

Administration 1

DLM Administration

Date of Publish: 2018-07-03

<http://docs.hortonworks.com>

Contents

Replication Concepts.....	4
HDFS cloud replication.....	4
Hive cloud replication.....	4
Non-support of replication of Hive-Managed tables written by Spark applications.....	5
Cloud replication guidelines and considerations.....	6
Pairing guidelines and considerations.....	6
Policy guidelines and considerations.....	7
Snapshot guidelines and considerations.....	8
Retention of Ranger policies, HDFS permissions, and ACLs.....	8
Hive replication bootstrap.....	8
Replicating Data.....	9
HDFS replication policy process overview.....	9
Replicating data on-premise to on-premise.....	9
Replicating data on-premise to cloud.....	9
Register cloud credentials.....	10
Pair clusters for replication.....	10
Create a replication policy.....	11
Edit an existing policy.....	13
View job status from the Policies page.....	14
View job status from the Overview page.....	15
View job status from the Notifications page.....	15
Update Cluster Endpoint.....	15
Failing Over Manually.....	16
Make the destination cluster the new source.....	16
Remove the Ranger deny policy.....	17
Activate a new destination cluster.....	18
DLM policy parameters.....	18
DLM version Information.....	19
Roles Required to Work with DLM.....	19
Infra Admin role.....	20
DataPlane Admin role.....	20
Roles Required for Installation and Troubleshooting.....	20
Data Lifecycle Manager Tasks and Required Roles.....	21
Tuning DLM Engine.....	21
Troubleshooting DLM.....	22

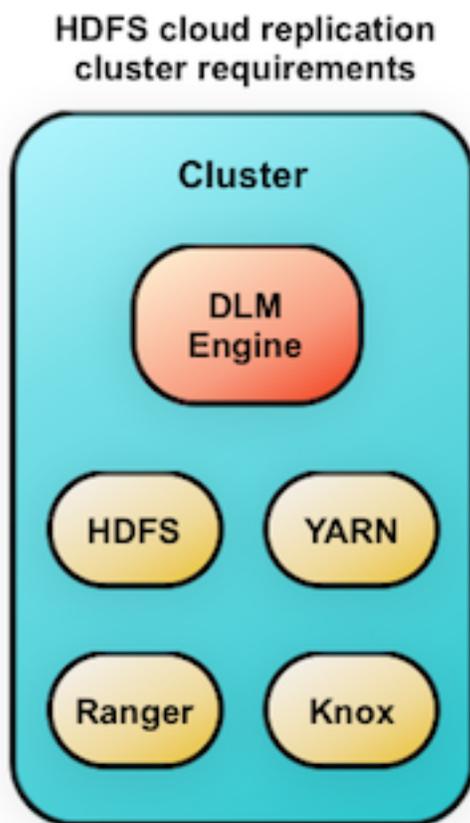
Ranger UI does not display deny policy items.....	22
Replication fails with TDE and non-TDE data.....	22
Hive data cannot be replicated.....	22
Instance of a policy stuck in a running state.....	23
Hive replication failure.....	23

Replication Concepts

HDFS cloud replication

DLM supports replication of HDFS data from cluster to cloud storage and vice versa. The replication policy runs on the cluster and either pushes or pulls the data from cloud storage.

The cluster can be an on-premise or IaaS cluster with data on local HDFS. The cluster requires HDFS, YARN, Ranger, Knox and Beacon services.



Hive cloud replication

DLM supports replication of the Hive database from a cluster with underlying HDFS to another cluster with cloud storage. It uses push-based replication, with the replication job running on the cluster with HDFS. Hive replication from cloud storage to HDFS is not supported.

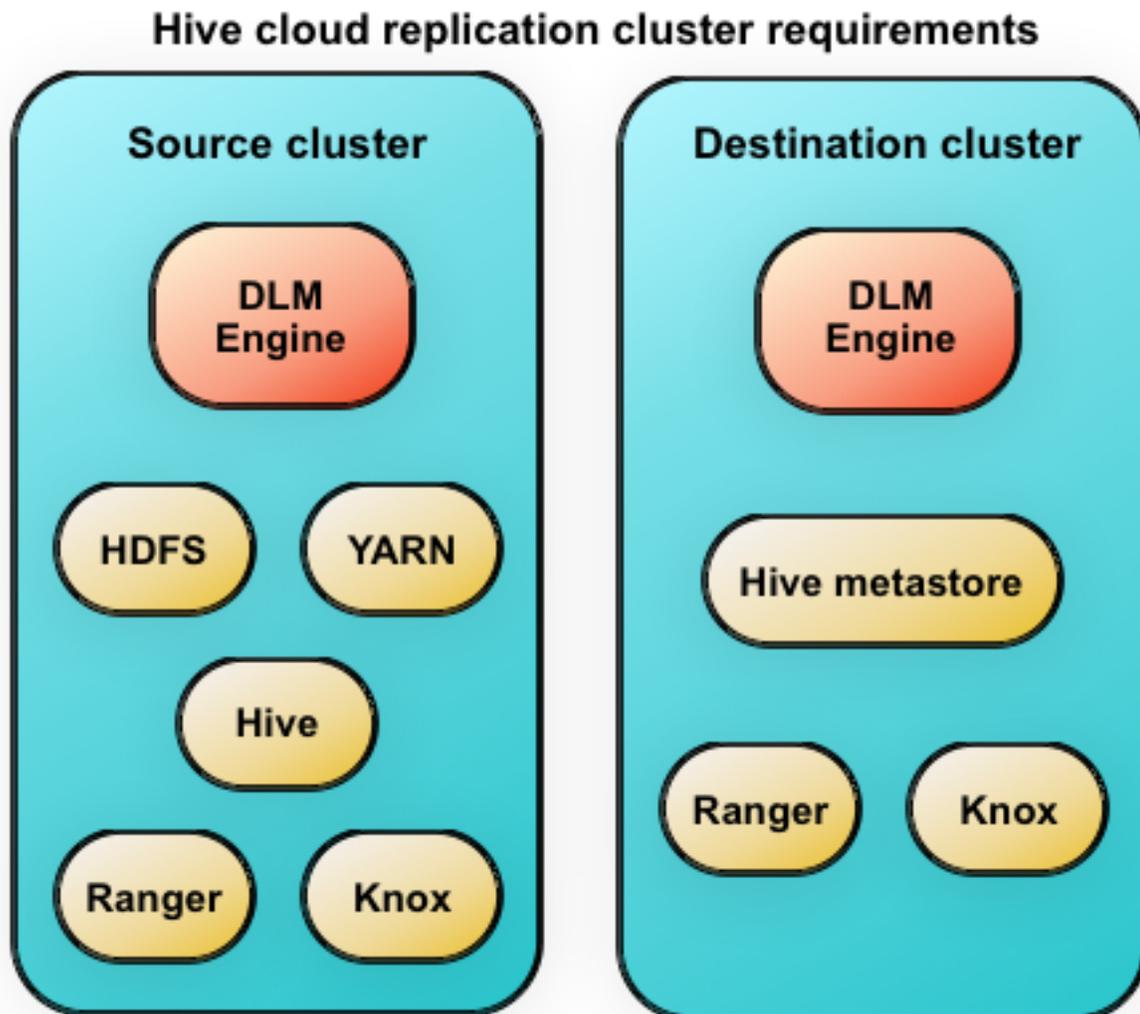
Hive stores its metadata in Hive Metastore, but the underlying data is stored in HDFS or cloud storage. In a Hadoop cluster with Hive service, the Hive warehouse directory can be configured with either HDFS or cloud storage.

- You can rename the dataset in the policy that is replicated.
- You can create a pull-based policy on the source cluster to move data from the target back into the source cluster Hive database.

- DLM does not manage Ranger policies and any PII/secure data that gets replicated from on-premise to S3. You must manage these items outside of DLM.
- Hive replication from an HDFS-based cluster to a cloud storage-based cluster requires the following:
 - Source cluster

The cluster with a Hive warehouse directory on local HDFS. This can be an on-premise cluster or an IaaS cluster with data on local HDFS. The required services are HDFS, YARN, Hive, Ranger, Knox and Beacon.
 - Destination cluster

The cluster with data on cloud storage. The cluster minimally requires Hive Metastore, Ranger, Knox and Beacon Services.



Non-support of replication of Hive-Managed tables written by Spark applications.

DLM Hive replication for Managed tables relies on replication events being published by Hive in Hive Metastore for every change that is made by Hive.

In case of External table replication, DLM replication does not rely on events being published and checks every table/partition directory for any new file that might have been added.

Important: Applications other than Hive do not always publish events for new data file addition to Managed tables. The list of such applications includes Spark. This can result in data loss if these applications write to a Managed table in HDP 2.6.5. External tables should be used for data written by such applications. While replication for External table has some overheads, it will capture files that have been added without any event generation as well.

Note: With Spark, the use of `hive.metastore.dml.events` is not supported in HDP. Spark should be treated as an application that does not reliably publish events for the changes.

Cloud replication guidelines and considerations

DLM supports replication of HDFS and Hive data between underlying HDFS and AWS S3 cloud storage.

Accessing AWS Services for Amazon S3

- Cloud bucket requirements
 - You need a cloud bucket with user credentials that you can enter in DLM, so DLM can access the bucket.
 - The bucket has to have enough space for the replicated data, and write permissions to copy the data.
 - The bucket needs to support cloud storage encryption types supported by DLM (SSE-S3 & SSE-KMS).
- Authentication
 - DLM supports access key and secret key authentication with AWS S3.
 - Unregistered credentials in DLM are credentials associated with a cluster node that does not have updated credentials.

An example of how this can arise is if a node was down when the credentials were changed on a bucket, and when the node is brought up it still has the old credentials.

- Impact of bucket changes
 - Changes made to a bucket configuration (secret/access keys, bucket name/endpoint, encryption type) can affect execution of the DLM policy and might require an update to DLM cloud credentials.

Credential changes are picked up by the next run of the policy. Any policies being run when the credential changes are made could fail, but succeeding runs will pick up the changes.

- Users can delete cloud credentials, but this triggers failures of any policies based on the deleted cloud credentials.

You must delete the DLM cloud policies associated with the deleted credentials and recreate the policies with the new credentials. You can view a list of policies associated with specific credentials on the Cloud Credentials page.

- Cloud encryption
 - When replicating data from cloud storage, the encryption algorithm specified by the user is used for validations on the replication policy.
 - When replicating data to cloud storage, the encryption algorithm and encryption key specified by the user are used for all the data written to the cloud storage.

This overrides any bucket level encryption set in the cloud provider.

- DLM does not allow replication of encrypted data to an unencrypted destination.

Pairing guidelines and considerations

You should take into consideration the following items when pairing clusters in DLM.

- 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.
- The HDFS nameservice for the source and destination clusters cannot be configured with the same name.
- Cluster security configurations must be symmetrical to pair clusters, including Kerberos, Ranger, Knox, etc.
- The DLM Engine must be configured as Ranger administrator on all DLM-enabled clusters.

Policy guidelines and considerations

You should take into consideration the following items when creating or modifying a replication policy.

Data security

- If using TDE for encryption, the entire source directory must be either encrypted or not encrypted, otherwise policy creation fails.
- If using an S3 cluster for your policy, your credentials must have been registered on the Cloud Credentials page.
- On destination clusters, the DLM Engine must have been granted write permissions on folders being replicated.
- 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 administrators have the ability to modify or delete folders or files that are viewable from the DLM UI.

Policy properties and settings

- 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 Skipped. If a job is consistently skipped, you might need to modify the frequency of the job.
- 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.

Cluster requirements

- The target folder or database on the destination cluster must either be empty or not exist prior to starting a new policy instance.
- The clusters you want to include in the replication policy must have been paired already.
- 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).

Hive restrictions

- ACID tables, external tables, storage handler-based tables (such as HBase), and column statistics are currently not replicated.
- 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.

Time considerations

- The first time you execute a job (an instance of a policy) with data that has not been previously replicated, Data Lifecycle Manager creates a new folder or database and bootstraps the data.

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. So you should plan the bootstrap accordingly.

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.

- 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.

Related Tasks

[Create a replication policy](#)

Snapshot guidelines and considerations

Before enabling snapshots on data to be replicated, consider the following items.

1. If the logged in Infra Admin user is an HDFS superuser, the source directory can be marked as snapshot enabled (snapshottable).
2. For every replication, DLM Engine creates a new snapshot on the source. DistCp then compares this snapshot with the previous snapshot to determine the list of files to be copied. Note that the actual file changes are copied, and not the snapshots themselves.
3. DLM Engine also creates snapshots on the destination HDFS after every replication. These snapshots are used to recover the destination HDFS state to a consistent state in case of failure.
4. DLM Engine also handles retention of snapshots on both source and target, configurable through Ambari.

Related Concepts

[Snapshot replication between HDP clusters](#)

Retention of Ranger policies, HDFS permissions, and ACLs

When a DLM replication job is run, data, metadata, and any Ranger policies that are associated with the replicated data are automatically exported to the target. The data on the destination is marked as read-only by adding a deny policy on the replicated data in Ranger in the destination cluster. This prevents accidental writes on the copy.

For on-premise to on-premise replications, the policies, permissions, and ACLs are retained and applied to the data on the target, except that the destination data is read-only.

For on-premise to cloud replication, the Ranger policies, permissions, and ACLs are stored in metadata files in cloud storage. Data in the cloud is protected using security features in the cloud environment.

Hive replication bootstrap

DLM allows you to replicate Hive databases from a source cluster to a target location on a destination cluster.

When you initiate the replication of Hive data, all of the data from the source location is copied to the destination. This bootstrapping of data can take hours to days, depending on factors such as the amount of data being copied and available network bandwidth. Subsequent replication jobs from the same source location to the same target on the destination are incremental, so only the changed data is copied.

If a bootstrap replication is interrupted, such as due to a network failure or an unrecoverable error, DLM automatically retries the job. If a retry succeeds, the replication job continues from the point at which it was interrupted. If the automatic retries are not successful, you must manually correct the problem before running the policy again. When you activate the policy again, the replication job resumes from the point at which it was suspended.

After the bootstrap replication succeeds, an incremental replication is automatically performed. This job synchronizes, between the source and destination clusters, any events that occurred during the bootstrap process. After the data is synchronized, the replicated data is ready for use on the destination.

External tables, ACID tables, storage handler-based tables (such as HBase), and column statistics are not currently replicated.

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.

HDFS replication policy process overview

1. The DLM App submits the replication policy to the DLM Engine on the destination cluster. The DLM Engine then schedules replication jobs at the specified frequency.
2. At the specific frequency, DLM Engine submits a DistCp job that runs on destination YARN, reads data from source HDFS, and writes to destination HDFS.
3. File length and checksums are used to determine changed files and validate that the data is copied correctly.
4. The Ranger policies for the HDFS directory are exported from source Ranger service and replicated to destination Ranger service.

DLM Engine also adds a deny policy on the destination Ranger service for the target directory so that the target is not writable.

Replicating data on-premise to on-premise

Before you can begin replicating data, you must pair your clusters and create a policy that specifies the data to replicate, the replication schedule, and other setting.

About this task

See the individual tasks linked below for considerations and tips when performing the tasks.

Before you begin

You must have the Infra Admin role to perform this set of tasks.

You must have registered clusters with DPS Platform prior to pairing them.

Procedure

1. [Pair clusters for replication](#)
Select the two clusters to use for replication and pair them so they can communicate.
2. [Create a replication policy](#)
Choose which cluster is source and which is destination, then set the schedule and other rules for replication jobs.
3. [View job status from the policies page](#).
Verify that the job starts and runs as expected.

Replicating data on-premise to cloud

The process for creating a replication job from on-premise to the cloud is similar to creating one for on-premise to on-premise. The primary difference is that you must register your cloud credentials with DLM, so DLM can access your cloud storage.

About this task

Attention: Replication of HDFS data from on-premise to cloud is a *Limited GA* feature in DPS 1.1. The HDFS data that you replicate to cloud requires security policies outside the Hadoop system, so you should work with Hortonworks support to ensure proper configuration of your environment. This does not apply to Hive replication to cloud.

See the individual tasks linked below for considerations and tips when performing the tasks.

Before you begin

You must have the Infra Admin role to perform this set of tasks.

Procedure

1. [Register cloud credentials with DLM.](#)

Enter the credentials for the bucket you want to replicate, so DLM can access the bucket.

2. [Create a replication policy.](#)

Choose which cluster is source and which is destination, then set the schedule and other rules for replication jobs.

3. [View job status.](#)

Verify that the job starts and runs as expected.

Register cloud credentials

If you plan to replicate data to a cloud account, you must register the cloud credentials so DLM can access your cloud account.

Procedure

1. In the navigation pane, click **Cloud Credentials**, and then click **Add**.

2. Enter or select the following information:

- Cloud Storage Type
- A unique name for the credential
- Authentication Type
 - Access & Secret Key or IAM Role
- If using Access & Secret Key, enter the keys

3. Submit the settings.

- If using Access & Secret Key, click **Validate**.

Using the validation feature is recommended to ensure that the S3 bucket keys are valid. If the keys are not valid, the DLM policy cannot execute a copy of data to the target S3 bucket.

- If using IAM Role, click **Save**.

4. Verify that your credentials display correctly on the Cloud Credentials page.

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.

About this task

Knox proxying can only be enabled on paired clusters running DLM Engine version 1.1 or higher.

Pairing clusters only establishes a communication relationship between two clusters. You determine which cluster of that pair is the source and which is the destination when you create the replication policy.

Before you begin

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

Procedure

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 the above steps to pair additional clusters.

Related Concepts

[How pairing works in Data Lifecycle Manager](#)

[Policy guidelines and considerations](#)

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.

About this task

- 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, DLM copies all of the data. The bootstrap process can take hours to days, depending on data size, so plan your time accordingly.

Before you begin

- You must use the DLM Infrastructure Admin role to perform this task.
- The target folder or database on the destination cluster must either be empty or not exist prior to starting a new policy instance.

Procedure

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

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

2. Click **Add > Policy**.
3. On the **General** page, enter or select the following information, and then click Select Source:
 - Policy Name
 - Description
 - Service: HDFS or Hive
4. On the **Select Source** page, enter or select the following information, and then click Select Destination:
 - Type: S3 or Cluster
 - Source Cluster (if Type=Cluster is selected)
 - Cloud Credential (if Type=S3 is selected)

You must have registered your credentials with DLM on the Cloud Credentials page.
 - Select a Folder Path (only if HDFS is selected)

TDE-enabled directories are identified by a lock icon. The entire source directory must be either encrypted or not encrypted, otherwise policy creation fails.
 - Enable snapshot based replication (only if HDFS is selected)

HDFS Admin role is required to enable snapshots.
 - Select Database (Only if Hive is selected)

TDE-enabled databases are identified by a lock icon.
5. On the **Select Destination** page, enter or select the following information, and then click Schedule:
 - Type: S3 or Cluster
 - Destination Cluster (if Type=Cluster is selected)
 - TDE Same Key (if Type=Cluster is selected)

Configures the policy to use the same TDE key for the source and destination.
 - Cloud Credential (if Type=S3 is selected)
6. On the **Schedule** page, select when you want the job to run, and then click Advanced Settings:

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

 - Start: On Schedule or From Now
 - Repeat
 - Start and End Dates
 - Start Time
7. Enter or select the **Advanced Settings**, and then click Create Policy:

Configuring Advanced Settings is optional.

 - Queue Name

If you are using Capacity Scheduler queues to limit resource consumption, enter the name of the YARN queue for the cluster to which the replication job will be submitted.
 - Maximum Bandwidth

You can adjust this setting so that each map task is throttled to consume only the specified bandwidth so that the net bandwidth used tends towards the specified value. The default value for the bandwidth is 1 MB per second.
 - Maximum Maps

Use this option to set the maximum number of map tasks (simultaneous copies) per replication job.

The Advanced Settings attributes are applied only during DLM replication jobs that are based on DistCp functionality.
8. Click Review and verify that the settings are correct.

After a policy is created, the policy name and the clusters associated with the policy cannot be modified.

- Click Submit. A message appears, stating that the submission was successful.
When the policy job runs, checks are performed to verify the copied data.

What to do next

View job status to verify that the replication job is running as intended.

Related Concepts

[Policy guidelines and considerations](#)

Related Information

[Using DistCp to Copy Files](#)

Edit an existing policy

You can edit some settings in your policies to better align with changing requirements. For example, you might want to change the frequency of a policy depending on the data size and importance of the data being replicated.

About this task

The Edit Replication Policy page is not available prior to DLM version 1.1.1.

- You can edit an existing policy, with the following restrictions:
 - Only non-expired policies in active or suspended state can be edited.
 - The start time cannot be modified if the policy has already started.
 - You cannot modify the policy name or the source or destination cluster.
- 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.

Before you begin

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

Procedure

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

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

- Locate the policy you want to edit and click



(Actions).

Status	Name	Source	Destination	Job	Duration	Last Good	
ACTIVE	contacts-data Every 20m	c1 /rest/contacts	c2 /rest/contacts	●●●	<1m	18m ago	⋮

← Actions

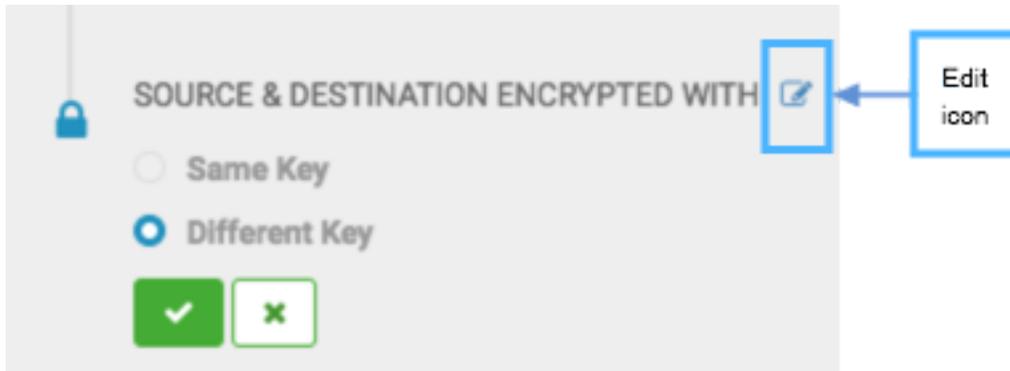
- Select **Edit** and then modify and save the policy.

The following options are available to edit:

- Frequency
- Start Date (if the policy has not yet run an initial job instance)
- End Date
- Start Time
- Queue Name
- Maximum Bandwidth
- Maximum Maps

If the Edit option does not display, verify that the policy status is active or suspended. Expired policies cannot be edited.

- To edit policy description and key selection, on the Policies Page, click the policy name. The Policy Settings display.
- Key selection is only available for policies that are replicating TDE-enabled data.
- Click the **Edit** icon next to Description or Source & Destination Encrypted With.



Clicking the Edit icon next to other items in Policy Settings opens the Edit Replication Policy wizard.

- Click the checkmark to save the change and close the edit option.

What to do next

View job status to verify that the replication job is running as intended.

View job status from the Policies page

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

Before you begin

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

About this task

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.

Procedure

- In the navigation pane, click **Policies**.
- Click the



or



icon to display the type of policies you want to see.

- 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.

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.

Procedure

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) sClick 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.

View job status from the Notifications page

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

Before you begin

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

Procedure

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.

Update Cluster Endpoint

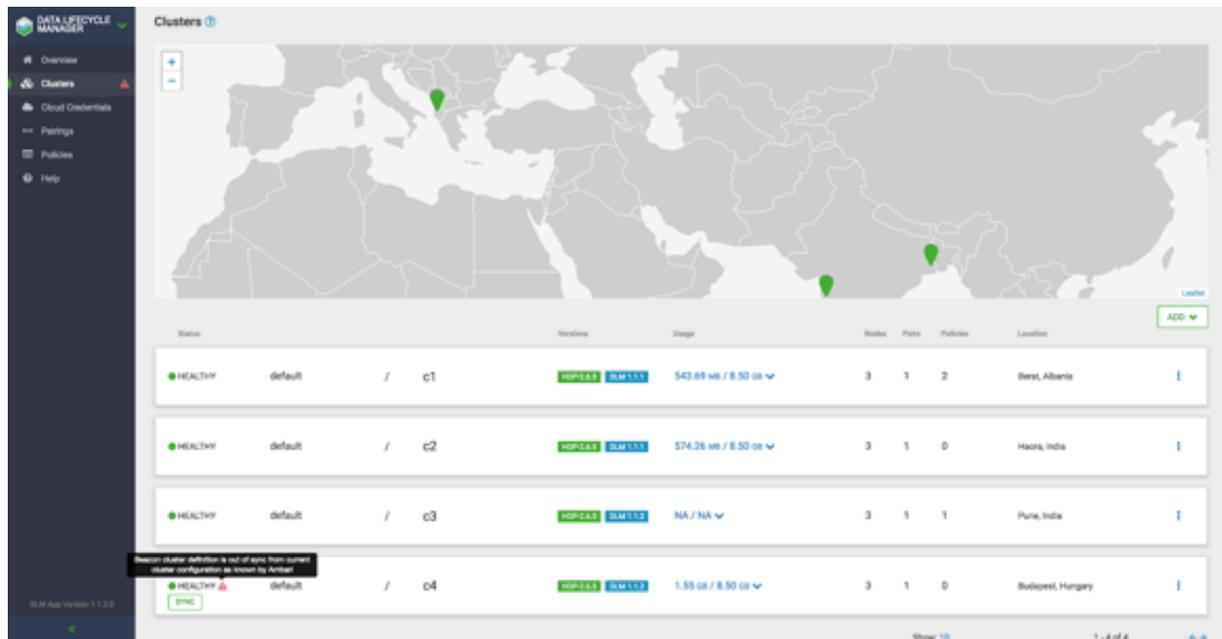
A DLM endpoint server is present for each cluster on the DataPlane Services that has DLM Engine installed. As an administrator, you can change the specific configurations on Ambari to update any cluster endpoint and ensures that it works with DLM.

Before you begin

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

Procedure

1. Log in to DataPlane services.
2. On the navigation pane, click **DATA LIFECYCLE MANAGER**.
3. Click **Clusters**.
You can view the DLM Engine clusters on the **Clusters** page.
4. Click the **Sync** button to synchronize the changes between the Ambari cluster and the DLM Engine.



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.

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.

Before you begin

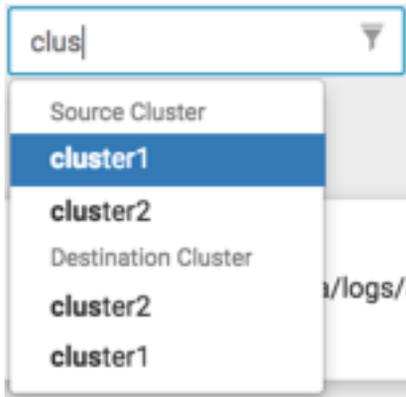
You must be logged in as Infra Admin to perform this task.

You need the name of the cluster that is offline.

Procedure

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.

What to do next

If the Ranger deny policy is enabled, remove the deny policy that is on the destination cluster.

Remove the Ranger deny policy

If the Ranger deny policy is enabled, you must remove the deny policy that is on the destination cluster so that DLM can access the target data to be retrieved.

Before you begin

You must be logged in as Ambari Admin to perform this task.

Procedure

1. 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.
2. 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” and remove the deny condition on the policies.

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.

Before you begin

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

Procedure

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.

DLM policy parameters

Add Policy Parameters

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.

Field	Description	Additional Information
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.
Enable snapshot based replication	Enables snapshot replication on the selected folder if you have the required permissions	When the job runs, snapshots are automatically created on the destination cluster and managed by DLM. HDFS Admin role is required to enable snapshots.
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	
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).
Maximum Maps	Sets the maximum number of map tasks (simultaneous copies) per replication job.	

DLM version Information

As a DLM Administrator, you can view various version-related details on the DLM user interface.

You can view the following details on the DLM UI:

- DLM Engine version
- HDP version on each cluster

Note: You can view the DLM Engine and HDP for each cluster on the following pages: Clusters, List Pairings, and Create Pairings.

Roles Required to Work with DLM

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.

Infra Admin role

The infrastructure administrator has access to DLM and administrative permissions to perform all actions in DLM.

- 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.

DataPlane Admin role

The DataPlane Admin has access to DPS Platform and administrative permissions to perform all actions in DPS Platform. A DataPlane Admin role is created during installation, so you can initially log in to DPS Platform.

The DPS Admin has the following capabilities and restrictions:

- 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 Admin, Data Steward, and so forth.

In addition to this role, there are app-specific admin roles that are needed to manage the enabled services. For more information, see the app-specific documentation.

Roles Required for Installation and Troubleshooting

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

See the Apache Ambari documentation for further details about these roles.

Ambari Admin

You must have the Ambari Admin role to perform many of the cluster and service troubleshooting actions in Ambari. The Ambari Admin cannot access DPS Platform or any enabled service in DPS Platform.

You can perform the following actions with the Ambari Admin role:

- Manage all aspects of Ambari
- Install HDP by using the Ambari installation wizard
- Install the DLM Engine (Beacon) management pack and configure the DLM Engine
- Start and stop the DLM Engine and troubleshoot cluster problems

Cluster Admin

You must have the Cluster Admin role to perform many of the cluster and service troubleshooting actions in Ambari. The Cluster Admin cannot access DPS Platform or any enabled service in DPS Platform.

You can perform the following actions with the Cluster Admin role:

- Manage a cluster, its hosts, and services
- Create clusters to be registered with DPS Platform
- Start and stop the DLM Engine and troubleshoot cluster problems

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	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		
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

Tuning DLM Engine

You can tune the DLM Engine for tasks such as running multiple concurrent policies and handling multiple files.

Run Multiple Concurrent Policies

Perform the following steps to run multiple concurrent policies in DLM:

1. Log in to Ambari.
2. Set the `beacon_quartz_thread_pool` property to a value greater than the number of policies required to run concurrently.

Handle Multiple Files

For the DLM Engine to handle multiple files that are listed, ensure that it has sufficient memory.

Troubleshooting DLM

To verify that your environment meets the requirements for DPS, see the DPS Support Matrices.

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>.

Replication fails with TDE and non-TDE data

HDFS Replication fails when some files are encrypted and some are unencrypted. If the source directory is unencrypted, but contains both encrypted and unencrypted subfolders, then replication jobs fail with checksum mismatch error.

Ensure that all folders in a source *root* directory have the same encryption setting (enabled/not enabled or same key).

Hive data cannot be replicated

If an initial Hive replication (bootstrap) fails in DLM, review the following possible causes and resolutions to try resolving the issue.

Notification events are missing in the meta store

REPL_EVENTS_MISSING_IN_METASTORE (20016)

Use the drop command to delete the target database and then resume the policy from the DLM App UI.

Target database is bootstrapped from some other path.

REPL_BOOTSTRAP_LOAD_PATH_NOT_VALID (20017)

Use the drop command to delete the target database and then resume the policy from the DLM App UI.

File is missing from both the source and CM path.

REPL_FILE_MISSING_FROM_SRC_AND_CM_PATH (20018)

Review the DLM Engine logs to locate the REPL DUMP directory, remove the directory, delete (drop) the target database, and then resume the policy from the DLM App UI.

Either the dump directory does not exist or it is not accessible

REPL_LOAD_PATH_NOT_FOUND (20019)

If the dump location does not exist, you can resume the policy and the DLM Engine creates a new dump.

If the directory is not accessible, you need to set the required permissions.

The source for the replication (repl.source.for) is not set in the database properties.

REPL_DATABASE_IS_NOT_SOURCE_OF_REPLICATION (20020)

On the source database, use DESC DATABASE EXTENDED <db_name> to determine if the parameter repl.source.for is set with the policy name.

If the policy is scheduled and the above parameter is not set, then set the parameter using ALTER DATABASE <db_name> SET DBPROPERTIES 'repl.source.for'='<policy_name>'.

Then resume the policy from the DLM App UI.

Instance of a policy stuck in a running state

If your policy is stuck in a running state because of some unknown exceptions, you must restart the DLM engine using Ambari. This process would in turn handle the failure scenarios.

Note: If a database failure is detected, you must first get the database service up and running.

Hive replication failure

Hive replication fails with an error message 'This operation is not allowed on source cluster: <ClusterOne>. Try it on target cluster: <ClusterTwo>'

If the Hive warehouse directory on target cluster is changed from HDFS to Cloud storage, you must 'Sync' the cluster in DLM UI. DLM UI must be aware about the cluster changes.