

Hortonworks DataPlane Service (DPS)

DLM Administration

(October 31, 2017)

Hortonworks DataPlane Service (DPS™): DLM Administration

Copyright © 2016-2017 Hortonworks, Inc. All rights reserved.

Please visit the [Hortonworks Data Platform](#) page for more information on Hortonworks technology. For more information on Hortonworks services, please visit either the [Support](#) or [Training](#) page. Feel free to [contact us](#) directly to discuss your specific needs.

Hortonworks reserves the right to change any products described herein at any time, and without notice. Hortonworks assumes no responsibility or liability arising from the use of products described herein, except as expressly agreed to in writing by Hortonworks.

Trademark: Hortonworks DataPlane Service and DPS are trademarks of Hortonworks, Inc. 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 HORTONWORKS, HORTONWORKS DOES NOT MAKE OR GIVE ANY REPRESENTATION, WARRANTY, OR COVENANT OF ANY KIND, WHETHER EXPRESS OR IMPLIED, IN CONNECTION WITH HORTONWORKS TECHNOLOGY OR RELATED SUPPORT PROVIDED IN CONNECTION THEREWITH. HORTONWORKS DOES NOT WARRANT THAT DPS WILL OPERATE UNINTERRUPTED OR THAT IT WILL BE FREE FROM DEFECTS OR ERRORS, THAT IT WILL PROTECT YOUR DATA FROM LOSS, CORRUPTION OR UNAVAILABILITY, OR THAT DPS WILL MEET ALL OF CUSTOMER'S BUSINESS REQUIREMENTS. WITHOUT LIMITING THE FOREGOING, AND TO THE MAXIMUM EXTENT PERMITTED BY APPLICABLE LAW, HORTONWORKS 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.

Open Source: Open source software licensing information for third party software used in connection with or included in the DataPlane Services software may be found in documentation accompanying the delivery of Hortonworks DataPlane Services software.

HORTONWORKS PROVIDES SEPARATELY LICENSED CODE, INCLUDING ALL OPEN SOURCE SOFTWARE, TO CUSTOMER WITHOUT WARRANTIES OF ANY KIND. HORTONWORKS DISCLAIMS ANY AND ALL EXPRESS AND IMPLIED WARRANTIES WITH RESPECT TO SEPARATELY LICENSED CODE, INCLUDING BUT NOT LIMITED TO IMPLIED WARRANTIES OF TITLE, NON-INFRINGEMENT, MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE.

Table of Contents

1. Hortonworks Data Lifecycle Manager Terminology	1
1.1. DPS Platform	1
1.2. Data Lifecycle Manager (DLM) Service	1
2. Replicating Data	3
2.1. Enable Snapshots on Directories	3
2.2. Pair Clusters for Replication	4
2.3. Create a Replication Policy	5
2.4. View Job Status	7
2.5. Pairing Guidelines and Considerations	8
2.6. Policy Guidelines and Considerations	9
2.7. How Pairing Works in Data Lifecycle Manager	10
2.8. How Policies Work in Data Lifecycle Manager	10
3. Understanding the UI	11
3.1. Overview Page	11
3.1.1. Cluster Health Panel	11
3.1.2. Policies Panel	11
3.1.3. Jobs Panel	12
3.1.4. Recent Issues Panel	13
3.1.5. Clusters Map	13
3.1.6. Issues & Updates Table	14
4. Roles Required to Work with DPS Services	16
4.1. Roles for DPS Services	16
4.1.1. DPS Admin	16
4.1.2. Infra Admin	16
4.1.3. Data Steward	17
4.2. Roles Required for Installation and Troubleshooting	17
4.2.1. Ambari Admin	17
4.2.2. Cluster Admin	17
4.3. Required Roles by Task	17
4.3.1. DPS Platform Tasks and Required Roles	18
4.3.2. Data Lifecycle Manager Tasks and Required Roles	18
4.3.3. Data Steward Studio	19
5. Replication Concepts	20
5.1. Snapshot Replication Between HDP Clusters	20
6. Configuring Security for DLM	22
6.1. Configure Authorization for DLM on Secured Clusters	22
6.2. Enable Authorization for DLM on Unsecured Clusters	22
7. Failing Over Manually	24
7.1. Make the Destination Cluster the New Source	24
7.2. Activate a New Destination Cluster	25
8. Troubleshooting DLM	27
8.1. Endpoint is unreachable	27
8.2. Ranger UI does not display deny policy items	27

1. Hortonworks Data Lifecycle Manager Terminology

You should be familiar with the terminology used in Hortonworks DataPlane Service (DPS) and in the Data Lifecycle Manager (DLM) service that interfaces with the DPS infrastructure.

1.1. DPS Platform

DPS Platform is a UI service platform from which you can manage and monitor various services on multiple Hortonworks Data Platform (HDP) Hadoop clusters. You can install DPS Platform on an HDP cluster or remote to the cluster.

Hortonworks DataPlane Service (DPS) service	<p>The family of components that include the DPS Platform service platform and all services that plug into it.</p> <p>An autonomous component in the DPS environment. DPS Platform is a service, as is each component that is enabled and managed through DPS Platform, such as Data Lifecycle Manager (DLM). DPS Platform and each of its plugin services must be installed as Docker containers.</p>
cluster	<p>A typical HDP Hadoop cluster. See the Cluster Planning guide for details.</p> <p>The cluster hosts the various Systems of Record (SoRs) for metadata (Apache Hive, Apache Atlas, Apache Ranger, HDFS, and so on) that DPS Platform and associated plugin services rely on. In an on-premise environment, a cluster often equates to a data center. However, a single data center can contain multiple HDP Hadoop clusters.</p>

1.2. Data Lifecycle Manager (DLM) Service

DLM is a UI service that is enabled through DPS Platform. From the DLM UI you can create and manage replication and disaster recovery policies and jobs.

DLM Engine	Also referred to as the Beacon engine, this is the replication engine that is required for Data Lifecycle Manager. The DLM Engine must be installed as a management pack on each cluster that is to be used in data replication jobs. The engine maintains, in a configured database, information about clusters and policies that are involved in replication.
data center	The facility that contains the computer, server, and storage systems and associated infrastructure, such as routers, switches, and so forth. Corporate data is stored, managed, and distributed from the data center.

	In an on-premise environment, a data center is often composed of a single Hadoop cluster.
policy	A set of rules applied to a replication relationship. The rules include which clusters serve as source and destination, the type of data to replicate, the schedule for replicating data, and so on.
job	An instance of a policy that is running or has run.

2. Replicating Data

After enabling the clusters for use with DLM, you must pair the clusters, and then create replication policies between the source and destination clusters. You can also enable snapshot functionality on the clusters.

Any source or destination cluster that you want to use in a replication relationship with Data Lifecycle Manager (DLM) must be managed by Apache Ambari and enabled for DLM through DPS Platform.

2.1. Enable Snapshots on Directories

You can use snapshot-enabled (snapshottable) HDFS directories for replication on the source and destination clusters. You must manually enable the directories for snapshots. Note that you do not need to create these directories for Apache Hive.

See [Replication Concepts](#) for an overview of how snapshots work and considerations when implementing snapshots.

About This Task

You must have HDFS superuser access on the source and destination clusters to perform this task.

After you create the snapshot-enabled directories and begin replication jobs, snapshots are created on the source and destination clusters. The snapshots are maintained in directories named `.snapshot`.

The three most recent snapshots are retained by default. Older snapshots are automatically deleted.

Steps

1. As HDFS superuser, log into a terminal on the *source* cluster.
2. Create a directory in which to store snapshots on the source:

```
$ hadoop dfs -mkdir <source directory path>
```

3. Set the source directory permissions to read/write/execute (777) for everyone:

```
hdfs dfs -chmod -R 777 <source directory path>
```

4. Enable snapshots on the source directory:

```
$ hdfs dfsadmin -allowSnapshot <source directory path>
```

5. If the Beacon (DLM Engine) user is not configured as HDFS superuser, then give ownership of the directory to the DLM Engine:

```
hadoop dfs -chown -R beacon <source directory path>
```

6. As HDFS superuser, log in to a console on the *destination* cluster and repeat Step 2 through Step 5 to create a target directory to store snapshots on the destination.

2.2. Pair Clusters for Replication

You must pair two clusters to communicate with one another before you can replicate data between the clusters. After the clusters are paired, you can create a replication policy, during which you establish which are the source and destination clusters in the replication relationship.

Prerequisites

Before pairing, clusters must have been enabled for use with Data Lifecycle Manager from the DPS Platform UI.

The clusters you want to pair must have symmetrical security configurations for Kerberos, LDAP, Ranger, Knox, HA, etc.

About This Task

You must be logged in with the DLM Infra Admin role to perform this task.

Steps

1. In the DLM navigation pane, click **Pairings**.
2. Click **Add Pairing**.

The Create Pairing page displays, showing the clusters that are enabled for replication.



Tip

You can place the cursor over the cluster name to display the cluster location.

3. Click one of the cluster names.

All clusters available to be paired with the cluster you selected display in a second column.



Tip

If a cluster displays but cannot be selected, it is already paired with the cluster you selected in the first column.

4. Click a cluster in the second column.
5. Click **Start Pairing**.
A progress bar displays.
6. Repeat steps to pair additional clusters.

Next Steps

[Create a Replication Policy \[5\]](#)

More Information

[Pairing Guidelines and Considerations \[8\]](#)

[How Pairing Works in Data Lifecycle Manager \[10\]](#)

2.3. Create a Replication Policy

You must create a policy to assign the rules for the replication job (instance of a policy) that you want to execute. You can set rules such as the type of data to replicate, the time and frequency of replication, the bandwidth allowed for a job, and so forth. During replication, data and associated file metadata or table structures or schemas are also replicated

Prerequisites

- The clusters you want to include in the replication policy must have been paired already.
- You must ensure that the clusters you select are healthy before you start a policy instance (job).
- On destination clusters, the DLM Engine must have been granted write permissions on folders being replicated.
- The target folder or database on the destination cluster must either be empty or not exist prior to starting a new policy instance.

About This Task

- You must use the DLM Infrastructure Admin role to perform this task.
- You cannot modify a policy after it is created.
To change a policy, you must create a new policy with the new settings.
- DLM does not support update of any cluster endpoints (HDFS, Hive, Ranger, or DLM Engine). If an endpoint must be modified, contact Hortonworks support for assistance.
- The first time you execute a job with data that has not been previously replicated, Data Lifecycle Manager creates a new folder or database and bootstraps the data.



Important

During a bootstrap operation, all data is replicated from the source cluster to the destination. As a result, the initial execution of a job can take a significant amount of time, depending on how much data is being replicated, network bandwidth, and so forth.

After initial bootstrap, data replication is performed incrementally, so only updated data is transferred. Data is in a consistent state only after incremental replication has captured any new changes that occurred during bootstrap.

Steps

1. In the DLM navigation pane, click **Policies**.

The Replication Policies page displays a list of any existing policies.

2. Click **Add>Policy**.
3. Enter or select the following information:

Field	Description	Additional Information
Policy Name	The policy name that will display in the UI	Maximum length of 64 characters. Spaces, dashes, and underscores are the only special characters allowed.
Description	Any useful information to identify the policy or its use	
Service	Hive or HDFS replication	For Hive replication, a corresponding Hive database structure must exist on the destination. For HDFS, the corresponding file system structure is created when the first replication job executes.
Source Cluster	The cluster that contains the data to be replicated	If the cluster you want is not listed, you need to enable the cluster for DLM.
Destination Cluster	The cluster to which the source data will be replicated	If the cluster you want is not listed, you need to enable the cluster for DLM.
Select a Folder Path (Only if HDFS is selected)	The HDFS directories available to browse and to select for replication	The Infra Admin role has read privileges, in the DLM UI only, for all HDFS directories on the source and destination clusters. Clusters must be paired before you can browse HDFS directories in DLM.
Select Database (Only if Hive is selected)	The internal or external databases available to browse and to select for replicated	The Infra Admin role has read privileges, in the DLM UI only, for all databases on the source and destination clusters.

4. Select how you want the job to run:

When setting the schedule, consider requirements such as RPO and RTO, network bandwidth, and so forth.

Field	Description	Additional Information
Repeat	How often you want the job to run	Choices are weeks, days, hours, or minutes. For a Hive replication policy, set the frequency so that changes are replicated often enough to avoid overly large copies.
Start and End Dates	The dates you want the job to start (required) and end (optional)	If you do not set an end date, the job runs at the set time and frequency until the job is manually cancelled.
Start Time	24-hour clock	

5. Enter or select the **Advanced Properties**:

Field	Description	Additional Information
Queue Name (Optional)	The YARN queue you want to use to prioritize job scheduling across multiple tenants	If no queue is entered, DLM defaults to the YARN queue identified in the Ambari View for YARN Capacity Scheduler. You can enter one queue name per policy.
Maximum Bandwidth (Optional)	The maximum bandwidth to be used when running a job based on this policy	Enables you to restrict the amount of data throughput to the specified value. Enter a number in megabytes per second (MBps).

6. Click **Review** and verify that the settings are correct.



Important

After a policy is created, it cannot be modified.

7. Click **Submit**.

A message appears, stating that the submission was successful.

Next Steps

Verify that the replication job is running as intended.

More Information

[Policy Guidelines and Considerations \[9\]](#)

[How Policies Work in Data Lifecycle Manager \[10\]](#)

2.4. View Job Status

You can check job status from several places in the DLM UI.

About This Task

You must use the DLM Infrastructure Admin role to perform this task.

Steps: View Job Status from the Policies Page

You can view the status and other information about policies and associated jobs from the Policies page. All jobs (policy instances) can be viewed from this page, regardless of status.

The Policies Page can display up to 200 policies.

1. In the navigation pane, click **Policies**.
2. Click the  HDFS or  Hive icon to display the type of policies you want to see.
3. Locate the policy associated with the job that you want to view by doing one of the following:
 - Browse the list to find the name of the policy.
 - Enter full or partial terms in the search field.
4. For the policy you located, click  in the Prev Jobs column to open or close the list of jobs associated with the policy.

A maximum of 10 jobs displays per page.

5. Click   to see the next or the previous list of jobs.

Steps: View Job Status from the Overview Page

The Overview page displays jobs that are either in progress or have not succeeded. While jobs are executing, they display in the list with a status of In Progress. If the job succeeds, it disappears from the list. Successful jobs can be viewed from the Policies page.

1. In the navigation pane, click **Overview**.
2. Browse the Issues & Updates list to locate the policy for the job you want status for.
3. View the Job Status column for the policy.
4. If the job did not succeed, click  next to the job status to view the job log.
5. Optionally, see information about previous job runs:
 - a. Click the dots in the Policy History column.

The policy displays in the Policies page.
 - b. Click the dots in the Prev Job column.

A list of jobs related to the selected policy displays, showing up to the last 10 jobs.

Steps: View Job Status from the Notifications Page

1. From any page in Data Lifecycle Manager, click  to display the last five job alerts.
2. From the Notifications dialog box, click **View All** to open the Notifications page, showing all previous notifications.

2.5. Pairing Guidelines and Considerations

- Pairings are bidirectional, so either cluster in a pair can serve as the source or as the destination in a replication policy.

You establish through the policy which cluster is the source and which is the destination.

- For pairing to succeed, host name resolution must work between all the nodes involved (DPS Platform host and all cluster hosts.)

For example, pairing in DLM fails if the DLM engines on the clusters being paired cannot resolve each other's host name.

- You can only pair clusters that have been registered with DPS Platform and enabled for use with Data Lifecycle Manager.
- Cluster security configurations must be symmetrical to pair clusters.
- The DLM Engine must be configured as Ranger administrator on all DLM-enabled clusters.
- Clusters must be paired before you can browse HDFS directories in DLM.

More Information

[Pair Clusters for Replication \[4\]](#)

[How Pairing Works in Data Lifecycle Manager \[10\]](#)

2.6. Policy Guidelines and Considerations

- When creating a schedule for a Hive replication policy, you should set the frequency so that changes are replicated often enough to avoid overly large copies.
- On the Create Policy page, the only requirement for clusters to display in the Source Cluster or Destination Cluster fields is that they are DLM-enabled.

You must ensure that the clusters you select are healthy before you start a policy instance (job).

- Any user with access to the DLM UI has the ability to browse, within the DLM UI, the folder structure of any clusters enabled for DLM.

Therefore, the DPS Admins and the Infra Admins can see folders, files, and databases in the DLM UI that they might not have access to in HDFS.

The DataPlane Admin and Infra Admin cannot view from the DLM UI the content of files on the source or destination clusters. Nor do these admins have the ability to modify or delete folders or files that are viewable from the DLM UI.

- The first time you execute a job (an instance of a policy) with data that has not been previously replicated, Data Lifecycle Manager replicates all data and associated metadata from the source cluster to the destination. As a result, the initial execution of a job can take a significant amount of time. After initial bootstrap, data replication is done incrementally, so only updated data is transferred.
- Policies cannot be modified after they are created.
- If you want to run a single job, you can either start the job and then cancel it manually.
- Achieving a one-hour Recovery Point Objective (RPO) depends on how you set up your replication jobs and the configuration of your environment:
 - Select data in sizes that replicate within 30 minutes.
 - Set replication frequency to 45 minutes or fewer.
 - Ensure that network bandwidth is sufficient, so that data can move fast enough to meet your RPO.
 - Consider the rate of change of data being replicated.
- Specify bandwidth per map, in MBps.

Each map is restricted to consume only the specified bandwidth. This is not always exact. The map throttles back its bandwidth consumption during a copy in such a way that the *net* bandwidth used tends towards the specified value.

- Ensure that the frequency is set so that a job finishes before the next job starts.

Jobs based on the same policy cannot overlap. If a job is not completed before another job starts, the second job does not execute and is given the status Ignored. If a job is consistently ignored, you might need to modify the frequency of the job.

More Information

[Create a Replication Policy \[5\]](#)

[How Policies Work in Data Lifecycle Manager \[10\]](#)

2.7. How Pairing Works in Data Lifecycle Manager

Pairing verifies compatibility and establishes communication between the clusters that you want to use as source or destination clusters in a replication relationship. Pairings are bi-directional, so either cluster in a pair can serve as the source or as the destination in a replication policy.

After pairing is complete, you can create a replication policy between the paired clusters, and you can establish through the policy which cluster is the source and which is the destination.

Pairings can only be performed on clusters that have been registered with Data Lifecycle Manager. If a cluster you want to use is not visible, you need to register it from the DPS Platform UI, logged in as DataPlane Admin.

More Information

[Pairing Guidelines and Considerations \[8\]](#)

2.8. How Policies Work in Data Lifecycle Manager

In Data Lifecycle Manager, you create policies to establish the rules you want applied to your replication and disaster recovery jobs. The policy rules you set can include which cluster is the source and which the destination, what data is replicated, what day and time the replication job occurs, the frequency of job execution, bandwidth restrictions, etc.

When scheduling how often you want a replication job to run, you should consider the recovery point objective (RPO) of the data being replicated; that is, what is the acceptable lag time between the active site and the replicated data on the destination. Data Lifecycle Manager supports a one-hour RPO: data is preserved up to one hour prior to the point of data recovery. To meet a one-hour RPO, you must consider how long it takes to replicate the selected data, how often the data is replicated, and network bandwidth capabilities.

As an example, if you have a set of data that you expect to take 15 minutes to replicate, then to meet a one-hour RPO, you would schedule the replication job to occur no more often than every 45 minutes, depending on network bandwidth.

More Information

[Policy Guidelines and Considerations \[9\]](#)

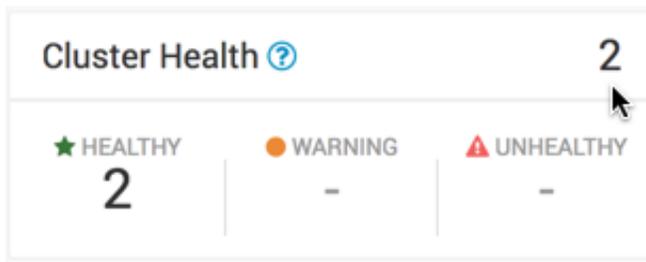
3. Understanding the UI

3.1. Overview Page

From the Overview page you can quickly identify any issues with, or verify the health of, the clusters, policies, or jobs in Data Lifecycle Manager (DLM).

- [Cluster Health Panel \[11\]](#)
- [Policies Panel \[11\]](#)
- [Jobs Panel \[12\]](#)
- [Recent Issues Panel \[13\]](#)
- [Clusters Map \[13\]](#)
- [Issues & Updates Table \[14\]](#)

3.1.1. Cluster Health Panel



You can use the Cluster Health panel of the Overview page to view the total number of clusters enabled for Data Lifecycle Manager, the number that are healthy, the number for which a warning is issued, and the number that are unhealthy .

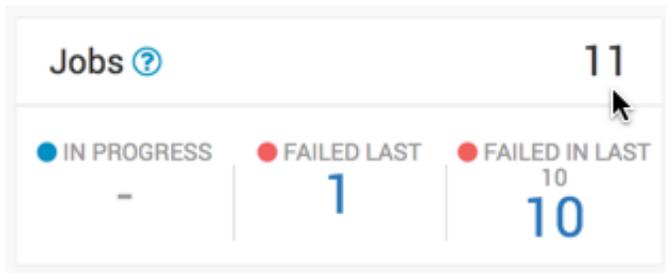
Healthy	Specifies the total number of clusters currently available to run replication jobs. The DLM Engine can be reached and all services are running.
Warning	Specifies the total number of clusters for which remaining disk capacity is less than 10%. If this value is greater than zero, you can click the number to open a table that displays the cluster name and remaining capacity.
Unhealthy	Specifies the total number of clusters for which at least one Apache Ambari service required for DLM (DLM Engine, HDFS, Apache Hive, or Apache Knox) is not started. If this value is greater than zero, you can click the number to open a table that displays the cluster name and the names of any Ambari services that have stopped.

3.1.2. Policies Panel

You can use the Policies panel of the Overview page to view the total number of policies in use and their status.

Active	Specifies policies with status of Submitted or Running. This item is not actionable
Suspended	Specifies policies that have been suspended by an administrator. This item is not actionable.
Unhealthy	Specifies policies associated with any cluster designated as Unhealthy in the Cluster Health panel. If the value is greater than zero, the number becomes clickable. You can click the number to display a table that contains the policy name, the names of the source and destination clusters, and which services are stopped on the source or destination cluster.

3.1.3. Jobs Panel



You can use the Jobs panel of the Overview page to view the total number of running and failed jobs and their status.

In Progress	Specifies the number of jobs with the status Running. If the value is greater than zero, the number becomes clickable. You can click the number to apply a filter to the Issues and Updates table, so that the table displays only in-progress jobs. The filter label Jobs: In Progress appears above the table. Running jobs display as a blue dot in the Policy History column of the table.
Failed Last	Specifies the last job that completed with status Failed. If the value is greater than zero, the number becomes clickable. You can click the number to apply a filter to the Issues and Updates table, so that the table displays only policies for which the last job had a status of Failed. The filter label Jobs: Failed Last appears above the table. Failed jobs display as a red dot in the Policy History column of the table.
Failed in Last 10	Indicates the number of policies for which at least one of the last 10 jobs completed with status Failed. If the value is greater than zero, the number becomes clickable. You can click the number to apply a filter to the Issues and Updates table, so that the table displays only policies for which at least one job failed out of the last 10 jobs. The filter label Jobs: Failed Last appears above the table. Failed jobs display as a red dot in the Policy History column of the table.

3.1.4. Recent Issues Panel

Shows the last four events with severity of warning, critical, or error. For each event, the panel shows the severity, the type, a message that includes the policy name and file icon, and the age of the event.

Severity icon	Displays in orange for warning and red for critical or error.
Event type	Displays in bold text above the event message. The type can be succeeded, deleted, or suspended.
Event message	Displays as text under the event name. When an event is associated with a policy or policy instance (job), then the message text contains two items: <ul style="list-style-type: none">• Policy name: You can click this term to navigate to the Policies page with a preset filter that displays information about only the selected policy.• File icon: You can move the cursor over the icon to display the text "View Log" and click the text to display log content for the associated policy or job.
Event age	Displays in numeric form how long ago from the current time the event occurred.

You can click View All at the bottom of the event list to navigate the browser to the Notifications page.

3.1.5. Clusters Map

Shows red, orange, and green markers that indicate the geolocation of each cluster.

- Red: At least one required service has stopped on the cluster.
- Orange: All required services are running but the remaining capacity on a cluster is less than 10%.
- Green: All required services are running and remaining disk capacity is greater than 10%.

You can move the cursor over a marker on the map displays a tooltip specifying the data center associated with the cluster, the cluster name, and the number of DLM policies that are associated with that cluster.

You can click a marker to open a panel showing the same information as in the tooltip, plus a Launch Ambari link. Clicking the link opens a new browser tab with the login page for the Ambari host for that cluster.

If the dot is red, the panel also displays a list of services that are in a Stopped state in Ambari.

3.1.6. Issues & Updates Table

The Issues & Updates table shows policies that have running jobs but at least one failed out of the most recent 10 jobs. You do not see any policy if its last 10 jobs were all successful. Table columns include the following:

Job Status	When the status of a job is Running, a status circle icon and progress bar display. For jobs that are not running, a status circle icon displays along with the text Success, Failed, or Ignored. You can move the cursor over a Failed status to see a "View Log" tooltip, which you can click to see the job log.
Source & Destination	The names of the source and destination clusters associated with the policy.
Service	Indicates whether the data being replicated is from HDFS or Hive.
Policy	The name assigned to the policy.
Policy History	Shows up to 10 job statuses as colored dots.

Color	Status	Description
Green	Succeeded	Job completed with no issues.
Red	Failed	Job did not complete.
Gray	Ignored	Job did not start because a previous instance of the policy was still in progress. Only one run of a job can be in progress at one time. If a job is consistently ignored, you might need to modify its frequency.

Clicking the colored dots navigates the browser to the Policies page, with the filter preset to show information only about the specified policy.

Transferred/Files	The amount of data transferred, in gigabytes, and the number of objects transferred, if available. When a job is running, the column displays In Progress.
Runtime	How long it took to run the most recent job.
Started	When the most recent job started.
Ended	When the most recent job ended.
Actions icon	<ul style="list-style-type: none"> • Abort Job: Aborts a running job. Enabled only when the job status is Running. • Re-run Job: Starts another instance of the policy. Disabled when a job is executing. • Delete Policy: Removes a policy from Data Lifecycle Manager. Delete cannot be undone. Always enabled.

- Suspend Policy: Suspends the policy and any job that is executing. Disabled when the policy status is Suspended.
- Activate Policy: Resumes a suspended policy. Disabled when the policy status is Running.

4. Roles Required to Work with DPS Services

To perform actions in DPS Platform and associated services, you must be an DPS administrator, an infrastructure administrator, or a data steward. In addition, to perform actions in Apache Ambari that impact DPS (such as creating clusters, changing configuration settings for services, and so forth), you must be an Ambari administrator or a cluster administrator.

Other roles might be required during installation, depending on your configuration. See the installation instructions for roles required during installation.

4.1. Roles for DPS Services

The following roles are available to perform actions in the DPS UI.

4.1.1. DPS Admin

- An DPS Admin role is created during installation, so you can initially log in to DPS Platform.
- Can access DPS Platform and perform all actions in DPS Platform related to clusters, users, and enabling services .
- Can access all services enabled with DPS Platform, and perform the same actions as each administrator role assigned to the enabled services, such as Infra (infrastructure) Admin, Data Steward, and so forth.

4.1.2. Infra Admin

- Can access DPS Platform service to manage clusters enabled for Data Lifecycle Manager (DLM).
- Cannot perform any other DPS Admin actions.
- Can access and perform all actions in Data Lifecycle Manager.
- Has read-only access to browse, within the DLM UI, the folder structure of any cluster enabled for DLM.
 - This access is available to all users logged in to DLM as the Infra Admin, even if they do not have permissions to view the directories or databases as a Linux user or Kerberos principal.
 - Cannot view the *content* of the source or copied files or databases from the DLM UI.
 - Cannot modify or delete folders or databases that you can view from the DLM UI.
- Can create replication policies for any folders or databases on clusters enabled for DLM, and replicate the data objects by using the DLM Engine.

This ability is available to all users logged in to DLM as the Infra Admin, even if the users do not have access to the target path on the target cluster.

4.1.3. Data Steward

- Can access the DPS Platform service to see and manage clusters enabled for Data Steward Studio.
- Can access Data Steward Studio and monitor all data in the DSS UI.

Interactions are only with Apache Atlas and Apache Ranger, so Atlas and Ranger authorization must be set up to enable read access.

- Cannot currently access any file contents from HDFS or Apache Hive.

4.2. Roles Required for Installation and Troubleshooting

You also need the Ambari Admin or Cluster Admin roles to install Hortonworks DPS, add clusters to Ambari, troubleshoot cluster issues, and so forth.

See [Apache Ambari Administration](#) for further details about these roles.

4.2.1. Ambari Admin

- Has full control over all aspects of Ambari.
- Can install HDP by using the Ambari installation wizard.
- Can install the DLM Engine (Beacon) management pack and configure the DLM Engine.
- Can start and stop the DLM Engine and troubleshoot cluster problems.
- *Cannot* access DPS Platform or any enabled service in DPS Platform.

4.2.2. Cluster Admin

- Has control over a cluster, its hosts, and services.
- Can create clusters to be registered with DPS Platform.
- Can start and stop the DLM Engine and troubleshoot cluster problems.
- *Cannot* access DPS Platform or any enabled service in DPS Platform.

4.3. Required Roles by Task

The following tables indicate the roles required to perform various tasks in DPS Platform and the associated services.

4.3.1. DPS Platform Tasks and Required Roles

The following table shows tasks you can perform that are related to DPS Platform and the roles required to perform the tasks.

Task	DataPlane Admin	Infra Admin	Data Steward	Ambari Admin	Cluster Admin
Install HDP using Ambari				X	
Install DPS Docker image	X				
Install service (DLM, DSS) Docker image	X				
Configure LDAP	X				X
Enable DPS services (DLM, DSS)	X				
Manage DPS users	X				
Manage DPS clusters (register, delete, etc.)	X				
Monitor clusters in DPS	X	X	X		
Create a cluster				X	X

4.3.2. Data Lifecycle Manager Tasks and Required Roles

The following table shows tasks you can perform that are related to DLM and the roles required to perform the tasks.

Task	DataPlane Admin	Infra Admin	Data Steward	Ambari Admin	Cluster Admin
Install DLM Engine MPack				X	
Enable DLM	X				
Log in to DLM		X			
Register clusters	X				
Pair clusters		X			
Browse HDFS folders		X			
Browse Hive databases		X			
Create or schedule a replication policy (HDFS or Hive)		X			
Manage a replication policy (suspend, delete, resume, etc.)		X			
Monitor a replication job (HDFS or Hive)		X			

Task	DataPlane Admin	Infra Admin	Data Steward	Ambari Admin	Cluster Admin
Create or schedule a DR policy (HDFS or Hive)		X			
Manage a DR policy (suspend, delete, resume, etc.)		X			
Monitor a DR job (HDFS or Hive)		X			
Monitor DLM alerts		X			
Access DLM log information		X			
Allocate DistCp jobs to YARN queue		X			
Allocate bandwidth		X			
Start or stop the DLM (Beacon) Engine in Ambari				X	X
Access Ambari for troubleshooting				X	X

4.3.3. Data Steward Studio

The following table shows actions you can perform that are related to DSS and the roles required to perform the actions.

Task	DataPlane Admin	Infra Admin	Data Steward	Ambari Admin	Cluster Admin
Enable DSS	X				
Log in to DSS		X			
Register clusters	X				
Start or stop the DSS Profiler engine				X	
Access Ambari for troubleshooting				X	X

More Information

[Apache Ambari Administration](#)

5. Replication Concepts

5.1. Snapshot Replication Between HDP Clusters

You can optionally enable HDFS snapshots for replication in Data Lifecycle Manager. Understanding how snapshots work, and some of the benefits and costs involved, can help you to decide whether or not to enable snapshot replication.

Understanding HDFS Snapshots

HDFS snapshots are read-only point-in-time copies of the filesystem, created from either a subtree of the filesystem or the entire filesystem. You can create a snapshot of any directory after the directory is snapshot enabled (snapshottable).

Blocks in datanodes are not copied during snapshot replication. The snapshot files record the block list and the file size. There is no data copying.

When you enable snapshots on a directory, all subdirectories are automatically enabled for snapshots as well. Snapshots cannot be nested. Snapshot operations are not allowed on a directory if one of its ancestors or descendants already contains snapshots.

When you create a snapshot of a directory, all content in that directory, including subdirectories, is included as part of the copy. There is no limit to the number of snapshot-enabled directories you can have. A snapshot-enabled directory can accommodate 65,536 simultaneous snapshots.

When snapshots are initially created, a directory named *.snapshot* is created on the source and destination clusters, under the directory being copied. All snapshots are retained within *.snapshot* directories. By default, the last three snapshots of a file or directory are retained. Snapshots older than the last three are automatically deleted.

Benefits of snapshots

Snapshot-based replication helps you to avoid unnecessarily copying renamed files and directories. If a large directory is renamed on the source side, a regular DistCp update operation sees the renamed directory as a new one and copies the entire directory.

Generating copy lists during incremental synchronization is more efficient with snapshots than using a regular DistCp update, which can take a long time to scan the whole directory and detect identical files. And because snapshots are read-only point-in-time copies between the source and destination, modification of source files during replication is not an issue, as it can be using other replication methods.

A snapshot cannot be modified. This protects the data against accidental or intentional modification, which is helpful in governance and in meeting disaster recovery (DR) requirements.

Considerations for using snapshots

There is a memory cost to enabling and maintaining snapshots. Tracking the modifications that are made relative to a snapshot increases the memory footprint on the NameNode and can therefore stress NameNode memory.

Because of the additional memory requirements, snapshot replication is recommended for situations in which it is most useful. Such circumstance might include: if you expect to do a lot of directory renaming, if the directory tree is very large, or if you expect changes to be made to source files while replication jobs execute.

Requirements for snapshot-based replication

You must have HDFS superuser privilege to enable or disable snapshot operations.

Replication using snapshots requires that the target filesystem data being replicated is identical to the source data for a given snapshot. *There must not be any modification to the data on the target.* Otherwise, the integrity of the snapshot cannot be guaranteed on the target and replication can fail in various ways.

If a directory contains snapshots but the directory is no longer snapshot-enabled, you must delete the snapshots prior to enabling the snapshot capability on the directory.

6. Configuring Security for DLM

DLM supports Kerberos for securing clusters. If you are working in an unsecured environment, authentication between DLM (Beacon) Engines uses basic authentication. For DLM authorization on secured clusters, there are some prerequisites that must be met. If you want to enable authorization on unsecured clusters, you must complete some manual configuration.

6.1. Configure Authorization for DLM on Secured Clusters

In addition to the security tasks you must complete for DPS, and to satisfy your environmental or corporate requirements, you must ensure the following are properly configured so that DLM replication jobs complete successfully on clusters with Kerberos enabled. No other special configuration is required for authorization and authentication with DLM on a secure cluster.

- HDFS, Hive, Knox, and Ranger are enabled in Ambari
- Ranger plugins are enabled for HDFS and Hive
- Clusters to be paired in DLM have identical configurations, including security
- Global LDAP is configured to share user-group mappings across clusters
- If using Kerberos with different KDCs, two-way trust is configured between the KDCs
- If using AD, there is no support for trust relationships across multiple domains or forests through domain and forest
- HTTPS is supported in DPS Platform, but is *not* supported in DLM in the first release.

6.2. Enable Authorization for DLM on Unsecured Clusters

By default, authorization for DLM is not enabled in an unsecured cluster. You can enable authorization by modifying a security file available with the DLM Engine on the HDP clusters. You can also enable authorization for individual users or groups.

Steps

1. On a node with the DLM Engine installed, navigate to `$DLM_HOME/conf`.
2. Open the file `beacon-security.xml`.
3. Add the following properties to the file:

```
<property>  
<name>beacon.authorization.enabled</name>
```

```
<value>>true</value>
</property>
<property>
<name>beacon.authorization.policy.file</name>
<value>policy-store.txt</value>
</property>
```

4. Enable authorization for individual users or groups by doing the following:

- a. In `$DLM_HOME/conf`, open the `policy-store.txt` file.

The file looks similar to the following:

```
##Policy Format
##r-READ, w-WRITE, u-UPDATE, d-DELETE
##Policy_Name; ;User_Name1:Operations_Allowed,
User_Name2:Operations_Allowed; ;Group_Name1:Operations_Allowed,
Group_Name2:
Operations_Allowed; ;Resource_Type1:Resource_Name,
Resource_Type2:Resource_Name
##
adminPolicy; ;admin:rwud; ;ROLE_ADMIN:rwud; ;cluster:*,policy:*,schedule:*,
event:*,logs:*
beaconPolicy; ;beacon:rwud; ;ROLE_ADMIN:rwud; ;cluster:*,policy:*,
schedule:*,event:*,logs:*
```

- b. Add the configuration settings for the users you want to authorize for access to DLM, following the `adminPolicy` and `beaconPolicy` examples in the file.

```
adminPolicy; ;admin:rwud; ;ROLE_ADMIN:rwud; ;cluster:*,policy:*,schedule:*,
event:*,logs:*
beaconPolicy; ;beacon:rwud; ;ROLE_ADMIN:rwud; ;cluster:*,policy:*,
schedule:*,event:*,logs:*
```

7. Failing Over Manually

If a source cluster used in a replication policy is offline and will not be brought online for an extended period, you should manually fail over the destination cluster to serve as the new source. After failover, the new source cluster will receive read and write requests. You might also want to designate a new destination cluster to which data will be copied from the new source.

7.1. Make the Destination Cluster the New Source

If the source cluster becomes unavailable for an extended period, you can configure the destination cluster to serve as the new source. Read and write requests from clients will then be redirected from the old source to the new source cluster.

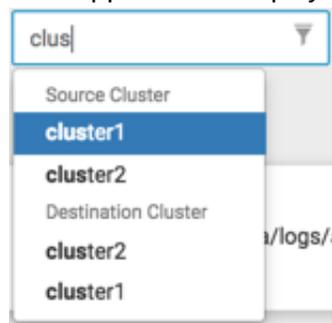
Prerequisites

You need the name of the cluster that is offline.

Steps

1. Log in to the DPS UI as Infra Admin.
2. Access the DLM UI by clicking the DPS icon in the upper left of the page and then clicking the Data Lifecycle Manager icon.
3. Identify the set of replication policies for which the offline cluster is the source in a replication relationship.
 - a. Click **Policies** in the navigation pane.
 - b. In the **Filter** field, type the name of the offline cluster.

A list appears that displays the cluster name as a source or a destination cluster.



- c. From the list, select the cluster name under Source Cluster.

The page content shows only the policies that use the selected cluster as the source for replication.

4. Delete all replication policies that use the offline cluster as the replication source.
 - a. At the end of each row in the policies list, click the **⋮** (Actions) icon.

- b. Click **Delete** in the drop-down menu, and then click **OK** to confirm deletion.

If a replication policy is in the process of running a job, the job aborts when you delete the policy.



Important

After a replication policy is deleted, it cannot be retrieved.

5. If the Ranger deny policy is enabled, remove the deny policy that is on the destination cluster.
 - a. Determine if the Ranger deny policy is enabled.
 - a. Navigate to the Ambari UI.
 - b. In the services list, click **DLM Engine**.
 - c. Click **Configs>Advanced**.
 - d. Scroll to the parameter `beacon.ranger.plugin.create.denypolicy` and verify if the Ranger Deny Policy is enabled or disabled.
 - b. If the Ranger Deny Policy is enabled, you must disable it.
 - a. Log in to the destination cluster, access Ranger, and then navigate to Ranger admin resource policies.
 - b. Identify Ranger policies that start with "`<sourcecluster>_beacon deny policy for`".

7.2. Activate a New Destination Cluster

If you have not prepared a cluster in advance to serve as an alternate destination in a failover scenario, then you must install the DLM Engine, configure the clusters for use by DLM, and pair the clusters before you can create new replication policies and begin copying data to the new destination.

Prerequisites

You must have the name of the cluster you want to configure as the new destination.

Steps

1. Identify the Ambari-managed cluster to use as the new destination.
2. Install the DLM Engine on the new destination, if it is not already installed.

[Installing DPS Services, Engines, and Agents](#)

3. Follow the instructions in [Setting Up the DPS Services](#) for the following tasks, as needed:
 - *Register Clusters with DPS*

- *Enable Services*

4. Pair the clusters you are using as source and destination, if they are not already paired.

[Pair Clusters for Replication](#)

5. Ensure that the HDFS folders or Hive databases to be copied either do not exist or are empty on the new destination cluster.

This is required prior to bootstrapping data from the source cluster to the destination cluster. Otherwise, the initial copy job fails.

6. Create and submit new replication policies between the source and destination clusters.

[Create a Replication Policy](#)

The first time a new policy is submitted, the entire contents of the source dataset is copied to the destination. Depending on the size of each dataset, these initial bootstrap copies can take a significant amount of time. After the initial copy, subsequent copies are incremental.

More Information

[Hortonworks DPS Installation](#)

8. Troubleshooting DLM

To verify that your environment meets the requirements for DPS, see [Planning for a DPS Installation](#) and the [Hortonworks Support Matrix](#).

8.1. Endpoint is unreachable

DLM does not support updating any cluster endpoints (HDFS, Hive, Ranger, or DLM Engine). If an endpoint must be modified, contact Hortonworks Support for assistance.

8.2. Ranger UI does not display deny policy items

If you need to view deny policy details related to a DLM replication policy, you need to use the Ranger UI. However, when a policy with deny conditions is created on Ranger-admin in a replication relationship, the Policy Details page in Ranger does not display the deny policy items. To make the policy visible, update the respective `service-def` with `enableDenyAndExceptionsInPolicies="true"` option.

Refer to section "2.2 Enhanced Policy model" in <https://cwiki.apache.org/confluence/display/RANGER/Deny-conditions+and+excludes+in+Ranger+policies>