

Administration 1

DLM Administration

Date of Publish: 2018-10-15



<http://docs.hortonworks.com>

Contents

Introduction.....	5
Purpose and scope.....	5
Audience and assumptions.....	5
Replication concepts.....	5
Data Lifecycle Manager terminology.....	5
Communication with HDP clusters.....	6
How Policies Work in Data Lifecycle Manager.....	7
UI overview.....	8
Cluster Health panel.....	8
Policies panel.....	9
Jobs panel.....	9
Recent Issues panel.....	10
Clusters map.....	10
Issues & Updates table.....	10
Preparing to setup replication policy.....	11
Roles required.....	11
Infrastructure Admin role.....	12
DataPlane Admin role.....	12
Roles Required for Installation and Troubleshooting.....	12
Data Lifecycle Manager Tasks and Required Roles.....	13
Add clusters.....	13
Cluster pairing.....	14
Pairing considerations.....	14
Cloud credentials.....	14
Register cloud credentials.....	14
Registering Amazon S3 cloud account.....	14
Considerations for Amazon S3.....	15
Registering WASB cloud account.....	15
Considerations for WASB.....	15
Data replication use cases.....	15
Replication of data using HDFS.....	16
HDFS cloud replication.....	16
On-premise to on-premise replication in HDFS.....	17
On-premise to Amazon S3 replication in HDFS.....	17
Amazon S3 to on-premise replication in HDFS.....	18
On-premise to WASB replication in HDFS.....	19
WASB to on-premise replication in HDFS.....	19
Replication of data using Hive.....	20
Hive cloud replication.....	20

Hive replication bootstrap.....	21
Non-support of replication of Hive-Managed tables written by Spark applications.....	22
On-premise to on-premise replication in Hive.....	22
On-premise to Amazon S3 replication in Hive.....	23
On-premise to WASB replication in Hive.....	24
Metadata replication.....	25
Ranger metadata.....	25
Atlas metadata.....	25
Snapshot replication between HDP clusters.....	26
Replication policy operations.....	27
Monitoring replication.....	27
Policies page.....	27
Overview page.....	28
Notifications page.....	28
Tuning replication policy (advanced options).....	29
Update replication policy.....	29
Browsing data directory.....	30
Cloud credentials operations.....	30
Update cloud credentials.....	30
Delete credentials.....	31
Unregistered credentials.....	31
Miscellaneous.....	31
Update Cluster Endpoint.....	31
Failing Over Manually.....	32
Make the destination cluster the new source.....	32
Remove the Ranger deny policy.....	33
Activate a new destination cluster.....	34
DLM policy parameters.....	34
DLM version Information.....	35
Tuning DLM Engine.....	35
Troubleshooting DLM.....	36
Ranger UI does not display deny policy items.....	36
Hive cloud replication is slow.....	36
Replication fails with TDE and non-TDE data.....	36
Hive data cannot be replicated.....	36
Instance of a policy stuck in a running state.....	37
Hive replication failure.....	37

Introduction

This *System Administrator Guide* provides information on how to setup, configure, and manage Data Lifecycle Manager (DLM) jobs to replicate data for small and large scale enterprise organisations.

Purpose and scope

This guide is intended to assist System Administrators to manage the complete data replication life cycle, using on-premise and cloud environments.

Audience and assumptions

The intended audience is composed of DLM System Administrators and community-based end-users. This document assumes that the reader has some experience installing and administering DataPlatform applications.

Replication concepts

Data Lifecycle Manager terminology

Data Lifecycle Manager (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 App or Service

The web UI that runs on the DPS platform host. The corresponding agent needs to be installed on the clusters.

DLM Engine

The agent required for DLM. Also referred to as the Beacon engine, this replication engine must be installed as a management pack on each cluster 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 HDP cluster. However, a single data center can contain multiple HDP clusters.

IaaS cluster

A full HDP cluster on cloud VMs with Apache services running, such as HDFS, YARN, Ambari, Hiveserver2, Ranger, Atlas, and DLM Engine. Replication behavior is similar to on-premise cluster replication.

The data is on local HDFS.

cloud data lake or data lake

An HDP cluster on the cloud, using VMs, with data retained on cloud storage. A cloud data lake requires minimal services for metadata and governance, such as Hive metastore, Ranger, Atlas, and DLM Engine.

cloud storage	The data is on the cloud. Any storage retained in a cloud account, such as Amazon S3 web service.
on-premise cluster	A full HDP cluster in a data center, with Apache services running, such as HDFS, Yarn, HMS, hiveserver2, Ranger, Atlas and Beacon. Replication behavior is similar to IaaS cluster replication. The data is on local HDFS.
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.
source cluster	The cluster that contains the source data that will be replicated to a destination cluster. Source data could be an HDFS dataset or a Hive database.
destination cluster	The cluster to which an HDFS dataset or Hive database is replicated.
target	The path on the destination cluster to which the HDFS dataset or Hive database is replicated.

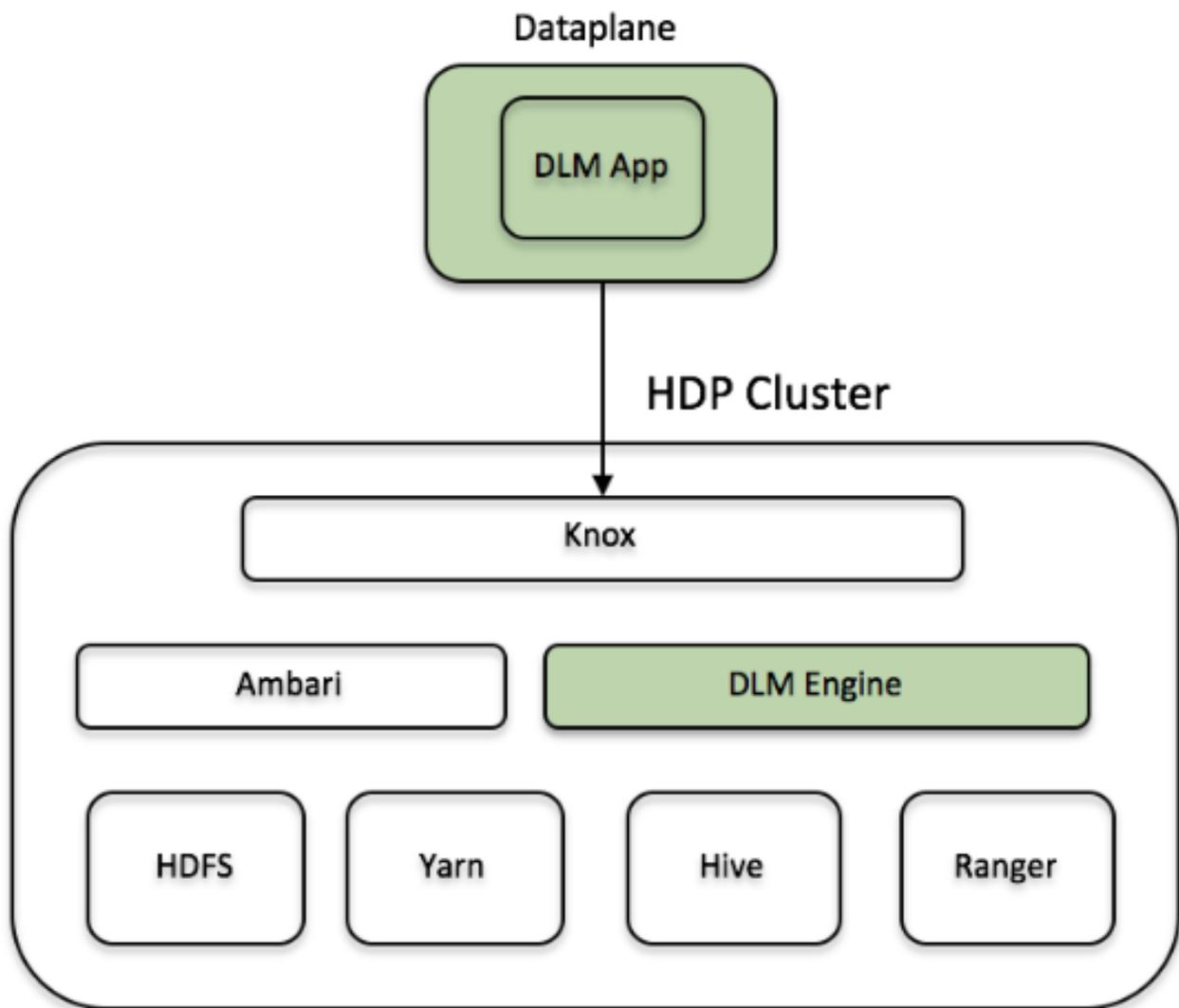
Related Concepts

[DP Platform terminology](#)

Communication with HDP clusters

DPS Platform and the DLM App communicate with the HDP cluster through Knox. Knox SSO is a required configuration for the DPS Platform host.

DLM replication also requires HDFS, YARN, Hive, and Ranger on the cluster. Knox Gateway is recommended to protect data being transferred between clusters.



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 is the destination, what data is replicated, what day and time the replication job occurs, the frequency of job execution, and bandwidth restrictions.

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.

- 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.

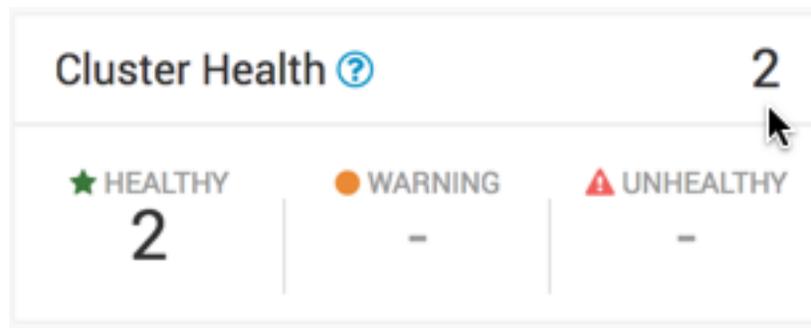
- Achieving a one-hour Recovery Point Objective (RPO) depends on how you set up your replication jobs and the configuration of your environment:
 - Select data in sizes that replicate within 30 minutes.
 - Set replication frequency to 45 minutes or less.
 - Ensure that network bandwidth is sufficient, so that data can move fast enough to meet your RPO.
 - Consider the rate of change of data being replicated.

UI overview

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.

You can investigate the issues associated with clusters that have a warning or unhealthy status by navigating to the Ambari web UI.



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.

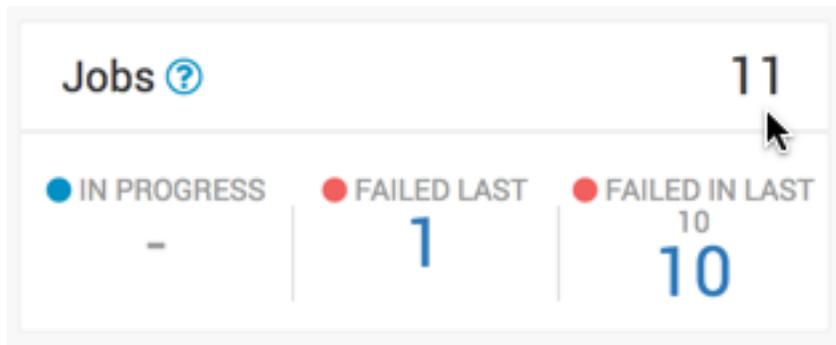
Policies panel

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.

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.

Recent Issues panel

This 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	<p>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:</p> <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.

Clusters map

The Clusters map indicates the geolocation of each cluster, using red, orange, and green markers on the map.

The colored markers indicate the following:

- 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.

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
-------------------	--

Source & Destination	Failed status to see a “View Log” tooltip, which you can click to see the job log.
Service	The names of the source and destination clusters associated with the policy.
Policy	Indicates whether the data being replicated is from HDFS or Hive.
Policy History	The name assigned to the policy. 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 was in progress. Only one run of a job can be in progress at a time. If a job is ignored, you might need to modify its configuration.

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. • Edit Policy: Allows editing of some policy settings. Disabled if a policy is expired. • 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.

Preparing to setup replication policy

Roles required

To perform actions in DPS Platform and associated services, you must be an DPS administrator, an infrastructure administrator.

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.

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

The DataPlane Admin has the following capabilities and restrictions:

- Can access DP Platform and perform all actions in DP Platform related to clusters, users, and enabling services.
- Can access all services enabled with DP 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

Add clusters

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

You must have a cluster registered with the Data Lifecycle Manager to which you replicate data. The cluster must have enough storage to accept data that gets replicated.

Cluster pairing

Before setting up replication between two hadoop clusters, the clusters need to be paired which will validate that the data and metadata can be replicated between two clusters. This will check the configurations of two clusters to communicate with each other.

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.

Pairing 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 LDAP, Kerberos, Ranger, Knox, etc.

Cloud credentials

Register cloud credentials

If you plan to replicate data to a storage cloud account, you must register the cloud credentials, so DLM can access your cloud account. There are two types of cloud storage account that are supported; **Amazon S3** and **WASB**.

Registering Amazon S3 cloud account

You must have valid Amazon S3 credentials to create a cloud account.

Procedure

1. In the DLM UI navigation pane, click **Cloud Credentials > Add**

2. Enter the details in the **Add Cloud Credential** window:

- **Cloud Storage Type** - Select the replication cloud account from the drop-down.
- **Name** - Provide a unique cloud credential name.
- **Authentication Type** - Select the authentication type as **Access Secret Key** from the drop-down.



Note: If you select **IAM Role**, click **Save** to proceed.

- **Access Key** - Enter the valid access key.
- **Secret Key** - Enter the valid secret key.

3. Click **Validate**.

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

Verify that your credentials are listed on the **Cloud Credentials** page.

Considerations 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).

Registering WASB cloud account

You must have valid credentials to create the WASB cloud account.

Procedure

1. Create a storage account in WASB. The **Access key** can be retrieved from the WASB storage account. The Access Key is used in DLM UI to set up the cloud account credentials.
2. In the DLM UI navigation pane, select **Cloud Credentials > Add**.
3. Enter the WASB cloud details in the **Add Cloud Credential** window:
 - **Cloud Storage Type** - Select the replication cloud account from the drop-down.
 - **Name** - Provide a unique cloud credential name.
 - **Storage Account Name** - Provide a name for the storage account.
 - **Access Key** - Paste the Access Key generated from the newly created storage account.
4. Click **Save** to save the changes. Verify that your credentials are listed on the **Cloud Credentials** page.

Considerations for WASB

- You need a cloud account with user credentials that you can enter in DLM, so DLM can access Blob containers.
- Create a WASB storage account using: <https://portal.azure.com>
- Blob containers must have enough space for the replicated data, and write permissions to copy the data. For more information about WASB, see <https://blogs.msdn.microsoft.com/cindygross/2015/02/04/understanding-wasb-and-hadoop-storage-in-azure/>

Data replication use cases

Replication of data using HDFS

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.

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.

File length and checksums are used to determine changed files and validate that the data is copied correctly.

The Ranger policies for the HDFS directory are exported from source Ranger service and replicated to destination Ranger service.

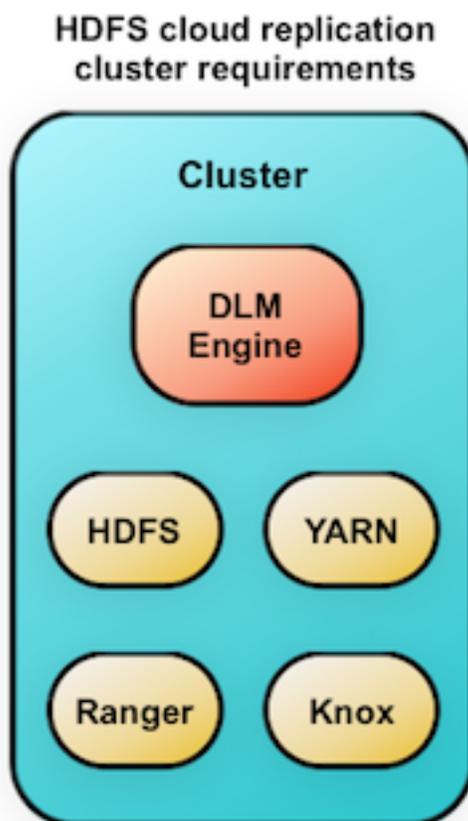
Atlas entities that related to HDFS directory are replicated. If no HDFS path entities are present within Atlas, they are created and then exported.

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

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.



The HDFS replication policy process overview consists of the following:

- 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.

- 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.
- File length and checksums are used to determine changed files and validate that the data is copied correctly.
- The Ranger policies for the HDFS directory are exported from source Ranger service and replicated to destination Ranger service.



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

- Atlas entities related to HDFS directory are replicated. If there are no HDFS path entities are present within Atlas, they are first created and then exported.

On-premise to on-premise replication in HDFS

Before you can begin replicating data using clusters, you must make sure that there are at least a couple of clusters that are registered in your DLM App instance.

Replication of data on-premise to on-premise in HDFS

You must create a replication policy that specifies the data to replicate, the replication schedule, and other settings.

About this task

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

Procedure

1. Pair clusters for replication. Select the two clusters to use for replication and pair them, so they can communicate with each other. For more information, see [Pair clusters for replication](#).
2. Create a replication policy.
3. Select **Policies** and click **Add Policy**. By default, **HDFS** is selected as the service in the **Create Replication Policy** page.
4. Enter the replication policy name and description.
5. Click **SELECT SOURCE** and select type and source cluster from the drop-down.
6. Provide the data replication folder path and click **SELECT DESTINATION**.
7. Select the destination type from the drop-down.
You must select another cluster available in the DLM App as your destination.
8. Select the path and click **VALIDATE**.
9. Once the validation is successful, click **SCHEDULE**.
10. Configure the job settings for the replication policy.
11. Click **ADVANCED SETTINGS** to set up the policy queue.
12. Click **CREATE POLICY**.

The data replication process is enabled.

View job status from the policies page. Verify that the job starts and runs as expected.

On-premise to Amazon S3 replication in HDFS

The process for creating a replication job from on-premise to Amazon S3 is similar to creating one for on-premise to on-premise. The primary difference is that, you must register your cloud account credentials with DLM App instance, so that DLM can access your cloud storage.

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.

Replication of data on-premise to Amazon S3 in HDFS

You must create a new replication policy to replicate data from on-premise Amazon S3 cloud storage.

About this task

Before you create a new replication policy, you must register Amazon S3 cloud account. For more information, see [Register cloud credentials](#).



Note: You can replicate data on-premise to Amazon S3 with a single cluster.

Procedure

1. Select **Policies** and click **Add Policy**. By default, **HDFS** is selected as the service in the **Create Replication Policy** page.
2. Enter the replication policy name and description.
3. Click **SELECT SOURCE** and select type and source cluster from the drop-down.
4. Provide the data replication folder path and click **SELECT DESTINATION**.
5. Select the destination type as **S3** and cloud credential from the drop-down.
6. Provide a folder path bucket_name/path and click **VALIDATE**.
7. Once the validation is successful, click **SCHEDULE**.
8. Configure the job settings for the replication policy.
9. Click **ADVANCED SETTINGS** to set up the policy queue.
10. Click **CREATE POLICY**.

The data replication process is enabled.

View job status from the policies page. Verify that the job starts and runs as expected.

Amazon S3 to on-premise replication in HDFS

You must have a cloud account registered in Amazon S3 before you perform data replication from Amazon S3 to on-premise.

Replication of data from Amazon S3 to on-premise in HDFS

You must create a new replication policy to replicate data from Amazon S3 cloud storage to on-premise.

About this task

Before you create a new replication policy, you must register Amazon S3 cloud account. For more information, see [Register cloud credentials](#).

Procedure

1. Select **Policies** and click **Add Policy**. By default, **HDFS** is selected as the service in the **Create Replication Policy** page.
2. Enter the replication policy name and description.
3. Click **SELECT SOURCE**.
4. Select type as **S3** and cloud credential from the drop-down and enter the S3 source path bucket_name/path.
5. Click **SELECT DESTINATION**.

Make sure you have one or more clusters in the DLM application.

6. Select type as cluster and destination cluster from the drop-down.
7. Enter the destination path and click **VALIDATE**.
8. Once the validation is successful, click **SCHEDULE**.
9. Configure the job settings for the replication policy.

10. Click **ADVANCED SETTINGS** to set up the policy queue.

11. Click **CREATE POLICY**.

The data replication process is enabled.

View job status from the policies page. Verify that the job starts and runs as expected.

On-premise to WASB replication in HDFS

The process for creating a data replication job from on-premise to WASB is similar to creating one for on-premise to on-premise. The primary difference is that, you must register your WASB cloud credentials with DLM App instance, so that DLM can access your WASB cloud storage. You must create a new data replication policy to replicate data from on-premise to WASB.

Replication of data on-premise to WASB in HDFS

You must create a new replication policy to replicate data from on-premise to WASB cloud account.

About this task

Before you create a new replication policy, you must register the WASB cloud account. For more information, see [Register cloud credentials](#).

Procedure

1. Select **Policies** and click **Add Policy**. By default, **HDFS** is selected as the service in the **Create Replication Policy** page.
2. Enter the replication policy name and description.
3. Click **SELECT SOURCE** and select type and source cluster from the drop-down.
4. Provide the data replication folder path and click **SELECT DESTINATION**.
5. Select the destination type as **WASB** and cloud credential from the drop-down.
6. Provide a folder path container_name/path and click **VALIDATE**.
7. Once the validation is successful, click **SCHEDULE**.
8. Configure the job settings for the replication policy.
9. Click **ADVANCED SETTINGS** to set up the policy queue.
10. Click **CREATE POLICY**.

The data replication process is enabled. View job status from the policies page. Verify that the job starts and runs as expected.

WASB to on-premise replication in HDFS

You must setup cloud storage account in WASB before you perform replication from WASB cloud storage to on-premise. Later, create a new replication policy to replicate data from WASB to on-premise.

Replication of data from WASB to on-premise in HDFS

You must create a new replication policy to replicate data from WASB cloud account to on-premise.

About this task

Before you create a new replication policy, you must register the WASB cloud account. For more information, see [Register cloud credentials](#).



Note: You must have a cluster registered with the Data Lifecycle Manager to which you replicate data. The cluster must have enough storage to accept data that gets replicated.

Procedure

1. Select **Policies** and click **Add Policy**. By default, **HDFS** is selected as the service in the **Create Replication Policy** page.
2. Enter the replication policy name and description.
3. Click **SELECT SOURCE**.
4. Select type as **WASB** and cloud credential from the drop-down and enter the path container_name/path for the WASB source.
5. Click **SELECT DESTINATION**.
You must have one or more clusters in the DLM application.
6. Select cluster type and destination cluster from the drop-down.
7. Enter the destination path and click **VALIDATE**.
8. Once the validation is successful, click **SCHEDULE**.
9. Configure the job settings for the replication policy.
10. Click **ADVANCED SETTINGS** to set up the policy queue.
11. Click **CREATE POLICY**.

The data replication process is enabled. View job status from the policies page. Verify that the job starts and runs as expected. For more information, see Viewing job status.

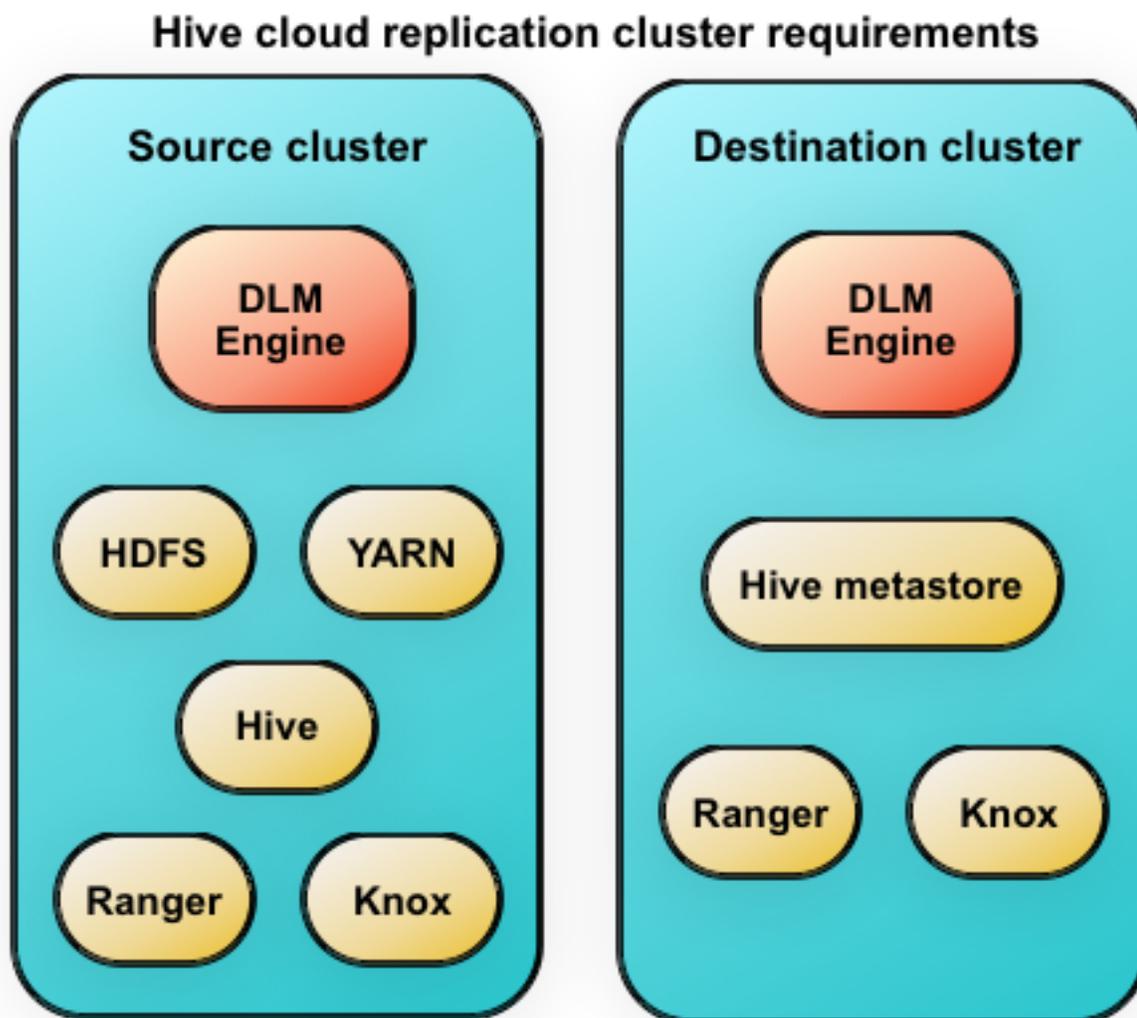
Replication of data using Hive

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.



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.

Functions such as User Defined Functions (UDF) in Hive are replicated. To enable this, UDFs have to be created using a syntax. An example of UDF creation syntax:

```
CREATE FUNCTION [db_name.]function_name AS class_name USING JAR|FILE|
ARCHIVE 'file_uri' [, JAR|FILE|ARCHIVE 'file_uri' ] ;
```

- 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.

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.

On-premise to on-premise replication in Hive

Before you can begin replicating data using clusters on Hive, you must make sure that there are at least a couple of clusters that are registered in your DLM App instance. The replication load happens on the target cluster.

Replication of data on-premise to on-premise in Hive

You must create a replication policy that specifies the data to replicate, the replication schedule, and other settings.

About this task

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

Procedure

1. Select **Policies** and click **Add Policy**. Select **HIVE** as the service in the **Create Replication Policy** page.
2. Enter the replication policy name and description.
3. Click **SELECT SOURCE** and select type and source cluster from the drop-down.
4. Provide the data replication folder path and click **SELECT DESTINATION**.
5. Select the destination type from the drop-down.
You must select another cluster available in the DLM App as your destination.
6. Select the path and click **VALIDATE**.
7. Once the validation is successful, click **SCHEDULE**.
8. Configure the job settings for the replication policy.
9. Click **ADVANCED SETTINGS** to set up the policy queue.
10. Click **CREATE POLICY**.

The data replication process is enabled.

View job status from the policies page.

Verify that the job starts and runs as expected.

On-premise to Amazon S3 replication in Hive

The process for creating a Hive data replication job from on-premise to Amazon S3 is similar to creating one for on-premise to on-premise. The primary difference is that, you must register your cloud account credentials with DLM App instance, so that DLM can access your cloud storage. The replication load happens on the source cluster.

Target cluster setup for Amazon S3

The target cluster for hive cloud replication should be setup on AWS instances with Hive warehouse directory on cloud AWS S3.

The target cluster is data lake cluster with metadata services such as HMS, Ranger, Atlas, and DLM engine. The following configurations are required in this target cluster:

```
hive.metastore.dml.events=false
```

```
hive.repl.cm.enabled=false
```

```
hive.warehouse.subdir.inherit.perms=false
```

```
hive.repl.replica.functions.root.dir=s3a://<bucket_name>/  
<replfunctions_path>
```

```
hive.metastore.warehouse.dir=s3a://<bucket_name>/<warehouse_path>
```

The cluster should have additional Amazon S3 credential configurations for access to Amazon S3 storage buckets. For more information, see https://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.6.5/bk_cloud-data-access/content/authentication-s3.html.

Replication of data on-premise to Amazon S3 in Hive

You must create a new data replication policy to replicate data from on-premise to Amazon S3.

About this task

Before you create a new replication policy, you must register Amazon S3 cloud account. For more information, see [Register cloud credentials](#).



Note: You can replicate data on-premise to Amazon S3 with a single cluster. The metastore must be running on the cloud. There is no requirement to run the HiveServer 2 on the cloud environment.



Important: You must setup target cluster before commencing the replication process.

Procedure

1. Select **Policies** and click **Add Policy**. Select **HIVE** as the service in the **Create Replication Policy** page.
2. Enter the replication policy name and description.
3. Click **SELECT SOURCE** and select type and source cluster from the drop-down.
4. Provide the data replication folder path and click **SELECT DESTINATION**.
5. Choose Type and Destination Cluster.
6. Enter Destination Database and select Cloud Credential from the drop-down.
7. Click **VALIDATE**.
8. Once the validation is successful, click **SCHEDULE**.

9. Configure the job settings for the replication policy.
10. Click **ADVANCED SETTINGS** to set up the policy queue.
11. Click **CREATE POLICY**.

The data replication process is enabled.

View job status from the policies page. Verify that the job starts and runs as expected.

On-premise to WASB replication in Hive

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

Target cluster setup for WASB

The target cluster for hive cloud replication should be setup on WASB instances with Hive warehouse directory on WASB cloud.

The target cluster is data lake cluster with metadata services such as HMS, Ranger, Atlas, and DLM engine. The following configurations are required in this target cluster:

```
hive.metastore.dml.events=false
```

```
hive.repl.cm.enabled=false
```

```
hive.warehouse.subdir.inherit.perms=false
```

```
hive.repl.replica.functions.root.dir=wasb://
<container_name>@<storage_account_name>.blob.core.windows.net/
<replfunctions_path>
```

```
hive.metastore.warehouse.dir=wasb://
<container_name>@<storage_account_name>.blob.core.windows.net/
<warehouse_path>
```

The cluster should have additional WASB credential configurations for access to WASB storage containers, refer to https://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.6.5/bk_cloud-data-access/content/authentication-wasb.html.

Replication of data on-premise to WASB in Hive

You must create a new data replication policy to replicate data from on-premise to WASB. You must setup target cluster before commencing the replication process.

Before you begin

You must register the WASB cloud account. For more information, see [Register cloud credentials](#).



Important: You must setup target cluster before commencing the replication process.

About this task

You can replicate data on-premise to WASB with a single cluster. The metastore must be running on the cloud. There is no requirement to run the HiveServer 2 on the cloud environment.

Procedure

1. Select **Policies** and click **Add Policy**. Select **HIVE** as the service in the **Create Replication Policy** page.
2. Enter the replication policy name and description.

3. Click **SELECT SOURCE** and select type and source cluster from the drop-down.
4. Provide the data replication folder path and click **SELECT DESTINATION**.
5. Choose Type and Destination Cluster.
6. Enter Destination Database and select Cloud Credential from the drop-down.
7. Click **VALIDATE**.
8. Once the validation is successful, click **SCHEDULE**.
9. Configure the job settings for the replication policy.
10. Click **ADVANCED SETTINGS** to set up the policy queue.
11. Click **CREATE POLICY**.

The data replication process is enabled.

View job status from the policies page. Verify that the job starts and runs as expected.

Metadata replication

This page provides information about various types of metadata replication.

Ranger metadata

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.

Atlas metadata

Atlas plugin within DLM is used to replicate Atlas metadata. It uses incremental export to move data across clusters, thereby optimizing the payload for speed and size.

Atlas entities replicated to target cluster are tagged with special classification. Tagging entities allows for easy access to the entities that are part of the available metadata due to replication.



Note: The lineage associated with the entities is not replicated.

On the source cluster, the entity's replicatedTo attribute is updated to indicate the cluster it is being replicated to and on the target cluster the entity's replicatedFrom attribute is modified to indicate its source. Since each cluster has its own identity, the entities that are part of replication are transformed such that, they appear to be native to the cluster they are going to reside within. This involves changing attributes that are indicative of their place of residence. In addition, within Atlas, new entities of type AtlasServer are created. This allows for a central place to access all the servers for which replication has been initiated. Replication audit logs can also be accessed here. Each audit entry has details of every export or import performed for that cluster.

When a DLM Atlas replication job is executed, any Atlas metadata associated with the dataset on source Atlas server, which is replicated, is exported from source, and imported in the target Atlas cluster. The associated replication policy must not be updated or modified during the course of the replication life cycle. You can perform Atlas replication on-premise to on-premise using both HDFS and Hive. You must make sure that there are at least two clusters that are

registered in your DLM App instance. And Atlas must be installed on source and target clusters. Optionally, using Ambari UI, you can verify if Atlas is installed on these clusters. While you create a new Atlas replication policy, do not select **Disable Atlas metadata replication** check-box.

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.

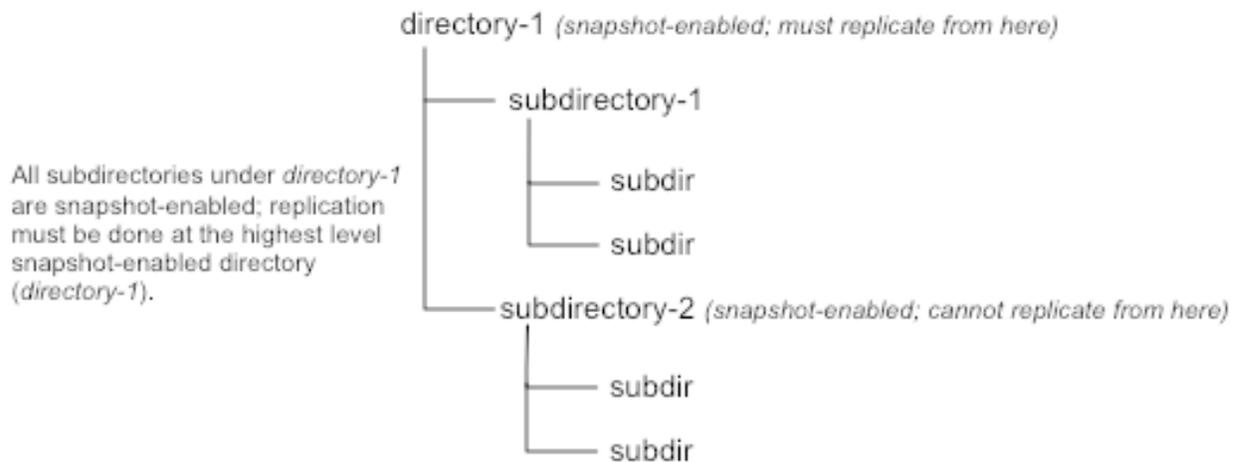
Understanding HDFS Snapshots

HDFS snapshots are read-only point-in-time copies of the filesystem. You can enable snapshots on the entire filesystem, or on a subtree of the filesystem. For DLM, you enable snapshots at a dataset level.

Enabling snapshots on a folder requires HDFS admin permissions, because it impacts the NameNode. When you enable snapshots, all subdirectories are automatically enabled for snapshots as well. So when you create a snapshot copy of a directory, all content in that directory, including subdirectories, is included as part of the copy. 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.

Snapshots must be taken on the highest-level parent directory that is snapshot-enabled. Snapshot operations are not allowed on a directory if one of its parent directories is already snapshot-enabled (snapshottable) or if descendants already contain snapshots. 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.

For example, in the directory tree image below, if *directory-1* is snapshot-enabled but you want to replicate *subdirectory-2*, you cannot select only *subdirectory-2* for replication. You must select *directory-1* for your replication policy.



There is no limit to the number of snapshot-enabled directories you can have. A snapshot-enabled directory can accommodate 65,536 simultaneous snapshots.

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 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.

Replication policy operations

Monitoring replication

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.

Policies page

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

Before you begin

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

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

2. Click the



or



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

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.

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) 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.

Notifications page

You can view the ongoing and completed activities on the notification page 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.

Tuning replication policy (advanced options)

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.

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.

Update replication 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	Jobs	Duration	Last Good	
ACTIVE	contacts-data Every 20m	c1 /test/contacts	c2 /test/contacts	●●●	<1m	18m ago	⋮

- 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

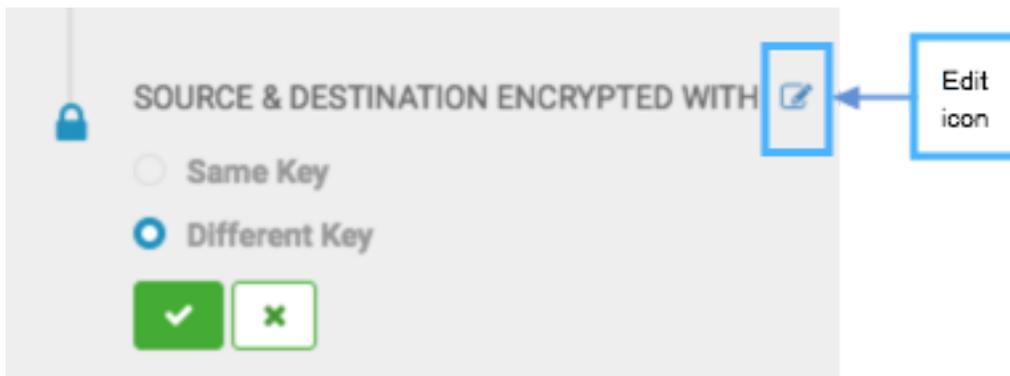
If the Edit option does not display, verify that the policy status is active or suspended. Expired policies cannot be edited.

4. 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.

5. 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.

6. 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.

Browsing data directory

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.

Cloud credentials operations

Update cloud credentials

- 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.

Delete credentials

- 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.

Unregistered credentials

- 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.

Miscellaneous

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

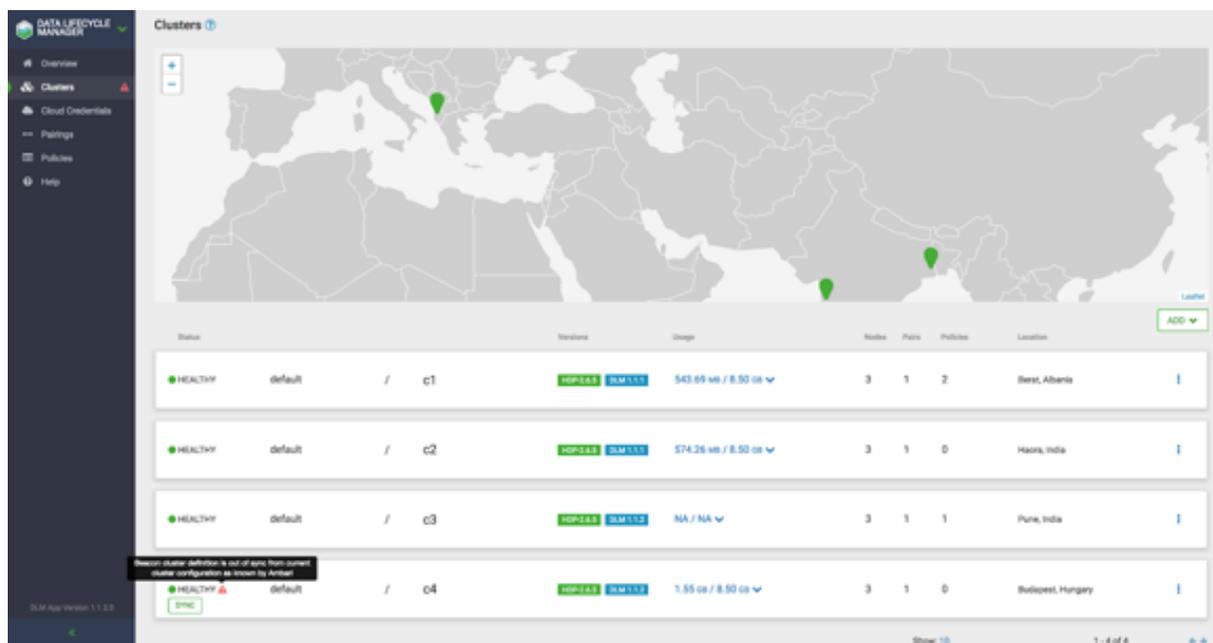
You must use the DLM Infrastructure Admin role to perform this task.

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.

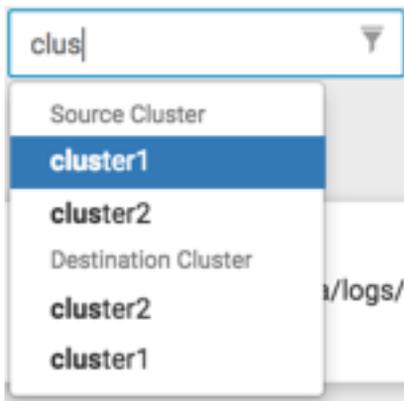
You need the name of the cluster that is offline.

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. For more information, see [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. For more information, see [Cluster pairing](#).
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.

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.
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.

Field	Description	Additional Information
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.

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>.

Hive cloud replication is slow

If the Hive cloud replication is slow, please refer to the documentation related to the target cluster setup instructions.

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.

DLM out of memory

DLM throws out of memory message after upgrading on the source machine.

You must increase heap to about 2GB and try again.