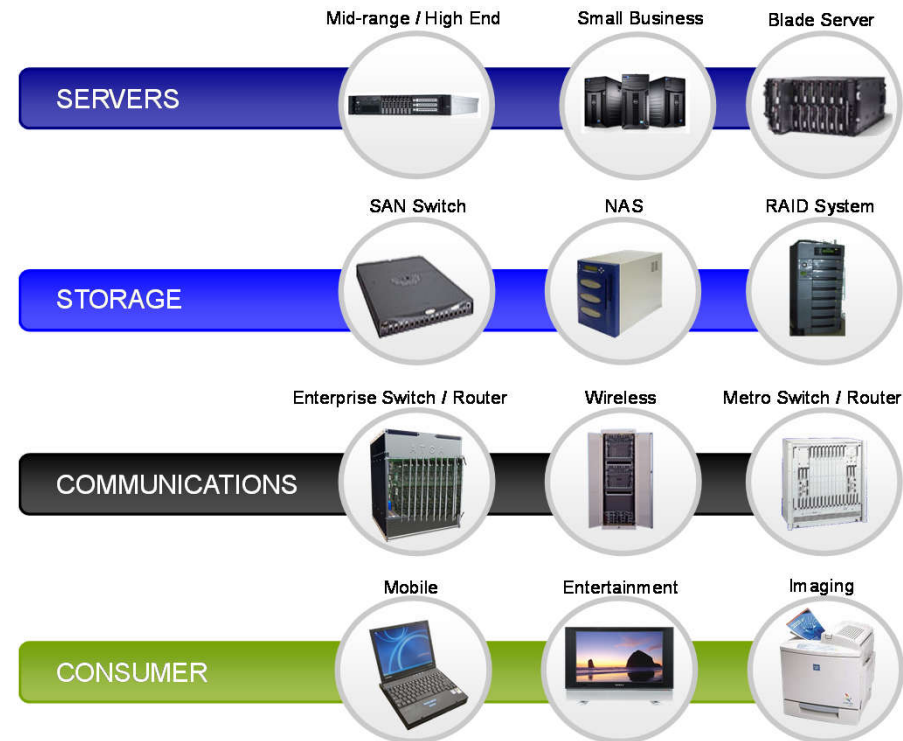# NVM
# PCIe Networked ~~Flash~~ Storage

## Peter Onufryk

## Microsemi Corporation

# PCI Express (PCIe)

- Specification defined by PCI-SIG
  - www.pcisig.com

- Packet-based protocol over serial links
  - Software compatible with PCI and PCI-X
  - Reliable, in-order packet transfer

- High performance and scalable from consumer to Enterprise
  - Scalable link speed (2.5 GT/s, 5.0 GT/s, 8.0 GT/s)
  - Scalable link width (x1, x2, x4, .... x32)

- Primary application is as an I/O interconnect

**SERVERS**
Mid-range / High End   Small Business   Blade Server

**STORAGE**
SAN Switch   NAS   RAID System

**COMMUNICATIONS**
Enterprise Switch / Router   Wireless   Metro Switch / Router

**CONSUMER**
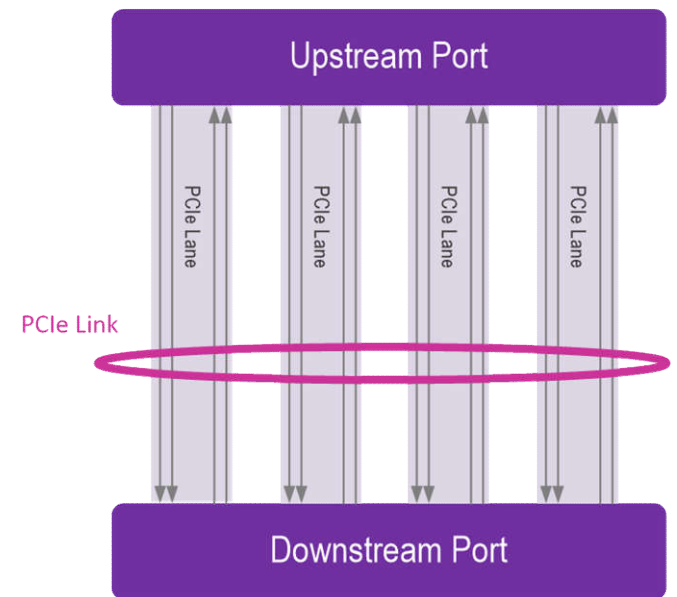Mobile   Entertainment   Imaging

# PCIe Characteristics

- Scalable speed
- Scalable width: x1, x2, x4, x8, x12, x16, x32
- Encoding
  - 8b10b: 2.5 GT/s and 5 GT/s
  - 128b/130b: 8 GT/s and 16 GT/s

| Generation | Raw Bit Rate | Bandwidth Per Lane Each Direction | Total x16 Link Bandwidth |
|---|---|---|---|
| Gen 1* | 2.5 GT/s | ~ 250 MB/s | ~ 8 GB/s |
| Gen 2* | 5.0 GT/s | ~500 MB/s | ~16 GB/s |
| Gen 3* | 8 GT/s | ~ 1 GB/s | ~ 32 GB/s |
| Gen 4 | 16 GT/s | ~ 2 GB/s | ~ 64 GB/s |

Note
*Source – PCI-SIG PCI Express 3.0 FAQ*



Upstream Port

PCIe Lane

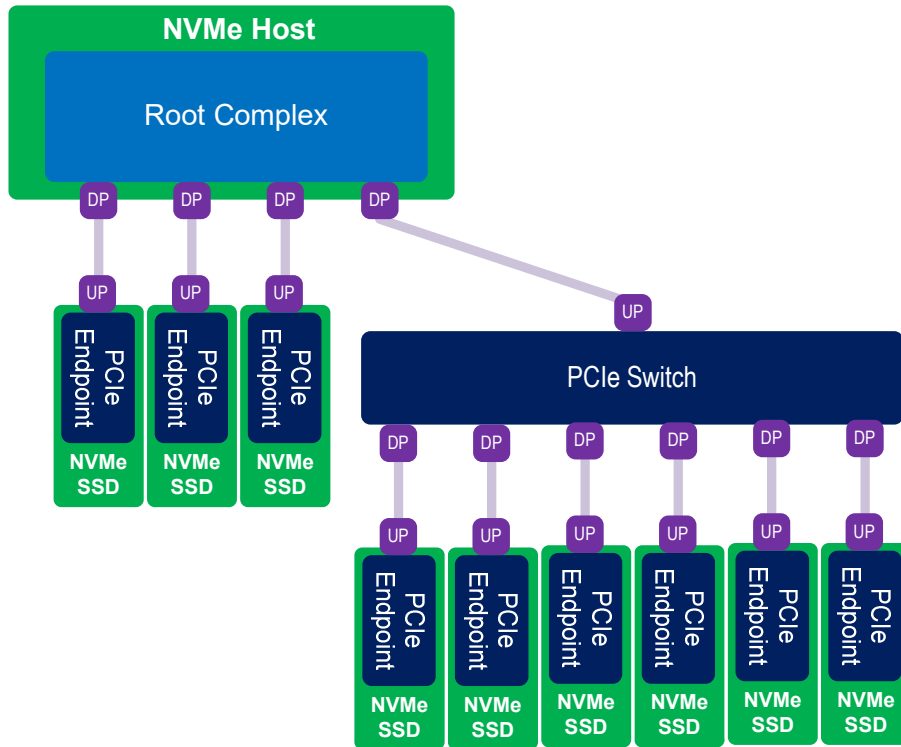PCIe Link

Downstream Port

# NVM Express (NVMe)

- **Two specifications**
  1. NVM Express (PCIe)
  2. NVM Express over Fabrics (RDMA and Fibre Channel)

- **Architected from the ground up for NVM**
  - Simple optimized command set
  - Fixed size 64 B commands and 16 B completions
  - Supports many-core processors without locking
  - No practical limit on the number of outstanding requests
  - Supports out-of-order data deliver
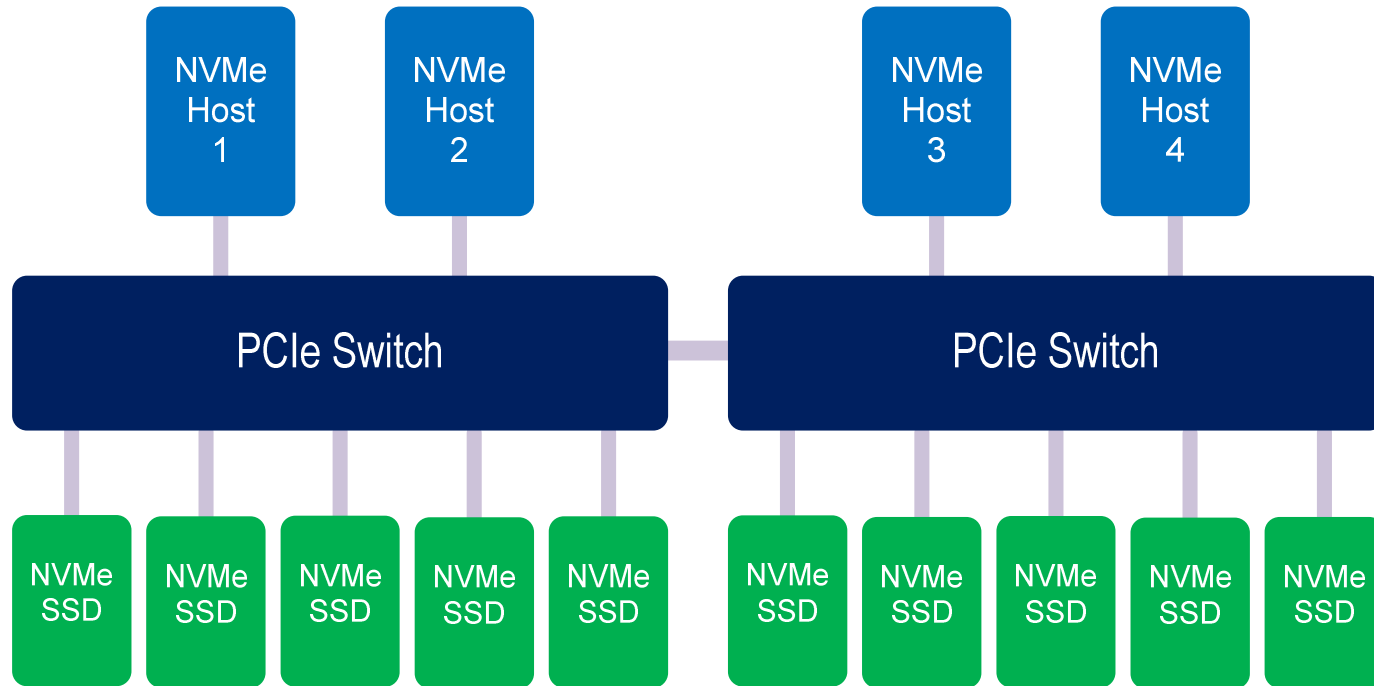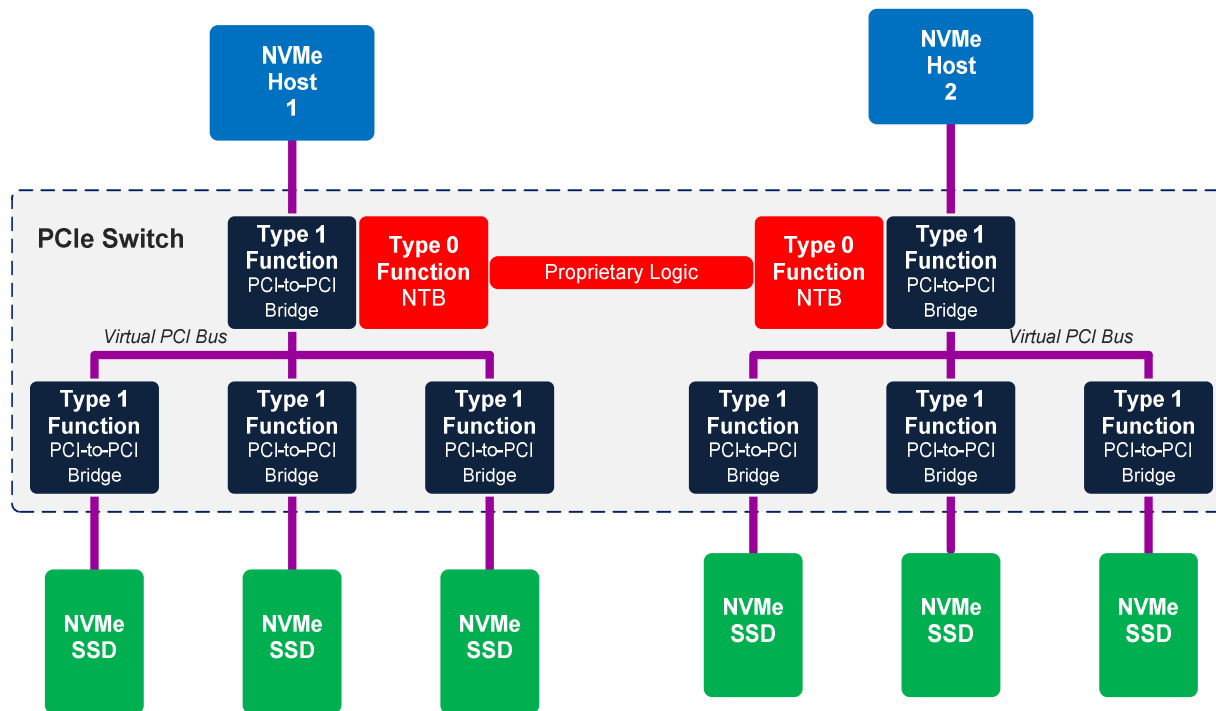


**PCIe SSD = NVMe SSD**

# PCIe and NVMe

**NVMe Host**

Root Complex

DP DP DP DP

UP UP UP UP

PCIe Endpoint — NVMe SSD
PCIe Endpoint — NVMe SSD
PCIe Endpoint — NVMe SSD

**PCIe Switch**

DP DP DP DP DP DP

UP UP UP UP UP UP

PCIe Endpoint — NVMe SSD
PCIe Endpoint — NVMe SSD
PCIe Endpoint — NVMe SSD
PCIe Endpoint — NVMe SSD
PCIe Endpoint — NVMe SSD
PCIe Endpoint — NVMe SSD

**Card**

**U.2**

**M.2**

# Ideal NVM Fabric

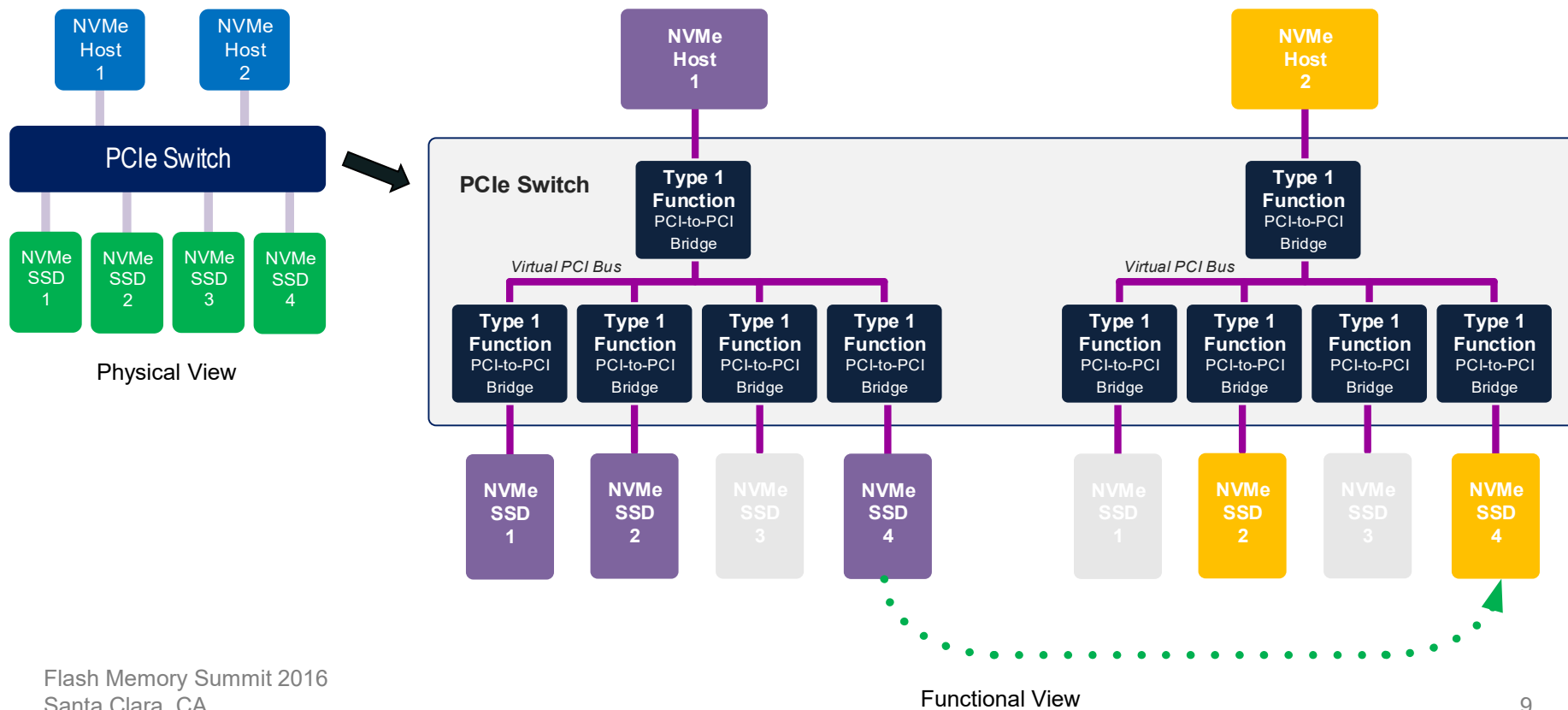| Property | Ideal Characteristic |
|---|---|
| Cost | Free |
| Complexity | None |
| Performance | High |
| Power consumption | None |
| Standards-based | Yes |
| Scalability | Infinite |

# PCIe Fabric

# Non-Transparent Bridging (NTB)

# Dynamic Partitioning



Physical View
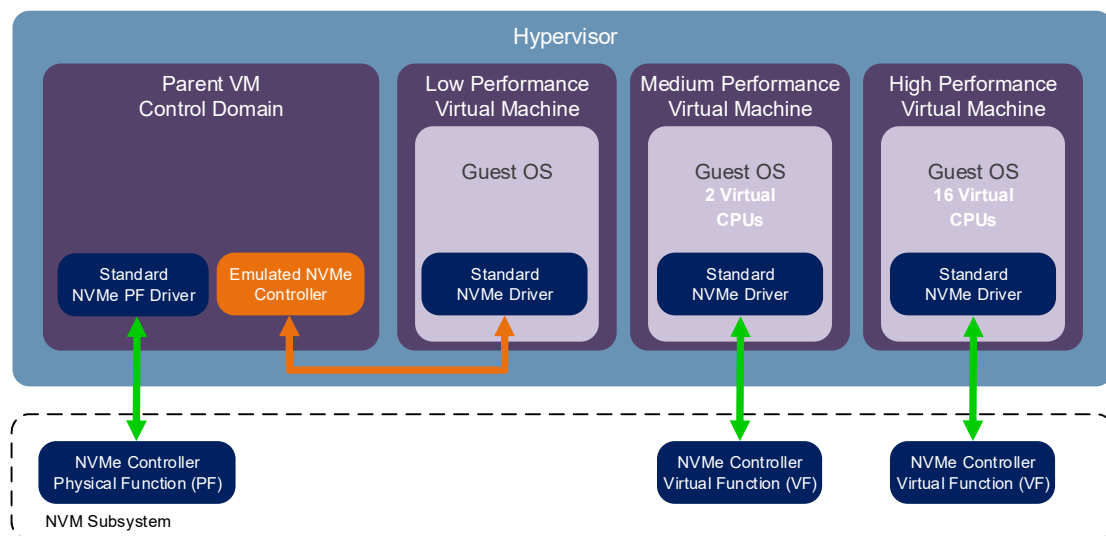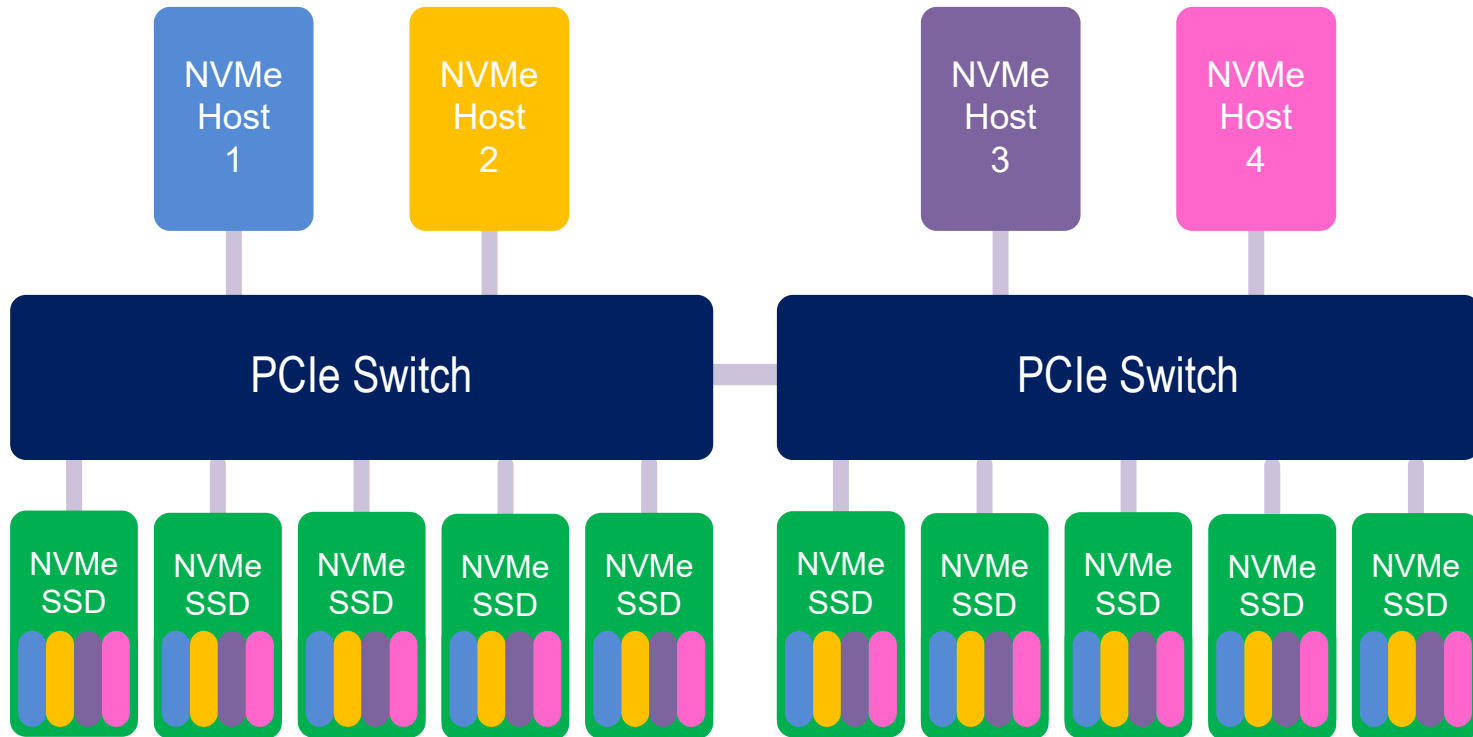
Functional View

# NVMe SR-IOV

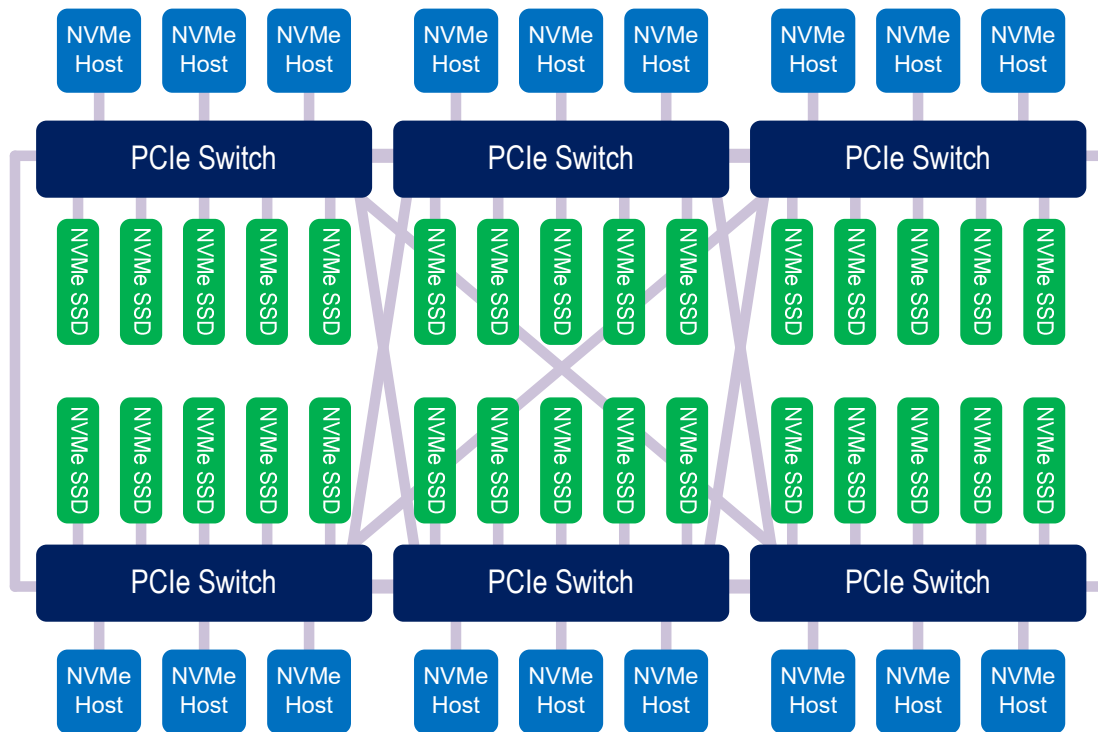# NVMe Single Root I/O Virtualization (SR-IOV)



- Datacenter service providers that host VMs win with high density, oversubscription, and differentiation
  - Multi-tenancy is the norm
  - Premium differentiator offering high speed storage
- Virtual machines are inherently mobile and hosts are inherently dynamic
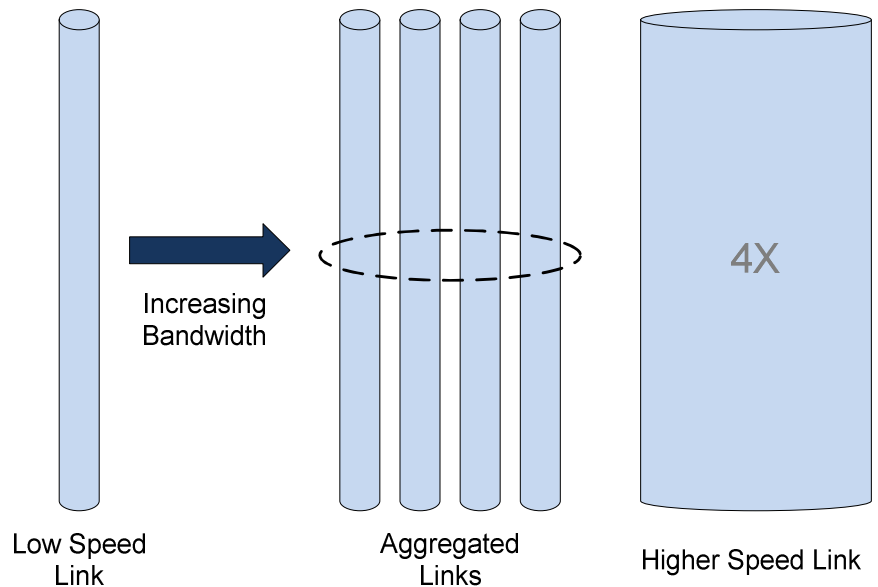
# Multi-Host I/O Sharing

# PCIe Fabric



- **Storage Functions**
  - Dynamic partitioning (drive-to-host mapping)
  - NVMe shared I/O (shared storage)
  - Ability to share other storage (SAS/SATA)

- **Host-to-Host Communications**
  - RDMA
  - Ethernet emulation

- **Manageability**
  - NVMe controller-to-host mapping
  - PCIe path selection
  - NVMe management

- **Fabric Resilience**
  - Supports link failover
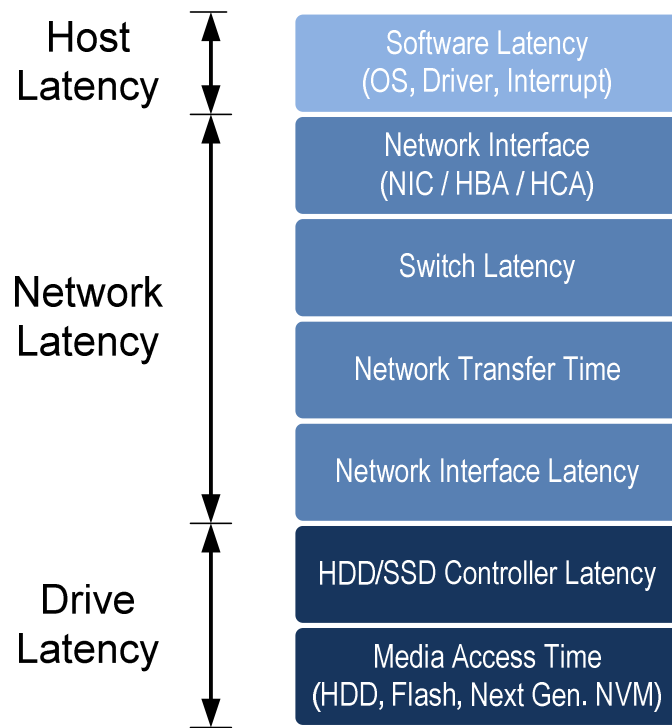  - Supports fabric manager failover

# Fabric Performance

- **A high performance fabric means:**
  - High bandwidth
  - Low latency
- **Increasing bandwidth is easy**
  - Aggregate parallel links
  - Increase link speed (fatter pipe)
- **Reducing latency is hard**
  - Transfer latency is typically a small component of overall latency
  - Other sources of latency:
    - Software (drivers)
    - Complex protocols
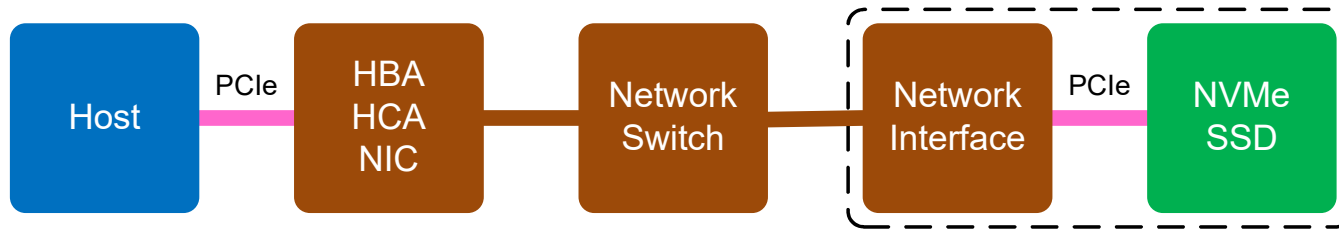    - Protocol translation
    - Fabric switches/hops

Increasing Bandwidth

4X

Low Speed Link

Aggregated Links

Higher Speed Link

# Latency and Next-Generation NVM



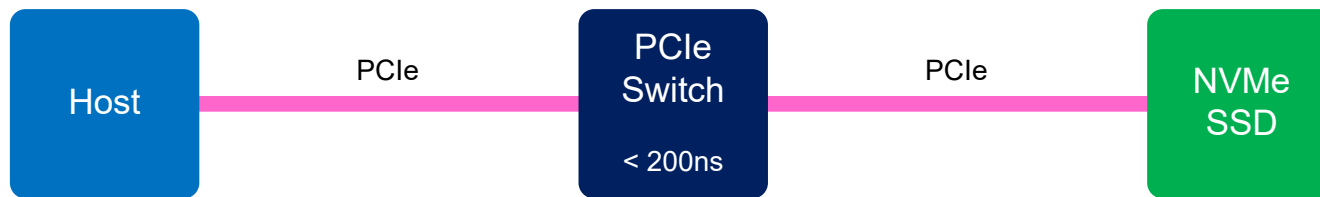| Host Latency | Software Latency (OS, Driver, Interrupt) |
|---|---|
| Network Latency | Network Interface (NIC / HBA / HCA) |
| | Switch Latency |
| | Network Transfer Time |
| | Network Interface Latency |
| Drive Latency | HDD/SSD Controller Latency |
| | Media Access Time (HDD, Flash, Next Gen. NVM) |

- Media Access Time
  - Hard drive – Milliseconds
  - NAND flash –Microseconds
  - Next-gen. NVM – Nanoseconds

# The PCIe Advantage

Host — PCIe — HBA HCA NIC — Network Switch — Network Interface — PCIe — NVMe SSD

**Other Flash Storage Networks**

Host — PCIe — PCIe Switch < 200ns — PCIe — NVMe SSD

**PCIe Fabric**

# PCIe Fabric Characteristics

| Property | Ideal Characteristic | PCIe Fabric | Notes |
|---|---|---|---|
| Cost | Free | Low | • PCIe built into virtually all hosts and NVMe drives |
| Complexity | None | Medium | • Builds on existing NVMe ecosystem with no changes<br>• PCIe fabrics are an emerging technology<br>• Requires PCIe SR-IOV drives for low-latency shared storage |
| Performance | High | High | • High bandwidth<br>• The absolute lowest latency |
| Power consumption | None | Low | • No protocol translation |
| Standards-based | Yes | Yes | • Works with standard hosts and standard NVMe SSDs |
| Scalability | Infinite | Limited | • PCIe hierarchy domain limited to 256 bus numbers<br>• PCIe has limited reach (cables)<br>• PCIe fabrics have limited scalability (less than 256 SSDs and 128 hosts) |

# Summary

- PCIe fabrics build on the existing PCIe and NVMe ecosystem
  - Work with standard NVMe SSDs and OS drivers
  - Leverage standard PCIe infrastructure

- PCIe fabrics are well suited for applications that require the absolute lowest latency and limited scalability
  - NVMe SSD sharing inside a rack
  - Small clusters

- PCIe fabrics are not well suited for long reach applications or where a high degree of scalability is required
  - NVMe over fabrics is well suited for these applications