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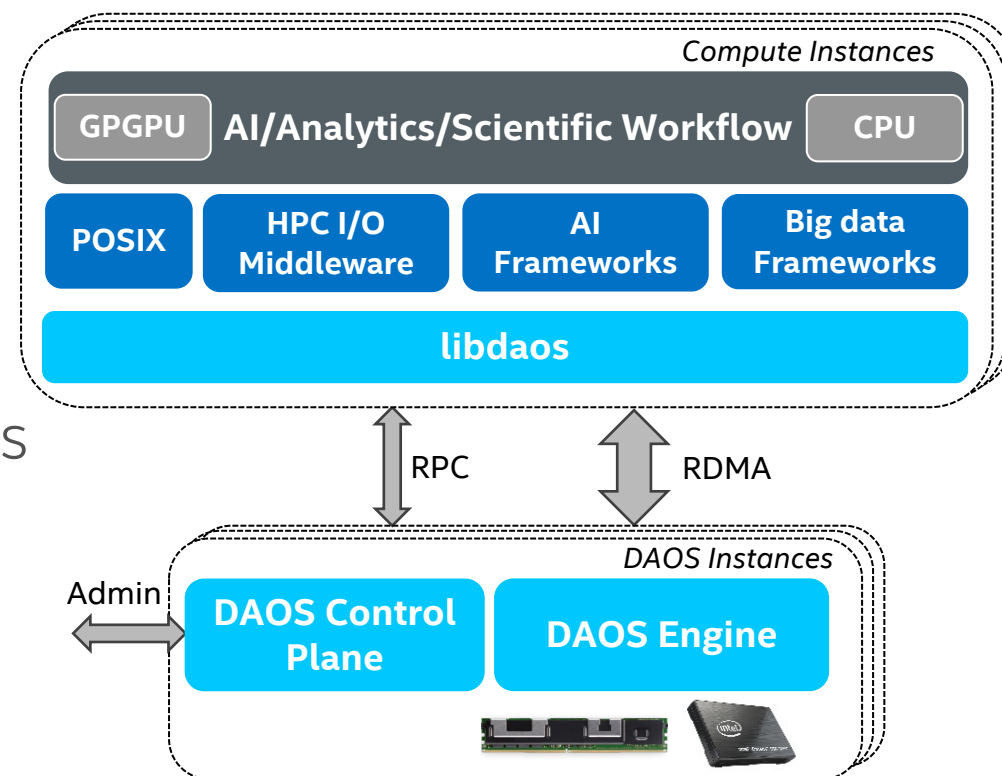
Intel does not control or audit third-party benchmark data or the web sites referenced in this document. You should visit the referenced web site and confirm whether referenced data are accurate.

# DAOS: Nextgen Open Storage Platform

- Fully Distributed multi-tenant global namespace

- Platform for innovation
  - Modular API and layering
  - Can leverage latest HW & SW technology

- Built for high performance
  - 10's μs latency, billions of IOPS, TB/s to PB/s

- Full userspace model
  - Run on-prem or in the cloud

- Growing open-source community

*Compute Instances*

| GPGPU | AI/Analytics/Scientific Workflow | CPU |

| POSIX | HPC I/O Middleware | AI Frameworks | Big data Frameworks |

**libdaos**

RPC    RDMA

*DAOS Instances*

Admin

**DAOS Control Plane**    **DAOS Engine**

# Middleware Ecosystem

**Compute Instances**

**AI/Analytics/Scientific Workflow**

GPGPU | CPU

| Apache Spark | Apache Hadoop | PyTorch |

| Legacy | MPI-IO | S3 | Block | TensorFlow | Hadoop Connector | Python | HDF5 | SEGY | FDB | ROOT | DAQ |

**libdfs (Parallel Filesystem)**

**Libdaos (key-value interface)**

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
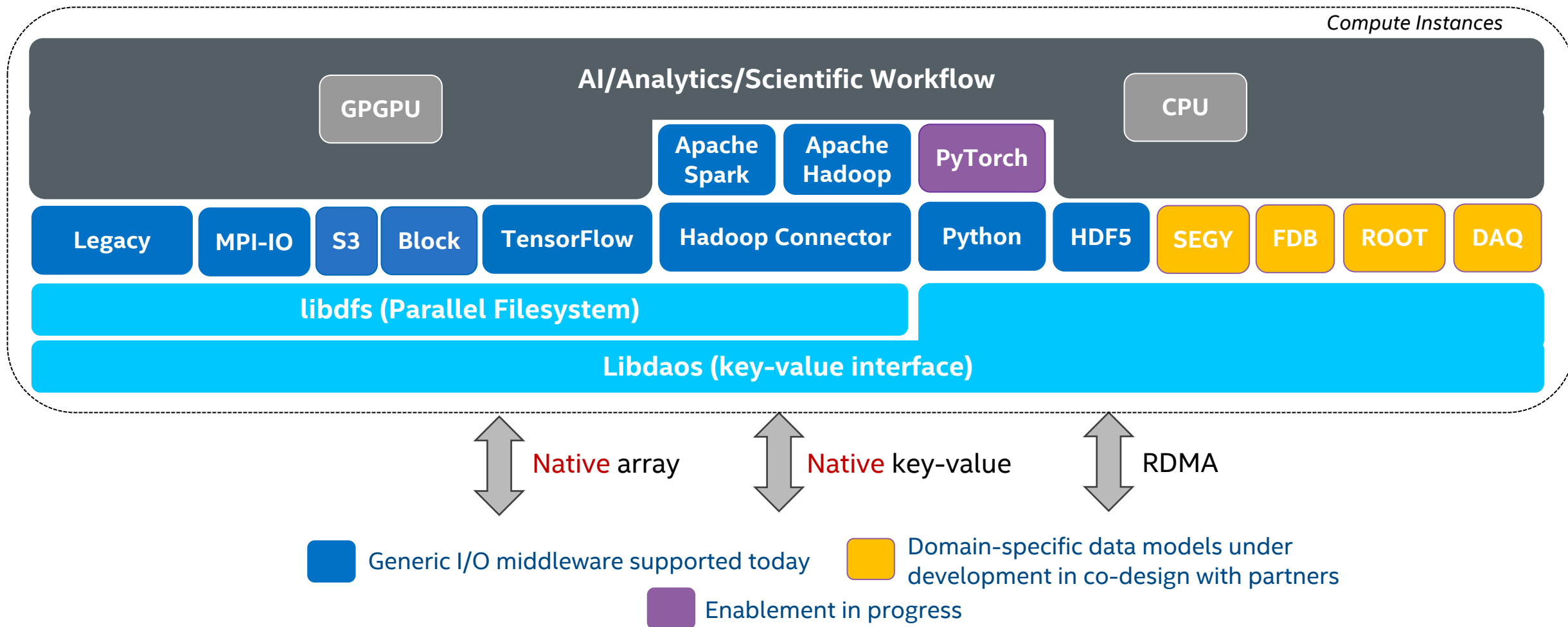
# APP: CosmoFlow

- A highly scalable Deep Learning application for Cosmology, using a 3D convolutional neural network trained on N-body cosmological simulation data to study dark matter in the universe
  - https://arxiv.org/pdf/1808.04728.pdf
  - Built over **TensorFlow**; uses **Horovod** for distributed training
- Community dataset: cosmoUniverse_2019_05_4parE_tf_v2
  - 1.7PiB
  - 524,288 samples for training
  - 65,536 samples for validation
  - Compressed TFRecord files

# APP: Workflow

# APP: TFRecord Loading

Compression ratio > 5x

| Read TFRecord File | Decompress File |
|---|---|
| <1 ms | >90 ms |

>90x ratio

Decompression coud be optimitized in future DAOS release with inline compression support with Intel® QAT

# The DAOS Exascale Storage Stack – Software Architecture

# DAOS Program Updates

DAOS to introduce two development paths:

- Stay on PMEM with transitioning to CXL2.0 Persistent memory third party products
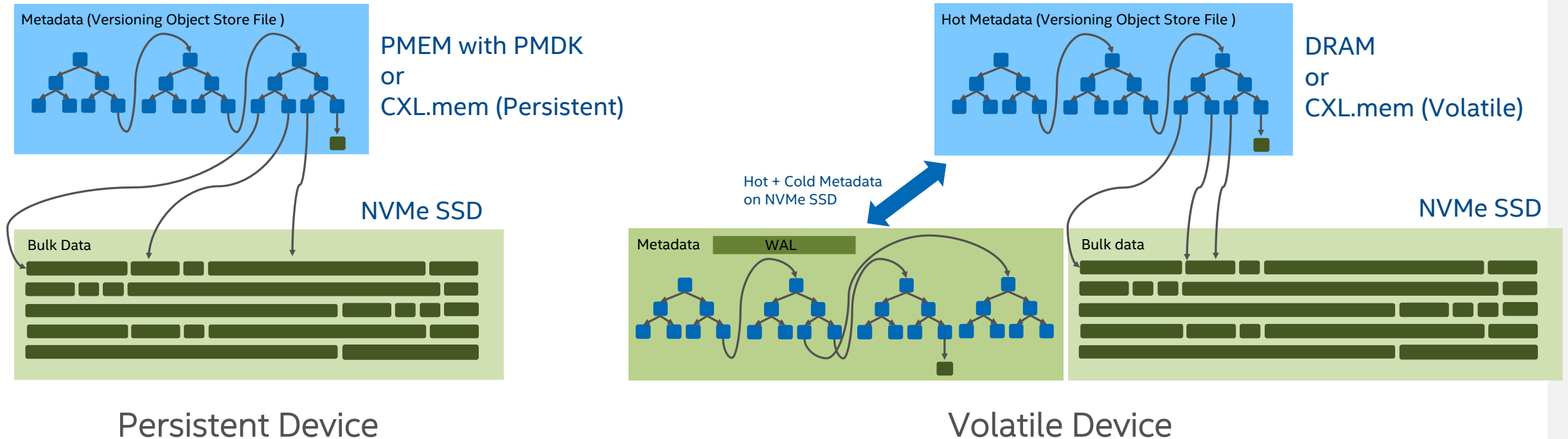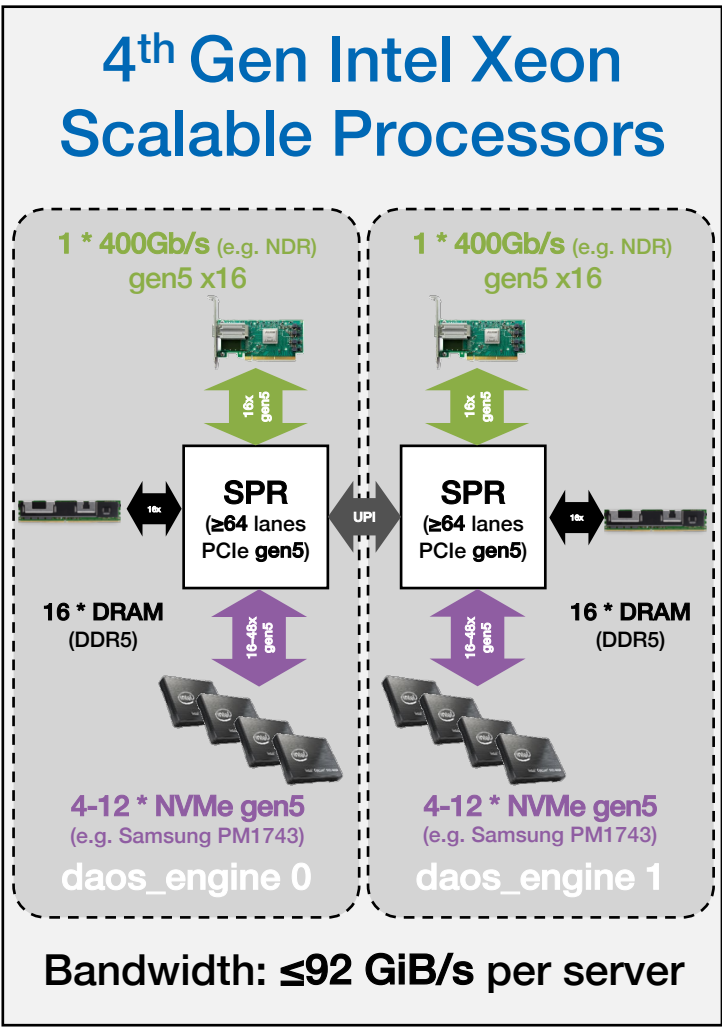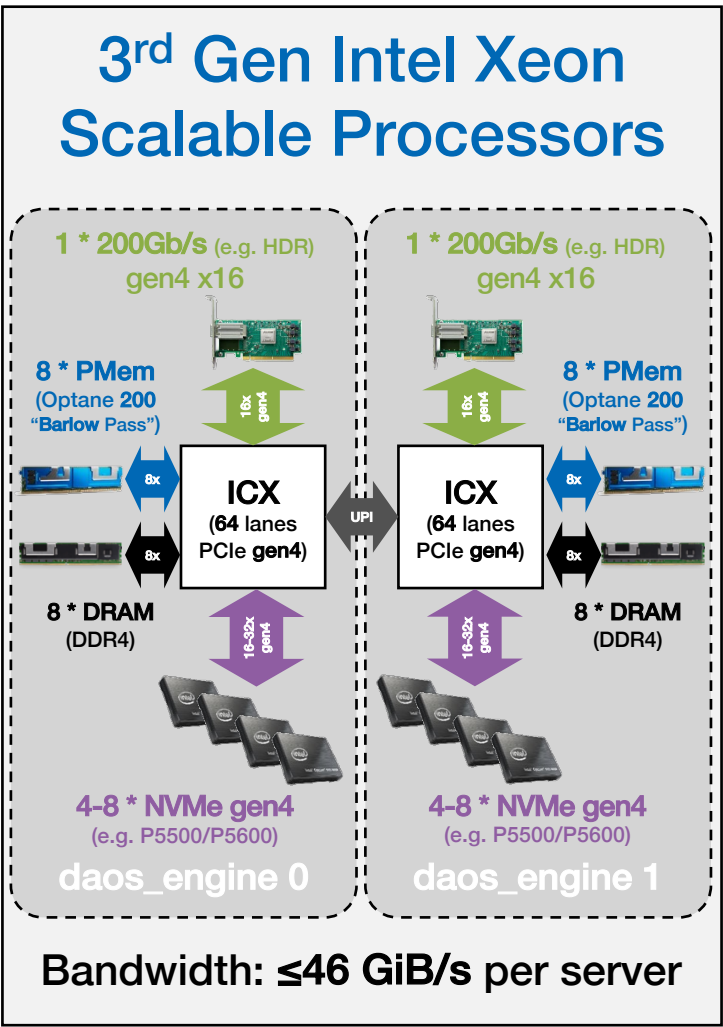- Implement DRAM / CXL.mem (volatile) + NVMe staging for metadata



Persistent Device

Volatile Device

# DAOS Servers on Intel Xeon Scalable Processors



## 2nd Gen Intel Xeon Scalable Processors

**1 * 100Gb/s** (e.g. EDR) gen3 x16

**6 * PMem** (Optane 100 "Apache Pass")

16x gen3

**CLX** (48 lanes PCIe gen3)

6x

UPI

**CLX** (48 lanes PCIe gen3)

6x

**6 * PMem** (Optane 100 "Apache Pass")

**1 * 100Gb/s** (e.g. EDR) gen3 x16

**6 * DRAM** (DDR4)

6x

6x

**6 * DRAM** (DDR4)

16-24x gen3

16-24x gen3

**4-6 * NVMe gen3** (e.g.P4500/P4600)

**4-6 * NVMe gen3** (e.g. P4500/P4600)

daos_engine 0

daos_engine 1

**Bandwidth: ≤23 GiB/s per server**

## 3rd Gen Intel Xeon Scalable Processors

**1 * 200Gb/s** (e.g. HDR) gen4 x16

**1 * 200Gb/s** (e.g. HDR) gen4 x16

**8 * PMem** (Optane 200 "Barlow Pass")

16x gen4

**ICX** (64 lanes PCIe gen4)

8x

UPI

**ICX** (64 lanes PCIe gen4)

8x

**8 * PMem** (Optane 200 "Barlow Pass")

16x gen4

**8 * DRAM** (DDR4)

8x

8x

**8 * DRAM** (DDR4)

16-32x gen4

16-32x gen4

**4-8 * NVMe gen4** (e.g. P5500/P5600)

**4-8 * NVMe gen4** (e.g. P5500/P5600)

daos_engine 0

daos_engine 1

**Bandwidth: ≤46 GiB/s per server**

## 4th Gen Intel Xeon Scalable Processors

**1 * 400Gb/s** (e.g. NDR) gen5 x16

**1 * 400Gb/s** (e.g. NDR) gen5 x16

16x gen5

**SPR** (≥64 lanes PCIe gen5)

1x

UPI

**SPR** (≥64 lanes PCIe gen5)

1x

16x gen5

**16 * DRAM** (DDR5)

**16 * DRAM** (DDR5)

16-48x gen5

16-48x gen5

**4-12 * NVMe gen5** (e.g. Samsung PM1743)

**4-12 * NVMe gen5** (e.g. Samsung PM1743)

daos_engine 0

daos_engine 1

**Bandwidth: ≤92 GiB/s per server**

# DAOS on Aurora



- **1024x** DAOS nodes, each with:
  - 2x Xeon 5320 CPUs
  - 512GB DRAM
  - 8TB Optane Persistent Memory 200
  - 244TB NVMe SSDs
  - 2x HPE Slingshot NIC
- **Usable** capacity
  - between 220PB and 249PB
    depending on redundancy level chosen



**DAOS Performance**
220 PB capacity @ EC16+2
≥ 25 TB/s

**Lustre Performance**
Grand – 100 PB @ 650 GB/s
Eagle – 100 PB @ 650 GB/s

**DAOS Nodes (DNs)**
Xeon servers
NVRAM and NVMe attached storage
DAOS service

Slingshot Fabric

System Service Nodes (SSNs)

User Access Nodes (UANs)

**Gateway Nodes (GNs)**
Xeon servers with no local storage
Access to external storage

Gateway nodes

**Scalable Storage Cluster (SSC)**
Xeon servers connected to JBOD
Lustre OSSs & MDSs

# IO500 SC22 Results

## IO500 SC22 List

| IO500 | 10 Node | Full | Historical | **Customize** | Download |

This is the SC22 IO500 list

| # ↑ | BOF | INSTITUTION | SYSTEM | STORAGE VENDOR | FILE SYSTEM TYPE | CLIENT NODES | TOTAL CLIENT PROC. | SCORE ↑ | BW (GIB/S) | MD (KIOP/S) |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | ISC21 | Pengcheng Laboratory | Pengcheng Cloudbrain-II on Atlas 900 | Pengcheng | MadFS | 512 | 36,864 | 36,850.40 | 3,421.62 | 396,872.82 |
| 2 | SC22 | Argonne National Laboratory | Aurora Storage | Intel | DAOS | 260 | 27,040 | 20,694.50 | 6,048.69 | 70,802.51 |
| 3 | SC22 | Sugon Cloud Storage Laboratory | ParaStor | Sugon | ParaStor | 10 | 2,560 | 8,726.42 | 718.11 | 106,042.93 |
| 4 | SC22 | SuPro Storteck | StarStor | SuPro Storteck | StarStor | 10 | 2,560 | 6,751.75 | 515.15 | 88,491.65 |
| 5 | SC22 | Tsinghua Storage Research Group | SuperStore | Tsinghua Storage Research Group | SuperFS | 10 | 1,200 | 5,517.73 | 179.60 | 169,515.95 |
| 6 | ISC22 | National Supercomputing Center in Jinan | Shanhe | PDSL | flashfs | 10 | 2,560 | 3,534.42 | 207.79 | 60,119.50 |
| 7 | SC22 | Cloudam HPC on OCI | HPC-OCI | Cloudam | BurstFS | 64 | 1,920 | 3,033.03 | 278.48 | 33,033.54 |
| 8 | SC21 | Huawei HPDA Lab | Athena | Huawei | OceanFS | 10 | 1,720 | 2,395.03 | 314.56 | 18,235.71 |
| 9 | SC21 | Olympus Lab | OceanStor Pacific | Huawei | OceanFS | 10 | 1,720 | 2,298.69 | 317.07 | 16,664.88 |
| 10 | SC21 | Huawei Cloud | | PDSL | Flashfs | 15 | 1,560 | 2,016.70 | 109.82 | 37,034.00 |

# DAOS Ecosystem



| | |
|---|---|
| Hardware Partners | Hewlett Packard Enterprise, Lenovo, intel, SUPERMICRO, inspur, Atos, NEC (Orchestrating a brighter world), QCT |
| Reseller Partners | COLFAX Customized Solutions, MEGWARE, RSC, E4 COMPUTER ENGINEERING, 2crsi, ScaleWorX |
| Software Development and 3rd party support | croit, SEAGATE, Brightskies, Google Cloud, Nettrix 宁畅, H3C 数字化解决方案领导者 |
| End Customers | JINR, Argonne NATIONAL LABORATORY, 瑞金 (Ruijin Hospital), UNIVERSITY OF CAMBRIDGE, lrz |

# 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 available at http://dug.daos.io

■ Upcoming BoFs at SC22

■ Intel landing page

- https://www.intel.com/content/www/us/en/high-performance-computing/daos.html