

# Accelerating High Performance Cluster Computing Through the Reduction of File System Latency

**David Fellingner**  
Chief Scientist, DDN Storage

## In large clusters, primarily on the cluster itself

- ✓ Lower latency interconnects
- ✓ More efficient message passing structures
- ✓ Higher performance processors & GPUs

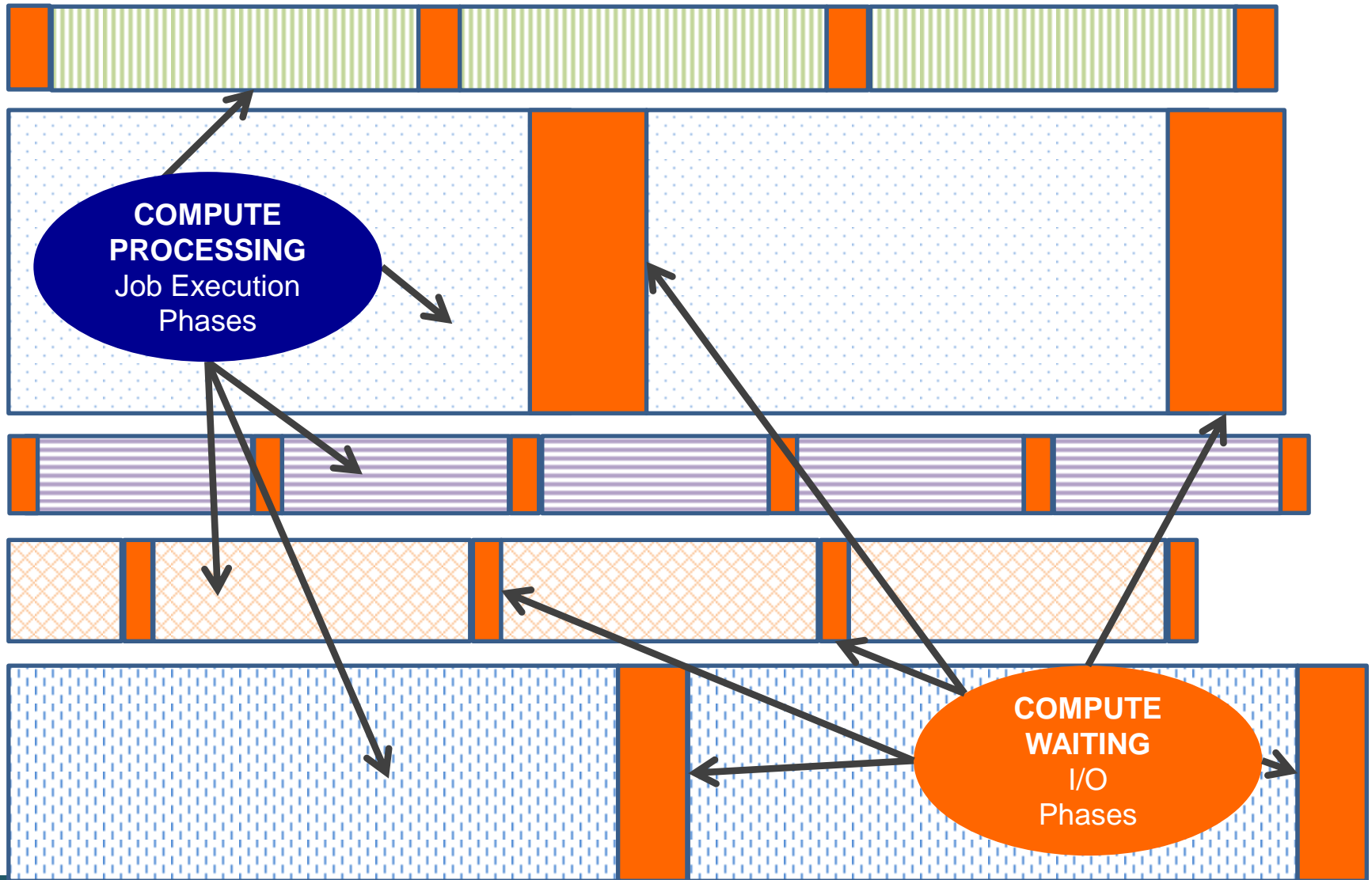
## Also, research & study on processing techniques to achieve true parallel processing operations

- Symmetric multi-processing vs. Efficient message passing



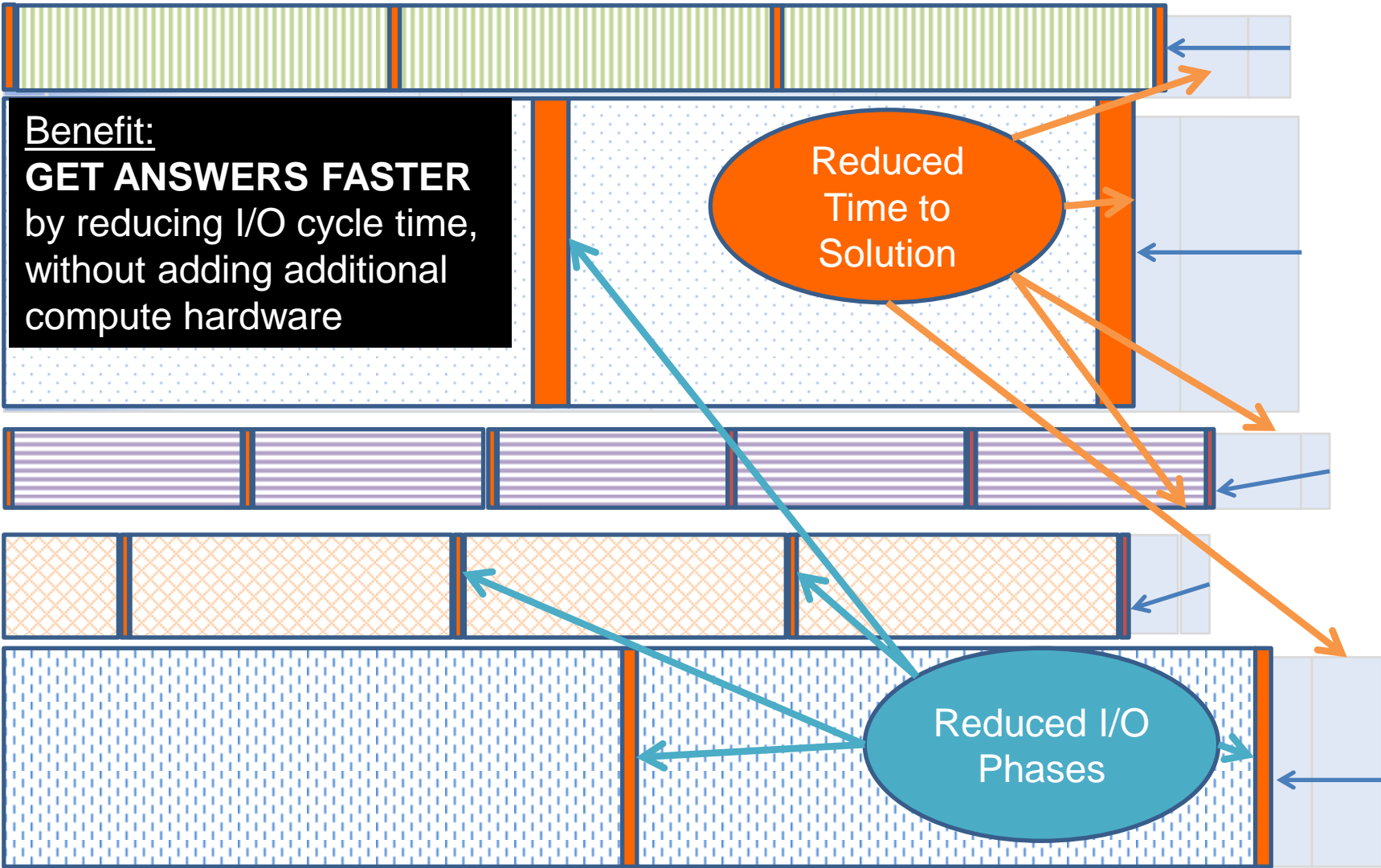
# Today's Challenge: Bulk I/O Latency

HPC Jobs

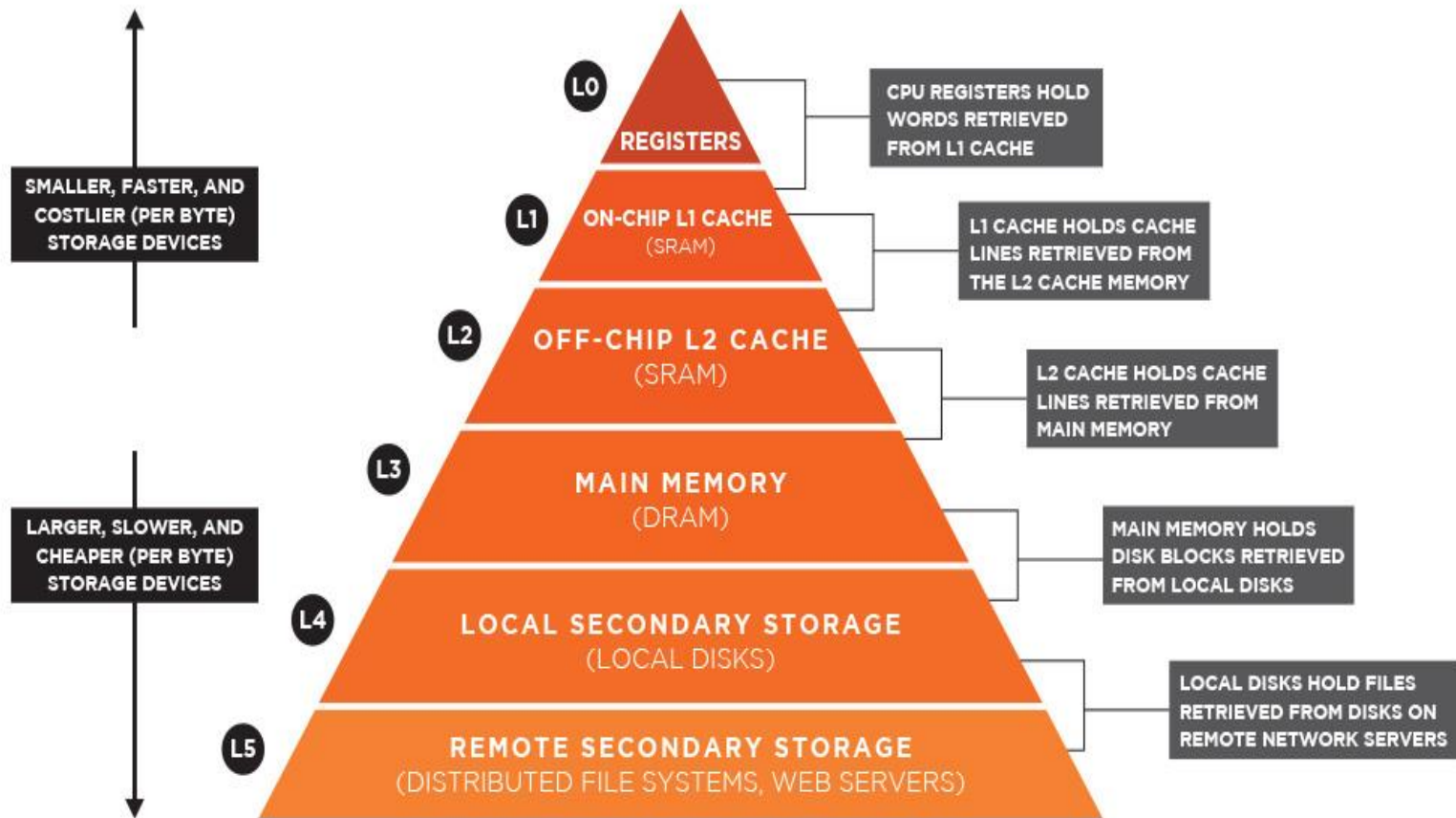


# What's Needed? Compressed I/O!

HPC Jobs

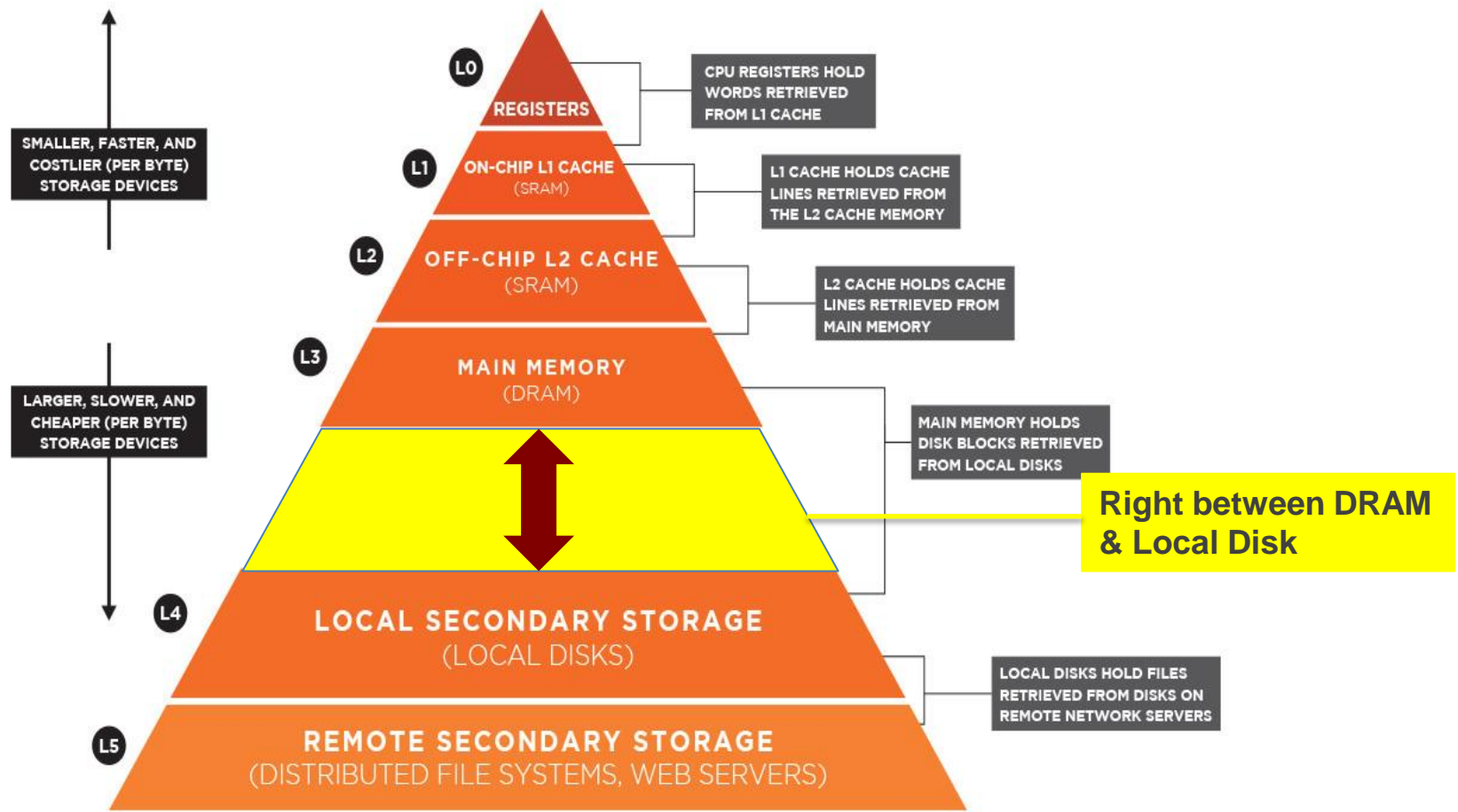


## AN EXAMPLE OF MEMORY HEIRARCHY

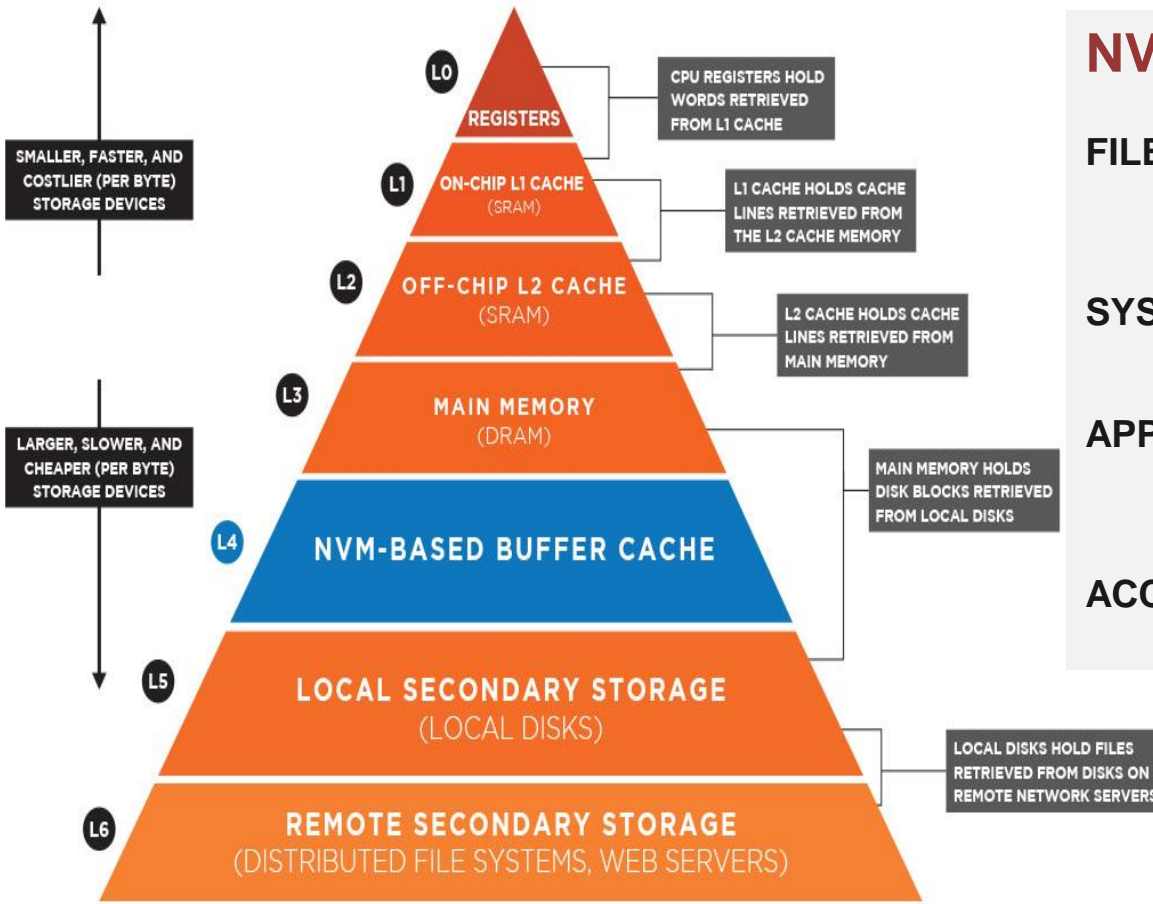


# The Emerging Storage Tier's Placement

## AN EXAMPLE OF MEMORY HEIRARCHY



## AN EXAMPLE OF MEMORY HEIRARCHY



## NVM-BASED BUFFER / CACHE

**FILE CACHE** between DRAM & local disk

- Performance
- Capacity
- Cost

**SYSTEM MANAGED** Resource

- Shared PFS acceleration, pinned libraries, common datasets, etc.

**APPLICATION MANAGED** Resource

- Allocated on a per-job basis, dedicated to a specific job or application, etc.
- Application “co-design”

**ACCELERATING** parallel file system I/O

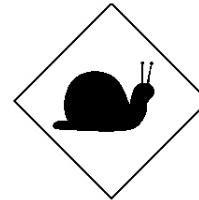
- POSIX, MPI, etc.

# Is NVM Memory Hardware Enough?

**Solid state storage offers high IOPS and low power,  
but in HPC . . .**



**File Systems  
Force Locks**



**Many Operations Are  
Forced To Piecewise  
Sequential Constructs**



**Internal FS Operations  
Are Single Threaded**



**Load Leveling  
Required To Extend  
SSD Life**



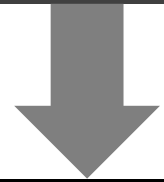
**Large Operations Limited  
By Port Bandwidths**

**SOFTWARE MUST BE DEVELOPED TO UTILIZE  
THE BENEFITS OF SSD HARDWARE**



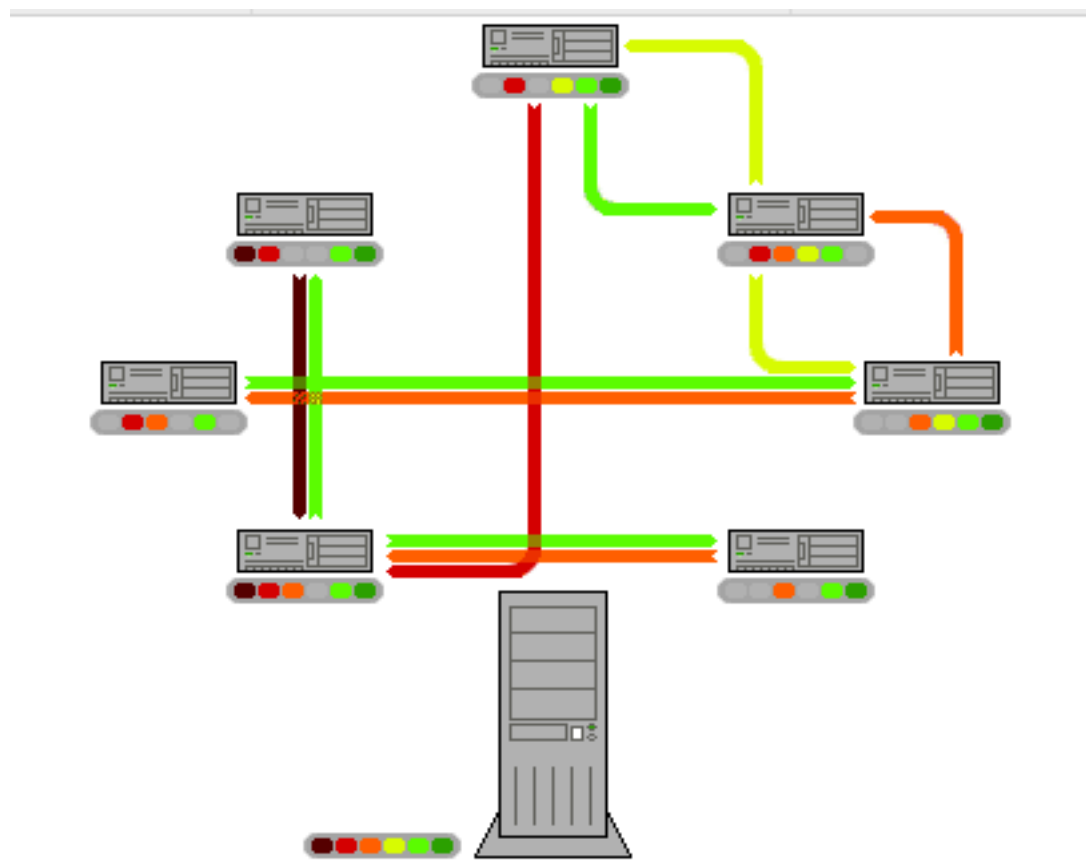
# Other Industry Approaches

**BEFORE:**  
**Single FTP Server**  
(regardless of performance & availability)



**NOW:**

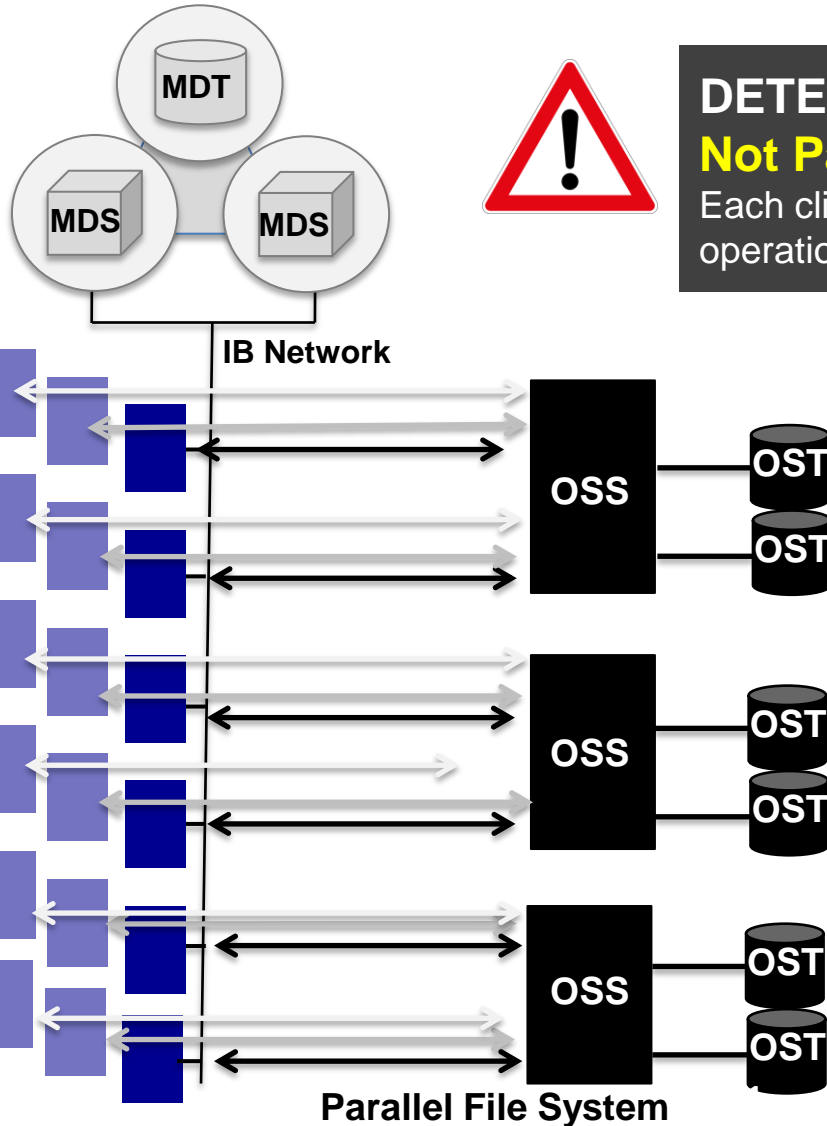
- Utilizing **parallelism** for network transfer efficiency
- **Data sharing**
- **P2P Grid**



# Why is HPC living in the “Internet 80’s”?



# Is a Parallel File System Really Parallel?

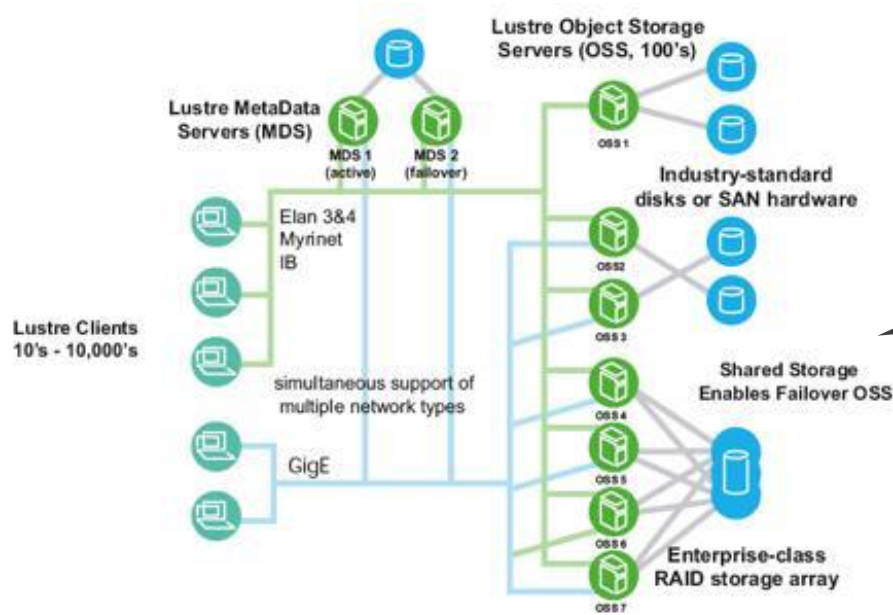


## DETERMINISTIC WRITE SCHEMA Not Parallel – It’s a Bottleneck!

Each client competes to talk to a single OSS (serial operation)

What’s  
Needed?

A WRITE ANYWHERE  
**FILE SYSTEM** is the next  
evolution



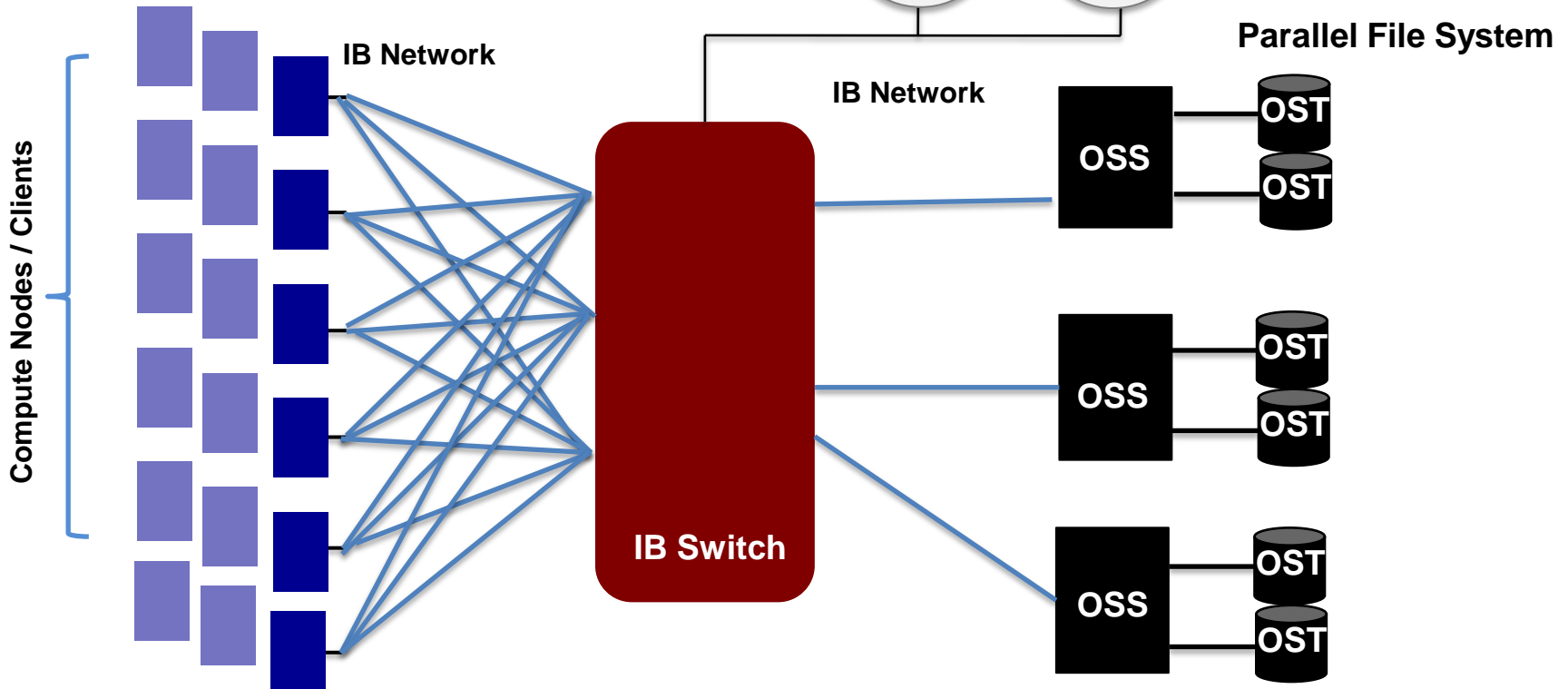
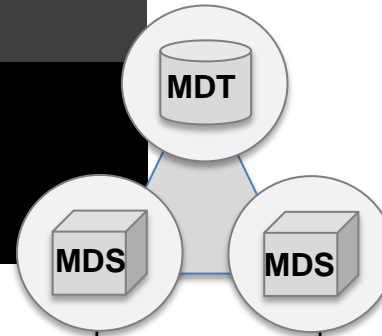
But . . .  
**Where's the  
FATs?**

- ① Ext4 extended write (I make a call to a Lustre client)
- ② Redirect request to metadata server
- ③ Metadata server returns OST number & iNode number, execute locks
- ④ Begin classic write operation: iNode to file loc table EXT 4 gathers blocks from garbage collection to extent list
- ⑤ Metadata server assignment of handle - Then lock is released

# IME Makes a Parallel File System Parallel

BEFORE: MDS DICTATED SERIAL ACCESS

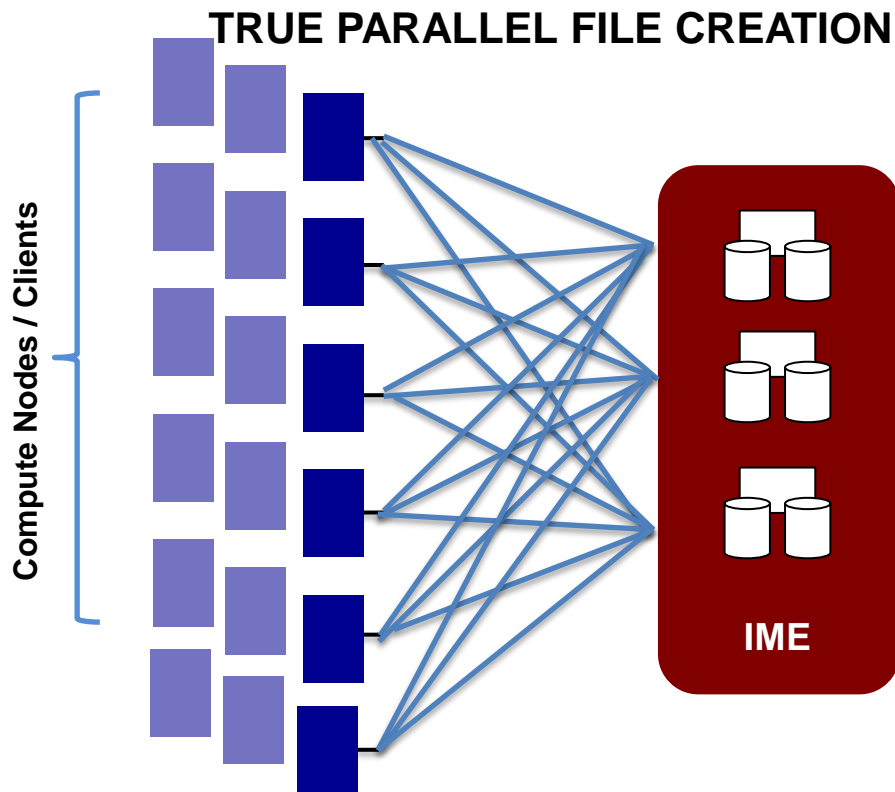
NOW: **PARALLEL DATA ACCESS**  
POSIX semantics & PFS bottleneck broken!



## The Magic of Write Anywhere

Now . . . **THE PFS IS PARALLEL!**

Every compute node can write file increments to every storage node along with metadata & erasure coding



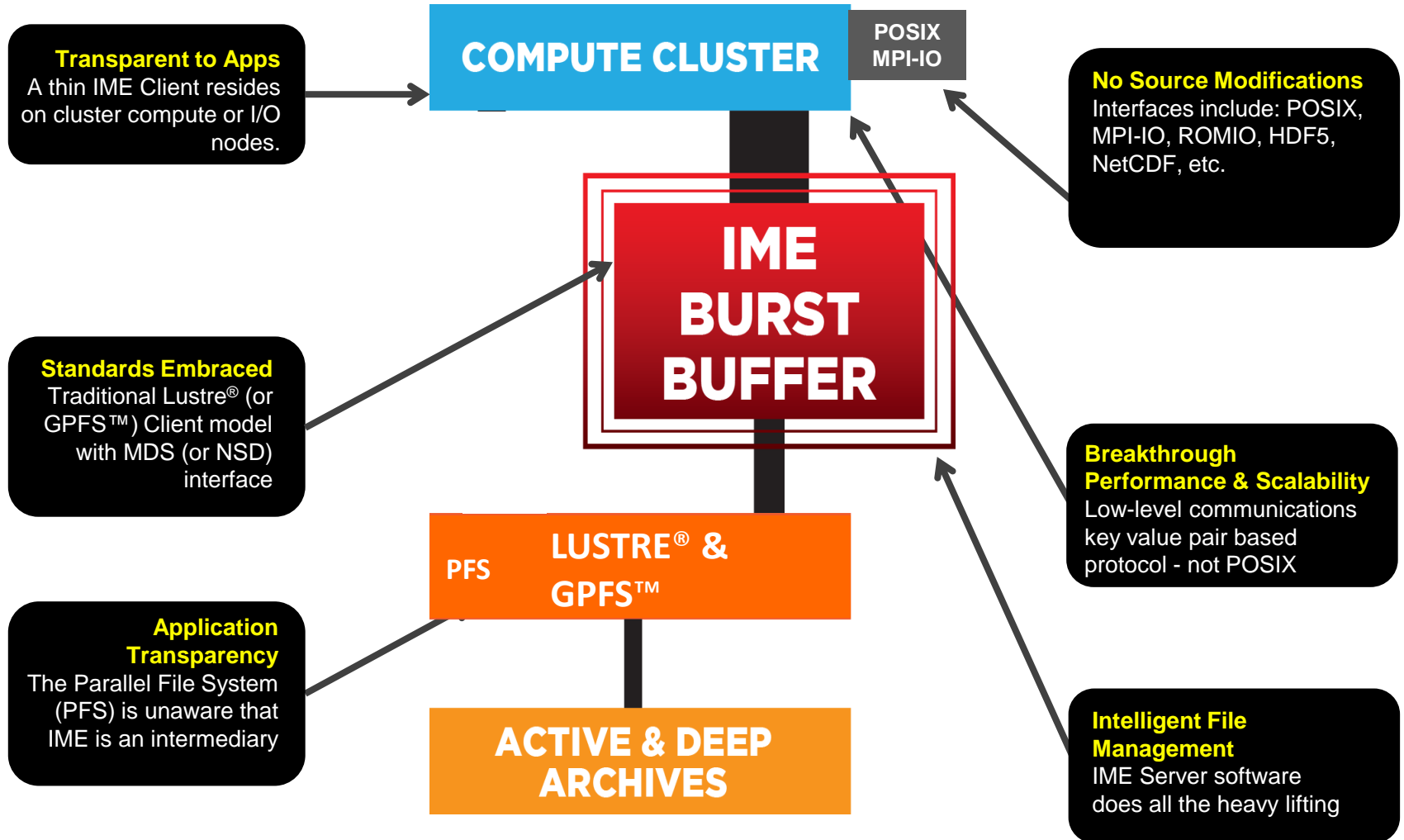
In IME, we've implemented a DHT.

Now, the compute is sharing files the same way Napster users do



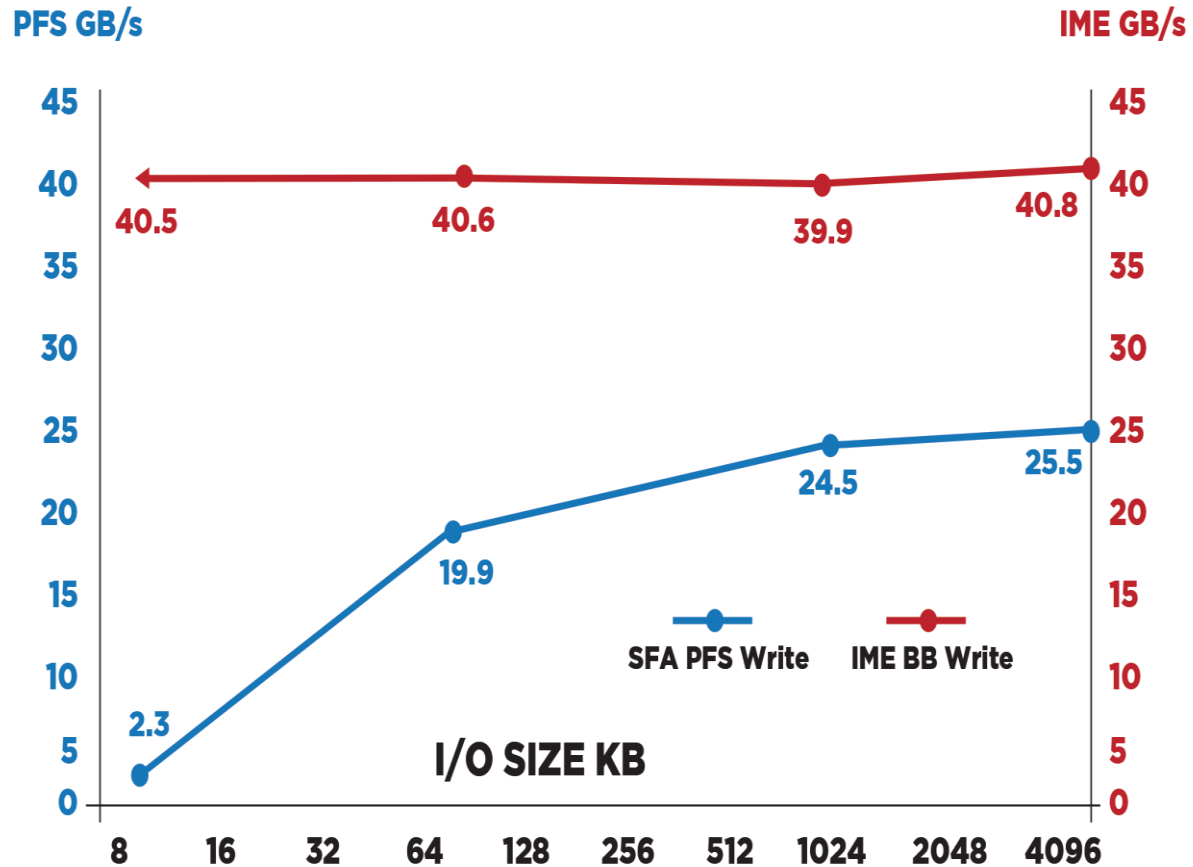
# Introducing IME<sup>®</sup>

## Technical Key Components and Operations



### DDN 32 x FDR GE CLIENTS - IOR MPIIO SINGLE SHARED FILE - NO SEGMENTS

**WRITES - PFS/BB COMPARISON -> DDN IME BB ADVANTAGE**



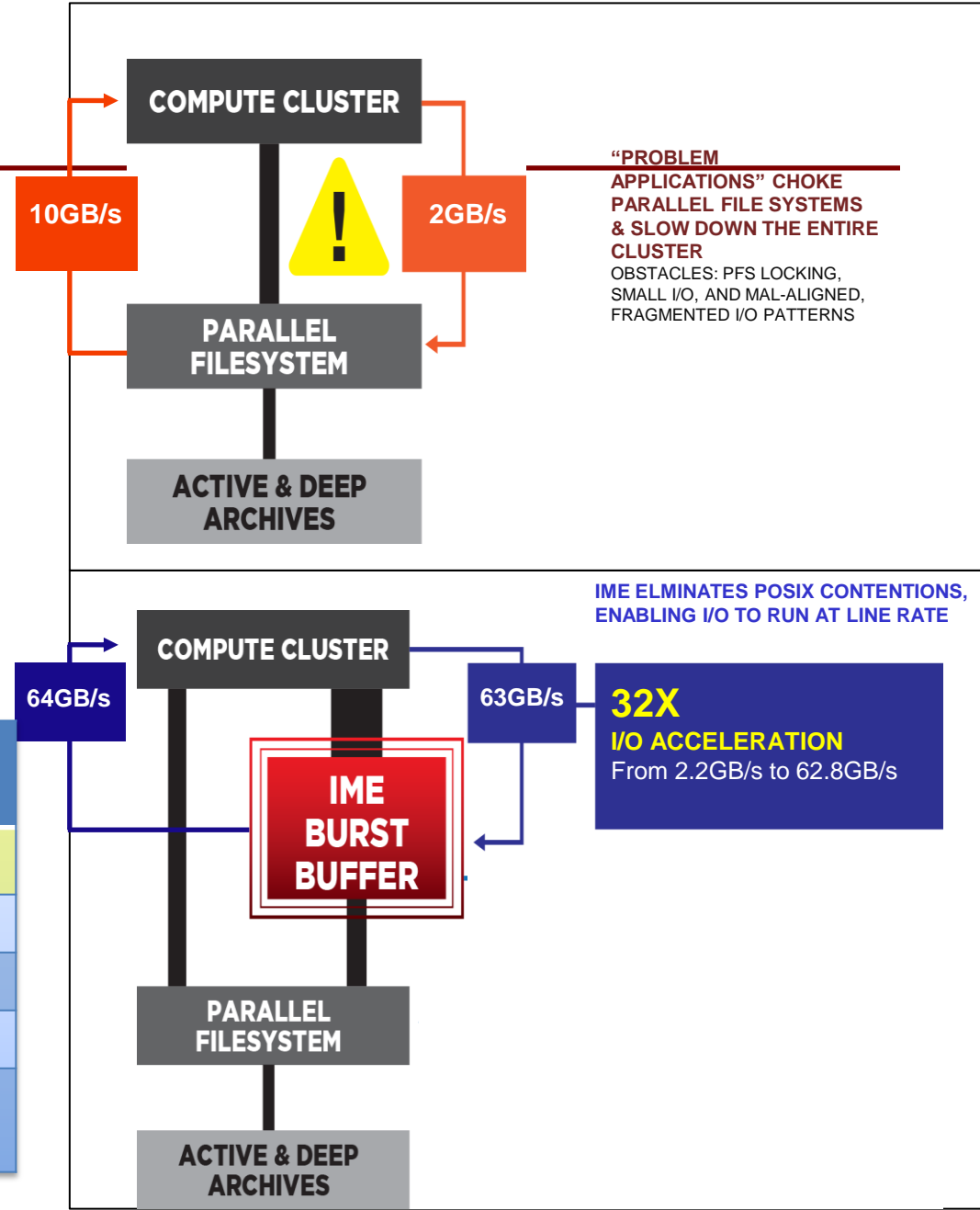


# IME<sup>®</sup> Benchmarking

## HACC\_IO @TACC Results

Description:  
HACC-IO is an HPC Cosmology Kernel

Particles per Process	Qty. Clients	IME Writes (GB/s)	IME Reads (GB/s)	PFS Writes (GB/s)	PFS Reads (GB/s)
34M	128	62.8	63.7	2.2	9.8
34M	256	68.9	71.2	4.6	6.5
34M	512	73.2	71.4	9.1	7.5
34M	1024	63.2	70.8	17.3	8.2
<b>IME Acceleration</b>		<b>3.7x-28x</b>	<b>6.5x-11x</b>		



HPC Cluster = Large Group of Data Users

✓ Why haven't we learned . . .

What the internet p2p guys have known for a long time?

✓ Learn to share!

The precious resources of network bandwidth & storage



# Thank You!

Keep in touch with us

## David Fellingner

Chief Scientist,  
DDN Storage

[dfellinger@ddn.com](mailto:dfellinger@ddn.com)



[sales@ddn.com](mailto: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)