

# Dell PowerVault ME5 Storage System Best Practices

## Abstract

This white paper highlights best practices for optimizing and deploying PowerVault ME5 (ME5012/5024/5084) and should be used in conjunction with other PowerVault ME5 manuals (Deployment guide, Admin Guide, Support Matrix etc.)

March 2023

## Revisions

Date	Description
March 2023	Initial release

## Acknowledgments

Author: Selim Selveroglu

Support: Dan Jobke

Joseph Catalanotti

Justin Oros

Martin Pritchard

Nigel Hart

Patrick Quah

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

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

This document may contain certain words that are not consistent with Dell's current language guidelines. Dell plans to update the document over subsequent future releases to revise these words accordingly.

This document may contain language from third party content that is not under Dell's control and is not consistent with Dell's current guidelines for Dell's own content. When such third party content is updated by the relevant third parties, this document will be revised accordingly.

Copyright © 27th of March 2023 Dell Inc. or its subsidiaries. All Rights Reserved. Dell Technologies, Dell, EMC, Dell EMC and other trademarks are trademarks of Dell Inc. or its subsidiaries. Other trademarks may be trademarks of their respective owners. [27-Mar-23] [White Paper] [H19551]

# Table of contents

Revisions.....	2
Acknowledgments.....	2
Table of contents .....	3
Executive summary.....	5
Intended Audience .....	5
Prerequisites.....	5
Related Documentation.....	5
Introduction .....	6
System Concepts.....	9
General Best Practices .....	12
Become Familiar with Manuals .....	12
Stay Up-to-date with Firmware.....	12
Always Use Supported Configurations.....	13
Host Information .....	14
Identifying Your Hosts Easily.....	14
How to Monitor Array Health .....	15
Configuring E-mail and SNMP Notifications.....	15
E-Mail Notifications.....	16
SNMP Notifications .....	16
How To Provision Virtual Pools Features.....	18
Thin Provisioning .....	18
Thin Provisioning Space Reclamation .....	20
Block Zeroing.....	20
Pool Balancing.....	21
Quick Rebuild .....	21
Modifying Virtual Volumes.....	21
Best Practices for High Availability .....	22
Volume Mapping.....	22
Direct Attach Cabling Example.....	22
Multipathing Configuration and Multipath Software.....	23
Snapshots.....	23
Dual Power Supplies .....	24
Fault Tolerance (Reverse) Cabling .....	24
SMART .....	25

Scrubbing .....	25
Autonomic Distributed Allocation Protection (ADAPT).....	25
Hot Spares.....	26
Virtual Volume Replication .....	27
<b>Best Practices for Performance .....</b>	<b>28</b>
Workload .....	28
Use SSDs for Randomly Access Data .....	28
PowerSizer .....	29
Dual pools.....	30
Linear vs Virtual.....	30
Magical Number 2 Method .....	31
Disk Groups in a Pool.....	31
Which RAID Level Should You Use? .....	32
ADAPT performance .....	32
Disk Count Per RAID Level .....	33
<b>Optimization of an existing System.....</b>	<b>34</b>
Volume Cache Options .....	34
<b>Automated Tiered Storage.....</b>	<b>36</b>
How Tiering Works? .....	36
Volume Tier Affinity Feature.....	37
Optimization in Tier Setup .....	38
<b>Best Practices for Firmware Updates .....</b>	<b>40</b>
Updating Disk-Drive Firmware .....	41

# Executive summary

This white paper highlights best practices for optimizing and deploying PowerVault ME5 (ME 5012/5024/5084) and should be used in conjunction with other PowerVault ME5 manuals (Deployment guide, Admin Guide, Support Matrix etc.) All manuals available from [Dell Support](#) web page.

---

Note: Features and recommendations in this document reflect current software functionality as of the latest firmware available at the time of publication. Features, functionality and GUI might vary with different storage system firmware levels.

---

## Intended Audience

This best practice document is intended for PowerVault ME5 storage administrators and presales & deployment engineers with previous SAN infrastructure and SAN storage knowledge.

## Prerequisites

Prerequisites for using this product include knowledge of:

- Storage system configuration
- SAN management
- Connectivity methods such as direct attached storage (DAS), Fibre Channel, and serial attached SCSI (SAS)
- Networking
- iSCSI and Ethernet protocols

## Related Documentation

In addition to this guide, other documents or materials for this product include:

- [Dell PowerVault ME5 Series Storage System Administrator's Guide](#)
- [Dell PowerVault ME5 Series Owner's Manual](#)
- [Dell PowerVault ME5 Series Storage System CLI Guide](#)
- [Dell PowerVault ME5 Series Storage System Deployment Guide](#)
- [Dell PowerVault ME5 Series Support Matrix](#)

# Introduction

The PowerVault ME5 (ME5) models referenced in this paper include PowerVault ME5012, ME5024 and ME5084.

ME5 brings the essential features and data services of mid-to-high end storage to small-medium businesses (SMBs). ME5 unleashes substantial performance, capacity and bandwidth gains over its predecessor system without compromising simplicity, features or business outcomes. Designed to meet several SMB data storage challenges, ME5 reduces complexity, drives more performance, supports higher capacity, ensures data protection and achieves budget stability when supporting a diverse SMB workload and application environment.

Designed with dual-active controllers, multi-core processing and five-9's availability, ME5 arrays deliver the performance of flash with the best economics of disk making it an ideal solution for SMB high-value workloads.

Using fast Intel Xeon processors, ME5 storage implements a dual-active controller architecture and delivers, 12GB/sec read and 10GB/sec write throughput and provides a 12Gb SAS backend protocol for rapid capacity expansion. ME5 Series storage implements a block architecture with VMware virtualization integration and concurrent support for native iSCSI, Fibre Channel, and SAS protocols. Additional storage capacity is added via Disk Array Enclosures (DAEs) while Distributed RAID (ADAPT) software delivers faster drive re-build times. And all ME5 Series arrays are managed by an integrated HTML5 web-based GUI – PowerVault Manager.

The two non-dense ME5 base arrays are 2U models and the dense ME5 array is a 5U model. Both models include dual controllers with, 16GB cache memory per controller.

Supported Front-end protocol options per controller are;

- 4x10Gb iSCSI Base-T
- 4x 10/25Gbit iSCSI
- 4x12Gb SAS
- 4x16/32Gbit FC

## PowerVault ME5012



- Up to 5.2 PB capacity
- 12 x 3.5" drive bays
- Up to 264 drives
- Multi-protocols – SAS, iSCSI, Fibre Channel
- Single/Dual Controller
- 12Gb SAS Backend
- All premium software features included

## PowerVault ME5024



- Up to 5.0 PB capacity
- 24 x 2.5" drive bays
- Up to 276 drives
- Multi-protocols – SAS, iSCSI, Fibre Channel
- Single/Dual Controller
- 12Gb SAS Backend
- All premium software features included

## PowerVault ME5084



- Up to 8.0 PB capacity
- 84 x 3.5" drive bays
- Up to 336 drives
- Multi-protocols – SAS, iSCSI, Fibre Channel
- Dual Controller
- 12Gb SAS Backend
- All premium software features included

### All Inclusive Softwares

- **ADAPT (Distributed RAID):** (like Dynamic Disk Pooling) reduces drive rebuild times
- **Thin Provisioning:** Allocate and consume physical storage capacity as needed in disk pools. Thin is virtual mode only.
- **SSD Read Cache:** Increase execution speed of applications by caching previously read data
- **IP & FC Remote Replication:** Safely replicate data to any global location that includes mirroring thin provisioned pools
- **Snapshots:** Easily recover files after accidental deletion or alteration with Redirect on Write snaps
- **3 Level Tiering:** Get great performance with less hardware expense
- **Volume Copy/Clones:** Enable seamless volume relocation and disk-based backup and recovery with a full, replicated copy of source data
- **Encryption (SEDs):** Render data useless to unauthorized users with drive-level encryption, even if the drive has been removed from the enclosure (internal key management included)
- **Virtualization Integrations:** VMware vSphere, vCenter SRM, Microsoft Hyper-V

---

**Note:** Softwares (except of ADAPT) requires virtual storage mode which will be covered in next section.

---



# System Concepts

## Virtual and Linear Storage

This product uses two different storage technologies that share a common user interface. One uses the virtual method while the other one uses the linear method.

**Virtual storage** (system default) is the most common selection and is recommended for most environments. Virtual storage allocates space in pages and allows data to be moved to improve system performance of the storage system. Virtual storage supports thin provisioning (with or without overcommitment), tiering, replication, and many other features not available to linear configurations. However, you cannot exceed 2 PB usable capacity in Virtual Storage.

Virtual storage is a method of mapping logical storage requests to physical storage (disks). It inserts a layer of virtualization such that logical host I/O requests are mapped onto pages of storage. Each page is then mapped onto physical storage. Within each page the mapping is contiguous, but there is no direct relationship between adjacent logical pages and their physical storage.

**Linear storage** is used in applications where performance and data workloads dictate that data be allocated on disks in a contiguous fashion with more predictable performance. Users in streaming media and video editing, High Performance Computing environments may prefer the performance and raw capacity available to a linear storage configuration. Features such as thin provisioning, snapshot, read cache, tiering, and replication are not available in a linear storage environment. Linear storage is similar to thick provisioning.

The linear method maps logical host requests directly to physical storage. In some cases, the mapping is one-to-one, while in most cases the mapping is across groups of physical storage devices, or slices of them. This linear method of mapping is highly efficient. The negative side of linear mapping is lack of flexibility. This makes it difficult to alter the physical layout after it is established.

### Notes:

- Virtual Mode is required for Thin Provisioning, Tiering, SSD Read Cache, Snapshots and Replication
- Both Virtual and Linear mode supports 8 PB raw capacity.
- ADAPT is supported in either Virtual or Linear Mode
- Once you select Virtual or Linear Mode you can't change it online.

### Pools

A pool is an aggregation of one or more disk groups that serves as a container for volumes. Virtual and linear storage systems both use pools. Dual controller systems consist of two pools. Each storage controller has ownership of a one pool.

In both virtual and linear storage, if the owning controller fails, the partner controller assumes temporary ownership of the pool and resources owned by the failed controller. If a fault-tolerant cabling configuration, with appropriate mapping and MPIO, is used to connect the controllers to hosts, LUNs for both controllers are accessible through the partner controller so I/O to volumes can continue without interruption.

## Disk Groups

A disk group is an aggregation of disks of the same type, using a specific RAID level that is incorporated as a component of a pool, for storing volume data. Disk groups are used in both virtual and linear storage environments. You can add virtual, linear, or read-cache disk groups to a pool.

## Volumes

A volume is a logical subdivision of a virtual or linear pool and can be mapped to host-based applications. A mapped volume provides addressable storage to a host (for example, a file system partition you create with your operating system or third-party tools). For more information about mapping, please check [Volume Mapping section](#) of this document.

## Volume groups

You can group a maximum of 1024 volumes (standard volumes, snapshots, or both) into a volume group. Doing so enables you to perform mapping operations for all volumes in a group at once, instead of for each volume individually. A volume can be a member of only one group. All volumes in a group must be in the same virtual pool. A volume group cannot have the same name as another volume group but can have the same name as any volume. A maximum of 256 volume groups can exist per system. If a volume group is being replicated, the maximum number of volumes that can exist in the group is 16.

## Thin Provisioning

Thin provisioning is a virtual storage feature that allows a system administrator to overcommit physical storage resources. This allows the host system to operate as though it has more storage available than is actually allocated to it. When physical resources fill up, the administrator can add physical storage by adding additional disk groups on demand.

## Automated Tiered Storage

Automated Tiered Storage is a virtual storage feature that automatically moves data residing in one class of disks to a more appropriate class of disks based on data access patterns, with no manual configuration necessary:

- Frequently accessed data can move to disks with higher performance.
- Infrequently accessed data can move to disks with lower performance and lower costs.

Each virtual disk group, depending on the type of disks it uses, is automatically assigned to one of the following tiers:

- **Performance**—This highest tier uses SSDs, which provide the best performance but also the highest cost.
- **Standard**—This middle tier uses enterprise-class spinning SAS disks, which provide good performance with mid-level cost and capacity.
- **Archive**—This lowest tier uses midline spinning SAS disks, which provide the lowest performance with the lowest cost and highest capacity.

## SSD read cache

SSD Read cache is a virtual storage feature. Unlike tiering, where a single copy of specific blocks of data resides in either spinning disks or SSDs, the Read Flash Cache (RFC) feature uses one SSD read-cache disk group per pool as a read cache for frequently accessed data only. Each read-cache disk group consists of one or two SSDs with a maximum usable capacity of 4TB. A separate copy of the data is also kept in spinning disks. Read-cache content is lost when a controller restart or failover occurs.

## ADAPT (Autonomic Distributed Allocation Protection)

ADAPT is a RAID-based data protection level that maximizes flexibility, provides built in spare capacity, and allows for very fast rebuilds, large storage pools, and simplified expansion. All disks in the ADAPT disk group must be the same type (enterprise SAS, for example), and in the same tier, but can have different capacities. ADAPT is shown as a RAID level in the management interfaces. For detailed information, please check [ADAPT section](#) of this document

# General Best Practices

This section has some general practices when administering ME5 storage arrays

## Become Familiar with Manuals

This document includes some information from manuals, for become familiar to your array reading all manuals are highly recommended. You can access ME5 public support page and documentation from [here](#)

## Stay Up-to-date with Firmware

For better performance, reliability and gaining new features, it's important to update your storage's firmware regularly. You can find most up to date firmware at <https://www.dell.com/support> web site. You can update your array's firmware in **PowerVault Manager -> Maintenance > Firmware > System..** For more details please check Admin Guide or this guide's "[Best Practices for firmware update](#)" section.

**DELL EMC** PowerVault Manager | ME5084

Dashboard  
Provisioning  
Settings  
Maintenance  
Storage  
Hardware  
Firmware  
About  
Support

### Firmware

System Disks

It is important to periodically check for new firmware updates that may be available for your system. Storage systems that are used by a replication set should run the same or compatible firmware versions. You can update the firmware in each controller module by loading a firmware file obtained from the enclosure vendor

Current Firmware Bundle	
Controller A (Currently using)	Controller B
ME5.1.0.1.0	ME5.1.0.1.0

Partner Firmware Update is Enabled ?

Upload the firmware bundle to install the bundle on controller A. You will have the opportunity to compare versions and activate the installed bundle using the table below.

# Always Use Supported Configurations

Always check ME5 Support Matrix for supported configurations, Operating systems and array rules. Do not risk your data, critical applications with unsupported configurations. Dell Technologies does not recommend or provide support for unsupported configurations. You can access latest support matrix from <https://www.dell.com/support> web site with your arrays array's service tag number or searching with your arrays model.

The screenshot displays the Dell support page for the PowerVault ME5024. At the top, there is a navigation bar with 'Support' and 'Product Support'. Below this, the product name 'PowerVault ME5024' is prominently displayed next to a small image of the device. A 'Change product' link is visible below the product name. The main content area features a search bar with the text 'Search PowerVault ME5024 Support Information' and a subtext 'Find articles, manuals and more to help support your product.' Below the search bar is a text input field with the placeholder 'What can we help you to find' and a search icon. To the left of the search bar, there are navigation tabs: 'Overview', 'Drivers & Downloads', 'Documentation' (which is selected), and 'Service Events'. Below the search bar, there is a 'Top Solutions' section. This section includes a 'See All' button and a list of three articles. The first article is titled 'PowerVault ME5: Access Key and PIN - We could not find the site you were searching, please verify if you have access to the site' and has a 'View Page' link. The second article is titled 'PowerVault ME5: Unable to Generate Access Key and PIN with Message "This page is only accessible to Dell business users"' and also has a 'View Page' link. The third article is titled 'ME5: Alert: Lost connectivity to the Dell SupportAssist backend and Event: Scheduler - SupportAssist is not registered with the Dell backend' and has a 'View Page' link. On the left side of the page, there is a sidebar with navigation links: 'TOP SOLUTIONS', 'MANUALS AND DOCUMENTS', 'REGULATORY INFORMATION', and 'VIDEOS'.

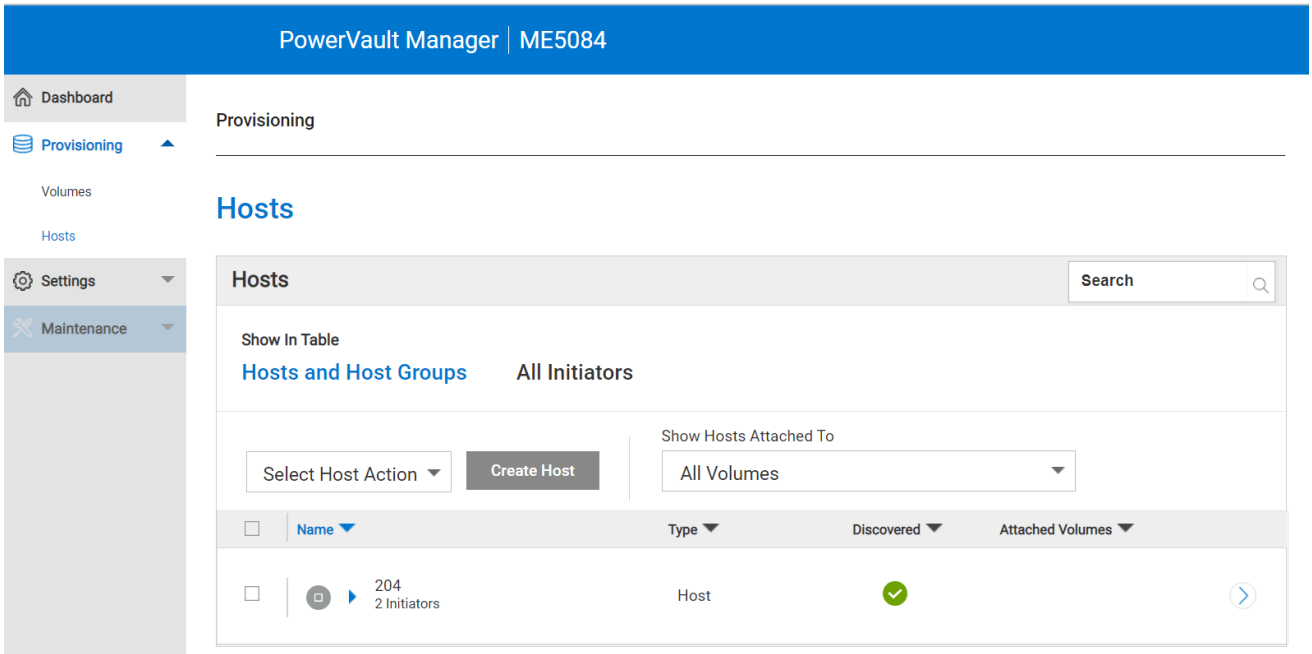
# Host Information

The Hosts block on PowerVault Manager GUI shows how many host groups, hosts, and initiators are defined in the system. An initiator identifies an external port to which the storage system is connected. The external port may be a port in an I/O adapter in a server, or a port in a network switch. A host is a user-defined set of initiators that represents a server. A host group is a user-defined set of hosts for ease of management. If the external port is a switch and there is no connection from the switch to an I/O adapter, then no host information will be shown.

## Identifying Your Hosts Easily

For easily identifying your hosts, it's highly recommended that using nicknames. recommended method for acquiring and renaming World-Wide Names (WWNs) is to connect one cable at a time and then rename the WWN to an identifiable name. You can change it via PowerVault Manager.

Select **Provisioning -> Hosts** on the left pane. Click arrow on the right-hand side and click pen figure on the next page, near the name.



## How to Monitor Array Health

PowerVault Manager dashboard can be used for monitor the system and see an overview of system status. Use the Alerts panel to monitor system health and performance issues and to track and acknowledge the resolution of these issues.

The screenshot shows the Dell EMC PowerVault Manager Alerts dashboard for ME5084. The dashboard includes a navigation menu on the left with options for Dashboard, Provisioning, Settings, and Maintenance. The main content area displays the Alerts panel, which is divided into four summary cards: System Health (OK), Health Alerts (0), Information Alerts (10), and Alerts to Acknowledge (12). Below these cards is a table titled 'Information Alerts (10)' with columns for Active Time, Severity, Component, and Acknowledged. The table lists three alerts: hostport\_A2 (16 days, Information), hostport\_B2 (16 days, Information), and hostport\_A3 (69 days, Information).

Active Time	Severity	Component	Acknowledged
16 days	Information	hostport_A2	
16 days	Information	hostport_B2	
69 days	Information	hostport_A3	

The ME5 storage systems can report their status through SNMP and/or Email. For more details please check [Administrator's Guide's](#) Dashboard Section

## Configuring E-mail and SNMP Notifications

The Notifications panel (Settings > Notifications) provides options to send system alert notifications to users through email, SNMP trap hosts, or a remote syslog server.

---

**Note:** You should enable at least one notification service to monitor the system.

---

## E-Mail Notifications

You can choose to be notified by email when system alerts occur. Alert notifications can be sent to a maximum of three email addresses. Weekly alerts concerning system health issues will also be sent until corrective action has been taken and the system health value has returned to OK.

Enter information in the text boxes to receive alert notifications. For details about panel options, see the on-screen tool tips. For information about SMTP notification parameters for events and managed logs, see the `set email-parameters` command in the [CLI Reference Guide](#).

---

**Note:** If the mail server is not on the local network, make sure that the gateway IP address was set in [Configuring controller network ports](#).

---

## SNMP Notifications

SNMP is a widely used network monitoring and control protocol. It is an application layer protocol that facilitates the exchange of management information between network devices. It is part of the Transmission Control Protocol/Internet Protocol (TCP/IP) protocol suite.

SNMP enables network administrators to manage network performance, find and solve network problems, and plan for network growth. Data is passed from SNMP agents reporting activity on each network device to the workstation console used to oversee the network. The agents return information contained in a Management Information Base (MIB), which is a data structure that defines what is obtainable from the device and what can be controlled (turned on and off, etc.).

The SNMP panel provides options to send alert notifications to SNMP trap hosts. You must enable SNMP for the system to send alert notifications to SNMP users. Enter information in the text boxes to receive alert notifications. For details about panel options, see the on-screen tool tips. See [Enabling or disabling system-management services](#) for more information on Admin Guide.



Dashboard

Provisioning

Settings

Network

Users

System

Notifications

Peer Connections

Maintenance

### Settings

## Notifications

Email

SNMP

Syslog

SMTP Server ?

Sender Email ?

Protocol ?

None

SMTP Port ?

Sender Password ?

Email Address 1 ?

Email Address 2

Email Address 3

# How To Provision Virtual Pools Features

This section outlines the best methods for optimizing virtual storage features such as Thin Provisioning, automated tiering for ME5 series.

## Thin Provisioning

Thin provisioning is a virtual storage feature that allows a system administrator to overcommit physical storage resources. This allows the host system to operate as though it has more storage available than is actually allocated to it. When physical resources fill up, the administrator can add physical storage by adding additional disk groups on demand.

Overcommit is enabled by default. The overcommit setting lets you oversubscribe the physical storage (that is, provision volumes in excess of the physical capacity). If you disable overcommit, you can provision virtual volumes only up to the available physical capacity. Overcommit is performed on a per-pool basis by using the Change Pool Settings option.

Maintenance

---

### Storage

A pool is an aggregation of one or more disk groups. When provisioning virtual storage, it is recommended that disks and provisioning are balanced between Pools A and B.

**Auto Storage Setup** This is optional. Apply a suggested configuration or expand a current configuration provided by the system. You can also add individual disk groups below.

Pool A	Size	Health	Available	Overcommit Size	
▼	166.2TB	✓	166.2TB	0	

Low Threshold

Middle Threshold

Pool Overcommit

▶ Disk Groups

To see or change overcommit settings go to : **PowerVault Manager -> Maintenance > Storage** click down arrow on the Pool A or Pool B and click Pen Symbol on right hand side, near Overcommit size value.

Each virtual pool has three thresholds for page allocation as a percentage of pool capacity. You can set the low and middle thresholds. The high threshold is automatically calculated based on the available capacity of the pool minus 200 GB of reserved space.

You can view and change settings that govern the operation of each virtual pool:

- **Low Threshold:** When this percentage of virtual pool capacity has been used, informational event 462 will be generated to notify the administrator. This value must be less than the Mid Threshold value. The default is 50 percent.
- **Mid Threshold:** When this percentage of virtual pool capacity has been used, event 462 will be generated to notify the administrator to add capacity to the pool. This value must be between the Low Threshold and High Threshold values. The default is 75 percent. If the pool is not overcommitted, the event will have Informational severity. If the pool is overcommitted, the event will have Warning severity.
- **High Threshold:** When this percentage of virtual pool capacity has been used, event 462 will be generated to alert the administrator to add capacity to the pool. This value is automatically calculated based on the available capacity of the pool minus 200 GB of reserved space. If the pool is not overcommitted, the event will have Informational severity. If the pool is overcommitted, the event will have Warning severity and the system will use write-through cache mode until virtual pool usage drops back below this threshold.
- **Pool Overcommit:** This check box controls whether overcommitting is enabled, and whether storage-pool capacity may exceed the physical capacity of disks in the system.

---

**Note:** If the pool size is 500 GB or smaller, or the middle threshold is relatively high or both, the high threshold may not guarantee 200 GB of reserved space in the pool. The controller will not automatically adjust the low and middle thresholds in such cases.

---

**Note:** If your system has a replication set, the pool might be unexpectedly overcommitted because of the size of the internal snapshots of the replication set.

---

**Note:** For more information about events, see the Event History panel (Maintenance > Support > Event History)

**Note:** If the pool is overcommitted and has exceeded its high threshold, its health will show as degraded in the Storage panel (**Maintenance > Storage**). If you try to disable overcommitment and the total space allocated to thin-provisioned volumes exceeds the physical capacity of their pool, an error will state that there is insufficient free disk space to complete the operation and overcommitment will remain enabled.

---

## Thin Provisioning Space Reclamation

The thin provisioning space reclamation primitive, also known as unmap, enables thin-provisioned datastores to be re-thinned to only consume the actual space they are consuming on the array. This frees up space on the array that has been deleted by ESXi, allowing thin-provisioned volumes to remain thin and reducing overall storage costs. Traditionally, the size of a thin-provisioned volume, as shown at the storage layer, reflects the maximum space consumption that occurred at some point since it was created. This is because ESXi did not inform the array that particular blocks of data had been deleted and no longer needed to be stored by the array. The T10 SCSI primitive unmap enables this information to be communicate to the array, through the SCSI storage stack. This unmap primitive is referred to as thin provisioning space reclamation by VMware.

With the release of vSphere 6.7, VMware updated the unmap API to run automatically in the background without user intervention as part of VMFS-6. This is dependent upon arrays utilizing 1 MB or smaller pages. PowerVault ME5 arrays utilize 4 MB pages, and therefore is incompatible with automatic unmap.

When a file is deleted on a Windows Server, the file pointer is deleted. However, the old data remains on the disk. Over time, the operating system overwrites the old data with new data.

For PowerVault volumes mapped to a Windows Server, the host passes a trim and unmap command to PowerVault when files are deleted. Within a few minutes, the PowerVault storage pool reflects the additional free capacity.

The ability to recover deleted disk space on PowerVault is a key benefit of thinly provisioned volumes. In cases where trim and unmap is not supported or disabled, reclaimed space appears as free in Windows, but not on the storage.

Windows Server and Hyper-V support trim and unmap natively with PowerVault given these conditions:

The Windows Server operating system must be version 2012 or newer (ME5 supports Server 2016 and newer).

Volumes must be basic disks that are formatted as NTFS volumes. Trim and unmap is not supported with other formats such as ReFS.

For more information related volume optimization on Windows Server operating System please visit [this link](#).

## Block Zeroing

Fault-tolerant virtual machines require VMDKs that are eager-zeroed thick. These VMDKs differ from standard thick or thin VMDKs in that the blocks are zeroed out when the VMDK is created. For large disks, this process can take a significant amount of time as each zero is written from the server to the array. Then, the array sends an acknowledgment of each write to the server, taking more time. With the block zeroing primitive, the ESXi host offloads to the ME5 array the task of zeroing out the blocks. The primitive also permits the host to continue creating the fault-tolerant virtual machine while the storage completes the zeroing task in the background. By offloading the block zeroing to the ME5 array, you can create fault-tolerant virtual machines much faster.

## Pool Balancing

In a storage system with two controller modules, try to balance the workload of the controllers. Each controller can own one virtual pool. Having the same number of disk groups and volumes in each pool will help balance the workload, increasing performance.

## Quick Rebuild

Quick rebuild is a method for reconstructing virtual disk groups that reduces the time that user data is less than fully fault-tolerant after a disk failure in a disk group. Taking advantage of virtual storage knowledge of where user data is written, quick rebuild only rebuilds the data stripes that contain user data.

Typically, storage is only partially allocated to volumes, so the quick-rebuild process completes significantly faster than a standard RAID rebuild. Data stripes that have not been allocated to user data are scrubbed in the background, using a lightweight process that allows future data allocations to be more efficient.

After a quick rebuild, a scrub starts on the disk group within a few minutes after the quick rebuild completes. Quick rebuild is only usable in ADAPT not in traditional RAID.

## Modifying Virtual Volumes

A virtual disk group requires the specification of a set of disks, RAID level, disk group type, pool target (A or B), and a name. If the virtual pool does not exist at the time of adding the disk group, the system will automatically create it. Multiple disk groups (up to 16) can be added to a single virtual pool.

You can expand a volume. If a virtual volume is not a secondary volume involved in replication, you can expand the size of the volume but not make it smaller. If a linear volume is neither the parent of a snapshot nor a primary or secondary volume, you can expand the size of the volume but not make it smaller. Because volume expansion does not require I/O to be stopped, the volume can continue to be used during expansion.

If there is not enough space in the pool's disk groups; the recommended method to expand the volume size is to add a new virtual disk group with the same RAID level, capacity disks, and physical number of disks as the existing virtual disk group in the same tier.

If storage array configured with ADAPT disk groups, instead of traditional RAID (RAID 5, RAID6 etc) you can even add 1 drive to ADAPT disk group for expanding volume. Please check [ADAPT section](#) of this paper for more information.

# Best Practices for High Availability

High availability is always advisable to protect assets in the event of a device failure. This section gives you some options and information to help you in the event of a failure.

## Volume Mapping

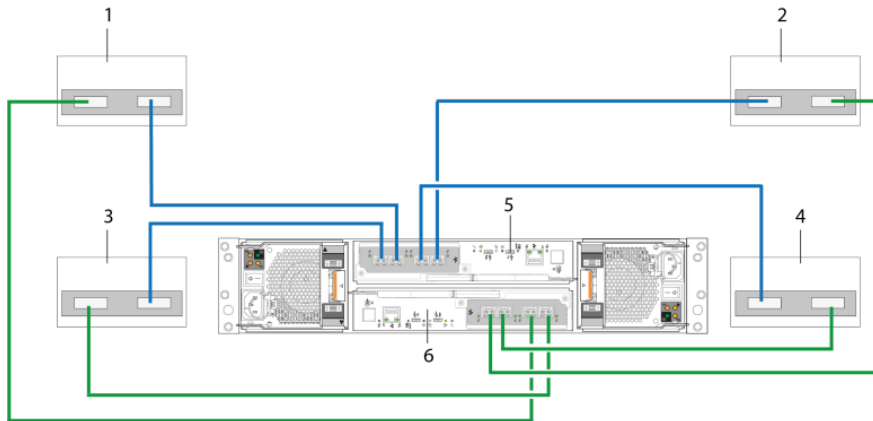
In both virtual and linear storage, if the owning controller fails, the partner controller assumes temporary ownership of the pool and resources owned by the failed controller. If a fault-tolerant cabling configuration, with appropriate mapping and MPIO, is used to connect the controllers to hosts, LUNs for both controllers are accessible through the partner controller so I/O to volumes can continue without interruption.

The best practice is to map volumes to two ports on each controller to take advantage of load balancing and redundancy.

To avoid multiple hosts mounting the volume and causing corruption, the hosts must be cooperatively managed, such as by using cluster software.

If multiple hosts mount a volume without being cooperatively managed, volume data is at risk for corruption. To control access by specific hosts, you can create an explicit mapping. An explicit mapping can use different access mode, LUN, and port settings to allow or prevent access by a host to a volume, overriding the default mapping. When an explicit mapping is deleted, the volume's default mapping takes effect.

## Direct Attach Cabling Example



ME5 Series 2U direct attach- four servers /one HBA per server / dual path

- 1- Server 1
- 2 – Server 2
- 3- Server 3
- 4 - Server 4
- 5 – Controller Module A
- 6 – Controller Module B

For more detailed information about host cabling please check the [System Deployment Guide](#)

## Multipathing Configuration and Multipath Software

ME5 systems comply with the SCSI-3 standard for Asymmetrical Logical Unit Access (ALUA). ALUA-compliant storage systems provide optimal and non-optimal path information to the host during device discovery. To implement ALUA, you must configure your servers to use multipath I/O (MPIO).

The Microsoft MPIO feature needs to be installed prior to connecting to ME5 storage arrays.

For Red Hat and Suse Linux Device Mapper multipath is required for multipath support.

Please check [PowerVault ME5 Support Matrix](#) for details.

Please check [Deployment Guide](#) for MPIO Setup

## Snapshots

The system can create snapshots of virtual volumes up to the maximum number supported by your system. Snapshots provide data protection by enabling you to create and save source volume data states at the point in time when the snapshot was created. Snapshots can be created manually, or you can schedule snapshot creation. After a snapshot has been created, the source volume can be expanded. To view the maximum number of snapshots for your system, see System configuration limits in [Administrator's Guide](#). When you reach the maximum number of snapshots for your system, before you can create a new snapshot, you must delete an existing snapshot.

You need to determine how you manage snapshots on virtual volumes in pools that have overcommit enabled. ME5 has two options for setting the frequency of snapshot management:

- 1) If you already maintain snapshots, then you probably do not need to change anything. The system automatically sets the limit at 10% of the pool and only notifies you if a threshold is crossed.
- 2) The set snapshot-space CLI command enables you to set the percent of the pool that can be used for snapshots (the snapshot space). Optionally, you can specify a limit policy to enact when the snapshot space reaches the percentage. You can set the policy to either notify you via the event log that the percentage has been reached (in which case the system continues to take snapshots, using the general pool space), or to notify you and trigger automatic deletion of snapshots. If automatic deletion is triggered, snapshots are deleted according to their configured retention priority. For more information, see the CLI documentation

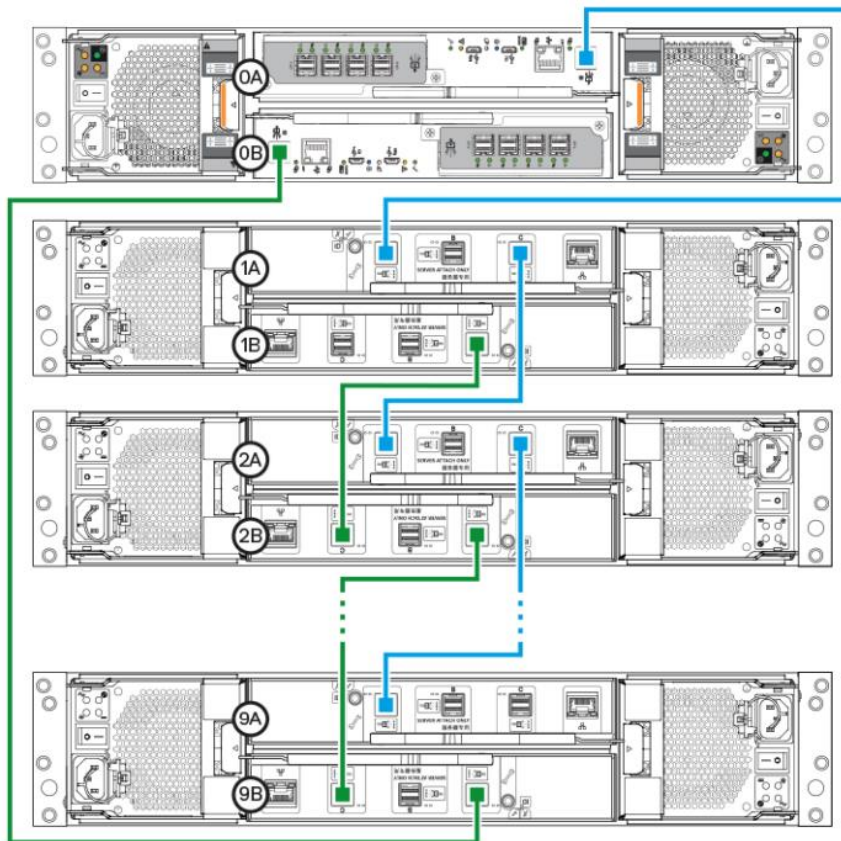
## Dual Power Supplies

All ME5 storage arrays ship with redundant power supplies. Each redundant power supply module requires power from an independent source or a rack power distribution unit with Uninterruptible Power Supply (UPS). 2U enclosures use standard AC power and the 5U84 enclosure requires high-line (high-voltage) AC power.

## Fault Tolerance (Reverse) Cabling

Reverse cabling allows any drive enclosure to fail—or be removed—while maintaining access to other enclosures. Fault tolerance and performance requirements determine whether to optimize the configuration for high availability or high performance when cabling.

The following figure shows the cabling configuration for a 2U controller enclosure with 2U expansion enclosures. The controller modules are identified as 0A and 0B, the IOMs in the first expansion enclosure are identified as 1A and 1B, and so on. Controller module 0A is connected to IOM 1A, with a chain of connections cascading down (blue). Controller module 0B is connected to the lower IOM (9B), of the last expansion enclosure, with connections moving in the opposite direction (green).



Please check [Deployment Guide](#) for other cabling examples



## SMART

Self-Monitoring Analysis and Reporting Technology (SMART) provides data that enables you to monitor disks and analyze why a disk failed. SMART is enabled by default and the system checks for SMART events one minute after a restart and every five minutes thereafter. SMART events are recorded in the event log.

Please check [Administrator's Guide](#) for how to modify SMART.

## Scrubbing

The system-level Disk Group Scrub option automatically checks all disk groups for disk defects. If this option is disabled, you can still perform a scrub on a selected disk group. Scrub analyzes the selected disk group to find and fix disk errors. It will fix parity mismatches for RAID 3, 5, 6, 50, and ADAPT; find but not fix mirror mismatches for RAID 1 and 10; and find media errors for all RAID levels.

Scrub can last over an hour, depending on the size of the disk group, the utility priority, and the amount of I/O activity. However, a manual scrub performed by Scrub Disk Group is typically faster than a background scrub performed by Disk Group Scrub. You can use a disk group while it is being scrubbed. When a scrub is complete, event 207 is logged and specifies whether errors were found and whether user action is required.

## Autonomic Distributed Allocation Protection (ADAPT)

ADAPT is a RAID-based data protection level that maximizes flexibility, provides built in spare capacity, and allows for very fast rebuilds, large storage pools, and simplified expansion. All disks in the ADAPT disk group must be the same type (enterprise SAS, for example), and in the same tier, but can have different capacities. ADAPT is shown as a RAID level in the management interfaces.

ADAPT disk groups use all available space to maintain fault tolerance, and data is spread evenly across all the disks. When new data is added, new disks are added, or the system recognizes that data is not distributed across disks in a balanced way, it moves the data to maintain balance across the disk group. Reserving spare capacity for ADAPT disk groups is automatic since disk space dedicated to sparing is spread across all disks in the system. In the case of a disk failure, data will be moved to many disks in the disk group, allowing for quick rebuilds and minimal disruption to I/O.

One of the key differences between ADAPT and traditional RAID groups is the width that arrays can be constructed. RAID 5 and 6 can be applied up to a width of 16 drives. However, whilst ADAPT widths are a minimum of 12, the maximum is 128 making the potential drive group width and therefore size much bigger than traditional R6. This has significant implications especially when potential topologies of a ME5084 are considered. With this ability one can consider how to layout the drive groups on a ME5084. For example, R6 would give 5 \* 16 disk groups and 4 spares.

If a disk fails in an ADAPT disk group, and the failed disk is replaced with a new disk in the same slot, the replacement disk will be added to the disk group automatically. All disks in the ADAPT disk group must be the same type (enterprise SAS, for example), but can have different capacities, provided the range of difference does not exceed a factor of two. For example, mixing a 600 GB disk and a 1.2 TB disk is acceptable; but mixing a 6 TB disk and a 16 TB disk could prove problematic. It is conceivable that a sizeable difference between mixed disk capacities (ratio greater than two) could prevent consuming space on disks due to insufficient distributed space required to support striping.

---

**Note:** Do not mix disks if the ratio of the largest disk to the smallest disk is greater than two.

---

Spare disks are not used by ADAPT disk groups since the RAID design provides built-in spare capacity that is spread across all disks in the disk group. In the case of a disk failure, data will be redistributed to many disks in the disk group, allowing for quick rebuilds and minimal disruption to I/O.

The system will automatically default to a target spare capacity that is the sum of the largest two disks in the ADAPT disk group, which is large enough to fully recover fault tolerance after loss of any two disks in the disk group. The actual spare capacity value can change depending on the current available spare capacity in the disk group. Spare capacity is determined by the system as disks are added to a disk group, or when disk groups are created, expanded, or rebalanced.

---

**Note:** If a disk fails in an ADAPT disk group and is replaced by a new disk in the same slot as the failed disk, the disk group automatically incorporates the replacement disk into the disk group.

---

For more information, please check the [PowerVault ME5 ADAPT whitepaper](#)

## Hot Spares

Spare disks are unused disks in your system that you designate to automatically replace a failed disk, restoring fault tolerance to disk groups in the system. Types of spares include:

**Dedicated spare:** Reserved for use by a specific linear disk group to replace a failed disk. Most secure way to provide spares for disk groups, but expensive to reserve a spare for each disk group.

**Global spare.** Reserved for use by any fault-tolerant disk group to replace a failed disk.

**Dynamic spare.** Available compatible disk that is automatically assigned to replace a failed disk in a fault-tolerant disk group.

When a disk fails, the system looks for a dedicated spare first. If it does not find a dedicated spare, it looks for a global spare. If it does not find a compatible global spare and the dynamic spares option is enabled, it takes any available compatible disk

---

**Note:** You cannot designate spares for ADAPT disk groups. For information on how ADAPT disk groups manage sparing Please check [Administrator's Guide](#) . A best practice is to designate spares for use if disks fail. Dedicating spares to disk groups is the most secure method, but it is also expensive to reserve spares for each disk group. Alternatively, you can enable dynamic spares or assign global spares.

---

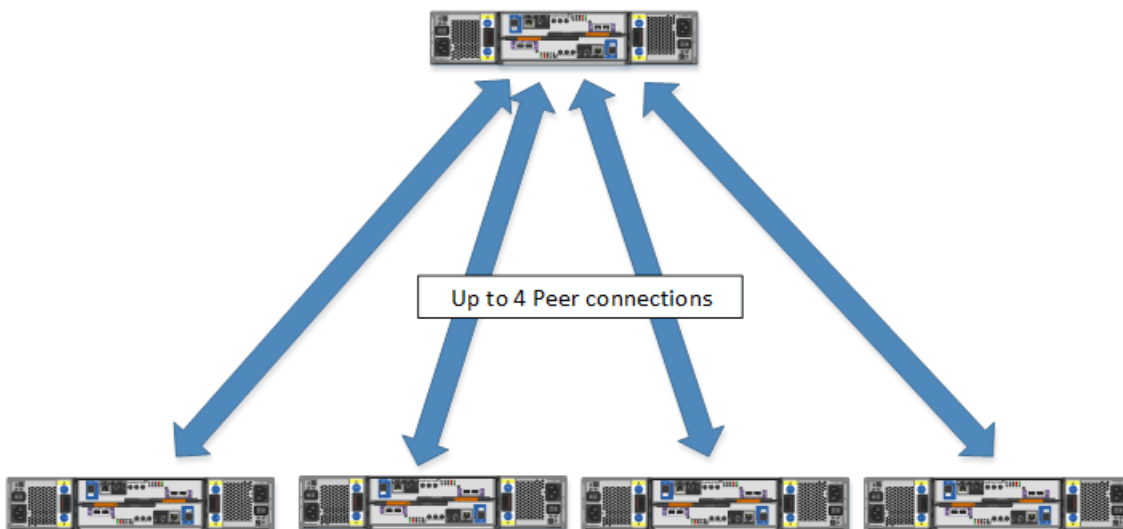
## Virtual Volume Replication

Replication is a feature for disaster recovery. This feature performs asynchronous replication of block-level data from a volume in a primary system to a volume in a secondary system by creating an internal snapshot of the primary volume and copying the changes to the data since the last replication to the secondary system via FC or iSCSI links.

The two associated standard volumes form a replication set, and only the primary volume (source of data) can be mapped for access by a server. Both must be connected through switches to the same fabric or network (no direct attach). The server accessing the replication set need only be connected to the primary system. If the primary system goes offline, a connected server can access the replicated data from the secondary system.

ME5 allows for up to four peer connections (replication partner), these are bi-directional connections, e.g.: a primary system can have up to four replication destinations (DR-sites), or the other way around, up to four prod systems replicate to the same DR site (many-to-one or one-to-many).

However, on a volume basis, a particular volume can only be part of one replica set, hence can only be replicated to one replication destination (one-to-one, not one-to-many).



For Cabling considerations of replication Please check [Deployment Guide](#). For other Replication prerequisites Please check [Administrator's Guide](#).

---

**Note:** Cross replication between ME4 and ME5 is supported. For more information, please refer to the Administrator's Guide.

---

# Best Practices for Performance

The secret to good performance is planning the configuration of the PowerVault system. The following section takes you through the things that should be considered from a performance perspective. What is typical in most system design exercises is that there are always compromises. This section should be used as guidance just from the performance perspective with the anticipation that in real life those compromises may mean it is not always possible to completely optimize for performance.

## Workload

The workload you plan to run on the system greatly influences the performance of the system. Random workloads deliver much slower bandwidth than sequential workloads. If possible, SSD's in random workload will mitigate these issues albeit at an increase in cost

## Use SSDs for Randomly Access Data

You can use SSDs for virtual disk groups. When combined with virtual disk groups that consist of other classes of disks, improved read and write performance is possible through automated tiered storage. Alternatively, you can use one or two SSDs in read-cache disk groups to increase read performance for pools without a Performance tier. The application workload of a system determines the percentage of SSDs of the total disk capacity that is needed for best performance.

Database indexes and Temporary Database files are good examples for SSD usage and you can get benefit of SSDs. Another great example for SSD usage is Virtual Desktop Infrastructures (VDI). For More details about VMware Horizon View VDI environments with ME5; please check [“Dell PowerVault ME5 Series: VMware Horizon VDI Best Practices”](#) document.

# PowerSizer

Utilize the [PowerSizer](#) to determine the achievable performance for a specific configuration. It is imperative to size the number and performance of drives beforehand using the [PowerSizer](#) as exceeding these parameters after deploying the system is not feasible.

### Workloads

Name: Workload 1 Pool Mapping: Single Flash Mode: Tiering

#### Tier Layout

Tier	Drive	Raid	Drive Count	Hot Spare Count
Extreme Performance	3.84TB SSD	RAID 5 (4+1)	5	1
Performance	1.2TB 10K SAS	RAID 6 (8+2)	20	1
Capacity	8TB NL-SAS	ADAPT (8+2)	12	0

Hot spares are not required and only available for drives used with traditional RAID.

RESET WORKLOAD

- 1** Workload - Workload 1 - To ensure the best performance with sequential workloads and RAID-5 or RAID-6 disk groups, use power-of-two data disks. For example, RAID 5 (4+1), RAID 5 (8+1), RAID 6 (2+2), RAID 6 (4+2).
- 1** Workload - Workload 2 - To ensure the best performance with sequential workloads and RAID-5 or RAID-6 disk groups, use power-of-two data disks. For example, RAID 5 (4+1), RAID 5 (8+1), RAID 6 (2+2), RAID 6 (4+2).
- ▲** For ME5084 please select high density and start with 28 drives and increments of 14 thereafter (28, 42, 56, 70, 84 etc). If adding a ME484 expansion that also needs to start with the initial 28 drives.

---

OVERVIEW
POWERPOINT

#### System Details

#### Workloads

##### WORKLOAD 1

Usable Capacity	Associated Pool(s)		Pool 1		Flash Mode	Tiering
89.63TiB	Rand R	Rand W	Seq R	Seq W		
Percent	80	20	0	0		
Size (KiB)	8	8	0	0		

##### WORKLOAD 2

Usable Capacity	Associated Pool(s)		Pool 2		Flash Mode	Tiering
89.63TiB	Rand R	Rand W	Seq R	Seq W		
Percent	80	20	0	0		
Size (KiB)	8	8	0	0		

#### System Details

#### Workloads

##### POOL 1

Usable Capacity	Raw Capacity	Performance Max	
89.63TiB	131.18TiB	13.054K IOPS / 101.98 MB/s	
Associated Workload	Flash Mode		
Workload 1 (Single)	Tiering		

Tier Layout	Tier	Drive	RAID	Hot Spare Count	Total Drive Count	Drive Groups	Usable Capacity	Performance
Extreme Performance		3.84TB SSD	RAID 5 (4+1)	1	6	1	13.96 TiB	10.937K IOPS
Performance		1.2TB 10K SAS	RAID 6 (8+2)	1	21	2	17.46 TiB	1.357K IOPS
Capacity		8TB NL-SAS	ADAPT (8+2)	0	12	1	58.2 TiB	0.76K IOPS

##### POOL 2

Usable Capacity	Raw Capacity	Performance Max	
89.63TiB	126.6TiB	13.054K IOPS / 101.98 MB/s	
Associated Workload	Flash Mode		
Workload 2 (Single)	Tiering		

Tier Layout	Tier	Drive	RAID	Hot Spare Count	Total Drive Count	Drive Groups	Usable Capacity	Performance
Extreme Performance		3.84TB SSD	RAID 5 (4+1)	0	5	1	13.96 TiB	10.937K IOPS
Performance		1.2TB 10K SAS	RAID 6 (8+2)	0	20	2	17.46 TiB	1.357K IOPS
Capacity		8TB NL-SAS	ADAPT (8+2)	0	12	1	58.2 TiB	0.76K IOPS

System Details

Workloads

Pools

## PowerVault ME5024

Usable Capacity 179.27TiB	Raw Capacity 257.78TiB	Performance Profile Recommended
Performance 26.108K IOPS / 203.96 MiB/s	Storage Mode Virtual	Total Drive Count 76
Rack Units 10U		

CAPACITY & DISKS

Internal Enclosure  
21 x 2.5" 1.2TB 10K SAS

Expansion Enclosure: (ME412)  
24 x 3.5" 8TB NL-SAS  
11 x 2.5" 3.84TB SSD

Expansion Enclosure: (ME424)  
20 x 2.5" 1.2TB 10K SAS

	SAS		SSD		NL-SAS
--	-----	--	-----	--	--------

Usable Capacity Percentage

19.49%

15.58%

64.93%

FRONTEND CONNECTIVITY

8 x 16Gb FC SFP

ADDITIONAL DETAILS

SizerID: pvault\_413583

Expansion: ME424 - 20 x 2.5" 1.2TB 10K SAS

Expansion: ME412 - 11 x 3.5" 8TB NL-SAS

Expansion: ME412 - 7 x 3.5" 8TB NL-SAS  
5 x 2.5" 3.84TB SSD

Expansion: ME412 - 6 x 2.5" 3.84TB SSD  
6 x 3.5" 8TB NL-SAS

Internal: 21 x 2.5" 1.2TB 10K SAS

## Dual pools

In order to make the best performance from the controllers, due to the ownership of drive groups by a given controller, best performance is achieved at a system level with at least 2 drive groups each being owned by the 2 controllers of the system.

## Linear vs Virtual

If the requirements are for consistent performance, linear is the best choice. Performance of Virtual is affected by the workload not just to the drive groups, but also the snaps and mirrors that may be in place in this system.

## Magical Number 2 Method

You can configure max “2” virtual pools in dual controller systems and each controller owns a pool. You need to balance disks between two pools. In other words, you need to use magical number “2” which means divide your disk quantity “2” for load balancing in each controller.

### **For example;**

If you have 13 x SSD drives and 25 x SAS drives and 49 NL-SAS drives in ME5 storage array, then::

Pool 1: 6 x SSD; 12 x SAS ; 24 x NL-SAS

Pool 2: 6 x SSD; 12 x SAS ; 24 x NL-SAS

And assign at least 1 x global spares for each type of drives.

## Disk Groups in a Pool

For better efficiency and performance, use similar disk groups per tier in a pool.

- Disk count balance: For example, with 20 disks, it is better to have two 8+2 RAID-6 disk groups than one 10+2 RAID-6 disk group and one 6+2 RAID-6 disk group.
- RAID balance: It is better to have two RAID-5 disk groups than one RAID-5 disk group and one RAID-6 disk group.
- In terms of the write rate, due to wide striping, tiers and pools are as slow as their slowest disk groups.
- All disks in a tier should be the same type. For example, use all 10K disks or all 15K disks in the Standard tier.
- Create more small disk groups instead of fewer large disk groups.
- Each disk group has a write queue depth limit of 100. This means that in write-intensive applications this architecture will sustain bigger queue depths within latency requirements.
- Using smaller disk groups will cost more raw capacity. For less performance-sensitive applications, such as archiving, bigger disk groups are desirable.

## Which RAID Level Should You Use?

There is not only one answer to this question. It depends on your environment. The following table describes the characteristics and use cases of each RAID level.

RAID level	Protection	Performance	Capacity	Application use cases	Suggested disk speed
RAID 1	Protects against up to one disk failure per mirror set	Great random I/O performance	Poor: 50% fault tolerance capacity loss	Databases, OLTP, Exchange Server	10K, 15K, 7K
RAID 5	Protects against up to one disk failure per RAID set	Good sequential I/O performance, moderate random I/O performance	Great: One-disk fault tolerance capacity loss	Big data, media and entertainment (ingest, broadcast, and past production)	10K, 15K, lower capacity 7K
RAID 6	Protects against up to two disk failures per RAID set	Moderate sequential I/O performance, poor random I/O performance	Moderate: Two disk fault tolerance capacity loss	Archive, parallel distributed file system	High capacity 7K

## ADAPT performance

The performance of an ADAPT array is similar to the performance of the RAID6 array, but there are some additional points to be considered with the deployment of an ADAP array type. Firstly, the number of drives in the drive group is different. R6 is 4-16 drives. However, ADAPT is 12-128 drives per drive group. What this means is that from a performance perspective a single ADAPT drive group can have significantly more spindles and therefore both sequential and random performance. As can be seen when using sizer, performance of the drive group will increase with additional members, up to the point of maximum performance of a given RAID controller.

Additionally ADAPT comes with 2 different settings for the array width 8+2 or 16+2. 16+2 can only be selected with a minimum of 20 drives. If you compare a 20-drive array with 16+2 or 8+2 you will find that both perform similarly on random performance as this is defined more by the number of spindles. However, for sequential performance, 16+2 will be faster than 8+2 and should be used for any sequential workload types.

ADAPT performance compared to R6 performance are very similar , however ADAPT is less consistent in sequential performance than R6. For most applications this is not a factor as this effect is small. However for latency sensitive application R6 may offer more consistent performance.



## Disk Count Per RAID Level

The following table shows recommended disk counts for RAID-6 and RAID-5 disk groups. Each entry specifies the total number of disks and the equivalent numbers of data and parity disks in the disk group. Note that parity is distributed among all the disks.

RAID Level	Total Disks	Data disks (equivalent)	Parity disks (equivalent)
RAID 6	4	2	2
	6	4	2
	10	8	2
RAID 5	3	2	1
	5	4	1
	6	8	1

To ensure best performance with sequential workloads and RAID-5 and RAID-6 disk groups, use a power-of-two data disks.

The controller breaks virtual volumes into 4-MB pages, which are referenced paged tables in memory. The 4-MB page is a fixed unit of allocation. Therefore, 4-MB units of data are pushed to a disk group. A write performance penalty is introduced in RAID-5 or RAID-6 disk groups when the stripe size of the disk group isn't a multiple of the 4-MB page.

- Example 1: Consider a RAID-5 disk group with five disks. The equivalent of four disks provide usable capacity, and the equivalent of one disk is used for parity. Parity is distributed among disks. The four disks providing usable capacity are the data disks and the one disk providing parity is the parity disk. The parity is distributed among all the disks, but conceiving of it in this way helps with the example.

---

**Note:** The number of data disks is a power of two (2, 4, and 8). The controller will use a 512-KB stripe unit size when the data disks are a power of two. This results in a 4-MB page being evenly distributed across two stripes. This is ideal for performance.

---

- Example 2: Consider a RAID-5 disk group with six disks. The equivalent of five disks now provides usable capacity. Assume the controller again uses a stripe unit of 512-KB. When a 4-MB page is pushed to the disk group, one stripe will contain a full page, but the controller must read old data and old parity from two of the disks in combination with the new data to calculate new parity. This is known as a read-modify-write, and it's a performance killer with sequential workloads. In essence, every page push to a disk group would result in a read-modify-write. To mitigate this issue, the controllers use a stripe unit of 64-KB when a RAID-5 or RAID-6 disk group isn't created with a power-of-two data disks. This results in many more full-stripe writes, but at the cost of many more I/O transactions per disk to push the same 4-MB page.

# Optimization of an existing System

The PowerVault ME5 storage system is designed to optimize itself where possible and therefore there are only a few options available to change performance on a configured system.

## Volume Cache Options

You can set options that optimize reads and writes performed for each volume. It is recommended that you use the default settings.

You can enable and disable the write-back cache for each volume. By default, volume write-back cache is enabled. Because controller cache is backed by supercapacitor technology, if the system loses power, data is not lost. For most applications, this is the preferred setting.

**CAUTION:** Only disable write-back caching if you fully understand how the host operating system, application, and adapter move data. Used incorrectly, write-back caching can hinder system performance.

## Using write-back or write through caching

When modifying a volume, you can change its write-back cache setting. Write-back is a cache-writing strategy in which the controller receives the data to be written to disks, stores it in the memory buffer, and immediately sends the host operating system a signal that the write operation is complete, without waiting until the data is written to the disk. Write-back cache mirrors all the data from one controller module cache to the other. Write-back cache improves the performance of write operations and the throughput of the controller.

When write-back cache is disabled, write-through becomes the cache-writing strategy. Using write-through cache, the controller writes the data to the disks before signaling the host operating system that the process is complete. Write-through cache has lower write throughput performance than write-back, but it is the safer strategy, with minimum risk of data loss on power failure. However, write-through cache does not mirror the write data because the data is written to the disk before posting command completion and mirroring is not required. You can set conditions that cause the controller to change from write-back caching to write-through caching. For more information, see [Changing system cache settings in Administrator's Guide](#).

In both caching strategies, active-active failover of the controllers is enabled.

---

**Note:** The best practice for a fault-tolerant configuration is to use write-back caching.

---

---

**Note:** Single controller system only use write-through

---

## Cache Optimization Mode

You can also change the optimization mode.

**Standard:** This controller cache mode of operation is optimized for sequential and random I/O and is the optimization of choice for most workloads. In this mode, the cache is kept coherent with the partner controller. This mode gives you high performance and high redundancy. This is the default.

**No-mirror:** In this mode of operation, the controller cache performs the same as the standard mode with the exception that the cache metadata is not mirrored to the partner. While this improves the response time of write I/O, it comes at the cost of redundancy. If this option is used, the user can expect higher write performance but is exposed to data loss if a controller fails.

**CAUTION:** Changing the cache optimization setting while I/O is active can cause data corruption or loss. Before changing this setting, quiesce I/O from all initiators.

## Optimizing read-ahead caching

You can optimize a volume for sequential reads or streaming data by changing its read-ahead cache settings.

You can change the amount of data read in advance. Increasing the read-ahead cache size can greatly improve performance for multiple sequential read streams.

- The **Adaptive** option works well for most applications: it enables adaptive read-ahead, which allows the controller to dynamically calculate the optimum read-ahead size for the current workload.
- The **Stripe** option sets the read-ahead size to one stripe. The controllers treat NRAID and RAID-1 disk groups internally as if they have a stripe size of 512 KB, even though they are not striped.
- Specific size options let you select an amount of data for all accesses.
- The **Disabled** option turns off read-ahead cache. This is useful if the host is triggering read ahead for what are random accesses. This can happen if the host breaks up the random I/O into two smaller reads, triggering read ahead.

---

**Note:** Only change read-ahead cache settings if you fully understand how the host operating system, application, and adapter move data so that you can adjust the settings accordingly.

---

# Automated Tiered Storage

Automated Tiered Storage is a virtual storage feature that automatically moves data residing in one class of disks to a more appropriate class of disks based on data access patterns, with no manual configuration necessary. Automated tiered storage operates as follows:

- Frequently accessed "hot" data can move to disks with higher performance
- Infrequently accessed "cool" data can move to disks with lower performance and lower costs.

Each virtual disk group, depending on the type of disks it uses, is automatically assigned to one of the following tiers:

**Performance**—This highest tier uses SSDs, which provide the best performance but also the highest cost.

**Standard**—This middle tier uses enterprise-class spinning SAS disks, which provide good performance with mid-level cost and capacity.

**Archive**—This lowest tier uses midline spinning SAS disks, which provide the lowest performance with the lowest cost and highest capacity.

When the status of a disk group in the Performance Tier becomes critical (CRIT), the system will automatically drain data from that disk group to disk groups using spinning disks in other tiers providing that they can contain the data on the degraded disk group. This occurs because similar wear across the SSDs is likely, so more failures may be imminent.

If a system only has one class of disk, no tiering occurs. However, automated tiered storage rebalancing happens when adding or removing a disk group in a different tier.

## How Tiering Works?

Auto-tiering uses a concept called "paging". User volumes are logically broken down into small, 4MB chunks called pages. Pages are ranked based upon an algorithm. The page rank is used to very efficiently select good pages to move between tiers. The result is that pages can be migrated between tiers automatically such that I/Os are optimized in real-time.

- Tiering algorithm runs every 5 seconds.
- Only 80 MB of data is migrated every five seconds to avoid degrading system throughput.
- Frequently accessed data moved up to higher performance disks
- Infrequently access data moved down to lower performance disks
- Pages are only migrated down if room needed for highly ranked page
- Single copy of specific blocks of data resides in either spinning drives or SSDs

## Volume Tier Affinity Feature

The volume tier affinity feature enables tuning the tier-migration algorithm for a virtual volume when creating or modifying the volume so that the volume data automatically moves to a specific tier, if possible. If space is not available in the preferred tier, another tier will be used. The three volume tier affinity settings are:

The screenshot displays the 'Overview' tab of a storage system interface. At the top, there are navigation tabs: 'Overview' (selected), 'Snapshots', 'Attached Hosts (0)', and 'Replications'. Below these are two blue buttons: 'Expand Volume' and 'Copy Volume'. A large blue heading 'Capacity' is followed by a progress bar. The progress bar shows 'USED' space as 0.0 KB / 0.0% and 'UNUSED' space as 199.9 GB / 100.0%. Below the capacity section is a 'Details' section with the following information:

Volume Name:	TEST
Serial Number:	00c0ff646498000012fde86101000000
WWN:	600C0FF00064649812FDE86101000000
Owner:	no affinity
Cache Write Policy:	archive
Optimization:	performance
Read-Ahead Size:	
Tier Affinity:	no affinity

**No Affinity** - This setting uses the highest available performing tiers first and only uses the Archive tier when space is exhausted in the other tiers. Volume data swaps into higher performing tiers based on the frequency of access and tier space availability.

**Performance** - This setting prioritizes volume data to the higher performing tiers. If no space is available, lower performing tier space is used. Performance affinity volume data swaps into higher tiers based upon frequency of access or when space is made available.

**Archive** - This setting prioritizes the volume data to the lowest tier of service. Volume data can move to higher performing tiers based on the frequency of access and available space in the tiers.

## Optimization in Tier Setup

In general, it is best to have two tiers instead of three tiers. The highest tier will nearly fill before using the lowest tier. The highest tier must be 95% full before the controller will evict cold pages to a lower tier to make room for incoming writes.

Typically, you should use tiers with SSDs and 10K/15K disks, or tiers with SSDs and 7K disks. An exception may be if you need to use both SSDs and faster spinning disks to hit a combination of price for performance, but you cannot hit your capacity needs without the 7K disks; this should be rare.

Recommended setting for Volume Tier Affinity is “**No Affinity**” for most configurations. This setting attempts to balance the frequency of data access, disk cost, and disk availability by moving the volume’s data to the appropriate tier.

If the virtual volume uses mostly random or burst low-latency workloads such as online transaction processing (OLTP), virtual desktop infrastructure (VDI), or virtualization environments, recommended setting is “**Performance**”. This setting keeps as much of the volume’s data in the performance tier for as long as possible.

If the virtual volume contains infrequently accessed workloads such as backup data or email archiving, recommended setting is “**Archive**”. This option keeps as much of the volume’s data in the archive tier for as long as possible.

Some workloads do not respond well to the tiering feature. This is when the workload offers no capability for the caching effects of the tier to improve performance. Examples of this are streaming environments where the host streams are never re accessed and therefore there is no improvement in performance. The workload therefore needs to be understood to see if it allows optimization with the tiering feature. Good performance is seen where the “working set” i.e. the data that is typically being used by the host application lives within the tier. The “working set” size and the amount of SSD storage are closely related and therefore increasing the amount of SSD storage in these cases will improve performance at the system level if this is the case.

## Gauging the percentage of life remaining for SSDs

An SSD can be written and erased a limited number of times. Through the SSD Life Left disk property, you can gauge the percentage of disk life remaining. This value is polled every 5 minutes. When the value decreases to 20%, an event is logged with Informational severity. This event is logged again with Warning severity when the value decreases to 5%, 2% or 1%, and 0%. If a disk crosses more than one percentage threshold during a polling period, only the lowest percentage will be reported. When the value decreases to 0%, the integrity of the data is not guaranteed. To prevent data integrity issues, replace the SSD when the value decreases to 5% of life remaining.

## All-Flash Array

The all-flash array feature, enabled by default, allows systems to run exclusively with disk groups that consist of SSDs, providing the ability to have a homogeneous SSD-only configuration. Systems using an all-flash array have one tier that consists solely of SSDs. If a system includes disk groups with spinning disks, the disk groups must be removed before the all-flash array feature can be used.

## SSD Read Cache

Unlike tiering, where a single copy of specific blocks of data resides in either spinning disks or SSDs, the Read Flash Cache (RFC) feature uses one SSD read-cache disk group per virtual pool as a read cache for frequently accessed data only. Each read-cache disk group consists of one or two SSDs with a maximum usable capacity of 4TB. A separate copy of the data is also kept in spinning disks. Read-cache content is lost when a controller restart or failover occurs. Taken together, these attributes have several advantages:

- The performance cost of moving data to read-cache is lower than a full migration of data from a lower tier to a higher tier.
- Read-cache does not need to be fault tolerant, potentially lowering system cost.
- Controller read cache is effectively extended by two orders of magnitude, or more.
- 

When a read-cache group consists of one SSD, it automatically uses NRAID. When a read-cache group consists of two SSDs, it automatically uses RAID 0.

## Full Disk Encryption (FDE)

A system and the FDE-capable disks in the system are initially unsecured but can be secured at any point as long as all disks in the system chain are FDE. Until the system is secured, FDE-capable disks function exactly like disks that do not support FDE.

Enabling FDE protection involves setting a passphrase and securing the system. Data that was present on the system before it was secured is accessible in the same way it was when it was unsecured. However, if a disk is transferred to an unsecured system or a system with a different passphrase, the data is not accessible.

FDE operates on a per-system basis, not a per-disk group basis. To use FDE, all disks in the system must be FDE-capable.

**CAUTION:** Do not change FDE configuration settings while running I/O. Temporary data unavailability may result. Also, the intended configuration change might not take effect.

---

**Note:** Be sure to record the passphrase as it cannot be recovered if lost.

---

# Best Practices for Firmware Updates

Controller modules, expansion modules, and disk drives contain firmware that operate them. As newer firmware versions become available, they may be installed at the factory or at a customer maintenance depot or they may be installed by storage-system administrators at customer sites.

- In the Alerts panel on the dashboard, verify that the system health is OK. If the system health is not OK, expand the view to see the active health alerts and resolve all problems before you update firmware. For information about Active Alerts, see [How to Monitor Array Health](#) section
- Run the `check firmware-upgrade-health` CLI command before upgrading firmware. This command performs a series of health checks to determine whether any conditions exist that must be resolved before upgrading firmware. Any conditions that are detected are listed with their potential risks. For information about this command, see the [CLI Reference Guide](#)
- If any unwritten cache data is present, firmware update will not proceed. Before you can update firmware, unwritten data must be removed from cache. See information about the clear cache command in the [CLI Reference Guide](#).

**CAUTION:** Removing unwritten data may result in data loss. Contact technical support for assistance.

- If a disk group is quarantined, resolve the problem that is causing it to be quarantined before updating firmware.
- To ensure success of an online controller firmware update, select a period of low I/O activity. This helps the update to complete as quickly as possible and avoids disruption to host and applications due to timeouts. Attempting to update a storage system that is processing a large, I/O-intensive batch job may cause hosts to lose connectivity with the storage system.
- Confirm PFU is enabled by clicking **Settings > System > Properties > Firmware Properties**.
- Do not perform a power cycle or controller restart during a firmware update. If the update is interrupted or there is a power failure, the module might become inoperative. If this occurs, contact technical support. The module might need to be returned to the factory for reprogramming.

When planning a controller firmware upgrade:

- Online firmware upgrades are performed while host I/O being processed, and this time frame performance of array can impact. Select appropriate time frame for upgrade operation especially when host I/O activity is low.
- Spare at least 30 minutes for firmware upgrades.
- Please ensure that both controller's management ports' ethernet connection available before start upgrade process.



---

**Note:** Please always check README-First.pdf document which can be found in the firmware zip file before doing any update.

---

---

**Note:** Expansion firmware is updated automatically with controller updates.

---

## Updating Disk-Drive Firmware

You can update disk-drive firmware by loading a firmware file obtained from Dell Support Web Site

A dual-ported disk drive can be updated from either controller. Disk Drive Firmware update is an OFFLINE process. Stop I/O to the storage system. During the update all volumes will be temporarily inaccessible to hosts. If I/O is not stopped, mapped hosts will report I/O errors. Volume access is restored after the update completes.

Please check [Administrator's Guide](#) for more details.