# DDN

# Accelerating Lustre with SSDs and NVMe

James Coomer, DDN

# DDN ExaScaler Software Components

| Data Management | Fast Data Copy | Tape | S3 (Cloud) | Object (WOS) |
| | | ExaScaler Data Management Framework | | |

| Management and Monitoring | DDN DirectMon | DDN ExaScaler Monitor | Intel IML |
| | DDN ExaScaler | | |

| Global Filesystem | DDN Clients | NFS/CIFS/S3 | DDN IME | Intel Hadoop |
| | DDN Lustre Edition with L2RC | | |
| | OFD Read Cache Layer | | |

| Local Filesystem | ldiskfs | OpenZFS | btrfs |
| | Intel DSS | | |

| Storage Hardware | DDN Block Storage with SFX Cache | 3rd Party HW |

Level 3 support provided by: | DDN | DDN & Intel | Intel HPDD | 3rd Party |

# DDN | ES14K
## Designed for Flash and NVMe

**Configuration Options**

- 72 SAS SSD or 48 NVMe
- SSDs or HDDs only
- HDDs with SSD caching
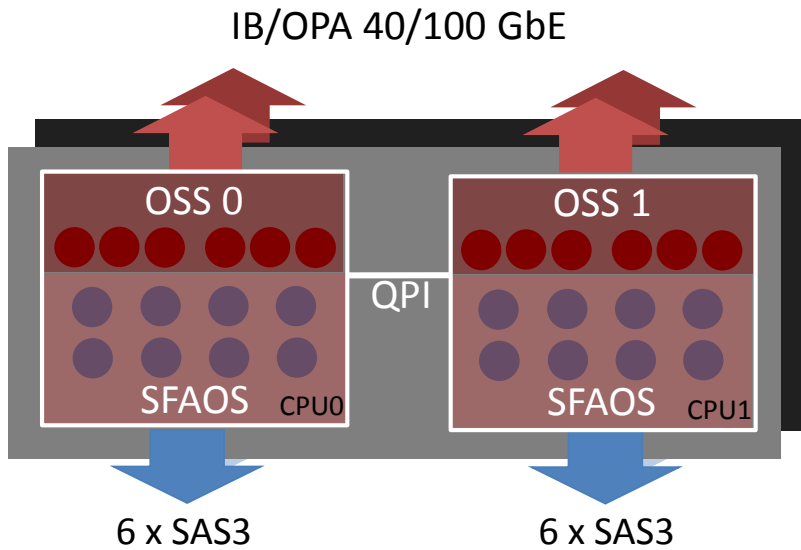- SSDs with HDD tier

**Connectivity**

- FDR/EDR
- OmniPath
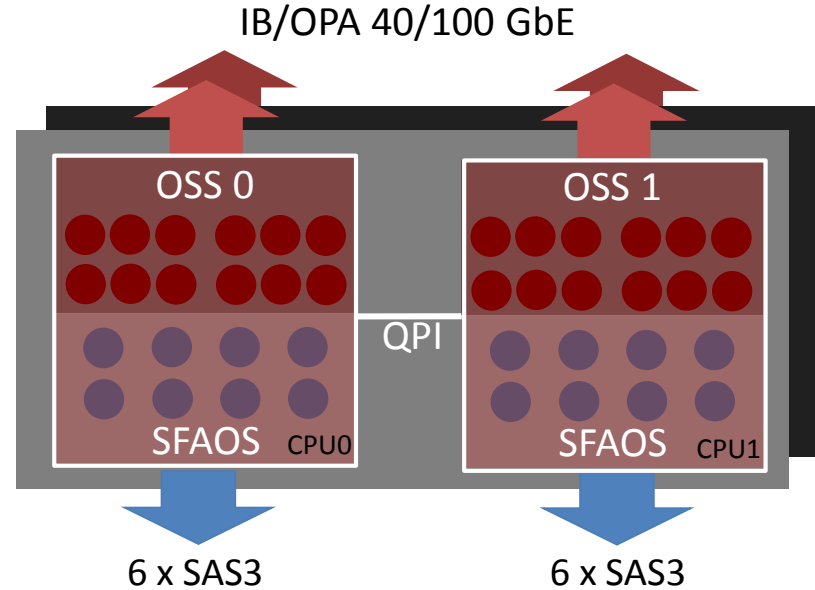- 40/100GbE

**Industry Leading Performance in 4U**

- Up to 40 GB/sec throughput
- Up to 6 million IOPS to cache
- Up to 3.5 million IOPS to storage
- 1PB+ capacity (with 16TB SSD)
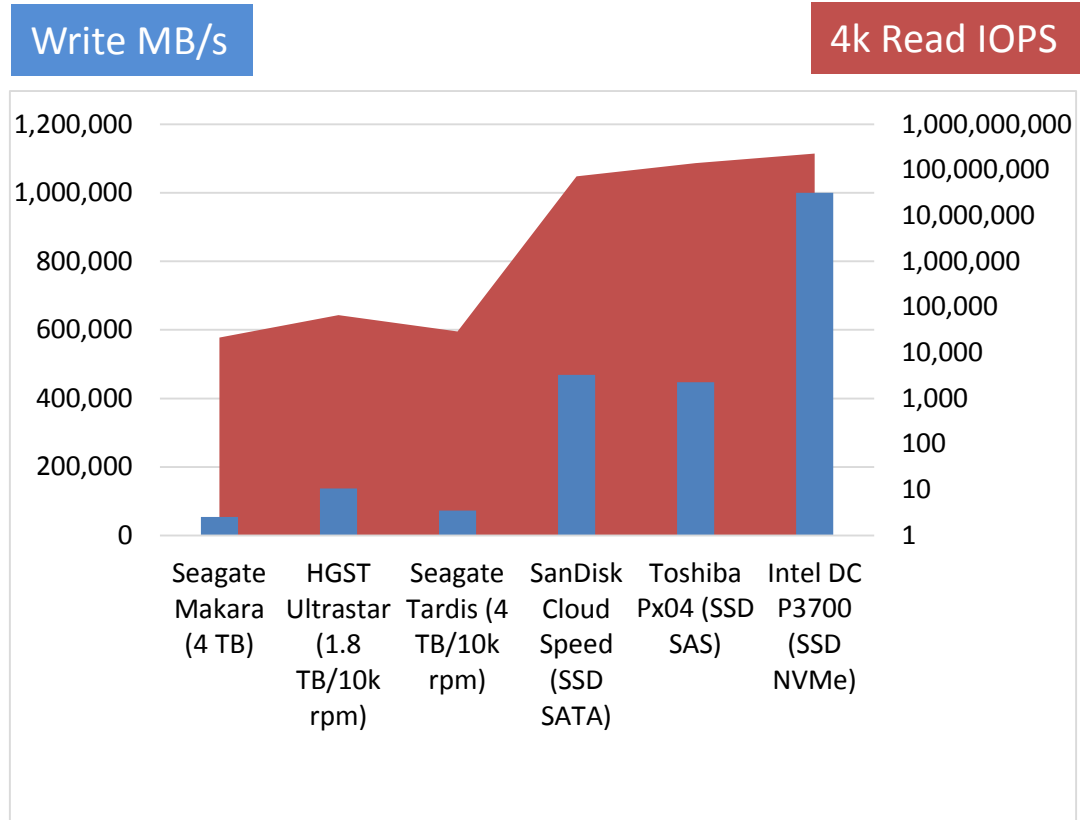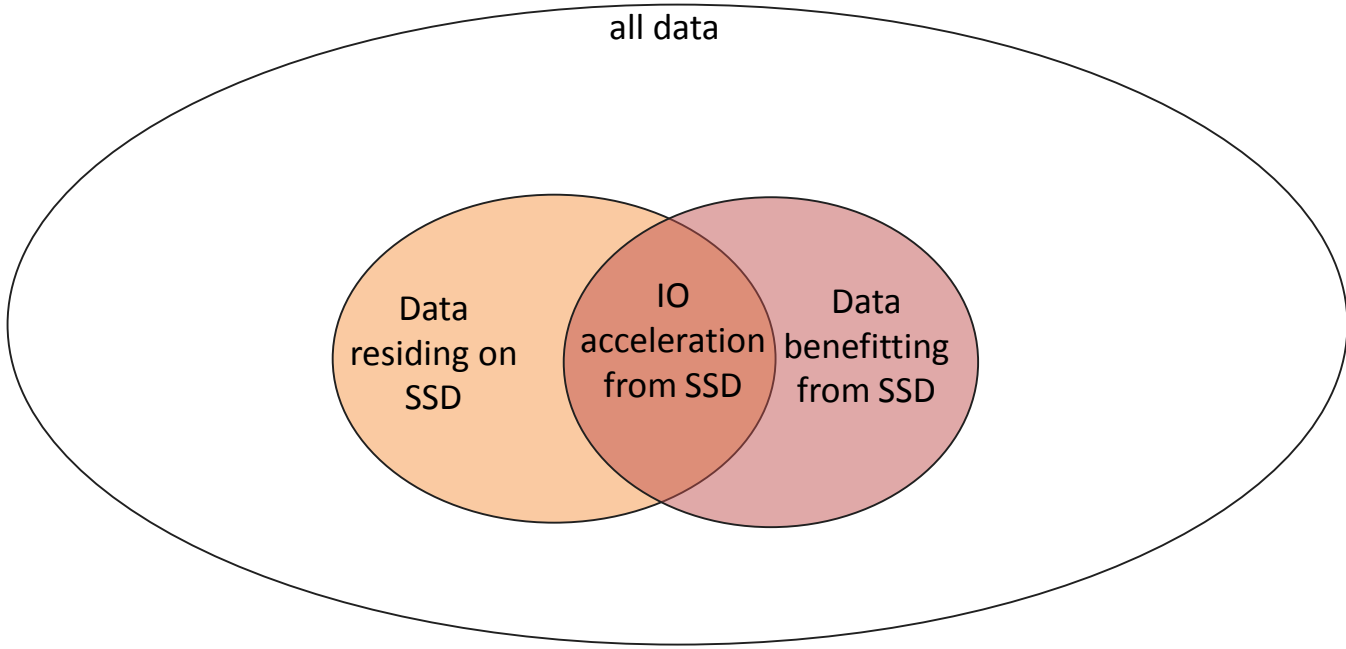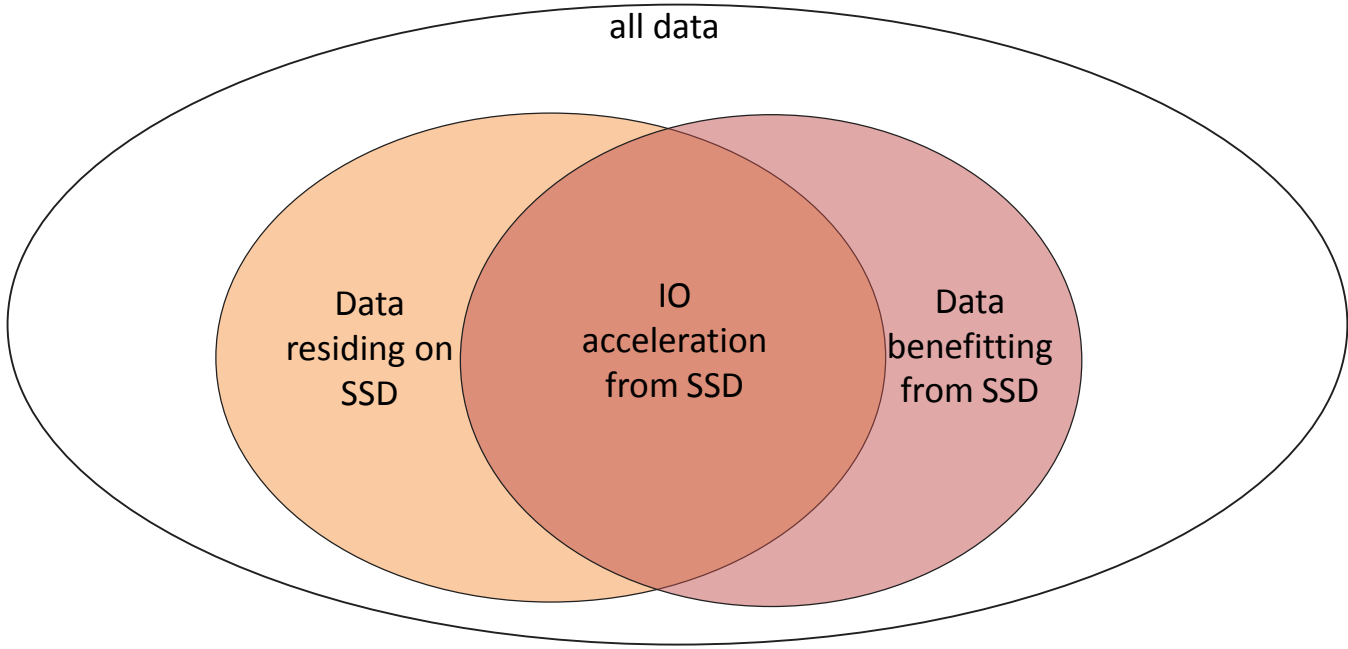- 100 millisecond latency

# ES14K Architecture

# Why SSD Cache?
Don't blow the power/space/management with spindles

**DDN**®

SSDs still pricey... So

▶ Optimise Data for SSDs

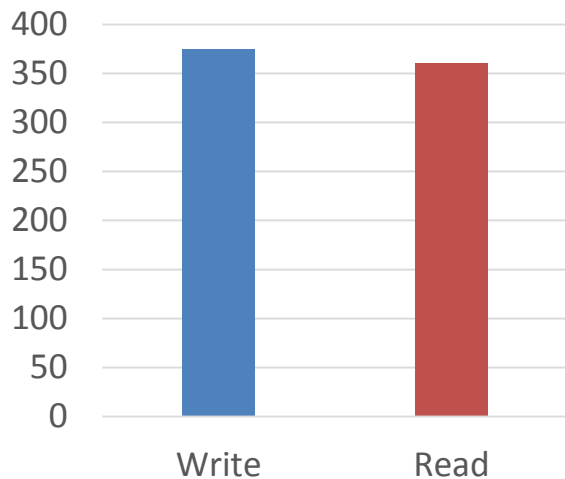▶ Optimise SSDs for Data



Write MB/s

4k Read IOPS

# SSD Options

**DDN®**

## 1 — All SSD Lustre

- Lustre HSM for Data Tiering to HDD namespace
- Generic Lustre I/O
- Millions of Read and Write IOPS

## 2 — SFX

- Block-Level Read Cache
- Instant Commit, DSS, fadvice()
- Millions of Read IOPS

## 3 — L2RC

- OSS-level Read cache
- Heuristics with **FileHeat**
- Millions of Read IOPS

## 4 — IME

- I/O level Write and Read Cache
- Transparent + hints
- 10s of Millions of Read/Write IOPS

OSS  OSS  OSS

OSS  OSS  OSS
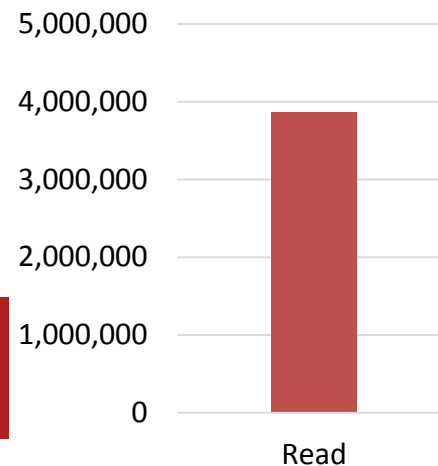
OSS  OSS  OSS

**IME**

OSS  OSS  OSS

# 1. Rack Performance: Lustre

**DDN**®

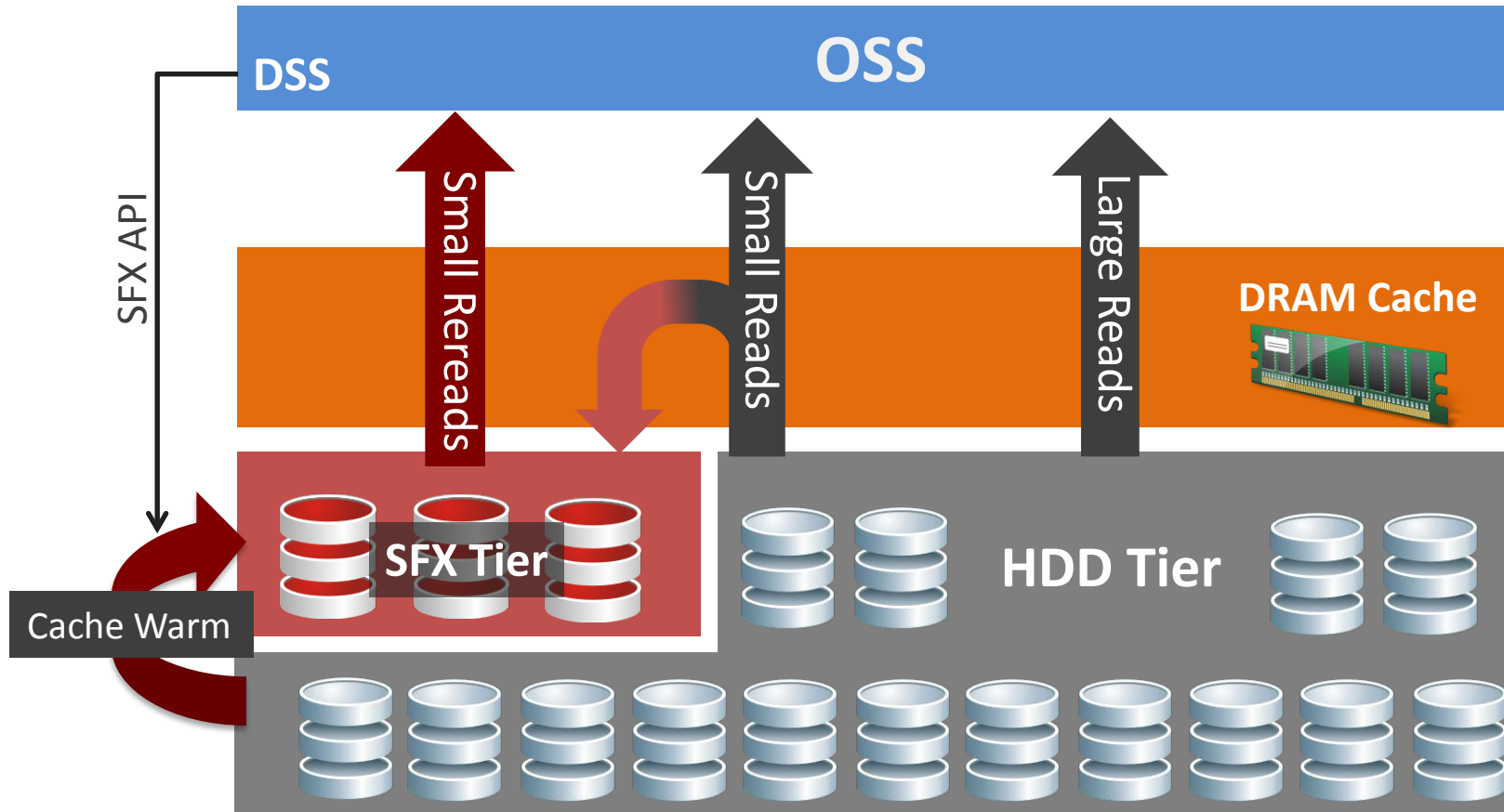**IOR File-per-Process (GB/s)**



**4k Random Read IOPS**

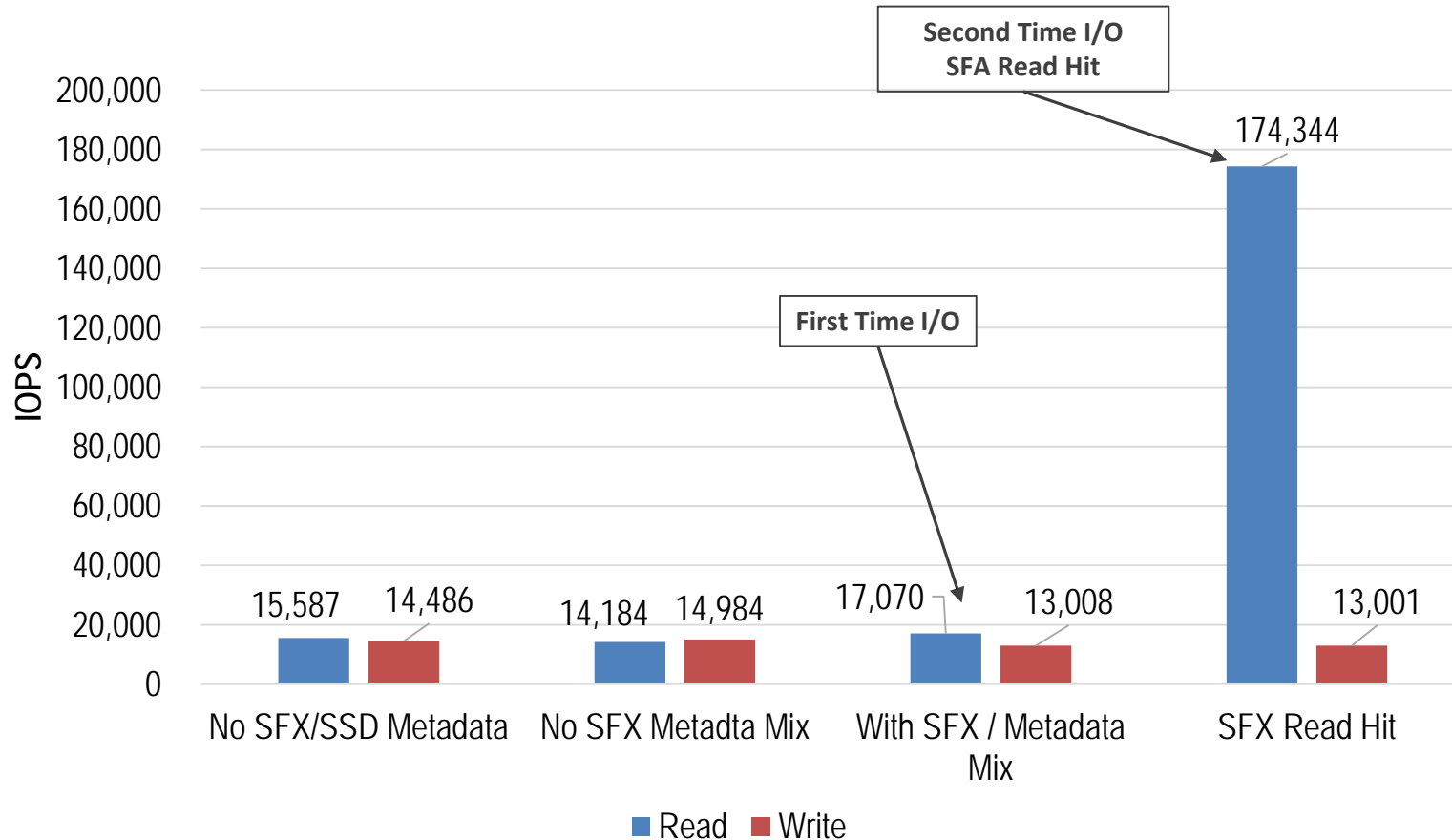up to 4PB Flash Capacity

4 Million IOPs

350GB/s Read and Write (IOR)

# 2. SFX & ReACT – Accelerating Reads
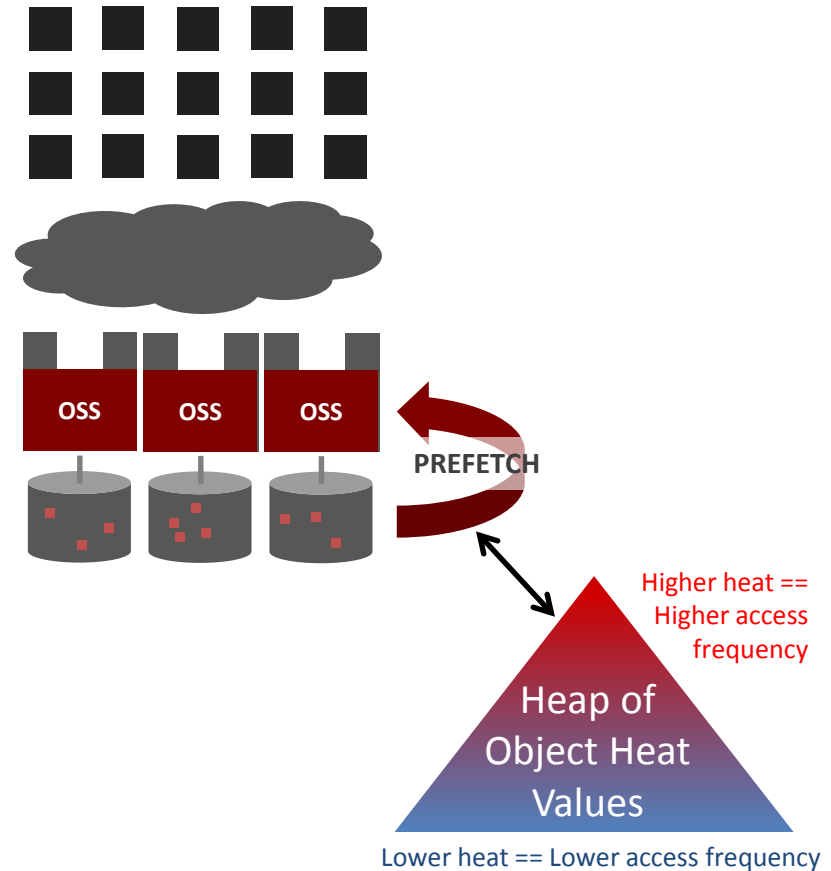Integrated with Lustre DSS

# 2. 4 KiB Random I/O

# 3. Lustre L2RC and File Heat

OSS-based Read Caching

- Uses SSDs (or SFA SSD pools) on the OSS as read cache
- Automatic prefetch management based on file heat
- File-heat is a relative (tunable) attribute that reflects file access frequency
- Indexes are kept in memory (worst case is 1 TB SSD for 10 GB memory)
- Efficient space management for the SSD cache space (4KB-1 MB extends)
- Full support for ladvice in Lustre

**OSS**   **OSS**   **OSS**

**PREFETCH**

Higher heat == Higher access frequency

Heap of Object Heat Values

Lower heat == Lower access frequency

# 3. File Heat Utility

- tune the arguments of file heat with proc interfaces
  ```
  /proc/fs/lustre/heat_period_second
  /proc/fs/lustre/heat_replacement_percentage
  ```
- Utils to get file heat values: `lfs heat_get <file>`
- Utils to set flags for file heat:
  ```
  lfs heat_set [--clear|-c] [--off|-o] [--on|-O] <file>
  ```
- Heat can be cleared by: `lfs heat_set --clear`
- Heat accounting of a file can be turned off by: `lfs heat_set --off`
- Heaps on OSTs which can be used to dump lists of FIDs sorted by heat:
  ```
  [root@server9-Centos6-vm01 cache]# cat
  /proc/fs/lustre/obdfilter/lustre-OST0000/heat_top
  [0x200000400:0x1:0x0] [0x100000000:0x2:0x0]: 0 740 0 775946240
  [0x200000400:0x9:0x0] [0x100000000:0x6:0x0]: 0 300 0 314572800
  [0x200000400:0x8:0x0] [0x100000000:0x5:0x0]: 0 199 0 208666624
  [0x200000400:0x7:0x0] [0x100000000:0x4:0x0]: 0 100 0 104857600
  [0x200000400:0x6:0x0] [0x100000000:0x3:0x0]: 0 100 0 104857600
  ```
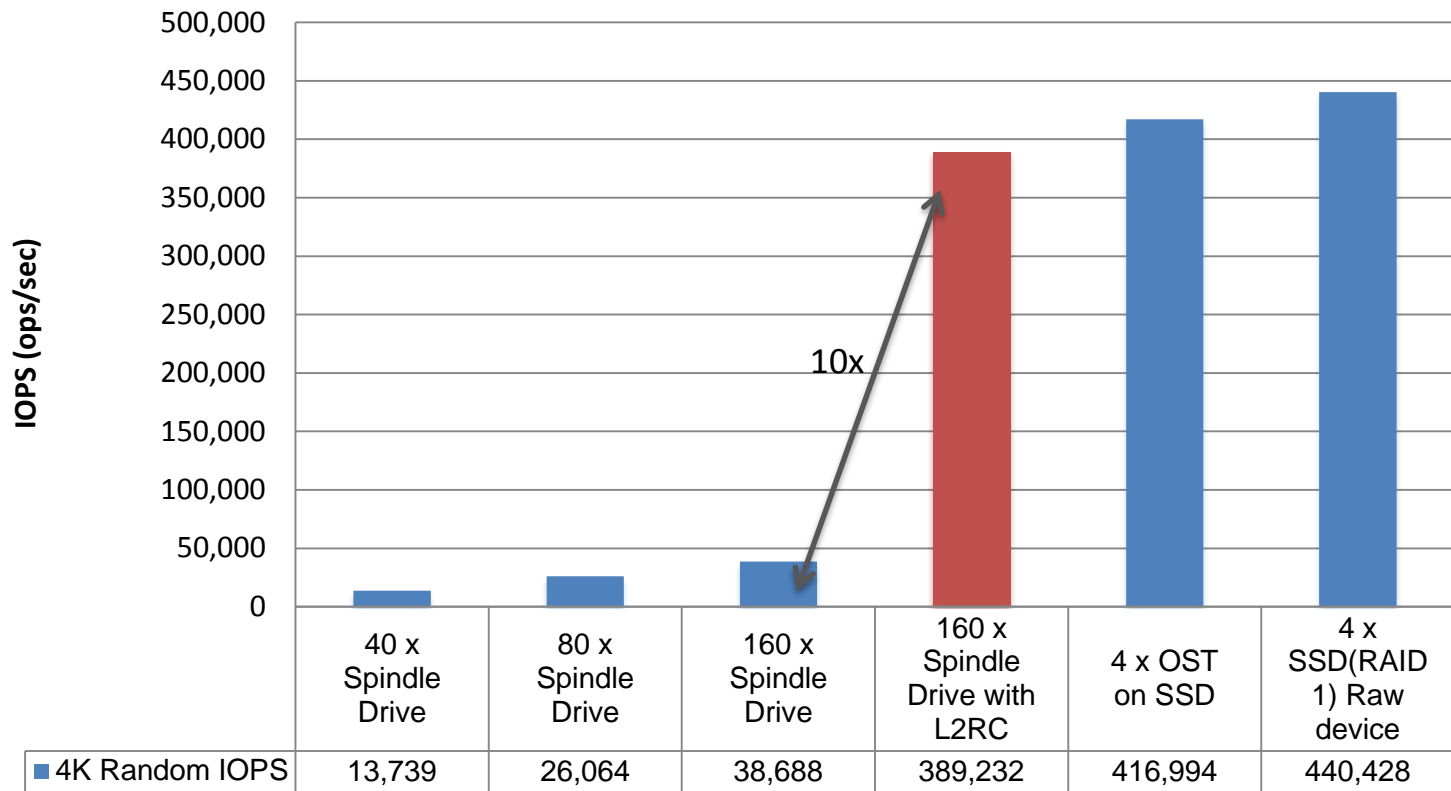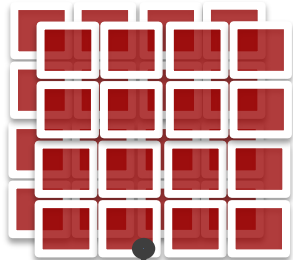
# 3. Random Read Performance with L2RC

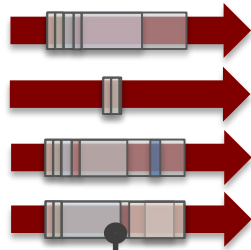**4KB Random Read IOPS (HDD/SSD based OST vs. OST & L2RC)**



| | 40 x Spindle Drive | 80 x Spindle Drive | 160 x Spindle Drive | 160 x Spindle Drive with L2RC | 4 x OST on SSD | 4 x SSD(RAID 1) Raw device |
|---|---|---|---|---|---|---|
| ■ 4K Random IOPS | 13,739 | 26,064 | 38,688 | 389,232 | 416,994 | 440,428 |

# DDN | IME

Application I/O Workflow



**COMPUTE**

**NVM TIER**

**LUSTRE**

Lightweight IME client intercepts application I/O. Places fragments into buffers + parity
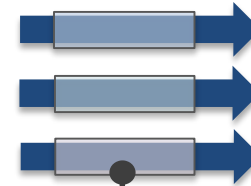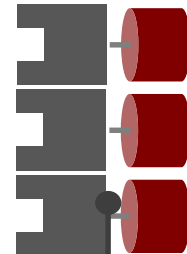
IME client sends fragments to IME servers

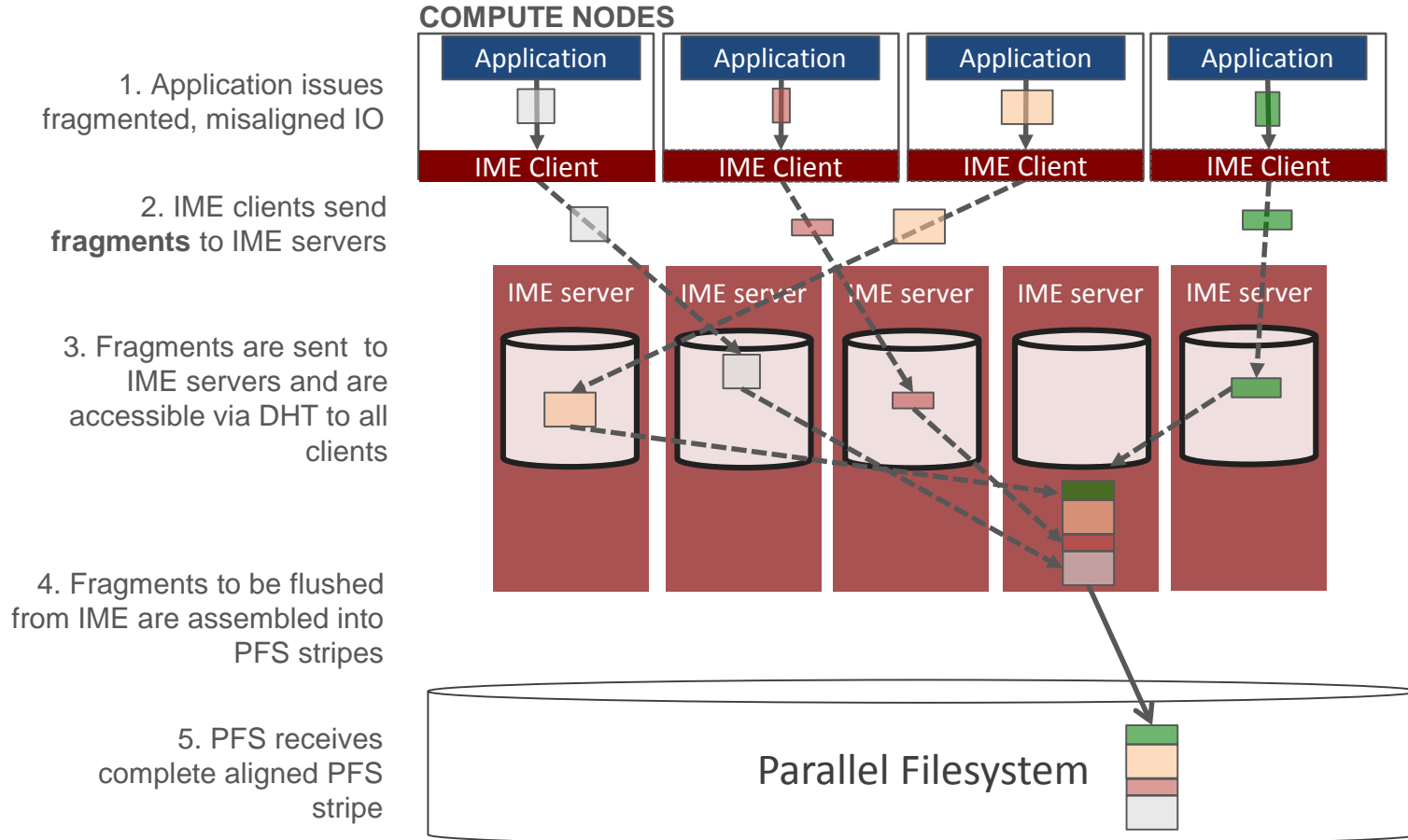IME servers write buffers to NVM and manage internal metadata

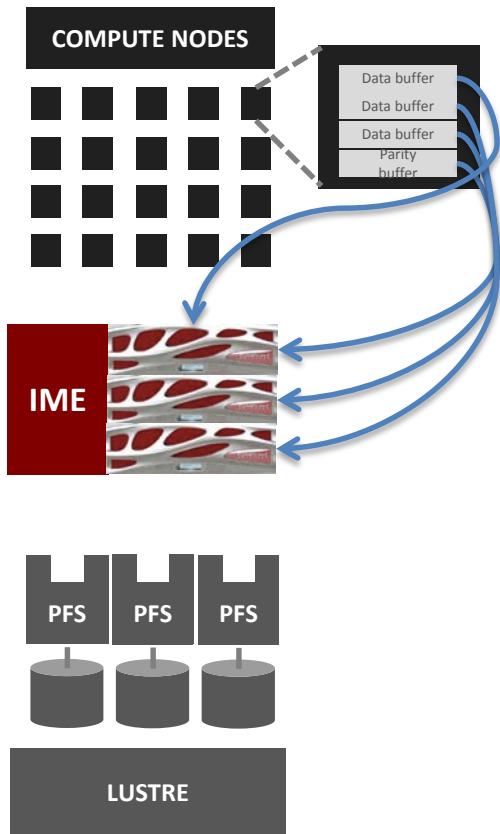IME servers write aligned sequential I/O to SFA backend

Parallel File system operates at maximum efficiency
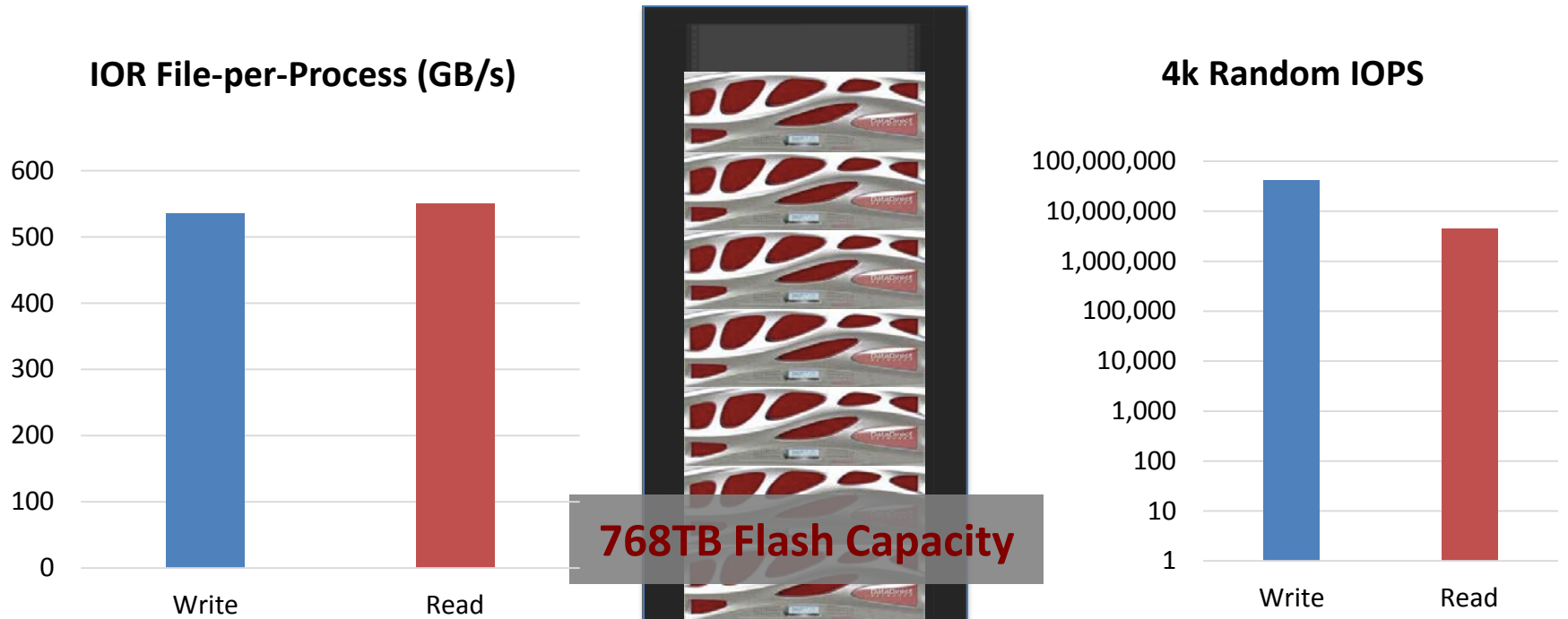
# 4. IME Write Dataflow

# 4. IME Erasure Coding



- Data protection against IME server or SSD Failure is optional
  - (the lost data is "just cache")
- Erasure Coding calculated at the Client
  - Great scaling with extremely high client count
  - Servers don't get clogged up
- Erasure coding does reduce useable Client bandwidth and useable IME capacity:
  - 3+1: 56Gb → 42Gb
  - 5+1: 56Gb → 47Gb
  - 7+1: 56Gb → 49Gb
  - 8+1: 56Gb → 50Gb

# Summary

- SSDs can today be seamlessly introduced into a Lustre Filesystem
  - **Modest** investment in SSDs
  - **Intelligent** policy-driven data moves the most appropriate blocks/files to SSD cache
  - **Block level and Lustre Object** Level data placement schemes
- IME is a ground-up NVM distributed cache which adds
  - **Write** Performance optimisation (not just read)
  - **Small, random** I/O optimisations
  - **Shared** (many-to-one) file optimisations
  - **Improved** SSD lifetime
  - Back-end **Lustre IO optimisation**