

CDP Private Cloud Base 7.0.3

Managing Clusters

Date published: 2020-11-30

Date modified: 2020-11-30

CLOUDBERA

<https://docs.cloudera.com/>

Legal Notice

© Cloudera Inc. 2024. All rights reserved.

The documentation is and contains Cloudera proprietary information protected by copyright and other intellectual property rights. No license under copyright or any other intellectual property right is granted herein.

Unless otherwise noted, scripts and sample code are licensed under the Apache License, Version 2.0.

Copyright information for Cloudera software may be found within the documentation accompanying each component in a particular release.

Cloudera software includes software from various open source or other third party projects, and may be released under the Apache Software License 2.0 (“ASLv2”), the Affero General Public License version 3 (AGPLv3), or other license terms. Other software included may be released under the terms of alternative open source licenses. Please review the license and notice files accompanying the software for additional licensing information.

Please visit the Cloudera software product page for more information on Cloudera software. For more information on Cloudera support services, please visit either the Support or Sales page. Feel free to contact us directly to discuss your specific needs.

Cloudera reserves the right to change any products at any time, and without notice. Cloudera assumes no responsibility nor liability arising from the use of products, except as expressly agreed to in writing by Cloudera.

Cloudera, Cloudera Altus, HUE, Impala, Cloudera Impala, and other Cloudera marks are registered or unregistered trademarks in the United States and other countries. All other trademarks are the property of their respective owners.

Disclaimer: EXCEPT AS EXPRESSLY PROVIDED IN A WRITTEN AGREEMENT WITH CLOUDERA, CLOUDERA DOES NOT MAKE NOR GIVE ANY REPRESENTATION, WARRANTY, NOR COVENANT OF ANY KIND, WHETHER EXPRESS OR IMPLIED, IN CONNECTION WITH CLOUDERA TECHNOLOGY OR RELATED SUPPORT PROVIDED IN CONNECTION THEREWITH. CLOUDERA DOES NOT WARRANT THAT CLOUDERA PRODUCTS NOR SOFTWARE WILL OPERATE UNINTERRUPTED NOR THAT IT WILL BE FREE FROM DEFECTS NOR ERRORS, THAT IT WILL PROTECT YOUR DATA FROM LOSS, CORRUPTION NOR UNAVAILABILITY, NOR THAT IT WILL MEET ALL OF CUSTOMER’S BUSINESS REQUIREMENTS. WITHOUT LIMITING THE FOREGOING, AND TO THE MAXIMUM EXTENT PERMITTED BY APPLICABLE LAW, CLOUDERA EXPRESSLY DISCLAIMS ANY AND ALL IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO IMPLIED WARRANTIES OF MERCHANTABILITY, QUALITY, NON-INFRINGEMENT, TITLE, AND FITNESS FOR A PARTICULAR PURPOSE AND ANY REPRESENTATION, WARRANTY, OR COVENANT BASED ON COURSE OF DEALING OR USAGE IN TRADE.

Contents

- Accessing the Cloudera Manager Admin Console.....7**
- Starting, Stopping, Refreshing, and Restarting a Cluster.....7**
- Pausing a Cluster in AWS.....9**
 - Shutting Down and Starting Up the Cluster.....9
- Renaming a Cluster.....10**
- Managing Hosts.....10**
 - Viewing Host Status..... 10
 - Adding a Host to a Cluster..... 11
 - Using the Add Hosts Wizard to Add Hosts.....12
 - Adding a Host by Installing the Packages Using Your Own Method..... 18
 - Parcels..... 19
 - Configuring Hosts..... 19
 - Viewing Host Role Assignments..... 19
 - Host Templates..... 20
 - Creating a Host Template.....20
 - Editing a Host Template.....20
 - Applying a Host Template to a Host..... 21
 - Hosts Disks Overview..... 21
 - Deleting Hosts.....22
 - Deleting a Host from Cloudera Manager.....22
 - Removing a Host From a Cluster..... 22
 - Stopping All the Roles on a Host..... 23
 - Starting All the Roles on a Host..... 23
 - Changing Hostnames..... 23
 - Moving a Host Between Clusters..... 25
 - Using Upgrade Domains to manage rolling restarts..... 26
 - Configuring Upgrade Domains..... 26
 - Adding or changing the Upgrade Domain for a single host..... 27
 - Putting all Hosts in an Upgrade Domain group into Maintenance Mode..... 27
 - Specifying Racks for Hosts..... 28
 - Performing Maintenance on a Cluster Host..... 28
 - Decommissioning Hosts..... 29
 - Recommissioning Hosts.....30
 - Tuning and Troubleshooting Host Decommissioning.....30
 - Maintenance Mode.....34
 - Viewing the Maintenance Mode Status of a Cluster..... 36
- Managing Roles.....36**
 - Role Instances..... 37
 - Adding a Role Instance..... 37

Starting, Stopping, and Restarting Role Instances.....	38
Decommissioning Role Instances.....	39
Recommissioning Role Instances.....	39
Deleting Role Instances.....	40
Configuring Roles to Use a Custom Garbage Collection Parameter.....	40
Role Groups.....	41
Creating a Role Group.....	41
Managing Role Groups.....	42
Default User Roles.....	42
Managing Cloudera Runtime Services.....	43
Adding a Service.....	43
Comparing Configurations for a Service Between Clusters.....	44
Starting a Cloudera Runtime Service on All Hosts.....	45
Stopping a Cloudera Runtime Service on All Hosts.....	46
Restarting a Cloudera Runtime Service.....	46
Rolling Restart.....	47
Aborting a Pending Command.....	50
Deleting Services.....	50
Renaming a Service.....	51
Configuring Maximum File Descriptors.....	51
Extending Cloudera Manager.....	52
Add-on Services.....	52
Core Configuration Service.....	54
Managing Cloudera Manager.....	54
Automatic Logout.....	54
Starting, Stopping, and Restarting the Cloudera Manager Server.....	56
Configuring Cloudera Manager.....	56
Configuring Cloudera Manager Server Ports.....	57
Configuring Network Settings for a Proxy Server.....	57
Moving the Cloudera Manager Server to a New Host.....	58
Migrating from the Cloudera Manager Embedded PostgreSQL Database Server to an External PostgreSQL Database.....	59
Step 1: Identify Roles that Use the Embedded Database Server.....	59
Step 2: Migrate Databases from the Embedded Database Server to the External PostgreSQL Database Server.....	61
Migrating from the Cloudera Manager External PostgreSQL Database Server to a MySQL/Oracle Database Server.....	65
Prerequisites.....	65
Migrate from the Cloudera Manager External PostgreSQL Database Server to a MySQL/Oracle Database Server.....	66
Managing Cloudera Manager Server Logs.....	68
Viewing the Cloudera Manager Server Logs.....	68
Setting the Cloudera Manager Server Log Location.....	68
Cloudera Manager Agents.....	69
Starting, Stopping, and Restarting Cloudera Manager Agents.....	69
Configuring Cloudera Manager Agents.....	71
Managing the Cloudera Manager Agent Logs.....	73
Overview of Parcels.....	74

Advantages of Parcels.....	74
Parcel Life Cycle.....	75
Parcel Locations.....	75
Managing Parcels.....	76
Viewing Parcel Usage.....	79
Parcel Configuration Settings.....	81
Managing Licenses.....	83
Accessing the License Page.....	83
Ending a CDP Private Cloud Base Trial.....	83
Upgrading from a CDP Private Cloud Base Trial to CDP Private Cloud Base.....	84
Renewing a License.....	84
Default User Roles.....	84
Other Cloudera Manager Tasks and Settings.....	85
Cloudera Management Service.....	86
Starting the Cloudera Management Service.....	88
Stopping the Cloudera Management Service.....	88
Restarting the Cloudera Management Service.....	88
Starting and Stopping Cloudera Management Service Roles.....	89
Configuring Management Service Database Limits.....	89
Performance Management.....	90
Optimizing Performance in Cloudera Runtime.....	90
Disable the tuned Service.....	90
Disabling Transparent Hugepages (THP).....	91
Setting the vm.swappiness Linux Kernel Parameter.....	92
Improving Performance in Shuffle Handler and IFile Reader.....	92
Tips and Best Practices for Jobs.....	93
Decrease Reserve Space.....	93
Choosing and Configuring Data Compression.....	93
Resource Management.....	94
Static Service Pools.....	95
Enabling and Configuring Static Service Pools.....	96
Disabling Static Service Pools.....	96
Linux Control Groups (cgroups).....	97
Data Storage for Monitoring Data.....	100
Configuring Service Monitor Data Storage.....	100
Configuring Host Monitor Data Storage.....	101
Viewing Host and Service Monitor Data Storage.....	101
Data Granularity and Time-Series Metric Data.....	101
Moving Monitoring Data on an Active Cluster.....	102
Host Monitor and Service Monitor Memory Configuration.....	102
Accessing Storage Using Amazon S3.....	103
Referencing S3 Credentials for YARN, MapReduce, or Spark Clients.....	104

Referencing Amazon S3 in URIs.....	105
Using Fast Upload with Amazon S3.....	106
Enabling Fast Upload using Cloudera Manager.....	106
Configuring and Managing S3Guard.....	107
Configuring S3Guard for Cluster Access to S3.....	107
Editing the S3Guard Configuration.....	108
Running the Prune Command Using Cloudera Manager Admin Console.....	109
Running the Prune Command Using the Cloudera Manager API.....	109
How to Configure a MapReduce Job to Access S3 with an HDFS Credstore.....	110
Importing Data into Amazon S3 Using Sqoop.....	111
Authentication.....	111
Sqoop Import into Amazon S3.....	113
S3Guard with Sqoop.....	115

Accessing Storage Using Microsoft ADLS Gen 2..... 116

Configuring OAuth in Data Hub.....	116
Configuring OAuth with core-site.xml.....	116
Configuring OAuth with the Hadoop CredentialProvider.....	117
Configuring Native TLS Acceleration.....	117
Importing Data into Microsoft Azure Data Lake Store (Gen1 and Gen2) Using Sqoop.....	118
Prerequisites.....	119
Authentication.....	119
Sqoop Import into ADLS.....	119

Accessing the Cloudera Manager Admin Console

After you create a Data Hub cluster using the Cloudera Management Console, you can access the Cloudera Manager Admin Console to manage, configure, and monitor the cluster and its Cloudera Runtime services.

About this task

To access the Cloudera Manager Admin Console:

Procedure

1. Open the Cloudera Management Console in a Web browser using the following link: `http://<cm-server-url>/`
2. Click the Data Hub Clusters service.
3. Click the name of the Data Hub cluster you want to manage.
The cluster details page displays.
4. Click the URL for Cloudera Manager.

Results


The Cloudera Manager Admin Console opens in a new browser tab. You do not need to login to the Cloudera Manager Admin Console.

Starting, Stopping, Refreshing, and Restarting a Cluster

Minimum Required Role: [Operator](#) (also provided by Configurator, Cluster Administrator, Limited Cluster Administrator, Full Administrator)

Complete the steps below to start, stop, refresh, and restart a cluster.

Starting a Cluster


1. On the HomeStatus tab, click  to the right of the cluster name and select Start.
2. Click Start that appears in the next screen to confirm. The Command Details window shows the progress of starting services.

When All services successfully started appears, the task is complete and you can close the Command Details window.



Note: The cluster-level Start action starts only Cloudera Runtime and other product services (Impala, Cloudera Search). It does not start the Cloudera Management Service. You must start the Cloudera Management Service separately if it is not already running.

Stopping a Cluster

1. On the HomeStatus tab, click  to the right of the cluster name and select Stop.
2. Click Stop in the confirmation screen. The Command Details window shows the progress of stopping services.

When All services successfully stopped appears, the task is complete and you can close the Command Details window.




Note: The cluster-level Stop action does not stop the Cloudera Management Service. You must stop the Cloudera Management Service separately.

Refreshing a Cluster

Runs a cluster refresh action to bring the configuration up to date without restarting all services. For example, certain masters (for example NameNode and ResourceManager) have some configuration files (for example, fair-scheduler.xml, mapred_hosts_allow.txt, topology.map) that can be refreshed. If anything changes in those files then a refresh can be used to update them in the master.




Important: If you have changed a configuration property that requires a redeployment of the client configurations, note that refreshing or restarting a cluster does not automatically re-deploy the client

configurations. A service or cluster displays a staleness icon () next to the cluster or service name that indicates that you must redeploy the client configuration. Click the icon to open the [Stale Configurations](#) page and follow the prompts to refresh the cluster and redeploy the client configuration.









Alternatively, after a restart is completed, you can select Deploy Client Configuration from the Actions menu for either a service or cluster.


Here is a summary of the operations performed in a refresh action:

 Refresh Cluster	Cluster 1	Finished	Mar 19, 2014 11:31:55 AM PDT	Mar 19, 2014 11:32:09 AM PDT
Successfully refreshed roles in the cluster.				

Command Progress

Completed 4 of 4 steps.


-  Run 1 steps in parallel
Successfully refreshed datanode allow/exclude lists.
[Details](#) 
-  Run 1 steps in parallel
Successfully refreshed ResourceManager.
[Details](#) 
-  Run 3 steps in parallel
Successfully refreshed NodeManager.
[Details](#) 
-  Run 3 steps in parallel
Refreshed Impala Daemon's Pools configuration and ACLs successfully.
[Details](#) 

To refresh a cluster, in the HomeStatus tab, click  to the right of the cluster name and select Refresh Cluster.


Restarting a Cluster



Important: If you have changed a configuration property that requires a redeployment of the client configurations, note that refreshing or restarting a cluster does not automatically re-deploy the client

configurations. A service or cluster displays a staleness icon () next to the cluster or service name that indicates that you must redeploy the client configuration. Click the icon to open the [Stale Configurations](#) page and follow the prompts to refresh the cluster and redeploy the client configuration.

Alternatively, after a restart is completed, you can select Deploy Client Configuration from the Actions menu for either a service or cluster.

1. On the HomeStatus tab, click  to the right of the cluster name and select Restart.

- Click Restart that appears in the next screen to confirm. If you have enabled high availability for HDFS, you can choose Rolling Restart instead to minimize cluster downtime. The Command Details window shows the progress of stopping services.

When All services successfully started appears, the task is complete and you can close the Command Details window.

Pausing a Cluster in AWS

Operator (also provided by Configurator, Cluster Administrator, Full Administrator)

If all data for a cluster is stored on EBS volumes, you can pause the cluster and stop your AWS EC2 instances during periods when the cluster will not be used. The cluster will not be available while paused and can't be used to ingest or process data, but you won't be billed by Amazon for the stopped EC2 instances. Provisioned EBS storage volumes will continue to accrue charges.



Important: Pausing a cluster requires using EBS volumes for all storage, both on management and worker nodes. Data stored on ephemeral disks will be lost after EC2 instances are stopped.

Shutting Down and Starting Up the Cluster

To pause an AWS cluster, follow the shutdown procedure. To restart the cluster after a pause, follow the startup procedure.


Operator (also provided by Configurator, Cluster Administrator, Full Administrator)


In the shutdown and startup procedures below, some steps are performed in the AWS console and some are performed in Cloudera Manager:

- For AWS actions, use one of the following interfaces:
 - AWS console
 - AWS CLI
 - AWS API
- For cluster actions, use one of the following interfaces:
 - The Cloudera Manager web UI
 - The Cloudera API start and stop commands

Shutdown procedure



To pause the cluster, complete the following steps:

- Navigate to the Cloudera Manager web UI.
- Stop the cluster.
 - On the HomeStatus tab, click  to the right of the cluster name and select Stop.
 - Click Stop in the confirmation screen. The Command Details window shows the progress of stopping services.

When All services successfully stopped appears, the task is complete and you can close the Command Details window.
- Stop the Cloudera Management Service.
 - On the HomeStatus tab, click  to the right of the service name and select Stop.
 - Click Stop in the next screen to confirm. When you see a Finished status, the service has stopped.
- In AWS, stop all cluster EC2 instances, including the Cloudera Manager host .

Startup procedure

To restart the cluster after a pause, the steps are reversed:

1. In AWS, start all cluster EC2 instances.
2. Navigate to the Cloudera Manager UI.
3. Start the Cloudera Management Service.
 - a. On the HomeStatus tab, click  to the right of the service name and select Start.
 - b. Click Start that appears in the next screen to confirm. When you see a Finished status, the service has started.
4. Start the cluster.
 - a. On the HomeStatus tab, click  to the right of the cluster name and select Start.
 - b. Click Start that appears in the next screen to confirm. The Command Details window shows the progress of starting services.

When All services successfully started appears, the task is complete and you can close the Command Details window.

Considerations after Restart

Since the cluster was completely stopped before stopping the EC2 instances, the cluster should be healthy upon restart and ready for use. You should be aware of the following about the restarted cluster:


- After starting the EC2 instances, Cloudera Manager and its agents will be running but the cluster will be stopped. There will be gaps in Cloudera Manager's time-based metrics and charts.
- EC2 instances retain their internal IP address and hostname for their lifetime, so no reconfiguration of CDH or Runtime is required after restart. The public IP and DNS hostnames, however, will be different. Elastic IPs can be configured to remain associated with a stopped instance at additional cost, but it isn't necessary to maintain proper cluster operation.

Renaming a Cluster

About this task

Minimum Required Role: [Full Administrator](#). This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

1. On the HomeStatus tab, click  to the right of the cluster name and select Rename Cluster.
2. Type the new cluster name and click Rename Cluster.

Managing Hosts

How to use Cloudera Manager to configure and manage the hosts in your clusters.

Viewing Host Status

You can view summary information about the hosts managed by Cloudera Manager. You can view information for all hosts, the hosts in a cluster, or individual hosts.

Viewing All Hosts

To display summary information about all the hosts managed by Cloudera Manager, click **All Hosts** in the left menu. The **All Hosts** page displays with a list of all the hosts managed by Cloudera Manager.

The list of hosts shows the overall status of the Cloudera Manager-managed hosts in your cluster.

- The information provided varies depending on which columns are selected. To change the columns, click the Columns: *n* Selected drop-down and select the checkboxes next to the columns to display.
- Click ▶ to the left of the number of roles to list all the role instances running on that host.
- Filter the hosts list by entering search terms (hostname, IP address, or role) in the search box separated by commas or spaces. Use quotes for exact matches (for example, strings that contain spaces, such as a role name) and brackets to search for ranges. Hosts that match any of the search terms are displayed. For example:

```
hostname[1-3], hostname8 hostname9, "hostname.example.com"
hostname.example.com "HDFS DataNode"
```

- You can also search for hosts by selecting a value from the facets in the Filters section at the left of the page. Click the Filters toggle to show or hide the Filters section.
- If the agent heartbeat and health status properties are configured as follows:
 - Send Agent heartbeat every *x*
 - Set health status to Concerning if the Agent heartbeats fail *y*
 - Set health status to Bad if the Agent heartbeats fail *z*

The value *v* for a host's Last Heartbeat facet is computed as follows:

- $v < x * y = \text{Good}$
- $v \geq x * y$ and $\leq x * z = \text{Concerning}$
- $v \geq x * z = \text{Bad}$

Viewing the Hosts in a Cluster

Do one of the following:

- Select Clusters *Cluster name* Hosts .
- In the Home screen, click  **Hosts** in a full form cluster table.

The **All Hosts** page displays with a list of the hosts filtered by the cluster name.

Viewing Individual Hosts

You can view detailed information about an individual host—resources (CPU/memory/storage) used and available, which processes it is running, details about the host agent, and much more—by clicking a host link on the **All Hosts** page.

Adding a Host to a Cluster

Steps to add hosts to a cluster.

Minimum Required Role: **Full Administrator**. This feature is not available when using Cloudera Manager to manage Data Hub clusters.

You can add one or more hosts to your cluster using the Add Hosts wizard, which installs the Oracle JDK, Cloudera Runtime, and Cloudera Manager Agent software. After the software is installed and the Cloudera Manager Agent is started, the Agent connects to the Cloudera Manager Server and you can use the Cloudera Manager Admin Console to manage and monitor Cloudera Runtime on the new host.

The Add Hosts wizard does not create roles on the new host; once you have successfully added the host(s) you can either add roles, one service at a time, or apply a host template, which can define role configurations for multiple roles.

**Important:**

- Unqualified hostnames (short names) must be unique in a Cloudera Manager instance. For example, you cannot have both *host01.example.com* and *host01.standby.example.com* managed by the same Cloudera Manager Server.
- All hosts in a single cluster must be running the same version of Cloudera Runtime.
- When you add a new host, you must install the same version of Cloudera Runtime to enable the new host to work with the other hosts in the cluster. The installation wizard lets you select the version of Cloudera Runtime to install, and you can choose a custom repository to ensure that the version you install matches the version on the other hosts.
- If you are managing multiple clusters, select the version of Cloudera Runtime that matches the version in use on the cluster where you plan to add the new host.
- When you add a new host, the following occurs:
 - YARN topology.map is updated to include the new host
 - Any service that includes topology.map in its configuration—Flume, Hive, Hue, Oozie, Solr, Spark, YARN—is marked stale

At a convenient point after adding the host you should restart the stale services to pick up the new configuration.

Use one of the following methods to add a new host:

Using the Add Hosts Wizard to Add Hosts

You can use the Add Hosts wizard to install Cloudera Runtime, Impala, and the Cloudera Manager Agent on a host.

Disable TLS Encryption or Authentication

If you have enabled TLS encryption or authentication for the Cloudera Manager Agents, you must disable both of them before starting the Add Hosts wizard. Otherwise, skip to the next step.

If you perform this step, then skip step 2 ([Alternate Method of Installing Cloudera Manager Agent without Disabling TLS](#)). If you skip step 1 and perform step 2, then continue to step 3 ([Using the Add Hosts Wizard to Add Hosts](#)).



Important: This step temporarily puts the existing cluster hosts in an unmanageable state; they are still configured to use TLS and so cannot communicate with the Cloudera Manager Server. Roles on these hosts continue to operate normally, but Cloudera Manager is unable to detect errors and issues in the cluster and reports all hosts as being in bad health. To work around this issue, you can manually install the Cloudera Manager Agent on the new host. See [Alternate Method of Installing Cloudera Manager Agent without Disabling TLS](#) on page 13.

1. From the Administration tab, select Settings.
2. Select the Security category.
3. Disable TLS by clearing the following options: Use TLS Encryption for Agents, and Use TLS Authentication of Agents to Server.
4. Click Save Changes to save the settings.
5. Log in to the Cloudera Manager Server host.
6. Restart the Cloudera Manager Server.

RHEL 7, SLES 12, Debian 8, Ubuntu 16.04 and higher

```
sudo systemctl restart cloudera-scm-server
```

RHEL 5 or 6, SLES 11, Debian 6 or 7, Ubuntu 12.04 or 14.04

```
sudo service cloudera-scm-server restart
```

Alternate Method of Installing Cloudera Manager Agent without Disabling TLS

If you have TLS encryption or authentication enabled in your cluster, you must either disable TLS during the installation, or install the Cloudera Manager Agent manually using the following procedure:

1. Copy the repository configuration file from an existing host in the cluster to the new host. For example:

OS	Command
RHEL	<pre>sudo scp mynode.example.com:/etc/yum.repos.d/cloudera-manager.repo /etc/yum.repos.d/cloudera-manager.repo</pre>
SLES	<pre>sudo scp mynode.example.com:/etc/zypp/zypper.conf/cloudera-cm.repo /etc/zypp/zypper.conf/cloudera-cm.repo</pre>
Ubuntu or Debian	<pre>sudo scp mynode.example.com:/etc/apt/sources.list.d/cloudera.list /etc/apt/sources.list.d/cloudera.list</pre>

2. Remove cached package lists and other transient data by running the following command:

OS	Command
RHEL	<pre>sudo yum clean all</pre>
SLES	<pre>sudo zypper clean --all</pre>
Ubuntu or Debian	<pre>sudo apt-get clean</pre>

3. Install the JDK package from the Cloudera Manager repository. Install the same version as is used on other cluster hosts. Only JDK 1.8 is supported:

Table 1: Oracle JDK 1.8

OS	Command
RHEL	<pre>sudo yum install jdk1.8.0_144-cloudera</pre>
SLES	<pre>sudo zypper install jdk1.8.0_144-cloudera</pre>
Ubuntu or Debian	<pre>sudo apt-get install jdk1.8.0_144-cloudera</pre>

Open JDK

RHEL

OpenJDK 8

```
sudo yum install java-1.8.0-openjdk-devel
```

OpenJDK 11

```
su -c yum install java-11-openjdk-devel
```

Ubuntu

OpenJDK 8

```
sudo apt-get install openjdk-8-jdk
```

OpenJDK 11

```
sudo apt-get install openjdk-11-jdk
```

SLES

OpenJDK 8

```
sudo zypper install java-1_8_0-openjdk-devel
```

OpenJDK 11

```
sudo zypper install java-11-openjdk-devel
```

4. Set up the TLS certificates using the same procedure that was used to set them up on other cluster hosts. See [Configuring TLS Encryption for Cloudera Manager Using Auto-TLS](#). If you have set up a custom truststore, copy that file from an existing host to the same location on the new host.

5. Install the Cloudera Manager Agent:

OS	Command
RHEL	<code>sudo yum install cloudera-manager-agent</code>
SLES	<code>sudo zypper install cloudera-manager-agent</code>
Ubuntu or Debian	<code>sudo apt-get install cloudera-manager-agent</code>

6. Copy the Cloudera Manager Agent configuration file from an existing cluster host that is already configured for TLS to the same location on the new host. For example:

```
sudo scp mynode.example.com:/etc/cloudera-scm-agent/config.ini /etc/cloudera-scm-agent/config.ini
```

7. Create and secure the file containing the password used to protect the private key of the Agent:

- a. Use a text editor to create a file called `agentkey.pw` that contains the password. Save the file in the `/etc/cloudera-scm-agent` directory.
- b. Change ownership of the file to root:

```
sudo chown root:root /etc/cloudera-scm-agent/agentkey.pw
```

- c. Change the permissions of the file:

```
sudo chmod 440 /etc/cloudera-scm-agent/agentkey.pw
```

8. Start the Agent on the new host:

```
sudo service cloudera-scm-agent start
```

9. Log in to Cloudera Manager and go to HostsAll Hosts page and verify that the new host is recognized by Cloudera Manager.

Add Hosts to an Existing Cluster

1. Click the Hosts tab.
2. Click the Add Hosts button.
3. Select Add hosts to cluster.
4. If the cluster uses Kerberos authentication, ensure that the Kerberos packages are installed on the new hosts. If necessary, use the package commands provided on the Add Hosts screen to install these packages.
5. Select the cluster where you want to add the host from the drop-down list.

6. Click Continue.

The Specify Hosts page displays. You can either add a new host to the cluster, or add an existing managed host to the cluster.

Do one of the following:

- Add a new host:
 - a. On the Specify Hosts page, enter a host name or pattern (click "using patterns" for more information) to search for new hosts to add to the cluster.

A list of matching hosts displays.
 - b. Select the hosts that you want to add.
 - c. Click Continue.
 - d. Select the Repository Location where Cloudera Manager can find the software to install on the new hosts. Select Public Cloudera Repository or Custom Repository and enter the URL of a custom repository available on your local network.
 - e. Click Continue.
 - f. Follow the instructions in the wizard to install the Oracle JDK.
 - g. Enter Login Credentials:
 1. Select root for the root account, or select Another user and enter the username for an account that has password-less sudo privileges.
 2. Select an authentication method:
 - If you choose password authentication, enter and confirm the password.
 - If you choose public-key authentication, provide a passphrase and path to the required key files.

You can modify the default SSH port if necessary.
 3. Specify the maximum number of host installations to run at once. The default and recommended value is 10. You can adjust this based on your network capacity.
 4. Click Continue.

The Install Agents page displays and Cloudera Manager installs the Agent software on the new hosts.
 5. When the agent installation finishes, click Continue.
 - Add an existing managed host:
 - a. Click the Currently Managed Hosts tab.

A list of hosts previously added to Cloudera Manager displays.
 - b. Select the hosts that you want to add to the cluster.
 - c. Click Continue.
- 7.** Cloudera Manager begins to install the Cloudera Runtime parcels.
- 8.** When the parcel installation finishes, click Continue.
- 9.** The Host Inspector runs and displays any problems with the hosts. Correct the problems before continuing.
- 10.** After correcting any problems, click Continue.
- 11.** To add role instances to the hosts now or select None to add them later.
- 12.** To add roles now:- a. Select an existing host template, or create a new one.
- b. To create a new host template, click the Create... button. The Create New Host Template screen opens.. See [Host Templates](#) on page 20 for details on how you select the role groups that define the roles that should run on a host. After you have created the template, it will appear in the list of host templates from which you can choose.
- c. Select the host template you want to use.
- d. By default Cloudera Manager will automatically start the roles specified in the host template on your newly added hosts. To prevent this, uncheck the option to start the newly-created roles.

13. When the wizard is finished, you can verify the Agent is connecting properly with the Cloudera Manager Server by clicking the Hosts tab and checking the health status for the new host. If the Health Status is Good and the value for the Last Heartbeat is recent, then the Agent is connecting properly with the Cloudera Manager Server.

If you did not specify a host template during the Add Hosts wizard, then no roles will be present on your new hosts until you add them. You can do this by adding individual roles under the Instances tab for a specific service, or by using a host template. See [Adding a Role Instance](#) for information about adding roles for a specific service. See [Host Templates](#) to create a host template that specifies a set of roles (from different services) that should run on a host.

Add New Hosts To Cloudera Manager

This option allows you to add a host, but without adding them to a specific cluster. Later, you can use these hosts to create new clusters or expand existing clusters.

1. Click the Hosts tab.
2. Click the Add Hosts button.
3. Select Add hosts to Cloudera Manager
4. Specify the hosts to add:
 - a. On the Specify Hosts page, enter a host name or pattern (click "using patterns" for more information) to search for new hosts to add to the cluster.

A list of matching hosts displays.
 - b. Select the hosts that you want to add.
 - c. Click Continue.
 - d. Select the Repository Location where Cloudera Manager can find the software to install on the new hosts. Select Public Cloudera Repository or Custom Repository and enter the URL of a custom repository available on your local network.
 - e. Click Continue.
 - f. Follow the instructions in the wizard to install the Oracle JDK.
 - g. Enter Login Credentials:
 1. Select root for the root account, or select Another user and enter the username for an account that has password-less sudo privileges.
 2. Select an authentication method:
 - If you choose password authentication, enter and confirm the password.
 - If you choose public-key authentication, provide a passphrase and path to the required key files.

You can modify the default SSH port if necessary.
 3. Specify the maximum number of host installations to run at once. The default and recommended value is 10. You can adjust this based on your network capacity.
 4. Click Continue.

The Install Agents page displays and Cloudera Manager installs the Agent software on the new hosts.
 5. When the agent installation finishes, click Continue.
5. If the cluster uses Kerberos authentication, ensure that the Kerberos packages are installed on the new hosts. If necessary, use the package commands provided on the Add Hosts screen to install these packages.
6. Select the Repository Location where Cloudera Manager can find the software to install on the new hosts. Select Public Cloudera Repository or Custom Repository and enter the URL of a custom repository available on your local network.
7. Follow the instructions in the wizard to install the Oracle JDK.

8. Enter Login Credentials:

- a. Select root for the root account, or select Another user and enter the username for an account that has password-less sudo privileges.
- b. Select an authentication method:
 - If you choose password authentication, enter and confirm the password.
 - If you choose public-key authentication, provide a passphrase and path to the required key files.

You can modify the default SSH port if necessary.

- c. Specify the maximum number of host installations to run at once. The default and recommended value is 10. You can adjust this based on your network capacity.
- d. Click Continue.

The Install Agents page displays and Cloudera Manager installs the Agent software on the new hosts.

- e. When the agent installation finishes, click Continue.

9. The Host Inspector runs and displays any problems with the hosts. Correct the problems before continuing.

10. After correcting any problems, click Continue.

Enable TLS Encryption or Authentication

If you previously enabled TLS security on your cluster, you must re-enable the TLS options on the Administration page and also configure TLS on each new host after using the Add Hosts wizard. Otherwise, you can ignore this step. For instructions, see [Configuring TLS Encryption for Cloudera Manager and Cloudera Runtime Using Auto-TLS](#).

Enable TLS/SSL for cluster Components

If you have previously enabled TLS/SSL on your cluster, and you plan to start these roles on this new host, make sure you install a new host certificate to be configured from the same path and naming convention as the rest of your hosts. Since the new host and the roles configured on it are inheriting their configuration from the previous host, ensure that the keystore or truststore passwords and locations are the same on the new host. For instructions on configuring TLS/SSL, see [Configuring TLS Encryption for Cloudera Manager and Cloudera Runtime Using Auto-TLS](#).

Enable Kerberos

If you have previously enabled Kerberos on your cluster:

1. Install the packages required to kinit on the new host (see the list in [Enabling Kerberos Authentication for Cloudera Runtime](#)).
2. If you have set up Cloudera Manager to manage krb5.conf, it will automatically deploy the file on the new host. Note that Cloudera Manager will deploy krb5.conf only if you use the Kerberos wizard. If you have used the API, you will need to manually perform the commands that the wizard calls.

If Cloudera Manager does not manage krb5.conf, you must manually update the file at /etc/krb5.conf.

Adding a Host by Installing the Packages Using Your Own Method

If you used a different mechanism to install the JDK, Cloudera Runtime, and Cloudera Manager Agent packages, you can use that same mechanism to install the JDK, Cloudera Runtime, Cloudera Manager Agent packages and then start the Cloudera Manager Agent.

1. Install the Oracle JDK, Cloudera Runtime, and Cloudera Manager Agent packages using your own method. For instructions on installing these packages, see [Install Cloudera Manager Server](#).
2. After installation is complete, start the Cloudera Manager Agent. For instructions, see [Starting, Stopping, and Restarting Cloudera Manager Agents](#) on page 69.
3. After the Agent is started, you can verify the Agent is connecting properly with the Cloudera Manager Server by clicking the Hosts tab and checking the health status for the new host. If the Health Status is Good and the value for the Last Heartbeat is recent, then the Agent is connecting properly with the Cloudera Manager Server.

4. If you have enabled TLS security on your cluster, you must enable and configure TLS on each new host. Otherwise, ignore this step.
 - a. Enable and configure TLS on each new host by specifying 1 for the `use_tls` property in the `/etc/cloudera-scm-agent/config.ini` configuration file.
 - b. Configure TLS security on the new hosts by following the instructions in [Configuring TLS Encryption for Cloudera Manager and Cloudera Runtime Using Auto-TLS](#).
5. If you have previously enabled TLS/SSL on your cluster, and you plan to start these roles on this new host, make sure you install a new host certificate to be configured from the same path and naming convention as the rest of your hosts. Since the new host and the roles configured on it are inheriting their configuration from the previous host, ensure that the keystore or truststore passwords and locations are the same on the new host. For instructions on configuring TLS/SSL, see [Configuring TLS Encryption for Cloudera Manager and Cloudera Runtime Using Auto-TLS](#).
6. If you have previously enabled Kerberos on your cluster:
 - a. Install the packages required to kinit on the new host (see the list in [Enabling Kerberos Authentication for Cloudera Runtime](#)).
 - b. If you have set up Cloudera Manager to manage `krb5.conf`, it will automatically deploy the file on the new host. Note that Cloudera Manager will deploy `krb5.conf` only if you use the Kerberos wizard. If you have used the API, you will need to manually perform the commands that the wizard calls.

If Cloudera Manager does not manage `krb5.conf`, you must manually update the file at `/etc/krb5.conf`.

Parcels

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

In the Parcels tab you can download, distribute, and activate available parcels to your cluster. You can use parcels to add new products to your cluster, or to upgrade products you already have installed.

Configuring Hosts

The Configuration tab lets you set properties related to parcels and to resource management, and also monitoring properties for the hosts under management.

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

The configuration settings you make here will affect all your managed hosts. You can also configure properties for individual hosts by clicking on the host in the **All Hosts** page, which will override the global properties set here).

To edit the default configuration properties for hosts, click the Configuration tab.

Related Information

[Modifying Configuration Properties Using Cloudera Manager](#)

Viewing Host Role Assignments

You can view the assignment of roles to hosts as follows:

1. In the left menu, click HostsRoles.
2. Click a cluster name or All Clusters.

Host Templates

The **Host Templates** page lets you create and manage host templates, which provide a way to specify a set of role configurations that should be applied to a host.

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Host templates let you designate a set of role groups that can be applied in a single operation to a host or a set of hosts. This significantly simplifies the process of configuring new hosts when you need to expand your cluster.



Important: A host template can only be applied on a host with a version of Cloudera Runtime that matches the Cloudera Runtime version running on the cluster to which the host template belongs.

You can create and manage host templates by clicking HostsHost Templates.

Templates are not required; Cloudera Manager assigns roles and role groups to the hosts of your cluster when you perform the initial cluster installation. However, if you want to add new hosts to your cluster, a host template can make this much easier.

If there are existing host templates, they are listed on the page, along with links to each role group included in the template.

If you are managing multiple clusters, you must create separate host templates for each cluster, as the templates specify role configurations specific to the roles in a single cluster. Existing host templates are listed under the cluster to which they apply.

- You can click a role group name to be taken to the Edit configuration page for that role group, where you can modify the role group settings.
- From the Actions menu associated with the template you can edit the template, clone it, or delete it.

Creating a Host Template

When you create a host template, you choose a name for the template and select appropriate role groups for each role.

About this task

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

1. Click HostsHost Templates.
2. From the **Host Templates** page, click Create.
The **Create New Host Template** pop-up window appears.
3. Type a name for the template.
4. For each role, select the appropriate role group. There may be multiple role groups for a given role type — you want to select the one with the configuration that meets your needs.
5. Click Create to create the host template.

Editing a Host Template

You can edit the name of a host template, in addition to any of the role group selections.

About this task

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

1. Click HostsHost Templates.
2. Pull down the Actions menu for the template you want to modify, and click Edit.
The **Edit Host Template** window appears. This page is identical to the Create New Host Template page. You can modify the template name or any of the role group selections.
3. Click OK when you have finished.

Applying a Host Template to a Host

You can use a host template to apply configurations for multiple roles in a single operation.

About this task

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

You can apply a template to a host that has no roles on it, or that has roles from the same services as those included in the host template. New roles specified in the template that do not already exist on the host will be added. A role on the host that is already a member of the role group specified in the template will be left unchanged. If a role on the host matches a role in the template, but is a member of a different role group, it will be moved to the role group specified by the template.

For example, suppose you have two role groups for a DataNode (DataNode Default Group and DataNode (1)). The host has a DataNode role that belongs to DataNode Default Group. If you apply a host template that specifies the DataNode (1) group, the role on the host will be moved from DataNode Default Group to DataNode (1).

However, if you have two instances of a service, such as MapReduce (for example, mr1 and mr2) and the host has a TaskTracker role from service mr2, you cannot apply a TaskTracker role from service mr1.

A host may have no roles on it if you have just added the host to your cluster, or if you decommissioned a managed host and removed its existing roles.

Also, the host must have the same version of CDH installed as is running on the cluster whose host templates you are applying.

If a host belongs to a different cluster than the one for which you created the host template, you can apply the host template if the "foreign" host either has no roles on it, or has only management roles on it. When you apply the host template, the host will then become a member of the cluster whose host template you applied. The following instructions assume you have already created the appropriate host template.


Procedure

1. Click HostsAll Hosts.
2. Select the host(s) to which you want to apply your host template.
3. From the Actions for Selected menu, select Apply Host Template.
4. In the pop-up window that appears, select the host template you want to apply.
5. Optionally you can have Cloudera Manager start the roles created per the host template. To enable this, check the box.
6. Click Confirm to initiate the action.

Hosts Disks Overview

How to view the status of all disks in a cluster.

In the left menu, click HostsDisks Overview to display an overview of the status of all disks in the deployment. The statistics exposed match or build on those in iostat, and are shown in a series of histograms that by default cover every physical disk in the system.

Adjust the endpoints of the time line to see the statistics for different time periods. Specify a filter in the box to limit the displayed data. For example, to see the disks for a single rack rack1, set the filter to: `logicalPartition = false` and `rackId = "rack1"` and click Filter. Click a histogram to drill down and identify outliers. Mouse over the graph and click  to display additional information about the chart.

Deleting Hosts

Minimum Required Role: [Full Administrator](#). This feature is not available when using Cloudera Manager to manage Data Hub clusters.

You can remove a host from a cluster in two ways:

- Delete the host entirely from Cloudera Manager.
- Remove a host from a cluster, but leave it available to other clusters managed by Cloudera Manager.

Both methods decommission the hosts, delete roles, and remove managed service software, but preserve data directories.

Deleting a Host from Cloudera Manager

To delete a host from Cloudera Manager, first decommission the host and then remove it.

About this task

Minimum Required Role: [Full Administrator](#). This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

1. In the Cloudera Manager Admin Console, go to Hosts All Hosts.
2. Select the hosts to delete.
3. Select Actions for SelectedHosts Decommission.
4. Stop the Agent on the host.
5. In the Cloudera Manager Admin Console, go to Hosts All Hosts.
6. Reselect the hosts you selected in Step 2.
7. Select Actions for SelectedRemove from Cloudera Manager.

Removing a Host From a Cluster

Removing a host from a cluster leaves the host managed by Cloudera Manager and preserves the Cloudera Management Service roles (such as the Events Server, Host Monitor, and so on).

About this task

Minimum Required Role: [Full Administrator](#). This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

1. In the Cloudera Manager Admin Console, click the Hosts tab.
2. Select the hosts to delete.
3. Select Actions for SelectedRemove From Cluster. The **Remove Hosts From Cluster** dialog box displays.
4. Leave the selections to decommission roles and skip removing the Cloudera Management Service roles. Click Confirm to proceed with removing the selected hosts.

Stopping All the Roles on a Host

You can stop all of the roles on a host from the **Hosts** page.

About this task

Minimum Required Role: **Operator** (also provided by Configurator, Cluster Administrator, Limited Cluster Administrator , Full Administrator)

Procedure

1. In the left menu, click ClustersHosts or HostsAll Hosts.
2. Select one or more hosts on which to stop all roles.
3. Select Actions for SelectedStop Roles on Hosts.

Starting All the Roles on a Host

You can start all the roles on a host from the **Hosts** page.

About this task

Minimum Required Role: **Operator** (also provided by Configurator, Cluster Administrator, Limited Cluster Administrator , Full Administrator)

Procedure

1. Click the Hosts tab.
2. Select one or more hosts on which to start all roles.
3. Select Actions for SelectedStart Roles on Hosts.

Changing Hostnames

After you have installed Cloudera Manager and created a cluster, you may need to update the names of the hosts running the Cloudera Manager Server or cluster services.

About this task

Minimum Required Role: **Full Administrator**. This feature is not available when using Cloudera Manager to manage Data Hub clusters.



Important:

- The process described here requires Cloudera Manager and cluster downtime.
- If any user-created scripts reference specific hostnames, those must also be updated.
- Due to the length and complexity of the following procedure, changing cluster hostnames is not recommended by Cloudera.

To update a deployment with new hostnames, follow these steps:

Procedure

1. Verify if TLS/SSL certificates have been issued for any of the services and make sure to create new TLS/SSL certificates in advance for services protected by TLS/SSL.

2. Export the Cloudera Manager configuration using one of the following methods:

- Open a browser and go to this URL `http://cm_hostname:7180/api/api_version/cm/deployment`. Save the displayed configuration.
- From terminal type:

```
$ curl -u admin:admin http://cm_hostname:7180/api/api_version/cm/deployment > cme-cm-export.json
```

If Cloudera Manager SSL is in use, specify the `-k` switch and the port number as 7183:

```
$ curl -k -u admin:admin https://cm_hostname:7183/api/api_version/cm/deployment > cme-cm-export.json
```

where `cm_hostname` is the name of the Cloudera Manager host and `api_version` is the correct version of the API for the version of Cloudera Manager you are using. For example, `http://tcdn5-1.ent.cloudera.com:7180/api/v40/cm/deployment`.

3. Stop all services on the cluster.

4. Stop the Cloudera Management Service.

5. Stop the Cloudera Manager Server.

6. Stop the Cloudera Manager Agents on the hosts that you want to change the hostname of.

7. Back up the Cloudera Manager Server database using `mysqldump`, `pg_dump`, or another preferred backup utility. Store the backup in a safe location.

8. Update names and principals:

- a) Update the target hosts using standard per-OS/name service methods (`/etc/hosts`, `dns`, `/etc/sysconfig/network`, `hostname`, and so on). Ensure that you remove the old hostname.

- b) If you are changing the hostname of the host running Cloudera Manager Server do the following:

1. Change the hostname per Step 8.a.
2. Update the Cloudera Manager hostname in `/etc/cloudera-scm-agent/config.ini` on all Agents.

- c) If the cluster is configured for Kerberos security, do the following:

1. Remove the old hostname cluster principals.

- If you are using an MIT KDC, remove old hostname cluster service principals from the KDC database using one of the following:

- Use the `delprinc` command within `kadmin.local` interactive shell.

OR

- From the command line:

```
kadmin.local -q "listprincs" | grep -E "(HTTP|hbase|hdfs|hive|ht
tpfs|hue|impala|mapred|solr|oozie|yarn|zookeeper) [^/]*/[^/]*@" >
cluster-princ.txt
```

Open `cluster-princ.txt` and remove any noncluster service principal entries. Make sure that the default `krbtgt` and other principals you created, or that were created by Kerberos by default, are not removed by running the following: `for i in `cat cluster-princ.txt`; do yes yes | kadmin.local -q "delprinc $i"; done.`

- For an Active Directory KDC, an AD administrator must manually delete the principals for the old hostname from Active Directory.

2. Start the Cloudera Manager database and Cloudera Manager Server.

3. Start the Cloudera Manager Agents on the newly renamed hosts. The Agents should show a current heartbeat in Cloudera Manager.

4. Within the Cloudera Manager Admin Console click the Hosts tab.

5. Select the checkbox next to the host with the new name.

6. Select ActionsRegenerate Keytab.

9. If one of the hosts that was renamed has a NameNode configured with high availability and automatic failover enabled, reconfigure the ZooKeeper Failover Controller znodes to reflect the new hostname.

a) Start ZooKeeper Servers.



Warning: All other services, and most importantly HDFS, and the ZooKeeper Failover Controller (FC) role within the HDFS, should not be running.

b) On one of the hosts that has a ZooKeeper Server role, run zookeeper-client.

1. If the cluster is configured for Kerberos security, configure ZooKeeper authorization as follows:

- Go to the HDFS service.
- Click the Instances tab.
- Click the Failover Controller role.
- Click the Process tab.
- In the Configuration Files column of the hdfs/hdfs.sh ["zkfc"] program, expand Show.
- Inspect core-site.xml in the displayed list of files and determine the value of the ha.zookeeper.auth property, which will be something like: digest:hdfs-fcs:TEbW2bgoODa96rO3ZTn7ND5fSOGx0h. The part after digest:hdfs-fcs: is the password (in the example it is TEbW2bgoODa96rO3ZTn7ND5fSOGx0h)
- Run the addauth command with the password:

```
addauth digest hdfs-fcs:TEbW2bgoODa96rO3ZTn7ND5fSOGx0h
```

2. Verify that the HA znode exists: ls /hadoop-ha.

3. Delete the HDFS znode: rmr /hadoop-ha/nameservice1.

4. If you are not running JobTracker in a high availability configuration, delete the HA znode: rmr /hadoop-ha.

c) In the Cloudera Manager Admin Console, go to the HDFS service.

d) Click the Instances tab.

e) Select ActionsInitialize High Availability State in ZooKeeper....

10. Update the Hive metastore:

a) Back up the Hive metastore database.

b) In the Cloudera Manager Admin Console, go to the Hive service.

c) Select ActionsUpdate Hive Metastore NameNodes and confirm the command.

11. Update the Database Hostname property for each of the cluster roles for which a database is located on the host being renamed. This is required for both Cloudera Management Service roles (Reports Manager, Activity Monitor, Navigator Audit and Metadata Server) and for cluster services such as Hue, Hive, and so on.

12. Start all cluster services.

13. Start the Cloudera Management Service.

14. Deploy client configurations.

Moving a Host Between Clusters

To move a host between clusters, you must first decommission the host, remove roles from the host, and complete other tasks.

About this task

Minimum Required Role: [Full Administrator](#). This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

1. Decommission the host.

2. Remove all roles from the host (except for the Cloudera Manager management roles).
3. Remove the host from the cluster but leave it available to Cloudera Manager.
4. Add the host to the new cluster.
5. Add roles to the host (optionally using one of the host templates associated with the new cluster).

Using Upgrade Domains to manage rolling restarts

Upgrade Domains allow to group cluster hosts for optimal performance during restarts and upgrades.

Upgrade Domains enable faster cluster restarts, faster Cloudera Runtime upgrades, and seamless OS patching & hardware upgrades across large clusters. Upgrade Domains provide an alternative to the default HDFS block placement policy, distributing data across a set of hosts (potentially larger than a single rack) that Cloudera Manager can upgrade/restart at once without compromising service and data availability. When you select Upgrade Domains as the block placement policy, you also assign an Upgrade Domain group to each DataNode host. The NameNode uses these groups to distribute blocks when writing data, and to orchestrate rolling restarts and upgrades. This feature is useful for very large clusters, or for clusters where rolling restarts happen frequently.

For example, if HDFS is configured with the default replication factor of 3, the NameNode places the replica blocks on DataNode hosts in 3 different Upgrade Domains and on at least two different racks.



Note:

- Cloudera recommends that you assign an approximately equal number of DataNode hosts to each Upgrade Domain.
- The number of Upgrade Domains in a cluster should be greater than or equal to the HDFS Replication Factor. When you perform a rolling restart on a cluster, all hosts in an Upgrade Domain group will be restarted simultaneously, followed by the hosts in each remaining Upgrade Domain group.
- You should create a sufficient number of Upgrade Domains so that the cluster can still function adequately when all the hosts in a single Upgrade Domain are taken offline. The appropriate number of Upgrade Domains depends on the workloads and capacity of the cluster and may require tuning for optimal performance.
- To take advantage of the improved rolling restart performance, Upgrade Domain groups should not duplicate rack assignments. The number of hosts in an Upgrade Domain group should be larger than the number of hosts in a rack.

Configuring Upgrade Domains

Steps to configure Upgrade domains.

About this task

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

1. Configure the Upgrade Domain for all hosts:
 - a) Click Hosts Configuration.
 - b) Search for the Host Upgrade Domain parameter.
 - c) Click the Add Host Overrides link for this parameter.
 - d) Enter the Upgrade Domain group name in the New Override Value field.
 - e) Select the hosts for this Upgrade Domain.
 - f) Click Add.

2. Set the HDFS Block Replica Placement Policy:
 - a) Open the Cloudera Manager Admin Console.
 - b) Go to the HDFS service for the cluster.
 - c) Click the Configuration tab.
 - d) Search for the HDFS Block Replica Placement Policy configuration parameter.
 - e) Select Upgrade Domains.
 - f) Click Save Changes.The Upgrade Domain assigned to each host displays in the Upgrade Domain column on the All Hosts page.
3. Restart the HDFS service.

Adding or changing the Upgrade Domain for a single host

Steps to add or change the Upgrade Domain for a single host.

About this task

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

1. Click HostsAll hosts .
2. Click on the host name.
3. Click the Configuration tab.
4. Search for the Host Upgrade Domain parameter.
5. Enter the name of the Upgrade Domain group for this host.
6. Click Save Changes.
7. Restart the HDFS Service.

Putting all Hosts in an Upgrade Domain group into Maintenance Mode

Steps to put hosts in an Upgrade Domain into Maintenance Mode.

About this task

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

1. In Cloudera Manager, select the cluster where you want to decommission hosts.
2. Click HostsAll Hosts.
3. In the Filters section, click Upgrade Domain.
4. Select an Upgrade Domain.
The All Hosts list now displays only the hosts belonging to the Upgrade Domain.
5. Select all of the hosts.
6. Click Actions for SelectedBegin Maintenance (Suppress Alerts/Decommission).
The Begin Maintenance (Suppress Alerts/Decommission) dialog box opens. The role instances running on the hosts display at the top. You can also use this dialog box to decommission the host.
7. Select the Take DataNode offline option to put the hosts into Maintenance Mode.
In this mode, alerts from the hosts are suppressed until the host exits Maintenance Mode. The events, however, are still logged. Hosts that are currently in Maintenance Mode display the icon.

8. Click Begin Maintenance.

The Host Decommission Command dialog box opens and displays the progress of the command.

Specifying Racks for Hosts

To get maximum performance, it is important to configure CDH so that it knows the topology of your network. Network locations such as hosts and racks are represented in a tree, which reflects the network “distance” between locations. HDFS will use the network location to be able to place block replicas more intelligently to trade off performance and resilience.

About this task

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

When placing jobs on hosts, CDH prefers within-rack transfers (where there is more bandwidth available) to off-rack transfers; the MapReduce and YARN schedulers use network location to determine where the closest replica is as input to a map task. These computations are performed with the assistance of rack awareness scripts.

Cloudera Manager includes internal rack awareness scripts, but you must specify the racks where the hosts in your cluster are located. If your cluster contains more than 10 hosts, Cloudera recommends that you specify the rack for each host. HDFS, MapReduce, and YARN will automatically use the racks you specify.

Cloudera Manager supports nested rack specifications. For example, you could specify the rack `/rack3`, or `/group5/rack3` to indicate the third rack in the fifth group. All hosts in a cluster must have the same number of path components in their rack specifications.

Procedure

1. Click the Hosts tab.
2. Check the checkboxes next to the host(s) for a particular rack, such as all hosts for `/rack123`.
3. Click Actions for Selected (*n*)Assign Rack, where *n* is the number of selected hosts.
4. Enter a rack name or ID that starts with a slash `/`, such as `/rack123` or `/aisle1/rack123`, and then click Confirm.
5. Optionally restart affected services. Rack assignments are not automatically updated for running services.

Performing Maintenance on a Cluster Host

You can perform minor maintenance on cluster hosts by using Cloudera Manager to manage the host decommission and recommission process.

In this process, you can specify whether to suppress alerts from the decommissioned host and, for hosts running the DataNode role, you can specify whether or not to replicate under-replicated data blocks to other DataNodes to maintain the cluster's replication factor. This feature is useful when performing minor maintenance on cluster hosts, such as adding memory or changing network cards or cables where the maintenance window is expected to be short and the extra cluster resources consumed by replicating missing blocks is undesirable.

You can also place hosts into Maintenance Mode, which suppresses unneeded alerts during a maintenance window but does not decommission the hosts.

To perform host maintenance on cluster hosts:

1. Decommission the hosts.
2. Perform the necessary maintenance on the hosts.
3. Recommission the hosts.

Decommissioning Hosts

Cloudera Manager manages the host decommission and recommission process and allows you the option to specify whether to replicate the data to other DataNodes, and whether or not to suppress alerts.

About this task

Decommissioning a host decommissions and stops all roles on the host without requiring you to individually decommission the roles on each service. Decommissioning applies to only to HDFS DataNode, MapReduce TaskTracker, YARN NodeManager, and HBase RegionServer roles. If the host has other roles running on it, those roles are stopped.



Note: Hosts with DataNodes and DataNode roles themselves can only be decommissioned if the resulting action leaves enough DataNodes commissioned to maintain the configured HDFS replication factor (by default 3). If you attempt to decommission a DataNode or a host with a DataNode in such situations, the decommission process will not complete and must be aborted.

Before you begin

Minimum Required Role: [Limited Operator](#) (also provided by Operator, Configurator, Cluster Administrator, Limited Cluster Administrator, or Full Administrator).

Procedure

To decommission one or more hosts:

1. If the host has a DataNode, and you are planning to replicate data to other hosts (for longer term maintenance operations or to permanently decommission or repurpose the host), perform the steps in [Tuning HDFS Prior to Decommissioning DataNodes](#).
2. In Cloudera Manager, select the cluster where you want to decommission hosts.
3. In the left menu, click Hosts>All Hosts.
4. Select the hosts that you want to decommission.
5. Select Actions for SelectedBegin Maintenance (Suppress Alerts/Decommission).

(If you are logged in as a user with the Limited Operator or Operator role, the menu item is labeled Decommission Host(s) and you will not see the option to suppress alerts.)

The Begin Maintenance (Suppress Alerts/Decommission) dialog box opens. The role instances running on the hosts display at the top.

6. To decommission the hosts and suppress alerts, select Decommission Host(s). When you select this option for hosts running a DataNode role, choose one of the following (if the host is not running a DataNode role, you will only see the Decommission Host(s) option):

- Decommission DataNodes

This option re-replicates data to other DataNodes in the cluster according to the configured replication factor. Depending on the amount of data and other factors, this can take a significant amount of time and uses a great deal of network bandwidth. This option is appropriate when replacing disks, repurposing hosts for non-HDFS use, or permanently retiring hardware.

- Take DataNode Offline

This option does not re-replicate HDFS data to other DataNodes until the amount of time you specify has passed, making it less disruptive to active workloads. After this time has passed, the DataNode is automatically recommissioned, but the DataNode role is not started. This option is appropriate for short-term maintenance tasks such as not involving disks, such as rebooting, CPU/RAM upgrades, or switching network cables.



Caution: Taking multiple DataNodes offline simultaneously increases the chances that some HDFS data may become unavailable during maintenance. Configuring the proper value for the Maintenance State Minimal Block Replication HDFS configuration property will avoid risking data availability.

7. Click Begin Maintenance.

The Host Decommission Command dialog box opens and displays the progress of the command.

Results



Note:

- You cannot start roles on a decommissioned host.
- When a DataNode is decommissioned, although HDFS data is replicated to other DataNodes, local files containing the original data blocks are not automatically removed from the storage directories on the host. If you want to permanently remove these files from the host to reclaim disk space, you must do so manually.

What to do next

Perform the necessary maintenance on the hosts.

Recommissioning Hosts

About this task

Only hosts that are decommissioned using Cloudera Manager can be recommissioned.

Before you begin

Minimum Required Role: [Operator](#) (also provided by Configurator, Cluster Administrator, Limited Cluster Administrator, Full Administrator)

Procedure

1. In Cloudera Manager, select the cluster where you want to recommission hosts.
2. In the left menu, click HostsAll Hosts.
3. Select the hosts that you want to recommission.
4. Select Actions for SelectedEnd Maintenance (Suppress Alerts/Decommission. The End Maintenance (Suppress Alerts/Decommission dialog box opens. The role instances running on the hosts display at the top.
5. To recommission the hosts, select Recommission Host(s).
6. Choose one of the following:
 - Bring hosts online and start all roles
All decommissioned roles will be recommissioned and started. HDFS DataNodes will be started first and brought online before decommissioning to avoid excess replication.
 - Bring hosts online
All decommissioned roles will be recommissioned but remain stopped. You can [restart the roles](#) later.
7. Click End Maintenance.

Results

The Recommission Hosts and Start Roles Command dialog box opens and displays the progress of recommissioning the hosts and restarting the roles

Tuning and Troubleshooting Host Decommissioning

Decommissioning a host decommissions and stops all roles on the host without requiring you to individually decommission the roles on each service. The decommissioning process can take a long time and uses a great deal of cluster resources, including network bandwidth. You can tune the decommissioning process to improve performance and mitigate the performance impact on the cluster.

You can use the Decommission and Recommission features to perform minor maintenance on cluster hosts using Cloudera Manager to manage the process.

Tuning HDFS Prior to Decommissioning DataNodes

When a DataNode is decommissioned, the NameNode ensures that every block from the DataNode will still be available across the cluster as dictated by the replication factor. This procedure involves copying blocks from the DataNode in small batches. If a DataNode has thousands of blocks, decommissioning can take several hours. Before decommissioning hosts with DataNodes, you should first tune HDFS:

About this task

Minimum Required Role: [Configurator](#) (also provided by Cluster Administrator, Limited Cluster Administrator, and Full Administrator)

Procedure

1. Run the following command to identify any problems in the HDFS file system:

```
hdfs fsck / -list-corruptfileblocks -openforwrite -files -blocks -locations 2>&1 > /tmp/hdfs-fsck.txt
```

2. Fix any issues reported by the fsck command. If the command output lists corrupted files, use the fsck command to move them to the lost+found directory or delete them:

```
hdfs fsck file_name -move
```

or

```
hdfs fsck file_name -delete
```

3. Raise the heap size of the DataNodes. DataNodes should be configured with at least 4 GB heap size to allow for the increase in iterations and max streams.
 - a) Go to the HDFS service page.
 - b) Click the Configuration tab.
 - c) Select ScopeDataNode.
 - d) Select CategoryResource Management.
 - e) Set the Java Heap Size of DataNode in Bytes property as recommended.

To apply this configuration property to other role groups as needed, edit the value for the appropriate role group.

4. Increase the replication work multiplier per iteration to a larger number (the default is 2, however 10 is recommended).
 - a) Select ScopeNameNode.
 - b) Expand the CategoryAdvanced category.
 - c) Configure the Replication Work Multiplier Per Iteration property to a value such as 10.

To apply this configuration property to other role groups as needed, edit the value for the appropriate role group.

5. Increase the replication maximum threads and maximum replication thread hard limits.
 - a) Select ScopeNameNode.
 - b) Expand the CategoryAdvanced category.
 - c) Configure the Maximum number of replication threads on a DataNode and Hard limit on the number of replication threads on a DataNode properties to 50 and 100 respectively. You can decrease the number of threads (or use the default values) to minimize the impact of decommissioning on the cluster, but the trade off is that decommissioning will take longer.

To apply this configuration property to other role groups as needed, edit the value for the appropriate role group.

- Restart the HDFS service.

Related Information

[Performance Considerations](#)

[Modifying Configuration Properties Using Cloudera Manager](#)

Tuning HBase Prior to Decommissioning DataNodes

To increase the speed of a rolling restart of the HBase service, set the Region Mover Threads property to a higher value.

Minimum Required Role: [Configurator](#) (also provided by Cluster Administrator, Limited Cluster Administrator, and Full Administrator)

This increases the number of regions that can be moved in parallel, but places additional strain on the HMaster. In most cases, Region Mover Threads should be set to 5 or lower.

Performance Considerations

Decommissioning a DataNode does not happen instantly because the process requires replication of a potentially large number of blocks. During decommissioning, the performance of your cluster may be impacted.

This section describes the decommissioning process and suggests solutions for several common performance issues.

Decommissioning occurs in two steps:

- The Commission State of the DataNode is marked as Decommissioning and the data is replicated from this node to other available nodes. Until all blocks are replicated, the node remains in a Decommissioning state. You can view this state from the NameNode Web UI. (Go to the HDFS service and select Web UI/NameNode Web UI.)
- When all data blocks are replicated to other nodes, the node is marked as Decommissioned.

Decommissioning can impact performance in the following ways:

- There must be enough disk space on the other active DataNodes for the data to be replicated. After decommissioning, the remaining active DataNodes have more blocks and therefore decommissioning these DataNodes in the future may take more time.
- There will be increased network traffic and disk I/O while the data blocks are replicated.
- Data balance and data locality can be affected, which can lead to a decrease in performance of any running or submitted jobs.
- Decommissioning a large numbers of DataNodes at the same time can decrease performance.
- If you are decommissioning a minority of the DataNodes, the speed of data reads from these nodes limits the performance of decommissioning because decommissioning maxes out network bandwidth when reading data blocks from the DataNode and spreads the bandwidth used to replicate the blocks among other DataNodes in the cluster. To avoid performance impacts in the cluster, Cloudera recommends that you only decommission a minority of the DataNodes at the same time.
- You can decrease the number of replication threads to decrease the performance impact of the replications, but this will cause the decommissioning process to take longer to complete.

Cloudera recommends that you add DataNodes and decommission DataNodes in parallel, in smaller groups. For example, if the replication factor is 3, then you should add two DataNodes and decommission two DataNodes at the same time.

Related Information

[Tuning HDFS Prior to Decommissioning DataNodes](#)

Troubleshooting Performance of Decommissioning

Several conditions can impact performance when you decommission DataNodes.

Open Files

Write operations on the DataNode do not involve the NameNode. If there are blocks associated with open files located on a DataNode, they are not relocated until the file is closed. This commonly occurs with:

- Clusters using HBase

- Open Flume files
- Long running tasks

To find open files, run the following command:

```
hdfs dfsadmin -listOpenFiles -blockingDecommission
```

The command returns output similar to the following example:

```
Client Host          Client Name          Open File Path
172.26.12.77        DFSCClient_NONMAPREDUCE_-698274460_1 /hbase/ol
dWALs/dn3.cloudera.com%2C22101%2C1540973344249.dn3.cloudera.com%
2C22101%2C1540973344249.regiongroup-0.154099857098
```

After you find the open files, perform the appropriate action to restart process to close the file. For example, major compaction closes all files in a region for HBase.

Alternatively, you may evict writers to those decommissioning DataNodes with the following command:

```
hdfs dfsadmin -evictWriters <datanode_host:ipc_port>
```

For example:

```
hdfs dfsadmin -evictWriters datanode1:20001
```

A block cannot be relocated because there are not enough DataNodes to satisfy the block placement policy.

For example, for a 10 node cluster, if the `mapred.submit.replication` is set to the default of 10 while attempting to decommission one DataNode, there will be difficulties relocating blocks that are associated with map/reduce jobs. This condition will lead to errors in the NameNode logs similar to the following:

```
org.apache.hadoop.hdfs.server.blockmanagement.BlockPlacementPolicyDefault: Not able to place enough replicas, still in need of 3 to reach 3
```

Use the following steps to find the number of files where the block replication policy is equal to or above your current cluster size:

1. Provide a listing of open files, their blocks, the locations of those blocks by running the following command:

```
hadoop fsck / -files -blocks -locations -openforwrite 2>&1 >
openfiles.out
```

2. Run the following command to return a list of how many files have a given replication factor:

```
grep repl= openfiles.out | awk '{print $NF}' | sort | uniq -c
```

For example, when the replication factor is 10, and decommissioning one:

```
egrep -B4 "repl=10" openfiles.out | grep -v '<dir>' | awk '/^
\\/{print $1}'
```

3. Examine the paths, and decide whether to reduce the replication factor of the files, or remove them from the cluster.

Maintenance Mode

Maintenance mode allows you to suppress alerts for a host, service, role, or an entire cluster. This can be useful when you need to take actions in your cluster (make configuration changes and restart various elements) and do not want to see the alerts that will be generated due to those actions.



Putting an entity into maintenance mode does not prevent events from being logged; it only suppresses the alerts that those events would otherwise generate. You can see a history of all the events that were recorded for entities during the period that those entities were in maintenance mode.

Explicit and Effective Maintenance Mode

When you enter maintenance mode on an entity (cluster, service, or host) that has subordinate entities (for example, the roles for a service) the subordinate entities are also put into maintenance mode. These are considered to be in *effective maintenance mode*, as they have inherited the setting from the higher-level entity.

For example:

- If you set the HBase service into maintenance mode, then its roles (HBase Master and all RegionServers) are put into effective maintenance mode.
- If you set a host into maintenance mode, then any roles running on that host are put into effective maintenance mode.

Entities that have been explicitly put into maintenance mode show the icon . Entities that have entered effective maintenance mode as a result of inheritance from a higher-level entity show the icon .

When an entity (role, host or service) is in effective maintenance mode, it can only be removed from maintenance mode when the higher-level entity exits maintenance mode. For example, if you put a service into maintenance mode, the roles associated with that service are entered into effective maintenance mode, and remain in effective maintenance mode until the service exits maintenance mode. You cannot remove them from maintenance mode individually.


Alternatively, an entity that is in effective maintenance mode can be put into explicit maintenance mode. In this case, the entity remains in maintenance mode even when the higher-level entity exits maintenance mode. For example, suppose you put a host into maintenance mode, (which puts all the roles on that host into effective maintenance mode). You then select one of the roles on that host and put it explicitly into maintenance mode. When you have the host exit maintenance mode, that one role remains in maintenance mode. You need to select it individually and specifically have it exit maintenance mode.



Entering Maintenance Mode

You can enable maintenance mode for a cluster, service, role, or host.

Minimum Required Role: [Configurator](#) (also provided by Cluster Administrator, Limited Cluster Administrator, and Full Administrator)

Putting a Cluster into Maintenance Mode


1. In the left menu, click Clusters<cluster name>.
2. Click the Actions menu () to the right of the cluster name and select Enter Maintenance Mode.
3. Confirm that you want to do this.


The cluster is put into explicit maintenance mode, as indicated by the  icon. All services and roles in the cluster are entered into effective maintenance mode, as indicated by the  icon.

Putting a Service into Maintenance Mode

1. In the left menu, click Clusters and select the service.

2. Click ActionsEnter Maintenance Mode.
3. Confirm that you want to do this.

The service is put into explicit maintenance mode, as indicated by the  icon. All roles for the service are entered

into effective maintenance mode, as indicated by the  icon.

Putting Roles into Maintenance Mode

1. In the left menu, click Clusters and select the service.
2. Click the Instances tab.
3. Select the role(s) you want to put into maintenance mode.
4. From the Actions for Selected menu, select Enter Maintenance Mode.
5. Confirm that you want to do this.

The roles will be put in explicit maintenance mode. If the roles were already in effective maintenance mode (because its service or host was put into maintenance mode) the roles will now be in explicit maintenance mode. This means that they will not exit maintenance mode automatically if their host or service exits maintenance mode; they must be explicitly removed from maintenance mode.

Putting Hosts into Maintenance Mode

1. In Cloudera Manager, select the cluster where you want to decommission hosts.
2. Click HostsAll Hosts.
3. Select the hosts that you want to put into Maintenance Mode.
4. Select Actions for SelectedBegin Maintenance (Suppress Alerts/Decommission).

The Begin Maintenance (Suppress Alerts/Decommission) dialog box opens. The role instances running on the hosts display at the top. You can also use this dialog box to decommission the host.

5. Deselect the Decommission Host(s) option to put the host into Maintenance Mode. In this mode, alerts from the hosts are suppressed until the host exits Maintenance Mode. The events, however, are still logged. Hosts that are

currently in Maintenance Mode display the  icon.


6. Click Begin Maintenance.

The Host Decommission Command dialog box opens and displays the progress of the command.


Exiting Maintenance Mode

When you exit maintenance mode, the maintenance mode icons are removed and alert notification resumes.

Exiting a Cluster from Maintenance Mode

1. Click  to the right of the cluster name and select Exit Maintenance Mode.
2. Confirm that you want to do this.

Exiting a Service from Maintenance Mode

1. Click  to the right of the service name and select Exit Maintenance Mode.
2. Confirm that you want to do this.

Exiting Roles from Maintenance Mode

1. Go to the services page that includes the role.
2. Go to the Instances tab.
3. Select the role(s) you want to exit from maintenance mode.

4. From the Actions for Selected menu, select Exit Maintenance Mode.
5. Confirm that you want to do this.

Taking Hosts out of Maintenance Mode

1. In Cloudera Manager, go to the cluster with the hosts you want to take out of Maintenance Mode.
2. Click Hosts>All Hosts.
3. Select the hosts that are ready to exit Maintenance Mode.
4. Select Actions for SelectedEnd Maintenance (Suppress Alerts/Decommission).

The End Maintenance (Suppress Alerts/Decommission) dialog box opens. The role instances running on the hosts display at the top.

5. Deselect the Recommission Host(s) option to take the host out of Maintenance Mode and re-enable alerts from the

hosts. Hosts that are currently in Maintenance Mode display the  icon on the All Hosts page.

6. Click End Maintenance.

Viewing the Maintenance Mode Status of a Cluster

For any cluster, you can view the components (service, roles, or hosts) that are in maintenance mode.

About this task


Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

1. From the Cloudera Manager Home page, select the cluster that you want to view the maintenance mode status for.
2. Click Actions View Maintenance Mode Status... .

This pops up a dialog box that shows the components in your cluster that are in maintenance mode, and indicates which are in effective maintenance mode as well as those that have been explicitly placed into maintenance mode.

From this dialog box you can select any of the components shown there and remove them from maintenance mode.

If individual services are in maintenance mode, you will see the maintenance mode icon  next to the Actions button for that service.



Note: The Actions button is not enabled if you are viewing status for a point of time in the past.

Managing Roles

When Cloudera Manager configures a service, it configures hosts in your cluster with one or more functions (called roles in Cloudera Manager) that are required for that service. The role determines which Hadoop daemons run on a given host. For example, when Cloudera Manager configures an HDFS service instance it configures one host to run the NameNode role, another host to run as the Secondary NameNode role, another host to run the Balancer role, and some or all of the remaining hosts to run DataNode roles.

Configuration settings are organized in role groups. A *role group* includes a set of configuration properties for a specific group, as well as a list of role instances associated with that role group. Cloudera Manager automatically creates default role groups.

For role types that allow multiple instances on multiple hosts, such as DataNodes, TaskTrackers, RegionServers (and many others), you can create multiple role groups to allow one set of role instances to use different configuration settings than another set of instances of the same role type. In fact, upon initial cluster setup, if you are installing on

identical hosts with limited memory, Cloudera Manager will (typically) automatically create two role groups for each worker role — one group for the role instances on hosts with only other worker roles, and a separate group for the instance running on the host that is also hosting master roles.

The HDFS service is an example of this: Cloudera Manager typically creates one role group (DataNode Default Group) for the DataNode role instances running on the worker hosts, and another group (HDFS-1-DATANODE-1) for the DataNode instance running on the host that is also running the master roles such as the NameNode, JobTracker, HBase Master and so on. Typically the configurations for those two classes of hosts will differ in terms of settings such as memory for JVMs.

Cloudera Manager configuration screens offer two layout options: classic and new. The new layout is the default; however, on each configuration page you can easily switch between layouts using the Switch to XXX layout link at the top right of the page.

Gateway Roles

A *gateway* is a special type of role whose sole purpose is to designate a host that should receive a client configuration for a specific service, when the host does not have any roles running on it. Gateway roles enable Cloudera Manager to install and manage client configurations on that host. There is no process associated with a gateway role, and its status will always be Stopped. You can configure gateway roles for HBase, HDFS, Hive, Kafka, MapReduce, Solr, Spark, Sqoop 1 Client, and YARN.

Related Information

[Cluster Configuration Overview](#)

Role Instances

Adding a Role Instance

About this task

Minimum Required Role: [Configurator](#) (also provided by Cluster Administrator, Limited Cluster Administrator, and Full Administrator)

After creating services, you can add role instances to the services. For example, after initial installation in which you created the HDFS service, you can add a DataNode role instance to a host where one was not previously running. Upon upgrading a cluster to a new version of Cloudera Runtime you might want to create a role instance for a role added in the new version.

Procedure

1. Go to the service for which you want to add a role instance. For example, to add a DataNode role instance, go to the HDFS service.
2. Click the Instances tab.
3. Click the Add Role Instances button.

4. Customize the assignment of role instances to hosts. The wizard evaluates the hardware configurations of the hosts to determine the best hosts for each role. The wizard assigns all worker roles to the same set of hosts to which the HDFS DataNode role is assigned. You can reassign role instances.

Click a field below a role to display a dialog box containing a list of hosts. If you click a field containing multiple hosts, you can also select All Hosts to assign the role to all hosts, or Custom to display the hosts dialog box.

The following shortcuts for specifying hostname patterns are supported:

- Range of hostnames (without the domain portion)

Range Definition	Matching Hosts
10.1.1.[1-4]	10.1.1.1, 10.1.1.2, 10.1.1.3, 10.1.1.4
host[1-3].company.com	host1.company.com, host2.company.com, host3.company.com
host[07-10].company.com	host07.company.com, host08.company.com, host09.company.com, host10.company.com

- IP addresses
- Rack name

Click the View By Host button for an overview of the role assignment by hostname ranges.

5. Click Continue.

6. In the Review Changes page, review the configuration changes to be applied.

Confirm the settings entered for file system paths. The file paths required vary based on the services to be installed. For example, you might confirm the NameNode Data Directory and the DataNode Data Directory for HDFS.

7. Click Continue.

Results

The wizard finishes by performing any actions necessary to prepare the cluster for the new role instances. For example, new DataNodes are added to the NameNode `dfs_hosts_allow.txt` file. The new role instance is configured with the default role group for its role type, even if there are multiple role groups for the role type. If you want to use a different role group, follow the instructions in the topic *Managing Role Groups* for moving role instances to a different role group.

Related Information

[Managing Role Groups](#)

Starting, Stopping, and Restarting Role Instances

About this task

Minimum Required Role: [Operator](#) (also provided by Configurator, Cluster Administrator, Limited Cluster Administrator, Full Administrator)

If the host for the role instance is currently decommissioned, you will not be able to start the role until the host has been recommissioned.



Important: Use Cloudera Manager to stop the Node Manager service. If it is stopped manually, it can cause jobs to fail.

Procedure

1. Go to the service that contains the role instances to start, stop, or restart.
2. Click the Instances tab.
3. Check the checkboxes next to the role instances to start, stop, or restart (such as a DataNode instance).
4. Select Actions for SelectedStart, Stop, or Restart, and then click Start, Stop, or Restart again to start the process. When you see a Finished status, the process has finished.

Related Information

[Rolling Restart](#)

Decommissioning Role Instances

You can remove a role instance such as a DataNode from a cluster while the cluster is running by decommissioning the role instance.

About this task

Minimum Required Role: [Operator](#) (also provided by Configurator, Cluster Administrator, Limited Cluster Administrator, Full Administrator)

When you decommission a role instance, Cloudera Manager performs a procedure so that you can safely retire a host without losing data. Role decommissioning applies to HDFS DataNode, MapReduce TaskTracker, YARN NodeManager, and HBase RegionServer roles.

Hosts with DataNodes and DataNode roles themselves can only be decommissioned if the resulting action leaves enough DataNodes commissioned to maintain the configured HDFS replication factor (by default 3). If you attempt to decommission a DataNode or a host with a DataNode in such situations, the decommission process will not complete and must be aborted.

A role will be decommissioned if its host is decommissioned.


To remove a DataNode from the cluster, you decommission the DataNode role as described here and then perform a few additional steps to remove the role. See the topic [Delete a DataNode](#).

Procedure

To decommission role instances:

1. If you are decommissioning DataNodes, perform the steps in the topic *Tuning HDFS Prior to Decommissioning DataNodes*.
2. Click the service instance that contains the role instance you want to decommission.
3. Click the Instances tab.
4. Check the checkboxes next to the role instances to decommission.
5. Select Actions for SelectedDecommission, and then click Decommission again to start the process.

Results

A Decommission Command pop-up displays that shows each step or decommission command as it is run. In the Details area, click  to see the subcommands that are run. Depending on the role, the steps may include adding the host to an "exclusions list" and refreshing the NameNode, JobTracker, or NodeManager; stopping the Balancer (if it is running); and moving data blocks or regions. Roles that do not have specific decommission actions are stopped.

You can abort the decommission process by clicking the Abort button, but you must recommit and restart the role.

The Commission State facet in the Filters list displays  Decommissioning while decommissioning is in progress, and  Decommissioned when the decommissioning process has finished. When the process is complete, a  is added in front of Decommission Command.

Related Information

[Tuning HDFS Prior to Decommissioning DataNodes](#)

Recommissioning Role Instances

About this task

Minimum Required Role: [Operator](#) (also provided by Configurator, Cluster Administrator, Limited Cluster Administrator, Full Administrator)

Procedure

1. Click the service that contains the role instance you want to recommission.
2. Click the Instances tab.
3. Check the checkboxes next to the decommissioned role instances to recommission.
4. Select Actions for SelectedRecommission, and then click Recommission to start the process. A Recommission Command pop-up displays that shows each step or recommission command as it is run. When the process is complete, a ✓ is added in front of Recommission Command.
5. Restart the role instance.

Deleting Role Instances

About this task

[Configurator](#) (also provided by Cluster Administrator, Full Administrator)

Deleting Role Instances

Procedure

1. Click the service instance that contains the role instance you want to delete. For example, if you want to delete a DataNode role instance, click an HDFS service instance.
2. Click the Instances tab.
3. Check the checkboxes next to the role instances you want to delete.
4. If the role instance is running, select Actions for SelectedStop and click Stop to confirm the action.
5. Select Actions for SelectedDelete. Click Delete to confirm the deletion.

Results



Note: Deleting a role instance does not clean up the associated client configurations that have been deployed in the cluster.

Configuring Roles to Use a Custom Garbage Collection Parameter

You can use Java configuration options to configure roles to use a custom garbage collection parameter.

Every Java-based role in Cloudera Manager has a configuration setting called Java Configuration Options for *role* where you can enter command line options. Commonly, garbage collection flags or extra debugging flags would be passed here. To find the appropriate configuration setting, select the service you want to modify in the Cloudera Manager Admin Console, then use the Search box to search for Java Configuration Options.

You can add configuration options for all instances of a given role by making this configuration change at the service level. For example, to modify the setting for all DataNodes, select the HDFS service, then modify the Java Configuration Options for DataNode setting.

To modify a configuration option for a given instance of a role, select the service, then select the particular role instance (for example, a specific DataNode). The configuration settings you modify will apply to the selected role instance only.

Related Information

[Modifying Configuration Properties Using Cloudera Manager](#)

Role Groups

Minimum Required Role: [Configurator](#) (also provided by Cluster Administrator, Limited Cluster Administrator , and Full Administrator)

A *role group* is a set of configuration properties for a role type, as well as a list of role instances associated with that group. Cloudera Manager automatically creates a default role group named *Role Type Default Group* for each role type. Each role instance can be associated with only a single role group.

Role groups provide two types of properties: those that affect the configuration of the service itself and those that affect monitoring of the service, if applicable (the Monitoring subcategory). Not all services have monitoring properties.

When you run the installation or upgrade wizard, Cloudera Manager configures the default role groups it adds, and adds any other required role groups for a given role type. For example, a DataNode role on the same host as the NameNode might require a different configuration than DataNode roles running on other hosts. Cloudera Manager creates a separate role group for the DataNode role running on the NameNode host and uses the default configuration for DataNode roles running on other hosts.

You can modify the settings of the default role group, or you can create new role groups and associate role instances to whichever role group is most appropriate. This simplifies the management of role configurations when one group of role instances may require different settings than another group of instances of the same role type—for example, due to differences in the hardware the roles run on. You modify the configuration for any of the service's role groups through the Configuration tab for the service. You can also override the settings inherited from a role group for a role instance.

If there are multiple role groups for a role type, you can move role instances from one group to another. When you move a role instance to a different group, it inherits the configuration settings for its new group.

Related Information

[Configuring Monitoring Settings](#)

[Overriding Configuration Properties](#)

Creating a Role Group

About this task

Minimum Required Role: [Configurator](#) (also provided by Cluster Administrator, Limited Cluster Administrator , and Full Administrator)

Procedure

1. Go to a service status page.
2. Click the Instances or Configuration tab.
3. Click Role Groups.
4. Click Create new group....
5. Provide a name for the group.
6. Select the role type for the group. You can select role types that allow multiple instances and that exist for the service you have selected.
7. In the Copy From field, select the source of the basic configuration information for the role group:
 - An existing role group of the appropriate type.
 - None.... The role group is set up with generic default values that are not the same as the values Cloudera Manager sets in the default role group, as Cloudera Manager specifically sets the appropriate configuration properties for the services and roles it installs. After you create the group you must edit the configuration to set missing properties (for example the TaskTracker Local Data Directory List property, which is not populated if you select None) and clear other validation warnings and errors.

Related Information

[Modifying Configuration Properties Using Cloudera Manager](#)

Managing Role Groups

About this task

Minimum Required Role: [Configurator](#) (also provided by Cluster Administrator, Limited Cluster Administrator, and Full Administrator)

Procedure

1. Go to a service status page.
2. Click the Instances or Configuration tab.
3. Click Role Groups.
4. Click the group you want to manage. Role instances assigned to the role group are listed.
5. Perform the appropriate procedure for the action:

- Rename
 - a. Click the role group name, and click Rename.
 - b. Specify the new name and click Rename.

- Delete

You cannot delete any of the default groups. The group must first be empty; if you want to delete a group you've created, you must move any role instances to a different role group.

- a. Click the role group name.
- b. Click Delete, and confirm by clicking Delete. Deleting a role group removes it from host templates.

- Move

- a. Select the role instance(s) to move.
- b. Select Actions for Selected Move To Different Role Group....
- c. In the pop-up that appears, select the target role group and click Move.

Related Information

[Managing Hosts](#)

Default User Roles

By default, Cloudera Manager ships with user roles that have privileges for all clusters managed by Cloudera Manager.

The following table describes the actions each user role can perform:

Permitted Operations	Auditor	Cluster Administrator	Configurator	Dashboard User	Full Administrator	Key Administrator	Limited Operator	Navigator Administrator	Operator	Read-Only	Replication Administrator	User Administrator
Apply policies to redact sensitive data		Y			Y							
Administer Cloudera Navigator					Y			Y				
Create, modify, and delete your own dashboards				Y	Y							
Manage user accounts and configuration of external authentication					Y							Y
See available hosts		Y			Y							
View and perform parcels operations		Y			Y							
Enter and exit Maintenance Mode		Y	Y		Y							
Edit the configuration of services and roles		Y	Y		Y							
View data in Cloudera Manager	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y
Start, stop, and restart KMS		Y	Y		Y	Y			Y			
Manage Full Administrator accounts					Y							
Decommission hosts		Y	Y		Y		Y		Y			
Create clusters		Y			Y							
View audit events	Y				Y			Y				
Create, update, or delete external account configuration					Y							Y
Configure HDFS Encryption, administer Key Trustee Server, and manage encryption keys					Y	Y						
Recommission hosts, and decommission and recommission roles		Y	Y		Y				Y			
Access all functionality that Cloudera Manager offers		Y			Y							
Create replication policies and snapshot policies					Y						Y	
Start, stop, and restart most clusters, services, and roles		Y	Y		Y				Y			

Managing Cloudera Runtime Services

Cloudera Manager service configuration features let you manage the deployment and configuration of Cloudera Runtime and managed services.

You can add new services and roles if needed, gracefully start, stop and restart services or roles, and decommission and delete roles or services if necessary. Further, you can modify the configuration properties for services or for individual role instances. You can also view past configuration changes and roll back to a previous revision. You can also generate client configuration files, enabling you to easily distribute them to the users of a service.

The topics in this chapter describe how to configure and use the services on your cluster. Some services have unique configuration requirements or provide unique features. See the documentation for an individual service for more information.

Adding a Service

After initial installation, you can use the **Add a Service** wizard to add and configure new service instances. For example, you may want to add a service such as Oozie that you did not select in the wizard during the initial installation.

About this task

Minimum Required Role: [Full Administrator](#). This feature is not available when using Cloudera Manager to manage Data Hub clusters.




Note:

The binaries for the following services are not packaged in Cloudera Runtime and must be installed individually before being adding the service:

- Accumulo
- Kafka
- Key Trustee KMS

If you do not add the binaries before adding the service, the service will fail to start.

Procedure

1. On the HomeStatus tab, click  to the right of the cluster name and select Add a Service. A list of service types display. You can add one type of service at a time.
2. Select a service and click Continue. If you are missing required binaries, a pop-up displays asking if you want to continue with adding the service.
3. Select the services on which the new service should depend. All services must depend on the same ZooKeeper service. Click Continue.
4. Customize the assignment of role instances to hosts. The wizard evaluates the hardware configurations of the hosts to determine the best hosts for each role. The wizard assigns all worker roles to the same set of hosts to which the HDFS DataNode role is assigned. You can reassign role instances.

Click a field below a role to display a dialog box containing a list of hosts. If you click a field containing multiple hosts, you can also select All Hosts to assign the role to all hosts, or Custom to display the hosts dialog box.

The following shortcuts for specifying hostname patterns are supported:

- Range of hostnames (without the domain portion)

Range Definition	Matching Hosts
10.1.1.[1-4]	10.1.1.1, 10.1.1.2, 10.1.1.3, 10.1.1.4
host[1-3].company.com	host1.company.com, host2.company.com, host3.company.com
host[07-10].company.com	host07.company.com, host08.company.com, host09.company.com, host10.company.com

- IP addresses
- Rack name

Click the View By Host button for an overview of the role assignment by hostname ranges.

5. Review and modify configuration settings, such as data directory paths and heap sizes and click Continue. The service is started.



Note: If you are adding the Ranger service, passwords for the Ranger Admin, Usersync, Tagsync, and KMS Keyadmin users must be a minimum of 8 characters long, with at least one alphabetic and one numeric character. The following characters are not valid: " ' \ ` ´ .

6. Click Continue then click Finish. You are returned to the home page.
7. Verify the new service is started properly by checking the health status for the new service. If the Health Status is Good, then the service started properly.

Comparing Configurations for a Service Between Clusters

You can compare the configuration settings for a particular service between two different clusters in a Cloudera Manager deployment.

About this task

Minimum Required Role: [Configurator](#) (also provided by Cluster Administrator, Limited Cluster Administrator, and Full Administrator)

Procedure

1. On the HomeStatus tab, click the name of the service you want to compare, or click the Clusters menu and select the name of the service.
2. Click the Configuration tab.

3. Click the drop-down menu above the Filters pane, and select from one of the options that begins Diff with...:
 - *service on cluster* - For example, HBASE-1 on Cluster 1. This is the default display setting. All properties are displayed for the selected instance of the service.
 - *service on all clusters* - For example, HBase on all clusters. All properties are displayed for all instances of the service.
 - Diff with *service on cluster* - For example, Diff with HBase on Cluster 2. Properties are displayed only if the values for the instance of the service whose page you are on differ from the values for the instance selected in the drop-down menu.
 - Diff with *service on all clusters* - For example, Diff with HBase on all clusters. Properties are displayed if the values for the instance of the service whose page you are on differ from the values for one or more other instances in the Cloudera Manager deployment.

The service's properties will be displayed showing the values for each property for the selected clusters. The filters on the left side can be used to limit the properties displayed.

You can also view property configuration values that differ between clusters across a deployment by selecting Non-uniform Values on the Configuration tab of the Cloudera Manager HomeStatus page.

Starting a Cloudera Runtime Service on All Hosts

Starting and Stopping Cloudera Runtime services.

About this task

Minimum Required Role: [Operator](#) (also provided by Configurator, Cluster Administrator, Limited Cluster Administrator, Full Administrator)

It is important to start and stop services that have dependencies in the correct order. For example, because MapReduce and YARN have a dependency on HDFS, you must start HDFS before starting MapReduce or YARN. The Cloudera Management Service and Hue are the only two services on which no other services depend; although you can start and stop them at anytime, their preferred order is shown in the following procedures.

The Cloudera Manager cluster actions start and stop services in the correct order. To start or stop all services in a cluster, follow the instructions in Starting, Stopping, Refreshing, and Restarting a Cluster.


Before you begin

The order in which to start services is:

1. Cloudera Management Service
2. ZooKeeper
3. HDFS
4. Solr
5. HBase
6. Key-Value Store Indexer
7. MapReduce or YARN
8. Hive
9. Impala
10. Oozie
11. Sqoop
12. Hue

Procedure

1. In the left menu, click Clusters and select a service.

2. Click  to the right of the service name and select Start.
3. Click Start in the next screen to confirm.
When you see a Finished status, the service has started.

Results



Note: If you are unable to start the HDFS service, it's possible that one of the roles instances, such as a DataNode, was running on a host that is no longer connected to the Cloudera Manager Server host, perhaps because of a hardware or network failure. If this is the case, the Cloudera Manager Server will be unable to connect to the Cloudera Manager Agent on that disconnected host to start the role instance, which will prevent the HDFS service from starting. To work around this, you can stop all services, abort the pending command to start the role instance on the disconnected host, and then restart all services again without that role instance.

Related Information

[Aborting a Pending Command](#)

Stopping a Cloudera Runtime Service on All Hosts

About this task


[Operator](#) (also provided by Configurator, Cluster Administrator, Full Administrator)

Before you begin

The order in which to stop services is:

1. Hue
2. Sqoop
3. Oozie
4. Impala
5. Hive
6. MapReduce or YARN
7. Key-Value Store Indexer
8. HBase
9. Flume
10. Solr
11. HDFS
12. ZooKeeper
13. Cloudera Management Service

Procedure

1. In the left menu, click Clusters and select a service.
2. Click  to the right of the service name and select Stop.
3. Click Stop in the next screen to confirm.
When you see a Finished status, the service has stopped.

Restarting a Cloudera Runtime Service


About this task

Operator (also provided by Configurator, Cluster Administrator, Full Administrator)

Before you begin




Important: If you have changed a configuration property that requires a redeployment of the client configurations, note that refreshing or restarting a cluster does not automatically re-deploy the client

configurations. A service or cluster displays a staleness icon () next to the cluster or service name that indicates that you must redeploy the client configuration. Click the icon to open the [Stale Configurations](#) page and follow the prompts to refresh the cluster and redeploy the client configuration.

Alternatively, after a restart is completed, you can select Deploy Client Configuration from the Actions menu for either a service or cluster.

Procedure

1. In the left menu, click Clusters and select a service.
2. Click  to the right of the service name and select Restart.
3. Click Start on the next screen to confirm.

Results

When you see a Finished status, the service has restarted.

What to do next

To restart all services, restart the cluster.

Rolling Restart

Minimum Required Role: **Operator** (also provided by Configurator, Cluster Administrator, Limited Cluster Administrator, Full Administrator)

Rolling restart allows you to conditionally restart the role instances of the following services to update software or use a new configuration:

- HBase
- HDFS
- Kafka
- Key Trustee Server
- MapReduce
- Oozie
- YARN
- ZooKeeper


If the service is not running, rolling restart is not available for that service. You can specify a rolling restart of each service individually.

If you have [HDFS High Availability](#) enabled, you can also perform a cluster-level rolling restart. At the cluster level, the rolling restart of worker hosts is performed on a host-by-host basis, rather than per service, to avoid all roles for a service potentially being unavailable at the same time. During a cluster restart, to avoid having your NameNode (and thus the cluster) be unavailable during the restart, Cloudera Manager forces a failover to the standby NameNode.

Job Tracker and Resource Manager High availability are not required for a cluster-level rolling restart. However, if you have JobTracker or ResourceManager high availability enabled, Cloudera Manager will force a failover to the standby JobTracker or ResourceManager.



Important: If you have changed a configuration property that requires a redeployment of the client configurations, note that refreshing or restarting a cluster does not automatically re-deploy the client

configurations. A service or cluster displays a staleness icon () next to the cluster or service name that indicates that you must redeploy the client configuration. Click the icon to open the [Stale Configurations](#) page and follow the prompts to refresh the cluster and redeploy the client configuration.

Alternatively, after a restart is completed, you can select Deploy Client Configuration from the Actions menu for either a service or cluster.

Performing a Service or Role Rolling Restart

You can initiate a rolling restart from either the Status page for one of the eligible services, or from the service's Instances page, where you can select individual roles to be restarted.

1. Go to the service you want to restart.
2. Do one of the following:
 - service - Select ActionsRolling Restart.
 - role -
 - a. Click the Instances tab.
 - b. Select the roles to restart.
 - c. Select Actions for SelectedRolling Restart.
3. In the pop-up dialog box, select the options you want:
 - Restart only roles whose configurations are stale
 - Restart only roles that are running outdated software versions
 - Which role types to restart
4. If you select an HDFS, HBase, MapReduce, or YARN service, you can have their worker roles restarted in batches. You can configure:
 - How many roles should be included in a batch - Cloudera Manager restarts the worker roles rack-by-rack in alphabetical order, and within each rack, hosts are restarted in alphabetical order. If you are using the default replication factor of 3, Hadoop tries to keep the replicas on at least 2 different racks. So if you have multiple racks, you can use a higher batch size than the default 1. But you should be aware that using too high batch size also means that fewer worker roles are active at any time during the upgrade, so it can cause temporary performance degradation. If you are using a single rack only, you should only restart one worker node at a time to ensure data availability during upgrade.
 - How long should Cloudera Manager wait before starting the next batch.

- The number of batch failures that will cause the entire rolling restart to fail (this is an advanced feature). For example if you have a very large cluster you can use this option to allow failures because if you know that your cluster will be functional even if some worker roles are down.



Note:

- HDFS - If you do not have HDFS high availability configured, a warning appears reminding you that the service will become unavailable during the restart while the NameNode is restarted. Services that depend on that HDFS service will also be disrupted. Cloudera recommends that you restart the DataNodes one at a time—one host per batch, which is the default.
- HBase
 - Administration operations such as any of the following should not be performed during the rolling restart, to avoid leaving the cluster in an inconsistent state:
 - Split
 - Create, disable, enable, or drop table
 - Metadata changes
 - Create, clone, or restore a snapshot. Snapshots rely on the RegionServers being up; otherwise the snapshot will fail.
 - To increase the speed of a rolling restart of the HBase service, set the Region Mover Threads property to a higher value. This increases the number of regions that can be moved in parallel, but places additional strain on the HMaster. In most cases, Region Mover Threads should be set to 5 or lower.
 - Another option to increase the speed of a rolling restart of the HBase service is to set the Skip Region Reload During Rolling Restart property to true. This setting can cause regions to be moved around multiple times, which can degrade HBase client performance.
- MapReduce - If you restart the JobTracker, all current jobs will fail.
- YARN - If you restart ResourceManager and ResourceManager HA is enabled, current jobs continue running; they do not restart or fail.
- ZooKeeper and Flume - For both ZooKeeper and Flume, the option to restart roles in batches is not available. They are always restarted one by one.

5. Click Confirm to start the rolling restart.

Performing a Cluster-Level Rolling Restart

You can perform a cluster-level rolling restart on demand from the Cloudera Manager Admin Console. A cluster-level rolling restart is also performed as the last step in a rolling upgrade when the cluster is configured with HDFS high availability enabled.

1. If you have not already done so, enable high availability. See [HDFS High Availability](#) for instructions. You do not need to enable automatic failover for rolling restart to work, though you can enable it if you want. Automatic failover does not affect the rolling restart operation.
2. For the cluster you want to restart select ActionsRolling Restart.
3. In the pop-up dialog box, select the services you want to restart. Please review the caveats in the preceding section for the services you elect to have restarted. The services that do not support rolling restart will simply be restarted, and will be unavailable during their restart.
4. If you select an HDFS, HBase, or MapReduce service, you can have their worker roles restarted in batches. You can configure:
 - How many roles should be included in a batch - Cloudera Manager restarts the worker roles rack-by-rack in alphabetical order, and within each rack, hosts are restarted in alphabetical order. If you are using the default replication factor of 3, Hadoop tries to keep the replicas on at least 2 different racks. So if you have multiple racks, you can use a higher batch size than the default 1. But you should be aware that using too high batch size also means that fewer worker roles are active at any time during the upgrade, so it can cause temporary performance degradation. If you are using a single rack only, you should only restart one worker node at a time to ensure data availability during upgrade.

- How long should Cloudera Manager wait before starting the next batch.
 - The number of batch failures that will cause the entire rolling restart to fail (this is an advanced feature). For example if you have a very large cluster you can use this option to allow failures because if you know that your cluster will be functional even if some worker roles are down.
5. Click Restart to start the rolling restart. While the restart is in progress, the Command Details page shows the steps for stopping and restarting the services.


Aborting a Pending Command

Minimum Required Role: [Operator](#) (also provided by Configurator, Cluster Administrator, Limited Cluster Administrator, Full Administrator)

Commands will time out if they are unable to complete after a period of time.

If necessary, you can abort a pending command. For example, this may become necessary because of a hardware or network failure where a host running a role instance becomes disconnected from the Cloudera Manager Server host. In this case, the Cloudera Manager Server will be unable to connect to the Cloudera Manager Agent on that disconnected host to start or stop the role instance which will prevent the corresponding service from starting or stopping. To work around this, you can abort the command to start or stop the role instance on the disconnected host, and then you can start or stop the service again.

To abort any pending command:

You can click the Recent Commands indicator (), which shows the number of commands that are currently running in your cluster (if any). This indicator is positioned above the Support link at the bottom of the left menu. Unlike the Commands tab for a role or service, this indicator includes all commands running for all services or roles in the cluster. In the **Running Commands** window, click Abort to abort the pending command.

To abort a pending command for a service or role:

1. In the left menu, click Clusters and select the service where the role instance you want to stop is located. For example, click ClustersHDFS Service if you want to abort a pending command for a DataNode.
2. Click the Instances tab.
3. In the list of instances, click the link for role instance where the command is running (for example, the instance that is located on the disconnected host).
4. Go to the Commands tab.
5. Find the command in the list of Running Commands and click Abort Command to abort the running command.

Related Information

[Viewing Running and Recent Commands](#)


Deleting Services

You can delete a service from the **Status** tab.

About this task

Minimum Required Role: [Full Administrator](#). This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

1. Stop the service.
2. On the HomeStatus tab, click  to the right of the service name and select Delete.

3. Click Delete to confirm the deletion. Deleting a service does not clean up the associated client configurations that have been deployed in the cluster or the user data stored in the cluster. For a given "alternatives path" (for example `/etc/hadoop/conf`) if there exist both "live" client configurations (ones that would be pushed out with deploy client configurations for active services) and ones that have been "orphaned" client configurations (the service they correspond to has been deleted), the orphaned ones will be removed from the alternatives database. In other words, to trigger cleanup of client configurations associated with a deleted service you must create a service to replace it. To remove user data, see the topic *Remove Cloudera Manager and User Data*.

Renaming a Service

A service is given a name upon installation, and that name is used as an identifier internally. However, Cloudera Manager allows you to provide a display name for a service, and that name will appear in the Cloudera Manager Admin Console instead of the original (internal) name.


About this task

Minimum Required Role: [Full Administrator](#). This feature is not available when using Cloudera Manager to manage Data Hub clusters.

**Note:**

The original service name will still be used internally, and may appear or be required in certain circumstances, such as in log messages or in the API.

Procedure

1. On the HomeStatus tab, click  to the right of the service name and select Rename.
2. Type the new name.
3. Click Rename.

The rename action is recorded as an Audit event.

When looking at Audit or Event search results for the renamed service, it is possible that these search results might contain either only the original (internal) name, or both the display name and the original name.

Configuring Maximum File Descriptors

You can set the maximum file descriptor parameter for all daemon roles. When not specified, the role uses whatever value it inherits from supervisor. When specified, configures soft and hard limits to the configured value.

About this task

Minimum Required Role: [Configurator](#) (also provided by Cluster Administrator, Limited Cluster Administrator, and Full Administrator)

Procedure

1. Go to a service.
2. Click the **Configuration** tab.
3. In the Search box, type `rlimit_fds`.
4. Set the Maximum Process File Descriptors property for one or more roles.
5. Enter a Reason for change, and then click Save Changes to commit the changes.
6. Restart the affected role instances.

Extending Cloudera Manager

In addition to the set of software packages and services managed by Cloudera Manager, you can also define and add new types of services using [custom service descriptors](#). When you deploy a custom service descriptor, the implementation is delivered in a Cloudera Manager [parcel](#) or other software package. For information on the extension mechanisms provided by Cloudera Manager for creating custom service descriptors and parcels, see [Cloudera Manager Extensions](#).

Add-on Services

Minimum Required Role: [Full Administrator](#). This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Cloudera Manager supports adding new types of services (referred to as an *add-on service*) to Cloudera Manager, allowing such services to leverage Cloudera Manager distribution, configuration, monitoring, resource management, and life-cycle management features. An add-on service can be provided by Cloudera or an independent software vendor (ISV). If you have multiple clusters managed by Cloudera Manager, an add-on service can be deployed on any of the clusters.



Note: If the add-on service is already installed and running on hosts that are not currently being managed by Cloudera Manager, you must first add the hosts to a cluster that's under management. See [Adding a Host to a Cluster](#) on page 11 for details.

Custom Service Descriptor Files

Integrating an add-on service requires a Custom Service Descriptor (CSD) file. A CSD file contains all the configuration needed to describe and manage a new service. A CSD is provided in the form of a JAR file.

Depending on the service, the CSD and associated software may be provided by Cloudera or by an ISV. The integration process assumes that the add-on service software (parcel or package) has been installed and is present on the cluster. The recommended method is for the ISV to provide the software as a parcel, but the actual mechanism for installing the software is up to the ISV. The instructions in [Installing an Add-on Service](#) on page 52 assume that you have obtained the CSD file from the Cloudera repository or from an ISV. It also assumes you have obtained the service software, ideally as a parcel, and have or will install it on your cluster either prior to installing the CSD or as part of the CSD installation process.

Configuring the Location of Custom Service Descriptor Files

The default location for CSD files is `/opt/cloudera/csd`. You can change the location in the Cloudera Manager Admin Console as follows:

1. Select AdministrationSettings.
2. Click the Custom Service Descriptors category.
3. Edit the Local Descriptor Repository Path property.
4. Enter a Reason for change, and then click Save Changes to commit the changes.
5. Restart Cloudera Manager Server:

RHEL 7 compatible, SLES, Ubuntu:

```
sudo systemctl restart cloudera-scm-server
```

RHEL 6 compatible:

```
sudo service cloudera-scm-server restart
```


Installing an Add-on Service

An ISV may provide its software in the form of a parcel, or they may have a different way of installing their software. If their software is not available as a parcel, then you must install their software before adding the CSD file. Follow the instructions from the ISV for installing the software. If the ISV has provided their software as a parcel, they may also have included the location of their parcel repository in the CSD they have provided. In that case, install the CSD first and then install the parcel.

Installing the Custom Service Descriptor File

1. Acquire the CSD file from Cloudera or an ISV.
2. Log on to the Cloudera Manager Server host, and place the CSD file under the [location configured](#) for CSD files.
3. Set the file ownership to `cloudera-scm:cloudera-scm` with permission 644.
4. Restart the Cloudera Manager Server:

```
service cloudera-scm-server restart
```

5. Log into the Cloudera Manager Admin Console and restart the Cloudera Management Service.
 - a. Do one of the following:
 - 1. Select Clusters Cloudera Management Service .
 - 2. Select ActionsRestart.
 - On the HomeStatus tab, click  to the right of Cloudera Management Service and select Restart.
 - b. Click Restart to confirm. The Command Details window shows the progress of stopping and then starting the roles.
 - c. When Command completed with *n/n* successful subcommands appears, the task is complete. Click Close.

Installing the Parcel



Note: It is not required that the Cloudera Manager server host be part of a managed cluster and have an agent installed. Although you initially copy the CSD file to the Cloudera Manager server, the Parcel for the add-on service will not be installed on the Cloudera Manager Server host unless the host is managed by Cloudera Manager.

If you have already installed the external software onto your cluster, you can skip these steps and proceed to [Adding an Add-on Service](#) on page 53.

1. Click Parcels in the left menu. If the vendor has included the location of the repository in the CSD, the parcel should already be present and ready for downloading. If the parcel is available, skip to [step 7](#).
2. Use one of the following methods to open the parcel settings page:
 - Navigation bar
 - a. Click the parcel icon in the top navigation bar or click Hosts and click the Parcels tab.
 - b. Click the Configuration button.
 - Menu
 - a. Select AdministrationSettings.
 - b. Select CategoryParcels.
3. In the Remote Parcel Repository URLs list, click the addition symbol to open an additional row.
4. Enter the path to the repository.
5. Enter a Reason for change, and then click Save Changes to commit the changes.
6. Click Parcels in the left navigation menu. The external parcel should appear in the set of parcels available for download.
7. Download, distribute, and activate the parcel.

Adding an Add-on Service

Add the service following the procedure in [Adding a Service](#).

Uninstalling an Add-on Service

1. Stop all instances of the service.
2. Delete the service from all clusters. If there are other services that depend on the service you are trying to delete, you must delete those services first.
3. Log on to the Cloudera Manager Server host and remove the CSD file.
4. Restart the Cloudera Manager Server:

```
service cloudera-scm-server restart
```

5. After the server has restarted, log into the Cloudera Manager Admin Console and restart the Cloudera Management Service.
6. Optionally remove the parcel.

Core Configuration Service

The Core Configuration Service allows you to create clusters without the HDFS service.

Minimum Required Role: [Configurator](#) (also provided by Cluster Administrator, Limited Cluster Administrator, and Full Administrator)

The Core Configuration service allows you to create more types of clusters without having to include the HDFS service. Previously, the HDFS service was required in many cases even when data was not being stored in HDFS because some services like Sentry and Spark required cluster-wide configuration files that Cloudera Manager deploys within the HDFS service. The Core Configuration service provides this configuration in a standalone fashion and thus eliminates the need for an HDFS service for certain types of clusters where no HDFS storage is required (e.g. Kudu, Kafka, or 'Compute' clusters using exclusively object storage like S3 or ADLS). The Core Configuration service is also useful when creating a Compute cluster that accesses data on an HDFS service located in the Base cluster.

You can add the Core Configuration service in the following ways:

- When installing a Cloudera Runtime cluster using the installation wizard you can select the Core Configuration from the list of services.
 - When adding a new cluster using the Add Cluster wizard, you can select the Core Configuration from the list of services.
 - You can add the Core Configuration service to an existing cluster that does not include HDFS:
1. Open the Cloudera Manager Admin Console and navigate to the cluster where you want to add the service.
 2. Click Actions > Add Service
 3. Select Core Configuration from the list of services.
 4. Click Continue.
 5. Follow the prompts on the screen to complete adding the service.



Important: You cannot use both the HDFS service and the Core Configuration service in the same cluster.

Managing Cloudera Manager

Automatic Logout

For security purposes, Cloudera Manager automatically logs out a user session after 30 minutes. You can change this session logout period.

Procedure

1. Click AdministrationSettings.
2. Click CategorySecurity.
3. Edit the Session Timeout property.
4. Enter a Reason for change, and then click Save Changes to commit the changes.

When the timeout is one minute from triggering, the user sees the following message:

Automatic Logout for Your Protection ✕

Due to inactivity, your current work session is about to expire. For your security, Cloudera Manager sessions automatically end after 30 minutes of inactivity.

Your current session will expire in **1 minute**.
Press any key or click anywhere to continue.

If the user does not click the mouse or press a key, the user is logged out of the session and the following message appears:

Automatic Log Out Due to Inactivity

You are now logged out of your account.

We hadn't heard from you for about 30 minute(s), so for your security Cloudera Manager automatically logged you out of your account. Log back in below to continue.

 Remember me

Starting, Stopping, and Restarting the Cloudera Manager Server

To start the Cloudera Manager Server:

```
sudo service cloudera-scm-server start
```

You can stop (for example, to perform maintenance on its host) or restart the Cloudera Manager Server without affecting the other services running on your cluster. Statistics data used by activity monitoring and service monitoring will continue to be collected during the time the server is down.

To stop the Cloudera Manager Server:

```
sudo service cloudera-scm-server stop
```

To restart the Cloudera Manager Server:

```
sudo service cloudera-scm-server restart
```

Configuring Cloudera Manager

From the Administration menu you can select options for configuring settings that affect how Cloudera Manager interacts with your clusters.

Settings

The Settings page provides a number of categories as follows:

- Performance - Set the Cloudera Manager Agent heartbeat interval.
- Advanced - Enable API debugging and other advanced options.
- Monitoring - Set Agent health status parameters. For configuration instructions, see the topic *Configuring Cloudera Manager Agents*.
- Security - Set TLS encryption settings to enable TLS encryption between the Cloudera Manager Server, Agents, and clients. For configuration instructions, see [Configuring TLS Encryption for Cloudera Manager Using Auto-TLS](#). You can also:
 - Set the realm for Kerberos security and point to a custom keytab retrieval script.
 - Specify session timeout and a "Remember Me" option.
- Ports and Addresses - Set ports for the Cloudera Manager Admin Console and Server.
- Other
 - Enable Cloudera usage data collection For configuration instructions, see *Managing Anonymous Usage Data Collection*.
 - Set a custom header color and banner text for the Admin console.
 - Set an "Information Assurance Policy" statement – this statement will be presented to every user before they are allowed to access the login dialog box. The user must click "I Agree" in order to proceed to the login dialog box.
 - Disable/enable the auto-search for the Events panel at the bottom of a page.
- Support
 - Configure diagnostic data collection properties. See *Diagnostic Data Collection*.
 - Configure how to access Cloudera Manager help documentation.
- External Authentication - Specify the configuration to use LDAP, Active Directory, or an external program for authentication.
- Parcels - Configure settings for parcels, including the location of remote repositories that should be made available for download, and other settings such as the frequency with which Cloudera Manager will check for new

parcels, limits on the number of downloads or concurrent distribution uploads. See [Overview of Parcels](#) on page 74 for more information.

- Network - Configure proxy server settings.
- Custom Service Descriptors - Configure custom service descriptor properties for [Cloudera Manager Add-on Services](#).

You can also configure the following:

- Alerts
- Users
- Kerberos
- License

See [Managing Licenses](#) on page 83.

- Language

You can change the language of the Cloudera Manager Admin Console User Interface through the language preference in your browser. Information on how to do this for the browsers supported by Cloudera Manager is shown under the Administration page. You can also change the language for the information provided with activity and health events, and for alert email messages by selecting Language, selecting the language you want from the drop-down list on this page, then clicking Save Changes.

Related Information

[Managing Anonymous Usage Data Collection](#)

[Diagnostic Data Collection](#)

[Alerts](#)

Configuring Cloudera Manager Server Ports

You can specify the ports used to access the Cloudera Manager Server using the Admin Console. You can also specify the port used by agents to connect to the Server.

About this task

Minimum Required Role: [Full Administrator](#). This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

1. Select AdministrationSettings.
2. Under the Ports and Addresses category, set the following options as described below:

Setting	Description
HTTP Port for Admin Console	Specify the HTTP port to use to access the Server using the Admin Console.
HTTPS Port for Admin Console	Specify the HTTPS port to use to access the Server using the Admin Console.
Agent Port to connect to Server	Specify the port for Agents to use to connect to the Server.

3. Click Save Changes.
4. Restart the Cloudera Manager Server.

Configuring Network Settings for a Proxy Server

How to configure a proxy server for connections to Cloudera Manager.

Minimum Required Role: [Full Administrator](#). This feature is not available when using Cloudera Manager to manage Data Hub clusters.

To configure a proxy server through which data is downloaded to and uploaded from the Cloudera Manager Server, do the following:

1. Select AdministrationSettings.
2. Click the Network category.
3. Configure proxy properties.
4. Enter a Reason for change, and then click Save Changes to commit the changes.

Moving the Cloudera Manager Server to a New Host

You can move the Cloudera Manager Server if either the Cloudera Manager database server or a current backup of the Cloudera Manager database is available.

Procedure

1. Identify a new host on which to install Cloudera Manager.
2. Install Cloudera Manager on a new host, using the method described in the topic *Install the Cloudera Manager Server Packages*.



Important:

- The Cloudera Manager version on the destination host must match the version on the source host.
 - Do not install the other components, such as CDH and databases.
3. Copy the entire contents of the `/var/lib/cloudera-scm-server/` directory on the old host to that same path on the new host. Ensure you preserve permissions and all file contents.
 4. Copy the entire contents of the local parcel directory on the old host to that same path on the new host. Ensure you preserve permissions and all file contents. (The default location is `/opt/cloudera/parcel-repo` but this can be configured with the Local Parcel Repository Path configuration property, under AdministrationSettings.)
 5. If the database server is not available:
 - a) Install the database packages on the host that will host the restored database. This could be the same host on which you have just installed Cloudera Manager or it could be a different host. If you used the embedded PostgreSQL database, install the PostgreSQL package as described in the topic *Managing the Embedded PostgreSQL Database*. If you used an external MySQL, PostgreSQL, or Oracle database, reinstall the database following the instructions in *Step 4: Install and Configure Databases*.
 - b) Restore the backed up databases to the new database installation.
 6. Update `/etc/cloudera-scm-server/db.properties` with the database name, database instance name, username, and password.
 7. Do the following on all cluster hosts:
 - a) In `/etc/cloudera-scm-agent/config.ini`, update the `server_host` property to the new hostname.
 - b) If you are replacing the Cloudera Manager database with a new database, and you are not using a backup of the original Cloudera Manager database, delete the `/var/lib/cloudera-scm-agent/cm_guid` file.
 - c) Restart the agent using the following command:

```
sudo service cloudera-scm-agent restart
```

8. Stop the Cloudera Manager server on the source host by running the following command:

```
service cloudera-scm-server stop
```

9. Copy any Custom Service Descriptor files for add-on services to the configured directory on the new Cloudera Manager host. The directory path is configured by going to AdministrationSettings and editing the Local Descriptor Repository Path property. The default value is `/opt/cloudera/csd`. See [Add-on Services](#).

10. Start the Cloudera Manager Server on the new (destination) host. Cloudera Manager should resume functioning as it did before the failure. Because you restored the database from the backup, the server should accept the running state of the Agents, meaning it will not terminate any running processes.

The process is similar with secure clusters, though files in `/etc/cloudera-scm-server` must be restored in addition to the database. See the *Security* documentation.

Migrating from the Cloudera Manager Embedded PostgreSQL Database Server to an External PostgreSQL Database

Cloudera Manager provides an embedded PostgreSQL database server for demonstration and proof of concept deployments when creating a cluster. To remind users that this embedded database is not suitable for production, Cloudera Manager displays the banner text: "You are running Cloudera Manager in non-production mode, which uses an embedded PostgreSQL database. Switch to using a supported external database before moving into production."

If, however, you have already used the embedded database, and you are unable to redeploy a fresh cluster, then you must migrate to an external PostgreSQL database.



Note: This procedure does not describe how to migrate to a database server other than PostgreSQL. Moving databases from one database server to a different type of database server is a complex process that requires modification of the schema and matching the data in the database tables to the new schema. It is strongly recommended that you engage with Cloudera Professional Services if you wish to perform a migration to an external database server other than PostgreSQL.

Prerequisites

Before migrating the Cloudera Manager embedded PostgreSQL database to an external PostgreSQL database, ensure that your setup meets the following conditions:

- The external PostgreSQL database server is running.
- The database server is configured to accept remote connections.
- The database server is configured to accept user logins using md5.
- No one has manually created any databases in the external database server for roles that will be migrated.



Note: To view a list of databases in the external database server (requires default superuser permission):

```
sudo -u postgres psql -l
```

- All health issues with your cluster have been resolved.

For details about configuring the database server, see the topic *Configuring and Starting the PostgreSQL Server*.



Important: Only perform the steps in *Configuring and Starting the PostgreSQL Server*. Do not proceed with the creation of databases as described in the subsequent section.

For large clusters, Cloudera recommends running your database server on a dedicated host. Engage Cloudera Professional Services or a certified database administrator to correctly tune your external database server.

Step 1: Identify Roles that Use the Embedded Database Server

Before you can migrate to another database server, you must first identify the databases using the embedded database server. When the Cloudera Manager Embedded Database server is initialized, it creates the Cloudera Manager database and databases for roles in the Management Services. The Installation Wizard (which runs automatically the first time you log in to Cloudera Manager) or Add Service action for a cluster creates additional databases for roles when run. It is in this context that you identify which roles are used in the embedded database server.

Procedure

1. Obtain and save the cloudera-scm superuser password from the embedded database server. You will need this password in subsequent steps:

```
head -1 /var/lib/cloudera-scm-server-db/data/generated_password.txt
```

2. Make a list of all services that are using the embedded database server. Then, after determining which services are not using the embedded database server, remove those services from the list. The scm database must remain in your list. Use the following table as a guide:

Table 2: Cloudera Manager Embedded Database Server Databases

Service	Role	Default Database Name	Default Username
Cloudera Manager Server		scm	scm
Cloudera Management Service	Activity Monitor	amon	amon
Hive	Hive Metastore Server	hive	hive
Hue	Hue Server	hue	7uu7uu7uhue
Cloudera Management Service	Navigator Audit Server	nav	nav
Cloudera Management Service	Navigator Metadata Server	navms	navms
Oozie	Oozie Server	oozie_oozie_server	oozie_oozie_server
Cloudera Management Service	Reports Manager	rman	rman
Sentry	Sentry Server	sentry	sentry

3. Verify which roles are using the embedded database. Roles using the embedded database server always use port 7432 (the default port for the embedded database) on the Cloudera Manager Server host.

For Cloudera Management Services:

- a. Select Cloudera Management Service > Configuration, and type "7432" in the Search field.
- b. Confirm that the hostname for the services being used is the same hostname used by the Cloudera Manager Server.



Note:

If any of the following fields contain the value "7432", then the service is using the embedded database:

- Activity Monitor
- Navigator Audit Server
- Navigator Metadata Server
- Reports Manager

For the Oozie Service:

- a. Select Oozie service > Configuration, and type "7432" in the Search field.
- b. Confirm that the hostname is the Cloudera Manager Server.

For Hive, Hue, and Sentry Services:

- a. Select the specific service > Configuration, and type "database host" in the Search field.
- b. Confirm that the hostname is the Cloudera Manager Server.
- c. In the Search field, type "database port" and confirm that the port is 7432.
- d. Repeat these steps for each of the services (Hive, Hue and Sentry).

4. Verify the database names in the embedded database server match the database names on your list (Step 2). Databases that exist on the database server and not used by their roles do not need to be migrated. This step is to confirm that your list is correct.



Note: Do not add the postgres, template0, or template1 databases to your list. These are used only by the PostgreSQL server.

```
psql -h localhost -p 7432 -U cloudera-scm -l
```

```
Password for user cloudera-scm: <password>
```

Name	Access	Owner	List of databases		
			Encoding	Collate	Ctype
amon		amon	UTF8	en_US.UTF8	en_US.U
TF8					
hive		hive	UTF8	en_US.UTF8	en_US.UT
F8					
hue		hue	UTF8	en_US.UTF8	en_US
.UTF8					
nav		nav	UTF8	en_US.UTF8	en_US.
UTF8					
navms		navms	UTF8	en_US.UTF8	en_US.U
TF8					
oozie_oozie_server		oozie_oozie_server	UTF8	en_US.UTF8	en_US.UT
F8					
postgres		cloudera-scm	UTF8	en_US.UTF8	en_US
.UTF8					
rman		rman	UTF8	en_US.UTF8	en_US.
UTF8					
scm		scm	UTF8	en_US.UTF8	en_US.U
TF8					
sentry		sentry	UTF8	en_US.UTF8	en_US.UT
F8					
template0		cloudera-scm	UTF8	en_US.UTF8	en_US
.UTF8	=c/"cloudera-scm"				
template1		cloudera-scm	UTF8	en_US.UTF8	en_US.UT
F8	=c/"cloudera-scm"				

(12 rows)

Results

You should now have a list of all roles and database names that use the embedded database server, and are ready to proceed with the migration of databases from the embedded database server to the external PostgreSQL database server.

What to do next

Proceed to Step 2: Migrate Databases from the Embedded Database Server to the External PostgreSQL Database Server.

Step 2: Migrate Databases from the Embedded Database Server to the External PostgreSQL Database Server

After you identify the roles that use the embedded database server, you can migrate from the Cloudera Manager embedded database server to the external PostgreSQL database server. When you migrate, you export the PostgreSQL user roles from the embedded database, import the PostgreSQL user roles into the external database, import the Cloudera Manager database on the external database server, and perform other tasks.

About this task

Minimum Required Role: [Operator](#) (also provided by Configurator, Cluster Administrator, Limited Cluster Administrator, Full Administrator)

While performing this procedure, ensure that the Cloudera Manager Agents remain running on all hosts. Unless otherwise specified, when prompted for a password use the cloudera-scm password.



Note: After completing this migration, you cannot delete the cloudera-scm postgres superuser unless you remove the access privileges for the migrated databases. Minimally, you should change the cloudera-scm postgres superuser password.

Procedure

1. In Cloudera Manager, stop the cluster services identified in the previous step as using the embedded database server. Be sure to stop the Cloudera Management Service as well. Also be sure to stop any services with dependencies on these services. The remaining CDH services will continue to run without downtime.



Note: If you do not stop the services from within Cloudera Manager before stopping Cloudera Manager Server from the command line, they will continue to run and maintain a network connection to the embedded database server. If this occurs, then the embedded database server will ignore any command line stop commands (Step 2) and require that you manually kill the process, which in turn causes the services to crash instead of stopping cleanly.

2. Navigate to Hosts > All Hosts, and make note of the number of roles assigned to hosts. Also take note whether or not they are in a commissioned state. You will need this information later to validate that your scm database was migrated correctly.
3. Stop the Cloudera Manager Server. To stop the server:

```
sudo service cloudera-scm-server stop
```

4. Obtain and save the embedded database superuser password (you will need this password in subsequent steps) from the generated_password.txt file:

```
head -1 /var/lib/cloudera-scm-server-db/data/generated_password.txt
```

5. Export the PostgreSQL user roles from the embedded database server to ensure the correct users, permissions, and passwords are preserved for database access. Passwords are exported as an md5sum and are not visible in plain text. To export the database user roles (you will need the cloudera-scm user password):

```
pg_dumpall -h localhost -p 7432 -U cloudera-scm -v --roles-only -f "/var/tmp/cloudera_user_roles.sql"
```

6. Edit /var/tmp/cloudera_user_roles.sql to remove any CREATE ROLE and ALTER ROLE commands for databases not in your list. Leave the entries for cloudera-scm untouched, because this user role is used during the database import.
7. Export the data from each of the databases on your list you created when you identified roles that use the embedded database server:

```
pg_dump -F c -h localhost -p 7432 -U cloudera-scm [database_name] > /var/tmp/[database_name]_db_backup-$(date +%m-%d-%Y).dump
```

The following is a sample data export command for the scm database:

```
pg_dump -F c -h localhost -p 7432 -U cloudera-scm scm > /var/tmp/scm_db_backup-$(date +%m-%d-%Y).dump
```

Password:

8. Stop and disable the embedded database server:

```
service cloudera-scm-server-db stop
chkconfig cloudera-scm-server-db off
```

Confirm that the embedded database server is stopped:

```
netstat -at | grep 7432
```

9. Back up the Cloudera Manager Server database configuration file:

```
cp /etc/cloudera-scm-server/db.properties /etc/cloudera-scm-server/db.pr
operties.embedded
```

10. Copy the file /var/tmp/cloudera_user_roles.sql and the database dump files from the embedded database server host to /var/tmp on the external database server host:

```
cd /var/tmp
scp cloudera_user_roles.sql *.dump <user>@<postgres-server>:/var/tmp
```

11. Import the PostgreSQL user roles into the external database server.

The external PostgreSQL database server superuser password is required to import the user roles. If the superuser role has been changed, you will be prompted for the username and password.



Note: Only run the command that applies to your context; do not execute both commands.

- To import users when using the default PostgreSQL superuser role:

```
sudo -u postgres psql -f /var/tmp/cloudera_user_roles.sql
```

- To import users when the superuser role has been changed:

```
psql -h <database-hostname> -p <database-port> -U <superuser> -f /var/tm
p/cloudera_user_roles.sql
```

For example:

```
psql -h pg-server.example.com -p 5432 -U postgres -f /var/tmp/cloudera_u
ser_roles.sql
```

```
Password for user postgres
```

12. Import the Cloudera Manager database on the external server. First copy the database dump files from the Cloudera Manager Server host to your external PostgreSQL database server, and then import the database data:



Note: To successfully run the `pg_restore` command, there must be an existing database on the database server to complete the connection; the existing database will not be modified. If the `-d <existing-database>` option is not included, then the `pg_restore` command will fail.

```
pg_restore -C -h <database-hostname> -p <database-port> -d <existing-database> -U cloudera-scm -v <data-file>
```

Repeat this import for each database.

The following example is for the scm database:

```
pg_restore -C -h pg-server.example.com -p 5432 -d postgres -U cloudera-scm -v /var/tmp/scm_server_db_backup-20180312.dump
```

```
pg_restore: connecting to database for restore
Password:
```

13. Update the Cloudera Manager Server database configuration file to use the external database server. Edit the `/etc/cloudera-scm-server/db.properties` file as follows:
- Update the `com.cloudera.cmf.db.host` value with the hostname and port number of the external database server.
 - Change the `com.cloudera.cmf.db.setupType` value from "EMBEDDED" to "EXTERNAL".
14. Start the Cloudera Manager Server and confirm it is working:

```
service cloudera-scm-server start
```

Note that if you start the Cloudera Manager GUI at this point, it may take up to five minutes after executing the start command before it becomes available.

In Cloudera Manager Server, navigate to Hosts > All Hosts and confirm the number of roles assigned to hosts (this number should match what you found in Step 2); also confirm that they are in a commissioned state that matches what you observed in Step 2.

15. Update the role configurations to use the external database hostname and port number. Only perform this task for services where the database has been migrated.
- For Cloudera Management Services:
 - Select Cloudera Management Service > Configuration, and type "7432" in the Search field.
 - Change any database hostname properties from the embedded database to the external database hostname and port number.
 - Click Save Changes.
 - For the Oozie Service:
 - Select Oozie service > Configuration, and type "7432" in the Search field.
 - Change any database hostname properties from the embedded database to the external database hostname and port number.
 - Click Save Changes.
 - For Hive, Hue, and Sentry Services:
 - Select the specific service > Configuration, and type "database host" in the Search field.
 - Change the hostname from the embedded database name to the external database hostname.
 - Click Save Changes.

16. Start the Cloudera Management Service and confirm that all management services are up and no health tests are failing.

17. Start all Services via the Cloudera Manager web UI. This should start all services that were stopped for the database migration. Confirm that all services are up and no health tests are failing.

18. On the embedded database server host, remove the embedded PostgreSQL database server:

a) Make a backup of the `/var/lib/cloudera-scm-server-db/data` directory:

```
tar czvf /var/tmp/embedded_db_data_backup-$(date +%m-%d-%Y).tgz /var/lib/cloudera-scm-server-db/data
```

b) Remove the embedded database package:

For RHEL/SLES:

```
rpm --erase cloudera-manager-server-db-2
```

For Debian/Ubuntu:

```
apt-get remove cloudera-manager-server-db-2
```

c) Delete the `/var/lib/cloudera-scm-server-db/data` directory.

Migrating from the Cloudera Manager External PostgreSQL Database Server to a MySQL/Oracle Database Server

Cloudera Manager provides an embedded PostgreSQL database server for demonstration and proof of concept deployments when creating a cluster. To remind users that this embedded database is not suitable for production, Cloudera Manager displays the banner text: "You are running Cloudera Manager in non-production mode, which uses an embedded PostgreSQL database. Switch to using a supported external database before moving into production."

If you have already used the embedded database, and you are unable to redeploy a fresh cluster, then you must migrate to an external PostgreSQL database.



Note: You can migrate to an external MySQL or Oracle database only after successfully migrating from the embedded PostgreSQL database server to the external PostgreSQL database server.

Prerequisites

Before migrating from the Cloudera Manager external PostgreSQL database to an external MySQL/Oracle database, ensure that your setup meets the following conditions:

- Configuration uses Cloudera Manager 5.15.0 or later on supported platforms.
- You must have a valid Cloudera Manager Enterprise license.
- If Cloudera Manager is secured, then you must import Kerberos account manager credentials and regenerate them.
- You must have a destination host installed with the supported database of choice (MySQL or Oracle). For details about installing and configuring MySQL for Cloudera, see the topic *Install and Configure MySQL for Cloudera Software*. For details about installing and configuring Oracle for Cloudera, see the topic *Install and Configure Oracle Database Software for Cloudera Software*.
- You have made configured target database hosts available.
- You have planned for cluster downtime during the migration process.
- You have a plan to follow service specific database migration instructions for services other than Cloudera Manager. Refer to the appropriate service migration documentation for your cluster setup.
- No one has manually created any databases in the external database server for roles that will be migrated.
- All health issues with your cluster are resolved.

For large clusters, Cloudera recommends running your database server on a dedicated host. Engage Cloudera Professional Services or a certified database administrator to correctly tune your external database server.

Migrate from the Cloudera Manager External PostgreSQL Database Server to a MySQL/Oracle Database Server

When you migrate from the Cloudera Manager External PostgreSQL database server to a MySQL or Oracle database server, you export the Cloudera Manager configuration, prepare the target database for Cloudera Manager, and complete other tasks.

About this task

Minimum Required Role: [Operator](#) (also provided by Configurator, Cluster Administrator, Limited Cluster Administrator, Full Administrator)

Procedure

1. Migrate from the embedded PostgreSQL database server to an external PostgreSQL database server as described in the topic *Migrating from the Cloudera Manager Embedded PostgreSQL Database Server to an External PostgreSQL Database*.



Important: Migrating directly from the Cloudera Manager embedded PostgreSQL to a MySQL/Oracle database is not supported. You must first migrate from the Cloudera Manager embedded PostgreSQL database server to the external PostgreSQL database server. After performing this migration, you can use this procedure to migrate from the external PostgreSQL database server to MySQL/Oracle database servers.

2. Export your Cloudera Manager Configuration. First, get the latest supported API version:

```
curl -u <admin_username>:<admin_password> "http://<cm_server_host>:7180/api/version"
```

```
curl -u <admin_username>:<admin_password> "http://<cm_server_host>:7180/api/<api_version> /cm/deployment" > <path_to_file>/cm-deployment.json
```

The following is an example of the API version command:

```
curl -u admin:admin "http://10.17.103.191:7180/api/v19/cm/deployment" > /root/cm-deployment.json
```



Note:

If you have Cloudera Manager with TLS for the Admin Console enabled, retrieve the certificate file and use curl with the `--cacert` option:

```
curl --cacert <certificate_file> -u admin:admin "https://<cm_server_host>:7183/api/version"
```

3. Preserve Cloudera Manager's GUID by running the following command to create a `/etc/cloudera-scm-server/uuid` file. On a host that has an agent, run:

```
sudo -u postgres psql -qtAX scm -c "select GUID from CM_VERSION" > uuid
```



Note: Check to confirm the name of your Cloudera Manager database in `/etc/cloudera-scm-server/db.properties`.

Then move the UUID file to Cloudera Manager server's `/etc/cloudera-scm-server` directory.

4. Stop the cluster and the Cloudera Management Services.

5. Stop the Cloudera Manager Server:

```
sudo service cloudera-scm-server stop
```



Note:

For RHEL/CentOS 7, use the `systemctl` option instead:

```
sudo service systemctl cloudera-scm-server stop
```

6. Prepare the target database for Cloudera Manager. For details, refer to the topics *Install and Configure MySQL for Cloudera Software* or *Install and Configure Oracle Database for Cloudera Software*.
7. The process directory (`/var/run/cloudera-scm-agent/process/`) must be cleaned out for all of the hosts that have agents running on them. The agent completes this cleanup with a server reboot. However, if a server reboot is not a viable option, use one of the following options to accomplish the same task.



Note: This "hard restart" works for all supported platforms except SLES 12.

- a. Stop the agent and supervisor:

```
sudo systemctl stop cloudera-scm-agent
```

- b. Confirm that the agent and supervisor process are stopped:

```
ps -ef | grep -i cmf-agent; ps -ef | grep -i supervisor
```

- c. Perform a clean start:

```
service cloudera-scm-agent next_start_clean
```

Alternatively, run the following command to view the start options available on your platform:

```
service cloudera-scm-agent clean_start
```

- d. Ensure that the process is empty:

```
ls -la /var/run/cloudera-scm-agent/process/
```

- Alternatively:

- a. Stop the agent and supervisor:

```
sudo systemctl stop cloudera-scm-agent
```

- b. Confirm that the agent and supervisor process are stopped:

```
ps -ef | grep -i cmf-agent; ps -ef | grep -i supervisor
```

- c. Move the existing `/var/run/cloudera-scm-agent/` directory:

```
mv /var/run/cloudera-scm-agent /var/run/cloudera-scm-agent-BU
```

The agent will recreate the directory. Delete the backed up copy after confirming that the migration was successful.

8. Start the Cloudera Manager server:

```
service cloudera-scm-server start
```

9. Log in to Cloudera Manager. Exit the installation wizard by clicking the product log in the upper-left corner to stop the wizard and return to the Cloudera Manager home page.

10. Upgrade the Cloudera Manager Enterprise License by navigating to Administration > Licenses and installing a valid Cloudera Manager license.
11. Restore the Cloudera Manager configuration:

```
curl -H "Content-Type: application/json" --upload-file <path_to_file>/cm-deployment.json -u <admin_username>:<admin_password> "http://<cm_server_host>:7180/api/<api_version>/cm/deployment?deleteCurrentDeployment=true"
```

The following example shows how to restore a Cloudera Manager configuration:

```
curl -H "Content-Type: application/json" --upload-file /root/cm-deployment.json -u admin:admin "http://172.31.113.146:7180/api/v19/cm/deployment?deleteCurrentDeployment=true"
```

12. Start the following: Cloudera Management Service, Host Monitor, and Services Monitor. Verify that all the services in the Cloudera Management Service started and are Healthy.
13. Select the Home > Status tab for the cluster(s) that you previously stopped, and in the Actions dropdown, select Start.

Managing Cloudera Manager Server Logs

You can use the Cloudera Manager Server logs to troubleshoot problems with Cloudera Manager .

Related Information

[Logs](#)

Viewing the Cloudera Manager Server Logs

To help you troubleshoot problems, you can view the Cloudera Manager Server log. You can view the logs in the **Logs** page or in specific pages for the log.

Procedure

1. In the left menu, click **Diagnostics** > **Logs**.
2. Next to **Sources**, select the Cloudera Manager Server checkbox and deselect the other options.
3. Adjust the search criteria and click **Search**.

What to do next

You can also view the raw Cloudera Manager Server log by logging in to the Cloudera Manager Server host and view the `/var/log/cloudera-scm-server/cloudera-scm-server.log` file.

Setting the Cloudera Manager Server Log Location

You can set the location of the Cloudera Manager Server log.

Procedure

1. Stop the Cloudera Manager Server:

```
sudo service cloudera-scm-server stop
```

2. Set the `CMF_VAR` environment variable in `/etc/default/cloudera-scm-server` to the new parent directory:

```
export CMF_VAR=/opt
```

3. Create `log/cloudera-scm_server` and run directories in the new parent directory and set the owner and group of all directories to `cloudera-scm`. For example, if the new parent directory is `/opt/`, do the following:

```
sudo su
cd /opt
mkdir log
chown cloudera-scm:cloudera-scm log
mkdir /opt/log/cloudera-scm-server
chown cloudera-scm:cloudera-scm log/cloudera-scm-server
mkdir run
chown cloudera-scm:cloudera-scm run
```

4. Restart the Cloudera Manager Server:

```
sudo service cloudera-scm-server start
```

Cloudera Manager Agents

The Cloudera Manager Agent is a Cloudera Manager component that works with the Cloudera Manager Server to manage the processes that map to role instances.

In a Cloudera Manager managed cluster, you can only start or stop role instance processes using Cloudera Manager. Cloudera Manager uses an open source process management tool called `supervisord`, that starts processes, takes care of redirecting log files, notifying of process failure, setting the effective user ID of the calling process to the right user, and so on. Cloudera Manager supports automatically restarting a crashed process. It will also flag a role instance with a bad health flag if its process crashes repeatedly right after start up.

The Agent is started by `init.d` at start-up. It, in turn, contacts the Cloudera Manager Server and determines which processes should be running. The Agent is monitored as part of Cloudera Manager's host monitoring. If the Agent stops heartbeating, the host is marked as having bad health.

One of the Agent's main responsibilities is to start and stop processes. When the Agent detects a new process from the Server heartbeat, the Agent creates a directory for it in `/var/run/cloudera-scm-agent` and unpacks the configuration. It then contacts `supervisord`, which starts the process.

cm_processes

To enable Cloudera Manager to run scripts in subdirectories of `/var/run/cloudera-scm-agent`, (because `/var/run` is mounted `noexec` in many Linux distributions), Cloudera Manager mounts a `tmpfs`, named `cm_processes`, for process subdirectories.

A `tmpfs` defaults to a max size of 50% of physical RAM but this space is not allocated until its used, and `tmpfs` is paged out to swap if there is memory pressure.

The lifecycle actions of `cm_processes` can be described by the following statements:

- Created when the Agent starts up for the first time with a new `supervisord` process.
- If it already exists without `noexec`, reused when the Agent is started using `start` and not recreated.
- Remounted if Agent is started using `clean_restart`.
- Unmounting and remounting cleans out the contents (since it is mounted as a `tmpfs`).
- Unmounted when the host is rebooted.
- Not unmounted when the Agent is stopped.

Related Information

[supervisord](#)

[tmpfs](#)

Starting, Stopping, and Restarting Cloudera Manager Agents

Starting Agents

To start Agents, the supervisord process, and all managed service processes, use the following command:

- Start

```
sudo systemctl start cloudera-scm-agent
```

Stopping and Restarting Agents

To stop or restart Agents while leaving the managed processes running, use one of the following commands:

- Stop

```
sudo systemctl stop cloudera-scm-agent
```

- Restart

```
sudo systemctl restart cloudera-scm-agent
```

Hard Stopping and Restarting Agents



Warning: The `hard_stop` and `hard_restart` commands stop all running managed service processes on the host(s) where the command is run.

To stop or restart Agents, the supervisord process, and all managed service processes, use one of the following commands:

- Hard Stop

RHEL 7, SLES 12, Ubuntu 16.04

```
sudo systemctl stop cloudera-scm-supervisord.service
```

RHEL 5 or 6, SLES 11, Debian 6 or 7, Ubuntu 12.04, 14.04

```
sudo systemctl stop cloudera-scm-agent
```

- Hard Restart

RHEL 7, SLES 12, Debian 8, Ubuntu 16.04 and higher

```
sudo systemctl stop cloudera-scm-supervisord.service
sudo systemctl restart cloudera-scm-agent
```

RHEL 5 or 6, SLES 11, Debian 6 or 7, Ubuntu 12.04, 14.04

```
sudo systemctl restart cloudera-scm-agent
```

Hard restart is useful for the following situations:

- You are upgrading Cloudera Manager and the supervisord code has changed between your current version and the new one. To properly do this upgrade you need to restart supervisor too.
- supervisord freezes and needs to be restarted.
- You want to clear out all running state pertaining to Cloudera Manager and managed services.

Checking Agent Status

To check the status of the Agent process, use the command:

```
sudo systemctl status cloudera-scm-agent
```

Configuring Cloudera Manager Agents

Minimum Required Role: [Full Administrator](#). This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Cloudera Manager Agents can be configured globally using properties you set in the Cloudera Manager Admin Console and by setting properties in Agent configuration files.

Configuring Agent Heartbeat and Health Status Options

You can configure the Cloudera Manager Agent heartbeat interval and timeouts to trigger changes in Agent health as follows:

1. Select AdministrationSettings.
2. Under the Performance category, set the following option:

Property	Description
Send Agent Heartbeat Every	The interval in seconds between each heartbeat that is sent from Cloudera Manager Agents to the Cloudera Manager Server. Default: 15 sec.

3. Under the Monitoring category, set the following options:

Property	Description
Set health status to Concerning if the Agent heartbeats fail	The number of missed consecutive heartbeats after which a Concerning health status is assigned to that Agent. Default: 5.
Set health status to Bad if the Agent heartbeats fail	The number of missed consecutive heartbeats after which a Bad health status is assigned to that Agent. Default: 10.

4. Click Save Changes.

Configuring the Host Parcel Directory



Important: If you modify the parcel directory location, make sure that all hosts use the same location. Using different locations on different hosts can cause unexpected problems.

To configure the location of distributed parcels:

1. Click Hosts in the top navigation bar.
2. Click the Configuration tab.
3. Select CategoryParcels.
4. Configure the value of the Parcel Directory property. The setting of the parcel_dir property in the Cloudera Manager Agent configuration file overrides this setting (see below).
5. Enter a Reason for change, and then click Save Changes to commit the changes.
6. [Restart](#) the Cloudera Manager Agent on all hosts.

Agent Configuration File

The Cloudera Manager Agent supports different types of configuration options in the `/etc/cloudera-scm-agent/config.ini` file. You must update the configuration on each host. After changing a property, restart the Agent:

```
sudo systemctl restart cloudera-scm-agent
```

Section	Property	Description
[General]	server_host, server_port, listening_port, listening_hostname, listening_ip	<p>Hostname and ports of the Cloudera Manager Server and Agent and IP address of the Agent.</p> <p>Also see Configuring Cloudera Manager Server Ports and Ports used by Cloudera Manager .</p> <p>The Cloudera Manager Agent configures its hostname automatically. You can also manually specify the hostname the Cloudera Manager Agent uses by updating the listening_hostname property. To manually specify the IP address the Cloudera Manager Agent uses, update the listening_ip property in the same file.</p> <p>To have a CNAME used throughout instead of the regular hostname, an Agent can be configured to use listening_hostname=CNAME. In this case, the CNAME should resolve to the same IP address as the IP address of the hostname on that machine. Users doing this will find that the host inspector will report problems, but the CNAME will be used in all configurations where that's appropriate. This practice is particularly useful for users who would like clients to use namenode.mycluster.company.com instead of machine1234.mycluster.company.com. In this case, namenode.mycluster would be a CNAME for machine1234.mycluster, and the generated client configurations (and internal configurations as well) would use the CNAME.</p>
	lib_dir	<p>Directory to store Cloudera Manager Agent state that persists across instances of the agent process and system reboots. The Agent UUID is stored here.</p> <p>Default: /var/lib/cloudera-scm-agent.</p>
	local_filesystem_whitelist	<p>The list of local filesystems that should always be monitored.</p> <p>Default: ext2,ext3,ext4.</p>
	log_file	<p>The path to the Agent log file. If the Agent is being started using the init.d script, /var/log/cloudera-scm-agent/cloudera-scm-agent.out will also have a small amount of output (from before logging is initialized).</p> <p>Default: /var/log/cloudera-scm-agent/cloudera-scm-agent.log.</p>
	max_collection_wait_seconds	<p>Maximum time to wait for all metric collectors to finish collecting data.</p> <p>Default: 10 sec.</p>
	metrics_url_timeout_seconds	<p>Maximum time to wait when connecting to a local role's web server to fetch metrics.</p> <p>Default: 30 sec.</p>
	parcel_dir	<p>Directory to store unpacked parcels.</p> <p>Default: /opt/cloudera/parcels.</p> <p>This property overrides the setting in Cloudera Manager. To use the recommended procedure, you must make sure that this property is commented out in each host config.ini file.</p>
	supervisord_port	<p>The supervisord port. A change takes effect the next time supervisord is restarted (not when the Agent is restarted).</p> <p>Default: 19001.</p>
	task_metrics_timeout_seconds	<p>Maximum time to wait when connecting to a local TaskTracker to fetch task attempt data.</p> <p>Default: 5 sec.</p>
[Security]	use_tls,verify_cert_file, client_key_file, client_keypw_file, client_cert_file	<p>Security-related configuration.</p> <p>See</p> <ul style="list-style-type: none"> • Configuring TLS Encryption for Cloudera Manager and CDH Using Auto-TLS • Adding a Host to a Cluster on page 11
[Cloudera]	mgmt_home	<p>Directory to store Cloudera Management Service files.</p> <p>Default: /usr/share/cmfd.</p>

Section	Property	Description
[JDBC]	cloudera_mysql_connector_jar, cloudera_oracle_connector_jar, cloudera_postgresql_jdbc_jar	Location of JDBC drivers. Default: <ul style="list-style-type: none"> MySQL - /usr/share/java/mysql-connector-java.jar Oracle - /usr/share/java/oracle-connector-java.jar PostgreSQL - /usr/share/cmfd/lib/postgresql-version-build.jdbc4.jar

Managing the Cloudera Manager Agent Logs

To help you troubleshoot problems, you can view the Cloudera Manager Agent logs. You can view the logs in the Logs page or in specific pages for the logs.

Viewing the Cloudera Manager Agent Logs

Use the procedure to view and search the logs from all Cloudera Manager agents managed by this instance of Cloudera Manager.

Procedure

1. In the left menu, click DiagnosticsLogs.
2. Click Select Sources to display the log source list.
3. Uncheck the All Sources checkbox.
4. Click ► to the left of Cloudera Manager and select the Agent checkbox.
5. Click Search.

What to do next

You can also view the Cloudera Manager Agent log at /var/log/cloudera-scm-agent/cloudera-scm-agent.log on the Agent hosts.

Setting the Cloudera Manager Agent Log Location

By default, the Cloudera Manager Agent log is stored in /var/log/cloudera-scm-agent/. If there is not enough space in that directory, you can change the location of the log file.

Procedure

1. Set the log_file property in the Cloudera Manager Agent configuration file:

```
log_file=/opt/log/cloudera-scm-agent/cloudera-scm-agent.log
```

2. Create log/cloudera-scm_agent directories and set the owner and group to cloudera-scm. For example, if the log is stored in /opt/log/cloudera-scm-agent, do the following:

```
sudo su
cd /opt
mkdir log
chown cloudera-scm:cloudera-scm log
mkdir /opt/log/cloudera-scm-agent
chown cloudera-scm:cloudera-scm log/cloudera-scm-agent
```

3. Restart the Agent:

```
sudo service cloudera-scm-agent restart
```

Overview of Parcels

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

A *parcel* is a binary distribution format containing the program files, along with additional metadata used by Cloudera Manager. The important differences between parcels and packages are:

- Parcels are self-contained and installed in a versioned directory, which means that multiple versions of a given parcel can be installed side-by-side. You can then designate one of these installed versions as the active one. With packages, only one package can be installed at a time so there is no distinction between what is installed and what is active.
- Parcels are required for rolling upgrades.
- You can install parcels at any location in the filesystem. They are installed by default in `/opt/cloudera/parcels`. In contrast, packages are installed in `/usr/lib`.
- When you install from the Parcels page, Cloudera Manager automatically downloads, distributes, and activates the correct parcel for the operating system running on each host in the cluster. All hosts that make up a logical cluster must run on the same major OS release to be covered by Cloudera Support. Cloudera Manager must run on the same major OS release as at least one of the clusters it manages, to be covered by Cloudera Support. The risk of issues caused by running different minor OS releases is considered lower than the risk of running different major OS releases. Cloudera recommends running the same minor release cross-cluster, because it simplifies issue tracking and supportability.



Important: Cloudera Manager manages parcels without the need for users to manipulate parcels in the filesystem. You might cause failures or unexpected behaviors in your cluster if you perform any of the following unsupported actions:

- Installing parcels within custom RPM packages and saving them to the Cloudera Manager parcel directory.
- Downloading parcels and manually placing them in the Cloudera Manager parcel directory.
- Manually adding, modifying, or deleting files within the root parcels directory or its subdirectories.

Parcels are available for CDH, for other managed services, and for Sqoop Connectors.

Advantages of Parcels

Because of their unique properties, parcels offer the following advantages over packages:

- Distribution of Cloudera Runtime as a single object - Instead of having a separate package for each component of Cloudera Runtime, parcels are distributed as a single object. This makes it easier to distribute software to a cluster that is not connected to the Internet.
- Internal consistency - All Cloudera Runtime components are matched, eliminating the possibility of installing components from different versions.
- Installation outside of `/usr` - In some environments, Hadoop administrators do not have privileges to install system packages. With parcels, administrators can install to `/opt`, or anywhere else.



Note: With parcels, the path to the Cloudera Runtime libraries is `/opt/cloudera/parcels/CDH/lib` instead of the usual `/usr/lib`. Do not link `/usr/lib/` elements to parcel-deployed paths, because the links can cause scripts that distinguish between the two paths to not work.

- Installation of Cloudera Runtime without `sudo` - Parcel installation is handled by the Cloudera Manager Agent running as root or another user, so you can install Cloudera Runtime without `sudo`.
- Decoupled distribution from activation - With side-by-side install capabilities, you can stage a new version of Cloudera Runtime across the cluster before switching to it. This allows the most time-consuming part of an upgrade to be done ahead of time without affecting cluster operations, thereby reducing downtime.

- Rolling upgrades - Using packages requires you to shut down the old process, upgrade the package, and then start the new process. Errors can be difficult to recover from, and upgrading requires extensive integration with the package management system to function seamlessly. With parcels, when a new version is staged side-by-side, you can switch to a new minor version by simply changing which version of Cloudera Runtime is used when restarting each process. You can then perform upgrades with rolling restarts, in which service roles are restarted in the correct order to switch to the new version with minimal service interruption. Your cluster can continue to run on the existing installed components while you stage a new version across your cluster, without impacting your current operations. Major version upgrades (for example, CDH 5 to CDH 6) require full service restarts because of substantial changes between the versions. Finally, you can upgrade individual parcels or multiple parcels at the same time.
- Upgrade management - Cloudera Manager manages all the steps in a CDH or Cloudera Runtime version upgrade. With packages, Cloudera Manager only helps with initial installation.
- Additional components - Parcels are not limited to Cloudera Runtime. Add-on service parcels are also available.
- Compatibility with other distribution tools - Cloudera Manager works with other tools you use for download and distribution, such as Puppet. Or, you can download the parcel to Cloudera Manager Server manually if your cluster has no Internet connectivity and then have Cloudera Manager distribute the parcel to the cluster.

Parcel Life Cycle

To enable upgrades and additions with minimal disruption, parcels have following phases:

- Downloaded - The parcel software is copied to a local parcel directory on the Cloudera Manager Server, where it is available for distribution to other hosts in any of the clusters managed by this Cloudera Manager Server. You can have multiple parcels for a product downloaded to your Cloudera Manager Server. After a parcel has been downloaded to the Server, it is available for distribution on all clusters managed by the Server. A downloaded parcel appears in the cluster-specific section for every cluster managed by this Cloudera Manager Server.
- Distributed - The parcel is copied to the cluster hosts, and components of the parcel are unpacked. Distributing a parcel does not upgrade the components running on your cluster; the current services continue to run unchanged. You can have multiple parcels distributed on your cluster. Distributing parcels does not require Internet access; the Cloudera Manager Agent on each cluster member downloads the parcels from the local parcel repository on the Cloudera Manager Server.
- Activated - Links to the parcel components are created. Activation does not automatically stop the current services or perform a restart. You can restart services after activation, or the system administrator can determine when to perform those operations.
- In Use - The parcel components on the cluster hosts are in use when you start or restart the services that use those components.
- Deactivated - The links to the parcel components are removed from the cluster hosts.
- Removed - The parcel components are removed from the cluster hosts.
- Deleted - The parcel is deleted from the local parcel repository on the Cloudera Manager Server.

Cloudera Manager detects when new parcels are available. You can configure Cloudera Manager to download and distribute parcels automatically. .

Parcel Locations

The default location for the local parcel directory on the Cloudera Manager Server is `/opt/cloudera/parcel-repo`. To change this location, follow the instructions in [Configuring Cloudera Manager Server Parcel Settings](#) on page 81.

The default location for the distributed parcels on managed hosts is `/opt/cloudera/parcels`. To change this location, set the `parcel_dir` property in `/etc/cloudera-scm-agent/config.ini` file of the Cloudera Manager Agent and restart the Cloudera Manager Agent or by following the instructions in [Configuring the Host Parcel Directory](#) on page 82.

Managing Parcels

Procedures for managing Parcels.

On the Parcels page in Cloudera Manager, you can manage parcel installation and activation and determine which parcel versions are running across your clusters. The Parcels page displays a list of parcels managed by Cloudera Manager. Cloudera Manager displays the name, version, and status of each parcel and provides available actions on the parcel.

Accessing the Parcels Page

Minimum Required Role: [Configurator](#) (also provided by Cluster Administrator, Limited Cluster Administrator, and Full Administrator)

Access the Parcels page by doing one of the following:

- Click the parcel icon in the top navigation bar.
- Click the Hosts in the top navigation bar, then the Parcels tab.

Use the selectors on the left side of the console to filter the displayed parcels:

- Location selector - View only parcels that are available remotely, only parcels pertaining to a particular cluster, or parcels pertaining to all clusters. When you access the Parcels page, the selector is set to Available Remotely.
- Error Status section of the Filters selector - Limit the list of displayed parcels by error status.
- Parcel Name section of the Filters selector - Limit the list of displayed parcels by parcel name.
- Status section of the Filters selector - Limit the list to parcels that have been distributed, parcels that have not been distributed (Other), or all parcels.

When you download a parcel, it appears in the list for each cluster managed by Cloudera Manager, indicating that the parcel is available for distribution on those clusters. Only one copy of the downloaded parcel resides on the Cloudera Manager Server. After you distribute the parcel, Cloudera Manager copies the parcel to the hosts in that cluster.

For example, if Cloudera Manager is managing two clusters, the rows in the All Clusters page list the information about the parcels on the two clusters. The Status column displays the current status of the parcels. The Version column displays version information about the parcel. Click the information icon to view the release notes for the parcel. The Actions column shows actions you can perform on the parcels, such as download, distribute, delete, deactivate, and remove from host.

The screenshot shows the Cloudera Manager interface for the Parcels page. The top navigation bar includes 'Clusters', 'Hosts', 'Diagnostics', 'Audits', 'Charts', 'Backup', and 'Administration'. The 'Hosts' tab is active. The page title is 'Parcels'. On the left, there are filters for 'Location' (Cluster 1, Cluster 2, All Clusters, Available Remotely) and 'Filters' (PARCEL NAME, STATUS). The main content area displays two tables of parcels for 'Cluster 1' and 'Cluster 2'.

Cluster	Parcel Name	Version	Status	Actions
Cluster 1	CDH 5	5.5.7-1.cdh5.5.7.p0.280	Downloaded	Distribute
		5.12.0-1.cdh5.12.0.p0.319	Available Remotely	Download
Cluster 2	CDH 5	5.5.7-1.cdh5.5.7.p0.280	Distributed, Activated	Deactivate
		5.12.0-1.cdh5.12.0.p0.319	Available Remotely	Download

Downloading a Parcel

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

1. Go to the Parcels page. In the Location selector, click *ClusterName* or Available Remotely. Parcels that are available for download display the Available Remotely status and a Download button.

If the parcel you want is not shown here—for example, you want to upgrade to a version of CDH that is not the most current version—you can make additional remote parcel repositories available. You can also configure the location of the local parcel repository and other settings. See [Parcel Configuration Settings](#) on page 81.

If a parcel version is too new to be supported by the Cloudera Manager version, the parcel appears with a red background and error message:

CDH 5	5.5.0-1.cdh5.5.0.p0.871	Available Remotely
<ul style="list-style-type: none"> • Local parcel error for parcel CDH-5.5.0-1.cdh5.5.0.p0.871-el6.parcel : The version 5.5.0-1.cdh5.5.0.p0.871 is too new to be supported. 		

Such parcels are also listed when you select the Error status in the Error Status section of the Filters selector.

2. Click the Download button of the parcel you want to download to your local repository. The status changes to Downloading.

After a parcel has been downloaded, it is removed from the Available Remotely page.

Distributing a Parcel

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Downloaded parcels can be distributed to the hosts in your cluster and made available for activation. Parcels are downloaded to the Cloudera Manager Server, so with multiple clusters, the downloaded parcels are shown as available to all clusters managed by the Cloudera Manager Server. However, you select distribution to a specific cluster's hosts on a cluster-by-cluster basis.

1. From the Parcels page, in the Location selector, select the cluster where you want to distribute the parcel, or select All Clusters. (The first cluster in the list is selected by default when you open the Parcels page.)
2. Click Distribute for the parcel you want to distribute. The status changes to Distributing. During distribution, you can:
 - Click the Details link in the Status column to view the Parcel Distribution Status page.
 - Click Cancel to cancel the distribution. When the Distribute action completes, the button changes to Activate, and you can click the Distributed status link to view the status page.

Distribution does not require Internet access; the Cloudera Manager Agent on each cluster member downloads the parcel from the local parcel repository hosted on the Cloudera Manager Server.

If you have a large number of hosts to which parcels must be distributed, you can control how many concurrent uploads Cloudera Manager performs. See [Parcel Configuration Settings](#) on page 81.

To delete a parcel that is ready to be distributed, click the triangle at the right end of the Distribute button and select Delete. This deletes the parcel from the local parcel repository.

Distributing parcels to the hosts in the cluster does not affect the current running services.

Activating a Parcel

Parcels that have been distributed to the hosts in a cluster are ready to be activated.

1. From the Parcels page, in the Location selector, choose *ClusterName* or All Clusters, and click the Activate button for the parcel you want to activate. This updates Cloudera Manager to point to the new software, which is ready to run the next time a service is restarted. A pop-up indicates which services must be restarted to use the new parcel.

2. Choose one of the following:

- Restart - Activate the parcel and restart services affected by the new parcel.
- Activate Only - Active the parcel. You can restart services at a time that is convenient. If you do not restart services as part of the activation process, you must restart them at a later time. Until you restart services, the current parcel continues to run.

3. Click OK.

Activating a new parcel also deactivates the previously active parcel for the product you just upgraded. However, until you restart the services, the previously active parcel displays a status of Still in use because the services are using that parcel, and you cannot remove the parcel until it is no longer being used.

If the parcel you activate updates the software for only a subset of services, even if you restart all of that subset, the previously active parcel displays Still in use until you restart the remaining services. For example, if you are running HDFS, YARN, Oozie, Hue, Impala, and Spark services, and you activate a parcel that updates only the Oozie service, the pop-up that displays instructs you to restart only the Oozie and Hue services. Because the older parcel is still in use by the HDFS, YARN, Impala, and Spark services, the parcel page shows that parcel as Still in use until you restart these remaining services.

Deactivating a Parcel

You can deactivate an active parcel; this updates Cloudera Manager to point to the previous software version, which is ready to run the next time a service is restarted. From the Parcels page, choose *ClusterName* or All Clusters in the Location selector, and click the Deactivate button on an activated parcel.

To use the previous version of the software, restart your services.




Important: If you originally installed from parcels, and one version of the software is installed (that is, no packages, and no previous parcels have been activated and started), when you attempt to restart after deactivating the current version, your roles will be stopped and will not be able to restart.

Removing a Parcel

From the Parcels page, in the Location selector, choose *ClusterName* or All Clusters, click the  to the right of an Activate button, and select Remove from Hosts.

Deleting a Parcel

From the Parcels page, in the Location selector, choose *ClusterName* or All Clusters, and click the  to the right of a Distribute button, and select Delete.



Warning:

Note: Do not remove the Cloudera Runtime parcel entirely from Cloudera Manager unless it is no longer in use on any other cluster managed by this instance of Cloudera Manager. To remove a parcel from specific managed hosts, select the Remove from Hosts option instead of the Delete option.

Changing the Parcel Directory

The default location of the parcel directory is `/opt/cloudera/parcels`. To relocate distributed parcels to a different directory, do the following:

1. Stop all services.
2. **Deactivate** all in-use parcels.
3. **Shut down** the Cloudera Manager Agent on all hosts.
4. Move the existing parcels to the new location.
5. **Configure** the host parcel directory.
6. **Start** the Cloudera Manager Agents.

7. [Activate](#) the parcels.
8. Start all services.

Troubleshooting

If you experience an error while performing parcel operations, click the red 'X' icons on the parcel page to display a message that identifies the source of the error.

If a parcel is being distributed but never completes, make sure you have enough free space in the [parcel download directories](#), because Cloudera Manager will try to download and unpack parcels even if there is insufficient space.

Viewing Parcel Usage

The Parcel Usage page shows parcels in current use in your clusters. In a large deployment, this makes it easier to keep track of different versions installed across the cluster, especially if some hosts were not available when you performed an installation or upgrade, or were added later. To display the Parcel Usage page:

1. Do one of the following:
 - Click the parcel icon in the top navigation bar.
 - Click Hosts in the top navigation bar and click the Parcels tab.
2. Click the Parcel Usage button.

This page only shows the usage of parcels, not components that were installed as packages. If you select a cluster running packages, the cluster is not displayed, and instead you see a message indicating the cluster is not running parcels.

The screenshot displays the Cloudera Manager interface for viewing parcel usage. At the top, there is a navigation bar with tabs for 'Hosts', 'Status', 'Configuration', 'Templates', 'Disks Overview', and 'Parcels'. The 'Parcels' tab is currently selected. Below the navigation bar, the 'Parcel Usage' section is shown. It features a 'Product' dropdown menu set to 'Cluster 1' and another dropdown menu set to 'CDH'. A legend below these menus shows a checked checkbox next to a blue square icon, representing 'CDH 5.1.0-1.cdh5.1.0.p0.460 (Active, 4)'. Below the legend, there are two unchecked checkboxes: 'No CDH processes running on this host' and 'Multiple product versions running on a single host'. To the right of the legend, a section titled 'Hosts with CDH processes running' displays a grid of four blue squares, indicating that four hosts in the cluster are running the specified CDH processes.

You can view parcel usage by cluster or by product.

You can also view just the hosts running only the active parcels, or just hosts running older parcels (not the currently active parcels), or both.

The host map at the right shows each host in the cluster, with the status of the parcels on that host. If the host is running the processes from the currently activated parcels, the host is indicated in blue. A black square indicates that a parcel has been activated, but that all the running processes are from an earlier version of the software. This occurs, for example, if you have not restarted a service or role after activating a new parcel. If you have individual hosts running components installed as packages, the square is empty.

Move the cursor over the grid icon to see the rack to which the hosts are assigned. Hosts on different racks are displayed in separate rows.

To view the exact versions of the software running on a given host, click the square representing the host. This displays the parcel versions installed on that host.

Hosts Status Configuration Templates Disks Overview Parcels

Parcel Usage

The screenshot shows the Cloudera Manager interface for Parcel Usage. On the left, the 'Product' dropdown is set to 'Cluster 1' and 'CDH'. Below it, there are checkboxes for 'CDH 5.1.0-1.cdh5...' (checked), 'No CDH processes running', and 'Multiple product versions ru...'. On the right, under 'Hosts with CDH processes running', there are four square icons. The first icon is a four-square grid, and a pop-up window is open over it. The pop-up displays the host name tcdn501-1.ent.cloudera.com and the title 'Versions used by running roles'. It lists 'CDH 5.1.0-1.cdh5.1.0.p0.460' with the status 'Active'. Below this, there are links for various roles: [Hive Metastore Server](#), [HiveServer2](#), [JobHistory Server](#), [NameNode](#), [Oozie Server](#), [ResourceManager](#), [SecondaryNameNode](#), and [Server](#). At the bottom of the pop-up, it says 'Other products in use by host'.

The pop-up lists the roles running on the selected host that are part of the listed parcel. Clicking a role opens the Cloudera Manager page for that role. It also shows whether the parcel is active or not.

If a host is running various software versions, the square representing the host is a four-square grid icon. When you move the cursor over that host, both the active and inactive components are shown. For example, in the image below, the older CDH parcel has been deactivated, but only the HDFS service has been restarted.

Hosts Status Configuration Templates Disks Overview **Parcels**

Parcel Usage

Product

Cluster 1

CDH

CDH 5.1.0-1.cdh5.0.1.p0.460

CDH 5.0.1-1.cdh5.0.1.p0.47

No CDH processes running

Multiple product versions running

Hosts with CDH processes running

tcdn501-1.ent.cloudera.com

Versions used by running roles

CDH 5.0.1-1.cdh5.0.1.p0.47 Inactive

[Hive Metastore Server](#) [HiveServer2](#) [Hue Server](#) [JobHistory Server](#)
[Oozie Server](#) [ResourceManager](#) [Server](#) [Sqoop 2 Server](#)

CDH 5.1.0-1.cdh5.1.0.p0.460 Active

[NameNode](#) [SecondaryNameNode](#)

Other products in use by host

Parcel Configuration Settings

You can configure where parcels are stored on the Cloudera Manager Server host, the URLs of parcel repositories, the properties of a proxy server through which parcels are downloaded, and where parcels distributed to cluster hosts are stored.

Configuring Cloudera Manager Server Parcel Settings

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

1. Use one of the following methods to open the parcel settings page:

- Navigation bar
 - a. Click the parcel icon in the top navigation bar or click Hosts and click the Parcels tab.
 - b. Click the Configuration button.
- Menu
 - a. Select Administration Settings .
 - b. Select Category Parcels .

2. Specify a property:

- Local Parcel Repository Path defines the path on the Cloudera Manager Server host where downloaded parcels are stored.
- Remote Parcel Repository URLs is a list of repositories that Cloudera Manager checks for parcels. Initially this points to the latest released CDH 5 and CDH 6 repositories, but you can add your own repository locations to

the list. Use this mechanism to add Cloudera repositories that are not listed by default, such as older versions of CDH. You can also use this to add your own custom repositories. The locations of the Cloudera parcel repositories are `https://archive.cloudera.com/product/parcels/version` , where *product* is a product name and *version* is a specific product version, latest, or the substitution variable {latest_supported}. The substitution variable appears after the parcel for the CDH version with the same major number as the Cloudera Manager version to enable substitution of the latest supported maintenance version of CDH.

To add a parcel repository:

- a. In the Remote Parcel Repository URLs list, click the addition symbol to open an additional row.
 - b. Enter the path to the repository.
3. Click Save Changes.

You can also:

- Set the frequency with which Cloudera Manager checks for new parcels.
- Configure a proxy to access to the remote repositories.
- Configure whether downloads and distribution of parcels should occur automatically when new ones are detected. If automatic downloading and distribution are not enabled (the default), go to the Parcels page to initiate these actions.
- Control which products can be downloaded if automatic downloading is enabled.
- Control whether to retain downloaded parcels.
- Control whether to retain old parcel versions and how many parcel versions to retain

You can tune the parcel distribution load on your network by configuring the bandwidth limits and the number of concurrent uploads. The defaults are up to 50 MiB/s aggregate bandwidth and 50 concurrent parcel uploads.

- Theoretically, the concurrent upload count (Maximum Parcel Uploads) is unimportant if all hosts have the same speed Ethernet. Fifty concurrent uploads is acceptable in most cases. However, if the server has more bandwidth (for example, 10 GbE, and the normal hosts are using 1 GbE), then the count is important to maximize bandwidth. It should be at least the difference in speeds (10x in this case).
- The bandwidth limit (Parcel Distribution Rate Limit) should be your Ethernet speed (in MiB/seconds) divided by approximately 16. You can use a higher limit if you have QoS configured to prevent starving other services, or if you can accept the risk associated with higher bandwidth load.

Configuring a Proxy Server

To configure a proxy server through which data and parcels are downloaded to and uploaded from the Cloudera Manager Server, do the following:

1. Select AdministrationSettings.
2. Click the Network category.
3. Configure proxy properties.
4. Enter a Reason for change, and then click Save Changes to commit the changes.

Configuring the Host Parcel Directory



Important: If you modify the parcel directory location, make sure that all hosts use the same location. Using different locations on different hosts can cause unexpected problems.

To configure the location of distributed parcels:

1. Click Hosts in the top navigation bar.
2. Click the Configuration tab.
3. Select Category Parcels .
4. Configure the value of the Parcel Directory property. The setting of the `parcel_dir` property in the Cloudera Manager agent configuration file overrides this setting.
5. Enter a Reason for change, and then click Save Changes to commit the changes.
6. Restart the Cloudera Manager Agent on all hosts.

Configuring Peer-to-Peer Distribution of Parcels

Cloudera Manager uses a peer-to-peer service to efficiently distribute parcels to cluster hosts. The service is enabled by default and is configured to run on port 7191. You can change this port number, and you can disable peer-to-peer distribution.

To modify peer-to-peer distribution of parcels:

1. Open Cloudera Manager and select **Hosts All Hosts Configuration**.
2. Change the value of the P2P Parcel Distribution Port property to the new port number.
Set the value to 0 to disable peer-to-peer distribution of parcels.
3. Enter a Reason for change, and then click **Save Changes** to commit the changes.

Managing Licenses

When you install Cloudera Manager, you can either upload your CDP Private Cloud Base license or select a 60-day trial version.

CDP Private Cloud Base offers the following two types of licenses:

- CDP Private Cloud Base Edition

When the license expires, you will no longer be able to access the Cloudera Manager Admin console until you upload a valid license. However, your clusters will continue to function and your data will remain intact. To obtain a CDP Private Cloud Base license, fill in the *Contact Us* form or call 866-843-7207

- CDP Private Cloud Base Edition Trial

The CDP Private Cloud Base Edition Trial is a free 60-day trial that does not require a license file. When the 60-day trial period expires, you will no longer be able to access the Cloudera Manager Admin console, though your clusters and data remain intact. You can obtain a CDP Private Cloud Base license to regain access to the Admin Console. To obtain a CDP Private Cloud Base license, fill in the *Contact Us* form or call 866-843-7207.

You can use a trial license only once; when the 60-day trial period expires or you have ended the trial, you cannot restart the trial.



Note: The only version available for trial installations is 7.1.1.

Accessing the License Page

To access the license page, click **AdministrationLicense**.

Minimum Required Role: **Cluster Administrator** (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

If you have a license installed, the license page indicates its status (for example, whether your license is currently valid) and displays the license details: the license owner, the license key, the license start date, the expiration date, and the deactivation date. Typically the expiration date is the same as the deactivation date, at which point the Admin Console is no longer accessible. If the license expires, your clusters and data are unaffected.

Ending a CDP Private Cloud Base Trial

If you are using the trial edition, the License page indicates when your license will expire. However, you can end the trial at any time (prior to expiration) as follows:

About this task

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

1. On the **License** page, click End Trial.
2. Confirm that you want to end the trial.
3. Restart the Cloudera Management Service, HBase, HDFS, and Hive services to pick up configuration changes.

Upgrading from a CDP Private Cloud Base Trial to CDP Private Cloud Base

You can upgrade your license from a trial license to a CDP Private Cloud Base license.

About this task

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

1. Purchase a CDP Private Cloud Base Edition license from Cloudera.
2. On the License page, click Update License.
3. Click the Select License File field.
4. Browse to the location of your license file, click the file, and click Open.
5. Click Upload.

Renewing a License

You can upload a license file to renew a CDP Private Cloud Base license.

Before you begin

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

1. Download the license file and save it locally.
2. In Cloudera Manager, go to the Home page.
3. Select AdministrationLicense.
4. Click Update License.
5. Browse to the license file you downloaded.
6. Click Upload.

Default User Roles

By default, Cloudera Manager ships with user roles that have privileges for all clusters managed by Cloudera Manager.

The following table describes the actions each user role can perform:

Permitted Operations	Auditor	Cluster Administrator	Configurator	Dashboard User	Full Administrator	Key Administrator	Limited Operator	Navigator Administrator	Operator	Read-Only	Replication Administrator	User Administrator
Apply policies to redact sensitive data		Y			Y							
Administer Cloudera Navigator					Y			Y				
Create, modify, and delete your own dashboards				Y	Y							
Manage user accounts and configuration of external authentication					Y							Y
See available hosts		Y			Y							
View and perform parcels operations		Y			Y							
Enter and exit Maintenance Mode		Y	Y		Y							
Edit the configuration of services and roles		Y	Y		Y							
View data in Cloudera Manager	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y
Start, stop, and restart KMS		Y	Y		Y	Y			Y			
Manage Full Administrator accounts					Y							
Decommission hosts		Y	Y		Y		Y		Y			
Create clusters		Y			Y							
View audit events	Y				Y			Y				
Create, update, or delete external account configuration					Y							Y
Configure HDFS Encryption, administer Key Trustee Server, and manage encryption keys					Y	Y						
Recommission hosts, and decommission and recommission roles		Y	Y		Y				Y			
Access all functionality that Cloudera Manager offers		Y			Y							
Create replication policies and snapshot policies					Y						Y	
Start, stop, and restart most clusters, services, and roles		Y	Y		Y				Y			

Other Cloudera Manager Tasks and Settings

From the Administration tab you can select options for configuring settings that affect how Cloudera Manager interacts with your clusters.

Settings

The Settings page provides a number of categories as follows:

- Performance - Set the Cloudera Manager Agent heartbeat interval.
- Advanced - Enable API debugging and other advanced options.
- Monitoring - Set Agent health status parameters.
- Security - Set TLS encryption settings to enable TLS encryption between the Cloudera Manager Server, Agents, and clients. You can also:
 - Set the realm for Kerberos security and point to a custom keytab retrieval script.
 - Specify session timeout and a "Remember Me" option.
- Ports and Addresses - Set ports for the Cloudera Manager Admin Console and Server.
- Other
 - Enable Cloudera usage data collection.
 - Set a custom header color and banner text for the Admin console.
 - Set an "Information Assurance Policy" statement – this statement will be presented to every user before they are allowed to access the login dialog box. The user must click "I Agree" in order to proceed to the login dialog box.
 - Disable/enable the auto-search for the Events panel at the bottom of a page.
- Support
 - Configure diagnostic data collection properties.
 - Configure how to access Cloudera Manager help files.
- External Authentication - Specify the configuration to use LDAP, Active Directory, or an external program for authentication.

- **Parcels** - Configure settings for parcels, including the location of remote repositories that should be made available for download, and other settings such as the frequency with which Cloudera Manager will check for new parcels, limits on the number of downloads or concurrent distribution uploads. See [Parcels](#) for more information.
- **Network** - Configure proxy server settings.
- **Custom Service Descriptors** - Configure custom service descriptor properties for Add-on services.

Alerts

See *Managing Alerts*.

Users

See *Cloudera Manager User Accounts*.

Kerberos

See *Enabling Kerberos Authentication for Cloudera Runtime*.

License

See *Managing Licenses*.

User Interface Language

You can change the language of the Cloudera Manager Admin Console User Interface through the language preference in your browser. Information on how to do this for the browsers supported by Cloudera Manager is shown under the Administration page. You can also change the language for the information provided with activity and health events, and for alert email messages by selecting Language, selecting the language you want from the drop-down list on this page, then clicking Save Changes.

Peers

See *Designating a Replication Source*.

Cloudera Management Service

The Cloudera Management Service is a set of roles used by Cloudera Manager to manage and monitor clusters.

The Cloudera Management Service implements various management features as a set of roles:

- **Host Monitor** - collects health and metric information about hosts
- **Service Monitor** - collects health and metric information about services and activity information from the YARN and Impala services
- **Event Server** - aggregates relevant Hadoop events and makes them available for alerting and searching
- **Alert Publisher** - generates and delivers alerts for certain types of events
- **Reports Manager** - generates reports that provide an historical view into disk utilization by user, user group, and directory, processing activities by user and YARN pool, and HBase tables and namespaces. This role is not added in Cloudera Express.




You can view the status of the Cloudera Management Service by doing one of the following:

- Select **Clusters Cloudera Management Service** .
- On the **HomeStatus** tab, in Cloudera Management Service table, click the Cloudera Management Service link.

Health Tests

Cloudera Manager monitors the health of the services, roles, and hosts that are running in your clusters using *health tests*. The Cloudera Management Service also provides health tests for its roles. Role-based health tests are enabled by default. For example, a simple health test is whether there's enough disk space in every NameNode data directory. A more complicated health test may evaluate when the last checkpoint for HDFS was compared to a threshold or whether a DataNode is connected to a NameNode. Some of these health tests also aggregate other health tests: in a distributed system like HDFS, it's normal to have a few DataNodes down (assuming you've got dozens of hosts), so we allow for setting thresholds on what percentage of hosts should color the entire service down.

Health tests can return one of three values: Good, Concerning, and Bad. A test returns Concerning health if the test falls below a warning threshold. A test returns Bad if the test falls below a critical threshold. The overall health of a service or role instance is a roll-up of its health tests. If any health test is Concerning (but none are Bad) the role's or service's health is Concerning; if any health test is Bad, the service's or role's health is Bad.

In the Cloudera Manager Admin Console, health tests results are indicated with colors: Good , Concerning , and Bad .

One common question is whether monitoring can be separated from configuration. One of the goals for monitoring is to enable it without needing to do additional configuration and installing additional tools (for example, Nagios). By having a deep model of the configuration, Cloudera Manager is able to know which directories to monitor, which ports to use, and what credentials to use for those ports. This tight coupling means that, when you install Cloudera Manager all the monitoring is enabled.

Metric Collection and Display

To perform monitoring, the Service Monitor and Host Monitor collects metrics. A *metric* is a numeric value, associated with a name (for example, "CPU seconds"), an entity it applies to ("host17"), and a timestamp. Most metric collection is performed by the Agent. The Agent communicates with a supervised process, requests the metrics, and forwards them to the Service Monitor. In most cases, this is done once per minute.


A few special metrics are collected by the Service Monitor. For example, the Service Monitor hosts an HDFS canary, which tries to write, read, and delete a file from HDFS at regular intervals, and measure whether it succeeded, and how long it took. Once metrics are received, they're aggregated and stored.

Using the Charts page in the Cloudera Manager Admin Console, you can query and explore the metrics being collected. Charts display *time series*, which are streams of metric data points for a specific entity. Each metric data point contains a timestamp and the value of that metric at that timestamp.

Some metrics (for example, `total_cpu_seconds`) are counters, and the appropriate way to query them is to take their rate over time, which is why a lot of metrics queries contain the `dt0` function. For example, `dt0(total_cpu_seconds)`. (The `dt0` syntax is intended to remind you of derivatives. The 0 indicates that the rate of a monotonically increasing counter should never have negative rates.)

Events, Alerts, and Triggers

An *event* is a record that something of interest has occurred – a service's health has changed state, a log message (of the appropriate severity) has been logged, and so on. Many events are enabled and configured by default.

An *alert* is an event that is considered especially noteworthy and is triggered by a selected event. Alerts are shown with an  badge when they appear in a list of events. You can configure the Alert Publisher to send alert notifications by email or by SNMP trap to a trap receiver.

A *trigger* is a statement that specifies an action to be taken when one or more specified conditions are met for a service, role, role configuration group, or host. The conditions are expressed as a *tsquery* statement, and the action to be taken is to change the health for the service, role, role configuration group, or host to either Concerning (yellow) or Bad (red).

Starting the Cloudera Management Service

How to start the Cloudera Management Service.

Before you begin

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

1. Do one of the following:
 - a. Select Clusters Cloudera Management Service .
 - b. Select ActionsStart.
 - On the HomeStatus tab, click the options menu to the right of Cloudera Management Service and select Start.
2. Click Start to confirm. The Command Details window shows the progress of starting the roles.

Results

When Command completed with *n/n* successful subcommands appears, the task is complete. Click Close.

Stopping the Cloudera Management Service

How to stop the Cloudera Management Service.

Before you begin

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

1. Do one of the following:
 - a. Select Clusters Cloudera Management Service .
 - b. Select ActionsStop.
 - On the HomeStatus tab, click the options menu to the right of Cloudera Management Service and select Stop.
2. Click Stop to confirm. The Command Details window shows the progress of stopping the roles.

Results

When Command completed with *n/n* successful subcommands appears, the task is complete. Click Close.

Restarting the Cloudera Management Service

How to restart the Cloudera Management Service.

Before you begin

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

1. Do one of the following:
 - a. Select Clusters Cloudera Management Service .
 - b. Select ActionsRestart.
 - On the HomeStatus tab, click the options menu to the right of Cloudera Management Service and select Restart.
2. Click Restart to confirm. The Command Details window shows the progress of restarting the roles.

Results

When Command completed with *n/n* successful subcommands appears, the task is complete. Click Close.

Starting and Stopping Cloudera Management Service Roles

Before you begin

Minimum Required Role: [Full Administrator](#). This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

1. Do one of the following:
 - Select Clusters Cloudera Management Service .
 - On the HomeStatus tab, in Cloudera Management Service table, click the Cloudera Management Service link.
2. Click the Instances tab.
3. Select a role.
4. Do one of the following:
 - Start: Select Actions for SelectedStart and click Start to confirm
 - Stop: Select Actions for SelectedStop and click Stop to confirm.

Results

When Command completed with *n/n* successful subcommands appears, the task is complete. Click Close.

Configuring Management Service Database Limits

Configuring database service limits lets you control the amount of retained monitoring data.

Before you begin

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

About this task

Each Cloudera Management Service role maintains a database for retaining the data it monitors. These databases (as well as the log files maintained by these services) can grow quite large. Limits on these data sets are configured when you create the management services, but you can modify these parameters through the Configuration settings in the Cloudera Manager Admin Console. For example, the Event Server lets you set a total number of events to store.

There are also settings for the logs that these various services create. You can throttle how big the logs are allowed to get and how many previous logs to retain.

Procedure

1. Do one of the following:
 - Select Clusters Cloudera Management Service .
 - On the HomeStatus tab, in Cloudera Management Service table, click the Cloudera Management Service link.
2. Click the Configuration tab.
3. Select Scope and then one of the following.
 - Host Monitor
 - Service Monitor
4. Select CategoryLog Files to view log file size properties.
5. Edit the appropriate properties.

To apply this configuration property to other role groups as needed, edit the value for the appropriate role group. See .
6. Click Save Changes.

Related Information

[Data Storage for Monitoring Data](#)

Performance Management

This section describes mechanisms and best practices for improving performance.

Optimizing Performance in Cloudera Runtime

This section provides solutions to some performance problems, and describes configuration best practices.



Important: Work with your network administrators and hardware vendors to ensure that you have the proper NIC firmware, drivers, and configurations in place and that your network performs properly. Cloudera recognizes that network setup and upgrade are challenging problems, and will do its best to share useful experiences.

Disable the tuned Service

If your cluster hosts are running RHEL/CentOS 7.x, disable the "tuned" service by running the following commands:

Procedure

1. Ensure that the tuned service is started:

```
systemctl start tuned
```

2. Turn the tuned service off:

```
tuned-adm off
```

3. Ensure that there are no active profiles:

```
tuned-adm list
```

The output should contain the following line:

```
No current active profile
```

4. Shutdown and disable the tuned service:

```
systemctl stop tuned
systemctl disable tuned
```

Disabling Transparent Hugepages (THP)

Most Linux platforms supported by CDH 5 include a feature called *transparent hugepages*, which interacts poorly with Hadoop workloads and can seriously degrade performance.

About this task

Symptom: top and other system monitoring tools show a large percentage of the CPU usage classified as "system CPU". If system CPU usage is 30% or more of the total CPU usage, your system may be experiencing this issue.

To see whether transparent hugepages are enabled, run the following commands and check the output:

```
$ cat defrag_file_pathname
$ cat enabled_file_pathname
```

- [always] never means that transparent hugepages is enabled.
- always [never] means that transparent hugepages is disabled.

To disable Transparent Hugepages, perform the following steps on all cluster hosts:

Procedure

1. (Required for hosts running RHEL/CentOS 7.x.) To disable transparent hugepages on reboot, add the following commands to the `/etc/rc.d/rc.local` file on all cluster hosts:

- RHEL/CentOS 7.x:

```
echo never > /sys/kernel/mm/transparent_hugepage/enabled
echo never > /sys/kernel/mm/transparent_hugepage/defrag
```

- RHEL/CentOS 6.x

```
echo never > /sys/kernel/mm/redhat_transparent_hugepage/defrag
echo never > /sys/kernel/mm/redhat_transparent_hugepage/enabled
```

- Ubuntu/Debian, OL, SLES:

```
echo never > /sys/kernel/mm/transparent_hugepage/defrag
echo never > /sys/kernel/mm/transparent_hugepage/enabled
```

Modify the permissions of the `rc.local` file:

```
chmod +x /etc/rc.d/rc.local
```

2. If your cluster hosts are running RHEL/CentOS 7.x, modify the GRUB configuration to disable THP:
 - a) Add the following line to the `GRUB_CMDLINE_LINUX` options in the `/etc/default/grub` file:

```
transparent_hugepage=never
```

- b) Run the following command:

```
grub2-mkconfig -o /boot/grub2/grub.cfg
```

3. Disable the tuned service, as described above.

You can also disable transparent hugepages interactively (but remember this will not survive a reboot).

To disable transparent hugepages temporarily as root:

```
# echo 'never' > defrag_file_pathname
# echo 'never' > enabled_file_pathname
```

To disable transparent hugepages temporarily using sudo:

```
$ sudo sh -c "echo 'never' > defrag_file_pathname"
$ sudo sh -c "echo 'never' > enabled_file_pathname"
```

Setting the vm.swappiness Linux Kernel Parameter

The Linux kernel parameter, `vm.swappiness`, is a value from 0-100 that controls the swapping of application data (as anonymous pages) from physical memory to virtual memory on disk. You can set the value of the `vm.swappiness` parameter for minimum swapping.

The higher the parameter value, the more aggressively inactive processes are swapped out from physical memory. The lower the value, the less they are swapped, forcing filesystem buffers to be emptied.

On most systems, `vm.swappiness` is set to 60 by default. This is not suitable for Hadoop clusters because processes are sometimes swapped even when enough memory is available. This can cause lengthy garbage collection pauses for important system daemons, affecting stability and performance.

Cloudera recommends that you set `vm.swappiness` to a value between 1 and 10, preferably 1, for minimum swapping on systems where the RHEL kernel is 2.6.32-642.el6 or higher.

To view your current setting for `vm.swappiness`, run:

```
cat /proc/sys/vm/swappiness
```

To set `vm.swappiness` to 1, run:

```
sudo sysctl -w vm.swappiness=1
```

To ensure persistence of the `vm.swappiness` value after reboot:

```
echo 'vm.swappiness=1' > /etc/sysctl.d/90-cloudera-swappiness.conf
```

Improving Performance in Shuffle Handler and IFile Reader

The MapReduce shuffle handler and IFile reader use native Linux calls, (`posix_fadvise(2)` and `sync_data_range`), on Linux systems with Hadoop native libraries installed.

Shuffle Handler

You can improve MapReduce shuffle handler performance by enabling shuffle readahead. This causes the TaskTracker or Node Manager to pre-fetch map output before sending it over the socket to the reducer.

- To enable this feature for YARN, set `mapreduce.shuffle.manage.os.cache`, to true (default). To further tune performance, adjust the value of `mapreduce.shuffle.readahead.bytes`. The default value is 4 MB.
- To enable this feature for MapReduce, set the `mapred.tasktracker.shuffle.fadvise` to true (default). To further tune performance, adjust the value of `mapred.tasktracker.shuffle.readahead.bytes`. The default value is 4 MB.

IFile Reader

Enabling IFile readahead increases the performance of merge operations. To enable this feature for either MRv1 or YARN, set `mapreduce.ifile.readahead` to true (default). To further tune the performance, adjust the value of `mapreduce.ifile.readahead.bytes`. The default value is 4MB.

Tips and Best Practices for Jobs

This section describes changes you can make at the job level.

Use the Distributed Cache to Transfer the Job JAR

Use the distributed cache to transfer the job JAR rather than using the `JobConf(Class)` constructor and the `JobConf.setJar()` and `JobConf.setJarByClass()` methods.

To add JARs to the classpath, use `-libjars jar1.jar2`. This copies the local JAR files to HDFS and uses the distributed cache mechanism to ensure they are available on the task nodes and added to the task classpath.

The advantage of this, over `JobConf.setJar`, is that if the JAR is on a task node, it does not need to be copied again if a second task from the same job runs on that node, though it will still need to be copied from the launch machine to HDFS.



Note: `-libjars` works only if your MapReduce driver uses ToolRunner. If it does not, you would need to use the DistributedCache APIs (Cloudera does not recommend this).

For more information, see item 1 in the blog post *How to Include Third-Party Libraries in Your MapReduce Job*.

Changing the Logging Level on a Job (MRv1)

You can change the logging level for an individual job. You do this by setting the following properties in the job configuration:

- `mapreduce.map.log.level`
- `mapreduce.reduce.log.level`

Valid values are NONE, INFO, WARN, DEBUG, TRACE, and ALL.

Example:

```
JobConf conf = new JobConf();
...

conf.set("mapreduce.map.log.level", "DEBUG");
conf.set("mapreduce.reduce.log.level", "TRACE");
...
```

Decrease Reserve Space

By default, the ext3 and ext4 filesystems reserve 5% space for use by the root user. This reserved space counts as Non DFS Used.

To view the reserved space use the `tune2fs` command:

```
# tune2fs -l /dev/sde1 | egrep "Block size:|Reserved block count"
Reserved block count: 36628312
Block size: 4096
```

The Reserved block count is the number of ext3/ext4 filesystem blocks that are reserved. The block size is the size in bytes. In this example, 150 GB (139.72 Gigabytes) are reserved on this filesystem.

Cloudera recommends reducing the root user block reservation from 5% to 1% for the DataNode volumes. To set reserved space to 1% with the `tune2fs` command:

```
# tune2fs -m 1 /dev/sde1
```

Choosing and Configuring Data Compression

For an overview of compression, see *Data Compression*.

Guidelines for Choosing a Compression Type

- GZIP compression uses more CPU resources than Snappy or LZO, but provides a higher compression ratio. GZip is often a good choice for cold data, which is accessed infrequently. Snappy or LZO are a better choice for hot data, which is accessed frequently.
- BZip2 can also produce more compression than GZip for some types of files, at the cost of some speed when compressing and decompressing. HBase does not support BZip2 compression.
- Snappy often performs better than LZO. It is worth running tests to see if you detect a significant difference.
- For MapReduce, if you need your compressed data to be splittable, BZip2 and LZO formats can be split. Snappy and GZip blocks are not splittable, but files with Snappy blocks inside a container file format such as SequenceFile or Avro can be split. Snappy is intended to be used with a container format, like SequenceFiles or Avro data files, rather than being used directly on plain text, for example, since the latter is not splittable and cannot be processed in parallel using MapReduce. Splittability is not relevant to HBase data.
- For MapReduce, you can compress either the intermediate data, the output, or both. Adjust the parameters you provide for the MapReduce job accordingly. The following examples compress both the intermediate data and the output. MR2 is shown first, followed by MR1.
- MRv2

```
hadoop jar hadoop-examples-.jar sort "-Dmapreduce.compress.map.output=true"
"-Dmapreduce.map.output.compression.codec=org.apache.hadoop.io.compress.GzipCodec"
"-Dmapreduce.output.compress=true"
"-Dmapreduce.output.compression.codec=org.apache.hadoop.io.compress.GzipCodec" -outKey
org.apache.hadoop.io.Text -outValue org.apache.hadoop.io.Text input output
```

- MRv1

```
hadoop jar hadoop-examples-.jar sort "-Dmapred.compress.map.output=true"
"-Dmapred.map.output.compression.codec=org.apache.hadoop.io.compress.GzipCodec"
"-Dmapred.output.compress=true"
"-Dmapred.output.compression.codec=org.apache.hadoop.io.compress.GzipCodec" -outKey
org.apache.hadoop.io.Text -outValue org.apache.hadoop.io.Text input output
```

Related Information

[Hadoop File Formats Support](#)

Resource Management

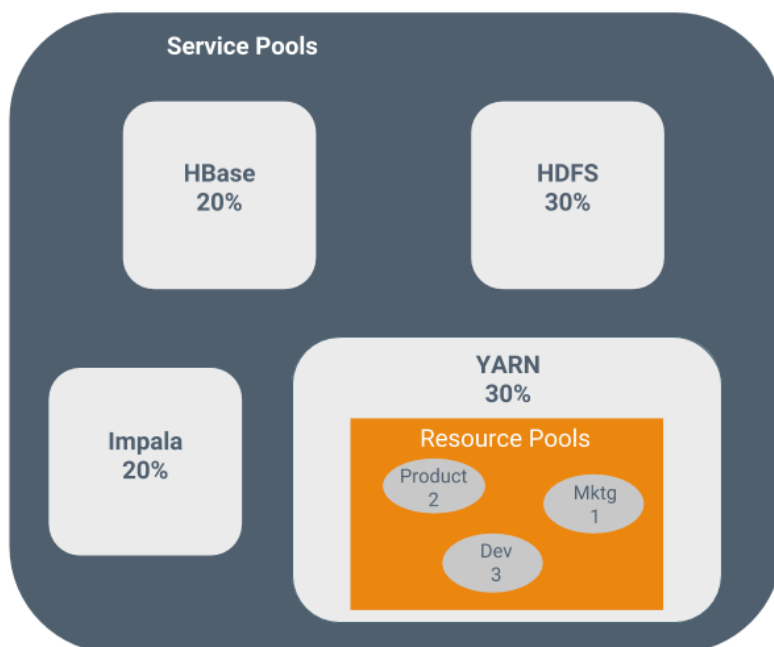
Resource management helps ensure predictable behavior by defining the impact of different services on cluster resources.

Use resource management to:

- Guarantee completion in a reasonable time frame for critical workloads.
- Support reasonable cluster scheduling between groups of users based on fair allocation of resources per group.
- Prevent users from depriving other users access to the cluster.

Statically allocating resources using cgroups is configurable through a single static service pool wizard. You allocate services as a percentage of total resources, and the wizard configures the cgroups.

For example, the following figure illustrates static pools for HBase, HDFS, Impala, and YARN services that are respectively assigned 20%, 30%, 20%, and 30% of cluster resources.



You can dynamically apportion resources that are statically allocated to YARN and Impala by using dynamic resource pools.

Depending on the version of CDH you are using, dynamic resource pools in Cloudera Manager support the following scenarios:

- **YARN** - YARN manages the virtual cores, memory, running applications, maximum resources for undeclared children (for parent pools), and scheduling policy for each pool. In the preceding diagram, three dynamic resource pools—Dev, Product, and Mktg with weights 3, 2, and 1 respectively—are defined for YARN. If an application starts and is assigned to the Product pool, and other applications are using the Dev and Mktg pools, the Product resource pool receives $30\% \times 2/6$ (or 10%) of the total cluster resources. If no applications are using the Dev and Mktg pools, the YARN Product pool is allocated 30% of the cluster resources.
- **Impala** - Impala manages memory for pools running queries and limits the number of running and queued queries in each pool.

Static Service Pools

Static service pools isolate the services in your cluster from one another, so that load on one service has a bounded impact on other services.

Services are allocated a static percentage of total resources—CPU, memory, and I/O weight—which are not shared with other services. When you configure static service pools, Cloudera Manager computes recommended memory, CPU, and I/O configurations for the worker roles of the services that correspond to the percentage assigned to each service. Static service pools are implemented per role group within a cluster, using Linux control groups (cgroups) and cooperative memory limits (for example, Java maximum heap sizes). Static service pools can be used to control access to resources by HBase, HDFS, Impala, MapReduce, Solr, Spark, YARN, and add-on services. Static service pools are not enabled by default.

**Note:**

- I/O allocation only works when short-circuit reads are enabled.
- I/O allocation does not handle write side I/O because cgroups in the Linux kernel do not currently support buffered writes.

Viewing Static Service Pool Status

Select Clusters *Cluster name* Static Service Pools. If the cluster has a YARN service, the Static Service Pools Status tab displays and shows whether resource management is enabled for the cluster, and the currently configured service pools.

Enabling and Configuring Static Service Pools

To enable and configure static service pools, you enter the percentage of resources to allocate to each service and then restart the cluster.

Before you begin

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

1. Select Clusters *Cluster name* Static Service Pools.
2. Click the Configuration tab.
The Step 1 of 4: Basic Allocation Setup page displays. In each field in the basic allocation table, enter the percentage of resources to give to each service. The total must add up to 100%.
3. Click Continue to proceed.
Step 2: Review Changes - The allocation of resources for each resource type and role displays with the new values as well as the values previously in effect. The values for each role are set by role group; if there is more than one role group for a given role type (for example, for RegionServers or DataNodes) then resources will be allocated separately for the hosts in each role group.
4. Take note of changed settings. If you have previously customized these settings, check these over carefully:
 - Click the **>** to the right of each percentage to display the allocations for a single service. Click **>** to the right of the Total (100%) to view all the allocations in a single page.
 - Click the Back button to go to the previous page and change your allocations.
5. When you are satisfied with the allocations, click Continue.
The Step 3 of 4: Restart Services page displays.
6. To apply the new allocation percentages, click Restart Now to restart the cluster. To skip this step, click Restart Later. If HDFS High Availability is enabled, you will have the option to choose a rolling restart.
7. Step 4 of 4: Progress displays the status of the restart commands. Click Finished after the restart commands complete.

After you enable static service pools, there are three additional tasks:

8. Delete everything under the local directory path on NodeManager hosts. The local directory path is configurable, and can be verified in Cloudera Manager with [YARN Configuration NodeManager Local Directories](#) .
9. Enable cgroups for resource management. You can enable cgroups in Cloudera Manager with [Yarn Configuration Use CGroups for Resource Management](#) .
10. If you are using the optional Impala scratch directory, delete all files in the Impala scratch directory. The directory path is configurable, and can be verified in Cloudera Manager with [Impala Configuration Impala Daemon Scratch Directories](#) .

Disabling Static Service Pools

To disable static service pools, disable cgroup-based resource management for all hosts in all clusters.

Before you begin

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

1. In the main navigation bar, click Hosts.
2. Click the Configuration tab.
3. Select ScopeResource Management.
4. Clear the Enable Cgroup-based Resource Management property.
5. Click Save Changes.
6. Restart all services.

Results

Static resource management is disabled, but the percentages you set when you configured the pools, and all the changed settings (for example, heap sizes), are retained by the services. The percentages and settings will also be used when you re-enable static service pools. If you want to revert to the settings you had before static service pools were enabled, follow the procedures in *Viewing and Reverting Configuration Changes*.

Linux Control Groups (cgroups)

Minimum Required Role: [Full Administrator](#). This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Cloudera Manager supports the Linux control groups (cgroups) kernel feature. With cgroups, administrators can impose per-resource restrictions and limits on services and roles. This provides the ability to allocate resources using cgroups to enable isolation of compute frameworks from one another. Resource allocation is implemented by setting properties for the services and roles.

Linux Distribution Support

Cgroups are a feature of the Linux kernel, and as such, support depends on the host's Linux distribution and version as shown in the following tables. If a distribution lacks support for a given parameter, changes to the parameter have no effect.

Table 3: RHEL-compatible

Distribution	CPU Shares	I/O Weight	Memory Soft Limit	Memory Hard Limit
Red Hat Enterprise Linux, CentOS, and Oracle Enterprise Linux 7	■	■	■	■
Red Hat Enterprise Linux, CentOS, and Oracle Enterprise Linux 6	■	■	■	■

Table 4: SLES

Distribution	CPU Shares	I/O Weight	Memory Soft Limit	Memory Hard Limit
SUSE Linux Enterprise Server 12	■	■	■	■
SUSE Linux Enterprise Server 11	■	■	■	■

Table 5: Ubuntu

Distribution	CPU Shares	I/O Weight	Memory Soft Limit	Memory Hard Limit
Ubuntu 16.04 LTS	■	■	■	■

Distribution	CPU Shares	I/O Weight	Memory Soft Limit	Memory Hard Limit
Ubuntu 16.04 LTS	■	■	■	■
Ubuntu 14.04 LTS	■	■	■	■
Ubuntu 12.04 LTS	■	■	■	■

Table 6: Debian

Distribution	CPU Shares	I/O Weight	Memory Soft Limit	Memory Hard Limit
Debian 7.1	■	■	■	■
Debian 7.0	■	■	■	■
Debian 6.0	■	■	■	■

Table 7: Oracle linux (OL)

Distribution	CPU Shares	I/O Weight	Memory Soft Limit	Memory Hard Limit
Oracle linux 7	■	■	■	■
Oracle linux 6	■	■	■	■

The exact level of support can be found in the Cloudera Manager Agent log file, shortly after the Agent has started. In the log file, look for an entry like this:

```
Found cgroups capabilities: {
  'has_memory': True,
  'default_memory_limit_in_bytes': 9223372036854775807,
  'writable_cgroup_dot_procs': True,
  'has_cpu': True,
  'default_blkio_weight': 1000,
  'default_cpu_shares': 1024,
  'has_blkio': True}
```

The `has_cpu` and similar entries correspond directly to support for the CPU, I/O, and memory parameters.

Resource Management with Control Groups

To use cgroups, you must enable cgroup-based resource management under the host resource management configuration properties. However, if you configure static service pools, this property is set as part of that process.

Before you begin

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

About this task

Cgroups-based resource management can be enabled for all hosts, or on a per-host basis.

Procedure

1. If you have upgraded from a version of Cloudera Manager older than Cloudera Manager 4.5, restart every Cloudera Manager Agent before using cgroups-based resource management:
 - a) Stop all services, including the Cloudera Management Service.
 - b) On each cluster host, run as root:
 - RHEL-compatible 7 and higher:

```
sudo service cloudera-scm-agent next_stop_hard
```

```
sudo service cloudera-scm-agent restart
```

- All other Linux distributions:

```
$ sudo service cloudera-scm-agent hard_restart
```

c) Start all services.

2. Click the Hosts tab.
3. Optionally click the link of the host where you want to enable cgroups.
- 4.
5. Select CategoryResource Management.
6. Select Enable Cgroup-based Resource Management.
7. Restart all roles on the host or hosts.

Limitations

- Role group and role instance override cgroup-based resource management parameters must be saved one at a time. Otherwise some of the changes that should be reflected dynamically will be ignored.
- The role group abstraction is an imperfect fit for resource management parameters, where the goal is often to take a numeric value for a host resource and distribute it amongst running roles. The role group represents a "horizontal" slice: the same role across a set of hosts. However, the cluster is often viewed in terms of "vertical" slices, each being a combination of worker roles (such as TaskTracker, DataNode, RegionServer, Impala Daemon, and so on). Nothing in Cloudera Manager guarantees that these disparate horizontal slices are "aligned" (meaning, that the role assignment is identical across hosts). If they are unaligned, some of the role group values will be incorrect on unaligned hosts. For example a host whose role groups have been configured with memory limits but that's missing a role will probably have unassigned memory.

Configuring Resource Parameters

After enabling cgroups, you can restrict and limit the resource consumption of roles (or role groups) on a per-resource basis.

All of these parameters can be found in the Cloudera Manager Admin Console, under the Resource Management category:

- CPU Shares - The more CPU shares given to a role, the larger its share of the CPU when under contention. Until processes on the host (including both roles managed by Cloudera Manager and other system processes) are contending for all of the CPUs, this will have no effect. When there is contention, those processes with higher CPU shares will be given more CPU time. The effect is linear: a process with 4 CPU shares will be given roughly twice as much CPU time as a process with 2 CPU shares.

Updates to this parameter are dynamically reflected in the running role.

- I/O Weight - The greater the I/O weight, the higher priority will be given to I/O requests made by the role when I/O is under contention (either by roles managed by Cloudera Manager or by other system processes).

This only affects read requests; write requests remain unprioritized. The Linux I/O scheduler controls when buffered writes are flushed to disk, based on time and quantity thresholds. It continually flushes buffered writes from multiple sources, not certain prioritized processes.

Updates to this parameter are dynamically reflected in the running role.

- Memory Soft Limit - When the limit is reached, the kernel will reclaim pages charged to the process if and only if the host is facing memory pressure. If reclaiming fails, the kernel may kill the process. Both anonymous as well as page cache pages contribute to the limit.

After updating this parameter, you must restart the role for changes to take effect.

- Memory Hard Limit - When a role's resident set size (RSS) exceeds the value of this parameter, the kernel will swap out some of the role's memory. If it is unable to do so, it will kill the process. The kernel measures memory consumption in a manner that does not necessarily match what the top or ps report for RSS, so expect that this limit is a rough approximation.

After updating this parameter, you must restart the role for changes to take effect.

Example: Protecting Production MapReduce Jobs from Impala Queries

Suppose you have MapReduce deployed in production and want to roll out Impala without affecting production MapReduce jobs. For simplicity, we will make the following assumptions:

- The cluster is using homogenous hardware
- Each worker host has two cores
- Each worker host has 8 GB of RAM
- Each worker host is running a DataNode, TaskTracker, and an Impala Daemon
- Each role type is in a single role group
- Cgroups-based resource management has been enabled on all hosts

Action	Procedure
CPU	<ol style="list-style-type: none"> 1. Leave DataNode and TaskTracker role group CPU shares at 1024. 2. Set Impala Daemon role group's CPU shares to 256. 3. The TaskTracker role group should be configured with a Maximum Number of Simultaneous Map Tasks of 2 and a Maximum Number of Simultaneous Reduce Tasks of 1. This yields an upper bound of three MapReduce tasks at any given time; this is an important detail for memory sizing.
Memory	<ol style="list-style-type: none"> 1. Set Impala Daemon role group memory limit to 1024 MB. 2. Leave DataNode maximum Java heap size at 1 GB. 3. Leave TaskTracker maximum Java heap size at 1 GB. 4. Leave MapReduce Child Java Maximum Heap Size for Gateway at 1 GB. 5. Leave cgroups hard memory limits alone. We'll rely on "cooperative" memory limits exclusively, as they yield a nicer user experience than the cgroups-based hard memory limits.
I/O	<ol style="list-style-type: none"> 1. Leave DataNode and TaskTracker role group I/O weight at 500. 2. Impala Daemon role group I/O weight is set to 125.

When you're done with configuration, restart all services for these changes to take effect. The results are:

1. When MapReduce jobs are running, all Impala queries together will consume up to a fifth of the cluster's CPU resources.
2. Individual Impala Daemons will not consume more than 1 GB of RAM. If this figure is exceeded, new queries will be cancelled.
3. DataNodes and TaskTrackers can consume up to 1 GB of RAM each.
4. We expect up to 3 MapReduce tasks at a given time, each with a maximum heap size of 1 GB of RAM. That's up to 3 GB for MapReduce tasks.
5. The remainder of each host's available RAM (6 GB) is reserved for other host processes.
6. When MapReduce jobs are running, read requests issued by Impala queries will receive a fifth of the priority of either HDFS read requests or MapReduce read requests.

Data Storage for Monitoring Data

The Service Monitor and Host Monitor roles in the Cloudera Management Service store time series data, health data, and Impala query and YARN application metadata.

Configuring Service Monitor Data Storage

The Service Monitor stores time series data and health data, Impala query metadata, and YARN application metadata.

By default, the data is stored in `/var/lib/cloudera-service-monitor/` on the Service Monitor host. You can change this by modifying the Service Monitor Storage Directory configuration (`firehose.storage.base.directory`). To change this configuration on an active system, see *Moving Monitoring Data on an Active Cluster*.

You can control how much disk space to reserve for the different classes of data the Service Monitor stores by changing the following configuration options:

- Time-series metrics and health data - Time-Series Storage (firehose_time_series_storage_bytes - 10 GB default, 10 GB minimum)
- Impala query metadata - Impala Storage (firehose_impala_storage_bytes - 1 GB default)
- YARN application metadata - YARN Storage (firehose_yarn_storage_bytes - 1 GB default)

For information about how metric data is stored in Cloudera Manager and how storage limits impact data retention, see *Data Granularity and Time-Series Metric Data*.

The default values are small, so you should examine disk usage after several days of activity to determine how much space is needed.

Configuring Host Monitor Data Storage

The Host Monitor stores time series data and health data.

By default, the data is stored in `/var/lib/cloudera-host-monitor/` on the Host Monitor host. You can change this by modifying the Host Monitor Storage Directory configuration. To change this configuration on an active system, follow the procedure in *Moving Monitoring Data on an Active Cluster*.

You can control how much disk space to reserve for Host Monitor data by changing the following configuration option:

- Time-series metrics and health data: Time Series Storage (firehose_time_series_storage_bytes - 10 GB default, 10 GB minimum)

For information about how metric data is stored in Cloudera Manager and how storage limits impact data retention, see *Data Granularity and Time-Series Metric Data*.

The default value is small, so you should examine disk usage after several days of activity to determine how much space they need. The Charts Library tab on the Cloudera Management Service page shows the current disk space consumed and its rate of growth, categorized by the type of data stored. For example, you can compare the space consumed by raw metric data to daily summaries of that data.

Viewing Host and Service Monitor Data Storage

The Cloudera Management Service page shows the current disk space consumed and its rate of growth, categorized by the type of data stored. For example, you can compare the space consumed by raw metric data to daily summaries of that data.

Procedure

1. Select ClustersCloudera Management Service.
2. Click the Charts Library tab.

Data Granularity and Time-Series Metric Data

The Service Monitor and Host Monitor store time-series metric data in a variety of ways.

When the data is received, it is written as-is to the metric store. Over time, the raw data is summarized to and stored at various data granularities. For example, after ten minutes, a summary point is written containing the average of the metric over the period as well as the minimum, the maximum, the standard deviation, and a variety of other statistics. This process is summarized to produce hourly, six-hourly, daily, and weekly summaries. This data summarization procedure applies only to metric data. When the Impala query and YARN application monitoring storage limit is reached, the oldest stored records are deleted.

The Service Monitor and Host Monitor internally manage the amount of overall storage space dedicated to each data granularity level. When the limit for a level is reached, the oldest data points at that level are deleted. Metric data for that time period remains available at the lower granularity levels. For example, when an hourly point for a particular time is deleted to free up space, a daily point still exists covering that hour. Because each of these data granularities consumes significantly less storage than the previous summary level, lower granularity levels can be retained for longer periods of time. With the recommended amount of storage, weekly points can often be retained indefinitely.

Some features, such as detailed display of health results, depend on the presence of raw data. Health history is maintained by the event store dictated by its retention policies.

Moving Monitoring Data on an Active Cluster

You can change where monitoring data is stored on a cluster.

Basic: Changing the Configured Directory

1. Stop the Service Monitor or Host Monitor.
2. Save your old monitoring data and then copy the current directory to the new directory (optional).
3. Update the Storage Directory configuration option (`firehose.storage.base.directory`) on the corresponding role configuration page.
4. Start the Service Monitor or Host Monitor.

Advanced: High Performance

For the best performance, and especially for a large cluster, Host Monitor and Service Monitor storage directories should have their own dedicated spindles. In most cases, that provides sufficient performance, but you can divide your data further if needed. You cannot configure this directly with Cloudera Manager; instead, you must use symbolic links.

For example, if all your Service Monitor data is located in `/data/1/service_monitor`, and you want to separate your Impala data from your time series data, you could do the following:

1. Stop the Service Monitor.
2. Move the original Impala data in `/data/1/service_monitor/impala` to the new directory, for example `/data/2/impala_data`.
3. Create a symbolic link from `/data/1/service_monitor/impala` to `/data/2/impala_data` with the following command:

```
ln -s /data/2/impala_data /data/1/service_monitor/impala
```

4. Start the Service Monitor.

Host Monitor and Service Monitor Memory Configuration

You can configure Java heap size and non-Java memory size. The memory recommended for these configuration options depends on the number of hosts in the cluster, the services running on the cluster, and the number of monitored entities.

Monitored entities are the objects monitored by the Service Monitor or Host Monitor. As the number of hosts and services increases, the number of monitored entities also increases.

In addition to the memory configured, the Host Monitor and Service Monitor use the Linux page cache. Memory available for page caching on the Host Monitor and Service Monitor hosts improves performance.

Configuring Memory Allocations

To configure memory allocations, determine how many entities are being monitored and then consult the tables below for required and recommended memory configurations.

About this task

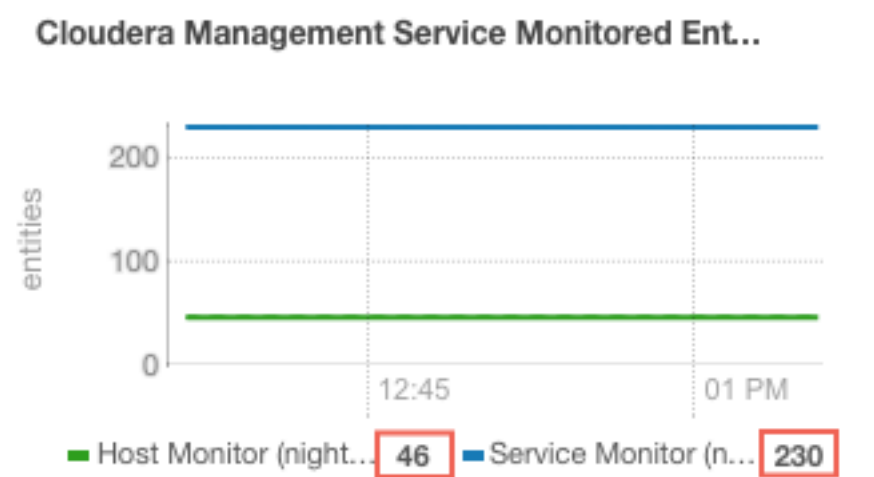
To determine the number of entities being monitored:

Procedure

1. Go to ClustersCloudera Management Service.

2. Locate the chart with the title Cloudera Management Service Monitored Entities.

The number of monitored entities for the Host Monitor and Service Monitor displays at the bottom of the chart. In the following example, the Host Monitor has 46 monitored entities and the Service Monitor has 230 monitored entities.



3. Use the number of monitored entities for the Host Monitor to determine its memory requirements and recommendations in the tables below.
4. Use the number of monitored entities for the Service Monitor to determine its memory requirements and recommendations in the tables below.

Clusters with HDFS, YARN, or Impala

Use the recommendations in this table for clusters where the only services having worker roles are HDFS, YARN, or Impala.

Number of Monitored Entities	Number of Hosts	Required Java Heap Size	Recommended Non-Java Heap Size
0-2,000	0-100	1 GB	6 GB
2,000-4,000	100-200	1.5 GB	6 GB
4,000-8,000	200-400	1.5 GB	12 GB
8,000-16,000	400-800	2.5 GB	12 GB
16,000-20,000	800-1,000	3.5 GB	12 GB

Clusters with HBase, Solr, Kafka, or Kudu

Use the recommendations when services such as HBase, Solr, Kafka, or Kudu are deployed in the cluster. These services typically have larger quantities of monitored entities.

Number of Monitored Entities	Number of Hosts	Required Java Heap Size	Recommended Non-Java Heap Size
0-30,000	0-100	2 GB	12 GB
30,000-60,000	100-200	3 GB	12 GB
60,000-120,000	200-400	3.5 GB	12 GB
120,000-240,000	400-800	8 GB	20 GB

Accessing Storage Using Amazon S3

Referencing S3 Credentials for YARN, MapReduce, or Spark Clients

If you have selected IAM authentication, no additional steps are needed. If you are not using IAM authentication, use one of the following three options to provide Amazon S3 credentials to clients.



Note: This method of specifying AWS credentials to clients does not completely distribute secrets securely because the credentials are not encrypted. Use caution when operating in a multi-tenant environment.

Programmatic

Specify the credentials in the configuration for the job. This option is most useful for Spark jobs.

Make a modified copy of the configuration files

Make a copy of the configuration files and add the S3 credentials:

1. For YARN and MapReduce jobs, copy the contents of the `/etc/hadoop/conf` directory to a local working directory under the home directory of the host where you will submit the job. For Spark jobs, copy `/etc/spark/conf` to a local directory under the home directory of the host where you will submit the job.
2. Set the permissions for the configuration files appropriately for your environment and ensure that unauthorized users cannot access sensitive configurations in these files.
3. Add the following to the `core-site.xml` file within the `<configuration>` element:

```
<property>
  <name>fs.s3a.access.key</name>
  <value>Amazon S3 Access Key</value>
</property>

<property>
  <name>fs.s3a.secret.key</name>
  <value>Amazon S3 Secret Key</value>
</property>
```

4. Reference these versions of the configuration files when submitting jobs by running the following command:

- YARN or MapReduce:

```
export HADOOP_CONF_DIR=path to local configuration directory
```

- Spark:

```
export SPARK_CONF_DIR=path to local configuration directory
```



Note: If you update the client configuration files from Cloudera Manager, you must repeat these steps to use the new configurations.

Reference the managed configuration files and add AWS credentials

This option allows you to continue to use the configuration files managed by Cloudera Manager. If you deploy new configuration files, the new values are included by reference in your copy of the configuration files while also maintaining a version of the configuration that contains the Amazon S3 credentials:

1. Create a local directory under your home directory.
2. Copy the configuration files from `/etc/hadoop/conf` to the new directory.
3. Set the permissions for the configuration files appropriately for your environment.

4. Edit each configuration file:
 - a. Remove all elements within the <configuration> element.
 - b. Add an XML <include> element within the <configuration> element to reference the configuration files managed by Cloudera Manager. For example:

```
<include xmlns="http://www.w3.org/2001/XInclude"
  href="/etc/hadoop/conf/hdfs-site.xml">
  <fallback />
</include>
```

5. Add the following to the core-site.xml file within the <configuration> element:

```
<property>
  <name>fs.s3a.access.key</name>
  <value>Amazon S3 Access Key</value>
</property>

<property>
  <name>fs.s3a.secret.key</name>
  <value>Amazon S3 Secret Key</value>
</property>
```

6. Reference these versions of the configuration files when submitting jobs by running the following command:
 - YARN or MapReduce:

```
export HADOOP_CONF_DIR=path to local configuration directory
```

- Spark:

```
export SPARK_CONF_DIR=path to local configuration directory
```

Example core-site.xml file:

```
<?xml version="1.0"?>
<?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
<configuration>
  <include xmlns="http://www.w3.org/2001/XInclude"
    href="/etc/hadoop/conf/core-site.xml">
    <fallback />
  </include>

  <property>
    <name>fs.s3a.access.key</name>
    <value>Amazon S3 Access Key</value>
  </property>

  <property>
    <name>fs.s3a.secret.key</name>
    <value>Amazon S3 Secret Key</value>
  </property>
</configuration>
```

Referencing Amazon S3 in URIs

By default, files are still placed on the local HDFS and not on S3 if the protocol is not specified in the URI. When you have added the Amazon S3 service, use one of the following options to construct the URIs to reference when submitting jobs:

- Amazon S3:

```
s3a://bucket_name/path
```

- HDFS:

```
hdfs://path
```

or

```
/path
```

Related Information

[Accessing Data Stored in Amazon S3 through Spark](#)

[Impala with Amazon S3](#)

Using Fast Upload with Amazon S3

Writing data to Amazon S3 is subject to limitations of the s3a OutputStream implementation, which buffers the entire file to disk before uploading it to S3. This can cause the upload to proceed very slowly and can require a large amount of temporary disk space on local disks.

You can configure a cluster to use the Fast Upload feature. This feature implements several performance improvements and has tunable parameters for buffering to disk (the default) or to memory, tuning the number of threads, and for specifying the disk directories used for buffering.

Related Information

[Hadoop-AWS module: Integration with Amazon Web Services](#)

Enabling Fast Upload using Cloudera Manager

Procedure

To enable Fast Upload for clusters managed by Cloudera Manager:

1. Go to the HDFS service.
2. Click the Configuration tab.
3. Search for "core-site.xml" and locate the Cluster-wide Advanced Configuration Snippet (Safety Valve) for core-site.xml property.
4. Add the fs.s3a.fast.upload property and set it to true.
5. Set any additional tuning properties in the Cluster-wide Advanced Configuration Snippet (Safety Valve) for core-site.xml configuration properties.
6. Click Save Changes.

Results

Cloudera Manager will indicate that there are stale services and which services need to be restarted.

Related Information

[Setting an Advanced Configuration Snippet for a Cluster](#)

[Restarting a Cloudera Runtime Service](#)

Configuring and Managing S3Guard

Minimum Required Role: [User Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Data written to Amazon S3 buckets is subject to the "eventual consistency" guarantee provided by Amazon Web Services (AWS), which means that data written to S3 may not be immediately available for queries and listing operations. This can cause failures in multi-step ETL workflows, where data from a previous step is not available to the next step. The S3Guard feature guarantees a consistent view of data stored in Amazon S3 by storing additional metadata in a table residing in an Amazon DynamoDB instance. Depending on the workload, this additional metadata store may also improve performance for Hive, Spark, and Impala jobs.

All processes that modify the S3 bucket that S3Guard is enabled for must use S3Guard. Since S3Guard works by logging metadata changes to an external database, modifying the bucket outside of S3Guard will cause the S3 data and the S3Guard database to go out of sync. This can cause issues such as S3A/S3Guard thinking that files are or are not present despite the bucket having different data.

To enable S3Guard, you set up an Amazon DynamoDB database from Amazon Web Services. Amazon charges an hourly rate for this service based on the capacity you provision.

When the data stored in S3 eventually becomes consistent (usually within 24 hours or less), the S3Guard metadata is no longer required and you can periodically prune the S3Guard metadata stored in the DynamoDB to clear older entries. Pruning can also reduce costs associated with the DynamoDB.

To configure S3Guard in your cluster, you must provide the following:

- Credentials for the Amazon S3 bucket.
- An instance of Amazon DynamoDB database provisioned from Amazon Web Services.
- The configured region for the DynamoDB database.
- A cluster managed by Cloudera Manager.

Pruning the S3Guard Metadata

Amazon charges for the amount of data stored in the DynamoDB and the bandwidth used for reads and writes to the database. To optimize costs and improve performance, you can remove stale metadata from the DynamoDB table by running the Prune command. Generally, data written to S3 becomes consistent after 24 hours or less, meaning that you only need to maintain metadata in DynamoDB for about one day. You can monitor the usage of DynamoDB using AWS tools to determine how often and when to prune the table.

Running the Prune command removes all metadata that is older than the age you specify with the S3Guard Metadata Pruning Age property in the S3Guard configuration. You can run this command from the Cloudera Manager Admin Console, or you can create a script to run the Prune command automatically using the Cloudera Manager API. Cloudera recommends that you run that script using a Linux cron job or other scheduling mechanism to regularly prune the metadata.

Configuring S3Guard for Cluster Access to S3

Procedure

1. Specify the AWS credentials for the Amazon S3 instance where you want to enable S3Guard.
 - Add a new AWS credential.
After adding the credential, the Edit S3Guard dialog box displays.
 - Use an existing AWS credential:
 - **a.** Go to Administration AWS Credentials .

- b. Locate the credential you want to use and click **Actions Edit S3Guard**.

The Edit S3Guard dialog box displays.

2. Select **Enable S3Guard**.
3. Edit the following S3Guard configuration properties:

Table 8: S3Guard Configuration Properties

Property	Description
Automatically Create S3Guard Metadata Table (fs.s3a.s3guard.ddb.table.create) API Name: s3guard_table_auto_create	When Yes is selected, the DynamoDB table that stores the S3Guard metadata is automatically created if it does not exist. When No is selected and the table does not exist, running the Prune command, queries, or other jobs on S3 will fail.
S3Guard Metadata Table Name (fs.s3a.s3guard.ddb.table) API Name: s3guard_table_name	The name of the DynamoDB table that stores the S3Guard metadata. By default, the table is named s3guard-metadata.
S3Guard Metadata Region Name (fs.s3a.s3guard.ddb.region) API Name: s3guard_region	The DynamoDB region to connect to for access to the S3Guard metadata. Set this property to a valid region.
Expand the Advanced section to configure the following properties:	
S3Guard Metadata Pruning Age (fs.s3a.s3guard.cli.prune.age) API Name: s3guard_cache_prune_age_ms	Maximum age for S3Guard metadata. Whenever the Prune command runs, entries in the S3Guard metadata cache older than this age will be deleted. You can enter this value in milliseconds, seconds, minutes, hours, or days.
S3Guard Metadata Table Read Capacity (fs.s3a.s3guard.ddb.table.capacity.read) API Name: s3guard_table_capacity_read	Provisioned throughput requirements, in capacity units, for read operations from the DynamoDB table used for the S3Guard metadata. This value is only used when creating a new DynamoDB table. After the table is created, you can monitor the throughput and adjust the read capacity using the DynamoDB AWS Management Console.
S3Guard Metadata Table Write Capacity (fs.s3a.s3guard.ddb.table.capacity.write) API Name: s3guard_table_capacity_write	Provisioned throughput requirements, in capacity units, for write operations to the DynamoDB table used for the S3Guard metadata. This value is only used when creating a new DynamoDB table. After the table is created, you can monitor the throughput and adjust the write capacity as needed using the DynamoDB AWS Management Console.

4. Click **Save**.

The Connect to Amazon Web Services dialog box displays.

5. To enable cluster access to S3 using the S3 Connector Service, click the **Enable for *Cluster Name*** link in the Cluster Access to S3 section.

Follow the prompts to add the S3 Connector Service.



Note: S3Guard is not supported for Cloud Backup and Restore and Cloudera Navigator Access to S3.

Editing the S3Guard Configuration

Procedure

To edit or disable the S3Guard configuration:

1. Click Administration AWS Credentials.
2. Locate the credential associated with the S3Guard configuration and click Actions Edit S3Guard . The Edit S3Guard dialog box displays.
3. Edit the S3Guard configuration. (To disable S3Guard for this credential, uncheck Enable S3Guard.)
4. Click Save.

Running the Prune Command Using Cloudera Manager Admin Console

Before you begin

Minimum Required Role: [Cluster Administrator](#) (also provided by Full Administrator) This feature is not available when using Cloudera Manager to manage Data Hub clusters.

Procedure

To prune the S3Guard metadata in the DynamoDB table using the Cloudera Manager Admin Console:

1. Go to Administration AWS Credentials .
2. Locate the credential associated with the S3 data and click Actions Run S3 Guard Prune Command .

Running the Prune Command Using the Cloudera Manager API

Cloudera recommends that you automate running the Prune command by creating a script that uses the Cloudera Manager API to run the command. You can run the command using a REST command, a Python script, or Java class. Configure the script using the Linux cron command or another scheduling mechanism to run on a regular schedule.

REST

See the Rest API documentation.

You can run the Prune command by issuing the following REST request:

```
curl -X POST -u username:password
  'Cloudera_Manager_server_URL:port_number/api/vAPI_version_number/externalAccounts/account/Credential_Name/commands/S3GuardPrune'
```

For example, the following request runs the S3Guard prune command on the data associated with the johnsmith credential. The response from Cloudera Manager is also displayed (within the curly brackets):

```
curl -X POST -u admin:admin 'http://clusterhost-1.gce.mycompany.com:7180/api/v16/externalAccounts/account/johnsmith/commands/S3GuardPrune'
{
  "id" : 322,
  "name" : "S3GuardPrune",
  "startTime" : "2017-03-20T23:35:55.453Z",
  "active" : true,
  "children" : {
    "items" : [ {
      "id" : 323,
      "name" : "HostS3GuardPrune",
      "startTime" : "2017-03-20T23:35:55.777Z",
      "active" : true,
      "hostRef" : {
        "hostId" : "ff988a15-3749-4178-b167-a60b15f91653"
      }
    }
  ]
}
```

Python

You can also use a Python script to run the Prune command. See the *aws.py* link under Related Information for the code and usage instructions.

Java

See the Javadoc.

How to Configure a MapReduce Job to Access S3 with an HDFS Credstore

Configure your MapReduce jobs to read and write to Amazon S3 using a custom password for an HDFS Credstore.

Procedure

1. Copy the contents of the `/etc/hadoop/conf` directory to a local working directory on the host where you will submit the MapReduce job. Use the `--dereference` option when copying the file so that symlinks are correctly resolved. For example:

```
cp -r --dereference /etc/hadoop/conf ~/my_custom_config_directory
```

2. Change the permissions of the directory so that only you have access:

```
chmod go-wrx -R my_custom_config_directory/
```

If you see the following message, you can ignore it:

```
cp: cannot open '/etc/hadoop/conf/container-executor.cfg' for reading: Permission denied
```

3. Add the following to the copy of the `core-site.xml` file in the working directory:

```
<property>
  <name>hadoop.security.credential.provider.path</name>
  <value>jceks://hdfs/user/username/awscreds.jceks</value>
</property>
```

4. Specify a custom Credstore by running the following command on the client host:

```
export HADOOP_CREDSTORE_PASSWORD=your_custom_keystore_password
```

5. In the working directory, edit the `mapred-site.xml` file:

- a) Add the following properties:

```
<property>
  <name>yarn.app.mapreduce.am.env</name>
  <value>HADOOP_CREDSTORE_PASSWORD=your_custom_keystore_password</value>
</property>

<property>
  <name>mapred.child.env</name>
  <value>HADOOP_CREDSTORE_PASSWORD=your_custom_keystore_password</value>
</property>
```

- b) Add `yarn.app.mapreduce.am.env` and `mapred.child.env` to the comma-separated list of values of the `mapreduce.job.redacted-properties` property. For example (new values shown bold):

```
<property>
  <name>mapreduce.job.redacted-properties</name>
  <value>fs.s3a.access.key,fs.s3a.secret
  .key,yarn.app.mapreduce.am.env,mapred.child.env</value>
</property>
```

6. Set the environment variable to point to your working directory:

```
export HADOOP_CONF_DIR=~/path_to_working_directory
```

7. Create the Credstore by running the following commands:

```
hadoop credential create fs.s3a.access.key
hadoop credential create fs.s3a.secret.key
```

You will be prompted to enter the access key and secret key.

8. List the credentials to make sure they were created correctly by running the following command:

```
hadoop credential list
```

9. Submit your job. For example:

- ls

```
hdfs dfs -ls s3a://S3_Bucket/
```

- distcp

```
hadoop distcp hdfs_path s3a://S3_Bucket/S3_path
```

- teragen (package-based installations)

```
hadoop jar /usr/lib/hadoop-mapreduce/hadoop-mapreduce-examples.jar teragen 100 s3a://S3_Bucket/teragen_test
```

- teragen (parcel-based installations)

```
hadoop jar /opt/cloudera/parcels/CDH/lib/hadoop-mapreduce/hadoop-mapreduce-examples.jar teragen 100 s3a://S3_Bucket/teragen_test
```

Importing Data into Amazon S3 Using Sqoop

Sqoop supports data import from RDBMS into Amazon S3.



Note: Sqoop import is supported only into the S3A (s3a:// protocol) filesystem.

Related Information

[Hadoop-AWS module: Integration with Amazon Web Services](#)

Authentication

You must authenticate to an S3 bucket using Amazon Web Service credentials. There are three ways to pass these credentials:

- Provide them in the configuration file or files manually.
- Provide them on the sqoop command line.
- Reference a credential store to "hide" sensitive data, so that they do not appear in the console output, configuration file, or log files.

Amazon S3 Block Filesystem URI example:

```
s3a://bucket_name/path/to/file
```

S3 credentials can be provided in a configuration file (for example, `core-site.xml`):

```
<property>
  <name>fs.s3a.access.key</name>
  <value>...</value>
</property>
<property>
  <name>fs.s3a.secret.key</name>
  <value>...</value>
</property>
```

You can also set up the configurations through Cloudera Manager by adding the configurations to the appropriate Advanced Configuration Snippet property.

Credentials can be provided through the command line:

```
sqoop import -Dfs.s3a.access.key=... -Dfs.s3a.secret.key=... --target-dir s3a://
```

For example:

```
sqoop import -Dfs.s3a.access.key=$ACCES_KEY -Dfs.s3a.secret.key=$SECRET_KEY
--connect $CONN --username $USER --password $PWD --table $TABLENAME --target-dir s3a://example-bucket/target-directory
```



Note: Entering sensitive data on the command line is inherently insecure. The data entered can be accessed in log files and other artifacts. Cloudera recommends that you use a credential provider to store credentials.

Using a Credential Provider to Secure S3 Credentials

You can run the `sqoop` command without entering the access key and secret key on the command line. This prevents these credentials from being exposed in the console output, log files, configuration files, and other artifacts. Running the command this way requires that you provision a credential store to securely store the access key and secret key. The credential store file is saved in HDFS.



Note: Using a Credential Provider does not work with MapReduce v1 (MRV1).

To provision credentials in a credential store:

1. Provision the credentials by running the following commands:

```
hadoop credential create fs.s3a.access.key -value access_key -provider jceks://hdfs/path_to_credential_store_file
hadoop credential create fs.s3a.secret.key -value secret_key -provider jceks://hdfs/path_to_credential_store_file
```

For example:

```
hadoop credential create fs.s3a.access.key -value foobar -provider jceks://hdfs/user/alice/home/keystores/aws.jceks
hadoop credential create fs.s3a.secret.key -value barfoo -provider jceks://hdfs/user/alice/home/keystores/aws.jceks
```

You can omit the `-value` option and its value. When the option is omitted, the command will prompt the user to enter the value.

2. Copy the contents of the `/etc/hadoop/conf` directory to a working directory.
3. Add the following to the `core-site.xml` file in the working directory:

```
<property>
```



```
<name>hadoop.security.credential.provider.path</name>
<value>jceks://hdfs/path_to_credential_store_file</value>
</property>
```

4. Set the HADOOP_CONF_DIR environment variable to the location of the working directory:

```
export HADOOP_CONF_DIR=path_to_working_directory
```

After completing these steps, you can run the sqoop command using the following syntax:

Import into a target directory in an Amazon S3 bucket while credentials are stored in a credential store file and its path is set in the core-site.xml.

```
sqoop import --connect $CONN --username $USER --password $PWD --table $TABLENAME --target-dir s3a://example-bucket/target-directory
```

You can also reference the credential store on the command line, without having to enter it in a copy of the core-site.xml file. You also do not have to set a value for HADOOP_CONF_DIR. Use the following syntax:

Import into a target directory in an Amazon S3 bucket while credentials are stored in a credential store file and its path is passed on the command line.

```
sqoop import -Dhadoop.security.credential.provider.path=jceks://hdfs/path-to-credential-store-file --connect $CONN --username $USER --password $PWD --table $TABLENAME --target-dir s3a://example-bucket/target-directory
```

Related Information

[Credential Management \(Apache Software Foundation\)](#)

Sqoop Import into Amazon S3

Import Data from RDBMS into an S3 Bucket

The --target-dir option must be set to the target location in the S3 bucket to import data from RDBMS into an S3 bucket.

Example command: Import data into a target directory in an Amazon S3 bucket.

```
sqoop import --connect $CONN --username $USER --password $PWD --table $TABLENAME --target-dir s3a://example-bucket/target-directory
```

Data from RDBMS can be imported into S3 as Sequence or Avro file format too.

Parquet import into S3 is also supported if the Parquet Hadoop API based implementation is used, meaning that the --parquet-configurator-implementation option is set to hadoop.

Example command: Import data into a target directory in an Amazon S3 bucket as Parquet file.

```
sqoop import --connect $CONN --username $USER --password $PWD --table $TABLENAME --target-dir s3a://example-bucket/target-directory --as-parquetfile --parquet-configurator-implementation hadoop
```

Import Data into S3 Bucket in Incremental Mode

The --temporary-rootdir option must be set to point to a location in the S3 bucket to import data into an S3 bucket in incremental mode.

Append Mode

When importing data into a target directory in an Amazon S3 bucket in incremental append mode, the location of the temporary root directory must be in the same bucket as the directory. For example: `s3a://example-bucket/temporary-rootdir` or `s3a://example-bucket/target-directory/temporary-rootdir`.

Example command: Import data into a target directory in an Amazon S3 bucket in incremental append mode.

```
sqoop import --connect $CONN --username $USER --password $PWD --table $TABLE_NAME --target-dir s3a://example-bucket/target-directory --incremental append --check-column $CHECK_COLUMN --last-value $LAST_VALUE --temporary-rootdir s3a://example-bucket/temporary-rootdir
```

Data from RDBMS can be imported into S3 in incremental append mode as Sequence or Avro file format. too

Parquet import into S3 in incremental append mode is also supported if the Parquet Hadoop API based implementation is used, meaning that the `--parquet-configurator-implementation` option is set to `hadoop`.

Example command: Import data into a target directory in an Amazon S3 bucket in incremental append mode as Parquet file.

```
sqoop import --connect $CONN --username $USER --password $PWD --table $TABLE_NAME --target-dir s3a://example-bucket/target-directory --incremental append --check-column $CHECK_COLUMN --last-value $LAST_VALUE --temporary-rootdir s3a://example-bucket/temporary-rootdir --as-parquetfile --parquet-configurator-implementation hadoop
```

Lastmodified Mode

When importing data into a target directory in an Amazon S3 bucket in incremental lastmodified mode, the location of the temporary root directory must be in the same bucket and in the same directory as the target directory. For example: `s3a://example-bucket/temporary-rootdir` in case of `s3a://example-bucket/target-directory`.

Example command: Import data into a target directory in an Amazon S3 bucket in incremental lastmodified mode.

```
sqoop import --connect $CONN --username $USER --password $PWD --table $TABLE_NAME --target-dir s3a://example-bucket/target-directory --incremental lastmodified --check-column $CHECK_COLUMN --merge-key $MERGE_KEY --last-value $LAST_VALUE --temporary-rootdir s3a://example-bucket/temporary-rootdir
```

Parquet import into S3 in incremental lastmodified mode is supported if the Parquet Hadoop API based implementation is used, meaning that the `--parquet-configurator-implementation` option is set to `hadoop`.

Example command: Import data into a target directory in an Amazon S3 bucket in incremental lastmodified mode as Parquet file.

```
sqoop import --connect $CONN --username $USER --password $PWD --table $TABLE_NAME --target-dir s3a://example-bucket/target-directory --incremental lastmodified --check-column $CHECK_COLUMN --merge-key $MERGE_KEY --last-value $LAST_VALUE --temporary-rootdir s3a://example-bucket/temporary-rootdir --as-parquetfile --parquet-configurator-implementation hadoop
```

Import Data into an External Hive Table Backed by S3

The AWS credentials must be set in the Hive configuration file (`hive-site.xml`) to import data from RDBMS into an external Hive table backed by S3. The configuration file can be edited manually or by using the advanced configuration snippets.

Both `--target-dir` and `--external-table-dir` options have to be set. The `--external-table-dir` has to point to the Hive table location in the S3 bucket.

Parquet import into an external Hive table backed by S3 is supported if the Parquet Hadoop API based implementation is used, meaning that the `--parquet-configurator-implementation` option is set to `hadoop`.

Example Commands: Create an External Hive Table Backed by S3

Create an external Hive table backed by S3 using HiveServer2:

```
sqoop import --connect $CONN --username $USER --password $PWD --table $TABLE_NAME --hive-import --create-hive-table --hs2-url $HS2_URL --hs2-user $HS2_USER --hs2-keytab $HS2_KEYTAB --hive-table $HIVE_TABLE_NAME --target-dir s3a://example-bucket/target-directory --external-table-dir s3a://example-bucket/external-directory
```

Create and external Hive table backed by S3 using Hive CLI:

```
sqoop import --connect $CONN --username $USER --password $PWD --table $TABLE_NAME --hive-import --create-hive-table --hive-table $HIVE_TABLE_NAME --target-dir s3a://example-bucket/target-directory --external-table-dir s3a://example-bucket/external-directory
```

Create an external Hive table backed by S3 as Parquet file using Hive CLI:

```
sqoop import --connect $CONN --username $USER --password $PWD --table $TABLE_NAME --hive-import --create-hive-table --hive-table $HIVE_TABLE_NAME --target-dir s3a://example-bucket/target-directory --external-table-dir s3a://example-bucket/external-directory --as-parquetfile --parquet-configurator-implementation hadoop
```

Example Commands: Import Data into an External Hive Table Backed by S3

Import data into an external Hive table backed by S3 using HiveServer2:

```
sqoop import --connect $CONN --username $USER --password $PWD --table $TABLE_NAME --hive-import --hs2-url $HS2_URL --hs2-user $HS2_USER --hs2-keytab $HS2_KEYTAB --target-dir s3a://example-bucket/target-directory --external-table-dir s3a://example-bucket/external-directory
```

Import data into an external Hive table backed by S3 using Hive CLI:

```
sqoop import --connect $CONN --username $USER --password $PWD --table $TABLE_NAME --hive-import --target-dir s3a://example-bucket/target-directory --external-table-dir s3a://example-bucket/external-directory
```

Import data into an external Hive table backed by S3 as Parquet file using Hive CLI:

```
sqoop import --connect $CONN --username $USER --password $PWD --table $TABLE_NAME --hive-import --target-dir s3a://example-bucket/target-directory --external-table-dir s3a://example-bucket/external-directory --as-parquetfile --parquet-configurator-implementation hadoop
```

S3Guard with Sqoop

The properties that enable S3Guard can be set through command line during Sqoop import.

Example command:

Import data into a target directory in Amazon S3 bucket and enable S3Guard.

```
sqoop import -Dfs.s3a.metadatastore.impl=org.apache.hadoop.fs.s3a.s3guard.DynamoDBMetadataStore -Dfs.s3a.s3guard.ddb.region=$BUCKET_REGION -Dfs.s3a.s3
```

```
guard.ddb.table.create=true --connect $CONN --username $USER --password $PWD
--table $TABLENAME --target-dir s3a://example-bucket/target-directory
```

Accessing Storage Using Microsoft ADLS Gen 2

These topics focused on Microsoft ADLS from the core Cloudera Enterprise documentation library can help you deploy, configure, manage, and secure clusters in the cloud. They are listed by broad category:

Note the following limitations:

- ADLS is not supported as the default filesystem. Do not set the default file system property (fs.defaultFS) to an abfss:// URI. You can use ADLS as secondary filesystem while HDFS remains the primary filesystem.
- Hadoop Kerberos authentication is supported, but it is separate from the Azure user used for ADLS authentication.
- Directory and file names should not end with a period. Paths that end in periods can cause inconsistent behavior, including the period disappearing. For more information, see [HADOOP-15860](#).

Configuring OAuth in Data Hub

To connect a DataHub cluster to ADLS Gen2 with OAuth, you must configure the Hadoop CredentialProvider or core-site.xml directly. Although configuring core-site.xml is convenient, it is insecure since the contents of core-site.xml are not encrypted. For this reason, Cloudera recommends using a credential provider.

Before you start, ensure that you have configured OAuth for Azure.

Configuring OAuth with core-site.xml

Before you begin

Configuring your OAuth credentials in core-site.xml is insecure. Cloudera recommends that you only use this method for development environments or other environments where security is not a concern.

Perform the following steps to connect your cluster to ADLS Gen2:

Procedure

1. In the Cloudera Manager Admin Console, search for the following property: Cluster-wide Advanced Configuration Snippet (Safety Valve) for core-site.xml. .
2. Add the following properties and values:

Table 9: OAuth Properties

Name	Value
fs.azure.account.auth.type	OAuth
fs.azure.account.oauth.provider.type	org.apache.hadoop.fs.azurebfs.oauth2.ClientCredsTokenProvider
fs.azure.account.oauth2.client.endpoint	Provide your tenant ID: https://login.microsoftonline.com/<Tenant_ID>/oauth2/token
fs.azure.account.oauth2.client.id	Provide your <Client_ID>
fs.azure.account.oauth2.client.secret	Provide your <Client_Secret>

What to do next

In addition, you can also provide account-specific keys. To do this, you need to add the following suffix to the key:

```
.<Account>.dfs.core.windows.net
```

Configuring OAuth with the Hadoop CredentialProvider

Before you begin

A more secure way to store your OAuth credentials is with the Hadoop CredentialProvider. When you submit a job, reference the CredentialProvider, which then supplies the OAuth information. Unlike the core-site.xml, the credentials are not stored in plain text.

The following steps describe how to create a credential provider and how to reference it when submitting jobs:

Procedure

1. Create a password for the Hadoop Credential Provider and export it to the environment:

```
export HADOOP_CREDSTORE_PASSWORD=password
```

2. Provision the credentials by running the following commands:

```
hadoop credential create dfs.adls.oauth2.client.id -provider jceks://hdfs/
user/USER_NAME/adls2keyfile.jceks -value client ID
hadoop credential create dfs.adls.oauth2.credential -provider jceks://h
dfs/user/USER_NAME/adls2keyfile.jceks -value client secret
hadoop credential create dfs.adls.oauth2.refresh.url -provider jceks://
hdfs/user/USER_NAME/adls2keyfile.jceks -value refresh URL
```

You can omit the -value option and its value and the command will prompt the user to enter the value.

For more details on the hadoop credential command, see [Credential Management \(Apache Software Foundation\)](#).

3. Export the password to the environment:

```
export HADOOP_CREDSTORE_PASSWORD=password
```

4. Reference the credential provider on the command line when you submit a job:

```
hadoop <command>
  -Ddfs.adls.oauth2.access.token.provider.type=ClientCredential \
  -Dhadoop.security.credential.provider.path=jceks://hdfs/user/USER_NAM
E/adls-cred.jceks \
  abfs[s]://<file_system>@<account_name>.dfs.core.windows.net/<path>/
<file_name>
```

Configuring Native TLS Acceleration

For ADLS Gen2, TLS is enabled by default using the Java implementation of TLS. For better performance, you can use the native OpenSSL implementation of TLS.

Perform the following steps to use the native OpenSSL implementation of TLS:

1. Verify the location of the OpenSSL libraries on the hosts with the following command:

```
whereis libssl
```

2. In the Cloudera Manager Admin Console, search for the following property: Gateway Client Environment Advanced Configuration Snippet (Safety Valve) for `hadoop-env.sh`.
3. Add the following parameter to the property:

```
HADOOP_OPTS="-Dorg.wildfly.openssl.path=<path to OpenSSL libraries> ${HADOOP_OPTS}"
```

For example, if the OpenSSL libraries are in `/usr/lib64`, add the following parameter:

```
HADOOP_OPTS="-Dorg.wildfly.openssl.path=/usr/lib64 ${HADOOP_OPTS}"
```

4. Save the change.
5. Search for the following property: HDFS Client Environment Advanced Configuration Snippet (Safety Valve) for `hadoop-env.sh`
6. Add the following parameter to the property:

```
HADOOP_OPTS="-Dorg.wildfly.openssl.path=<path to OpenSSL libraries> ${HADOOP_OPTS}"
```

For example, if the OpenSSL libraries are in `/usr/lib64`, add the following parameter:

```
HADOOP_OPTS="-Dorg.wildfly.openssl.path=/usr/lib64 ${HADOOP_OPTS}"
```

7. Save the change.
8. Restart the stale services.
9. [Deploy the client configurations.](#)
10. Verify that you configured native TLS acceleration successfully by running the following command from any host in the cluster:

```
hadoop fs -ls abfss://<container>@<account>.dfs.core.windows.net/
```

A message similar to the following should appear:

```
org.wildfly.openssl.SSL init
INFO: WFOPENSSL0002 OpenSSL Version OpenSSL 1.0.1e-fips 11 Feb 2013
```

The message may differ slightly depending on your operating system and OpenSSL version.

Importing Data into Microsoft Azure Data Lake Store (Gen1 and Gen2) Using Sqoop

Microsoft Azure Data Lake Store (ADLS) is a cloud object store designed for use as a hyper-scale repository for big data analytic workloads. ADLS acts as a persistent storage layer for CDH clusters running on Azure.

There are two generations of ADLS, Gen1 and Gen2. You can use Apache Sqoop with both generations of ADLS to efficiently transfer bulk data between these file systems and structured datastores such as relational databases. For more information on ADLS Gen 1 and Gen 2, see:

- [Microsoft ADLS Gen1 documentation](#)
- [Microsoft ADLS Gen2 documentation](#)

You can use Sqoop to import data from any relational database that has a JDBC adaptor such as SQL Server, MySQL, and others, to the ADLS file system.



Note: Sqoop export from the Azure files systems is not supported.

The major benefits of using Sqoop to move data are:

- It leverages RDBMS metadata to get the column data types
- It ensures fault-tolerant and type-safe data handling
- It enables parallel and efficient data movement

Prerequisites

The configuration procedure presumes that you have already set up an Azure account, and have configured an ADLS Gen1 store or ADLS Gen2 storage account and container. See the following resources for information:

- [Microsoft ADLS Gen1 documentation](#)
- [Microsoft ADLS Gen2 documentation](#)
- [Hadoop Azure Data Lake Support](#)

Authentication

To connect a cluster to ADLS with OAuth, you must configure the Hadoop CredentialProvider or core-site.xml directly. Although configuring the core-site.xml is convenient, it is insecure, because the contents of core-site.xml configuration file are not encrypted. For this reason, Cloudera recommends using a credential provider. For more information, see [Configuring OAuth in CDH](#).

You can also pass the credentials by providing them on the Sqoop command line as part of the import command.

```
sqoop import
-Dfs.azure.account.auth.type=...
-Dfs.azure.account.oauth.provider.type=...
-Dfs.azure.account.oauth2.client.endpoint=...
-Dfs.azure.account.oauth2.client.id=...
-Dfs.azure.account.oauth2.client.secret=...
```

For example:

```
sqoop import
-Dfs.azure.account.oauth2.client.endpoint=https://login.microsoftonline.com/$TENANT_ID/oauth2/token
-Dfs.azure.account.oauth2.client.id=$CLIENT_ID
-Dfs.azure.account.oauth2.client.secret=$CLIENT_SECRET
-Dfs.azure.account.auth.type=OAuth
-Dfs.azure.account.oauth.provider.type=org.apache.hadoop.fs.azurebfs.oauth2.ClientCredsTokenProvider
```

Sqoop Import into ADLS

To import data into ADLS from diverse data sources, such as a relational database, enter the Sqoop import command on the command line of your cluster. Make sure that you specify the Sqoop connection to the data source you want to import.

If you want to enter a password for the data source, use the `-P` option in the connection string. If you want to specify a file where the password is stored, use the `--password-file` option.

```
sqoop import
-Dfs.azure.account.auth.type=...
-Dfs.azure.account.oauth.provider.type=...
-Dfs.azure.account.oauth2.client.endpoint=...
-Dfs.azure.account.oauth2.client.id=...
-Dfs.azure.account.oauth2.client.secret=...
```

```
--connect... --username... --password... --table... --target-dir... --split-  
by...
```

ABFS example:

```
sqoop import  
-Dfs.azure.account.oauth2.client.endpoint=https://login.microsoftonline  
.com/$TENANT_ID/oauth2/token  
-Dfs.azure.account.oauth2.client.id=$CLIENT_ID  
-Dfs.azure.account.oauth2.client.secret=$CLIENT_SECRET  
-Dfs.azure.account.auth.type=OAuth  
-Dfs.azure.account.oauth.provider.type=org.apache.hadoop.fs.azurebfs.oauth2.  
ClientCredsTokenProvider  
--connect $CONN --username $USER --password $PWD --table $TABLENAME --targe  
t-dir abfs://$CONTAINER$ACCOUNT.dfs.core.windows.net/$TARGET-DIRECTORY --s  
plit-by $COLUMN_NAME
```

ADLS example:

```
sqoop import  
-Dfs.adl.oauth2.refresh.url=https://login.windows.net/$TENANT_ID/oauth2/toke  
n  
-Dfs.adl.oauth2.client.id=$CLIENT_ID  
-Dfs.adl.oauth2.credential=$CLIENT_SECRET  
-Dfs.adl.oauth2.access.token.provider.type=ClientCredential  
--connect $CONN --username $USER --password $PWD --table $TABLENAME --targe  
t-dir adl://$TARGET-ADDRESS/$TARGET-DIRECTORY --split-by $COLUMN_NAME
```