# Optimizing storage stacks for AI
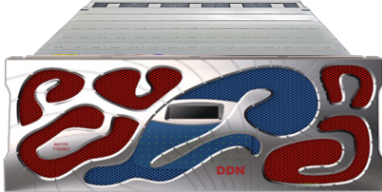
Spectrum Scale CIUK UG

December, 2018

Sven Oehme – Chief Research Officer DDN

# DDN SFA | ALL-FLASH AND HYBRID BLOCK STORAGE PLATFORMS

| 200NV | 400NV | 7990 | 14KX | 18K |
|---|---|---|---|---|
| 23GB/s | 42GB/s | 23GB/s | 60GB/s | 92GB/s |
| 1M IOP/s | 3M IOP/s | 1M IOP/s | 4M IOP/s | 3.2M IOP/s |
| 24 NVME Slots | 24 NVME Slots | | 48 NVMe Slots | 48 NVMe Slots |
| | | Up to 450 SSD/HDD | Up to 1872 SSD/HDD | Up to 1872 SSD/HDD |
| EDR IB (4), OPA (2) | EDR IB (8), OPA (4) | EDR IB (4), OPA (2) | EDR IB (12\|8) | EDR IB (16), OPA (8) |
| FC32 (8), FC (8) | | FC16 (8) | OPA (4), FC16 (24) | |
| NEW 2018 | COMING 2019 | NEW 2018 | | COMING 2019 |

# DDN | GRIDScaler

Massively Scalable NAS & Parallel File Storage Appliance

| | GS200NV | GS400NV | GS7990 | GS14KX | GS18K |
|---|---|---|---|---|---|
| GRIDScaler v4 | ✓ | ✓* | ✓ | ✓ | ✓* |
| v4 upgrade to v5 | ✓ | ✓* | ✓ | ✓ | ✓* |
| GRIDScaler v5 | ✓ | ✓* | ✓ | ✓ | ✓* |

▶ Easy to deploy, All-in-One Appliance for All Flash Array with HDD, archive and cloud tiering options

▶ Scale-out building blocks architecture
- Configurations scale from <100 TB to PBs of storage and 10s of TBs/sec of performance

▶ Flash Centric Architecture - custom embedded fabric delivers optimal SSD performance

▶ Feature-Rich, Enterprise Grade Quality and High Availability with no single point of failure

▶ Simple, Intuitive but Powerful DDN Insight monitoring solution

# Optimizations for GRIDScaler V5

# Optimizations for GRIDScaler V5

▶ Updated device drivers, OS and Scale tuning parameters and SFA multi-queue LUN support

▶ Embedded systems can now achieve up to 1.2 Million random 4k read IOPS

▶ External SFA14KX NSD Server performance went from 1.25 Million to 2.96 Million*

▶ This enhancements were used to produce the SpecSFS 2014 record publications

* test was using external NSD Server. all numbers are measured from network attached clients with GPFSPERF using one 100 GB file per client during random 4k reads using O_DIRECT

# Platform Optimizations help significantly
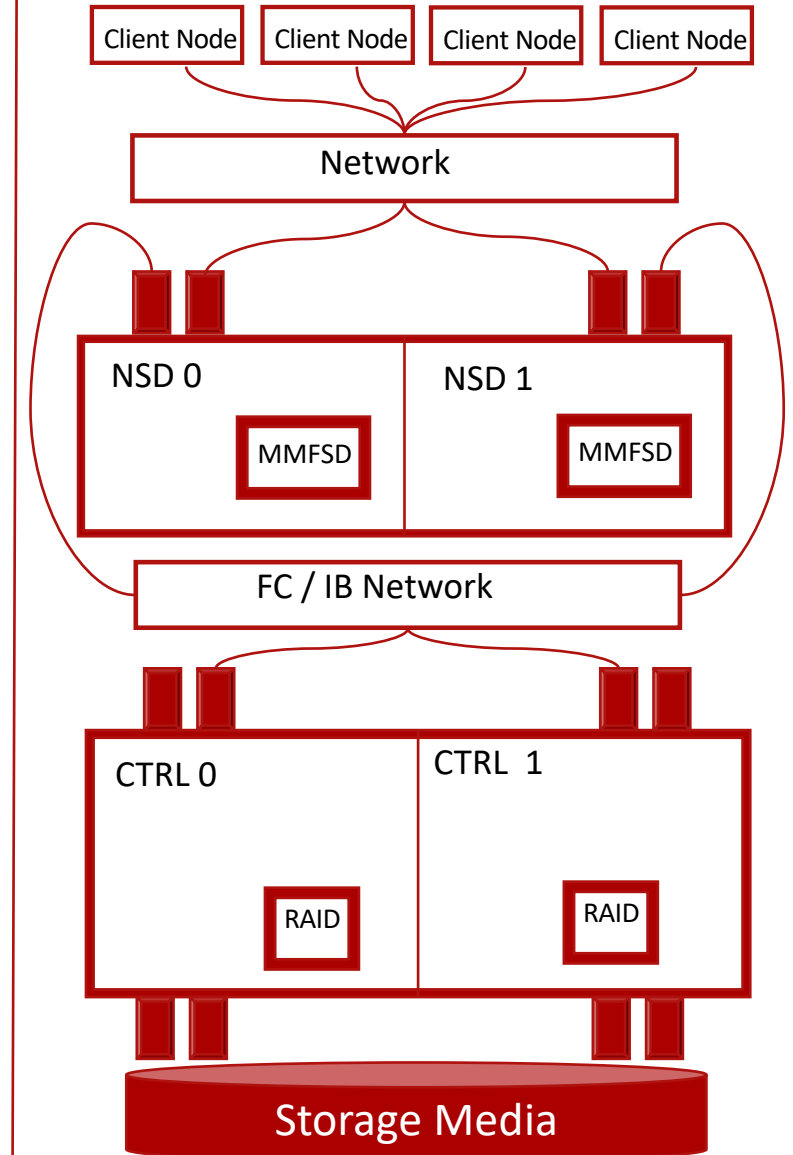
# ES200NV | LOW LATENCY DESIGNED-IN

| CLIENT | → | HCA/NIC | → | SWITCH | → | HCA/NIC | → | SERVER | → | FILESYSTEM | → | HBA | → | SAN SWITCH | → | HBA | → | STORAGE |

## IO PATHS

**TRADITIONAL**

**Components Simplified with SFA™ Embedded Appliances**

## SFA EMBEDDED FILESERVICE

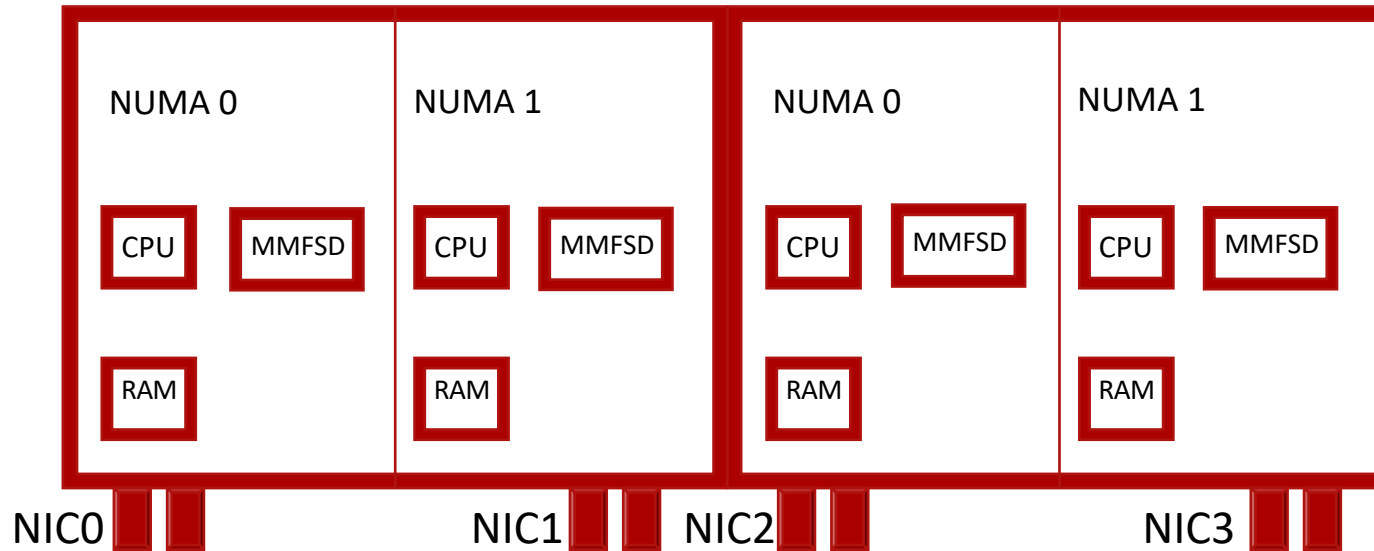| CLIENT | → | HCA/NIC | → | SWITCH | → |

SFA400NV

## SCALER™
APPLIANCES

# The fastest network hop is the one you can avoid

# Why all this work, what's to gain ?

► Remote NUMA region HW access in SW is one of the biggest issue to achieve HW capable performance targets

► even just a 2 NUMA Zone system (e.g. modern Intel 2 socket system) has significant overhead as without optimization on the SW, 50% of the access is remote, as larger the number of NUMA nodes as more overhead , each IBM power or modern AMD CPU has 2 NUMA nodes. So a 2 socket Power 8 system has 4 NUMA zones and a 75% chance your data is on the wrong side.

► Databases developers have spend years to optimize their SW stack to be NUMA aware, storage stacks are trying to catch up. On databases tests have show between 2-4x improvements with proper memory placement, for Storage the benefit can be even greater as it typically interacts with HW beyond memory that is NUMA dependent (e.g. HBA's or HCA's)

► Remote HW access significant increases latency and causes very unpredictable performance

► Linear scaling with increased core counts gets eliminated by contention on interconnects or lock overhead requiring synchronization between NUMA regions

# SFA NUMA awareness



The system is perfectly balanced across numa nodes, which allows affinitizing of mmfsd threads to memory, core and network for lowest latency and consistent scaling

# DIO Random 4k writes into a 100GB files

```
/usr/lpp/mmfs/samples/perf/gpfsperf write rand /target/sven-100g
  recSize 4K nBytes 100G fileSize 100G
  nProcesses 1 nThreadsPerProcess 1
  file cache flushed before test
  using direct I/O
  offsets accessed will cycle through the same file segment
  not using shared memory buffer
  not releasing byte-range token after open
  no fsync at end of test
    Data rate was 34659.88 Kbytes/sec, Op Rate was 8461.89 Ops/sec, Avg Latency
was 0.118 milliseconds, thread utilization 1.000, bytesTransferred 1039802368
```

# DIO Random 4k reads from a 100GB files (exceeds all cache by 4x)

```
/usr/lpp/mmfs/samples/perf/gpfsperf read rand /target/sven-100g
  recSize 4K nBytes 100G fileSize 100G
  nProcesses 1 nThreadsPerProcess 1
  file cache flushed before test
  using direct I/O
  offsets accessed will cycle through the same file segment
  not using shared memory buffer
  not releasing byte-range token after open
    Data rate was 21763.50 Kbytes/sec, Op Rate was 5313.36 Ops/sec, Avg Latency
was 0.188 milliseconds, thread utilization 1.000, bytesTransferred 652910592
```

# DIO Random 4k reads from a 100GB files (exceeds all cache by 4x) ETH

```
/work/oehmes/bin/gpfsperf read rand -r 4k -n 100g -th 1 -dio -millis 5000
/ai200g/test.sven
/work/oehmes/bin/gpfsperf read rand /ai200g/test.sven
  recSize 4K nBytes 100G fileSize 100G
  nProcesses 1 nThreadsPerProcess 1
  file cache flushed before test
  using direct I/O
  offsets accessed will cycle through the same file segment
  not using shared memory buffer
  not releasing byte-range token after open
    Data rate was 14888.68 Kbytes/sec, Op Rate was 3634.93 Ops/sec, Avg Latency
was 0.275 milliseconds, thread utilization 1.000, bytesTransferred 74448896
```
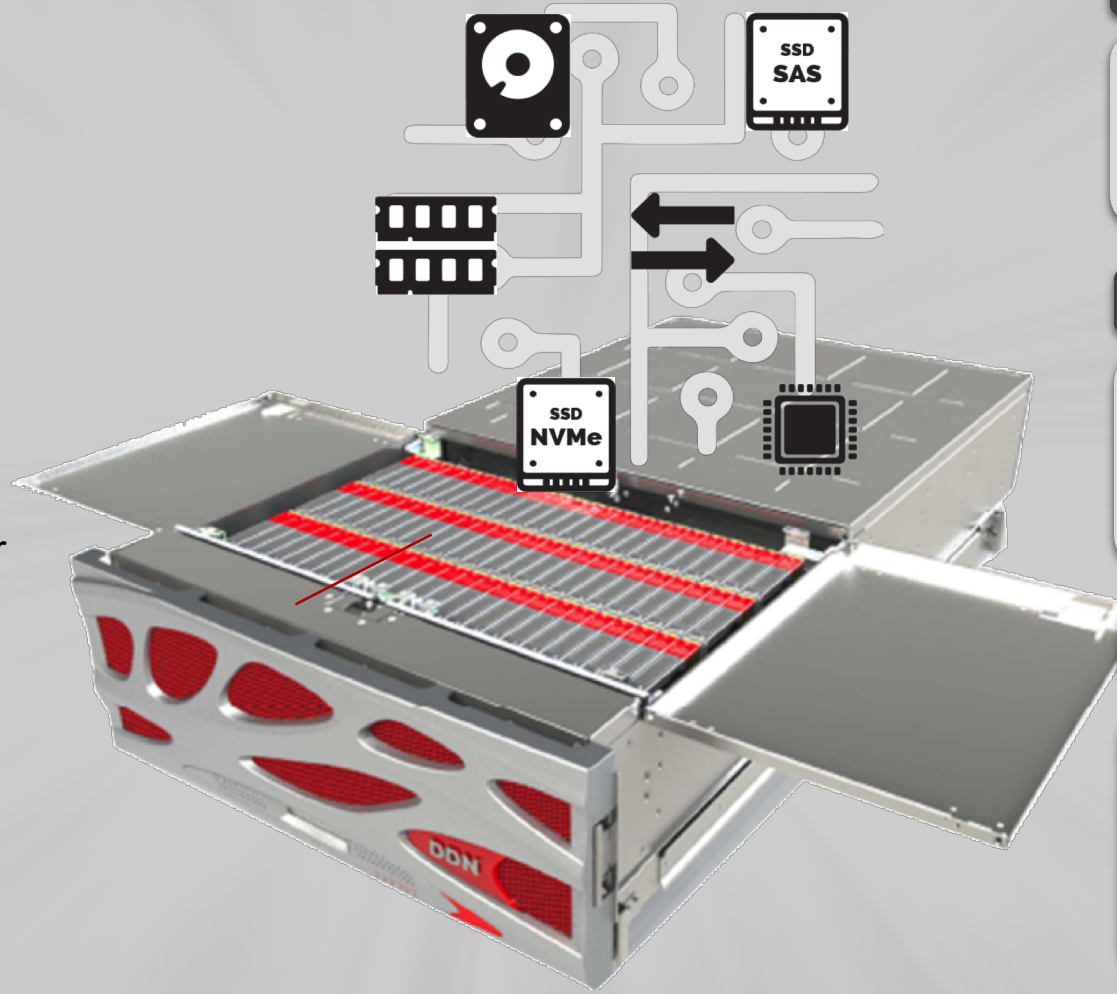
# World record SpecSFS 2014 with GRIDScaler*

*world record has been broken with an 8 Storage System setup - we use ONE !

# DDN SFA14KX
*Fastest, Densest and Simplest at Scale*

**Low Latency, Highly Efficient Architecture**

- All in one integrated design with expansion capability
- Dual Redundant Controllers
- 72 Drive High-Density 2.5" Enclosure with NVMe support for 48 2.5" dual ported NVMe
- Optimized Building Block for BW or IOPs
- Support for up to 20 SS9012 12Gb/s 90 drive Enclosures

## Flexible Connectivity

- ► 10/40/100GbE
- ► IB and OmniPath
- ► 16/32Gb FC

## Industry Leading Performance
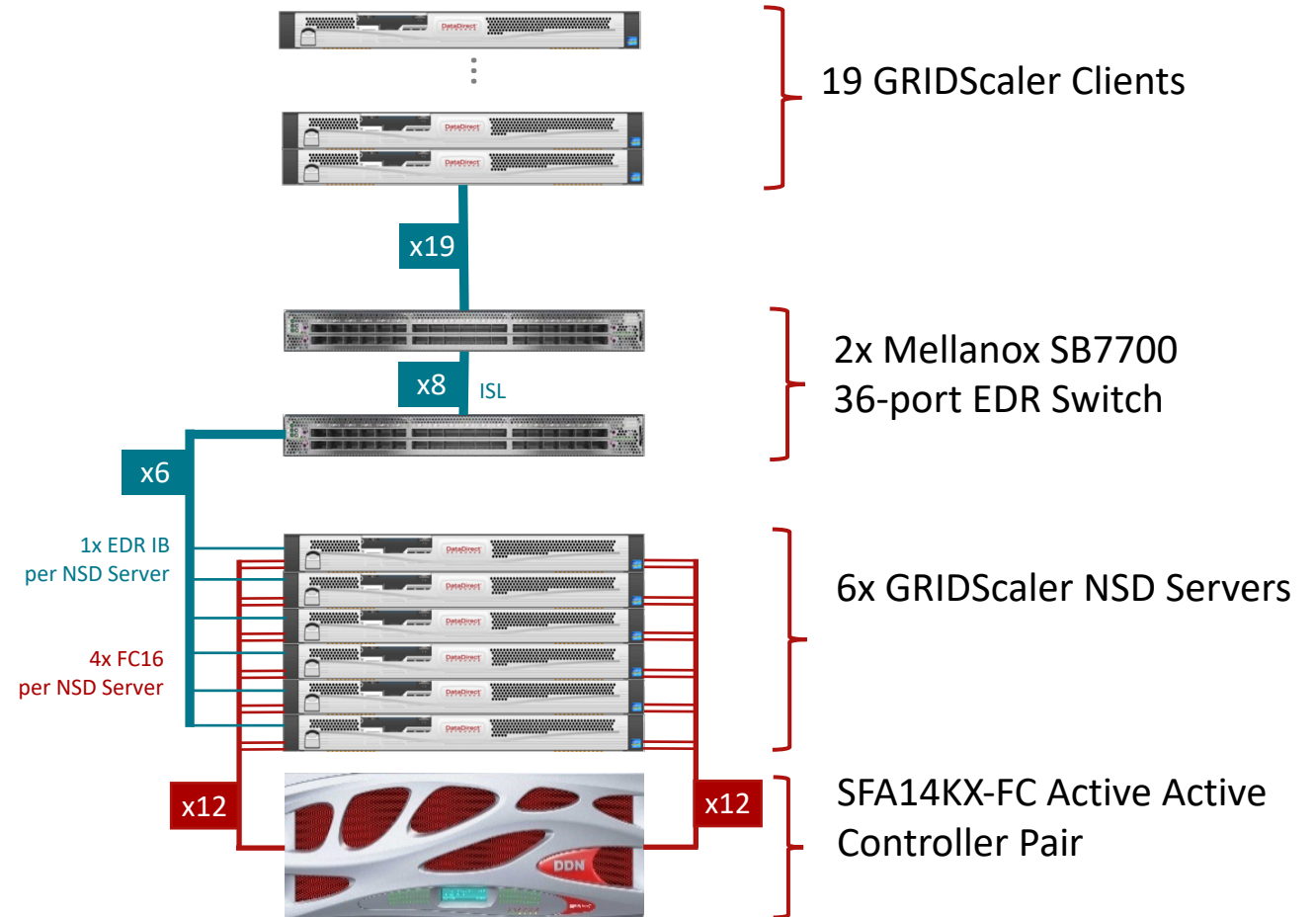
- ► 72 SAS SSD or 48 NVMe
- ► Up to 60 GB/sec throughput
- ► Up to 4 million IOPS

## Best Data Protection

- ► Fully Declustered RAID
- ► Higher Data Availability
- ► Flexible Pool Management
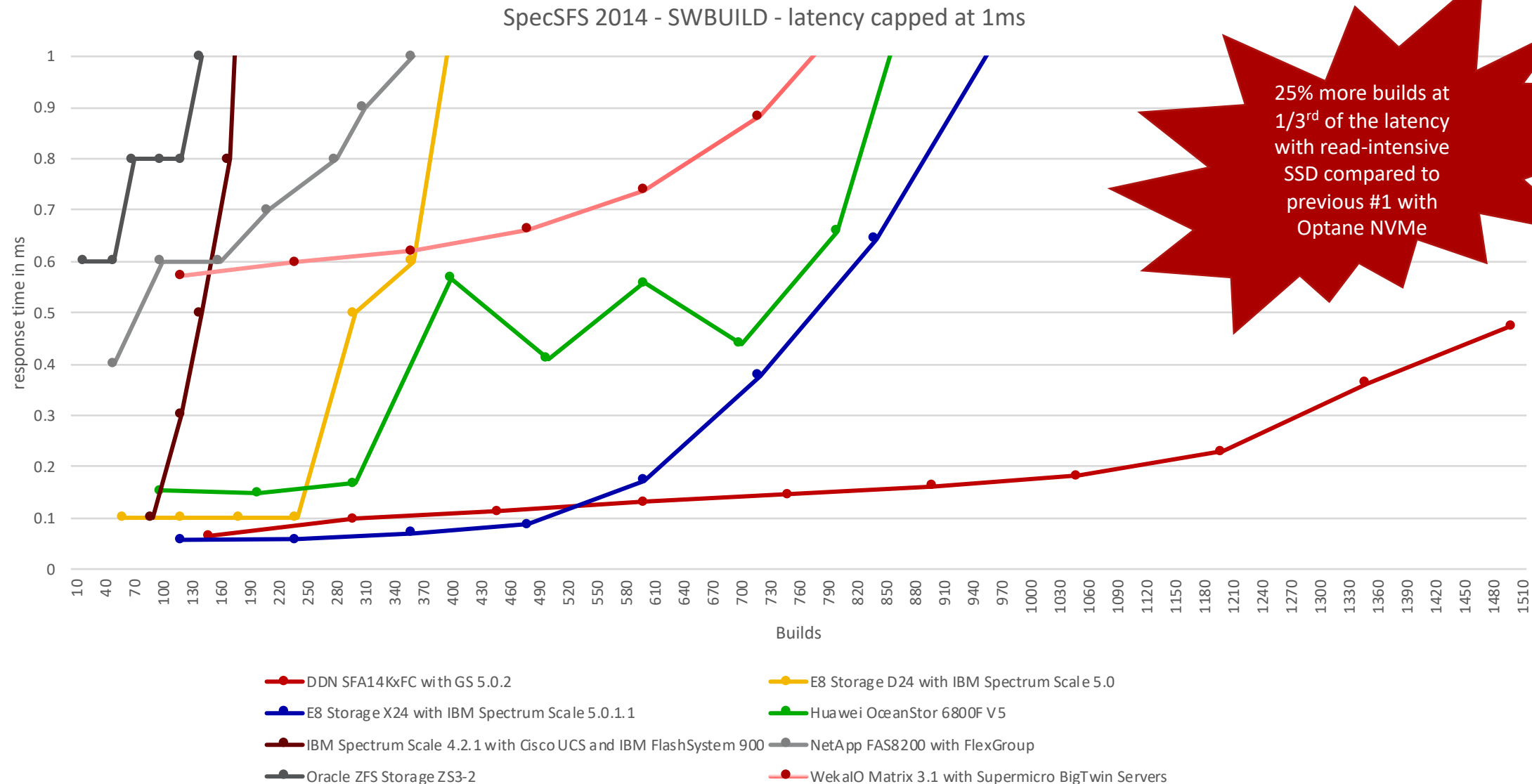- ► Optimized for both Random and Sequential IO

# DDN SFA14KX with GRIDScaler

▶ With the SFA14KX and GRIDScaler parallel filesystem, DDN gains pole position for SPEC SFS

▶ DDN's SFA14KX running SFAOS with Declustered RAID and connecting to 6 GRIDScaler servers Sustains 25% more builds at 1/3rd of the Overall roundtrip latency with read-intensive SSD compared to previous #1 with Optane NVMeof - the next nearest competitor

19 GRIDScaler Clients

x19

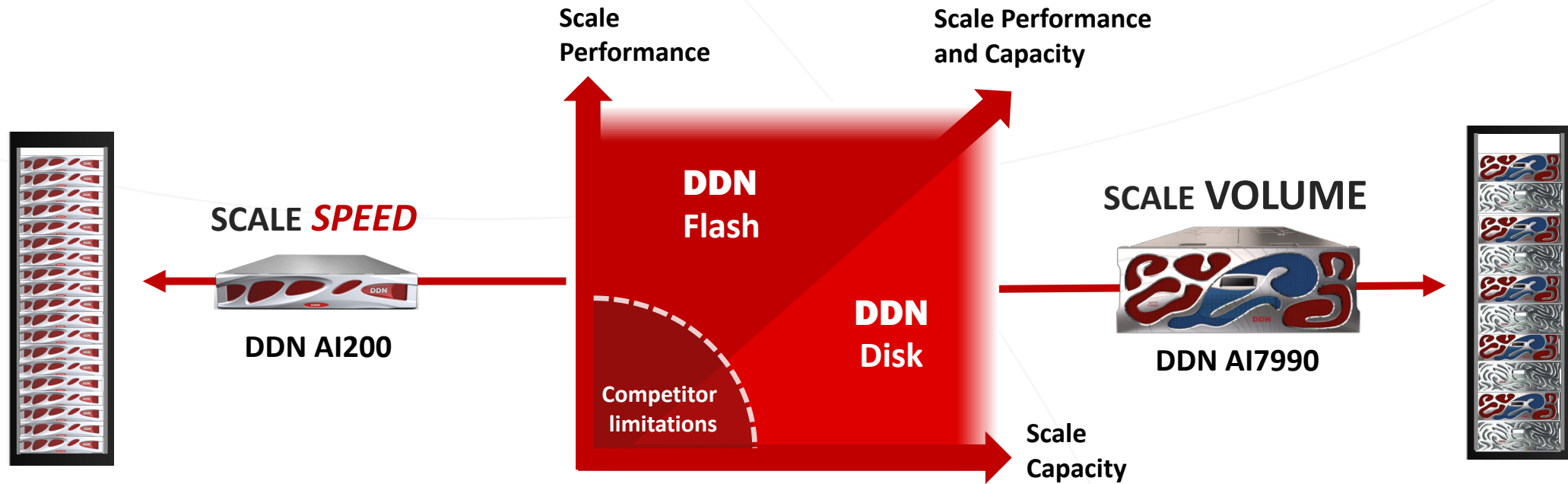2x Mellanox SB7700 36-port EDR Switch

x8    ISL

x6

1x EDR IB per NSD Server

4x FC16 per NSD Server

6x GRIDScaler NSD Servers

x12                    x12

SFA14KX-FC Active Active Controller Pair

**System Benchmarked for SPEC SFS**

# SpecSFS 2014 – SWBUILD compare to other vendors



SpecSFS 2014 - SWBUILD - latency capped at 1ms

response time in ms

Builds

25% more builds at 1/3rd of the latency with read-intensive SSD compared to previous #1 with Optane NVMe

- DDN SFA14KxFC with GS 5.0.2
- E8 Storage X24 with IBM Spectrum Scale 5.0.1.1
- IBM Spectrum Scale 4.2.1 with Cisco UCS and IBM FlashSystem 900
- Oracle ZFS Storage ZS3-2
- E8 Storage D24 with IBM Spectrum Scale 5.0
- Huawei OceanStor 6800F V5
- NetApp FAS8200 with FlexGroup
- WekaIO Matrix 3.1 with Supermicro BigTwin Servers

DDN A³I Solutions: Turnkey, integrated and optimized for NVIDIA DGX-1 and HP Apollo 6500

# SCALE UP, SCALE OUT OR SCALE BOTH



Scale
Performance

Scale Performance
and Capacity

SCALE *SPEED*

**DDN**
Flash

SCALE **VOLUME**

DDN AI200

**DDN**
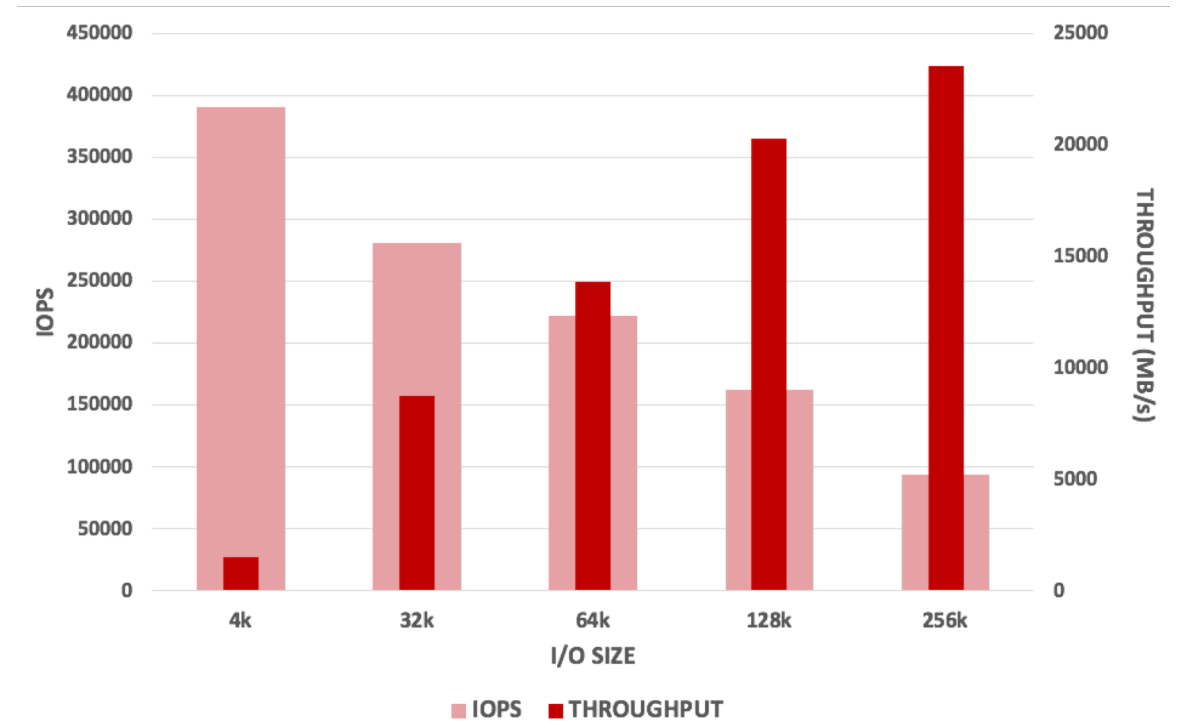Disk

DDN AI7990

Competitor
limitations

Scale
Capacity

# DDN A³I SOLUTIONS TO A SINGLE CONTAINER ON DGX-1

## 23 GB/s and 395K IOPS to a single container*

DDN A³I parallel storage client demonstrates over 23 GB per second and over 395K IOPS to a single container on DGX-1.

Typical deep learning codes perform IO using 128K size for which DDN delivers over 20 GB/s of sustained performance.



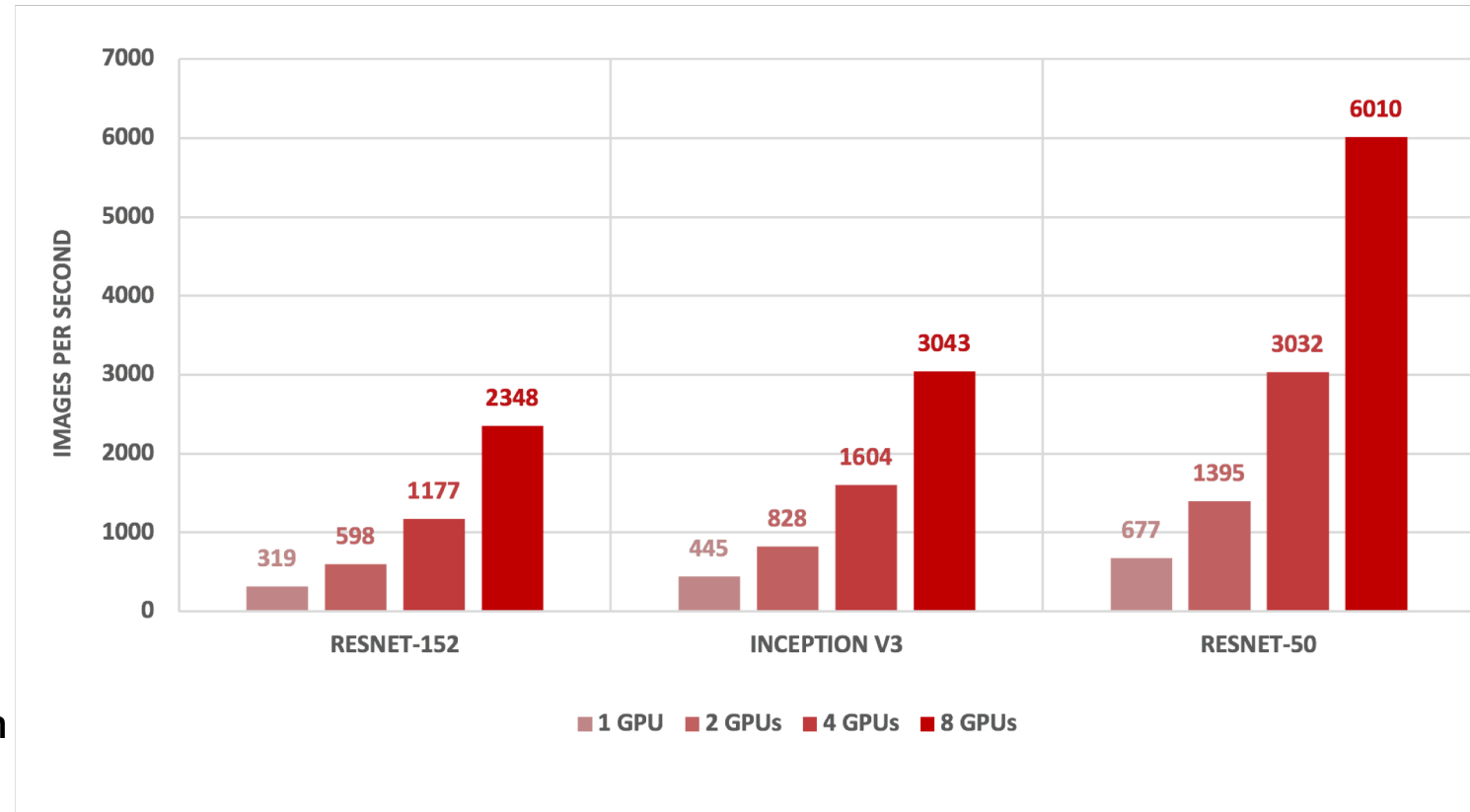*numbers are with a single AI200 and was limited by client side performance of single DGX client

# DDN A³I SOLUTIONS TENSORFLOW TRAINING PERFORMANCE

## Fast, Consistent, Linear AI and DL Performance

DL Training application performance scales linearly using multiple GPUs on DGX-1 with DDN parallel storage.
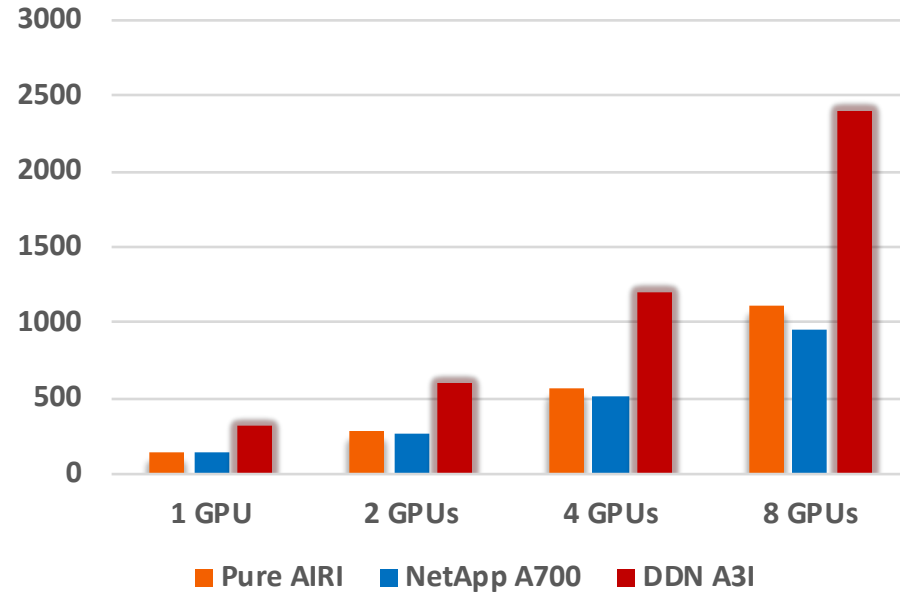
Parallel storage performance and shared architecture magnify end-to-end DL workflow acceleration.

Extensive application interoperability and performance testing has been engaged by DDN in close collaboration with NVIDIA and customers.
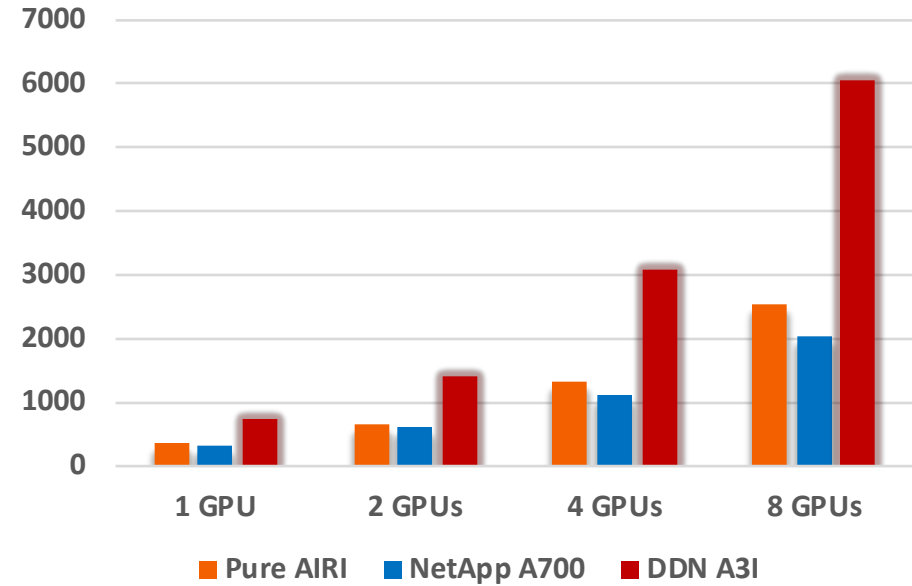
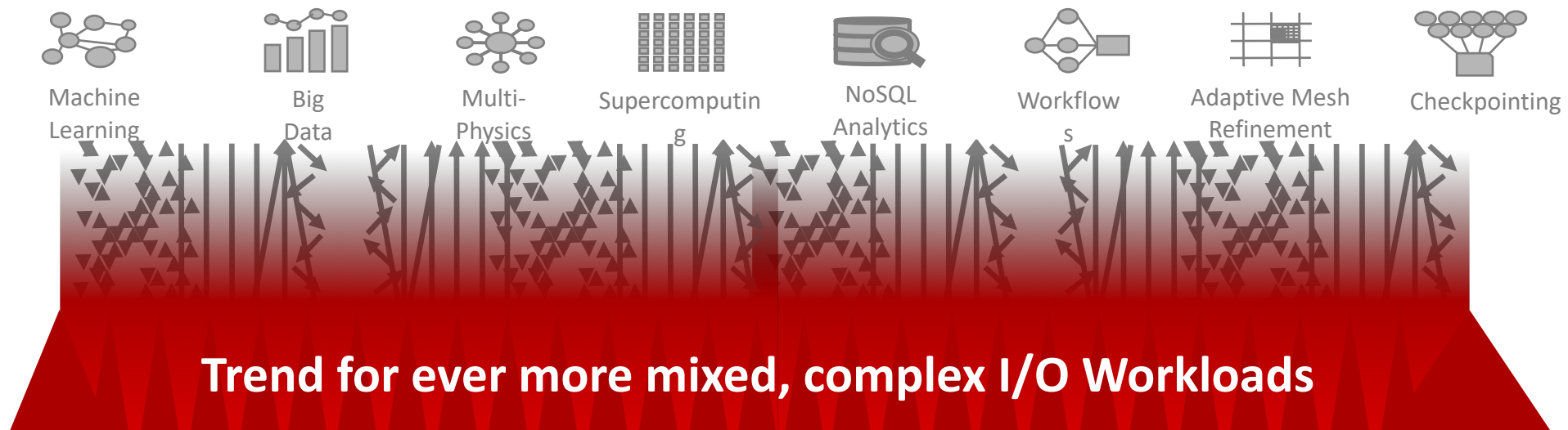# DDN A³I SOLUTIONS LEADS PERFORMANCE FOR AI AND DL



## ResNet-152

## ResNet-50

Pure AIRI    NetApp A700    DDN A3I

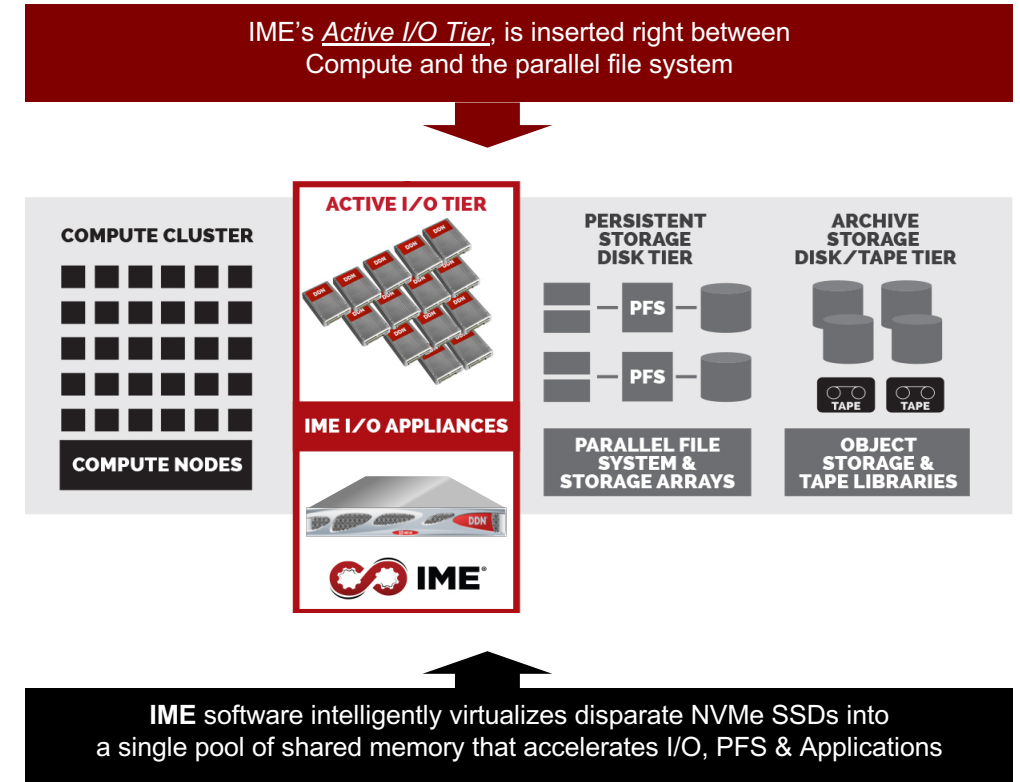"In the Resnet-152 and Resnet-50 tests, the AI200 tested faster than competing Pure, NetApp and Dell EMC systems."

# Challenges in I/O Performance and Behavior

► Newer applications need to operate on byte addressable data

► Significant shift from sequential to random I/O

► Multifold increase of metadata to data ratio

► Average data sizes are less homogeneous and are now fractions or multiples of previous workloads. gap between small and large data seems to wide (bytes on one end , GB's on the other end of the spectrum)

► Interactive, outcome and event driven analytics are driven by latency rather than bandwidth

Machine Learning    Big Data    Multi-Physics    Supercomputing    NoSQL Analytics    Workflows    Adaptive Mesh Refinement    Checkpointing

**Trend for ever more mixed, complex I/O Workloads**

# WHAT IS IME?

► Scale-Out Flash Cache Layer using NVMe SSDs inserted between compute cluster and Parallel File System (PFS)

- IME is configured as CLUSTER with multiple NVMe servers
- All compute nodes can access cache data on IME

► Accelerates difficult IO patterns: small/random/shared file/high concurrency due to thin SW IO management layer

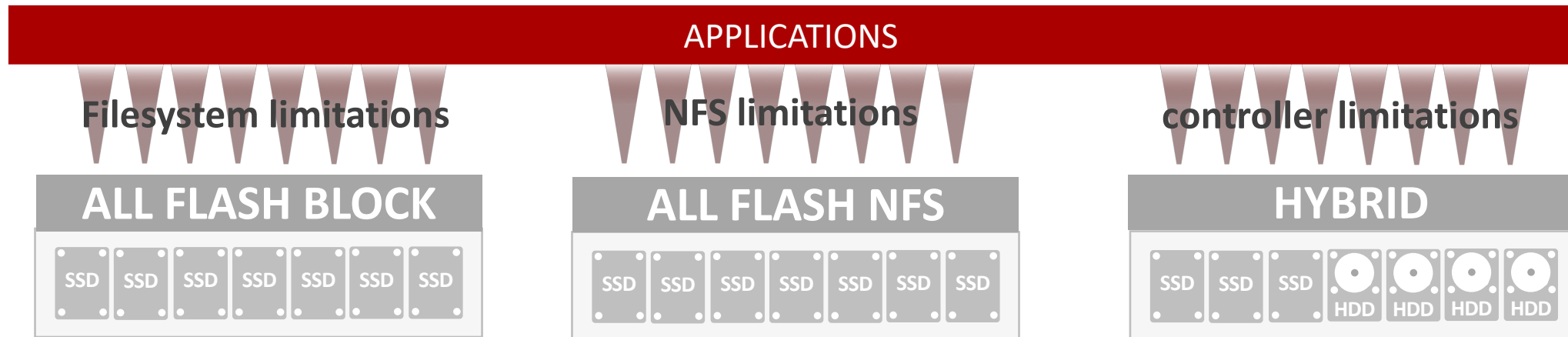► configured as scale-out massive cache layer with huge IO bandwidth and IOPs



IME's *Active I/O Tier*, is inserted right between Compute and the parallel file system

COMPUTE CLUSTER

COMPUTE NODES

ACTIVE I/O TIER

IME I/O APPLIANCES

IME

PERSISTENT STORAGE DISK TIER

PFS
PFS

PARALLEL FILE SYSTEM & STORAGE ARRAYS

ARCHIVE STORAGE DISK/TAPE TIER

TAPE    TAPE

OBJECT STORAGE & TAPE LIBRARIES

**IME** software intelligently virtualizes disparate NVMe SSDs into a single pool of shared memory that accelerates I/O, PFS & Applications

# Expansion in Active Data Volumes requires a new economics for fast data at scale

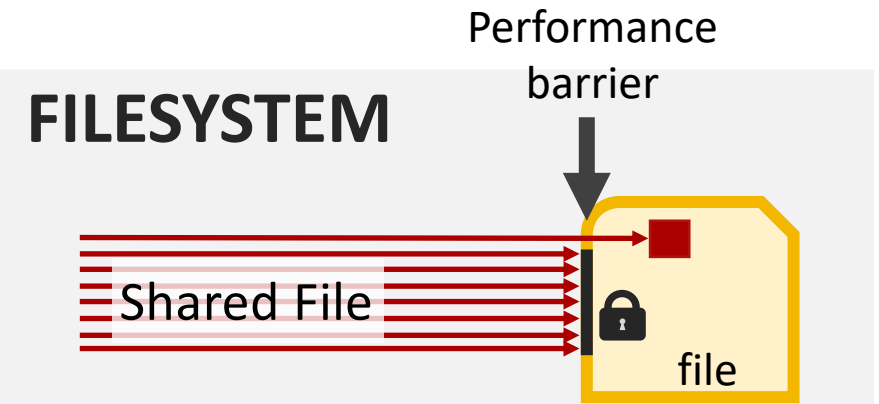**All-Flash block** doesn't solve the problem. **Block IOPs ≠ File IOPs**

**All-Flash NFS** too slow and **too expensive** for real at-scale data problems

Traditional **Hybrid Approach doesn't enable flash at scale** – still limited by the storage controller

APPLICATIONS

**Filesystem limitations**

**NFS limitations**

**controller limitations**

ALL FLASH BLOCK

SSD SSD SSD SSD SSD SSD SSD

ALL FLASH NFS

SSD SSD SSD SSD SSD SSD SSD

HYBRID

SSD SSD SSD HDD HDD HDD HDD

# IME enables new levels of filesystem performance

▶ Parallel File systems can exhibit extremely poor performance for shared file IO due to internal lock management as a result of managing files in large lock units
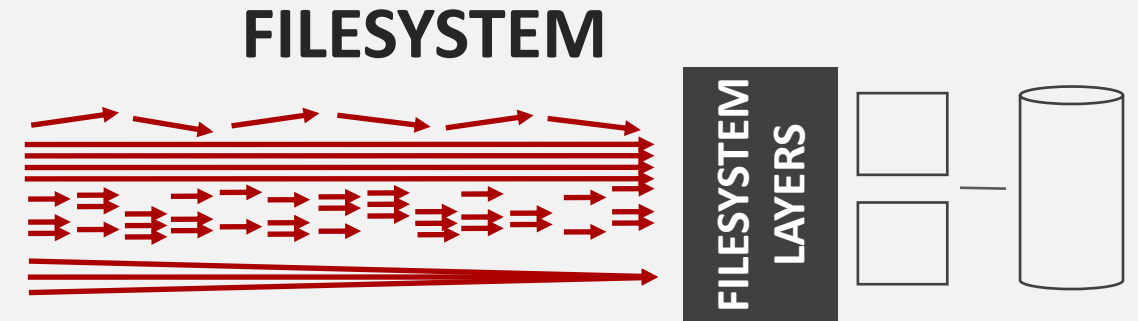
**FILESYSTEM**

Performance barrier

Shared File

file

▶ IME eliminates contention by managing IO fragments directly, and coalescing IO's prior to flushing to the parallel file system
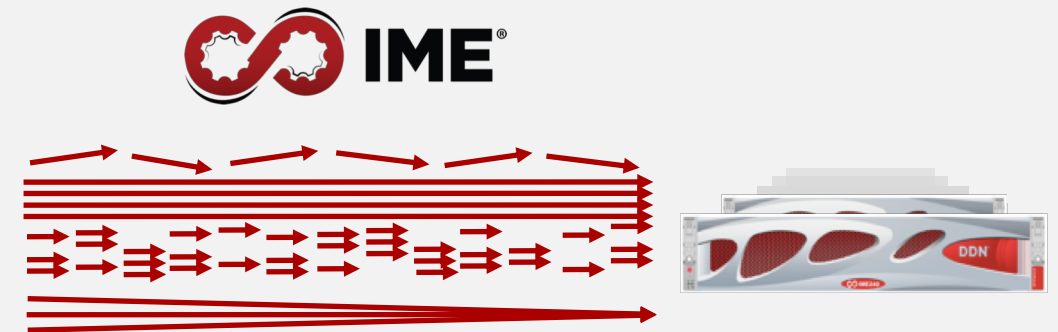
**IME**®

Shared File

file

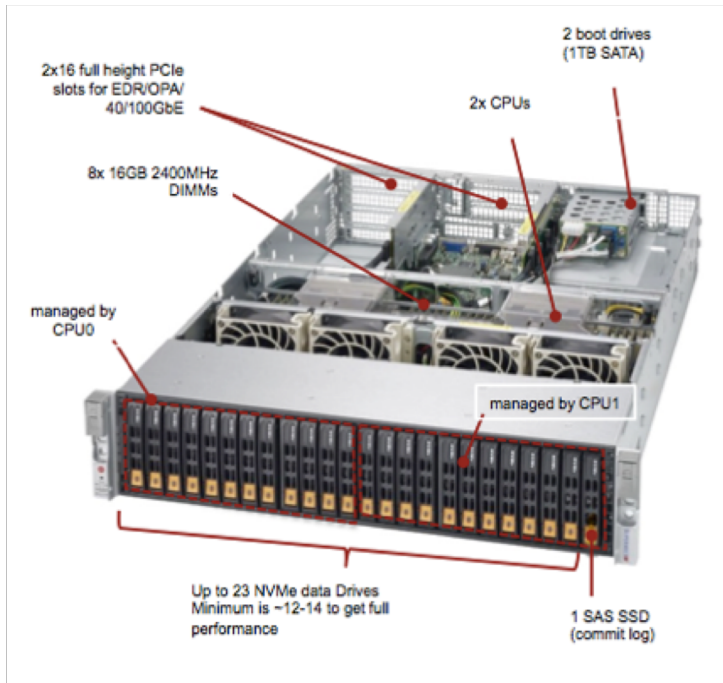# IME enables new levels of filesystem performance

▶ Thick File system SW layers and traditional data layout severely restricts performance for tough workloads

**FILESYSTEM**



▶ IME's lean write anywhere, fully parallel IO completely removes the barriers that prevent your application seeing full performance

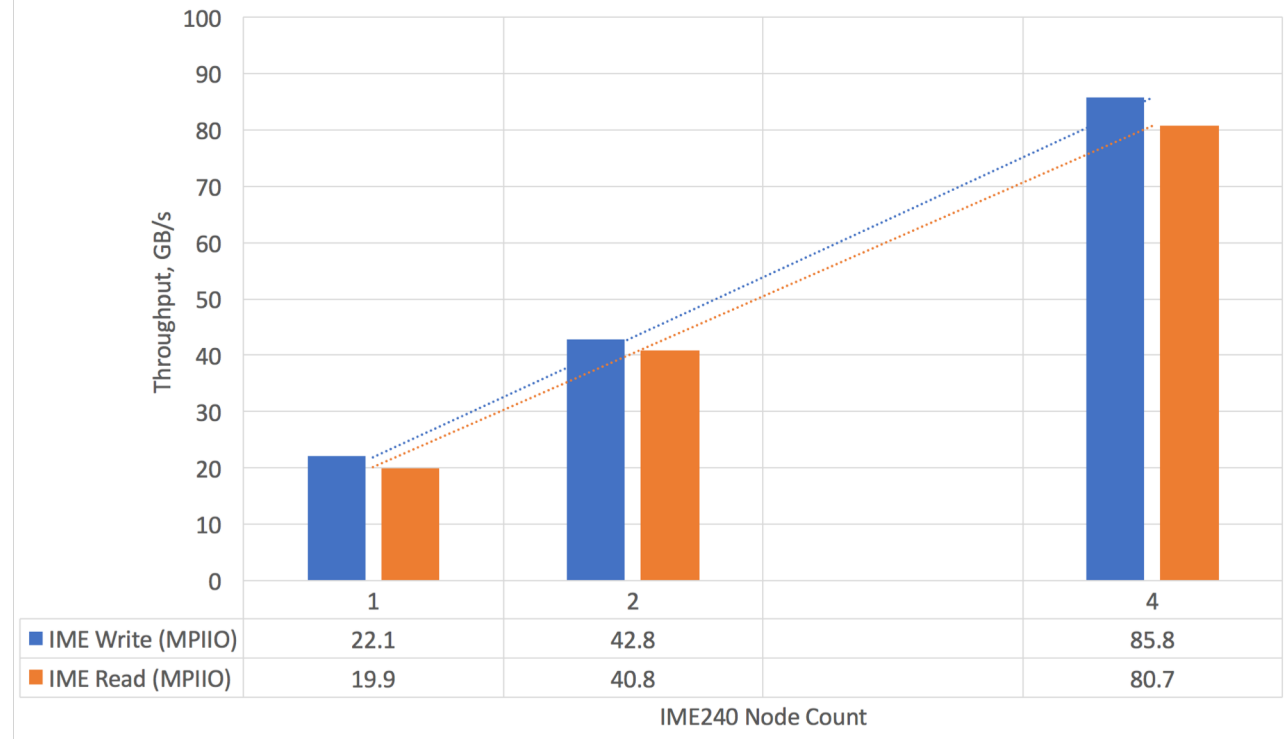# IME Performance Scalability & R/W Parity
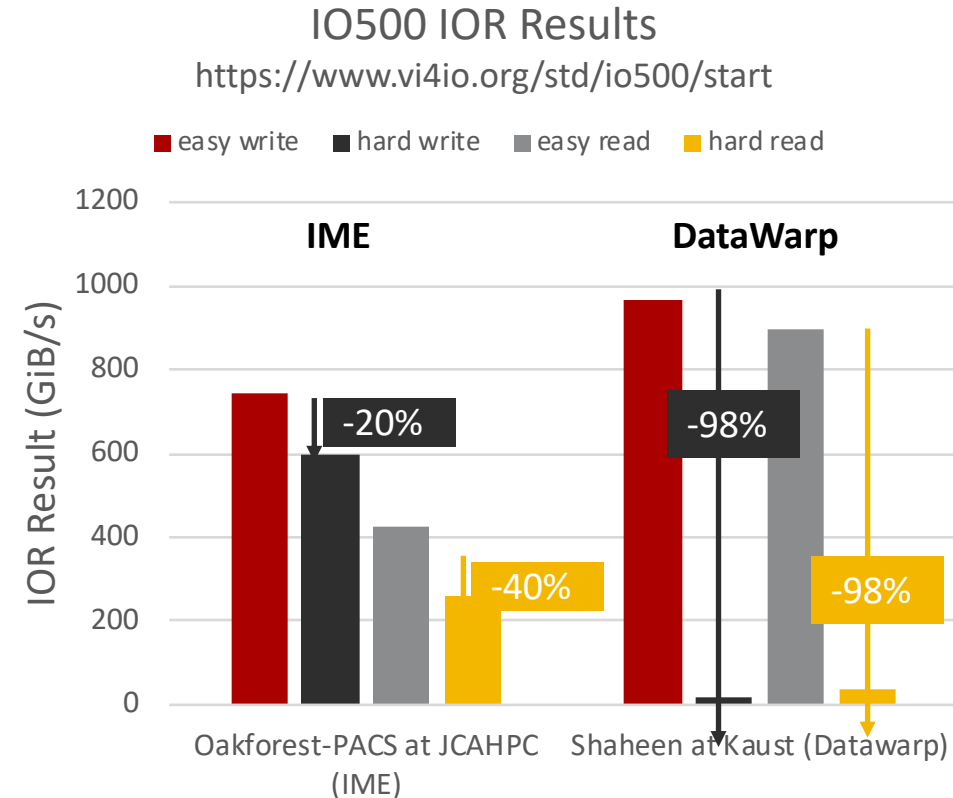


IME 240



- 2x16 full height PCIe slots for EDR/OPA/ 40/100GbE
- 2x CPUs
- 2 boot drives (1TB SATA)
- 8x 16GB 2400MHz DIMMs
- managed by CPU0
- managed by CPU1
- Up to 23 NVMe data Drives Minimum is ~12-14 to get full performance
- 1 SAS SSD (commit log)

## IME240 Sequential Throughput - File-per-Process
### IOR, 20x NVMe drives, 32 Clients, IB FDR



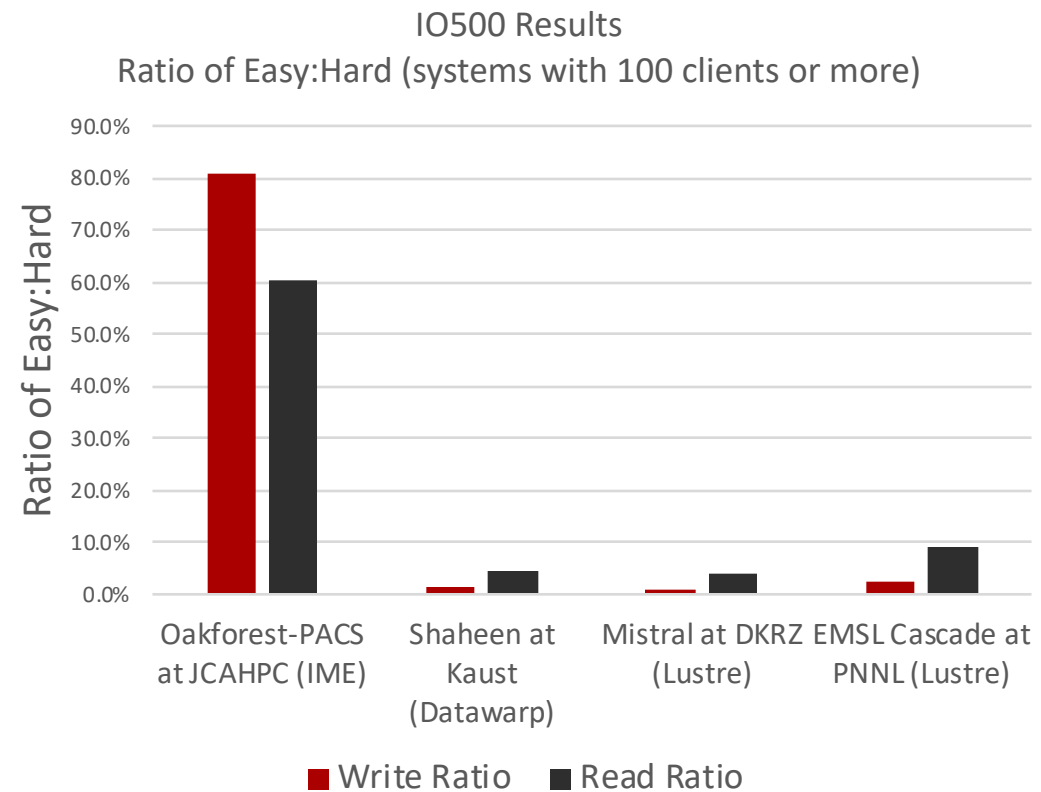| IME240 Node Count | 1 | 2 | 4 |
|---|---|---|---|
| ■ IME Write (MPIIO) | 22.1 | 42.8 | 85.8 |
| ■ IME Read (MPIIO) | 19.9 | 40.8 | 80.7 |

# APPLICATION EFFICIENCY FOR THE REAL WORLD

► IME's datapath is designed to deliver the potential of flash to the application

► Other Burst Buffers use a conventional filesystem which severely limits the ability to deliver flash performance

► The IO500 uses "Easy" and "Hard" IOR benchmarks

- IOR easy. You can set the parameters to be whatever you would like. You can use any of the modules such as HDF5 or MPI-IO. Typically people maximize performance by doing file-per-process and large aligned IO.

- IOR hard. We enforce a particular set of parameters. Specifically, the IOs are 47008 bytes each interspersed in a single shared file. Your only control is to specify how many writes each thread does.

► *Anyone can get good performance with enough equipment with the easy benchmark. Good Performance with the Hard Benchmark requires a new approach*
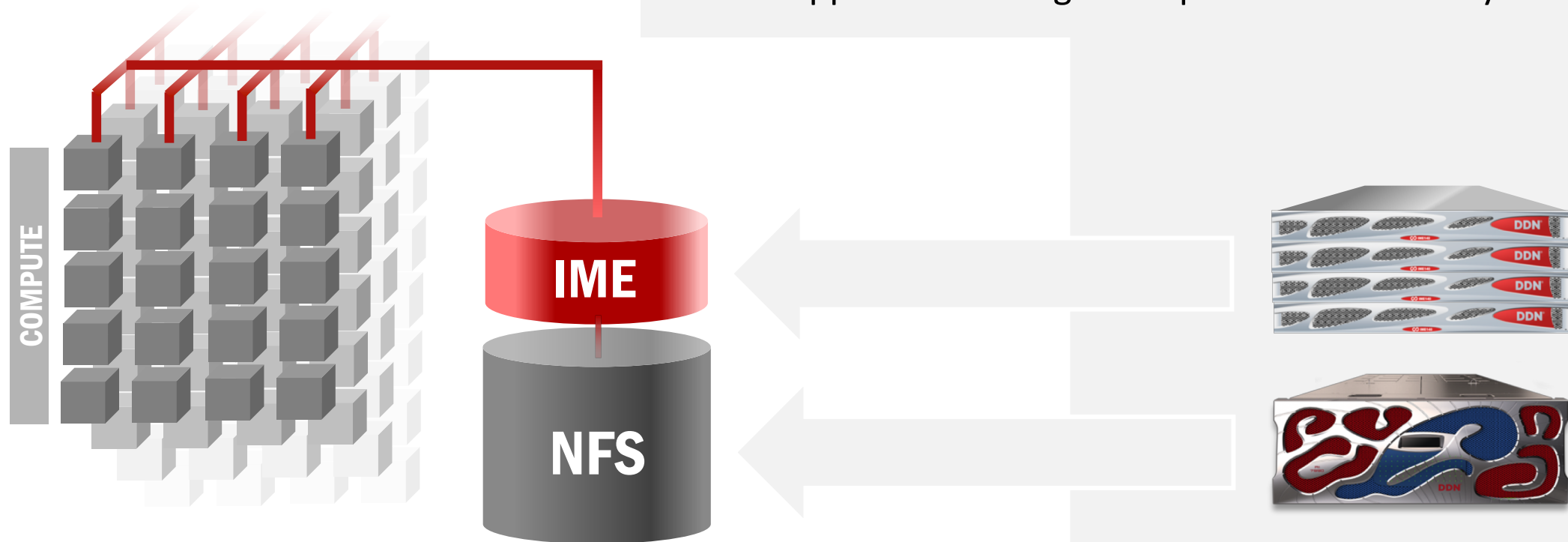


IO500 IOR Results
https://www.vi4io.org/std/io500/start

■ easy write  ■ hard write  ■ easy read  ■ hard read

IME          DataWarp
-20%    -98%
-40%    -98%

IOR Result (GiB/s)

Oakforest-PACS at JCAHPC (IME)      Shaheen at Kaust (Datawarp)

# APPLICATION EFFICIENCY FOR THE REAL WORLD

► Extracting results from IO500 where the client count is 100 nodes or more

► Filesystem options show huge degradation when the IO patterns is tough.

► Only IME is able to present Flash to the applications efficiently

IO500 Results
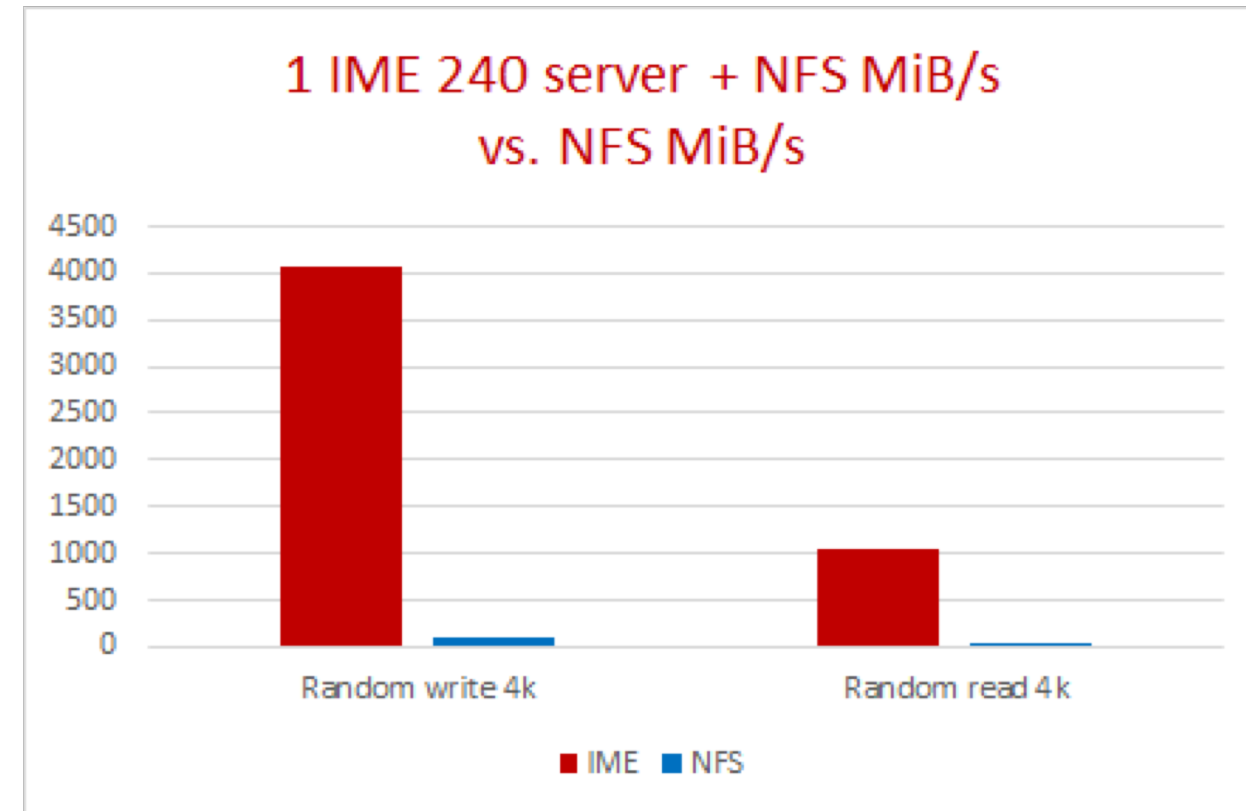Ratio of Easy:Hard (systems with 100 clients or more)

# IME – Burst buffer for NFS

► Brings scale-out Flash native performance to NFS access

► Shield NFS server from "tough" IO

► Increase IO throughput from NFS hardware

► Zero application changes  - replace NFS mount by IME mount

COMPUTE

IME

NFS

# IME – Burst buffer for NFS

IME with NFS

▶ Brings scale out Flash native performance to NFS Systems

▶ Removes complexity associated with Parallel Filesystems

▶ Shield NFS server for "bad" IO

▶ Increase IO throughput on top of NFS hardware

▶ No application changes - replace NFS mount by IME mount



1 IME 240 server + NFS MiB/s vs. NFS MiB/s

**DDN.COM/A3I**

# Thank You!

Keep in touch with us.

sales@ddn.com

@ddn_limitless

company/datadirect-networks

9351 Deering Avenue
Chatsworth, CA 91311

1.800.837.2298
1.818.700.4000

DDN STORAGE

ddn.com