# SCA/HPCAsia 2026 Tutorial: Accelerating HPC Application I/O with Fast Node-Local Storage

**Chen Wang**

chen.wang@ntu.edu.sg

Assistant Professor

Nanyang Technological University

**Jae-Sung Yeom**

yeom1@llnl.gov

Computer Scientist

Lawrence Livermore National Laboratory

*Contributers: Hariharan Devarajan, Cameron Stanavige,*
*Michael Brim, Kathryn Mohror, Michela Taufer, Ian Lumsden*

Tutorial link: https://wangchen.io/unifyfs-dyad-tutorial

# Schedule

**Tutorial Link:**

https://wangchen.io/unifyfs-dyad-tutorial

| Time | Event |
|---|---|
| 13:30 - 14:00 | Introduction and Motivation |
| 14:00 - 14:30 | Deep Dive: UnifyFS |
| 14:30 - 14:50 | Break |
| 14:50 - 15:20 | Deep Dive: DYAD |
| 15:20 - 15:50 | DYAD Walkthrough |
| 15:50 - 16:15 | UnifyFS Walkthrough |
| 16:15 - 16:30 | Wrap-up and Q&A |

**Cluster account request:**

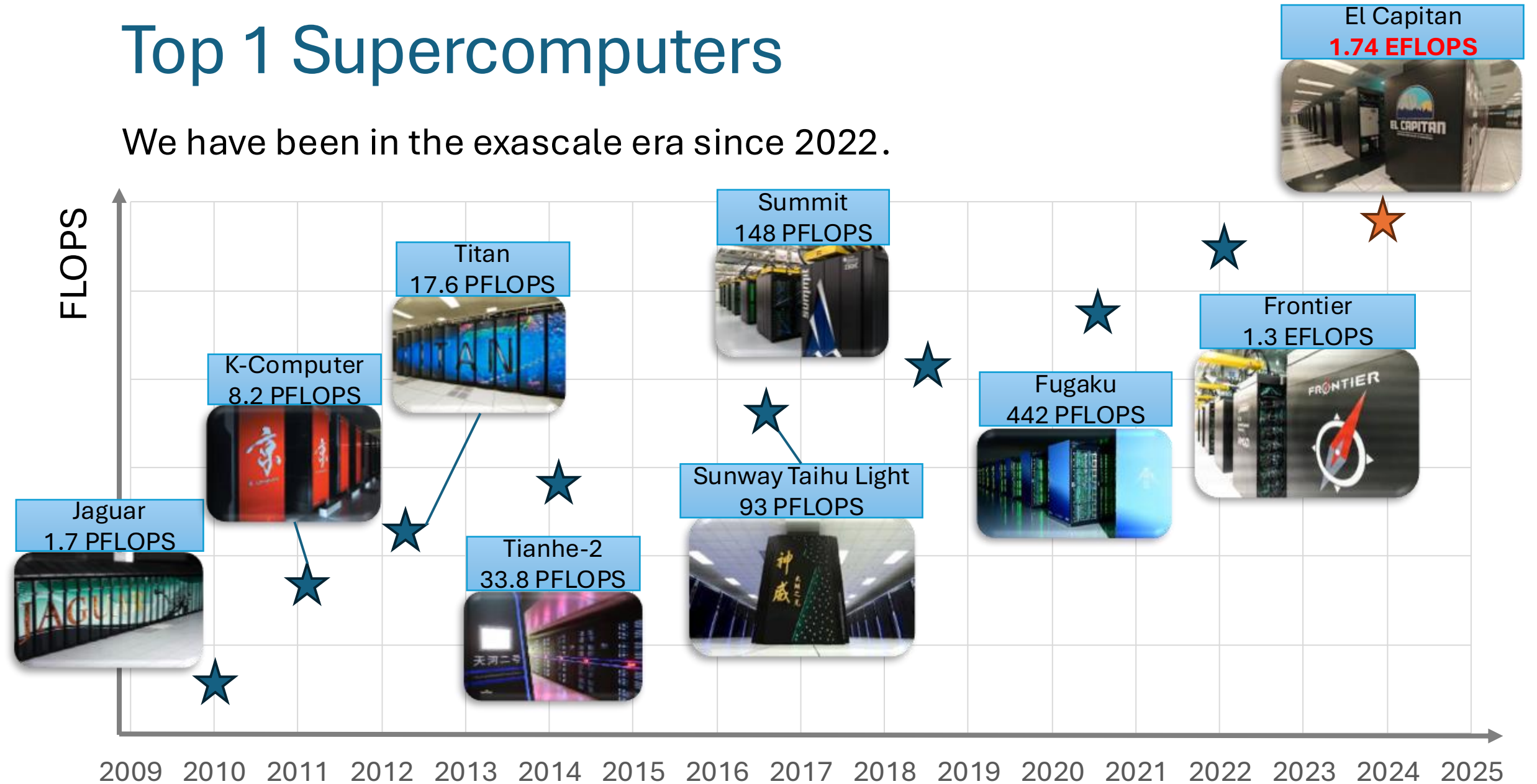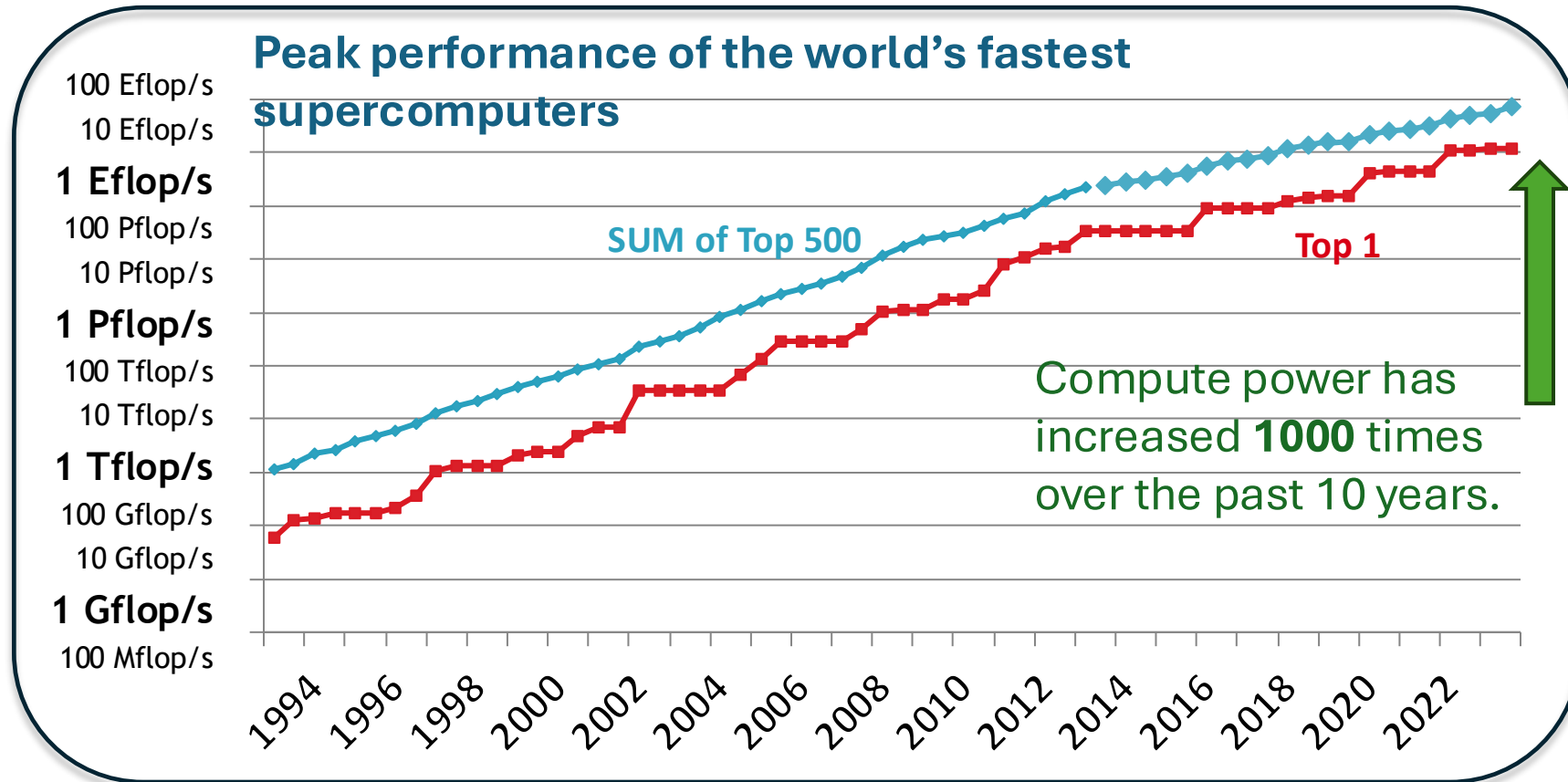https://forms.gle/Ya1WRVn9cxdRSMmXA



Scan me!

```
ssh -p 2223 username@13.215.163.223
password: P@ssw0rd
```

# Top 1 Supercomputers

We have been in the exascale era since 2022.



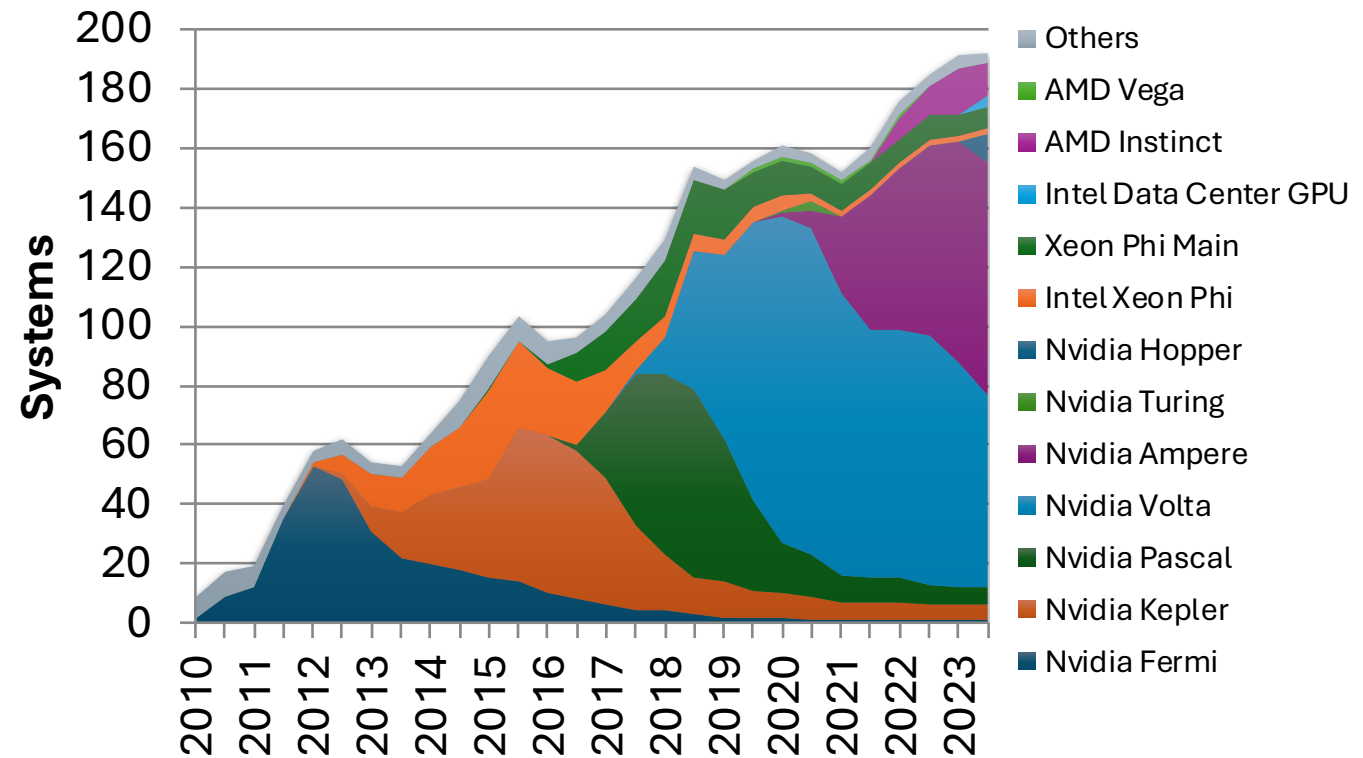El Capitan
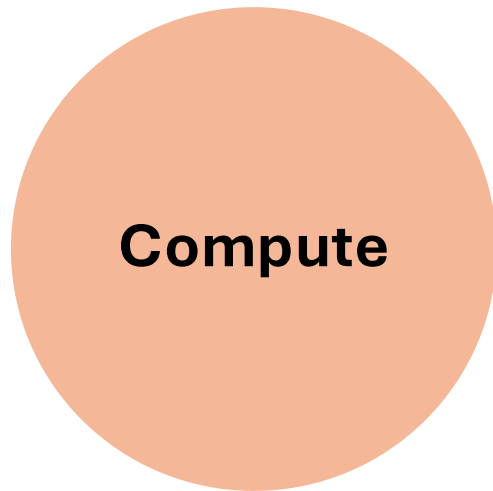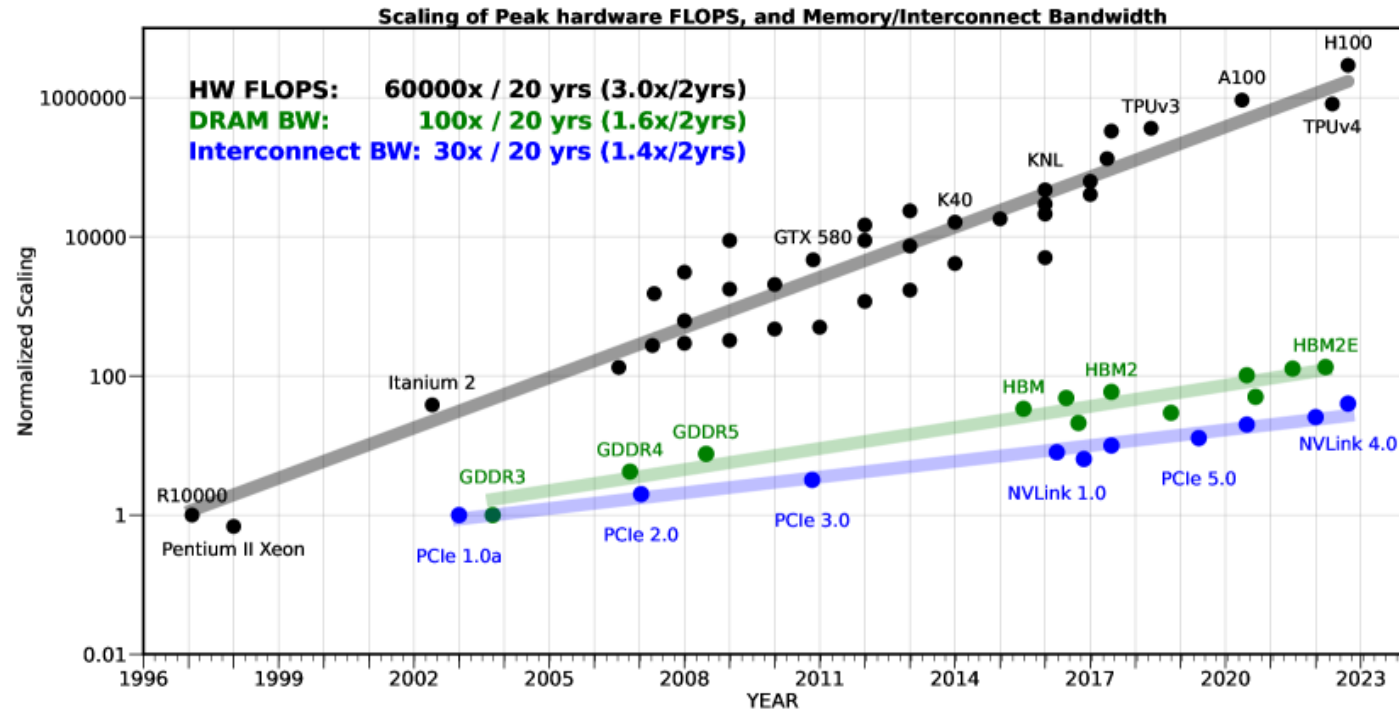**1.74 EFLOPS**

Summit
148 PFLOPS

Titan
17.6 PFLOPS

K-Computer
8.2 PFLOPS

Frontier
1.3 EFLOPS

Fugaku
442 PFLOPS

Jaguar
1.7 PFLOPS

Sunway Taihu Light
93 PFLOPS

Tianhe-2
33.8 PFLOPS

FLOPS

2009  2010  2011  2012  2013  2014  2015  2016  2017  2018  2019  2020  2021  2022  2023  2024  2025

Tutorial link: https://wangchen.io/unifyfs-dyad-tutorial

# Compute Power Has Increased Significantly

**Peak performance of the world's fastest supercomputers**

100 Eflop/s
10 Eflop/s
**1 Eflop/s**
100 Pflop/s
10 Pflop/s
**1 Pflop/s**
100 Tflop/s
10 Tflop/s
**1 Tflop/s**
100 Gflop/s
10 Gflop/s
**1 Gflop/s**
100 Mflop/s

SUM of Top 500

Top 1

Compute power has increased **1000** times over the past 10 years.

1994 1996 1998 2000 2002 2004 2006 2008 2010 2012 2014 2016 2018 2020 2022

*Data Source: Top500.org*

Tutorial link: https://wangchen.io/unifyfs-dyad-tutorial

# GPUs Have Become the Primary Workhorse

**Compute**



Over 1/3 of top 500 systems have accelerators

Legend (top to bottom): Others, AMD Vega, AMD Instinct, Intel Data Center GPU, Xeon Phi Main, Intel Xeon Phi, Nvidia Hopper, Nvidia Turing, Nvidia Ampere, Nvidia Volta, Nvidia Pascal, Nvidia Kepler, Nvidia Fermi

X-axis: 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, 2023
Y-axis: Systems (0–200)

# AI and Memory Wall



**GPU FLOPS vs. Memory Bandwidth**

The performance gap is expected to grow at 50% per year.

*Gholami, Amir, Zhewei Yao, Sehoon Kim, Coleman Hooper, Michael W. Mahoney, and Kurt Keutzer. "Ai and memory wall." IEEE Micro 44, no. 3 (2024): 33-39.*

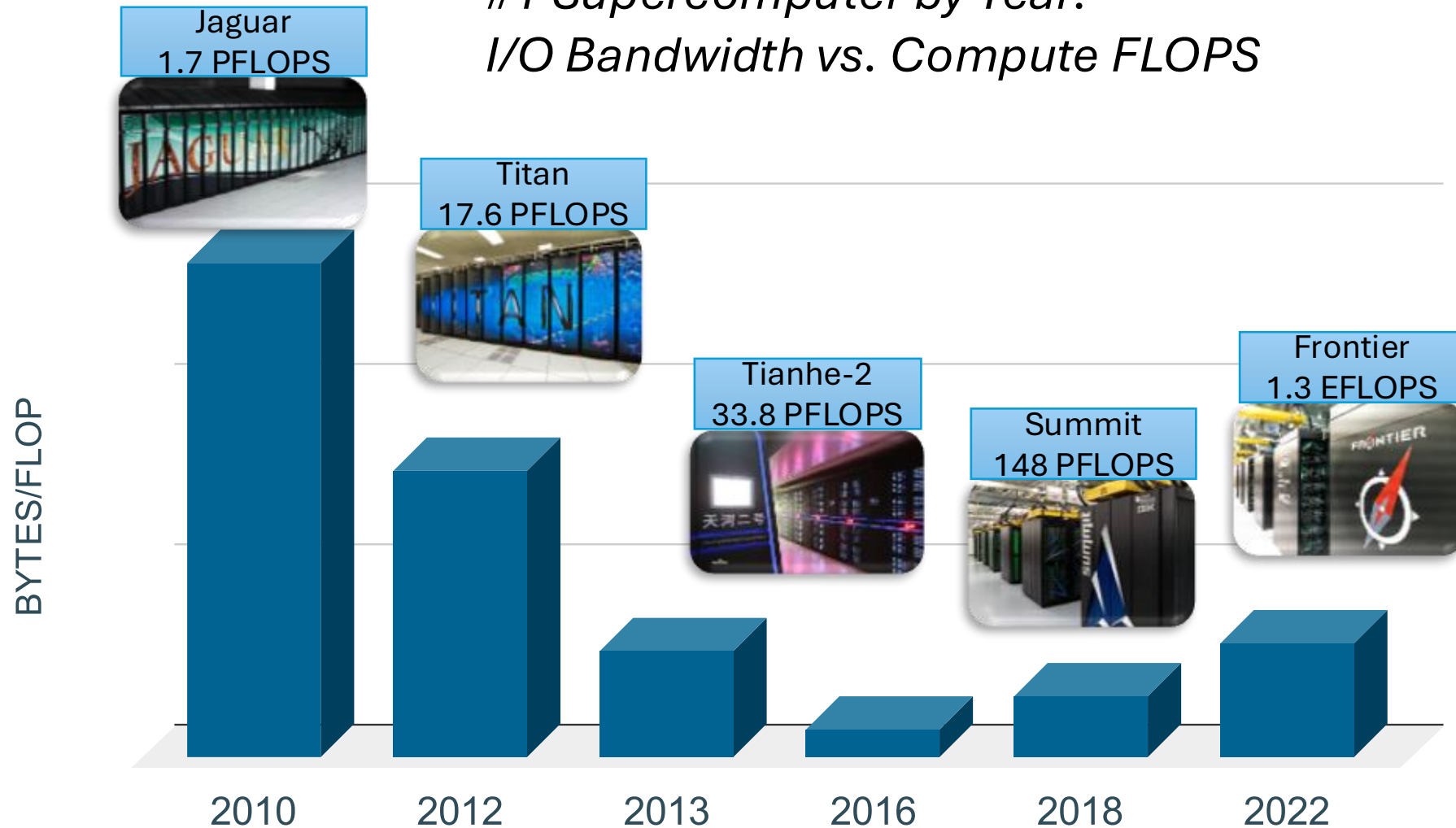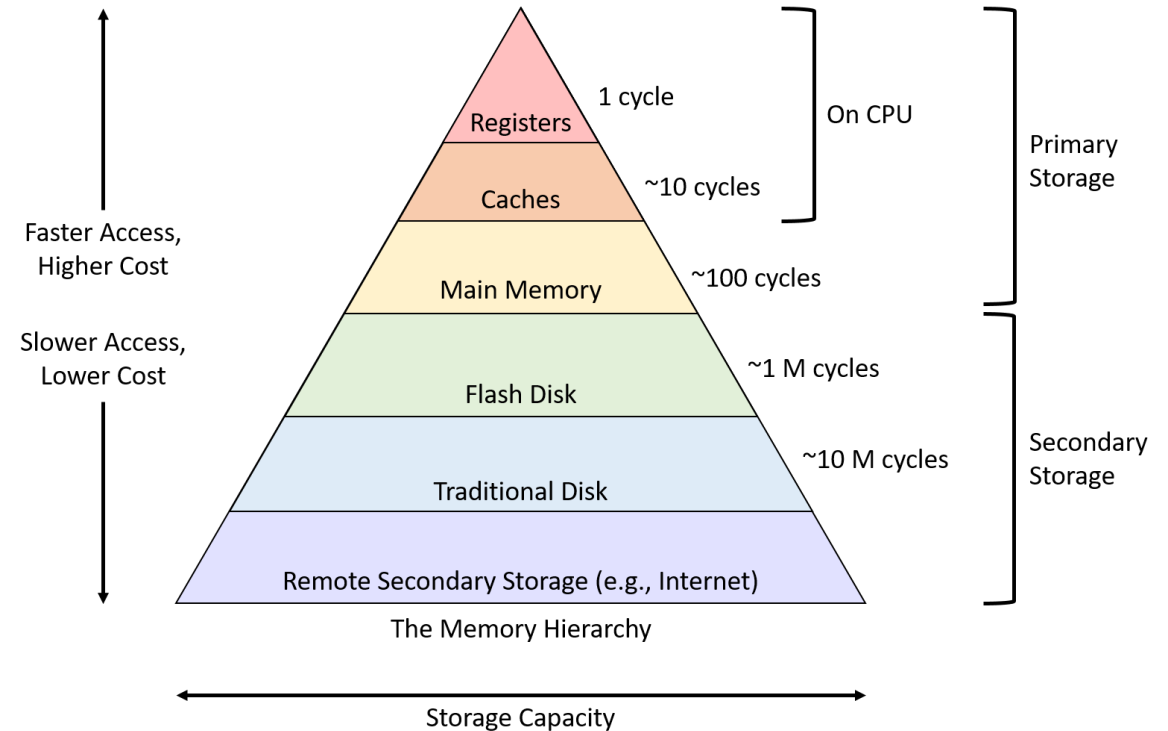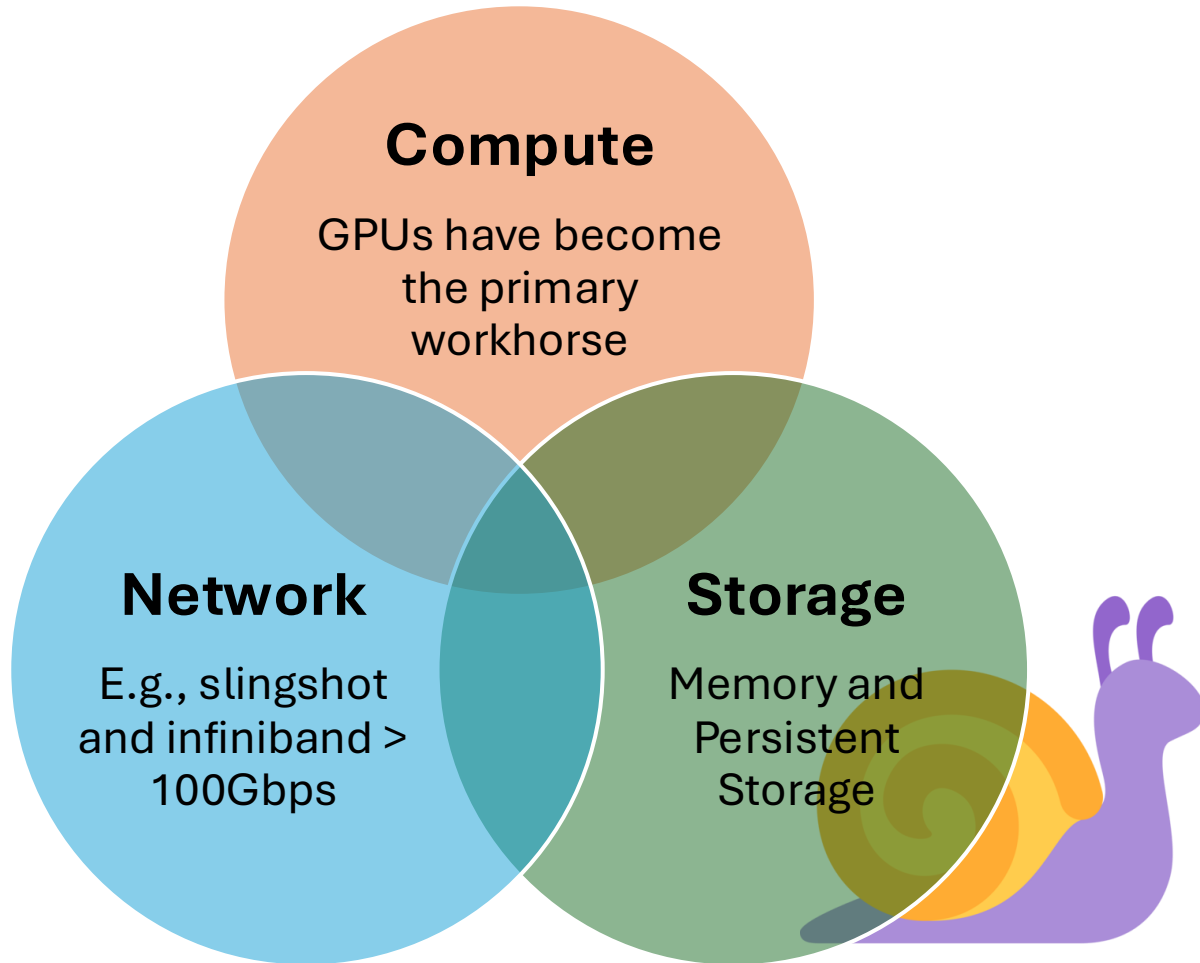# AI and Memory Wall



**Transformer Size vs. Memory Capacity**



**The Evolution of GPT Models**

*Gholami, Amir, Zhewei Yao, Sehoon Kim, Coleman Hooper, Michael W. Mahoney, and Kurt Keutzer. "Ai and memory wall." IEEE Micro 44, no. 3 (2024): 33-39.*

Tutorial link: https://wangchen.io/unifyfs-dyad-tutorial

# I/O is even Slower

*#1 Supercomputer by Year:*
*I/O Bandwidth vs. Compute FLOPS*

BYTES/FLOP

Jaguar
1.7 PFLOPS

Titan
17.6 PFLOPS

Tianhe-2
33.8 PFLOPS

Summit
148 PFLOPS

Frontier
1.3 EFLOPS

2010　　2012　　2013　　2016　　2018　　2022

Tutorial link: https://wangchen.io/unifyfs-dyad-tutorial

# The Deep Storage Hierarchy

**Compute**

GPUs have become the primary workhorse

**Network**

E.g., slingshot and infiniband > 100Gbps

**Storage**

Memory and Persistent Storage

Faster Access, Higher Cost

Slower Access, Lower Cost

| | |
|---|---|
| Registers | 1 cycle |
| Caches | ~10 cycles |
| Main Memory | ~100 cycles |
| Flash Disk | ~1 M cycles |
| Traditional Disk | ~10 M cycles |
| Remote Secondary Storage (e.g., Internet) | |

On CPU

Primary Storage

Secondary Storage

The Memory Hierarchy

Storage Capacity

*The entire storage hierarchy is getting deeper and more complex, and the boundary between memory and storage is steadily blurring.*

Tutorial link: https://wangchen.io/unifyfs-dyad-tutorial

# I/O Subsystem Inefficiency

A study of 4 million jobs over four years on two LLNL systems shows that

- on average, jobs which performing I/O spread I/O activities across 78.8% of their runtime.

- less than 22% write-intensive jobs perform efficient writes.

- HPC jobs are no longer write dominated

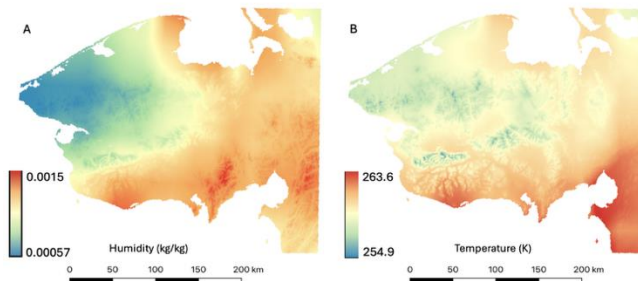*Paul, Arnab K., et al. "Understanding HPC Application I/O behavior using system level statistics." 2020 IEEE HiPC.*



**Percentage I/O (Write vs. Read) by User**

# Case Study: Energy Exascale Earth System Model (E3SM)



| File | Size |
|------|------|
| Surface (I) | 188 GB |
| Forcing (I) | 1.4 TB |
| History (O) | 134 GB |
| Restart (O) | 4.2 TB |

I:Input   O:Output

A high-resolution (1kmx1km, previously 10kmx10km) land simulation over Alaska (21.6 Million land grid cell)
Used three supercomputers: Perlmutter, Summit (#1 from 2018-2020),  and Frontier (#1 2022-2023, #2 now)



*2025 CCGRID SCALE Challenge Finalist: Dali Wang, Chen Wang, Qinglei Cao, and et.al. "Scaling Ultrahigh-Resolution E3SM Land Model for Leadership-Class Supercomputers".*

Tutorial link: https://wangchen.io/unifyfs-dyad-tutorial

# Case Study: Energy Exascale Earth System Model (E3SM)

Strong scaling results on Frontier.

- Up to 1200 nodes with I/O.
  - **Bottleneck**: 76,800 processes concurrently write to a single file.
- Up to 4000 nodes (nearly half of the Frontier) without I/O.
  - Bottleneck: initialization phase
- *Note this is only a 5-day test simulation.*
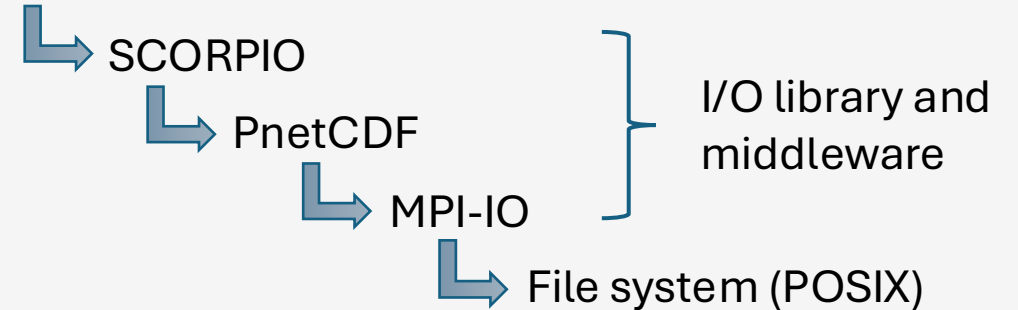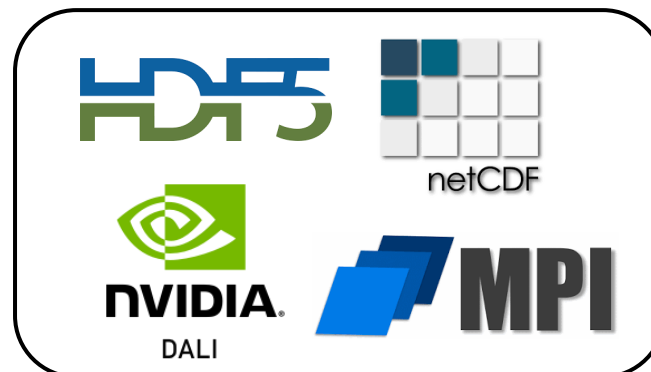- *We encountered both the scalability issue and the I/O bottleneck.*



I/O can take **up to 80%** of total execution time.

286s on I/O

*2025 CCGRID SCALE Challenge Finalist: Dali Wang, Chen Wang, Qinglei Cao, and et.al. "Scaling Ultrahigh-Resolution E3SM Land Model for Leadership-Class Supercomputers".*

Tutorial link: https://wangchen.io/unify-fs-dyad-tutorial

# How I/O Works in HPC

**HPC Application**



**The E3SM Example:**

E3SM
↳ SCORPIO
    ↳ PnetCDF
        ↳ MPI-IO
            ↳ File system (POSIX)

I/O library and middleware

**I/O Library and middleware**



**HPC File System**



Tutorial link: https://wangchen.io/unifyfs-dyad-tutorial

# Case Study: Energy Exascale Earth System Model (E3SM)



Legend: ■ I/O  ■ Compute

Y-axis: Time (seconds), 0 to 800
X-axis: Number of nodes — 150, 300, 600, 1200, 1200 (Tuned) (76,800 cores)

286s on I/O
81s on I/O

76,800 processes concurrently open/read/write a single file, causing significant congestions.
→ Delegate all I/O to one aggregator per node. Operate on one file per node.

After tuning:
- I/O time: 286s → 81s.  (3.5x)
- Write bandwidth → ~300GB/s. This is still far away from the system peak performance (5TB/s).

Tutorial link: https://wangchen.io/unifyfs-dyad-tutorial

*2025 CCGRID SCALE Challenge Finalist: Dali Wang, Chen Wang, Qinglei Cao, and et.al. "Scaling Ultrahigh-Resolution E3SM Land Model for Leadership-Class Supercomputers".*

# Optimizing HPC File Systems?

**HPC Application**



**I/O Library and middleware**



**HPC File System**



Tutorial link: https://wangchen.io/unifyfs-dyad-tutorial

# Optimizing HPC File Systems?

Traditional HPC file systems are **global resources shared by all users and jobs**. They are static and unable to adapt to different workloads, making it basically impossible to optimize for a single job. This limitation affects all applications.
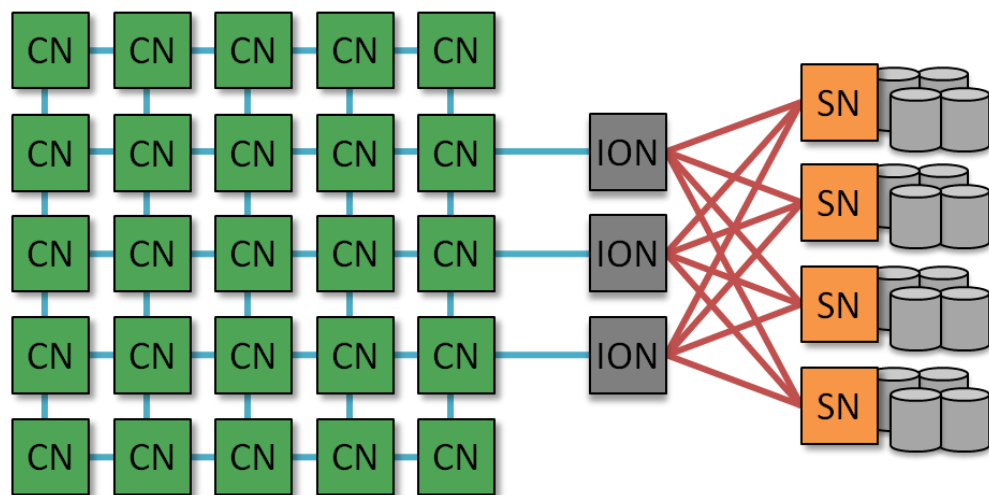
# Burst Buffers: Yet Another Storage Layer



"Tape is Dead. Disk is Tape. Flash is Disk." (at CIRD'07)

Jim Gray
(1998 Turing Award Winner)

CN: Compute Node; ION: I/O Node; SN: Storage Node

## Historical price of computer memory and storage

This data is expressed in US dollars per terabyte (TB), adjusted for inflation. "Memory" refers to random access memory (RAM), "disk" to magnetic storage, "flash" to special memory used for rapid data access and rewriting, and "solid state" to solid-state drives (SSDs).



Data source: John C. McCallum (2023); U.S. Bureau of Labor Statistics (2024)     OurWorldinData.org/technological-change | CC BY
Note: For each year, the time series shows the cheapest historical price recorded until that year. This data is expressed in constant 2020 US$.

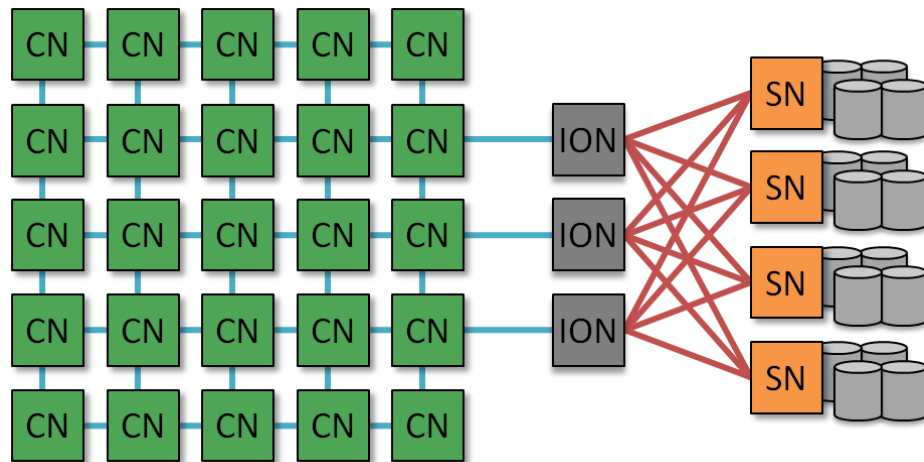*Figures courtesy of Glenn K. Lockwood.*
*https://blog.glennklockwood.com/2017/03/reviewing-state-of-art-of-burst-buffers.html*

Tutorial link: https://wangchen.io/unifyfs-dyad-tutorial

# Burst Buffers: Yet Another Storage Layer

Node-local burst buffer:

- Attach one SSD to each compute node.
  - Scales linearly.
- **Only accessiable from the attached node**.
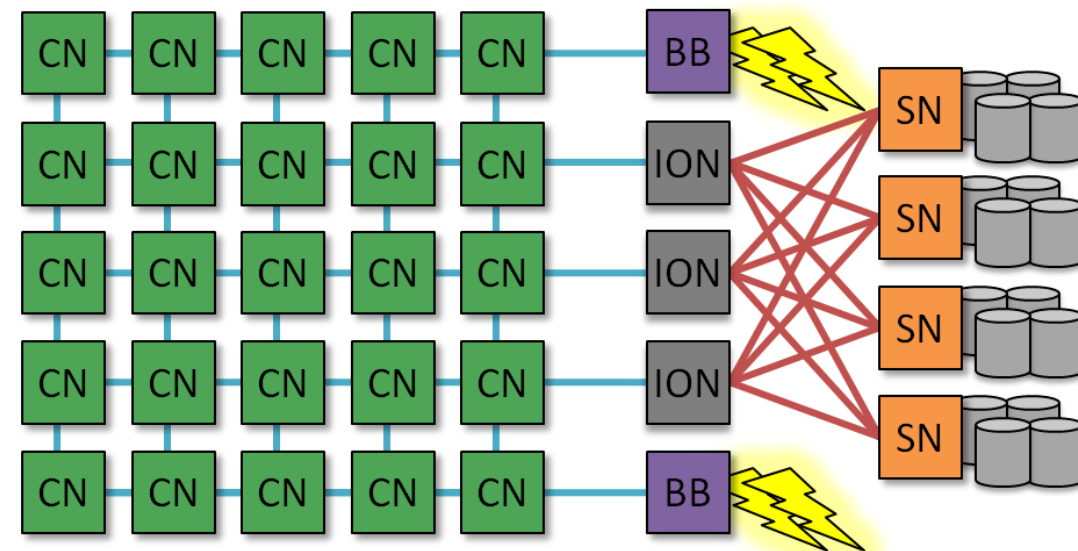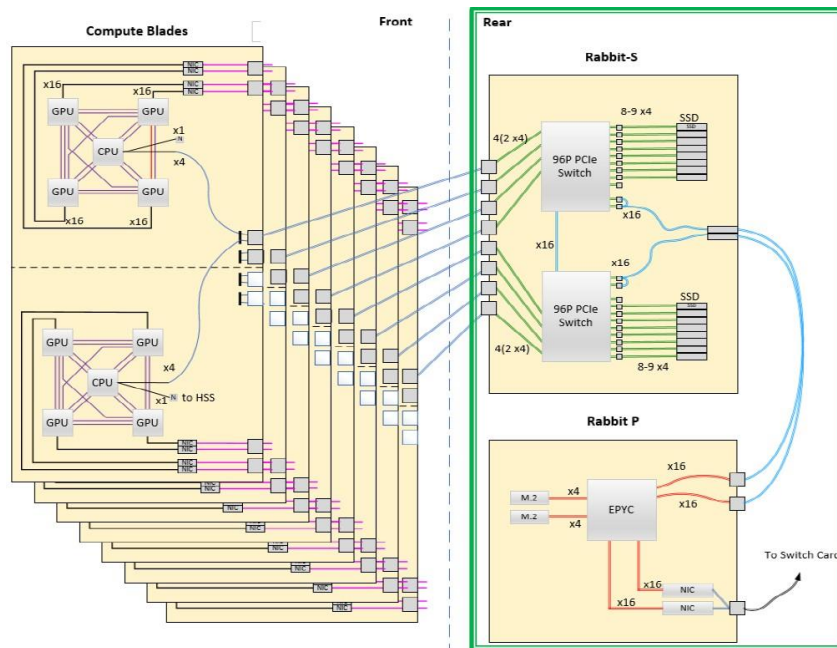  - Users need to manage data transfers across layers and between nodes.

CN: Compute Node; ION: I/O Node; SN: Storage Node

Node-local Burst Buffer

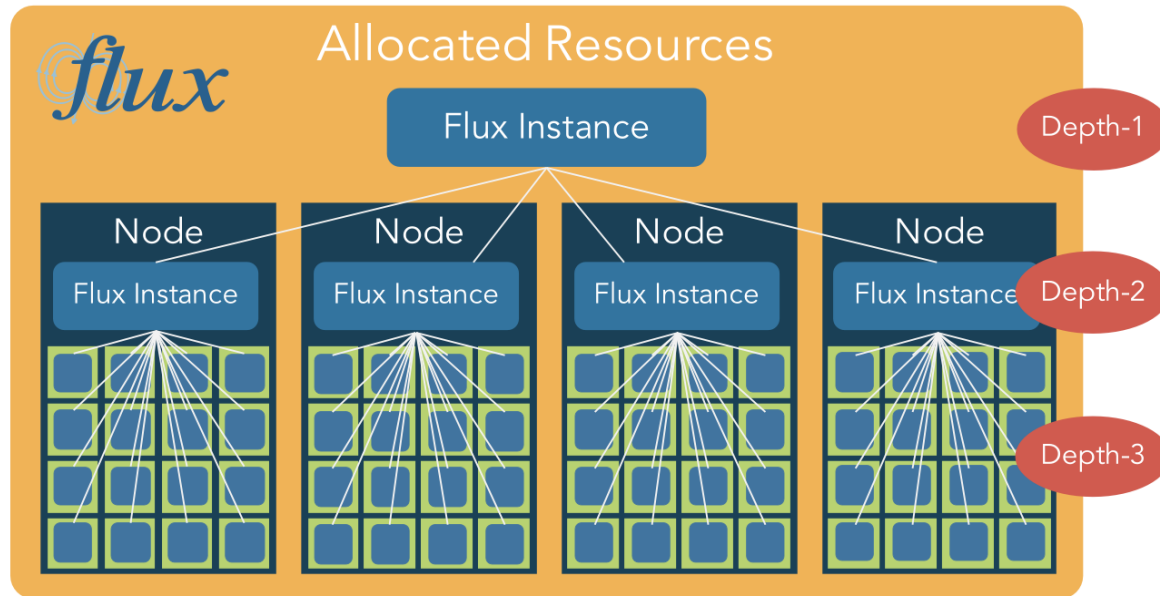# Burst Buffers: Yet Another Storage Layer

The "Rabbit" way:

- Each Rabbit node consists of *N* SSDs and one processor.
- Two Rabbit nodes sit in each rack; Each is directly connected to all compute nodes within the same rack.
- Provide a shared address space to all compute nodes.
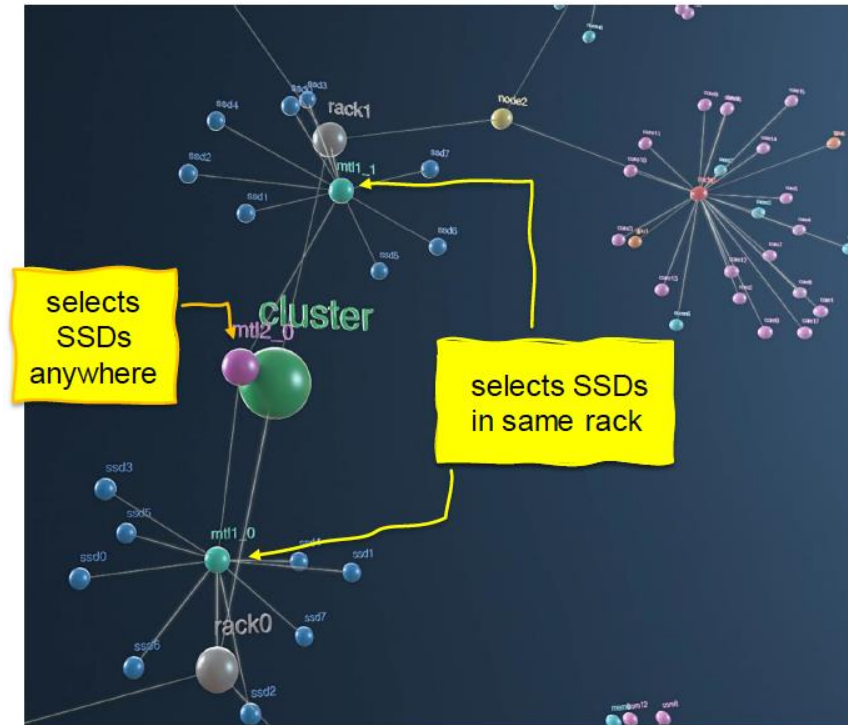


BB: Burst Buffer Node

# Flux hierarchically manages of resources and jobs for scalability



- Hierarchical and modular management
- Mitigates the centralized scheduler bottleneck
- Deployed at El Capitan
  - Ranked 1st in top500 supercomputers

# Flux's graph-based resource representation addresses challenges with Rabbit storage system



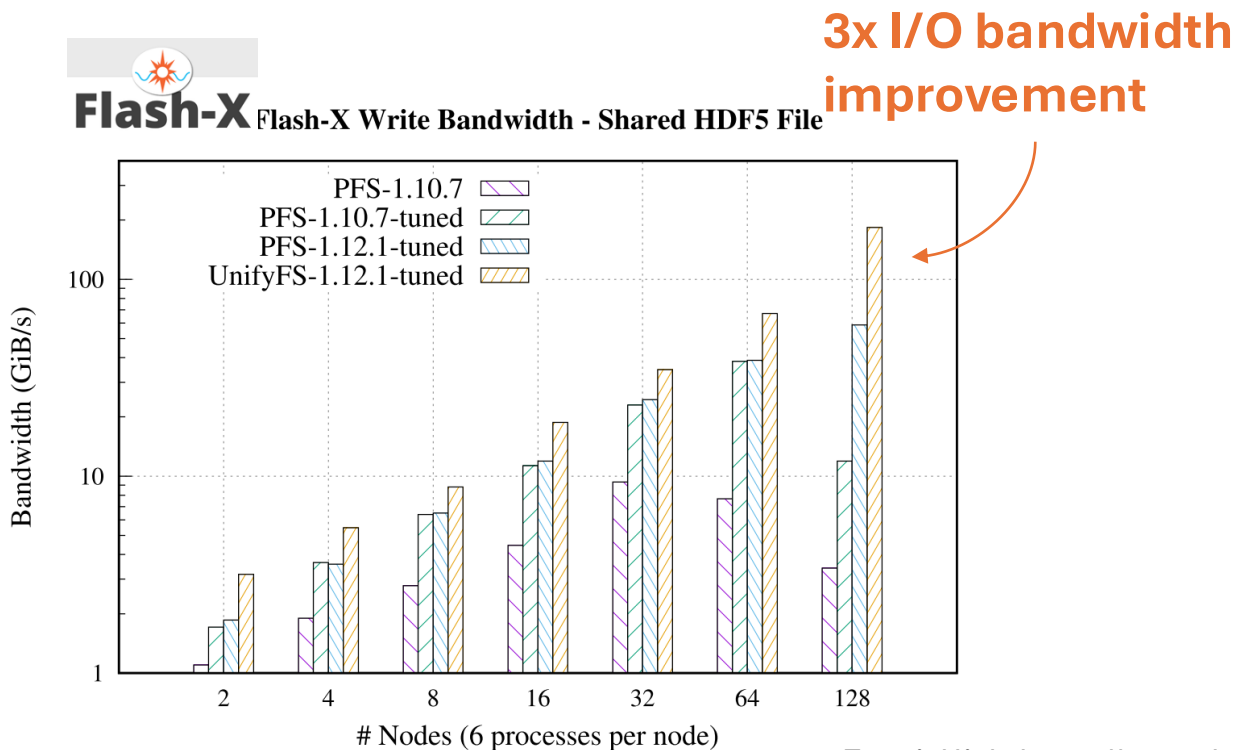- Better suited for dynamic and heterogenous resources
- Rabbit is a near-compute disaggregate storage system that can be dynamically allocated per job as either shared or node-local.

Tutorial link: https://wangchen.io/unifyfs-dyad-tutorial

# UnifyFS and DYAD

**UnifyFS**: A specialized burst buffer parallel file system for supercomputers. Designed **for write-heavy HPC applications**.
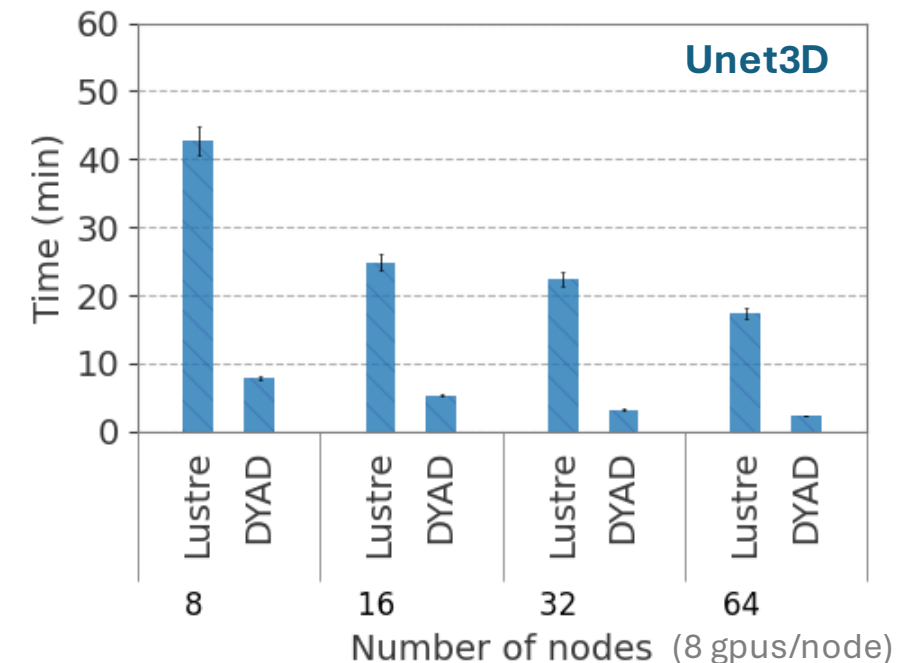
- https://github.com/LLNL/UnifyFS

**3x I/O bandwidth improvement**

**DYAD:** is a data streamer optimized for deep learning (DL) training.

- https://github.com/flux-framework/dyad

**Up to 8x training time improvement**



Flash-X Write Bandwidth - Shared HDF5 File



Tutorial link: https://wangchen.io/unifyfs-dyad-tutorial

# Thank You!

**UnifyFS**: A specialized burst buffer parallel file system for supercomputers. Our goal is to provide easy, portable, and fast support for I/O-intensive applications.

- https://github.com/LLNL/UnifyFS
- *Michael Brim, Adam Moody, Seung-Hwan Lim, Ross Miller, Swen Boehm, Cameron Stanavige, Kathryn Mohror, Sarp Oral, "UnifyFS: A User-level Shared File System for Unified Access to Distributed Local Storage," 37th IEEE International Parallel & Distributed Processing Symposium (IPDPS), May 2023.*

**DYAD:** is a data streamer optimized for deep learning (DL) training and scientific workflows.

- https://github.com/flux-framework/dyad
- DYAD data movement coordination strategy reduces network contention and thus improving data movement
- DYAD employs a hierarchical sample discovery technique to isolate metadata accesses, and that improves lookup throughput
- DYAD's novel streaming RPC over RDMA protocol boosts inter-node access
- DYAD optimizes large-scale DL workloads by 8.2x as compared to state-of-the-art solutions.

# Call for Papers

**7th International Workshop on Extreme-Scale Storage and Analysis (ESSA 2026)**, held in conjunction with IPDPS (New Orleans, May 2026)!

_Deadline: Feb. 6, 2026._



# ESSA
## IPDPS
### 7th Workshop on Extreme-Scale Storage and Analysis
Held in conjunction with IPDPS 2026, New Orleans, LA, USA

- Paper submission deadline: **February 6, 2026**
- Acceptance notification: **February 27, 2026**
- Camera-ready deadline: March 6, 2026
- Workshop date: May 26, 2026

Final extension

https://sites.google.com/view/essa-2026

---

Special Issue

# New Advances in Parallel and Distributed Computing

**Message from the Guest Editors**

This Special Issue focuses on emerging advances that enhance the performance, scalability, and intelligence of parallel and distributed computing systems. It emphasizes new technologies and paradigms in high-performance computing (HPC) in the era of AI, particularly in the context of the convergence of HPC, cloud, and edge computing. Topics of interest include, but are not limited to, the following:

- Advanced architectures for the HPC–cloud–edge computing continuum.
- Parallel and distributed algorithms for large-scale AI and large language models (LLMs).
- Data-driven optimization techniques for complex scientific and AI workflows.
- Cross-facility resource management and scheduling.
- Workflow management frameworks for distributed and federated infrastructures.
- Performance modeling, optimization, and benchmarking of distributed systems.
- Disaggregated storage and high-throughput data flow optimization.
- Green and sustainable computing strategies for data-intensive workloads.
- Integration of HPC, quantum, and AI accelerators into distributed platforms.
- Emerging applications, including digital twins, autonomous systems, IoT, and 6G networking.

---

## Electronics

**an Open Access Journal by MDPI**

**Impact Factor 2.6
CiteScore 6.1**

mdpi.com/si/267658

mdpi.com/journal/
electronics