

NEXT GENERATION DATA PROTECTION WITH ONEFS SNAPSHOTIQ

ABSTRACT

Most file systems are a thin layer of organization on top of a block device and cannot efficiently address data at large scale. This paper focuses on OneFS, a modern file system that meets the unique needs of Big Data. OneFS includes SnapshotIQ, a native data protection capability, which enables enterprises to reduce storage costs and footprint and increase data availability, without sacrificing management simplicity.

April 2017

TABLE OF CONTENTS

SNAPSHOTS OVERVIEW.....3

 Data Protection with SnapshotIQ 4

 Reading from a Snapshot..... 8

 Painting Algorithm 8

 Taking a Snapshot 9

 Writing to a Snapshot..... 10

 Deleting a Snapshot..... 11

 Restoring a Snapshot..... 12

 User Driven File Recovery 13

 Snapshot scheduling 15

 Snapshot Performance..... 16

 SnapshotIQ Licensing 18

 Snapshot Reserve..... 18

 Snapshot Overhead 19

SNAPSHOT BEST PRACTICES 19

SNAPSHOT CONSIDERATIONS..... 20

SNAPSHOT AND ONEFS FEATURE INTEGRATION 20

 File clones 23

CONCLUSION 24

INTRODUCTION

Information technology managers across most areas of commerce are grappling with the challenges presented by explosive file data growth, which significantly raises the cost and complexity of storage environments.

This proliferation of unstructured data has left traditional storage architectures unable to satisfy the demands of this growth and has necessitated the development of a new generation of storage technologies. Additionally, broader data retention requirements, regulatory compliance, tighter availability service level agreements (SLAs) with internal/external customers, and cloud and virtualization initiatives are only serving to compound this issue.

Snapshots Overview

The availability and protection of data can be effectively described in terms of a continuum:

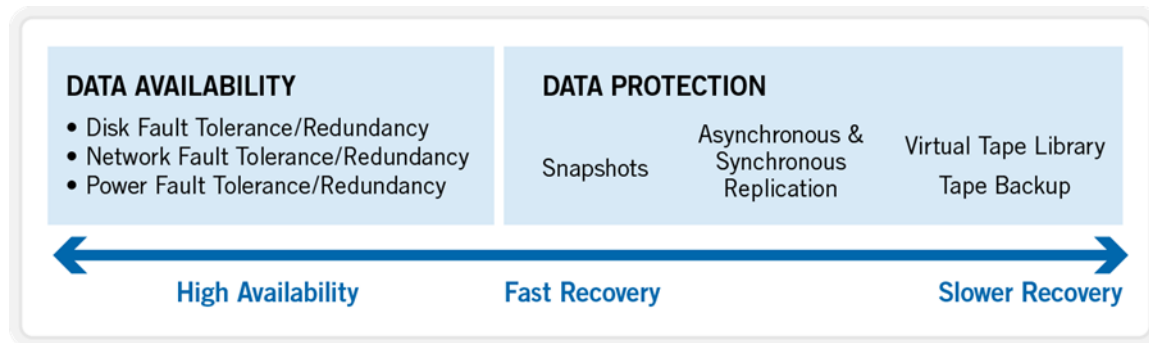


Figure 1: Data Protection Continuum

At the beginning of the continuum sits high availability. This requirement is usually satisfied by redundancy and fault tolerant designs. The goal here is continuous availability and the avoidance of downtime by the use of redundant components and services.

Further along the continuum lie the data recovery approaches in order of decreasing timeliness. These solutions typically include point-in-time snapshots for fast recovery, followed by replication, and, finally, backup to tape or a virtual tape library.

The following diagram illustrates how the core components of the Isilon data protection portfolio align with the notion of an availability and protection continuum and associated recovery objectives.

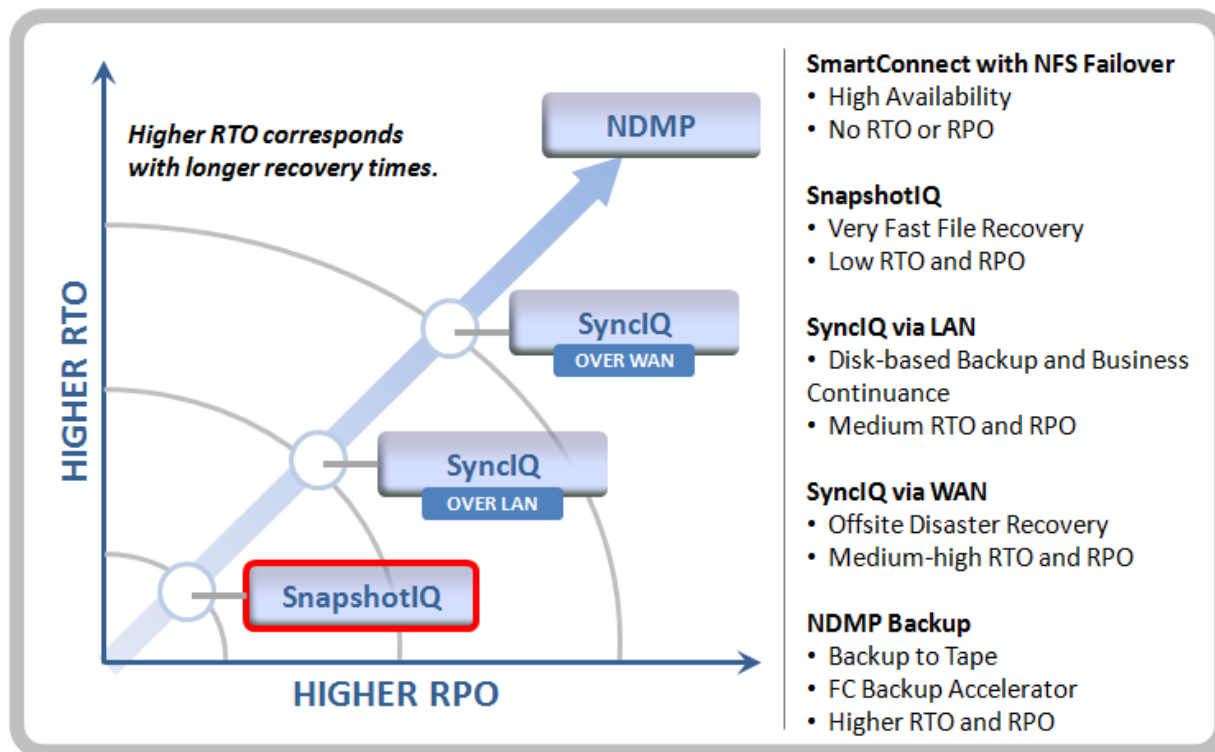


Figure 2: Isilon Data Protection technology alignment with protection continuum

The recovery time objective (RTO) of a snapshot can be very small and the recovery point objective (RPO) is also highly flexible with the use of rich policies and granular schedules.

Data Protection with SnapshotIQ

Isilon SnapshotIQ™ can take read-only, point-in-time copies (snapshots) of any directory or subdirectory within OneFS. When a snapshot is taken, it preserves the exact state of a file system at that instant, which can then be accessed later. This immutable, point-in-time copy has a variety of applications. For example, snapshots can be used to make consistent backups, or to restore files which were inadvertently changed or deleted. Snapshots are also for quickly identifying filesystem changes.

Typically, only a small percentage of the file system changes in a given period of time, so most storage systems do not keep fully redundant copies. Instead, generally only the differences from the currently file system state are stored.

An Isilon OneFS snapshot is basically a logical pointer to data that is stored on a cluster at a particular point in time. Each snapshot references a specific directory under OneFS, and includes all the files stored in that directory and its subdirectories. If the data referenced by a snapshot is modified, the snapshot stores a physical copy of the data that was modified. Snapshots are either created according to user configuration, or are automatically generated by OneFS to facilitate system operations.

Snapshots can be identified and located either by a unique name or a system generated snapshot ID. A snapshot name is specified by a user and assigned to the virtual directory which contains the snapshot. A snapshot ID is a numerical identifier that OneFS automatically assigns to a snapshot.

OneFS Snapshots are highly scalable and typically take less than one second to create. They create little performance overhead, regardless of the level of activity of the file system, the size of the file system, or the size of the directory being copied. Also, only the changed blocks of a file are stored in a snapshot, thereby ensuring highly-efficient storage capacity utilization. User access to the available snapshots is via a special hidden 'portal' directory under each file system directory.

Isilon SnapshotIQ can also create up to twenty thousand snapshots on a cluster. This large quantity provides a substantial benefit over the majority of other snapshot implementations, because it allows the snapshot intervals to be far shorter, and hence offer more granular recovery point objectives (RPOs).

Snapshot Architecture

SnapshotIQ has several fundamental differences as compared to most snapshot implementations. The most significant of these are:

- **Directory based:** OneFS snapshots are per-directory based. This is in contrast to the traditional approach, where snapshots are taken at a file system or volume boundary.
- **Logical snapshot process:** Since OneFS manages and protects data at the file-level, there is no inherent, block-level indirection layer for snapshots to use. Instead, OneFS takes copies of files, or pieces of files (logical blocks and inodes) in what's termed a logical snapshot process.
- **Snapshot space allocation:** There is no requirement for reserved space for snapshots in OneFS. Snapshots can use as much or little of the available file system space as desirable. A snapshot reserve can be configured if preferred, although this will be an accounting reservation rather than a hard limit.

The process of capturing a snapshot in OneFS is near instantaneous. However, there is a small amount of snapshot preparation work that has to occur. The moment a snapshot is taken, it consumes zero space until any file writes occur to the data it contains.

Any changes to a dataset are then recorded in the pertinent snapshot inodes, which contain only referral ('ditto') records, until any of the logical blocks they reference are altered or another snapshot is taken. In order to reconstruct data from a particular snapshot, OneFS will look through all of the more recent versions snapshot tracking files (STFs) until it reaches HEAD (current version). In so doing, it will systematically find all the changes and recreate the point-in-time view of that dataset.

The process of capturing a snapshot in OneFS is near instantaneous. However, there is a small amount of snapshot preparation work that has to occur. The moment a snapshot is taken, it consumes zero space until any file writes occur to the data it contains.

Any changes to a dataset are then recorded in the pertinent snapshot inodes, which contain only referral ('ditto') records, until any of the logical blocks they reference are altered or another snapshot is taken. In order to reconstruct data from a particular snapshot, OneFS will look through all of the more recent versions snapshot tracking files (STFs) until it reaches HEAD (current version). In so doing, it will systematically find all the changes and recreate the point-in-time view of that dataset.

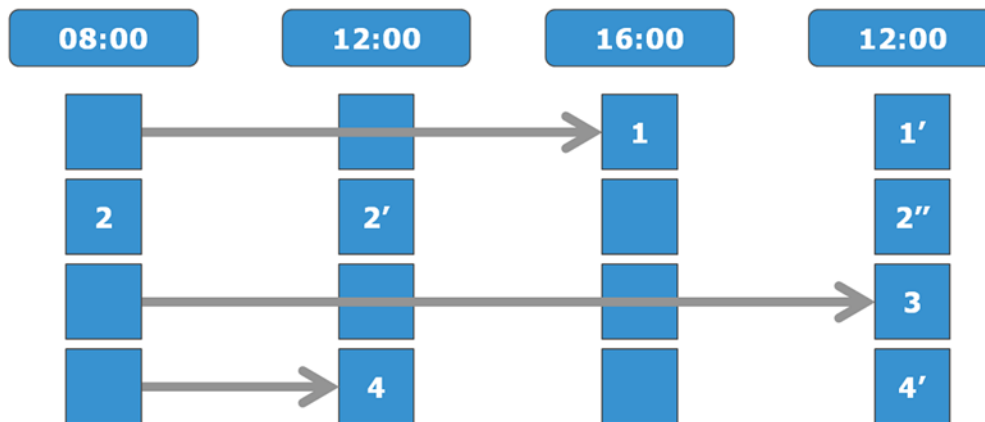


Figure 3: Snapshot Read Chain

In this example, in order to reconstruct file 1-2-3-4 as it appeared at 08:00, SnapshotIQ would need to read snapshots 12:00, 16:00 (4PM) and 20:00 (8PM).

In extreme cases, OneFS might have to 'paint' through hundreds or thousands of snapshots in order to get the correct older version.

Snapshot Tracking Files

Snapshot tracking files (STF) are the main data structure associated with a snapshot. A snapshot tracking file has three major purposes:

- Indicating which snapshots are active.
- Storing snapshot attributes, such as usage data, creation time, and root directory paths.
- Recording a list of LINs modified in the snapshot, which can be freed when the snapshot is deleted.

Snapshot Tracking Files are a special file type with several unique characteristics, and are involved in the full snapshot lifecycle, including the creation, storing any changes, and deletion of snapshots.

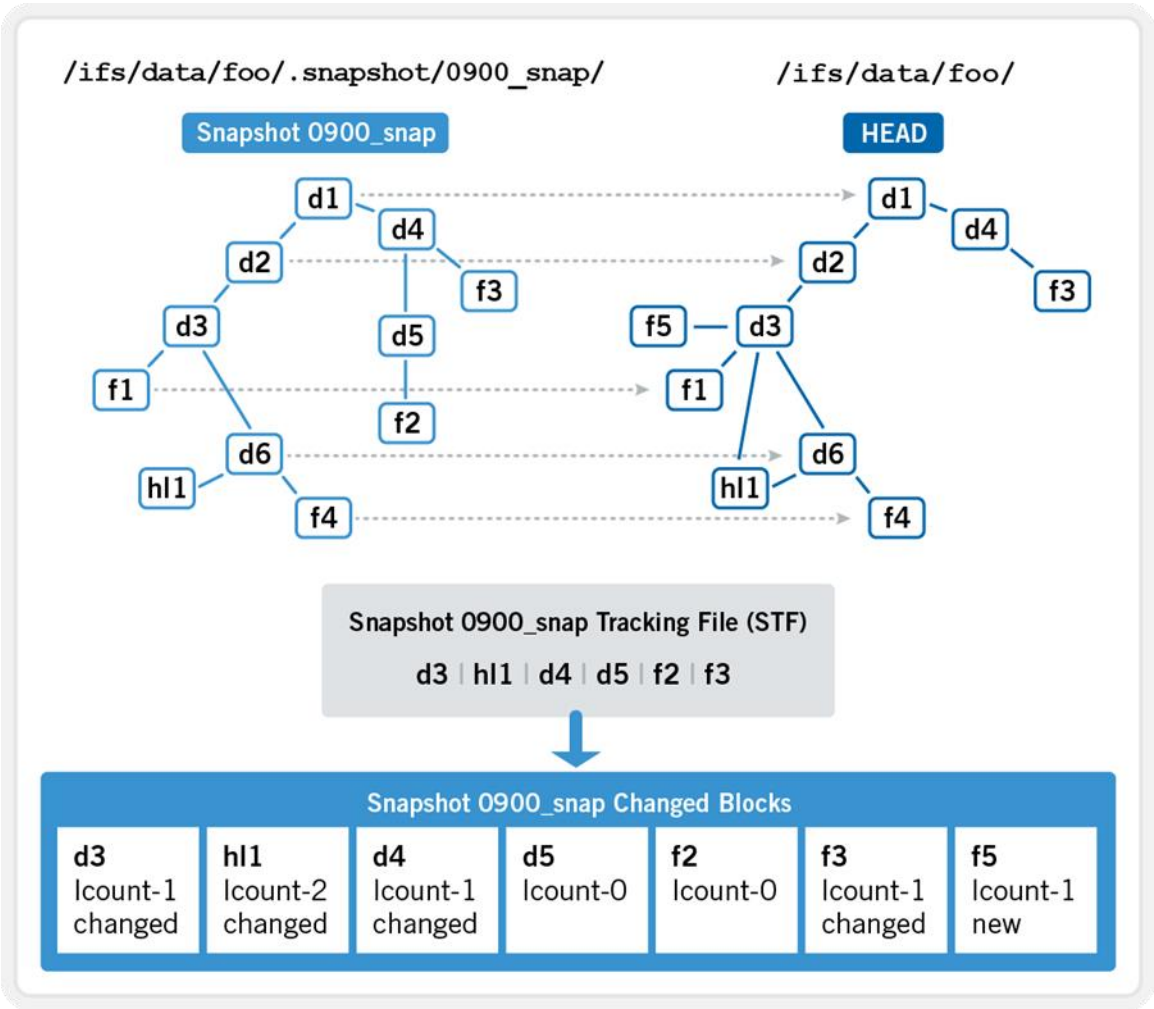


Figure 4: Snapshot Change Tracking

LIN Table

Finding the next version of a particular file is a fast operation due to the structure of the LIN (logical inode) table. While the LIN table is frequently called a table, it's actually a B-Tree. It's sorted by (LIN,Version), so finding the next newer version of a LIN is an inexpensive operation.

A LIN is fundamentally the identity of an object and possesses the following attributes:

- Directories point to LINs.
- The LIN Table maps a LIN to blocks on disk for the inode.
- The Inode contains the map of data ranges to blocks on disk.
- The LIN table maps the LIN, SnapID that identifies a file version to the location on disk that houses the metadata about that file.

Given a LIN:SnapID pairing, finding the next higher SnapID is a fast and efficient process, and the HEAD version is always represented by the biggest SnapID

LIN	Snap ID	Inodes
1:abcd:1234	98	1,1,8;2,3,9;23,7,16
1:abcd:1234	100	2,3,5;5,7,2;13,3,1
2:0021:abcd	MAX	...
2:0031:1cf8	MAX	...

/ifs/data

uid =2300
gid =2300
mtime =whenever
last_snap_id = 100
governing_snaps = { 100 }

foo 2:0021:abcd / MAX
file 1:abcd:1234 / MAX

Figure 5: LIN Table Example

In the example above, 1:abcd:1234/MAX is not in the LIN Tree. This suggests that the file was deleted sometime after snapshot 100 was taken.

Copy or Redirect on Write

SnapshotIQ uses both copy on write (CoW) and redirect on write (RoW) strategies for its differential snapshots, and utilizes the most appropriate method for a given situation. Both have pros and cons, and OneFS dynamically picks which flavor to use in order to maximize performance and keep overhead to a minimum.

With copy on write, as the name suggests, a new write to HEAD results in the old blocks being copied out to the snapshot version first. Although this incurs a double write penalty, it results in less fragmentation of the HEAD file, which is better for cache prefetch, etc. Typically, CoW is most prevalent in OneFS, and is primarily used for small changes, inodes and directories.

Redirect on write, on the other hand, avoids the double write penalty by writing changes to a snapshot protected file directly to another free area of the filesystem. However, the flip side to this is increased file fragmentation. Since file contiguity not maintained by virtue of writing changes to other filesystem regions, RoW in OneFS is used for more substantial changes such as deletes and large sequential writes.

Note: There is no reserved space requirement for snapshots in OneFS. Snapshots can use as much or little of the available file system space as desirable. A snapshot reserve can be configured, if preferred.

Accessing Snapshots via the .snapshot Directory

The .snapshot directory is the namespace entry point for accessing snapshots. It is also referred to as a 'portal' since it provides a path from the HEAD version of the file system into the associated snapshot(s). Much like the UNIX '.' and '..' directory notations, the .snapshot directory is a virtual entry in that it does not have an associated record in the directory B-tree, much like '.' and '..'.

Within the .snapshot directory, there is a virtual entry for every snapshot which either governs the parent directory or was taken on a descendant of the directory. The virtual entries consist of the parent directory's LIN and the snapshot's ID. This is implemented by listing the snapshot names directory, and filtering the entries which do not apply. If a process opens an entry, it gets a snapshot version of the portal's parent. The .snapshot portal can also be disabled entirely, have access restricted to the root and subdirectories, and be renamed to 'snapshot', rather than '.snapshot'.

Snapshot Management

Reading from a Snapshot

The illustration below shows three snapshots of a given file with logical inode (LIN) value “1:abcd:1234”. In this case, blocks 3 and 4 were changed after the first snapshot (Snap_ID 98) was taken and before the second (Snap_ID 100), and blocks 0 and 4 were changed after the second snapshot was taken.

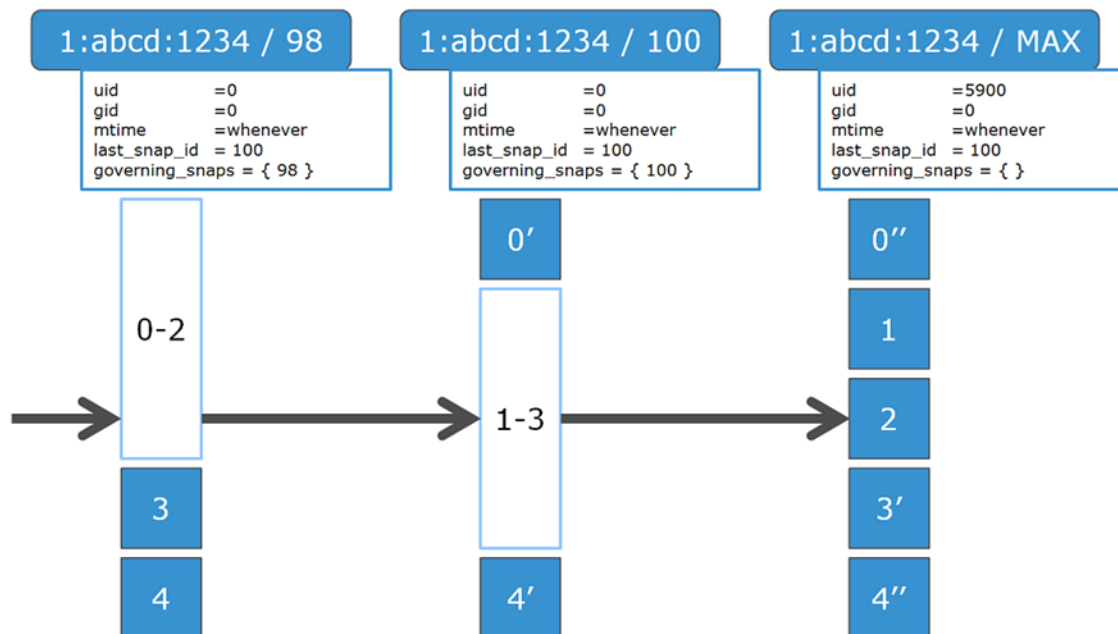


Figure 6: Reading from a Snapshot

When the data is not in the snapshot, the block tree of the inode on the snapshot doesn't point to a real data block. Instead it has a flag marking it as a "Ditto Block". A Ditto-block means that the data is the same as the next newer version of the file, so OneFS will automatically look forward to find the newer version of the block.

The arrow represents the logic for reading block 2 from snapshot 98. Since it wasn't changed in snapshot 98, the read has to fall forward to snapshot 100. It wasn't changed there either, so the write falls forward to the head version of the file.

The performance implications here should be clear: If you have thousands of snapshots of the same unchanged file, reading from the oldest snapshot can potentially be somewhat slow.

Painting Algorithm

When a file is written to, the system needs to do a small amount of work to determine if the file is part of a snapshot. If so, a copy of the old data needs to be kept. This is done via a process known as the "painting algorithm".

The system keeps the most recent snapshot ID in a cluster-wide global variable. When a file is modified, OneFS looks first at the file's last_snap_id. If the last_snap_id is not the most recent snap_id, there is a likelihood that the governing_snaps information in the file is out of date. In this case, OneFS recursively searches the parent directories until it finds up-to-date information, and then uses the correct directory's governing-snaps information. For example:

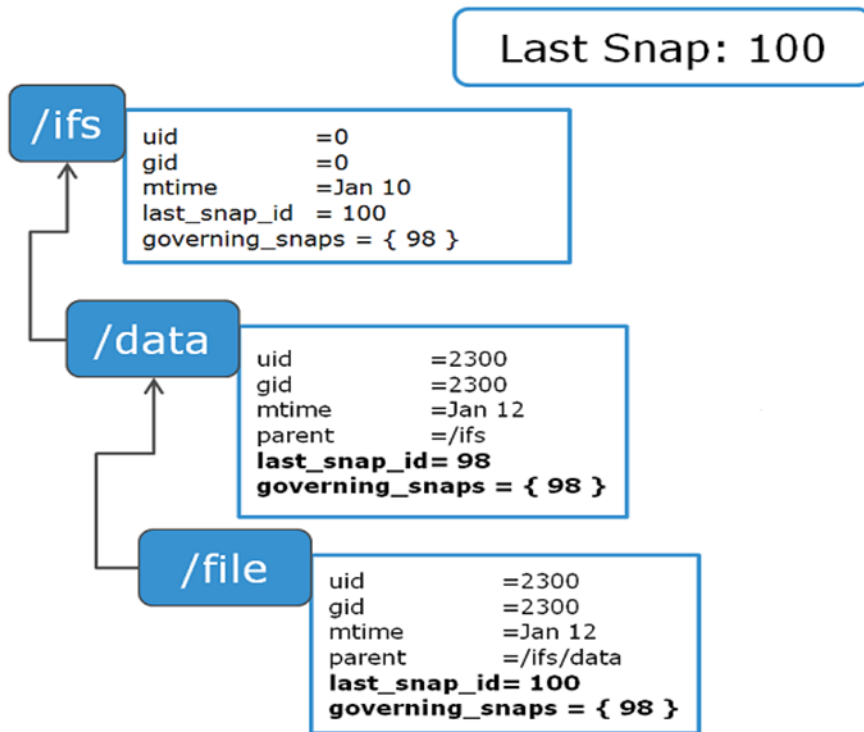


Figure 7: SnapshotIQ Painting Algorithm

In the instance above, /ifs has a governing_snaps ID of “98”. This implies that snapshots 98 and 99 were taken on another directory; possibly /ifs/home.

Taking a Snapshot

There are two methods to create user snapshots:

- On-demand.
- From a preconfigured snapshot schedule.

The first, on-demand, method is initiated via the simple ‘Capture a new snapshot’ button in the ‘snapshot summary’ tab of the OneFS user interface.

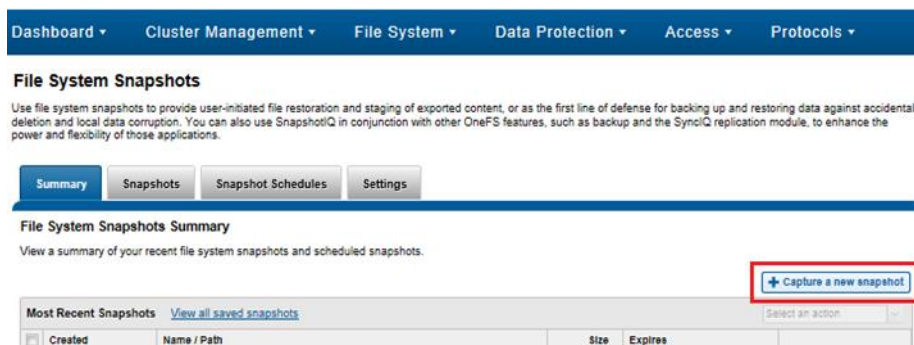


Figure 8: Taking an On-demand Snapshot

On the backend, creating a snapshot involves synchronizing changes across several levels, from the write-back cache (coalescer) to the snapshot tracking file (STF). The main elements of this process are:

- 1. The snapshot acquires an exclusive snapshot lock.
- 2. Coalescers are coordinated and suspended across all nodes.
- 3. The write lock is upgraded to exclusive.
- 4. The maximum snapshot ID is incremented.
- 5. The snapshot (and associated minisnaps) is generated up to the snapshot root.
- 6. A snapshot tracking file is created and updated.
- 7. If requested, a snapshot alias is created.

Writing to a Snapshot

This following figure illustrates the process of updating a file’s metadata, by way of a simple UID change. When this occurs, and OneFS detects that it needs to duplicate the old data, it copies the old blocks to the snapshot. In this case, there was no LIN 1:abcd:1234/98 prior to the write, so any attempt to read this LIN would have fallen forward to the head version.

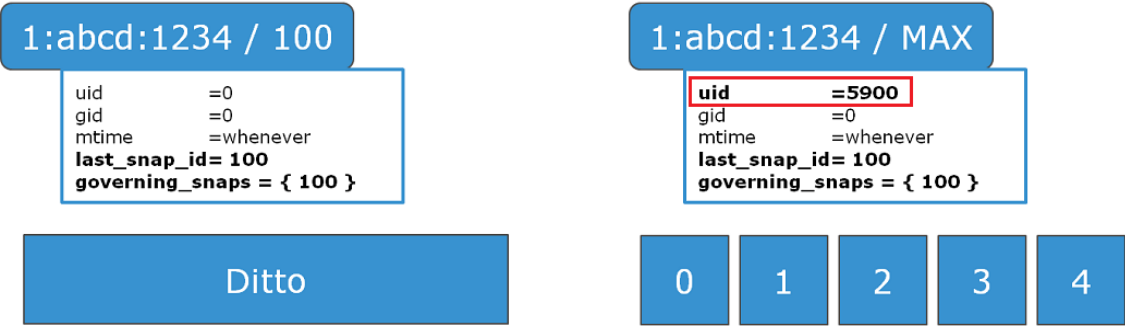


Figure 9: Changing the UID of a File in a Snapshot – Part 1

When a file’s metadata is changed, OneFS needs to create a copy of the data, and increment the Snap_ID. Since there is no actual data written, there’s no need to copy over any old data blocks to the new snapshot, so ditto-blocks are used instead. However, the first time a file is changed after a snapshot is taken, the LIN is recorded in the snapshot tracking file (STF). This STF lets OneFS efficiently know what data can be removed when a snapshot is deleted.

The file being written to is unchanged since the last snapshot. So the data blocks from the head version of the file are copied to the snapshot version, replacing the ditto-blocks, and the head version is overwritten with the new data. This is an example of copy on write (COW).

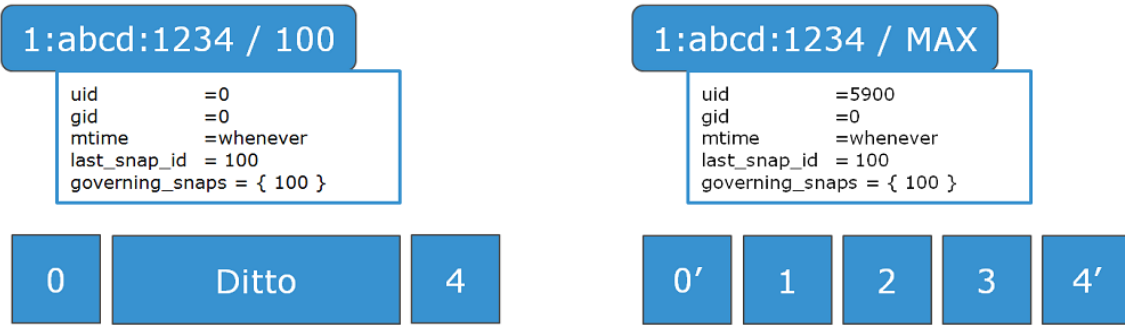


Figure 10: Changing the UID of a File in a Snapshot – Part 2

Deleting a Snapshot

When snapshots are manually deleted, OneFS marks the appropriate snapshot IDs and queue a job engine job to affect their removal. The SnapshotDelete job is queued immediately but the job engine will typically wait a minute or so to actually start running it. During this interval, the snapshot will be listed as 'delete pending'.

A similar procedure occurs with expired snapshots. Here, the snapshot daemon is responsible for checking expiration of snapshots and marking them for deletion. The daemon performs the check every 10-seconds. The job is queued to delete a snapshot completely and then it is up to the job engine to schedule it. The SnapshotDelete job might run immediately (after a min or so of wait) if the job engine determines that the job is runnable and there are no other contending jobs with higher priority running at that moment. For SnapshotDelete, it is only run if the cluster is in a fully available state, i.e., no drives/nodes are down.

The most efficient method for deleting multiple snapshots simultaneously is to process older through newer, and SnapshotIQ will automatically attempt to orchestrate deletes in this manner. A SnapshotDelete job engine schedule can also be defined so snapshot deletes only occur during desired times.

On deletion of a snapshot, OneFS immediately simply modifies some of the tracking data and the snapshot disappears from view. However, the actual behind-the-scenes clean up of the snapshot can involve a fair amount of work, which is performed in the SnapshotDelete job.

In the example below, snapshot ID 100 is being deleted. To accomplish this, any changes will likely need to be moved to the prior snapshot (ID 98), because that snapshot will no longer be able to read forward.

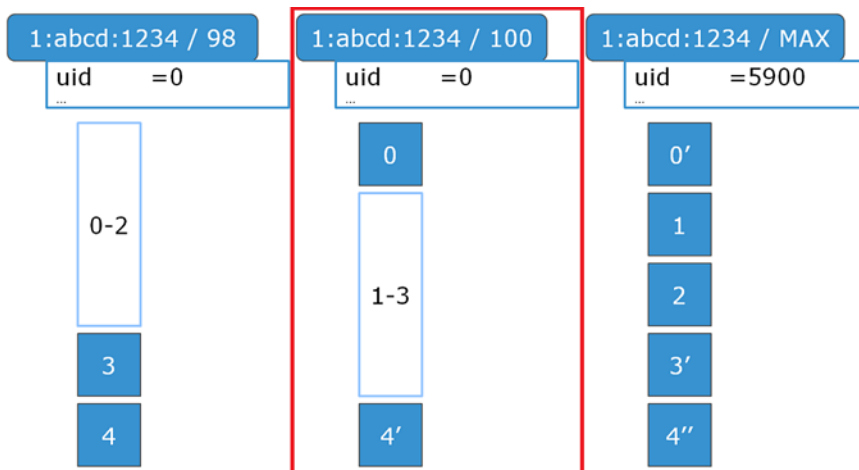


Figure 11: Deleting a Snapshot – Part 1

Snapshot 100 has two changed blocks: block 0 and block 4. Snapshot 98 was changed after snapshot 98 was taken, so block 4 can be deleted, but block 0 needs to be moved over to snapshot 98.

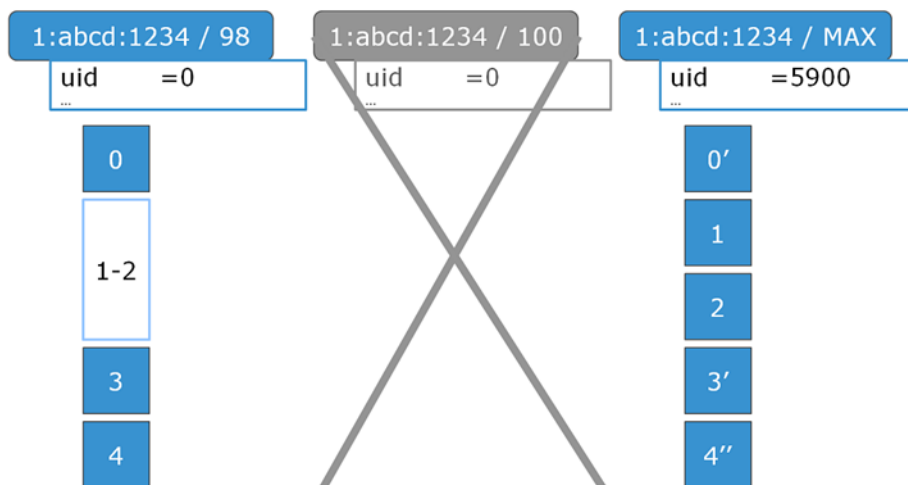


Figure 12: Deleting a Snapshot – Part 2

The oldest snapshot can be deleted very quickly. An ordered deletion is the deletion of the oldest snapshot of a directory, and is a recommended best practice for snapshot management. An unordered deletion is the removal of a snapshot that is not the oldest in a directory, and can often take approximately twice as long to complete and consume more cluster resources than ordered deletions.

Restoring a Snapshot

There are three main methods for restoring data from a snapshot:

- Copying specific files and directories directly from the snapshot
- Cloning a file from the snapshot
- Reverting the entire snapshot via the SnapRevert job

Copying a file from a snapshot duplicates that file, which roughly doubles the amount of storage space it consumes. Even if the original file is deleted from HEAD, the copy of that file will remain in the snapshot.

Cloning a file from a snapshot also duplicates that file. However, unlike a copy, a clone does not consume any additional space on the cluster, unless either the clone or original file is modified. File cloning is covered in more detail towards the end of this paper.

However, the most efficient of these approaches is the SnapRevert job, which automates the restoration of an entire snapshot to its top level directory. This is invaluable for quickly reverting to a previous, known-good recovery point, for example in the event of virus or malware outbreak. The SnapRevert job can be run from the Job Engine WebUI, and requires adding the desired snapshot ID.

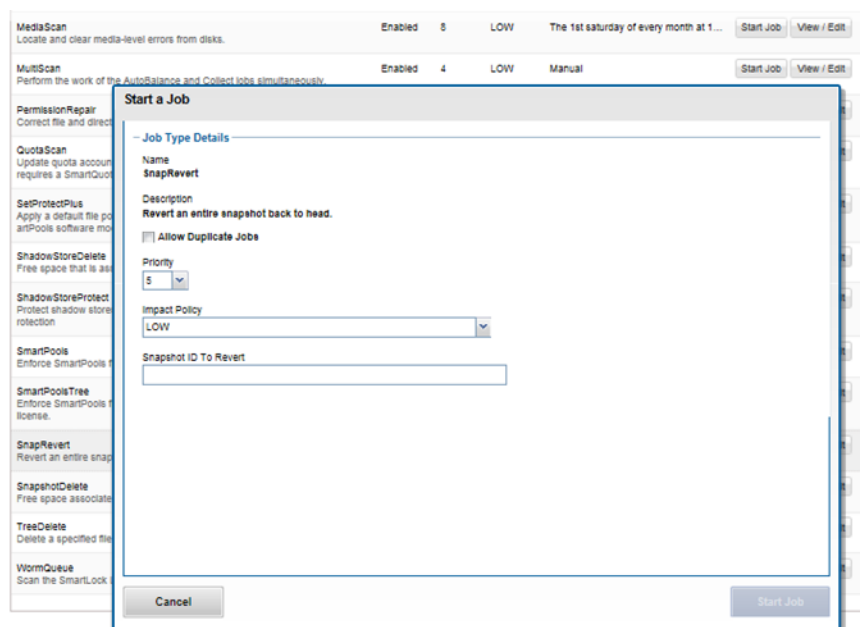


Figure 13: Running a Snapshot Revert Job

Before a snapshot is reverted, SnapshotIQ creates a snapshot of the data that is being replaced. This enables the snapshot revert to be undone later, if necessary.

Additionally, individual files, rather than entire snapshots, can also be restored in place using the `isi_file_revert` command line utility. This can help drastically simplify virtual machine management and recovery.

User Driven File Recovery

With the appropriate access credentials and permissions, NFS and SMB users can view and recover data from OneFS snapshots. The snapshots are accessed via the `.snapshot` directory, as described previously in this paper. The following screenshot from a Windows client shows the list of snapshots available on a OneFS SMB share:

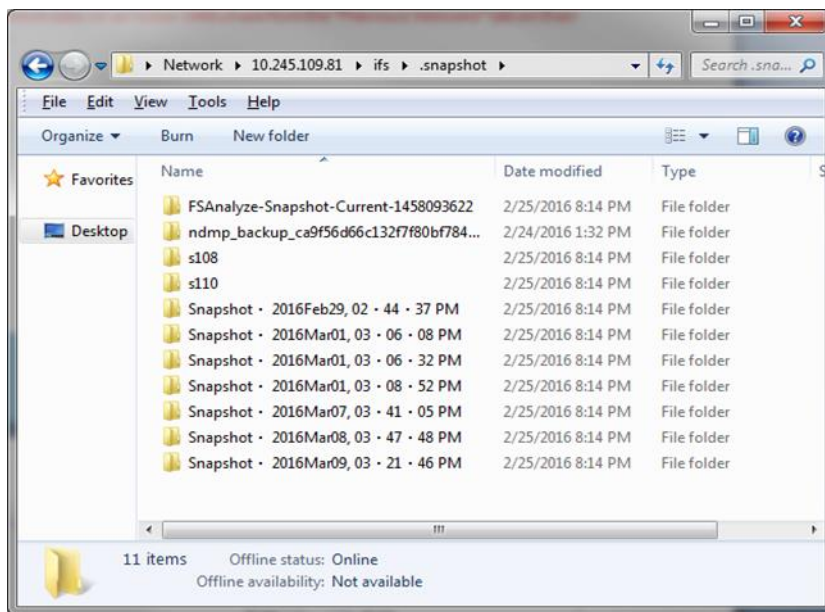


Figure 14: Accessing snapshots from the `.snapshot` portal

In the example below, a user accidentally deletes a file ``/ifs/data/foo/bar.txt`` at 9.10am and notices it's gone a couple of minutes later. By accessing the 9am snapshot, the user is able to recover the deleted file themselves at 9.14am, by copying it directly from the snapshot directory

``/ifs/data/foo/.snapshot./0900_snap/bar.txt`` back to its original location at

``/ifs/data/foo/bar.txt``.

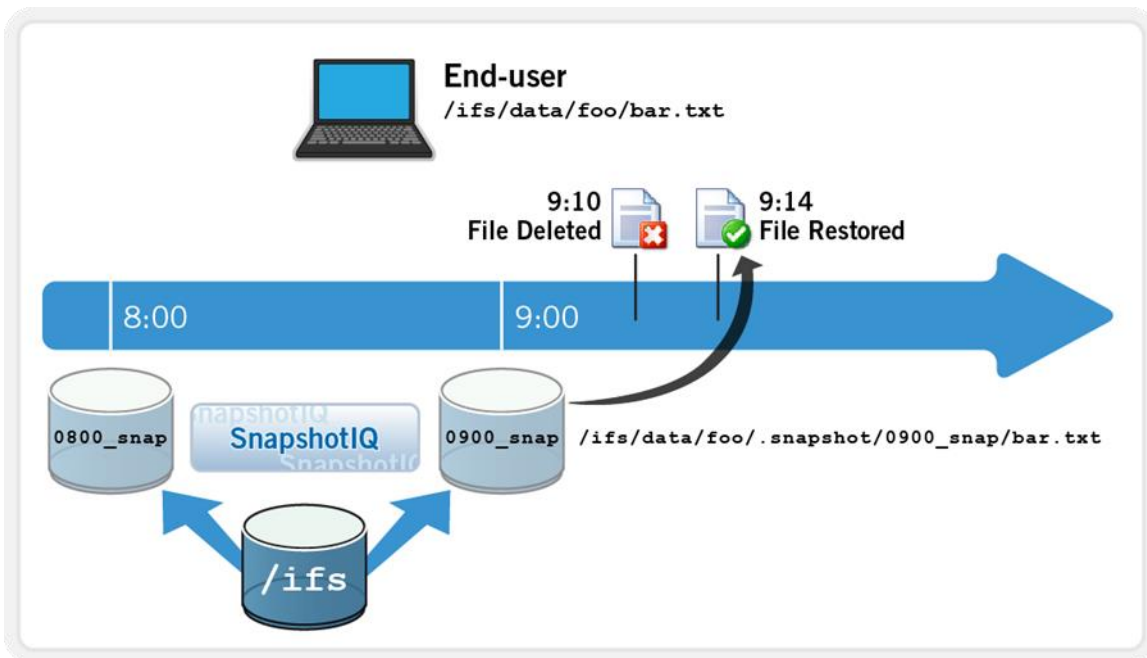


Figure 15: User driven file recovery with SnapshotIQ

SnapshotIQ integration with Windows Volume Snapshot Manager allows Windows users a simple way to restore data on an Isilon SMB share from the “Previous Versions” tab on their desktop.

Snapshot Aliases

As mentioned earlier, a snapshot can be referenced by a common name, in addition to a system defined snapshot ID. Snapshot names must be unique, and often incorporate a date, timestamp, and other context in their nomenclature.

Most users and applications are interested in referencing a snapshot relative to the HEAD version of data. Common naming schemes include, for example, “latest-weekly” or “latest-monthly,” for which the underlying snapshot ID changes as time goes by. Snapshot aliases provide for snapshot name indirection and enable the administrator to create user-friendly names, which can point at a different underlying snapshot over time.

Aliases are added and modified using the `--alias` option to the snapshot create and modify CLI subcommands, or via the snapshot schedule configuration pages of the graphical WebUI.

The screenshot shows the 'Snapshot Schedules' configuration page. The 'Create a Snapshot Schedule' form is visible. The 'Create an Alias' section is highlighted with a red box, showing the 'Yes' radio button selected and the 'Alias Name' field populated with 'Snapshot Schedule 132096734_Alias'.

Figure 16: Creating a Snapshot Alias

Only one alias per snapshot is allowed, which is enforced in the file system. Aliases also have an alias snapshot ID. This is provided as a convenient reference and to make interfaces consistent when referencing snapshots using IDs.

Snapshot scheduling

Snapshot scheduling allows cluster administrators to automatically generate snapshots according to a pre-defined itinerary. OneFS snapshot schedules can be configured at daily, weekly, monthly or yearly intervals, with single or multiple job frequency per schedule, and down to a per-minute granularity. OneFS snapshot schedules can be configured from either the CLI or the WebUI.

The screenshot shows the 'Create a Snapshot Schedule' web interface. It includes a 'Create a snapshot schedule' button at the top right. The form has several sections:

- Schedule Name:** A text field containing 'Snapshot Schedule 132096734'.
- Naming Pattern for Generated Snapshots:** A text field containing 'ScheduleName_duration_%Y-%m-%d-%H-%M'.
- Directory Path:** A text field containing '/ifs/data' and a 'Browse...' button.
- Snapshot Frequency:** A dropdown menu set to 'Weekly'.
- Generate snapshot every:** A text field containing '1' and a unit dropdown set to 'week(s)'.
- Generate snapshots on:** A grid of checkboxes for days of the week: Monday, Tuesday, Wednesday, Thursday, Friday, Saturday, and Sunday. The 'Sunday' checkbox is checked.
- Generate one snapshot per specified day:** A radio button that is selected.
- Generate snapshot at:** A dropdown menu set to '12:00 PM'.
- Generate multiple snapshots per specified day:** A radio button that is not selected.

Figure 17: Configuring a Weekly Snapshot Schedule

Similarly, automatic snapshot deletion can be configured per defined schedule at an hourly through yearly range.

The following table provides a recommended snapshot schedule for both ordered and unordered deletion configurations:

Deletion Type	Snapshot Frequency	Snapshot Time	Snapshot Expiration	Max Snaps Retained
Ordered deletion (for mostly static data)	Every hour	Start at 12:00 AM End at 11.59 AM	1 month	720
Unordered deletion (for frequently modified data)	Every other hour	Start at 12:00 AM, end at 11.59 PM	1 day	27
	Every day	At 12:00AM	1 week	
	Every week	Saturday at 12:00 AM	1 month	
	Every month	First Saturday of month at 12:00 AM	3 months	

Figure 18: Snapshot Schedule Recommendations

Note: A snapshot schedule cannot span multiple days. For example, you cannot specify to begin generating snapshots at 5:00 PM Monday and end at 5:00 AM Tuesday. To continuously generate snapshots for a period greater than a day, two individual snapshot schedules are required.

In order to generate snapshots from 5:00 PM Monday to 5:00 AM Tuesday, for example, create one schedule that generates snapshots from 5:00 PM to 11:59 PM on Monday, and another schedule that generates snapshots from 12:00 AM to 5:00 AM on Tuesday.

Configuring SnapshotIQ

The SnapshotIQ pages of the OneFS WebUI contains a 'settings' tab, as shown below. This provides a number of global snapshot settings, including:

- The ability to enable and disable the snapshot service
- Control of auto-creation and deletion of expired snapshots
- Per-protocol and complete control of snapshot visibility and accessibility

Summary Snapshots Snapshot Schedules **Settings**

File System Snapshot Settings

Snapshot Service: ☒ **On**

☒ **Auto-create Snapshots**
Create scheduled snapshots automatically

☒ **Auto-delete Snapshots**
Delete expired snapshots automatically

☐ **Off**
Turning snapshot service off will disable the creation of any snapshot and the automatic deletion

Global Visibility & Accessibility: ☒ **On**

NFS Visibility & Accessibility

☒ NFS Root Directory Accessible

☒ NFS Root Directory Visible

☒ NFS Sub-directories Accessible

SMB Visibility & Accessibility

☒ SMB Root Directory Accessible

☒ SMB Root Directory Visible

☒ SMB Sub-directories Accessible

Local Visibility & Accessibility

☒ Local Root Directory Accessible

☒ Local Root Directory Visible

☒ Local Sub-directories Accessible

☐ **Off**
Disabling visibility & accessibility will turn off any restoration from snapshots

Figure 19: Global Snapshot Configuration Settings

Snapshot Performance

Creating a snapshot in OneFS is a relatively instantaneous process. It typically takes substantially less than a second to perform, depending on the amount of snapshot prep work that has to happen first.

First, OneFS' write-back cache (coalescer) is paused and any uncommitted writes are flushed to allow the file system to be quiesced for the short period of time required to create the snapshot. Next, a marker is placed at the top-level directory inode for a particular snapshot and a unique snapshot ID is assigned.

Once this is done, the coalescer resumes and writes continue as normal, and any changes to HEAD (current version) are recorded in the snapshot inodes when any of the logical blocks they reference are altered - until another snapshot is taken.

Note that the moment a snapshot is created, it essentially consumes zero space until file creates, deletes, modifies and truncates start occurring in the structure underneath the marked top-level directory.

Snapshot Locks

As the name suggests, the exclusive snapshot lock which synchronizes the process of creating and deleting snapshots. If a snapshot has one or more locks applied to it, the snapshot cannot be deleted and is referred to as a locked snapshot. If the duration period of a locked snapshot expires, OneFS will not delete the snapshot until all locks on the snapshot have been deleted.

OneFS applies snapshot locks to ensure that snapshots generated by OneFS applications are not deleted prematurely. For this reason, it is recommended that you do not delete snapshot locks or modify the duration period of snapshot locks. A limited number of locks can be applied to a snapshot at a time. If you create snapshot locks, the limit for a snapshot might be reached, and OneFS could be unable to apply a snapshot lock when necessary. For this reason, it is recommended that you do not create snapshot locks.

The OneFS CLI command `'isi snapshot locks delete <snap_ID>'` allows you to remove existing snapshot locks, if necessary.

SnapshotIQ Monitoring and Reporting

Once a snapshot create request has completed, or has been terminated, a report is available. This can be accessed from the WebUI by navigating to **Data Protection > Snapshots > Saved Snapshots**, and selecting 'View Details' action button on the desired line item.

OneFS provides a variety of information about snapshots, including the total amount of space consumed by all snapshots. The following information is displayed in the Saved Snapshots area of the WebUI:

- **SnapshotIQ Status**
Indicates whether a SnapshotIQ license has been activated on the cluster.
- **Total Number of Saved Snapshots**
Indicates the total number of snapshots that exist on the cluster.
- **Total Number of Snapshots Pending Deletion**
Indicates the total number of snapshots that were deleted on the cluster since the last snapshot delete job was run. The space consumed by the deleted snapshots is not freed until the snapshot delete job is run again.
- **Total Number of Snapshot Aliases**
Indicates the total number of snapshot aliases that exist on the cluster.
- **Capacity Used by Saved Snapshots**
Indicates the total amount of space consumed by all snapshots.

Summary

Snapshots

Snapshot Schedules

Settings

Saved Snapshots

Manage and view your existing file system snapshots, and manually capture new snapshots as you need them.

SnapshotIQ Status: **Licensed & Active** [Versions/licensing info](#)

Total Number of Saved Snapshots: 10

Total Number of Snapshots Pending Deletion: 1

Total Number of Snapshot Aliases: 0

Capacity Used by Saved Snapshots: 6 GB

View snapshots | [View snapshot aliases](#) [+ Capture a new snapshot](#)

Created	Name / Path	Size	Expires	Select an action
2016-03-08 15:47:52	Snapshot: 2016Mar08, 03:47:48 PM Path: /ifs	14 MB	Never	Hide details Delete

Snapshot Details

Snapshot Name: **Snapshot: 2016Mar08, 03:47:48 PM** [Edit](#)

Snapshot ID: 136

Directory Path: /ifs

Size: 14 MB

Snapshot Creation Date: 2016-03-08 15:47:52

Snapshot Expiration: -- [Edit](#)

[Close](#) ✕

Figure 20: Viewing Snapshot Details

For SnapRevert and SnapshotDelete jobs, the Job Engine framework provides comprehensive run time and completion reporting for each individual job instance.

While a snapshot related job is underway, its status is available at a glance via the progress column in the active jobs table. Additional progress information is provided in an Active Job Details status update, which includes an estimated completion percentage based on the number of logical inodes (LINs) that have been counted and processed.

SnapshotIQ Licensing

SnapshotIQ is included as a core component of Dell EMC Isilon OneFS but requires a valid product license key in order to activate. This license key can be purchased through your Isilon account team. To create and manage snapshots, you must activate a SnapshotIQ license on the cluster. An unlicensed cluster will show a SnapshotIQ warning until a valid product license has been purchased and applied to the cluster.

License keys can be easily added via the ‘Activate License’ section of the OneFS WebUI, accessed by navigating via Cluster Management > Licensing.

Note: Some applications, such as SyncIQ and InsightIQ, must generate snapshots to function, but do not require an active SnapshotIQ license. By default, these snapshots are automatically deleted when OneFS no longer needs them. However, if you activate a SnapshotIQ license, you can retain these snapshots. Snapshots generated by other modules can still be viewed without a SnapshotIQ license.

Snapshot Reserve

There is also no requirement for reserved space for snapshots in OneFS. Snapshots can use as much or little of the available file system space as desirable and necessary.

A snapshot reserve can be configured if preferred, although this will be an accounting reservation rather than a hard limit, and is not a recommend best practice. If desired, snapshot reserve can be set via the OneFS command line interface (CLI) by running the ‘isi snapshot settings modify –reserve’ command.

For example, the following command will set the snapshot reserve to 20%:

```
# isi snapshot settings modify --reserve 20
```

Note: Snapshot reserve does not limit the amount of space that snapshots can consume on the cluster. Snapshots can consume a greater percentage of storage capacity specified by the snapshot reserve.

Additionally, when using Isilon SmartPools, snapshots can be stored on a different disk tier than the one the original data resides on.

Snapshot Overhead

The following guidelines can be used to help calculate snapshot usage overhead:

1. For a given directory, decide how frequently snapshots are taken, and calculate how much the primary data changes within the frequency interval
2. Decide how many snapshots are retained at any given time.

Overhead is the product of multiplying the rate of change (1) and number of snapshots to retain (2). For example, consider the following scenario:

- Snapshots under `/ifs/data/example` are taken every day and the data change rate is 5% daily.
- Three copies of the snapshots are retained at any given moment.

In this case, overhead is $5\% \times 3 \text{ (copies)} = 15\%$ for that directory only.

Within OneFS, file system snapshot overhead is very low and, in most cases, marginal. It can be complex to calculate and changes based on the number of files and file sizes in the given data set. However, even in the most extreme cases (such as many small files) the file system snapshot overhead is typically only 1-3%. For example, the file system snapshot overhead for a data set with 1 million files of 1K size, 5% daily change rate, and 3 snapshots retained, is 1%.

Sometimes, however, a cluster has many old snapshots that take up a lot of space. Reasons for this may include:

- The cluster is configured to take a lot of snapshots, but is not configured to have the snapshots expire as quickly as they are created.
- There is a device that is down or smartfailed on the cluster (in other words, the cluster is in a “degraded protection” state). You cannot run a SnapshotDelete job (or any job other than FlexProtect or FlexProtectLin) while the cluster is in this state. This could cause the number of snapshots to increase without being noticed.
- There is no SnapshotIQ license on the cluster, so you cannot delete old snapshots that still exist.

Snapshot Best Practices

For optimal cluster performance, Isilon recommends observing the following SnapshotIQ best practices. Please note that some of this information may be covered elsewhere in this paper.

- Configure the cluster to take fewer snapshots, and for the snapshots to expire more quickly, so that less space will be consumed by old snapshots. Take only as many snapshots as you need, and keep them active for only as long as you need them.
- Integration with Windows Volume Snapshot Manager allows Windows clients a method to restore from “Previous Versions”
- Snapshots are easily managed using flexible policies and schedules.
- Using SmartPools, snapshots can physically reside on a different disk tier than the original data.
- Recommend limiting snapshot creation to 1,024 per directory.
- The default snapshot limit is 20,000 per cluster.
- Avoid creating snapshots of directories that are already referenced by other snapshots. If you create a snapshot of a root directory, that snapshot counts towards the total number of snapshots for any subdirectories of the root directory. For example, if you create 500 snapshots of `/ifs/data` and 500 snapshots of `/ifs/data/media`, you have created 1000 snapshots of `/ifs/data/media`.
- It is recommended that you do not create more than 1000 hard links per file in a snapshot to avoid performance degradation.
- Always attempt to keep directory paths as shallow as possible. The deeper the depth of directories referenced by snapshots, the greater the performance degradation.
- If you have a lot of old, unneeded snapshots taking up space, and you do not have a SnapshotIQ license, contact Isilon Technical Support to discuss your options and provide assistance on deleting the old snapshots.
- Creating snapshots of directories higher on a directory tree will increase the amount of time it takes to modify the data referenced by the snapshot and require more cluster resources to manage the snapshot and the directory. However, creating snapshots of directories lower on directories trees will require more snapshot schedule, which can be difficult to manage.
- It is recommended that you do not create snapshots of `/ifs`. The main reason not to take a snapshot of `/ifs` is actually space: because we store `.ifsvar` under `/ifs`, a snapshot of `/ifs` will snapshot `/ifs/.ifsvar` as well. The files within `ifsvar` are generally highly protected (8x mirrored), and many of them are written quite often. This can waste a lot of space, especially on a small cluster.

- Avoid taking snapshots of `/ifs`, `/ifs/data`, and `/ifs/home` in favor of more specific targets when possible. In particular, avoid taking nested snapshots, redundant snapshots, or overly scoped snapshots. For example, if you schedule snapshots of `/ifs/data` and `/ifs/data/foo` and `/ifs/data/foo/bar`, consider taking snapshots of only the intermediate or most granularly scoped part (`/ifs/data/foo` or `/ifs/data/foo/bar`).
- The recommendation is not to disable the snapshot delete job, since this prevents unused disk space from being freed and can also cause performance degradation.
- If you need to delete snapshots and there are down or smartfailed components, or the cluster is in an otherwise degraded state, contact Isilon Technical Support for assistance.
- If the system clock is set to a time zone other than Coordinated Universal Time (UTC), SnapshotIQ modifies snapshot duration periods to match Daylight Savings Time (DST). Upon entering DST, snapshot durations are increased by an hour to adhere to DST; when exiting DST, snapshot durations are decreased by an hour to adhere to standard time.
- If you intend on reverting snapshots for a directory, it is recommended that you create SnapRevert domains for those directories while the directories are empty. Creating a domain for a directory that contains less data takes less time.
- Delete snapshots in order, beginning with the oldest. Where possible, avoid deleting snapshots from the middle of a time range. Newer snapshots are mostly pointers to older snapshots, and they look larger than they really are. Removing the newer snapshots will not free up much space. Deleting the oldest snapshot ensures you will actually free up the space. You can determine snapshot order (if not by name or date) by using the `isi snapshot list -l` command. The snapshot IDs (first column) are non-conserved, serial values.
- Create several snapshot schedules for a single directory, and then assign different snapshot duration periods for each schedule. Ensure that all snapshots are created at the same time when possible.
- Do not delete SyncIQ snapshots (snapshots with names that start with SIQ), unless the only remaining snapshots on the cluster are SyncIQ snapshots, and the only way to free up space is to delete those SyncIQ snapshots. Deleting SyncIQ snapshots resets the SyncIQ policy state, which requires a reset of the policy and potentially a full sync or initial diff sync. A full sync or initial diff sync could take many times longer than a regular snapshot-based incremental sync.

NOTE: If the oldest snapshot is a SyncIQ snapshot and you do not want to delete it, you can delete the next-oldest snapshot. Even if the next-oldest snapshot is in the same path, you will still free up some amount of data.

Snapshot Considerations

As discussed earlier, snapshots aren't free. There's always trade-off between cluster resource consumption (CPU, memory, disk), the potential for data fragmentation, and the benefit of increased data availability, protection, and recovery.

- Snapshots are created at the directory-level instead of the volume-level, thereby providing improved granularity.
- There is no requirement for reserved space for snapshots in OneFS. Snapshots can use as much or little of the available file system space as desirable.
- Quotas can be used to calculate a file and directory count that includes snapshot revisions, provided the quota is configured to include snaps in its accounting via the `--snaps=true` configuration option.
- The Lincount job will also include snapshot revisions of LINs in its count.
- Files with alternate data streams or resource forks are fully supported by SnapshotIQ.
- The SmartDedupe job will automatically ignore (and not deduplicate) file system snapshots.
- SnapshotDelete will only run if the cluster is in a fully available state, i.e., no drives or nodes are down.
- Snapshots of file clones, and shadow stores in general, are not allowed, since shadow stores have no hard links.
- Snapshot data will not be containerized by the OneFS Storage Efficiency for Healthcare PACS feature.
- A snapshot schedule cannot span multiple days: To generate snapshots from 5:00 PM Monday to 5:00 AM Tuesday, create one schedule that generates snapshots from 5:00 PM to 11:59 PM on Monday, and another schedule that generates snapshots from 12:00 AM to 5:00 AM on Tuesday.
- If a directory is moved, you cannot revert any snapshots of that directory which were taken prior to its move.

Snapshot and OneFS Feature Integration

The Isilon IQ Snapshots feature allows an administrator to create a frozen, point-in-time view of OneFS, while allowing normal file system modifications to continue without interruption. Efforts are made to ensure that Snapshots functionality consumes minimal system resources (disk space, CPU, etc.) while providing maximum flexibility to the system administrator and users.

As we have seen, Snapshots can be used on their own to provide functionality like user-initiated file restoration and staging of exported content, or in conjunction with other OneFS features such as Backup and SyncIQ to enhance the power and flexibility of those applications.

CloudPools and Snapshots

CloudPools, the OneFS cloud tiering product, is able to archive data sets that have associated snapshots. However, archiving snapshotted files to the cloud will not result in space savings on the cluster until all the snapshots taken prior to archiving have either expired or been deleted.

Also, as part of a DR process, the on-cluster stub files after a cloud archiving has completed can be snapshotted and replicated to another cluster.

Changelist Job and Snapshots

One of the classes of Job Engine jobs utilizes a 'changelist', rather than a full LIN-based scan, in order to discover its scope of work. The changelist approach analyzes two snapshots to find the LINs which changed (delta) between the snapshots, and from there determines the exact changes.

SyncIQ replication and the File System Analyze cluster analytics are good examples of a job that leverages snapshot deltas and the ChangelistCreate mechanism.

NDMP and Snapshots

The NDMP Snapshot Management Extension Interface leverages the extensibility of NDMP v4 to define a mechanism and protocol for controlling primary storage file system images commonly referred to as snapshots.

Specifically this interface supports the management of automated and manual snapshot creation, snapshot deletion, and general directory browsing as well as full snapshot recovery and selective file recovery. This interface provides functionality allowing snapshots to be used to implement near-line data protection solutions that offer faster backup and recovery times compared to traditional tape based secondary storage.

InsightIQ and Snapshots

The InsightIQ analytics engine provides statistics on the amount of space that snapshots consume on a cluster.

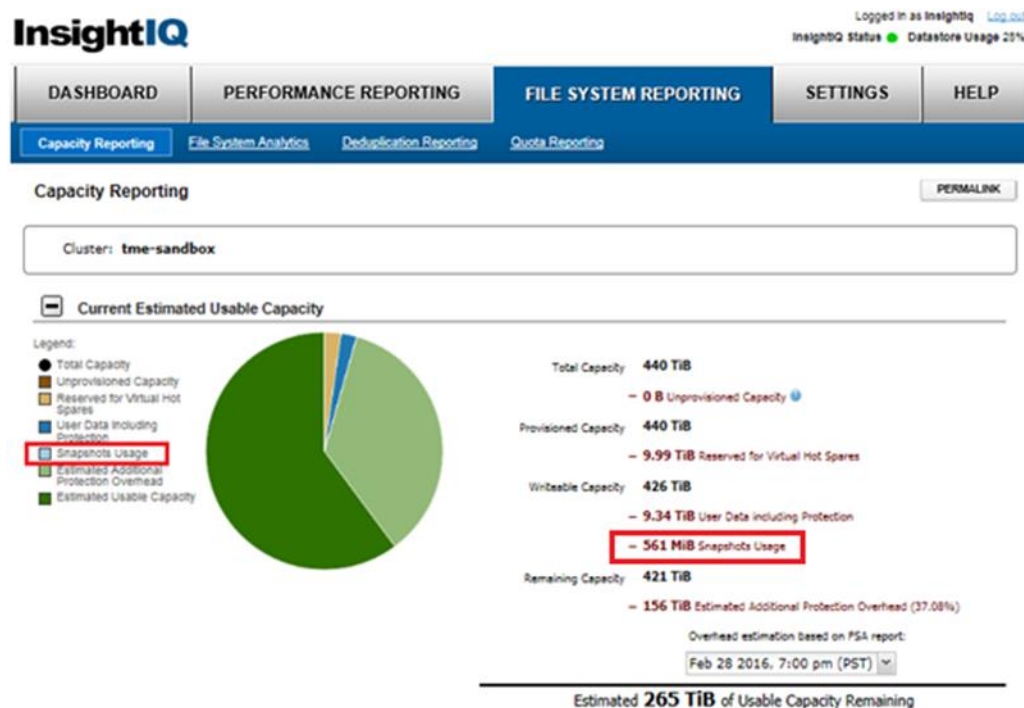


Figure 21: InsightIQ Snapshot Usage Data

SmartPools and Snapshots

When using Isilon SmartPools, snapshots can be stored on a different disk tier than the one the original data resides on. For example, the snapshots taken on a performance aligned tier can be physically housed on a more cost effective archive tier.

For example, the snapshots taken on a performance aligned tier can be physically housed on a more cost effective archive tier.

The snapshot storage target setting is applied to each file version by SmartPools. When a snapshot is taken, the storage pool setting is simply preserved. This means that the snapshot will initially be written to the default data pool and then moved. The SmartPools job subsequently finds the snapshot version and moves it to the intended pool during the next scheduled SmartPools job run.

Figure 22: Saving Snapshots to a different storage tier

SyncIQ Replication and Snapshots

SyncIQ also leverages snapshots for the consistency points required to facilitate replication, failover, and failback between Isilon clusters. This means that only the changes between the source and target datasets need to be replicated between the two clusters, allowing for efficient replication and granular recovery objectives. The snapshots generated by SyncIQ can also be used for archival purposes on the target cluster.

Source cluster snapshots

SyncIQ creates snapshots on the source cluster to ensure that a consistent point-in-time image is replicated and that unaltered data is not sent to the target cluster. Before running a replication job, SyncIQ takes a snapshot of the source directory. It then replicates data according to the snapshot rather than the current state of the cluster, allowing users to modify source-directory files while ensuring that an exact point-in-time image of the source directory is replicated.

For example, if a replication job of `/ifs/home/user1/` starts at 10:00 PM and finishes at 10:20 PM, and `/ifs/home/user1/file` is modified at 10:10 PM, the modifications are not reflected on the target cluster, even if `/ifs/home/user1/file` is not replicated until 10:15 PM.

SyncIQ can also replicate data according to either an on-demand or scheduled snapshot generated directly by SnapshotIQ. If data is replicated using a SnapshotIQ snapshot, SyncIQ does not generate another snapshot of the source directory. This method can be useful if you want to replicate identical copies of data to multiple Isilon clusters.

SyncIQ generates source snapshots to ensure that replication jobs do not transfer unmodified data. When a job is created for a replication policy, SyncIQ checks whether it is the first job created for the policy. If not, SyncIQ compares the snapshot generated for the earlier job with the snapshot generated for the new job.

SyncIQ replicates only data that has changed since the last time a snapshot was generated for the replication policy. When a replication job is completed, SyncIQ deletes the previous source-cluster snapshot and retains the most recent snapshot until the next job is run.

Target cluster snapshots

When a replication job is run, SyncIQ generates a snapshot on the target cluster to facilitate failover operations. When the next replication job is created for the replication policy, the job creates a new snapshot and deletes the old one.

If a SnapshotIQ license has been activated on the target cluster, you can configure a replication policy to generate additional snapshots that remain on the target cluster even as subsequent replication jobs run.

SyncIQ generates target snapshots to enable failover on the target cluster regardless of whether a SnapshotIQ license has been configured on the target cluster. Failover snapshots are generated when a replication job completes. SyncIQ retains only one failover snapshot per replication policy, and deletes the old snapshot after the new snapshot is created.

If a SnapshotIQ license has been activated on the target cluster, you can configure SyncIQ to generate archival snapshots on the target cluster that are not automatically deleted when subsequent replication jobs run. Archival snapshots contain the same data as the snapshots that are generated for failover purposes. However, you can configure how long archival snapshots are retained on the target cluster. You can access archival snapshots the same way that you access other snapshots generated on a cluster.

File clones

In complement to SnapshotIQ's read-only snapshots, OneFS also provides the ability to create writable clones of files. OneFS File Clones provides a rapid, efficient method for provisioning multiple writable copies of files. Common blocks are shared between the original file and clone, providing space efficiency and offering similar performance and protection levels across both. This mechanism is ideal for the rapid provisioning and protection of virtual machine files and is integrated with VMware's linked cloning and block and file storage APIs. This utilizes the OneFS shadow store metadata structure, which is able to reference physical blocks, references to physical blocks, and nested references to physical blocks. This is in contrast to snapshots which, as we have seen, have ditto records which reference other versions.

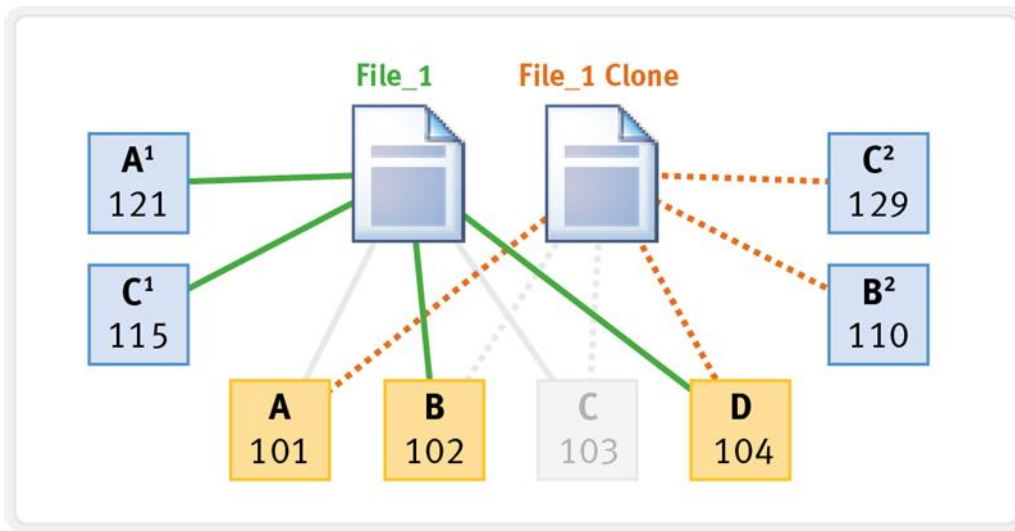


Figure 23: File Clones

Shadow stores are hidden files, which behave differently than other files. For example:

- Reading shadow-store references might be slower than reading data directly. Specifically, reading non-cached shadow-store references is slower than reading non-cached data. Reading cached shadow-store references takes no more time than reading cached data.
- When files that reference shadow stores are replicated to another Isilon cluster or backed up to a Network Data Management Protocol (NDMP) backup device, the shadow stores are not transferred to the target Isilon cluster or backup device. The files are transferred as if they contained the data that they reference from shadow stores. On the target Isilon cluster or backup device, the files consume the same amount of space as if they had not referenced shadow stores.
- When OneFS creates a shadow store, OneFS assigns the shadow store to a storage pool of a file that references the shadow store. If you delete the storage pool that a shadow store resides on, the shadow store is moved to a pool occupied by another file that references the shadow store.
- OneFS does not delete a shadow-store block immediately after the last reference to the block is deleted. Instead, OneFS waits until the ShadowStoreDelete job is run to delete the unreferenced block. If a large number of unreferenced blocks exist on the cluster, OneFS might report a negative deduplication savings until the ShadowStoreDelete job is run.
- Shadow stores are protected at least as much as the most protected file that references it. For example, if one file that references a shadow store resides in a storage pool with +2 protection and another file that references the shadow store resides in a storage pool with +3 protection, the shadow store is protected at +3.
- Quotas account for files that reference shadow stores as if the files contained the data referenced from shadow stores; from the perspective of a quota, shadow-store references do not exist. However, if a quota includes data protection overhead, the quota does not account for the data protection overhead of shadow stores.

Note: Clones cannot contain alternate data streams (ADS). If you clone a file that contains alternate data streams, the clone will not contain the alternate data streams.

Conclusion

SnapshotIQ is designed to help you provide highly efficient, cost-effective, low recovery objective data protection. Once a baseline snapshot has been established, only changes to blocks that make up a file are reflected in updates to the current version of snapshots. This allows highly efficient snapshot storage utilization. Additionally, since snapshots are an integral part of the OneFS filesystem, there is no need to pre-allocate dedicated snapshot reserve space.

One of the largest IT costs associated with backup and restore is the sheer number of help desk calls from end-users who accidentally delete a file or directory. To reduce these costs, SnapshotIQ can be used to empower end-users by enabling them to easily find and restore their own accidentally deleted data – without any IT intervention. In this way, business users remain productive and impact on IT staff can be significantly reduced.

SnapshotIQ transcends the limits of traditional approaches by reliably distributing a highly scalable number of snapshots across multiple, highly available Isilon scale-out storage nodes. The result is a remarkably simple, scalable and efficient data backup and recovery capability that allows you to meet the demanding data availability requirements of your business.

TAKE THE NEXT STEP

Contact your Dell EMC sales representative or authorized reseller to learn more about how Isilon scale-out NAS storage solutions can benefit your organization.

[Shop Dell EMC Isilon](#) to compare features and get more information.



Dell, EMC and other trademarks are trademarks of Dell Inc. or its subsidiaries.

© Copyright 2017 Dell Inc. All rights reserved. Published in the USA. 4/17 White Paper H15048.2.

Dell EMC believes the information in this document is accurate as of its publication date. The information is subject to change without notice.