



Enhancing SSD Control of NVMe Devices for Hyperscale Applications

Luca Bert - Seagate

Chris Petersen - Facebook



Agenda

- Introduction & overview (Luca)
- Problem statement & proposed solution (Chris)
- SSD implication and proto work (Luca)
- Test results and analysis (All)
- Conclusion (All)

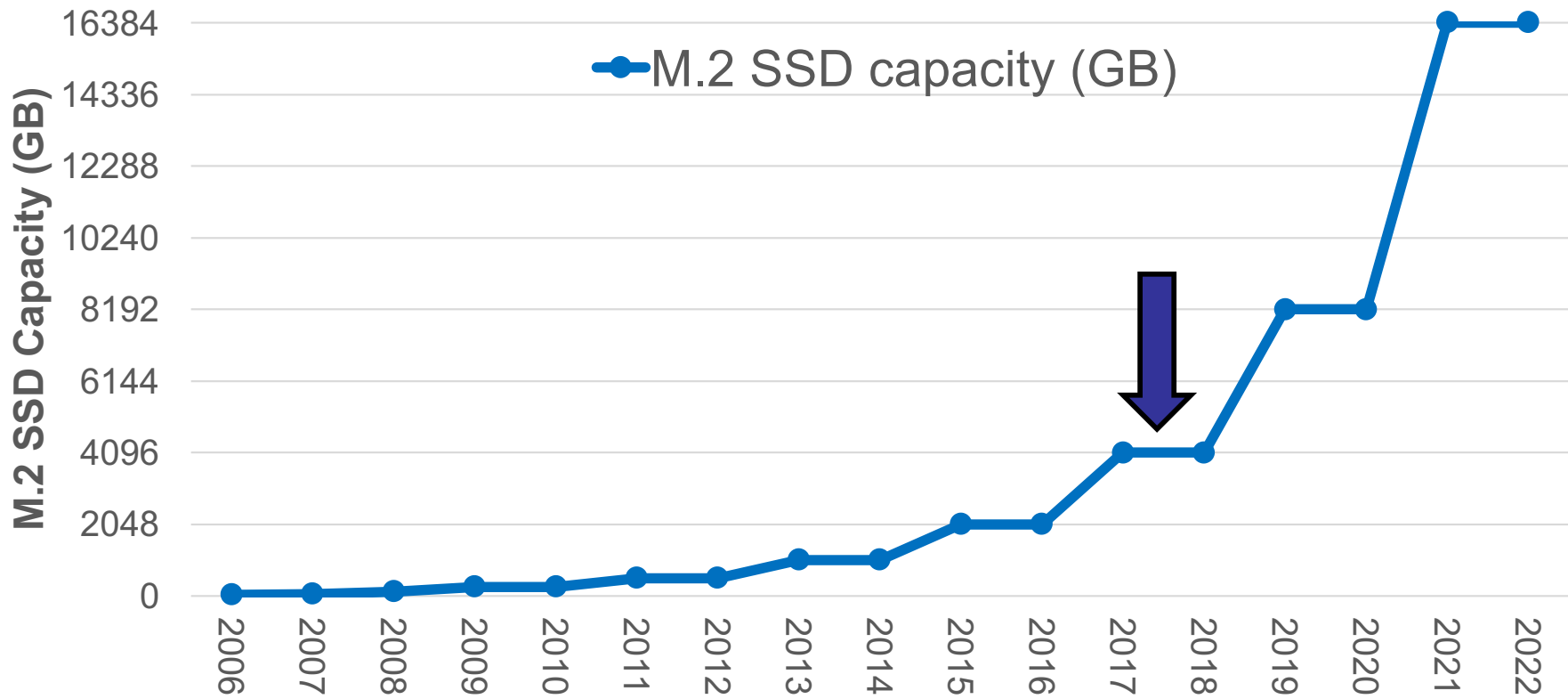


Problem Statement & Proposed Solution



Flash Memory Summit

NAND Flash Trends

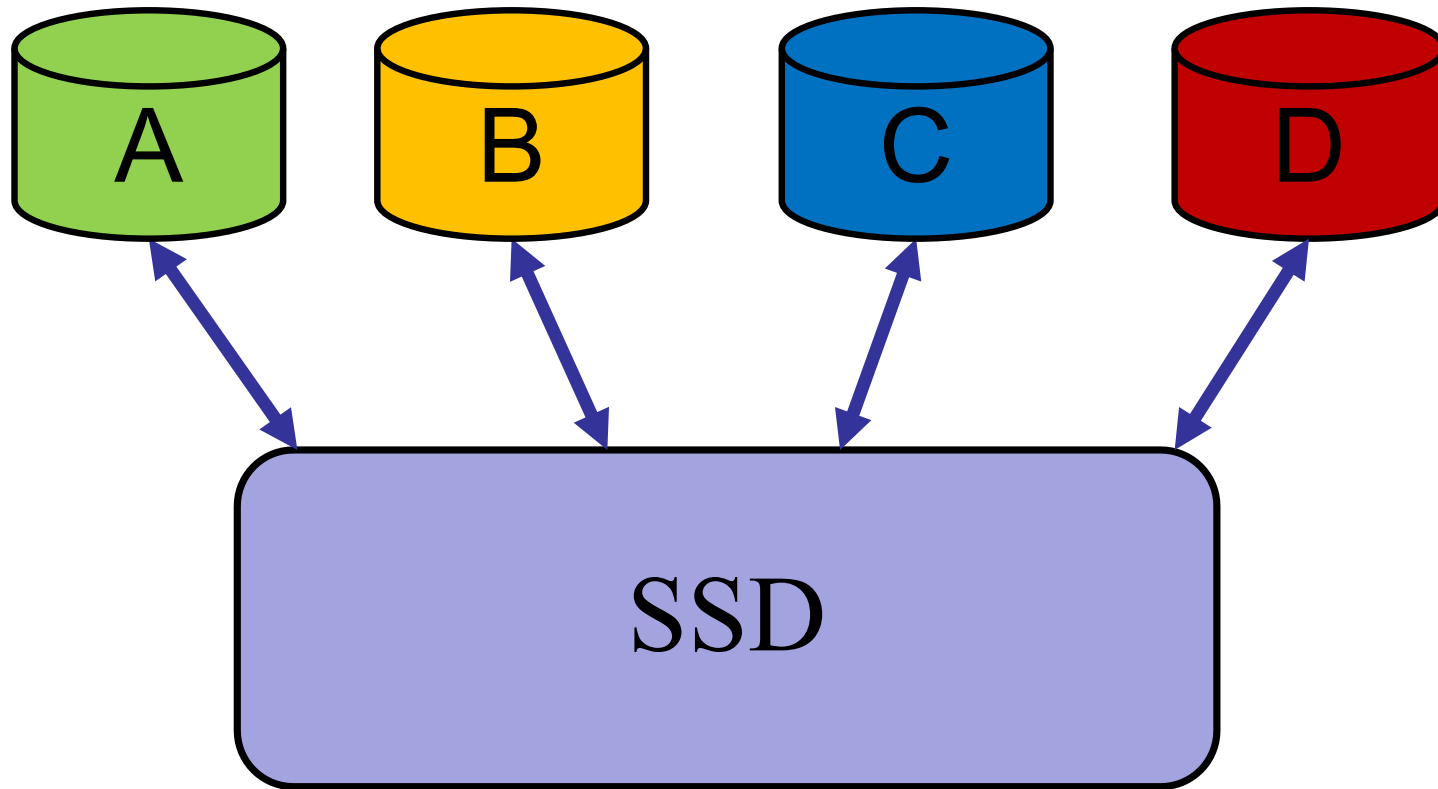


facebook SEAGATE



Flash Memory Summit

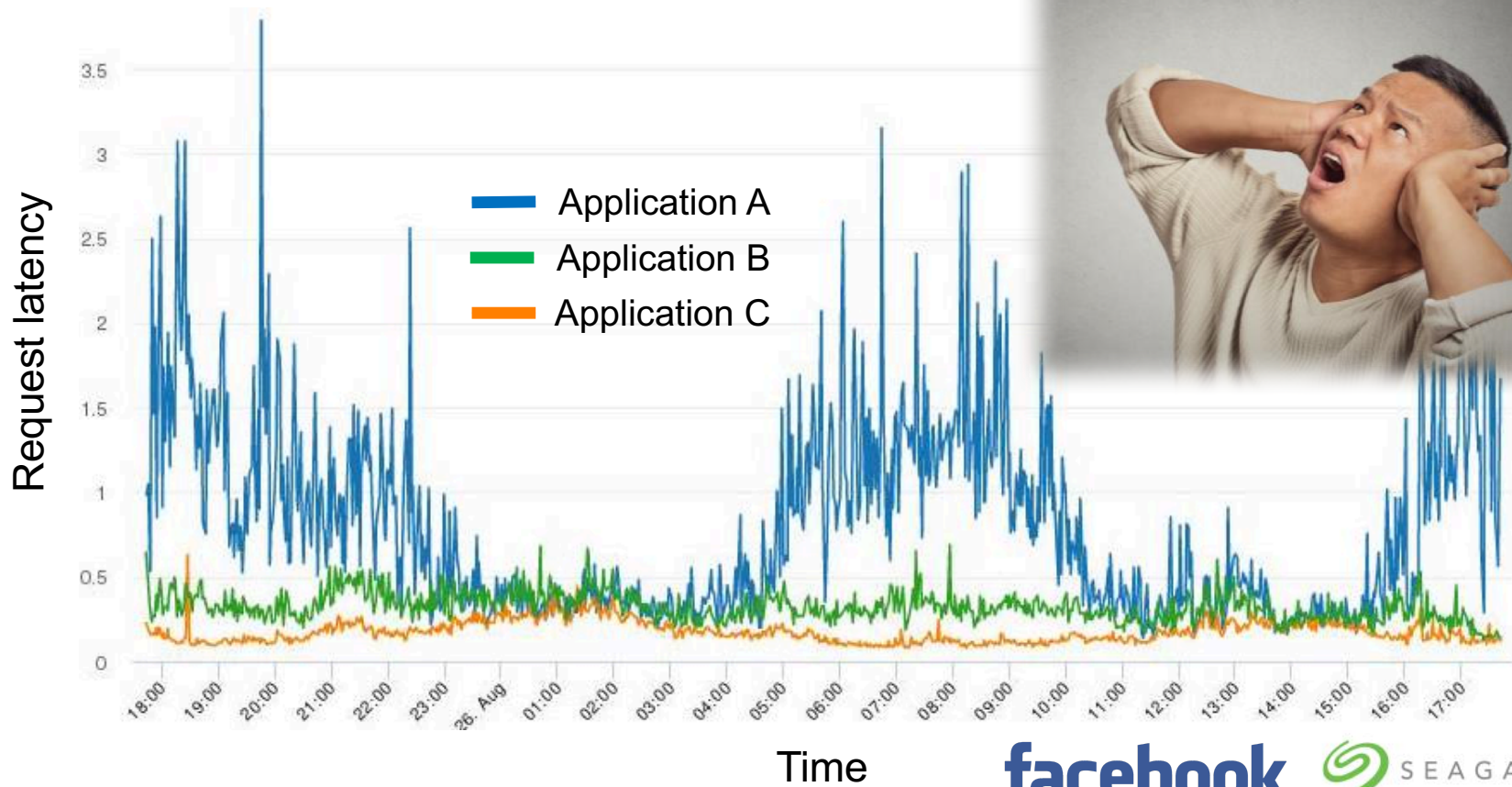
SSD Capacity = Shared Resource





Flash Memory Summit

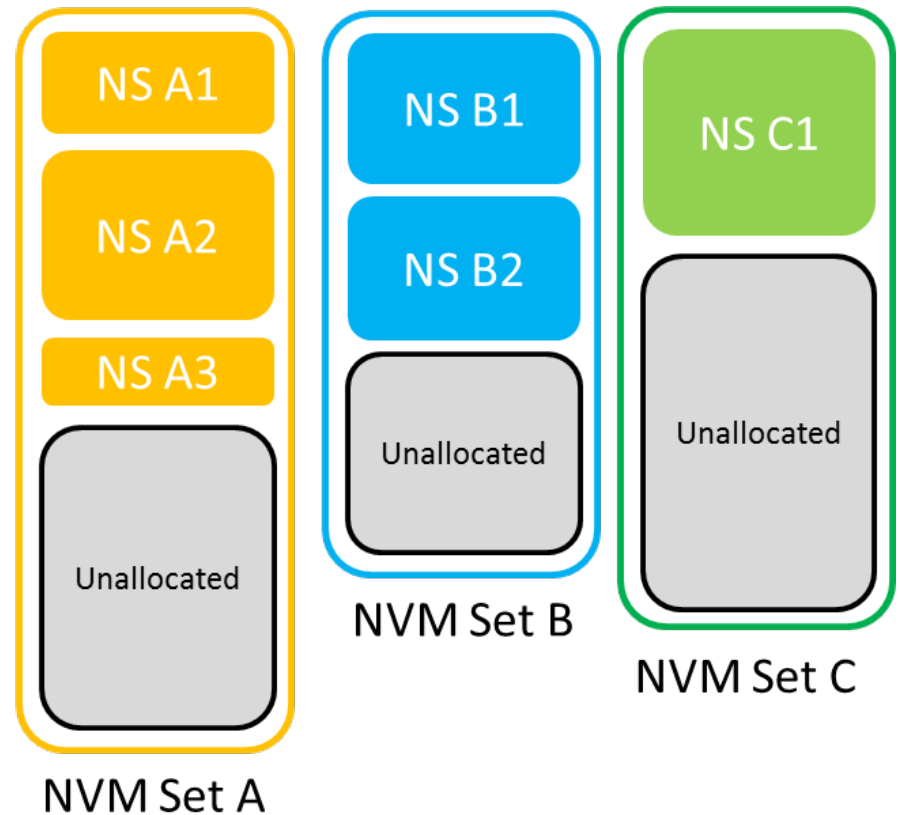
Noisy Neighbors





Introducing I/O Determinism

- New NVMe feature called: I/O Determinism
- SSD is configured as multiple NVM Sets (e.g. A, B, C)
- NVM Sets are QoS isolated regions
 - A write to Set A does not impact a read to Set B or C
- One or more namespaces are allocated from an NVM Set





SSD Implications and Prototype Work

SSD Data Layout Impact

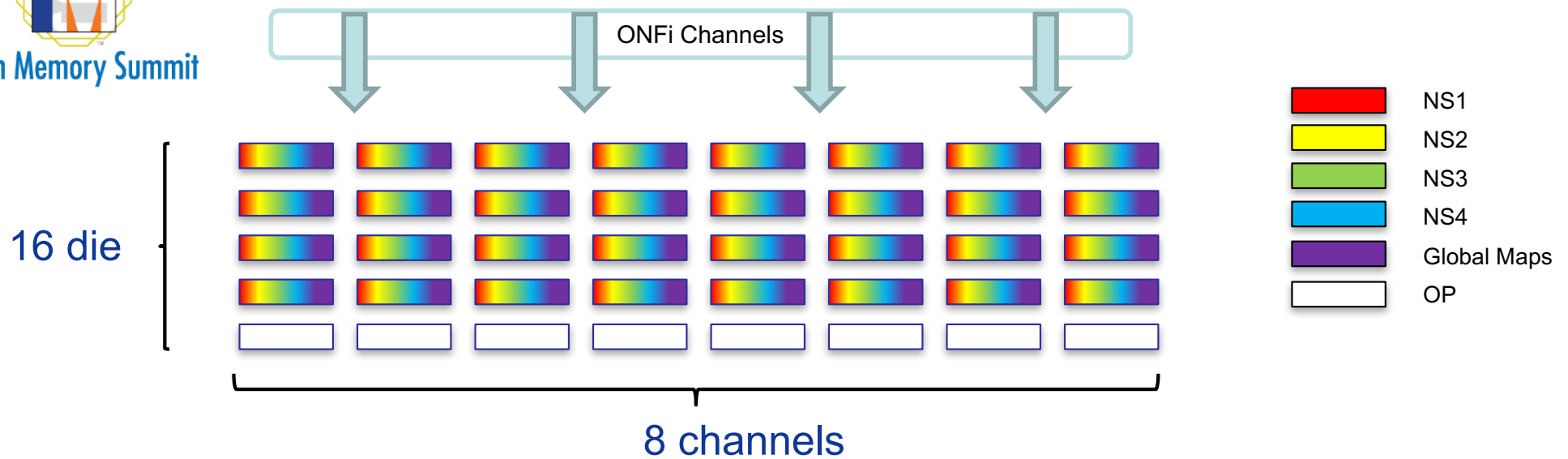


Flash Memory Summit

Key Hurdles to Address

- Physical resource segregation
 - Die have to be allocated to VM Sets
- “Localization of the Globals”
 - **Data:** Maps, Journals, Logs, Metadata... are global and need to be redesigned to be local to VM Sets
 - **Processes:** Same for GC and other background processes
 - **Resources:** OP, cache, memory... can be local or global

Baseline: Conventional SSD with 4 NS



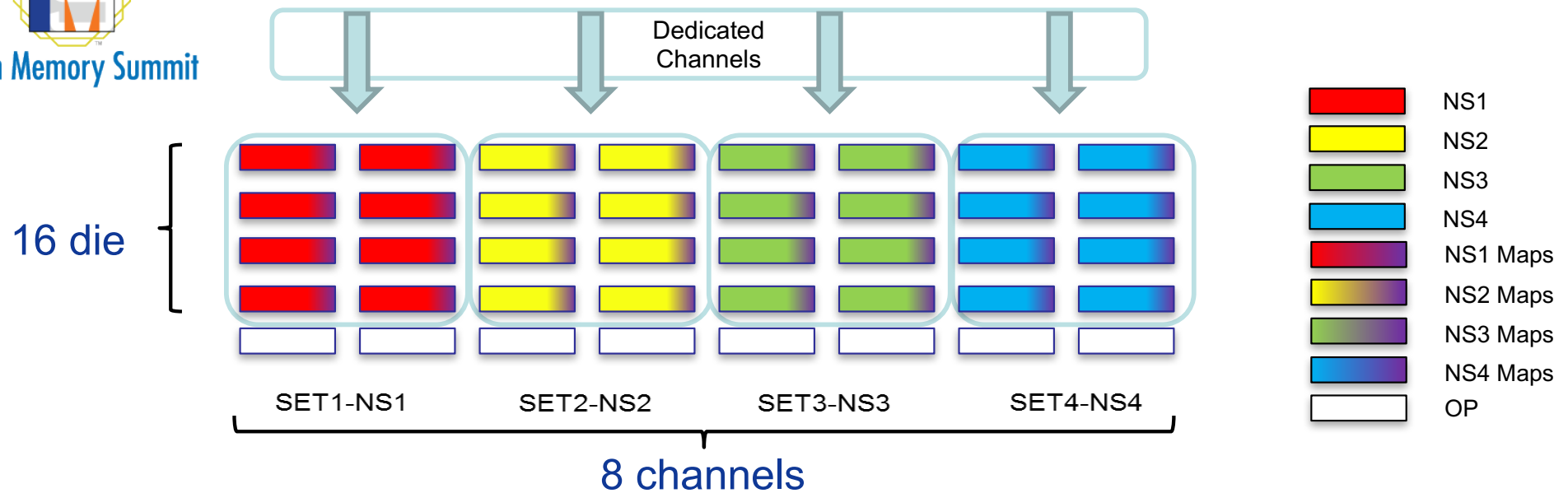
Pros

- Simplest configuration
- No logical – physical affiliation
- Most balanced model

Cons

- IO pattern conflicts

IOD: Vertical Sets (VS)



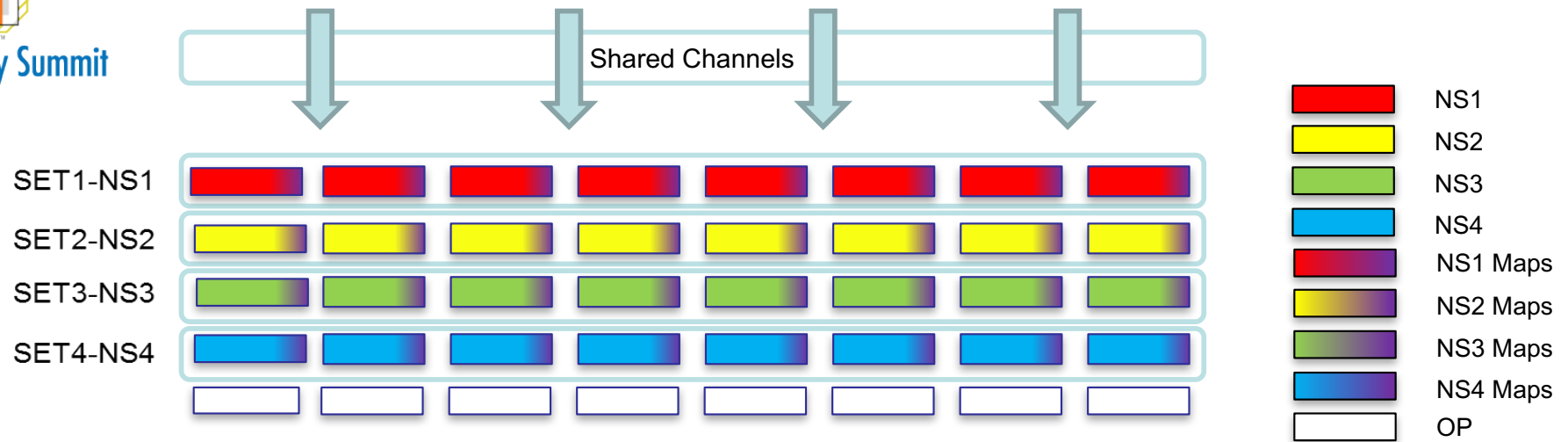
Pros

- Few dedicated channels for each Set
- No conflict in channel or Die
- Every Set runs in parallel resources

Cons

- 1/4 BW possible per Set (limited by number of dedicated channels)

IOD: Horizontal Sets (HS)



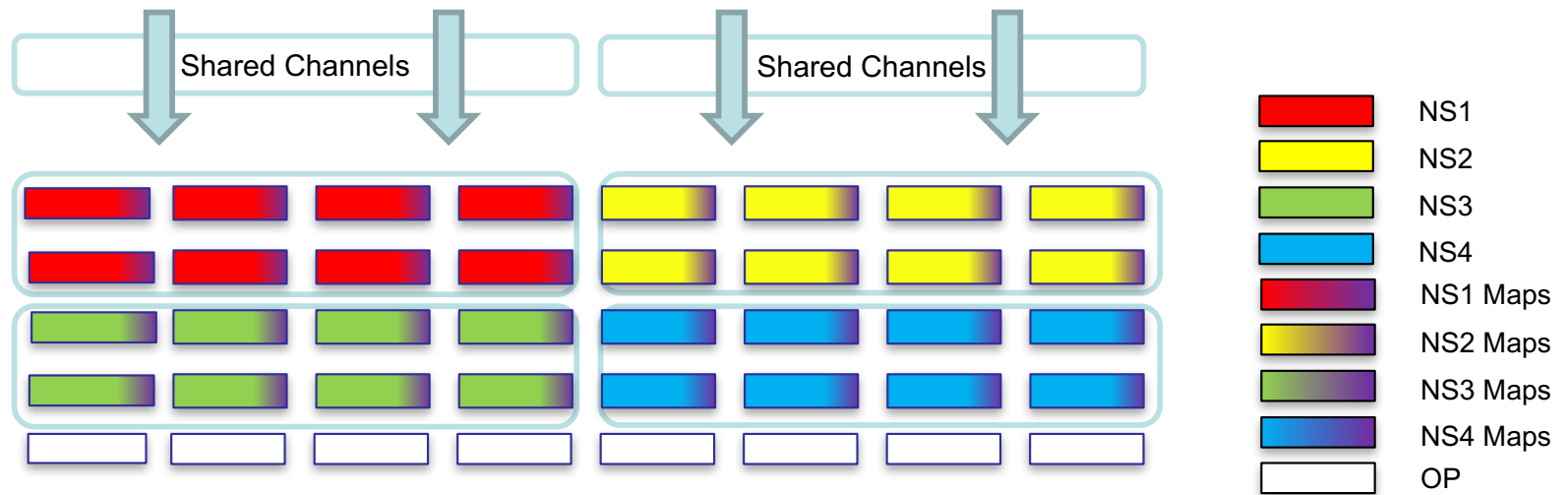
Pros

- All channels are utilized by all Sets
- Full BW possible with single active Set (or shared BW with all Sets)
- No Conflict in Die

Cons

- No Dedicated channels or queues for Set

IOD: Mixed Sets (MS)



Pros

- Few dedicated channels for group of Sets
- No conflict in Die
- Groups of Sets runs in parallel resources

Cons

- 1/2 of the BW possible per Set (limited by number of dedicated/shared channels)



Flash Memory Summit

Areas Not Covered in Proto Work

- R and W path are HW assisted
 - No complete separation between them
 - Conflicts do happen
- Cache memory is shared across VS
 - Conflicts happen but limited by test nature
- R and W have same priority
 - No priority to any specific IO
- No Erase Suspend



Test Results and Analysis



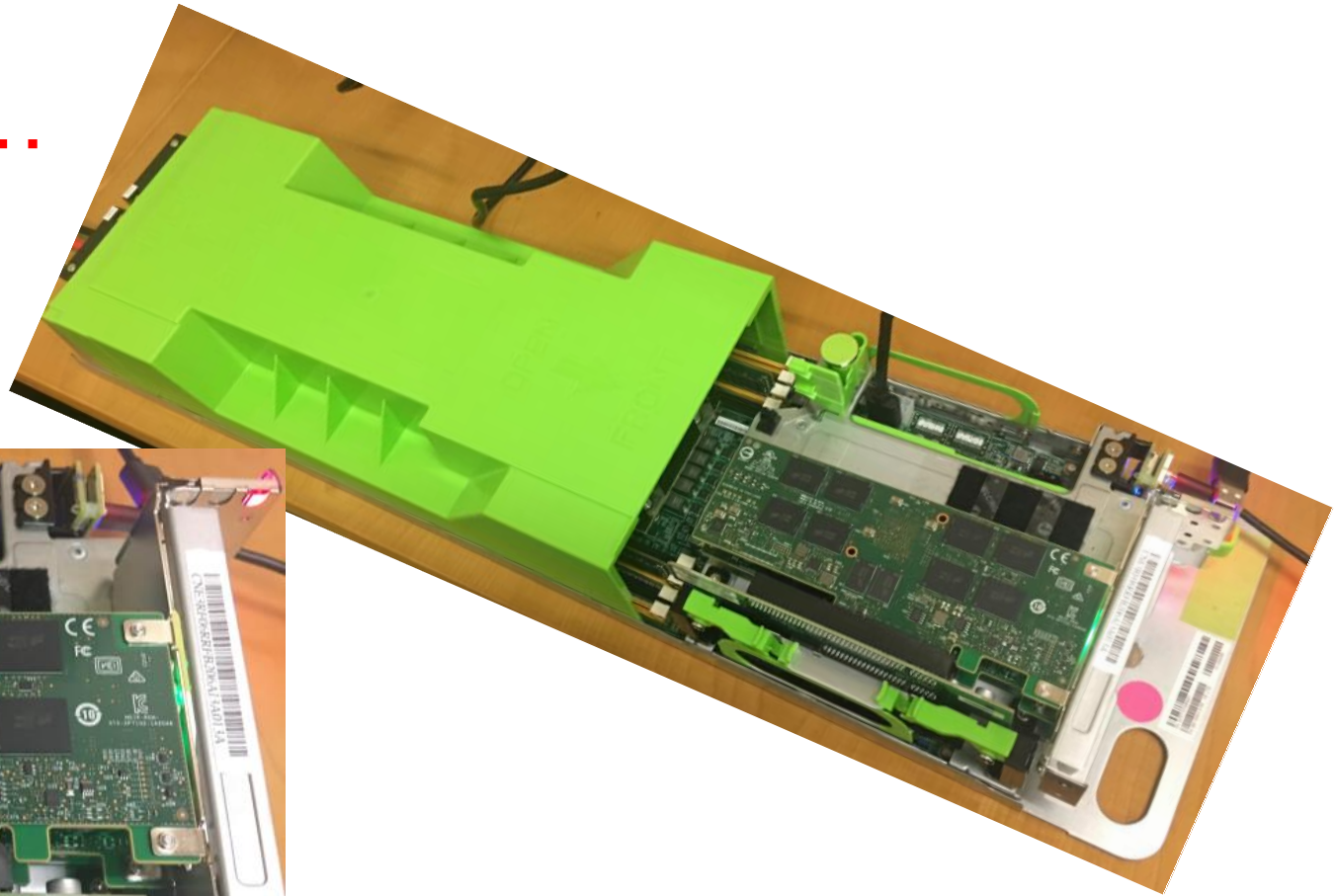
Configuration Under Test

- Facebook OCP Leopard System
 - 2-socket Intel[®] Xeon[®] v4 CPU
 - 32GB DRAM
- STX Nytro XP7102 (2TB eMLC),
 - 128 dies on 8 channels, ONFi @400 MT/s
- Tested with Baseline (No IOD), VSets and Hsets
- Writes = 32KB SW @QD4
- Reads = 4KB RR @QD8



Flash Memory Summit

PICS...



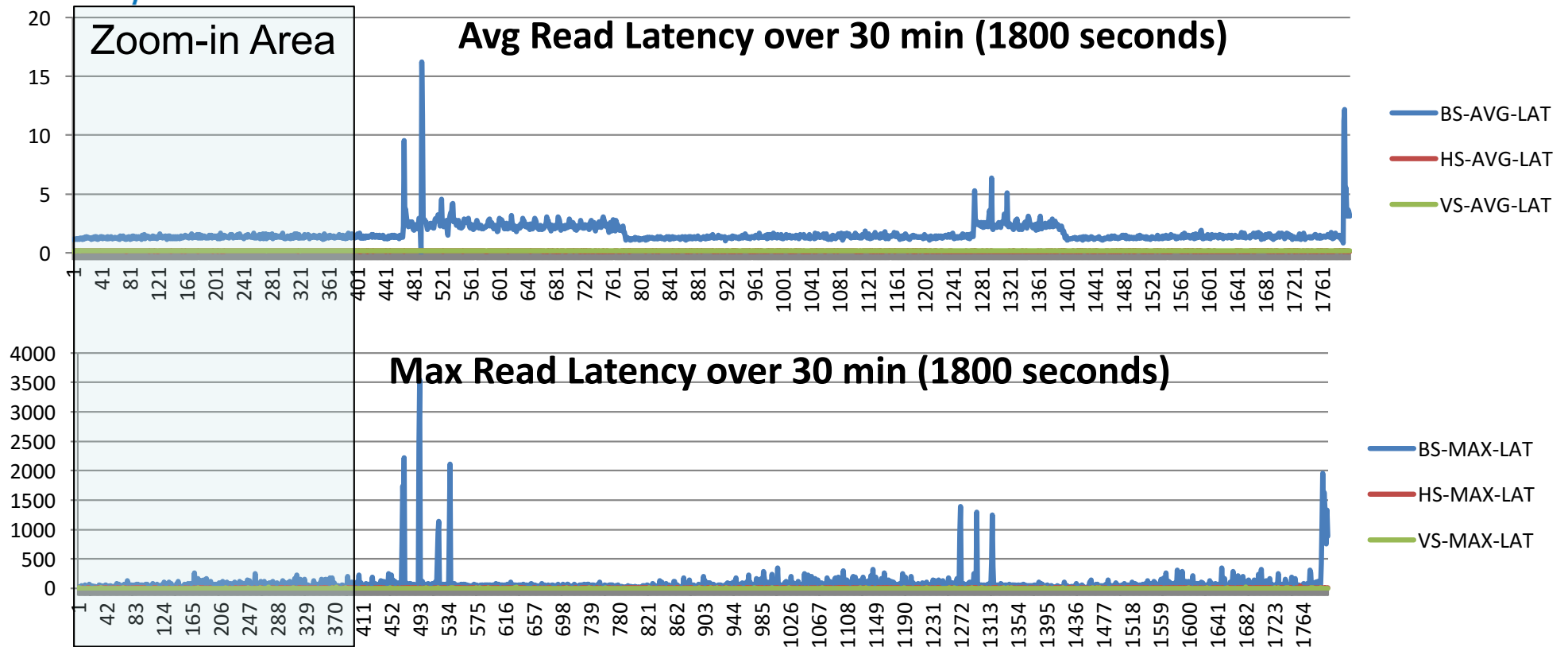


“The Noisy Neighbor”

- One Set under heavy writes, the other under heavy reads
 - Write is 32K Seq Write with QD=4
 - Read is 4K Random Read with QD=8
- Expected Outcome: minimal to no impact on Reads from the “noisy” Writes

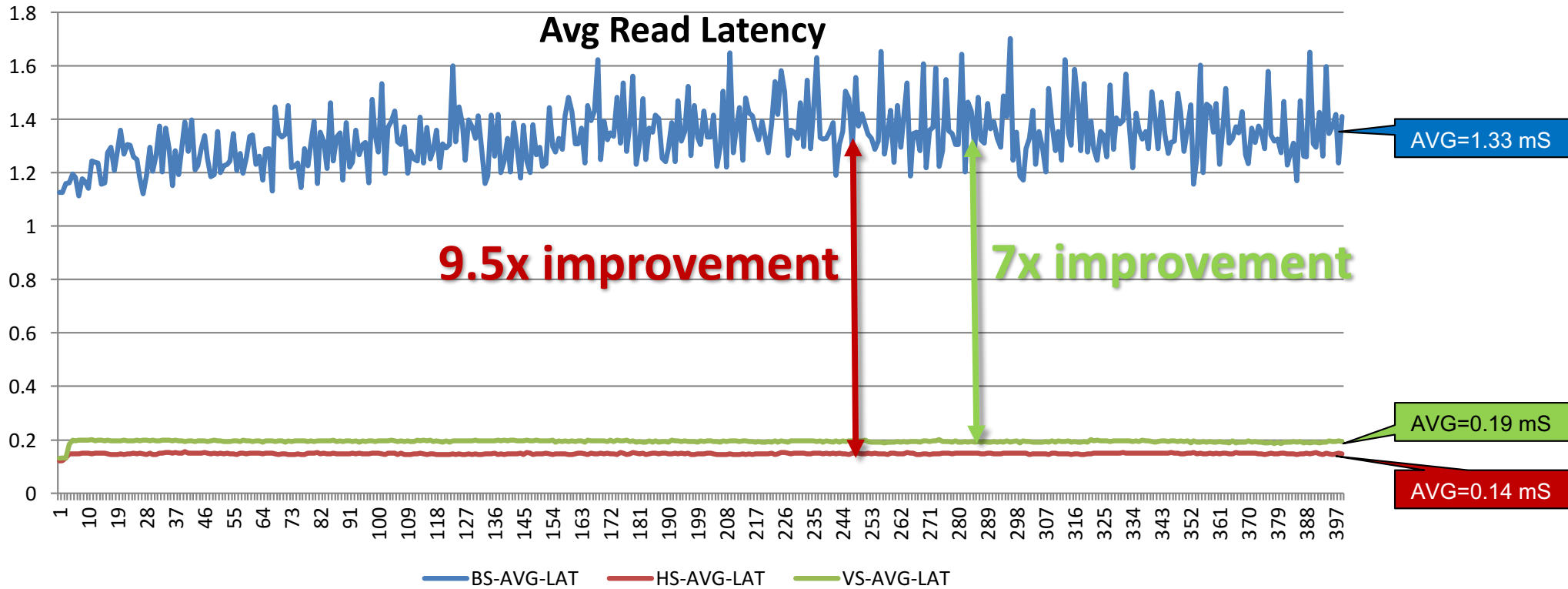


WR on NS1 + RD on NS2



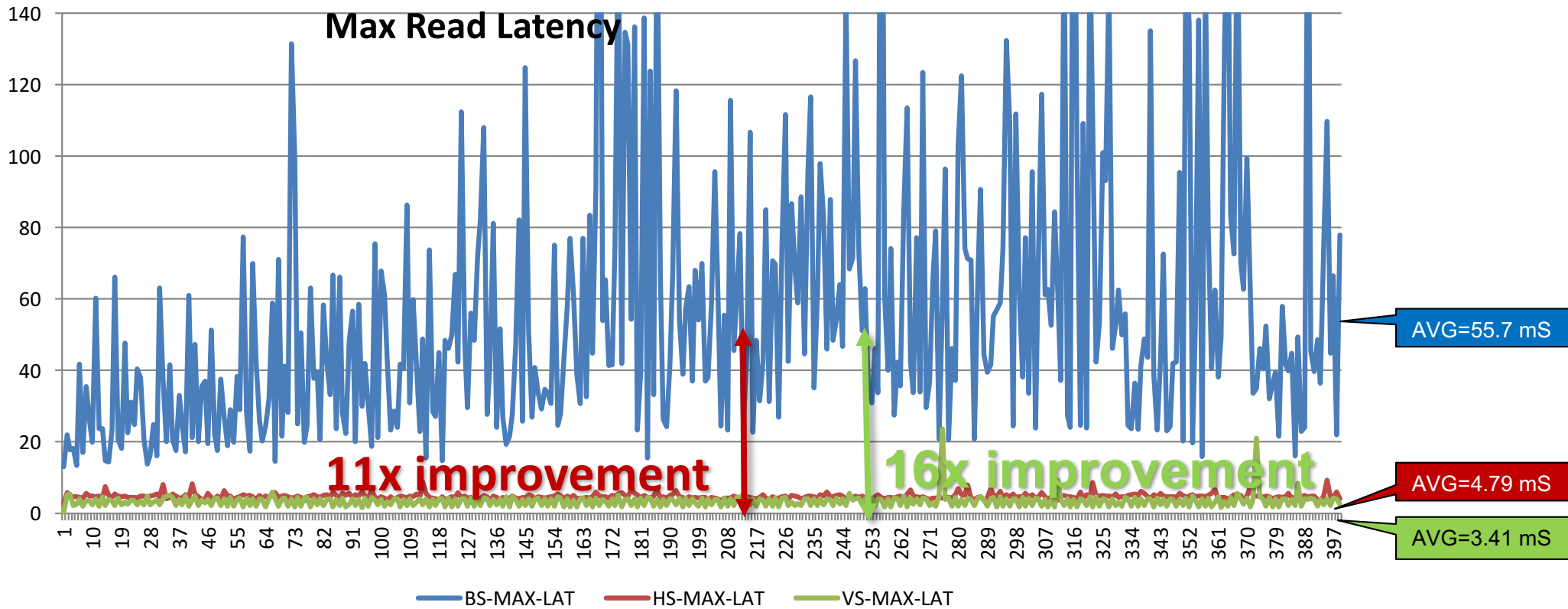


WR on NS1 + RD on NS2

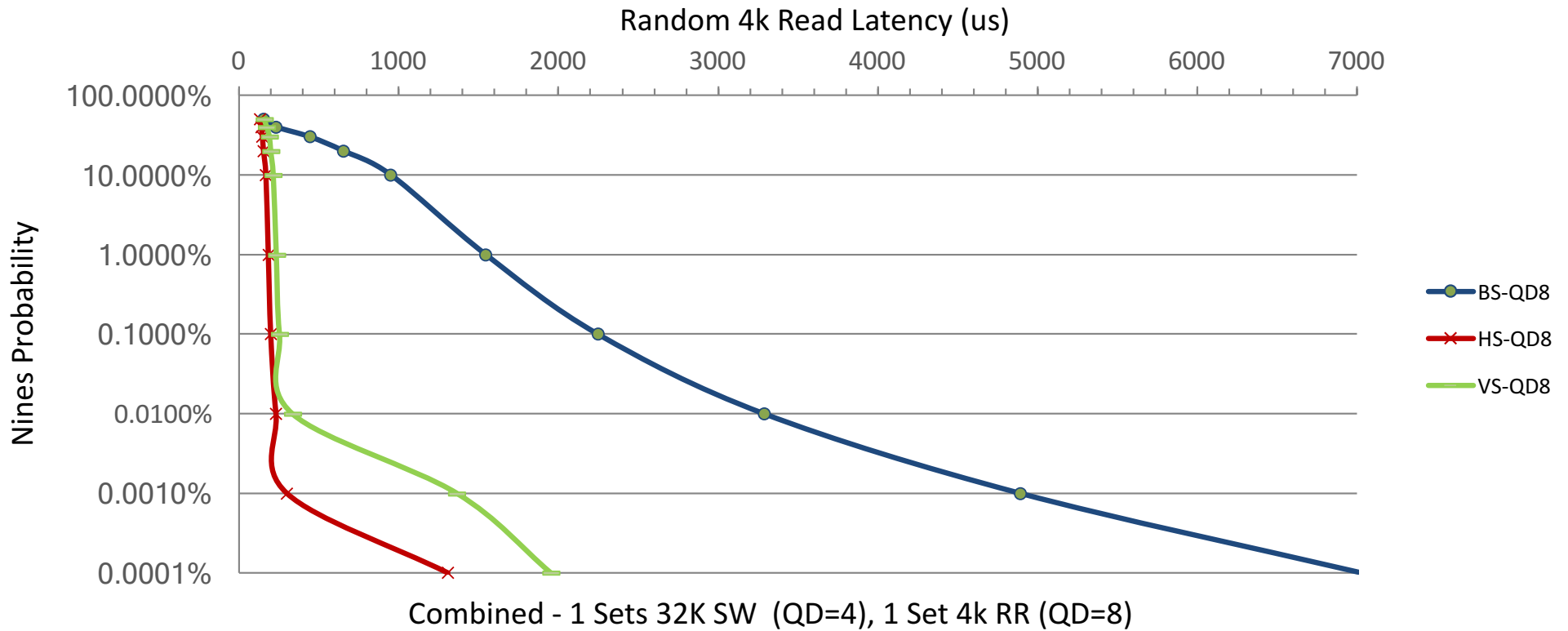




WR on NS1 + RD on NS2



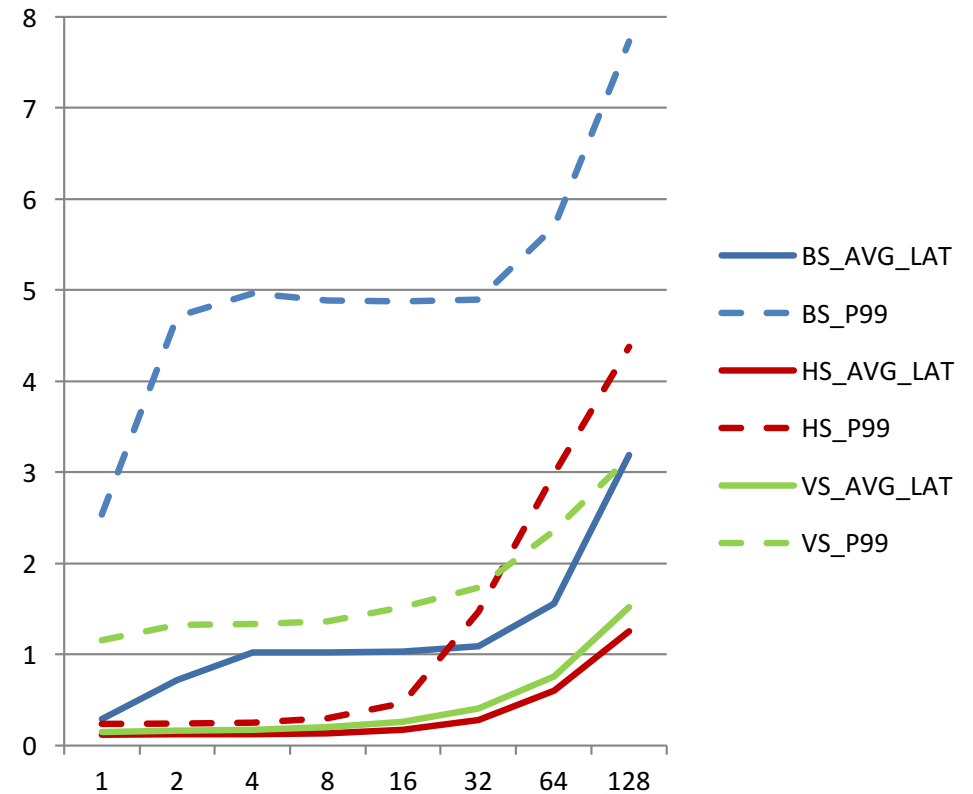
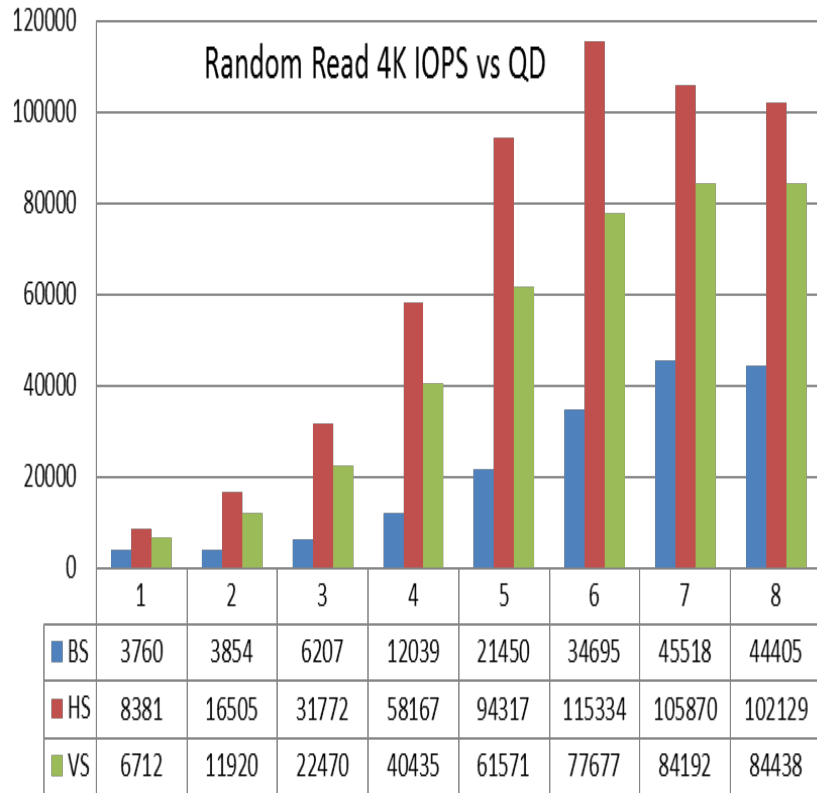
WR on NS1 + RD on NS2





Flash Memory Summit

WR on NS1 + RD on NS2





Conclusions

Conclusions

- Average values: IOD is behaving largely as expected
- Max values: there is more under the hood to address but even proto show huge advantages
- IOD is all about die separation but Controller, DRAM, ONFi, and other resources play key roles under specific conditions
- Most of all... improvements are so high that there may be other usage models outside Hyperscalers



Flash Memory Summit

facebook.



Visit Seagate Booth #505

Learn about Seagate's portfolio of SSDs,
flash solutions and system level products for
every segment.

www.seagate.com/Nytro