DataWorks

Quick Start

MORE THAN JUST CLOUD | C-) Alibaba Cloud

Quick Start

This guide describes how to quickly perform data development and O&M operations.

NOTE:

If this is the first time you are using DataWorks, make sure that you have prepared an account and configured the project roles and project according to steps in **Preparation**. Then, go to the **DataWorks console** page and click **Enter Workspace** after a project to go to the **Data Development** page of DataWorks to start data development.

Generally, DataWorks project space data development and O&M involve the following operations:

- Step 1: Create a table and upload data
- Step 2: Create a flow
- Step 3: Create a synchronization task
- Step 4: Set the period and dependency
- Step 5: Perform O&M and log troubleshooting

A general process is shown in the following figure:



This section uses the creation of the tables bank_data and result_table as an example to describe how to create a table and upload data. The table of bank_data is used to store business data, while the

result_table is used to store results after data analysis.

Instructions

Create bank_data

Log on to the project, and click Data Development > New > Create Table.



Enter the table creation statements, and click **OK**.For more information on table creation SQL syntax, see MaxCompute-based table creation, view, and deletion.

The statements used for table creation in this example are as follows:

```
CREATE TABLE IF NOT EXISTS bank_data
(
age BIGINT COMMENT 'age',
job STRING COMMENT 'job type',
marital STRING COMMENT 'marital status',
education STRING COMMENT 'educational level',
default STRING COMMENT 'credit card ownership',
housing STRING COMMENT 'mortgage',
loan STRING COMMENT 'loan',
contact STRING COMMENT 'contact information',
month STRING COMMENT 'month',
day_of_week STRING COMMENT 'day of the week',
duration STRING COMMENT 'Duration',
campaign BIGINT COMMENT 'contact times during the campaign',
pdays DOUBLE COMMENT 'time interval from the last contact',
previous DOUBLE COMMENT 'previous contact times with the customer',
poutcome STRING COMMENT 'marketing result',
emp_var_rate DOUBLE COMMENT 'employment change rate',
cons_price_idx DOUBLE COMMENT 'consumer price index',
cons conf idx DOUBLE COMMENT 'consumer confidence index',
euribor3m DOUBLE COMMENT 'euro deposit rate',
nr_employed DOUBLE COMMENT 'number of employees',
y BIGINT COMMENT 'has time deposit or not'
);
```

After the table is created, click **Table Query** in the left-side navigation bar and enter the table name for search.



Create result_table

Click Data Development > New > Create Table.

On the Create Table page, enter the table creation statements, and click **OK**. The statements used for table creation are as follows:

```
CREATE TABLE IF NOT EXISTS result_table
(
education STRING COMMENT 'educational level',
num BIGINT COMMENT 'number of people'
);
```

After the table is created, click **Table Query** in the left-side navigation bar and enter the table name for search.

Upload local data to bank_data

DataWorks supports the following operations:

- Upload data in local text files to a table in the workspace.
- Use the data integration module to import business data from multiple different data sources to the workspace.

NOTE:

In this section, local files are used as the data source. Local text file uploads have the following restrictions:

- File type: Only .txt and .csv files are supported.
- File size: The file size cannot exceed 10 MB.
- Operation objects: Partition and non-partition tables can be imported, but Chinese partition values are not supported.

Using the import of the local file **banking.txt** to DataWorks as an example, the instruction is as follows:

Click Import > Import Local Data.



Select a local data file, configure the import information, and click Next.

	Selected fi	iles: bai	nking.txt Only .txt	.csv and .lo	g files	are s	upported								
	Delim	iter: 💿	Comma 🔶	◎ 自定义											
Ori	ginal character	set: G	BK												
	First line is t	itle:	Yes												
4	blue-collar	married	basic.4y	unknown	yes	no	cellular	aug	thu	210	1	999	0	nonexistent	1
3	technician	married	unknown	no	no	no	cellular	nov	fri	138	1	999	0	nonexistent	ļ
8	management	single	university.degree	по	yes	no	cellular	jun	thu	339	3	6	2	success	
9	services	married	high.school	no	no	no	cellular	apr	fri	185	2	999	0	nonexistent	
5	retired	married	basic.4y	no	yes	no	cellular	aug	fri	137	1	3	1	success	
0	management	divorced	basic.4y	no	yes	no	cellular	jul	tue	68	8	999	0	nonexistent	
7	blue-collar	married	basic.4y	no	yes	no	cellular	may	thu	204	1	999	0	nonexistent	
39	blue-collar	divorced	basic.9y	по	yes	no	cellular	may	fri	191	1	999	0	nonexistent	

Enter at least two letters to search for the table by name. Select the table to which the data is to be imported, for example, bank_data.

To create a new table, click **Create Table**.

Import loca	al data)
Table:	ba				Create	e Table
Matching:	bank_data Match by position	Match by name				
Target field			Source field			
				Prev	Import	Cancel

Select the field matching method ("Match by Position" is used in this example), and click **Import**.

able: ban	k_data			Create Table
Natching: 💿 Ma	tch by position 🔘 Ma	tch by name		
Target field	So	urce field		
age			\$	
job			\$	
marital			\$	
education			\$	
default			\$	
housing			\$	
loan			\$	

After the file is imported, the system displays a data import success or failure prompt.

Other data import methods

Create a data synchronization task

Applicability:

Data saved in multiple source types, including RDS, MySQL, SQL Server, PostgreSQL, MaxCompute, OCS, DRDS, OSS, Oracle, FTP, dm, HDFS, and MongoDB.

For information on DataWorks data synchronization task creation, see Create a data synchronization task.

Upload a local file

Applicability:

The file size cannot exceed 10 MB, and only .txt and .csv files are supported. Only non-partitioned tables are supported.

For information on DataWorks local file uploads, see the Upload local data to bank_data section.

Use Tunnel commands to upload files

Applicability:

Local files and other resource files are larger than 10 MB.

Using the Tunnel commands provided by the MaxCompute **Client** to upload or download data, you can upload a local data file to a partitioned table.

For details, see Tunnel command operations.

Use DataX open-source tools

Applicability:

DataX is applicable to the import of local data in batches. The imported data must have a twodimensional table structure. This method can be used in other scenarios that described above. For details, see DataX.

For more information about DataX open-source tools, go to the DataX open-source website.

Subsequent steps

You have learned how to create a table and upload data. You can go to the next tutorial for further study. This tutorial shows you how to create a flow for further data analysis and computing in the project space. For details, see Create a flow for data analysis.

The data development function of DataWorks supports the graphic design of data analysis flows and processes data and forms mutual dependencies through flow tasks and inner nodes. Currently, it supports multiple task types, including ODPS_SQL, data synchronization, OPEN_MR, SHELL, machine learning, and virtual nodes. For details about the use of each task type, see Task type description.

This section uses the creation of a flow task named "work" as an example to show how to create nodes in a flow, configure dependencies, and conveniently design and display steps and sequences for data analysis. We briefly describe how to use the data development function for further data analysis and computing in the workspace.

Prerequisites

You have prepared the business data table bank_data, the data it contains, and the result_table in the workspace according to Create a table and upload data.

Instruction

Create a flow

log on to the project, and click **Data Development** > **New** > **Create Task**.



Select the relevant content in the dialog box and specify the task type as Flow task.

Create task		×
*Task type:	Workflow task O Node task	
*Name:	work	
*OSchedule type:	Manual scheduling eriodic scheduling	
Description:		
Select directory:	1	
	> Task development	
	•	
	Create Cance	I

Note : Once selected, the scheduling attribute cannot be changed.

9	P New 🕶	🖺 Save	Submit	🔁 Test run	D Full Screen	🕑 Import 🕶
	work	×				
	Nodes	6				
	Data Proc	cess				
	OPEN_N	WR				
	ODPS_S	QL				
	ODPS_N	ИR				
	Data SY	NC				
	Algorith	ım				
	Script	t				
	SHELI	L				
	Contro	bl				
	Virtua	I				

Create a node and dependency on the flow canvas

This section shows how to create a virtual node "start" and an odps_sql node "insert_data", and to configure "insert_data" to depend on "start".

Note :

- As a control-type node, the virtual node does not affect the data during flow operation and is only used for O&M control of downstream nodes.
- When a virtual node is depended on by other nodes and its status is manually set to failure by the O&M personnel, its downstream nodes that have not yet run cannot be triggered. This prevents further propagation of erroneous upstream data during the O&M process. For details, see the section on virtual nodes in Task type description.

In summary, we recommend that you create a virtual node as the root node to control the whole flow when designing a flow.

Double-click the virtual node, and enter the node name" start" .

Create node			×
*Name :	start		
*Type :	Virtual		
Description :	Enter description		
		Create	Cancel

Double-click **ODPS_SQL** and enter the node name "insert_data" .

Create node			×
*Name :	insert_data		
*Type :	ODPS_SQL		
Description :	Enter description		
		Create	Cancel

Click the start note, and draw a line between start and insert_data to make insert_data dependent on start.



Edit the code in ODPS_SQL

This section describes how to use the SQL code in the ODPS_SQL node **insert_data** to query the quantity of mortgages for individual persons with different education backgrounds, and save the results for analysis or display by subsequent nodes. The SQL statements are as follows. For details about the syntax, see the MaxCompute documentation.

```
INSERT OVERWRITE TABLE result_table --Insert data to result_table
SELECT education
, COUNT(marital) AS num
FROM bank_data
WHERE housing = 'yes'
AND marital = 'single'
GROUP BY education
```

Run and debug ODPS_SQL

After editing the SQL statements in the insert_data node, click Save to prevent code loss.

Click **Run** to view the operations logs and results.

I work ×
← Back ③ Run ① Stop 🗄 Format ③ Cost Estimate
1 INSERT OVERWRITE TABLE result_tableInsert data to result_table
² SELECT education
3 , COUNT(marital) AS num
4 FROM bank_data
5 WHERE housing = 'yes'
6 AND marital = 'single'
7 GROUP BY education
Log
TableSink_REL887317: 8 (min: 8, max: 8, avg: 8)
reader dumps:
StreamLineRead_REL887314: (min: 0, max: 0, avg: 0)
201/-10-19 1/:20:05 INFO ====================================
2017 10 10 17,20,05 TNEO Evit code of the Shall command 0
2017-10-19 17:20:05 INFO Exit code of the Shell command 0
2017-10-19 17:20:05 INFO Exit code of the Shell command 0 2017-10-19 17:20:05 INFO Invocation of Shell command completed 2017-10-19 17:20:05 INFO Shell run successfully!
2017-10-19 17:20:05 INFO Exit code of the Shell command 0 2017-10-19 17:20:05 INFO Invocation of Shell command completed 2017-10-19 17:20:05 INFO Shell run successfully! 2017-10-19 17:20:05 INFO Current task status: FINISH

Then, click Table Query on the left to query data in the table.

Task	All p	roject: 🔻	Q () 🗄	Œ	New 🕶	🖺 Save	🕜 Submit		
deve	~ 🚘	ODPS表			work	×			
elopmen	0	<pre>result_tail</pre>	ble testbyxilin	÷	Back	∋ Run	(I) Stop	88 Format	٩
-				1	INSERT O	VERWRITE	TABLE result	_tableInse	er
Scri				2	SELECT e	ducation			
pt d				3	, CO	UNT(marit	al) AS num		
evel				4	FROM ban	k_data ,	,		
opn				6	WHERE NO	using = ;	yes ,,,		
nent				7	AND CROID BY	maritai -	single		
					GROOP DI	educatio.			
Res									
Resourc	4		ŀ						
Resource	•	Column	► Information						
Resource Fur	•	Column	► Information	Log					
Resource Functio	∢ Filte	Column	Information	Log		TableSin	k_REL887317:	8 (min: 8, ma	DX :
Resource Function	✓ Filte	Column er column	Information	Log	reader	TableSin dumps:	k_REL887317:	8 (min: 8, ma	X:
Resource Function Ta	Filte	Column er column Column nam	Information	Log	reader	TableSin dumps: StreamLi	k_REL887317: neRead_REL883	8 (min: 8, ma 7314: (min: 0, m	ix : ma
Resource Function Table	 Filte □ 	Column er column Column nam education	Information	Log : 0K 2017-	reader 10-19 17:	TableSin dumps: StreamLi 20:05 INFO	k_REL887317: neRead_REL887	8 (min: 8, ma 7314: (min: 0, n	ma
Resource Function Table que	 Filte □ □ 	Column er column Column nam education num	e Column type STRING BIGINT	Log OK 2017- 2017-	reader 10-19 17: 10-19 17:	TableSin dumps: StreamLi 20:05 INFO 20:05 INFO	k_REL887317: neRead_REL883 	8 (min: 8, ma 7314: (min: 0, m f the Shell com	ma ma
Resource Function Table query	 Filt □ □ □ 	Column er column Column nam education num	e Column type STRING BIGINT	Log OK 2017- 2017- 2017- 2017-	reader 10-19 17: 10-19 17: 10-19 17:	TableSin dumps: StreamLi 20:05 INFO 20:05 INFO 20:05 INFO	k_REL887317: neRead_REL883 Exit code of Invocat	8 (min: 8, ma 7314: (min: 0, m f the Shell com ion of Shell co	sma ma ma
Resource Function Table query	Filte	Column er column Column nam education num	e Column type STRING BIGINT	Log OK 2017- 2017- 2017- 2017- 2017- 2017-	reader 10-19 17: 10-19 17: 10-19 17: 10-19 17: 10-19 17:	TableSin dumps: StreamLi 20:05 INFO 20:05 INFO 20:05 INFO 20:05 INFO 20:05 INFO 20:05 INFO	k_REL887317: neRead_REL88: Exit code o Invocat: Shell run su	8 (min: 8, ma: 7314: (min: 0, n f the Shell com ion of Shell co uccessfully! k status: FINIS	ma ma ima ima

Save and submit a flow

After running and debugging the ODPS_SQL node "insert_data", return to the flow page, and save and submit the whole flow.



Subsequent steps

Now, you know how to create, save, and submit a flow. Continue to the next tutorial for further study. This tutorial shows you how to create a synchronization task to export data to data sources of different types. For details, see Create a synchronization task to export results.

Currently data source types supported by the data synchronization jobs include: MaxCompute, RDS (MySQL, SQL Server, PostgreSQL), Oracle, FTP, ADS, OSS, OCS, and DRDS.



Taking the data synchronization from RDS to MaxCompute for example, the detailed descriptions can be found below:

Step 1: Create a data table

For details on how to create a MaxCompute table, refer to Create a Table.

Step 2: Create a data source

Note:

Only users with the project administrator role are allowed to create a new data source.

Preparation

Currently only the China East 1 (Hangzhou) region is supported as a RDS data source, and the Beijing region is not yet supported. In addition, when the RDS data sources in the Hangzhou region cannot be connected to during testing, you should add a whitelist of data synchronization server IP addresses onto your RDS:

10.152.69.0/24, 10.153.136.0/24, 10.143.32.0/24, 120.27.160.26, 10.46.67.156, 120.27.160.81, 10.46.64.81, 121.43, 110.160, 10.117.39.238, 121.43, 112.137, 10.117.28.203, 118.178.84.74, 10.27.63.41, 118.178.56.228, 10.27.63.60, 118.178.59.233, 10.27.63.38, 118.178.142.154, 10.27.63.15

The specific steps are as follows:

Go to Alibaba Cloud Dataplus platform > DataWorks Kit > Console as a developer, click the Enter Work Zone in the action bar of the corresponding project.

Click **Manage Projects** in the top menu bar, and then click **Manage Data Sources** in the left navigation bar.

Click New Data Source.

-) Alibaba big data plat	form coolshell_demo	🚽 Data Dev	elopment Data M	Management	Operation Center	Other 🔺	yangyi.pt@ 👻	English 🗸
III Project Configuration	Data Source Management				:	Project Management		ew Data Source
Project Member Manage	Enter name to search O	Search					J.	
Data Source Management	Data Source Name	Data Source Type	Link Info			Data Source Description	Op	eration
Scheduling Resource List MaxCompute Config	odps_first	odps	ODPS Enderstand American ODPS Project Manager Access Id	ternise adjusel y soa addelly dena (Sall B	eren/agei	connection from odps calc e	engine 121	
	coolshell_ads	ads	Connext Tel	illan kakar berjing Skalit	Labalyana ang 10		del	lete edit

Fill in the configuration items in the "New Data Source" pop-up box.

* Data Source Name :	Enter a data source name	
Data Source Description :	Enter a data source description	
* Data Source Type :	rds v mysql v	
* RDS Instance ID :	Enter the RDS instance ID	
* RDS Instance Purchaser ID:	enter the RDS Instance Purchaser ID	
	How to find the ID of the RDS instance purchaser, clickhere	
* Database Name :	Enter the RDS database name	
* User name :	Enter the RDS user name	
* Password :	Enter the RDS password	
	You need to add the data source to the RDS whitelist to connect it successfully, Click here to view how to add an entry to the whitelist.	

Specific descriptions of the configuration items in the figure above are as follows:

- Data source name: A data source name may consist of letters, numbers, and underscores. It must begin with a letter or an underscore and cannot exceed 30 characters in length.
- Data source descriptions: A brief description of the data source. The description should not exceed 1,024 characters in length.
- Data source type: The data source type selected currently (RDS>MySQL>RDS).
- RDS instance ID: The ID of the MySQL data source RDS instance.
- RDS instance purchaser ID: The purchaser ID of the MySQL data source RDS instance.

Note: If you have selected the JDBC form to configure the data source, the format of the JDBC connection information is: jdbc:mysql://IP:Port/database.

- Database name: The database name of the data source.
- User name/password: The user name and password of the database.

Click Test Connectivity.

If the test result is connected successfully, click the **Save** button to save the configuration information.

For detailed configurations of other types of data sources (MaxCompute, RDS, Oracle, FTP, ADS, OSS, OCS, and DRDS), see Data Source Configuration.

Step 3: Create a new job

Take the "wizard mode" new task as an example.

1. In the data integration interface, click on the left navigation bar to synchronize tasks;

Click the wizard mode in the interface to get to the task configuration page.

New Synchronization Tasks :



Step 4: Configure the data synchronization job

The synchronization job node includes five configuration items: **"Select Data Source and Target"**, **"Field Mapping"**, **"Channel Control"** and **Preview & Save**.

Select the data source

Select Data Source (The data source has been created in Step 2), and then select the data table.

1 Select Source S	elect Target	Field Mapping	Channel Control	Preview & Save	9
ou may need to select the source type of data, * Data Source :	it can be your own dw_log_detail_rd	independent databas s (mysql)	e server, or RDS in Alib	aba Cloud, see su	pport data source typ
* Table:	`adm_user_meas	ures' X		~	
Data Filter:	New Data Source	e + Γ(createtime,'%Y-%	m-%d')='\${ct}'		
Split Key:	device				
	l l	Preview Data 🗸			

Extraction Filtering: You can specify the WHERE filter based on the corresponding SQL syntax (You do not need to specify the WHERE keyword). The WHERE filter will be used as a condition of incremental synchronization.

The WHERE filter is used for source data filtering. The specified column, table, and WHERE filter are concatenated to create an SQL command for data extraction. The WHERE filter can be used for full synchronization and incremental synchronization. Specific descriptions are as follows:

• Full synchronization:

Full synchronization is usually executed when data is imported for the first time. You do not need to configure the WHERE filter. You can set the WHERE filter limit to 10 to avoid a large data size during tests.

 Incremental synchronization: In the actual service scenario, incremental synchronization usually synchronizes the data generated on the current day. Before compiling the WHERE filter, you usually need to first determine the field that describes the increment (timestamp) in the table. For example, if in Table A, the field that describes the increment is "creat_time", you need to compile "creat_time>\$yesterday" in the WHERE filter and assign a value to the parameter in parameter configuration.

Splitting key: If the data synchronization job is RDS/Oracle/MaxCompute, the splitting key configuration will be displayed on the page. **only supports integer fields.** During data reading, the data will be split based on the configured fields to achieve concurrent reading, improving data synchronization efficiency. The splitting key configuration item will only be displayed when the synchronization job is for importing RDS/Oracle data into MaxCompute.

If the source is a MySQL data source, the data synchronization job also supports databaseand table-based data importing (on a premise that the table structure must be consistent, no matter whether the data is stored in the same database or different databases).

Database- and table-based data importing supports the following scenarios:

Multiple tables in the same database: Click **Search Table** to search for the tables and add the tables you want to synchronize.

Multiple tables in different databases: First click **Add** to select the source database, and then click **Search Table** to add the tables.

Select the data target

Click **rapid establishment of table** and you will be able to convert the tabulation statements of the source table to DDL statements conforming to the MaxCompute SQL syntax to create a target table. After making the necessary selections, click **Next**.

Select Source Se	elect Target		Channel Control	Preview 8	Save
ou may need to select the destination type of data	, it can be your own	i independent data	base server, or RDS in Alibal	oa Cloud, s	see support the data target type
* Data Source :	odps_first (odps))		\sim	
* Table:	my_region			\checkmark	rapid establishment of table
* Partition:	pt	=	\${bdp.system.bizdate}		0
cleansing rules:	 write before cl 	eaning with availa	ole data Insert Overwrite		
	oformer reserva	itions have been in	cluded in the data Insert Inte	D	

Partition information: Partitioning helps you to easily search for the special columns introduced by some data. By specifying the partition, you can quickly locate the desired data. Constant partitions and variable partitions are supported.

Clearing rules:

Clear existing data before writing: Before data importing, all the data in the table or partition should be cleared, which is equivalent to **Insert Overwrite**.

Keep existing data before writing: No data needs to be cleared before data importing. New data is always appended with each run, which is equivalent to **Insert into**.

Assign values to parameters in the parameter configuration, as shown in the figure below:

System parameter config	guration 😡	Schedu
\${bdp.system.bizda	te} yyyyMMdd	uling contigu
User-defined parameter	configuration 😡	ration
ct	\$[yyyy-mm-dd-1]	Paramete
		r configurati

Field Mapping

You need to configure the field mapping relationships. The **Source Table Fields** on the left correspond one to one with the **Target Table Fields** on the right.

	Sel	ect Source	Select Target	3 Field Mapping Ch	annel Control	5 Preview & Save	
ou may	need to configure the	source table and com	I the destination table plete the mapping by	mapping relationship, conne peer mapping。 data synchro	ct the fields to be a nization document	synchronized via the connection, or you It	can
	Source Table Field	Туре		Target Table Field	Туре	Auto Mapping	
	device	VARCHAR	•	e device	STRING	Auto Layout	
	pv	BIGINT	•	pv .	BIGINT		
	UV	BIGINT	•	uv uv	BIGINT		
	createtime	DATETIME	•	createtime	DATETIME		
	New Row +						

- Add/Delete: Click **Add a Line** to add a single field. Move and hover the cursor on a line above, and click the **Delete** icon, and you will delete the current field.

Writing method for user-defined variables and constants:To import a constant or variable to a field in the MaxCompute table, you only need to click the **Insert** button and enter the value of the constant or variable enclosed in single quotation marks. For example, for the

'\${yesterday}' variable, you can then assign a value to the variable using the parameter configuration component, such as yesterday=\$[yyyymmdd].

Channel Control

The **Channel Control** is used to configure the maximum speed of the job and the dirty data check rules, as shown in the figure:

Select Source Se	elect Target	Field Mapping	4 Channel Control	Preview	& Save
You can configure the transfer rate of the job and th	he number of error le	ogs to control the ent	ire data synchronizatior	process	, data synchronization document
* Maximum Speed Rate :	10MB/s			\sim	0
Incorrect records more than :	Dirty data num	ber range, allow di	rty data default		number, to end task 🛛 🕥

- The maximum speed of the job refers to the speed of the current data synchronization job, with a maximum value of 10 MB/s supported (The channel traffic measured value is the measured value of the data synchronization job, and does not represent the actual traffic of the network interface card).

Dirty data check rules (available for writing data to RDS and Oracle):

- When the number of error records (that is the volume of dirty data) exceeds the configured quantity, the data synchronization job ends.

Preview & Save

When you complete the above configuration, click **next** to preview, if correct, click **save**, as shown below:

Select Source	Select Target	Field Mapping	Channel Control	Preview & Save
Please confirm and save the configured i	nformation that you c	an test to run or config	ure the scheduling prop	erties, data synchronization document
Select Source				Edit
* Data Sour	ce: dw_log_detail_	rds		
* Ta	ble: `adm_user_me	asures		0
Data Fil	ter: DATE_FORM	AT(createtime;%Y-	%m-%d')='\${ct}'	
Split H Select Target	ey: Unfilled			Edit
		Previous Save	1	

Step 5: Submit the data synchronization job and test the workflow

Click the top menu bar to submit the job.

After the job is submitted successfully, click Test Run.

Because some createtime values in the source table in this example are 2017-01-04, while the scheduling time parameters used in the configuration are \$[yyyy-mm-dd-1] and \${bdp.system.bizdate}, we set the partition value of the target table to 20170104 to assign the value of 2017-01-04 to the createtime parameter in the test. The 2017-01-04 should be selected as the business time in the test, as shown in the figure below:

Instance name:	data_sync_2017_03_09	
*Business date:	2017-01-04	
f the selected bu	iness date is before yesterday, the task will be executed immediately.	
f the selected bu	iness date is yesterday, the task will be executed at the specified time.	

After the test task is triggered successfully, you can click **Go to O&M Center** to view the task progress.



Task name :	data sync
Current status :	Run successfully
Status description :	Instance run successfully
Application name :	coolshell_demo
Job type :	Data Synchronization
Regular time :	2017-01-05 00:00:00
Start time :	2017-03-09 18:14:07
End time :	2017-03-09 18:14:40
Owner :	yangyi.pt@aliyun-test.com

View the synchronized data.

1 read m	ny_region ;			
Log	Results[1] ×			
No.	device	pv	uv	createtime pt
1	android	937	73	2017-01-04 20:51 20170104
2	iphone	428	31	2017-01-04 20:49 20170104
3	macintosh	830	107	2017-01-04 20:51 20170104
4	unknown	4124	444	2017-01-04 20:51 20170104
5	windows_pc	5650	649	2017-01-04 20:51 20170104

DataWorks provides powerful scheduling capabilities including time-based or dependency-based

task trigger mechanisms to perform **tens of millions** of tasks accurately and punctually each day based on DAG relationships. It supports scheduling by minute, hour, day, week and month. For details, see Scheduling configuration description.

This section uses write_result created in Create a synchronization task as an example and configures the scheduling period to weekly, to explain the scheduling configurations and task O&M functions of DataWorks.

Instruction

Configure the scheduling attribute of a synchronization task

Go to the **Data Development > Task Development** page, and double-click the synchronization task you want to configure (write_result), then click **Scheduling Configuration** to configure the **scheduling attribute** of the task, as shown in the following figure:

- Basic attribute	es 🕨		Î	Sche
- Scheduling at	tribute 👻 —			duling
Scheduling status:	Frozen) configurati
Auto	🔲 open ?)		g
icuy.				Para
Activation date:	1970-01-01	to 2116-10-20		meter co
*Scheduling period:	Day	*		nfiguratio
*Specific time:	00	♣ : 00 ♣		2

The configuration parameters are described below:

- Scheduling status: When this parameter is selected, the task is paused.
- Error retry: When this parameter is selected, error retry is enabled.
- Start date: The date on which the task takes effect, which can be set based on actual needs.
- Scheduling period: The operating period of the task, which can be set by month, week, day, hour, and minute. For example, a task can be scheduled weekly.
- Specific time: The specific operating time of the task. For example, you can set up the task to

run at 02:00 each Tuesday.

Configure the dependency attribute of a synchronization task

After configuring the scheduling attribute of a task, you can configure its dependency attribute, as shown in the following figure:

Dependency a	attribute 🔻 ———		
Project:			
Upstream task:	Enter a keyword	to query upstr	re Q
Project name	Task name	Owner	Actions
	work	shu	Delete

- Cro	ss_cvcle dependency 👻
OIC	33-cycle dependency +
۲	Not dependent on the previous scheduling period
\bigcirc	Self-dependent; operation can continue after the conclusion of the previous scheduling period
\bigcirc	Operation can continue after the conclusion of the previous downstream task scheduling period
0	Operation can continue after the conclusion of the previous custom task scheduling period

You can configure an upstream dependency for a task. In this way, even if the scheduled time of an instance of the current task is reached, the task can run only after the instance of its upstream task is completed.

The configuration in the above figure indicates that instances of the current task will be triggered only after the instance of the upstream task write_result is finished. You can enter "work" in the upstream task to configure an upstream task for write_result.

If no upstream task is configured, the current task is triggered by the project by default. Therefore, by default, the upstream task of the current task is project_start in the scheduling system. By default, a project_start task is created as a root task for each project.

Submit a synchronization task

Save the synchronization task "write_result", and click **Submit** to submit it to the scheduling system, as shown in the following figure:

운 New 🕶 🖹 Save 🕜 Submit 🗖 Test run (고) Full Screen 🕑 Import 🕶	⊖ Go to O&M
Image: select Source Select Target	- Basic attributes > * Sreeduling attribute > * Scheduling Frozen status: *
Select Target	Activation 1970-01-01 to 2116-10-20 the construction of the constr
Pre-import statement: Unfilled	*Specific 00 🔶 : 00 🔶
Prepare statements after import: Unfilled	Dependency attribute Cross-cycle dependency Cross-cycle dependency Not dependent on the previous scheduling period Self-dependent; operation can continue after the conclusion of the previous scheduling period Operation can continue after the conclusion of the

The system automatically generates an instance for the task at each time point according to the scheduling attribute configuration and periodically runs the task from the second day only after a task is submitted to a scheduling system.

NOTE:

If a task is submitted after 23:30, the scheduling system automatically generates instances for the task and periodically runs the task from the third day.

Subsequent steps

Now you know how to set the scheduling attribute and dependency of a synchronization task. Continue to the next tutorial for further study. This tutorial shows you how to perform periodic O&M for submitted tasks and view the log troubleshooting results. For details, see Perform periodic O&M and view log troubleshooting results.

In the previous operations, you have set a synchronization task to be run at 02:00 every Tuesday. After the task is submitted, you can view the automatic operation results in the scheduling system from the second day. Now, how can we check whether the instance schedule and dependency are as expected? DataWorks provides three triggering methods: test run, data population, and periodic running. Details about the three methods are as follows:

Test run: The task is triggered manually. If you need to check the timing and operation of a single task, test run is recommended.

Data population: The task is triggered manually. This method applies if you need to check the timing and dependencies of multiple tasks or re-execute data analysis and computing from a root task.

Periodic running: The task is triggered automatically. After a task is submitted successfully, the scheduling system automatically generates task instances at different time points starting from 00:00 of the second day. It checks whether upstream instances of each instance have run successfully at the scheduled time. If all the upstream instances have run successfully at the scheduled time, the current instance runs automatically without manual intervention.

NOTE:

The scheduling system periodically generates instances based on the same rules that apply in both manual and automatic triggering modes.

The period can be set to monthly, weekly, daily, hourly, or even by minute. The scheduling system always generates an instance for the task on the specified day or at the specified time.

The scheduling system only regularly runs the instance on the specified date and generates operation logs.

Instances rather than on the specified date are not run, and their statuses are directly changed to "Successful" when the running conditions are met. Therefore, no running logs are generated.

The following instructions show how to configure these three triggering methods.

Test run

Manually trigger the test run

Click the Test Run button on the flow page.

•	Data Inte	egration	Data Develo	pment	Data	Management
[+] New ▼	🖺 Save		Test run	D Full	Screen	➢ Import ◄
🔝 work	×					

As promoted on the page, click **Confirm** and **Run**.

Cancel OK Test run Instance name: work_2017_10_20 *Business date: 2017-10-19 If the selected business date is before yesterday, the task will be executed immediately. If the selected business date is yesterday, the task will be executed at the specified time.	his operation may affe	ct the data output by cyclically scheduled tasks. Proceed with caution!	
Test run Instance name: work_2017_10_20 *Business date: 2017-10-19 If the selected business date is before yesterday, the task will be executed immediately. If the selected business date is yesterday, the task will be executed at the specified time.			Cancel OK
Instance name: work_2017_10_20 *Business date: 2017-10-19 If the selected business date is before yesterday, the task will be executed immediately. If the selected business date is yesterday, the task will be executed at the specified time.	Fest run		>
Instance name: work_2017_10_20 *Business date: 2017-10-19 * If the selected business date is before yesterday, the task will be executed immediately. * If the selected business date is yesterday, the task will be executed at the specified time.			
*Business date: 2017-10-19	Instance name:	work_2017_10_20	
* If the selected business date is before yesterday, the task will be executed immediately. * If the selected business date is yesterday, the task will be executed at the specified time.	*Business date:	2017-10-19	
* If the selected business date is yesterday, the task will be executed at the specified time.	* If the selected busi	ess date is before yesterday, the task will be executed immediately.	
	* If the selected busin	ess date is yesterday, the task will be executed at the specified time.	

Click Go to O&M Center to view the task operation status.

Workflow task test run	×
Workflow task test run triggered. Go to the O&M center to view progress.	
	Cancel Go to O&M Center

View the information and operation logs of the test instance

Click the task name to view the instance DAG. In the instance DAG view, right-click an instance to view its dependencies and detailed information. Also, you can terminate or re-run the instance. In the instance DAG view, double-click an instance and a dialog box appears, showing the task attributes, running logs, operation logs, and code.

Note:

In test run mode, the task is triggered manually. The task runs immediately as long as the set time is reached, regardless of the instance' s upstream dependencies.

According to the previously mentioned instance generation rules, set up the task write_result to run at 02:00 each Tuesday. If the business date of test run is Monday (business date = running date -1), the instance runs at 02:00. If not, the instance status is changed to "Successful" at 2:00 and no logs are generated.

Data population

Manually trigger data population

If you need to check the timing and dependency of **multiple tasks** or re-execute data analysis and computing from a root task, go to the **O&M Center > Task List > Task Scheduling** page and click **Data Population Task** to run multiple tasks of a specific period of time.

Instruction

Log on to the O&M Center > Task Scheduling page and enter the task name.

Select the task query results and click the Data Population button.

Set the business date of the data population as May 11, 2017 to May 12, 2017, select the insert_data and write_result node tasks, and click **OK**.

Click View Data Population Results.

View the information and operation logs of the data population instance

On the **Data Population Instance** page, find the task instance: Click the task name to view the instance DAG. In the instance DAG view, right-click an instance to view its dependencies and detailed information. Also, you can terminate or re-run the instance. In the instance DAG view, double-click an instance and a dialog box appears, showing the task attributes, running logs, operation logs, and code.

Note:

Data population task instances are dependent on instances from the previous day. For example, for a data population task within the period from September 15, 2017 to September 18, 2017, if the instance on the 15th is failed to run, the instance on the 16th is not run.

According to the previously mentioned instance generation rules, set up the task write_result to run at 02:00 each Tuesday. If the business date selected during data population is Monday (service date = running date -1), the instance runs at 02:00. If not, the instance status is changed to "Successful" at 02:00 and no logs are generated.

Periodic automatic run

In periodic automatic run mode, the scheduling system automatically triggers tasks according to all task scheduling configurations. Therefore, no operation portal is provided on the page. You can view the instance information and operation logs in either of the following methods:

Go to the **O&M Center > Task Scheduling** page, select parameters such as service date or running date, search instances corresponding to the task write_result, and then right-click on an instance to view its information and operation logs.

Click the task name to view the instance DAG. In the instance DAG view, right-click an instance to view its dependencies and detailed information. Also, you can terminate or rerun the instance. In the instance DAG view, double-click an instance and a dialog box appears, showing the task attributes, running logs, operation logs, and code.

Note:

If the initial status of a task instance is "Not Run", when the scheduled time is reached, the scheduling system checks whether all the upstream instances are successful.

The instance is triggered only when all of its upstream instances are successful and its scheduled time is reached.

For an instance in Not Run status, check that all its upstream instances are successful and its scheduled time has been reached.