switchover from a faulty node to its standby node. Reversely, automatic switchback from the standby node to the primary node is called failback.
FDW	A foreign data wrapper (FDW) is a SQL interface provided by Postgres. It is used to access big data objects stored in remote data so that DBAs can integrate data from unrelated data sources and store them in public schema in the database.
freeze	An operation automatically performed by the AutoVacuum Worker process when transaction IDs are exhausted. GaussDB(DWS) records transaction IDs in row headings. When a transaction reads a row, the transaction ID in the row heading and the actual transaction ID are compared to determine whether this row is explicit. Transaction IDs are integers containing no symbols. If exhausted, transaction IDs are recalculated outside of the integer range, causing the explicit rows to become implicit. To prevent such a problem, the freeze operation marks a transaction ID as a special ID. Rows marked with these special transaction IDs are explicit to all transactions.
GDB	As a GNU debugger, GDB allows you to see what is going on 'inside' another program while it executes or what another program was doing the moment that it crashed. GDB can perform four main kinds of things (make PDK functions stronger) to help you catch bugs in the act: Starts your program, specifying anything that might affect its behavior. Stops a program in a specific condition. Checks what happens when a program stops.
	Modifies the program content to rectify the fault and proceeds with the next one.
GDS	General Data Service (GDS). To import data to GaussDB(DWS), you need to deploy the tool on the server where the source data is stored so that DNs can use this tool to obtain data.
GIN index	Generalized inverted index (GIN) is used for handling cases where the items to be indexed are composite values, and the queries to be handled by the index need to search for element values that appear within the composite items.

Term	Description
GNU	The GNU Project was publicly announced on September 27, 1983 by Richard Stallman, aiming at building an OS composed wholly of free software. GNU is a recursive acronym for "GNU's Not Unix!". Stallman announced that GNU should be pronounced as Guh-NOO. Technically, GNU is similar to Unix in design, a widely used commercial OS. However, GNU is free software and contains no Unix code.
gsql	GaussDB(DWS) interaction terminal. It enables you to interactively type in queries, issue them to GaussDB(DWS), and view the query results. Queries can also be entered from files. gsql supports many meta commands and shell-like commands, allowing you to conveniently compile scripts and automate tasks.
GTM	Global Transaction Manager (GTM) manages the status of transactions.
GUC	Grand unified configuration (GUC) includes parameters for running databases, the values of which determine database system behavior.
НА	High availability (HA) is a solution in which two modules operate in primary/standby mode to achieve high availability. This solution helps to minimize the duration of service interruptions caused by routine maintenance (planned) or sudden system breakdowns (unplanned), improving the system and application usability.
НВА	Host-based authentication (HBA) allows hosts to authenticate on behalf of all or some of the system users. It can apply to all users on a system or a subset using the Match directive. This type of authentication can be useful for managing computing clusters and other fairly homogenous pools of machines. In all, three files on the server and one on the client must be modified to prepare for host-based authentication.
HDFS	Hadoop Distributed File System (HDFS) is a subproject of Apache Hadoop. HDFS is highly fault tolerant and is designed to run on low-end hardware. The HDFS provides high-throughput access to large data sets and is ideal for applications having large data sets.
server	A combination of hardware and software designed for providing clients with services. This word alone refers to the computer running the server OS, or the software or dedicated hardware providing services.
advanced package	Logical and functional stored procedures and functions provided by GaussDB(DWS).

Term	Description
isolation	One of the ACID features of database transactions. Isolation means that the operations inside a transaction and data used are isolated from other concurrent transactions. The concurrent transactions do not affect each other.
relational database	A database created using a relational model. It processes data using methods of set algebra.
archive thread	A thread started when the archive function is enabled on a database. The thread archives database logs to a specified path.
failover	The automatic substitution of a functionally equivalent system component for a failed one. The system component can be a processor, server, network, or database.
environment variable	An environment variable defines the part of the environment in which a process runs. For example, it can define the part of the environment as the main directory, command search path, terminal that is in use, or the current time zone.
checkpoint	A mechanism that stores data in the database memory to disks at a certain time. GaussDB(DWS) periodically stores the data of committed and uncommitted transactions to disks. The data and redo logs can be used for database restoration if a database restarts or breaks down.
encryption	A function hiding information content during data transmission to prevent the unauthorized use of the information.
node	Cluster nodes (or nodes) are physical and virtual severs that make up the GaussDB(DWS) cluster environment.
error correction	A technique that automatically detects and corrects errors in software and data streams to improve system stability and reliability
process	An instance of a computer program that is being executed. A process may be made up of multiple threads of execution. Other processes cannot use a thread occupied by the process.
PITR	Point-In-Time Recovery (PITR) is a backup and restoration feature of GaussDB(DWS). Data can be restored to a specified point in time if backup data and WAL logs are normal.
record	In a relational database, a record corresponds to data in each row of a table.
cluster	A cluster is an independent system consisting of servers and other resources, ensuring high availability. In certain conditions, clusters can implement load balancing and concurrent processing of transactions.

K~O

Table 16-3 Terms K to O

Term	Description
LLVM	LLVM is short for Low Level Virtual Machine. Low Level Virtual Machine (LLVM) is a compiler framework written in C++ and is designed to optimize the compile-time, link-time, run-time, and idle-time of programs that are written in arbitrary programming languages. It is open to developers and compatible with existing scripts.
	GaussDB(DWS) LLVM dynamic compilation can be used to generate customized machine code for each query to replace original common functions. Query performance is improved by reducing redundant judgment conditions and virtual function invocation, and by making local data more accurate during actual queries.
LVS	Linux Virtual Server (LVS), a virtual server cluster system, is used for balancing the load of a cluster.
logical replication	Data synchronization mode between primary and standby databases or between two clusters. Different from physical replication which replays physical logs, logical replication transfers logical logs between two clusters or synchronizes data through SQL statements in logical logs.
logical log	Logs recording database changes made through SQL statements. Generally, the changes are logged at the row level. Logical logs are different from physical logs that record changes of physical pages.
logical decoding	Logic decoding is a process of extracting all permanent changes in database tables into a clear and easy-to-understand format by decoding Xlogs.
logical replication slot	In a logical replication process, logic replication slots are used to prevent Xlogs from being reclaimed by the system or VACUUM . In GaussDB(DWS), a logical replication slot is an object that records logical decoding positions. It can be created, deleted, read, and pushed by invoking SQL functions.
МРР	Massive Parallel Processing (MPP) refers to cluster architecture that consists of multiple machines. The architecture is also called a cluster system.
MVCC	Multi-Version Concurrency Control (MVCC) is a protocol that allows a tuple to have multiple versions, on which different query operations can be performed. A basic advantage is that read and write operations do not conflict.
NameNode	The NameNode is the centerpiece of a Hadoop file system, managing the namespace of the file system and client access to files.

Term	Description
Node Group	In GaussDB(DWS), a Node Group refers to a DN set, which is a sub-cluster. Node Groups can be classified into Storage Node Groups, which store local table data; and Computing Node Groups, which perform aggregation and join for queries.
NULL FIRST	In SQL, NULL FIRST is a clause used to explicitly control the positions of NULL values during sorting. It is usually used in the ORDER BY statement. It puts NULL values before non-NULL values. When NULL FIRST is used in ORDER BY, all rows with NULL values are displayed first in the result set, followed by the rows with non-null values, regardless of whether the sorting order is ASC or DESC.
OLAP	Online analytical processing (OLAP) is the most important application in the database warehouse system. It is dedicated to complex analytical operations, helps decision makers and executives to make decisions, and rapidly and flexibly processes complex queries involving a great amount of data based on analysts' requirements. In addition, the OLAP provides decision makers with query results that are easy to understand, allowing them to learn the operating status of the enterprise. These decision makers can then produce informed and accurate solutions based on the query results.
ОМ	Operations Management (OM) provides management interfaces and tools for routine maintenance and configuration management of the cluster.
ORC	Optimized Row Columnar (ORC) is a widely used file format for structured data in a Hadoop system. It was introduced from the Hadoop HIVE project.
client	A computer or program that accesses or requests services from another computer or program.
free space management	A mechanism for managing free space in a table. This mechanism enables the database system to record free space in each table and establish an easy-to-search data structure, accelerating operations (such as INSERT) performed on the free space.
cross-cluster	In GaussDB(DWS), users can access data in other DBMS through foreign tables or using an Extension Connector. Such access is cross-cluster.
junk tuple	A tuple that is deleted using the DELETE and UPDATE statements. When deleting a tuple, GaussDB(DWS) only marks the tuples that are to be cleared. The Vacuum thread will then periodically clear these junk tuples.

Term	Description
column	An equivalent concept of "field". A database table consists of one or more columns. Together they describe all attributes of a record in the table.
logical node	Multiple logical nodes can be installed on the same node. A logical node is a database instance.
schema	A collection of database objects that define the logical structure, such as tables, views, sequences, stored procedures, synonyms, indexes, clusters, and database links.
schema file	A SQL file that determines the database structure.

P~T

Table 16-4 Terms P to T

Term	Description
Page	Minimum memory unit for row storage in the GaussDB(DWS) relational object structure. The default size of a page is 8 KB. By default, the page size is determined during database initialization and cannot be dynamically changed.
PostgreSQL	An open-source DBMS developed by volunteers all over the world. PostgreSQL is not controlled by any companies or individuals. Its source code can be used for free.
Postgres-XC	Postgres-XC is an open source PostgreSQL cluster to provide write-scalable, synchronous, multi-master PostgreSQL cluster solution.
Postmaster	A thread started when the database service is started. It listens to connection requests from other nodes in the cluster or from clients.
	After receiving and accepting a connection request from the standby server, the primary server creates a WAL Sender thread to interact with the standby server.
RHEL	Red Hat Enterprise Linux (RHEL)
redo log	A log that contains information required for performing an operation again in a database. If a database is faulty, redo logs can be used to restore the database to its original state.
SCTP	The Stream Control Transmission Protocol (SCTP) is a transport-layer protocol defined by Internet Engineering Task Force (IETF) in 2000. The protocol ensures the reliability of datagram transport based on unreliable service transmission protocols by transferring SCN narrowband signaling over IP network.

Term	Description		
savepoint	A savepoint marks the end of a sub-transaction (also known as a nested transaction) in a relational DBMS. The process of a long transaction can be divided into several parts. After a part is successfully executed, a savepoint will be created. If later execution fails, the transaction will be rolled back to the savepoint instead of being totally rolled back. This is helpful for recovering database applications from complex errors. If an error occurs in a multi-statement transaction, the application can recover by rolling back to the save point without terminating the entire transaction.		
session	A task created by a database for a connection when an application attempts to connect to the database. Sessions are managed by the session manager. They execute initial tasks to perform all user operations.		
shared- nothing architecture	A distributed computing architecture, in which none of the nodes share CPUs or storage resources. This architecture has good scalability.		
SIMD	Single Instruction, Multiple Data (SIMD) is a parallel computing technique where a single instruction operates on multiple data elements simultaneously. It significantly improves the performance of compute-intensive tasks.		
SLES	SUSE Linux Enterprise Server (SLES) is an enterprise Linux OS provided by SUSE.		
SMP	Symmetric multiprocessing (SMP) lets multiple CPUs run on a computer and share the same memory and bus. To ensure an SMP system achieves high performance, an OS must support multi-tasking and multi-thread processing. In databases, SMP means to concurrently execute queries using the multi-thread technology, efficiently using all CPU resources and improving query performance.		
SQL	Structure Query Language (SQL) is a standard database query language. It consists of DDL, DML, and DCL.		
SSL	Secure Socket Layer (SSL) is a network security protocol introduced by Netscape. It is a security protocol based on the TCP and IP communications protocols and uses the public key technology. SSL supports a wide range of networks and provides three basic security services, all of which use the public key technology. SSL ensures the security of service communication through the network by establishing a secure connection between the client and server and then sending data through this connection.		
convergence ratio	Downlink to uplink bandwidth ratio of a switch. A high convergence ratio indicates a highly converged traffic environment and severe packet loss.		

Term	Description		
ТСР	Transmission Control Protocol (TCP) sends and receives data through the IP protocol. It splits data into packets for sending, and checks and reassembles received package to obtain original information. TCP is a connection-oriented, reliable protocol that ensures information correctness in transmission.		
trace	A way of logging to record information about the way a program is executed. This information is typically used by programmers for debugging purposes. System administrators and technical support can diagnose common problems by using software monitoring tools and based on this information.		
escape character	In a database, an escape character is a special character used to escape other characters. It enables the database system to identify and correctly process characters that have special meanings and may conflict with SQL syntax (such as quotation marks and special characters) as common characters. Common escape characters include backslashes (\) and double quotation marks (").		
full backup	Backup of the entire database cluster.		
full synchronizati on	A data synchronization mechanism specified in the GaussDB(DWS) HA solution. Used to synchronize all data from the primary server to a standby server.		
log file	A file to which a computer system writes a record of its activities.		
transaction	A logical unit of work performed within a DBMS against a database. A transaction consists of a limited database operation sequence, and must have ACID features.		
data	A representation of facts or directives for manual or automatic communication, explanation, or processing. Data includes constants, variables, arrays, and strings.		
data redistribution	A process whereby a data table is redistributed among nodes after users change the data distribution mode.		
data distribution	A mode in which table data is split and stored on each database instance in a distributed system. Table data can be distributed in hash, replication, or random mode. In hash mode, a hash value is calculated based on the value of a specified column in a tuple, and then the target storage location of the tuple is determined based on the mapping between nodes and hash values. In replication mode, tuples are replicated to all nodes. In random mode, data is randomly distributed to the nodes.		
data partitioning	A division of a logical database or its constituent elements into multiple parts (partitions) whose data does not overlap based on specified ranges. Data is mapped to storage locations based on the value ranges of specific columns in a tuple.		

Term	Description			
database	A collection of data that is stored together and can be accessed, managed, and updated. Data in a view in the database can be classified into the following types: numerals, full text, digits, and images.			
DB instance	A database instance consists of a process in GaussDB(DWS) and files controlled by the process. GaussDB(DWS) installs multiple database instances on one physical node. GTM, CM, CN, and DN installed on cluster nodes are all database instances. A database instance is also called a logical node.			
database HA	GaussDB(DWS) provides a highly reliable HA solution. Every logical node in GaussDB(DWS) is identified as a primary or standby node. Only one GaussDB(DWS) node is identified as primary at a time. When the HA system is deployed for the first time, the primary server synchronizes all data from each standby server (full synchronization). The HA system then synchronizes only data that is new or has been modified from each standby server (incremental synchronization). When the HA system is running, the primary server can receive data read and write operation requests and the standby servers only synchronize logs.			
database file	A binary file that stores user data and the data inside the database system.			
data flow operator	An operator that exchanges data among query fragments. By their input/output relationships, data flows can be categorized into Gather flows, Broadcast flows, and Redistribution flows. Gather combines multiple query fragments of data into one. Broadcast forwards the data of one query fragment to multiple query fragments. Redistribution reorganizes the data of multiple query fragments and then redistributes the reorganized data to multiple query fragments.			
data dictionary	A reserved table within a database which is used to store information about the database itself. The information includes database design information, stored procedure information, user rights, user statistics, database process information, database increase statistics, and database performance statistics.			
deadlock	Unresolved contention for the use of resources.			
index	An ordered data structure in the database management system. An index accelerates querying and the updating of data in database tables.			
statistics	Information that is automatically collected by databases, including table-level information (number of tuples and number of pages) and column-level information (column value range distribution histogram). Statistics in databases are used to estimate the cost of execution plans to find the plan with the lowest cost.			

Term	Description	
stop word	In computing, stop words are words which are filtered out before or after processing of natural language data (text), saving storage space and improving search efficiency.	

U~Z

Table 16-5 Terms U to Z

Term	Description		
unlogged table	Unlogged tables are a special type of tables that do not record write-ahead logs (WALs) for data operations. These table can significantly improve write performance in some scenarios. However, they are automatically truncated upon conflicts, OS restart, forcible restart, power-off, or unexpected shutdown, which may cause data loss. The content of an unlogged table is not replicated to standby servers. Any indexes created on an unlogged table are not automatically logged as well. If the UNLOGGED parameter is specified in the CREATE TABLE syntax, an unlogged table is created.		
Vacuum	A thread that is periodically started up by a database to clear junk tuples. Multiple Vacuum threads can be started concurrently by setting a parameter.		
V2 table	A V2 table refers to a table whose colversion defined in the CREATE TABLE syntax is 2.0 during table creation, indicating that each column of the column-store table is combined and stored in a file named relfilenode.C1.0 , and data is stored on the local disk. For storage-compute coupled clusters, if colversion is not specified, the created column-store table is a V2 table by default.		
V3 table	A V3 table refers to a table whose colversion defined in the CREATE TABLE syntax is 3.0 during table creation, indicating that each column of the column-store table is stored in a file named C1_field.0 , and data is stored in the OBS file system. For storage-compute coupled clusters, if colversion is not specified, the created column-store table is a V3 table by default.		
verbose	The VERBOSE option specifies the information to be displayed.		
WAL	Write-ahead logging (WAL) is a standard method for logging a transaction. Corresponding logs must be written into a permanent device before a data file (carrier for a table and index) is modified.		

Term	Description			
WAL Receiver	A thread created by the standby server during database duplication. The thread is used to receive data and commands from the primary server and to tell the primary server that the data and commands have been acknowledged. Only one WAL receiver thread can run on one standby server.			
WAL Sender	Name of a thread created by the primary node after the primary node receives a connection request from the standby node during database replication. This thread is used to send data and commands to standby servers and to receive responses from the standby servers. Multiple WAL Sender threads may run on one primary server. Each WAL Sender thread corresponds to a connection request initiated by a standby server.			
WAL Writer	A thread for writing redo logs that are created when a database is started. This thread is used to write logs in the memory to a permanent device, such as a disk.			
WLM	The WorkLoad Manager (WLM) is a module for controlling and allocating system resources in GaussDB(DWS).			
Xlog	A transaction log. A logical node can have only one Xlog file.			
xDR	X detailed record. It refers to detailed records on the user and signaling plans and can be categorized into charging data records (CDRs), user flow data records (UFDRs), transaction detail records (TDRs), and data records (SDRs).			
network backup	Network backup provides a comprehensive and flexible data protection solution to Windows, UNIX, and Linux platforms. Network backup can back up, archive, and restore files, folders, directories, volumes, and partitions on a computer.			
Predicate	A predicate is a logical expression used to filter or limit data. It is essentially a conditional expression that returns a Boolean value (TRUE, FALSE, or NULL). Predicates are used in WHERE, CHECK, and JOIN clauses. For example, salary > 5000 in WHERE salary > 5000, is a predicate.			
predicate column	A predicate column is a column referenced in a predicate expression. For example, salary in salary > 5000 is the predicate column. If the predicate column in a query contains an index, the query optimizer may use the index to accelerate data filtering. This is called "index scan". GaussDB(DWS) collects statistics on predicate columns (supported by clusters of version 9.1.0.100 and later) to help generate better query plans.			
physical node	A physical machine or device.			
materialized view	A materialized view is a special database object. It pre-computes complex query results and stores them in the database to accelerate queries.			

Term	Description		
system catalog	A table storing meta information about the database. The meta information includes user tables, indexes, columns, functions, and the data types in a database.		
pushdown	GaussDB(DWS) is a distributed database, where CN can send a query plan to multiple DNs for parallel execution. This CN behavior is called pushdown. It achieves better query performance than extracting data to CN for query.		
compression	Data compression, source coding, or bit-rate reduction involves encoding information that uses fewer bits than the original representation. Compression can be either lossy or lossless. Lossless compression reduces bits by identifying and eliminating statistical redundancy. No information is lost in lossless compression. Lossy compression reduces bits by identifying and removing unnecessary or unimportant information. The process of reducing the size of a data file is commonly referred as data compression, although its formal name is source coding (coding done at the source of the data, before it is stored or transmitted).		
consistency	One of the ACID features of database transactions. Consistency is a database status. In such a status, data in the database must comply with integrity constraints.		
metadata	Data that provides information about other data. Metadata describes the source, size, format, or other characteristics of data. In database columns, metadata explains the content of a data warehouse.		
atomicity	One of the ACID features of database transactions. Atomicity means that a transaction is composed of an indivisible unit of work. All operations performed in a transaction must either be committed or uncommitted. If an error occurs during transaction execution, the transaction is rolled back to the state when it was not committed.		
Zhparser	Zhparser is an extension for full-text search of Chinese. It provides Chinese word segmentation, enabling GaussDB(DWS) to process Chinese text search and analysis requests more efficiently.		
online scale- out	Online scale-out means that data can be saved to the database and query services are not interrupted during redistribution in GaussDB(DWS).		
dirty page	A page that has been modified and is not written to a permanent device.		
incremental backup	Incremental backup stores all files changed since the last valid backup.		

Term	Description		
incremental synchronizati on	A data synchronization mechanism in the GaussDB(DWS) HA solution. Only data modified since the last synchronization is synchronized to the standby server.		
Host	A node that receives data read and write operations in the GaussDB(DWS) HA system and works with all standby servers. At any time, only one node in the HA system is identified as the primary server.		
thesaurus	Standardized words or phrases that express document themes and are used for indexing and retrieval.		
dump file	A specific type of the trace file. A dump is typically a one-time output of diagnostic data in response to an event, whereas a trace tends to be continuous output of diagnostic data.		
resource pool	Resource pools used for allocating resources in GaussDB(DWS). By binding a user to a resource pool, you can limit the priority of the jobs executed by the user and resources available to the jobs.		
tenant	A database service user who runs services using allocated computing (CPU, memory, and I/O) and storage resources. Service level agreements (SLAs) are met through resource management and isolation.		
minimum restoration point	A method used by GaussDB(DWS) to ensure data consistency. During startup, GaussDB(DWS) checks consistency between the latest WAL logs and the minimum restoration point. If the record location of the minimum restoration point is greater than that of the latest WAL logs, the database fails to start.		

1 Hybrid Data Warehouse

17.1 Introduction to Hybrid Data Warehouse

Hybrid: A hybrid data warehouse offers both large-scale data query and analysis capabilities, as well as low-cost, high-concurrency, high-performance, and low-latency transaction processing capabilities.

□ NOTE

- To use hybrid data warehouse capabilities, choose the storage-compute coupled architecture when you create a GaussDB(DWS) cluster on the console and ensure the vCPU to memory ratio is 1:4 when setting up cloud disk flavors. For more information, see Data Warehouse Flavors.
- When setting up a GaussDB(DWS) cluster, make sure to have a vCPU to memory ratio of 1:8 for standard data warehouses and a ratio of 1:4 for hybrid data warehouses. You can distinguish a standard data warehouse from a real-time data warehouse by comparing their vCPU to memory ratios.

A hybrid data warehouse needs to work with data sources, such as upstream databases or applications, to insert, upsert, and update data in real time. The data warehouse should also be able to query data shortly after it was imported.

Currently, the existing row-store and column-store tables in a conventional GaussDB(DWS) data warehouse cannot meet real-time data import and query requirements. Row-store tables have strong real-time import capabilities and support highly concurrent updates, but their disk usage is high and query efficiency is low. Column-store tables have high data compression ratio and good OLAP query performance, but do not support concurrent updates. Concurrent import will cause severe lock conflicts.

To solve these problems, we use column storage to reduce the disk usage, support highly concurrency updates, and improve query speed. GaussDB(DWS) hybrid data warehouses use HStore tables to achieve high performance during real-time data import and query, and have the transaction processing capabilities required for traditional OLTP scenarios.

The HStore tables uniquely support single and small-batch real-time IUD operations, as well as regular large-batch import. Data can be queried immediately after being imported. You can deduplicate traditional indexes (such

as primary keys) and accelerate point queries. You can further accelerate OLAP queries through partitioning, multi-dimensional dictionaries, and partial sorting. Strong data consistency can be ensured for transactions with heavy workloads, such as TPC-C.

□ NOTE

- Only clusters 8.2.0.100 and later support the HStore tables of the hybrid data warehouse.
- HStore Opt tables are recommended for 9.1.0 and later versions. HStore tables can be replaced by HStore_opt tables for better performance, except in scenarios requiring high performance without micro-batch updates.
- The hybrid data warehouse is used for both production and analysis. It is applicable to
 hybrid transaction and analysis scenarios. It can be deployed in single-node or cluster
 mode. For how to create a hybrid data warehouse, see Creating a GaussDB(DWS)
 Storage-Compute Coupled Cluster.
- Hot and cold data management is supported for HStore tables. For details, see Best Practices of Hot and Cold Data Management. This function is supported only by cluster versions 8.2.0.101 and later.
- HStore is a table type designed for the hybrid data warehouse and is irrelevant to the SQL parameter hstore.

Differences from Standard Data Warehouses

Hybrid data warehouses and standard data warehouses are two types of GaussDB(DWS) data warehouses with different specifications and usage. For details, see **Table 17-1**.

Table 17-1 Comparison between hybrid and standard data warehouses

Туре	Standard Data Warehouse (Compute-Storage Coupled Architecture with 1:8 vCPU to Memory Ratio)	Hybrid Data Warehouse (Compute-Storage Coupled Architecture with 1:4 vCPU to Memory Ratio)
Application scenario	Converged data analysis using OLAP. It is used in sectors such as finance, government and enterprise, e-commerce, and energy.	Real-time data import + Hybrid analysis. Real- time upstream data import + Real-time query after data import. It is mainly used in scenarios that have high requirements on real-time data import, such as e- commerce and finance.

Туре	Standard Data Warehouse (Compute-Storage Coupled Architecture with 1:8 vCPU to Memory Ratio)	Hybrid Data Warehouse (Compute-Storage Coupled Architecture with 1:4 vCPU to Memory Ratio)
Advantage	It is cost-effective and widely used. Cost effective, both hot and cold data analysis supported, elastic storage and compute capacities. Hybrid load, high daimport performance. It achieves high querefficiency and high data compression rate that are equivalent to those of column storage. It can also process transactions traditional OLTP scenarios.	
Features Excellent performance in interactive analysis and offline processing of mass data, as well as complex data mining.		It supports highly concurrent update operations on massive amounts of data and can achieve high query efficiency. It achieves high performance when processing a large amount of data in scenarios like high-concurrency import and latency-sensitive queries.
SQL syntax	Highly compatible with SQL syntax	Compatible with column-store syntax
GUC parameter	You can configure a wide variety of GUC parameters to tailor your data warehouse environment.	It is compatible with standard data warehouse GUC parameters and supports hybrid data warehouse tuning parameters.

Technical Highlights

• Transaction consistency

Data can be retrieved for queries immediately after being inserted or updated. After concurrent updates, data is strongly consistent, and there will not be incorrect results caused by wrong update sequence.

High query performance
 In complex OLAP queries, such as multi-table correlation, the data warehouse achieves high performance through comprehensive distributed query plans

and distributed executors. It also supports complex subqueries and stored procedures.

Quick import

There will not be lock conflicts on column-store CUs. High-concurrency update and import operations are supported. The concurrent update performance can be over 100 times higher than before in general scenarios.

• High compression

Column storage can achieve a high compression ratio. Data is stored in the column-store primary table through MERGE can be compressed to greatly reduce disk usage and I/O.

• Query acceleration

You can deduplicate traditional indexes (such as primary keys) and accelerate point queries. You can further accelerate OLAP queries through partitioning, multi-dimensional dictionaries, and partial sorting.

Comparison Between Row-store, Column-store, and HStore Tables

Table 17-2 Comparison between row-store, column-store, and HStore tables

Table Type	Row-Store	Column-Store	HStore
Data storage mode	The attributes of a tuple are stored nearby.	The values of an attribute are stored nearby in the unit of CU.	Data is stored in the column-store primary tables as CUs. Updated columns and data inserted in small batches is serialized and then stored in a newly designed delta table.
Data write	Row-store compression has not been put into commercial use. Data is stored as it is, occupying a large amount of disk space.	In row storage, data with the same attribute value types is easy to compress. Data write consumes much fewer I/O resources and less disk space.	Data inserted in batches is directly written to CUs, which are as easy to compress as column storage. Updated columns and data inserted in small batches are serialized and then compressed. They will also be periodically merged to primary table CUs.

Table Type	Row-Store	Column-Store	HStore
Data update	Data is updated by row, avoiding CU lock conflicts. The performance of concurrent updates (UPDATE/UPSERT/DELETE) is high.	The entire CU needs to be locked even if only one record in it is updated. Generally, concurrent updates (UPDATE/UPSERT/DELETE) are not supported.	CU lock conflicts can be avoided. The performance of concurrent updates (UPDATE/UPSERT/DELETE) is higher than 60% of the rowstore update performance.
Data read	Data is read by row. An entire row needs to be retrieved even if only one column in it needs to be accessed. The query performance is low.	When data is read by column, only the CU of a column needs to be accessed. CUs can be easily compressed, occupying less I/O resources, and achieve high read performance.	Data in a column- store primary table is read by column. Updated columns and data inserted in small batches are deserialized and then retrieved. After data is merged to the primary table, the data can be read as easily as that in column storage.
Advantage	The concurrent update performance is high.	The query performance is high, and the disk space usage is small.	The concurrent update performance is high. After data merges, the query and compression performance are the same as those of column storage.
Disadvantag e	A large amount of disk space is occupied, and the query performance is low.	Generally, concurrent updates are not supported.	A background permanent thread is required to clear unnecessary HStore table data after merge. Data is merged to the primary table CUs and then cleared. This operation is irrelevant to the SQL syntax MERGE .

Table Type	Row-Store	Column-Store	HStore
Application scenario	 OLTP transactions with frequent update and deletion operations Point queries (simple queries that are based on indexes and return a small amount of data) 	 OLAP query and analysis A large volume of data is imported, and is rarely updated or deleted after the import. 	 Data is concurrently imported to the database in real time. High-concurrency update and import; and high- performance query

17.2 Support and Constraints

A hybrid data warehouse is compatible with all column-store syntax.

Table 17-3 Supported syntax

Syntax	Supported
CREATE TABLE	Yes
CREATE TABLE LIKE	Yes
DROP TABLE	Yes
INSERT	Yes
COPY	Yes
SELECT	Yes
TRUNCATE	Yes
EXPLAIN	Yes
ANALYZE	Yes
VACUUM	Yes
ALTER TABLE DROP PARTITION	Yes
ALTER TABLE ADD PARTITION	Yes
ALTER TABLE SET WITH OPTION	Yes
ALTER TABLE DROP COLUMN	Yes
ALTER TABLE ADD COLUMN	Yes

Syntax	Supported
ALTER TABLE ADD NODELIST	Yes
ALTER TABLE CHANGE OWNER	Yes
ALTER TABLE RENAME COLUMN	Yes
ALTER TABLE TRUNCATE PARTITION	Yes
CREATE INDEX	Yes
DROP INDEX	Yes
DELETE	Yes
Other ALTER TABLE syntax	Yes
ALTER INDEX	Yes
MERGE	Yes
SELECT INTO	Yes
UPDATE	Yes
CREATE TABLE AS	Yes

Constraints

- 1. To use HStore tables, use the following parameter settings, or the performance of HStore tables will deteriorate significantly:
 - autovacuum_max_workers_hstore=3, autovacuum_max_workers=6, and autovacuum=true
- 2. In version 8.2.1 and later, you can now clear the dirty data from column-store indexes. This is especially beneficial when dealing with frequent data updates and imports into the database. By efficiently managing the index space, it improves both the import and query performance.
- 3. When using HStore asynchronous sorting, pay attention to the following:
 - DML operations on certain data may be blocked during asynchronous sorting. The maximum blocking granularity is the row threshold for asynchronous sorting. This function is not recommended for frequent DML operations.
 - Automatic asynchronous sorting and column-store VACUUM cannot be used together. If the autovacuum process meets the conditions for column-store VACUUM, asynchronous sorting is skipped and will wait for the next trigger. In some cases, column-store VACUUM may be continuously triggered due to a high volume of DML operations, which means asynchronous sorting will never be triggered.

17.3 Hybrid Data Warehouse Syntax

17.3.1 CREATE TABLE

Function

Create an HStore table in the current database. The table will be owned by the user who created it.

In a hybrid data warehouse, you can use DDL statements to create HStore tables. To create an HStore table, set **enable_hstore** to **on** and set **orientation** to **column**.

To enhance performance, GaussDB(DWS) 9.1.0 and later versions have optimized HStore tables and kept the old ones for compatibility purposes. The optimized tables are known as HStore Opt tables. HStore tables can be replaced by HStore Opt tables for better performance, except in scenarios requiring high performance without micro-batch updates.

You cannot convert a created HStore table to an HStore_opt table. Create an HStore Opt table if you want to use it.

◯ NOTE

- To use hybrid data warehouse capabilities, choose the storage-compute coupled architecture when you create a GaussDB(DWS) cluster on the console and ensure the vCPU to memory ratio is 1:4 when setting up cloud disk flavors. For more information, see Data Warehouse Flavors.
- When setting up a GaussDB(DWS) cluster, make sure to have a vCPU to memory ratio
 of 1:8 for standard data warehouses and a ratio of 1:4 for hybrid data warehouses. You
 can distinguish a standard data warehouse from a real-time data warehouse by
 comparing their vCPU to memory ratios.

Precautions

- When creating an HStore table, ensure that the database GUC parameter settings meet the following requirements:
 - autovacuum is set to on.
 - The value of autovacuum max workers hstore is greater than 0.
 - The value of autovacuum_max_workers is greater than that of autovacuum max workers hstore.
- To create an HStore table, you must have the **USAGE** permission on schema cstore.
- The table-level parameters enable_delta and enable_hstore_opt cannot be enabled at the same time. The parameter enable_delta is used to enable delta for common column-store tables and conflicts with enable_hstore_opt.
- Each HStore table is bound to a delta table. The OID of the delta table is recorded in the **reldeltaidx** field in **pg_class**. (The **reldelta** field is used by the delta table of the column-store table).

Syntax

Differences Between Delta Tables

Table 17-4 Differences between the delta tables of HStore and column-store tables

Туре	Column-Store Delta Table	HStore Delta Table	HStore Opt Delta Table
Table struct ure	Same as that defined for the column-store primary table.	Different from that defined for the primary table.	Different from the definitions of the primary table and but same as the definitions of the HStore table.
Functi onality	Used to temporarily store a small batch of inserted data. After the data size reaches the threshold, the data will be merged to the primary table. In this way, data will not be directly inserted to the primary table or generate a large number of small CUs.	Persistently stores UPDATE, DELETE, and INSERT information. It is used to restore the memory structure that manages concurrent updates, such as the memory update chain, in the case of a fault.	Persistently stores UPDATE, DELETE, and INSERT information. It is used to restore the memory structure that manages concurrent updates, such as the memory update chain, in the case of a fault. It is further optimized compared with HStore.
Weak ness	If data is not merged in a timely manner, the delta table will grow large and affect query performance. In addition, the table cannot solve lock conflicts during concurrent updates.	The merge operation depends on the background AUTOVACUUM.	The merge operation depends on the background AUTOVACUUM.

Туре	Column-Store Delta Table	HStore Delta Table	HStore Opt Delta Table
Specification differences	Concurrent requests in the same CU are not supported. It is applicable to the scenario where there are not many concurrent updates.	 Insertion and update restrictions: MERGE INTO does not support concurrent updates of the same row or repeated updates of the same key. Concurrent UPDATE or DELETE operations on the same row are not supported. Otherwise, an error is reported. Index and query restrictions: Indexes do not support array condition filtering, IN expression filtering, partial indexes, or expression indexes. Indexes cannot be invalidated. Table structure and operation restrictions: Ensure that the tables to be exchanged are HStore tables during partition exchange or relfilenode operations. The distribution column cannot be modified 	 Insertion and update restrictions: MERGE INTO does not support concurrent updates of the same row or repeated updates of the same key. Concurrent updates or deletions of the same row is not supported. hstore_opt does not support cross-partition upserts. Index and query restrictions: Bitmap indexes are supported. Global dictionaries are supported. bitmap_column s must be specified during table creation and cannot be modified after being set. The opt version does not support transparent parameter transmission during SMP streaming. In multi-table join queries that require partition pruning, avoid using replicated tables or setting query_dop.

Туре	Column-Store Delta Table	HStore Delta Table	HStore Opt Delta Table	
		using the UPDATE command. You are not advised to modify the partition column using the UPDATE command. (No error is reported, but the performance is poor.)	 Table structure and operation restrictions: Distribution columns and partition columns cannot be modified using UPDATE. The enable_hstore_o pt attribute must be set when the table is created and cannot be changed after being set. 	
Data import sugges tions	it is recommended involving micro-bat no data updates, yo	. For optimal data import, query performance, and space utilization, it is recommended to choose the HStore Opt table. In scenarios involving micro-batch copying with high performance demands and no data updates, you can choose the HStore table.		
	The performanc	 Similarities between HStore and HStore Opt tables: The performance of importing data using UPDATE is poor. You 		
	When using DEI	 are advised to use UPSERT to import data. When using DELETE to import data, use index scanning. The JDBC batch method is recommended. 		
	the data volume	 Use MERGE INTO to import data records to the database when the data volume exceeds 1 million per DN and there is no concurrent data. 		
	Do not modify of	or add data in cold partiti	ons.	

Туре	Column-Store Delta Table	HStore Delta Table	HStore Opt Delta Table			
Point query sugges tions	2. Similarities between HStore and HStore Opt tables: Create a level-2 partition on the column where the equal-version filter condition is most frequently used and distinct values are evenly distributed.					
	 Accelerating ind effect. You are a If the data type 	 Suggestions on using HStore tables for point queries: Accelerating indexes other than primary keys may have poor effect. You are advised not to enable index acceleration. If the data type is numeric or strings less than 16 bytes, Turbo acceleration is recommended. 				
	columns involve query, use the C you are advised five index colum • For all string col indexes can be s	filter columns not in leve d in the filter criteria are B-tree index. If the colum to use the GIN index. Do	basically fixed in the ins change continuously, not select more than at filtering, bitmap stion. The number of			
	 Specify columns partition column If the number of index scanning of this case, you ar 	that can be filtered by times. If returned data records expressions and the significantly enhance advised to use the GUC to test the performance	me range as the sceeds 100,000 per DN, ance performance. In parameter			

Parameters

• IF NOT EXISTS

If **IF NOT EXISTS** is specified, a table will be created if there is no table using the specified name. If there is already a table using the specified name, no error will be reported. A message will be displayed indicating that the table already exists, and the database will skip table creation.

table name

Specifies the name of the table to be created.

The table name can contain a maximum of 63 characters, including letters, digits, underscores (_), dollar signs (\$), and number signs (#). It must start with a letter or underscore (_).

column_name

Specifies the name of a column to be created in the new table.

The column name can contain a maximum of 63 characters, including letters, digits, underscores (_), dollar signs (\$), and number signs (#). It must start with a letter or underscore (_).

data_type

Specifies the data type of the column.

• LIKE source_table [like_option ...]

Specifies a table from which the new table automatically copies all column names and their data types.

The new table and the original table are decoupled after creation is complete. Changes to the original table will not be applied to the new table, and scans on the original table will not be performed on the data of the new table.

Columns copied by **LIKE** are not merged with the same name. If the same name is specified explicitly or in another **LIKE** clause, an error will be reported.

HStore tables can be inherited only from HStore tables.

WITH ({ storage_parameter = value } [, ...])

Specifies an optional storage parameter for a table.

ORIENTATION

Specifies the storage mode (time series, row-store, or column-store) of table data. This parameter cannot be modified once it is set. For HStore tables, use the column storage mode and set **enable_hstore** to **on**.

Options:

- TIMESERIES indicates that the data is stored in time series.
- COLUMN indicates that the data is stored in columns.
- ROW indicates that table data is stored in rows.

Default value: ROW

COMPRESSION

Specifies the compression level of the table data. It determines the compression ratio and time. Generally, a higher compression level indicates a higher compression ratio and a longer compression time, and vice versa. The actual compression ratio depends on the distribution characteristics of loading table data.

Options:

This parameter is available only for HStore tables and column-store tables. Value: **LOW** (default value), **MIDDLE**, or **HIGH**.

COMPRESSLEVEL

Specifies table data compression rate and duration at the same compression level. This divides a compression level into sub-levels, providing you with more choices for compression ratio and duration. As the value becomes greater, the compression rate becomes higher and duration longer at the same compression level. The parameter is only valid for time series tables and column-store tables.

Value range: 0 to 3
Default value: **0**MAX BATCHROW

Specifies the maximum number of rows in a storage unit during data loading. The parameter is only valid for time series tables and column-store tables.

Value range: 10000 to 60000

Default value: **60000**- PARTIAL_CLUSTER_ROWS

Specifies the number of records to be partially clustered for storage during data loading. The parameter is only valid for time series tables and column-store tables.

Value range: 600000 to 2147483647

Default value: 4,200,000

enable delta

Specifies whether to enable delta tables in column-store tables. This parameter cannot be enabled for HStore tables.

Default value: **off** enable hstore

Specifies whether to create a table as an HStore table (based on column-store tables). The parameter is only valid for column-store tables. This parameter is supported by version 8.2.0.100 or later clusters.

Default value: off

If this parameter is enabled, the following GUC parameters must be set to ensure that HStore tables are cleared.

autovacuum=true, autovacuum_max_workers=6, autovacuum_max_workers_hstore=3.

enable disaster cstore

Specifies whether fine-grained DR will be enabled for column-store tables. This parameter only takes effect on column-store tables whose COLVERSION is 2.0 and cannot be set to **on** if **enable_hstore** is **on**. This parameter is supported by version 8.2.0.100 or later clusters.

Default value: off

<u>A</u> CAUTION

Before enabling this function, set the GUC parameter **enable_metadata_tracking** to **on**. Otherwise, fine-grained DR may fail to be enabled.

SUB PARTITION COUNT

Specifies the number of level-2 partitions. This parameter specifies the number of level-2 partitions during data import. This parameter is configured during table creation and cannot be modified after table creation. You are not advised to set the default value, which may affect the import and query performance.

Value range: 1 to 1024

Default value: 32

DELTAROW_THRESHOLD

Specifies the maximum number of rows (SUB_PARTITION_COUNT x DELTAROW_THRESHOLD) to be imported to the delta table.

Value range: 0 to 60000 Default value: **60000**

- COLVERSION

Specifies the version of the storage format. HStore tables support only version 2.0, and **enable_hstore_opt** tables support versions 2.0 and 3.0.

Options:

1.0: Each column in a column-store table is stored in a separate file. The file name is **relfilenode.C1.0**, **relfilenode.C2.0**, **relfilenode.C3.0**, or similar.

2.0: All columns of a column-store table are combined and stored in a file. The file is named **relfilenode.C1.0**.

Default value: **2.0** enable binlog

Specifies whether to enable the binlog function for the HStore Opt table. This parameter is supported only by clusters of version 9.1.0 or later.

Value range: on and off

Default value: off

enable_binlog_timestamp

Determines whether to enable the binlog function with timestamps for HStore Opt tables. This parameter and **enable_binlog** cannot be enabled at the same time. Only clusters of 9.1.0.200 and later versions support this parameter.

Value range: on and off

Default value: **off**DISTRIBUTE BY

Specifies how the table is distributed or replicated between DNs.

Options:

HASH (column_name): Each row of the table will be placed into all the DNs based on the hash value of the specified column.

TO { GROUP groupname | NODE (nodename [, ...]) }

TO GROUP specifies the Node Group in which the table is created. Currently, it cannot be used for HDFS tables. **TO NODE** is used for internal scale-out tools.

PARTITION BY

Specifies the initial partition of an HStore table.

secondary part column

Specifies the name of a level-2 partition column in a column-store table. Only one column can be specified as the level-2 partition column. This parameter applies only to HStore Opt column-store tables. This parameter is supported only by clusters of version 9.1.0 and later. V3 tables do not support this parameter and will use hashbucket pruning.

■ NOTE

- The column specified as a level-2 partition column cannot be deleted or modified.
- The level-2 partition column can be specified only when a table is created.
 After a table is created, the level-2 partition column cannot be modified.
- You are not advised to specify a distribution column as a level-2 partition column.
- The level-2 partition column determines how the table is logically split into hash partitions on DNs, which enhances the query performance for that column.

secondary_part_num

Specifies the number of level-2 partitions in a column-store table. This parameter applies only to HStore Opt column-store tables. This parameter is supported only by clusters of version 9.1.0 and later. V3 tables do not support this parameter and will use hashbucket pruning.

Value range: 1 to 32 Default value: **8**

- MOTE
 - This parameter can be specified only when secondary_part_column is specified.
 - The number of level-2 partitions can be specified only when a table is created and cannot be modified after the table is created.
 - You are not advised to change the default value, which may affect the import and query performance.

Example

Create a simple HStore Opt table.

```
CREATE TABLE warehouse_t1
                        INTEGER
  W_WAREHOUSE_SK
                                        NOT NULL,
  W WAREHOUSE ID
                        CHAR(16)
                                        NOT NULL.
  W_WAREHOUSE_NAME
                          VARCHAR(20)
  W WAREHOUSE SQ FT
                          INTEGER
  W_STREET_NUMBER
                        CHAR(10)
  W_STREET_NAME
                       VARCHAR(60)
                      CHAR(15)
  W_STREET_TYPE
  W_SUITE_NUMBER
                        CHAR(10)
  W CITY
                   VARCHAR(60)
  W_COUNTY
                     VARCHAR(30)
  W STATE
                   CHAR(2)
  W_ZIP
                  CHAR(10)
  W_COUNTRY
                      VARCHAR(20)
  W_GMT_OFFSET
                      DECIMAL(5,2)
)WITH(ORIENTATION=COLUMN, ENABLE_HSTORE_OPT=ON);
CREATE TABLE warehouse_t2 (LIKE warehouse_t1 INCLUDING ALL);
```

17.3.2 INSERT

Function

Insert one or more rows of data into an HStore table.

- To use hybrid data warehouse capabilities, choose the storage-compute coupled architecture when you create a GaussDB(DWS) cluster on the console and ensure the vCPU to memory ratio is 1:4 when setting up cloud disk flavors. For more information, see Data Warehouse Flavors.
- When setting up a GaussDB(DWS) cluster, make sure to have a vCPU to memory ratio
 of 1:8 for standard data warehouses and a ratio of 1:4 for hybrid data warehouses. You
 can distinguish a standard data warehouse from a real-time data warehouse by
 comparing their vCPU to memory ratios.

Precautions

- If the data to be inserted at a time is greater than or equal to the value of the table-level parameter **DELTAROW_THRESHOLD**, the data is directly inserted into the primary table to generate a compression unit (CU).
- If the data to be inserted is smaller than **DELTAROW_THRESHOLD**, a record of the type **I** will be inserted into the delta table. The data will be serialized and stored in the **values** field of the record.
- CUIDs are allocated to the data in the delta table and the primary table in a unified manner.
- The data inserted into the delta table depends on **AUTOVACUUM** to merge to primary table CUs.

Syntax

Parameters

table_name

Specifies the name of the target table.

Value range: an existing table name

AS

Specifies an alias for the target table *table_name*. *alias* indicates the alias name.

column_name

Specifies the name of a column in a table.

query

Specifies a query statement (**SELECT** statement) that uses the query result as the inserted data.

Example

Create the reason t1 table.

```
-- Create the reason_t1 table.

-- Create the reason_t1 table.

CREATE TABLE reason_t1

(

TABLE_SK INTEGER ,

TABLE_ID VARCHAR(20) ,

TABLE_NA VARCHAR(20)

)WITH(ORIENTATION=COLUMN, ENABLE HSTORE OPT=ON);
```

17.3.3 **DELETE**

Function

Delete data from an HStore Opt table.

◯ NOTE

- To use hybrid data warehouse capabilities, choose the storage-compute coupled architecture when you create a GaussDB(DWS) cluster on the console and ensure the vCPU to memory ratio is 1:4 when setting up cloud disk flavors. For more information, see Data Warehouse Flavors.
- When setting up a GaussDB(DWS) cluster, make sure to have a vCPU to memory ratio of 1:8 for standard data warehouses and a ratio of 1:4 for hybrid data warehouses. You can distinguish a standard data warehouse from a real-time data warehouse by comparing their vCPU to memory ratios.

Precautions

- To delete all the data from a table, you are advised to use the TRUNCATE syntax to improve performance and reduce table bloating.
- If a single record is deleted from an HStore Opt table, a record of the type D
 will be inserted into the delta table. The memory update chain will also be
 updated to manage concurrency.
- If multiple records are deleted from an HStore Opt table at a time, a record of the type **MD** will be inserted for the consecutive deleted records in each CU.
- In concurrent deletion scenarios, operations on the same CU will get queued in traditional column-store tables and result in low performance. For HStore Opt tables, the operations can be concurrently performed, and the deletion performance can be more than 100 times that of column-store tables.
- The syntax is fully compatible with column storage. For more information, see the **UPDATE** syntax.

Syntax

```
DELETE FROM [ ONLY ] table_name [ * ] [ [ AS ] alias ]
    [ USING using_list ]
    [ WHERE condition ]
```

Parameters

ONLY

If **ONLY** is specified, only that table is deleted. If **ONLY** is not specified, this table and all its sub-tables are deleted.

table_name

Specifies the name (optionally schema-qualified) of a target table.

Value range: an existing table name

alias

Specifies the alias for the target table.

Value range: a string. It must comply with the naming convention.

using_list

Specifies the **USING** clause.

condition

Specifies an expression that returns a value of type boolean. Only rows for which this expression returns **true** will be deleted.

Example

Create the reason_t2 table.

```
CREATE TABLE reason_t2
(

TABLE_SK INTEGER ,

TABLE_ID VARCHAR(20) ,

TABLE_NA VARCHAR(20)
)WITH(ORIENTATION=COLUMN, ENABLE_HSTORE_OPT=ON);
INSERT INTO reason_t2 VALUES (1, 'S01', 'StudentA'),(2, 'T01', 'TeacherA'),(3, 'T02', 'TeacherB');
```

Use the WHERE condition for deletion.

```
DELETE FROM reason_t2 WHERE TABLE_SK = 2;
DELETE FROM reason_t2 AS rt2 WHERE rt2.TABLE_SK = 2;
```

Use the **IN** syntax for deletion.

DELETE FROM reason_t2 WHERE TABLE_SK in (1,3);

17.3.4 UPDATE

Function

Update specified data in an HStore Opt table.

□ NOTE

- To use hybrid data warehouse capabilities, choose the storage-compute coupled architecture when you create a GaussDB(DWS) cluster on the console and ensure the vCPU to memory ratio is 1:4 when setting up cloud disk flavors. For more information, see Data Warehouse Flavors.
- When setting up a GaussDB(DWS) cluster, make sure to have a vCPU to memory ratio
 of 1:8 for standard data warehouses and a ratio of 1:4 for hybrid data warehouses. You
 can distinguish a standard data warehouse from a real-time data warehouse by
 comparing their vCPU to memory ratios.

Precautions

- Similar to column storage, the UPDATE operation on an HStore Opt table in the current version involves DELETE and INSERT. You can configure a global GUC parameter to control the lightweight UPDATE of HStore Opt. In the current version, the lightweight UPDATE is disabled by default.
- In concurrent update scenarios, operations on the same CU will cause lock conflicts in traditional column-store tables and result in low performance. For

HStore Opt tables, the operations can be concurrently performed, and the update performance can be more than 100 times that of column-store tables.

Syntax

Parameters

• plan_hint clause

Following the keyword in the /*+ */ format, hints are used to optimize the plan generated by a specified statement block. For details, see Hint-based Tuning.

table name

Name (optionally schema-qualified) of the table to be updated.

Value range: an existing table name

alias

Specifies the alias for the target table.

Value range: a string. It must comply with the naming convention.

expression

Specifies a value assigned to a column or an expression that assigns the value.

DEFAULT

Sets the column to its default value.

The value is **NULL** if no specified default value has been assigned to it.

from_list

A list of table expressions, allowing columns from other tables to appear in the **WHERE** condition and the update expressions. This is similar to the list of tables that can be specified in the **FROM** clause of a **SELECT** statement.

NOTICE

Note that the target table must not appear in the **from_list**, unless you intend a self-join (in which case it must appear with an alias in the **from list**).

condition

An expression that returns a value of type **boolean**. Only rows for which this expression returns **true** are updated.

Example

Create the **reason_update** table.

Insert data to the reason_update table.

INSERT INTO reason_update VALUES (1, 'S01', 'StudentA'),(2, 'T01', 'TeacherA'),(3, 'T02', 'TeacherB');

Perform the UPDATE operation on the **reason_update** table.

UPDATE reason_update SET TABLE_NA = 'TeacherD' where TABLE_SK = 3;

17.3.5 **UPSERT**

Function

HStore Opt is compatible with the **UPSERT** syntax. You can add one or more rows to a table. When a row duplicates an existing primary key or unique key value, the row will be ignored or updated.

∩ NOTE

- To use hybrid data warehouse capabilities, choose the storage-compute coupled architecture when you create a GaussDB(DWS) cluster on the console and ensure the vCPU to memory ratio is 1:4 when setting up cloud disk flavors. For more information, see Data Warehouse Flavors.
- When setting up a GaussDB(DWS) cluster, make sure to have a vCPU to memory ratio
 of 1:8 for standard data warehouses and a ratio of 1:4 for hybrid data warehouses. You
 can distinguish a standard data warehouse from a real-time data warehouse by
 comparing their vCPU to memory ratios.

Precautions

- The **UPSERT** statement of updating data upon conflict can be executed only when the target table contains a primary key or unique index.
- Similar to column storage, an update operation performed using **UPSERT** on an HStore Opt table in the current version involves DELETE and INSERT.
- In concurrent **UPSERT** scenarios, operations on the same CU will cause lock conflicts in traditional column-store tables and result in low performance. For HStore Opt tables, the operations can be concurrently performed, and the upsert performance can be more than 100 times that of column-store tables.

Syntax

Table 17-5 UPSERT syntax

Syntax	Update Data Upon Conflict	Ignore Data Upon Conflict
Syntax 1: No index is specified.	INSERT INTO ON DUPLICATE KEY UPDATE	INSERT IGNORE INSERT INTO ON CONFLICT DO NOTHING

Syntax	Update Data Upon Conflict	Ignore Data Upon Conflict
Syntax 2: The unique key constraint can be inferred from the specified column name or constraint name.	INSERT INTO ON CONFLICT() DO UPDATE SET INSERT INTO ON CONFLICT ON CONSTRAINT con_name DO UPDATE SET	INSERT INTO ON CONFLICT() DO NOTHING INSERT INTO ON CONFLICT ON CONSTRAINT con_name DO NOTHING

Parameters

In syntax 1, no index is specified. The system checks for conflicts on all primary keys or unique indexes. If a conflict exists, the system ignores or updates the corresponding data.

In syntax 2, a specified index is used for conflict check. The primary key or unique index is inferred from the column name, the expression that contains column names, or the constraint name specified in the **ON CONFLICT** clause.

• Unique index inference

Syntax 2 infers the primary key or unique index by specifying the column name or constraint name. You can specify a single column name or multiple column names by using an expression. Example: column1, column2, column3

• **UPDATE** clause

The **UPDATE** clause can use **VALUES(colname)** or **EXCLUDED.colname** to reference inserted data. **EXCLUDED** indicates the rows that should be excluded due to conflicts.

- WHERE clause
 - The WHERE clause is used to determine whether a specified condition is met when data conflict occurs. If yes, update the conflict data. Otherwise, ignore it.
 - Only syntax 2 of Update Data Upon Conflict can specify the WHERE clause, that is, INSERT INTO ON CONFLICT(...) DO UPDATE SET WHERE.

Example

Create table reason upsert and insert data into it.

```
CREATE TABLE reason_upsert
(
    a int primary key,
    b int,
    c int
)WITH(ORIENTATION=COLUMN, ENABLE_HSTORE_OPT=ON);
INSERT INTO reason_upsert VALUES (1, 2, 3);
```

Ignore conflicting data.

INSERT INTO reason_upsert VALUES (1, 4, 5),(2, 6, 7) ON CONFLICT(a) DO NOTHING;

Update conflicting data.

INSERT INTO reason_upsert VALUES (1, 4, 5), (3, 8, 9) ON CONFLICT(a) DO UPDATE SET b = EXCLUDED.b, c = EXCLUDED.c;

17.3.6 MERGE INTO

Function

The **MERGE INTO** statement is used to conditionally match data in a target table with that in a source table. If data matches, **UPDATE** is executed on the target table; if data does not match, **INSERT** is executed. You can use this syntax to run **UPDATE** and **INSERT** at a time for convenience.

- To use hybrid data warehouse capabilities, choose the storage-compute coupled architecture when you create a GaussDB(DWS) cluster on the console and ensure the vCPU to memory ratio is 1:4 when setting up cloud disk flavors. For more information, see Data Warehouse Flavors.
- When setting up a GaussDB(DWS) cluster, make sure to have a vCPU to memory ratio
 of 1:8 for standard data warehouses and a ratio of 1:4 for hybrid data warehouses. You
 can distinguish a standard data warehouse from a real-time data warehouse by
 comparing their vCPU to memory ratios.

Precautions

In concurrent **MERGE INTO** scenarios, the update operations triggered on the same CU will cause lock conflicts in traditional column-store tables and result in low performance. For HStore Opt tables, the operations can be concurrently performed, and the **MERGE INTO** performance can be more than 100 times that of column-store tables.

Syntax

Parameters

INTO clause

Specifies the target table that is being updated or has data being inserted.

table_name

Specifies the name of the target table.

alias

Specifies the alias for the target table.

Value range: a string. It must comply with the naming convention.

USING clause

Specifies the source table, which can be a table, view, or subquery.

ON clause

Specifies the condition used to match data between the source and target tables. Columns in the condition cannot be updated. The **ON** association condition can be **ctid**, **xc_node_id**, or **tableoid**.

• WHEN MATCHED clause

Performs **UPDATE** if data in the source table matches that in the target table based on the condition.

□ NOTE

Distribution columns, system catalogs, and system columns cannot be updated.

• WHEN NOT MATCHED clause

Specifies that the INSERT operation is performed if data in the source table does not match that in the target table based on the condition.

□ NOTE

- An INSERT clause can contain only one VALUES.
- The sequence of WHEN NOT MATCHED and WHEN NOT MATCHED clauses can be exchanged. One of them can be omitted, but they cannot be omitted at the same time.
- Two **WHEN MATCHED** or **WHEN NOT MATCHED** clauses cannot be specified at the same time.

Example

Create a target for **MERGE INTO**.

CREATE TABLE target(a int, b int)WITH(ORIENTATION = COLUMN, ENABLE_HSTORE_OPT = ON); INSERT INTO target VALUES(1, 1),(2, 2);

Create a data source table.

CREATE TABLE source(a int, b int)WITH(ORIENTATION = COLUMN, ENABLE_HSTORE_OPT = ON); INSERT INTO source VALUES(1, 1),(2, 2),(3, 3),(4, 4),(5, 5);

Run the MERGE INTO command.

MERGE INTO target t
USING source s
ON (t.a = s.a)
WHEN MATCHED THEN
UPDATE SET t.b = t.b + 1
WHEN NOT MATCHED THEN
INSERT VALUES (s.a, s.b) WHERE s.b % 2 = 0;

17.3.7 SELECT

Function

Reads data from an HStore Opt table.

□ NOTE

- To use hybrid data warehouse capabilities, choose the storage-compute coupled architecture when you create a GaussDB(DWS) cluster on the console and ensure the vCPU to memory ratio is 1:4 when setting up cloud disk flavors. For more information, see Data Warehouse Flavors.
- When setting up a GaussDB(DWS) cluster, make sure to have a vCPU to memory ratio
 of 1:8 for standard data warehouses and a ratio of 1:4 for hybrid data warehouses. You
 can distinguish a standard data warehouse from a real-time data warehouse by
 comparing their vCPU to memory ratios.

Precautions

- Currently, neither column-store tables and HStore Opt tables support the **SELECT FOR UPDATE** syntax.
- When a SELECT query is performed on an HStore Opt table, the system will scan the data in column-store primary table CUs, the delta table, and the update information in each row in the memory. The three types of information will be combined before returned.
- In the scenario where data is queried using the primary key index or unique index:

For traditional column-store tables, the unique index stores both the data location information (blocknum, offset) of the row-store Delta table and the data location information (cuid, offset) of the column-store primary table. After the data is merged to the primary table, a new index tuple will be inserted, and the index will keep bloating.

For HStore Opt tables, global CUIDs are allocated in a unified manner. Therefore, only cuid and offset are stored in index tuples. After data is merged, no new index tuples will be generated.

Syntax

```
[ WITH [ RECURSIVE ] with_query [, ...] ]
SELECT [/*+ plan_hint */] [ ALL | DISTINCT [ ON ( expression [, ...] ) ] ]
{* | {expression [ [ AS ] output_name ]} [, ...] }
[ FROM from_item [, ...] ]
[ WHERE condition ]
[ GROUP BY grouping_element [, ...] ]
[ HAVING condition [, ...] ]
[ { UNION | INTERSECT | EXCEPT | MINUS } [ ALL | DISTINCT ] select ]
[ ORDER BY {expression [ [ ASC | DESC | USING operator ] | nlssort_expression_clause ] [ NULLS { FIRST | LAST } ]} [, ...] ]
[ { LIMIT { count | ALL } ] [ OFFSET start [ ROW | ROWS ] ] } | { LIMIT start, { count | ALL } } ]
```

Parameters

• DISTINCT [ON (expression [, ...])]

Removes all duplicate rows from the **SELECT** result set.

ON (expression [, ...]) is only reserved for the first row among all the rows with the same result calculated using given expressions.

SELECT list

Indicates columns to be queried. Some or all columns (using wildcard character *) can be queried.

You may use the **AS output_name** clause to give an alias for an output column. The alias is used for the displaying of the output column.

• **FROM** clause

Indicates one or more source tables for **SELECT**.

The **FROM** clause can contain the following elements:

• WHERE clause

The **WHERE** clause forms an expression for row selection to narrow down the query range of **SELECT**. The condition is any expression that evaluates to a result of Boolean type. Rows that do not satisfy this condition will be eliminated from the output.

In the **WHERE** clause, you can use the operator (+) to convert a table join to an outer join. However, this method is not recommended because it is not the standard SQL syntax and may raise syntax compatibility issues during platform migration. There are many restrictions on using the operator (+):

• GROUP BY clause

Condenses query results into a single row all selected rows that share the same values for the grouped expressions.

HAVING clause

Selects special groups by working with the **GROUP BY** clause. The **HAVING** clause compares some attributes of groups with a constant. Only groups that matching the logical expression in the **HAVING** clause are extracted.

ORDER BY clause

Sorts data retrieved by **SELECT** in descending or ascending order. If the **ORDER BY** expression contains multiple columns:

Example

```
Create the reason select table and insert data into the table.
```

Perform the GROUP BY operation.

SELECT COUNT(*), r_reason_sk FROM reason_select GROUP BY r_reason_sk;

Perform the HAVING filtering operation.

SELECT COUNT(*) c,r_reason_sk FROM reason_select GROUP BY r_reason_sk HAVING c > 1;

Perform the ORDER BY operation.

SELECT * FROM reason_select ORDER BY r_reason_sk;

17.3.8 ALTER TABLE

Function

Modify a table, including modifying the definition of a table, renaming a table, renaming a specified column in a table, adding or updating multiple columns, and changing a column-store table to an HStore Opt table.

□ NOTE

- To use hybrid data warehouse capabilities, choose the storage-compute coupled architecture when you create a GaussDB(DWS) cluster on the console and ensure the vCPU to memory ratio is 1:4 when setting up cloud disk flavors. For more information, see Data Warehouse Flavors.
- When setting up a GaussDB(DWS) cluster, make sure to have a vCPU to memory ratio
 of 1:8 for standard data warehouses and a ratio of 1:4 for hybrid data warehouses. You
 can distinguish a standard data warehouse from a real-time data warehouse by
 comparing their vCPU to memory ratios.

Precautions

- You can set enable_hstore by using ALTER to change a column-store table to an HStore table, or to change it back. If enable_delta is set to on, enable hstore cannot be set to on.
- For some ALTER operations (such as modifying column types, merging partitions, adding NOT NULL constraints, and adding primary key constraints), HStore tables need to merge data to the primary table and then perform ALTER, which may cause extra performance overhead. The impact on performance depends on the data volume in the delta table.
- When you add a column, do not use ALTER to specify other operations (for example, modifying the column type). An ALTER statement with only the ADD COLUMN parameter can achieve high performance, because it does not require FULL MERGE.
- The storage parameter ORIENTATION cannot be modified.

Modifying Table Attributes

Syntax:

ALTER TABLE [IF EXISTS] <table_name> SET ({ENABLE_HSTORE = ON} [, ...]);

To change a column-store table to an HStore table, run the following command:

CREATE TABLE alter_test(a int, b int) WITH(ORIENTATION = COLUMN); ALTER TABLE alter_test SET (ENABLE_HSTORE = ON);

NOTICE

To use HStore tables, set the following parameters, or the HStore performance will deteriorate severely. The recommended settings are as follows:

autovacuum_max_workers_hstore=3, autovacuum_max_workers=6, autovacuum=true

Adding a Column

Syntax:

ALTER TABLE [IF EXISTS] <table_name> ADD COLUMN <new_column> <data_type> [DEFAULT <default_value>];

Example:

Create the alter_test2 table and add a column to it.

CREATE TABLE alter_test2(a int, b int) WITH(ORIENTATION = COLUMN,ENABLE_HSTORE_OPT = ON); ALTER TABLE alter_test ADD COLUMN c int;

∩ NOTE

When adding a column, you are not advised to use **ALTER** to specify other operations in the same SQL statement.

Renaming

Syntax:

ALTER TABLE [IF EXISTS] <table_name> RENAME TO <new_table_name>;

Example:

Create table alter_test3 and rename it as alter_new.

CREATE TABLE alter_test3(a int, b int) WITH(ORIENTATION = COLUMN,ENABLE_HSTORE_OPT = ON); ALTER TABLE alter_test3 RENAME TO alter_new;

17.4 Hybrid Data Warehouse Functions

hstore_light_merge(rel_name text)

Description: This function is used to manually perform lightweight cleanup on HStore tables and holds the level-3 lock of the target table.

Return type: int

Example:

SELECT hstore_light_merge('reason_select');

hstore_full_merge(rel_name text, partitionName text)

Description: This function is used to manually perform full cleanup on HStore tables. The second input parameter is optional and is used to specify a single partition for operations.

Return type: int

NOTICE

- This operation forcibly merges all the visible operations of the delta table to the primary table, and then creates an empty delta table. During this period, this operation holds the level-8 lock of the table.
- The duration of this operation depends on the amount of data in the delta table. You must enable the HStore clearing thread to ensure unnecessary data in the HStore table is cleared in a timely manner.
- The second parameter partitionName is only supported by clusters of version 9.1.0 and later. However, these versions do not allow calling this function via call because it lacks reload capability.

Example:

SELECT hstore_full_merge('reason_select', 'part1');

pgxc_get_small_cu_info(rel_name text, row_count int)

Description: Obtains the small CU information of the target table. The second parameter **row_count** is optional and indicates the small CU threshold. If the number of live tuples in a CU is fewer than the threshold, the CU is considered as a small CU. The default value is **200**. This function is supported only by clusters of version 8.2.1.200 or later.

Return type: record

Return value:

node_name: DN name.

part name: partition name. This column is empty for non-partitioned tables.

zero_cu_count: number of 0 CUs. If all data in a CU is deleted, the CU is called 0 CU.

small_cu_count: number of small CUs. When a CU has live data that is less than the threshold, the CU is called a small CU.

total_cu_count: total number of CUs.

sec_part_cu_num: number of CUs in each level-2 partition. This column is displayed only when **secondary_part_column** is specified. This field is available only in clusters of version 8.3.0 or later.

It should be noted that a CU may contain multiple columns.

Example:

gs_hstore_compaction(rel_name text, row_count int)

Description: Merges small CUs of the target table. The second parameter **row_count** is optional and indicates the small CU threshold. If the number of live tuples in a CU is fewer than the threshold, the CU is considered as a small CU. The default value is **100**. This function is supported only by version 8.2.1.200 or later.

Return type: int

Return value: **numCompactCU**, which indicates the number of small CUs to be merged.

- A CU may contain multiple columns.
- The partition name cannot be input in the function. Currently, a single partition cannot be specified in this function.

Example:

SELECT qs_hstore_compaction('hs', 10);

pgxc_get_hstore_delta_info(rel_name text)

Description: This function is used to obtain the delta table information of the target table, including the delta table size and the number of **INSERT**, **DELETE**, and **UPDATE** records. This function is supported only by clusters of version 8.2.1.100 or later.

Return type: record

Return value:

node name: DN name.

part_name: partition name. This column is set to non-partition table if the table
is not a partitioned table.

live_tup: number of live tuples.

n_ui_type: number of records with a type of *ui* (small CU combination and upsert insertion through update). An **ui** record represents a single or batch insertion. This parameter is supported only by 8.3.0.100 and later versions.

n_i_type: number of records whose type is **i** (insert). An **i** record indicates one insertion, which can be single insertion or batch insertion.

n_d_type: number of records whose type is **d** (delete). One **d** record indicates one deletion, which can be single deletion or batch deletion.

n_x_type: number of records whose type is **x** (deletions generated by update).

n_u_type: number of records whose type is **u** (lightweight update).

n_m_type: number of records whose type is **m** (merge).

data_size: total size of the delta table (including the size of the index and toast data on the delta table).

Example:

```
SELECT * FROM pgxc get hstore delta info('hs part');
node_name | part_name | live_tup | n_ui_type | n_i_type | n_d_type | n_x_type | n_u_type | n_m_type |
data_size
dn_1
                         2 |
                                 0 |
                                         2 |
                                                 0 |
                                                                 0 |
         | p1
                                                         0 [
                                                                          0 1
                                                                                8192
dn 1
                         2 |
                                 0 1
                                         2 |
                                                 0 |
                                                         0 1
                                                                 0 |
                                                                          0 |
                                                                                8192
         | p2
dn_1
                         2 |
                                 0 |
                                                                  0 |
                                                                                8192
         | p3
(3 rows)
```

pgxc_get_binlog_sync_point(rel_name text, slot_name text, checkpoint bool, node id int)

Description: Obtains the synchronization point information corresponding to a slot from the **pg_binlog_slots** system catalog. This function is applicable only to tables with binlog or binlog timestamp enabled. This function is supported only by clusters of version 9.1.0.200 or later.

Return type: record

Return value:

node_name: DN name

node_id: node ID

last_sync_point: last synchronization point

latest_sync_point: latest synchronization point

xmin: xmin corresponding to the synchronization point

Example:

pgxc_get_binlog_changes(rel_name text, node_id int, start_csn bigint, end_csn bigInt)

Description: Obtains the incremental data of the target table within the specified synchronization point range on a specified DN. If **node_id** is set to **0**, all DNs are specified. This function is applicable only to tables with binlog or binlog timestamp enabled. This function is supported only by clusters of version 9.1.0.200 or later.

Return type: record

Return value:

gs_binlog_sync_point: synchronization point

gs_binlog_event_sequence: sequence in the same transaction

gs_binlog_event_type: binlog type

gs_binlog_timestamp_us: timestamp of the binlog record. For the binlog table whose **enable_binlog_timestamp** is **false**, this column is empty.

value columns: data of each user field in the target table

Example:

```
SELECT * FROM pgxc_get_binlog_changes('hstore_binlog_source', 0, 0, 9999999999);
gs_binlog_sync_point | gs_binlog_event_sequence | gs_binlog_event_type | gs_binlog_timestamp_us | c1 | c2
| c3
                                                | 1731570520900211 | 100 | 1 | 1
          10516 |
                                211
                                                     1731570520904425 | 100 | 1 | 1
          10517
                                3 | d
                                                      1731570520909055 | 200 | 1 | 1
          10518 |
                                2 | 1
                                                      1731570520914102 | 200 | 1 | 1
          10519 I
                                3 | B
          10519 |
                                4 | U
                                                      1731570520914154 | 200 | 2 | 1
```

pgxc_register_binlog_sync_point(rel_name text, slot_name text, node_id int, end_csn bigInt, checkpoint bool, xmin bigint)

Description: Registers synchronization points and can be used only for tables with binlog or binlog timestamp enabled. This function is supported only by clusters of version 9.1.0.200 or later.

Return type: int

Return value: number of nodes that are successfully registered

Example:

pgxc_consumed_binlog_records(rel_name text, node_id int)

Description: Obtains the consumption status of the target table on a specified DN. This function can be used only for tables with binlog or binlog timestamp enabled. This function is supported only by clusters of version 9.1.0.200 or later.

Return type: int

Return value: If **0** is returned, the binlog of the target table is not completely consumed (including all slots and checkpoint synchronization points). If **1** is returned, the binlog of the target table is completely consumed.

Example:

```
SELECT * FROM pgxc_consumed_binlog_records('hstore_binlog_source',0);
pgxc_consumed_binlog_records
------

1
(1 row)
```

pgxc_get_binlog_cursor_by_timestamp(rel_name text, timestamp timestampTz, node_id int)

Description: Obtains information about the first binlog record after a specified time point in the target table. This function can be used only for tables with the binlog timestamp enabled.

This function is supported only by clusters of version 9.1.0.200 or later.

Return type: record

Return value:

node name: DN name

node_id: node ID

latest sync point: latest synchronization point

binlog_sync_point: synchronization point of the first binlog record after the time

point

binlog_timestamp_us: timestamp of the first binlog record after the time point

binlog_xmin: xmin recorded in the first binlog after the time point

Example:

```
SELECT * FROM pgxc_get_binlog_cursor_by_timestamp('hstore_binlog_source','2024-11-14 15:48:40.900211+08', 0);
node_name | node_id | latest_sync_point | binlog_sync_point | binlog_timestamp_us | binlog_xmin
```

	+	+	+
	-1051926843	10532	10516 1731570520900211 10510
dn_1		10532	10518 1731570520909055 10510
(2 rows)		

pgxc_get_binlog_cursor_by_syncpoint(rel_name text, csn int8, node_id int)

Description: Obtains the first binlog record after a specified synchronization point on the target table. This function can be used only for tables with the binlog timestamp enabled.

This function is supported only by clusters of version 9.1.0.200 or later.

Return type: record

Return value:

node_name: DN name

node_id: node ID

latest_sync_point: latest synchronization point

binlog_sync_point: synchronization point of the first binlog record after the time

point

binlog_timestamp_us: timestamp of the first binlog record after the time point

binlog_xmin: xmin recorded in the first binlog after the time point

Example:

pgxc_get_binlog_consume_progress(rel_name text, node_id int)

Description: Obtains the consumption progress of the target table. This function can be used only for tables with the binlog timestamp enabled.

This function is supported only by clusters of version 9.1.0.200 or later.

Return type: record

Return value:

node_name: DN name

node_id: node ID

slot name: slot name

checkpoint: indicates whether the location is a checkpoint.

latest_consumed_timestamp: timestamp of the latest consumed binlog

latest timestamp: latest timestamp in all binlogs

latest_consumed_csn: CSN of the latest consumed binlog

latest_csn: latest CSN in all binlogs

unconsumed_binlog_count: number of binlogs that have not been consumed

Example:

pgxc_get_cstore_dirty_ratio(rel_name text, partition_name)

Description: This function is used to obtain the cu, delta, and cudesc dirty page rates and sizes of the target table on each DN. Only HStore Opt tables are supported.

The **partition_name** parameter is optional. If the partition name is specified, only the information about the partition is returned. If the partition name is not specified and the table is a primary table, the information about all partitions is returned. It is supported only by clusters of version 9.1.0.100 or later.

Return type: record

Return value:

node_name: DN name

database_name: name of the database where the table is located

rel_name: primary table name

part_name: partition name

cu_dirty_ratio: dirty page rate of CU files

cu_size: CU file size

delta dirty ratio: dirty page rate of the delta table

delta size: delta table size

cudesc_dirty_ratio: dirty page rate of the cudesc table

cudesc_size: cudesc table size

Example:

```
SELECT * FROM pgxc_get_cstore_dirty_ratio('hs_opt_part');
node_name | database_name | rel_name | partition_name | cu_dirty_ratio | cu_size | delta_dirty_ratio
| delta_size | cudesc_dirty_ratio | cudesc_size
-----+----+-----+-----+----
dn_1 | postgres | public.hs_opt_part | p1
                                              0 | 0 |
                                                                          0 | 16384
         0 | 24576
       | postgres | public.hs_opt_part | p2
                                                        0 |
                                                              0 |
                                                                          0 |
                                                                                16384
         0 | 24576
dn_1
       | postgres | public.hs_opt_part | p3
                                                        0 |
                                                              0 |
                                                                          0 |
                                                                               16384
        0 | 24576
```

dn_1	postgres public.hs_opt_part p4	0 0	0 16384
	0 24576		
dn_1	postgres public.hs_opt_part other	0 1105920	0 524288
	0 40960		

17.5 Hybrid Data Warehouse Binlog

17.5.1 Subscribing to Hybrid Data Warehouse Binlog

Binlog Usage

The HStore table within the GaussDB(DWS) hybrid data warehouse offers binlog to facilitate the capture of database events. This enables the export of incremental data to third-party components like Flink. By consuming binlog data, you can synchronize upstream and downstream data, improving data processing efficiency.

Unlike traditional MySQL binlog, which logs all database changes and focuses on data recovery and replication. The GaussDB(DWS) hybrid data warehouse binlog is optimized for real-time data synchronization, recording DML operations—INSERT, DELETE, and UPDATE (UPSERT)—while excluding DDL operations.

GaussDB(DWS) Binlog has the following advantages:

- Table-level on-demand switch: enables or disables binlog for specific tables as needed.
- Full incremental integrated consumption: supports full synchronization followed by real-time incremental consumption after a Flink task is started.
- Cleanup upon consumption: allows asynchronous clearing of incremental data after consumption, reducing space usage.

With Flink's real-time processing capabilities and Binlog, you can build a hybrid data warehouse efficiently without additional components like Kafka. The architecture is streamlined, and data flows efficiently, driven by Flink SQL.

Constraints and Limitations

- 1. Currently, only 8.3.0.100 and later versions support HStore and HStore Opt to record binlogs. V3 tables are in the trial commercial use phase. Before using them, contact technical support for evaluation.
- 2. Binlog requires a primary key, an HStore or HStore-opt table, and supports only hash distribution.
- 3. Binlog tables log DML operations like INSERT, DELETE, and UPDATE (UPSERT), excluding DDL operations.
- 4. Binlog tables do not support INSERT OVERWRITE, altering distribution columns, enabling Binlog on temporary tables, or partition operations like EXCHANGE, MERGE, and SPLIT PARTITION.
- 5. Users can perform the following DDL operations, but these will reset incremental data and synchronization details.
 - ADD COLUMN, DROP COLUMN, SET TYPE, and TRUNCATE
- 6. The system waits for binlog consumption before further scaling. The default wait time is 1 hour. Timeouts or errors will cause the scaling process to fail.

- 7. The system waits for the consumption of binlog records before the **VACUUM FULL** operation is performed on a binlog table. The default wait time is 1
 hour. Timeouts or errors will cause the **VACUUM FULL** process to fail.
 Additionally, even if VACUUM FULL is executed for a partition table, a level-7
 lock is added to the primary table of the partition, which blocks the insertion, update, or deletion of the entire table.
- 8. Binlog tables are backed up as standard HStore tables. Post-restoration, you must restart data synchronization as incremental data and sync details are reset.
- 9. The Binlog timestamp function is supported. This function can be enabled by activating **enable_binlog_timestamp**. Only the HStore and HStore Opt tables support this function. This constraint is supported only in 9.1.0.200 and later versions.

Binlog Formats and Principles

Table 17-6 binlog fields

Field	Туре	Description
gs_binlog_sync _point	BIGINT	Binlog system field, which indicates the synchronization point. In common GTM mode, the value is unique and ordered.
gs_binlog_even t_sequence	BIGINT	Binlog system field, which indicates the sequence of operations of the same transaction type.
gs_binlog_even t_type	CHAR	Binlog system field, which indicates the operation type of the current record. The options are as follows:
		I refers to INSERT, indicating that a new record is inserted into the current binlog.
		d refers to DELETE, indicating that a record is deleted from the current binlog.
		B refers to BEFORE_UPDATE, indicating that the current binlog is a record before the update.
		U refers to AFTER_UPDATE, indicating that the current binlog is a record after the update.
gs_binlog_time stamp_us	BIGINT	System field of Binlog, indicating the timestamp when the current record is saved to the database.
		This field is available only when the Binlog timestamp function is enabled. If the Binlog timestamp function is disabled, this field is left blank. Only 9.1.0.200 and later versions support this function.

Field	Туре	Description
user_column_1	User column	User-defined data column
user_column_n	User column	User-defined data column

NOTE

- For each UPDATE (or UPSERT-triggered update), two binlog records—BEFORE_UPDATE and AFTER_UPDATE—are created. BEFORE_UPDATE verifies the accuracy of data processed by third-party components like Flink.
- During UPDATE and DELETE operations, the GaussDB(DWS) hybrid data warehouse generates BEFORE_UPDATE and DELETE binlogs without querying or populating all user columns, enhancing database import efficiency.
- Enabling binlog for an HStore table in the GaussDB(DWS) hybrid data warehouse is in fact the process of creation of a supplementary table. This table includes three system columns gs_binlog_event_sync_point, gs_binlog_event_event_sequence, and gs_binlog_event_type, and a value column that serializes all user columns.
- When the enable_binlog_timestamp parameter is enabled, binlog records are retained until the TTL expires, causing extra space overhead proportional to the data volume updated within the TTL. When enable_binlog is enabled, binlogs can be cleared asynchronously once consumed by downstream processes, significantly reducing space usage. Only 9.1.0.200 and later versions support this function.

Enabling Binlog

You can specify the table-level parameter **enable_binlog** when creating an HStore Opt table to enable binlogs.

```
CREATE TABLE hstore_binlog_source (
c1 INT PRIMARY KEY,
c2 INT,
c3 INT
) WITH (
ORIENTATION = COLUMN,
enable_hstore_opt=true,
enable_binlog=on,
binlog_ttl = 86400
);
```

■ NOTE

- Binlog recording begins only after a synchronization point is registered for the task, not during the initial data import. Once binlog synchronization in Flink is activated, it periodically acquires the synchronization point and incremental data, then registers the synchronization point.
- The binlog_ttl parameter defaults to 86,400 seconds and is optional. If a registered synchronization point exceeds this TTL without undergoing incremental synchronization, it will be cleared. Subsequently, binlogs before the oldest synchronization point are asynchronously deleted to free up space.
- Space overhead: For a table with common binlog enabled, if incremental data can be consumed by downstream processes in a timely manner, the space can be cleared and reclaimed promptly.

Run the **ALTER** command to enable the binlog function for an existing HStore table

```
CREATE TABLE hstore_binlog_source (
    c1 INT PRIMARY KEY,
    c2 INT,
    c3 INT
) WITH (
    ORIENTATION = COLUMN,
    enable_hstore_opt=true
);
ALTER TABLE hstore_binlog_source SET (enable_binlog=on);
```

Querying Binlogs

You can use the system functions provided by GaussDB(DWS) to query the binlog information of the target table on a specified DN and check whether the binlog is consumed by downstream processes.

```
-- Simulate Flink to call a system function to obtain the synchronization point. The parameters indicate the table name, slot name, whether the point is a checkpoint, and target DN (0 indicates all DNs). select * from pg_catalog.pgxc_get_binlog_sync_point('hstore_binlog_source', 'slot1', false, 0); select * from pg_catalog.pgxc_get_binlog_sync_point('hstore_binlog_source', 'slot1', true, 0); -- Incremental binlogs are generated after additions, deletions, and modifications. INSERT INTO hstore_binlog_source VALUES(100, 1, 1); delete hstore_binlog_source where c1 = 100; INSERT INTO hstore_binlog_source VALUES(200, 1, 1); update hstore_binlog_source set c2 = 2 where c1 = 200; -- Simulate Flink to call a system function to query the binlog of a specified CSN range. The parameters indicate the table name, target DN (0 indicates all DNs), start CSN point, and end CSN point. select * from pgxc_get_binlog_changes('hstore_binlog_source', 0, 0, 9999999999);
```

```
      postgres=# select * from pgxc_get_binlog_changes('hstore_binlog_source', 0, 0, 9999999999);

      gs_binlog_sync_point | gs_binlog_event_sequence | gs_binlog_event_type | gs_binlog_timestamp_us | c1 | c2 | c3

      10241 | 2 | I
      | 100 | 1 | 1

      10242 | 3 | 4 | I
      | 100 | 1 | 1

      10243 | 4 | I
      | 100 | 1 | 1

      10245 | 5 | B
      | 100 | 1 | 1

      10245 | 6 | U
      | 100 | 100 | 1

      (5 rows)
```

Two **INSERT** operations generate two records with **gs_binlog_event_type** as **I**. The **DELETE** operation generates a record whose type is **d**. The **UPDATE** operation generates a **B** record for **BeforeUpdate** and a **U** record for **AfterUpdate**, indicating the values before and after the update.

You can call the system function pgxc_consumed_binlog_records to check whether the binlogs of the target table are consumed by all slots. The parameters indicate the target table name and target DN (0 indicates all DNs).

```
-- Simulate Flink to call the system function to register a synchronization point. The parameters indicate the table name, slot name, registered point, whether the point is a checkpoint, and xmin corresponding to the point (provided when the synchronization point is obtained). select pgxc_register_binlog_sync_point('hstore_binlog_source', 'slot1', 0, 9999999999, false, 100); select pgxc_register_binlog_sync_point('hstore_binlog_source', 'slot1', 0, 9999999999, true, 100); -- Check whether all binlogs in the table are consumed. If 1 is returned, all binlogs have been consumed by downstream slots. select * from pgxc_consumed_binlog_records('hstore_binlog_source',0);
```

Enabling the Binlog Timestamp Function

If you need to read binlogs generated after a specified time point, specify the table-level parameter **enable_binlog_timestamp** when creating an HStore table

to enable the binlog timestamp function of the HStore table. Only 9.1.0.200 and later versions support this function.

```
CREATE TABLE hstore_binlog_source(
    c1 INT PRIMARY KEY,
    c2 INT,
    c3 INT
) WITH (
    ORIENTATION = COLUMN,
    enable_hstore_opt=true,
    enable_binlog_timestamp =on,
    binlog_ttl = 86400
);
```

□ NOTE

- Binlog recording begins only after a synchronization point is registered for the task, not during the initial data import. Once the binlog timestamp is enabled, the system periodically acquires the synchronization point and incremental data, then registers the synchronization point.
- Binlog_ttl is an optional parameter. If not set, the default value is **86400** seconds (i.e., data is retained for one day by default). If the timestamp of the binlog record is greater than the current TTL, the binlog record will be deleted asynchronously.
- Space overhead: For a table with the binlog timestamp enabled, the binlog records recorded in the auxiliary table are retained until the TTL expires. This results in extra space overhead, which is proportional to the amount of data updated and imported into the database within the TTL.

Query the binlog on the table where the binlog timestamp function is enabled.

```
postgres=# select * from pgxc_get_binlog_changes('hstore_binlog_source', 0, 0 , 9999999999);

gs_binlog_sync_point | gs_binlog_event_sequence | gs_binlog_event_type | gs_binlog_timestamp_us | c1 | c2 | c3

10516 | 2 | I | 17315705209000211 | 100 | 1 | 1

10517 | 3 | d | 1731570520900425 | 100 | 1 | 1

10518 | 2 | I | 1731570520904425 | 100 | 1 | 1

10519 | 3 | B | 1731570520914102 | 200 | 1 | 1

10519 | 4 | U | 1731570520914154 | 200 | 2 | 1

(5 rows)
```

Convert **gs_binlog_timestamp_us** from the BigInt type to a readable timestamp.

select to_timestamp(1731569598408661/1000000);

To obtain the first binlog information of the target table after the specified time point on each DN (if the value is empty, no binlog exists after the time point).

select * from pgxc_get_binlog_cursor_by_timestamp('hstore_binlog_source','2024-11-14 15:33:18.40866+08', 0):

```
postgres=# select * from pgxc_get_binlog_cursor_by_timestamp('hstore_binlog_source','2024-11-14 15:48:40.900211+08', 0);
node_name | node_id | latest_sync_point | binlog_sync_point | binlog_timestamp_us | binlog_xmin

dn_2 | -1051926843 | 10532 | 10516 | 1731570520900211 | 10510
dn_1 | -1300059100 | 10532 | 10518 | 1731570520909055 | 10510
(2 rows)
```

Obtain the consumption progress of the table for which the binlog timestamp function is enabled.

The returned fields indicate the timestamp of the latest consumed binlog, the latest timestamp on the binlog, the CSN point of the latest consumed binlog, the latest CSN point on the binlog, and the number of unconsumed binlog records.

-- Simulate Flink to call the system function to register a synchronization point. The parameters indicate the table name, slot name, registered point, whether the point is a checkpoint, and **xmin** corresponding to the point (provided when the synchronization point is obtained). select pgxc_register_binlog_sync_point('hstore_binlog_source', 'slot1', 0, 9999999999, false, 100); select pgxc_register_binlog_sync_point('hstore_binlog_source', 'slot1', 0, 9999999999, true, 100); -- Query the consumption progress of each slot in the target table. select * from pgxc_get_binlog_consume_progress('hstore_binlog_source', 0);

```
postgres# select * from pax get binlog consume progress('hstore binlog source', 0); node_name | node_id | slot_name | checkpoint | latest_consumed_timestamp | latest_timestamp | latest_consumed_csn | latest_csn | unconsumed_binlog_count | dn_1 | .1300059100 | slot1 | f | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 10519 | 10519 | 0 | dn_1 | .1300059100 | slot1 | t | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 10519 | 10519 | 0 | dn_2 | .1051926843 | slot1 | f | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 10517 | 10517 | 0 | dn_2 | .1051926843 | slot1 | t | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 10517 | 10517 | 0 | dn_2 | .1051926843 | slot1 | t | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 10517 | 10517 | 0 | dn_2 | .1051926843 | slot1 | t | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 10517 | 10517 | 0 | dn_2 | .1051926843 | slot1 | t | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-11-14 15:48:40+08 | 2024-
```

Preventing DML from Generating Binlogs

You can set the session-level parameter **enable_generate_binlog** to **off** to control the DML of the current session. When a table for which binlog is enabled is imported to the database, no binlog record is generated.

17.5.2 Real-Time Binlog Consumption by Flink

Precautions

- Currently, only 8.3.0.100 and later versions support HStore and HStore Opt tables to record binlogs.
- V3 HStore tables do not support binlogs, whereas only V3 HStore Opt tables offer binlog support. V3 is currently in trial commercial use and should be thoroughly evaluated before deployment.
- The Binlog function is only supported for HStore and HStore Opt tables in GaussDB(DWS). These tables must have primary keys and one of parameters **enable_binlog** and **enable_binlog_timestamp** must be set to **on**.
- The name of the consumed binlog table cannot contain special characters, such as periods (.) and double quotation marks (").
- If multiple tasks consume binlog data of a single table, ensure that **binlogSlotName** of each task is unique.
- For maximum consumption speed, match task concurrency with the number of DNs in your GaussDB(DWS) cluster.
- If you use the sink capability of **dws-connector-flink** to write binlog data, pay attention to the following:
 - To ensure the data write sequence on DNs, set **connectionSize** to **1**.
 - If the primary key is updated on the source end or Flink is required for aggregation calculation, set ignoreUpdateBefore to false. Otherwise, you are not advised to set ignoreUpdateBefore to false (the default value is true).

Real-Time Binlog Consumption by Flink

Use DWS Connector to consume binlogs in real time. For details, see **DWS-Connector**.

If full data has been synchronized to the target end using other synchronization tools, and only incremental synchronization is required, you can call the following system function to update the synchronization points.

SELECT * FROM pg_catalog.pgxc_register_full_sync_point('table_name', 'slot_name');

Source Table DDL

The source autonomously assigns the appropriate Flink RowKind type (INSERT, DELETE, UPDATE_BEFORE, or UPDATE_AFTER) to each data row based on the operation type. This mechanism facilitates the synchronization of table data in a mirrored way, akin to the Change Data Capture (CDC) feature in MySQL and PostgreSQL databases.

```
CREATE TABLE test_binlog_source (
a int,
b int,
c int,
primary key(a) NOT ENFORCED
) with (
'connector' = 'dws',
'url' = 'jdbc:gaussdb://ip:port/gaussdb',
'binlog' = 'true',
'tableName' = 'test_binlog_source',
'binlogSlotName' = 'slot',
'username'='xxx',
'password'='xxx')
```

Binlog Parameters

The following table describes the parameters involved in binlog consumption.

Table 17-7 Parameters

Parameter	Description	Data Type	Default Value
binlog	Specifies whether to read binlog information.	Boolean	false
binlogSlotName	Slot, which serves as an identifier. Multiple Flink tasks can simultaneously consume binlog data of the same table, so each task's binlogSlotName must be unique.	String	Name of the Flink mapping table
binlogBatchRead- Size	Rows of binlog data read in batches.	Integer	5000
fullSyncBinlogBat- chReadSize	Rows of binlog data fully read.	Integer	50000
binlogReadTimeout	Timeout for incrementally consuming binlog data, in milliseconds.	Integer	600000
fullSyncBinlogRead- Timeout	Timeout for fully consuming binlog data, in milliseconds.	Long	1800000

Parameter	Description	Data Type	Default Value
binlogSleepTime	Sleep duration when no real- time binlog data is consumed, in milliseconds. The sleep duration with consecutive read failures is binlogSleepTime * failures, up to binlogMaxSleepTime. The value is reset after successful data read.		500
binlogMaxSleepTim e	Maximum sleep duration when no real-time binlog data is consumed, in milliseconds.	Long	10000
binlogMaxRetryTim es	Maximum number of retries after a binlog data consumption error.	Integer	1
binlogRetryInterval	Interval between retries after a binlog data consumption error, in milliseconds. Sleep duration during retry, which is calculated as binlogRetryInterval * (1~binlogMaxRetryTimes) + Random(100). The unit is millisecond.	Long	100
binlogParallelNum	Number of threads for consuming binlog data. This parameter is valid only when task concurrency is less than the number of DNs in the GaussDB(DWS) cluster.	Integer	3
connectionPoolSize	Number of connections in the JDBC connection pool.	Integer	5
needRedistribution	Determines compatibility with expansion redistribution. To ensure compatibility, upgrade the kernel to the corresponding version. If the kernel is an older version, set this parameter to false. If set to true, restart-strategy of Flink cannot be set to none.	Boolean	true

Parameter	Description	Data Type	Default Value
newSystemValue	Indicates whether to use the new system field when reading binlog data. (The kernel needs to be upgraded to the corresponding version. If the kernel is an older version, set this parameter to false.)	Boolean	true
checkNodeChangeI nterval	Interval for detecting node changes. This parameter is valid only when needRedistribution is set to true .	Long	10000
connectionSocket- Timeout	Timeout interval for connection processing, in milliseconds. It can also be considered as the timeout interval for executing SQL statements on the client. The default value is 0 , which means that the timeout interval is not set.	Integer	0
binlogIgnoreUpda- teBefore	Determines whether to filter out before_update records in binlogs and whether to return only primary key information for delete records. This parameter is supported only in 9.1.0.200 and later versions.	Boolean	false
binlogStartTime	Sets the time point from which binlogs are consumed can be set using the format yyyy-MM-dd hh:mm:ss. enable_binlog_timestamp must be enabled for the table. This parameter is supported only in 9.1.0.200 and later versions.	String	N/A
binlogSyncPointSize	Specifies the size of the synchronization point range for incrementally reading binlogs. This can control data flushing if the data volume is too large. This parameter is supported only in 9.1.0.200 and later versions.	Integer	5000

Data Synchronization Example

On GaussDB(DWS):

When creating a binlog table, set **enable_hstore_binlog_table** to true. You can run the **show enable_hstore_binlog_table** command to query the binlog table.

-- Source table (generating binlogs)

CREATE TABLE test_binlog_source(a int, b int, c int, primary key(a)) with(orientation=column, enable_hstore_opt=on, enable_binlog=true);

-- Target table

CREATE TABLE test_binlog_sink(a int, b int, c int, primary key(a)) with(orientation=column, enable_hstore_opt=on);

On Flink:

Run the following commands to perform complete data synchronization:

```
-- Create a mapping table for the source table.
CREATE TABLE test_binlog_source (
  a int,
 b int,
 c int,
  primary key(a) NOT ENFORCED
) with (
  'connector' = 'dws',
  'url' = 'jdbc:gaussdb://ip:port/gaussdb',
  'binlog' = 'true',
  'tableName' = 'test_binlog_source',
  'binlogSlotName' = 'slot',
  'username'='xxx',
  'password'='xxx');
-- Create a mapping table for the target table:
CREATE TABLE test_binlog_sink (
  a int,
  b int,
 c int,
  primary key(a) NOT ENFORCED
) with (
  'connector' = 'dws',
  'url' = 'jdbc:gaussdb://ip:port/gaussdb',
  'tableName' = 'test_binlog_sink',
  'ignoreUpdateBefore'='false',
  'connectionSize' = '1',
  'username'='xxx',
  'password'='xxx');
INSERT INTO test_binlog_sink select * from test_binlog_source;
```

Example of Using Java Programs

Create a source table and a target table.

```
-- source
create table binlog_test_source(a int, b int, c int, primary key(a)) with(orientation=column,
enable_hstore_opt=on, enable_binlog=true);
-- sink
create table binlog_test_sink(a int, b int, c int, primary key(a)) with(orientation=column,
enable_hstore_opt=on, enable_binlog=true);
```

Demo program:

```
public class BinlogDemo {

//Name of the binlog table
```

```
private static final String BINLOG_TABLE_NAME = "binlog_test_source";
  //Slot name of the binlog table
  private static final String BINLOG_SLOT_NAME = "binlog_test_slot";
  //Name of the table to be written
  private static final String SINK_TABLE_NAME = "binlog_test_sink";
  public static void main(String[] args) throws Exception {
     DwsConfig dwsConfig = buildDwsConfig();
     DwsClient dwsClient = new DwsClient(dwsConfig);
     TableSchema sourceTableSchema =
dwsClient.getTableSchema(TableName.valueOf(BINLOG_TABLE_NAME));
     TableSchema sinkTableSchema = dwsClient.getTableSchema(TableName.valueOf(SINK_TABLE_NAME));
     // Columns to be written
     List<String> sinkColumns = sinkTableSchema.getColumnNames();
     // Thread pool
     DwsConnectionPool dwsConnectionPool = new DwsConnectionPool(dwsConfig);
     //Queue for storing data
     BlockingQueue<BinlogRecord> queue = new LinkedBlockingQueue<>();
     //Columns to be synchronized
     List<String> sourceColumnNames = sourceTableSchema.getColumnNames();
     BinlogReader binlogReader = new BinlogReader(dwsConfig, queue, sourceColumnNames,
dwsConnectionPool);
     //Start the read task.
     binlogReader.start();
     binlogReader.getRecords();
     while (binlogReader.isStart()) {
       try {
          while (!queue.isEmpty() && !binlogReader.hasException()) {
             // Read data.
             BinlogRecord record = queue.poll();
             if (Objects.isNull(record)) {
               continue;
             BinlogRecordType type = BinlogRecordType.toBinlogRecordType(record.getType());
             List<Object> columnValues = record.getColumnValues();
             if (BinlogRecordType.INSERT.equals(type) || BinlogRecordType.UPDATE AFTER.equals(type)) {
               Operate upsert = dwsClient.write(sinkTableSchema);
               for (int i = 0; i < sinkColumns.size(); i++) {
                  upsert.setObject(i, columnValues.get(i), false);
               upsert.commit();
             } else if (BinlogRecordType.DELETE.equals(type) ||
BinlogRecordType.UPDATE BEFORE.equals(type)) {
               Operate delete = dwsClient.delete(sinkTableSchema);
               for (int i = 0; i < sinkColumns.size(); i++) {
                  String field = sinkColumns.get(i);
                  if (!sinkTableSchema.isPrimaryKey(field)) {
                     continue;
                  delete.setObject(i, columnValues.get(i), false);
                delete.commit();
            }
          binlogReader.checkException();
       } catch (Exception e) {
          throw new DwsClientException(ExceptionCode.GET_BINLOG_ERROR, "get binlog has error", e);
```

```
private static DwsConfig buildDwsConfig() {
    //Initialize configuration information. (Only necessary parameters are listed. For more information
about the configuration, see the document.)
    TableConfig tableConfig = new TableConfig().withBinlog(true)
        .withNewSystemValue(true)
        .withNeedRedistribution(false)
        .withBinlogSlotName(BINLOG_SLOT_NAME);
return DwsConfig.builder()
        .withUrl("Link information")
        .withUsername("Username")
        .withPassword ("Password")
        .withPassword ("Password")
        .withBinlogTableName(BINLOG_TABLE_NAME)
        .withTableConfig(BINLOG_TABLE_NAME, tableConfig)
        .build();
}
```