# DAOS 2.4 and Beyond

Johann Lombardi, Senior Principal Engineer, AXG, Intel
6th DAOS User Group, Dallas, Nov 22, 2022

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.
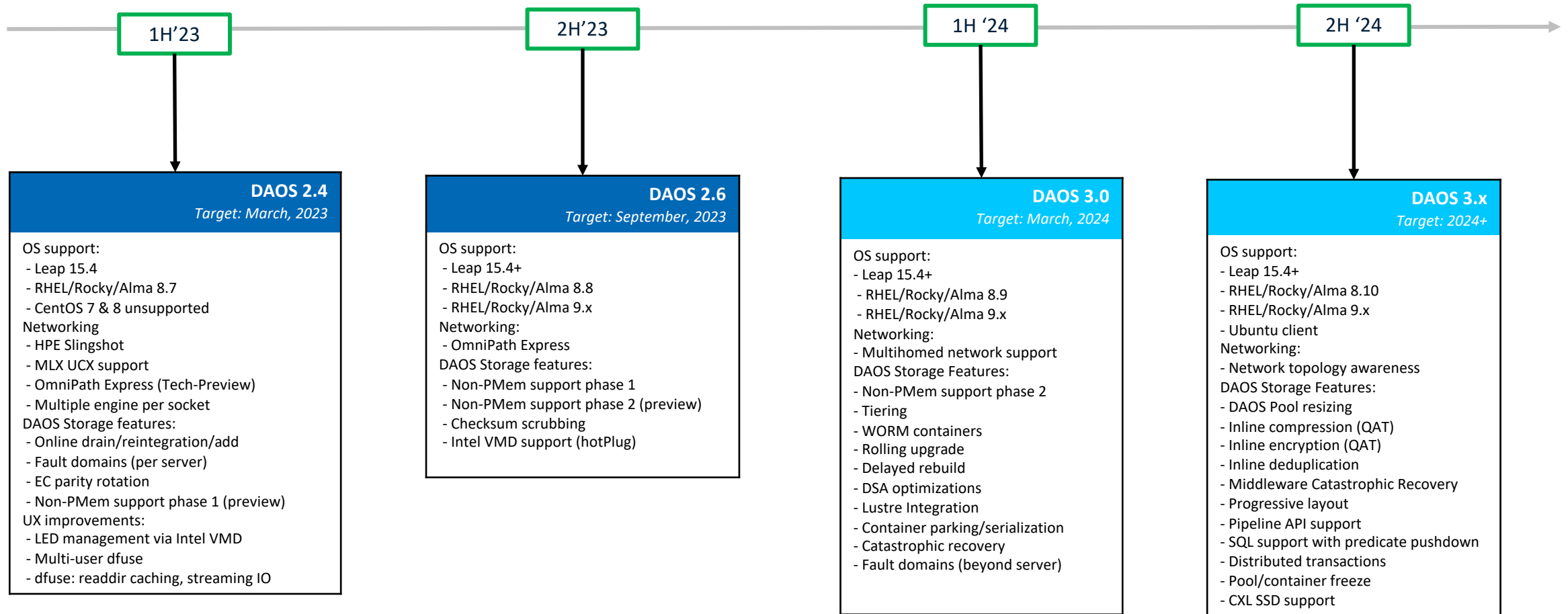
Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, 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.
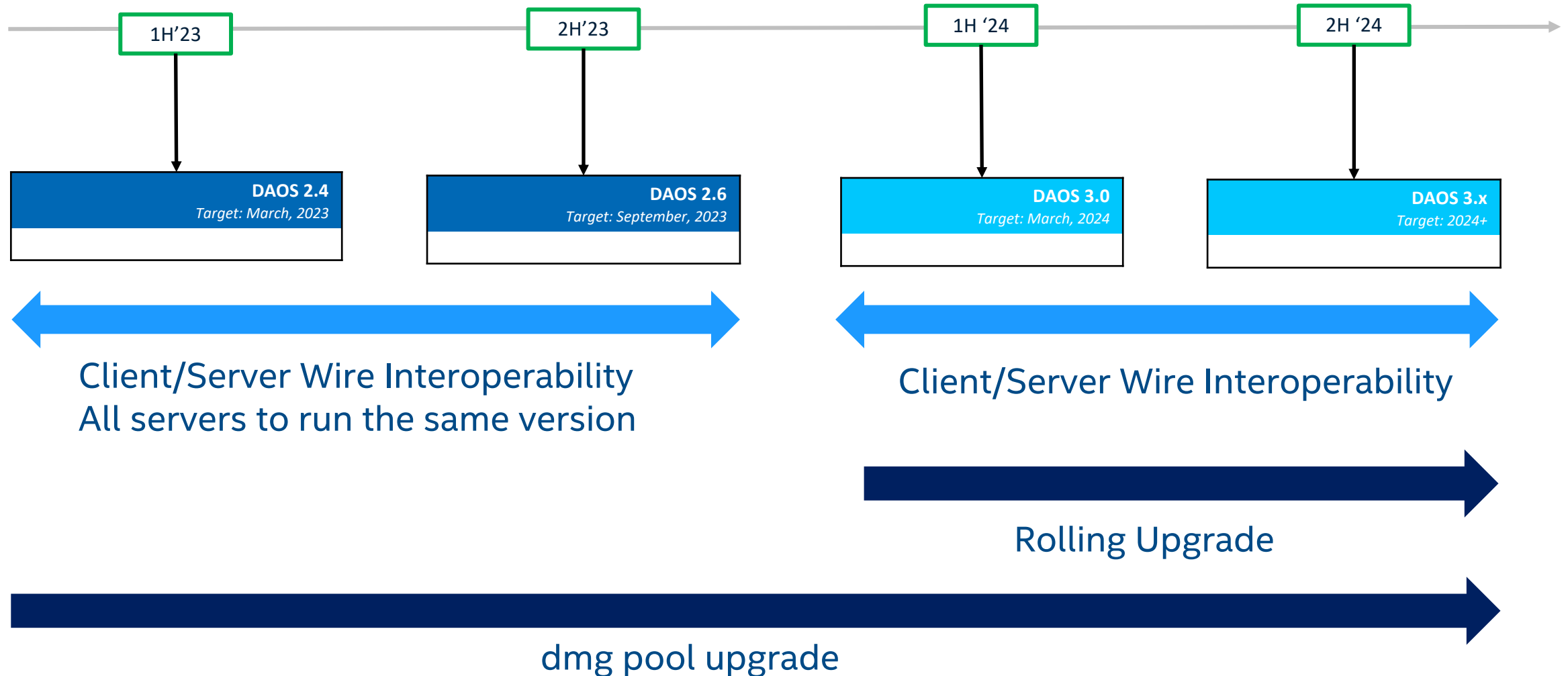
# DAOS Upcoming Community Releases

**1H'23**     **2H'23**     **1H '24**     **2H '24**

## DAOS 2.4
*Target: March, 2023*

OS support:
- Leap 15.4
- RHEL/Rocky/Alma 8.7
- CentOS 7 & 8 unsupported

Networking
- HPE Slingshot
- MLX UCX support
- OmniPath Express (Tech-Preview)
- Multiple engine per socket

DAOS Storage features:
- Online drain/reintegration/add
- Fault domains (per server)
- EC parity rotation
- Non-PMem support phase 1 (preview)

UX improvements:
- LED management via Intel VMD
- Multi-user dfuse
- dfuse: readdir caching, streaming IO

## DAOS 2.6
*Target: September, 2023*

OS support:
- Leap 15.4+
- RHEL/Rocky/Alma 8.8
- RHEL/Rocky/Alma 9.x

Networking:
- OmniPath Express

DAOS Storage features:
- Non-PMem support phase 1
- Non-PMem support phase 2 (preview)
- Checksum scrubbing
- Intel VMD support (hotPlug)

## DAOS 3.0
*Target: March, 2024*

OS support:
- Leap 15.4+
- RHEL/Rocky/Alma 8.9
- RHEL/Rocky/Alma 9.x

Networking:
- Multihomed network support

DAOS Storage Features:
- Non-PMem support phase 2
- Tiering
- WORM containers
- Rolling upgrade
- Delayed rebuild
- DSA optimizations
- Lustre Integration
- Container parking/serialization
- Catastrophic recovery
- Fault domains (beyond server)

## DAOS 3.x
*Target: 2024+*

OS support:
- Leap 15.4+
- RHEL/Rocky/Alma 8.10
- RHEL/Rocky/Alma 9.x
- Ubuntu client

Networking:
- Network topology awareness

DAOS Storage Features:
- DAOS Pool resizing
- Inline compression (QAT)
- Inline encryption (QAT)
- Inline deduplication
- Middleware Catastrophic Recovery
- Progressive layout
- Pipeline API support
- SQL support with predicate pushdown
- Distributed transactions
- Pool/container freeze
- CXL SSD support

*NOTE: All information provided in this roadmap is subject to change without notice.*

intel

# Interoperability / Upgradability

# Network: Buffer Mgmt Improvements (2.4)

- Enable MR (memory registration) cache
  - Required for CXI support
  - Dfuse changes to be MR cache-friendly in progress
- OFI MULTI_RECV support
  - Faced issue with buffer exhaustion at larger scale
  - Change Mercury to support OFI multi-receive API
    - Post large buffers where incoming messages are appended
    - Protocol changes since no tag support
  - Supported by most providers

# Network: UCX Support (2.4)

- Mellanox UCX library
- Future path for MLX HW support:
  - Shipped as part of MOFED packages
  - Better scalability on MLX fabric via DC & UD support
  - Native multi-rail support
  - GDS support on MLX/NVIDIA HW
- Shoud eventually replace OFI verbs provider
- Feature preview in 2.2 and official support in 2.4
- More work to be done to demonstrate scalability

**DAOS**

**CART**

**Mercury**

**OFI** **UCX**

intel.

# Network: Multi-Homed Network (3.0)

# Network: Multi-Homed Network (3.0)

# Backend: Metadata on SSD Requirements (2.4+)

- Incremental changes
  - Deliver incremental functionality w/o waiting for another 4y
- Offer continuity to our current and future customers
- Maintain performance leadership
- Turn this into an opportunity to:
  - Strenghten the DAOS open-source community
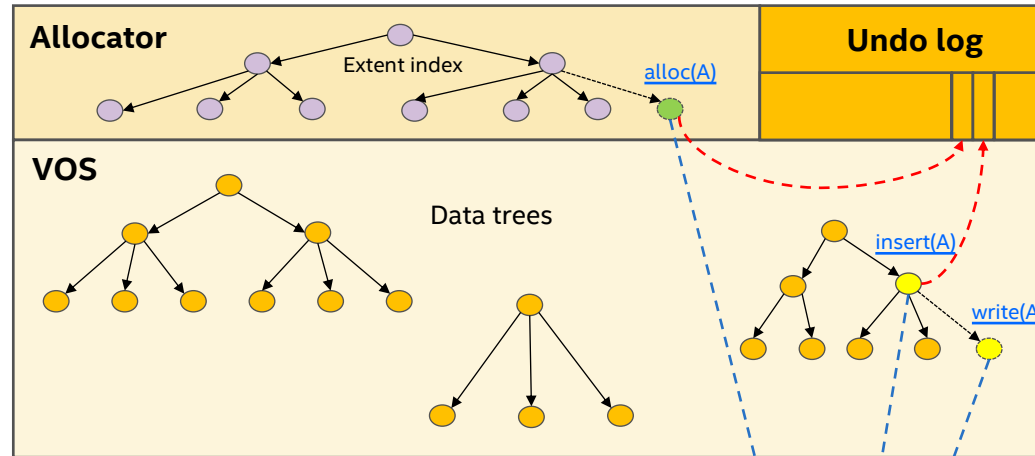  - Run DAOS on a wider range of hardware
  - Broaden the DAOS market
- https://daosio.atlassian.net/wiki/spaces/DC/pages/11196923911/Metadata+on+SSDs

# Backend: Persistent Memory



MMAP File

PMEM with PMDK
CXL.mem SSD

Data Blob

# Backend: Volatile Memory



MMAP File

DRAM/tmpfs

Synchronous Commit

Asynchronous Checkpointing

Meta Blob

WAL

Data Blob

# Backend: New (Meta)Data Structures



**md.mem**

Allocator — Extent index — alloc(A)

Undo log

Undo log is **off-heap**, it has no persistent mirror

VOS — Data trees — insert(A) — write(A)

**md.blob**

Allocator — Extent index

WAL

WAL is on-heap **(or in a separate blo**b),  it has no ephemeral mirror.

VOS — Data trees

# Backend: Undo vs Redo



WAL (redo)

. . . . . .

. operation = set_val(30)
. address = 0x2004
. size = 4

. operation = allocate
. address = 0x6000
. size = 64

. operation = set_val(5)
. address = 0x6000
. size = 4

. . . . . .

Overwrite value

x= 2

addr is 0x2004
x = 30

Allocate and assign value

y

addr = alloc()
addr is 0x6000
size is 64
y = 5

UMEM(undo)

. . . . . .

. operation = set_val(2)
. address = 0x2004
. size = 4

. operation = allocate
. address = 0x6000
. size = 64

. . . . . .

intel.

13

# Backend: Stack Layering (2.2)

# Backend: Stack Layering (2.4+)

# Erasure Code: Parity Rotation (2.4)



IOR Easy Read

IOR Easy Write
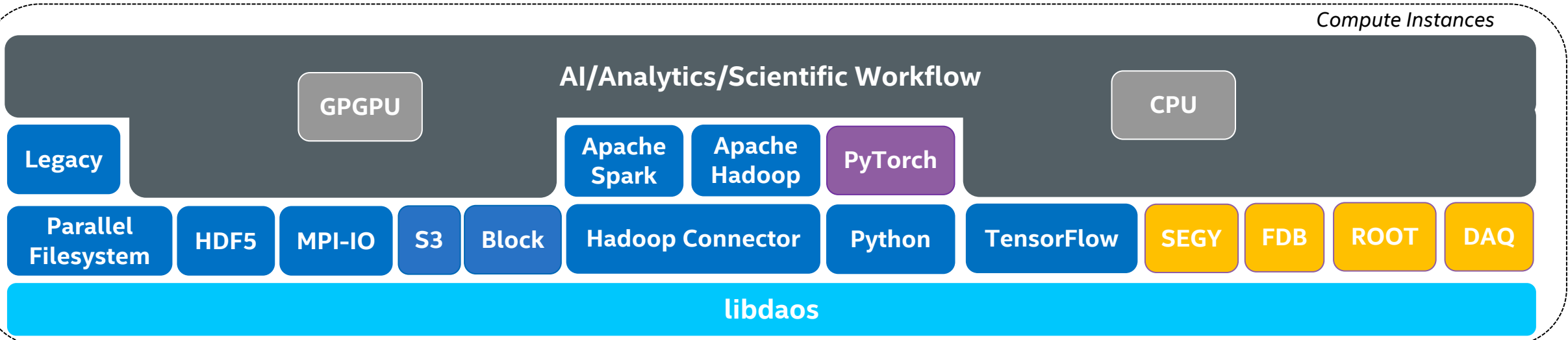
# Middleware: State of Affairs



**Compute Instances**

**AI/Analytics/Scientific Workflow**

GPGPU · CPU

Legacy · Apache Spark · Apache Hadoop · PyTorch

Parallel Filesystem · HDF5 · MPI-IO · S3 · Block · Hadoop Connector · Python · TensorFlow · SEGY · FDB · ROOT · DAQ

**libdaos**

Native array · Native key-value · RDMA

Generic I/O middleware supported today

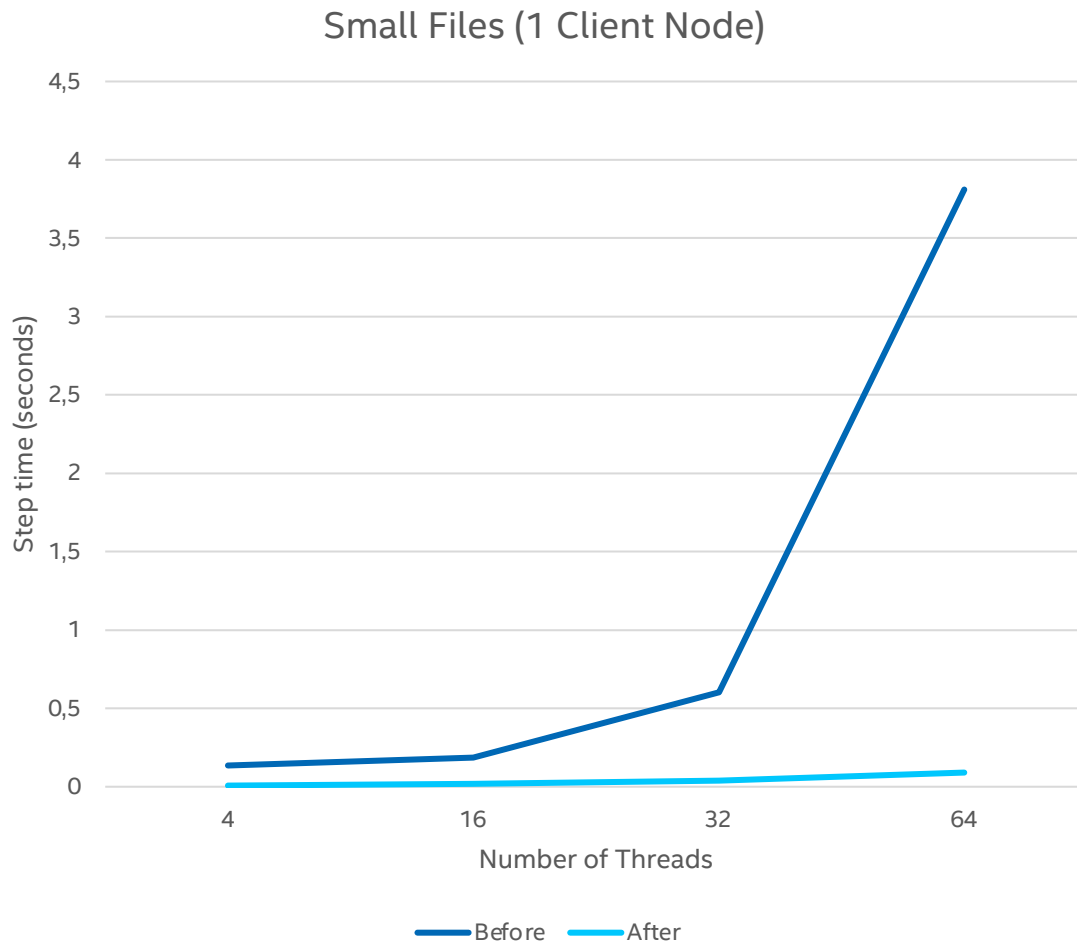Domain-specific data models under development in co-design with partners

Enablement in progress

# Middleware: dfuse & IL Improvements (2.4)

- Multi-user dfuse
  - Mountpoint owned by root and accessible by all users
  - Unix permissions apply
  - Can be used with IL, but no global handle
- Interception of streaming functions
  - f{open,read,write,close,…}
  - No caching
- Aggressive dfuse caching
  - Readdir and other operations
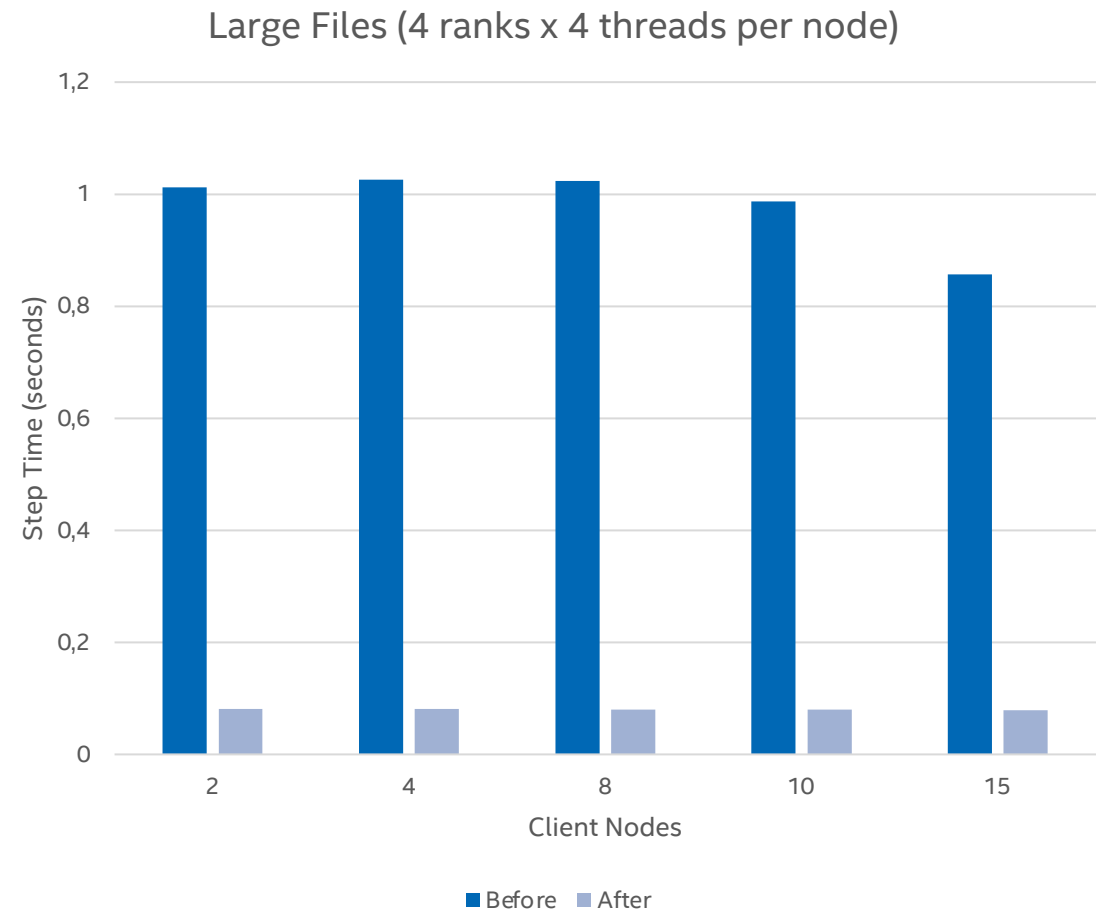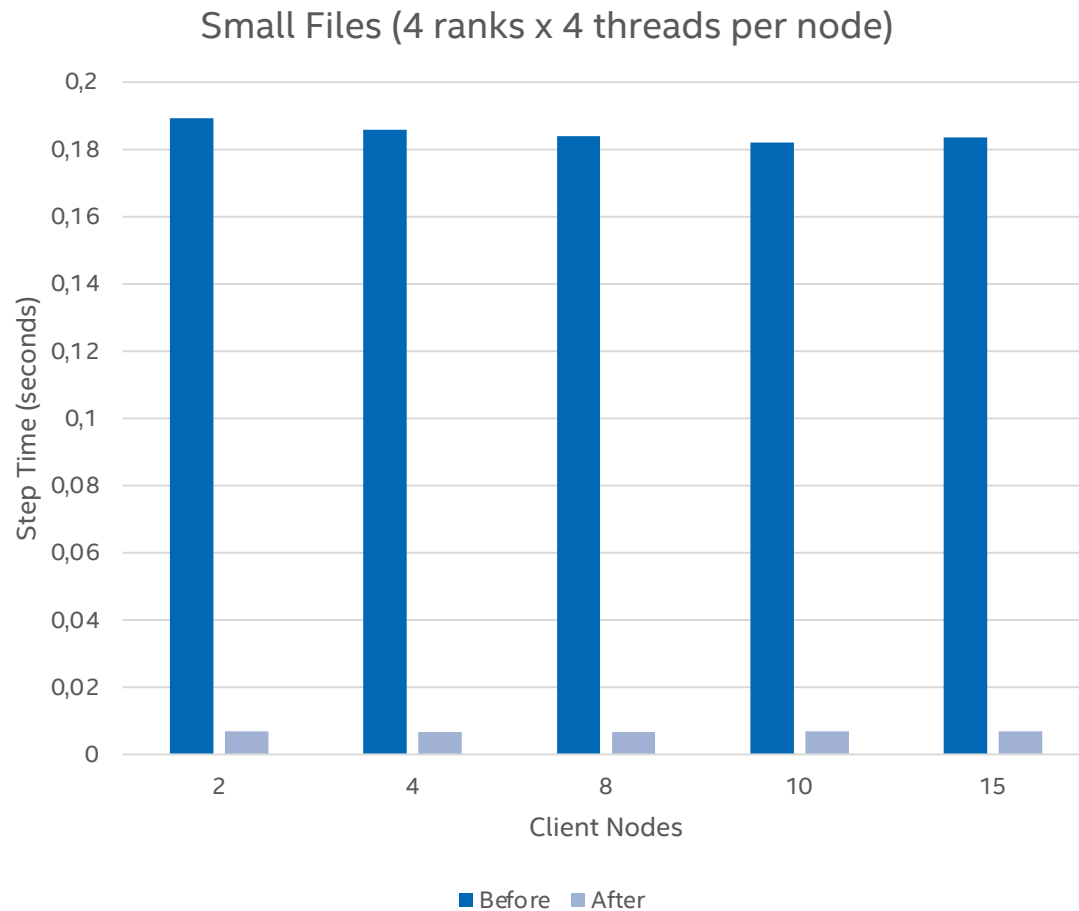
# Middleware: dfuse Caching AI Benchmark

- Benchmark to simulate some AI / ML workloads
  - 2 Datasets: large number of 4k files, small number of 4m files
  - x MPI ranks, each having y threads, all reading the dataset at the same time
    - Stat, open, read, close
- Measure on wolf using dfuse:
  - Before results: all default options
  - After results: do all the optimizations for aggressive caching
- Benchmark from LRZ: Durillo Barrionuevo, Juan; Hammer, Nicolay

intel.

# Middleware: dfuse Caching 1 Client

## Small Files (1 Client Node)

Step time (seconds) vs Number of Threads (4, 16, 32, 64)

— Before — After

## Large Files (1 Client Node)

Step time (seconds) vs Number of Threads (4, 16, 32, 64)

— Before — After

The information on this page is subject to the use and disclosure restrictions provided on the second page to this document.

intel. 20

# Middleware: dfuse Caching Many Clients



Small Files (4 ranks x 4 threads per node)

Large Files (4 ranks x 4 threads per node)

intel.

# Middleware: WORM Containers (3.0)

- Write Once Read Many Containers

- Optimizations for read-only datasets

  - Aggressive caching

  - Read-optimized layout

  - Size optimizations for immutable files

  - Indexing

- Use cases

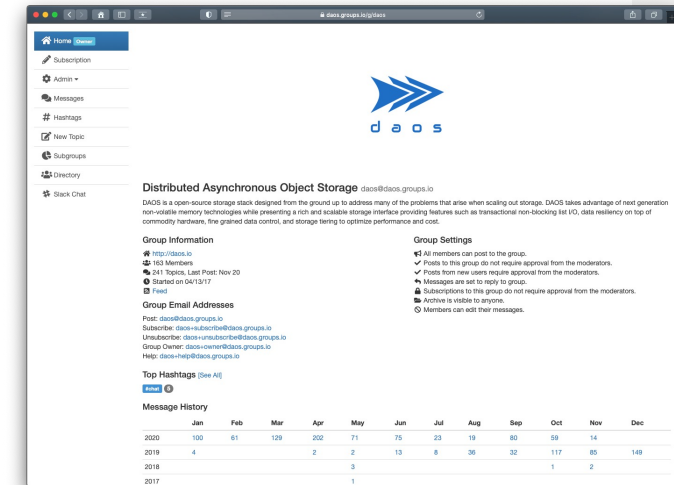  - Training/verification datasets
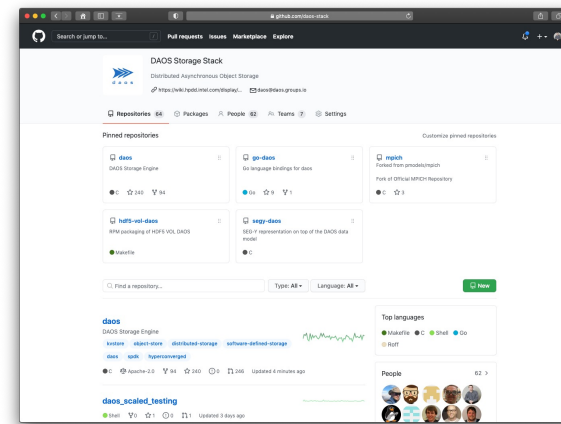
# Direction: Full Userspace POSIX Support

- Intercept all POSIX calls, including metadata operations
  - File open/create/stat/close and derivative
  - mmap via userfaultfd
- Collaborative caching
  - Use shared memory for data & metadata caching
  - Accessible by all ranks running on the same node
  - Size of cache configurable

# Direction: Scale-out Active Storage

- **DAOS pipeline API**
  - Offload data-intensive processing to storage
  - Pre-defined or user-defined (ubpf)
- **Leverage HW acceleration**
  - computational storage devices
  - accelerators/smartNICs
- **Many use cases**
  - POSIX find(1)
  - SQL query / MariaDB prototype using predicate push-down
  - In-place data filtering/pre-processing/transformation for AI frameworks
  - Calculate max/min/sum/avg … or searching for specific pattern/value on metadata/data

**Compute**

Node #1

Node #2

Reduce

Filter

Node #N

Transform

**Performance Tier**

DAOS System

Container 1

Container 2

Container 3

SSDs

DAOS node #1

SSDs

DAOS node #Y

# Resources



- Open-source Community
  - 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
- 6th DAOS User Group (DUG'22)
  - Recordings will be available at http://dug.daos.io
- DAOS BoF Community at SC'22
- Intel landing page
  - https://www.intel.com/content/www/us/en/high-performance-computing/daos.html





intel.