# Principles of HPC I/O

**Phil Carns**

**carns@mcs.anl.gov**

Mathematics and Computer Science Division

Argonne National Laboratory
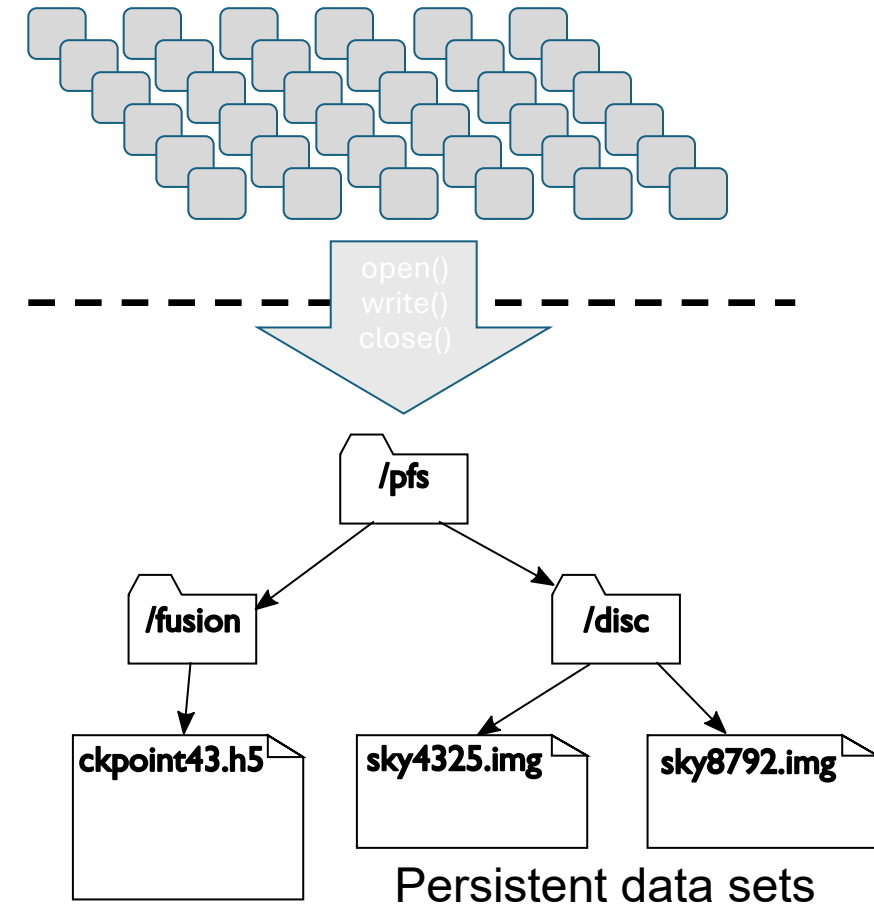
August 7, 2025

# What is HPC I/O?

HPC I/O: storing and retrieving persistent scientific data on a high-performance computing platform

– Data is usually stored on a **parallel file system** that has been optimized to rapidly store and access enormous volumes of data.

– *This is an important job! Valuable CPU time is wasted if applications spend too long waiting for data.*

– It also means that parallel file systems are quite specialized and have some unusual properties.

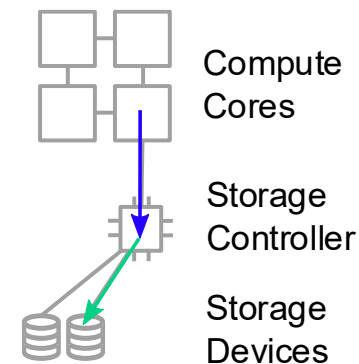Today's lectures are really all about the proper care and feeding of exotic parallel file systems.

Scientific application processes

open()
write()
close()
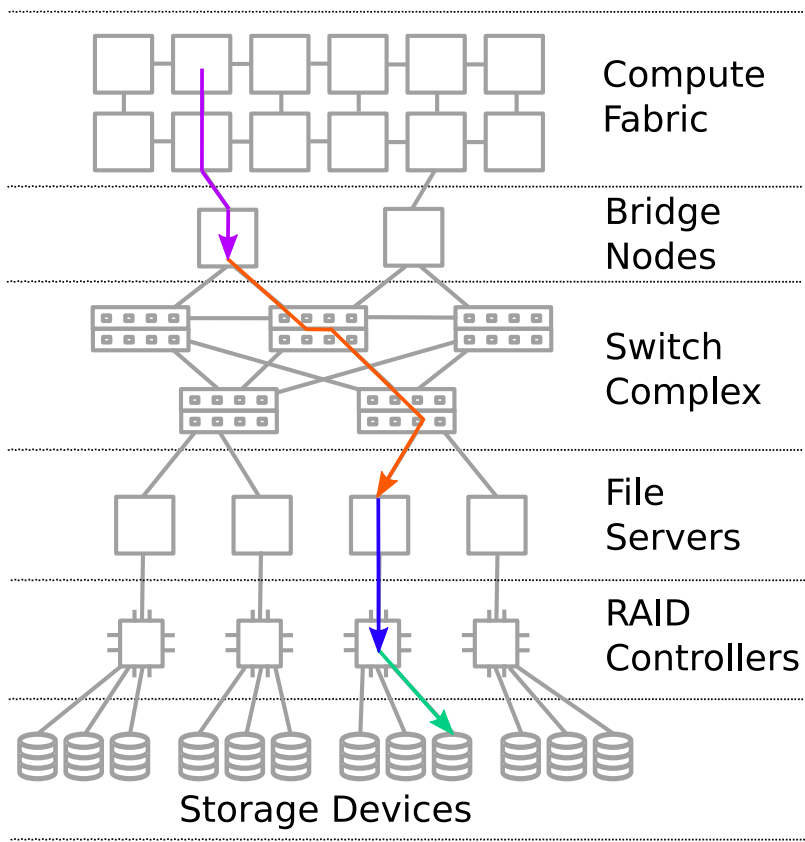
/pfs

/fusion

/disc

ckpoint43.h5

sky4325.img

sky8792.img

Persistent data sets

Hands on exercises:
https://github.com/radix-io/hands-on

extremecomputingtraining.anl.gov

ATPESC2025

Argonne
NATIONAL LABORATORY

# A look under the hood

**Workstation (laptop) storage path**

Compute Cores

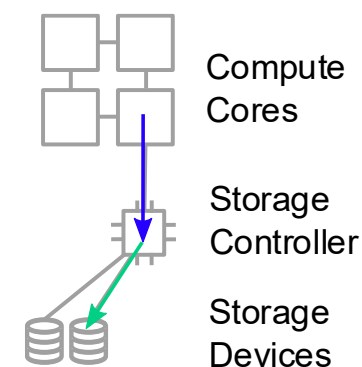Storage Controller

Storage Devices

- A typical workstation/laptop has only one storage device.
- The path between applications and storage is *short*.
- Properties:
  - Low latency
  - Low bandwidth

Hands on exercises:
https://github.com/radix-io/hands-on

ATPESC2025

Argonne
NATIONAL LABORATORY

# A look under the hood

**HPC system storage path**

Compute
Fabric

Bridge
Nodes

Switch
Complex

File
Servers

RAID
Controllers

Storage Devices
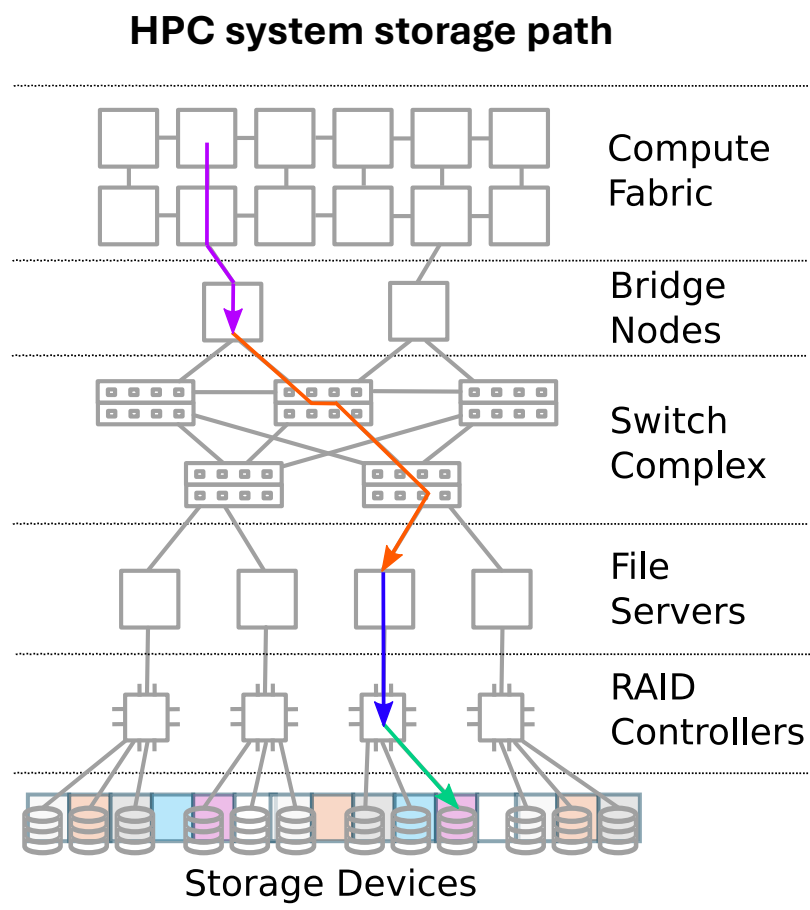
**Workstation (laptop) storage path**

Compute
Cores

Storage
Controller

Storage
Devices

- In contrast, an HPC storage system manages many (e.g., **thousands** of) disaggregated devices.
- Paths between applications and storage devices are quite long, but numerous.
- Properties:
  - High latency
  - High bandwidth

Hands on exercises:
https://github.com/radix-io/hands-on

extremecomputingtraining.anl.gov

ATPESC2025

ARGONNE
NATIONAL LABORATORY

# Striping / Layout

**HPC system storage path**



Compute Fabric

Bridge Nodes

Switch Complex

File Servers
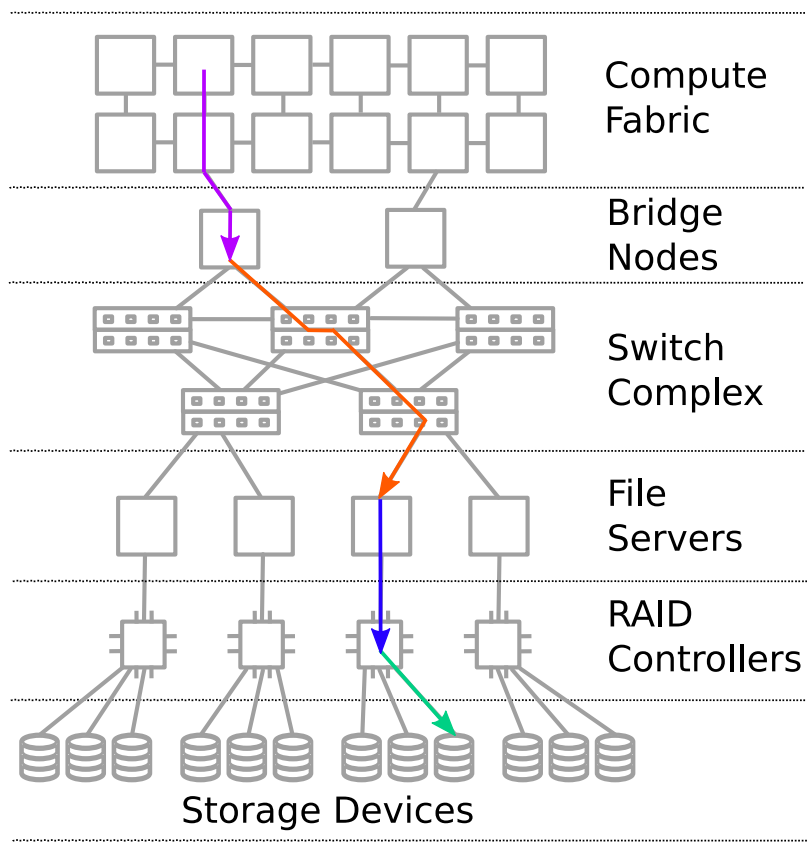
RAID Controllers

Storage Devices

- Large files (or sets of files) are not generally stored on a single storage device.
- They are distributed across multiple servers (and then each server further distributes across storage devices).
- This is referred to as **data layout** or **striping**.
- Different file systems use different striping strategies.
- It can usually be tuned to better suit your application.

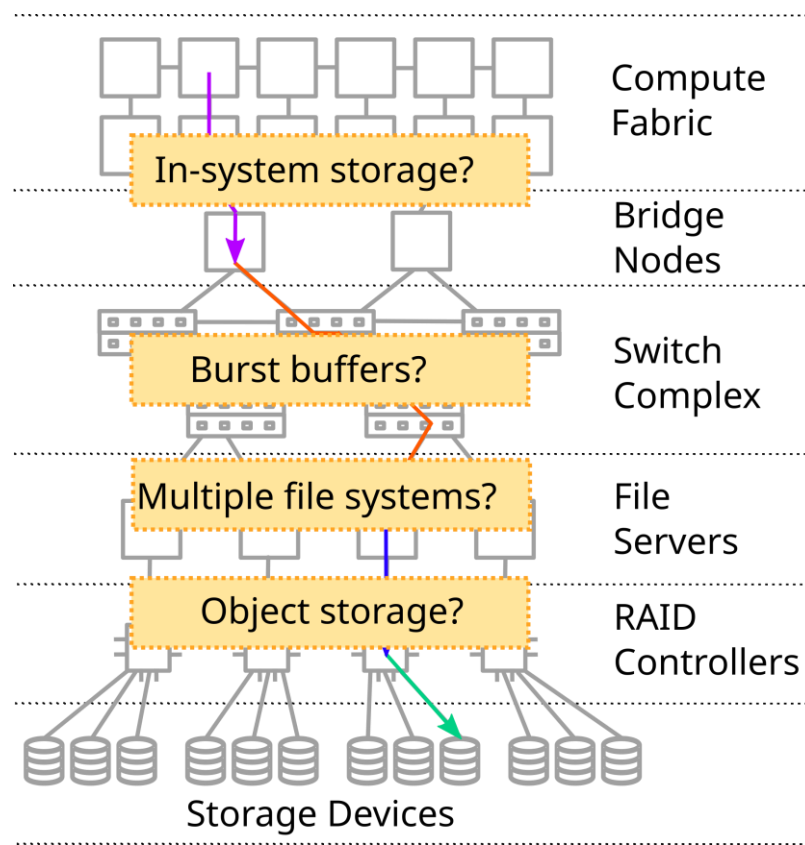Example of a single logical file striped across all available servers and storage devices

Hands on exercises:
https://github.com/radix-io/hands-on

extremecomputingtraining.anl.gov

Argonne
NATIONAL LABORATORY

# Is that all?

**HPC system storage path**

Compute Fabric

Bridge Nodes

Switch Complex

File Servers

RAID Controllers

Storage Devices

Hands on exercises:
https://github.com/radix-io/hands-on

extremecomputingtraining.anl.gov

Argonne
NATIONAL LABORATORY

# Is that all?

**HPC system storage path**



Compute Fabric

In-system storage?

Bridge Nodes

Burst buffers?

Switch Complex

Multiple file systems?

File Servers

Object storage?

RAID Controllers

Storage Devices

Each HPC storage system is a special unicorn. Some systems have:

- **In-system storage**: low latency but not shared
- **Burst buffers**: high performance with limited capacity
- **Multiple file systems**: storage systems optimized for different kinds of data
- **Object stores**: alternative methods of organizing data

Don't worry. The tools and techniques that we will teach today will help to tame this complexity. The important thing to know for now is *why* HPC storage systems need specialized techniques.

ATPESC2025

Hands on exercises:
https://github.com/radix-io/hands-on

extremecomputingtraining.anl.gov

Argonne
NATIONAL LABORATORY

# Presenting storage to HPC applications

A parallel file system can be accessed just like any other file system:

- open() / close() / read() / write() for binary data
- fopen() / fclose() / fprintf() for text data
- Various language-specific bindings

- Data is organized in a hierarchy of directories and files.

- We call this API the "POSIX interface"; it is standardized across all UNIX-like systems.

- This API works, and is great for compatibility, but it was created 50 years ago before the rise of parallel computing.

```c
fd = open("foo.dat", O_CREAT|O_WRONLY, 0600);
if(fd < 0) {
    perror("open");
    return(-1);
}

ret = write(fd, &buffer, 8);
printf("wrote %d bytes\n", ret);

close(fd);
```
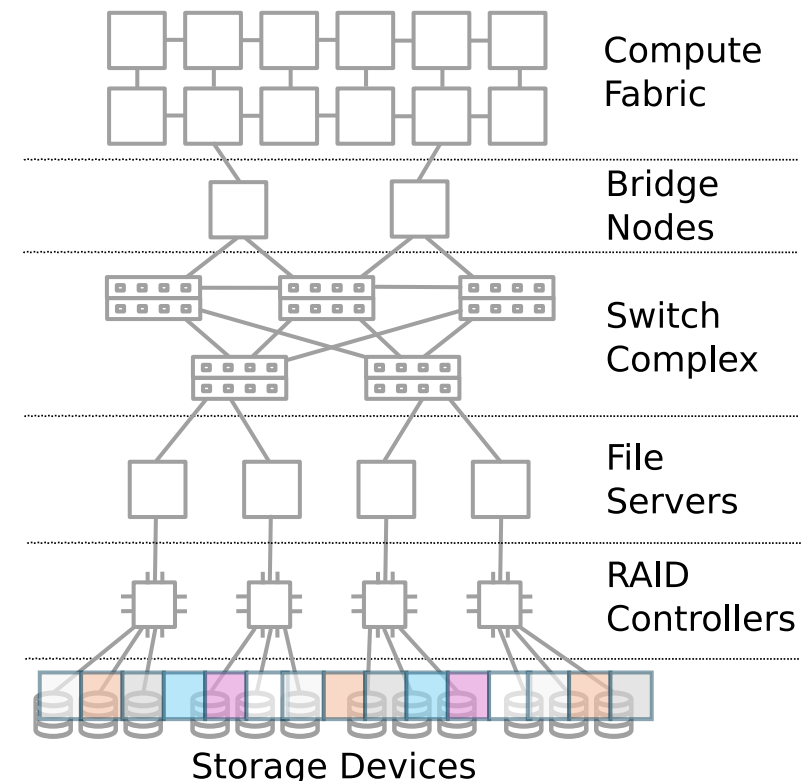
- The API has no concept of parallel access; semantics for that are largely undefined.

- Files are unstructured streams of bytes.

- File descriptors are stateful and unique to each process.

ATPESC2025

Hands on exercises:
https://github.com/radix-io/hands-on

extremecomputingtraining.anl.gov

Argonne
NATIONAL LABORATORY

# Why is it difficult to access files concurrently with POSIX?
## Example 1: writing different parts of the same file

- Consider a case in which two ranks write data simultaneously to different parts of a file.

- In this example, we have a big gap (32 MiB) between them. Assume we are writing reasonably large chunks to optimize bandwidth.

Rank 0: lseek(0); write(256 KiB);
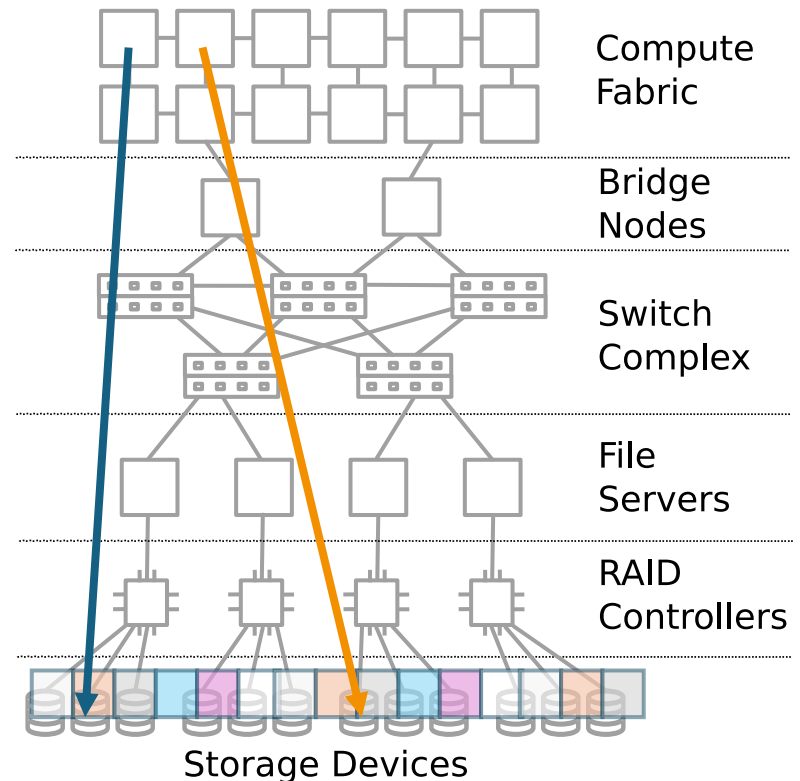Rank 1: lseek(32 MiB); write(256 KiB);

Compute Fabric

Bridge Nodes

Switch Complex

File Servers

RAID Controllers

Storage Devices

Hands on exercises:
https://github.com/radix-io/hands-on

extremecomputingtraining.anl.gov

Argonne
NATIONAL LABORATORY

# Why is it difficult to access files concurrently with POSIX?
## Example 1: writing different parts of the same file

- Consider a case in which two ranks write data simultaneously to different parts of a file.

- In this example, we have a big gap (32 MiB) between them. Assume we are writing reasonably large chunks to optimize bandwidth.

- ☑ The writes probably map to different servers and devices.

- ☑ There is no device contention, and both I/O operations can be executed at the same time.

- ❓ There is also **no coordination** of which path each write will take, though, and eventually you will want access adjacent data...

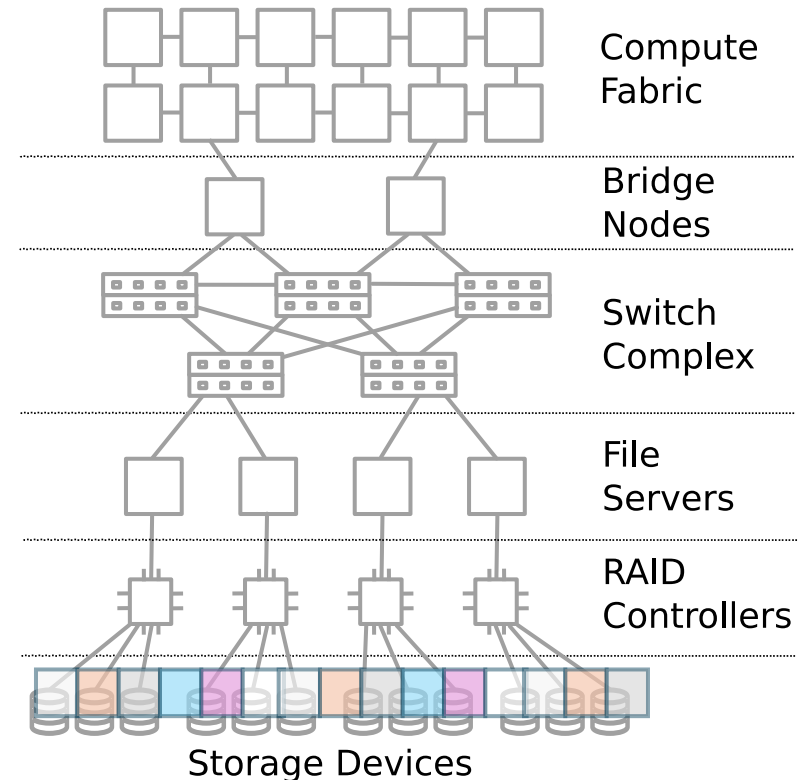Rank 0: lseek(0); write(256 KiB);
Rank 1: lseek(32 MiB); write(256 KiB);



Compute Fabric

Bridge Nodes

Switch Complex

File Servers

RAID Controllers

Storage Devices

Hands on exercises:
https://github.com/radix-io/hands-on

extremecomputingtraining.anl.gov

# Why is concurrent access hard?
## Example 2: writing adjacent parts of the same file

- Consider a case in which two ranks write data simultaneously to different parts of a file.

- In this case the writes still don't overlap, but they access adjacent bytes.

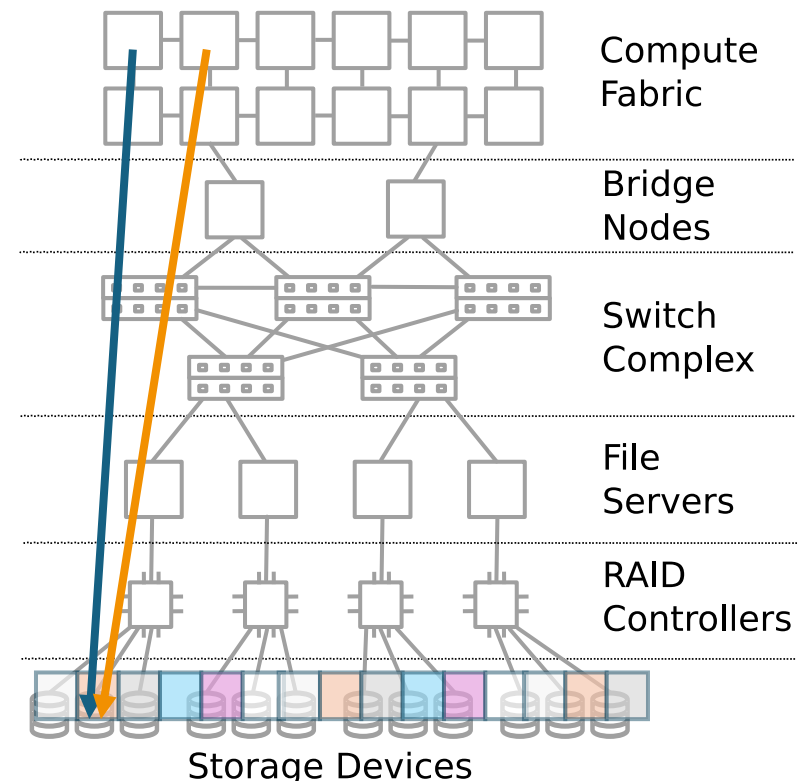Rank 0: lseek(0); write(256 KiB);

Rank 1: lseek(256 KiB); write(256 KiB);



Compute Fabric

Bridge Nodes

Switch Complex

File Servers

RAID Controllers

Storage Devices

Hands on exercises:
https://github.com/radix-io/hands-on

extremecomputingtraining.anl.gov

Argonne
NATIONAL LABORATORY

# Why is concurrent access hard?
## Example 2: writing adjacent parts of the same file

- Consider a case in which two ranks write data simultaneously to different parts of a file.

- In this case the writes still don't overlap, but they access adjacent bytes.

- ❓ The writes are likely to access the same server and storage device because of locality.

- ❌ Counterintuitively, this **almost certainly causes conflicting access**. The granularity of caching and locking has no relationship to the access size.

- ❌ Uncoordinated adjacent access can cause "false sharing" and serialize I/O operations that should have proceeded in parallel.

Rank 0: lseek(0); write(256 KiB);
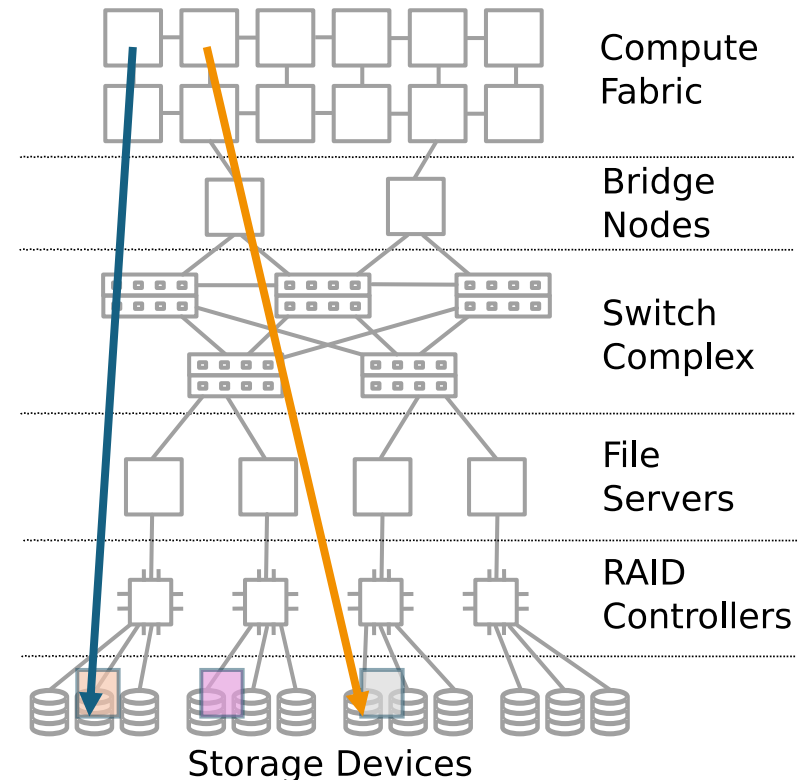Rank 1: lseek(256 KiB); write(256 KiB);



Compute Fabric

Bridge Nodes

Switch Complex

File Servers

RAID Controllers

Storage Devices

Hands on exercises:
https://github.com/radix-io/hands-on

extremecomputingtraining.anl.gov

# Why is concurrent access hard?
## Example 3: writing separate files

- Consider a case in which two ranks write data simultaneously to different files.

- There is no possibility of I/O conflict. That should be good, right?

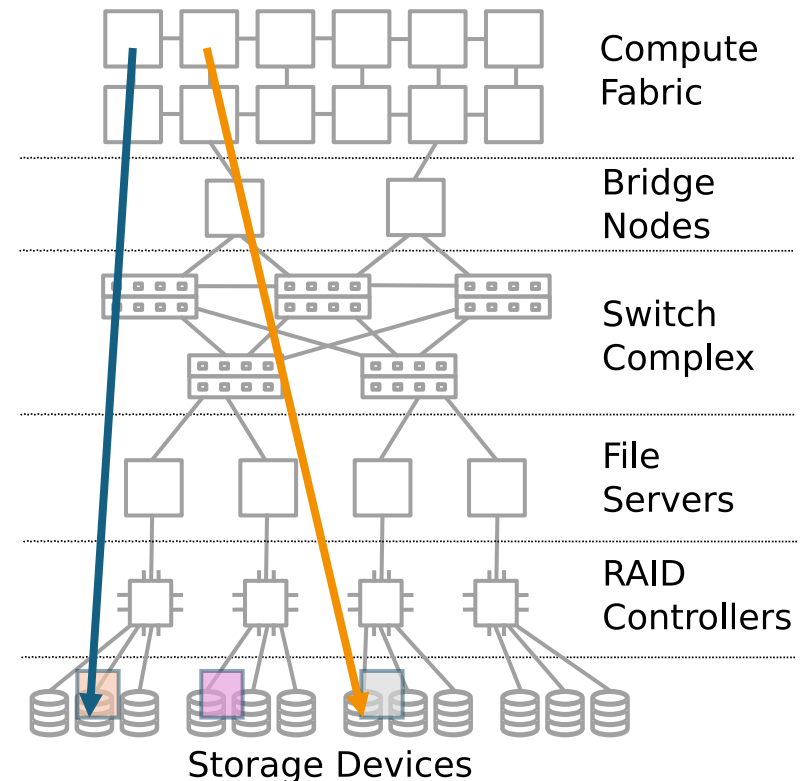Rank 0: open("a"); write(256 KiB);

Rank 1: open("b"); write(256 KiB);



Compute Fabric

Bridge Nodes

Switch Complex

File Servers

RAID Controllers

Storage Devices

# Why is concurrent access hard?
## Example 3: writing separate files

- Consider a case in which two ranks write data simultaneously to different files.

- There is no possibility of I/O conflict. That should be good, right?

- ☑ The writes are indeed issued to independent servers and storage devices. This probably works well at small scale.

Rank 0: open("a"); write(256 KiB);
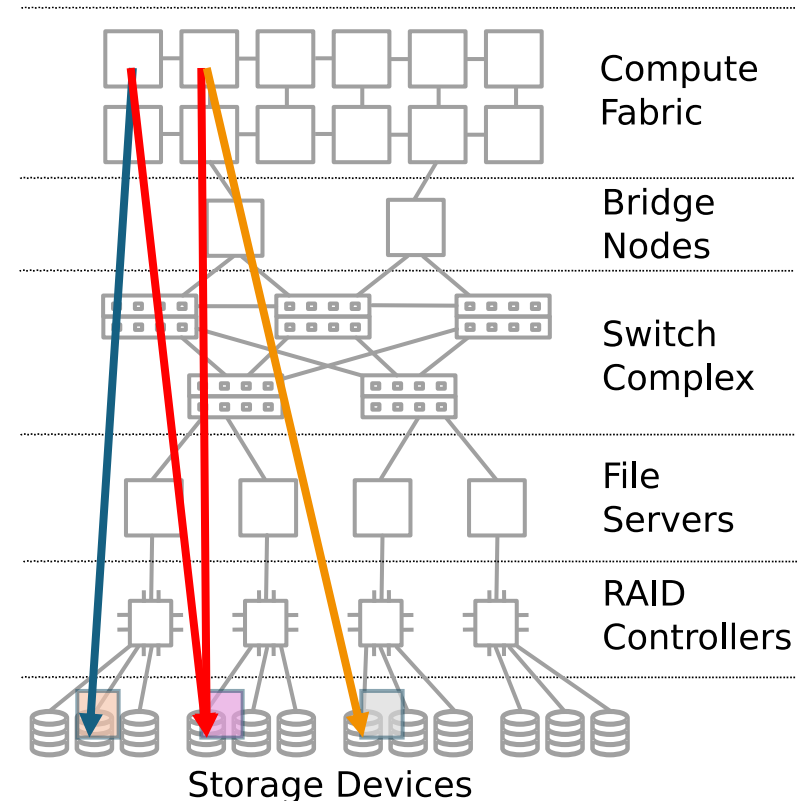
Rank 1: open("b"); write(256 KiB);



Compute Fabric

Bridge Nodes

Switch Complex

File Servers

RAID Controllers

Storage Devices

Hands on exercises:
https://github.com/radix-io/hands-on

extremecomputingtraining.anl.gov

# Why is concurrent access hard?
## Example 3: writing separate files

- Consider a case in which two ranks write data simultaneously to different files.

- There is no possibility of I/O conflict. That should be good, right?

- ☑ The writes are indeed issued to independent servers and storage devices.  This probably works well at small scale.

- ☒ Directories are hierarchical, though, so processes will conflict at open() time to coordinate access to the parent directory.  This problem gets progressively worse at scale.

- ☒ It also makes the file system spend more time managing metadata (names, permissions, attributes) than performing productive data transfer.

- ☒ It also (eventually) places more burden on the user to manage files.

Hands on exercises:
https://github.com/radix-io/hands-on

Rank 0: open("a"); write(256 KiB);
Rank 1: open("b"); write(256 KiB);



Compute Fabric

Bridge Nodes

Switch Complex

File Servers

RAID Controllers

Storage Devices

extremecomputingtraining.anl.gov

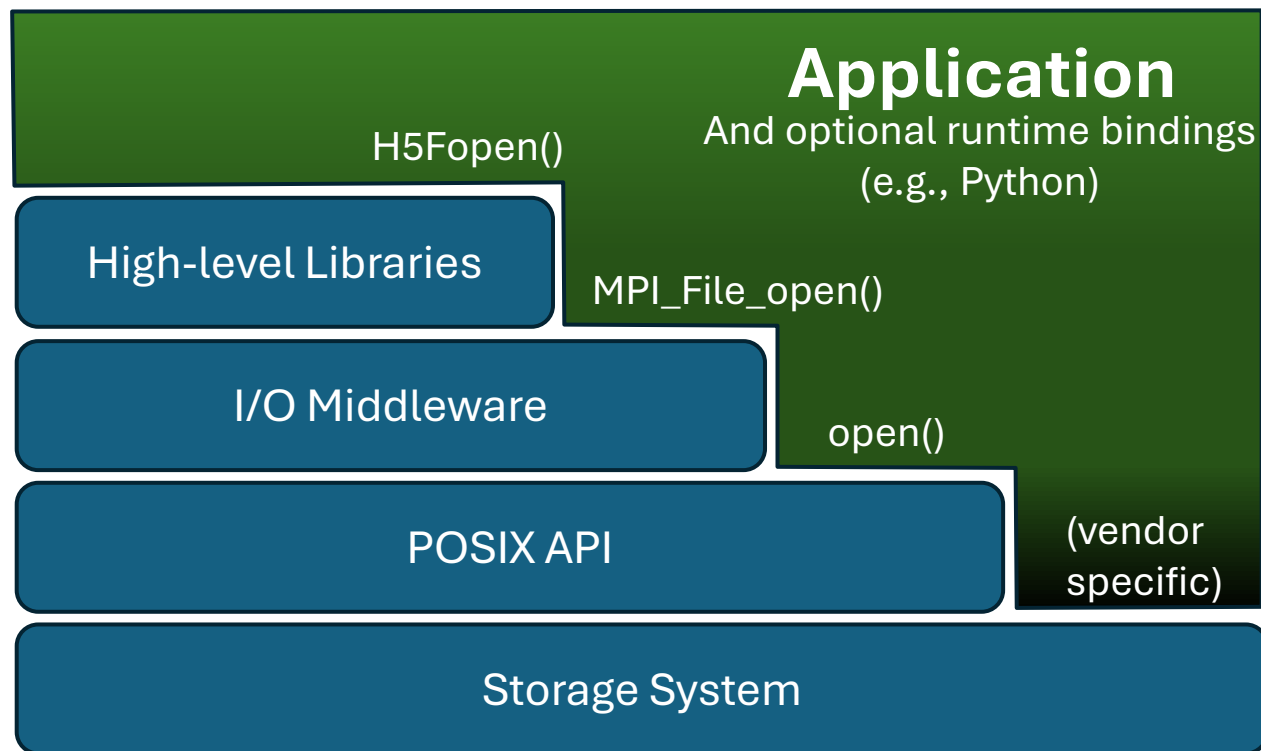# Why is concurrent access hard?
## The common theme

There is a common underlying problem in each of the preceding examples:

**Fundamentally, the sequential POSIX API cannot describe a complete, coordinated parallel access pattern to the file system.**

Because each process issues I/O operations independently, the storage system must service each one in isolation (even if there are thousands or even millions in flight). There isn't much opportunity to aggregate or structure the flow of data.

Hands on exercises:
https://github.com/radix-io/hands-on

extremecomputingtraining.anl.gov

Argonne
NATIONAL LABORATORY

# Help is on the way: Application-level I/O interface layers

**Application**
And optional runtime bindings (e.g., Python)

H5Fopen()

High-level Libraries

MPI_File_open()

I/O Middleware

open()

POSIX API

(vendor specific)

Storage System

- High level interfaces **translate between scientific data concepts (hierarchical, multi-dimensional, multi-variable, with rich metadata) and file system constructs (directories and opaque files)** so that you don't have to. Benefits include:

  - Performance portability

  - Expressive interfaces

  - Future proofing

  - Interoperability

- Don't worry if you have to use POSIX; we can help with that too!

- We will discuss all of these layers, and how to use them effectively, in depth today.

Hands on exercises:
https://github.com/radix-io/hands-on

**ATPESC**2025

**Argonne**
NATIONAL LABORATORY

# One final warning: even if you do everything right, … performance can still be surprising



- – Thousands of storage devices will *never* perform perfectly at the same time.

- – You are sharing storage with many other users across multiple HPC systems.

- – You are also sharing storage with remote transfers, tape archives, and more.

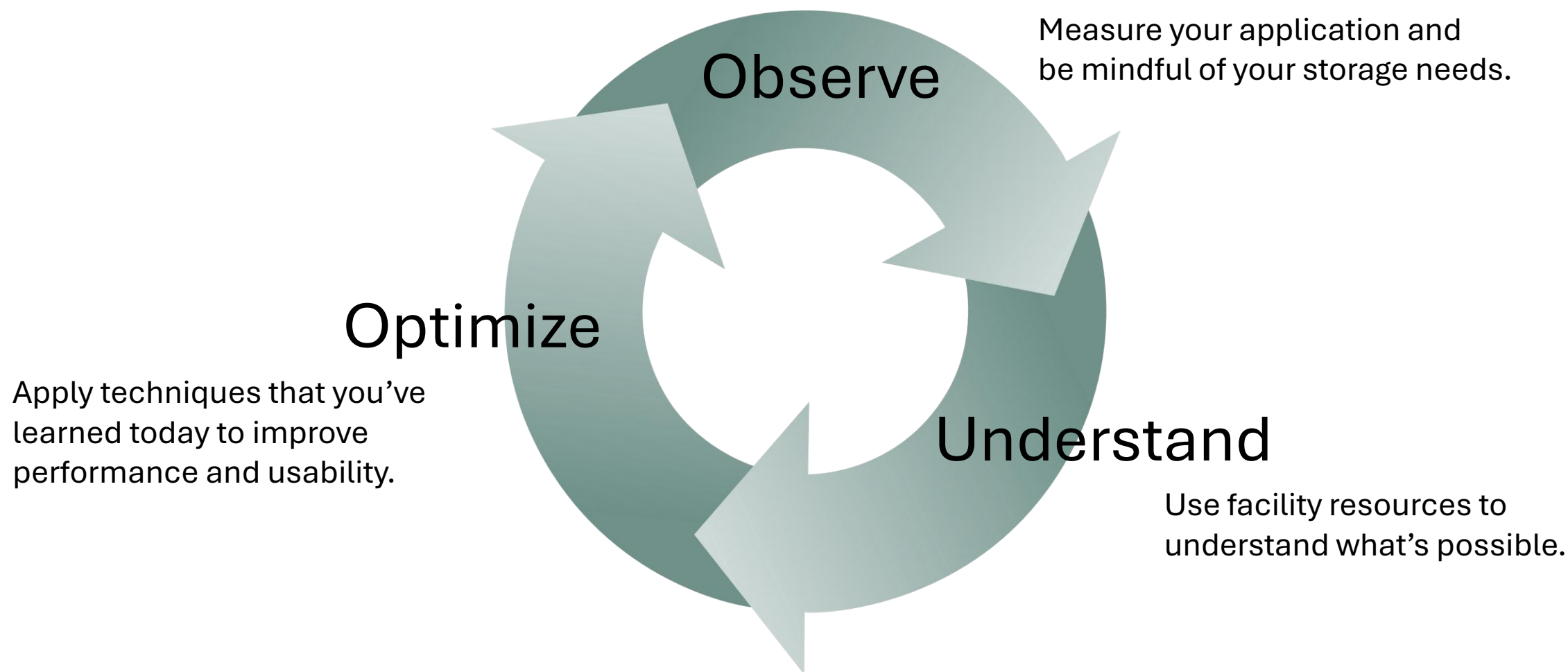Compute nodes belong exclusively to you during a job allocation, but the storage system does not. ***Storage performance varies in ways that are fundamentally different from compute performance.***

Best practice: take multiple samples of I/O performance to fully understand subtle changes. Want a demonstration of what variance you can expect? See the hands-on/variance exercise.

ATPESC2025

Hands on exercises:
https://github.com/radix-io/hands-on

extremecomputingtraining.anl.gov

Argonne
NATIONAL LABORATORY

# Parting thoughts on the Principles of HPC I/O: HPC I/O optimization is an ongoing process



**Observe**

Measure your application and be mindful of your storage needs.

**Understand**

Use facility resources to understand what's possible.

**Optimize**

Apply techniques that you've learned today to improve performance and usability.

ATPESC 2025

Hands on exercises:
https://github.com/radix-io/hands-on

extremecomputingtraining.anl.gov

Argonne
NATIONAL LABORATORY

# HPC storage systems



We mentioned earlier that each HPC system has its own unique storage architecture.

What kinds of high-performance storage systems can you find at the Office of Science compute facilities?

Let's take a quick tour!

Hands on exercises:
https://github.com/radix-io/hands-on

extremecomputingtraining.anl.gov

Argonne
NATIONAL LABORATORY

# HPC storage systems

**Oak Ridge (OLCF)**        Lawrence Berkeley (NERSC)        Argonne (ALCF)

Facility documentation is your friend!  Some of the systems that we discuss will change over time.

- Orion parallel file system (Lustre)
  - Multiple storage tiers are transparently combined:
    - Metadata tier: ~10 PiB capacity
    - Performance tier: ~11 PiB capacity, ~10 TiB/s performance
    - Capacity tier: ~680 PiB capacity, ~5 TiB/s performance
  - File layouts: Orion uses a sophisticated progressive file layout. Unlike some Lustre file systems, it is **not** recommended for you to manually change file stripe settings on Orion!

- Node-local SSDs on Frontier
  - NVMe drives provide 5 GiB/s read and 2 GiB/s write per node.
  - Excellent for latency, but they are not shared between nodes or accessible outside of your job.
  - You must explicitly request them in your job script.
  - You must explicitly stage data on and off of them.

# HPC storage systems

Oak Ridge (OLCF)          **Lawrence Berkeley (NERSC)**          Argonne (ALCF)

- Perlmutter Scratch parallel file system (Lustre)
  - Unlike the Orion file system at OLCF, the Perlmutter scratch file system consists of a single all-flash (high-performance) tier of storage.
  - 35 PiB capacity, 5 TiB/s throughput
  - File layouts: 1 server per file by default
    - This layout is great for small files accessed by single processes, but you'll want to increase the striping if you need to access large files in parallel.

# HPC storage systems

Oak Ridge (OLCF)          Lawrence Berkeley (NERSC)          **Argonne (ALCF)**

- Flare and Eagle parallel file systems (Lustre)
  - Similar caveats as Polaris Scratch: you may need to tune the Lustre settings for your use case; the default settings are only optimal for small, non-shared files.
- Polaris also has local NVMe drives
  - Similar caveats as Frontier local disks: you need to think about how to stage data.
- Aurora also has a DAOS storage system
  - 230 PiB capacity, 30 TiB/s throughput.
  - Under the covers it is a distributed object store.
    - You can access it like a conventional file system or use libraries like HDF5 or PyTorch that have been adapted to access DAOS natively.
  - We'll learn more about how to use it this afternoon!  DAOS is remarkably fast but has a little bit more of a learning curve.

ATPESC2025

Hands on exercises:
https://github.com/radix-io/hands-on

extremecomputingtraining.anl.gov

Argonne
NATIONAL LABORATORY

# Thank you!

# Any questions before we move on to the next presentation?

# ARGONNE TRAINING PROGRAM ON EXTREME-SCALE COMPUTING