



Choosing the right NVMe server storage platform for your application

Stuart Campbell
Principal Platform Architect
E8 Storage

Santa Clara, CA
August 2017



Abstract

- This presentation is based on E8 Storage's experience in selecting and qualifying a HA NVMe server hardware platform for a storage application, with 24x U.2 dual ported NVMe SSDs.
- All NVMe platforms are not equal, and each application will have different requirements.
- This presentation will look at a number of capabilities and features of NVMe server storage platforms, and help guide selection based on the requirements of the application.

The logo for the Flash Memory Summit, featuring a stylized sunburst icon in yellow and orange above the text "Flash Memory" in blue and "SUMMIT" in white on a blue background.

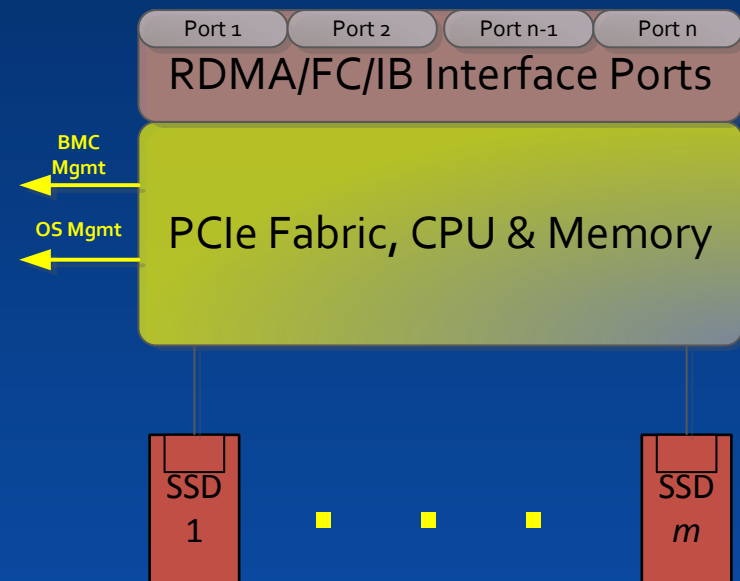
Flash Memory Summit Agenda

- Platform Selection
- SSD Selection
- PCIe Hot-plug & Management
- SSD Management
- System Power
- System Thermals



Platform Selection – Introduction

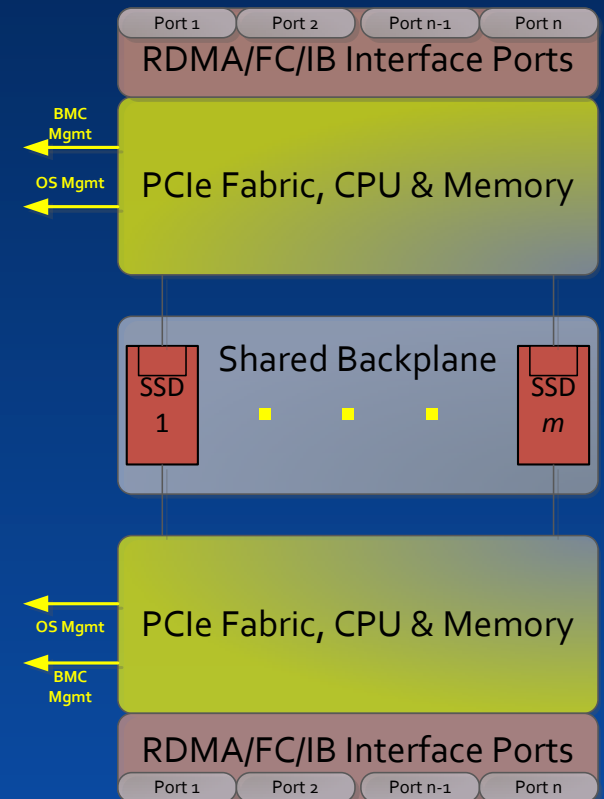
- Basic solution requirements – answer following questions:
 - HA or single server node?
 - Serviceability requirements?
 - PCIe connectivity?
 - CPU requirements?
 - Memory requirements?
 - Power-loss protection?
 - Management capability?
 - Boot device?





Platform Selection – HA

- Why HA?
 - Increased Reliability, Availability and Serviceability (RAS)
- Per node requirements:
 - Power isolation and control
- Shared platform requirements:
 - SSD power control
 - Thermal monitoring & control
 - No SPOF
 - Isolated clocking (or SRIS), resets, power
 - “Low complexity” common components





Platform Selection – Serviceability

- Serviceability
 - What parts are FRUs or customer replaceable?
 - Can all serviceable parts be removed without tools?
 - Can be more challenging in higher density systems
- Mechanical
 - Is there a limit for physical size, e.g. 1U / 2U versus 4U?
 - Simpler installations, less shipping/installation weight restrictions
- Off-the-shelf
 - As much as possible use off-the-shelf components including standard NICs
 - Enables easier migration to different standards (OCP, etc.)



Platform Selection – PCIe Connectivity

- Depends on root-complex capability
 - Each RC has specific number of ports/lanes
 - PCIe switch(es) required if insufficient lanes/ports on root-complex
 - Allocate sufficient PCIe lanes to SSDs to match required system bandwidth performance
- Ensure external interfaces are not bottleneck
 - 50 / 100GbE RDMA, IB or FC
 - Match PCIe NIC/HBA connectivity to SSDs for balanced end-to-end performance



Platform Selection – CPU/Memory

- Required CPU capability depends on application
 - In-path processing of commands and/or data requires more CPU power
 - Additional data protection (e.g. RAID) processing requires more CPU power
 - CPU core/frequency limits defined by architecture and power/thermal
- Required memory configuration depends on application
 - Best to populate all memory channels to provide adequate bandwidth
 - Use pcm (<https://github.com/opcm/pcm>) to monitor memory usage
 - Every data copy uses bandwidth, so application architecture is critical to achieving maximum performance



Platform Selection – Misc.

- Power loss protection
 - Battery backup or NVDIMM
 - Requires application architecture support

- Onboard Management
 - 1GbE or 10GbE – number of ports
 - Dedicated or shared IPMI

- Boot device
 - Use industry standard – M.2 NVMe is now common
 - Use multiple devices and RAID on motherboard



SSD Selection Criteria

- Direct
 - Performance
 - Media Type
 - Endurance
 - Single/Dual Port

- Indirect
 - Pricing
 - Support
 - Warranty
 - Vendor Relationships



SSD Selection – Performance / Media

- Need to select for both Bandwidth & IOPs
- What performance can the platform and/or application deliver?
 - No need to use higher performing SSD than SW/HW can use
 - Suitable application can allow for SSD performance to be shared
- What performance does the solution need?
 - May depend on end-user requirements
 - Requires flexibility in price/performance tradeoff
- What latency does the application add?
 - Does it benefit from Storage Class Memory?



SSD Selection – Endurance

- What is the expected (or guaranteed to customer) amount of write data?
 - Ensure accurate estimation of expected system performance, and customer expectation on SSD lifetime
 - Ensure selected endurance level can meet performance requirements
- What is the write amplification in the solution?
 - Can have large impact on endurance, optimize where possible
- Does the SSD support dynamic over-provisioning?
 - Depends on SSD vendor, single SKU reduces inventory requirements



SSD Selection – Single/Dual Port

- Mostly driven by solution requirements (HA or not)
- Most datacenter SSDs support single or dual port in same SKU
 - Reduces qualification required
- Dual-port performance requires active-active application architecture
 - Required for full performance from dual port, even for IOPs
- Dual port requires application synchronization of I/Os between ports
 - Similar to previous SAS/FC HA implementations



PCIe Hot-plug & Management

- Hot-plug & Resiliency
- PCIe Enumeration
- Device Identification
- Error Monitoring
- Performance Monitoring



Hot Plug & Resiliency Support

- Hardware Support
 - Hot-plug capable hardware, electrically and with notification capability
- Firmware Support
 - Hot-plug aware BIOS for resource allocation
- OS Support
 - Hot-plug aware OS/Drivers for re-enumeration and graceful removal
- Advanced Features
 - Downstream Port Containment (DPC) in switches
 - Enhanced DPC in CPU



PCIe Enumeration

- To support hot-plug, pre-allocate sufficient resources in BIOS or use advanced PCIe switch/fabric features
 - PCI Busses
 - Memory (prefetchable, non-prefetchable)
 - I/O space is typically not required, and insufficient for large number of SSDs
- Fixed bus numbers and slot for physical locations makes it much easier to debug (see next slide)



Slot & NVMe Device Identification

- Note multiple references in dmesg to same device
- Example shows device being removed

```
[1310.866368] pciehp 0000:80:03.0:pcie04: Card not present on Slot(1)
[1310.866380] pciehp 0000:80:03.0:pcie04: slot(1): Link Down event
[1311.169637] nvme 0000:86:00.0: Failed status: 0xffffffff, reset controller
[1311.170005] blk_update_request: I/O error, dev nvme6n1, sector 2725422000
[1311.173279] nvme 0000:86:00.0: Removing after probe failure status: -19
[1311.173294] nvme6n1: detected capacity change from 2000398934016 to 0
```



PCIe Error Monitoring

- AER Registers
 - Monitor PCIe AER registers on both sides of the link
 - Clear errors after bootup
 - Correctable errors should be zero or very low rate
 - Uncorrectable errors shouldn't occur on stable system
 - Regularly monitor, clear and count errors
- Link Status Registers
 - Monitor for correct link speed and width
- Use vendor-specific advanced error counters where available



PCIe Performance Monitoring

- Advanced performance monitoring tools available from CPU/SOC or switch vendors
- Can be used to check performance on every link to look for issues
- Can validate packet size and other parameters
- Ask for platform firmware (BMC) to expose these



SSD Management

- Multiple standards for I2C / SMBus SSD management
 - NVMe-MI
 - Enterprise SSD Form Factor 1.0a
 - Key Capabilities
 - Temperature (for platform management)
 - Model / Serial Number
 - Often available over IPMI – can be used remotely & during manufacturing
 - Requires HA support if HA platform

- Platform Features
 - Power control over IPMI
 - LED control over IPMI



System Power Capability

- Ensure PSUs are rated for max system power consumption
 - Based on application power usage, and may vary
 - U.2 SSDs can consume up to 25W
 - In an example 24 SSD system, max power is up to 600W
 - Difficult to get all SSDs consuming 25W, depends on upstream connectivity and data path width
 - SSD power cap feature is useful and may be necessary
 - Throttle CPU if necessary to keep within power budget
 - CPU utilization can be large factor
 - Determine if 100V / 110V operation is required and architect accordingly

- Use Instrumented PSUs and subcomponents
 - Some systems have this built in, and expose via IPMI



System Thermal Capability

- Related to power utilization, and to architecture
 - Location of fans, CPUs, SSDs etc. is important
- Test and validate based on product requirements
 - e.g. Ashrae A2/A3/A4, or specific range such as 5 - 35C
- Determine if CPU/Memory/SSD throttling is acceptable
 - For normal use and fan fail cases
 - Develop FRU replacement strategy
- Use system thermal sensors to monitor
 - CPUs, Memory, Add-in-Cards, SSDs, available via IPMI



Questions?

stuart@e8storage.com