Data Lake Insight

Best Practices

Issue 01

Date 2023-07-21





Copyright © Huawei Technologies Co., Ltd. 2023. 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

HUAWEI 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 "Vul. Response Process". For details about the policy, see the following website: https://www.huawei.com/en/psirt/vul-response-process

For enterprise customers who need to obtain vulnerability information, visit: https://securitybulletin.huawei.com/enterprise/en/security-advisory

Contents

1 Overview	1
2 Data Migration	2
2.1 Overview	2
2.2 Migrating Data from Hive to DLI	4
2.3 Migrating Data from MRS Kafka to DLI	13
2.4 Migrating Data from Elasticsearch to DLI	21
2.5 Migrating Data from RDS to DLI	29
2.6 Migrating Data from GaussDB(DWS) to DLI	36
3 Data Analysis	44
3.1 Analyzing Driving Behavior Data	44
3.2 Converting Data Format from CSV to Parquet	54
3.3 Analyzing E-commerce BI Reports	57
3.4 Analyzing DLI Billing Data	64
3.5 Using DLI Flink SQL to Analyze e-Commerce Business Data in Real Time	68
3.6 Interconnecting Yonghong BI with DLI to Submit Spark Jobs	83
3.6.1 Preparing for Yonghong BI Interconnection	83
3.6.2 Adding Yonghong BI Data Source	84
3.6.3 Creating Yonghong BI Data Set	88
3.6.4 Creating a Chart in Yonghong Bl	90
3.7 Interconnecting FineBl with DLI Trino	93
3.8 Interconnecting Power BI with DLI Trino	104
4 Connections	119
4.1 Configuring the Connection Between a DLI Queue and a Data Source in a Private Network	119
4.2 Configuring the Connection Between a DLI Queue and a Data Source in the Internet	125
A Change History	131

1 Overview

This document gives you best practices for data migration and analysis, helping you better use DLI for large-scale data analysis and processing.

Data Migration

You can use **Cloud Data Migration Service** (CDM) to easily migrate data from other cloud services or service platforms to DLI. You can refer to the following best practices:

- Migrating Data from Hive to DLI
- Migrating Data from MRS Kafka to DLI
- Migrating Data from Elasticsearch to DLI
- Migrating Data from RDS to DLI
- Migrating Data from GaussDB(DWS) to DLI

Data Analysis

DLI is widely used to analyze massive amounts of log data and in extract, transform, and load (ETL) processes, giving you great insight into data of a wide range of industries. You can refer to the following best practices of data analysis:

- Analyzing Driving Behavior Data
- Converting Data Format from CSV to Parquet
- Analyzing E-commerce BI Reports
- Analyzing DLI Billing Data

2 Data Migration

2.1 Overview

This section describes how you can migrate data to DLI in an efficient way. You can use **Cloud Data Migration Service** (CDM) to migrate data from other cloud services or platforms to DLI.

DLI is a serverless data processing and analysis service. It processes streaming data and batch data and supports interactive analysis. Its high-scalability framework supports the convergence of batch and streaming data analysis, and provides real-time, efficient, and diversified compute resources for TB-to EB-level data processing.

Best Practices of Data Migration

- You can migrate Hive data to DLI. For details, see Migrating Data from Hive to DLI.
- You can migrate Kafka data to DLI. For details, see Migrating Data from MRS Kafka to DLI.
- You can migrate Elasticsearch data to DLI. For details, see Migrating Data from Elasticsearch to DLI.
- You can migrate RDS data to DLI. For details, see Migrating Data from RDS to DLI.
- You can migrate GaussDB(DWS) data to DLI. For details, see Migrating Data from GaussDB(DWS) to DLI.

Data Type Mapping

If you migrate data from other cloud services or platforms to DLI, data types need to be converted and source and destination data must be mapped by type. **Table 2-1** lists the mapping relationships.

Table 2-1 Data type mapping

MySQL	Hive	DWS	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	DECIMAL	NUMERIC	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	BYTEA	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
TIME	Not supported (use String instead)	TIME	DATE	TIME	TIME	Not supporte d (use String instead)
DATETI ME	TIMESTA MP	TIMESTAM P	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

MySQL	Hive	DWS	Oracle	Postgre SQL	Hologre s	DLI Spark
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)	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

■ NOTE

If a service does not support a standard data type, you can use the recommended data type.

2.2 Migrating Data from Hive to DLI

This section describes how to use the CDM data synchronization function to migrate data from MRS Hive to DLI. Data of other MRS Hadoop components can be bidirectionally synchronized between CDM and DLI.

Prerequisites

• You have created a DLI SQL queue.

A CAUTION

When you create a queue, set its Type to For SQL.

- 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 details about how to create a cluster, see
 Creating a CDM Cluster.

□ NOTE

- If the destination data source is an on-premises database, you need the Internet or Direct Connect. When using the Internet, ensure that an EIP has been bound to the CDM cluster, the security group of CDM allows outbound traffic from the host where the off-cloud 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 MRS or GaussDB(DWS) on a cloud, the network must meet the following requirements:
 - i. If the CDM cluster and the cloud service are in different regions, a public network or a dedicated connection is required for enabling communication between the CDM cluster and the cloud service. If the Internet is used for communication, ensure 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, VPC, subnet, and security group, they 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 rules and security group rules.

For details about how to configure routes, see **Configure routes**. For details about how to configure security groups, see section **Security Group Configuration Examples**.

iii. The cloud service instance and the CDM cluster belong to the same enterprise project. If they do not, you can modify the enterprise project of the workspace.

In this example, the VPC, subnet, and security group of the CDM cluster are the same as those of the MRS cluster.

Step 1: Prepare Data

- Create a Hive table in the MRS cluster and insert data in the table.
 - Log in to MRS Manager by referring to Accessing FusionInsight Manager.
 - On MRS Manager, click **System** in the top navigation pane. On the page displayed, choose **Permission** > **Role** from the left navigation pane. On the displayed page, configure 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.

Role of Create Role

Role of Create Role of Create Role

Role of Create Rol

Figure 2-1 Creating a Hive role

For details about how to create a role, see Creating a Role.

- c. On the MRS Manager console, click **System** in the top navigation pane. On the displayed page, choose **Permission** > **User** from the left navigation pane. 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 password of the current user and enter it again.
 - iv. User Group and Primary Group: Select supergroup.
 - v. **Role**: Select the role created in **b** and the **Manager_viewer** role.

🖐 | FusionInsight Manager Homepage Cluster 🕶 Hosts O&M Audit Tenant Resources System User > Create * Usemame: System Human-Machine
Machine-Machine . User Type: . Confirm Password • Role Security Policy supergroup × · Domain and Mutual Trust supergroup Primary Group: Add Clear All Create Role hivetestrole × Manager_viewer ×

Figure 2-2 Creating a Hive User

d. Download and install the Hive client by referring to **Installing an MRS Client**. For example, the Hive client is installed in the **/opt/hiveclient** directory on the active MRS node.

OK Cancel

e. Go to the client installation directory as user root.

For example, run the cd /opt/hiveclient command.

f. Run the following command to set environment variables:

source bigdata_env

g. Run the following command to authenticate the user created in **c** as Kerberos authentication has been enabled for the current cluster:

kinit <Username in c>

Example: kinit hivetestusr

h. Run the following command to connect to Hive:

beeline

Create a table and insert data into it.

Run the following statement to create a table:

create table user_info(id string,name string,gender string,age int,addr string);

Run the following statements to 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 R").

insert into table user_info(id,name,gender,age,addr) values ("12005000202","B","male",20,"City B")

□ NOTE

In the preceding example, data is migrated by creating a table and inserting data. To migrate an existing Hive database, run the following commands to obtain Hive database and table information:

 Run the following command on the Hive client to obtain database information:

show databases

• Switch to the Hive database from which data needs to be migrated.

use Hive database name

 Run the following command to display information about all tables in the database:

show tables

 Run the following command to query the creation statement of the Hive table:

show create table table name

The queried table creation statements must be processed to comply with the DLI table creation syntax before being executed.

- Create a database and table on DLI.
 - a. 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 editing window to create a database, for example, the migrated DLI database **testdb**: For details about the syntax for creating a DLI database, see **Creating a Database**.

create database testdb:

b. Create a table in the database.

■ NOTE

You need to edit the table creation statement obtained by running **show create table** *hive table name* in MRS Hive to ensure the statement complies with the table creation syntax of DLI. For details about 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

- Create a CDM connection to MRS Hive.
 - a. Create a connection to link CDM to the data source MRS Hive.
 - Log in to the CDM console, 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 Hive** and click **Next**.

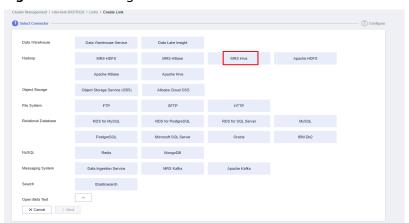


Figure 2-3 Selecting the MRS Hive connector

iii. Configure the connection. The following table describes the required parameters.

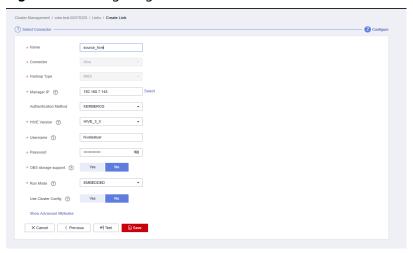
	Table 2-2 MRS	Hive	connection	configurations
--	---------------	------	------------	----------------

Parameter	Value
Name	Name of the MRS Hive data source, for example, source_hive
Manager IP	Click Select next to the text box and select the MRS Hive cluster. The Manager IP address is automatically specified.
Authenticatio n Method	Set this parameter to KERBEROS if Kerberos authentication is enabled for the MRS cluster. Set this parameter to SIMPLE if the MRS cluster is a common cluster. In this example, set this parameter to KERBEROS .

Parameter	Value
Hive Version	Set this parameter to the Hive version you have selected during MRS cluster creation. If the current Hive version is 3.1.0, set this parameter to HIVE_3_X.
Username	Name of the MRS Hive user created on c
Password	Password of the MRS Hive user

Retain default values for other parameters.

Figure 2-4 Configuring the connection to MRS Hive



- iv. Click **Save** to complete the configuration.
- b. Create a connection to link CDM to DLI.
 - Log in to the CDM console, 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.

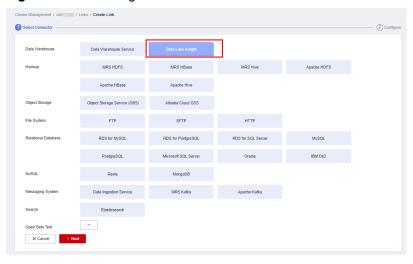
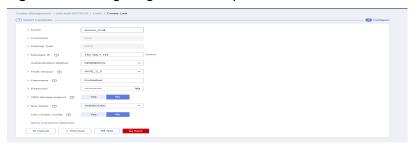


Figure 2-5 Selecting the DLI connector

iii. Configure the connection parameters.

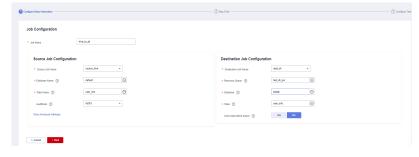
Figure 2-6 Configuring connection parameters



After the configuration is complete, click **Save**.

- 2. Create a CDM migration job.
 - a. Log in to the CDM console, 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, choose the **Table/File Migration** tab and click **Create Job**.
 - c. On the **Create Job** page, specify job information.

Figure 2-7 Configuring the CDM job



i. **Job Name**: Name of the data migration job, for example, **hive_to_dli**

ii. Set parameters required for **Source Job Configuration**.

Table 2-3 Source job configuration parameters

Parameter	Value
Source Link Name	Select the name of the data source created in 1.a.
Database Name	Select the name of the MRS Hive database you want to migrate to DLI. For example, the default database.
Table Name	Select the name of the Hive table. In this example, a database created on DLI and the user_info table are selected.
readMode	In this example, HDFS is selected.
	Two read modes are available: HDFS and JDBC. By default, the HDFS mode is used. If you do not need to use the WHERE condition to filter data or add new fields on the field mapping page, select the HDFS mode.
	The HDFS mode shows good performance, but in this mode, you cannot use the WHERE condition to filter data or add new fields on the field mapping page.
	The JDBC mode allows you to use the WHERE condition to filter data or add new fields on the field mapping page.

For details about parameter settings, see From Hive.

iii. Set parameters required for **Destination Job Configuration**.

Table 2-4 Destination job configuration parameters

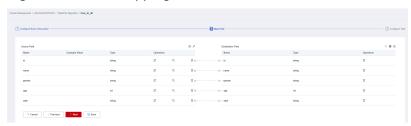
Parameter	Value
Destination Link Name	Select the DLI data source connection created in 1.b.
Resource Queue	Select a created DLI SQL queue.
Database	Select a created DLI database. In this example, database testdb created in Create a database and table on DLI is selected.
Table	Select the name of a table in the database. In this example, table user_info created in Create a database and table on DLI is created.

Parameter	Value
Clear data before import	Whether to clear data in the destination table before data import. In this example, set this parameter to No .
	If this parameter is set to Yes , data in the destination table will be cleared before the task is started.

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.
 - If the field mapping is incorrect, you can drag the fields to adjust the mapping.
 - If the type is automatically created at the migration destination, you need to configure the type and name of each field.
 - CDM supports field conversion during the migration. For details, see Field Conversion.

Figure 2-8 Field mapping



4. Click **Next** and set task parameters. Generally, retain the default values of all parameters.

In this step, you can configure the following optional functions:

- 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 details about 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: Specify 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 so that 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 2-9 Job progress and execution result

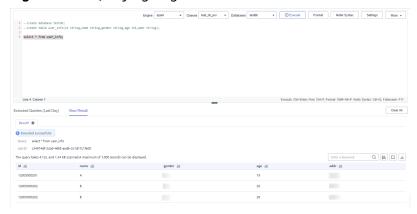


Step 3: Query Results

After the migration job is complete, log in to the DLI management console and click **SQL Editor**. In the displayed page, set **Engine** to **spark**, **Queue** to the created SQL queue, and **Database** to the database created in **a**. Execute the following query statement and check whether the Hive table data has been migrated to the **user_info** table:

select * from user_info;

Figure 2-10 Querying migrated data



2.3 Migrating Data from MRS Kafka to DLI

This section describes how to use the CDM data synchronization function to migrate data from MRS Kafka to DLI.

Prerequisites

• You have created a DLI SQL queue. For details about how to create a DLI queue, see **Creating a Queue**.



When you create a queue, set its **Type** to **For SQL**.

- You have created an MRS security cluster that contains the Kafka component.
 For details about 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 details about how to create a cluster, see Creating a CDM Cluster.

□ NOTE

- If the destination data source is an on-premises database, you need the Internet or
 Direct Connect. When using the Internet, ensure that an EIP has been bound to the
 CDM cluster, the security group of CDM allows outbound traffic from the host
 where the off-cloud 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 MRS or GaussDB(DWS), the network must meet the following requirements:
 - i. If the CDM cluster and the cloud service are in different regions, a public network or a dedicated connection is required for enabling communication between the CDM cluster and the cloud service. If the Internet is used for communication, ensure 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, VPC, subnet, and security group, they 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 rules and security group rules.

For details about how to configure routes, see **7. Configure routes**. For details about how to configure security groups, see section **Security Group Configuration Examples**.

iii. The cloud service instance and the CDM cluster belong to the same enterprise project. If they do not, you can modify the enterprise project of the workspace.

In this example, the VPC, subnet, and security group of the CDM cluster are the same as those of the MRS cluster.

Step 1: Prepare Data

- Create a Kafka topic for the MRS cluster and send messages to the topic.
 - Log in to MRS Manager by referring to Accessing FusionInsight Manager.
 - On MRS Manager, click **System** in the top navigation pane. On the page displayed, choose **Permission** > **User** from the left navigation pane. On the displayed page, configure 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 password of the current user and enter it again.
 - iv. User Group and Primary Group: Select kafkaadmin.
 - v. Role: Select Manager_viewer.

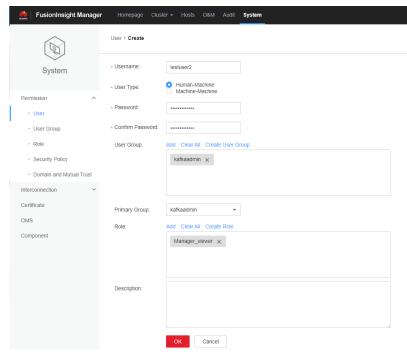


Figure 2-11 Creating a Kafka user

- c. On the MRS Manager console, choose Cluster > Name of the desired cluster > Service > ZooKeeper > Instance. On the displayed page, obtain the IP address of the ZooKeeper instance.
- d. On the MRS Manager console, 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 Kafka client is installed in the **/opt/kafkaclient** directory on the active MRS node.
- f. Go to the client installation directory as user **root**.
 - Example command: cd /opt/kafkaclient
- g. Run the following command to set environment variables:

source bigdata_env

h. Run the following command to authenticate the user created in **b** since Kerberos authentication has been enabled for the cluster:

kinit <Username in b>

Example command: kinit testuser2

- i. Run the following command to 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 deployed is that of the ZooKeeper instance obtained in **c**.
- j. Run the following command to 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 deployed in that 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 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 editing window to create a database. The following example creates the migrated DLI database **testdb**. For details about the syntax for creating a DLI database, see **Creating a Database**.

create database testdb;

 Create a table in the database. For details about the table creation syntax, see Creating a DLI Table Using the DataSource Syntax. CREATE TABLE testdlitable(value STRING);

Step 2: Migrate Data

- Create a CDM connection to MRS Kafka.
 - a. Create a connection to link CDM to the data source MRS Kafka.
 - Log in to the CDM console, 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**.

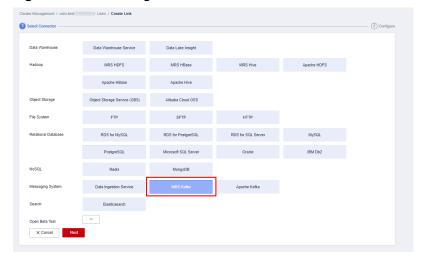


Figure 2-12 Selecting the MRS Kafka connector

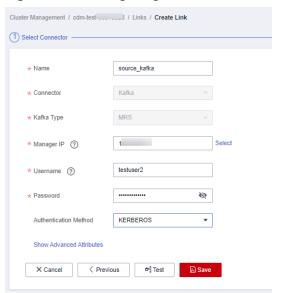
iii. Configure the connection. The following table describes the required parameters.

Table 2-5 MRS Kafka connection configurations

Parameter	Value
Name	Name of the MRS Kafka data source, for example, source_kafka.
Manager IP	Manager IP address of the cluster. The value is automatically specified after you click Select next to the text box and select the MRS Kafka cluster.
Username	Name of the MRS Kafka user created in b .
Password	Password of the MRS Kafka user.
Authenticatio n Method	KERBEROS if Kerberos authentication is enabled for the MRS cluster; SIMPLE if the MRS cluster is a common cluster
	In this example, set this parameter to KERBEROS .

For more details about the parameters, see Link to Kafka.

Figure 2-13 Configuring the MRS Kafka connection



- iv. Click **Save** to complete the configuration.
- b. Create a connection to link CDM to DLI.
 - Log in to the CDM console, 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.

Figure 2-14 Selecting the DLI connector

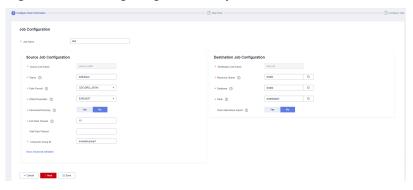
iii. Configure the connection parameters. For details about parameter settings, see **Link to DLI**.

Figure 2-15 Configuring connection parameters



- iv. After the configuration is complete, click **Save**.
- 2. Create a CDM migration job.
 - Log in to the CDM console, 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, choose the **Table/File Migration** tab and click **Create Job**.
 - c. On the **Create Job** page, specify job information.

Figure 2-16 Configuring the CDM job



- i. Job Name: Name of the data migration job, for example, test
- ii. Set parameters required for **Source Job Configuration**.

Table 2-6 Source job configuration parameters

Parameter	Value
Source Link Name	Select the name of the data source created in 1.a.
Topics	Name of the topics you want to migrate to DLI. You can select one or more topics. Example: kafkatopic.
Data Format	Select the message format as needed. In this example, CDC (DRS_JSON) 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 EARLIEST . Available values are as follows:
	Latest: Maximum offset, indicating that the latest data will be extracted
	Earliest: Minimum offset, indicating that the earliest data will be extracted
	Submitted: Data that has been submitted
	Time Range: Data within a specified time range
Permanent Running	Whether a job runs permanently. In this example, set this parameter to No .
Pull Data Timeout	Maximum minutes allowed for a continuous data pulling. In this example, set this parameter to 15 .
Wait Data Timeout	(Optional) Maximum seconds allowed for waiting data reading. In this example, leave this parameter empty.
Consumer Group ID	Consumer group ID. The default Kafka message group ID example-group1 is used.

For details about parameter settings, see From Apache Kafka.

iii. Set parameters required for **Destination Job Configuration**.

Parameter	Value
Destination Link Name	Select the DLI data source connection created in 1.b.
Resource Queue	Select a created DLI SQL queue.
Database	Select a created DLI database. In this example, database testdb created in Create a database and table on DLI is selected.
Table	Select the name of a table in the database. In this example, table testdlitable created in Create a database and table on DLI is selected.
Clear data before import	Whether to clear data in the destination table before data import. In this example, set this parameter to No .
	If this parameter is set to Yes , data in the destination table will be cleared before the task is started.

Table 2-7 Destination job configuration 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.
 - If the field mapping is incorrect, you can drag the fields to adjust the mapping.
 - If the type is automatically created at the migration destination, you need to configure the type and name of each field.
 - CDM supports field conversion during the migration. For details, see Field Conversion.

Figure 2-17 Field mapping



4. Click **Next** and set task parameters. Generally, retain the default values of all parameters.

In this step, you can configure the following optional functions:

- 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 details about 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: Specify 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. In this example, retain the default value No so that 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 2-18 Job progress and execution result



Step 3: Query Results

After the migration job is complete, log in to the DLI management console and click **SQL Editor**. In the displayed page, set **Engine** to **spark**, **Queue** to the created SQL queue, and **Database** to the database created in **a**. Execute the following query statement and check whether the Kafka table data has been migrated to the **testdlitable** table:

select * from testdlitable;

2.4 Migrating Data from Elasticsearch to DLI

This section describes how to use the CDM data synchronization function to migrate data from a CSS Elasticsearch cluster to DLI. Data in a self-built Elasticsearch cluster can also be bidirectionally synchronized between CDM and DLI.

Prerequisites

• You have created a DLI SQL queue. For details about how to create a DLI queue, see **Creating a Queue**.



When you create a queue, set its Type to For SQL.

- You have created a CSS Elasticsearch cluster. For details about how to create a CSS cluster, see Creating an Elasticsearch Cluster in Non-Security Mode.
 In this example, the version of the created CSS cluster is 7.6.2, and security mode is disabled for the cluster.
- You have created a CDM cluster. For details about how to create a CDM cluster, see Creating a CDM Cluster.

- If the destination data source is an on-premises database, you need the Internet or
 Direct Connect. When using the Internet, ensure that an EIP has been bound to the
 CDM cluster, the security group of CDM allows outbound traffic from the host
 where the off-cloud 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 CSS on a cloud, the network must meet the following requirements:
 - i. If the CDM cluster and the cloud service are in different regions, a public network or a dedicated connection is required for enabling communication between the CDM cluster and the cloud service. If the Internet is used for communication, ensure 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, VPC, subnet, and security group, they 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 rules and security group rules.

For details about how to configure routes, see **Configure routes**. For details about how to configure security groups, see section **Security Group Configuration Examples**.

iii. The cloud service instance and the CDM cluster belong to the same enterprise project. If they do not, you can modify the enterprise project of the workspace.

In this example, the VPC, subnet, and security group of the CDM cluster are the same as those of the CSS cluster.

Step 1: Prepare Data

- Create an index for the CSS cluster and import data.
 - Log in to the CSS management console and choose Clusters > Elasticsearch from the navigation pane on the left.
 - On the Clusters page, click Access Kibana in the Operation column of the created CSS cluster.
 - c. In the navigation pane of Kibana, choose **Dev Tools**. The **Console** page is displayed.
 - 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": {
                 "type": "text",
                 "analyzer": "ik_smart"
        },
        "size": {
                 "type": "keyword"
        }
    }
}
```

e. Run the following command to 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":"2018 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 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 editing window to create a database, for example, the migrated DLI database **testdb**: For details about the syntax for creating a DLI database, see **Creating a Database**. create database testdb;

 Create a table in the database. For details about 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 to MRS Hive.
 - a. Create a connection to link CDM to the data source CSS.
 - Log in to the CDM console, 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.

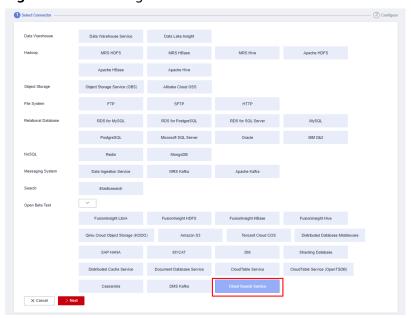


Figure 2-19 Selecting the CSS connector

iii. Configure the connection. The following table describes the required parameters. For details about parameter settings, see Link to Elasticsearch/CSS.

Table 2-8 CSS data source configuration

Parameter	Value.
Name	Name of the CSS data source, for example, source_css.
Elasticsearch Server List	Click Select next to the text box and select the CSS cluster. The Elasticsearch server list is automatically displayed.
Security mode Authenticatio n	If you have enabled the security mode for the CSS cluster, set this parameter to Yes . Otherwise, set this parameter to No . In this example, set this parameter to No .

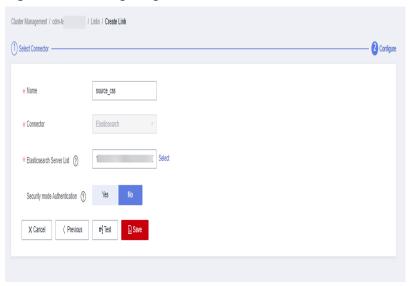


Figure 2-20 Configuring the CSS connection

- iv. Click **Save** to complete the configuration.
- b. Create a connection to link CDM to DLI.
 - Log in to the CDM console, 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 Data Lake Insight and click Next.

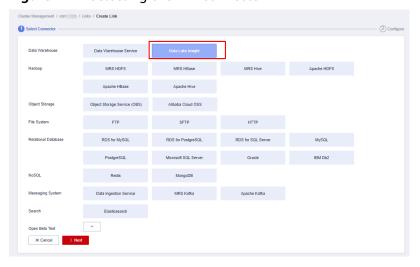


Figure 2-21 Selecting the DLI connector

ii. Configure the connection parameters. For details about parameter settings, see **Link to DLI**.

Figure 2-22 Configuring connection parameters



- iv. After the configuration is complete, click Save.
- 2. Create a CDM migration job.
 - Log in to the CDM console, 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, choose the **Table/File Migration** tab and click **Create Job**.
 - c. On the **Create Job** page, specify job information.

Figure 2-23 Configuring the CDM job



- i. Job Name: Name of the data migration job, for example, css_to_dli
- ii. Set parameters required for Source Job Configuration.

Table 2-9 Source job configuration parameters

Parameter	Value
Source Link Name	Select the name of the data source created in 1.a.
Index	Select the Elasticsearch index created for the CSS cluster. In this example, the my_test index created in Create an index for the CSS cluster and import data is used.
	The index can contain only lowercase letters.
Туре	Elasticsearch type, which is similar to the table name of a relational database. The type name can contain only lowercase letters. Example: _doc.

For details about other parameters, see From Elasticsearch or CSS.

iii. Set parameters required for **Destination Job Configuration**.

Table 2-10 Destination job configuration parameters

Parameter	Value
Destination Link Name	Select the DLI data source connection created in 1.b.
Resource Queue	Select a created DLI SQL queue.
Database	Select a created DLI database. In this example, database testdb created in Create a database and table on DLI is selected.
Table	Select the name of a table in the database. In this example, table tablecss created in Create a database and table on DLI is created.
Clear data before import	Whether to clear data in the destination table before data import. In this example, set this parameter to No .
	If this parameter is set to Yes , data in the destination table will be cleared before the task is started.

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.
 - If the field mapping is incorrect, you can drag the fields to adjust the mapping.
 - If the type is automatically created at the migration destination, you need to configure the type and name of each field.
 - CDM supports field conversion during the migration. For details, see Field Conversion.

Figure 2-24 Field mapping



4. Click **Next** and set task parameters. Generally, retain the default values of all parameters.

In this step, you can configure the following optional functions:

 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 details about 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: Specify 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 so that 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 2-25 Job progress and execution result

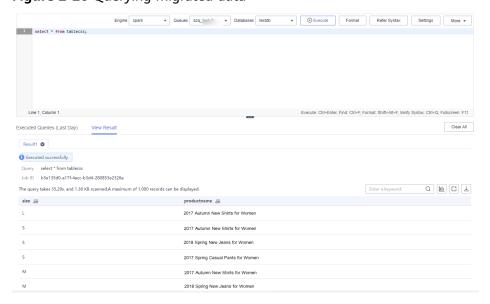


Step 3: Query Results

After the migration job is complete, log in to the DLI management console and click **SQL Editor**. In the displayed page, set **Engine** to **spark**, **Queue** to the created SQL queue, and **Database** to the database created in **a**. Execute the following query statement and check whether the CSS table data has been migrated to the **tablecss** table:

select * from tablecss;

Figure 2-26 Querying migrated data



2.5 Migrating Data from RDS to DLI

This section describes how to use the CDM data synchronization function to migrate data from an RDS DB instance to DLI. Data in other relational databases can also be bidirectionally synchronized between CDM and DLI.

Prerequisites

• You have created a DLI SQL queue. For details about how to create a DLI queue, see **Creating a Queue**.

<u>A</u> CAUTION

When you create a queue, set its Type to For SQL.

- You have created an RDS for MySQL DB instance. For details about how to create an RDS cluster, see Buy a DB 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 details about how to create a CDM cluster, see Creating a CDM Cluster.

Ⅲ NOTE

- If the destination data source is an on-premises database, you need the Internet or Direct Connect. When using the Internet, ensure that an EIP has been bound to the CDM cluster, the security group of CDM allows outbound traffic from the host where the off-cloud 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 on a cloud, the network must meet the following requirements:
 - i. If the CDM cluster and the cloud service are in different regions, a public network or a dedicated connection is required for enabling communication between the CDM cluster and the cloud service. If the Internet is used for communication, ensure 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, VPC, subnet, and security group, they 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 rules and security group rules.

For details about how to configure routes, see **Configure routes**. For details about how to configure security groups, see section **Security Group Configuration Examples**.

iii. The cloud service instance and the CDM cluster belong to the same enterprise project. If they do not, you can modify the enterprise project of the workspace.

In this example, the VPC, subnet, and security group of the CDM cluster are the same as those of the RDS for MySQL DB instance.

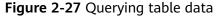
Step 1: Prepare Data

- Create databases and tables on the RDS for MySQL DB instance.
 - a. Log in to the RDS console. On the displayed page, locate the target DB instance and choose **More** > **Log In** in the **Operation** column.
 - b. On the displayed login page, 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 of rest parameters. Then, click **OK**.
 - d. In the **Operation** column of row where the created database locates, click **SQL Window** and enter the following statement to create a table:

e. Run the following statements to insert data to the created table: insert into tabletest VALUES ('123','abc');

```
insert into tabletest VALUES ('456','efg');
insert into tabletest VALUES ('789','hij');
```

f. Run the following statement to query table data: select * from tabletest;





- Create a database and table on DLI.
 - a. 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 editing window to create a database, for example, the migrated DLI database **testdb**: For details about the syntax for creating a DLI database, see **Creating a Database**.

create database testdb;

b. In SQL Editor, select testdb for Database and run the following table creation statement to create a table in the database. For details about the table creation syntax, see Creating a DLI Table Using the DataSource Syntax.

create table tabletest(id string,name string);

Step 2: Migrate Data

- Create a CDM connection to MRS Hive.
 - a. Create a connection to the RDS database.

- Log in to the CDM console, 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 crating a connection to RDS for MySQL, upload the MySQL driver. Choose the **Links** tab and click **Driver Management**. The **Driver Management** page is displayed.
- Download the MySQL driver to your local PC by referring to Managing Drivers and decompress the driver package to obtain the JAR file.
 - For example, download the **mysql-connector-java-5.1.48.zip** package and decompress 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** to upload the driver file obtained in 1.a.iii.
- v. On the **Driver Management** page, click **Back** to return to the **Links** tab. Click **Create Link**, select **RDS for MySQL**, and click **Next**.
- vi. Configure the connection. The following table describes the required parameters.

Table 2-11 Connection parameters

Parameter	Value
Name	Name of the RDS data source, for example, source_rds
Database Server	Click Select 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 entered after you select the database server.
Database Name	Name of the RDS DB instance you want to migrate. The testrdsdb database created in c is used in this example.
Username	Username used for accessing the database. This account must have the permissions required to read and write data tables and metadata.
	In this example, the default user root for creating the RDS for MySQL DB instance is used.
Password	Password of the user.

For other parameters, retain the default values. For details, see **Link to Relational Databases**. Click **Save** to complete the configuration.

Figure 2-28 Configuring the connection to the RDS for MySQL DB instance

- b. Create a connection to the DLI.
 - Log in to the CDM console, 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 Data Lake Insight and click Next.

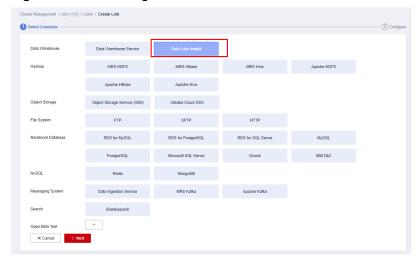


Figure 2-29 Selecting the DLI connector

 Create a connection to link CDM to DLI. For details about parameter settings, see Link to DLI.

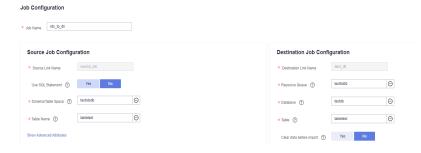
Figure 2-30 Selecting the DLI connector



After the configuration is complete, click Save.

- 2. Create a CDM migration job.
 - Log in to the CDM console, 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, choose the **Table/File Migration** tab and click **Create Job**.
 - c. On the **Create Job** page, specify job information.

Figure 2-31 Configuring the migration job



- i. Job Name: Name of the data migration job, for example, rds_to_dli
- ii. Set parameters required for **Source Job Configuration**.

Table 2-12 Source job configuration parameters

Parameter	Value	
Source Link Name	Select the name of the data source created in 1.a.	
Use SQL Statement	When Use SQL Statement is set to Yes , enter an SQL statement here. CDM exports data based on the SQL statement.	
	In this example, set this parameter to No .	
Schema/Table Space	Select the name of the RDS for MySQL DB instance you want to migrate to DLI. For example, the testrdsdb database.	

Parameter	Value
Table Name	Name of the table you want to migrate. In this example, use tabletest created in d .

For details about parameter settings, see **From PostgreSQL/SQL Server**.

iii. Set parameters required for **Destination Job Configuration**.

Table 2-13 Destination job configuration parameters

Parameter	Value	
Destination Link Name	Select the DLI data source connection.	
Resource Queue	Select a created DLI SQL queue.	
Database	Select a created DLI database. In this example, database testdb created in Create a database and table on DLI is selected.	
Table	Select the name of a table in the database. In this example, table tabletest created in Create a database and table on DLI is created.	
Clear data before import	Whether to clear data in the destination table before data import. In this example, set this parameter to No .	
	If this parameter is set to Yes , 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.
 - If the field mapping is incorrect, you can drag the fields to adjust the mapping.
 - If the type is automatically created at the migration destination, you need to configure the type and name of each field.
 - CDM supports field conversion during the migration. For details, see Field Conversion.

Figure 2-32 Field mapping



v. Click **Next** and set task parameters. Generally, retain the default values of all parameters.

In this step, you can configure the following optional functions:

- 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 details about 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: Specify 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 so that dirty data is not recorded.
- vi. Click **Save and Run**. On the **Job Management** page, you can view the job execution progress and result.

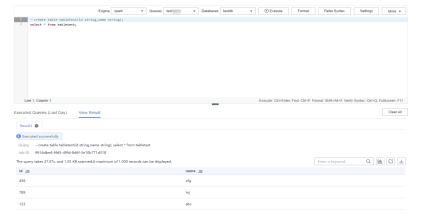
Figure 2-33 Job progress and execution result



Step 3: Query Results

After the migration job is complete, log in to the DLI management console and click **SQL Editor**. In the displayed page, set **Engine** to **spark**, **Queue** to the created SQL queue, and **Database** to the database created in **Create a database and table on DLI**. Execute the following query statement and check whether the table data has been migrated to the **tabletest** table: select * from tabletest;

Figure 2-34 Querying data in the table



2.6 Migrating Data from GaussDB(DWS) to DLI

This section describes how to use the CDM data synchronization function to migrate data from GaussDB(DWS) to DLI.

Prerequisites

• You have created a DLI SQL queue. For details about how to create a DLI queue, see **Creating a Queue**.

CAUTION

When you create a queue, set its **Type** to **For SQL**.

- You have created a GaussDB(DWS) cluster. For details about how to create a GaussDB(DWS) cluster, see Creating a Cluster.
- A CDM cluster has been created for migration. For details about how to create a CDM cluster, see **Creating a CDM Cluster**.

- If the destination data source is an on-premises database, you need the Internet or
 Direct Connect. When using the Internet, ensure that an EIP has been bound to the
 CDM cluster, the security group of CDM allows outbound traffic from the host
 where the off-cloud 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 on a cloud, the network must meet the following requirements:
 - i. If the CDM cluster and the cloud service are in different regions, a public network or a dedicated connection is required for enabling communication between the CDM cluster and the cloud service. If the Internet is used for communication, ensure 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, VPC, subnet, and security group, they 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 rules and security group rules.

For details about how to configure routes, see **Configure routes**. For details about how to configure security groups, see section **Security Group Configuration Examples**.

iii. The cloud service instance and the CDM cluster belong to the same enterprise project. If they do not, you can modify the enterprise project of the workspace.

In this example, the VPC, subnet, and security group of the CDM cluster are the same as 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**.

- Connect to the default database gaussdb of a 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 network address is used for connection, set this parameter to Public Network Address or Public Network Access Domain Name. If a private network address is used for connection, set this parameter to Private Network Address or Private Network Access Domain Name. For details, see Obtaining the Cluster Connection Address. If an ELB is used for connection, set this parameter to ELB Address.
 - dbadmin: Default administrator username used during cluster creation
 - -W: Default password of the administrator
- Run the following command to create the **testdwsdb** database: CREATE DATABASE testdwsdb;
- d. Run the following command to exit the **gaussdb** database and connect to **testdwsdb**:

```
\q
gsql -d testdwsdb -h Connection address of the GaussDB(DWS) cluster -U dbadmin -p 8000 -W
password -r
```

e. Run the following commands to create a table and import data to the table.

```
Run the following command to create a table:

CREATE TABLE table1 (id int, a char(6), b varchar(6),c varchar(6));

Run the following statements to 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 2-35 Querying data in the table

- Create a database and table on DLI.
 - a. 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 editing window to create a database, for example, the migrated DLI database **testdb**: For details about the syntax for creating a DLI database, see **Creating a Database**. create database testdb;

b. In **SQL Editor**, select **testdb** for **Database** and run the following table creation statement to create a table in the database: For details about

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 to MRS Hive.
 - a. Create a connection to the GaussDB(DWS) database.
 - Log in to the CDM console, 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. The following table describes the required parameters.

Table 2-14 GaussDB(DWS) data source configuration

Parameter	Value
Name	Name of the GaussDB(DWS) data source, for example, source_dws .
Database Server	Click Select next to the text box to select the name of the created GaussDB(DWS) cluster.
Port	Port number of the GaussDB(DWS) database. The default value is 8000 .
Database Name	Name of the GaussDB(DWS) database you want to migrate The testdwsdb database created in Create a database and table in the GaussDB(DWS) cluster is used in this example.
Username	Username used for accessing the database. This account must have the permissions required to read and write data tables and metadata.
	In this example, the default administrator dbadmin specified when you create the GaussDB(DWS) database is used.
Password	Password of the GaussDB(DWS) database user.

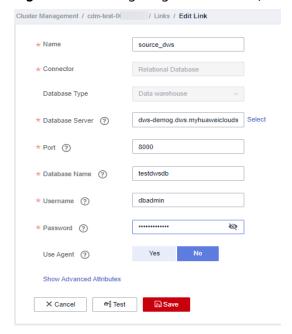


Figure 2-36 Configuring the GaussDB(DWS) connection

For other parameters, retain the default values. For details, see **Link to Relational Databases**. Click **Save** to complete the configuration.

- b. Create a connection to the DLI.
 - Log in to the CDM console, 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 Data Lake Insight and click Next.

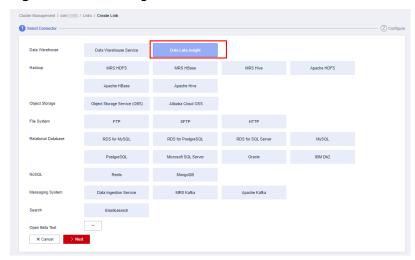


Figure 2-37 Selecting the DLI connector

 Create a connection to link CDM to DLI. For details about parameter settings, see Link to DLI.

Figure 2-38 Selecting the DLI connector



After the configuration is complete, click Save.

- 2. Create a CDM migration job.
 - Log in to the CDM console, 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, choose the **Table/File Migration** tab and click **Create Job**.
 - c. On the **Create Job** page, specify job information.

Figure 2-39 Configuring the migration job



- i. Job Name: Name of the data migration job, for example, test
- ii. Set parameters required for **Source Job Configuration**.

Table 2-15 Source job configuration parameters

Parameter	Value
Source Link Name	Select the name of the data source created in 1.a.
Use SQL Statement	When Use SQL Statement is set to Yes , enter an SQL statement here. CDM exports data based on the SQL statement. In this example, set this parameter to No .

Parameter	Value	
Schema/Table Space	Name of the schema or tablespace from which data will be extracted. This parameter is displayed when Use SQL Statement is set to No . Click the icon next to the text box to go to the page for selecting a schema or directly enter a schema or tablespace.	
	In this example, no schema is created in Create a database and table in the GaussDB(DWS) cluster. In this case, set this parameter to the default value public.	
	If the desired schema or tablespace is not displayed, check whether the login account has the permissions required to query metadata.	
	NOTE The parameter value can contain wildcard characters (*), which is used to export all databases whose names start with a certain prefix or end with a certain suffix. For example:	
	SCHEMA* indicates that all databases whose names starting with SCHEMA are exported.	
	*SCHEMA indicates that all databases whose names ending with SCHEMA are exported.	
	SCHEMA indicates that all databases whose names containing SCHEMA are exported.	
Table Name	Name of the table you want to migrate. In this example, table1 created in Create a database and table in the GaussDB(DWS) cluster is used.	

For details about parameter settings, see **From a Relational Database**.

iii. Set parameters required for **Destination Job Configuration**.

Table 2-16 Destination job configuration parameters

Parameter	Value	
Destination Link Name	Select the DLI data source connection.	
Resource Queue	Select a created DLI SQL queue.	
Database	Select a created DLI database. In this example, database testdb created in Create a database and table on DLI is selected.	

Parameter	Value
Table	Select the name of a table in the database. In this example, table tabletest created in Create a database and table on DLI is created.
Clear data before import	Whether to clear data in the destination table before data import. In this example, set this parameter to No .
	If this parameter is set to Yes , 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.
 - If the field mapping is incorrect, you can drag the fields to adjust the mapping.
 - If the type is automatically created at the migration destination, you need to configure the type and name of each field.
 - CDM supports field conversion during the migration. For details, see Field Conversion.

Figure 2-40 Field mapping



v. Click **Next** and set task parameters. Generally, retain the default values of all parameters.

In this step, you can configure the following optional functions:

- 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 details about 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: Specify 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** so that dirty data is not recorded.

vi. Click **Save and Run**. On the **Job Management** page, you can view the job execution progress and result.

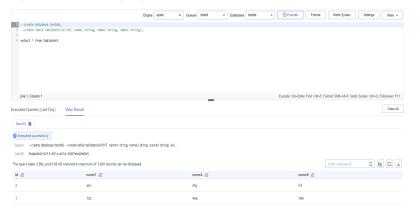
Figure 2-41 Job progress and execution result



Step 3: Query Results

After the migration job is complete, log in to the DLI management console and click **SQL Editor**. In the displayed page, set **Engine** to **spark**, **Queue** to the created SQL queue, and **Database** to the database created in **Create a database and table on DLI**. Execute the following query statement and check whether the table data has been migrated to the **tabletest** table: select * from tabletest;

Figure 2-42 Querying data in the table



3 Data Analysis

3.1 Analyzing Driving Behavior Data

Application Scenarios

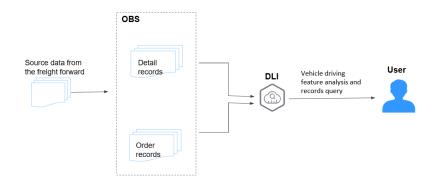
Cloud computing and big data provide companies with data analysis and mining capabilities required in the Internet of Vehicle (IoV) field, helping companies or department of motor vehicles manage and analyze vehicle and driving behavior data quickly and scientifically.

Solution Architecture

DLI can query the records of vehicle driving features based on the detail records and freight order data regularly reported by the freight forwarder.

Data Types describes the data types used by DLI to record the data.

Figure 3-1 Solution Overview



Process

To use DLI to analyze driving behavior data, perform the following steps:

Step 1: Uploading Data. Upload the data to OBS.

Step 2: Analyzing Data. Use DLI to query the data.

Example Code

Download the **data package** for sample data and detailed SQL statements.

Solution Advantages

- Free of data migration: DLI can interconnect with multiple data sources. You only need to create SQL tables and map data sources.
- Easy to use: You can use standard SQL statements to compile metric analysis logic without paying attention to the complex distributed computing platform.
- Pay-per-use: Log analysis is scheduled periodically based on time-critical requirements. There is a long idle period between every two scheduling operations. DLI uses the pay-per-use billing mode, which effectively reduces your costs.

Resource Planning and Costs

Table 3-1 Resource planning and costs

Resource	Description	Cost
OBS	You need to create an OBS bucket and upload data to OBS for data analysis using DLI.	You will be charged for using the following OBS resources:
		• Storage Fee for storing static website files in OBS.
		 Request Fee for accessing static website files stored in OBS.
		 Traffic Fee for using a custom domain name to access OBS over the public network.
		The actual fee depends on the size of the stored file, the number of user access requests, and the traffic volume. Estimate the fee based on your service requirements.
DLI	Before creating a SQL job, you need to purchase a queue. When	For example, if you purchase a pay-peruse queue, you will be billed based on the number of CUHs used by the queue.
	using queue resources, you are billed based on the CUH of the queue.	Usage is billed by the hour. For example, 58 minutes of usage will be rounded to the hour. CUH pay-per-use billing = Unit price x Number of CUs x Number of hours.

Data Types

Detail records

Detail records include the regularly reported location records and data of alarms triggered by abnormal driving behavior.

Table 3-2 Detail records

Field	Data Type	Description
driverID	string	Driver ID
carNumber	string	License plate number
latitude	double	Latitude
longitude	double	Longitude
speed	int	Speed
direction	int	Direction
siteName	string	Site name
time	timestamp	Report time of the records
isRapidlySpeedup	int	Whether the vehicle rapidly speeds up. 1 indicates that the vehicle suddenly speeds up, and 0 indicates that the vehicle does not.
isRapidlySlowdown	int	Whether the vehicle suddenly slows down.
isNeutralSlide	int	Whether the vehicle is coasting.
isNeutralSlideFinished	int	Whether vehicle coasting has stopped.
neutralSlideTime	bigint	Time length of vehicle coasting.
isOverspeed	int	Whether the vehicle is speeding.
isOverspeedFinished	int	Whether the vehicle stops speeding.
overspeedTime	bigint	Duration of the vehicle speeding
isFatigueDriving	int	Whether fatigue driving occurs.

Field	Data Type	Description
isHthrottleStop	int	Whether the driver revs the engine in neutral.
isOilLeak	int	Abnormal oil consumption

Order data

Order data refers to the records of freight orders.

Table 3-3 Order data

Field	Data Type	Description
orderNumber	string	Order ID
driverID	string	Driver ID
carNumber	string	License plate number
customerID	string	Customer ID
sourceCity	string	Departure
targetCity	string	Destination
expectArriveTime	timestamp	Expected delivery time
time	timestamp	Time when a record is generated.
action	string	Event type, including creating an order, dispatching goods, delivering packages, and signing orders.

Step 1: Uploading Data

Upload the data to OBS for data analysis using DLI.

- Download OBS Browser+. For details about the download address, see Object Storage Service Tool Guide.
- 2. Install OBS Browser+. For details about the installation procedure, see **Object Storage Service Tool Guide**.
- Log in to OBS Browser+. OBS Browser+ supports two login modes: AK login (using access keys) or authorization code login. For details about the login procedure, see Object Storage Service Tool Guide.
- 4. Upload data using the OBS browser+.
 - Start the OBS Browser+, click **Create Bucket** on the homepage. Select a region and enter a bucket name (for example, **DLI-demo**). After the bucket is created, return to the bucket list and click **DLI-demo**. OBS Browser+ supports

upload by dragging. You can drag one or more files or folders from a local path to the object list of a bucket or a parallel file system on OBS Browser+. You can even drag a file or folder directly to a specified folder on OBS Browser+.

Obtain the test data by downloading the **Best_Practice_01.zip** file and decompressing it. Perform the following operations:

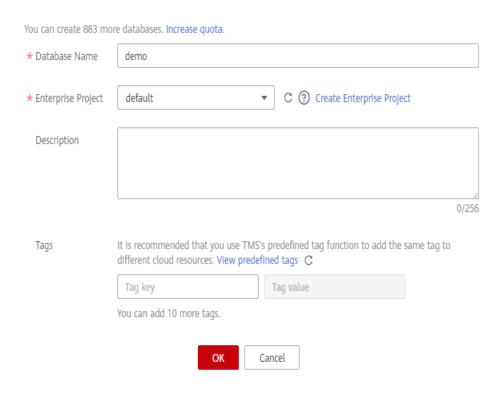
- Detail records: Upload the detail-records folder in the Data directory to the root directory of the OBS bucket.
- Order data: Upload the order-records folder in the Data directory to the root directory of the OBS bucket.

Step 2: Analyzing Data

Use DLI to query the data for analysis.

- 1. Creating a Database and a Table
 - a. On the homepage of the management console, choose Service List >
 Analytics > Data Lake Insight.
 - b. On the DLI console, click **SQL Editor**.
 - c. In the left pane of the SQL Editor, select the **Databases** tab and click to create the **demo** database.

Figure 3-2 Creating a database Create Database



NOTE

Database Name cannot be set to **default** because **default** is the built-in database.

d. Choose the **demo** database, and enter the following SQL statement in the editing box:

```
create table detail_records(
 driverID String,
 carNumber String,
 latitude double.
 longitude double,
 speed int,
 direction int,
 siteName String,
 time timestamp,
 isRapidlySpeedup int,
 isRapidlySlowdown int,
 isNeutralSlide int,
 isNeutralSlideFinished int,
 neutralSlideTime long,
 isOverspeed int,
 isOverspeedFinished int,
 overspeedTime long,
 isFatigueDriving int,
 isHthrottleStop int,
 isOilLeak int) USING CSV OPTIONS (PATH 'obs://dli-demo/detail-records/');
```

□ NOTE

Replace the file path in the preceding statement with the actual OBS path where the detail records are stored.

e. Click Execute to create the detail_records table. See Figure 3-3.

Figure 3-3 Creating the detail_records table



f. Run the following SQL statements to create the **event_records** table in the **demo** database. The operation is similar to **1.d** and **1.e**.

```
create table event_records(
 driverID String,
 carNumber String,
 latitude double,
 longitude double,
 speed int,
 direction int,
 siteName String,
 time timestamp,
 isRapidlySpeedup int,
 isRapidlySlowdown int,
 isNeutralSlide int,
 isNeutralSlideFinished int,
 neutralSlideTime long,
 isOverspeed int,
 isOverspeedFinished int,
 overspeedTime long,
 isFatigueDriving int,
 isHthrottleStop int,
 isOilLeak int)
```

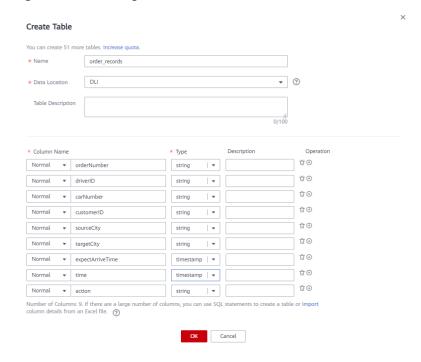
g. Run the following SQL statements to extract the alarm and event data from the detail records and insert it into the **event records** table.

```
insert into table event_records
(select *
from detail_records
where isRapidlySpeedup > 0
OR isRapidlySlowdown > 0
OR isNeutralSlide > 0
OR isNeutralSlideFinished > 0
OR isOverspeed > 0
OR isOverspeedFinished > 0
OR isFatigueDriving > 0
OR isHthrottleStop > 0
OR isOilLeak > 0)
```

h. Use another method to create the **order_records** table.

On the left of the SQL job editor, click the **Databases** tab and click the demo database. Click the plus icon (+) on the right of **Table** to create a table, and set **Data Location** to **DLI**. Set the column types according to **Order data**.

Figure 3-4 Creating the order_records table



i. Import the OBS data to the order_records table. Choose Data Management > Databases and Tables. Click the demo database to go to the table management page. In the Operation column of the order_records table, choose More > Import. Set File Format to CSV, the data storage path to obs://DLI-demo/order-records/, and retain default values for the rest parameters. Click OK.

The default timestamp format is **yyyy-MM-dd HH:mm:ss**. To use other formats, select **Advanced Settings** and enter the desired timestamp format (not modified in this example).

Import Data Database Name demo Table Name order_records * File Format CSV Set the options in Advanced Settings as required. DLI supports the read of CSV data that is not compressed or compressed by gzip. default * Path obs://DLI-demo/order-records/ Advanced Settings Cancel

Figure 3-5 Importing table data

2. Querying Data

a. Run the following SQL statements to query the alarm events of all drivers in a certain time period.

You can save the frequently-used query statements as a template by clicking **More** > **Save as Template** in the upper right corner of the editing window. The template is available for future use or can be modified in the SQL editor again.

Choose **Job Templates** > **SQL Templates** and click the **Custom Templates** tab. In the **Operation** column of the target template, click **Execute** to switch to the SQL editor. You can modify it as needed.

```
select
 driverID.
 carNumber,
 sum(isRapidlySpeedup) as rapidlySpeedupTimes,
 sum(isRapidlySlowdown) as rapidlySlowdownTimes,
 sum(isNeutralSlide) as neutralSlideTimes,
 sum(neutralSlideTime) as neutralSlideTimeTotal,
 sum(isOverspeed) as overspeedTimes,
 sum(overspeedTime) as overspeedTimeTotal,
 sum(isFatigueDriving) as fatigueDrivingTimes,
 sum(isHthrottleStop) as hthrottleStopTimes,
 sum(isOilLeak) as oilLeakTimes
from
 event_records
where
 time >= "2017-01-01 00:00:00"
 and time <= "2017-02-01 00:00:00"
group by
 driverID
 carNumber
order by
 rapidlySpeedupTimes desc,
 rapidlySlowdownTimes desc,
 neutralSlideTimes desc,
 neutralSlideTimeTotal desc,
 overspeedTimes desc,
 overspeedTimeTotal desc,
 fatigueDrivingTimes desc,
 hthrottleStopTimes desc,
 oilLeakTimes desc
```

In the query result, click to view graphical results.

- Set **Graph Type** to the bar chart.
- Set **X-AXIS** to **driverID**.
- Set Y-AXIS to rapidlySpeedupTimes.
- Set **Results** to **10**.

The command output is as follows:

Figure 3-6 Rapid acceleration



b. Run the following SQL statement to query the detailed record of a driver in a certain time period.

```
select

*
from
event_records
where
driverID = "panxian1000005"
and time >= "2017-01-01 00:00:00"
and time <= "2017-02-01 00:00:00"
```

In the query result, click to view graphical results.

- Set Graph Type to the bar chart.
- Set X-AXIS to driverID.
- Set **Y-AXIS** to **speed**.
- Set Results to 10.

The command output is as follows:

Figure 3-7 Speeding record



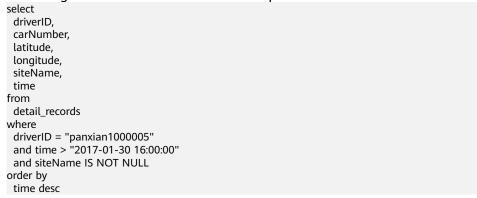
c. Run the following SQL statement to guery the order information.



Figure 3-8 Order information



d. Run the following SQL statement to query a vehicle's driving feature according to the driver ID and time of departure.



In the query result, click to view graphical results.

- Set **Graph Type** to the bar chart.
- Set X-AXIS to time.
- Set Y-AXIS to latitude.
- Set **Results** to **10**.

The command output is as follows:

Figure 3-9 Driving information

3.2 Converting Data Format from CSV to Parquet

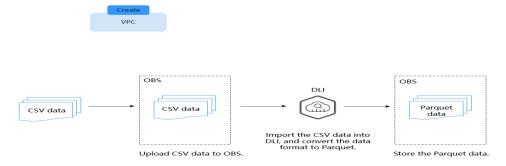
Application Scenarios

Parquet is a columnar storage substrate created for simpler data analysis. This format can speed up queries by allowing only the required columns to be read and calculated. In addition, Parquet is built to support efficient compression schemes, which maximizes the storage efficiency on disks. Using DLI, you can easily convert data format form CSV to Parquet.

Solution Overview

Upload CSV data to an OBS bucket, convert CSV data into Parquet data with DLI, and store the converted Parquet data to OBS.

Figure 3-10 Solution overview



Process

To use DLI to convert CSV data into Parquet data, perform the following steps:

Step 1: Creating and Uploading Data. Upload data to your OBS bucket.

Step 2: Using DLI to Convert CSV Data into Parquet Data. Import CSV data to DLI and convert it into Parquet data.

Solution Advantages

The query performance is improved.

If you have text-based data files or tables in an HDFS and are using Spark SQL to query data, converting data format to Parquet can improve the query performance by about 30 times (or more in some cases), despite of the time consumed during the conversion.

• Storage is saved.

Parquet is built to support efficient compression schemes, which maximizes the storage efficiency on disks. With Parquet, the storage cost can be reduced by about 75%.

Resource Planning and Costs

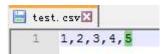
Table 3-4 Resource planning and costs

Resource	Description	Cost	
OBS	You need to create an OBS bucket and upload data to OBS for data analysis using DLI.	You will be charged for using the following OBS resources:	
		• Storage Fee for storing static website files in OBS.	
		 Request Fee for accessing static website files stored in OBS. 	
		 Traffic Fee for using a custom domain name to access OBS over the public network. 	
		The actual fee depends on the size of the stored file, the number of user access requests, and the traffic volume. Estimate the fee based on your service requirements.	
DLI	Before creating a SQL job, you need to purchase a queue. When	For example, if you purchase a pay-per- use queue, you will be billed based on the number of CUHs used by the queue.	
	using queue resources, you are billed based on the CUH of the queue.	Usage is billed by the hour. For example, 58 minutes of usage will be rounded to the hour. CUH pay-per-use billing = Unit price x Number of CUs x Number of hours.	

Step 1: Creating and Uploading Data

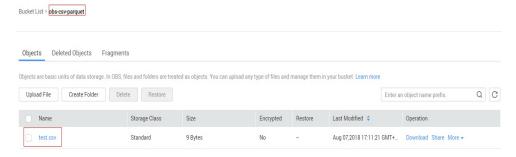
1. Create a CSV file. See **test.csv** in **Figure 3-11**.

Figure 3-11 Creating a test.csv file



2. In the OBS management console, create a bucket, name it **obs-csv-parquet**, and upload the **test.csv** file to the bucket.

Figure 3-12 Uploading CSV data to OBS



3. Create a bucket and name it **obs-parquet-data** to store the converted parquet data.

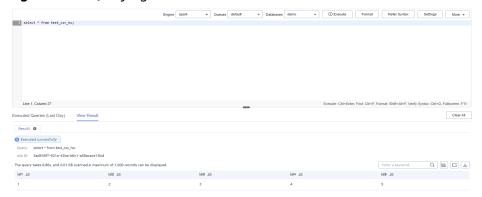
Step 2: Using DLI to Convert CSV Data into Parquet Data

- 1. Go to the DLI console, click **SQL Editor** in the navigation pane.
- 2. In the left pane of the SQL editor, click the **Databases** tab. Click 🕀, create a database, and name it **demo**.
- 3. In the SQL editing window, set **Engine** to **spark**, **Queue** to **default**, and **Database** to **demo**. Execute the following statement to create table **test_csv_hw** to import the data in the **test.csv** file from OBS.

```
create table test_csv_hw(id1 int, id2 int, id3 int, id4 int, id5 int)
using csv
options(
path 'obs://obs-csv-parquet/test.csv'
)
```

4. In the SQL editing window, query data in the **test csv hw** table.

Figure 3-13 Querying data



5. In the SQL job editing window, create a table to store the OBS data in Parquet format and name the table **test_parquet_hw**.

```
create table `test_parquet_hw` (`id1` INT, `id2` INT, `id3` INT, `id4` INT, `id5` INT) using parquet options ( path 'obs://obs-parquet-data/' )
```

Ⅲ NOTE

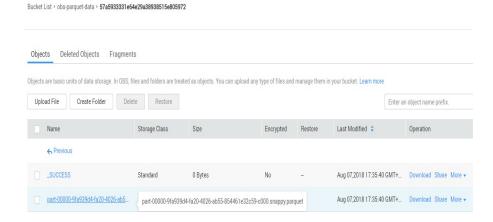
You do not need to specify a file because no Parquet file exists in this OBS bucket before the data is converted.

6. In the SQL editing window, execute the following statement to convert the CSV data to Parquet format and store the data in the specified OBS folder:

insert into test_parquet_hw select * from test_csv_hw

7. Check the result. OBS automatically created a file for saving the result.

Figure 3-14 Parquet data saved in a file in OBS



3.3 Analyzing E-commerce BI Reports

Application Scenarios

As a self-operated e-commerce company in China, the *X* mall has developed hundreds of millions of loyal users and accumulated massive amounts of authentic data while maintaining high-speed development. How to use the BI tool to find business opportunities from historical data is a key issue in the precision marketing of big data applications. It is also the core technology required for intelligent upgrade of all e-commerce platforms.

This case uses HUAWEI CLOUD DLI, GaussDB(DWS), and Yonghong BI to analyze data features of users and offerings based on the real user, product, and comment data (anonymized) of the mall, providing high-quality information for marketing decision-making, advertising recommendation, credit rating, brand monitoring, and user behavior prediction.

Process

To use DLI to analyze e-commerce data, perform the following steps:

Step 1: Uploading Data. Upload the data to OBS for data analysis using DLI.

Step 2: Analyzing Data. Use DLI to query the data for analysis.

Data Types

To protect user privacy and data security, all sampled data is anonymized.

User data

Table 3-5 User data

Field	Data Type	Description	Value
user_id	int	User ID	Anonymized
age	int	Age group	-1 indicates that the user age is unknown.
gender	int	Gender	0: Male1: Female2: Confidential
rank	Int	User level	Sequenced list of user level. The higher the user level, the larger the number.
register_tim e	string	User registration date	Unit: day

Product data

Table 3-6 Product data

Field	Data Type	Description	Value
product_id	int	Product No.	Anonymized
a1	int	Attribute 1	Enumerated value. The value -1 indicates unknown.
a2	int	Attribute 2	Enumerated value. The value -1 indicates unknown.
a3	int	Attribute 3	Enumerated value. The value -1 indicates unknown.
category	int	Category ID	Anonymized
brand	int	Brand ID	Anonymized

Comment data

Table 3-7 Comment data

Field	Data Type	Description	Value
deadline	string	End time	Unit: day

Field	Data Type	Description	Value
product_id	int	Product No.	Anonymized
comment_num	int	Segments of accumulated comment count	 0: No comment 1: One comment 2: 2 to 10 comments 3: 11-50 comments 4: More than 50 comments
has_bad_comm ent	int	Whether there is negative feedback.	0: No; 1: Yes.
bad_comment_ rate	float	Dissatisfaction rate	Proportion of the negative feedback.

• Action data

Table 3-8 Action data

Data Type	Description	Value
int	User ID	Anonymized
int	Product No.	Anonymized
string	Time of action	-
string	Module ID	Anonymized
string	 Browse (refers to the offering details page) Add to cart Remove from cart Place an order Follow 	-
	int int string string	int User ID int Product No. string Time of action string Module ID string • Browse (refers to the offering details page) • Add to cart • Remove from cart • Place an order

Step 1: Uploading Data

Upload the data to OBS for data analysis using DLI.

- Download OBS Browser+. For details about the download address, see Object Storage Service Tool Guide.
- Install OBS Browser+. For details about the installation procedure, see Object Storage Service Tool Guide.

- Log in to OBS Browser+. OBS Browser+ supports two login modes: AK login (using access keys) or authorization code login. For details about the login procedure, see Object Storage Service Tool Guide.
- 4. Upload data using the OBS Browser+.

On the OBS Browser+ page, click **Create Bucket**. Select a region and enter a bucket name (for example, **DLI-demo**). After the bucket is created, return to the bucket list and click **DLI-demo**. OBS Browser+ supports upload by dragging. You can drag one or more files or folders from a local path to the object list of a bucket or a parallel file system on OBS Browser+. You can even drag a file or folder directly to a specified folder on OBS Browser+.

Obtain the test data by downloading the **Best_Practice_04.zip** file, decompressing it, and uploading the **Data** folder to the root directory of the OBS bucket. The test data directory is as follows:

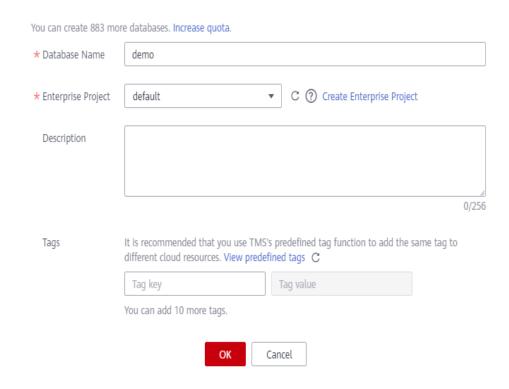
- data/JData_User: Data in the user table
- data/JData_Product: Data in the product table
- data/JData_Product/JData_Comment: Data in the comment table
- data/JData_Action: Data the action table

Step 2: Analyzing Data

- 1. Creating a Database and a Table
 - On the top menu bar of the portal page, choose Products > Analytics > Data Lake Insight (DLI).
 - b. Create a demo database. On the DLI console, choose Job Management
 SQL Jobs. Click the created job on the displayed page to go to the SQL Editor page.
 - c. In the left pane of the SQL Editor, select the **Databases** tab and click to create the **demo** database. For details, see **Figure 3-15**.

Figure 3-15 Creating a database

Create Database



□ NOTE

The **default** database is a built-in database. You cannot create a database named **default**.

d. Choose the **demo** database, and enter the following SQL statement in the editing box:

```
create table user(
    user_id int,
    age int,
    gender int,
    rank int,
    register_time string
) USING csv OPTIONS (path "obs://DLI-demo/data/JData_User")
```

□ NOTE

The file path in the preceding SQL statement is the actual OBS path for storing data.

- e. Click **Execute** to create the user information table user.
- f. Create the **product**, **comment**, and **action** tables in the same way.
 - Product data

```
create table product(
product_id int,
a1 int,
a2 int,
a3 int,
category int,
brand int
) USING csv OPTIONS (path "obs://DLI-demo/data/JData_Product")
```

Comment table

```
create table comment(
    deadline string,
    product_id int,
    comment_num int,
    has_bad_comment int,
    bad_comment_rate float
) USING csv OPTIONS (path "obs://DLI-demo/data/JData_Comment")
```

Action table

```
create table action(
user_id int,
product_id int,
time string,
model_id string,
type string
) USING csv OPTIONS (path "obs://DLI-demo/data/JData_Action");
```

2. Querying Data

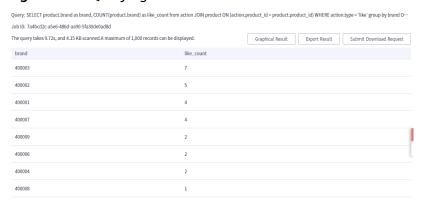
You can save common query statements as templates on the **Template Management** page for later use. For details, see **SQL Template Management** in *Data Lake Insight User Guide*.

- Top 10 products with the most likes
 - i. Run the following SQL statement to analyze the top 10 products with the most likes.

```
SELECT
product.brand as brand,
COUNT(product.brand) as like_count
from
action
JOIN product ON (action.product_id = product.product_id)
WHERE
action.type = 'like'
group by
brand
ORDER BY like_count desc
limit
10
```

ii. Click **Execute**. The execution results are displayed, as shown in **Figure 3-16**.

Figure 3-16 Querying results



iii. Click to view the result in a chart.

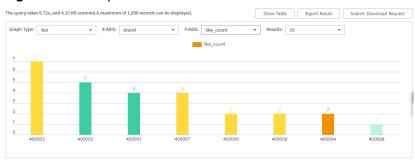


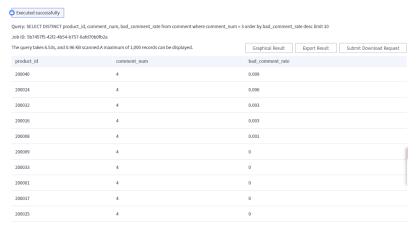
Figure 3-17 Graphical results

- Top 10 worst-rated products
 - i. Run the following SQL statement to analyze the top 10 worst-rated products:

```
SELECT
DISTINCT product_id,
comment_num,
bad_comment_rate
from
comment
where
comment_num > 3
order by
bad_comment_rate desc
limit
10
```

ii. Click **Execute**. The execution results are displayed, as shown in **Figure 3-18**.

Figure 3-18 Querying results



iii. Click to view the result in a chart.

Figure 3-19 Graphical result

You can also analyze data for age distribution, gender ratio, offering evaluation, purchase number, and browsing statistics of users.

3.4 Analyzing DLI Billing Data

Application Scenarios

You can analyze DLI billing data (account information has been masked) on the big data analysis platform of DLI, find possible optimization, and figure out some measures to reduce costs for using DLI.

Analysis Process

Perform the following steps to analyze billing data and reduce costs:

Step 1: Obtaining Consumption Data. Obtain billing data of an account.

Step 2: Analyzing Billing Data and Reducing Costs. Analyze the consumption data, find the resources or users with high expenditure, and provide optimization measures to reduce cost.

Resources and Costs

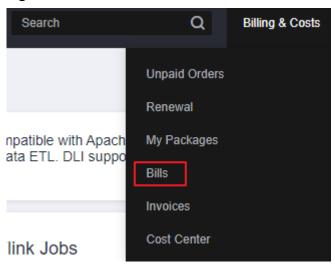
Table 3-9 Resource planning and costs

Resource	Description	Cost
Resource DLI	Description DLI is a big data analytics platform on HUAWEI CLOUD. You are billed for using storage and compute resources. DLI supports three billing modes: yearly/monthly, package, and pay-per-use.	You can run SQL jobs, Flink jobs, and Spark jobs on DLI. For SQL jobs, you are billed for both storage and compute resources. Compute resources can be billed based on a yearly/monthly basis or pay-peruse. If you choose the yearly/monthly billing mode, fees are deducted based on the subscription period. This billing mode is recommended for its preferential price and exclusive compute resources within the
		 In pay-per-use mode, fees are deducted by hour. You can choose either billing by CUH or by the amount of data scanned. Billing by CUH is recommended, for you can have exclusive resources and clear costing. In addition, you can purchase and use packages. Billing for CUH used = Number of CUs x Usage duration x Unit price. The unit of the usage duration is hour. If the duration is less than one hour, it is rounded to one hour.
		 Billing for the amount of data scanned = Amount of data scanned during SQL statement execution x Unit price. If a computing task times out or fails, no fee is charged for the task. For Flink and Spark jobs, you will be billed for compute resources only. The billing rules are same to those of SQL jobs. For details, see Price Calculator.

Step 1: Obtaining Consumption Data

- 1. Obtain billing details.
 - a. Log in to the DLI console.
 - b. Click **Billing & Costs** on the upper right corner of the page. Choose **Bills**.

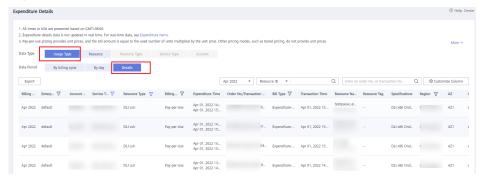
Figure 3-20 Bills



On the Dashboard page of the Billing Center, click Expenditure Details.
 On the displayed page, set Data Type to Usage Type and Data Period to Details. Set time to the billing cycle you want.

In the title row of the displayed table, set **Service Type** to **Data Lake Insight (DLI)** and **Resource Type** to **DLI cuh**. Click **Export**. On the **Export**page, configure **Export Content** and **Period** as you need, and click **Export**. The **Export History** page is displayed.

Figure 3-21 DLI Bills



d. On the **Export History** page, wait until the file status changes to **Successful**. Click **Download**.

Step 2: Analyzing Billing Data and Reducing Costs

- 1. Analyze billing details.
 - Upload the billing details downloaded in Step 1: Obtaining Consumption Data to the created OBS bucket.

- b. Create a table on DLI.
 - i. Log in to the DLI console. In the navigation pane, choose **SQL Editor**. Select **spark** for **Engine**, and select the queue and database. In this example, the default queue and database are used.
 - ii. The downloaded file contains information such as time and usage. Create a table on DLI based on these table headers. For details, see the following example.

```
CREATE TABLE `spending` (
 account_period string,
 EnterpriseProject string,
 EnterpriseProjectID string,
 accountID string,
 product_type_code string,
 product_type string,
 product_code string,
 product_name string,
 product_id string,
 mode string,
 time1 string,
 use_start string,
 use_end string,
 orderid string,
 ordertime string,
 resource_type string,
 resource_id string,
 resouce_name string,
 tag string,
 skuid string,
 `c22name` STRING,
`c23name` STRING,
 `c24name` STRING,
 `c25name` STRING,
 `c26name` STRING,
 `c27name` STRING.
 `c28name` STRING,
 `c29name` STRING,
 size STRING,
 `c31name` STRING,
 `c32name` STRING,
 `c33name` STRING,
 `c34name` STRING.
 `c35name` STRING,
 `amount` STRING,
 `c37name` STRING,
 `c38name` STRING,
 `c39name` STRING,
 `c40name` STRING,
 `c41name` STRING,
 `c42name` STRING,
 `c43name` STRING,
 `c44name` STRING,
 `c45name` STRING,
 `c46name` STRING,
 `c47name` STRING,
 `c48name` STRING,
 `c49name` STRING,
 `c50name` STRING,
`c51name` STRING,
 `c52name` STRING,
 `c53name` STRING,
 `c54name` STRING
) USING csv options (
 path 'obs://xxx/Spendings(ByTransaction)_20200501_20200531.csv',
 header true)
```

c. Query **resource_id** and **resource_name** with the highest amount within the period.

The following statement shows the amount charged for using the SQL and Flink queues.

select resource_id, resouce_name, sum(size)
as usage, sum(amount)
as sum_amount
from spending
group by resource_id, resouce_name
order by sum_amount desc

Figure 3-22 Query results

resource_id	resouce_name	usage	sum_amount
d91d4616-b10c-471a-820d-e676e6c514b4	sql	5264	1842.399999999896
8163cc27-89ce-4bac-aa85-38cb753ee425	flink	5264	1842.39999999896
9bd0736b-f8ca-4bfb-b3e7-0e391ef7dd8b	null	48	14.3999999999999
dd3a12ff-c0af-4ad1-bbc1-858bf4d3661c	ditest	32	11.2
f8265ef5-eb5f-4eff-b8d6-9ca91ed20009	test	16	5.6

d. Run the following statements to analyze the usage periods of SQL and Flink resources:

select * from spending where resource_id = 'd91d4616-b10c-471a-820d-e676e6c5f4b4' order by ordertime

The SQL queue was billed each hour from May 14 2020 17:00:00 GMT +08:00 to May 28, 2020 10:00:00 GMT+08:00.

Similarly, the Flink queue was continuously used from May 14, 2020 17:00:00 GMT+08:00 to May 28 2020 10:00:00 GMT+08:00.

2. Suggestion for reducing the cost

You can change the SQL and Flink queues to yearly/monthly queues for lower costs. If you are sure about the number of CUHs required for a job, you can purchase a package to reduce the cost.

DLI helps you to analyze billing details of your enterprise to quickly find the unreasonable expenses and control costs. You can also use DLI to reduce your cost on HUAWEI CLOUD.

3.5 Using DLI Flink SQL to Analyze e-Commerce Business Data in Real Time

Application Scenarios

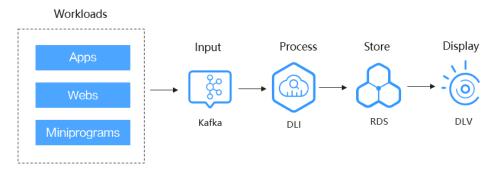
Online shopping is very popular for its convenience and flexibility. e-Commerce platform can be accessed via an array of methods, such as visiting the websites, using shopping apps, and accessing through mini-programs. A large volume of statistics data such as the real-time access volume, number of orders, and number of visitors needs to be collected and analyzed on each e-commerce platform every day. These data needs to be displayed in an intuitive way and updated in real time to help managers learn about data changes in a timely manner and adjust marketing policies accordingly. How can we efficiently and quickly collect statistics based on these metrics?

Assume the order information of each offering is written into Kafka in real time. The information includes the order ID, channel (websites or apps), order creation time, amount, actual payment amount after discount, payment time, user ID, username, and region ID. We need to collect statistics on such information based on metrics of each sales channel in real time, store the statistics in a database, and display the statistics on screens.

Solution Overview

The following figure gives you an overview to user DLI Flink to analyze and process real-time e-commerce business data and sales data of all channels.

Figure 3-23 Solution overview



Process

To analyze real-time e-commerce data with DLI Flink, perform the following steps:

Step 1: Creating Resources. Create resources required for creating jobs belong to your account, including VPC, DMS, DLI, and RDS.

Step 2: Obtaining the DMS Connection Address and Creating a Topic. Obtain the connection address of the DMS Kafka instance and create a DMS topic.

Step 3: Creating an RDS Database Table. Obtain the private IP address of the RDS DB instance and log in to the instance to create an RDS database and MySQL table.

Step 4: Creating an Enhanced Datasource Connection. Create an enhanced datasource connection for the queue and test the connectivity between the queue and the RDS instance and the queue and the DMS instance, respectively.

Step 5: Creating and Submitting a Flink Job. Create a DLI Flink OpenSource SQL iob and run it.

Step 6: Querying the Result. Query the Flink job results and display the results on a screen in DLV.

Solution Advantages

- Cross-source analysis: You can perform association analysis on sales summary data of each channel stored in OBS. There is no need for data migration.
- SQL only: DLI has interconnected with multiple data sources. You can create tables using SQL statements to complete data source mapping.

Resource Planning and Costs

Table 3-10 Resource planning and costs

Resource	Description	Cost
OBS	You need to create an OBS bucket and upload data to OBS for data analysis using DLI.	 You will be charged for using the following OBS resources: Storage Fee for storing static website files in OBS. Request Fee for accessing static website files stored in OBS. Traffic Fee for using a custom domain name to access OBS over the public network. The actual fee depends on the size of the stored file, the number of user access requests, and the traffic volume. Estimate the fee based on your service requirements.
DLI	Before creating a SQL job, you need to purchase a queue. When using queue resources, you are billed based on the CUH of the queue.	For example, if you purchase a pay-per-use queue, you will be billed based on the number of CUHs used by the queue. Usage is billed by the hour. For example, 58 minutes of usage will be rounded to the hour. CUH pay-per-use billing = Unit price x Number of CUs x Number of hours.
VPC	You can customize subnets, security groups, network ACLs, and assign EIPs and bandwidths.	The VPC service is free of charge. EIPs are required if your resources need to access the Internet. EIP supports two billing modes: pay-per-use and yearly/monthly. For more, see VPC Billing Description.
DMS Kafka	Kafka provides premium instances with computing, storage, and exclusive bandwidth resources.	Kafka supports two billing modes: pay-per-use and yearly/monthly. Billing items include Kafka instances and Kafka disk storage space. For details, see DMS for Kafka Billing Description.
RDS MySQL	RDS for MySQL provides online cloud database services.	You are billed for RDS DB instances, database storage, and backup storage (optional). For details, see RDS Billing Description.

Resource	Description	Cost
DLV	DLV adapts to a wide range of on-premise and cloud data sources, and provides diverse visualized components for you to quickly customize your data screens.	If you use the DLV service, you will be charged for the purchased yearly/monthly DLV package.

Example Data

• Order details wide table

Field	Data Type	Description
order_id	string	Order ID.
order_channel	string	Order channel (websites or apps)
order_time	string	Time
pay_amount	double	Order amount
real_pay	double	Actual amount paid
pay_time	string	Payment time
user_id	string	User ID
user_name	string	Username
area_id	string	Region ID

• Result table: real-time statistics of the total sales amount in each channel

Field	Data Type	Description
begin_time	varchar(32)	Start time for collecting statistics on metrics
channel_code	varchar(32)	Channel code
channel_name	varchar(32)	Channel
cur_gmv	double	Gross merchandises value (GMV) of the day

Field	Data Type	Description
cur_order_user_count	bigint	Number of users who settled the payment in the day
cur_order_count	bigint	Number of orders paid on the day
last_pay_time	varchar(32)	Latest settlement time
flink_current_time	varchar(32)	Flink data processing time

Step 1: Creating Resources

Create VPC, DMS, RDS, DLI, and DLV resources listed in Table 3-11.

Table 3-11 Cloud resources required

Resource	Description	Instructions
VPC	A VPC manages network resources on the cloud.	Creating a VPC and Subnet
	The network planning is described as follows:	
	The VPCs specified for the Kafka and MySQL instances must be the same.	
	The VPC network segment where the Kafka and MySQL instances belong cannot conflict with that of the DLI queue.	
DMS Kafka	In this example, the DMS for Kafka instance is the data source.	Getting Started with DMS for Kafka
RDS MySQL	In this example, an RDS for MySQL instance provides the cloud database service.	Getting Started with RDS for MySQL
DLI	DLI provides real-time data analysis. Create a general-purpose queue that uses dedicated resources in yearly/monthly or pay-per-use billing mode. Otherwise, an enhanced network connection cannot be created.	Creating a Queue
DLV	DLV displays the result data processed by the DLI queue in real time.	Creating Screens

Step 2: Obtaining the DMS Connection Address and Creating a Topic

 Hover the mouth on the Service List icon and choose Distributed Message Service in Application. The DMS console is displayed. On the DMS for Kafka page, locate the Kafka instance you have created.

Figure 3-24 Kafka instances



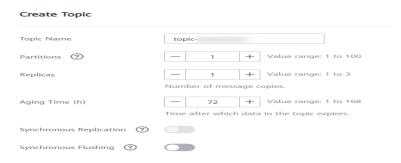
2. The instance details page is displayed. Obtain the **Instance Address (Private Network)** in the **Connection** pane.

Figure 3-25 Connection address



3. Create a topic and name it trade_order_detail_info.

Figure 3-26 Creating a topic



Configure the required topic parameters as follows:

- Partitions: Set it to 1.
- Replicas: Set it to 1.
- Aging Time: Set it to 72 hours.
- Synchronous Flushing: Disable this function.

Step 3: Creating an RDS Database Table

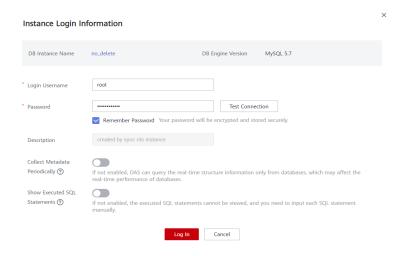
 Log in to the console, hover your mouse over the service list icon and choose Relational Database Service in Databases. The RDS console is displayed. On the Instances page, locate the created DB instance and view its floating IP address.

Figure 3-27 Viewing the floating IP address



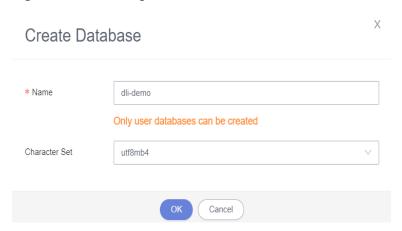
 Click More > Log In in the Operation column. On the displayed page, enter the username and password for logging in to the instance and click Test Connection. After Connection is successful is displayed, click Log In.

Figure 3-28 Logging in to an Instance



Click Create Database. In the displayed dialog box, enter database name dlidemo. Then, click OK.

Figure 3-29 Creating a database

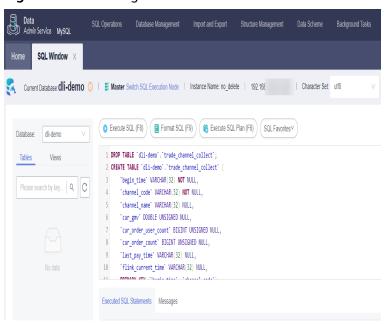


 Choose SQL Operation > SQL Query and run the following SQL statement to create a MySQL table for test (Example Data describes the fields):

```
DROP TABLE 'dli-demo'. 'trade_channel_collect';
CREATE TABLE 'dli-demo'. 'trade_channel_collect' (
    'begin_time' VARCHAR(32) NOT NULL,
    'channel_code' VARCHAR(32) NOT NULL,
    'channel_name' VARCHAR(32) NULL,
    'cur_gmv' DOUBLE UNSIGNED NULL,
    'cur_order_user_count' BIGINT UNSIGNED NULL,
    'cur_order_count' BIGINT UNSIGNED NULL,
    'last_pay_time' VARCHAR(32) NULL,
    'flink_current_time' VARCHAR(32) NULL,
```

```
PRIMARY KEY ('begin_time', 'channel_code')
) ENGINE = InnoDB
DEFAULT CHARACTER SET = utf8mb4
COLLATE = utf8mb4_general_ci
COMMENT = 'Real-time statistics on the total sales amount of each channel';
```

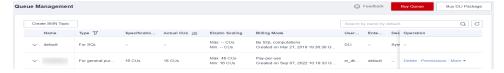
Figure 3-30 Creating a table



Step 4: Creating an Enhanced Datasource Connection

 On the management console, hover the mouse on the service list icon and choose Analytics > Data Lake Insight. The DLI management console is displayed. Choose Resources > Queue Management to query the created DLI queue.

Figure 3-31 Queue list



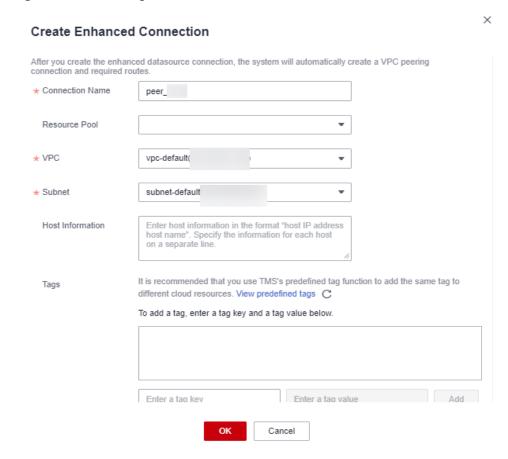
 In the navigation pane of the DLI management console, choose Global Configuration > Service Authorization. On the displayed page, select VPC Administrator, and click Update to grant the DLI user the permission to access VPC resources. The permission is used to create a VPC peering connection.

Figure 3-32 Updating agency permissions



- 3. Choose **Datasource Connections**. On the displayed **Enhanced** tab, click **Create**. Configure the following parameters, and click **OK**.
 - **Connection Name**: Enter a name.
 - **Resource Pool**: Select the general-purpose queue you have created.
 - VPC: Select the VPC where the Kafka and MySQL instances are.
 - Subnet: Select the subnet where the Kafka and MySQL instances are.

Figure 3-33 Creating an enhanced datasource



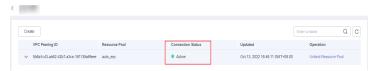
The status of the created datasource connection is **Active** in the **Enhanced** tab.

Click the name of the datasource connection. On the details page, the connection status is **ACTIVE**.

Figure 3-34 Connection status



Figure 3-35 Details



- 4. Test whether the queue can connect to RDS for MySQL and DMS for Kafka instances, respectively.
 - On the Queue Management page, locate the target queue. In the Operation column, click More > Test Address Connectivity.

Figure 3-36 Testing address connectivity



b. Enter the connection address of the DMS for Kafka instance and the private IP address of the RDS for MySQL instance to test the connectivity. If the test is successful, the DLI queue can connect to the Kafka and MySQL instances.

Figure 3-37 Testing address connectivity



If the test fails, modify the security group rules of the VPC where the Kafka and MySQL instances are to allow DLI queue access on ports 9092 and 3306. You can obtain the network segment of the queue on its details page.

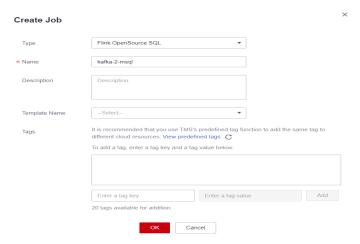
Figure 3-38 VPC security group rules



Step 5: Creating and Submitting a Flink Job

- In the navigation pane on the left, choose Job Management > Flink Jobs.
 Click Create Job.
 - Type: Select Flink OpenSource SQL.
 - Name: Enter a name.

Figure 3-39 Creating a Flink Job



2. Click **OK**. The job editing page is displayed. The following is a simple SQL statement. You need to modify some parameter values based on the RDS and DMS instance information.

```
******************
create table trade_channel_collect(
 begin_time string, --Start time of statistics collection
                       -- Channel ID
 channel_code string,
 channel_name string, -- Channel name
                       -- GMV
 cur amv double.
 cur_order_user_count bigint, -- Number of payers
 cur_order_count bigint, -- Number of orders paid on the day last_pay_time string, -- Latest settlement time
 flink_current_time string,
 primary key (begin_time, channel_code) not enforced
) with (
 "connector.type" = "idbc",
 "connector.url" = "jdbc:mysql://xxxx:3306/xxxx", -- MySQL connection address, in JDBC format
 "connector.table" = "xxxx",
                                  -- MySQL table name
 "connector.driver" = "com.mysql.jdbc.Driver",
 'pwd_auth_name'= 'xxxxx', -- Name of the datasource authentication of the password type created
on DLI. If datasource authentication is used, you do not need to set the username and password for
 "connector.write.flush.max-rows" = "1000",
 "connector.write.flush.interval" = "1s"
);
-- Temporary intermediate table
                              ********
create view tmp order detail
select *
  , case when t.order_channel not in ("webShop", "appShop", "miniAppShop") then "other"
       else t.order_channel end as channel_code -- Redefine channels. Only four enumeration values
are available: webShop, appShop, miniAppShop, and other.
  , case when t.order_channel = "webShop" then _UTF16"Website"
       when t.order_channel = "appShop" then _UTF16"Shopping App"
       when t.order_channel = "miniAppShop" then _UTF16" Miniprogram"
       else _UTF16"Other" end as channel_name -- Channel name
from (
  select *
     , row_number() over(partition by order_id order by order_time desc ) as rn -- Remove duplicate
order data
     , concat(substr("2021-03-25 12:03:00", 1, 10), " 00:00:00") as begin_time
     , concat(substr("2021-03-25 12:03:00", 1, 10), " 23:59:59") as end_time
  from trade_order_detail
  where pay_time >= concat(substr("2021-03-25 12:03:00", 1, 10), " 00:00:00") --Obtain the data of
the current day. To accelerate running, 2021-03-25 12:03:00 is used to replace
cast(LOCALTIMESTAMP as string).
  and real_pay is not null
where t.rn = 1;
-- Collect data statistics by channel.
insert into trade_channel_collect
select
   begin_time --Start time of statistics collection
  , channel_code
  , channel_name
  , cast(COALESCE(sum(real_pay), 0) as double) as cur_gmv -- GMV
  , count(distinct user_id) as cur_order_user_count -- Number of payers
  , count(1) as cur_order_count -- Number of orders paid on the day
  , max(pay_time) as last_pay_time -- Settlement time
  , cast(LOCALTIMESTAMP as string) as flink_current_time -- Current time of the flink task
from tmp_order_detail
where pay_time >= concat(substr("2021-03-25 12:03:00", 1, 10), " 00:00:00")
group by begin_time, channel_code, channel_name;
```

□ NOTE

Job logic

- 1. Create a Kafka source table to read consumption data from a specified Kafka topic.
- 2. Create a result table to write result data into MySQL through JDBC.
- 3. Implement the processing logic to collect statistics on each metric.

 To simplify the final processing logic, create a view to preprocess the data.
 - Use over window condition and filters to remove duplicate data (the top N method is used). In addition, the built-in functions concat and substr are used to set 00:00:00 as the start time and 23:59:59 of the same day as the end time, and to collect statistics on orders paid later than 00:00:00 on the day. (To facilitate data simulation, replace cast(LOCALTIMESTAMP as string) with 2021-03-25 12:03:00.)
 - Based on the channels of the order data, the built-in condition function is used to set channel_code and channel_name to obtain the field information in the source table and the values of begin_time, end_time, channel_code, and channel_name.
- 4. Collect statistics on the required metrics, filter the data as required, and write the results to the result table.
- 3. Select the created general-purpose queue and submit the job.

Queue 1.12 ★ Flink Version **UDF** Jar --Select--+? ★ CUs 2 ★ Job Manager CUs 1 +★ Parallelism Task Manager Configu... ★ OBS Bucket Save Job Log Alarm Generation upo... Enable Checkpointing Auto Restart upon Exc... Idle State Retention Time -

Figure 3-40 Flink OpenSource SQL Job

4. Wait until the job status changes to **Running**. Click the job name to view the details.

Figure 3-41 Job status



5. Use the Kafka client to send data to a specified topic to simulate real-time data streaming.

For details, see Connecting to an Instance Without SASL.

Figure 3-42 Real-time data streaming



6. Run the following command:

sh kafka_2.11-2.3.0/bin/kafka-console-producer.sh --broker-list *KafKa connection address* --topic *Topic name*

Example data is as follows:

```
{"order id":"202103241000000001", "order channel":"webShop", "order time":"2021-03-24 10:00:00",
"pay_amount":"100.00", "real_pay":"100.00", "pay_time":"2021-03-24 10:02:03", "user_id":"0001",
"user_name":"Alice", "area_id":"330106"}
{"order_id":"202103241606060001", "order_channel":"appShop", "order_time":"2021-03-24 16:06:06",
"pay_amount":"200.00", "real_pay":"180.00", "pay_time":"2021-03-24 16:10:06", "user_id":"0001",
"user_name":"Alice", "area_id":"330106"}
{"order_id":"202103251202020001", "order_channel":"miniAppShop", "order_time":"2021-03-25
12:02:02", "pay_amount":"60.00", "real_pay":"60.00", "pay_time":"2021-03-25 12:03:00",
"user_id":"0002", "user_name":"Bob", "area_id":"330110"} {"order_id":"202103251505050001", "order_channel":"qqShop", "order_time":"2021-03-25 15:05:05", "pay_amount":"500.00", "real_pay":"400.00", "pay_time":"2021-03-25 15:10:00", "user_id":"0003",
"user_name":"Cindy", "area_id":"330108"}
{"order_id":"202103252020200001", "order_channel":"webShop", "order_time":"2021-03-24 20:20:20", "pay_amount":"600.00", "real_pay":"480.00", "pay_time":"2021-03-25 00:00:00", "user_id":"0004",
"user_name":"Daisy", "area_id":"330102"}
{"order_id":"202103260808080001", "order_channel":"webShop", "order_time":"2021-03-25 08:08:08",
 'pay_amount":"300.00", "real_pay":"240.00", "pay_time":"2021-03-25 08:10:00", "user_id":"0004",
"user_name":"Daisy", "area_id":"330102"}
{"order_id":"202103261313130001", "order_channel":"webShop", "order_time":"2021-03-25 13:13:13",
"user_name":"Daisy", "area_id":"330102"}
{"order_id":"202103270606060001", "order_channel":"appShop", "order_time":"2021-03-25 06:06:06",
"pay_amount":"50.50", "real_pay":"50.50", "pay_time":"2021-03-25 06:07:00", "user_id":"0001", "user_name":"Alice", "area_id":"330106"}
{"order_id":"202103270606060002", "order_channel":"webShop", "order_time":"2021-03-25 06:06:06",
"pay_amount":"66.60", "real_pay":"66.60", "pay_time":"2021-03-25 06:07:00", "user_id":"0002",
"user_name":"Bob", "area_id":"330110"}
{"order_id":"202103270606060003", "order_channel":"miniAppShop", "order_time":"2021-03-25
06:06:06", "pay_amount":"88.80", "real_pay":"88.80", "pay_time":"2021-03-25 06:07:00",
"user_id":"0003", "user_name":"Cindy", "area_id":"330108"}
{"order_id":"202103270606060004", "order_channel":"webShop", "order_time":"2021-03-25 06:06:06",
"pay_amount":"99.90", "real_pay":"99.90", "pay_time":"2021-03-25 06:07:00", "user_id":"0004",
"user_name":"Daisy", "area_id":"330102"}
```

7. In the navigation pane on the left, choose **Job Management** > **Flink Jobs**, and click the job submitted in **3**. On the job details page, view the number of processed data records.

Figure 3-43 Job details



Step 6: Querying the Result

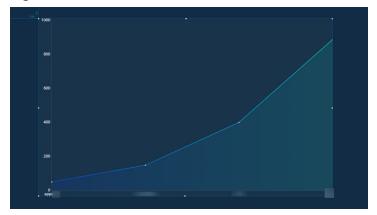
 Log in to the MySQL instance by referring to 2 and run the following SQL statement to query the result data processed by the Flink job: SELECT * FROM `dli-demo`.`trade_channel_collect`;

Figure 3-44 Querying results



 Log in to the DLV console, configure a DLV screen, and run SQL statements to query data in the RDS for MySQL instance to display the data on the screen.
 For details, see Editing Screens.

Figure 3-45 DLV screen



3.6 Interconnecting Yonghong BI with DLI to Submit Spark Jobs

3.6.1 Preparing for Yonghong BI Interconnection

Scenario

Prepare for the interconnection between Yonghong BI system and DLI.

Procedure

- Step 1 (Optional) In the upper left corner of the Huawei Cloud management console, click Service List and choose Analytics > Data Lake Insight. On the Overview page displayed, find the Common Links area on the right, and click SDK Download. On the DLI SDK DOWNLOAD page displayed, download a DLI JDBC driver, for example, dli-jdbc-1.1.0-jar-with-dependencies-jdk1.7.jar. For details, see Downloading the JDBC Driver Package.
- **Step 2** The AK/SK and token authentication modes can be used for JDBC authentication. The AK/SK authentication mode is recommended.
- **Step 3** Contact Yonghong customer service personnel to obtain the username and password of the Yonghong SaaS production environment.
- **Step 4** Log in to the Yonghong SaaS production environment and enter the username and password.

----End

3.6.2 Adding Yonghong BI Data Source

Scenario

Add the DLI data source to the Yonghong SaaS production environment.

Procedure

Step 1 On the homepage of the Yonghong SaaS production environment, click **Create Connection** from the left navigation tree. See **Figure 3-46**.

Figure 3-46 Adding a connection



Step 2 On the **New Connection Type** page, choose **GENERIC** for the type of the new connection. See **Figure 3-47**.

Does public allow control to sourch.

Does public above control to the super cord to sourch.

Does public above control to the super cord to sourch.

Does public above control to the super cord to sourch.

Does public above control to the super cord to sourch.

Does public above control to the super cord to sourch.

Does public above control to the super cord to sourch.

Does public above control to the super cord to sourch.

Does public above control to the super cord to sourch.

Does public above control to the super cord to sourch.

Does public above control to the super cord to sourch.

Does public above control to the super cord to sourch.

Does public above control to the super cord to sourch.

Does public above control to the super cord to sourch.

Does public above control to the super cord to sourch.

Does public above control to the super cord to sourch.

Does public above control to the super cord to sourch.

Does public above control to the super cord to sourch.

Does public above control to the super cord to sourch.

Does public above control to the super cord to sourch.

Does public above control to sou

Figure 3-47 Choosing a new connection type

Step 3 Configure the new connection. See Figure 3-48.

In Driver, enter com.huawei.dli.jdbc.DliDriver.

In **URL**, select **Self-defined Protocol**. Enter the URL of the DLI JDBC driver. For details about the URL format and the attributes, see **Table 3-12** and **Table 3-13**, respectively.

Ⅲ NOTE

- In Schema, you can optionally enter the name of the database to be accessed. If you enter the name, only tables in the database are displayed during data set creation. If you do not enter the name, tables in all databases are displayed during data set creation. For details about how to create a data set, see Creating Yonghong BI Data Set.
- Retain default values of other parameters. You do not need to select **Request Login**.

Figure 3-48 Configuring the new connection

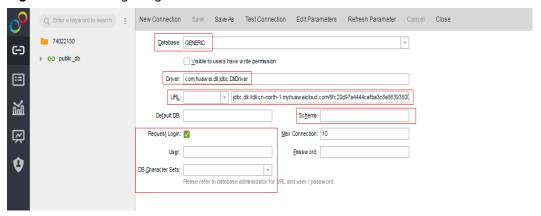


Table 3-12 Database connection parameters

Parameter	Description
URL	The URL format is as follows:
	jdbc:dli:// <endpoint>/<projectid>? <key1>=<val1>;<key2>=<val2></val2></key2></val1></key1></projectid></endpoint>
	NOTE
	 endpoint indicates the domain name of DLI. For details, see Regions and Endpoints.
	 projectId indicates the project ID, which can be obtained from the My Credentials page of the public cloud platform.
	The question mark (?) is followed by other configuration items. Each configuration item is listed in the "key=value" format. Semicolons (;) are used to separate configuration items. For details, see Table 3-13.

Table 3-13 Attribute-related configuration items

Attribute (key)	Mandatory	Defaul t Value (value)	Description
queuename	Yes	-	Queue name of DLI.
databasename	No	-	Default database to be accessed. If this parameter is not specified in the URL, you need to use db.table (for example, select * from dbother.tabletest) to access tables in the database.
authentication mode	Yes	token	Authentication method, which can be token or aksk . Value aksk is recommended during the interconnection with Yonghong BI system.
accesskey	This parameter must be configured if authentication mode is set to aksk.	-	For details, see Preparing for Yonghong BI Interconnection.

Attribute (key)	Mandatory	Defaul t Value (value)	Description
secretkey	This parameter must be configured if authentication mode is set to aksk.	-	For details, see Preparing for Yonghong BI Interconnection.
regionname	This parameter must be configured if authentication mode is set to aksk.	-	For details, see Regions and Endpoints .
servicename	This parameter must be configured if authentication mode is set to aksk.	-	servicename=dli
dli.sql.checkNo ResultQuery	No	false	Whether to allow invoking the executeQuery API to execute statements (for example, DDL) that do not return results. Value false indicates that invoking of the executeQuery
			API is allowed. • Value true indicates that invoking of the executeQuery API is not allowed. NOTE If dli.sql.checkNoResultQuery is set to false , non-query statements will be executed twice.

Step 4 On the tool bar of the displayed page, click **Test Connection**. After the test is complete, click **Save**. Enter the data source name, and save the data source.

□ NOTE

Currently, you are not allowed to save the data source to the root directory. Therefore, you can only save the data source to an existing folder.

----End

3.6.3 Creating Yonghong BI Data Set

Scenario

Create a DLI data set in the Yonghong SaaS production environment.

Procedure

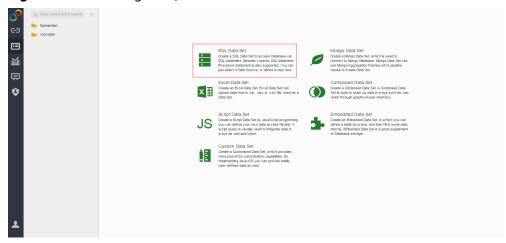
Step 1 On the home page of the Yonghong SaaS production environment, click **Create Data Set** in the left navigation tree. See **Figure 3-49**.

Figure 3-49 Creating a data set



Step 2 On the displayed page, click SQL Data Set. See Figure 3-50.

Figure 3-50 Creating a SQL data set



Step 3 On the displayed page, select the added DLI data source from the **Connection** drop-down list box. See **Figure 3-51**.

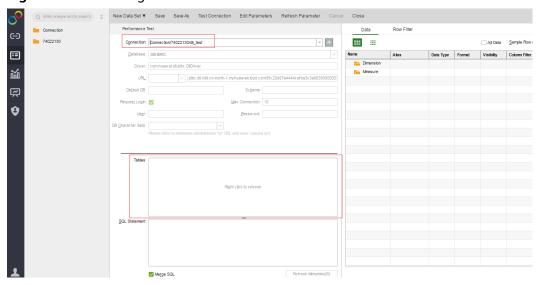
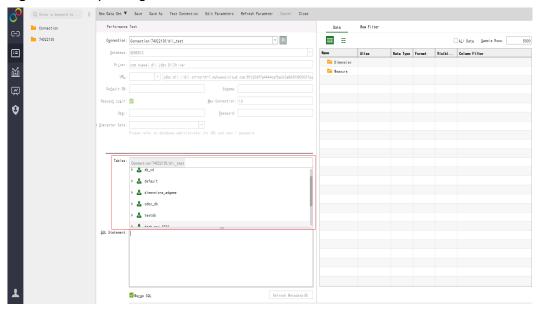


Figure 3-51 Selecting a data source

Step 4 In the **Table** area on the left pane, right-click and choose **Update** to update tables. All databases and their tables are listed in the area. **Figure 3-52** shows the page displayed when **Table Structure** is not configured during connection creation.

Figure 3-52 Updating tables



Step 5 In the SQL Statement area on the left pane, enter the select * from table_name command to query tables. On the Preview Data page on the right pane, click

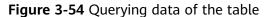
. Metadata of the table, including fields and field types, is displayed. See Figure 3-53.

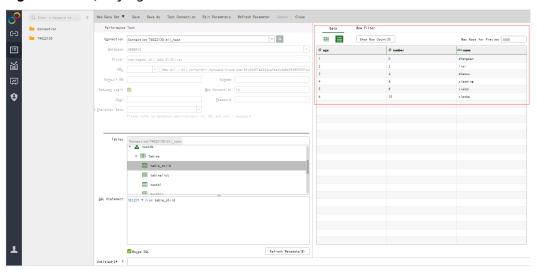
Consection
To Consection
To Page 1988

The Data Served Serve Serve

Figure 3-53 Querying the table

Step 6 Click on the right pane to query data details. See **Figure 3-54**.





Step 7 On the tool bar of the displayed page, click **Save**.

----End

3.6.4 Creating a Chart in Yonghong BI

Scenario

Create a chart in the Yonghong SaaS production environment.

Procedure

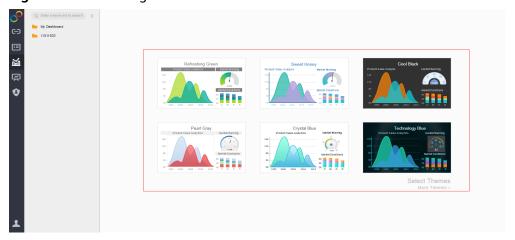
Step 1 On the home page of the Yonghong SaaS production environment, click **Create Dashboard** in the left navigation tree. See **Figure 3-55**.

Figure 3-55 Creating a dashboard



Step 2 Select a theme. See Figure 3-56.

Figure 3-56 Selecting a theme



Step 3 In this example, the Refreshing Green theme is selected. On the left pane, select the created data set from the drop-down list box and choose a table as the data source (for example, **table_child**). Metadata (including fields and field types) of the table is displayed in the lower part of the **Data** column. See **Figure 3-57**.

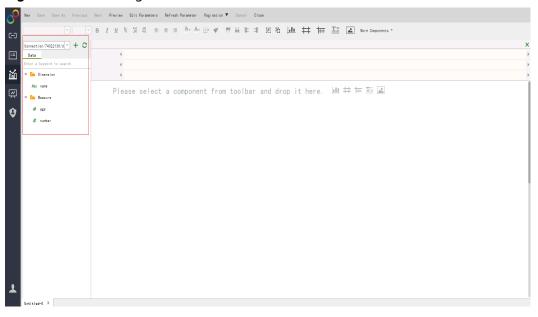
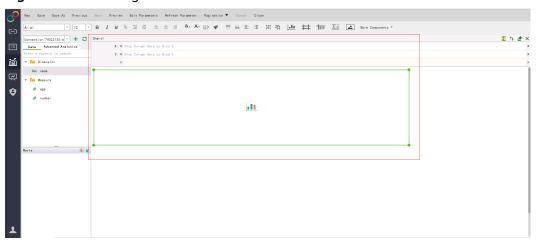


Figure 3-57 Selecting a connection for the table

Step 4 On the report creation page, chart, table, matrix, and list filtering components are available. For example, if you want to create a chart, click and drag it to the editing area. See **Figure 3-58**.

Figure 3-58 Creating a chart



Step 5 In **X**, choose **name**. In **Y**, choose **age**. Drag them to the corresponding area, and the system automatically generates a bar chart. See **Figure 3-59**.

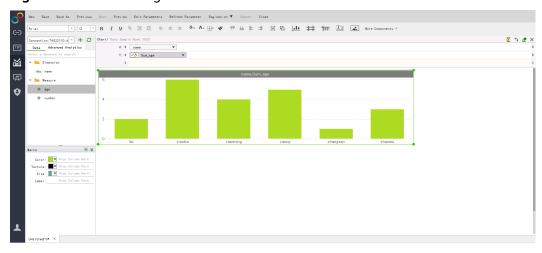


Figure 3-59 Generating a chart

Step 6 On the tool bar of the displayed page, click **Save**.

----End

3.7 Interconnecting FineBI with DLI Trino

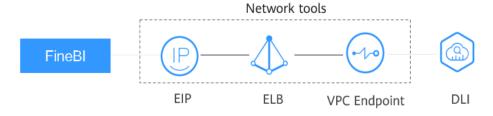
FineBI is a BI tool for big data analytics developed by FanRuan Software. It provides business personnel and data analysts with data exploration capabilities such as data management, editing, and visualization. Huawei Cloud DLI integrates data analysis and processing. SQL jobs using the Trino interactive engine are more suitable for interactive analysis and query. It provides FineBI with efficient engine compute capabilities and effective high-quality data for subsequent data statistics and analysis, helping enterprises make data decisions.

This section describes how to interconnect FineBI with DLI.

Solution Overview

This solution uses VPCEP to connect FinBI to DLI.

Figure 3-60 Architecture



Constraints

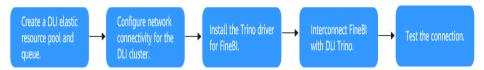
- Trino engine queues support only HTTPS connections.
- When the Trino engine is used, the created SQL queues cannot be scaled in or out.

To adjust the CU size of a queue, you need to first **delete** the queue in the elastic resource pool and then **create** a queue with Trino as the engine and with an appropriate CU size in the elastic resource pool.

- The DLI Trino engine is in the open beta test (OBT) phase. If you need it, contact customer service to apply for it.
 - The DLI Trino engine is available in the following regions: **CN North-Beijing4**, **CN East-Shanghai1**, **CN-Hong Kong**, **AP-Bangkok**, **AP-Singapore**, and **AF-Johannesburg**.
- Currently, only foreign tables created using the Hive syntax can be used for FineBI interconnection.

Process

Figure 3-61 Process of interconnecting DLI with FlineBI



To interconnect FineBI with Huawei Cloud DLI, perform the following steps:

- 1. Creating an Elastic Resource Pool and Queue
- 2. Configuring Network Connectivity for the DLI Cluster
- 3. Installing the Trino Driver for FineBI
- 4. Interconnecting FineBI with DLI Trino
- 5. Testing the Connection

Solution Advantages

- As a next-gen BI tool for self-service big data analytics, FineBI provides enterprises with one-stop solutions for enterprise business intelligence, such as multi-source data collection, self-service exploratory analysis, multi-screen support, and enterprise-level management and control.
- Huawei Cloud DLI provides convergent data analysis and processing capabilities. DLI can interconnect with multiple data sources and map data sources by creating tables using SQL statements. You can use standard SQL statements to compile metric analysis logic without paying attention to the complex distributed computing platform.
- FineBI interconnects with Huawei Cloud DLI for real-time data ingestion, efficient data processing, and good data visualization. DLI can connect to FineBI from multiple data sources. Fine BI can display DLI data in charts and reports, making data more intuitive and improving decision-making accuracy and efficiency.

Resource Planning and Costs

Table 3-14 Resource planning and costs

Resource	Description	Cost
OBS	DLI needs to use OBS buckets to store logs.	You will be charged for using the following OBS resources:
		• Storage Fee for storing static website files in OBS.
		 Request Fee for accessing static website files stored in OBS.
		 Traffic Fee for using a custom domain name to access OBS over the public network.
		The actual fee depends on the size of the stored file, the number of user access requests, and the traffic volume. Estimate the fee based on your service requirements.
DLI	In this example, an elastic resource pool is used to create SQL jobs.	When using a DLI elastic resource pool, you are billed based on the CUH of the elastic resource pool.
VPCEP	Used to enable the network connection between FineBI and DLI.	For details about the billing for VPCEP, see Billing .
ELB	Elastic Load Balance (ELB) distributs access traffic to multiple backend servers based on distribution policies.	For details about the billing for ELB, see Billing .
EIP	Provides independent public IP addresses and bandwidth for Internet access.	For details about the billing for EIP, see Billing.

Step 1: Creating an Elastic Resource Pool and Queue

- **Step 1** Log in to the DLI management console.
- **Step 2** In the navigation pane on the left, choose **Resources** > **Resource Pool**.
- **Step 3** On the **Resource Pool** page, click **Buy Resource Pool** in the upper right corner.
- **Step 4** On the displayed page, set the following parameters:

Table 3-15 Parameters

Parameter	Description
Billing Mode	Pay-per-use/Yearly/Monthly
Region	Select a region near you to ensure the lowest latency possible.
Project	Each region corresponds to a project.
Name	Name of the elastic resource pool.
CU Range	The maximum and minimum CUs allowed for the elastic resource pool.
Description	Description of the elastic resource pool.
CIDR Block	CIDR block the elastic resource pool belongs. If DLI enhanced datasource connections are required, the CIDR block of the elastic resource pool cannot overlap with that of the data source. The CIDR block of the elastic resource pool cannot be changed after being set. Recommended CIDR blocks: 10.0.0.0~10.255.0.0/16~19 172.16.0.0~172.31.0.0/16~19 192.168.0.0~192.168.0.0/16~19
Enterprise Project	If the created elastic resource pool belongs to an enterprise project, select the enterprise project.
Required Duration	You must specify the Required Duration if Billing Mode is set to Yearly/Monthly . The longer the subscription duration is, the more discounts you can get. If Auto renew is selected, monthly subscriptions are renewed each month. And yearly subscriptions are renewed each year.
Tag	Tags used to identify cloud resources.

- **Step 5** Click **Buy** and confirm the configurations.
- **Step 6** Wait until the status of the elastic resource pool changes to **Available**. The elastic resource pool is successfully created.
- **Step 7** Add a SQL queue to the elastic resource pool and select Trino as the execution engine.

□ NOTE

When buying a SQL queue, select Trino as the execution engine.

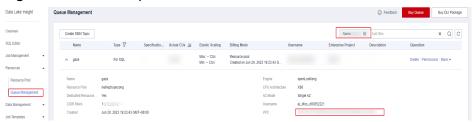
----End

Step 2: Configuring Network Connectivity for the DLI Cluster

Step 1 On the **Queue Management** page of the created DLI queue, view the VPC endpoint information of the queue.

- 1. On the DLI console, choose **Resources** > **Queue Management**, and view the VPC endpoint information about one minute after the queue is created.
- 2. Locate the created queue and click in front of the queue name to obtain the VPC endpoint information of the queue.

Figure 3-62 VPC endpoint information



Step 2 Create a VPC endpoint.

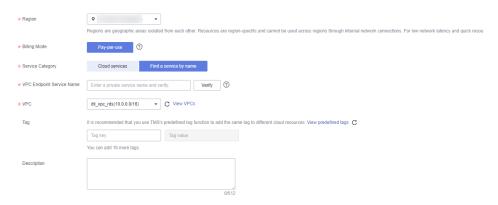
- 1. Log in to the VPC Endpoint management console.
- 2. On the VPC Endpoints page displayed, click Buy VPC Endpoint.
- 3. Set Service Category to Find a service by name.
- 4. In the **VPC Endpoint Service Name** field box, enter the **obtained** VPC endpoint information, excluding the port.

Example:

The VPC endpoint information of the queue is xxx.3a715f69-b1b0-45d0-bc4a-d917137bcd08:18090.

Enter xxx.3a715f69-b1b0-45d0-bc4a-d917137bcd08 in the field box.

Figure 3-63 Buy VPC Endpoint page



Step 3 Obtain the IP address of the VPC endpoint.

- In the navigation pane of the VPCEP console, choose VPC Endpoint > VPC Endpoints.
- Click the ID of the VPC endpoint and view the node IP address on the Summary tab page.

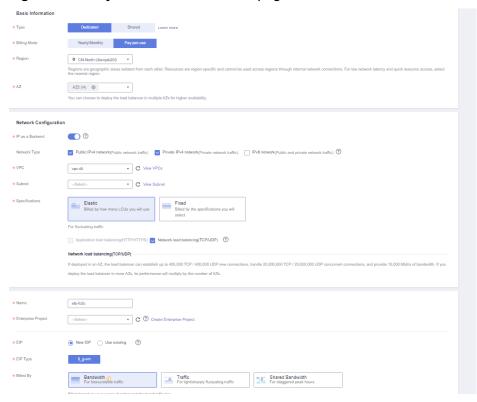
Figure 3-64 IP address of the VPC endpoint



Step 4 Create an ELB.

- Log in to the ELB console.
- 2. Click **Buy Elastic Load Balancer** and then configure the parameters.

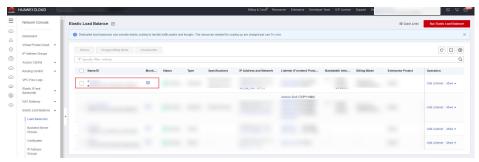
Figure 3-65 Buy Elastic Load Balancer page



Step 5 Obtain the service IP address of the ELB.

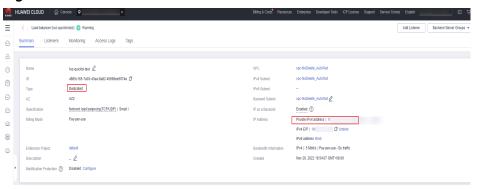
1. On the ELB console, choose **Elastic Load Balance** > **Load Balancers**.

Figure 3-66 Load balancer list



Click the ID of the created load balancer. On the Summary tab page, view the load balancer information, and record the IPv4 EIP address.

Figure 3-67 Dedicated load balancer



Step 6 Create a datasource connection.

- 1. Log in to the DLI management console.
- 2. In the left navigation pane, choose **Datasource Connections**.

On the **Enhanced** tab page displayed, click **Create**. In the **Create Enhanced Connection** dialog box, enter a connection name in **Connection Name**, set **Resource Pool** to the elastic resource pool that contains the Trino engine queue created in **step 1**, and configure **VPC** and **Subnet**. For details about the parameters, see **Table 3-16**.

Figure 3-68 Creating an enhanced datasource connection

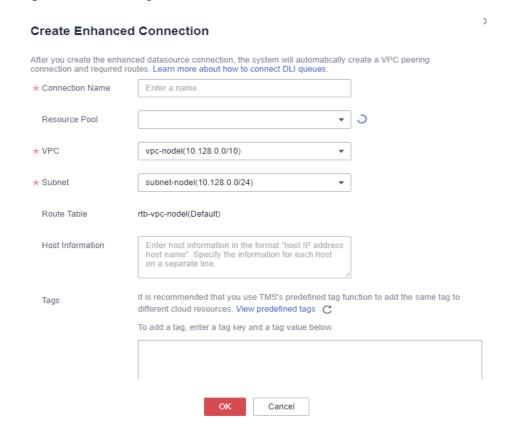


Table 3-16 Parameters

Paramete r	Description
Connectio n Name	 Name of the datasource connection to be created The name can contain only digits, letters, and underscores (_) but cannot be left blank. Enter a maximum of 64 characters.
Resource Pool	It binds an elastic resource pool or queue that uses a datasource connection. This parameter is optional. Only dedicated queues charged in yearly/monthly or pay-peruse billing mode can be bound to elastic resource pools. In regions where this function is available, an elastic resource pool with the same name is created by default for the yearly/monthly or pay-per-use dedicated queue created in "Creating a Queue." NOTE Before using an enhanced datasource connection, you must bind a queue to the connection and ensure that the VPC peering connection is in the Active state.
VPC	VPC used by the destination data source
Subnet	Subnet used by the destination data source
Route Table	Route table of the subnet NOTE - 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. - 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.
Tags	Tags used to identify cloud resources.

3. Click OK.

- Check whether the datasource connection is successfully created.
 Click the name of the created datasource connection to view its connection status. If the status is **Active**, the datasource connection is successfully created.
- **Step 7** Add a backend server group as the VPC backend. The cross-VPC backend IP address is the IP address of the purchased VPC endpoint.
 - On the ELB console, choose Elastic Load Balance > Backend Server Groups.
 On the page displayed, click Create Backend Server Group.
 - Select the created load balancer for Load Balancer and click Next. On the Backend Servers tab page, click Next. On the Confirm page, click Create Now.

< │ Create Backend Server Group Configure Routing Policy — Add Backend Server 3 Confirm * Load Balancing Type ▼ C View load balancers * Load Balancer -Select-* Backend server group name server_group-3315 * Backend Protocol Weighted round robin Weighted least connections Source IP hash Sticky Session ? Slow Start ? Description

Figure 3-69 Creating a backend server group

 On the Backend Servers tab page, click Add Backend Server in the Operation column of the created backend server group to add backend servers.

Figure 3-70 Cross-VPC backend IP address and service port



Step 8 Verify that the network connection between VPCEP and DLI is normal.

On the **IP** as **Backend Servers** tab page, check whether the health check result is normal. If yes, the network connection is normal.

Figure 3-71 Successful network connection



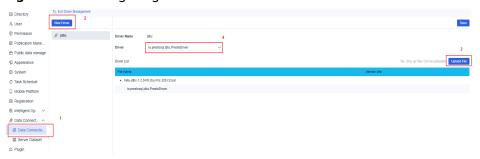
----End

Step 3: Installing the Trino Driver for FineBI

- 1. Install FineBI.
- Install the Trino driver for FineBI.
 Visit Trino to download the Trino driver JAR file.

On the FineBI console, choose **Management System > Data Connection Management**, click **New Driver**, click **Upload File**, and upload the downloaded driver package.

Figure 3-72 Configuring the Presto driver



Step 4: Interconnecting FineBI with DLI Trino

Configure the interconnection between FineBI and DLI.

- 1. On the FineBI management console, choose **Data Connection Management** > **Create Data Connection** > **Other** > **Other** D**BC**.
- 2. Enter information about the data connection.
 - Enter the data connection name.
 - b. Set **Driver** to **Custom** and select **io.prestosql.jdbc.PrestoDriver** as the driver.
 - c. Enter the cross-VPC backend IP address for **Host** and enter the service port in **Port**. For details, see **Creating a backend server group**.
 - d. Enter the username and password. The username is in the format of Account name/Username/Project ID. For details about how to obtain a project ID, see Obtaining a Project ID. If a primary account is used for connection, Account name and Username are both the account name.
 - e. Example data connection URL: jdbc:presto://{ip}/dli/default? SSL=true
 - □ NOTE

In the URL, **SSL=true** indicates that backend requests use HTTPS. Currently, Trino engine queues support only HTTPS connections.

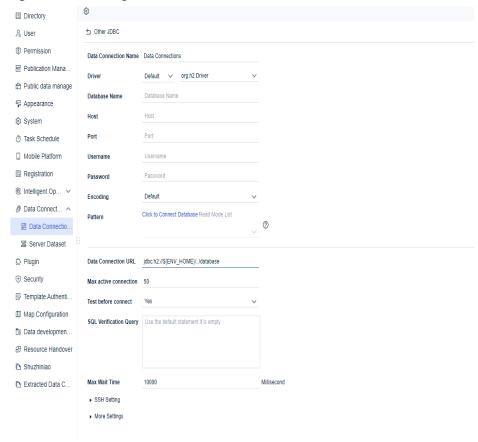


Figure 3-73 Configuration information

Step 5: Testing the Connection

Click **Test Connection** in the upper right corner of the FineBI data connection management page. If the connection is successful, you can use the connection to query DLI tables for BI report analytics.

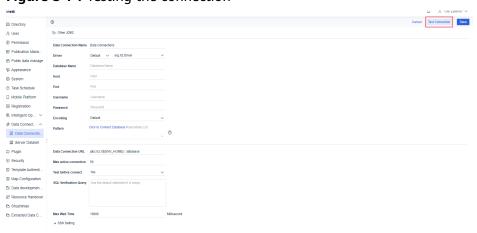


Figure 3-74 Testing the connection

Related Operations

Trino supports the SQL syntax. For details about the Trino SQL syntax, see
 Trino SQL Syntax. Currently, the Trino engine supports only the SELECT query operation.

3.8 Interconnecting Power BI with DLI Trino

Application Scenarios

Power BI is a unified, scalable self-service and enterprise business intelligence (BI) platform. You can use it to connect to and visualize any data, and seamlessly integrate visual objects into your daily applications.

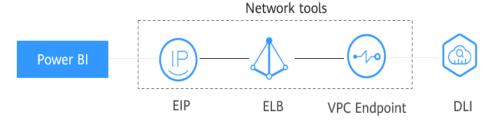
Huawei Cloud DLI provides Power BI with standard, effective, and high-quality data through converged data analysis and processing for subsequent data statistics and analysis, helping enterprises make data decisions.

For more information about Power BI, see Power BI.

Solution Overview

This solution uses VPCEP to connect Power BI to DLI.

Figure 3-75 Architecture



Constraints

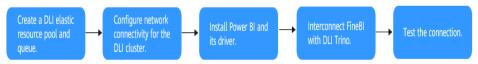
- Trino engine queues support only HTTPS connections.
- When the Trino engine is used, the created SQL queues cannot be scaled in or out.

To adjust the CU size of a queue, you need to first **delete** the queue in the elastic resource pool and then **create** a queue with Trino as the engine and with an appropriate CU size in the elastic resource pool.

- The DLI Trino engine is in the open beta test (OBT) phase. If you need it, contact customer service to apply for it.
 - The DLI Trino engine is available in the following regions: **CN North-Beijing4**, **CN East-Shanghai1**, **CN-Hong Kong**, **AP-Bangkok**, **AP-Singapore**, and **AF-Johannesburg**.
- Currently, only foreign tables created using the Hive syntax can be used for Power BI interconnection.

Process

Figure 3-76 Process of interconnecting Power BI with DLI Trino



To interconnect DLI Trino with Power BI, perform the following steps:

- 1. Creating a DLI Elastic Resource Pool and Queue
- 2. Configuring Network Connectivity for the DLI Cluster
- 3. Installing Power BI and Its Driver
- 4. Interconnecting Power BI with DLI Trino
- 5. **Testing the Connection**

Solution Advantages

- This BI tool for big data analytics provides data exploration capabilities by converting data into various forms, such as charts, tables, maps, and dashboards, making data more intuitive and easy to understand.
- Huawei Cloud DLI provides convergent data analysis and processing capabilities. DLI can interconnect with multiple data sources and map data sources by creating tables using SQL statements. DLI provides powerful data exploration capabilities to deeply explore data potentials through data filtering, sorting, and grouping.
- Power BI interconnects with Huawei Cloud DLI for real-time data ingestion, high data accuracy, efficient data processing, and good data visualization. After DLI interconnects with Power BI, data from different data sources can be integrated. DLI supports big data processing. It can process a large amount of data and mine the potential of the data. Power BI can quickly visualize the data to improve the efficiency and precision of data analytics.

Resource Planning and Costs

Table 3-17 Resource planning and costs

Resource	Description	Cost
OBS	DLI needs to use OBS buckets to store logs.	You will be charged for using the following OBS resources:
		• Storage Fee for storing static website files in OBS.
		 Request Fee for accessing static website files stored in OBS.
		 Traffic Fee for using a custom domain name to access OBS over the public network.
		The actual fee depends on the size of the stored file, the number of user access requests, and the traffic volume. Estimate the fee based on your service requirements.
DLI	In this example, an elastic resource pool is used to create SQL jobs.	When using a DLI elastic resource pool, you are billed based on the CUH of the elastic resource pool.
VPCEP	Used to enable the network connection between FineBI and DLI.	For details about the billing for VPCEP, see Billing .
ELB	Elastic Load Balance (ELB) distributs access traffic to multiple backend servers based on distribution policies.	For details about the billing for ELB, see Billing .
EIP	Provides independent public IP addresses and bandwidth for Internet access.	For details about the billing for EIP, see Billing .

Step 1: Creating an Elastic Resource Pool and Queue

- **Step 1** Log in to the DLI management console.
- **Step 2** In the navigation pane on the left, choose **Resources** > **Resource Pool**.
- **Step 3** On the **Resource Pool** page, click **Buy Resource Pool** in the upper right corner.
- **Step 4** On the displayed page, set the following parameters:

Table 3-18 Parameters

Parameter	Description
Billing Mode	Pay-per-use/Yearly/Monthly
Region	Select a region. Select a region near you to ensure the lowest latency possible.
Project	Each region corresponds to a project.
Name	Name of the elastic resource pool.
CU Range	The maximum and minimum CUs allowed for the elastic resource pool.
Description	Description of the elastic resource pool.
CIDR Block	CIDR block the elastic resource pool belongs. If DLI enhanced datasource connections are required, the CIDR block of the elastic resource pool cannot overlap with that of the data source. The CIDR block of the elastic resource pool cannot be changed after being set. Recommended CIDR blocks: 10.0.0.0~10.255.0.0/16~19 172.16.0.0~172.31.0.0/16~19 192.168.0.0~192.168.0.0/16~19
Enterprise Project	If the created elastic resource pool belongs to an enterprise project, select the enterprise project.
Required Duration	You must specify the Required Duration if Billing Mode is set to Yearly/Monthly . The longer the subscription duration is, the more discounts you can get. If Auto renew is selected, monthly subscriptions are renewed each month. And yearly subscriptions are renewed each year.
Tags	Tags used to identify cloud resources.

- **Step 5** Click **Buy** and confirm the configurations.
- **Step 6** Wait until the status of the elastic resource pool changes to **Available**. The elastic resource pool is successfully created.
- **Step 7** Add a SQL queue to the elastic resource pool. The selected engine is Trino.

□ NOTE

When buying a SQL queue, select Trino as the execution engine.

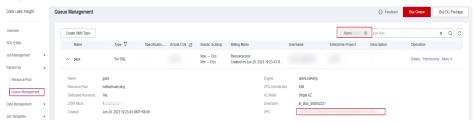
----End

Step 2: Configuring Network Connectivity for the DLI Cluster

Step 1 On the **Queue Management** page of the created DLI queue, view the VPC endpoint information of the queue.

- 1. On the DLI console, choose **Resources** > **Queue Management**.
- 2. Locate the created queue and click in front of the queue name to obtain the VPC endpoint information of the queue.

Figure 3-77 VPC endpoint information



Step 2 Create a VPC endpoint.

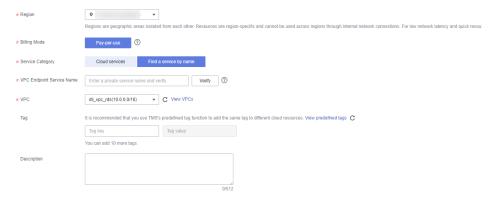
- 1. Log in to the VPC Endpoint management console.
- 2. Click **Buy VPC Endpoint**. The **Buy VPC Endpoint** page is displayed.
- 3. Set Service Category to Find a service by name.
- In the VPC Endpoint Service Name field box, enter the obtained VPC endpoint information, excluding the port.

Example:

The VPC endpoint information of the queue is xxx.3a715f69-b1b0-45d0-bc4a-d917137bcd08:18090.

Enter xxx.3a715f69-b1b0-45d0-bc4a-d917137bcd08 in the field box.

Figure 3-78 Buy VPC Endpoint page



Step 3 Obtain the IP address of the VPC endpoint.

- In the navigation pane of the VPCEP console, choose VPC Endpoint > VPC Endpoints.
- Click the ID of the VPC endpoint and view the node IP address on the Summary tab page.

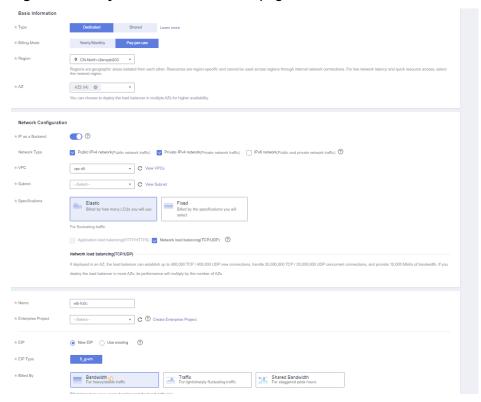
Figure 3-79 IP address of the VPC endpoint



Step 4 Create an ELB.

- 1. Log in to the ELB console.
- 2. Click **Buy Elastic Load Balancer** and then configure the parameters.

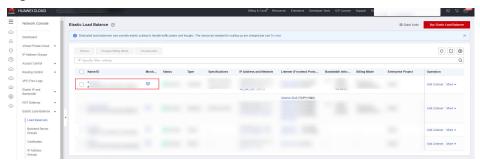
Figure 3-80 Buy Elastic Load Balancer page



Step 5 Obtain the service IP address of the ELB.

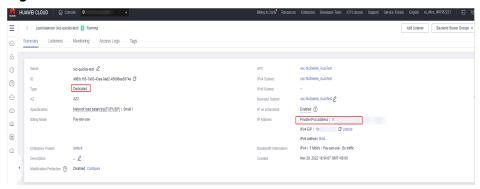
1. On the ELB console, choose **Elastic Load Balance** > **Load Balancers**.

Figure 3-81 Load balancer list



Click the ID of the created load balancer. On the Summary tab page, view the load balancer information, and record the IPv4 EIP address.

Figure 3-82 Dedicated load balancer



Step 6 Create a datasource connection.

- 1. Log in to the DLI management console.
- 2. In the left navigation pane, choose **Datasource Connections**.

On the **Enhanced** tab page displayed, click **Create**. In the **Create Enhanced Connection** dialog box, enter a connection name in **Connection Name**, set **Resource Pool** to the elastic resource pool that contains the Trino engine queue created in **step 1**, and configure **VPC** and **Subnet**. For details about the parameters, see **Table 3-19**.

Figure 3-83 Creating an enhanced datasource connection

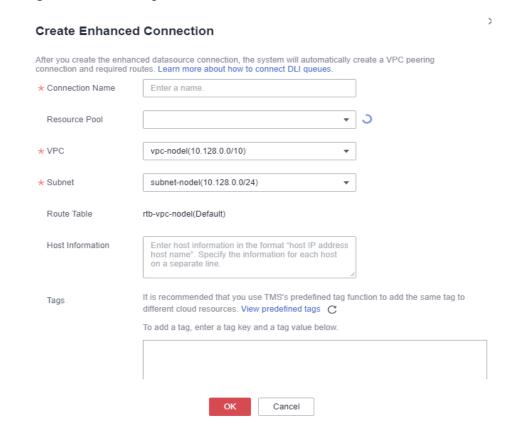


Table 3-19 Parameters

Paramete r	Description
Connectio n Name	 Name of the datasource connection to be created The name can contain only digits, letters, and underscores (_) but cannot be left blank. Enter a maximum of 64 characters.
Resource Pool	It binds an elastic resource pool or queue that uses a datasource connection. This parameter is optional. Only dedicated queues charged in yearly/monthly or pay-peruse billing mode can be bound to elastic resource pools. In regions where this function is available, an elastic resource pool with the same name is created by default for the yearly/monthly or pay-per-use dedicated queue created in "Creating a Queue." NOTE Before using an enhanced datasource connection, you must bind a queue to the connection and ensure that the VPC peering connection is in the Active state.
VPC	VPC used by the destination data source
Subnet	Subnet used by the destination data source
Route Table	Route table of the subnet NOTE - 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. - 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.
Tags	Tags used to identify cloud resources.

3. Click OK.

- Check whether the datasource connection is successfully created.
 Click the name of the created datasource connection to view its connection status. If the status is **Active**, the datasource connection is successfully created.
- **Step 7** Add a backend server group as the VPC backend. The cross-VPC backend IP address is the IP address of the purchased VPC endpoint.
 - On the ELB console, choose Elastic Load Balance > Backend Server Groups.
 On the page displayed, click Create Backend Server Group.
 - Select the created load balancer for Load Balancer and click Next. On the Backend Servers tab page, click Next. On the Confirm page, click Create Now.

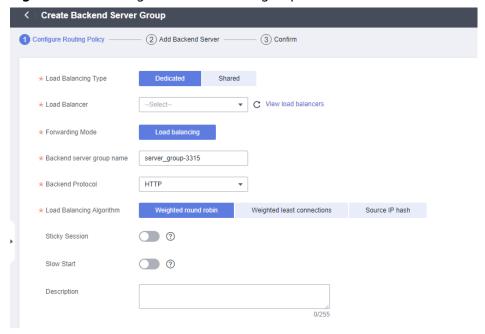


Figure 3-84 Creating a backend server group

 On the Backend Servers tab page, click Add Backend Server in the Operation column of the created backend server group to add backend servers.

Step 8 Verify that the network connection between VPCEP and DLI is normal.

On the **IP as Backend Servers** tab page, check whether the health check result is normal. If yes, the network connection is normal.

----End

Step 3: Installing Power BI and Its Driver

Step 1 Install Power BI of the desktop version.

Download the Power BI desktop version.

Step 2 Install the openLooKeng ODBC driver.

Before installing the **driver**, ensure that you have the administrator rights.

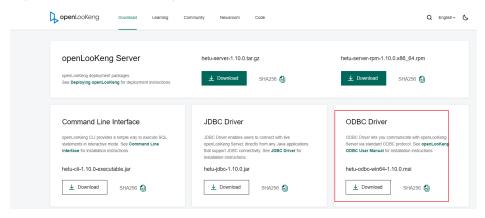
- Double-click the .msi installation package. The welcome page is displayed. Click Next.
- 2. The second page is the user agreement. Accept the terms and click **Next**.
- 3. On the third page, select an installation mode. You are advised to select **Complete**.
- 4. On the fourth page, select an installation path and click **Next**.
- 5. After the preceding installation settings are complete, click **Install** on the last page to start the installation.

□ NOTE

During the installation, the CLI is displayed to show the process of installing the driver components. After the installation is complete, the CLI is automatically closed. The openLooKeng ODBC driver is installed.

In the dialog box displayed, use DSN for new installation if you have configured the user DSN for an earlier version and click **Finish**.

Figure 3-85 openLooKeng ODBC driver



----End

Step 4: Interconnecting Power BI with DLI Trino

- 1. Stop the ODBC service.
 - a. Run the following command to go to the C:\Program Files \openLooKeng\openLooKeng ODBC Driver 64-bit\odbc_gateway \mycat\bin directory:
 - cd C:\Program Files\openLooKeng\openLooKeng ODBC Driver 64-bit\odbc_gateway\mycat\bin
 - Run the following command to stop the ODBC service: mycat.bat stop
- 2. Replace the JDBC driver.
 - Copy the JDBC JAR file obtained from driver to the C:\ProgramFiles \openLooKeng\openLooKeng ODBC Driver 64-bit\odbc_gateway \mvcat\lib directory.
 - b. Delete the existing **hetu-jdbc-1.0.1.jar** file from the directory.
- 3. Edit the protocol prefix of the ODBC server.xml file.
 - Change the property value of the **server.xml** file in the **C:\Program Files** **openLooKeng\openLooKeng ODBC Driver 64- bit\odbc_gateway\mycat \conf** directory from **property name="jdbcUrlPrefix">jdbc:lk://
 to property name="jdbcUrlPrefix">jdbc:presto://
 property>.**
- 4. Configure the connection mode of using the username and password.
 - Create the **jdbc_param.properties** file in a user-defined path, for example, **D:**, and add the following content to the file:

SSL=true user={Account name}/{Username}/{Project ID} password={Password}

- Restart the ODBC service.
 - a. Run the following command to go to the C:\Program Files
 \openLooKeng\openLooKeng ODBC Driver 64-bit\odbc_gateway
 \mvcat\bin directory:

cd C:\Program Files\openLooKeng\openLooKeng ODBC Driver 64-bit\odbc_gateway\mycat\bin

- b. Run the following command to restart the ODBC service: mycat.bat restart
- 6. Set up ODBC data sources (64-bit).

Enter **odbc** in the control panel of the Windows OS to search for the ODBC management program. Click **Set up ODBC data sources (64-bit)**.

Figure 3-86 Clicking Set up ODBC data sources (64-bit)



7. Add a driver.

In the dialog box displayed, click **Add**. In the **Create New Data Source** dialog box, click **openLookeng ODBC 1.9 Driver** and click **Finish**.

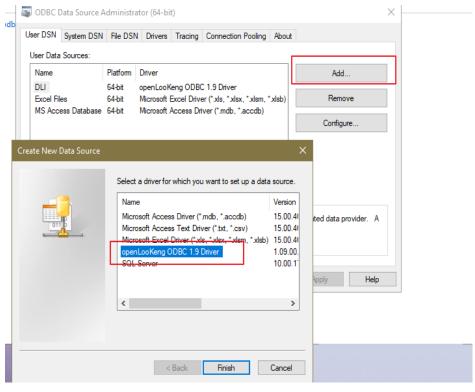
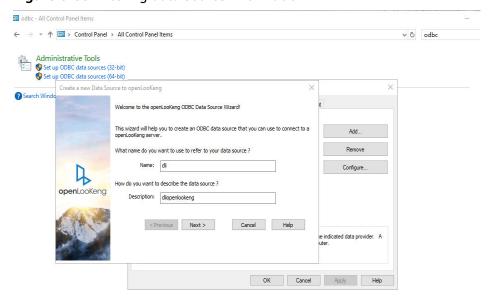


Figure 3-87 Adding a driver

Enter data source information.
 Enter the name and description by referring to Figure 3-88 and click Next.

Figure 3-88 Entering data source information



- 9. Configure related information.
 - Connect URL: indicates the address for accessing openLooKeng. Enter the value in the format of *IP address.Port*. For details about how to obtain the IP address and port, see Obtaining the IP address and port.

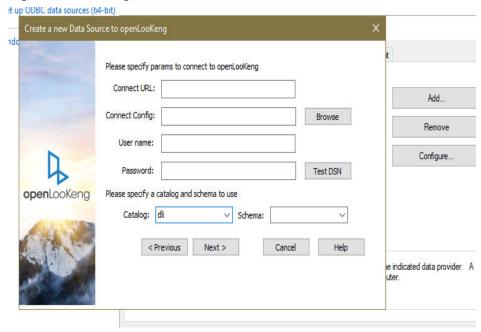
Connect Config: Select the jdbc_param.properties file in 4.

Example content of the properties file:
SSL=true
user=xxx_d21/xx_352221/xxxxacc00a2e2
password=xxxx12*

In the URL, **SSL=true** indicates that backend requests use HTTPS. Currently, Trino engine queues support only HTTPS connections.

- Set Catalog to dli.

Figure 3-89 Configuration information



Step 5: Testing the Connection

1. Click **Test DSN**. If the connection is successful, click **Next** > **Finish** to complete the test.

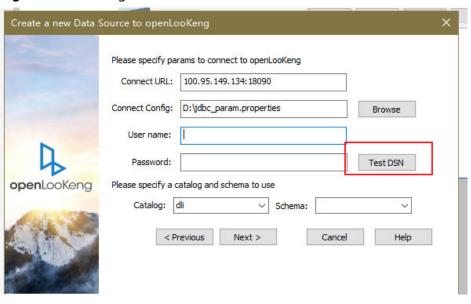
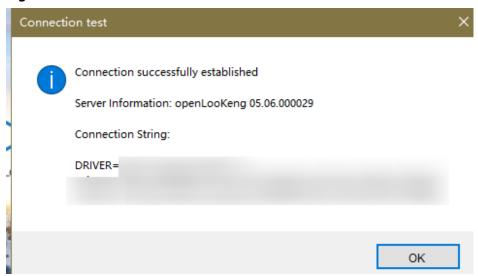


Figure 3-90 Testing the connection

Figure 3-91 Successful connection



Use Power BI to interconnect with DLI. Choose Get data > All > ODBC > Connect.

You need to enter the password for logging in to the DLI console for the first connection.

Get Data

All

All

Other

Certified Connectors Template Apps

Connect Cancel

Figure 3-92 Connecting to Power BI

Related Operations

Trino supports the SQL syntax. For details about the Trino SQL syntax, see
 Trino SQL Syntax. Currently, the Trino engine supports only the SELECT query operation.

4 Connections

4.1 Configuring the Connection Between a DLI Queue and a Data Source in a Private Network

Background

If your DLI jobs need to connect to external data sources, for example, MRS, RDS, CSS, Kafka, or GaussDB(DWS), you need to enable the network between DLI and the external data sources. DLI enhanced datasource connection uses VPC peering to directly connect the VPC networks of the destination data sources for point-to-point data exchanges.

This section provides a guide to help you connect to data sources. You can also refer to this section to rectify connection faults.

Development Process

Figure 4-1 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.

Table 4-1 Reference for creating clusters of other data sources

Service Name	Reference Documents
RDS	Getting Started with RDS for MySQL
GaussDB(DWS)	Creating a GaussDB(DWS) Cluster
DMS Kafka	Creating a Kafka Instance CAUTION When you create the instance, do not enable Kafka SASL_SSL.
CSS	Creating a CSS Cluster
MRS	Creating an MRS Cluster

CAUTION

- 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 4-2 Data source information to be obtained

Dat a Sour ce	Obtain Method
DMS Kafk a	 On the Kafka management console, click an instance name on the DMS for Kafka page. Basic information of the Kafka instance is displayed.
	 In the Connection pane, obtain the Instance Address (Private Network) value. In the Network pane, obtain the VPC and subnet of the instance.
	3. In the Network pane, obtain the security group of the instance.
RDS	On the Instances page of the RDS console, click the target DB instance name. In the displayed page, locate the Connection Information pane and obtain the Floating IP Address , VPC , Subnet , Database Port , and Security Group .

Dat a Sour ce	Obtain Method
CSS	 On the CSS management console, choose Clusters > Elasticsearch. On the displayed page, click the name of the created CSS cluster to view basic information.
	On the Cluster Information page, obtain the Private Network Address, VPC, Subnet, and Security Group.
Gaus sDB(DWS	1. 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.
)	2. On the Basic Information tab, locate the Database Attributes pane and obtain the private IP address and port number of the DB instance. In the Network pane, obtain the VPC, subnet, and security group information.

Dat	Obtain Method	
a Sour ce		
MRS	An MRS 3.x cluster is used as an example.	
HBa se	 Log in to the MRS management console, click a cluster name on the Clusters > Active Clusters page to view basic information. 	
	2. On the dashboard, obtain VPC, subnet, and security group from the Basic Information 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 root in SSH mode. For details, see Logging In to an ECS. 	
	 d. Run the following command to obtain MRS hosts information. Copy and save the information. cat /etc/hosts 	
	An example query result is as follows:	
	[root§node-master1kOno ~]# 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-master3tVbG.mrs-v08w.com node-master3tVbG.mrs-v08w.com. 192.168.0.238 node-group-1ySwO.mrs-v08w.com node-group-1ySwO.mrs-v08w.com. 192.168.0.123 node-master1kOno.mrs-v08w.com node-master1kOno.mrs-v08w.com. 192.168.0.124 node-group-1zKgA.mrs-v08w.com node-group-1zKgA.mrs-v08w.com. 192.168.0.71 node-master2qlhC.mrs-v08w.com node-master2qlhC.mrs-v08w.com. 192.168.0.71 node-master2qlhC.mrs-v08w.com node-group-1yRpv.mrs-v08w.com.	

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 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.

To obtain the security group information, go to the management console of the data source service and follow the steps provided in **Step 1: Obtain the Floating IP Address, Port Number, and Security Group of an External Data Source**.

4. In the **Inbound Rules** tab, add a rule to allow access from the queue network segment.

For details about how to set the inbound rule parameters, see **Table 4-3**.

Figure 4-2 Adding an inbound rule

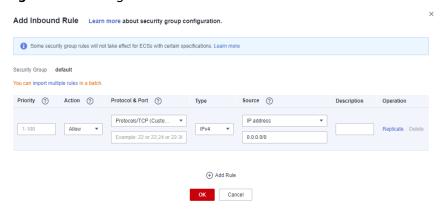


Table 4-3 Inbound rule parameters

Parameter	Description	Example
Priority	The security group rule priority.	1
	The priority value ranges from 1 to 100. The default value is 1, indicating the highest priority. A smaller value indicates a higher priority of a security group rule.	
Action	Action of the security group rule.	Select Allow .

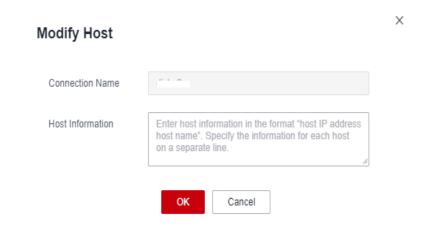
Parameter	Description	Example
Protocol &Port	 Network protocol: The value can be All, TCP, UDP, ICMP, or GRE. Port: Port or port range over which the traffic can reach your instance. The port ranges from 1 to 65535. 	In this example, select TCP. Leave the port blank or set it to the data source port obtained in Step 1: Obtain the Floating IP Address, Port Number, and Security Group of an External Data Source.
Туре	Type of IP addresses.	IPv4
Source	Allow access from IP addresses or instances in another security group.	In this example, enter the queue network segment 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

- 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 displayed dialog box, 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 4-3 Modifying host information



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 obtained IP address and port number of the data source in the address box, and click **Test**. If the queue passes the test, it can access the data source.

For MRS HBase, use **ZooKeeper IP address:ZooKeeper port** or **ZooKeeper host information:ZooKeeper port** for the test.

4.2 Configuring the Connection Between a DLI Queue and a Data Source in the Internet

Scenario

This section provides instructions to enable network connectivity for DLI queues to be accessed from the Internet. You can configure SNAT rules and add routes to the public network to enable communications between a queue and the Internet.

Procedure

Figure 4-4 Configuration process



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 4-5 Creating a VPC



Step 2: Create a Dedicated Queue

In this example, you will create a pay-per-use queue that uses dedicated resources.



The billing mode of the queue must be **Yearly/Monthly** or **Pay-per-use**. (If you select **Pay-per-use**, select **Dedicated Resource Mode** after you select a queue type.)

Enhanced datasource connections can be created only for yearly/monthly resources or pay-per-use resources in dedicated resource mode.

- 1. Log in to the DLI management console.
- Click **Buy Queue** in the upper left corner on the homepage page. On the displayed page, specify specifications and other required parameters.
 For details about the parameters for purchasing a queue, see **Creating a** Queue.

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 4-6 Creating an enhanced datasource connection

Create Enhanced Connection

After you create the enhanced datasource connection, the system will automatically create connection and required routes.



Step 4: Buy an EIP

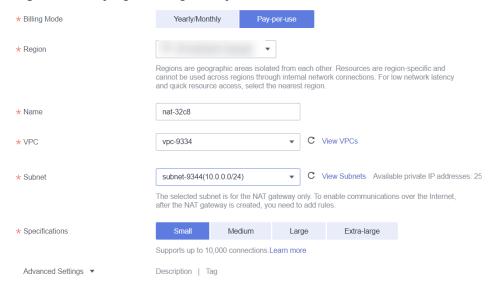
- 1. Log in to the **EIPs** page of the network console, click **Buy EIP**.
- 2. In the displayed page, configure the parameters as required. For details about how to set the parameters, see **Buy EIP**.

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.
- Click Buy Public NAT Gateway and configure the required parameters.
 For details, see Buying a Public NAT Gateway.

Figure 4-7 Buying a NAT gateway



3. Click **Next**, confirm the configurations, and click **Submit**.

◯ NOTE

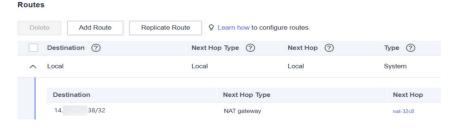
During the configuration, you need to set **VPC** to the one created in **Step 1: Create a VPC**.

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 4-8 Viewing the route



Step 3 Add an SNAT rule.

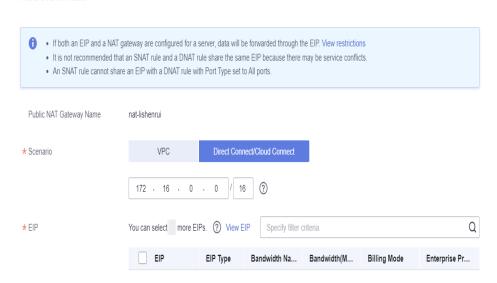
You need to add SNAT rules for the new NAT gateway to allow the hosts in the subnet to communicate with the Internet.

- 1. Click the name of the created NAT gateway on the **Public NAT Gateways** page of the network console.
- 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 4-9 Adding an SNAT rule

Add SNAT Rule



6. Click OK.

----End

Step 6: Adding 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**.

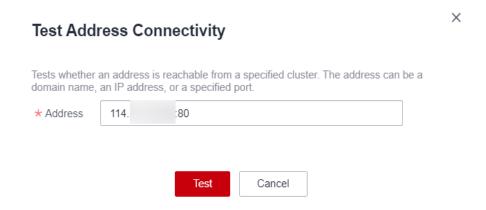
Figure 4-10 Adding route information for test



Step 7: Testing 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 4-11 Testing address connectivity



A Change History

Released On	What's New
2023-10-09	 Modified the following sections: Added the regions that support the Trino engine to Interconnecting FineBI with DLI Trino and Interconnecting Power BI with DLI Trino.
2023-08-19	Added the following section: Using DLI Flink SQL to Analyze e-Commerce Business Data in Real Time
2023-07-21	Added the following sections: Interconnecting FineBI with DLI Trino Interconnecting Power BI with DLI Trino
2023-03-09	Adjusted the document structure and moved the content related to DLI data development to <i>Data Lake Insight Development Guide</i> .
2023-03-02	Modified the prerequisites in Migrating Data from MRS Kafka to DLI.
2023-01-18	Added the description of configuring migration job scenarios to Migrating Data from MRS Kafka to DLI.
2023-01-06	Optimized the procedure for adding an inbound rule to the security group of the external data source to allow access from the DLI queue in Configuring the Connection Between a DLI Queue and a Data Source in a Private Network.
2022-10-31	Optimized the following sections and added information about solution advantages, process guidance, and resource planning and costs. • Analyzing Driving Behavior Data • Converting Data Format from CSV to Parquet

Released On	What's New
2022-09-27	Added Configuring the Connection Between a DLI Queue and a Data Source in the Internet.