

Elastic MapReduce

Data Migration

Product Documentation



Copyright Notice

©2013-2019 Tencent Cloud. All rights reserved.

Copyright in this document is exclusively owned by Tencent Cloud. You must not reproduce, modify, copy or distribute in any way, in whole or in part, the contents of this document without Tencent Cloud's the prior written consent.

Trademark Notice



All trademarks associated with Tencent Cloud and its services are owned by Tencent Cloud Computing (Beijing) Company Limited and its affiliated companies. Trademarks of third parties referred to in this document are owned by their respective proprietors.

Service Statement

This document is intended to provide users with general information about Tencent Cloud's products and services only and does not form part of Tencent Cloud's terms and conditions. Tencent Cloud's products or services are subject to change. Specific products and services and the standards applicable to them are exclusively provided for in Tencent Cloud's applicable terms and conditions.

Contents

Data Migration

- Getting Started

- Using COS to Migrate Data

- How to View COS Info

- Using DistCp to Migrate Data

 - Brief Description

 - Network Connectivity

 - Copy

 - Must Knows

Data Migration

Getting Started

Last updated : 2019-07-26 17:53:19

EMR currently allows you to migrate your data in two ways: 1. Migrating your data to COS; 2. Copying EMR cluster data to your self-built HDFS using DistCP. This method requires your self-built cluster is connect with the EMR cluster network.

Using COS to Migrate Data

Last updated : 2019-07-26 17:53:25

- Migrating non-HDFS files

If your source file is a non-HDFS file, upload it to COS via COS Console or API, and then analyze it in the EMR cluster.

- Migrating HDFS files

- i. Get the COS migration tool

Click [here](#) to download the migration tool. For more migration tools, see [here](#).

- ii. Configuring the tool

All configuration files are stored in the conf directory of the tool directory. Copy the core-site.xml file of the HDFS cluster to be synced to conf, which contains the configuration information of the NameNode. Edit the configuration file cos_info.conf by including your appid, bucket, region, and key information.

- iii. Command parameter descriptions

```
-ak <ak> the cos secret id
-appid,--appid <appid> the cos appid
-bucket,--bucket <bucket_name> the cos bucket name
-cos_info_file,--cos_info_file <arg> the cos user info config default is ./conf/cos_info.conf
-cos_path,--cos_path <cos_path> the absolute cos folder path
-h,--help print help message
-hdfs_conf_file,--hdfs_conf_file <arg> the hdfs info config default is ./conf/core-site.xml
-hdfs_path,--hdfs_path <hdfs_path> the hdfs path
-region,--region <region> the cos region. legal value cn-south, cn-east, cn-north, sg
-sk <sk> the cos secret key
-skip_if_len_match,--skip_if_len_match skip upload if hadoop file length match cos
```

- iv. Executing data migration

```
# All operations must be performed in the tool directory. If both configuration files and comma
nd line parameters are set, the latter will prevail
```

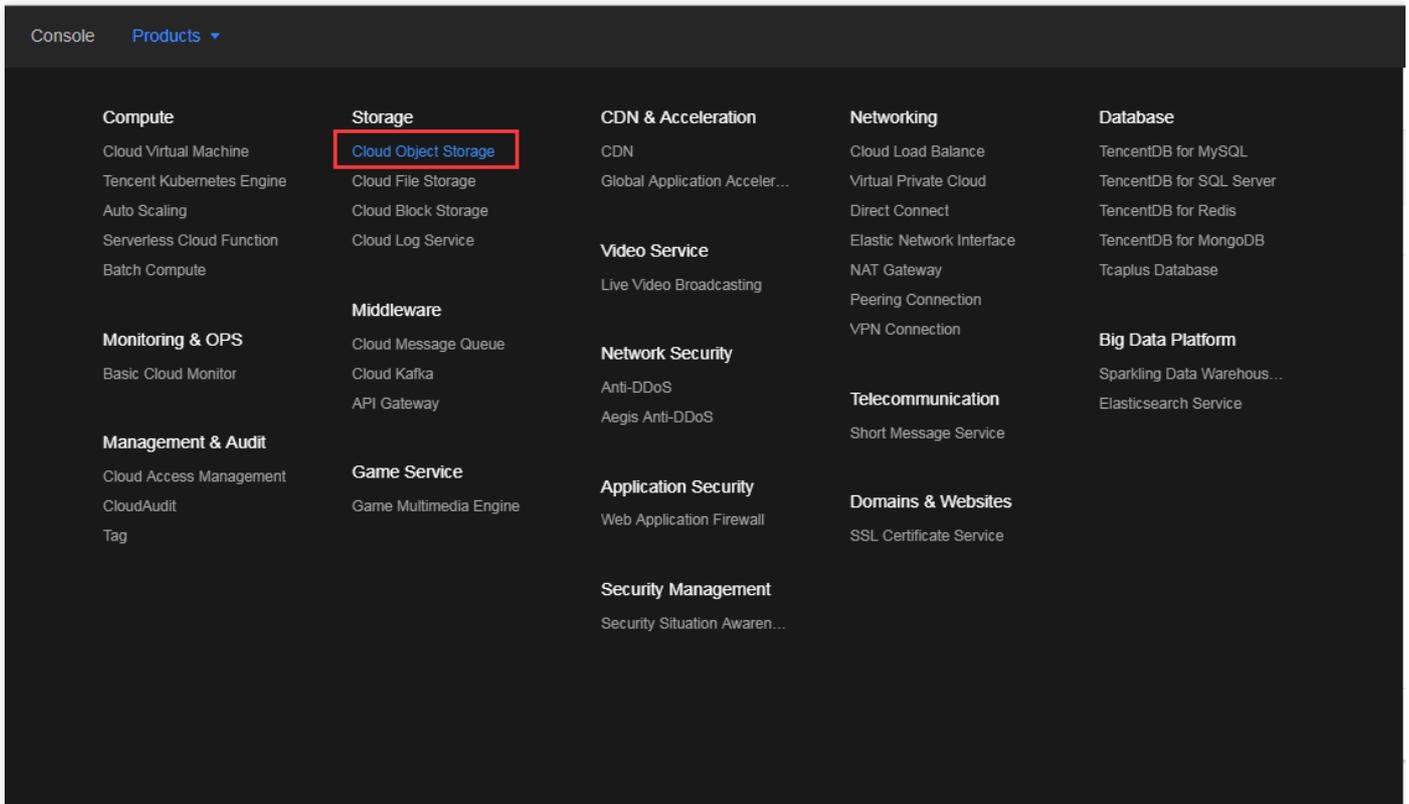

data in Hadoop.

- Files that already exist in COS are overwritten by default in case of repeated upload, unless you explicitly specify the `-skip_if_len_match` parameter, which indicates to skip files if they have the same length as the existing files.
- The COS path is always considered as a directory, and files that are eventually copied from HDFS will be stored in this directory.

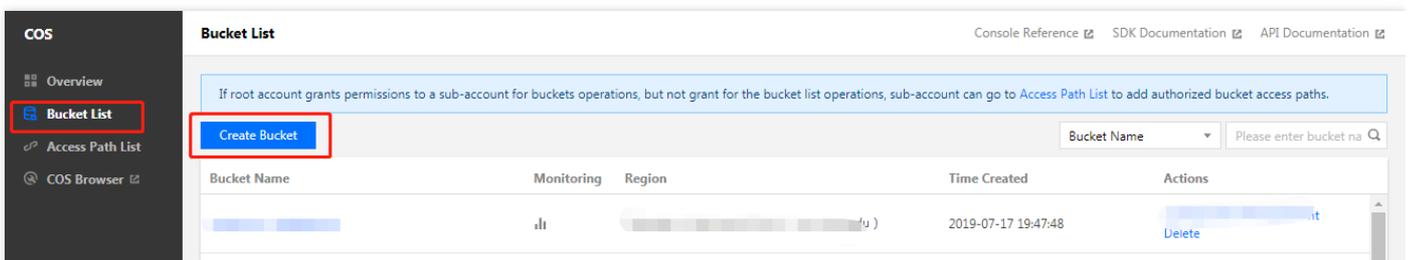
How to View COS Info

Last updated : 2019-07-26 17:53:34

- Log in to your Tencent Cloud account, go to the console and select Cloud Object Storage



- Select a bucket



- View the bucket's domain name

20190718test01-1258469122

Objects Basic Configuration **Domain Management** Permission Management

Default CDN Acceleration Domain [Edit](#)

Status Off

[Learn more](#)

User-defined Acceleration Domain

Domain	Acceleration Region	Origin Type	Origin-pull Authentication	CDN Authentication	CNAME	State	Actions
Add Domain							

Note: Make sure that the added domain name has been filed and the corresponding CNAME has been set on the DNS provider website. For more information or help, please refer to [Learn more](#)

- Query Secretkey

Bucket List

If root account grants permissions to a sub-account for buckets operations, but not grant for the bucket list operations, sub-account can go to [Access Path List](#) to add authorized bucket access paths.

[Create Bucket](#) Bucket Name Please enter bucket name

Bucket Name	Monitoring	Region	Time Created	Actions
20190717-1258469122		Chengdu (Mainland China) (ap-chengdu)	2019-07-17 19:47:48	Configuration Management Delete
20190718test01-1258469122		Guangzhou (Mainland China) (ap-guangzhou)	2019-07-18 11:55:11	Configuration Management Delete
211-test-1258469122		Guangzhou (Mainland China) (ap-guangzhou)	2019-06-24 17:23:42	Configuration Management Delete
211-test-chongqing-1258469122		ap-chongqing (Mainland China) (ap-chongqing)	2019-06-25 17:14:19	Configuration Management Delete
beijing-1258469122		ap-beijing (Mainland China) (ap-beijing)	2019-01-18 15:53:02	Configuration Management Delete

Using DistCp to Migrate Data

Brief Description

Last updated : 2019-07-26 17:53:45

DistCp (distributed copy) is a tool used for large inter/intra-cluster copying. It uses MapReduce to effect its distribution, error handling and recovery, and reporting. It expands a list of files and directories into input to map tasks, each of which will copy a partition of the files specified in the source list. It is a file migration tool that comes with Hadoop.

Network Connectivity

Last updated : 2019-08-12 20:04:18

Migrating Files in a Local Self-built HDFS to EMR

The migration of files in a local self-built HDFS to an EMR cluster requires a Direct Connect line for network connectivity. You can contact the developers for assistance.

Migrating Files in a Self-built HDFS in CVM to EMR

- If the network where the CVM instance resides and the one where the EMR cluster resides are in the same VPC, the files can be transferred freely.
- Otherwise, a peering connection is required for network connectivity.

Using a Peering Connection

IP range 1: Subnet A 192.168.1.0/24 in VPC1 of Guangzhou.

IP range 2: Subnet B 10.0.1.0/24 in VPC2 of Beijing.

1. Log in to the Tencent Cloud Console and select **Products** > **Networking** > **Virtual Private Cloud** in the navigation bar.
2. In the VPC Console, go to the "Peering Connections" page, select the region "Guangzhou" in the list at the top, select "VPC1", and click **Create**.

The screenshot shows the Tencent Cloud console interface for Peering Connections. The region is set to South China (Guangzhou) and the VPC list is expanded to show all VPCs. A table of existing peering connections is visible, with one connection highlighted.

ID/Name	M...	Status	Local Region	Local VPC	Peer Region	Peer account	Peer VPC	Ban...	Serv...	Billing Mode	Operation
peering-connection-1		Connect...	South China (...	vpc-e0krgxj7 farley_test_vpc(...	South China (...	My Account	vpc-apsw6eh9	Unlimi ted	Gold	Free	Delete

3. Go to the peering connection creation page.

- Enter the name of the peering connection as shown in box 1, for example, PeerConn.
- Enter the local region and network information as shown in box 2, for example, Guangzhou and VPC1.
- Enter the account of the opposite network as shown in box 3. If the two networks in Guangzhou and Beijing are under the same account, select **My account**; otherwise, select "Another account". It should be noted here that if both the local network and the opposite network are in the same region (e.g., Guangzhou), the communication is free of charge, and there is no need to select the upper limit for bandwidth as shown in box 5; otherwise, fees will be incurred and the upper limit for bandwidth can be set.
- Enter the opposite region and network as shown in box 4, for example, Beijing and VPC2.

Create a peering connection ✕

Name

Local Region

Local network

Peer account type My Account Other accounts

Peer Region

Peer network

Bandwidth Cap No restriction

Billing method Free

4. A peering connection between VPCs under the same account takes effect immediately after creation; otherwise, it can take effect only after the opposite account accepts it.
5. Configure the local and opposite route tables for the peering connection.

i. Go to the **Subnets** page in the VPC Console.

ID/Name	Network	CIDR	IPv6 CIDR	Availability Zone	Associated route table	Subnet broadcast	Cloud VPC	Available IP	Default S...	Operation
subnet-44azita4 PrivateLink...	vpc-ezt5qmz Default-VPC		-	e ①	rtb-80ccr7fw default	<input type="checkbox"/>	0	4055	No	Delete Change route table
subnet-8f48cy... subnet-fair...	vpc-e0krqx7 farley_test_vpc		-	Guangzhou Zo...	rtb-b0ex0vqo	<input type="checkbox"/>	0	253	No	Delete Change route table
subnet-h3f7gtay emr_eric_t...	vpc-r5l8g7xf emr_eric		-	Guangzhou Zo...	rtb-46hnmj46 default	<input type="checkbox"/>	2	251	No	Delete Change route table
subnet-om0zb... ...	vpc-34mkwjhd			Guangzhou Zo...	rtb-8gxp30u default	<input type="checkbox"/>	1	249	No	Delete Change route table
subnet-5weh8lis subnet-fair...	vpc-e0krqx7 farley_test_vpc		-	Guangzhou Zo...	rtb-b0ex0vqo	<input type="checkbox"/>	0	252	No	Delete Change route table

ii. Click "Associated Route Tables" (as shown in the red box above) of the subnet of the local network (e.g., subnet VPC1 in Guangzhou) to enter the route table details page.

iii. Click **Edit** to edit the routing policy. Enter the opposite CIDR (for example, the CIDR of VPC2 in Beijing is 10.0.0.1/24) for the destination, select **Peering connection** for the next hop type, and select the created peering connection (PeerConn) for the next hop.

Destination	Next hop type	Next hop	Notes	Enable routing	Operation
Local	Local	Local	Released by the system by default, ...	<input checked="" type="checkbox"/>	ⓘ

iv. You've configured the route table from Guangzhou VPC1 to Beijing VPC2 in the previous steps. Now you need to repeat the steps above to configure the route table from Beijing VPC2 to Guangzhou VPC1.

v. After the route tables are configured, different VPC IP ranges can communicate to each other.

Copy

Last updated : 2019-07-26 17:54:18

```
# Copy the specified folder from one cluster to another
hadoop distcp hdfs://nn1:9820/foo/bar hdfs://nn2:9820/bar/foo

# Copy the specified file
hadoop distcp hdfs://nn1:9820/foo/a hdfs://nn1:9820/foo/b hdfs://nn2:9820/bar/foo

# If too many files were specified, use -f parameter to separate them.
```

Must Knows

Last updated : 2019-07-26 17:54:23

1. For the commands above, the source and destination versions must be the same.
2. You may fail to copy the source file if another client is writing data to it; you will fail to rewrite a file if it is being copied to the destination; you will fail to copy the source file if it was being moved, and you will see an error saying `FileNotFoundException`.