



# INTEGRATING & TROUBLESHOOTING VMWARE WITH EMC'S MIDRANGE CLARION: DEEP DIVE

**EMC<sup>2</sup>** EMC Proven Professional Knowledge Sharing 2012



Jason L. Gates

Systems Engineer - Storage & Virtualization

Presidio Networked Solutions

[jgates@presidio.com](mailto:jgates@presidio.com)

<http://www.linkedin.com/in/mrjasongates>

**EMC<sup>2</sup>**

## Contents

Audience .....	3
Overview .....	3
Troubleshooting Connectivity VMware ESX Version 4, 3.x with EMC CLARiiON .....	4
Configuring Multipathing and Failover for High Availability.....	6
ISCSI Troubleshooting .....	8
SCSI Reservations .....	9
LUN Layout Considerations.....	10
Conclusion.....	11
References .....	12

Disclaimer: The views, processes or methodologies published in this article are those of the author. They do not necessarily reflect EMC Corporation's views, processes or methodologies.

## **Audience**

This article is intended for system administrators, VMware engineers, and storage administrators. Readers are expected to be familiar with VMware ESX/ESXi hosts, basic CLARiiON® system operations, fabric switches, basic networking, VCenter Server, and EMC Navisphere/Unisphere Manager.

## **Overview**

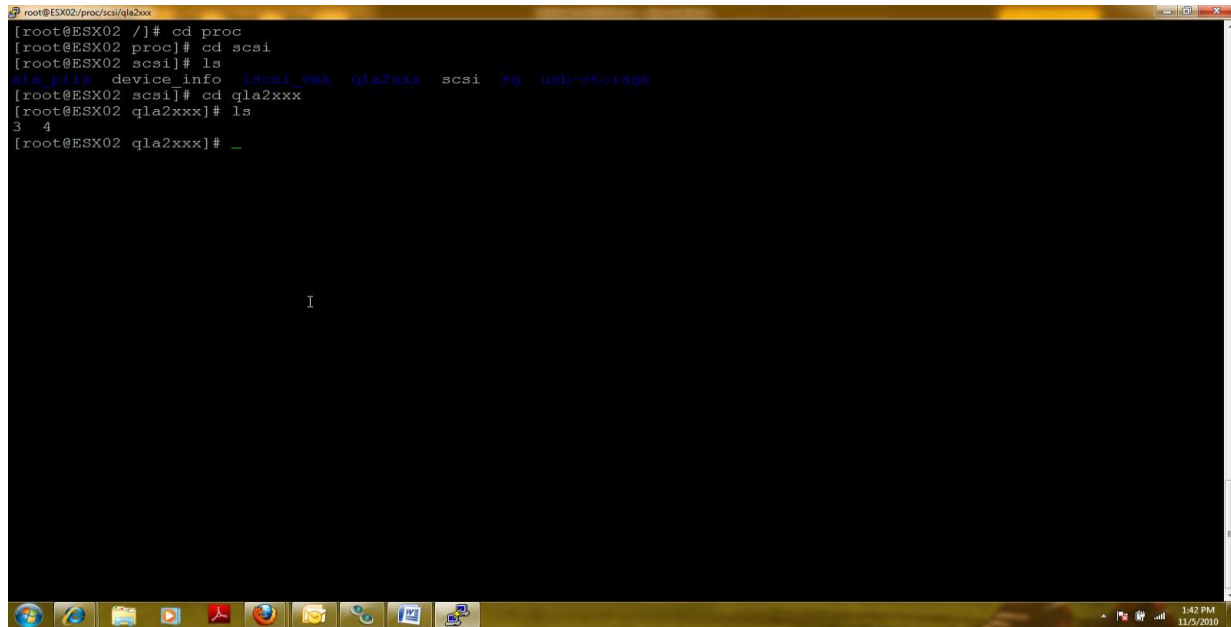
Many in the EMC Proven™ community have questions about best practices and the advanced skills that it takes to troubleshoot problems at the Systems Engineer level when it involves VMware and CLARiiON all in one solution/configuration. Problems do arise and resolving these problems takes an engineer who has a good understanding of the VMware hosts and expected behavior of the CLARiiON at the SCSI level.

What kind of problems arise? LUN trespassing, performance issues, connectivity, iSCSI and SCSI reservation issues. What happens when these type of problems occur? The storage administrator is pointing fingers at the VMware server and the VMware engineer is pointing fingers at the storage subsystem.

Common questions are, “What is the best path policy?” “What invokes a trespass?” “How does native multipath behave?” “What causes SCSI Reservations and how to resolve?” “What are best practices to configure iSCSI with the CLARiiON?” “How should my LUNs be carved out based on VM (Virtual machine) needs?” This article covers troubleshooting skills, configuration methods, and best practices anyone can utilize on the world’s number one mid-range storage sub-system and VMware’s leading virtualization software.

## Troubleshooting Connectivity VMware ESX Version 4, 3.x with EMC CLARiiON

Most seasoned VMware administrators and storage administrators understand that ESX hosts access CLARiiON LUNs via the Fibre Channel protocol or iSCSI at the block level. Basics such as zoning, LUN masking, creating storage groups, registering hosts, and assigning host/servers to storage groups are steps that need to occur for storage to be available, which the majority of EMC professionals are comfortable with. What happens when all the above is done and there is still no storage available to our ESX hosts? Reboot? Rebooting is a common workaround for problems with storage but sometimes this is not an acceptable workaround, so we must dig deeper to really understand what is happening. I'm a strong supporter of using ESX command line to troubleshoot issues, so anytime there is a fibre channel (FC) connectivity problem, I use the command line to delve deeper.



```
root@ESX02/proc/qla2xxx
[root@ESX02 /]# cd proc
[root@ESX02 proc]# cd scsi
[root@ESX02 scsi]# ls
qla qlx device_info local_vmk qla2xxx scsi sg usb-storage
[root@ESX02 scsi]# cd qla2xxx
[root@ESX02 qla2xxx]# ls
3 4
[root@ESX02 qla2xxx]# _
```

Figure 1 Using SSH as in this display, `proc/scsi/qla2xxx` is the directory that stores very useful information about the HBA. (In your environment, the HBA nomenclature could be different)

```
root@ESX02:/proc/scsi/qla2xxx
[root@ESX02 /]# cd proc
[root@ESX02 proc]# cd scsi
[root@ESX02 scsi]# ls
isa_piix device info laci1_vmk qla2xxx scsi sg usb-storage
[root@ESX02 scsi]# cd qla2xxx
[root@ESX02 qla2xxx]# more 4
QLogic PCI to Fibre Channel Host Adapter for QLE8042:
    Firmware version 4.04.09 [IP] [Multi-ID] [84XX] , Driver version 8.02.01-k1-vmw38
BIOS version 2.04
FCODE version 2.02
EFI version 2.00
Flash FW version 4.04.00
ISP: ISP8432
Request Queue = 0x1b0cb000, Response Queue = 0x1b14c000
Request Queue count = 4096, Response Queue count = 512
Total number of interrupts = 3929
    Device queue depth = 0x20
Number of free request entries = 4096
Number of mailbox timeouts = 0
Number of ISP aborts = 0
Number of loop resyncs = 1
Host adapter:loop state = <DEAD>, flags = 0x105a83
Dpc flags = 0x40180c0
MBX flags = 0x0
Link down Timeout = 030
Port down retry = 005
Login retry count = 008
Execution throttle = 1024
ZIO mode = 0x6, ZIO timer = 1
Commands retried with dropped frame(s) = 0
Product ID = 0000 0000 0000 0000

NPIV Supported : Yes
Max Virtual Ports = 63

SCSI Device Information:
scsi-qlal-adapter-node=2001001b32a00e1a:000000:0;
scsi-qlal-adapter-port=2101001b32a00e1a:000000:0;

FC Target-Port List:

FC Port Information:
[root@ESX02 qla2xxx]#
```

Figure 2 Using the command *more 4*, we see the loop state of this HBA is “Dead” This indicates a bad HBA, HBA cable, or perhaps a GBIC on the switch. Sometimes it is difficult to trace cables and wires so this is a quick way to confirm connectivity.

```
root@ESX02:/proc/scsi/qla2xxx
Flash FW version 4.04.00
ISP: ISP8432
Request Queue = 0x1b013000, Response Queue = 0x1b094000
Request Queue count = 4096, Response Queue count = 512
Total number of interrupts = 1477435
    Device queue depth = 0x20
Number of free request entries = 3953
Number of mailbox timeouts = 0
Number of ISP aborts = 0
Number of loop resyncs = 6
Host adapter:loop state = <READY>, flags = 0x145ac3
Dpc flags = 0x10000
MBX flags = 0x0
Link down Timeout = 030
Port down retry = 005
Login retry count = 008
Execution throttle = 1024
ZIO mode = 0x6, ZIO timer = 1
Commands retried with dropped frame(s) = 0
Product ID = 0000 0000 0000 0000

NPIV Supported : Yes
Max Virtual Ports = 63

SCSI Device Information:
scsi-qla0-adapter-node=2000001b32800e1a:680001:0;
scsi-qla0-adapter-port=2100001b32800e1a:680001:0;

FC Target-Port List:
scsi-qla0-target-0=5006016930213fad;
scsi-qla0-target-1=5006016a3ce00ebd;
scsi-qla0-target-2=50060e80143a0811;

FC Port Information:
--More-- (0%)
```

Figure 3 I replaced the HBA cable. Now, the Target-Port List has the WWPN of the CLARiON listed, which indicates connectivity is OK. The state of the HBA is changed to “Ready”.

```
root@ESX02:/proc/scsi# cd /
root@ESX02 /# cd /proc
root@ESX02 /proc# cd /proc/scsi
root@ESX02 /proc/scsi# cat /proc/scsi/scsi
Attached devices:
Host: scsi5 Channel: 00 Id: 00 Lun: 00
  Vendor: TEAC      Model: DV-28E-V      Rev: C.AB
  Type:   CD-ROM
  ANSI SCSI revision: 05
Host: scsi0 Channel: 00 Id: 00 Lun: 00
  Vendor: VMware   Model: Virtual disk  Rev: 1.0
  Type:   Direct-Access
  ANSI SCSI revision: 02
Host: scsi5 Channel: 00 Id: 01 Lun: 00
  Vendor: EMC      Model: Celerra      Rev: 0002
  Type:   Direct-Access
  ANSI SCSI revision: 05
Host: scsi5 Channel: 00 Id: 02 Lun: 00
  Vendor: NETAPP   Model: LUN           Rev: 7320
  Type:   Direct-Access
  ANSI SCSI revision: 04
Host: scsi5 Channel: 00 Id: 03 Lun: 00
  Vendor: DGC      Model: RAID 5       Rev: 0226
  Type:   Direct-Access
  ANSI SCSI revision: 04
Host: scsi5 Channel: 00 Id: 04 Lun: 00
  Vendor: DGC      Model: RAID 5       Rev: 0226
  Type:   Direct-Access
  ANSI SCSI revision: 04
Host: scsi5 Channel: 00 Id: 05 Lun: 00
  Vendor: DGC      Model: LUNZ         Rev: 0429
  Type:   Direct-Access
  ANSI SCSI revision: 04
Host: scsi5 Channel: 00 Id: 06 Lun: 00
  Vendor: Generic  Model: STORAGE DEVICE Rev: 9412
  Type:   Direct-Access
  ANSI SCSI revision: 02
root@ESX02 /proc/scsi#
```

Figure 4 Connectivity is now intact. We can also confirm that SCSI devices are attached from `proc/scsi` and `cat` the file `scsi`. The highlighted entry is a CLARiiON LUN, because the vendor type is “DGC”, (Data General Corporation.)

### Configuring Multipathing and Failover for High Availability

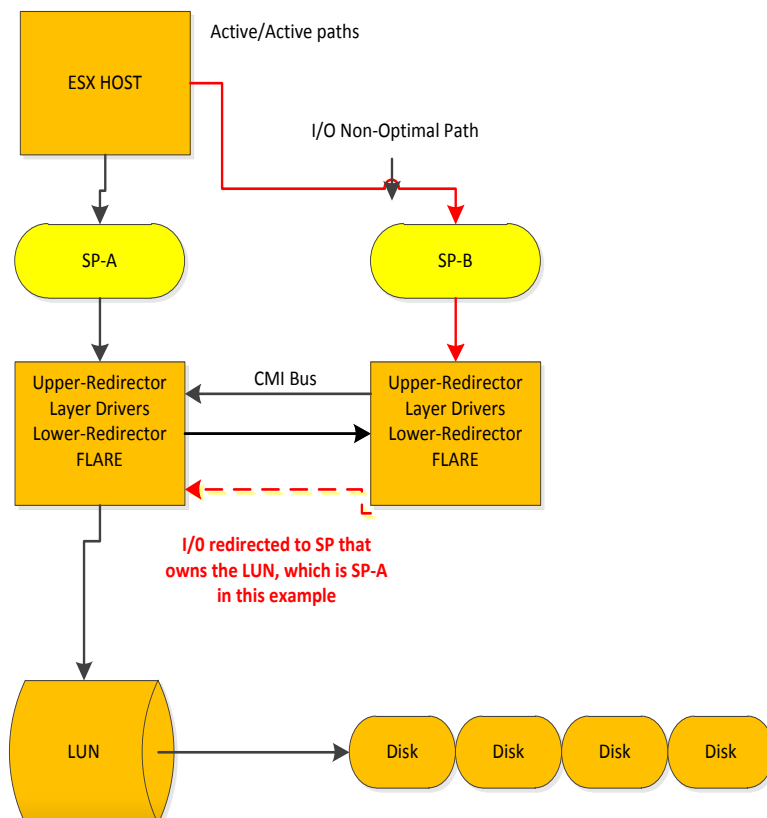
Multipathing and load balancing increase the level of availability for applications running on ESX servers. CLARiiON supports nondisruptive upgrade (NDU) operations for VMware’s native failover software and EMC’s popular failover software, PowerPath. The CLARiiON is an active/passive array and the policy for all paths should be set as MRU if using ESX 3 to avoid path thrashing. Path thrashing occurs when a server cannot access a LUN or access is very slow. In most cases two or more servers are attempting to access the LUN through different storage processors and as a result, the LUN is **never** truly available. Let’s discuss this in detail.

The storage processors are independent computers that access shared storage and algorithms are used to determine how concurrent access is handled. When using the MRU policy, there is no concept of preferred path; in this case, the preferred path can be disregarded. The MRU policy uses the most recent path to the disk until this path becomes unavailable. If two ESX servers are connected, with path one from HBA1 to SPA, and path two from HBA0 to SPB, a single LUN configured as a VMFS volume can be accessed by multiple ESX servers. Using MRU, each HBA should be cross-zoned to each SP on the CLARiiON array. The ESX host will only use one active path for I/O and thus makes use of only one HBA. In this configuration, the

host will continue to use the active path until it finds a failure and then will failover to the alternate path (which could be via the same or different SP). Keep in mind, native failover is not dynamic I/O balancing software; a manual restore is needed to take the original path again or it will continue to use the alternate path even after the original path recovers.

With VSphere 4.0 and above, failover and multipathing can be configured for Asymmetric Logical Unit Access (ALUA) on the CLARiiON array with FLARE 28 and higher. I highly recommend this setting. This configuration is quite different from the active/passive configuration. To understand multipathing behavior of VMware at the host level, we must understand expected behavior on the CLARiiON backend with active/active configuration. Asymmetric Active/Active introduces initiator Failover Mode (Failover mode 4,) where initiators are permitted to send I/O to a LUN regardless of which SP actually owns the LUN.

ALUA uses SCSI 3 commands that are part of the standard SCSI SPC-3 specification to determine I/O paths. Case in point, 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. The redirection is done through internal communication (upper-redirect) between the SP's; a LUN trespass is not needed and the redirection is transparent to the hosts. The diagram below illustrates ALUA behavior:



The following event code will be logged in the SP(s) event log:

“UpperRedirector 711b0000 IO is being forwarded to the other SP for processing”

Use of ALUA failover mode has additional benefits when combined with VSphere native multipathing software, providing automatic restore of LUNs to its preferred paths after a fault has been repaired. However, this only applies to "Fixed" policy configuration. This does not apply to "Round Robin" policy. So what is the major difference between Fixed policy and Round Robin? The Fixed policy supports auto-restore of failed paths but NO load balancing. The Round Robin policy performs load balancing but NO auto-restore of failed paths. The ALUA feature also eliminates path thrashing in Active/Active cluster configurations, and avoids unavailability when booting from a SAN LUN during a path failure.

## ISCSI Troubleshooting

Internet Small Computer System Interface (iSCSI) is a cost-effective, easy-to-manage storage solution that many companies deploy into production. The iSCSI protocol allows clients (called *initiators*) to send SCSI commands to SCSI storage devices (*targets*). Configuring the software iSCSI initiator is straight-forward with both versions of ESX using VMware's iSCSI configuration guide. Most iSCSI issues are network related—poor connections, switch issues, hops, routing issues, and so forth. These kinds of problems can severely impact the read & write performance of storage attached to the ESX software through the iSCSI initiator.

If your network is congested or has very high periods of intensive I/O activities that negatively affect your ESX environment, a workaround could be to disable the delayed acknowledgement (ACK) feature on your VMWare server via a configuration option. ESX uses delayed ACK to increase efficiency in both the network and the hosts by sending less than one ACK acknowledgment segment per data segment received.

Congestion can manifest itself as a delay, timeout, or packet loss and during this “congestion” period the CLARiiON will only re-transmit **one** lost data segment at a time. The ACK feature coupled with the CLARiiON only re-transmitting one data segment would slow read performance to a halt in a congested network and frequent timeouts will be reported in the ESX kernel logs, virtual machine guest would timeout, and so on. To modify the delayed ack settings on a specific discovery target, do the following:

- 1) Select the **Static Discovery** tab
- 2) Select the **Server Address** tab
- 3) Click **Settings**
- 4) Click **Advanced**
- 5) Uncheck **DelayedAck**

*Note: As a rule of thumb, this setting should be tested to confirm this is ideal for your applications and network before deploying into production.*

## **SCSI Reservations**

All VMware clusters use the SCSI protocol to manage disks on a shared bus. When a reserve command is issued by a host bus adapter (HBA), the command allows a HBA to obtain or maintain ownership of a SCSI device. A device that is reserved refuses all commands from all other host bus adapters except the one that initially reserved it, the initiator. In the VMware world, VMFS data store-level operations require a lock or reservation to ensure there is no data corruption, because two hosts cannot write to the same device at the same time; hence, the need for reservations.

So what happens when host A has a reservation on LUN 1 and host B attempts to access LUN 1? The answer is a “reservation conflict” is reported by Host B. A reservation conflict means that host B is trying to access the LUN and is not being allowed to, because host A has already reserved the LUN for its own access. Some administrators believe that the CLARiiON creates the SCSI reservations. This is incorrect; reservations are not created by the array itself. The CLARiiON simply responds to the commands sent by the ESX hosts. When the CLARiiON fails to release the reservation, the request may not have come through (hardware, firmware, pathing problems, and so forth). For example, virtual machines may experience I/O failures and blue screens due to too many SCSI reservation conflicts. To resolve reservations conflicts, the following should be done or considered:

- 1) Confirm all HBA drivers are up-to date.
- 2) Determine which HBA has the lock; remove the HBA cable from fabric switch, wait 10 seconds, then reconnect. This forces a Port Login and Fabric Login which would clear the reserve.
- 3) Do not place Exchange, SQL, or Oracle VM's on the same LUN(s), because these types of applications change data frequently.

- 4) Reduce heavy amount of VM snapshots.
- 5) Confirm failover settings are correct on all hosts accessing shared storage.
- 6) As a last resort, the Storage Processor that owns the LUN may need to be rebooted or reboot the ESX host that placed the reservation on the LUN.

Reservation conflicts should be expected in an ESX clustered environment, but a high number of conflicts could cause unexpected outages. I strongly suggest keeping all systems—including the CLARiiON—on latest codes and patches.

## **LUN Layout Considerations**

When designing LUN layouts, each situation is different; however there are some basic guidelines to follow. Select RAID 5 protection on Fibre Channel drives for boot volumes and servers such as DNS or large file servers because low I/O activity is expected from these type of virtual machines. RAID 5 writes require parity updates to each disk; however each disk can be read from independently.

SATA-II drives should be used for storing archived data and could be RAID 5 or RAID 1 protected.

For write intensive application data volumes and logs of databases, use RAID 1/0 on fast Fibre Channel drives such as 146GB 15k or flash drives. Also in a RAID 1/0 group, I would advise to span the drives across two separate buses if feasible. The purpose of creating the RAID Group this way is to place data and mirrors on two separate enclosures. In the event of an enclosure failure, the other enclosure would still be online and maintain access to the data or the mirrored data. See the diagram below:

Order of Disks Example for Raid 1/0 six disk RAID Group

First Data1 1\_0\_0

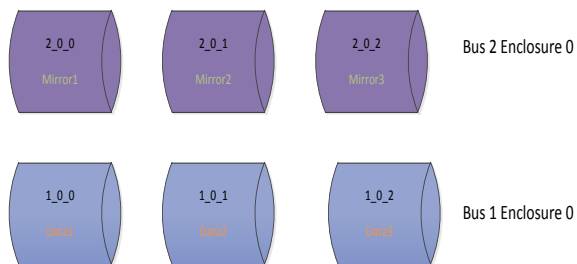
Second Mirror1 2\_0\_0

Third Data2 1\_0\_1

Fourth Mirror2 2\_0\_1

Fifth Data3 1\_0\_2

Sixth Mirror3 2\_0\_2



*Note: VMware Vmfs-3 volumes are already aligned to 64KB during creation. However, if using Windows 2003/2000, the volume will need to be aligned at the virtual machine level using diskpart. This is not required on Windows 2008. The CLARiiON formats disks in blocks of 128 per disk, which is equivalent to a 64 KB of data that is written to a disk from write cache.*

## Conclusion

EMC's mid-range storage leader, CLARiiON, and VMware provide a robust, industry-leading Information Lifecycle solution. These technologies complement each other well when configured properly using best practices. I sincerely hope that this article will be a great asset to EMC Proven Professionals and the community in general when troubleshooting common problems in a VMware/EMC CLARiiON environment.

## References

EMC CLARiiON Integration with VMware ESX Server - White Paper

EMC's Host Connectivity Guide for VMware - Best Practices Guide

Fibre Chanel SAN Configuration Guide -

[www.vmware.com/pdf/vsphere4/r40/vsp\\_40\\_san\\_cfg.pdf](http://www.vmware.com/pdf/vsphere4/r40/vsp_40_san_cfg.pdf)

ISCSI SAN configuration Guide -

[www.vmware.com/pdf/vsphere4/r40/vsp\\_40\\_iscsi\\_san\\_cfg.pdf](http://www.vmware.com/pdf/vsphere4/r40/vsp_40_iscsi_san_cfg.pdf)

Performance for VSphere - [www.vmware.com/pdf/Perf\\_Best\\_Practices\\_vSphere4.0.pdf](http://www.vmware.com/pdf/Perf_Best_Practices_vSphere4.0.pdf)

Next Generation Best Practices for Storage and VMware - <http://virtualgeek.typepad.com>

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.