## Overview

This module provides information on known best practices for configuring disk groups, sizing disk groups, and improving their availability and performance. The module begins by describing the factors that go into disk group configuration and sizing, including availability, performance, and cost. After each factor is described, the best practices for optimizing that factor are discussed. Next, a formula for determining the number of disks in a disk group is included. Finally, information is presented about current utilities for documenting the current configuration and for determining storage capacity requirements.

# Objectives

After completing this module, you should be able to:

- Describe the key considerations that go into configuring disk groups.

- List the best practices for configuring disk groups with respect to availability, performance, and cost.

- Describe the factors that go into sizing disk groups.

- Use the formula for calculating disk counts needed for the given requirements of a disk group.

- Identify the utilities available for documenting the storage system configuration, finding troubleshooting information, and for calculating storage capacity requirements.

# Available documentation and tools

The following white papers on the TechBB, Cybrary, or Tiger Team e-Room give details about best practices for configuring and sizing disk groups:

■ *EVA Availability Best Practices — Robust Availability Configuration*. The focus is on optimizing availability of the array. This paper is geared toward the service engineer and is not available to the customer.

■ *EVA Best Practices — Cost, Performance, and Availability*. The focus is on all three areas and their tradeoffs. This paper is geared to and available to the customer.

■ *Configuring Disk Groups and Virtual Disks in a StorageWorks Enterprise Virtual Array*

■ *Sizing Disk Groups in a StorageWorks Enterprise Virtual Array*

The following utilities are also available:

■ EVA Documentation Tool (EVA-DT) — HP Storage Tools Portal

■ EVA Troubleshooting Tool — NSS Online Products (Boise) and HP Storage Tools Portal

■ Disk Group Sizing Utility — TechBB or Cybrary

■ Disk Groups and Virtual Disks Calculator — TechBB

These utilities are described later in this module.

# Disk group configuration overview

Given that the customer must decide on the number and organization of disk groups in the storage system, the configuration goal is to optimize the number of disk groups with respect to the following three factors:

- Availability requirements

- Performance requirements

- Cost of ownership of the storage

It is not always possible to simultaneously optimize a configuration for all three factors. This may yield contradictory advice, since one of the three choices must yield to demands made by another. For example, VRAID0 is the best choice from a cost standpoint, since all of the storage is available for user data. But from an availability standpoint, VRAID1 is a much better choice, although storage utilization is only 50%.

Other tradeoffs are minor by comparison, but must sometimes still be made. There is no best choice in these situations, since it will depend on the needs of a particular environment.

---

**Note**

The material in this discussion comes largely from the white papers, *EVA Availability Best Practices — Robust Availability Configuration*, and *EVA Best Practices — Cost, Performance, and Availability.*

---

# Best practices for optimizing availability

The EVA is designed to be fully redundant with no single point of failure. There are, however, supported configurations that result in multiple drives on a single shelf being part of the same disk group. These configurations may result in data unavailability due to multidisk failures or human error. This section describes how to minimize these failures and to prevent unexpected downtime by using what is called a Robust Availability Configuration.

Factors that go into optimizing for availability include VRAID type, protection level, the number of disk groups, and drive replacement procedures. Best practices for each of these areas are described in the subtopics.

## VRAID type

Each VRAID type provides a different level of availability:

■   VRAID0 — Provides no availability protection. If one disk fails, there is no parity protection and the data will be lost.

■   VRAID1 — Provides the highest availability protection. If one disk fails, the data is available an exact copy is available on another disk. This provides the best protection for configuration below 8 disk drive enclosures. It is only 50% efficient because two copies of the data are stored.

■   VRAID5 — Provides the most efficient protection from the standpoint of disk utilization because it is 80% efficient (parity occupies 20% of disk space). However, it is less efficient in read performance if a disk failure occurs because lost and data must be computed from the parity and other disks. Some failures, such as an enclosure failures, will not be fully protected in configurations with less than 8 drive enclosures.

**Availability best practice:**

**Do not use VRAID0 because it does not provide protection from a disk failure. For best performance and availability, use VRAID1. For best storage efficiency and cost, but at lower performance and availability protection, use VRAID5.**

## Protection level

The protection level is reserved (assigned) capacity used to rebuild the data on a failed disk. Conceptually, space is reserved during disk group creation to handle zero, one, or two disk failures. The space reserved is specific to a particular group, and cannot span group boundaries.

The algorithm for reserving protected space is to find the largest disk in the disk group, multiply the result by 0, 2, or 4 (depending on a protection level of None, Single, or Double), and then remove that capacity from free space and distribute it across all disks in the disk group. Unlike traditional arrays, it does not reserve physical disks; all disks remain in use. The reason for finding the largest disk is that even though there may only be a few large disks in a group, they must be protected as well as the smaller ones, so space must be preserved. That size must be doubled, since the recovery algorithm for VRAID1 failures involve transferring the data to two new disks.

Having one or two large disks in a group means that extra space will be reserved. For example, if a disk group of 168 disks consisted of all 36GB disks with the exception of two 72GB disks, a protection level of Double would reserve 288GB of space (72GB x 4), even though the vast majority of the disks require only 144 of protection (36GB x 4).

VRAID1 mirror partners can only use the capacity of the smallest disk in a mirror pair. If, for example, a 36GB drive was mirrored with a 72GB drive, only 36GB of data could be used for mirroring on the 72GB drive, potentially wasting large amounts of space.

**Availability best practice: Do not mix disk sizes in a single disk group.**

The sparing algorithm attempts to use unassigned capacity in the disk group first and will only use the reserved spare space if there is not adequate unassigned capacity. This maximizes the availability of the array. By using unassigned capacity for recovery first, the array maintains the reserved spare space if another drive fails. This provides the highest level of data availability.

**Availability best practice: Always configure spare space for every array group. Do not rely on unassigned capacity to spare a failed drive.**

# Number of disk groups

Although the EVA offers numerous levels of data protection and redundancy, a catastrophic failure can result in loss of a disk group. This is extremely unlikely, and requires multiple simultaneous disk failures of disks in the same RSS. In spite of this very low probability, installations that demand the ultimate in data availability might consider creating two separate disk groups.

Although two groups will result in a slightly higher cost of ownership and potentially lower performance, the increase in availability may be the right decision for a very high availability application.

In order for two disk groups to prevent data loss, each disk group must contain sufficient independent information to reconstruct the entire data set. A practical example of this is a database that contains both data and log files. In this instance, placing the data files in one group and duplexing the log files (a typical feature of the database) to both the data file group and another group ensures that loss of an entire disk group will not prevent recovering the data.

> **Availability best practice: For critical database applications, consider placing data and log files in separate disk groups.**

# Drive failure and replacement

Although following the above rules will protect against loss of data access in the event of a shelf failure, there are specific steps that must be taken to help continue this protection after a disk fails.

When a disk fails, the EVA will rebuild the failed disk data through a process known as "sparing". This sparing action will rebuild the original level of redundancy, but may place two members of the redundancy set on the same shelf. To restore the disk group to the original configuration, specific steps must be followed:

> **Availability best practice:**
>
> a.  **Wait for the sparing to be completed. This will be signaled by an entry in the event log (VCS versions 2.002 and above).**
>
> b.  **Remove the failed disk from the shelf and replace with a new one.**
>
> c.  **Add the new disk into the original disk group.**

Following this action, the EVA will initiate a leveling operation to evenly distribute data across all disks in the disk group, which implicitly results in restoring the original configuration.

When a disk is inserted into a shelf, there will be some transient activity on the back-end fibre bus. In order to keep this from causing a false indication of excessive errors, insertion of multiple disks should be done carefully and slowly, with a pause between inserting disks.

> **Availability best practice: After inserting a disk drive into a shelf, wait 60 seconds before inserting another disk.**

## Robust Availability Configuration

Beginning with VCS V2.0 firmware, physically configuring the array to contain eight or more disk shelves and then creating disk groups that contain a multiple of eight disks will help maximize data availability following a hardware failure on the array, including a shelf failure. The firmware will distribute data for virtual disks to maximize data availability if shelf failure occurs. This is the VCS Robust Availability Configuration algorithm (formerly known as the *Non-Stop VRAID Configuration*).

The robust configuration requires the disks in the array to be organized vertically so that the disk mechanisms reside in the same slots in each shelf. It also requires that the disk groups created on the array contain a multiple of eight drives. When creating a disk group, the user specifies the size of the disk group but the firmware must select the mechanisms to be contained in the disk group. The firmware will allocate the mechanisms from the shelves vertically. This physical allocation of disks to the disk group sets up the array so that the RSSs in the disk group will not contain more than one disk per shelf.

### Robust Availability Configuration rules

The Robust Availability Configuration rules are the following:

- The array **must** contain eight or more disk shelves per controller pair.

- Disk mechanisms should be placed in the shelves so the same slots are used in all shelves.

- The array should contain a multiple of eight disks, for example, 8, 16, 24, 32.

- Disk groups **must** contain a multiple of eight disks.

- Allow the VCS firmware to choose the mechanisms to include in the disk group. It will pick disks vertically, ensuring no shelf contains two disks for the same RSS.

- The array must have functioning EMUs.

Given the fact that a Robust Availability Configurations requires a minimum of eight disk shelves per controller pair:
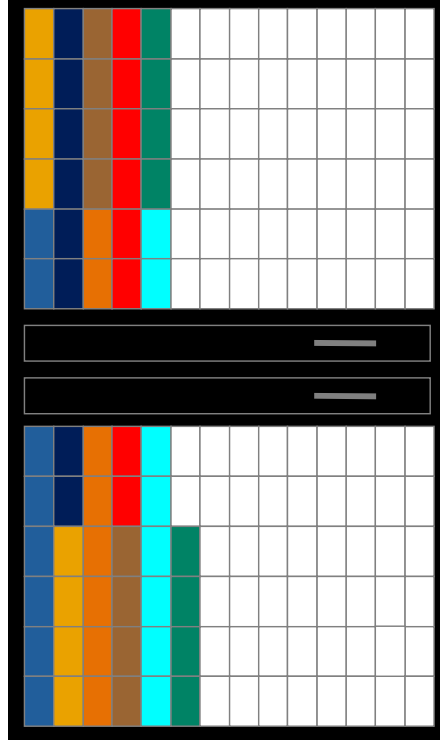
- The 2C2D, 2C4D, 2C6D, 8C8D, and eva3000 configurations cannot be put in a Robust Availability Configuration.

    > **Note**
    > An 8C8D is actually four 2C2D configurations in a single rack.

- Availability for these configurations can be maximized by creating only VRAID1 virtual disks.

- The VCS firmware maximizes data availability by mirroring data between different shelves even if the array configurations contains fewer than eight disk shelves.

The following shows a 2C12D with 64 drives in a Robust Availability Configuration. Each color represents a set of eight disks and RSSs.



**Note**

The firmware allocates drives from the bottom up, rotating from left to right by disk group.

You should note that a disk group with 64 disks cannot spare completely if a shelf fails. This is because five or six drives will need spares, and if Double protection is specified, only four drives maximum are allowed for sparing. There is no assurance that unassigned capacity can provide spare space. Ways to avoid this situation and provide maximum availability are provided later.

## Altering the configuration

When an array is configured in a Robust Availability Configuration, the firmware sets up the array so that there are not two disks from the same RSS residing in the same disk shelf. When a VRAID1 virtual disk is created, both copies of any given block of data contained in that VRAID1 virtual disk are stored on a "married" pair of disks. A married pair of mechanisms will always reside in the same RSS. For a VRAID5 virtual disk, all of the data in any given VRAID5 stripe (five blocks, four data and one parity) will always reside on different disks in the same RSS. Combining a Robust Availability Configuration and the method the VCS firmware uses to protect data within an RSS maximizes availability by preventing two critical blocks of data from residing in the same disk shelf. Even if a shelf failure occurs and multiple disk mechanisms become unavailable simultaneously, the worst that happens is that VRAID0 data becomes unavailable.

If a join or split operation has occurred for an RSS, the resulting RSS(s) will still follow the rules for a block of data in a VRAID1 or the data and parity in a VRAID5 stripe always residing in the same RSS. The firmware, however, cannot ensure that two disks from the same RSS will not reside in the same disk shelf. Therefore, you cannot ensure the array is still in a Robust Availability Configuration following an RSS combine or split. This occurs because following a disk mechanism failure or addition of new mechanisms to a disk group, disks may move from one RSS to another as the firmware attempts to reach the optimum RSS size of eight members.

⚠ **WARNING**

Any time even a single disk mechanism leaves or joins a disk group, the Robust Availability Configuration of the array may change and the array may no longer be in a Robust Availability Configuration.

There is a process that you can use to check the configuration in Appendix A of the white paper, *EVA Availability Best Practices — Robust Availability Configuration*.

## Maximizing availability

VRAID1 and VRAID5 virtual disks maximize data availability following a disk mechanism failure and a Robust Availability Configuration maximizes data availability if a disk shelf fails. The problem is how to maximize data availability following a shelf failure. By configuring the disk groups on the array in such a manner that there will always be enough spare space in each disk group to spare **all** the mechanisms in the failed shelf, you can maximize the array's ability to provide data availability even if a disk fails **after** a disk shelf failure.

To provide the highest availability of your disk groups after a shelf failure, you should adhere to specific rules called Availability Envelope rules. These rules maximize availability but can severely limit the number of disks in the disk group. But by adhering to these rules, full data protection will be maximized after sparing has completed following a disk shelf failure on the array.

If Single protection is specified (two drives for sparing):

■ Never put more than one disk per shelf in the disk group if VRAID1 virtual disks are going to be created in the disk group.

■ Never put more than two disks per shelf in the disk group if the disk group will only contain only VRAID5 virtual disks.

**Example**

With a 2C18D, you can create a disk group that contains up to 32 disks and still follow these rules.

If Double protection is specified (four drives for sparing):

- Never put more than two disks per shelf in the disk group if VRAID1 virtual disks are going to be created in the disk group.

- Never put more than four disks per shelf in the disk group if the disk group will only contain VRAID5 virtual disks.

  **Example**

  With a 2C18D, you can create a disk group that contains up to 72 disks and still follow these rules.

By adhering to these rules, data availability will be maximized by ensuring that data will still be available if a second disk failure follows any other single failure, including the failure of an entire disk shelf (sparing must complete before the second failure occurs).

## Adding to or removing drives from a disk group

Any time disks are added to or removed from a disk group, the Robust Availability Configuration becomes suspect and must be checked and repaired if necessary. There are two best practices to consider when adding or removing drives.

**Availability best practice: If additional disks are added to or removed from an existing disk group, the disk group must be inspected to be sure it is still in a Robust Availability Configuration.**

**Availability best practice: An alternative to adding disks to an existing disk group is to create a new disk group and set it up to follow the Robust Availability Configuration rules.**

# Best practices for optimizing performance

Performance and price are generally incompatible. The underlying virtualization technology goes a long way towards masking this, but there are still cases where higher performance comes at the cost of both price and availability. The following sections describe the factors that contribute to obtaining the best performance.

## Disk count

Since most random access applications are limited by the physical disk speeds, increasing the numbers of disks accessed by a LUN will translate directly to an increased performance potential. With the high transfer rates of modern disk drives, maximum sequential performance can be attained with only a few disks. Having additional disks does allow multiple sequential streams, or even intermixed random and sequential streams, to co-exist with minimal interaction, so for most applications, there will be a direct relationship between the number of disk drives and I/O performance.

**Performance best practice: Fill the array with as many disk drives as possible.**

## Number of disk groups

For typical workloads, an increased number of disk drives under a LUN imply increased performance potential, and since a LUN can only exist within a single disk group, having a single disk group maximizes the performance capability. Similarly, the larger numbers of disk drives associated with a single disk group imply less interaction among multiple I/O streams, as well as the ability to share disk resources as the I/O load shifts from one LUN to another.

**Performance best practice: Use a single disk group.**

## Disk RPM

For applications that perform large block sequential I/O, such as data warehousing and decision support, disk RPM has little or no effect on performance. As such, large capacity 10K RPM disks make the most sense.

For applications that issue small block random I/O, such as interactive databases, file and print servers, and mail servers, higher RPM disk drives offer a substantial performance advantage. Workloads such as these can see gains of 30% to 40% in the request rate when changing from 10K to 15K RPM disks. Although it seems contradictory to use 10K RPM disks for better performance in these circumstances, there are instances where it may make sense.

Although not guaranteed, it is a likely assumption that 15K RPM disks cost more than the equivalent capacity 10K drives. Since the gain from a 15K drive is in the range of 30% to 40%, if the 15K drives are more than 30-40% more expensive than the 10K drives, then it makes sense to purchase a larger number of 10K drives. Although a 15K RPM disk offers higher performance than the equivalent capacity 10K RPM disk, the lower cost of the 10K disk means that more can be purchased for the same total amount, and the increased number of disks translates directly into higher overall performance. Therefore, there are times when a "lower performance" disk such as a 10K RPM can, for the same dollar amount, yield not only higher performance, but also more capacity (with an attendant lower cost of ownership).

> **Performance best practice: Using 15K RPM disks is generally best, but carefully consider cost and quantity tradeoffs between 10K and 15K RPM disks.**

## Write cache mirroring

The purpose of cache mirroring is to prevent data loss in the unlikely event of a cache board failure. In operation, writes from the host are sent to the controller that has the LUN online. That controller will store the write data in its non-volatile cache, and then send that data across a 2Gb mirror port to the other controller. The second controller will store a copy of the data in its cache, and then return completion to the original controller. The original controller will then signal completion of the request to the host. Since there are two independent copies of the data maintained in separate caches, failure of a cache memory board does not result in loss of the data.

Because there is a data copy operation involved in this mirroring operation, there is an impact on the performance of writes. The amount of the impact depends on many factors, including the percentage of write operations, the size of the write requests, and the overall write request rate as presented by the various hosts.

From a performance perspective, it is possible to obtain significant gains in performance when a LUN is created with cache mirroring disabled. The clear disadvantage of disabling cache mirroring is, of course, the possibility of data loss if a cache board that contains data yet to be written to disk fails. There are applications that are not concerned with this possibility, such as databases that are reloaded every night, or other applications where loss of write data is a secondary concern in relation to performance, and these applications may be candidates for disabling write cache mirroring.

> **Performance best practice: Under certain, carefully considered circumstances, disabling write cache mirroring will result in significantly increased write performance.**

## Mixed disk speeds

Although it is esthetically preferable to have all disks in a group be of the same type, there are instances where mixed drive speeds may be unavoidable. From a performance standpoint, 15K RPM disks will perform 30% to 40% faster than 10K RPM disks with a small block, random access workload. Mixing drive speeds within a single group would then result in host level performance that will vary depending on which drive processes the I/O request. It is important to note that even when mixing drives of different speeds in the same disk group, the disk group does not slow down to the speed of the slowest drive in the group.

Although variations in performance are generally not desirable, the overall performance of a LUN is dependent on the number of physical disks accessed by that LUN. If drives of differing speeds are separated into different disk groups, then it follows that LUNs in each disk group will have fewer drives. As such, the performance of each LUN will be lower that it would have been with all drives in the same group. Since there is almost always some amount of I/O load imbalance between different LUNs, the total performance of two LUNs in separate disk groups will be less than the total performance of those same two LUNs in a single large disk group.

> **Performance best practice: Drives with different performance characteristics may be placed in the same disk group, which will result in higher system performance than if they were in separate disk groups.**

## Mixed disk capacities

In a manner similar to disk speeds, disks with different capacities would seem to be better placed in separate disk groups. As with the preceding advice however, more disks will usually produce better performance, so placing disks of differing capacities in a single disk group will result in overall higher performance of the LUNs in that group.

There is a minor issue dealing with leveling on the EVA however. The EVA will attempt to ensure that the amount on each physical disk drive is proportional to that drive's contribution to the overall capacity. This means that larger drives will have more data on them than smaller drives. As an example, a 72 GB disk will have twice as much user data on it as a 36 GB drive. In a random access type of application, this implies that the larger drives will have twice as much I/O as the smaller drives, resulting in an I/O load imbalance at the disk drive level.

Although more I/O to a drive implies a higher response time, this must be weighed against the fact that for a given I/O load, more drives under a LUN equate to a lower per-drive I/O rate. This is important when there is a load imbalance between LUNs, since one LUN in a disk group may be at the maximum sustainable I/O rate, while a LUN in a different disk group might be operating at a fraction of its potential.

As an example, assume that both 36 GB and 72 GB disks are capable of 150 requests per second at a reasonable response time. If an EVA were configured with 84 of each of these drives in a separate disk group, then a single LUN in a single disk group would be capable of 84 x 150, or slightly over 14,000 IOPS. Although the total of both groups would be twice that value, this would only happen if there were a perfect balance of I/O requests between the two groups. Since this rarely happens in practice, combining all disks in the same group would allow the I/O rate of a single LUN to increase over that of a LUN in a single, smaller disk group. The total I/O rate may not be doubled, but the overall realizable rate would be much higher than that of two smaller disk groups.

> **Performance best practice: Drives with different capacities may be placed in the same disk group and will usually result in higher performance than if they were in separate disk groups.**

## Read cache

One of the parameters that may be set at either LUN creation or dynamically at a later date is read caching. The parameter affects both random access read caching and sequential (prefetch) caching, although the algorithms and cache usage are completely different.

Both algorithms are designed to come into play only when they will have a positive effect on performance. Random access caching will be enabled and disabled automatically as the I/O workload changes, while prefetch caching will only come into play if a sequential read stream is detected. Because of this dynamic response to changing I/O workload conditions, cache efficiency is maintained at a high level, and there will be no negative impact on either cache usage or performance in the presence of an I/O workload that is "cache unfriendly".

Since there is no negative impact of leaving cache enabled, and there is always a chance of a performance gain through caching, read cache should always be left enabled.

> **Performance best practice: Always leave read caching enabled on a LUN.**

## LUN balancing

Although both controllers in an EVA can access all physical disks, a LUN is online to only one controller at a time. Because of this, a single LUN can only use the cache of a single controller, may be constrained by the processing capability of a single controller, and can only use two out of the four host ports on the EVA. As such, it makes sense to ensure that each controller in an EVA pair has an equal share of the I/O load.

Although the default controller preference is established via the element manager, this is only a preference, and is usually controlled instead by the host operating system. Because of this , attention should be paid at the operating system level to ensure that the I/O load on both controllers is reasonably balanced.

> **Performance best practice: Always attempt to balance LUNs between the two controllers on an EVA based on the I/O load.**

# Best practices for optimizing cost

The cost of ownership of the storage is the cost per MB (or GB, TB, or more) of the entire storage subsystem. It is obtained by dividing the total cost of the storage by the useable data as seen by the customer. Items affecting cost of ownership are the protection level, number of disk groups, disk quantity, and disk types.

## Protection level

As described in the Availability section, having one or two large disks in a group means that extra space will be reserved. For example, if a disk group of 168 disks consisted of all 36GB disks with the exception of two 72GB disks, a protection level of 2 would reserve 288GB of space (72GB x 4), even though the vast majority of the disks require only 144 of protection (36GB x 4). As such, cost considerations dictate that a disk group should consist of disks that are all the same size.

In addition, VRAID1 mirror partners can only use the capacity of the smallest disk in a mirror pair. If, for example, a 36GB drive was mirrored with a 72GB drive, only 36GB of data could be used for mirroring on the 72GB drive, potentially wasting large amounts of space.

> **Cost of ownership best practice: Do not mix disk sizes in a single disk group.**

Since protected space cannot be shared across disk groups, having multiple disk groups can result in excessive protected space. As an example, consider a 2C12D with 168 disks. With a single disk group and a protection level of 2, there will be a total of 4 disks worth of space set aside. With two disk groups, each group will have 4 disks worth of reserved space, or a total of 8 disks worth of space reserved. As more and more groups are added, the amount of reserved space increased, reducing the user data capacity, and increasing the effective cost of ownership.

> **Cost of ownership best practice: Use a single disk group.**

## Number of disk groups

Spare (unassigned) capacity is the space in a disk group that remains after space has been reserved for the protection level. It is used for creation of virtual disks, snapshots, and snapclones, as well as temporary space used by the EVA for certain internal operations.

Spare capacity decreases when virtual disks, snapshots or snapclones are created, or when physical disks are removed. Spare capacity increases when virtual disks or snapshots are deleted, or when additional physical disks are added to the disk group.

Spare capacity is also used by the system to reconstruct data from redundancy information (VRAID1 and VRAID5 only) due to a physical disk failure. If there is spare capacity available, the system will use this before accessing the space reserved by the protection level.

Spare capacity, like protected space, exists within a disk group, and cannot be shared across groups. If there are two or more disk groups, this can lead to what is known as *stranded* capacity. This happens when it is desired to create a LUN in a disk group, and although the total spare capacity of all disk groups in the storage system is sufficient for the creation of this LUN, there is insufficient capacity in any one of the individual groups to create it. The only solution for this is a massive reconfiguration of the storage, or the addition of disks to increase the spare capacity.

**Cost of ownership best practice: Use a single disk group.**

## Disk quantity

Since a fixed price is paid for the controller and supporting infrastructure, it makes sense to amortize this cost over as much storage as possible. To this end, it is reasonable to use as many disks as possible in a single configuration.

**Cost of ownership best practice: Fill the storage system with as many disk drives as possible.**

## Disk types

When looking at disks, larger disks usually offer better price per capacity. Although prices continuously change, it is reasonable to assume that at any point in time, you can purchase more capacity for the same price when using 72GB drives than with 36GB drives. Similarly, higher performance drives, such as 15K rpm drives, are generally more expensive than their lower performance 10K rpm counterparts.

**Cost of ownership best practice: Use lower performance, larger capacity disks wherever possible.**

# Best practices summary

All of the preceding recommendations can be summarized in tabular format. This not only makes it relatively easy to choose between the various possibilities, it also highlights the fact that many of the "best practice" recommendations contradict each other. In many cases, there is no correct choice, since the best one depends on what the goal is; cost, availability, or performance. Note also that in some cases, a choice has no impact.

Here is how the first line of the table would be read: to optimize cost and availability, do not mix disk capacities in a disk group, but it is acceptable to mix disk capacities when optimizing performance.

|  | Availability | Performance | Cost |
|---|---|---|---|
| Mixed disk capacities in a disk group | Avoid mixing | Acceptable to mix | Avoid mixing |
| Number of disk groups | 1-2 | 1 | 1 |
| Number of disks in a group | Multiple of 8 (Robust Availability Configuration) | Use maximum | Use maximum |
| Total number of disks | Multiple of 8 (Robust Availability Configuration) | Use maximum | Use maximum |
| Higher performance disks |  | OK to use | Avoid using |
| Write cache mirroring | Enable | Maybe disable |  |
| Mixed disk speeds in a disk group |  | Acceptable to mix |  |
| Read cache |  | Enable |  |
| LUN balancing |  | Always balance |  |

# Disk group sizing

When an EVA subsystem is deployed, one of the configuration questions that must be resolved is the number of physical disks required to deliver the desired *usable capacity*, that is, the capacity of the virtual disk as seen by the hosts to which it is presented.

The raw physical storage capacity available in an EVA subsystem is consumed for a number of purposes. The first and most obvious is the storage of the data written by the operating system and applications. Some of the physical storage, however, is used to store the information that makes the subsystem's fault tolerance and virtualization features possible.

---

**Note**

The material in this discussion comes largely from the white paper, *Sizing Disk Groups in a StorageWorks Enterprise Virtual Array by Compaq*.

---

## Disk group sizing factors

The following factors impact the translation of raw capacity into usable capacity:

1.  Hardware versus software storage representations

2.  System metadata overheads

3.  VRAID redundancy overheads

4.  Spare capacity

5.  Snapshot working space

Therefore, sizing a disk group involves determining how much usable capacity is required, accounting for the fixed-system overheads, then factoring in the variable overheads associated with the VRAID type, spare capacity, and snapshot activity.

### Hardware versus software storage representations

Most hardware storage capacity (including that of the EVA) is quoted using a decimal representation of bytes. Many operating systems, including those from Microsoft, use a binary representation (power of 2). Whereas 1 gigabyte is 1,000,000,000 bytes as a decimal representation, the binary representation is $2^{30}$, or 1,073,741,824 bytes. Similarly, 1 terabyte is 1,000,000,000,000 bytes decimal, whereas the binary representation is $2^{40}$, or 1,099,627,776 bytes.

The difference means that 1 physical gigabyte must be consumed for every 0.93 software gigabyte of usable capacity delivered, and that 1 physical terabyte must be consumed for every 0.91 software terabyte of usable capacity delivered. For EVA drives, that can mean using either fourteen or sixteen 72GB drives. The disk sizing formula accounts for this difference.

## System metadata overheads

The subsystem stores its configuration, the tables that map virtual disks to specific physical disks blocks, and other system metadata on the disk drives. At disk group creation, the EVA keeps a minimum of five copies of metadata. Thereafter, metadata is duplicated in up to a maximum of 16 disk groups. It has been determined that approximately 0.2% of the physical capacity is consumed for this purpose. The disk sizing formula accounts for this overhead.

## VRAID redundancy overheads

The parity used to protect VRAID1 and VRAID5 data from a disk failure also requires storage. VRAID1 data is stored twice, and thus consumes two blocks of physical capacity for every block of usable capacity. VRAID5 data stores one block of parity for every four blocks of data, and thus consumes 1.25 blocks of physical capacity for every block of usable capacity. VRAID0 data has no parity protection, and thus consumes only one block of physical capacity for every block of usable capacity.

## Spare capacity

When a disk drive fails, the subsystem reconstructs the missing VRAID1 and VRAID5 data, drawing from unused physical capacity in the affected disk group. To ensure that sufficient spare capacity is available to reconstruct all of the affected data, the system reserves physical capacity equivalent to twice the largest disk in the disk group for each level of disk failure protection selected.

For example, if your disk group has five 36GB disks and five 72GB disks, and you need double protection, four 72GB drives are set aside for spare capacity. If you need single protection, two 72GB drives are set aside.

> **Note**
>
> Spare capacity is not used unless available capacity in the disk group is depleted. Spare capacity is **true** spare space reserved for times when the disk group is nearly full.

## Snapshot working space

An installation that creates snapshots or snapclones on a regular basis needs free capacity from which these new disks can draw. Snapclones and standard snapshots consume the same physical capacity as the original virtual disk. This additional capacity is consumed as part of the operation that creates the snapclone or standard snapshot.

Virtually Capacity-Free Snapshots (Vsnaps) consume physical capacity only when new data is written and capacity is allocated from the disk group in 1GB increments. When the Vsnap is first taken, both the Vsnap and the original target disk contain the same data, and thus share the same physical storage. Each time a block on the original target disk is written subsequent to the snapshot operation, the snapshot contents are preserved by copying the original data to a new location, after which the original target virtual disk block is updated with the new write data.

Therefore, a Vsnap consumes physical capacity as a function of the rate at which the original target disk is changed. A Vsnap for a virtual disk that changes at a slow rate consumes little additional storage; a Vsnap for a virtual disk that is completely rewritten subsequent to the snapshot operation consumes the same physical capacity as the original virtual disk.

An estimate of the usable capacity required for a Vsnap is calculated in two ways:

1.  Use an estimate of the rate at which new data will be written to the Vsnap's original virtual disk, and the length of time the Vsnap will be in existence:

    Virtual Disk Write Rate x Snapshot Duration

    For example, if new data is being written to a VRAID5 virtual disk at 10 MB/s, and the Vsnap is to exist for two hours, the additional VRAID5 usable capacity consumed by the Vsnap will be:

    10MB/s x 7200s = 72GB

    The above applies no matter how many snapshots exist for a specific target virtual disk, as long as the Vsnaps themselves are not being actively written. The original data need only be copied once, so in this case the estimating rule is insensitive to the number of snapshots.

2.  Use the 80/20 rule, that is, usable capacity for a Vsnap will be 20% of the total capacity of the disk group.

# Formula for determining disk count

A formula has been determined to define the approximate relationship between hardware disk capacity, number of disks, and usable capacity for each of the three VRAID types. To solve for the number of disks, supply the usable capacity for each of the three VRAID types, the disk drive capacity, and the protection level desired.

The formula assumes disk groups that follow the guidelines described in the topic called Disk Group Structuring. In particular, the groups are assumed to consist of a multiple of six or eight drives of identical or similar capacity.

The formula to determine disk count is:

DiskCount $\cong$
 ((UsableV0 x 538) + (UsableV5 x 673) + (UsableV1 x 1076)) ÷ (DiskCap x 476) + (ProtLevel x 2)

Where each of the variables is defined as follows:

| | |
|---|---|
| DiskCap | Disk drive capacity in hardware GB |
| DiskCount | Integer number of disk drives |
| ProtLevel | 0 for None, 1 for Single, 2 for Double |
| UsableV0 | Desired usable VRAID0 capacity in software GB |
| UsableV1 | Desired usable VRAID1 capacity in software GB |
| UsableV5 | Desired usable VRAID5 capacity in software GB |

For example, if you need 1TB of usable VRAID0, 1TB of usable VRAID5, and 1TB of usable VRAID1, all using 72GB disks, and you need double protection level, you can calculate the disk count you need:

DiskCount $\cong$ ((1000GB x538) + (1000GB x 673) + (1000GB x 1076)) ÷ (72GB x 476) + (2 x 2)

$\cong$ 66.7 + 4

$\cong$ 70.7 disks

$\cong$ 71 disks

---

**Note**

Remember to convert to the same units, for example, 1TB = 1000GB.

---

The above formulas are approximations because of the dynamic nature of physical capacity use in a virtualized environment. The variance inherent in dynamic algorithms results in a small variability in the outcome. Most of that variability has been accounted for in the formulas, and while not exact, the results should be generally applicable.

## Disk group sizing procedure summary

To determine the number of physical disks required for the disk group:

1. Determine the physical capacity of the drives to be used, and the disk failure protection level required for the disk group.

2. Determine the number of virtual disks that are to reside in the disk group, the usable capacity required for each, and the VRAID type required for each.

3. Determine the amount of additional usable capacity required for snapshot working space for each virtual disk that will have snapshots. The VRAID type of the snapshots will match the original.

4. Sum the total usable capacity required for each VRAID type. This is the sum of the usable capacity for each virtual disk of that type and its corresponding snapshot working space (if any).

5. Solve for the number of disks using the formula.

6. Add margins appropriate to the circumstances. Subsystem management flexibility improves significantly as free capacity grows.

# EVA documentation tool

The EVA tools program is designed to provide a set of powerful, interoperable support tools for the field, support centers, and escalation centers. It is based on the XP tools program that contains documentation (XPDT), configuration (XPCT), and performance tools (XPPT). At present, the EVA tools program consists of the Windows-based EVA documentation tool (EVA-DT), V1.20.0000, supporting VCS versions 2.X and 3.X. There are also plans to build a configuration tool (EVA-CT), performance tool (EVA-PT), and capacity planning tool (EVA-CPT).

## General description

The purpose of EVA-DT is to give field support a way to acquire EVA configuration information and documentation. Instead of looking into raw data files and displays, this tool will allow a field engineer to create a visual representation of the data and present it in a variety of ways. The tool will be able to browse all areas of a configuration and produce graphical images of the array from a physical and logical perspective. The configuration can then be analyzed and described to the customer, and professional documentation will be available.
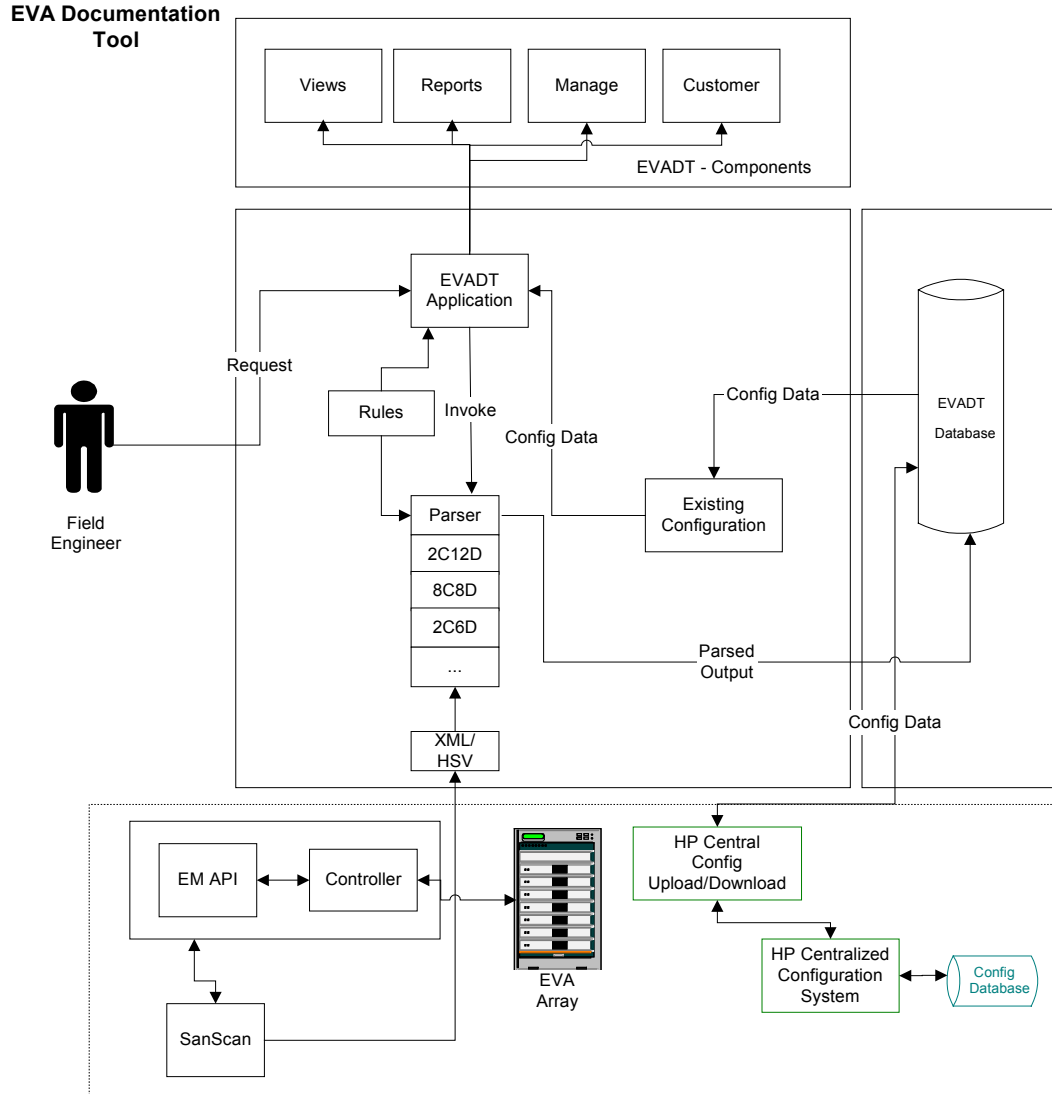
## Functionality

EVA-DT is a Windows-based tool that extracts array configuration information provided by SANscan into both logical and physical graphical views of an array. The tool will provide customizable reports of the array data as well as offer differing views of the array based on a range of criteria.

EVA-DT will support all the EVA configurations depending on the hardware version available in the configuration file. The following are some other considerations about the tool:

- It does not interact with the array or Command View EVA

- It is not an online tool

- It does not offer new solution, but helps the field to propose new solutions manually

- It does not integrate with other EVA support tools

# Block diagram

The following block diagram gives an overview of how the tool works.

**EVA Documentation Tool**

| Views | Reports | Manage | Customer |

EVADT - Components

EVADT Application

Field Engineer

Request

Rules

Invoke

Config Data

Config Data

EVADT Database

Existing Configuration

Parser
2C12D
8C8D
2C6D
...

Parsed Output

Config Data

XML/ HSV

EM API — Controller

EVA Array

SanScan

HP Central Config Upload/Download

HP Centralized Configuration System

Config Database

The following topics describe key components of the diagram.

## Users

The target users for EVA-DT are:

- Pre-sales Technical Consultants (TC) — Helps customer choose EVA array hardware and software.

- Field Engineers — Account Support Engineer (ASE) or Customer Engineer (CE).

- Escalation personnel

- Consulting and Integration (C&I)

## Configurations

The current supported configurations are generated as XML/text files from the SANscan generator. The EVA-DT tool will parse these configurations and dump the data into the database.

Supported configurations are:

- Model 2C6D, 2C6D-A, 2C6D-B

- Model 2C12D, 2C12D-A, 2C12D-B

- Model (2) 2C6D + Expansion Model 0C6D, (2) 2C6D-A + Expansion Model 0C6D-A, (2) 2C6D-B + Expansion Model 0C6D-B

- Model (2) 2C12D + Expansion Model 0C12D, (2) 2C12D-A + Expansion Model 0C12D-A, (2) 2C12D-B + Expansion Model 0C12D-B

In general, EVA-DT should support all the EVA configurations depending on the hardware version available in the configuration file.

## SANscan generated files

The SANscan application will interact with Command View EVA (HSV Element Manager) through the Command View EVA API to collect all the data related to the EVA configured array. The data files generated by SANscan are the inputs for EVA-DT.

## EVA-DT tool

When a **new** configuration is selected, the tool will parse the selected configuration file based on the rules and dump the data into a database. After the data is available in the data repository, the user can view all the provided features (views, reports, support) related to the configured array.

When an **existing** configuration is selected, the tool will provide all the available features related to that configuration like physical view, associate View and reports.

## Hardware and software requirements for installation

Before installation, you should have the following hardware:

- Compatible Windows-based PC

  - Pentium 2, 3, or 4 processor

  - Minimum processor speed of 400MHz

  - Minimum 128MB of main memory

  - Minimum display resolution of 1024 x 768

  - Minimum 100MB of hard disk storage

You also need the following software:

- Windows 2000 Professional SP1/SP2/SP3, Windows XP Home, or Windows XP Professional

- NET Framework V1.0 or above (V1.1 recommended)

- Microsoft Data Access Components (MDAC) V2.6 SP1 or above

- MS Jet V4.0 SP3 (included in Windows 2000 or XP)

- Microsoft Internet Explorer V5.5 SP1 or above

---

**Note**

MS Jet is part of Windows 2000 or XP. If a user is using some other operating system, he may need to install MS Jet. MS Jet applications currently available through the download page include Windows 9X, NT, and Millenium.

---

# Installing EVA-DT

You must do the following to install EVA-DT:

- Install SANscan. This is the program that creates the EVA configuration files for input to EVA-DT.

  For information on downloading SANscan, go to the SANscan download page at **http://storagetools.lvld.hp.com/fibre/sanscan/download.asp**

  Here, you can download SANscan V1.4.0.

- Go to the EVA-DT download page at **http://storagetools.lvld.hp.com/eva/evadoctool/download.asp** and do the following:

  - Install .NET framework, if necessary

  - Install MDAC V2.6 SP1, if necessary

  - Install EVA-DT executable (also installs VS Flex Grid for .NET and MetaDraw V3.0)

  - Review the Read Me file

    **Note**

    To get the latest EVA tools, you should register at the Tools portal at **http://15.2.237.51/.**

# Running EVA-DT

Before running the actual EVA-DT application, you must run SANscan on any host or PC that talks to the EVA. SANscan uses the Command View EVA API to collect information on the array. The resulting configuration file is input to EVA-DT.

**Note**

For configurations that have more than one EVA attached to an SMA, SANscan will generate multiple files. Each of these files will need to be run through EVA-DT to generate configuration files.

You can invoke EVA-DT in two ways:

- By clicking on the desktop icon to present the GUI

- From the command line prompt

## Using the command line

The primary purpose of using the command line is to interact with other tools such as SANXpert. To use the command line, use the following syntax:

```
EVADT.exe  -i = "<Input File>"  -o = "<Output path>"
[xml] [-w=y]
```

The following describes the parameters:

<Input File> = absolute path of the input file from SANscan

<Output Path> = path for the output folder

[xml] = Optional parameter to generate sortable HTML configuration files

[-w=y] = Optional parameter for overwriting the output file

### Example

```
EVADT.exe -i = "c:\EVA1_10_2_1_18.xml"  -o = "c:\text".
```
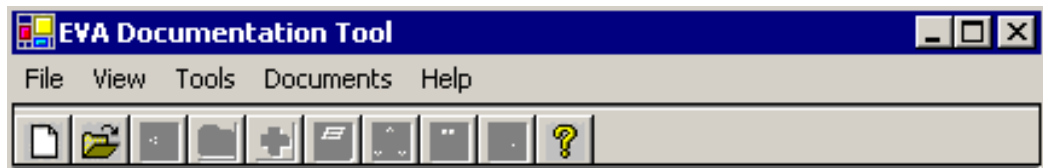
The XML file is parsed and the XML and HTML files are placed in \text.

---

**Note**

When using the command line, be sure you are in the folder where evadt is resident (probably c:\Program Files\Hewlett-Packard\EVA Documentation Tool).

---

## Using the GUI

When you click on the desktop icon, you have access to all of the EVA-DT operations from the main menu.

Menu options include:

- File options

  - New — Parses a new configuration (.xml)

    When you select *File, New*, you are asked to supply the XML file created by SANscan to create configuration files (XML and HTML).

  - Open — Opens an existing configuration by using *File, Open*Close — Closes the open configuration

  - Edit Title — Allows editing of Associative View title

  - Copy to Clipboard — Copies the physical view of the currently displayed configuration to the clipboard

  - Save Graphic — Saves the Associative View in a .bmp, .jpg, .emf, or .wmf format

  - Page Setup — Displays page options for printer

  - Print — Prints the views or reports

  - Exit — Exits the EVA-DT application

- View

  - Associative View — Displays the logical relationships of the array

  - Array Global Report — Provides information about the entire array

  - Physical Disk Report — Shows information about the disks

  - Detail Report — Shows columns of array information by logical units

  - Disk Group Report — Shows information about different groups

  - Virtual Disk Report — Shows information about virtual disks

  - Host Report — Shows information about the hosts

- Tools

  - Logging — Turn logging on or off

  - Toolbar — Toggle display of toolbar

  - Options — Allows changes to the display of the physical view

■ Documents — Generates the reports in an HTML format

■ Help

- EVA-DT Online User Guide (HP Intranet only)

- Feedback (HP Intranet only)

- Download Page (HP Intranet only)
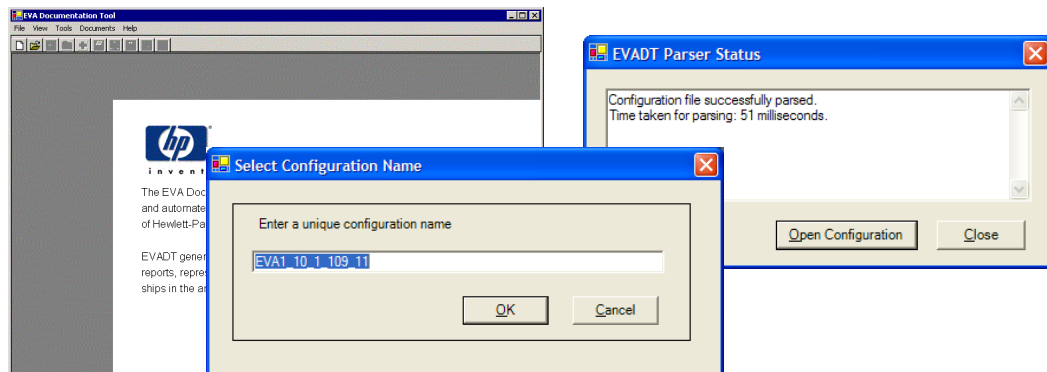
- About EVA-DT — Version information about EVA-DT

## Sample screens

The following are some of the main screens you will see when using EVA-DT.

### Creating a new configuration

When you click on the EVA-DT icon on the desktop, the EVA Documentation Tool main screen displays. Use the following steps to create a new configuration:

1. Select *File* → *New* from the menu or click the *New* icon.

2. Browse to and open the **.xml** configuration file you want. The Select Configuration Name dialog box displays.

3. Enter a unique configuration name and click *OK*. The EVA-DT parser creates the configuration. This may take a minute for large files.

4. Click *Open Configuration*. The Physical View for the configuration displays.

## Physical View

The Physical View displays the actual physical arrangement of the components of an array. When you open a configuration, the Physical View appears in the main window of EVA-DT. You can modify the physical viewing characteristics by using *Tools → Options*.
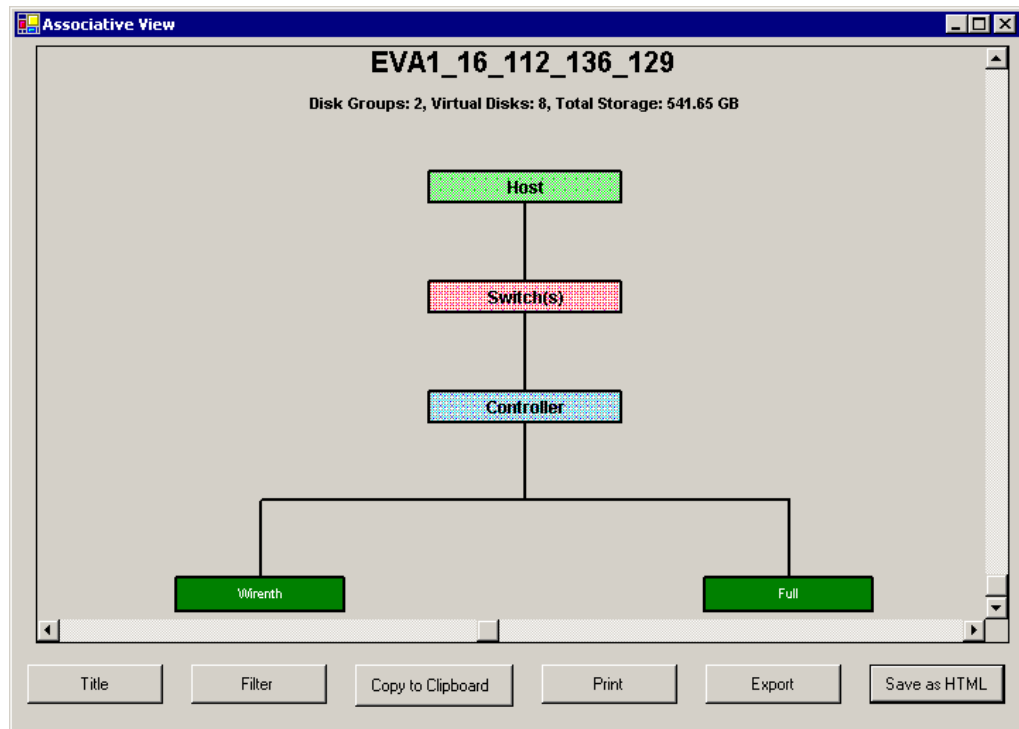
The following are components in the physical view, of which several are clickable:

■ System Information — Displays beneath the title showing the Total Storage and Used Storage.

■ Cache (clickable) — Displays Controller Cache information.

■ Controller (clickable) — Displays hardware information about the controller.

■ Port (clickable) — Displays the Controller Port information (Host Ports and Device Ports).

■ Individual Disks — Numbered 1-14 or marked E for empty. A disk displaying an L is loaded, but its attributes and properties are unknown. Placing the cursor on the disk displays the Bay number and Shelf number. Disks are color-coded to indicate their group, as shown in the legend.

■ Backend Loop Switch — There can be four loop switches that form part of the array, two above the controllers, and two below.

■ Shelf number — Displays the shelf number on either side of the shelf. Numbers range from 1-13.

## Associative View

The Associative View displays the logical relationships of the components of an array. In the EVA-DT main window, select *View → Associative View*, or click on the icon in the toolbar.

Functions available through buttons at the bottom of the Associative View include:
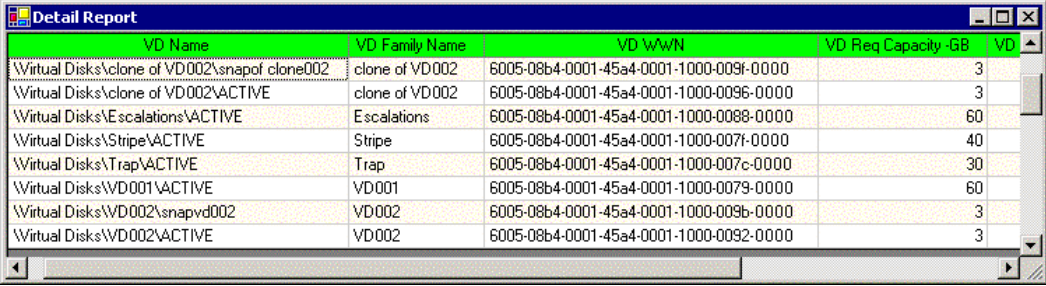
- Title — Allows editing of the title of the Associative View. Not saved to the database.

- Filter — Allows filtering of the information displayed in the Associative View (not yet implemented).

- Copy to Clipboard — Copies the Associative View to the clipboard in **.emf** or **.wmf** format so you can paste it into other applications.

- Print — Prints the Associative View to the default printer.

- Export — Saves the Associative View as a graphic (.bmp, .jpg, .emf, .wmf).

- Save as HTML — Saves the Associative View in an HTML format.

- Close — Closes the Associative View window.

There are several components in the Associative View, of which several are clickable.

- Disk Group (clickable) — Displays disk group information.

- Controller (clickable) — Displays controller information.

- Switches — Switches with the array.

- Host (clickable) — Displays host information.

## Reports

There are several available reports in EVA-DT, of which one is the Detail Report. Select *View → Detail Report* to display this report.



The detail report displays columns of array information delineated by logical units, showing the relationship between virtual disks, disk group, LUN, and the host.

---

**Note**

For all reports, you can sort the data by clicking on the column headers and you can drag the columns wherever you want.

---

# EVA troubleshooting tool

The EVA troubleshooting tool, version 01.00.00, is an online help application that allows you to drill down into EVA service documentation to diagnose and solve onsite or remote problems. It allows you to access such items as troubleshooting information, hardware information, error messages, manuals, and best practices. This information is point in time, that is, newer copies of information or documentation need to be updated as the product changes.

> **Note**
> This tool is available to HP storage-trained specialists or customers, however, customers should use the tool under HP supervision.

## Installing the troubleshooting tool

Installation requirements for the troubleshooting tools include the following:

- Internet Explorer 5.0 or higher

- Minimum 128MB of main memory

- Minimum display resolution of 1024 x 768

- Minimum 40MB of hard disk storage

To install and open the tool, do the following:

- Install the tool from the EVA Troubleshooting Tool Review Page at **http://tce-web.boi.hp.com/prod_port/review/nsas_eva_ts_tool.htm**Choose the second link down to install the tool (CE Version)

- Open the tool by clicking on the desktop icon

# Using the troubleshooting tool

When you click on the desktop icon, you will see the following screen.



The main page is divided into two areas:

- Table of Contents on the left. This is the default but you can also use the Search button to search for a specific item.
- Topic contents on the right

Use the screen as you would an online help system.

# Storage calculation utilities

A few utilities are available for you to determine storage capacities for your Enterprise storage system configuration. These tools enforce the current maximums and minimums known for any Enterprise storage system while providing appropriate factors for virtualization overhead. The TechBB for the EVA contains access to these utilities.

## Disk group sizing utility

This Microsoft Excel spreadsheet allows you to determine the net storage capacity of an Enterprise storage system resulting from the selection of up to 16 disk groups, their VRAID mixes, and their protection levels. The impact of virtualization overhead is included through a VCS impact factor, and you are allowed to select the disk size (36GB or 72GB).

## Disk groups and virtual disks calculator

This Microsoft Excel spreadsheet allows you to interactively choose an EVA configuration with a number of drives and capacity per drive. For this configuration, you interactively select the number of virtual disks per disk group (up to 16), and for each of those virtual disks (up to 512), select the usable capacity of the virtual disk and the VRAID type to determine the capacity used and the resulting free capacity. Virtualization overhead is accounted for in the calculations.

# Learning check

1.  What are the three factors that must be considered when configuring disk groups?

    ....................................................................................................................................

    ....................................................................................................................................

    ....................................................................................................................................

2.  List the best practices with respect to mixing disk capacities in a disk group.

    ....................................................................................................................................

    ....................................................................................................................................

    ....................................................................................................................................

3.  List the best practices with respect to the number of disks in a disk group.

    ....................................................................................................................................

    ....................................................................................................................................

    ....................................................................................................................................

4.  List the five factors that determine the number of disks required to deliver a given usable capacity in an Enterprise storage system.

    ....................................................................................................................................

    ....................................................................................................................................

    ....................................................................................................................................

    ....................................................................................................................................

5.  To deliver 2TB of usable VRAID0 and 1TB of usable VRAID5, and using 72GB drives with single protection, what is the approximate disk count that you would need?

    ....................................................................................................................................

    ....................................................................................................................................

    ....................................................................................................................................

HP Restricted