

DAOS at Exascale

DAOS Users Group 2022



U.S. DEPARTMENT OF
ENERGY

Hewlett Packard
Enterprise

intel.

Kevin Harms
Argonne Leadership Computing Facility



Aurora

Leadership Computing Facility
Exascale Supercomputer

Peak Performance
 ≥ 2 Exaflops DP

Intel GPU
**Intel® Data Center
GPU Max**

Intel Xeon Processor
**Xeon Intel® Xeon®
CPU Max**

Platform
HPE Cray-Ex

Compute Node

2 Xeon Intel® Xeon® CPU Max processors
6 Intel® Data Center GPU Max Node Unified Memory Architecture
8 fabric endpoints

GPU Architecture

Intel XeHPC architecture
High Bandwidth Memory Stacks

Node Performance

>130 TF

System Size

>9,000 nodes

Aggregate System Memory

>10 PB aggregate System Memory

System Interconnect

HPE Slingshot 11
Dragonfly topology with adaptive routing

Network Switch

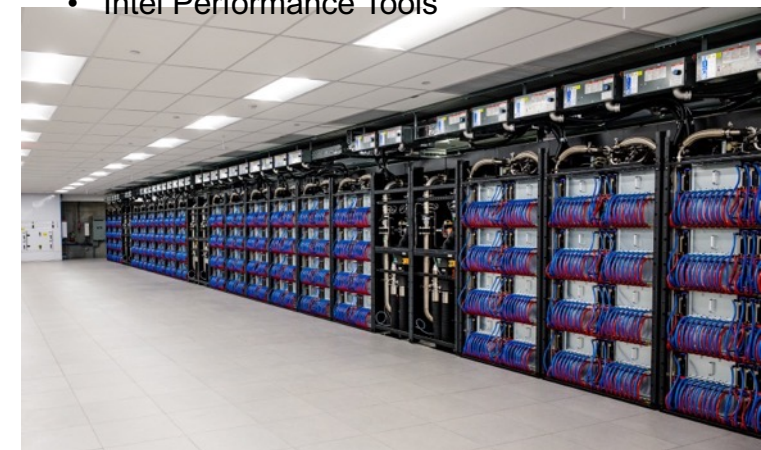
25.6 Tb/s per switch (64 200 Gb/s ports)
Links with 25 GB/s per direction

High-Performance Storage

220 PB
 ≥ 25 TB/s DAOS bandwidth

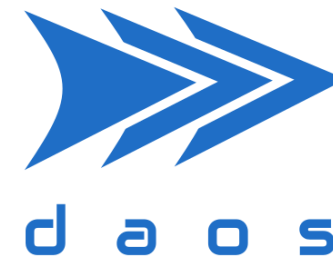
Software Environment

- C/C++
- Fortran
- SYCL/DPC++
- OpenMP offload
- Kokkos
- RAJA
- Intel Performance Tools

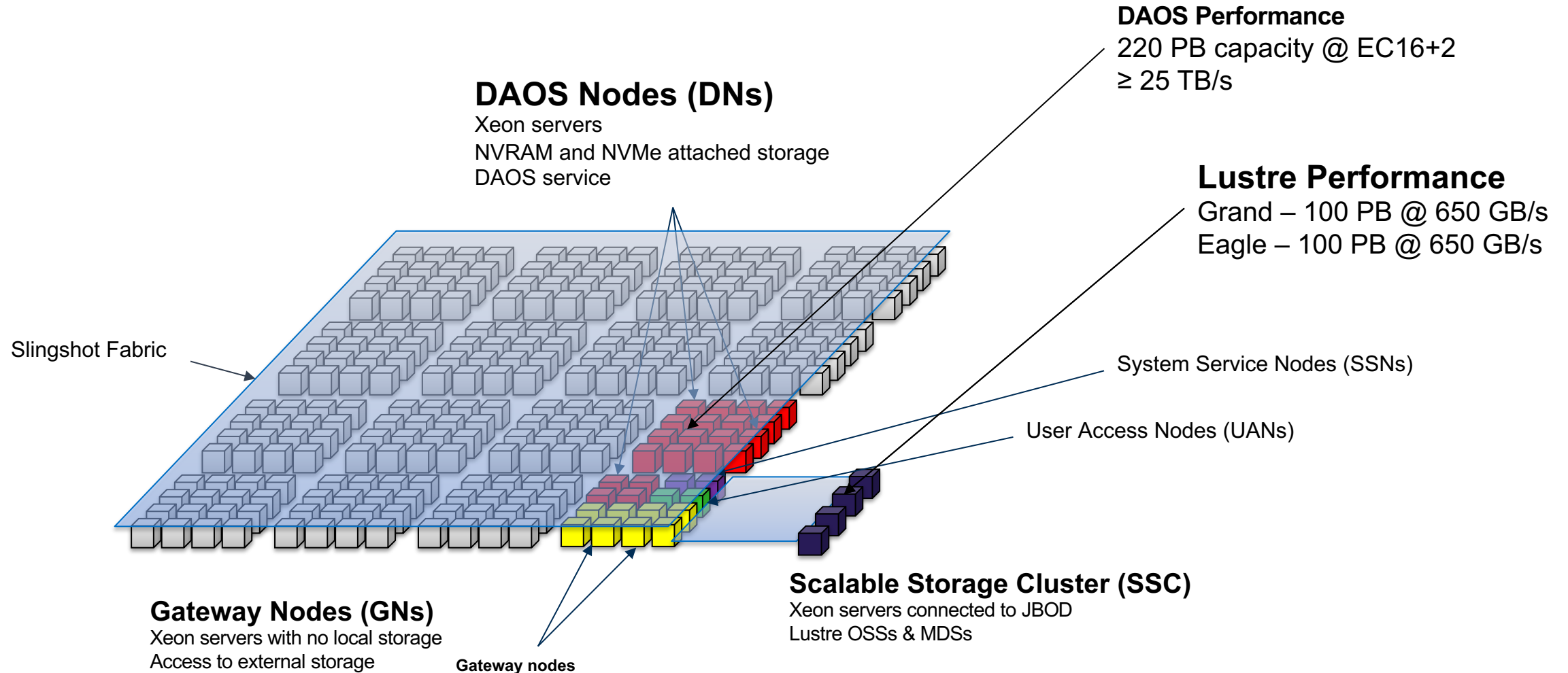


ALCF and DAOS

- Argonne Leadership Computing Facility and Intel started a collaboration on DAOS in 2015
- Collaboration on design and features related to Aurora
- Part of Non-recurring Engineering (NRE) of Aurora
 - Support for multiple simultaneous libfabric providers
 - Application optimizations for DAOS
 - Optimized object placement
 - Catastrophic Recovery



Aurora Overview

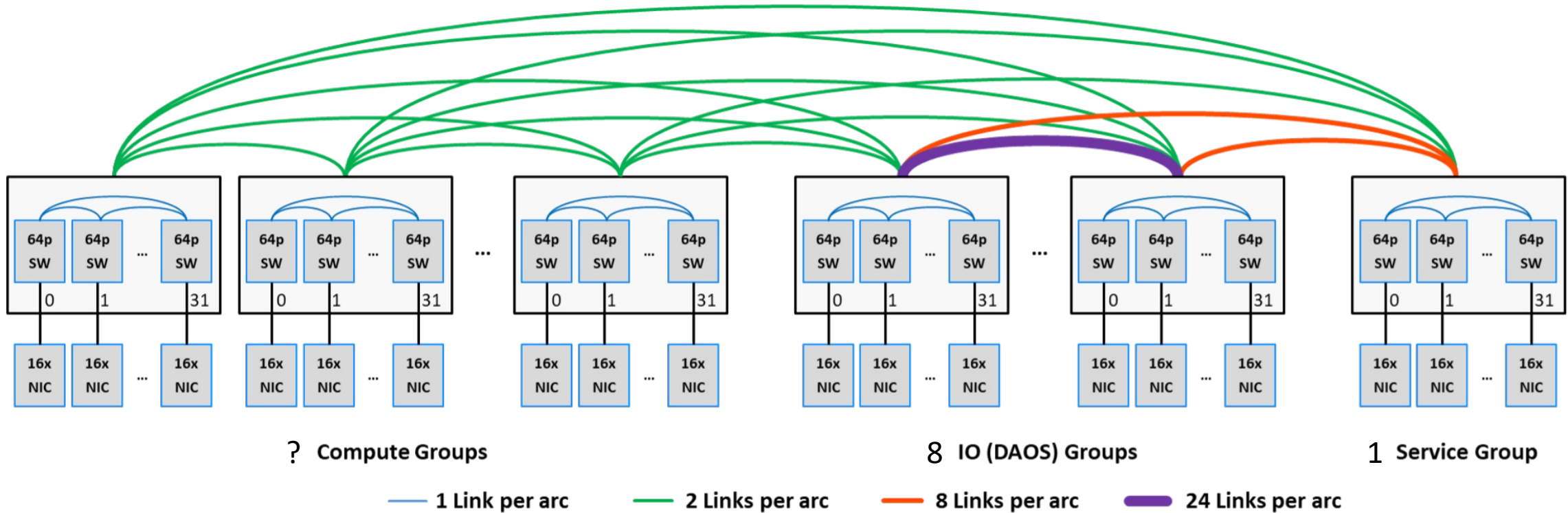


DAOS Node Details

- Intel Coyote Pass System
 - (2) Xeon 5320 CPU (Ice Lake)
 - (16) 32GB DDR4 DIMMs
 - (16) 512GB Intel Optane Persistent Memory 200
 - (16) 15.3TB Samsung PM1733
 - (2) HPE Slingshot NIC
- 1024 Total Servers
 - Each node will run 2 DAOS engines
 - 2048 DAOS engines



Aurora Network Architecture



- Increased DAOS inter-group bandwidth
 - Support rebuilding and inter-server communication
 - Prevent DAOS server traffic interfering with application communication
- Increased bandwidth to service group
 - Support off-cluster access and data-movement

Aurora DAOS Status

- *Note: All work being done by DAOS testing team*
- Initial hardware validation of DAOS nodes completed
 - Each server operating as expected
- Initial per-dragonfly group testing
 - Run automated test system scaling DAOS servers up to full dragonfly group size
 - Run soak testing on system
 - Using gateway nodes or other DAOS nodes as clients
- Scale-up testing
 - Running automated testing on multi-dragonfly group scale
 - Working through various network and DAOS issues
 - Captured in DAOS Jira

Sunspot

<https://www.alcf.anl.gov/support-center/aurora/getting-started-sunspot>

- ALCF's Test and Development system
 - Think of it as a baby Aurora
- Two compute racks / groups
 - 128 compute nodes
- DAOS deployment
 - 20 DAOS nodes
 - Identical server configuration to Aurora
 - Allows running EC16+2 – 18 nodes with 2 nodes for failover
- First? production environment for DAOS at ALCF
 - Follow pool and container usage plan for Aurora
 - 1 pool per project
 - Pool allocated to ~60-80% of targets
 - ACL limits pool to project members
 - Users create containers
 - Suggested default data protection of EC16+2 on containers
- Examine storage ratio of metadata to data

IO-500 Results

- IO-500 SC22 BoF submission
— <https://io500.org>
- **IO500: The High-Performance Storage Community**
— Tuesday, 15 November 2022 - 5:15pm - 6:45pm
— D174



<https://io500.org>

Join the Aurora Team

- Looking for a post-doc to work on DAOS
 - ALCF's performance engineering group is looking for a Postdoctoral Appointee to perform research and development on the open source DAOS storage system, in the context of the upcoming exascale platforms, and Aurora in particular.
 - Three areas of interest for study are:
 - new opportunities for applications to optimize I/O that isn't oriented around file access. DAOS provides very low latency access and the possibility allows applications to write data in a more "read-optimized" format with minimal penalty versus write-optimized formats.
 - DAOS supports a prototype "active storage" interface, and exploration of some HPC type workloads (like pointer chasing, lookup tables, etc.)
 - With the proliferation of CPUs and accelerators with significant dedicated high performance memory, the DAOS client should provide a mechanism to utilize device memory with direct-to-NIC memory movement bypassing CPU memory.
- https://argonne.wd1.myworkdayjobs.com/Argonne_Careers/job/Lemont-IL-USA/Postdoctoral-Appointee---Exascale-Storage-using-DAOS_414419

Acknowledgements

This research used resources of the Argonne Leadership Computing Facility, which is a DOE Office of Science User Facility supported under Contract DE-AC02-06CH11357.