# WeData Data Development Platform

# Quick Start

Service Notice

This document provides an overview of the as-is details of Tencent Cloud's products and services in their entirety or part. The descriptions of certain products and services may be subject to adjustments from time to time.

The commercial contract concluded by you and Tencent Cloud will provide the specific types of Tencent Cloud products and services you purchase and the service standards. Unless otherwise agreed upon by both parties, Tencent Cloud does not make any explicit or implied commitments or warranties regarding the content of this document.

Contact Us

We are committed to providing personalized pre-sales consultation and technical after-sale support. Don't hesitate to contact us at 4009100100 or 95716 for any inquiries or concerns.

# Contents

# Quick Start
# Overall Introduction

Last updated: 2025-04-18 15:44:55

> ⓘ **Note**
> The product features involved in this case are only the commonly used features within the module. If you need a more detailed understanding of all the features of the module, please view WeData Operation Guide.
> Besides the document, you can also learn the content of "Quick Start" through Tencent Cloud WeData Big Data Development and Governance Training Camp.

## Background

This document is the basic usage documentation of Tencent Cloud WeData. The target is to help you quickly understand WeData and have a basic concept of the entire process of business data processing.

This document takes the order data synchronization analysis scenario as an example, concatenates the features of each module such as data table structure design, data integration, data development, data quality, and data service, and helps you complete your initial experience of the WeData end-to-end process.

## Learning through This Document, You Can Understand the Following Content:

- Understand the entire process of business data development.
- Understand the role of each product module and upstream and downstream collaboration.
- Learn about basic concepts of data table structure design.
- Master the offline data synchronization process.
- Master the offline data development process.
- Master the data quality inspection process.
- Master the data service development process.

## Roles and Division of Roles Involved in This Documentation

- **Enterprise administrator:**
  - Responsible for registering and authenticating Tencent Cloud accounts.
  - Responsible for building the network environment.
  - Responsible for purchasing various cloud resources, including: EMR, WeData, MySQL, data service resources.
  - Responsible for creating sub-accounts, projects, and data tables.
  - Responsible for adding sub-accounts and binding data sources in WeData.
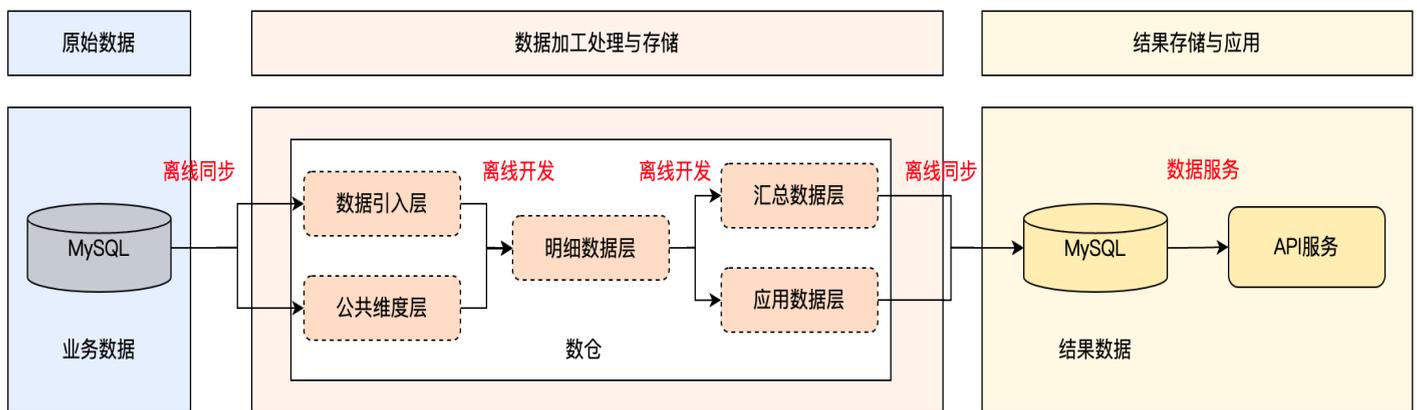- **data developer**

○ Responsible for data table structure design.

○ Responsible for the Data Integration module.

○ Responsible for the data development module.

○ Responsible for the data quality module.

○ Responsible for the data service module.

# Scenario Introduction

In a certain e-commerce platform, the business side **expects to learn about the sales performance of different categories in different cities through analysis of order data** so that it can adjust the operational promotion strategy targeting different cities.

# Holistic Picture of the Process

Throughout the entire product chain of WeData products, stages such as raw data research, data processing and storage, and result storage and application will be involved.

# Preparations

Last updated：2025-04-18 15:45:59

Before starting WeData data development, we need to make the following preparations first:

> **Note**
> The following steps involve resource purchasing and payment, which need to be performed by an enterprise administrator.

## Flowchart

**The specific operations include:**

| Steps | Description |
| --- | --- |
| Sign up for a Tencent Cloud account | • Sign up for a Tencent Cloud account.<br>• Real-name authentication (enterprise authentication recommended).<br>• Create a Tencent Cloud sub-account. |
| Prepare the network environment | • Create a new VPC.<br>• Bind the subnet.<br>• Apply for a public IP.<br>• Purchase a Public NAT Gateway, bind it to the VPC, and bind the public IP. |
| Prepare the engine resource environment (Taking Tencent Cloud EMR as an example) | • Purchase an EMR cluster.<br>• Purchase WeData, including integration resources and scheduling resources.<br>• Create a project and bind EMR, bind integration resources, and scheduling resources.<br>• In the WeData project, add a user. |
| Prepare the business data resource environment (Taking TencentDB for MySQL as an example) | • Purchase MySQL and initialize the business data.<br>• In the WeData project, add a data source. |

## Signing up for Tencent Cloud account

All the cloud resources involved in this tutorial are purchased through a Tencent Cloud account. Please use the same main Tencent Cloud account. If you already have a Tencent Cloud account and have completed enterprise authentication, please skip this step.

Role: Enterprise Administrator.

# Signing up for Tencent Cloud account

Go to Tencent Cloud registration page. You can register by scanning the QR code with WeChat or using an email. For the detailed registration process, please refer to Tencent Cloud Registration Guide.

# Enterprise Identity Verification Guide

After completing the registration, you need to go through enterprise authentication. The available authentication methods are as follows. For the detailed registration process, please refer to Tencent Cloud Enterprise Authentication.

| Authentication Method | Authentication Duration | Note |
|---|---|---|
| WeChat Public Platform Authentication | Instant completion | Enterprises that have registered a WeChat Official Account and completed WeChat identity verification can use this method for immediate authentication. |
| Enterprise legal person WeChat scan code authentication | Instant completion | Use the personal WeChat of the enterprise legal person to scan the code for authentication. After the legal person authorizes via WeChat scan, the authentication is completed. |
| Enterprise legal person Face Recognition authentication | Instant completion | Use the personal WeChat of the enterprise legal person to scan the code for Face Recognition. After passing Face Recognition, the authentication is completed. |
| Tencent Cloud recharge authentication | 1 business day | Transfer a small verification amount (less than 1 RMB) generated randomly by the system from the enterprise bank account (the amount will be added to the balance). Once Tencent Cloud receives the transfer, the authentication is completed. |
| Enterprise remittance authentication | 1–5 business days | Enter the enterprise bank account information. After Tencent Cloud successfully transfers the money, enter the transferred amount to complete the authentication. |

# Create a Tencent Cloud sub-account

1. Go to the **Tencent Cloud console >** User List **> Create User > Rapid creation.** Modify user permissions, click the **Edit Icon** for User Permissions.
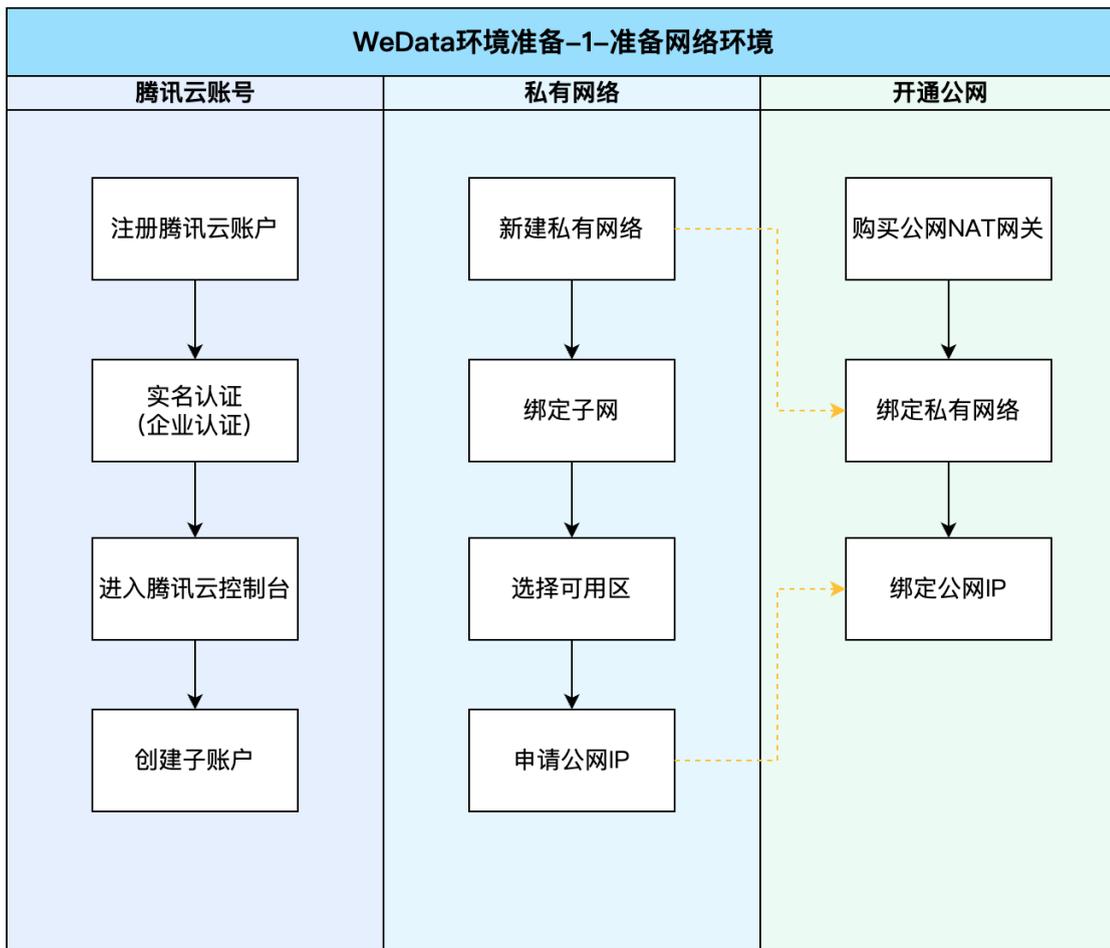
2. Enter wedata, click the search icon. After selecting all policies, click **OK**.

# Prepare the network environment

This tutorial involves multiple cloud resources. To ensure network connectivity, you need to set up a VPC environment.

- **Role:** Enterprise Administrator.
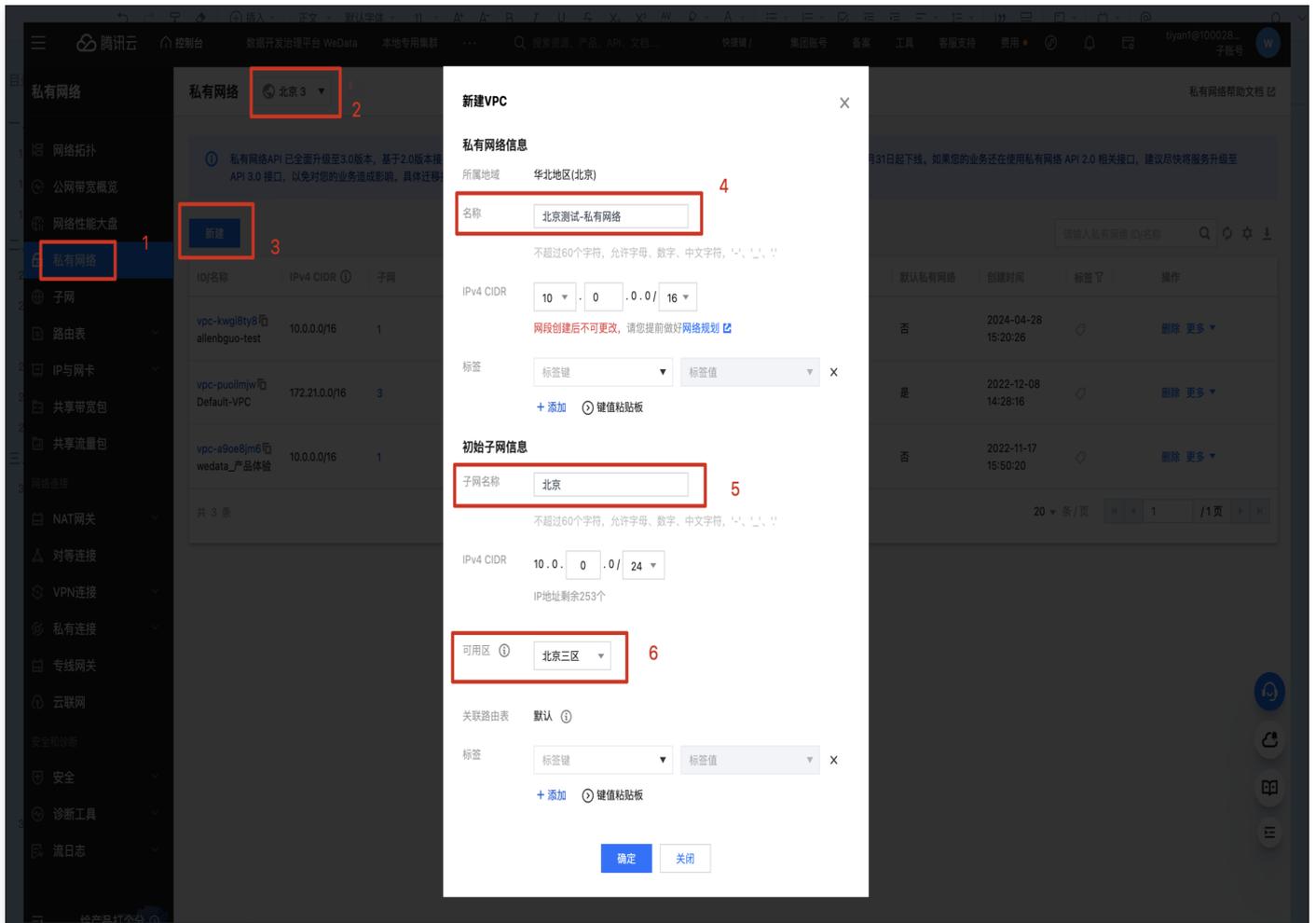- **Account:** Tencent Cloud primary account.
- **Steps:**

## Create a VPC

1. Log in to Tencent Cloud **VPC Console**. At the top of the **VPC** page, select the VPC's region, for example, select Beijing, and click **Create**.

2. Enter the new VPC interface, fill in the VPC information and initial subnet information. After completing the form, click **OK**.

   ○ VPC Name: You can name it anything for easy distinction. Example: Beijing – VPC.

   ○ Subnet Name: You can name it anything for easy distinction. It is recommended to match the optional zone below. Example: Beijing Zone 3.

   ○ Optional Zone: You can choose any zone, for example, Beijing Zone 3. When purchasing other resources later, if this zone is not available, you can add a subnet in the VPC.

> ⓘ **Note**
> - Choose the Beijing region. This is just an example. It is recommended to select a region that is closer to you.
> - All resources purchased later in this tutorial will be in the Beijing region, so please choose carefully.

## Purchase a NAT Gateway

1. Navigate to the Tencent Cloud Public Network NAT Gateway page, at the top of the NAT Gateway page, select the VPC region. For example: select Beijing, then click **Create.**

2. If you do not have a NAT Gateway, please go to the purchase page. After selecting the configuration, click **Enable Now**. Verify the bill, and complete the payment to enable it.

  ○ Region: Select Beijing,

  ○ VPC: Select the newly created VPC,

  ○ EIP: Select a new EIP. If you have already applied for a public IP, you can bind it here directly.

# Prepare the engine resource environment

WeData, as a data development and management platform, needs to be bound with Tencent Cloud Big Data Suite as the data storage and computing engine, such as Tencent Cloud EMR, DLC, TCHouse, etc.

In this tutorial, EMR is used as an example to introduce WeData's data synchronization and development process. Therefore, we need to purchase an EMR environment on Tencent Cloud first.

- **Role:** Enterprise Administrator.
- **Account:** Tencent Cloud primary account.
- **Steps:**

## Purchase EMR

1. Enter the Tencent Cloud  EMR purchase page , first select the software configuration. After selection, click  Next .

   ○ Region: Select North China – Beijing

   ○ Scenarios: Default scenario

   ○ Deployment Components: Choose Hive-3.1.3. In this tutorial, Hive is used as the storage and computing engine.

2. In the second step, select the region and hardware configuration. After selection, click **Next** .

○ Cluster Network: Select the newly created VPC.

○ Availability Zone: Select the availability zone where the subnet of the VPC is located. If it is not available, return to the **VPC** page to bind the subnet.

○ Security Group: Here, the default is to create a new security group.

3. Enter the node configuration page, expand the details, set the number of nodes, and use the default configuration.

4. Enter the basic configuration interface, set the server password, check **auto-renewal** and **Terms of Service**, then click **Purchase Now**. Verify the bill and complete the payment to activate. You can refer to Password Setting Format.



## Purchase WeData

1. Go to the Tencent Cloud WeData Purchase Page, complete the quick configuration, click **Purchase Now**, verify the bill, and complete the payment to activate.

   ○ Region: Select Beijing. It is recommended to choose a region closer to you.

   ○ Product Version: Select the professional version. For details, refer to WeData Version Differences.

   ○ Scheduling Resources: Select the test specification. For detailed information about scheduling resources, refer to Scheduling Resources Billing Overview.

   ○ Scheduling Resources Network: Select the newly created VPC.

   ○ Configuration Scheme: Select the basic specification. For detailed information about scheduling resources, refer to Integration Resources Billing Overview.

   ○ Network: Select the newly created VPC.

Tencent Cloud

# WeData数据开发治理平台 | 返回产品详情

**快速配置**　自定义配置

**购买须知**

温馨提示　本页面提供产品版本、调度及集成资源快速配置方案；更多版本及资源配置选择，请前往 自定义配置

**选择配置**

地域　　2　| **北京** | 广州 | 美国硅谷 | 上海 | 新加坡 | 上海金融 | 北京金融 | 香港 |

选择产品版本服务、调度及集成资源所在地域（了解详情 ☒），处于不同地域的云产品间网络不互通，创建成功后不可切换地域，请您谨慎选择。更多地域，请选择 自定义配置

产品版本　3

| 专业版 | 企业版 |
|---|---|
| • 完善数据开发与运维<br>• 基础数据治理能力 | • 智能高效数据开发与运维<br>• 全链路数据与成本治理 |

更多版本功能对比　了解详情 ☒

调度资源　调度资源用于调度离线开发任务（包括 SQL 类开发任务、Shell 任务、数据质量检测任务、元数据采集任务等），了解详情 ☒

配置方案　4

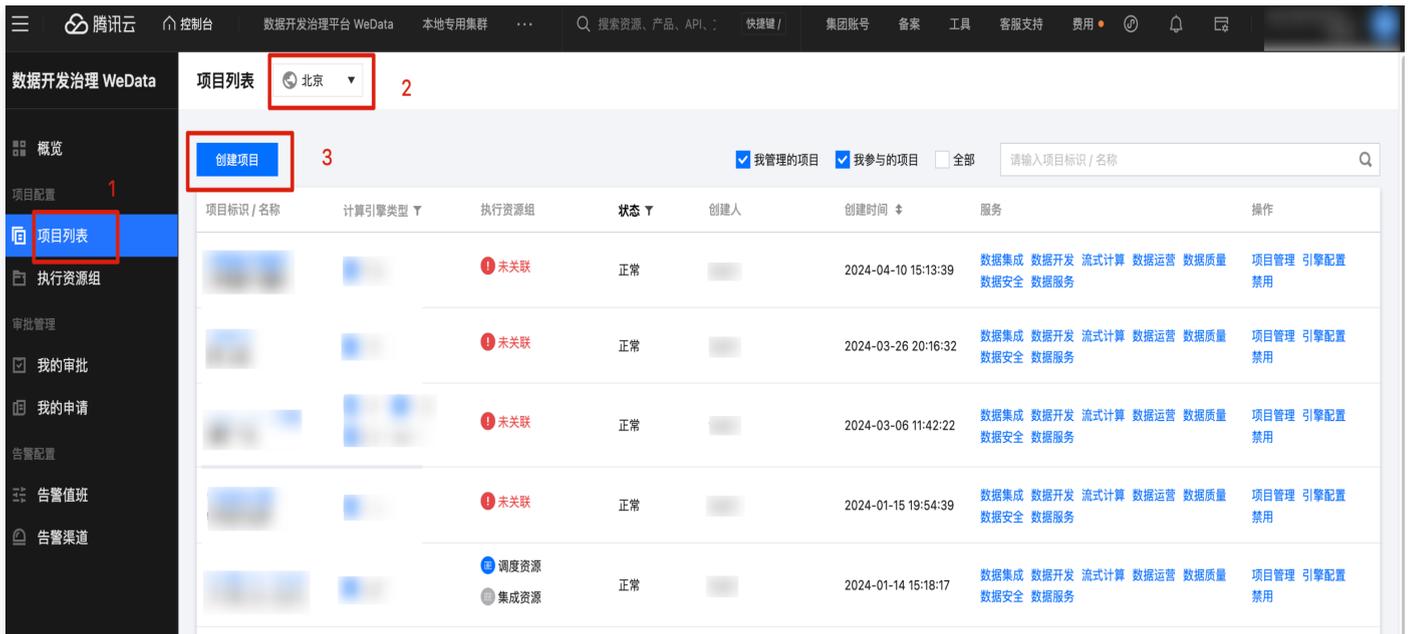| 测试规格 | 基础规格 | 普及规格 |
|---|---|---|
| • 适合测试，体验的场景<br>• 最大8并发实例数<br>• 100GB硬盘 | • 适合任务量与并发小的场景<br>• 最大16并发实例数<br>• 400GB硬盘 | • 适合任务量与并发适中的场景<br>• 最大32并发实例数<br>• 500GB硬盘 |

## Create a project in WeData

1. Log in to Tencent Cloud WeData Console, click **Project List** on the left-side menu, enter the Project List interface, select the top region as **Beijing**, then click **Create Project**.

2. Enter the project creation interface, select and fill in the relevant information, then click **Confirm** to complete the project creation.

- **Creation Method:** Select to create and configure the project.
- **Basic Information:**
  - Project Identification: Enter any text for easy distinction, e.g., test_bj_project.
  - Project Name: Enter any text for easy distinction, e.g., Beijing Test Project.
- **Configure Storage and Computing Engine:** Engine Region: Select Beijing.
- **Engine Type:** Select **EMR.**
  - EMR Cluster: Bind it directly from the dropdown. If you've completed the purchase in step 2.3.1, the EMR cluster name will be displayed here.
  - Account: The default is root.
  - Password: The one set during **Purchase EMR**.
  - Connectivity: Please click **Testing.**
  - Yarn Resource Queue: The default is `default`.
  - Engine Metadata Collection: Yes.

- ○ **Scheduling Resources:** Select **Immediately Associate,** check the previously created scheduling resources for binding **,** and the available resource groups will be displayed here.
- ○ **Integration Resources:** Select **Immediately Associate,** check the previously created integration resources for binding **,** and the available resource groups will be displayed here.

3. After the project is created, you can configure the account and add members by clicking **Project Management/Storage and Computing Engine Configuration** and **Project Management/Member Management.**



4. **Storage and Computing Engine Configuration:** Enter the storage and computing engine configuration interface, set EMR as the **primary account.**

5. **Member Configuration:** In the Member and Role Management page, click **Add,** enter the Add Member page, and add the Tencent Cloud sub-account as the project administrator.

> ⊘ **Note**
> This sub-account can perform subsequent data synchronization and data development operations.

# Prepare the business data resource environment

In this tutorial, we simulate an e-commerce mall order data synchronization and analysis scenario, so we need to prepare the original data of the e-commerce mall.

In this tutorial, Tencent Cloud MySQL is used as an example to introduce WeData's data synchronization process. Therefore, we need to purchase a MySQL database on Tencent Cloud first.

- **Role:** Enterprise Administrator.
- **Account:** Tencent Cloud primary account.
- **Steps:**

## Purchase MySQL

1. Go to the TencentDB for MySQL purchase page, complete the quick configuration, click **Purchase Now**, verify the bill, and complete the payment to activate.

   ○ Region: Select Beijing (this is just an example, you may choose a region closer to you).

   ○ Architecture: Select **Single-node**

   ○ Availability Zone: Select **Beijing Zone 3** , choose the availability zone where the subnet of the VPC is located.

○ Instance specifications: Select **Basic**



○ Network: Select the newly created VPC mentioned above.

○ Security Group: Select the default security group created earlier.

- ○ Character Set: Select **UTF-8**
- ○ Root Password: Set the root user's password.

2. After completing the settings, click **Next,** to confirm the configuration information.



3. Click **Buy Now** , verify the bill and complete the payment to activate.

# Initialize business data in MySQL

1. Enter the Tencent Cloud **database MySQL Console** , click the left-side menu **Instance List** to access the MySQL instance list page. Select the region at the top to **Beijing** , then click **Login to** .



2. Enter the DMC login interface, input your account and password, then click **Login to** .



3. In the DMC interface, select the top menu **Create** > **Create Database** .

4. Enter the Create Database interface, click **Create New Database** to access the new database creation page, fill in the database name, suggested to be 'emall'. After completing, click **OK**.



5. In the DMC interface, click the top menu **SQL Window** > **SQL** to enter the SQL interface. Quickly create tables by executing SQL statements.

   ○ Copy the following table creation SQL statements one by one. Each time you copy an SQL statement, click **Execute**. After execution, clear the SQL content before copying the next statement.

○ The specific table creation statements are as follows:

 ○ **Create City Table (cities)**

```sql
-- Create cities table in MySQL
CREATE TABLE cities (
    city_id INT NOT NULL AUTO_INCREMENT,
    city_name VARCHAR(50) NOT NULL,
    PRIMARY KEY (city_id)
) ENGINE=InnoDB DEFAULT CHARSET=utf8mb4;

-- Insert data
INSERT INTO cities (city_name) VALUES
('Beijing 市'),
('Shanghai 市'),
('Guangzhou 市'),
('Shenzhen 市'),
('Chengdu 市'),
('武汉市'),
('Nanjing 市'),
('杭州市'),
('Chongqing 市'),
('西安市'),
('苏州市'),
('Tianjin 市'),
('郑州市'),
('长沙市'),
('青岛市'),
('沈阳市');
```

○ **Create Product Category Table (categories)**

```sql
-- Create categories table
CREATE TABLE categories (
    category_id INT NOT NULL AUTO_INCREMENT,
    category_name VARCHAR(50) NOT NULL,
    PRIMARY KEY (category_id)
) ENGINE=InnoDB DEFAULT CHARSET=utf8mb4;

-- Insert data
INSERT INTO categories (category_name) VALUES
('电子产品'),
('家用电器'),
('服装鞋帽'),
('食品饮料'),
('图书音像'),
('运动户外'),
('家居建材'),
('母婴用品'),
('汽车用品');
```

○ **Create Product Table (products)**

```sql
-- Create products table
CREATE TABLE products (
    product_id INT NOT NULL AUTO_INCREMENT,
    category_id INT NOT NULL,
    product_name VARCHAR(100) NOT NULL,
    PRIMARY KEY (product_id),
    FOREIGN KEY (category_id) REFERENCES categories(category_id)
) ENGINE=InnoDB DEFAULT CHARSET=utf8mb4;

-- Insert data
INSERT INTO products (category_id, product_name) VALUES
(1, '智能手机'),
(1, '笔记本电脑'),
(1, '平板电脑'),
(2, '空调'),
(2, '洗衣机'),
(3, '男士外套'),
(3, '女士裙子'),
(4, '碳酸饮料'),
(4, '矿泉水'),
(5, '现代小说'),
(5, '历史书籍'),
(6, '跑步鞋'),
(6, '瑜伽垫'),
```

```
(7, '实木家具'),
(7, '床上用品'),
(8, '婴儿奶粉'),
(8, '儿童玩具');
```

○ **Create Order Table (orders)**

```sql
-- Create orders table
CREATE TABLE orders (
    order_id INT NOT NULL AUTO_INCREMENT,
    product_id INT NOT NULL,
    quantity INT NOT NULL CHECK (quantity > 0),
    unit_price DECIMAL(10, 2) NOT NULL,
    amount DECIMAL(10, 2) NOT NULL,
    order_time DATETIME NOT NULL,
    shipping_city_id INT NOT NULL,
    shipping_address TEXT NOT NULL,
    PRIMARY KEY (order_id),
    FOREIGN KEY (product_id) REFERENCES products(product_id),
    FOREIGN KEY (shipping_city_id) REFERENCES cities(city_id)
) ENGINE=InnoDB DEFAULT CHARSET=utf8mb4;

-- Inserting data
INSERT INTO orders (product_id, quantity, unit_price, amount,
order_time, shipping_city_id, shipping_address) VALUES
(1, 1, 4999.00, 4999.00, '2024-04-01 10:00:00', 1, 'Beijing Haidian
District'),
(2, 1, 6999.00, 6999.00, '2024-04-02 11:00:00', 2, 'Shanghai Pudong
New District'),
(3, 2, 3999.00, 7998.00, '2024-04-03 12:00:00', 3, 'Guangzhou Tianhe
District'),
(4, 1, 5999.00, 5999.00, '2024-04-04 13:00:00', 4, 'Shenzhen Nanshan
District'),
(5, 1, 999.00, 999.00, '2024-04-05 14:00:00', 5, 'Chengdu Wuhou
District'),
(6, 1, 699.00, 699.00, '2024-04-06 15:00:00', 6, 'Wuhan Jianghan
District'),
(7, 1, 2999.00, 2999.00, '2024-04-07 16:00:00', 7, 'Nanjing Gulou
District'),
(8, 1, 3999.00, 3999.00, '2024-04-08 17:00:00', 8, 'Hangzhou West Lake
District'),
(9, 1, 4999.00, 4999.00, '2024-04-09 18:00:00', 9, 'Chongqing Yuzhong
District'),
(10, 1, 1999.00, 1999.00, '2024-04-10 19:00:00', 10, 'Xi'an Beilin
District');
```

# Bind MySQL in WeData

1. Log in to Tencent Cloud **WeData Console**, click on the left menu **Project List**, select the top region as Beijing, in the corresponding project operation column, click **Project Management.**



2. In the Data Source Management interface, click **New Data Source**, select the data source type as MySQL, click **Next**.

3. Enter the New MySQL Data Source page, fill in the relevant information, and click **Save** .

- ○ Data Source Name:
  - ○ Data Source Name: Fill in as required for easy identification, for example: bj_mall.
  - ○ Display Name: Fill in as required for easy identification, for example: Beijing—Test—Mall.
- ○ Instance Information:
  - ○ Region: Select **Beijing.**
  - ○ Select Instance: Select from the drop-down list.
- ○ Table Connection Information:
  - ○ Database Name: emall
  - ○ Username: root
  - ○ Password: Enter the password set when purchasing MySQL mentioned above.

All preparations are now complete. We will now officially start the data synchronization and data development sections of this tutorial. Subsequent operations can be performed using a Tencent Cloud Sub-account.

# Data Table Structure Design

Last updated：2025-04-18 15:46:32

## Business Research

### Source Data Storage Location

Through researching the existing technical architecture of the mall system, it was found that the data is stored in a MySQL database.
Here, it is assumed that the TencentDB for MySQL database is used.

### Target Business Scenario Analysis

By analyzing the target business scenarios: sales performance in various cities and categories, we need to obtain the following tables:

- Order Table: (At this point, ignore the design of subtables such as order details, assuming the order table includes product ID, product quantity, product price, shipping address, order time, etc.).
- Product Table: (At this point, ignore the design of subtables such as SKU, assuming the product table includes product ID, product category, etc.).
- City Table: (Assuming the geographic location coding table only goes down to the city level, the city table includes: city code, city name).
- Product Category Table: (Assuming there is only one level of category, the category table includes: category code, category name).

### Actual Structure

The following are the actual structures of the researched Order Table and Product Table:

### 1. Order Table (orders)

| Field name | Field Type | Field Length | Field Description | Sample Code |
|---|---|---|---|---|
| order_id | INT | 10 | Order ID, primary key, auto-increment | 10001 |
| product_id | INT | 10 | Product ID, foreign key | 1001 |
| quantity | INT | 5 | Quantity of Goods, positive integer | 2 |
| unit_price | DECIMAL(10,2) | – | Unit Price of Goods, round to two decimal places | 99.99 |
| amount | DECIMAL(10,2) | – | Subtotal Amount of Goods, which is quantity multiplied by the unit price | 199.98 |

| order_time | DATETIME | – | Order Time, accurate to the minute | "2024-04-04 10:30:00" |
| shipping_city_id | INT | 10 | Shipping Address City ID, foreign key | 1101 |
| shipping_address | TEXT | – | Shipping Address, includes province, city, district, and detailed address | "Chaoyang District, Beijing, XXX Community" |

## 2. Product Table (products)

| Field Code | Field Type | Field Length | Field Description | Sample Code |
| --- | --- | --- | --- | --- |
| product_id | INT | 10 | Product ID,primary key,Auto-increment | 1001 |
| category_id | INT | 10 | Category ID,Foreign key | 101 |
| product_name | VARCHAR(100) | – | Product Name | "Smartphone" |

## 3. City Table (cities)

| Field Code | Field Type | Field Length | Field Description | Sample Code |
| --- | --- | --- | --- | --- |
| city_id | INT | 10 | City Code,primary key,Auto-increment | 1101 |
| city_name | VARCHAR(50) | – | City Name | "Beijing City" |

## Product Category Table (categories)

| Field Code | Field Type | Field Length | Field Description | Sample Code |
| --- | --- | --- | --- | --- |
| category_id | INT | 10 | Category ID,primary key,Auto-increment | 1 |
| category_name | VARCHAR(50) | – | Category Name | "Electronic Products" |

# Architecture Design

Based on business scenario needs, the final business output involves Data Warehouse Layering and Data Table Structure.

## Model Specification

Model specifications help teams unify data warehouse design rules, streamline the data development process, better accumulate data assets, and lay a foundation for building data services and data marts. During the design of data warehouse model specifications, multiple categories are included, such as data domain and principal domain.

In this scenario, the core objective is DataInLong and the data development process, and therefore detailed teaching of data model specifications is not covered in this tutorial.

Below are examples of model specifications related to this scenario:

| Category | Chinese Description | English Name |
|---|---|---|
| Business Category | Sales | trade |
| Data Domain | Order<br>Product | order<br>product |
| Business process | Order Creation | ordercreate |
| Subject Domain | Product | product |
| Dimension | Date<br>Region<br>Category | date<br>city<br>category |
| Metrics | Sales Volume<br>Sales Quantity | amount<br>quantity |

# Data Warehouse Layering

## 1. Data Ingestion Layer ODS

Import raw data that hasn't undergone any processing into the data warehouse. The table structure in the ODS layer is consistent with the table structure in the original data system.

Therefore, we need to create 4 Hive tables based on the raw data (table creation operations are not needed here; this will be covered in subsequent lessons), with table structures identical to the MySQL source data tables.

The naming of the four tables is as follows:

- Order Table: ods_order_order
- Product Table: ods_product_product
- Category Table: ods_product_category
- City Table: ods_order_city

> ⓘ **Note**
> The suggested naming convention is: ods_{data domain}_{self-definition content}.

## 2. Common Dimension Layer DIM

This section focuses on data synchronization logic, temporarily ignoring the design of the dimension layer.

> ⓘ **Note**
>
> It is recommended to redundantly store the field attributes from the dimension tables in the detailed data tables.

## 3. Detailed Data Layer DWD

Build the most granular detailed data table. It is advisable to redundantly store some fields in this table to reduce the association between the detailed data table and the dimension table.

- Build a detailed table: Table creation is not required here; it will be covered in subsequent lessons.
- Product Sales Detail Report: dwd_trade_order_ordercreate_productsales.

> ⓘ **Note**
>
> The suggested naming convention is: dwd_{business category}_{data domain}_{business process}_{self-definition content}.

| Field Code | Field Type | Field Length | Field Description | Sample Code |
|---|---|---|---|---|
| order_id | INT | 10 | Order ID, primary key | 10001 |
| product_id | INT | 10 | Product ID | 1001 |
| category_id | INT | 10 | Category ID | 101 |
| category_name | STRING | 50 | Category Name | "Electronic Products" |
| product_name | STRING | 50 | Product Name | "Smartphone" |
| quantity | INT | 5 | Quantity of Goods, positive integer | 2 |
| unit_price | DECIMAL | 10,2 | Unit Price of Goods, round to two decimal places | 99.99 |
| amount | DECIMAL | 10,2 | Subtotal Amount of Goods, which is quantity multiplied by the unit price | 199.98 |
| order_time | DATETIME | – | Order Time, accurate to the minute | "2024-04-04 10:30:00" |
| shipping_city_id | INT | 10 | Shipping Address City ID, foreign key | 1101 |
| shipping_city_ | STRIN | 50 | City Name | "Beijing City" |

| name | G | | | |
| --- | --- | --- | --- | --- |
| shipping_address | TEXT | – | Shipping Address, includes province, city, district, and detailed address | "Chaoyang District, Beijing, XXX Community" |
| pt_date | STRING | 50 | Partition Field | "2024-04-01" |

## Additional Notes: Hive Table Partitioning

### Partitioned Table Overview

Partitioning is an essential database optimization technique that improves performance, simplifies management, reduces costs, and enhances data availability and security by dividing datasets into smaller, logically independent parts.

Partitioning is especially important in big data scenarios.

### Benefits of Hive Table Partitioning

The benefits of partitioning storage in Hive can be reflected in several aspects:

| Advantage | Description |
| --- | --- |
| Improve Query Performance | By storing data in different partitions, queries can target specific partitions, avoiding the need to scan the entire table's data and significantly reducing query time. |
| Optimize Data Management | Partitioning is a logical way of organizing data, making it easier to maintain, clean, and perform bulk operations such as backup and recovery. |
| Horizontal Scaling | Partitions can horizontally distribute data storage pressure, physically distributing data to different storage units, enhancing system scalability. |
| Reduce Data Skew | In cases of uneven data distribution, partitioning can prevent data skew by avoiding situations where some partitions have too much data while others have too little. |
| Data Isolation | Partitions can be used for data isolation. For example, data can be divided into different partitions based on time, facilitating version control and historical data management. |
| Reduce Data Loading Time | During data loading or ETL processes, data can be loaded into specific partitions more quickly without operating on the entire table. |
| Save Storage Space | Partitions can help delete or archive old partition data, thereby saving storage space. |
| Parallel Processing | Partition tables can better leverage Hadoop's MapReduce parallel processing capabilities, as queries can be executed in parallel on different partitions. |

| Data Security and Access Control | Partitions can be used to implement finer-grained data access control, such as setting stricter access permissions for certain partitions. |
|---|---|
| Maintain Data Integrity | Partitions can ensure data integrity because each partition can have its own data integrity constraints. |
| Support Hot and Cold Data Layering | By partitioning, data can be classified into "hot data" and "cold data" based on its frequency of use, and different storage strategies can be applied. |
| Simplify Data ETL Process | During data extraction, transformation, and loading, partitions can simplify data organization and processing workflows. |
| Improve Data Availability | Partitions can improve data availability, as the unavailability of one partition does not affect access to other partitions. |

Therefore, plan partition fields as early as possible when creating Hive tables.

## 4. Summary Data Layer DWS

Construct summary indicator data tables with business-use granularity.
- Build a summary table: Table creation is not required here; it will be covered in subsequent lessons.
- Daily Product Sales Summary Table: dws_trade_order_productsales_1d.
- The suggested naming convention is: dws_{business category}_{data domain}_{self-definition content}_{time period}.

| Field Code | Field Type | Field Length | Field Description | Sample Code |
|---|---|---|---|---|
| order_date | DATE | – | Date | 2021-04-01 |
| city_id | INT | 10 | City ID | 1 |
| category_id | INT | 10 | Category ID | 1 |
| city_name | STRING | 50 | City Name | " Beijing " |
| category_name | STRING | 50 | Category Name | "Electronic Products" |
| quantity | INT | 10 | Total Product Sales | 100 |
| amount | DECIMAL | (10, 2) | Total Product Sales Volume | 9999.99 |
| pt_date | STRING | 50 | Partition Field | "2021-04-01" |

# 5. Application Data Layer (ADS)

Build an indicator table for final business analysis requirements. Since this scenario is relatively simple, it will be temporarily ignored here.

# Data Integration

Last updated：2025-04-18 15:48:35

We will sync the raw data to the data warehouse in this step.

## Add a New Data Source

### Raw Data Source: MySQL

We have bound the data source to the project. It can be ignored here.

### Target Data Source: Hive

After you bind the storage-compute engine EMR, the system will collect the Hive data source in the EMR cluster within 10 minutes. Therefore, it is not necessary for you to actively bind the Hive data source.

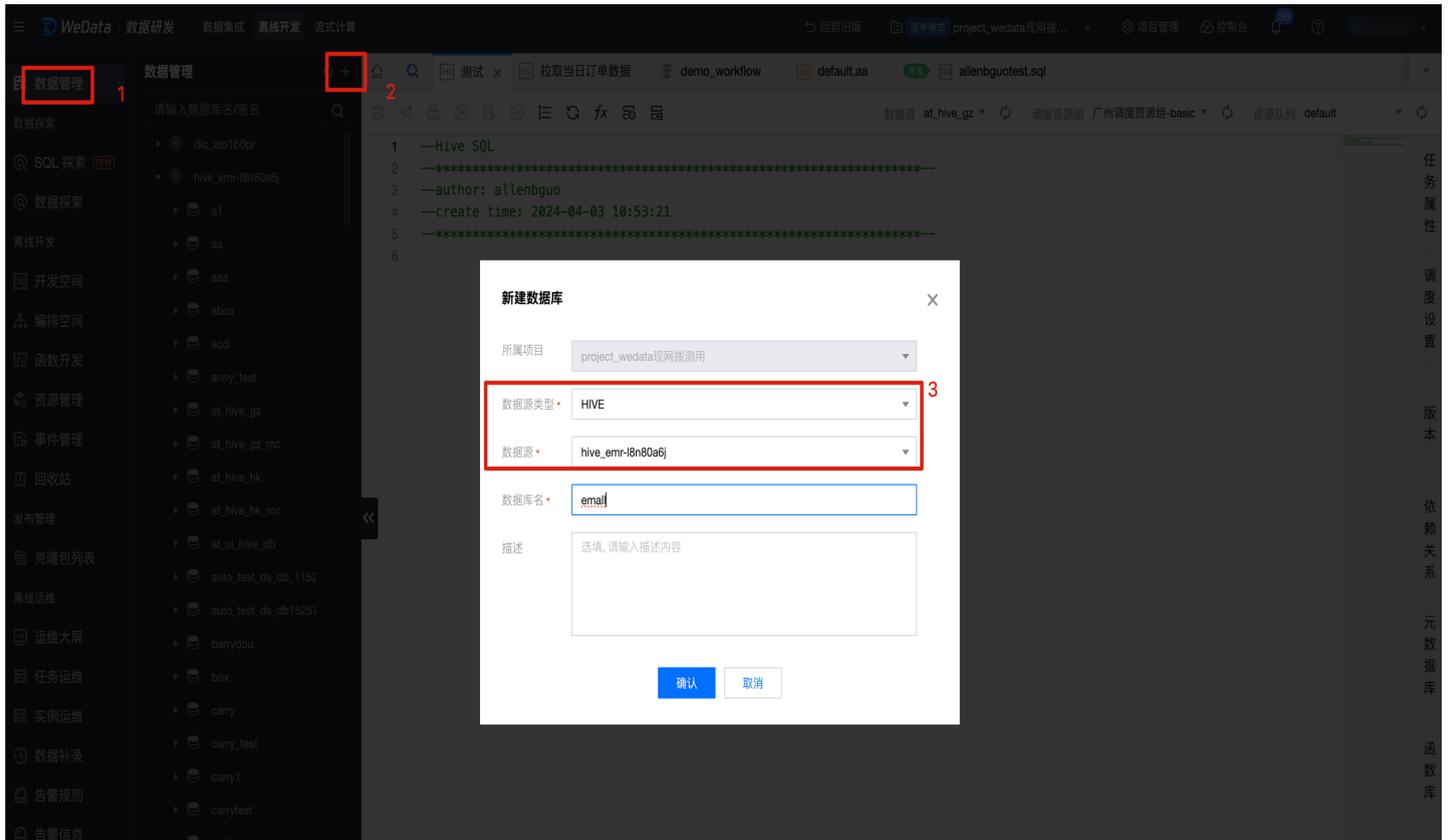However, we need to **create a database** in the Hive data source to store the collected original data.

## Creating a Database

Enter **Data Development** > **Data Management**, click + to create a database, select and fill in the required content, and after completing, click **OK**.

- Select data source type: Hive.
- Data source selection: hive_emr-XXX.

> ⓘ **Note**
>
> After you bind the storage-compute engine EMR, the system will collect the Hive data source in the EMR cluster within 10 minutes.
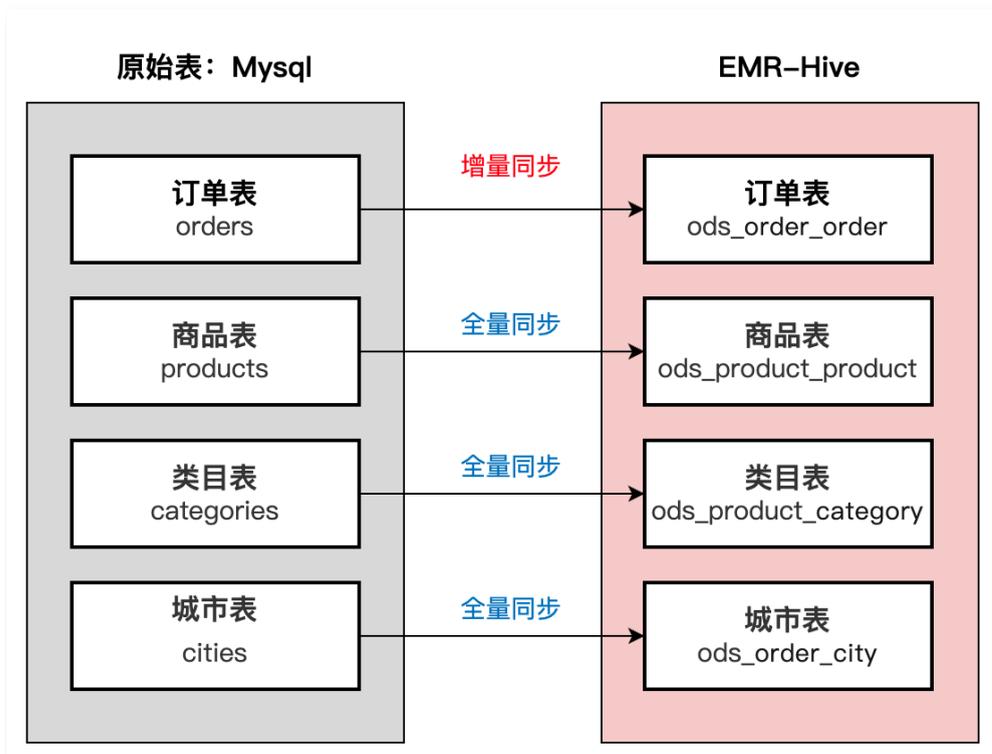
# Offline Synchronization Task Design

We will now create an offline synchronization task to synchronize the original data in the MySQL data source to the Hive table in the EMR cluster.

We have already known from the above operations that we need to synchronize 4 raw data tables, which are:

| No. | Table Name | Raw Data Source: MySQL | <Target Data Source: Hive Table Name> |
|-----|------------|------------------------|----------------------------------------|
| 1 | Order Table | orders | ods_order_order |
| 2 | Product Table | products | ods_product_product |
| 3 | Category Table | categories | ods_product_category |
| 4 | City Table | cities | ods_order_city |

The task development plan is designed as follows:

## Additional Notes

Difference between full synchronization and incremental synchronization:

| Name | Description |
|---|---|
| | **Definition:** Full synchronization refers to the process where the system transmits all data from two databases or data warehouses during each sync operation. |
| | **Use cases:** Full synchronization is usually suitable for small amounts of data or infrequent data change situations, as well as during initial sync or data migration. |
| Full synchronization | **Strengths:**<br>• Simple and easy to implement: No need to track data changes, just directly copy all data.<br>• Integrity: Underwrite the integrity and consistency of data, because all data is resynchronized. |
| | **Drawbacks:**<br>• High time and resource consumption: A large amount of data needs to be transmitted, which takes a long time and occupies bandwidth.<br>• High cost: For large data volume, more storage and computing resources may be needed. |
| Incremental synchronization | **Definition:** Incremental synchronization only synchronizes the data that has changed since the last synchronization, rather than all data. |

Use cases: Suitable for environments with large volumes of data or frequent updates.

**Strengths:**
- High efficiency: Synchronize only the changed data, saving time and bandwidth.
- Low cost: Reduces the need for storage and computing resources.
- Real-timeness: Can reflect the latest state of data more quickly.

**Drawbacks:**
- High complexity: A mechanism is required to track and record data changes.
- Consistency issues may exist: If problems occur during the sync process, data inconsistency may result.

**Summary:**
The choice of sync method depends on the specific application scenario and requirements. For small amounts of data with infrequent changes, full synchronization may be simpler and more efficient. Conversely, for large amounts of data with frequent updates, incremental synchronization can significantly improve efficiency and reduce costs. In practical applications, these two strategies are sometimes used in conjunction, for example, regularly conducting full synchronization to ensure data integrity, while using incremental synchronization in daily operations to improve efficiency.
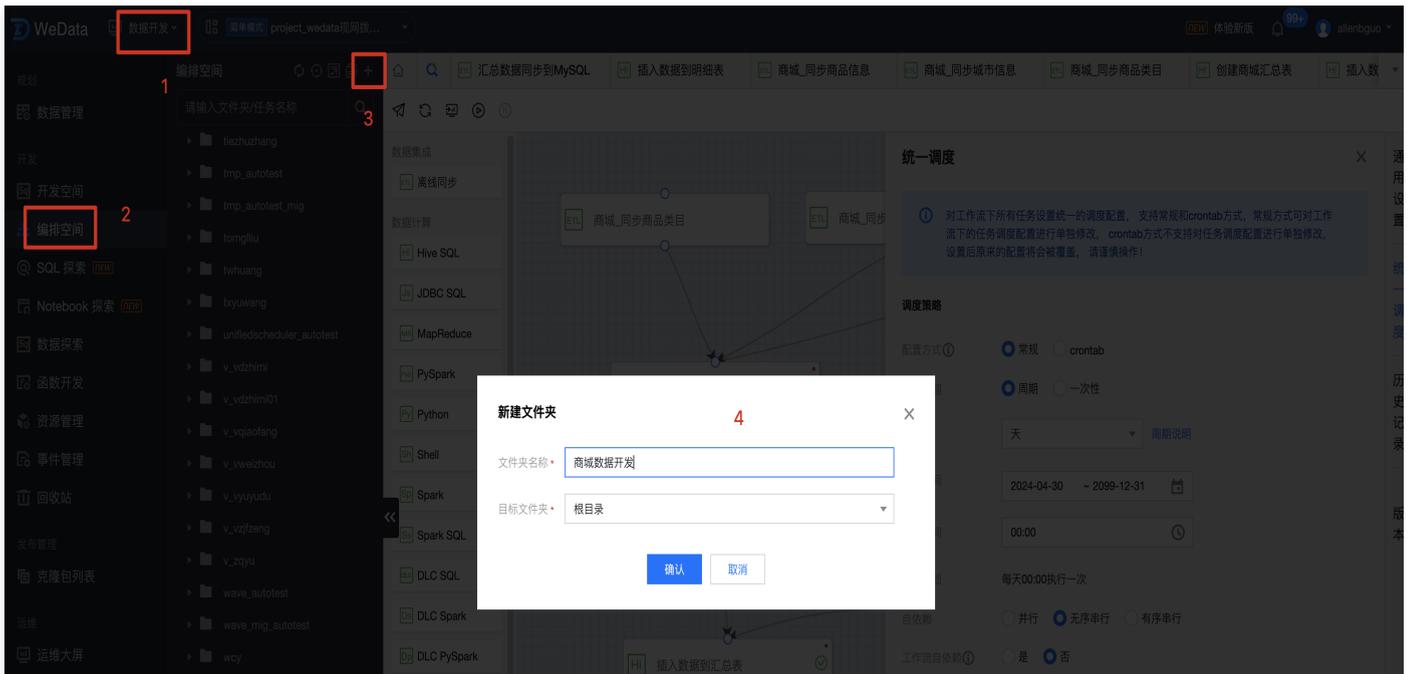
# Offline Synchronization Task Development

## Creating Workflow

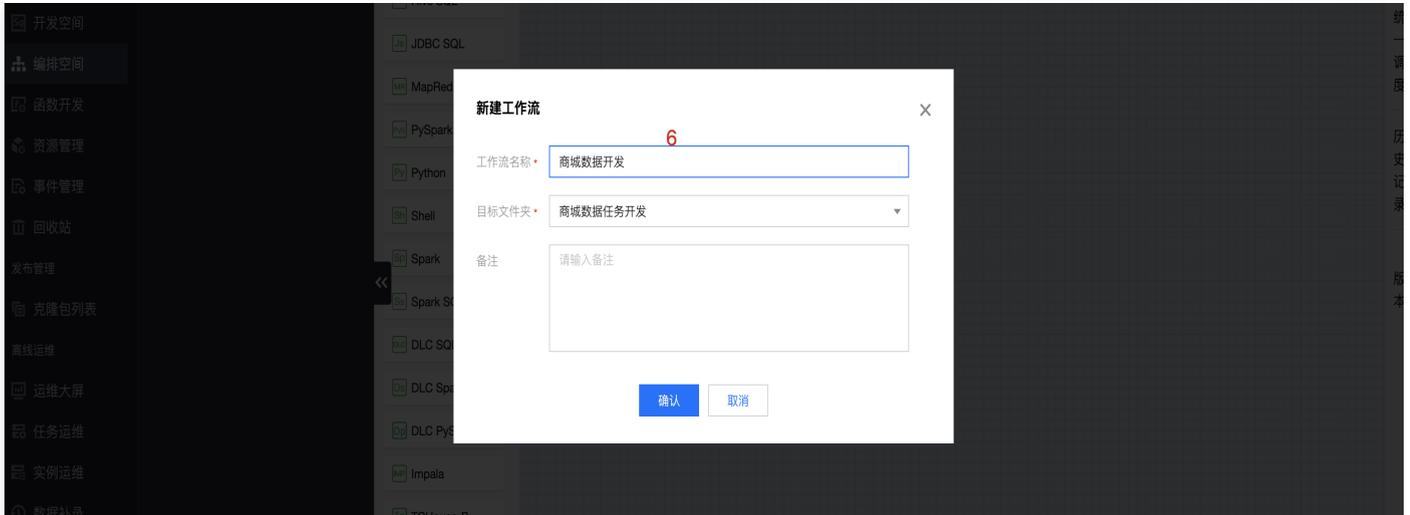1. Enter the offline development page from the project list .

2. Select offline **Data Development** > **Orchestration Space**, click + to create a new folder (named: Mall Data Task Development), and store subsequent development tasks.



3. Find the **created folder** > Right-click **create workflow**.
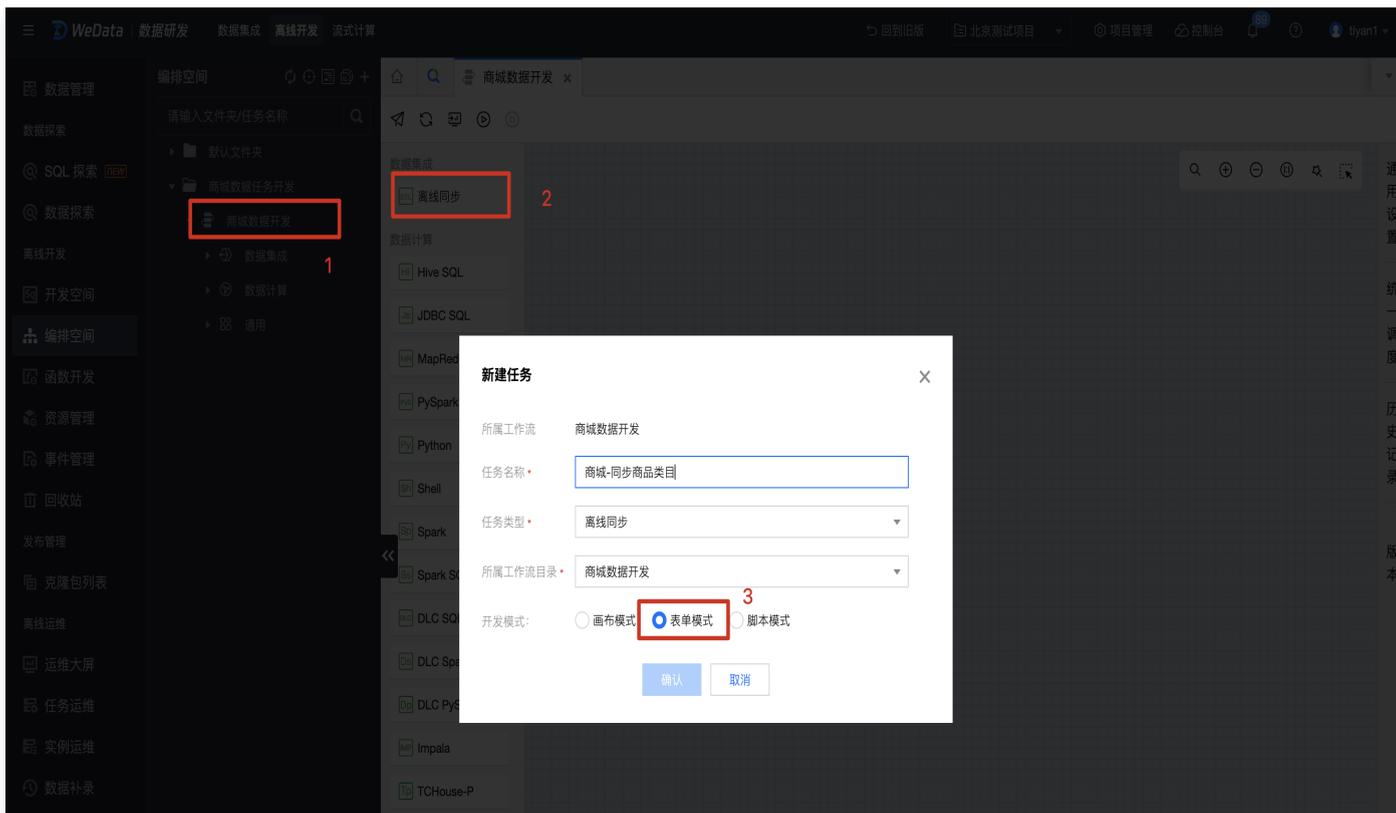
4. Create a new workflow named Mall Data Development, select the corresponding target folder, and click **Confirm**.



## Synchronize Category Table

First, we will synchronize the category table from MySQL to the Hive table.

1. Find the newly created **Mall Data Task Development** > **Mall Data Development** in the Orchestration Space. Click **Offline Sync.** Select the configuration mode (Task Name: Mall_Synchronize Product Category, Development Mode: Select **Form Mode** here), and click **Confirm.**

2. Configure the raw data source.

2.1 Select the database and table where the raw data table is stored. Please select the MySQL data source added in the previous step. In the following text, when we introduce order table sync, we will introduce how to achieve incremental data synchronization by setting filter conditions.
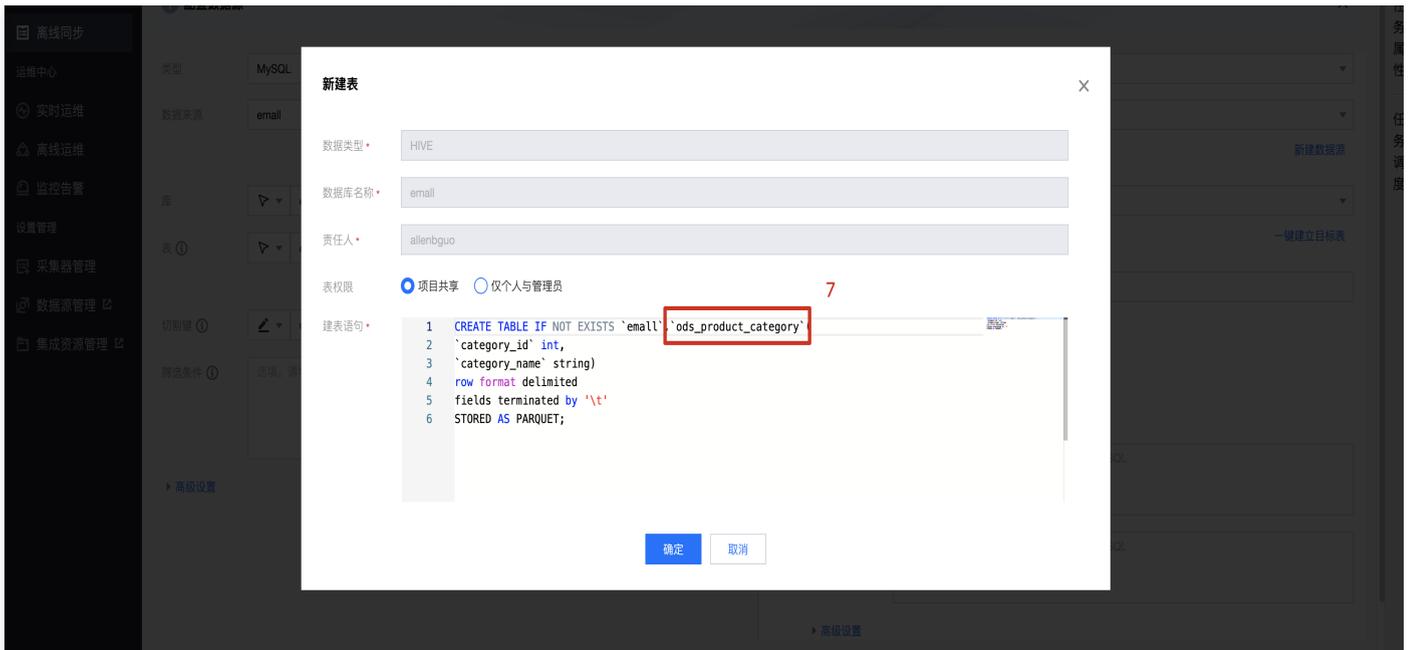
> ⚠ **Notes:**
>
> Since the category table is basic information and the amount of data is generally small, we do not need to set filter conditions (i.e., Where statements) here.

2.2 Configure the target data source. Select the Hive table where you need to store data. Please select the Hive data source and database added in the previous step.

- ○ Type: Hive
- ○ Data destination: search hive_emr
- ○ Database: emall (the database created in the previous step)

2.3 Establish a target table. Here, use **create target tables with one click** to replicate the MySQL Table Structure.
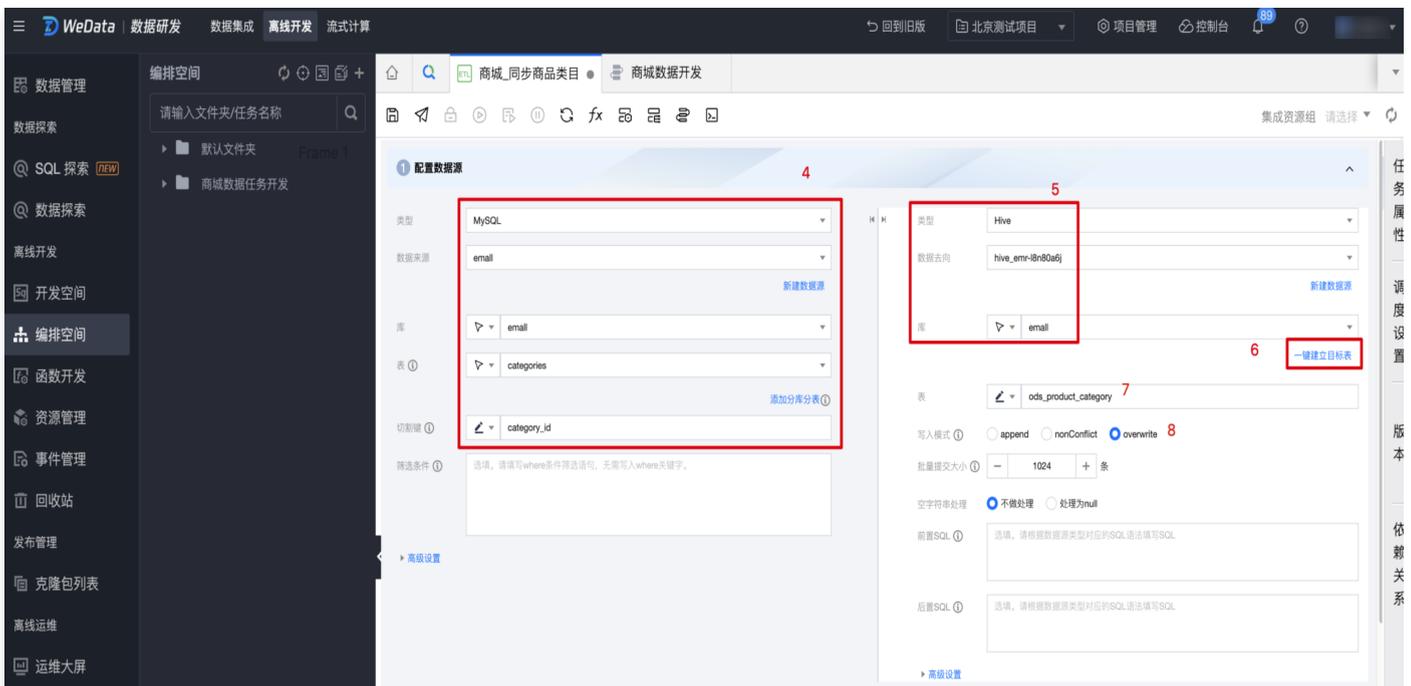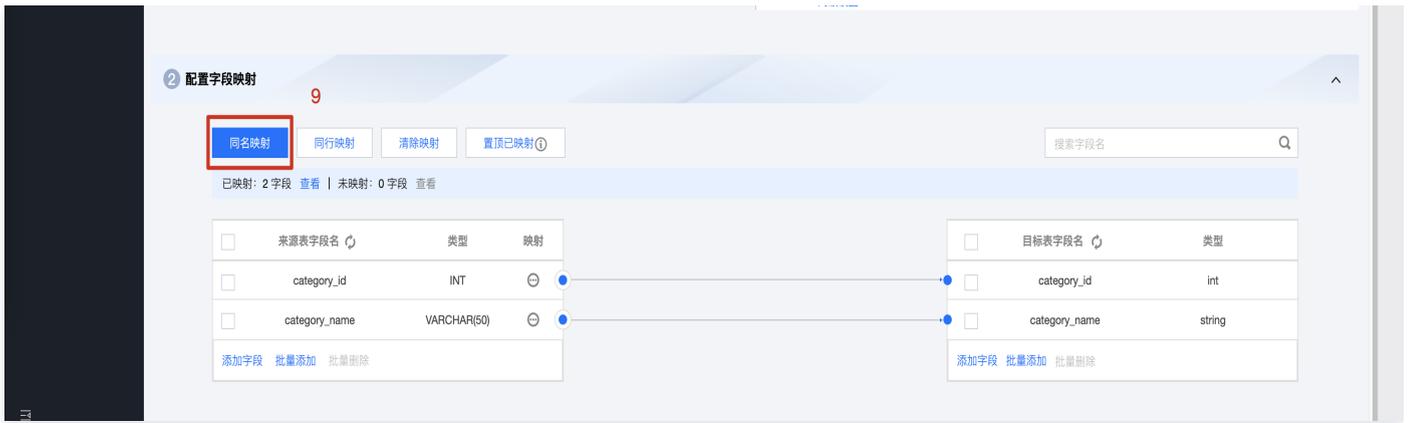
> ⚠ **Notes:**
>
> Please modify the table creation statement in the pop-up box. Modify table name:
> ods_product_category.

2.4 Selection table: ods_product_category

2.5 Since the category table is basic information, here we select **overwrite,** that is, overwrite update
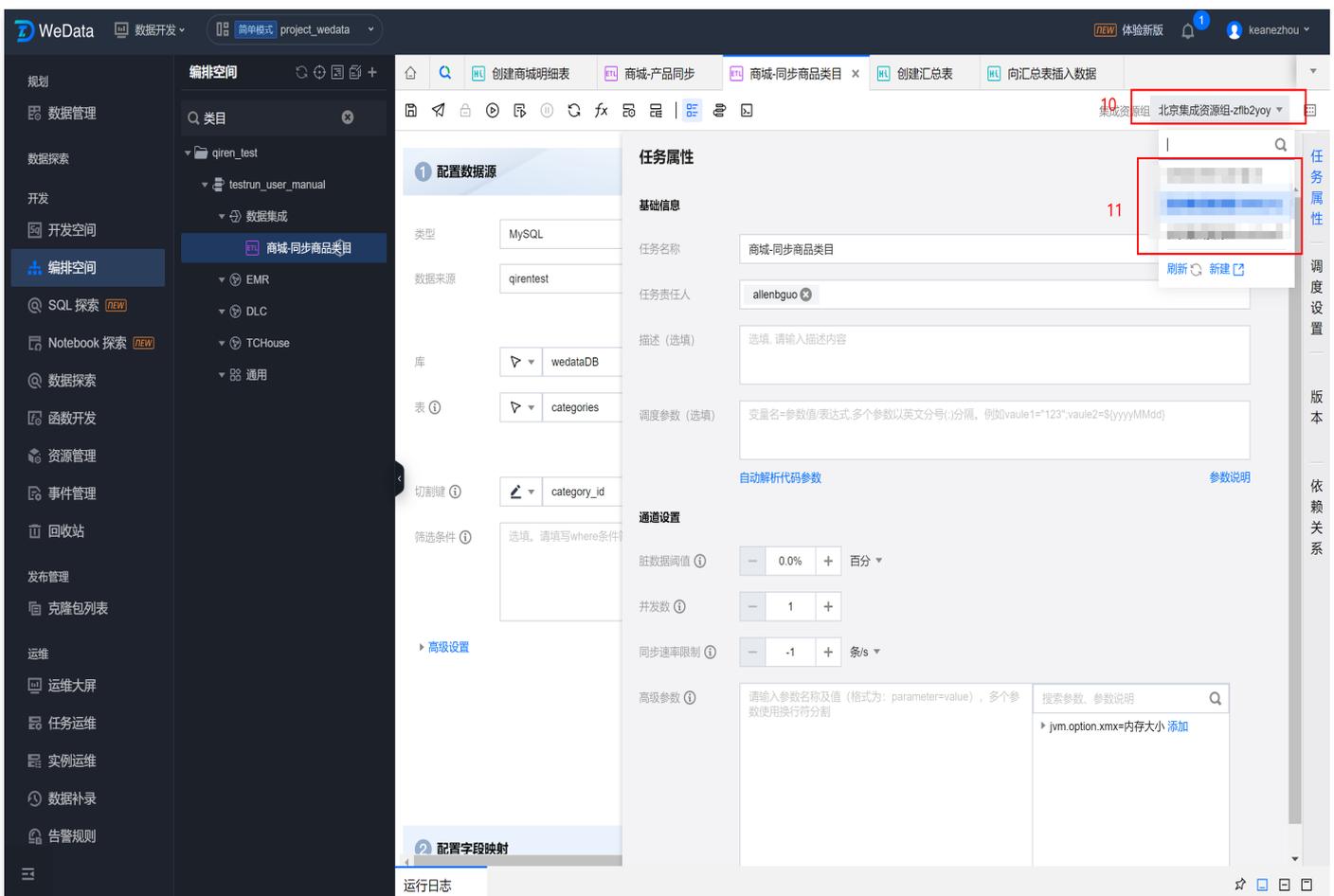  every time.



3. Configure field mappings. Here, we need to map the fields of the source table to the fields of the
  target table one-to-one. Since our table structures are identical, we can use **same name mapping.**

4. Set task properties. Select the integration resource group required at runtime.
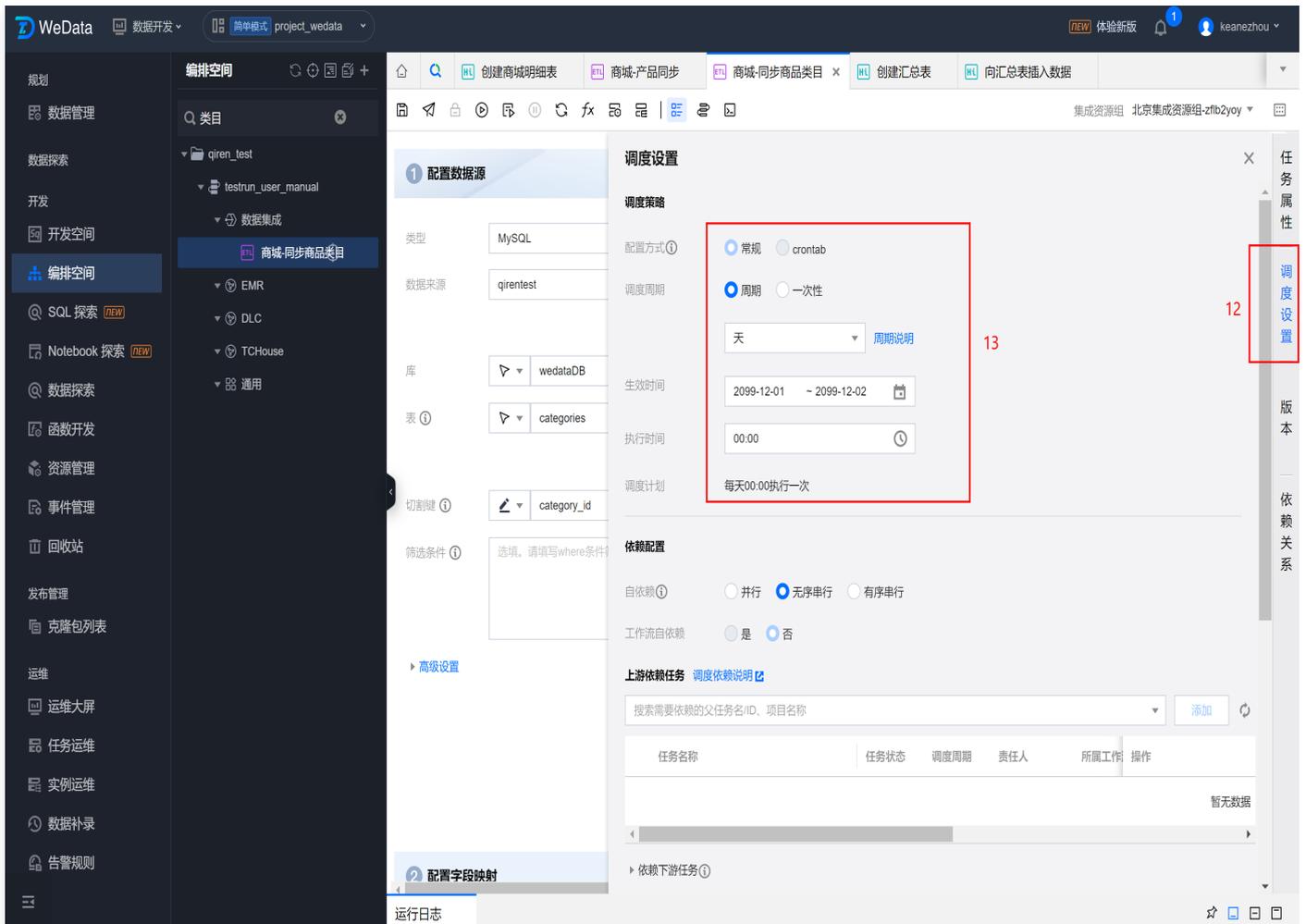
> ⚠ **Notes:**
> Here, it is essential to ensure network connectivity among the integration resource group, MySQL instance, and EMR cluster. Tencent Cloud resources in the same region must be purchased.



5. Set task scheduling. Here, we need to set the running strategy of the task. Since the frequency of changes in the category table is small, we set a daily sync once in the wee hours.

- Scheduling method: Periodic scheduling.
- Effective date: Default

- Scheduling cycle: Day
- Execution time: 00:00



6. Set task attributes:

7. after the above steps are completed, please save data promptly.

8. You can perform a dry run once before the official submission. At this point, the system will also check the integrity of the configuration and network connectivity. Once the detection passes, the system will start running immediately. The running log will appear at the bottom of the page.
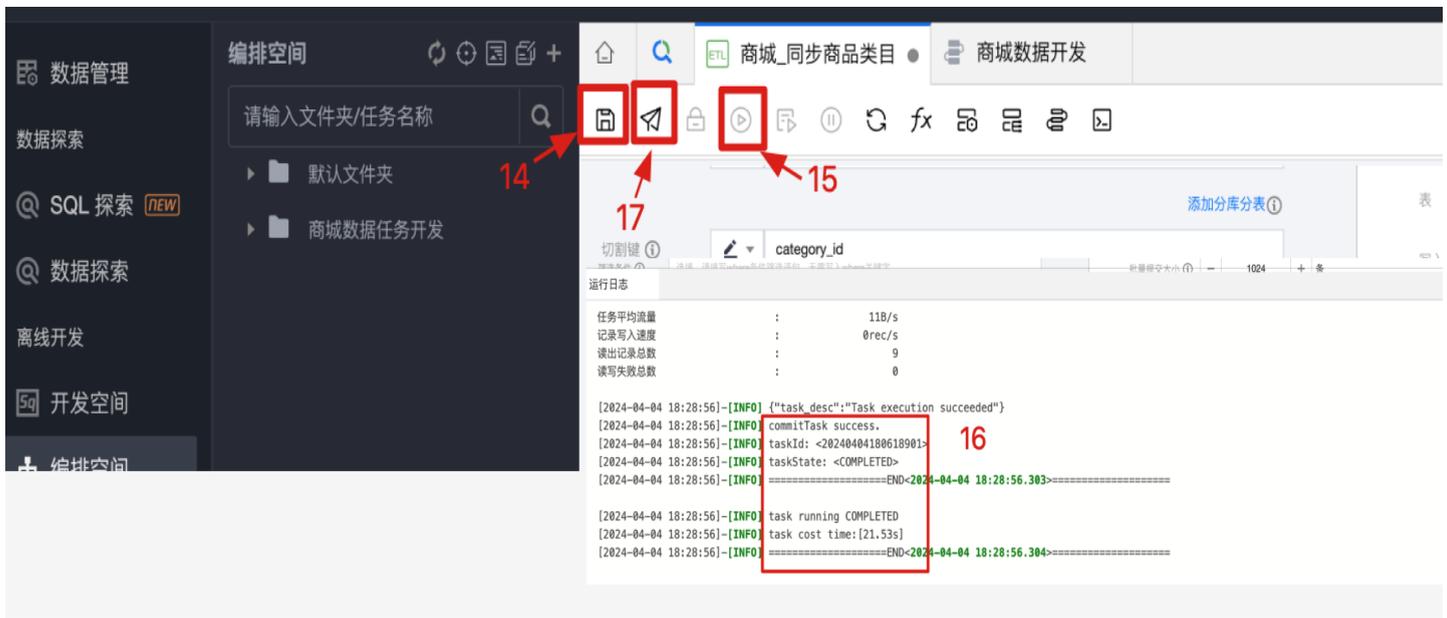
9. The logs and progress of the trial run can be seen at the bottom of the page. When Success or Completed appears, it indicates that the trial run is successful.

10. Submit the task to the scheduling resource server. Once the set running time is reached, the task will automatically start running. At this point, the configuration of the data integration task is completed.

> ⚠ **Notes:**
> Upon submission, the system will automatically detect the integrity and network connectivity of the configuration.

> If a notification about network connectivity issues appears, do not ignore it. Check immediately whether the networks of the integration resource group, MySQL instance, and EMR cluster are reachable.



By completing steps 1 – 17, you have completed synchronizing the category table from the MySQL database to the EMR cluster. And every midnight, WeData will automatically pull full data for overwrite update.

## Synchronizing City Table

We will now create the second offline synchronization task to synchronize the City Table from MySQL to the Hive table.

Since you have completed a sync task of a data table, you should have a certain understanding of data synchronization. Here we add a step: develop a synchronization strategy.

The synchronization strategy generally includes:

- **Offline sync or real-time synchronization?**

  When the business has small requirements for access latency and real-timeness, we generally select offline sync.

- **Additional Notes**

Difference between real-time synchronization and offline sync:

| Method | Description |
| --- | --- |
| Real-time Sync | **Definition:** Real-time synchronization refers to data being transmitted to the target system almost immediately after a change occurs in the source system. |
| | **Use cases:** Applications requiring high real-timeness of data, such as financial transactions, online collaboration tools. |

| | **Strengths:** |
|---|---|
| | • Real-timeness: Data changes can be immediately reflected in the target system, reducing the time window of data inconsistency. |
| | • Data consistency: Due to the fast synchronization speed, problems caused by data inconsistency can be reduced. |
| | **Drawbacks:** |
| | • High system resource consumption: A continuous network connection and relatively high system resources are required to maintain real-timeness. |
| | • High cost: Real-time synchronization may need more complex technical support and higher O&M costs. |
| | • Complexity: Implementing real-time synchronization requires more complex logic to handle data conflicts and synchronization status. |
| Offline sync | **Definition:** Offline sync refers to data not being transmitted immediately after a change occurs in the source system, but rather batch-transmitted at a specific time point or under certain conditions. |
| | **Use cases:** Environments where data has low real-time requirement or network conditions are unstable, such as data backup on mobile devices, regular synchronization of specific enterprise data, etc. |
| | **Strengths:** |
| | • Low system resource consumption: Synchronization can be arranged according to network and system resource conditions to reduce the need for real-time resources. |
| | • Low cost: Compared with real-time synchronization, offline sync has lower O&M costs and technical support requirements. |
| | • Flexibility: The time and frequency of synchronization can be arranged according to your actual needs. |
| | **Drawbacks:** |
| | • Delay: There may be a delay in data sync, and the latest data changes may not be reflected immediately. |
| | • Data consistency risk: If the sync interval is relatively long, it may result in data inconsistency. |

**Summary:**

The choice of sync method depends on business requirements, the importance of data, network conditions, and cost budget. Real-time synchronization is suitable for scenarios with extremely high requirements for data real-timeness, while offline sync is applicable to environments that can tolerate a certain amount of data delay. In some cases, these two strategies can also be used in conjunction with each other. For example, use offline sync when network conditions are

poor and real-time synchronization when network conditions are good to balance real-timeness and cost.

- **Should incremental or full synchronization policy be selected?**

In real business scenarios, incremental synchronization is generally selected. Full synchronization is only selected during data table initialization.
In this tutorial, slice fields such as date are not set in city, category, and Product Table. We depend on selecting full synchronization.
In following steps, when synchronizing the order table, we will introduce how to set up incremental synchronization.

- **If it is offline sync, then how often do you select for synchronization, per day or per hour?**

Synchronization frequency needs to be determined according to business needs. The smaller the frequency, the larger the resource consumption.
In this tutorial, we all select to sync once every day at wee hours.
The synchronization strategy of the City Table is identical to that of the category table. Please refer to steps 1 – 17 in the category table and repeat the operation.

> ⚠ **Notes:**
> - **The following steps are all marked with the step numbers in the diagram.**
> - Step 3: Task name: Mall_Synchronize city information
> - Step 4: Selection table: cities;
> - Step 6, 7: The table name needs to be modified to: ods_order_city;
> - Steps 10 – 13: These four procedures are easy to ignore;
> - Step 15: Regardless of how familiar you are with the operation, remember to run it once before submission to ensure the task runs accurately.

## Synchronize Product Table

Create the third offline synchronization task below, synchronizing the Product Table from MySQL to the Hive table.
Still think before reoperating.

| No. | Issue | Conclusion |
|---|---|---|
| 1 | Offline sync or real-time synchronization? | Offline |
| 2 | Should incremental or full synchronization policy be selected? | Full<br><br>ⓘ **Note**<br>Actually, it should be incremental. Here, we select full for teaching purposes. |

| 3 | If it is offline sync, then how often do you select for synchronization, per day or per hour? | Sync at midnight every day |
| --- | --- | --- |

Based on the above thinking, the synchronization strategies of the Product Table, category table and City Table are also identical. Please continue to refer to steps 1 – 17 in the category table and repeat the operation.

> ⚠ **Notes:**
> - **The following steps are all marked with the step numbers in the diagram.**
> - Step 3: Task name: Mall_Synchronize product information.
> - Step 4: Selection table: products.
> - Step 6, 7: The table name needs to be modified to: ods_product_product.
> - Steps 10 – 13: These four procedures are easy to ignore.
> - Step 15: Regardless of how familiar you are with the operation, remember to run it once before submission to ensure the task runs accurately.

## Synchronize Order Table

Create the fourth offline synchronization task below, synchronizing the order table from MySQL to the Hive table.
Still think before reoperating.

| No. | Issue | Conclusion |
| --- | --- | --- |
| 1 | Offline sync or real–time synchronization? | Offline |
| 2 | Should incremental or full synchronization policy be selected? | Incremental<br><br>> ⓘ **Note**<br>> Since the order table often has a large data volume, it is more suitable for incremental synchronization, that is, synchronizing the order data of the previous day every day. |
| 3 | If it is offline sync, then how often do you select for synchronization, per day or per hour? | Sync at midnight every day |

Through the above thinking, the synchronization strategy of the order table is different from that of the above three tables. Here we have selected incremental synchronization, which we will focus on introducing in this step.

## Incremental sync logic:

There is a special field in the original data order table: order_time, which will change as time goes by. Therefore, we can use the order creation time as a partition and ensure the pulling of real-time incremental data daily based on order_time.

To dynamically compare **task running time** (represented as ${yyyy-MM-dd}) with **order_time** during scheduling

For example: When the running time is 2024-04-01 00:10, then ${yyyy-MM-dd} = 2024-04-01, and meanwhile ${yyyy-MM-dd**-1d**} = 2024-03-31

Therefore, we can use date(order_time) = '${yyyy-MM-dd**-1d**}' to represent the data of yesterday.

## Create offline synchronization task

First, please continue to refer to steps 1 – 13 in the category table and repeat the operation.

> ⚠ **Notes:**
> - **The following serial numbers are all marked with the step numbers in the diagram.**
> - Step 4: Selection table: orders;
> - Step 6, 7: The table name needs to be modified to: ods_order_order;
> - Do not perform operation steps 14 – 17 (do not submit). We need to modify the following configuration.

## Operation Demonstration Screenshot:

1. Open the Configure Data Source page. Fill in the **filter conditions** on the left. Select **write mode** on the right as **overwrite.**

- Filter conditions: Fill in the corresponding filter statement according to the data type. This statement will serve as the filter conditions for the data to be synchronized.

- Write mode

   1.1 Append: Retain original data and append new rows.

   1.2 nonConflict: Error reported on data conflict.

   1.3 Overwrite: Delete the original data and rewrite it.

2. Open the **Create Table** page, modify the table name to ods_order_order, and add the partition field PARTITIONED BY (pt_date date) at the same time.



3. On the **Configure Field Mappings** page, click **Add Fields**. The field name is date(order_time) and the type is function. Click on the small circle on the right. Drag and drop with the mouse to the target field pt_date on the right to **establish mapping relationship**.

4. Click on the **Task Scheduling** button on the right task bar. Find **Execution Time** on the Task Scheduling page and modify the execution time to: 00:10.



Create an order table SQL:

```
--One-click creation of order table
CREATE TABLE IF NOT EXISTS `emall`.`ods_order_order`(
`order_id` int,
`product_id` int,
`quantity` int,
`unit_price` decimal(10,2),
`amount` decimal(10,2),
`order_time` timestamp,
```

```
`shipping_city_id` int,
`shipping_address` string)
PARTITIONED BY (pt_date date)
row format delimited
fields terminated by '\t'
STORED AS PARQUET;
```

So far, we have completed all offline synchronization tasks from raw data tables to Hive tables. And every midnight, WeData will automatically conduct full/incremental data synchronization.

# Summary

Now you have completed the study of the Data Integration part. Now summarize:

| No. | Step Name |
|-----|-----------|
| 1 | **Confirm the original data table and the data target table**<br>Base Table: Read: Data Source<br>Target Table: Write: Data Destination |
| 2 | **Confirm offline sync or real-time synchronization**<br>According to business needs, if not necessary, offline sync is available to reduce resource consumption. |
| 3 | **Confirm incremental synchronization or full synchronization**<br>Generally, full synchronization is required during data initialization, and incremental synchronization is used for periodic synchronization.<br>During incremental synchronization, it is required to set filter conditions to ensure no overlap when pulling data. |
| 4 | **Confirm network environment interoperability**<br>Three environments are involved in the sync process:<br>1. original database instance<br>2. Integration Resource Group<br>3. EMR cluster<br><br>⚠ **Notes:**<br>Must ensure that the integration resource group can access the original database instance and EMR cluster. |

Below we will learn part of offline development, that is, performing data processing in the Hive table of the EMR cluster.

# Offline development

Last updated：2025-04-18 15:49:07

Through the previous steps, we have synchronized all the raw data to the Hive tables in the EMR cluster.

- However, these data are all in their raw structures and cannot be directly used for business purposes.
- Combining the content from the **Data table structure design** step, we have analyzed the business requirements and divided the data warehouse hierarchy.

Next, we will complete the generation and data processing of the detail and summary tables through data development.

## Offline development task design



## Offline development task development

### Detail Table Development

Completing the detail table development mainly includes the following 4 steps:

### Create a shopping mall detail table

1. In the **Orchestration Space** module, under any folder in the shopping mall detail table development workflow directory, click the **Hive SQL icon** under Data Calculation, and create a new **Hive SQL** task. Name the task 'Create a shopping mall detail table', task type is Hive SQL, then click **Confirm**.

2. After confirmation, the Hive SQL development script page will pop up, follow these steps to complete the Hive SQL development task.

- Data Source: Select hive_emr-XXX.

- Scheduling Resource Group: Select the resource group we purchased, it is recommended to be in the same domain.

- Development Script: Write the table creation SQL script for the shopping mall detail table on the script page.

- Scheduling Settings: In the right taskbar, click **Scheduling Settings**, select the scheduling cycle as **One-time**, and default the execution event.

  ○ One-time: This Hive SQL executes only once according to the execution time.

  ○ Cycle: This Hive SQL executes periodically according to the execution time.

- **Click Save** button: Save this Hive SQL task.

- **Click Run** button: Execute the task once to verify the script task.

- **Click Submit** button: Officially submit this task to the scheduling resource server. The task will be executed according to the scheduling cycle at the specified time.

3. **Create Shopping Mall Detail Table HiveQL Statement:**

```sql
--Create Detail Table HiveQL Statement
CREATE TABLE
  IF NOT EXISTS emall.dwd_trade_order_ordercreate_productsales (
    order_id INT COMMENT 'Order ID, primary key',
    product_id INT COMMENT 'Product ID',
    product_name STRING COMMENT 'Product Name',
    category_id INT COMMENT 'Product Category ID',
    category_name STRING COMMENT 'Product Category Name',
    quantity INT COMMENT 'Quantity, positive integer',
    unit_price DECIMAL(10, 2) COMMENT 'Unit Price, with two decimal places',
    amount DECIMAL(10, 2) COMMENT 'Subtotal Amount, i.e., quantity times unit
price',
    order_time TIMESTAMP COMMENT 'Order Time, accurate to the minute',
    shipping_city_id INT COMMENT 'Shipping City ID, foreign key',
    shipping_city_name STRING COMMENT 'City Name',
    shipping_address STRING COMMENT 'Shipping Address, including province, city,
district, and detailed address'
  )
  COMMENT 'Product Sales Detail Table, records the sales details of each order'
  PARTITIONED BY (pt_date STRING)
```

```
    row format delimited fields terminated by '\t'
    STORED AS PARQUET;
```

By completing steps 7 – 15, you have created a Hive table in the Hive data source within the EMR cluster.

## Write data to the Detail Table

Next, we will start writing data into the detail table:

Please repeat steps 7 – 15, paying attention to the following:

> ⚠ **Note:**
> - **The following steps correspond to the step numbers indicated in the illustrations**
> - In the same workflow, create a new HiveQL node
> - Step 8: Name it: Insert data into the Detail Table;
> - Steps 9 and 10: These two steps are often overlooked;
> - Step 11: The HiveQL statement is as follows;
> - Step 12:
>
>   Change the schedule cycle to: **Cycle**
>   Set the execution time to: 01:00
>   Description: This task needs to run once daily
> - Step 14: No matter how familiar you are with the operation, remember to test run before and after submission to ensure the task runs correctly.
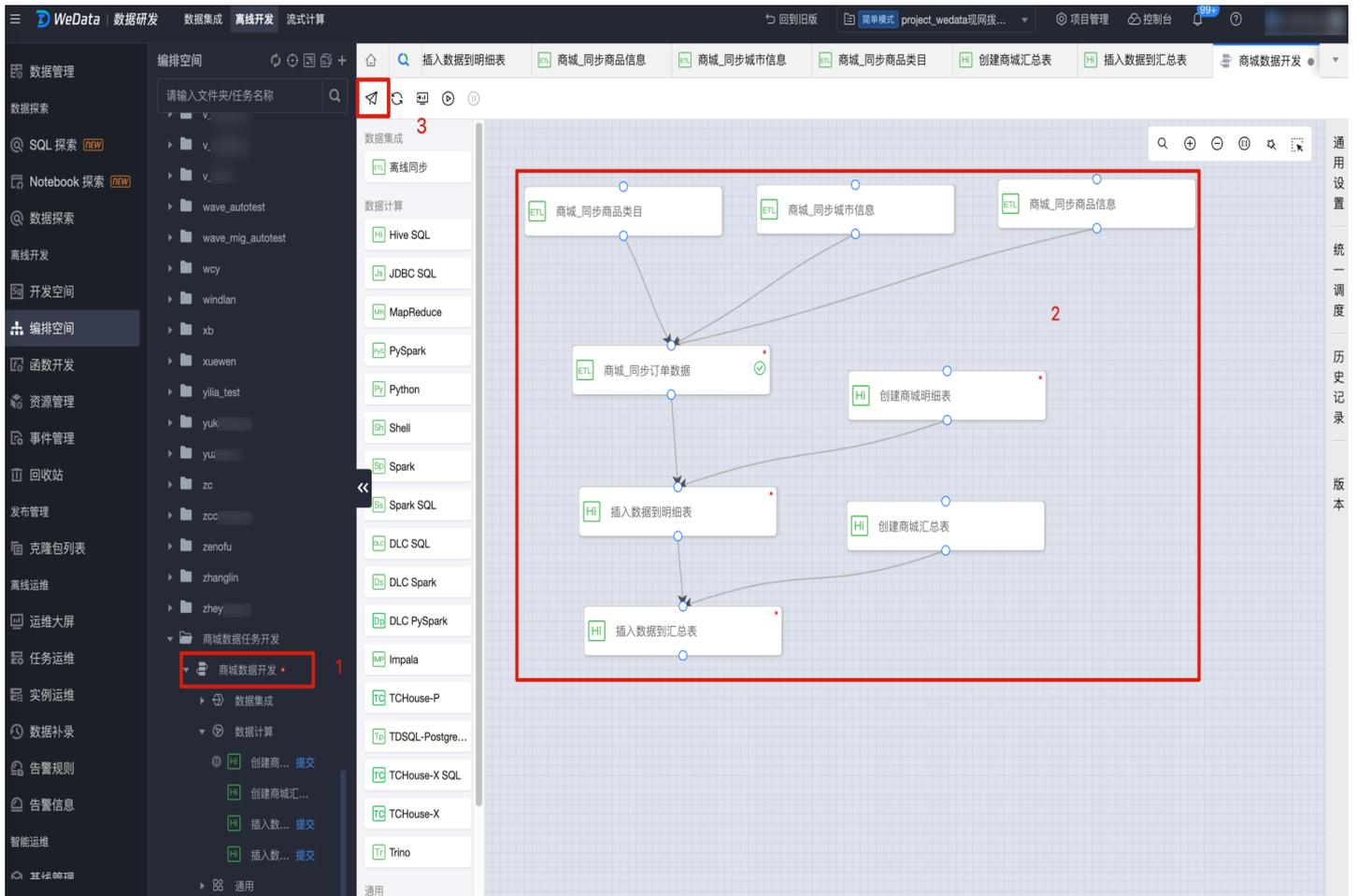
HiveQL statement to insert data into the Detail Table:

```
--Insert data into the Detail Table HiveQL statement
SET hive.exec.dynamic.partition.mode=nonstrict;
INSERT INTO TABLE emall.dwd_trade_order_ordercreate_productsales PARTITION
(pt_date)
SELECT
    o.order_id,
    o.product_id,
    p.product_name,
    p.category_id,
    ca.category_name,
    o.quantity,
    o.unit_price,
    o.amount,
    o.order_time,
    o.shipping_city_id,
    ci.city_name,
    o.shipping_address,
    o.pt_date
FROM
```

```
    emall.ods_order_order o
JOIN
    emall.ods_product_product p ON o.product_id = p.product_id
JOIN
    emall.ods_product_category ca ON p.category_id = ca.category_id
JOIN
    emall.ods_order_city ci ON o.shipping_city_id = ci.city_id
WHERE o.pt_date = '${yyyy-MM-dd-1d}';
```

At this point, we have completed the development task for the Detail Table. After the data is synchronized from the original table to the Hive cluster at midnight, the system will automatically link the four tables and compile the data into the Detail Table.

> ⓘ **Note**
>
> There are some redundant fields in the Detail Table at this point, which is intended to minimize joins when processing the Summary Table to improve computational efficiency.

## Summary Table Development

### Create the Mall Summary Table

Next, we will start developing the Summary Table
First, please repeat steps 5 – 6 to create a new workflow: Mall Summary Table Development.
Next, repeat steps 7 – 15 to create the summary table. Pay attention to the following points:

> ⚠ **Note:**
> - **The following steps correspond to the step numbers indicated in the illustrations**
> - Create a new HiveQL node in the same workflow;
> - Step 8: Name it: Create Mall Summary Table;
> - Steps 9 and 10: These two steps are often overlooked;
> - Step 11: The table creation SQL statement is as follows;
> - Step 12:
>
>   Set the scheduling cycle to:**Once**
>   Set the execution time to: Default
>   Description: This task only needs to run once
> - Step 14: No matter how familiar you are with the operation, remember to test run before and after submission to ensure the task runs correctly.

**Create the emall summary table HiveQL statement:**

```
--Create summary table HiveQL statement
CREATE TABLE IF NOT EXISTS emall.dws_trade_order_productsales_1d (
    order_date DATE COMMENT 'Date of statistics, primary key',
    city_id INT COMMENT 'City ID',
```

```
    city_name STRING COMMENT 'City name',
    category_id INT COMMENT 'Product category ID',
    category_name STRING COMMENT 'Product category name',
    quantity INT COMMENT 'Total product sales, positive integer',
    amount DECIMAL(10, 2) COMMENT 'Total product sales in monetary terms, two
decimal places kept'
)
COMMENT 'Daily summary table of product sales situations'
PARTITIONED BY (pt_date STRING)
row format delimited fields terminated by '\t'
STORED AS PARQUET;
```

## Write data into the summary table

Next, repeat steps 7 – 15 to insert data into the summary table. Pay attention to the following points:

> ⚠️ **Note:**
> - **The following steps correspond to the step numbers indicated in the illustrations**
> - Create a new HiveQL node in the same workflow;
> - Step 8: Name it: Insert Data into the Summary Table;
> - Steps 9 and 10: These two steps are often overlooked;
> - Step 11: The HiveQL statement is as follows;
> - Step 12:
>   Change the schedule cycle to: **Cycle**
>   Set the execution time to: 01:00
>   Description: This task needs to run once daily
> - Step 14: No matter how familiar you are with the operation, remember to test run before and after submission to ensure the task runs correctly.

**Insert data into the Summary Table HiveQL statement:**

```
--Insert data into the Summary Table HiveQL statement
SET hive.exec.dynamic.partition.mode=nonstrict;
INSERT INTO TABLE emall.dws_trade_order_productsales_1d PARTITION (pt_date)
SELECT
    p.pt_date AS order_date,
    p.shipping_city_id,
    p.shipping_city_name,
    p.category_id,
    p.category_name,
    SUM(p.quantity) AS quantity,
    SUM(p.amount) AS amount,
    p.pt_date
FROM
```

```
    emall.dwd_trade_order_ordercreate_productsales p
WHERE p.pt_date = '${yyyy-MM-dd-1d}'
GROUP BY
    p.pt_date,
    p.shipping_city_id,
    p.shipping_city_name,
    p.category_id,
    p.category_name;
```

We have completed the offline development tasks for the detail table and summary table. Every early morning, WeData will automatically perform the calculation tasks for the detail and summary tables.

## Establish dependencies and submit

1. **Double-click** the mall data development workflow. A workflow canvas will pop up, displaying all the tasks under the workflow. Establish the dependencies among the workflows in sequence. Once completed, click the submit button.

- Establish dependencies in sequence:
  - Mall_Synchronize_Product_Category → Mall_Synchronize_Order_Data.
  - Mall_Synchronize_City_Information → Mall_Synchronize_Order_Data.
  - Mall_Synchronize_Product_Information → Mall_Synchronize_Order_Data.
  - Mall_Synchronize_Order_Data → Insert Data into the Detail Table.
  - Create Mall Detail Table → Insert Data into the Detail Table.
  - Insert Data into the Detail Table → Insert Data into the Summary Table.
  - Create Mall Summary Table → Insert Data into the Summary Table.
- Finally, submit the workflow. The tasks will run in an orderly manner according to the set dependencies.

# Offline Development Task Ops

You can view the running status of workflows or offline tasks in Task Ops.

# Data Quality

Last updated：2025-04-18 15:51:17

In this step, we will complete quality monitoring of the data tables in the data warehouse to prevent dirty data from being transmitted downstream.

## Quality Monitoring Task Design

When the following fields in the detailed table are empty, it will cause a serious impact on the summary table:

- Monitoring table: dwd_trade_order_ordercreate_productsales.
- Monitoring fields: amount, order_date.
- Monitoring logic: Depend on the completion of the detailed task and then automatically detect whether there are null values.

## Quality Monitoring Task Development

Quality task development mainly includes the following 7 steps:

## Null Value Detection Task

### Step 1: Select Monitoring Fields

1. Click the **Data Quality** module, enter the **Data Monitoring** page, and click the **Add Rules to Multiple Tables** button again.
   - Add rules to multiple tables: Support setting monitoring rules for multiple tables or more fields at one time.



2. Click **Monitor Multiple Fields**, then click the **Add Fields** button to start adding the fields that need to be monitored.

3. On the add field page, select the add method as **manually add**, select **data source** as hive_emr-XXX, select **database** as hive database emall. After completion of the selection, tables and fields below will be updated. Select table dwd_trade_order_ordercreate_productsales and corresponding fields amount, order_date, and click **save**.



4. After successful saving, the page will refresh to display the selected tables and fields. Click the **Next** button to start configuring monitoring rules.

## Step 2: Configure Monitoring Rules

1. Configure monitoring rules for the previously selected tables and fields. Select **rule template** as system template, **select template** as number of empty fields, fill in **rule name** as value not empty, set **trigger condition** as equal to or greater than 1, set **trigger level** as high. After the settings are completed, click **next** button to start configuring execution policies and subscription information.

   ○ Rule template: WeData has already built-in more than 50 system templates, which can be used directly here.

   ○ Select a template: On the right, you can **view template descriptions**.

   ○ Trigger condition: means that when the count of empty values is equal to or greater than 1, immediately interrupt downstream tasks and send an alarm.

# Procedure 3: Set Execution Policy

1. Click **Rule Name,** batch-select all rules, and then click **Batch Set Execution Strategy** button to configure execution policies.



2. **Select Execution Mode** as Associated Production Scheduling, **Select Task** as Insert Data into Detail Table, and click **Save** button after the settings are completed.

   ○ Associate with production scheduling: means associating quality tasks with data development tasks. Only after the associated tasks are completed will this quality monitoring task be executed. Since we select **Insert data into detail table** here, data integrity will be detected immediately after data is inserted into the detail table.

   ○ The execution engine, computational resource, and execution resource are consistent with the selection in the above text.

   ○ Select task: the data development task that needs to be associated.

## Procedure 4: Set Subscription Notifications

1. Click **rule name,** batch-select all rules, and click **Batch Set Subscription Information** button to configure subscription information.

   ○ Set up subscription notifications: Decide what method will be used to send a message reminder when an exception is detected.



2. **Subscription Configuration**: Select **email** and **Short Message Service,** and select the recipient as XXX.

## Procedure 5: Set Detection Scope

1. Click **Rule Name**, batch-select all rules, and then click **Batch Set Detection Scope** button to configure the detection scope.

   ○ Set detection scope: Set which data scope to detect.



2. Set **Detection Scope** as conditional scanning, fill in content: order_date = ${yyyy-MM-dd-1}, and click **Save** button.

   ○ Conditional scanning: Only check the newly produced data of each day according to the filled-in conditions, rather than performing a full check every day. Because the larger the amount of monitoring data, the larger the resource consumption.

   ○ Related reference information can be viewed on the right.

3. The above indicates that the configuration is all completed. Click **Complete** button.



## Procedure 6: Task Trial Run

1. Click **Trial Run** button on the right side of the page to configure, select the scheduling time as the trial run time, and click **Start Trial Run** button again.

   - execution engine, computational resource and execution resource should be consistent with the above context. **Note network connectivity.**

   - Trial run: Before releasing the task, you can run it once to detect whether there are errors in the task.

   - Start trial run: Just wait for pending execution and monitoring.

     - During the trial run, it should be ensured that there is data in the detailed table. If the trial run has been performed during both data integration and data development processes, there should be data in the detailed table.

## Procedure 7: Task Release

1. Return to the Data Monitoring page, select the **View Rules by Table** page, and click the **Monitoring Status** button to enable monitoring.
   - View Rules by Table: Filters out rules based on data source, database, and data table.



# Quality Monitoring Task Operation and Maintenance

You can view the task running result of the quality monitoring task in Ops management.

# Data services

Last updated：2024-09-05 16:45:56

> ⓘ **Note**
>
> The data service module is a **non-essential step** in the data development process. If your business does not involve this aspect, you **can skip** this module.
> If you need to experience this module's tutorial, please contact the enterprise administrator to complete the environment preparation.

According to the initial scheme design of the tutorial, the business side expects to understand the sales performance of different categories in different regions every day.
In Chapter Five, we have created a Mall summary table, which records product sales quantity and total sales amount by day, city, and category.
Regarding the data application of the Shopping Mall Detail Table, there are multiple ways:

1. Using SQL to perform direct data extraction and view daily sales performance is quite simple and will not be covered in this tutorial.

2. Displaying data report views using BI software (e.g., BI). For relevant documentation, please refer to the BI official website. This tutorial will not cover this method.

3. Using WeData data service features to produce API services for external systems to call. This part will be the focus of this tutorial.

Next, we will introduce how to use WeData data service features to produce API services for external systems to call.
Before developing the data service interface, we need to understand some technical information.

> ⓘ **Note**
>
> - **Hive**'s feature is big data processing. It is mainly used for handling large-scale structured data, especially in data warehouse and big data analysis scenarios. However, its query speed is relatively slow, making it unsuitable for APIs that require rapid response. Therefore, it is generally not used as backend data storage for APIs. The summarized data computed from big data must be imported into an engine with better query performance before providing API queries.
> - **MySQL** can be used for applications that require rapid response and transaction support. It has fast query speed and supports complex queries and transaction processing.
>
> Therefore, we use MySQL as the underlying data storage for the API.

Next, we need to transfer the computed data to the MySQL database.

## Environment preparations

**Role:** Enterprise Administrator.

> ⓘ **Note**

> The following operations (environment and resource preparation) involve resource purchase and paid content, which need to be performed by the Enterprise Administrator.

# Purchase Data Service Resources

1. Go to the WeData Console and enter the Execution Resource Group page, select the region as **Beijing**, click **Data Service Resource Group** to enter the page, then click the **Create** button.
   - Region: In this tutorial, select resources in the **Beijing** region.



2. After clicking the **Create** button, you will enter the **Data Service Resource Group Purchase** page. In **Resource Configuration**, select the region as **Beijing**, select the **Network** as the VPC created in step Create a New VPC, and select the specifications as **Test Specifications**.
   - Specifications: For detailed information about service resources, please see Service Resource Billing Explanation.

3. In **Associated Project Space**, select **Associate Project** as Immediate Association, associate with the project created in step Create a Project in WeData , and click the **Purchase Now** button.



# Purchase API Gateway Resources

1. Go to the Tencent Cloud Official Website > Cloud Native Gateway . If you are using it for the first time, you will need to enable permissions by clicking the **Use Now** button.

2. Log in to the Tencent Cloud Console > **Instance List**, go to the Gateway **Instance** list, select the region as **Beijing,** and click the **Create** button.



3. In the configuration selection, select the region as **Beijing,** select the instance specification as **Basic Edition,** and select the network as the VPC created in step **Create a New VPC**.

- Availability Zone: The availability zone where the subnet in the VPC is located.

## Create Summary Table

Create a data summary table in Tencent Cloud MySQL. For specific steps, see Summary Table Development .
The table creation SQL statement is as follows:

```
-- In MySQL, create a summary table
CREATE TABLE emall.dws_trade_order_productsales_1d (
    order_date DATE NOT NULL,
    city_id INT NOT NULL,
    category_id INT NOT NULL,
    city_name VARCHAR(50) NOT NULL,
    category_name VARCHAR(50) NOT NULL,
    quantity INT NOT NULL,
    amount DECIMAL(10, 2) NOT NULL,
    pt_date VARCHAR(50) NOT NULL
);
```

## Data Service Interface Design

**Interface Input Parameters:**

| Field Name | Associated Fields | Mandatory or Not | Format |
|---|---|---|---|
| Date | order_date | Mandatory | Date |
| Region | city_id | Optional | City Code, INT |
| Category | category_id | Optional | Category Code, INT |

**Interface Output:**

Here we output all fields of the data table.

# Data Service Interface Development

## Sync data to MySQL

First, we need to sync the data from the Hive mall summary table to the MySQL summary table.

### Step 1: Create a data synchronization task

> ⓘ **Note**
>
> In the **WeData > Data Development > Mall Development** workflow, add an offline synchronization
> node. For detailed steps, refer to Offline Synchronization Task Development .

1. In the orchestration space, find the created **Mall Data Task Development** > **Mall Data Development** >
   click **Offline Synchronization**, select configuration mode, click **Confirm.**

- Task Name: Sync summary data to MySQL.
- Development Mode: Choose **Form Mode.**

2. Configure **Data Source** and **Data Destination** in order.

> ⚠ **Note:**
> - Data Flow Direction: Hive → MySQL.
> - It is not possible to create MySQL tables in one click in WeData, so you need to create the data table in MySQL first.
> - Hive table filter condition: order_date = '${yyyy−MM−dd−1d}'.
> - MySQL write mode: overwrite.
> - Scheduling Settings:
>   - Scheduling Period: Choose a period.
>   - Set the execution time to: 01:00.

## Step 2: Suggest dependencies and submit

Establish an association with the upstream task for the data synchronization task just created. For specific operations, please refer to Establish Dependencies and Submit .
The suggested task dependency relationship after association is as shown below:



Upon completing the above steps, the data will automatically sync from the Hive table to MySQL.

## Create Interface Service

## Step 1: Create a Service

1. Enter the **Data Service > Service Development** page, click the **Create New Folder** button. Configure the folder name as **Shopping Mall API**. Select the target folder as the root directory. Click the **OK** button.



2. Then click **+ sign**, choose **Create New API** and click to enter the Create New API page. Select the configuration method as **Wizard Mode**, then configure **Authentication Method, Associated Gateway** and **Service Resource Group** sequentially. Click the **OK** button.

- Authentication Method: No Authentication.
- Associated Gateway: Select the API Gateway purchased in the step Purchasing API Gateway Resource.
- Service Resource Group: Select the data service purchased in the step Purchase Data Service Resources.

## Step 2: Set Data Source

1. Select **MySQL** data source, and choose the **Summary Table** created in the previous step.



## Step 3: Set Input and Output Parameters

1. Configure sequentially **Request Parameters, Response Parameters.**

- Request parameters:
  - Set three parameters: date, city_id, category_id.
  - Bind the database field.
  - Adjust the parameter types.
  - Keep date as the required field, others as not required.
- Response parameters:
  - Set all fields in the table as output fields.
  - Quick operation: After selecting the bound field, the parameter name will be automatically brought out.
  - Adjust the parameter types.



2. Configure **Sorting Parameters.**

- Sorting parameters: Set city_id, category_id to ascending sort.

| 排序参数 | 序号 | 10 | 字段名称 | 排序方式 | 操作 |
|---|---|---|---|---|---|
| | 1 | | city_id | 正序 | 编辑 删除 |
| | 2 | | category_id | 正序 | 编辑 删除 |

+ 添加参数  + 批量添加  - 清空

**高级配置**

| 超时时间 ⓘ | 5 | s |
|---|---|---|
| 单次最大返回<br>条数限制 ⓘ | 2000 | 条 |
| 最大每秒请求<br>数(QPS) ⓘ | 200 | 次/秒 |

## Step 4: API Testing

1. Click the **Test** button, in the pop-up test window, fill in the **mandatory parameter** date, e.g., 2024–04–01, then click the **Initiate a call** button, wait for the response, and check the output result.

## Step 5: Submit

1. After passing the test, click the **Submit** button. After submission, the interface is available.



# View the Data Service Interface

After the Data Service API is submitted, you can view this service in the service list. In this list, you can manage the API, such as setting visibility, setting alarms, decommissioning, and other operations.

# Business Research and Model Design (Optional)

Last updated: 2024-08-24 17:36:18
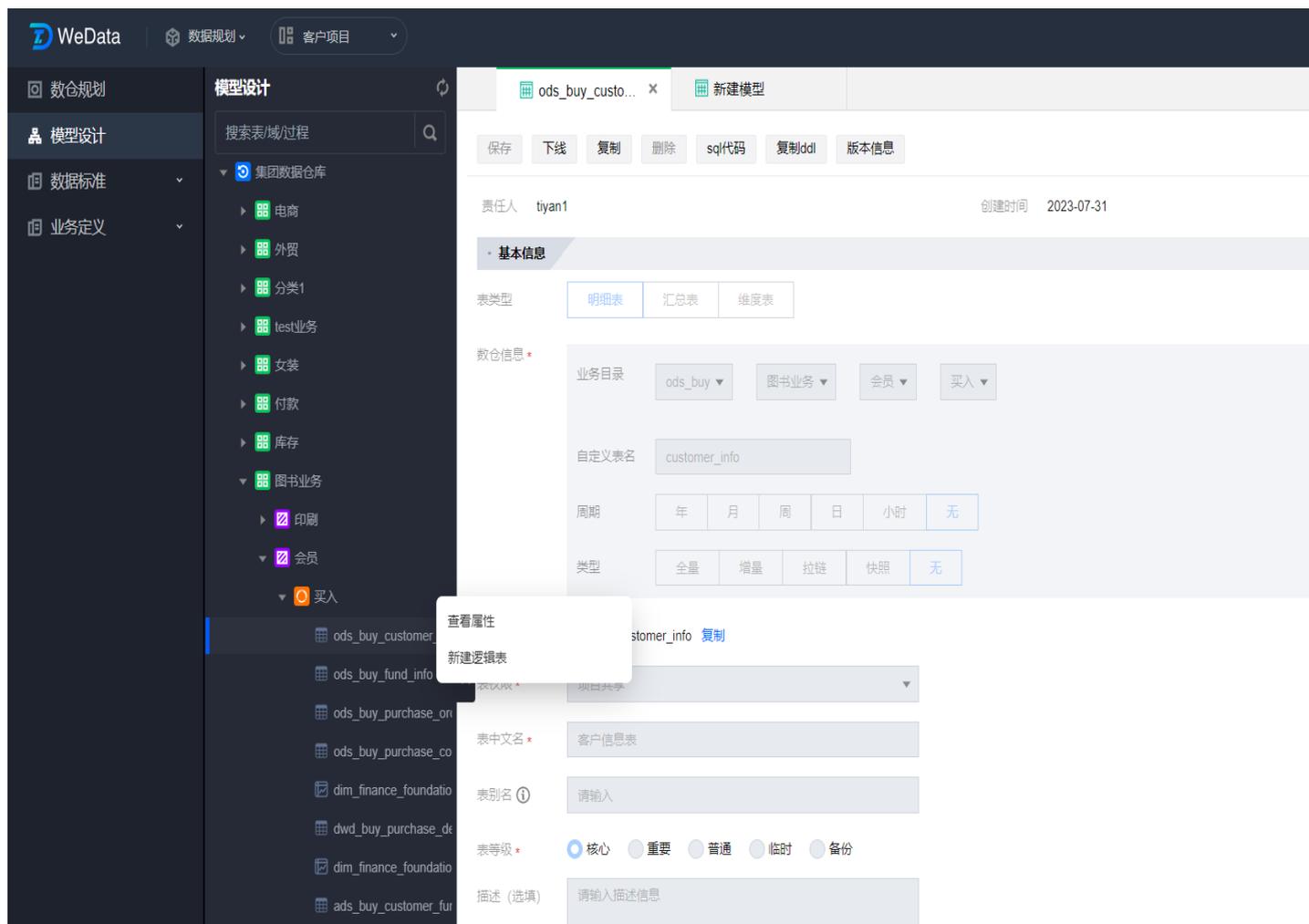
## Building a Data Warehouse

Business research will involve building the data warehouse framework based on business and data dimensions, using the concept of layering, categorization, and domain-based abstraction Definition:

- The Data Hierarchy Definition includes ODS, DWD, DWS, ADS, and DIM layers, which are mapped through logical layering and physical database associations.
- Business Type Definition includes business classification, subject domain, and business process, managing data objects with custom business catalog management.
- After Definition, the data warehouse architecture will be automatically generated. The subsequent process of defining model metrics dimensions will rely on the overall data warehouse architecture for Definition management.



## Model Design

After building the data warehouse, the designer will define the logical model and physical model based on data features and business scenarios. During the logical table definition, standardized naming is performed according to the data warehouse architecture. Additionally, metadata and value range standards can be bound during field configuration to complete the standardized definition process

The model design process will consider both data-driven (Bottom-up) and business-driven (Top-down) approaches:

1. For the data-driven dimension, the designer first needs to synchronize raw data from the production source system to the interface layer. After cleansing and transforming, the raw data in the interface layer will form the detail table, also known as the fact table, which stores the finest granularity data. The detail table is the source for metric statistics, and its fields will be bound with basic metrics and dimension conditions.

2. From the business-driven perspective, based on business scenarios, designers need to define an aggregation layer and market layer. The summary table in the aggregation layer will store aggregated metric data under different dimension conditions and will be used as the target table for derived metrics, forming a one-to-one binding.

3. The defined analysis dimensions will create dimension tables that store attribute hierarchy data for dimensions. These tables will be one-to-one bound with common dimensions and can be auto-generated during dimension definition.

After completing the logical model design, the physical model can be generated through the publishing action, linking the design and development processes. If some physical models have already been created, they can be reverse-imported to generate logical models, completing the design phase.
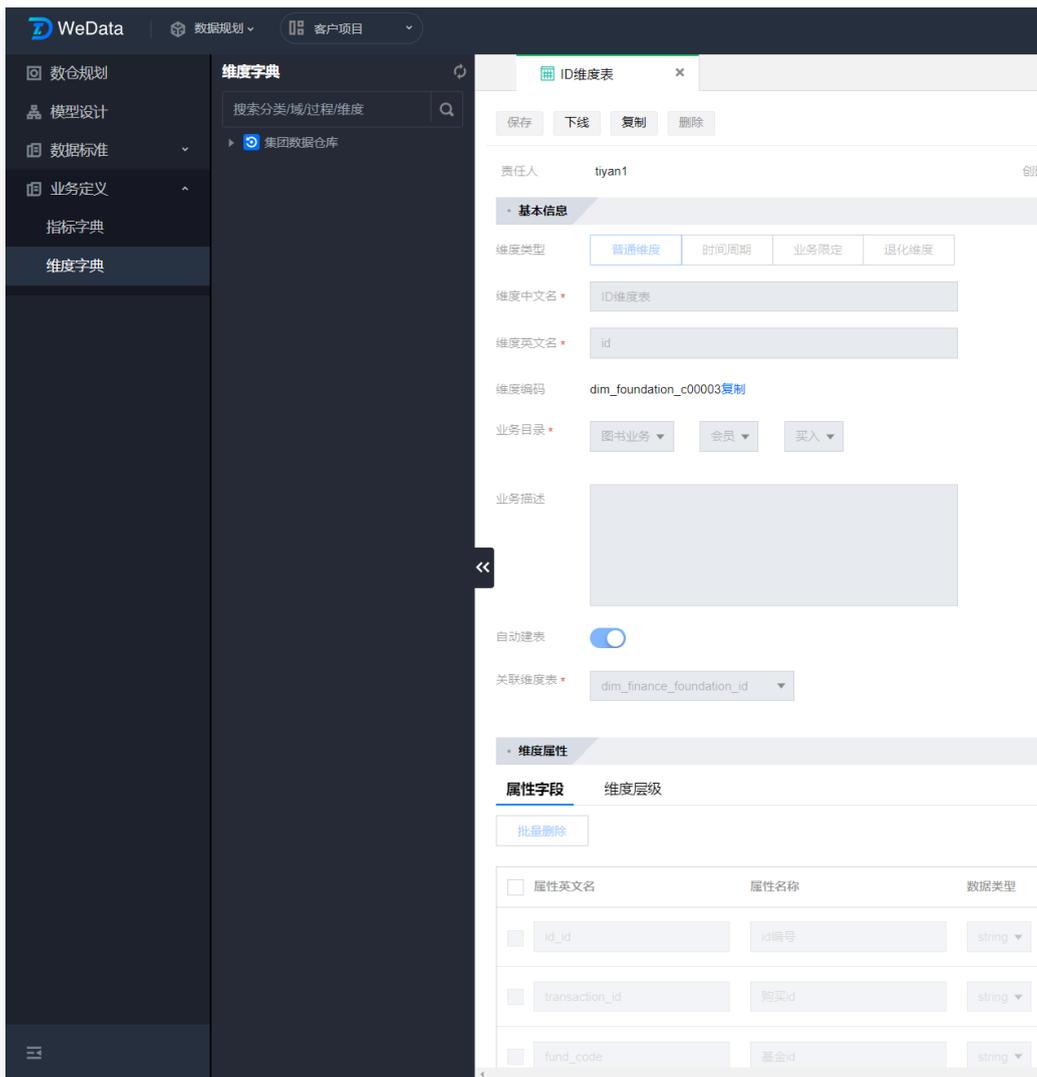
# Business Definition

During business research, the designer needs to abstractly define indicators and dimensions based on business scenarios:

1. Indicator Definition



Indicators are divided into two categories: basic metrics and derived metrics

1.1 Basic metrics are measurements that do not include dimensional conditions. They require definitions of their basic attributes, statistical calibers, and unit precisions. Derived metrics will inherit the unit of the basic metrics. The data for basic metrics comes from a specific field in the detail table, so they need to be associated in the indicator definition.

1.2 Derived metrics can be defined by adding dimensional conditions to basic metrics for specific characteristic ranges, such as user growth in a certain channel or product type. They can also be the result of combined calculations of multiple derived metrics under the same dimensional conditions, such as growth rate. Once defined, derived metrics are bound to a field in a summary table, facilitating indicator production.

2. Dimension Definition

Dimensions can be classified into the following categories:

2.1 Common dimension: This can be understood as the group by condition in SQL. A common dimension uniquely corresponds to a dimension table, associated during dimension modeling

2.2 Business Constraint: Also known as a modifier, it is used to filter tag characteristics from business dimensions

2.3 Time period: Time-based limiting conditions

2.4 Degenerate Dimension: Dimensions reverted to the fact table. This usually happens when a dimension has no other content besides the primary key, even though it's a legitimate dimension key. Reverting it to the fact table reduces the number of associations and improves query performance.
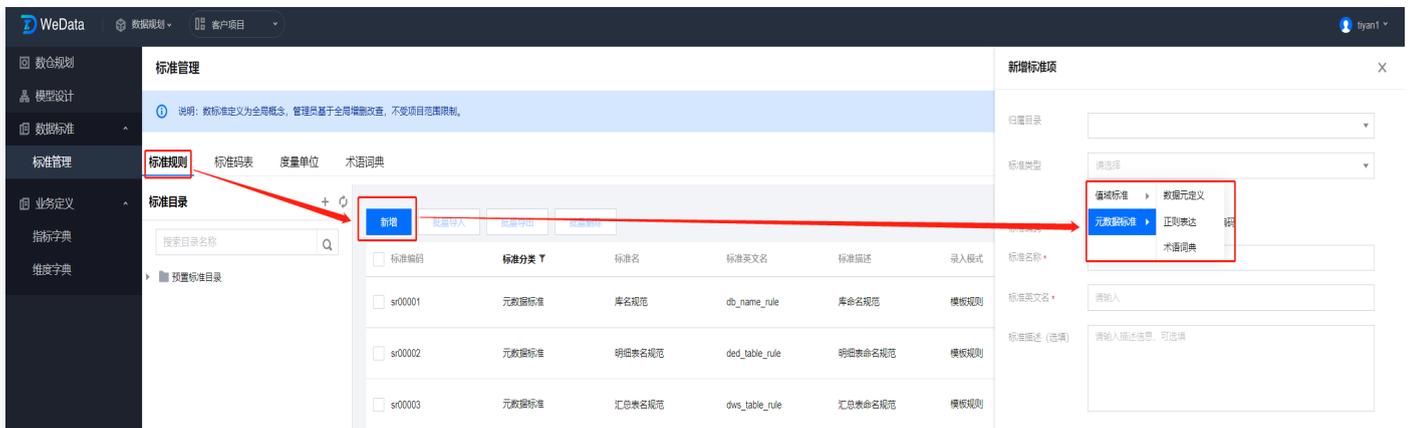
Indicators and dimensions need to be defined and published sequentially to establish associations with the table model and be referenced by subsequent derived definitions, guiding the implementation of indicator production.

# Data Standard

Model Metrics Definition: During the development and production process, operations must adhere to a unified data standard. Therefore, business objects need rules to be defined. Standards management
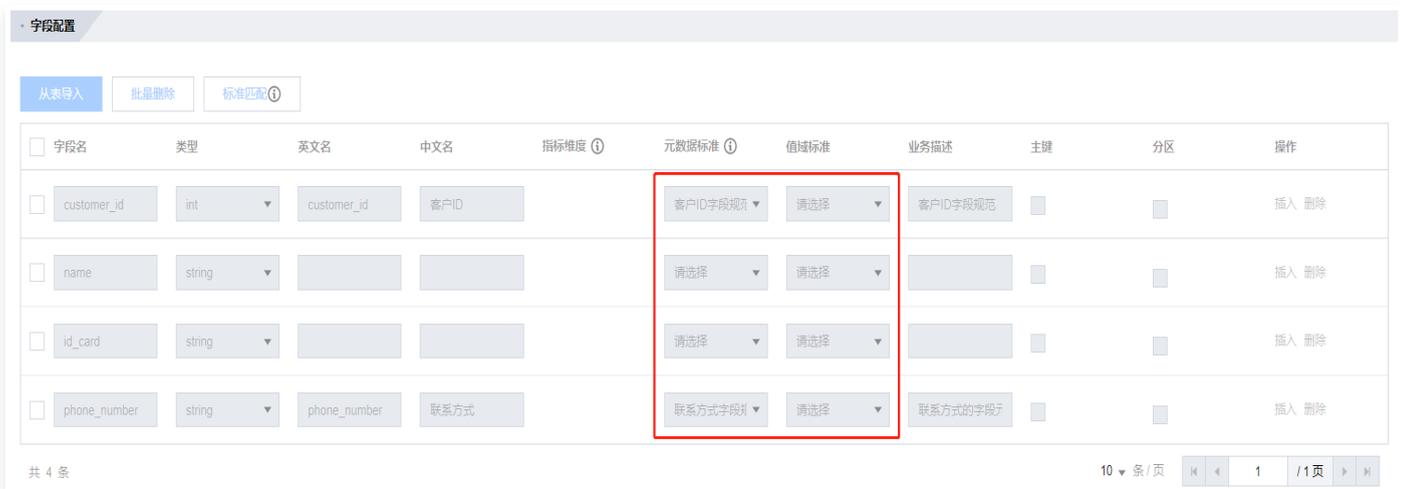
involves defining standards across the following four modules:
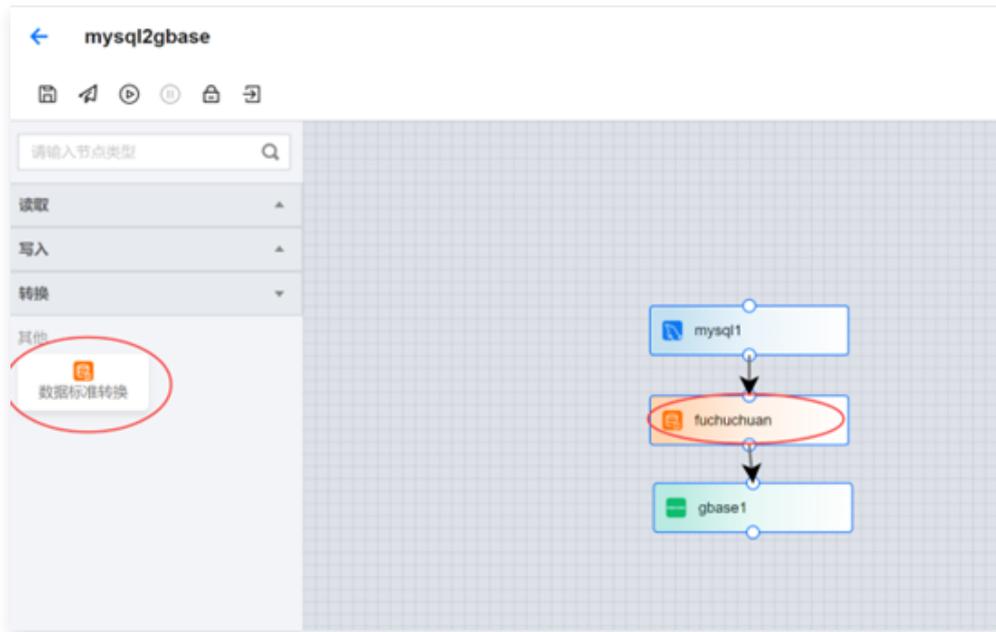
1. Definition Standard Rules



You can define standards at the table, field, and indicator levels. Metadata standards define the naming and type specifications of business objects, and value range standards define the characteristics of value ranges.

After rule release, you can perform association binding in the model design/physical model fields:



You can also use ETL tasks for standard conversion tasks during subsequent data development processes:
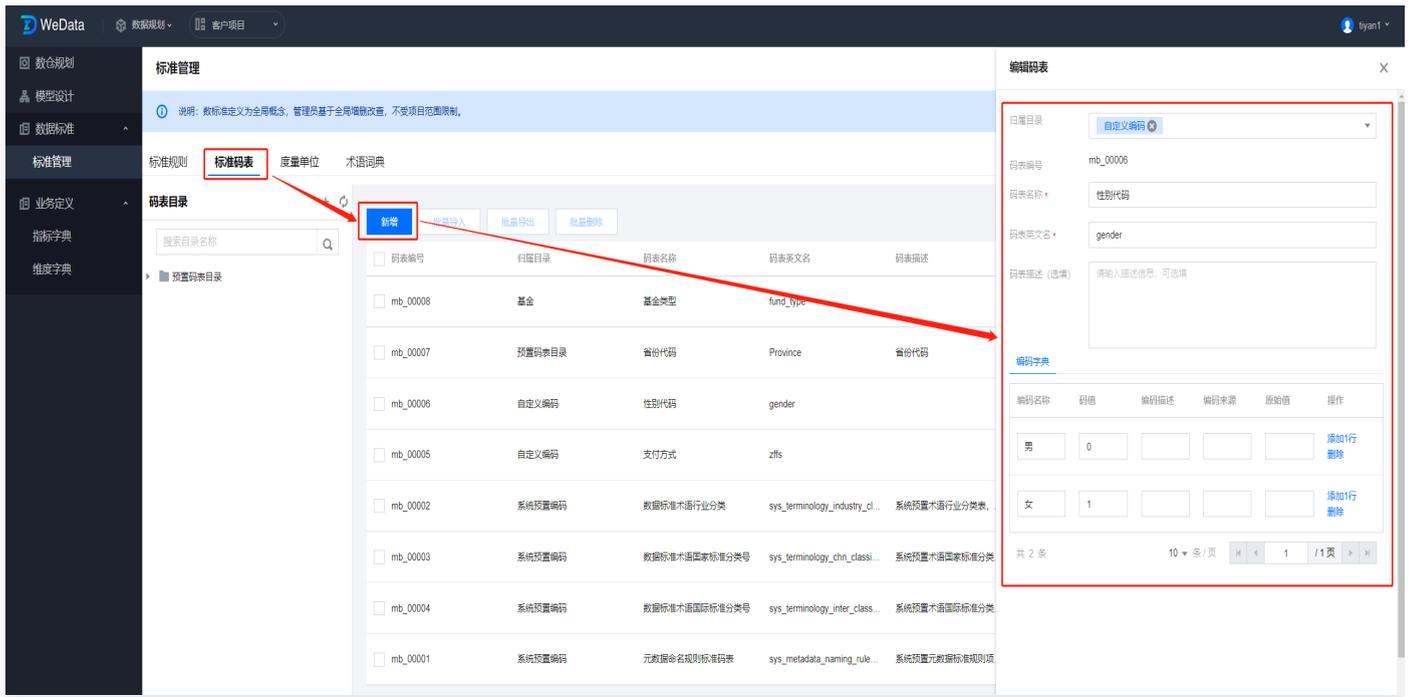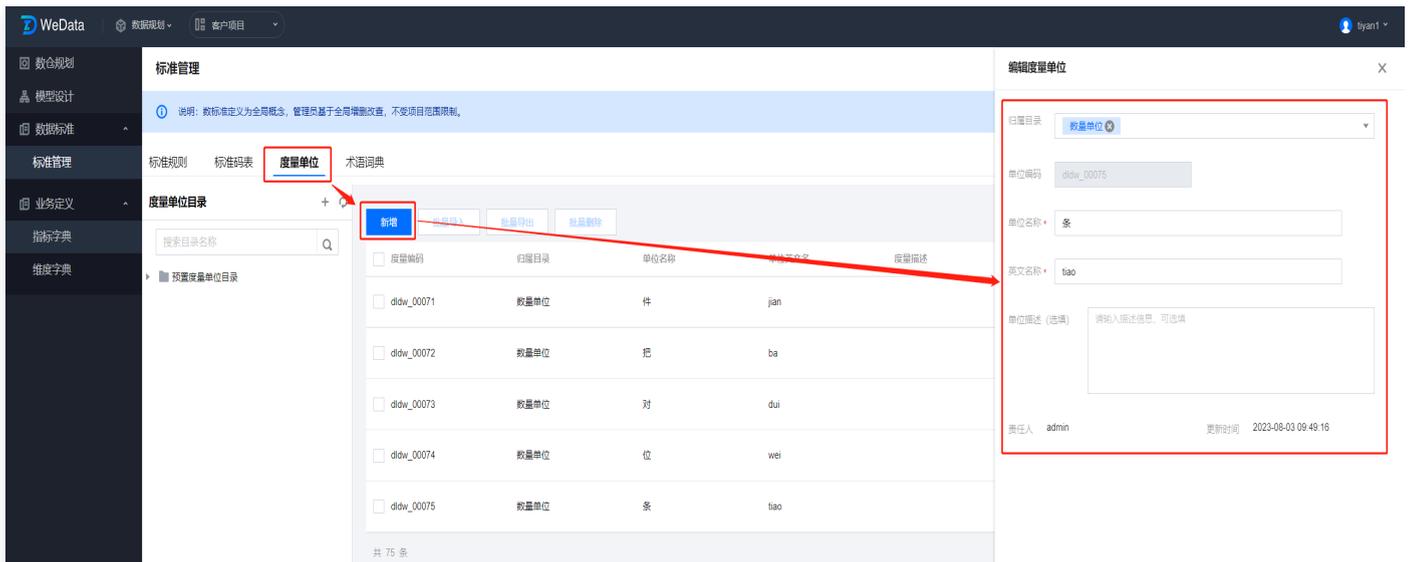
## Configure Conversion Rules:



2. Definition Standard Encoding

For data enumeration types, management requires standard encoding. After encoding is defined and released, it can be referenced in standard rules.

## 3. Definition Measurement Unit

When defining indicators, measurement units will be used. The system presets common units. For custom units, you can define them in this module.



## 4. Definition Terminology Dictionary

Industry standard metadata will be defined in bulk in the terminology dictionary. After definition release, it can be referenced in standard rules.