# DDN Update for 2017

PCCC

**DataDirect Networks Japan, Inc**

橋爪信明

Dec, 2016

# Why IME?
cater better for real-world IO patterns

IO Benchmarking
for acceptance

ddn.com

# Why IME?
cater better for real-world IO patterns



All the time after that

4

ddn.com

# DDN のFlash(IME)に対する取り組み

従来のHPC向け並列ファイルシステムでは解決出来ないIOワークロード（Big Data, AIなど）に向けてFlashを活用した新しいIOシステムを提案

- **2014 : Infinite Memory Engine(IME)発表**

- **2015 : IMEを大規模スパコンに提案**

- **2016 : IMEを国内3サイトに提供**

- **2017 : 小さいIMEを展開する予定**

ddn.com

# DDN | IME
## Application I/O Workflow

**Compute**

**∞ IME™**

**SFA**

**Diverse, high concurrency applications**

**Fast Data NVM & SSD**

**Persistent Data (Disk)**

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

**DDN™ STORAGE**

ddn.com

# DDN | EXAScaler & Lustre Case Studies
# JCAHPC System

DDN STORAGE | 東京大学 THE UNIVERSITY OF TOKYO | JCAHPC | 筑波大学 University of Tsukuba | (intel)
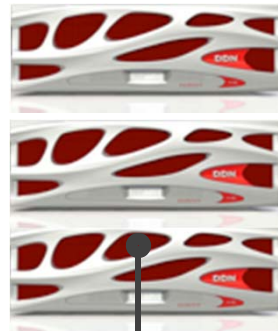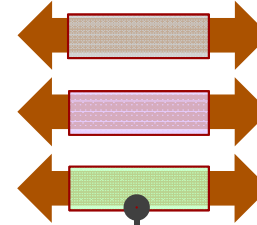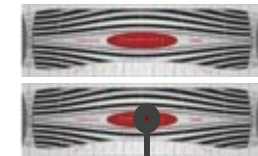
May 10, 2016

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

John Russell

- **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

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

Knights Landing Die Photo

DDN.COM

# DDN | EXAScaler & Lustre Case Studies Reedbush Supercomputer System

**DDN® STORAGE**

**iTC** 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
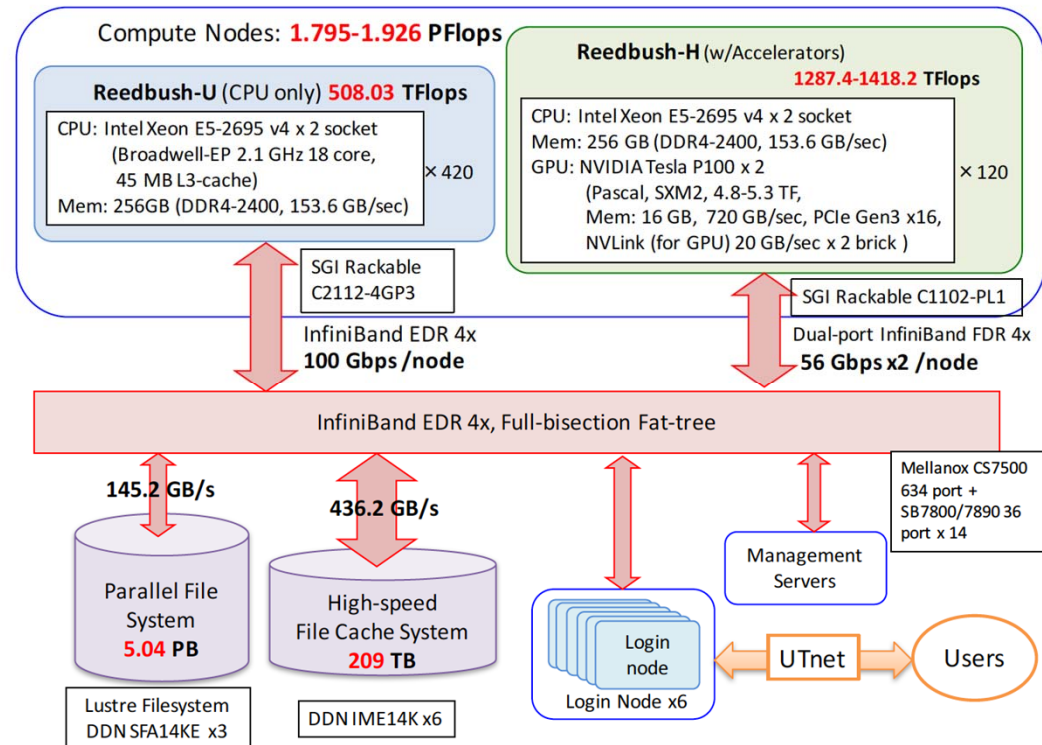
- **University of Tokyo**

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

Compute Nodes: **1.795-1.926** PFlops

**Reedbush-U** (CPU only) **508.03** TFlops

CPU: Intel Xeon E5-2695 v4 x 2 socket
(Broadwell-EP 2.1 GHz 18 core,
45 MB L3-cache)
Mem: 256GB (DDR4-2400, 153.6 GB/sec)

× 420

SGI Rackable C2112-4GP3

InfiniBand EDR 4x
**100 Gbps /node**

**Reedbush-H** (w/Accelerators)
**1287.4-1418.2** TFlops

CPU: Intel Xeon E5-2695 v4 x 2 socket
Mem: 256 GB (DDR4-2400, 153.6 GB/sec)
GPU: NVIDIA Tesla P100 x 2
(Pascal, SXM2, 4.8-5.3 TF,
Mem: 16 GB, 720 GB/sec, PCIe Gen3 x16,
NVLink (for GPU) 20 GB/sec x 2 brick )

× 120

SGI Rackable C1102-PL1

Dual-port InfiniBand FDR 4x
**56 Gbps x2 /node**

InfiniBand EDR 4x, Full-bisection Fat-tree

**145.2 GB/s**

**436.2 GB/s**

Mellanox CS7500
634 port +
SB7800/7890 36
port x 14

Parallel File System
**5.04 PB**

High-speed File Cache System
**209 TB**

Login node

Management Servers

UTnet

Users

Lustre Filesystem
DDN SFA14KE x3

DDN IME14K x6

Login Node x6

DDN.COM

# DDN | EXAScaler & Lustre Case Studies
# Kyoto University Supercomputer System



- **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 = 3.05TFops 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
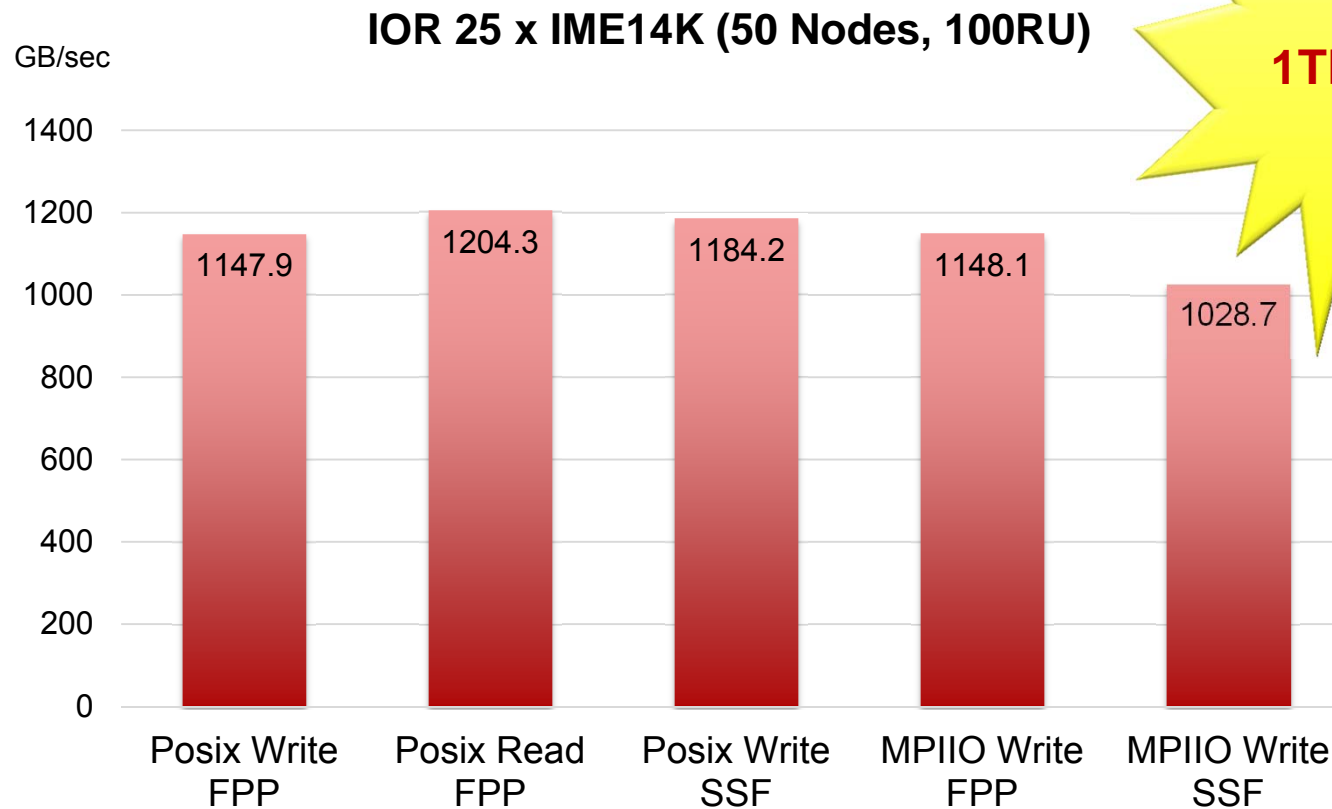- Memory capacity = (32+8GB) x 482 ➜ 18.8 TB

**Storage**

**DataDirect NETWORKS  SFA12K**
- Disk capacity = 3.0 PB
- Bandwidth = 24 GB/sec

高速通信網 **InfiniBand FDR/QDR**

**DDN**  DDN.COM

# IMEシステム@Japan

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

DDN STORAGE

ddn.com

# IME Performance, OPA, IOR on OFP



**IOR 25 x IME14K (50 Nodes, 100RU)**

GB/sec

1TB/sec

Bar chart values:
- Posix Write FPP: 1147.9
- Posix Read FPP: 1204.3
- Posix Write SSF: 1184.2
- MPIIO Write FPP: 1148.1
- MPIIO Write SSF: 1028.7

Y-axis: 0, 200, 400, 600, 800, 1000, 1200, 1400

FPP : File Per Process
SSF : Single Shared File

ddn.com

# FPP vs SSF

- **FPP(File Per Process)**
  - プロセス毎に独立したファイルを使用
  - メリット：並列ファイルシステムが得意、スケールしやすい
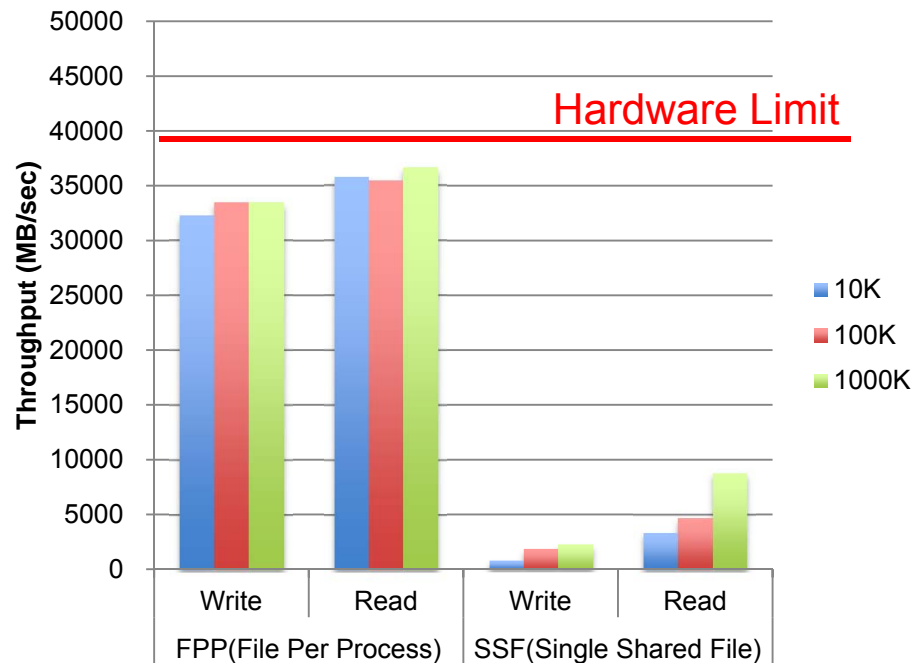  - デメリット：超並列（何万・何十万プロセス）実行時など膨大なファイルによってメタデータ性能ネックになる
- **SSF(Single Shared File)**
  - 全てのプロセスで単一ファイルを共有
  - HDF5, NetCDFなどを利用
  - メリット：メタデータ性能ネックにならない
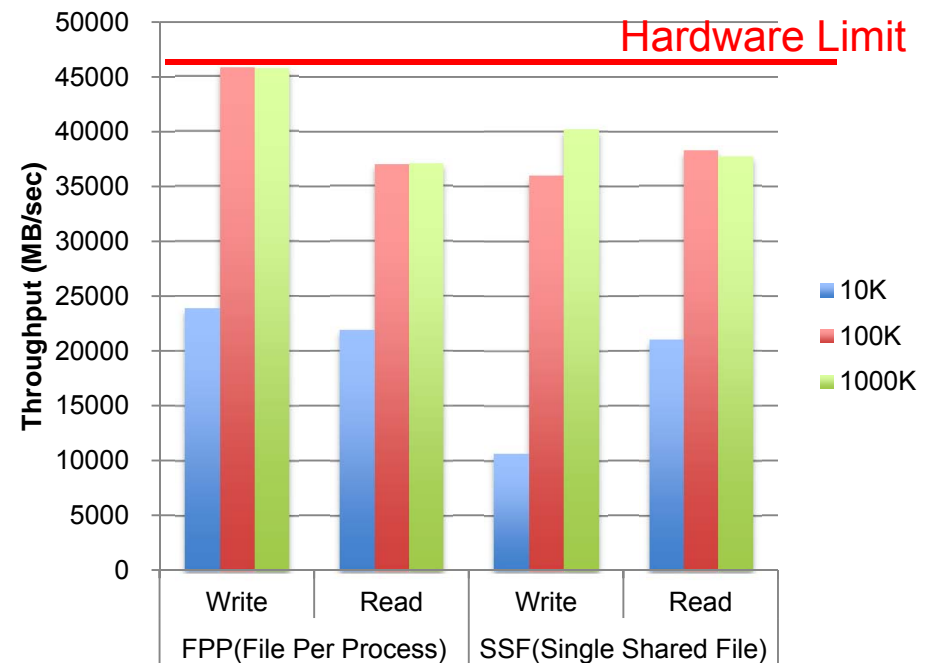  - デメリット：FPPより場合によっては複雑

ddn.com

# Lustre vs IME (IOR POSIX)
# 32 clients, 512 process, 3.3TB File Size



IOR(Lustre, POSIX)

IOR(IME, POSIX)

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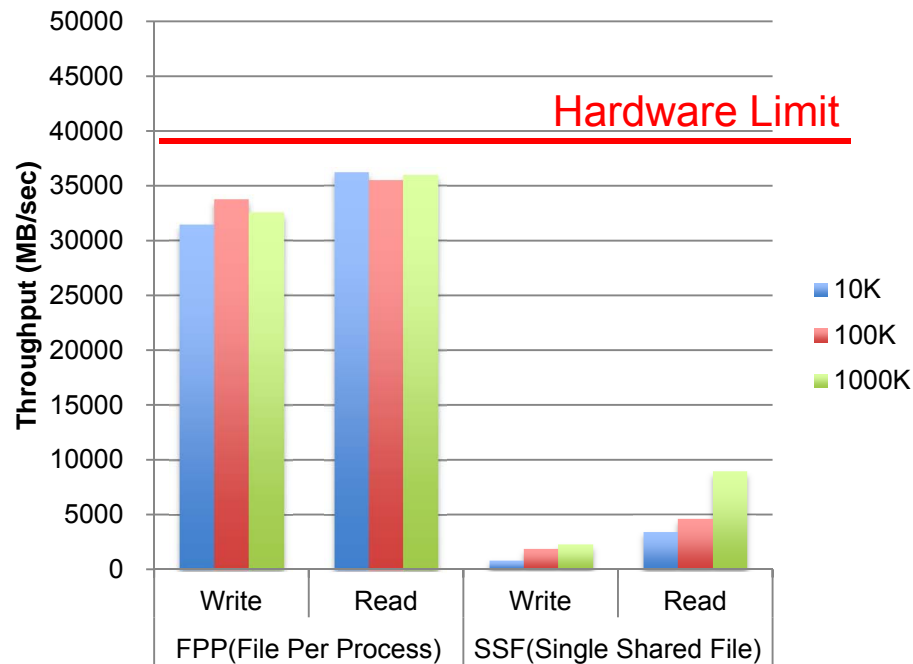SSP I/O Efficiency ~5%(Write) ~22%(Read)

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

ddn.com

DDN STORAGE
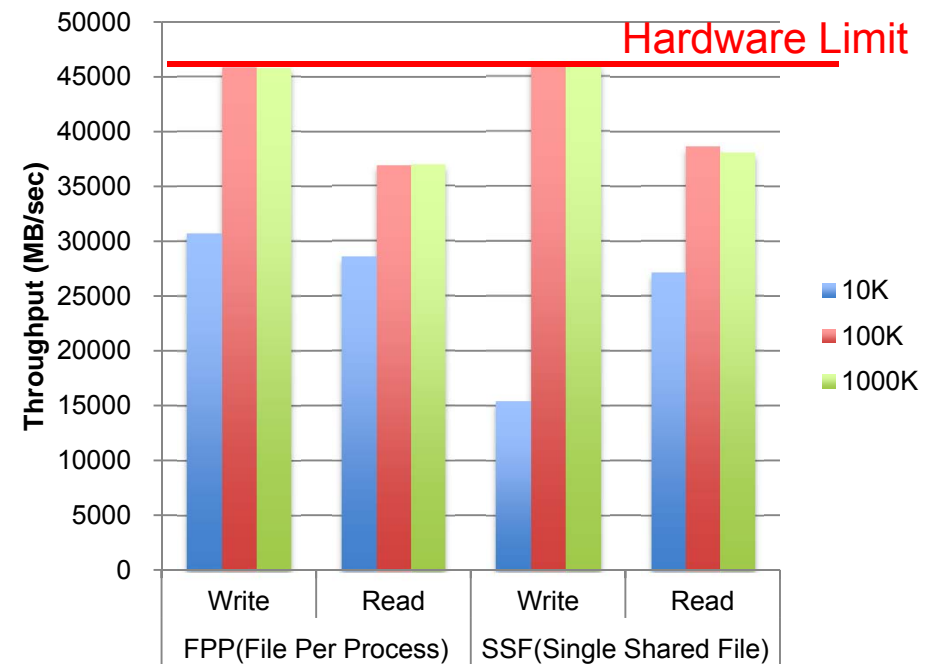
# Lustre vs IME (IOR MPI-IO)
# 32 clients, 512 process, 3.3TB File Size

**IOR(Lustre, MPIIO)**



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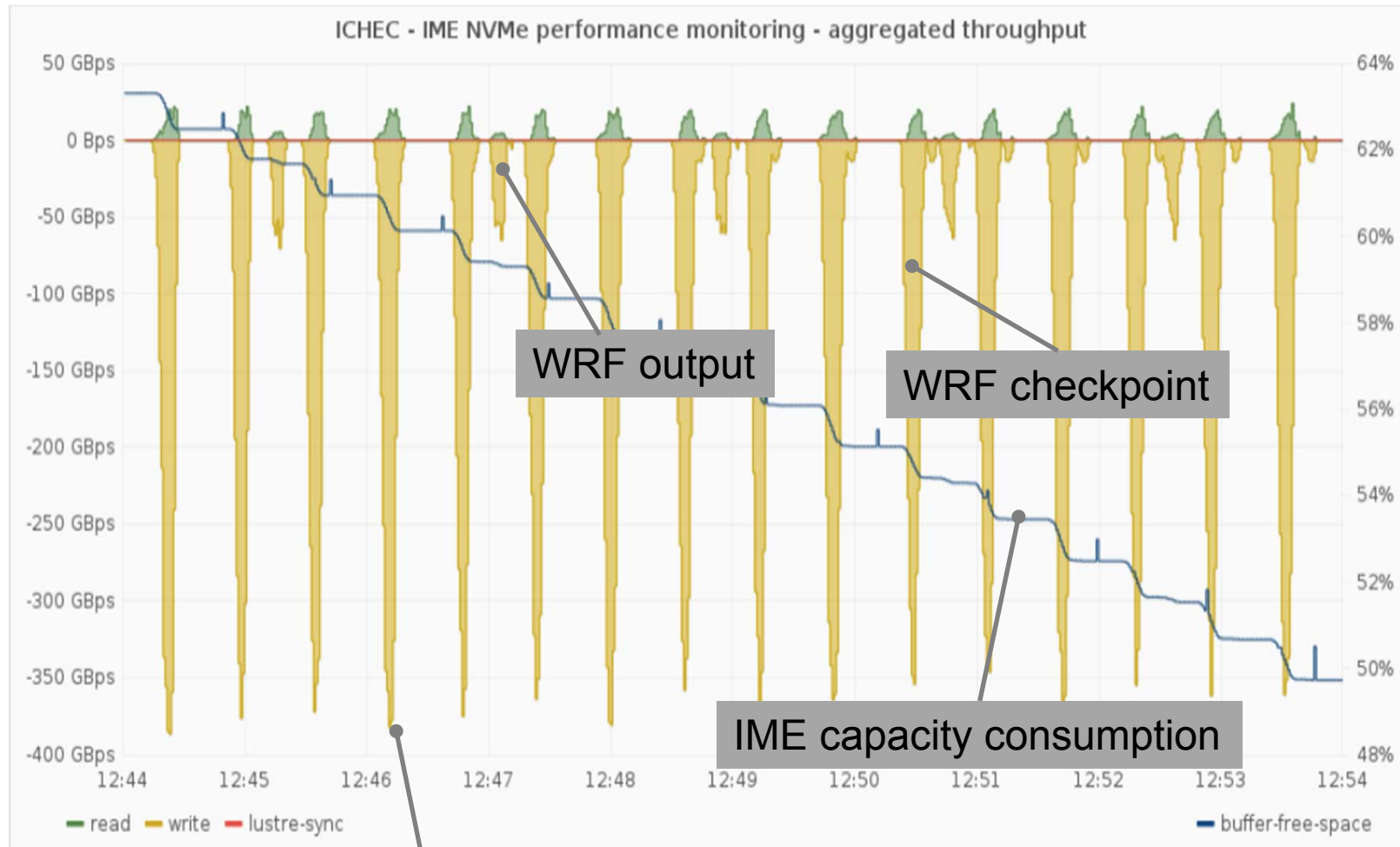SSP I/O Efficiency ~5%(Write) ~22%(Read)

**IOR(IME, MPIIO)**



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

ddn.com

**DDN STORAGE**

# WRF on IME

48 job, total 240 compute node – throughput – focus, 5 minute window



ICHEC - IME NVMe performance monitoring - aggregated throughput

Labels on chart: WRF output, WRF checkpoint, IME capacity consumption, ~390GB/s

Legend: read, write, lustre-sync, buffer-free-space

# WRF on Lustre

48 job, total 240 compute node – throughput – focus, 5 minute window



ICHEC - Lustre OST obdfilter - aggregated throughput

~100GB/s

ddn.com

# WRF on IME at scale
## Application runtime speedup over 1.5x

ddn.com

# DDN | DDN IME Platforms

‣ **IME120 (SuperMicro Platform)**
  • Runs IME server software
  • 1P Broadwell, 1x EDR/OPA 1RU
  • 6 2.5" NVMe SSDs
  • 64-128 GB DRAM @2400
  • 4.6-12TB IME Capacity
  • Max 10GB/s throughput per server

‣ **IME240 (SuperMicro Platform)**
  • Runs IME server software
  • 2P Broadwell, 2x EDR/OPA, 2 RU
  • 20 2.5" NVMe SSDs
  • 128 – 256 GB DRAM @2400
  • 9.6-40TB IME Capacity
  • Max 20 GB/s throughput per server

‣ **IME14K Appliance**
  • Two controllers in 4 RU
  • 2P Haswell / Broadwell, EDR/OPA
  • 48 NVMe SSDs (V1.0.0) and 72 SAS SSDs (1H'16)
  • 38.4-86.4TB IME Capacity
  • 50 GB/s raw throughput per appliance

*Note: Erasure coding topology will impact achievable peak performance!*

ddn.com

# Thank You!

Keep in touch with us

sales@ddn.com

@ddn_limitless

company/datadirect-networks

2929 Patrick Henry Drive
Santa Clara, CA 95054

1.800.837.2298
1.818.700.4000

**DataDirect**™
NETWORKS

ddn.com