



Dell Fluid File System

A Dell Technology White Paper

Version 2.0

Storage Solutions Engineering

Dell Product Group

March 2012

THIS TECHNOLOGY WHITE PAPER IS FOR INFORMATIONAL PURPOSES ONLY, AND MAY CONTAIN TYPOGRAPHICAL ERRORS AND TECHNICAL INACCURACIES. THE CONTENT IS PROVIDED AS IS, WITHOUT EXPRESS OR IMPLIED WARRANTIES OF ANY KIND.

© 2012 Dell Inc. All rights reserved. Reproduction of this material in any manner whatsoever without the express written permission of Dell Inc. is strictly forbidden. For more information, contact Dell.

Dell, the *DELL* logo, and the *DELL* badge, and *PowerVault* are trademarks of Dell Inc. *Microsoft* and *Windows* are either trademarks or registered trademarks of Microsoft Corporation in the United States and/or other countries. *Symantec*, *BackupExec*, and *NetBackup* are trademarks of Symantec Corporation or its affiliates in the U.S. and other countries. *Commvault* and *Simpana* are registered trademarks of Commvault Systems, Inc. Other trademarks and trade names may be used in this document to refer to either the entities claiming the marks and names or their products. Dell Inc. disclaims any proprietary interest in trademarks and trade names other than its own.

March 2012

Contents

- 1 Abstract 4
- 2 Dell Fluid File System Overview..... 5
- 3 Architecture..... 6
 - 3.1 Dell Fluid File System Components 6
 - 3.2 Fluid File System Clustering.....7
 - 3.3 Distributed File System7
 - 3.4 Global Namespace..... 9
 - 3.5 Scalability 9
- 4 Performance11
 - 4.1 Optimization for large and small file sizes11
 - 4.2 Optimized Caching11
 - 4.3 Write Optimization11
 - 4.4 Balanced I/O11
 - 4.5 Load Balancing12
 - 4.6 Resource Optimization12
- 5 Data Integrity.....13
 - 5.1 Cache Mirroring and Metadata13
 - 5.2 High availability13
- 6 Data Protection.....14
 - 6.1 Snapshots.....14
 - 6.2 Backups.....15
 - 6.3 Replication.....15
- 7 Product Integration.....16
- Appendix A Fluid File System Solutions17
 - A.1.1 Dell PowerVault NX350017
 - A.1.2 Dell EqualLogic FS750018

1 Abstract

Traditional approaches to handling file data growth have proven to be costly, hard to manage, and difficult to scale effectively and efficiently. Dell™ Fluid File System is designed to go beyond the limitations of traditional file systems with a flexible architecture that enables organizations to scale out non-disruptively. It addresses challenges that organizations face by allowing them to gain control of their data, reduce complexity, and meet growing data demands over time.

The Fluid File System architecture is open-standards based, supports industry standard protocols, and provides innovative features relating to high availability, performance, efficient data management, data integrity, and data protection. As a core component of the Dell Fluid Data architecture, Fluid File System brings differentiated value to the various Dell storage offerings. It is a network attached storage (NAS) file system accessed using CIFS and NFS protocols, but it has features and enhancements that make it unique, as discussed in the remaining sections of this document.

2 Dell Fluid File System Overview

The relentless growth of unstructured and file data is accelerating the need for network file storage systems. Organizations coping with data growth are confronted with several challenges:

- Data silos prevent easy access to vital business information.
- Data migration, backup, and disaster recovery are complex, consuming administrative time and resources.
- Meeting data growth by deploying more and more storage systems increases both the administrative burden and capital expenditure at a time when businesses need to run lean.
- Traditional file systems have scalability limitations that make them unwieldy for organizations with rapidly expanding file data.

Dell Fluid File System is the result of Dell's focus on offering superior technology that enables our customers to meet these types of critical enterprise IT challenges. Built on intellectual property acquired from Exanet, Ltd., Fluid File System has been developed by Dell and has a dedicated roadmap for future enhancements. It provides a consistent file system across all Dell storage platforms.

Fluid File System is an enterprise-class distributed file system that provides customers with the tools necessary to manage file data in an efficient and simple manner. It removes the scalability limitations associated with traditional file systems, and it supports both scale-out performance and scale-up capacity expansion, all within a single namespace for ease of administration. An optimal combination of performance and scalability makes Fluid File System an excellent choice for a wide range of use cases and deployment environments.

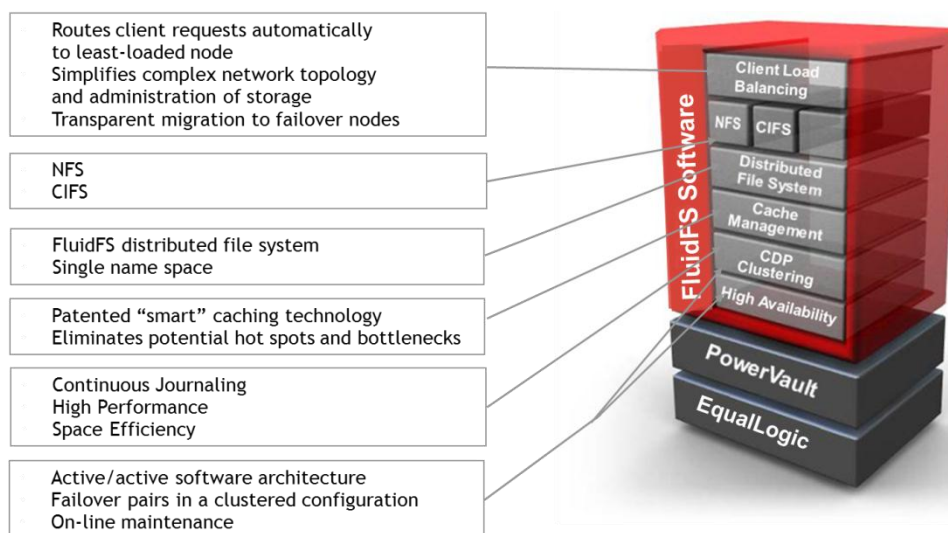


Figure 1 Fluid File System Software Stack

3 Architecture

The Dell Fluid File System architecture was designed from the ground up to provide solutions that address the demands and challenges customers face every day. This section will cover the fundamental features that allow Fluid File System to perform exceptionally well and allow growth.

3.1 Dell Fluid File System Components

The components that make up the solution are the Fluid File System software running on two or more front-end controllers, a backup power supply, and block-based back-end storage subsystems. The client-facing controllers host the Fluid File System software and provide the file sharing infrastructure. The controllers are connected in redundant pairs to enable data access and provide failover capability.

Figure 2 shows the logical architecture of Fluid File System. The different components are discussed in further detail, later in this document.

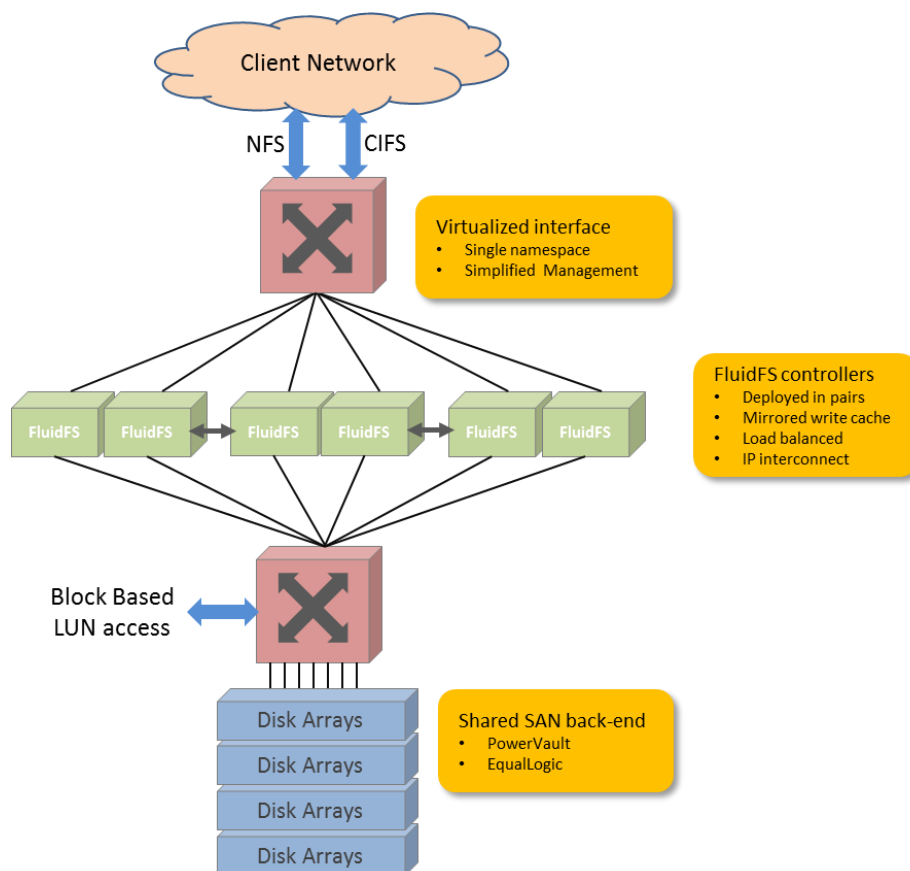


Figure 2 Fluid File System Architecture

3.2 Fluid File System Clustering

The Fluid File System architecture is inherently highly available by virtue of an underlying cluster technology that consists of multiple controllers working together, monitoring each other, and providing automatic failover capabilities. The basic implementation is a pair of controllers in a cluster, but it can be scaled to multiple controller pairs, depending on client workload characteristics. To achieve data distribution and maintain high availability, each controller in a cluster has access to all other controllers in the cluster through a dedicated and redundant interconnect network.

Dell has integrated and designed the highly available Fluid File System cluster architecture to incorporate the following criteria:

- **No single point of failure** — All critical system components including hardware and software are redundant. Write cache is mirrored between controllers to provide availability and prevent data loss. Multiple network paths to each controller shield against network failures.
- **Automatic recovery** — Fluid File System continuously monitors all hardware and software components and, in the event of failure, maintains data availability without manual intervention.
- **Self-healing** — A cluster enables each member controller to monitor its peer. If a controller detects a service failure on a peer controller, it tries to restart the controller before initiating a failover.

3.3 Distributed File System

Dell Fluid File System is a distributed file system, meaning that processing is distributed across multiple pieces of hardware and software while being transparent to the end user. This effectively uses all available resources in a balanced manner increasing the overall system utilization. The distributed file system consists of a collection of multiple instances of coupled modules, each performing a particular function:

- **Front-end service** — Acts as a protocol converter, translating between client-side protocol requests and the internal file system requests. Dedicated front-end services are available for CIFS and NFS protocols.
- **Store agent** — Handles all I/O to the back-end, block-based storage subsystem.
- **File System Daemon (FSD)** — Responsible for all client requests and is the fundamental component of the Fluid File System.

The File System Daemon is used to manage file access. The front-end service uses FSDs to open sessions and interact with the Fluid File System. The FSDs manage the file system data comprising actual data and the metadata (information about the data). Each FSD owns the metadata for the files it creates, while allowing direct data access by other FSDs. Each of the FSDs sees and accesses the entire file system and has dedicated system resources that make it extremely efficient in responding to client requests.

When a client creates and modifies files, free blocks from available storage can be allocated to any FSD. As a client deletes data, the system frees the storage blocks, returning them to the Fluid File System's overall storage capacity. 0 shows the different FSD modules and how they interact with each other.

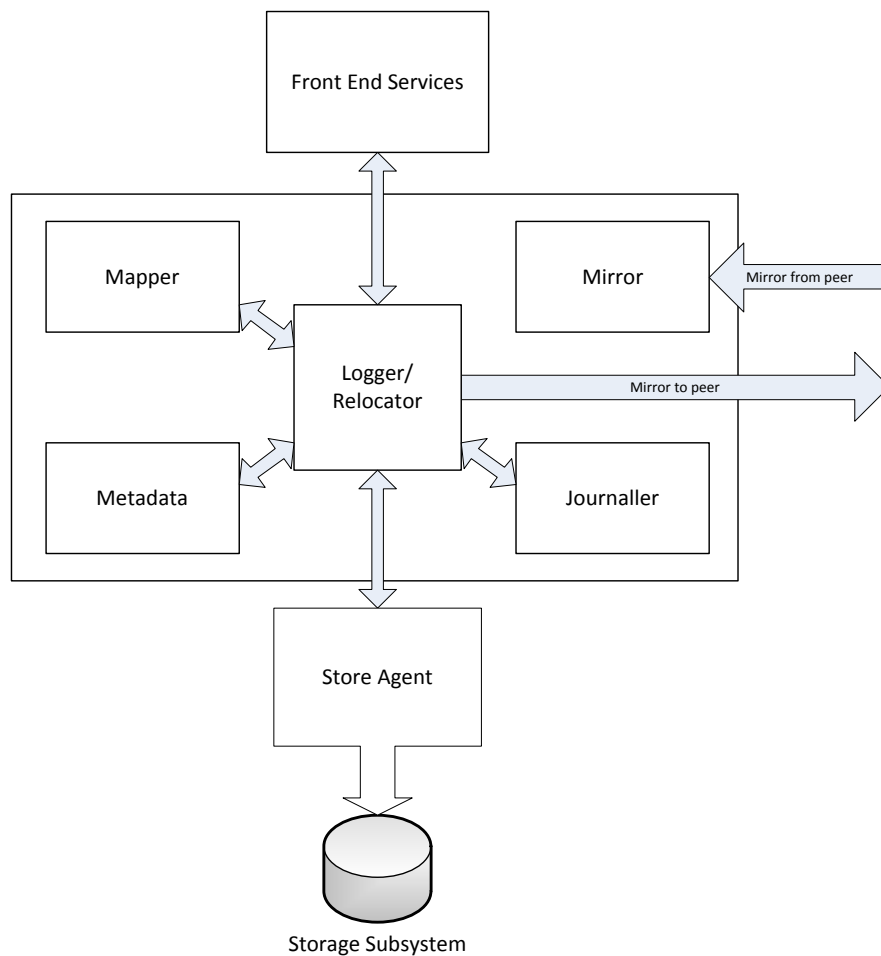


Figure 3 File System Daemon (FSD) Modules

Logger — Performs the two major tasks of managing the read/write cache and data relocation. The cache manages all current active data and metadata.

Relocator — Background relocation process that sends the written data to the back-end, block-based storage subsystem. The relocator applies coalescing algorithms to accelerate disk I/O. This aspect is explained further in the Performance section of this document.

Mapper — Serves as a pointer repository that tracks the location of data distributed across the system.

Metadata — A plain repository for file metadata, such as ownership, permission, and size. Metadata access in Fluid File System is highly optimized and will be discussed in detail in the Performance section of this document.

Mirror — Maintains an image of the write cache of the paired instance.

Journaler — Constantly journals the metadata to disk during normal operation for consistent file system recovery in the event of a failure. The actual data and metadata are also journaled in certain scenarios, which is discussed in detail in the Data Integrity section of this document.

3.4 Global Namespace

A Global Namespace allows clients to access the entire storage capacity as a single entity without knowing the intricacies of the physical infrastructure, and it is key to the effectiveness of a distributed clustered file system.

A Fluid File System cluster is accessed and managed as a single NAS system, regardless of how many controllers are in the cluster. After controllers are discovered and added to the cluster, there are no controller-specific operations for the administrator to manage. NAS volumes are virtual entities that span the underlying capacity presented from the back-end, block-based storage subsystem and provide a context for additional inherent features. When new back-end storage capacity is added, the available space for the NAS volumes can be increased dynamically. Additionally, the NAS volumes can be resized non-disruptively.

In a Fluid File System cluster, Global Namespace is achieved by using one or more Virtual IP (VIP) addresses to access the entire NAS file system. This means that as the Fluid File System platform scales, customers don't need to worry about managing multiple mounts or redesigning applications to accommodate a fragmented namespace.

3.5 Scalability

Fluid File System provides the ability to scale up and scale out. Scaling up simply increases the storage capacity, without disruption, as the data storage needs grow (see Figure 4). Scaling out grows the cluster by adding additional controller pairs as well as additional back-end, block-based storage subsystems, depending upon the solution being implemented. Figure 5 and Figure 6 show the scale out options available in Fluid File System.

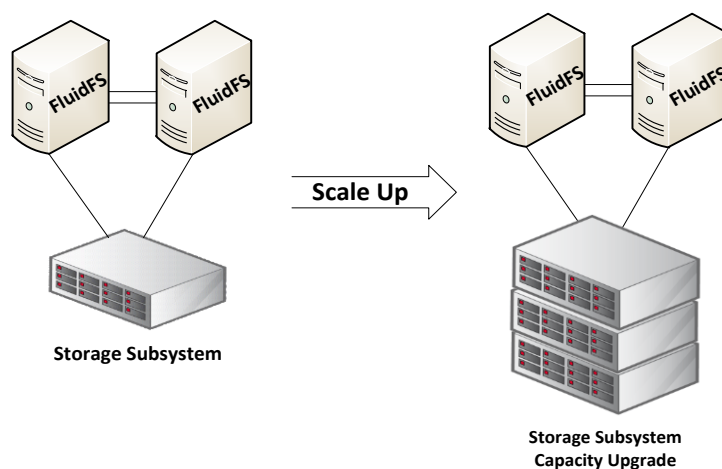


Figure 4 Scale up capacity

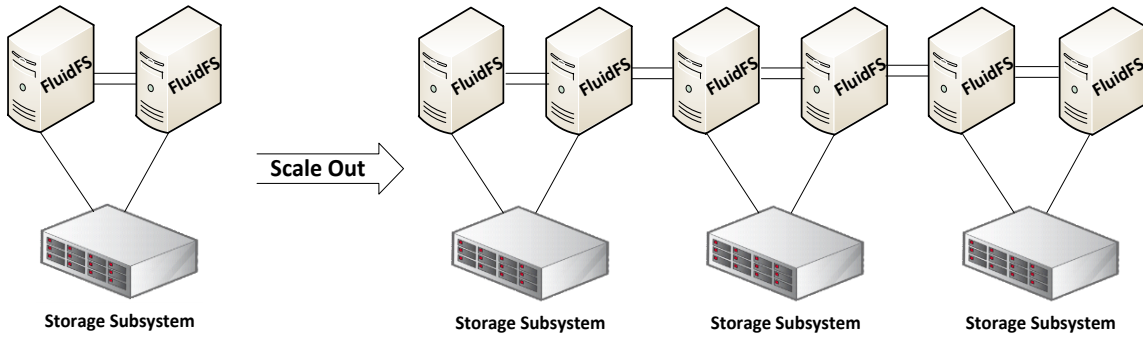


Figure 5 (Scale Out - Controller and Storage)

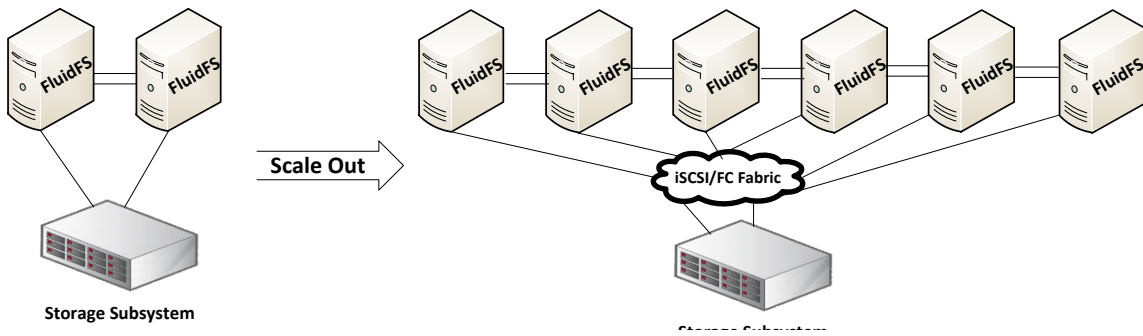


Figure 6 (Scale Out - Controller Only)

4 Performance

Traditional NAS file systems are serviced by a single controller, and hence the performance attributes of a single file system are restricted by the available resources within a controller. However, Fluid File System — with its distributed and clustered architecture, along with the global namespace — ensures that all File System Daemons (FSD) across all controllers actively engage in servicing client I/O requests and deliver high performance. Performance can be scaled further by adding controller pairs to the cluster.

4.1 Optimization for large and small file sizes

Fluid File System is optimized for both large and small file sizes to ensure performance, reliability, and capacity efficiency associated with specialized workloads. For large files, Fluid File System improves performance and minimizes fragmentation by distributing NAS data intelligently across all available back-end, block-based storage coalesced at 1MB. Files smaller than 4KB are stored along with the metadata. This allows metadata and the actual data to be read with a single disk I/O operation, which improves read access times and overall file system performance.

4.2 Optimized Caching

The Fluid File System cache is organized as a pool of 4KB pages and is used for data as well as metadata. Data is evicted from cache based on the Least Recently Used (LRU) algorithm. Fluid File System maintains separate LRUs for data and metadata, thus ensuring metadata is retained longer in cache. This allows Fluid File System to deliver high metadata performance, which eliminates one of the major bottlenecks of traditional NAS systems.

In addition, Fluid File System adapts to read-intensive, write-intensive, and mixed workloads by maintaining separate LRUs for read and write, as well as auto-tuning the size of the shared read/write cache at all times. Each FSD reads and caches the data that is accessed by the clients connected to it. All subsequent access to the same data is serviced from cache, reducing back-end disk operations and improving the response time.

4.3 Write Optimization

Fluid File System aggregates small files and then stripes the data across the available back-end storage for more efficient write operations. This process — also known as “write coalescing” — takes a random access pattern and converts it into a sequential disk operation that yields much higher throughput. Additionally, all client writes are acknowledged after being written to the cache of the local controller and mirrored to the cache of the paired controller, thereby avoiding the latency associated with disk access. The data is later asynchronously de-staged to disk.

4.4 Balanced I/O

Fluid File System stripes data across the available back-end storage capacity to improve performance. The data placement algorithm also takes into account the capacity and the percentage used of the underlying block storage to ensure that it is both load and capacity balanced.

4.5 Load Balancing

Client access to Fluid File System is load balanced across all the NAS controllers, as well as the network interfaces within these controllers, for higher performance. The load is balanced between the interfaces within a controller using one of two industry standard algorithms — Alternate Load Balancing (ALB) or Link Aggregation Protocol (LACP). ALB is a MAC address-based balancing mechanism that does not require any configuration on the network switch. LACP, which is also known as 802.3ad, is another available option and is supported by most major manufacturers' managed switches; some configuration is required on the switch.

For client access within the same subnet, Fluid File System uses a process called the “arper” to balance client connections across all available network interfaces within the cluster. For load balancing with multiple subnets, Fluid File System supports multiple solutions, such as additional virtual IPs and DNS Round Robin. Figure 7 shows load balancing in a Fluid File System cluster on a single subnet.

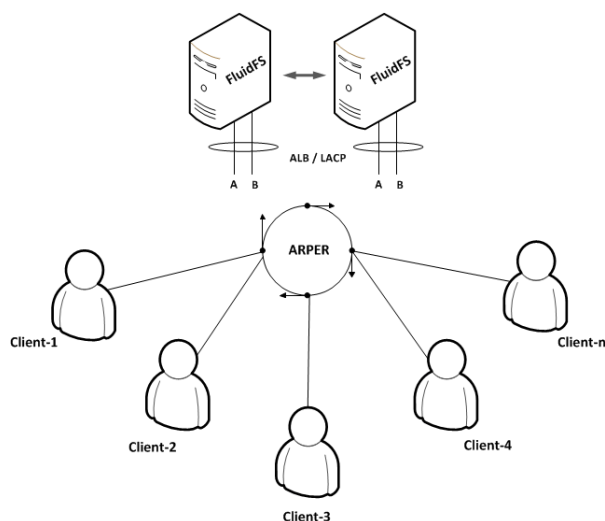


Figure 7 Load Balancing in a Fluid File System Cluster

4.6 Resource Optimization

Fluid File System actively uses all controllers for I/O and has no passive controllers or idle resources. Because all controllers in a Fluid File System cluster support active I/O, organizations benefit from high intrinsic performance without the need to manually distribute application load across multiple storage controllers. Load balancing sends client requests automatically to the controller with the least-current workload.

5 Data Integrity

In addition to providing outstanding performance in NAS environments, Fluid File System also provides a series of mechanisms to provide a high level of integrity and system resiliency for data at rest and in transit.

5.1 Cache Mirroring and Metadata

During normal operation, the write cache, which includes the data and the metadata, is mirrored between the controller pairs in the Fluid File System cluster. Additionally, important metadata is journaled to back-end storage capacity by a journaling process that runs continuously. This ensures file system consistency, in the event of a failure.

When a controller fails, cache mirroring is not possible, and the surviving controller journals the data and metadata to the storage presented by the back-end storage subsystem. This ensures that the file system remains consistent and that all data is protected in case of additional failures.

In the event of a power outage, the write cache is journaled to a temporary staging location within the controller. This ensures that all I/O in flight remains consistent and that there is no data loss. In addition, this ensures data consistency is not compromised, irrespective of the duration of the outage, as seen by traditional battery backed systems, which will lose data if power is lost before the write completes.

These mechanisms provide resiliency and redundancy for a multitude of failure scenarios to ensure that data is always intact and accessible.

5.2 High availability

In a Fluid File System cluster, any single controller can fail without affecting data availability or causing data loss — even if write operations were in flight. Cross-cluster reliability is achieved through a variety of mechanisms, including a high speed cluster interconnect, write cache mirroring, failsafe journaling, and data integrity checks to ensure data store consistency.

Fluid File System monitors the health of the server platform, including temperature and power condition, to ensure cluster reliability and maximize data availability in cases of hardware or software failures. If failures occur, hardware components in the storage subsystem are redundant and hot-swappable.

Each controller receives its power from the power grid and a dedicated backup power supply (BPS), which is regularly monitored to ensure that the BPS maintains a minimum level of power for normal operation. The BPS has sufficient battery power to allow the controllers to execute their shutdown procedures and use the cache as NVRAM. The BPS also provides enough time to write all the data from the cache to disk.

6 Data Protection

Dell Fluid File System enables data protection within a single system, across systems, and to external NAS repositories. This section will discuss some of the features and benefits of using Fluid File System to store and protect your data.

6.1 Snapshots

Snapshots provide the first level of data protection by providing the ability to recover data instantly. Dell Fluid File System provides the ability take point-in-time snapshots of the entire NAS volume with no impact to user access. Each NAS volume has its own snapshot policy to allow greater granularity.

Fluid File System incorporates redirect-on-write snapshots, instead of the copy-on-write solutions typical of other file systems. Redirect-on-write requires only one I/O operation, thereby preventing performance degradation. Figure 8 shows the redirect-on-write mechanism and the different stages of the snapshot process.

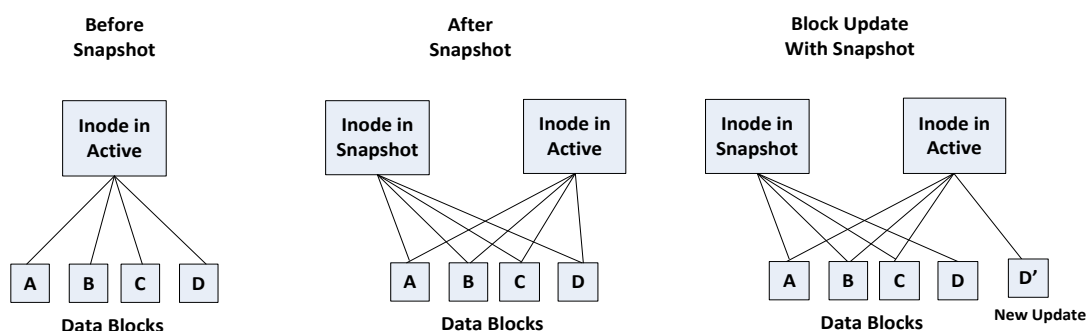


Figure 8 Redirect-on-Write Snapshot Process

Snapshots are available to users as a read-only copy of the file system, which allows them to restore their documents in a simple manner without helpdesk or administrator intervention. Administrators also can easily restore very large data sets (terabyte scale) as a whole to a particular point in time. This eliminates long file copies and the need for free space for the recovery process.

6.2 Backups

Dell Fluid File System supports standard backup software using Network Data Management Protocol (NDMP) with no changes required to existing backup workflows. Dell has partnered with industry leaders to provide comprehensive backup solutions that integrate with Fluid File System. Currently supported backup software includes:

- Symantec™ BackupExec™
- Symantec NetBackup™
- CommVault® Simpana®

See www.dell.com for any other backup solutions that are supported on Fluid File System.

6.3 Replication

Fluid File System allows fast and reliable snapshot-based replication of any number of volumes to a partner. After the initial synchronization, only incremental changes are replicated, which improves network bandwidth utilization. This replication is native to Fluid File System and does not require any additional hardware. The data is always consistent on the partner site and available as read-only.

In addition to data, NAS configurations (volumes, exports, etc.) are replicated. This reduces administrative burden and enables continuous access to data in the case of a disaster or site failure to assure business continuity.

Replication is bi-directional, meaning that the same system can host both source and destination volumes. In addition, the direction can be reversed without requiring a full resynchronization. Fluid File System also supports “one-to-many” and “many-to-one” replication between NAS systems using unique volumes.

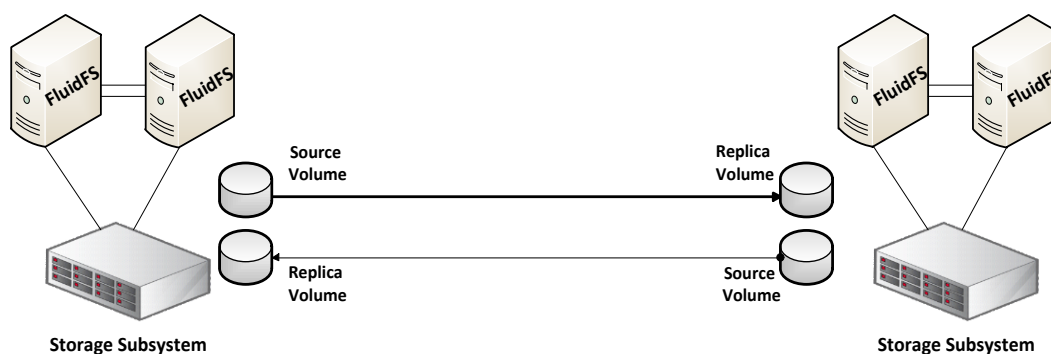


Figure 9 Bi-directional Replication Between Clusters

7 Product Integration

Fluid File System is being implemented as a core component in a number of Dell storage solutions that serve the needs of different customer scenarios and data workloads. Each solution has unique feature, function and value propositions, and each is differentiated by the type of controllers being used in the Fluid File System cluster and the architecture of the underlying block-based back-end storage subsystem. This approach to design, coupled with the distributed and clustered architecture of Fluid File System, results in solutions that can easily leverage future technology enhancements.

The current products that incorporate Fluid File System technology are:

- Dell PowerVault™ NX3500 — A unified storage platform that provides an easy-to-manage and cost-effective solution for both file- and block-based applications.
- Dell EqualLogic™ FS7500 — A unified storage platform that offers high performance scale-out capability. The grid architecture native to both Fluid File System and EqualLogic storage systems provides linear performance scalability in line with capacity growth.

More information about these solutions is located in the Appendix.

Appendix A Fluid File System Solutions

A.1.1 Dell PowerVault NX3500

The PowerVault NX3500 is the first in a series of products based upon Dell Fluid File System that delivers enterprise-class file services to Microsoft® Windows® and Linux clients. It works with PowerVault MD32x0i and MD36x0i storage arrays, providing affordable unified storage with iSCSI, CIFS and NFS access to block and file data. The NX3500 lowers the barriers posed by traditional clustered file system implementations by reducing deployment complexity and offering clustered file systems benefits such as high availability, load balancing etc.

Organizations can use the PowerVault NX3500 to consolidate user data as well as other file and block applications into a single, easy-to-manage unified storage system with best-of-breed data management and scaling capabilities. The PowerVault NX3500's scale-up architecture delivers a flexible, load-balanced pool of high performance storage, making it easy to grow capacity up while avoiding the scalability constraints and challenges of managing separate block and file systems. With dual active-active file controllers and backup power supply, the PowerVault NX3500 gives you data protection and excellent performance with no single point of failure. More information about the NX3500 is available at Dell.com/NX3500.



Dell PowerVault NX3500 and Dell Fluid File System Technical Specifications

Feature	Max Value (2-controller system)
Max system size	576 TB
Max file size	4 TB
Max files	~32 billion
Number of directories	~34 billion
Max NAS volumes	512
Max snapshots per volume	512
Max snapshots per NX3500 system	10,000
Memory per NX3500 system	24 GB (12 GB) per controller
Max LUNs	40
File name length	255 bytes
Max NFS mounts	1024
Max CIFS shares	1024
Max Quota rules per NX3500 system (user quotas)	65,536
Max quota rules per volume	256
Max block level replication policies	256
Max directory depth	1,024

A.1.2 Dell EqualLogic FS7500

The EqualLogic FS7500 is a high-performance solution that enables organizations to easily configure and manage iSCSI, CIFS, and NFS storage from a single interface. Its unique, Fluid File System-based architecture lets organizations scale both capacity and performance and pay as they grow. As storage needs grow and change, block and file capacity can be modified without disrupting existing applications and storage systems. A single file system can be expanded up to the capacity of the EqualLogic back end (up to 509TB usable storage). NAS service can be configured and added to EqualLogic arrays that have been deployed quickly and efficiently. The EqualLogic FS7500 includes a file-based snapshot capability (separate from iSCSI snapshots). Users can restore previous versions of files from a directory of these snapshots themselves, without contacting IT.

A dual active/active controller architecture and sizable onboard cache give the EqualLogic FS7500 outstanding performance. Each controller contains 24GB memory protected by a backup power supply. The EqualLogic FS7500 supports all new and existing EqualLogic arrays running a current version of the EqualLogic firmware. A dual active/active controller architecture and sizable onboard cache give the EqualLogic FS7500 outstanding performance. Each system provides 48GB of battery protected cache, and traffic is automatically load balanced across all controllers.

The EqualLogic FS7500 supports all new and existing EqualLogic arrays. A single FS7500 system can support up to eight EqualLogic PS Arrays with the ability to add another FS7500 system into the same namespace to improve file performance. As with all Dell EqualLogic products, the FS7500's features, software licensing and future firmware enhancements are included in the base price.

Dell EqualLogic FS7500 and Dell Fluid File System Technical Specifications

Feature	Dell EqualLogic FS7500 with Dell Fluid File System
Max system size	509 TB
Max file size	4 TB
Max files	~64 billion
Number of directories	~34 billion
Max NAS file systems	256 per 2-controller FS7500 system, 512 per 4-controller FS7500 system
Max snapshots per NAS File system	512
Max snapshots	10,000 per 2-controller system or 4-controller solution
Memory per FS7500 2-controller system	48 GB/24 GB per controller
File name length	255 bytes
Max NFS mounts	1024 per 2-controller FS7500 system, 2048 per 4-controller FS7500 solution
Max CIFS shares	1024 per 2-controller FS7500 system, 2048 per 4-controller FS7500 solution
Max Quota rules per FS7500 system (user quotas)	100,000
Max quota rules per volume	512
Max directory depth	512