# DAOS Next Generation Storage

EXASCALE COMPUTING PROJECT

Kevin Harms (ALCF)
Mohamad Chaawari (Intel)
Johann Lombardi (Intel)
John Mellor-Crummey (Rice)
Yumeng Liu (Rice)

RICE UNIVERSITY

intel.

Argonne
NATIONAL LABORATORY

NNSA
National Nuclear Security Administration

U.S. DEPARTMENT OF ENERGY | Office of Science

# Agenda

- Introductions

- DAOS Overview - Johann

- DFS Overview - Mohamad

- Accelerating I/O in HPCToolkit using DAOS - Yumeng

- Discussion - Kevin

# DAOS Next Generation Storage

2022 ECP Community BOF Days

# 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 History

# 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**

# Storage Pooling

| Pool 1 | ⬛ | Project Apollo | 100PB usable | 20TB/s | 200M IOPS |
|--------|-----|----------------|--------------|--------|-----------|
| Pool 2 | 🟨 | Project Gemini | 10PB usable | 2TB/s | 20M IOPS |
| Pool 3 | 🟪 | Project Mercury | 30TB usable | 80GB/s | 2M IOPS |

**DAOS System**

*DAOS Nodes (DN)*

# Dataset Management



POSIX Namespace



File-per-process
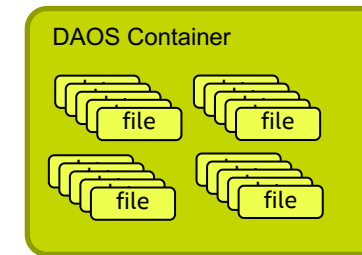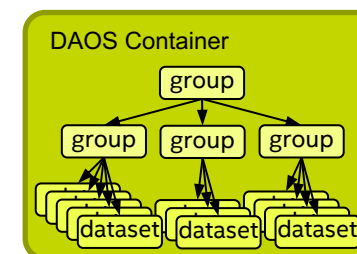
- New data model to unwind 30+y of file-based management

- Introduce notion of dataset: DAOS container

  - Basic unit of storage

  - Containers have a type

  - POSIX containers can have billions of files/directories

  - Many parameters

  - ACLs

- Advanced container query capabilities



HDF5 « File »



Key-value store



S3 Bucket



Graph

# Storage Management

# DAOS Objects



I/O Middleware View

Container (eg POSIX)
- root
  - dir
  - dir
  - dir
    - file
    - file
    - file

Mapping →

DAOS Layout View

Container
- obj1
- obj2
- obj3
- obj4
- obj5
- obj10
- obj20

object →

128-bit object identifier

Array

Multidimensional Array

Key-value Store

Multi-level Key-value Store

# Software Ecosystem

# 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 Community Roadmap



| | | | | | |
|---|---|---|---|---|---|
| **1H'21** | **2H'21** | **1H'22** | **2H'22** | **1H'23** | **2H'23** |

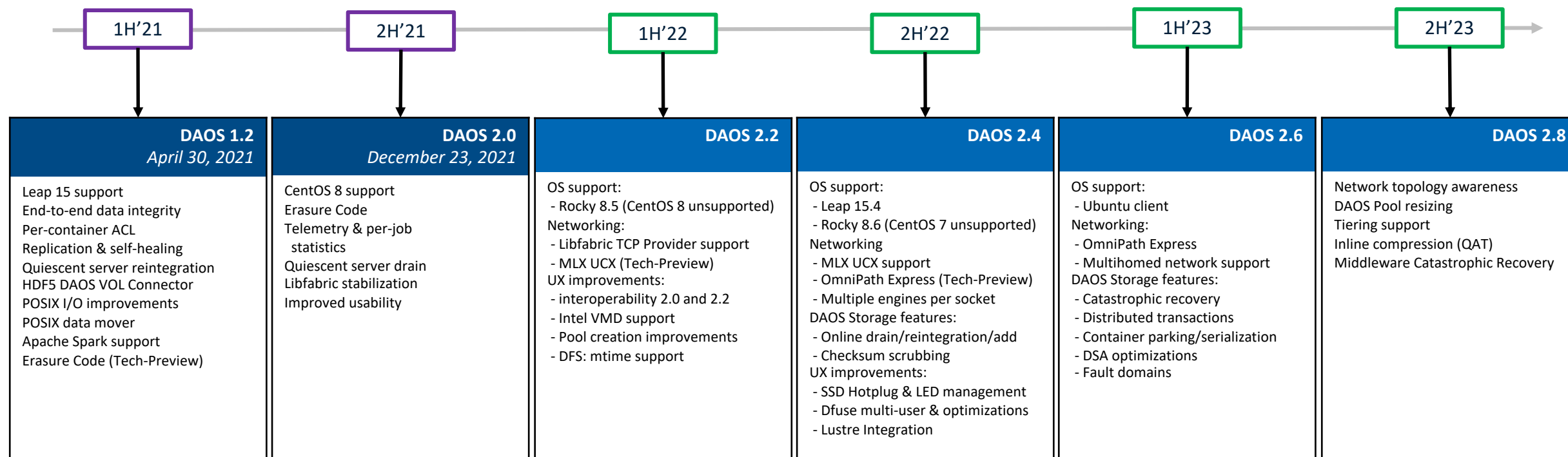| **DAOS 1.2**<br>*April 30, 2021* | **DAOS 2.0**<br>*December 23, 2021* | **DAOS 2.2** | **DAOS 2.4** | **DAOS 2.6** | **DAOS 2.8** |
|---|---|---|---|---|---|
| Leap 15 support<br>End-to-end data integrity<br>Per-container ACL<br>Replication & self-healing<br>Quiescent server reintegration<br>HDF5 DAOS VOL Connector<br>POSIX I/O improvements<br>POSIX data mover<br>Apache Spark support<br>Erasure Code (Tech-Preview) | CentOS 8 support<br>Erasure Code<br>Telemetry & per-job<br> statistics<br>Quiescent server drain<br>Libfabric stabilization<br>Improved usability | OS support:<br> - Rocky 8.5 (CentOS 8 unsupported)<br>Networking:<br> - Libfabric TCP Provider support<br> - MLX UCX (Tech-Preview)<br>UX improvements:<br> - interoperability 2.0 and 2.2<br> - Intel VMD support<br> - Pool creation improvements<br> - DFS: mtime support | OS support:<br> - Leap 15.4<br> - Rocky 8.6 (CentOS 7 unsupported)<br>Networking<br> - MLX UCX support<br> - OmniPath Express (Tech-Preview)<br> - Multiple engines per socket<br>DAOS Storage features:<br> - Online drain/reintegration/add<br> - Checksum scrubbing<br>UX improvements:<br> - SSD Hotplug & LED management<br> - Dfuse multi-user & optimizations<br> - Lustre Integration | OS support:<br> - Ubuntu client<br>Networking:<br> - OmniPath Express<br> - Multihomed network support<br>DAOS Storage features:<br> - Catastrophic recovery<br> - Distributed transactions<br> - Container parking/serialization<br> - DSA optimizations<br> - Fault domains | Network topology awareness<br>DAOS Pool resizing<br>Tiering support<br>Inline compression (QAT)<br>Middleware Catastrophic Recovery |

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

# 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
- 5th DAOS User Group (DUG'21)
  - Recordings available at http://dug.daos.io
- Upcoming BoF at ISC'22
- Intel landing page
  - https://www.intel.com/content/www/us/en/high-performance-computing/daos.html

# DAOS NEXT GENERATION STORAGE

Johann Lombardi, Mohamad Chaarawi (Intel)
Kevin Harms (Argonne National Laboratory)

May 11, 2022

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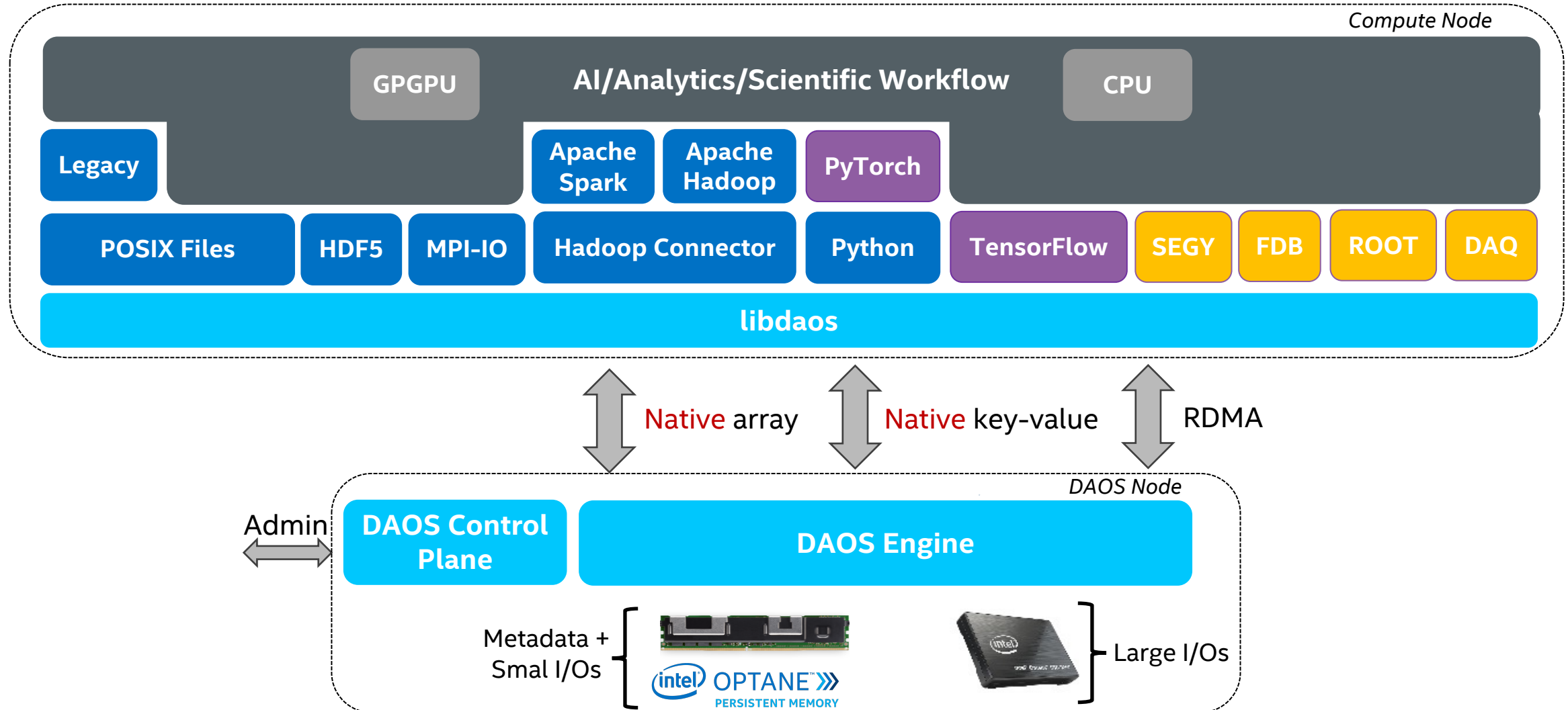
Cost reduction scenarios described are intended as examples of how a given Intel-based product, in the specified circumstances and configurations, may affect future costs and provide cost savings.  Circumstances will vary.  Intel does not guarantee any costs or cost reduction.

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.
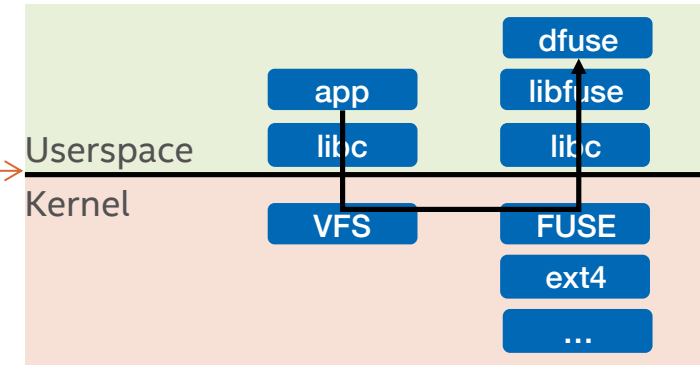
# POSIX I/O Support



- User space DFS library with an API like POSIX.
  - Requires application changes (new API)
  - Kernel Bypass, no client cache
- DFUSE plugin to support POSIX API
  - No application changes
  - Fuse Kernel Supports data (wb and ra) & metadata caching (stat, open, etc.)
- DFUSE + IL
  - No application changes, runtime LD_PRELOAD
  - Kernel Bypass for raw data IO only.

# How to use DFS?

- You should have access to a pool (identified by a string label).

- Create a POSIX container with the daos tool:

  - `daos cont create mypool --label=mycont --type=POSIX`

  - Alternatively, you can programmatically create a container to use directly in your application (if you are using DFS and changing your app).

- Open the DFS mount:

  - `dfs_connect (mypool, mycont, O_RDWR, .. &dfs);`

  - `dfs_disconnect (dfs);`

# DFS API

| POSIX | DFS |
|---|---|
| mkdir(), rmdir() | dfs_mkdir(), dfs_rmdir() |
| open(), close(), access() | dfs_open(), dfs_release(),dfs_lookup() |
| pwritev(), preadv() | dfs_read/write() |
| {set,get,list,remove}xattr() | dfs_{set,get,list,remove}xattr |
| stat(), fstat() | dfs_stat(),ostat() |
| readdir() | dfs_readdir() |
| ... | ... |

- Mostly 1-1 mapping from POSIX API to DFS API.
- Instead of File & Directory descriptors, use DFS objects.
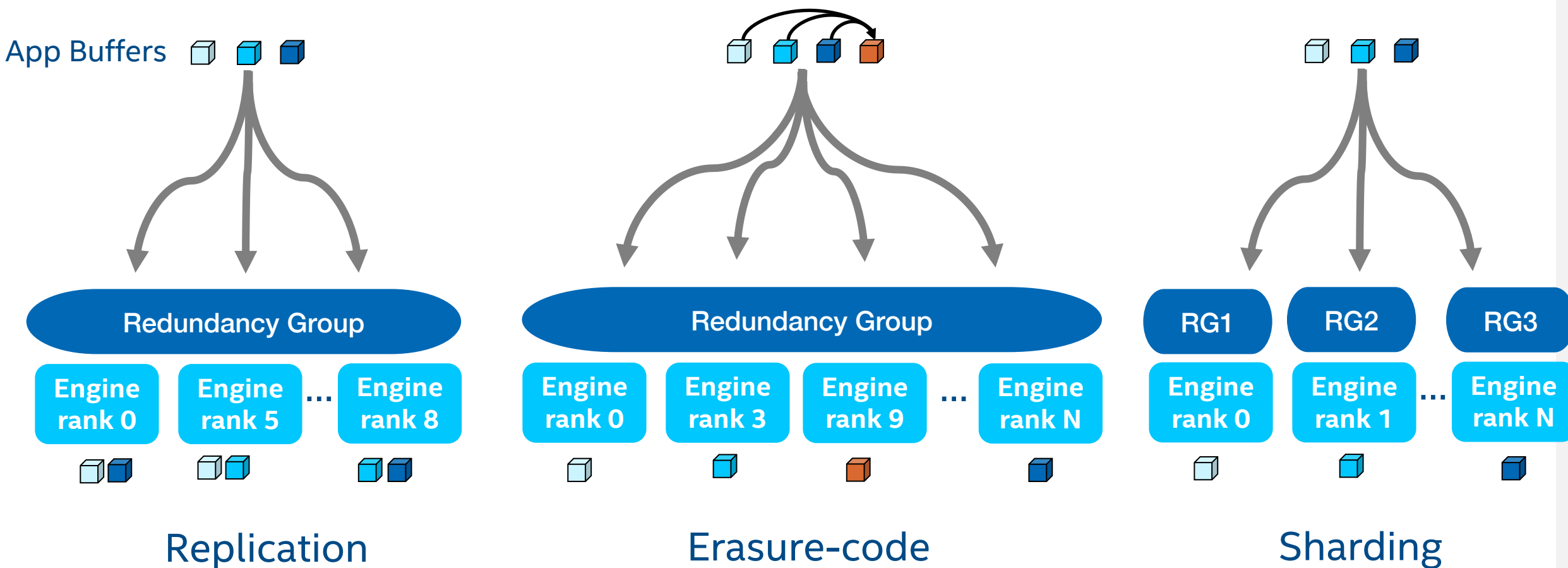- All calls need the DFS mount which is usually done once initialization.

# DFUSE

- To mount an existing POSIX container with dfuse, run the following command:

    - `dfuse --pool mypool --container mycont -m /mnt/dfuse`

    - No one can access your container / mountpoint unless access is provided on the pool and container (through ACLs).

    - Multi-user concurrent support is planned for DAOS 2.4

- Now you can access files / directories as any namespace in the container, and applications can run without any modifications (the easy path).

- Interception Library:

    - This library works in conjunction with dfuse and allow to interception of POSIX I/O calls and issue the I/O operations directly from the application context through libdaos without any application changes.

    - This provides kernel-bypass for I/O. To use this set the LD_PRELOAD to point to the shared library in the DOAS install dir

        - LD_PRELOAD=/path/to/daos/install/lib/libioil.so

# To Intercept or Not Intercept

- The dfuse interception library bypasses the fuse kernel and sends IO to DFS directly

  - No caching (no writeback / readahead)

  - Works well for bulk IO.

  - Does not work well for very tiny IO (order of bytes) where dfuse without interception takes advantage of the kernel write back and read ahead caching to reduce IO over the wire.

intel. 8

# DAOS Data Model: Distribution & Fault Tolerance



App Buffers

Redundancy Group

Engine rank 0    Engine rank 5    ...    Engine rank 8

Replication

Redundancy Group

Engine rank 0    Engine rank 3    Engine rank 9    ...    Engine rank N

Erasure-code

RG1    RG2    RG3

Engine rank 0    Engine rank 1    ...    Engine rank N

Sharding

# DAOS Object Class

- Each DAOS object is identified by an object ID which encode the object class.

- The object class controls the redundancy and sharding of the object and allows the client to do algorithmic placement when accessing the object.

  - Redundancy:

    - None: S

    - N-way Replication: RP_n

    - Erasure Code: EC_dPp (d = data, p = parity)

  - Sharding:

    - Single target/group:

      - S1, RP_nG1, EC_dPpG1

    - Widely / max sharding:

      - SX, RP_nGX, EC_dPpGX

intel.

# Container Level Redundancy Setting

- On Container create, one can set the redundancy factor property (`--properties=rf:n`):
  - 0 (default): no redundancy
    - One can create objects with higher redundancy after
  - 1: tolerate 1 failure domain
    - At least RP2 or ECxP1
  - Up to 4
- DAOS automatically selects an object class for files and directories using the redundancy factor.
  - Widely striped, EC oclass for files (SX, EC_xP1GX, EC_xP2GX, etc.)
  - Single stripe, RP oclass for directories (S1, RP_2GX, RP_3GX, etc.)

# Setting the Object Class

- Users can select the object class if the default is not appropriate.
  - Default works well if files are large, directories are small (not many entries)
- DFS API allows users to pass an object class when creating an object
  - Passing 0 would use the default
- DAOS tool allows also to create a **new** file with a specific oclass.
  - `daos fs set-attr --path=/mnt/dfuse/testfile --oclass=RP_2G1`
- DAOS tool allows changing the object class on a directory to use that oclass for all entries created in that directory (does not change the oclass of the directory itself)
  - `mkdir /mnt/dfuse/dir1`
  - `daos fs set-attr --path=/mnt/dfuse/dir1 --oclass=RP_2G1`

# Striping Considerations

- Widely striped objects
  - Best IO BW (read/write) for single shared file, Highest IOPS for entry creation in single shared directory
  - Slow file stat (file size), ls (directory listing)
    - Must query all targets
    - Caching with dfuse helps in this case – limited to single node or read-only
- Small stripe objects
  - Limited size and bandwidth to number of targets
  - Good for FPP or DPP
  - Fast stat, ls (not querying all targets on storage system)

# Redundancy Considerations

- Erasure Code performance characterization for IO:
  - Best for large IO access patterns
  - Full stripe write: 33% performance hit
  - Partial stripe write: 66% performance hit
  - Read should be the same.
- Replication:
  - Best for metadata objects (directories)
  - Small files (<= 16k)
  - Write IOPS: n x slower
  - Read IOPS: equal or better – more shards to serve concurrent requests

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

intel.    14

# Challenges / Future Work

- Selecting correct object class still not straightforward:
  - Auto object class selection works fine but is not ideal in some situations.
- Extend (POSIX) container create options to allow users to provide hints as to what their expected dataset looks like:
  - small/medium/large directories
  - small/medium/large files
  - Based on those hints, the auto object class selection can be tuned to select better object class for performance considerations in such scenarios
- Extend DAOS FS tool to accept hints too for setting object class on parent directories or new files.
- *Likely DAOS 2.4 features*

# HPCToolkit Funding Acknowledgments

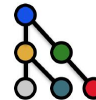# Outline

- HPCToolkit Project Overview
- Why DAOS in HPCToolkit?
- An I/O abstraction layer for HPCToolkit
- Integration experience
    - DAOS usage
    - Choice of parameters
    - Challenges
    - Differences from POSIX
- Early results
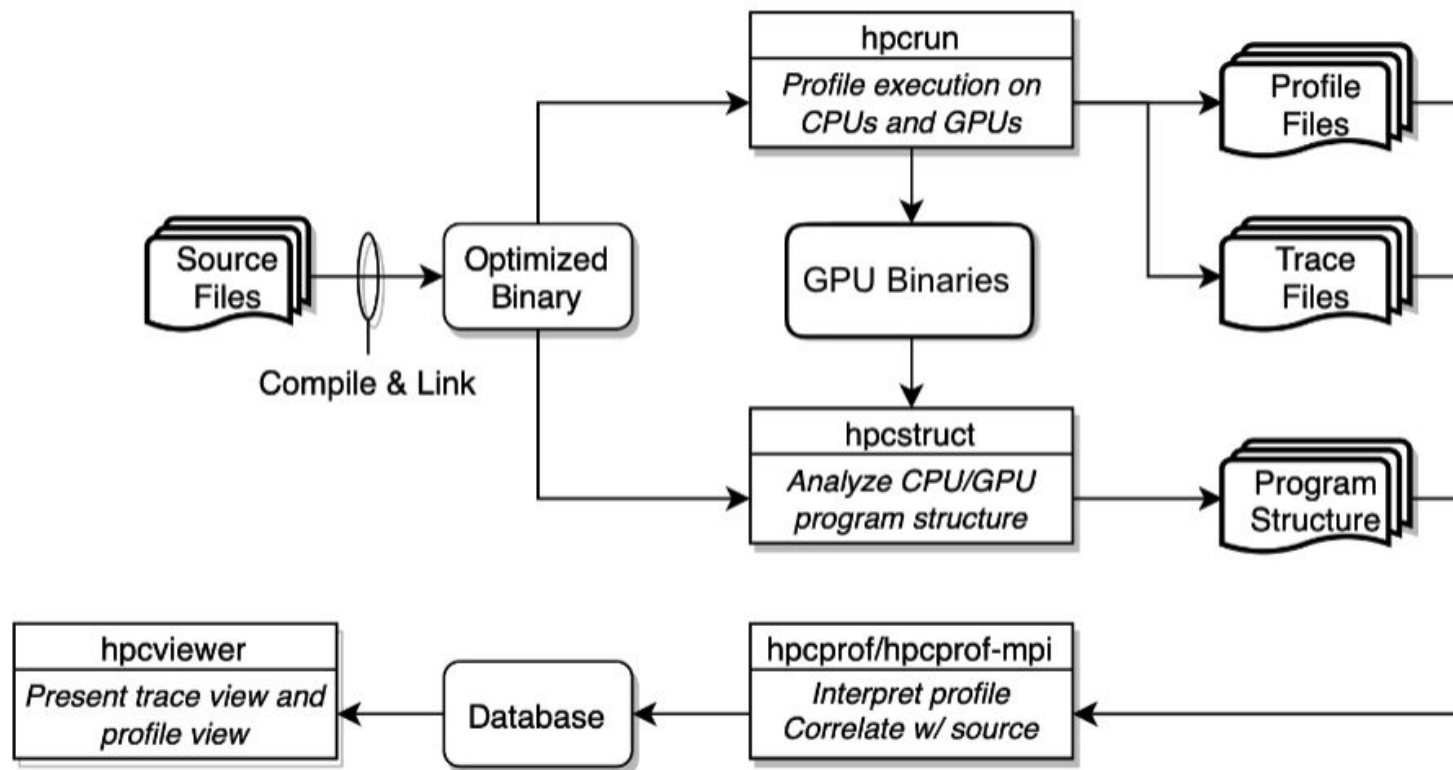
# HPCToolkit Project Overview

- **Goals**
  - Develop effective tools to measure, attribute, analyze, and understand performance of applications, libraries, tools, and system software on extreme-scale parallel systems
  - Develop new strategies and mechanisms to measure and analyze performance and resource utilization of GPU-accelerated compute nodes
  - Influence the evolution of hardware and software ecosystems to enable better tools
- **Team**
  - Lead Institution: Rice University (HPCToolkit performance tools)
    - PI: Prof. John Mellor-Crummey
    - Research staff: Dr. Laksono Adhianto, Dr. Mark Krentel, Dr. Xiaozhu Meng, Dr. Scott Warren
    - Grad students: Keren Zhou, Jonathon Anderson, Yumeng Liu, Aaron Cherian, Dejan Grubisic
  - Subcontractor: University of Wisconsin – Madison (Dyninst binary analysis toolkit)
    - Lead: Prof. Barton Miller

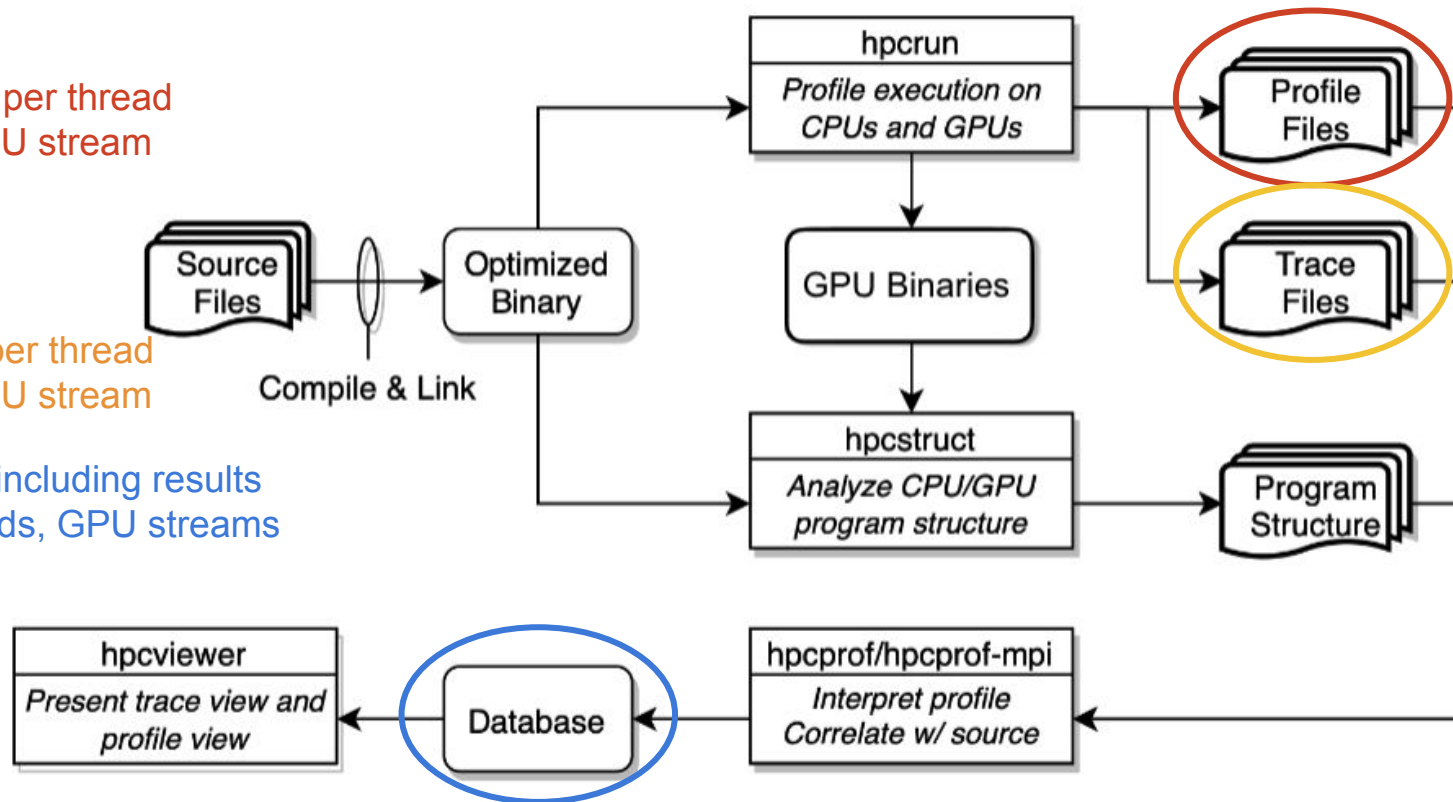# HPCToolkit Project Overview

# HPCToolkit Project Overview



One profile per thread and per GPU stream

One trace per thread and per GPU stream

Three files including results for all threads, GPU streams and traces

# Why DAOS in HPCToolkit? - I

**Concern**: **Opening millions of small measurement files in parallel**

- Without DAOS: metadata server can be a bottleneck
    - Solution:
        - merge profiles and traces for one rank or even one node in one file during writing
        - separate the profiles and traces on the fly in post-processing


- With DAOS: no default metadata besides 128bits
    - We can keep our original code: one profile and one trace per thread and per GPU stream

# Why DAOS in HPCToolkit? - II

**Concern**: **Parallel assembly of a file with irregular structure**

- All processes and threads know their offsets for writing

- <u>Without DAOS</u>:
  - Overhead for maintaining page cache consistency between multiple processes on different nodes

- <u>With DAOS</u>:
  - No cache, no extra overhead
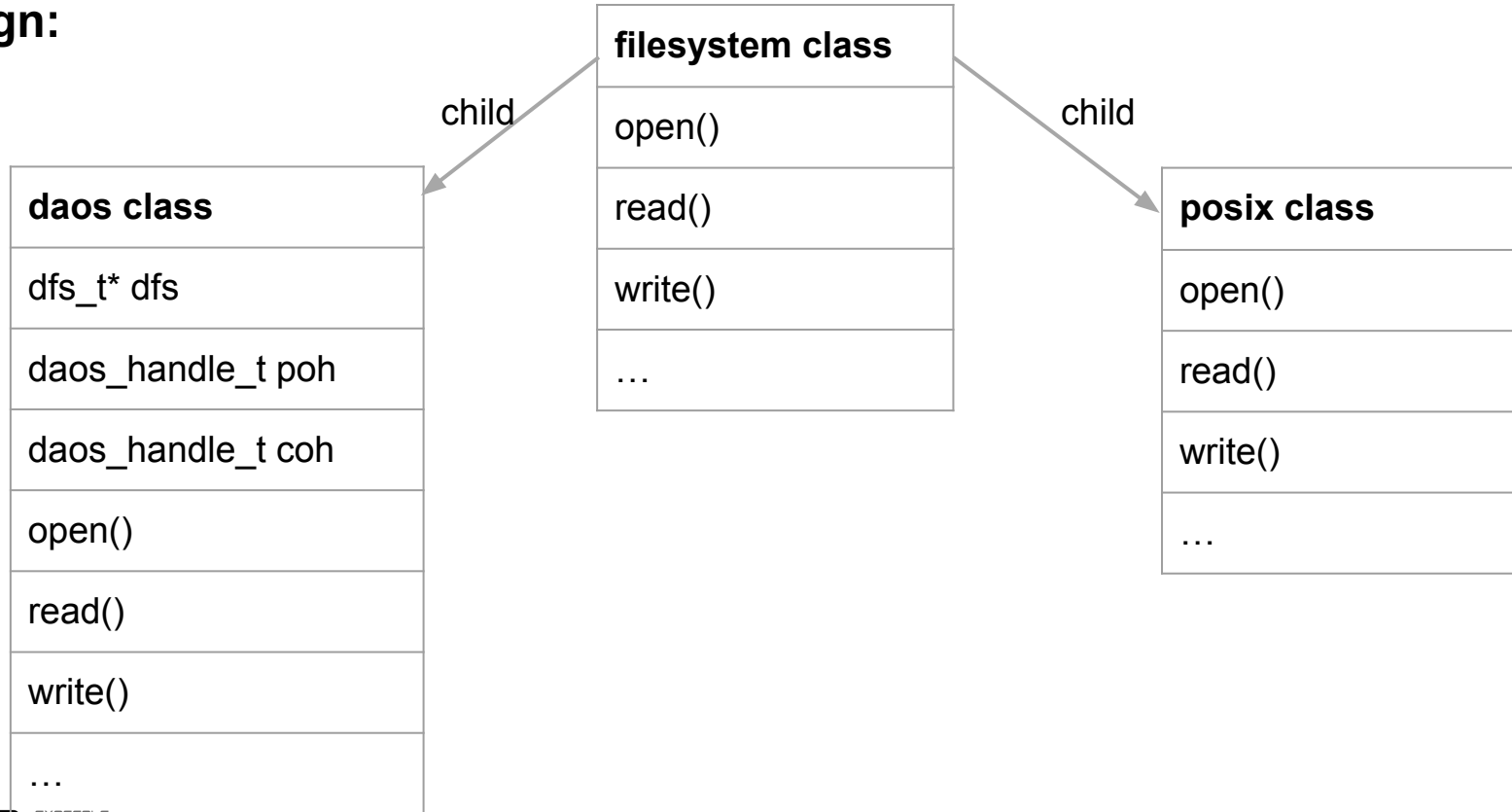
# An I/O abstraction layer for HPCToolkit

**Goals:**

- Freedom to choose any I/O option (DAOS, POSIX, Lustre …) for any file
  - Example:
    - `hpcprof -o daos://<POOL>/<CONT>/database-dir measurement-dir`
    - We use DAOS for database and use POSIX for measurement files at the same time
- Easy integration of other I/O options in the future
  - Example: HPE Rabbit Near Node Storage (Livermore's El Capitan), Lustre ...
- Details of I/O are invisible to other parts of the software
  - Example: they just call `io->write(...)` without concern for underlying I/O implementations

# An I/O abstraction layer for HPCToolkit
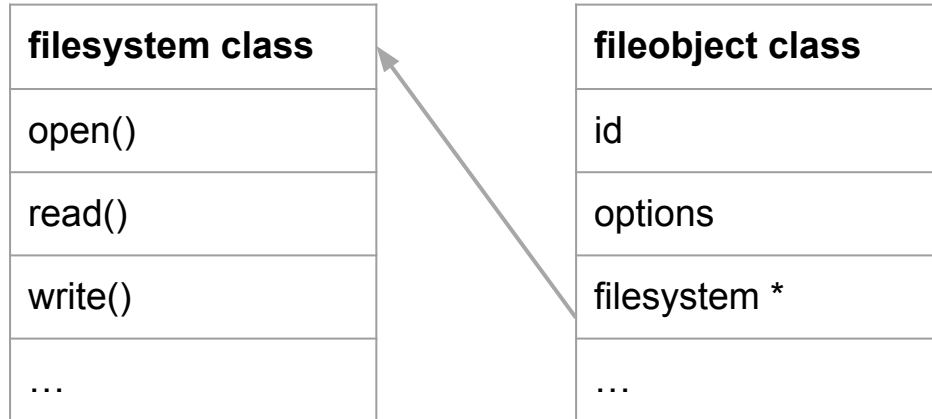
**Design:**



**filesystem class**
- open()
- read()
- write()
- …

child →

**daos class**
- dfs_t* dfs
- daos_handle_t poh
- daos_handle_t coh
- open()
- read()
- write()
- …

child →

**posix class**
- open()
- read()
- write()
- …

# An I/O abstraction layer for HPCToolkit

**Design:**

| filesystem class |
| --- |
| open() |
| read() |
| write() |
| … |

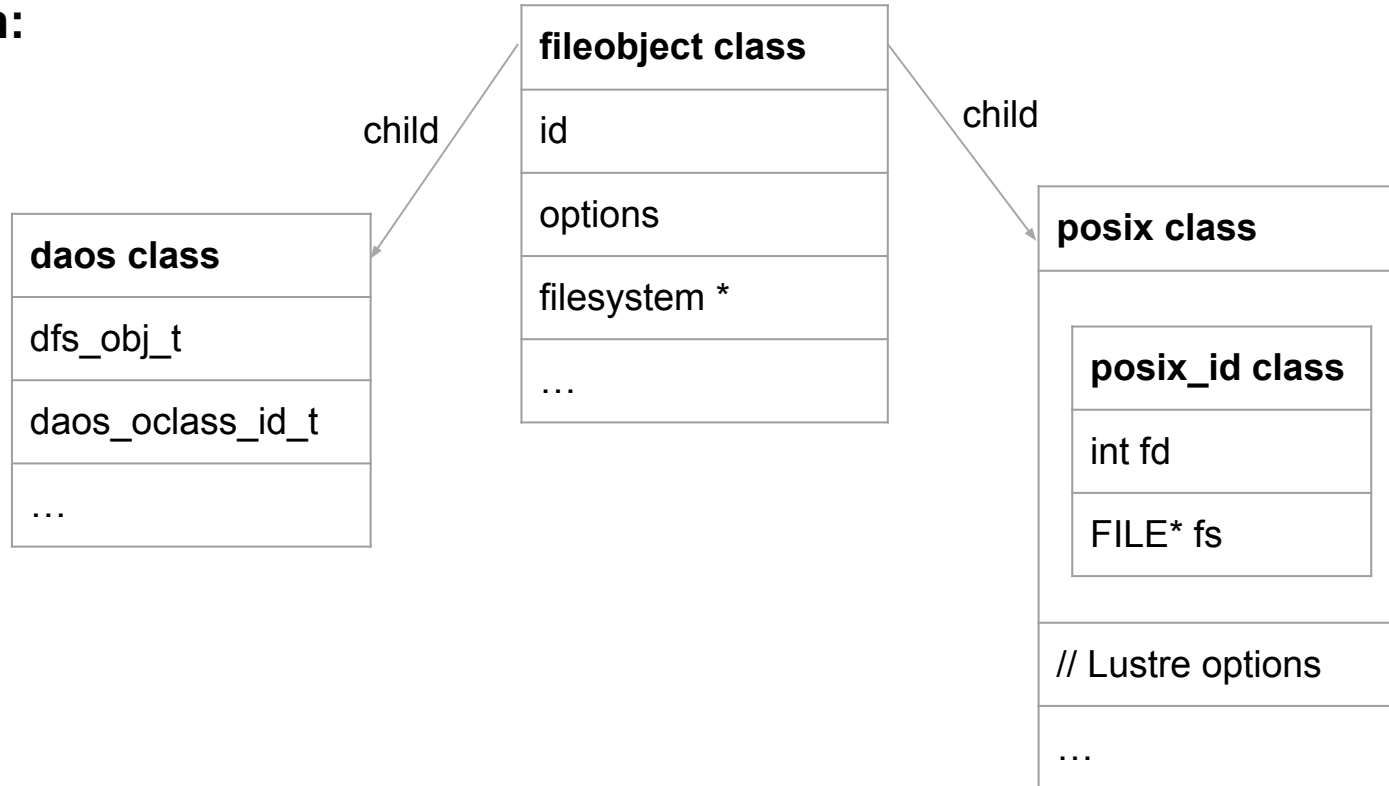| fileobject class |
| --- |
| id |
| options |
| filesystem * |
| … |

Pointer to the filesystem containing this object

# An I/O abstraction layer for HPCToolkit

**Design:**

# An I/O abstraction layer for HPCToolkit

**Simplified example of calling open():**

```
hpctio_obj_t * hpctio_obj_open(const char *path, int flags, mode_t md, int
wmode, int sizetype, hpctio_sys_t * sys)
{
    hpctio_obj_t * obj = (hpctio_obj_t *) malloc(sizeof(hpctio_obj_t));

    obj->sys_ptr = sys;
    obj->opt_ptr = sys->create_obj_options(wmode, sizetype);
    obj->id = sys->open(path, flags, md, opt, sys->params);

    return obj;
}
```

# An I/O abstraction layer for HPCToolkit

**Revisit goals:**

✓ Freedom to choose any I/O option (DAOS, POSIX, Lustre …) for any file
  - By defining different filesystem objects

✓ Easy integration of other I/O options in the future
  - Just define another children pair of filesystem class and fileobject class

✓ Details of I/O are invisible to other parts of the software
  - After defining filesystem, pass it in when opening the file object, no details exposed

# Integration experience - DAOS usage

|  | bypass OS kernel for read/write | bypass OS kernel for metadata operations | Use more than one pool and one container at a time |
|---|---|---|---|
| DFuse |  |  |  |
| DFuse + IL | YES |  |  |
| DFS | YES | YES | YES |

We expect better performance with DFS API directly

# Integration experience - choice of parameters

- Object Class
    - Profile files and Trace files: **OC_S1**
        - no striping: several MB, relatively small
        - no redundancy: not concerned about file loss, ok to lose a few profiles out of many
    - Analysis files: **OC_EC_2P1GX (Aurora: OC_EC_16P2GX)**
        - maximum striping:
            - MB - GB, relatively large
            - written in parallel by multiple processes and threads
        - erasure coding: don't want to lose analysis results

- Chunk size
    - default 1MB

# Integration experience - challenges

No matching DAOS function for every POSIX function

- `fread, fwrite, fseek` …
  - `fread`: We reworked our code to avoid it by using `read_at`
  - `fwrite`: DAOS only has `write_at`, manage our own buffer and cursor
  - `fseek`: We reworked our code to avoid it by using `read_at` and `write_at`
- `write`
  - Multiple processes append to the same log file (debug info for developers)
  - With `write_at`, we need to share a buffer or a cursor between processes
  - We reworked so that each process creates its own log file
- C++ filesystem abstractions (`directory_iterator`, `exists`, `remove_all` …)
  - We could use them if we chose DFuse + IL
  - With DFS API, we needed to reimplement these abstractions using our I/O abstraction layer

# Integration experience - differences from POSIX

- Extra initialization and finalization steps
  - Each process initializes its own DAOS pool and container handles
  - We set up the output directory and begin recording trace data BEFORE MPI initialization
- Need to manage file system accesses carefully
  - Before: use POSIX anywhere without any direct coordination
  - Now: pass around I/O abstractions to access the correct file system
    - e.g., `hpctio_sys_mkdir(const char *path, mode_t md, **hpctio_sys_t* sys**)`
- Caching semantics
  - Objects created were only visible through command line after cache refreshed
  - Using `--disable-caching` can help, but may hurt performance

# DAOS Testbed: Presque on JLSE

## 8 Clients

- Intel Xeon Gold 6140 (Skylake) CPU @ 2.30Ghz with 18 cores
- 192 GB DDR4
- Mellanox ConnectX-5

## 4 Servers

- Dual Intel Xeon Gold 6248 (Cascade Lake) CPU @ 2.5 GHz with 22 cores
- 384GB DDR4-2666 DRAM
- 1.5TB Intel Optane
- 12TB Intel P4510 NVMe SSDs
- 2 Mellanox ConnectX-5

# Collecting data for DAOS experiments with `hpcprof`

Measurement files from AMG2006 on Theta

- 1 CPU metric
- 8 nodes, 512 ranks, 2048 threads
- profile files: ~948MB
    - 512 files: ~1.5MB each
    - 1544 files: ~120KB each
- trace files: ~22GB
    - 2056 files: ~11MB each

# Early Results

- Run hpcprof on 1 node of Presque with 1 MPI rank, 36 threads
  - POSIX: 4m56s
  - DAOS: 1m50s

- Using DAOS I/O is 2.69x faster than POSIX I/O
- With POSIX, data may still be in the disk cache, may not be in stable storage
- We expect to have larger benefits from using DAOS on Aurora

# Discussion

# Acknowledgements