# Applying DDN
# to
# Machine Learning

Jean-Thomas Acquaviva

jacquaviva@ddn.com

# Learning from What?

**Image data**

*Facial recognition*
*Action recognition*
*Object detection and recognition*
*Handwriting and character recognition*
*Aerial images*
*...*

**Anomaly data**
Time series

**Biological data**

*Human*
*Animal*
*Plant*
*Microbe*
*Drug Discovery*

**Multivariate data**

*Financial*
*Weather*
*Census*
*Transit*
*Internet*
*Games*
*Other multivariate*

**Text data**

*Reviews*
*News articles*
*Messages*
*Twitter and tweets*
*Social network*

**Signal data**

*Electrical*
*Motion-tracking*
*Other signals*

**Physical data**

*High-energy physics*
*Systems*
*Astronomy*
*Earth science*
*...*

**Sound data**

*Music*
*Speech data*

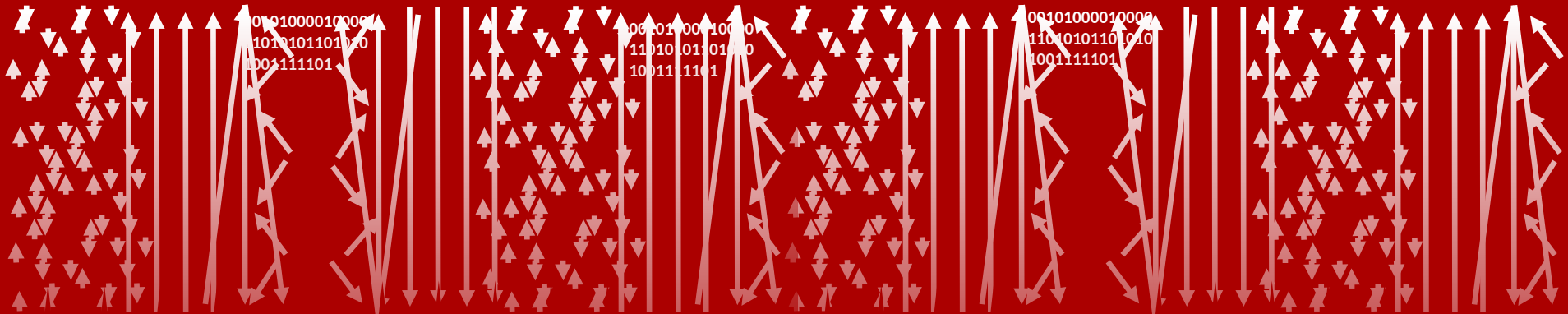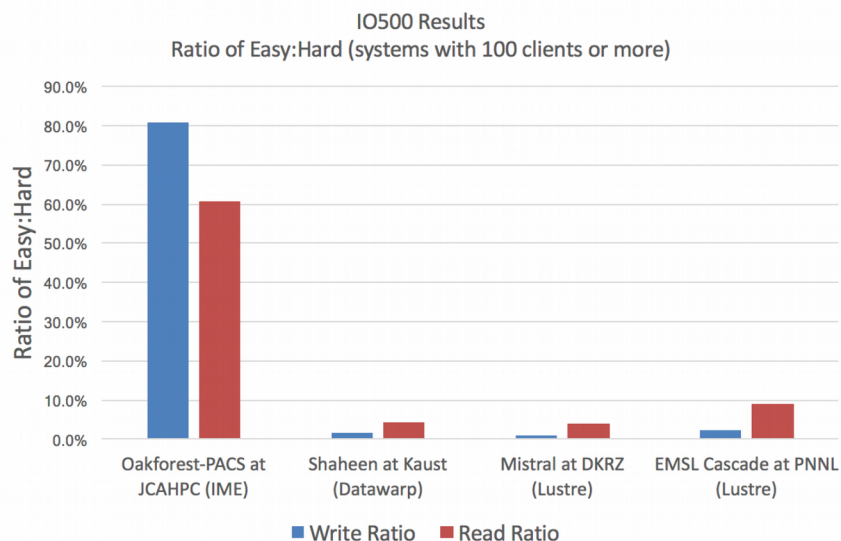| Type of data | Supported operations |
|---|---|
| discret quantitative data | Calculations, equality / difference, inferiority / superiority |
| Continuous quantitative data | Calculations, equality / difference, inferiority / superiority |
| nominal qualitative data | Equality / difference |
| ordinal qualitative data | Equality / difference, inferiority / superiority |

Machine Learning

Big Data

NoSQL Analytics

IO Characteristics: Read, Random, High Throughput per Client, File and IO Sizes between a few kb and a few MB
**Training Sets typically larger than local caches**

# Diversity of Load:  IO500

Detailed write

| Rank | System | Institution | Filesystem | Client Nodes | Score | BW | MD | Easy Write | Hard Write | Hard vs. Easy | Easy Read | Hard Read | Hard vs. Easy |
|------|--------|-------------|------------|--------------|-------|------|--------|------------|------------|---------------|-----------|-----------|---------------|
| | | | | | | GiB/s | kIOP/s | GiB/s | GiB/s | | GiB/s | GiB/s | |
| 1 | Oakforest-PACS | JCAHPC | IME | 2048 | 101.48 | 471.25 | 19.04 | 742.38 | 600.28 | 80.9% | 427.41 | 258.93 | 60.6% |
| 2 | Shaheen | Kaust | DataWarp | 300 | 70.9 | 151.53 | 33.17 | 969.45 | 15.55 | 1.6% | 894.76 | 39.09 | 4.4% |
| 3 | Shaheen | Kaust | Lustre | 1000 | 41 | 54.17 | 31.03 | 333.03 | 1.44 | 0.4% | 220.62 | 81.38 | 36.9% |



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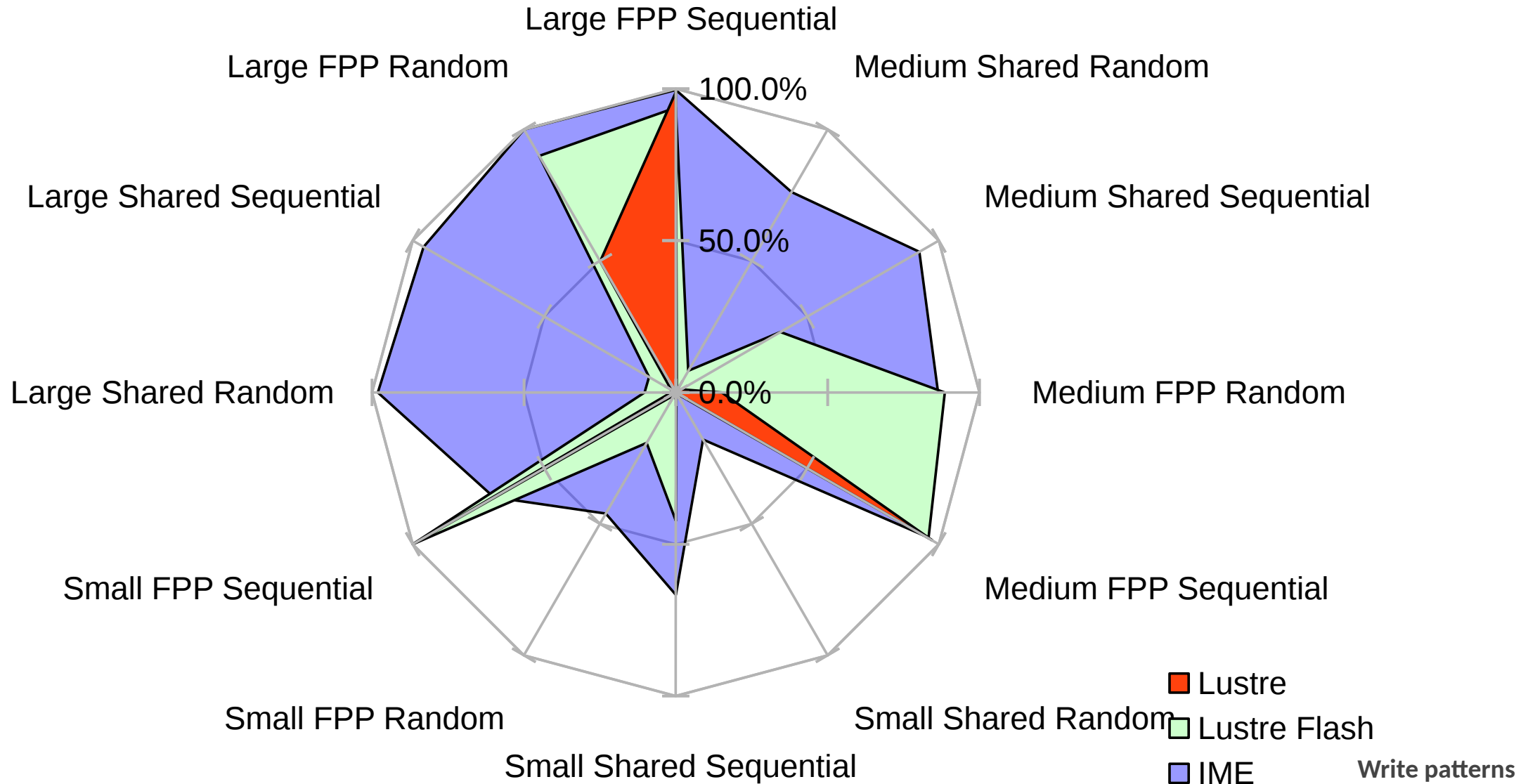Ratio of Easy:Hard (systems with 100 clients or more)

- ▶ KAUST BurstBuffer and Lustre at DKRZ show massive falls in IO performance
- ▶ Small DDN Lustre based on 12K at PNNL shows a similar pattern
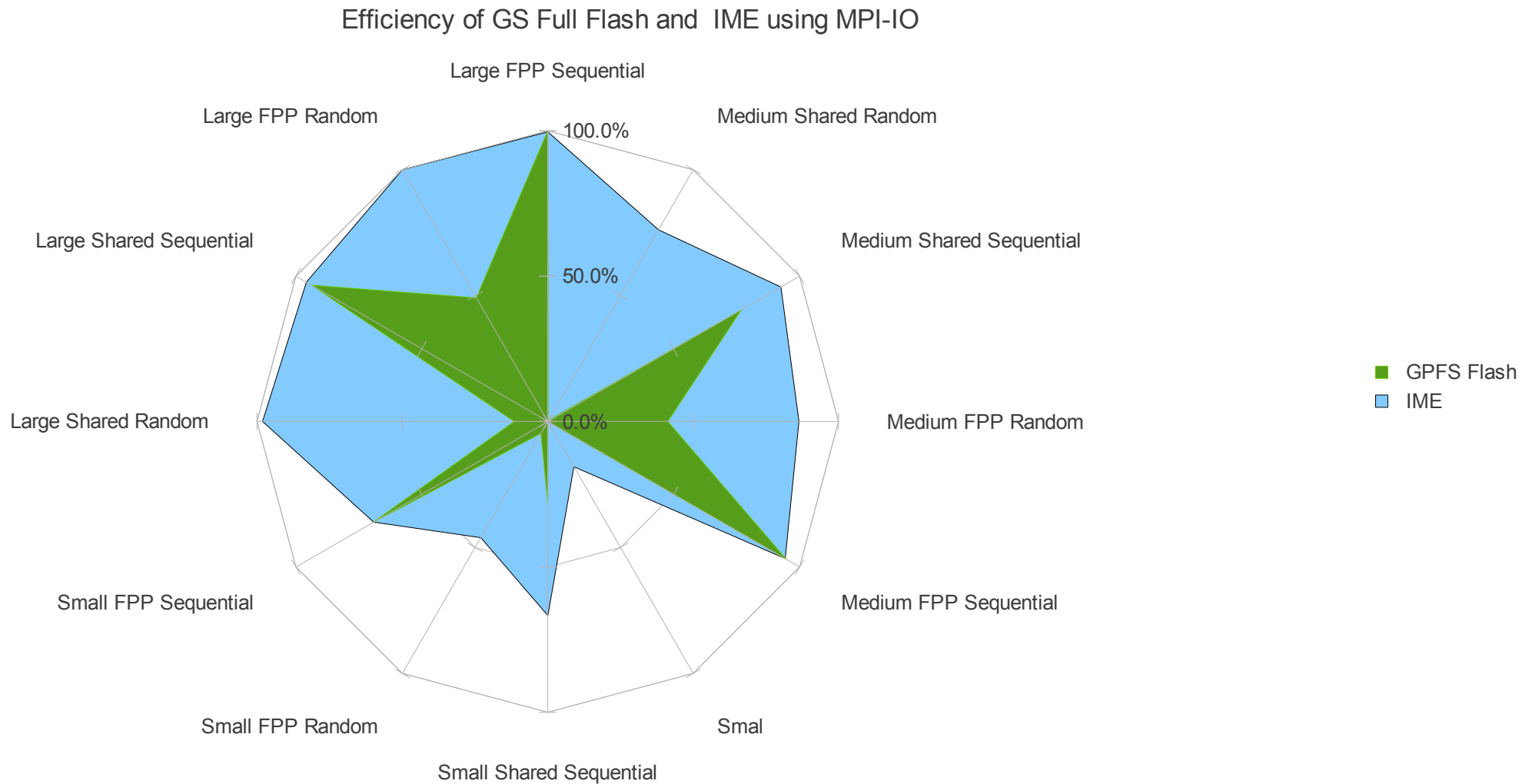- ▶ IME has a order of magnitude better ratio between easy and hard

# Acknowledging multi-criteria performance metrics

| I/O Granularity | I/O control plane Pattern | I/O Data plane Pattern | |
|---|---|---|---|
| Large (>= 1MB) | File Per Process ( = share nothing) | Sequential | **IO500 Easy !** |
| Large | File Per Process | Random | |
| Large | Single Shared File | Sequential | |
| Large | Single Shared File | Random | |
| Small | File Per Process | Sequential | |
| Small | File Per Process | Random | |
| Small (47008 Bytes) | Single Shared File | Sequential | **IO500 Hard !** |
| Small | Single Shared File | Random | |

# IO500 to a comprehensive picture: DDN Flash native vs Lustre



Large FPP Sequential

Medium Shared Random

Large FPP Random

Medium Shared Sequential

Large Shared Sequential

100.0%

Medium FPP Random

Large Shared Random

50.0%

0.0%

Medium FPP Sequential

Small FPP Sequential

Small FPP Random

Small Shared Random

Small Shared Sequential

Write patterns

- Lustre
- Lustre Flash
- IME

DDN Storage | © 2018 DDN Storage

# IO500 to a comprehensive picture: DDN Flash native E vs GridScaler



Efficiency of GS Full Flash and IME using MPI-IO

- Large FPP Sequential
- Medium Shared Random
- Large FPP Random
- Medium Shared Sequential
- Large Shared Sequential
- Medium FPP Random
- Large Shared Random
- Medium FPP Sequential
- Small FPP Sequential
- Smal
- Small FPP Random
- Small Shared Sequential

100.0%
50.0%
0.0%

- GPFS Flash
- IME

**Write patterns**

# Example: EXAScaler DGX Solution (hardware view)



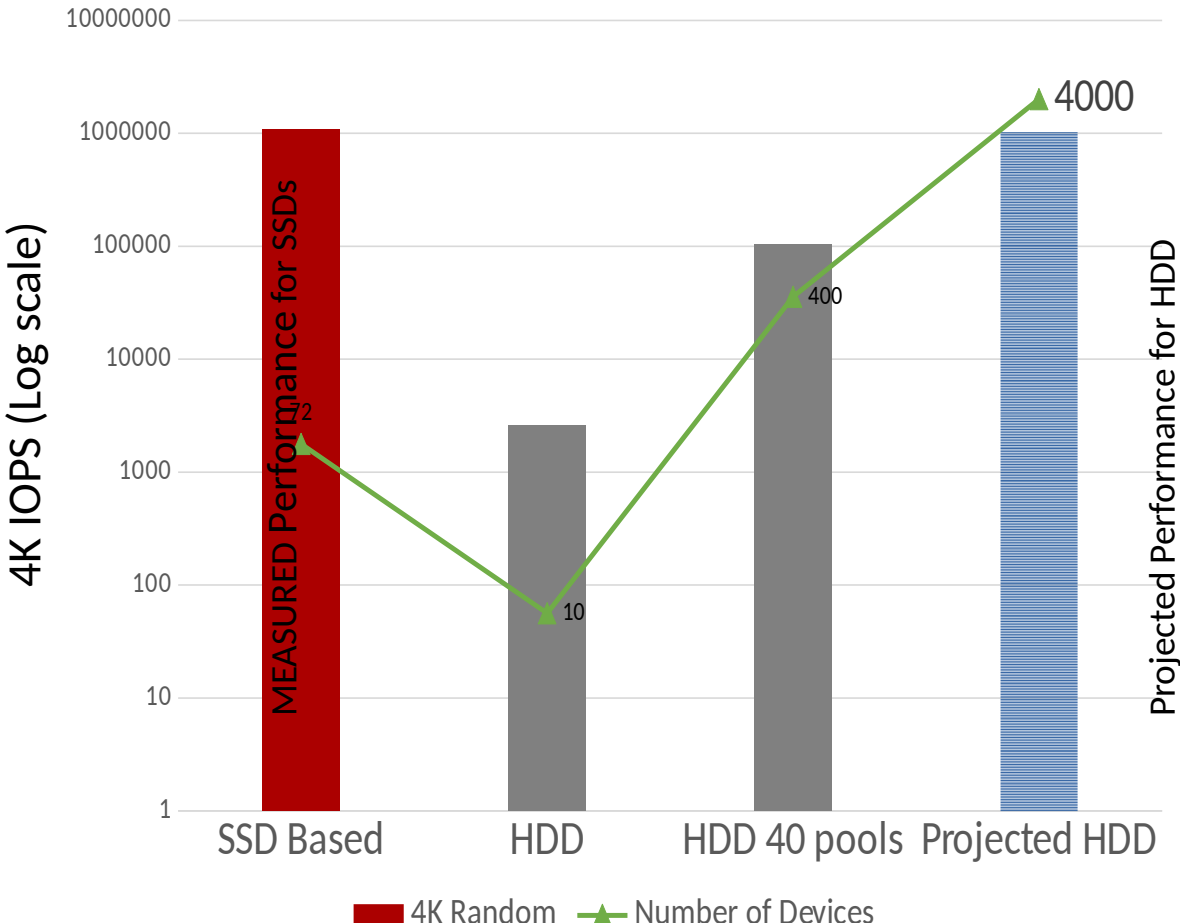72 x 1.2 TB NVMe

ES14KXE 72 SSDs

# Platform DDN ES14KXE Full Flash: 1M IOPS – 40 GB/s

- ES14KX ALL Flash active-active controllers deliver 1M file IOPs – the equivalent of 4000 HDDs

- Scale IOPs further in the namespace with additional controllers

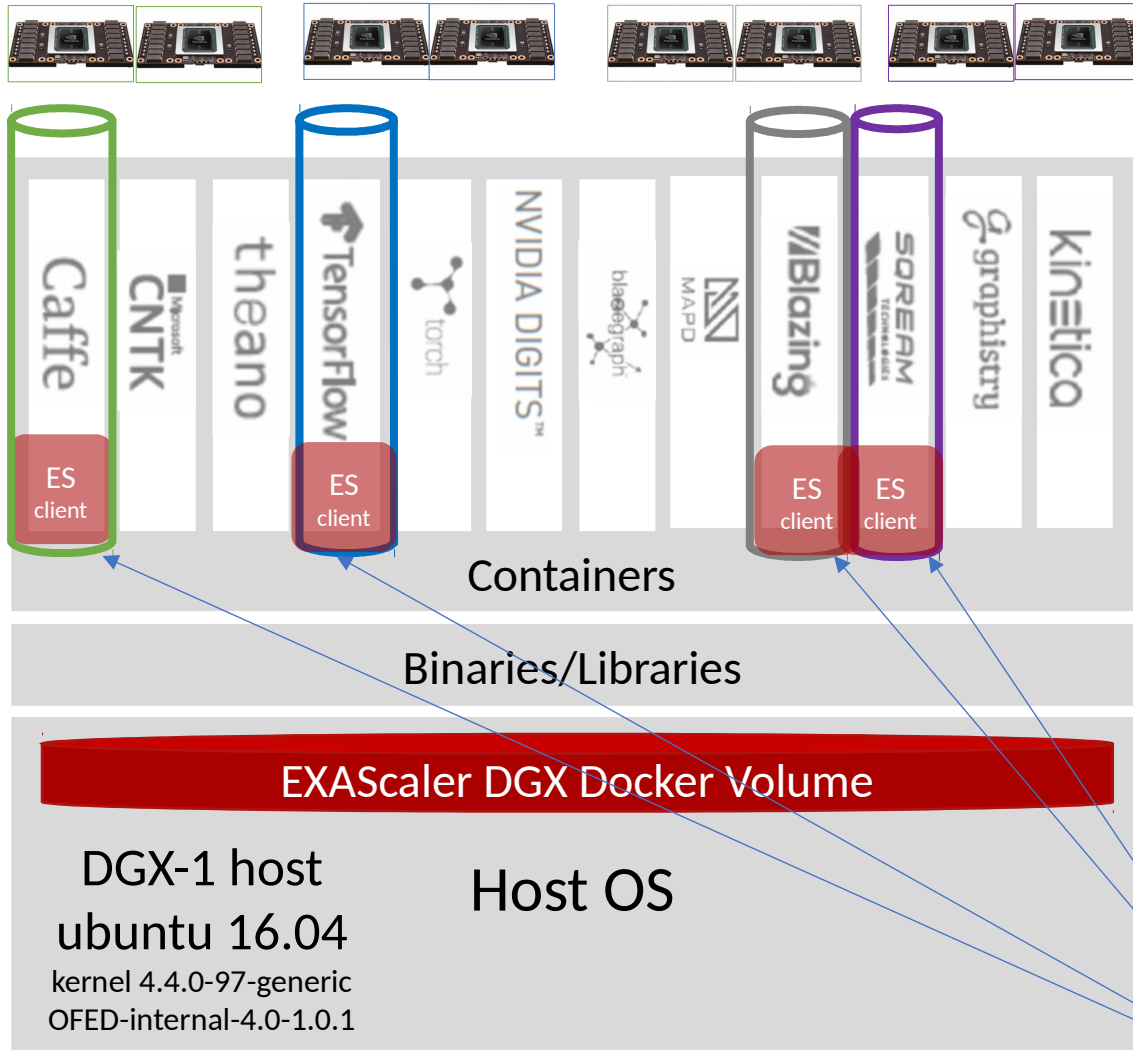- Augment Flash with HDD at scale with up to 1680 HDDs per controller

Random Read 4K IOPS on ES14KX (All Flash vs. HDD)



4K IOPS (Log scale)

MEASURED Performance for SSDs

Projected Performance for HDD

10000000
1000000
100000
10000
1000
100
10
1

2

10

400

4000

SSD Based    HDD    HDD 40 pools    Projected HDD

■ 4K Random    ▲ Number of Devices

# WHAT PFS FOR AI APPLICATIONS?

| Feature | Importance for AI | GPFS | Lustre |
|---|---|---|---|
| Shared Metadata Operations | **High** - training data are usually curated into a single directory | ✘ Lower than 10K (minimal improvements with v5) | ✔ Up to 200K |
| Support for high-performance mmap() I/O Calls | **High** - many AI applications use mmap() calls | ✘ Extremely poor | ✔ Strong |
| Container Support | **High** - most AI applications are containerized | ✘ Poor (network complexity & root issues) | ✔ Available |
| Data Isolation for Containers | **Medium/High** – important for shared environments | ✘ Not available today | ✔ Available |
| Data-on-Metadata (small file support) | **Medium/High** – depends on data set | ✘ DOM only for files smaller than 3.4k | ✔ DOM is highly tunable |
| Unique Metadata Operations | **Medium** - depends on Installation Size and Application Workflow | ✔ Highly scalable | ✔ Highly scalable with DNE 1/2 |

# Example: EXAScaler DGX Solution (host part)



Tesla P100 NVLINK (170 Tflps)

- Integrated Flash Parallel File System Access via TCP or IB

- Extreme Data Access Rates for concurrent DGX Containers

**EXAScaler DGX Docker Volume**

- Lustre ES3.2 kernel modules compiled for Ubuntu kernel and host's OFED
- Lustre userspace tools
- scripts for Lustre mount/umount

Resource isolation: Ios/GPUs/NIC/memory/namespaces For the application/SW suite
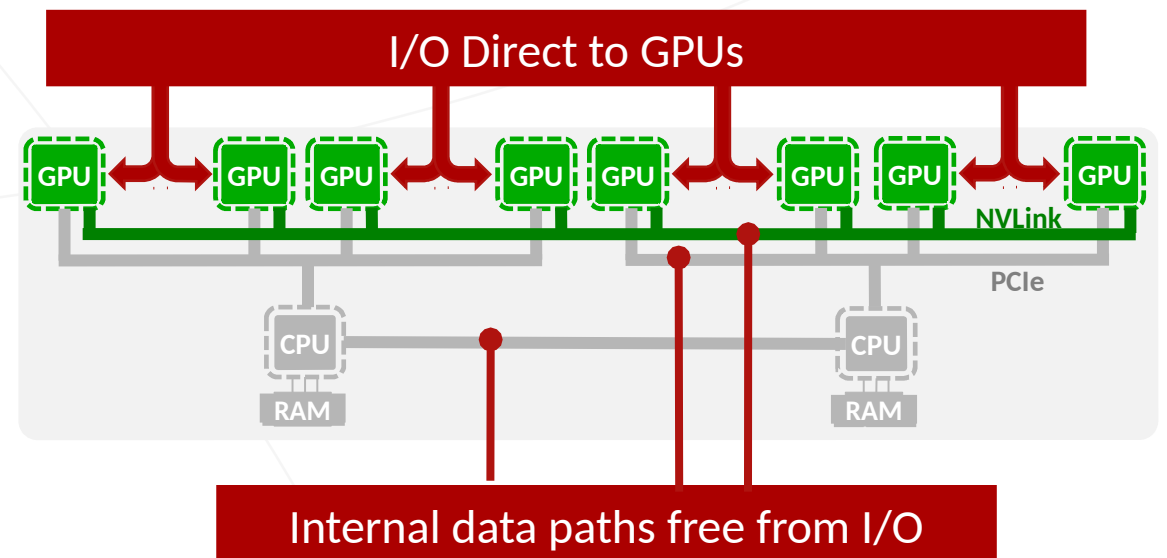
## CONTAINER PINNING

DDN's EXASCALER for DGX manages I/O-paths optimally through DGX-1 to maximize performance to your AI application and keep IO traffic from consuming internal data paths

## mmap() support

LMDB is key to manage file in several framework (caffé). LMDB relies on mmap() hence page fault to trigger kernel internal I/O call (read_ahead)
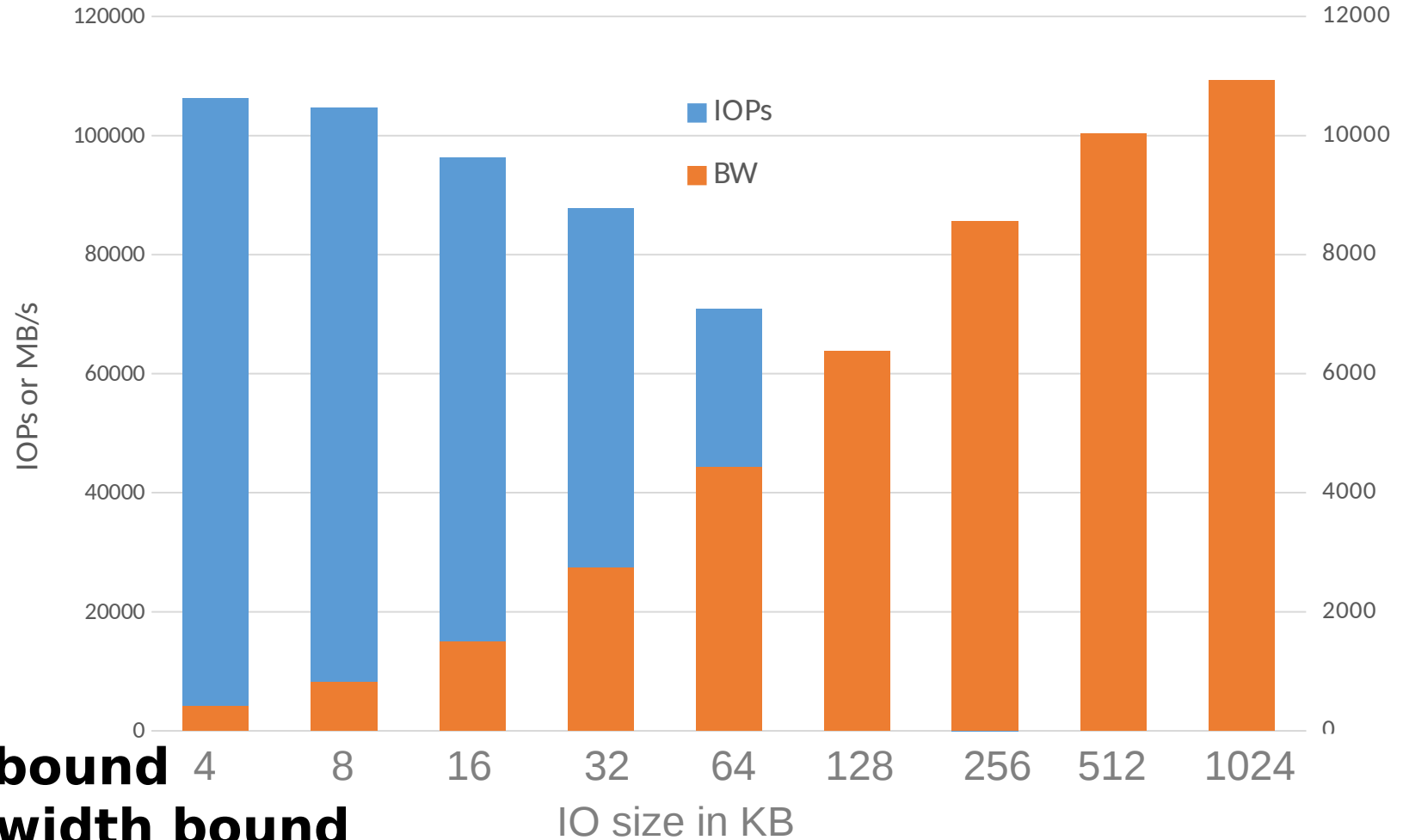
I/O Direct to GPUs

GPU GPU GPU GPU GPU GPU GPU GPU

NVLink

PCIe

CPU CPU

RAM RAM

Internal data paths free from I/O

# CONTAINER PINNING OPTIMIZATION

## Random Read (4k) IOPs



Legend:
- Directed I/O
- Unoptimized

Y-axis: IOPs — 0, 50,000, 100,000, 150,000, 200,000, 250,000, 300,000

X-axis: Container Count — 1, 2, 4, 8

**FROM IOPS TO BANDWIDTH**

Single Container Performance (Lustre/SSD)

IOPs
BW

IOPs or MB/s

IO size in KB

4  8  16  32  64  128  256  512  1024
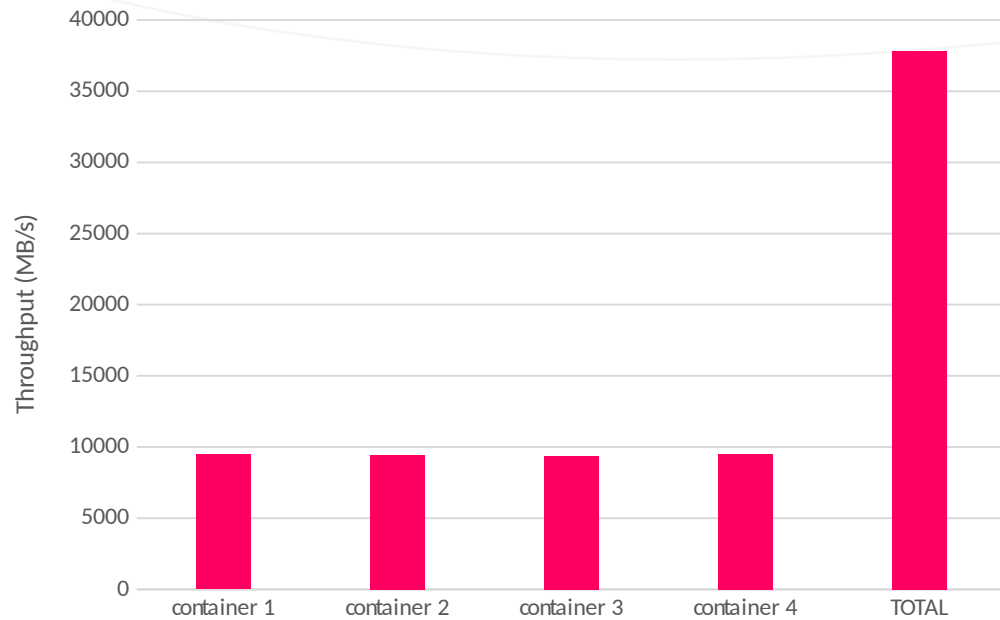
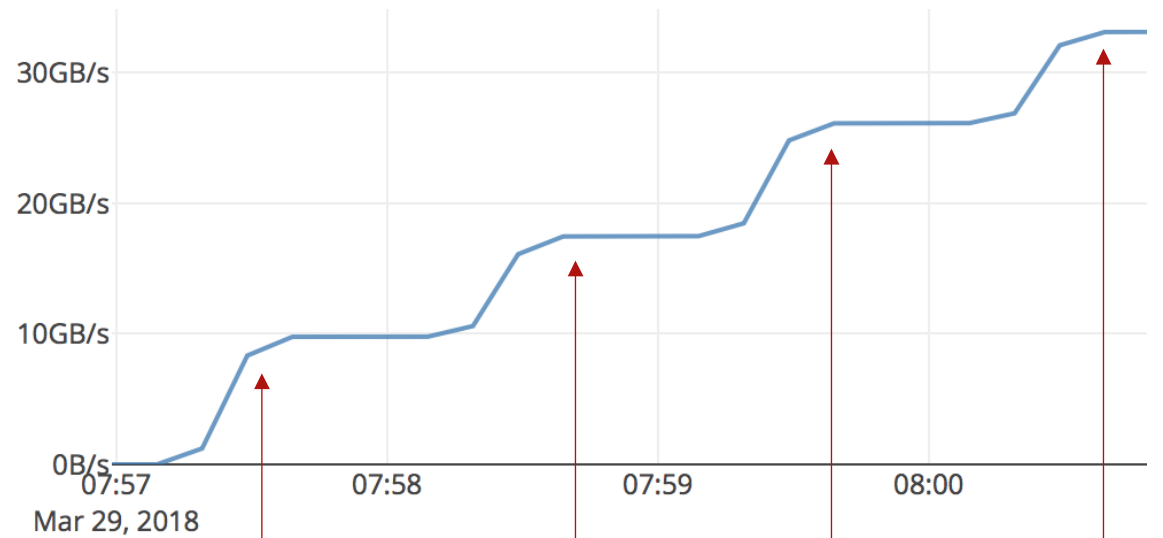**4KB random read is IOPS bound**
**1MB random read is Bandwidth bound**

# SCALING UP WORKLOADS IN DGX-1

Container Throughput



Overall Throughput



1 container     2 containers     3 containers     4 containers

# Remove I/O burden from data scientist shoulders

→ I/O no longer the limiting factor
→ Saturation of the network
→ 250 KIOPS on a DGX-1
→ 1 Millions IOPS with 4 DGX-1

## Single DGX-1

**38GB/s
>250  KIOPs**

# Bringing HPC technologies and know-how to analytics = x12

SCALE
SPEED

SCALE
VOLUME

**ULTIMATE FLEXIBILITY**

40GB/s in 4RU

1PB in 4RU

**1TB/s &
10M IOPS**

**12x More Answers/Minute for
Your AI at Any Scale!**

**Over 20PB**

DDN
STORAGE

ddn.com