



E&G ECO ENGINEERING

# DAOS: SCALE-OUT SOFTWARE-DEFINED STORAGE FOR HPC/BIG DATA/AI CONVERGENCE

Johann Lombardi, Principal Engineer

# 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 computer system 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.

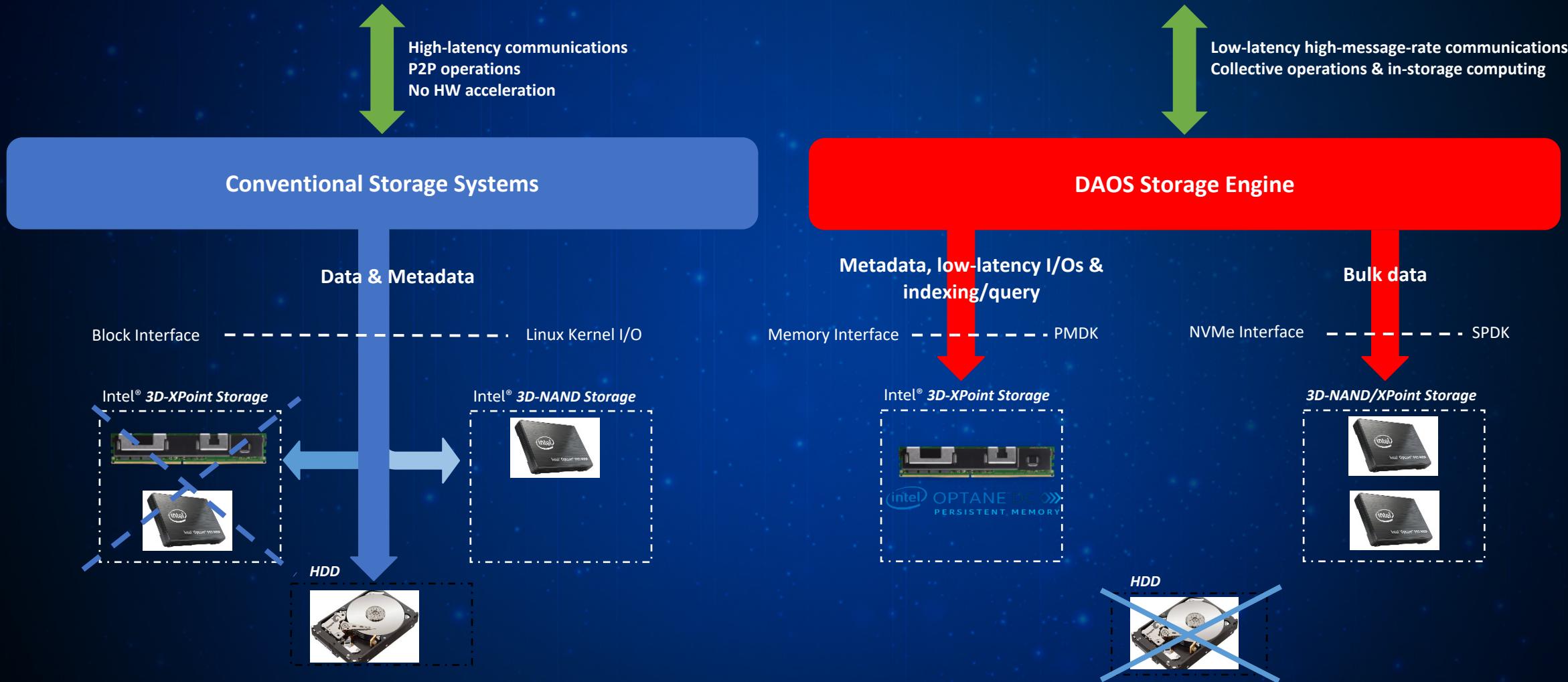
Cost reduction scenarios described are intended as examples of how a given Intel-based product, in the specified circumstances and configurations, may affect future costs and provide cost savings. Circumstances will vary. Intel does not guarantee any costs or cost reduction.

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.

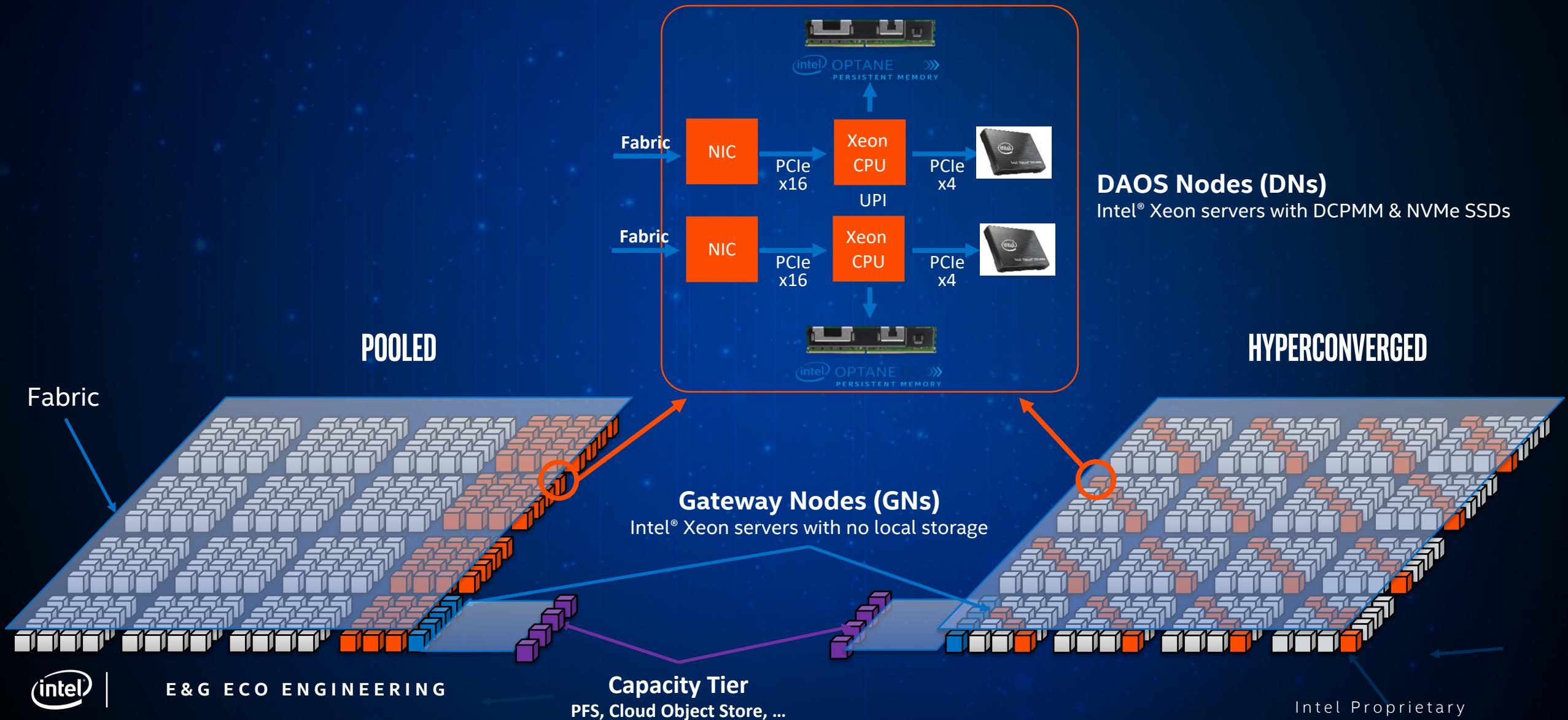
Intel, the Intel logo, Intel Optane and Intel Xeon are trademarks of Intel Corporation in the U.S. and/or other countries.

\*Other names and brands may be claimed as property of others.

# DAOS ARCHITECTURE



# DAOS DEPLOYMENTS



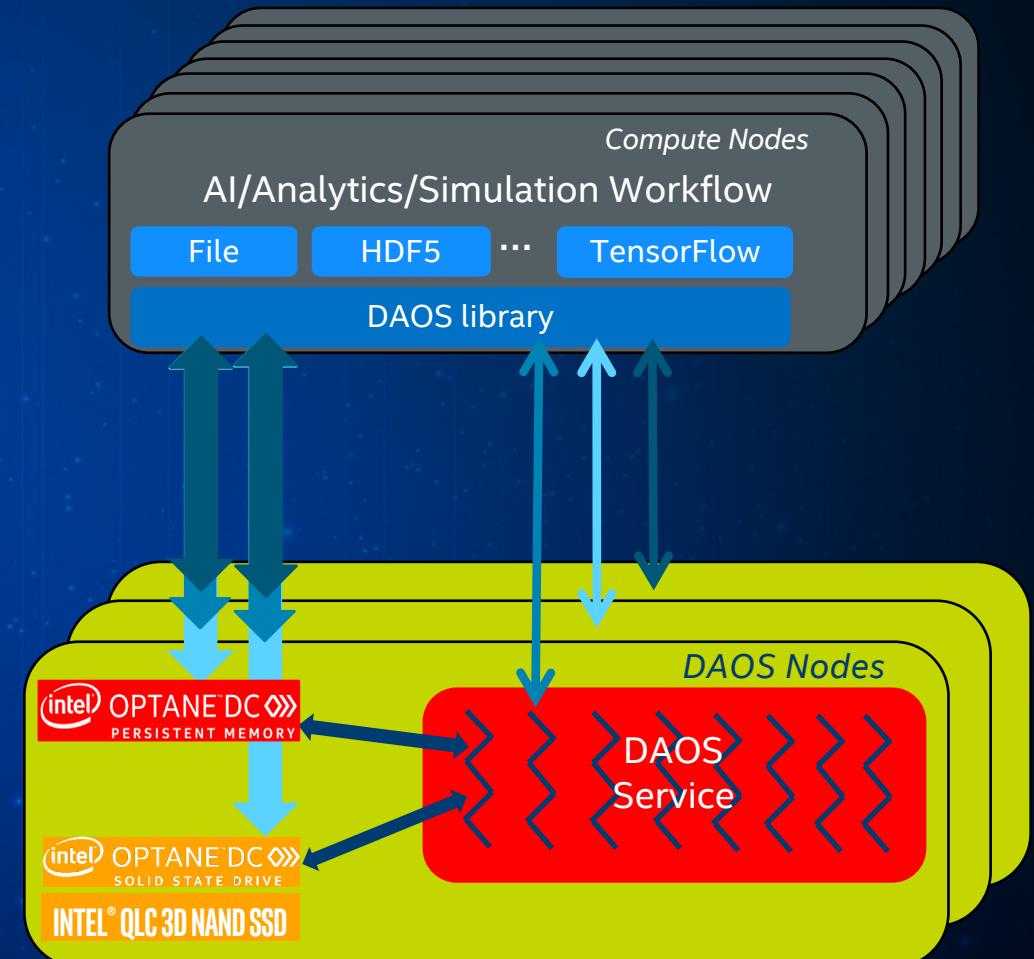
# DAOS TIER ANATOMY

## DAOS Tier

- Globally accessible from any compute nodes
- Large capacity (100's PB)

## DAOS Nodes

- COTS Intel® Xeon servers running the DAOS service
- RNIC attached for communications
  - Support multiple RNICs per server to sustain backend storage IOPS/bandwidth
- Mix of storage technologies attached
  - Intel® Optane™ DC Persistent Memory (DCPMM)
  - NVMe SSD (\*NAND, Intel® Optane™ SSDs)



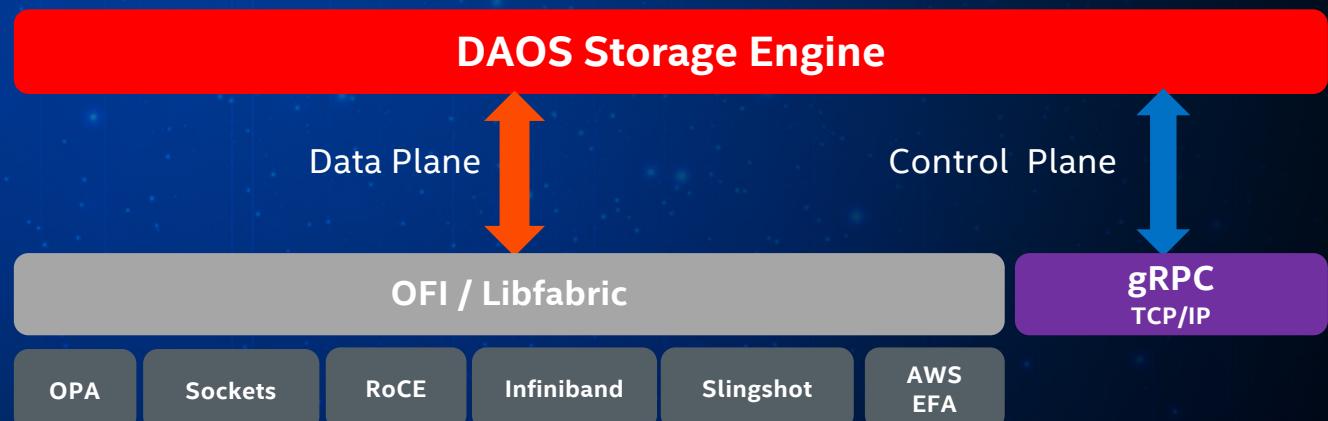
# NETWORK SUPPORT

Performance-critical I/O path over libfabric

- Low-latency messaging
  - End-to-end in userspace
- Native support for RDMA
  - True zero-copy I/O
- Non-blocking
- Scalable collective communications

Out-of-band channel for administration

- Manage hardware, service & pools
- Telemetry & troubleshooting
- Secured with TLS & certificate



# STORAGE VIRTUALIZATION & MULTI-TENANCY

## Distributed storage reservation

- Intel® Optane™ DC Persistent Memory (DCPMM)
- NVMe SSD

## Predictable capacity

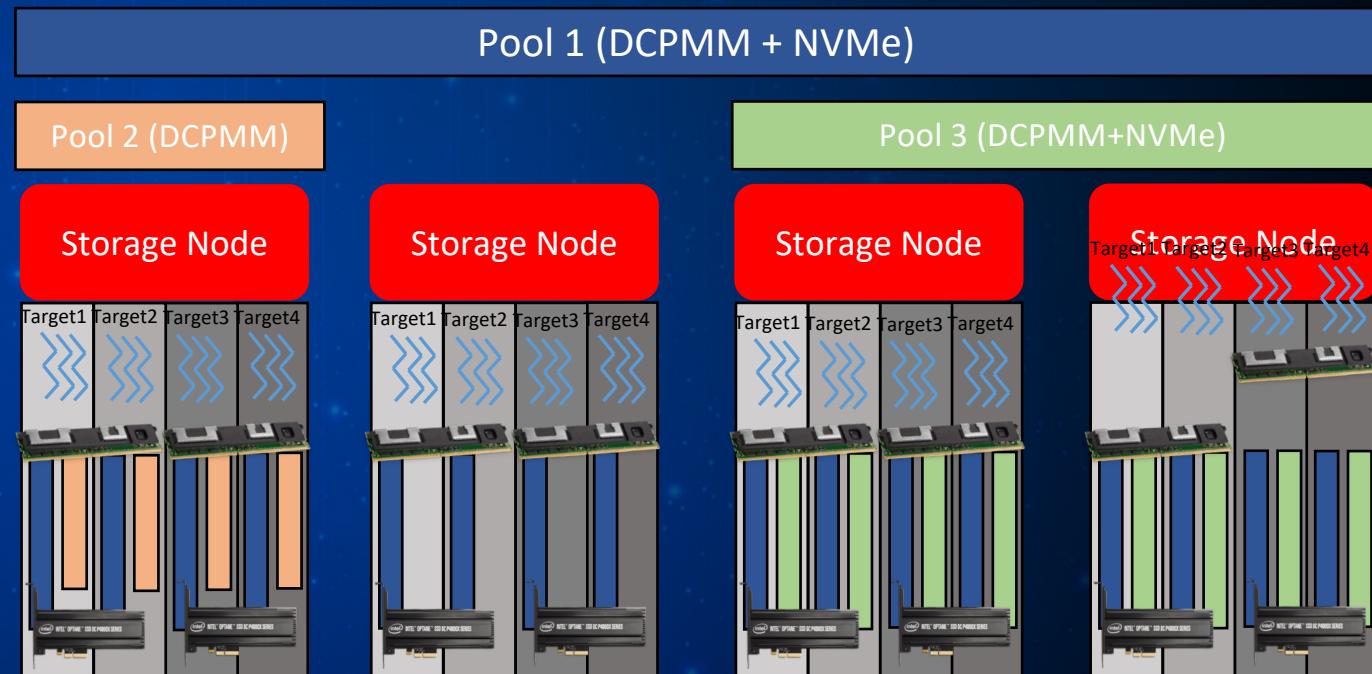
- Can be resized
- Can be extended to span more servers

## Multi-tenancy

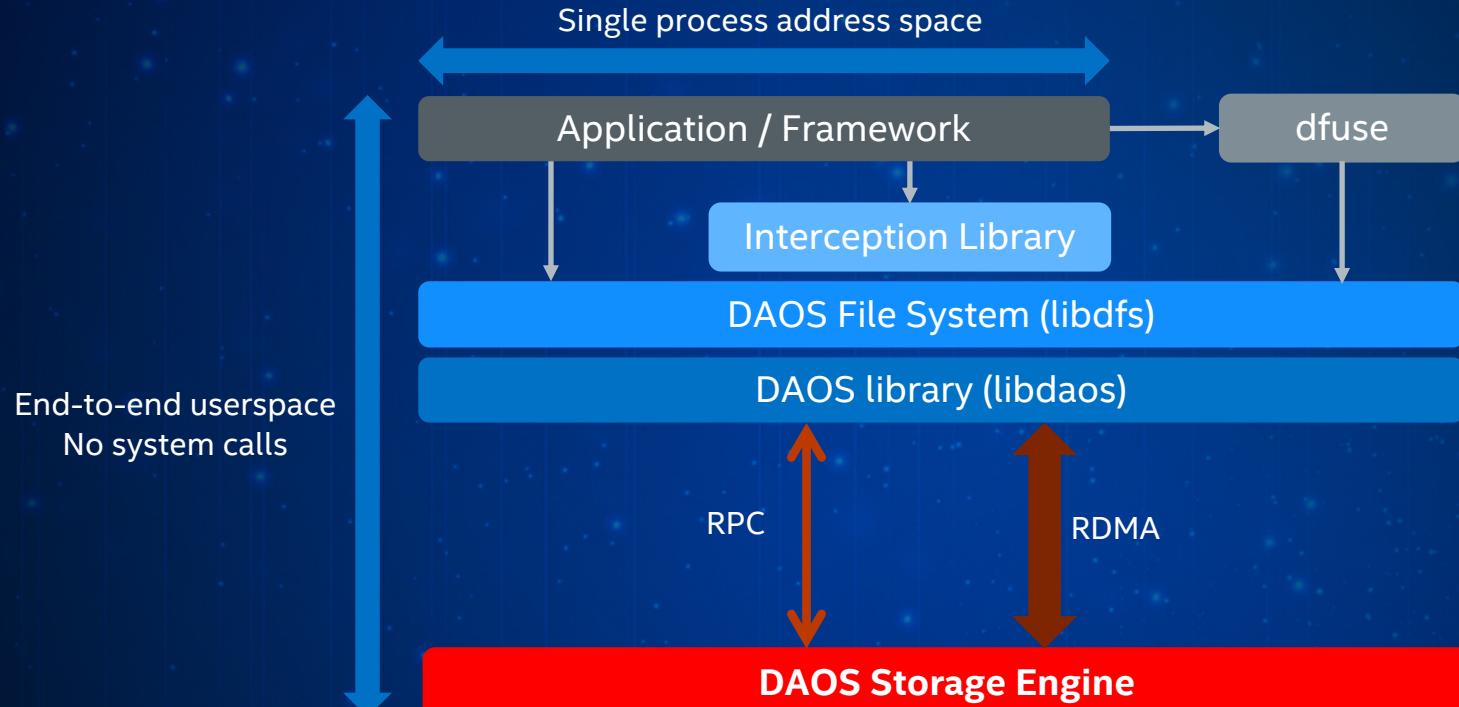
- NFSv4-type ACLs

Typically 1 pool = 1 project

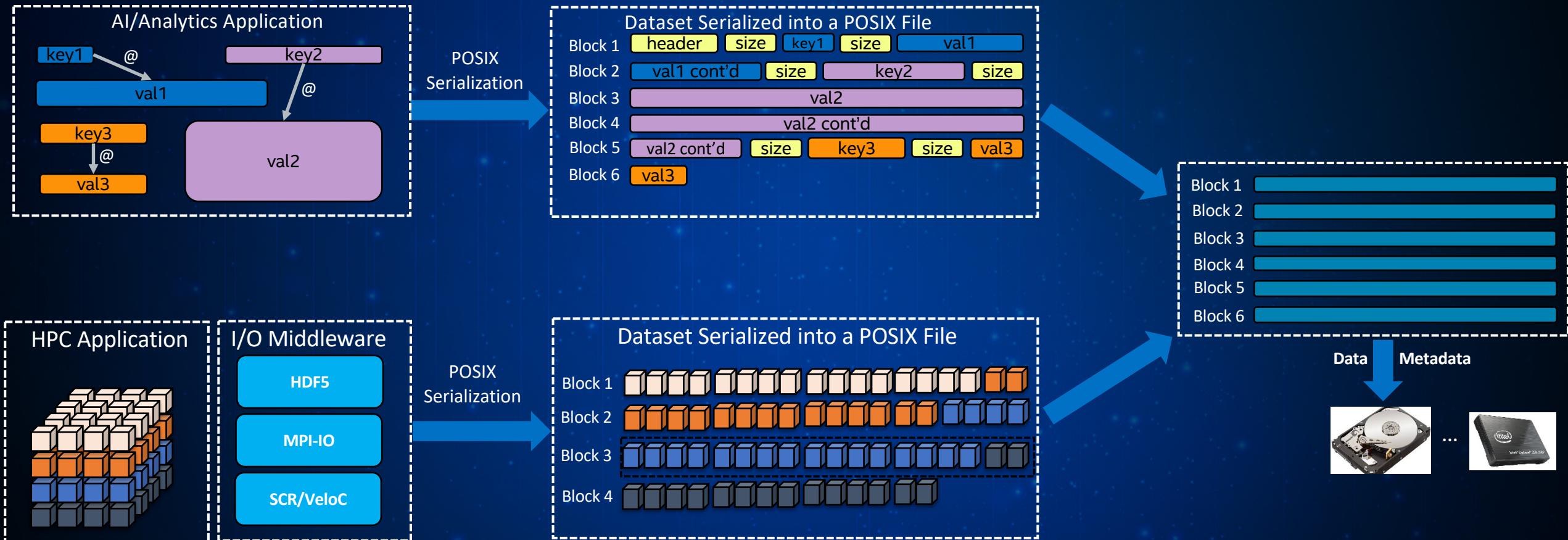
- Can have a single pool or 100's



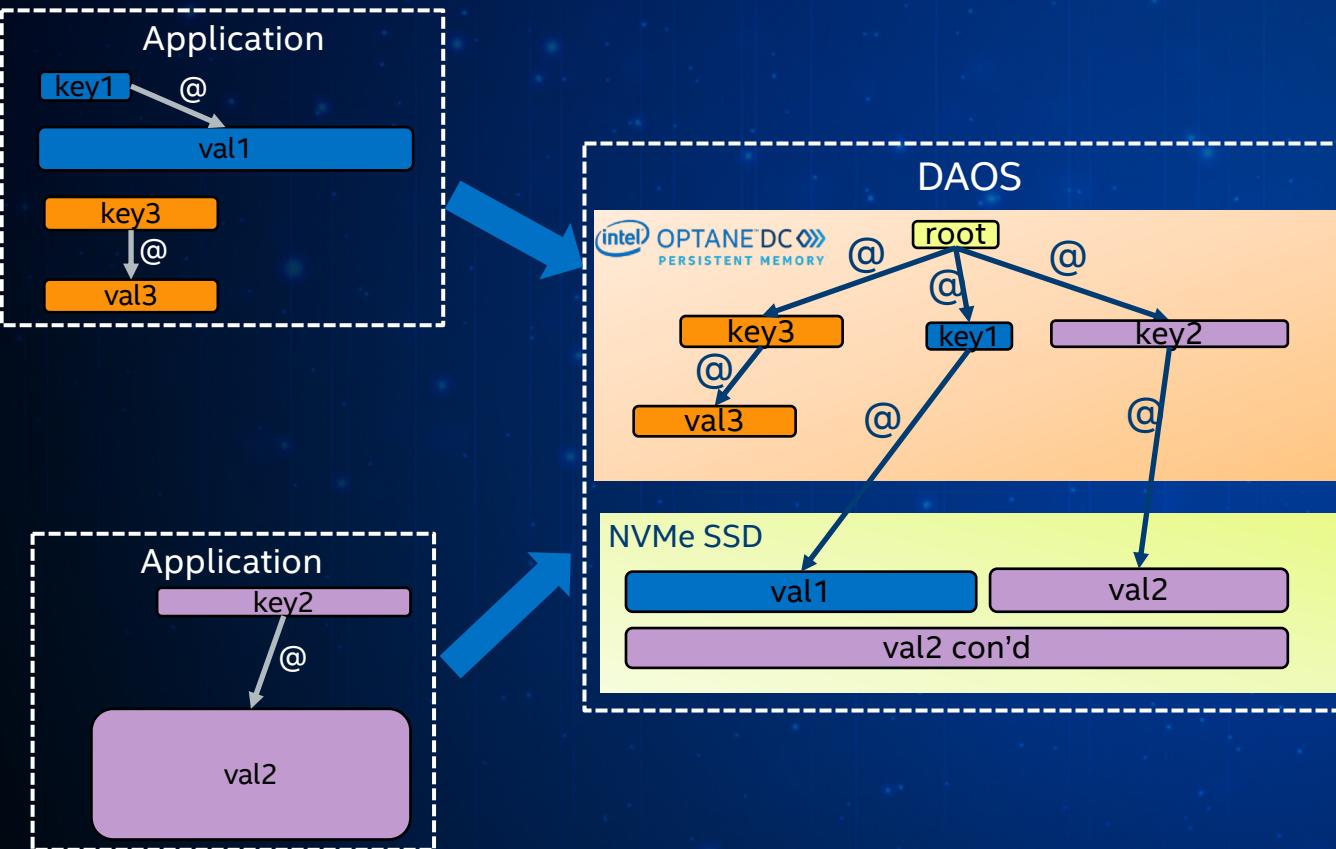
# POSIX I/O SUPPORT



# POSIX I/O LIMITATIONS

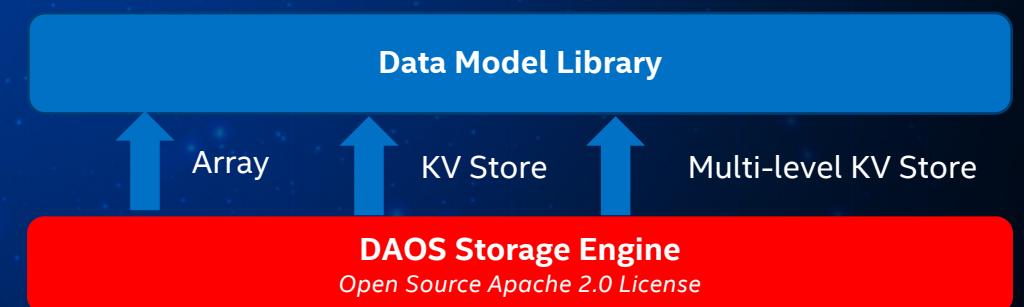


# DAOS API

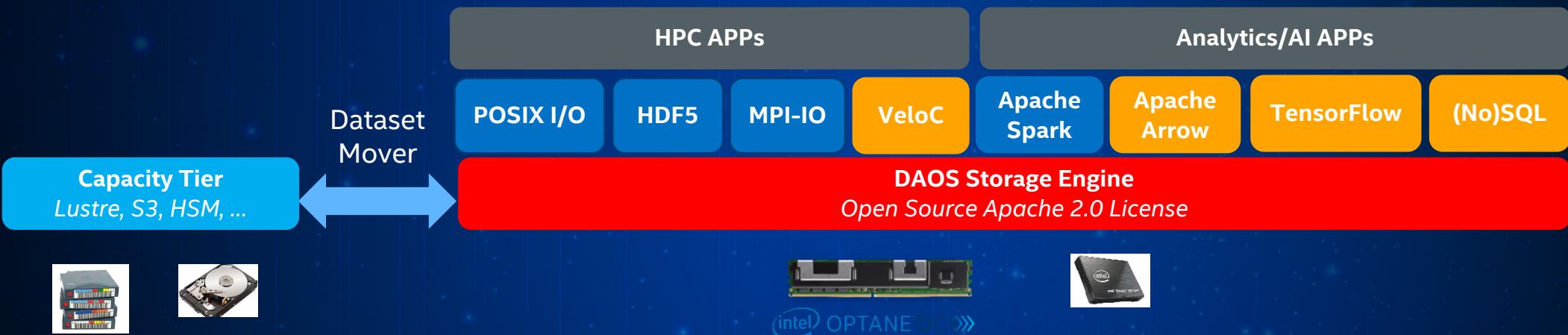


Native support for structured, semi-structured & unstructured data models

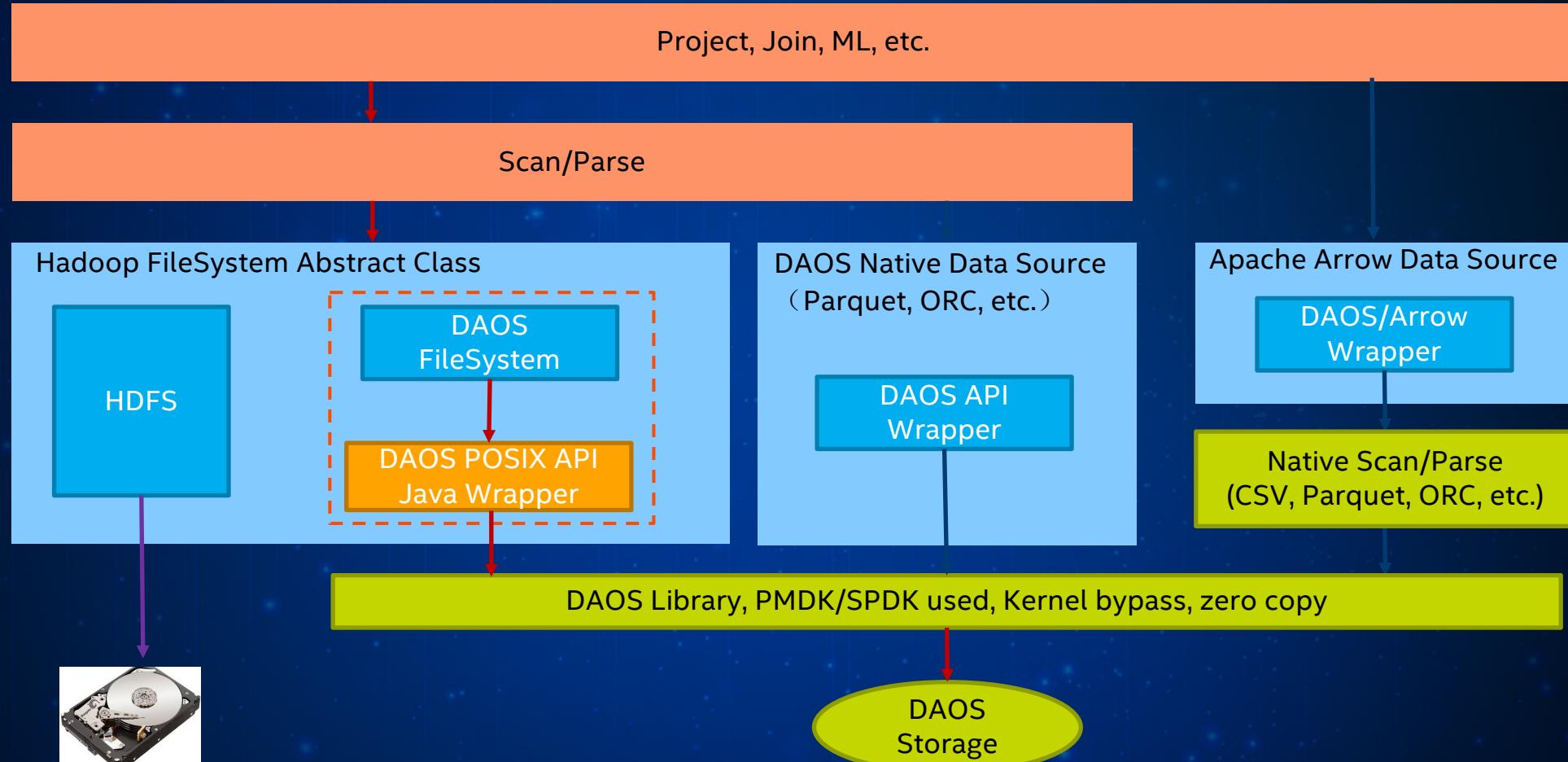
- Built on top of DCPMM
- Unconstrained by POSIX serialization
- Data access time orders of magnitude faster ( $\mu\text{s}$ )
- Scalable concurrent updates & high IOPS
- Enable in-storage computing



# APPLICATION INTERFACE



# DAOS & BIG DATA / AI



intel OPTANE DC PERSISTENT MEMORY



# DAOS: PRIMARY STORAGE FOR AURORA



## Aurora DAOS configuration

- Capacity: **230PB**
- Bandwidth **>25TB/s**

"The Argonne Leadership Computing Facility will be the first major production deployment of the DAOS storage system as part of Aurora, the first US exascale system coming in 2021. The DAOS storage system is designed to provide the levels of metadata operation rates and bandwidth required for I/O extensive workloads on an exascale-level machine."

**Susan Coghlan, ALCF-X Project Director/Exascale Computing Systems Deputy Director**



# DAOS COMMUNITY ROADMAP

