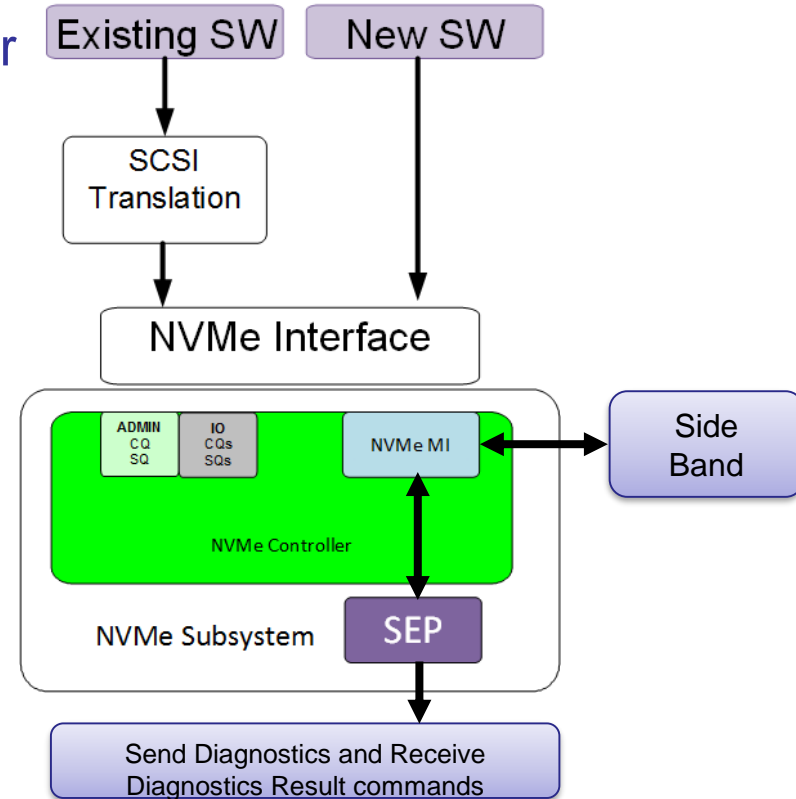


NVMe Enclosure Management and Dual Port Drive

Parag Maharana

SES over NVMe Overview

- Inband SES management via NVMe driver
- Out of band SES management via NVMeMI
- SEP can be exposed as a Namespace or as a separate PCIe function or both
- Existing applications can use SCSI interface to access SES pages based on SCSI to NVMe Translation
- New application can access directly via NVMe admin queue pair

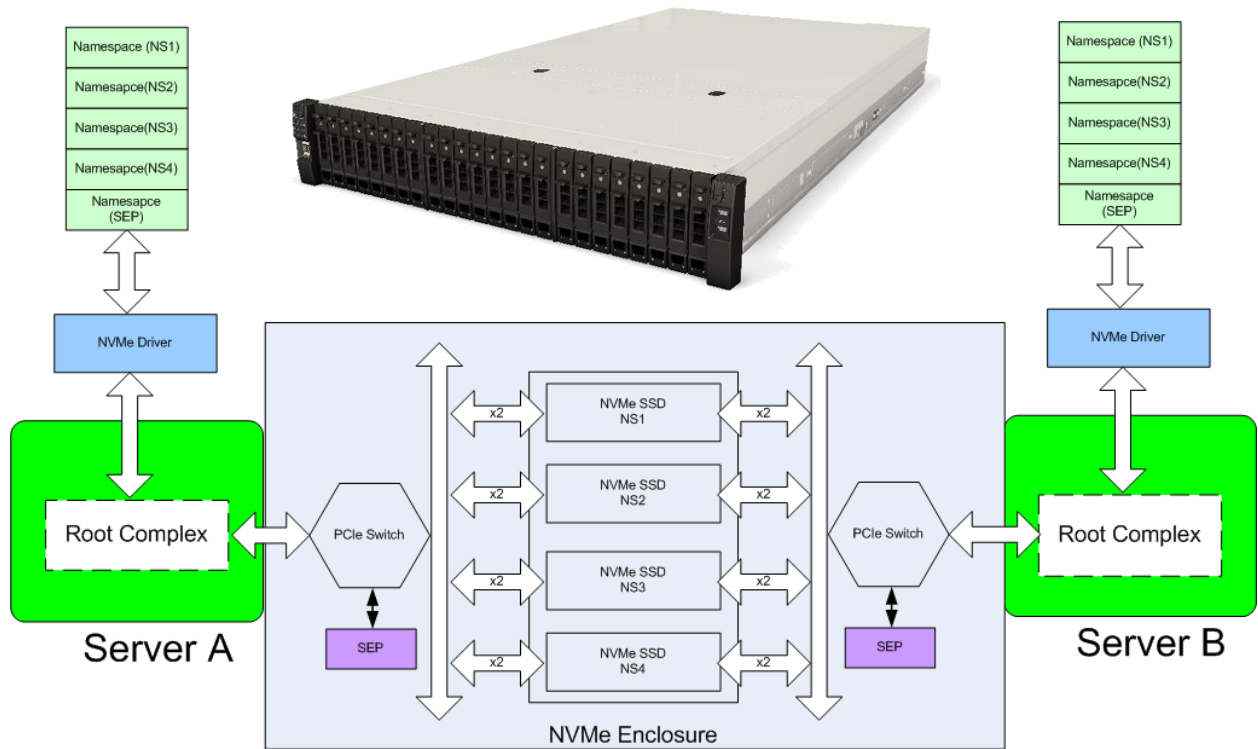


Enclosure Services Management usage model

- Main goal is to use standard SCSI Enclosure Services tools on NVMe. This will leverage investment on existing management tools for SES
- Also generic sg_ses tool in Linux that can be used to manage NVMe enclosure services, can manage NVMe enclosure
- When sg_ses open SEP namespace (e.g. /dev/nvme0ns1) then NVMe driver can translate/convert send and receive diagnostic commands to NVMe admin command that will be send to NVMe HBA or NVMe enclosure
- NVMe Admin commands are issued in-band through PCIe.
- Linux generic SES management tool will work as it is by selecting NVMe controller namespace “sg_ses /dev/nvme0ns1”

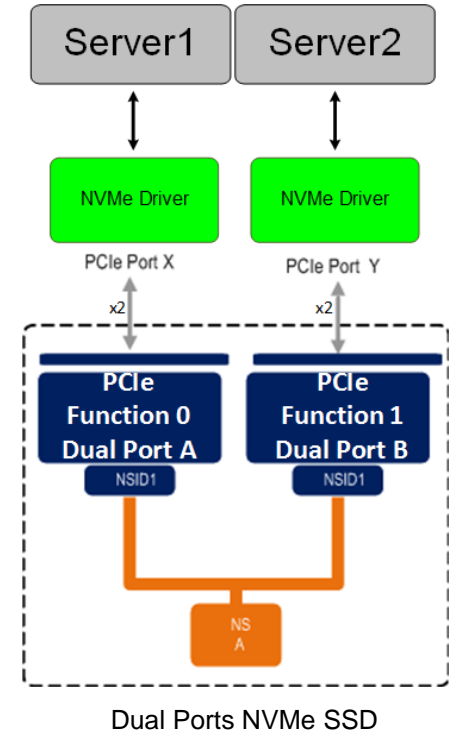
SES over NVMe for Multi-Host Topology

- Multi-Host Topology
- Directly connected to root complex of the servers
- Each Namespace will be expose through a separate PCIe functions
- Each namespace will be handed by separate driver instance



Dual Ports NVMe Drive

- Single port U.2 NVMe drives are PCIe x4
- Dual Ports U.2 NVMe drives are each PCIe x2
- Dual port drives shared namespace across both ports
- Dual ports can have dedicated namespace per port (mainly for boot)
- Dual port drives required reservation to allow active-active support from two or more servers on shared namespace
- NVMe Subsystem reset will reset entire drive



NVMe Persistent Reserve and Release

- Reservation functionalities are compatible with T10 persistent reservations
- Dual port drive can support any one type of Reservation. Typical reservation type is Write Exclusive with Registrant
- Dual port drive need to support 4 PR calls i.e. register, acquire, release and report
- The reservations and registrations can be persisted across all type of resets including NVM Subsystem Resets. This can be enabled via NVMe PTPLS Register
- Reservation can only be cleared by preempt or preempt and abort commands. Also drive can updates PR log pages and AENs
- Dual port drive should handle Host Identifier as well as change of Host Identifier (HOSTID) in case of server replaced or drive swapped

Reservation Type	Reservation Holder		Registrant		Non-Registrant		Reservation Holder Definition
	Read	Write	Read	Write	Read	Write	
Write Exclusive	Y	Y	Y	N	Y	N	One reservation holder
Exclusive Access	Y	Y	N	N	N	N	One reservation holder
Write Exclusive – Registrants only	Y	Y	Y	Y	Y	N	One reservation holder
Exclusive Access – Registrants only	Y	Y	Y	Y	Y	N	One reservation holder
Write Exclusive – All Registrants	Y	Y	Y	Y	Y	N	all Registrants are Reservation Holder
Exclusive Access – All Registrants	Y	Y	Y	Y	N	N	all Registrants are Reservation Holder

NVMe I/O Command	Operation/Action
Reservation Register	<ol style="list-style-type: none"> Register a reservation key Unregister a reservation key Replace a reservation key
Reservation Acquire	<ol style="list-style-type: none"> Acquire a reservation on a namespace Preempt reservation held on a namespace Preempt and abort a reservation held on a namespace
Reservation Release	<ol style="list-style-type: none"> Release a reservation held on a namespace Clear a reservation held on a namespace
Reservation Report	<ol style="list-style-type: none"> Retrieve reservation status data structure <ol style="list-style-type: none"> Type of reservation held on the namespace (if any) Persist through power loss state Reservation status, Host ID, reservation key for each

Conclusion

- NVMe Dual port drives fits well in NVMe enclosures
- NVMe enclosures can direct connect to server using PCIe
- Also NVMe enclosures can work well with fabrics transport
- NVMe enclosure will provide a fault tolerant infrastructure using redundant path using Dual port drives
- Intelligence NVMe enclosures can function as NVMe appliances or AFA
- SES over NVMe will provide a common/standard mechanism to control all non NVMe devices in enclosure (e.g. Fans, Power supplies, LED, etc)

Thank You! Questions?

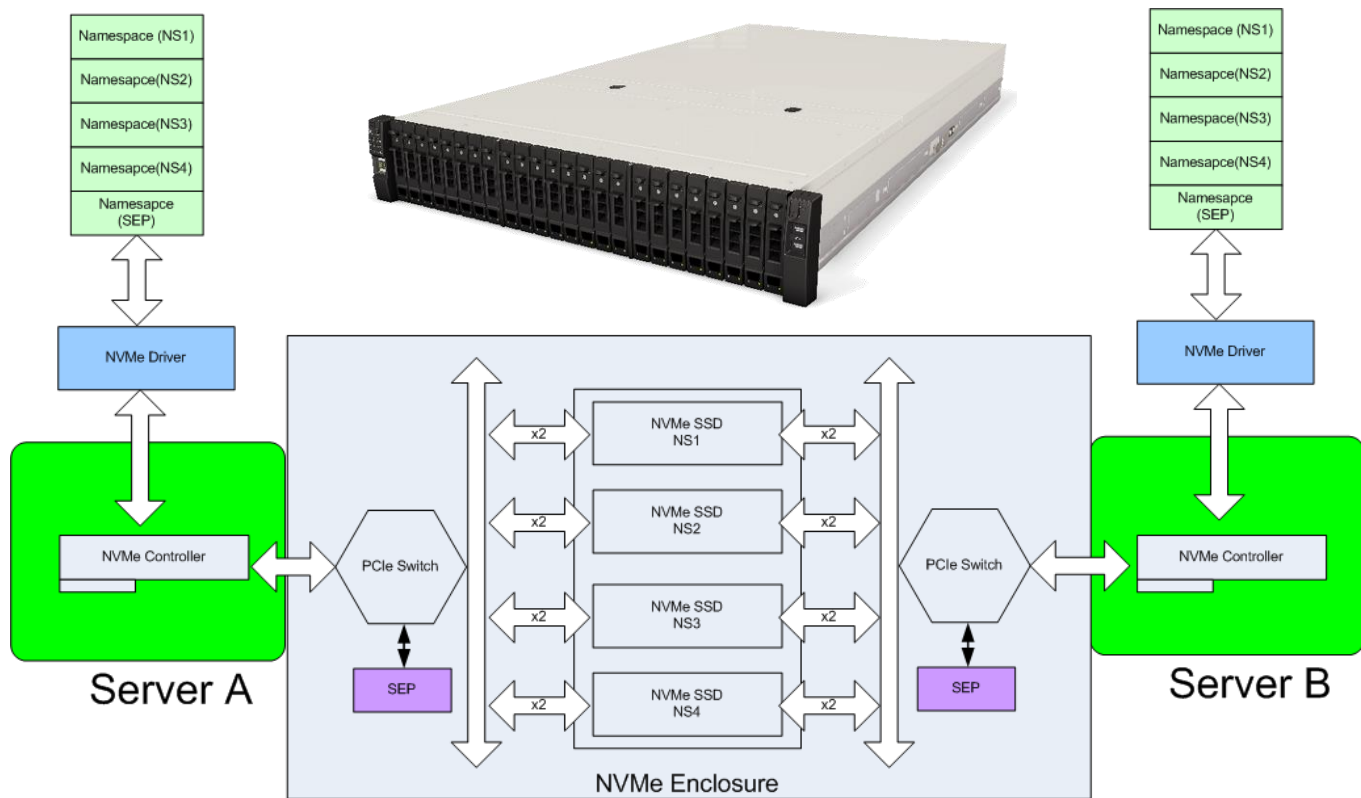
A large, white, stylized letter 'S' logo is positioned on the left side of a green horizontal banner. The 'S' is composed of thick, rounded strokes and is partially cut off by the left edge of the banner.

Visit Seagate Booth #505

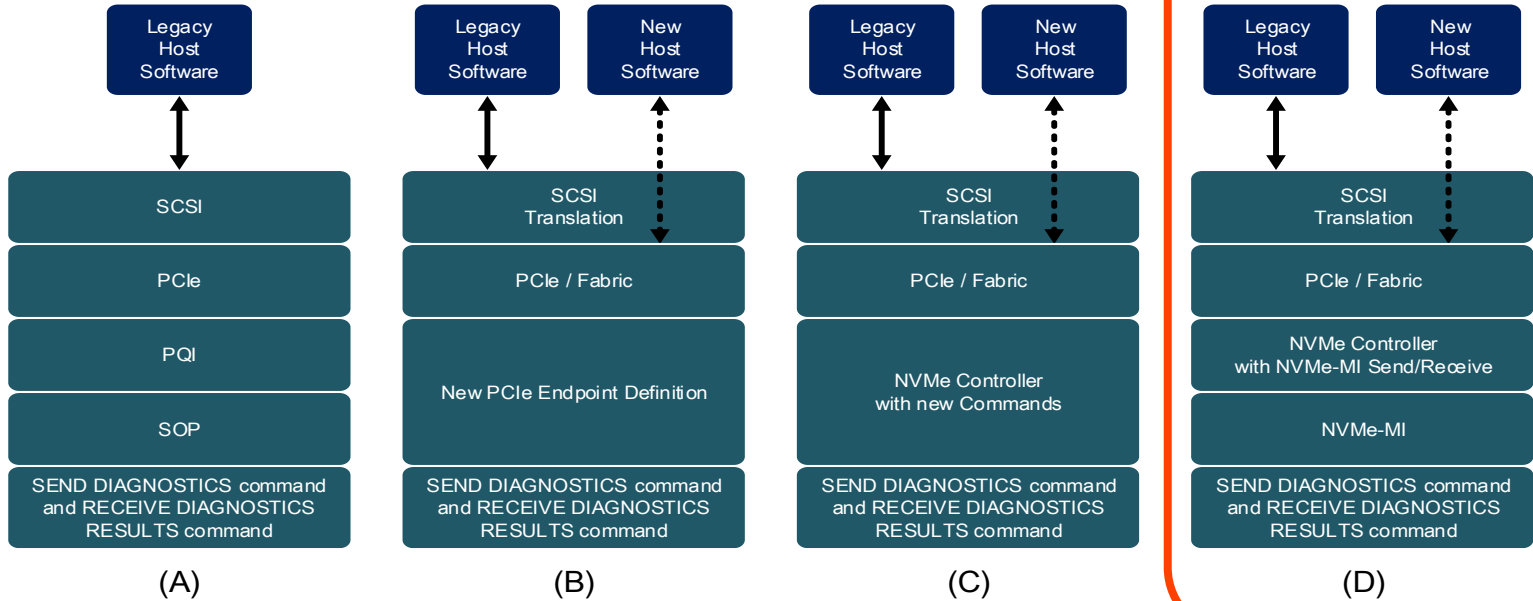
Learn about Seagate's ever-expanding portfolio of SSDs, Flash solutions and system level products for every segment

SES over NVMe for Multi-Host Topology

- Multi-Host Topology with NVMe aggregator Controller at Host



Possible Options (there may be others)



NVMe workgroup selected option 'D' for NVMeMI pass-through model