# **Infinite Memory Engine (IME)**

# Partial Non-Deterministic I/O System for Exascale

Jean-Thomas Acquaviva

Teratec Forum Palaiseau, France

June, 2015

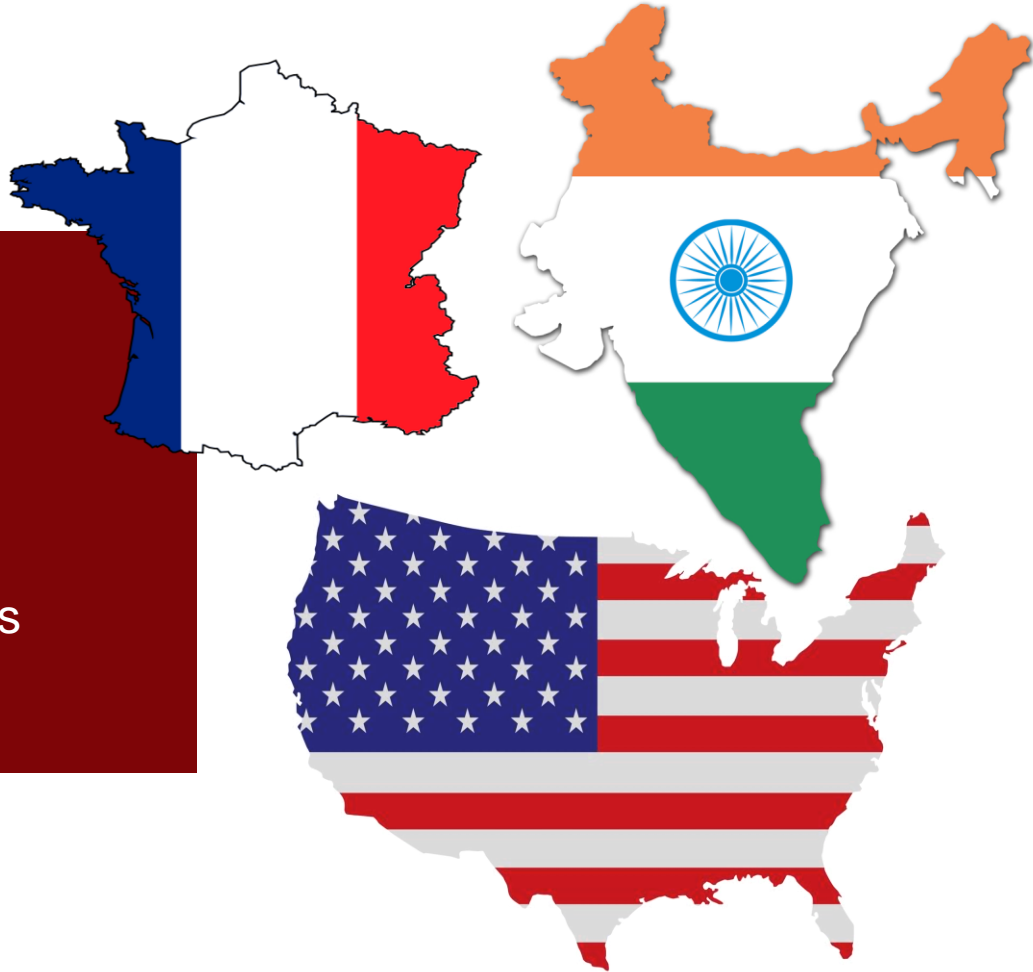ddn.com

# Who's IME?

**Lead Architect: Paul Nowoczynski**

**R&D team on 3 continents**
ₗMaryland, US
ₗParis, France
ₗPune, India
First project developed in the new Paris Advanced Technical Center

# Why IME?
## Bandwidth coarse estimation

**Bandwidth needs next-gen pre-Exascale systems**

**Rules of thumb**

**1/ Checkpointing less than 6 minutes per hours**
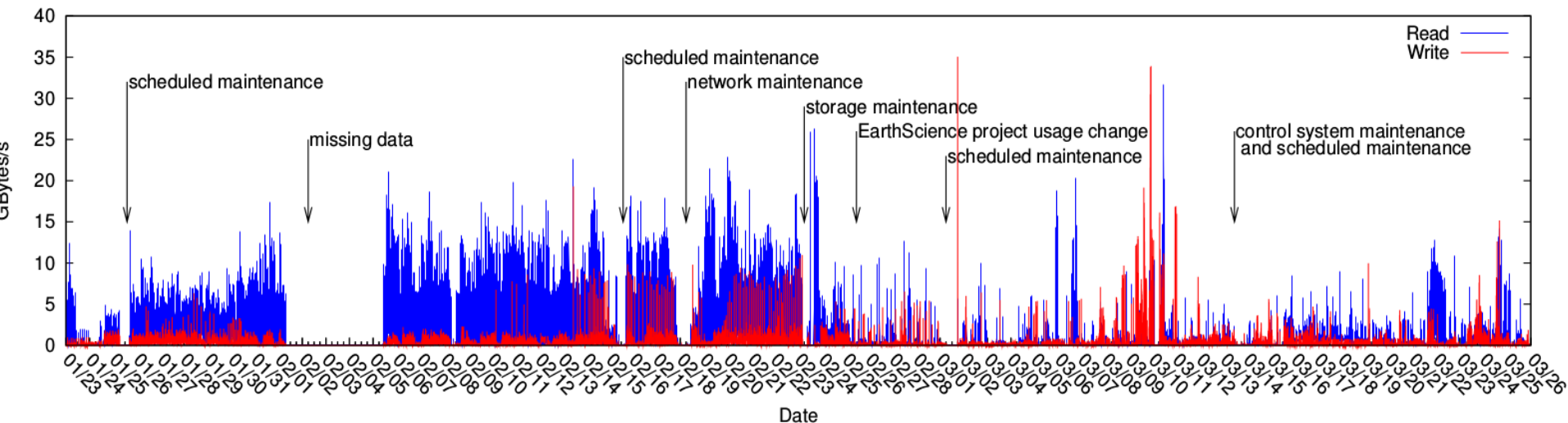**2/ Checkpointing means draining half of system memory**

Pre-Exascale system:
**4 Petabyte → bandwidth requirement 5.6 TB/s**

Oakridge,
**Teng Wang, Weikuan Yu et al. " An Efficient Distributed Burst Buffer for Linux"**

**DataDirect**™
N E T W O R K S

ddn.com

# Why IME?
## Bandwidth detailed view



**99% of the time the IO sub-system is stressed bellow 30% of its bandwidth**
**70% of the time the system is stress under 5% of its peak bandwidth**

Argone lab.
*P. Carns, K. Harms et al, Understanding and Improving Computational Science Storage Access through Continuous Characterization*

ddn.com

# Why IME?
## economics

**SSD reshuffle the parameters**
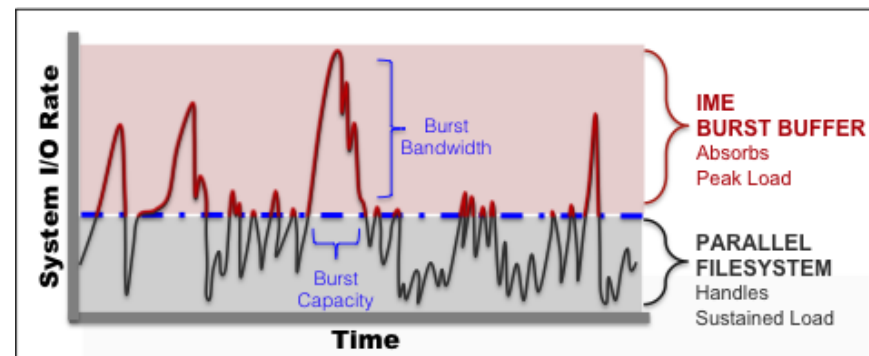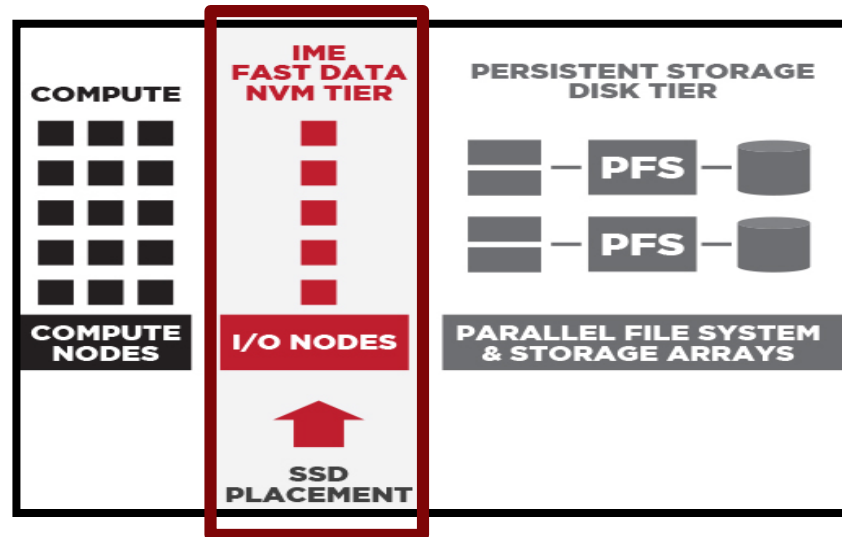**Latency / 40 :   4ms → 0,1 ms**
**Bandwidth x 3: 150 → 450 MB/s**
**Capacity / 8 :    8 → 1TB.**

**Cost x 10 $ 0,05/Gbit → $0.04**

**+ SSD deprecation rate !!**

**What can we do with a costly high bandwidth low latency technology ?**

ddn.com

# What is IME?
## Distributed virtually shared coherent array of SSDs



**Hardware tiering**

Tier 0
Tier 1
Tier 2
Tier 3
Long-Term Archive ($/GB/Watt)

**System architecture**

Compute Nodes — Compute Nodes
Compute Nodes — Compute Nodes
Compute Nodes — Compute Nodes

ION | ION | ION | ION | ION | ION | ION | ION

IME

Switch

PFS | PFS | PFS | PFS

SFA | SFA | SFA | SFA

# IME New Data Flow, New software stack



**Compute Nodes**

**2** IME Client sends buffers or fragments to IME Servers

**3** IME Servers write buffers to NVM and manage internal metadata. IME performs data coalescence and full-stripe alignment

I/O

**IME Servers**

**4** IME Servers write aligned sequential I/O to DDN SFA or other storage array backend

**Parallel Filesystem Storage**

IME Client

**Application**

**1** IME Client intercepts application I/O. Places fragments into buffers with corresponding parity buffers

DataDirect™ NETWORKS

ddn.com

# IME means Scalability and Reliability
Alleviate lock pressure: byte addressable, no page lock

**Every IO is recorded as a frag:**
 **File ID**
 **Address**
 **Length**

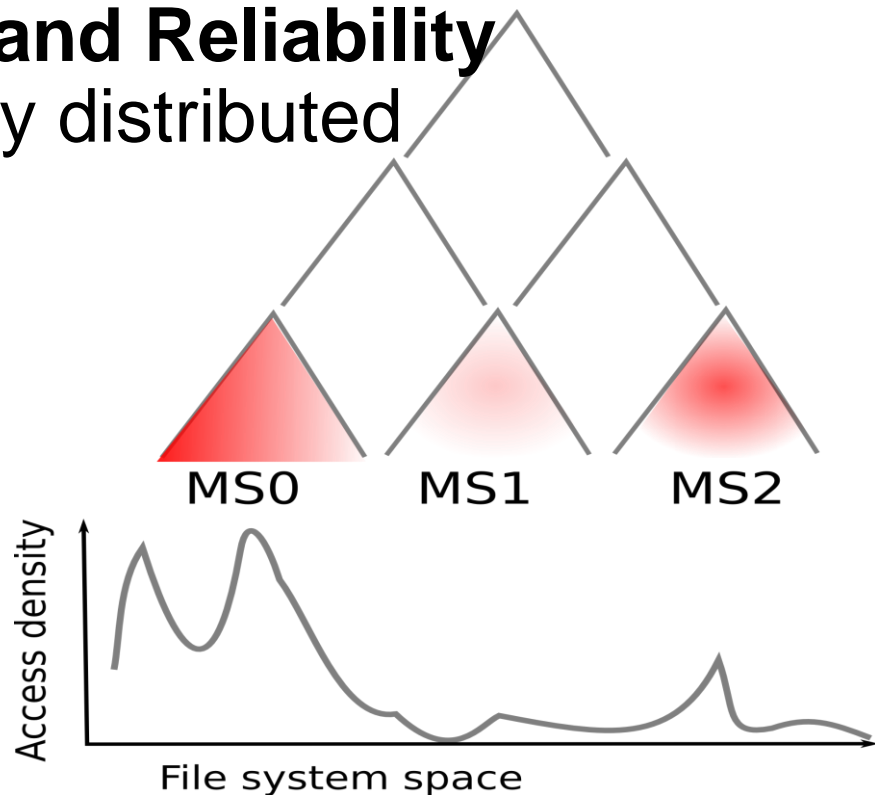**Frags are aggregated in large network buffer**

**Frags are weakly ordered, weak Posix compliance**

# IME means Scalability and Reliability
## Metadata directory is fully distributed



**No centralized directory**
**Distribution independent from structure**
                **→ no hot directory issue**
**Common knowledge ensured through hash**
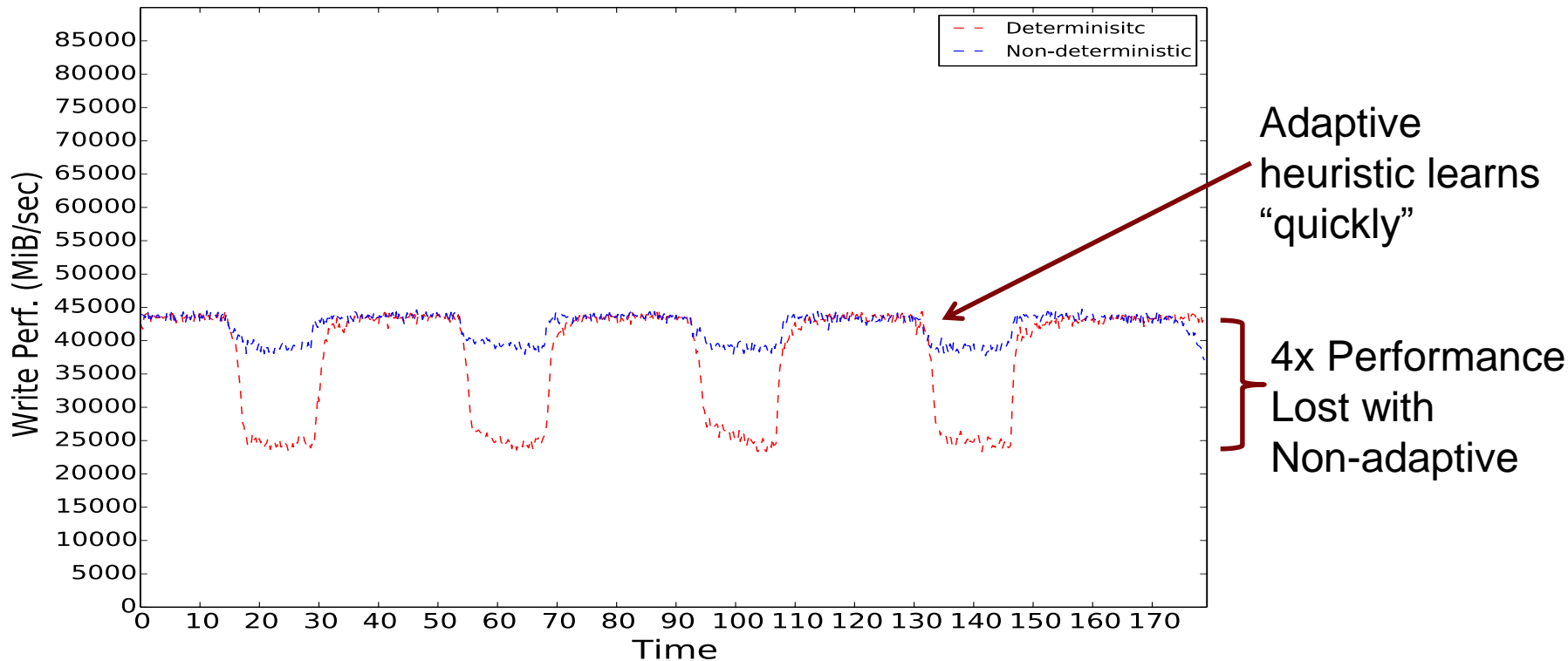
                **→ CRUSH like hash function**

Santa Cruz U., Supercomputing 2006
**Sage Weil, Scott Brandt et al.** *CRUSH: Controlled, Scalable, Decentralized Placement of Replicated Data*

ddn.com

DataDirect™
N E T W O R K S

# From fault tolerance to resilience
## Leverage server fault tolerance for load balancing



Adaptive heuristic learns "quickly"

4x Performance Lost with Non-adaptive

# IME means Scalability and Reliability
## Parity and data protection is not going to scale

**IME software only**
**Erasure coding not scalable on server side**
**Done on the IME client side**
 - **Vector instructions**
 - **Heavily multithreaded**

Pittsburgh Supercomputing Center., Supercomputing 2006
**Paul Nowoczynski et al.** Zest Checkpoint Storage System for Large Supercomputers

DataDirect™
N E T W O R K S

ddn.com

# IME Fault Recovery

- SSD Failure
  - Data recovery responsibility belongs to IME server with failed device
  - Recovered data are rewritten to remaining SSDs

- IME server failure
  - Recovery performed by remaining servers
    - Aka "Node Ejection"
  - Fully declustered
    - Remaining servers share rebuild processing

# IME Boost

**Accelerates applications**
**→ especially small or mal-aligned I/O**

**No page lock**
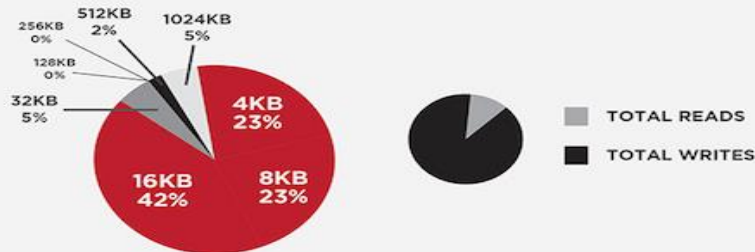**Client coalescing:**
  **→ Bulk RPC transfer**
  **→ Save SSD write cycle**
**Server flush scheme**
  **→ tune for PFS parameters**
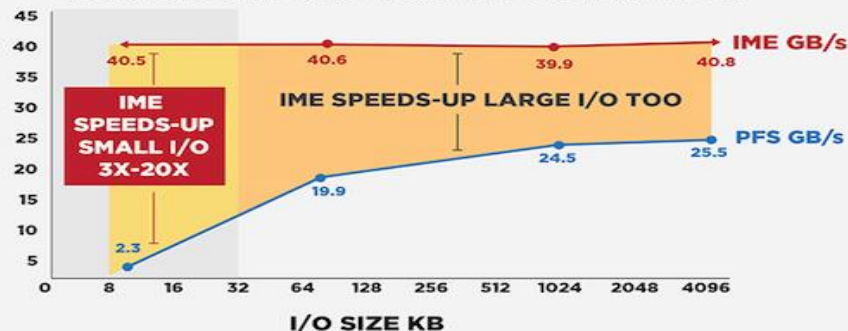


**WRITE DISTRIBUTION FOR MULTI-DISCIPLINARY HPC CLUSTER**
EVEN LARGE HPC SITES DRIVE A LOT OF SMALL I/O

256KB 0%
512KB 2%
1024KB 5%
128KB 0%
32KB 5%
4KB 23%
16KB 42%
8KB 23%

TOTAL READS
TOTAL WRITES

**90% OF ALL I/O IN TYPICAL HPC DATACENTERS IS <32KB IN SIZE**

**HOW IME HELPS**
ENABLES HIGHER PEAK BANDWIDTH THAN DISK-BASED PFS, ESPECIALLY FOR SMALL I/OS

IME GB/s
40.5 40.6 39.9 40.8
IME SPEEDS-UP SMALL I/O 3X-20X
IME SPEEDS-UP LARGE I/O TOO
PFS GB/s
19.9 24.5 25.5
2.3

I/O SIZE KB

**DataDirect**™
**N E T W O R K S**

ddn.com

# IME Boost
## IOR SSF 4k interleaved writes (SC'13)



*3 IOR Checkpoints to IME*

**SC'13 IME peak bandwidth (~50GB/s)**

*Migration of 3rd checkpoint to PFS*

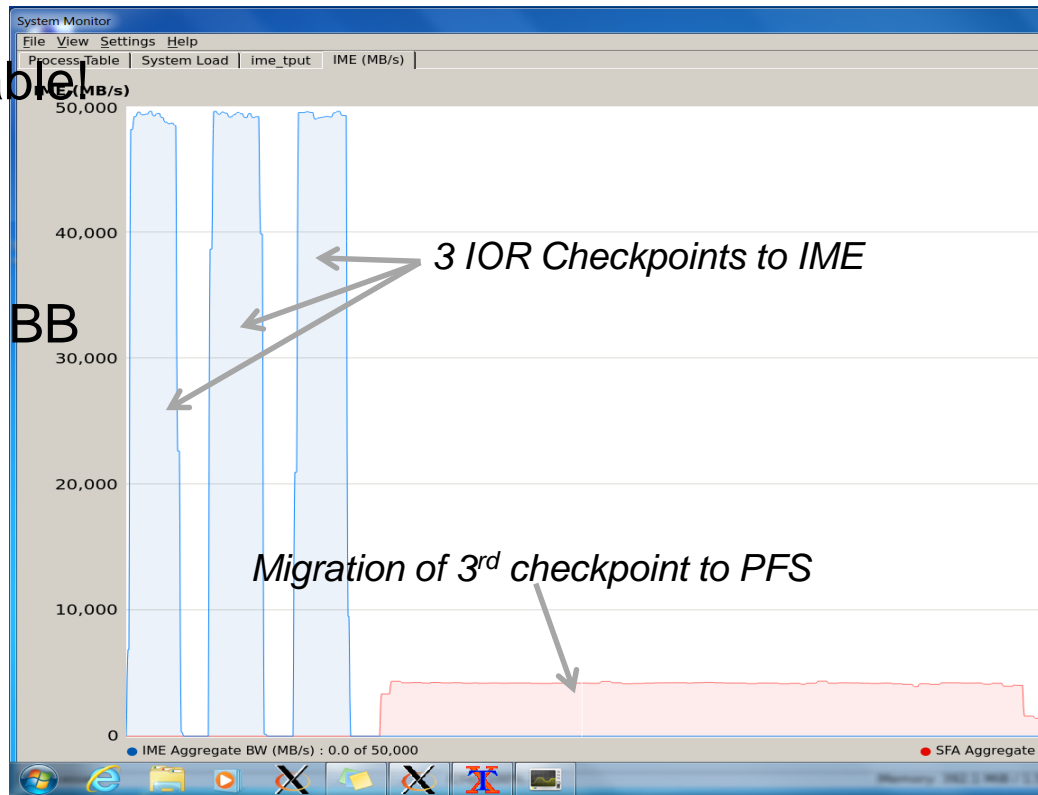**SC'13 PFS peak bandwidth (~4GB/s)**

ddn.com

# IME Boost
## Reduce data movement

Copying all data to PFS is not desirable!

- Reduce data movement

- 3$^{rd}$ checkpoint should be freed from BB following copy to PFS

- ***Explicit management is required***



*3 IOR Checkpoints to IME*

*Migration of 3$^{rd}$ checkpoint to PFS*

ddn.com

# Future: Node Local NVM

Node local provides a new set of challenges.

•Affiliation of file components to specific compute nodes must be expressed

•Will a job actually run on the compute nodes where data has been staged-in?
  –At large issue in HPC deployments

•Blur the distinction between client and server.
  –Performance instructiveness in compute node

ddn.com

# IME take away

- Extra tier of SSD in storage hierarchy
    - → Re-design the software stack to address real issues
    - → Resilience and scalability
- Keep data as close as possible from compute node
- Pluggable with other IO storage technology
    - → IO libraries
    - → Object storage
- Pave the way for future evolutions
    - → A component of the Exascale storage stack

**Thanks !**

# IME Approach

## Distributed Hash Table + Log Structuring