Data Lake Insight

## **User Guide**

 Issue
 01

 Date
 2025-07-11





HUAWEI TECHNOLOGIES CO., LTD.

#### Copyright © Huawei 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 Technologies Co., Ltd.

#### **Trademarks and Permissions**

NUAWEI and other Huawei trademarks are trademarks 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 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, guarantees 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.

## **Security Declaration**

#### Vulnerability

Huawei's regulations on product vulnerability management are subject to the *Vul. Response Process.* For details about this process, visit the following web page:

https://www.huawei.com/en/psirt/vul-response-process

For vulnerability information, enterprise customers can visit the following web page: <u>https://securitybulletin.huawei.com/enterprise/en/security-advisory</u>

## **Contents**

1 DLI Job Development Process	1
2 Preparations	5
- 2.1 Configuring DLI Agency Permissions	5
2.2 Creating an IAM User and Granting Permissions	11
2.3 Configuring a DLI Job Bucket	13
3 Creating an Elastic Resource Pool and Queues Within It	.15
3.1 Overview of DLI Elastic Resource Pools and Queues	15
3.2 Creating an Elastic Resource Pool and Creating Queues Within It	22
3.3 Managing Elastic Resource Pools	30
3.3.1 Viewing Basic Information	31
3.3.2 Managing Permissions	33
3.3.3 Binding a Queue	35
3.3.4 Setting CUs	36
3.3.5 Modifying Specifications	38
3.3.6 Managing Tags	42
3.3.7 Adjusting Scaling Policies for Queues in an Elastic Resource Pool	44
3.3.8 Viewing Scaling History	50
3.3.9 Allocating to an Enterprise Project	51
3.4 Managing Queues	52
3.4.1 Viewing Basic Information About a Queue	52
3.4.2 Queue Permission Management	53
3.4.3 Allocating a Queue to an Enterprise Project	56
3.4.4 Creating an SMN Topic	57
3.4.5 Managing Queue Tags	58
3.4.6 Setting Queue Properties	60
3.4.7 Enabling Spark Native Operator Optimization	67
3.4.8 Testing Address Connectivity	69
3.4.9 Deleting a Queue	70
3.4.10 Enabling Elastic Scaling for a Queue in a Non-Elastic Resource Pool	70
3.4.11 Creating a Scheduled Elastic Scaling Task for a Queue in a Non-Elastic Resource Pool	73
3.4.12 Changing the CIDR Block of a Queue in a Non-Elastic Resource Pool	77
3.5 Example Use Case: Creating an Elastic Resource Pool and Running Jobs	78

3.6 Example Use Case: Configuring Scaling Policies for Queues in an Elastic Resource Pool	84
3.7 Creating a Non-Elastic Resource Pool Queue (Deprecated, Not Recommended)	89
4 Creating a Data Directory, Database, and Table	94
4.1 Understanding Data Catalogs, Databases, and Tables	94
4.2 Creating a Data Catalog, Database, and Table on the DLI Console	
4.3 Viewing Table Metadata	108
4.4 Managing Data Catalogs on the DLI Console	109
4.4.1 Configuring Data Catalog Permissions on the DLI Console	109
4.5 Managing Database Resources on the DLI Console	112
4.5.1 Configuring Database Permissions on the DLI Console	112
4.5.2 Deleting a Database on the DLI Console	118
4.5.3 Changing the Database Owner on the DLI Console	119
4.5.4 Managing Tags	120
4.6 Managing Table Resources on the DLI Console	123
4.6.1 Configuring Table Permissions on the DLI Console	123
4.6.2 Deleting a Table on the DLI Console	131
4.6.3 Changing the Table Owner on the DLI Console	132
4.6.4 Importing OBS Data to a DLI Table	133
4.6.5 Exporting DLI Table Data to OBS	136
4.6.6 Previewing Table Data on the DLI Console	140
4.7 Creating and Using LakeFormation Metadata	140
4.7.1 Connecting DLI to LakeFormation	140
4.7.2 Permission Policies and Supported Actions for LakeFormation Resources	164
5 Data Import and Migration	172
5.1 Overview	172
5.2 Migrating Data from External Data Sources to DLI	173
5.2.1 Overview of Data Migration Scenarios	173
5.2.2 Using CDM to Migrate Data to DLI	175
5.2.3 Example Typical Scenario: Migrating Data from Hive to DLI	182
5.2.4 Example Typical Scenario: Migrating Data from Kafka to DLI	191
5.2.5 Example Typical Scenario: Migrating Data from Elasticsearch to DLI	199
5.2.6 Example Typical Scenario: Migrating Data from RDS to DLI	206
5.2.7 Example Typical Scenario: Migrating Data from GaussDB(DWS) to DLI	214
6 Configuring DLI to Read and Write Data from and to External Data Sources.	222
6.1 Configuring DLI to Read and Write External Data Sources	222
6.2 Configuring the Network Connection Between DLI and Data Sources (Enhanced Datasource Connection)	223
6.2.1 Overview of Enhanced Datasource Connections	223
6.2.2 Creating an Enhanced Datasource Connection	226
6.2.3 Establishing a Network Connection Between DLI and Resources in a Shared VPC	232
6.2.4 Common Development Methods for DLI Cross-Source Analysis	234
6.3 Using DEW to Manage Access Credentials for Data Sources	237

6.4 Using DLI Datasource Authentication to Manage Access Credentials for Data Sources	240
6.4.1 Overview	240
6.4.2 Creating a CSS Datasource Authentication	243
6.4.3 Creating a Kerberos Datasource Authentication	245
6.4.4 Creating a Kafka_SSL Datasource Authentication	
6.4.5 Creating a Password Datasource Authentication	252
6.4.6 Datasource Authentication Permission Management	255
6.5 Managing Enhanced Datasource Connections	256
6.5.1 Viewing Basic Information About an Enhanced Datasource Connection	257
6.5.2 Enhanced Connection Permission Management	257
6.5.3 Binding an Enhanced Datasource Connection to an Elastic Resource Pool	258
6.5.4 Unbinding an Enhanced Datasource Connection from an Elastic Resource Pool	
6.5.5 Adding a Route for an Enhanced Datasource Connection	
6.5.6 Deleting the Route for an Enhanced Datasource Connection	262
6.5.7 Modifying Host Information in an Elastic Resource Pool	263
6.5.8 Enhanced Datasource Connection Tag Management	
6.5.9 Deleting an Enhanced Datasource Connection	
6.6 Example Typical Scenario: Connecting DLI to a Data Source on a Private Network	267
6.7 Example Typical Scenario: Connecting DLI to a Data Source on a Public Network	273
7 Configuring an Agency to Allow DLI to Access Other Cloud Services	281
7.1 DLI Agency Overview	
7.2 Creating a Custom DLI Agency	
7.3 Agency Permission Policies in Common Scenarios	
7.4 Example of Configuring DLI Agency Permissions in Typical Scenarios	295
8 Submitting a SQL Job on the DLI Management Console	296
9.1 Creating and Submitting a SOL Jab	
8.1 Creating and Submitting a SQL Job	
8.2 Example of a Typical Scenario: Analyzing OBS Data Using a Spark SQL Job	
<ul><li>8.1 Creating and Submitting a SQL Job.</li><li>8.2 Example of a Typical Scenario: Analyzing OBS Data Using a Spark SQL Job.</li><li>8.3 Exporting SQL Job Results.</li></ul>	301 312
<ul> <li>8.1 Creating and Submitting a SQL Job.</li> <li>8.2 Example of a Typical Scenario: Analyzing OBS Data Using a Spark SQL Job.</li> <li>8.3 Exporting SQL Job Results.</li> <li>8.4 Creating a SQL Inspection Rule.</li> </ul>	
<ul> <li>8.1 Creating and Submitting a SQL Job.</li> <li>8.2 Example of a Typical Scenario: Analyzing OBS Data Using a Spark SQL Job.</li> <li>8.3 Exporting SQL Job Results.</li> <li>8.4 Creating a SQL Inspection Rule.</li> <li>8.5 Setting the Priority for a SQL Job.</li> </ul>	
<ul> <li>8.1 Creating and Submitting a SQL Job</li></ul>	
<ul> <li>8.1 Creating and Submitting a SQL Job.</li> <li>8.2 Example of a Typical Scenario: Analyzing OBS Data Using a Spark SQL Job.</li> <li>8.3 Exporting SQL Job Results.</li> <li>8.4 Creating a SQL Inspection Rule.</li> <li>8.5 Setting the Priority for a SQL Job.</li> <li>8.6 Querying Logs for SQL Jobs.</li> <li>8.7 Managing SQL Jobs.</li> </ul>	
<ul> <li>8.1 Creating and Submitting a SQL Job.</li> <li>8.2 Example of a Typical Scenario: Analyzing OBS Data Using a Spark SQL Job.</li> <li>8.3 Exporting SQL Job Results.</li> <li>8.4 Creating a SQL Inspection Rule.</li> <li>8.5 Setting the Priority for a SQL Job.</li> <li>8.6 Querying Logs for SQL Jobs.</li> <li>8.7 Managing SQL Jobs.</li> <li>8.8 Viewing a SQL Execution Plan.</li> </ul>	
<ul> <li>8.1 Creating and Submitting a SQL Job</li></ul>	
<ul> <li>8.1 Creating and Submitting a SQL Job.</li> <li>8.2 Example of a Typical Scenario: Analyzing OBS Data Using a Spark SQL Job.</li> <li>8.3 Exporting SQL Job Results.</li> <li>8.4 Creating a SQL Inspection Rule.</li> <li>8.5 Setting the Priority for a SQL Job.</li> <li>8.6 Querying Logs for SQL Jobs.</li> <li>8.7 Managing SQL Jobs.</li> <li>8.8 Viewing a SQL Execution Plan.</li> <li>8.9 Creating and Managing SQL Job Templates.</li> <li>8.9.1 Creating a SQL Job Template.</li> </ul>	
<ul> <li>8.1 Creating and Submitting a SQL Job</li></ul>	
<ul> <li>8.1 Creating and Submitting a SQL Job</li></ul>	301 312 316 324 325 327 330 331 332 334 335
<ul> <li>8.1 Creating and Submitting a SQL Job</li> <li>8.2 Example of a Typical Scenario: Analyzing OBS Data Using a Spark SQL Job</li> <li>8.3 Exporting SQL Job Results</li> <li>8.4 Creating a SQL Inspection Rule</li> <li>8.5 Setting the Priority for a SQL Job</li> <li>8.6 Querying Logs for SQL Jobs</li> <li>8.7 Managing SQL Jobs</li> <li>8.8 Viewing a SQL Execution Plan</li> <li>8.9 Creating and Managing SQL Job Templates</li> <li>8.9.1 Creating a SQL Job Template</li> <li>8.9.2 Developing and Submitting a SQL Job Using a SQL Job Template</li> <li>8.9.3 TPC-H Sample Data in the SQL Templates Preset on DLI</li> <li>9 Developing a DLI SQL Job in DataArts Studio</li> </ul>	301 312 316 324 325 327 330 331 332 334 335 <b>339</b>
<ul> <li>8.1 Creating and submitting a SQL Job</li></ul>	301 312 316 324 325 327 330 331 332 334 335 <b>339</b> <b>352</b>
<ul> <li>8.1 Creating and submitting a SQL Job</li></ul>	301 312 316 324 325 327 330 331 331 332 334 335 <b>339</b> <b>339</b> 352

10.3 APIs Supported By the DLI JDBC Driver	360
11 Submitting a Flink Job on the DLI Management Console	363
11.1 Flink Job Overview	
11.2 Creating a Flink OpenSource SQL Job	364
11.3 Creating a Flink Jar Job	
11.4 Configuring Flink Job Permissions	
11.5 Managing Flink Jobs	394
11.5.1 Viewing Flink Job Details	395
11.5.2 Setting the Priority for a Flink Job	
11.5.3 Enabling Dynamic Scaling for Flink Jobs	
11.5.4 Querying Logs for Flink Jobs	404
11.5.5 Common Operations of Flink Jobs	407
11.6 Managing Flink Job Templates	
11.7 Adding Tags to a Flink Job	
12 Submitting a Spark Job on the DLI Management Console	423
12.1 Creating a Spark Job	423
12.2 Example of a Typical Scenario: Reading and Querying OBS Data Using a Spark Jar Job	433
12.3 Setting the Priority for a Spark Job	447
12.4 Querying Logs for Spark Jobs	448
12.5 Managing Spark Jobs	450
12.6 Managing Spark Job Templates	451
13 Developing a DLI Spark Job in DataArts Studio	453
14 Submitting a Spark Job Using a Notebook Instance	460
15 Submitting a Spark Jar Job Using Livy	475
16 Monitoring DLI Using Cloud Eye	481
17 Using CTS to Audit DLI	493
18 Permissions Management	497
18.1 Overview	497
18.2 Creating a Custom Policy	502
18.3 DLI Resources	508
18.4 DLI Request Conditions	509
18.5 Common Operations Supported by DLI System Policy	
19 Common DLI Management Operations	515
19.1 Enhancing the Job Runtime Environment Using a Custom Image	515
19.2 Managing DLI Global Variables	520
19.3 Managing Program Packages of Jar Jobs	522
19.3.1 Package Management Overview	522
19.3.2 Creating a DLI Package	524
19.3.3 Configuring DLI Package Permissions	527

19.3.4 Changing the DLI Package Owner	530
19.3.5 Managing DLI Package Tags	531
19.3.6 DLI Built-in Dependencies	532
19.4 Managing DLI Resource Quotas	. 560

## DLI Job Development Process

This chapter walks you through on how to develop a DLI job.

#### **Creating an IAM User and Granting Permissions**

- To manage fine-grained permissions for your DLI resources using IAM, create an IAM user and grant them permissions to DLI if you are an enterprise user. For details, see **Creating an IAM User and Granting Permissions**.
- When using DLI for the first time, you need to update the DLI agency according to the console's guidance so that DLI can use other cloud services and perform resource O&M operations on your behalf. The agency includes permissions to obtain IAM user information, access and use VPCs, CIDR blocks, routes, and peering connections, and send notifications via SMN in case of job execution failure.

For more information on the specific permissions included in the agency, refer to **Configuring DLI Agency Permissions**.

#### **Creating Compute Resources and Metadata Required for Running Jobs**

 Before submitting a job using DLI, you need to create an elastic resource pool and create queues within it. This will provide the necessary compute resources for running the job. For how to create an elastic resource pool and create queues within it, see Overview of DLI Elastic Resource Pools and Queues.

Alternatively, you can enhance DLI's computing environment by creating custom images. Specifically, to enhance the functions and performance of Spark and Flink jobs, you can create custom images by downloading the base images provided by DLI and adding dependencies (files, JAR files, or software) and private capabilities required for job execution. This changes the container runtime environment for the jobs.

For example, you can add a Python package or C library related to machine learning to a custom image to help you extend functions. For how to create a custom image, see **Enhancing the Job Runtime Environment Using a Custom Image**.

• DLI metadata is the basis for developing SQL and Spark jobs. Before executing a job, you need to define databases and tables based on your business scenario.

#### D NOTE

source.

Flink allows for dynamic data types, enabling the definition of data structures at runtime without the need for predefined metadata.

- Define your data structures, including data catalogs, databases, and tables. For details, see Creating a Data Directory, Database, and Table.
- Create a bucket to store temporary data generated during job running, such as job logs and job results. For details, see Configuring a DLI Job Bucket.
- Configure the permission to access metadata. For details, see
   Configuring Database Permissions on the DLI Console and
   Configuring Table Permissions on the DLI Console.

#### Importing Data to DLI

- DLI allows you to analyze and query data stored in OBS without the need to migrate it. Simply upload your data to OBS and use DLI for data analysis.
- Migrate data from various sources to DLI for central storage and processing.
   For how to migrate data to DLI, see Using CDM to Migrate Data to DLI.
- After the data migration, you can submit jobs.
  Cross-source access can reduce data duplication and latency when real-time access and processing of data from different sources is required for service

needs. The prerequisites for cross-source access are that DLI can communicate with the data source network and DLI can obtain the access credentials to the data

- Configure network connection between DLI and the data source by referring to Configuring the Network Connection Between DLI and Data Sources (Enhanced Datasource Connection).
- Manage data source credentials.
  - You can use DLI's datasource authentication to manage the authentication information for accessing a specified datasource.

This applies to SQL jobs and Flink 1.12 jobs. For details, see Using DLI Datasource Authentication to Manage Access Credentials for Data Sources.

 You can also use DEW to manage access credentials for data sources and use a custom agency to authorize DLI to access DEW.

This applies to Spark 3.3.1 or later and Flink 1.15 or later.

For details, see Using DEW to Manage Access Credentials for Data Sources and Configuring an Agency to Allow DLI to Access Other Cloud Services.

#### Submitting a Job Using DLI

• DLI offers a serverless service that integrates stream processing, batch processing, and interactive analytics. It supports various job types to meet different data processing needs.

Job Type	Description	Use Case
SQL job	This type is suitable for scenarios where standard SQL statements are used for querying. It is typically used for querying and analyzing structured data. For details, see <b>Creating and</b> <b>Submitting a SQL Job</b> .	It applies to scenarios such as data warehouse query, report generation, and online analytical processing (OLAP).
Flink job	<ul> <li>This type is specifically designed for real-time data stream processing, making it ideal for scenarios that require low latency and quick response. It is well-suited for real-time monitoring and online analysis.</li> <li>Flink OpenSource job: DLI provides standard connectors and various APIs to facilitate quick integration with other data systems. For details, see Creating a Flink OpenSource SQL Job.</li> <li>Flink Jar job: allows you to submit Flink jobs compiled into JAR files, providing greater flexibility and customization capabilities. It is suitable for complex data processing scenarios that require user-defined functions (UDFs) or specific library integration. The Flink ecosystem can be utilized to implement advanced stream processing logic and status management. For details, see Creating a Flink Jar Job.</li> </ul>	It applies to scenarios that require quick response, such as real- time data monitoring and real-time recommender systems. Flink Jar jobs are suitable for data analysis scenarios that require custom stream processing logic, complex state management, or integration with specific libraries.
Spark job	Compute jobs can be submitted through interactive sessions or batch processing. Jobs are submitted to queues created within an elastic resource pool, simplifying resource management and job scheduling. It supports multiple data sources and formats, providing rich data processing capabilities, including but not limited to SQL queries and machine learning. For details, see <b>Creating a Spark Job</b> .	It is suitable for large- scale data processing and analysis, such as machine learning training, log analysis, and large-scale data mining.

 Table 1-1
 Job types supported by DLI

• Manage program packages of Jar jobs.

DLI allows you to submit Flink or Spark jobs compiled as JAR files, which contain the necessary code and dependency information for executing the job. These files are used for specific data processing tasks such as data query, analysis, and machine learning. You can manage program packages required for jobs on the DLI console.

To submit a Spark Jar or Flink Jar job, you must first upload the program package to OBS, create a program package in DLI, and then submit the program package, data, and job parameters to run the job. For details, see Managing Program Packages of Jar Jobs.

#### **NOTE**

For Spark 3.3.1 or later and Flink 1.15 or later, when creating a Jar job, you can directly configure the program package in OBS. Program packages cannot be read from DLI.

#### Using Cloud Eye to Monitor DLI

You can query DLI monitoring metrics and alarms through Cloud Eye management console or APIs.

For example, you can monitor the resource usage and job status of a DLI queue. For details about DLI metrics, see **Monitoring DLI Using Cloud Eye**.

#### Using CTS to Audit DLI

With CTS, you can log operations related to DLI, making it easier to search, audit, and trace in the future. For the supported operations, see **Using CTS to Audit DLI**.

# **2** Preparations

## 2.1 Configuring DLI Agency Permissions

#### Scenarios of dli\_management\_agency

Before using DLI, configure DLI agency permissions. This section describes the scenarios and steps for configuring DLI agency permissions (**dli\_management\_agency**).

• If you use DLI for the first time, configure DLI agency permissions by referring to this section.

DLI needs to work with other cloud services. You must grant DLI basic operation permissions of these services so that DLI can access them and perform resource O&M operations on your behalf.

• If you are currently using the older DLI agency **dli\_admin\_agency**, follow the instructions in this section to upgrade it to the newer version, **dli\_management\_agency**.

To balance practical business needs with the risk of excessive delegation, DLI upgraded its system agency to achieve more granular control over permissions. The previous **dli\_admin\_agency** was upgraded to **dli\_management\_agency**, which includes permissions for accessing IAM user information, datasource operations, and message notifications. This effectively prevents uncontrolled permission issues related to the services associated with DLI. After the upgrade, the DLI agency is more flexible and more suitable for scenario-based agency customization for medium- and large-sized enterprises.

After agency permissions are configured, the **dli\_management\_agency** agency is generated on the **Agencies** page of the IAM console. Do not delete this default system agency. Otherwise, the permissions included in the agency will be automatically revoked. The system cannot obtain IAM user information, access network resources required by datasource connections, or access SMN to send notifications.

#### **Notes and Constraints**

- Only the tenant account or a member account of user group **admin** can authorize the service.
- DLI authorization (**dli\_management\_agency**) needs to be conducted by project. The permissions of required agencies must be updated separately in each project. This means you need to switch to the corresponding project and then update the agency by following the instructions provided in this section.

#### Updating DLI Agency Permissions (dli\_management\_agency)

- 1. In the navigation pane of the DLI console, choose **Global Configuration** > **Service Authorization**.
- 2. On the displayed page, select permissions for scenarios.

Click 🗖 on a permission card to view its detailed permission policies. **Table 2-1** describes these agencies.

Use Case	Agency	Description
Basic usage	IAM ReadOnlyAccess	To authorize IAM users who have not logged in to DLI, you need to obtain their information. So, the permissions contained in the <b>IAM ReadOnlyAccess</b> policy are required.
		<b>IAM ReadOnlyAccess</b> is a global policy. Make sure you select this policy. If you do not select it, all its permissions will become invalid in all regions, and the system cannot obtain IAM user information.
Datasource	DLI Datasource Connections Agency Access	Permissions to access and use VPCs, subnets, routes, and VPC peering connections
O&M	DLI Notification Agency Access	Permissions to send notifications through SMN when a job fails to be executed

Table 2-1	Permissions	contained	in	the dli	managemen	t agency	v ad	encv
	1 6111113516113	contraintea		the dd	_managemen	_agene	, ~~	circy

#### **NOTE**

Among the permissions contained in **dli\_management\_agency**:

- The authorization scope of the IAM ReadOnlyAccess policy covers all global service resources in all regions.
  - If you select this policy when updating a DLI agency in any region, this policy's permissions apply to the projects in all regions.
  - If you do not select this policy when updating an agency in any project, this policy's permissions will be revoked from all regions. This means that all projects cannot obtain IAM user information.
- The authorization scope of the DLI Datasource Connections Agency Access and DLI Notification Agency Access policies covers the project resources in specified regions.

These policies' permissions only apply to projects for which these policies are selected and the DLI agency permissions are updated. Projects for which these policies are not selected do not have the permissions required in datasource scenarios and the permission to send notifications using SMN.

**Example 1: Configure Permissions for DLI Usage, Datasource Connection, and O&M Scenarios in Project A** and **Example 2: Configure Permissions for DLI Usage, Datasource Connection, and O&M Scenarios in Project B** demonstrate the agency permission differences caused by updating a DLI agency for different projects in a region.

3. Select the policies to be included in **dli\_management\_agency** and click **Update**.

,	Sasic Usage
	IAM ReadOnlyAccess
	N Datasource
	DLI Datasource Connections Agency Access     Permissions to access and use VPCs, subnets, ro
	∧ 08M
	DLI Notification Agency Access
	e service authorization has succeeded, an agency named dii_management_agency on IAM will be created. Go to the agency list to view the details.
,	tes
	coo nly the tenant account or sub-accounts under User Group admin can perform authorization.
	a set defend the second deserved deserved a second

#### Figure 2-1 Updating agency permissions

- 4. View and understand the notes for updating the agency and click **OK**. The DLI agency permissions are updated.
  - The system upgrades your dli\_admin\_agency to dli\_management\_agency.
  - To maintain compatibility with existing job agency permission requirements, dli\_admin\_agency will still be listed in the IAM agency list even after the update.

- Do not delete the agency created by the system by default.

#### **Follow-Up Operations**

In addition to the permissions provided by **dli\_management\_agency**, you need to create an agency on the IAM console and add information about the new agency to the job configuration for scenarios like allowing DLI to read and write data from and to OBS to transfer logs, or allowing DLI to access DEW to obtain data access credentials. For details, see **Creating a Custom DLI Agency** and **Agency Permission Policies in Common Scenarios**.

- When Flink 1.15, Spark 3.3.1 (Spark general queue scenario), or a later version is used to execute jobs, you need to create an agency on the IAM console.
- If the engine version is earlier than Flink 1.15, **dli\_admin\_agency** is used by default during job execution. If the engine version is earlier than Spark 3.3.1, user authentication information (AK/SK and security token) is used during job execution.

This means that jobs whose engine versions are earlier than Flink 1.15 or Spark 3.3.1 are not affected by the update of agency permissions and do not require custom agencies.

#### Common service scenarios where you need to create an agency:

- Data cleanup agency required for clearing data according to the lifecycle of a table and clearing lakehouse table data. You need to create a DLI agency named dli\_data\_clean\_agency on IAM and grant permissions to it. You need to create an agency and customize permissions for it. However, the agency name is fixed to dli\_data\_clean\_agency.
- Tenant Administrator permissions are required to access data from OBS to execute Flink jobs on DLI, for example, obtaining OBS data sources, log dump (including bucket authorization), checkpointing enabling, and job import and export.
- The AK/SK required by DLI Flink jobs is stored in DEW. To allow DLI to access DEW data during job execution, you need to create an agency to delegate the permissions to operate on DEW data to DLI.
- To allow DLI to access DLI catalogs to retrieve metadata when executing jobs, you need to create a new agency that grants DLI catalog data operation permissions to DLI. This will enable DLI to access DLI catalogs on your behalf.
- Cloud data required by DLI Flink jobs is stored in LakeFormation. To allow DLI to access catalogs to retrieve metadata during job execution, you need to create an agency to delegate the permissions to operate on catalog data to DLI.

When creating an agency, you cannot use the default agency names **dli\_admin\_agency**, **dli\_management\_agency**, or **dli\_data\_clean\_agency**. It must be unique.

For more information about custom agency operations, see **Creating a Custom DLI Agency** and **Agency Permission Policies in Common Scenarios**.

## Example 1: Configure Permissions for DLI Usage, Datasource Connection, and O&M Scenarios in Project A

- **Operation instruction**: A DLI user upgrades **dli\_admin\_agency** to **dli\_management\_agency** for project A in **CN North-Beijing4**.
  - a. On the DLI management console, switch to project A in the **CN North-Beijing4** region and choose **Global Configuration** > **Service Authorization**.
  - b. Select the policies under **Basic Usage**, **Datasource**, and **O&M**.

Figure 2-2 Updating a DLI agency for project A in CN North-Beijing4

Basic Usage		
IAM ReadOnlyAccess Permissions to obtain IAM user information	۵ ۲	
Datasource		
DLI Datasource Connections Agency Ac Permissions to access and use VPCs, subnets	ess वि	
<ul> <li>○ 0&amp;M</li> </ul>		

- c. Click **Update**.
- **Permission description**: The agency permissions are updated for project A in **CN North-Beijing4**.
  - The IAM ReadOnlyAccess policy's permissions are granted to global service resources, meaning that all regions and projects have these permissions.
  - The DLI Datasource Connections Agency Access and DLI Notification Agency Access policies contain only regional permissions, meaning that their permissions only apply to project A in CN North-Beijing4.

Example of Agency Permissions for Project A in CN North- Beijing4	Example of Agency Permissions for Project B in CN North- Beijing4
dli_management_agency	dli_management_agency
The new agency contains the following policy permissions:	The new agency contains the following policy permissions:
IAM ReadOnlyAccess	<ul> <li>IAM ReadOnlyAccess</li> </ul>
DLI Datasource Connections     Agency Access	
DLI Notification Agency Access	

## Example 2: Configure Permissions for DLI Usage, Datasource Connection, and O&M Scenarios in Project B

**Operation instruction**: To assign the permissions of the **DLI Datasource Connections Agency Access** and **DLI Notification Agency Access** policies to project B in the **CN North-Beijing4** region, perform the following operations to update the agency permissions of project B in **CN North-Beijing4**:

- On the DLI management console, switch to project B in the CN North-Beijing4 region and choose Global Configuration > Service Authorization.
- 2. Select the policies under **Basic Usage**, **Datasource**, and **O&M**.

#### 

When updating the agency permissions for project B, you will need to select the **IAM ReadOnlyAccess** policy, as its authorization scope covers global service resources. If you deselect it and update the agency, all its permissions will become invalid in all regions and projects.



nagement-related Agency Settings	(Agency Name: dli_management_agency)
Basic Usage	
IAM ReadOnlyAccess Permissions to obtain IAM user information	۵
Datasource	
DLI Datasource Connections Agency Access Permissions to access and use VPCs, subnets, ro	
08M	
DLI Notification Agency Access	۵

#### 3. Click **Update**.

#### Permission description:

The authorization scope of the **DLI Datasource Connections Agency Access** and **DLI Notification Agency Access** policies covers the project resources in specified regions. Following updates to agency permissions, project B in **CN North-Beijing4** has the permission to obtain IAM user information, perform datasource operations, and send message notifications.

Example of Agency Permissions in Region A	Example of Agency Permissions in Region B	
dli_management_agency	dli_management_agency	
The new agency contains the following policy permissions:	The new agency contains the following policy permissions:	
IAM ReadOnlyAccess	IAM ReadOnlyAccess	
DLI Datasource Connections Agency Access	DLI Datasource Connections Agency Access	
DLI Notification Agency Access	• DLI Notification Agency Access	

## 2.2 Creating an IAM User and Granting Permissions

You can use Identity and Access Management (IAM) to implement fine-grained permissions control on DLI resources. For details, see **Overview**.

If your cloud account does not need individual IAM users, then you may skip over this section.

This section describes how to create an IAM user and grant DLI permissions to the user. **Figure 2-4** shows the procedure.

#### Prerequisites

Before assigning permissions to user groups, you should learn about system policies and select the policies based on service requirements. For details about system permissions supported by DLI, see **DLI System Permissions**.

#### **Process Flow**



Figure 2-4 Process for granting DLI permissions

#### 1. Create a user group and grant permissions to it.

Create a user group on the IAM console, and assign the **DLI ReadOnlyAccess** permission to the group.

#### 2. Create an IAM user.

Create a user on the IAM console and add the user to the group created in 1.

3. Log in and verify permissions.

Log in to the management console using the newly created user, and verify the user permissions.

- Choose Service List > Data Lake Insight. The DLI management console is displayed. If you can view the queue list on the Queue Management page but cannot buy DLI queues by clicking Buy Queue in the upper right corner (assume that the current permission contains only DLI ReadOnlyAccess), the DLI ReadOnlyAccess permission has taken effect.
- Choose any other service in Service List. If a message appears indicating that you have insufficient permissions to access the service, the DLI ReadOnlyAccess permission has already taken effect.

#### More

- For how to create an IAM user, see Creating an IAM User.
- For how to create a custom policy, see **Creating a Custom Policy**.
- For how to modify a user policy, see **Modifying or Deleting a Custom Policy**.

## 2.3 Configuring a DLI Job Bucket

Before using DLI, you need to configure a DLI job bucket. The bucket is used to store temporary data generated during DLI job running, such as job logs and results.

Configure a DLI job bucket on the **Global Configuration** > **Project** page of the DLI management console.

#### **NOTE**

If you have enabled the function to save job results to a DLI job bucket for your SQL queue, make sure to configure the DLI job bucket before submitting SQL jobs. Failure to do so may result in SQL jobs not being submitted successfully. For details, refer to **How Do I Check if Job Result Saving to a DLI Job Bucket Is Enabled for a SQL Queue?** 

#### Preparations

Before the configuration, create an **OBS bucket** or **parallel file system (PFS)**.

In big data scenarios, you are advised to create a PFS. PFS is a high-performance file system provided by OBS, with access latency in milliseconds. PFS can achieve a bandwidth performance of up to TB/s and millions of IOPS, which makes it ideal for processing high-performance computing (HPC) workloads.

For details about PFS, see "Parallel File System Feature Guide" in the *Object Storage Service User Guide*.

#### Notes

- Do not use the OBS bucket for other purposes.
- The OBS bucket must be set and modified by the main account. Member users do not have the permission.
- If the bucket is not configured, you will not be able to view job logs.
- You can create lifecycle rules to automatically delete objects or change storage classes for objects that meet specified conditions.
- Inappropriate modifications of the job bucket may lead to loss of historical data.

#### Procedure

- In the navigation pane of the DLI console, choose Global Configuration > Project.
- 2. On the **Project** page, click and next to **Job Bucket** to configure bucket information.

#### Figure 2-5 Project

	Data Lake Insight	Project
 &	Overview SOL Editor	Job Bucket
$\bigcirc$	SQL EURO	Bucket Name: dli-cn-north-7-330e068af1334c9782f4226acc00a2e2
ത	Job Management 🔻	This bucket is used to store temporary data generated by DLI, such as job logs and job results. Do not use this bucket for other purposes. If you do not
0	Resources 🔻	create this bucket, you will not be able to view job logs. You can use the main account to set and modify the bucket. Sub-users do not have modification
$\odot$	Data Management 👻	permissions. You can set a lifecycle rule to periodically delete objects in a bucket or change its storage class. Exercise caution when you modify the lifecycle rule prevent historical data being deleted by microke.
	Job Templates 🔻	incycle fulle to prevent instorical data being detecto by finstake.
٢	Datasource Connections	
(BS)	Global Configuration	
6	Global Variables	
	Project	•
	Service	
	Authorization	
	Intelligent Tuning	

- 3. Click  $\bigcirc$  to view available buckets.
- 4. Select the bucket for storing the temporary data of the DLI job and click **OK**. Temporary data generated during DLI job running will be stored in the OBS bucket.

#### Figure 2-6 Setting the job bucket

Set Job Bu	ıcket	×
★ Job Bucket	This path must start with "obs://".	B.
	OK	

# **3** Creating an Elastic Resource Pool and Queues Within It

## 3.1 Overview of DLI Elastic Resource Pools and Queues

DLI compute resources are the foundation to run jobs. This section describes the modes of DLI compute resources and queue types.

#### What Are Elastic Resource Pools and Queues?

Before we dive into the compute resource modes of DLI, let us first understand the basic concepts of elastic resource pools and queues.

• An **elastic resource pool** is a pooled management mode for DLI compute resources, which can be seen as a collection of DLI compute resources. DLI supports the creation of multiple queues within an elastic resource pool, and these queues can share the resources in the elastic resource pool.

SQL jobs	Spark jobs	Flink jobs
Elast Unified resource management Tenant isolation	ic resource pool Time-based Job-I scaling policy resource.	evel Auto scaling
Kubernetes cluster 1		Kubernetes cluster N
Node 1 Executer pod 1 Container Container	Node N Pod 1N Container	Pod 1N Container
x86 server ARM64 server	GPU-accelerated server	FPGA-accelerated server
Huawei Cloud networking and storage se	ervices (EVS, OBS, SFS, V	PC, ELB, NAT, and more)

Figure 3-1 Elastic resource pool architecture

Elastic resource pools have the following advantages:

- Unified resource management
  - You can manage multiple internal clusters and schedule jobs in a unified manner. The scale of compute resources can reach million vCPUs.
  - Elastic resource pools can be deployed across multiple AZs to support cross-AZ high availability.
- Tenant resource isolation

Resources of different queues are isolated to reduce the impact on each other.

- Time-based on-demand elasticity
  - Minute-level scaling to cope with traffic peaks and resource requirements.
  - You can queue priorities and CU quotas at different times to improve resource utilization.
- Job-level resource isolation (not implemented currently and will be supported in later versions)

You can run SQL jobs on independent Spark instances, reducing mutual impacts between jobs.

 Automatic scaling (not implemented currently and will be supported in later versions)

Queue quotas are automatically updated in real time based on queue loads and priorities.

The advantages of using elastic resource pools include:

Di me nsi on	No Elastic Resource Pool	Elastic Resource Pool
Exp ans ion dur ati on	You will need to spend several minutes manually scaling out.	No manual intervention is required, as dynamic scale out can be done in seconds.
ResResources cannot be shared among queues.ceFor example, if queue 1 has 10 unused CUs and queue 2 requires more resources due to heavy load, queue 2 cannot utilize the resources of queue 1. It has to be scaled up.		Multiple queues added to the same elastic resource pool can share CU resources, enhancing resource utilization.
	When creating a datasource connection, you need to assign non-overlapping network segments to each queue, consuming a significant amount of VPC network segments.	You can centrally assign network segments to multiple queues in an elastic resource pool, thereby simplifying datasource configuration.
Res our ce all oca tio n	You cannot set priorities when scaling out multiple queues concurrently. If there are insufficient resources, some queues will fail to be scaled out.	You can set the priority for each queue in an elastic resource pool based on the peak and off-peak hours of the current service to ensure reasonable resource allocation.

• **Queues** are the basic units of compute resources that are actually used and allocated in DLI. You can create different queues for different jobs or data processing tasks, and allocate and adjust resources for these queues as needed. To learn more about the queue types in DLI, refer to **DLI Queue Types**.

#### **NOTE**

DLI elastic resource pools are physically isolated, while queues within the same elastic resource pool are logically isolated.

You are advised to create separate elastic resource pools for testing and production scenarios to ensure the independence and security of resource management through physical isolation.

#### **DLI Compute Resource Modes**

DLI offers three compute resource management modes, each with unique advantages and use cases.

#### Figure 3-2 DLI compute resource modes



- Elastic resource pool mode: a pooled management mode for compute resources that provides dynamic scaling capabilities. Queues within the same elastic resource pool share compute resources. By setting up a reasonable compute resource allocation policy for queues, you can improve compute resource utilization and meet resource demands during peak hours.
  - Use cases: suitable for scenarios with significant fluctuations in business volume, such as periodic data batch processing tasks or real-time data processing needs.
  - Supported queue types: for SQL (Spark), for SQL (HetuEngine), and for general purpose. To learn more about the queue types in DLI, refer to DLI Queue Types.

D NOTE

General-purpose queues and SQL queues in elastic resource pool mode do not support cross-AZ deployment.

 Usage: first create an elastic resource pool, then create queues within the pool and allocate compute resources. Queues are associated with specific jobs and data processing tasks.

For how to buy an elastic resource pool and create queues within it, see **Creating an Elastic Resource Pool and Creating Queues Within It**.

#### • Global sharing mode:

Global sharing mode is a compute resource allocation mode that allocates resources based on the actual amount of data scanned in SQL queries. It does not support specifying or reserving compute resources.

The **default** queue, which is preset by DLI, is the compute resource for global sharing mode, and the resource size is allocated on demand. Users who are unsure of the data size or occasionally need to process data can use the **default** queue to run jobs.

- Use cases: suitable for testing jobs or scenarios with low resource consumption.

- Supported queue types: Only the preset **default** queue in DLI is the compute resource for global sharing mode.

The **default** queue is typically used by users who are new to DLI but it may lead to resource contention and prevent you from getting the resources you need for your jobs, as its resources are shared among all users. You are advised to use self-built queues to run production jobs.

Usage: The default queue is only applicable to submitting SQL jobs.
 When submitting SQL jobs on the DLI management console, select the default queue.

#### • Non-elastic resource pool mode (discarded and not recommended):

The previous-gen of DLI's compute resource management mode is no longer recommended due to its lack of flexibility.

Non-elastic resource pool mode provides fixed-specification compute resources that are purchased and exclusively used, and cannot be dynamically adjusted according to demand, which may result in resource waste or insufficient resources during peak hours.

DLI Compute Resource Mode	Supported Queue Type	Resource Characteristic	Use Case
Elastic resource pool mode	For SQL (Spark) For SQL (HetuEngine ) For general purpose	Resources are shared among multiple queues for a single user. Resources are dynamically allocated and can be flexibly adjusted.	Suitable for scenarios with significant fluctuations in business demand, where resources need to be flexibly adjusted to meet peak and off-peak demands.
Global sharing mode	default queue	Resources are shared among multiple queues for multiple users. You are billed on a pay-per-use basis. Resources cannot be reserved.	Suitable for temporary or testing projects where data size is uncertain or data processing is only required occasionally.
Non-elastic resource pool mode (discarded, not recommend ed)	For SQL For general purpose	Resources are exclusively used by a single user and a single queue. Resources cannot be dynamically adjusted and may remain idle.	Discarded and not recommended.

Table 3-1 DLI of	compute resource	modes and s	supported	queue types

To help you understand the use cases for different DLI compute resource modes, we can compare purchasing DLI compute resources to using car services:

• The elastic resource pool mode can be compared to "renting a car" where you can dynamically adjust the scale of resources based on actual needs.

This mode is suitable for scenarios with significant fluctuations in business demand, allowing for flexible adjustment of resources based on peak and offpeak demands to optimize costs.

• The global sharing mode can be compared to "taking a taxi" where you only pay for the actual amount of data used.

This mode is suitable for scenarios where data size is uncertain or data processing is only required occasionally, allowing for on-demand use of resources without the need to pre-purchase or reserve resources.

#### **Elastic Resource Pool Scaling**

Creating or deleting queues within an elastic resource pool triggers elastic resource scaling.

Scaling in an elastic resource pool may affect nodes containing shuffle data, leading to the recomputation of Spark tasks. This causes automatic retries for Spark and SQL jobs, and if the retries exceed the limit, the job execution fails, requiring you to rerun the job.

#### **NOTE**

- Spark 2.3 jobs need to be upgraded to a later Spark version to support dynamic scale-in of the jobs while they are running.
- Spark Streaming and Flink jobs cannot be scaled in while they are running. To perform a scale-in, suspend the jobs or migrate them to another elastic resource pool.

#### DLI Queue Types

DLI is divided into three queue types: **default** queue, for SQL, and for general purpose. You can choose the most suitable queue type based on your business scenario and job characteristics.

• default queue:

The **default** queue is a preset queue that is shared among all users.

The **default** queue does not support specifying the size of resources and resources are allocated on-demand during job execution, with billing based on the actual amount of data scanned.

As resources of the **default** queue are shared among all users, there may be resource contention during use, and it cannot be guaranteed that resources will be available for every job execution.

The **default** queue is suitable for small-scale or temporary data processing needs. For important jobs or jobs that require guaranteed resources, you are advised to buy an elastic resource pool and create queues within it to execute jobs.

• For SQL:

For SQL queues are used to execute SQL jobs and supports specifying engine types including Spark and HetuEngine.

This type of queues is suitable for businesses that require fast data query and analysis, as well as regular cache clearing or environment resetting.

#### For general purpose:

For general purpose queues are used to execute Spark jobs, Flink OpenSource SQL jobs, and Flink Jar jobs.

It is suitable for complex data processing, real-time data stream processing, or batch data processing scenarios.

#### **Use Cases of Elastic Resource Pools**

Queues in an elastic resource pool are recommended, as they offer the flexibility to use resources with high utilization as needed. This part describes common use cases of elastic resource pools.

#### Resources too fixed to meet a range of requirements.

The quantities of compute resources required for jobs change in different time of a day. If the resources cannot be scaled based on service requirements, they may be wasted or insufficient. **Figure 3-3** shows the resource usage during a day.

- After ETL jobs are complete, no other jobs are running during 04:00 to 07:00 in the early morning. The resources could be released at that time.
- From 09:00 to 12:00 a.m. and 02:00 to 04:00 p.m., a large number of ETL report and job queries are queuing for compute resources.



#### Figure 3-3 Fixed resources

#### Resources are isolated and cannot be shared.

A company has two departments, and each run their jobs on a DLI queue. Department A is idle from 08:00 to 12:00 a.m. and has remaining resources, while department B has a large number of service requests during this period and needs more resources to meet the requirements. Since the resources are isolated and cannot be shared between department A and B, the idle resources are wasted.



#### Figure 3-4 Resource waste due to resource isolation

Elastic resource pools can be accessed by different queues and automatically scaled to improve resource utilization and handle resource peaks.

You can use elastic resource pools to centrally manage and allocate resources. Multiple queues can be bound to an elastic resource pool to share the pooled resources.

## **3.2 Creating an Elastic Resource Pool and Creating Queues Within It**

An elastic resource pool offers compute resources (CPU and memory) required for running DLI jobs, which can adapt to the changing demands of services.

You can create multiple queues within an elastic resource pool. These queues are associated with specific jobs and data processing tasks, and serve as the basic unit for resource allocation and usage within the pool. This means queues are specific compute resources required for executing jobs.

Queues within an elastic resource pool can be shared to execute jobs. This is achieved by correctly setting the queue allocation policy. This enhances queue utilization. This section describes how to create an elastic resource pool and create queues within it.

#### **NOTE**

DLI elastic resource pools are physically isolated, while queues within the same elastic resource pool are logically isolated.

You are advised to create separate elastic resource pools for testing and production scenarios to ensure the independence and security of resource management through physical isolation.

#### Notes and Constraints

Table 3-2	2 Notes	and	constraints	on	elastic	resource	pools

ltem	Description
Resource specifications	• An elastic resource pool currently supports up to 32,000 CUs.
	• Minimum CUs of a queue that can be created in an elastic resource pool:
	<ul> <li>General purpose queue: 4 CUs</li> </ul>
	<ul> <li>SQL queue: Spark SQL queue: 8 CUs; HetuEngine SQL queue: 96 CUs</li> </ul>
Billing mode	<ul> <li>Billing mode of the elastic resource pool, which can be pay-per-use or yearly/monthly.</li> </ul>
	<ul> <li>You cannot change the billing mode once the elastic resource pool is created.</li> </ul>
	• Currently, you can only change the specifications of yearly/ monthly elastic resource pools.
	• For a pay-per-use elastic resource pool, <b>Dedicated</b> <b>Resource Mode</b> is selected by default. The resource pool is billed by the calendar hour since it is created.
Managing elastic resource	<ul> <li>You cannot change the region of an elastic resource pool once the pool is created.</li> </ul>
pools	• Flink 1.10 or later jobs can run in elastic resource pools.
	<ul> <li>The CIDR block of an elastic resource pool cannot be changed once set.</li> </ul>
	• You can view only the scaling history of an elastic resource pool within 30 days.
	Elastic resource pools cannot directly access the Internet.
Associating an	Associating an elastic resource pool with a queue:
elastic resource pool with a	<ul> <li>Only pay-per-use queues (including dedicated queues) can be associated with elastic resource pools.</li> </ul>
queue	<ul> <li>No resources are frozen.</li> </ul>

Item	Description
Elastic resource pool scaling	• Changes to elastic resource pool CUs can occur when setting the CU, adding or deleting queues in an elastic resource pool, or modifying the scaling policies of queues in an elastic resource pool, or when the system automatically triggers elastic resource pool scaling. However, in some cases, the system cannot guarantee that the scaling will reach the target CUs as planned.
	<ul> <li>If there are not enough physical resources, an elastic resource pool may not be able to scale out to the desired target size.</li> </ul>
	<ul> <li>The system does not guarantee that an elastic resource pool will be scaled in to the desired target size. The system checks the resource usage before scaling in the elastic resource pool to determine if there is enough space for scaling in. If the existing resources cannot be scaled in according to the minimum scaling step, the pool may not be scaled in successfully or only partially.</li> </ul>
	The scaling step may vary depending on the resource specifications, usually 16 CUs, 32 CUs, 48 CUs, 64 CUs, etc.
	For example, if the elastic resource pool has a capacity of 192 CUs and the queues in the pool are using 68 CUs due to running jobs, the plan is to scale in to 64 CUs.
	When executing a scaling in task, the system determines that there are 124 CUs remaining and scales in by the minimum step of 64 CUs. However, the remaining 60 CUs cannot be scaled in any further. Therefore, after the elastic resource pool executes the scaling in task, its capacity is reduced to 128 CUs.

#### **Creating an Elastic Resource Pool**

- **Step 1** In the navigation pane on the left, choose **Resources** > **Resource Pool**.
- **Step 2** On the displayed page, click **Buy Resource Pool** in the upper right corner.
- **Step 3** On the displayed page, set the following parameters:

Parameter	Description
Region	Select a region where you want to buy the elastic resource pool. A region refers to the location of the physical data center of an elastic resource pool. Select a region near you to ensure the lowest latency possible.
Project	Each region corresponds to a project.

Table 3-3 Parameters

Parameter	Description
Name	Name of the elastic resource pool.
	• Only numbers, letters, and underscores (_) are allowed. The value cannot contain only numbers or start with an underscore (_) or number.
	• The value can contain a maximum of 128 characters.
	<b>NOTE</b> The elastic resource pool name is case-insensitive. Uppercase letters will be automatically converted to lowercase letters.
Туре	Basic edition: offers 16 CUs to 64 CUs
	<ul> <li>This edition is suitable for testing scenarios with low resource consumption and low requirements for resource reliability and availability.</li> </ul>
	<ul> <li>High reliability and availability are not supported.</li> </ul>
	<ul> <li>Job priorities cannot be set.</li> </ul>
	• Standard edition: offers at least 64 CUs This edition offers powerful computing capabilities, high availability, and flexible resource management. It is suitable for large-scale computing tasks and business scenarios with long-term resource planning needs.
CU Range	The maximum and minimum CUs allowed for the elastic resource pool.
	CU settings are used to control the maximum and minimum CU ranges for elastic resource pools to avoid unlimited resource scaling.
	In <b>CU Range</b> , set the minimum CUs on the left and the maximum CUs on the right.
	• The total minimum CUs of all queues in an elastic resource pool must be no more than the minimum CUs of the pool.
	• The maximum CUs of any queue in an elastic resource pool must be no more than the maximum CUs of the pool.
	An elastic resource pool should at least ensure that all queues in it can run with the minimum CUs and should try to ensure that all queues in it can run with the maximum CUs.
	The specifications (yearly/monthly CUs) of an elastic resource pool are equal to the minimum CUs allocated during creation. This means that when the elastic resource pool is first created, the actual CUs will be equal to the specifications, which is also the minimum CUs.
Description	Description of the elastic resource pool

Parameter	Description
CIDR Block	CIDR block the elastic resource pool belongs to. If you use an enhanced datasource connection, this CIDR block cannot overlap that of the data source. <b>Once set, this CIDR block</b> <b>cannot be changed.</b>
	Recommended CIDR block:
	10.0.0-10.255.0.0/16-19
	172.16.0.0–172.31.0.0/16–19
	192.168.0.0–192.168.0.0/16–19
IPv6	If IPv6 is enabled, an IPv6 CIDR block will be automatically allocated to the elastic resource pool.
	IPv6 can only be enabled during creation.
	• After enabling, both IPv4 and IPv6 addresses will be available, allowing for both private and public network access. IPv6 cannot be used independently.
	• The IPv6 CIDR block cannot be specified; it is allocated by the system.
	• After IPv6 is enabled, it cannot be disabled.
Enterprise Project	If the created elastic resource pool belongs to an enterprise project, select the enterprise project.
	Enterprise projects let you manage cloud resources and users by project.
	For how to set enterprise projects, see <b>Enterprise Management User Guide</b> .
	<b>NOTE</b> This parameter is displayed only for users who have enabled the Enterprise Management Service.

Parameter	Description
Tags	Tags used to identify cloud resources. A tag includes the tag key and tag value. If you want to use the same tag to identify multiple cloud resources, that is, to select the same tag from the drop-down list box for all services, you are advised to create predefined tags on the Tag Management Service (TMS).
	If your organization has configured tag policies for DLI, add tags to resources based on the policies. If a tag does not comply with the tag policies, resource creation may fail. Contact your organization administrator to learn more about tag policies.
	For details, see Tag Management Service User Guide.
	NOTE
	• A maximum of 20 tags can be added.
	<ul> <li>Only one tag value can be added to a tag key.</li> </ul>
	• The key name in each resource must be unique.
	• Tag key: Enter a tag key name in the text box.
	<b>NOTE</b> A tag key can contain a maximum of 128 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed, but the value cannot start or end with a space or start with <b>_sys_</b> .
	• Tag value: Enter a tag value in the text box.
	NOTE A tag value can contain a maximum of 255 characters. Only letters, numbers, spaces, and special characters ( .:+-@) are allowed.

- **Step 4** Click **Buy** and confirm the configurations.
- **Step 5** Click **Pay**. Wait until the status of the elastic resource pool changes to **Available**. The elastic resource pool is successfully created.
- **Step 6** Refer to **Example Use Case: Creating an Elastic Resource Pool and Running Jobs** and **Example Use Case: Configuring Scaling Policies for Queues in an Elastic Resource Pool** to perform subsequent operations as needed.

----End

#### **Creating Queues Within an Elastic Resource Pool**

Create one or more queues within an elastic resource pool to run jobs. This section describes how to create a queue within an elastic resource pool.

Queues added to an elastic resource pool are not billed separately, but billed as billing items of the elastic resource pool.

Creating a queue within an elastic resource pool will trigger changes of elastic resource CUs.

- **Step 1** In the navigation pane on the left, choose **Resources** > **Resource Pool**.
- **Step 2** Locate the elastic resource pool in which you want to create queues and click **Add Queue** in the **Operation** column.

**Step 3** On the **Add Queue** page, set basic queue parameters based on the table below.

Parameter	Description
Name	Name of the queue to add
Туре	• For SQL: The queue is used to run SQL jobs.
	• For general purpose: The queue is used to run Spark and Flink jobs.
Engine	If <b>Type</b> is <b>For SQL</b> , the queue engine can be <b>Spark</b> or <b>HetuEngine</b> .
	If <b>HetuEngine</b> is selected, the minimum number of CUs of the SQL queue cannot be fewer than 96 CUs.
	To use HetuEngine to submit SQL jobs, you need to configure a DLI job bucket. For details, see <b>Configuring a DLI Job Bucket</b> .
Enterprise Project	Select the enterprise project the queue belongs to. Queues under different enterprise projects can be added to an elastic resource pool.
	Enterprise projects let you manage cloud resources and users by project.
	For how to set enterprise projects, see <b>Enterprise Management User Guide</b> .
	<b>NOTE</b> This parameter is available only for users who have subscribed to the Enterprise Management Service.
Description	Description about the queue.

 Table 3-4 Basic parameters for adding a queue
Parameter	Description
Tags	Tags used to identify cloud resources. A tag includes the tag key and tag value. If you want to use the same tag to identify multiple cloud resources, that is, to select the same tag from the drop-down list box for all services, you are advised to create predefined tags on the Tag Management Service (TMS).
	If your organization has configured tag policies for DLI, add tags to resources based on the policies. If a tag does not comply with the tag policies, resource creation may fail. Contact your organization administrator to learn more about tag policies.
	For details, see Tag Management Service User Guide.
	NOTE
	• A maximum of 20 tags can be added.
	Only one tag value can be added to a tag key.
	• The key name in each resource must be unique.
	• Tag key: Enter a tag key name in the text box.
	NOTE A tag key can contain a maximum of 128 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed, but the value cannot start or end with a space or start with <b>_sys</b>
	• Tag value: Enter a tag value in the text box.
	NOTE A tag value can contain a maximum of 255 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed.

**Step 4** Click **Next**. On the displayed page, configure a scaling policy for the queue in the elastic resource pool.

**Figure 3-5** Configuring a scaling policy when adding a queue

) Basic Configuration —	Elastic Resources		
View scaling policies of	of all queues ir		
1. The priority ranges fr 2. A new policy overwri 3 You can only set the	rom 1 to 100. If you do not set the priority for a specific per tes the default policy. period to burys in (start time end time) format	riod, the default value is 1.	
4. The total minimum C 5. The maximum CUs of	CUs of all queues in an elastic resource pool cannot be more of any queue in an elastic resource pool cannot be more th	re than the minimum CUs of the pool. nan the maximum CUs of the pool.	
4.The total minimum C 5.The maximum CUs of Priority	Sus of all queues in an elastic resource pool cannot be more th of any queue in an elastic resource pool cannot be more th Period	re than the minimum CUs of the pool. aan the maximum CUs of the pool. Min CU	Max CU
4.The total minimum C 5.The maximum CUs of Priority	2Us of all queues in an elastic resource pool cannot be mo of any queue in an elastic resource pool cannot be more th Period 00	re than the minimum CUs of the pool. an the maximum CUs of the pool. Min CU - 16 +	Max CU

Click **Create** to add a scaling policy with varying priority, period, minimum CUs, and maximum CUs. The parameters of each scaling policy are:

Paramete r	Description
Priority	Priority of the scaling policy in the current elastic resource pool. A larger value indicates a higher priority. You can set a number ranging from 1 to 100.
Period	<ul> <li>Time segment when the policy takes effect. It can be set only by hour. The start time is on the left, and the end time is on the right.</li> <li>The time range includes the start time but not the end time, that is, [start time, end time). For example, if you set <b>Period</b> to <b>01</b> and <b>17</b>, the scaling policy takes effect at 01:00 a.m. till 05:00 p.m.</li> <li>The periods of scaling policies with different priorities should not overlap.</li> </ul>
Min CU	<ul> <li>Minimum number of CUs allowed by the scaling policy.</li> <li>In any time segment of a day, the total minimum CUs of all queues in an elastic resource pool cannot be more than the minimum CUs of the pool.</li> <li>If the minimum CUs of the queue is less than 16 CUs, both Max. Spark Driver Instances and Max. Prestart Spark Driver Instances set in the queue properties do not apply. Refer to Setting Queue Properties.</li> <li>For a HetuEngine SQL queue, there must be at least 96 CUs.</li> </ul>
Max CU	Maximum number of CUs allowed by the scaling policy. In any time segment of a day, the maximum CUs of any queue in an elastic resource pool cannot be more than the maximum CUs of the pool.

- The first scaling policy is the default policy, and its **Period** parameter configuration cannot be deleted or modified.
- Flink jobs cannot trigger automatic scaling of queues in an elastic resource pool.
- **Step 5** Click **OK**. View all queues and scaling policies added to the elastic resource pool by referring to **Adjusting Scaling Policies for Queues in an Elastic Resource Pool**.

----End

# 3.3 Managing Elastic Resource Pools

# 3.3.1 Viewing Basic Information

After creating an elastic resource pool, you can check and manage it on the management console.

This section describes how to check basic information about an elastic resource pool on the management console, including the VPC CIDR block, IPv6 CIDR block, and creation time.

# Procedure

- 1. Log in to the DLI management console.
- 2. In the navigation pane on the left, choose **Resources** > **Resource Pool**.
- 3. On the displayed page, select the elastic resource pool you want to check.
  - In the upper right corner of the list page, click <sup>129</sup> to customize the columns to display and set the rules for displaying the table content and the **Operation** column.
  - In the search box above the list, you can filter the required elastic resource pool by name or tag.
- 4. Click  $\cong$  to expand the basic information card of the elastic resource pool and view detailed information about the pool.

These details include the name of the pool, the user who created it, the date it was created, whether IPv6 is enabled, and the VPC CIDR block. If IPv6 is enabled, the subnet's IPv6 CIDR block will also be displayed.

For the definitions of actual CUs, used CUs, CU range, and specifications (yearly/monthly CUs) of the elastic resource pool, refer to Actual CUs, Used CUs, CU Range, and Specifications (Yearly/Monthly CUs) of an Elastic Resource Pool.

Figure 3-6 Basic information about an elastic resource pool

^	For SQL	16 CUs	16 CUs	Max: Cl Min: CL	Us Js	
Name					Engine	Spark
Default Version	2.4.5				Supported Versions	2.4.5
Max CUs	16				CPU Architecture	x86
Dedicated Resource	Yes				AZ Mode	Single AZ
CIDR Block	172.16.0.0/16				Username	dli_function
Created	Mar 06, 2023 18:	09:55 GMT+08:00				

# Actual CUs, Used CUs, CU Range, and Specifications (Yearly/Monthly CUs) of an Elastic Resource Pool

- Actual CUs: number of CUs that can be allocated in the elastic resource pool.
  - When there is no queue in the resource pool, the actual CUs are equal to the minimum CUs when the elastic resource pool is created.
  - When there are queues in the resource pool, the formula for calculating actual CUs is:
    - Actual CUs = min{sum(maximum CUs of the queue), maximum CUs of the elastic resource pool}.

- The calculation result must be a multiple of 16 CUs. If it is not exactly divisible by 16, round up to the nearest multiple.
- Example of actual CU allocation:

In **Table 3-6**, the calculation process for the actual allocation of CUs in an elastic resource pool is as follows:

- i. Calculate the sum of maximum CUs of the queues: sum(maximum CUs) = 32 + 56 = 88 CUs.
- ii. Compare the sum of maximum CUs of the queues with the maximum CUs of the elastic resource pool and take the smaller value: min{88 CUs, 112 CUs} = 88 CUs.
- iii. Check if 88 CUs is a multiple of 16 CUs. Since 88 is not divisible by 16, round up to 96 CUs.

Scenario	Resource	CU Range
New elastic resource pool: 64–112 CUs	Elastic resource pool	64–112 CUs
Queues A and B are created within the elastic	Queue A	16-32 CUs
resource pool. The CU ranges of the two queues are:	Queue B	16–56 CUs
• CU range of queue A: 16–32 CUs		
CU range of queue B: 16–56 CUs		

Table 3-6 Example of actual CU allocation of an elastic resource pool

• **Used CUs**: CUs that have been used by jobs or tasks. These CU resources may be currently engaged in computing tasks and therefore temporarily unavailable.

#### **NOTE**

The CUs used by HetuEngine match the actual CUs.

- **CU range**: CU settings are used to control the maximum and minimum CU ranges for elastic resource pools to avoid unlimited resource scaling.
  - The total minimum CUs of all queues in an elastic resource pool must be no more than the minimum CUs of the pool.
  - The maximum CUs of any queue in an elastic resource pool must be no more than the maximum CUs of the pool.
  - An elastic resource pool should at least ensure that all queues in it can run with the minimum CUs and should try to ensure that all queues in it can run with the maximum CUs.
  - When an elastic resource pool is scaled up, the minimum value of the CU range is adjusted in synchronization with the pool's specifications (yearly/ monthly CUs). After the specifications are modified, the minimum value of the CU range will be updated to match the new specification (yearly/ monthly CUs).

 Specifications (Yearly/Monthly CUs): The minimum CUs selected during elastic resource pool purchase are elastic resource pool specifications. Specifications are exclusive to yearly/monthly elastic resource pools. Billing for these specifications is yearly/monthly, while any usage exceeding the specifications is billed on a pay-per-use basis.

# **3.3.2 Managing Permissions**

Administrators can assign permissions of different operation scopes to users for each elastic resource pool.

# Precautions

- The administrator and elastic resource pool owner have all permissions, which cannot be set or modified by other users.
- When you set resource pool permissions for a new user, ensure that the region of the user group to which the user belongs has the Tenant Guest permission. For details about the Tenant Guest permission and how to apply for the permission, see Creating a User Group and Assigning Permissions and System Permissions.

# Procedure

- Step 1 In the navigation pane on the left of the DLI console, choose Resources > Resource Pool.
- Step 2 Select the desired elastic resource pool and choose More > Permissions in the Operation column. The User Permissions area displays the list of users who have permissions of elastic resource pools.

You can assign permissions to new users, modify permissions for users who already have some permissions of elastic resource pools, and revoke all permissions of a user on a pool.

• Assign permissions to a new user.

A new user does not have permissions on the elastic resource pool.

- a. Click **Set Permission** in the **Operations** column on **User Permissions** page. The **Set Permission** dialog box is displayed.
- b. Set **Username** to the name of the desired IAM user, and select the required permissions for the user.
- c. Click **OK** to.

 Table 3-7 describes the related parameters.

#### Figure 3-7 Managing permissions

Grant Pern	nission			×
* Username	Enter a username.			
Select the permis	sions to be granted to the user			
Select all				
Update		Resources	Delete	
Modify Queu	e Specifications	Grant Permission	Revoke Permission	
View Other U	Jser's Permissions			
		OK Cancel		

## Table 3-7 Parameters

Parameter	Description
Username	Name of the user you want to grant permissions to <b>NOTE</b>
	The username must be an existing IAM username and has logged in to the DLI management console.
Select the permissions to be	<ul> <li>Update: Update the description of an elastic resource pool.</li> </ul>
granted to the user	<ul> <li>Resources: Add queues, delete queues, and configure scaling policies for queues in an elastic resource pool.</li> </ul>
	<ul> <li>Delete: Delete the elastic resource pool.</li> </ul>
	<ul> <li>Modify Specifications: Change the specifications of a yearly/monthly elastic resource pool.</li> </ul>
	<ul> <li>Grant Permission: Grant the elastic resource pool permissions to other users.</li> </ul>
	<ul> <li>Revoke Permission: Revoke the permissions that other users have but you cannot revoke the owner's permissions.</li> </ul>
	<ul> <li>View Other User's Permissions: View the elastic resource pool permissions of other users.</li> </ul>

- To assign or revoke permissions of a user who has some permissions on the elastic resource pool, perform the following steps:
  - a. In the list under **User Permissions**, select the user whose permissions need to be modified and click **Set Permission** in the **Operation** column.
  - b. In the displayed **Set Permission** dialog box, modify the permissions of the user. **Table 3-7** lists the detailed permission descriptions.

If **Set Permission** is gray, you are not allowed to change permissions on this elastic resource pool. You can apply to the administrator, elastic resource pool owner, or other authorized users for granting and revoking permissions.

#### Figure 3-8 Managing permissions

Grant Perr	nission			×
* Username	testpooljwx			
Select the permis	ssions to be granted to the user			
Select all				
Update		Resources	V Delete	
Modify Queu	ue Specifications	Grant Permission	Revoke Permission	
View Other	User's Permissions			
		OK Cancel		

- c. Click **OK**.
- To revoke all permissions of a user on an elastic resource pool, perform the following steps:

In the list under **User Permissions**, select the desired user whose permissions and click **Set Permission** in the **Operation** column. Click **Yes** in the **Revoke Permission** dialog box.

----End

# 3.3.3 Binding a Queue

#### Scenario

If you want a queue to use resources in an elastic resource pool, bind the queue to the pool.

You can click **Associate Queue** on the **Resource Pool** page to bind a queue to an elastic resource pool, or bind a queue on the **Queue Management** page.

#### **NOTE**

Elastic resource pools support only Flink 1.10 or later. If jobs using Flink 1.7 run on a queue that is bound to an elastic resource pool, errors may occur due to incompatibility.

## Prerequisites

- Both the elastic resource pool and queue are available.
- The queue you want to bind must be a dedicated queue in pay-per-use billing mode.
- No resources are frozen.
- Only queues under the same enterprise project can be bound to an elastic resource pool.

## Associating a Queue

- **Step 1** In the navigation pane on the left, choose **Resources** > **Resource Pool**.
- **Step 2** Locate the row that contains the desired elastic resource pool, click **More** in the **Operation** column, and select **Associate Queue**.
- **Step 3** In the displayed dialog box, select the desired queue and click **OK**.

----End

# Allocating a Queue to an Elastic Resource Pool

- **Step 1** In the navigation pane on the left, choose **Resources** > **Queue Management**.
- **Step 2** Locate the target queue and choose **More** > **Bind Resource Pool** in the **Operation** column.
- Step 3 Select the desired elastic resource pool and click OK.

----End

# 3.3.4 Setting CUs

CU settings are used to control the maximum and minimum CU ranges for elastic resource pools to avoid unlimited resource scaling.

For example, an elastic resource pool has a maximum of 256 and two queues, and each queue must have at least 64 CUs. If you want to add another queue that needs at lest 256 CUs to the elastic resource pool, the operation is not allowed due to the maximum CUs of the elastic resource pool.

# Precautions

- In any time segment of a day, the total minimum CUs of all queues in an elastic resource pool cannot be more than the minimum CUs of the pool.
- In any time segment of a day, the maximum CUs of any queue in an elastic resource pool cannot be more than the maximum CUs of the pool.
- When you change the minimum CUs of a created elastic resource pool, ensure that the value is no more than the current CU value. Otherwise, the modification fails.

# Setting CUs

- **Step 1** In the navigation pane on the left, choose **Resources** > **Resource Pool**.
- **Step 2** Locate the row that contains the desired elastic resource pool, click **More** in the **Operation** column, and select **Set CUs**.
- **Step 3** In the **Set CUs** dialog box, set the minimum CUs on the left and the maximum CUs on the right. Click **OK**.

----End

# How Do I Increase the Minimum Value of the CU Range?

The minimum value of the CU range of an elastic resource pool is less than or equal to the actual CUs of the elastic resource pool. To adjust the minimum value of the CU range to be greater than the current CUs, increase the actual CUs first.

The following operations apply when the number of target CUs is less than or equal to the maximum value of the CU range. If the number of target CUs is greater than the maximum value of the CU range, increase the maximum CUs of the elastic resource pool first.

 For a yearly/monthly elastic resource pool, you can increase its actual CUs by adjusting the maximum CUs of queues or add queues to the pool. Afterwards, you can modify the specifications to match or exceed the target CUs. This will ensure that the actual CUs match the provided specifications. Finally, set the minimum value of the CU range to be equal to the target CUs.

• For a pay-per-use elastic resource pool, you can increase its actual CUs by adjusting the maximum CUs of queues or add queues to the pool. Afterwards, set the minimum value of the CU range to match the target CUs.

#### Example

• Example 1: For a yearly/monthly elastic resource pool, the number of actual CUs is 64, the CU range is from 64 to 96, the number of CUs provided by the specifications is 64, and the target CU range is from 80 to 96.

Procedure

a. You can increase the actual CUs of an elastic resource pool by adjusting the maximum CUs of queues in or adding queues to the elastic resource pool.

When the total number of maximum CUs of queues in the elastic resource pool is greater than its actual CUs, scale-out is triggered for the elastic resource pool. Actual CUs after scale-out = min (Total number of maximum CUs of queues, Maximum value of the CU range)

- b. Change the specifications of the elastic resource pool to 80 CUs. After the change, the minimum value of the CU range is automatically changed to the number of CUs provided by the specifications.
- Example 2: For a pay-per-use elastic resource pool, the number of actual CUs is 64, the CU range is from 64 to 96, and the target CU range is from 80 to 96. Procedure
  - a. You can increase the actual CUs of an elastic resource pool by adjusting the maximum CUs of queues in or adding queues to the elastic resource pool.

When the total number of maximum CUs of queues in the elastic resource pool is greater than its actual CUs, scale-out is triggered for the elastic resource pool. Actual CUs after scale-out = min (Total number of maximum CUs of queues, Maximum value of the CU range)

b. Adjust the CU range to the target one (80 to 96).

#### D NOTE

- The adjustment of the CU range, specifications change, and CU setting of an elastic resource pool take effect on the next hour.
- The adjustment of the actual CUs of an elastic resource pool by adding queues takes effect immediately.

# How Do I Decrease the Maximum Value of the CU Range?

The minimum value of the CU range of an elastic resource pool is less than or equal to the actual CUs of the elastic resource pool. To adjust the maximum value of the CU range to be less than the current CUs, decrease the actual CUs first.

• For a yearly/monthly elastic resource pool, you can decrease its actual CUs by adjusting the maximum CUs of queues or delete queues from the pool. Afterwards, you can modify the specifications to be equal to or less than the target CUs. This will ensure that the actual CUs match the provided

specifications. Finally, set the minimum value of the CU range to be equal to the target CUs.

• For a pay-per-use elastic resource pool, you can decrease its actual CUs by adjusting the maximum CUs of queues or delete queues from the pool. Afterwards, set the minimum value of the CU range to match the target CUs.

#### Example

• Example 1: For a yearly/monthly elastic resource pool, the number of actual CUs is 96, the CU range is from 64 to 128, the number of CUs provided by the specifications is 96, and the target CU range is from 64 to 80. Procedure

# a. You can decrease the actual CUs of an elastic resource pool by deceasing the maximum CUs of queues in or deleting queues from the elastic resource pool.

When the total number of maximum CUs of queues in the elastic resource pool is less than its actual CUs, scale-in is triggered for the elastic resource pool. Actual CUs after scale-in = min (Total number of maximum CUs of queues, Maximum value of the CU range)

- b. Change the specifications of the elastic resource pool to 80 CUs. After the change, the minimum value of the CU range is automatically changed to the number of CUs provided by the specifications.
- Example 2: For a pay-per-use elastic resource pool, the number of actual CUs is 96, the CU range is from 64 to 128, and the target CU range is from 64 to 80.

#### Procedure

a. You can decrease the actual CUs of an elastic resource pool by deceasing the maximum CUs of queues in or deleting queues from the elastic resource pool.

When the total number of maximum CUs of queues in the elastic resource pool is less than its actual CUs, scale-in is triggered for the elastic resource pool. Actual CUs after scale-in = min (Total number of maximum CUs of queues, Maximum value of the CU range)

b. Adjust the CU range to the target one (64 to 80).

#### D NOTE

- The adjustment of the CU range, specifications change, and CU setting of an elastic resource pool take effect on the next hour.
- The adjustment of the actual CUs of an elastic resource pool by adding queues takes effect immediately.

# 3.3.5 Modifying Specifications

## Scenario

You can modify the specifications of an elastic resource pool to adjust the resource configuration and billing mode based on the actual resource usage requirements of your services. This helps you efficiently use resources and optimize costs.

If the number of CUs of an yearly/monthly elastic resource pool is within the specified range (yearly/monthly CUs), then the resource pool is billed on a yearly/

monthly basis. Any excess CUs beyond the specified range (yearly/monthly CUs) are billed by CUH. You can adjust the specifications to make billing more cost-effective based on actual CU usage.

For example, if the current specification for the elastic resource pool (yearly/ monthly CUs) is 64 CUs, but most of the time the actual usage exceeds 128 CUs, the 64 CUs will be billed annually/monthly, while the excess 64 CUs will be billed by CUH. To achieve more cost-effective billing, you can change the specification of the elastic resource pool to 128 CUs. Once the specification change is successful, the entire 128 CU range will be billed yearly/monthly, resulting in a more costeffective billing.

In short, you can change the billing mode of the resources that exceed the specification (yearly/monthly CUs) from pay-per-use to yearly/monthly. This helps you efficiently use resources and optimize costs.

## Precautions

Currently, only yearly/monthly elastic resource pools support specification (yearly/ monthly CUs) changes.

## Concepts

The specification change of an elastic resource pool depends on the actual CUs of the resource pool.

- Actual CUs: currently allocated available CUs by the elastic resource pool.
  - When there is no queue in the resource pool, the actual CUs are equal to the minimum CUs when the elastic resource pool is created.
  - When there are queues in the resource pool, the formula for calculating actual CUs is:
    - Actual CUs = min{sum(maximum CUs of the queue), maximum CUs of the elastic resource pool}.
    - The calculation result must be a multiple of 16 CUs. If it cannot be exactly divided by 16 CUs, round up to the nearest multiple.
- **CU range**: CU settings mainly control the maximum and minimum CU ranges for elastic resource pool scaling to avoid unlimited resource expansion risks. When expanding the specifications of an elastic resource pool, the minimum value of the CU range is linked to the specifications (yearly/monthly CUs) of the elastic resource pool. After changing the specifications of the elastic resource pool, the minimum value of the CU range is modified to match the specifications (yearly/monthly CUs).
- **Specifications (yearly/monthly CUs)**: The minimum value of the CU range selected when purchasing an elastic resource pool is the elastic resource pool specifications. Specifications are unique to yearly/monthly elastic resource pools. The specification part is billed on a yearly/monthly basis, while parts beyond the specifications are billed on a pay-per-use basis.

# **Checking Before Specification Change (Expansion)**

Before changing the specifications (expansion), check whether the **actual CUs** are **greater than or equal to** the target CUs of the new specifications.

If the actual CUs are fewer than the target CUs, you need to increase the maximum CUs of the queue or add more queues to adjust the actual CUs.

Example: A yearly/monthly elastic resource pool has 64 actual CUs, a CU range of 64–96, and specifications of 64 CUs. The planned target specification is 80 CUs.

Procedure

1. Increase the actual CUs of the elastic resource pool to 80 by adjusting the **maxCU** of existing queues or adding new queues within the current elastic resource pool.

When the total maximum CUs of all queues exceeds the actual CUs of the elastic resource pool, the **actual CUs will increase**. The adjusted actual CUs = min(total maximum CUs of queues, maximum CU limit of the elastic resource pool). (Changes to the CU range of a queue take effect at the next full hour.)

2. Once the **actual CUs** have been increased to 80, proceed with the specification change operation to update the elastic resource pool's specifications to 80 CUs. (Elastic resource pool specification changes take effect at the next full hour.)

After the specification change, the minimum CU in the CU range of the elastic resource pool will also align with the actual CUs.

# Scaling Out

- 1. In the navigation pane on the left of the console, choose **Resources** > **Resource Pool**.
- 2. Locate the elastic resource pool you want to scale out, click **More** in the **Operation** column, and select **Modify Yearly/Monthly CU**.
- 3. On the **Modify Yearly/Monthly CU** page, set **Operation** to **Scale-out** and specify the number of CUs you want to add.

#### Figure 3-9 Scaling out

< Modify Yearly/Monthly CU

Name	
Billing Mode	Yearly/Monthly
Everine d	D 02, 0000 00:50:50 OMT: 00:00
Expired	Dec 03, 2026 23.59.59 GMT+08.00
Before Modification	64 CUs
Operation	Scale-out Scale-in
Amount (CUs)	- <b>1</b> 6 +
After Modification	80 CH s
Alter Moullication	00 003

- 4. Confirm the changes and click **OK**.
- Choose Job Management > SQL Jobs to view the status of the SCALE\_POOL SQL job.

If the job status is **Scaling**, the elastic resource pool is scaling out. Wait until the job status changes to **Finished**.

# **Scaling In**

#### **NOTE**

By default, the minimum number of CUs is **16**. That is, when the specifications of an elastic resource pool are **16 CUs**, you cannot scale the pool down.

- 1. In the navigation pane on the left, choose **Resources** > **Resource Pool**.
- 2. Locate the elastic resource pool you want to scale in, click **More** in the **Operation** column, and select **Modify Yearly/Monthly CU**.
- 3. On the **Modify Yearly/Monthly CU** page, set **Operation** to **Scale-in** and specify the number of CUs you want to decrease.

#### Figure 3-10 Scaling in

<   Modify Year	ly/Monthly CU
Name	
Billing Mode	Yearly/Monthly
Expired	Dec 03, 2026 23:59:59 GMT+08:00
Before Modification	64 CUs
Operation	Scale-out Scale-in
Amount (CUs)	_   16   <b>+</b>
After Modification	48 CUs

- 4. Confirm the changes and click **OK**.
- 5. Choose **Job Management** > **SQL Jobs** to view the status of the SCALE\_POOL SQL job.

If the job status is **Scaling**, the elastic resource pool is scaling in. Wait until the job status changes to **Finished**.

# 3.3.6 Managing Tags

## **Tag Management**

A tag is a key-value pair that you can customize to identify cloud resources. It helps you to classify and search for cloud resources. A tag consists of a tag key and a tag value.

If you use tags in other cloud services, you are advised to create the same tag (key-value pairs) for cloud resources used by the same business to keep consistency.

DLI supports the following two types of tags:

- Resource tags: non-global tags created on DLI
- Predefined tags: global tags created on Tag Management Service (TMS).

For more information about predefined tags, see **Tag Management Service User Guide**.

If your organization has configured tag policies for DLI, add tags to resources based on the policies. If a tag does not comply with the tag policies, resource creation may fail. Contact your organization administrator to learn more about tag policies.

 $\times$ 

DLI allows you to add, modify, or delete tags for queues.

- **Step 1** In the left navigation pane of the DLI console, choose **Resources** > **Resource Pool**.
- **Step 2** In the **Operation** column of the queue, choose **More** > **Tags**.
- **Step 3** The tag management page is displayed, showing the tag information about the current queue.
- **Step 4** Click **Add/Edit Tag**. The **Add/Edit Tag** dialog is displayed. Enter a tag and a value, and click **Add**.

Figure 3-11 Adding/Editing tags

#### Add/Edit Tag

It is recommended that you use TMS's predefined tag function to add the same tag to different cloud resources. View predefined tags C

To add a tag, enter a tag key and a tag value below.

Enter a tag value	Add
Entor a tag tatao	
	Enter a tag value

Cancel

ок

Parame ter	Description					
Tag key	You can specify the tag key in either of the following ways:					
	• Click the text box and select a predefined tag key from the drop- down list.					
	To add a predefined tag, you need to create one on TMS and then select it from the <b>Tag key</b> drop-down list. You can click <b>View</b> <b>predefined tags</b> to go to the <b>Predefined Tags</b> page of the TMS console. Then, click <b>Create Tag</b> in the upper corner of the page to create a predefined tag.					
	For details, see <b>Creating Predefined Tags</b> in <i>Tag Management</i> Service User Guide.					
	• Enter a tag key in the text box.					
	<b>NOTE</b> A tag key can contain a maximum of 128 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed, but the value cannot start or end with a space or start with <b>_sys_</b> .					
Tag	You can specify the tag value in either of the following ways:					
value	• Click the text box and select a predefined tag value from the drop- down list.					
	• Enter a tag value in the text box.					
	<b>NOTE</b> A tag value can contain a maximum of 255 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed.					

#### Table 3-8 Tag parameters

## **NOTE**

- A maximum of 20 tags can be added.
- Only one tag value can be added to a tag key.
- The key name in each resource must be unique.

#### Step 5 Click OK.

**Step 6** (Optional) To delete a tag, locate the row where the tag locates in the tag list and click **Delete** in the **Operation** column to delete the tag.

#### ----End

# **3.3.7 Adjusting Scaling Policies for Queues in an Elastic Resource Pool**

Multiple queues can be added to an elastic resource pool. For how to add a queue, see **Creating an Elastic Resource Pool and Creating Queues Within It**. You can configure the number of CUs you want based on the compute resources used by DLI queues during peaks and troughs and set priorities for the scaling policies to ensure stable running of jobs.

# Precautions

• You are advised to implement fine-grained management of resource pools for stream and batch processing jobs by placing Flink real-time stream jobs and SQL batch processing jobs in separate elastic resource pools.

Flink real-time stream jobs can run stably without forced scale-in, thus avoiding job interruption and system instability.

SQL batch processing jobs are placed in independent resource pools, which can scale out and in more flexibly, significantly enhancing the success rate and operational efficiency of scaling operations.

- In any time segment of a day, the total minimum CUs of all queues in an elastic resource pool cannot be more than the minimum CUs of the pool.
- In any time segment of a day, the maximum CUs of any queue in an elastic resource pool cannot be more than the maximum CUs of the pool.
- The periods of scaling policies cannot overlap.
- The period of a scaling policy can only be set by hour and specified by the start time and end time. For example, if you set the period to **00-09**, the time range when the policy takes effect is [00:00, 09:00). The period of the default scaling policy cannot be modified.
- In any period, compute resources are preferentially allocated to meet the minimum number of CUs of all queues. The remaining CUs (total CUs of the elastic resource pool – total minimum CUs of all queues) are allocated in accordance with the scaling policy priorities.
- After the queue is scaled out, the system starts billing you for the added CUs. So, if you do not have sufficient requirements, scale in your queue to release unnecessary CUs to save cost.

· · · ·	-
Scenario	CUs
An elastic resource pool has a maximum number of 256 CUs for queue A and queue B. The scaling policies are as follows:	<ul> <li>From 00:00 a.m. to 09:00 a.m.:</li> <li>1. The minimum CUs are allocated to the two queues. Queue A has 32 CUs, and queue B has 64 CUs. There are 160 CUs remaining</li> </ul>
00:00–09:00; minimum CU: 32; maximum CU: 64	<ol> <li>The remaining CUs are allocated based on priority. Since queue B has a</li> </ol>
<ul> <li>Queue B: priority 10; time period: 00:00–09:00; minimum CU: 64; maximum CU: 128</li> </ul>	higher priority than queue A, 64 CUs will be allocated to queue B first, followed by the allocation of 32 CUs to queue A.

Table 3-9 CU allocation (without jobs)

Scenario	CUs
<ul> <li>An elastic resource pool has a maximum number of 96 CUs for queue A and queue B. The scaling policies are as follows:</li> <li>Queue A: priority 5; period: 00:00–09:00; minimum CU: 32; maximum CU: 64</li> <li>Queue B: priority 10; time period: 00:00–09:00; minimum CU: 64; maximum CU: 128</li> </ul>	<ul> <li>From 00:00 a.m. to 09:00 a.m.:</li> <li>1. The minimum CUs are allocated to the two queues. Queue A has 32 CUs, and queue B has 64 CUs. There are no remaining CUs.</li> <li>2. The allocation is complete.</li> </ul>
<ul> <li>An elastic resource pool has a maximum number of 128 CUs for queue A and queue B. The scaling policies are as follows:</li> <li>Queue A: priority 5; period: 00:00-09:00; minimum CU: 32; maximum CU: 64</li> <li>Queue B: priority 10; time period: 00:00-09:00; minimum CU: 64; maximum CU: 128</li> </ul>	<ul> <li>From 00:00 a.m. to 09:00 a.m.:</li> <li>1. The minimum CUs are allocated to the two queues. Queue A has 32 CUs, and queue B has 64 CUs. There are 32 CUs remaining.</li> <li>2. The remaining 32 CUs are preferentially allocated to queue B.</li> </ul>
<ul> <li>An elastic resource pool has a maximum number of 128 CUs for queue A and queue B. The scaling policies are as follows:</li> <li>Queue A: priority 5; period: 00:00-09:00; minimum CU: 32; maximum CU: 64</li> <li>Queue B: priority 5; time period: 00:00-09:00; minimum CU: 64; maximum CU: 128</li> </ul>	<ul> <li>From 00:00 a.m. to 09:00 a.m.:</li> <li>1. The minimum CUs are allocated to the two queues. Queue A has 32 CUs, and queue B has 64 CUs. There are 32 CUs remaining.</li> <li>2. The two queues have the same priority, the remaining 32 CUs are randomly allocated to the two queues.</li> </ul>

Scenario	Actual CUs of Elastic Resourc e Pool	CUs Alloca ted to Queue A	CUs Alloca ted to Queue B	Allocation Description
Queues A and B are added to the elastic resource pool. The scaling policies are as follows: • Oueue A:	192 CUs	64 CUs	128 CUs	If the actual CUs of the elastic resource pool are greater than or equal to the sum of the maximum CUs of the two queues, the maximum CUs are allocated to both queues.
period: 00:00– 09:00; minimum CU: 32; maximum CU: 64 • Queue B: period: 00:00– 09:00; minimum CU:	96 CUs	32 CUs	64 CUs	The elastic resource pool preferentially meets the minimum CUs of the two queues. After the minimum CUs are allocated to the two queues, no CUs are allocatable.
64; maximum CU: 128	128 CUs	32 CUs to 64 CUs	64 CUs to 96 CUs	The elastic resource pool preferentially meets the minimum CUs of the two queues. That is, 32 CUs are allocated to queue A, 64 CUs are allocated to queue B, and the remaining 32 CUs are available. The remaining CUs are allocated based on the queue load and priority. The actual CUs of the queue change within the range listed.

Table 3-10 CU allocation (with jobs)

# Managing Queues

- **Step 1** In the navigation pane on the left, choose **Resources** > **Resource Pool**.
- **Step 2** Locate the target elastic resource pool and click **Queue MGMT** in the **Operation** column. The **Queue Management** page is displayed.
- **Step 3** View the queues added to the elastic resource pool.

Parameter	Description
Name	Name of the queue to add
Туре	Queue type • For SQL • For general purpose
Period	The start and end time of the queue scaling policy. This time range includes the start time but not the end time, that is, [start time, end time).
Min CUs	Minimum number of CUs allowed by the scaling policy.
Max CUs	Maximum number of CUs allowed by the scaling policy.
Priority	Priority of the scaling policy for a queue in the elastic resource pool. The priority ranges from 1 to 100. A smaller value indicates a lower priority.
Engine	For a queue running SQL jobs, the engine is Spark. For a queue for general purpose, the engine can be Spark or Flink, but it is displayed by in this page.
Created	Time when a queue is added to the elastic resource pool
Enterprise Project	Enterprise project the queue belongs to. Queues under different enterprise projects can be added to an elastic resource pool.
Owner	User who added this queue
Operation	<ul><li>Edit: Modify or add a scaling policy.</li><li>Delete: Delete the queue.</li></ul>

 Table 3-11
 Queue parameters

Figure 3-12 Managing queues

Queue Managemen	t(test									Add Queue
									Enter a name.	QMC
Name	Type	Period	Min CUs	Max CUs	Priority (?)	Engine	Created	Enterprise Project	Owner	Operation
period_queue_1	For general pur	(00.00, 17:00) [00.00, 24:00)	16 16	16 64	1	-	May 18, 2022 16:15:37 GMT+0	default	r 0	Edit Delete
fink	For SQL	(00:00,24:00)	16	16	1	spark	May 10, 2022 10:14:50 GMT+0	cefault		Edit Delete

- **Step 4** Locate the target queue and click **Edit** in the **Operation** column.
- **Step 5** In the displayed **Queue Management** pane, perform the following operations as needed:

Cancel

Edit Queue							
Elastic Resources	The priority ran priority, time se	ges from 1 to 100. gments, max and r	If you do not set the p nin CUs. You can only	priority for a spe y set the period	ecific period, the de to hours in [start ti	fault value is 1.The new me,end time) format.	priority overrides the d
	Priority	Period			Min CU	Max CU	
	1	00	v 24		- 16	+ 64	+
	1	00	<b>•</b> 17	¥	- 16	+ 16	+ Delete
	It is recommon	lad that you use Ti	JC's predefined too f	Cri	eate	areat cloud recourses 1	for medafined tage
Tags	To add a tag, er	nter a tag key and	a tag value below.		ne same ray to uni	erent cioud resources. Y	iew predenned tags
	Enter a tao ke	ΡV	Enter a tao va	lue	Add		
	Entor d tag ta	-1					

Figure 3-13 Editing scaling policies for a queue

- Add: Click Create to add a scaling policy. Set Priority, Period, Min CU, and Max CU, and click OK.
- Modify: Modify parameters of an existing scaling policy and click OK.
- **Delete**: Locate the row that contains the scaling policy you want, click **Delete** and click **OK**.

#### D NOTE

The Priority and Period parameters must meet the following requirements:

- **Priority**: The default value is **1**. The value ranges from 1 to 100. A larger value indicates a higher priority.
- Period:
  - You can only set the period to hours in [start time,end time) format.
  - For example, if the **Period** to **01** and **17**, the scaling policy takes effect at 01:00 a.m. till 05:00 p.m.
  - The periods of scaling policies with different priorities cannot overlap.
- Max CUs and Min CUs:
  - In any time segment of a day, the total minimum CUs of all queues in an elastic resource pool cannot be more than the minimum CUs of the pool.
  - In any time segment of a day, the maximum CUs of any queue in an elastic resource pool cannot be more than the maximum CUs of the pool.
- **Step 6** After you finish the settings, click statistics icon in the upper right corner of the queue list to view all scaling policies of all queue in the elastic resource pool.

Figure 3-14 Viewing statistic graphics







Step 7 View the scaling task generated when the scaling starts. Go to Job Management > SQL Jobs and view the jobs of the SCALE\_QUEUE type.

----End

# 3.3.8 Viewing Scaling History

# Scenario

If you added a queue to or deleted one from an elastic resource pool, or you scaled an added queue, the CU quantity of the elastic resource pool may be changed. You can view historical CU changes of an elastic resource pool on the console.

# Prerequisites

Currently, you can only view the historical records generated within the last 30 days on the console.

# Viewing Scaling History of an Elastic Resource Pool

- 1. In the navigation pane on the left, choose **Resources** > **Resource Pool**.
- 2. Select the desired elastic resource pool and choose **More** > **Expansion History** in the **Operation** column.
- 3. On the displayed page, select a duration to view the CU usage.

You can view the number of CUs before and after a scaling, and the target number of CUs.

The historical records can be displayed in charts or tables. Click  $\exists$  in the upper right corner to switch the display.

For example, the scaling is abnormal according to **Viewing Scaling History**, and the **Figure 3-17** shows that the target number of CUs is 80, the original number of CUs is 64, and the scaled number of CUs is 64. The scaling fails.



Figure 3-16 Scaling history in a chart

Figure 3-17 Scaling history in a table

				Date	Aug 05, 2022 06:05:45	Aug 05, 2022 12:05:41
Max CUs	Min CUs	Target CUs	CUs before expansion	CUs after expansion	Status 🗸	Operation time
128	64	128	80	80	• fail	Aug 05, 2022 09:43:26 GMT+08:00
80	64	80	64	80	success	Aug 05, 2022 09:35:29 GMT+08:00
80	64	80	64	64	• fail	Aug 05, 2022 09:29:53 GMT+08:00
80	64	80	64	64	fail	Aug 05, 2022 09:20:40 GMT+08:00
80	64	64	80	64	success	Aug 05, 2022 09:17:45 GMT+08:00

# 3.3.9 Allocating to an Enterprise Project

You can create enterprise projects matching the organizational structure of your enterprises to centrally manage cloud resources across regions by project. Then you can create user groups and users with different permissions and add them to enterprise projects.

DLI allows you to select an enterprise project when creating an elastic resource pool. This section describes how to bind an elastic resource pool to and modify an enterprise project.

**NOTE** 

Modifying the enterprise project of an elastic resource pool will modify the enterprise projects of the queues in the elastic resource pool.

Only queues under the same enterprise project can be bound to an elastic resource pool.

## Prerequisites

Data Lake Insight

User Guide

You have logged in to the Enterprise Project Management Service console and created an enterprise project by referring to **Creating an Enterprise Project**.

# **Binding an Enterprise Project**

When creating an elastic resource pool, you can select a created enterprise project for **Enterprise Project**.

Alternatively, you can click **Create Enterprise Project** to go to the Enterprise Project Management Service console to create an enterprise project and check existing ones.

For how to create a queue, see **Creating an Elastic Resource Pool and Creating Queues Within It**.

# Modifying an Enterprise Project

You can modify the enterprise project bound to a created cluster as needed.

- 1. Log in to the DLI management console.
- 2. In the navigation pane on the left, choose **Resources** > **Resource Pool**.
- 3. In the elastic resource pool list, locate the elastic resource pool for which you want to modify the enterprise project, click **More** in the **Operation** column, and select **Allocate to Enterprise Project**.
- 4. In the **Modify Enterprise Project** dialog box displayed, select an enterprise project.

Alternatively, you can click **Create Enterprise Project** to go to the Enterprise Project Management Service console to create an enterprise project and check existing ones.

5. After the modification, click **OK** to save the enterprise project information of the elastic resource pool.

# **Related Operations**

For details about how to modify the enterprise project of a queue, see **Allocating a Queue to an Enterprise Project**.

# 3.4 Managing Queues

# 3.4.1 Viewing Basic Information About a Queue

This section walks you through how to view basic information about a queue on the management console, including the engine type and version.

# Procedure

- 1. Log in to the DLI management console.
- 2. In the navigation pane on the left, choose **Resources** > **Queue Management**.
- 3. On the displayed page, locate the queue whose basic information you want to view.
  - In the upper right corner of the list page, click <sup>(2)</sup> to customize the columns to display and set the rules for displaying the table content and the **Operation** column.

- In the search box above the list, you can filter the required queue by name or tag.
- 4. Click ^ to view queue details.

The parameters of the queue are:

- **Engine**: type of the engine that executes jobs on the queue.
- **Default Version**: default version of the execution engine or the version the system will use if no specific version is specified.
- Supported Versions: all versions supported by the execution engine. By viewing the supported versions of the queue, you can find out which versions of the execution engine can be used to process jobs on the queue.

#### Figure 3-18 Basic queue information

^	For SQL	16 CUs	16 CUs	Max: CUs Min: CUs	
Name				Engine	Spark
Default Version	2.4.5			Supported Versions	2.4.5
Max CUs	16			CPU Architecture	x86
Dedicated Resource	Yes			AZ Mode	Single AZ
CIDR Block	172.16.0.0/16			Username	dli_function
Created	Mar 06, 2023 18:0	9:55 GMT+08:00			

# 3.4.2 Queue Permission Management

Administrators and queue owners have full operation permissions on queues. They can grant operation permissions to other users based on service needs. This ensures that users can execute their jobs independently without any impact on the performance of other users' job execution. This section describes how to manage queue permissions.

## **Operation Precautions**

- The administrator and queue owner have all permissions, which cannot be set or modified by other users.
- When setting queue permissions for a new user, ensure that the region of the user group to which the user belongs has the **Tenant Guest** permission.

For details about the **Tenant Guest** permission and how to apply for the permission, see **System Permissions** and **Creating a User Group** in *Identity and Access Management User Guide*.

# Operations

- Step 1 On the top menu bar of the DLI management console, click Resources > Queue Management.
- Step 2 Select the queue to be configured and choose Manage Permissions in the Operation column. The User Permission Info area displays the list of users who have permissions on the queue.

You can grant queue permissions to new users, modify permissions for users who have some permissions on a queue, and revoke all permissions for a user on a queue.

• Grant permissions to a new user.

A new user does not have permissions on the queue.

- a. Click **Set Permission** on the right of **User Permissions** page. The **Set Permission** dialog box is displayed.
- b. Specify **Username** and select corresponding permissions.
- c. Click **OK**.

Table 3-12 describes the related parameters.

#### Figure 3-19 Queue permission granting

Grant Perm	ission		
* Username	Enter a username.		
Select the permis	ssions to be granted to the user		
Select all			
Delete Queues		Submit Jobs	Terminate Job
Grant Permission		Revoke Permission	View Other User's Permissions
Restart Que	ues	Modify Queue Specifications	
		<b>OK</b> Cancel	

Table 3-12 Parameter description

Parameter	Description
Username	Name of the authorized user.
	<b>NOTE</b> The username is an existing IAM user name and has logged in to the DLI management console.

Parameter	Description
Permission Settings	<ul> <li>Delete Queues: This permission allows you to delete the queue.</li> </ul>
	<ul> <li>Submit Jobs: This permission allows you to submit jobs using this queue.</li> </ul>
	<ul> <li>Terminate Jobs: This permission allows you to terminate jobs submitted using this queue.</li> </ul>
	<ul> <li>Grant Permission: This permission allows you to grant queue permissions to other users.</li> </ul>
	<ul> <li>Revoke Permission: This permission allows you to revoke the queue permissions that other users have but cannot revoke the queue owner's permissions.</li> </ul>
	<ul> <li>View Other User's Permissions: This permission allows you to view the queue permissions of other users.</li> </ul>
	<ul> <li>Restart Queues: This permission allows you to restart queues.</li> </ul>
	<ul> <li>Modify Queue Specifications: This permission allows you to modify queue specifications.</li> </ul>

- To grant or revoke permissions for a user who already has certain permissions on a queue, perform the following steps:
  - a. In the list under **User Permission Info** for a queue, select the user whose permissions need to be modified and click **Set Permission** in the **Operation** column.
  - b. In the displayed **Set Permission** dialog box, modify the permissions of the current user. **Table 3-12** lists the detailed permission descriptions.

If all options under **Set Permission** are gray, you are not allowed to change permissions on this queue. You can apply to the administrator, queue owner, or other authorized users for queue permission granting and revoking.

#### Figure 3-20 Setting queue permissions

Set Permission		
* Username		
Select the permissions to be granted to	the user	
Select all		
Delete Queues	Submit Jobs	Terminate Job
Grant Permission	Revoke Permission	View Other User's Permission
Restart Queues	Modify Queue Specifications	
	OK	

c. Click **OK**.

# • To revoke all permissions for a user on a queue, perform the following steps:

In the user list under **Permission Info**, select the user whose permission needs to be revoked and click **Revoke Permission** under **Operation**. In the **Revoke Permission** dialog box, click **OK**. All permissions on this queue are revoked.

----End

# 3.4.3 Allocating a Queue to an Enterprise Project

You can create enterprise projects matching the organizational structure of your enterprises to centrally manage cloud resources across regions by project. Then you can create user groups and users with different permissions and add them to enterprise projects.

DLI allows you to select an enterprise project when creating a queue. This section describes how to bind a DLI queue to and modify an enterprise project.

#### **NOTE**

Currently, enterprise projects can be modified only for queues that have not been added to elastic resource pools.

## Prerequisites

You have logged in to the Enterprise Project Management Service console and created an enterprise project by referring to **Creating an Enterprise Project**.

# **Modifying an Enterprise Project**

You can change the enterprise project associated with a queue that has been created.

- 1. Log in to the DLI management console.
- 2. In the navigation pane on the left, choose **Resources** > **Queue Management**.
- 3. In the queue list, locate the queue for which you want to modify the enterprise project, click **More** in the **Operation** column, and select **Modify Enterprise Project**.
- 4. In the **Modify Enterprise Project** dialog box displayed, select an enterprise project.

Alternatively, you can click **Create Enterprise Project** to go to the Enterprise Project Management Service console to create an enterprise project and check existing ones.

5. After the modification, click **OK** to save the enterprise project information of the queue.

# **Related Operations**

For details about how to modify the enterprise project of an elastic resource pool, see **Allocating to an Enterprise Project**.

# 3.4.4 Creating an SMN Topic

## Scenario

Once you have created an SMN topic, you can easily subscribe to it by going to the **Topic Management** > **Topics** page of the SMN console. You can choose to receive notifications via SMS or email. After the subscription is successful, if a job fails, the system automatically sends a message to the subscription endpoint you specified.

- If a job fails within 1 minute of submission, a message notification is not triggered.
- If a job fails after 1 minute of submission, the system automatically sends a message to the subscriber terminal you specified.

## Procedure

1. On the **Resources** > **Queue Management** page, click **Create SMN Topic** on the upper left side. The **Create SMN Topic** dialog box is displayed.

Figure 3-21 Creating an SMN topic

#### Create SMN Topic

After creating a message notification topic, you can subscribe to it in different ways. After the subscription succeeds, any job failure will automatically be sent to your subscription endpoints.

Select	All queues (default)	•
	Ok	Cancel

2. Select a queue and click **OK**.

## D NOTE

- You can select a single queue or all queues.
- If you create a topic for a queue and another topic for all queues, the SMN of all queues does not include the message of the single queue.
- After a message notification topic is created, you will receive a message notification only when a Spark job created on the subscription queue fails.

#### Figure 3-22 Successfully created a topic

# Topic Created

011c99a26ae84a1bb963a75e7637d3fd\_all\_dli\_topic created successfully. You can add subscription to this topic in Topic Management of SMN service.



3. Click **Topic Management** in to go to the **Topic Management** page of the SMN service.

#### Figure 3-23 Topic management page

imple Message lotification	Topi	a ()							O Feedback + Grade Topic
unboard							All projects	• Enter a name.	Q Search by Tag. (6)
opic Management .		Masse		unin (3)		Enterprise Project	Display Name	Operation	
Topics		DLI, Fink, Jinko		umann.co-earth-4050411ffa4002570286ec090ec091832.0L1,fink_info		default.		Publish Mi	mage   Add Subscription   More +
Subscriptions			39R	PRCPH	052dLaparkists.dk	default.	dL0:103e3ae80055148c0062082dLspatiant.retRution	Publish Mi	ssage   Add Subscription   More +
Message Templates			Rutepic	p11371	052cl.sporkiest222	default.	d1.0.1933e3ae800490348e80982082df.spackae4222.rot#cation	Publish Me	mage   Add Subscription   Mare +
			Cope	ant on	0121E,101300023,	default	di_0.10100etae00010110108.0098208234_seepast923_authoriae	Publick Ma	scage   Add Industrythin   Mare +
			.dli_topic	encer.	Hell10,wo000935,6	default.	d1,0465eac11554056ac7c0a55ba46a018,xxd80525,54,red8cation	Public M	mage   Add Subscription   More +
		end-text		Trens.		default.		Publish Mi	ssage   Add Subscription   More +
		And Olis-CaliBack Teple		ant yet		default.	Yesh Oles-Califfack Tepic	Publish Me	sage Add Subscription Mare +
		antiko -		ant on		default		Publick Ma	scage   Add Subscription   Mare +
		Myttepic		and the second se		default.		Public M	ssage   Add Subscription   More +

- 4. In the **Operation** column of the topic, click **Add Subscription**. Select **Protocol** to determine the subscription mode.
  - If you select SMS for the subscription protocol, you need to enter the mobile number for receiving the confirmation SMS message in the Endpoint text box.
  - If you select **Email**, enter the email address for receiving the confirmation email in the **Endpoint** text box.

#### Figure 3-24 Adding a subscription

Add Subscrip	tion	
Topic Name	0c1093e9ae800f361f48	8c0096200b2df_sparktest_dli_topic
★ Protocol	SMS	•
* Endpoint ⑦	Endpoints	Description
	+ Add Endpoint	
		OK Cancel

- 5. After you click the link in the SMS message or email, you will receive a message indicating that the subscription is successful.
- 6. Go to the **Subscriptions** page of SMN, and check that subscription status is **Confirmed**.

# 3.4.5 Managing Queue Tags

## Tag Management

A tag is a key-value pair that you can customize to identify cloud resources. It helps you to classify and search for cloud resources. A tag consists of a tag key and a tag value.

If you use tags in other cloud services, you are advised to create the same tag (key-value pairs) for cloud resources used by the same business to keep consistency.

 $\times$ 

DLI supports the following two types of tags:

- Resource tags: non-global tags created on DLI
- Predefined tags: global tags created on Tag Management Service (TMS).
   For more information about predefined tags, see Tag Management Service User Guide.

If your organization has configured tag policies for DLI, add tags to resources based on the policies. If a tag does not comply with the tag policies, resource creation may fail. Contact your organization administrator to learn more about tag policies.

DLI allows you to add, modify, or delete tags for queues.

- **Step 1** In the navigation pane of the DLI management console, choose **Resources** > **Queue Management**.
- Step 2 In the Operation column of the queue, choose More > Tags.
- **Step 3** The tag management page is displayed, showing the tag information about the current queue.
- **Step 4** Click **Add/Edit Tag** to switch to the **Add/Edit Tag** dialog box. Enter a tag and a value, and click **Add**.

Figure 3-25 Adding/Editing tags

#### Add/Edit Tag

Х

It is recommended that you use TMS's predefined tag function to add the same tag to different cloud resources. View predefined tags C

To add a tag, enter a tag key and a tag value below.

Enter a tag key	Enter a tag value	Add
10 tags available for addition.		

Cancel



OK

Parame ter	Description				
Tag key	You can specify the tag key in either of the following ways:				
	<ul> <li>Click the text box and select a predefined tag key from the drop- down list.</li> </ul>				
	To add a predefined tag, you need to create one on TMS and then select it from the <b>Tag key</b> drop-down list. You can click <b>View</b> <b>predefined tags</b> to go to the <b>Predefined Tags</b> page of the TMS console. Then, click <b>Create Tag</b> in the upper corner of the page to create a predefined tag.				
	For details, see <b>Creating Predefined Tags</b> in <i>Tag Management</i> Service User Guide.				
	• Enter a tag key in the text box.				
	NOTE				
	A tag key can contain a maximum of 128 characters. Only letters, numbers, spaces, and special characters (_:+-@) are allowed, but the value cannot start or end with a space or start with _ <b>sys</b>				
Tag	You can specify the tag value in either of the following ways:				
value	• Click the text box and select a predefined tag value from the drop- down list.				
	• Enter a tag value in the text box.				
	<b>NOTE</b> A tag value can contain a maximum of 255 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed.				

#### Table 3-13 Tag parameters

## **NOTE**

- A maximum of 20 tags can be added.
- Only one tag value can be added to a tag key.
- The key name in each resource must be unique.

#### Step 5 Click OK.

**Step 6** (Optional) To delete a tag, locate the row where the tag resides in the tag list and click **Delete** in the **Operation** column to delete the tag.

----End

# **3.4.6 Setting Queue Properties**

# Scenario

DLI allows you to set properties for queues.

You can currently set the following property parameters:

• Spark driver parameters: Set them to improve the scheduling efficiency of queues.

- Set Result Saving Policy: Set it to determine whether to save job results for queues to a DLI job bucket.
- Enable Spark Native operator optimization: Enable the Spark Native engine feature to improve Spark SQL job performance and reduce CPU and memory consumption.

This section describes how to set queue properties on the management console.

# Notes and Constraints

- Queue properties can only be set for SQL queues of the Spark engine in an elastic resource pool of the standard edition.
- You cannot set queue properties in batches.
- The constraints on different queue properties vary. For details, see Table 3-14.

Property	Phase Where You Can Set This Property	Constraint	Helpful Link
Spark driver parameter s of a queue	You can set this property only after a queue is created.	For a queue in an elastic resource pool, if the minimum number of CUs of the queue is less than 16 CUs, both <b>Max. Spark Driver Instances</b> and <b>Max. Prestart Spark Driver</b> <b>Instances</b> set in the queue properties do not apply.	Procedu re
Job result saving policy yroperty only after a queue is created.		After the <b>Set Result Saving Policy</b> is enabled, that is, job results are saved to a DLI job bucket, you must configure the DLI job bucket before submitting SQL jobs. Failure to do so may result in SQL jobs not being submitted successfully.	Procedu re

Table 3-14 Constraints on different queue properties

Property	Phase Where You Can Set This Property	Constraint	Helpful Link
Enabling Spark Native operator optimizati on	<ul> <li>Creating queues within an elastic resource pool</li> <li>Setting queue propertie s after queue creation</li> </ul>	<ul> <li>For created queues, if you change the Spark Native setting (enabled/disabled) through the DLI management console or API, you need to restart the queue for the modification to take effect.</li> <li>To enable the Spark Native engine for a queue in an elastic resource pool, the following conditions must be met simultaneously: <ul> <li>Type of an elastic resource pool: Standard</li> <li>Type of a queue: For SQL</li> <li>Spark version: Spark 3.3.1 or later</li> </ul> </li> <li>For the default queue, when Spark 3.3.1 or later is used, Spark Native is disabled by default.</li> <li>To disable Spark Native for a job, configure spark.gluten.enabled=false in the job parameters to disable Spark Native at the job level.</li> </ul>	Enabling Spark Native Operato r Optimiz ation

# Procedure

- 1. In the navigation pane of the DLI management console, choose **Resources** > **Queue Management**.
- 2. Locate the queue for which you want to set properties, click **More** in the **Operation** column, and select **Set Property**.
- 3. Go to the queue property setting page and set property parameters. **Table 3-15** describes the property parameters.

Propert y Type	Property	API Parameter	Description	Value Range
Spark driver type	Max. Spark Driver Instances	computeEn gine.maxIn stance	Maximum number of Spark drivers can be started on this queue, including the Spark driver that is prestarted and the Spark driver that runs jobs.	<ul> <li>For a 16-CU queue, the value is 2.</li> <li>For a queue that has more than 16 CUs, the value range is [2, queue CUs/16].</li> <li>If the minimum number of CUs of the queue is less than 16 CUs, this configuration item does not apply.</li> </ul>
	Max. Prestart Spark Driver Instances	computeEn gine.maxPr efetchInsta nce	Maximum number of Spark drivers can be prestarted on this queue. When the number of Spark drivers that run jobs exceeds the value of <b>Max.</b> <b>Concurrency</b> <b>per Instance</b> , the jobs are allocated to the Spark drivers that are prestarted.	<ul> <li>For a 16-CU queue, the value range is 0 to 1.</li> <li>For a queue that has more than 16 CUs, the value range is [2, queue CUs/16].</li> <li>If the minimum number of CUs of the queue is less than 16 CUs, this configuration item does not apply.</li> </ul>

Table 3-15 Queue properties

Propert y Type	Property	API Parameter	Description	Value Range
	Max. Concurrenc y per Instance	job.maxCo ncurrent	Maximum number of jobs can be concurrently executed by a Spark driver. When the number of jobs exceeds the value of this parameter, the jobs are allocated to other Spark drivers.	1-32
Propert y Type	Property	API Parameter	Description	Value Range
-------------------------	--------------------------------	--------------------------------------	---	-------------
Job result saving	Set Result Saving Policy	job.saveJob ResultToJo bBucket	Whether to save job results to a DLI job bucket.	N/A
policy			This parameter is only available for Spark SQL queues.	
			Once enabled, this feature cannot be disabled, and job results of a queue are consistently saved to the DLI job bucket you configure.	
			Before enabling this feature, make sure you have configured a DLI job bucket. For how to configure a DLI job bucket, see <b>Configuring</b> a DLI Job Bucket.	
			To check if you have enabled the function to save SQL job results to a DLI job bucket, refer to How Do I Check if Job Result Saving to a DLI Job Bucket Is Enabled for a SQL Queue?	
			You are advised to enable the feature of saving job results to DLI job buckets to more effectively	

Propert y Type	Property	API Parameter	Description	Value Range
			manage and store SQL job query results.	
Enablin g Spark Native operato r optimiz ation	DLI Spark Native Acceleratio n	computeEn gine.spark. nativeEnab led	Enables the Spark Native engine to improve Spark SQL job performance and reduce CPU and memory usage.	Enabled or disabled
			For more information, see Enabling Spark Native Operator Optimization.	

4. Click OK.

# How Do I Check if Job Result Saving to a DLI Job Bucket Is Enabled for a SQL Queue?

- Method 1: View the result path on the SQL job details page.
  - a. Log in to the DLI console. In the navigation pane on the left, choose **Job Management** > **SQL Jobs**.
  - b. Locate a desired SQL job and click  $\checkmark$  next to the queue name to expand the job details.
  - c. Check the result path in the job details.
    - If the result path displays a custom DLI job bucket, the feature of saving job results to a DLI job bucket is enabled for the queue running the job.
    - If the result path is not displayed in the job details, the feature is not enabled for the queue running the job.
- Method 2: Check whether the feature of saving job results to a job bucket is enabled in the SQL queue properties.
  - a. Log in to the DLI management console. In the navigation pane on the left, choose **Resources** > **Queue Management**.
  - b. Locate the queue for which you want to set properties, click **More** in the **Operation** column, and select **Set Property**.
  - c. On the **Set Property** dialog box that appears, check whether the feature is enabled.

# 3.4.7 Enabling Spark Native Operator Optimization

## Scenario

Spark Native is a core component of Apache Spark designed to enhance the performance of Spark SQL computations. By utilizing vectorized C++ acceleration libraries, it accelerates the performance of Spark operators. Enabling Spark Native can improve the performance of Spark SQL jobs, reducing CPU and memory consumption.

After enabling Spark Native in a queue, it currently supports optimization for Scan and Filter operators.

• Scan: The Scan operator is typically triggered by query statements, such as select \* from test\_table.

The following conditions support enabling Native:

- Hive tables and datasource tables in Parquet format
- Datasource tables in ORC format
- **Filter**: The Filter operator is typically triggered by **WHERE** clauses, such as **select \* from test\_table where id =** *xxx*.

#### D NOTE

Using the **EXPLAIN** statement, you can view the types of operators triggered by SQL commands, for example, **Explain select \* from test\_table**.

This section describes how to enable Spark Native operator optimization.

#### Notes and Constraints

- To enable the Spark Native engine for a queue in an elastic resource pool, the following conditions must be met simultaneously:
  - Type of an elastic resource pool: Standard
  - Type of a queue: **For SQL**
  - Spark version: Spark 3.3.1 or later
- For the **default** queue, when Spark 3.3.1 or later is used, Spark Native is disabled by default.
- To disable Spark Native for a job, configure **spark.gluten.enabled=false** in the job parameters to disable Spark Native at the job level.

#### **Enabling Spark Native Operator Optimization**

 When creating a SQL queue in an elastic resource pool, you can enable Spark Native.

Enable Spark Native acceleration on DLI.

For details, see **Creating an Elastic Resource Pool and Creating Queues Within It**.

- For SQL queues in an existing elastic resource pool, you can enable Spark Native by setting queue properties.
  - a. In the navigation pane on the left of the DLI management console, choose **Resources** > **Queue Management**.

- b. Locate the queue for which you want to set properties, click **More** in the **Operation** column, and select **Set Property**.
- c. Go to the queue property setting page and set property parameters. **Table 3-16** describes the property parameters.

D NOTE

For created queues, if you change the Spark Native setting (enabled/disabled) through the DLI management console or API, you need to restart the queue for the modification to take effect.

#### Table 3-16 Queue properties

Property	Description	Example Value
DLI Spark Native Acceleration	Enables the Spark Native engine to improve Spark SQL job performance and reduce CPU and memory usage.	Enabled

d. Click OK.

## **Disabling Spark Native Operator Optimization**

- Disable Spark Native for SQL queues in an elastic resource pool.
  - a. In the navigation pane on the left of the DLI management console, choose **Resources** > **Queue Management**.
  - b. Locate the queue for which you want to set properties, click **More** in the **Operation** column, and select **Set Property**.
  - c. Go to the queue property setting page and set property parameters. **Table 3-17** describes the property parameters.

Table 3-17	Queue	properties
------------	-------	------------

Property	Description	Example Value
DLI Spark Native Acceleration	Enables the Spark Native engine to improve Spark SQL job performance and reduce CPU and memory usage.	Disabled

- d. Click OK.
- Disable Spark Native for a specified job when a queue has Spark Native enabled.

After Spark Native is enabled for a SQL queue, if you want to disable Spark Native for a particular job running in the queue,

add **spark.gluten.enabled=false** to the parameter settings of the SQL job to disable Spark Native.

# 3.4.8 Testing Address Connectivity

DLI's address connectivity testing feature can be used to verify network connectivity between DLI queues and destination addresses.

This feature is typically utilized for reading and writing external data sources. Once a datasource connection is configured, the communication capability between the DLI queue and the bound peer address is verified.

# Testing the Address Connectivity Between a Queue and the Data Source

- 1. Log in to the DLI management console. In the navigation pane on the left, choose **Resources** > **Queue Management**.
- 2. On the **Queue Management** page, locate the row containing the target queue, click **More** in the **Operation** column, and select **Test Address Connectivity**.
- 3. On the **Test Address Connectivity** page, enter the address to be tested. The domain name and IP address are supported, and the port number can be specified.

You can input the data source address in the following formats: IPv4 address; IPv4 address + Port number; Domain name; Domain name + Port number.

- · IPv4 address: 192.168.x.x
- · IPv4 + Port number: 192.168.x.x:8080
- · Domain name: domain-xxxxx.com
- · Domain name + Port number: domain-xxxxx.com:8080

#### Figure 3-26 Testing address connectivity

Tests whether an address is reachable from a specified cluster. The address can be a domain name, an IP address, or a specified port.	
* Address Enter IP address or Domain	

- 4. Click Test.
  - If the test address is reachable, a message is displayed on the page, indicating that the address is reachable.
  - If the test address is unreachable, the system displays a message indicating that the address is unreachable. Check the network configurations and retry. Network configurations include the VPC peering and the datasource connection. Check whether they have been activated.

# **Related Operations**

Why Is a Datasource Connection Successfully Created But the Network Connectivity Test Fails?

# 3.4.9 Deleting a Queue

You can delete a queue based on actual conditions.

#### **NOTE**

- This operation will fail if there are jobs in the **Submitting** or **Running** state on this queue.
- Deleting a queue does not cause table data loss in your database.

#### Procedure

- Step 1 In the navigation pane on the left of the DLI management console, choose Resources > Queue Management.
- **Step 2** Locate the row where the target queue locates and click **Delete** in the **Operation** column.

Figure 3-27 Deleting a Queue

Delete		×
Are you sure you want to delete the sele g state on this queue. Collapse 🔺	ected queue? This operation will fail if the	ere are jobs in the Submitting or Runnin
Name	Specifications	Created
testsql	16 CUs	Feb 11, 2022 17:36:57 GMT+08:00
	Yes No	

#### **NOTE**

If **Delete** in the **Operation** column is gray, the current user does not have the permission of deleting the queue. You can apply to the administrator for the permission.

**Step 3** In the displayed dialog box, click **OK**.

----End

# 3.4.10 Enabling Elastic Scaling for a Queue in a Non-Elastic Resource Pool

# Prerequisites

Newly created pay-per-use queues need to run jobs before they can be scaled in or out.

#### D NOTE

The operations described in this section only apply to queues in non-elastic resource pools.

## Notes and Constraints

- Queues with 16 CUs do not support scale-out or scale-in.
- Queues with 64 CUs do not support scale-in.
- Only queues whose billing mode is Pay-per-use/By CUH and Pay-per-use/ Dedicated resource mode support elastic scaling.
- If **Status of queue xxx is assigning, which is not available** is displayed on the **Elastic Scaling** page, the queue can be scaled only after the queue resources are allocated.
- If there are not enough physical resources, a queue may not be able to scale out to the desired target size.
- The system does not guarantee that a queue will be scaled in to the desired target size. Typically, the system checks the resource usage before scaling in the queue to determine if there is enough space for scaling in. If the existing resources cannot be scaled in according to the minimum scaling step, the queue may not be scaled in successfully or only partially.

The scaling step may vary depending on the resource specifications, usually 16 CUs, 32 CUs, 48 CUs, 64 CUs, etc.

For example, if the queue size is 48 CUs and job execution uses 18 CUs, the remaining 30 CUs do not meet the requirement for scaling in by the minimum step of 32 CUs. If a scaling in task is executed, it will fail.

# **Scaling Out**

If the current queue specifications do not meet service requirements, you can add the number of CUs to scale out the queue.

#### **NOTE**

Scale-out is time-consuming. After you perform scale-out on the **Elastic Scaling** page of DLI, wait for about 10 minutes. The duration is related to the CU amount to add. After a period of time, refresh the **Queue Management** page and check whether values of **Specifications** and **Actual CUs** are the same to determine whether the scale-out is successful. Alternatively, on the **Job Management** page, check the status of the **SCALE\_QUEUE** SQL job. If the job status is **Scaling**, the queue is being scaled out.

The procedure is as follows:

- 1. In the navigation pane on the left of the DLI management console, choose **Resources > Queue Management**.
- 2. Select the queue to be scaled out, click **More** in the **Operation** column, and select **Elastic Scaling**.
- 3. On the displayed page, select **Scale-out** for **Operation** and set the scale-out amount.

#### Figure 3-28 Scale-out

Name	sparktest	
Billing Mode	Pay-per-use	
Actual CUs	16 CUs	
Opération	Scale-out	Scale-in
Amount (CUs)	- 16 +	
Final CU Count	32 CUs	

4. Click OK.

## Scaling In

If the current queue specifications are too much for your computing service, you can reduce the number of CUs to scale in the queue.

#### **NOTE**

- Scale-in is time-consuming. After you perform scale-in on the **Elastic Scaling** page of DLI, wait for about 10 minutes. The duration is related to the CU amount to reduce. After a period of time, refresh the **Queue Management** page and check whether values of **Specifications** and **Actual CUs** are the same to determine whether the scale-in is successful. Alternatively, on the **Job Management** page, check the status of the **SCALE\_QUEUE** SQL job. If the job status is **Scaling**, the queue is being scaled in.
- By default, the minimum number of CUs is **16**. That is, when the queue specifications are **16 CUs**, you cannot scale in the queue.

The procedure is as follows:

- 1. In the navigation pane on the left of the DLI management console, choose **Resources > Queue Management**.
- 2. Select the queue to be scaled in, click **More** in the **Operation** column, and select **Elastic Scaling**.
- 3. On the displayed page, select **Scale-in** for **Operation** and set the scale-in amount.

#### Figure 3-29 Manual scale-in

Name	createoutcuqueue	
Billing Mode	Pay-per-use	
Actual CUs	128 CUs	
Opération	Scale-out	Scale-in
Amount (CUs)	- 16 +	
Final CU Count	112 CUs	

4. Click OK.

# 3.4.11 Creating a Scheduled Elastic Scaling Task for a Queue in a Non-Elastic Resource Pool

## Scenario

When services are busy, you might need to use more compute resources to process services in a period. After this period, you do not require the same amount of resources. If the purchased queue specifications are small, resources may be insufficient during peak hours. If the queue specifications are large, resources may be wasted.

DLI provides scheduled tasks for elastic scale-in and -out in the preceding scenario. You can set different queue sizes (CUs) at different time or in different periods based on your service period or usage and the existing queue specifications to meet your service requirements and reduce costs.

#### **NOTE**

The operations described in this section only apply to queues in non-elastic resource pools.

# Precautions

- Newly created queues need to run jobs before they can be scaled in or out.
- Scheduled scaling tasks are available only for a queue with more than 64 CUs. That is, the minimum specifications of a queue are 64 CUs.
- A maximum of 12 scheduled tasks can be created for each queue.
- When each scheduled task starts, the actual start time of the specification change has a deviation of 5 minutes. It is recommended that the task start time be at least 20 minutes earlier than the time when the queue is actually used.
- The interval between two scheduled tasks must be at least 2 hours.
- Changing the specifications of a queue is time-consuming. The time required for changing the specifications depends on the difference between the target

specifications and the current specifications. You can view the specifications of the current queue on the **Queue Management** page.

- If a job is running in the current queue, the queue may fail to be scaled in to the target CU amount value. Instead, it will be scaled in to a value between the current queue specifications and the target specifications. The system will try to scale in again 1 hour later until the next scheduled task starts.
- If a scheduled task does not scale out or scale in to the target CU amount value, the system triggers the scaling plan again 15 minutes later until the next scheduled task starts.

## **Creating Periodic Task**

- If only scale-out or scale-in is required, you need to create only one task for changing specifications. Set the **Task Name**, **Final CU Count**, and **Executed** parameters. For details, see **Table 3-18**.
- To set both scale-out and scale-in parameters, you need to create two periodic tasks, and set the **Task Name**, **Final CU Count**, and **Executed** parameters. For details, see **Table 3-18**.

The procedure is as follows:

- 1. In the navigation pane on the left of the DLI management console, choose **Resources > Queue Management**.
- 2. Locate the queue for which you want to schedule a periodic task for elastic scaling, and choose **More** > **Schedule CU Changes** in the **Operation** column.
- 3. On the displayed page, click **Create Periodic Task** in the upper right corner.
- 4. On the Create Periodic Task page, set the required parameters. Click OK.

#### Figure 3-30 Creating a periodic task

#### Create Periodic Task

★ Task Name	Enter a name.
Enable Task	
Validity Period	Select a date and time.
Total Number	32 CUs
★ Final CU Count	— 16 + CUs
Repeat	<ul> <li>Select all</li> <li>Monday</li> <li>Tuesday</li> <li>Wednesday</li> <li>Thursday</li> <li>Friday</li> <li>Saturday</li> </ul>
* Executed	14:27 OK Cancel

#### Table 3-18 Parameters

Param eter	Description
Task Name	<ul> <li>Enter the name of the periodic task.</li> <li>Only numbers, letters, and underscores (_) are allowed. The value cannot contain only numbers or start with an underscore (_) or be left unspecified.</li> <li>The value can contain a maximum of 128 characters.</li> </ul>
Enable Task	Whether to enable periodic elastic scaling. The task is enabled by default. If disabled, the task will not be triggered on time.
Validit y Period	<ul> <li>Time segment for executing the periodic task. The options include</li> <li>Date and Time. If there is no time segment restriction, leave this parameter empty, indicating that the task takes effect permanently. If you need to specify the time segment for the task to take effect, set this parameter based on the service requirements.</li> <li>NOTE</li> <li>The start time of the Validity Period must be later than the current system time.</li> <li>If only scale-out is configured, the system does not automatically scale in after the Validity Period expires. You need to manually modify or configure a periodic scale-in task. Similarly, if only scale-out task. That is, a scheduled scaling task is executed at a time.</li> <li>If both scale-out and scale-in are configured, the system scales in or out resources based on the configured queue specifications within the validity period. After the validity period expires, the system retains the last configured queue specifications.</li> </ul>
Actual CUs	Queue specifications before scale-in or scale-out.
Final CUs	<ul> <li>Specifications after the queue is scaled in or out.</li> <li>NOTE <ul> <li>By default, the maximum specifications of a queue are 512 CUs.</li> <li>The minimum queue specifications for scheduled scaling are 64 CUs. That is, only when Actual CUs are more than 64 CUs, the scheduled scaling can be performed.</li> <li>The value of Actual CUs must be a multiple of 16.</li> </ul> </li> </ul>

Param eter	Description		
Repeat	Time when scheduled scale-out or scale-in is repeat. Scheduled tasks can be scheduled by week in <b>Repeat</b> .		
	• By default, this parameter is not configured, indicating that the task is executed only once at the time specified by <b>Executed</b> .		
	<ul> <li>If you select all, the plan is executed every day.</li> </ul>		
	• If you select some options of <b>Repeat</b> , the plan is executed once a week at all specified days.		
	NOTE		
	<ul> <li>You do not need to set this parameter if you only need to perform scale-in or scale-out once.</li> </ul>		
	<ul> <li>If you have set scaling, you can set <b>Repeat</b> as required. You can also set the repeat period together with the validity period.</li> </ul>		
Execut	Time when scheduled scale-out or scale-in is performed		
ed	• When each scheduled task starts, the actual start time of the specification change has a deviation of 5 minutes. It is recommended that the task start time be at least 20 minutes earlier than the time when the queue is actually used.		
	<ul> <li>The interval between two scheduled tasks must be at least 2 hours.</li> </ul>		

After a periodic task is created, you can view the specification change of the current queue and the latest execution time on the page for scheduling CU changes.

Alternatively, on the **Queue Management** page, check whether the **Specifications** change to determine whether the scaling is successful.

You can also go to the **Job Management** page and check the status of the **SCALE\_QUEUE** job. If the job status is **Scaling**, the queue is being scaled in or out.

# Modifying a Scheduled Task

If a periodic task cannot meet service requirements anymore, you can modify it on the **Schedule CU Changes** page.

- 1. In the navigation pane on the left of the DLI management console, choose **Resources > Queue Management**.
- Locate the queue for which you want to schedule a periodic task for elastic scaling, and choose More > Schedule CU Changes in the Operation column.
- 3. On the displayed page, click **Modify** in the **Operation** column. In the displayed dialog box, modify the task parameters as needed.

# **Deleting a Scheduled Task**

If you do not need the task anymore, delete the task on the **Schedule CU Changes** page.

- 1. In the navigation pane on the left of the DLI management console, choose **Resources > Queue Management**.
- 2. Locate the queue for which you want to schedule a periodic task for elastic scaling, and choose **More** > **Schedule CU Changes** in the **Operation** column.
- 3. On the displayed page, click **Delete** in the **Operation** column. In the displayed dialog box, click **OK**.

# 3.4.12 Changing the CIDR Block of a Queue in a Non-Elastic Resource Pool

If the CIDR block of the DLI queue conflicts with that of the user data source, you can change the CIDR block of the queue.

If the queue whose CIDR block is to be modified has jobs that are being submitted or running, or the queue has been bound to enhanced datasource connections, the CIDR block cannot be modified.

**NOTE** 

The operations described in this section only apply to queues in non-elastic resource pools.

#### Procedure

#### 

Currently, you can modify the CIDR block of for queues whose billing mode is **Yearly/Monthly** or **Pay-per-use/Dedicated resource mode**.

- 1. In the navigation pane on the left of the DLI management console, choose **Resources > Queue Management**.
- 2. Select the queue to be modified and click **Modify CIDR Block** in the **Operation** column.

Figure 3-31 Modifying a CIDR block

Modify C	IDR Block
Name	testaaa
CIDR Block	172 . 16 . 0 . 0       /       18 ▼         If you need to use DLI enhanced datasource connections, the CIDR block entered here cannot be the same as that of the data source.         Recommended CIDR blocks:         10 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
	OK Cancel

3. Enter the required CIDR block and click **OK**. After the CIDR block of the queue is successfully changed, wait for 5 to 10 minutes until the cluster to which the queue belongs is restarted and then run jobs on the queue.

Recommended CIDR block: 10.0.0.0-10.255.0.0/8-24 172.16.0.0-172.31.0.0/12-24 192.168.0.0-192.168.0.0/16-24

# **3.5 Example Use Case: Creating an Elastic Resource Pool and Running Jobs**

This section walks you through the procedure of adding a queue to an elastic resource pool and binding an enhanced datasource connection to the elastic resource pool.



Figure 3-32 Process of creating an elastic resource pool

 Table 3-19
 Procedure

Step	Description	Reference
Create an elastic resource pool	Create an elastic resource pool and configure basic information, such as the billing mode, CU range, and CIDR block.	Creating an Elastic Resource Pool and Creating Queues Within It

Step	Description	Reference
Add a queue to the elastic resource pool	<ul> <li>Add the queue where your jobs will run on to the elastic resource pool. The operations are as follows:</li> <li>1. Set basic information about the queue, such as the name and type.</li> <li>2. Configure the scaling policy of the queue, including the priority, period, and the maximum and minimum CUs allowed for scaling.</li> </ul>	Creating an Elastic Resource Pool and Creating Queues Within It Adjusting Scaling Policies for Queues in an Elastic Resource Pool
(Optional) Create an enhanced datasource connection.	If a job needs to access data from other data sources, for example, GaussDB(DWS) and RDS, you need to create a datasource connection. The created datasource connection must be bound to the elastic resource pool.	Creating an Enhanced Datasource Connection
Run a job.	Create and submit the job as you need.	Managing SQL Jobs Flink Job Overview Creating a Spark Job

# Step 1: Create an Elastic Resource Pool

- 1. Log in to the DLI management console. In the navigation pane on the left, choose **Resources** > **Resource Pool**.
- 2. On the displayed page, click **Buy Resource Pool** in the upper right corner.
- 3. On the displayed page, set the following parameters:
  - **Name**: Enter the name of the elastic resource pool. For example, **pool\_test**.
  - **CU range**: Minimum and maximum CUs of the elastic resource pool.
  - **CIDR Block**: Network segment of the elastic resource pool. For example, **172.16.0.0/18**.
  - Set other parameters as required.

	5
<   Buy Resource Po	01
Billing Mode	Yearly/Monthly Pay-per-use
Region	Resources are region-specific and cannot be used across regions through internal network connections. For low network latency and quick access, select the nearest region.
Project	C
* Name	pool_test
CU range	-         64         +         -         64         +         ∅           ✓         Dedicated Resource Mode         ⑦         ∅ </th
Description	0056
CIDR Block	172 · 16 · 0 · 0 / 18 ·
	An estimated 4096 CUs are available.
	This network segment cannot be changed after having been set. To use enhanced datasource connections, make sure that the network segment of the resource pool does not overlap with that of the datasources.
	Recommended CIDR blocks: 10.0.0.~10.255.0.0/16~19,172.16.0.0~172.31.0.0/16~19,192.168.0.0~192.168.0.0/16~19
* Enterprise Project	default   C (2) Create Enterprise Project

Figure 3-33 Creating an elastic resource pool

For details about how to create an elastic resource pool, see **Creating an Elastic Resource Pool and Creating Queues Within It**.

- 4. Click **Buy**. Confirm the configuration and click **Pay**.
- 5. Go to the **Resource Pool** page to view the creation status. If the status is **Available**, the elastic resource pool is ready for use.

# Step 2: Add a Queue to the Elastic Resource Pool

- 1. In the **Operation** column of the created elastic resource pool, click **Add Queue**.
- 2. Specify the basic information about the queue. The configuration parameters are as follows:
  - Name: Queue name
  - Type: Queue type In this example, select For general purpose.
     For SQL: The queue is used to run Spark SQL and HetuEngine jobs.
     For general purpose: The queue is used to run Flink and Spark Jar jobs.
  - Set other parameters as required.

#### Figure 3-34 Creating a queue

<   Add Queue(pool	I_test)
Basic Configuration —	(2) Elastic Resources
* Name	general_test
<b>*</b> Туре	For SQL For general purpose
* Enterprise Project	default C 🕐 Create Enterprise Project
Description	
	It is recommended that you use TMS's predefined tag function to add the same tag to different cloud resources. View predefined tags $$ C
	To add a tag, enter a tag key and a tag value below.
Tags	
	Enter a tag key Enter a tag value Add
	10 tags available for addition.

3. Click Next. On the displayed page, set Min CU to 64 and Max CU to 64.

Figure 3-35 Set scaling policy for the queue

<   Add Queue(pool_t	< Add Queue[pool_test]					
Basic Configuration	Basic Configuration — 2 Elastic Resources					
View scaling policies of all o	ueues in pool_test. o 100. If you do not set the priority for a specific period, the de	fault value is 1.The new priority overrides the default p	riority, time segments, max and min Cl	Js. You can only set the period		
to hours in [start time, end ti Priority	Period	Min CU	Max CU			
1	00 * - 24 *	- 64 +	- 64 +			
		Create				

4. Click **OK**. The queue is added.

# (Optional) Step 3: Create an Enhanced Datasource Connection

In this example, a datasource connection is required to connect to RDS. You need to create a datasource connection. If your job does not need to connect to an external data source, skip this step.

1. Log in to the RDS console and create an RDS DB instance.

For details, see **Buying an RDS for MySQL Instance**.

- 2. Click **Create Database**. In the dialog box that appears, enter database name **test2**. Then, click **OK**.
- Locate the row that contains the test2 database, click Query SQL Statements in the Operation column. On the displayed page, enter the following statement to create table tabletest2. Click Execute SQL. The table creation statement is as follows: CREATE TABLE `tabletest2` (
  - `id` int(11) unsigned,
  - `name` VARCHAR(32)
  - ) ENGINE = InnoDB DEFAULT CHARACTER SET = utf8mb4;

- 4. On the RDS console, choose **Instances** form the navigation pane. Click the name of a created RDS DB instance to view its basic information.
- 5. In the **Connection Information** pane, obtain the floating IP address, database port, VPC, and subnet.
- Click the security group name. In the Inbound Rules tab, add a rule to allow access from the CIDR block of the elastic resource pool. For example, if the CIDR block of the elastic resource pool is 172.16.0.0/18 and the database port is 3306, set the rule Priority to 1, Action to Allow, Protocol to TCP and Port to 3306, Type to IPv4, and Source to 172.16.0.0/18.

Click **OK**. The security group rule is added.

- 7. Log in to the DLI management console. In the navigation pane on the left, choose **Datasource Connections**. On the displayed page, click **Create** in the **Enhanced** tab.
- 8. In the displayed dialog box, set the following parameters:
  - Connection Name: Name of the enhanced datasource connection
  - Resource Pool: Select the elastic resource pool created in Step 1: Create an Elastic Resource Pool.

**NOTE** 

If you cannot decide the elastic resource pool in this step, you can skip this parameter, go to the **Enhanced** tab, and click **More** > **Bind Resource Pool** in the **Operation** column of the row that contains this datasource connection after it is created.

- **VPC**: Select the VPC of the RDS DB instance obtained in **5**.
- **Subnet**: Select the subnet of the RDS DB instance obtained in **5**.
- Set other parameters as you need.

Click **OK**. Click the name of the created datasource connection to view its status. You can perform subsequent steps only after the connection status changes to **Active**.

- Click Resources > Queue Management, select the target queue, for example, general\_test. In the Operation column, click More and select Test Address Connectivity.
- 10. In the displayed dialog box, enter *Floating IP address:Database port* of the RDS database in the **Address** box and click **Test** to check whether the database is reachable.

#### Step 4: Run a Job

Run a Flink SQL jab on a queue in an elastic resource pool.

- 1. On the DLI management console, choose **Job Management** > **Flink Jobs**. On the **Flink Jobs** page, click **Create Job**.
- 2. In the **Create Job** dialog box, set **Type** to **Flink SQL** and **Name** to **testFlinkSqUob**. Click **OK**.
- 3. On the job editing page, set the following parameters:

#### Figure 3-36 Creating a Flink SQL job

testFlinkSqUo	x		
testFlinkS ID: 14173	nglobb (tour) Joan Type Print SGL		Start Save Save As
Check Seman	ics Debug Format Save as Template Theme Settings Help		
1 CREATE 2 type 3 regi 4 user	SIN: SINAH cer_info (id DH, new SINH) HZTH (	* Queue UDF Jar	
5 pass 6 db_u	and • , ' 1 • "9",	★ CUs	- 2 + ®
8 ); 9 INSERT	_mane = Tabletett2"	\star Job Manager CUs	- 1 + Pa
10 car_ 11 SELECT 12 13.	lefo	+ Parallelism	- 1 + (?)
13 'abo		Task Manager Configu	
		★ OBS Bucket	dion- 6
		Save Job Log	≥ at
		Alarm Generation upo.	
		Enable Checkpointing	
		Checkpoint Interval	- 30 + s
		Checkpoint Mode	Exactly once +
		Auto Hestart upon Exc	
		Dirty Data Policy	-Select-
-	<ul> <li>In Step 2: Add a Queue to the Elastic Resource Parallelistic Resource Resou</li></ul>	ogs an ollowir ed.	nd grant access ng is an
	<pre>type = "rds", region = "", /* Change the value to the current region ID. */ 'pwd_auth_name'="xxxxx", // Name of the datasource authentica created on DLI. If datasource authentication is used, you do not nee password for the job. db_url = "mysql://192.168.x.x:3306/test2", /* The format is mysql: number of the RDS database/database name. */ table_name = "tabletest2" /* Table name in RDS database */ ); INSERT INTO car_info SELECT 13, 'abc';</pre>	tion of t ed to set :// <i>floatii</i>	the password type t the username and <i>ng IP address:port</i>
Clic	k <b>Check Semantic</b> and ensure that the SOL stateme	nt na	ses the check

- 4. Click **Check Semantic** and ensure that the SQL statement passes the check. Click **Save**. Click **Start**, confirm the job parameters, and click **Start Now** to execute the job.
- 5. Wait until the job is complete. The job status changes to **Completed**.
- 6. Log in to the RDS console, click the name of the RDS DB instance. On the displayed page, click the name of the created database, for example, **test2**, and click **Query SQL Statements** in the **Operation** column of the row that containing the **tabletest2** table.
- 7. On the displayed page, click **Execute SQL**. Check whether data has been written into the RDS table.

#### Figure 3-37 Query result

12 × SQL Window ×		
Master Switch SQL Execution Node   Instance Name: no_delete   192.168.168.56.3386   Character Set. ut18		Save Executed SQL Statements
Execute SOL (F8)     (# Format SOL (F9)     (# Execute SOL Plan (F6)     SOL Favorites v		SQL Input Prompt (1) 🚺 Full Screen 💥
1 select * from 'tabletest2'		
Executed SQL Statements Messages Result Set1 ×		Overwrite Mode 💿
The following is the execution result set of select * from 'tabletest2'.	O This object has no primary key and cannot be estiled or exported in SOL format.	Copy Row Copy Column v Column Settings v
id	name	
1 13	abc	

# 3.6 Example Use Case: Configuring Scaling Policies for Queues in an Elastic Resource Pool

# Scenario

A company has multiple departments that perform data analysis in different periods during a day.

- Department A requires a large number of compute resources from 00:00 a.m. to 09:00 a.m. In other time segments, only small tasks are running.
- Department B requires a large number of compute resources from 10:00 a.m. to 10:00 p.m. Some periodical tasks are running in other time segments during a day.

In the preceding scenario, you can add two queues to an elastic resource pool: queue **test\_a** for department A, and queue **test\_b** for department B. You can add scaling policies for 00:00-09:00 and 10:00-23:00 respectively to the **test\_a** and **test\_b** queues. For jobs in other periods, you can modify the default scaling policy.

Que ue	Period	Priorit y	CUs	Defa ult Peri od	Def ault Prio rity	Default CUs	Remarks
test_ a	[00:00 , 09:00)	20	Minimu m CU: 64 Maximu m CU: 128	The time seg men ts beyo nd [00: 00, 09:0 0)	5	Minimum CU: 16 Maximum CU: 32	Jobs of department A

 Table 3-20 Scaling policy

Que ue	Period	Priorit y	CUs	Defa ult Peri od	Def ault Prio rity	Default CUs	Remarks
test_ b	[10:00 , 23:00)	20	Minimu m CU: 64 Maximu m CU: 128	The time seg men ts beyo nd [10: 00, 23:0 0)	5	Minimum CU: 32 Maximum CU: 64	Jobs of department B

# Precautions

• You are advised to implement fine-grained management of resource pools for stream and batch processing jobs by placing Flink real-time stream jobs and SQL batch processing jobs in separate elastic resource pools.

Flink real-time stream jobs can run stably without forced scale-in, thus avoiding job interruption and system instability.

SQL batch processing jobs are placed in independent resource pools, which can scale out and in more flexibly, significantly enhancing the success rate and operational efficiency of scaling operations.

- In any time segment of a day, the total minimum CUs of all queues in an elastic resource pool cannot be more than the minimum CUs of the pool.
- In any time segment of a day, the maximum CUs of any queue in an elastic resource pool cannot be more than the maximum CUs of the pool.
- The periods of scaling policies cannot overlap.
- The period of a scaling policy can only be set by hour and specified by the start time and end time. For example, if you set the period to **00-09**, the period when the policy takes effect is [00:00, 09:00). The period of the default scaling policy cannot be modified.
- In any period, compute resources are preferentially allocated to meet the minimum number of CUs of all queues. The remaining CUs (maximum CUs of the elastic resource pool total minimum CUs of all queues) are allocated in accordance with the scaling policy priorities.
  - Scaling policies with smaller priority values are prior to over those with larger priority values.
  - If the scaling policies of two queues have the same priority, resources are randomly allocated to a queue. If there are remaining resources, they are randomly allocated until there is no more left.

Scenario	CUs			
<ul> <li>An elastic resource pool has a maximum number of 256 CUs for queue A and queue B. The scaling policies are as follows:</li> <li>Queue A: priority 5; period: 00:00–09:00; minimum CU: 32; maximum CU: 128</li> <li>Queue B: priority 10; time period: 00:00–09:00; minimum CU: 64; maximum CU: 128</li> </ul>	<ul> <li>From 00:00 a.m. to 09:00 a.m.:</li> <li>1. The minimum CUs are allocated to the two queues. Queue A has 32 CUs, and queue B has 64 CUs. There are 160 CUs remaining.</li> <li>2. The remaining CUs are allocated based on the priorities. Queue B is prior to queue A. Therefore, queue B gets 64 CUs, and queue A has 96 CUs.</li> </ul>			
<ul> <li>An elastic resource pool has a maximum number of 96 CUs for queue A and queue B. The scaling policies are as follows:</li> <li>Queue A: priority 5; period: 00:00–09:00; minimum CU: 32; maximum CU: 64</li> <li>Queue B: priority 10; time period: 00:00–09:00; minimum CU: 64; maximum CU: 128</li> </ul>	<ul> <li>From 00:00 a.m. to 09:00 a.m.:</li> <li>1. The minimum CUs are allocated to the two queues. Queue A has 32 CUs, and queue B has 64 CUs. There are no remaining CUs.</li> <li>2. The allocation is complete.</li> </ul>			
<ul> <li>An elastic resource pool has a maximum number of 128 CUs for queue A and queue B. The scaling policies are as follows:</li> <li>Queue A: priority 5; period: 00:00-09:00; minimum CU: 32; maximum CU: 64</li> <li>Queue B: priority 10; time period: 00:00-09:00; minimum CU: 64; maximum CU: 128</li> </ul>	<ul> <li>From 00:00 a.m. to 09:00 a.m.:</li> <li>1. The minimum CUs are allocated to the two queues. Queue A has 32 CUs, and queue B has 64 CUs. There are 32 CUs remaining.</li> <li>2. The remaining 32 CUs are all preferentially allocated to queue B.</li> </ul>			
<ul> <li>An elastic resource pool has a maximum number of 128 CUs for queue A and queue B. The scaling policies are as follows:</li> <li>Queue A: priority 5; period: 00:00-09:00; minimum CU: 32; maximum CU: 64</li> <li>Queue B: priority 5; time period: 00:00-09:00; minimum CU: 64; maximum CU: 128</li> </ul>	<ul> <li>From 00:00 a.m. to 09:00 a.m.:</li> <li>1. The minimum CUs are allocated to the two queues. Queue A has 32 CUs, and queue B has 64 CUs. There are 32 CUs remaining.</li> <li>2. The two queues have the same priority, the remaining 32 CUs are randomly allocated to the two queues.</li> </ul>			

# Setting a Scaling Policy

- Step 1 Log in to the DLI management console and create an elastic resource pool. Set the minimum and maximum number of CUs of the pool to 128 and 256 respectively. For details, see Creating an Elastic Resource Pool and Creating Queues Within It.
- **Step 2** Choose **Resources** > **Resource Pool**. Locate the row that contains the created elastic resource pool, and click **Queue MGMT** in the **Operation** column.
- **Step 3** Refer to **Creating an Elastic Resource Pool and Creating Queues Within It** to create the **test\_a** queue and set the scaling policy.
  - 1. Set the priority of the default scaling policy to 5, **Min CU** to **16**, and **Max CU** to **32**.
  - 2. Click create to add a scaling policy. Set the priority to **20**, **Period** to **00--09**, **Min CU** to **64**, and **Max CU** to **128**.

Figure 3-38 Adding a scaling policy for queue test\_a

< Add Queue				
(1) Basic Configuration —	Elastic Resources			
View scaling policies of The priority ranges fro to hours in (start time,	if all queues in	he default value is 1. The new priority overrides the default prior	rity, time segments, max and min	CUs. You can only set the period
Priority	Period	Min CU	Max CU	
5	00 = 24 =	- 16 +	- 32 +	
20	00 • 09 •	- 64 +	- 128 +	Delete
Create				

**Step 4** View the scaling policy on the **Queue Management** page of the specific elastic resource pool.

Figure 3-39 Viewing the new scaling policy

										Enter a name.		2 🛛	<u>  </u> [	3
Name	Туре	Period	Min CUs	Max CUs	Priority 🕐	Engine	Created	Enterprise Project	Owner		Operation			
lest_a	For SQL	(00:00,09:00) (00:00,24:00)	64 16	128 32	20 5	spark	May 18, 2022 17:00:13 GMT+0	default			Edit Delete			

Click to view graphical statistics of priorities and CU settings for all time segments.





- **Step 5** Refer to **Creating an Elastic Resource Pool and Creating Queues Within It** to create the **test\_b** queue and set the scaling policy.
  - 1. Set the priority of the default scaling policy to **5**, **Min CU** to **32**, and **Max CU** to **64**.
  - 2. Click create to add a scaling policy. Set the priority to **20**, **Period** to **10--23**, **Min CU** to **64**, and **Max CU** to **128**.

Figure 3-41 Adding a scaling policy for queue test\_b

<   Add Queue(;				
1 Basic Configuration —	2 Elastic Resources			
View scaling policies of a	Il queues in pool_test.			
The priority ranges from 1 to hours in [start time,end	1 to 100. If you do not set the priority for a specific perio I time) format.	d, the default value is 1. The new priority overrides the default prior	ority, time segments, max and min CUs. You ca	n only set the period
Priority	Period	Min CU	Max CU	
5	00 v - 24 v	- 32 +	- 64 +	
20	10 v - 23 v	- 64 +	- 128 + Del	ete
		Create		

**Step 6** View the scaling policy on the **Queue Management** page of the specific elastic resource pool.

#### Figure 3-42 Viewing the new scaling policy

< Queue Managem	ent(pool									Add Ouese
									Enter a name	Q M C
Name	Type	Period	Min CUs	Max CUs	Priority (2)	Engine	Created	Enterprise Project	Owner	Operation
test_a	For SQL	[00:00,09:00) [00:00,24:00)	64 16	128 32	20 5	spark	May 18, 2022 17:47:40 GMT+0	default		Edit   Delete
test_b	For SQL	[00:00,24:00) [10:00,23:00)	32 64	64 128	5 20	spark	May 18, 2022 17:46:59 GMT+0	default	0	Edit   Delete

Click to view graphical statistics on priorities and CU settings of the two queues for all time segments.



Figure 3-43 Scaling policies of both queues

----End

# 3.7 Creating a Non-Elastic Resource Pool Queue (Deprecated, Not Recommended)

Queues in the non-elastic resource pool mode are the previous-gen of resource management for DLI. It involved purchasing and releasing resources based on usage demands, requiring estimation of resource needs before making purchases.

Queues in an elastic resource pool are recommended, as they offer the flexibility to use resources with high utilization as needed. For how to buy an elastic resource pool and create queues within it, see **Creating an Elastic Resource Pool and Creating Queues Within It**.

#### D NOTE

- If you use a sub-account to create a queue for the first time, log in to the DLI management console using the main account and keep records in the DLI database before creating a queue.
- It takes 6 to 10 minutes for a job running on a new queue for the first time.
- After a queue is created, if no job is run within one hour, the system releases the queue.
- Queues with 16 CUs do not support scale-out or scale-in.
- Queues with 64 CUs do not support scale-in.

# **Notes and Constraints**

ltem	Description
Resource type	Queue types:
	<ul> <li>default queue: A queue named default is preset in DLI, where you can use resources on demand. You are billed based on the amount of data scanned in each job (unit: GB).</li> </ul>
	<ul> <li>For SQL: Spark SQL jobs can be submitted to SQL queues.</li> </ul>
	<ul> <li>For general purpose: The queue is used to run Spark programs, Flink SQL jobs, and Flink Jar jobs.</li> </ul>
	<ul> <li>You cannot change the queue type once a queue is purchased. To use another queue type, purchase a new queue.</li> </ul>
Managing	• The billing mode of a queue cannot be changed.
queues	• The region of a queue cannot be changed.
	• When creating a queue in the non-elastic resource pool mode, you can only select cross-AZ active-active for yearly/ monthly queues and pay-per-use dedicated queues. The price of a cross-AZ queue is double that of a single-AZ queue.
	DLI queues cannot access the Internet.
Queue scaling	• Queues with 16 CUs do not support scale-out or scale-in.
	Queues with 64 CUs do not support scale-in.
	• Newly created queues need to run jobs before they can be scaled in or out.

Table 3-22 Notes and	constraints on queues
----------------------	-----------------------

# Procedure

- 1. You can create a queue on the **Overview**, **SQL Editor**, or **Queue Management** page.
  - In the upper right corner of the **Overview** page, click Purchase Queue.
  - To create a queue on the **Queue Management** page:
    - i. In the navigation pane on the left of the DLI management console, choose **Resources** > **Queue Management**.
    - ii. In the upper right corner of the **Queue Management** page, click **Buy Queue** to create a queue.
  - To create a queue on the **SQL Editor** page:
    - i. In the navigation pane of the DLI management console, click **SQL Editor**.

- ii. Click **Queues**. On the tab page displayed, click  $\bigcirc$  on the right to create a queue.
- 2. On the **Buy Queue** page displayed, set the parameters according to **Table 3-23**.

Table	3-23	Parameters	5
-------	------	------------	---

Paramet er	Description
Billing Mode	<b>Pay-per-use</b> . Billing for CUH used = Number of CUs x Usage duration x Unit price. You are billed for used CUs on an hourly basis (rounded up to the nearest hour). The pay-per-use billing mode is adopted. You are advised to purchase the CUH package to enjoy preferential price.
Region	Select a region. Select a region near you to ensure the lowest latency possible.
Project	Each region corresponds to a project.
Name	<ul> <li>Queue name</li> <li>Only numbers, letters, and underscores (_) are allowed. The value cannot contain only numbers, start with an underscore (_), or be left unspecified.</li> <li>The value can contain a maximum of 128 characters.</li> <li>NOTE The queue name is case-insensitive. Uppercase letters will be automatically converted to lowercase letters.</li> </ul>
Туре	<ul> <li>For SQL: compute resources used for SQL jobs.</li> <li>For general purpose: compute resources used for Spark and Flink jobs.</li> <li>NOTE         When a dedicated queue is not in use, its resources are still reserved and not released. This means that resources are constantly being held, regardless of whether the queue is being utilized. By using a dedicated queue, resources can be guaranteed to be available whenever jobs are submitted. You can create enhanced datasource connections for a dedicated queue.     </li> <li>When buying a pay-per-use queue, you have the option to choose a dedicated queue. Dedicated queues are billed 24/7 regardless of whether they are used.</li> </ul>

Paramet er	Description
AZ Mode	Available only when <b>Dedicated Resource Mode</b> is selected for <b>Type</b> .
	The deployment mode for the DLI queue. Select dual-AZ if you require high availability.
	<ul> <li>Currently, only SQL queues in dedicated resource mode support the dual-AZ policy.</li> </ul>
	<ul> <li>Dual-AZ improves data availability by creating a duplicate queue in the second AZ, but at an increased cost (twice as much as that of single AZ mode).</li> </ul>
	• This is a one-time configuration and cannot be changed later.
Specifica tions	The compute nodes' total number of CUs. One CU equals one vCPU and 4 GB of memory. DLI automatically assigns CPU and memory resources to each compute node, and the client does not need to know how many compute nodes are being used.
Enterpris e Project	If the created queue belongs to an enterprise project, you can select the corresponding enterprise project.
	Enterprise projects let you manage cloud resources and users by project.
	For how to set an enterprise project, see <b>Enterprise</b> Management User Guide.
	<b>NOTE</b> This parameter is displayed only for users who have enabled the Enterprise Management Service.
Descripti on	Description of the queue to be created. The value can contain a maximum of 128 characters.
Advance d	In the <b>Queue Type</b> area, select <b>Dedicated Resource Mode</b> and then click <b>Advanced Settings</b> .
Settings	• <b>Default</b> : The system automatically configures the parameter.
	• Custom CIDR Block: Enter a CIDR block range. If DLI enhanced datasource connection is used, the CIDR block of the DLI queue cannot overlap with that of the data source.
	Recommended CIDR blocks:
	10.0.0–10.255.0.0/8–24
	172.16.0.0-172.31.0.0/12-24
	192.168.0.0–192.168.0.0/16–24
	enhanced. When running other jobs, select <b>Basic</b> .

Paramet er	Description			
Tags	Tags used to identify cloud resources. A tag includes the tag key and tag value. If you want to use the same tag to identify multiple cloud resources, that is, to select the same tag from the drop-down list box for all services, you are advised to create predefined tags on the Tag Management Service (TMS).			
	If your organization has configured tag policies for DLI, add tags to resources based on the policies. If a tag does not comply with the tag policies, resource creation may fail. Contact your organization administrator to learn more about tag policies.			
	For details, see Tag Management Service User Guide.			
	NOTE			
	• A maximum of 20 tags can be added.			
	<ul> <li>Only one tag value can be added to a tag key.</li> </ul>			
	<ul> <li>The key name in each resource must be unique.</li> </ul>			
	• Tag key: Enter a tag key name in the text box.			
	<b>NOTE</b> A tag key can contain a maximum of 128 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed, but the value cannot start or end with a space or start with <b>sys</b> .			
	• Tag value: Enter a tag value in the text box.			
	<b>NOTE</b> A tag value can contain a maximum of 255 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed.			

3. Click **Buy** to confirm the configuration.

When creating a queue for the first time, select **Agree to the above privacy agreements** and click **OK**.

4. Confirm the configuration and click **Submit**.

If the queue name already exists, the system displays a message indicating that the queue name already exists when you click **Submit**. In this case, click **Previous** to go back to modify the queue name.

5. Once created, you can check and use it on the **Queue Management** page.

#### **NOTE**

It takes 6 to 10 minutes for a job running on a new queue for the first time.

# **4** Creating a Data Directory, Database, and Table

# 4.1 Understanding Data Catalogs, Databases, and Tables

Databases and tables are the basis for developing SQL and Spark jobs. Before running a job, you need to define databases and tables based on your service scenarios.

#### **NOTE**

Flink allows for dynamic data types, enabling the definition of data structures at runtime without the need for predefined metadata.

# Data Catalog

A data catalog is a metadata management object that can contain multiple databases.

DLI currently supports DLI data catalogs and LakeFormation data catalogs.

For how to create a database and a table in a DLI data catalog, see **Creating a Data Catalog**, **Database**, **and Table on the DLI Console**.

For how to create and use a LakeFormation metadata file, see **Creating and Using LakeFormation Metadata**.

#### Database

A database is a repository of data organized, stored, and managed on computer storage devices according to data structures. Databases are typically used to store, retrieve, and manage structured data, consisting of multiple data tables that are interrelated through keys and indexes.

#### Table

Tables are one of the most important components of a database, consisting of rows and columns. Each row represents a data item, while each column represents

a property or feature of the data. Tables are used to organize and store specific types of data, making it possible to query and analyze the data effectively.

A database is a framework, and tables are its essential content. A database contains one or more tables.

You can create databases and tables on the management console or using SQL statements. For details about using SQL statements, see **Creating a Database**, **Creating an OBS Table**, and **Creating a DLI Table**. This section describes how to create a database and a table on the management console.

#### **NOTE**

When creating a database or table, you need to grant permissions to other users so that they can view the new database or table.

## **Table Metadata**

Metadata is data used to define the type of data. It primarily describes information about the data itself, including its source, size, format, or other characteristics. In database fields, metadata is used to interpret the content of a data warehouse.

When creating a table, metadata is defined by three columns: column name, type, and column description.

#### Table Types Supported by DLI

#### • DLI table

DLI tables are data tables stored in a DLI data lake. They can store structured, semi-structured, and unstructured data.

Data in DLI tables is stored internally within the DLI service, resulting in better query performance. This makes it suitable for time-sensitive services, such as interactive queries.

In the navigation pane on the left of the DLI console, choose **Data Management** > **Databases and Tables**. On the displayed page, click the name of a database. In the displayed table list, tables whose **Type** is **MANAGED** are DLI tables.

#### OBS table

Data in OBS tables is stored in the OBS service, which is suitable for latencyinsensitive services, such as historical data statistics and analysis.

An OBS table stores data in the form of objects. Each object contains data and related metadata.

In the navigation pane on the left of the DLI console, choose **Data Management** > **Databases and Tables**. On the displayed page, click the name of a database. In the displayed table list, tables whose **Type** is **EXTERNAL** and **Storage Location** is **OBS** are OBS tables.

#### • View table

A view table is a virtual table that does not store actual data. Instead, it dynamically generates data based on the defined query logic. Views are typically used to simplify complex queries or provide customized data views for different users or applications.

A view table can be created based on one or multiple tables, providing a flexible way to display data without affecting the storage and organization of the underlying data.

In the navigation pane on the left of the DLI console, choose **Data Management** > **Databases and Tables**. On the displayed page, click the name of a database. In the displayed table list, tables whose **Type** is **VIEW** are view tables.

**NOTE** 

A view can only be created using SQL statements. You cannot create a view on the **Create Table** page. Table or view information in a view cannot be modified. Otherwise, the query may fail.

#### • Datasource table

A datasource table is a data table that can be queried and analyzed across multiple data sources. This type of tables can integrate data from varying data sources and provide a unified data view.

Datasource tables are typically used in data warehouse and data lake architectures, allowing users to perform complex queries across multiple data sources.

In the navigation pane on the left of the DLI console, choose **Data Management** > **Databases and Tables**. On the displayed page, click the name of a database. In the displayed table list, tables whose **Type** is **EXTERNAL** and **Storage Location** is not **OBS** are datasource tables.

## Notes and Constraints on Databases and Tables

ltem	Description
Database	<ul> <li>default is the database built in DLI. You cannot create a database named default.</li> </ul>
	• DLI supports a maximum of 50 databases.
Data table	• DLI supports a maximum of 5,000 tables.
	<ul> <li>DLI supports the following table types:</li> </ul>
	- MANAGED: Data is stored in a DLI table.
	- <b>EXTERNAL</b> : Data is stored in an OBS table.
	<ul> <li>View: A view can only be created using SQL statements.</li> </ul>
	<ul> <li>Datasource table: The table type is also EXTERNAL.</li> </ul>
	<ul> <li>You cannot specify a storage path when creating a DLI table.</li> </ul>

 Table 4-1 Notes and constraints on DLI resources

ltem	Description
Data import	<ul> <li>Only OBS data can be imported to DLI or OBS.</li> <li>You can import data in CSV, Parquet, ORC, JSON, or Avro format from OBS to tables created on DLI.</li> <li>To import data in CSV format to a partitioned table, place the partition column in the last column of the data source.</li> </ul>
	<ul> <li>The encoding format of imported data can only be UTF-8.</li> </ul>
Data export	• Data in DLI tables (whose table type is <b>MANAGED</b> ) can only be exported to OBS buckets, and the export path must contain a folder.
	• The exported file is in JSON format, and the text format can only be UTF-8.
	• Data can be exported across accounts. That is, after account B authorizes account A, account A has the permission to read the metadata and permission information of account B's OBS bucket as well as the read and write permissions on the path. Account A can export data to the OBS path of account B.

# Table Management Page

From the **Data Management** page, click the database name or **Tables** in the **Operation** column to switch to the table management page.

The displayed page lists all tables created in the current database. You can view the table type, data storage location, and other information. Tables are listed in chronological order by default, with the most recently created tables displayed at the top.

# 4.2 Creating a Data Catalog, Database, and Table on the DLI Console

• A data catalog is a metadata management object that can contain multiple databases.

DLI currently supports DLI data catalogs and LakeFormation data catalogs.

- DLI data catalog: The data catalog service provided by DLI, which is used to store and manage metadata in a data lake. The default name of a DLI data catalog is dli.
- LakeFormation data catalog: LakeFormation provides unified metadata management. You need to create a connection to a LakeFormation catalog on the DLI management console to access the catalog stored within a LakeFormation instance. After connecting DLI to the default LakeFormation instance and authorizing access to LakeFormation

resources, you can use LakeFormation metadata during DLI job development.

- A database is a repository of data organized, stored, and managed on computer storage devices according to data structures.
- A table is one of the most essential components of a database. It is composed of rows and columns, with each column regarded as a field. The values within each field represent a specific type of data.

The database is a framework and the table contains data content. A database has one or more tables.

You can create databases and tables on the management console or using SQL statements.

For details about how to create a database and table using SQL statements, see **Creating a Database**, **Creating an OBS Table**, and **Creating a DLI Table**.

This section describes how to create a data catalog, database, and table on the management console.

#### **NOTE**

- Views can be created only by using SQL statements, not through the **Create Table** page.
- For Hudi tables created using SQL statements, you need to configure Hive synchronization parameters before they can be checked in the databases and tables on the DLI management console.

Why Is a Hudi Table Not Displayed on the DLI Console?

## Precautions

When creating a data catalog, database, or table, you need to grant permissions to other users so that they can view the new data catalog, database, or table. For details, see **Common Operations Supported by DLI System Policy**.

#### Creating a Data Catalog

The DLI management console provides the DLI data catalog by default. You can also follow this section's instructions to create a connection to a LakeFormation catalog on the DLI management console. Once created, the LakeFormation catalog will be displayed under the data catalog list on the DLI management console.

- 1. Before creating a connection to a LakeFormation catalog on DLI, ensure that a data catalog has been created on the LakeFormation management console.
  - a. Log in to the LakeFormation management console.
  - b. In the navigation pane on the left, choose **Metadata** > **Catalog**.
  - c. On the displayed page, click **Create**.

Set catalog instance parameters as needed.

For parameter settings and descriptions, see Creating a Catalog.

d. Once created, you can view information about the created catalog on the **Catalog** page.

DLI can only connect to the default LakeFormation instance. Set the instance in LakeFormation as the default to ensure successful connection.

#### 2. Create a data catalog on the DLI management console.

On the DLI management console, you need to create a connection to a LakeFormation catalog to enable access to the catalog stored within a LakeFormation instance when submitting jobs on DLI.

You can create data catalog connections on three pages of the DLI management console, and the created data catalog connections will be visible under the **Catalog** tab of the **SQL Editor** page.

- On the Catalog tab of the SQL Editor page, click 
   to create a connection to a LakeFormation catalog.
- On the Flink job editing page, click 
   next to Catalog Name to create a connection to a LakeFormation catalog. (Only Flink 1.17 or later supports configuring data catalogs.)
- On the Spark job editing page, click 

   next to Catalog Name to create a connection to a LakeFormation catalog. (Only Spark 3.3.1 supports configuring data catalogs.)

#### **NOTE**

You can only create one mapping for each data catalog in LakeFormation.

For example, a user creates a mapping named **catalogMapping1** in DLI, which corresponds to the data catalog **catalogA** in LakeFormation. Once created, you cannot create a mapping to **catalogA** in the same project space.

Take creating a data catalog connection on the **Catalog** tab of the **SQL Editor** page as an example:

- a. Log in to the DLI management console.
- b. In the navigation pane on the left, choose SQL Editor.
- c. On the SQL editor page, select a data catalog under **Catalog**.
- d. Click  $\textcircled{\textcircled{}}$  to create a data catalog.
- e. In the Create Catalog dialog box, set data catalog parameters.

#### Table 4-2 Data catalog parameters

Parameter	Mand atory	Description
External Catalog Name	Yes	Catalog name of the default LakeFormation instance.
Туре	Yes	Currently, the only available option is <b>LakeFormation</b> .
		This option is fixed and does not need to be selected.

Parameter	Mand atory	Description
Catalog Name	Yes	Catalog mapping name used in DLI. When running SQL statements, you need to specify the catalog mapping to identify the external metadata to be accessed. You are advised to set this parameter to the same value as <b>External Catalog Name</b> .
		Currently, DLI can only connect to the data catalog of the default LakeFormation instance.
Description	No	Description of the data catalog.

- f. Click **OK**.
- g. After the data catalog is created, the connection status of the data catalog is displayed in the data catalog list.
  - Blinking indicates that the data catalog is being created.
  - indicates that the data catalog has been created and the data catalog connection has been activated.
  - indicates that the data catalog fails to be created. In this case, it is advised to delete the current data connection and create a data catalog again.

#### Creating a Database

- **Step 1** You can create a database on either the **Data Management** page or the **SQL Editor** page.
  - To create a database on the **Data Management** page:
    - a. On the left of the management console, choose **Data Management** > **Databases and Tables**.
    - b. In the upper right corner of the **Databases and Tables** page, click **Create Database** to create a database.
  - To create a database on the **SQL Editor** page:
    - a. On the left of the management console, click **SQL Editor**.
    - b. In the navigation pane on the left, click  $\textcircled{\textcircled{}}$  next to **Databases**.
- **Step 2** In the displayed **Create Database** dialog box, specify **Name** and **Description** by referring to **Table 4-3**.
# Figure 4-1 Creating a database

Create Databas	se	×
You can create 1 more	databases. Increase quota.	
★ Databases	testdb	
Description		
★ Enterprise Project	0/* Enterprise Project  C ⑦ Create Enterprise Project  It is recommended that you use TMS's predefined tag function to add the same tag to differen cloud resources. View predefined tags ⑦ To add a tag, enter a tag key and a tag value below.	128 t
Tags		
	Enter a tag key     Enter a tag value     Add       10 tags available for addition.     OK     Cancel	

# Table 4-3 Description

Paramete r	Description
Database Name	<ul> <li>Only numbers, letters, and underscores (_) are allowed. The value cannot contain only numbers or start with an underscore (_).</li> </ul>
	• The database name is case insensitive and cannot be left blank.
	• The value can contain a maximum of 128 characters.
	<b>NOTE</b> The <b>default</b> database is a built-in database. You cannot create the <b>default</b> . database.
Enterprise Project	If the created queue belongs to an enterprise project, you can select the corresponding enterprise project.
	Enterprise projects let you manage cloud resources and users by project.
	For how to set enterprise projects, see <b>Enterprise Management</b> <b>User Guide</b> .
	<b>NOTE</b> This parameter is displayed only for users who have enabled the Enterprise Management Service.
Descriptio n	Description of a database.

Paramete r	Description
Tags	Tags used to identify cloud resources. A tag includes the tag key and tag value. If you want to use the same tag to identify multiple cloud resources, that is, to select the same tag from the drop-down list box for all services, you are advised to create predefined tags on the Tag Management Service (TMS).
	If your organization has configured tag policies for DLI, add tags to resources based on the policies. If a tag does not comply with the tag policies, resource creation may fail. Contact your organization administrator to learn more about tag policies.
	For details, see Tag Management Service User Guide.
	NOTE
	• A maximum of 20 tags can be added.
	<ul> <li>Only one tag value can be added to a tag key.</li> </ul>
	The key name in each resource must be unique.
	• Tag key: Enter a tag key name in the text box.
	<b>NOTE</b> A tag key can contain a maximum of 128 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed, but the value cannot start or end with a space or start with <b>_sys</b>
	• Tag value: Enter a tag value in the text box.
	<b>NOTE</b> A tag value can contain a maximum of 255 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed.

#### Step 3 Click OK.

After a database is created, you can view and select the database for use on the **Databases and Tables** page or **SQL Editor** page.

----End

# Creating a Table

Before creating a table, ensure that a database has been created.

**Step 1** You can create a table on either the **Databases and Tables** page or the **SQL Editor** page.

Datasource connection tables, such as View tables, HBase (CloudTable/MRS) tables, OpenTSDB (CloudTable/MRS) tables, GaussDB(DWS) tables, RDS tables, and CSS tables, cannot be created. You can use SQL to create views and datasource connection tables. For details, see **Data Lake Insight SQL Syntax Reference**.

- To create a table on the **Data Management** page:
  - a. On the left of the management console, choose **Data Management** > **Databases and Tables**.

- On the Databases and Tables page, select the database for which you want to create a table. In the Operation column, click More > Create Table to create a table in the current database.
- To create a table on the **SQL Editor** page:
  - a. On the left of the management console, click **SQL Editor**.
  - b. In the navigation pane of the displayed **SQL Editor** page, click **Databases**. You can create a table in either of the following ways:
    - Click a database name. In the Tables area, click on the right to create a table in the current database.
- **Step 2** In the displayed **Create Table** dialog box, set parameters as required.
  - If you set Data Location to DLI, set related parameters by referring to Table 4-4.

#### Figure 4-2 Creating a DLI table

Create Table			×
You can create 53 n	nore tables. Increase quota.		
* Table Name	testtable		
* Data Location	DLI		
Description			
		0/256	
			Add Column Excel Import
Column Type	Column	Туре	Description
Normal	▼ col1	string •	
		OK Cancel	

• If you set **Data Location** to **OBS**, set related parameters by referring to **Table 4-4** and **Table 4-5**.

When there are both a folder and a file with the same name in the OBS directory, creating an OBS table pointing to that path will prioritize the file over the folder.

# Figure 4-3 Creating an OBS table

Create Table					×
You can create 53 n	nore tables. Increase quota.				
★ Table Name	testtable				
* Data Location	OBS •				
★ Data Format	CSV 💌				
* Storage Path	This path must start with "obs://".				
Description					
			0/256		
				Add Column	Excel Import
Column Type	Column	Туре		Description	
Normal	▼ col1	string	•		
Advanced Settings					
		OK Cancel			

# Table 4-4 Common parameters

Paramet er	Description	Exampl e
Table Name	<ul> <li>Only numbers, letters, and underscores (_) are allowed. The value cannot contain only numbers or start with an underscore (_).</li> </ul>	table01
	<ul> <li>The table name is case insensitive and cannot be left unspecified.</li> </ul>	
	<ul> <li>The table name can contain the dollar sign (\$). An example value is <b>\$test</b>.</li> </ul>	
	<ul> <li>The value can contain a maximum of 128 characters.</li> </ul>	
Data Location	Data storage location. Currently, DLI and OBS are supported.	DLI
Descripti on	Description of the table.	-
Column Type	Available values: Normal or Partition	Normal

Paramet er	Description	Exampl e
Column	Name of a column in a table. The column name must contain at least one letter and can contain underscores (_). It cannot contain only numbers. You can select <b>Normal</b> or <b>Partition</b> . Partition columns are dedicated to partition tables. User data is partitioned to improve query efficiency. <b>NOTE</b> The column name is case-insensitive and must be unique.	name
Туре	<ul> <li>Data type of a column. This parameter corresponds to Column Name.</li> <li>string: The data is of the string type.</li> <li>int: Each integer is stored on four bytes.</li> <li>date: The value ranges from 0000-01-01 to 9999-12-31.</li> <li>double: Each number is stored on eight bytes.</li> <li>boolean: Each value is stored on one byte.</li> <li>decimal: The valid bits are positive integers between 1 to 38, including 1 and 38. The decimal digits are integers less than 10.</li> <li>smallint/short: The number is stored on eight bytes.</li> <li>bigint/long: The number is stored on eight bytes.</li> <li>timestamp: The data indicates a date and time. The value can be accurate to six decimal points.</li> <li>float: Each number is stored on one byte. Only OBS tables support this data type.</li> </ul>	string
Column Descripti on	Description of a column.	-
Operatio n	<ul> <li>Add Column</li> <li>Delete</li> <li>NOTE         <ul> <li>If the table to be created includes a great number of columns, you are advised to use SQL statements to create the table or import column information from the local EXCEL file.</li> </ul> </li> </ul>	-

Paramete r	Description	Example
Data Format	DLI supports the following data formats:	CSV
	<ul> <li>Parquet: DLI can read non- compressed data or data that is compressed using Snappy and gzip.</li> </ul>	
	<ul> <li>CSV: DLI can read non-compressed data or data that is compressed using gzip.</li> </ul>	
	<ul> <li>ORC: DLI can read non-compressed data or data that is compressed using Snappy.</li> </ul>	
	<ul> <li>JSON: DLI can read non-compressed data or data that is compressed using gzip.</li> </ul>	
	<ul> <li>Avro: DLI can read uncompressed Avro data.</li> </ul>	
Storage Path	Enter or select an OBS path. The path can be a file or folder.	obs://obs1/ sampledata.csv
	<ul> <li>When creating an OBS table, you must specify a folder as the path. If a file is specified, data cannot be imported.</li> </ul>	
	<ul> <li>When there are both a folder and a file with the same name in the OBS directory, importing data pointing to that path will prioritize the file over the folder.</li> </ul>	
Table Header: No/Yes	This parameter is valid only when <b>Data</b> <b>Format</b> is set to <b>CSV</b> . Whether the data source to be imported contains the table header.	-
	Click <b>Advanced Settings</b> and select the checkbox next to <b>Table Header:</b> <b>No</b> . If the checkbox is selected, the table header is displayed. If the checkbox is deselected, no table header is displayed.	

Table 4-5 Parameter description when Data Location is set to OBS

Paramete r	Description	Example
User- defined Delimiter	This parameter is valid only when <b>Data</b> Format is set to CSV and you select User-defined Delimiter.	Comma (,)
	The following delimiters are supported:	
	– Comma (,)	
	- Vertical bar ( )	
	– Tab character (\t)	
	<ul> <li>Others: Enter a user-defined delimiter.</li> </ul>	
User- defined Quotation	This parameter is valid only when <b>Data</b> Format is set to CSV and you select User-defined Quotation Character.	Single quotation mark (')
Character	The following quotation characters are supported:	
	<ul> <li>Single quotation mark (')</li> </ul>	
	<ul> <li>Double quotation marks (")</li> </ul>	
	<ul> <li>Others: Enter a user-defined quotation character.</li> </ul>	
User- defined Escape Character	This parameter is valid only when <b>Data</b> Format is set to CSV and you select User-defined Escape Character. The following escape characters are	Backslash (\)
	supported:	
	- DdCKSIdSII (\) Others: Enter a user defined escape	
	character.	
Date Format	This parameter is valid only when <b>Data</b> Format is set to CSV or JSON.	2000-01-01
	This parameter specifies the format of the date in the table and is valid only <b>Advanced Settings</b> is selected. The default value is <b>yyyy-MM-dd</b> . For definition of characters involved in the date pattern, see Table 3 in <b>Importing</b> <b>Data to the Table</b> .	

Paramete r	Description	Example
Timestam p Format	This parameter is valid only when <b>Data</b> Format is set to CSV or JSON.	2000-01-01 09:00:00
	This parameter specifies the format of the timestamp in the table and is valid only <b>Advanced Settings</b> is selected. The default value is <b>yyyy-MM-dd</b> <b>HH:mm:ss</b> . For definition of characters involved in the time pattern, see Table 3 in <b>Importing Data to the Table</b> .	

#### Step 3 Click OK.

After a table is created, you can view and select the table for use on the **Data Management** page or **SQL Editor** page.

----End

# **Related Operations**

After a table is created, you can import data from other OBS buckets to the table.

For details about how to import data, see Importing OBS Data to a DLI Table.

# 4.3 Viewing Table Metadata

# **Metadata Description**

- Metadata is used to define data types. It describes information about the data, including the source, size, format, and other data features. In database fields, metadata interprets data content in the data warehouse.
- When you create a table, metadata is defined, consisting of the column name, type, and description.
- The **Metadata** page displays information about the target table, including **Column Name**, **Column Type**, **Data Type**, and **Description**.

# Procedure

You can view metadata on either the **Data Management** page or the **SQL Editor** page.

- To view metadata on the **Data Management** page:
  - a. On the left of the management console, choose **Data Management** > **Databases and Tables**.
  - b. On the displayed **Data Management** page, click the name of the database where the target table whose data you want to export resides to switch to the **Manage Tables** page.
  - c. Click **More** in the **Operation** column of the target table and select **View Properties**. In the **Metadata** tab, view the metadata of the table.

- To view metadata on the **SQL Editor** page:
  - a. On the left of the management console, click **SQL Editor**.
  - b. In the navigation pane of the displayed **SQL Editor** page, click **Databases**.
  - c. Click the corresponding database name to view the tables in the database.
  - d. Click  $\equiv$  on the right of the table and choose **View Properties** from the shortcut menu. On the **Metadata** tab page, view the metadata of the table.

# 4.4 Managing Data Catalogs on the DLI Console

# 4.4.1 Configuring Data Catalog Permissions on the DLI Console

# Scenario

- DLI data catalogs support authorization on the DLI console or authentication through IAM. By setting permissions, you can grant varying data catalog permissions to different users.
- Administrators and data catalog owners (users who create data catalogs on DLI) have all data catalog permissions by default. You do not need to set permissions for them, and other users cannot modify their data catalog permissions.
- When setting data catalog permissions for a new user, ensure that the region of the user group the user belongs to has the **Tenant Guest** permission.

For details about the **Tenant Guest** permission and how to apply for the permission, see **System Permissions** and **Creating a User Group** in *Identity and Access Management User Guide*.

# Precautions

- Data catalog permissions are non-inherited permissions, meaning they apply only to the current data catalog. Databases and tables within the data catalog cannot inherit any permissions from it.
- Data catalog owners or users with the **Assign catalog access to specified users** permission can grant permissions to data catalogs.
- If you create a data catalog with the same name after deleting an existing one, the permissions will not be inherited and must be regranted to users.

# **Data Catalog Permissions**

You can also use IAM to grant data catalog permissions to specified users. For details about data catalog permissions, see **Table 4-6**.

Operation	Permission Set ( service:resourc e:action )	Authorization on the DLI Console	API-based Authoriza tion	IAM- based Authoriz ation
Unbinding a data catalog	dli:catalog:unbin d	Supported	Supported	Supporte d
Querying data catalog binding details	dli:catalog:get	Supported	Supported	Supporte d
Granting permissions	dli: catalog: grantPrivilege	Supported	Supported	Supporte d
Revoking permissions	dli: catalog: revokePrivilege	Supported	Supported	Supporte d
Viewing permissions of other users	dli: catalog: showPrivileges	Supported	Supported	Supporte d
Binding a data catalog	dli:catalog:bind	Not supported	Supported	Supporte d
Querying the data catalog binding list	dli:catalog:list	Not supported	Supported	Supporte d

#### Table 4-6 Data catalog permissions

# Granting Data Catalog Permissions to a New User on the DLI Management Console

Grant permissions to a new user or project that previously did not have permissions on this data catalog.

- 1. In the navigation pane on the left of the management console, choose **SQL Editor**.
- 2. On the displayed **Catalog** tab, locate the data catalog you want to view, click ≡ , and select **Permissions**.
- 3. On the displayed page, click **Grant Permission** in the upper right corner. In the dialog box that appears, enter the username you want to grant permissions to and select required permissions. For details about the permissions, see **Table 4-7**.

#### Figure 4-4 Granting permissions on a data catalog to a user

Grant Perr	nission		×	
★ Username	Enter a username.			
Select the permi	ssions to be granted to the user			
Select all				
Query the d	etails of bound catalogs	Unbind catalogs	Assign catalog access to specified users	
Revoke cata	alog access from specified users	View catalog access rights for other users		
			Cancel	

### Table 4-7 Parameter descriptions

Parameter	Description
Username	Name of an IAM user you wish to grant permissions to. <b>NOTE</b> The username must be an existing IAM username and has been used to log in to the DLI management console.
Permission	<ul> <li>Selecting a permission grants it to a user, while deselecting a permission revokes it from the user.</li> <li>Data catalog permissions are non-inherited permissions, meaning they apply only to the current data catalog.</li> <li>Databases and tables within the data catalog cannot inherit any permissions from it.</li> <li>Unbind catalogs: permission to unbind the data catalog from DLI.</li> <li>Query the details of bound catalogs: permission to view data catalog binding information. This permission is required if you need to use the data catalog when submitting jobs.</li> <li>Assign catalog access to specified users: permission to grant permissions on a data catalog to specified users.</li> <li>Revoke catalog access from specified users: permission to revoke permissions on a data catalog from specified users.</li> <li>View catalog access rights for other users: permission to view the permissions of other users in the current data catalog.</li> </ul>

4. Click **OK**.

# Modifying Permissions on a Data Catalog

If a user has certain permissions on a data catalog, you can modify or revoke those permissions for the user.

#### **NOTE**

If the options in **Set Permission** are gray, the corresponding account does not have the permission to modify the data catalog. You can request the **Assign catalog access to specified users** and **Revoke catalog access from specified users** permissions from administrators, data catalog owners, or other authorized users with permission-granting permissions.

- 1. On the **User Permissions** page, locate the user you wish to set permissions for.
  - If the user is an IAM user, you can set permissions for it.
  - If the user is already an administrator, you can only view the permissions information.
- 2. In the **Operation** column of the IAM user or project, click **Set Permission**. The **Set Permission** dialog box appears.

For details about data catalog permissions, see Table 4-7.

3. Select or deselect the permissions and click **OK**.

# 4.5 Managing Database Resources on the DLI Console

# 4.5.1 Configuring Database Permissions on the DLI Console

### Scenario

- By setting permissions, you can assign varying database permissions to different users.
- The administrator and database owner have all permissions, which cannot be set or modified by other users.
- When setting database permissions for a new user, ensure that the region of the user group to which the user belongs has the **Tenant Guest** permission. For details about the Tenant Guest permission and how to apply for the permission, see **Creating a User Group and Assigning Permissions** and **System Permissions**.

## Precautions

- By the rules in **Common Operations Supported by DLI System Policy**, you need to grant the current user the permission to view the databases of an administrator or another user.
- Lower-level objects automatically inherit permissions granted to upper-level objects. The hierarchical relationship is database > table > column.
- The database owner, table owner, and **authorized** users can assign permissions on the database and tables.
- Columns can only inherit the query permission. For details about Inheritable Permissions, see Configuring Database Permissions on the DLI Console.
- The permissions can be revoked only at the initial level to which the permissions are granted. You need to grant and revoke permissions at the same level. You need to grant and revoke permissions at the same level. For example, after you are granted the insertion permission on a database, you

can obtain the insertion permission on the tables in the database. Your insertion permission can be revoked only at the database level.

• If you create a database with the same name as a deleted database, the database permissions will not be inherited. In this case, you need to grant the database permissions to users or projects.

For example, user A is granted with the permission to delete the **testdb** database. Delete the database and create another one with the same name. You need to grant user A the deletion permission of the **testdb** database again.

### **Viewing Database Permissions**

- On the left of the management console, choose Data Management > Databases and Tables.
- 2. Locate the row where the target database resides and click **Manage Permissions** in the **Operation** column.

Permissions can be granted to new users or projects, modified for users or projects with existing permissions, or revoked from a user or project.

### Granting Permissions to a New User or Project

Here, the new user or project refers to a user or a project that does not have permissions on the database.

- 1. Click a database you need. In the displayed **Database Permission Management** page, click **Grant Permission** in the upper right corner.
- 2. In the displayed dialog box, select **User** or **Project**, enter the username or select the project to be authorized, and select the required permissions. For details about the permissions, see **Table 4-8**.

Figure 4-5 Granting database permissions to a user

Grant Permission						~
★ Authorization Object	User	Project				
★ Username	Enter a usernam	ie.				
Select the permissions to be Non-Inherited Permi	e granted to the u	ser				
Select all View Table View Table Creation St Insert Access Metadata Revoke Permission	atement	Select Table Drop Table Overwrite Table View Other Use	e er's Permissions	View Table Ir Rename Tabl Add Column Grant Permis	ifo e sion	
Column Permissions				E	nter a column name.	Q
Name			Permission Type			
id			Select			
name			Select			
score			Select			
		ОК	Cancel			

Grant Permission					×
* Authorization Object User Pro	oject				
* Project ID	•				
Select the permissions to be granted to the user					
Non-Inherited Permissions					
Select all					
Select Table	View Table Info		View Tabl	le Creation Statement	
Drop Table	Rename Table		Insert		
Overwrite Table	Add Column		Access M	etadata	
View Other User's Permissions	Grant Permission		Revoke Pe	ermission	
Column Permissions				Enter a column name.	Q
Name		Permission Type			
id		Select			
name		Select			
score		Select			
	ОК	ancel			

# Figure 4-6 Granting database permissions on a project

#### Table 4-8 Parameters

Parameter	Description
Authorizatio n Object	Select <b>User</b> or <b>Project</b> .
Username/ Project Name	<ul> <li>If you select User, enter the IAM username when adding a user to the database.</li> </ul>
	The username is an existing IAM user name and has logged in to the DLI management console.
	• If you select <b>Project</b> , select the project to be authorized in the current region.
	NOTE When <b>Project</b> is selected:
	<ul> <li>If you set Non-inheritable Permissions, you cannot view tables in the corresponding database within the project.</li> </ul>
	<ul> <li>If you set Inheritable Permissions, you can view all tables in the database within the project.</li> </ul>

Parameter	Description			
Non- Inherited	Select a permission to grant it to the user, or deselect a permission to revoke it.			
Permissions	Non-inherited permissions apply only to the current database.			
	<ul> <li>The following permissions are applicable to both user and project authorization:</li> </ul>			
	<ul> <li>Drop Database: This permission allows you to delete the current database.</li> </ul>			
	<ul> <li>Create Table: This permission allows you to create tables in the current database.</li> </ul>			
	<ul> <li>Create View: This permission allows you to create views in the current database.</li> </ul>			
	<ul> <li>Execute SQL EXPLAIN: This permission allows you to execute an EXPLAIN statement and view information about how this database executes a query.</li> </ul>			
	<ul> <li>Create Role: This permission allows you to create roles in the current database.</li> </ul>			
	<ul> <li>Delete Role: This permission allows you to delete roles of the current database.</li> </ul>			
	<ul> <li>View Role: This permission allows you to view the role of the current user.</li> </ul>			
	<ul> <li>Bind Role: This permission allows you to bind roles to the current database.</li> </ul>			
	<ul> <li>Unbind Role: This permission allows you to bind roles from the current database.</li> </ul>			
	<ul> <li>View All Binding Relationships: This permission allows you to view the binding relationships between all roles and users.</li> </ul>			
	<ul> <li>Create Function: This permission allows you to create a function in the current database.</li> </ul>			
	<ul> <li>Delete Function: This permission allows you to delete functions from the current database.</li> </ul>			
	<ul> <li>View All Functions: This permission allows you to view all functions in the current database.</li> </ul>			
	<ul> <li>View Function Details: This permission allows you to view details about the current function.</li> </ul>			
	• The following permissions can only be granted to users:			
	<ul> <li>View All Tables: This permission allows you to view all tables in the current database.</li> </ul>			
	<b>NOTE</b> If this permission of a specific database is not granted, all tables in the database will not be displayed.			
	<ul> <li>View Database: This permission allows you to view the information about the current database.</li> </ul>			

Parameter	Description
	<b>NOTE</b> If this permission is not granted, the database will not be displayed.

Parameter	Description
Inherited Permissions	Select a permission to grant it to the user, or deselect a permission to revoke it.
	Inherited permissions are applicable to the current database and all its tables. However, only the query permission is applicable to table columns.
	The following permissions can be granted to both user and project.
	• <b>Drop Table</b> : This permission allows you to delete tables in a database.
	• <b>Select Table</b> : This permission allows you to query data of the current table.
	• View Table Information: This permission allows you to view information about the current table.
	• <b>Insert</b> : This permission allows you to insert data into the current table.
	• Add Column: This permission allows you to add columns to the current table.
	• <b>Overwrite</b> : This permission allows you to insert data to overwrite the data in the current table.
	• <b>Grant Permission</b> : This permission allows you to grant database permissions to other users or projects.
	• <b>Revoke Permission</b> : This permission allows you to revoke the permissions of the database that other users have but cannot revoke the database owner's permissions.
	• Add Partition to Partition Table: This permission allows you to add a partition to a partition table.
	• <b>Delete Partition from Partition Table</b> : This permission allows you to delete existing partitions from a partition table.
	• <b>Configure Path for Partition</b> : This permission allows you to set the path of a partition in a partition table to a specified OBS path.
	• <b>Rename Table Partition</b> : This permission allows you to rename partitions in a partition table.
	• <b>Rename Table</b> : This permission allows you to rename tables.
	• <b>Restore Table Partition</b> : This permission allows you to export partition information from the file system and save the information to metadata.
	• View All Partitions: This permission allows you to view all partitions in a partition table.
	• View Other Users' Permissions: This permission allows you to query other users' permission on the current database.

#### 3. Click **OK**.

# Modifying Permissions for an Existing User or Project

For a user or project that has some permissions on the database, you can revoke the existing permissions or grant new ones.

#### **NOTE**

If the options in **Set Permission** are gray, the corresponding account does not have the permission to modify the database. You can apply to the administrator, database owner, or other authorized users for granting and revoking permissions of databases.

- 1. In the **User Permission Info** list, find the user whose permission needs to be set.
  - If the user is an IAM user, you can set permissions for it.
  - If the user is already an administrator, you can only view the permissions information.

In the **Project Permission Info** list, locate the project for which you want to set permissions and click **Set Permission**.

2. In the **Operation** column of the IAM user or project, click **Set Permission**. The **Set Permission** dialog box is displayed.

For details about the permissions of database users or projects, see Table 4-8.

3. Click **OK**.

#### **Revoking All Permissions of a User or Project**

Revoke all permissions of a user or a project.

• In the user list under **User Permission Info**, locate the row where the target IAM user resides and click **Revoke Permission** in the **Operation** column. In the displayed dialog box, click **OK**. In this case, the user has no permissions on the database.

#### **NOTE**

If a user is an administrator, **Revoke Permission** is gray, indicating that the user's permission cannot be revoked.

• In the **Project Permission Info** area, select the project whose permissions need to be revoked and click **Revoke Permission** in the **Operation** column. After you click **OK**, the project does not have any permissions on the database.

# 4.5.2 Deleting a Database on the DLI Console

You can delete an unused database from the DLI console: when a database is no longer needed, such as after a test database has completed testing; if a database has errors or anomalies that cannot be fixed; when there is a need to reorganize the data structure, such as by modifying table designs; or if a database is idle and has no practical use.

This section describes how to delete a database on the DLI management console.

# Precautions

- You are not allowed to delete databases or tables that are being used for running jobs.
- The administrator, database owner, and users with the database deletion permission can delete the database.

#### D NOTE

If a database or table is deleted, it cannot be recovered. Exercise caution when performing this operation.

#### **Deleting a Database**

- On the left of the management console, choose Data Management > Databases and Tables.
- 2. Locate the row where the target database locates and click **More** > **Drop Database** in the **Operation** column.

#### **NOTE**

You cannot delete databases that contain tables. To delete a database containing tables, delete the tables first.

3. In the displayed dialog box, click **Yes**.

# 4.5.3 Changing the Database Owner on the DLI Console

In practical use, developers create databases and tables, which are then handed over to testers for testing. Once testing is complete, the databases and tables are handed over to O&M personnel for experience. In this scenario, ownership of the data can be transferred to another owner by changing the database owner.

### Changing the Database Owner

You can change the owner of a database on either the **Data Management** page or the **SQL Editor** page.

- On the **Data Management** page, change the database owner.
  - a. On the left of the management console, choose **Data Management** > **Databases and Tables**.
  - b. On the **Databases and Tables** page, locate the database you want and click **More** > **Modify Database** in the **Operation** column.
  - c. In the displayed dialog box, enter a new owner name (an existing username) and click **OK**.
- Change the database owner on the **SQL Editor** page.
  - a. On the left of the management console, click SQL Editor.
  - b. In the navigation tree on the left, click **Databases**, click  $\equiv$  on the right of the database you want to modify, and choose **Modify Database** from the shortcut menu.
  - c. In the displayed dialog box, enter a new owner name (an existing username) and click **OK**.

# 4.5.4 Managing Tags

# **Tag Management**

A tag is a key-value pair that you can customize to identify cloud resources. It helps you to classify and search for cloud resources. A tag consists of a tag key and a tag value. If you use tags in other cloud services, you are advised to create the same tag (key-value pairs) for cloud resources used by the same business to keep consistency.

If your organization has configured tag policies for DLI, add tags to resources based on the policies. If a tag does not comply with the tag policies, resource creation may fail. Contact your organization administrator to learn more about tag policies.

DLI supports the following two types of tags:

- Resource tags: non-global tags created on DLI
- Predefined tags: global tags created on Tag Management Service (TMS).
   For more information about predefined tags, see Tag Management Service User Guide.

This section describes how to add, modify, and delete tags for databases and tables.

# Database Tags

- Step 1 In the navigation pane on the left, choose Data Management > Databases and Tables.
- **Step 2** Locate the row that contains the target database, and click **More** > **Tags** in the **Operation** column.
- **Step 3** The tag management page is displayed, and the tags (if there are) are displayed.
- **Step 4** On the displayed page, click **Add/Edit Tag**. The **Add/Edit Tag** dialog box is displayed.

Enter a tag key and a tag value in the text boxes and click **Add**.

 $\times$ 

#### Figure 4-7 Adding/Editing tags

#### Add/Edit Tag

It is recommended that you use TMS's predefined tag function to add the same tag to different cloud resources. View predefined tags  $\,$  C

To add a tag, enter a tag key and a tag value below.

test = test 🛛 🛞		
dli	001	Add
19 tags available for addition.		
	OK Cancel	

#### Table 4-9Tag parameters

Parameter	Description			
Tag key	You can specify the tag key in either of the following ways:			
	• Click the text box for tag key and select a predefined tag key from the drop-down list.			
	To add a predefined tag, you need to create one on TMS and then select it from the <b>Tag key</b> drop-down list. You can click <b>View predefined tags</b> to go to the <b>Predefined Tags</b> page of the TMS console. Then, click <b>Create Tag</b> in the upper corner of the page to create a predefined tag			
	For details, see <b>Creating Predefined Tags</b> in <i>Tag Management</i> Service User Guide.			
	• Enter a tag key in the text box.			
	<b>NOTE</b> A tag key can contain a maximum of 128 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed, but the value cannot start or end with a space or start with <b>_sys_</b> .			
Tag value	You can specify the tag value in either of the following ways:			
	• Click the tag value text box and select a predefined tag value from the drop-down list.			
	• Enter a tag value in the text box.			
	NOTE A tag value can contain a maximum of 255 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed.			

 $\times$ 

tag to

Add

#### **NOTE**

- A maximum of 20 tags can be added.
- Only one tag value can be added to a tag key.
- The key name in each resource must be unique.
- Step 5 Click OK. The database tag is added.

To delete a tag, click **Delete** in the **Operation** column of the target tag.

----End

### Table Tags

- Step 1 In the navigation pane on the left, choose Data Management > Databases and Tables.
- **Step 2** Click a database name to view the tables in the database.
- **Step 3** Locate the row that contains the target table and click **More** > **Tag** in the **Operation** column.
- **Step 4** The tag management page is displayed, and the tags (if there are) are displayed.
- **Step 5** On the displayed page, click **Add/Edit Tag**. The **Add/Edit Tag** dialog box is displayed.

Enter a tag key and a tag value in the text boxes and click Add

#### Figure 4-8 Adding/Editing tags

Add/Edit Tag
It is recommended that you use TMS's predefined tag function to add the same different cloud resources. View predefined tags $$
To add a tag, enter a tag key and a tag value below.
gff = ztt 💿
di 001

19 tags available for addition.



Parameter	Description			
Tag key	<ul> <li>You can specify the tag key in either of the following ways:</li> <li>Click the text box for tag key and select a predefined tag key from the drop-down list. To add a predefined tag, you need to create one on TMS and then select it from the <b>Tag key</b> drop-down list. You can click <b>View predefined tags</b> to go to the <b>Predefined Tags</b> page of the TMS console. Then, click <b>Create Tag</b> in the upper corner of the page to create a predefined tag.</li> </ul>			
	For details, see <b>Creating Predefined Tags</b> in <i>Tag Management Service User Guide</i> .			
	• Enter a tag key in the text box.			
	<b>NOTE</b> A tag key can contain a maximum of 128 characters. Only letters, numbers, spaces, and special characters (_:+-@) are allowed, but the value cannot start or end with a space or start with <b>_sys_</b> .			
Tag value	You can specify the tag value in either of the following ways:			
	• Click the tag value text box and select a predefined tag value from the drop-down list.			
	• Enter a tag value in the text box.			
	<b>NOTE</b> A tag value can contain a maximum of 255 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed.			

Table 4-10	Tag	parameters
------------	-----	------------

### D NOTE

- A maximum of 20 tags can be added.
- Only one tag value can be added to a tag key.
- The key name in each resource must be unique.

**Step 6** Click **OK**. The table tag is added.

To delete a tag, click **Delete** in the **Operation** column of the target tag.

----End

# 4.6 Managing Table Resources on the DLI Console

# 4.6.1 Configuring Table Permissions on the DLI Console

# **Operation Scenario**

- By setting permissions, you can assign varying table permissions to different users.
- The administrator and table owner have all permissions, which cannot be set or modified by other users.

• When setting table permissions for a new user, ensure that the region of the user group the user belongs to has the **Tenant Guest** permission. For details about the Tenant Guest permission and how to apply for the permission, see **Permissions Policies** and **Creating a User Group and Assigning Permissions** in the *Identity and Access Management User Guide*.

# Precautions

- By the rules in **Common Operations Supported by DLI System Policy**, you need to authorize a user to view tables in a database of the owner account.
- If you create a table with the same name as a deleted table, the table permissions will not be inherited. In this case, you need to grant the table permissions to users or projects.

For example, user A is granted with the permission to delete the **testTable** table. Delete the table and create another one with the same name. You need to grant user A the deletion permission of the **testTable** table again.

# **Viewing Table Permissions**

- On the left of the management console, choose Data Management > Databases and Tables.
- 2. Click the database name in the table whose authority is to be set. The **Table Management** page of the database is displayed.
- 3. Locate the row where the target table resides and click **Manage Permissions** in the **Operation** column.

Та	ble Management / bucket_multibk				Grant Permission
	User Permission Info				Enter a username. Q C
Username Inherited Permissions		Inherited Permissions	Non-Inherited Permissions	Column Permissions	Operation
	Add Columns/Add Partition to Partition Table/Delete Parti				Set Permission   Revolve Permission
			ALL		Set Permission   Revole Permission
	Project Permission Info				Please enter the project Q C
	Project Name	Inherited Permissions	Non-Inherited Permissions	Column Permissions	Operation
		Drop Table/Select Table/Insert/Add Columns/Overwrite/Gr			Set Permission   Revolve Permission
		Drop Table/Select Table/Insert/Add Columns/Overwrite/Gr			Set Permission   Revolve Permission

#### Figure 4-9 Table permission management

Permissions can be granted to new users or projects, modified for users or projects with existing permissions, or revoked from a user or project.

# Granting Permissions to a New User or a Project

Here, the new user or project refers to a user or a project that does not have permissions on the database.

- 1. Click the table you need. In the displayed table permissions page, click **Grant Permission** in the upper right corner.
- 2. In the displayed **Grant Permission** dialog box, select the required permissions.
  - For details about the DLI table permissions, see **Table 4-11**.

 $\times$ 

×

#### Figure 4-10 Granting DLI table permissions to a user

Grant Permission					
* Authorization Object	User	Project			
★ Username	Enter a usernam	e.			
Select the permissions to	be granted to the us	er			
Non-inheritable Permis	ssions				
Select all					
View Table		Sele	ect Table		View Table Info
View Table Creation S	Statement	Dro	p Table		Rename Table
Insert		Ove	rwrite Table		Add Column
Access Metadata		View	v Other User's Permissions		Grant Permission
Revoke Permission					
			OK Cancel		

### Figure 4-11 Granting DLI table permissions to a project

Grant Permission					
* Authorization Object	User	Project			
★ Project	(	•			
Select the permissions to t	be granted to the us	er			
Non-inheritable Permis	sions				
Select all					
Select Table		View Table Info		View Table Creation Statement	
Drop Table		Rename Table		Insert	
Overwrite Table		Add Column		Access Metadata	
View Other User's Per	missions	Grant Permission		Revoke Permission	
		OK	Cancel		

 Table 4-11
 Parameter description

Parameter	Description		
Authorizati on Object	Select <b>User</b> or <b>Project</b> .		
Username/ Project	<ul> <li>If you select User, enter the IAM username when granting table permissions to the user.</li> <li>NOTE         The username is an existing IAM user name and has logged in to the DLI management console.     </li> </ul>		
	<ul> <li>If you select <b>Project</b>, select the project to be authorized in the current region.</li> <li><b>NOTE</b>         If you select <b>Project</b>, you can only view information about the authorized tables and their databases.     </li> </ul>		

Parameter	Description				
Non- inheritable	Select a permission to grant it to the user, or deselect a permission to revoke it.				
Permissions	• The following permissions are applicable to both user and project authorization:				
	<ul> <li>Select Table: This permission allows you to query data of the current table.</li> </ul>				
	<ul> <li>View Table Information: This permission allows you to view information about the current table.</li> </ul>				
	<ul> <li>View Table Creation Statement: This permission allows you to view the statement for creating the current table.</li> </ul>				
	<ul> <li>Drop Table: This permission allows you to delete the current table.</li> </ul>				
	<ul> <li>Rename Table: Rename the current table.</li> </ul>				
	<ul> <li>Insert: This permission allows you to insert data into the current table.</li> </ul>				
	<ul> <li>Overwrite: This permission allows you to insert data to overwrite the data in the current table.</li> </ul>				
	<ul> <li>Add Column: Add columns to the current table.</li> </ul>				
	<ul> <li>Grant Permission: The current user can grant table permissions to other users.</li> </ul>				
	<ul> <li>Revoke Permission: The current user can revoke the table's permissions that other users have but cannot revoke the table owner's permissions.</li> </ul>				
	<ul> <li>View Other Users' Permissions: This permission allows you to query other users' permission on the current table.</li> </ul>				
	The partition table also has the following permissions:				
	<ul> <li>Delete Partition: This permission allows you to delete existing partitions from a partition table.</li> </ul>				
	<ul> <li>View All Partitions: This permission allows you to view all partitions in a partition table.</li> </ul>				
	<ul> <li>The following permissions can only be granted to users:</li> </ul>				
	<ul> <li>View Table: This permission allows you to display the current table.</li> </ul>				

- For details about the OBS table permissions, see **Table 4-12**.

 $\times$ 

 $\times$ 

# Figure 4-12 Granting OBS table permissions to a user

Grant Permission						
* Authorization Object	User	Project				
* Username	Enter a userna	me.				
Select the permissions to b	be granted to the	user				
Non-inheritable Permis	sions					
Select all						
View Table		Sele	ct Table	View T	able Info	
View Table Creation S	tatement	Drop	Table	Renan	ne Table	
Insert		Ove	rwrite Table	Add Co	olumn	
Access Metadata		Viev	Other User's Permissions	Grant	Permission	
Revoke Permission						
			OK Cancel			

### Figure 4-13 Granting OBS table permissions to a project

Grant Permission	1		
* Authorization Object	User	Project	
* Project		•	
Select the permissions to be	e granted to the us	er	
Non-inheritable Permiss	ions		
Select all			
Select Table		View Table Info	View Table Creation Statement
Drop Table		Rename Table	Insert
Overwrite Table		Add Column	Access Metadata
View Other User's Perr	nissions	Grant Permission	Revoke Permission
		ОК Са	ncel

 Table 4-12
 Parameter description

Paramete r	Description
Authoriza tion Object	Select <b>User</b> or <b>Project</b> .
Username /Project	<ul> <li>If you select User, enter the IAM username when granting table permissions to the user.         NOTE         The username is an existing IAM user name and has logged in to the DLI management console.     </li> <li>If you select Project, select the project to be authorized in the current region.         NOTE         If you select Project, you can only view information about the authorized tables and their databases.     </li> </ul>

Paramete r	Description
Non- inheritabl	Select a permission to grant it to the user, or deselect a permission to revoke it.
e Permissio	• The following permissions are applicable to both user and project authorization:
115	<ul> <li>View Table Creation Statement: This permission allows you to view the statement for creating the current table.</li> </ul>
	<ul> <li>View Table Information: This permission allows you to view information about the current table.</li> </ul>
	<ul> <li>Select Table: This permission allows you to query data of the current table.</li> </ul>
	<ul> <li>Drop Table: This permission allows you to delete the current table.</li> </ul>
	- Rename Table: Rename the current table.
	<ul> <li>Insert: This permission allows you to insert data into the current table.</li> </ul>
	<ul> <li>Overwrite: This permission allows you to insert data to overwrite the data in the current table.</li> </ul>
	<ul> <li>Add Column: This permission allows you to add columns to the current table.</li> </ul>
	<ul> <li>Grant Permission: This permission allows you to grant table permissions to other users or projects.</li> </ul>
	<ul> <li>Revoke Permission: This permission allows you to revoke the table's permissions that other users or projects have but cannot revoke the table owner's permissions.</li> </ul>
	<ul> <li>View Other Users' Permissions: This permission allows you to query other users' permission on the current table.</li> </ul>
	The partition table also has the following permissions:
	<ul> <li>Add Partition: This permission allows you to add a partition to a partition table.</li> </ul>
	<ul> <li>Delete Partition: This permission allows you to delete existing partitions from a partition table.</li> </ul>
	<ul> <li>Configure Path for Partition: This permission allows you to set the path of a partition in a partition table to a specified OBS path.</li> </ul>
	<ul> <li>Rename Table Partition: This permission allows you to rename partitions in a partition table.</li> </ul>
	<ul> <li>Restore Table Partition: This permission allows you to export partition information from the file system and save the information to metadata.</li> </ul>

Paramete r	Description
	<ul> <li>View All Partitions: This permission allows you to view all partitions in a partition table.</li> </ul>
	<ul> <li>The following permissions can only be granted to users:</li> <li>View Table: This permission allows you to view the current table.</li> </ul>

- For details about the view permissions, see Table 4-13.

#### D NOTE

A view can be created only by using SQL statements. You cannot create a view on the **Create Table** page.

#### Figure 4-14 Granting view permissions to a user

Grant Permission					
* Authorization Object	User	Project			
* Username	Enter a username.				
Select the permissions to be	e granted to the user				
Non-Inherited Permi	ssions				
Select all					
View Table		View	Table Info		View Table Creation Statement
Drop Table		Select	t Table		Rename Table
View Other User's Pern	nissions	Grant	Permission		Revoke Permission
			ОК	Cancel	

### Figure 4-15 Granting view permissions to a project

Grant Permission			
* Authorization Object	User	Project	
* Enterprise Project		•	
Select the permissions to b	e granted to the use	er.	
Non-Inherited Permi	issions		
Select all			
View Table Info		View Table Creation Statement	t Drop Table
Select Table		Rename Table	View Other User's Permissions
Grant Permission		Revoke Permission	
		<b>OK</b> Cancel	

Parameter	Description
Authorizatio n Object	Select <b>User</b> or <b>Project</b> .
Username/ Project	<ul> <li>If you select User, enter the IAM username when adding a user to the database.</li> <li>NOTE         <ul> <li>The username is an existing IAM user name and has logged in to the DLI management console.</li> </ul> </li> <li>If you select Project, select the project to be authorized in the current region.</li> <li>NOTE         <ul> <li>If you select Project, you can only view information about the authorized tables and their databases.</li> </ul> </li> </ul>
Non- inheritable Permissions	<ul> <li>Select a permission to grant it to the user, or deselect a permission to revoke it.</li> <li>The following permissions are applicable to both user and project authorization: <ul> <li>View Table Information: This permission allows you to view information about the current table.</li> <li>View Table Creation Statement: This permission allows you to view the statement for creating the current table.</li> <li>Drop Table: This permission allows you to delete the current table.</li> <li>Select Table: This permission allows you to query data of the current table.</li> <li>Rename Table: Rename the current table.</li> <li>Grant Permission: The current user or project can grant table permissions to other users or project can revoke the table's permissions that other users</li> </ul> </li> </ul>
	<ul> <li>can revoke the table's permissions that other users or projects have but cannot revoke the table owner's permissions.</li> <li>View Other Users' Permissions: This permission allows you to query other users' permission on the current table.</li> <li>Only applicable to <ul> <li>View Table: This permission allows you to view</li> </ul> </li> </ul>

 Table 4-13 Parameter description

### 3. Click OK.

# Modifying Permissions for an Existing User or Project

For a user or project that has some permissions on the database, you can revoke the existing permissions or grant new ones.

#### **NOTE**

If all options under **Set Permission** are gray, you are not allowed to change permissions on this table. You can apply to the administrator, table owner, or other authorized users for granting and revoking table permissions.

- 1. In the **User Permission Info** list, find the user whose permission needs to be set.
  - If the user is an IAM user and is not the owner of the table, you can set permissions.
  - If the user is an administrator or table owner, you can only view permissions.

In the **Project Permission Info** list, locate the project for which you want to set permissions and click **Set Permission**.

- 2. In the **Operation** column of the IAM user or project, click **Set Permission**. The **Set Permission** dialog box is displayed.
  - For details about DLI table user or project permissions, see **Table 4-11**.
  - For details about OBS table user or project permissions, see **Table 4-12**.
  - For details about View table user or project permissions, see Table 4-13.
- 3. Click **OK**.

#### **Revoking All Permissions of a User or Project**

Revoke all permissions of a user or a project.

• In the user list under **User Permission Info**, locate the row where the target IAM user resides and click **Revoke Permission** in the **Operation** column. In the displayed dialog box, click **OK**. In this case, the user has no permissions on the table.

#### **NOTE**

In the following cases, **Revoke Permission** is gray, indicating that the permission of the user cannot be revoked.

- The user is an administrator.
- The IAM user is the owner of the table.
- The IAM user has only inheritable permissions.
- In the **Project Permission Info** area, select the project whose permissions need to be revoked and click **Revoke Permission** in the **Operation** column. After you click **OK**, the project does not have any permissions on the table.

#### **NOTE**

If a project has only inheritable permissions, **Revoke Permission** is gray, indicating that the permissions of the project cannot be revoked.

# 4.6.2 Deleting a Table on the DLI Console

You can delete an unused data table from the DLI console: when a data table is no longer needed, such as after a test data table has completed testing; if a data table has errors or anomalies that cannot be fixed; when there is a need to reorganize the data structure, such as by modifying table designs; or if a data table is idle and has no practical use.

This section describes how to delete a data table on the DLI management console.

# Precautions

- Databases or tables that are being used for running jobs cannot be deleted.
- Only administrators, table owners, and users with table deletion permission can delete tables.

#### **NOTE**

Deleted data tables cannot be restored. Exercise caution when performing this operation.

# Deleting a Table

You can delete a table on either the **Data Management** page or the **SQL Editor** page.

- On the **Data Management** page
  - a. In the navigation pane on the left of the console, choose **Data Management > Databases and Tables**.
  - b. Locate the database whose tables you want to delete and click its name.
  - c. On the displayed page, select the table you want to delete, click **More** in the **Operation** column, and select **Drop Table**.
  - d. In the displayed dialog box, click **Yes**.
- On the **SQL Editor** page
  - a. In the navigation pane on the left of the console, choose **SQL Editor**.
  - b. In the middle pane, click the **Databases** tab. Click the name of the database whose tables you want to delete.
  - c. Click  $\equiv$  next to the table to delete and select **Delete**.
  - d. In the **Drop Table** dialog box that appears, click **OK**.

# 4.6.3 Changing the Table Owner on the DLI Console

In practical use, developers create databases and tables, which are then handed over to testers for testing. Once testing is complete, the databases and tables are handed over to O&M personnel for experience. In this scenario, ownership of the data can be transferred to another owner by changing the table owner.

# Changing the Table Owner

- In the navigation pane on the left of the console, choose Data Management
   > Databases and Tables.
- 2. Click the name of the database whose tables you want to modify.
- 3. Locate the table for which you want to change the owner, click **More** in the **Operation** column, and select **Modify Owner**.

4. In the displayed **Modify Owner** dialog box, enter a new owner name (an existing username) and click **OK**.

# 4.6.4 Importing OBS Data to a DLI Table

This section describes how to import data stored in OBS to a table on the DLI console. Data can be imported to both DLI tables (table type: MANAGED) and OBS tables (table type: EXTERNAL).

# Precautions

- Only one path can be specified during data import. The path cannot contain commas (,).
- To import data in CSV format to a partitioned table, place the column to be partitioned in the last column of the data source.
- You are advised not to concurrently import data in to a table. If you concurrently import data into a table, there is a possibility that conflicts occur, leading to failed data import.
- The imported file can be in CSV, Parquet, ORC, JSON, and Avro format. The encoding format must be UTF-8.

# Prerequisites

The data to be imported has been stored on OBS.

# Procedure

- **Step 1** You can import data on either the **Data Management** page or the **SQL Editor** page.
  - To import data on the **Data Management** page:
    - a. On the left of the management console, choose **Data Management** > **Databases and Tables**.
    - b. Click the name of the database corresponding to the table where data is to be imported to switch to the table management page.
    - c. Locate the row where the target table resides and choose **More** > **Import** in the **Operation** column. The **Import** dialog box is displayed.

#### Figure 4-16 Importing data

Table Name	Type V Owner						
		Storage Loc V	Rize	Data Path	Created	Last Accessed	Operation
	MANAGED	DU	0.8		Jan 05, 2023 20:31:25 GMT+08:00	Jan 10, 2023 19:21:42 GMT+06:00	Permissions More +
	MANAGED	DU	716 B		Jan 06, 2023 20:31:25 GMT+00:00	Jan 29, 2023 11:34:02 GMT-00:00	Drop Table Modify Owner
	MANAGED	DLI	0.8		Mar 20, 2023 11:45:05 GMT+05:00	Mar 20, 2023 11:45:08 GMT+08:00	Import
	MANAGED	DU	0.0		Jan 10, 2023 16:37:56 OMT+00:00	Jan 10, 2023 16:37:56 GMT+00:00	Export
	MANAGED	DLI	0 8		Jan 10, 2023 10:39:19 GMT+0E:00	Jan 10, 2023 18:39:19 GMT+08:00	Tage
	EXTERNAL	685	0.8		Mar 18, 2023 17:29:35 GMT+08:00	Mar 18, 2023 17:29:35 GMT+08:00	Permissions   More +
	EXTERNAL	089	0.0		Jan 17, 2023 11:08:31 GMT+08:00	Jan 17, 2023 11:10:54 GMT+00:00	Permissions   More +
	EXTERNAL	OBS	0 8		Jan 17, 2023 11:12:29 GMT+08:00	Jan 17, 2023 11:12:29 GMT+08:00	Permissions   More +
	EXTERNAL	089	0.0		Jan 17, 2023 11:13:27 GMT+00:00	Jan 17, 2023 11:13:27 OMT+08:00	Permissions   More +
	EXTERNAL	CHN	0 8		Mar 10, 2023 10:44 22 GMT+08:00	Mar 10, 2023 10:44:22 GMT+0E:00	Permissions More +

- To import data on the **SQL Editor** page:
  - a. On the left of the management console, click **SQL Editor**.
  - b. In the navigation tree on the left of **SQL Editor**, click **Databases** to see all databases. Click the database where the target table belongs. The table list is displayed.

c. Click  $\equiv$  on the right of the table and choose **Import** from the shortcut menu. The **Import** page is displayed.

Databases Queues	· 🗸		
Databases /		1	SELECT doc_id FR
Table(58)	⊕ C	2	SELECT name FROM
Enter a name.	Q		
External(26)			
-	≡		
	Delete		
	Import		
	View Pr	operties	
	=		

Figure 4-17 SQL editor - importing data

**Step 2** In the **Import** dialog box, set the parameters based on **Table 4-14**.

Table	4-14	Description	
Tuble	<b>T</b> 1 <b>T</b>	Description	

Parameter	Description	Example
Databases	Database where the current table is located.	-
Table Name	Name of the current table.	-
Queues	Queue where the imported data will be used	-
File Format	Format of the data source file to be imported. The CSV, Parquet, ORC, JSON, and Avro formats are supported. Encoding format. Only UTF-8 is supported.	CSV

Parameter	Description	Example
Path	You can directly enter a path or click and select an OBS path. If no bucket is available, you can directly switch to the OBS management console and create an OBS bucket. The path can be a file or folder.	obs://DLI/ sampledat a.csv
	<ul> <li>When creating an OBS table, you must specify a folder as the directory. If a file is specified, data import will fail.</li> <li>When there are both a folder and a file with</li> </ul>	
	the same name in the OBS directory, importing data pointing to that path will prioritize the file over the folder.	
Table Header: No/Yes	This parameter is valid only when <b>File Format</b> is set to <b>CSV</b> . Whether the data source to be imported contains the table header.	-
	Click <b>Advanced Settings</b> and select the checkbox next to <b>Table Header: No</b> . If the checkbox is selected, the table header is displayed. If the checkbox is deselected, no table header is displayed.	
User-defined Delimiter	This parameter is valid only when <b>File Format</b> is set to <b>CSV</b> and you select <b>User-defined Delimiter</b> .	Default value:
	The following delimiters are supported:	comma (,)
	• Comma (,)	
	• Vertical bar ( )	
	• Tab character (\t)	
	Others: Enter a user-defined delimiter.	
User-defined Quotation Character	This parameter is valid only when <b>File Format</b> is set to <b>CSV</b> and <b>User-defined Quotation</b> <b>Character</b> is selected.	Default value: double
	The following quotation characters are supported:	quotation
	Single quotation mark (')	()
	<ul> <li>Double quotation marks (")</li> </ul>	
	Others: Enter a user-defined quotation character.	
User-defined Escape Character	This parameter is valid only when <b>File Format</b> is set to <b>CSV</b> and you select <b>User-defined Escape Character</b> .	Default value: backslash
	The following escape characters are supported:	(\)
	• Backslash (\)	
	• Others: Enter a user-defined escape character.	

Parameter	Description	Example
Date Format	This parameter is valid only when <b>File Format</b> is set to <b>CSV</b> or <b>JSON</b> .	2000-01-0 1
	This parameter specifies the format of the date in the table and is valid only <b>Advanced Settings</b> is selected. The default value is <b>yyyy-MM-dd</b> . For definition of characters involved in the date pattern, see Table 3 in <b>Importing Data to the</b> <b>Table</b> .	
Timestamp Format	This parameter is valid only when <b>File Format</b> is set to <b>CSV</b> or <b>JSON</b> .	2000-01-0 1 09:00:00
	This parameter specifies the format of the timestamp in the table and is valid only <b>Advanced Settings</b> is selected. The default value is <b>yyyy-MM-dd HH:mm:ss</b> . For definition of characters involved in the time pattern, see Table 3 in <b>Importing Data to the Table</b> .	
Error Records Path	This parameter is valid only when <b>File Format</b> is set to <b>CSV</b> or <b>JSON</b> .	obs://DLI/
	The parameter specifies the error data is stored in the corresponding OBS path and is valid only <b>Advanced Settings</b> is selected.	

#### Step 3 Click OK.

**Step 4** You can view the imported data in either of the following ways:

**NOTE** 

Currently, only the first 10 records are displayed.

- Choose Data Management > Databases and Tables in the navigation pane of the console. Locate the row that contains the database where the target table belongs and click More > View Properties in the Operation column. In the displayed dialog box, click the Preview tab to view the imported data.
- On the **Databases** tab of the **SQL Editor**, click the database name to go to the table list. Click  $\equiv$  on the right of a table name and choose **View Properties** from the shortcut menu. In the displayed dialog box, click **Preview** to view the imported data.
- Step 5 (Optional) View the status and execution result of the importing job on the Job Management > SQL Jobs page.

----End

# 4.6.5 Exporting DLI Table Data to OBS

You can export data from a DLI table to OBS. During the export, a folder is created in OBS or the content in the existing folder is overwritten.
#### Precautions

- The exported file can be in JSON format, and the text format can only be UTF-8.
- Data in DLI tables (table type: MANAGED) can only be exported to OBS buckets, and the export path must be a folder.
- Data can be exported across accounts. That is, after account B authorizes account A, account A can export data to the OBS path of account B if account A has the permission to read the metadata and permission information about the OBS bucket of account B and read and write the path.

#### Procedure

- **Step 1** You can export data on either the **Data Management** page or the **SQL Editor** page.
  - To export data on the **Data Management** page:
    - a. On the left of the management console, choose **Data Management** > **Databases and Tables**.
    - b. Click the name of the database corresponding to the table where data is to be exported to switch to the **Manage Tables** page.
    - c. Select the corresponding table (DLI table) and choose **More** > **Export** in the **Operation** column. The **Export Data** page is displayed.
  - To export data on the **SQL Editor** page:
    - a. On the left of the management console, click **SQL Editor**.
    - b. In the navigation tree on the left, click **Databases** to see all databases. Click the database name corresponding to the table to which data is to be exported. The tables are displayed.
    - c. Click  $\equiv$  on the right of the managed table (DLI table) whose data is to be exported, and choose **Export** from the shortcut menu.

#### Figure 4-18 Exporting the table

Databases Que	ues · 🗸				
Databases /		1			
Table(58)	⊕ C				
Enter a name.	Q				
+ External(26)					
Managed(32)					
+	Ξ				
- +	Dele	te			
- +	Impo	ort			
+ Export					
- +	View	Properties			
( <del>+</del> )	_	-	1, Column 1		

**Step 2** In the displayed **Export Data** dialog box, specify parameters by referring to **Table 4-15**.

Export Result			×
★ Databases	db1		
★ Table Name	src		
★ Data Format	json		
Queues		•	
Compression Format	bzip2	•	
★ Storage Path	This path must start with "obs://"		?
★ Export Mode (?)	New OBS directory	Existing OBS directory (Overwirtten)	
Table Header	No Yes		
	ОК	Cancel	

#### Figure 4-19 Exporting data

Paramet er	Description				
Databas es	Database where the current table is located.				
Table Name	Name of the current table.				
Data Format	Format of the file storing data to be exported. Formats other than JSON will be supported in later versions.				
Queue	Select a queue.				
Compres sion	Compression format of the data to be exported. The following compression formats are supported:				
Format	• none				
	• bzip2				
	• deflate				
	• gzip				
Storage	• Enter or select an OBS path.				
Path	• The export path must be a folder that does not exist in the OBS bucket. Specifically, you need to create a folder in the target OBS directory.				
	<ul> <li>The folder name cannot contain the special characters of \/:*?</li> <li>"&lt;&gt; , and cannot start or end with a dot (.).</li> </ul>				
Export	Storage mode of the data to be exported.				
Mode	• <b>New OBS directory</b> : If the specified export directory exists, an error is reported and the export operation cannot be performed.				
	• Existing OBS directory (Overwritten): If you create a file in the specified directory, the existing file will be overwritten.				
Table Header: No/Yes	Whether the data to be exported contains the table header.				

#### Table 4-15 Parameter description

#### Step 3 Click OK.

- Step 4 (Optional) You can view the job status (indicated by Status), statements (indicated by Statement), and other information about exporting jobs on the SQL Jobs page.
  - 1. Select **EXPORT** from the **Job Type** drop-down list box and specify the time range for exporting data. The jobs meeting the requirements are displayed in the job list.
  - 2. Click  $\checkmark$  to view details about an exporting job.

----End

# 4.6.6 Previewing Table Data on the DLI Console

The **Preview** page displays the first 10 records in the table.

#### Procedure

You can preview data on either the **Data Management** page or the **SQL Editor** page.

- To preview data on the **Data Management** page:
  - a. On the left of the management console, choose **Data Management** > **Databases and Tables**.
  - b. On the displayed **Data Management** page, click the name of the database where the target table whose data you want to export resides to switch to the **Manage Tables** page.
  - c. Click **More** in the **Operation** column of the target table and select **View Properties**.
  - d. Click the **Preview** tab to preview the table data.
- To preview data on the **SQL Editor** page:
  - a. On the left of the management console, click **SQL Editor**.
  - b. In the navigation pane of the displayed **SQL Editor** page, click **Databases**.
  - c. Click the corresponding database name to view the tables in the database.
  - d. Click  $\equiv$  on the right of the corresponding table, choose **View Properties** from the list menu, and click the **Preview** tab to preview the data of the table.

# 4.7 Creating and Using LakeFormation Metadata

# 4.7.1 Connecting DLI to LakeFormation

#### Scenario

LakeFormation is an enterprise-level, all-in-one lakehouse construction service that provides unified metadata management capabilities. It supports seamless integration with various compute engines and big data cloud services, enabling you to efficiently build data lakes and operate related businesses, thereby accelerating the extraction of business data value.

In Spark and SQL job scenarios, LakeFormation can be connected to achieve unified metadata management. This section will guide you through the steps to configure the data connection between DLI and LakeFormation.

For the Spark syntax of LakeFormation, see Spark Syntax Reference.

For the Flink syntax of LakeFormation, see Flink Syntax Reference.

#### Notes

To use this function, which is currently in the whitelist, submit a request by choosing **Service Tickets** > **Create Service Ticket** in the upper right corner of the management console.

Connecting DLI to LakeFormation depends on the availability of the LakeFormation service. To understand the availability scope of LakeFormation, refer to **Global Products and Services**.

#### Procedure

#### Figure 4-20 Procedure

[	LakeFormation ma	nagement console	 	DLI r	management consol	e		
	Create a LakeFormation instance.	← a catalog.	Create a DLI catalog.	<b>→</b>	Authorize DLI to use LakeFormation resources.	<b>_</b>	Create and submit a SQL job and use LakeFormation metadata.	
1								

#### Notes and Constraints

• **Table 4-16** lists the DLI queue and engine types that allow you to connect DLI to LakeFormation to obtain metadata.

For the engine type and version of a queue, see **Viewing Basic Information About a Queue**.

**Table 4-16** Queue and engine types that allow for the connection of DLI toLakeFormation for metadata retrieval

Queue Type	Engine Type and Version
default queue	• Spark 3.3. <i>x</i> : can connect to LakeFormation to obtain metadata.
	• HetuEngine 2.1.0: can connect to LakeFormation to obtain metadata.
For SQL	• Spark 3.3. <i>x</i> : can connect to LakeFormation to obtain metadata.
	• HetuEngine 2.1.0: can connect to LakeFormation to obtain metadata.
For general purpose	Flink job scenario: Flink 1.15 or later supports integration with LakeFormation to obtain metadata when using queues within an elastic resource pool.

- DLI can only connect to the default LakeFormation instance. Set the instance in LakeFormation as the default to ensure successful connection.
- DLI can read data in Avro, JSON, Parquet, CSV, ORC, Text, and Hudi formats from LakeFormation.
- LakeFormation manages the permissions for databases and tables in the data catalog of LakeFormation.

• After connecting DLI to LakeFormation, the original databases and tables in DLI will be moved to the data directory of DLI.

#### Step 1: Create a LakeFormation Instance for Metadata Storage

LakeFormation instances provide basic resources for metadata management. DLI can only connect to the default LakeFormation instance.

#### 1. Creating an instance

- a. Log in to the LakeFormation management console.
- b. Click **Buy Now** or **Buy Instance** in the upper right corner of the page.

If there are no instances available on the page, **Buy Now** is displayed. If there are any LakeFormation instances on the page, **Buy Instance** is displayed.

c. Set LakeFormation instance parameters as needed to complete instance creation.

In this example, we create a pay-per-use shared instance.

For parameter settings and descriptions, see **Creating a LakeFormation Instance**.

#### 2. Setting a LakeFormation instance as the default

- a. View the value of **Default Instance** in the **Basic Information** area.
  - **true**: The instance is the default.
  - **false**: The instance is not the default.
- b. To set an instance as the default, click **Set as Default** in the upper right corner of the page.
- c. In the dialog box that appears, select I understand the consequences of changing the default instance and have still decided to perform this operation. and click OK.

#### D NOTE

DLI can currently only connect to the default LakeFormation instance. Changing an instance to the default may impact the services that use LakeFormation. Exercise caution when performing this operation.

#### Step 2: Create a Catalog on the LakeFormation Management Console

A data catalog is a metadata management object that can contain multiple databases. You can create and manage multiple catalogs in LakeFormation to isolate metadata of different external clusters.

- 1. Log in to the LakeFormation management console.
- 2. In the navigation pane on the left, choose **Metadata** > **Catalog**.
- 3. On the displayed page, click **Create**.

Set catalog instance parameters as needed.

For parameter settings and descriptions, see **Creating a Catalog**.

4. Once created, you can view information about the created catalog on the **Catalog** page.

#### Step 3: Create a Data Catalog on the DLI Management Console

On the DLI management console, you need to create a link to the catalog to access the catalog stored in the LakeFormation instance.

#### **NOTE**

- DLI can only connect to the default LakeFormation instance. Set the instance in LakeFormation as the default to ensure successful connection.
- You can only create one mapping for each data catalog in LakeFormation. For example, a user creates a mapping named **catalogMapping1** in DLI, which corresponds to the data catalog **catalogA** in LakeFormation. Once created, you cannot create a mapping to **catalogA** in the same project space.
- 1. Log in to the DLI management console.
- 2. In the navigation pane on the left, choose **SQL Editor**.
- 3. The **Catalog** tab of the SQL editor is displayed.
- 4. Click  $\bigcirc$  to create a data catalog.
- 5. In the **Create Catalog** dialog box, set data catalog parameters.

 Table 4-17 Data catalog parameters

Parameter	Mand atory	Description
External Catalog Name	Yes	Catalog name of the default LakeFormation instance.
Туре	Yes	Currently, the only available option is <b>LakeFormation</b> . This option is fixed and does not need to be selected.
Catalog Name	Yes	Catalog mapping name used in DLI. When running SQL statements, you need to specify the catalog mapping to identify the external metadata to be accessed. You are advised to set this parameter to the same value as <b>External Catalog Name</b> . Currently, DLI can only connect to the data catalog of the default LakeFormation
		instance.
Description	No	Description of the data catalog.

6. Click OK.

#### Step 4: Authorize to Use LakeFormation Resources

#### • SQL job scenarios

Before submitting a SQL job, you need to authorize DLI to access LakeFormation resources such as metadata, databases, tables, columns, and functions, to ensure that the job can access required data and resources during execution. **Supported Actions for LakeFormation SQL Resources** describes LakeFormation resources and corresponding permissions.

To use LakeFormation resources, you need to separately complete IAM finegrained authorization and LakeFormation SQL resource authorization.

 IAM fine-grained authorization for LakeFormation: Authorize DLI to use LakeFormation APIs.

IAM offers multiple methods to manage the permissions of users, groups, and roles to access resources. You can create policies on the IAM console to define which users or roles can call LakeFormation APIs. Then, attach these policies to the specified users or roles.

Method 1: Role-based authorization

A coarse-grained authorization strategy that defines permissions by job responsibility. Only a limited number of service-level roles are available for authorization.

For example, grant users the read-only permission to query LakeFormation metadata resources by referring to LakeFormation Permissions Management.

Alternatively, grant all operation permissions on LakeFormation-related metadata resources.

Example:

```
{
    "Version": "1.1",
    "Statement": [
        {
            "Effect": "Allow",
            "Action": [
                "lakeformation:table:*",
                "lakeformation:catalog:*",
                "lakeformation:function:*",
                "lakeformation:transaction:*",
                "lakeformation:policy:describe",
                "lakeformation:credential:describe"
            ]
        }
    ]
}
```

Method 2: Policy-based fine-grained authorization

A policy defines the permissions required to perform actions on a specific cloud resource under certain conditions.

For LakeFormation permissions policies, see LakeFormation Permissions and Supported Actions.

For how to grant permissions, see **Creating a User and Granting LakeFormation Permissions**.

 LakeFormation SQL resource authorization: authorizes users to use specific LakeFormation resources (such as metadata, databases, tables, columns, and functions).

Users are allowed to access specific resources. This controls access to LakeFormation data and metadata.

Two methods:

 Method 1: Authorize access to resources on the LakeFormation management console.

For LakeFormation SQL resource permissions, see **Data Permissions**.

Method 2: Run the GRANT SQL statement on the DLI management console.

The **GRANT** statement is used for authorization in SQL.

You can run the **GRANT** statement to grant users or roles the permission to access databases, tables, columns, and functions.

**Supported Actions for LakeFormation SQL Resources** describes authorization policies for LakeFormation resources.

#### D NOTE

Currently, you cannot authorize access to catalogs by running the **GRANT** statement on the DLI console. To authorize access, use **method 1**.

- Spark Jar, Flink OpenSource SQL, and Flink Jar job scenarios:
  - Method 1: Use agency authorization. Before using Spark 3.3.1 or later and Flink 1.15 to run jobs, you need to create an agency on the IAM console and add the new agency information when configuring the job.

For agency permission examples, see **Creating a Custom DLI Agency** and **Agency Permission Policies in Common Scenarios**.

- Method 2: Use DEW for authorization.
  - You have granted the IAM user the IAM and LakeFormation permissions. For details, see IAM authorization in SQL job scenarios.
  - A shared secret has been created on the DEW console and the secret value has been stored. For details, see **Creating a Shared Secret**.
  - An agency has been created and authorized for DLI to access DEW. The agency must have been granted the following permissions:
    - Permission of the ShowSecretVersion interface for querying secret versions and secret values in DEW: csms:secretVersion:get.
    - Permission of the ListSecretVersions interface for listing secret versions in DEW: csms:secretVersion:list.
    - Permission to decrypt DEW secrets: kms:dek:decrypt

For agency permission examples, see **Creating a Custom DLI Agency** and **Agency Permission Policies in Common Scenarios**.

#### Step 5: Use LakeFormation Metadata During DLI Job Development

After connecting DLI to the default LakeFormation instance and authorizing access to LakeFormation resources, you can use LakeFormation metadata during DLI job development.

• DLI SQL:

For the SQL syntax of LakeFormation, see **Data Lake Insight Spark SQL** Syntax Reference. When running a SQL job, you can select the catalog where the SQL statement is located on the console, as shown in **Figure 4-21**, or specify **catalogName** in the SQL statement. **catalogName** is the mapping name of the data catalog on the DLI console.

Figure 4-21 Selecting a data catalog on the SQL editor page

Catalog Queues	Tei …				Engine	Spark	~	Queues	default	~	Catalog	hive	~	Databases	default	~
	0 Đ	1	show tables	]	(											
Enter a name.	Q															
_																
🗏 hive	Ξ															

#### D NOTE

- When connecting DLI to a LakeFormation instance, you need to specify the OBS path for storing the database when creating it.
- When connecting DLI to a LakeFormation instance, you cannot set table lifecycle and versioning when creating a table.
- When connecting DLI to a LakeFormation instance, the LOAD DATA statement does not support datasource tables, and partitions must be specified if the statement is used to import data into a partitioned table.
- If the databases and tables created on the LakeFormation console contain Chinese characters, operations on them cannot be performed on DLI.
- When connecting DLI to a LakeFormation instance, you cannot specify filter criteria to delete partitions.
- When connecting DLI to a LakeFormation instance, you cannot create Truncate Datasource or Hive foreign tables.
- DLI does not currently support the use of LakeFormation row filter criteria.
- When DLI reads binary data for console display, it converts the binary data to Base64.
- DLI does not currently support authorizing access to LakeFormation paths.
- DLI Spark Jar:

This part explains how to configure LakeFormation metadata when submitting a Spark Jar job on the DLI management console.

- Example Spark Jar
  - SparkSession spark = SparkSession.builder() .enableHiveSupport() .appName("java\_spark\_demo") .getOrCreate();

spark.sql("show databases").show();

- Configuring a Spark Jar job on the DLI management console
  - (Recommended) Method 1: Configure a Spark Jar job's access to LakeFormation metadata using the parameters (such as agency and metadata source) provided on the console

When creating or editing a Spark Jar job, refer to **Table 4-18** to configure the job's access to LakeFormation metadata.

Table 4-18	Configuring a	ı Spark Jar	job's access	to LakeForm	ation
metadata					

Paramete r	Description	Exampl e Value
Spark Version	Spark 3.3. <i>x</i> or later can connect to LakeFormation.	3.3.1
Agency	Before using Spark 3.3.1 or later to run jobs, you need to create an agency on the IAM console and add the new agency information. Once set, the system automatically adds the following agency configuration to your job: spark.dli.job.agency.name= <i>agency</i>	-
	For agency permission examples, see Creating a Custom DLI Agency and Agency Permission Policies in Common Scenarios.	
Access Metadata	Whether to allow the Spark job to access metadata.	Yes
Metadata Source	Type of metadata the Spark job accesses. In this scenario, select <b>lakeformation</b> . Once set to <b>lakeformation</b> , the system automatically adds the following configurations to your job to load LakeFormation dependencies. spark.sql.catalogImplementation=hive spark.hadoop.hive- ext.dlcatalog.metastore.client.enable=true spark.hadoop.hive- ext.dlcatalog.metastore.session.client.class=com.huawei.cl oud.dalf.lakecat.client.hiveclient.LakeCatMetaStoreClient og // Load LakeFormation dependencies. spark.driver.extraClassPath=/usr/share/extension/dli/ spark-jar/lakeformation/* spark.executor.extraClassPath=/usr/share/extension/dli/ spark-jar/lakeformation/* You can also set the metadata source in the <b>Spark Arguments(conf)</b> parameter. This means if you set the metadata source in both <b>Metadata Source</b> and <b>Spark</b> <b>Arguments(conf)</b> , the system preferentially uses the information configured in <b>Spark Arguments(conf)</b> . You are advised to set the metadata source through <b>Metadata Source</b> .	LakeFo rmatio n

Paramete r	Description	Exampl e Value
Catalog Name	Name of the data catalog the Spark job accesses. The data catalog created on the DLI management console is selected, that is, the mapping between DLI and the data catalog of the default LakeFormation instance. This data catalog is connected to the data catalog of the default LakeFormation instance. To specify other LakeFormation instances, configure the LakeFormation instances and data catalogs to be connected in <b>Spark Arguments( conf)</b> . For details, see <b>Method 2: Configure</b> <b>a Spark Jar job's access to LakeFormation</b> <b>metadata using the Spark Arguments( conf) parameter</b> .	-
	Once set to LakeFormation, the system automatically adds the following configuration to your job to connect to the data catalog of the default LakeFormation instance: spark.hadoop.lakecat.catalogname.default=lfcatalog You can also set the data catalog name in the Spark Arguments(conf) parameter. This means if you set the data catalog name in both Catalog Name and Spark Arguments(conf), the system preferentially uses the information	
	configured in <b>Spark Arguments(conf)</b> . You are advised to set the data catalog name through <b>Catalog Name</b> .	

Paramete r	Description	Exampl e Value
Spark Argument s(conf)	You can also set the metadata source and data catalog name in the <b>Spark</b> <b>Arguments(conf)</b> parameter. This means if you set the metadata source and data catalog name in both <b>Metadata Source</b> and <b>Catalog Name</b> , as well as <b>Spark</b> <b>Arguments(conf)</b> , the system preferentially uses the information configured in <b>Spark Arguments(conf)</b> .	-
	<ul> <li>To access a Hudi data table, add the following configurations to the Spark Arguments(conf) parameter: spark.sql.extensions=org.apache.spark.sql.hudi.Hoodie SparkSessionExtension spark.hadoop.hoodie.write.lock.provider=org.apache.h udi.lakeformation.LakeCatMetastoreBasedLockProvid- er</li> </ul>	
	<ul> <li>To access a Delta data table, add the following configurations to the Spark Arguments(conf) parameter: spark.sql.catalog.spark_catalog=org.apache.spark.sql.d elta.catalog.DeltaCatalog spark.sql.extensions=io.delta.sql.DeltaSparkSessionEx- tension</li> </ul>	

#### Method 2: Configure a Spark Jar job's access to LakeFormation metadata using the Spark Arguments(--conf) parameter

When creating or editing a Spark Jar job, configure the following information in the **Spark Arguments(--conf)** parameter on the job configuration page to access LakeFormation metadata:

spark.sql.catalogImplementation=hive

spark.hadoop.hive-ext.dlcatalog.metastore.client.enable=true

spark.hadoop.hiveext.dlcatalog.metastore.session.client.class=com.huawei.cloud.dalf.lakecat.client.hiveclient.L akeCatMetaStoreClient

spark.sql.extensions=org.apache.spark.sql.hudi.HoodieSparkSessionExtension // Hudi is supported, which is optional.

spark.hadoop.hoodie.write.lock.provider=org.apache.hudi.lakeformation.LakeCatMetastoreB asedLockProvider // Hudi is supported, which is optional.

// Access using an agency with the OBS and LakeFormation permissions. You are advised to configure the minimum permission set.

spark.dli.job.agency.name=agencyForLakeformation

// ID of the LakeFormation instance to be accessed, which can be viewed on the LakeFormation console. This parameter is optional. If unset, the default LakeFormation instance is accessed.

spark.hadoop.lakeformation.instance.id=xxx

// Name of the catalog to be accessed on the LakeFormation side, which can be viewed on the LakeFormation console. This parameter is optional. If unset, the default value is hive.

spark.hadoop.lakecat.catalogname.default=lfcatalog

// Load LakeFormation dependencies.

spark.driver.extraClassPath=/usr/share/extension/dli/spark-jar/lakeformation/\* spark.executor.extraClassPath=/usr/share/extension/dli/spark-jar/lakeformation/\*

• DLI Flink OpenSource SQL

Example 1: Connecting DLI to LakeFormation using an agency

Create a Flink OpenSource SQL job and	d set the following parameters:
---------------------------------------	---------------------------------

Parameter	Description	Exampl e Value
Flink Version	Flink 1.15 or later can connect to LakeFormation.	1.15
Agency	Before using Flink 1.15 or later to run jobs, you need to create an agency on the IAM console and add the new agency information. Once set, the system automatically adds the following agency configuration to your job: flink.dli.job.agency.name= <i>agency</i> For agency permission examples, see Creating a Custom DLI Agency and Agency Permission Policies in Common Scenarios.	-
Enable Checkpoint ing	Select it.	Select it.
Runtime Configurati on	<ul> <li>Type of metadata the Flink job accesses. In this scenario, select lakeformation. flink.dli.job.catalog.type=lakeformation</li> <li>Name of the data catalog the Flink job accesses. flink.dli.job.catalog.name=[Catalog name in LakeFormation]</li> <li>The data catalog created on the DLI management console is selected, that is, the mapping between DLI and the data catalog of the default LakeFormation instance. This data catalog is connected to the data catalog of the default LakeFormation instance.</li> </ul>	-

For the catalog parameters in the example, see **Table 4-19**.

Parameter	Description	Mandatory	Example Value
type	Catalog type	Yes	Fixed at hive
hive-conf-dir	hive-conf path, which is fixed at <b>/opt/flink/</b> conf.	Yes	Fixed at /opt/flink/ conf

Table 4-19	Catalog	parameters	in the	example	Flink O	penSource	SOL
	catatog	parameters		example		pensource	~~-

Parameter	Description	Mandatory	Example Value
default- database	Default database name	No	Default database

```
CREATE CATALOG hive
WITH
 (
  'type' = 'hive',
  'hive-conf-dir' = '/opt/flink/conf', -- Fixed at /opt/flink/conf
  'default-database'='default'
 );
USE CATALOG hive;
CREATE TABLE IF NOT EXISTS
 dataGenSource612 (user_id string, amount int)
WITH
 (
  'connector' = 'datagen',
  'rows-per-second' = '1',
'fields.user_id.kind' = 'random',
  'fields.user_id.length' = '3'
 );
CREATE table IF NOT EXISTS
 printSink612 (user_id string, amount int)
WITH
 ('connector' = 'print');
INSERT INTO
 printSink612
SELECT
FROM
dataGenSource612;
```

- Example 2: Connecting DLI to LakeFormation using DEW

#### Create a Flink OpenSource SQL job and set the following parameters:

Parameter	Description	Exampl e Value
Flink Version	Flink 1.15 or later can connect to LakeFormation.	1.15
Agency	Before using Flink 1.15 or later to run jobs, you need to create an agency on the IAM console and add the new agency information. Once set, the system automatically adds the following agency configuration to your job:	
	For agency permission examples, see <b>Creating</b> a Custom DLI Agency and Agency Permission Policies in Common Scenarios.	
Enable Checkpoint ing	Select it.	Select it.

Parameter	Description	Exampl e Value
Runtime Configurati on	<ul> <li>Type of metadata the Flink job accesses. In this scenario, select lakeformation. flink.dli.job.catalog.type=lakeformation</li> <li>Name of the data catalog the Flink job accesses. flink.dli.job.catalog.name=[Catalog name in LakeFormation]</li> <li>The data catalog created on the DLI management console is selected, that is, the mapping between DLI and the data catalog of the default LakeFormation instance. This data catalog is connected to the data catalog of the default LakeFormation instance.</li> </ul>	-

For the catalog parameters in the example, see **Table 4-20**.

Set **properties.catalog.lakeformation.auth.identity.util.class** to **com.huawei.flink.provider.lakeformation.FlinkDewIdentityGenerator** and configure DEW.

Table 4-20 Catalog parameters in the example Flink OpenSource SQL
(using DEW)

Parameter	Description	Mandatory	Example Value
type	Catalog type	Yes	Fixed at hive
hive-conf-dir	hive-conf path, which is fixed at <b>/opt/flink/</b> conf.	Yes	Fixed at <b>/opt/flink/</b> conf
default- database	Default database name	No	If unset, the default database is used.
properties.cat alog.lakecat.a uth.identity.ut il.class	Authentication information acquisition class	Yes	Mandatory when DEW is used, which is fixed at com.huawei.flink.p rovider.lakeformati on.FlinkDewIdentit yGenerator.

Parameter	Description	Mandatory	Example Value
properties.cat alog.dew.proj ectId	ID of the project DEW belongs to. The default value is the ID of the project where the Flink job is.	Yes	Mandatory when DEW is used
properties.cat alog.dew.end point	Endpoint of the DEW service to be used.	Yes	Mandatory when DEW is used. Example: kms.xxx.com
properties.cat alog.dew.csm s.secretName	Name of the shared secret in DEW's secret management.	Yes	Mandatory when DEW is used
properties.cat alog.dew.csm s.version	Version number of the shared secret in DEW's secret management.	Yes	Mandatory when DEW is used
properties.cat alog.dew.acce ss.key	Enter the key corresponding to the <b>access.key</b> value in DEW's secret.	Yes	Mandatory when DEW is used
properties.cat alog.dew.secr et.key	Enter the key corresponding to the <b>secret.key</b> value in DEW's secret.	Yes	Mandatory when DEW is used

CREATE CATALOG myhive WITH ( 'type' = 'hive', 'hive-conf-dir' = '/opt/flink/conf', 'default-database'='default', -- The following is the DEW configuration. Change the parameter values based on site requirements. 'properties.catalog.lakeformation.auth.identity.util.class' = 'com.huawei.flink.provider.lakeformation.FlinkDewIdentityGenerator', 'properties.catalog.dew.endpoint'='kms.xxx.com', 'properties.catalog.dew.csms.secretName'='obsAksK', 'properties.catalog.dew.access.key' = 'myak', 'properties.catalog.dew.secret.key' = 'mysk', 'properties.catalog.dew.projectId'='330e068af1334c9782f4226xxxxxxx', 'properties.catalog.dew.csms.version'='v9' USE CATALOG myhive; create table IF NOT EXISTS dataGenSource\_dew612(

user\_id string,

);

```
amount int
) with (
'connector' = 'datagen',
'rows-per-second' = '1',
'fields.user_id.kind' = 'random',
'fields.user_id.length' = '3'
);
create table IF NOT EXISTS printSink_dew612(
user_id string,
amount int
) with (
'connector' = 'print'
);
```

insert into printSink\_dew612 select \* from dataGenSource\_dew612;

 Example 3: Connecting DLI to LakeFormation using an agency to write data to a Hudi table

Parameter	Description	Exampl e Value
Flink Version	Flink 1.15 or later can connect to LakeFormation.	1.15
Agency	Before using Flink 1.15 or later to run jobs, you need to create an agency on the IAM console and add the new agency information. Once set, the system automatically adds the following agency configuration to your job: flink.dli.job.agency.name= <i>agency</i> For agency permission examples, see <b>Creating</b>	-
	a Custom DLI Agency and Agency Permission Policies in Common Scenarios.	
Enable Checkpoint ing	Select it.	Select it.
Runtime Configurati	• Type of metadata the Flink job accesses. In this scenario, select <b>lakeformation</b> .	-
	<ul> <li>flink.dli.job.catalog.type=lakeformation</li> <li>Name of the data catalog the Flink job accesses.</li> <li>flink.dli.job.catalog.name=[Catalog name in LakeFormation]</li> </ul>	
	The data catalog created on the DLI management console is selected, that is, the mapping between DLI and the data catalog of the default LakeFormation instance. This data catalog is connected to the data catalog of the default LakeFormation instance.	

Create a Flink OpenSource SQL job and set the following parameters:

For the catalog parameters in the example, see **Table 4-21**.

Parameter	Description	Mandatory	Example Value
type	Catalog type	Yes	<b>hudi</b> for a Hudi table
hive-conf-dir	hive-conf path. The value is fixed at <b>/opt/flink/</b> <b>conf</b> .	Yes	Fixed at <b>/opt/flink/</b> conf
default- database	Default database name	No	Default database
mode	The value can be 'hms' or 'non- hms'. • hms indicates that the created Hudi catalog uses Hive Metastore to store metadata.	Yes	Fixed at <b>hms</b>
	<ul> <li>non-hms indicates that Hive Metastore is not used to store metadata.</li> </ul>		

 Table 4-21
 Parameters for configuring a catalog of the Hudi type

 Table 4-22
 Connector parameters for a Hudi sink table

Parameter	Description	Mandatory	Example Value
connector	Flink connector type If set to <b>hudi</b> , the sink table is a Hudi table.	Yes	hudi
path	Basic path of the table. If the path does not exist, the system will create it.	Yes	Refer to the values configured in the sample code.
hoodie.datas ource.write.re cordkey.field	Unique key field name of the Hoodie table	No	Set <b>order_id</b> to a unique key.

Parameter	Description	Mandatory	Example Value
EXTERNAL	Whether the table is foreign	Yes Mandatory for the Hudi table and must be set to <b>true</b> .	true

```
CREATE CATALOG hive_catalog
 WITH (
   'type'='hive',
   'hive-conf-dir' = '/opt/flink/conf',
  'default-database'='test'
 );
USE CATALOG hive_catalog;
create table if not exists genSource618 (
 order id STRING,
 order name STRING,
 price INT,
 weight INT
) with (
 'connector' = 'datagen',
 'rows-per-second' = '1',
 'fields.order_id.kind' = 'random',
 'fields.order_id.length' = '8',
 'fields.order_name.kind' = 'random',
 'fields.order_name.length' = '5'
);
CREATE CATALOG hoodie_catalog
 WITH (
   'type'='hudi',
   'hive.conf.dir' = '/opt/flink/conf',
   'mode'='hms' -- supports 'dfs' mode that uses the DFS backend for table DDLs
persistence
 );
CREATE TABLE if not exists hoodie_catalog.`test`.`hudiSink618` (
 `order_id` STRING PRIMARY KEY NOT ENFORCED,
 `order_name` STRING,
 `price` INT,
 `weight` INT,
 `create_time` BIGINT,
 `create date` String
 PARTITIONED BY (create_date) WITH (
)
 'connector' = 'hudi',
 'path' = 'obs://xxx/catalog/dbtest3/hudiSink618',
 'hoodie.datasource.write.recordkey.field' = 'order_id',
 'write.precombine.field' = 'create_time',
 'EXTERNAL' = 'true' -- must be set
);
insert into hoodie_catalog.`test`.`hudiSink618`
select
 order_id,
 order_name,
 price,
 weight,
 UNIX_TIMESTAMP() as create_time,
 FROM_UNIXTIME(UNIX_TIMESTAMP(), 'yyyyMMdd') as create_date
from genSource618;
```

#### • DLI Flink Jar

- Example 1: Connecting DLI to LakeFormation using an agency
  - i. Develop a Flink Jar program, compile and upload the JAR file to OBS. In this example, the file is uploaded to the **obs://obs-test/dlitest/** directory.

The sample code is as follows:

In this example, random data is generated using the DataGen table and then output to the Print result table.

For other connector types, see List of Connectors Supported by Flink 1.15.

package com.huawei.test;

```
import org.apache.flink.api.java.utils.ParameterTool;
import org.apache.flink.contrib.streaming.state.RocksDBStateBackend;
import org.apache.flink.runtime.state.filesystem.FsStateBackend;
import org.apache.flink.streaming.api.CheckpointingMode;
import org.apache.flink.streaming.api.environment.CheckpointConfig;
import org.apache.flink.streaming.api.environment.StreamExecutionEnvironment;
import org.apache.flink.table.api.EnvironmentSettings;
import org.apache.flink.table.api.bridge.java.StreamTableEnvironment;
import org.slf4j.Logger;
import org.slf4j.LoggerFactory;
import java.text.SimpleDateFormat;
@SuppressWarnings({"deprecation", "rawtypes", "unchecked"})
public class GenToPrintTaskAgency {
  private static final Logger LOGGER =
LoggerFactory.getLogger(GenToPrintTaskAgency.class);
  private static final String datePattern = "yyyy-MM-dd_HH-mm-ss";
  public static void main(String[] args) {
     LOGGER.info("Start task.");
     ParameterTool paraTool = ParameterTool.fromArgs(args);
     String checkpointInterval = "180000000";
     // set up execution environment
     StreamExecutionEnvironment env =
StreamExecutionEnvironment.getExecutionEnvironment();
     EnvironmentSettings settings = EnvironmentSettings.newInstance()
          .inStreamingMode().build();
     StreamTableEnvironment tEnv = StreamTableEnvironment.create(env, settings);
env.getCheckpointConfig().setCheckpointingMode(CheckpointingMode.EXACTLY_ONCE);
     env.getCheckpointConfig().setCheckpointInterval(Long.valueOf(checkpointInterval));
     env.getCheckpointConfig().enableExternalizedCheckpoints(
          CheckpointConfig.ExternalizedCheckpointCleanup.RETAIN_ON_CANCELLATION);
     SimpleDateFormat dateTimeFormat = new SimpleDateFormat(datePattern);
     String time = dateTimeFormat.format(System.currentTimeMillis());
     RocksDBStateBackend rocksDbBackend =
          new RocksDBStateBackend(
               new FsStateBackend("obs://obs/xxx/testcheckpoint/" + time), true);
     env.setStateBackend(rocksDbBackend);
     String createCatalog = "CREATE CATALOG lf_catalog WITH (\n" +
              type' = hive', n'' +
              'hive-conf-dir' = '/opt/hadoop/conf'\n" +
          ");";
     tEnv.executeSql(createCatalog);
     String dataSource = "CREATE TABLE if not exists
lf_catalog.`testdb`.`dataGenSourceJar618_1` (\n" +
          " user id string,\n" +
```

```
" amount int\n" +
             ") WITH (\n" +
" 'connector' = 'datagen',\n" +
             '' 'rows-per-second' = '1',\n" +
'' 'fields.user_id.kind' = 'random',\n" +
             " 'fields.user_id.length' = '3'\n" +
             ")";
/*testdb is a custom database.*/
      tEnv.executeSql(dataSource);
      String printSink = "CREATE TABLE if not exists lf_catalog.`testdb`.`printSinkJar618_1`
(\n" +
              " user_id string,\n" +
              " amount int\n" +
             ") WITH ('connector' = 'print')";
      tEnv.executeSql(printSink);
/*testdb is a custom database.*/
      String query = "insert into lf_catalog.`test`.`printSinkJar618_1` " +
"select * from lf_catalog.`test`.`dataGenSourceJar618_1`";
      tEnv.executeSql(query);
   }
}
```

ii. Create a Flink Jar job and set the following parameters:

Paramete r	Description	Exampl e Value
Flink Version	Flink 1.15 or later can connect to LakeFormation.	1.15
Agency	Before using Flink 1.15 or later to run jobs, you need to create an agency on the IAM console and add the new agency information. Once set, the system automatically adds the following agency configuration to your job:	-
	flink.dli.job.agency.name= <i>agency</i>	
	For agency permission examples, see Creating a Custom DLI Agency and Agency Permission Policies in Common Scenarios.	

Paramete r	Description	Exampl e Value
Runtime Configurat ion	<ul> <li>Type of metadata the Flink job accesses. In this scenario, select lakeformation. flink.dli.job.catalog.type=lakeformation</li> <li>Name of the data catalog the Flink job accesses. flink.dli.job.catalog.name=[Catalog name in LakeFormation]</li> <li>The data catalog created on the DLI management console is selected, that is, the mapping between DLI and the data catalog of the default LakeFormation</li> </ul>	-
	instance. This data catalog is connected to the data catalog of the default LakeFormation instance.	

#### - Example 2: Connecting DLI to LakeFormation using DEW

i. Develop a Flink Jar program, compile and upload the JAR file to OBS. In this example, the file is uploaded to the **obs://obs-test/dlitest/** directory.

The sample code is as follows:

In this example, random data is generated using the DataGen table and then output to the Print result table.

For other connector types, see **List of Connectors Supported by Flink 1.15**.

package com.huawei.test;

import org.apache.flink.api.java.utils.ParameterTool; import org.apache.flink.contrib.streaming.state.RocksDBStateBackend; import org.apache.flink.runtime.state.filesystem.FsStateBackend; import org.apache.flink.streaming.api.CheckpointingMode; import org.apache.flink.streaming.api.environment.CheckpointConfig; import org.apache.flink.streaming.api.environment.StreamExecutionEnvironment; import org.apache.flink.table.api.EnvironmentSettings; import org.apache.flink.table.api.bridge.java.StreamTableEnvironment; import org.slf4j.Logger; import org.slf4j.LoggerFactory; import java.text.SimpleDateFormat; @SuppressWarnings({"deprecation", "rawtypes", "unchecked"}) public class GenToPrintTaskDew { private static final Logger LOGGER = LoggerFactory.getLogger(GenToPrintTaskAgency.class); private static final String datePattern = "yyyy-MM-dd\_HH-mm-ss"; public static void main(String[] args) { LOGGER.info("Start task."); ParameterTool paraTool = ParameterTool.fromArgs(args); String checkpointInterval = "180000000"; // set up execution environment StreamExecutionEnvironment env =



ii. Create a Flink Jar job and set the following parameters:

Paramete r	Description	Exampl e Value
Flink Version	Flink 1.15 or later can connect to LakeFormation.	1.15

Paramete r	Description	Exampl e Value
Agency	Before using Flink 1.15 or later to run jobs, you need to create an agency on the IAM console and add the new agency information. Once set, the system automatically adds the following agency configuration to your job:	
	For agency permission examples, see Creating a Custom DLI Agency and Agency Permission Policies in Common Scenarios.	
Runtime Configurat ion	<ul> <li>Type of metadata the Flink job accesses. In this scenario, select lakeformation. flink.dli.job.catalog.type=lakeformation</li> <li>Name of the data catalog the Flink job accesses. flink.dli.job.catalog.name=[Catalog name in LakeFormation] The data catalog created on the DLI management console is selected, that is,</li> </ul>	-
	the mapping between DLI and the data catalog of the default LakeFormation instance. This data catalog is connected to the data catalog of the default LakeFormation instance.	

#### - Example 3: Flink Jar jobs supporting Hudi tables

i. Develop a Flink Jar program, compile and upload the JAR file to OBS. In this example, the file is uploaded to the **obs://obs-test/dlitest/** directory.

The sample code is as follows:

In this example, random data is generated using the DataGen table and then output to the Hudi result table.

For other connector types, see **List of Connectors Supported by** Flink 1.15.

package com.huawei.test;

import org.apache.flink.api.java.utils.ParameterTool; import org.apache.flink.contrib.streaming.state.RocksDBStateBackend; import org.apache.flink.runtime.state.filesystem.FsStateBackend; import org.apache.flink.streaming.api.CheckpointingMode; import org.apache.flink.streaming.api.environment.CheckpointConfig; import org.apache.flink.streaming.api.environment.StreamExecutionEnvironment; import org.apache.flink.table.api.EnvironmentSettings; import org.apache.flink.table.api.bridge.java.StreamTableEnvironment; import org.slf4j.Logger;

import org.slf4j.LoggerFactory;

import java.io.IOException; import java.text.SimpleDateFormat; public class GenToHudiTask4 { private static final Logger LOGGER = LoggerFactory.getLogger(GenToHudiTask4.class); private static final String datePattern = "yyyy-MM-dd\_HH-mm-ss"; public static void main(String[] args) throws IOException { LOGGER.info("Start task."); ParameterTool paraTool = ParameterTool.fromArgs(args); String checkpointInterval = "30000"; // set up execution environment StreamExecutionEnvironment env = StreamExecutionEnvironment.getExecutionEnvironment(); EnvironmentSettings settings = EnvironmentSettings.newInstance() .inStreamingMode().build(); StreamTableEnvironment tEnv = StreamTableEnvironment.create(env, settings); env.getCheckpointConfig().setCheckpointingMode(CheckpointingMode.EXACTLY\_ONCE); env.getCheckpointConfig().setCheckpointInterval(Long.valueOf(checkpointInterval)); env.getCheckpointConfig().enableExternalizedCheckpoints( CheckpointConfig.ExternalizedCheckpointCleanup.RETAIN\_ON\_CANCELLATION); SimpleDateFormat dateTimeFormat = new SimpleDateFormat(datePattern); String time = dateTimeFormat.format(System.currentTimeMillis()); RocksDBStateBackend rocksDbBackend = new RocksDBStateBackend( new FsStateBackend("obs://xxx/jobs/testcheckpoint/" + time), true); env.setStateBackend(rocksDbBackend); String catalog = "CREATE CATALOG hoodie\_catalog\n" + " WITH (\n" + " 'type'='hudi',\n" + ... 'hive.conf.dir' = '/opt/hadoop/conf',\n" + " 'mode'='hms'\n" + ")"; tEnv.executeSql(catalog); String dwsSource = "CREATE TABLE if not exists genSourceJarForHudi618\_1 (\n" + order\_id STRING,\n" + " order\_name STRING,\n" + " price INT,\n" + " weight INT\n" + ") WITH (\n" + " 'connector' = 'datagen',\n" + " 'rows-per-second' = '1',\n" + " 'fields.order\_id.kind' = 'random',\n" + " 'fields.order\_id.length' = '8',n" + " 'fields.order\_name.kind' = 'random',\n" + " 'fields.order\_name.length' = '8'\n" + ")"; tEnv.executeSql(dwsSource); /\*testdb is a custom database.\*/ String printSinkdws = "CREATE TABLE if not exists hoodie\_catalog. `testdb`.`hudiSinkJarHudi618\_1` (\n" + " order\_id STRING PRIMARY KEY NOT ENFORCED,\n" + " order\_name STRING,\n" + " price INT,\n" + " weight INT,\n" + " create\_time BIGINT,\n" + " create\_date String\n" + ") WITH (" + "'connector' = 'hudi',\n" + "'path' = 'obs://xxx/catalog/dbtest3/hudiSinkJarHudi618\_1',\n" + "'hoodie.datasource.write.recordkey.field' = 'order\_id',\n" + "'EXTERNAL' = 'true'\n" + ")";

tEnv.executeSql(printSinkdws); /\*testdb is a custom database.\*/ String query = "insert into hoodie\_catalog.`testdb`.`hudiSinkJarHudi618\_1` select\n" + " order\_id,\n" + " price,\n" + " weight,\n" + " UNIX\_TIMESTAMP() as create\_time,\n" + " FROM\_UNIXTIME(UNIX\_TIMESTAMP(), 'yyyyMMdd') as create\_date\n" + " from genSourceJarForHudi618\_1"; tEnv.executeSql(query); } }

Parameter	Description	Mandator y	Example Value
connector	Flink connector type	Yes	hudi
	If set to <b>hudi</b> , the sink table is a Hudi table.		
path	Basic path of the table. If the path does not exist, the system will create it.	Yes	Refer to the values configured in the sample code.
hoodie.datas ource.write.r ecordkey.fiel d	Unique key field name of the Hoodie table	No	Set <b>order_id</b> to a unique key.
EXTERNAL	Whether the table is foreign	Yes Mandator y for the Hudi table and must be set to <b>true</b> .	true

Table 4-23 Connector	parameters	for a	Hudi	sink	table
----------------------	------------	-------	------	------	-------

#### ii. Create a Flink Jar job and set the following parameters:

Paramete r	Description	Exampl e Value
Flink Version	Flink 1.15 or later can connect to LakeFormation.	1.15

Paramete r	Description	Exampl e Value
Agency	Before using Flink 1.15 or later to run jobs, you need to create an agency on the IAM console and add the new agency information. Once set, the system automatically adds the following agency configuration to your job:	
	flink.dli.job.agency.name= <i>agency</i>	
	For agency permission examples, see Creating a Custom DLI Agency and Agency Permission Policies in Common Scenarios.	
Runtime Configurat	<ul> <li>Type of metadata the Flink job accesses. In this scenario, select lakeformation.</li> </ul>	-
ion	flink.dli.job.catalog.type=lakeformation	
	<ul> <li>Name of the data catalog the Flink job accesses. flink.dli.job.catalog.name=[Catalog name in LakeFormation]</li> </ul>	
	The data catalog created on the DLI management console is selected, that is, the mapping between DLI and the data catalog of the default LakeFormation instance. This data catalog is connected to the data catalog of the default LakeFormation instance.	

# 4.7.2 Permission Policies and Supported Actions for LakeFormation Resources

#### Supported Actions for LakeFormation SQL Resources

For the list of actions supported by DLI for SQL resource authentication, refer to **Data Permissions List**.

 Table 4-24 lists the supported actions for LakeFormation SQL resources.

Resource Type	Permission Type
Database	ALL
	ALTER
	DROP

 Table 4-24 Supported actions for LakeFormation SQL resources

Resource Type	Permission Type
	DESCRIBE
	LIST_TABLE
	LIST_FUNC
	CREATE_TABLE
	CREATE_FUNC
Table/View	ALL
	ALTER
	DROP
	DESCRIBE
	UPDATE
	INSERT
	SELECT
	DELETE
Column	SELECT
Function	ALL
	ALTER
	DROP
	DESCRIBE
	EXEC

# LakeFormation Permission Policies (Spark)

Table 4-25 LakeFormation	permission	policies
--------------------------	------------	----------

Туре	SQL Statement	Permission to Authenticate Access to Metadata Using IAM	Permission to Authenticate Access to SQL Resources
DDL	ALTER DATABASE	database:describe	database:DESCRIBE
statement		database:alter	database:ALTER

Туре	SQL Statement	Permission to Authenticate Access to Metadata Using IAM	Permission to Authenticate Access to SQL Resources
	ALTER TABLE	database:describe table:describe table:alter database:create	database:DESCRIBE table:DESCRIBE table:ALTER database:CREATE_TAB LE column:SELECT or table:SELECT
	ALTER VIEW	database:describe table:describe table:alter	database:DESCRIBE table:DESCRIBE column:SELECT table:ALTER
	CREATE DATABASE	database:describe database:create	database:DESCRIBE catalog:CREATE_DATA BASE
	CREATE OR REPLACE FUNCTION (CREATE)	database:describe function:create	database:DESCRIBE database:CREATE_FU NC
	CREATE OR REPLACE FUNCTION (REPLACE)	database:describe function:describe function:alter	database:CREATE_FU NC database:DESCRIBE function:DESCRIBE function:ALTER
	CREATE TABLE	database:describe table:describe table:create	database:DESCRIBE database:CREATE_TAB LE
	CREATE VIEW	database:describe table:describe table:drop table:create	database:CREATE_TAB LE table:DESCRIBE (source\target) table:DROP(target) column:SELECT
	DROP DATABASE	database:describe database:drop	database:DESCRIBE database:DROP

Туре	SQL Statement	Permission to Authenticate Access to Metadata Using IAM	Permission to Authenticate Access to SQL Resources
	DROP FUNCTION	database:describe function:describe function:drop	database:DESCRIBE function:DESCRIBE function:DROP
	DROP TABLE	database:describe table:describe credential:describe table:drop	database:DESCRIBE table:DESCRIBE table:DROP
	DROP VIEW	database:describe table:describe table:drop	database:DESCRIBE table:DESCRIBE(target \source) table:DROP(target)
	REPAIR TABLE	database:describe table:describe credential:describe table:alter	database:DESCRIBE table:DESCRIBE table:ALTER table:SELECT
	TRUNCATE TABLE	database:describe table:describe table:alter	database:DESCRIBE table:DESCRIBE table:SELECT table:UPDATE
DML statement	INSERT TABLE	database:describe table:describe table:alter credential:describe	database:DESCRIBE table:DESCRIBE table:ALTER table:INSERT column:SELECT or table:SELECT
	LOAD DATA	database:describe table:describe credential:describe	database:DESCRIBE table:DESCRIBE table:UPDATE table:ALTER table:SELECT
DR statement	SELECT	database:describe table:describe credential:describe	database:DESCRIBE table:DESCRIBE column:SELECT
	EXPLAIN	Depends on the SQL statement.	Depends on the SQL statement.

Туре	SQL Statement	Permission to Authenticate Access to Metadata Using IAM	Permission to Authenticate Access to SQL Resources
Auxiliary statement	ANALYZE TABLE	database:describe table:describe credential:describe table:alter	database:DESCRIBE table:DESCRIBE table:SELECT table:ALTER
	DESCRIBE DATABASE	database:describe	database:DESCRIBE
	DESCRIBE FUNCTION	database:describe function:describe	database:DESCRIBE function:DESCRIBE
	DESCRIBE QUERY	database:describe table:describe	database:DESCRIBE table:DESCRIBE table:SELECT
	DESCRIBE TABLE	database:describe table:describe	database:DESCRIBE table:DESCRIBE
	REFRESH TABLE	database:describe table:describe credential:describe	database:DESCRIBE table:DESCRIBE table:SELECT
	REFRESH FUNCTION	database:describe function:describe	database:DESCRIBE function:DESCRIBE
	SHOW COLUMNS	database:describe table:describe	database:DESCRIBE table:DESCRIBE
	SHOW CREATE TABLE	database:describe table:describe	database:DESCRIBE table:DESCRIBE
	SHOW DATABASES	database:describe	catalog:LIST_DATABA SE database:DESCRIBE
	SHOW FUNCTIONS	database:describe function:describe	database:DESCRIBE
	SHOW PARTITIONS	database:describe table:describe	database:DESCRIBE table:DESCRIBE

Туре	SQL Statement	Permission to Authenticate Access to Metadata Using IAM	Permission to Authenticate Access to SQL Resources
	SHOW TABLE EXTENDED	database:describe table:describe	catalog:LIST_DATABA SE database:DESCRIBE table:DESCRIBE database:LIST_TABLE
	SHOW TABLES	database:describe table:describe	catalog:LIST_DATABA SE database:LIST_TABLE database:DESCRIBE
	SHOW TBLPROPERTIES	database:describe table:describe	database:DESCRIBE table:DESCRIBE
	SHOW VIEWS	database:describe table:describe	catalog:LIST_DATABA SE database:LIST_TABLE database:DESCRIBE

### LakeFormation Permission Policies (HetuEngine)

**Table 4-26** Reference for configuration LakeFormation permissions usingHetuEngine syntax

Туре	Syntax	LakeFormation Permission Required for SQL Authentication	LakeFormation Permission Required for Metadata API Calling
Sche ma	create schema	catalog:CREATE_DATABA SE	catalog:CREATE_DATABASE catalog:DESCRIBE
	show schemas	catalog:LIST_DATABASE	catalog:LIST_DATABASE
	drop schema	database:DROP	catalog:LIST_DATABASE database:DESCRIBE database:DROP
	alter schema set location/ owner	database:ALTER	catalog:LIST_DATABASE database:DESCRIBE database:ALTER

Туре	Syntax	LakeFormation Permission Required for SQL Authentication	LakeFormation Permission Required for Metadata API Calling
	desc schema	database:LIST_DATABASE	database:LIST_DATABASE database:DESCRIBE
Table	create table	database:CREATE_TABLE	database:DESCRIBE database:CREATE_TABLE
	create table as select	database:CREATE_TABLE Source table: SELECT (or column:SELECT)	database:DESCRIBE database:CREATE_TABLE table:DESCRIBE (source table) table:select (source table)
	show create table	table:DESCRIBE	table:DESCRIBE table:select
	select from table	table:SELECT (or column:SELECT)	table:DESCRIBE table:SELECT (or column:SELECT)
	insert into table	table:INSERT table:SELECT (or column:SELECT)	table:DESCRIBE table:ALTER
	alter table	table:ALTER	table:DESCRIBE table:ALTER
	show tables	database:LIST_TABLE	catalog:LIST_DATABASE database:LIST_TABLE
	drop table	table:DROP	table:DESCRIBE table:DROP
	truncate table	table:DELETE	table:DESCRIBE
	desc table	table:DESCRIBE	catalog:LIST_DATABASE table:DESCRIBE
	comment	table:ALTER	table:DESCRIBE table:ALTER
view	create view	database:CREATE_TABLE Source table: SELECT (or column:SELECT)	database:CREATE_TABLE table:DESCRIBE (source table) table:select (source table)

Туре	Syntax	LakeFormation Permission Required for SQL Authentication	LakeFormation Permission Required for Metadata API Calling
	drop view	table:DROP	table:DESCRIBE table:DROP
	alter view	table:ALTER	table:DESCRIBE table:ALTER (table:SELECT)
	select from view	table:DESCRIBE (source table and view) table:select (source table	table:DESCRIBE (source table and view) table:select (source table
		and view)	and view)
	show views	database:LIST_TABLE	catalog:LIST_DATABASE database:LIST_TABLE table:DESCRIBE
	show create view	table:DESCRIBE	table:DESCRIBE
colum n	show columns	table:SELECT (or column:SELECT)	catalog:LIST_DATABASE table:DESCRIBE table:SELECT (or column:SELECT)
	select [column] from table	table:SELECT (or column:SELECT)	table:DESCRIBE table:SELECT (or column:SELECT)
stats	show stats	table:SELECT (or column:SELECT)	table:DESCRIBE table:SELECT (or column:SELECT)
	analyze	table:INSERT table:SELECT (or column:SELECT)	table:DESCRIBE table:ALTER

# **5** Data Import and Migration

# 5.1 Overview

#### Importing Data to an OBS Table

DLI enables direct access to data stored in OBS for query and analysis, eliminating the need for data migration.

To begin using DLI for data analysis, just import your local data into OBS.

#### Migrating Data to DLI

To centrally analyze and manage scattered data from various systems, you can utilize migration tools like Cloud Data Migration (CDM) to migrate the data to DLI. Once the migration is complete, you can submit DLI jobs to analyze the data.

CDM supports multiple types of data sources, such as databases, data warehouses, and files. You can configure data source migration tasks on the GUI to enhance the efficiency of data migration and integration.

For details, see Migrating Data from External Data Sources to DLI.



Figure 5-1 Migrating data to DLI
### Configuring DLI to Read and Write External Data Sources

If you prefer not to import data into OBS or DLI tables, DLI offers cross-source access capabilities, allowing you to connect to data sources for analysis without the need to migrate the data.

For details, see **Configuring DLI to Read and Write External Data Sources**.

### 5.2 Migrating Data from External Data Sources to DLI

### 5.2.1 Overview of Data Migration Scenarios

To centrally analyze and manage scattered data from various systems, you can utilize migration tools like Cloud Data Migration (CDM) to migrate the data to DLI. Once the migration is complete, you can submit DLI jobs to analyze the data.

CDM supports multiple types of data sources, such as databases, data warehouses, and files. You can configure data source migration tasks on the GUI to enhance the efficiency of data migration and integration.

Figure 5-2 Migrating data to DLI



### **Common Migration Scenarios and Solutions**

Data Type	Migration Tool	Solution
Hive	CDM	Example Typical Scenario: Migrating Data from Hive to DLI
Kafka	CDM	Example Typical Scenario: Migrating Data from Kafka to DLI
Elasticsearch	CDM	Example Typical Scenario: Migrating Data from Elasticsearch to DLI

Table 5-1 Common migration scenarios and solutions

Data Type	Migration Tool	Solution
RDS	CDM	Example Typical Scenario: Migrating Data from RDS to DLI
GaussDB(DWS)	CDM	Example Typical Scenario: Migrating Data from GaussDB(DWS) to DLI

### Data Type Mapping

Refer to **Table 5-2** for data type mapping between data sources and destinations during data migration between DLI and other cloud services and platforms. This will aid in data type conversion and mapping.

**Table 5-2** Data type mapping

MySQL	Hive	GaussDB( DWS)	DB( Oracle Postgre SQL		Hologre s	DLI Spark
CHAR	CHAR	CHAR	CHAR	CHAR	CHAR	CHAR
VARCH AR	VARCHAR	VARCHAR	VARCHAR	VARCHA R	VARCHA R	VARCHA R/ STRING
DECIMA L	DECIMA DECIMAL		NUMERIC	NUMERI C	DECIMA L	DECIMAL
INT	INT	INTEGER	NUMBER	INTEGER	INTEGER	INT
BIGINT BIGINT		BIGINT	NUMBER	BIGINT	BIGINT	BIGINT/ LONG
TINYINT	TINYINT	SMALLINT	NUMBER	SMALLI NT	SMALLI NT	TINYINT
SMALLI NT	SMALLINT	SMALLINT	NUMBER	SMALLI NT	SMALLI NT	SMALLIN T/SHORT
BINARY	BINARY BINARY		RAW	BYTEA	BYTEA	BINARY
VARBIN ARY	BINARY	BYTEA	RAW	BYTEA	BYTEA	BINARY
FLOAT	FLOAT	FLOAT4	FLOAT	DOUBLE	FLOAT4	FLOAT
DOUBL E	DOUBLE	FLOAT8	FLOAT	REAL/ DOUBLE	FLOAT8	DOUBLE
DATE	DATE	TIMESTAM P	DATE	DATE	DATE	DATE

MySQL	Hive	GaussDB( DWS)	Oracle	Postgre SQL	Hologre s	DLI Spark
TIME	TIME Not supported (use String instead)		DATE	TIME	TIME	Not supporte d (use String instead)
DATETI ME	DATETI TIMESTA ME MP		TIME	TIME	TIMESTA MP	TIMESTA MP
TINYINT TINYINT		BOOLEAN	Not supporte d	TINYINT	BOOLEA N	BOOLEA N
Not support ed (use TEXT instead)	Not supported (use String instead)	Not supported (use TEXT instead)	Not supporte d (use VARCHAR instead)	Not supporte d (use TEXT instead)	Not supporte d (use TEXT instead)	ARRAY
Not Not support supported ed (use (use TEXT String instead) instead)		Not supported (use TEXT instead)	Not supporte d (use VARCHAR instead)	Not supporte d (use TEXT instead)	Not supporte d (use TEXT instead)	МАР
Not support ed (use TEXT instead)	Not supported (use String instead)	Not supported (use TEXT instead)	Not supporte d (use VARCHAR instead)	Not supporte d (use TEXT instead)	Not supporte d (use TEXT instead)	STRUCT

### 

If a service does not support a standard data type, you can use the recommended data type.

### 5.2.2 Using CDM to Migrate Data to DLI

On its GUI, CDM enables you to create data migration tasks from multiple data sources to a data lake.

This section describes how to use CDM to migrate data from data sources to DLI.

Figure 5-3 Process of migrating data to DLI using CDM

			CDM	ma	anagement console					DLI management console
a	Create CDM cluster.	¥	Create a connection between a data source and CDM.	•	Create a data connection between CDM and DLI.	•	Create a data migration job.	÷	Monitor the job execution status.	Review the data migration results.

### Step 1: Create a CDM Cluster

A CDM cluster is used to execute data migration jobs that migrate data from data sources to DLI.

- 1. Log in to the CDM management console.
- 2. Click **Buy CDM Cluster** in the upper right corner. On the displayed page, set CDM cluster parameters.
  - You are advised to set the region, VPC, subnet, security group, and enterprise project to be the same as those of the data source and DLI.
  - Once a cluster is created, its specifications cannot be modified. If you need to use higher specifications, you will need to create a new cluster.

For how to set CDM cluster parameters, see Creating a Cluster.

- 3. Click **Buy Now**.
- 4. Confirm the settings and click **Submit**. The system starts to create a CDM cluster. You can check the creation progress on the **Cluster Management** page.

### Step 2: Create a Data Connection Between the Data Source and CDM

This step uses the MySQL data source as an example to describe how to create a data connection between the data source and CDM.

- Step 1 In the navigation pane on the left of the CDM console, choose Cluster Management. Locate the cdm-aff1 cluster created in Step 1: Create a CDM Cluster.
- **Step 2** Click **Job Management** in the **Operation** column.
- **Step 3** Click the **Link Management** tab then **Create Link**.

Data Warehouse	Data Warehouse Service	Data Lake Insight	MRS ClickHouse	
Hadoop	MRS HDFS	Apache HDFS	MRS HBase	Apache HBas
	MRS Hive	Apache Hive	MRS Hudi	
Object Storage	Object Storage Service (OBS)			
File System	FTP	SFTP	HTTP	
Relational Database	RDS for MySQL	MySQL	RDS for PostgreSQL	PostgreSQL
	RDS for SQL Server	Microsoft SQL Server	Oracle	
NoSQL	Redis	MongoDB		
Messaging System	Data Ingestion Service	MRS Kafka	Apache Kafka	
Search	Elasticsearch			
Open Beta Test	^			
X Cancel > Next				

Figure 5-4 Selecting a connector

### **Step 4** Select **RDS for MySQL** and click **Next** to set the link parameters.

Click **Show Advanced Attributes** to view more optional parameters. For details, see **Link to an RDS for MySQL/MySQL Database**. Retain the default values for the optional parameters and set the mandatory parameters based on **Table 5-3**.

Parameter	Description	Example Value		
Name	Name Enter a unique link name.			
Database Server	IP address or domain name of the MySQL database	-		
Port	Port number of the MySQL database	3306		
Database Name	Name of the MySQL database	sqoop		
Username	User who has the read, write, and delete permissions on the MySQL database	admin		
Password	Password of the user	-		
Use Local API	Whether to use the local API of the database for acceleration. (The system attempts to enable the <b>local_infile</b> system variable of the MySQL database.)	Yes		
Use Agent	This parameter does not need to be set as the agent function will be unavailable soon.	-		
local_infile Character Set	When using local_infile to import data to MySQL, you can set the encoding format.	utf8		
Driver Version	Before connecting CDM to a relational database, you need to upload the JDK 8 .jar driver of the relational database. Download the MySQL driver 5.1.48 from https:// downloads.mysql.com/archives/c- j/, obtain mysql-connector- java-5.1.48.jar, and upload it.	-		

 Table 5-3 MySQL link parameters

**Step 5** Click **Test** to check whether the parameters are correctly configured. If the test is successful, click **Save** to create the link and return to the **Links** page.

----End

### Step 3: Create a Data Connection Between CDM and DLI

- Step 1 In the navigation pane on the left of the CDM console, choose Cluster Management. Locate the cdm-aff1 cluster created in Step 1: Create a CDM Cluster.
- **Step 2** Click **Job Management** in the **Operation** column.
- **Step 3** Click the **Link Management** tab then **Create Link**.

### Figure 5-5 Selecting a connector

Data Warehouse Service	Data Lake Insight	MRS ClickHouse				
MRS HDFS	Apache HDFS	MRS HBase	Apache HBase			
MRS Hive	Apache Hive	MRS Hudi				
Object Storage Service (OBS)						
FTP	SFTP	HTTP				
RDS for MySQL	MySQL	RDS for PostgreSQL	PostgreSQL			
RDS for SQL Server	Microsoft SQL Server	Oracle				
Redis	MongoDB					
Data Ingestion Service	MRS Kafka	Apache Kafka				
Elasticsearch						
Open Beta Test						
	Data Warehouse Service MRS HDFS MRS Hive Object Storage Service (OBS) FTP RDS for MySQL RDS for SQL Server Redis Data Ingestion Service Elasticsearch	Data Warehouse Service     Data Lake Insight       MRS HDFS     Apache HDFS       MRS Hive     Apache Hive       Object Storage Service (OBS)        FTP     SFTP       RDS for MySQL     MySQL       RDS for SQL Server     Microsoft SQL Server       Redis     MongoDB       Data Ingestion Service     MRS Kafka	Data Warehouse Service     Data Lake Insight     MRS ClickHouse       MRS HDFS     Apache HDFS     MRS HBase       MRS Hive     Apache Hive     MRS Hudi       Object Storage Service (OBS)         FTP     SFTP     HTTP       RDS for MySQL     MySQL     RDS for PostgreSQL       RDS for SQL Server     Microsoft SQL Server     Oracle       Redis     MRS Kafka     Apache Kafka       Elasticsearch			

**Step 4** Select **Data Lake Insight** for **Data Warehouse** and click **Next**. On the displayed page, set DLI link parameters.

Click **Show Advanced Attributes** to view optional parameters. For details, see **Link to DLI**. Retain the default values for the optional parameters and set the mandatory parameters based on **Table 5-4**.

Fable 5-4 DL	I link	parameters
--------------	--------	------------

Parameter	Description	Example Value
Name	Link name, which should be defined based on the data source type, so it is easier to remember what the link is for	dlilink

Parameter	Description	Example Value
Access Key	AK/SK required for authentication during	-
Secret Key	You need to create an access key for the current account and obtain an AK/SK pair.	-
	<ol> <li>Log in to the management console, move the cursor to the username in the upper right corner, and select My Credentials from the drop-down list.</li> </ol>	
	2. On the <b>My Credentials</b> page, choose Access Keys, and click Create Access Key. See Figure 5-6.	
	Figure 5-6 Clicking Create Access Key	
	3. Click <b>OK</b> and save the access key file as prompted. The access key file will be saved to your browser's configured download location. Open the <b>credentials.csv</b> file to view <b>Access Key Id</b> and <b>Secret Access Key</b> .	
	<ul> <li>Only two access keys can be added for each user.</li> <li>To ensure access key security, the access key is automatically downloaded only when it is generated for the first time and cannot be obtained from the management console</li> </ul>	
	later. Keep them properly.	
Project ID	Project ID of the region where DLI is A project is a group of tenant resources, and an account ID corresponds to the current account. The IAM ID corresponds to the current user. You can view the project IDs, account IDs, and user IDs in different regions	-
	<ul><li>on the corresponding pages.</li><li>1. Register with and log in to the management console.</li></ul>	
	2. Hover the cursor on the username in the upper right corner and select <b>My Credentials</b> from the drop-down list.	
	3. On the <b>API Credentials</b> page, obtain the account name, account ID, IAM username, and IAM user ID, and obtain the project and its ID from the project list.	

**Step 5** Click **Test** to check whether the parameters are correctly configured. If the test is successful, click **Save** to create the link and return to the **Links** page.

----End

### Step 4: Create a Data Migration Job on CDM

After establishing the data connection between the data source and CDM, as well as between CDM and DLI, you need to create a data migration job to migrate the data from the data source to DLI.

- Step 1 On the Cluster Management page, locate the cdm-aff1 cluster created in Step 1: Create a CDM Cluster.
- **Step 2** Click **Job Management** in the **Operation** column.
- **Step 3** Click the **Table/File Migration** tab, click **Create Job**, and configure basic job information.
  - Job Name: Enter a unique job name, for example, mysql2dli.
  - Source job parameters
    - **Source Link Name**: Select the MySQL link **mysqllink**.
    - Use SQL: Select No.
    - Schema/Tablespace: Select the MySQL database from which the table is to be exported.
    - **Table Name**: Select the table from which data is to be exported.
    - Retain the default values for other optional parameters. For details, see From MySQL.
  - Destination job parameters
    - **Destination Link Name**: Select the DLI link **dlilink**.
    - **Schema/Tablespace**: Select the schema to which data is to be imported.
    - Auto Table Creation: Select Auto creation. If the table specified by Table Name does not exist, CDM automatically creates the table in DLI.
    - **Table Name**: Select the table to which data is to be imported.
    - Advanced Attributes > Extend Field Length: Select Yes. The encoding techniques used for storing Chinese characters in MySQL and DLI differ, and the required lengths also vary. In UTF-8 encoding, a Chinese character can take up to three bytes. If this parameter is set to Yes, DLI will automatically create tables with character fields that are three times the length of the original table to avoid errors caused by insufficient character field length in DLI tables.
    - Retain the default values for other optional parameters. For details, see To DWS.
- **Step 4** Click **Next**. The **Map Field** tab page is displayed. CDM automatically maps the fields in the source and destination data tables. You need to check if the field mapping is correct.
  - If the mapping is incorrect, click the row where the field is located and hold down the left mouse button to drag the field to adjust the mapping.

- When importing data into DLI, you must manually choose the distribution columns. You are advised to select the distribution columns based on the following principles:
  - a. Use the primary key as the distribution column.
  - b. If multiple data segments are combined as primary keys, specify all primary keys as the distribution column.
  - c. In the scenario where no primary key is available, if no distribution column is selected, DWS uses the first column as the distribution column by default. As a result, data skew risks exist.
- To convert the content of the source fields, perform the operations in this step. For details, see **Converting Fields**. In this example, field conversion is not required.
- **Step 5** Click **Next** and set task parameters. Typically, retain the default values for all parameters.

In this step, you can configure the following optional features:

- **Retry Upon Failure**: If the job fails to be executed, you can determine whether to automatically retry. Retain the default value **Never**.
- **Group**: Select the group to which the job belongs. The default group is **DEFAULT**. On the **Job Management** page, jobs can be displayed, started, or exported by group.
- **Schedule Execution**: Determine whether to automatically execute the job at a scheduled time. Retain the default value **No**.
- **Concurrent Extractors**: Enter the number of concurrent extractors. An appropriate value improves migration efficiency. For details, see **Performance Tuning**. Retain the default value **1**.
- Write Dirty Data: Specify this parameter if data that fails to be processed or filtered out during job execution needs to be written to OBS for future viewing. Before writing dirty data, create an OBS link on the CDM console. Retain the default value **No**, meaning dirty data is not recorded.
- **Step 6** Click **Save and Run**. CDM starts to execute the job immediately.

----End

### **Step 5: View Data Migration Results**

This step describes how to view a job's execution results and its historical information within the past 90 days, including the number of written rows, read rows, written bytes, written files, and log information.

- Viewing the status of the migration job on CDM
  - a. On the **Cluster Management** page, locate the **cdm-aff1** cluster created in **Step 1: Create a CDM Cluster**.
  - b. Click Job Management in the Operation column.
  - c. Locate the mysql2dli job created in Step 4: Create a Data Migration Job on CDM and check its execution status. If the job status is Succeeded, the migration is successful.
- Viewing data migration results on DLI

- a. After the CDM migration job is complete, log in to the DLI management console.
- b. In the navigation pane on the left, choose **SQL Editor**.

Set **Engine** to **Spark**, **Queues** to the created SQL queue, and **Databases** to the created database. Run the following DLI table query statement to check whether the MySQL data has been successfully migrated to the DLI table: select \* from *tablename*;

# 5.2.3 Example Typical Scenario: Migrating Data from Hive to DLI

This section describes how to use CDM's data synchronization to migrate data from MRS Hive to DLI. All other MRS Hadoop component data can be synchronized bidirectionally with DLI using CDM.

### Prerequisites

• You have created a DLI SQL queue.

### 

Set Type to For SQL when buying a queue.

- You have created an MRS security cluster that contains the Hive component.
  - In this example, the MRS cluster and component versions are as follows:
    - Cluster version: MRS 3.1.0
    - Hive version: 3.1.0
    - Hadoop version: 3.1.1
  - In this example, Kerberos authentication is enabled when the MRS cluster is created.
- You have created a CDM cluster. For how to create a CDM cluster, see Creating a CDM Cluster.

### D NOTE

- To connect the cluster to an on-premises database as the destination data source, you can use either Internet or Direct Connect. If the Internet is used, make sure that an EIP has been bound to the CDM cluster, the security group of CDM allows outbound traffic from the host where the on-premises data source is located, the host where the data source is located can access the Internet, and the connection port has been enabled in the firewall rules.
- To successfully connect to cloud services like MRS and GaussDB (DWS) as data sources, the following requirements must be met:

i. If the CDM cluster and the cloud service are in different regions, they must be connected through either the Internet or Direct Connect. If the Internet is used, make sure that an EIP has been bound to the CDM cluster, the host where the data source is located can access the Internet, and the port has been enabled in the firewall rules.

ii. If the CDM cluster and the cloud service are in the same region, instances in the same VPC, subnet, and security group can communicate with each other by default. If the CDM cluster and the cloud service are in the same VPC but in different subnets or security groups, you must configure routing and security group rules.

For how to configure routing rules, see **Configure routes**. For how to configure security group rules, see **Security Group Configuration Examples**.

iii. The cloud service instance and the CDM cluster belong to the same enterprise project. If they do not, change the enterprise project of the workspace.

In this example, the VPC, subnet, and security group of the CDM cluster match those of the MRS cluster.

### Step 1: Prepare Data

- Create a Hive table in the MRS cluster and insert data into the table.
  - a. Log in to MRS Manager by referring to Accessing FusionInsight Manager.
  - b. Click **System** and choose **Permission** > **Role**. On the displayed page, set the following parameters:
    - Role Name: Enter a role name, for example, hivetestrole.
    - Configure Resource Permission: Select the MRS cluster name and then Hive. Select Hive Admin Privilege.

### Figure 5-7 Creating a Hive role

	n FusionInsight Manager	Homepage Cluster - Hosts	O&M Audit Tenant Resources System				
	FOI	Role > Create Role					
	System	Role Name: Configure Resource Permission:	hivetestrole All resources > mrs_test_00378328 > Hive View Name				
	Permission ^	compare recourse remained.					
ļ	User Group	[	Hive Read Write Privileges				
1	Role     Security Policy	Description:					
	Domain and Mutual Trust Interconnection						
	Certificate		OK Cancel				
	OMS Component						

For how to create a role, see Creating a Role.

- c. Click **System** and choose **Permission** > **User**. On the displayed page, set the following parameters:
  - i. **Username**: Enter a username. In this example, enter **hivetestusr**.
  - ii. User Type: Select Human-Machine.
  - iii. **Password** and **Confirm Password**: Enter the current user's password and confirm it by entering it again.
  - iv. User Group and Primary Group: Select supergroup.
  - v. **Role**: Select the role created in **b** and the **Manager\_viewer** role.

Figure 5-8 Creating a Hive user

Kanager	Homepage Clust	ter - Hosts O&M Audit Tenant Resources System
	User > Create	
System	* Username:	hivetestusr
	* User Type:	Human-Machine     Machine-Machine
Permission ^	* Password:	
User Group	* Confirm Password:	
- Role	User Group:	Add Clear All Create User Group
<ul> <li>Security Policy</li> </ul>		supergroup ×
Domain and Mutual Trust		
Interconnection ~		
Certificate	Primary Group:	Supergroup
OMS	Dala:	Add. Clear All. Create Belo
Component	Role.	
		melestore X managel_nener X
	Description:	
		OK Cancel

- d. Download and install the Hive client by referring to Installing an MRS Client. For example, the current Hive client is installed in the /opt/ hiveclient directory of the active MRS node.
- e. Go to the client installation directory as user **root**.

Example: cd /opt/hiveclient

f. Configure environment variables.

### source bigdata\_env

g. Authenticate the user created in **c** as Kerberos authentication has been enabled for the cluster:

kinit Username in c

Example: kinit hivetestusr

h. Connect to Hive.

### beeline

i. Create a table and insert data into it.

Create a table.

create table user\_info(id string,name string,gender string,age int,addr string);

#### Insert data into the table.

insert into table user\_info(id,name,gender,age,addr) values("12005000201","A","Male",19,"City A");

insert into table user\_info(id,name,gender,age,addr) values("12005000202","B","Male",20,"City B");

insert into table user\_info(id,name,gender,age,addr) values("12005000202","B","Male",20,"City B");

### D NOTE

In the preceding example, data is migrated by creating a table and inserting data into the table. To migrate existing Hive databases and table data, obtain information about them.

• Run the following command on the Hive client to obtain database information:

#### show databases

• Switch to the Hive database to migrate.

use Hive database name

- Displays information about all tables in the current database. show tables
- Query the Hive table creation statement.

show create table Hive table name

The table creation statement obtained from the query needs to be modified to match DLI's syntax before running it on DLI.

- Create a database and table on DLI.
  - a. Log in to the DLI management console. In the navigation pane on the left, choose **SQL Editor**. On the displayed page, set **Engine** to **Spark** and **Queues** to the created SQL queue.

Create a database, for example, **testdb**. For the syntax to create a DLI database, see **Creating a Database**.

create database testdb;

b. Create a table in the database.

You need to modify the table creation statement obtained by running **show create table hive\_table\_name** to comply with DLI's table creation syntax. For the table creation syntax, see **Creating a DLI Table Using the DataSource Syntax**.

create table user\_info(id string,name string,gender string,age int,addr string);

### Step 2: Migrate Data

- 1. Create a CDM connection.
  - a. Create a connection to link CDM to MRS Hive.
    - i. Log in to the CDM console. In the navigation pane on the left, choose **Cluster Management**. On the displayed page, locate the created CDM cluster and click **Job Management** in the **Operation** column.
    - On the Job Management page, click the Links tab. On this tab page, click Create Link. On the displayed page, select MRS Hive and click Next.

Data Warehouse	Data Warehouse Service	Data Lake Insight			
Hadoop	MRS HDFS	MRS HBase	MRS Hive	Apache HDFS	
	Apache HBase	Apache Hive			
Object Storage	Object Storage Service (OBS)	Alibaba Cloud OSS			
File System	FTP	SFTP	HTTP		
Relational Database	RDS for MySQL	RDS for PostgreSQL	RDS for SQL Server	MySQL	
	PostgreSQL	Microsoft SQL Server	Oracle	IBM Db2	
NoSQL	Redis	MongoDB			
Messaging System	Data Ingestion Service	MRS Kafka	Apache Kafka		
Search	Elasticsearch				
Open Beta Test	^				

Figure 5-9 Selecting the MRS Hive connector

iii. Configure the connection as follows:

Parameter	Value
Name	Name of the MRS Hive data source, for example, <b>source_hive</b> .
Manager IP	Manager IP address, which is automatically filled in after you click <b>Select</b> next to the text box and select the MRS Hive cluster.
Authenticatio n Method	Set it to <b>KERBEROS</b> if Kerberos authentication is enabled for the MRS cluster or to <b>SIMPLE</b> if the MRS cluster is a common cluster. In this example, set it to <b>KERBEROS</b> .
Hive Version	Set it to the Hive version you selected when creating the MRS cluster. Select <b>HIVE_3_X</b> if the current Hive version is 3.1.0.
Username	Name of the MRS Hive user created in <b>c</b> .
Password	Password of the MRS Hive user.

Table 5-	5 MRS	Hive	connection	configurat	ions
Table J		TINC	connection	connigurat	10113

Retain the default values for other parameters.

elect Connector			- 2 Configu
* Name	source_hive		
* Connector	Hive		
Hadoop Type	MRS		
* Manager IP 🕜	192.168.7.145	Select	
Authentication Method	KERBEROS	•	
* HIVE Version (?)	HIVE_3_X	•	
* Username 🕜	hivetestusr		
* Password		8	
* OBS storage support ⑦	Yes No		
* Run Mode 🕥	EMBEDDED	•	
Use Cluster Config 🏼 🕜	Yes No		
Show Advanced Attributes			

Figure 5-10 Configuring the connection to MRS Hive

- iv. Click Save.
- b. Create a connection to link CDM to DLI.
  - i. Log in to the CDM console. In the navigation pane on the left, choose **Cluster Management**. On the displayed page, locate the created CDM cluster and click **Job Management** in the **Operation** column.
  - ii. On the Job Management page, click the Links tab, and click Create Link. On the displayed page, select Data Lake Insight and click Next.

Data Warehouse	Data Warehouse Service	Data Lake Insight		
fadoop	MRS HDFS	MRS HBase	MRS Hive	Apache HDFS
	Apache HBase	Apache Hive		
Dbject Storage	Object Storage Service (OBS)	Alibaba Cloud OSS		
ile System	FTP	SFTP	HTTP	
Relational Database	RDS for MySQL	RDS for PostgreSQL	RDS for SQL Server	MySQL
	PostgreSQL	Microsoft SQL Server	Oracle	IBM Db2
loSQL	Redis	MongoDB		
lessaging System	Data Ingestion Service	MRS Kafka	Apache Kafka	
Search	Elasticsearch			

Figure 5-11 Selecting the DLI connector

iii. Set the connection parameters.

Figure 5-12 Setting connection parameters

ielect Connector			🙆 Configu
* Name	source_hive		
* Connector	Hive		
* Hadoop Type	MRS		
* Manager IP	192.168.7.145	Select	
Authentication Method	KERBEROS		
* HIVE Version (2)	HIVE_3_X		
* Username 💮	hivetestusr		
* Password	40		
* OIIX storage support	Yes No		
* Run Mode 🕐	EMBEDDED		
Use Cluster Config	Yes No		
Show Advanced Allebades			

Click Save.

- 2. Create a CDM migration job.
  - a. Log in to the CDM console. In the navigation pane on the left, choose **Cluster Management**. On the displayed page, locate the created CDM cluster and click **Job Management** in the **Operation** column.
  - b. On the **Job Management** page, click the **Table/File Migration** tab. On the displayed tab, click **Create Job**.
  - c. On the **Create Job** page, set job parameters.

Figure 5-13 Setting CDM job parameters

Configure Easis Information		(Z) v	Map Field		(3) Configure Tax
Job Configuration					
* Job Name	hive_io_di				
Source Job Configuration	on		Destination Job Configura	tion	
* Source Link Name	source_hive •		Destination Link Name	dest_di •	
* Database Name 🕥	orbut 😔		* Resource Queue	test_di_jex	
+ Table Name 🕥	user_inte 🖂		• Calabase (1)	teskib 🛛	
madiliode 💮	HDFS *		* Table ()	user_info	
Show Advanced Attributes			Ciear data before import 🛞	Ves No	
× Cancel ) Need					

- i. **Job Name**: Name of the data migration job, for example, **hive\_to\_dli**.
- ii. Set the parameters in the Source Job Configuration area as follows:

Table 5-6 Source job parameters

Parameter	Value
Source Link Name	Select the name of the data source created in <b>1.a</b> .
Database Name	Select the name of the MRS Hive database you want to migrate to DLI. For example, the <b>default</b> database.
Table Name	Select the name of the Hive table. In this example, a database created on DLI and the <b>user_info</b> table are selected.

Parameter	Value
readMode	In this example, set it to <b>HDFS</b> . The options are described as follows:
	There are two read modes available: HDFS and JDBC. HDFS is used by default. If you do not need to use the WHERE condition to filter data or add fields on the field mapping page, select <b>HDFS</b> .
	When using the HDFS mode, data reading performance is good, but it does not support filtering data using WHERE conditions or adding fields on the field mapping page.
	The JDBC mode allows you to use WHERE conditions to filter data or add fields on the field mapping page.

For details about parameter settings, see **From Hive**.

iii. Set the parameters in the **Destination Job Configuration** area as follows:

Parameter	Value
Destination Link Name	Select the DLI data source connection created in <b>1.b</b> .
Resource Queue	Select a created DLI SQL queue.
Database Name	Select a created DLI database. In this example, database <b>testdb</b> created in <b>Create a database</b> <b>and table on DLI</b> is selected.
Table Name	Select the name of a table in the database. In this example, table <b>user_info</b> created in <b>Create a database and table on DLI</b> is selected.
Clear data before import	Whether to clear data in the destination table before data import. In this example, set it to <b>No</b> . If set to <b>Yes</b> , data in the destination table will be cleared before the task is started.

Table 5-7 Destination	job	parameters
-----------------------	-----	------------

For details about parameter settings, see **To DLI**.

- 3. Click **Next**. The **Map Field** page is displayed. CDM automatically matches the source and destination fields.
  - You can drag any unmatched fields to match them.

- If the type is automatically created at the migration destination, you need to configure the type and name of each field.
- CDM allows for field conversion during migration. For details, see Field Conversion.

### Figure 5-14 Field mapping

Juster Management J. cdm-leat-60378326	/ Table/Tile Migration / hive_do_dli								
() Configure Basic Information —					Ø =	ış Field			- (3) Configure Task
Source Field					0/	Destination Field			<b>∓ e</b> e
Name	Example Value	Туре	Operation			Namo	Туре	Operation	
м		string	2	Q	Ø	-)+ 1 H	shing	T	
name		string	2	Q	Ø	-ini ram	string	π	
gender		stolog	8	Q	<b>B</b> 0	-> pender	string	π	
494		in .	8	Q	<b>B</b> 0	- apa	int	π	
addr		shing	8	Q	<b>U</b> 0	- addr	ating	π	
× Cancel < Pre-	rous → Nect 🗈 Sava								

4. Click **Next** and set task parameters. Typically, retain the default values for all parameters.

In this step, you can configure the following optional features:

- Retry Upon Failure: If the job fails to be executed, you can determine whether to automatically retry. Retain the default value Never.
- Group: Select the group to which the job belongs. The default group is DEFAULT. On the Job Management page, jobs can be displayed, started, or exported by group.
- Scheduled Execution: For how to configure scheduled execution, see Scheduling Job Execution. Retain the default value No.
- **Concurrent Extractors**: Enter the number of extractors to be concurrently executed. Retain the default value **1**.
- Write Dirty Data: Set this parameter if data that fails to be processed or filtered out during job execution needs to be written to OBS. Before writing dirty data, create an OBS link. You can view the data on OBS later. Retain the default value No, meaning dirty data is not recorded.
- 5. Click **Save and Run**. On the **Job Management** page, you can view the job execution progress and result.

Figure 5-15 Job progress and execution result

TableTile Migration	Entire	08 Migration Links	Agents Settings									
() Corato Job	Run	🖞 Delete				C Feedback	pot 🖸 Impot	Schedule	* Al statutes	* Job name	• Jab name or link type	Q C
0 / Þ T	¢	Nome 🕫	Link Details	Created By 20	Last Execution Time 32	Daration 32	Write Statistics	Status	Group Name	Operation		
Enter a group name.	1	Mingla di	source_hive-dest_di	100378325	Mar 25, 2022 20:41 10 GMT+05:00	10x	Written rows: 3	Succeeded	DEFAULT	Ran   Historical Record   Edit	More +	
Groups												
DEFAULT												

### Step 3: Query Results

Once the migration job is complete, check whether the Hive table data has been migrated to the **user\_info** table. Specifically, do as follows: Log in to the DLI management console and choose **SQL Editor**. On the displayed page, set **Engine** to **Spark**, **Queues** to the created SQL queue, and **Databases** to the database created in **a**. Then, execute the following query statement: select \* from user\_info;

### Figure 5-16 Querying migrated data

	Engine	spark •	Queues	test_cli_twx	* Datat	bases [	tesidb +	O Execut	le Format	Refer Syntax	Settings	More +
1rrate database testob) 1rrate tables con_info[id str 2 solart from user_info]	tag,name string,gender string,nge lar,bader stri	a);	,									
Line 4, Column 1				-				Execute: Ctrl	+Enter, Find: CM+F, F	ormat: Shift+Alt+F, Veri	fy Syntax: Old+Q,	Fullscreen: F11
Executed Queries (Last Day) Vie	w Result											Clear All
Deecuted successfully     Query select * from user_info     Job ID c34974df-2cbd-46fd-aed6-2c	1d17c1fe00											
The query takes 4.12s, and 1.34 KB scann	red.A maximum of 1,000 records can be displayed.									Enter a keyword.	Q	<u>u</u> Cí <u>4</u>
id J≘	name ↓Ξ	gender ↓⊟					age JΞ			addr .¦≘		
12005000201	A						19					
12005000202	8						20					
12005000202	8						20					

## 5.2.4 Example Typical Scenario: Migrating Data from Kafka to DLI

This section describes how to use CDM's data synchronization to migrate data from MRS Kafka to DLI.

### Prerequisites

 You have created a DLI SQL queue. For how to create a DLI queue, see Creating a Queue.

### 

Set **Type** to **For SQL** when buying a queue.

- You have created an MRS security cluster that contains the Kafka component. For how to create an MRS cluster, see **Purchasing a Custom Cluster**.
  - In this example, the version of the MRS cluster is 3.1.0.
  - You have enabled Kerberos authentication for the MRS cluster.
- You have created a CDM cluster. For how to create a CDM cluster, see Creating a CDM Cluster.

### D NOTE

- To connect the cluster to an on-premises database as the destination data source, you can use either Internet or Direct Connect. If the Internet is used, make sure that an EIP has been bound to the CDM cluster, the security group of CDM allows outbound traffic from the host where the on-premises data source is located, the host where the data source is located can access the Internet, and the connection port has been enabled in the firewall rules.
- To successfully connect to cloud services MRS and GaussDB (DWS) as data sources, the following requirements must be met:

i. If the CDM cluster and the cloud service are in different regions, they must be connected through either the Internet or Direct Connect. If the Internet is used, make sure that an EIP has been bound to the CDM cluster, the host where the data source is located can access the Internet, and the port has been enabled in the firewall rules.

ii. If the CDM cluster and the cloud service are in the same region, instances in the same VPC, subnet, and security group can communicate with each other by default. If the CDM cluster and the cloud service are in the same VPC but in different subnets or security groups, you must configure routing and security group rules.

For how to configure routing rules, see **Configure routes**. For how to configure security group rules, see **Security Group Configuration Examples**.

iii. The cloud service instance and the CDM cluster belong to the same enterprise project. If they do not, change the enterprise project of the workspace.

In this example, the VPC, subnet, and security group of the CDM cluster match those of the MRS cluster.

### Step 1: Prepare Data

- Create a Kafka topic for the MRS cluster and send messages to the topic.
  - a. Log in to MRS Manager by referring to Accessing FusionInsight Manager.
  - b. Click **System** and choose **Permission** > **User**. On the displayed page, set the following parameters:
    - i. **Username**: Enter a username. In this example, enter **testuser2**.
    - ii. User Type: Select Human-Machine.
    - iii. **Password** and **Confirm Password**: Enter the current user's password and confirm it by entering it again.
    - iv. User Group and Primary Group: Select kafkaadmin.
    - v. Role: Select Manager\_viewer.

SusionInsight Manager	Homepage Clus	ter + Hosts O&M Audit <b>System</b>
	User > Create	
System	• Username:	testuser2
	* User Type:	Human-Machine     Machine-Machine
• User	* Password:	
User Group	* Confirm Password:	
Role	User Group:	Add Clear All Create User Group
Security Policy		kafkaadmin 🗙
Domain and Mutual Trust		
Interconnection ~		
Certificate	Primary Group:	kafkaadmin 👻
OMS	Role:	Add Clear All Create Role
Component		Manager_viewer 🗙
	Description:	
		OK Cancel

Figure 5-17 Creating a Kafka user

- c. On MRS Manager, choose Cluster > Name of the desired cluster > Service > ZooKeeper > Instance. On the displayed page, obtain the IP address of the ZooKeeper instance.
- On MRS Manager, choose Cluster > Name of the desired cluster > Service > Kafka > Instance. On the displayed page, obtain the IP address of the Kafka instance.
- e. Download and install the Kafka client by referring to **Installing an MRS Client**. For example, the current Hive client is installed in the **/opt/ kafkaclient** directory of the active MRS node.
- f. Go to the client installation directory as user **root**.

Example: cd /opt/kafkaclient

g. Configure environment variables.

### source bigdata\_env

h. Authenticate the user created in **b** as Kerberos authentication has been enabled for the cluster:

kinit Username in b

Example: kinit testuser2

i. Create a Kafka topic named **kafkatopic**. kafka-topics.sh --create --zookeeper *IP address 1 of the node where the ZooKeeper role is*2181,*IP address 2 of the node where the ZooKeeper role is*.2181,*IP address 3 of the node where the ZooKeeper role is*.2181/kafka --replication-factor 1 --partitions 1 --topic kafkatopic

In this command, *IP address of the node where the ZooKeeper role is* is the IP address of the ZooKeeper instance obtained in **c**.

j. Send a test message to **kafkatopic**.

kafka-console-producer.sh --broker-list *IP address 1 of the node where the Kafka role is*:21007;*IP address 2 of the node where the Kafka role is*:21007;*IP address 3 of the node where the Kafka role is*:21007 --topic kafkatopic --producer.config /opt/kafkaclient/Kafka/kafka/config/ producer.properties In this command, *IP address of the node where the Kafka role is* is the IP address of the Kafka instance obtained in **d**.

The content of the test message is as follows: {"PageViews":5, "UserID":"4324182021466249494", "Duration":146,"Sign":-1}

- Create a database and table on DLI.
  - a. Log in to the DLI management console. In the navigation pane on the left, choose **SQL Editor**. On the displayed page, set **Engine** to **Spark** and **Queues** to the created SQL queue.

Create a database, for example, **testdb**. For the syntax to create a DLI database, see **Creating a Database**. create database testdb;

 b. Create a table in the database. For the table creation syntax, see Creating a DLI Table Using the DataSource Syntax. CREATE TABLE testdlitable(value STRING);

### Step 2: Migrate Data

- 1. Create a CDM connection.
  - a. Create a connection to link CDM to MRS Kafka.
    - i. Log in to the CDM console. In the navigation pane on the left, choose **Cluster Management**. On the displayed page, locate the created CDM cluster and click **Job Management** in the **Operation** column.
    - ii. On the Job Management page, click the Links tab and click Create Link. On the displayed page, select MRS Kafka and click Next.



iii. Configure the connection as follows:

Parameter	Value
Name	Name of the MRS Kafka data source, for example, <b>source_kafka</b> .

Table 5-8 MRS Kafka connection configurations

Parameter	Value
Manager IP	Manager IP address, which is automatically filled in after you click <b>Select</b> next to the text box and select the MRS Kafka cluster.
Username	Name of the MRS Kafka user created in <b>b</b> .
Password	Password of the MRS Kafka user.
Authenticatio n Method	Set it to <b>KERBEROS</b> if Kerberos authentication is enabled for the MRS cluster or to <b>SIMPLE</b> if the MRS cluster is a common cluster. In this example, set it to <b>KERBEROS</b> .

For more details about the parameters, see Link to Kafka.

Select Connector			
* Name	source_kafka		
* Connector	Kafka		
★ Kafka Type	MRS		
* Manager IP	1		Select
* Username (?)	testuser2		
* Password	•••••	Ø	
Authentication Method	KERBEROS	•	
Show Advanced Attributes			

Figure 5-19 Configuring the MRS Kafka connection

- iv. Click Save.
- b. Create a connection to link CDM to DLI.
  - i. Log in to the CDM console. In the navigation pane on the left, choose **Cluster Management**. On the displayed page, locate the created CDM cluster and click **Job Management** in the **Operation** column.
  - ii. On the Job Management page, click the Links tab, and click Create Link. On the displayed page, select Data Lake Insight and click Next.

Data Warehouse	Data Warehouse Service	Data Lake Insight			
fadoop	MRS HDFS	MRS HBase	MRS Hive	Apache HDFS	
	Apache HBase	Apache Hive			
Dbject Storage	Object Storage Service (OBS)	Alibaba Cloud OSS			
file System	FTP	SFTP	HTTP		
Relational Database	RDS for MySQL	RDS for PostgreSQL	RDS for SQL Server	MySQL	
	PostgreSQL	Microsoft SQL Server	Oracle	IBM Db2	
VoSQL	Redis	MongoDB			
vlessaging System	Data Ingestion Service	MRS Kafka	Apache Kafka		
Search	Elasticsearch				

Figure 5-20 Selecting the DLI connector

iii. Set the connection parameters. For details about parameter settings, see Link to DLI.

Figure 5-21 Setting connection parameters

Select Connector	uan-en-eou/racar / Lans / L <b>rear Lan</b>	Configure
* Name	dest_di	
* Connector	DU *	
* AK (?)		
* SK 🕐	······ · · · · · · · · · · · · · · · ·	
* Project ID	05/ 1	
× Cancel	C Previous 0° Test	

- iv. Click Save.
- 2. Create a CDM migration job.
  - a. Log in to the CDM console. In the navigation pane on the left, choose **Cluster Management**. On the displayed page, locate the created CDM cluster and click **Job Management** in the **Operation** column.
  - b. On the **Job Management** page, click the **Table/File Migration** tab. On the displayed tab, click **Create Job**.
  - c. On the **Create Job** page, set job parameters.

Figure 5-22 Setting CDM job parameters

Configure Sasic Information		2) Map From		(3) Contigure Task
Job Configuration	nt			
Source Job Configuration		Destination Job Configuration	n	
* Source Link Name	source_kafea	<ul> <li>Destination Link Name</li> </ul>	dest_fil	
* Topica 🕐	kathatopic	* Resource Queue	Test68	
* Data Pormat	000,0R5_3000 +	A Database (7)	Inskib O	
* Offset Parameter 🕥	EARLIEST *	A Table (2)	International C	
* Permanent Running 💮	Yes No	Ciear data before import	Yan No	
* Pul Data Timeout 💮	13			
Walt Data Treeval				
* Consumer Group ID	ecample.group1			
Show Advanced Attributes				
× Cencel > Not Save				

- i. Job Name: Name of the data migration job, for example, test.
- ii. Set the parameters in the **Source Job Configuration** area as follows:

Parameter	Value
Source Link Name	Select the name of the data source created in <b>1.a</b> .
Topics	Name of the topics you want to migrate to DLI. You can select one or more topics. Example: <b>kafkatopic</b> .
Data Format	Select the message format as needed. In this example, <b>CDC (DRS_JSON)</b> is selected, indicating that the source data will be parsed in DRS_JSON format.
Offset Parameter	Initial offset when data is pulled from Kafka. In this example, select <b>EARLIEST</b> . The options are:
	• Latest: Maximum offset, meaning that the latest data will be pulled.
	• <b>Earliest</b> : Minimum offset, meaning that the earliest data will be pulled.
	• <b>Submitted</b> : The submitted data is pulled.
	• <b>Time Range</b> : Data within a time range is pulled.
Permanent Running	Whether a job runs permanently. In this example, set it to <b>No</b> .
Pull Data Timeout	Maximum minutes allowed for a continuous data pulling. In this example, set it to <b>15</b> .
Wait Data Timeout	(Optional) Maximum seconds allowed for waiting data reading. In this example, leave this parameter blank.
Consumer Group ID	Consumer group ID. The default Kafka message group ID <b>example-group1</b> is used.

Table 5-9 Source job parameters

For details about parameter settings, see From Apache Kafka.

iii. Set the parameters in the **Destination Job Configuration** area as follows:

Parameter	Value
Destination Link Name	Select the DLI data source connection created in <b>1.b</b> .
Resource Queue	Select a created DLI SQL queue.
Database Name	Select a created DLI database. In this example, database <b>testdb</b> created in <b>Create a database</b> <b>and table on DLI</b> is selected.
Table Name	Select the name of a table in the database. In this example, table <b>testdlitable</b> created in <b>Create a database and table on DLI</b> is selected.
Clear data before import	Whether to clear data in the destination table before data import. In this example, set it to <b>No</b> . If set to <b>Yes</b> , data in the destination table will be cleared before the task is started.

Table 5-10 Destination job parameters

For details about parameter settings, see **To DLI**.

- 3. Click **Next**. The **Map Field** page is displayed. CDM automatically matches the source and destination fields.
  - You can drag any unmatched fields to match them.
  - If the type is automatically created at the migration destination, you need to configure the type and name of each field.
  - CDM allows for field conversion during migration. For details, see Field Conversion.

### Figure 5-23 Field mapping



4. Click **Next** and set task parameters. Typically, retain the default values for all parameters.

In this step, you can configure the following optional features:

- Retry Upon Failure: If the job fails to be executed, you can determine whether to automatically retry. Retain the default value Never.
- Group: Select the group to which the job belongs. The default group is DEFAULT. On the Job Management page, jobs can be displayed, started, or exported by group.
- Scheduled Execution: For how to configure scheduled execution, see
   Scheduling Job Execution. Retain the default value No.
- **Concurrent Extractors**: Enter the number of extractors to be concurrently executed. Retain the default value **1**.

- Write Dirty Data: Set this parameter if data that fails to be processed or filtered out during job execution needs to be written to OBS. Before writing dirty data, create an OBS link. You can view the data on OBS later. Retain the default value No, meaning dirty data is not recorded.
- 5. Click **Save and Run**. On the **Job Management** page, you can view the job execution progress and result.

Figure 5-24 Job progress and execution result

Table/Tile Migration Entire I	28 Migration Links Agen	ta Settingo									
🕢 Create Juli	C Debelle				C Parellant C Eq.	et Direct	202milute	* All statuses	* Job name	<ul> <li>Jub name or link type</li> </ul>	Q C
0 / P 0 (	Ramo (II	Link Defails	Created By 28	Last Execution Time 20	Duration 28	Wite Materica	Status	Group Marrie	Operation		
Enter a group name. Q	net	source_kafka=dest_dli		Apr 07, 2922 95:42:25 GMT+08:09	15m 15a	Virillen rave: 3	Gutteeded	OFFAULT	Ram   Historical Record   Edit	More +	
Groupe											
DEFAULT											

### Step 3: Query Results

Once the migration job is complete, check whether the Kafka table data has been migrated to the **testdlitable** table. Specifically, do as follows: Log in to the DLI management console and choose **SQL Editor**. On the displayed page, set **Engine** to **Spark**, **Queues** to the created SQL queue, and **Databases** to the database created in **a**. Then, execute the following query statement: select \* from testdlitable;

## 5.2.5 Example Typical Scenario: Migrating Data from Elasticsearch to DLI

This section describes how to use CDM's data synchronization to migrate data from a CSS Elasticsearch cluster to DLI. Data in a self-built Elasticsearch cluster can also be bidirectionally synchronized with DLI using CDM.

### Prerequisites

 You have created a DLI SQL queue. For how to create a DLI queue, see Creating a Queue.

Set **Type** to **For SQL** when buying a queue.

 You have created a CSS Elasticsearch cluster. For how to create a CSS cluster, see Creating a CSS Cluster.

In this example, the version of the created CSS cluster is 7.6.2, and the cluster is a non-security cluster.

 You have created a CDM cluster. For how to create a CDM cluster, see Creating a CDM Cluster.

### D NOTE

- To connect the cluster to an on-premises database as the destination data source, you can use either Internet or Direct Connect. If the Internet is used, make sure that an EIP has been bound to the CDM cluster, the security group of CDM allows outbound traffic from the host where the on-premises data source is located, the host where the data source is located can access the Internet, and the connection port has been enabled in the firewall rules.
- To successfully connect to the cloud service CSS as a data source, the following requirements must be met:

i. If the CDM cluster and the cloud service are in different regions, they must be connected through either the Internet or Direct Connect. If the Internet is used, make sure that an EIP has been bound to the CDM cluster, the host where the data source is located can access the Internet, and the port has been enabled in the firewall rules.

ii. If the CDM cluster and the cloud service are in the same region, instances in the same VPC, subnet, and security group can communicate with each other by default. If the CDM cluster and the cloud service are in the same VPC but in different subnets or security groups, you must configure routing and security group rules.

For how to configure routing rules, see **Configure routes**. For how to configure security group rules, see **Security Group Configuration Examples**.

iii. The cloud service instance and the CDM cluster belong to the same enterprise project. If they do not, change the enterprise project of the workspace.

In this example, the VPC, subnet, and security group of the CDM cluster match those of the CSS cluster.

### Step 1: Prepare Data

- Create an index for the CSS cluster and import data.
  - a. Log in to the CSS management console and choose **Clusters** > **Elasticsearch** from the navigation pane on the left.
  - b. On the **Clusters** page, locate the created CSS cluster and click **Access Kibana** in the **Operation** column.
  - c. In the navigation pane of Kibana, choose **Dev Tools**.
  - d. On the displayed **Console** page, run the following command to create index **my\_test**:

```
PUT /my_test
{
    "settings": {
        "number_of_shards": 1
    },
    "mappings": {
            "properties": {
            "productName": {
               "type": "text",
               "analyzer": "ik_smart"
            },
            "size": {
               "type": "keyword"
            }
        }
    }
}
```

e. Import data to the my\_test index. POST /my\_test/\_doc/\_bulk {"index":{}} {"productName":"2017 Autumn New Shirts for Women", "size":"L"} {"index":{}}
{"productName":"2017 Autumn New Shirts for Women", "size":"M"}
{"index":{}}
{"productName":"2017 Autumn New Shirts for Women", "size":"S"}
{"index":{}}
{"productName":"2018 Spring New Jeans for Women", "size":"M"}
{"index":{}}
{"productName":"2018 Spring New Jeans for Women", "size":"S"}
{"index":{}}
{"productName":"2017 Spring Casual Pants for Women", "size":"L"}
{"index":{}}
{"productName":"2017 Spring Casual Pants for Women", "size":"S"}

If **errors** is **false** in the command output, the data is imported.

- Create a database and table on DLI.
  - a. Log in to the DLI management console. In the navigation pane on the left, choose **SQL Editor**. On the displayed page, set **Engine** to **Spark** and **Queues** to the created SQL queue.

Create a database, for example, **testdb**. For the syntax to create a DLI database, see **Creating a Database**.

create database testdb;

 b. Create a table in the database. For the table creation syntax, see Creating a DLI Table Using the DataSource Syntax. create table tablecss(size string, productname string);

### Step 2: Migrate Data

- 1. Create a CDM connection.
  - a. Create a connection to link CDM to the data source CSS.
    - i. Log in to the CDM console. In the navigation pane on the left, choose **Cluster Management**. On the displayed page, locate the created CDM cluster and click **Job Management** in the **Operation** column.
    - On the Job Management page, click the Links tab, and click Create Link. On the displayed page, select Cloud Search Service and click Next.

lect Connector				
ata Warehouse	Data Warehouse Service	Data Lake Insight		
fadoop	MRS HDFS	MRS HBase	MRS Hive	Apache HDFS
	Apache HBase	Apache Hive		
bject Storage	Object Storage Service (OBS)	Alibaba Cloud OSS		
ile System	FTP	SFTP	HTTP	
telational Database	RDS for MySQL	RDS for PostgreSQL	RDS for SQL Server	MySQL
	PostgreSQL	Microsoft SQL Server	Oracle	IBM Db2
loSQL	Redis	MongoDB		
fessaging System	Data Ingestion Service	MRS Kafka	Apache Kafka	
iearch	Elasticsearch			
)pen Beta Test	~			
	FusionInsight LibrA	FusionInsight HDFS	FusionInsight HBase	FusionInsight Hive
	Qiniu Cloud Object Storage (KODO)	Amazon S3	Tencent Cloud COS	Distributed Database Middleware
	SAP HANA	MYCAT	DM	Sharding Database
	Distributed Cache Service	Document Database Service	CloudTable Service	CloudTable Service (OpenTSDB)
	Cassandra	DMS Kafka	Cloud Search Service	
× Cancel > Next				

Figure 5-25 Selecting the CSS connector

iii. Configure the connection. For details about parameter settings, see Link to Elasticsearch/CSS.

Parameter	Value
Name	Name of the CSS data source, for example, <b>source_css</b> .
Elasticsearch Server List	Elasticsearch server list, which is automatically displayed after you click <b>Select</b> next to the text box and select the CSS cluster.
Security Mode Authenticatio n	If you have enabled the security mode for the CSS cluster, set this parameter to <b>Yes</b> . Otherwise, set this parameter to <b>No</b> . In this example, set it to <b>No</b> .

Table 5-11 CSS data source configuration

elect Connector		2
* Name	source_css	
* Connector	Elasicsearch v	
* Elasticsearch Server List	Select	
Security mode Authentication	Yês No	
X Cancel <pre></pre>	of Test	

Figure 5-26 Configuring the CSS connection

- iv. Click Save.
- b. Create a connection to link CDM to DLI.
  - i. Log in to the CDM console. In the navigation pane on the left, choose **Cluster Management**. On the displayed page, locate the created CDM cluster and click **Job Management** in the **Operation** column.
  - ii. On the Job Management page, click the Links tab, and click Create Link. On the displayed page, select Data Lake Insight and click Next.

er Management / cdm /	Links / Create Link				
Select Connector					
Data Warehouse	Data Warehouse Service	Data Lake Insight			
Hadoop	MRS HDFS	MRS HBase	MRS Hive	Apache HDFS	
	Apache HBase	Apache Hive			
Object Storage	Object Storage Service (OBS)	Alibaba Cloud OSS			
File System	FTP	SFTP	HTTP		
Relational Database	RDS for MySQL	RDS for PostgreSQL	RDS for SQL Server	MySQL	
	PostgreSQL	Microsoft SQL Server	Oracle	IBM Db2	
NoSQL	Redis	MongoDB			
Messaging System	Data Ingestion Service	MRS Kafka	Apache Kafka		
Search	Elasticsearch				
Open Beta Test X Cancel					

**Figure 5-27** Selecting the DLI connector

iii. Set the connection parameters. For details about parameter settings, see Link to DLI.

Figure 5-28 Setting connection parameters

sect Connector		0
* Name	source_css	
* Connector	Elasticsearch v	
* Elasticsearch Server List 🕐	1 Select	
Security mode Authentication (2)	Yes No	
X Cancel < Previous	ort Test	

- iv. Click Save.
- 2. Create a CDM migration job.
  - a. Log in to the CDM console. In the navigation pane on the left, choose **Cluster Management**. On the displayed page, locate the created CDM cluster and click **Job Management** in the **Operation** column.
  - b. On the **Job Management** page, click the **Table/File Migration** tab. On the displayed tab, click **Create Job**.
  - c. On the **Create Job** page, set job parameters.

Figure 5-29 Setting CDM job parameters

* Job Name Css_to_dli	1				
Source Job Cor	nfiguration		Destination Job Conf	iguration	
* Source Link Name	source_css		* Destination Link Name	dest_di	
* Index 📀	my_test	Θ	* Resource Queue (?)	bz	Θ
* Type 🌘	_doc	Θ	* Database (?)	testdb	Θ
Show Advanced Altribute	8		* Table (?)	tablecss	Θ
			Clear data before import	Yes No	

- i. Job Name: Name of the data migration job, for example, css\_to\_dli.
- ii. Set the parameters in the Source Job Configuration area as follows:

Table 5-12 Source job parameters

Parameter	Value
Source Link Name	Select the name of the data source created in <b>1.a</b> .
Index	Select the Elasticsearch index created for the CSS cluster. In this example, the <b>my_test</b> index created in <b>Create an index for the CSS cluster and import data</b> is used.
	Only lowercase letters are allowed.
Туре	Elasticsearch type, which is similar to the table name of a relational database. Only lowercase letters are allowed. Example: <b>_doc</b> .

For details about other parameters, see From Elasticsearch or CSS.

### iii. Set the parameters in the **Destination Job Configuration** area as follows:

Parameter	Value
Destination Link Name	Select the DLI data source connection created in <b>1.b</b> .
Resource Queue	Select a created DLI SQL queue.
Database Name	Select a created DLI database. In this example, the database <b>testdb</b> created in <b>Create a database and table on DLI</b> is selected.
Table Name	Select the name of a table in the database. In this example, the table <b>tablecss</b> created in <b>Create a database and table on DLI</b> is used.
Clear data before import	Whether to clear data in the destination table before data import. In this example, set it to <b>No</b> . If set to <b>Yes</b> , data in the destination table will be cleared before the task is started.

Table 5-13 Destination	job	parameters
------------------------	-----	------------

For details about parameter settings, see **To DLI**.

- 3. Click **Next**. The **Map Field** page is displayed. CDM automatically matches the source and destination fields.
  - You can drag any unmatched fields to match them.
  - If the type is automatically created at the migration destination, you need to configure the type and name of each field.
  - CDM allows for field conversion during migration. For details, see Field Conversion.

### Figure 5-30 Field mapping

Source Field				(	∋ ∥	Destination Field		<u></u>
Type	Name	Example Value	Operation			Name	Type	Operation
string	productName		8	Q	<u>ت</u> و	> productname	sting	Ψ
keyword	size	L	8	Q	ii و ا	⊳ size	sting	ΰ
				(	∋ ∥			🕁 🕲 🖸

4. Click **Next** and set task parameters. Typically, retain the default values for all parameters.

In this step, you can configure the following optional features:

- **Retry Upon Failure**: If the job fails to be executed, you can determine whether to automatically retry. Retain the default value **Never**.
- Group: Select the group to which the job belongs. The default group is DEFAULT. On the Job Management page, jobs can be displayed, started, or exported by group.

- Scheduled Execution: For how to configure scheduled execution, see Scheduling Job Execution. Retain the default value No.
- **Concurrent Extractors**: Enter the number of extractors to be concurrently executed. Retain the default value **1**.
- Write Dirty Data: Set this parameter if data that fails to be processed or filtered out during job execution needs to be written to OBS. Before writing dirty data, create an OBS link. You can view the data on OBS later. Retain the default value No, meaning dirty data is not recorded.
- 5. Click **Save and Run**. On the **Job Management** page, you can view the job execution progress and result.

Figure 5-31 Job progress and execution result

Table/File Migration	Entire	DB Migration Links A	gents Settings						
@ Create Job	💰 Run	C Delete			Feedback     Expert	C Import Schedule	• All statuses	• Job name • Job name or link type Q	С
0 / > 0	¢	Name JE	Link Details	Created By JE Last Execution Time JE	Duration JE Write Sta	diatica Slatua	Group Name	Operation	
Enter a group name.	Q	0_00_00_0	source_css-dest_dl	e_dici_d003 Apr 11, 2022 19:29:39 0MT+08:00	Ten 195 Wittlen ro	rvis: 7 O Succeeded	DEPAULT	Run   Hatorical Record   Edit   More +	
Groups									
DEFAULT									

### Step 3: Query Results

Once the migration job is complete, check whether the CSS table data has been migrated to the **tablecss** table. Specifically, do as follows: Log in to the DLI management console and choose **SQL Editor**. On the displayed page, set **Engine** to **Spark**, **Queues** to the created SQL queue, and **Databases** to the database created in **a**. Then, execute the following query statement: select \* from tablecss;

Figure 5-32 Querying	migrated data
----------------------	---------------

Engine spark 👻 Qu	eues bzq_ · *	Databases testdb	Execute     Forma	t Refer Syntax	Settings More •
1 select * from tablecss;					
Line 1 Column 1			Execute: Ctrl+Enter Find: Ctrl+	-F Format Shift+Alt+F Verify S	wntax: Ctri+Q Euliscreen: E11
		-			
Executed Queries (Last Day) View Result					Clear All
Resulti 0					
Executed successfully					
Query select * from tablecss					
Job ID b5e135d0-a17f-4acc-b3d4-280853e2326a					
The query takes 35.29s, and 1.36 KB scanned.A maximum of 1,000 records can be	displayed.			Enter a keyword.	Q 🔟 🖸 🛓
size ↓Ξ	productname ↓Ξ				
L	2017 Autumn New Shirts	for Women			
S	2017 Autumn New Shirts	s for Women			
S	2018 Spring New Jeans	for Women			
S	2017 Spring Casual Pant	ts for Women			
м	2017 Autumn New Shirts	s for Women			
M	2018 Spring New Jeans	for Women			

## 5.2.6 Example Typical Scenario: Migrating Data from RDS to DLI

This section describes how to use CDM's data synchronization to migrate data from an RDS DB instance to DLI. Data in other relational databases can also be bidirectionally synchronized with DLI using CDM.

### Prerequisites

• You have created a DLI SQL queue. For how to create a DLI queue, see Creating a Queue.

### 

Set **Type** to **For SQL** when buying a queue.

- You have created an RDS for MySQL DB instance. For how to create an RDS cluster, see **Buying an RDS for MySQL Instance**.
  - In this example, the RDS DB engine is MySQL.
  - In this example, the DB engine version is 5.7.
- You have created a CDM cluster. For how to create a CDM cluster, see Creating a CDM Cluster.

### **NOTE**

- To connect the cluster to an on-premises database as the destination data source, you can use either Internet or Direct Connect. If the Internet is used, make sure that an EIP has been bound to the CDM cluster, the security group of CDM allows outbound traffic from the host where the on-premises data source is located, the host where the data source is located can access the Internet, and the connection port has been enabled in the firewall rules.
- If the data source is RDS or MRS, the network must meet the following requirements:

i. If the CDM cluster and the cloud service are in different regions, they must be connected through either the Internet or Direct Connect. If the Internet is used, make sure that an EIP has been bound to the CDM cluster, the host where the data source is located can access the Internet, and the port has been enabled in the firewall rules.

ii. If the CDM cluster and the cloud service are in the same region, instances in the same VPC, subnet, and security group can communicate with each other by default. If the CDM cluster and the cloud service are in the same VPC but in different subnets or security groups, you must configure routing and security group rules.

For how to configure routing rules, see **Configure routes**. For how to configure security group rules, see **Security Group Configuration Examples**.

iii. The cloud service instance and the CDM cluster belong to the same enterprise project. If they do not, change the enterprise project of the workspace.

In this example, the VPC, subnet, and security group of the CDM cluster match those of the RDS for MySQL DB instance.

### Step 1: Prepare Data

- Create a database and table on the RDS for MySQL DB instance.
  - a. Log in to the RDS management console. In the navigation pane on the left, choose **Instances**. On the displayed page, locate the target DB instance and click **Log In** in the **Operation** column.
  - b. In the displayed login dialog box, enter the correct username and password and click **Log In**.

- c. On the **Databases** page, click **Create Database**. In the displayed dialog box, enter **testrdsdb** as the database name and retain default values for the rest parameters. Then, click **OK**.
- d. In the **Operation** column of row where the created database is, click **SQL Window** and enter the following statement to create a table:

CREATE TABLE tabletest ( `id` VARCHAR(32) NOT NULL, `name` VARCHAR(32) NOT NULL, PRIMARY KEY (`id`)

- ENGINE = InnoDB
- DEFAULT CHARACTER SET = utf8mb4;
- e. Insert data into the table. insert into tabletest VALUES ('123','abc'); insert into tabletest VALUES ('456','efg'); insert into tabletest VALUES ('789','hij');
- f. Query table data. select \* from tabletest;

### Figure 5-33 Querying table data

Home SQL Window X				
Current Database testralsabi 📀	🗿 Maatar Switch SQL Execution Nada    Indem	e Name: rds-lest-00376328   192.168.0.157.3306   Charac	ther Set 1485 V	Save Executed SQL Statements ()
Database Nethodo ∨ Tobia Viens Phone search by k_   0, C F ⊠ tableted	Insects 302 (FT)      Faced 502 (FT)	$\overline{\mathbf{S}}$ Kanada 1924 Pila (Pili) ( 1924 Pilavilla $\overline{\mathbf{v}}$ )		SCL vysd Provyd 🛞 🌑 Full Sones 🗶
	Executed SQL Statements Messages Result Set	×		🔘 Overveille Made 🛇
	The following in the execution result set of select * from tabletest.		Click on the cell to edit the data. After adding or editing, you need to submit and save the changes.	Copy Rose Copy Column v Column Settings v
		M	nanoa	
		142	4H	
	2	454	**8	
	3	789	NJ	

- Create a database and table on DLI.
  - a. Log in to the DLI management console. In the navigation pane on the left, choose **SQL Editor**. On the displayed page, set **Engine** to **Spark** and **Queues** to the created SQL queue.

Create a database, for example, **testdb**. For the syntax to create a DLI database, see **Creating a Database**.

create database testdb;

b. On the **SQL Editor** page, set **Databases** to **testdb** and run the following table creation statement to create a table in the database. For the table creation syntax, see **Creating a DLI Table Using the DataSource Syntax**. create table tabletest(id string,name string);

### Step 2: Migrate Data

- 1. Create a CDM connection.
  - a. Create a connection to the RDS database.
    - i. Log in to the CDM console. In the navigation pane on the left, choose **Cluster Management**. On the displayed page, locate the created CDM cluster and click **Job Management** in the **Operation** column.
    - ii. If this is your first time creating a connection to RDS for MySQL, you need to upload the MySQL driver. Click the **Links** tab and click **Driver Management**.
    - Download the MySQL driver to your local PC by referring to Managing Drivers and extract the driver package to obtain the JAR file.
For example, download the **mysql-connector-java-5.1.48.zip** package and extract it to obtain the driver file **mysql-connector-java-5.1.48.jar**.

- iv. Return to the Driver Management page. Locate the MYSQL driver and click Upload in the Operation column. In the Import Driver File dialog box, click Select File and upload the driver file obtained in 1.a.iii.
- v. On the **Driver Management** page, click **Back** in the lower left corner to return to the **Links** tab. Click **Create Link**, select **RDS for MySQL**, and click **Next**.
- vi. Configure the connection as follows:

Parameter	Value
Name	Name of the RDS data source, for example, <b>source_rds</b> .
Database Server	Click <b>Select</b> next to the text box and click the name of the created RDS DB instance. The database server address is automatically entered.
Port	Port number of the RDS DB instance. The value is automatically filled in after you select the database server.
Database Name	Name of the RDS DB instance you want to migrate. In this example, the <b>testrdsdb</b> database created in <b>c</b> is used.
Username	Username used for accessing the database. This user must have the permissions to read and write data tables and metadata.
	In this example, the default user <b>root</b> for creating the RDS for MySQL DB instance is used.
Password	Password of the user.

 Table 5-14 Connection parameters

For other parameters, retain the default values. For details, see Link to Relational Databases. Click Save.

elect Connector		Configu
When you create a data	tabase link for the first time, upload the required driver on the Driver Management page or this page.	
* Name	source_rds	
* Connector	Relational Database v	
Database Type	MySQL ~	
* Database Server 🕐	Select	
Port (?)	3306	
* Database Name (	testrdadb	
* Username (?	root	
* Password ⑦	······································	
Use Local API	Yes No	
Use Agent (?)	Yes No	
	mysol-connector-iava-5.1.48.iar Upload I Copy from SETP	

Figure 5-34 Configuring the connection to the RDS for MySQL DB

b. Create a connection to DLI.

instanco

- i. Log in to the CDM console. In the navigation pane on the left, choose **Cluster Management**. On the displayed page, locate the created CDM cluster and click **Job Management** in the **Operation** column.
- ii. On the Job Management page, click the Links tab, and click Create Link. On the displayed page, select Data Lake Insight and click Next.

er Management / cdm	/ Links / Create Link				(
			-		
Data Warehouse	Data Warehouse Service	Data Lake Insight			
Hadoop	MRS HDFS	MRS HBase	MRS Hive	Apache HDFS	
	Apache HBase	Apache Hive			
Object Storage	Object Storage Service (OBS)	Alibaba Cloud OSS			
File System	FTP	SFTP	HTTP		
Relational Database	RDS for MySQL	RDS for PostgreSQL	RDS for SQL Server	MySQL	
	PostgreSQL	Microsoft SQL Server	Oracle	IBM Db2	
NoSQL	Redis	MongoDB			
Messaging System	Data Ingestion Service	MRS Kafka	Apache Kafka		
Search	Elasticsearch				
Open Beta Test	^				
× Cancel > Ne	xt				

Figure 5-35 Selecting the DLI connector

i. Create a connection to link CDM to DLI. For details about parameter settings, see Link to DLI.

Figure 5-36 Selecting the DLI connector

Cluster Management	/ cdm.4es4.60378328 / Links / Create Link	2 Configure
* Name	dest_dli	
* Connector	DLI v	
* AK 🕐		
* sк 🕜	·······	
* Project ID	05 f	
X Cancel	C Previous Of Test Save	

Click Save.

- 2. Create a CDM migration job.
  - a. Log in to the CDM console. In the navigation pane on the left, choose **Cluster Management**. On the displayed page, locate the created CDM cluster and click **Job Management** in the **Operation** column.
  - b. On the **Job Management** page, click the **Table/File Migration** tab. On the displayed tab, click **Create Job**.
  - c. On the **Create Job** page, set job parameters.

Figure 5-37 Configuring the migration job

Job Configuration				
* Job Name rds_to_dli				
Source Job Configuration	in	Destinatio	on Job Configuration	
* Source Link Name Sou	rce_rds	* Destination	Link Name dest_di	
Use SQL Statement (?)	Yes No	* Resource Q	ueue (?) test0402	Θ
* Schema/Table Space (?) tes	rdsdb 🖂	* Database (	() testdb	Θ
* Table Name (?) tab	etest 😔	* Table 🕐	tabletest	Θ
Show Advanced Attributes		Clear data b	efore import (?) Yes No	

- i. Job Name: Name of the data migration job, for example, rds\_to\_dli.
- ii. Set the parameters in the **Source Job Configuration** area as follows:

Table 5-15 Source job parameters

Parameter	Value
Source Link Name	Select the name of the data source created in <b>1.a</b> .
Use SQL Statement	When set to <b>Yes</b> , enter a SQL statement. CDM exports data based on the statement.
	In this example, set it to <b>No</b> .
Schema/Table Space	Select the name of the RDS for MySQL database you want to migrate to DLI. For example, the <b>testrdsdb</b> database.
Table Name	Name of the table you want to migrate. In this example, use <b>tabletest</b> created in <b>d</b> .

For details about parameter settings, see **From PostgreSQL/SQL Server**.

iii. Set the parameters in the **Destination Job Configuration** area as follows:

Parameter	Value
Destination Link Name	Select the DLI data source connection.
Resource Queue	Select a created DLI SQL queue.
Database Name	Select a created DLI database. In this example, the database <b>testdb</b> created in <b>Create a database and table on DLI</b> is selected.
Table Name	Select the name of a table in the database. In this example, the table <b>tabletest</b> created in <b>Create a database and table on DLI</b> is selected.
Clear data before import	Whether to clear data in the destination table before data import. In this example, set it to <b>No</b> . If set to <b>Yes</b> , data in the destination table will be cleared before the task is started.

Table 5-16         Destination	job	parameters
--------------------------------	-----	------------

For details about parameter settings, see **To DLI**.

- iv. Click **Next**. The **Map Field** page is displayed. CDM automatically matches the source and destination fields.
  - You can drag any unmatched fields to match them.
  - If the type is automatically created at the migration destination, you need to configure the type and name of each field.
  - CDM allows for field conversion during migration. For details, see **Field Conversion**.

#### Figure 5-38 Field mapping

()∞	nfgure Basic Information	e Baic Internation						(3) Configure Task			
	Source Field					•	,	Destination Field			
	Name	Example Velae	Туря	Operation				Name	Туре	Operation	
	16		VARCHAR(32)	8	Q	0	0)	- H	prints	ά.	
	name		WRCHAR(32)	8	Q	0	0)	rane	string	Ω.	
						•	1			0	<b>@</b> •
	× Cancel ( Previous	> Next Save									

v. Click **Next** and set task parameters. Typically, retain the default values for all parameters.

In this step, you can configure the following optional features:

- **Retry Upon Failure**: If the job fails to be executed, you can determine whether to automatically retry. Retain the default value **Never**.
- **Group**: Select the group to which the job belongs. The default group is **DEFAULT**. On the **Job Management** page, jobs can be displayed, started, or exported by group.
- **Scheduled Execution**: For how to configure scheduled execution, see **Scheduling Job Execution**. Retain the default value **No**.
- **Concurrent Extractors**: Enter the number of extractors to be concurrently executed. Retain the default value **1**.
- Write Dirty Data: Set this parameter if data that fails to be processed or filtered out during job execution needs to be written to OBS. Before writing dirty data, create an OBS link. You can view the data on OBS later. Retain the default value No, meaning dirty data is not recorded.
- vi. Click **Save and Run**. On the **Job Management** page, you can view the job execution progress and result.

Figure 5-39 Job progress and execution result

Table#Tile Migration Entre	DB Migration Links 7	Agents Settings								
🕀 Crivate Job 🌐 🖓 Run	2 Delete				G Feedback	3 Expert 2 im	port Schedule	<ul> <li>Al statutes</li> </ul>	• Job name • Job name or link type	QC
⊙≠⊳∓ <	Name JE	Link Details	Created By JE	Last Execution Time 48	Duration 48	Write Statistics	Status	Group Name	Operation	
Enter a group name. Q	inte_te_di	source_rds-dest_dli	el_dics_6603	Apr 02, 2022 15:45:32 GMT+08:00	1m 27s	Wintten rows: 3	Succeeded	DEFAULT	Run   Historical Record   Edt   Mare +	
Groups										
DEFAULT										

### **Step 3: Query Results**

Once the migration job is complete, check whether the RDS for MySQL table data has been migrated to the **tabletest** table. Specifically, do as follows: Log in to the DLI management console and choose **SQL Editor**. On the displayed page, set **Engine** to **Spark**, **Queues** to the created SQL queue, and **Databases** to the database created in **Create a database and table on DLI**. Then, execute the following query statement: select \* from tabletest;

Figure 5-40 Querying data in the table

Engine spark 👻 Queues test	▼ Databases testdb	Execute Format Refer Syntax Settings More
1create table tabletest(id string,name string);		
<ul> <li>Maker - Yrom Ebureett;</li> </ul>		
Line 1, Column 1		Execute: Ctrl+Enter, Find: Ctrl+F, Format: Shift+Alt+F, Verify Syntax: Ctrl+Q, Fullscreen: F11
xecuted Queries (Last Day) View Result		Clear All
Executed successfully		
Querycreate table tabletest(id string,name string); select * from tabletest		
Job ID 9616dbe4-9fd5-49fd-8d6f-3e10b771d318		
The query takes 27.07s, and 1.05 KB scanned.A maximum of 1,000 records can be displayed.		Enter a keyword. Q
id J≣	name J≣	
456	efg	
789	hij	
123	abc	

# 5.2.7 Example Typical Scenario: Migrating Data from GaussDB(DWS) to DLI

This section describes how to use CDM's data synchronization to migrate data from GaussDB(DWS) to DLI.

### Prerequisites

• You have created a DLI SQL queue. For how to create a DLI queue, see Creating a Queue.

### 

Set Type to For SQL when buying a queue.

- You have created a GaussDB(DWS) cluster. For how to create a GaussDB(DWS) cluster, see Creating a Cluster.
- You have created a CDM cluster. For how to create a CDM cluster, see Creating a CDM Cluster.

**NOTE** 

- To connect the cluster to an on-premises database as the destination data source, you can use either Internet or Direct Connect. If the Internet is used, make sure that an EIP has been bound to the CDM cluster, the security group of CDM allows outbound traffic from the host where the on-premises data source is located, the host where the data source is located can access the Internet, and the connection port has been enabled in the firewall rules.
- If the data source is GaussDB(DWS) or MRS, the network must meet the following requirements:

i. If the CDM cluster and the cloud service are in different regions, they must be connected through either the Internet or Direct Connect. If the Internet is used, make sure that an EIP has been bound to the CDM cluster, the host where the data source is located can access the Internet, and the port has been enabled in the firewall rules.

ii. If the CDM cluster and the cloud service are in the same region, instances in the same VPC, subnet, and security group can communicate with each other by default. If the CDM cluster and the cloud service are in the same VPC but in different subnets or security groups, you must configure routing and security group rules.

For how to configure routing rules, see **Configure routes**. For how to configure security group rules, see **Security Group Configuration Examples**.

iii. The cloud service instance and the CDM cluster belong to the same enterprise project. If they do not, change the enterprise project of the workspace.

In this example, the VPC, subnet, and security group of the CDM cluster match those of the GaussDB(DWS) cluster.

### Step 1: Prepare Data

- Create a database and table in the GaussDB(DWS) cluster.
  - a. Connect to the existing GaussDB(DWS) cluster by referring to Using the gsql CLI Client to Connect to a Cluster.

- b. Connect to the default database **gaussdb** of the GaussDB(DWS) cluster. gsql -d gaussdb -h *Connection address of the GaussDB(DWS) cluster* -U dbadmin -p 8000 -W *password* -r
  - gaussdb: Default database of the GaussDB(DWS) cluster.
  - Connection address of the DWS cluster: If a public address is used, set it to Public Network Address or Public Network Access Domain Name. If a private address is used, set it to Private Network Address or Private Network Access Domain Name. For details, see Obtaining the Cluster Connection Address. If an ELB is used, set it to the ELB address.
  - dbadmin: Default administrator username used during cluster creation.
  - -W: Default password of the administrator.
- c. Create the **testdwsdb** database. CREATE DATABASE testdwsdb;
- d. Exit the **gaussdb** database and connect to **testdwsdb**.

gsql -d testdwsdb -h *Connection address of the GaussDB(DWS) cluster* -U dbadmin -p 8000 -W *password* -r

e. Create a table and import data into it.

Create a table. CREATE TABLE table1(id int, a char(6), b varchar(6),c varchar(6));

Insert data into the table. INSERT INTO table1 VALUES(1,'123','456','789'); INSERT INTO table1 VALUES(2,'abc','efg','hif');

f. Query the table data to verify that the data is inserted. select \* from table1;

Figure 5-41 Querying data in the table

test	dwsdb=>	select '	* from	table1;
id	a	Ь	c	
	+	+	+	
1	123	456	789	
2	abc	efg	hif	
(2 ro	ows)			

- Create a database and table on DLI.
  - a. Log in to the DLI management console. In the navigation pane on the left, choose **SQL Editor**. On the displayed page, set **Engine** to **Spark** and **Queues** to the created SQL queue.

Create a database, for example, **testdb**. For the syntax to create a DLI database, see **Creating a Database**.

create database testdb;

b. On the **SQL Editor** page, set **Databases** to **testdb** and run the following table creation statement to create a table in the database. For the table creation syntax, see **Creating a DLI Table Using the DataSource Syntax**. create table tabletest(id INT, name1 string, name2 string, name3 string);

### Step 2: Migrate Data

- 1. Create a CDM connection.
  - a. Create a connection to the GaussDB(DWS) database.
    - i. Log in to the CDM console. In the navigation pane on the left, choose **Cluster Management**. On the displayed page, locate the created CDM cluster and click **Job Management** in the **Operation** column.
    - ii. On the Job Management page, click the Links tab, and click Create Link. On the displayed page, select Data Warehouse Service and click Next.
    - iii. Configure the connection as follows:

Parameter	Value
Name	Name of the GaussDB(DWS) data source, for example, <b>source_dws</b> .
Database Server	Click <b>Select</b> next to the text box and select the name of the created GaussDB(DWS) cluster.
Port	Port number of the GaussDB(DWS) database, which is <b>8000</b> by default.
Database Name	Name of the GaussDB(DWS) database you want to migrate. In this example, the <b>testdwsdb</b> database created in <b>Create a database and table</b> <b>in the GaussDB(DWS) cluster</b> is used.
Username	Username used for accessing the database. This user must have the permissions to read and write data tables and metadata.
	In this example, the default administrator <b>dbadmin</b> specified when you create the GaussDB(DWS) database is used.
Password	Password of the GaussDB(DWS) database user.

Table 5-17 GaussDB(DWS) data source configuration

Cluster Management / cdm-test-0	/ Links / Edit Link
* Name	source_dws
* Connector	Relational Database
Database Type	Data warehouse
* Database Server	dws-demog.dws.myhuaweiclouds Select
* Port ⑦	8000
* Database Name  ?	testdwsdb
* Username (?)	dbadmin
* Password (?)	····· &
Use Agent (?)	Yes No
Show Advanced Attributes	
X Cancel	🖹 Save

Figure 5-42 Configuring the GaussDB(DWS) connection

For other parameters, retain the default values. For details, see Link to Relational Databases. Click Save.

- b. Create a connection to the DLI.
  - i. Log in to the CDM console. In the navigation pane on the left, choose **Cluster Management**. On the displayed page, locate the created CDM cluster and click **Job Management** in the **Operation** column.
  - ii. On the Job Management page, click the Links tab, and click Create Link. On the displayed page, select Data Lake Insight and click Next.

Data Warehouse	Data Warehouse Service	Data Lake Insight	7		
Hadoop	MRS HDFS	MRS HBase	MRS Hive	Apache HDFS	
	Apache HBase	Apache Hive			
Object Storage	Object Storage Service (OBS)	Alibaba Cloud OSS			
File System	FTP	SFTP	HTTP		
Relational Database	RDS for MySQL	RDS for PostgreSQL	RDS for SQL Server	MySQL	
	PostgreSQL	Microsoft SQL Server	Oracle	IBM Db2	
NoSQL	Redis	MongoDB			
Messaging System	Data Ingestion Service	MRS Kafka	Apache Kafka		
Search	Elasticsearch				
Open Beta Test	^				

Figure 5-43 Selecting the DLI connector

i. Create a connection to link CDM to DLI. For details about parameter settings, see Link to DLI.

Figure 5-44 Selecting the DLI connector

uster Management	/ cdm-lest-00378328 / Linits / Create Link	Darter
Select Connector		- Contigure
* Name	dest_dii	
* Connector	DU *	
* AK 🕐		
* SK 🕜		
* Project ID	05/ f	
× Cancel	<     Previous     O <sup>2</sup> Test     Save	

Click Save.

- 2. Create a CDM migration job.
  - a. Log in to the CDM console. In the navigation pane on the left, choose **Cluster Management**. On the displayed page, locate the created CDM cluster and click **Job Management** in the **Operation** column.
  - b. On the **Job Management** page, click the **Table/File Migration** tab. On the displayed tab, click **Create Job**.
  - c. On the **Create Job** page, set job parameters.

Figure 5-45 Configuring the migration job

Configure Seals Information	D Map Trees (3) Contigue	t Tank
Job Configuration * Job Turne ****		
Source Job Configuration	Destination Job Configuration	
* Source Link Hame zource_dun	* Destination Link Hame dout_cili	
Use SQL Statement (2) Yes No	* Resource Gause (3) Set405	
* Schema/Table Space () politic ()	* Catalase () Saturb	
• Table Name (i) Table1	+ Table (7) Isolatest (9)	
Show Advanced Addulan	Clear data before import	
× Center → Inst G Sem		

- i. Job Name: Name of the data migration job, for example, test.
- ii. Set the parameters in the **Source Job Configuration** area as follows:

 Table 5-18
 Source job parameters

Parameter	Value
Source Link Name	Select the name of the data source created in <b>1.a</b> .
Use SQL Statement	When set to <b>Yes</b> , enter a SQL statement. CDM exports data based on the statement. In this example, set it to <b>No</b> .

Parameter	Value
Schema/Table Space	Name of the schema or tablespace from which data will be extracted. This parameter is available when <b>Use SQL Statement</b> is set to <b>No</b> . Click the icon next to the text box to select a schema or tablespace or directly enter a schema or tablespace.
	In this example, set this parameter to the default value <b>public</b> as there is no schema created in <b>Create a database and table in the GaussDB(DWS) cluster</b> .
	If there are no schemas or tablespaces available, check if the account has the permission to query metadata.
	<b>NOTE</b> The parameter value can contain wildcard characters (*), which allows for the export of all databases with names starting or ending with a certain prefix or suffix, respectively. For example:
	SCHEMA* indicates that all databases with names starting with SCHEMA are exported.
	<b>*SCHEMA</b> indicates that all databases with names ending with SCHEMA are exported.
	<b>*SCHEMA*</b> indicates that all databases with names containing <b>SCHEMA</b> are exported.
Table Name	Name of the table you want to migrate. In this example, <b>table1</b> created in <b>Create a database</b> and table in the GaussDB(DWS) cluster is used.

For details about parameter settings, see **From a Relational Database**.

iii. Set the parameters in the **Destination Job Configuration** area as follows:

Parameter	Value
Destination Link Name	Select the DLI data source connection.
Resource Queue	Select a created DLI SQL queue.
Database Name	Select a created DLI database. In this example, the database <b>testdb</b> created in <b>Create a database and table on DLI</b> is used.

Table	5-19	Destination	iob	parameters
iable	5 .5	Destination	100	parameters

Parameter	Value
Table Name	Select the name of a table in the database. In this example, the table <b>tabletest</b> created in <b>Create a database and table on DLI</b> is used.
Clear data before import	Whether to clear data in the destination table before data import. In this example, set it to <b>No</b> . If set to <b>Yes</b> , data in the destination table will be cleared before the task is started

For details about parameter settings, see **To DLI**.

- iv. Click **Next**. The **Map Field** page is displayed. CDM automatically matches the source and destination fields.
  - You can drag any unmatched fields to match them.
  - If the type is automatically created at the migration destination, you need to configure the type and name of each field.
  - CDM allows for field conversion during migration. For details, see Field Conversion.

### Figure 5-46 Field mapping

(1) Configure Basic Information	n Bans Information 🕘 Vice Print								
Source Field					0/	Destination Field			<b>⊽ e c</b>
Name	Example Value	Type	Operation			Name	Type	Operation	
		ыт	8	Q	<b>п</b>	- 4	н.	α.	
1 C C		CH4R(E)	8	Q	₫ ₀	-i nerel	phing	α	
ь.		WACHING	8	Q	<b>п</b>	name2	string	α	
e		WRCHAR(6)	8	Q	a	-i name3	string	α	
× Cancel C Previous	Ned Stee								

v. Click **Next** and set task parameters. Typically, retain the default values for all parameters.

In this step, you can configure the following optional features:

- **Retry Upon Failure**: If the job fails to be executed, you can determine whether to automatically retry. Retain the default value **Never**.
- **Group**: Select the group to which the job belongs. The default group is **DEFAULT**. On the **Job Management** page, jobs can be displayed, started, or exported by group.
- Scheduled Execution: For how to configure scheduled execution, see Scheduling Job Execution. Retain the default value No.
- **Concurrent Extractors**: Enter the number of extractors to be concurrently executed. Retain the default value **1**.
- Write Dirty Data: Set this parameter if data that fails to be processed or filtered out during job execution needs to be written to OBS. Before writing dirty data, create an OBS link. You can view the data on OBS later. Retain the default value No, meaning dirty data is not recorded.

vi. Click **Save and Run**. On the **Job Management** page, you can view the job execution progress and result.

Figure 5-47 Job progress and execution result

Table File Migration	Entire DB Migration Links	Agents Settings								
() Create Job	di Run				C Feedback	toot 🚺 Import	Schedule	• Al statutes	• Job name • Job name or link type	Q C
0 / Þ T	C Name 42	Link Details	Created By 40	Last Execution Time 48	Duration 42	Write Statistics	Status	Group Name	Operation	
Enter a group name.	Q 🗆 🚥	source_dvs-dest_di	00378325	Apr 66, 2022 17:34:50 GMT+08:00	104	Weller rove 2	Succeeded	DEFAULT	Run   Historical Record   Edit   More +	
Grosps										
DEFAULT										

### **Step 3: Query Results**

Once the migration job is complete, check whether the GaussDB(DWS) table data has been migrated to the **tabletest** table. Specifically, do as follows: Log in to the DLI management console and choose **SQL Editor**. On the displayed page, set **Engine** to **Spark**, **Queues** to the created SQL queue, and **Databases** to the database created in **Create a database and table on DLI**. Then, execute the following query statement: select \* from tabletest;

#### Figure 5-48 Querying data in the table

	Engine spark + Q	ueues testdi +	Databases	testdb 👻	Execute	Format	Refer Syntax	Settings	More +
Increase database testion; Increase database testion; Increase table testistes(16 bit, same string, name string); I select * from tabletest	Engine spark • Q	ucues testil +	Databases	testdb •	© Execute	Format	Refer Syntax	Settings	More +
Line 1, Colume 1 Executed Queries (Last Day) Vew Result Result 0					Execute: Ctrl+Ente	t, Find: Clil+F, Fo	mat Shift+All+F, Verify	Syntax Ctrl+Q, F	Fullscreen: F11 Clear Al
Executed successfully Ourrycreate database testiditycreate table tabletestijd INT, name1 string, name2 s     Do ID fods48c6-4211-401c-a2214-1087be50be50b5  The name table stract wide 196 KR screened & maximum of 1.000 records can be disclosed	tring, name3 string); sel						Enter a lensered	0	
id.j≘ name1.j≣		name2 JE			nam	ic8 (≘			
2 abc		efg			hif				
1 123		456			789				

# **6** Configuring DLI to Read and Write Data from and to External Data Sources

# 6.1 Configuring DLI to Read and Write External Data Sources

To read and write external data sources when running DLI jobs, two conditions must be met:

- Establish network connectivity between DLI and the external data source to ensure that the DLI queue is connected to the data source network.
- Securely store the access credentials for the data source to ensure authentication security and facilitate secure DLI access to the data source.

This section describes how to configure DLI to read and write external data sources.

- Configure network connection between DLI and the data source by referring to Configuring the Network Connection Between DLI and Data Sources (Enhanced Datasource Connection).
- Manage credentials for DLI to access data sources.
  - Spark 3.3.1 or later and Flink 1.15 or later jobs accessing data sources using datasource connections
    - You are advised to use Data Encryption Workshop (DEW) to store authentication information of data sources, addressing data security, key security, and complex key management issues.

For details, see Using DEW to Manage Access Credentials for Data Sources.

- To manage data source access credentials using DEW, you also need to create a DLI agency to grant DLI access to read access credentials for other services (DEW).
- When SQL and Flink 1.12 jobs access data sources using datasource connections, use DLI's datasource authentication feature to manage data

source access credentials. For details, see Using DLI Datasource Authentication to Manage Access Credentials for Data Sources.

# 6.2 Configuring the Network Connection Between DLI and Data Sources (Enhanced Datasource Connection)

### **6.2.1 Overview of Enhanced Datasource Connections**

### Why Create Enhanced Datasource Connections?

In cross-source data analysis scenarios, DLI needs to connect to external data sources. However, due to the different VPCs between the data source and DLI, the network cannot be connected, which results in DLI being unable to read data from the data source. DLI's enhanced datasource connection feature enables network connectivity between DLI and the data source.

This section will introduce a solution for cross-VPC data source network connectivity:

- Creating an enhanced datasource connection: Establish a VPC peering connection to connect DLI and the data source's VPC network.
- Testing network connectivity: Verify the connectivity between the queue and the data source's network.

For details about the data sources that support cross-source access, see **Common Development Methods for DLI Cross-Source Analysis**.

### 

In cross-source development scenarios, there is a risk of password leakage if datasource authentication information is directly configured. You are advised to use Data Encryption Workshop (DEW) to store authentication information of data sources when Spark 3.3.1 or later and Flink 1.15 or later jobs access data sources using datasource connections. This will help you address issues related to data security, key security, and complex key management. For details, see Using DEW to Manage Access Credentials for Data Sources.

### **Notes and Constraints**

ltem	Description
Use case	<ul> <li>Datasource connections cannot be created for the default queue.</li> </ul>
	• Flink jobs can directly access DIS, OBS, and SMN data sources without using datasource connections.
	• When compute resources in non-elastic resource pools are used, enhanced datasource connections can only be created for yearly/monthly and pay-per-use dedicated queues.
Permission	• VPC Administrator permissions are required for enhanced connections to use VPCs, subnets, routes, VPC peering connections. You can set these permissions by referring to Service Authorization.
Usage	• If you use an enhanced datasource connection, the CIDR block of the elastic resource pool or queue cannot overlap with that of the data source.
	• Only queues bound with datasource connections can access datasource tables.
	• Datasource tables do not support the preview function.

 Table 6-1 Notes and constraints on enhanced datasource connections

ltem	Description
Connectivity check	• When checking the connectivity of datasource connections, the notes and constraints on IP addresses are:
	<ul> <li>The IP address must be valid, which consists of four decimal numbers separated by periods (.). The value ranges from 0 to 255.</li> </ul>
	<ul> <li>During the test, you can add a port after the IP address and separate them with colons (:). The port can contain a maximum of five digits. The value ranges from 0 to 65535.</li> <li>For example, <b>192.168.</b>xx.xx or <b>192.168.</b>xx.xx.<b>8181</b>.</li> </ul>
	• When checking the connectivity of datasource connections, the notes and constraints on domain names are:
	<ul> <li>The domain name can contain 1 to 255 characters.</li> <li>Only letters, numbers, underscores (_), and hyphens</li> <li>(-) are allowed.</li> </ul>
	<ul> <li>The top-level domain name must contain at least two letters, for example, .com, .net, and .cn.</li> </ul>
	<ul> <li>During the test, you can add a port after the domain name and separate them with colons (:). The port can contain a maximum of five digits. The value ranges from 0 to 65535.</li> <li>Example: example.com:8080</li> </ul>

### **Cross-Source Analysis Process**

To use DLI for cross-source analysis, you need to create a datasource connection to connect DLI to the data source, and then develop jobs to access the data source.

Figure 6-1 Cross-source analysis flowchart



### Helpful Links

• Creating an enhanced datasource connection on the management console

Creating an Enhanced Datasource Connection

- Creating an enhanced datasource connection using APIs Creating an Enhanced Datasource Connection
- Practice of creating an enhanced datasource connection

- Example Typical Scenario: Connecting DLI to a Data Source on a Private Network
- Example Typical Scenario: Connecting DLI to a Data Source on a Public Network

### 6.2.2 Creating an Enhanced Datasource Connection

### Scenario

Create an enhanced datasource connection for DLI to access, import, query, and analyze data of other data sources.

For example, to connect DLI to the MRS, RDS, CSS, Kafka, or GaussDB(DWS) data source, you need to enable the network between DLI and the VPC of the data source.

Create an enhanced datasource connection on the console.

### Notes and Constraints

able 6-2 Notes and	constraints on	enhanced da	atasource co	nnections
--------------------	----------------	-------------	--------------	-----------

ltem	Description
Use case	<ul> <li>Datasource connections cannot be created for the default queue.</li> </ul>
	<ul> <li>Flink jobs can directly access DIS, OBS, and SMN data sources without using datasource connections.</li> </ul>
	<ul> <li>When compute resources in non-elastic resource pools are used, enhanced datasource connections can only be created for yearly/monthly and pay-per-use dedicated queues.</li> </ul>
Permission	• VPC Administrator permissions are required for enhanced connections to use VPCs, subnets, routes, VPC peering connections. You can set these permissions by referring to Service Authorization.
Usage	<ul> <li>If you use an enhanced datasource connection, the CIDR block of the elastic resource pool or queue cannot overlap with that of the data source.</li> </ul>
	<ul> <li>Only queues bound with datasource connections can access datasource tables.</li> </ul>
	• Datasource tables do not support the preview function.

ltem	Description
Connectivity check	• When checking the connectivity of datasource connections, the notes and constraints on IP addresses are:
	<ul> <li>The IP address must be valid, which consists of four decimal numbers separated by periods (.). The value ranges from 0 to 255.</li> </ul>
	<ul> <li>During the test, you can add a port after the IP address and separate them with colons (:). The port can contain a maximum of five digits. The value ranges from 0 to 65535.</li> <li>For example, <b>192.168</b>.<i>xx</i>.<i>xx</i> or <b>192.168</b>.<i>xx</i>.<i>xx</i>.<b>8181</b>.</li> </ul>
	• When checking the connectivity of datasource connections, the notes and constraints on domain names are:
	<ul> <li>The domain name can contain 1 to 255 characters.</li> <li>Only letters, numbers, underscores (_), and hyphens</li> <li>(-) are allowed.</li> </ul>
	<ul> <li>The top-level domain name must contain at least two letters, for example, .com, .net, and .cn.</li> </ul>
	<ul> <li>During the test, you can add a port after the domain name and separate them with colons (:). The port can contain a maximum of five digits. The value ranges from 0 to 65535.</li> <li>Example: example.com:8080</li> </ul>

### Process

### Figure 6-2 Enhanced datasource connection creation flowchart



### Prerequisites

- An elastic resource pool or queue has been created.
- You have obtained the VPC, subnet, private IP address, port, and security group information of the external data source.
- The security group of the external data source has allowed access from the CIDR block of the elastic resource pool or queue.

### Procedure

### Step 1 Create an Enhanced Datasource Connection

1. Log in to the DLI management console.

- 2. In the navigation pane on the left, choose **Datasource Connections**.
- On the displayed Enhanced tab, click Create.
   Set parameters based on Table 6-3.

### Table 6-3 Parameters

Parameter	Description
Connection Name	<ul> <li>Name of the created datasource connection.</li> <li>Only letters, numbers, and underscores (_) are allowed. The parameter must be specified.</li> <li>A maximum of 64 characters are allowed.</li> </ul>
Resource Pool	This parameter is optional when you create an enhanced datasource connection. However, you must bind an elastic resource pool to the enhanced datasource connection before using it. The status of the enhanced datasource connection's VPC peering connection is <b>active</b> . Used to bind an elastic resource pool or queue that uses a datasource connection. If you use compute resources in non-elastic resource pools, DLI will create a resource pool with the same name for your yearly/monthly or pay-per-use dedicated queue (only yearly/monthly or pay-per-use dedicated queues can be bound) after the elastic resource pool function is enabled. Here, you can select the corresponding resource pool and bind it to an enhanced datasource connection. <b>NOTE</b> Before using an enhanced datasource connection, ensure that the created VPC peering connection is in the <b>Active</b> state.
VPC	VPC used by the data source.
Subnet	Subnet used by the data source. If the selected subnet has IPv6 enabled, the enhanced datasource connection you create will also support IPv6, and you can subsequently add routes for IPv6 addresses.
Route Table	<ul> <li>Route table of the subnet.</li> <li>NOTE <ul> <li>The route table is associated with the subnet used by the destination data source, which is not the table containing the route you add by Manage Route in the Operation column. The route you add on the Manage Route page is contained in the route table associated with the subnet used by the queue to be bound.</li> <li>The subnet used by the destination data source must be different from that used by the queue to be bound. Otherwise, a segment conflict occurs.</li> </ul> </li> </ul>

Parameter	Description
Host Information	In this text field, you can configure the mapping between host IP addresses and domain names so that jobs can only use the configured domain names to access corresponding hosts. This parameter is optional.
	For example, when accessing the HBase cluster of MRS, you need to configure the host name (domain name) and IP address of the ZooKeeper instance. Enter one record in each line in the format of <i>IP address Host name Domain name</i> .
	Example:
	192.168.0.22 node-masterxxx1.com
	192.168.0.23 node-masterxxx2.com
	For details about how to obtain host information, see <b>How</b> <b>Do I Obtain MRS Host Information?</b> .
Tags	Tags used to identify cloud resources. A tag includes the tag key and tag value. If you want to use the same tag to identify multiple cloud resources, that is, to select the same tag from the drop-down list box for all services, you are advised to create predefined tags on the Tag Management Service (TMS).
	If your organization has configured tag policies for DLI, add tags to resources based on the policies. If a tag does not comply with the tag policies, resource creation may fail. Contact your organization administrator to learn more about tag policies.
	For details, see <b>Tag Management Service User Guide</b> .
	NOTE
	<ul> <li>A maximum of 20 tags can be added.</li> </ul>
	<ul> <li>Only one tag value can be added to a tag key.</li> </ul>
	<ul> <li>The key name in each resource must be unique.</li> </ul>
	– Tag key: Enter a tag key name in the text box.
	<b>NOTE</b> A tag key can contain a maximum of 128 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed, but the value cannot start or end with a space or start with <b>_sys_</b> .
	- Tag value: Enter a tag value in the text box.
	NOTE A tag value can contain a maximum of 255 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed.

4. Click OK.

After the creation is complete, the enhanced datasource connection is in the **Active** state, indicating that the connection is successfully created.

#### Step 2 Security Group Where the Data Source Belongs Allows Access from the CIDR Block of the Elastic Resource Pool

1. On the DLI management console, obtain the network segment of the elastic resource pool or queue.

Choose **Resources** > **Queue Management** from the left navigation pane. On the page displayed, locate the queue on which jobs are running, and click the button next to the queue name to obtain the CIDR block of the queue.

- 2. Log in to the VPC console and find the VPC the data source belongs to.
- 3. On the network console, choose Virtual Private Cloud > Network Interfaces. On the Network Interfaces tab page displayed, search for the security group name, click More in the Operation column, and select Change Security Group.
- 4. In the navigation pane on the left, choose **Access Control** > **Security Groups**.
- 5. Click the name of the security group to which the external data source belongs.
- 6. Click the **Inbound Rules** tab and add a rule to allow access from the CIDR block of the queue. See **Figure 6-3**.

Configure the inbound rule parameters according to **Table 6-4**.

#### Figure 6-3 Adding an inbound rule

Add Inbound	d Rule Learn	more about security group co	nfiguration.			
Some secur	rity group rules will no	t take effect for ECSs with certain spe	cifications. Learn mor	e		
Security Group s	sg-705b tiple rules in a batch.					
Priority (?)	Action (?)	Protocol & Port (?)	Туре	Source (?)	Description	Operation
1-100	Allow •	Protocols/TCP (Custo   Example: 22 or 22,24 or 22-3	IPv4 •	IP address 0.0.0.0/0	•	Replicate Delete
			+ Add Rule	cel		

 Table 6-4 Inbound rule parameters

Parameter	Description	Example Value
Priority	Priority of a security group rule.	1
	The priority value ranges from 1 to 100. The default value is <b>1</b> , indicating the highest priority. A smaller value indicates a higher priority of a security group rule.	
Action	Action of the security group rule.	Allow

Parameter	Description	Example Value
Protocol & Port	<ul> <li>Network protocol. The value can be All, TCP, UDP, ICMP, or GRE.</li> <li>Port: Port or port</li> </ul>	In this example, select <b>TCP</b> . Leave the port blank or set it to the data source port.
	range over which the traffic can reach your instance. The port ranges from 1 to 65535.	
Туре	Type of IP addresses.	IPv4
Source	Allows access from IP addresses or instances in another security group.	In this example, enter the obtained queue CIDR block.
Description	Supplementary information about the security group rule. This parameter is optional.	_

### Step 3 Test the Connectivity Between the DLI Queue and the Data Source

- 1. Obtain the private IP address and port number of the data source.
  - Take the RDS data source as an example. On the **Instances** page, click the target DB instance. On the page displayed, locate the **Connection Information** pane and view the private IP address. In the **Connection Information** pane, locate the **Database Port** to view the port number of the RDS DB instance.
- 2. In the navigation pane of the DLI management console, choose **Resources** > **Queue Management**.
- 3. Locate the queue bound with the enhanced datasource connection, click **More** in the **Operation** column, and select **Test Address Connectivity**.
- 4. Enter the data source connection address and port number to test the network connectivity.

Format: IP address.Port number

### 

Before testing the connection, ensure that the security group of the external data source has allowed access from the CIDR block of the queue.

**Figure 6-4** Testing the network connectivity between the queue and the data source

Test Add	Test Address Connectivity			
Tests whether domain name,	an address an IP addr	is reachable from ess, or a specified	a specified cluster. The address can be a port.	
★ Address	192.	57: 3		
		Test	Cancel	
End				

### **Related Operations**

• Why Is a Datasource Connection Successfully Created But the Network Connectivity Test Fails?

## 6.2.3 Establishing a Network Connection Between DLI and Resources in a Shared VPC

### VPC Sharing Overview

VPC sharing allows sharing VPC resources created in one account with other accounts using Resource Access Manager (RAM). For example, account A can share its VPC and subnets with account B. After accepting the share, account B can view the shared VPC and subnets and use them to create resources.

For more information about VPC sharing, see **VPC Sharing** in *Virtual Private Cloud User Guide*.

### DLI Use Cases

An enterprise IT management account creates a VPC and subnets and shares them with other service accounts to facilitate centralized configuration of VPC security policies and orderly resource management.

Service accounts use the shared VPC and subnets to create resources and want to use DLI to submit jobs and access resources in the shared VPC. To do this, they need to establish a network connection between DLI and the resources in the shared VPC.

For example, account A is the enterprise IT management account and the owner of VPC resources. It creates the VPC and subnets and shares them with service account B.

Account B is a service account that uses the shared VPC and subnets to create resources and uses DLI to access them.

### Prerequisites

- Account A has been configured with a DLI agency, which includes the DLI Datasource Connections Agency Access permission. This will grant the necessary permissions to access and use VPCs, subnets, routes, and VPC peering connections. For details, see Configuring DLI Agency Permissions.
- Account A, as the resource owner, has created a VPC and subnets and designated account B as the principal.

For details, see **Creating a Resource Share**.

### Establishing a Network Connection Between DLI and Resources in a Shared VPC

**Step 1** Account A creates an enhanced datasource connection.

- 1. Log in to the DLI management console using account A.
- 2. In the navigation pane on the left, choose **Datasource Connections**.
- 3. On the displayed **Enhanced** tab, click **Create**.

Set parameters based on **Table 6-5**.

 Table 6-5 Parameters for creating an enhanced datasource connection

Parameter	Description
Connection Name	Name of the datasource connection to be created
Resource Pool	You do not need to set this parameter in this scenario.
VPC	VPC shared by account A to account B
Subnet	Subnet shared by account A to account B
Route Table	You do not need to set this parameter in this scenario.
Host Information	You do not need to set this parameter in this scenario.
Tags	Tags used to identify cloud resources. A tag includes the tag key and tag value.

- 4. Click **OK**.
- **Step 2** Account A grants account B access to the enhanced datasource connection created in **Step 1**.
  - 1. In the enhanced datasource connection list, locate the row containing the newly created one, click **More** in the **Operation** column, and select **Manage Permission** from the drop-down list.
  - 2. In the displayed **Permissions** dialog box, select **Grant Permission** for **Set Permission**, enter the ID of the project account B belongs to in **Project ID**, and click **OK**.
- **Step 3** Account B binds a DLI elastic resource pool to the shared enhanced datasource connection.

- 1. Log in to the DLI management console using account B.
- 2. In the navigation pane on the left, choose **Datasource Connections**.
- 3. On the displayed **Enhanced** tab, locate the row containing the enhanced datasource connection shared by account A, click **More** in the **Operation** column, and select **Bind Resource Pool** from the drop-down list.
- 4. In the displayed **Bind Resource Pool** dialog box, select the created elastic resource pool for **Resource Pool** and click **OK**.

If there is no elastic resource pool available, create one by referring to **Creating an Elastic Resource Pool and Creating Queues Within It**.

**Step 4** Account B tests the network connectivity between the elastic resource pool and resources in the VPC.

#### **NOTE**

If there are resources in the shared VPC, ensure that the security group the resources belong to has allowed access to the elastic resource pool's CIDR block.

1. Obtain the private IP address and port number of the data source in the shared VPC.

Take the RDS data source as an example. On the **Instances** page, click the target DB instance. On the displayed page, locate the **Connection Information** pane and view the private IP address. In the **Connection Information** pane, locate the **Database Port** to view the port number of the RDS DB instance.

- 2. In the navigation pane of the DLI management console, choose **Resources** > **Queue Management**.
- 3. Locate the queue under the elastic resource pool bound with the enhanced datasource connection, click **More** in the **Operation** column, and select **Test Address Connectivity**.
- 4. Enter the data source connection address and port number to test the network connectivity.

If the address is reachable, it means that account B has established a network connection between the DLI resource and the resources in the shared VPC. Account B can then submit jobs to the elastic resource pool's queue and access the resources in the shared VPC.

----End

### 6.2.4 Common Development Methods for DLI Cross-Source Analysis

### **Cross-Source Analysis**

If DLI needs to access external data sources, you need to establish enhanced datasource connections to enable the network between DLI and the data sources, and then develop different types of jobs to access the data sources. This is the process of DLI cross-source analysis.

This section describes how to develop data sources supported by DLI for crosssource analysis.

### Notes

- Flink jobs can directly access DIS, OBS, and SMN data sources without using datasource connections.
- You are advised to use enhanced datasource connections to connect DLI to data sources.

### **Cross-Source Analysis Development Methods**

**Table 6-6** lists the data sources supported by DLI and the corresponding development methods.

Service	Spark SQL Job	Spark Jar Job	Flink OpenSource SQL Job	Flink Jar Job
CloudTable HBase	<ul> <li>Create an HBase association table</li> <li>Insert data</li> <li>Query data</li> </ul>	<ul> <li>Scala sample code</li> <li>PySpark sample code</li> <li>Java sample code</li> </ul>	<ul> <li>HBase source table</li> <li>HBase result table</li> <li>HBase dimension table</li> </ul>	-
CloudTable OpenTSDB	<ul> <li>Create an OpenTSDB association table</li> <li>Insert data</li> <li>Query data</li> </ul>	<ul> <li>Scala sample code</li> <li>PySpark sample code</li> <li>Java sample code</li> </ul>	-	-
CSS	<ul> <li>Create a CSS association table</li> <li>Insert data</li> <li>Query data</li> </ul>	<ul> <li>Scala sample code</li> <li>PySpark sample code</li> <li>Java sample code</li> </ul>	• Elasticsearch result table	-

 Table 6-6 Syntax reference for cross-source analysis

Service	Spark SQL Job	Spark Jar Job	Flink OpenSource SQL Job	Flink Jar Job
DCS Redis	<ul> <li>Create a DCS association table</li> <li>Insert data</li> <li>Query data</li> </ul>	<ul> <li>Scala sample code</li> <li>PySpark sample code</li> <li>Java sample code</li> </ul>	<ul> <li>Redis source table</li> <li>Redis result table</li> <li>Redis dimension table</li> </ul>	Flink job sample
DDS	<ul> <li>Create a DDS association table</li> <li>Insert data</li> <li>Query data</li> </ul>	<ul> <li>Scala sample code</li> <li>PySpark sample code</li> <li>Java sample code</li> </ul>	-	-
DMS	-	-	<ul> <li>Kafka source table</li> <li>Kafka result table</li> </ul>	-
GaussDB(DW S)	<ul> <li>Create a GaussDB(D WS) association table</li> <li>Insert data</li> <li>Query data</li> </ul>	<ul> <li>Scala sample code</li> <li>PySpark sample code</li> <li>Java sample code</li> </ul>	<ul> <li>GaussDB(D WS) source table</li> <li>GaussDB(D WS) result table</li> <li>GaussDB(D WS) dimension table</li> </ul>	Flink job sample
MRS HBase	<ul> <li>Create an HBase association table</li> <li>Insert data</li> <li>Query data</li> </ul>	<ul> <li>Scala sample code</li> <li>PySpark sample code</li> <li>Java sample code</li> </ul>	<ul> <li>HBase source table</li> <li>HBase result table</li> <li>HBase dimension table</li> </ul>	Flink job sample

Service	Spark SQL Job	Spark Jar Job	Flink OpenSource SQL Job	Flink Jar Job
MRS Kafka	-	-	<ul> <li>Kafka source table</li> <li>Kafka result table</li> </ul>	<ul> <li>Flink job sample</li> </ul>
MRS OpenTSDB	<ul> <li>Create an OpenTSDB association table</li> <li>Insert data</li> <li>Query data</li> </ul>	<ul> <li>Scala sample code</li> <li>PySpark sample code</li> <li>Java sample code</li> </ul>	-	_
RDS for MySQL	<ul> <li>Create an RDS association table</li> <li>Insert data</li> <li>Query data</li> </ul>	<ul> <li>Scala sample code</li> <li>PySpark sample code</li> <li>Java sample code</li> </ul>	• MySQL CDC source table	_
RDS PostGre	<ul> <li>Create an RDS association table</li> <li>Insert data</li> <li>Query data</li> </ul>	<ul> <li>Scala sample code</li> <li>PySpark sample code</li> <li>Java sample code</li> </ul>	Postgres CDC source table	-

# 6.3 Using DEW to Manage Access Credentials for Data Sources

When using DLI to submit jobs that involve reading and writing data from external sources, it is crucial to securely access these sources by properly storing their access credentials. This ensures the authentication of the data source and enables secure access by DLI. DEW is a comprehensive cloud-based encryption service that addresses data security, key security, and complex key management issues. This section describes how to use DEW to store authentication information for a data source.

For details, see Data Encryption Workshop (DEW).

### Creating a Shared Secret in DEW

This example describes how to configure a credential for accessing RDS DB instances in a DLI job and store the credential in DEW.

- 1. Log in to the DEW management console.
- 2. In the navigation pane on the left, choose **Cloud Secret Management Service** > **Secrets**.
- 3. Click **Create Secret**. On the displayed page, configure basic secret information.
  - Secret Name: Enter a secret name. In this example, the name is secretInfo.
  - Secret Value: Enter the username and password for logging in to the RDS for MySQL DB instance.
    - The key in the first line is **MySQLUsername**, and the value is the username for logging in to the DB instance.
    - The key in the second line is MySQLPassword, and the value is the password for logging in to the DB instance.

#### Figure 6-5 Secret Value

Secret key/value	Plaintext	
MySQLUsername		Delet
MvSQLPassword		Delet

4. Set other parameters as required and click **OK**.

### Using the Secret Created in DEW in a DLI Job

This part uses a Flink job as an example to describe how to use credentials created in DEW.

```
WITH (
'connector' = 'jdbc',
'url? = 'jdbc:mysql://MySQLAddress:MySQLPort/flink',--flink is the MySQL database where the orders table
locates.
'table-name' = 'orders',
'username' = 'MySQLUsername', -- Shared secret in DEW whose name is secretInfo and version is v1. The
key MySQLUsername defines the secret value. The value is the user's sensitive information.
'password' = 'MySQLPassword, -- Shared secret in DEW whose name is secretInfo and version is v1. The
key MySQLPassword defines the secret value. The value is the user's sensitive information.
'password' = 'MySQLPassword, -- Shared secret in DEW whose name is secretInfo and version is v1. The
key MySQLPassword defines the secret value. The value is the user's sensitive information.
'sink.buffer-flush.max-rows' = '1',
'dew.endpoint'='kms.xxxx.com', -- Endpoint information for the DEW service being used
'dew.csms.secretName'='secretInfo', --Name of the DEW shared secret
```

'dew.csms.decrypt.fields'='username,password', --The **password** field value must be decrypted and replaced using DEW secret management. 'dew.csms.version'='v1' );

### **Related Operations**

For how to use a DLI agency to obtain access credentials, see **Table 6-7**.

Туре	Helpful Link	Description
Flink job	Flink OpenSource SQL Jobs Using DEW to Manage Access Credentials	Guideline for using DEW to manage and access credentials for Flink OpenSource SQL jobs. When writing the output data of Flink jobs to MySQL or GaussDB(DWS), set attributes such as the username and password in the connector.
	Flink Jar Jobs Using DEW to Acquire Access Credentials for Reading and Writing Data from and to OBS	Guideline for Flink Jar jobs to acquire an AK/SK to read and write data from and to OBS.
	Obtaining Temporary Credentials for Flink Job Agencies	DLI provides a common interface to obtain temporary credentials for Flink job agencies set by users during job launch. The interface encapsulates the obtained temporary credentials for the job agency in the <b>com.huaweicloud.sdk.core.auth.</b> <b>BasicCredentials</b> class. Guideline for obtaining a temporary credential for a Flink job agency.
Spark job	Spark Jar Jobs Using DEW to Acquire Access Credentials for Reading and Writing Data from and to OBS	Guideline for Spark Jar jobs to acquire an AK/SK to read and write data from and to OBS.
	Obtaining Temporary Credentials for Spark Job Agencies	Guideline for obtaining a temporary credential for a Spark Jar job agency.

Table 6-7 Guidelines for configuring DLI agency permissions in specific scenarios

# 6.4 Using DLI Datasource Authentication to Manage Access Credentials for Data Sources

### 6.4.1 Overview

### What Is Datasource Authentication?

When analyzing across multiple sources, you are advised not to configure authentication information directly in a job as it can lead to password leakage. Instead, you are advised to use either Data Encryption Workshop (DEW) or datasource authentication provided by DLI to securely store data source authentication information.

• DEW is a comprehensive cloud-based encryption service that addresses data security, key security, and complex key management issues. You are advised to use DEW to store authentication information for data sources.

You are advised to use DEW to store authentication information of data sources when Spark 3.3.1 or later and Flink 1.15 or later jobs access data sources using datasource connections. This will help you address issues related to data security, key security, and complex key management. For details, see **Using DEW to Manage Access Credentials for Data Sources**.

 Datasource authentication is used to manage authentication information for accessing specified data sources. After datasource authentication is configured, you do not need to repeatedly configure data source authentication information in jobs, improving data source authentication security while enabling DLI to securely access data sources.

When SQL and Flink 1.12 jobs access data sources using datasource connections, use DLI's datasource authentication feature to manage data source access credentials.

This section describes how to use datasource authentication provided by DLI.

### Notes and Constraints

ltem	Description
Use case	<ul> <li>Only Spark SQL and Flink OpenSource SQL 1.12 jobs support datasource authentication.</li> </ul>
	<ul> <li>Flink jobs can use datasource authentication only on queues created after May 1, 2023.</li> </ul>

Table 6-8 Notes and constraints on datasource authentication

Item	Description
Datasource authentication type	<ul> <li>DLI supports four types of datasource authentication.</li> <li>Select an authentication type specific to each data source.</li> </ul>
	<ul> <li>CSS: applies to 6.5.4 or later CSS clusters with the security mode enabled.</li> </ul>
	<ul> <li>Kerberos: applies to MRS security clusters with Kerberos authentication enabled.</li> </ul>
	<ul> <li>Kafka_SSL: applies to Kafka with SSL enabled.</li> </ul>
	<ul> <li>Password: applies to GaussDB(DWS), RDS, DDS, and DCS.</li> </ul>

### **Datasource Authentication Types**

DLI supports four types of datasource authentication. Select an authentication type specific to each data source.

- CSS: applies to 6.5.4 or later CSS clusters with the security mode enabled. During the configuration, you need to specify the username, password, and authentication certificate of the cluster and store the information in DLI through datasource authentication so that DLI can securely access CSS data sources. For details, see Creating a CSS Datasource Authentication.
- Kerberos: applies to MRS security clusters with Kerberos authentication enabled. During the configuration, you need to specify MRS cluster authentication credentials, including the krb5.conf and user.keytab files. For details, see Creating a Kerberos Datasource Authentication.
- Kafka\_SSL: applies to Kafka with SSL enabled. During the configuration, you need to specify the KafkaTruststore path and password. For details, see **Creating a Kafka\_SSL Datasource Authentication**.
- Password: applies to GaussDB(DWS), RDS, DDS, and DCS data sources. During the configuration, you need to store the passwords of the data sources in DLI. For details, see **Creating a Password Datasource Authentication**.

### Jobs That Can Connect to Data Sources Through Datasource Authentication

Different types of jobs can connect to data sources through different types of datasource authentication.

- For details about the data sources that Spark SQL jobs can connect to through datasource authentication and their constraints, see Table 6-9.
- For details about the data sources that Flink OpenSource SQL 1.12 jobs can connect to through datasource authentication and their constraints, see Table 6-10.

**Table 6-9** Data sources that Spark SQL jobs can connect to through datasource authentication

Datasource Authentication Type	Data Source	Notes and Constraints
CSS	CSS	The CSS cluster version must be 6.5.4 or later. The security mode has been enabled for the CSS cluster.
Password	GaussDB(DWS), RDS, DDS, and Redis	-

**Table 6-10** Data sources that Flink OpenSource SQL 1.12 jobs can connect to through datasource authentication

Table Type	Datasource Authenticati on Type	Data Source	Notes and Constraints
Sourc e	Kerberos	HBase	Kerberos authentication has been enabled for the MRS cluster.
table		Kafka	Kerberos authentication has been enabled for MRS Kafka.
	Kafka_SSL	Kafka	SASL_SSL authentication has been enabled for DMS Kafka.
			SASL authentication has been enabled for MRS Kafka.
			SSL authentication has been enabled for MRS Kafka.
	Password	GaussDB(DWS), RDS, and Redis	-
Result table	Kerberos	HBase	Kerberos authentication has been enabled for the MRS cluster.
		Kafka	Kerberos authentication has been enabled for MRS Kafka.
	Kafka_SSL	Kafka	SASL_SSL authentication has been enabled for DMS Kafka.
			SASL authentication has been enabled for MRS Kafka.
			SSL authentication has been enabled for MRS Kafka.

Table Type	Datasource Authenticati on Type	Data Source	Notes and Constraints
	Password	GaussDB(DWS), RDS, CSS, and Redis	-
Dime nsion	Kerberos	HBase	Kerberos authentication has been enabled for the MRS cluster.
table	Password	GaussDB(DWS), RDS, and Redis	-

### 6.4.2 Creating a CSS Datasource Authentication

### Scenario

Create a CSS datasource authentication on the DLI console to store the authentication information of the CSS security cluster to DLI. This will allow you to access to the CSS security cluster without having to configure a username and password in SQL jobs.

Create a datasource authentication for a CSS security cluster on the DLI console.

### Notes

A CSS security cluster has been created and has met the following conditions:

- The cluster version is 6.5.4 or later.
- The security mode has been enabled for the cluster.

### Procedure

- 1. Download the authentication credential of the CSS security cluster.
  - a. Log in to the CSS management console and choose **Clusters** > **Elasticsearch**.
  - b. On the **Clusters** page displayed, click the cluster name.
  - c. On the **Cluster Information** page displayed, find the security mode and download the certificate of the CSS security cluster.
- 2. Upload the authentication credential to the OBS bucket.
- 3. Create a datasource authentication.
  - a. Log in to the DLI management console.
  - b. Choose **Datasource Connections**. On the page displayed, click **Datasource Authentication**.
  - c. Click Create.Configure CSS authentication parameters according to Table 6-11.

#### Table 6-11Parameters

Parameter	Description	
Authentica tion Certificate	<ul> <li>Name of the datasource authentication information to be created.</li> <li>The name can contain only digits, letters, and underscores (_), but cannot contain only digits or start with an underscore (_).</li> </ul>	
	<ul> <li>The length of the database name cannot exceed 128 characters.</li> </ul>	
	<ul> <li>It is recommended that the name contain the CSS security cluster name to distinguish security authentication information of different clusters.</li> </ul>	
Туре	Select <b>CSS</b> .	
Username	Username for logging in to the security cluster.	
Password	The password of the security cluster	
Certificate Path	Enter the OBS path to which the security certificate is uploaded, that is, the OBS bucket address in <b>2</b> .	

#### Figure 6-6 Creating a datasource authentication-CSS

### **Create Authentication**

Туре	CSS	•	
★ Authentication Certificate	Enter a name for the certificate.		
★ Username	Enter a username.		?
* Password	Enter a password.	Ø	?
★ Certificate Path	obs://	ß	
	OK Cancel		

4. Create a table to access the CSS cluster.

When creating a table, associate the table with the created datasource authentication to access the CSS cluster.

For example, when using Spark SQL to create a table for accessing the CSS cluster, configure **es.certificate.name** to set the datasource authentication name and then connect to the CSS security cluster.

X
Use Spark SQL to create a table for accessing the CSS cluster by referring to **Creating a DLI Table and Associating It with CSS**.

# 6.4.3 Creating a Kerberos Datasource Authentication

# Scenario

Create a Kerberos datasource authentication on the DLI console to store the authentication information of the data source to DLI. This will allow you to access to the data source without having to configure a username and password in SQL jobs.

### D NOTE

- When Kerberos authentication is enabled for MRS Kafka but SSL authentication is disabled, create a Kerberos authentication. When creating a table, configure **krb\_auth\_name** to associate the datasource authentication.
- If Kerberos authentication and SSL authentication are both enabled for MRS Kafka, you need to create Kerberos and Kafka\_SSL authentications. When creating a table, configure krb\_auth\_name and ssl\_auth\_name to associate the datasource authentications.
- Datasource authentication is not required when Kerberos authentication is disabled but SASL authentication is enabled for MRS Kafka (for example, when a username and a password are used for PlainLoginModule authentication).
- When Kerberos authentication is disabled but SSL authentication is enabled for MRS Kafka, you need to create a Kafka\_SSL authentication. When creating a table, configure **ssl\_auth\_name** to associate the datasource authentication.
- When Kerberos authentication is disabled but SASL authentication and SSL authentication are enabled for MRS Kafka, you need to create a Kafka\_SSL authentication. When creating a table, configure **ssl\_auth\_name** to associate the datasource authentication.

# Data Sources Supported by Kerberos Datasource Authentication

**Table 6-12** lists the data sources supported by Kerberos datasource authentication.

Job Type	Table Type	Data Source	Notes and Constraints
Flink OpenSource SQL	Source table	HBase	Kerberos authentication has been enabled for the MRS cluster.
		Kafka	Kerberos authentication has been enabled for MRS Kafka.
	Result table	HBase	Kerberos authentication has been enabled for the MRS cluster.
		Kafka	Kerberos authentication has been enabled for MRS Kafka.

Table 6-12 Data sources supported	d by Kerberos	datasource	authentication
-----------------------------------	---------------	------------	----------------

Job Type	Table Type	Data Source	Notes and Constraints
	Dimension table	HBase	Kerberos authentication has been enabled for the MRS cluster.

# Procedure

- 1. Download the authentication credential of the data source.
  - a. Log in to MRS Manager.
  - b. Choose **System** > **Permission** > **User**.
  - c. Click **More**, select **Download Authentication Credential**, save the file, and decompress it to obtain the **keytab** and **krb5.conf** files.
- 2. Upload the authentication credential to the OBS bucket.
- 3. Create a datasource authentication.
  - a. Log in to the DLI management console.
  - b. Choose **Datasource Connections**. On the page displayed, click **Datasource Authentication**.
  - c. Click Create.

Configure Kerberos authentication parameters according to Table 6-13.

Parameter	Description
Туре	Select <b>Kerberos</b> .
Authenticatio n Certificate	<ul> <li>Name of the datasource authentication to be created.</li> <li>The name can contain only digits, letters, and underscores (_), but cannot contain only digits or start with an underscore (_).</li> <li>The name can contain a maximum of 128 characters.</li> <li>It is recommended that the name contain the MRS security cluster name to distinguish security authentication information of different clusters.</li> </ul>
Username	Username for logging in to the security cluster.
krb5_conf Path	OBS path to which the <b>krb5.conf</b> file is uploaded. <b>NOTE</b> The <b>renew_lifetime</b> configuration item under <b>[libdefaults]</b> must be removed from <b>krb5.conf</b> . Otherwise, the "Message stream modified (41)" error may occur.
keytab Path	OBS path to which the <b>user.keytab</b> file is uploaded.

### Table 6-13 Parameters

Create Authentication				
Туре	Kerberos	•		
* Authentication Certificate	Enter a name for the certificate.			
★ Username	Enter a username.		?	
★ krb5_conf Path	obs://	ß		
★ keytab Path	obs://	ß		
	OK Cancel			

Figure 6-7 Creating a datasource authentication – Kerberos

4. Create a table to access the MRS cluster.

When creating a data source, associate the data source with the created datasource authentication to access the data source.

 Table 6-14 lists the fields used to associate with the datasource authentication during table creation.

**Table 6-14** Fields that are used to associate with Kerberos datasourceauthentication during table creation

Job Type	Dat a Sou rce	Para meter	Ma nda tory	Data Type	Description
Flink OpenS ource SQL	HBa se	krb_a uth_n ame	No	String	This field is used to associate datasource authentications when source, result, and dimension tables are created.

Job Type	Dat a Sou rce	Para meter	Ma nda tory	Data Type	Description
	Kafk a	krb_a uth_n ame	No	String	This field is used to associate datasource authentications when source and result tables are created.
					Name of the created Kerberos datasource authentication.
					If SASL_PLAINTEXT and Kerberos authentication are both used, you need to configure the following parameters:
					<ul> <li>'properties.sasl.mechanism' = 'GSSAPI'</li> </ul>
					<ul> <li>'properties.security.protocol' = 'SASL_PLAINTEXT'</li> </ul>

For details about how to create a table, see *Data Lake Insight Syntax Reference*.

 Flink OpenSource SQL Syntax Reference: Creating an HBase Source Table

# 6.4.4 Creating a Kafka\_SSL Datasource Authentication

# Scenario

Create a Kafka\_SSL datasource authentication on the DLI console to store the Kafka authentication information to DLI. This will allow you to access to Kafka instances without having to configure a username and password in SQL jobs.

## D NOTE

- When Kerberos authentication is enabled for MRS Kafka but SSL authentication is disabled, create a Kerberos authentication. When creating a table, configure **krb auth name** to associate the datasource authentication.
- If Kerberos authentication and SSL authentication are both enabled for MRS Kafka, you need to create Kerberos and Kafka\_SSL authentications. When creating a table, configure **krb\_auth\_name** and **ssl\_auth\_name** to associate the datasource authentications.
- Datasource authentication is not required when Kerberos authentication is disabled but SASL authentication is enabled for MRS Kafka (for example, when a username and a password are used for PlainLoginModule authentication).
- When Kerberos authentication is disabled but SSL authentication is enabled for MRS Kafka, you need to create a Kafka\_SSL authentication. When creating a table, configure **ssl\_auth\_name** to associate the datasource authentication.
- When Kerberos authentication is disabled but SASL authentication and SSL authentication are enabled for MRS Kafka, you need to create a Kafka\_SSL authentication. When creating a table, configure **ssl\_auth\_name** to associate the datasource authentication.

# Data Sources Supported by Kafka\_SSL Datasource Authentication

 Table 6-15 lists the data sources supported by Kafka\_SSL datasource authentication.

Јоb Туре	Table Type	Data Source	Notes and Constraints
Flink OpenSource SQL	Source table and result	Kafka	SASL_SSL authentication has been enabled for DMS Kafka.
table	table		SASL authentication has been enabled for MRS Kafka.
			SSL authentication has been enabled for MRS Kafka.

Table 6-15 Data sources supported by Kafka\_SSL datasource authentication

# Procedure

- 1. Download the authentication credential.
  - DMS Kafka
    - i. Log in to the DMS (for Kafka) console and click a Kafka instance to access its details page.
    - ii. In the connection information, find the SSL certificate and click **Download**.

Decompress the downloaded **kafka-certs** package to obtain the **client.jks** and **phy\_ca.crt** files.

- MRS Kafka
  - i. Log in to MRS Manager.

- ii. Choose System > Permission > User.
- iii. Click **More**, select **Download Authentication Credential**, save the file, and decompress it to obtain the truststore file.
- 2. Upload the authentication credential to the OBS bucket.
- 3. Create a datasource authentication.
  - a. Log in to the DLI management console.
  - b. Choose **Datasource Connections**. On the page displayed, click **Datasource Authentication**.
  - c. Click Create.

Configure Kafka authentication parameters according to **Table 6-16**.

Table 6-16 Parameters

Parameter	Description
Туре	Select <b>Kafka_SSL</b> .
Authenticatio n Certificate	<ul> <li>Name of the datasource authentication to be created.</li> <li>The name can contain only digits, letters, and underscores (_), but cannot contain only digits or start with an underscore (_).</li> <li>The name can contain a maximum of 128 characters.</li> </ul>
Truststore Path	<ul> <li>OBS path to which the SSL truststore file is uploaded.</li> <li>For MRS Kafka, enter the OBS path of the <b>Truststore.jks</b> file.</li> <li>For DMS Kafka, enter the OBS path of the <b>client.jks</b> file.</li> </ul>
Truststore Password	Truststore password.
Keystore Path	OBS path to which the SSL keystore file (key and certificate) is uploaded.
Keystore Password	Keystore (key and certificate) password.
Key Password	Password of the private key in the keystore file.

Figure 6-8	Creating a	a datasource	authentication -	Kafka_SSI
------------	------------	--------------	------------------	-----------

# **Create Authentication**

•		
1	ч	r
1	~	`

Туре	Kafka_SSL	•
★ Authentication Certificate	Enter a name for the certificate.	
★ Truststore Path	obs://	Þ
Truststore Password	Enter the password.	Ø
Keystore Path	obs://	Þ
Keystore Password	Enter the password.	Ø
Key Password	Enter the password.	Ø
	OK Cancel	

4. Access Kafka with SASL\_SSL authentication enabled.

When creating a data source, associate the data source with the created datasource authentication to access the data source.

**Table 6-17** lists the fields used to associate with the datasource authentication during table creation.

Table 6-17 Fields that are used to associate with Kafka_SSL	. datasource
authentication during table creation	

Paramet er	Ma nda tor y	Dat a Typ e	Description
ssl_auth_ name	No	Stri ng	This field is used to associate datasource authentications when source, result, and dimension tables are created.
			Name of the created Kafka_SSL datasource authentication. This configuration is used when SSL is configured for Kafka.
			<ul> <li>If only SSL is used, configure the following parameter: 'properties.security.protocol '= 'SSL';</li> </ul>
			<ul> <li>If SASL_SSL is used, configure the following parameters:</li> </ul>
			<ul><li>- 'properties.security.protocol' = 'SASL_SSL',</li></ul>
			<ul> <li>'properties.sasl.mechanism' ='GSSAPI or PLAIN'</li> </ul>
			<ul> <li>'properties.sasl.jaas.config' =         <ul> <li>'org.apache.kafka.common.security.plain.Plai</li> <li>nLoginModule required username=\"xxx\"</li>             password=\"xxx\";'</ul></li> </ul>

For details about how to create a table, see *Data Lake Insight Syntax Reference*.

- Flink OpenSource SQL Syntax Reference: Creating a Kafka Source Table

# 6.4.5 Creating a Password Datasource Authentication

# Scenario

Create a password datasource authentication on the DLI console to store passwords of the GaussDB(DWS), RDS, DCS, and DDS data sources to DLI. This will allow you to access to the data sources without having to configure a username and password in SQL jobs.

# Data Sources Supported by Password Datasource Authentication

 Table 6-18 lists the data sources supported by password datasource authentication.

Job Type	Table Type	Data Source	
Spark SQL	-	GaussDB(DWS), RDS, DDS, and Redis	
Flink OpenSource	Source table	GaussDB(DWS), RDS, and Redis	
SQL	Result table	GaussDB(DWS), RDS, CSS, and Redis	
	Dimension table	GaussDB(DWS), RDS, and Redis	

**Table 6-18** Data sources supported by password datasource authentication

# Procedure

- 1. Create a datasource authentication.
  - a. Log in to the DLI management console.
  - b. Choose **Datasource Connections**. On the page displayed, click **Datasource Authentication**.
  - c. Click **Create**.

Configure authentication parameters according to Table 6-19.

### Table 6-19 Parameters

Paramet er	Description
Туре	Select Password.
Authentic ation Certificat e	<ul> <li>Name of the datasource authentication to be created.</li> <li>The name can contain only digits, letters, and underscores (_), but cannot contain only digits or start with an underscore (_).</li> <li>The name can contain a maximum of 128 characters.</li> </ul>
Usernam e	Username for accessing the data source.
Password	Password for accessing the data source.

# Figure 6-9 Creating a datasource authentication – Password

# **Create Authentication**

`	/
2	٢

Туре	Password	
* Authentication Certificate	Enter a name for the certificate.	
Username	Enter a username.	?
* Password	Enter a password.	?
	OK Cancel	

2. Access the data source.

When creating a data source, associate the data source with the created datasource authentication to access the data source.

**Table 6-20** lists the fields used to associate with the datasource authentication during table creation.

Table 6-20 Fields that are used to associate with password datasource
authentication during table creation

Job Type	Parame ter	Man dato ry	Data Type	Description
Spark SQL	passwd auth	No	String	Name of datasource authentication. It is applicable to GaussDB(DWS), RDS, DDS, and Redis data sources.
Flink OpenSo urce SQL	pwd_au th_nam e	No	String	This field is used to associate datasource authentications when source, result, and dimension tables are created.
				Set <b>pwd_auth_name</b> to the name of the password datasource authentication. If this parameter is set, you do not need to configure a username and a password of the data source in SQL jobs.

For details about how to create a table, see *Data Lake Insight Syntax Reference*.

 Flink OpenSource SQL Syntax Reference: Creating a GaussDB(DWS) Source Table

# 6.4.6 Datasource Authentication Permission Management

# Scenario

Grant permissions on a datasource authentication to users so multiple user jobs can use the datasource authentication without affecting each other.

### Notes

- The administrator and the datasource authentication owner have all permissions. You do not need to set permissions for them, and their datasource authentication permissions cannot be modified by other users.
- When setting datasource authentication permissions for a new user, ensure that the user group to which the user belongs has the **Tenant Guest** permission.

For details about the **Tenant Guest** permission and how to apply for the permission, see **System Permissions** and **Creating a User Group** in *Identity and Access Management User Guide*.

# **Granting Permissions on Datasource Connections**

- 1. Log in to the DLI management console.
- 2. Choose **Datasource Connections**. On the page displayed, click **Datasource Authentication**.
- 3. Locate the row containing the datasource authentication to be authorized and click **Manage Permission** in the **Operation** column. The **User Permissions** page is displayed.
- 4. Click **Grant Permission** in the upper right corner of the page. On the **Grant Permission** dialog box displayed, grant permissions on this datasource authentication to other users.

### Figure 6-10 Granting permissions on datasource connections

#### Grant Permission

* Username	Enter a username.		
Select the permi	ssions to be granted to	the user	
Select all			
Access		Update	Delete
Grant Permi	ssion	Revoke Permission	View Other User's Permissions
			7
		Ok Cancel	

Parameter	Description			
Username	Name of the IAM user to whom permissions on the datasource connection are to be granted. <b>NOTE</b> The username is the name of an existing IAM user.			
Select the permissions to be granted to the user	<ul> <li>Access: This permission allows you to access the datasource connection.</li> <li>Update: This permission allows you to update the</li> </ul>			
	<ul> <li>datasource connection.</li> <li>Delete: This permission allows you to delete the datasource connection.</li> </ul>			
	• Grant Permission: This permission allows you to grant the datasource connection permission to other users.			
	• Grant Permission: This permission allows you to revoke the datasource connection permission to other users. However, you cannot revoke the permissions of the datasource connection owner.			
	• View Other User's Permissions: This permission allows you to view the datasource connection permissions of other users.			

**Table 6-21** Permission granting parameters

# Modifying the Permissions of Current User

- 1. Log in to the DLI management console.
- 2. Choose **Datasource Connections**. On the page displayed, click **Datasource Authentication**.
- 3. Locate the row containing the datasource authentication to be authorized and click **Manage Permission** in the **Operation** column. The **User Permissions** page is displayed.
- 4. Click **Set Permission** in the **Operation** column to modify the permissions of the current user. **Table 6-21** lists the detailed permission descriptions.

D NOTE

- If all options under **Set Permission** are gray, you are not allowed to change permissions on this datasource connection. You can apply to the administrator, group owner, or other users who have the permission to grant permissions for the permissions to grant and revoke the datasource authentication permissions.
- To revoke all permissions of the current user, click **Revoke Permission** in the **Operation** column. The IAM user will no longer have any permission on the datasource authentication.

# 6.5 Managing Enhanced Datasource Connections

# 6.5.1 Viewing Basic Information About an Enhanced Datasource Connection

After creating an enhanced datasource connection, you can view and manage it on the management console.

This section describes how to view basic information about an enhanced datasource connection on the management console, including the enhanced datasource connection's host information, IPv6 support, and more.

# Procedure

- 1. Log in to the DLI management console.
- 2. In the navigation pane on the left, choose **Datasource Connections**.
- 3. On the displayed **Enhanced** tab, locate the enhanced datasource connection whose basic information you want to view.
  - In the upper right corner of the list page, click <sup>(2)</sup> to customize the columns to display and set the rules for displaying the table content and the **Operation** column.
  - In the search box above the list, you can filter the required enhanced datasource connection by name or tag.
- 4. Click  $\checkmark$  to expand details about the enhanced datasource connection.

You can view the following information:

- IPv6 Support: If you selected a subnet with IPv6 enabled when creating the enhanced datasource connection, then your enhanced datasource connection will support IPv6.
- Host Information: When accessing an MRS HBase cluster, you need to configure the host name (domain name) and the corresponding IP address of the instance. For details, see Modifying Host Information in an Elastic Resource Pool.

# **6.5.2 Enhanced Connection Permission Management**

# Scenario

Enhanced connections support user authorization by project. After authorization, users in the project have the permission to perform operations on the enhanced connection, including viewing the enhanced connection, binding a created resource pool to the enhanced connection, and creating custom routes. In this way, the enhanced connection can be used across projects. Grant and revoke permissions to and from a user for an enhanced connection.

### **NOTE**

- If the authorized projects belong to different users in the same region, you can use the user account of the authorized projects to log in.
- If the authorized projects belong to the same user in the same region, you can use the current account to switch to the corresponding project.

# Use Cases

Project B needs to access the data source of project A. The operations are as follows:

- For Project A:
  - a. Log in to DLI using the account of project A.
  - b. Create an enhanced datasource connection **ds** in DLI based on the VPC information of the corresponding data source.
  - c. Grant project B the permission to access the enhanced datasource connection **ds**.
- For Project B:
  - a. Log in to DLI using the account of project B.
  - b. Bind the enhanced datasource connection **ds** to a queue.
  - c. (Optional) Set host information and create a route.

After creating a VPC peering connection and route between the enhanced datasource connection of project A and the queue of project B, you can create a job in the queue of project B to access the data source of project A.

# Procedure

- 1. Log in to the DLI management console.
- 2. In the left navigation pane, choose **Datasource Connections**.
- 3. On the **Enhanced** tab page displayed, locate the desired enhanced connection, click **More** in the **Operation** column, and select **Manage Permission**.
  - Granting permission
    - i. In the **Permissions** dialog box displayed, select **Grant Permission** for **Set Permission**.
    - ii. Enter the project ID.
    - iii. Click **OK** to grant the resource pool operation permission to the project.
  - Revoking permission
    - i. In the **Permissions** dialog box displayed, select **Revoke Permission** for **Set Permission**.
    - ii. Select a project ID.
    - iii. Click **OK** to revoke the resource pool operation permission from the specified project.

# 6.5.3 Binding an Enhanced Datasource Connection to an Elastic Resource Pool

# Scenario

To connect other resource pools to data sources through enhanced datasource connections, bind enhanced datasource connections to resource pools on the **Enhanced** tab page.

# Constraints

- Enhanced datasource connections support only dedicated pay-per-use resource pools and queues.
- The CIDR block of the DLI queue to be bound with a datasource connection cannot overlap with that of the data source.
- The **default** queue preset in the system cannot be bound with a datasource connection.

# Procedure

- 1. Log in to the DLI management console.
- 2. In the left navigation pane, choose **Datasource Connections**.
- 3. On the **Enhanced** tab page displayed, bind an enhanced datasource connection to an elastic resource pool:
  - a. Locate your desired enhanced datasource connection, click **More** in the **Operation** column, and select **Bind Resource Pool**.
  - b. In the **Bind Resource Pool** dialog box, select the resource pool to be bound for **Resource Pool**.
  - c. Click **OK**.
- 4. View the connection status on the **Enhanced** tab page.
  - After an enhanced datasource connection is created, the status is Active, but it does not indicate that the queue is connected to the data source. Go to the queue management page to check whether the data source is connected. The procedure is as follows:
    - i. In the navigation pane on the left, choose **Resources** > **Queue Management**. On the page displayed, locate a desired queue.
    - ii. Click More in the Operation column and select Test Address Connectivity.
    - iii. Enter the IP address and port number of the data source.
  - On the details page of an enhanced datasource connection, you can view information about the VPC peering connection.
    - VPC peering ID: ID of the VPC peering connection created in the cluster to which the queue belongs.

A VPC peering connection is created for each queue bound to an enhanced datasource connection. The VPC peering connection is used for cross-VPC communication. Ensure that the security group used by the data source allows access from the CIDR block of the DLI queue, and do not delete the VPC peering connection during the datasource connection.

Status of the VPC peering connection:

The status of a datasource connection can be **Creating**, **Active**, or **Failed**.

If the connection status is **Failed**, click  $\checkmark$  on the left to view the detailed error information.

### Figure 6-11 Viewing details



# 6.5.4 Unbinding an Enhanced Datasource Connection from an Elastic Resource Pool

# Scenario

Unbind an enhanced datasource connection from an elastic resource pool that does not need to access a data source through an enhanced datasoruce connection.

# Constraints

If the status of the VPC peering connection created for binding an enhanced datasource connection to an elastic resource pool is **Failed**, the elastic resource pool cannot be unbound.

- 1. Log in to the DLI management console.
- 2. In the left navigation pane, choose **Datasource Connections**.
- 3. On the **Enhanced** tab page displayed, use either of the following methods to unbind an enhanced datasource connection from an elastic resource pool:
  - Method 1:
    - i. Locate your desired enhanced datasource connection, click **More** in the **Operation** column, and select **Unbind Resource Pool**.
    - ii. In the **Unbind Resource Pool** dialog box, select the resource pool to be unbound for **Resource Pool**.
    - iii. Click OK.
  - Method 2:
    - i. Click your desired enhanced datasource connection in the list.
    - ii. Locate your desired resource pool and click **Unbind Resource Pool** in the **Operation** column.
    - iii. Click OK.

# 6.5.5 Adding a Route for an Enhanced Datasource Connection

# Scenario

A route is configured with the destination, next hop type, and next hop to determine where the network traffic is directed. Routes are classified into system routes and custom routes.

After an enhanced connection is created, the subnet is automatically associated with the default route. You can add custom routes as needed to forward traffic destined for the destination to the specified next hop.

### **NOTE**

- When an enhanced connection is created, the associated route table is the one associated with the subnet of the data source.
- The route to be added in the **Add Route** dialog box must be one in the route table associated with the subnet of the resource pool.
- The subnet of the data source must be different from that used by the resource pool. Otherwise, a network segment conflict occurs.

- 1. Log in to the DLI management console.
- 2. In the left navigation pane, choose **Datasource Connections**.
- 3. On the **Enhanced** tab page displayed, locate the row containing the enhanced connection to which a route needs to be added, and add the route.
  - Method 1:
    - i. On the **Enhanced** tab page displayed, locate the enhanced datasource connection to which a route needs to be added and click **Manage Route** in the **Operation** column.
    - ii. Click Add Route.
    - iii. In the **Add Route** dialog box, enter the route information. For details about the parameters, see **Table 6-22**.
    - iv. Click OK.
  - Method 2:
    - i. On the **Enhanced** tab page displayed, locate the enhanced datasource connection to which a route needs to be added, click **More** in the **Operation** column, and select **Add Route**.
    - ii. In the **Add Route** dialog box, enter the route information. For details about the parameters, see **Table 6-22**.
    - iii. Click OK.

Parameter	Description					
Route Name	Name of a custom route, which is unique in the same enhanced datasource scenario. The name can contain 1 to 64 characters. Only digits, letters, underscores (_), and hyphens (-) are allowed.					
IP Address	You can add an IPv4 or IPv6 address.					
Туре	If your data source has IPv6 enabled and the current enhanced datasource connection supports IPv6, you can select IPv6 routes when adding a route table.					
	You can check whether the current enhanced datasource connection supports IPv6 in its basic information. For details, see Viewing Basic Information About an Enhanced Datasource Connection.					
	The route IP address example is as follows:					
	• IPv4 address: <b>192.168.2.0/24</b> .					
	<ul> <li>IPv6 address: 2407:c080:802:be7::/64.</li> </ul>					
IP Address	Custom route CIDR block. The CIDR block of different routes can overlap but cannot be the same.					
	Do not add the CIDR blocks <b>100.125</b> . <i>xx</i> . <i>xx</i> and <b>100.64</b> . <i>xx</i> . <i>xx</i> to prevent conflicts with the internal CIDR blocks of services such as SWR. This can lead to failure of the enhanced datasource connection.					

Table 6-22 Parameters for adding a custom route

4. After adding a route, you can view the route information on the route details page.

# 6.5.6 Deleting the Route for an Enhanced Datasource Connection

# Scenario

Delete a route that is no longer used.

# Constraints

A custom route table cannot be deleted if it is associated with a subnet.

- 1. Log in to the DLI management console.
- 2. In the left navigation pane, choose **Datasource Connections**.
- 3. On the **Enhanced** tab page displayed, locate the row containing the enhanced connection from which the route needs to be deleted, and delete the route.
  - Method 1:

- i. On the **Enhanced** tab page displayed, locate the enhanced connection from which the route needs to be deleted and click **Manage Route** in the **Operation** column.
- ii. Locate the route to be deleted and click **Delete** in the **Operation** column.
- iii. In the dialog box displayed, click **OK**.
- Method 2:
  - i. On the **Enhanced** tab page displayed, locate the enhanced connection from which the route needs to be deleted, click **More** in the **Operation** column, and select **Delete Route**.
  - ii. In the **Delete Route** dialog box displayed, confirm the route information.
  - iii. Click **Yes**.

# 6.5.7 Modifying Host Information in an Elastic Resource Pool

# Scenario

Host information is the mapping between host IP addresses and domain names. After you configure host information, jobs can only use the configured domain names to access corresponding hosts. After a datasource connection is created, you can modify the host information.

When accessing the HBase cluster of MRS, you need to configure the host name (domain name) and IP address of the instance.

# Constraints

You have obtained the MRS host information by referring to **How Do I Obtain** MRS Host Information?

# **Modifying Host Information**

- 1. Log in to the DLI management console.
- 2. In the left navigation pane, choose **Datasource Connections**.
- 3. On the **Enhanced** tab page displayed, locate the enhanced datasource connection to be modified, click **More** in the **Operation** column, and select **Modify Host**.
- 4. In the **Modify Host** dialog box displayed, enter the obtained host information.

Enter host information in the format of *Host IP address Host name*. Information about multiple hosts is separated by line breaks.

Example:

192.168.0.22 node-masterxxx1.com

192.168.0.23 node-masterxxx2.com

Obtain the MRS host information by referring to **How Do I Obtain MRS Host Information?** 

5. Click OK.

# How Do I Obtain MRS Host Information?

## • Method 1: View MRS host information on the management console.

To obtain the host name and IP address of an MRS cluster, for example, MRS 3.*x*, perform the following operations:

- a. Log in to the MRS management console.
- b. On the **Active Clusters** page displayed, click your desired cluster to access its details page.
- c. Click the **Components** tab.
- d. Click ZooKeeper.
- e. Click the **Instance** tab to view the corresponding service IP addresses. You can select any service IP address.
- f. Modify host information by referring to **Modifying Host Information**.

### **NOTE**

If the MRS cluster has multiple IP addresses, enter any service IP address when creating a datasource connection.

# • Method 2: Obtain MRS host information from the /etc/hosts file on an MRS node.

- a. Log in to any MRS node as user root.
- b. Run the following command to obtain MRS hosts information. Copy and save the information.

### cat /etc/hosts

### Figure 6-12 Obtaining hosts information



- c. Modify host information by referring to **Modifying Host Information**.
- Method 3: Log in to FusionInsight Manager to obtain host information.
  - a. Log in to FusionInsight Manager.
  - b. On FusionInsight Manager, click **Hosts**. On the **Hosts** page, obtain the host names and service IP addresses of the MRS hosts.

### Figure 6-13 FusionInsight Manager

🟥 🛛 FusionInsight Manager	Homepage Cluster - Hosts O&M Au	dit Tenant Resources System			<b>O</b> o <b>O</b> o <b>O</b>	0 0 :   🗐 0   Hello, admin 🗸
(h)	Hosts					Host View Role View
Hosts	Add More * ExpertAll				All types	• Advanced Search V C @
	Host Name JE	Management IP Addr JE Service IP Address JE	Running Status JE	CPU Usage ↓Ξ	Memory(GB) ↓=	Disk(GB) JE
Holli	server-2110081635-0001		<ul> <li>Normal</li> </ul>		11/30.9	- 34.87/280.47
Resource Overview	server-2110081635-0002		Normal	20%	9.35/30.9	- 31.35/280.47
	server-2110081635-0003		<ul> <li>Normal</li> </ul>	21%	14.04/30.9	- 24.45/280.47
	#server-2110082001-0017		<ul> <li>Normal</li> </ul>	73%	43.89/56.49	157.05/414.76
	server-2110062001-0018		<ul> <li>Normal</li> </ul>	- 85%	46.39/56.49	144.89/414.76
	server-2110052001-0019		<ul> <li>Normal</li> </ul>	59%	28.45/95.49	70 2/297.14

c. Modify host information by referring to **Modifying Host Information**.

# 6.5.8 Enhanced Datasource Connection Tag Management

# Scenario

A tag is a key-value pair customized by users and used to identify cloud resources. It helps users to classify and search for cloud resources. A tag consists of a tag key and a tag value.

If you use tags in other cloud services, you are advised to create the same tag keyvalue pairs for cloud resources used by the same business to keep consistency.

If your organization has configured tag policies for DLI, add tags to resources based on the policies. If a tag does not comply with the tag policies, resource creation may fail. Contact your organization administrator to learn more about tag policies.

DLI supports the following two types of tags:

- Resource tags: non-global tags created on DLI.
- Predefined tags: global tags created on Tag Management Service (TMS).
   For more information about predefined tags, see Tag Management Service User Guide.
- DLI allows you to add, modify, or delete tags for datasource connections.

- 1. In the left navigation pane of the DLI management console, choose **Datasource Connections**.
- 2. In the **Operation** column of the link, choose **More** > **Tags**.
- 3. The tag management page is displayed, showing the tag information about the current connection.
- 4. Click **Add/Edit Tag**. The **Add/Edit Tag** dialog is displayed. Add or edit tag keys and values and click **OK**.

### Figure 6-14 Adding/Editing a tag

# Add/Edit Tag

 $\times$ 

It is recommended that you use TMS's predefined tag function to add the same tag to different cloud resources. View predefined tags  $\ C$ 

To add a tag, enter a tag key and a tag value below.

Enter a tag key	E	Enter a tag value	Add
0 tags available for addition.			
	OK	Cancel	

## Table 6-23Tag parameters

Param eter	Description
Tag key	You can perform the following operations:
	• Click the text box and select a predefined tag key from the drop-down list.
	To add a predefined tag, you need to create one on TMS and then select it from the <b>Tag key</b> drop-down list. You can click <b>View predefined tags</b> to go to the <b>Predefined Tags</b> page of the TMS console. Then, click <b>Create Tag</b> in the upper corner of the page to create a predefined tag.
	For details, see <b>Creating Predefined Tags</b> in <i>Tag Management Service User Guide</i> .
	• Enter a tag key in the text box.
	<b>NOTE</b> A tag key can contain a maximum of 128 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed, but the value cannot start or end with a space or start with <b>_sys_</b> .
Tag	You can perform the following operations:
value	• Click the text box and select a predefined tag value from the drop-down list.
	Enter a tag value in the text box.
	NOTE A tag value can contain a maximum of 255 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed.

- 5. Click **OK**.
- 6. (Optional) To delete a tag, locate the row where the tag locates in the tag list and click **Delete** in the **Operation** column to delete the tag.

# 6.5.9 Deleting an Enhanced Datasource Connection

# Scenario

Delete an enhanced datasource connection that is no longer used on the console.

# Procedure

- 1. Log in to the DLI management console.
- 2. In the left navigation pane, choose Datasource Connections.
- 3. On the **Enhanced** tab page displayed, locate the enhanced datasource connection to be deleted and click **Delete** in the **Operation** column.
- 4. In the dialog box displayed, click **Yes**.

# 6.6 Example Typical Scenario: Connecting DLI to a Data Source on a Private Network

# Scenario

Typically, connecting DLI to a data source on a private network means connecting it to a Huawei Cloud service, such as MRS, RDS, CSS, Kafka, or GaussDB(DWS). DLI's enhanced datasource connections use VPC peering connections to directly connect DLI to VPC networks of destination data sources.

This section describes how to connect DLI to a data source on a private network using an enhanced datasource connection.

If the network is disconnected when an enhanced datasource connection is created, you can rectify the fault based on the overall process and procedure in this section.

# **Overall Process**

Figure 6-15 Configuration process of an enhanced datasource connection



# Prerequisites

• You have created a queue. For details about how to create a queue, see **Creating a Queue**.

# 

The queue billing mode must be **Pay-per-use**, and **Dedicated Resource Mode** must be selected after you select a queue type.

Enhanced datasource connections can be created only for pay-per-use resources in dedicated resource mode.

• A cluster of the external data source has been created. You can select a data source as needed.

Service Name	Reference Documents
RDS	Buying an RDS for MySQL Instance
GaussDB(DWS)	Creating a GaussDB(DWS) Cluster
DMS Kafka	Creating a Kafka Instance
CSS	Creating a CSS Cluster
MRS	Creating an MRS Cluster

Table 6-24 Reference for creating clusters of other data sources
--

# 

- The CIDR block of the DLI queue bound with a datasource connection cannot overlap with the CIDR block of other data sources.
- Datasource connections cannot be bound with the **default** queue.

# Step 1: Obtain the Floating IP Address, Port Number, and Security Group of an External Data Source

Table 6-25 Data source information to be obtained

Dat a Sour ce	Obtain Method
DMS Kafk a	<ol> <li>On the Kafka management console, click an instance name on the DMS for Kafka page. Basic information of the Kafka instance is displayed.</li> </ol>
	<ol> <li>In the Connection pane, obtain the Instance Address (Private Network) value. In the Network pane, obtain the VPC and subnet of the instance.</li> </ol>
	3. In the <b>Network</b> pane, obtain the security group of the instance.

Dat a Sour ce	Obtain Method
RDS	On the <b>Instances</b> page of the RDS console, click the target DB instance name. In the displayed page, locate the <b>Connection Information</b> pane and obtain the <b>Floating IP Address</b> , <b>VPC</b> , <b>Subnet</b> , <b>Database Port</b> , and <b>Security Group</b> .
CSS	<ol> <li>On the CSS management console, choose Clusters &gt; Elasticsearch. On the displayed page, click the name of the created CSS cluster to view basic information.</li> <li>On the Cluster Information page, obtain the Private Network Address, VPC, Subnet, and Security Group.</li> </ol>
Gaus sDB( DWS )	<ol> <li>On the GaussDB(DWS) management console, choose Clusters. On the displayed page, click the name of the created GaussDB(DWS) cluster to view basic information.</li> <li>On the Basic Information tak leasts the Connection Information.</li> </ol>
,	2. On the <b>Basic information</b> tab, locate the <b>Connection Information</b> pane and obtain the private IP address and port number of the DB instance. In the <b>Network</b> pane, obtain the VPC, subnet, and security group information.

Dat a Sour ce	Obtain Method			
MRS	An MRS 3.x cluster is used as an example.			
HBa se	<ol> <li>Log in to the MRS management console, click a cluster name on the Clusters &gt; Active Clusters page to view basic information.</li> </ol>			
	2. On the dashboard, obtain VPC, subnet, and security group from the <b>Basic Information</b> pane.			
	3. The ZooKeeper instance and its port of the MRS cluster are required for creating a job that connects DLI to MRS HBase. You need to obtain the host information of the MRS cluster.			
	a. Log in to MRS Manager by referring to Accessing FusionInsight Manager. On MRS Manager, choose Cluster > Name of the desired cluster > Services > ZooKeeper. Click the Instance tab and obtain the ZooKeeper host information such as the host name and service IP address.			
	b. On MRS Manager, choose Cluster and click the name of the desired cluster. Choose Services > ZooKeeper. Click the Configurations tab and select All Configurations, search for the clientPort parameter, and obtain its value, that is, the ZooKeeper port number.			
	c. Log in to any MRS node as user <b>root</b> in SSH mode. For details, see Logging In to an ECS.			
	<ul> <li>Run the following command to obtain MRS hosts information.</li> <li>Copy and save the information.</li> <li>cat /etc/hosts</li> </ul>			
	An example query result is as follows:			
	<pre>[root@node-master1kUno ~]# cat /etc/hosts ::1 localhost localhost.localdomain localhost6 localhost6.localdomain6 127.0.0.1 localhost localhost.localdomain localhost4 localhost4.localdomain4 10.10.10.10 hadoop.hadoop.com 10.10.10.10 manager 192.168.0.22 node-master3tVb6.mrs-v08w.com node-group-1ySw0.mrs-v08w.com. 192.168.0.238 node-group-1ySw0.mrs-v08w.com node-group-1ySw0.mrs-v08w.com. 192.168.0.154 node-group-1ySw0.mrs-v08w.com node-master1kOno.mrs-v08w.com. 192.168.0.154 node-group-1yLgA.mrs-v08w.com node-master2qLhC.mrs-v08w.com. 192.168.0.71 node-master2qLhC.mrs-v08w.com node-group-1yRpv.mrs-v08w.com. 192.168.0.7</pre>			

# Step 2: Obtain the CIDR Block of the DLI Queue

On the DLI management console, choose **Resources** > **Queue Management** from

the navigation pane. Locate the queue you have created, and click  $\stackrel{\checkmark}{\sim}$  next to the queue name to view the CIDR block of the queue.

# Step 3: Add a Rule to the Security Group of the External Data Source to Allow Access from the DLI Queue

- 1. Log in to the VPC console.
- 2. In the navigation pane on the left, choose **Access Control** > **Security Groups**.
- 3. Click the name of the security group to which the external data source belongs.

Obtain the security group name of the data source on the management console of the data source by referring to **Step 1: Obtain the Floating IP Address, Port Number, and Security Group of an External Data Source**.

4. On the **Inbound Rules** tab, add a rule to allow access from the queue network segment.

Set the inbound rule parameters based on Table 6-26.

### Figure 6-16 Adding an inbound rule

Add Inbound	d Rule Learn	more about security group co	nfiguration.			
Some secur	ity group rules will not	take effect for ECSs with certain spe	cifications. Learn mor	e		
Security Group d	lefault tiple rules in a batch.					
Priority (?)	Action (?)	Protocol & Port (?)	Туре	Source (?)	Description	Operation
1-100	Allow <b>v</b>	Protocols/TCP (Custo	IPv4 v	IP address         •           0.0.0.0/0         •		Replicate Delete
			+ Add Rule			
			OK Can	cel		

### Table 6-26 Inbound rule parameters

Parameter	Description	Example
Priority	The security group rule priority.	1
	The priority value ranges from 1 to 100. The default value is <b>1</b> , indicating the highest priority. A smaller value indicates a higher priority of a security group rule.	
Action	Action of the security group rule.	Select <b>Allow</b> .

Parameter	Description	Example	
Protocol &Port	<ul> <li>Network protocol: The value can be All, TCP, UDP, ICMP, or GRE.</li> </ul>	In this example, select <b>TCP</b> . Leave the port blank or set it to the data source port obtained in <b>Step 1: Obtain the Floating</b> <b>IP Address, Port Number, and</b> <b>Security Group of an External</b> <b>Data Source</b> .	
	• Port: Port or port range over which the traffic can reach your instance. The port ranges from 1 to 65535.		
Туре	Type of IP addresses.	IPv4	
Source	Allow access from IP addresses or instances in another security group.	In this example, enter the queue CIDR block obtained in Step 2: Obtain the CIDR Block of the DLI Queue.	
Description	Supplementary information about the security group rule. This parameter is optional.	_	

# Step 4: Create an Enhanced Datasource Connection

- 1. Log in to the DLI management console. In the navigation pane on the left, choose **Datasource Connections**. On the displayed page, click **Create** in the **Enhanced** tab.
- 2. In the displayed dialog box, set the following parameters:
  - **Connection Name**: Name of the enhanced datasource connection
  - **Resource Pool**: Select the target DLI queue. (Queues that are not added to a resource pool are displayed in this list.)
  - VPC: VPC of the data source obtained in Step 1: Obtain the Floating IP Address, Port Number, and Security Group of an External Data Source
  - Subnet: Subnet of the data source obtained in Step 1: Obtain the Floating IP Address, Port Number, and Security Group of an External Data Source
  - Set other parameters as you need.
- 3. Click **OK**. Click the name of the created datasource connection to view its status. You can perform subsequent steps only after the connection status changes to **Active**.
- 4. To connect to MRS HBase, you need to add MRS host information. The procedure is as follows:
  - a. On the **Datasource Connections** page, click the **Enhanced** tab and locate the row that contains the created enhanced datasource connection. Click **More** > **Modify Host** in the **Operation** column.
  - b. In the dialog box that appears, enter the MRS HBase host information obtained in **Step 1: Obtain the Floating IP Address, Port Number, and**

×

**Security Group of an External Data Source** to the **Host Information** box.

### Figure 6-17 Modifying host information

### Modify Host

Connection Name		
Host Information	Enter host information in the format "host IP address host name". Specify the information for each host on a separate line.	
	OK Cancel	

c. Click **OK**.

# **Step 5: Test Network Connectivity**

- Choose Resources > Queue Management from the left navigation pane, locate the target queue. In the Operation column, click More > Test Address Connectivity.
- In the displayed dialog box, enter the IP address and port number of the data source obtained in Step 1: Obtain the Floating IP Address, Port Number, and Security Group of an External Data Source in the address box and click Test. If the queue passes the test, it can access the data source.

**NOTE** 

For MRS HBase, use **ZooKeeper IP address:ZooKeeper port** or **ZooKeeper host** information:ZooKeeper port for the test.

# 6.7 Example Typical Scenario: Connecting DLI to a Data Source on a Public Network

### Scenario

A public network data source is one that is accessible over the Internet and has a public IP address. By connecting DLI to the public network, you can access these data sources.

This section explains how to connect DLI to a public network by setting up SNAT rules and configuring routing information.

# Procedure



# Step 1: Create a VPC

Log in to the VPC console and create a VPC. The created VPC is used for NAT to access the public network.

For details about how to create a VPC, see Creating a VPC.

## Figure 6-19 Creating a VPC

Basic Information	
Region	·
	Regions are geographic areas isolated from each other. Resources are region-specific and cannot latency and quick resource access, select the nearest region.
Name	vpc-9334
IPv4 CIDR Block	
	Recommended:10.0.0.0/8-24 ( Select ) 172.16.0.0/12-24 ( Select ) 192.168.0.0/16-24 ( Select )

# Step 2: Create an Elastic Resource Pool and Add Queues Within It

- 1. Log in to the DLI management console.
- 2. In the navigation pane on the left, choose **Resources** > **Resource Pool**.
- 3. On the displayed page, click **Buy Resource Pool** in the upper right corner.
- 4. On the displayed page, set the parameters.

## Table 6-27 Parameter descriptions

Parameter	Description
Region	Select a region where you want to buy the elastic resource pool.
Project	Project uniquely preset by the system for each region
Name	Name of the elastic resource pool
Specifications	Specifications of the elastic resource pool
CU Range	The maximum and minimum CUs allowed for the elastic resource pool
CIDR Block	CIDR block the elastic resource pool belongs to. If you use an enhanced datasource connection, this CIDR block cannot overlap that of the data source. <b>Once set, this CIDR block cannot be changed.</b>

Parameter	Description
Enterprise Project	Select an enterprise project for the elastic resource pool.

- 5. Click **Buy**.
- 6. Click **Submit**.
- 7. In the elastic resource pool list, locate the pool you just created and click **Add Queue** in the **Operation** column.
- 8. Set the basic parameters listed below.

Table 6-28 Basic p	arameters for adding a queue
	<b>–</b> • •

Parameter	Description
Name	Name of the queue to add
Туре	<ul> <li>Type of the queue</li> <li>To execute SQL jobs, select For SQL.</li> <li>To execute Flink or Spark jobs, select For general purpose.</li> </ul>
Engine	SQL queue engine. The options are <b>Spark</b> and <b>HetuEngine</b> .
Enterprise Project	Select an enterprise project for the elastic resource pool.

9. Click **Next** and configure scaling policies for the queue.

Click **Create** to add a scaling policy with varying priority, period, minimum CUs, and maximum CUs.

Figure 6-20 shows the scaling policy configured in this example.

Figure 6-20 Configuring a scaling policy when adding a queue

Basic Configuration —	2 Elastic Resources		
View scaling policies of	all queues in		
1. The priority ranges fro	om 1 to 100. If you do not set the priority for a specific	eriod, the default value is 1.	
2.A new policy overwrite	es the default policy.		
4 The total minimum Cl	Is of all queues in an elastic resource pool cannot be i	ore than the minimum CUs of the pool	
5.The maximum CUs of	f any queue in an elastic resource pool cannot be more	than the maximum CUs of the pool.	
Priority	Period	Min CU	Max CU
1	00 ~ 24 ~	- 16 +	- 16 +

Paramet er	Description	Example Value
Priority	Priority of the scaling policy in the current elastic resource pool. A larger value indicates a higher priority. In this example, only one scaling policy is configured, so its priority is set to <b>1</b> by default.	1
Period	The first scaling policy is the default policy, and its <b>Period</b> parameter configuration cannot be deleted or modified. The period for the scaling policy is from 00	00-24
	to 24.	
Min CU	Minimum number of CUs allowed by the scaling policy	16
Max CU	Maximum number of CUs allowed by the scaling policy	64

Table 6-29	Scaling	policy	parameters
------------	---------	--------	------------

10. Click **OK**.

# Step 3: Create an Enhanced Datasource Connection Between the Queue and a VPC

- 1. In the navigation pane of the DLI management console, choose **Datasource Connections**.
- 2. In the **Enhanced** tab, click **Create**.
  - Enter the connection name, select the created queue, VPC, and subnet, and enter the host information (optional).

## Figure 6-21 Creating an enhanced datasource connection

# **Create Enhanced Connection**

After you create the enhanced datasource connection, the system will automatically create a connection and required routes.

* Connection Name	dli_peer_0927	
Resource Pool		•
* VPC	vpc-9334(10.0.0/8)	•
* Subnet	subnet-9344(10.0.0/24)	•

# Step 4: Buy an EIP

- 1. Log in to the **EIPs** page of the network console, click **Buy EIP**.
- In the displayed page, configure the parameters as required.
   For details about how to set the parameters, see Buy EIP.

# Step 5: Configure a NAT Gateway

**Step 1** Create a NAT gateway.

- 1. Log in to the console and search for **NAT Gateway** in the Service List. The **Public NAT Gateways** page of the network console is displayed.
- 2. Click **Buy Public NAT Gateway** and configure the required parameters. For details, see **Buying a Public NAT Gateway**.

### Figure 6-22 Buying a NAT gateway

★ Billing Mode	Yearly/Monthly Pay-per-use	
* Region	<b>v</b>	
	Regions are geographic areas isolated from e cannot be used across regions through interna and quick resource access, select the nearest	each other. Resources are region-specific and al network connections. For low network latency t region.
★ Name	nat-32c8	
* VPC	vpc-9334 💌	C View VPCs
* Subnet	subnet-9344(10.0.0.0/24)	C View Subnets Available private IP addresses: 2
	The selected subnet is for the NAT gateway o after the NAT gateway is created, you need to	nly. To enable communications over the Internet, add rules.
* Specifications	Small Medium Larg	ge Extra-large
	Supports up to 10,000 connections.Learn more	re
Advanced Settings <	Description   Tag	

3. Click Next, confirm the configurations, and click Submit.

### **NOTE**

During the configuration, set **VPC** to the one created in **Step 1: Create a VPC**.

**Step 2** Add a route.

In the navigation pane on the left of the network console, choose Virtual **Private Cloud** > **Route Tables**. After a NAT gateway instance is created, a route to that gateway is automatically created. Click the route table name to view the automatically created route.

The destination address is the public IP address you want to access, and the next hop is the NAT gateway.

### Figure 6-23 Viewing the route

Routes					
Delete	Add Route	Replicate Route Q	Learn how to configur	re routes.	
Destir	nation (?)	Next Ho	р Туре	Next Hop  ?	Туре (?)
∧ Local		Local		Local	System
Deat	tination		Next Hen Time		Next Hen
Desi	unation		Next Hop Type		Next Hop
14.	.38/32		NAT gateway		nat-32c8

### Step 3 Add an SNAT rule.

You need to add SNAT rules for the new NAT gateway to allow the hosts in the subnet to communicate with the Internet.

- 1. Click the name of the created NAT gateway on the **Public NAT Gateways** page of the network console.
- 2. On the **SNAT Rules** tab, click **Add SNAT Rule**.

For details, see Adding an SNAT Rule.

3. Scenario: Select Direct Connect/Cloud Connect.

- 4. **Subnet**: Select the subnet where the queue you want to connect locates.
- 5. **EIP**: Select the target EIP.

### Figure 6-24 Adding an SNAT rule

Add SNAT Rule

<ul> <li>If both an EIP and a NAT g</li> <li>It is not recommended that</li> <li>An SNAT rule cannot share</li> </ul>	ateway are configured for a an SNAT rule and a DNAT r an EIP with a DNAT rule wi	server, data will be fo ule share the same E th Port Type set to Al	rwarded through th EIP because there r I ports.	e EIP. View restriction nay be service conflic	ns cts.	
Public NAT Gateway Name	nat-lishenrui					
* Scenario	VPC	Direct Connect	Cloud Connect			
	172 · 16 · 0	· 0 / 16	?			
* EIP	You can select more El	Ps. (?) View EIP	Specify filter of	riteria.		Q
	EIP	EIP Type E	Bandwidth Na	Bandwidth(M	Billing Mode	Enterprise Pr

6. Click OK.

----End

# Step 6: Add a Custom Route

Add a custom route for the enhanced datasource connection you have created. Specify the route information of the IP address you want to access.

For details, see **Custom Route Information**.

Figure 6-25 Adding route information for test

Add Route		×
★ Route Name		
★ IP Address	14 . 0 24	
	OK Cancel	

# Step 7: Test the Connectivity to the Public Network

Test the connectivity between the queue and the public network. Click **More** > **Test Address Connectivity** in the **Operation** column of the target queue and enter the public IP address you want to access.

Figure 6-26 Testing address connectivity

# Test Address Connectivity

 $\times$ 

Tests whether an address is reachable from a specified cluster. The address can be a domain name, an IP address, or a specified port.

* Address	114.	:80		
		Test	Cancel	
# Configuring an Agency to Allow DLI to Access Other Cloud Services

# 7.1 DLI Agency Overview

# What Is an Agency?

Cloud services often interact with each other, with some of which dependent on other services. You can create an agency to delegate DLI to use other cloud services and perform resource O&M on your behalf.

For example, the AK/SK required by DLI Flink jobs is stored in DEW. To allow DLI to access DEW data during job execution, you need to provide an IAM agency to delegate the permissions to perform operations on DEW data to DLI.

Figure 7-1 DLI service agency



# **DLI Agencies**

Before using DLI, you are advised to set up DLI agency permissions to ensure the proper functioning of DLI.

- By default, DLI provides the following agencies: dli\_admin\_agency, dli\_management\_agency, and dli\_data\_clean\_agency. The names of these agencies are fixed, but the permissions contained in them can be customized. In other scenarios, you need to create custom agencies. For details about the agencies, see Table 7-1.
- DLI upgrades **dli\_admin\_agency** to **dli\_management\_agency** to meet the demand for fine-grained agency permissions management. The new agency has the necessary permissions for datasource operations, notifications, and user authorization operations. For details, see **Configuring DLI Agency Permissions**.
- To use Flink 1.15, Spark 3.3.1 (Spark general queue scenario), or a later version to execute jobs, perform the following operations:

Create an agency on the IAM console and add the agency information to the job configuration. For details, see **Creating a Custom DLI Agency**.

- Common scenarios for creating an agency: DLI is allowed to read and write data from and to OBS, dump logs, and read and write Flink checkpoints. DLI is allowed to access DEW to obtain data access credentials and access catalogs to obtain metadata.
- You cannot use the default agency names dli\_admin\_agency, dli\_management\_agency, or dli\_data\_clean\_agency. It must be unique.
- If the engine version is earlier than Flink 1.15, **dli\_admin\_agency** is used by default during job execution. If the engine version is earlier than Spark 3.3.1, user authentication information (AK/SK and security token) is used during job execution.

This means that jobs whose engine versions are earlier than Flink 1.15 or Spark 3.3.1 are not affected by the update of agency permissions and do not require custom agencies.

• To maintain compatibility with existing job agency permission requirements, **dli\_admin\_agency** will still be listed in the IAM agency list even after the update.

#### D NOTE

- Only the tenant account or a member account of user group **admin** can authorize the service.
- Do not delete the agency created by the system by default.

Table 7-1	DLI	agencies
-----------	-----	----------

Permission	Туре	Description
dli_admin_agency	Default agency	This agency has been discarded and is not recommended. Upgrade the agency to <b>dli_management_agency</b> as soon as possible.
		For details about how to update an agency, see <b>Configuring DLI Agency Permissions</b> .
dli_management_a gency	Default agency	DLI system agency, which is used to delegate operation permissions to DLI so that DLI can use other cloud services and perform resource O&M operations on your behalf. This agency grants permissions for datasource operations, message notifications, and user authorization operations. For details about the permissions of an agency, see Table 7-2.
dli_data_clean_age ncy	Default agency, which needs to be authorized by users	Data cleanup agency, which is used to clean up data according to the lifecycle of a table and clean up lakehouse table data. You need to create a DLI agency named <b>dli_data_clean_agency</b> on IAM and grant permissions to it.
		You need to create an agency and customize permissions for it. However, the agency name is fixed to <b>dli_data_clean_agency</b> .
		For details about the permission policies of an agency, see <b>Agency Permission</b> <b>Policies in Common Scenarios</b> .

Permission	Туре	Description
Other custom agencies	Custom agency	When using Flink 1.15, Spark 3.3, or a later version to execute jobs, create an agency on the IAM console and add new agency information to the job configuration. For details, see <b>Creating a Custom DLI Agency</b> .
		Common scenarios for creating an agency: DLI is allowed to read and write data from and to OBS to transfer logs. DLI is allowed to access DEW to obtain data access credentials and access catalogs to obtain metadata.
		You cannot use the default agency names dli_admin_agency, dli_management_agency, or dli_data_clean_agency. It must be unique.
		For details about the permission policies of an agency, see <b>Agency Permission</b> <b>Policies in Common Scenarios</b> .

Table 7-2 Permissions cor	ntained in the dli	_management_	_agency agency
---------------------------	--------------------	--------------	----------------

Policy	Description
IAM ReadOnlyAccess	To authorize IAM users who have not logged in to DLI, you need to obtain their information. So, the permissions contained in the <b>IAM ReadOnlyAccess</b> policy are required.
DLI Datasource Connections Agency Access	Permissions to access and use VPCs, subnets, routes, and VPC peering connections
DLI Notification Agency Access	Permissions to send notifications through SMN when a job fails to be executed

# 7.2 Creating a Custom DLI Agency

When Flink 1.15, Spark 3.3, or a later version is used to execute jobs and the required agency is not included in the DLI system agency **dli\_management\_agency**, you need to create an agency on the IAM console and add information about the new agency to the job configuration. **dli\_management\_agency** contains the permissions required for datasource operations, message notifications, and user authorization operations. For other agency permission requirements, you need to create custom DLI agencies. For details about **dli\_management\_agency**, see **DLI Agency Overview**.

This section walks you through how to create a custom agency, complete service authorization, and add information about the new agency to the job configuration.

# **DLI Custom Agency Scenarios**

Scenario	Agency Name	Description	Permission Policy
Allowing DLI to clear data according to the lifecycle of a table	dli_data_cl ean_agenc y	Data cleanup agency, which is used to clean up data according to the lifecycle of a table and clean up lakehouse table data. You need to create an agency and customize permissions for it. However, the agency name is fixed to <b>dli_data_clean_agency</b> .	Data Cleanup Agency Permission Configuratio n
Allowing DLI to read and write data from and to OBS to transfer logs	Custom	For DLI Flink jobs, the permissions include downloading OBS objects, obtaining OBS/ GaussDB(DWS) data sources (foreign tables), transferring logs, using savepoints, and enabling checkpointing. For DLI Spark jobs, the permissions allow downloading OBS objects and reading/writing OBS foreign tables.	Permission Policies for Accessing and Using OBS
Allowing DLI to obtain data access credentials by accessing DEW	Custom	DLI jobs use DEW-CSMS' secret management.	Permission to Use DEW's Encryption Function
Allowing DLI to access DLI catalogs to retrieve metadata	Custom	DLI accesses catalogs to retrieve metadata.	Permission to Access DLI Catalog Metadata
Allowing DLI to access LakeFormation catalogs to retrieve metadata	Custom	DLI accesses LakeFormation catalogs to retrieve metadata.	Permission to Access LakeFormati on Catalog Metadata

Гable	7-3	DLI	custom	agency	scenarios

# Procedure

Figure 7-2 Process of creating a custom agency



# Notes and Constraints

- The custom agency name cannot be the default agency names dli\_admin\_agency, dli\_management\_agency, or dli\_data\_clean\_agency. It must be unique.
- The name of the agency that allows DLI to clear data according to the lifecycle of a table must be **dli\_data\_clean\_agency**.
- Custom agencies can be configured only for jobs executed by engines like Flink 1.15, Spark 3.3.1 (Spark general queue scenario), or later versions.
- Once the agency permissions are updated, your dli\_admin\_agency will be upgraded to dli\_management\_agency. This new agency will have the necessary permissions for datasource operations, notifications, and user authorization operations. If you have other agency permission needs, you will need to create custom agencies. For details about dli\_management\_agency, see DLI Agency Overview.
- Common scenarios for creating an agency: DLI is allowed to read and write data from and to OBS, dump logs, and read and write Flink checkpoints. DLI is allowed to access DEW to obtain data access credentials and access catalogs to obtain metadata. For details about agency permissions in these scenarios, see Agency Permission Policies in Common Scenarios.

# Step 1: Create a Cloud Service Agency on the IAM Console and Grant Permissions

- 1. Log in to the management console.
- 2. In the upper right corner of the page, hover over the username and select **Identity and Access Management**.
- 3. In the navigation pane of the IAM console, choose Agencies.
- 4. On the displayed page, click **Create Agency**.
- 5. On the **Create Agency** page, set the following parameters:
  - Agency Name: Enter an agency name, for example, dli\_obs\_agency\_access.
  - Agency Type: Select Cloud service.
  - Cloud Service: This parameter is available only when you select Cloud service for Agency Type. Select Data Lake Insight (DLI) from the dropdown list.
  - Validity Period: Select Unlimited.

# - **Description**: You can enter **Agency with OBS OperateAccess permissions**. This parameter is optional.

Figure 7-3 Creating an agency

* Agency Name	dli_obs_agency_access
* Agency Type	<ul> <li>Account         Delegate another Huawei Cloud account to perform operations on your resources.     </li> <li>Cloud service         Delegate a cloud service to access your resources in other cloud services.     </li> </ul>
* Cloud Service	Data Lake Insight (DLI)
* Validity Period	Unlimited •
Description	OBS OperateAccess
	Next Cancel

- 6. Click Next.
- 7. Click the agency name. On the displayed page, click the **Permissions** tab. Click **Authorize**. On the displayed page, click **Create Policy**.
- 8. Configure policy information.
  - a. Enter a policy name, for example, **dli-obs-agency**.
  - b. Select JSON.
  - c. In the **Policy Content** area, paste a custom policy.

In this example, the permissions allow access and usage of OBS in various scenarios. For DLI Flink jobs, this includes downloading OBS objects, obtaining OBS/GaussDB(DWS) data sources (foreign tables), transferring logs, using savepoints, and enabling checkpointing. For DLI Spark jobs, the permissions allow downloading OBS objects and reading/writing OBS foreign tables.

For how to configure common agency permissions for Flink jobs, see **Agency Permission Policies in Common Scenarios**.

```
"Version": "1.1",
"Statement": [
{
"Effect": "Allow",
"Action": [
"obs:bucket:GetBucketPolicy",
"obs:bucket:GetLifecycleConfiguration",
"obs:bucket:GetBucketLocation",
```



- d. Enter a policy description as required.
- 9. Click Next.
- 10. On the **Select Policy/Role** page, select **Custom policy** from the first dropdown list and select the custom policy created in **8**.

#### Figure 7-4 Selecting the created custom policy

Assign selec	teo permissions to oil_obs_agency_access.			
View Sel	ected (1) Copy Permissions from Another Project		Custom policy	٣
•	Policy/Role Name	Туре		
•	palcyvdd834 	Custom policy		
<b>⊻</b> ~	di-dos-agency 	Custom policy		

11. Click Next. On the Select Scope page, set the authorization scope.

For details about authorization operations, see **Creating a User Group and Assigning Permissions**.

- All resources: IAM users will be able to use all resources, including those in enterprise projects, region-specific projects, and global services under your account based on assigned permissions.
- Global services: Global services are deployed in all regions, so you can access them seamlessly without having to switch between regions. However, global services cannot be authorized based on regional projects. They include services like Object Storage Service (OBS) and Content Delivery Network (CDN). After authorization, you can use the selected services based on your permissions.
- Region-specific projects: IAM users will be able to use resources in the selected region-specific projects based on assigned permissions.

 Enterprise projects: IAM users will be able to use resources in the selected enterprise projects based on assigned permissions. For example, an enterprise project may contain resources that are deployed in different regions. After you associate the enterprise project with the IAM users, they can access the resources in this enterprise project based on the assigned permissions.

In this example, the custom policy is an OBS agency. So, select **Global services**. If a DLI agency is used, you are advised to select **Region-specific projects**.

12. Click OK.

It takes 15 to 30 minutes for the authorization to be in effect.

## Step 2: Set Agency Permissions for a Job

When Spark 3.3.1, Flink 1.15, or a later version is used to execute jobs, you need to add information about the new agency to the job configuration.

Otherwise, If you do not specify an agency for Spark 3.3.1 jobs, the jobs cannot use OBS. If you do not specify an agency for a Flink 1.15 job, checkpointing cannot be enabled, savepoints cannot be used, logs cannot be transferred, and data sources such as OBS and GaussDB(DWS) cannot be used.

#### 

- You can only specify an agency for Flink 1.15 and Spark 3.3.1 jobs running on queues in an elastic resource pool.
- After specifying an agency for a job, be careful when modifying the permissions granted to the agency. Any changes made may impact the job's normal operation.
- Specifying an agency for a Flink Jar job
  - a. Log in to the DLI console. In the navigation pane, choose **Job Management** > **Flink Jobs**.
  - b. Select a desired job and click Edit in the Operation column.
  - c. In the job configuration area on the right, configure agency information.
    - Flink Version: Select 1.15.
    - Agency: Select the agency created in Step 1: Create a Cloud Service Agency on the IAM Console and Grant Permissions.

In this example, set it to **dli\_obs\_agency\_access**.

* Queue	~	
* Flink Version	( 1.15 v	
Catalog Name		
* Application ⑦		View Built-in Dependencies
Main Class	Default Manually assign	
* Class Name		
Class Arguments	Enter a class argument (Separate multiple class arguments with spaces).	
JAR Package Dependencies	Enter OBS paths, separating each path by pressing Enter.	
Other Dependencies (?)	Enter OBS paths, separating each path by pressing Enter.	View Built-in Dependencies
Јор Туре	Basic Image	
Agency	dli_obs_agency_access ~	

Figure 7-5 Specifying an agency for a Flink Jar job

- Specifying an agency for a Flink OpenSource SQL job
  - a. Log in to the DLI console. In the navigation pane, choose **Job Management** > **Flink Jobs**.
  - b. Select a desired job and click **Edit** in the **Operation** column.
  - c. In the job configuration area on the right, configure agency information.
    - On the **Running Parameters** tab, ensure that the selected Flink version is **1.15**.
    - Agency: Select the agency created in Step 1: Create a Cloud Service Agency on the IAM Console and Grant Permissions.

In this example, set it to **dli\_obs\_agency\_access**.

* Queue		Running
* Flink Version	( 1.15 V )	Para
UDF Jar	Enter OBS paths, separating each path by pressing Enter.	neters Runtime Conf
Catalog Name	~ (e	D iguration
Agency	dli_obs_agency_access ~	

Figure 7-6 Specifying an agency for a Flink OpenSource SQL job

- Specifying an agency for a Spark job
  - a. Log in to the DLI console. In the navigation pane, choose **Job Management** > **Spark Jobs**.
  - b. Select the target job and click **Edit** in the **Operation** column.
  - c. In the job configuration area on the right, configure agency information.
    - Spark Version: Make sure to select 3.3.1.
    - Agency: Select the agency created in Step 1: Create a Cloud Service Agency on the IAM Console and Grant Permissions.

In this example, set it to **dli\_obs\_agency\_access**.

Figure 7-7 Specifying an agency for a Spark job

<	Fill Form	Write API					
	Select a Que	ue					
	Queues					~	
	Spark Version		3.3.1			~	
	Application (	Configuration					
	Application					Ð	
	Agency		dli_obs_age	ncy_access		~	

# 7.3 Agency Permission Policies in Common Scenarios

This section provides agency permission policies for common scenarios, which can be used to configure agency permission policies when you customize your permissions. The "Resource" in the agency policy should be replaced according to specific needs.

## **Data Cleanup Agency Permission Configuration**

Application scenario: Data cleanup agency, which is used to clean up data according to the lifecycle of a table and clean up lakehouse table data. You need to create an agency and customize permissions for it. However, the agency name is fixed to **dli\_data\_clean\_agency**.

#### **NOTE**

Set the authorization scope of an agency as follows:

- For an OBS agency, select **Global services**.
- For a DLI agency, select Region-specific projects.

```
{
   "Version": "1.1",
   "Statement": [
     {
        "Effect": "Allow",
         "Action": [
           "obs:object:GetObject",
            "obs:object:DeleteObject",
            "obs:bucket:HeadBucket",
           "obs:bucket:ListBucket".
            "obs:object:PutObject"
        1
     }
  ]
}
   "Version": "1.1",
   "Statement": [
     {
        "Effect": "Allow",
        "Action": [
            "dli:table:showPartitions",
           "dli:table:select",
           "dli:table:dropTable",
            "dli:table:alterTableDropPartition"
        ]
     }
  ]
```

# Permission Policies for Accessing and Using OBS

Application scenario: For DLI Flink jobs, the permissions include downloading OBS objects, obtaining OBS/GaussDB(DWS) data sources (foreign tables), transferring logs, using savepoints, and enabling checkpointing. For DLI Spark jobs, the permissions allow downloading OBS objects and reading/writing OBS foreign tables.

```
"Version": "1.1",
"Statement": [
{
"Effect": "Allow",
"Action": [
"obs:bucket:GetBucketPolicy",
"obs:bucket:GetLifecycleConfiguration",
```



# Permission to Use DEW's Encryption Function

Application scenario: DLI Flink and Spark jobs use DEW-CSMS' secret management.



# Permission to Access DLI Catalog Metadata

Application scenario: DLI Flink and Spark jobs are authorized to access DLI metadata.

```
"Version": "1.1",
"Statement": [
{
"Effect": "Allow",
"Action": [
"dli:table:showPartitions",
"dli:table:alterTableAddPartition",
"dli:table:alterTableAddColumns",
```

"dli:table:alterTableRenamePartition", "dli:table:delete", "dli:column:select", "dli:database:dropFunction", "dli:table:insertOverwriteTable", "dli:table:describeTable", "dli:database:explain", "dli:table:insertIntoTable", "dli:database:createDatabase", "dli:table:alterView", "dli:table:showCreateTable", "dli:table:alterTableRename", "dli:table:compaction", "dli:database:displayAllDatabases", "dli:database:dropDatabase", "dli:table:truncateTable", "dli:table:select", "dli:table:alterTableDropColumns", "dli:table:alterTableSetProperties", "dli:database:displayAllTables", "dli:database:createFunction", "dli:table:alterTableChangeColumn", "dli:database:describeFunction", "dli:table:showSegments", "dli:database:createView", "dli:database:createTable", "dli:table:showTableProperties", "dli:database:showFunctions", "dli:database:displayDatabase", "dli:table:alterTableRecoverPartition", "dli:table:dropTable", "dli:table:update", "dli:table:alterTableDropPartition" 1

# Permission to Access LakeFormation Catalog Metadata

} ] }

Application scenario: DLI Spark jobs are authorized to access LakeFormation metadata.



# 7.4 Example of Configuring DLI Agency Permissions in Typical Scenarios

Туре	Helpful Link	Description
Flink job	Flink OpenSource SQL Jobs Using DEW to Manage Access Credentials	Guideline for using DEW to manage and access credentials for Flink OpenSource SQL jobs. When writing the output data of Flink jobs to MySQL or GaussDB(DWS), set attributes such as the username and password in the connector.
	Flink Jar Jobs Using DEW to Acquire Access Credentials for Reading and Writing Data from and to OBS	Guideline for Flink Jar jobs to acquire an AK/SK to read and write data from and to OBS.
	Obtaining Temporary Credentials for Flink Job Agencies	DLI provides a common interface to obtain temporary credentials for Flink job agencies set by users during job launch. The interface encapsulates the obtained temporary credentials for the job agency in the <b>com.huaweicloud.sdk.core.auth.</b> <b>BasicCredentials</b> class.
		Guideline for obtaining a temporary credential for a Flink job agency.
Spark job	Spark Jar Jobs Using DEW to Acquire Access Credentials for Reading and Writing Data from and to OBS	Guideline for Spark Jar jobs to acquire an AK/SK to read and write data from and to OBS.
	Obtaining Temporary Credentials for Spark Job Agencies	Guideline for obtaining a temporary credential for a Spark Jar job agency.

 Table 7-4 Guidelines for configuring DLI agency permissions in specific scenarios

# 8 Submitting a SQL Job on the DLI Management Console

# 8.1 Creating and Submitting a SQL Job

# Introduction

DLI offers a SQL editor for executing data query operations using SQL statements.

DLI's SQL editor supports SQL:2003, is compatible with SparkSQL, and allows batch execution of SQL statements. Additionally, commonly used syntax in the job editing window is highlighted in various colors. Both single-line and multi-line comments are supported (starting with "--", followed by the comment). For more detailed syntax descriptions, refer to **Data Lake Insight SQL Syntax Reference**.

This section describes how to create and submit a SQL job using the DLI SQL editor.

### Notes

- Before submitting a SQL job, configure a DLI job bucket. The bucket is used to store temporary data generated by DLI, such as job logs. You cannot view job logs if you choose not to create the bucket.
  - For details about how to configure a job bucket, see Configuring a DLI
     Job Bucket. The job bucket name is set by default.
  - On the OBS management console, you can configure lifecycle rules for a bucket to automatically delete objects within it or change object storage classes on a regular basis. For details, see Configuring a Lifecycle Rule.

If you have enabled the function to save job results to a DLI job bucket for your SQL queue, make sure to configure the DLI job bucket before submitting SQL jobs. Failure to do so may result in SQL jobs not being submitted successfully. For details, refer to **How Do I Check if Job Result Saving to a DLI Job Bucket Is Enabled for a SQL Queue?** 

# Creating and Submitting a SQL Job Using the SQL Editor

1. Log in to the DLI management console. In the navigation pane on the left, choose **SQL Editor**.

#### **NOTE**

On the SQL editor page, the system prompts you to create an OBS bucket to store temporary data generated by DLI jobs. In the **Set Job Bucket** dialog box, click **Setting**. On the page displayed, click the edit button in the upper right corner of the job bucket card. In the displayed **Set Job Bucket** dialog box, enter the job bucket path and click **OK**.

2. Above the SQL job editing window, set the parameters required for running a SQL job, such as the queue and database. For how to set the parameters, refer to **Table 8-1**.

Button & Drop- Down List	Description
Engine	<ul> <li>SQL jobs support the Spark and HetuEngine engines.</li> <li>Spark is suitable for offline analysis.</li> <li>HetuEngine is suitable for interactive analysis.</li> <li>For more information about DLI engines, see DLI Compute</li> </ul>

#### Table 8-1 Setting SQL job parameters

Button & Drop- Down List	Description
Queues	Resource queue used to execute SQL jobs.
	A queue determines the compute resources accessible to a job during its operation within an elastic resource pool. Every queue is allocated with specific resources, known as CUs, whose configuration significantly impacts the job's performance and execution efficiency.
	Before submitting a job, assess its resource needs and select an appropriate queue.
	SQL jobs can only be executed on SQL queues.
	If no queue is available, you can create a queue or use the <b>default</b> queue.
	• For details about how to create a queue, see <b>Creating an</b> <b>Elastic Resource Pool and Creating Queues Within It</b> .
	• The <b>default</b> queue serves well for temporary or testing scenarios characterized by indeterminate data volumes or infrequent data processing needs.
Catalog	A data catalog is a metadata management object that can contain multiple databases.
	For more information about data catalogs, see Understanding Data Catalogs, Databases, and Tables.
	You can create and manage multiple catalogs in DLI to isolate different metadata.
Databases	Select a database from the drop-down list box.
	If no database is available, the <b>default</b> database is displayed.
	For how to create a database, see <b>Creating a Data Catalog</b> , <b>Database</b> , and <b>Table on the DLI Console</b> .
	If you have specified a database where tables are located in SQL statements, the database you choose here does not apply.
Settings	Add parameters and tags.
	<b>Parameter Settings</b> : Set parameters in key/value format for SQL jobs. For details, see <b>Data Lake Insight SQL Syntax Reference</b> .
	<b>Tags</b> : Assign key-value pairs as tags to a SQL job.

3. Create a database and a table.

Create them in advance by referring to **Creating a Data Catalog, Database, and Table on the DLI Console**. For example, create a table named **qw**.

4. In the SQL job editing window, enter the following SQL statement: SELECT \* FROM qw.qw LIMIT 10; Alternatively, you can double-click the table name **qw**. The query statement is automatically entered in the SQL job editing window.

DLI offers a range of SQL templates that come with use cases, code examples, and usage guides. You can also use these templates to quickly implement your service logic. For more information about templates, see **Creating a SQL Job Template**.

- 5. On top of the editing window, click **More** > **Verify Syntax** to check whether the SQL statement is correct.
  - a. If the verification fails, check the SQL statement syntax by referring to **Data Lake Insight SQL Syntax Reference**.
  - b. If the syntax verification is successful, click **Execute**. Read and agree to the privacy agreement. Click **OK** to execute the SQL statement.

Once successfully executed, you can check the execution result on the **View Result** tab below the SQL job editing window.

6. View job execution results.

On the **View Result** tab, click <sup>[]</sup> to display execution results in a chart. Click

to switch back to the table view.

You can view a maximum of 1,000 data records on this View Result tab. To

view more or full data, click  $\Box$  to export the data to OBS.

**NOTE** 

- If no column of the numeric type is displayed in the execution result, the result cannot be represented in charts.
- You can view the data in a bar chart, line chart, or fan chart.
- In the bar chart and line chart, the X axis can be any column, while the Y axis can only be columns of the numeric type. The fan chart displays the corresponding legends and indicators.

### **Setting SQL Job Parameters**

Click **Settings** in the upper right corner of the **SQL Editor** page. You can set parameters and tags for the SQL job.

• Parameter Settings: Assign key-value pairs as parameter settings.

For details, see Data Lake Insight SQL Syntax Reference.

• **Tags**: Assign key-value pairs as tags to a SQL job.

Parameter	Default Value	Description
spark.sql.files.maxRec ordsPerFile	0	Maximum number of records to be written into a single file. If the value is zero or negative, there is no limit.

#### Table 8-2 Parameters for SQL job running

Parameter	Default Value	Description
spark.sql.autoBroadca stJoinThreshold	20971520 0	Maximum size, in bytes, of the table that displays all working nodes when a connection is executed. You can set this parameter to -1 to disable the display. NOTE Currently, only configuration units that store tables analyzed using the ANALYZE TABLE COMPUTE statistics noscan command and file-based data source tables that calculate statistics directly from data files are supported.
spark.sql.shuffle.partit ions	200	Default number of partitions used to filter data for join or aggregation.
spark.sql.dynamicPart itionOverwrite.enable d	false	When set to <b>false</b> , DLI will delete all partitions that meet the conditions before overwriting them. For example, if there is a partition named <b>2021-01</b> in a partitioned table and you use the <b>INSERT</b> <b>OVERWRITE</b> statement to write data to the <b>2021-02</b> partition, the data in the <b>2021-01</b> partition will also be overwritten. When set to <b>true</b> , DLI will not delete partitions in advance, but will overwrite partitions with data written during runtime.
spark.sql.files.maxPart itionBytes	13421772 8	Maximum number of bytes to be packed into a single partition when a file is read.
spark.sql.badRecordsP ath	-	Path of bad records.
dli.sql.sqlasync.enable d	true	Whether DDL and DCL statements are executed asynchronously. The value <b>true</b> indicates that asynchronous execution is enabled.
dli.sql.job.timeout	-	Job running timeout interval, in seconds. If the job times out, it will be canceled.

# More Common Functions of the SQL Editor

• Switching to the SparkUI page to view the SQL statement execution process

The SQL editor allows you to switch to the SparkUI to view the SQL statement execution process.

- You can view only the latest 100 job records on DLI's SparkUI.

 If a job is running on the **default** queue or is a synchronization one, you cannot switch to the SparkUI to view the SQL statement execution process.

#### **NOTE**

When you execute a job on a created queue, the cluster is restarted. It takes about 10 minutes. If you click **SparkUI** before the cluster is created, an empty **projectID** will be cached. The SparkUI page cannot be displayed. You are advised to use a dedicated queue so that the cluster will not be released. Alternatively, wait for a while after the job is submitted (the cluster is created), and then check **SparkUI**.

#### • Archiving SQL run logs

On the **Executed Queries (Last Day)** tab of the **SQL Editor** page, click **More** and select **View Log** in the **Operation** column of the SQL job. The system automatically switches to the OBS path where logs are stored. You can download logs as needed.

**NOTE** 

The **View Log** button is not available for synchronization jobs and jobs running on the **default** queue.

#### • SQL Editor shortcuts

Table 8-3 Keyboard shortcuts

Shortcut	Description
Ctrl+Enter	Execute SQL statements. You can run SQL statements by pressing <b>Ctrl+R</b> or <b>Ctrl + Enter</b> on the keyboard.
Ctrl+F	Search for SQL statements. You can press Ctrl+F to search for a required SQL statement.
Shift+Alt+F	Format SQL statements. You can press <b>Shift + Alt + F</b> to format a SQL statement.
Ctrl+Q	Syntax verification. You can press <b>Ctrl + Q</b> to verify the syntax of SQL statements.
F11	Full screen. You can press <b>F11</b> to display the SQL Job Editor window in full screen. Press <b>F11</b> again to leave the full screen.

# 8.2 Example of a Typical Scenario: Analyzing OBS Data Using a Spark SQL Job

DLI allows you to use data stored on OBS. You can create OBS tables on DLI to access and process data in your OBS bucket.

This section describes how to create an OBS table on DLI, import data to the table, and insert and query table data.

# Prerequisites

- You have created an OBS bucket. For details, see *Object Storage Service User Guide*. In this example, the OBS bucket name is **dli-test-021**.
- You have created a DLI SQL queue. For details, see Creating a Queue.
   Note: When you create the DLI queue, set Type to For SQL.

# Preparations

#### Creating a Database on DLI

- 1. Log in to the DLI management console and click **SQL Editor**. On the displayed page, set **Engine** to **spark** and **Queue** to the created SQL queue.
- Enter the following statement in the SQL editing window to create the testdb database. For details about the syntax for creating a DLI database, see Creating a Database. create database testdb;

The following operations in this section must be performed for the **testdb** database.

# DataSource and Hive Syntax for Creating an OBS Table on DLI

The main difference between DataSource syntax and Hive syntax lies in the range of table data storage formats supported and the number of partitions supported. For the key differences in creating OBS tables using these two syntax, refer to **Table 8-4**.

Synt ax	Data Types	Partitioning	Number of Partitions
Data Sourc e	ORC, PARQUET, JSON, CSV, and AVRO	You need to specify the partitioning column in both CREATE TABLE and PARTITIONED BY statements. For details, see <b>Creating a Single-</b> <b>Partition OBS Table Using</b> <b>DataSource Syntax</b> .	A maximum of 7,000 partitions can be created in a single table.
Hive	TEXTFILE, AVRO, ORC, SEQUENCEFILE, RCFILE, and PARQUET	Do not specify the partitioning column in the CREATE TABLE statement. Specify the column name and data type in the PARTITIONED BY statement. For details, see <b>Creating an OBS Table</b> <b>Using Hive Syntax</b> .	A maximum of 100,000 partitions can be created in a single table.

Table 8-4 Syntax differences
------------------------------

For details about the DataSource syntax, see **Creating an OBS Table Using the DataSource Syntax**.

For details about the Hive syntax, see **Creating an OBS Table Using the Hive Syntax**.

### Creating an OBS Table Using the DataSource Syntax

The following describes how to create an OBS table for CSV files. The methods of creating OBS tables for other file formats are similar.

- Create a non-partitioned OBS table.
  - Specify an OBS file and create an OBS table for the CSV data.
    - i. Create the **test.csv** file containing the following content and upload the **test.csv** file to the **root** directory of OBS bucket **dli-test-021**: Jordon,88,23 Kim,87,25 Henry,76,26
    - Log in to the DLI management console and choose SQL Editor from the navigation pane on the left. In the SQL editing window, set Engine to spark, Queue to the SQL queue you have created, and Database to testdb. Run the following statement to create an OBS table:

CREATE TABLE testcsvdatasource (name STRING, score DOUBLE, classNo INT ) USING csv OPTIONS (path "obs://dli-test-021/test.csv");

#### 

If you create an OBS table using a specified file, you cannot insert data to the table with DLI. The OBS file content is synchronized with the table data.

iii. Run the following statement to query data in the **testcsvdatasource** table.

select \* from testcsvdatasource;

#### Figure 8-1 Query results

	Engine spark +	Queues	testdi +	Databases	testdb +	• Execute	Format	Refer Syntax	Settings	More +
1CREATE TABLE testcsvdatasource (name STRING, score DOUBLE, classio I) 2	(T) USING csv OPTIONS (path	obs://dli	-test-021/test.csv	51						
3 select * from testcsvdatasource; 4										
5										
Line 3, Column 1			-			Execute: Ctrl+E	inter, Find: Chi+F, F	ormat Shift+Alt+F, Verify	Syntax: Chi+Q, F	ullscreen: F11
Executed Queries (Last Day) View Result										Clear All
Result1 🔕										
Executed successfully										
Query select * from testcsvdatasource										
Job ID bb4de710-48fb-4f94-8781-7b73e9b14ceb										
The query takes 4.00s, and 0.04 KB scanned.A maximum of 1,000 records can be displa	yed.							Enter a keyword.	Q	: Ľ ±
name 4E	score ↓Ξ				dassNo 🕽	=				
Jordon	88				23					
Kim	87				25					
Henry	76				26					

iv. Open the **test.csv** file on the local PC, add **Aarn,98,20** to the file, and replace the original **test.csv** file in the OBS bucket.

Jordon,88,23
Kim,87,25
Henry,76,26
Aarn.98.20

v. In the DLI **SQL Editor**, query the **testcsvdatasource** table for **Aarn,98,20**. The result is displayed. select \* from testcsvdatasource;

#### Figure 8-2 Query results

	Engine spark • Queues testdi	Databases testds      •      OExecute     F	ormat Refer Syntax Settings More •
1 select * from testcsvdatasource;			
Line 1, Column 33		Execute: Ctrl+Enter, Find:	Chi+F, Format: Shift+Alt+F, Verity Syntax: Chi+Q, Fullscreen: F11
xecuted Queries (Last Day) View Result		-	Clear All
Result 0			
Executed successfully			
Query select * from testcsvdatasource			
Job ID 9717ae98-a928-497f-bf82-708572484a94			
The query takes 1.49s, and 0.05 KB scanned A maximum of 1,000 records can be disp	/ayed.		Enter a keyword. Q
name (Ξ	score ↓Ξ	classNo ↓Ξ	
Jordon	88	23	
Kim	87	25	
Henry	76	26	

- Specify an OBS directory and create an OBS table for CSV data.
  - The specified OBS data directory does not contain files you want to import to the table.
    - 1) Create the file directory **data** in the **root** directory of the OBS bucket **dli-test-021**.
    - 2) Log in to the DLI management console and click SQL Editor. On the displayed page, set Engine to spark, Queue to the created SQL queue, and Database to testdb. Run the following statement to create OBS table testcsvdata2source in the testdb database on DLI: CREATE TABLE testcsvdata2source (name STRING, score DOUBLE, classNo INT) USING csv OPTIONS (path "obs://dli-test-021/data");
    - 3) Run the following statement to insert table data: insert into testcsvdata2source VALUES('Aarn','98','20');
    - Run the following statement to query data in the testcsvdata2source table: select \* from testcsvdata2source;

#### Figure 8-3 Query results

	Engine spark • Queues testd	I • Databases testob • OExecute	Format Refer Syntax Settings More *
1 select * from testcsvdata2source;			
Line 1, Column 33		Execute: CM+Enter	; Find: Old+F; Format: Shill+All+F; Vetify Syntax: Old+O; Fullscreen: F11
Executed Queries (Last Day) View Result			Clear All
Result1 O			
Executed successfully			
Query select * from testcsvdata2source			
Job ID abda6f6c-1cfb-4e06-aee9-bcdf36e6e329			
The query takes 4.96s, and 0.01 KB scanned.A maximum of 1,000 records can be displayed	ed.		Enter a keyword. Q
name JE	score ↓Ξ	classNo ↓≡	
áam.	91	20	

5) Refresh the **obs://dli-test-021/data** directory of the OBS bucket and query the data. A CSV data file is generated, and the data is added to the file.

#### Figure 8-4 Query results

lých / data / <b>1563485e2-6412615</b> Ø								
Objects Deleted Objects Fragments	Objects Deleted Objects Fragments							
Upload Object Create Folder Duints More	•					Enter an object name prefix. Q		
Name	Storage Class	Size ⑦ J⊟	Encrypted	Restoration Status	Last Modified 🕥 🐙	Operation		
4- Back								
	Standard	0 byte	No	-	Apr 22, 2022 09:17:41 GMT+08:00	Download Share More +		
parl-00000-1988409ec7e5402007988727815691e7-c000.cev	Standard	13 bytes	No	-	Apr 22, 2022 09:17:40 GMT+00.00	Download   Share   More +		

- The specified OBS data directory contains files you want to import to the table.
  - Create file directory data2 in the root directory of the OBS bucket dli-test-021. Create the test.csv file with the following content and upload the file to the obs://dli-test-021/data2 directory: Jordon,88,23 Kim,87,25 Henry,76,26
  - 2) Log in to the DLI management console and click SQL Editor. On the displayed page, set Engine to spark, Queue to the created SQL queue, and Database to testdb. Run the following statement to create OBS table testcsvdata3source in the testdb database on DLI: CREATE TABLE testcsvdata3source (name STRING, score DOUBLE, classNo INT) USING csv OPTIONS (path "obs://dli-test-021/data2");
  - 3) Run the following statement to insert table data: insert into testcsvdata3source VALUES('Aarn','98','20');
  - Run the following statement to query data in the testcsvdata3source table: select \* from testcsvdata3source;

#### Figure 8-5 Query results

name 🚛	score ↓≣	dassNo ↓Ξ
Jordon	88	23
Kim	87	25
Henry	76	26
Aam	98	20

5) Refresh the **obs://dli-test-021/data2** directory of the OBS bucket and query the data. A CSV data file is generated, and the data is added to the file.

#### Figure 8-6 Query results

jets / 6078663/684592 Ø							
Objects Deleted Objects Fragments							
Upload Object Create Folder Deloto More	•					Enter an object name prefix. Q	
Name	Storage Class	Size ③ J⊟	Encrypted	Restoration Status	Last Modified 🛞 🐙	Operation	
4. Back							
	Standard	0 byte	No	-	Apr 22, 2022 10:31:44 GMT+00:00	Download   Share   More +	
part-00000-4370bba3a09e4015b2b3ada0M3ccF4-c000.csv	Standard	13 bytes	No	-	Apr 22, 2022 10:31:44 GMT+08.00	Download   Share   More +	

• Create an OBS partitioned table

- Create a single-partition OBS table
  - i. Create file directory **data3** in the **root** directory of the OBS bucket **dli-test-021**.
  - ii. Log in to the DLI management console and click SQL Editor. On the displayed page, set Engine to spark, Queue to the created SQL queue, and Database to testdb. Run the following statement to create OBS table testcsvdata4source using data in the specified OBS directory obs://dli-test-021/data3 and partition the table on the classNo column.

CREATE TABLE testcsvdata4source (name STRING, score DOUBLE, classNo INT) USING csv OPTIONS (path "obs://dli-test-021/data3") PARTITIONED BY (classNo);

iii. Create the classNo=25 directory in the obs://dli-test-021/data3 directory of the OBS bucket. Create the test.csv file based on the following file content and upload the file to the obs://dli-test-021/ data3/classNo=25 directory of the OBS bucket. Jordon,88,25 King 725

Kim,87,25 Henry,76,25

iv. Run the following statement in the SQL editor to add the partition data to OBS table **testcsvdata4source**:

PARTITION (classNo = 25) LOCATION 'obs://dli-test-021/data3/classNo=25';

 Run the following statement to query data in the classNo=25 partition of the testcsvdata4source table: select \* from testcsvdata4source where classNo = 25;

#### Figure 8-7 Query results

name J≣	score ↓Ξ	classNo ↓Ξ
Jordon	88	25
Kim	87	25
Henry	76	25

vi. Run the following statement to insert the following data to the **testcsvdata4source** table:

insert into testcsvdata4source VALUES('Aarn','98','25'); insert into testcsvdata4source VALUES('Adam','68','24');

vii. Run the following statement to query data in the **classNo=25** and **classNo=24** partitions of the **testcsvdata4source** table:

### 

When a partitioned table is queried using the where condition, the partition must be specified. Otherwise, the query fails and "DLI.0005: There should be at least one partition pruning predicate on partitioned table" is reported.

ALTER TABLE testcsvdata4source ADD

select \* from testcsvdata4source where classNo = 25;

#### Figure 8-8 Query results

Jordon         88         25           Kim         87         25           Heny         76         25           Aam         98         25	name ↓Ξ	score ↓Ξ	classNo ↓Ξ
Kim         87         25           Henry         76         25           Aarn         98         25	Jordon	88	25
Henry         76         25           Aam         98         25	Kim	87	25
Aarn 98 25	Henry	76	25
	Aam	98	25

select \* from testcsvdata4source where classNo = 24;

#### Figure 8-9 Query results

name ↓Ξ	score ↓Ξ	classNo ↓Ξ
Adam	68	24

viii. In the **obs://dli-test-021/data3** directory of the OBS bucket, click the refresh button. Partition files are generated in the directory for storing the newly inserted table data.

#### Figure 8-10 classNo=25 file on OBS

Objects / data3 / 1c08c2797a4a4aada6 / classNo+25 🗇	Sijelis / dista] / tells2797talvalandalis / elassilio-23 of							
Objects Deleted Objects Fragments								
Objects are basic units of data stange. In OBS, files and folders are treated You can use OBS Browser+ to move an object to any other folder in this bud Upload Object Create Folder Detitie More v	Optim are basic with of dial strongs in OSE. Bits and bolins are brotect as algobs. Any the type can be applied and macapet is a basic Lean news       Tacks and OBE Showshin and algobs that and the basic.       Optim an applied basic     Charter and point and point and the basic.       Optim an applied basic     Charter and point and							
Name JE	Storage Class 4	Size (=	Encrypted 4E	Restoration Status JE	Last Modified UF	Operation		
4. End								
part-60000-1c00c2797a4a4aada6b6cb257c961157.c000.csv	Standard	10 bytes	No		Apr 22, 2022 11:41:34 GMT+68:00	Download   Share   More +		

#### Figure 8-11 classNo=24 file on OBS

bjess / data/ / 58/20070208003860 / CaseNo-24 ()							
Objects Deleted Objects Fragments							
Objects are basic units of data storage. In OBS, files and folders are treated as You can use OBS Browser+ to move an object to any other folder in this bucket	Objects are basic with of data strongs. In OSE, this and folders are treated as adjects. Any the type can be optioned and managed in a basist. Learn more Via can set OSE Strongers in some an adject to any other taken in this basist.						
Upload Object Create Folder Deinie More +						Enter an object name prefix. Q	
Name 4E	Storage Class 4E	Size (E	Encrypted JE	Restoration Status JE	Last Modified 4F	Operation	
4 Bad							
part-40000-fclic2xeb732x843dtia50b0dxeb11327b1.c000.cev	Standard	10 bytes	No	-	Apr 22, 2022 11:43:58 GMT+08:00	Download   Share   More +	

- Create an OBS table partitioned on multiple columns.
  - i. Create file directory **data4** in the **root** directory of the OBS bucket **dli-test-021**.
  - ii. Log in to the DLI management console and click SQL Editor. On the displayed page, set Engine to spark, Queue to the created SQL queue, and Database to testdb. Run the following statement to create OBS table testcsvdata5source using data in the specified OBS directory obs://dli-test-021/data4 and partition the table on classNo and dt columns.

CREATE TABLE testcsvdata5source (name STRING, score DOUBLE, classNo INT, dt varchar(16)) USING csv OPTIONS (path "obs://dli-test-021/data4") PARTITIONED BY (classNo,dt);

- iii. Run the following statements to insert the following data into the testcsvdata5source table: insert into testcsvdata5source VALUES('Aarn','98','25','2021-07-27'); insert into testcsvdata5source VALUES('Adam','68','25','2021-07-28');
- iv. Run the following statement to query data in the classNo partition of the testcsvdata5source table: select \* from testcsvdata5source where classNo = 25;

#### Figure 8-12 Query results

name ↓≘	score 4⊟	classNo ↓⊟	dt ↓≘
Aam	98	25	2021-07-27
Adam	68	25	2021-07-28

v. Run the following statement to query data in the **dt** partition of the **testcsvdata5source** table:

select \* from testcsvdata5source where dt like '2021-07%';

#### Figure 8-13 Query results

name ↓≣	score ↓Ξ	classNo ↓≣	dt ↓≣
Aam	98	25	2021-07-27
Adam	68	25	2021-07-28

- vi. Refresh the **obs://dli-test-021/data4** directory of the OBS bucket. The following data files are generated:
  - File directory 1: obs://dli-test-021/data4/xxxxx/classNo=25/ dt=2021-07-27

#### Figure 8-14 Query results

Objecks / data / 1-6022/65489654cd9654 / classNar-25 / dx-20121.07.27 🗗							
Objects Deleted Objects Fragments							
Upload Object Create Folder Deicle More +						Enter an object name prefix. Q	
Name 4E	Storage Class JE	Size JE	Encrypted JE	Restoration Status 4E	Last Modified UF	Operation	
fn Back							
part-00000-e022a/5469/54bc8x548/4ecc86x8988.c000.csv	Standard	10 bytes	No	-	May 07, 2022 19:41:24 GMT+08:00	Download Share More +	

 File directory 2: obs://dli-test-021/data4/xxxxx/classNo=25/ dt=2021-07-28

#### Figure 8-15 Query results

cts / dats4 / fed53241bd5841cc3cts / daes04e=25 / <b>de=2621.07.28 (1</b>								
Deleted Objects Fragments								
Upload Object Create Folder Dukte More +						Enter an object name prefix. Q		
Name 4≣	Storage Class ↓⊞	Size JE	Encrypted 48	Restoration Status JE	Last Modified JF	Operation		
6 Back								
part-00000-fed63241bb5641cc3cb4d6703d0ecbd.c000.cav	Standard	10 bytes	No	-	May 07, 2022 19:41:24 GMT+08:00	Download   Share   More +		

vii. Create the partition directory classNo=24 in obs://dli-test-021/ data4, and then create the subdirectory dt=2021-07-29 in classNo=24. Create the test.csv file using the following file content and upload the file to the obs://dli-test-021/data4/classNo=24/ dt=2021-07-29 directory.

Jordon,88,24,2021-07-29 Kim,87,24,2021-07-29 Henry,76,24,2021-07-29

viii. Run the following statement in the SQL editor to add the partition data to OBS table **testcsvdata5source**: ALTER TABLE

testcsvdata5source ADD PARTITION (classNo = 24,dt='2021-07-29') LOCATION 'obs://dli-test-021/data4/ classNo=24/dt=2021-07-29';

ix. Run the following statement to query data in the **classNo** partition of the **testcsvdata5source** table:

select \* from testcsvdata5source where classNo = 24;

#### Figure 8-16 Query results

name ↓Ξ	score ↓⊟	classNo ↓⊟	dt ↓≣
Jordon	88	24	2021-07-29
Kim	87	24	2021-07-29
Henry	76	24	2021-07-29

x. Run the following statement to query all data in July 2021 in the **dt** partition:

select \* from testcsvdata5source where dt like '2021-07%';

#### Figure 8-17 Query results

name ↓Ξ	score ↓Ξ	classNo ↓⊟	dt ↓≘
Jordon	88	24	2021-07-29
Kim	87	24	2021-07-29
Henry	76	24	2021-07-29
Aam	98	25	2021-07-27
Adam	68	25	2021-07-28

## Creating an OBS Table Using Hive Syntax

The following describes how to create an OBS table for TEXTFILE files. The methods of creating OBS tables for other file formats are similar.

- Create a non-partitioned OBS table.
  - a. Create file directory data5 in the root directory of the OBS bucket dlitest-021. Create the test.txt file based on the following file content and upload the file to the obs://dli-test-021/data5 directory: Jordon,88,23 Kim,87,25 Henry,76,26
  - b. Log in to the DLI management console and click SQL Editor. On the displayed page, set Engine to spark, Queue to the created SQL queue, and Database to testdb. Run the following Hive statement to create an OBS table using data in obs://dli-test-021/data5/test.txt and set the row data delimiter to commas (,):

CREATE TABLE hiveobstable (name STRING, score DOUBLE, classNo INT) STORED AS TEXTFILE LOCATION 'obs://dli-test-021/data5' ROW FORMAT DELIMITED FIELDS TERMINATED BY ',';

#### **NOTE**

**ROW FORMAT DELIMITED FIELDS TERMINATED BY ','** indicates that records are separated by commas (,).

c. Run the following statement to query data in the **hiveobstable** table: select \* from hiveobstable;

#### Figure 8-18 Query results

name ↓Ξ	score ↓Ξ	classNo ↓Ξ
Jordon	88	23
Kim	87	25
Henry	76	26

- d. Run the following statements to insert data into the table: insert into hiveobstable VALUES('Aarn','98','25'); insert into hiveobstable VALUES('Adam','68','25');
- e. Run the following statement to query data in the table to verify that the data has been inserted: select \* from hiveobstable;

#### Figure 8-19 Query results

name ↓Ξ	score ↓Ξ	classNo ↓Ξ
Adam	68	25
Aam	98	25
Jordon	88	23
Kim	87	25
Henry	76	26

f. In the **obs://dli-test-021/data5** directory, refresh the page and query the data. Two files are generated containing the newly inserted data.

#### Figure 8-20 Query results

Objects / data5 🗇						
Objects Deleted Objects Fragments						
Upload Object Create Folder Dutots More +						Enter an object name prefix. Q
Name JE	Storage Class JE	Size (1)	Encrypted J≣	Restoration Status JE	Last Modified UF	Operation
← Back						
parl-00000-27dcf27c-6750-4254-9f16-dac973ed133e-c000	Standard	13 bytes	No	-	Apr 22, 2022 15:59:49 GMT+08:00	Download Share More +
part-00000-a86a9d1c-44c1-4115-9645-cbb7b7041a78-c000	Standard	13 bytes	No	-	Apr 22, 2022 15:59:39 GMT+00:00	Download   Share   More +

#### Create an OBS Table Containing Data of Multiple Formats

- a. Create file directory data6 in the root directory of the OBS bucket dlitest-021. Create the test.txt file based on the following file content and upload the file to the obs://dli-test-021/data6 directory: Jordon,88-22,23:21 Kim,87-22,25:22 Henry,76-22,26:23
- b. Log in to the DLI management console and click SQL Editor. On the displayed page, set Engine to spark, Queue to the created SQL queue, and Database to testdb. Run the following Hive statement to create an OBS table using data stored in obs://dli-test-021/data6.
   CREATE TABLE hiveobstable2 (name STRING, hobbies ARRAY<string>, address map<string,string>) STORED AS TEXTFILE LOCATION 'obs://dli-test-021/data6' ROW FORMAT DELIMITED FIELDS TERMINATED BY ',' COLLECTION ITEMS TERMINATED BY '-' MAP KEYS TERMINATED BY ':';

#### 

- **ROW FORMAT DELIMITED FIELDS TERMINATED BY ','** indicates that records are separated by commas (,).
- **COLLECTION ITEMS TERMINATED BY '-'** indicates that the second column **hobbies** is in array format. Elements are separated by hyphens (-).
- MAP KEYS TERMINATED BY ':' indicates that the address column is in the key-value format. Key-value pairs are separated by colons (:).
- c. Run the following statement to query data in the **hiveobstable2** table: select \* from hiveobstable2;

#### Figure 8-21 Query results

name J≘	hobbies ↓Ξ	address ↓Ξ
Jordon	["88-22"]	{"23":"21"}
Kim	['87-22']	{"25":"22"}
Henry	['76-22']	{"26":"23"}

• Create a partitioned OBS table.

- Create file directory data7 in the root directory of the OBS bucket dlia. test-021.
- Log in to the DLI management console and click **SQL Editor**. On the b. displayed page, set Engine to spark, Queue to the created SQL queue, and **Database** to **testdb**. Run the following statement to create an OBS table using data stored in obs://dli-test-021/data7 and partition the table on the classNo column:

CREATE TABLE IF NOT EXISTS hiveobstable3(name STRING, score DOUBLE) PARTITIONED BY (classNo INT) STORED AS TEXTFILE LOCATION 'obs://dli-test-021/data7' ROW FORMAT DELIMITED FIELDS TERMINATED BY ',';

# A CAUTION

You can specify the partition key in the **PARTITIONED BY** statement. Do not specify the partition key in the CREATE TABLE IF NOT EXISTS statement. The following is an incorrect example:

CREATE TABLE IF NOT EXISTS hiveobstable3(name STRING, score DOUBLE, classNo INT) PARTITIONED BY (classNo) STORED AS TEXTFILE LOCATION 'obs://dli-test-021/data7';

- Run the following statements to insert data into the table: C. insert into hiveobstable3 VALUES('Aarn','98','25'); insert into hiveobstable3 VALUES('Adam','68','25');
- Run the following statement to query data in the table: d. select \* from hiveobstable3 where classNo = 25;

Figure 8-22 Query results

name ↓Ξ	score ↓Ξ	classNo ↓Ξ
Adam	68	25
Aarn	98	25

Refresh the **obs://dli-test-021/data7** directory. A new partition directory e. **classno=25** is generated containing the newly inserted table data.

#### Figure 8-23 Query results

Objects / data7 / classno=25 🗗					
Objects Deleted Objects Fragments					
Upload Object Create Folder Delete More +					Enter an object name prefix. Q
Nome JE Storage Class JE	Size JE	Encrypted JE	Restoration Status JE	Last Modified 47	Operation
≪ Back					
parl.00001-4216628-8866-4et3-ab63-5e706041e447.c000 Standard	10 bytes	No	-	Apr 22, 2022 19:50:08 GMT+08:00	Download Share More •
parl-0000-63777694-6203-4619-a223-2244cb29387.c000 Standard	10 bytes	No	-	Apr 22, 2022 16:49:59 GMT+08:00	Download   Share   More +

- f. Create partition directory classno=24 in the obs://dli-test-021/data7 directory. Create the **test.txt** file using the following file content and upload the file to the **obs://dli-test-021/data7/classno=24** directory: Jordon,88,24 Kim,87,24 Henry,76,24
- Run the following statement in the SQL editor to add the partition data q. to OBS table hiveobstable3: ALTER TABLE hiveobstable3 ADD

PARTITION (classNo = 24) LOCATION 'obs://dli-test-021/data7/classNo=24';

h. Run the following statement to query data in the **hiveobstable3** table: select \* from hiveobstable3 where classNo = 24;

#### Figure 8-24 Query results

name ↓Ξ	score 4≡	classNo ↓Ξ
Jordon	88	24
Kim	87	24
Henry	76	24

## FAQs

• **Q1**: What should I do if the following error is reported when the OBS partition table is queried?

DLI.0005: There should be at least one partition pruning predicate on partitioned table `xxxx`.`xxxx`.;

**Cause**: The partition key is not specified in the query statement of a partitioned table.

Solution: Ensure that the where condition contains at least one partition key.

• **Q2**: What should I do if "DLI.0007: The output path is a file, don't support INSERT...SELECT error" is reported when I use a DataSource statement to insert data in a specified OBS directory into an OBS table and the execution fails?

The statement is similar to the following: CREATE TABLE testcsvdatasource (name string, id int) USING csv OPTIONS (path "obs://dli-test-021/ data/test.csv");

**Cause**: Data cannot be inserted if a specific file is used in the table creation statement. For example, the OBS file **obs://dli-test-021/data/test.csv** is used in the preceding example.

**Solution**: Replace the OBS file to the file directory. You can insert data using the INSERT statement. The preceding example statement can be modified as follows:

CREATE TABLE testcsvdatasource (name string, id int) USING csv OPTIONS (path "obs://dli-test-021/ data");

• Q3: What should I do if the syntax of a Hive statement used to create a partitioned OBS table is incorrect? For example, the following statement creates an OBS table partitioned on **classNo**: CREATE TABLE IF NOT EXISTS testtable(name STRING, score DOUBLE, classNo INT) PARTITIONED BY (classNo) STORED AS TEXTFILE LOCATION 'obs://dli-test-021/data7';

**Cause**: Do not specify the partition key in the list following the table name. Specify the partition key in the **PARTITIONED BY** statement.

**Solution**: Specify the partition key in **PARTITIONED BY**. For example: CREATE TABLE IF NOT EXISTS testtable(name STRING, score DOUBLE) PARTITIONED BY (classNo INT) STORED AS TEXTFILE LOCATION 'obs://dli-test-021/data7';

# 8.3 Exporting SQL Job Results

Store the data results of analyzed SQL jobs in a specified location in the desired format.

By default, DLI stores SQL job results in its job bucket. You can also download job results to a local host or export job results to a specified OBS bucket.

# Exporting Job Results to the DLI Job Bucket

DLI specifies a default OBS bucket for storing job results. You can configure the bucket information on the **Global Configuration** > **Project** page of the DLI management console. Once a job is complete, the system automatically stores its results to this bucket.

The following conditions must be met if you want to read job results from the DLI job bucket:

- You have configured the job bucket on the Global Configuration > Project page of the DLI management console by referring to Configuring a DLI Job Bucket.
- You have submitted a service ticket to request the whitelisting of the feature that allows writing job results to buckets.
- The user who executes jobs has been granted read and write permissions either on the job bucket or on the **jobs/result** path of the job bucket.

For details, see **Creating a Custom Bucket Policy**.

For how to obtain job results from the DLI job bucket, see "Object Management" > **Downloading Objects** in *Object Storage Service User Guide*.

# Exporting Job Results to a Specified Location in Another Bucket

In addition to storing job results in the default bucket, you can also export them to a specified location in another bucket, increasing the flexibility of job result management and making it easier to organize and manage them.

On the console, you can only view a maximum of 1,000 job results. To view additional results, you can export them to an OBS path. The procedure is as follows:

You can export job results on either the **SQL Jobs** or the **SQL Editor** page.

- SQL Jobs page: In the navigation pane on the left, choose Job Management
   SQL Jobs. On the displayed page, locate the row containing a desired job, click More in the Operation column, and select Export Result.
- SQL Editor page: In the navigation pane on the left, choose SQL Editor. On

the displayed page, once query statements are successfully executed, click in next to the **View Result** tab to export job results.

#### D NOTE

- If there are no numerical columns in the query results, job results cannot be exported.
- Ensure that the user who exports job results has the read and write permissions on the OBS bucket.

Parameter	Mandator y	Description	
Data Format	Yes	Choose a data format for the job results you want to export. The options include <b>json</b> and <b>csv</b> .	
Queues	Yes	Select the queue where the job is executed. SQL jobs can only be executed on SQL queues.	
Compression Format	No	Compression format of the data to be exported. The options are:	
		• deflate	
Storage Path	Yes	Path in an OBS bucket where the job results are exported	
		<ul> <li>If Export Mode is set to New OBS directory, then You need to manually enter a directory name and ensure that the directory name does not exist. Otherwise, the system returns an error message and the export operation cannot be performed.</li> <li>NOTE The folder name cannot contain special characters (\/:*?"&lt;&gt; ) and cannot start or end with a period (.).</li> <li>For example, after selecting the storage path obs://bucket/src1/, you need to manually enter a directory name to change the path to obs://bucket/src1/ src2/ and ensure that the src2 directory name does not exist under src1.</li> <li>The job result export path is obs://</li> </ul>	
		<ul> <li>bucket/src1/src2/test.csv.</li> <li>If Export Mode is set to Existing OBS directory (Overwritten), then After selecting a bucket path, the job results are exported to that path. If there are files with the same name, they will be automatically overwritten.</li> <li>For example, if you select obs://bucket/</li> </ul>	
		<pre>src1/ as the bucket path, then The job result export path is obs:// bucket/src1/test.csv.</pre>	

Parameter	Mandator y	Description
Export Mode	Yes	• New OBS directory If you select this mode, a new folder path is created and the job results are saved to this path. This mode is used when you need to save exported results to a new location, making it easier to manage and track job results.
		If you select this option, you must manually enter an export directory in <b>Storage Path</b> and ensure that the directory must not exist. If the directory already exists, the system displays an error message and the export operation cannot be performed.
		<ul> <li>Existing OBS directory (Overwritten): When exporting job results, you can choose an existing file path as the output directory. If there is a file with the same name in that path, it will be automatically overwritten by the new exported job result file. This mode is used when you only need to save a single job result file in the same path, and you do not need to keep old job results.</li> </ul>
Number of Results	No	Number of results to be exported If you do not specify or set it to <b>0</b> , all results will be exported.
Table Header	No	Whether the job results to be exported contain table headers

# **Exporting Job Results to a Local Host**

You can download the results of asynchronous DDL and QUERY statements to a local directory. By default, you can download a maximum of 1,000 data records to a local host.

The procedure is as follows:

 Locate the row containing a desired job whose asynchronous DDL or QUERY statement has been successfully executed, click More in the Operation column, and select Submit Download Request. In the displayed dialog box, click OK. After a few seconds, the Submit Download Request button would change to Download.

	9.			J				-1		
sqi	QL Jobs 🕐 😳 🖗									
					Date	Feb 21, 2022 14:43:45	Select a date and the	me. 🟥	Search by statement by default.	Q @ C
		Queues 🏹	Engine 🍞	Username	Туре 🍞	Status 🏹	Query	Duration 7	Created 7	Operation
	~	testsql	spark	-	QUERY	Finished	Q01: Price summary r	16.08s	Feb 22, 2022 09:35:58 GMT+08:00	Edit Spark U More 🔺
	~	default	spark		QUERY	Finished	Q01: Price summary r	14.28s	Feb 22, 2022 09:11:22 GMT+08:00	Re-execute
	~	default	spark	-	QUERY	Finished	Q01: Price summary r	47.04s	Feb 22, 2022 09:00:29 GMT+08:00	Download Request
										View Result
										Export Result
										Export Log

Figure 8-25 Selecting Submit Download Request

2. Click **Download** to download the results to your local host.

# 8.4 Creating a SQL Inspection Rule

# What Is SQL Inspection?

There are numerous SQL engines in the big data field, which bring diversity to the solutions but also expose some issues such as varying quality of SQL input statements, difficult SQL problem localization, and excessive resource consumption by large SQL statements.

Poor quality SQL can have unforeseeable impacts on data analysis platforms, affecting system performance or platform stability.

DLI offers this feature to allow you to create inspection rules for the Spark SQL engine. This helps prevent common issues like large or low-quality SQL statements by providing pre-event information, blocking, and in-event circuit breaker. You do not need to change your SQL submission method or syntax, so it is easy to implement without affecting service operations.

- You can configure SQL inspection rules in a visualized manner and also have the ability to query and modify these rules.
- During query execution and response, each SQL engine proactively inspects SQL statements based on the rules.
- Administrators can choose to display hints on, intercept, or block SQL statements. The system logs SQL inspection events in real time for SQL audit. O&M engineers can analyze the logs, assess the quality of SQL statements on the live network, identify potential risks, and take preventive measures.

This section describes how to create a SQL inspection rule to enhance SQL defense capabilities.

# Notes and Constraints of DLI SQL Inspection Rules

- SQL inspection is only supported by Spark 3.3.*x* or later.
- Only one inspection rule can be created for an action or a queue.
- Each rule can be associated with a maximum of 50 SQL queues.
- A maximum of 1,000 rules can be created for each project.
#### **Creating a SQL Inspection Rule**

You can create SQL inspection rules for specified SQL queues on the **SQL Inspector** page. The system will prompt, block, or perform circuit breaking on SQL requests that trigger the rules.

#### **NOTE**

When creating or modifying a SQL inspection rule, evaluate the appropriateness of enabling the rules and setting the threshold based on the service scenario. This will help avoid the negative impact of unreasonable inspection rules blocking or performing circuit breaking on relevant SQL requests on the services.

- 1. Log in to the DLI management console.
- 2. In the navigation pane on the left, choose **Global Configuration** > **SQL Inspector**.
- 3. On the displayed **SQL Inspector** page, click **Create Rule** in the upper right corner. In the **Create Rule** dialog box, set parameters based on the table below.

Parameter	Description				
Rule Name	Name of a SQL inspection rule				
System Rules	Select an inspection rule. For details about the system nspection rules supported by DLI, see <b>SQL Inspection</b> System Rules That DLI Supports.				
Queues	select the queues the rules are bound to.				
Description	Enter a rule description.				
Rule Action	Actions that the current SQL inspection rule supports.				
	SQL rules support the following types of actions:				
	• Info: Record logs and provide a hint for handling the SQL request. If the rule has parameters, you need to configure the threshold.				
	• <b>Block</b> : Intercept the SQL request that meets the rule. If the rule has parameters, you need to configure the threshold.				
	• <b>Circuit Breaker</b> : Perform circuit breaking on the SQL request that meets the inspection rules. If the rule has parameters, you need to configure the threshold.				

Table 8-6	Parameters	for	creating	a SOL	inspection	rule
	ruruncters	101	creating	u JQL	inspection	ruic

4. Click OK.

View the added inspection rule on the **SQL Inspector** page. The rule takes effect dynamically.

To modify a rule, click **Modify** in its **Operation** column.

#### SQL Inspection System Rules That DLI Supports

This part describes the system inspection rules supported by DLI. For details, see **Table 8-7**.

- Default system rules are SQL inspection rules automatically created by DLI when a queue is created. These rules are bound to the queue and cannot be deleted.
- Default system rules include Scan files number, Scan partitions number, Shuffle data(GB), Count(distinct) occurrences, and Not in<Subquery>.
- Only one inspection rule can be created for an action or a queue.
- For every supported action, a default system rule is created. For example, when a queue is created, a Scan files number rule is automatically created for both the Info and Block actions.
- Different engine versions support different inspection rules.

To view the engine version of a queue, choose **Resources** > **Queue Management** in the navigation pane on the left, select the queue, doubleclick the pane at the bottom of the page, and check the value of **Default Version**.

#### Figure 8-26 Viewing the engine version of a queue



Rule ID	Rule Name	Descript ion	Ty p e	A p lic ab le En gi ne	Action	Value	Def aul t Sys te m Rul e	Exampl e SQL Statem ent	Supp orted Engin e Versi on
dynami c_0001	Scan files number	Maximu m number of files to be scanned	D yn a m ic	Sp ar k H et uE ng in e	Info Blo ck	Value range: 1– 20000 00 Defaul t value: 20000 0	Yes	N/A	Spark 3.3.1
dynami c_0002	Scan partitio ns number	Maximu m number of partition s involved in the operatio ns (select, delete, update, and alter) that can be performe d on a table	D yn a m ic	Sp ar k	Info Blo ck	Value range: 1– 50000 0 Defaul t value: 5000	Yes	select * from Partition ed table	Spark 3.3.1
runnin g_0002	Memor y used(M B)	Peak memory usage of the SQL stateme nt	R u ni n g	Sp ar k	Circ uit Bre ake r	Value range: 1– 83886 08	No	N/A	Spark 3.3.1

 Table 8-7 System inspection rules supported by DLI

Rule ID	Rule Name	Descript ion	Ty p e	A p lic ab le En gi ne	Act ion	Value	Def aul t Sys te m Rul e	Exampl e SQL Statem ent	Supp orted Engin e Versi on
runnin g_0003	Run time(S)	Maximu m running duration of the SQL stateme nt	R u ni g	Sp ar k	Circ uit Bre ake r	Unit: second Value range: 1– 43200	No	N/A	Spark 3.3.1
runnin g_0004	Scan data(G B)	Maximu m amount of data to be scanned	R u ni n g	Sp ar k	Circ uit Bre ake r	Unit: GB Value range: 1– 10240	No	N/A	Spark 3.3.1
runnin g_0005	Shuffle data(G B)	Maximu m amount of data to be shuffled	R u ni g	Sp ar k	Circ uit Bre ake r	Unit: GB Value range: 1– 10240 Defaul t value: <b>2048</b>	Yes	N/A	Spark 3.3.1 Spark 2.4.5

Rule ID	Rule Name	Descript ion	Ty p e	A p lic ab le En gi ne	Action	Value	Def aul t Sys te m Rul e	Exampl e SQL Statem ent	Supp orted Engin e Versi on
static_0 001	Count( distinct ) occurre nces	Maximu m number of occurren ces of count(di stinct) in the SQL stateme nt	St at ic	Sp ar k	Info Blo ck	Value range: 1–100 Defaul t value: <b>10</b>	Yes	SELECT COUNT( DISTINC T deviceId ), COUNT( DISTINC T collDevi ceId) FROM table GROUP BY deviceN ame, collDevi ceName , collCurr entVersi on;	Spark 3.3.1

Rule ID	Rule Name	Descript ion	Ty p e	A p lic ab le En gi ne	Action	Value	Def aul t Sys te m Rul e	Exampl e SQL Statem ent	Supp orted Engin e Versi on
static_0 002	Not in <sub query&gt;</sub 	Check whether <b>not in</b> <b><subque< b=""> <b>ry&gt;</b> is used in the SQL stateme nt.</subque<></b>	St at ic	Sp ar k	Info Blo ck	Value range: Yes, No Defaul t value: <b>Yes</b>	Yes	SELECT * FROM Orders o WHERE Orders. Order_I D not in (Select Order_I D FROM HeldOrd ers h where h.order_i d = o.order_i d);	Spark 3.3.1
static_0 003	Join occurre nces	Maximu m number of joins in the SQL stateme nt	St at ic	Sp ar k	Info Blo ck	Value range: 1–50	No	SELECT name, text FROM table_1 JOIN table_2 ON table_1.I d = table_2.I d	Spark 3.3.1

Rule ID	Rule Name	Descript ion	Ty p e	A p lic ab le En gi ne	Action	Value	Def aul t Sys te m Rul e	Exampl e SQL Statem ent	Supp orted Engin e Versi on
static_0 004	Union occurre nces	Maximu m number of union all times in the SQL stateme nt	St at ic	Sp ar k	Info Blo ck	Value range: 1–100	No	select * from tables t1 union all select * from tables t2 union all select * from tables t3	Spark 3.3.1
static_0 005	Subque ry nesting layers	Maximu m number of subquery nesting layers	St at ic	Sp ar k	Info Blo ck	Value range: 1–20	No	select * from ( with temp1 as (select * from tables) select * from temp1);	Spark 3.3.1
static_0 006	Sql size(KB )	Maximu m size of a SQL file	St at ic	Sp ar k	Info Blo ck	Unit: KB Value range: 1– 1024	No	N/A	Spark 3.3.1

Rule ID	Rule Name	Descript ion	Ty p e	A p lic ab le En gi ne	Action	Value	Def aul t Sys te m Rul e	Exampl e SQL Statem ent	Supp orted Engin e Versi on
static_0 007	Cartesi an product	Limitatio n of Cartesia n products when multiple tables are being associate d	St at ic	Sp ar k	Info Blo ck	Value range: 0–1	No	select * from A,B;	Spark 3.3.1

# 8.5 Setting the Priority for a SQL Job

#### Scenario

In actual job running, it is necessary to prioritize and ensure the normal running of important and urgent tasks due to their varying levels of importance and urgency. This requires providing the necessary compute resources for their normal operations.

DLI offers a feature to set job priorities for each SQL job, which prioritizes the allocation of compute resources to higher priority jobs when resources are limited.

#### Notes

- Priorities cannot be set for jobs running in queues within an elastic resource pool of the basic edition.
- You can assign a priority level of 1 to 10 for each job, with a larger value indicating a higher priority. Compute resources are preferentially allocated to high-priority jobs. That is, if compute resources required for high-priority jobs are insufficient, compute resources for low-priority jobs are reduced.
- Jobs running on a SQL queue have a default priority level of 3.
- To change the priority for a job, you must first stop the job, change the priority level, and then submit the job for the modification to take effect.

#### Setting the Priority for a SQL Job

Click **Settings**. In the **Parameter Settings** area, configure the following parameter. *x* indicates the priority value. **spark.sql.dli.job.priority=x** 

- 1. Log in to the DLI management console.
- 2. In the navigation pane on the left, choose **Job Management** > **SQL Jobs**.
- 3. Locate the row containing the job for which you want to set the priority and click **Edit** in the **Operation** column.
- 4. Click **Settings**. In the **Parameter Settings** area, configure the **spark.sql.dli.job.priority** parameter.

Figure 8-27	Example	configuration	for	a SQ	L job
-------------	---------	---------------	-----	------	-------

Execute     Fe	ormat	Refer Syntax	Settings	More 👻
Parameter Se	ttings			🕼 Quick Links
spark.sql.dli	.job.priority	× 8		ច
+				
Tags				
Кеу		Value		Ū
+		^		

## 8.6 Querying Logs for SQL Jobs

#### Scenario

DLI job buckets are used to store temporary data generated during DLI job running, such as job logs and results.

This section describes how to configure a bucket for DLI jobs on the DLI console and obtain SQL job logs.

#### Notes

- To avoid disordered job results, do not use the OBS bucket configured for DLI jobs for any other purposes.
- DLI jobs must be set and modified by the main account as IAM users do not have required permissions.
- You cannot view the logs for DLI jobs before configuring a bucket.
- You can configure lifecycle rules to periodically delete objects from buckets or change storage classes of objects.
- Exercise caution when modifying the job bucket, as it may result in the inability to retrieve historical data.

#### Prerequisites

Before the configuration, create an OBS bucket or parallel file system (PFS). In big data scenarios, you are advised to create a PFS. PFS is a high-performance file system provided by OBS, with access latency in milliseconds. PFS can achieve a bandwidth performance of up to TB/s and millions of IOPS, which makes it ideal for processing high-performance computing (HPC) workloads.

For details about PFS, see "Parallel File System Feature Guide" in the *Object Storage Service User Guide*.

#### **Configuring a Bucket for DLI Jobs**

- 1. In the navigation pane of the DLI console, choose **Global Configuration** > **Project**.
- 2. On the **Project** page, click and next to **Job Bucket** to configure bucket information.

#### Figure 8-28 Project

Lake Insigl	Project	
view		
L Editor	Job Bu	ucket
b Management	→ Bucket	Name: rain
esources	✓ This buc logs and create th	cket is used to store temporary data generated by DLI, such as job d job results. Do not use this bucket for other purposes. If you do not this bucket you will not be able to view job loos. You can use the main
ata Management	<ul> <li>account</li> <li>permissi</li> </ul>	t to set and modify the bucket. Sub-users do not have modification sions. You can set a lifecycle rule to periodically delete objects in a
ob Templates	✓ bucket o lifecycle	or change its storage class. Exercise caution when you modify the e rule to prevent historical data being deleted by mistake.
atasource Connection		
Global Configuration	^	
Global Variables		
SQL Inspector	<	
Project		
Service Authorization		

- 3. Click 🗁 to view available buckets.
- 4. In the displayed **OBS** dialog box, click the name of a bucket or search for and click a bucket name and then click **OK**. In the **Set Job Bucket** dialog box, click **OK**.

Temporary data generated during DLI job running will be stored in the OBS bucket.

Figure 8-29 Setting the job bucket

Set Job Bu	cket	>
1 Parallel fil	e buckets are recommended.	×
★ Job Bucket	rain	
		Cancel OK

#### **Querying Logs for SQL Jobs**

- Log in to the DLI console. In the navigation pane on the left, choose Job Management > SQL Jobs.
- 2. Select the SQL job whose jobs you want to query, click **More** in the **Operation** column, and select **View Log**.

The system automatically switches to the log path of the DLI job bucket.

3. On the **Files** tab, select the log file of the desired date and time and click **Download** in the **Operation** column to download the file to your local host.

#### Figure 8-30 Downloading SQL job logs

les / jobs / logs / 1 / 2024-09-18_10-48-55 /								
Files Fragments								
You can use OBS Browser+ to mov Preview Objects in OBS from My Br Upload File Create F	You can use QBS Browser+ to move a file to any other folder in this parallel file system. For security reasons, files cannot be previewed online when you access them from a browser. To preview files online, see How Preview Opects in QBS from My Browser?							
C Enter a file name prefi	х.							
Name	Storage Class	Size	Last Modified	Operation				
	Standard	36.55 KB	Sep 18, 2024 14:59:00 GMT+08:00	Download Share More 🗸				
🕑 📑 spark.log	Standard	1.29 KB	Sep 18, 2024 14:34:00 GMT+08:00	Download Share More V				

## 8.7 Managing SQL Jobs

#### Viewing Basic Job Information on the SQL Jobs Page

The **SQL Jobs** page displays all SQL jobs, which may span multiple pages if there are many jobs. You can navigate to a specific page as needed. DLI allows you to view jobs in all statuses. Jobs in the job list are sorted by creation time in descending order by default.

Parameter	Description
Queues	Name of the queue to which a job belongs
Engine	<ul> <li>SQL jobs support the Spark and HetuEngine engines.</li> <li>Spark: displays jobs whose execution engine is Spark.</li> <li>HetuEngine: displays jobs whose execution engine is HetuEngine.</li> </ul>
Username	Name of the user who executed the job.

 Table 8-8 SQL Job management parameters

Parameter	Description
Туре	<ul> <li>Job type. The following types are supported:</li> <li>IMPORT: A job that imports data to DLI</li> <li>EXPORT: A job that exports data from DLI</li> <li>DCL: Conventional DCLs and operations related to queue permissions</li> <li>DDL:Conventional DDLs, including creating and deleting databases and tables</li> <li>QUERY: A job that queries data by running SQL statements</li> <li>INSERT: A job that inserts data by running SQL statements</li> <li>UPDATE: A job that updates data.</li> <li>DELETE: A job that deletes a SQL job.</li> </ul>
	<ul> <li>DATA_MIGRATION: A job that migrates data.</li> <li>RESTART_QUEUE: A job that restarts a queue.</li> <li>SCALE_QUEUE: A job that changes queue specifications, including sale-out and scale-in.</li> </ul>
Status	Job status. Possible values are as follows: • Submitting • Running • Finished • Canceled • Failed • Scaling
Query	SQL statements for operations such as exporting and creating tables You can click <sup>1</sup> to copy the query statement.
Duration	Running duration of a job
Created	Time when a job is created. Jobs can be displayed in ascending or descending order of the job creation time.

Parameter	Description
Operation	• Edit: Edit the job.
	Cancel
	<ul> <li>You can terminate a job only when the job is in Submitting or Running status.</li> </ul>
	<ul> <li>A job whose status is Finished, Failed, or Canceled cannot be terminated.</li> </ul>
	<ul> <li>If the Cancel button is gray, you are not allowed to perform this operation.</li> </ul>
	• <b>Re-execute</b> : Execute the job again.
	• <b>SparkUI</b> : Display the Spark job execution page.
	NOTE
	• When you execute a job on a created queue, the cluster is restarted. It takes about 10 minutes. If you click <b>SparkUI</b> before the cluster is created, an empty <b>projectID</b> will be cached. The SparkUI page cannot be displayed. You are advised to use a dedicated queue so that the cluster will not be released. Alternatively, wait for a while after the job is submitted (the cluster is created), and then check <b>SparkUI</b> .
	• SparkUI can currently only display the latest 100 jobs.
	• In addition to the preceding operations, the following operations are available for QUERY jobs and asynchronous DDL jobs.
	<ul> <li>Submit Download Request: Download the results of asynchronous DDL and QUERY statements to a local directory. For details, see Exporting Job Results to a Local Host.</li> </ul>
	<ul> <li>View Result: View the job running result.</li> </ul>
	<ul> <li>Export Result: Export the job running result to the created OBS bucket.</li> </ul>
	• In addition to the preceding operations, the EXPORT job also includes the following operations:
	– Download
	• <b>View Log</b> : Save job logs to the temporary OBS bucket created by DLI.
	<b>NOTE</b> The <b>View Log</b> button is not available for synchronization jobs and jobs running on the default queue.

#### Viewing Job Details

On the **SQL Jobs** page, you can click  $\checkmark$  in front of a job record to view details about the job.

Job details vary with job types. The job details vary depending on the job types, status, and configuration options. The following describes how to load data, create

a table, and select a job. For details about other job types, see the information on the management console.

- Load data (job type: IMPORT) include the following information: queue, job ID, username, type, status, execution statement, running duration, creation time, end time, parameter settings, label, number of results, scanned data, number of scanned data, number of error records, storage path, data format, database, table, table header, separator, reference character, escape character, date format, timestamp format, total CPU used, and output bytes.
- **Create table** (job type: DDL) include the following information: queue, job ID, username, type, status, execution statement, running duration, creation time, end time, parameter settings, tags, number of results, scanned data, and database.
- Select (job type: QUERY) include the following information: queue, job ID, username, type, status, execution statement, running duration, creation time, end time, parameter setting, label, number of results (results of successful executions can be exported), and scanned data, username, result status (results of successful tasks can be viewed. Failure causes of failed tasks are displayed), database, total CPU used, and output bytes.

#### **NOTE**

- Total CPU Used (Core x ms): total CPU used during job execution.
- **Output Bytes**: number of output bytes after the job is executed.

#### Searching for a Job

On the **SQL Jobs** page, you can search jobs with any of the following operations.

- Select a queue name.
- Select an engine.
- Set the date range.
- Enter a username, statement, tag, or job ID.
- Select the creation time in ascending or descending order.
- Select a job type.
- Select a job status.
- Select the job execution duration in ascending or descending order.

#### Terminating a SQL Job

On the **SQL Jobs** page, you can click **Terminate** in the **Operation** column to stop a submitting or running job.

## 8.8 Viewing a SQL Execution Plan

A SQL execution plan is a logical flowchart of a database query that shows how a database management system executes a specific SQL query. The execution plan details the steps needed to execute the query, such as table scans, index lookups, join operations (for example, inner join, outer join), sorting, and aggregation. Viewing an execution plan can help analyze query performance, identify potential performance bottlenecks, understand the query's execution logic, and use this

information to adjust the query or database structure to improve SQL query efficiency.

This section describes how to view a SQL execution plan on the DLI management console.

#### Notes and Constraints

- You can only view SQL execution plans for Spark 3.3.*x* or later queues and HetuEngine queues.
- You can only view a SQL execution plan after a SQL job is executed.
- You can only view the SQL execution plan for SQL jobs that have reached the **Finished** state.
- Make sure you have authorized DLI to use OBS buckets for saving the SQL execution plans of user jobs.
- SQL execution plans are stored in paid storage buckets for DLI jobs. The system does not automatically delete them. You are advised to configure the bucket lifecycle and specify rules to regularly delete or migrate unused SQL execution plans. Refer to Configuring a DLI Job Bucket.

#### Procedure

- 1. Log in to the DLI management console.
- 2. In the navigation pane on the left, choose **Job Management** > **SQL Jobs**.
- 3. Select the SQL job you want to query.
- 4. At the bottom of the page, click the name of the job you selected to view its details.

In the details area, click **Expand** next to **SQL Execution Plan**. The system queries the SQL execution plan of the job from the DLI job bucket and displays the plan on the console.

If the SQL execution plan in the DLI job bucket is deleted, the plan may not be displayed because the source file is missing.

Figure 8-31	Viewing	a SQL	execution	plar
-------------	---------	-------	-----------	------

	Cant	OUERY	Eloiohad	calast # from test, data, size	1000 0 160	Aug 16, 2024 12:21:08 CMT - 08:00
Sq.	Spark	QUERT	Pinished	select - from test_data_size	Imin U. Tos	Aug 16, 2024 13:21:06 GM1+06:00
	Spark	QUERY	Finished	select * from test_data_size	20.50s	Aug 16, 2024 12:58:11 GMT+08:00
	Spark	DDL	Canceled	CREATE EXTERNAL TABLE IF NOT EXISTS tpch.nation	3h 0 min4.25s	Aug 16, 2024 12:55:21 GMT+08:00
	Spark	QUERY	S Finished	select * from tpch.orders	3min 59.86s	Aug 16, 2024 11:31:02 GMT+08:00
d: 03d42c5f-1fe1-	-4d06-b250-5e41aad860ea					
d: 03d42c5f-1fe1-	-4d06-b250-5e41aad860ea					
d: 03d42c5f-1fe1-	-4006-b250-5e41aad860ea					
d: 03d42c5f-1fe1-	4d06-b250-5e41aad860ea					
d: 03d42c5f-1fe1-	-4d06-b250-5e41aad860ea					
d: 03d42c5f-1fe1-	-4d06-b250-5e41aad860ea				0	-

## 8.9 Creating and Managing SQL Job Templates

## 8.9.1 Creating a SQL Job Template

To facilitate SQL operation execution, DLI allows you to customize query templates or save the SQL statements in use as templates. After templates are saved, you do not need to compile SQL statements. Instead, you can directly perform the SQL operations using the templates.

SQL templates include sample templates and custom templates. The default sample template contains 22 standard TPC-H query statements, which can meet most TPC-H test requirements. For details, see **TPC-H Sample Data in the SQL Templates Preset on DLI**.

#### **NOTE**

In the navigation pane on the left, choose **Job Templates** > **SQL Templates**. In the upper right corner of the displayed page, click **Settings**. In the displayed **Settings** dialog box, choose whether to display templates by group.

If you enable **Display by Group**, the display options are **Expand the first group**, **Expand all**, and **Collapse all**.

#### Creating a SQL Job Template

You can create a template on either the **Job Templates** or the **SQL Editor** page.

- To create a template on the Job Templates page:
  - a. On the left of the management console, choose **Job Templates** > **SQL Templates**.
  - b. On the **SQL Templates** page, click **Create Template** to create a template.

Enter the template name, SQL statement, and description information. For details, see **Table 8-9**.

#### Figure 8-32 Creating a template

		-
Creeke	Toma	
Create	lem	Diate
cicace		Place

You can create 89 more templates. Increase quota.

* Name	Enter a template name.						
* Statement	Enter the SQL statem	ents for this template					
				11			
				0/10,000			
Description							
				0/256			
				-/			
Group	Use Existing	Use new	Do not use				
* Group Name	Select			•			
A Group Marrie	Detect						
	O	K Cancel					

#### Table 8-9 Parameter description

Parameter	Description
Name	Indicates the template name.
	<ul> <li>A template name can contain only digits, letters, and underscores (_), but cannot start with an underscore (_) or contain only digits. It cannot be left empty.</li> </ul>
	<ul> <li>The template name can contain a maximum of 50 characters.</li> </ul>
Statement	SQL statement to be saved as a template.
Description	Description of the template you create.
Group	Use existing
	Use new
	Do not use
Group Name	If you select <b>Use existing</b> or <b>Use new</b> , you need to enter the group name.

- c. Click **OK**.
- To create a template on the **SQL Editor** page:
  - a. On the left of the management console, click **SQL Editor**.
  - b. In the SQL job editing area of the displayed **SQL Editor** page, click **More** in the upper right corner and choose **Save as Template**.

Enter the template name, SQL statement, and description information. For details, see **Table 8-9**.

c. Click **OK**.

#### Submitting a SQL Job Using a Template

Perform the template operation as follows:

- 1. On the left of the management console, choose **Job Templates** > **SQL Templates**.
- 2. On the **SQL Templates** page, select a template and click **Execute** in the **Operation** column. The **SQL Editor** page is displayed, and the corresponding SQL statement is automatically entered in the SQL job editing window.
- 3. In the upper right corner of the SQL job editing window, Click **Execute** to run the SQL statement. After the execution is complete, you can view the execution result below the current SQL job editing window.

#### Searching for a SQL Job Template

On the **SQL Templates** page, you can enter the template name keyword in the search box on the upper right corner to search for the desired template.

#### Modifying a SQL Job Template

Only custom templates can be modified. To modify a template, perform the following steps:

- **Step 1** On the **SQL Templates** page, locate the target template and click **Modify** in the **Operation** column.
- **Step 2** In the displayed **Modify Template** dialog box, modify the template name, statement, and description as required.
- Step 3 Click OK.

----End

#### **Deleting a Template**

On the **SQL Templates** page, select one or more templates to be deleted and click **Delete** to delete the selected templates.

# 8.9.2 Developing and Submitting a SQL Job Using a SQL Job Template

DLI allows you to create custom templates or save currently used SQL statements as templates for quick and convenient SQL operations. Once a template is saved, you can execute SQL operations directly through the template without the need to write SQL statements.

The system offers sample templates that include various standard TPC-H query statements. You can choose to use one of these templates or create a custom template to create a SQL job.

This example shows how to use a TPC-H sample template to develop and submit a SQL job.

#### Procedure

- 1. Log in to the DLI management console.
- 2. In the navigation pane on the left, choose **Job Templates** > **SQL Templates**.
- 3. On the displayed **Sample Templates** tab, find a sample template that matches your service scenario under **tpchQuery** and click **Execute** in the **Operation** column.
- 4. In the upper part of the editing window, set **Engine** to **Spark**, **Queues** to **default**, and **Databases** to **default**, and click **Execute**.

	Engine Spark • Queues default • Databases default • OExecute Format Refer Syntax Settings More •
1	Q01: Price summary report query.
2	Quiry on a single table lineitem within a certain period of time, statistics on all kinds of commodities that have been paid and already shipped, including business volume billing, shipping, discounts,
3	Features: Single table query operations with grouping, sorting, and aggregation coexist. This query will read 95% to 97% of the data on the table.
4	SELECT 1_returnflag,
5	1_linestatus,
6	sum(l_quantity) AS sum_qty,
7	<pre>sum(1_extendedprice) AS sum_base_price,</pre>
8	<pre>sum(l_extendedprice * (1 - 1_discount)) AS sum_disc_price,</pre>
9	<pre>sum(1_extendedprice * (1 - 1_discount) * (1 + 1_tax)) AS sum_charge,</pre>
10	avg(l_quantity) AS avg_aty,
11	avg(l_extendedprice) AS avg_price,
12	avg(1_discount) AS avg_disc,
13	count(*) AS count_order
14	FRON tpch.lineitem
15	WHERE 1_shipdate <= DATE "1998-12-01" - INTERVAL "90" DAY
16	renn BV 1 satusfise 1 lisatteur
17	na 6 Oshum 37. Evenite: OtlaE Earnet-Ohila Mac Merilio Sudore: OtlaE Carmet-Ohila Mac Merilio Sudore: OtlaE Difference: Ett
	Execute Curvenies / Print Curve, Pointait Shill PAILER, Pointait

5. Check the query result on the **View Result** tab below the editing window.

Executed Queries (Last [	Day) View Result								Clear All
Result O									
Executed successfully									
Query Q01: Price s Job ID 7dcda5f5-edfc	ummary report query Qu I-419b-b1e8-e6bafdc26d18	iry on a single table line	eitem within a certain period of ti	me,					
The query takes 47.04s, an	nd 70.87 MB scanned.A ma:	kimum of 1,000 records	can be displayed.					Enter a keyword.	QLLLL
l_returnflag ↓≡	I_linestatus ↓≣	sum_qty ↓≣	sum_base_price ↓≣	sum_disc_price ↓Ξ	sum_charge ↓≣	avg_qty ↓≣	avg_price ↓≣	avg_disc ↓≣	count_order ↓≣
A	F	3774200	5320753880.6900215	5054096266.682835	5256751331.449283	25.53758711685	36002.12382901429	0.0501445970634	147790
N F 95257 133737795.83999997 127132372.65119997 132286291.22544508 25.30066407062 35521.3269163346 0.0493944223107 3765								3765	
N	0	7459297	10512270008.89988	9986238338.38475	10385578376.585411	25.54553767123	36000.9246880132	0.0500959589041	292000
R	F	3785523	5337950526.469969	5071818532.941963	5274405503.049363	25.52594385742	35994.0292140307	0.0499892785618	148301

This example uses the **default** queue and database preset in the system as an example. You can also run the command in a self-created queue and database.

For details, see **Creating a Queue**. For how to create a database, see **Creating a Database**.

### 8.9.3 TPC-H Sample Data in the SQL Templates Preset on DLI

#### **TPC-H Sample Data**

TPC-H is a test set developed by the Transaction Processing Performance Council (TPC) to simulate decision-making support applications. It is widely used in academia and industry to evaluate the performance of decision-making support technology. This business test has higher requirements on vendors, because it can comprehensively evaluate the overall business computing capability. With universal business significance, is widely used in analysis of bank credit, credit card, telecom operation, tax, as well as tobacco industry decision-making analysis.

The TPC-H benchmark test is derived from the TPC-D standard, which was established by the TPC organization in 1994 for decision support system testing. TPC-H implements a 3NF data warehouse that contains eight basic relationships, with a data volume range from 1 GB to 3 TB. The TPC-H benchmark test includes 22 queries (Q1 to Q22). The main evaluation indicator is the response time of each query (from submission to result return). The unit of the TPC-H benchmark test is the query number per hour (QphH@size). H indicates the average number of complex queries per hour. **size** indicates the size of database, which reflects the query processing capability of the system. TPC-H can evaluate key performance parameters that other tests cannot evaluate, because it is modeled based on the actual production and operation environment. In a word, the TPC-H standard by TPC meets the test requirements of data warehouse and motivate vendors and research institutes to stretch the limit of this technology.

In this example, DLI directly queries the TPC-H dataset on OBS. DLI has generated a standard TPC-H-2.18 dataset of 100 MB which is uploaded to the tpch folder on OBS. The read-only permission is granted to you to facilitate query operations.

#### **TPC-H Test and Metrics**

TPC-H test is divided into three sub-tests: data loading test, Power test, and Throughput test. Data loading indicates the process of setting up a test database, and the loading test is to test the data loading ability of DBMS. The first test is data loading test that tests data loading time, which is time-consuming. The second test is Power test, also called raw query. After data loading test is complete, the database is in the initial state without any other operation, especially the data in the buffer is not tested. Power test requires that the 22 queries be executed once in sequence and a pair of RF1 and RF2 operations be executed at the same time. The third test is Throughput test, the core and most complex test, more similar to the actual application environment. With multiple query statement groups and a pair of RF1 and RF2 update flows, Throughput test pose greater pressure on the SUT system than Power test does.

The basic data in the test is related to the execution time (the time of each data loading step, each query execution, and each update execution), based on which you can calculate the data loading time, Power@Size, Throughput@Size, qphH@Size and \$/QphH@Size.

Power@Size is the result of the Power test, which is defined as the reciprocal of the geometric average value of the query time and change time. The formula is as follows:

$$\frac{3600 * SF}{\sum_{i=1}^{24} \prod_{i=1}^{i=22} QI(i,0) * \prod_{j=1}^{j=2} RI(j,0)}$$
TPC-H Power@Size =

Size indicates the data size. SF is the scaling factor of data scale. QI (i, 0) indicates the time of the ith query, in seconds. R (I j, 0) is the update time of RFj, in seconds.

Throughput@Size is the Throughput test result, which is defined as the reciprocal of the average value of all query execution time. The formula is as follows:

QphH@Size =  $\sqrt{Power}$  @ Size \* Throughput @ Size

#### Service Scenario

You can use the built-in TPC-H test suite of DLI to perform interactive query without uploading data.

#### Advantages of DLI Built-in TPC-H

- You can log in to DLI and get permission to run SQL statements without creating tables or import data.
- The 22 preset TPC-H SQL query templates with rich functions meet the requirements of most business scenarios. You do not need to download TPC-H query statements, which saves your time and energy.
- Data Lake gives you brand-new experience of serverless DLI product within the minimum time.

#### Precautions

When a sub-account uses the TPC-H test suite, the main account needs to grant the sub-account the OBS access permission and the permission to view the main account table. If the master account has not logged in to DLI, the sub-account needs to have the permissions to create databases and tables in addition to the preceding permissions.

#### Developing and Submitting a SQL Job Using the TPC-H Sample Template

DLI allows you to create custom templates or save currently used SQL statements as templates for quick and convenient SQL operations. Once a template is saved, you can execute SQL operations directly through the template without the need to write SQL statements.

- 1. Log in to the DLI management console.
- On the DLI management console, choose Job Templates > SQL Templates in the navigation pane on the left. On the displayed Sample Templates tab, click the plus sign next to tpchQuery to find the Q1 Price summary report guery template. Then, click Execute in the

**Q1\_Price\_summary\_report\_query** template. Then, click **Execute** in the **Operation** column.

SQL Temp	plates			Set Property 🛛 🖗 Quick Links	Create Template
Samp	ple Templates Custom Tem	nplates			
				Template Name	QC
Name		Description	Statement		Operation
<b>—</b> tp	pchQuery				
Q	21_Price_summary_report_query	Query on a single table lineitem within a	Q01: Price summary report query Quiry on a single table	lineitem within a certain period of time, statisti	Execute
Q	22_The_lowest_cost_supplier_analy	In the given area, for a specified part (pa	Q02: Minimum cost supplier analysis In the given area, fo	or the specified parts (parts of a certain type an	Execute
Q	23_Shipping_priority_analysis	This query analyzes the shipping priority	QQ03: Shipping priority analysis Description: Query the to	p 10 orders that have not yet shipped Shippi	Execute
4	24_Order_priority_check_analysis	Query to obtain order priority statistics	Q04: Analysis of order priority check Query to obtain orde	r priority statistics. Calculate the number of or	Execute
Q	25_Analysis_of_the_number_of_loc	Query the number of local suppliers who	Q5: Analysis of the number of local suppliers This query li	sts the amount of income obtained through loc	Execute
Q	26_Analysis_of_forecasted_income	This query determines the amount of rev	Q6: Analysis of forecasted income changes This query dete	ermines the amount of revenue increase that is	Execute

3. In the upper part of the editing window, set **Engine** to **Spark**, **Queues** to **default**, and **Databases** to **default**, and click **Execute**.

	Engine spark 💌 Queues default 🔍 Database default 🔍 💽 Execute Format Set Property More 🔻
_	
1	Q01: Price summary report query.
2	Quiry on a single table lineitem within a certain period of time, statistics on all kinds of commodities that have been paid and already
3	Features: Single table query operations with grouping, sorting, and aggregation coexist. This query will read 95% to 97% of the data on
4	SELECT 1_returnflag,
5	l_linestatus,
6	sum(1 quantity) AS sum qty,
7	<pre>sum(l_extendedprice) AS sum base price,</pre>
8	<pre>sum(1_extendedprice * (1 - 1_discount)) AS sum_disc_price,</pre>
9	<pre>sum(1_extendedprice * (1 - l_discount) * (1 + l_tax)) AS sum_change,</pre>
10	<pre>avg(1_quantity) AS avg_qty,</pre>
11	<pre>avg(1_extendedprice) AS avg_price,</pre>
12	avg(l_discount) AS avg_disc,
13	count(*) AS count_order
14	FROM tpch.lineitem
15	WHERE 1_shipdate <= DATE "1998-12-01" - INTERVAL "90" DAY
16	GROUP BY 1_returnflag, 1_linestatus
17	ORDER BY 1 returnflag. 1 linestatus:

4. Check the query result on the **View Result** tab below the editing window.

Executed Queries (Last D	Day) View Result								Clear All
Result1 O									
Executed successfully									
Query Q01: Price si Job ID 7dcda5f5-edfd	ummary report query Q I-419b-b1e8-e6bafdc26d18	uiry on a single table lir 8	eitem within a certain period of t	ime,					
The query takes 47.04s, ar	nd 70.87 MB scanned.A ma	aximum of 1,000 records	can be displayed.					Enter a keyword.	Q HL LÍ 🕹
l_returnflag ↓≣	I_linestatus ↓≡	sum_qty ↓≣	sum_base_price ↓≡	sum_disc_price ↓≡	sum_charge ↓≣	avg_qty ↓≡	avg_price ↓≣	avg_disc ↓⊞	count_order ↓≣
A	F	3774200	5320753880.6900215	5054096266.682835	5256751331.449283	25.53758711685	36002.12382901429	0.0501445970634	147790
N	F	95257	133737795.83999997	127132372.65119997	132286291.22944508	25.30066401062	35521.3269163346	0.0493944223107	3765
N	0	7459297	10512270008.89988	9986238338.38475	10385578376.585411	25.54553767123	36000.9246880132	0.0500959589041	292000
R	F	3785523	5337950526.469969	5071818532.941963	5274405503.049363	25.52594385742	35994.0292140307	0.0499892785618	148301

This example uses the **default** queue and database preset in the system as an example. You can also run the command in a self-created queue and database.

# **9** Developing a DLI SQL Job in DataArts Studio

#### Scenario

Huawei Cloud DataArts Studio provides a one-stop data governance platform that integrates with DLI for seamless data integration and development, enabling enterprises to manage and control their data effectively.

This section walks you through how to develop a DLI SQL job in DataArts Studio.

Figure 9-1 Process of developing a DLI SQL job in DataArts Studio

#### **Development Process**



- table. See Step 1: Create a Database and Table.3. Import service data: Submit SQL scripts to import service data. See Step 2:
- Calculate and Process Service Data.Query and analyze data: Submit SQL scripts to analyze service data, for
- example, querying daily sales. See **Step 3: Query and Analyze Sales Data**.
- 5. Orchestrate a job: Orchestrate data processing and analysis scripts into a pipeline. DataArts Studio executes all nodes based on the orchestrated pipeline sequence. See **Step 4: Orchestrate a Job**.
- 6. Test job runs: Test if jobs can run properly. See **Step 5: Test Job Running**.
- 7. Configure job scheduling and monitoring: Set job scheduling attributes and monitoring rules. See **Step 6: Set Periodic Job Scheduling** and **Related Operations**.

а

#### **Prepare Environments**

- Prepare a DLI resource environment.
  - Configure a DLI job bucket.

Before using DLI, you need to configure a DLI job bucket. The bucket is used to store temporary data generated during DLI job running, such as job logs and results.

For details, see **Configuring a DLI Job Bucket**.

- Create an elastic resource pool and create a SQL queue within it.

An elastic resource pool offers compute resources (CPU and memory) required for running DLI jobs, which can adapt to the changing demands of services.

You can create multiple queues within an elastic resource pool. These queues are associated with specific jobs and data processing tasks, and serve as the basic unit for resource allocation and usage within the pool. This means queues are specific compute resources required for executing jobs.

Queues within an elastic resource pool can be shared to execute jobs. This is achieved by properly setting the queue allocation policy. This improves queue utilization.

For details, see **Creating an Elastic Resource Pool and Creating Queues** Within It.

#### • Prepare a DataArts Studio resource environment.

- Buy a DataArts Studio instance.

Buy a DataArts Studio instance before submitting a DLI job using DataArts Studio.

For details, see **Buying a DataArts Studio Basic Package**.

#### - Access the DataArts Studio instance's workspace.

i. After buying a DataArts Studio instance, click Access.

#### Figure 9-2 Accessing a DataArts Studio instance

	Enterprise Project: default		
Version	Enterprise	Billing Mode	Yearly/Monthly
Created	Oct 23, 2019 10:02:13 GMT+08:00	Name	
Expires	Oct 22, 2020 23:59:59 GMT+08:00	Instance ID	9ab2da986bf34d70b62c323850
Order No.	-	Status	Valid
Description	🖉		
A	ccess   Renew	le Buy	♥ Upgrade

Click the Workspaces tab to access the data development page.
 By default, a workspace named default is created for the user who has purchased the DataArts Studio instance, and the user is assigned the administrator role. You can use the default workspace or create one.

For how to create a workspace, see **Creating and Managing a Workspace**.

	Back to Instar	nce List			
Dashboard	Workspaces	Roles	Industry	Assets	Tags
Create Work	space				
Q Select a pr	roperty or enter a keyv	vord.			
Name/ID 👙	Status	÷ Mo	ode ≑	Description	÷
	📀 En	able Sii	mple		

#### Figure 9-3 Accessing the DataArts Studio instance's workspace

#### Figure 9-4 Accessing DataArts Studio's data development page

Data			
Development	AI Owne ⊘ C ↓Ξ :		
Overview	Enter a keyword		Create Script
	Cital disciplicate. Gr	- 1773	Develop, debug, and run scripts online, or run developed scripts in jobs.
Data			
Development		<u> </u>	+ SQL + Hive SQL + DLI SQL + DWS SQL + Spark SQL + Spark Python + Flink SQL
Davelon Projet	the second se		+ RDS SQL + Presto SQL + HetuEngine SQL + ClickHouse SQL + Impala SQL + Shell
Develop Script			
Develop Job			+ Python + Dons SQL
·			
Monitoring			
Overview			Create Job
Job Monitoring	< .		cleate vob
Monitor Instance			Drag and drop the nodes on a canvas and connect them to easily develop jobs.
Monitor instance			+ Create Job + Create Data Micration Job
Monitor PatchData			
Duty Schedules			
Manage			
Notification			Create Data Connection
Manage Backup			
Conception 1 Patron		a St	Compute data source mormation and create data connections, through which you can access data sources when developing scripts and pos-
Operation History		9	+ Create Data Connection

#### Step 1: Create a Database and Table

#### Step 1 Develop SQL scripts for creating databases and tables.

Databases and tables are the basis for developing SQL jobs. Before running a job, you need to define databases and tables based on your service scenarios.

This part describes how to develop a SQL script to create databases and tables.

- In the left navigation pane of DataArts Factory, choose Data Development > Develop Script.
- 2. In the **Create Script** area on the right, click + **DLI SQL**.

#### Figure 9-5 Creating a DLI SQL script

	Create Script
·	Develop, debug, and run scripts online, or run developed scripts in jobs.
	+ ISQL + Hive SQL + DLI SQL + DWS SQL
	+ Spark SQL + Spark Python + Flink SQL + RDS SQL
	+ Presto SQL + HetuEngine SQL + ClickHouse SQL
	+ Impala SQL + Shell + Python + Doris SQL

3. On the script editing page, enter sample code for creating a database and table.

"SQL -- Create a database.
 CREATE DATABASE IF not EXISTS supermarket\_db;-- Create product dimension table.
 CREATE TABLE IF not EXISTS supermarket*db.products ( product*id INT, productname STRING, category STRING, price DECIMAL(10,2) ) using parquet;
 - Create a transaction table. *(product*id INT, -- Product ID productname STRING, -- Product name category STRING, -- Product category price DECIMAL(10,2) -- Unit price)
 CREATE TABLE IF not EXISTS supermarket*db.transactions (transaction*id INT, productid INT, quantity INT, dt STRING ) using parquet partitioned by (dt);
 - Create a sales analysis table. *(transaction*id INT, -- Transaction ID productid INT, -- Product ID quantity INT, -- Quantity dt STRING -- Date)
 CREATE TABLE IF not EXISTS supermarket*db.analyze (transaction*id INT, product*id INT, product*name STRING, quantity INT, dt STRING ) using parquet partitioned by (dt);
 - *(transaction*id INT, -- Transaction ID product*id INT, -- Product* ID product*id INT, -- Product* name quantity INT, -- Quantity dt STRING -- Date)

- 4. Click **Save** to save the SQL script. In this example, the script is named **create\_tables**.
- 5. Click **Submit** to run the script to create a database and table.

#### Step 2 Create a SQL job running script.

 In the left navigation pane of DataArts Factory, choose Data Development > Develop Job.

# Figure 9-6 Creating a job

Data				
Development	All Owne	⊚ C 1≣ ≣		
Overview	Enter a keyword.	Q		Create Script
			□	Develop, debug, and run scripts online, or run developed scripts in jobs.
Data				+ SOI + Hwe SOI + DU SOI + DWS SOI
Development			- 1/20	
Develop Script				+ Spark SQL + Spark Python + Flink SQL + RDS SQL
Develop Job				+ Presto SQL + HetuEngine SQL + ClickHouse SQL
				+ Impala SQL + Shell + Python + Doris SQL
Monitoring		_		
Overview	<	_		
Job Monitoring				
t de sites la stances				Create Job
MONITOR INSTANCE	-	_		Drag and drop the nodes on a canvas and connect them to easily develop jobs.
Monitor PatchData		_	=	
Duty Schedules				+ Create Job + Create Data Migration Job

2. Click **Create Job**. In the dialog box that appears, edit job information. In this example, the SQL job is named **job\_create\_tables**.

Create Job A maximum of 100	0,000 jobs can be created. You can create 99,982 more jobs.
* Job Name	job_creat_tables
Јор Туре	Batch processing     Real-time processing
Mode	Pipeline      Single task
Select Directory	/Jobs/   (Jobs/
Owner (?)	× 🕀
Priority	High  Medium Low
Agency ⑦	Select an agency.
Log Path	obs://dlf-log-330e068af1334c9782f4226acc00a2e2/
	<ul> <li>I agree to create OBS bucket obs://dlf-log- 330e068af1334c9782f4226acc00a2e2/. This bucket is used only for storing run logs of DLF jobs.</li> </ul>
	To change the log path, go to the WorkSpaces page.
	OK Cancel

#### Figure 9-7 Editing job information

- 3. On the job development page, drag the DLI SQL node to the canvas and click the node to edit its properties.
  - **SQL or Script**: Select **SQL script** in this example. Then, select the script created in **Step 2.2** for **SQL script**.
  - **Database Name**: Select the database configured in the SQL script.
  - Queue Name: Select the SQL queue created in Create an elastic resource pool and create a SQL queue within it.

For more property parameters, see **Parameters of DLI SQL nodes**.

↓ sob_creat_tables ×	+					
Node Library ^	💾 Save 🔥 Submit 🄓 Unlock 🔒 Loci	k 🕨 Test	👲 Clear 🛛 🗧 Full Sc	reen 📑 Export 🕲 Refresh 👳 Monitor	∬ Base	eline Ta
MRS Kafka		To Execute		DLI SQL		
Kafka Client ROMA FDI Job				Properties * Node Name		^
A TICS Job				creat_tables * SQL or Script		
Compute & Analytics 🔺				<ul> <li>SQL Statement</li> <li>SQL script</li> </ul>		
DLI Flink Job DLI SQL	creat_lables			* SQL script	⊕ <b>₽</b>	₫
DLI Spark DWS SQL				Script Parameter C		
seet 😨				* Database Name		0
MRS Spark MRS Hive SQL SQL						+
() ()				* Queue Name		
MRS Presto MRS SQL HetuEngine					$\oplus$	0

Figure 9-8 Editing the properties of a DLI SQL node

4. Once the properties are edited, click **Save** to save the configuration.

#### Step 3 Configure job scheduling.

As the database and table only need to be created once, only one-time scheduling is configured in this example.

- 1. Left-click the blank area of the canvas.
- 2. Click **Scheduling Setup** and select **Run once**. (The job will be scheduled only once and will not be automatically scheduled later.)

#### Figure 9-9 Configuring job scheduling

호 Clear 🗧 Full Screen 📑 Export 💿 Refresh 🖙 Monitor 🔑 Baseline Task	Link
Scheduling Setup	Basic Info
Run once Run periodically Event-based	Scheduling Setup
Dry Run Task Groups Do not select	Parameter Setup

3. After configuring the scheduling, click **Execute**. Click **Go to O&M Center** to view the job status.

----End

#### Step 2: Calculate and Process Service Data

#### Step 1 Develop a SQL script for importing service data.

This part describes how to submit a SQL script to import service data.

- 1. In the left navigation pane of DataArts Factory, choose Data Development > Develop Script.
- 2. In the **Create Script** area on the right, click + **DLI SQL**.

#### Figure 9-10 Creating a DLI SQL script

#### Create Script

Develop, debug, and run scripts online, or run developed scripts in jobs. </>> + DWS SQL + DLI SQL + SQL + Hive SQL + Spark SQL + Flink SQL + RDS SQL + Spark Python + Presto SQL + HetuEngine SQL + ClickHouse SQL + Impala SQL + Doris SQL + Shell + Python

3. On the script editing page, enter sample code for analyzing data. SQL: Data in actual services is typically from other data sources. This example simplifies the data

import logic and simulates the insertion of product data. INSERT INTO supermarketdb products (productid, productname, category, price) VALUES (1001, 'Shampoo', 'Daily necessities', 39.90), (1002, 'Toothpaste', 'Daily necessities', 15.90) (1003, 'Instant noodles', 'Food', 4.50) (1004, 'Coke', 'Beverage', 3.50); -- Data in actual services is typically from other data sources. This example simplifies the data import logic and simulates the insertion of transaction records. INSERT INTO supermarketdb.transactions (transactionid, productid, quantity, dt) VALUES (1, 1001, 50, '2024-11-01'), -- 50 bottles of shampoo were sold. (2, 1002, 100, '2024-11-01'), -- 100 tubes of toothpaste were sold.

(3, 1003, 30, '2024-11-02'), -- 30 packs of instant noodles were sold. (4, 1004, 24, '2024-11-02'); -- 24 bottles of Coke were sold.

-- Simulate supermarket business analysis and query the transaction records of a certain product. INSERT INTO supermarketdb.analyze SELECT t.transactionid, t.productid, p.productname, t.quantity, t.dt

FROM supermarketdb.transactions t

JOIN supermarketdb.products p ON t.productid = p.productid WHERE t.dt = '2024-11-01';

- Click Save to save the SQL script. In this example, the script is named 4. job\_process\_data.
- Click **Submit** to execute the script. 5.

#### Step 2 Create a SQL job.

In the left navigation pane of DataArts Factory, choose **Data Development** > 1. Develop Job.

#### Figure 9-11 Creating a job

Data Development Overview Data Development Develop Script Develop Job Monitoring	I I I I I I I I I I I I I I I I I I I	> III 3 000	• [2] .	Create Script Develop, debug, and + SOL + Spark SOL + Presto SOL + Impata SOL	un scripts online, or run developed scripts in jobs. + Hive SQL + DU SQL + DWS SQL + Spark Python + Flink SQL + RDS SQL + HetuEngine SQL + ClickHouse SQL + Shell + Python + Dons SQL
Overview					
Monitor Instance Monitor PatchData Duty Schedules	and the second sec			Create Job Drag and drop the not + Create Job	tes on a canvas and connect them to easily develop jobs.
	•				

2. Click **Create Job**. In the dialog box that appears, edit job information. In this example, the SQL job is named **job\_process\_data**.

* Job Name	job_process_data	
Јор Туре	Batch processing     Real-time processing	
Mode	Pipeline     Single task	
Select Directory	/Jobs/	÷
Owner ⑦		× 🕀
Priority	● High ○ Medium ○ Low	
Agency ⑦	Select an agency.	(+)
Log Path	obs://dlf-log-330e068af1334c9782f4226acc00a2e2/	
	I agree to create OBS bucket obs://dlf-log- 330e068af1334c9782f4226acc00a2e2/. This bucket for storing run logs of DLF jobs.	t is used only
	To change the log path, go to the WorkSpaces page.	

Figure 9-12 Editing job information

- 3. On the job development page, drag the DLI SQL node to the canvas and click the node to edit its properties.
  - SQL or Script: Select SQL script in this example. Then, select the script created in Step 1 for SQL script.
  - **Database Name**: Select the database configured in the SQL script.
  - Queue Name: Select the SQL queue created in Create an elastic resource pool and create a SQL queue within it.
  - Environment variable: The DLI environment variable is optional.

The description of the parameter added in this example is as follows:

spark.sql.optimizer.dynamicPartitionPruning.enabled = true

- This parameter is used to control whether to enable dynamic partition pruning. Dynamic partition pruning can help reduce the amount of data that needs to be scanned and improve query performance when executing SQL queries.
- When set to true, dynamic partition pruning is enabled. SQL automatically detects and deletes partitions that do not meet the WHERE clause conditions during query. This is useful for tables that have a large number of partitions.

- If SQL queries contain a large number of nested left join operations and the table has a large number of dynamic partitions, a large number of memory resources may be consumed during data parsing. As a result, the memory of the driver node is insufficient and there are frequent Full GCs.
- To avoid such issues, you can disable dynamic partition pruning by setting this parameter to **false**.

However, disabling this optimization may reduce query performance. Once disabled, Spark does not automatically prune the partitions that do not meet the requirements.

For more property parameters, see **Parameters of DLI SQL nodes**.

📱 Save 🔥 Submit 🔓 Unlock	🔒 Lock ㅣ	▶ Test	🔶 Clear	🕤 Full Scree	n 🚺 Base	line Task L	inl∈
		Execute					Node P
		Properties				^	roperties
		* Node Name					line
		process_data					eageli
		* SQL or Script					лfo
	_	SQL Statement	💿 ऽହା	_ script			
process_data		* SQL script					
		process_data			⊕₽	$\underline{\mathcal{O}}$	
		Script Parameter	С				

Figure 9-13 Editing the properties of a DLI SQL node

4. Once the properties are edited, click **Save** to save the configuration.

----End

#### Step 3: Query and Analyze Sales Data

#### Develop a SQL script for analyzing and processing data.

This part describes how to submit a SQL script to analyze data.

- In the left navigation pane of DataArts Factory, choose Data Development > Develop Script.
- 2. In the **Create Script** area on the right, click + **DLI SQL**.

#### Figure 9-14 Creating a DLI SQL script

	Create Script	
•	Develop, debug, and run scripts online, or run developed scripts in jobs.	
	+ BQL + Hive SQL + DLI SQL + DWS SQL	
	+ Spark SQL + Spark Python + Flink SQL + RDS	S SQL
	+ Presto SQL + HetuEngine SQL + ClickHouse SQL	
	+ Impala SQL + Shell + Python + Doris SQL	L

- On the script editing page, enter sample code for analyzing data.
   Query daily sales
   SELECT transaction\_id, productid, productname, quantity, dt FROM supermarket\_db.analyze WHERE dt = '2024-11-01';
- 4. Click **Save** to save the SQL script. In this example, the script is named **select\_analyze\_data**.
- 5. Click **Submit** to execute the script.

#### Step 4: Orchestrate a Job

- 1. Create a DLI SQL node named **select\_analyze\_data** in the **job\_process\_data** job. Then, click the node to edit its properties.
  - SQL or Script: Select SQL script in this example. Then, select the script created in Step 1 for SQL script.
  - **Database Name**: Select the database configured in the SQL script.
  - Queue Name: Select the SQL queue created in Create an elastic resource pool and create a SQL queue within it.

For more property parameters, see **Parameters of DLI SQL nodes**.

Figure 9-15 Editing the properties of a DLI SQL node

	🐻 Execute	Enter:	Node P
		Properties ^	roperties
		* Node Name	lin
		select_analyze_data	eagel
$\bigcirc$		* SQL or Script	nfo
		SQL Statement () SQL script	
process_data	select_analyze_data	* SQL script	
		select_analyze_data 💿 🔉 🖉	
		Script Parameter C	

- 2. Once the properties are edited, click Save to save the configuration.
- 3. Orchestrate the two nodes into a pipeline. DataArts Studio executes all nodes based on the orchestrated pipeline sequence. Then, click **Save** and **Submit** in the upper left corner.

#### Step 5: Test Job Running

After orchestrating the job, click **Test** to test it.

After testing the job, open the **select\_analyze\_data** SQL script file and click **Execute** to query and analyze sales details.

If query results meet your expectation, go to **Step 6**: **Set Periodic Job Scheduling** to set periodic job scheduling.

Figure 9-16 Running the select\_analyze\_data script



#### Step 6: Set Periodic Job Scheduling

- In the left navigation pane of DataArts Factory, choose Data Development > Develop Job.
- 2. Double-click job\_process\_data.
- 3. Click the Scheduling Setup tab on the right.
- 4. Select **Run periodically** and set scheduling properties.

In this example, the job's scheduling policy starts at 10:15:00 on November 22, 2024. The first scheduled time is 10:20:00 on the same day. With a daily scheduling interval, the job automatically runs at 10:20:00 a.m. each day, executing the nodes based on the orchestrated pipeline sequence.

#### Figure 9-17 Configuring job scheduling

Save	🔥 Submit	🔓 Unlock	Lock	▶ Test		👲 Clear 🛛 📓 F	full Screen 📑 E	xport 📀 Refresh	👳 Monitor 🛛	N Bas	eline Task	Link
				76 Execute	Ð	Scheduling \$	Setup					Basic In
						Scheduling Ty Run once	Pe Run periodically Nation	C Event-based				fo Scheduling Setup
						Scheduling Pr	roperties				^	Param
	<u> </u>	→ <u>©</u>	)			* From	Nov 22, 2024 10:2	2:00		##	to	eter Setup
pr	ocess_data	select_analyze	e_data			l	Valid permanent	ime. Iy				Version

5. Click **Save**, **Submit**, and **Execute** in sequence to complete periodic scheduling configuration.

For more job scheduling settings, see **Setting Up Scheduling for a Job**.

#### **Related Operations**

#### • Configure job monitoring.

DataArts Studio monitors the status of batch jobs.

This type of job is a pipeline that consists of one or more nodes and is scheduled as a whole.

In the left navigation pane of DataArts Factory, choose **Monitoring** > **Job Monitoring**. On the displayed **Batch Jobs** tab, view the scheduling status, scheduling interval, and scheduling start time of batch jobs.

For details, see Monitoring a Batch Job.

#### Figure 9-18 Configuring job monitoring

Batch Jobs	Real-Time Jobs									
Execute	Stop Scheduling	Configure Notification	Add Job Tag	Export All Data	Exact match	Job Name 🔻	(	C Filter jobs by notification co	▼ In-Progr	es 🜒 🖡 🛛 Please select j 🔻
										Start Date – End Dale 📋 C 💿
	Name	Scheduling	Scheduli 7	Sched 7 Next Plan Time	÷	Owner	Last Modified By	Last Modified $\frac{A}{\Psi}$	Last Instan	La Operation
. v	)00_01   W	🜖 in-Progr	Run periodic	1 days Mar 30, 2024 00;	10:00 GMT+0	$q \log (1/2^{1/2})$	$(der (0))^{1/2}$	Jan 10, 2024 09:20:20 GMT+0	Successful	Ma: Pause   Stop Scheduling   More 🕶

#### • Configure instance monitoring.

Each time a job is executed, a job instance record is generated.

In the left navigation pane of DataArts Factory, choose **Monitoring** > **Monitor Instance**. On the instance monitoring page that appears, you can view job instance information and perform more operations on the instances as needed.

For more information, see Instance Monitoring.

#### FAQ

# • If a DataArts Studio job fails, and the logs provided by DataArts Studio are not detailed enough, what should I do? Where can I find more specific logs?

You can locate the DLI job ID through the logs provided by DataArts Studio, and then find the specific job in the DLI console using the DLI job ID.

#### Figure 9-19 Monitoring log file



Once you find the specific job in the DLI console, click archived logs to view detailed logs.

#### Figure 9-20 Entering the job ID

Data Lake Insight 🔍	SQL Jobs 💿							😛 Feedlack	Onute Job
Overview SQL 5/Hor Jub Management	-Select- Queses	V OCL X DOL X DELETE X	ll ∨ Leet t ey Type Stat	v Q	200 (D. 70:10004-40:a-4000-00:a5-00040:0540:a44 × Add fram	r Duration (t)	Created (t)	Operation	× 0 0
SQL Jobs First Jobs	A 7.36	5perk	ouenr 😋	Trained		34min 41.32s	Peb 11, 2025 08:47:58-09/7+08:00	Diff SparkU More v	
Sperk Jobs V Reserves V Data Normalistic V Data Serviciation V Datascons Connections Obtain Configuration V	Ourses Datortame Oury Crasel Parameter Settings Nurster of Houston Press Tribute Detriferens Output Byes Presst Format				Jan D Type Donatron Kanael Sacone Data Canalig Tabat (2014 Sacone) Read Yahn SQL Samatane Panel	1001004-00-400-000-004004044	119172712712712712712711111111111111111	1010 0000000000 ()	
	Tatal Resords: 1							10 × 1 < 1	> au 🗔

You can also search for the job in the DLI console using the node name or job name provided by DataArts Studio.

#### Figure 9-21 Node name or job name

	node_name=analyze_dataplan_time=202	×
da1657bf-aa66-4e4c-bd0e-2b4221b5d701	4_11_20_19_17_55;dgc_job_id=370A8F48 A41045179E4AC301CADA66D7vhOwaQU	
INSERT	.job_name=job_DLI_SQL;sqlStep=1;works	
43.63s	pace=test- 001:updateTime=1732101508000:	
2024/11/20 19:23:08 GMT+08:00	001,0000000,	
node_name=analyze_data;plan_time=2024_	11_20_19_17_55;dgc_job_id=370A8F48A41045179E4A0	C301CADA66D7vhOwaQUJ;job_name.
0.93 KB		
0		
supermarket_db		
0		

# • What should I do if I encounter permission errors when running complex DLI jobs?

DLI needs to work with other cloud services. You must grant DLI basic operation permissions of these services so that DLI can access them and perform resource O&M operations on your behalf.

For more information, see **Configuring DLI Agency Permissions**.

# **10** Submitting a SQL Job Using JDBC

# 10.1 Downloading and Installing the JDBC Driver Package

#### Scenario

To connect to DLI, JDBC is utilized. You can obtain the JDBC installation package from Maven or download the JDBC driver file from the DLI management console.

This section describes how to connect to DLI and submit a SQL job using JDBC.

#### **Obtaining the Server Connection Address**

The address for connecting to DLI is in the format of **jdbc:dli**://*<endPoint>/ <projectId>*. So, you need to obtain the endpoint and project ID.

Obtain the DLI endpoint from **Regions and Endpoints**. Specifically, log in to the public cloud, hover over your username in the upper right corner, and choose **My Credentials** from the shortcut menu. You can obtain the project ID on the displayed **API Credentials** tab page.

#### Downloading and Installing the JDBC Driver

#### D NOTE

Once JDBC 2.X has undergone function reconstruction, query results can only be accessed from DLI job buckets. To utilize this feature, certain conditions must be met:

- On the DLI management console, choose **Global Configuration** > **Project** to configure the job bucket.
- Starting May 2024, new users can directly use DLI's function to write query results into buckets without needing to whitelist it.

For users who started using DLI before May 2024, to use this function, they must submit a service ticket to whitelist it.

Method 1: Adding the JDBC driver using the Maven central repository
 The Maven central repository is part of the Apache Maven project that

provides Java libraries and frameworks.
When the JDBC retrieval method is not specified, the default approach is to add the JDBC driver using the Maven central repository.

Use Maven to add the Maven configuration item on which **huaweicloud-dlijdbc** depends. (This is the default operation and does not need to be configured separately.)

<dependency>

<groupId>com.huawei.dli</groupId> <artifactId>huaweicloud-dli-jdbc</artifactId> <version>*x.x.*</version> </dependency>

• Method 2: Obtaining the JDBC driver using Maven to configure the Huawei image source

When using Maven to manage project dependencies, you can modify the **settings.xml** file to configure the Huawei image source to obtain the JDBC driver.

```
<mirror>
<id>huaweicloud</id>
<mirrorOf>*</mirrorOf>
<url>https://mirrors.huaweicloud.com/repository/maven/</url>
</mirror>
```

- Method 3: Downloading the JDBC driver file from the DLI management console
  - a. Log in to the DLI management console.
  - b. Click **SDK Download** in the **Common Links** area on the right of the **Overview** page.
  - c. On the **DLI SDK DOWNLOAD** page, select a driver and download it.

Click huaweicloud-dli-jdbc-x.x.x to download a JDBC driver package.

**NOTE** 

The JDBC driver package is named **huaweicloud-dli-jdbc-***<version>.zip*. It can be used in all versions of all platforms (such as Linux and Windows) and depends on JDK 1.7 or later.

d. The downloaded JDBC driver package contains **.bat** (Windows) or **.sh** (Linux/Mac) scripts, which are used to automatically install the JDBC driver to the local Maven repository.

You can choose a script based on your OS to install the JDBC driver.

- Windows: Double-click the .bat file or run the file in the CLI.
- Linux/Mac: Run the .sh script.

# Authentication

You need to be authenticated when using JDBC to create DLI driver connections.

JDBC currently supports two authentication modes: AK/SK-based and token-based. Token-based authentication is only supported by **dli-jdbc-1**.*x*. AK/SK-based authentication is recommended.

#### • (Recommended) Generating an AK/SK

- a. Log in to the DLI management console.
- b. Hover over the username in the upper right corner and select **My Credentials** from the shortcut list.

- c. The **Projects** area is displayed on the **API Credentials** tab page by default. Choose **Access Keys** in the navigation pane on the left.
- d. Click **Create Access Key**. In the dialog box that appears, set **Login Password** and **SMS Verification Code**.
- e. Click **OK** to download the certificate.
- f. Once the certificate is downloaded, you can obtain the AK and SK information in the **credentials** file.

#### **NOTE**

Hard coding AKs and SKs or storing them in code in plaintext poses significant security risks. You are advised to store them in encrypted form in configuration files or environment variables and decrypt them when needed to ensure security.

#### • Obtaining a token

When using token-based authentication, you need to obtain the user token and configure the token information in the JDBC connection parameters. You can obtain the token as follows:

a. Send *POST https://<IAM\_Endpoint>/v3/auth/tokens*. To obtain the IAM endpoint and region name in the message body, see Regions and Endpoints.

Here is an example request:

```
NOTE
```

Replace the content in italic in the sample code with the actual values. For details, see **Identity and Access Management API Reference**.

```
{
 "auth": {
  "identity": {
    "methods": [
     "password"
    ],
    "password": {
     "user": {
       "name": "username",
       "password": "password",
      "domain": {
        "name": "domainname"
      }
     }
   }
  },
   "scope": {
    "project": {
     "id": "0aa253a31a2f4cfda30eaa073fee6477" //Assume that project_id is
0aa253a31a2f4cfda30eaa073fee6477.
   }
  }
}
3
```

b. After the request is processed, the value of **X-Subject-Token** in the response header is the token value.

# 10.2 Connecting to DLI and Submitting SQL Jobs Using JDBC

# Scenario

In Linux or Windows, you can connect to the DLI server using JDBC.

## D NOTE

- Jobs submitted to DLI using JDBC are executed on the Spark engine.
- Once JDBC 2.X has undergone function reconstruction, query results can only be accessed from DLI job buckets. To utilize this feature, certain conditions must be met:
  - On the DLI management console, choose **Global Configuration** > **Project** to configure the job bucket.
  - Starting May 2024, new users can directly use DLI's function to write query results into buckets without needing to whitelist it.

For users who started using DLI before May 2024, to use this function, they must submit a service ticket to whitelist it.

DLI supports 13 data types. Each type can be mapped to a JDBC type. If JDBC is used to connect to the server, you must use the mapped Java type. **Table 10-1** describes the mapping relationships.

DLI Data Type	JDBC Туре	Java Type
INT	INTEGER	java.lang.Integer
STRING	VARCHAR	java.lang.String
FLOAT	FLOAT	java.lang.Float
DOUBLE	DOUBLE	java.lang.Double
DECIMAL	DECIMAL	java.math.BigDecimal
BOOLEAN	BOOLEAN	java.lang.Boolean
SMALLINT/SHORT	SMALLINT	java.lang.Short
TINYINT	TINYINT	java.lang.Short
BIGINT/LONG	BIGINT	java.lang.Long
TIMESTAMP	TIMESTAMP	java.sql.Timestamp
CHAR	CHAR	Java.lang.Character
VARCHAR	VARCHAR	java.lang.String
DATE	DATE	java.sql.Date

#### Table 10-1 Data type mapping

# Prerequisites

Before using JDBC, perform the following operations:

1. Getting authorized.

DLI uses the Identity and Access Management (IAM) to implement finegrained permissions for your enterprise-level tenants. IAM provides identity authentication, permissions management, and access control, helping you securely access your HUAWEI CLOUD resources.

With IAM, you can use your HUAWEI CLOUD account to create IAM users for your employees, and assign permissions to the users to control their access to specific resource types.

Currently, roles (coarse-grained authorization) and policies (fine-grained authorization) are supported. For details about permissions and authorization operations, see the **Data Lake Insight User Guide**.

- Create a queue. Choose Resources > Queue Management. On the page displayed, click Buy Queue in the upper right corner. On the Buy Queue page displayed, select For general purpose for Type, that is, the compute resources of the Spark job.
  - **NOTE**

If the user who creates the queue is not an administrator, the queue can be used only after being authorized by the administrator. For details about how to assign permissions, see **Queue Permission Management**.

## Procedure

- **Step 1** Install JDK 1.7 or later on the computer where JDBC is installed, and configure environment variables.
- Step 2 Obtain the DLI JDBC driver package huaweicloud-dli-jdbc-<version>.zip by referring to Downloading and Installing the JDBC Driver Package. Decompress the package to obtain huaweicloud-dli-jdbc-<version>-jar-with-dependencies.jar.
- **Step 3** On the computer using JDBC, add **huaweicloud-dli-jdbc-1.1.1-jar-withdependencies.jar** to the **classpath** path of the Java project.
- **Step 4** DLI JDBC provides two authentication modes, namely, token and AK/SK, to connect to DLI. For how to obtain the token and AK/SK, see **Authentication**.
- Step 5 Run the Class.forName() command to load the DLI JDBC driver.

#### Class.forName("com.huawei.dli.jdbc.DliDriver");

**Step 6** Call the GetConnection method of DriverManager to create a connection.

#### Connection conn = DriverManager.getConnection(String url, Properties info);

JDBC configuration items are passed using the URL. For details, see **Table 10-2**. JDBC configuration items can be separated by semicolons (;) in the URL, or you can dynamically set the attribute items using the Info object. For details, see **Table 10-3**.

Param eter	Description
url	The URL format is as follows:
	jdbc:dli:// <endpoint>/projectId? <key1>=<val1>;<key2>=<val2></val2></key2></val1></key1></endpoint>
	<ul> <li>EndPoint indicates the DLI domain name. ProjectId indicates the project ID.</li> <li>To obtain the endpoint corresponding to DLI, see Regions and Endpoints. To obtain the project ID, log in to the public cloud, move the mouse on the account, and click My Credentials from the endpoint.</li> </ul>
	<ul> <li>Other configuration items are listed after ? in the form of key=value. The configuration items are separated by semicolons (;). They can also be passed using the Info object.</li> </ul>
Info	The Info object passes user-defined configuration items. If Info does not pass any attribute item, you can set it to null. The format is as follows: info.setProperty ("Attribute item", "Attribute value").

#### Table 10-2 Database connection parameters

#### Table 10-3 Attribute items

ltem	Mandatory	Defau lt Value	Description	Supporte d dli-jdbc
queuename	Yes	-	Queue name of DLI.	dli- jdbc-1.x dli- jdbc-2.x
databasenam e	No	-	Name of a database.	dli- jdbc-1.x dli- jdbc-2.x
authenticatio nmode	No	token	Authentication mode. Currently, token- and AK/SK-based authentication modes are supported.	dli- jdbc-1.x
accesskey	Yes	-	AK that acts as the authentication key. For how to obtain the AK, see Authentication.	dli- jdbc-1.x dli- jdbc-2.x
secretkey	Yes	-	SK that acts as the authentication key. For how to obtain the SK, see Authentication.	dli- jdbc-1.x dli- jdbc-2.x

ltem	Mandatory	Defau lt Value	Description	Supporte d dli-jdbc
regionname	This parameter must be configured if <b>authenticat</b> <b>ionmode</b> is set to <b>aksk</b> .	-	Region name. For details, see <b>Regions and Endpoints</b> .	dli- jdbc-1.x dli- jdbc-2.x
token	This parameter must be configured if <b>authenticat</b> <b>ionmode</b> is set to <b>token</b> .	-	Token-based authentication. For details, see Authentication.	dli- jdbc-1.x
charset	No	UTF-8	JDBC encoding mode.	dli- jdbc-1.x dli- jdbc-2.x
usehttpproxy	No	false	Whether to use the access proxy.	dli- jdbc-1.x
proxyhost	This parameter must be configured if <b>usehttppro</b> <b>xy</b> is set to <b>true</b> .	-	Access proxy host.	dli- jdbc-1.x dli- jdbc-2.x
proxyport	This parameter must be configured if <b>usehttppro</b> <b>xy</b> is set to <b>true</b> .	_	Access proxy port.	dli- jdbc-1.x dli- jdbc-2.x

Item	Mandatory	Defau lt Value	Description	Supporte d dli-jdbc
dli.sql.checkN oResultQuery	No	false	<ul> <li>Whether to allow invoking the executeQuery API to execute statements (for example, DDL) that do not return results.</li> <li>Value false indicates that</li> </ul>	dli- jdbc-1.x dli- jdbc-2.x
			invoking of the executeQuery API is allowed.	
			<ul> <li>Value true indicates that invoking of the executeQuery API is not allowed.</li> </ul>	
jobtimeout	No	300	End time of the job submission. Unit: second	dli- jdbc-1.x dli- jdbc-2.x
directfetchthr eshold	No	1000	Check whether the number of returned results exceeds the threshold based on service requirements.	dli- jdbc-1.x
			The default threshold is <b>1000</b> .	

**Step 7** Create a Statement object, set related parameters, and submit Spark SQL to DLI.

```
Statement statement = conn.createStatement();
```

statement.execute("SET
dli.sql.spark.sql.forcePartitionPredicatesOnPartitionedTable.enabled=true");

#### statement.execute("select \* from tb1");

Step 8 Obtain the result.

#### ResultSet rs = statement.getResultSet();

**Step 9** Display the result.

```
while (rs.next()) {
int a = rs.getInt(1);
int b = rs.getInt(2);
}
```

**Step 10** Close the connection.

conn.close();

----End

# Example

#### **NOTE**

- Hard-coded or plaintext AK and SK pose significant security risks. To ensure security, encrypt your AK and SK, store them in configuration files or environment variables, and decrypt them when needed.
- In this example, the AK and SK stored in the environment variables are used. Specify the environment variables **System.getenv("AK")** and **System.getenv("SK")** in the local environment first.

```
import java.sql.*;
import java.util.Properties;
```

public class DLIJdbcDriverExample {

```
public static void main(String[] args) throws ClassNotFoundException, SQLException {
   Connection conn = null;
   try {
      Class.forName("com.huawei.dli.jdbc.DliDriver");
     String url = "jdbc:dli://<endpoint>/<projectId>?databasename=db1;queuename=testqueue";
      Properties info = new Properties();
     info.setProperty("authenticationmode", "aksk");
     info.setProperty("regionname", "<real region name>");
     info.setProperty("accesskey", "<System.getenv("AK")>");
info.setProperty("secretkey", "<System.getenv("SK")>");
     conn = DriverManager.getConnection(url, info);
     Statement statement = conn.createStatement();
     statement.execute("select * from tb1");
      ResultSet rs = statement.getResultSet();
     int line = 0:
      while (rs.next()) {
        line ++;
        int a = rs.getInt(1);
        int b = rs.getInt(2);
        System.out.println("Line:" + line + ":" + a + "," + b);
     }
     statement.execute("SET dli.sql.spark.sql.forcePartitionPredicatesOnPartitionedTable.enabled=true");
     statement.execute("describe tb1");
      ResultSet rs1 = statement.getResultSet();
     line = 0:
     while (rs1.next()) {
        line ++;
        String a = rs1.getString(1);
        String b = rs1.getString(2);
        System.out.println("Line:" + line + ":" + a + "," + b);
     3
   } catch (SQLException ex) {
   } finally {
     if (conn != null) {
        conn.close();
     }
   }
}
```

# 10.3 APIs Supported By the DLI JDBC Driver

The DLI JDBC driver supports multiple APIs of the JDBC standard, but some APIs cannot be invoked by users. For example, when transaction-related API **prepareCall** is invoked, the **SQLFeatureNotSupportedException** exception is reported. For details about the APIs, see the JDBC official website https://docs.oracle.com/javase/8/docs/api/java/sql/package-summary.html.

# **Supported APIs**

The following tables list the APIs supported by the DLI JDBC driver and provide remarks on possible incompatibilities with the JDBC standard.

- Common signatures supported by Connection APIs
  - Statement createStatement()
  - PreparedStatement prepareStatement(String sql)
  - void close()
  - boolean isClosed()
  - DatabaseMetaData getMetaData()
  - PreparedStatement prepareStatement(String sql, int resultSetType, int resultSetConcurrency)
- Common signatures supported by Driver APIs
  - Connection connect(String url, Properties info)
  - boolean acceptsURL(String url)
  - DriverPropertyInfo[] getPropertyInfo(String url, Properties info)
- Common signatures supported by Connection APIs
  - String getColumnClassName(int column)
  - int getColumnCount()
  - int getColumnDisplaySize(int column)
  - String getColumnLabel(int column)
  - String getColumnName(int column)
  - int getColumnType(int column)
  - String getColumnTypeName(int column)
  - int getPrecision(int column)
  - int getScale(int column)
  - boolean isCaseSensitive(int column)
- Common signatures supported by Statement APIs
  - ResultSet executeQuery(String sql)
  - int executeUpdate(String sql)
  - boolean execute(String sql)
  - void close()
  - int getMaxRows()
  - void setMaxRows(int max)
  - int getQueryTimeout()
  - void setQueryTimeout(int seconds)
  - void cancel()
  - ResultSet getResultSet()
  - int getUpdateCount()
  - boolean isClosed()
- Common signatures supported by PreparedStatement APIs

- void clearParameters()
- boolean execute()
- ResultSet executeQuery()
- int executeUpdate()
- PreparedStatement Set methods
- Common signatures supported by ResultSet APIs
  - int getRow()
  - boolean isClosed()
  - boolean next()
  - void close()
  - int findColumn(String columnLabel)
  - boolean wasNull()
  - Get methods
- Common signatures supported by DatabaseMetaData APIs
  - ResultSet getCatalogs()

**NOTE** 

- DLI does not have the concept of Catalog, so an empty ResultSet is returned.
- ResultSet getColumns(String catalog, String schemaPattern, String tableNamePattern, String columnNamePattern)
- Connection getConnection()
- getTables(String catalog, String schemaPattern,String tableNamePattern, String types[])

**NOTE** 

This method does not use the **Catalog** parameter, and schemaPattern corresponds to the database concept of DLI.

- ResultSet getTableTypes()
- ResultSet getSchemas()
- ResultSet getSchemas(String catalog, String schemaPattern)

# **11** Submitting a Flink Job on the DLI Management Console

# **11.1 Flink Job Overview**

DLI supports two types of Flink jobs:

- Flink OpenSource SQL job:
  - It is fully compatible with Flink of the community edition, ensuring that jobs can run smoothly on these Flink versions.
  - DLI Flink has expanded the support for connectors based on Flink of the community edition, supporting Redis and GaussDB(DWS) as new data source types. With this expansion, you can now utilize a wider range of data source types, providing greater flexibility and convenience when working with datasets.
  - Flink OpenSource SQL jobs are ideal for scenarios where stream processing logic can be defined and executed through SQL statements. This simplifies stream processing, allowing developers to focus more on implementing service logic.

For how to create a Flink OpenSource SQL job, see **Creating a Flink OpenSource SQL Job**.

- Flink Jar job:
  - DLI allows you to submit Flink jobs compiled into JAR files, providing higher flexibility and customization capabilities. It is applicable to scenarios where complex data processing is required.
  - If the connectors provided by Flink of the community edition cannot meet specific needs, you can use Jar jobs to implement custom connectors or data processing logic.
  - It is ideal for scenarios where user-defined functions (UDFs) or specific library integration are required. You can use the Flink ecosystem to implement advanced stream processing logic and status management.

For how to create a Flink Jar job, see Creating a Flink Jar Job.

# 11.2 Creating a Flink OpenSource SQL Job

This section describes how to create a Flink OpenSource SQL job.

DLI Flink OpenSource SQL jobs are fully compatible with the syntax of Flink provided by the community. In addition, Redis and GaussDB(DWS) data source types are added based on the community connector. For the syntax and constraints of Flink SQL DDL, DML, and functions, see **Table API & SQL**.

- For the Flink OpenSource SQL 1.15 syntax, see Flink OpenSource SQL 1.15 Syntax.
- For the Flink OpenSource SQL 1.12 syntax, see Flink OpenSource SQL 1.12 Syntax.

# Prerequisites

- You have prepared the source and sink streams.
- A datasource connection has been created to enable the network between the queue where the job is about to run and external data sources.
  - For details about the external data sources that can be accessed by Flink jobs, see Common Development Methods for DLI Cross-Source Analysis.
  - For how to create a datasource connection, see Configuring the Network Connection Between DLI and Data Sources (Enhanced Datasource Connection).

On the **Resources** > **Queue Management** page, locate the queue you have created, click **More** in the **Operation** column, and select **Test Address Connectivity** to check if the network connection between the queue and the data source is normal. For details, see **Testing Address Connectivity**.

# Precautions

Before creating jobs and submitting tasks, you are advised to enable CTS to record DLI operations for queries, audits, and tracking. **Using CTS to Audit DLI** lists DLI operations that can be recorded by CTS.

For how to enable CTS and view trace details, see **Cloud Trace Service Getting Started**.

# Creating a Flink OpenSource SQL Job

- Step 1 In the left navigation pane of the DLI management console, choose Job Management > Flink Jobs. The Flink Jobs page is displayed.
- Step 2 In the upper right corner of the Flink Jobs page, click Create Job.

Figure 11-1 Creating a Flink OpenSource SQL job

Create Job		×
Туре	Flink OpenSource SQL	
* Name	testFlinkJob	
Description	Description	
Template Name	-Select	
Tags	It is recommended that you use TMS's predefined tag function to add the same tag to different cloud resources. View predefined tags C To add a tag, enter a tag key and a tag value below.	
	Enter a tag key     Enter a tag value     Add	
	10 tags available for addition.           OK         Cancel	

# **Step 3** Set job parameters.

#### Table 11-1 Job parameters

Parameter	Description
Туре	Set <b>Type</b> to <b>Flink OpenSource SQL</b> . You will need to start jobs by compiling SQL statements.
Name	Job name. The value can contain up to 57 characters. Only letters, digits, hyphens (-), and underscores (_) are allowed. The job name must be globally unique.
Description	Job description. It can contain up to 512 characters.
Template Name	You can select a sample template or a custom job template. For details about templates, see Managing Flink Job Templates.

Parameter	Description
Tags	Tags used to identify cloud resources. A tag includes the tag key and tag value. If you want to use the same tag to identify multiple cloud resources, that is, to select the same tag from the drop-down list box for all services, you are advised to create predefined tags on the Tag Management Service (TMS).
	If your organization has configured tag policies for DLI, add tags to resources based on the policies. If a tag does not comply with the tag policies, resource creation may fail. Contact your organization administrator to learn more about tag policies.
	For details, see Tag Management Service User Guide.
	NOTE
	• A maximum of 20 tags can be added.
	Only one tag value can be added to a tag key.
	The key name in each resource must be unique.
	• Tag key: Enter a tag key name in the text box.
	<b>NOTE</b> A tag key can contain a maximum of 128 characters. Only letters, numbers, spaces, and special characters (_::+-@) are allowed, but the value cannot start or end with a space or start with <b>_sys_</b> .
	• Tag value: Enter a tag value in the text box.
	NOTE A tag value can contain a maximum of 255 characters. Only letters, numbers, spaces, and special characters ( .:+-@) are allowed.

- **Step 4** Click **OK** to enter the editing page.
- **Step 5** Edit an OpenSource SQL job.

Enter detailed SQL statements in the statement editing area. For details about SQL statements, see **Data Lake Insight Flink OpenSource SQL Syntax Reference**.

#### Step 6 Click Check Semantics.

- You can **Start** a job only after the semantic verification is successful.
- If verification is successful, the message "The SQL semantic verification is complete. No error." will be displayed.
- If verification fails, a red "X" mark will be displayed in front of each SQL statement that produced an error. You can move the cursor to the "X" mark to view error details and change the SQL statement as prompted.

#### D NOTE

Flink 1.15 does not support syntax verification.

#### **Step 7** Set job running parameters.

Parameter	Description
Queue	Resource queue used to execute Flink jobs. A queue determines the compute resources accessible to a job during its operation within an elastic resource pool. Every queue is allocated with specific resources, known as CUs, whose configuration significantly impacts the job's performance and execution efficiency.
	Before submitting a job, assess its resource needs and select an appropriate queue. Flink OpenSource SQL jobs support selecting <b>For general</b> <b>purpose</b> queues.
Flink Version	<ul> <li>Flink version used for job running. Flink versions have varying feature support.</li> <li>If you choose to use Flink 1.15, make sure to configure the agency information for the cloud service that DLI is allowed to access in the job.</li> <li>For the syntax of Flink 1.15, see Flink OpenSource SQL 1.15</li> <li>Usage and Flink OpenSource SQL 1.15 Syntax.</li> <li>For the syntax of Flink 1.12, see Flink OpenSource SQL 1.12</li> <li>Syntax.</li> <li>NOTE <ul> <li>You are advised not to use Flink of different versions for a long time.</li> <li>Doing so can lead to code incompatibility, which can negatively impact job execution efficiency.</li> <li>Doing so may result in job execution failures due to conflicts in dependencies. Jobs rely on specific versions of libraries or components.</li> </ul> </li> </ul>
UDF Jar	<ul> <li>UDF JAR file, which contains UDFs that can be called in subsequent jobs.</li> <li>There are the following ways to manage UDF JAR files:</li> <li>Upload packages to OBS: Upload Jar packages to an OBS bucket in advance and select the corresponding OBS path.</li> <li>Upload packages to DLI: Upload JAR files to an OBS bucket in advance and create a package on the Data Management &gt; Package Management page of the DLI management console. For details, see Creating a DLI Package.</li> <li>For Flink 1.15 or later, only OBS packages can be selected when creating jobs, and DLI packages are not supported.</li> </ul>
Agency	If you choose Flink 1.15 or later to execute your job, you can create a custom agency to allow DLI to access other services. For how to create a custom agency, see <b>Creating a Custom DLI Agency</b> .

Parameter	Description	
Resource Configuration	DLI offers various resource configuration templates based on different Flink engine versions.	
Version	Compared with the v1 template, the v2 template does not support the setting of the number of CUs. The v2 template supports the setting of <b>Job Manager Memory</b> and <b>Task</b> <b>Manager Memory</b> .	
	<b>v1</b> : applicable to Flink 1.12, 1.13, and 1.15.	
	<b>v2</b> : applicable to Flink 1.13, 1.15, and 1.17.	
	You are advised to use the parameter settings of v2.	
	For details about the parameters of v1, see Table 11-3.	
	For details about the parameters of v2, see <b>Table 11-4</b> .	

<b>Table 11-3</b> Resource specification parameters of V	Table 11-3	Resource	specification	parameters	of v1
--	------------	----------	---------------	------------	-------

Parameter	Description	
CUs	Sum of the number of compute units and job manager CUs of DLI. CU is also the billing unit of DLI. One CU equals 1 vCPU and 4 GB.	
	The value is the number of CUs required for job running and cannot exceed the number of CUs in the bound queue.	
	NOTE When Task Manager Config is selected, elastic resource pool queue management is optimized by automatically adjusting CUs to match Actual CUs after setting Slot(s) per TM.	
	CUs = Actual number of CUs = max[Job Manager CPUs + Task Manager CPU, (Job Manager Memory + Task Manager Memory/4)]	
	<ul> <li>Job Manager CPUs + Task Manager CPUs = Actual TMs x CU(s) per TM + Job Manager CUs.</li> </ul>	
	<ul> <li>Job Manager Memory + Task Manager Memory = Actual TMs x Memory per TM + Job Manager Memory</li> </ul>	
	<ul> <li>If Slot(s) per TM is set, then: Actual TMs = Parallelism/Slot(s) per TM.</li> </ul>	
	<ul> <li>If Slot(s) per TM is not set, then: Actual TMs = (CUs – Job Manager CUs)/CU(s) per TM.</li> </ul>	
	• If <b>Memory per TM</b> and <b>Job Manager Memory</b> in the optimization parameters are not set, then: Memory per TM = CU(s) per TM x 4. Job Manager Memory = Job Manager CUs x 4.	
	• The parallelism degree of Spark resources is jointly determined by the number of Executors and the number of Executor CPU cores.	
Job Manager CUs	Number of CUs of the management unit.	

Parameter	Description	
Parallelism	Number of tasks concurrently executed by each operator in a job. <b>NOTE</b> This value cannot be greater than four times the compute units (number of CUs minus the number of job manager CUs).	
Task Manager Config	<ul> <li>Whether Task Manager resource parameters are set</li> <li>If selected, you need to set the following parameters: <ul> <li>CU(s) per TM: Number of resources occupied by each Task Manager.</li> <li>Slot(s) per TM: Number of slots contained in each Task Manager.</li> </ul> </li> <li>If not selected, the system automatically uses the default values. <ul> <li>CU(s) per TM: The default value is 1.</li> <li>Slot(s) per TM: The default value is (Parallelism x CU(s) per TM)/(CUs – Job Manager CUs).</li> </ul> </li> </ul>	
OBS Bucket	OBS bucket to store job logs and checkpoint information. If the OBS bucket you selected is unauthorized, click <b>Authorize</b> .	
Save Job Log	<ul> <li>Whether job running logs are saved to OBS. The logs are saved in the following path: <i>Bucket name</i>/jobs/logs/<i>Directory starting with the job ID</i>.</li> <li>CAUTION <ul> <li>You are advised to configure this parameter. Otherwise, no run log is generated after the job is executed. If the job fails, the run log cannot be obtained for fault locating.</li> <li>If this option is selected, you need to set the following parameters:</li> <li>OBS Bucket: Select an OBS bucket to store job logs. If the OBS bucket you selected is unauthorized, click Authorize.</li> </ul> </li> <li>NOTE <ul> <li>If Enable Checkpointing and Save Job Log are both selected, you only need to authorize OBS once.</li> </ul> </li> </ul>	
Alarm on Job Exception	<ul> <li>Whether to notify users of any job exceptions, such as running exceptions or arrears, via SMS or email.</li> <li>If this option is selected, you need to set the following parameters:</li> <li>SMN Topic</li> <li>Select a custom SMN topic. For how to create a custom SMN topic, see Creating a Topic.</li> </ul>	

Parameter	Description
Enable Checkpointing	Whether to enable job snapshots. If this function is enabled, jobs can be restored based on the checkpoints.
	If this option is selected, you need to set the following parameters:
	• <b>Checkpoint Interval</b> : interval for creating checkpoints, in seconds. The value ranges from 1 to 999999, and the default value is <b>30</b> .
	• <b>Checkpoint Mode</b> can be set to either of the following values:
	<ul> <li>At least once: Events are processed at least once.</li> </ul>
	- <b>Exactly once</b> : Events are processed only once.
	OBS Bucket: Select an OBS bucket to store your checkpoints. If the OBS bucket you selected is unauthorized, click Authorize.  The sheelyneint path is <i>Bucket</i> pame/iebs/checkpoint/
	Directory starting with the job ID.
	NOTE If Enable Checkpointing and Save Job Log are both selected, you only need to authorize OBS once.
Auto Restart upon Exception	Whether automatic restart is enabled. If enabled, jobs will be automatically restarted and restored when exceptions occur.
	If this option is selected, you need to set the following parameters:
	• <b>Max. Retry Attempts</b> : maximum number of retries upon an exception. The unit is times/hour.
	<ul> <li>Unlimited: The number of retries is unlimited.</li> </ul>
	<ul> <li>Limited: The number of retries is user-defined.</li> </ul>
	• <b>Restore Job from Checkpoint</b> : This parameter is available only when <b>Enable Checkpointing</b> is selected.
Idle State Retention Time	Clears intermediate states of operators such as <b>GroupBy</b> , <b>RegularJoin</b> , <b>Rank</b> , and <b>Depulicate</b> that have not been updated after the maximum retention time. The default value is 1 hour.
Dirty Data Policy	Policy for processing dirty data. The following policies are supported: <b>Ignore</b> , <b>Trigger a job exception</b> , and <b>Save</b> .
	If you set this field to <b>Save</b> , the <b>Dirty Data Dump Address</b> must be set. Click the address box to select the OBS path for storing dirty data.
	This parameter is available only when a DIS data source is used.

Parameter	Description
Parallelism	Number of tasks concurrently executed by each operator in a job. NOTE
	<ul> <li>The minimum parallelism must not be less than 1. The default value is 1.</li> <li>This value cannot be greater than four times the compute units.</li> </ul>
	(number of CUs minus the number of JobManager CUs).
Job Manager	Number of CPU cores available for Job Manager.
CPU	The default value is <b>1</b> . The minimum value cannot be less than 0.5.
Job Manager	Number of memory available for Job Manager.
Memory	The default value is 4 GB. The minimum size cannot be less than 2 GB (2048 MB). The default unit is GB, which can be set to GB or MB.
Task Manager	Number of CPU cores available for Task Manager.
CPU	The default value is <b>1</b> . The minimum value cannot be less than 0.5.
Task Manager	Number of memory available for Task Manager.
Memory	The default value is 4 GB. The minimum size cannot be less than 2 GB (2048 MB). The default unit is GB, which can be set to GB or MB.
Slot(s) per TM	Number of parallel tasks that a single Task Manager can support. Each task slot can execute one task in parallel. Increasing task slots enhances the parallel processing capacity of the Task Manager but also increases resource consumption.
	The number of task slots is linked to the CPU count of the Task Manager since each CPU can offer one task slot.
	By default, a single TM slot is set to <b>1</b> . The minimum parallelism must not be less than 1.
OBS Bucket	OBS bucket to store job logs and checkpoint information. If the OBS bucket you selected is unauthorized, click <b>Authorize</b> .

 Table 11-4 Resource specification parameters of v2

Parameter	Description
Save Job Log	Whether job running logs are saved to OBS. The logs are saved in the following path: <i>Bucket name/jobs/logs/Directory starting with the job ID</i> .
	<b>CAUTION</b> You are advised to configure this parameter. Otherwise, no run log is generated after the job is executed. If the job fails, the run log cannot be obtained for fault locating.
	If this option is selected, you need to set the following parameters:
	<b>OBS Bucket</b> : Select an OBS bucket to store job logs. If the OBS bucket you selected is unauthorized, click <b>Authorize</b> .
	NOTE If Enable Checkpointing and Save Job Log are both selected, you only need to authorize OBS once.
Alarm on Job Exception	Whether to notify users of any job exceptions, such as running exceptions or arrears, via SMS or email.
	If this option is selected, you need to set the following parameters:
	SMN Topic
	Select a custom SMN topic. For how to create a custom SMN topic, see <b>Creating a Topic</b> .
Enable Checkpointing	Whether to enable job snapshots. If this function is enabled, jobs can be restored based on the checkpoints.
	If this option is selected, you need to set the following parameters:
	• <b>Checkpoint Interval</b> : interval for creating checkpoints, in seconds. The value ranges from 1 to 999999, and the default value is <b>30</b> .
	• <b>Checkpoint Mode</b> can be set to either of the following values:
	<ul> <li>At least once: Events are processed at least once.</li> </ul>
	<ul> <li>Exactly once: Events are processed only once.</li> </ul>
	<ul> <li>OBS Bucket: Select an OBS bucket to store your checkpoints. If the OBS bucket you selected is unauthorized, click Authorize.</li> <li>The checkpoint path is <i>Bucket name/jobs/checkpoint/ Directory starting with the job ID</i>.</li> </ul>
	NOTE If Enable Checkpointing and Save Job Log are both selected, you only need to authorize OBS once.

Parameter	Description
Auto Restart upon Exception	Whether automatic restart is enabled. If enabled, jobs will be automatically restarted and restored when exceptions occur.
	If this option is selected, you need to set the following parameters:
	• <b>Max. Retry Attempts</b> : maximum number of retries upon an exception. The unit is times/hour.
	<ul> <li>Unlimited: The number of retries is unlimited.</li> </ul>
	<ul> <li>Limited: The number of retries is user-defined.</li> </ul>
	• <b>Restore Job from Checkpoint</b> : This parameter is available only when <b>Enable Checkpointing</b> is selected.
Idle State Retention Time	Clears intermediate states of operators such as <b>GroupBy</b> , <b>RegularJoin</b> , <b>Rank</b> , and <b>Depulicate</b> that have not been updated after the maximum retention time. The default value is 1 hour.
Dirty Data Policy	Policy for processing dirty data. The following policies are supported: <b>Ignore</b> , <b>Trigger a job exception</b> , and <b>Save</b> .
	If you set this field to <b>Save</b> , the <b>Dirty Data Dump Address</b> must be set. Click the address box to select the OBS path for storing dirty data.
	This parameter is available only when a DIS data source is used.

**Step 8** (Optional) Set the runtime configuration as required. For details about related parameters, see **How Do I Optimize the Performance of a Flink Job?** 

Figure 11-2 Runtime configuration



You can set compute resource specification parameters on the **Runtime Configuration** tab of Flink jobs, and the parameter values have a higher priority than the specified values.

 Table 11-5 describes the parameter mapping.

#### D NOTE

In Flink 1.12, you are advised to set compute resource specification parameters based on the configuration method on the console. Using automatic parameter settings may result in discrepancies in actual CU statistics.

Table 11-5 Mapping between compute resource specification parameters on the
console and those on the Runtime Configuration tab

Runtime Configuration	Compute Resource Specificat ion Paramete r of v1	Compute Resource Specificati on Parameter of v2	Description
kubernetes.jobmanag er.cpu	Job Manager CUs	Job Manager CPU	Number of CPU cores available for Job Manager. The default value is <b>1</b> . The minimum value cannot be less than 0.5.
kubernetes.taskmana ger.cpu	CU(s) per TM	Task Manager CPU	Number of CPU cores available for Task Manager. The default value is <b>1</b> . The minimum value cannot be less than 0.5.
jobmanager.memory.p rocess.size	-	Job Manager Memory	Number of memory available for Job Manager. The default value is 4 GB. The minimum size cannot be less than 2 GB (2048 MB). The default unit is GB, which can be set to GB or MB.
taskmanager.memory. process.size	-	Task Manager Memory	Number of memory available for Task Manager. The default value is 4 GB. The minimum size cannot be less than 2 GB (2048 MB). The default unit is GB, which can be set to GB or MB.

#### Step 9 Click Save.

**Step 10** Click **Start**. On the displayed **Start Flink Jobs** page, confirm the job specifications and the price, and click **Start Now** to start the job.

After the job is started, the system automatically switches to the **Flink Jobs** page, and the created job is displayed in the job list. You can view the job status in the **Status** column. Once a job is successfully submitted, its status changes from **Submitting** to **Running**. After the execution is complete, the status changes to **Completed**.

If the job status is **Submission failed** or **Running exception**, the job fails to submit or run. In this case, you can hover over the status icon in the **Status** column of the job list to view the error details. You can click I to copy these details. Rectify the fault based on the error information and resubmit the job.

#### **NOTE**

Other buttons are as follows:

- Save As: Save the created job as a new job.
- **Static Stream Graph**: Provide the static concurrency estimation function and stream graph display function. See Figure 11-4.
- **Simplified Stream Graph**: Display the data processing flow from the source to the sink. See **Figure 11-3**.
- Format: Format the SQL statements in the editing box.
- Set as Template: Set the created SQL statements as a job template.
- Theme Settings: Set the theme related parameters, including Font Size, Wrap, and Page Style.
- Help: Redirect to the help center to provide you with the SQL syntax for stream jobs.

----End

## Simplified Stream Graph

On the OpenSource SQL job editing page, click Simplified Stream Graph.

**NOTE** 

Simplified stream graph viewing is only supported in Flink 1.12 and Flink 1.10.

#### Figure 11-3 Simplified stream graph



# Static Stream Graph

On the OpenSource SQL job editing page, click Static Stream Graph.

### D NOTE

- Simplified stream graph viewing is only supported in Flink 1.12 and Flink 1.10.
- If you use a UDF in a Flink OpenSource SQL job, it is not possible to generate a static stream graph.

The Static Stream Graph page also allows you to:

- Estimate concurrencies. Click **Estimate Concurrencies** on the **Static Stream Graph** page to estimate concurrencies. Click **Restore Initial Value** to restore the initial value after concurrency estimation.
- Zoom in or out the page.
- Expand or merge operator chains.
- You can edit **Parallelism**, **Output rate**, and **Rate factor**.
  - **Parallelism**: indicates the number of concurrent tasks.
  - Output rate: indicates the data traffic of an operator. The unit is piece/s.
  - Rate factor: indicates the retention rate after data is processed by operators. Rate factor = Data output volume of an operator/Data input volume of the operator (Unit: %)

Figure 11-4 Static stream graph

Stopped] ID: 102726 Jab Type: Fliet OpenSource SOL	Start Save	Save As Static Stream	Graph
Check Semantics Simplified Stream Graph Format Save as Template Theme Settings Help			î a
1 CRAT THE order ( order_thatrig, 1 order_thatrig, 2 order_thatrig, 3 order_thatrig, 4 order_thatrig, 5 pay_emont double, 4 real_spatial, 5 pay_emont double, 4 real_spatial, 5 pay_emont double, 5 real_spatial, 5 pay_emont double, 5 real_spatial, 5 pay_emont double, 6 real_spatial, 8 real_spatial, 9 real_spatial, 9 real_spatial, 10 re	CUIS     CUIS     CUIS     Paratleism     Task Manager Cong	Competence         •           112         •           -Select-         •           -         2         +         ①           -         1         +         ●           -         1         +         ●	uning Persenters Runtime Configuration
18 "format" - 'fson'			
Static Stream Graph Save	Edit Estimate Concu	irrencies Restore Initial Valu	ye X
Vortex (4 Operators)           Source: TableSourceScantzable-(Selsuc, catalog, education), education (selsuc, education), education (sels	le+[default_catalog default_data m 1 scords/s)	base o	

# 11.3 Creating a Flink Jar Job

A Flink Jar job involves developing a custom application Jar package based on Flink's capabilities and submitting it to a DLI queue for execution.

To create a Flink Jar job, you need to write and build your own application Jar package. This is suitable for users who require stream data processing and are proficient in Flink's secondary development capabilities.

This section describes how to create a Flink Jar job on the DLI management console.

## Prerequisites

• When you use a Flink Jar job to access other external data sources, such as OpenTSDB, HBase, Kafka, GaussDB(DWS), RDS, CSS, CloudTable, DCS Redis,

and DDS, you need to create a datasource connection to connect the job running queue to the external data source.

- For details about the external data sources that can be accessed by Flink jobs, see Common Development Methods for DLI Cross-Source Analysis.
- For how to create a datasource connection, see Configuring the Network Connection Between DLI and Data Sources (Enhanced Datasource Connection).

On the **Resources** > **Queue Management** page, locate the queue you have created, click **More** in the **Operation** column, and select **Test Address Connectivity** to check if the network connection between the queue and the data source is normal. For details, see **Testing Address Connectivity**.

- To run a Flink Jar job, you need to build your custom application code into a JAR file and upload it to the OBS bucket that has already been created.
- Flink dependencies have been built in the DLI server and security hardening has been performed based on the open-source community version. To avoid dependency package compatibility issues or log output and dump issues, be careful to exclude the following files when packaging:
  - Built-in dependencies (or set the package dependency scope to **provided** in Maven or SBT)
  - Log configuration files (example, **log4j.properties**/**logback.xml**)
  - JAR package for log output implementation (example, **log4j**).

# Precautions

Before creating jobs and submitting tasks, you are advised to enable CTS to record operations associated with DLI for later query, audit, and backtrack operations. To view the DLI operations that can be recorded by CTS, see Using CTS to Audit DLI.

For how to enable CTS and view trace details, see **Cloud Trace Service Getting Started**.

# Creating a Flink Jar Job

- Step 1 In the left navigation pane of the DLI management console, choose Job Management > Flink Jobs. The Flink Jobs page is displayed.
- **Step 2** In the upper right corner of the **Flink Jobs** page, click **Create Job**.

Figure 11-5 Creating a Flink Jar job

Create Job		×
Туре	Flink Jar 💌	
* Name	Enter a name.	
Description	Description	
Tags	It is recommended that you use TMS's predefined tag function to add the same tag to different cloud resources. View predefined tags C To add a tag, enter a tag key and a tag value below.	
	Enter a tag key         Enter a tag value         Add	
	10 tags available for addition.           OK         Cancel	

**Step 3** Specify job parameters.

 Table 11-6 Job configuration information

Paramet er	Description
Туре	Select <b>Flink Jar</b> .
Name	Job name. The value can contain up to 57 characters. Only letters, digits, hyphens (-), and underscores (_) are allowed. <b>NOTE</b> The job name must be globally unique.
Descripti on	Job description. It can contain up to 512 characters.

Paramet er	Description
Tags	Tags used to identify cloud resources. A tag includes the tag key and tag value. If you want to use the same tag to identify multiple cloud resources, that is, to select the same tag from the drop-down list box for all services, you are advised to create predefined tags on the Tag Management Service (TMS).
	If your organization has configured tag policies for DLI, add tags to resources based on the policies. If a tag does not comply with the tag policies, resource creation may fail. Contact your organization administrator to learn more about tag policies.
	For details, see Tag Management Service User Guide.
	NOTE
	• A maximum of 20 tags can be added.
	<ul> <li>Only one tag value can be added to a tag key.</li> </ul>
	• The key name in each resource must be unique.
	• Tag key: Enter a tag key name in the text box.
	<b>NOTE</b> A tag key can contain a maximum of 128 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed, but the value cannot start or end with a space or start with <b>_sys_</b> .
	• Tag value: Enter a tag value in the text box.
	<b>NOTE</b> A tag value can contain a maximum of 255 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed.

**Step 4** Click **OK** to enter the editing page.

Step 5 Select a queue.

**Step 6** Configuring Flink Jar Job parameters

★ Queue	-Select-			*	
* Application (?)	-Select				View Built-in Dependencies
Main Class	Default	Manually	assign		
	Default main class	is specified by	the Manifest	file of the	application.
Class Arguments	Enter a class arg arguments with s	jument (Separa spaces).	te multiple cl	ass	
				4	
JAR Package Dependenc	Select			•	View Built-in Dependencies
Other Dependencies (?)	Select			•	
Job Type	Basic	Image			
Flink Version	Select			*	Select a queue first.
Runtime Configuration	Enter arguments Press Enter to se	using the key = eparate multiple	= value forma e key-value pa	at. airs.	
				11	

Figure 11-6 Configuring Flink Jar Job parameters

## Table 11-7 Parameters

Parameter	Description	
Queue	Select a queue where you want to run your job.	
Flink Version	Flink version used for job running. Flink versions have varying feature support.	
	If you choose to use Flink 1.15, make sure to configure the agency information for the cloud service that DLI is allowed to access in the job.	
	For the syntax of Flink 1.15, see <b>Flink OpenSource SQL 1.15</b> Usage and Flink OpenSource SQL 1.15 Syntax.	
	For the syntax of Flink 1.12, see <b>Flink OpenSource SQL 1.12</b> <b>Syntax</b> .	
	<b>NOTE</b> You are advised not to use Flink of different versions for a long time.	
	<ul> <li>Doing so can lead to code incompatibility, which can negatively impact job execution efficiency.</li> </ul>	
	<ul> <li>Doing so may result in job execution failures due to conflicts in dependencies. Jobs rely on specific versions of libraries or components.</li> </ul>	

Parameter	Description
Application	Select a Jar job package. There are the following ways to manage JAR files:
	<ul> <li>Upload packages to OBS: Upload Jar packages to an OBS bucket in advance and select the corresponding OBS path.</li> </ul>
	<ul> <li>Upload packages to DLI: Upload JAR files to an OBS bucket in advance and create a package on the Data Management &gt; Package Management page of the DLI management console. For details, see Creating a DLI Package.</li> </ul>
	For Flink 1.15 or later, you can only select packages from OBS, instead of DLI.
Main Class	The name of the JAR package to be loaded, for example, <b>KafkaMessageStreaming</b> .
	<ul> <li>Default: Specified based on the Manifest file in the JAR package.</li> </ul>
	• <b>Manually assign</b> : You must enter the class name and confirm the class arguments (separated by spaces).
	When a class belongs to a package, the main class path must contain the complete package path, for example, <b>packagePath.KafkaMessageStream-ing</b> .
Class Arguments	List of arguments of a specified class. The arguments are separated by spaces.
	Flink parameters support replacement of non-sensitive global variables. For example, if you add the global variable windowsize in Global Configuration > Global Variables, you can add the - windowsSize {{windowsize}} parameter for the Flink Jar job.
JAR Package Dependenc	Select a user-defined package dependency. The dependent program packages are stored in the classpath directory of the cluster.
ies	There are the following ways to manage JAR files:
	<ul> <li>Upload packages to OBS: Upload Jar packages to an OBS bucket in advance and select the corresponding OBS path.</li> </ul>
	<ul> <li>Upload packages to DLI: Upload JAR files to an OBS bucket in advance and create a package on the Data Management &gt; Package Management page of the DLI management console. For details, see Creating a DLI Package.</li> </ul>
	For Flink 1.15 or later, you can only select packages from OBS, instead of DLI.
	When creating a JAR file for a Flink Jar job, you do not need to upload existing built-in dependency packages to avoid package information conflicts.
	For details about built-in dependency packages, see <b>DLI Built-in</b> <b>Dependencies</b> .

Parameter	Description		
Other Dependenc	User-defined dependency files. Other dependency files need to be referenced in the code.		
ies	There are the following ways to manage dependency files:		
	• Upload packages to OBS: Upload dependency files to an OBS bucket in advance and select the corresponding OBS path.		
	<ul> <li>Upload packages to DLI: Upload dependency files to an OBS bucket in advance and create a package on the Data Management &gt; Package Management page of the DLI management console. For details, see Creating a DLI Package.</li> </ul>		
	For Flink 1.15 or later, you can only select packages from OBS, instead of DLI.		
	You can add the following command to the application to access the corresponding dependency file. In the command, <b>fileName</b> indicates the name of the file to be accessed, and <b>ClassName</b> indicates the name of the class that needs to access the file. ClassName.class.getClassLoader().getResource("userData/fileName")		
Job Type	Image type used for creating a Flink Jar job. It is used to specify the image type of the DLI container cluster.		
	• <b>Basic</b> : base image provided by DLI, which is selected by default.		
	• Image: Select the image name and image version. Image set on the Software Repository for Container (SWR) console. For details, see Enhancing the Job Runtime Environment Using a Custom Image.		
Agency	If you choose Flink 1.15 or later to execute your job, you can create a custom agency to allow DLI to access other services.		
	For how to create a custom agency, see <b>Creating a Custom DLI Agency</b> .		

Parameter	r Description		
Runtime Configurati	User-defined optimization parameters. The parameter format is <b>key=value</b> .		
on	Flink optimization parameters support replacement non-sensitive global variable. For example, if you create global variable <b>phase</b> in <b>Global Configuration</b> > <b>Global Variables</b> , optimization parameter <b>table.optimizer.agg-phase.strategy={{phase}}</b> can be added to the Flink Jar job.		
	Flink 1.15 supports minimal submission of Flink Jar jobs. Enable this by configuring <b>flink.dli.job.jar.minimize-</b> <b>submission.enabled=true</b> in the runtime optimization		
	NOTE Minimal submission means Flink only submits the necessary job dependencies, not the entire Flink environment. By setting the scope of non-Connector Flink dependencies (starting with <b>flink-</b> ) and third-party libraries (like Hadoop, Hive, Hudi, and MySQL-CDC) to <b>provided</b> , you ensure these dependencies are excluded from the Jar job, avoiding conflicts with Flink core dependencies.		
	Only Flink 1.15 supports minimal submission of Flink Jar jobs.		
	• For Flink-related dependencies, use the <b>provided</b> scope by adding < <b>scope&gt;provided</b> in the dependencies, especially for non-Connector dependencies under the <b>org.apache.flink</b> group starting with <b>flink</b>		
	<ul> <li>For dependencies related to Hadoop, Hive, Hudi, and MySQL-CDC, also use the provided scope by adding <scope>provided</scope> in the dependencies.</li> </ul>		
	<ul> <li>In the Flink source code, only methods marked with @Public or @PublicEvolving are intended for user invocation. DLI guarantees compatibility with these methods.</li> </ul>		

**Step 7** Set compute resource specification parameters.

	* Flink Version	-Select-   Select a queue first.
	Runtime Configuration	Enter arguments using the key = value format. Press Enter to separate multiple key-value pairs.
8		
	* CUs	2 + Total number of CUs occupied by a job. The value must be the same
	★ Job Manager CUs	- 1 +
	★ Parallelism	1 + Default maximum number of operators that can be concurrently ex- times the number of CUs. However, it is recommended that you set
	Task Manager Configuration	
	Save Job Log	
	★ OBS Bucket	dli-ae-ad-1-0acce4eece805a022f06c007f6c086b8
	Alarm Generation upon J	
	Auto Restart upon Except	

Figure 11-7 Configuring job parameters

DLI offers various resource configuration templates based on different Flink engine versions.

Compared with the v1 template, the v2 template does not support the setting of the number of CUs. The v2 template supports the setting of **Job Manager Memory** and **Task Manager Memory**.

v1: applicable to Flink 1.12, 1.13, and 1.15.

**v2**: applicable to Flink 1.13, 1.15, and 1.17.

You are advised to use the parameter settings of v2.

For details about the parameters of v1, see Table 11-8.

For details about the parameters of v2, see Table 11-9.

Parameter Description			
CUs	One CU consists of one vCPU and 4 GB of memory. The number of CUs ranges from 2 to 10000.		
	When Task Manager Config is selected, elastic resource pool queue management is optimized by automatically adjusting CUs to match Actual CUs after setting Slot(s) per TM.		
	CUs = Actual number of CUs = max[Job Manager CPUs + Task Manager CPU, (Job Manager Memory + Task Manager Memory/4)]		
	<ul> <li>Job Manager CPUs + Task Manager CPUs = Actual TMs x CU(s) per TM + Job Manager CUs.</li> </ul>		
	<ul> <li>Job Manager Memory + Task Manager Memory = Actual TMs x Memory per TM + Job Manager Memory</li> </ul>		
	• If <b>Slot(s) per TM</b> is set, then: Actual TMs = Parallelism/Slot(s) per TM.		
	<ul> <li>If Slot(s) per TM is not set, then: Actual TMs = (CUs – Job Manager CUs)/CU(s) per TM.</li> </ul>		
	<ul> <li>If Memory per TM and Job Manager Memory in the optimization parameters are not set, then: Memory per TM = CU(s) per TM x 4. Job Manager Memory = Job Manager CUs x 4.</li> </ul>		
	• The parallelism degree of Spark resources is jointly determined by the number of Executors and the number of Executor CPU cores.		
Job Manager CUs	Number of CUs for the job management unit.		
Parallelism	Number of tasks concurrently executed by each operator in a job.		
	<ul> <li>The value cannot exceed four times the number of compute units (CUs         – Job Manager CUs).</li> </ul>		
	<ul> <li>Set this parameter to a value greater than that configured in the code to avoid job submission failures.</li> </ul>		
Task	Whether Task Manager resource parameters are set		
Manager Config	<ul> <li>If this option is selected, you need to set the following parameters:</li> </ul>		
	<ul> <li>CU(s) per TM: Number of resources occupied by each Task Manager.</li> </ul>		
	<ul> <li>Slot(s) per TM: Number of slots contained in each Task Manager.</li> </ul>		
	<ul> <li>If not selected, the system automatically uses the default values.</li> </ul>		
	- <b>CU(s) per TM</b> : The default value is <b>1</b> .		
	<ul> <li>Slot(s) per TM: The default value is (Parallelism x CU(s) per TM)/(CUs – Job Manager CUs).</li> </ul>		

Table 11-8 Parameters descriptions of	v1
---------------------------------------	----

Parameter	Description
Save Job Log	Whether to save the job running logs to the OBS bucket. CAUTION
	after the job is executed. If the job is abnormal, the run log cannot be obtained for fault locating.
	If this option is selected, you need to set the following parameters:
	<b>OBS Bucket</b> : Select an OBS bucket to store job logs. If the selected OBS bucket is not authorized, click <b>Authorize</b> .
Enable Checkpoint ing	Checkpoints are used to periodically save the job state. Enabling checkpointing allows for the quick recovery of a specific job state in case of system failure.
	There are two ways to enable checkpointing in DLI:
	<ul> <li>Configure checkpoint-related parameters in the job code, suitable for Flink 1.15 or earlier.</li> </ul>
	• Enable checkpointing on the Jar job configuration page of the DLI management console, suitable for Flink 1.15 or later.
	For Flink 1.15, do not configure checkpoint-related parameters both in the job code and the Jar job configuration page. The configurations in the job code have higher priority. Duplicate configurations may lead to the use of incorrect checkpoint paths during abnormal restarts, causing recovery failures or data inconsistencies.
	After selecting <b>Enable Checkpointing</b> , set the following parameters to enable checkpointing:
	Checkpoint Interval: The interval between checkpoints, in seconds.
	• <b>Checkpoint Mode</b> : Select a mode for checkpoints. The options are:
	<ul> <li>At least once: Events are processed at least once.</li> </ul>
	<ul> <li>Exactly once: Events are processed exactly once.</li> </ul>
	CAUTION
	• After selecting <b>Enable Checkpointing</b> , you need to set <b>OBS Bucket</b> to save the checkpoint information. The default checkpoint save path is <i>Bucket name</i> /jobs/checkpoint/ <i>Directory with job ID prefix</i> .
	• Once checkpointing is enabled, do not set checkpoint parameters in the job code, as the parameters configured in the job code have a higher priority than those configured on the job configuration page. Duplicate configurations may cause the job to use incorrect checkpoint paths during abnormal restarts, resulting in recovery failures or data inconsistencies.
	• After enabling checkpointing, if <b>Auto Restart on Exception</b> and <b>Restore</b> <b>Job from Checkpoint</b> are both selected, you do not need to set <b>Checkpoint Path</b> . The system will automatically determine the path based on the <b>Enable Checkpointing</b> configuration.

Parameter	Description
Alarm on Job	Whether to notify users of any job exceptions, such as running exceptions or arrears, via SMS or email.
Exception	If this option is selected, you need to set the following parameters:
	SMN Topic
	Select a custom SMN topic. For how to create a custom SMN topic, see <b>Creating a Topic</b> .
Auto Restart	Whether automatic restart is enabled. If enabled, jobs will be automatically restarted and restored when exceptions occur.
upon Excontion	If this option is selected, you need to set the following parameters:
Exception	• <b>Max. Retry Attempts</b> : maximum number of retries upon an exception. The unit is times/hour.
	<ul> <li>Unlimited: The number of retries is unlimited.</li> </ul>
	<ul> <li>Limited: The number of retries is user-defined.</li> </ul>
	• <b>Restore Job from Checkpoint</b> : Restore the job from the saved checkpoint.
	If you select this parameter, you also need to set <b>Checkpoint</b> <b>Path</b> .
	<b>Checkpoint Path</b> : Select a path for storing checkpoints. This path must match that configured in the application package. Each job must have a unique checkpoint path, or, you will not be able to obtain the checkpoint.
	NOTE
	<ul> <li>If you also select Enable Checkpointing, you do not need to set Checkpoint Path. The system will automatically determine the path based on the Enable Checkpointing configuration.</li> </ul>
	<ul> <li>If you do not select Enable Checkpointing, you need to set Checkpoint Path.</li> </ul>

# Table 11-9 Parameters descriptions of v2

Parameter	Description		
Parallelism	Number of tasks concurrently executed by each operator in a job. <b>NOTE</b>		
	<ul> <li>The minimum parallelism must not be less than 1. The default value is</li> <li>1.</li> </ul>		
	<ul> <li>This value cannot be greater than four times the compute units (number of CUs minus the number of Job Manager CUs).</li> </ul>		
Job Manager CPU	Number of CPU cores available for Job Manager. The default value is <b>1</b> . The minimum value cannot be less than 0.5.		

Parameter	Description
Job Manager Memory	Number of memory available for Job Manager.
	The default value is 4 GB. The minimum size cannot be less than 2 GB (2048 MB). The default unit is GB, which can be set to GB or MB.
Task	Number of CPU cores available for Task Manager.
Manager CPU	The default value is <b>1</b> . The minimum value cannot be less than 0.5.
Task	Number of memory available for Task Manager.
Manager Memory	The default value is 4 GB. The minimum size cannot be less than 2 GB (2048 MB). The default unit is GB, which can be set to GB or MB.
Slot(s) per TM	Number of parallel tasks that a single Task Manager can support. Each task slot can execute one task in parallel. Increasing task slots enhances the parallel processing capacity of the Task Manager but also increases resource consumption.
	The number of task slots is linked to the CPU count of the Task Manager since each CPU can offer one task slot.
	By default, a single TM slot is set to <b>1</b> . The minimum parallelism must not be less than 1.
Save Job	Whether to save the job running logs to the OBS bucket.
Log	<b>CAUTION</b> You are advised to select this parameter. Otherwise, no run log is generated after the job is executed. If the job is abnormal, the run log cannot be obtained for fault locating.
	If this option is selected, you need to set the following parameters:
	<b>OBS Bucket</b> : Select an OBS bucket to store job logs. If the selected OBS bucket is not authorized, click <b>Authorize</b> .
Parameter	Description
-----------------------------	--
Enable Checkpoint ing	Checkpoints are used to periodically save the job state. Enabling checkpointing allows for the quick recovery of a specific job state in case of system failure.
	There are two ways to enable checkpointing in DLI:
	<ul> <li>Configure checkpoint-related parameters in the job code, suitable for Flink 1.15 or earlier.</li> </ul>
	• Enable checkpointing on the Jar job configuration page of the DLI management console, suitable for Flink 1.15 or later.
	For Flink 1.15, do not configure checkpoint-related parameters both in the job code and the Jar job configuration page. The configurations in the job code have higher priority. Duplicate configurations may lead to the use of incorrect checkpoint paths during abnormal restarts, causing recovery failures or data inconsistencies.
	After selecting <b>Enable Checkpointing</b> , set the following parameters to enable checkpointing:
	<ul> <li>Checkpoint Interval: The interval between checkpoints, in seconds.</li> </ul>
	• <b>Checkpoint Mode</b> : Select a mode for checkpoints. The options are:
	<ul> <li>At least once: Events are processed at least once.</li> </ul>
	<ul> <li>Exactly once: Events are processed exactly once.</li> </ul>
	CAUTION
	<ul> <li>After selecting Enable Checkpointing, you need to set OBS Bucket to save the checkpoint information. The default checkpoint save path is Bucket name/jobs/checkpoint/Directory with job ID prefix.</li> </ul>
	• Once checkpointing is enabled, do not set checkpoint parameters in the job code, as the parameters configured in the job code have a higher priority than those configured on the job configuration page. Duplicate configurations may cause the job to use incorrect checkpoint paths during abnormal restarts, resulting in recovery failures or data inconsistencies.
	• After enabling checkpointing, if <b>Auto Restart on Exception</b> and <b>Restore</b> <b>Job from Checkpoint</b> are both selected, you do not need to set <b>Checkpoint Path</b> . The system will automatically determine the path based on the <b>Enable Checkpointing</b> configuration.
OBS Bucket	OBS bucket to store job logs and checkpoint information. If the OBS bucket you selected is unauthorized, click <b>Authorize</b> .
Alarm on Job	Whether to notify users of any job exceptions, such as running exceptions or arrears, via SMS or email.
Exception	If this option is selected, you need to set the following parameters:
	SMN Topic
	Select a custom SMN topic. For how to create a custom SMN topic, see <b>Creating a Topic</b> .

Parameter	Description
Auto Restart	Whether automatic restart is enabled. If enabled, jobs will be automatically restarted and restored when exceptions occur.
upon	If this option is selected, you need to set the following parameters:
Exception	• <b>Max. Retry Attempts</b> : maximum number of retries upon an exception. The unit is times/hour.
	<ul> <li>Unlimited: The number of retries is unlimited.</li> </ul>
	<ul> <li>Limited: The number of retries is user-defined.</li> </ul>
	<ul> <li>Restore Job from Checkpoint: Restore the job from the saved checkpoint.</li> <li>If you select this parameter, you also need to set Checkpoint Path.</li> </ul>
	<b>Checkpoint Path</b> : Select a path for storing checkpoints. This path must match that configured in the application package. Each job must have a unique checkpoint path, or, you will not be able to obtain the checkpoint.
	NOTE
	<ul> <li>If you also select Enable Checkpointing, you do not need to set Checkpoint Path. The system will automatically determine the path based on the Enable Checkpointing configuration.</li> </ul>
	<ul> <li>If you do not select Enable Checkpointing, you need to set Checkpoint Path.</li> </ul>

You can set compute resource specification parameters in the **Runtime Configuration** of Flink jobs, and the parameter values have a higher priority than the specified values.

 Table 11-10 describes the parameter mapping.

#### 

In Flink 1.12, you are advised to set compute resource specification parameters based on the configuration method on the console. Using automatic parameter settings may result in discrepancies in actual CU statistics.

Runtime Configuration	Compute Resource Specificat ion Paramete r of v1	Compute Resource Specificati on Parameter of v2	Description
kubernetes.jobmanag er.cpu	Job Manager CUs	Job Manager CPU	Number of CPU cores available for Job Manager. The default value is <b>1</b> . The minimum value cannot be less than 0.5.
kubernetes.taskmana ger.cpu	CU(s) per TM	Task Manager CPU	Number of CPU cores available for Task Manager. The default value is <b>1</b> . The minimum value cannot be less than 0.5.
jobmanager.memory.p rocess.size	-	Job Manager Memory	Number of memory available for Job Manager. The default value is 4 GB. The minimum size cannot be less than 2 GB (2048 MB). The default unit is GB, which can be set to GB or MB.
taskmanager.memory. process.size	-	Task Manager Memory	Number of memory available for Task Manager. The default value is 4 GB. The minimum size cannot be less than 2 GB (2048 MB). The default unit is GB, which can be set to GB or MB.

**Table 11-10** Mapping between compute resource specification parameters on the console and those in the Runtime Configuration

- **Step 8** Click **Save** on the upper right of the page.
- Step 9 Click Start in the upper right corner. On the displayed Start Flink Job page, confirm the job specifications and the price, and click Start Now to start the job. After the job is started, the system automatically switches to the Flink Jobs page, and the created job is displayed in the job list. You can view the job status in the Status column.
  - Once a job is successfully submitted, its status changes from **Submitting** to **Running**. After the execution is complete, the status changes to **Completed**.
  - If the job status is **Submission failed** or **Running exception**, the job fails to submit or run. In this case, you can hover over the status icon in the **Status**

column of the job list to view the error details. You can click  $\Box$  to copy these details. Rectify the fault based on the error information and resubmit the job.

----End

# **11.4 Configuring Flink Job Permissions**

#### Scenario

- You can isolate Flink jobs allocated to different users by setting permissions to ensure data query performance.
- The administrator and job creator have all permissions, which cannot be set or modified by other users.
- When setting job permissions for a new user, ensure that the region of the user group to which the user belongs has the **Tenant Guest** permission. For details about the Tenant Guest permission and how to apply for the permission, see **Permissions Policies** and **Creating a User Group and Assigning Permissions** in the *Identity and Access Management User Guide*.

#### **Flink Job Permission Operations**

- On the left of the DLI management console, choose Job Management > Flink Jobs.
- Select the job to be configured and choose More > Permissions in the Operation column. The User Permissions area displays the list of users who have permissions on the job.

Flink Jobs ①					C Peerbeck E Export Job	23 Import Job   19 Quick Links (	Video Tatorial Greate Job	Aanage Edge Authentication Code
Start Stop Delote			NI Types 🔹	All statures +	V, O Add 10		×	Q Search by Tag. V C
_ 10 J≣ Name Type 7/	Status 🏆	Description	Username	Crea	ated	Started	Duration	Operation
199273 Eventstate Firek SQL	•			Nev	24, 2021 11:52:14 GMT+08:00			Edit Start More -
- 199274 t Flok 5QL				Nev	24, 2021 14:11:11 GMT+08:00			Filekul
D 192176 c Flink SQL	•			Aug	13, 2021 09:56:38 GMT+08:00			Delete
D 197000 c Flink SQL	•			Oct	27, 2021 11:29:30 GMT+08:00	Oct 27, 2021 14:50:59 GMT+08:00		Modify Name and Description
- 197020 t Filmk SQL	•			Oct	27, 2021 14:38:18 GMT+08:00			Import Savepoint
- 196056 E Film Edge S				Oct	12, 2021 2138:02 GMT+08:00			Tripper Savepoint Manage Permission
- 192748 c Flok lar	•			Sep	02, 2021 2041:16 GMT+08:00			Furthers
D 194424 t Flink OpenS	HC •			Sep	13, 2021 16:21:10 GMT+00:00			LOR STAT MONEY
193235 t Fink Opera	H			Aug	28. 2021 19:20:17 GMT+08:00			Edit   Stort   More +
193234 t Flok Jar				per la constante de	28, 2021 19:05:54 GMT+08:00			Edit   Start   More +
10 • Total Records: 11 < 1 2 > Go 1 •								

You can assign queue permissions to new users, modify permissions for users who have some permissions of a queue, and revoke all permissions of a user on a queue.

- Assign permissions to a new user.

A new user does not have permissions on the job.

- i. Click Grant Permission on the right of User Permissions page. The Grant Permission dialog box is displayed.
- ii. Specify **Username** and select corresponding permissions.
- iii. Click **OK**.

 Table 11-11 describes the related parameters.

#### Figure 11-8 Granting permissions

Grant Pern	nission		
* Username	Enter a username.		
Select the perm	issions to be granted to the use	r	
Select all			
Get Job Det	tails	Modify Jobs	Delete Jobs
Start job		Stop job	Export Job
Grant Perm	ission	Revoke Permission	View Other User's Permissions
		Ok Cancel	]

Table 11-11	Permission	parameters
-------------	------------	------------

Parameter	Description		
Username	Name of the user you want to grant permissions to. <b>NOTE</b> This username must be an existing IAM username. In addition, the user can perform authorization operations only after logging in to the Huawei Cloud platform.		
Permissions	Select all: All permissions are selected.		
to be granted to	• <b>View Job Details</b> : This permission allows you to view the job details.		
the user	• <b>Modify Job</b> : This permission allows you to modify the job.		
	• <b>Delete Job</b> : This permission allows you to delete the job.		
	• <b>Start Job</b> : This permission allows you to start the job.		
	• <b>Stop Job</b> : This permission allows you to stop the job.		
	• <b>Export Job</b> : This permission allows you to export the job.		
	• <b>Grant Permission</b> : This permission allows you to grant job permissions to other users.		
	• <b>Revoke Permission</b> : This permission allows you to revoke the job permissions that other users have but cannot revoke the job creator's permissions.		
	• View Other User's Permissions: This permission allows you to view the job permissions of other users.		

- To assign or revoke permissions of a user who has some permissions on the job, perform the following steps:
  - i. In the list under **User Permissions** for a job, select the user whose permissions need to be modified and click **Set Permission** in the **Operation** column.

ii. In the displayed **Set Permission** dialog box, modify the permissions of the current user. **Table 11-11** lists the detailed permission descriptions.

If all options under **Set Permission** are gray, you are not allowed to change permissions on this job. You can apply to the administrator, job creator, or other authorized users for job permission granting and revoking.

- iii. Click **OK**.
- To revoke all permissions of a user on a job, perform the following steps:

In the list under **User Permissions** for a job, locate the user whose permissions need to be revoked, click **Revoke Permission** in the **Operation** column, and click **Yes**. After this operation, the user does not have any permission on the job.

#### Flink Job Permissions

- View Job Details
  - Tenants and the admin user can view and operate all jobs.
  - Subusers and users with the read-only permission can only view their own jobs.

#### D NOTE

If another user grants any permission other than the job viewing permission to a subuser, the job is displayed in the job list, but the details cannot be viewed by the subuser.

• Start Job

You must have the permission to submit and start jobs.

• Stop Job

You must have the permission to stop queues and jobs.

- Delete Job
  - If a job can be deleted, you can delete the job if you were granted this permission.
  - If a job cannot be deleted, the system stops the job before you delete it.
     For details about how to stop a job, see Stop Job. In addition, you must have the permission to delete the job.
- Create Job
  - By default, sub-users cannot create jobs.
  - To create a job, you must have this permission. Currently, only the admin user has the permission to create jobs. In addition, the user must have the permission of the related package group or package used by the job.

#### • Modify Job

When modifying a job, you need to have the permission to update the job and the permission to the package group or package used by the job belongs.

# 11.5 Managing Flink Jobs

# 11.5.1 Viewing Flink Job Details

After creating a Flink job, you can check the basic information, job details, task list, and execution plan of the job on the DLI console.

This section describes how to check information about a Flink job.

Туре	Description	Helpful Link
Basic information	Includes the job ID, job type, job execution status, and more.	Viewing Basic Information
Job details	Includes SQL statements and the parameter settings for Flink Jar jobs.	Viewing Details
Job monitoring	You can use Cloud Eye to check job data input and output details.	Viewing Monitoring Information
Task list	You can view details about each task running on a job, including the task start time, number of received and transmitted bytes, and running duration.	Viewing the Task List
Execution plan	You can understand the operator flow direction of a running job.	Viewing the Execution Plan

Table 11-12 Viewing Flink job information

# Viewing Basic Information

In the navigation pane of the DLI console, choose **Job Management** > **Flink Jobs**. The **Flink Jobs** page displays all Flink jobs. You can check basic information about any Flink jobs in the list.

Table 11-13 Basic in	formation about a Flink job	)
----------------------	-----------------------------	---

Parameter	Description
ID	ID of a submitted Flink job, which is generated by the system by default.
Name	Name of the submitted Flink job.
Туре	<ul> <li>Type of the submitted Flink job, which includes:</li> <li>Flink SQL</li> <li>Flink Jar</li> <li>Flink OpenSource SQL</li> </ul>

Description
Job status, which is subject to the console.
Description of the submitted Flink job.
Name of the user who submits the job.
Time when the job was created.
Time when the Flink job started to run.
Time consumed by job running.
<ul> <li>Edit: Edit a created job.</li> <li>Start: Start and run a job.</li> <li>More         <ul> <li>FlinkUI: Selecting this will display the Flink job execution page.</li> <li>NOTE</li></ul></li></ul>

# Viewing Details

This section describes how to view job details. After you create and save a job, you can click the job name to view job details, including SQL statements and parameter settings. For a Jar job, you can only view its parameter settings.

- Step 1In the left navigation pane of the DLI management console, choose JobManagement > Flink Jobs. The Flink Jobs page is displayed.
- **Step 2** Click the name of your desired job.

On the displayed page, you can view the job details, SQL statements, job configuration information, task list, execution plan, commit logs, run logs, log list, and job tags.

----End

#### **Viewing Monitoring Information**

You can use Cloud Eye to view details about job data input and output.

- Step 1In the left navigation pane of the DLI management console, choose JobManagement > Flink Jobs. The Flink Jobs page is displayed.
- **Step 2** Click the name of the job you want. The job details are displayed.

Click **Job Monitoring** in the upper right corner of the page to switch to the Cloud Eye console.

#### Figure 11-9 Monitoring a Job



The following table describes monitoring metrics related to Flink jobs.

Name	Description
Flink Job Data Read Rate	Displays the data input rate of a Flink job for monitoring and debugging. Unit: record/s.
Flink Job Data Write Rate	Displays the data output rate of a Flink job for monitoring and debugging. Unit: record/s.
Flink Job Total Data Read	Displays the total number of data inputs of a Flink job for monitoring and debugging. Unit: records
Flink Job Total Data Write	Displays the total number of output data records of a Flink job for monitoring and debugging. Unit: records
Flink Job Byte Read Rate	Displays the number of input bytes per second of a Flink job. Unit: byte/s
Flink Job Byte Write Rate	Displays the number of output bytes per second of a Flink job. Unit: byte/s

Name	Description
Flink Job Total Read Byte	Displays the total number of input bytes of a Flink job. Unit: byte
Flink Job Total Write Byte	Displays the total number of output bytes of a Flink job. Unit: byte
Flink Job CPU Usage	Displays the CPU usage of Flink jobs. Unit: %
Flink Job Memory Usage	Displays the memory usage of Flink jobs. Unit: %
Flink Job Max Operator Latency	Displays the maximum operator delay of a Flink job. The unit is <b>ms</b> .
Flink Job Maximum Operator Backpressure	Displays the maximum operator backpressure value of a Flink job. A larger value indicates severer backpressure. <b>0</b> : OK <b>50</b> : low
	<b>100</b> : high

----End

# Viewing the Task List

You can view details about each task running on a job, including the task start time, number of received and transmitted bytes, and running duration.

#### D NOTE

If the value is **0**, no data is received from the data source.

- Step 1 In the left navigation pane of the DLI management console, choose Job Management > Flink Jobs. The Flink Jobs page is displayed.
- **Step 2** Click the name of the job you want. The job details are displayed.
- **Step 3** On **Task List** and view the node information about the task.

#### Figure 11-10 Task list

test			C Job Monitoring	Edit	Start More 👻
Job Detail Task List Execution Plan	Commit Logs Run Log	Tags			
					С
Name Duration	Max C Task	Status Back Delay	Sent R Sent B	Receiv Receiv	Started Ended
Sou EKA 34d 20	1 000	-	33 9.469	0 0B	Jan 18
Map 34d 20	1 000	48	33 9.469	33 9.479	Jan 18
Sink 34d 20		131	0 0 B	33 9.48 MB	Jan 18

View the operator task list. The following table describes the task parameters.

Parameter	Description	
Name	Name of an operator.	
Duration	Running duration of an operator.	
Max Concurrent Jobs	Number of parallel tasks in an operator.	
Task	<ul> <li>Operator tasks are categorized as follows:</li> <li>The digit in red indicates the number of failed tasks.</li> <li>The digit in light gray indicates the number of canceled tasks.</li> <li>The digit in yellow indicates the number of tasks that are being canceled.</li> <li>The digit in green indicates the number of finished tasks.</li> <li>The digit in blue indicates the number of running tasks.</li> <li>The digit in sky blue indicates the number of tasks that are being deployed.</li> <li>The digit in dark gray indicates the number of tasks in a queue.</li> </ul>	
Status	Status of an operator task.	
Back Pressure Status	<ul> <li>Working load status of an operator. Available options are as follows:</li> <li>OK: indicates that the operator is in normal working load.</li> <li>LOW: indicates that the operator is in slightly high working load. DLI processes data quickly.</li> <li>HIGH: indicates that the operator is in high working load. The data input speed at the source end is slow.</li> </ul>	
Delay	Duration from the time when source data starts being processed to the time when data reaches the current operator. The unit is millisecond.	
Sent Records	Number of data records sent by an operator.	
Sent Bytes	Number of bytes sent by an operator.	
Received Bytes	Number of bytes received by an operator.	
Received Records	Number of data records received by an operator.	
Started	Time when an operator starts running.	
Ended	Time when an operator stops running.	

#### Table 11-15 Parameter descriptions

----End

# Viewing the Execution Plan

You can view the execution plan to understand the operator stream information about the running job.

- Step 1 In the left navigation pane of the DLI management console, choose Job Management > Flink Jobs. The Flink Jobs page is displayed.
- **Step 2** Click the name of the job you want. The job details are displayed.
- **Step 3** Click the **Execution Plan** tab to view the operator flow direction.

#### Figure 11-11 Execution plan



Click a node. The corresponding information is displayed on the right of the page.

- Scroll the mouse wheel to zoom in or out.
- The stream diagram displays the operator stream information about the running job in real time.

----End

# 11.5.2 Setting the Priority for a Flink Job

#### Scenario

In actual job running, it is necessary to prioritize and ensure the normal running of important and urgent tasks due to their varying levels of importance and urgency. This requires providing the necessary compute resources for their normal operations.

DLI offers a feature to set job priorities for each Flink job, which prioritizes the allocation of compute resources to higher priority jobs when resources are limited.

#### **NOTE**

You can set the priority for Flink 1.12 or later jobs.

#### Notes

- Priorities cannot be set for jobs running in queues within an elastic resource pool of the basic edition.
- You can assign a priority level of 1 to 10 for each job, with a larger value indicating a higher priority. Compute resources are preferentially allocated to high-priority jobs. That is, if compute resources required for high-priority jobs are insufficient, compute resources for low-priority jobs are reduced.
- Flink jobs running on a general-purpose queue have a default priority level of 5.

- The job priority change will only be in effect once the job has been stopped, edited, and resubmitted.
- You can set the priority for Flink jobs only after enabling dynamic scaling by setting flink.dli.job.scale.enable to true. For details, see Enabling Dynamic Scaling for Flink Jobs.
- To change the priority for a job, you must first stop the job, change the priority level, and then submit the job for the modification to take effect.

# Setting the Priority for a Flink OpenSource SQL Job

- 1. Log in to the DLI management console.
- 2. In the navigation pane on the left, choose **Job Management > Flink Jobs**.
- 3. Select the job for which you want to set the priority and click **Edit** in the **Operation** column.
- 4. On the far right of the displayed page, click **Runtime Configuration**.
- 5. Enter statements in the text box to enable dynamic scaling and then set the job priority.

**NOTE** 

To set the priority for Flink jobs, you must first enable dynamic scaling by setting **flink.dli.job.scale.enable** to **true**.

For more parameter settings, see Enabling Dynamic Scaling for Flink Jobs.

```
flink.dli.job.scale.enable=true
flink.dli.job.priority=x
```

Figure 11-12 Example configuration for a Flink OpenSource SQL job

Start Save	Save As	Static Stream Graph	
flink.dli.job.scale.enable=true flink.dli.job.priority=8		ß	Running Parameters
			Runtime Configuration

# Setting the Priority for a Flink Jar Job

Enter the following statement in the **Runtime Configuration** text box, where *x* indicates the priority value:

#### flink.dli.job.priority=*x*

- 1. Log in to the DLI management console.
- 2. In the navigation pane on the left, choose **Job Management** > **Flink Jobs**.
- 3. Select the job for which you want to set the priority and click **Edit** in the **Operation** column.

4. In the **Runtime Configuration** text box, enter the following statements to enable dynamic scaling and set the job priority:

#### **NOTE**

To set the priority for Flink jobs, you must first enable dynamic scaling by setting **flink.dli.job.scale.enable** to **true**.

For more parameter settings, see **Enabling Dynamic Scaling for Flink Jobs**. flink.dli.job.scale.enable=true

flink.dli.job.priority=x

Figure 11-13 Example configuration for a Flink Jar job

Runtime Configuration	flink.dli.job.scale.enable=true flink.dli.job.priority=8	
		4

# **11.5.3 Enabling Dynamic Scaling for Flink Jobs**

#### Scenario

In actual job operations, the compute resources required by a job vary depending on the data volume. As a result, compute resources are wasted when the volume is small and are insufficient when the volume is large.

DLI provides dynamic scaling to dynamically adjust the compute resources used by a job based on the job load, such as the data input and output volume, data input and output rate, and backpressure, to improve resource utilization.

After enabling dynamic scaling for Flink jobs, the system will adjust resource allocation based on the actual resource requirements of the Flink jobs. When the remaining pod resources in the elastic resource pool are sufficient to meet the minimum resource requirements of the job, the system will automatically reduce the number of nodes where the job is running. This ensures efficient job running while improving resource utilization.

#### **NOTE**

Currently, dynamic scaling can only be enabled for Flink 1.12 jobs.

#### Notes

- During dynamic scaling of a Flink job, if queue resources are preempted and the remaining resources are insufficient for starting the job, the job may fail to be restored.
- When the resources that can be used by a Flink job are dynamically scaled in or out, the background job needs to be stopped and then restored from the savepoint. So, the job cannot process data before the restoration is successful.
- Savepoints need to be triggered during scaling. So, you must configure an OBS bucket, save logs, and enable checkpointing.
- Do not set the scaling detection period to a small value to avoid frequent job start and stop.
- The restoration duration of a scaling job is affected by the savepoint size. If the savepoint size is large, the restoration may take a long time.

• To adjust the configuration items of dynamic scaling, you need to stop the job, edit the job, and submit the job for the modification to take effect.

#### Procedure

Dynamic scaling applies to Flink OpenSource SQL and Flink Jar jobs.

- 1. Log in to the DLI management console.
- 2. In the navigation pane on the left, choose **Job Management** > **Flink Jobs**.
- 3. Select the job for which you want to enable dynamic scaling and click **Edit** in the **Operation** column.
  - For a Flink OpenSource SQL job, click **Runtime Configuration** on the right to configure dynamic scaling parameters.
  - For a Flink Jar job, click the **Runtime Configuration** box to configure dynamic scaling parameters.

Parameter	Default Value	Description
flink.dli.job.scale.ena ble	false	Whether to enable dynamic scaling, that is, whether to allow DLI to adjust the resources used by jobs based on job loads and job priorities.
		If this parameter is set to <b>false</b> , the function is disabled.
		If this parameter is set to <b>true</b> , the function is enabled.
		The default value is <b>false</b> .
flink.dli.job.scale.inte rval	30	Interval for checking whether to scale the resources for the current job, in minutes. The default value is <b>30</b> . For example, <b>30</b> indicates that the job is checked every 30 minutes to determine whether to scale in or out the resources used by the job. Note: This configuration is effective only when dynamic scaling is enabled.
flink.dli.job.cu.max	Initial CU value	Maximum number of CUs that can be used by the current job during dynamic scaling. If this parameter is not set, the default value is the initial total number of CUs of the job.
		Note: The value of this parameter cannot be smaller than the total number of CUs configured by the user. In addition, this parameter is effective only when dynamic scaling is enabled.

Table 11-16 Dynamic scaling parameters

Parameter	Default Value	Description
flink.dli.job.cu.min	2	Minimum number of CUs that can be used by the current job during dynamic scaling. The default value is <b>2</b> .
		Note: The value of this parameter cannot be greater than the total number of CUs configured by the user. In addition, this parameter is effective only when dynamic scaling is enabled.

# **11.5.4 Querying Logs for Flink Jobs**

## Scenario

DLI job buckets are used to store temporary data generated during DLI job running, such as job logs and results.

This section describes how to configure a bucket for DLI jobs on the DLI console and obtain Flink job logs.

#### Notes

- To avoid disordered job results, do not use the OBS bucket configured for DLI jobs for any other purposes.
- DLI jobs must be set and modified by the main account as IAM users do not have required permissions.
- You cannot view the logs for DLI jobs before configuring a bucket.
- You can configure lifecycle rules to periodically delete objects from buckets or change storage classes of objects.
- Exercise caution when modifying the job bucket, as it may result in the inability to retrieve historical data.

# Prerequisites

Before the configuration, create an OBS bucket or parallel file system (PFS). In big data scenarios, you are advised to create a PFS. PFS is a high-performance file system provided by OBS, with access latency in milliseconds. PFS can achieve a bandwidth performance of up to TB/s and millions of IOPS, which makes it ideal for processing high-performance computing (HPC) workloads.

For details about PFS, see "Parallel File System Feature Guide" in the *Object Storage Service User Guide*.

# Configuring a Bucket for DLI Jobs

In the navigation pane of the DLI console, choose Global Configuration > Project.

2. On the **Project** page, click and next to **Job Bucket** to configure bucket information.

#### Figure 11-14 Project

Data Lake Insig	ht a	Project
Overview		
SQL Editor		Job Bucket
Job Management	$\sim$	Bucket Name: rain
Resources	~	This bucket is used to store temporary data generated by DLI, such as job logs and job results. Do not use this bucket for other purposes. If you do not create this bucket, you will not be able to view job logs. You can use the main
Data Management	$\sim$	account to set and modify the bucket. Sub-users do not have modification permissions. You can set a lifecycle rule to periodically delete objects in a
Job Templates	$\sim$	bucket or change its storage class. Exercise caution when you modify the lifecycle rule to prevent historical data being deleted by mistake.
Datasource Connectio	ns	
Global Configuration	^	
Global Variables		
SQL Inspector	<	
Project		
Service Authorization		

- 3. Click  $\square$  to view available buckets.
- 4. In the displayed **OBS** dialog box, click the name of a bucket or search for and click a bucket name and then click **OK**. In the **Set Job Bucket** dialog box, click **OK**.

Temporary data generated during DLI job running will be stored in the OBS bucket.

#### Figure 11-15 Setting the job bucket

Set Job Bu	cket		×
Parallel fil	e buckets are recommended.	×	
★ Job Bucket	rain		
		Cancel	

## **Viewing Commit Logs**

You can check commit logs to locate commit faults.

- **Step 1** In the navigation pane of the DLI console, choose **Job Management** > **Flink Jobs**.
- **Step 2** Click the name of the Flink job whose commit logs you want to check.
- Step 3 Click the Commit Logs tab and check the job commit process.

#### Figure 11-16 Commit logs





# **Viewing Run Logs**

You can check run logs to locate job running faults.

- **Step 1** In the navigation pane of the DLI console, choose **Job Management > Flink Jobs**.
- **Step 2** Click the name of the Flink job whose commit logs you want to check.
- **Step 3** Click the **Run Log** tab and check the JobManager and TaskManager information of the running job.

#### Figure 11-17 Run logs

ID: 234396 Job Type: Flink OpenSource SQL		
	Run Log	

JobManager and TaskManager information is updated every minute. By default, the run logs generated in the last minute are displayed.

If you have configured an OBS bucket to store job logs, you can access it to download and check historical logs.

#### **NOTE**

For details about how to upload files to OBS, see **Uploading an Object** in *Object Storage Service Getting Started*.

If the job is not running, you cannot check Task Manager information.

----End

# Viewing the Log List

You can view historical job files by viewing the Flink job log list.

- **Step 1** In the navigation pane of the DLI console, choose **Job Management** > **Flink Jobs**.
- **Step 2** Click the name of your desired job.

**Step 3** On the **Logs** tab page, choose **Job Manager** and **Task Manager** separately to view their logs.

Figure 11-18 Flink job log list

datagen ID: 81403 Job Type: Flink OpenSource SQL	
	Logs
Name	
jobmanager.log jobmanager.log 2023-06-21_10.0.log	
gc-2023-06-21_10-36-16.log.0.current	
jobmanager.out	
jobmanager.err	

----End

# 11.5.5 Common Operations of Flink Jobs

After creating a job, you can manage it by performing various operations such as editing its basic information, starting or stopping it, and importing or exporting it.

## **Editing a Job**

You can edit a created job, for example, by modifying the SQL statement, job name, job description, or job configurations.

- Step 1 In the left navigation pane of the DLI management console, choose Job Management > Flink Jobs. The Flink Jobs page is displayed.
- **Step 2** In the row where the job you want to edit locates, click **Edit** in the **Operation** column to switch to the editing page.
- **Step 3** Edit the job as required.

For details, see Creating a Flink OpenSource SQL Job and Creating a Flink Jar Job.

----End

#### Starting a Job

You can start a saved or stopped job.

- Step 1 In the left navigation pane of the DLI management console, choose Job Management > Flink Jobs. The Flink Jobs page is displayed.
- **Step 2** Use either of the following methods to start jobs:
  - Starting a single job
     Select a job and click Start in the Operation column.
     Alternatively, you can select the row where the job you want to start locates and click Start in the upper left of the job list.

• Batch starting jobs

Select the rows where the jobs you want to start locate and click **Start** in the upper left of the job list.

After you click **Start**, the **Start Flink Jobs** page is displayed.

**Step 3** On the **Start Flink Jobs** page, confirm the job information and price. If they are correct, click **Start Now**.

After a job is started, you can view the job execution result in the **Status** column.

----End

## Stopping a Job

You can stop a job in the **Running** or **Submitting** state.

- Step 1 In the left navigation pane of the DLI management console, choose Job Management > Flink Jobs. The Flink Jobs page is displayed.
- **Step 2** Stop a job using either of the following methods:
  - Stopping a job

Locate the row that contains the job to be stopped, click **More** in the **Operation** column, and select **Stop**.

Alternatively, you can select the row where the job you want to stop locates and click **Stop** in the upper left of the job list.

• Batch stopping jobs

Locate the rows containing the jobs you want to stop and click **Stop** in the upper left of the job list.

Step 3 In the displayed Stop Job dialog box, click OK to stop the job.

Figure 11-19 Stopping a job

Are you sure you want to stop the following 1 jobs?

Name	Status	Description
	Running	
Trigger Savepoint	1	OK Cancel

#### D NOTE

- Before stopping a job, you can trigger a savepoint to save the job status information. When you start the job again, you can choose whether to restore the job from the savepoint.
- If you select **Trigger savepoint**, a savepoint is created. If **Trigger savepoint** is not selected, no savepoint is created. By default, the savepoint function is disabled.
- The lifecycle of a savepoint starts when the savepoint is triggered and stops the job, and ends when the job is restarted. The savepoint is automatically deleted after the job is restarted.

When a job is being stopped, the job status is displayed in the **Status** column of the job list. The details are as follows:

- **Stopping**: indicates that the job is being stopped.
- **Stopped**: indicates that the job is stopped successfully.
- **Stop failed**: indicates that the job failed to be stopped.

----End

## Deleting a Job

If you do not need to use a job, perform the following operations to delete it. A deleted job cannot be restored. Therefore, exercise caution when deleting a job.

- Step 1 In the left navigation pane of the DLI management console, choose Job Management > Flink Jobs. The Flink Jobs page is displayed.
- **Step 2** Perform either of the following methods to delete jobs:
  - Deleting a single job

Locate the row containing the job you want to delete and click **More** > **Delete** in the **Operation** column.

Alternatively, you can select the row containing the job you want to delete and click **Delete** in the upper left of the job list.

• Deleting jobs in batches

Select the rows containing the jobs you want to delete and click **Delete** in the upper left of the job list.

Step 3 Click Yes.

----End

#### Exporting a Job

You can export the created Flink jobs to an OBS bucket.

This mode is applicable to the scenario where a large number of jobs need to be created when you switch to another region, project, or user. In this case, you do not need to create a job. You only need to export the original job, log in to the system in a new region or project, or use a new user to import the job.

#### **NOTE**

When switching to another project or user, you need to grant permissions to the new project or user. For details, see **Configuring Flink Job Permissions**.

- **Step 1** In the left navigation pane of the DLI management console, choose **Job Management > Flink Jobs**. The **Flink Jobs** page is displayed.
- **Step 2** Click **Export Job** in the upper right corner. The **Export Job** dialog box is displayed.

Figure 11-20 Exporting a job



- Step 3 Select the OBS bucket where the job is stored. Click Next.
- Step 4 Select job information you want to export.

By default, configurations of all jobs are exported. You can enable the **Custom Export** function to export configurations of the desired jobs.

**Step 5** Click **Confirm** to export the job.

----End

#### Importing a Job

You can import the Flink job configuration file stored in the OBS bucket to the **Flink Jobs** page of DLI.

This mode is applicable to the scenario where a large number of jobs need to be created when you switch to another region, project, or user. In this case, you do not need to create a job. You only need to export the original job, log in to the system in a new region or project, or use a new user to import the job.

To import a self-created job, use the job creation function.

For details, see **Creating a Flink OpenSource SQL Job** and **Creating a Flink Jar Job**.

#### **NOTE**

- When switching to another project or user, you need to grant permissions to the new project or user. For details, see **Configuring Flink Job Permissions**.
- Only jobs whose data format is the same as that of Flink jobs exported from DLI can be imported.
- Step 1 In the left navigation pane of the DLI management console, choose Job Management > Flink Jobs. The Flink Jobs page is displayed.
- **Step 2** Click **Import Job** in the upper right corner. The **Import Job** dialog box is displayed.
- **Step 3** Select the complete OBS path of the job configuration file to be imported. Click **Next**.
- Step 4 Configure the same-name job policy and click next. Click Next.
  - Select **Overwrite job of the same name**. If the name of the job to be imported already exists, the existing job configuration will be overwritten and the job status switches to **Draft**.

- If **Overwrite job of the same name** is not selected and the name of the job to be imported already exists, the job will not be imported.
- **Step 5** Ensure that **Config File** and **Overwrite Same-Name Job** are correctly configured. Click **Confirm** to import the job.

----End

## Modifying the Name and Description of a Flink Job

You can change the job name and description as required.

- Step 1 In the left navigation pane of the DLI management console, choose Job Management > Flink Jobs. The Flink Jobs page is displayed.
- Step 2 In the Operation column of the job whose name and description need to be modified, choose More > Modify Name and Description. The Modify Name and Description dialog box is displayed. Change the name or modify the description of a job.
- Step 3 Click OK.

----End

## **Triggering a Savepoint**

Before stopping a job, you can trigger a savepoint to save the status information of your job. When you restart the job, you can choose whether to quickly recover it from the most recent savepoint.

- **Step 1** In the navigation pane of the DLI console, choose **Job Management > Flink Jobs**.
- **Step 2** Locate the job you want to stop, click **More** in the **Operation** column, and select **Trigger Savepoint**. In the displayed dialog box, select a save path.
- Step 3 Click OK.

----End

**NOTE** 

- You can click **Trigger Savepoint** for jobs in the **Running** status to save the job status.
- The lifecycle of a savepoint starts when the savepoint is triggered and stops the job, and ends when the job is restarted. The savepoint is automatically deleted after the job is restarted.

#### Importing to a Savepoint

Flink jobs can be restored based on imported savepoints.

- **Step 1** In the navigation pane of the DLI console, choose **Job Management > Flink Jobs**.
- **Step 2** Locate the job you want to stop, click **More** in the **Operation** column, and select **Import Savepoint**. In the displayed dialog box, select a save path.
- Step 3 Click OK.

----End

Х

# **Runtime Configuration**

You can configure job exception alarms and restart options by selecting **Runtime Configuration**.

#### **NOTE**

This configuration is only available for Flink OpenSource SQL jobs and Flink Jar jobs.

- 1. Locate the desired Flink job, click **More** in the **Operation** column, and select **Runtime Configuration**.
- 2. In the **Runtime Configuration** dialog box, set the following parameters:

Figure 11-21 Runtime configuration

## **Runtime Configuration**

Name	testflinkjar		
Alarm Generation upo			
* SMN Topic	DLI_fink_info Configure Topic	•	
Auto Restart upon Exc	. 🔽		
Max. Retry Attempts	Unlimited	Limited	
Restore Job from Che			
★ Checkpoint Path	dli-test-obs01/		ß
	The checkpoint path in the application pac path for each job mu checkpoint cannot be	must be the same kage. Note that th st be unique. Othe e obtained.	as that you set e checkpoint rwise, the
	ОК	Cancel	

#### Table 11-17 Running parameters

Parameter	Description
Name	Job name.

Parameter	Description
Alarm Generation upon Job	Whether to report job exceptions, for example, abnormal job running or exceptions due to an insufficient balance, to users via SMS or email.
Exception	If this option is selected, you need to set the following parameters:
	SMN Topic
	Select a custom SMN topic. For how to create a custom SMN topic, see <b>Creating a Topic</b> .
Auto Restart upon Exception	Whether to enable automatic restart. If this function is enabled, any job that has become abnormal will be automatically restarted.
	If this option is selected, you need to set the following parameters:
	• Max. Retry Attempts: maximum number of retries upon an exception. The unit is times/hour.
	<ul> <li>Unlimited: The number of retries is unlimited.</li> </ul>
	<ul> <li>Limited: The number of retries is user-defined.</li> </ul>
	<ul> <li>Restore Job from Checkpoint: Restore the job from the saved checkpoint.</li> </ul>
	This parameter cannot be configured for Flink SQL jobs or Flink OpenSource SQL jobs.
	If this parameter is selected, you need to set <b>Checkpoint Path</b> for Flink Jar jobs.
	<b>Checkpoint Path</b> : Select the checkpoint saving path. The checkpoint path must be the same as that you set in the application package. Note that the checkpoint path for each job must be unique. Otherwise, the checkpoint cannot be obtained.

# **11.6 Managing Flink Job Templates**

Flink templates include sample templates and custom templates. You can modify an existing sample template to meet the actual job logic requirements and save time for editing SQL statements. You can also customize a job template based on your habits and methods so that you can directly invoke or modify the template in later jobs.

Flink template management provides the following functions:

- Flink SQL Sample Template
- Flink OpenSource SQL Sample Template
- Custom Templates
- Creating a Template

- Creating a Job Based on a Template
- Modifying a Template
- Deleting a Template

## Flink SQL Sample Template

The template list displays existing sample templates for Flink SQL jobs. **Table 1** describes the parameters in the template list.

The scenarios of sample templates can be different, which are subject to the console.

Para meter	Description
Name	Template name. The value can contain up to 64 characters. Only letters, digits, hyphens (-), and underscores (_) are allowed.
Descri ption	Description of a template. It can contain 0 to 512 characters.
Opera tion	<b>Create Job</b> : Create a job directly by using the template. After a job is created, the system switches to the <b>Edit</b> page under <b>Job Management</b> .

Table 11-18 Parameters in the Flink SQL sample template list

# Flink OpenSource SQL Sample Template

The template list displays existing sample templates for Flink SQL OpenSource jobs. **Table 1** describes the parameters in the template list.

Para meter	Description
Name	Template name. The value can contain up to 64 characters. Only letters, digits, hyphens (-), and underscores (_) are allowed.
Descri ption	Description of a template. It can contain 0 to 512 characters.
Opera tion	<b>Create Job</b> : Create a job directly by using the template. After a job is created, the system switches to the <b>Edit</b> page under <b>Job Management</b> .

Table 11-19 Parameters in the sample template list for Flink OpenSource SQL jobs

The existing sample templates apply to the following scenarios:

- Create a wide table of order information from the dimension table of address information.
- Generate statistics on indicators such as the daily transaction amount, number of orders, and number of paid users in real time.

• Generate statistics on offerings with the highest real-time click-through rate.

# **Custom Templates**

The custom template list displays all Jar job templates. **Table 1** describes parameters in the custom template list.

Parameter	Description
Name	Template name. The value can contain up to 64 characters. Only letters, digits, hyphens (-), and underscores (_) are allowed.
Description	Description of a template. It can contain 0 to 512 characters.
Created	Time when a template is created.
Updated	Latest time when a template is modified.
Operation	<ul> <li>Edit: Modify a template that has been created.</li> <li>Create Job: Create a job directly by using the template. After a job is created, the system switches to the Edit page under Job Management.</li> <li>More: <ul> <li>Delete: Delete a created template.</li> <li>Tags: View or add tags.</li> </ul> </li> </ul>

Table 11-20 Parar	neters in the	e custom	template	list
-------------------	---------------	----------	----------	------

# **Creating a Template**

You can create a template using any of the following methods:

- Creating a template on the **Template Management** page
  - a. In the left navigation pane of the DLI management console, choose **Job Templates** > **Flink Templates**.
  - b. Click **Create Template** in the upper right corner of the page. The **Create Template** dialog box is displayed.
  - c. Specify **Name** and **Description**.

Create Templa	ate	>
Туре	Flink SQL 🗸	
★ Name	Enter a name.	
Description	Description	
Tags	It is recommended that you use TMS's predefined tag function to add the same tag different cloud resources. View predefined tags $C$ To add a tag, enter a tag key and a tag value below.	to
	Enter a tag key Enter a tag value	Add
	10 tags available for addition.           OK         Cancel	

Figure 11-22 Creating a Flink template

 Table 11-21
 Template parameters

Parame ter	Description
Туре	<ul><li>Template type</li><li>Flink SQL job template</li><li>Flink OpenSource SQL job template</li></ul>
Name	Template name. The value can contain up to 64 characters. Only letters, digits, hyphens (-), and underscores (_) are allowed. <b>NOTE</b> The template name must be unique.
Descript ion	Description of a template. It can contain 0 to 512 characters.

Parame ter	Description
Tags	Tags used to identify cloud resources. A tag includes the tag key and tag value. If you want to use the same tag to identify multiple cloud resources, that is, to select the same tag from the drop-down list box for all services, you are advised to create predefined tags on the Tag Management Service (TMS).
	If your organization has configured tag policies for DLI, add tags to resources based on the policies. If a tag does not comply with the tag policies, resource creation may fail. Contact your organization administrator to learn more about tag policies.
	For details, see Tag Management Service User Guide.
	NOTE
	• A maximum of 20 tags can be added.
	<ul> <li>Only one tag value can be added to a tag key.</li> </ul>
	<ul> <li>The key name in each resource must be unique.</li> </ul>
	• Tag key: Enter a tag key name in the text box.
	<b>NOTE</b> A tag key can contain a maximum of 128 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed, but the value cannot start or end with a space or start with _ <b>sys</b>
	• Tag value: Enter a tag value in the text box.
	NOTE A tag value can contain a maximum of 255 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed.

d. Click **OK** to enter the editing page.

The **Table 11-22** describes the parameters on the template editing page.

Table 11-22 Template parameters

Parameter	Description
Name	You can modify the template name.
Description	You can modify the template description.
Saving Mode	• <b>Save Here</b> : Save the modification to the current template.
	• <b>Save as New</b> : Save the modification as a new template.
SQL statement editing area	In the area, you can enter detailed SQL statements to implement business logic. For how to compile SQL statements, see <b>Data Lake Insight SQL Syntax</b> <b>Reference</b> .

Parameter	Description	
Save	Save the modifications.	
Create Job	Use the current template to create a job.	
Format	Format SQL statements. After SQL statements are formatted, you need to compile SQL statements again.	
Theme Settings	Change the font size, word wrap, and page style (black or white background).	

- e. In the SQL statement editing area, enter SQL statements to implement service logic. For how to compile SQL statements, see **Data Lake Insight SQL Syntax Reference**.
- f. After the SQL statement is edited, click **Save** in the upper right corner to complete the template creation.
- g. (Optional) If you do not need to modify the template, click **Create Job** in the upper right corner to create a job based on the current template. For how to create a job, see **Creating a Flink Jar Job**.
- Creating a template based on an existing job template
  - a. In the left navigation pane of the DLI management console, choose **Job Templates** > **Flink Templates**. Click the **Custom Templates** tab.
  - b. In the row where the desired template is located in the custom template list, click **Edit** under **Operation** to enter the **Edit** page.
  - c. After the modification is complete, set **Saving Mode** to **Save as New**.
  - d. Click **Save** in the upper right corner to save the template as a new one.
- Creating a template using a created job
  - a. In the left navigation pane of the DLI management console, choose Job Management > Flink Jobs. The Flink Jobs page is displayed.
  - b. Click **Create Job** in the upper right corner. The **Create Job** page is displayed.
  - c. Specify parameters as required.
  - d. Click **OK** to enter the editing page.
  - e. After the SQL statement is compiled, click **Set as Template**.
  - f. In the **Set as Template** dialog box that is displayed, specify **Name** and **Description** and click **OK**.
- Creating a template based on the existing job
  - a. In the left navigation pane of the DLI management console, choose **Job Management** > **Flink Jobs**. The **Flink Jobs** page is displayed.
  - b. In the job list, locate the row where the job that you want to set as a template resides, and click **Edit** in the **Operation** column.
  - c. After the SQL statement is compiled, click **Set as Template**.
  - d. In the **Set as Template** dialog box that is displayed, specify **Name** and **Description** and click **OK**.

# Creating a Job Based on a Template

You can create jobs based on sample templates or custom templates.

- 1. In the left navigation pane of the DLI management console, choose **Job Templates** > **Flink Templates**.
- 2. In the sample template list, click **Create Job** in the **Operation** column of the target template. For how to create a job, see **Creating a Flink OpenSource SQL Job** and **Creating a Flink Jar Job**.

# Modifying a Template

After creating a custom template, you can modify it as required. The sample template cannot be modified, but you can view the template details.

- In the left navigation pane of the DLI management console, choose Job Templates > Flink Templates. Click the Custom Templates tab.
- 2. In the row where the template you want to modify is located in the custom template list, click **Edit** in the **Operation** column to enter the **Edit** page.
- 3. In the SQL statement editing area, modify the SQL statements as required.
- 4. Set Saving Mode to Save Here.
- 5. Click **Save** in the upper right corner to save the modification.

## Deleting a Template

You can delete a custom template as required. The sample templates cannot be deleted. Deleted templates cannot be restored. Exercise caution when performing this operation.

- In the left navigation pane of the DLI management console, choose Job Templates > Flink Templates. Click the Custom Templates tab.
- 2. In the custom template list, select the templates you want to delete and click **Delete** in the upper left of the custom template list.

Alternatively, you can delete a template by performing the following operations: In the custom template list, locate the row where the template you want to delete resides, and click **More** > **Delete** in the **Operation** column.

3. In the displayed dialog box, click Yes.

# 11.7 Adding Tags to a Flink Job

A tag is a key-value pair customized by users and used to identify cloud resources. It helps users to classify and search for cloud resources. A tag consists of a tag key and a tag value.

DLI allows you to add tags to Flink jobs. You can add tags to Flink jobs to identify information such as the project name, service type, and background. If you use tags in other cloud services, you are advised to create the same tag key-value pairs for cloud resources used by the same business to keep consistency.

DLI supports the following two types of tags:

• Resource tags: indicate non-global tags created on DLI.

Predefined tags: global tags created on Tag Management Service (TMS).
 For more information about predefined tags, see Tag Management Service User Guide.

If your organization has configured tag policies for DLI, add tags to resources based on the policies. If a tag does not comply with the tag policies, resource creation may fail. Contact your organization administrator to learn more about tag policies.

This section includes the following content:

- Managing a Job Tag
- Searching for a Job by Tag

#### Managing a Job Tag

DLI allows you to add, modify, or delete tags for jobs.

- Step 1 In the left navigation pane of the DLI management console, choose Job Management > Flink Jobs. The Flink Jobs page is displayed.
- Step 2 Click the name of the job to be viewed. The Job Details page is displayed.
- **Step 3** Click **Tags** to display the tag information about the current job.

Figure 11-23 Managing a job tag



- Step 4 Click Add/Edit Tag to open to the Add/Edit Tag dialog box.
- **Step 5** Configure the tag parameters in the **Add/Edit Tag** dialog box.

Figure 11-24 Adding a tag

Add/Edit Tag						
It is recommended that you use TMS's predefined tag function to add the same tag to different cloud resources. View predefined tags $C$						
To add a tag, enter a tag key and a tag value below.						
Enter a tag key     Enter a tag value     Add						
10 tags available for addition.						
OK Cancel						

Parame ter	Description
Tag key	You can perform the following operations:
	<ul> <li>Click the text box and select a predefined tag key from the drop- down list.</li> </ul>
	To add a predefined tag, you need to create one on TMS and then select it from the <b>Tag key</b> drop-down list. You can click <b>View</b> <b>predefined tags</b> to go to the <b>Predefined Tags</b> page of the TMS console. Then, click <b>Create Tag</b> in the upper corner of the page to create a predefined tag.
	For details, see <b>Creating Predefined Tags</b> in <i>Tag Management</i> Service User Guide.
	• Enter a tag key in the text box.
	<b>NOTE</b> A tag key can contain a maximum of 128 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed, but the value cannot start or end with a space or start with <b>_sys_</b> .
Tag value	You can perform the following operations:
	• Click the text box and select a predefined tag value from the drop- down list.
	• Enter a tag value in the text box.
	<b>NOTE</b> A tag value can contain a maximum of 255 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed.

#### Table 11-23Tag parameters

#### **NOTE**

- A maximum of 20 tags can be added.
- Only one tag value can be added to a tag key.
- The key name in each resource must be unique.

#### Step 6 Click OK.

**Step 7** (Optional) In the tag list, locate the row where the tag you want to delete resides, click **Delete** in the **Operation** column to delete the tag.

----End

# Searching for a Job by Tag

If tags have been added to a job, you can search for the job by setting tag filtering conditions to quickly find it.

- Step 1 In the left navigation pane of the DLI management console, choose Job Management > Flink Jobs. The Flink Jobs page is displayed.
- **Step 2** In the upper right corner of the page, click the search box and select **Tags**.

Figure 11-25 Searching for jobs by tag



- **Step 3** Choose a tag key and value as prompted. If no tag key or value is available, create a tag for the job. For details, see Managing a Job Tag.
- **Step 4** Choose other tags to generate a tag combination for job search.
- **Step 5** Click search icon. The target job will be displayed in the job list.

----End

# **12** Submitting a Spark Job on the DLI Management Console

# 12.1 Creating a Spark Job

DLI Spark jobs provide fully managed Spark computing services.

On the **Overview** page, click **Create Job** in the upper right corner of the **Spark Jobs** tab or click **Create Job** in the upper right corner of the **Spark Jobs** page. The Spark job editing page is displayed.

On the Spark job editing page, a message is displayed, indicating that a temporary DLI data bucket will be created. The created bucket is used to store temporary data generated by DLI, such as job logs and job results. You cannot view job logs if you choose not to create the bucket. You can **configure a lifecycle rule** to periodically delete objects in a bucket or transit objects between different storage classes. The bucket will be created and the default bucket name is used.

If you do not need to create a DLI temporary data bucket and do not want to receive this message, select **Do not show again** and click **Cancel**.

# Prerequisites

- You have uploaded the dependencies to the corresponding OBS bucket on the **Data Management > Package Management** page.
- Before creating a Spark job to access other external data sources, such as OpenTSDB, HBase, Kafka, GaussDB(DWS), RDS, CSS, CloudTable, DCS Redis, and DDS, you need to create a datasource connection to enable the network between the job running queue and external data sources.
  - For details about the external data sources that can be accessed by Spark jobs, see Common Development Methods for DLI Cross-Source Analysis.
  - For how to create a datasource connection, see Configuring the Network Connection Between DLI and Data Sources (Enhanced Datasource Connection).

On the **Resources** > **Queue Management** page, locate the queue you have created, click **More** in the **Operation** column, and select **Test** 

Address Connectivity to check if the network connection between the queue and the data source is normal. For details, see Testing Address Connectivity.

## Procedure

 In the left navigation pane of the DLI management console, choose Job Management > Spark Jobs. The Spark Jobs page is displayed.

Click **Create Job** in the upper right corner. In the job editing window, you can set parameters in **Fill Form** mode or **Write API** mode.

The following uses the **Fill Form** as an example. In **Write API** mode, refer to the **Data Lake Insight API Reference** for parameter settings.

- 2. Select a queue.
  - a. **Queues**: Select a queue from the drop-down list.
  - b. **Spark Version**: Select a Spark version from the drop-down list. The latest version is recommended.

#### D NOTE

You are advised not to use Spark of different versions for a long time.

- Doing so can lead to code incompatibility, which can negatively impact the job execution efficiency.
- Doing so may result in job execution failures due to conflicts in dependencies. Jobs rely on specific versions of libraries or components.
- 3. Configure the application.

Table 12-1	Application	configuration	parameters
------------	-------------	---------------	------------

Parameter	Description	
Application	Select the package to be executed. The value can be <b>.jar</b> or <b>.py</b> .	
	There are the following ways to manage JAR files:	
	<ul> <li>Upload packages to OBS: Upload JAR files to an OBS bucket in advance and select the corresponding OBS path.</li> </ul>	
	<ul> <li>Upload packages to DLI: Upload JAR files to an OBS bucket in advance and create a package on the Data Management &gt; Package Management page of the DLI management console. For details, see Creating a DLI Package.</li> </ul>	
	For Spark 3.3. <i>x</i> or later, you can only select packages in OBS paths.	
Parameter	Description	
---------------------------	---	--
Agency	Before using Spark 3.3.1 or later (Spark general queue scenario) to run jobs, you need to create an agency on the IAM console and add the new agency information. For details, see <b>Customizing DLI Agency Permissions</b> .	
	Common scenarios for creating an agency: DLI is allowed to read and write data from and to OBS to transfer logs. DLI is allowed to access DEW to obtain data access credentials and access catalogs to obtain metadata.	
	For details, see <b>Configuring DLI Agency Permissions</b> .	
Main Class ( class)	Enter the name of the main class. When the application type is <b>.jar</b> , the main class name cannot be empty.	
Application Parameters	User-defined parameters. Separate multiple parameters by <b>Enter</b> .	
	These parameters can be replaced with global variables. For example, if you create a global variable <b>batch_num</b> on the <b>Global Configuration</b> > <b>Global Variables</b> page, you can use <b>{{batch_num}}</b> to replace a parameter with this variable after the job is submitted.	

#### 4. Configure the job.

#### Table 12-2 Job configuration parameters

Parameter	Description
Job Name (name)	Set a job name.

Parameter	Description
Spark Argument	Enter a parameter in the format of <b>key=value</b> . Press Enter to separate multiple key-value pairs.
s(conf)	These parameters can be replaced with global variables. For example, if you create a global variable <b>custom_class</b> on the <b>Global Configuration</b> > <b>Global Variables</b> page, you can use "spark.sql.catalog"={{custom_class}} to replace a parameter with this variable after the job is submitted.
	NOTE
	<ul> <li>The JVM garbage collection algorithm cannot be customized for Spark jobs.</li> </ul>
	• If the Spark version is <b>3.1.1</b> , configure <b>Spark parameters (conf)</b> to select a dependency module. For details about the example configuration, see <b>Table 12-3</b> .
	If you select <b>3.3.1</b> for <b>Spark Version</b> , you can configure compute resource specification parameters in <b>Spark</b> <b>Argument(conf)</b> . Note that the configuration priority of <b>Spark Argument(conf)</b> is higher than that of <b>Resource</b> <b>Specifications(Optional)</b> in <b>Advanced Settings</b> .
	Table 12-4 describes the parameter mapping.
	<b>NOTE</b> When configuring compute resource specification parameters in <b>Spark</b> <b>Argument(conf)</b> , you can use the unit M, G, or K. If the unit is not specified, the default unit is byte.
Access Metadata	Choose whether to enable access to metadata for a Spark job. Set it to <b>Yes</b> if you need to configure the metadata type accessed by the job. DLI metadata is accessed by default. When set to <b>Yes</b> , you also need to set <b>Metadata Source</b> .
Retry upon	Indicates whether to retry a failed job
Failure	If you select <b>Ves</b> you need to set the following parameters:
	Maximum Datrice: Maximum number of retry times. The
	maximum value is <b>100</b> .

	Table	12-3	Spark	Parameter	(conf)	configuration
--	-------	------	-------	-----------	--------	---------------

Datasource	Example Value
CSS	spark.driver.extraClassPath=/usr/share/extension/dli/ spark-jar/datasource/css/*
	spark.executor.extraClassPath=/usr/share/extension/dli/ spark-jar/datasource/css/*
GaussDB(DWS)	spark.driver.extraClassPath=/usr/share/extension/dli/ spark-jar/datasource/dws/*
	spark.executor.extraClassPath=/usr/share/extension/dli/ spark-jar/datasource/dws/*

Datasource	Example Value	
HBase	spark.driver.extraClassPath=/usr/share/extension/dli/ spark-jar/datasource/hbase/*	
	spark.executor.extraClassPath=/usr/share/extension/dli/ spark-jar/datasource/hbase/*	
OpenTSDB	spark.driver.extraClassPath=/usr/share/extension/dli/ spark-jar/datasource/opentsdb/*	
	spark.executor.extraClassPath=/usr/share/extension/dli/ spark-jar/datasource/opentsdb/*	
RDS	spark.driver.extraClassPath=/usr/share/extension/dli/ spark-jar/datasource/rds/*	
	spark.executor.extraClassPath=/usr/share/extension/dli/ spark-jar/datasource/rds/*	
Redis	spark.driver.extraClassPath=/usr/share/extension/dli/ spark-jar/datasource/redis/*	
	spark.executor.extraClassPath=/usr/share/extension/dli/ spark-jar/datasource/redis/*	
Mongo	spark.driver.extraClassPath=/usr/share/extension/dli/ spark-jar/datasource/mongo/*	
	spark.executor.extraClassPath=/usr/share/extension/dli/ spark-jar/datasource/mongo/*	

**Table 12-4** Mapping between compute resource specification parameters onthe console and Spark Argument(--conf)

Console Parameter	Spark Argument( conf)	Description	Notes and Constraint s
Executor Memory Complete executor memory = spark.executor.m emory + spark.executor.m emoryOverhead	spark.executor.me mory	Executor memory, which is configurable.	-

Console Parameter	Spark Argument( conf)	Description	Notes and Constraint s
	spark.executor.me moryOverhead	Amount of off-heap memory for each executor in a Spark application. This parameter is not configurable. spark.executor.memor yOverhead=spark.exec utor.memory * spark.executor.memor yOverheadFactor	The minimum value is 384 MB. That is, when the value of <b>spark.exec</b> utor.mem ory multiplied by <b>spark.exec</b> utor.mem oryOverhe adFactor is less than 384 MB, the system automatic ally sets the value to 384 MB.
	spark.executor.me moryOverheadFac tor	This parameter determines the ratio of off-heap memory allocation to on-heap memory allocation. The default value is <b>0.1</b> for Spark applications run with the JAR file and <b>0.4</b> for those run with Python. This parameter is configurable.	The priority of <b>spark.exec</b> <b>utor.mem</b> <b>oryOverhe</b> <b>adFactor</b> is higher than that of <b>spark.kub</b> <b>ernetes.m</b> <b>emoryOve</b> <b>rheadFact</b> <b>or</b> .
Executor Cores	spark.executor.cor es	Number of executor cores, which is configurable.	-
Executors	spark.executor.inst ances	Number of executors, which is configurable.	-
Driver Cores	spark.driver.cores	Number of driver cores, which is configurable.	-

Console Parameter	Spark Argument( conf)	Description	Notes and Constraint s
Driver Memory Complete driver	spark.driver.memo ry	Driver memory, which is configurable.	-
Complete driver memory = spark.driver.mem ory + spark.edriver.me moryOverhead	spark.driver.memo ryOverhead	Amount of off-heap memory for each driver in a Spark application. This parameter is not configurable. spark.driver.memoryO verhead= spark.driver.memory * spark.driver.memoryO verheadFactor	The minimum value is 384 MB. That is, when the value of <b>spark.driv</b> <b>er.memor</b> <b>y</b> multiplied by <b>spark.driv</b> <b>er.memor</b> <b>yOverhea</b> <b>dFactor</b> is less than 384 MB, the system automatic ally sets the value to 284 MB
	spark.driver.memo ryOverheadFactor	This parameter determines the ratio of off-heap memory allocation to on-heap memory allocation. The default value is <b>0.1</b> for Spark applications run with the JAR file and <b>0.4</b> for those run with Python. This parameter is configurable.	The priority of <b>spark.driv</b> er.memor yOverhea dFactor is higher than that of <b>spark.kub</b> ernetes.m emoryOve rheadFact or.

Console Parameter	Spark Argument( conf)	Description	Notes and Constraint s
-	spark.kubernetes. memoryOverhead Factor	Amount of memory allocated outside the memory assigned to Spark executors. The default value is <b>0.1</b> for Spark applications run with the JAR file and <b>0.4</b> for those run with Python. This parameter is configurable.	The priority of spark.exec utor.mem oryOverhe adFactor and spark.driv er.memor yOverhea dFactor is higher than that of spark.kub ernetes.m emoryOve rheadFact or.

#### 5. (Optional) Configure dependencies.

#### Table 12-5 Dependency configuration parameters

Parameter	Description	
JAR Package Dependencies ( jars)	JAR file on which the Spark job depends. You can enter the JAR file name or the OBS path of the JAR file in the format of <b>obs:</b> // <i>Bucket name</i> / <i>Folder path</i> / <i>JAR file</i> <i>name</i> .	
Python File Dependencies ( py-files)	py-files on which the Spark job depends. You can enter the Python file name or the corresponding OBS path of the Python file. The format is as follows: <b>obs://Bucket</b> <i>name</i> / <i>Folder name</i> / <i>File name</i> .	
Other Dependencies ( files)	Other files on which the Spark job depends. You can enter the name of the dependency file or the corresponding OBS path of the dependency file. The format is as follows: <b>obs://Bucket name/Folder name/</b> <i>File name</i> .	
Group Name	If you select a group when creating a package, you can select all the packages and files in the group. For how to create a package, see <b>Creating a DLI Package</b> .	
	Spark 3.3. <i>x</i> or later does not support group information configuration.	

- 6. Set the following parameters in advanced settings:
  - Select Dependency Resources: For details about the parameters, see Table 12-6.
  - Configure Resources: For details about the parameters, see Table 12-7.

**NOTE** 

The parallelism degree of Spark resources is jointly determined by the number of Executors and the number of Executor CPU cores.

# Maximum number of tasks that can be concurrently executed = Number of Executors x Number of Executor CPU cores

You can properly plan compute resource specifications based on the compute CUs of the queue you have purchased.

Note that Spark tasks need to be jointly executed by multiple roles, such as **driver** and **executor**. So, the number of executors multiplied by the number of executor CPU cores must be less than the number of compute CUs of the queue to prevent other roles from failing to start Spark tasks. For more information about roles for Spark tasks, see **Apache Spark**.

#### Calculation formula for Spark job parameters:

- Number of CUs = Actual number of CUs = Max{(Driver Cores + Executors x Executor Cores), [(Driver Memory + Executors x Executor Memory)/4]}
- Memory = Driver Memory + (Executors x Executor Memory)

Table 12-6 Pa	arameters for selecting dependency resources

Parameter	Description		
modules	If the Spark version is <b>3.1.1</b> , you do not need to select a module. Configure <b>Spark parameters (conf)</b> .		
	Dependency modules provided by DLI for executing datasource connection jobs. To access different services, you need to select different modules.		
	CloudTable/MRS HBase: sys.datasource.hbase		
	DDS: sys.datasource.mongo		
	CloudTable/MRS OpenTSDB: sys.datasource.opentsdb		
	DWS: sys.datasource.dws		
	RDS MySQL: sys.datasource.rds		
	RDS PostGre: sys.datasource.rds		
	DCS: sys.datasource.redis		
	CSS: sys.datasource.css		
	DLI internal modules include:		
	• sys.res.dli-v2		
	• sys.res.dli		
	sys.datasource.dli-inner-table		
Resource	JAR package on which the Spark job depends.		
Package	Spark 3.3. <i>x</i> or later does not support this parameter. Configure resource package information in <b>jars</b> , <b>pyFiles</b> , and <b>files</b> .		

Parameter	Description		
Resource Specifications	Select a resource specification from the drop-down list box. The system provides three resource specification options for you to choose from.		
	Resource specifications involve the following parameters:		
	Executor Memory		
	Executor Cores		
	Executors		
	Driver Cores		
	Driver Memory		
	If modified, your modified settings of the items are used.		
Executor Memory	Customize the configuration item based on the selected resource specifications.		
	Memory of each Executor. It is recommended that the ratio of Executor CPU cores to Executor memory be 1:4.		
Executor Cores	Number of CPU cores of each Executor applied for by Spark jobs, which determines the capability of each Executor to execute tasks concurrently.		
Executors	Number of Executors applied for by a Spark job		
Driver Cores	Number of CPU cores of the driver		
Driver Memory	Driver memory size. It is recommended that the ratio of the number of driver CPU cores to the driver memory be 1:4.		

**Table 12-7** Resource specification parameters

 If you select 3.3.1 for Spark Version, you can configure compute resource specification parameters in Spark Argument(--conf). Note that the configuration priority of Spark Argument(--conf) is higher than that of Resource Specifications(Optional) in Advanced Settings.

 Table 12-4 describes the parameter mapping.

**NOTE** 

- When configuring compute resource specification parameters in **Spark Argument(--conf)**, you can use the unit M, G, or K. If the unit is not specified, the default unit is byte.
- Spark 3.3.1 or later includes notes and constraints on the compute resource specifications for jobs. For details, see **Table 12-8**.

#### 

If the compute resource specification is set too high, beyond the resource allocation capacity of the cluster or project, the job may fail to run due to resource request failures.

Table 12-8 value ranges of compute resources specifications				
Parameter	Elastic Resource Pool of Standard Edition After Modification	Elastic Resource Poo of Basic Edition		
Executor Memory	450 MB to 64 GB	450 MB to 16 GB		
Executor Cores	0 to 16	0 to 4		
Executors	Unlimited	Unlimited		
Driver Cores	0 to 16	0 to 4		
Driver Memory	450 MB to 64 GB	450 MB to 16 GB		
Job CU Quota	Unlimited	Unlimited		

Table 12-8	Value	ranges of	compute	resources	specifications
	value	runges or	compute	resources	specifications

Click **Execute** in the upper right corner of the Spark job editing page. 7.

After the message "Batch processing job submitted successfully" is displayed, you can view the status and logs of the submitted job on the Spark Jobs page.

#### **NOTE**

During the Spark job submission process, if the job fails to acquire resources successfully for an extended period, the job status will change to Failed after waiting for approximately 3 hours, indicating that the session has exited. For details about Spark job statuses, see Viewing Basic Information.

### 12.2 Example of a Typical Scenario: Reading and Querying OBS Data Using a Spark Jar Job

#### Scenario

DLI is fully compatible with open-source Apache Spark and allows you to import, query, analyze, and process job data by programming. This section describes how to write a Spark program to read and query OBS data, compile and package the code, and submit it to a Spark Jar job.

#### **Environment Preparations**

Before you start, set up the development environment.

ltem	Description
OS	Windows 7 or later
JDK	JDK 1.8.
IntelliJ IDEA	This tool is used for application development. The version of the tool must be 2019.1 or other compatible versions.
Maven	Basic configurations of the development environment. Maven is used for project management throughout the lifecycle of software development.

 Table 12-9
 Spark Jar job development environment

#### **Development Process**

The following figure shows the process of developing a Spark Jar job.

Figure 12-1 Development process



 Table 12-10 Process description

No	Phase	Softw are Portal	Description
1	Create a queue for general use.	DLI consol e	The DLI queue is created for running your job.
2	Upload data to an OBS bucket.	OBS consol e	The test data needs to be uploaded to your OBS bucket.
3	Create a Maven project and configure the POM file.	IntelliJ IDEA	Write your code by referring to the sample code for reading data from OBS.
4	Write code.		
5	Debug, compile, and pack the code into a Jar package.		

No	Phase	Softw are Portal	Description
6	Upload the Jar package to OBS and DLI.	OBS consol e DLI consol e	You can upload the generated Spark JAR package to an OBS directory and DLI program package.
7	Create a Spark Jar Job.	DLI consol e	The Spark Jar job is created and submitted on the DLI console.
8	Check the job execution result.	DLI consol e	You can check the job status and run logs.

#### Step 1: Create a Queue for General Purpose

Create a queue before submitting Spark jobs. In this example, we will create a general-purpose queue named **sparktest**.

- 1. Log in to the DLI management console.
- 2. In the navigation pane on the left, choose **Resources** > **Resource Pool**.
- 3. On the displayed page, click **Buy Resource Pool** in the upper right corner.
- 4. On the displayed page, set the parameters.

In this example, we will buy the resource pool in the **CN East-Shanghai2** region. **Table 12-11** describes the parameters.

Table 12-11 Parameters

Parameter	Description	Example Value
Region	Select a region where you want to buy the elastic resource pool.	CN East-Shanghai2
Project	Project uniquely preset by the system for each region	Default
Name	Name of the elastic resource pool	dli_resource_pool
Specificatio ns	Specifications of the elastic resource pool	Standard
CU Range	The maximum and minimum CUs allowed for the elastic resource pool	64-64

Parameter	Description	Example Value
CIDR Block	CIDR block the elastic resource pool belongs to. If you use an enhanced datasource connection, this CIDR block cannot overlap that of the data source. <b>Once</b> <b>set, this CIDR block cannot be</b> <b>changed.</b>	172.16.0.0/19
Enterprise Project	Select an enterprise project for the elastic resource pool.	default

- 5. Click **Buy**.
- 6. Click **Submit**.
- 7. In the elastic resource pool list, locate the pool you just created and click **Add Queue** in the **Operation** column.
- 8. Set the basic parameters listed below.

Parameter	Description	Example Value
Name	Name of the queue to add	dli_queue_01
Туре	<ul> <li>Type of the queue</li> <li>To execute SQL jobs, select For SQL.</li> <li>To execute Flink or Spark jobs, select For general purpose.</li> </ul>	For SQL jobs, select <b>For SQL</b> . For other scenarios, select <b>For general</b> <b>purpose</b> .
Engine	SQL queue engine. The options are <b>Spark</b> and <b>HetuEngine</b> .	Spark
Enterprise Project	Select an enterprise project.	default

#### Table 12-12 Basic parameters for adding a queue

9. Click **Next** and configure scaling policies for the queue.

Click **Create** to add a scaling policy with varying priority, period, minimum CUs, and maximum CUs.

Figure 12-2 shows the scaling policy configured in this example.



	Add Queue( )			
$\oslash$	Basic Configuration ——— (	2 Elastic Resources		
	View scaling policies of all queue 1.The priority ranges from 1 to 10 2.A new policy overwrites the def 3.You can only set the period to 1 4.The total minimum CUs of all of The maximum CUs of and one	s in 10. If you do not set the priority for a specific period, the default value ault policy. useues in an elastic resource pool cannot be more than the minimum in an adaptic resource pool cannot be more than the minimum (* 11 in an adaptic resource pool cannot be more than the minimum (* 11)	is 1. CUs of the pool.	
	Priority	Period	Min CU	Max CU
	1	00 - 24 -	- 16 +	- 16 +
			Create	

 Table 12-13
 Scaling policy parameters

Paramet er	Description	Example Value
Priority	Priority of the scaling policy in the current elastic resource pool. A larger value indicates a higher priority. In this example, only one scaling policy is configured, so its priority is set to <b>1</b> by default.	1
Period	The first scaling policy is the default policy, and its <b>Period</b> parameter configuration cannot be deleted or modified. The period for the scaling policy is from 00 to 24.	00–24
Min CU	Minimum number of CUs allowed by the scaling policy	16
Max CU	Maximum number of CUs allowed by the scaling policy	64

10. Click **OK**.

#### Step 2: Upload Data to OBS

- Create the **people.json** file containing the following content: {"name":"Michael"} {"name":"Andy", "age":30} {"name":"Justin", "age":19}
- 2. Log in to the OBS console. On the **Buckets** page, click the name of the OBS bucket you created. In this example, the bucket name is **dli-test-obs01**.
- 3. On the **Objects** tab, click **Upload Object**. In the displayed dialog box, upload the **people.json** file to the root directory of the OBS bucket.
- 4. In the root directory of the OBS bucket, click **Create Folder** to create a folder and name it **result**.
- 5. Click the **result** folder, click **Create Folder** on the displayed page to create a folder and name it **parquet**.

#### Step 3: Create a Maven Project and Configure the pom Dependency

This step uses IntelliJ IDEA 2020.2 as an example.

1. Start IntelliJ IDEA and choose **File** > **New** > **Project**.

#### Figure 12-3 Creating a project

	<u>F</u> i	le <u>E</u> dit	<u>V</u> iew	<u>N</u> avigate	<u>C</u> ode	Analy <u>z</u> e	<u>R</u> efa	ctor	<u>B</u> uild	R <u>u</u> n	<u>T</u> ools	VC <u>s</u>	<u>W</u> indow	<u>H</u> elp
Μv		<u>N</u> ew												
ŕ	Þ	<u>0</u> pen						Pro	oject f	rom E	xisting	I Sour	rces	
ect		Open <u>R</u> ec	ent				►	Pro	oject f	rom V	ersion	Contr		
roj		Close Pr	roject					Мо	dule					
<u>م</u>	ر کر	Se <u>t</u> tings				Ctrl+Alt	+S	Мо	dule fr	om Ex	isting	Sourc	es	
		Project	Struct	ture	Ctrl+	Alt+Shift	+S 🤇	) Jav	va Clas					
		File Pro	opertie	es			► Í	🛃 Кот	tlin Fi	le/Cl	ass			
		<u>S</u> ave All				Ctrl	+S 🤇	🛛 Gro	bovy Cl	ass				
	G	Reload A	ll fro	om Disk		Ctrl+Alt	+Y 🕯	Fi	Le					
		Invalida	ate Cao	ches / Res	tart		É	Scı	ratch F	ile	Ctrl+	Alt+S	Shift+In	sert
		Manage I	DE Set	ttings			► C	Pac	ckage					
		New Proj	jects S	Settings			► C	Py	thon Pa	ckage				
		Export						FXI	1L File					
	_	Doint						<b>p</b> ao	ckage-i	nfo.j	ava			

2. Choose Maven, set Project SDK to 1.8, and click Next.

#### Figure 12-4 Creating a project

🖳 New Project	
📕 Java	Project \$DK: 📑 1.8 java version "1.8.0_202"
🗬 Gradle	
<pre>Image: Image: Ima</pre>	
	Previous Next Cancel Helm

3. Set the project name, configure the storage path, and click **Finish**.

#### Figure 12-5 Creating a project



In this example, the Maven project name is **SparkJarObs**, and the project storage path is **D:\DLITest\SparkJarObs**.

4. Add the following content to the **pom.xml** file.



<dependency> <groupId>org.apache.spark</groupId> <artifactId>spark-sql\_2.11</artifactId> <version>2.3.2</version> </dependency>

</dependencies>

#### Figure 12-6 Modifying the pom.xml file



5. Choose **src** > **main** and right-click the **java** folder. Choose **New** > **Package** to create a package and a class file.

Figure 12-7 Creating a package

8	<u>F</u> ile <u>E</u> dit <u>V</u> iew	<u>N</u> avigate <u>C</u> ode Analy	<u>z</u> e <u>R</u> efactor <u>B</u> u:	ild R <u>u</u> n <u>T</u> ools VC <u>S</u>	<u>§ W</u> indow <u>H</u> elp SparkJarObs – pom.xml (SparkJarObs) – Admini
s	parkJarObs $ angle$ src $ angle$ m	ain 👌 🖿 java			
ct	🔲 Project 👻	🕀 🛨 🗢 –	🎢 pom.xml (Spa	arkJarObs) $ imes$	
òje	🗸 📄 SparkJarObs		1 xml</th <th>version="1.0" enc</th> <th>coding="UTF-8"?&gt;</th>	version="1.0" enc	coding="UTF-8"?>
Ъ	> 🖿 .idea		2 🔤 <proje< td=""><td></td><td></td></proje<>		
2	🗸 🖿 src				
2	∽ 🖿 main			xsi:schemaLoca	ation="http://maven.apache.org/POM/4.0.0 http://maven.ap
	<ul> <li>java</li> <li>java</li> <li>reso</li> <li>test</li> <li>mpom.xml</li> <li>SparkJar0</li> <li>IIII External Lib</li> <li>Scratches and</li> </ul>	New Cut Copy D Paste Find Usages Find in Path Replace in Path Analyze	5 <m< th=""><th>odelVersion&gt;4.0.0 Ctrl+X Ctrl+X Ctrl+V Alt+F7 Ctrl+Shift+F Ctrl+Shift+R</th><th>G/ImdalVarion&gt;</th></m<>	odelVersion>4.0.0 Ctrl+X Ctrl+X Ctrl+V Alt+F7 Ctrl+Shift+F Ctrl+Shift+R	G/ImdalVarion>
					<sup>6</sup> / <sub>4</sub> Python File <sup>m</sup> l <sup>1</sup> / <sub>4</sub> Jupyter Notebook <sup>#</sup> / <sub>4</sub> HTML File <sup>#</sup> / <sub>4</sub> Stylesheet <sup>#</sup> / <sub>4</sub> Styleshret <sup>#</sup> / <sub>4</sub> JupsCoript File

Set the package name as you need. In this example, set **Package** to **com.huawei.dli.demo**. Then, press **Enter**.

Create a Java Class file in the package path. In this example, the Java Class file is **SparkDemoObs**.

Figure 12-8 Creating a Java class file

8	<u>F</u> ile <u>E</u> dit <u>V</u> iew <u>N</u> avigate <u>C</u> ode Analy	<u>z</u> e <u>R</u> efa	ctor <u>B</u> uild R <u>u</u> n <u>T</u> ools VC <u>S</u> <u>W</u> indow <u>H</u> elp Sp	arkJarObs – SparkDemoObs.java – Administrator
S	parkJarObs $ angle$ src $ angle$ main $ angle$ java $ angle$ com $ angle$ hua			
sct		🕒 Spa	rkDemoObs.java ×	
, oj	✓ 📑 SparkJarObs D:\DLITest\SparkJarOb		<pre>package com.huawei.dli.demo;</pre>	
ē	> 🖿 .idea			
1	🗸 🖿 src			
_	✓ 🖿 main			
	🗸 🖿 java			
	✓ ➡ com.huawei.dli.demo			
	💿 SparkDemoObs			
	resources			
	> 🖿 test			
	🛃 SparkJarObs.iml			
	> III External Libraries			
	🏹 Scratches and Consoles			

#### Step 4: Write Code

Code the **SparkDemoObs** program to read the **people.json** file from the OBS bucket, create the temporary table **people**, and query data.

For the sample code, see **Sample Code**.

Import dependencies.

 import org.apache.spark.sql.Dataset;
 import org.apache.spark.sql.Row;
 import org.apache.spark.sql.SaveMode;
 import org.apache.spark.sql.SparkSession;

import static org.apache.spark.sql.functions.col;

- 2. Create Spark session **spark** using the AK and SK of the current account. SparkSession spark = SparkSession
  - .builder()
    - .config("spark.hadoop.fs.obs.access.key", "xxx")

.appName("java\_spark\_demo")
.getOrCreate();

- Replace xxx of "spark.hadoop.fs.obs.access.key" with the AK of the account.
- Replace yyy of "spark.hadoop.fs.obs.secret.key" with the SK of the account.

For details about how to obtain the AK and SK, see **How Do I Obtain the AK/SK Pair?** 

3. Read the **people.json** file from the OBS bucket.

dli-test-obs01 is the name of the sample OBS bucket. Replace it with the
actual OBS bucket name.
Dataset<Row> df = spark.read().json("obs://dli-test-obs01/people.json");
df.printSchema();

- Create temporary table **people** to read data. df.createOrReplaceTempView("people");
- Query data in the **people** table. Dataset<Row> sqlDF = spark.sql("SELECT \* FROM people"); sqlDF.show();
- Export people table data in Parquet format to the result/parquet directory of the OBS bucket. sqlDF.write().mode(SaveMode.Overwrite).parquet("obs://dli-test-obs01/result/parquet"); spark.read().parquet("obs://dli-test-obs01/result/parquet").show();
- Disable the **spark** session. spark.stop();

#### Step 5: Debug, compile, and pack the code into a JAR package.

1. Double-click **Maven** in the tool bar on the right, and double-click **clean** and **compile** to compile the code.

After the compilation is successful, double-click **package**.

# Dit Up: Bur Bondard Do Analyse Brother Bull Am [ost V g mone but Bondards: Sendersembaches (so v) and (so v) and

Figure 12-9 Compiling and packaging

The generated JAR package is stored in the **target** directory. In this example, **SparkJarObs-1.0-SNAPSHOT.jar** is stored in **D:\DLITest\SparkJarObs\target**.

Figure 12-10 Exporting the JAR file

#### Step 6: Upload the JAR Package to OBS and DLI

• Spark 3.3 or later:

You can only set the **Application** parameter when creating a Spark job and select the required JAR file from OBS.

- a. Log in to the OBS console and upload the JAR file to the OBS path.
- Log in to the DLI console. In the navigation pane, choose Job Management > Spark Jobs.
- c. Locate the row containing a desired job and click **Edit** in the **Operation** column.
- d. Set **Application** to the OBS path in **a**.

Figure 12-11	Configuring	the application

Select a Queue		
* Queues		•
* Spark Version	3.3.1	•
Job Configurations		
Job Name(name)	Enter a name.	
* Application	obs://	6

• Versions earlier than Spark 3.3:

Upload the JAR file to OBS and DLI.

- a. Log in to the OBS console and upload the JAR file to the OBS path.
- b. Upload the file to DLI for package management.

- i. Log in to the DLI management console and choose **Data Management** > **Package Management**.
- ii. On the **Package Management** page, click **Create** in the upper right corner.
- iii. In the Create Package dialog, set the following parameters:
  - 1) Type: Select JAR.
  - 2) **OBS Path**: Specify the OBS path for storing the package.
  - 3) Set **Group** and **Group Name** as required for package identification and management.
- iv. Click **OK**.

Create Package					>
Туре	JAR	PyFile	File	ModelFile	
* OBS Path	Specify each p	arameter on a s	eparate line.		
Group	Use existin	g Use	new	Do not use	
★ Group Name				•	
Tags	It is recommended different cloud re	ed that you use T sources. View pi er a tag key and	'MS's predefin redefined tags a tag value be	ed tag function to a C elow.	dd the same tag to
	Enter a tag key	1	Enter	a tag value	Add
	20 tags available	for addition.	Cancel		

Figure 12-12 Creating a package

#### Step 7: Create a Spark Jar Job

- Log in to the DLI console. In the navigation pane, choose Job Management > Spark Jobs.
- 2. On the Spark Jobs page, click Create Job.
- 3. On the displayed page, configure the following parameters:
  - Queue: Select the created queue. For example, select the queue sparktest created in Step 1: Create a Queue for General Purpose.
  - Select a supported Spark version from the drop-down list. The latest version is recommended.
  - Job Name (--name): Name of the Spark Jar job. For example, SparkTestObs.

- Application: Select the package uploaded in Step 6: Upload the JAR Package to OBS and DLI. For example, select SparkJarObs-1.0-SNAPSHOT.jar.
- Main Class (--class): The format is program package name + class name.
   For example, com.huawei.dli.demo.SparkDemoObs.

You do not need to set other parameters.

For more information about Spark JAR job submission, see **Creating a Spark Job**.

4. Click **Execute** to submit the Spark Jar job. On the **Spark Jobs** page, check the status of the job you submitted.

#### Figure 12-13 Job status

				Search by name by default.		
Job ID	Nome	Gueues V	Username	Status T Created	Last Modified	Operation
				<ul> <li>Starting</li> </ul>		

#### **Step 8: Check Job Execution Result**

- 1. On the **Spark Jobs** page, check the status of the job you submitted. The initial status is **Starting**.
- 2. If the job is successfully executed, the job status is **Finished**. Click **More** in the **Operation** column and select **Driver Logs** to check the run log.

Figure 12-14 Selecting Diver Logs







- 3. If the job is successfully executed, go to the **result/parquet** directory in the OBS bucket to check the generated **parquet** file.
- 4. If the job fails to be executed, click **More** in the **Operation** column and select **Driver Logs** to check detailed error information.

For example, the following figure shows that when you create the Spark Jar job, you did not add the package path to the main class name.



Figure 12-16 Error information

In the **Operation** column, click **Edit**, change the value of **Main Class** to **com.huawei.dli.demo.SparkDemoObs**, and click **Execute** to run the job again.

#### Follow-up Guide

- If you want to use Spark Jar jobs to access other data sources, see Using Spark Jobs to Access Data Sources of Datasource Connections.
- If you want to create a database and table using a Spark Jar job, see Using Spark Jobs to Access DLI Metadata.

#### Sample Code

#### **NOTE**

Hard-coded or plaintext **access.key** and **secret.key** pose significant security risks. To ensure security, encrypt your AK and SK, store them in configuration files or environment variables, and decrypt them when needed.

package com.huawei.dli.demo;

```
import org.apache.spark.sql.Dataset;
import org.apache.spark.sql.Row;
import org.apache.spark.sql.SaveMode;
import org.apache.spark.sql.SparkSession;
import static org.apache.spark.sql.functions.col;
public class SparkDemoObs {
  public static void main(String[] args) {
     SparkSession spark = SparkSession
           .builder()
           .config("spark.hadoop.fs.obs.access.key", "xxx")
           .config("spark.hadoop.fs.obs.secret.key", "yyy")
           .appName("java_spark_demo")
           .getOrCreate();
     // can also be used --conf to set the ak sk when submit the app
     // test json data:
     // {"name":"Michael"}
     // {"name":"Andy", "age":30}
// {"name":"Justin", "age":19}
     Dataset<Row> df = spark.read().json("obs://dli-test-obs01/people.json");
     df.printSchema();
     // root
     // |-- age: long (nullable = true)
     // |-- name: string (nullable = true)
     // Displays the content of the DataFrame to stdout
     df.show();
     // +----+
     // | age| name|
     // +----+
     // |null|Michael|
     // | 30| Andy|
// | 19| Justin|
     // +----+
     // Select only the "name" column
     df.select("name").show();
     // +----+
     // | name|
     // +----+
     // |Michael|
     // | Andy
     // | Justin|
     // +----+
     // Select people older than 21
     df.filter(col("age").gt(21)).show();
     // +---+
     // |age|name|
     // +---+----
     // | 30|Andy|
     // +---+
     // Count people by age
     df.groupBy("age").count().show();
     // +----+
     // | age|count|
```

```
// +----+
  // | 19| 1|
  // |null| 1|
  // | 30| 1|
  // +----+
  // Register the DataFrame as a SQL temporary view
  df.createOrReplaceTempView("people");
  Dataset<Row> sqlDF = spark.sql("SELECT * FROM people");
  sqlDF.show();
  // +----+
  // | age| name|
  // +----+
  // |null|Michael|
  // | 30| Andy|
  // | 19| Justin|
  // +----+
  sqlDF.write().mode(SaveMode.Overwrite).parquet("obs://dli-test-obs01/result/parquet");
  spark.read().parquet("obs://dli-test-obs01/result/parquet").show();
   spark.stop();
}
```

# 12.3 Setting the Priority for a Spark Job

#### Scenario

In actual job running, it is necessary to prioritize and ensure the normal running of important and urgent tasks due to their varying levels of importance and urgency. This requires providing the necessary compute resources for their normal operations.

DLI offers a feature to set job priorities for each Spark job, which prioritizes the allocation of compute resources to higher priority jobs when resources are limited.

#### **NOTE**

}

You can set the priority for Spark 2.4.5 or later jobs.

#### Notes

- Priorities cannot be set for jobs running in queues within an elastic resource pool of the basic edition.
- You can assign a priority level of 1 to 10 for each job, with a larger value indicating a higher priority. Compute resources are preferentially allocated to high-priority jobs. That is, if compute resources required for high-priority jobs are insufficient, compute resources for low-priority jobs are reduced.
- Spark jobs running on a general-purpose queue have a default priority level of 3.
- To change the priority for a job, you must first stop the job, change the priority level, and then submit the job for the modification to take effect.

#### Procedure

Enter the following statement in the **Spark Arguments(--conf)** text box, where *x* indicates the priority value:

spark.dli.job.priority=x

- 1. Log in to the DLI management console.
- 2. In the navigation pane on the left, choose **Job Management** > **Spark Jobs**.
- 3. Select the job for which you want to set the priority and click **Edit** in the **Operation** column.
- 4. Set the **spark.dli.job.priority** parameter in the **Spark Arguments(--conf)** text box.

Figure 12-17 Example configuration for a Spark job

spark.dli.jo	ob.priority=8	

# 12.4 Querying Logs for Spark Jobs

Spark Arguments(--conf)

#### Scenario

DLI job buckets are used to store temporary data generated during DLI job running, such as job logs and results.

This section describes how to configure a bucket for DLI jobs on the DLI console and obtain Spark job logs.

#### Notes

- To avoid disordered job results, do not use the OBS bucket configured for DLI jobs for any other purposes.
- DLI jobs must be set and modified by the main account as IAM users do not have required permissions.
- You cannot view the logs for DLI jobs before configuring a bucket.
- You can configure lifecycle rules to periodically delete objects from buckets or change storage classes of objects.
- Exercise caution when modifying the job bucket, as it may result in the inability to retrieve historical data.
- Spark log splitting rules:
  - Split by size: By default, each log file has a maximum size of 128 MB.
  - Split by time: A new log file is automatically created every hour.

#### Prerequisites

Before the configuration, create an OBS bucket or parallel file system (PFS). In big data scenarios, you are advised to create a PFS. PFS is a high-performance file system provided by OBS, with access latency in milliseconds. PFS can achieve a

bandwidth performance of up to TB/s and millions of IOPS, which makes it ideal for processing high-performance computing (HPC) workloads.

For details about PFS, see "Parallel File System Feature Guide" in the *Object Storage Service User Guide*.

#### Configuring a Bucket for DLI Jobs

- 1. In the navigation pane of the DLI console, choose **Global Configuration** > **Project**.
- 2. On the **Project** page, click and next to **Job Bucket** to configure bucket information.

#### Figure 12-18 Project

ata Lake Insigl	ht a	Project
verview		
QL Editor		Job Bucket
b Management	$\sim$	Bucket Name: rain
sources	~	This bucket is used to store temporary data generated by DLI, such as job logs and job results. Do not use this bucket for other purposes. If you do not create this bucket, you will not be able to view job logs. You can use the main
ita Management	$\sim$	account to set and modify the bucket. Sub-users do not have modification permissions. You can set a lifecycle rule to periodically delete objects in a
b Templates	~	bucket or change its storage class. Exercise caution when you modify the lifecycle rule to prevent historical data being deleted by mistake.
tasource Connection	ns	
obal Configuration	^	
Global Variables		
SQL Inspector	<	
Project		
Service Authorization		

- 3. Click  $\square$  to view available buckets.
- 4. In the displayed **OBS** dialog box, click the name of a bucket or search for and click a bucket name and then click **OK**. In the **Set Job Bucket** dialog box, click **OK**.

Temporary data generated during DLI job running will be stored in the OBS bucket.

Figure 12-19 Setting the job bucket

Set Job Bucket	×
Parallel file buckets are recommended.	×
★ Job Bucket	
	Cancel OK

#### **Querying Logs for Spark Jobs**

- Log in to the DLI console. In the navigation pane on the left, choose Job Management > Spark Jobs.
- 2. Select the Spark job whose jobs you want to query, click **More** in the **Operation** column, and select **View Log**.

The system automatically switches to the log path of the DLI job bucket.

3. On the **Files** tab, select the log file of the desired date and time and click **Download** in the **Operation** column to download the file to your local host.

Figure 12-20 Downloading Spark job logs

Files / jobs / logs /	/ 2024-09-18_10-48-	55 /	Ō	
Files Fragments				
You can use OBS Browser+ to mo Preview Objects in OBS from My E Upload File Create	ve a file to any other folder in this pa Browser? Folder Delete M	rallel file system. For security rea ore $\checkmark$	isons, files cannot be previewed online when you access them	from a browser. To preview files online, see How
S C Enter a file name pre	efix.			
Name	Storage Class	Size	Last Modified	Operation
	Standard	36.55 KB	Sep 18, 2024 14:59:00 GMT+08:00	Download Share More $\vee$
Spark.log	Standard	1.29 KB	Sep 18, 2024 14:34:00 GMT+08:00	Download Share More ~

# 12.5 Managing Spark Jobs

#### Viewing Basic Information

On the **Overview** page, click **Spark Jobs** to go to the SQL job management page. Alternatively, you can click **Job Management** > **Spark Jobs**. The page displays all Spark jobs. If there are a large number of jobs, they will be displayed on multiple pages. DLI allows you to view jobs in all statuses.

Parameter	Description
Job ID	ID of a submitted Spark job, which is generated by the system by default.
Name	Name of a submitted Spark job.
Queues	Queue where the submitted Spark job runs
Username	Name of the user who executed the Spark job
Status	<ul> <li>Job status. The following values are available:</li> <li>Starting: The job is being started.</li> <li>Running: The job is being executed.</li> <li>Failed: The session has exited.</li> <li>Finished: The session is successfully executed.</li> <li>Restoring: The job is being restored.</li> </ul>
Created	Time when a job is created. Jobs can be displayed in ascending or descending order of the job creation time.
Last Modified	Time when a job is completed.

Parameter	Description					
Operation	• <b>Edit</b> : You can modify the current job configuration and re- execute the job.					
	• <b>SparkUI</b> : After you click this button, the Spark job execution page is displayed.					
	NOTE					
	• The SparkUI page cannot be viewed for jobs in the <b>Starting</b> state.					
	<ul> <li>Currently, only the latest 100 job information records are displayed on the SparkUI of DLI.</li> </ul>					
	• <b>Terminate Job</b> : Cancel a job that is being started or running.					
	• <b>Re-execute</b> : Run the job again.					
	• Archive Log: Save job logs to the temporary bucket created by DLI.					
	• <b>Commit Log</b> : View the logs of submitted jobs.					
	• Driver Log: View the logs of running jobs.					

#### **Re-executing a Job**

On the **Spark Jobs** page, click **Edit** in the **Operation** column of the job. On the Spark job creation page that is displayed, modify parameters as required and execute the job.

#### Searching for a Job

On the **Spark Jobs** page, select **Status** or **Queues**. The system displays the jobs that meet the filter condition in the job list.

#### **Terminating a Job**

On the **Spark Jobs** page, choose **More** > **Terminate Job** in the **Operation** column of the job that you want to stop.

### 12.6 Managing Spark Job Templates

#### Scenario

You can modify a sample template to meet the Spark job requirements, saving time for editing SQL statements.

Currently, the cloud platform does not provide preset Spark templates. You can customize Spark job templates. This section describes how to create a Spark job template.

#### **Creating a Management Job Template**

To create a Spark job template, you can save a Spark job information as a template.

- 1. In the left navigation pane of the DLI management console, choose **Job Templates** > **Spark Templates**. The **Spark Templates** page is displayed.
- 2. Configure job parameters by referring to **Creating a Spark Job**.
- 3. After you finish editing the job, click **Save as Template**.
- 4. Enter the template name and other information as you need.
- 5. Set a template group for future management.
- 6. Click OK.

# **13** Developing a DLI Spark Job in DataArts Studio

Huawei Cloud DataArts Studio provides a one-stop data governance platform that integrates with DLI for seamless data integration and development, enabling enterprises to manage and control their data effectively.

This section describes how to develop a DLI Spark job using DataArts Factory of DataArts Studio.

#### Procedure

- 1. Obtain a demo JAR file of the Spark job and associate it with DataArts Factory on the DataArts Studio console.
- 2. On the DataArts Studio console, create a DataArts Factory job and submit the Spark job through the DLI Spark node.

#### **Environment Preparations**

- Prepare a DLI resource environment.
  - Configure a DLI job bucket.

Before using DLI, you need to configure a DLI job bucket. The bucket is used to store temporary data generated during DLI job running, such as job logs and results.

For details, see **Configuring a DLI Job Bucket**.

- Prepare a JAR file and upload it to an OBS bucket.

The Spark job code used in this example comes from the Maven repository (download address: https://repo.maven.apache.org/maven2/org/apache/spark/spark-examples\_2.10/1.1.1/spark-examples\_2.10-1.1.1.jar). This Spark job is used to calculate the approximate value of  $\pi$ .

After obtaining the JAR file of the Spark job code, upload the JAR file to the OBS bucket. In this example, the storage path is **obs://dlfexample/spark-examples\_2.10-1.1.1.jar**.

 Create an elastic resource pool and create general-purpose queues within it. An elastic resource pool offers compute resources (CPU and memory) required for running DLI jobs, which can adapt to the changing demands of services.

You can create general-purpose queues within an elastic resource pool to submit Spark jobs. These queues are associated with specific jobs and data processing tasks, and serve as the basic unit for resource allocation and usage within the pool. This means queues are specific compute resources required for executing jobs.

For details, see **Creating an Elastic Resource Pool and Creating Queues** Within It.

#### • Prepare a DataArts Studio resource environment.

- Buy a DataArts Studio instance.

Buy a DataArts Studio instance before submitting a DLI job using DataArts Studio.

For details, see **Buying a DataArts Studio Basic Package**.

- Access the DataArts Studio instance's workspace.
  - i. After buying a DataArts Studio instance, click Access.

Figure 13-1 Accessing a DataArts Studio instance

Enterprise Project: default						
Version	Enterprise	Billing Mode	Yearly/Monthly			
Created	Oct 23, 2019 10:02:13 GMT+08:00	Name				
Expires	Oct 22, 2020 23:59:59 GMT+08:00	Instance ID	9ab2da986bf34d70b62c323850			
Order No.	-	Status	Valid			
Description	- <u>/</u>					
	ccess   Renew	舍 Buy	Upgrade •			

ii. Click the **Workspaces** tab to access the data development page.

By default, a workspace named **default** is created for the user who has purchased the DataArts Studio instance, and the user is assigned the administrator role. You can use the default workspace or create one.

For how to create a workspace, see **Creating and Managing a Workspace**.

	Back to Instar	nce List			
Dashboard	Workspaces	Roles	Industry	Assets	Tags
Create Works	space				
Q Select a pro	operty or enter a keyv	vord.			
Name/ID 👙	Status	÷ Mo	ode ≑	Description	÷
	📀 En	able Sir	nple		

Figure 13-2 Accessing the DataArts Studio instance's workspace

#### Figure 13-3 Accessing DataArts Studio's data development page

Data	0 🛛 🗑 📾 💭 🔍		
Development	AL Owne O C JE :		
Overview	Enter a keyword.		Create Script
		· (7)	Develop, debug, and run scripts online, or run developed scripts in jobs.
Data Development			+ SQL + Hive SQL + DLI SQL + DWS SQL + Spark SQL + Spark Python + Filmk SQL
Develop Script			+ RDS SQL + Presto SQL + HetuEngine SQL + ClickHouse SQL + Impala SQL + Shell
Develop Job			+ Python + Doris SQL
Monitoring			
Overview			Create Job
Job Monitoring	< .		Press and date the andre as a constra and connect them to easily devolve inte
Monitor Instance			ends and and the model of a carrier and connect metric carry develop year.
Monitor PatchData			+ Create Job + Create Data Migration Job
Duty Schedules			
Manage	8		
Notification			Create Data Connection
Manage Backup		· Cor	Configure data source information and create data connections, through which you can access data sources when developing scripts and jobs.
Operation History		Cest	+ Create Data Connection

#### Step 1: Obtain the Spark Job Code

- **Step 1** After obtaining the JAR file of the Spark job code, upload the JAR file to the OBS bucket. The storage path is **obs://dlfexample/spark-examples\_2.10-1.1.1.jar**.
- Step 2 On the DataArts Studio console, locate a workspace and click DataArts Factory.
- **Step 3** In the navigation pane on the left, choose **Configuration** > **Manage Resource**.
- Step 4 On the displayed page, click Create Resource, create a resource named sparkexample on DataArts Factory, and associate it with the JAR file obtained in Step 1.

								_	
Manage Resource									
Enter a keyword. Q 🕀 C ᠄	Create Re	esource							
Resources	Name	Ту	pe	Storage Path		Depended P	Created By		Cr
	c	Create Resou	rce					×	
	*	Name	spark-exa	mple					
		Туре	jar				•		
		Resource Location (	OBS				•		
		Main JAR Package	obs://xxx/	oox.jar			B		
	D	Depended Package 🤅	) ±						
	*	Select Directory	/Resource	s/			$\oplus$		
	D	Description							
							// 0/256		
				ОК	Cancel				

Figure 13-4 Creating a resource

----End

#### Step 2: Submit a Spark Job

You need to create a job in DataArts Factory and submit the Spark job using the DLI Spark node.

Step 1 In the navigation pane on the left, choose Data Development > Develop Job. In the displayed job list, locate the target directory, right-click it, and select Create Job. In the dialog box that appears, set Job Name to job\_DLI\_Spark and set other parameters as needed.

×

#### Figure 13-5 Creating a job

#### Create Job

A maximum of 100,000 jobs can be created. You can create 99,315 more jobs.

* Job Name	job_DLI_Spark
Job Type	Batch processing     Real-time processing
Mode	Pipeline      Single task
Select Directory	/Jobs/ (+)
Owner ⑦	★ ⊕
Priority	● High ○ Medium ○ Low
Agency ⑦	Select an agency.
Log Path	obs://dlf-log-62099355b894428e8916573ae635f1f9/
	I agree to create OBS bucket obs://dlf-log- 62099355b894428e8916573ae635f1f9/. This bucket is used only for storing run logs of DLF jobs.
	To change the log path, go to the WorkSpaces page.
	OK Cancel

**Step 2** Go to the job development page, drag the DLI Spark node to the canvas, and click the node to configure its properties.

-	- <del>.</del>	DLI Spark		Node
		Properties	^	Properti
		* Node Name		es
		DLI_Spark_2032		lineag
		* Job Name		gelnfo
DLI_Spark_2032	[	* DLI Queue		
		•	•	
		Spark Version		
			•	
		Job Running Resource		
		8-core 32 GB memory 💌	] ~	
	E	Major Job Class 🕜		
		* Spark program resource package ⑦		
			$(\neq)$	
		* Resource Type ⑦		<b>O</b> Data
		OBS path OLI program package		IA r



Description of key properties:

- **DLI Queue**: Select a DLI queue.
- Job Running Resource: Maximum CPU and memory resources that can be used by a DLI Spark node.
- Major Job Class: major class of a DLI Spark node. In this example, the major class is org.apache.spark.examples.SparkPi.
- Spark program resource package: Select the resources created in Step 4.

**Step 3** Click  $\triangleright$  to test the job.

Figure 13-7 Job logs (for reference only)

Logs

```
[INFO] [Oct 09,2018 10:22:43 GMT +08:00] : The job starts to run.
[INFO] [Oct 09,2018 10:22:53 GMT +08:00] : Node DLI_Spark started to run.
```

**Step 4** If there are no errors in the logs, save and submit the job.

----End

# **14** Submitting a Spark Job Using a Notebook Instance

Notebook is an interactive data analysis and mining module that has been deeply optimized based on the open-source JupyterLab. It provides online development and debugging capabilities for writing and debugging model training code. After connecting DLI to a notebook instance, you can write code and develop jobs using Notebook's web-based interactive development environment, as well as flexibly perform data analysis and exploration. This section describes how to submit a DLI job using a notebook instance.

For how to perform operations on Jupyter Notebook, see **Jupyter Notebook Documentation**.

Use notebook instances to submit DLI jobs in scenarios involving online development and debugging. You can perform data analysis and exploration seamlessly, without the need to set up a development environment.

#### Notes

- To use this function, which is currently in the whitelist, submit a request by choosing **Service Tickets** > **Create Service Ticket** in the upper right corner of the management console.
- Deleting an elastic resource pool on the DLI management console will not delete the associated notebook instances. If you no longer need the notebook instances, log in to the ModelArts management console to delete them.

#### Procedure

- Create an elastic resource pool and create general-purpose queues within it. To create a notebook instance on DLI, first create an elastic resource pool and create a general-purpose queue within the pool. So the queue can offer compute resources required to run DLI jobs. See Step 1: Create an Elastic Resource Pool and Create General-Purpose Queues Within It.
- 2. Create a VPC and security group.

After configuring the elastic resource pool, the pool will prepare the components required for the notebook instance. See **Step 2: Create a VPC and Security Group**.
3. Create an enhanced datasource connection, which will be used to connect the DLI elastic resource pool to a notebook instance.

See Step 3: Create an Enhanced Datasource Connection.

- Prepare a custom image.
   See Step 4: Register a ModelArts Custom Image.
- 5. Create a custom agency, which will be used to access a notebook instance. See **Step 5: Create a DLI Custom Agency**.
- Create a notebook instance in the DLI elastic resource pool.
   See Step 6: Create a Notebook Instance in the DLI Elastic Resource Pool.
- 7. Configure the notebook instance to access DLI or LakeFormation metadata.
  - (Optional) Configuring the Notebook Instance to Access DLI Metadata
  - (Optional) Configuring the Notebook Instance to Access LakeFormation Metadata
- 8. Write and debug code in JupyterLab.

On the JupyterLab home page, you can edit and debug code in the **Notebook** area. See **Step 8: Use the Notebook Instance to Write and Debug Code**.

## Notes and Constraints

- To submit a DLI job using a notebook instance, you must have a generalpurpose queue within an elastic resource pool.
- Each elastic resource pool is associated with a unique notebook instance.
- Temporary data generated during the running of notebook jobs is stored in DLI job buckets in a parallel file system.
- Manage notebook instances on the ModelArts management console. For details, see **Managing Notebook Instances**.
- Notebook instances are used for code editing and development, and associated queues are used for job execution.

To change the queue associated with a notebook instance, perform related operations on the ModelArts management console.

## Step 1: Create an Elastic Resource Pool and Create General-Purpose Queues Within It

- 1. Create an elastic resource pool.
  - a. Log in to the DLI management console. In the navigation pane on the left, choose **Resources** > **Resource Pool**.
  - b. On the displayed page, click **Buy Resource Pool** in the upper right corner.
  - c. On the displayed page, set the parameters based on **Creating an Elastic Resource Pool and Creating Queues Within It**.
    - **CU range**: Reserve over 16 CUs.
    - CIDR Block: Make sure the CIDR block differs from the following ones:

172.18.0.0/16, 172.16.0.0/16, 10.247.0.0/16

- d. Click **Buy**.
- e. Click **Submit**. Wait until the elastic resource pool changes to the **Available** state.

## 2. Create a general-purpose queue within the elastic resource pool.

- a. Locate the elastic resource pool in which you want to create queues and click **Add Queue** in the **Operation** column.
- b. On the Add Queue page, configure basic information about the queue. For details about the parameters, see Creating an Elastic Resource Pool and Creating Queues Within It.

## Set Type to For general purpose.

- c. Click **Next**. On the displayed page, configure a scaling policy for the queue.
- d. Click **OK**.

## Step 2: Create a VPC and Security Group

- Create a VPC.
  - a. Log in to the VPC management console and click **Create VPC** in the upper right corner of the page.
  - b. On the **Create VPC** page, set the parameters as prompted.

For details about the parameters, see **Creating a VPC**.

Make sure not to set IPv4 CIDR Block to any of the following ones:

172.18.0.0/16, 172.16.0.0/16, 10.247.0.0/16

## • Create a security group.

- a. On the network console, access the **Security Groups** page.
- b. Click **Create Security Group** in the upper right corner.

On the displayed page, set security group parameters as prompted.

For details about the parameters, see Creating a Security Group.

Ensure that the security group allows TCP ports **8998** and **30000–32767** to pass through the CIDR block of the DLI elastic resource pool.

## Step 3: Create an Enhanced Datasource Connection

- 1. Log in to the DLI management console.
- 2. In the navigation pane on the left, choose **Datasource Connections**.
- 3. On the displayed **Enhanced** tab, click **Create**.

Set parameters based on Table 6-3.

When creating an enhanced datasource connection:

- Resource Pool: Select the elastic resource pool created in Step 1: Create an Elastic Resource Pool and Create General-Purpose Queues Within It.
- VPC: Select the VPC created in Step 2: Create a VPC and Security Group.

## Step 4: Register a ModelArts Custom Image

Based on the preset MindSpore image provided by ModelArts and the ModelArts CLI, you can load the image creation template and modify a Dockerfile to create an image. Then, register the image.

For details about the ModelArts CLI, see **ma-cli image Commands for Building Images**.

 Base image address: swr.{endpoint}/atelier/pyspark\_3\_1\_1:develop-remotepyspark\_3.1.1-py\_3.7-cpu-ubuntu\_18.04-x86\_64uid1000-20230308194728-68791b4

Replace *endpoint* (region name) with the actual one.

For example, the endpoint of AP-Singapore is **ap-southeast-3.myhuaweicloud.com**.

The combined base image address is **swr.apsoutheast-3.myhuaweicloud.com/atelier/pyspark\_3\_1\_1:develop-remotepyspark\_3.1.1-py\_3.7-cpu-ubuntu\_18.04-x86\_64uid1000-20230308194728-68791b4**.

• For how to create and register a custom image on ModelArts, see **Creating a Custom Image Using Dockerfile**.

## Step 5: Create a DLI Custom Agency

Create a DLI custom agency, which will be used to access a notebook instance. For details, see **Creating a Custom DLI Agency**.

Make sure the agency includes the following permissions: **ModelArts FullAccess**, **DLI FullAccess**, **OBS Administrator**, and IAM permission to pass agencies to cloud services.

If using role/policy-based authorization, grant the IAMiam:agencies:\* permission.

"Version": "1.1", "Statement": [ { "Effect": "Allow", "Action": [ "iam:agencies:\*" ] }, { "Effect": "Deny", "Action": [ "iam:agencies:update\*", "iam:agencies:delete\*", "iam:agencies:create\*" ] } ] }

## Step 6: Create a Notebook Instance in the DLI Elastic Resource Pool

## **NOTE**

Log in to the ModelArts management console. In the navigation pane on the left, choose **System Management** > **Permission Management**. On the displayed page, check if the access authorization for ModelArts is configured. The new agency must include the IAM permission to pass agencies to cloud services. For details about permission policies, see **Step 5: Create a DLI Custom Agency**.

- **Step 1** On the DLI elastic resource pool page, preset DLI resource information required for creating a notebook instance.
  - 1. Log in to the DLI management console. In the navigation pane on the left, choose **Resources** > **Resource Pool**.
  - On the displayed page, locate the elastic resource pool created in Step 1: Create an Elastic Resource Pool and Create General-Purpose Queues Within It.
  - 3. Click More in the Operation column and select Notebook (New).
  - 4. In the slide-out panel, click **Create Notebook**. In the dialog box that appears, set the following parameters:
    - Image: Select the image registered in Step 4: Register a ModelArts Custom Image.
    - Queue: Select the queue created in Step 1: Create an Elastic Resource Pool and Create General-Purpose Queues Within It.
    - **Spark Version**: **3.3.1** is recommended.
    - Enhanced: Select the enhanced datasource connection created in Step 3: Create an Enhanced Datasource Connection.

**Figure 14-1** Presetting DLI resource information required for creating a notebook instance

mage	pyspark_3_1_1	~
	Can't find the image version you want? Create	e
Queue		~
Spark	3.3.1	~
/ersion		
Enhanced		~
Enhanced		~

- 5. Click **OK**. The instance creation page is displayed.
- **Step 2** On the displayed page, set notebook instance parameters.
  - Create a notebook instance.
     For details about the parameters, see Creating a Notebook Instance.

Set the parameters as follows:

- Image: Select the image registered in Step 4: Register a ModelArts Custom Image.
- VPC Access: Enable VPC access.

D NOTE

Contact customer support to enable the VPC access function for the notebook instance.

Select the security group created in **Step 2: Create a VPC and Security Group**. The security group must allow TCP ports **8998** and **30000–32767** to pass through the CIDR block of the DLI elastic resource pool.

Click Create.

**Step 3** Connect the notebook instance to DLI.

- 1. In the notebook instance list, locate the notebook instance and click **Open** in the **Operation** column to access the notebook instance page.
- 2. On the notebook instance page, click **connect** in the upper right corner to connect to DLI.

Figure 14-2 Connecting to DLI



3. In the **Connect Cluster** dialog box, configure job running information.

Connect Cluster	×
*Service Type	DLI -
* Pool Name	
*Queue Name	notebook 👻
* Spark Version	•
Advanced Settings	
User Defined Image	image_org/image_name:tag
Spark Arguments( conf)	Enter a parameter in the key=value format and separate multiple parameters by pressing Enter
Jar Package	Enter obs path of file and separate
Dependencies(jars)	multiple path by pressing Enter
Python File Dependencies(py- files)	Enter obs path of file and separate multiple path by pressing Enter
Other	Enter obs path of file and separate
Dependencies( files) Resource Config	multiple path by pressing Enter
thereare boining	
Driver Memory	1 GB 👻

## Table 14-1 Connect Cluster

Figure 14-3 Connect Cluster

Paramete r	Description	Example Value
Service Type	Name of the service to connect	DLI

Paramete r	Description	Example Value
Pool Name	Elastic resource pool of the queue where the notebook job is running	In this example, set this parameter to the elastic resource pool created in Step 1: Create an Elastic Resource Pool and Create General- Purpose Queues Within It.
Queue Name	Queue where the notebook job is running	In this example, set this parameter to the queue created in <b>Step 1</b> : <b>Create an Elastic</b> <b>Resource Pool and</b> <b>Create General-</b> <b>Purpose Queues</b> <b>Within It.</b>
Spark Version	Spark version	Only Spark 3.3.1 currently supports submitting DLI jobs using notebook instances.
Spark Argument s(conf)	Allows you to configure custom parameters for the DLI job.	See Table 14-2.

Table 14-2 Common Spark parameters

Parameter	Description
spark.dli.job.agenc	Name of the agency for the DLI job
y.name	When Flink 1.15, Spark 3.3, or a later version is used to execute jobs, you need to add information about the new agency to the job configuration.
	Example configuration:
	In this example, set this parameter to <b>dli_notebook</b> .
	spark.dli.job.agency.name= <i>dli_notebook</i>
spark.sql.session.st	Configuration item for accessing metadata
ate.builder	Example configuration: Set this parameter to access DLI metadata.
	spark.sql.session.state.builder=org.apache.spark.sql.hiv e.DliLakeHouseBuilder

Parameter	Description
spark.sql.catalog.cl ass	Different data sources and metadata management systems
	Example configuration: Set this parameter to access DLI metadata.
	spark.sql.catalog.class=org.apache.spark.sql.hive.DliLak eHouseCatalog
spark.dli.metaAcce ss.enable	Enables or disables access to DLI metadata. spark.dli.metaAccess.enable=true

4. Click **connect**. When the **connect** button in the upper right corner changes to the queue name and the dot before the name turns green, the connection is successful. Then, you can execute the notebook job.

Figure 14-4 Notebook instance connected

notebook_general	þ
+ ^ ~ ©	í

5. Click **connect** to test the connection.

----End

Once the notebook instance is initialized, you can perform online data analysis on it. Instance initialization typically takes about 2 minutes.

When you run SQL statements in the notebook instance, a Spark job is started in DLI, and the results are displayed in the instance.

## Step 7: Configure the Notebook Instance to Access DLI Metadata

Before running a job, you need to configure the notebook instance to access DLI or LakeFormation metadata.

- (Optional) Configuring the Notebook Instance to Access DLI Metadata
- (Optional) Configuring the Notebook Instance to Access LakeFormation Metadata

## Step 8: Use the Notebook Instance to Write and Debug Code

After the notebook instance is connected to the DLI queue, you can edit and debug code in the **Notebook** area.

You can choose to submit a job using the notebook instance or through the **Spark Jobs** page of the DLI management console.

- For notebook-related operations, see JupyterLab Overview and Common Operations.
- For how to upload data to a notebook instance, see Uploading Files to JupyterLab.

• For how to download data from a notebook instance, see **Downloading a File from JupyterLab to a Local PC**.

## (Optional) Configuring the Notebook Instance to Access DLI Metadata

After connecting the notebook instance to DLI, you need to configure access to metadata if you plan to submit DLI jobs using the notebook instance. This section describes how to configure access to DLI metadata.

For how to configure the notebook instance to access LakeFormation metadata, see **(Optional) Configuring the Notebook Instance to Access LakeFormation Metadata**.

- 1. Specify a notebook image.
- 2. Create a custom agency to authorize DLI to use DLI metadata and OBS.

For how to create a custom agency, see **Creating a Custom DLI Agency**. Make sure the custom agency contains the following permissions:

Scenario	Agency Name	Use Case	Permission Policy
Allowing DLI to read and write data from and to OBS to transfer logs	Custom	For DLI Flink jobs, the permissions include downloading OBS objects, obtaining OBS/ GaussDB(DWS) data sources (foreign tables), transferring logs, using savepoints, and enabling checkpointing. For DLI Spark jobs, the permissions allow downloading OBS objects and reading/writing OBS foreign tables.	Permission Policies for Accessing and Using OBS
Allowing DLI to access DLI catalogs to retrieve metadata	Custom	DLI accesses catalogs to retrieve metadata.	Permission to Access DLI Catalog Metadata

 Table 14-3 DLI custom agency scenarios

## 3. Confirm access to DLI metadata.

- a. Log in to the ModelArts console and choose **Development Workspace** > **Notebook**.
- b. Create a notebook instance. When the instance is **Running**, click **Open** in the **Operation** column.
- c. On the displayed JupyterLab page, choose **File** > **New** > **Terminal**. The **Terminal** page appears.

## Figure 14-5 Accessing the Terminal page

$\sim$	File Edit View Run	Kernel Git Tabs Settings Help
	New	Console
	New Launcher	Ctrl+Shift+L 🔲 Notebook
Ø	Open from Path	s_ Terminal

d. Run the following commands to go to the Livy configuration directory and view the Spark configuration file:

## cd /home/ma-user/livy/conf/

## vi spark-defaults.conf

Ensure that the **spark.dli.user.catalogName=dli** configuration item exists. This item is used to access DLI metadata.

It is the default configuration item.

## Figure 14-6 Disabling default access to DLI metadata

File Edit	View	Run	Kerne	el Git	Tabs	Settings	Help					
+	Đ	±	G	${\rm I}$	Π.	Untitled13.ip	bynb	٠	S. Termin	al 1	×	:
Filter file	es by na	me		Q	spa	rk.yarn.	isPython	=true				
<b>I</b>					spa	irk.pyspa	rk.pytho	n=pythor	13			
Name	•		Last M	odified	<b>###</b> s	park.dli	.user.ca	talogNam	e=dli			
.mode	larts		6 da	ays ago		•		0				

- e. Use Notebook to edit a job.
  - For notebook-related operations, see JupyterLab Overview and Common Operations.
  - For how to upload data to a notebook instance, see Uploading Files to JupyterLab.
  - For how to download data from a notebook instance, see Downloading a File from JupyterLab to a Local PC.

## (Optional) Configuring the Notebook Instance to Access LakeFormation Metadata

After connecting the notebook instance to DLI, you need to configure access to metadata if you plan to submit DLI jobs using the notebook instance. This section describes how to configure access to LakeFormation metadata.

For how to configure the notebook instance to access DLI metadata, see **(Optional) Configuring the Notebook Instance to Access DLI Metadata**.

- 1. Connect DLI to LakeFormation.
  - a. For details, see **Connecting DLI to LakeFormation**.
- 2. Specify a notebook image.
- 3. Create a custom agency to authorize DLI to use LakeFormation metadata and OBS.

For how to create a custom agency, see **Creating a Custom DLI Agency**. Make sure the custom agency contains the following permissions:

Scenario	Agency Name	Use Case	Permission Policy
Allowing DLI to read and write data from and to OBS to transfer logs	Custom	For DLI Flink jobs, the permissions include downloading OBS objects, obtaining OBS/ GaussDB(DWS) data sources (foreign tables), transferring logs, using savepoints, and enabling checkpointing. For DLI Spark jobs, the permissions allow downloading OBS objects and reading/writing OBS foreign tables.	Permission Policies for Accessing and Using OBS
Allowing DLI to Custom access LakeFormation catalogs to retrieve metadata		DLI accesses LakeFormation catalogs to retrieve metadata.	Permission to Access LakeFormat ion Catalog Metadata

 Table 14-4 DLI custom agency scenarios

## 4. On the notebook instance page, set Spark parameters.

- a. Select the queue of the DLI notebook image, click **connect**, and set Spark parameters.
  - spark.sql.catalogImplementation=hive

spark.hadoop.hive-ext.dlcatalog.metastore.client.enable=true

spark.hadoop.hive-

ext.dlcatalog.metastore.session.client.class=com.huawei.cloud.dalf.lakecat.client.hiveclient.LakeCa tMetaStoreClient

spark.hadoop.lakecat.catalogname.default=lfcatalog // Specify the catalog to access.

spark.dli.job.agency.name=agencyForLakeformation // The agency must have the necessary

permissions on LakeFormation and OBS and must be delegated to DLI.

spark.driver.extraClassPath=/usr/share/extension/dli/spark-jar/lakeformation/\*

spark.executor.extraClassPath=/usr/share/extension/dli/spark-jar/lakeformation/\*

spark.sql.extensions=org.apache.spark.sql.hudi.HoodieSparkSessionExtension

spark.hadoop.hoodie.support.write.lock=org.apache.hudi.lakeformation.LakeCatMetastoreBasedL ockProvider

Table 1	4-5	Parameter	description
---------	-----	-----------	-------------

Parameter	Ma nda tory	Example Value	Configuration Scenario
spark.sql.catalogImplem entation	Yes	hive	Type of catalog used to store and manage metadata

Parameter	Ma nda tory	Example Value	Configuration Scenario			
spark.hadoop.hive- ext.dlcatalog.metastore. client.enable	Yes	true	Mandatory when LakeFormation metadata access is enabled			
spark.hadoop.hive- ext.dlcatalog.metastore. session.client.class	Yes	com.huawei.cl oud.dalf.lakec at.client.hivecl ient.LakeCat MetaStoreClie nt	Mandatory when LakeFormation metadata access is enabled			
spark.hadoop.lakecat.ca talogname.default	No	lfcatalog	Name of the LakeFormation data directory to access The default value is <b>hive</b> .			
spark.dli.job.agency.nam e	Yes	User-defined agency name	<ul> <li>User-defined agency name</li> <li>For how to create a custom agency, see Creating a Custom DLI Agency.</li> <li>For DLI metadata agency permissions, see Permission to Access LakeFormation Catalog Metadata.</li> </ul>			
spark.driver.extraClassPa th	Yes	/usr/share/ extension/dli/ spark-jar/ lakeformation /*	Loading of the LakeFormation dependency package			
spark.executor.extraClas sPath	Yes	/usr/share/ extension/dli/ spark-jar/ lakeformation /*	Loading of the LakeFormation dependency package			

Parameter	Ma nda tory	Example Value	Configuration Scenario
spark.sql.extensions	No	org.apache.sp ark.sql.hudi.H oodieSparkSe ssionExtensio n	Mandatory in Hudi scenarios
spark.hadoop.hoodie.su pport.write.lock	No	org.apache.hu di.lakeformati on.LakeCatMe tastoreBasedL ockProvider	Mandatory in Hudi scenarios

- 5. Disable the default access to DLI metadata and use LakeFormation metadata.
  - a. Log in to the ModelArts management console and choose **DevEnviron** > **Notebook**.
  - b. Create a notebook instance. When the instance is **Running**, click **Open** in the **Operation** column.
  - c. On the displayed JupyterLab page, choose File > New > Terminal. The Terminal page appears.

## Figure 14-7 Accessing the Terminal page

$\bowtie$	File	Edit	View	Run	Kernel	Git	Tabs	Settin	gs	Help
	N	lew						•	۶.,	Console
	N	lew Lau	incher			(	Ctrl+Shi	ft+L		Notebook
0	C	nen fro	om Path.						\$_	Terminal

d. Run the following commands to go to the Livy configuration directory and modify the Spark configuration file to disable the default access to DLI metadata:

## cd /home/ma-user/livy/conf/

## vi spark-defaults.conf

Use # to comment out **spark.dli.user.catalogName=dli** to disable the default access to DLI metadata.



e. Use Notebook to edit a job.

Run the spark.sql statement to access LakeFormation metadata and Hudi tables.

Figure 14-9 Accessing LakeFormation metadata

+ 10	± C 🕸		Untitle	d11.	ipynb			٠										
Filter files by na	ime Q	8	+	Ж	Ō	Ĉ	۲		C	**	Code		~ (	D	git			
<b>III</b> /		40 ms	[1]	spa	rk.sc	1("5	how	tabl	es")									
Name 🔺	Last Modified			Sta	rting	s Spa	irk a	ppli	catio	on								
modelarts 🖬	5 days ago	1		ID				YAI	RN Ap	plicat	tion ID	Kir	nd St	tate	Spark UI	Driver log	User	Current session?
Untitled Fol	4 days ago			0	fef7d	a65-2	79f-4	276-8	8dc-et	5952b	c965c6	pyspa	rk i	idle	Link	Link	None	*
📃 cwk.ipynb	11 hours ago			Spa	rkSes	sion	ave	ilab	le as	s 'sp	ark'.							
🖪 lftest.ipynb	4 days ago			Dat	aFran	ne[na	mesp	ace:	stra	ing,	table	Name:	strin	ng,	isTempor	ary: bool	ean]	
📕 Untitled.ipy	4 days ago	5.9	0	spa	rk.sc	1(5	how	tabl	es").	show	()							
📕 Untitled1.ip	2 days ago																	
🖪 Untitled10.i	11 hours ago		-	+		+-					+		+					
🔲 Untitled11.i	seconds ago			Ina	mespa	ace			tal	oleNa	me is	Tempor	ary					
📕 Untitled2.ip	2 days ago			i.	defau	lt				te	st	fa	lse					
🖪 Untitled3.ip	a day ago				defau	1t		h	udi_r	th1	bl	<b>0</b> a	lakefo	orme	tion table			
📕 Untitled4.ip	a day ago			i	defau	ilt		hudi	mor	tbl	rt	fa	lse					
📕 Untitled5.ip	18 hours ago			E	defau defau	lt  lt h	udi	hud	i_cow	non.	_1	fa fa	lse					

# **15** Submitting a Spark Jar Job Using Livy

## Introduction to DLI Livy

DLI Livy is an Apache Livy-based client tool used to submit Spark jobs to DLI.

## Preparations

- Create an elastic resource pool and create queues within it. When creating a queue, select **General-purpose**, which is the compute resources used to run Spark jobs. For details, see **Creating a Queue**.
- Prepare a Linux ECS for installing DLI Livy.
  - Enable ports 30000 to 32767 and port 8998 on the ECS. For details, see
     Adding a Security Group Rule.
  - Install JDK on the ECS. JDK 1.8 is recommended. Configure Java environment variable JAVA\_HOME.
  - View the ECS details to obtain its private IP address.
- Use an enhanced datasource connection to connect the DLI queue to the VPC where the Livy instance is located. For details, see Enhanced Datasource Connection.

## Step 1: Download and Install DLI Livy

## **NOTE**

The software package used in the following operations is **apache-livy-0.7.2.0107bin.tar.gz**. Replace it with the latest one.

- **Step 1 Download** the DLI Livy software package.
- **Step 2** Use WinSCP to upload the obtained software package to the prepared ECS directory.
- Step 3 Log in to ECS as user root and perform the following steps to install DLI Livy:
  - Run the following command to create an installation directory: mkdir Livy installation directory

For example, to create the **/opt/livy** directory, run the **mkdir /opt/livy** command. The following operations use the **/opt/livy** installation directory as an example. Replace it as required.

2. Run the following command to decompress the software package to the installation directory:

tar --extract --file apache-livy-0.7.2.0107-bin.tar.gz --directory /opt/livy -strip-components 1 --no-same-owner

- 3. Run the following commands to change the configuration file name:
  - cd /opt/livy/conf

mv livy-client.conf.template livy-client.conf

- mv livy.conf.template livy.conf
- mv livy-env.sh.template livy-env.sh
- mv log4j.properties.template log4j.properties
- mv spark-blacklist.conf.template spark-blacklist.conf
- touch spark-defaults.conf

----End

## Step 2: Modify the DLI Livy Configuration File

Step 1 Upload the specified DLI Livy JAR package to the OBS bucket directory.

- 1. Log in to OBS console and create a directory for storing the DLI Livy JAR package in the specified OBS bucket, for example: **obs://bucket/livy/jars/**.
- 2. Go to the installation directory of the ECS where the DLI-Livy tool has been installed in **Step 3.1**, obtain Livy JAR packages, and upload them to the OBS bucket directory created in **Step 1.1**:

For example, if the installation path is **/opt/livy**, the JAR packages you need to upload are as follows: /opt/livy/rsc-jars/livy-api-0.7.2.0107.jar

/opt/livy/rsc-jars/livy-api-0.7.2.0107.jar /opt/livy/rsc-jars/livy-rsc-0.7.2.0107.jar /opt/livy/repl\_2.11-jars/livy-repl\_2.11-0.7.2.0107.jar /opt/livy/repl\_2.11-jars/livy-repl\_2.11-0.7.2.0107.jar

- **Step 2** Modify the DLI Livy configuration file.
  - 1. Run the following command to modify the **/opt/livy/conf/livy-client.conf** configuration file:

## vi /opt/livy/conf/livy-client.conf

Add the following content to the file and modify the configuration items as required: # Set the private IP address of the ECS, which can be obtained by running the ifconfig command. livy.rsc.launcher.address = X.X.X # Set the ports enabled on the ECS. livy.rsc.launcher.port.range = 30000~32767

2. Run the following command to modify the **/opt/livy/conf/livy.conf** configuration file:

## vi /opt/livy/conf/livy.conf

Add the following content to the file and modify the configuration items as required:

```
livy.server.port = 8998
livy.spark.master = yarn
```

livy.server.contextLauncher.custom.class=org.apache.livy.rsc.DliContextLauncher livy.server.batch.custom.class=org.apache.livy.server.batch.DliBatchSession livy.server.interactive.custom.class=org.apache.livy.server.interactive.DliInteractiveSession livy.server.sparkApp.custom.class=org.apache.livy.utils.SparkDliApp

livy.server.recovery.mode = recovery livy.server.recovery.state-store = filesystem # Change the following file directory of DLI Livy as needed: livy.server.recovery.state-store.url = file:///opt/livy/store/

livy.server.session.timeout-check = true livy.server.session.timeout = 1800s livy.server.session.state-retain.sec = 1800s

livy.dli.spark.version = 2.3.2 livy.dli.spark.scala-version = 2.11

# Enter the OBS bucket path that stores the Livy JAR file. livy.repl.jars = obs://bucket/livy/jars/livy-core\_2.11-0.7.2.0107.jar, obs://bucket/livy/jars/livyrepl\_2.11-0.7.2.0107.jar livy.rsc.jars = obs://bucket/livy/jars/livy-api-0.7.2.0107.jar, obs://bucket/livy/jars/livyrsc-0.7.2.0107.jar

 Run the following command to modify the /opt/livy/conf/sparkdefaults.conf configuration file:

### vi /opt/livy/conf/spark-defaults.conf

Add the following content to the file. Set the parameters based on **Table 15-1**.

# The following parameters can be overwritten when a job is submitted.
spark.yarn.isPython=true
spark.pyspark.python=python3

# Enter the production environment URL of DLI. spark.dli.user.uiBaseAddress=https://console.huaweicloud.com/dli/web # Set the region where the queue is located. spark.dli.user.regionName=XXXX

# Set the DLI endpoint address. spark.dli.user.dliEndPoint=XXXX

# Enter the name of the created DLI queue. spark.dli.user.queueName=XXXX

# Set the project ID used for submitting a job. spark.dli.user.projectId=XXXX

Table 15-1	Mandatory	parameters	in spark-	-defaults.conf
------------	-----------	------------	-----------	----------------

Parameter	Description
spark.dli.user.r	Name of the region where the DLI queue is.
egionName	Obtain the region name from <b>Regions and Endpoints</b> .
spark.dli.user.d	Endpoint where the DLI queue is located.
liEndPoint	Obtain the endpoint from <b>Regions and Endpoints</b> .
spark.dli.user.q ueueName	Queue name.

Parameter	Description
spark.dli.user.a	User's AK/SK. The user must have Spark job permissions.
ccess.key	For details, see <b>Permissions Management</b> .
spark.dli.user.s ecret.key	For how to obtain the AK/SK, see <b>Obtaining the AK/SK</b> .
spark.dli.user.p	Project ID. Obtain it by referring to <b>Obtaining a Project</b>
rojectId	ID.

The following parameters are optional. Set them based on the parameter description and site requirements. For details about these parameters, see **Spark Configuration**.

Spark Job Parameter	Spark Batch Processing Parameter	Remarks
spark.dli.user.file	file	Not required for connecting to the notebook tool
spark.dli.user.class Name	class_name	Not required for connecting to the notebook tool
spark.dli.user.scTy pe	sc_type	Same as the native Livy configuration
spark.dli.user.args	args	Same as the native Livy configuration
spark.submit.pyFil es	python_files	Same as the native Livy configuration
spark.files	files	Same as the native Livy configuration
spark.dli.user.mod ules	modules	-
spark.dli.user.ima ge	image	Custom image used for submitting a job. This parameter is available for container clusters only and is not set by default.
spark.dli.user.auto Recovery	auto_recovery	-
spark.dli.user.max RetryTimes	max_retry_times	-

Table 15-2 Optional parameters in spark-defaults.conf

Spark Job Parameter	Spark Batch Processing Parameter	Remarks
spark.dli.user.cata logName	catalog_name	To access metadata, set this parameter to <b>dli</b> .

----End

## Step 3: Start DLI Livy

**Step 1** Run the following command to go to the DLI Livy installation directory:

Example: cd /opt/livy

Step 2 Run the following command to start DLI Livy:

## ./bin/livy-server start

----End

## Step 4: Submit a Spark Job to DLI Using DLI Livy

The following demonstrates how to submit a Spark job to DLI using DLI Livy and by running the curl command.

**Step 1** Upload the JAR file of the developed Spark job program to the OBS directory.

For example, upload **spark-examples\_2.11-XXXX.jar** to the **obs://bucket/path** directory.

## **NOTE**

To write the output data of a Spark Jar job to OBS, AK/SK is required for accessing OBS. To ensure the security of AK/SK data, you can use Data Encryption Workshop (DEW) and Cloud Secret Management Service (CSMS) for unified management of AK/SK, effectively avoiding sensitive information leakage and business risks caused by hard-coded or plaintext configuration of programs.

For details, see **Obtaining Temporary Credentials from a Spark Job's Agency for Accessing Other Cloud Services**.

**Step 2** Log in to the ECS server where DLI Livy is installed as user **root**.

Step 3 Run the curl command to submit a Spark job request to DLI using DLI Livy.

## **NOTE**

ECS\_IP indicates the private IP address of the ECS where DLI Livy is installed.

```
curl --location --request POST 'http://ECS_IP.8998/batches' \

--header 'Content-Type: application/json' \

--data '{

"driverMemory": "3G",

"driverCores": 1,

"executorMemory": "2G",

"executorCores": 1,

"numExecutorS": 1,

"args": [

"1000"
```

```
],

"file": "obs://bucket/path/spark-examples_2.11-XXXX.jar",

"className": "org.apache.spark.examples.SparkPi",

"conf": {

"spark.dynamicAllocation.minExecutors": 1,

"spark.executor.instances": 1,

"spark.dynamicAllocation.initialExecutors": 1,

"spark.dynamicAllocation.maxExecutors": 2

}
```

----End

## **16** Monitoring DLI Using Cloud Eye

## Description

This section describes metrics reported by DLI to Cloud Eye as well as their namespaces and dimensions. You can use the management console or APIs provided by Cloud Eye to query the metrics of the monitored object and alarms generated for DLI.

## Namespace

SYS.DLI

## Metric

Table 16-1 DLI metrics

Metric ID	Name	Descripti on	Valu e Ran ge	Unit	Con vers ion Rule	Monito red Object	Monitoring Period (Raw Data)
queue_cu _num	Queue CU Usage	Displays the number of CUs applied by the user queue	≥ 0	Count	N/A	Queues	5 minutes

Metric ID	Name	Descripti on	Valu e Ran ge	Unit	Con vers ion Rule	Monito red Object	Monitoring Period (Raw Data)
queue_jo b_launchi ng_num	Numbe r of Jobs Being Submitt ed	Displays the number of jobs in the Submitti ng state in the user queue.	≥ 0	Count	N/A	Queues	5 minutes
queue_jo b_running _num	Numbe r of Runnin g Jobs	Displays the number of running jobs in the user queue.	≥ 0	Count	N/A	Queues	5 minutes
queue_jo b_succeed _num	Numbe r of Finishe d Jobs	Displays the number of complete d jobs in the user queue.	≥ 0	Count	N/A	Queues	5 minutes
queue_jo b_failed_n um	Failed Jobs	Displays the number of failed jobs in the user queue.	≥ 0	Count	N/A	Queues	5 minutes
queue_jo b_cancell ed_num	Numbe r of Cancele d Jobs	Displays the number of canceled jobs in the user queue.	≥ 0	Count	N/A	Queues	5 minutes

Metric ID	Name	Descripti on	Valu e Ran ge	Unit	Con vers ion Rule	Monito red Object	Monitoring Period (Raw Data)
queue_all oc_cu_nu m	Allocat ed CUs (queue)	Displays the CU allocatio n for user queues.	≥ 0	Count	N/A	Queues	5 minutes
queue_mi n_cu_nu m	Minimu m CUs for Queue	Displays the minimu m number of CUs for a user queue.	≥ 0	Count	N/A	Queues	5 minutes
queue_m ax_cu_nu m	Maxim um CUs for Queue	Displays the maximu m number of CUs for a user queue.	≥ 0	Count	N/A	Queues	5 minutes
queue_pri ority	Queue Priority	Displays the priority of a user queue.	1– 100	N/A	N/A	Queues	5 minutes
queue_cp u_usage	Queue CPU Usage	Displays the CPU usage of user queues.	0- 100	%	N/A	Queues This metric applies only to queues in non- elastic resourc e pools.	5 minutes

Metric ID	Name	Descripti on	Valu e Ran ge	Unit	Con vers ion Rule	Monito red Object	Monitoring Period (Raw Data)
queue_dis k_usage	Queue Disk Usage	Displays the disk usage of user queues.	0– 100	%	N/A	Queues This metric applies only to queues in non- elastic resourc e pools.	5 minutes
queue_dis k_used	Max Disk Usage	Displays the maximu m disk usage of user queues.	0– 100	%	N/A	Queues This metric applies only to queues in non- elastic resourc e pools.	5 minutes
queue_m em_usage	Queue Memor y Usage	Displays the memory usage of user queues.	0– 100	%	N/A	Queues This metric applies only to queues in non- elastic resourc e pools.	5 minutes
queue_m em_used	Used Memor y	Displays the memory usage rate of the user queues.	≥ 0	МВ	N/A	Queues This metric applies only to queues in non- elastic resourc e pools.	5 minutes

Metric ID	Name	Descripti on	Valu e Ran ge	Unit	Con vers ion Rule	Monito red Object	Monitoring Period (Raw Data)
queue_jo b_launchi ng_max_d uration	Longest Job Submis sion	The longest submitte d job that is still in progress at the sampling time (includin g SQL, Flink, and Spark jobs).	≥ 0	Secon ds	N/A	Queues	5 minutes This metric is an instantaneou s sampling metric (non- continuous sampling), used to record the longest submitted jobs that are still in progress at the moment of sampling, specifically those in the <b>Submitting</b> or <b>Starting</b> state. It does not serve as a statistical metric for all jobs. Data statistics for historical jobs or completed jobs are not included. It is only applicable for monitoring the status of queues.

Metric ID	Name	Descripti on	Valu e Ran ge	Unit	Con vers ion Rule	Monito red Object	Monitoring Period (Raw Data)
queue_sql _job_runn ing_max_ duration	Longest SQL Job	The longest running SQL job that is still in progress at the sampling time.	≥ 0	Secon ds	N/A	Queues	5 minutes This metric is an instantaneou s sampling metric (non- continuous sampling), used to record the longest running SQL jobs that are still in progress at the moment of sampling, specifically those in the <b>Running</b> state. It does not serve as a statistical metric for all jobs. Data statistics for historical jobs or completed jobs are not included. It is only applicable for monitoring the status of queues.

Metric ID	Name	Descripti on	Valu e Ran ge	Unit	Con vers ion Rule	Monito red Object	Monitoring Period (Raw Data)
queue_sp ark_job_r unning_m ax_durati on	Longest Spark Job	The longest running Spark job that is still in progress at the sampling time.	≥ 0	Secon ds	N/A	Queues	5 minutes This metric is an instantaneou s sampling metric (non- continuous sampling), used to record the longest running Spark jobs that are still in progress at the moment of sampling, specifically those in the <b>Running</b> state. It does not serve as a statistical metric for all jobs. Data statistics for historical jobs or completed jobs are not included. It is only applicable for monitoring the status of queues.

Metric ID	Name	Descripti on	Valu e Ran ge	Unit	Con vers ion Rule	Monito red Object	Monitoring Period (Raw Data)
flink_read _records_ per_secon d	Flink Job Data Read Rate	Displays the data input rate of a Flink job for monitori ng and debuggin g.	≥ 0	recor d/s	N/A	Flink jobs	10 seconds
flink_writ e_records _per_seco nd	Flink Job Data Write Rate	Displays the data output rate of a Flink job for monitori ng and debuggin g.	≥ 0	recor d/s	N/A	Flink jobs	10 seconds
flink_read _records_t otal	Flink Job Total Data Read	Displays the total number of data inputs of a Flink job for monitori ng and debuggin g.	≥ 0	recor d/s	N/A	Flink jobs	10 seconds
flink_writ e_records _total	Flink Job Total Data Write	Displays the total number of output data records of a Flink job for monitori ng and debuggin g.	≥ 0	recor d/s	N/A	Flink jobs	10 seconds

Metric ID	Name	Descripti on	Valu e Ran ge	Unit	Con vers ion Rule	Monito red Object	Monitoring Period (Raw Data)
flink_read _bytes_pe r_second	Flink Job Byte Read Rate	Displays the number of input bytes per second of a Flink job.	≥ 0	byte/ s	102 4(IE C)	Flink jobs	10 seconds
flink_writ e_bytes_p er_second	Flink Job Byte Write Rate	Displays the number of output bytes per second of a Flink job.	≥ 0	byte/ s	102 4(IE C)	Flink jobs	10 seconds
flink_read _bytes_tot al	Flink Job Total Read Byte	Displays the total number of input bytes of a Flink job.	≥ 0	byte/ s	102 4(IE C)	Flink jobs	10 seconds
flink_writ e_bytes_t otal	Flink Job Total Write Byte	Displays the total number of output bytes of a Flink job.	≥ 0	byte/ s	102 4(IE C)	Flink jobs	10 seconds
flink_cpu_ usage	Flink Job CPU Usage	Displays the CPU usage of Flink jobs.	0– 100	%	N/A	Flink jobs	10 seconds
flink_me m_usage	Flink Job Memor y Usage	Displays the memory usage of Flink jobs.	0- 100	%	N/A	Flink jobs	10 seconds

Metric ID	Name	Descripti on	Valu e Ran ge	Unit	Con vers ion Rule	Monito red Object	Monitoring Period (Raw Data)
flink_max _op_laten cy	Flink Job Max Operat or Latency	Displays the maximu m operator delay of a Flink job. The unit is <b>ms</b> .	≥ 0	ms	N/A	Flink jobs	10 seconds
flink_max _op_back pressure_l evel	Flink Job Maxim um Operat or Backpre ssure	Displays the maximu m operator backpres sure value of a Flink job. A larger value indicates severer backpres sure. 0: OK 50: low 100: high	0- 100	N/A	N/A	Flink jobs	10 seconds
elastic_re source_po ol_cpu_us age	CPU Usage of Elastic Resourc e Pool	Displays the CPU usage of elastic resource pools.	0- 100	%	N/A	Elastic resourc e pools	5 minutes
elastic_re source_po ol_mem_ usage	Memor y Usage of Elastic Resourc e Pool	Displays the memory usage of elastic resource pools.	0- 100	%	N/A	Elastic resourc e pools	5 minutes

Metric ID	Name	Descripti on	Valu e Ran ge	Unit	Con vers ion Rule	Monito red Object	Monitoring Period (Raw Data)
elastic_re source_po ol_disk_us age	Disk Usage of Elastic Resourc e Pool	Displays the disk usage of elastic resource pools.	0- 100	%	N/A	Elastic resourc e pools	5 minutes
elastic_re source_po ol_disk_m ax_usage	Maxim um Disk Usage of Elastic Resourc e Pool	Displays the maximu m disk usage of elastic resource pools.	0– 100	%	N/A	Elastic resourc e pools	5 minutes
elastic_re source_po ol_cu_nu m	CU Usage of Elastic Resourc e Pool	Displays the CU usage of elastic resource pools.	≥ 0	Count	N/A	Elastic resourc e pools	5 minutes
elastic_re source_po ol_alloc_c u_num	Allocat ed CUs of Elastic Resourc e Pool	Displays the CU allocatio n of elastic resource pools.	≥ 0	Count	N/A	Elastic resourc e pools	5 minutes
elastic_re source_po ol_min_cu _num	Minimu m CUs of Elastic Resourc e Pool	Displays the minimu m number of CUs of elastic resource pools.	≥ 0	Count	N/A	Elastic resourc e pools	5 minutes

Metric ID	Name	Descripti on	Valu e Ran ge	Unit	Con vers ion Rule	Monito red Object	Monitoring Period (Raw Data)
elastic_re source_po ol_max_c u_num	Maxim um CUs of Elastic Resourc e Pool	Displays the maximu m number of CUs of elastic resource pools.	≥ 0	Count	N/A	Elastic resourc e pools	5 minutes

## Dimension

## Table 16-2 Dimension

Кеу	Value
queue_id	Queue
flink_job_id	Flink job

## Viewing DLI Monitoring Metrics on Cloud Eye

- 1. Search for Cloud Eye on the management console.
- In the navigation pane on the left of the Cloud Eye console, click Cloud Service Monitoring > Data Lake Insight.
- 3. Select a queue to view its information.

## **17** Using CTS to Audit DLI

With CTS, you can log operations related to DLI, making it easier to search, audit, and trace in the future.

Operation	Resource Type	Trace Name
Creating a database	database	createDatabase
Deleting a database	database	deleteDatabase
Changing the database owner	database	alterDatabaseOwner
Creating a table	table	createTable
Deleting tables	table	deleteTable
Exporting table data	table	exportData
Importing table data	table	importData
Changing the owner of a table	table	alterTableOwner
Creating a queue	queue	createQueue
Deleting a queue	queue	deleteQueue
Granting permissions to a queue	queue	queueAuthorize
Modifying the CIDR block of a queue	queue	replaceQueue
Restarting a queue	queue	queueActions
Scaling out/in a queue	queue	queueActions
Submitting a job (SQL)	queue	submitJob
Canceling a job (SQL)	jobs	cancelJob

Table 17-1 DLI operations that can be recorded by (
---

Operation	Resource Type	Trace Name
Granting DLI the permission to access OBS buckets	obs	authorizeObsBuckets- ForStream
Checking SQL syntax	jobs	checkSQL
Deleting a job	jobs	deleteStreamJob
Creating a Flink OpenSource SQL job	jobs	createStreamSqlJob
Updating a Flink OpenSource SQL job	jobs	updateStreamSqlJob
Deleting Flink jobs in batches	jobs	deleteStreamJobs
Stopping a Flink job	jobs	stopStreamJobs
Submitting a Flink job	jobs	submitStreamJobs
Creating a Flink Jar job	jobs	createStreamJarJob
Updating a Flink Jar job	jobs	updateStreamJarJob
Viewing Flink jobs	jobs	checkStreamJob
Importing a savepoint	jobs	dealSavepoint
Purchasing CUH packages	order	orderPackage
Granting data permissions	authorization	dataAuthorize
Granting permissions on other projects	authorization	projectDataAuthorize
Exporting query results	jobs	storeJobResult
Saving a SQL template	template	createTemplate
Updating a SQL template	template	updateTemplate
Deleting a SQL template	template	deleteTemplates
Creating a Flink template	template	createStreamTemplate
Updating a Flink template	template	updateStreamTemplate
Viewing Flink templates	template	checkStreamTemplate
Deleting a Flink template	template	deleteStreamTemplate
Creating a datasource authentication and uploading a certificate	datasourceauth	uploadAuthInfo

Operation	Resource Type	Trace Name
Updating datasource authentication information	datasourceauth	updateAuthInfo
Deleting datasource authentication information	datasourceauth	deleteAuthInfo
Uploading a resource package	resource	uploadResources
Deleting a resource package	resource	deleteResource
Creating an enhanced datasource connection	edsconnection	createConnection
Deleting an enhanced datasource connection	edsconnection	deleteConnection
Creating a basic datasource connection	edsconnection	createConnection
Deleting a basic datasource connection	edsconnection	deleteConnection
Binding a queue	edsconnection	associateQueueToCon- nection
Unbinding a queue	edsconnection	disassociateQueueTo- Connection
Modifying host information	edsconnection	updateHostInfo
Adding a route	edsconnection	addRoute
Deleting a route	edsconnection	deleteRoute
Creating a batch processing job	jobs	createBatch
Canceling a batch processing job	jobs	cancelBatch
Creating a global variable	variable	createGlobalVariable
Deleting a global variable	variable	deleteGlobalVariable
Modifying a global variable	variable	updateGlobalVariable

For how to enable CTS and view trace details, see **Cloud Trace Service Getting Started**.

For details about key fields in the CTS trace structure, see **Trace StructureTrace Structure** in the *Cloud Trace Service User Guide*.
# **18** Permissions Management

# 18.1 Overview

DLI has a comprehensive permission control mechanism and supports fine-grained authentication through Identity and Access Management (IAM). You can create policies in IAM to manage DLI permissions. You can use both the DLI's permission control mechanism and the IAM service for permission management.

# **Application Scenarios of IAM Authentication**

When using DLI on the public cloud, enterprise users need to manage DLI resources (queues) used by employees in different departments, including creating, deleting, using, and isolating resources. In addition, data of different departments needs to be managed, including data isolation and sharing.

DLI uses IAM for refined enterprise-level multi-tenant management. IAM provides identity authentication, permissions management, and access control, helping you securely access to your public cloud resources.

With IAM, you can use your Huawei Cloud account to create IAM users for your employees, and assign permissions to the users to control their access to specific resource types. For example, some software developers in your enterprise need to use DLI resources but must not delete them or perform any high-risk operations. To achieve this result, you can create IAM users for the software developers and grant them only the permissions required for using DLI resources.

## **NOTE**

For a new user, you need to log in for the system to record the metadata before using DLI.

IAM is free to use, and you only need to pay for the resources in your account. For more information about IAM, see IAM Service Overview.

If your Huawei Cloud account does not need individual IAM users for permissions management, skip over this section.

# **DLI System Permissions**

 Table 18-1 lists all the system-defined roles and policies supported by DLI.

Type: There are roles and policies.

- Roles: A type of coarse-grained authorization mechanism that defines permissions related to user responsibilities. Only a limited number of service-level roles are available. When using roles to grant permissions, you also need to assign other roles on which the permissions depend. However, roles are not an ideal choice for fine-grained authorization and secure access control.
- Policies: A type of fine-grained authorization mechanism that defines permissions required to perform operations on specific cloud resources under certain conditions. This mechanism allows for more flexible policy-based authorization, meeting requirements for secure access control. For example, you can grant DLI users only the permissions for managing a certain type of ECSs.

For details about the system policies you need to perform common SQL operations, see **Common Operations Supported by DLI System Policy**.

Role/Policy Name	Description	Category	Dependency
DLI FullAccess	Full permissions for DLI.	System- defined policy	This role depends on other roles in the same project.
	policy		<ul> <li>Creating a datasource connection: VPC ReadOnlyAccess</li> </ul>
			<ul> <li>Creating yearly/ monthly resources: BSS Administrator</li> </ul>
			<ul> <li>Creating a tag: TMS FullAccess and EPS EPS FullAccess</li> </ul>
			<ul> <li>Using OBS for storage: OBS OperateAccess</li> </ul>
			<ul> <li>Creating an agency:</li> <li>Security</li> <li>Administrator</li> </ul>

 Table 18-1
 DLI system permissions

Role/Policy Name	Description	Category	Dependency
DLI ReadOnlyAcce ss	Read-only permissions for DLI. With read-only permissions, you can use DLI resources and perform operations that do not require fine-grained permissions. For example, create global variables, create packages and package groups, submit jobs to the default queue, create tables in the default database, create datasource connections, and delete datasource connections.	System- defined policy	None
Tenant Administrator	<ul> <li>Tenant administrator</li> <li>Job execution permissions for DLI resources. After a database or a queue is created, the user can use the ACL to assign rights to other users.</li> <li>Scope: project-level service</li> </ul>	System- defined role	None
DLI Service Administrator	<ul> <li>DLI administrator.</li> <li>Job execution permissions for DLI resources. After a database or a queue is created, the user can use the ACL to assign rights to other users.</li> <li>Scope: project-level service</li> </ul>	System- defined role	None

#### For details, see Creating an IAM User and Granting Permissions, How Do I Create an IAM user? and How Do I Modify a User Policy?

## **DLI Permission Types**

**Table 18-2** lists the DLI service permissions. For details about the resources that can be controlled by DLI, see **Table 18-7**.

Permission Type	Subtype	Console Operations	SQL Syntax
Queue Permissions	Queue management permissions	For details, see Queue Permission Management.	None
	Queue usage permission		
Data Permissions	Database permissions	For details, see Configuring	For details, see <b>Data</b> Permissions List.
	Table permissions	Database Permissions on the DLI Console and	
	Column permissions	Configuring Table Permissions on the DLI Console.	
Job Permissions	Flink job permissions	For details, see Configuring Flink Job Permissions.	None
Package Permissions	Package group permissions	For details, see Configuring DLI	None
	Package permissions	Package Permissions.	
Datasource Connection Permissions	Datasource connection permissions	For details, see Datasource Authentication Permission Management.	None

Table	18-2	DLI	permission	types
-------	------	-----	------------	-------

# Examples

An Internet company mainly provides game and music services. DLI is used to analyze user behaviors and assist decision making.

As shown in Figure 18-1, the Leader of the Basic Platform Team has applied for a Tenant Administrator account to manage and use public cloud services. The Leader of the Basic Platform Team creates a subaccount with the DLI Service Administrator permission to manage and use DLI, as the Big Data Platform **Team** requires DLI for data analysis. The **Leader of the Basic Platform Team** creates a **Queue A** and assigns it to **Data Engineer A** to analyze the gaming data. A **Queue B** is also assigned to **Data Engineer B** to analyze the music data. Besides granting the queue usage permission, the **Leader of the Basic Platform Team** grants data (except the database) management and usage permissions to the two engineers.





The **Data Engineer A** creates a table named **gameTable** for storing game prop data and a table named **userTable** for storing game user data. The music service is a new service. To explore potential music users among existing game users, the **Data Engineer A** assigns the query permission on the **userTable** to the **Data Engineer B**. In addition, **Data Engineer B** creates a table named **musicTable** for storing music copyrights information.

Table 18-3 describes the queue and data permissions of Data Engineer A and Data Engineer B.

User	Data Engineer A (game data analysis)	Data Engineer B (music data analysis)
Queues	Queue A (queue usage permission)	Queue B (queue usage permission)
Data (Table)	gameTable (table management and usage permission)	musicTable (table management and usage permission)
	userTable (table management and usage permission)	userTable (table query permission)

Table	18-3	Permission	description
iubic	10.5	1 CITII J JIOT	acochption

D NOTE

The queue usage permission includes job submitting and terminating permissions.

# **18.2 Creating a Custom Policy**

Custom policies can be created as a supplement to the system policies of DLI. You can add actions to custom policies. For the actions supported for custom policies, see .

You can create custom policies in either of the following two ways:

- Visual editor: Select cloud services, actions, resources, and request conditions without the need to know policy syntax.
- JSON: Create a policy in the JSON format from scratch or based on an existing policy.

For details, see **Creating a Custom Policy**. This section describes common DLI custom policies.

## **Policy Field Description**

The following example assumes that the authorized user has the permission to create tables in all databases in all regions:

```
Version": "1.1",
    "Statement": [
    {
        "Effect": "Allow",
        "Action": [
          "dli:database:createTable"
    ],
        "Resource": [
          "dli:*:*:database:*"
    ]
    }
]
```

Version

**1.1** indicates a fine-grained permission policy that defines permissions required to perform operations on specific cloud resources under certain conditions.

Effect

The value can be **Allow** and **Deny**. If both **Allow** and **Deny** are found in statements, the **Deny** overrides the **Allow**.

Action

Specific operation on a resource. A maximum of 100 actions are allowed, as shown in **Figure 18-2**.

#### Figure 18-2 DLI actions

Policy View	v	isual editor	JSON								
* Policy Content		Allow		C Data Lake Insight	G Actions: 0		C AI		C (Optional) Add request condition	Ē	Ū
		Select all	Enter a keyword.					Q			
		∩ □ Read	Only 14 in total								
		diitable     Query a	a showPartitions Il partitions.	Cuery t	le showCreateTable the table creation statement.	Cuery job deta	lis		dli:group:getGroup Query resource package group details.		
		Cuery to	e:showTableProperties able configurations.	Cuery t	le describeTable the table structure.	<ul> <li>dii:database.c</li> <li>Query databas</li> </ul>	lisplayDatabase es.		dli database showFunctions Query a function.		
		Gitable Select t	a select able	Cuery a	abase displayAlDatabases all databases.	dEcolumn sel     Select column	ect		dit resource getResource Query resource package details.		
		dii data     Query a	ibase displayAllTables il tables.	Cuery a	le showSegments a data segment.						
		Y 🗌 Reed	Write 50 in total								
		Y 🗌 Listor	nly 4 in total								
		Y 🗌 Permi	issions 17 in total								

#### D NOTE

- The format is *Service name*:*Resource type*:*Action*, for example, **dli:queue:submit\_job**.
- Service name: product name, such as dli, evs, and vpc. Only lowercase letters are allowed. Resource types and operations are not case-sensitive. You can use an asterisk (\*) to represent all operations.
- *Resource type*: For details, see **Table 18-7**.
- *Action*: action registered in IAM.
- Condition

Conditions determine when a policy takes effect. A condition consists of a condition key and operator.

A condition key is a key in the **Condition** element of a statement. There are global and service-level condition keys.

- Global condition keys (prefixed with g:) apply to all actions. For details, see condition key description in Policy Syntax.
- Service-level condition keys apply only to operations of the specific service.

An operator is used together with a condition key to form a complete condition statement. For details, see **Table 18-4**.

IAM provides a set of DLI predefined condition keys. The following table lists the predefined condition keys of DLI.

Condition Key	Ту ре	Operator	Description
g:CurrentTime	Glo bal	Date and time	Time when an authentication request is received <b>NOTE</b> The time is expressed in the format defined by <b>ISO 8601</b> , for example, <b>2012-11-11T23:59:59Z</b> .
g:MFAPresent	Glo bal	Boolean	Whether multi-factor authentication is used during user login
g:UserId	Glo bal	String	ID of the current login user
g:UserName	Glo bal	String	Current login user
g:ProjectName	Glo bal	String	Project that you have logged in to
g:DomainName	Glo bal	String	Domain that you have logged in to
g:ResourceTag	Glo bal	StringEquals	Resource tag value.

Table 18-4 DLI request conditions

Resource

The format is *Service name*.*Region*.*Domain ID*.*Resource type*.*Resource path*. The wildcard (\*) indicates all options. For details about the resource types and path, see **Table 18-7**.

Example:

dli:\*:\*:queue:\* indicates all queues.

# Creating a Custom Policy

You can set actions and resources of different levels based on scenarios.

1. Define an action.

The format is *Service name:Resource type:Action*. You can use wildcards \*. Example:

Table 18-5 Action

Action	Description
dli:queue:submit_job	Submission operations on a DLI queue
dli:queue:*	All operations on a DLI queue
dli:*:*	All operations on all DLI resource types

For more information about the relationship between operations and system permissions, see **Common Operations Supported by DLI System Policy**.

2. Define a resource.

The format is *Service name*.*Region*.*Domain ID*.**Resource type**:**Resource path**. The wildcard (\*) indicates all resources. The five fields can be flexibly set. Different levels of permission control can be set for resource paths based on scenario requirements. If you need to set all resources of the service, you do not need to specify this field. For details about the definition of Resource, see **Table 18-6**. For details about the resource types and paths in Resource, see **Table 18-7**.

Table	18-6	Resource
-------	------	----------

Resource	Description
DLI:*:*:table:databases.dbname.t ables.*	DLI, any region, any account ID, all table resources of database <b>dbname</b>
DLI:*:*:database:databases.dbna me	DLI, any region, any account ID, resource of database <b>dbname</b>

Resource	Description
DLI:*:*:queue:queues.*	DLI, any region, any account ID, any queue resource
DLI:*:*:jobs:jobs.flink.1	DLI, any region, any account ID, Flink job whose ID is 1

Table 18-7 D	I resources a	nd their paths
--------------	---------------	----------------

Resource Type	Resource Names	Path
elasticresou rcepool	DLI elastic resource pool	elasticresourcepools.name
queue	DLI queue	queues.queuename
database	DLI database	databases.dbname
table	DLI table	databases.dbname.tables.tbname
column	DLI column	databases.dbname.tables.tbname.columns.c olname
jobs	DLI Flink job	jobs.flink.jobid
resource	DLI package	resources.resourcename
group	DLI package group	groups.groupname
datasource auth	DLI cross- source authentication information	datasourceauth.name
edsconnect ions	Enhanced datasource connection	edsconnections. Connection ID
variable	DLI global variable	variables.name
sqldefendr ule	SQL inspector rule	sqldefendes.*

- Specific resources:

## Figure 18-3 Specific resources

Scope	Global serv	ices Project-level :	services					
	Project-level ser	vices, such as ECS and VPC, i	can be deployed and accessed in specific regions.					
Policy View	Visual edito	r ISON						
* Policy Content	Select Existing P	ilicy/Role						
	~ •	llow 🔷 🥥	Data Lake Insight	🛛 S dlidatabase:createTable	Specific resources	Optional) Add request condition		00
	Resource	s 🖲 Specific 🛛 🔿 All						
	database	Any Specify a	isource path					
		DUrmdatabasedatabase	es.dbname.tables.tbname.columns.colname				27	Ì
		Add Resource Path						
	A Add December	iner.						

- All resources: all resources of the service

```
Figure 18-4 All resources
```



3. Combine all the preceding fields into a JSON file to form a complete policy. You can set multiple actions and resources. You can also create a policy on the visualized page provided by IAM. For example:

Create a policy that grants users the permission to create and delete databases, submit jobs for any queue, and delete tables under any account ID in any region of DLI.



# **Example Custom Policies**

• Example 1: Allow policies

{

}

- Allow users to create tables across all databases in all regions:

```
"Version": "1.1",
"Statement": [
    {
        "Effect": "Allow",
        "Action": [
            "dli:database:createTable"
        ],
        "Resource": [
            "dli:*:*:database:*"
        ]
    }
]
```

 Allow users to query column col in the table tb of the database db in the region where the user is located:

```
"Version": "1.1",
"Statement": [
{
"Effect": "Allow",
"Action": [
"dli:column:select"
],
```

```
"Resource": [
     "dli:xxx:*:column:databases.db.tables.tb.columns.col"
    ]
    }
]
```

• Example 2: Deny policies

A deny policy must be used together with other policies. That is, a user can set a deny policy only after being assigned some operation permissions. Otherwise, the deny policy does not take effect.

If the permissions assigned to a user contain both Allow and Deny actions, the Deny actions take precedence over the Allow actions.

 Deny users to create or delete databases, submit jobs (except the default queue), or delete tables.

```
"Version": "1.1",
   "Statement": [
     {
        "Effect": "Deny",
"Action": [
           "dli:database:createDatabase",
           "dli:database:dropDatabase",
           "dli:queue:submitJob",
           "dli:table:dropTable"
        ],
"Resource": [
           "dli:*:*:database:*",
           "dli:*:*:queue:*",
           "dli:*:*:table:*"
        1
     }
  ]
}
Deny users to submit jobs in the demo queue.
   "Version": "1.1",
   "Statement": [
     {
         "Effect": "Deny",
         "Action": [
            "dli:queue:submitJob"
        1,
         "Resource": [
            "dli:*:*:queue:queues.demo"
        1
     }
  ]
```

• Example 3: Tag authentication. You need to specify an action and bind it to a condition, and specify the key and value of **g:ResourceTag**.

**Condition g: ResourceTag** indicates a resource with the *key=value* tag. Only the operation on this resource contained in the policy action list is allowed.

The key is case insensitive. Fuzzy match for the value is not supported.

```
"Version": "1.1",
"Statement": [
{
"Effect": "Allow",
"Action": [
"dli:database:dropDatabase",
"dli:table:select",
"dli:database:createTable",
```

{

```
"dli:table:dropTable"
],
"Condition": {
"StringEquals": {
"g:ResourceTag/key": [
"value"
]
}
}
}
```

# **18.3 DLI Resources**

}

A resource is an object that exists within a service. You can select DLI resources by specifying their paths.

Resource Type	Resource Names	Path
elasticresou rcepool	DLI elastic resource pool	elasticresourcepools.name
queue	DLI queue	queues.queuename
database	DLI database	databases.dbname
table	DLI table	databases.dbname.tables.tbname
column	DLI column	databases.dbname.tables.tbname.columns.coln ame
jobs	DLI Flink job	jobs.flink.jobid
resource	DLI package	resources.resourcename
group	DLI package group	groups.groupname
datasourcea uth	DLI cross-source authentication information	datasourceauth.name
edsconnecti ons	Enhanced datasource connection	edsconnections. Connection ID
variable	DLI global variable	variables.name
sqldefendrul e	SQL inspector rule	sqldefendes.*

Table 18-8 DLI resources and their paths

# **18.4 DLI Request Conditions**

Request conditions are useful in determining when a custom policy takes effect. A request condition consists of a condition key and operator. Condition keys are either global or service-level and are used in the Condition element of a policy statement. **Global condition keys** (starting with **g**:) are available for operations of all services, while service-level condition keys (starting with a service name such as **dli**) are available only for operations of a specific service. An operator is used together with a condition key to form a complete condition statement.

IAM provides a set of DLI predefined condition keys. The following table lists the predefined condition keys of *DLI*.

Condition Key	Тур е	Operator	Description
g:CurrentTime	Glo bal	Date and time	Time when an authentication request is received <b>NOTE</b> The time is expressed in the format defined by <b>ISO 8601</b> , for example, <b>2012-11-11T23:59:59Z</b> .
g:MFAPresent	Glo bal	Boolean	Whether multi-factor authentication is used during user login
g:UserId	Glo bal	String	ID of the current login user
g:UserName	Glo bal	String	Current login user
g:ProjectName	Glo bal	String	Project that you have logged in to
g:DomainName	Glo bal	String	Domain that you have logged in to
g:ResourceTag	Glo bal	StringEquals	Resource tag value.

Table 18-9 DLI request conditions

# **18.5 Common Operations Supported by DLI System Policy**

**Table 18-10** lists the common operations supported by each system policy of DLI. Choose proper system policies according to this table. For details about the SQL statement permission matrix in DLI in terms of permissions on databases, tables, and roles, see **Data Permission List**.

Res our ce	Operation	Description	DLI FullAcces s	DLI ReadOnl yAccess	Tenant Administ rator	DLI Service Admini strator
Qu eue	DROP_QUE UE	Deleting a Queue	$\checkmark$	×	$\checkmark$	$\checkmark$
	SUBMIT_JO B	Submitting a job	$\checkmark$	×	$\checkmark$	$\checkmark$
	CANCEL_JO B	Terminating a Job	~	×	$\checkmark$	$\checkmark$
	RESTART	Restarting a queue	$\checkmark$	×	$\checkmark$	$\checkmark$
	GRANT_PRI VILEGE	Granting permissions to a queue	$\checkmark$	×	$\checkmark$	$\checkmark$
	REVOKE_PRI VILEGE	Revoking permissions to a queue	$\checkmark$	×	√	√
	SHOW_PRIV ILEGES	Viewing the queue permissions of other users	$\checkmark$	×	$\checkmark$	~
Dat aba	DROP_DATA BASE	Deleting a database	$\checkmark$	×	$\checkmark$	$\checkmark$
se	CREATE_TAB LE	Creating a table	$\checkmark$	×	$\checkmark$	$\checkmark$
	CREATE_VIE W	Creating a view	$\checkmark$	×	$\checkmark$	$\checkmark$
	EXPLAIN	Explaining the SQL statement as an execution plan	√	×	$\checkmark$	$\checkmark$
	CREATE_RO LE	Creating a role	$\checkmark$	×	$\checkmark$	$\checkmark$
	DROP_ROLE	Deleting a role	$\checkmark$	×	$\checkmark$	$\checkmark$
	SHOW_ROL ES	Displaying a role	$\checkmark$	×	$\checkmark$	$\checkmark$

 Table 18-10 Common operations supported by each system permission

Res our ce	Operation	Description	DLI FullAcces s	DLI ReadOnl yAccess	Tenant Administ rator	DLI Service Admini strator
	GRANT_ROL E	Binding a role	$\checkmark$	×	$\checkmark$	$\checkmark$
	REVOKE_RO LE	Unbinding a role	$\checkmark$	×	$\checkmark$	$\checkmark$
	SHOW_USE RS	Displaying the binding relationships between all roles and users	$\checkmark$	×	$\checkmark$	~
	GRANT_PRI VILEGE	Granting permissions to the database	$\checkmark$	×	$\checkmark$	$\checkmark$
	REVOKE_PRI VILEGE	Revoking permissions to the database	$\checkmark$	×	$\checkmark$	$\checkmark$
	SHOW_PRIV ILEGES	Viewing database permissions of other users	$\checkmark$	×	$\checkmark$	~
	DISPLAY_AL L_TABLES	Displaying tables in a database	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
	DISPLAY_DA TABASE	Displaying databases	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
	CREATE_FU NCTION	Creating a function	$\checkmark$	×	$\checkmark$	$\checkmark$
	DROP_FUN CTION	Deleting a function	$\checkmark$	×	$\checkmark$	$\checkmark$
	SHOW_FUN CTIONS	Displaying all functions	$\checkmark$	×	$\checkmark$	$\checkmark$
	DESCRIBE_F UNCTION	Displaying function details	$\checkmark$	×	$\checkmark$	$\checkmark$
Tab le	DROP_TABL E	Deleting tables	$\checkmark$	×	$\checkmark$	$\checkmark$

Res our ce	Operation	Description	DLI FullAcces s	DLI ReadOnl yAccess	Tenant Administ rator	DLI Service Admini strator
	SELECT	Querying tables	$\checkmark$	×	$\checkmark$	$\checkmark$
	INSERT_INT O_TABLE	Inserting table data	$\checkmark$	×	$\checkmark$	$\checkmark$
	ALTER_TABL E_ADD_COL UMNS	Adding a column	$\checkmark$	×	$\checkmark$	~
	INSERT_OVE RWRITE_TA BLE	Overwriting a table	$\checkmark$	×	$\checkmark$	$\checkmark$
	ALTER_TABL E_RENAME	Renaming a table	$\checkmark$	×	$\checkmark$	$\checkmark$
	ALTER_TABL E_ADD_PAR TITION	Adding partitions to the partition table	$\checkmark$	×	$\checkmark$	~
	ALTER_TABL E_RENAME_ PARTITION	Renaming a table partition	$\checkmark$	×	$\checkmark$	~
	ALTER_TABL E_DROP_PA RTITION	Deleting partitions from a partition table	√	×	√	~
	SHOW_PAR TITIONS	Displaying all partitions	$\checkmark$	×	$\checkmark$	$\checkmark$
	ALTER_TABL E_RECOVER _PARTITION	Restoring table partitions	$\checkmark$	×	$\checkmark$	~
	ALTER_TABL E_SET_LOCA TION	Setting the partition path	$\checkmark$	×	$\checkmark$	~
	GRANT_PRI VILEGE	Granting permissions to the table	$\checkmark$	×	$\checkmark$	$\checkmark$
	REVOKE_PRI VILEGE	Revoking permissions to the table	√	×	√	√

Res our ce	Operation	Description	DLI FullAcces s	DLI ReadOnl yAccess	Tenant Administ rator	DLI Service Admini strator
	SHOW_PRIV ILEGES	Viewing table permissions of other users	$\checkmark$	×	√	~
	DISPLAY_TA BLE	Displaying a table	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
	DESCRIBE_T ABLE	Displaying table information	$\checkmark$	×	$\checkmark$	$\checkmark$
Elas tic reso urc	DROP	Deleting an elastic resource pool	$\checkmark$	×	$\checkmark$	$\checkmark$
e poo l	RESOURCE_ MANAGEME NT	Managing an elastic resource pool	$\checkmark$	×	$\checkmark$	~
	SCALE	Scaling an elastic resource pool	$\checkmark$	×	$\checkmark$	√
	UPDATE	Updating an elastic resource pool	$\checkmark$	×	$\checkmark$	$\checkmark$
	CREATE	Creating an elastic resource pool	$\checkmark$	×	$\checkmark$	$\checkmark$
	SHOW_PRIV ILEGES	Viewing elastic resource pool permissions of other users	$\checkmark$	×	$\checkmark$	$\checkmark$

Res our ce	Operation	Description	DLI FullAcces s	DLI ReadOnl yAccess	Tenant Administ rator	DLI Service Admini strator
	GRANT_PRI VILEGE	Granting elastic resource pool permissions	$\checkmark$	×	$\checkmark$	$\checkmark$
	REVOKE_PRI VILEGE	Retrieving elastic resource pool permissions	$\checkmark$	×	$\checkmark$	~
Enh anc ed dat aso	BIND_QUEU E	Binding an enhanced datasource connection to a queue	×	×	×	×
urc e con nec tion		It is only used to grant permissions across projects.				

# **19** Common DLI Management Operations

# 19.1 Enhancing the Job Runtime Environment Using a Custom Image

# Scenario

To enhance the functions and performance of Spark and Flink jobs, you can create custom images by downloading the base images provided by DLI and adding dependencies (files, JAR files, or software) and private capabilities required for job execution. This changes the container runtime environment for the jobs.

For example, you can add a Python package or C library related to machine learning to a custom image to help you extend functions.

#### **NOTE**

To use the custom image function, you need to have basic knowledge of **Docker**.

# **Notes and Constraints**

- The base images provided by DLI must be used to create custom images.
- You cannot modify the DLI components and directories in the base images.
- Only Spark Jar and Flink Jar jobs are supported.

# **Use Process**



- 1. Obtain DLI base images.
- 2. Use Dockerfile to pack dependencies (files, JAR files, or software) required for job execution into the base image to create a custom image.
- 3. Publish the custom image to SoftWare Repository for Container (SWR).
- 4. On the DLI job editing page, select the created image and run the job.
- 5. Check the job execution status.

## **Obtaining DLI Base Images**

Select the base image of the same type as the architecture of the queue.

For the CPU architecture type of a queue, see **Viewing Basic Information About** a **Queue**.

Image Type	Architecture	URL
General image	x86	swr.ap- southeast-3.myhuaweicloud.com/dli- public/spark_general- x86_64:3.3.1-2.3.8.1120250109929356 803819072.202501141605
General image	Arm	swr.ap- southeast-3.myhuaweicloud.com/dli- public/spark_general- aarch64:3.3.1-2.3.8.1120250109929356 803819072.202501141605
Notebook image	x86	swr.ap- southeast-3.myhuaweicloud.com/dli- public/spark_notebook- x86_64:3.3.1-2.3.8.1120250109929356 803819072.202501141605
Notebook image	Arm	swr.ap- southeast-3.myhuaweicloud.com/dli- public/spark_notebook- aarch64:3.3.1-2.3.8.1120250109929356 803819072.202501141605

Table 19-1 Obtaining DLI base images

# Creating a Custom Image

The following describes how to package TensorFlow into an image to generate a custom image with TensorFlow installed. Then, you can use the image to run jobs in DLI.

**Step 1** Prepare the container environment.

For details, see "Step 1: Install the Container Engine" in **Uploading an Image Through a Container Engine Client**. **Step 2** Log in to the prepared container environment as user **root** and run a command to obtain the base image.

In this example, the Spark base image is used and downloaded to the container image environment in **Step 1** by running the following command:

docker pull Address for downloading the base image

For details about the address, see **Obtaining DLI Base Images**.

- Step 3 Access SWR.
  - 1. Log in to the SWR management console.
  - 2. In the navigation pane on the left, choose **Dashboard** and click **Generate**

**Login Command** in the upper right corner. On the displayed page, click  $\square$  to copy the login command.

- 3. Run the login command on the VM where the container engine is installed.
- Step 4 Create an organization. If an organization has been created, skip this step.
  - 1. Log in to the SWR management console.
  - 2. In the navigation pane on the left, choose **Organization Management**. On the displayed page, click **Create Organization** in the upper right corner.
  - 3. Enter the organization name and click OK.
- **Step 5** Write a Dockerfile.

#### vi Dockerfile

Pack TensorFlow into the image as follows:

```
ARG BASE_IMG=swr.xxx/dli-public/spark_general-
x86_64:3.3.1-2.3.8.1120250109929356803819072.202501141605//Replace xxx with the URL of the base
image.
FROM ${BASE_IMG} as builder
USER omm //Run this command as user omm.
```

RUN set -ex && \ mkdir -p /home/omm/.pip && \ pip3 install tensorflow==2.4.0 Copy the content to the base image. USER omm

The following steps are included:

- 1. Set the available repository address of pip.
- 2. Use pip3 to install the TensorFlow algorithm package.
- 3. Copy the content in the temporary image builder where the algorithm package is installed to the base image (this step is to reduce the image size) to generate the final custom image.
- **Step 6** Use Dockerfile to generate a custom image.

Format of the image packaging command: docker build -t [*Custom organization name*]/[*Custom image name*]: [*Image version*] --build-arg BASE\_IMG= [*DLI base image path*] -f Dockerfile .

The DLI base image path is the image path in **Obtaining DLI Base Images**.

The following is an example:

docker build -t mydli/spark:2.4 --build-arg BASE\_IMG=swr.xxx/dli-public/spark\_generalx86\_64:3.3.1-2.3.8.1120250109929356803819072.202501141605 -f Dockerfile .

**Step 7** Add a tag to the custom image.

**docker tag** [Organization name]/[Image name]:[Image version][Image repository address]/[Organization name]/[Image name:version] in **Step 6** 

The following is an example: docker tag mydli/spark:2.4 swr.xxx/testdli0617/spark:2.4.5.tensorflow

**Step 8** Upload the custom image.

**docker push** [Image repository address]/[Organization name]/[Image name:Version]

Set [Image repository address]/[Organization name]/[Image name:Version] the same as those in **Step 7**.

The following is an example: docker push swr.xxx/testdli0617/spark:2.4.5.tensorflow

- Step 9 When submitting a Spark or Flink JAR job in DLI, select a custom image.
  - Open the Spark job or Flink job editing page on the management console, select the uploaded and shared image from the custom image list, and run the job.

If you select a non- shared image, the system displays a message indicating that the image is not authorized. You can use the image only after it is authorized. Click **Authorize** as prompted. Set other job execution parameters and execute the job.

fill form Write API					
Select a	a Queue				
* Queue	e	cce_general			<b>*</b>
Job Con	nfigurations				
Job Na	lame(name)	Enter a job nam	ie.		
* Applica	cation				Ð
Main C	Class(class)	This parameter	cannot be left blank	if the Application is a	JAR package.
Applica	cation Parameters	Specify each pa	rameter on a separa	te line.	
Spark A	Arguments(conf)	Enter argument	s using the key = va	lue format. Press Ente	r to separate
			the pairs.		
					10
Job Typ	ype	Basic	Al-enhanced	Image	
		Select	Ψ.	Select	
	c.	Configure Image	Use Custom Imag		
JAR Pa	ackage Dependencies (jars)				*
Datas	n Ela Desendencias ( nu filas)				
-30101	in the Dependencies (py-mes)				•
Other	r Dependencies (files)				*
Group	o Name				
		_			
Access	s Metadata	Yes	No		
Retry u	upon Failure	Yes	No		
Advant	nced Settings	Skip	Configure		

Figure 19-2 Selecting a custom image on the DLI Spark job editing page

Fill Form Write API		Execute Save as Template
oourranic(name)	restotapidonaliconitoswia.ioooo	
* Application		arkPi.jar
Main Class(class)	t i	
Application Parameters	5	
Spark Arguments(conf)	spark.kubernetes.executor.deleteOnTe	ermination=false
Job Type	Basic Al-enhanced	Image
	ei_	•
	The ei_dlics	not authorized.
	Authorize	
	Configure ImageUse Custom Image	

## Figure 19-3 Authorizing a Spark job image

Figure 19-4 Selecting a custom image on the DLI Flink Jar job editing page

test_s				
ID: 93229 Job Type: Fli	nk Jar			
Queue	xie_container_general			
Application ?	WindowJoin.jar	View Built-in Dependencies		
Main Class	Default Manually assign			
	Default main class is specified by the Manifest file of the	application.		
Class Arguments	Enter a class argument (Separate multiple class arguments with spaces).			
JAR Package Dependenc	-Select-	View Built-in Dependencies		
Other Dependencies ⑦	-Select			
Job Type	Basic Image			
	Irx/Irxflink   Configure ImageUse Custom Image	•		
Flink Version	1.10 💌	]		
Runtime Configuration				

• Specify the image parameter in job parameters on API to use a custom image to run a job. For details about Spark jobs, see **Creating a Batch Processing Job**. For details about Flink Jar jobs, see **Creating a Flink Jar Job**.

----End

# **19.2 Managing DLI Global Variables**

# What Is a Global Variable?

DLI allows you to set variables that are frequently used during job development as global variables on the DLI management console. This avoids repeated definitions during job editing and reduces development and maintenance costs. Global variables can be used to replace long and difficult variables, simplifying complex parameters and improving the readability of SQL statements.

This section describes how to create a global variable.

# Creating a Global Variable

- 1. In the navigation pane of the DLI console, choose **Global Configuration** > **Global Variables**.
- 2. On the **Global Variables** page, click **Create** in the upper right corner to create a global variable.

Parameter	Description
Variable	Name of the created global variable.
Value	Global variable value.

3. After creating a global variable, use **{{xxxx}}** in the SQL statement to replace the parameter value set as the global variable. **xxxx** indicates the variable name. For example, if you set global variable **abc** to represent the table name, replace the actual table name with **{{abc}}** in the table creation statement.

create table {{table\_name}} (String1 String, int4 int, varchar1 varchar(10)) partitioned by (int1 int,int2 int,int3 int)

#### **NOTE**

Do not use global variables in **OPTIONS** of the table creation statements.

#### **Related operations:**

- Modifying a global variable

On the **Global Variables** page, locate a desired variable and click **Modify** in the **Operation** column.

#### D NOTE

If there are multiple global variables with the same name in the same project under an account, delete the redundant global variables to ensure that the global variables are unique in the same project. In this case, all users who have the permission to modify the global variables can change the variable values.

#### - Deleting a global variable

On the **Global Variables** page, click **Delete** in the **Operation** column of a variable to delete the variable value.

#### D NOTE

- If there are multiple global variables with the same name in the same project under an account, delete the global variables created by the user first. If there are only unique global variables, all users who have the delete permission can delete the global variables.
- After a variable is deleted, the variable cannot be used in SQL statements.

#### **Permissions Management for Global Variables**

You can assign different users different global variables through permission settings. The administrator and owners of global variables have all permissions. You do not need to set permissions for them, and their global variable permissions cannot be modified by other users.

When setting global variables for a new user, the region hosting the user's group must have the **Tenant Guest** permission. For details about the **Tenant Guest** permission and how to apply for the permission, see **System-defined Permissions** and **Creating a User Group and Assigning Permissions**.

- Granting permissions on a global variable to a user
  - a. In the navigation pane on the left of the DLI console, choose Global Configuration > Global Variables. On the displayed page, locate a desired global variable and click Set Permission in the Operation column. On the displayed User Permissions page, you can grant, set, and revoke permissions on the global variable to users.
  - b. Click Grant Permission in the upper right corner.

Figure 19-5 Granting permissions of a global variable to a user

#### Grant Permission

* Username	Enter a username.		
Select the permis	sions to be granted to the user		
Select all			
Update		Delete	Grant Permission
Revoke Perr	nission	View Other User's Permissions	
		OK Cancel	

Parameter	Description
Username	Name of the IAM user who is granted permissions <b>NOTE</b> This username must be an existing IAM username.
Permissions	• <b>Update</b> : This permission allows you to update the global variable.
	• <b>Delete</b> : This permission allows you to delete the global variable.
	• <b>Grant Permission</b> : This permission allows you to grant permissions of the global variable to other users.
	• <b>Revoke Permission</b> : This permission allows you to revoke the global variable permissions that other users have but cannot revoke the global variable owner's permissions.
	• View Other User's Permissions: This permission allows you to view the global variable permissions of other users.

 Table 19-3 Global variable parameters

## • Granting permissions on a global variable to a user

On the **User Permissions** page, locate a desired IAM user and click **Set Permission** in the **Operation** column. **Table 19-3** lists the permission parameters.

If all permission options are grayed out, it means you do not have the authority to modify the permissions on this global variable. You can request the modification permission from users who have authorization, such as the administrator or group owners.

## • Revoking permissions on a global variable from a user

On the **User Permissions** page, locate a desired IAM user and click **Revoke Permission** in the **Operation** column. Once the revoke operation is complete, the IAM user will no longer have any permissions on the global variable.

# 19.3 Managing Program Packages of Jar Jobs

# 19.3.1 Package Management Overview

Before running DLI jobs, UDF JAR files or Jar job packages need to be uploaded to the cloud platform for unified management and maintenance.

There are two ways to manage packages:

- (Recommended) Upload packages to OBS: Upload Jar packages to an OBS bucket in advance and select the OBS path when configuring a job.
- The DLI package function will soon be discontinued. Upload packages to DLI: Upload Jar packages to an OBS bucket in advance, create a package on the

**Data Management** > **Package Management** page of the DLI management console, and select the DLI package when configuring a job.

This section describes how to upload and manage packages on the DLI management console.

#### **NOTE**

- The DLI package function will soon be discontinued. When using Spark 3.3.1 or later or Flink 1.15 or later to run jobs, you are advised to select packages stored in OBS.
- When packaging Spark or Flink Jar jobs, do not upload the dependency packages that the platform already has to avoid conflicts with the built-in dependency packages of the platform. Refer to **DLI Built-in Dependencies** for built-in dependency packages.

# **Notes and Constraints**

ltem	Description
Package	• A package can be deleted, but a package group cannot be deleted.
	• The following types of packages can be uploaded:
	– <b>JAR</b> : JAR file
	<ul> <li>PyFile: User Python file</li> </ul>
	– File: User file
	- ModelFile: User AI model file

#### Table 19-4 Notes and constraints on package usage

## Package Management Page

#### Table 19-5 Package management parameters

Parameter	Description
Group Name	Name of the group to which the package belongs. If the package is not grouped, is displayed.
Package Name	Name of a package.
Owner	Name of the user who uploads the package.
Туре	<ul> <li>Type of a package. The following package types are supported:</li> <li>JAR: JAR file</li> <li>PyFile: User Python file</li> <li>File: User file</li> <li>ModelFile: User AI model file</li> </ul>

Parameter	Description		
Status	Status of the package to be created.		
	Uploading: The file is being uploaded.		
	• Finished: The resource package has been uploaded.		
	Failed: The resource package upload failed.		
Created	Time when a package is created.		
Updated	Time when the package is updated.		
Operation	Manage Permissions: Manage user permissions for a package.		
	Delete: Delete the package.		
	More:		
	• <b>Modify Owner</b> : Modify the owner of the package.		
	• <b>Tags</b> : Add or edit package tags.		

# 19.3.2 Creating a DLI Package

DLI allows you to submit program packages in batches to the general-use queue for running.

## D NOTE

If you need to update a package, you can use the same package or file to upload it to the same location (in the same group) on DLI to overwrite the original package or file.

# Prerequisites

All software packages must be uploaded to OBS for storage in advance.

# Procedure

- On the left of the management console, choose Data Management > Package Management.
- 2. On the **Package Management** page, click **Create** in the upper right corner to create a package.
- 3. In the displayed **Create Package** dialog box, set related parameters by referring to **Table 19-6**.

## Figure 19-6 Creating a package

Create Package					×
Туре	JAR	PyFile	File	ModelFile	
★ OBS Path	Specify each p	parameter on a s	eparate line.		<b>D</b>
Group	Use existir	ng Use	e new	Do not use	
★ Group Name	test001			•	
Tags	It is recommended	ed that you use ⊺ esources. View p	TMS's predefine	ed tag function to a	add the same tag to
	To add a tag, en	ter a tag key and	a tag value be	elow.	
	Enter a tag ke	У	Enter	a tag value	Add
	10 tags available	e for addition.			
		ОК	Cancel		

## Table 19-6 Parameter description

Paramete r	Description
Package Type	<ul> <li>The following package types are supported:</li> <li>JAR: JAR file</li> <li>PyFile: User Python file</li> <li>File: User file</li> <li>ModelFile: User AI model file</li> </ul>
Package File Path	<ul> <li>Select the OBS path of the corresponding packages.</li> <li>NOTE <ul> <li>The packages must be uploaded to OBS for storage in advance.</li> <li>Only files can be selected.</li> </ul> </li> </ul>
Group Policy	You can select <b>Use existing group</b> , <b>Use new group</b> , or <b>No grouping</b> .

Paramete r	Description		
Group	• Use existing group: Select an existing group.		
Name	• Use new group: Enter a custom group name.		
	• No grouping: No need to select or enter a group name.		
	NOTE		
	<ul> <li>If you select a group, the permission management refers to the permissions of the corresponding package group.</li> </ul>		
	<ul> <li>If no group is selected, the permission management refers to the permissions of the corresponding package.</li> </ul>		
	For details about how to manage permissions on package groups and packages, see Managing Permissions on Packages and Package Groups.		
Tag	Tags used to identify cloud resources. A tag includes the tag key and tag value. If you want to use the same tag to identify multiple cloud resources, that is, to select the same tag from the drop-down list box for all services, you are advised to create predefined tags on the Tag Management Service (TMS).		
	If your organization has configured tag policies for DLI, add tags to resources based on the policies. If a tag does not comply with the tag policies, resource creation may fail. Contact your organization administrator to learn more about tag policies.		
	For details, see Tag Management Service User Guide.		
	NOTE		
	• A maximum of 20 tags can be added.		
	<ul> <li>Only one tag value can be added to a tag key.</li> </ul>		
	<ul> <li>The key name in each resource must be unique.</li> </ul>		
	• Tag key: Enter a tag key name in the text box.		
	<b>NOTE</b> A tag key can contain a maximum of 128 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed, but the value cannot start or end with a space or start with _ <b>sys</b>		
	• Tag value: Enter a tag value in the text box.		
	<b>NOTE</b> A tag value can contain a maximum of 255 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed.		

4. Click OK.

After a package is created, you can view and select the package for use on the **Package Management** page.

Once a job is executed, you can free up DLI storage space by promptly deleting the job package from the package management page if it is no longer needed.

# **19.3.3 Configuring DLI Package Permissions**

By configuring permissions, you can grant different package groups or packages to various users, ensuring that job efficiency remains unaffected and job performance is maintained.

- Administrators and package group owners have all permissions for the package group. No permission settings are required, and other users cannot modify their package group permissions.
- Administrators and package owners have all permissions for the package. No permission settings are required, and other users cannot modify their package permissions.
- Package groups are used to manage packages with consistent behavior, so they support granting related permissions to package groups, but do not support granting individual permissions to packages within a package group.
- When an administrator assigns package group or package permissions to a new user, the region of the administrator's user group must have Tenant Guest permissions.

For details about the **Tenant Guest** permission and how to apply for the permission, see **System Permissions** and **Creating a User Group** in *Identity and Access Management User Guide*.

# **Configuring Permissions on Package Groups or Packages**

- On the Data Management > Package Management page, locate the package group or package whose permissions you want to grant and click Manage Permission in the Operation column.
- 2. On the displayed **User Permissions** page, click **Grant Permission** in the upper right corner of the page. In the dialog box that appears, enter a username in **Username**, select required permissions, and click **OK**.

#### 

- If you select a group when creating a package, you can manage permissions of the corresponding program package group.
- If you select **No grouping** when creating a package, you can manage permissions of the corresponding package.
- Granting permissions on package groups

Figure 19-7 Granting permissions on package groups

Grant	Permission	

* Username	Enter a username.			
Select the permi	ssions to be granted to the	user		
Select all				
Use Group		Update Group		Get Group
Delete Grou	p	Grant Permission		Revoke Permission
View Other	User's Permissions			
		Ok	Cancel	

Parameter	Description
Username	Name of the authorized IAM user. NOTE The username is the name of an existing IAM user.
Select the permissions	• <b>Use Group</b> : This permission allows you to use the package of this group.
to be granted to the user	• <b>Update Group</b> : This permission allows you to update the packages in the group, including creating a package in the group.
	• <b>Query Group</b> : This permission allows you to query the details of a package in a group.
	• <b>Delete Group</b> : This permission allows you to delete the package of the group.
	• <b>Grant Permission</b> : This permission allows you to grant group permissions to other users.
	• <b>Revoke Permission</b> : This permission allows you to revoke the group permissions that other users have but cannot revoke the group owner's permissions.
	• View Other User's Permissions: This permission allows you to view the group permissions of other users.

 Table 19-7 Permission parameters

- Granting permissions on packages

#### Figure 19-8 Granting permissions on package groups

#### Grant Permission

* Username	Enter a username.		
Select the permi	ssions to be granted to the	user	
Select all			
Use Resource	2	Update Resource	Get Resource
Delete Resou	irce	Grant Permission	Revoke Permission
View Other	Jser's Permissions		
		Ok Canc	el

#### Table 19-8 Permission parameters

Parameter	Description
Username	Name of the authorized IAM user.
	<b>NOTE</b> The username is the name of an existing IAM user.

Parameter	Description
Select the permissions to be granted to the user	• <b>Use Package</b> : This permission allows you to use the package.
	<ul> <li>Update Package: This permission allows you to update the package.</li> </ul>
	• <b>Query Package</b> : This permission allows you to query the package.
	• <b>Delete Package</b> : This permission allows you to delete the package.
	<ul> <li>Grant Permission: This permission allows you to grant package permissions to other users.</li> </ul>
	• <b>Revoke Permission</b> : This permission allows you to revoke the package permissions that other users have but cannot revoke the package owner's permissions.
	• View Other User's Permissions: This permission allows you to view the package permissions of other users.

# Modifying Permissions on Package Groups or Packages

- 1. On the **Data Management** > **Package Management** page, locate the desired package group or package and click **Manage Permission** in the **Operation** column.
- 2. On the displayed **User Permissions** page, click **Set Permission** in the **Operation** column of the IAM user for whom you want to modify the permissions.

 Table 19-7 and Table 19-8 list the detailed permission descriptions.

- If you set **Group** to **Use existing** or **Use new** when creating a package, you will modify the permissions on the group you selected.
- If you set **Group** to **Do not use** when creating a package, you will modify the permissions on the package.

If **Set Permission** on the **User Permissions** page is grayed out, you do not have the permission to modify the permissions on the package group or package.

You can request the **Grant Permission** and **Revoke Permission** permissions on package groups or packages from users who have authorization privileges, such as administrators or group owners.

# **Revoking Permissions on Package Groups or Packages**

DLI allows you to revoke permissions on package groups or packages with just one click.

 On the Data Management > Package Management page, locate the desired package group or package and click Manage Permission in the Operation column. 2. On the displayed **User Permissions** page, click **Revoke Permission** in the **Operation** column of the IAM user for whom you want to revoke the permissions.

Once the permissions are revoked, the IAM user does not have any permissions on the package group or package.

**NOTE** 

- If you set **Group** to **Use existing** or **Use new** when creating a package, you will revoke the permissions on the group you selected.
- If you set **Group** to **Do not use** when creating a package, you will revoke the permissions on the package.

# 19.3.4 Changing the DLI Package Owner

DLI allows you to change the owner of a package group or package.

- Log in to the DLI management console and choose Data Management > Package Management.
- 2. On the **Package Management** page, locate the package whose owner you want to change, click **More** in the **Operation** column, and select **Modify Owner**.
  - If the package has been grouped, you can change the owner of the group or package by selecting Group or Packages for Select Type and entering the new owner's username in Username.

Figure 19-9 Modifying the package owner

## Modify Owner

Group Name			
Name			
Select Type	Group	Resource	
Username			
		<b>OK</b> Cancel	

- If the package has not been grouped, change its owner directly.

#### Figure 19-10 Modifying the owner of a package

Modify Owner		×
Name	SparkJarObs-1.0-SNAPSHOT.jar	
★ Username	j	
	OK	

#### Table 19-9 Description

Parameter	Description
Group Name	• If you select a group when creating a package, the name of the group is displayed.
	<ul> <li>If no group is selected when creating a package, this parameter is not displayed.</li> </ul>
Name	Name of a package.
Select Type	<ul> <li>If you select a group when creating a package, you can change the owner of the group or package.</li> </ul>
	<ul> <li>If no group is selected when creating a package, this parameter is not displayed.</li> </ul>
Username	Name of the package owner.
	<b>NOTE</b> The username is the name of an existing IAM user.

3. Click **OK**.

# 19.3.5 Managing DLI Package Tags

Tags are key-value pairs that you can define to identify cloud resources. They assist you in categorizing and searching for cloud resources. A tag consists of a key and a value.

DLI allows you to add tags to package groups or packages.

- Log in to the DLI management console and choose Data Management > Package Management.
- 2. On the **Package Management** page, locate the desired package, click **More** in the **Operation** column, and select **Tags**.
- 3. On the page that appears, click **Add/Edit Tag** in the upper left corner.
- 4. In the Add/Edit Tag dialog box, set parameters.

Table 19-10	Tag parameters	
-------------	----------------	--

Param eter	Description
Tag key	You can specify the tag key in either of the following ways:
	<ul> <li>Click the text box and select a predefined tag key from the drop-down list.</li> </ul>
	To add a predefined tag, you need to create one on TMS and then select it from the tag key drop-down list. You can click <b>View predefined tags</b> to go to the <b>Predefined Tags</b> page of the TMS console. Then, click <b>Create Tag</b> in the upper corner of the page to create a predefined tag.
	For details, see <b>Creating Predefined Tags</b> in <i>Tag Management</i> Service User Guide.
	• Enter a tag key in the text box.
	<b>NOTE</b> A tag key can contain a maximum of 128 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed, but the value cannot start or end with a space or start with _ <b>sys</b>
Tag	You can specify the tag value in either of the following ways:
value	• Click the text box and select a predefined tag value from the drop-down list.
	Enter a tag value in the text box.
	NOTE A tag value can contain a maximum of 255 characters. Only letters, numbers, spaces, and special characters (:+-@) are allowed.

#### **NOTE**

- A maximum of 20 tags can be added.
- Only one tag value can be added to a tag key.
- The key name in each resource must be unique.
- 5. Click OK.
- 6. (Optional) To delete a tag, locate the tag in the tag list and click **Delete** in its **Operation** column.

# 19.3.6 DLI Built-in Dependencies

DLI built-in dependencies are provided by the platform by default. In case of conflicts, you do not need to upload them when packaging JAR packages of Spark or Flink Jar jobs.
## Spark 3.1.1 Dependencies

Table	19-11	Spark	3.1.1	dependencies
		opaire	<b>U</b>	acpenacheres

Dependency		
accessors-smart-1.2.jar	hive-shims- scheduler-3.1.0- h0.cbu.mrs.321.r10.jar	metrics-graphite-4.1.1.jar
activation-1.1.1.jar	hive-spark-client-3.1.0- h0.cbu.mrs.321.r10.jar	metrics-jmx-4.1.1.jar
aggdesigner- algorithm-6.0.jar	hive-standalone- metastore-3.1.0- h0.cbu.mrs.321.r10.jar	metrics-json-4.1.1.jar
aircompressor-0.16.jar	hive-storage-api-2.7.2.jar	metrics-jvm-4.1.1.jar
algebra_2.12-2.0.0- M2.jar	hive-vector-code- gen-3.1.0- h0.cbu.mrs.321.r10.jar	minlog-1.3.0.jar
annotations-17.0.0.jar	hk2-api-2.6.1.jar	netty-3.10.6.Final.jar
ant-1.10.9.jar	hk2-locator-2.6.1.jar	netty-all-4.1.86.Final.jar
ant-launcher-1.10.9.jar	hk2-utils-2.6.1.jar	netty- buffer-4.1.86.Final.jar
antlr4-runtime-4.8-1.jar	hppc-0.7.2.jar	netty- codec-4.1.86.Final.jar
antlr-runtime-3.5.2.jar	httpclient-4.5.6.jar	netty-codec- dns-4.1.86.Final.jar
aopalliance-1.0.jar	httpcore-4.4.10.jar	netty-codec- haproxy-4.1.86.Final.jar
aopalliance- repackaged-2.6.1.jar	istack-commons- runtime-3.0.8.jar	netty-codec- http2-4.1.86.Final.jar
apiguardian- api-1.1.0.jar	ivy-2.5.0.jar	netty-codec- http-4.1.86.Final.jar
arpack_combined_all-0. 1.jar	jackson- annotations-2.13.2.jar	netty-codec- memcache-4.1.86.Final.ja r
arrow-format-2.0.0.jar	jackson-core-2.13.2.jar	netty-codec- mqtt-4.1.86.Final.jar
arrow-memory- core-2.0.0.jar	jackson-core-asl-1.9.13- atlassian-4.jar	netty-codec- redis-4.1.86.Final.jar
arrow-memory- netty-2.0.0.jar	jackson- databind-2.13.2.2.jar	netty-codec- smtp-4.1.86.Final.jar

Dependency		
arrow-vector-2.0.0.jar	jackson-dataformat- yaml-2.13.2.jar	netty-codec- socks-4.1.86.Final.jar
asm-5.0.4.jar	jackson-datatype- jsr310-2.11.2.jar	netty-codec- stomp-4.1.86.Final.jar
audience- annotations-0.5.0.jar	jackson-mapper- asl-1.9.13-atlassian-4.jar	netty-codec- xml-4.1.86.Final.jar
automaton-1.11-8.jar	jackson-module-jaxb- annotations-2.13.2.jar	netty- common-4.1.86.Final.jar
avatica-1.22.0.jar	jackson-module- scala_2.12-2.13.2.jar	netty- handler-4.1.86.Final.jar
avatica-core-1.16.0.jar	jaeger-client-1.6.0.jar	netty-handler- proxy-4.1.86.Final.jar
avatica- metrics-1.16.0.jar	jaeger-core-1.6.0.jar	netty-handler-ssl- ocsp-4.1.86.Final.jar
avatica-server-1.16.0.jar	jaeger-thrift-1.6.0.jar	netty- resolver-4.1.86.Final.jar
avro-1.8.2.jar	jaeger- tracerresolver-1.6.0.jar	netty-resolver- dns-4.1.86.Final.jar
avro-ipc-1.8.2.jar	jakarta.activation- api-1.2.1.jar	netty-resolver-dns- classes- macos-4.1.86.Final.jar
avro-mapred-1.8.2.jar	jakarta.annotation- api-1.3.5.jar	netty-resolver-dns- native- macos-4.1.86.Final-osx- aarch_64.jar
java-sdk- bundle-1.11.856.jar	jakarta.el-3.0.3.jar	netty-resolver-dns- native- macos-4.1.86.Final-osx- x86_64.jar
base64-2.3.8.jar	jakarta.el-api-3.0.3.jar	netty- transport-4.1.86.Final.jar
bcpkix-jdk15on-1.69.jar	jakarta.inject-2.6.1.jar	netty-transport-classes- epoll-4.1.86.Final.jar
bcprov-jdk15on-1.69.jar	jakarta.servlet- api-4.0.3.jar	netty-transport-classes- kqueue-4.1.86.Final.jar
bcutil-jdk15on-1.69.jar	jakarta.validation- api-2.0.2.jar	netty-transport-native- epoll-4.1.86.Final-linux- aarch_64.jar

Dependency		
bonecp-0.8.0.RELEASE.ja r	jakarta.ws.rs-api-2.1.6.jar	netty-transport-native- epoll-4.1.86.Final-linux- x86_64.jar
breeze_2.12-1.0.jar	jakarta.xml.bind- api-2.3.2.jar	netty-transport-native- kqueue-4.1.86.Final-osx- aarch_64.jar
breeze- macros_2.12-1.0.jar	jamon-runtime-2.4.1.jar	netty-transport-native- kqueue-4.1.86.Final-osx- x86_64.jar
caffeine-2.8.1.jar	janino-3.0.16.jar	netty-transport-native- unix- common-4.1.86.Final.jar
calcite-core-1.22.0.jar	JavaEWAH-0.3.2.jar	netty-transport- rxtx-4.1.86.Final.jar
calcite-druid-1.19.0.jar	java-sdk-core-3.0.12.jar	netty-transport- sctp-4.1.86.Final.jar
calcite-linq4j-1.22.0.jar	javassist-3.25.0-GA.jar	netty-transport- udt-4.1.86.Final.jar
cats-kernel_2.12-2.0.0- M4.jar	javax.activation- api-1.2.0.jar	nimbus-jose-jwt-8.19.jar
checker-qual-3.5.0.jar	javax.annotation- api-1.3.2.jar	objenesis-2.5.1.jar
chill_2.12-0.9.5.jar	javax.inject-1.jar	okhttp-3.14.9.jar
chill-java-0.9.5.jar	javax.jdo-3.2.0-m3.jar	okio-1.17.2.jar
classmate-1.5.1.jar	java-xmlbuilder-1.1.jar	opencsv-2.3.jar
commons- beanutils-1.9.4.jar	javax.servlet-api-3.1.0.jar	opentelemetry- api-1.16.0.jar
commons-cli-1.2.jar	javax.transaction- api-1.3.jar	opentelemetry- context-1.16.0.jar
commons-codec-1.15.jar	javax.ws.rs-api-2.1.1.jar	opentelemetry- semconv-1.16.0-alpha.jar
commons- collections-3.2.2.jar	javolution-5.5.1.jar	opentracing- api-0.33.0.jar
commons- compiler-3.0.16.jar	jaxb-api-2.2.11.jar	opentracing- noop-0.33.0.jar
commons- compress-1.21.jar	jaxb-runtime-2.3.2.jar	opentracing- tracerresolver-0.1.8.jar

Dependency		
commons- configuration2-2.1.1.jar	jboss- logging-3.4.1.Final.jar	opentracing- util-0.33.0.jar
commons- crypto-1.0.0-20191105.j ar	jboss- threads-2.3.3.Final.jar	orc-core-1.6.8.jar
commons- daemon-1.0.13.jar	jcip-annotations-1.0-1.jar	orc-mapreduce-1.6.8.jar
commons-dbcp-1.4.jar	jcl-over-slf4j-1.7.36.jar	orc-shims-1.6.8.jar
commons- dbcp2-2.6.0.jar	jcodings-1.0.57.jar	orc-tools-1.6.7- h0.cbu.mrs.321.r10.jar
commons- digester-2.1.jar	jdo-api-3.2.jar	oro-2.0.8.jar
commons- httpclient-3.1.jar	jersey-client-2.34.jar	osgi-resource- locator-1.0.3.jar
commons-io-2.8.0.jar	jersey-common-2.34.jar	paranamer-2.8.jar
commons-lang-2.4.jar	jersey-container- servlet-2.34.jar	parquet- column-1.12.2.jar
commons-lang-2.6.jar	jersey-container-servlet- core-2.34.jar	parquet- common-1.12.2.jar
commons-lang3-3.10.jar	jersey-hk2-2.34.jar	parquet- encoding-1.12.2.jar
commons- logging-1.2.jar	jersey-server-2.34.jar	parquet-format- structures-1.12.2.jar
commons- math3-3.4.1.jar	jets3t-0.9.4-1.0.0.jar	parquet- hadoop-1.12.2.jar
commons-net-3.1.jar	jettison-1.1.jar	parquet-hadoop- bundle-1.12.0-ei-2.0.jar
commons- pool2-2.6.1.jar	jetty- http-9.4.41.v20210516.jar	parquet- jackson-1.12.2.jar
commons-text-1.10.0.jar	jetty- io-9.4.41.v20210516.jar	postgresql-42.3.5.jar
commons- validator-1.7.jar	jetty- rewrite-9.4.43.v20210629.j ar	protobuf-java-2.5.0.jar
compress-lzf-1.0.3.jar	jetty- security-9.4.41.v20210516 .jar	py4j-0.10.9.jar

Dependency		
core-1.1.2.jar	jetty- server-9.4.41.v20210516.j ar	pyrolite-4.30.jar
curator-client-2.13.0.jar	jetty- servlet-9.4.41.v20210516.j ar	re2j-1.1.jar
curator- framework-2.13.0.jar	jetty- util-9.4.41.v20210516.jar	RoaringBitmap-0.9.0.jar
curator- recipes-2.13.0.jar	jetty-util- ajax-9.4.41.v20210516.jar	scala-collection- compat_2.12-2.1.1.jar
datanucleus-api- jdo-4.2.4.jar	jetty- webapp-9.4.41.v20210516 .jar	scala-compiler-2.12.16.jar
datanucleus- core-4.1.17.jar	jetty- xml-9.4.41.v20210516.jar	scala-library-2.12.16.jar
datanucleus-rdbms- fi-4.1.19-302022.jar	JLargeArrays-1.5.jar	scala-parser- combinators_2.12-1.1.2.ja r
derby-10.14.2.0.jar	jline-3.21.0.jar	scala-reflect-2.12.16.jar
disruptor-3.4.2.jar	joda-time-2.10.5.jar	scala-xml_2.12-1.2.0.jar
dli-catalog- client-2.3.7-20240108.0 90504-101.jar	jodd-core-3.5.2.jar	secComponentApi-1.1.8.j ar
dli-catalog-hive3- client-2.3.7-20240108.0 90513-100.jar	jodd-util-6.0.0.jar	serializer-2.7.2.jar
dli-catalog-hive- extension-2.3.7-202401 08.090517-100.jar	joni-2.1.43.jar	shapeless_2.12-2.3.3.jar
dnsjava-2.1.7.jar	jpam-1.1.jar	shims-0.9.0.jar
dropwizard-metrics- hadoop-metrics2- reporter-0.1.2.jar	jsch-0.1.72.jar	sketches-core-0.9.0.jar
error_prone_annotation s-2.18.0.jar	json-20210307.jar	slf4j-api-1.7.30.jar
esdk-obs-java- optimized-3.22.10.2.jar	json4s-ast_2.12-3.7.0- M5.jar	slf4j-log4j12-1.7.25.jar
esri-geometry- api-2.2.0.jar	json4s-core_2.12-3.7.0- M5.jar	snakeyaml-1.30.jar

Dependency		
fastutil-6.5.6.jar	json4s-jackson_2.12-3.7.0- M5.jar	snappy-java-1.1.8.2.jar
flatbuffers-java-1.9.0.jar	json4s-scalap_2.12-3.7.0- M5.jar	spark-avro_2.12-3.1.1- h1.cbu.dli.20230607.r1.jar
generex-1.0.2.jar	json-path-2.4.0.jar	spark-catalyst_2.12-3.1.1- h1.cbu.dli.20230607.r1.jar
glassfish-corba- omgapi-4.2.2.jar	json-smart-2.3.jar	spark-core_2.12-3.1.1- h1.cbu.dli.20230607.r1.jar
gson-2.8.9.jar	jsr305-3.0.0.jar	spark-graphx_2.12-3.1.1- h1.cbu.dli.20230607.r1.jar
gson-fire-1.8.5.jar	JTransforms-3.1.jar	spark-hive_2.12-3.1.1- h1.cbu.dli.20230607.r1.jar
guava-14.0.1.jar	jul-to-slf4j-1.7.36.jar	spark- kubernetes_2.12-3.1.1- h1.cbu.dli.20230607.r1.jar
guice-3.0.jar	kafka-clients-2.8.0.jar	spark-kvstore_2.12-3.1.1- h1.cbu.dli.20230607.r1.jar
guice- assistedinject-3.0.jar	kerb-admin-2.0.2.jar	spark- launcher_2.12-3.1.1- h1.cbu.dli.20230607.r1.jar
guice-servlet-4.0.jar	kerb-client-2.0.2.jar	spark-mllib_2.12-3.1.1- h1.cbu.dli.20230607.r1.jar
hadoop- annotations-3.1.1- h0.cbu.mrs.313.r9.jar	kerb-common-2.0.2.jar	spark-mllib- local_2.12-3.1.1- h1.cbu.dli.20230607.r1.jar
hadoop-archives-3.3.1- h0.cbu.mrs.321.r10.jar	kerb-core-2.0.2.jar	spark-network- common_2.12-3.1.1- h1.cbu.dli.20230607.r1.jar
hadoop-auth-3.3.1- h0.cbu.mrs.321.r16.jar	kerb-crypto-2.0.2.jar	spark-network- shuffle_2.12-3.1.1- h1.cbu.dli.20230607.r1.jar
hadoop-3.3.1- h0.cbu.mrs.321.r16.jar	kerb-identity-2.0.2.jar	spark-quota- manager_2.12-3.1.1-2.3.7 .dli-SNAPSHOT.jar
hadoop-client-3.1.1- h0.cbu.mrs.313.r9.jar	kerb-server-2.0.2.jar	spark-repl_2.12-3.1.1- h1.cbu.dli.20230607.r1.jar
hadoop-common-3.3.1- h0.cbu.mrs.321.r10.jar	kerb-simplekdc-2.0.2.jar	spark-sketch_2.12-3.1.1- h1.cbu.dli.20230607.r1.jar

Dependency		
hadoop-distcp-3.3.1- h0.cbu.mrs.321.r10.jar	kerb-util-2.0.2.jar	spark-sql_2.12-3.1.1- h1.cbu.dli.20230607.r1.jar
hadoop-hdfs-3.3.1- h0.cbu.mrs.321.r16.jar	kerby-asn1-2.0.2.jar	spark- streaming_2.12-3.1.1- h1.cbu.dli.20230607.r1.jar
hadoop-hdfs- client-3.3.1- h0.cbu.mrs.321.r10.jar	kerby-config-2.0.2.jar	spark-tags_2.12-3.1.1- h1.cbu.dli.20230607.r1.jar
hadoop-3.1.1-52.1.jar	kerby-pkix-2.0.2.jar	spark-unsafe_2.12-3.1.1- h1.cbu.dli.20230607.r1.jar
hadoop-mapreduce- client-common-3.1.1- h0.cbu.mrs.313.r9.jar	kerby-util-2.0.2.jar	spark- uquery_2.12-3.1.1-2.3.7.d li-SNAPSHOT.jar
hadoop-mapreduce- client-core-3.1.1- h0.cbu.mrs.313.r9.jar	kerby-xdr-2.0.2.jar	spire_2.12-0.17.0-M1.jar
hadoop-mapreduce- client-jobclient-3.1.1- h0.cbu.mrs.313.r9.jar	kotlin-stdlib-1.4.21.jar	spire- macros_2.12-0.17.0- M1.jar
hadoop-mapreduce- client-nativetask-3.3.1- h0.cbu.mrs.321.r10.jar	kotlin-stdlib- common-1.4.21.jar	spire- platform_2.12-0.17.0- M1.jar
hadoop-registry-3.3.1- h0.cbu.mrs.321.r10.jar	kryo-shaded-4.0.2.jar	spire-util_2.12-0.17.0- M1.jar
hadoop-shaded- guava-1.1.1.jar	kubernetes- client-5.4.1-20211025.jar	sqlline-1.3.0.jar
hadoop-shaded- protobuf_3_7-1.1.1.jar	kubernetes-model- admissionregistra- tion-5.4.1-20211025.jar	ST4-4.0.4.jar
hadoop-yarn-api-3.1.1- h0.cbu.mrs.313.r9.jar	kubernetes-model- apiextensions-5.4.1-20211 025.jar	stax2-api-4.2.1.jar
hadoop-yarn- client-3.1.1- h0.cbu.mrs.313.r9.jar	kubernetes-model- apps-5.4.1-20211025.jar	stax-api-1.0.1.jar
hadoop-yarn- registry-3.3.1- h0.cbu.mrs.321.r10.jar	kubernetes-model- autoscaling-5.4.1-202110 25.jar	stream-2.9.6.jar
hbase-asyncfs-2.4.14- h0.cbu.mrs.321.r10.jar	kubernetes-model- batch-5.4.1-20211025.jar	streamingClient

Dependency		
hbase-client-2.4.14- h0.cbu.mrs.321.r10.jar	kubernetes-model- certificates-5.4.1-2021102 5.jar	streamingClient010
hbase-common-2.4.14- h0.cbu.mrs.321.r10.jar	kubernetes-model- common-5.4.1-20211025.j ar	swagger- annotations-2.2.8.jar
hbase-hadoop2- compat-2.4.14- h0.cbu.mrs.321.r10.jar	kubernetes-model- coordination-5.4.1-20211 025.jar	tephra-api-0.6.0.jar
hbase-hadoop- compat-2.4.14- h0.cbu.mrs.321.r10.jar	kubernetes-model- core-5.4.1-20211025.jar	tephra-core-0.6.0.jar
hbase-http-2.4.14- h0.cbu.mrs.321.r10.jar	kubernetes-model- discovery-5.4.1-20211025. jar	tephra-hbase- compat-1.0-0.6.0.jar
hbase-logging-2.4.14- h0.cbu.mrs.321.r10.jar	kubernetes-model- events-5.4.1-20211025.jar	threetenbp-1.3.5.jar
hbase-metrics-2.4.14- h0.cbu.mrs.321.r10.jar	kubernetes-model- extensions-5.4.1-2021102 5.jar	threeten-extra-1.5.0.jar
hbase-metrics- api-2.4.14- h0.cbu.mrs.321.r10.jar	kubernetes-model- flowcontrol-5.4.1-202110 25.jar	tink-1.6.0.jar
hbase-procedure-2.4.14- h0.cbu.mrs.321.r10.jar	kubernetes-model- metrics-5.4.1-20211025.ja r	token-provider-2.0.2.jar
hbase-protocol-2.4.14- h0.cbu.mrs.321.r10.jar	kubernetes-model- networking-5.4.1-202110 25.jar	tomcat-servlet- api-8.5.61.jar
hbase-protocol- shaded-2.4.14- h0.cbu.mrs.321.r10.jar	kubernetes-model- node-5.4.1-20211025.jar	transaction-api-1.1.jar
hbase- replication-2.4.14- h0.cbu.mrs.321.r10.jar	kubernetes-model- policy-5.4.1-20211025.jar	twill-api-0.6.0- incubating.jar
hbase-server-2.4.14- h0.cbu.mrs.321.r10.jar	kubernetes-model- rbac-5.4.1-20211025.jar	twill-common-0.6.0- incubating.jar
hbase-shaded- gson-4.1.4.jar	kubernetes-model- scheduling-5.4.1-2021102 5.jar	twill-core-0.6.0- incubating.jar

Dependency				
hbase-shaded- jersey-4.1.4.jar	kubernetes-model- storageclass-5.4.1-202110 25.jar	twill-discovery-api-0.6.0- incubating.jar		
hbase-shaded- jetty-4.1.4.jar	leveldbjni- all-1.8-20191105.jar	twill-discovery- core-0.6.0-incubating.jar		
hbase-shaded- miscellaneous-4.1.4.jar	libfb303-0.9.3.jar	twill-zookeeper-0.6.0- incubating.jar		
hbase-shaded- netty-4.1.4.jar	libthrift-0.14.1- ei-311001.jar	univocity- parsers-2.9.1.jar		
hbase-shaded- protobuf-4.1.4.jar	log4j-1.2.17-cloudera1.jar	us-common-1.0.66.jar		
hbase-unsafe-4.1.4.jar	log4j-api-2.17.1.jar	velocity-1.7.jar		
hbase- zookeeper-2.4.14- h0.cbu.mrs.321.r10.jar	log4j-rolling- appender-20131024-2017. jar	velocity-engine- core-2.3.jar		
hibernate- validator-6.2.5.Final.jar	logging- interceptor-3.14.9.jar	wildfly-client- config-1.0.1.Final.jar		
HikariCP-2.6.1.jar	luxor-cluster-quota- manager- transport_2.12-2.3.7-2023 1226.034700-559.jar	wildfly- common-1.5.2.Final.jar		
hive-classification-3.1.0- h0.cbu.mrs.321.r10.jar	luxor- encrypt-2.3.7-20231226.0 34423-1046.jar	woodstox-core-5.4.0.jar		
hive-common-3.1.0- h0.cbu.mrs.321.r10.jar	luxor- fs3-2.3.7-20231226.03443 8-1039.jar	xalan-2.7.2.jar		
hive-exec-3.1.0- h0.cbu.mrs.321.r10- core.jar	luxor-obs- fs3-2.3.7-20231226.03444 3-1038.jar	xbean-asm7- shaded-4.15.jar		
hive-llap-client-2.3.3- ei-12-20210120.005053 -2.jar	luxor- rpc_2.12-2.3.7-20231226.0 34653-560.jar	xercesImpl-2.12.2.jar		
hive-llap- common-3.1.0- h0.cbu.mrs.321.r10.jar	luxor-scc- adapter-2.3.7-20231226.0 34418-1045.jar	xml-apis-1.4.01.jar		
hive-llap-tez-3.1.0- h0.cbu.mrs.321.r10.jar	luxor- transport-2.3.7-20231226. 034433-1038.jar	xnio-api-3.8.4.Final.jar		

Dependency				
hive-metastore-3.1.0- h0.cbu.mrs.321.r10.jar	lz4-java-1.7.1.jar	xz-1.5.jar		
hive-serde-3.1.0- h0.cbu.mrs.321.r10.jar	machinist_2.12-0.6.8.jar	zjsonpatch-0.3.0.jar		
hive-service-rpc-3.1.0- h0.cbu.mrs.321.r10.jar	macro- compat_2.12-1.1.1.jar	zookeeper-3.5.6- ei-302002.jar		
hive-shims-0.23-3.1.0- h0.cbu.mrs.321.r10.jar	memarts-ccsdk-1.0.jar	zookeeper-jute-3.5.6- ei-302002.jar		
hive-shims-3.1.0- h0.cbu.mrs.321.r10.jar	memory-0.9.0.jar	zstd-jni-1.4.9-1.jar		
hive-shims- common-3.1.0- h0.cbu.mrs.321.r10.jar	metrics-core-4.1.1.jar	-		

## Spark 2.4.5 Dependencies

 Table 19-12
 Spark 2.4.5
 dependencies

Dependency				
JavaEWAH-1.1.7.jar	httpclient-4.5.6.jar	lucene- queryparser-7.7.2.jar		
RoaringBitmap-0.7.45.jar	httpcore-4.4.10.jar	lucene- sandbox-7.7.2.jar		
ST4-4.3.1.jar	ivy-2.4.0.jar	luxor- encrypt-2.0.0-2022062 3.010726-213.jar		
accessors-smart-1.2.jar	jackson- annotations-2.11.4.jar	luxor- fs3-2.0.0-20220623.01 0750-209.jar		
activation-1.1.1.jar	jackson-core-2.11.4.jar	luxor-obs- fs3-2.0.0-20220623.01 0756-209.jar		
aircompressor-0.16.jar	jackson-core-asl-1.9.13- atlassian-4.jar	luxor- rpc_2.11-2.0.0-202206 23.010737-182.jar		
alluxio-2.3.1-luxor- SNAPSHOT-client.jar	jackson-databind-2.11.4.jar	luxor- transport-2.0.0-202206 23.010744-71.jar		

Dependency		
annotations-17.0.0.jar	jackson-dataformat- yaml-2.11.4.jar	lz4-java-1.7.1.jar
antlr-2.7.7.jar	jackson-datatype- jsr310-2.11.2.jar	machinist_2.11-0.6.1.ja r
antlr-runtime-3.4.jar	jackson-jaxrs-base-2.10.3.jar	macro- compat_2.11-1.1.1.jar
antlr4-runtime-4.8-1.jar	jackson-jaxrs-json- provider-2.10.3.jar	metrics-core-3.1.5.jar
aopalliance-1.0.jar	jackson-mapper-asl-1.9.13- atlassian-4.jar	metrics- graphite-3.1.5.jar
aopalliance- repackaged-2.4.0-b34.jar	jackson-module-jaxb- annotations-2.10.3.jar	metrics- jmx-4.1.12.1.jar
apache-log4j- extras-1.2.17.jar	jackson-module- paranamer-2.11.4.jar	metrics-json-3.1.5.jar
arpack_combined_all-0.1 .jar	jackson-module- scala_2.11-2.11.4.jar	metrics-jvm-3.1.5.jar
arrow-format-0.12.0.jar	jakarta.activation- api-1.2.1.jar	minlog-1.3.0.jar
arrow-memory-0.12.0.jar	jakarta.xml.bind- api-2.3.2.jar	mssql- jdbc-6.2.1.jre7.jar
arrow-vector-0.12.0.jar	janino-3.0.9.jar	netty- all-4.1.51.Final.jar
asm-5.0.4.jar	java-util-1.9.0.jar	nimbus-jose- jwt-8.19.jar
audience- annotations-0.5.0.jar	java-xmlbuilder-1.1.jar	objenesis-2.5.1.jar
automaton-1.11-8.jar	javassist-3.18.1-GA.jar	okhttp-3.14.9.jar
avro-1.8.2.jar	javax.annotation-api-1.2.jar	okio-1.17.2.jar
avro-ipc-1.8.2.jar	javax.inject-1.jar	opencsv-2.3.jar
avro-mapred-1.8.2.jar	javax.inject-2.4.0-b34.jar	opencsv-4.6.jar
java-sdk- bundle-1.11.856.jar	javax.servlet-api-3.1.0.jar	opencv-4.3.0-2.jar
base64-2.3.8.jar	javax.ws.rs-api-2.0.1.jar	orc-core-1.6.8- nohive.jar
bcpkix-jdk15on-1.66.jar	javolution-5.3.1.jar	orc-mapreduce-1.6.8- nohive.jar
bcprov-jdk15on-1.67.jar	jaxb-api-2.2.11.jar	orc-shims-1.6.8.jar

Dependency		
bonecp-0.8.0.RELEASE.ja r	jcip-annotations-1.0-1.jar	oro-2.0.8.jar
breeze- macros_2.11-0.13.2.jar	jcl-over-slf4j-1.7.30.jar	osgi-resource- locator-1.0.1.jar
breeze_2.11-0.13.2.jar	jdo-api-3.0.1.jar	paranamer-2.8.jar
calcite-avatica-1.2.0- incubating.jar	jersey-client-2.23.1.jar	parquet- column-1.12.2.jar
chill-java-0.9.3.jar	jersey-common-2.23.1.jar	parquet- common-1.12.2.jar
chill_2.11-0.9.3.jar	jersey-container- servlet-2.23.1.jar	parquet- encoding-1.12.2.jar
commons- beanutils-1.9.4.jar	jersey-container-servlet- core-2.23.1.jar	parquet-format- structures-1.12.2.jar
commons-cli-1.2.jar	jersey-guava-2.23.1.jar	parquet- hadoop-1.12.2.jar
commons-codec-1.15.jar	jersey-media-jaxb-2.23.1.jar	parquet-hadoop- bundle-1.6.0.jar
commons- collections-3.2.2.jar	jersey-server-2.23.1.jar	parquet- jackson-1.12.2.jar
commons- collections4-4.2.jar	jets3t-0.9.4.jar	postgresql-42.2.14.jar
commons- compiler-3.0.9.jar	jettison-1.1.jar	protobuf-java-2.5.0.jar
commons- compress-1.4.1.jar	jetty- http-9.4.34.v20201102.jar	py4j-0.10.7.jar
commons- configuration2-2.1.1.jar	jetty-io-9.4.34.v20201102.jar	pyrolite-4.13.jar
commons- crypto-1.0.0-20191105.ja r	jetty- security-9.4.34.v20201102.ja r	re2j-1.1.jar
commons- daemon-1.0.13.jar	jetty- server-9.4.34.v20201102.jar	scala- compiler-2.11.12.jar
commons- dbcp2-2.7.0.jar	jetty- servlet-9.4.34.v20201102.jar	scala- library-2.11.12.jar
commons- httpclient-3.1.jar	jetty- util-9.4.34.v20201102.jar	scala-parser- combinators_2.11-1.1. 2.jar

Dependency		
commons-io-2.5.jar	jetty-util- ajax-9.4.34.v20201102.jar	scala- reflect-2.11.12.jar
commons-lang-2.6.jar	jetty- webapp-9.4.34.v20201102.ja r	scala- xml_2.11-1.0.5.jar
commons-lang3-3.5.jar	jetty- xml-9.4.34.v20201102.jar	secComponentApi-1.0. 6.jar
commons- logging-1.2.jar	joda-time-2.9.3.jar	shapeless_2.11-2.3.2.ja r
commons- math3-3.4.1.jar	jodd-core-3.5.2.jar	shims-0.7.45.jar
commons-net-3.1.jar	json-20200518.jar	slf4j-api-1.7.30.jar
commons- pool2-2.8.0.jar	json-io-2.5.1.jar	slf4j-log4j12-1.7.30.jar
commons-text-1.3.jar	json-sanitizer-1.2.1.jar	snakeyaml-1.26.jar
compress-lzf-1.0.3.jar	json-smart-2.3.jar	snappy-java-1.1.8.2.jar
core-1.1.2.jar	json4s-ast_2.11-3.5.3.jar	solr-core-7.7.2.jar
crypter-0.0.6.jar	json4s-core_2.11-3.5.3.jar	solr-solrj-7.7.2.jar
curator-client-4.2.0.jar	json4s-jackson_2.11-3.5.3.jar	spark- avro_2.11-2.4.5.0100-2 .0.0.dli-20220617.0855 36-9.jar
curator- framework-4.2.0.jar	json4s-scalap_2.11-3.5.3.jar	spark- avro_2.11-4.0.0.jar
curator-recipes-2.7.1.jar	jsp-api-2.1.jar	spark- catalyst_2.11-2.4.5.010 0-2.0.0.dli-20220617.0 85405-16.jar
datanucleus-api- jdo-3.2.6.jar	jsr305-1.3.9.jar	spark- core_2.11-2.4.5.0100-2 .0.0.dli-20220617.0853 27-16.jar
datanucleus- core-3.2.10.jar	jta-1.1.jar	spark- graphx_2.11-2.4.5.010 00.dli-20220617.0853 36-16.jar

Dependency		
datanucleus- rdbms-3.2.9.jar	jtransforms-2.4.0.jar	spark- hive_2.11-2.4.5.0100-2. 0.0.dli-20220617.0854 23-16.jar
derby-10.14.2.0.jar	jts-core-1.16.1.jar	spark- kubernetes_2.11-2.4.5. 0100-2.0.0.dli-2022061 7.085519-16.jar
dnsjava-2.1.7.jar	jul-to-slf4j-1.7.30.jar	spark- kvstore_2.11-2.4.5.010 0-2.0.0.dli-20220617.0 85249-16.jar
ecj-3.21.0.jar	junit-4.11.jar	spark- launcher_2.11-2.4.5.01 00-2.0.0.dli-20220617. 085435-16.jar
ehcache-3.3.1.jar	kerb-admin-1.0.1.jar	spark-mllib- local_2.11-2.4.5.0100- 2.0.0.dli-20220617.085 349-16.jar
expiringmap-0.5.9.jar	kerb-client-1.0.1.jar	spark- mllib_2.11-2.4.5.0100- 2.0.0.dli-20220617.085 342-16.jar
fastutil-8.2.3.jar	kerb-common-1.0.1.jar	spark-network- common_2.11-2.4.5.01 00-2.0.0.dli-20220617. 085254-16.jar
flatbuffers-java-1.9.0.jar	kerb-core-1.0.1.jar	spark-network- shuffle_2.11-2.4.5.010 0-2.0.0.dli-20220617.0 85300-16.jar
fst-2.50.jar	kerb-crypto-1.0.1.jar	spark- om_2.11-2.4.5.0100-2. 0.0.dli-20220617.0853 16-16.jar
generex-1.0.2.jar	kerb-identity-1.0.1.jar	spark- repl_2.11-2.4.5.0100-2. 0.0.dli-20220617.0854 30-16.jar

Dependency		
geronimo- jcache_1.0_spec-1.0- alpha-1.jar	kerb-server-1.0.1.jar	spark- sketch_2.11-2.4.5.0100 -2.0.0.dli-20220617.08 5243-16.jar
gson-2.2.4.jar	kerb-simplekdc-1.0.1.jar	spark- sql_2.11-2.4.5.0100-2.0 .0.dli-20220617.08541 4-16.jar
guava-14.0.1.jar	kerb-util-1.0.1.jar	spark- streaming_2.11-2.4.5.0 1000.dli-20220617.08 5359-16.jar
guice-4.0.jar	kerby-asn1-1.0.1.jar	spark- tags_2.11-2.4.5.0100-2 .0.0.dli-20220617.0853 22-16.jar
guice-servlet-4.0.jar	kerby-config-1.0.1.jar	spark- unsafe_2.11-2.4.5.0100 -2.0.0.dli-20220617.08 5311-16.jar
hadoop- annotations-3.1.1- ei-302002.jar	kerby-pkix-1.0.1.jar	spark- uquery_2.11-2.4.5.010 0-2.0.0.dli- SNAPSHOT.jar
hadoop-auth-3.1.1- ei-302002.jar	kerby-util-1.0.1.jar	spark- yarn_2.11-2.4.5.0100-2 .0.0.dli-20220617.0855 31-16.jar
hadoop-3.1.1- ei-302002.jar	kerby-xdr-1.0.1.jar	spire- macros_2.11-0.13.0.jar
hadoop-client-3.1.1- ei-302002.jar	kryo-shaded-4.0.2.jar	spire_2.11-0.13.0.jar
hadoop-common-3.1.1- ei-302002.jar	kubernetes- client-5.4.1-20211025.jar	stax-api-1.0-2.jar
hadoop-hdfs-3.1.1- ei-302002.jar	kubernetes-model- admissionregistra- tion-5.4.1-20211025.jar	stax2-api-3.1.4.jar
hadoop-hdfs- client-3.1.1-ei-302002.jar	kubernetes-model- apiextensions-5.4.1-2021102 5.jar	stream-2.7.0.jar
hadoop-3.1.1-46.jar	kubernetes-model- apps-5.4.1-20211025.jar	stringtemplate-3.2.1.ja r

Dependency		
hadoop-mapreduce- client-common-3.1.1- ei-302002.jar	kubernetes-model- autoscaling-5.4.1-20211025. jar	threeten- extra-1.5.0.jar
hadoop-mapreduce- client-core-3.1.1- ei-302002.jar	kubernetes-model- batch-5.4.1-20211025.jar	tink-1.6.0.jar
hadoop-mapreduce- client-jobclient-3.1.1- ei-302002.jar	kubernetes-model- certificates-5.4.1-20211025.j ar	token- provider-1.0.1.jar
hadoop-minikdc-3.1.1- ei-302002.jar	kubernetes-model- common-5.4.1-20211025.jar	tomcat-api-9.0.39.jar
hadoop-yarn-api-3.1.1- ei-302002.jar	kubernetes-model- coordination-5.4.1-2021102 5.jar	zookeeper-jute-3.5.6- ei-302002.jar
hadoop-yarn- client-3.1.1-ei-302002.jar	kubernetes-model- core-5.4.1-20211025.jar	tomcat-el- api-9.0.39.jar
hadoop-yarn- common-3.1.1- ei-302002.jar	kubernetes-model- discovery-5.4.1-20211025.ja r	tomcat- jasper-9.0.39.jar
hadoop-yarn- registry-3.1.1- ei-302002.jar	kubernetes-model- events-5.4.1-20211025.jar	tomcat-jasper- el-9.0.39.jar
hadoop-yarn-server- applicationhistoryser- vice-3.1.1-ei-302002.jar	kubernetes-model- extensions-5.4.1-20211025.j ar	tomcat-jsp- api-9.0.39.jar
hadoop-yarn-server- common-3.1.1- ei-302002.jar	kubernetes-model- flowcontrol-5.4.1-20211025. jar	tomcat-juli-9.0.39.jar
hadoop-yarn-server- resourcemanager-3.1.1- ei-302002.jar	kubernetes-model- metrics-5.4.1-20211025.jar	tomcat-servlet- api-9.0.39.jar
hadoop-yarn-server- web-proxy-3.1.1- ei-302002.jar	kubernetes-model- networking-5.4.1-20211025. jar	tomcat-util-9.0.39.jar
hamcrest-core-1.3.jar	kubernetes-model- node-5.4.1-20211025.jar	tomcat-util- scan-9.0.39.jar
hive- common-1.2.1-2.0.0.dli- 20220528.090500-402.ja r	kubernetes-model- policy-5.4.1-20211025.jar	univocity- parsers-2.7.3.jar

Dependency		
hive- exec-1.2.1-2.0.0.dli-2022 0528.090521-401.jar	kubernetes-model- rbac-5.4.1-20211025.jar	zstd-jni-1.4.9-1.jar
hive- metastore-1.2.1-2.0.0.dli -20220528.090509-402.j ar	kubernetes-model- scheduling-5.4.1-20211025.j ar	validation- api-1.1.0.Final.jar
hive- shims-0.23-1.2.1-2.0.0.dli -20220528.090445-403.j ar	kubernetes-model- storageclass-5.4.1-20211025 .jar	velocity-1.7.jar
hive- shims-1.2.1-2.0.0.dli-202 20528.090455-403.jar	leveldbjni- all-1.8-20191105.jar	woodstox- core-5.0.3.jar
hive-shims- common-1.2.1-2.0.0.dli- 20220528.090441-404.ja r	libfb303-0.9.3.jar	xbean-asm6- shaded-4.8.jar
hive-shims- scheduler-1.2.1-2.0.0.dli- 20220528.090450-403.ja r	libthrift-0.12.0.jar	xercesImpl-2.12.0.jar
hk2-api-2.4.0-b34.jar	log4j-1.2.17-cloudera1.jar	xml-apis-1.4.01.jar
hk2-locator-2.4.0-b34.jar	log4j-rolling- appender-20131024-2017.ja r	xz-1.0.jar
hk2-utils-2.4.0-b34.jar	logging- interceptor-3.14.9.jar	zjsonpatch-0.3.0.jar
hppc-0.7.2.jar	lucene-analyzers- common-7.7.2.jar	zookeeper-3.5.6- ei-302002.jar
htrace-core4-4.2.0- incubating-1.0.0.jar	lucene-core-7.7.2.jar	-

## Spark 2.3.2 Dependencies

 Table 19-13
 Spark 2.3.2
 dependencies

Dependency		
accessors-smart-1.2.jar	HikariCP-java7-2.4.12.jar	logging- interceptor-3.14.4.jar

Dependency		
activation-1.1.1.jar	hive- common-1.2.1-2.1.0.dli-2 0201111.064115-91.jar	luxor- encrypt-2.1.0-20201106. 065437-53.jar
aircompressor-0.8.jar	hive- exec-1.2.1-2.1.0.dli-20201 111.064444-91.jar	luxor- fs3-2.1.0-20201106.065 612-53.jar
alluxio-2.3.1-luxor- SNAPSHOT-client.jar	hive- metastore-1.2.1-2.1.0.dli- 20201111.064230-91.jar	luxor-obs- fs3-2.1.0-20201106.065 616-53.jar
antlr-2.7.7.jar	hk2-api-2.4.0-b34.jar	luxor- rpc_2.11-2.1.0-2020110 6.065541-53.jar
antlr4-runtime-4.8-1.jar	hk2-locator-2.4.0-b34.jar	luxor-rpc- protobuf2-2.1.0-202011 06.065551-53.jar
antlr-runtime-3.4.jar	hk2-utils-2.4.0-b34.jar	lz4-java-1.7.1.jar
aopalliance-1.0.jar	hppc-0.7.2.jar	machinist_2.11-0.6.1.jar
aopalliance- repackaged-2.4.0-b34.jar	htrace-core4-4.2.0- incubating-1.0.0.jar	macro- compat_2.11-1.1.1.jar
apache-log4j- extras-1.2.17.jar	httpclient-4.5.4.jar	metrics-core-3.1.5.jar
arpack_combined_all-0.1.j ar	httpcore-4.4.7.jar	metrics- graphite-3.1.5.jar
arrow-format-0.8.0.jar	ivy-2.4.0.jar	metrics-jmx-4.1.12.1.jar
arrow-memory-0.8.0.jar	j2objc-annotations-1.3.jar	metrics-json-3.1.5.jar
arrow-vector-0.8.0.jar	jackson- annotations-2.10.0.jar	metrics-jvm-3.1.5.jar
asm-5.0.4.jar	jackson-core-2.10.0.jar	minlog-1.3.0.jar
audience- annotations-0.5.0.jar	jackson-core-asl-1.9.13- atlassian-4.jar	mssql-jdbc-6.2.1.jre7.jar
automaton-1.11-8.jar	jackson- databind-2.10.0.jar	netty-3.10.6.Final.jar
avro-1.7.7.jar	jackson-dataformat- yaml-2.10.0.jar	netty-all-4.1.51.Final.jar
avro-ipc-1.7.7.jar	jackson-datatype- jsr310-2.10.3.jar	nimbus-jose-jwt-8.19.jar
avro-ipc-1.7.7-tests.jar	jackson-jaxrs- base-2.10.3.jar	objenesis-2.1.jar

Dependency		
avro-mapred-1.7.7- hadoop2.jar	jackson-jaxrs-json- provider-2.10.3.jar	okhttp-3.14.4.jar
java-sdk- bundle-1.11.271.jar	jackson-mapper- asl-1.9.13-atlassian-4.jar	okio-1.17.2.jar
base64-2.3.8.jar	jackson-module-jaxb- annotations-2.10.3.jar	opencsv-2.3.jar
bcpkix-jdk15on-1.66.jar	jackson-module- paranamer-2.10.0.jar	opencsv-4.6.jar
bcprov-jdk15on-1.66.jar	jackson-module- scala_2.11-2.10.0.jar	opencv-4.3.0-2.jar
bonecp-0.8.0.RELEASE.jar	jakarta.activation- api-1.2.1.jar	orc-core-1.4.4-nohive.jar
breeze_2.11-0.13.2.jar	jakarta.xml.bind- api-2.3.2.jar	orc-mapreduce-1.4.4- nohive.jar
breeze- macros_2.11-0.13.2.jar	janino-3.0.8.jar	oro-2.0.8.jar
calcite-avatica-1.2.0- incubating.jar	javacpp-1.5.4.jar	osgi-resource- locator-1.0.1.jar
calcite-core-1.2.0- incubating.jar	javacpp-1.5.4-linux- x86_64.jar	paranamer-2.8.jar
calcite-linq4j-1.2.0- incubating.jar	javacv-1.5.4.jar	parquet- column-1.8.3.jar
checker-qual-2.11.1.jar	JavaEWAH-1.1.7.jar	parquet- common-1.8.3.jar
chill_2.11-0.8.4.jar	javassist-3.18.1-GA.jar	parquet- encoding-1.8.3.jar
chill-java-0.8.4.jar	javax.annotation- api-1.2.jar	parquet-format-2.3.1.jar
commons- beanutils-1.9.4.jar	javax.inject-1.jar	parquet- hadoop-1.8.3.jar
commons-cli-1.2.jar	javax.inject-2.4.0-b34.jar	parquet-hadoop- bundle-1.6.0.jar
commons- codec-2.0-20130428.2021 22-59.jar	javax.servlet-api-3.1.0.jar	parquet- jackson-1.8.3.jar
commons- collections-3.2.2.jar	javax.ws.rs-api-2.0.1.jar	parquet-format-2.3.1.jar

Dependency		
commons- collections4-4.2.jar	java-xmlbuilder-1.1.jar	parquet- hadoop-1.8.3.jar
commons- compiler-3.0.8.jar	javolution-5.3.1.jar	parquet-hadoop- bundle-1.6.0.jar
commons- compress-1.4.1.jar	jaxb-api-2.2.11.jar	parquet- jackson-1.8.3.jar
commons- configuration2-2.1.1.jar	jcip-annotations-1.0-1.jar	postgresql-42.2.14.jar
commons- crypto-1.0.0-20191105.jar	jcl-over-slf4j-1.7.26.jar	protobuf-java-2.5.0.jar
commons- daemon-1.0.13.jar	jdo-api-3.0.1.jar	py4j-0.10.7.jar
commons-dbcp-1.4.jar	jersey-client-2.23.1.jar	pyrolite-4.13.jar
commons-dbcp2-2.7.0.jar	jersey-common-2.23.1.jar	re2j-1.1.jar
commons- httpclient-3.1.jar	jersey-container- servlet-2.23.1.jar	RoaringBitmap-0.5.11.ja r
commons-io-2.5.jar	jersey-container-servlet- core-2.23.1.jar	scala- compiler-2.11.12.jar
commons-lang-2.6.jar	jersey-guava-2.23.1.jar	scala-library-2.11.12.jar
commons-lang3-3.5.jar	jersey-media- jaxb-2.23.1.jar	scalap-2.11.0.jar
commons-logging-1.2.jar	jersey-server-2.23.1.jar	scala-parser- combinators_2.11-1.1.0.j ar
commons-math3-3.4.1.jar	jets3t-0.9.4.jar	scala-reflect-2.11.12.jar
commons-net-2.2.jar	jetty- http-9.4.31.v20200723.jar	scala-xml_2.11-1.0.5.jar
commons-pool-1.5.4.jar	jetty- io-9.4.31.v20200723.jar	secComponentApi-1.0.5 c.jar
commons-pool2-2.8.0.jar	jetty- security-9.4.31.v2020072 3.jar	shapeless_2.11-2.3.2.jar
commons-text-1.3.jar	jetty- server-9.4.31.v20200723.j ar	slf4j-api-1.7.30.jar
compress-lzf-1.0.3.jar	jetty- servlet-9.4.31.v20200723. jar	slf4j-log4j12-1.7.30.jar

Dependency		
core-1.1.2.jar	jetty- util-9.4.31.v20200723.jar	snakeyaml-1.24.jar
curator-client-4.2.0.jar	jetty-util- ajax-9.4.31.v20200723.jar	snappy-java-1.1.7.5.jar
curator- framework-4.2.0.jar	jetty- webapp-9.4.31.v2020072 3.jar	spark- catalyst_2.11-2.3.2.0101 -2.1.0.dli-20201111.073 826-143.jar
curator-recipes-2.7.1.jar	jetty- xml-9.4.31.v20200723.jar	spark- core_2.11-2.3.2.01010. dli-20201111.073836-13 4.jar
datanucleus-api- jdo-3.2.6.jar	joda-time-2.9.3.jar	spark- graphx_2.11-2.3.2.0101- 2.1.0.dli-20201111.0738 47-129.jar
datanucleus- core-3.2.10.jar	jodd-core-4.2.0.jar	spark- hive_2.11-2.3.2.01010. dli-20201111.073854-13 2.jar
datanucleus- rdbms-3.2.9.jar	json-20200518.jar	spark- kubernetes_2.11-2.3.2.0 101-2.1.0.dli-20201111. 073916-85.jar
derby-10.12.1.1.jar	json4s-ast_2.11-3.2.11.jar	spark- kvstore_2.11-2.3.2.0101- 2.1.0.dli-20201111.0739 33-127.jar
dnsjava-2.1.7.jar	json4s- core_2.11-3.2.11.jar	spark- launcher_2.11-2.3.2.010 1-2.1.0.dli-20201111.07 3940-127.jar
ehcache-3.3.1.jar	json4s- jackson_2.11-3.2.11.jar	spark- mllib_2.11-2.3.2.0101-2. 1.0.dli-20201111.073946 -127.jar
eigenbase- properties-1.1.5.jar	json-sanitizer-1.2.1.jar	spark-mllib- local_2.11-2.3.2.0101-2. 1.0.dli-20201111.073953 -127.jar

Dependency		
error_prone_annotations- 2.3.4.jar	json-smart-2.3.jar	spark-network- common_2.11-2.3.2.010 1-2.1.0.dli-20201111.07 3959-127.jar
failureaccess-1.0.1.jar	jsp-api-2.1.jar	spark-network- shuffle_2.11-2.3.2.0101- 2.1.0.dli-20201111.0740 07-127.jar
fastutil-8.2.3.jar	jsr305-3.0.2.jar	spark- om_2.11-2.3.2.01010.dl i-20201111.074019-125. jar
ffmpeg-4.3.1-1.5.4.jar	jta-1.1.jar	spark- repl_2.11-2.3.2.0101-2.1. 0.dli-20201111.074028- 125.jar
ffmpeg-4.3.1-1.5.4-linux- x86_64.jar	jtransforms-2.4.0.jar	spark- sketch_2.11-2.3.2.0101- 2.1.0.dli-20201111.0740 35-125.jar
flatbuffers-1.2.0-3f79e055 .jar	jul-to-slf4j-1.7.26.jar	spark- sql_2.11-2.3.2.0101-2.1. 0.dli-20201111.074041- 126.jar
generex-1.0.2.jar	junit-4.11.jar	spark- streaming_2.11-2.3.2.01 01-2.1.0.dli-20201111.0 74100-123.jar
geronimo- jcache_1.0_spec-1.0- alpha-1.jar	kerb-admin-1.0.1.jar	spark- tags_2.11-2.3.2.0101-2.1 .0.dli-20201111.074136- 123.jar
gson-2.2.4.jar	kerb-client-1.0.1.jar	spark- tags_2.11-2.3.2.0101-2.1 .0.dli-20201111.074141- 124-tests.jar
guava-29.0-jre.jar	kerb-common-1.0.1.jar	spark- unsafe_2.11-2.3.2.0101- 2.1.0.dli-20201111.0741 44-123.jar

Dependency		
guice-4.0.jar	kerb-core-1.0.1.jar	spark- uquery_2.11-2.3.2.0101- 2.1.0.dli-20201111.0749 06-210.jar
guice-servlet-4.0.jar	kerb-crypto-1.0.1.jar	spark- yarn_2.11-2.3.2.0101-2.1 .0.dli-20201111.074151- 123.jar
hadoop- annotations-3.1.1- ei-302002.jar	kerb-identity-1.0.1.jar	spire_2.11-0.13.0.jar
hadoop-auth-3.1.1- ei-302002.jar	kerb-server-1.0.1.jar	spire- macros_2.11-0.13.0.jar
hadoop-3.1.1- ei-302002.jar	kerb-simplekdc-1.0.1.jar	ST4-4.3.1.jar
hadoop-client-3.1.1- ei-302002.jar	kerb-util-1.0.1.jar	stax2-api-3.1.4.jar
hadoop-common-3.1.1- ei-302002.jar	kerby-asn1-1.0.1.jar	stax-api-1.0-2.jar
hadoop-hdfs-3.1.1- ei-302002.jar	kerby-config-1.0.1.jar	stream-2.7.0.jar
hadoop-hdfs-client-3.1.1- ei-302002.jar	kerby-pkix-1.0.1.jar	stringtemplate-3.2.1.jar
hadoop-3.1.1-41.jar	kerby-util-1.0.1.jar	token-provider-1.0.1.jar
hadoop-mapreduce- client-common-3.1.1- ei-302002.jar	kerby-xdr-1.0.1.jar	univocity- parsers-2.5.9.jar
hadoop-mapreduce- client-core-3.1.1- ei-302002.jar	kryo-shaded-3.0.3.jar	validation- api-1.1.0.Final.jar
hadoop-mapreduce- client-jobclient-3.1.1- ei-302002.jar	kubernetes- client-4.9.2-20200804.jar	woodstox-core-5.0.3.jar
hadoop-minikdc-3.1.1- ei-302002.jar	kubernetes- model-4.9.2-20200804.jar	xbean-asm5- shaded-4.4.jar
hadoop-yarn-api-3.1.1- ei-302002.jar	kubernetes-model- common-4.9.2-20200804. jar	xercesImpl-2.12.0.jar
hadoop-yarn-client-3.1.1- ei-302002.jar	leveldbjni- all-1.8-20191105.jar	xml-apis-1.4.01.jar

Dependency		
hadoop-yarn- common-3.1.1- ei-302002.jar	libfb303-0.9.3.jar	xz-1.0.jar
hadoop-yarn- registry-3.1.1- ei-302002.jar	libthrift-0.12.0.jar	zjsonpatch-0.3.0.jar
hadoop-yarn-server- common-3.1.1- ei-302002.jar	listenablefuture-9999.0- empty-to-avoid-conflict- with-guava.jar	zookeeper-3.5.6- ei-302002.jar
hadoop-yarn-server-web- proxy-3.1.1-ei-302002.jar	log4j-1.2.17-cloudera1.jar	zookeeper-jute-3.5.6- ei-302002.jar
hamcrest-core-1.3.jar	log4j-rolling- appender-20131024-201 7.jar	zstd-jni-1.4.4-11.jar

### Flink 1.15 Dependencies

Obtain information about the Flink 1.15 dependencies from the logs of a Flink job.

- 1. Check the logs of a Flink job.
  - a. Log in to the DLI console. In the navigation pane on the left, choose **Job Management** > **Flink Jobs**.
  - b. Click the name of the desired job. On the displayed page, click the **Run Log** tab.
  - c. Check the latest run logs. For more logs, check the OBS bucket where the job logs are stored.
- 2. Search for dependency information in the logs.

Search for **Classpath:** in the logs to check the dependencies.

### Flink 1.12 Dependencies

Table	19-14	Flink	1.12	depend	lencies
-------	-------	-------	------	--------	---------

Dependency		
bcpkix-jdk15on-1.60.jar	flink-json-1.12.2- ei-313001- dli-2022011002.jar	libtensorflow-1.12.0.jar
bcprov-jdk15on-1.60.jar	flink- kubernetes_2.11-1.12.2- ei-313001- dli-2022011002.jar	log4j-1.2-api-2.17.1.jar

Dependency		
clickhouse-jdbc-0.3.1- ei-313001-SNAPSHOT.jar	flink-metrics- prometheus_2.11-1.12.2- ei-313001- dli-2022011002.jar	log4j-api-2.17.1.jar
commons-codec-1.9.jar	flink-obs-hadoop- fs-2.0.0-20220226.03442 1-73.jar	log4j-core-2.17.1.jar
commons- configuration-1.7.jar	flink-s3-fs- hadoop-1.12.2.jar	log4j-slf4j-impl-2.17.1.jar
dataflow-fs- obs-2.0.0-20220226.0344 02-190.jar	flink-shaded- zookeeper-3.6.3- ei-313001- SNAPSHOT.jar	luxor- encrypt-2.0.0-20220405. 072004-199.jar
deeplearning4j- core-0.9.1.jar	flink-sql-avro-1.12.2- ei-313001- dli-2022011002.jar	luxor- fs3-2.0.0-20220405.0720 25-195.jar
deeplearning4j- nlp-0.9.1.jar	flink-sql-avro-confluent- registry-1.12.2- ei-313001- dli-2022011002.jar	luxor-obs- fs3-2.0.0-20220405.0720 30-195.jar
deeplearning4j- nn-0.9.1.jar	flink-table_2.11-1.12.2- ei-313001- dli-2022011002.jar	manager-hadoop- security- crypter-8.1.3-313001- SNAPSHOT.jar
ejml-cdense-0.33.jar	flink-table- blink_2.11-1.12.2- ei-313001- dli-2022011002.jar	manager- wc2frm-8.1.3-313001- SNAPSHOT.jar
ejml-core-0.33.jar	guava-18.0.jar	mrs-obs- provider-3.1.1.49.jar
ejml-ddense-0.33.jar	guava-26.0-jre.jar	nd4j-api-0.9.1.jar
ejml-dsparse-0.33.jar	hadoop-hdfs- client-3.1.1- ei-302002.jar	nd4j-native-0.9.1.jar
ejml- experimental-0.33.jar	hadoop-3.1.1-46.jar	nd4j-native-api-0.9.1.jar
ejml-fdense-0.33.jar	hadoop- plugins-8.1.3-313001- SNAPSHOT.jar	nd4j-native- platform-0.9.1.jar
ejml-simple-0.33.jar	httpasyncclient-4.1.2.jar	okhttp-3.14.8.jar
ejml-zdense-0.33.jar	httpclient-4.5.3.jar	okio-1.14.0.jar

Dependency		
elsa-3.0.0-M7.jar	httpcore-4.4.4.jar	ranger-obs- client-0.1.1.jar
flink-changelog- json-1.12.2-ei-313001- dli-2022011002.jar	httpcore-nio-4.4.4.jar	secComponentApi-1.0.5.j ar
flink-csv-1.12.2- ei-313001- dli-2022011002.jar	java-xmlbuilder-1.1.jar	slf4j-api-1.7.26.jar
flink-dist_2.11-1.12.2- ei-313001- dli-2022011002.jar	jna-4.1.0.jar	tensorflow-1.12.0.jar

## Flink 1.10 Dependencies

Г

For details about the sample code for a Flink 1.10 program, see Writing Data to OBS Using Flink Jar.

Only queues created after December 2020 can use the Flink 1.10 dependencies.

Dependency		
bcpkix-jdk15on-1.60.jar	esdk-obs-java-3.20.6.1.jar	java-xmlbuilder-1.1.jar
bcprov-jdk15on-1.60.jar	flink-cep_2.11-1.10.0.jar	jna-4.1.0.jar
commons-codec-1.9.jar	flink-cep- scala_2.11-1.10.0.jar	libtensorflow-1.12.0.jar
commons- configuration-1.7.jar	flink-dist_2.11-1.10.0.jar	log4j-over-slf4j-1.7.26.jar
deeplearning4j- core-0.9.1.jar	flink- python_2.11-1.10.0.jar	logback-classic-1.2.3.jar
deeplearning4j- nlp-0.9.1.jar	flink-queryable-state- runtime_2.11-1.10.0.jar	logback-core-1.2.3.jar
deeplearning4j- nn-0.9.1.jar	flink-sql- client_2.11-1.10.0.jar	nd4j-api-0.9.1.jar
ejml-cdense-0.33.jar	flink-state-processor- api_2.11-1.10.0.jar	nd4j-native-0.9.1.jar
ejml-core-0.33.jar	flink-table_2.11-1.10.0.jar	nd4j-native-api-0.9.1.jar
ejml-ddense-0.33.jar	flink-table- blink_2.11-1.10.0.jar	nd4j-native- platform-0.9.1.jar

 Table 19-15
 Flink 1.10
 dependencies

Dependency		
ejml-dsparse-0.33.jar	guava-26.0-jre.jar	okhttp-3.14.8.jar
ejml- experimental-0.33.jar	hadoop-3.1.1-41.jar	okio-1.14.0.jar
ejml-fdense-0.33.jar	httpasyncclient-4.1.2.jar	secComponentApi-1.0.5.j ar
ejml-simple-0.33.jar	httpclient-4.5.3.jar	slf4j-api-1.7.26.jar
ejml-zdense-0.33.jar	httpcore-4.4.4.jar	tensorflow-1.12.0.jar
elsa-3.0.0-M7.jar	httpcore-nio-4.4.4.jar	-

## Flink 1.7.2 Dependencies

For details about the sample code of a Flink 1.7.2 program, see **luxor-demo\dli-flink-demo** in **DLI examples**.

Dependency		
bcpkix-jdk15on-1.60.jar	esdk-obs-java-3.1.3.jar	httpcore-4.4.4.jar
bcprov-jdk15on-1.60.jar	flink-cep_2.11-1.7.0.jar	httpcore-nio-4.4.4.jar
commons-codec-1.9.jar	flink-cep- scala_2.11-1.7.0.jar	java-xmlbuilder-1.1.jar
commons- configuration-1.7.jar	flink-dist_2.11-1.7.0.jar	jna-4.1.0.jar
deeplearning4j- core-0.9.1.jar	flink- gelly_2.11-1.7.0.jar	libtensorflow-1.12.0.jar
deeplearning4j-nlp-0.9.1.jar	flink-gelly- scala_2.11-1.7.0.jar	log4j-over-slf4j-1.7.21.jar
deeplearning4j-nn-0.9.1.jar	flink-ml_2.11-1.7.0.jar	logback-classic-1.2.3.jar
ejml-cdense-0.33.jar	flink- python_2.11-1.7.0.jar	logback-core-1.2.3.jar
ejml-core-0.33.jar	flink-queryable-state- runtime_2.11-1.7.0.jar	nd4j-api-0.9.1.jar
ejml-ddense-0.33.jar	flink-shaded- curator-1.7.0.jar	nd4j-native-0.9.1.jar
ejml-dsparse-0.33.jar	flink-shaded-hadoop2- uber-1.7.0.jar	nd4j-native-api-0.9.1.jar

Table 19-16 Flink 1.7.2 dependencies

Dependency			
ejml-experimental-0.33.jar	flink- table_2.11-1.7.0.jar	nd4j-native- platform-0.9.1.jar	
ejml-fdense-0.33.jar	guava-26.0-jre.jar	okhttp-3.14.8.jar	
ejml-simple-0.33.jar	hadoop-3.1.1-41-2020 1014.085840-4.jar	okio-1.14.0.jar	
ejml-zdense-0.33.jar	httpasyncclient-4.1.2.ja r	slf4j-api-1.7.21.jar	
elsa-3.0.0-M7.jar	httpclient-4.5.12.jar	tensorflow-1.12.0.jar	
log4j-api-2.16.0.jar	log4j-core-2.16.0.jar	log4j-api-2.8.2.jar	
log4j-core-2.8.2.jar	-	-	

# **19.4 Managing DLI Resource Quotas**

## What Is a Quota?

A quota limits the quantity of a resource available to users, thereby preventing spikes in the usage of the resource.

You can also request for an increased quota if your existing quota cannot meet your service requirements.

### How Do I View My Quotas?

- 1. Log in to the management console.
- 2. Click  $\bigcirc$  in the upper left corner and select a region and a project.
- In the upper right corner of the page, choose Resources > My Quotas. The Service Quota page is displayed.

#### Figure 19-11 My quotas

Billing Center	Resources	
My Resources		
My Quotas		
Open Beta Tests		
My Marketplace		

4. View the used and total quota of each type of resources on the displayed page.

If a quota cannot meet service requirements, increase a quota.

## How Do I Apply for a Higher Quota?

- 1. Log in to the management console.
- In the upper right corner of the page, choose Resources > My Quotas. The Service Quota page is displayed.

### Figure 19-12 My quotas

Billing Center	Resources
My Resources	
My Quotas	
Open Beta Tests	
My Marketplace	

- 3. Click Increase Quota.
- On the Create Service Ticket page, configure parameters as required. In the Problem Description area, fill in the content and reason for adjustment.
- 5. Select the agreement and click **Submit**.