

# DAOS Foundation

ISC 2024 Meeting



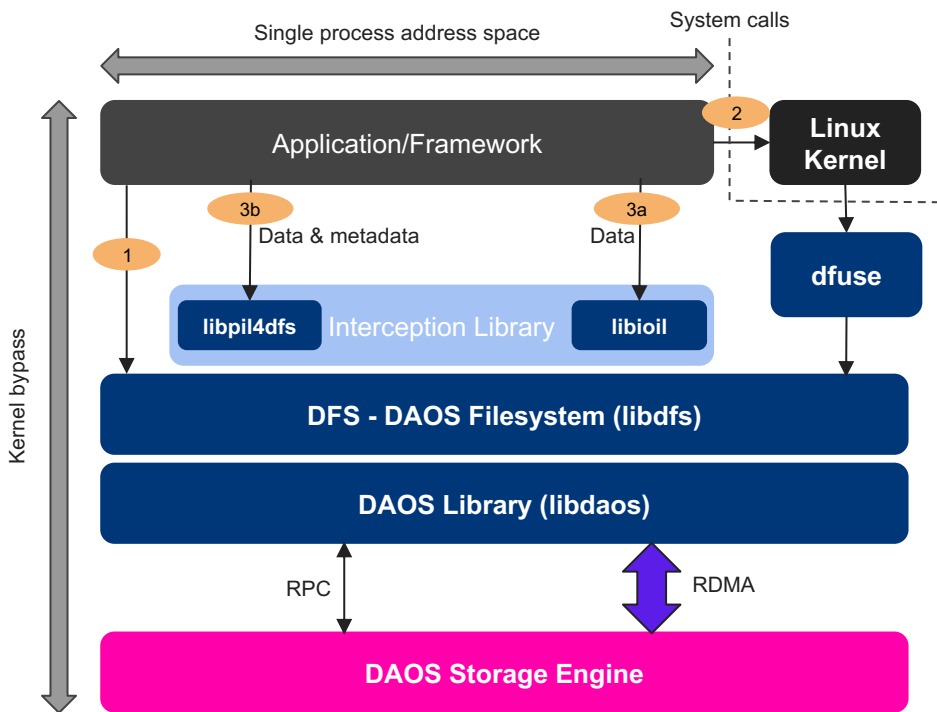
<https://foundation.daos.io>



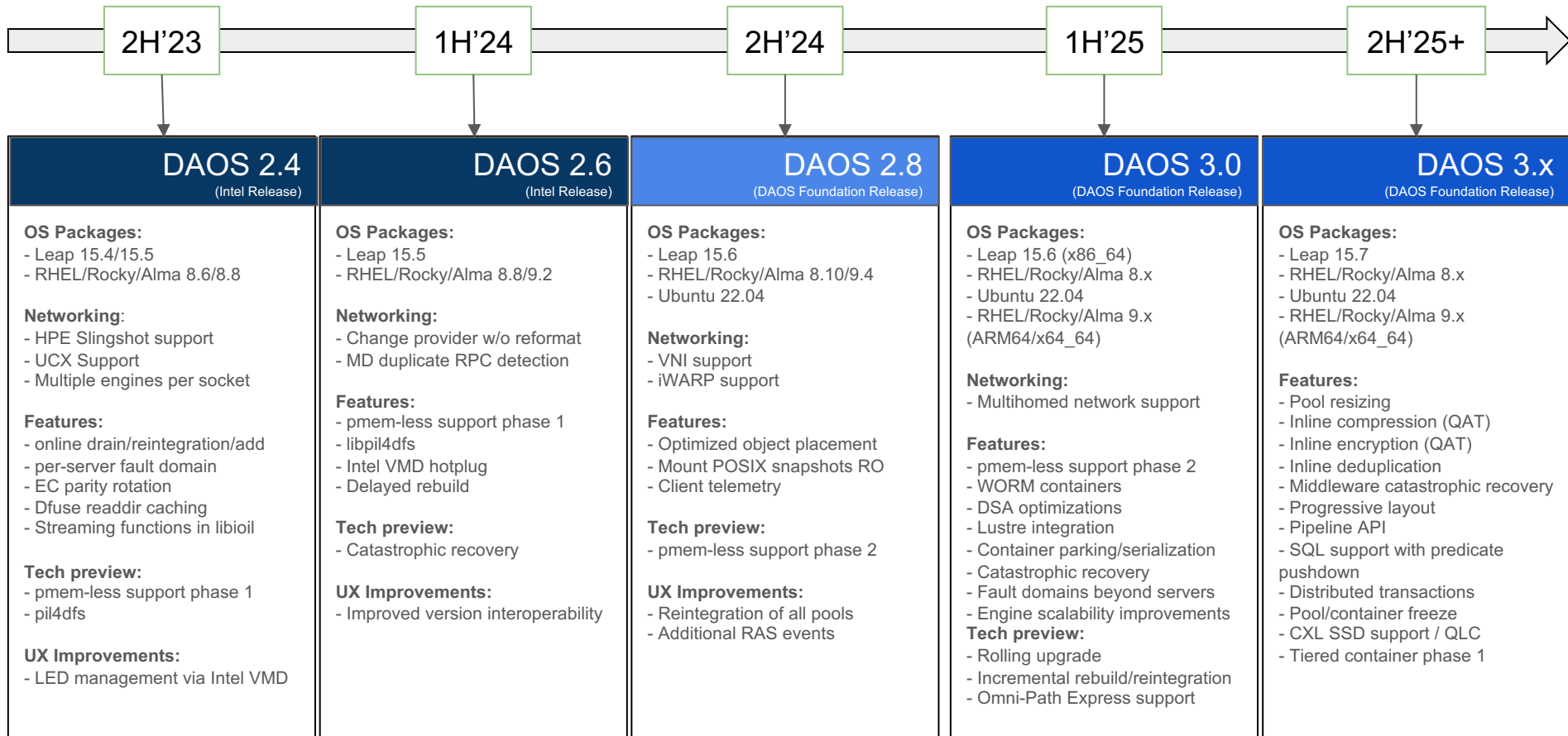
# Roadmap



# POSIX Support & Interception



1. Userspace DFS library with API like POSIX
    - **Require** application changes
    - Low latency & high concurrency
    - No caching
  2. DFUSE daemon to support POSIX API
    - **No** application changes
    - VFS mount point & high latency
    - Caching by Linux kernel
  3. DFUSE + Interception library
    - **No** application changes
    - 2 flavors using LD\_PRELOAD
- 3a libioil
- (f)read/write interception
  - Metadata via dfuse
- 3b libpil4dfs
- Data & metadata interception
  - Aim at delivering same performance as #1 w/o any application change
  - Mmap & binary execution via fuse



## 2.8 Plan

- 2.8 milestones
  - Test builds
  - Feature freeze
  - Code freeze
  - Release candidates
  - Release: shooting for Q4'24
- Generate more release candidates
  - Give opportunity for more community testing
- Train model

# Networking

- Slingshot VNI support
  - Isolate jobs running on the fabric with virtual
  - All jobs should still access the multi-tenant storage system
  - Block of VNIs allocated to DAOS at initialization time
- iWARP support
  - RDMA over tcp/ip
  - Supported via libfabric verbs provider
  - To be tested/validated

# Client Telemetry

- Infrastructure already landed for 2.6
- Leverage initial contribution from HPE
- Restructured by Intel and Google
- Client metrics accessible via multiple ways:
  - Prometheus endpoint exported by agent
  - Exported to csv file when job terminates
- Exploring alternatives:
  - Darshan
  - Dump to a “metrics” container?

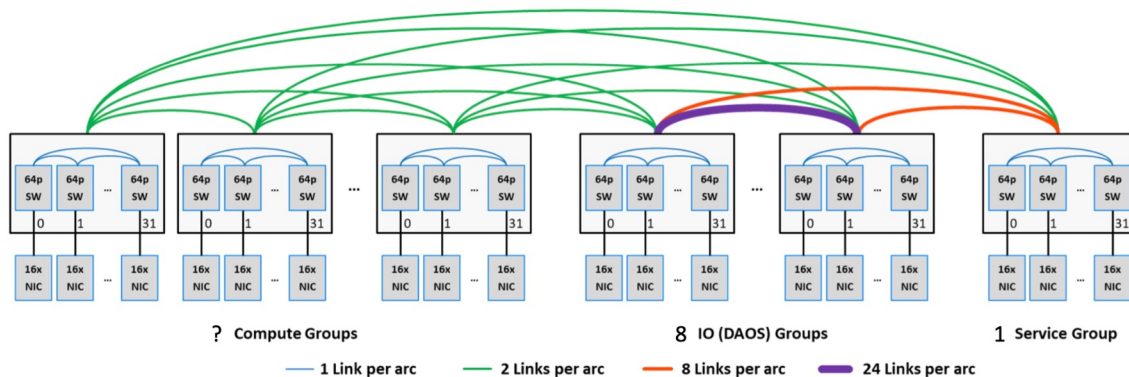
# POSIX Snapshot

- Snapshot already supported for any container type via daos utility
  - Identified by epoch number
- Add ability to access a snapshot of a POSIX container
- libdfs
  - Mount libdfs and “switch” mountpoint to a snapshot
  - Must have no active dfs objects to perform the switch
- dfuse
  - All snapshots accessible via .snapshot directory under the root
  - Accessible read-only
  - Can copy files/subtree from .snapshots/ back to main namespace
- Can delete snapshots via daos utility
  - Will automatically disappear from .snapshots directory



# Optimized Object Placement

- Placement aware of the network topology
  - E.g. DragonFly fabric
- Performance domains complementary to fault domains
- Different strategies
  - Spread redundancy groups as widely as possible across different performance domains
  - Spread redundancy groups across as fewer performance domains as possible

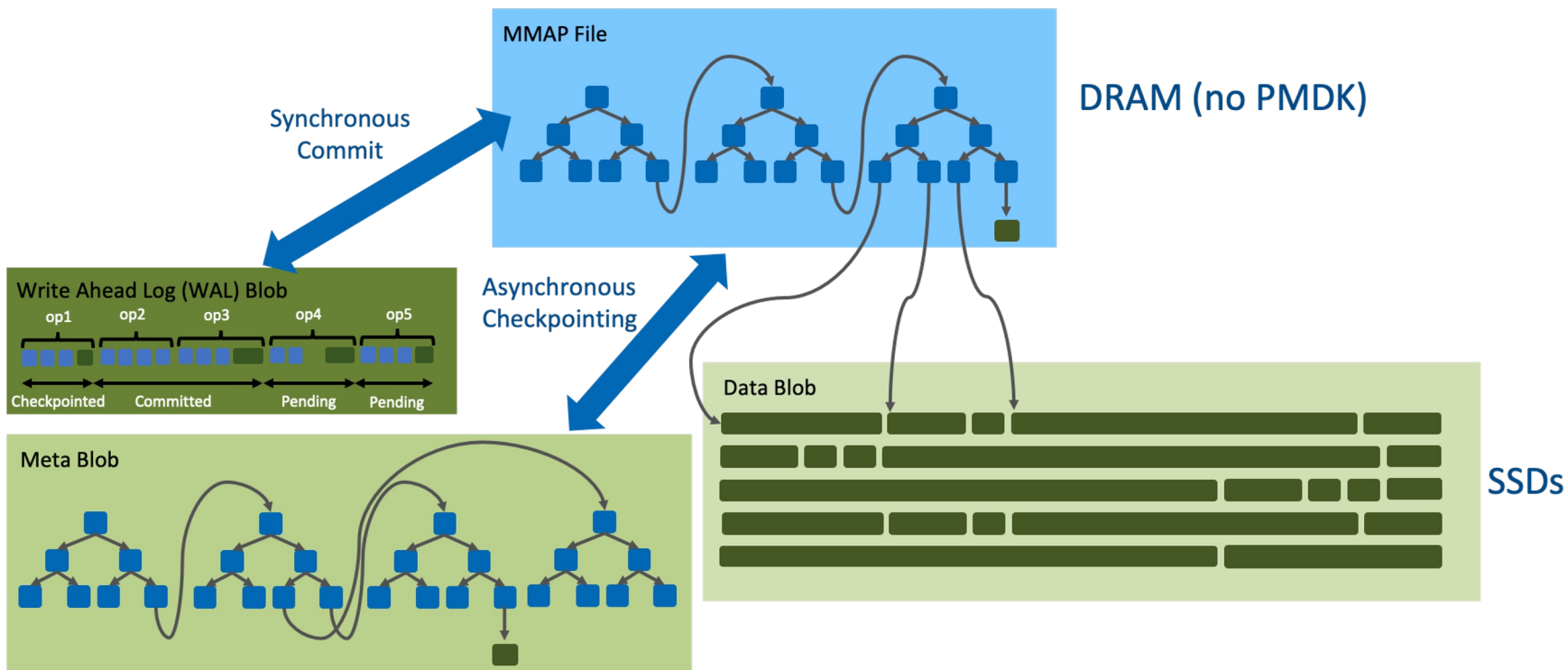


Aurora Network Topology  
(from Kevin Harms' presentation at DUG'22)

# UX Improvements

- Reintegration of all pools
  - Upon exclusion of an engine from a pool, it must be reintegrated manually
  - No automatic reintegration (bad idea, right?)
  - Must be reintegrate from each pool individually
    - `dmg pool reintegrate pool1 -rank=20; dmg pool reintegrate pool2 -rank=20; ...`
  - Provide a way to reintegrate an engine from all the pools it was excluded in a single command line
- Additional RAS events
  - Generate notifications when some operations are completed on the cluster
  - E.g. a new pool or container is created

# Pmem-less Support Phase 2 (tech preview)



# Pmem-less Support Phase 2 (tech preview)

- Phase 1 GA in 2.6
  - Meta blob size = mmap size
  - Work to reduce VOS memory footprint
  - ~600M 4K files per TB of RAM (using S1)
- Phase 2
  - Meta blob size > mmap size
  - Memory bucket allocator with dynamic eviction
    - 64MB bucket flushed to metadata blob
    - Ability to flush to data blob (3.0 scope)
  - Object flattening and eviction (3.0 scope)
    - Small objects flattened in metadata blob
    - Fetch entirely on cache miss