

EMC CLARiiON Asymmetric Active/Active Feature (ALUA)

A Detailed Review

Abstract

This white paper provides an overview of the EMC[®] CLARiiON[®] Asymmetric Active/Active feature. It highlights the configuration, best practices, implementation details of EMC PowerPath[®], and native multipathing software when the host initiator is configured in Asymmetric Active/Active mode.

March 2008

Copyright © 2007, 2008 EMC Corporation. All rights reserved.

EMC believes the information in this publication is accurate as of its publication date. The information is subject to change without notice.

THE INFORMATION IN THIS PUBLICATION IS PROVIDED “AS IS.” EMC CORPORATION 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 EMC software described in this publication requires an applicable software license.

For the most up-to-date listing of EMC product names, see EMC Corporation Trademarks on EMC.com

All other trademarks used herein are the property of their respective owners.

Part Number H2890.2

EMC CLARiiON
Asymmetric Active/Active Feature (ALUA)
A Detailed Review

Table of Contents

Executive summary	4
Introduction	4
Audience	4
Terminology	4
Overview of the Asymmetric Active/Active feature	5
Configuration.....	5
Implementation	5
Layered drivers impact.....	6
Explicit vs. implicit trespass	7
Configuring for ALUA.....	7
Configuring the host.....	7
Using PowerPath with ALUA.....	7
Using native failover with ALUA	8
Configuring the CLARiiON	9
CLI commands to configure Active/Active.....	10
CLARiiON Failover Modes	10
Getting the ALUA state from a CLARiiON.....	10
Storage system support	11
Performance considerations	12
Best practices	12
Impact of the Asymmetric Active/Active feature.....	12
Benefits of CLARiiON Asymmetric Active/Active.....	12
Improved customer experience	12
Supports standard SCSI multipath interfaces	14
Impact of the Redirector.....	14
Back-end failure masking	14
Comparison of failover methodologies	15
Explicit and implicit failover software	15
PowerPath.....	16
Native failover software (MPIO and others).....	17
Nondisruptive upgrade (NDU).....	17
Manual trespass.....	17
Path, HBA, and switch failure	18
SP failure.....	18
Conclusion	18
References	19

Executive summary

In FLARE® release 26, EMC introduced the Asymmetric Active/Active feature for CLARiiON® storage systems. Asymmetric Active/Active provides a powerful new way for CLARiiON storage systems to present logical storage units (LUNs) to hosts, and eliminates the need for hosts to use the LUN ownership model. This changes how a host manages multiple communication paths to LUNs on the array (commonly referred to as *path management*) by permitting I/O to flow to either or both storage processors. This white paper discusses the benefits and implementation details of this feature.

Introduction

This white paper provides a technical overview of the Asymmetric Active/Active feature on the CLARiiON, and discusses how PowerPath® and native MPIO (Multi-Pathed Input/Output) software work with this feature.

Prior to release 26, all CLARiiON storage systems used the standard active/passive feature in which one SP *owns* the LUN, and all I/O to that LUN is sent to that SP. If all paths from a host to the SP fail, host-based path-management software adjusts the I/O path by issuing a *trespass* command. This causes the storage system to change the ownership of the LUN to the peer SP, and I/O is then sent to the peer SP.

CLARiiON Asymmetric Active/Active is a feature that introduces a new initiator Failover Mode (Failover Mode 4). When configured as Failover Mode 4, initiators can send I/O to a LUN regardless of which SP owns the LUN. While this feature allows a trespass-like command, explicit trespass behavior is not required.

This white paper replaces the paper *EMC CLARiiON Asymmetric Active/Active Feature*.

Audience

This white paper is intended for customers, partners, and EMC field personnel who want a better understanding about the implementation, benefits, and configuration of the CLARiiON Asymmetric Active/Active feature and the potential impact in their storage environment.

Terminology

- **ALUA** – Asymmetric Logic Unit Access.
- **CLARiiON LUNs** — Logical subdivisions of RAID groups in a CLARiiON storage system. These are volumes that are presented to hosts.
- **CMI (CLARiiON Messaging Interface)** — Redundant PCI Express connections that provide communication between the two SPs.
- **Failover Mode** — Determines how the array responds to I/O requests that are directed to LUNs on the non-owning SP.
- **Non-optimal path** — A path that is ready to do I/O, but that may not yield the best performance.
- **Optimal path** — A path that is ready to do I/O and will yield the best performance.
- **Preferred bit** — This bit represents that a given SP is the default owner of a LUN
- **SP (storage processor)** — CLARiiON controller.
- **Target port group** — A set of target ports that are in either primary SP ports or secondary SP ports.
- **Trespass** — A command that allows an SP or its peer to take ownership of the LUN.

Overview of the Asymmetric Active/Active feature

CLARiiON Asymmetric Active/Active is based on the Asymmetric Logical Unit Access (ALUA) standard. ALUA uses SCSI 3 primary commands that are part of the standard SCSI SPC-3 specification (not a CLARiiON-specific implementation) to determine I/O paths. In dual-SP systems, like a CLARiiON, I/O can be routed through either SP. For example, if I/O for a LUN is sent to an SP that does not own the LUN, that SP redirects the I/O to the SP that *does* own the LUN. This redirection is done through internal communication within the storage system. It is transparent to the host, which is not aware that the I/O was processed by the other SP. Hence, a trespass is not required when I/O is sent to the non-owning (or non-optimal) SP.

Configuration

Dual-SP storage systems that support ALUA define a set of *target port groups* for each LUN. One target port group is defined for the SP that currently owns the LUN, and the other target port group is for the SP that does not own the LUN. The standard ALUA commands allow the host failover software to determine the state of a LUN's path.

The REPORT TARGET PORT GROUP command reports the following three attributes:

- Preferred - indicates whether this port group is the default (preferred) port group.
- Asymmetric access state - indicates the state of the port group. Port group states include Active/Optimal, Active/Non-optimal, Standby, and Unavailable, and are discussed below.
- Attribute - indicates whether the current asymmetric access state was explicitly set by a SET TARGET PORT GROUP command or was implicitly set or changed by the storage system.

The SET TARGET PORT GROUP command allows the access state (Active/Optimal, Active/Non-optimal, Standby, and Unavailable) of each port group to be set or changed. Access states are defined as below:

- Active/Optimal – best performing path, does not require upper level redirection in order to complete I/O
- Active/Non-optimal – required upper level redirection to complete I/O
- Standby – state is not supported.
- Unavailable – returned for port groups on a SP that is down. It is not settable via set target port groups.

Target port group commands are implemented in the ALUA layer of the storage system. However, it is actually host-based path-management software that executes the commands and manages the paths. This is similar to the traditional path-management mechanisms. However, ALUA has standardized the mechanisms, whereas in the past they were vendor-specific.

Implementation

FLARE 26 includes a new redirector driver that improves CLARiiON availability; this driver consists of an upper and lower redirector for each SP. The upper redirector sits closer to the host connection, while the lower redirector sits closer to the CLARiiON back end. The layered drivers such as Navisphere® Quality of Service Manager (NQM), SnapView™, and so forth sit between the upper and lower redirectors as shown in Figure 1.

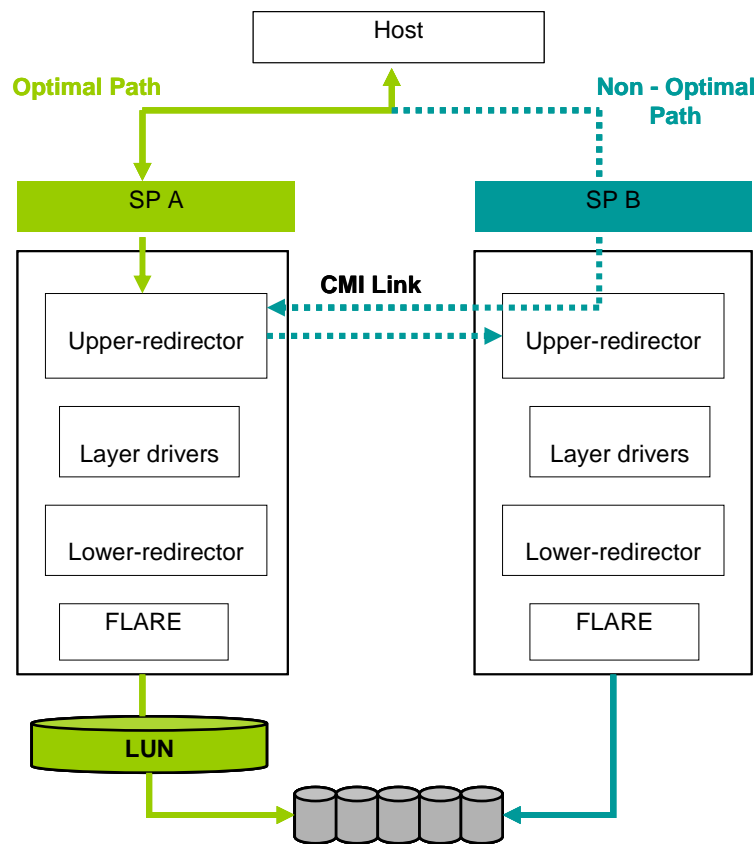


Figure 1. CLARiiON Asymmetric Active/Active feature

CLARiiON's Asymmetric Active/Active feature supports two target port groups: optimal and non-optimal. An optimal target port group represents the ports that belong to the current owner of the LUN. A non-optimal port group represents the ports owned by an SP that is not the current owner of the LUN. This is implemented in the redirector driver.

I/O is accepted on all ports. In event of front-end path failure I/O may be received by an SP that does not own the LUN. The upper redirector of the SP that does not own the LUN routes the I/O to the owning SP through the internal CMI channel. The I/O request is serviced by the SP that owns the LUN. Note that an I/O that is received by an SP has to be *acknowledged* over the same path by that SP.

Load balancing across ports on one SP works just like it does in PowerPath multipathing software. Load balancing across optimal and non-optimal paths is not recommended.

Layered drivers impact

The behavior of the layered drivers (MetaLUNs, LUN Migration, SnapView, MirrorView™ and SAN Copy™) does not change with Asymmetric Active/Active. For example, if a LUN is trespassed during a LUN Migration, the destination LUN is still trespassed by the array and the synchronization continues from a checkpoint.

With SnapView, the source and replica (snap or clone) are always owned by the same SP. Hence, if the source or SnapView replica trespasses, the source or the SnapView replica will still follow its replica or source, respectively. However, depending on the failover software deployed on the host, the *likelihood* of

trespassing the source or SnapView replica is reduced in ALUA mode. See the “Avoids LUN ownership thrashing situations that may occur in SnapView and cluster configurations” section for more information.

For MirrorView/S, if a primary LUN is trespassed, the secondary trespasses immediately. For MirrorView/A, if a MirrorView primary LUN is trespassed, the secondary MirrorView LUN trespasses during an update or at the start of an update.

For local and remote SAN Copy sessions, if the source or destination LUN trespasses, the session needs to be restarted on the peer SP since the SAN Copy initiators do not use ALUA mode. As a result, MirrorView and SAN Copy will behave the same as before.

Explicit vs. implicit trespass

An explicit trespass is the result of an external command from a user or the failover software. When an SP receives this command (from the failover software or a user issuing the LUN trespass in Navisphere), LUN ownership is transferred to that SP. These trespasses can be issued by PowerPath path-management software (using CLARiiON proprietary commands) or other ALUA-compatible path-management software (using the SET TARGET PORT GROUP command).

An implicit trespass is the result of software controls within the storage system. For example, an implicit trespass occurs when the amount of I/O transferred across the non-optimal path exceeds the optimal path I/O by a certain amount (threshold). The software uses counters to keep track of optimal and non-optimal path I/O. When it detects that the non-optimal path has received 128,000 more I/Os than the optimal path, it initiates a trespass.

Configuring for ALUA

There are two options for failover software on the host: PowerPath, and native OS failover software (such as MPIO) that is compliant with ALUA. The following sections discuss both of these options.

To ensure that application data is highly available, the host must be configured to withstand a single point of failure, including a failure in the host bus adapter (HBA), Fibre cable, or failover software. The *EMC CLARiiON Open Systems Configuration Guide* (on Powerlink®) outlines the attach methodologies that CLARiiON supports.

Configuring the host

Many hosts (see Table 1), using their native failover software, can take advantage of ALUA mode on CLARiiON. Upcoming FLARE releases in 2008 will support additional configurations, as noted.

Redundancy for all points of a configuration is essential for optimal high availability (including maximum data accessibility and server uptime). This means that each server must have at least two HBAs (or NICs in an iSCSI environment), and is connected to both SPs. Two switches provide independent, discrete paths to the storage system from the server. If more than two switches are employed, redundant switches should be connected via interswitch links for fabric failover purposes.

Each HBA can be configured to see an SP port (if direct connected) or zoned to see both SPs. This protects your configuration against the loss of an HBA, cable, switch, or SP, and takes advantage of the failover capabilities and load-balancing features of your path-management software.

Having multiple active paths to a LUN available from a server ensures that you can use path-management software load-balancing algorithms to avoid becoming *path bound* (or restricted to a single port).

Using PowerPath with ALUA

PowerPath version 5.1 is the first ALUA compliant release. Ensure that your version is 5.1 or later and consult the E-Lab™ Navigator on Powerlink for various host considerations.

PowerPath load balances across optimal paths. If PowerPath detects all optimal paths have failed, PowerPath initiates a trespass to change the LUN's ownership.

Table 1. PowerPath with native MPIO software

	Native with Active/Passive	Native with ALUA	PowerPath with Active/Passive	PowerPath with ALUA
Longhorn	No	Yes	See Note 1	See Note 1
W2K3	No	No	Yes	Yes
Win2K	No	No	Yes	Yes
HP-UX 11i v1 and v2	Yes	No	Yes	Yes
HP-UX 11iv3	No	Yes (11.31.0709)	No	No
Solaris 9/10	Yes	Yes ¹	Yes	Yes ¹
Linux (RH & SuSE)	Yes	Yes (SLES 10.1 SLES 9 SP4 RH 5.1, 46)	Yes	Yes
AIX	No	No	Yes	No
VMware	Yes	No	No	No

Note 1: Available in the first half of 2008 when used with PowerPath 5.1 SP2. Consult the *EMC Support Matrix* for details.

Definitions for columns:

- Native with Active/Passive: Standard Active/Passive failover (*not* ALUA) is provided when using the noted operating system (OS) alone. In this case, the OS issues trespass commands to enable alternate paths.
- Native with ALUA: ALUA features are provided when using the noted OS alone (PowerPath is not required).
- PowerPath with Active/Passive: Standard Active/Passive failover (*not* ALUA) is provided by PowerPath with the noted OS; PowerPath issues trespass commands to enable alternate paths.
- PowerPath with ALUA: ALUA features are provided when the specified PowerPath release is used with the noted OS.

Using native failover with ALUA

MPIO and other native host-based failover applications can work with ALUA if they are ALUA-compliant. Please note the following:

- For HPUX 11i v3.1, patch 11.31.0709 is required for native failover support with ALUA.
- For Solaris 9, MPIO requires StorEdge SAN Foundation Software 4.4.12.
- Solaris 10 update 3 is required for native failover support with ALUA.

It is important to consult the E-Lab Navigator on Powerlink for up-to-date host considerations. Table 1 lists support for native failover software on various operating systems. Native failover software often does not support features such as autorestore, load balancing, and SAN boot. Please refer to the *EMC Support Matrix* for more information about which software supports these features.

¹ ALUA mode is not supported with Sun Cluster at this time.

Configuring the CLARiiON

To configure Asymmetric Active/Active, the host initiators must be configured with Failover Mode 4. You can set its Failover Mode to **4** using Navisphere Manager or CLI. (For more information about Failover Mode see the “CLARiiON Failover Modes” section.)

As shown in Figure 2, the **Failover Mode** pull-down menu (in the **Group Edit Initiators** dialog box) includes **4** as an option. Since Failover mode is an *initiator* option (rather than *storage group* option) both ALUA hosts and non-ALUA hosts (hosts not configured with a Failover Mode 4) can be attached to the same LUN.

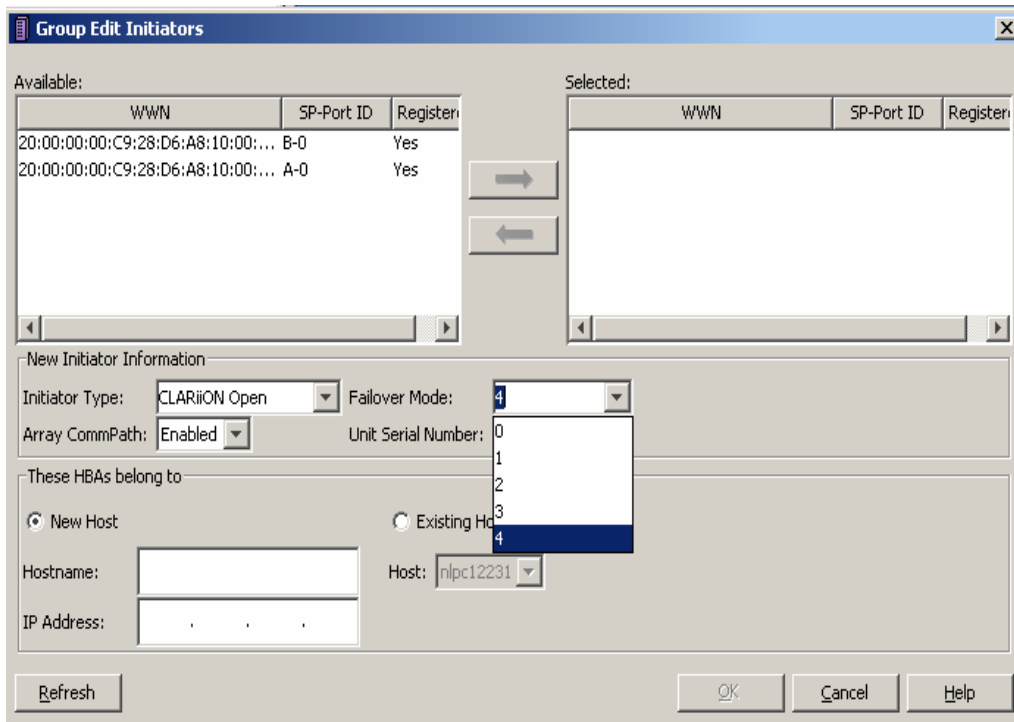


Figure 2. Group Edit Initiators dialog box in Navisphere displaying new Failover Mode for Asymmetric Active/Active

CLI commands to configure Active/Active

Command to set failover mode to 4 (Asymmetric Active/Active) on the host:

```
naviseccli -h <SP_IP_Address> -user a -password a -scope 0 storagegroup -sethost -ip < IP Address> -failovermode 4
```

Command to set failover mode to 4 (Active/Active) for HBAUID:

```
naviseccli -h <SP_IP_Address> -user a -password a -scope 0 storagegroup -setpath -hbauid xxxxx -sp a -spport xxxxx -failovermode 4
```

Command to display default failover mode value:

```
naviseccli -h <SP IP Address> -user a -password a -scope 0 port -list -failovermode
```

Note: After configuring devices for ALUA (Failover Mode 4), a host reboot is required.

CLARiiON Failover Modes

A CLARiiON LUN trespass may be initiated by the storage system, or by the path-management software residing on the server. The Failover Mode of the initiator specifies how the CLARiiON should respond to a trespass condition. CLARiiON supports five distinct Failover Modes depending on the operating system type attached:

- Failover Mode= 0. Auto trespass mode; any media access to the non-owning SP is rejected.
- Failover Mode= 1. Passive Not Ready; a command failure when I/O is sent to a non-owning SP.
- Failover Mode= 2. (DMP mode), Quiet Trespass on I/O to non-owning SP.
- Failover Mode= 3. Passive Always Ready; some commands (for example, Test Unit Ready) return Passive Always Ready status.
- Failover Mode= 4. Asymmetric Active/Active.

Getting the ALUA state from a CLARiiON

The Navisphere CLI **getlun** command displays additional information about the LUN connected to an Asymmetric Active/Active host. An excerpt of the command is shown as follows:

```
# naviseccli -h [SPipaddress] getlun 0
• Read Requests SPA:      236
• Read Requests SPB:     1480
• Write Requests SPA:     426
• Write Requests SPB:     627
• LUN Busy Ticks SPA:     273
• LUN Busy Ticks SPB:     297
• LUN Idle Ticks SPA:     0
• LUN Idle Ticks SPB:     0
• Number of arrivals with non-zero queue SPA: 398
• Number of arrivals with non-zero queue SPB: 398
• Sum queue lengths by arrivals SPA:         398
• Sum queue lengths by arrivals SPB:         398
• Explicit Trespasses: 5800
• Explicit Trespasses SPA: 2346
```

- Explicit Trespases SPB: 3454
- Implicit Trespases: 320
- Implicit Trespases SPA: 214
- Implicit Trespases SPB: 106

The meaning of the output for this command has not changed; however with Failover Mode 4, a LUN reports statistics for both SPA and SPB. Also, the output displays explicit and implicit trespases.

The **LUN Properties** dialog box (shown in Figure 3) has been enhanced to display statistics for LUNs that are connected to initiators with Failover Mode 4. These statistics are available when **Statistics Logging** is enabled on the storage system. The statistics include the number of reads and writes routing through the Optimal and Non-optimal Paths. The Redirector ReAssignment value records the number of implicit trespases by Asymmetric Active/Active since **Statistics Logging** was enabled.

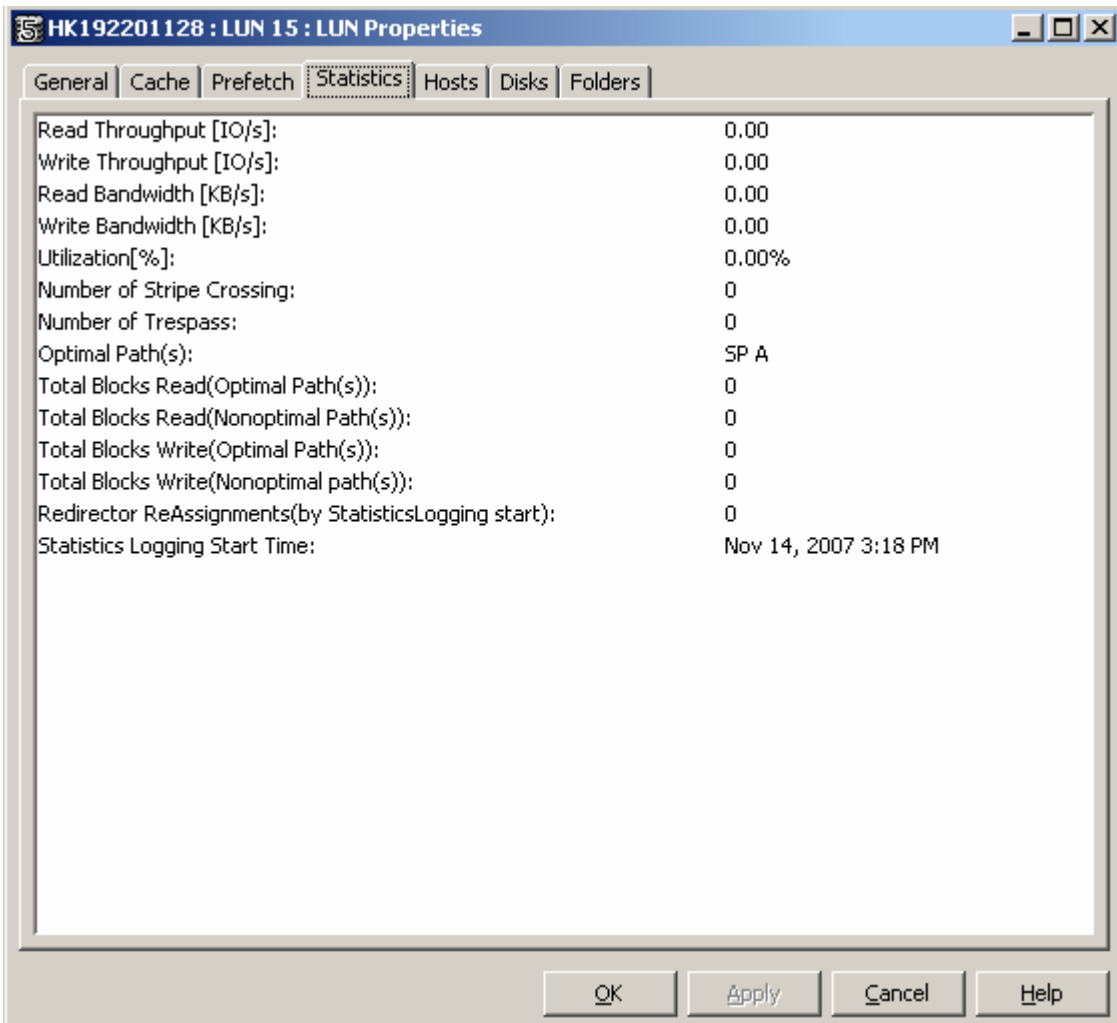


Figure 3. LUN Properties dialog box statistics

Storage system support

The CX3 and CX (700/500/300) storage systems running FLARE release 26 are ALUA compliant. The ALUA functionality is available for both Fibre Channel and iSCSI attach hosts if they are connected using Failover Mode 4 and have ALUA-aware failover software that supports FC and/or iSCSI connections. AX systems do not support ALUA.

EMC CLARiiON

Asymmetric Active/Active Feature (ALUA)

A Detailed Review

Performance considerations

There will be some performance impact if I/O is routed to the non-owning SP (through a non-optimal path). Non-optimal paths are slower; for this reason this is not the preferred method for normal access to the storage. New statistics are available on the array to help you determine if I/O is flowing through the optimal or non-optimal paths. Refer to the *EMC CLARiiON Best Practices for Fibre Channel Storage* white paper for details on performance impact when running I/O through non-optimal paths.

Best practices

- Balance LUN ownership between the two SPs.
- To avoid performance impact when I/O is routed to the non-owning SP of a LUN, EMC recommends that you configure failover software so that it only load balances across the active-optimal path for a given LUN. PowerPath does this by default.
- After SAN changes (component failures, replacements) that may cause I/O paths to change, ensure that hosts are still using optimal paths to their LUNs.
- After an NDU operation, ensure that all LUNs are returned to the default owner. PowerPath does this automatically.
- In case of failure or a performance issue, ensure I/O is routing through the optimal path.

Impact of the Asymmetric Active/Active feature

Asymmetric Active/Active is a request-forwarding implementation that honors the storage system LUN ownership feature (only one controller owns a given LUN); however, it allows I/O to route through either controller. The controller that is not the current owner of the LUN redirects the I/O to the controller that owns the LUN using internal communication paths within the storage system.

There is no benefit of ALUA in event of SP failure as the redirector driver is not accessible to redirect the I/O to a non-owning SP. So, in the event of SP failure, ALUA works as PNR (Failover mode 2) mode.

Benefits of CLARiiON Asymmetric Active/Active

Improved customer experience

The CLARiiON Asymmetric Active/Active feature:

- **Avoids unavailability of boot from SAN during path failure**
ALUA avoids issues when the BIOS attempts to boot from non-optimal path and there is no failover software available (because the system is booting).

In the past, when the boot server could not get to an SP (for example, all paths to that SP had failed before the operating system boot), the user had to manually trespass the LUN to the other SP for the server to boot successfully. With the request forwarding method, users do not need to explicitly trespass the LUN to the other SP. If the HBA boot BIOS is configured so that it can issue an I/O to the surviving SP, I/O will route through the upper redirector to the owning SP and boot the operating system successfully.

For more information, consult the host connectivity guides on E-Lab Navigator on Powerlink that explain how to configure the HBA boot BIOS.

- **Reduces data unavailable situations due to misconfiguration of the host**
Occasionally a user may misconfigure a host in such a way that an application sends an I/O to a non-optimal path (meaning the SP that does not own the LUN). In this case, depending on the failover software installed on the host, the CLARiiON storage system does not return an I/O error

condition. Instead, due to the request forwarding feature, the I/O is routed to the SP that owns the LUN.

Furthermore, when the CLARiiON addresses a change in the access of storage it automatically adjusts the optimal path setting for a LUN. (This “implicit trespass” is discussed in the “Explicit vs. implicit trespass” section.) This automatic adjustment is extremely beneficial in larger environments where the chances of misconfiguration are higher, and in environments where access of a LUN may vary over time.

- **Avoids LUN ownership thrashing situations that may occur in SnapView and cluster configurations**

For Active/Active cluster configurations, if the LUNs are shared and written to by multiple hosts, ALUA avoids the trespass of LUNs between the two SPs in a one-path-per-SP host configuration.

For SnapView, since the clone or snapshot must be owned by the same SP that owns the source LUN, in a case where a production server is writing to the source LUN while the snapshot/clone is mounted and written to by a backup server, a one-path-per-SP configuration on both the production and backup hosts (where each server has a path to the other SP than its peer) can cause path thrashing with servers in non-ALUA mode.

For both cluster and SnapView configurations, with the introduction of the ALUA standard, the LUN will not trespass back and forth between the two SPs, but will be owned by the SP through which maximum I/O requests for that LUN are received by a given host in ALUA mode.

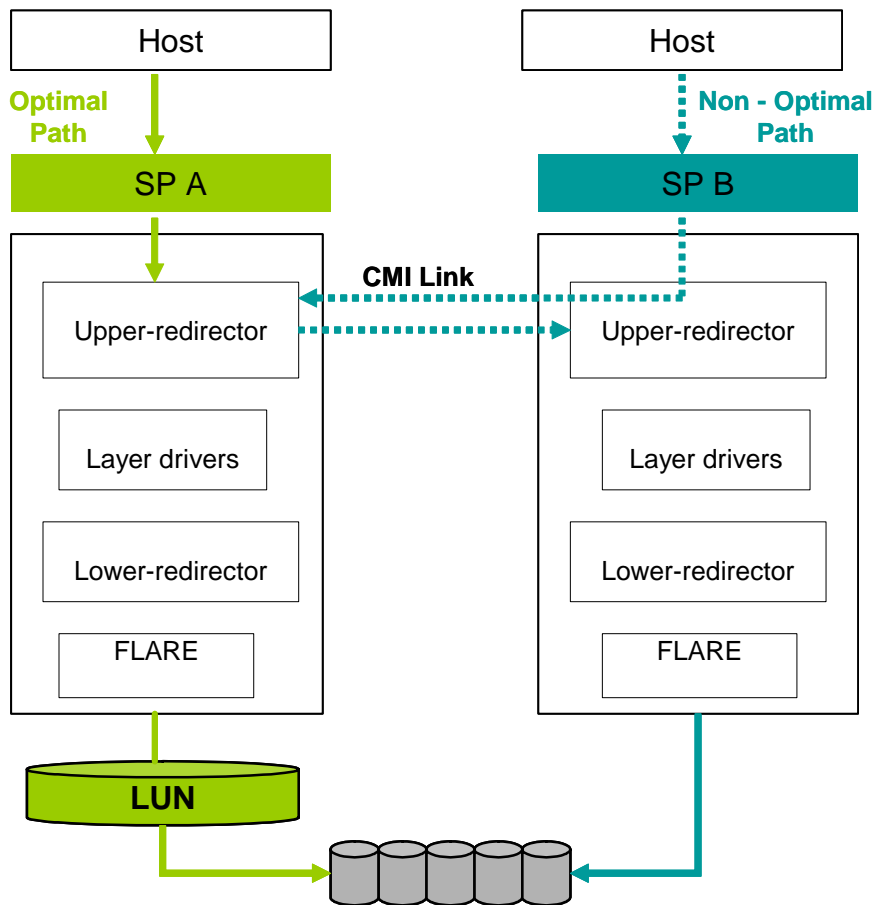


Figure 4. Impact of ALUA with SnapView and cluster configurations

Supports standard SCSI multipath interfaces

ALUA standardizes the implementation of OS vendors' native multipathing and other failover software. Host failover software does not need to contend with CLARiiON-specific trespass commands since CLARiiON implements the ALUA SCSI standard.

Impact of the Redirector

In addition to front-end pathing advantages, the new FLARE release 26 Redirector provides benefits for all attach (traditional PNR as well as ALUA-based) types by providing redirection services for the "back end" as well as the front.

Back-end failure masking

In the case of a back-end failure (for example, an LCC failure) on the SP that is the current owner of the LUN, using ALUA's request forwarding feature, I/O is routed via the lower redirector to the peer SP with the stable back end. The I/O acknowledgement is sent through the SP that owns the LUN, using the lower redirector, as shown in Figure 5. No intervention of failover software is needed on the host, thus masking certain CLARiiON back-end failures.

Note that in the following example, the LUN will be trespassed by FLARE to the SP that can access the LUN at the back end. The host continues to send I/O to that LUN through the peer SP. Once the back-end error is corrected, the LUN is trespassed back to the previous SP, and any I/O for that LUN will not be redirected to the peer SP. As a result, the host sees a minimum delay in I/O during the trespass operations.

This benefits CLARiiON layered applications such as SnapView, MirrorView, and SAN Copy since it isolates those layers from having any knowledge that redirection is taking place. The back-end fault masking feature is provided to all Failover Mode types.

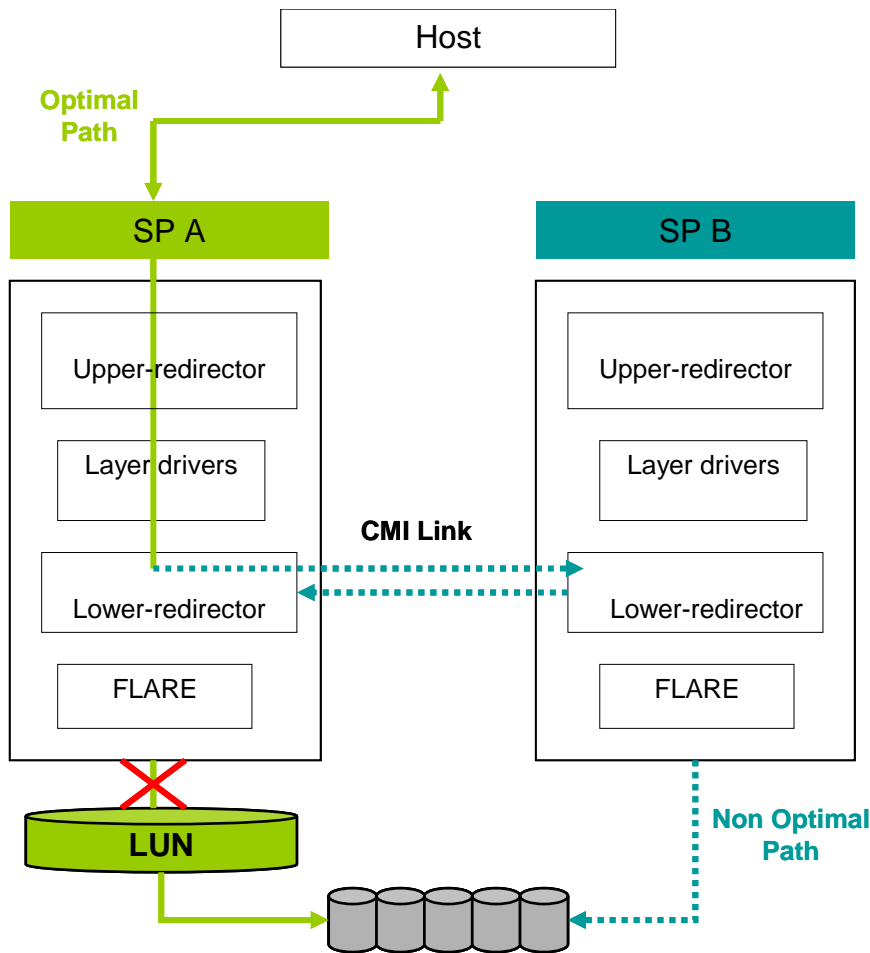


Figure 5. Masking back-end failures without failover software

Comparison of failover methodologies

Path-management software is generally a server-based application that interacts with the storage system to automate and manage multiple paths to data devices. Several options of path-management software are supported by the CLARiiON storage system. As is explained in the “CLARiiON Failover Modes” section, CLARiiON can be configured to run in several different Failover Modes with different path-management software.

EMC PowerPath is host-based software that provides path management. PowerPath works with several storage systems, on several operating systems, with both Fibre Channel and iSCSI data channels. Implementation details and various features of PowerPath and respective native failover software (MPIO) are discussed next.

Explicit and implicit failover software

Explicit and implicit failover *software* should not be confused with the explicit and implicit trespass *commands* mentioned earlier.

Failover software that supports *explicit* ALUA commands monitors and adjusts the active paths between optimal and non-optimal using the REPORT TARGET PORT GROUPS and SET TARGET PORT GROUPS commands, respectively. The SET TARGET PORT GROUPS command, when issued on one controller, trespasses the LUN to the other controller. PowerPath has the same net effect as an explicit

trespass, although it uses traditional CLARiiON proprietary commands to initiate a trespass. Longhorn MPIO (Windows Server 2008) will support explicit ALUA commands.

It is possible for host multipathing software to support *implicit* ALUA functionality. The host could redirect traffic to the non-optimal path for some reason that is not apparent to the storage system. In this case, the storage system must implicitly trespass once it detects enough I/O on the non-optimal paths. To get the optimal path information, the failover software uses the REPORT TARGET PORT GROUP command. This failover software does not provide autorestore capability. Without the autorestore capability, after an NDU, all LUNs could end up on the same SP. The initially supported HPUX 11iv3 worked in this way, although a patch has made this obsolete, and it is not expected to be encountered in the field. Cluster software could also make use of this functionality.

PowerPath

PowerPath works with the storage system to intelligently manage I/O paths. Under normal conditions, PowerPath only issues I/O to LUNs on optimal paths (paths to the owning SP). If optimal paths have failed, PowerPath issues a trespass command; changes the LUN ownership to the other SP; and redefines the non-optimal path as the optimal path.

A unique feature of PowerPath is that it uses a vendor-unique query to the peer SP to determine if the failed path is due to a failed SP. In that case, PowerPath issues a trespass immediately without checking path by path. This is an important advantage of PowerPath over MPIO.

PowerPath is explicit failover software since it load balances only across the optimal paths. If all optimal paths fail, PowerPath will issue a trespass and make the non-optimal paths optimal, and load balance across them. PowerPath also supports the Asymmetric Active/Active feature for operating systems that do not have native failover software that supports ALUA. For example, a Windows 2003 host can be configured in ALUA mode with the 5.1 version of PowerPath that supports ALUA.

PowerPath supports the following features:

- PowerPath will load-balance all I/O only across the optimal paths, as opposed to load balancing I/O across optimal and non-optimal paths, which can lead to lower performance. PowerPath supports various policies for load balancing.
- PowerPath has an autorestore capability to restore LUNs to their default SPs after an SP returns to health following a failure/NDU/cable or switch failure. This ensures even workload balancing of LUNs across SPs and more predictable performance.
- PowerPath supports device prioritization and proactive path testing.
- PowerPath version 5.1 and later support both ALUA and non-ALUA mode LUNs on the same host attached to multiple arrays.
- PowerPath has ALUA-specific I/O retry optimization to improve failover time in the event of SP failures and non-disruptive upgrades (NDU).
- For each LUN, the modes (ALUA versus non-ALUA) are displayed in the PowerPath display along with CLARiiON nice-names
- PowerPath supports Boot on SAN.

A PowerPath license (full functionality) is embedded in the PowerPath package for AX series. For CX300 and CX3-10 models, the PowerPath license (full functionality) is bundled with the storage system.

PowerPath SE is bundled free with CX3-20 and higher models of CLARiiON. PowerPath SE provides single HBA support and limited load-balancing options. To get full functionality PowerPath support with CX3-20 and higher models, full PowerPath software can be purchased. Full PowerPath has more flexible options for attach (multiple HBA) and load-balancing policies, together with unique features as mentioned previously.

Native failover software (MPIO and others)

To many people, native failover software means “MPIO.” MPIO is a commonly used name for the host-side interface that handles multipathed LUNs. However, MPIO implementations on different vendor operating systems are not the same API. Some operating systems require storage vendors to provide libraries to implement their MPIO framework, while some operating systems include libraries (such as Solaris MPIO).

MPIO framework can be multi-mode. Some vendors support Active/Active, Active/Passive, and ALUA. HP-UX 11iv3 and MS Windows Server 2008 (Longhorn) limit support to ALUA or Active/Active.

CLARiiON supports various native MPIO software, as noted in Table 1.

Native failover systems (ALUA and non-ALUA) often do not offer all the features that PowerPath offers. These should be noted before a commitment to native failover software is made. Some of the features not offered in failover systems include:

- Windows MPIO (iSCSI) has no autorestore.
- HP-UX PVLlinks “classic” (before 11iv2) has no load balancing.
- HP-UX 11iv3.1 MPIO will offer load balancing but no autorestore.
- LINUX MPIO “classic” RH4 and SUSE 9 don’t support boot from SAN.
- LINUX ALUA-based MPIO has no autorestore and cannot detect a hung SP. There are no timers in Linux. PowerPath can detect a hung SP as PowerPath times and checks the peer SP.
- AIX MPIO and Solaris MPxIO don’t support autorestore. EMC PowerPath provides autorestore capability. This ensures even workload balancing of LUNs across SPs and more predictable performance.

Nondisruptive upgrade (NDU)

NDU requires a reboot of each SP in turn. For the most part, ALUA does not change how the storage system works in an NDU. Using PowerPath with Active/Passive mode (Failover Mode 1), PowerPath trespasses the LUN to the other SP while the owner SP is rebooting. After the NDU operation completes, the PowerPath autorestores the LUN to the SP that is the default owner of the LUN.

In Asymmetric Active/Active mode, during the SP reboot the failover software detects the failed paths and redirects I/O to the SP that is up and running. When the NDU finishes, CLARiiON relies on the failover software to interrogate the preferred bit returned by the REPORT PORT GROUPS command to move LUNs back to the default owner – thus restoring the original path structure.

PowerPath trespasses a LUN to the peer SP during an NDU operation. After an NDU completes, PowerPath autorestores the LUN to its default owner. This is true whether the initiators are configured with Failover Mode 1 (Active/Passive) or 4 (ALUA). Note that PowerPath does not use the preferred bit to perform the autorestore but uses CLARiiON proprietary commands to issue a trespass.

Manual trespass

In the case of Active/Passive mode (Failover Mode 1), when a manual trespass is issued (using Navisphere Manager or CLI), subsequent I/O for that LUN is rejected over the SP on which the manual trespass was issued. This would result in a unit attention condition for the host; failover software would detect that error and re-route the I/O to the SP that owns the LUN.

Any trespass operation, automatic or manual, causes the ownership of the LUN to change and a **unit attention** to be sent to connected hosts. These hosts, being ALUA aware, would interpret the unit attention and query the current status using appropriate Target Port Group commands.

PowerPath and native MPIO that is ALUA aware would immediately act on the unit attention by routing all I/O to the optimal path so that I/O does not go through the non-optimal path.

So, if you manually trespass a LUN, both PowerPath and MPIO software will continue to use optimal paths. As per the CLARiiON ALUA implicit trespass mechanism, changes in I/O balance may cause the LUN to trespass implicitly (as discussed in “Explicit vs. implicit trespass” section).

Path, HBA, and switch failure

If a host is configured with Failover Mode 1 and all the paths to the SP that owns a LUN fail, the LUN is trespassed to the other SP by the host’s failover software.

With Failover Mode 4, in the case of a path, HBA, or switch failure, when I/O routes to the non-owning SP, the LUN may not trespass immediately (depending on the failover software on the host). If the LUN is not trespassed, FLARE will trespass the LUN to the SP that receives the most I/O requests to that LUN based on the implicit trespass mechanism.

SP failure

In case of an SP failure for a host configured as Failover Mode 1, the failover software trespasses the LUN to the surviving SP.

With Failover Mode 4, if an I/O arrives from an ALUA initiator on the SP that does not own the LUN (non-optimal), failover software or FLARE initiates an internal trespass operation. This operation changes ownership of the target LUN to the surviving SP, since its peer SP is dead. Hence, the host (failover software) must have access to the secondary SP so that it can issue an I/O under these circumstances.

Conclusion

The Asymmetric Active/Active feature within FLARE release 26 increases configuration flexibility and investment protection for a changing environment, and provides increased availability for systems.

Asymmetric Active/Active is an emerging SCSI standard. It supports host failover methods that adhere to this specification. Hosts can now avoid unwanted trespassing in certain scenarios, as I/O can be redirected to the SP that owns the LUN

The back-end fault masking feature is provided to all Failover Mode types. In the case of a back-end failure on the SP that is the current owner of the LUN, I/O is routed via the lower redirector of the peer SP with the stable back end. This avoids trespassing, which can often impede replication operations.

All ports can be used to access the same LU simultaneously. The Asymmetric Active/Active multipathing feature is a software enhancement to the current base software package. It will be upward compatible with existing Failover Modes. It presents an Asymmetric Active/Active model, allowing host I/O to a LUN over all ports based on the optimal and non-optimal path as reported by REPORT TARGET PORT GROUP.

Asymmetric Active/Active is part of the SPC-3 SCSI standard and is a new selectable Failover Mode just like the current modes of 0 (Auto-trespass), 1 (Passive Not Ready), 2 (DMP), and 3 (Passive Always Ready). Note that the Asymmetric Active/Active model can be applicable for any given attach to the array but is not the default system behavior. A PNR initiator and an Asymmetric Active/Active initiator can connect through the same physical SP port and have different failover behavior where PNR behaves the same as in legacy systems.

References

- *EMC CLARiiON High Availability (HA) – Best Practice Planning*