# Snapshot Policies in Replication Manager (Preview)

Date published: 2022-02-25
Date modified: 2022-02-25

# Legal Notice

# Contents

# What's New

A snapshot is a set of metadata information and a point-in-time backup of HDFS files and HBase tables. You can create snapshot policies for HDFS directories and HBase tables in registered classic clusters and SDX Data Lake clusters to take snapshots at regular intervals. Before you create an HDFS snapshot policy for an HDFS directory, you must enable snapshots for the directory in Cloudera Manager.

After a snapshot policy takes a snapshot of the HDFS directory or HBase table, you can perform the following tasks:
- Restore the snapshot to the original directory or table using the **Restore Snapshot** option.
- Restore a directory or table to a different directory or table using the **Restore Snapshot as** option.

# Snapshots page in Replication Manager

The **Snapshots** page in Replication Manager shows the list of snapshot policies with its current status, the classic cluster and the directories and tables that reside on it for which you are taking snapshots using the policy, the configured schedule to take snapshots, and the last policy run. You can edit, suspend, or delete the snapshot policy, and run the snapshot policy on demand on the **Snapshots** page.

When you click a snapshot on the **Snapshots** page, the **Snapshot History** tab appears. Along with the tab, the page also shows the classic cluster name, service name as HDFS or HBase, the directories, paths, or tables chosen for the snapshot policy, number of times the policy ran, and the snapshot location. You can also perform actions such editing the policy, suspend or resume, and remove the snapshot policy.

The **Snapshot History** tab shows the list of snapshots, current status, timestamp of the previous snapshot policy run, the configured schedule for the snapshot, and the time the snapshot expires. After the snapshot exceeds the maximum number specified in **Create Snapshot Policy > Keep** field, the least recent snapshot is deleted by Cloudera Manager and the snapshot is marked as **Expired**. It appears in the **Expired Snapshot** list on the **Snapshots History** tab. Use the **Search Snapshot History** field to search for snapshots. You can view the logs for a snapshot, restore a snapshot, or delete a snapshot.

# Snapshots

A snapshot is a set of metadata information and a point-in-time backup of HDFS files and HBase tables. Snapshots save the state of all or part of the HDFS filesystem or a set of HBase tables at a specified point in time. You can create snapshot policies to take snapshots of HDFS directories and HBase tables in the registered classic clusters at regular intervals.

Snapshots are efficient because they copy only the changed data. You can schedule the snapshot policies to take snapshots at regular intervals or on demand. You can restore the data in these snapshots to a version of the data from any previous snapshot in Replication Manager.

Cloudera Manager and Replication Manager support the following snapshots:

- **HDFS snapshots.** These are point-in-time backup of HDFS directories or the entire filesystem without data cloning. These snapshots appear on the filesystem as read-only directories that can be accessed just like other ordinary directories. Snapshots improve data replication performance and prevent errors caused by changes to a source directory. Before you create a HDFS snapshot policy for an HDFS directory in Replication Manager, you must enable snapshots for the directory in the cluster's Cloudera Manager. Enable snapshots for the HDFS replication source directories to improve HDFS replication performance. Enable snapshot for Hive warehouse directories to improve Hive replications. For more information about snapshots and how to enable HDFS snapshots, see HDFS snapshots.
- **HBase snapshots.** These are point-in-time backups of tables without making data copies and with minimal impact on RegionServers. Snapshots are enabled by default. For more information, see HBase snapshots.

You can use snapshots to restore a directory or table (to a previous version) to the original location or to a different location. For example, you enable snapshots for an HDFS directory **/hive/tmp** in Cloudera Manager on Monday. You then create a snapshot policy **hive-tmp-snapshot** for the directory in Replication Manager where you set a daily schedule for taking snapshots. After a few days, you notice that a user deleted some files by mistake on Thursday. To restore the deleted files, you can restore the snapshot taken on Wednesday that has the deleted files. To restore the files, go to the **Replication Manager > Snapshots Policies** page, click the snapshot policy to view the snapshot history. Click **Restore Snapshot** action for the snapshot taken on Wednesday. This action triggers an HDFS replication policy job to restore the data.

**Note:** When a file or directory is deleted, the snapshot ensures that the file blocks corresponding to the deleted files or directories are not removed from the file system. Only the metadata is modified to reflect the deletion.

# HDFS snapshots

You can schedule taking HDFS snapshots for replication in the Replication Manager. 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. In Replication Manager, you take snapshots at a dataset level. Understanding how snapshots work and some of the benefits and costs involved can help you to decide whether or not to enable snapshots.

To improve the performance and consistency of HDFS replications, enable the HDFS replication source directories for snapshots, and for Hive replications, enable the Hive warehouse directory for snapshots. For more information, see HDFS snapshots.

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

Take snapshots 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.

For example, in the following directory tree image, 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.

## Requirements and benefits of HDFS snapshots

You might want to consider the benefits and memory cost of using snapshots. Verify the requirements before you enable snapshots.

**Requirements**

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 be no 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.

**Benefits**

Snapshot-based replication helps you to avoid unnecessary copying of 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, because 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.

**Memory cost**

There is a memory cost to enable and maintain 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 where 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 run.

# HBase snapshots

HBase snapshots are point-in-time backups of the table without making data copies, and with minimal impact on RegionServers. Snapshots are enabled by default. You can create HBase snapshot policies in Replication Manager to take snapshots of a table at regular intervals.

HBase snapshots support enables you to take a snapshot of a table without much impact on RegionServers, because snapshot, clone, and restore operations do not involve data copying. For more information, see [HBase snapshots](#).

# Snapshot policies

You can create an HDFS snapshot policy for a snapshottable HDFS directory  A directory is *snapshottable* after it has been enabled for snapshots, or because a parent directory is enabled for snapshots in Cloudera Manager. Subdirectories of a snapshottable directory are included in the snapshot. You can also create HBase snapshot policies for HBase tables. Snapshots for HBase tables are enabled by default. Snapshot policies take snapshots of the directory or table at regular intervals.

## Enabling and taking snapshots in Cloudera Manager

Before you take snapshots (in Cloudera Manager) or create snapshot policies (in Replication Manager) for HDFS directories, you must enable snapshots for the directories in Cloudera Manager.

1. To enable snapshots for HDFS directories, navigate to the directory on the **Cloudera Manager > Clusters > HDFS service > File Browser** tab, and click **Enable Snapshots**.

   **Note:** If you enable snapshots for a directory, you cannot enable snapshots for its parent directory. Snapshots can be taken only on directories that have snapshots enabled.

   **Tip:** To disable snapshots for a directory that has snapshots enabled, click **Disable Snapshots**. Ensure that the snapshots of the directory are deleted before you disable snapshots for the directory.

2. To take a snapshot of a directory or table, perform the following steps:
   a. Navigate to the directory or folder.
   b. Click **Take Snapshot** in the drop-down menu next to the full file path.
   c. Specify a unique name for the snapshot.

   The snapshot name with the timestamp of its creation appears in the snapshot list.

3. Click  **Actions > Delete** to delete a snapshot. After you delete a snapshot, you can restore it, if required.

**After you finish:** You can create a snapshot policy to schedule taking snapshots of the directory at regular intervals in Replication Manager.

# Creating snapshot policies in Replication Manager

After you enable snapshots for an HDFS directory in Cloudera Manager, you can create an HDFS snapshot policy to schedule taking snapshots of the directory at regular intervals in Replication Manager. You can also create an HBase snapshot policy to take snapshots of HBase tables at regular intervals.

1. In the Replication Manager, go to the **Snapshots** page.

2. Click **Create Snapshot.**

3. In the **Create Snapshot Policy** wizard, enter or choose the following options.

4. **Source Cluster** - Choose the classic cluster where the directory or table resides for which you want to take a snapshot.

5. **Service** - Choose **HDFS** to choose one or more HDFS directories or **HBase** to choose the required HBase tables.

6. **Snapshot name** - Enter a name for the snapshot which is unique in the selected cluster.

   **Important:** Ensure that the snapshot policy name neither contains the characters **% . ; / \** nor any character that is not ASCII printable, which includes the ASCII characters less than 32 and the ASCII characters that are greater than or equal to 127.

7. If you choose **HDFS**, enter the following options as required:

   ○ **HDFS Directory Path** - Choose one or more directory paths. The directory path appears here after you enable snapshots for the directory in Cloudera Manager.

   ○ **Schedule** - Choose a schedule to take snapshots of the directory. In the **Keep** field, enter the maximum number of snapshots to be retained as currently available snapshots by the Replication Manager. When a snapshot exceeds the maximum number, the least recent snapshot is deleted by Cloudera Manager and the snapshot is marked as an **Expired**. It appears in the **Expired Snapshot** list on the **Snapshots History** tab.

   ○ **Alerts** - Choose one or more options to generate alerts during various state changes in the replication workflow. You can choose to generate alerts during **On Failure**, **On Start**, **On Success**, and **On Abort** of the replication job.

8. If you choose **HBase**, enter the following options as required:

   ○ **HBase Table** - Choose one or more tables.

   ○ **Schedule** - Choose a schedule to take snapshots of the directory. In the **Keep** field, enter the maximum number of snapshots to be retained as currently

available snapshots by the Replication Manager. When a snapshot exceeds the maximum number, the least recent snapshot is deleted by Cloudera Manager and the snapshot is marked as an **Expired**. It appears in the **Expired Snapshot** list on the **Snapshots History** tab.

- ○ **Alerts** - Choose one or more options to generate alerts during various state changes in the replication workflow. You can choose to generate alerts during **On Failure**, **On Start**, **On Success**, and **On Abort** of the replication job.

9. Click **Create**.
   The snapshot policy appears on the **Snapshots** page.

**Tip:** The snapshot policy appears on the **Cloudera Manager > Replication > Snapshot** page as well.

## Snapshot policy operations

On the **Snapshots** page, you can edit, suspend, or delete the snapshot policy. You can run the snapshot policy on demand. When you click a snapshot policy, you can view the snapshot history.

Use the **Actions** menu to perform the following policy operations on a snapshot policy:

1. Click **Edit** to make changes to the snapshot policy in the **Edit Snapshot Policy** wizard. The next policy run uses the changed configuration to take snapshots.

2. Click **Suspend** to pause taking snapshots. You can resume taking snapshots when required.

3. Click **Remove** to remove the snapshot policy. Replication Manager removes the snapshot policy and all the snapshots taken by the policy.

4. Click the snapshot name to view the snapshots taken by a snapshot policy on the **Snapshot History** tab.

# View snapshot history

When you click a snapshot on the **Snapshots** page, the **Snapshot History** tab appears. Along with the tab, the page also shows the classic cluster name, service name as HDFS or HBase, the directories, paths, or tables chosen for the snapshot policy, number of times the policy ran, and the snapshot location. You can also perform actions on the snapshot policy or on a specific snapshot on this page.

The following options are available in the **Actions** menu for the snapshot policy:

- **Edit** - Opens the **Edit Snapshot Policy** wizard where you can edit the snapshot as required.

- **Suspend** - Pauses the snapshot policy. You can resume the snapshot policy run when required.

- **Remove** - Removes the snapshot policy and all the snapshots taken by the snapshot policy.

The **Snapshot History** tab shows the list of snapshots, current status, timestamp of the previous snapshot policy run, the configured schedule for the snapshot, and the time the snapshot expires. After the snapshot exceeds the maximum number specified in **Create Snapshot Policy > Keep** field, the least recent snapshot is deleted by Cloudera Manager and the snapshot is marked as **Expired**. It appears in the **Expired Snapshot** list on the **Snapshots History** tab. Use the **Search Snapshot History** field to search for snapshots.

The **Snapshot History** tab lists the currently available snapshots (**Finished** status), next scheduled snapshot policy run (**Scheduled Next** status), suspended snapshots (**Suspended** status), and expired snapshots (appears under the **Expired Snapshots** list).
You can perform the following actions on the currently available snapshots:

- **View Log** - Shows the log details in a window. You can download the log file to your machine as a text file, copy the contents, or open the log window in a new tab.

- **Restore Snapshot** - Opens the **Restore Snapshot Policy** wizard to restore the snapshot to its original location. For more information, see [Restoring snapshots](#).

- **Restore Snapshot As** - Opens the **Restore Snapshot Policy** wizard to restore the snapshot to another location. For more information, see [Restoring snapshots](#).

- **Delete Snapshot** - Deletes the snapshot.

# Restoring snapshots in Replication Manager

In Replication Manager, you can choose to restore the data to a version from a previous snapshot. You can also choose to restore it to its original location or to another location.

**Before you begin:** Ensure you have adequate disk space for restored snapshots in the source cluster.

1. Click the snapshot policy on the **Snapshots** page.

2. On the **Snapshot History** tab, click the **Actions** menu for the required snapshot and choose one of the following options:
   - **Restore Snapshot -** Choose to restore the snapshot to the same directory or table. This action overwrites existing data.
   - **Restore Snapshot As** - Choose to restore the snapshot to another specified directory or table. If the specified directory or table exists or has similar data, this action overwrites it.

3. In the **Restore Snapshot** wizard, enter the following details:

4. **Source Cluster** - Shows the source cluster of the snapshot.

5. **Service** - Shows the service for the snapshot.

6. **Recover Data from Snapshot** - Shows the timestamp of the scheduled snapshot policy.

7. The following options appear for HDFS snapshots:
   - **Directory path to restore** - Shows the HDFS directory of the snapshot that is to be restored.
   - **The new directory path** - Enter the new location to save the snapshot. This option appears when you choose **Restore Snapshot As.**

   **Caution:** If the directory or table entered already exists, that directory is overwritten.

8. The following options appear for HBase snapshots:
   - **Table to restore** - Shows the table to restore.
   - **Restore as** - Enter a name of the table to which the snapshot is to be restored. This option appears when you choose **Restore Snapshot As.**

   **Caution:** If the directory or table entered already exists, that directory is overwritten.

9. **Use the HDFS 'copy' command** - Choose this option to copy the contents of the snapshot as a subdirectory or as files in the target directory. This restore option takes time to complete and does not require credentials in a secure cluster.

10. **Use DistCp / MapReduce -** Choose this option to merge the target directory with the contents of the source snapshot. This restore operation is faster and requires credentials (**Run As Username**) in secure clusters.

---

**Important:** Before you restore HDFS snapshots using DistCp, ensure that the user name specified in the **Run as Username** in the HDFS snapshot policy is part of the superuser user-group.

To add a user to the supergroup user-group, run the following commands:

```
sudo adduser [***username***]
sudo usermod -a -G hdfs [***username***]
sudo usermod -a -G hadoop [***username***]
id -Gn [***username***]
kinit -kt /cdep/keytabs/hdfs.keytab hdfs
hdfs dfs -mkdir /user/[***username***]
hdfs dfs -chown
[***username***]:[***username***]/user/[***username***]
sudo groupadd supergroup
sudo usermod -a -G supergroup [***username***]
hdfs groups [***username***]
```

---

Choose the following options as required:

| Option | Description |
|---|---|
| **MapReduce Service** | Choose a MapReduce or YARN service in the cluster. DistributedCopy (distcp) restores the directories and it increases the speed of restoration. |
| **Scheduler Pool** | Enter the name of a resource pool in the field. The value you enter is used by the MapReduce Service you specified when Cloudera Manager runs the MapReduce job for the replication.<br><br>The job specifies the value using If a MapReduce or YARN service is not available, a normal copy is performed.one of these properties:<br>    ● MapReduce – Fair scheduler: **mapred.fairscheduler.pool**<br>    ● MapReduce – Capacity scheduler: **queue.name**<br>    ● YARN – **mapreduce.job.queuename** |

| | |
|---|---|
| **Run as Username** | Uses the default username. Enter the username for kerberized clusters. |
| **Log Path** | Enter an alternate path for the log files. |
| **Maximum Map Slots** | The maximum number of map tasks (simultaneous copies) per replication job. The default value is 20. |
| **Maximum Bandwidth** | The maximum bandwidth in MB for a mapper. Each map task is throttled to consume only the specified bandwidth such that the net used bandwidth tends towards the specified value. The default is 100 MB. |
| **Replication Strategy** | Choose **Static** to distribute the file replication tasks among the mappers statically. Static replication distributes the file replication tasks among the mappers up front to achieve a uniform distribution based on the file sizes. **Dynamic** replication distributes the file replication tasks in small sets to the mappers, and as each mapper completes its tasks, it dynamically acquires and processes the next unallocated set of tasks. The default replication strategy is **Dynamic**. |
| **Delete Policy** | You can choose how to handle the replicated data in the destination cluster if the files are deleted on the source and also determine the handling of files in the destination location that are unrelated to the source. |
| **Preserve** | You can choose to preserve the block size, replication count, permissions (including ACLs), and extended attributes (XAttrs) as they exist on the source file system, or use the settings as configured on the destination file system. |
| **Abort on Error** | Aborts the restore job if an error appears. |
| **Skip Checksum Checks** | Does not validate the checksum checks on the copied files. By default, checksums are checked for the copied files. |

11. Click **Restore**.