

Data Hub

Accessing Clusters

Date published: 2019-12-17

Date modified: 2023-06-27

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

Enabling admin and user access to Data Hubs.....	4
Understanding Data Hub cluster details.....	4
Accessing Cloudera Manager, cluster UIs, and endpoints.....	7
Accessing Data Hub cluster via SSH.....	8
Set workload password.....	10
Finding your workload user name.....	11
Retrieving keytabs for workload users.....	11
Running workloads.....	12

Enabling admin and user access to Data Hubs

Data Hub resource roles can be assigned on the scope of a specific Data hub cluster.

When you grant access to admins and users of a Data Hub, consider the following guidelines:

- Any user or group that needs to access a specific Data Hub needs the EnvironmentUser role at the scope of the environment where that Data Hub is running.
- A user with the DataHubCreator (or higher) account role can create Data Hubs.
- The user who creates a Data Hub gets the Owner role for that Data Hub.
- The Owner of the Data Hub cluster can grant others access to the cluster. The following roles can be assigned:
 - Owner - Grants the permission to manage the Data Hub cluster in CDP and delete it. It does not grant any cluster-level access (such as access to Cloudera Manager).
 - DataHubAdmin (Technical Preview) - Grants administrative rights over the Data Hub cluster.

The roles are described in detail in *Resource roles*. The steps for assigning the roles are described in *Assigning resource roles to users* and *Assigning resource roles to groups*.

Related Information

[Resource roles](#)

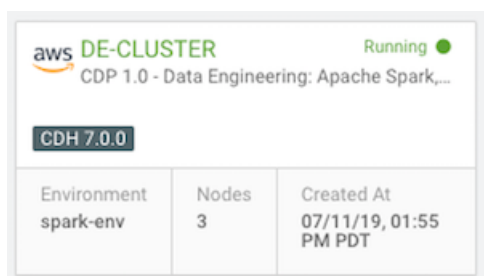
[Assigning resource roles to users](#)

[Assigning resource roles to groups](#)

Understanding Data Hub cluster details

To access information about Data Hub clusters, navigate to the Data Hub Clusters service or to the Management Console service > Data Hub Clusters.

Each cluster is represented by a tile that includes basic information about the cluster:



Click the cluster tile to display more information:

Data Hubs / long-run-dh / Event History

The screenshot shows the 'long-run-dh' cluster details page. At the top, there's a 'Stop' button and an 'Actions' dropdown. Below this is a summary bar with four columns: STATUS (Running), NODES (8), CREATED AT (04/29/21, 10:08 AM CDT), and CLUSTER TEMPLATE (7.2.8 - Data Engineering: HA: Apache Spark, Apache Hive, Apache Oozie). The 'Environment Details' section shows NAME (perf-long-run-env), DATA LAKE (perf-long-run-dl), CREDENTIAL (perf-long-run-env-cred), REGION (us-west-2), and AVAILABILITY ZONE (us-west-2a). The 'Services' section lists various services like CM-UI, Data Analytics Studio, HUE, Job History Server, Livy Server, Name Node, Queue Manager, Resource Manager, Spark History Server, and Zeppelin. The 'Cloudera Manager Info' section shows CM URL, CM VERSION (7.4.0), PLATFORM VERSION (7.2.8-1.cd7.2.8.p0.11560957), and LOGS. At the bottom, there's an 'Event History' section with filters for Show All, Autoscale, and Cluster, and a 'DOWNLOAD' button.

The summary bar includes the following information:

Item	Description
Status	Current cluster status. When a cluster is healthy, the status is Running.
Nodes	The current number of cluster nodes.
Created at	The date and time when the cluster was created. The format is MM/DD/YY, Time AM/PM Timezone. For example: 06/20/19, 7:56 AM PDT.
Cluster Template	The name of the cluster template used to create the cluster.

Environment Details

You can find information about the cluster cloud provider environment under Environment Details:

Item	Description
Cloud Provider	The logo of the cloud provider on which the cluster is running.
Name	The name of the environment used to create the cluster.
Data Lake	The name of a Data Lake to which the cluster is attached.
Credential	The name of the credential used to create the cluster.
Region	The region in which the cluster is running in the cloud provider infrastructure.
Availability Zone	The availability zone within the region in which the cluster is running.

Services

In the Services section, you will find links to cluster UIs. The exact content depends on what components are running on your cluster.


Cloudera Manager Info

The Cloudera Manager Info section provides the following information:

Item	Description
CM URL	Link to the Cloudera Manager web UI.
CM Version	The Cloudera Manager version which the cluster is currently running.
Platform Version	The Cloudera Runtime version which the cluster is currently running.

Event History and other tabs

Under Cloudera Manager, the Event History tab is displayed by default. You can also click the other tabs to view additional cluster information.

Item	Description
Event History	Shows events logged for the cluster, with the most recent event at the top. The Download option allows you to download the event history.
Hardware	This section includes information about your cluster instances: instance names, instance IDs, instance types, their status, fully qualified domain names (FQDNs), and private and public IPs. If you click on the  , you can access more information about the instance, storage, image, and packages installed on the image.
Tags	This section lists keys and values of the user-defined tags, in the same order as you added them.
Endpoints	This section includes the endpoints for various cluster services.
Recipes	This section includes recipe-related information. For each recipe, you can see the host group on which a recipe was executed, recipe name, and recipe type.
Repository Details	This section includes Cloudera Manager and Cloudera Runtime repository information, as you provided when creating a cluster.
Image Details	This section includes information about the prewarmed or base image that was used for the cluster nodes.
Network	This section includes information about the names of the network and subnet in which the cluster is running and the links to related cloud provider console.
Cloud Storage	This section provides information about the base storage locations used for YARN and Zeppelin.
Database	This section provides information about any external managed database you might have created for the cluster.
Telemetry	This section provides information about logging, metering, cluster log collection, and other analytics.

Actions Menu

Click Show Cluster Template on the Actions menu to review the template used in cluster creation. Click Show CLI Command to review the CDP CLI command used to create the cluster (which you can copy to create similar clusters via the CDP CLI). Select Manage Access to manage access to the cluster.

You can also perform some basic Data Hub management functions from the Actions menu, such as resizing, retrying, and repairing the cluster, as well renewing the host certificate.

Accessing Cloudera Manager, cluster UIs, and endpoints

Cluster UIs and endpoints can be accessed from cluster details.

Required role: EnvironmentUser at the scope of the environment where the Data Hub is running, but Cloudera Manager access is read-only. EnvironmentAdmin grants a limited administrator role in Cloudera Manager. DatahubAdmin or the Owner of the Data Hub can access cluster details, but access to Cloudera Manager is read-only.

To access cluster UIs and endpoints navigate to the Data Hub Clusters service and click the tile for your cluster. This opens the cluster details page, which lists the URLs for the cluster UIs and endpoints:

Data Hubs / [aws-ec2-us-west-2](#) / Endpoints

Click the URL for the service that you would like to access and you will be logged in automatically with your CDP credentials. All of the UIs and endpoints are accessible via the Knox gateway. The URLs listed connect you to a chosen service via Knox, and Knox securely passes your CDP credentials.

Credentials to use for logging in

The following table lists the credentials to use to access clusters:

Method	URL	Credentials
Cloudera Manager web UI	Access from the URL listed in cluster details > Services section.	You do not need to provide any credentials. You are automatically logged in with your CDP credentials. When accessing a Data Hub cluster via Cloudera Manager, you assume the Configurator role.
All cluster web UIs	Access from the URLs listed in cluster details.	You do not need to provide any credentials. You are automatically logged in with your CDP credentials.

Method	URL	Credentials
Data Analytics Studio (DAS)	Access from the URLs listed in cluster details.	<p>Access DAS with your workload user name and workload password.</p> <p>When accessing CDP for the first time, you must set a workload password. For more instructions on how to set your workload password, refer to Set or reset workload password.</p> <p>For instructions on how to find your workload user name, refer to Finding your workload user name.</p>
All cluster endpoints	Access by using the API endpoint listed in cluster details > Endpoints tab.	<p>Access all cluster API endpoints (such as JDBC and ODBC) with your workload user name and workload password.</p> <p>When accessing CDP for the first time, you must set a workload password. For more instructions on how to set your workload password, refer to Set or reset workload password.</p> <p>For instructions on how to find your workload user name, refer to Finding your workload user name.</p> <p>For information on how to set up a connection from a business intelligence tool such as Tableau, refer to Configuring JDBC for Impala and Configuring ODBC for Impala.</p>

Security exception

The first time you access the UIs, your browser will attempt to confirm that the SSL Certificate is valid. Since CDP automatically generates a certificate with self-signed CA, your browser will warn you about an untrusted connection and ask you to confirm a security exception. Depending on your browser, perform the steps below to proceed:

Browser	Steps
Firefox	Click Advanced > Click Add Exception... > Click Confirm Security Exception
Safari	Click Continue
Chrome	Click Advanced > Click Proceed...

Accessing Data Hub cluster via SSH

You can use SSH to access cluster nodes via a command line client.

Method	Credentials
Root SSH access to cluster VMs	<p>Required role: None</p> <p>CDP administrators can access cluster nodes as cloudbreak user with the SSH key provided during environment creation.</p> <p>On Mac OS, you can use the following syntax to SSH to the VM:</p> <pre>ssh -i "privatekey.pem" cloudbreak@publicIP</pre> <p>For example:</p> <pre>ssh -i "testkey-kp.pem" cloudbreak@90.101.0.132</pre> <p>On Windows, you can access your cluster via SSH by using an SSH client such as PuTTY. For more information, refer to How to use PuTTY on Windows.</p>
Non-root SSH access to cluster VMs	<p>Required role: Any user who has access to the environment (EnvironmentUser, DataSteward, and EnvironmentAdmin) can access Data Hubs via SSH.</p> <p>All authorized users can access cluster nodes via SSH using either a private key that is paired with the user's public key, or with their workload user name and workload password.</p> <p>For SSH access through a workload user name and password:</p> <p>When accessing CDP for the first time, you must set a workload password. The password also needs to be reset each time you are added to a new environment.</p> <p>For more information about workload passwords and instructions for setting/resetting it, refer to Set or Reset Workload Password.</p> <p>For instructions on how to find your workload user name, refer to Finding Your Workload User Name.</p> <p>On Mac OS, you can use the following syntax to SSH to the VM:</p> <pre>\$ ssh workload-user@publicIP</pre> <p>For example:</p> <pre>\$ ssh jsmith@190.101.0.132</pre> <p>To SSH to a cluster using the private key file that pairs with the public key associated with a user, use the ssh utility:</p> <pre>\$ ssh -i path-to-private-key-file user@nodeIPAddress</pre> <p>For example:</p> <pre>\$ ssh -i ~/.ssh/my-private.key jsmith@192.12.141.12</pre> <p>On Windows, you can access your cluster via SSH by using an SSH client such as PuTTY. For more information, refer to How to use PuTTY on Windows.</p>

Set workload password

A workload password is used to access Data Hub clusters via SSH, endpoints such as JDBC/ODBC, and some UIs such as DAS.

Required role: All users can manage their workload passwords from the account management page. All users can manage their workload password from CDP CLI, but this action requires an API access key, which can only be generated by users with the IAMUser role. As a CDP administrator or PowerUser, you can manage the workload password for all user accounts.

The workload password is independent from your SSO password.

To set or reset the password:

1. Sign in to the CDP web interface.
2. Navigate to Management Console > User Management.
3. Search for your user name, and then click your user name:

User Management

Users Identity Providers

Search: dbialek Type: All Identity Provider: All Clear all Actions

Type	Name	Email	Identity Provider	Workload User Name	Access Keys
	Dominika Bialek	dbialek@cloudera.com.staging	cloudera-ss	csso_dbialek	1 active

Displaying 1 - 1 of 1 < 1 > 25 / page


4. Click Set Password for User:

Users / Dominika Bialek

Name	Dominika Bialek
Email	dbialek@cloudera.com.staging
Workload User Name	csso_dbialek
CRN	crn:altus:iam:us-west-1:9d74eee4-1cad-45d7-b645-7ccf9edbb73d:user:1fa...
Identity Provider	cloudera-ss
Last Interactive Login	10/24/2019 9:54 AM PDT
Profile Management	View profile
FreelPA Password	Set FreelPA Password


Access Keys Roles Resources Groups

- In the dialog box that appears, enter the workload password twice:

* Password 

* Confirm Password

Environment

All 

Select environments you want to set or update your password in. Leave blank to update your password in all environments.

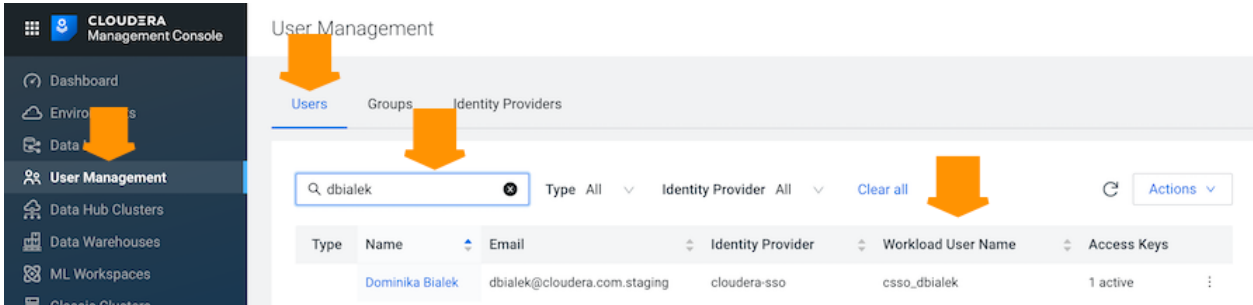
[Set Workload Password](#)

- Click Set Workload Password.

Finding your workload user name

Once you have reset your workload password, locate your workload user name.

To check your workload user name, select Management Console > User Management > Users, find your user name, and then find your Workload User Name:



User Management

Users Groups Identity Providers

Search: dbialek Type: All Identity Provider: All Clear all Actions

Type	Name	Email	Identity Provider	Workload User Name	Access Keys
	Dominika Bialek	dbialek@cloudera.com.staging	cloudera-ssso	csso_dbialek	1 active

Retrieving keytabs for workload users

A keytab file stores long-term keys for a principal in Kerberos. You can generate a keytab either through the Management Console user interface or the CDP CLI.

About this task

Required roles: All users can retrieve their keytabs from the account management page. All users can retrieve their keytabs from CDP CLI, but this action requires an API access key, which can only be generated by users with the IAMUser role. As a CDP administrator or PowerUser, you can retrieve the keytab for all user accounts.

You may need to generate a keytab for a workload user in certain Data Hub use cases, for example long-running Spark streaming jobs, which require a keytab as a long-lived credential; or NiFi flows requiring a keytab to write data into HBase.



Note: Keytabs are scoped to an environment, whereas workload passwords are the same for every environment. A keytab is, however, tied to the workload password. If you change the workload password, you must retrieve a new keytab. When you change a workload password, retrieve the keytab only after the user sync operation is complete. For more information on user sync, see *Assigning resources to users*.

Procedure

You can retrieve a keytab either in the Management Console or in the CDP CLI:

- Management Console:
 - a. Click User ManagementUsers and then search for and select the Name of the user that you want to get a keytab for.
 - b. Click ActionsGet Keytab.
 - c. Select the environment in which the Data Hub cluster is running and then click Download.
 - d. Save the keytab file in a location of your choice.

Once you have downloaded the keytab file, you can copy it to the machine on which the cluster runs and use the keytab to authenticate as the workload user principal, or point to the keytab file when running a Spark job or other job that requires a keytab.

- CDP CLI:
 - a. Keytab retrieval (get-keytab) is within the environments module. Run `cdp environments get-keytab help` for more information. You will need to pass the environment name and an actor CRN:

```
cdp environments get-keytab \  
--environment-name=EnvironmentName \  
--actor-crn=ActorCrn
```

- b. The output of the command is a base64-encoded representation of a keytab. The contents of the output must be base64 decoded and saved to a file for it to work as a keytab.



Note: There are ways to generate keytabs with utilities outside of CDP, such as `ipa-getkeytab` or `ktutil`. Cloudera recommends against using these methods as they may not work as expected. For example, `ipa-getkeytab` creates a keytab that may work but only temporarily.

Related Information

[Assigning resources to users](#)

[CLI client setup](#)

Running workloads

Once your cluster is running, refer to the Cloudera Runtime and Data Hub documentation for information on how to run workloads using services such as Hive, Spark, Impala, Hue, or Kafka.

Related Information

[Data Access](#)

[Data Science](#)

[Streaming](#)

[Flow Management](#)

[Operational Database](#)