



The Network Appliance Enterprise Storage Architecture: System and Data Availability

Michael J. Marchi & Andy Watson | Network Appliance | TR-3065

TECHNICAL REPORT

Network Appliance, a pioneer and industry leader in data storage technology, helps organizations understand and meet complex technical challenges with advanced storage solutions and global data management strategies.

Table of Contents

1. Overview	3
2. System Availability	3
3. Data Availability	8
4. Summary	12
5. Appendix	12
6. References	13

1. Overview

Information is a critical business asset and customers today require continuous availability to data. Enterprise storage solutions must furnish a high degree of protection for corporate data, provide near- continuous data access, and incorporate procedures to correct problems with minimal business impact.

The Network Appliance™ mission is to enable continuous data access throughout the enterprise. Network Appliance accomplishes this in two ways:

- System availability of greater than 99.99 percent.
- Unparalleled data availability and recoverability.

Most vendors today only focus on system availability. While extremely important, this does not address data availability emergencies that can occur. Some examples include:

- Data corruption occurring within an application (e.g., Oracle®).
- UNIX® user simultaneously edits file being read and/or written by Windows® user, causing file corruption or crashed Windows application.
- Major software upgrades failing or corrupting data.
- Critical file(s) being accidentally deleted or incorrectly modified.
- Natural or man-made disaster(s) (ex. flood, fire, earthquake, etc.).

The normal recovery mechanism for such emergencies is to recover a previous instance of the data set by reloading from tape. This means that data is unavailable during the recovery period. This is unacceptable in today's business environment. Thus, this paper discusses both system and data availability.

2. System Availability

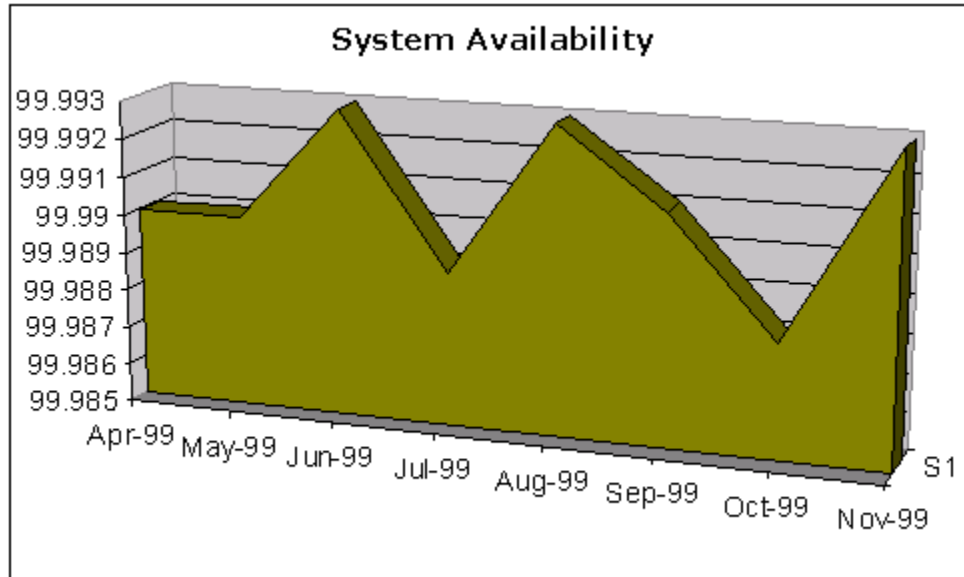
Availability is typically measured as a percentage of total uptime available over the course of a year. For example, a 99.99 percent availability requirement translates into 53 minutes of downtime, whereas a 99.9 percent availability requirement means 8.8 hours of downtime per year. Downtime can include planned maintenance and unscheduled outages.

Availability Classification	Level of Availability (%)	Annual Downtime
Continuous Processing	100%	0 minutes
Fault Tolerant	99.999%	5 minutes
Fault Resilient	99.99%	53 minutes
High Availability	99.9%	8.8 hours
Normal Commercial Availability (single node)	99 - 99.5%	87.6 - 43.8 hours

Source: *Dataquest Perspective - High Availability*

By the above standards, the Network Appliance Enterprise Storage Architecture provides a "fault resilient" level of system availability. This is demonstrated with average system availability statistics as measured across the Network Appliance installed based of over 9,000 storage appliances which is shown to be on average 99.99%.

Network Appliance Inc.



Average system availability measured across Network Appliance installed base of storage appliances for fiscal year 2000.

To better understand how Network Appliance achieves such high average system availability it is important to understand the major causes of system failures.

IT managers say the major causes of system failure are, in order of frequency:

Major Cause of System Failures (in order of frequency)
Software defects/failures
Planned administrative downtime
Operator error
Hardware outage/maintenance
Building/site disaster (fire)
Metropolitan disaster (storm, flood, etc...)

Source: The Gartner Group

Network Appliance effectively achieves "fault resilient" system availability by excelling in each of the areas listed above. The appliance approach improves reliability because it performs a single function very well. A general purpose computer has many features and applications which make it impossible to test all possible usage patterns. Network Appliance storage appliances can be tested much more thoroughly because they only do one thing.

The Network Appliance Storage Architecture is driven by a robust, tightly-coupled, multi-tasking, real-time microkernel (Data ONTAP™ software). This pre-tuned compact kernel minimizes complexity and improves reliability. In fact, Data ONTAP software is less than 2% the total size of general purpose operating systems.

The Network Appliance approach also helps improve overall application availability in that file system operations normally run on general purpose application file servers are no longer executed improving general application server availability. This is a clear differentiation when compared to conventional storage subsystems. In these examples, the odds of application server downtime are increased due to the 100% dependency on the application server's OS and file system software for all I/O operations. This contrasts significantly with Network Appliance deployment options, which allow for multiple application servers, such that the failure of any one of those application servers does not preclude the other application server(s) from accessing the data. This is an added benefit not measured in Network Appliance "fault resilient" availability.

Network Appliance storage appliances utilize proven high-volume, industry-standard hardware components, which help drive high hardware reliability. The most common components to fail are disk drives, followed by power supplies and fans. Network Appliance storage appliances utilize redundant disks (RAID), power supplies and fans for system units and shelves to protect customers against these common component failures.

The Data ONTAP kernel utilizes the robust WAFL™ (Write Anywhere File Layout) file system. WAFL and RAID were designed together to avoid the performance problems that most file systems experience with RAID and to ensure the highest level of reliability. RAID is integrated into the WAFL file system as opposed to other approaches (ex. Veritas Volume Manager sitting on top of Solaris) which eliminates operator errors, OS and application software release mismatches, patch level mismatches, etc.

Network Appliance storage appliances use RAID-4 parity protection for all data stored in the disk subsystem. In the event that any disk drive fails, the data on the failed drive is reconstructed to a global "hot spare" disk drive. While reconstruction occurs, requests for data from the failed disk are served by reconstructing the data "on the fly" with no interruption in file service. RAID-4 provides the benefit of dynamic file system and RAID group expansion with just a single command.

Some vendors will state that running RAID-4, RAID-5 or RAID-S is unsafe due to the fact that a double-disk failure within a single RAID group will cause data loss. This is often stated for one of two reasons:

- The vendor's RAID-5 or RAID-S performance is significantly slower than Network Appliance RAID-4 performance.
- The vendor wants to sell RAID 0+1 (striping + mirroring) which is twice the number of disks.

While it is true that a double-disk failure within a single RAID group will cause data loss, one must look at the probability of such an event happening. The mean time to data loss (MTTDL) for a single RAID group containing N data disks can be calculated according to the equation [Gibson1992] enclosed in the Appendix.

Let's examine four cases using 84 18-GB disks (1.512 TB raw capacity) connected to the Network Appliance high-end F760 storage appliance:

- F760, 84 18-GB disks, 1 volume, twelve 7 drive RAID groups.
- F760, 84 18-GB disks, 1 volume, six 14 drive RAID groups.
- F760, 84 18-GB disks, 3 volumes, four 7 drive RAID groups per volume.
- F760, 84 18-GB disks, 3 volumes, two 14 drive RAID groups per volume.

# Volumes	# RAID Groups per Volume	# Disks per RAID Group	MTTDL in Years Per Volume	MTTDL in Years for Entire Storage Appliance
1	12	7	348,460	348,460
1	6	14	99,560	99,560
3	4	7	1,045,378	348,460
3	2	14	298,679	99,560

Based on these results one can see that the odds of a double-disk failure are extremely rare.

RAID protection means that the chance of a double disk failure where data might be lost is measured in terms of tens of thousands of years. After months or years of use, a few blocks on a disk will go bad. Therefore, a few media read errors on disk blocks over time is normal for a disk. ONTAP will retry reading a disk if there is a media error. If the read error persists, ONTAP will do the following: recalculate the data by reading other disks in the RAID group, reap the bad block to another area of the disk, and store the correct data in the remapped block. As files are read and re-read over time, ONTAP will automatically reap blocks as necessary. Some files are not read for months at a time, or are never re-read. In this case the areas of the disk where these files exist may never be accessed. ONTAP has a special feature, RAID scrubbing, which forces a read on every disk block. Even if a user never reads a given file, NetApp ensures that it will be read by RAID scrubbing which forces a read on every disk block. If a media error is detected the block will be remapped. This avoids the situation where a disk block could go bad over a period of months or years and normally not be repaired.

To protect against this scenario, Network Appliance storage appliances routinely verify all data stored in the file system using RAID scrubbing - something not provided by other vendors. By default, this occurs once per week, early on Sunday morning, though this can be rescheduled or suppressed altogether. During this process, all data blocks are read in parallel. If a media error is encountered, the bad block is recomputed and the data is rewritten to a spare block.

The WAFL file system uses non-volatile RAM (NVRAM) to keep a log of NFS requests it has processed since the last consistency point. (NVRAM is special memory with batteries that allow it to store data even when system power is off.) If, for some reason, a power failure were to occur, this would force an unclean shutdown. When the storage appliance reboots following a power failure, the storage appliance finds the current consistent state on the disks and replays the outstanding requests from the log. The file system can ignore any writes that were in progress when the storage appliance lost power because it knows the written blocks are unallocated in the last consistent image.

AutoSupport is a proactive service provided with Network Appliance storage appliances to ensure the highest level of system and data availability. Storage appliances will periodically send email messages to the Network Appliance Global Support Center (GSC) regarding the health and operation of a customer's storage appliance. The messages cover a wide variety of events and factors. The GSC staff reviews and archives all AutoSupport messages. This is done to check history on specific systems, or to scan for other types of symptoms. Selected types of AutoSupport messages automatically result in the creation of a case by Network Appliance on behalf of the customer. The GSC staff will work on the case just as if the customer had initiated a call themselves. The hours of coverage and response times are the same for cases created from AutoSupport and from the customer directly.

The Network Appliance storage appliance hardware itself provides several features to enhance system and data availability. For example, it has a watchdog timer to detect certain software

failures, environmental monitoring, low mean time to repair and redundancy in failure-prone components such as main memory (ECC), disks (RAID), power supplies and fans. The robust Data ONTAP software, is based on a simple, message-passing kernel that has fewer failure modes than general purpose operating systems. These features combine to demonstrate average system availability greater than 99.99 percent.

While the measured availability on average is "fault resilient," the storage appliance does not tolerate failures of main system components, such as the system board. To eliminate a single point of failure, Network Appliance offers Cluster Failover as an option. Cluster Failover provides hardware redundancy without adding complexity.

Data availability is also affected by planned downtime. Planned downtime typically occurs at predetermined times, such as once a month or quarter. With conventional file servers and storage subsystems, time must be planned for activities which include: backup, software maintenance, hardware maintenance, application/database upgrade(s), operating system upgrade(s), etc... The storage appliance architecture minimizes, and almost completely eliminates, the need for planned downtime. Let's walk through some causes of planned downtime and examine how the Network Appliance Enterprise Storage Architecture addresses each of these:

- Strategic Research Corporation states that the average annual aggregate user productivity loss due to disk grooming (for the purpose of increasing file system size on server storage) is 3,129 hours. Conventional storage subsystems require complex reconfiguration or scheduled downtime to accomplish file system and RAID group expansion. Network Appliance addresses this requirement by allowing system administrators to add disk storage and dynamically expand file systems and RAID groups with a single command and with zero downtime. Also, logical partitions and shares within file systems can be dynamically expanded with zero downtime.
- Data management tasks, such as the spreading of data out across spindles to get better performance, is automatically managed by the Network Appliance storage appliance. It can be thought of as a self-tuning automobile in this respect requiring minimal intervention, if any.
- Operating system upgrades normally require hours of downtime. With Network Appliance, upgrades are installed while the system is operational and serving data. A simple reboot of the storage appliance is then scheduled at your preferred time, which only takes about 90 seconds.
- In many cases data must be taken offline to ensure a safe, consistent backup. Or, a particular application must be put into "hot backup mode" which affects overall system performance while the backup is occurring. With Network Appliance, a safe, consistent backup can be achieved from a Snapshot™ copy of the file system with zero downtime (Snapshots discussed further in the Data Availability section). Specific applications, such as a RDBMS, can be put in "hot backup mode" for a few seconds while the Snapshot copy is taken to ensure 100% data availability. Many competitive offerings require that the data be replicated prior to being backed up. Network Appliance believes this is not cost effective or simple to manage.

Operator error often translates into unplanned downtime, which is more disruptive than planned downtime. The Network Appliance Enterprise Storage Architecture greatly minimizes the chance of operator error given that there are simply few tasks that ever need to be done. For example, RAID is built into the file system so no setup or ongoing configuration is required. No data management involvement, other than expanding file systems, need be done. User authentication

is done via a Windows NT® Primary Domain Controller or NIS server. For those few administrative tasks, an easy-to-use Web-based interface is provided. For Windows administrators, the storage appliance takes advantage of Windows NT User and Server Manager.

3. Data Availability

The Network Appliance Enterprise Storage Architecture is designed to provide a level of data availability never before seen in the industry. Most vendors only focus on system availability, which Network Appliance believes is not satisfactory given today's business environment. Increased data availability translates into higher revenue, profit and productivity. Historically, organizations have paid significant premiums for hardware and software to obtain high system availability, which often does not produce increased data availability.

3.1. The SecureShare® Feature: Data Integrity and Protection

Many organizations are struggling with fast growing UNIX and Windows data requirements. Most engineering organizations utilize both Windows and UNIX systems for software development, chip design, etc. Most IT organizations utilize both Windows and UNIX application servers and would like the ability to consolidate this data to ensure that data integrity is protected. The Network Appliance Enterprise Storage Architecture allows organizations to consolidate both Windows and UNIX data within the same file system for true data sharing with data integrity never previously possible using traditional, uniprotocol-centric file service technologies.

SecureShare software is the Network Appliance patent pending multiprotocol file sharing technology, which provides:

- Kernel-integrated cross-platform lock enforcement.
- Performance-optimized cross-protocol Windows oplock management.
- Cross-platform data integrity.
- Cross-platform Change-Notify.

The SecureShare feature enables UNIX and Windows-based applications to concurrently access and update shared files, with the integrity and cache coherency of the shared data being protected by Network Appliance enforced locking and file-open semantics. The SecureShare feature coordinates and manages all lock types used with multiprotocol file sharing.

For example, if a Windows user is editing a file and a UNIX user accesses the file using NFS (assuming this user has read/write access to the file) the UNIX user will be granted read-only access to the file. Thus, the UNIX user cannot edit, change and corrupt the file while it is being edited by a Windows user.

Data consolidation with true data sharing offers tremendous flexibility to an IT organization. For example, let's assume an IT organization today utilizes Sun Microsystems Solaris-based HTTP servers accessing Web content stored on a Network Appliance storage appliance. One year from now the IT organization decides to move from Sun Solaris HTTP servers to Windows NT HTTP servers. With Network Appliance SecureShare technology, the Windows NT HTTP servers can dynamically be added to the environment with zero downtime and true data sharing. Conventional storage subsystems require that the Web content be replicated to a different volume to enable access from the Windows NT HTTP servers. This requires that not only the data be replicated, but also that two sets of data continuously be updated as web content changes within an organization.

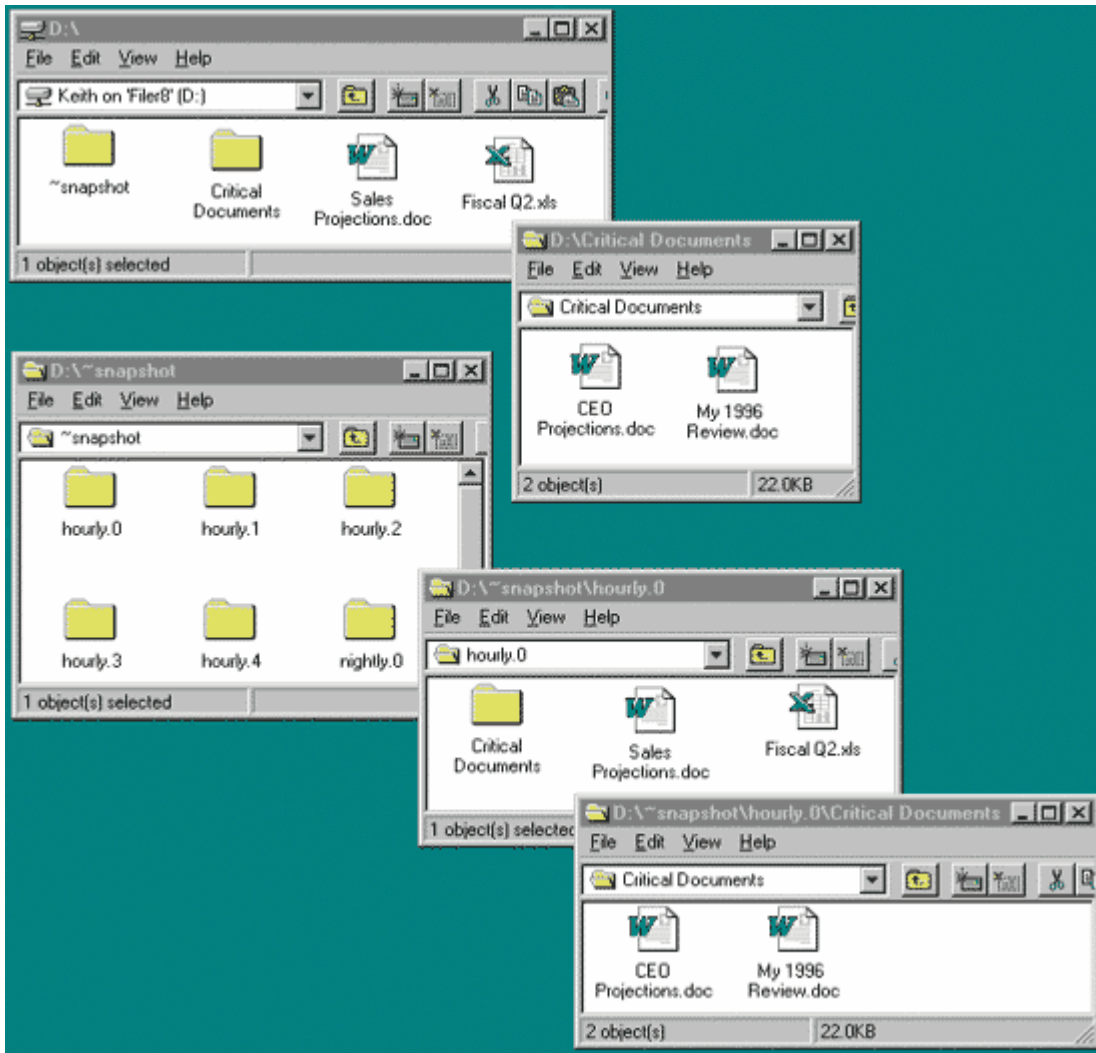
In summary, SecureShare provides a level of data integrity and protection not before seen.

3.2. The Snapshot Feature: Instantaneous File Recovery

Accidental deletion of a critical file(s) usually results in a user calling an IT Helpdesk and requesting that a system administrator restore the file(s) from tape. This result is a productivity loss for the user until the file is restored. It also requires that a system administrator spend valuable time going to the data center, loading a tape and retrieving the file from tape. This is not uncommon.

In fact, a recent study by Strategic Research states that the average site restores files from tape 144 times per year.

The Snapshot technology enables users to instantaneously recover accidentally deleted files without having to call an IT Helpdesk. This results in productivity gains for the user and less strain on already stressed IT organizations.



Example of file recovery utilizing Snapshot technology

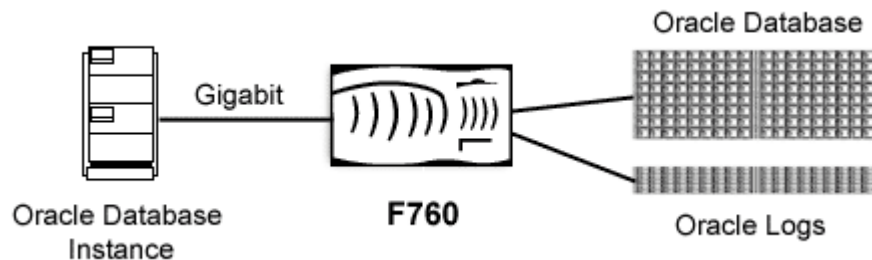
3.3. SnapRestore™ Software: Instantaneous File System Recovery

Time to recovery has become an important measurement in many IT organizations. Some dire situations (e.g., the discovery of corruption in a database) may require the full restoration of a previously saved state leaving data unavailable for an extended period of time. Strategic Research reports that the average site does two full file system restorations per year. SnapRestore software allows a file system to revert back to a previous point in time providing a level of data availability never before seen.

The Snapshot feature allows a file system to be frozen in time. SnapRestore software allows a file system to revert to the state and contents of a previous Snapshot. The system administrator may select any of the up-to-twenty existing Snapshot copies to revert the file system back to.

Some examples of how IT organizations can benefit from SnapRestore technology include:

Data corruption within a database (e.g., Oracle, Sybase, Informix, etc.) is discovered. Normal recovery mechanisms require restoring the damaged portion of the database from tape. SnapRestore software eliminates this time consuming option. SnapRestore software enables the IT administrator to quickly revert the file system back to a previous state when the database was consistent. The log files are replayed and users are again accessing data. Time to recovery is now three minutes (to revert the file system back) plus the log replay time. Compare this to reloading the entire damaged portion of the database from tape which could take a day or longer.



Many Network Appliance customers store application binaries on our storage appliances so that upgrades and patch updates are done in one place as opposed to having the binaries installed on each individual application server requiring each to be upgraded and patched. With SnapRestore software, IT organizations now see an even greater benefit in utilizing Network Appliance storage appliances.

For example, let's assume an IT organization is doing a major application software upgrade (i.e., Oracle, Microsoft Exchange, Informix, Rational ClearCase, Cadence, etc.) and something goes wrong. Upon completing the upgrade the data conversion does not work causing data corruption. Or, things simply don't work for one reason or another. With SnapRestore software, the previous environment can quickly be restored within three minutes without having to reinstall the previous release of software and data from tape. This is a level of data availability that translates into higher revenue, profit and productivity.

3.4. SnapMirror® Software: Cost Effective Automated File System Replication

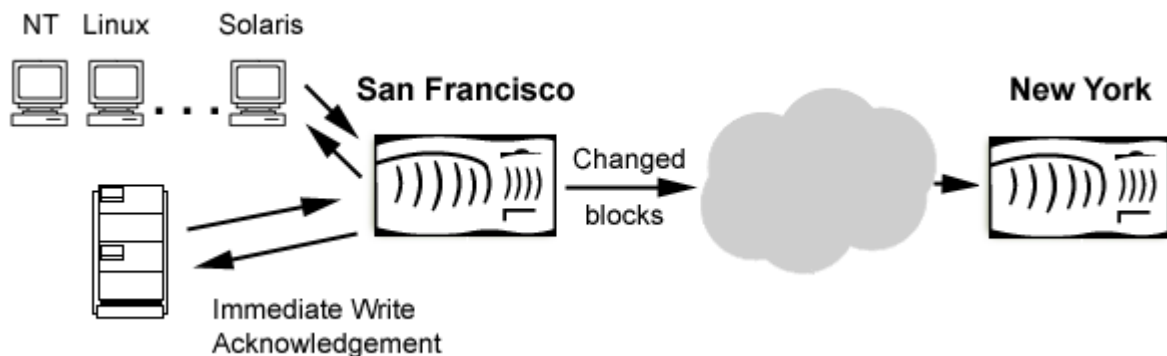
Most IT organizations have business continuance teams in place to help plan for natural or man-made disaster(s) (ex. flood, fire, earthquake, etc.). Most IT organizations today archive data to tape and send the tapes to an offsite location. However, time recovery in the case of disaster is days which is unacceptable in a world where data must be accessed 7x24x365. Thus, many companies are planning for and deploying real time data replication technology.

Network Appliance Inc.

SnapMirror software leverages the WAFL Snapshot capability to provide an automated file system replication facility. Using SnapMirror technology, a storage appliance can replicate one or more file systems to a partner storage appliance, keeping the target file system synchronized with Snapshot copies that are created automatically on the source file system. The target of a SnapMirror replication scenario can be located virtually any distance from the source. It can be in the same building as the source storage appliance or on the other side of the world.

SnapMirror software offers two key advantages over conventional replication products:

- Because SnapMirror software leverages the WAFL file system design, only changed 4-KB blocks are sent from the source to destination. This is significantly less overhead than other offerings which must replicate at an entire disk track due to them having no knowledge of the file system.
- No impact on performance. Write operations reaching the source storage appliance are acknowledged immediately providing sub-10ms response time to end users.



Because time to recovery has become such an important measurement, many IT organizations are looking at different application data sets within their enterprise and determining the minimum time to recovery requirements for each data set. For example, lack of data access to ERP data at the end of a quarter could easily result in revenue impact of \$100,000 or more per hour to a company, whereas the impact of an individual not being able to get to a non-critical file in his/her home directory would be minimal.

Another example would be an E-Commerce site where company revenue is tied to online transactions. Lack of data access for online transactions could result in millions of dollars of lost revenue per hour. However, data measuring individual customer buying habits stored in a data warehouse for marketing research purposes would not result in millions of dollars of lost revenue per hour.

Many IT organizations are choosing to deploy the Network Appliance SnapMirror technology to automatically replicate their most critical data and keep it online where time to recovery is critical. The following diagram illustrates a data replication solution for mirroring the most critical data from one storage appliance and replicating it online using SnapMirror software to a second storage appliance.

File Setup Wizard - Network Services Page 3 of 6

DNS Domain: If your site uses Domain Name Service (DNS), enter the domain name. If your site does not use DNS, leave this field blank. ?

DNS Servers: If your site uses DNS, enter the IP address for one or two DNS servers that the filer should use for host name lookups. If your site does not use DNS, leave these fields blank. ?

NIS Domain: If your site uses Network Information Service (NIS), enter the NIS domain name. If your site does not use NIS, leave this field blank. ?

Routing Gateway: Enter the name or IP address of the primary gateway to use for routing outbound network traffic. If the filer does not need routing to reach its clients, leave this field blank. ?

Email Gateway: Enter the name or IP address of a computer running the Simple Mail Transport Protocol (SMTP) that the filer can use for mail about problems it encounters. ?

Cancel <- Back Next ->

4. Summary

The Network Appliance mission is to enable continuous data access throughout the enterprise. Network Appliance accomplishes this in two ways:

- System availability of greater than 99.99 percent.
- Unparalleled data availability and recoverability.

Data availability - especially for disaster recovery or rapid return to an uncorrupted database condition - can be vastly enhanced by the use of Network Appliance storage appliances within an enterprise. Given the Network Appliance Enterprise Storage Architecture, traditional data management practices should be re-examined. In many cases they may no longer be necessary at all.

5. Appendix

The mean time to data loss for a single RAID group containing N data disks can be calculated according to the following equation from [Gibson1992]:

$$MTTDL = \frac{MTBF^2}{N \times (N + 1) \times MTTR}$$

where

- MTTDL = the Mean Time To Data Loss (i.e., double-disk failure in a RAID group).
 MTBF = the Mean Time Between Failure of any disk.
 MTTR = the Mean Time To Repair (i.e., reconstruct the failed disk onto a spare).
 N = the Number of data disks in the RAID group.

Network Appliance storage appliances support multiple RAID groups within single volumes. Most customers today utilize 7 or 14 drive RAID groups within large volumes. The formula changes to reflect this, where

- G = the number of RAID groups (requiring an additional G-1 parity disks, one per RAID group).

Using this formula, the probability of data loss becomes approximately:

$$\text{MTTDL} = \frac{\text{MTBF}^2}{N \times (N / G + 1) \times \text{MTTR}_G}$$

MTTR is smaller in this case because reconstruction time is a function of the amount of data in the RAID group. Experiments conducted on the following RAID group sizes shows the approximate MTTR values as:

F760: 84 18 GB disk drives using 7 drive RAID groups.

$$\text{MTTR} \sim 3.425 \times (6 + 1) + 15 = 62.95 \text{ (~39 minutes)} = 0.65 \text{ hours}$$

F760: 84 18 GB disk drives using 14 drive RAID groups.

$$\text{MTTR} \sim 3.425 \times (13 + 1) + 15 = 62.95 \text{ (~63 minutes)} = 1.05 \text{ hours}$$

6. References

[Gibson1992]

Garth A. Gibson, Redundant Disk Arrays: Reliable, Parallel Secondary Storage, MIT Press, Cambridge, Massachusetts, 1992.