



How Facebook and Microsoft Successfully Leverage NVMe[™] Cloud Storage

Sponsored by NVM Express[™] organization, the owner of NVMe[™], NVMe-oF[™] and NVMe-MI[™] standards

Speakers

Ross Stenfort



Lee Prewitt





NVMe[™] In The Real World

Ross Stenfort, Hardware System Engineer

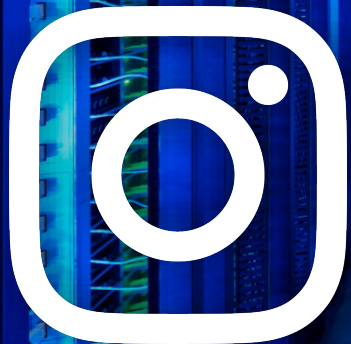
Facebook



facebook

**Facebook's mission is to give people the power to build
community and bring the world closer together.**

Facebook @ Scale



1 Billion



1.3 Billion



2.7 Billion

PAPILLION, NE



LOS LUNAS, NM



PRINEVILLE, OR



FOREST CITY, NC



FORT WORTH, TX



NEWTON COUNTY, GA



NEW ALBANY, OH



ODENSE, DENMARK



LULEÅ, SWEDEN



ALTOONA, IA



CLONEE, IRELAND



HENRICO, VA



EAGLE MOUNTAIN, UT



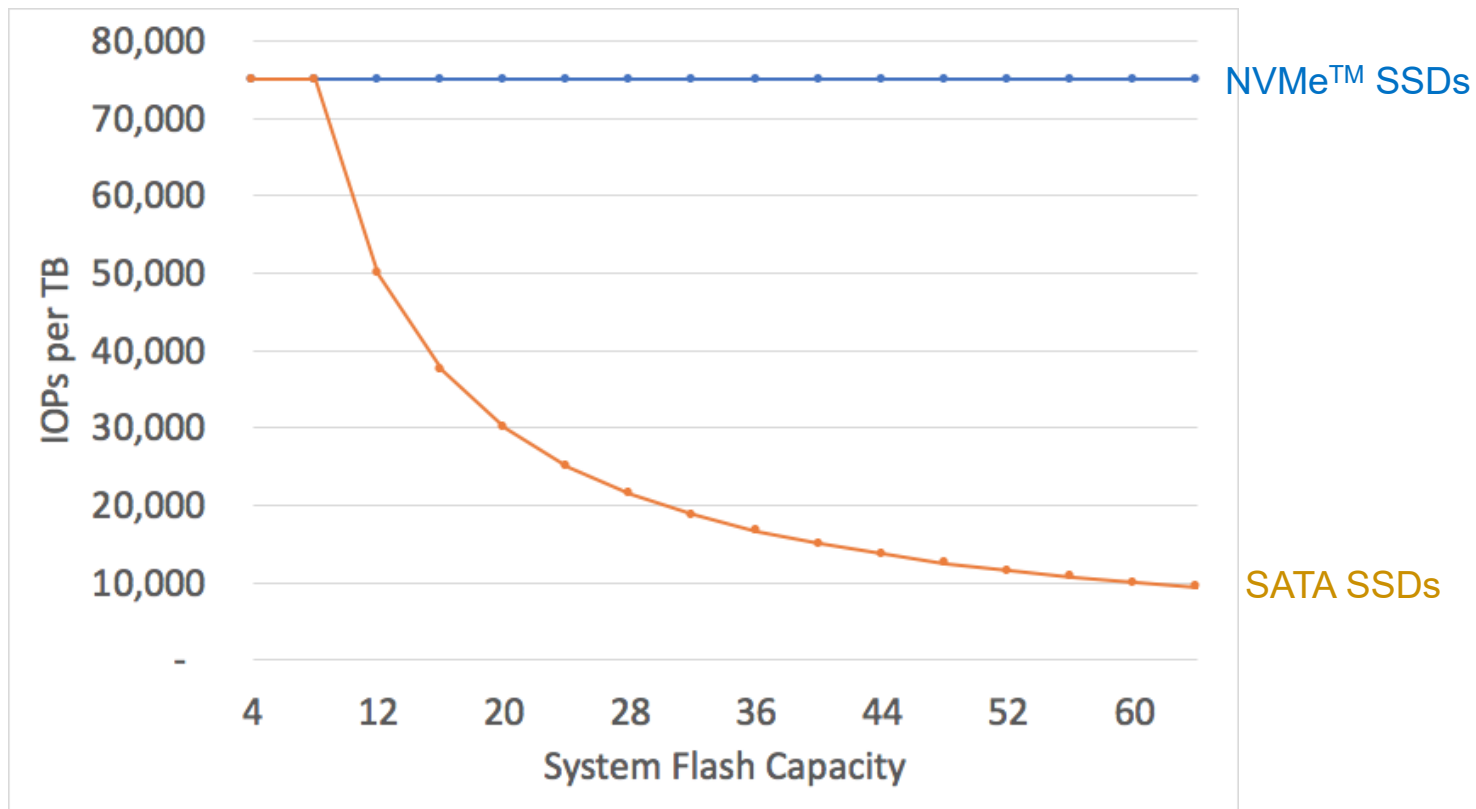
HUNTSVILLE, AL



SINGAPORE



Hyperscale Requires IOPS to Scale with Capacity



NVMe™ De-Allocate: Challenges and Improvements

➤ NVMe™ De-allocate

- Goal: It's a hint from the system to the SSD that the system is no longer tracking certain LBAs
- Good
 - Reduces Write Amplification
 - Improves performance/endurance
- Bad
 - Latency spikes due to De-allocate blocking Read/ Write

➤ Old Solution

- Tune De-Allocate size on a system
- Problem: The optimized de-allocate size varies based on supplier. Thus which supplier should I optimize for?

➤ Improved solution

- NVMe 1.4 allows the SSD to advertise its preferred De-allocation size
 - If NSFEAT bit 4 = 0x1 then Namespace Preferred Deallocation Granularity (NPDG) is valid
- This allows systems to be optimized standard mechanisms.



Flash Memory Summit

nvm
EXPRESS®

Managing at Scale (1 of 3)

- Challenge: Hyperscale Requires Debug with no physical access to the SSD.
- Challenge#1: Restricted access for vendor unique tools
- Solution:
 - NVMe™ CLI – Open source with active industry contribution and updates
 - <https://github.com/linux-nvme/nvme-cli>
 - Vendor-unique CLI plugin that pulls and reports the logs in a common format

Challenge#2: How do I get the debug information needed to resolve the issue

Solution: Telemetry

- This allows SSD providers to get remote debug information to resolve issues
- Different data areas allows for different levels of debugging



Flash Memory Summit

nvm
EXPRESS®

Managing at Scale (2 of 3)

- Background: The amount of data written to a SSD may exceed the endurance of the SSD given the expected lifetime of the SSD. Given a fixed amount of write bandwidth a low the capacity SSD will wear out faster than a higher capacity SSD. Examples of applications where this can occur are logging and caching.
- Challenge/ Real World Example:
 - Application only needs 256 GB but will use all the SSD capacity
 - Application write rate is high enough that it will wear out the 256 GB SSD
 - Application write rate scales per TB: Thus increasing capacity will not keep the SSD from wearing out.
- Solution: Namespace Management
 - Allows a 512 GB SSD to be configured as a 256GB SSD with double the endurance of a 256 GB SSD
 - Thus the application view is a 256GB with double the endurance



Flash Memory Summit

nvm
EXPRESS®

Managing at Scale (3 of 3)

- Challenge: How many blocks in my SSD have data and how many do not? If I de-allocate some blocks how many blocks really contain data? What is the effective over provisioning from a performance perspective?
- Solution:
 - Namespace Utilization (NUSE)
 - Allows user to determine the number of LBAs that actually contain data.



Industry Challenge

➤ Security challenges are growing

- NVM Express™ supports SECURITY_SEND/ RECIEVE will allows for security protocols to be tunneled into NVM Express
- There is even an open source tool for NVMe™ Opal security:
<https://github.com/Drive-Trust-Alliance/sedutil>
- Secure Boot is also a common security requirement. This is a process that ensures the firmware running on the device is from the manufacture and not some other source.



➤ Problem/ Industry call to action:

- There is no standard way to know if secure boot failed
- If firmware on a device is compromised, how is this identified vs any other type of failure?



Flash Memory Summit

nvm
EXPRESS®



Flash Memory Summit



NVMe™ at Hyper-Scale

Lee Prewitt, Principle Hardware Program Manager

Azure CSI - Microsoft

- Azure at a Glance
- Why NVMe™?
- Issues at Scale
 - Form factors
 - Need to allow for “rot in place”
 - Need for remote debugging
 - Need for security



2 Million

miles
intra-datacenter fiber

72+

Tb per second
backbone

54

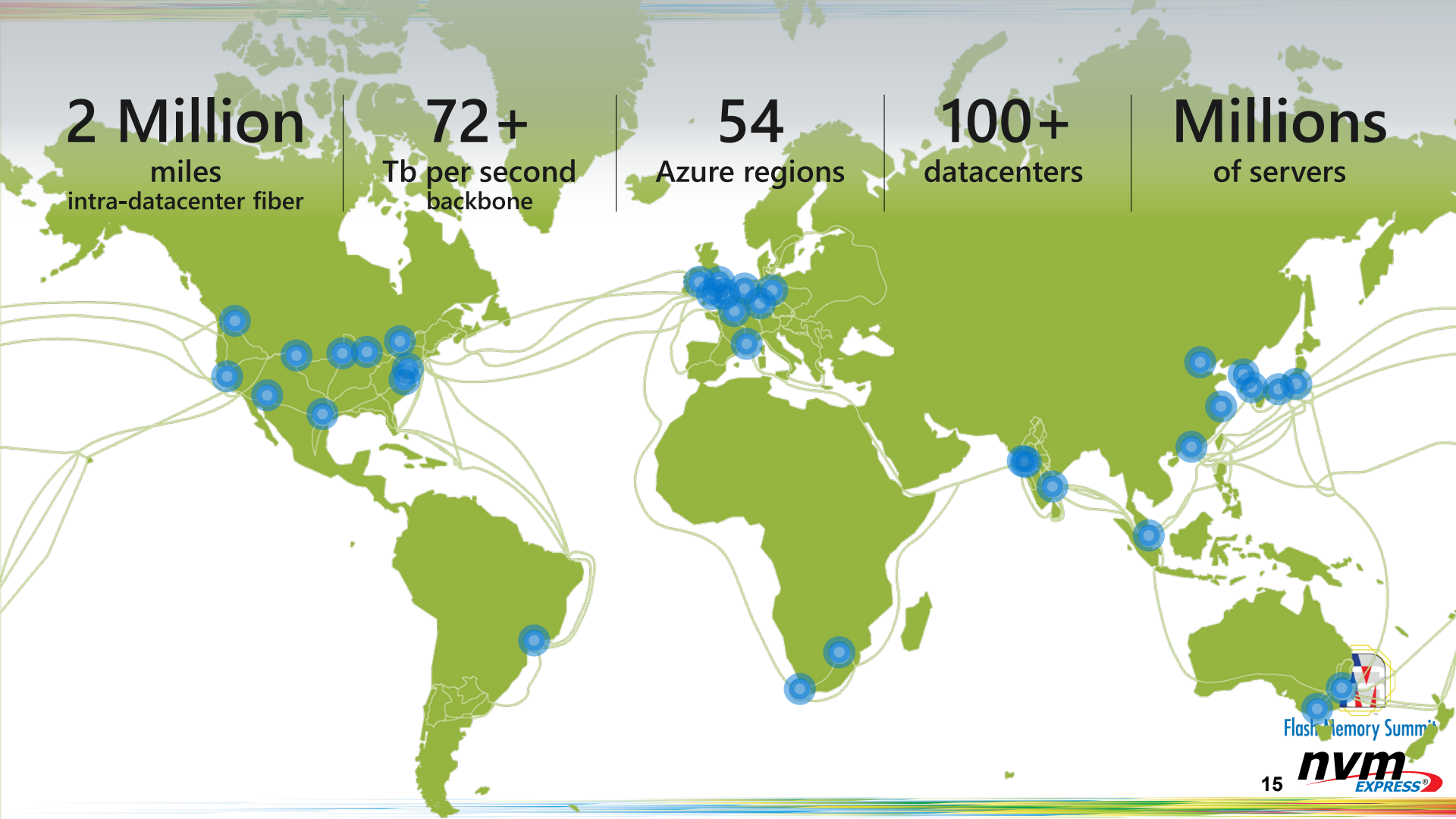
Azure regions

100+

datacenters

Millions

of servers

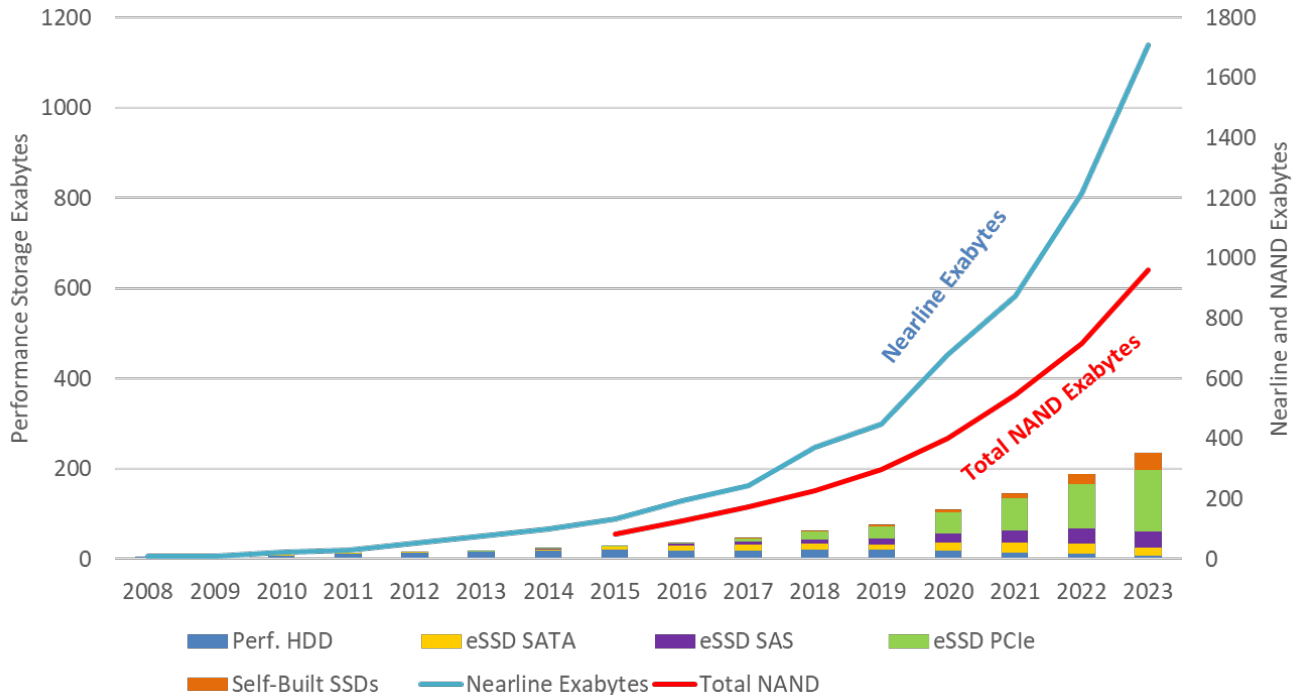


Flash Memory Summit



Why NVMe™? - Exploding Storage Growth

TRENDFOCUS

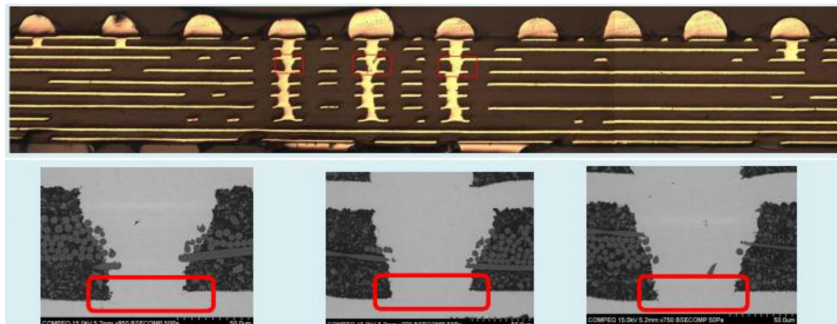


Flash Memory Summit



Issues at Scale – Form Factors

- m.2 has run its course
 - Power and thermal constraints
 - Fragile PCB and connector
 - Not hot-swappable
- E1.L and E1.S are here to replace it
 - Built from the ground up for datacenter use cases
- **Good news is that they support NVMe™ too!**



E1.L (SFF-TA-1007)

- Density Optimized
- 318.75 x 38.4 mm
- Supports > 40W
- Up to 48 Standard NAND sites



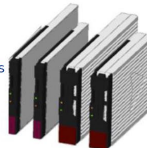
E1.S (SFF-TA-1006)

- 111.5 x 31.5 mm
- Up to 12 Standard NAND sites
- Supports > 12W



E3 (SFF-TA-1008)

- Ultra high-performance applications
- (104.9/142.2) x 78mm
- Supports up to 70W
- Up to 48 Standard NAND sites

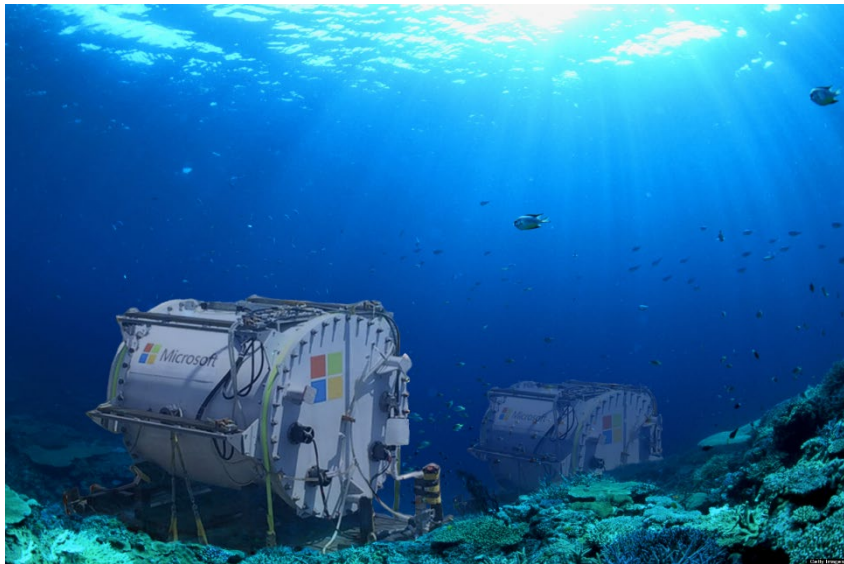


EDSFF Advantages

- **Same** Protocol: NVMe
- **Same** Interface: PCIe
- **Same** Connector: SFF-TA-1002
- **Same** Pinout and Functions

Different Usages
Same Expectations!

Issues at Scale – Need to allow for “Rot in Place”



Use the Endurance and Performance metrics for auto tiering

- Allows for fitting the workload to the device
- Allows for the ability to adjust the temperature of the data over time
- Allow for 5 - 7 year device service life

Zoned Name Spaces for QLC

- Reduce WAF due to large sequential writes
- Reduce DRAM due to large indirection unit
- Reduce overprovisioning due to minimal garbage collection



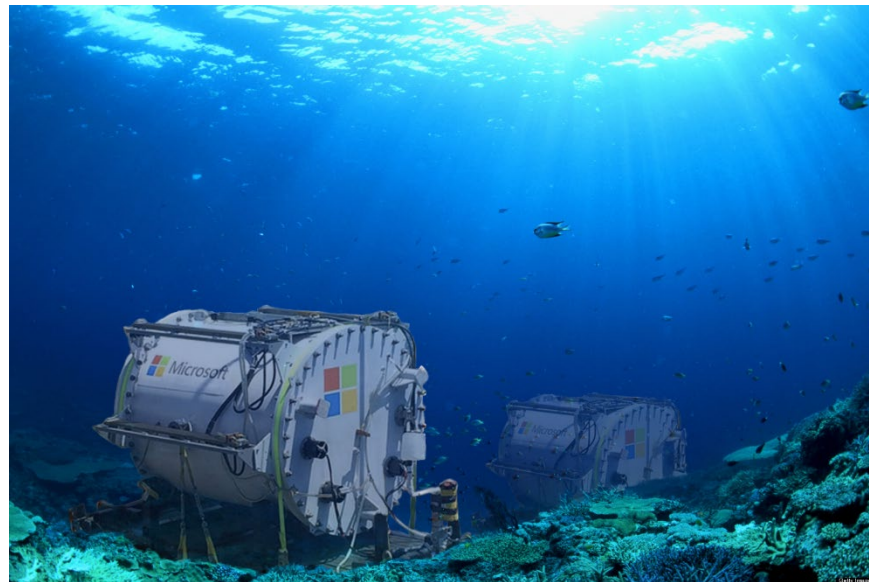
Flash Memory Summit

nvm
EXPRESS®

Issues at Scale – Need for Remote Debugging

- Timestamp
 - Drive events correlated to system (BIOS and OS) events
- Telemetry
 - Host initiated - IO failures
 - Drive Initiated - Firmware panic?
- SMART
 - Both standard and vendor unique collected once an hour
 - Hey SSD IHVs. How many terabytes would you like to see?

Caveat: Any data that leaves the datacenter must be in human readable form!



Flash Memory Summit

nvm
EXPRESS®

Issues at Scale – Need for Security



eDrive on Windows

- Opal v2 plus IEEE 1667 secure silo

Hardware Root of Trust

- Secure boot
- Signed firmware
- Cerberus

Device Hardening

- Pen and Fuzz testing
- Locking of debug ports and vendor unique commands



Flash Memory Summit

nvm
EXPRESS®

Questions?



Flash Memory Summit

nvm
EXPRESS®

