



DDN Update Oct, 2016

JLUG2016

Oct, 2016

DataDirect Networks Japan, Inc

Nobu Hashizume

**2016年はExaFlops Systemに向けてDDNが
提唱するStorage SystemのFirst Versionを
日本で提供できた年になりました**



ちなみに去年のJLUGではIMEのコンセプト
のお話しとSFA14Kのご紹介を話していまし
た





#1 IN HPC LEADERSHIP & INNOVATION
NVMe. SSD. FILE SYSTEM. ARCHIVE. CLOUD.

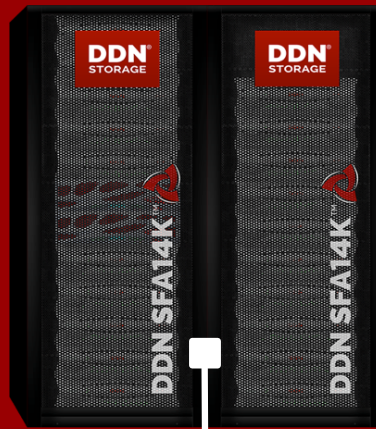
DDN END-TO-END DATA LIFECYCLE MANAGEMENT



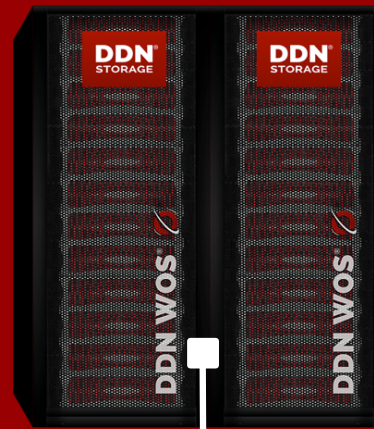
FAST DATA & COMPUTE



PERSISTENT DATA & FILE SYSTEM



LIVE ARCHIVE & COLLABORATION



CLOUD

- SPEED UP YOUR APPS 1000 TIMES!
- ACCELERATE YOUR LUSTRE & GPFS
- WORLD'S FASTEST BURST BUFFER

- HYBRID STORAGE
- EDR IB, 100GB ETHERNET
- MIX SSD AND DISK
- EMBED APPS AND F/S
- 60GB/S & 6M IOPS
- PCI-E FABRIC SPEED

- FASTEST OBJ. STORAGE
- HIGHEST SECURITY
- LOWEST COST ARCHIVE
- MULTI-SITE
- GLOBAL COLLABORATION
- INTELLIGENT TIERING

- HYBRID PUBLIC CLOUD
- GPFS™ AND LUSTRE® BRIDGE
- OPEN SOURCE ACCESS



DDN | EXAScaler & Lustre Case Studies

JCAHPC System



JCAHPC



筑波大学
University of Tsukuba

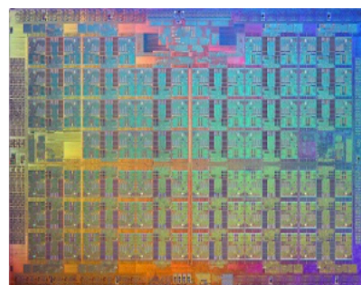


- **University of Tokyo & University of Tsukuba**
- **25 PF System with 8208 KNL Nodes provided by Fujitsu**
- **I/O System by DDN**
 - ▶ **Intel Omnipath**
 - ▶ 26 PB ExaScaler/Lustre @ 400 GB/sec
 - ▶ 1 PB of IME Burst Buffer with NVMe @ 1400 GB/sec

May 10, 2016

Japan Unveils Details of 25 PFLOPS Machine to be Operational in December 2016

John Russell



Knights Landing Die Photo

It's a good day to be Intel, Data Direct Networks (DDN), and Fujitsu. The Joint Center for Advanced High Performance Computing (JCAHPC) in Japan today released the details of its next generation supercomputer – Oakforest-PACs – which will deliver 25 PFLOPS, use Intel's Xeon Phi (Knights Landing) manycore processors and Omni-Path Architecture, be built by Fujitsu, and be operational in December 2016.

When fired up, the Oakforest-PACS will be the fastest supercomputer system in Japan for the moment. Twenty-five

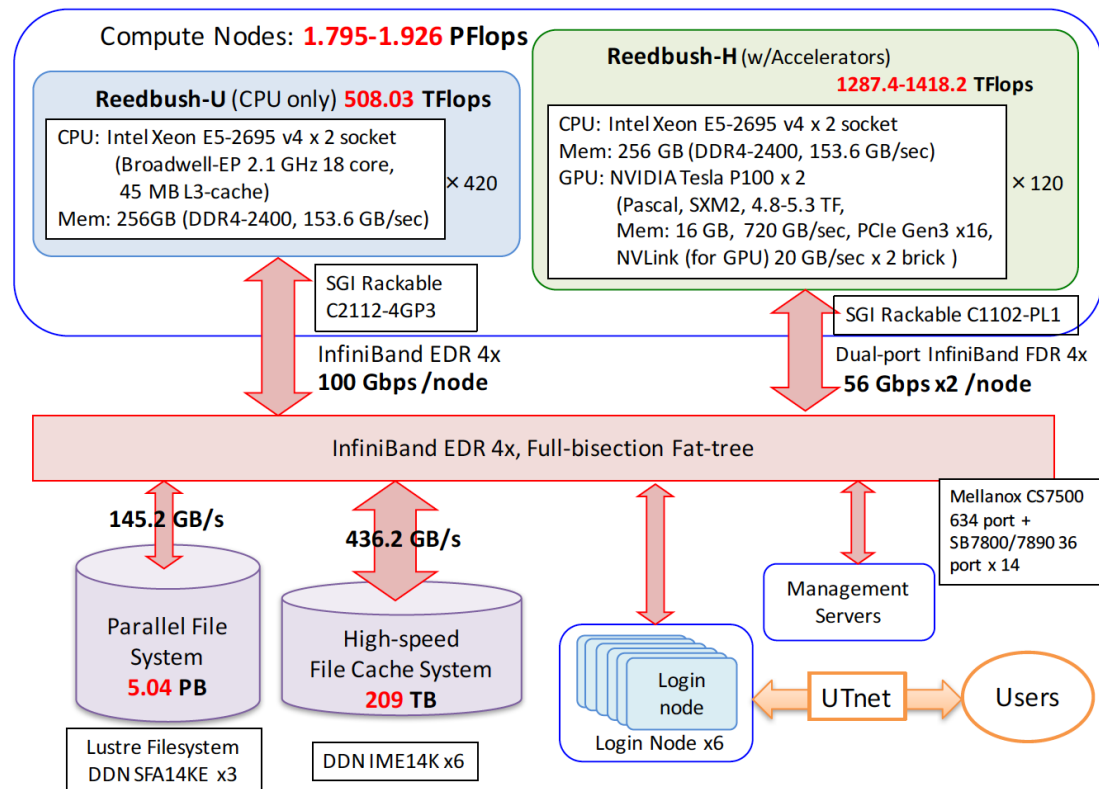
DDN | EXAScaler & Lustre Case Studies

Reedbush Supercomputer System



東京大学情報基盤センタースーパーコンピューティング部門
Supercomputing Division, Information Technology Center
The University of Tokyo

- **University of Tokyo**
- **I/O System by DDN**
 - ▶ 5 PB ExaScaler/Lustre @ 145.2 GB/sec
 - ▶ 200TB of IME Burst Buffer with NVMe @ 436.2 GB/sec



DDN | EXAScaler & Lustre Case Studies

Kyoto University Supercomputer System



京都大学学術情報メディアセンター
Academic Center for Computing and Media Studies, Kyoto University

- Kyoto University
- I/O System by DDN
 - ▶ 24PB ExaScaler/Lustre @ 150 GB/sec
 - ▶ 230TB of IME Burst Buffer with NVMe @ 240 GB/sec

Camphor 2 (System A)

CRAY XC40

intel Xeon Phi KNL 68cores 1.4GHz x 1 /node
#nodes = 1,800
#total cores = 68 cores x 1,800 → 122,400 cores
Peak performance = 3.05TFlops x 1,800 → 5.48 PFlops
Memory capacity = (96+16 GB) x 1,800 → 196.9 TB
Burst buffer = 230 TB, 200 GB/sec DATAWARP

Storage

DataDirect NETWORKS ExaScaler (SFA14K)

Disk capacity = 24 PB
Bandwidth = 150 GB/sec
(Oct. 2016 - Mar. 2018 : 16 PB, 100GB/sec)
Burst buffer = 230 TB, 240 GB/sec IME

高速通信網 InfiniBand EDR/FDR

Laurel 2 (System B)

CRAY CS400 2820XT

intel Xeon Broadwell 18cores 2.1GHz x 2 /node
#nodes = 850
#total cores = 36 cores x 850 → 30,600 cores
Peak performance = 1.21 TFlops x 850 → 1.03 PFlops
Memory capacity = 128 GB x 850 → 106.3 TB

Cinnamon 2 (System C)

CRAY CS400 4840X

intel Xeon Haswell 18cores 2.3GHz x 4 /node
#nodes = 16
#total cores = 72 cores x 16 → 1,152 cores
Peak performance = 2.65 TFlops x 16 → 42.4 TFlops
Memory capacity = 3 TB x 16 → 48.0 TB

高速通信網 Omni-Path

Camellia (System E)

CRAY XC30 with MIC

intel Xeon Ivy Bridge 10cores 2.5GHz x 1 /node
intel Xeon Phi KNC 60cores 1.053GHz, x 1 /node
#nodes = 482
#total cores = (10+60cores) x 482 → 33,740 cores
Peak performance = 1.21 TFlops x 482 → 0.58 PFlops
Memory capacity = (32+8GB) x 482 → 18.8 TB

Storage

DataDirect NETWORKS SFA12K

Disk capacity = 3.0 PB
Bandwidth = 24 GB/sec

高速通信網 InfiniBand FDR/QDR

IMEシステム@Japan

システム	IME	ExaScaler(Lustre)
JCAHPC Oakforest-PACS	物理960TB 25 x IME14K-OPA	物理32PB 10 x SFA14KXE-OPA
東大 Reedbush	物理230TB 6 x IME14K-EDR	物理6.25PB 3 x SFA14KE-EDR
京大 ACCMS2	物理230TB 6 x IME14K-OPA	物理24PB(P1: 16PB, P2: 8PB) 3 x SFA14K-EDR (P1) 2 x SFA14KXE-EDR (P2)
合計	物理1.42PB	物理54.25PB (62.25PB)



WORLD'S FASTEST HYBRID SSD & DISK, EMBEDDED & HYPER-CONVERGED

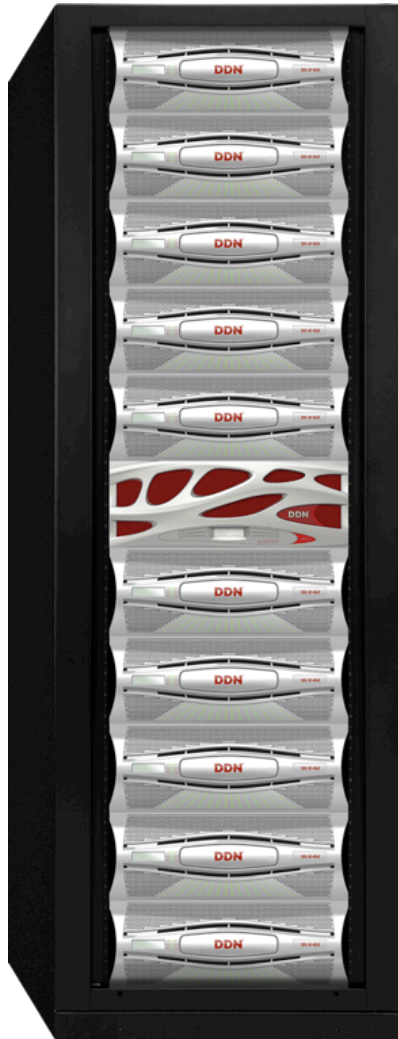
- Delivering 60GB/s and 6 Million IOPS in 4u
- NVMe and SSD Boosts I/O and Shrinks Latency
- Speeds Up SFA[®], IME[®], Lustre[®], GPFS[™]
- Most Advanced Embedded PCI-e Fabric
- EDR IB, 100Gb/s Ethernet, 12Gb/s SAS, Omni-Path

***THE FUTURE OF HPC STORAGE
IS NOW SHIPPING!***



DDN | SFA14K

Design Criteria



Overcome Intel 20GB/s Bottleneck

- No Effective Bandwidth Increases since Sandy Bridge
- PCIe-3 maximum speed has not increased

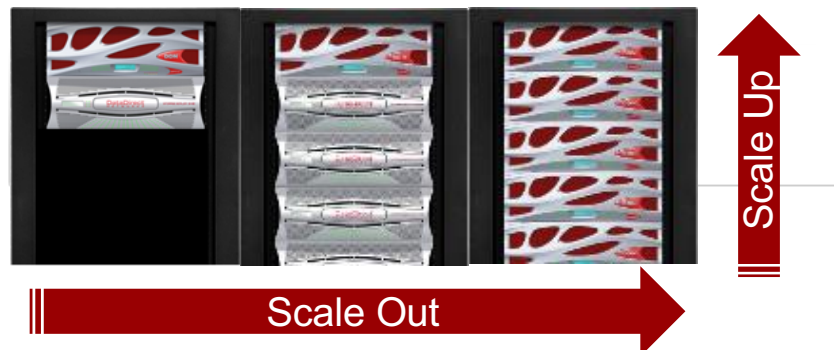
Accommodate New Interconnect Technologies

- Mellanox EDR, Intel Omni-Path, PCIe-3 x16, 100 Gb/s Ethernet, 32 Gb/s FC
- Adaptive Motherboard Design to Encompass Multiple Configurations

Integrate Cutting Edge Solid State Technologies

- 12Gb/s SAS Dual Port, 25W PCIe, NVMe Dual Port
- Most efficient IOPs, \$/capacity, \$/performance

Flexible Configurations for Multiple Use Cases



SFA14K

Most Versatile Next-Generation Storage Platform

SFA Powered By SSD & Disk

Maximize Bandwidth and IOPS Efficiency

- Data analytics acceleration
- Large and small file performance
- Up to 60GB/s bandwidth
- Support for controller clustering

Embedded Solutions

Parallel File System Appliance: [x]Scaler

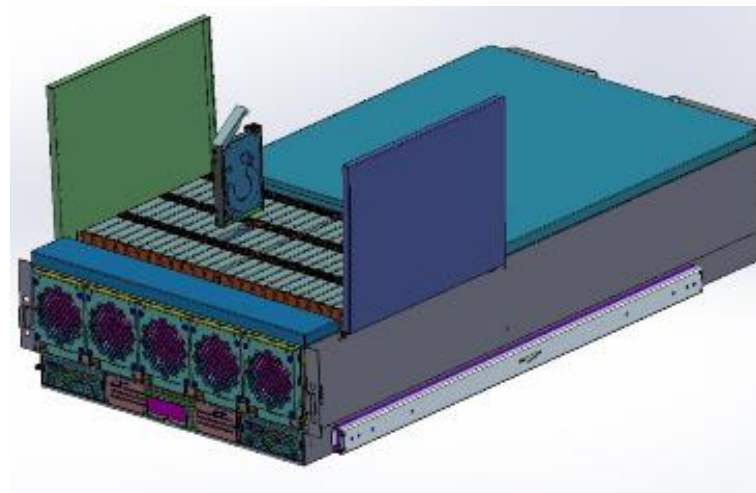
- Run in the controller head
- GPFS or Lustre
- Scale IOPS, throughput, capacity
- Up to 38GB/s bandwidth via file system

IME Application Acceleration

Embedded IME Technology

- Lower latency than flash arrays
- Application acceleration
- IO provisioning
- 50-60GB/s of buffer burst speed

- All in one integrated design with expansion capability
- Dual Redundant Controllers Intel Haswell / Broadwell Processor
- 72 Drive High Density 2.5" Enclosure
- SSDs become an Integral part of the platform
- 48 2.5" dual ported PCIe/SAS drives and 24 2.5" SAS 12 Gb/S dual ported drives
- Optimized Building Block for BW or IOPs
- Direct Connection for 10 – 8412 12Gb/S 84 drive Enclosures



Front/Side Isometric View



Right Side View

Design Specs

Appliance	4RU Active/Active Dual-Controller <ul style="list-style-type: none"> - 48x PCIe NVMe/SAS 2.5" SSD Slots - 24x 12Gb/s SAS 2.5" SSD/HDD Slots - Power Fail-Safe 	
Internal Fabric External Expansion	Switched 600+ lanes PCIe Gen3 24 Quad 12Gb/s SAS	
CPU – Memory Controller	Intel® Xeon® Processor E5 v3 (Haswell) Up to 2048GB DDR4 2133 Memory	
Platform Performance	6M IOPS /4U 60 GB/s /4U	60M IOPS /Rack 600 GB/s /Rack
Host & Network Connectivity	32Gb/s FC EDR IB 100Gig E Intel OmniPath	
Storage Scalability	72 Internal 2.5" & up to 20 direct attached 84 slot 2.5"/3.5" expansion enclosures	
External Clustering	Controller Grid Fabric Interconnect ready	
F/W and Software Options	SFAOS w SFX, Grid- & ExaScaler Embedded, IME, Open Platform	





WORLD'S MOST ADVANCED APPLICATION AWARE I/O ACCELERATION SOFTWARE



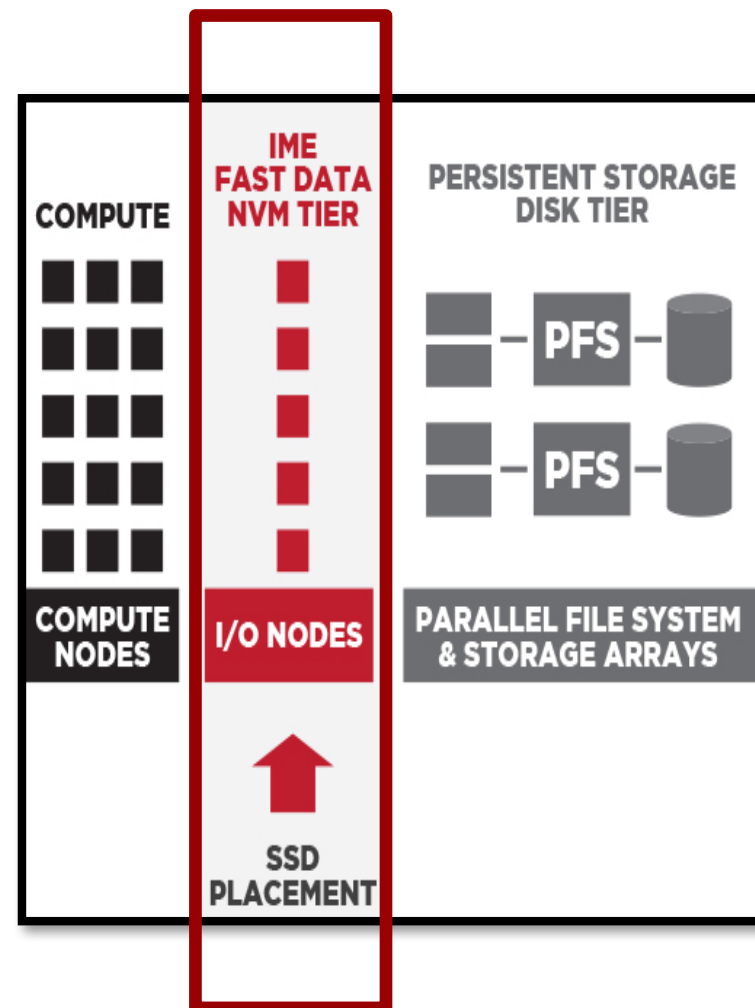
- Accelerate Your Applications Up to 1000 Times!
- Eliminate The Randomness of Your Workflows
- Boost Your Lustre® and GPFS™ File System Speed
- The Fastest and Most Reliable Burst Buffer



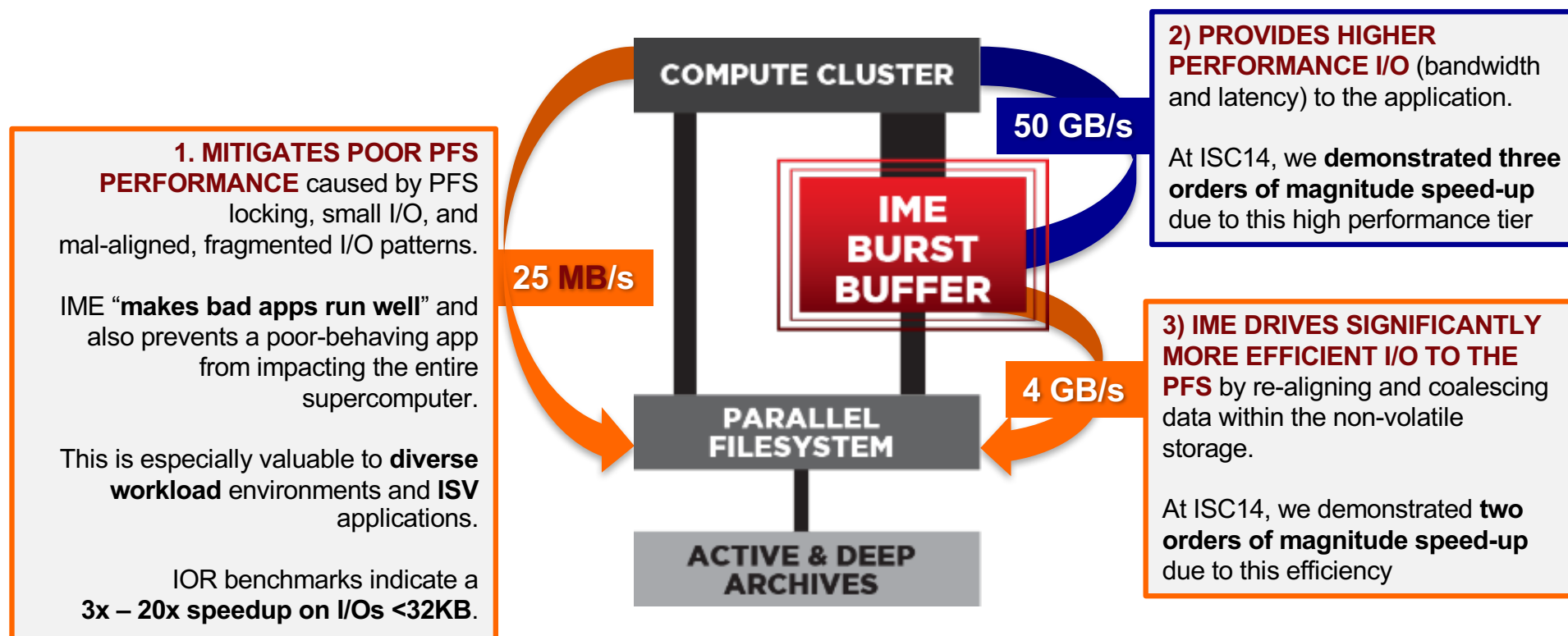
What is IME?



IME creates a new application-aware fast data tier that resides right between compute and the parallel file system to accelerate I/O, reduce latency and provide greater operational and economic efficiency

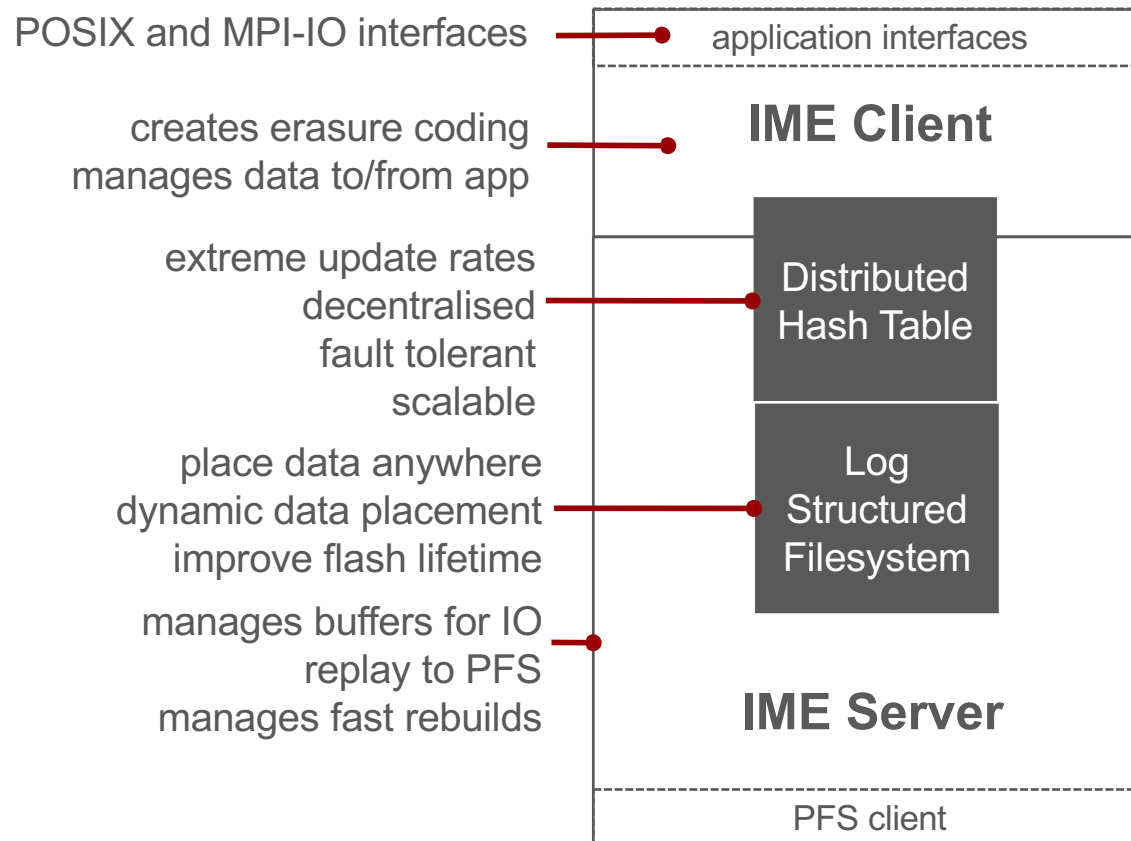


Where IME Provides Value



S3D Turbulent Flow Model

IME Architecture



IME Architecture

POSIX and MPI-IO interfaces

creates erasure coding
manages data to/from app

extreme update rates
decentralised
fault tolerant
scalable

place data anywhere
dynamic data placement
improve flash lifetime

manages buffers for IO
replay to PFS
manages fast rebuilds

application interfaces

IME Client

Distributed
Hash Table

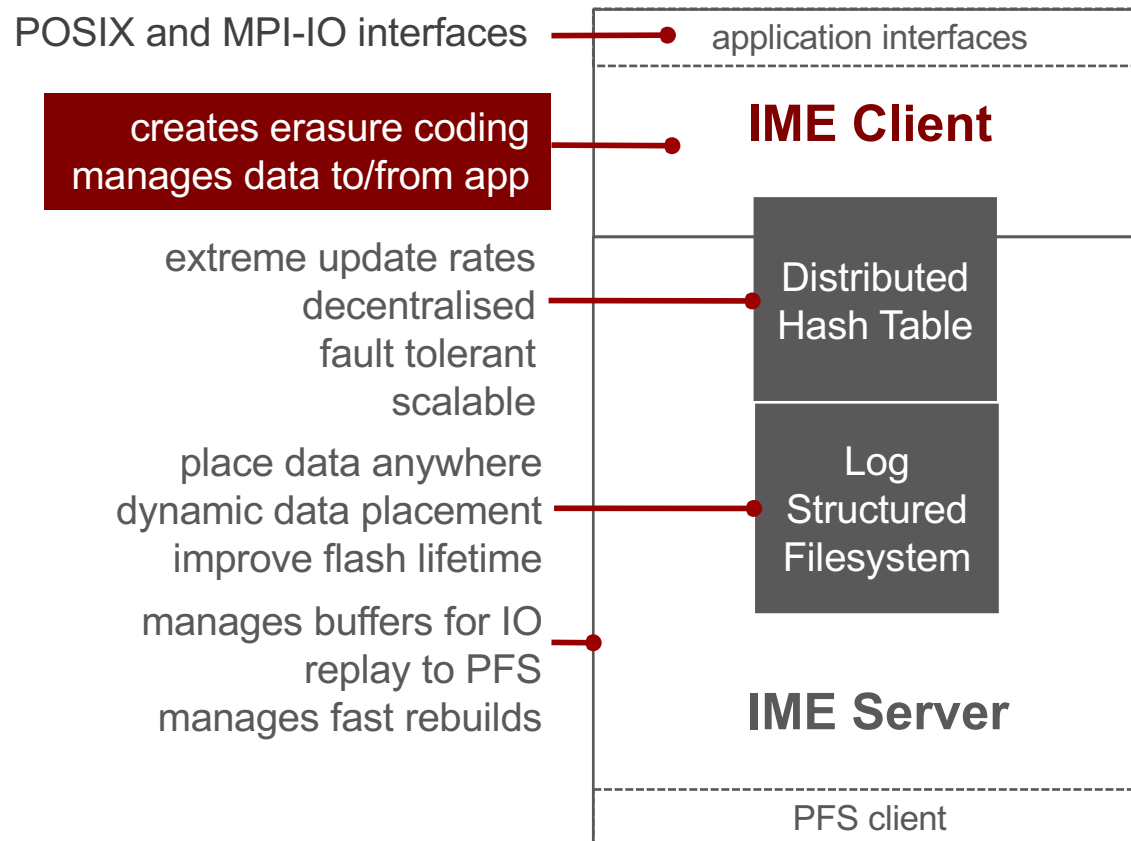
Log
Structured
Filesystem

IME Server

PFS client

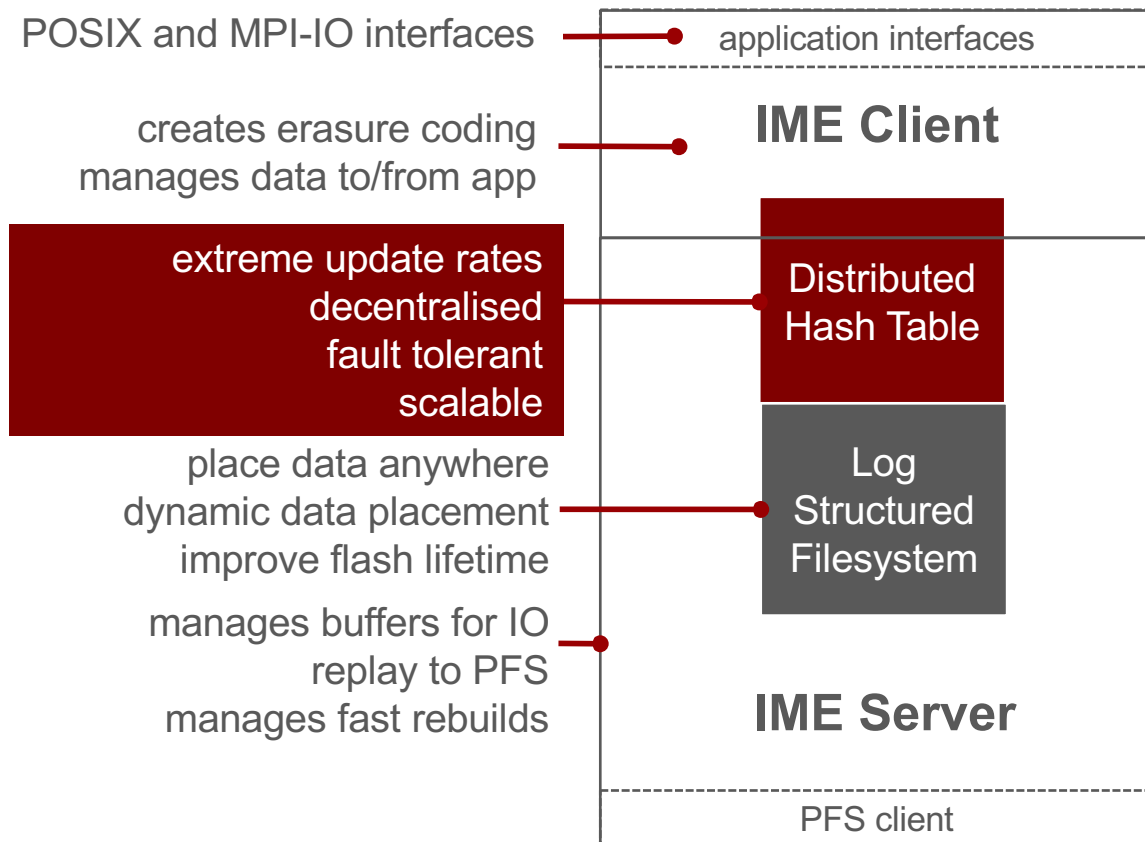
- ▶ No application changes required
- ▶ Applications must be recompiled/relinked to utilise IME

IME Architecture



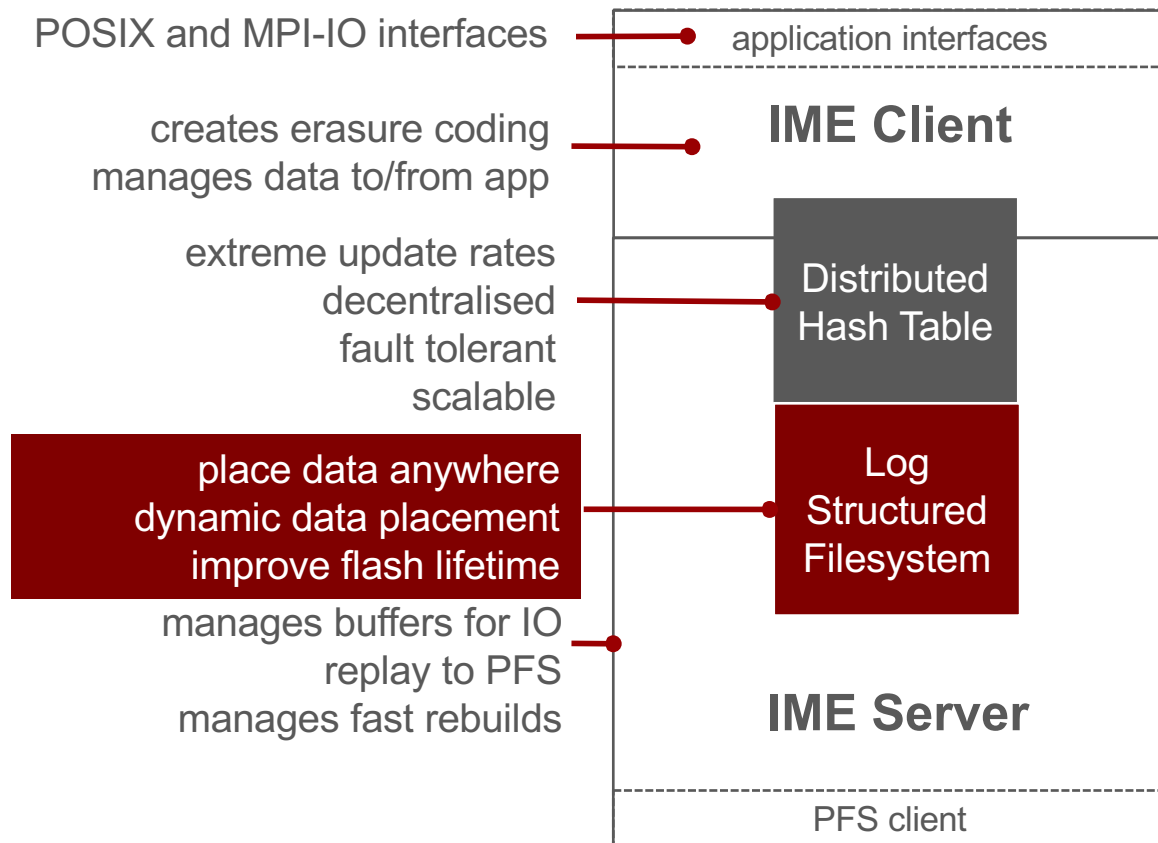
- ▶ **Client buffers IO fragments and sends to IME servers**
- ▶ **also (optionally) creates parity buffers according to RAID scheme chosen**
- ▶ **sends application and data to IME servers such that loss of a server or SSD does not result in data loss**
- ▶ **Heuristics in IME clients can intelligently pre-fetch IME-resident data to applications**

IME Architecture



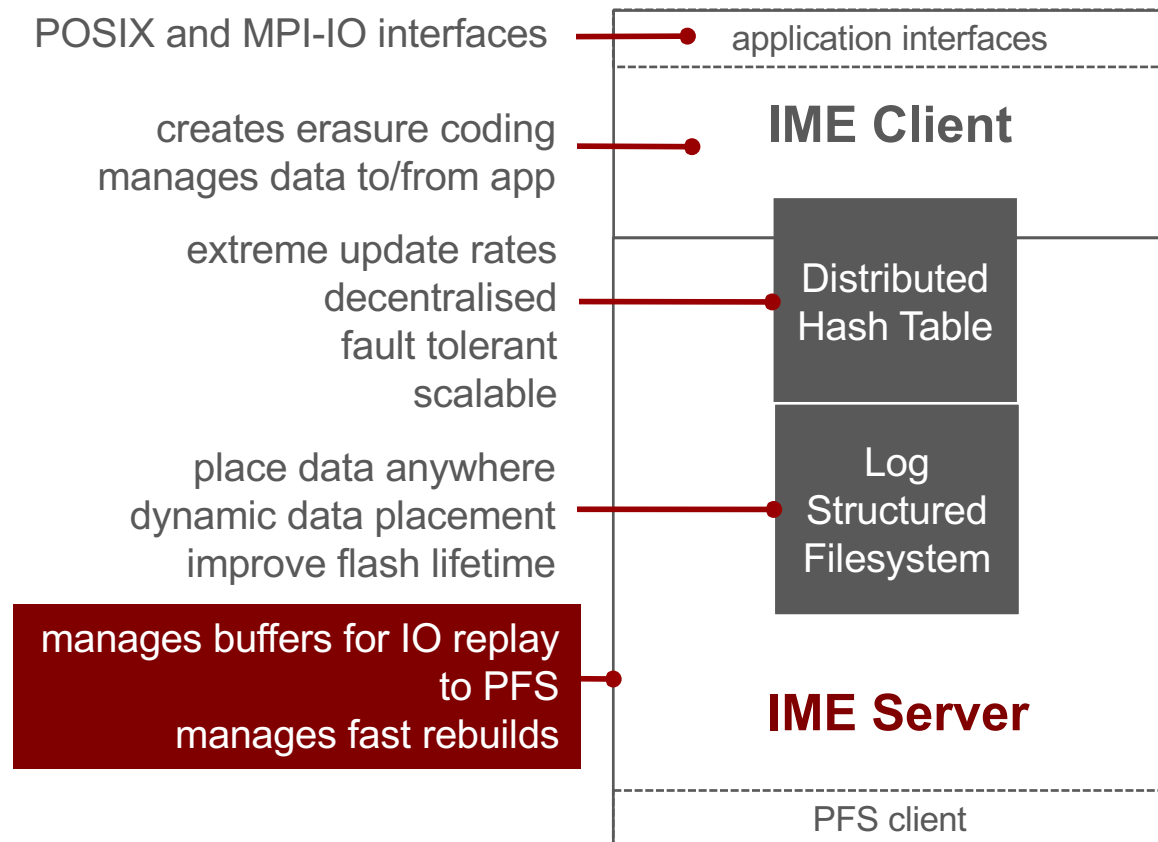
- ▶ **DHT is at the core of IME**
- ▶ **distributed index manages locations of files and objects**
- ▶ **Unique routing and data placement properties differentiate IME DHT from other DHTs**
- ▶ **Fast $O(1)$ routing algorithm built on high-performance, non-cryptographic hash algorithms**
- ▶ **Load is uniformly distributed across DHT nodes**

IME Architecture



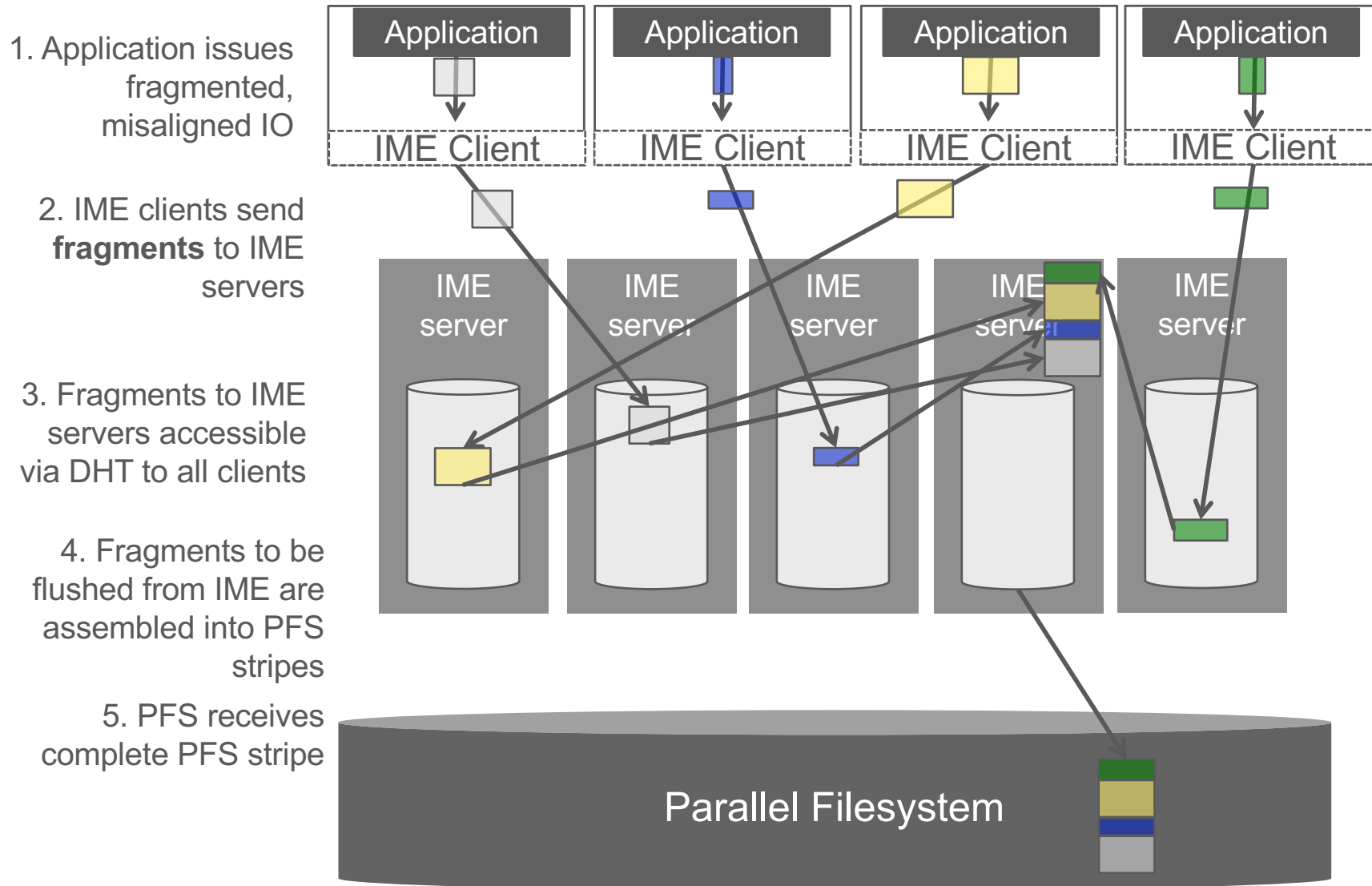
- ▶ **Data can be placed anywhere within IME – any server, any SSD**
- ▶ **IOPs to SSDs are minimised and optimised for NAND Flash**

IME Architecture

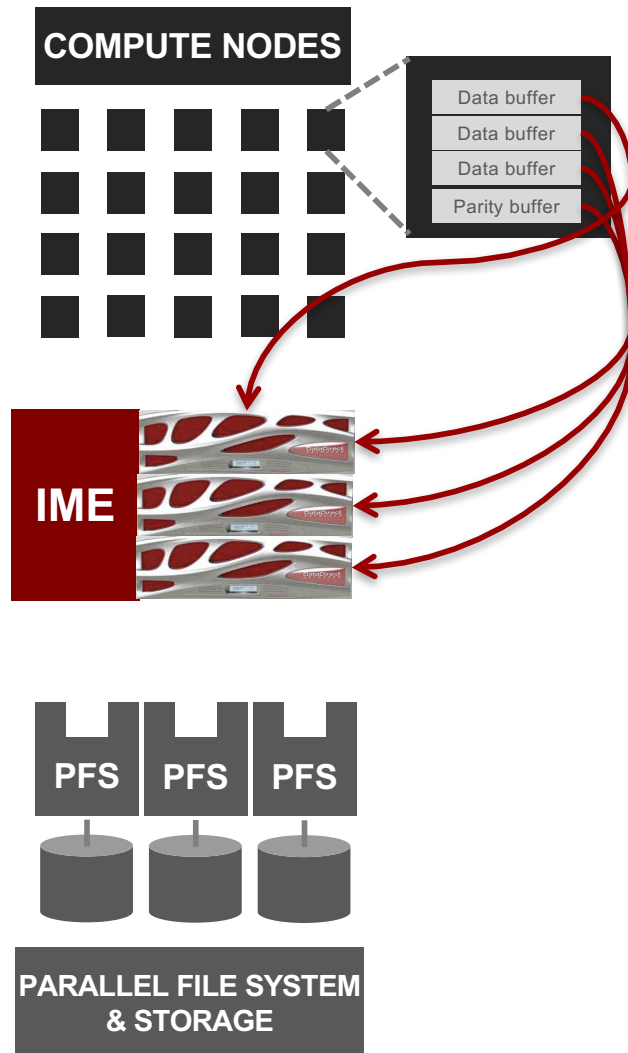


- ▶ **Actively reorganize I/Os to coalesce small block I/Os into large chunks and perform full stripe alignment before submission to the back end, Read acceleration is almost the reverse process**
- ▶ **Protect data by disseminating erasure coded chunks across the IME appliance cluster**
- ▶ **Pre-fetch data into IME from PFS using out-of-band commands and APIs**

IME Dataflow



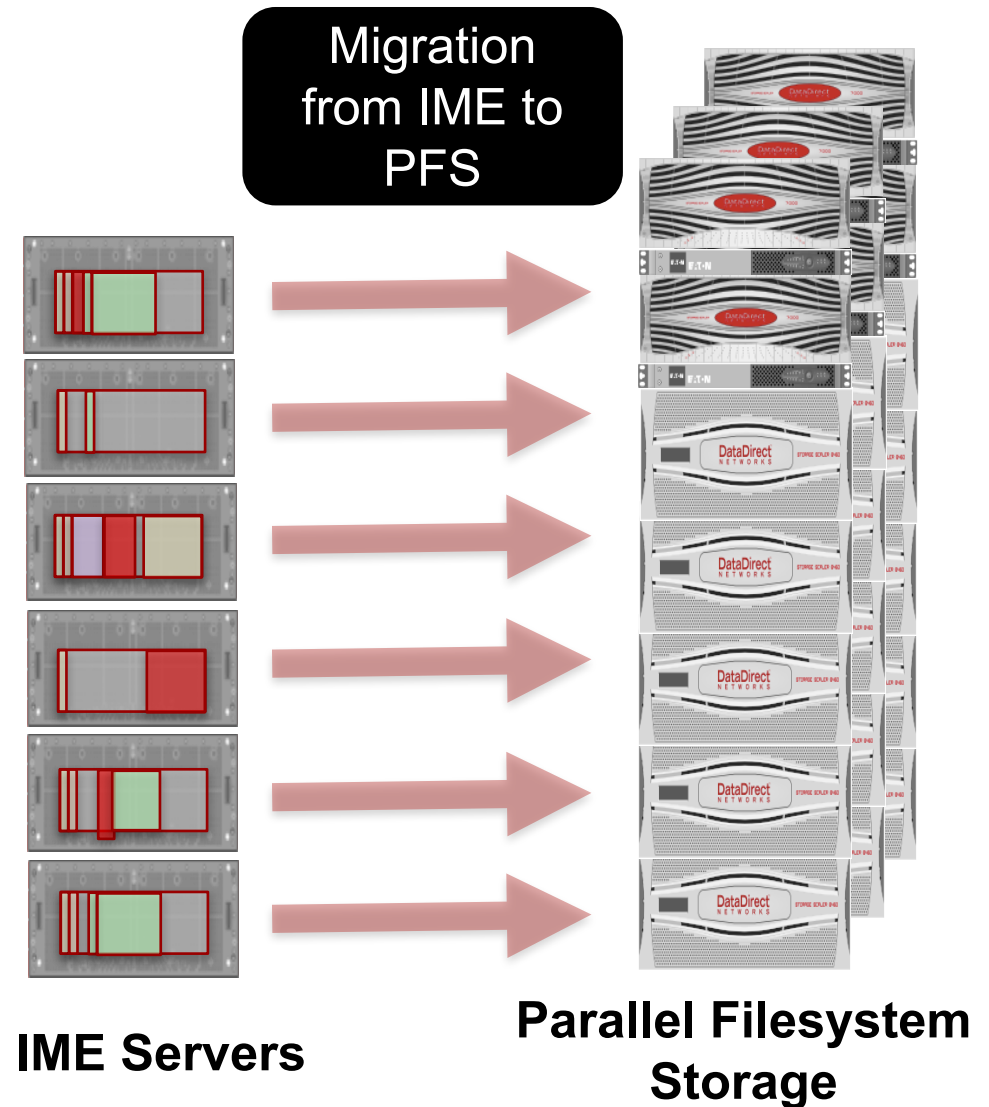
IME Erasure Coding



- ▶ **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

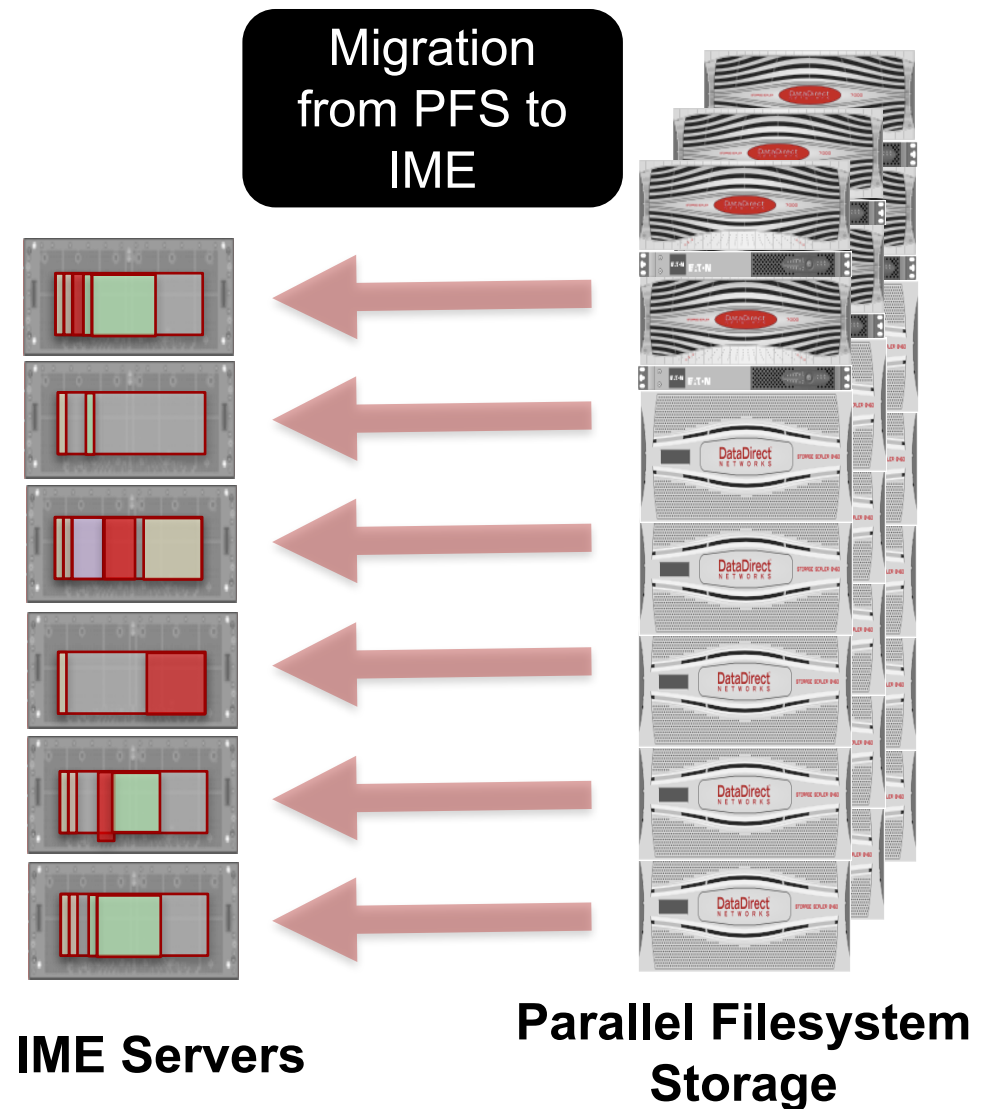
Data Migration Policy – Egress (to PFS)

- **Two modes:**
 - ▶ **Fully dynamic**
 - Data are flushed to PFS as they arrive in IME
 - .. Essentially an LRU
 - ▶ **Explicit**
 - Files are not migrated to the PFS by default
 - Command is required from user or scheduler
 - Shell command interface
 - These files *may* be flushed if IME has reached capacity



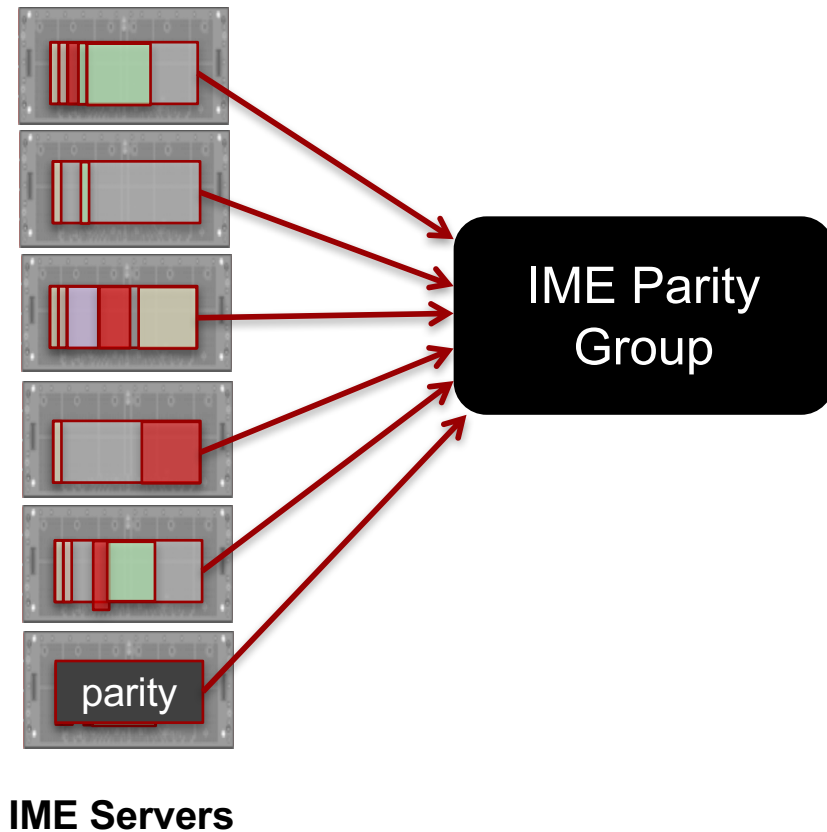
Data Migration Policy – Ingress (from PFS)

- **Fully parallel**
 - ▶ All IME servers participate in ingress procedures
- **Explicit only in IME 1.0**
 - ▶ Shell command interface
 - ▶ Issued by administrator or scheduler
- **Reads of non-cached data are transparently 'passed through to PFS'**
 - ▶ Data is **not** stored on IME SSD



IME Fault Tolerance

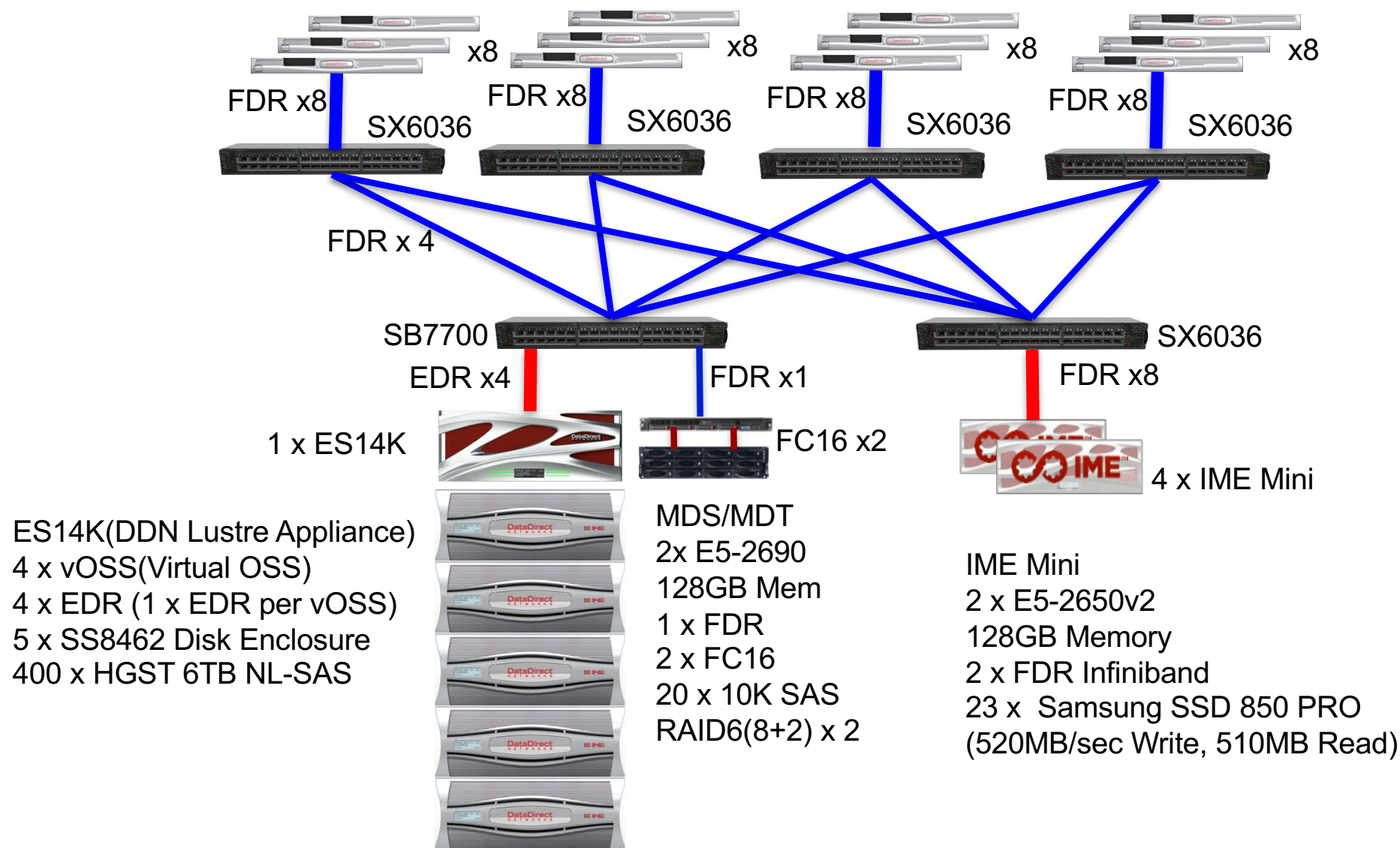
- **Across IME servers**
 - ▶ **no** RAID within the IME server
- **Flexible based on configuration and workload**
 - ▶ N+M schemes
 - ▶ Replication
 - Small files
 - ▶ May be disabled
 - Per-file basis
- **Very High Performance**
 - ▶ Parity group may encapsulate non-contiguous fragments
 - ▶ fully non-blocking
 - ▶ No 'read-modify-writes' or network locking required



27

Benchmark Results

ベンチマーク環境



ハードウェア、ソフトウェア構成

▶ Hardware

- Lustre OSS/OST
 - ES14K (SFA14KE/5 x SS8462, 400 x 6TB NL-SAS)
- Mellanox EDR/FDR Switch (2:1 Blocking Configuration)
- 1 x Lustre MDS/MDT
 - 1 x SuperMicro Server(2 x E5-2690, 128GB Memory, 1 x FDR)
 - 1 x EF4024(FC16, 10Krpm 600GB SAS x 20, RAID6(8+2) x 2)
- 32 x Lustre Client
 - 16 x Dell R630 (2 x E5-2650 v2, 64GB Memory, 1x FDR)
 - 16 x Dell R620 (2 x E5-2640 v3, 64GB Memory, 1 x FDR)

▶ Software Stack

- ES3.0 Beta(Lustre Server)
- CentOS6.7 (Client)
- Mellanox OFED-3.2-2.0.0
- Lustre-2.5.41.ddn11

ベンチマーク概要

▶ NeRSCの調達ベンチマークルールを使用

- <http://www.nersc.gov/users/computational-systems/cori/nersc-8-procurement/trinity-nersc-8-rfp/nersc-8-trinity-benchmarks>

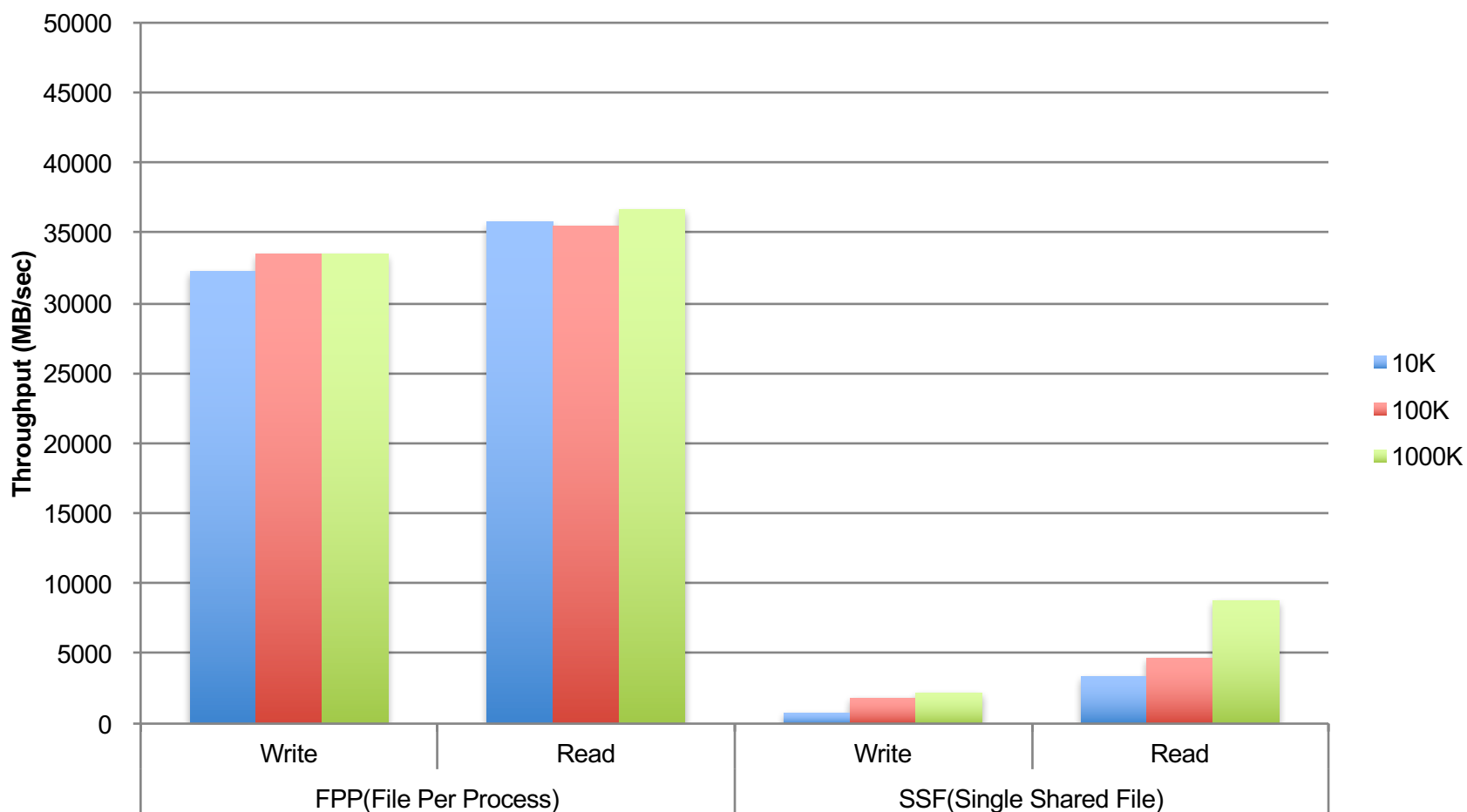
▶ IOR

- POSIX I/O file per process(FPP)
- POSIX I/O single shared file(SSF)
- MPI/IO file per process(FPP)
- MPI/IO single shared file(SSF)

IOR (POSIX, Lustre)

32 clients, 512 process, 3.3TB File Size

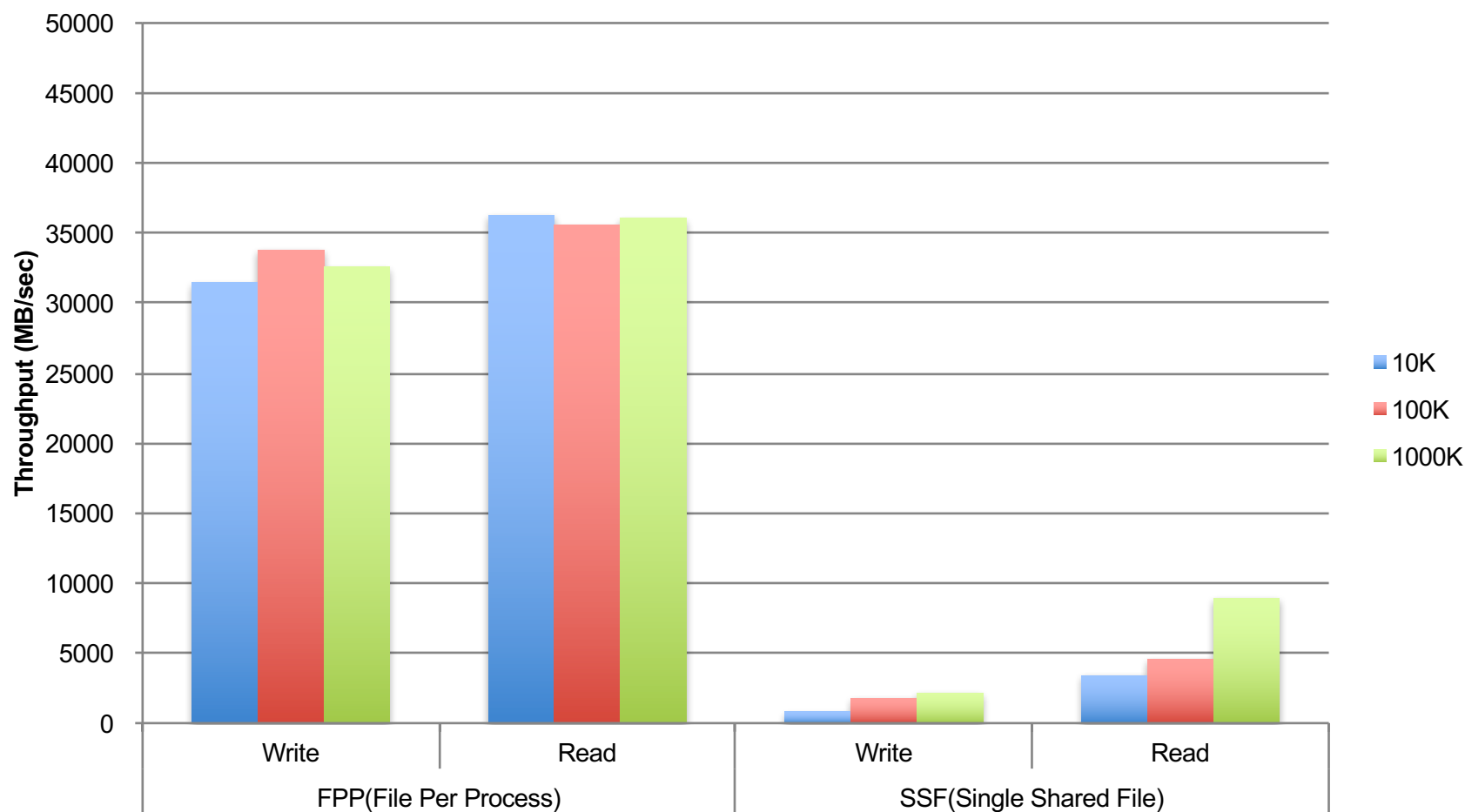
IOR(Lustre, POSIX, SFA14KE, 400 x NL-SAS)



IOR (MPIIO, Lustre)

32 clients, 512 process, 3.3TB File Size

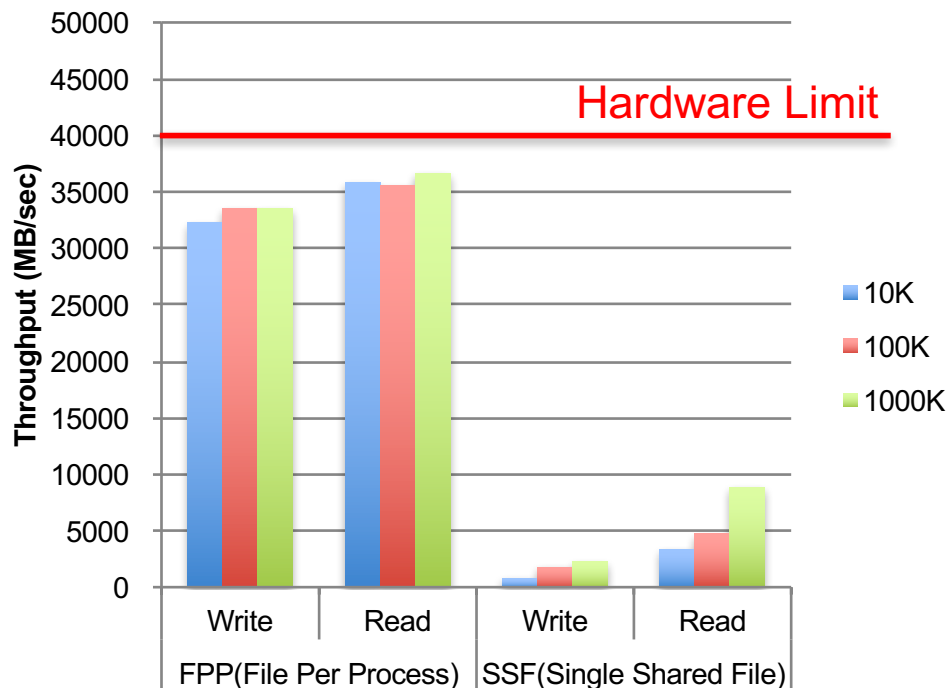
IOR(Lustre, MPIIO, SFA14KE, 400 x NL-SAS)



IOR (POSIX, Lustre and IME)

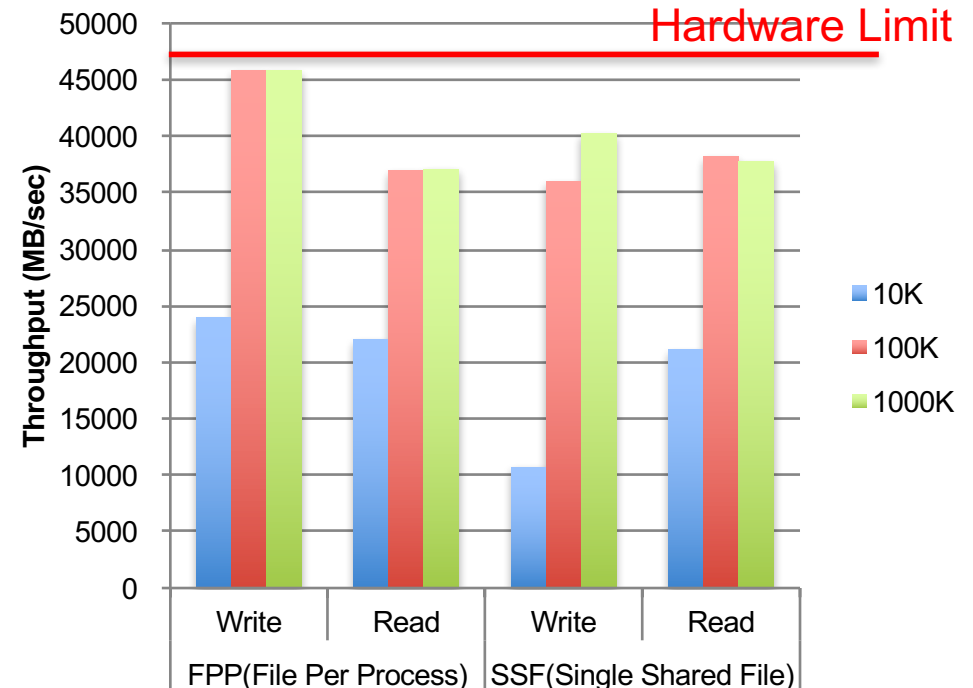
32 clients, 512 process, 3.3TB File Size

IOR(Lustre, POSIX)



FPP I/O Efficiency ~84%(Write) ~90%(Read)
 SSP I/O Efficiency ~5%(Write) ~22%(Read)

IOR(IME, POSIX)

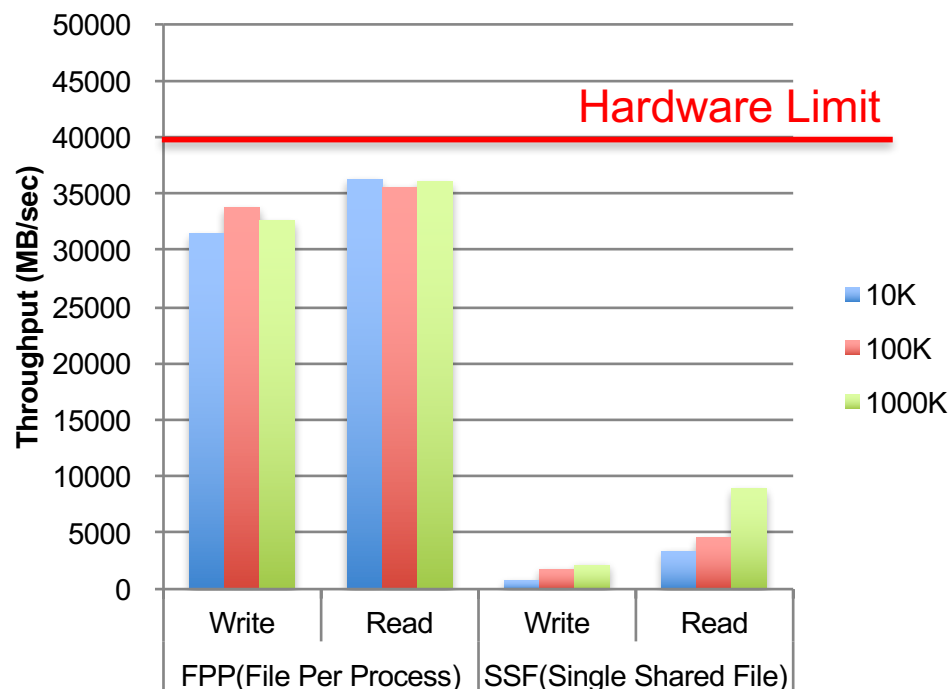


FPP I/O Efficiency ~97%(Write) ~78%(Read)
 SSP I/O Efficiency ~85%(Write) ~81%(Read)
 Still under optimizations

IOR (MPIIO, Lustre and IME)

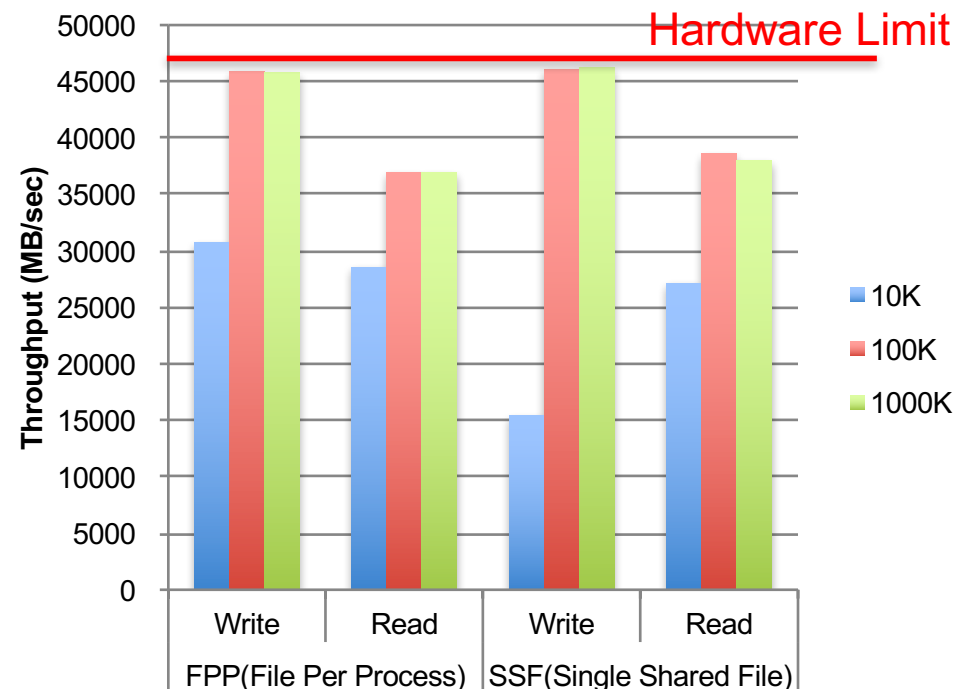
32 clients, 512 process, 3.3TB File Size

IOR(Lustre, MPIIO)



FPP I/O Efficiency ~84%(Write) ~90%(Read)
 SSP I/O Efficiency ~5%(Write) ~22%(Read)

IOR(IME, MPIIO)



FPP I/O Efficiency ~97%(Write) ~78%(Read)
 SSP I/O Efficiency ~97%(Write) ~81%(Read)
 Still under optimizations

Thank You!

Keep in touch with us



sales@ddn.com



2929 Patrick Henry Drive
Santa Clara, CA 95054



[@ddn_limitless](https://twitter.com/ddn_limitless)



1.800.837.2298
1.818.700.4000



[company/datadirect-networks](https://www.linkedin.com/company/datadirect-networks)

DataDirect[™]
NETWORKS

© 2014 DataDirect Networks, Inc. * Other names and brands may be claimed as the property of others.
Any statements or representations around future events are subject to change.

ddn.com