# DAOS: Storage Innovations Driven by Intel® Optane™

Liang.zhen@intel.com

Principal Engineer, Intel

intel®

# Notices and Disclaimers

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration.

No product or component can be absolutely secure.

Tests document performance of components on a particular test, in specific systems. Differences in hardware, software, or configuration will affect actual performance. For more complete information about performance and benchmark results, visit http://www.intel.com/benchmarks .

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products.   For more complete information visit http://www.intel.com/benchmarks .

Intel Advanced Vector Extensions (Intel AVX) provides higher throughput to certain processor operations. Due to varying processor power characteristics, utilizing AVX instructions may cause a) some parts to operate at less than the rated frequency and b) some parts with Intel® Turbo Boost Technology 2.0 to not achieve any or maximum turbo frequencies. Performance varies depending on hardware, software, and system configuration and you can learn more at http://www.intel.com/go/turbo.

Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice.

Cost reduction scenarios described are intended as examples of how a given Intel-based product, in the specified circumstances and configurations, may affect future costs and provide cost savings.  Circumstances will vary.  Intel does not guarantee any costs or cost reduction.

Intel does not control or audit third-party benchmark data or the web sites referenced in this document. You should visit the referenced web site and confirm whether referenced data are accurate.
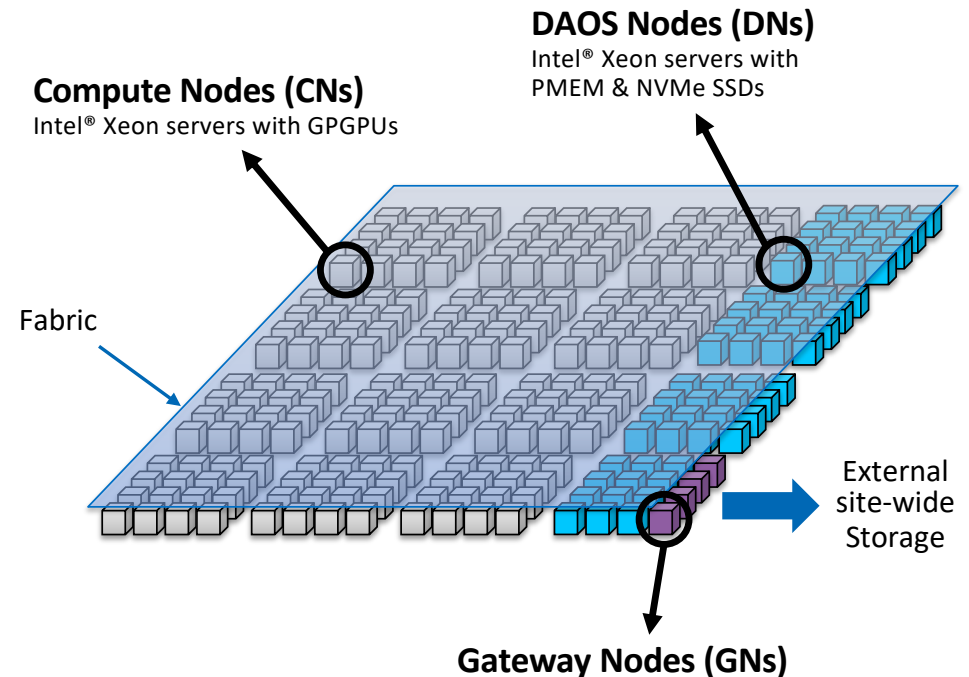
# DAOS overview

intel.

# What is DAOS?
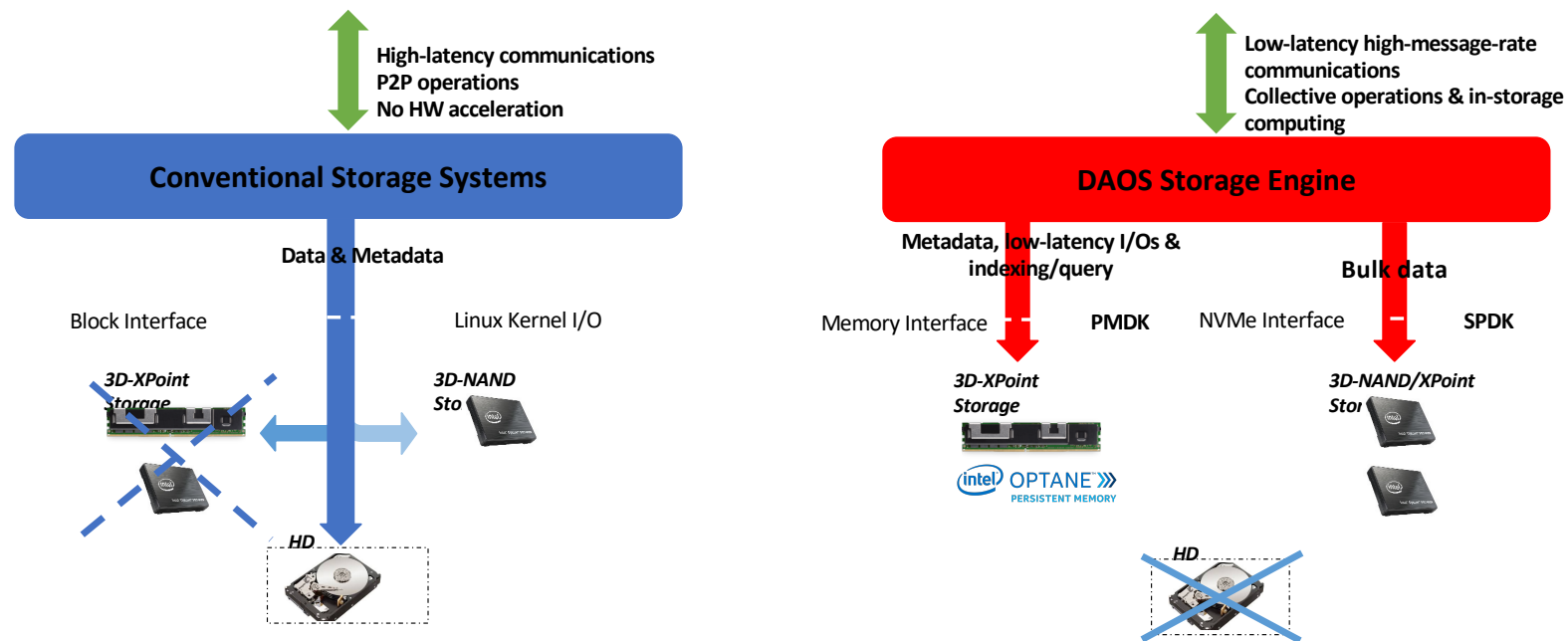
- Storage pools globally accessible over the fabric

- Overcomes industry bottlenecks by leveraging Intel® Optane™ Persistent Memory and NVMe SSDs

- Delivers exceptionally high bandwidth and IOPS, meeting the demands of HPC and AI

- Strong distributed consistency (Database like)

- Tightly integrated with Applications

- Can be deployed as either a standalone file system, or a performance tier integrated with existing storage systems

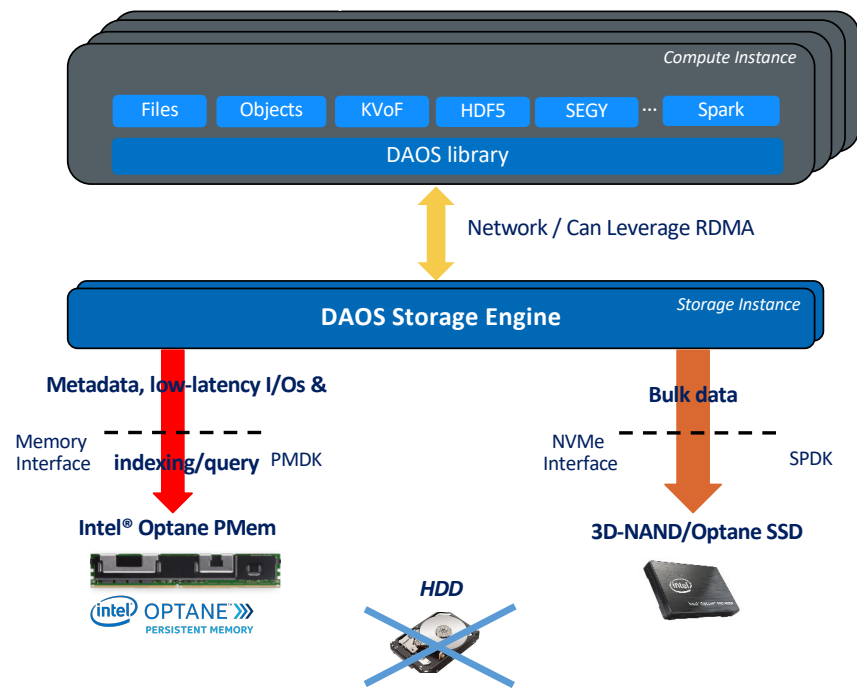**DAOS Nodes (DNs)**
Intel® Xeon servers with PMEM & NVMe SSDs

**Compute Nodes (CNs)**
Intel® Xeon servers with GPGPUs

Fabric

External site-wide Storage

**Gateway Nodes (GNs)**

**Generational leap forward in storage performance**

# DAOS architecture

**High-latency communications
P2P operations
No HW acceleration**

## Conventional Storage Systems

**Data & Metadata**

Block Interface                    Linux Kernel I/O

*3D-XPoint
Storage*

*3D-NAND
Sto*

*HD*

**Low-latency high-message-rate
communications
Collective operations & in-storage
computing**

## DAOS Storage Engine

**Metadata, low-latency I/Os &
indexing/query**

**Bulk data**

Memory Interface          **PMDK**      NVMe Interface          **SPDK**

*3D-XPoint
Storage*

(intel) OPTANE⟫⟫
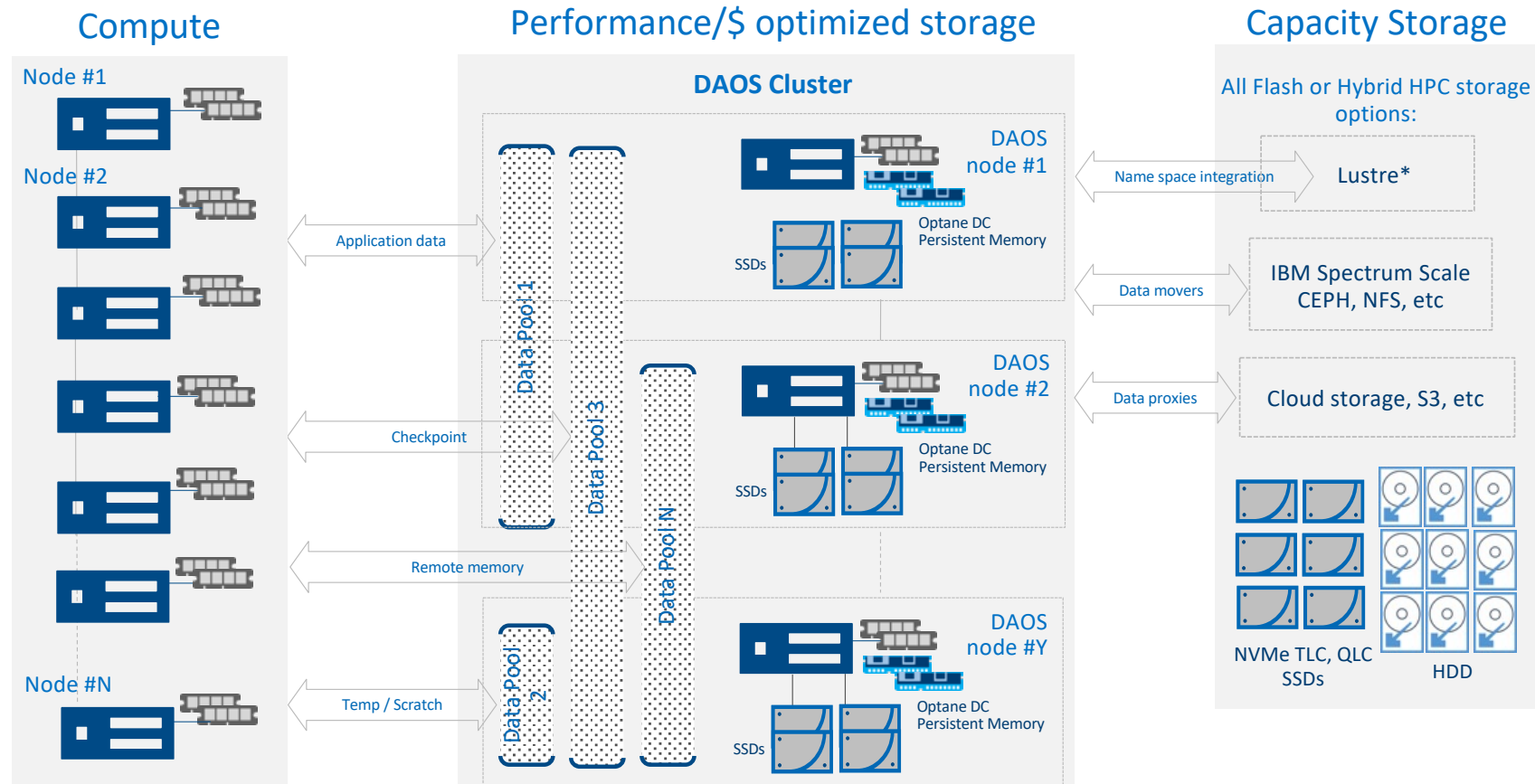PERSISTENT MEMORY

*3D-NAND/XPoint
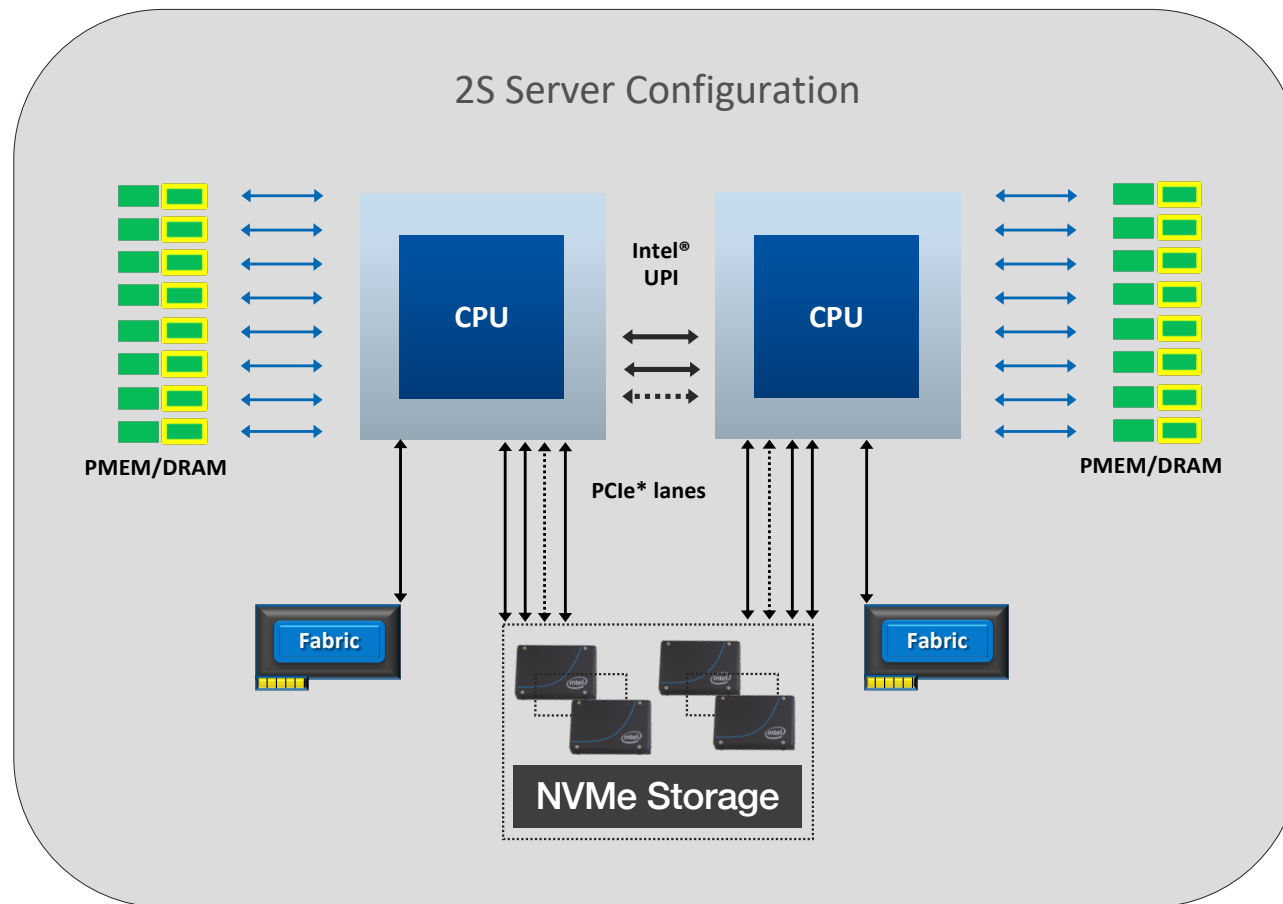Stor*

*HD*

# DAOS stack overview



- High **throughput/IOPS @arbitrary** alignment/size (2M IOPS in 1U)

- **Low-latency fine-grained** I/O

- Data access time **orders of magnitude faster** (μs vs ms)

- Highly **scalable**

- Operates in **userspace**

- **Rich** storage semantics

- Support **smooth migration** with support in common frameworks such as Apache Spark*, MPI-IO, HDF, POSIX, etc
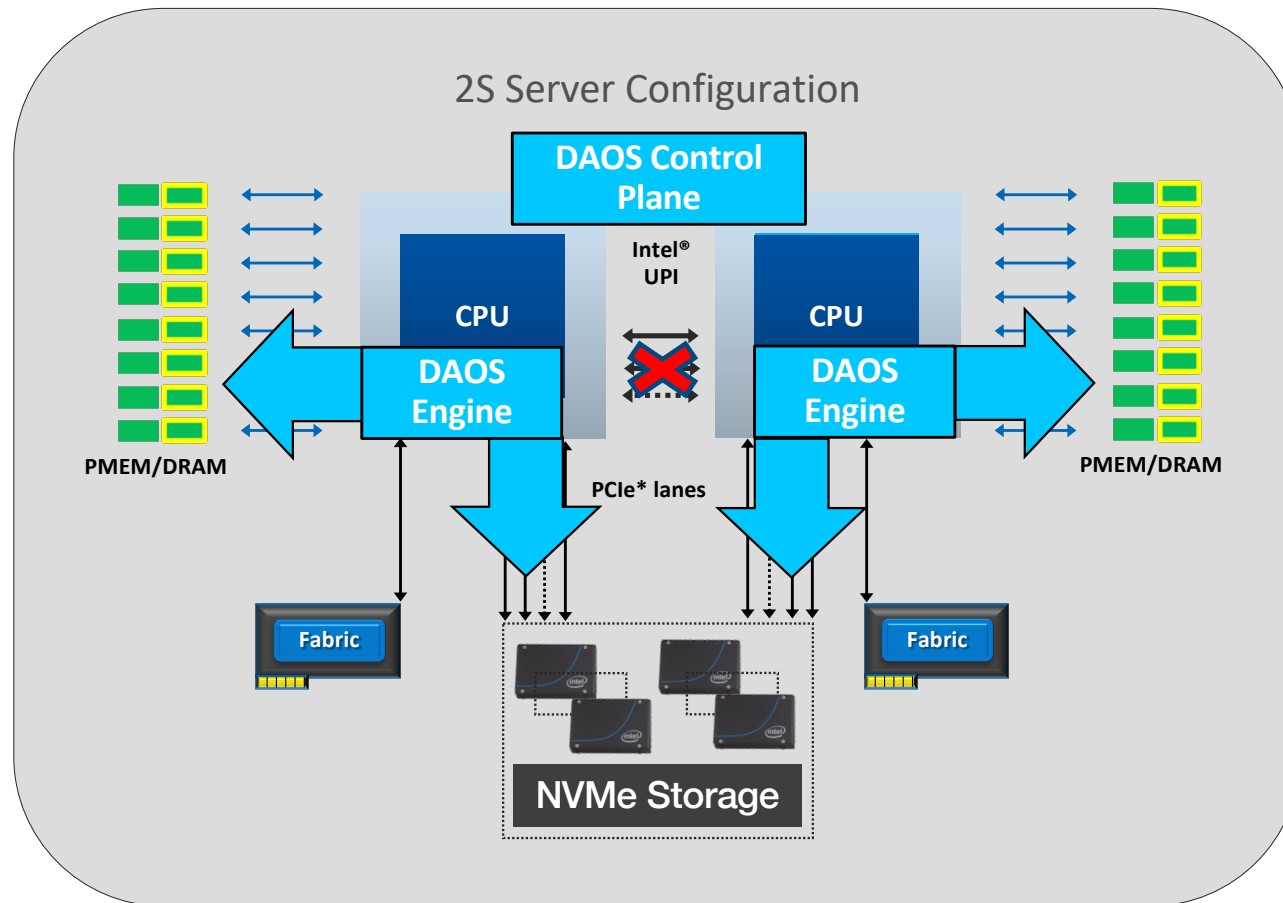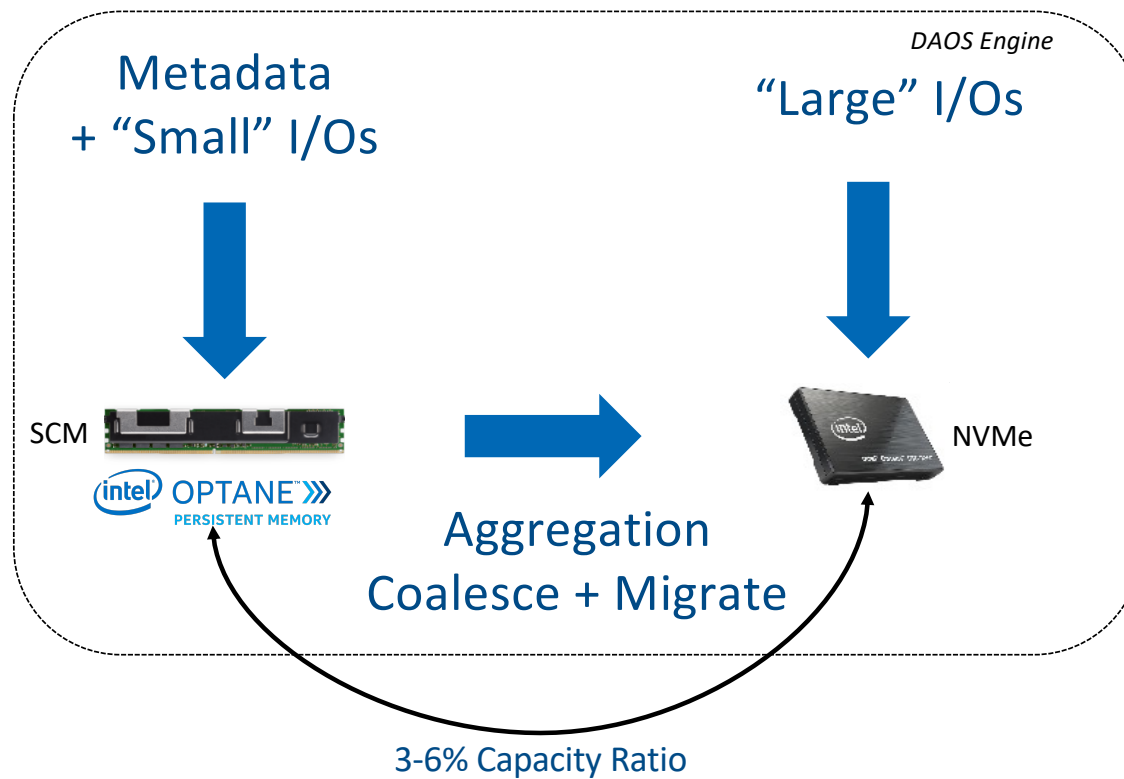
# DAOS in the Overall Cluster Architecture

**Compute**

**Performance/$ optimized storage**
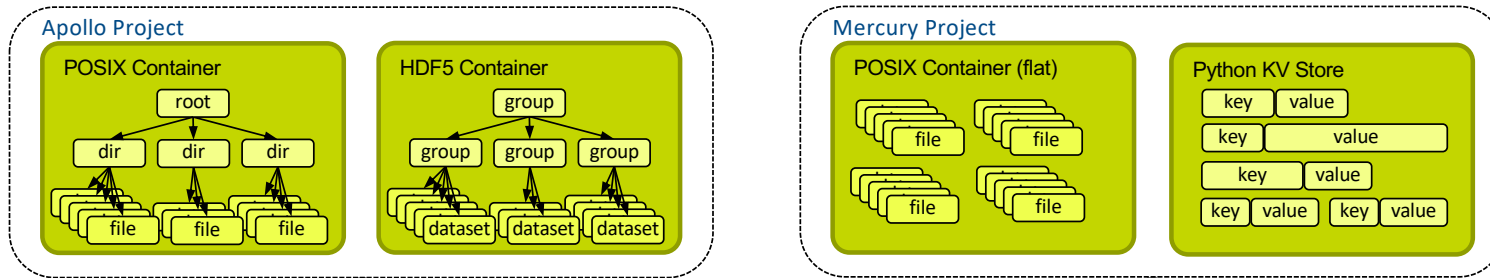
**Capacity Storage**

Node #1

Node #2

Node #N

**DAOS Cluster**

Data Pool 1

Data Pool 3

Data Pool N

Data Pool 2

Application data

Checkpoint

Remote memory

Temp / Scratch

DAOS node #1

Optane DC Persistent Memory

SSDs

DAOS node #2

Optane DC Persistent Memory

SSDs

DAOS node #Y

Optane DC Persistent Memory

SSDs

Name space integration

Data movers

Data proxies

All Flash or Hybrid HPC storage options:

Lustre*

IBM Spectrum Scale CEPH, NFS, etc

Cloud storage, S3, etc

NVMe TLC, QLC SSDs

HDD

# DAOS Node Design



2S Server Configuration

PMEM/DRAM

Intel® UPI

CPU    CPU

PCIe* lanes

PMEM/DRAM

Fabric    Fabric

NVMe Storage

# DAOS Node Design



2S Server Configuration

DAOS Control Plane

Intel® UPI

CPU · CPU

DAOS Engine · DAOS Engine

PMEM/DRAM · PMEM/DRAM

PCIe* lanes

Fabric · Fabric

NVMe Storage

# Engine: Media Management



DAOS Engine

Metadata + "Small" I/Os

"Large" I/Os

SCM

NVMe

Aggregation
Coalesce + Migrate

3-6% Capacity Ratio

# DAOS Data Model: Storage Pooling



| | | | | | |
|---|---|---|---|---|---|
| Pool 1 | 🟦 | Project Apollo | 100PB usable | 20TB/s | 200M IOPS |
| Pool 2 | 🟧 | Project Gemini | 10PB usable | 2TB/s | 20M IOPS |
| Pool 3 | 🟪 | Project Mercury | 30TB usable | 80GB/s | 2M IOPS |

# DAOS Data Model: Objects

**I/O Middleware View**

Container (eg POSIX)
- root
  - dir
  - dir
  - dir
    - file
    - file
    - file

Mapping →

**DAOS Layout View**

Container
- obj1
- obj2  obj3  obj4
- obj5  obj10  obj20

object →

128-bit object identifier

Array

Multidimensional Array

key @ → val
key @ → val
key @ → val
val

Key-value Store

key1 key2 key3 @ → val
key1 key2 key3 @ → val []
key1 key2 key3 @ → val

Multi-level Key-value Store

# DAOS Data Model: Distribution & Fault Tolerance

App Buffers

App Buffers

App Buffers

**Redundancy Group**

| Engine rank 0 | Engine rank 5 | ... | Engine rank 8 |

Replication

**Redundancy Group**

| Engine rank 0 | Engine rank 3 | Engine rank 9 | ... | Engine rank N |

Erasure-code

| RG1 | RG2 | RG3 |

| Engine rank 0 | Engine rank 1 | ... | Engine rank N |

Sharding

# Algorithmic object placement

- No explicit object layout
  - Expensive to maintain
  - Extra round-trip to get the layout
- Object class
  - Data protection method
  - Data distribution requirements
  - Pre-defined attributes, identified by 16-bit integer (class ID)

- Object ID
  - Object class ID + 96-bit ID
- Algorithmic object placement
  - determines where those shards will be stored on the physical system (node, target) based on pool configuration and OID

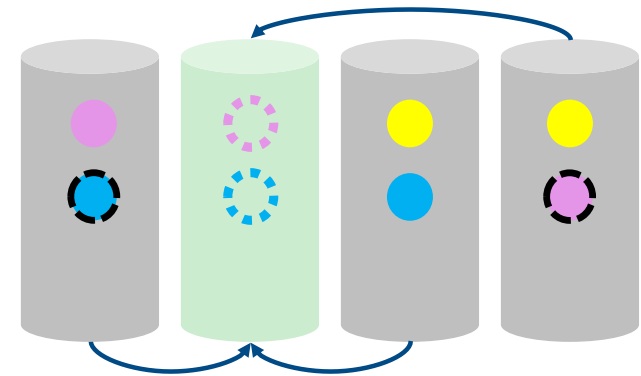Object ID → PLACEMENT

# Automated Exclusion & Self-Healing

- Health monitoring

  - Detects failed nodes via SWIM

  - Failed nodes are automatically evicted

- Failure recovery ("rebuild")

  - All other surviving nodes are notified of the failure

  - Impacted objects are automatically determined and reconstructed on surviving nodes

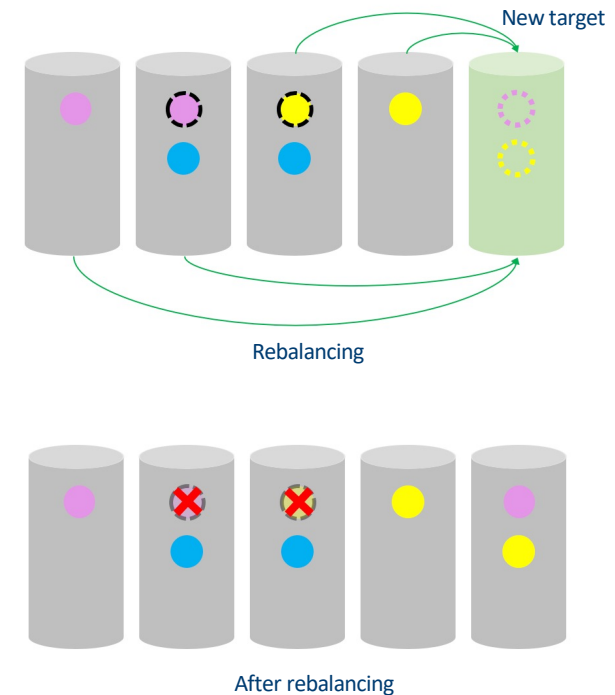# Reintegration ("reintegrate")

- Restore previously evicted servers/devices to active service

- Data moves back from fallback locations to reintegrated capacity

- Fallback devices service read/write requests until reintegration is complete

- Writes during reintegration also go to reintegrating device to ensure consistency

# Extension ("extend")
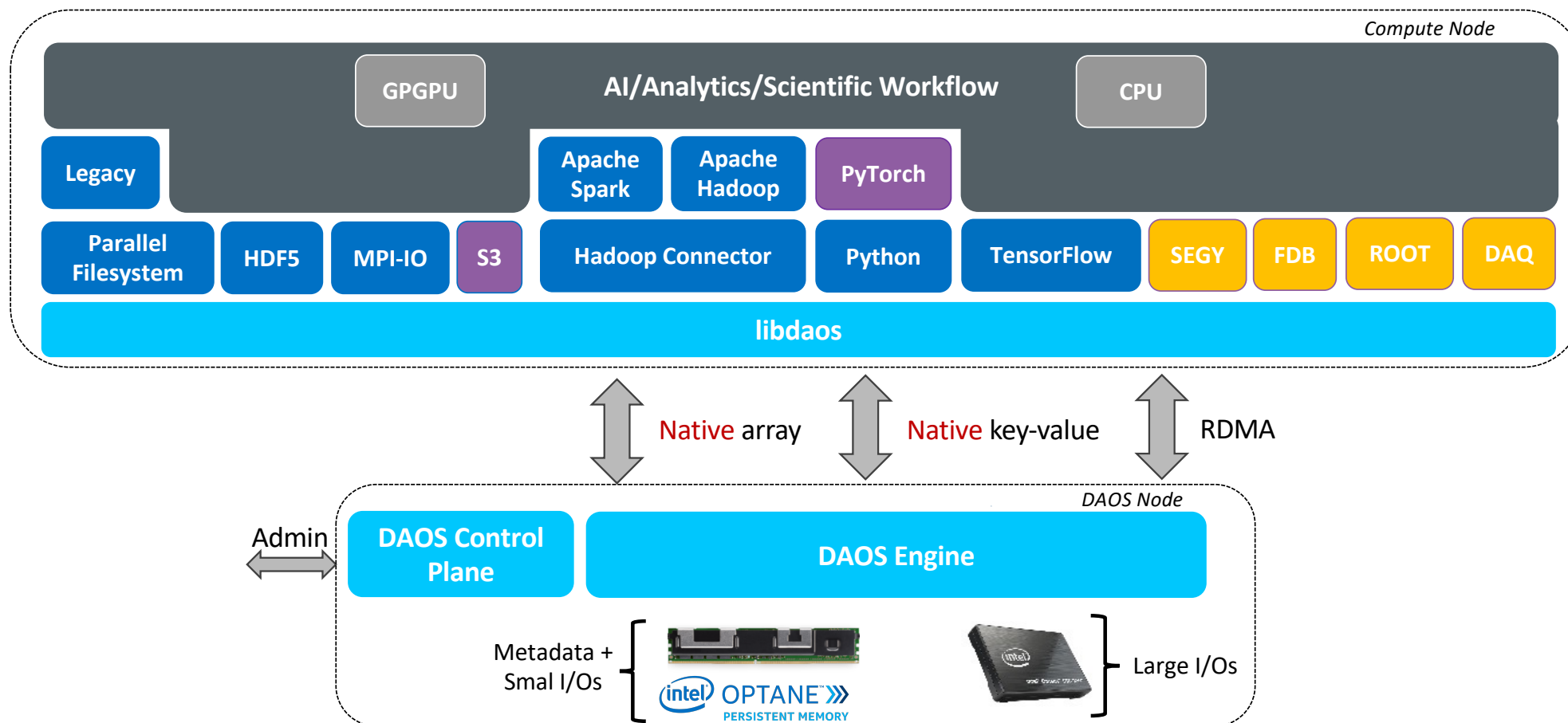


New target

Rebalancing

After rebalancing

- Add new servers/devices to the storage system

- Data rebalances automatically across new capacity

- Original devices service read/write requests until extension is complete

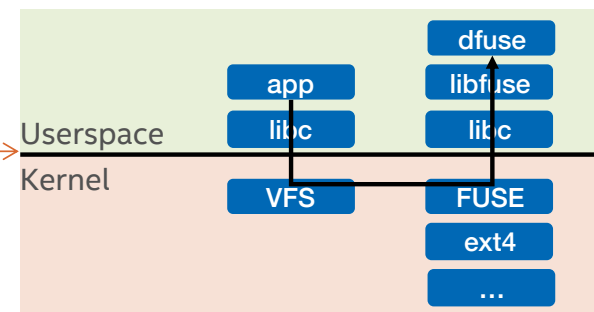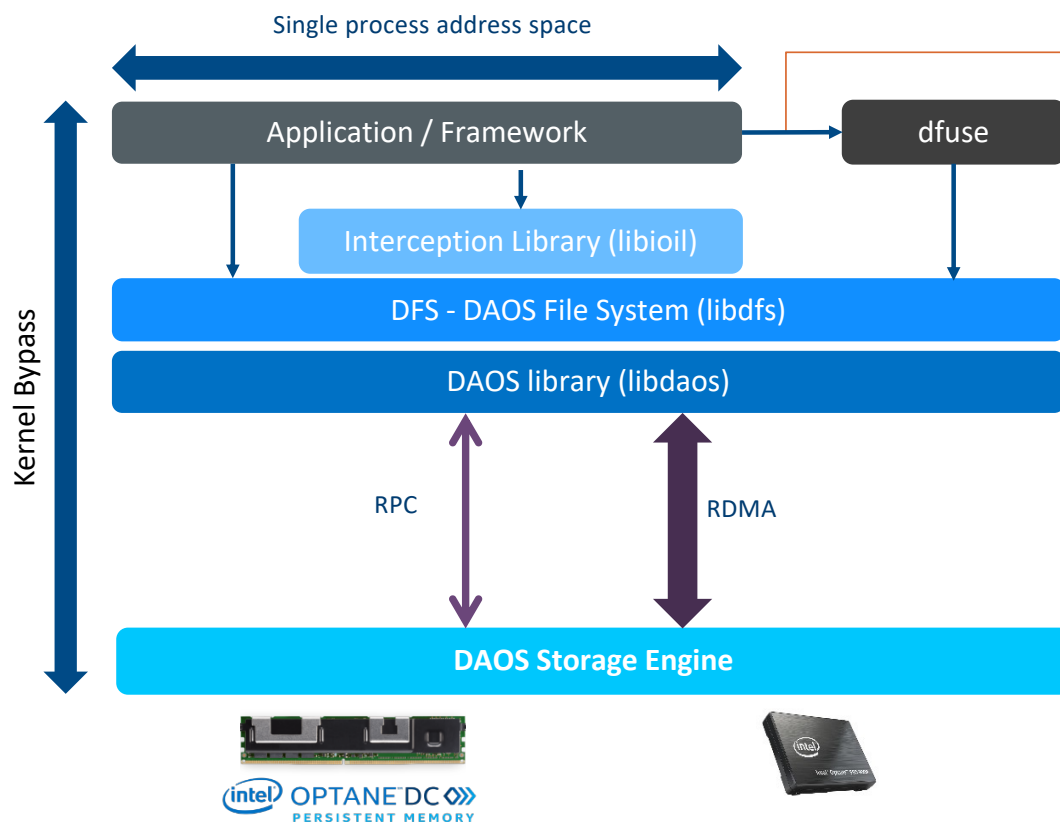- Writes during extension also go to new devices to ensure consistency

# DAOS Software Ecosystem

intel.

# DAOS Ecosystem

Generic I/O middleware supported today

Domain-specific data models under development in co-design with partners

Enablement in progress

*Compute Node*

AI/Analytics/Scientific Workflow

GPGPU

CPU

Apache Spark

Apache Hadoop

PyTorch

Legacy

Parallel Filesystem

HDF5

MPI-IO

S3

Hadoop Connector

Python

TensorFlow

SEGY

FDB

ROOT

DAQ

libdaos

Native array

Native key-value

RDMA

*DAOS Node*

Admin

DAOS Control Plane

DAOS Engine

Metadata + Smal I/Os

intel OPTANE ››› PERSISTENT MEMORY

Large I/Os

intel

# POSIX I/O Support

Single process address space



- User space DFS library with an API like POSIX.
  - Requires application changes (new API)
- DFUSE plugin to support POSIX API
  - No application changes
  - Limited performance
- DFUSE + IL
  - No application changes, runtime LD_PRELOAD
  - Good raw data I/O performance, limited metadata performance

# DFS API

| POSIX | DFS |
|---|---|
| mkdir(), rmdir() | dfs_mkdir(), dfs_rmdir() |
| open(), close(), access() | dfs_open(), dfs_release(),dfs_lookup() |
| pwritev(), preadv() | dfs_read/write() |
| {set,get,list,remove}xattr() | dfs_{set,get,list,remove}xattr |
| stat(), fstat() | dfs_stat(),ostat() |
| readdir() | dfs_readdir() |
| ... | ... |

Mostly 1-1 mapping from POSIX API to DFS API.

Instead of File & Directory descriptors, use DFS objects.

All calls need a DFS mount which is usually done on initialization with the pool / container access handles.

# PyDAOS Primer

- python module primarily written in C
  - Expose DAOS key-value store objects as a python dictionary
    - Support python iterator, direct assignments, …
    - Bulk insert/retrieve
  - Other data structures are under consideration (see later)
- Python objects allocated by PyDAOS:
  - are **persistent**
    - identified by a string name
  - are immediately **visible** upon creation
    - to any process running on the same or a different node.
  - have a **very low memory footprint** since the actual content is stored remotely
    - This allows to manipulate gigantic datasets way bigger than the amount of memory available on the node
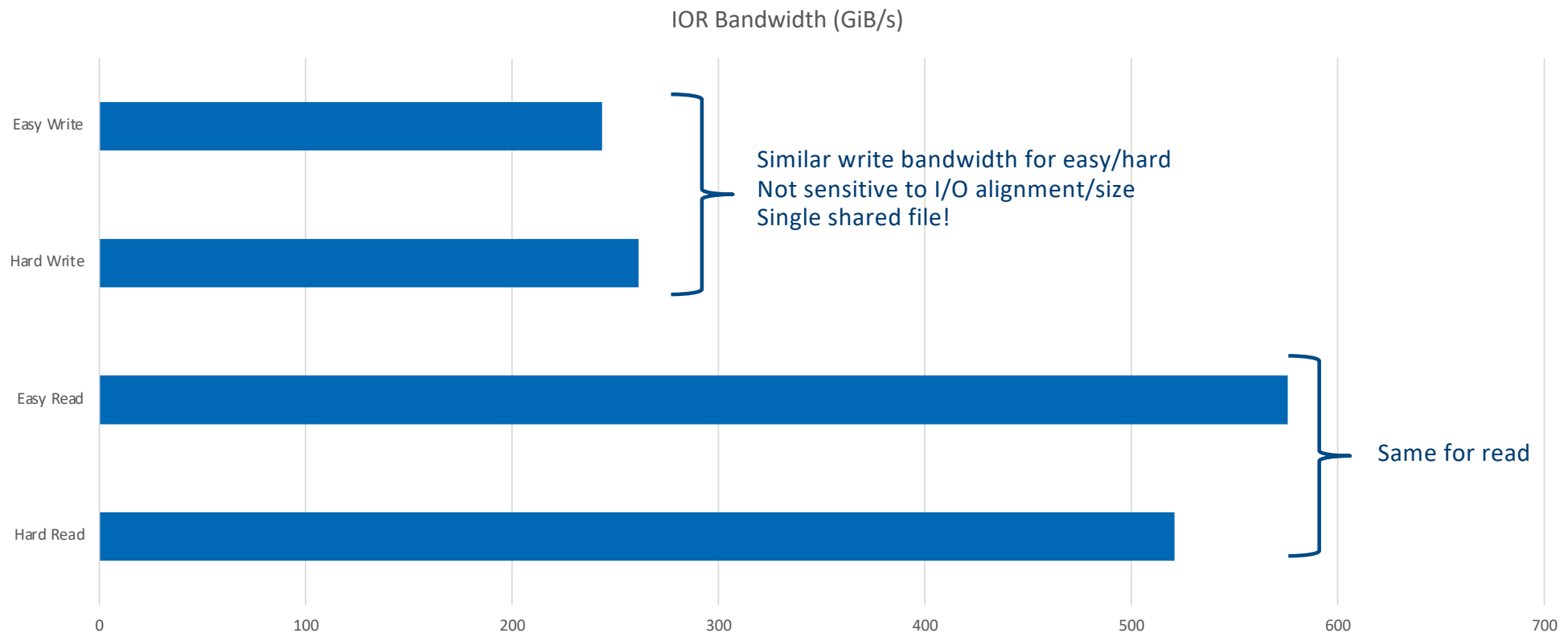
# TensorFlow Integration

- Done via TensorFlow-IO (see https://www.tensorflow.org/io)

- Initial integration at the DFS level

  - Use TF-IO filesystem API

  - Compatible with POSIX container

  - Provide full OS bypass

- Development completed & testing underway

  - https://github.com/daos-stack/tensorflow-io-daos/tree/devel

  - https://github.com/daos-stack/tensorflow-io-daos/blob/devel/docs/daos_tf_docs.md
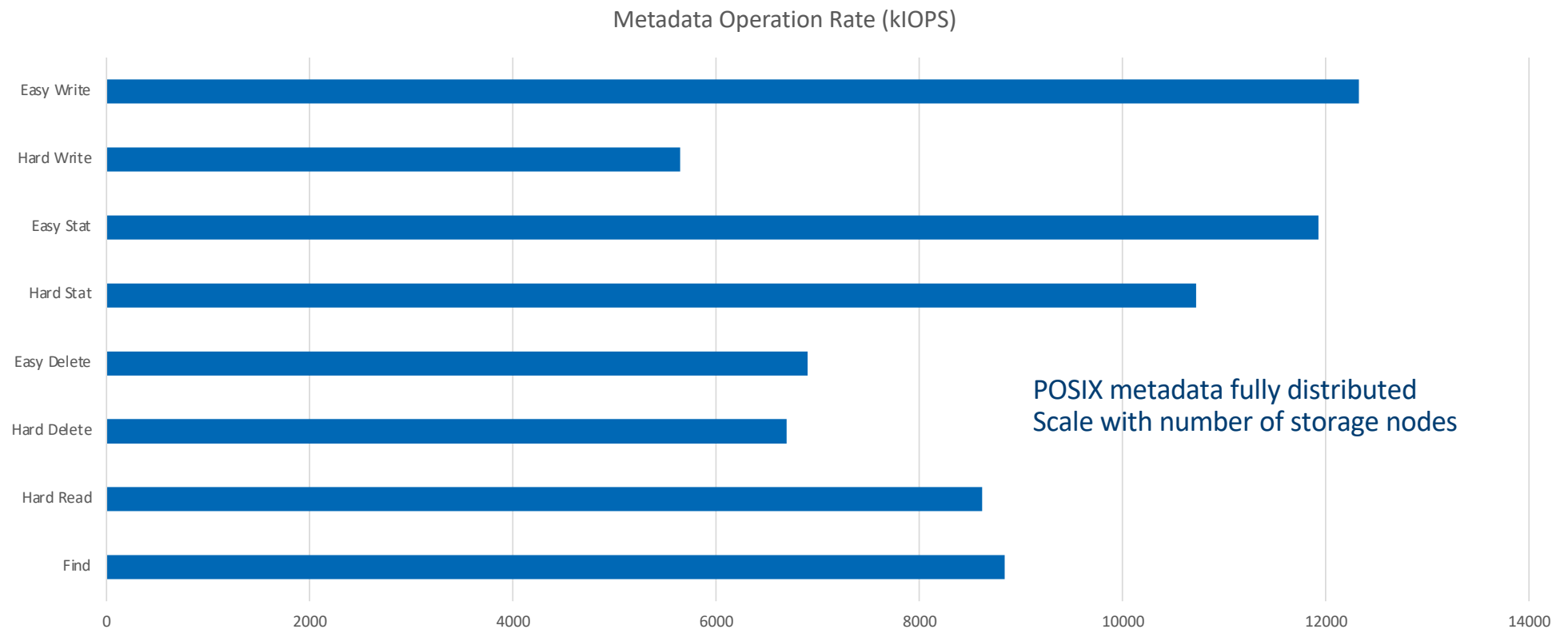
# DAOS Performance

intel.

# DAOS Bandwidth on IO500

IO⁵⁰⁰

IOR Bandwidth (GiB/s)



Similar write bandwidth for easy/hard
Not sensitive to I/O alignment/size
Single shared file!

Same for read

intel.

# DAOS Metadata Performance on IO500

IO$^{500}$

Metadata Operation Rate (kIOPS)



POSIX metadata fully distributed
Scale with number of storage nodes

# Resources



- ▪ Community Resources
  - Github: https://github.com/daos-stack/daos
  - Online doc: http://daos.io
  - Mailing list & slack: https://daos.groups.io
  - YouTube channel: http://video.daos.io
- ▪ 5$^{th}$ DAOS User Group (DUG'21)
  - Recordings available at http://dug.daos.io
- ▪ Intel landing page
  - https://www.intel.com/content/www/us/en/high-performance-computing/daos.html

# Q & A

intel.