

HIGH AVAILABILITY & DATA PROTECTION WITH EMC ISILON SCALE-OUT NAS

Abstract

This white paper gives a detailed look at the challenges organizations face as they deal with the deluge of digital content and unstructured data and the growing importance of data protection. It details how the EMC Isilon OneFS architecture provides high availability and data protection needed to meet these challenges.

December 2012

Copyright © 2012 EMC Corporation. All Rights Reserved.

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

The information in this publication is provided "as is." EMC Corporation makes no representations or warranties of any kind with respect to the information in this publication, and specifically disclaims implied warranties of merchantability or fitness for a particular purpose.

Use, copying, and distribution of any EMC software described in this publication requires an applicable software license.

For the most up-to-date listing of EMC product names, see EMC Corporation Trademarks on EMC.com.

All other trademarks used herein are the property of their respective owners.

Part Number **H10588.2**

Table of Contents

Introduction	5
High availability and reliability at the petabyte scale	5
Areal density and drive rebuild times	5
Silent data corruption	6
Data protection continuum	6
High availability with EMC Isilon	7
Isilon scale-out architecture	7
OneFS architectural overview	8
Safe writes.....	9
Cluster group management	9
Concurrency and locking	10
File layout	10
Flexible protection	11
Failure domains and resource pools.....	12
Automatic partitioning	12
Virtual hot spares	13
OneFS fault tolerance	13
File system journal.....	13
Proactive device failure.....	14
Isilon data integrity.....	14
Protocol checksums.....	14
Dynamic sector repair	14
MediaScan	14
IntegrityScan	15
Fault isolation.....	15
Accelerated drive rebuilds	15
Isilon data protection.....	16
High availability and data protection strategies.....	16
The Isilon high availability and data protection suite.....	17
Connection load balancing and failover	18
SmartConnect	18
Snapshots	19
SnapshotIQ.....	19
SnapshotIQ architecture	20
Snapshot scheduling	21
Snapshot deletes	21
Snapshot restore	22
File clones.....	22
Replication	23

SyncIQ	23
SyncIQ linear restore	24
SyncIQ replica protection	25
SyncIQ failover and failback	25
Archive.....	26
SmartLock	26
Nearline, VTL and tape backup	27
Backup Accelerator	27
Backup from snapshots	27
Parallel streams	28
NDMP	28
Direct NDMP model	29
Remote NDMP model	30
Incremental backups	31
Direct access recovery	31
Directory DAR	31
Certified backup applications.....	31
Summary.....	32
Isilon acronyms glossary	33
About EMC Isilon	33
Contact Isilon.....	34

Introduction

Today, organizations of all sizes across the full spectrum of the business arena are facing a similar problem: An explosion in the sheer quantity of file-based data they are generating and, by virtue, are forced to manage. This proliferation of unstructured data, often dubbed 'big 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.

High availability and reliability at the petabyte scale

Once data sets grow into the hundreds of terabytes on up, a whole new level of availability, management and protection challenges arise. At this magnitude, given the law of large numbers with the sheer quantity of components involved, there will almost always be one or more components in a degraded state at any point in time within the storage infrastructure. As such, guarding against single points of failure and bottlenecks becomes a critical and highly complex issue. Other challenges that quickly become apparent at the petabyte scale include the following:

- **File System Limitations**
 - How much capacity and how many files can a file system accommodate?
- **Disaster recovery**
 - How do you duplicate the data off site and then how do you retrieve it?
- **Scalability of Tools**
 - How do you take snapshots of massive data sets?
- **Software Upgrades and Hardware refresh**
 - How do you upgrade software and replace outdated hardware with new?
- **Performance Issues**
 - How long will searches & treewalks take with large, complex datasets?
- **Backup and Restore**
 - How do you back up a large dataset and how long will it take to restore?

Given these challenges, the requirement for a new approach to file storage is clear. Fortunately, when done correctly, scale-out NAS can fulfill this need.

Areal density and drive rebuild times

In today's world of large capacity disk drives, the probability that secondary device failures will occur has increased dramatically. Areal density, the amount of written information on the disk's surface in bits per square inch, continues to outstrip Moore's

law. However, the reliability and performance of disk drives are not increasing at the same pace, and this is compounded by the growing amount of time it takes to rebuild drives.

Large capacity disks, such as the current three and four terabyte SATA drives, require much longer drive reconstruction times, since each subsequent generation of disk still has the same number of heads and actuators servicing increased density platters—currently up to one terabyte per platter and with an areal density of 635Gb/inch. This significantly raises the probability of a multiple drive failure scenario.

Silent data corruption

Another threat that needs to be addressed, particularly at scale, is the looming specter of hardware induced corruption. For example, when CERN tested the data integrity of standard disk drives they discovered some alarming findings. To do this, they built a simple write and verify application which they ran across a pool of three thousand servers, each with a hardware RAID controller. After five weeks of testing, they found in excess of five hundred instances of silent data corruption spread across seventeen percent of the nodes - after having previously thought everything was fine. Under the hood, the hardware RAID controller only detected a handful of the most blatant data errors and the rest passed unnoticed.

Suffice to say, this illustrates two inherent data protection requirements: First, the need for an effective, end-to-end data verification process to be integral to a storage device in order to detect and mitigate such instances of silent data corruption. Second, the requirement for regular and reliable backups as the linchpin of a well-founded data protection plan.

Data protection continuum

The availability and protection of data can be usefully illustrated in terms of a continuum:

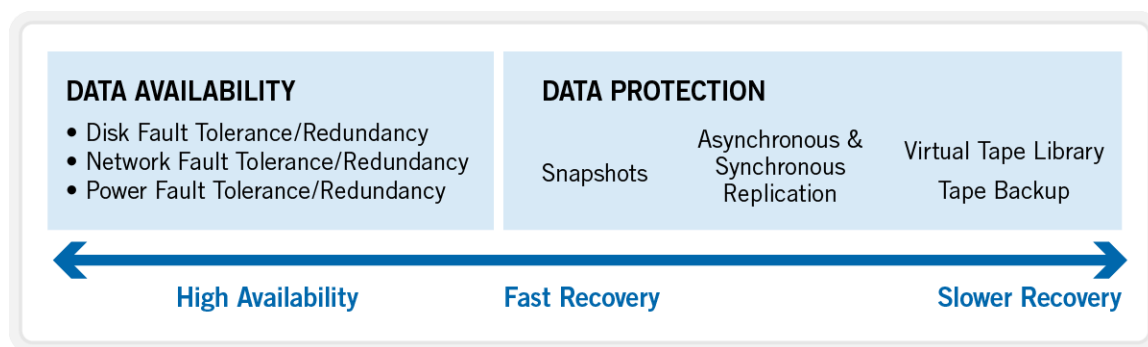


Figure 1: Data Protection Continuum

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 a form of point-in-time snapshots for fast recovery, followed by synchronous and asynchronous replication. Finally, backup to tape or a virtual tape library sits at the end of the continuum, providing insurance against large scale data loss, natural disasters and other catastrophic events.

High availability with EMC Isilon

As we will see, EMC® Isilon® OneFS® takes a holistic approach to ensuring that data is consistent and intact - both within the file system, and when exiting the cluster via a network interface. Furthermore, the Isilon clustering technology is uncompromisingly designed to simplify the management and protection of multi-petabyte datasets.

Isilon scale-out architecture

An Isilon cluster is built on a highly redundant and scalable architecture, based upon the hardware premise of shared nothing. The fundamental building blocks are platform nodes, of which there are anywhere from three to one hundred and forty four nodes in a cluster. Each of these platform nodes contain CPU, memory, disk and I/O controllers in an efficient 2U or 4U rack-mountable chassis. Redundant Infiniband (IB) adapters provide a high speed back-end cluster interconnect—essentially a distributed system bus - and each node houses a fast, battery-backed file system journal device. With the exception of the IB controller, journal card and an LCD control front panel, all of a node's components are standard enterprise commodity hardware.

These Isilon nodes contain a variety of storage media types and densities, including SAS & SATA hard disk drives (HDDs), solid-state drives (SSDs), and a configurable quantity of memory. This allows customers to granularly select an appropriate price, performance and protection point to accommodate the requirements of specific workflows or storage tiers.

Highly available storage client access is provided via multiple 1 or 10Gb/s Ethernet interface controllers within each node, and across a variety of file- and block-based protocols including NFS, SMB/CIFS and iSCSI.

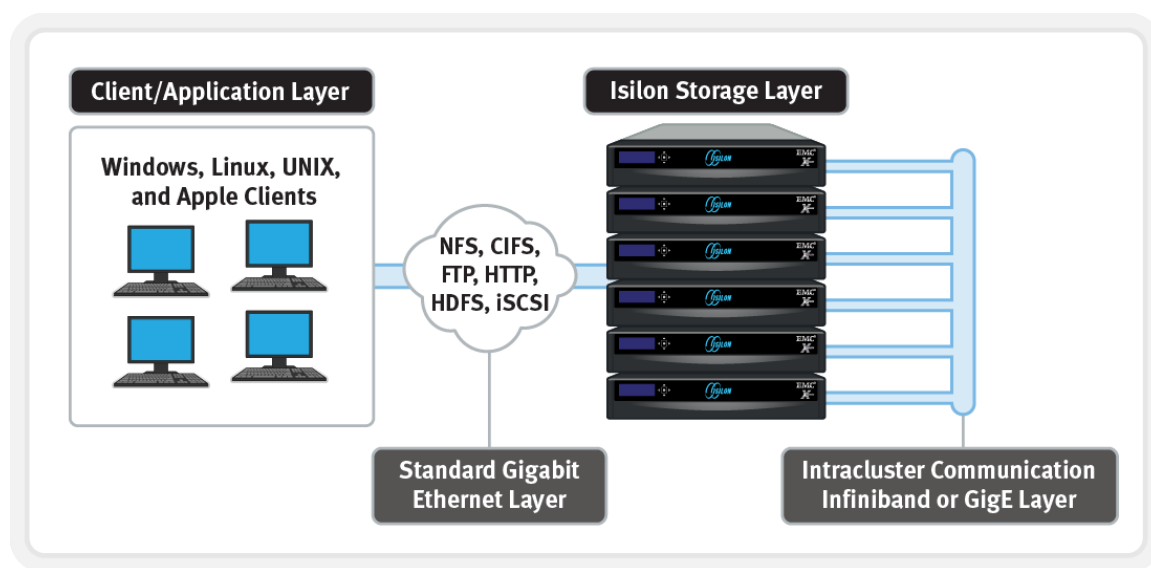


Figure 2: Isilon Scale-out NAS Architecture

OneFS architectural overview

OneFS collapses the traditional elements of the storage stack—data protection, volume manager, file system, etc.—into a single, unified software layer (see Figure 3 below). This allows for a highly extensible file system that affords unparalleled levels of protection and availability.

Built atop FreeBSD's UNIX implementation, availability and resilience are integral to OneFS from the lowest level on up. For example, unlike BSD, OneFS provides mirrored volumes for the root and /var file systems via the Isilon Mirrored Device Driver (IMDD), stored on flash drives. OneFS also automatically saves last known good boot partitions for further resilience.

On the network side, the Isilon logical network interface (LNI) framework provides a robust, dynamic abstraction for easily combining and managing differing interfaces, enabling network resilience. Multiple network interfaces can be trunked together with Link Aggregation Control Protocol (LACP) and Link Aggregation and Link Failover (LAGG) to provide bandwidth aggregation in addition to client session failover and general network resilience.

Within the cluster, every disk within each node is assigned both a Globally Unique Identifier (GUID) and logical drive number and is subdivided into 32MB cylinder groups comprised of 8KB blocks. Each cylinder group is responsible for tracking, via a bitmap, whether its blocks are used for data, inodes or other metadata constructs. The combination of node number, logical drive number and block offset comprise a block or inode address and fall under the control of the aptly named Block Allocation Manager (BAM).

In addition to block and inode allocation, the BAM also handles file layout and locking and abstracts the details of OneFS distributed file system from the kernel and userspace. The BAM never actually touches the disk itself, instead delegating tasks to the local and remote block manager elements respectively on the appropriate nodes. The Remote Block Manager (RBM) is essentially a Remote Procedure Call (RPC) protocol that utilizes the Socket Direct Protocol (SDP) over redundant Infiniband for

reliable, ultra low-latency back-end cluster communication. These RBM messages—everything from cluster heartbeat pings to distributed locking control - are then processed by a node's Local Block Manager via the Device Worker Thread (DWT) framework code.

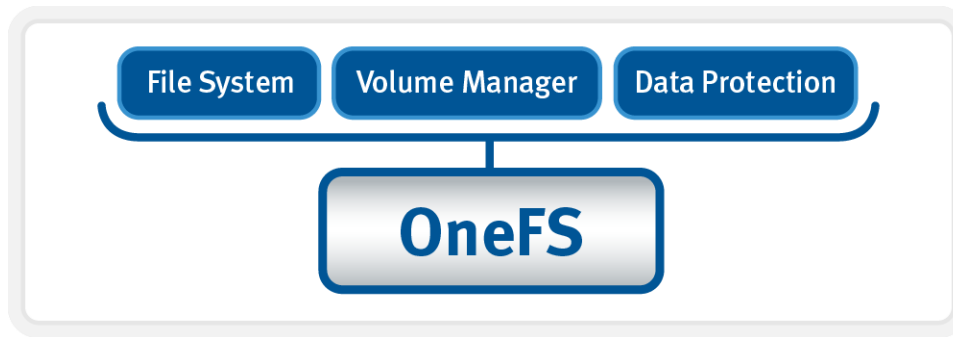


Figure 3: OneFS Collapsed Stack Storage Architecture

Safe writes

For write operations, where coherency is vital, the BAM first sets up a transaction. Next it uses a 2-phase commit protocol (2PC) over the RBM to guarantee the success of an atomic write operation across all participant nodes. This is managed via the BAM Safe Write (BSW) code path. The 2PC atomically updates multiple disks across the 2PC participant nodes, using their NVRAM journals for transaction logging. The write path operates as follows:

1. Client performs a transactional write.
 - Block is written to Non-Volatile Random Access Memory (NVRAM) journal; memory buffer is pinned.
 - Rollback data is maintained.
2. Transaction commits.
 - NVRAM data is pinned; memory buffer is dirty.
 - Rollback data can now be discarded.
 - Top level operation is complete.
3. OneFS asynchronously flushes dirty buffers to disk at some point.
 - Placed into the writeback cache.
 - NVRAM data still required and memory buffer discarded.
4. Journal approaches full or timeout and issues disk writeback cache flush.
 - This occurs relatively infrequently.
5. Cache flush complete.
 - NVRAM data discarded for writes that were returned prior to flush.

Cluster group management

Cluster coherence and quorum is handled by OneFS Group Management Protocol (GMP). The challenge is combining the various elements—performance, coherency,

client access protocols - across multiple heads. The GMP is built on several distributed algorithms and strictly adheres to Brewer's Theorem, which states that it is impossible for a distributed computer system to simultaneously guarantee all three of the following; consistency, availability and partition tolerance. OneFS does not compromise on either consistency or availability.

Given this, a quorum group comprising more than half of a cluster's nodes must be active and responding at any given time. In the event that a node is up and responsive but not a member of the quorum group, it is forced into a read-only state.

OneFS employs this notion of a quorum to prevent "split-brain" conditions that might possibly result from a temporary cluster division. The quorum also dictates the minimum number of nodes required to support a given data protection level. For example, seven or more nodes are needed for a cluster to support an N+3 configuration. This allows for a simultaneous loss of three nodes while still maintaining a quorum of four nodes, allowing the cluster to remain operational.

The group management protocol keeps track of the state of all the nodes and drives that are considered part of the cluster. Whenever devices are added or removed from the cluster, either proactively or reactively, a group change is broadcast, the group ID is incremented and any uncommitted journal write transactions are resolved.

Concurrency and locking

OneFS employs a distributed lock manager that utilizes a proprietary hashing algorithm to orchestrate coherent locking on data across all nodes in a storage cluster. The design is such that a lock coordinator invariably ends up on a different node than the initiator and either shared or exclusive locks are granted as required. The same distributed lock manager mechanism is used to orchestrate file system structure locks as well as protocol and advisory locks across the entire cluster. OneFS also provides support for delegated locks (i.e. SMB opportunistic locks and NFSv4 delegations) and also byte-range locks.

File layout

OneFS is a single file system providing one vast, scalable namespace—free from multiple volume concatenations or single points of failure. As such, all nodes access the same structures across the cluster using the same block addresses and all directories are inode number links emanating from the root inode.

The way data is laid out across the nodes and their respective disks in a cluster is fundamental to OneFS functionality. As mentioned previously, OneFS uses an 8KB block size, and sixteen of these blocks are combined to create a 128KB stripe unit. Files are striped across nodes allowing files to use the resources (spindles and cache) of up to twenty nodes, based on per-file policies.

The layout decisions are made by the BAM on the node that initiated a particular write operation using the 2PC described above. The BAM Safe Write (BSW) code takes the cluster group information from GMP and the desired protection policy for the file and makes an informed decision on where best to write the data blocks to ensure the file is properly protected. To do this, the BSW generates a write plan, which comprises all the steps required to safely write the new data blocks across the protection group. Once complete, the BSW will then execute this write plan and guarantee its successful completion.

All files, inodes and other metadata structures (B-trees, etc) within OneFS are either mirrored up to eight times or parity protected, with the data spread across the various disk cylinder groups of multiple nodes. Parity protection uses an N+M scheme with N representing the number of nodes—the stripe width—and M the number of parity blocks. This is described in more detail within the 'Flexible Protection' chapter below.

OneFS will not write files at less than the desired protection level, although the BAM will attempt to use an equivalent mirrored layout if there is an insufficient stripe width to support a particular forward error correction (FEC) protection level.

Flexible protection

OneFS is designed to withstand multiple simultaneous component failures (currently four) while still affording unfettered access to the entire file system and dataset. Data protection is implemented at the file system level and, as such, is not dependent on any hardware RAID controllers. This provides many benefits, including the ability add new data protection schemes as market conditions or hardware attributes and characteristics evolve. Since protection is applied at the file-level, a OneFS software upgrade is all that's required in order to make new protection and performance schemes available.

OneFS employs the popular Reed-Solomon erasure coding algorithm for its parity protection calculations. Protection is applied at the file-level, enabling the cluster to recover data quickly and efficiently. Inodes, directories and other metadata are protected at the same or higher level as the data blocks they reference. Since all data, metadata and FEC blocks are striped across multiple nodes, there is no requirement for dedicated parity drives. This both guards against single points of failure and bottlenecks and allows file reconstruction to be a highly parallelized process. Today, OneFS provides N+1 through N+4 parity protection levels, providing protection against up to four simultaneous component failures respectively. A single failure can be as little as an individual disk or, at the other end of the spectrum, an entire node.

OneFS also supports several hybrid protection schemes. These include N+2:1 and N+3:1, which protect against two drive failures or one node failure, and three drive failures or one node failure, respectively. These protection schemes are particularly useful for high density node configurations, where each node contains up to thirty six, multi-terabyte SATA drives. Here, the probability of multiple drives failing far surpasses that of an entire node failure. In the unlikely event that multiple devices have simultaneously failed, such that the file is "beyond its protection level", OneFS will re-protect everything possible and report errors on the individual files affected to the cluster's logs.

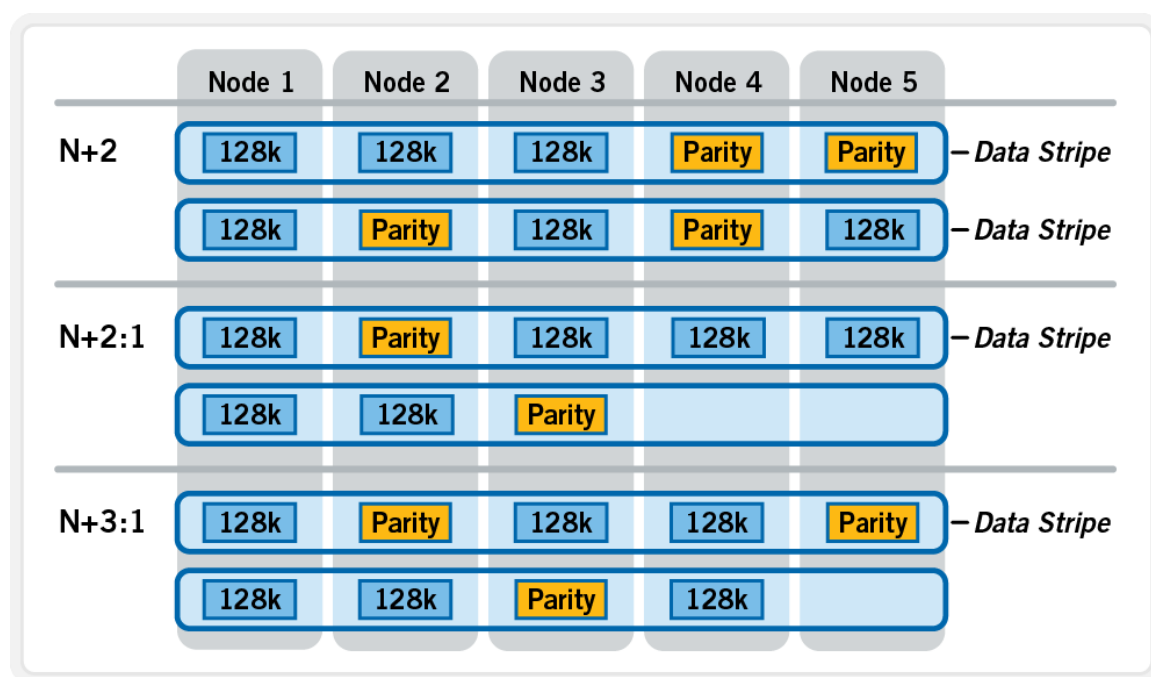


Figure 4: OneFS Hybrid Parity Protection Schemes (N+M:x)

As mentioned earlier, OneFS also provides a variety of mirroring options ranging from 2x to 8x, allowing from two to eight mirrors of the specified content. Metadata, for example, is mirrored at one level above FEC by default. For example, if a file is protected at N+1, its associated metadata object will be 3x mirrored.

Striped, distributed metadata coupled with continuous auto-balancing affords OneFS truly linear performance characteristics, regardless of fullness of file system. Both metadata and file data are spread across the entire cluster keeping the cluster balanced at all times.

Failure domains and resource pools

Data tiering and management in OneFS is handled by Isilon SmartPools™ software. From a data protection point of view, SmartPools facilitates the subdivision of large numbers of high-capacity, homogeneous nodes into smaller, more Mean Time to Data Loss (MTTDL)-friendly disk pools. For example, an 80-node nearline cluster would typically run at N+4 protection level. However, partitioning it into four, twenty node disk pools would allow each pool to run at N+2, thereby lowering the protection overhead and improving data utilization without any net increase in management overhead.

Automatic partitioning

In keeping with the goal of storage management simplicity, OneFS will automatically calculate and partition the cluster into pools of disks or 'node pools' which are optimized for both MTTDL and efficient space utilization. This means that protection level decisions, such as the 80-node cluster example above, are not left to the customer—unless desired.

With Automatic Provisioning, every set of equivalent node hardware is automatically divided into node pools comprising up to twenty nodes and six drives per node. These node pools are protected by default at N+2:1, and multiple pools can then be combined into logical tiers and managed using Isilon SmartPools file pool policies. By subdividing a node's disks into multiple, separately protected pools, nodes are significantly more resilient to multiple disk failures than previously possible.

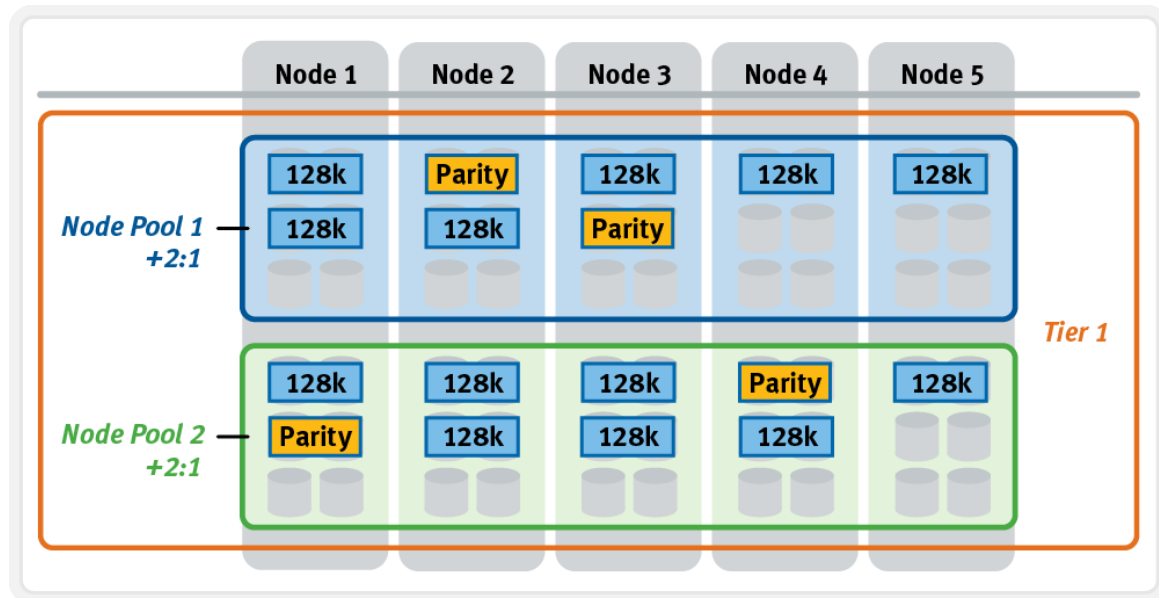


Figure 5: SmartPools Automatic Provisioning

Virtual hot spares

SmartPools also provides a virtual hot spare option, if desired. This functionality allows space to be reserved in a disk pool, equivalent to up to four full drives. This virtual hot spare pool can be immediately utilized for data re-protection in the event of a drive failure.

From a data availability and management point of view, SmartPools also applies storage tiering concepts to disk pools, allowing the storage and movement of data according to rich file policies or attributes. As such, SmartPools facilitates the automated alignment of data with the appropriate class of storage according to its business value, performance profile, and availability requirements. An Isilon cluster can thereby provide multiple storage pools, each supporting a range of availability SLAs within a single, highly scalable and easily managed file system. This resource pool model aligns beautifully with the current IT trend of private and hybrid cloud initiatives.

OneFS fault tolerance

File system journal

Every Isilon node is equipped with a dual-battery backed 512MB NVRAM card, which guards that node's file system journal. Each journal is used by OneFS as stable storage, and guards write transactions against sudden power loss or other

catastrophic events. The journal protects the consistency of the file system and the battery charge lasts up to three days. Since each member node of an Isilon cluster contains an NVRAM controller, the entire OneFS file system is therefore fully journaled.

Proactive device failure

OneFS will proactively remove, or SmartFail, any drive that reaches a particular threshold of detected Error Correction Code (ECC) errors, and automatically reconstruct the data from that drive and locate it elsewhere on the cluster. Both SmartFail and the subsequent repair process are fully automated and hence require no administrator intervention.

Isilon data integrity

Isilon “ISI” Data Integrity (IDI) is the OneFS process that protects file system structures against corruption via 32-bit CRC checksums. All Isilon blocks, both for file and metadata, utilize checksum verification. Metadata checksums are housed in the metadata blocks themselves, whereas file data checksums are stored as metadata, thereby providing referential integrity. All checksums are recomputed by the initiator, the node servicing a particular read, on every request.

In the event that the recomputed checksum does not match the stored checksum, OneFS will generate a system alert, log the event, retrieve and return the corresponding parity block to the client and attempt to repair the suspect data block.

Protocol checksums

In addition to blocks and metadata, OneFS also provides checksum verification for Remote Block Management (RBM) protocol data. As mentioned above, the RBM is a unicast, RPC-based protocol developed by Isilon for use over the back-end cluster interconnect. Checksums on the RBM protocol are in addition to the Infiniband hardware checksums provided at the network layer, and are used to detect and isolate machines with certain faulty hardware components and exhibiting other failure states.

Dynamic sector repair

OneFS includes a Dynamic Sector Repair (DSR) feature whereby bad disk sectors can be forced by the file system to be rewritten elsewhere. When OneFS fails to read a block during normal operation, DSR is invoked to reconstruct the missing data and write it to either a different location on the drive or to another drive on the node. This is done to ensure that subsequent reads of the block do not fail. DSR is fully automated and completely transparent to the end-user. Disk sector errors and Cyclic Redundancy Check (CRC) mismatches use almost the same mechanism as the drive rebuild process.

MediaScan

MediaScan’s role within OneFS is to check disk sectors and deploy the above DSR mechanism in order to force disk drives to fix any sector ECC errors they may encounter. Implemented as one of the phases of the OneFS job engine, MediaScan is run automatically based on a predefined schedule. Designed as a low-impact,

background process, MediaScan is fully distributed and can thereby leverage the benefits of Isilon's unique parallel architecture.

IntegrityScan

IntegrityScan, another component of the OneFS job engine, is responsible for examining the entire file system for inconsistencies. It does this by systematically reading every block and verifying its associated checksum. Unlike traditional 'fsck' style file system integrity checking tools, IntegrityScan is designed to run while the cluster is fully operational, thereby removing the need for any downtime. In the event that IntegrityScan detects a checksum mismatch, a system alert is generated and written to the syslog and OneFS automatically attempts to repair the suspect block.

The IntegrityScan phase is run manually if the integrity of the file system is ever in doubt. Although this process may take several days to complete, the file system is online and completely available during this time. Additionally, like all phases of the OneFS job engine, IntegrityScan can be prioritized, paused or stopped, depending on the impact to cluster operations.

Fault isolation

Because OneFS protects its data at the file-level, any inconsistencies or data loss is isolated to the unavailable or failing device—the rest of the file system remains intact and available.

For example, a ten node, S200 cluster, protected at n+2, sustains three simultaneous drive failures—one in each of three nodes. Even in this degraded state, I/O errors would only occur on the very small subset of data housed on all three of these drives. The remainder of the data striped across the other two hundred and thirty-seven drives would be totally unaffected. Contrast this behavior with a traditional RAID6 system, where losing more than two drives in a RAID-set will render it unusable and necessitate a full restore from backups.

Similarly, in the unlikely event that a portion of the file system does become corrupt (whether as a result of a software or firmware bug, etc) or a media error occurs where a section of the disk has failed, only the portion of the file system associated with this area on disk will be affected. All healthy areas will still be available and protected.

As mentioned above, referential checksums of both data and meta-data are used to catch silent data corruption (data corruption not associated with hardware failures). The checksums for file data blocks are stored as metadata, outside the actual blocks they reference, and thus provide referential integrity.

Accelerated drive rebuilds

The time that it takes a storage system to rebuild data from a failed disk drive is crucial to the data reliability of that system. With the advent of four terabyte drives, and the creation of increasingly larger single volumes and file systems, typical recovery times for multi-terabyte drive failures are becoming multiple days or even weeks. During this MTDL period, storage systems are vulnerable to additional drive failures and the resulting data loss and downtime.

Since OneFS is built upon a highly distributed architecture, it's able to leverage the CPU, memory and spindles from multiple nodes to reconstruct data from failed drives in a highly parallel and efficient manner. Because Isilon is not bound by the speed of any particular drive, OneFS is able to recover from drive failures extremely quickly and this efficiency grows relative to cluster size. As such, a failed drive within an Isilon cluster will be rebuilt an order of magnitude faster than hardware RAID-based storage devices. Additionally, OneFS has no requirement for dedicated 'hot-spare' drives.

Isilon data protection

To effectively protect a file system that is hundreds of terabytes or petabytes in size requires an extensive use of multiple data availability and data protection technologies. As mentioned above, the demand for storage is continuing to grow exponentially and all predictions suggest it will continue to expand at a very aggressive rate for the foreseeable future.

In tandem with this trend, the demand for ways to protect and manage that storage also increases. Today, several strategies for data protection are available and in use. As mentioned earlier, if data protection is perceived as a continuum, at the beginning lies high availability. Without high availability technologies such as drive, network and power redundancy, data loss and its subsequent recovery would be considerably more prevalent.

Historically, data protection was always synonymous with tape backup. However, over the past decade, several technologies like replication, synchronization and snapshots, in addition to disk based backup (such as nearline storage and VTL), have become mainstream and established their place within the data protection realm. Snapshots offer rapid, user-driven restores without the need for administrative assistance, while synchronization and replication provide valuable tools for business continuance and offsite disaster recovery.

The Isilon data management suite spans the breadth of the data protection continuum and throughout the course of this paper we will examine the constituent parts in more detail.

High availability and data protection strategies

At the core of every effective data protection strategy lies a solid business continuance plan. All enterprises need an explicitly defined and routinely tested plan to minimize the potential impact to the workflow when a failure occurs or in the event of a natural disaster. There are a number ways to address data protection and most enterprises adopt a combination of these methods, to varying degrees.

Among the primary approaches to data protection are fault tolerance, redundancy, snapshots, replication (local and/or geographically separate), and backups to nearline storage, VTL, or tape.

Some of these methods are biased towards cost efficiency but have a higher risk associated with them, and others represent a higher cost but also offer an increased level of protection. Two ways to measure cost versus risk from a data protection point of view are:

- **Recovery Time Objective (RTO):** RTO is the allotted amount of time within a Service Level Agreement (SLA) to recover data. For example, an RTO of four hours means data must be restored and made available within four hours of an outage.
- **Recovery Point Objective (RPO):** RPO is the acceptable amount of data loss that can be tolerated per an SLA. With an RPO of 30-minutes, this is the maximum amount of time that can elapse since the last backup or snapshot was taken.

The Isilon high availability and data protection suite

Data Protection—Described in detail earlier, at the heart of OneFS is FlexProtect. This unique, software based data protection scheme allows differing levels of protection to be applied in real time down to a per-file granularity, for the entire file system, or at any level in between.

Redundancy— As we have seen, Isilon’s clustered architecture is designed from the ground-up to support the following availability goals:

- No single point of failure
- Unparalleled levels of data protection in the industry
- Tolerance for multi-failure scenarios
- Fully distributed single file system
- Pro-active failure detection and pre-emptive, fast drive rebuilds
- Flexible data protection
- Fully journalled file system
- High transient availability

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.

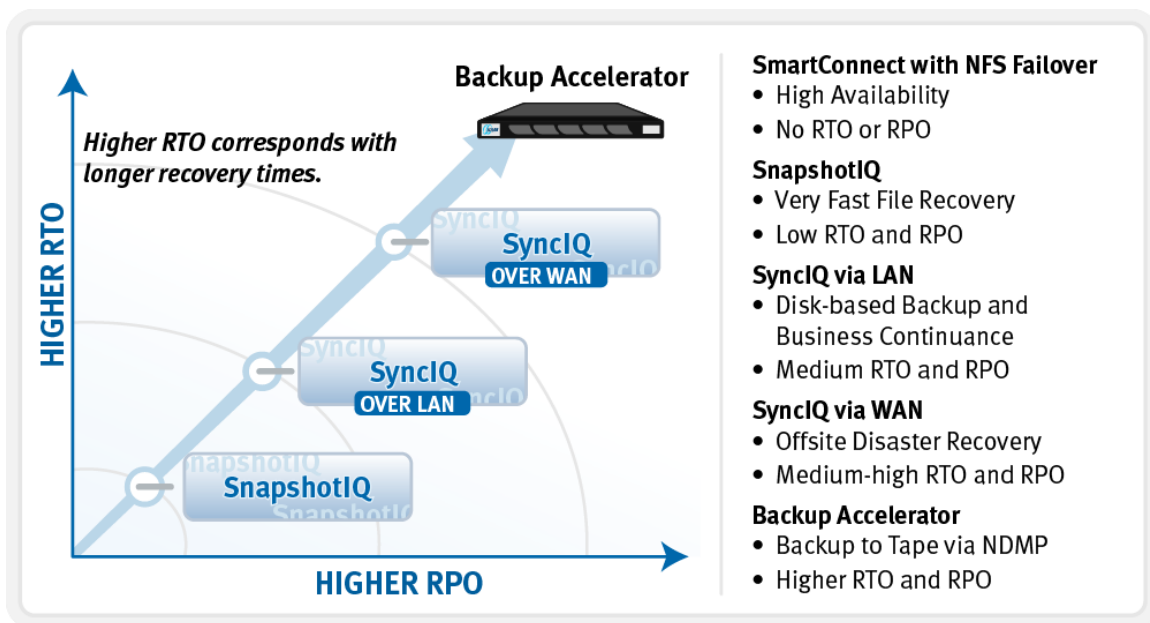


Figure 6: Isilon Data Protection technology alignment with protection continuum

Connection load balancing and failover

SmartConnect

As mentioned previously, at the leading edge of the data protection continuum lies high availability. This not only includes disk, CPU, and power redundancy, but also network resilience. EMC Isilon SmartConnect™ software contributes to data availability by supporting dynamic NFS failover and failback for Linux and UNIX clients. This ensures that when a node failure occurs, all in-flight reads and writes are handed off to another node in the cluster to finish its operation without any user or application interruption. Windows clients also benefit by easily being able to remount an SMB share using any other available node in the cluster.

During failover, clients are evenly redistributed across all remaining nodes in the cluster, ensuring minimal performance impact. If a node is brought down for any reason, including a failure, the virtual IP addresses on that node is seamlessly migrated to another node in the cluster. When the offline node is brought back online, SmartConnect automatically rebalances the NFS clients across the entire cluster to ensure maximum storage and performance utilization. For periodic system maintenance and software updates, this functionality allows for per-node rolling upgrades affording full-availability throughout the duration of the maintenance window.

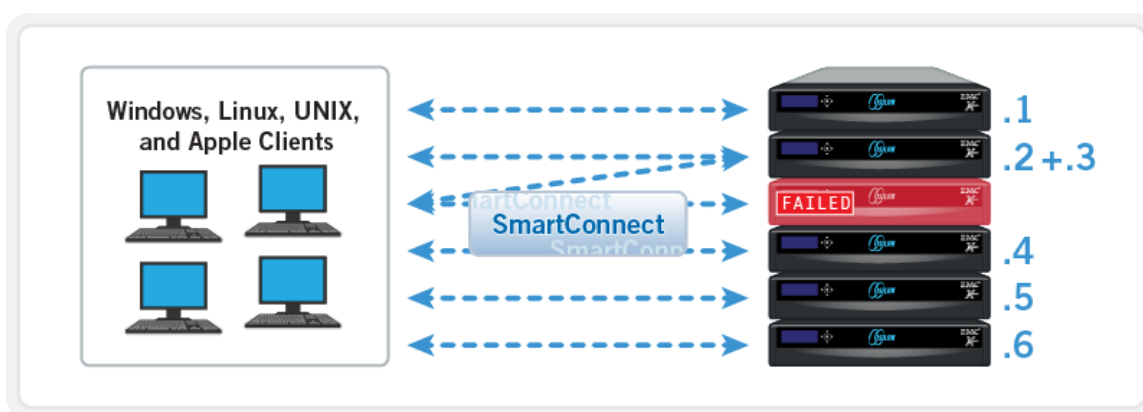


Figure 7: Seamless Client Failover with SmartConnect

Snapshots

SnapshotIQ

Next along the high availability and data protection continuum are snapshots. The RTO of a snapshot can be very small and the RPO is also highly flexible with the use of rich policies and schedules. Isilon SnapshotIQ™ software can take read-only, point-in-time copies of any directory or subdirectory within OneFS.

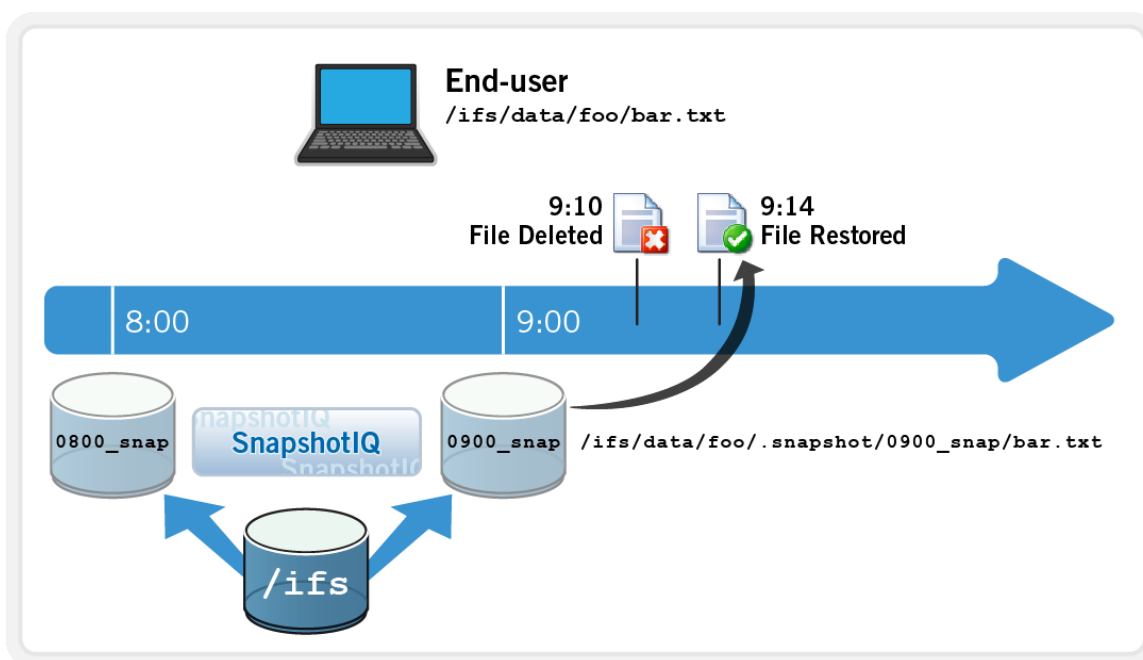


Figure 8: User Driven File Recovery with SnapshotIQ

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 when updating the snapshots, thereby ensuring highly-efficient snapshot storage utilization. User access to the

available snapshots is via a /.snapshot hidden directory under each file system directory.

Isilon SnapshotIQ can also create unlimited snapshots on a cluster. This provides a substantial benefit over the majority of other snapshot implementations because the snapshot intervals can be far more granular and hence offer improved RPO time frames.

SnapshotIQ architecture

SnapshotIQ has several fundamental differences as compared to most snapshot implementations. The most significant of these are, first, that 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. Second, 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.

The process of taking a snapshot in OneFS is relatively instantaneous. However, there is a small amount of snapshot preparation work that has to occur. First, the coalescer is paused and any existing write caches flushed in order for the file system to be quiesced for a short period of time. Next, a marker is placed at the top-level directory inode for a particular snapshot and a unique snapshot ID is assigned. Once this has occurred, the coalescer resumes and writes continue as normal. Therefore, the moment a snapshot is taken, it essentially consumes zero space until file creates, delete, modifies and truncates start occurring in the structure underneath the marked top-level directory.

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 iterate 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 'paint' the point-in-time view of that dataset.

OneFS uses both Copy on Write (CoW) and Redirect on Write (RoW) strategies for its differential snapshots and utilizes the most appropriate method for any given situation. Both have advantages and disadvantages and OneFS dynamically picks which flavor to use in order to maximize performance and keep overhead to a minimum. Typically, CoW is most prevalent, and is primarily used for small changes, inodes and directories. RoW, on the other hand, is adopted for more substantial changes such as deletes and large sequential writes.

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. Additionally, 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.

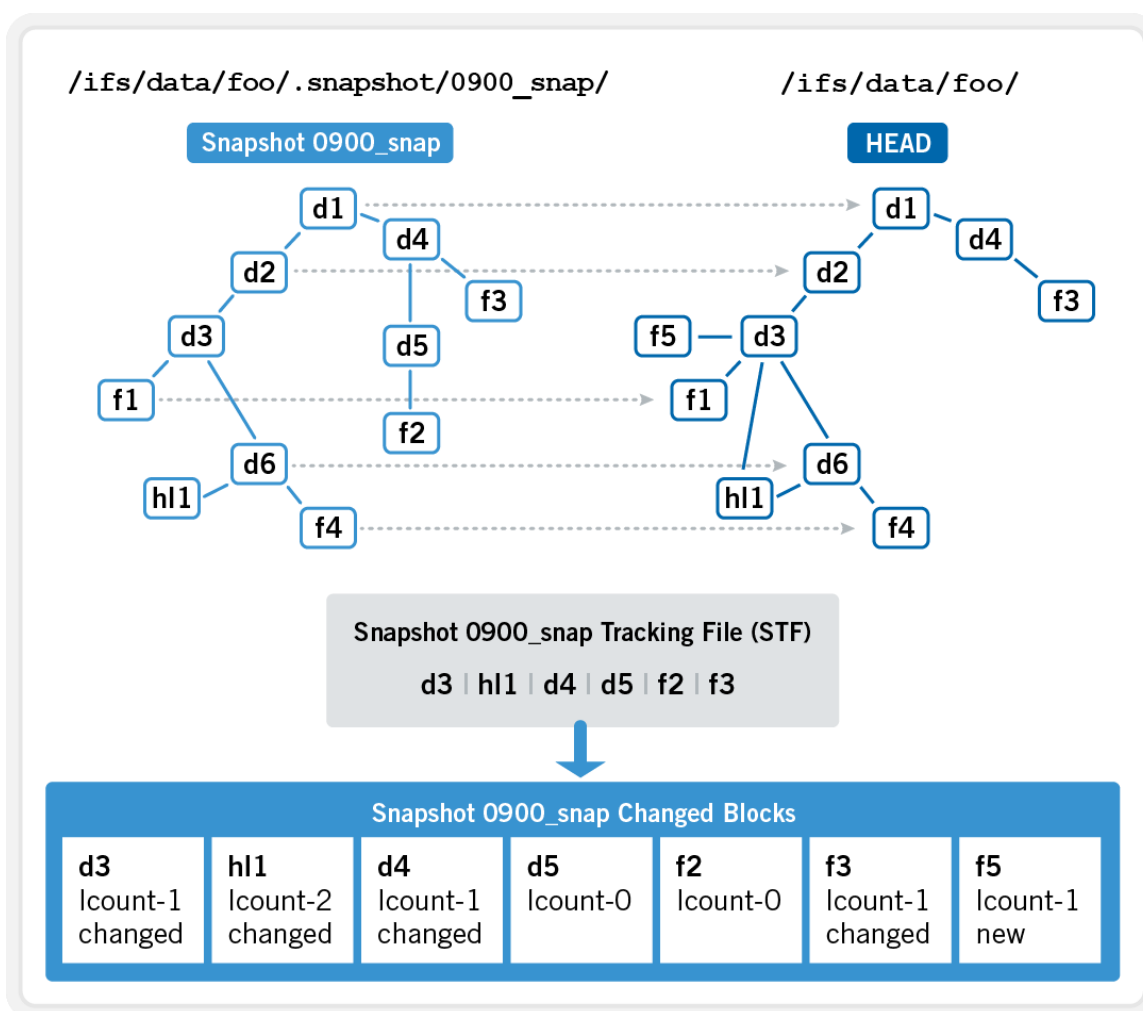


Figure 9: Snapshot Change Tracking

Snapshot scheduling

Snapshot schedules are configured at a daily, weekly, monthly or yearly interval, with single or multiple job frequency per schedule, down to a per-minute granularity. And automatic deletion can be configured per defined schedule at an hourly through yearly range.

Snapshot deletes

When snapshots are manually deleted, OneFS will mark 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 marked 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 then queued to delete a snapshot completely and then it is up to the job engine to schedule it. The 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 jobs with higher priority running at the moment. For

SnapshotDelete, it is only run if the group is in a pristine 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.

In summary, SnapshotIQ affords the following benefits:

- 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.
- 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.
- Up to 1,024 snapshots can be created per directory, and there is no hard limit of snapshots at the cluster-level (although this can be configured via a sysctl).
- The default snapshot limit is 2048 per cluster, but this is a soft limit which can be adjusted via a sysctl. However, depending on the rate of data change, cluster performance may be impacted when the number of snapshots exceeds 4000 or so.

Snapshot restore

For simple, efficient snapshot restoration, SnapshotIQ provides SnapRevert functionality. Using the Job Engine for scheduling, a SnapRevert job automates the restoration of an entire snapshot to its top level directory. This is invaluable for quickly and efficiently reverting to a previous, known-good recovery point, for example in the event of virus or malware outbreak. Additionally, individual files, rather than entire snapshots, can also be restored in place using FileRevert functionality. This can help drastically simplify virtual machine management and recovery.

File clones

OneFS File Clones provides a rapid, efficient method for provisioning multiple read/write copies of files and iSCSI LUNs. 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.

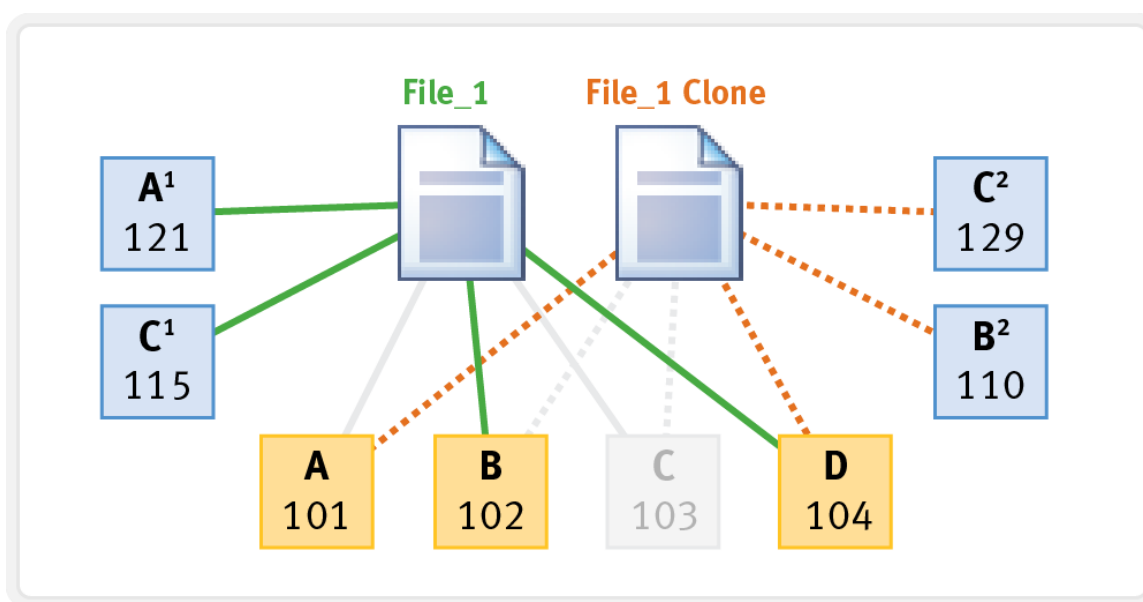


Figure 10: File Clones

Replication

SyncIQ

While snapshots provide an ideal solution for infrequent or smaller-scale data loss occurrences, when it comes to catastrophic failures or natural disasters, a second, geographically separate copy of a dataset is clearly beneficial. Here, a solution is required that is significantly faster and less error-prone than a recovery from tape, yet still protects the data from localized failure.

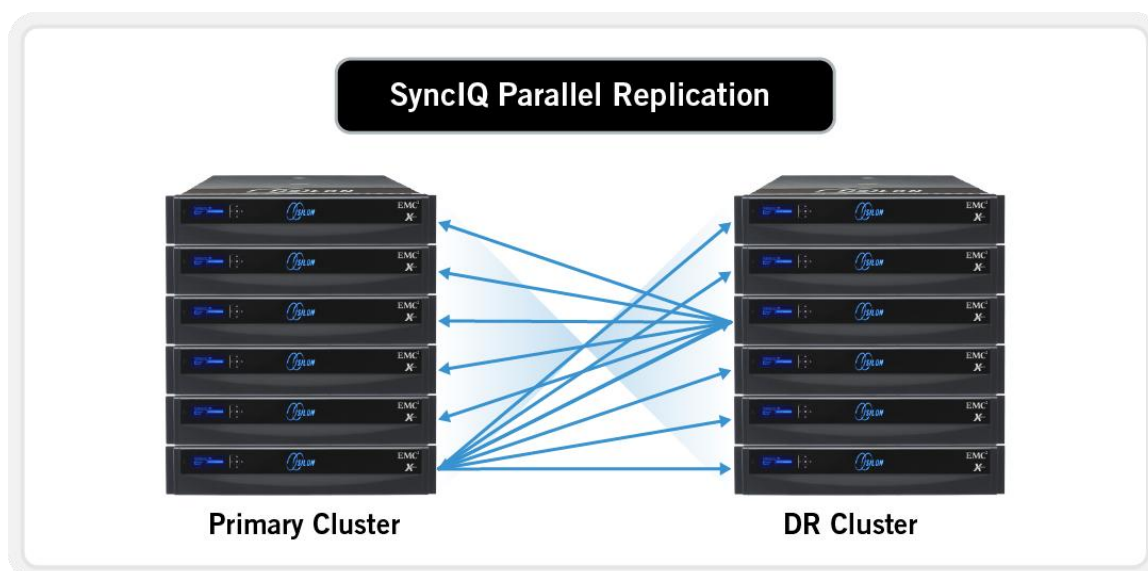


Figure 11: Disaster Recovery with SyncIQ

Isilon SyncIQ™ software delivers high-performance, asynchronous replication of unstructured data to address a broad range of recovery point objectives (RPO) and recovery time objectives (RTO). This enables customers to make an optimal tradeoff between infrastructure cost and potential for data loss if a disaster occurs. SyncIQ does not impose a hard limit on the size of a replicated file system so will scale linearly with an organization's data growth up into the multiple petabyte ranges.

SyncIQ is easily optimized for either LAN or WAN connectivity in order to replicate over short or long distances, thereby providing protection from both site-specific and regional disasters. Additionally, SyncIQ utilizes a highly-parallel, policy-based replication architecture designed to leverage the performance and efficiency of clustered storage. As such, aggregate throughput scales with capacity and allows a consistent RPO over expanding data sets.

There are two basic implementations of SyncIQ:

- The first is utilizing SyncIQ to replicate to a local target cluster within a datacenter. The primary use case in this scenario is disk backup and business continuance.
- The second implementation uses SyncIQ to replicate to a remote target cluster, typically located in a geographically separate datacenter across a WAN link. Here, replication is typically utilized for offsite disaster recovery purposes.

In either case, a secondary cluster synchronized with the primary production cluster can afford a substantially improved RTO and RPO than tape backup and both implementations have their distinct advantages. And SyncIQ performance is easily tuned to optimize either for network bandwidth efficiency across a WAN or for LAN speed synchronization. Synchronization policies may be configured at the file-, directory- or entire file system-level and can either be scheduled to run at regular intervals or executed manually.

SyncIQ linear restore

Leveraging OneFS SnapshotIQ infrastructure, the Linear Restore functionality of SyncIQ is able to detect and restore (commit) consistent, point in time, block-level changes between cluster replication sets, with a minimal impact on operations and a granular RPO. This 'change set' information is stored in a mirrored database on both source and target clusters and is updated during each incremental replication job, enabling rapid failover and failback RTOs.

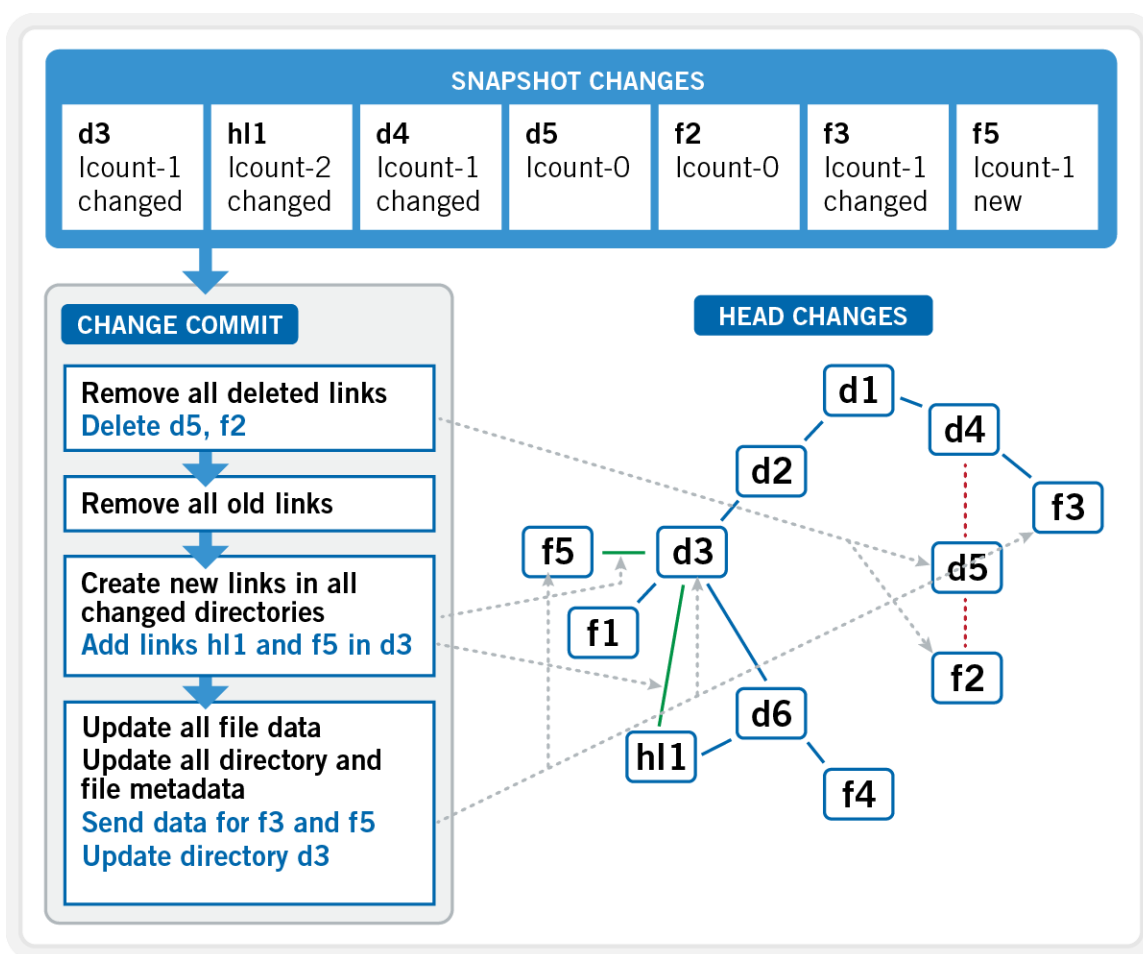


Figure 12: The SyncIQ Linear Restore Change Commit Mechanism

SyncIQ replica protection

All writes outside of the synchronization process itself are disabled on any directory that is a target for a specific SyncIQ job. However, if the association is broken between a target and a source, the target may then return to a writable state. Subsequent resolution of a broken association will force a full resynchronization to occur at the next job run. As such, restricted writes prevent modification, creation, deletion, linking or movement of any files within the target path of a SyncIQ job. Therefore, replicated disaster recovery (DR) data is protected within and by its SyncIQ container or restricted-writer domain, until a conscious decision is made to bring it into a writeable state.

SyncIQ failover and fallback

In the event that a primary cluster becomes unavailable, SyncIQ provides the ability to failover to a mirrored, DR cluster. During such a scenario, the administrator makes the decision to redirect client I/O to the mirror and initiates SyncIQ failover on the DR cluster. Users will continue to read and write to the DR cluster while the primary cluster is repaired.

Once the primary cluster becomes available again, the administrator may decide to revert client I/O back to it. To achieve this, the administrator initiates a SyncIQ

failback prep process which synchronizes any incremental changes made to the DR cluster back to the primary.

Failback is divided into three distinct phases:

1. First, the prep phase readies the primary to receive changes from the DR cluster by setting up a restricted writer domain and then restoring the last known good snapshot.
2. Next, upon successful completion of failback prep, a final failback differential sync is performed.
3. Lastly, the administrator commits the failback, which restores the primary cluster back to its role as the source and relegates the DR cluster back to a target again.

In addition to the obvious unplanned failover and failback, SyncIQ also supports controlled, proactive cluster failover and failback. This provides two major benefits:

- The ability to validate and test DR procedures and requirements
- Performing planned cluster maintenance.

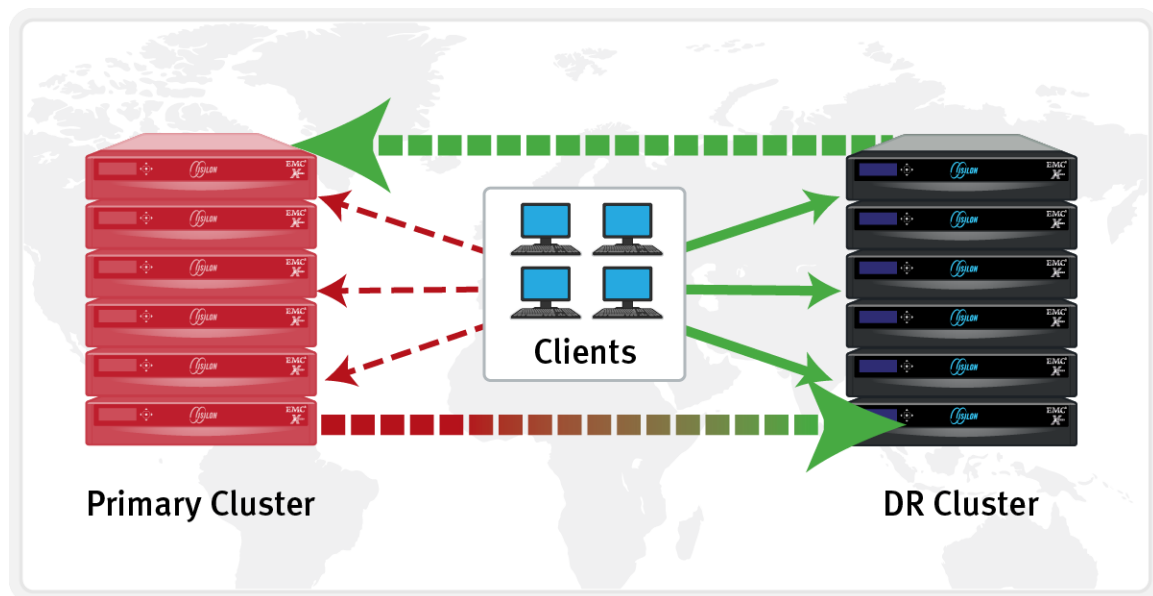


Figure 13: SyncIQ Automated Data Failover and Failback

Archive

As we have seen, Isilon SyncIQ software (described above) enables the simple creation of a secure remote archive. Additionally, SmartPools (OneFS tiering module) also facilitates the creation and management of a dedicated local archive pool within a cluster for data retention and high availability purposes.

SmartLock

OneFS utilizes Isilon SmartLock™ software to provide immutable storage for data. Based on a write once, read many (WORM) locking capability, SmartLock ensures tamper-proof archiving of critical data sets for disaster recovery and regulatory

compliance purposes. Configured at the directory-level, SmartLock delivers simple to manage secure data containers that remain locked for a configurable duration or indefinitely. Additionally, SmartLock satisfies the regulatory compliance demands of stringent data retention policies, including SEC 17a-4.

Nearline, VTL and tape backup

At the trailing end of the protection continuum lies traditional backup and restore—whether to tape or disk. This is the bastion of any data protection strategy and usually forms the crux of a ‘data insurance policy’. With high RPO and RTOs, often involving a retrieval of tapes from secure, offsite storage, tape backup is typically the mechanism of last resort for data recovery in the face of a disaster.

Backup Accelerator

Isilon provides the ability to perform large-scale backup and restore functions across massive, single-volume data sets—while leveraging an enterprise’s existing, SAN-based tape and VTL infrastructure. This is enabled by the Backup Accelerator (BA) node, which features a quad-port 4GB/s Fibre Channel card, quad-core processors, and 8GB of RAM.

A single Backup Accelerator can concurrently stream backups at 480MB/s, or 1.7TB/hour, across its four Fibre Channel ports. Additionally, as data grows, multiple Backup Accelerator nodes can be added to a single cluster to support a wide range of RPO/RTO windows, throughput requirements and backup devices.

Backup from snapshots

In addition to the benefits provided by SnapshotIQ in terms of user recovery of lost or corrupted files, it also offers a powerful way to perform backups while minimizing the impact on the file system.

Initiating backups from snapshots affords several substantial benefits. The most significant of these is that the file system does not need to be quiesced, since the backup is taken directly from the read-only snapshot. This eliminates lock contention issues around open files and allows users full access to data throughout the duration of the backup job.

SnapshotIQ also automatically creates an alias which points to the latest version of each snapshot on the cluster, which facilitates the backup process by allowing the backup to always refer to that alias. Since a snapshot is by definition a point-in-time (PIT) copy, by backing up from a snapshot, the consistency of the file system or sub-directory is maintained.

This process can be further streamlined by using the Network Data Management Protocol (NDMP) snapshot capability to create a snapshot as part of the backup job, then delete it upon successful completion of the backup.

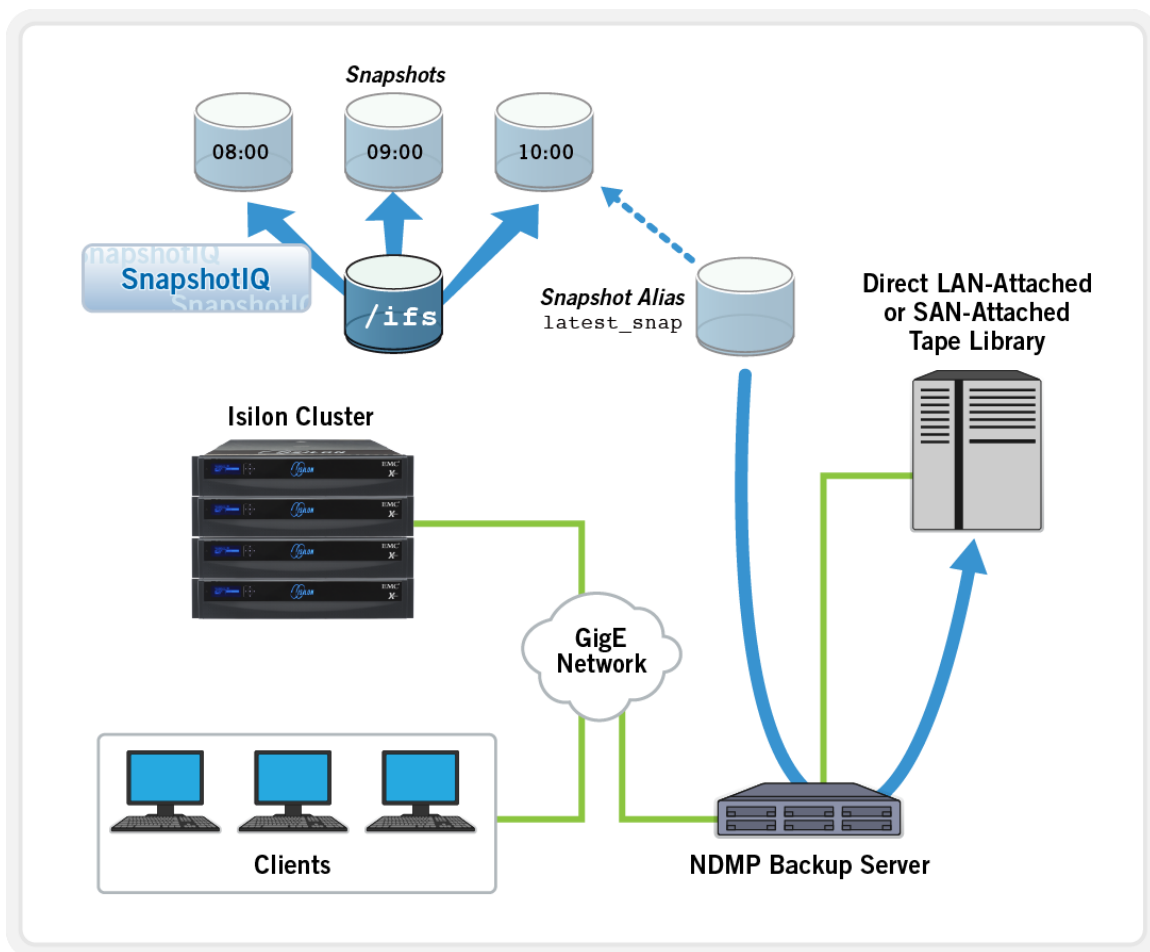


Figure 14: Backup Using SnapshotIQ

Parallel streams

Isilon's distributed architecture allows backups to be spread across multiple network streams from the cluster, which can significantly improve performance. This is achieved by dividing the root file system into several paths based on the number of nodes in the cluster and the structure of sub-directories under the file system root. For example, if the file system on a four-node cluster can be segregated logically among four sub-directories, each of these sub-directories can be backed up as a separate stream, one served from each node.

NDMP

OneFS facilitates performant backup and restore functionality via its support of the ubiquitous Network Data Management Protocol (NDMP). NDMP is an open-standard protocol that provides interoperability with leading data-backup products and Isilon supports both NDMP versions 3 and 4. The OneFS NDMP module includes the following functionality:

- Full and incremental backups and restores using NDMP
- Direct Access Restore/Directory Direct Access Restore (DAR/DDAR), single-file restores, and three-way backups

- Restore-to-arbitrary systems
- Seamless integration with access control lists (ACLs), alternate data streams and resource forks
- Selective File Recovery
- Replicate then backup

While some backup software vendors may support backing up OneFS over CIFS and NFS, the advantages of using NDMP include:

- Increased performance
- Retention of file attributes and security and access controls
- Backups utilize automatically generated snapshots for point-in-time consistency.
- Extensive support by backup software vendors

OneFS provides support for NDMP version 4, and both direct NDMP (referred to as 2-way NDMP), and remote NDMP (referred to as 3-way NDMP) topologies.

Direct NDMP model

This is the most efficient model and results in the fastest transfer rates. Here, the data management application (DMA) uses NDMP over the Ethernet front-end network to communicate with the Backup Accelerator. On instruction, the Backup Accelerator, which is also the NDMP tape server, begins backing up data to one or more tape devices which are attached to it via Fibre Channel.

The Backup Accelerator is an integral part of the Isilon cluster and communicates with the other nodes in the cluster via the internal InfiniBand network. The DMA, a separate server, controls the tape library's media management. File History, the information about files and directories, is transferred from the Backup Accelerator via NDMP to the DMA, where it is maintained in a catalog.

Direct NDMP is the fastest and most efficient model for backups with OneFS and obviously requires one or more Backup Accelerator nodes to be present within a cluster.

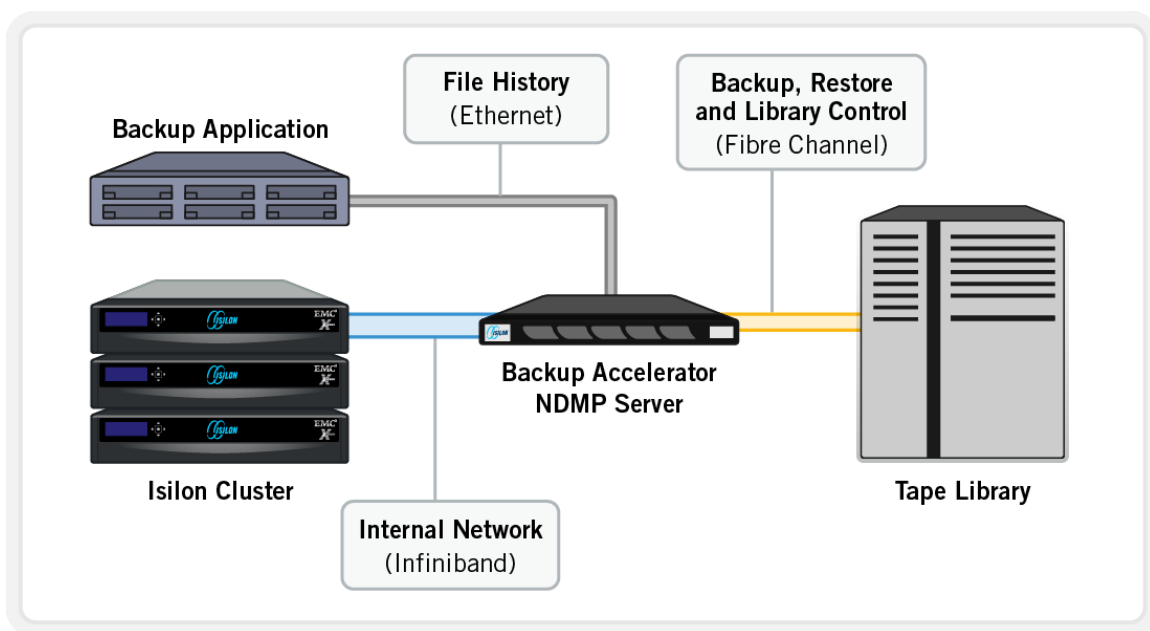


Figure 15: Recommended Two-way NDMP with Backup Accelerator

Remote NDMP model

In the remote NDMP scenario, there is no Backup Accelerator present. In this case, the DMA uses NDMP over the LAN to instruct the cluster to start backing up data to the tape server - either connected via Ethernet or directly attached to the DMA host. In this model, the DMA also acts as the Backup/Media Server.

During the backup, file history is transferred from the cluster via NDMP over the LAN to the backup server, where it is maintained in a catalog. In some cases, the backup application and the tape server software both reside on the same physical machine.

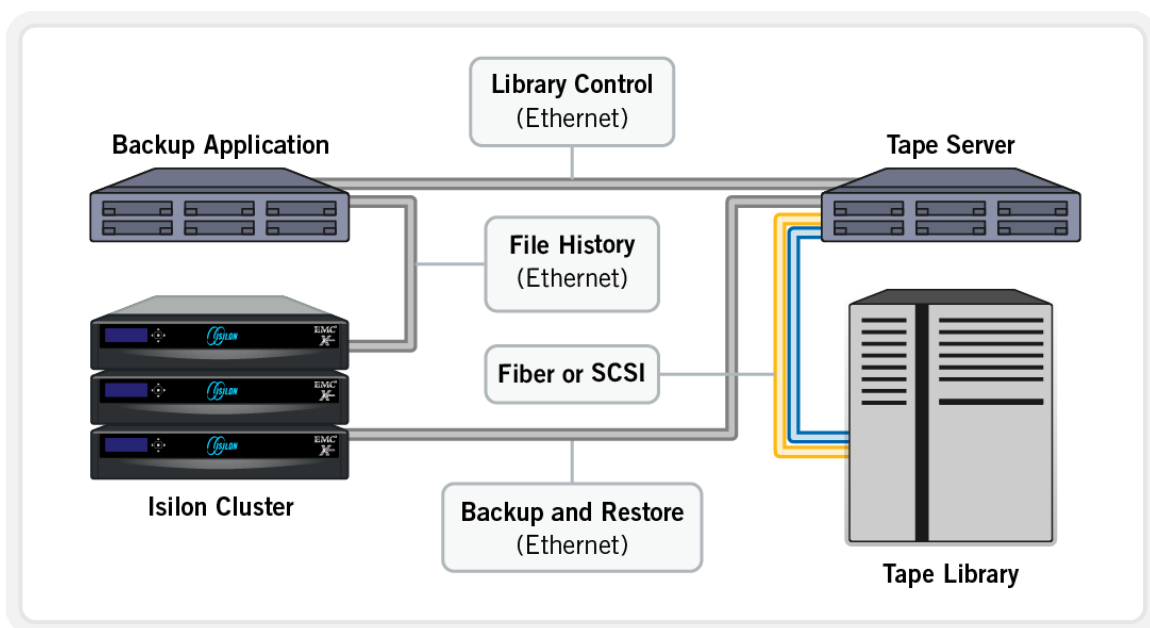


Figure 16: Remote Three-way NDMP Backup

Incremental backups

Isilon OneFS accommodates the range of full, incremental and token-based backups. In standard DR nomenclature, Level 0 indicates a full backup, and levels 1-9 are incrementals. Any level specified as 1-9 will back up all the files that have been modified since the previous lower level backup.

Token-based incremental backups are also supported. There are achieved by configuring the data management application (DMA) to maintain a timestamp database and to pass the reference time token on to the cluster for use during each incremental backup. This method does not rely on level based incremental backups, as described above, at all.

Direct access recovery

OneFS provides full supports for Direct Access Recovery (DAR). Direct Access Recovery allows the NDMP server to go directly to the location of a file within an archive and quickly recover that file. As such, it eliminates the need to scan through vast quantities of data typically spread across multiple tapes in an archive set, in order to recover a single file. This capability uses the offset information that is contained in the file history data passed to the DMA at backup time.

Directory DAR

Isilon OneFS NDMP also supports Directory DAR (DDAR), an extension of DAR. DDAR allows the NDMP server to go directly to the location of a directory within an archive and quickly recover all files/directories contained within the directory tree hierarchy. Clearly, both DAR and DDAR provide an improved RTO for smaller scale data recovery from tape.

OneFS NDMP offers Selective File Recovery—the ability to recovering a subset of files within a backup archive. Also supported is the ability to restore to alternate path locations.

Certified backup applications

Isilon OneFS is certified with a wide range of leading enterprise backup applications, including:

- Symantec NetBackup
- EMC NetWorker
- IBM TSM
- CommVault Simpana
- Quest Software (formerly Bakbone) NetVault
- Atempo Time Navigator

OneFS is also certified to work with the EMC Cloud Tiering Appliance to simplify data migration and with EMC DataDomain appliance products for deduplicated backup and archiving.

Summary

Organizations of all sizes around the globe are dealing with a deluge of digital content and unstructured data that is driving massive increases in storage needs. As these enterprise datasets continue to expand to unprecedented sizes, data protection has never been more crucial. A new approach is needed to meet the availability, protection and performance requirements of this era of 'big data'.

EMC Isilon enables organizations to linearly scale capacity and performance to over 20 petabytes, 106GB per second and 1.6 million SPECsfs2008 CIFS file operations per second. Moreover, they can do this within a single file system—one which is both simple to manage and highly available and redundant, as we have seen. Built on commodity hardware and powered by the revolutionary OneFS distributed file system, Isilon scale-out NAS solutions deliver the following key tenets:

- Unparalleled levels of data protection
- No single point of failure
- Fully distributed single file system
- Industry leading tolerance for multi-failure scenarios
- Pro-active failure detection and pre-emptive, fast drive rebuilds
- Flexible, file-level data protection
- Fully journalled file system
- Extreme transient availability

Isilon acronyms glossary

BAM	Block Allocation Manager	LAGG	Link Aggregation
BAT	Block Allocation Type	LBM	Local Block Manager
BH	Block History	LIN	Logical Inode
BSD	Berkeley Software Distribution UNIX	MDS	Mirrored Data Structure
BSW	BAM Safe Write	NFS	Network File System
CIFS	Common Internet File System	PiT	Point in Time snapshot
CoW	Copy on Write snapshot	POSIX	Portable Operating System Interface for UNIX
DFM	Directory Format Manager	RBM	Remote Block Manager
DSR	Dynamic Sector Repair	RoW	Redirect on Write snapshot
DWT	Device worker Thread	SDP	Sockets Direct Protocol
FEC	Forward Error Correction	SMB	Server Message Block
IDI	Isilon data integrity	TXN	Transaction Code
IFM	Inode Format Manager	VFS	Virtual File System
IMDD	Isilon Mirrored Device Driver	VOPs	Vnode Operations
iSCSI	Internet SCSI		
LACP	Link Aggregation Control Protocol		

About EMC Isilon

Isilon, a division of EMC, is the global leader in scale-out NAS. We deliver powerful yet simple solutions for enterprises that want to manage their data, not their storage. Isilon products are simple to install, manage and scale, at any size and, unlike traditional enterprise storage, Isilon stays simple no matter how much storage is added, how much performance is required, or how business needs change in the future. We're challenging enterprises to think differently about their storage, because when they do, they'll recognize there's a better, simpler way. Learn what we mean at www.isilon.com.

Contact Isilon

<http://www.isilon.com>

505 1st Avenue South, Seattle, WA 98104

Toll-Free: 877-2-ISILON • Phone: +1-206-315-7602

Fax: +1-206-315-7501 • Email: sales@isilon.com