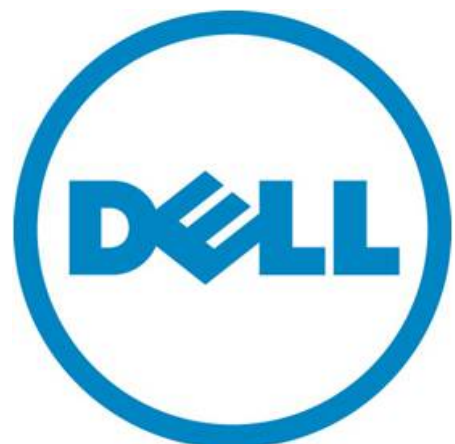


Dell PowerVault MD3600f/MD3620f Design Guide for Disaster Recovery Using Remote Replication



THIS 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.

© 2011 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, *PowerConnect*, 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. 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.

June 2011

Abstract

This guide explains how to design a disaster recovery solution using the Remote Replication feature of Dell MD36X0F Storage Array. Various topics that must be considered before designing a robust disaster recovery solution are discussed. In addition, guidelines, best practices and rules of thumb are provided.

Contents

Abstract	iii
Dell PowerVault MD3600F/MD3620F (MD36X0F) Remote Replication Overview	3
Benefits of PowerVault MD Remote Replication	3
Design Considerations	4
What is the motivation to implement data replication?	4
Achieve Compliance	4
Centralized Backup	4
Disaster Recovery	4
What is the cost of implementing a DR solution?	5
Storage	5
SAN	5
LAN	5
WAN	6
Software Licenses	6
What is a realistic implementation time?	6
How much storage capacity is needed?	6
Prioritize and characterize the data	7
Is the recovery site data center equipped for Remote Replication?	7
What is the distance between primary and recovery sites?	7
How much data link bandwidth do I need?	8
Is failback part of the solution?	8
What about backup procedures?	9
Does the DR solution proposal meet both current and future storage needs?	9
Which replication mode should be used?	9
Implementation Guidelines	10
MD36X0F Storage Array Sizing	10
Sizing the primary storage array	10
Sizing the Recovery Storage array	11
Network Sizing	12
Gather I/O information	12
Latency	12
Bandwidth	14
Buffer Credits	14

- Other Network Considerations..... 14
- Application Considerations..... 14
 - I/O Block Size..... 14
 - File System/Database Layout..... 15
 - Temporary/Scratch Files and Tablespaces 15
- Enhance your Disaster Recovery Solution 15
 - Implementing Snapshot and Virtual Disk Copy in a DR solution 16
 - Snapshot..... 16
 - Virtual Disk Copy..... 16
 - Backups..... 17
- Testing Considerations 17
 - Performance Testing 17
 - Testing for Failover and Failback 17
- Services 18
- Conclusion 18
- Appendix A Remote Replication on Databases for Disaster Recovery 19
 - What needs to be replicated in a database?..... 19
 - Method 1—Replicate Everything..... 19
 - Method 2—Replicate Log Files Only..... 19
 - The Consistency Group 20
 - Other Database Considerations 20
- Appendix B Disaster Recovery Terminology 22

Figures

- Figure 1. Distance between sites..... 13

Dell PowerVault MD3600F/MD3620F (MD36X0F) Remote Replication Overview

Remote Replication (RR) is the MD36x0f array feature that provides the ability to remotely move data from one MD36x0f storage array to another. It is an optional premium feature that requires activation, configuration, and in normal replication operations, administration and control; the latter two tasks do not require lots of time and manpower resources and they can be automated through scripting.

Benefits of PowerVault MD Remote Replication

- **SAN-based replication** - Data replication between the primary virtual disk and the secondary virtual disk is managed by the MD36X0f storage arrays, and is transparent to host machines and applications.
- **Replication modes** - RR supports *synchronous*, *asynchronous*, and *asynchronous with write order consistency* replication modes. These modes help administrators choose the replication method that best meets protection, distance or performance requirements.
- **Dynamic replication mode switching without suspending the replication** - Users can switch between replication modes at any time. This enables administrators to accommodate changing application and bandwidth requirements without sacrificing protection.
- **Ease of use** - IT administrators can enable replication through the Modular Disk Management System (MDSM) GUI or MDSM CLI. Once enabled, RR functions can be managed using the MDSM GUI, the CLI, or MDSM vCenter plug-in. These tools are available free of charge.
- **Multiple replication relationships** - A replication connection is not limited to a single primary RR system and a single remote RR system. In general, each MD36X0f with the optional RR premium feature installed and activated may be either a primary system or a secondary system, or both. Up to 16 replication relationships are supported.
- **Suspend and resume** - RR provides the ability to suspend replication by explicitly using the Suspend command, and under certain circumstances to automatically suspend replication (consistency group suspension, asynchronous replication when link bandwidth exceeded).
- **Role reversal/fail over** - RR allows reversing the roles of primary and secondary virtual disks to recover from a disaster.
- **Read-only replica access (includes Snapshot creation)** - RR enables the remote data to be utilized prior to a disaster without sacrificing protection of the primary site data.

For more information on the functional aspects of Dell PowerVault MD36X0F Remote Replication, refer to the *Dell PowerVault MD36X0F Remote Replication Functional Guide*. The rest of the paper covers the factors that need to be considered before implementing a Remote Replication solution.

NOTE: Refer to Appendix B for definitions of key terms used in this paper.

Design Considerations

There are some key critical design issues that need to be addressed before designing a disaster recovery (DR) solution with Remote Replication:

- What is the motivation to implement Remote Replication?
- How much will it cost to implement a Remote Replication solution? What infrastructure do I need?
- What are the recovery objectives—Recovery Point Objectives (RPOs) and Recovery time Objectives (RTOs)—and Business Continuity plans should a disaster occur?
- How will my decisions affect future growth and current backup procedures?
- What is the recommended distance between the primary and secondary sites?
- How much data must be replicated? Should you replicate all the data, or should you select only a portion of the data for replication?

What is the motivation to implement data replication?

There are three valid reasons to replicate data (keep in mind that data replication is only one component of a larger solution), but each has a different set of requirements and integration tasks.

- Achieve Compliance
- Centralized Backup
- Disaster Recovery

Achieve Compliance

Many compliance initiatives are driven by legislation, such as HIPPA and Sarbanes/Oxley (SOX) in the United States and Basel (BAL) in Europe; consequently they have high executive priority due to increased accountability.

Both SOX and HIPPA compliance include a requirement for timely disaster recovery. Depending on the size of the company and the complexity of its operations, this requirement could be satisfied either by having regular tape back-ups maintained offsite, or by real-time replication of transaction records and no more than a few minutes of downtime after a disaster.

Centralized Backup

Remote Replication can be a component of a centralized backup solution. If a business operates in several locations but only has one backup facility available, Remote Replication is a viable technique to centralize backup.

Disaster Recovery

A centralized backup facility can be used as a disaster recovery site as well, but additional measures beyond replication for backup must be taken into account to implement a true disaster recovery capability.

NOTE: If remote data replication is not required, using Snapshot and Virtual Disk Copy to create physical point-in-time copies within one storage array may be an alternative to replication between storage arrays. These replication tools are particularly useful for application testing, data migration and data mining.

What is the cost of implementing a DR solution?

There are a number of factors that contribute to the overall cost of a DR solution, such as storage, distance, SAN, LAN, WAN etc. The most important factor is the distance data is to be moved. For a given data capacity, the cost of the solution increases with the distance between primary MD36X0F storage arrays and the secondary MD36X0F storage array used to replicate data.

The distance requirements depend on the level of protection. For example, if the data is replicated to a recovery location in the same building, the solution will provide protection from failures of the primary MD36X0F, the computer room where the other MD36X0F is located, and the building floor which houses the computer room. This approach may very well protect against the majority of failures the organization is likely to encounter.

On the other hand, if the data is replicated to another building in the same city, the solution will provide additional protection against a disaster involving the entire building where the primary MD36X0F is located. But this solution will cost more because the link between the two systems will cost more. Likewise, if the recovery site is located in a separate and distant city, it will provide even more protection, but the cost will increase accordingly. The solution design should explore this issue, and Dell can provide valuable assistance in evaluating the tradeoff between distance-related costs and level of protection.

Storage

The total cost of a DR solution includes a significant storage component. Not only does the number of storage arrays required at least double (and the storage consumed within them), but there is an increase in infrastructure resources required for the SANs as well. MD36X0F premium features also add to the overall cost (Remote Replication is required for each storage array; Snapshot and Virtual Disk Copy are recommended for both).

SAN

A SAN is required at both the primary and the recovery site.

While these SANs may already be in place, special consideration should be given to their design to ensure that there is no single point of failure between the replicated storage arrays. This may mean that additional switches are needed. While cost-cutting at the recovery site might be a temptation, the DR solution will not deliver needed results if the recovery site is not configured properly—especially if it will remain online for an extended period of time. As discussed earlier, it is possible that the primary site will not be available for failback for an extended time after a disaster, and maybe never.

Also, if a failure takes the recovery site replication capability down and a primary site failure or disaster occurs before replication is resumed, there is no way the DR solution can meet its RPOs and RTOs. The only way to eliminate this very real exposure is to assure that the recovery site is configured for high availability—comparable to the primary site. Completing an extensive BIA will help to catch and address potential issues like this.

LAN

To design a truly resilient solution, additional components may be needed for the LAN. Again, this may involve additional switches, routers and/or hubs. Additional management capabilities may be needed as well to ensure the LAN is properly managed and monitored.

WAN

This is typically the most expensive component of the solution—primarily because it is a recurring cost. *While WAN costs can be reduced by implementing the solution with a lower bandwidth than is recommended by your Dell representative, the solution will not perform as designed and will not deliver expected results, thus exposing the entire project to failure and associated sunken costs.* It is absolutely critical to support the solution with the recommended network bandwidth.

Software Licenses

Since the purpose of a DR solution is to facilitate rapid failover of applications to a remote recovery site, the organization may need to obtain software licenses to support the failover process. This may add cost to the solution.

What is a realistic implementation time?

Planning, designing, implementing and testing a DR solution is not something that can happen in a week. A generous amount of time should be allocated just for measurements and testing alone. To reduce the complexity and speed implementation time of the overall project, Dell recommends that the DR solution be divided into phases such as those suggested below, and managed one phase at a time.

Phase 1—The first phase of data replication is to provide appropriate data protection. By itself, data replication is not disaster recovery, but it is a crucial step to achieving disaster recovery, and can be implemented in a relatively short amount of time, especially if professional services are employed. Further, it can provide other benefits, such as enhancing compliance and moving toward centralizing backup.

Phase 2—Beyond data protection, disaster recovery requires that host servers failover to connect to the data in the recovery site so that people connect to these servers over networks to operate them.

Phase 3—Time and labor must also be spent to test the complete solution.

How much storage capacity is needed?

This is a critical question. The answer depends on:

- How much data is kept online today, and how fast it is expected to grow?
- How much of the current data needs to be replicated?
- The number of virtual disk pairs and images (copies of virtual disks) that need to be maintained.

Keep in mind that each virtual disk and each physical image requires as much storage capacity as the source data virtual disk. In contrast, a snapshot typically requires only a fraction of additional capacity. Snapshot should be considered where appropriate to create logical images in order to reduce the cost of storage capacity

It is necessary to size both the primary and recovery site storage arrays for initial implementation—accounting for all source virtual disks and images—and then to add additional capacity to each system for future data growth. Failure to anticipate growth can cause the disaster recovery solution to become obsolete quickly. It is important to understand that disk capacities and associated performance requirements often grow at a rate of 10% to 15% on a quarterly basis!

Prioritize and characterize the data

Prioritizing and characterizing the data can help make storage capacity decisions easier and more effective. Prioritize the data to be replicated by RPO and RTO. For most organizations, different types of data have different levels of business value, and for DR purposes, they should have different recovery point objectives and different recovery time objectives.

Types of business data typically range from simple user files and directories to complex database tables, logs and indices. Each data type has its own set of requirements for integrity and proper recovery. It is important to understand the types of data to be replicated and the associated requirements for successful recovery.

Assuming all data has equal value and trying to replicate equally for disaster recovery purposes is usually not realistic. Most businesses find that setting several different RPOs and RTOs is appropriate. Data replication should be prioritized by data value (which translates to RPO) versus time to recover (RTO). This is where the business impact analysis (BIA) becomes very important.

It is important that diagrams of data structures, data schemas and data relationships be documented at the recovery site so they can be quickly accessed after a disaster (if necessary). These diagrams must include all the components of files, directories and databases that are to be protected with the DR solution. In addition, diagrams are required that detail, both logically and physically, how the data is laid out within the storage array and how it is connected to the servers and applications.

Understanding the current data layout may help to uncover potential problems before recovery is needed, and will reduce the likelihood of related issues arising unexpectedly during or after recovery.

If any of the relationships are going to change, or if a data migration is planned, this information should be factored into the DR solution plan so adequate time and resources can be allocated for these tasks.

Is the recovery site data center equipped for Remote Replication?

A key step in designing a DR solution is to make an assessment of the recovery site data center. A successful DR solution cannot tolerate surprises when recovery is attempted. Some relevant things to consider include:

- What equipment is available?
- Is the equipment adequate to handle the recovery and resume production of targeted applications after they fail over?
- Does the equipment need to be upgraded?
- What kind of performance can be expected when recovery is finished?
- What relevant software is licensed to the site?
- Does the software need to be upgraded?
- What power and cooling facilities are available?
- How will required personnel get to the recovery site?
- What facilities are available to service other personnel who are not onsite?
- Will the network have to be shared, lowering its effective speed?

What is the distance between primary and recovery sites?

One of the most important steps in designing a DR solution is planning for the distance that replicated data must travel between sites. With longer distances, the laws of physics become an issue and the

technical challenges of remote replication become more complex. For example, if the distance is less than a mile, the infrastructure will be completely different than one that covers fifty miles. Implementation and operating costs will be significantly different as well. Understanding and addressing these challenges is crucial to the design and implementation of a successful solution.

For some industries, such as banking, government regulations govern the minimum distance between sites. For others, the cost of providing adequate network bandwidth is the primary consideration. In some cases, special region-specific conditions should be taken into account as well. For example, in southern California the recovery site should be on a different tectonic plate than the primary site. But in London the recovery site may only need to be blocks away to enable recovery from a terrorist attack or a fire. The potential for flooding and the need to put the sites on different power grids are two other factors that should be considered.

How much data link bandwidth do I need?

Another important factor in a successful DR solution is the bandwidth requirement for data replication. Without adequate bandwidth, the disaster recovery solution will fail.

The data link bandwidth requirement encompasses the amount of data that needs to be replicated as well as the link speed required to provide acceptable performance. Since the amount of data that will be replicated will likely grow at a hefty rate, the link speed must be dynamically adaptable to meet future requirements. Some questions that should be asked are:

- What is the minimum bandwidth recommended by Dell to provide adequate performance for the distance between sites and the amount of data that is to be replicated initially?
- What is the expected rate of growth of the replicated data?
- What minimum bandwidth does Dell recommend for the expected data growth?
- If the data is being replicated to an existing site, how much bandwidth is available now? Is it being shared? If so, what portion is available for this replication? In the future?
- Are there plans to increase the bandwidth? If not, what is required to get it planned?

If the data link drops for any reason, replication is suspended until the link is re-established. This requires that the virtual disk pairs be resynchronized (fully replicated again) before they can be used for recovery, which may take a significant amount of time, based on a number of factors. A good DR plan will provide enough network bandwidth to minimize resynchronization time.

Is failback part of the solution?

The goal of a DR solution is to provide failover of selected applications to a recovery site that is not critically affected by the disaster. To restore full business operations after the disaster, failback to the primary site or failover to a third site (depending on the status of the primary site) may be required. This is an important consideration that should be factored into the DR solution design.

Failback introduces new requirements for the DR solution design that affect both implementation and testing of the solution. For example, failback adds the potential of losing transactions beyond the disaster recovery window. Any data with an RPO less than zero are subject to some amount of loss during disaster recovery failover *and* during failback. Failback also causes a second window of application downtime since it takes a finite amount of time to restore any application. The organization must be prepared to handle these realities.

The simplest and preferred way to implement failback is to reverse the pre-disaster process of replication. In cases where a planned failover is needed (e.g., evacuating an area that is in the path of an approaching hurricane), shutting down the servers at the primary site, suspending the replication, and then reversing the direction of replication is an acceptable way to failover quickly and facilitate rapid failback.

Failback is more complicated if failover is triggered by an unplanned outage, a complete copy of the recovery database must be restored at the primary site before resuming operations. This is because servers at the primary site will likely be running when the failure occurs, so the two images of the database will not be synchronized. There is no practical way to merge the two images and guarantee data integrity—hence one site or the other must be designated as the master and the other must have its data overwritten with the master image.

A major concern is the amount of time it will take to perform the failback replication which, in some cases, could take days or weeks—especially if the database has increased its size substantially since the failover occurred. Regular reviews of failback time requirements will ensure that failback time objectives are met. Reducing the time to failback may require physically moving a fresh backup copy from the recovery site to the primary site. Once the database has been restored, replication the logs and applying them is accomplished in a much shorter timeframe.

There is no guarantee that the original primary site will survive the disaster such that failback can actually be carried out. Therefore, a complete DR solution should plan for the possibility of never resuming operations at the original primary site.

What about backup procedures?

A DR solution will certainly affect backup procedures for the data involved in disaster recovery, and may have even broader reach. To get a handle on this, it is a good idea to characterize all data into backup classes that are either local or remote, then determine by data type if it is desirable to maintain local backup or move to remote. This analysis will help size required storage arrays and contribute to the final DR solution design.

In most cases it is preferable to do remote backups so the backup images are not stored at the primary site where they can be rendered useless by contingencies and disasters. Remote backup is a requirement for most DR scenarios. Archives using tape or other permanent media will still be needed, and should be part of the overall DR solution.

Does the DR solution proposal meet both current and future storage needs?

It is imperative to keep current and future storage needs in mind when designing a DR solution. Some questions to ask are:

- Is the data already stored on a MD36X0F storage array? Will it have to be migrated from another storage platform? Are firmware levels up to date?
- What about future storage needs? Most likely growth will require additional capacity, but how much? Will additional MD36X0F storage arrays be required in the future?
- To keep costs down, is it better to start with small MD36X0Fs and upgrade each to larger configurations as growth takes place?

Which replication mode should be used?

The answer to this often confusing question is largely determined by the distance between the primary and recovery sites. MD36X0F synchronous replication is recommended for relatively short distances to

provide maximum throughput. It is appropriate when the data link runs at fibre channel speeds. If Synchronous Replication distances exceed the standard fibre channel limit of 10 km (6.2 miles), Remote Replication requires additional equipment certification to extend beyond the fibre channel boundaries. If the distance between sites exceeds synchronous replication distances, asynchronous replication or asynchronous replication with write order consistency is the only answer.

The advantage of asynchronous replication and asynchronous replication with write order consistency is enabling longer distances and/or lower link speeds for the DR solution. However, the data link for asynchronous replication and asynchronous replication with write order consistency must provide sufficient bandwidth and I/O rate to support a successful Remote Replication implementation. Also, they must be faithfully monitored and managed to ensure they are running smoothly and efficiently.

A unique feature of MD36X0F Remote Replication is that replication modes can be switched dynamically—from synchronous replication to asynchronous replication and then to asynchronous replication with write order consistency, and vice versa. With this feature, all the modes can be tested for a performance comparison (assuming synchronous replication is even a possibility).

Implementation Guidelines

MD36X0F Storage Array Sizing

Sizing the primary storage array

Size the primary storage array as usual to meet production capacity and performance requirements, but then modify the sizing and configuration to address the following:

- Double the write I/O rates of the virtual disks that will be replicated. Replication requires the MD36X0F to write every I/O twice (to itself and to a secondary system at the recovery site), and this takes additional resources within the storage array. The additional resources will help to compensate for the performance impact that will occur when virtual disk synchronization is being performed.
- Increase the read rate of all virtual disks by 25% to account for replication latencies, synchronization and overhead. The data transfer rate, replication priority, and network latencies all influence this, but 25% is a reasonable rule of thumb.
- MD36X0F Remote Replication reserves the last port of each controller for its own use. This reduces the overall bandwidth available between the MD36X0F and the SAN. If the last port is already being used, it must be reconfigured to support Remote Replication and its associated virtual disks must be reassigned to another port.

Above all, ensure that the storage array is sized to provide enough horsepower and spindles to sustain expected performance growth requirements as well as the desired remote replication performance.

If other MD36X0F replication features are used while remote replication is in progress, they will also impact performance and should be taken into account when sizing the primary storage array. The same holds true for host-based data replication.

The replication repositories for Remote Replication should not be placed on the same array as the data being replicated. This is a best practice. They should be located on a high performance/highly available array that is less busy.

Sizing the Recovery Storage array

The remote storage array presents a more complex sizing effort than the primary storage array. Some questions to ask about the recovery storage array include:

- Is it to be used for more than recovery purposes?
- What is the performance expectation in the event of failover?
- Are more or fewer applications going to be running on it (compared to the primary storage array) after failover?
- How will post-failover backup impact it?
- Will it become a new primary storage array for remote replication to another site (e.g., to facilitate failback)?

If the recovery virtual disks are expected to perform as well as the primary virtual disks, then the same sizing for production work will apply to the remote storage array as the primary. Don't forget that one port on each controller will be reserved for Remote Replication, and don't forget to compensate for performance impacts caused by MD36X0F data replication features (as mentioned in the Primary Storage array section above).

As a best practice, the capacity sizing for production work at the recovery site needs to be multiplied by a factor of at least 3.2. This is to provide additional capacity as follows:

- 1X is for the continuous replication of production virtual disks that takes place prior to failover.
- 2X provides additional capacity for a second copy of point-in-time virtual disks to enable database roll forward capability in case a corruption of the database at the primary site is replicated to the recovery site.
- 3X adds more capacity to replicate the source virtual disks for testing or online backup.
- 0.2X provides incremental capacity to enable Snapshot for the replicated virtual disks.
- 3.2X enables all of the above (i.e., it provides capacity for the Remote Replication, a point-in-time copy for fall-back, space to copy the database for testing or backup, and capacity for Snapshot repositories).

In other words, if 10 TB of production data at the primary site is to be replicated and put into production at the recovery site after failover, the recovery storage array should have at least 32 TB of capacity—with enough spindles to meet performance expectations.

It is not unusual to have a smaller storage array at the recovery site as some of the applications running at the primary side may not be necessary in the event of a site failure. This is why it is so important to determine what data will be replicated and what applications will actually need to run at the remote site.

If, after failover, the recovery site is to replicate data to a remote site for failback or to cascade replication, then the recovery storage array must be sized accordingly. It will, in fact, function as a primary storage array in the new data replication scheme, and therefore must have its sizing altered according to the rules of thumb provided in the Primary Storage array section above.

If the data stored on the primary system is being replicated to a bunker site instead of a recovery site, performance is not as important as price. This is the ONLY case where SATA drives may be considered for remote replication. A bunker site is typically used for data mining or centralized backup, and the appropriate sizing rules apply.

Network Sizing

One of the most important aspects of designing a successful DR solution is properly sizing the network that connects the primary and recovery sites. If bandwidth is inadequate, required throughput will never be achieved. And if latency, which increases with distance, is excessive, RPO may not be obtainable. Furthermore, if the cost of obtaining the required bandwidth is excessive, the solution may not be approved by executive management.

One very important matter to consider in designing the network is the amount of time that will be required to perform initial synchronization as well as ongoing resynchronizations. Additional bandwidth may be needed to get virtual disk pairs synchronized in a reasonable timeframe. If the network is to be shared, sizing for synchronization is especially important because Remote Replication could potentially use all available bandwidth or, worst-case, not have enough bandwidth.

When sizing bandwidth, it is best to use actual performance measurements.

In some cases, it may be appropriate to ship a copy of the database to the recovery site and load it there rather than attempting synchronization of the entire database via replication. This approach may provide faster synchronization at the recovery site, and will reduce network bandwidth requirements. To accomplish this, the database is first transferred to the recovery site.

Once the database has been loaded at the recovery site, the logs are replicated from the primary site to the recovery site, and then the logs are applied to the database through a roll forward process. As long as logs continue to be replicated and applied, the database stays current with changes to the primary site database. Since only the logs are replicated, and not the database, there is a net reduction in network bandwidth requirement.

Gather I/O information

Determine I/O information on a virtual disk basis as this information is required to determine the minimum bandwidth required for successful remote replication. It is also useful for setting performance expectations if the recovery site is anticipated to deliver lower performance for any reason (network constraints, slower servers, older storage, etc.).

There are several utilities available to collect this information, including IOSTAT, PERFMON, and MD Storage Manager. Required I/O information per virtual disk includes:

- I/Os per second
- Read/write ratio
- I/O block size
- Variances between normal and peak workloads
- Total byte counts (read and write)

Ideally these measurements will be taken over a long period of time and include minimum, maximum, and average for each virtual disk to be replicated.

Latency

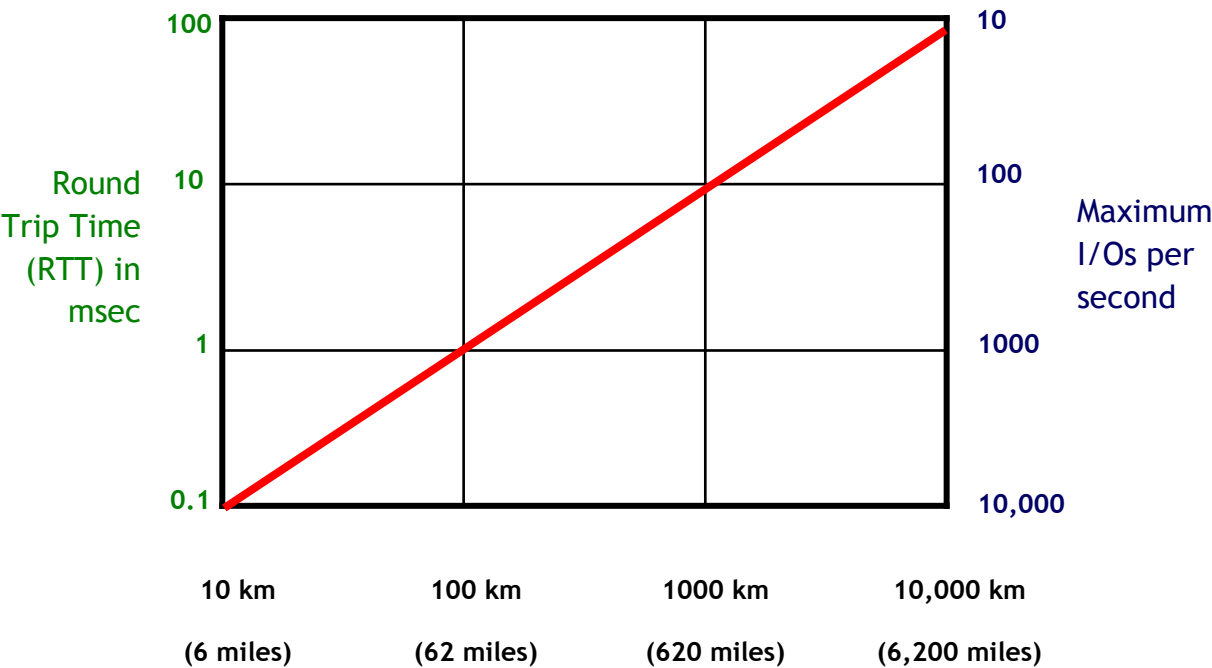
Latency is the time it takes a replication I/O and its acknowledgment to traverse a network link. The longer the distance between primary and recovery sites, the more time it takes to send the data and receive an acknowledgement. Distance becomes a limiting factor for the number of I/Os that can be

sent per second. *Therefore, latency is the controlling factor for I/O rate that can be supported in a disaster recovery solution.*

To illustrate, if it takes 2 milliseconds (msec) to send an I/O to a recovery site 125 miles (200 km) apart over one link and receive an acknowledgment (this is called round trip time or RTT), then a maximum of 500 I/Os can be processed in one second using that link ($500 \times 2 \text{ msec} = 1 \text{ second}$). If distance between sites increases to 1000 km (620 miles), then RTT increases to 10 msec and the I/O rate for that link is reduced to 100 I/Os per second. In these examples, additional latency that results from the communications infrastructure (switches, routers, firewalls, etc.) is not included.

The following figure illustrates how round trip time (RTT) increases while best-case maximum I/O rate decreases with distance between sites involved in a disaster recovery solution.

Figure 1. Distance between sites



If the link has sufficient bandwidth available and if the application to be replicated allows independent I/Os, then the Remote Replication I/Os may be multiplexed across the link, which will increase the effective I/O rate of the solution. However, databases typically require synchronization of transactions for consistency; therefore, the benefits of multiplexing may be limited when replicating databases.

At the MD36X0F level, asynchronous replication with write order consistency is used to sequence I/Os to the recovery site in the original database I/O order. The telecommunications equipment connecting the two sites must have in-order packet delivery turned on to ensure write order consistency.

Note: Consult with your telecommunications provider for more detail on this critical topic.

Bandwidth

Network bandwidth defines how much data can be sent through the data link. Note that if the data being transmitted has been converted to TCP/IP, there is a significant increase in overhead that reduces effective bandwidth—up to 50% for low link speeds!

Buffer Credits

Buffer credits help to optimize replication performance for extended distance SANs. Buffer credits allow multiple data frames to be in flight simultaneously over a single fibre channel link.

Specifying the correct number of buffer credits ensures that the fibre channel link is running at optimal efficiency in terms of data delivery. When longer distances are involved, the appropriate number of buffer credits becomes important in order to prevent a networking term called *drooping*. Drooping occurs when data cannot be sent because the sending system is waiting for the opportunity to send more data through the link. For fibre channel, drooping can be caused by not having enough buffer credits.

The optimal number of buffer credits is a function of bandwidth and distance. A 2 KB fibre channel frame (the standard frame size) traveling at the speed of light will be approximately 2 km in length when using a 2 Gb link. That same frame will be 4 km long using a 1Gb link and 1 km when using a 4 Gb link. Frame length is used to calculate the number of buffer credits required to optimize replication performance (divide distance between sites by the frame length).

Note: Consult with your SAN extension provider for more detail on buffer credits.

Other Network Considerations

Since the network is so important to the overall solution, it is important that all required information be obtained before the design is finalized. A detailed SAN/LAN/WAN diagram should be developed so any discrepancies or problems can be discovered early on during implementation.

Required network information includes:

- Switches involved (number, layout, type)
- Actual data path to be used for the interconnect between storage arrays
- Will the WAN/LAN connection be dedicated or shared?
- How much of the available bandwidth will be used by other applications or users?
- Are there any planned changes that could affect available bandwidth?
- What network monitoring tools are available? Can they be made available during implementation?
- How will performance be monitored?
- How will a network failure affect replication? Will remote replication be suspended? If so, how long?

Application Considerations

I/O Block Size—As a rule of thumb, make I/O block size as large as possible while keeping performance acceptable. Databases and file systems are designed with standard block sizes, which result from balancing the size of a block with server memory caching efficiency. Databases offer the ability to change block size, but the change is a disruptive process.

Another performance best practice is to ensure that data blocks line up on a storage segment boundary. This is normally not an issue for UNIX, but it is frequently encountered in Microsoft Windows

environments. For more information on this topic, a Microsoft document titled: “How to Align Exchange I/O with Storage Track Boundaries” is located at:

<http://www.microsoft.com/technet/prodtechnol/exchange/guides/StoragePerformance/0e24eb22-fbd5-4536-9cb4-2bd8e98806e7.mspx>.

File System/Database Layout

Applications that make use of a file system for file access may have additional overhead that should be considered when designing and building the DR solution. Specifically, a Journaled File System (JFS) will require additional I/Os since an I/O is staged in the file system journal before it is written to the main area on disk. A JFS may result in increased cost for the DR solution if it requires additional network bandwidth. *Whenever possible, avoid file systems and use raw devices to reduce this overhead as much as possible.* This may not be important with smaller databases, but it becomes increasingly important as databases grow.

Remember that partitions are not the same as virtual disks. In Windows environments, it is important to optimize the placement of partitions on virtual disks.

Temporary/Scratch Files and Tablespaces

Whenever possible, temporary/scratch files and temporary/scratch tablespaces should be assigned to a separate virtual disk that is NOT replicated. These temporary data holders are typically used for reports or data conversions and are not part of the data that needs to be replicated to the recovery site. If replicated, they will translate to additional I/Os (and possibly additional bandwidth cost) with no advantage in the event of a failover. The database or system administrator should be able to determine if this condition exists and remedy it.

For the same reasons, swap file is another storage object that should not be replicated.

Please refer to [Appendix A](#) for recommendations on Implementing Remote Replication for Databases.

Enhance your Disaster Recovery Solution

If the database needs to be recovered quickly at the primary site in the event of database corruption, a quick recovery solution should be included with the broader DR solution. A quick recovery solution typically invokes Snapshot and Virtual Disk Copy.

Some of the questions that need to be asked when considering quick recovery are:

- What are a reasonable RPO and RTO for quick recovery from database corruption?
- Can these objectives be achieved using Snapshot and Virtual Disk Copy?
- Can the MD36X0F storage array support the additional overhead and performance impact?

A quick recovery solution typically follows the following sequence of events:

1. The database is placed into hot backup mode.
2. A Snapshot is made of the database.
3. The database is put back into normal mode.
4. A Virtual Disk Copy is made of the Snapshot.
5. When the Virtual Disk Copy is complete, the Snapshot is deleted.

In the event of database corruption, the database is dismounted by the database administrator and the Virtual Disk Copy image is mounted in its place. In effect, the database is recovered to the point in time when the Snapshot was taken. Log files are used to roll the database forward and reassert database integrity.

Implementing Snapshot and Virtual Disk Copy in a DR solution

Snapshot

A Snapshot can provide—in seconds—a point-in-time copy that typically requires only a fraction (the default is 20%) of the capacity of the data being copied. This is very useful for non-disruptive backups because the database need only be quiesced (put in a consistent state) for a relatively short period of time before normal operations can resume.

Snapshots are also useful in multi-step updates where point-in-time copies are created at each step. Should a subsequent update step fail, an earlier Snapshot is used to restore the last known good recovery point. Prior to disk-based Snapshots, the only way to keep a point-in-time image was to back up the entire database to tape and keep the image around in case a restore was required. This required a significant amount of time (to back up the data and verify its integrity) and restores took hours. Tape-based point-in-time copies are no longer effective for this purpose (however, they remain useful for archiving).

A Snapshot of the database should be maintained at the recovery site to ensure that there is always a viable image there. *Best practice is to use a rolling pair of Snapshots where a second Snapshot is made before resynchronizing the databases.* Once the databases have been synchronized, the older Snapshot is deleted to improve overall storage array performance.

If Snapshot is to be integrated into the DR solution (which is *highly* recommended), there can be a storage array performance impact while a Snapshot is being made and is subsequently active. However, the impact may be reduced by adding drives to the array configuration. *A rule of thumb is to add 15% overhead for reads and 25% overhead for writes to the database.* If the Snapshot is made at the recovery side, these numbers should have minimal effect on the overall solution (unless the recovery storage array is being used for other work at the recovery site).

The performance impact of Snapshot must be considered when sizing the storage arrays.

Virtual Disk Copy

Virtual Disk Copy is another MD36X0F premium feature that can enhance a DR solution. Virtual Disk Copy can make a copy of a virtual disk and make it available as a different virtual disk (a clone). This is different from a Snapshot because the new image is a complete, physical copy that can be manipulated separately from the original data. There are many uses for this new image, including data mining, reporting, testing, and data migration.

While the virtual disk is being copied, the original image must remain stable and unmodified. Because of this, it is a best practice to make a Snapshot of the original database, and then use the Snapshot as the source for the Virtual Disk Copy. Since Virtual Disk Copy reads every block of the original image, there will be an additional impact on performance while the copy is being made and this should be sized when designing the solution.

Backups

Backups are an important component of a DR solution design. Not only are backups a critical part of the business, they will require additional MD36X0F capacity and premium features.

The first question to ask is where the backups will occur and be stored? In many cases, they will be made and maintained at the primary site to provide the ability to recover from primary site failures when failover to the recovery site is not required.

Backups may also be created at the recovery site. Some reasons for choosing this approach are:

- Backups may already be stored offsite from the primary location.
- The additional workload caused by backups (on servers, storage, network and SAN) is removed from the primary site, allowing more consistent and possibly higher performance there.

Testing Considerations

Performance Testing

Verifying the performance impact of replication, along with associated Snapshot and Virtual Disk Copy operations, to support disaster recovery is a very important aspect of the overall solution.

The ability to obtain metrics to verify acceptable performance for both the storage arrays and the application is a critical component of the solution design. These metrics should include:

- Average and maximum time to post a transaction
- Average and maximum network bandwidth utilization
- Average and maximum I/O rates and response times

The solution should be tested in a lab environment before implementation. This urgent requirement must be addressed when designing the overall solution.

Testing for Failover and Failback

Testing failover and failback can be a significant challenge and needs to be included when defining the DR solution. In many cases, executive management will not want to do the testing, primarily because of the time it will take to failback to the primary site. But these procedures must be tested to verify they will work and to develop efficiency enhancements. Simulating failover is fairly easy and can be tested without a lot of effort. However, simulating failback can be a much more complicated task and will require additional equipment—including additional storage capacity or storage arrays—which is why failback is often left out of the overall solution.

If failback is not needed, it is much easier to test the DR solution because failover testing proceeds without taking down the primary systems. Failback testing requires that the primary site be used as a recovery site; hence normal production cannot be continued with the primary systems while failback tests are conducted.

If failback is required, a modified approach may possibly address the testing issue discussed above. The modified approach is to treat failback as a second failover and test it only after the first failover has occurred. In other words, delay testing failback until recovery has been completed at the secondary

site. This involves starting the replication process over again from the secondary site that has now taken on primary status.

The downside of this approach is that it could take longer to failover to the original primary site than the failover/failback approach since the new failover process will have to be implemented and tested first. If resources at the secondary site do not provide adequate performance, this extra time may be unacceptable. Planning ahead for this modified approach will minimize the time required to get the original primary site up and running again.

In order to safely and adequately test failover and failback, multiple servers and large amounts of additional storage are required—which adds to the cost of the overall solution. One way to lessen this cost is to use equipment allocated to another project for the tests and then deploy it for the intended purpose after the planned DR testing exercise has finished.

Keep in mind, however, that DR testing is never complete! There must be an ongoing executive commitment to keep the DR solution current and to do regular testing to ensure that it remains a viable means to business continuity.

Services

Dell offers remote replication implementation service that can reduce the time required to plan, design, implement and test a DR solution. *While this service has a price, it can actually lower the final solution cost by reducing total project time and, more importantly, minimizing the risk of achieving a successful solution.* See your Dell representative for more details.

Conclusion

Dell MD36X0F Remote Replication provides the technical capability to implement a disaster recovery solution. However, it is only one component of an effective solution. The solution must be planned, designed, implemented and tested to address all the topics covered in this paper.

Appendix A Remote Replication on Databases for Disaster Recovery

What needs to be replicated in a database?

There are several ways to replicate a database for disaster recovery. Each method has advantages and disadvantages—understanding them will help in designing the best solution.

The two most popular database replication methods are (1) replicate everything (2) replicate log files only.

Method 1—Replicate Everything

With this approach, the entire database and the log files are replicated. The advantage here is that massive database updates will normally be handled without additional procedures or resources. However, there are several disadvantages; this approach is more susceptible to problems, it leads to a more complex solution, and it requires more network bandwidth.

Typical sequence of events for replication a database and logs:

1. Establish replication for all database virtual disks
2. Verify that virtual disk pairs are fully synchronized (in optimal state)

Then, on a regular basis:

3. Put the primary database into hot backup mode
4. Suspend replication to the recovery site
5. Create a Snapshot of the replicated image at the recovery site
6. Resume replication between sites
7. Exit database hot backup mode at the primary site and resume normal operation
8. Using the Snapshot, backups and data migration tasks can now be performed at the recovery site
9. *Do not delete the Snapshot!*

For details on hot back up mode refer to *Dell PowerVault MD36X0F/MD3620f Remote Replication Functional Guide*.

It is best to always have at least one Snapshot available at the recovery site in case database corruption at the primary site is replicated to the recovery site. *Best practice is to keep several copies of the database at the recovery site for multiple recovery points. Snapshot enables a quick and easy point-in-time process to accomplish this.*

Method 2—Replicate Log Files Only

Using this approach, the entire database is replicated to the recovery storage array initially, and then only the log files are replicated thereafter (until it becomes necessary to replicate the database again). The logs are applied to the database at the remote site. This reduces the bandwidth required for replication, but also requires a server at the recovery site to apply the logs. If a massive database change is made, the entire database should be copied over again—which takes time and additional bandwidth temporarily. One significant advantage of this approach is that the two database images are truly separate from each other. In the event of database corruption at the primary site, the database is still intact at the recovery location.

Typical sequence of events for replication only log files:

1. Establish replication for all database virtual disks
2. Suspend or remove replication of the database when synchronization is complete
3. Verify that virtual disk pairs are fully synchronized (in optimal state)
4. Continue replication the log files

Then, on a regular basis:

5. Put the primary database into hot backup mode
6. Suspend replication to the recovery site
7. Create a Snapshot of the log files at the recovery site
8. Resume replication of log files between sites
9. Exit database hot backup mode at the primary site and resume normal operation
10. Using the Snapshot, apply the log files to the recovery site database

The log Snapshots may or may not be removed once the log files have been applied. *Best practice retains at least one Snapshot image of the log files.*

It is also recommended that the log file space on the primary storage array be able to retain a minimum of 24 hours of logs in case problems occur at the recovery site (such as a failure that prevents the log files from being applied). It is also important to make sure that the log files are not deleted on the primary side until they have been applied at the recovery side. This can be accomplished with scripts.

RTO and RPO will also influence the DR solution design. Does the data need to be synchronized at all times or can images be sent over the data link in batch mode on a regular basis (e.g., hourly or at a shift change)? Can the organization afford to lose an hour or more of work? Must the DR solution provide an RPO within a single transaction? All these questions must be addressed during the RTO/RPO evaluation and answers must be provided in the design.

The Consistency Group

An MD36X0F storage array will support one consistency group. The consistency group is very important when the data to be replicated spans more than a single virtual disk. Designing a DR solution correctly involves determining which virtual disks should be included in the consistency group.

In a database environment, the best rule is to place the database virtual disks and the log files into the consistency group if “replicate everything” is the approach selected. Otherwise, place the log files into the consistency group if “replicate log files only” is chosen.

Other Database Considerations

It may be necessary to turn off replication during major database changes, such as bulk loads, no-log operations and rebuilds. Typically these types of changes require that logging be turned off to speed up the process. If only log files are being replicated for disaster recovery purposes, a new copy of the database must be replicated before log replication is started again. Be sure to optimize the database before starting the replication process (row-chaining, data block size, optimizer, statistics). If the changes are extensive, they may affect every block of data in the database, requiring that every block be replicated. A much more efficient approach is to suspend replication of the database, make the changes, and then full copy the source virtual disks.

Consider multiplexing the log files (if this is supported by the database)—it can reduce the performance impact. Most databases allow for multiple copies of log files. If the log files are multiplexed (two or three copies of the data instead of just one) and only one log image is replicated, the primary system will have nearly the same performance *with* replication as it does without it.

Note: It is important to create Snapshots of the database environment just after a link failure and before Remote Replication resynchronization begins. This will establish a consistent recovery point for the database prior to the link failure. Likewise, Snapshots should also be made after Remote Replication resynchronization completes to establish a new (consistent) recovery point for the database.

Appendix B Disaster Recovery Terminology

Disaster Recovery (DR)

Disaster recovery is the business process of recovering from a disaster, whether that disaster is caused by a manmade or natural event. DR requires the restore, recovery and restart of a business application and related processes after a disaster occurs. Disasters may arise from a multitude of sources, and each should have a documented plan of recovery. Types of disaster vary greatly, depending on geography, business environment and political considerations. A site in the Midwestern United States should plan for disasters caused by tornadoes and flooding, but a site in Hawaii should be concerned with typhoons, earthquakes, volcanic eruptions and tsunamis. Both geographies need to consider extended power loss, employee actions and fire, but individual sites may or may not be concerned about a terrorist attack. Each potential source of disaster needs to be addressed with a documented and tested disaster recovery plan. Every company or organization that wants to survive a disaster should have such a plan in place.

Business Continuity (BC)

Business continuity is the term that describes the ability to maintain normal business operations with minimal or no downtime—either planned or unplanned. Downtime can result from various causes besides disaster. For example, taking a database offline for backup is an example of planned downtime. Having a server crash and take down the order entry system is an example of unplanned downtime. Neither case involves a disaster, but both impact the business operations of an organization. BC can have a variety of implications for an organization, depending on its size and operating model. For example, a small business with ten employees that is only open during normal business hours will have a completely different set of BC standards than a 24x7 international company with thousands of employees and customers. Taking the order entry database offline for backup can be done quite easily by the smaller company during off-hours, but it is completely unacceptable for the international company. This company can use a non-disruptive Snapshot to make an image of the database for backup, ensuring that their BC requirements are met.

Disaster Avoidance (DA)

Disaster avoidance is the term for preventing disasters from occurring. For example, implementing a tight access control system can prevent users from accessing critical data that should not be modified. Likewise, protecting a facility with a fire extinguishing system is a way to reduce or prevent damage by fire. DA techniques like these should be part of every DR strategy. One dimension of DA that needs to be addressed is external factors that may impact the company's BC capability. For example, if an automobile manufacturer uses only one external supplier for brakes, what happens if that supplier shuts down? DA issues are usually discovered during a comprehensive contingency planning exercise.

Recovery Point Objective (RPO)

Recovery point objective is the maximum time interval between a disaster and the last time data was recoverable. It translates directly to the amount of data that can be lost when a disaster occurs. RPO should generally be minimized; however, it may vary by data type, the criticalness, the value, and the cost of maintenance of the data. For example, the component ordering system and the plant operations system of a manufacturing company will usually have more aggressive RPOs than other systems. In contrast, data warehousing applications and user file and print data typically have less aggressive RPOs than order entry systems.

The RPO for a given application may be a day, a shift, an hour or thirty minutes—again depending on the attributes of the data mentioned above. Typically data accessibility is maintained on a sliding time scale. For example, a company may want to keep multiple generations of critical data in the event of a corruption so it can roll back to a previous image, then roll forward to the current point in time. In some of the more critical environments (stock exchanges for example), the RPO needs to be set so that every completed trade is recoverable.

Recovery Time Objective (RTO)

Recovery time objective establishes the maximum amount of time allowed to recover data from a recovery point after a disaster occurs and should generally be minimized. But again, RTO should be set by the data attributes described above under RPO. For example, a mail order business will want to recover its order entry system as quickly as possible (minutes) but let its data warehouse application take up to 24 hours to recover. A determining factor to consider here is the cost of low RTO compared to the cost of not having the data available.

Network Recovery Objective (NRO)

Network recovery objective defines the amount of time it will take to actually get the remote recovery location up and running. This requires examining the overall network environment and determining what the requirements are for a complete switchover to the remote site. While this may not involve storage arrays directly (other than using the network for data replication or replication), it is important to be aware of NRO.

All three of the objectives defined above (RPO, RTO, and NRO) depend on one another and must be considered when designing a complete disaster recovery solution. For example, it is pointless to set a RTO at 30 minutes if the NRO is two hours. It may be appropriate to spend more resources and money to get the NRO down to one hour and spend less on the RTO, allowing it to meet the NRO at one hour. In designing a disaster recovery solution, you should ask two questions for each objective:

If we spend more (or less) on this objective, can we improve our bottom line disaster recovery time? If so, is it worth the additional cost (or savings)?

Contingency Plan (CP)

Contingency planning is the process of exploring a series of “what if?” questions that relate to disaster recovery. Planning for events that “might” or “could” happen is critical to an organization’s business continuity. If a fire breaks out in the main computer room, what happens? How does the fire get put out? How will this affect the operation of the company? What steps can be taken to reduce or eliminate this potential problem and its impact? How much will it cost and how much is it worth to the company?

A comprehensive contingency plan will invariably uncover more weak points than originally expected, but it will help a company continue business operations when more obscure problems are encountered (or at least a plan has been developed to reduce the exposure).

High Availability (HA)

High availability means that a system which enables business operations is able to continue running in the event of a single component failure. If a server is used for a business critical application, it needs to be upgraded to a highly available configuration or it needs to be configured in a server cluster that can take over the application in the event the server fails. Fortunately, MD36X0F storage arrays can be configured to run at different availability levels while providing high performance.

HA is important, not only in the data center, but in other areas of a company as well. For example, if a single source is used to power the manufacturing floor and that source fails, manufacturing is essentially shut down. The same applies to wide area networks—especially with the trend toward globalization. If a single WAN service provider is used, what happens if the connection is lost?

HA is achieved by designing in redundancy of the various components, subsystems and systems that constitute the data center. HA levels range from simple redundancy at the component level (HBAs for example) to duplex systems that protect against system failure. The goal of HA is to eliminate single points of failure so systems will continue to operate through unplanned failures.

Fault Containment and Isolation (FCI)

Fault containment and isolation prevents single failures from causing cascading events that result in even more failures. Using the MD36X0F to illustrate successful implementation of FCI, if a single drive fails in a properly configured system using RAID and switched drives, that failure will not affect any other components in the storage array and the application will continue to run.

The same concept applies throughout the data center. If a single network switch fails, the network should be resilient enough to continue functioning by bypassing the failed switch. Looking at FCI from a higher level, if a failure shuts down the development system, it should not shut down the order entry system as well. Sadly, many of history's notable failures were caused by a series of minor events that, if they had been contained and isolated properly, would never have occurred or would not have been as catastrophic.

Business Impact Analysis (BIA)

A business impact analysis is the output of assessing the impacts of unplanned events and contingencies on the business. It is an assessment that compares the cost of implementing specific DR solutions to the impact of not having a DR solution—considering financial effect, company reputation, lost customers, etc. In short, it is a comparison of the cost of DR against the expected losses. If the order entry system is brought down because of a site failure, how much revenue will be lost (both immediate and long-term) due to lost customers and how much will it cost to prevent such a failure? Site failures are not always caused by information technology. Having a vendor shut down for a week because of a hurricane and consequently running out of a critical manufacturing component is an example of how non-technical external events can affect a company's bottom line.