



Wednesday, Dec 1<sup>st</sup> 2021  
Webinar`

# Agenda

2

- Welcome message from leadership
- Hermes overview
- Brief information about the Beta release
- Configuration and Deployment guide
- IOR on Hermes demo
- Frequently asked questions



Short message  
by the PI  
Dr. Xian-He Sun



## What it is

- A multi-tiered I/O buffering system for HPC
- NSF OCI CSSI Framework 1835764

## The Hermes System

## Who we are

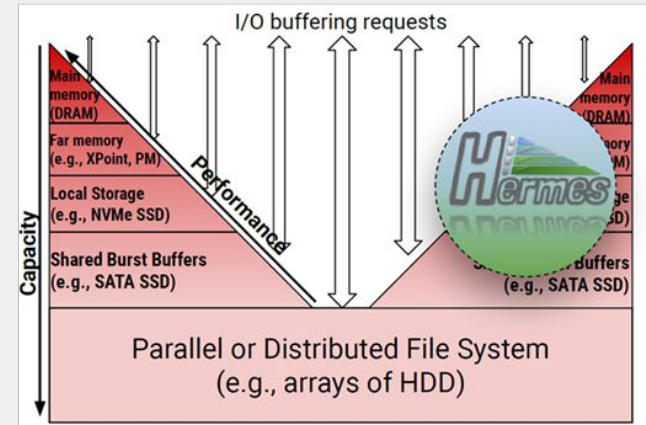
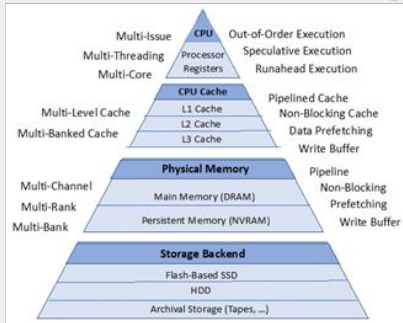
- Illinois Tech
- The HDF Group and UIUC, LBNL

## Motivation

- I/O is a performance bottleneck
- Memory/storage devices with different characteristics become available
- Current Bust Buffer system does not work well
- The 80-20 rule

## Solution

The Hermes System





Effective



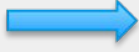
Fully parallel  
System optimization

Easy to use



It is a part of the HDF5 system

Automatic



User does not need to do any code change or system pre-setting  
Support heterogenous environments

Smart



Carry user information to all levels  
Hermes performance and management tools

Tested

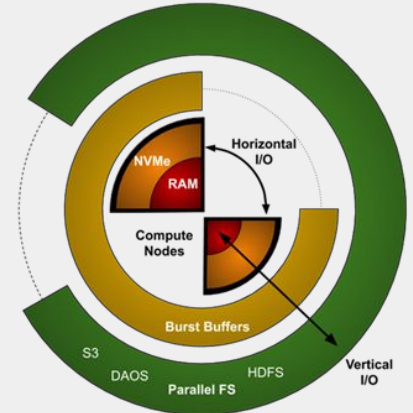
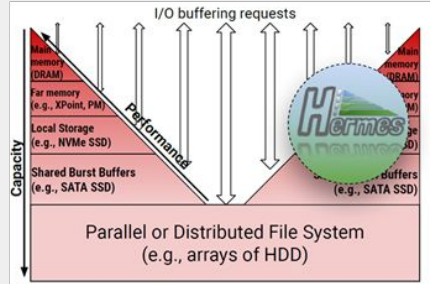


Tested on different applications and for different servers

Available



The Beta version is released





Short message from  
NSF OCI Program Director  
**Dr. Seung-Jong “Jay” Park**



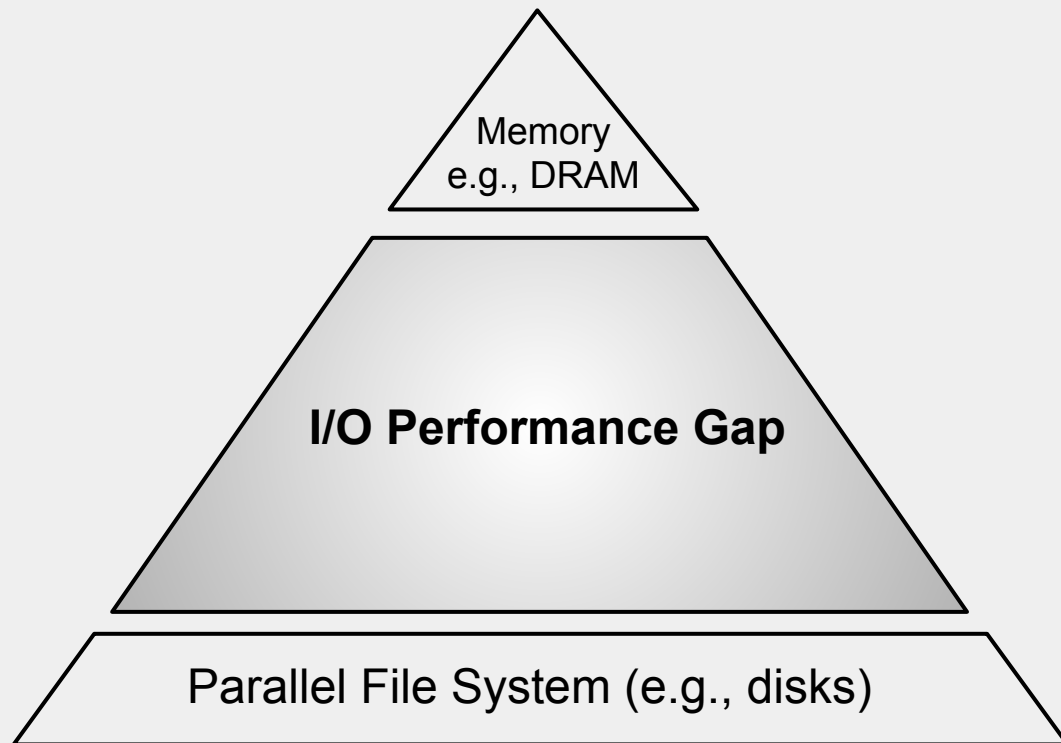


# Project Overview

# I/O Performance Gap

8

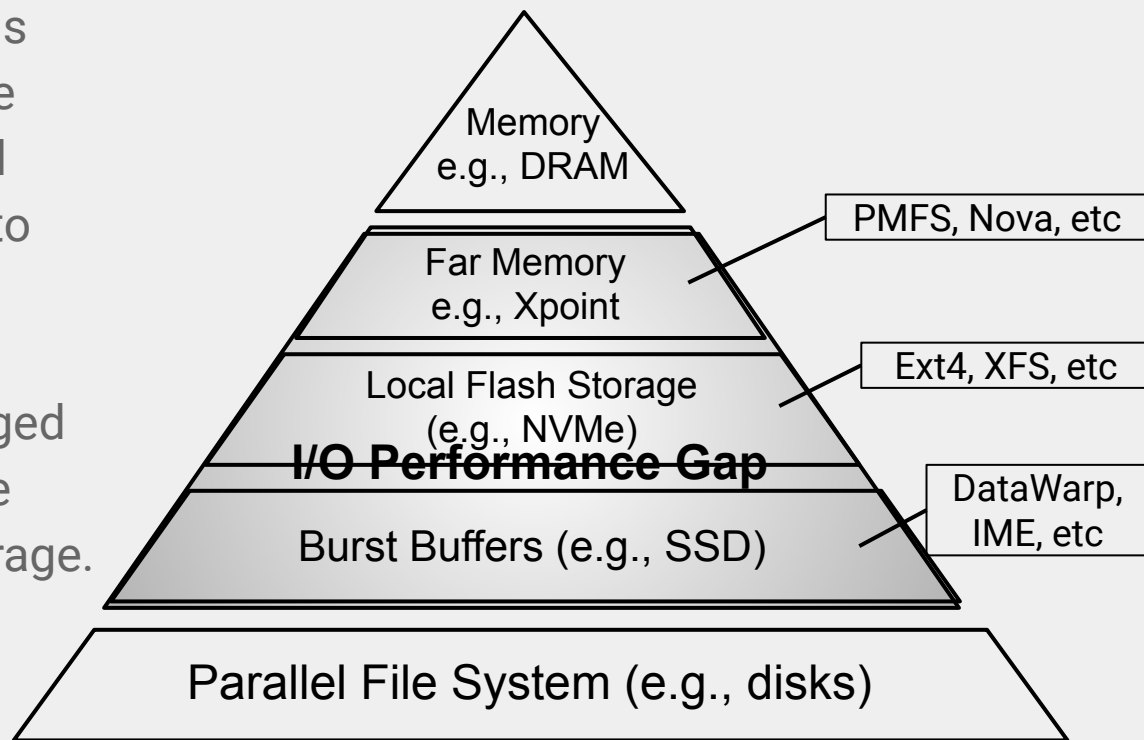
- Traditionally memory systems and storage demonstrate wildly different performance.
  - Access latency
  - Bandwidth
  - Data representation
- Applications experience performance degradation due to slow remote access to storage.





- Modern storage system designs include multiple tiers of storage organized in a **deep distributed storage hierarchy**. The goal is to mask the I/O gap.
- Each system is independently designed, deployed, and managed making very difficult to reap the benefits of the hierarchical storage.

Ideally, the presence of multiple tiers of storage should be **transparent** to applications without having to sacrifice **I/O performance**.

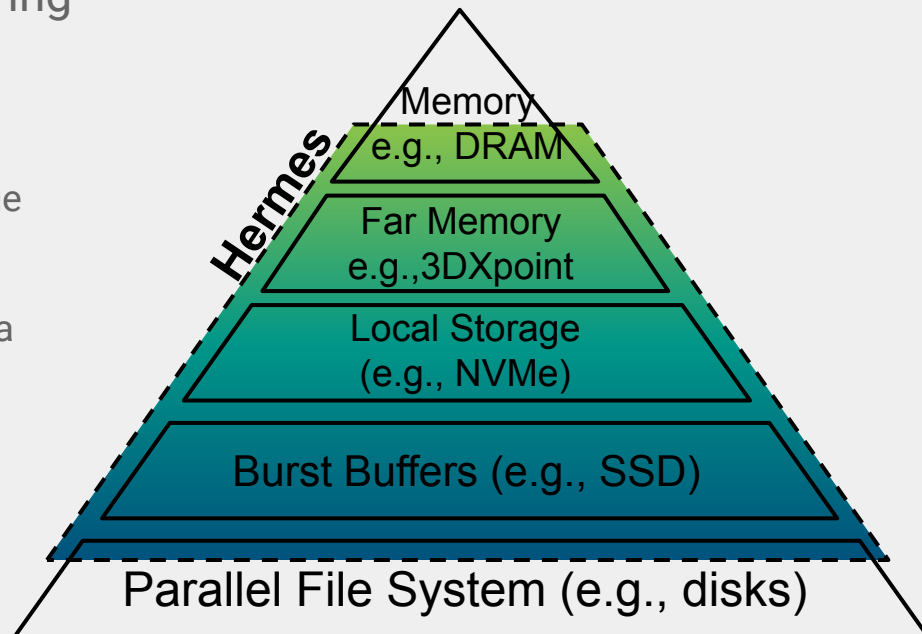




# Hermes Overview

10

- A new, multi-tiered, distributed buffering system that:
  - Enables, manages, and supervises I/O operations in the Deep Distributed Storage Hierarchy (DDSH).
  - Offers selective and dynamic layered data placement.
  - Is modular, extensible, and performance-oriented.
  - Supports a wide variety of applications (scientific, BigData, etc.,).

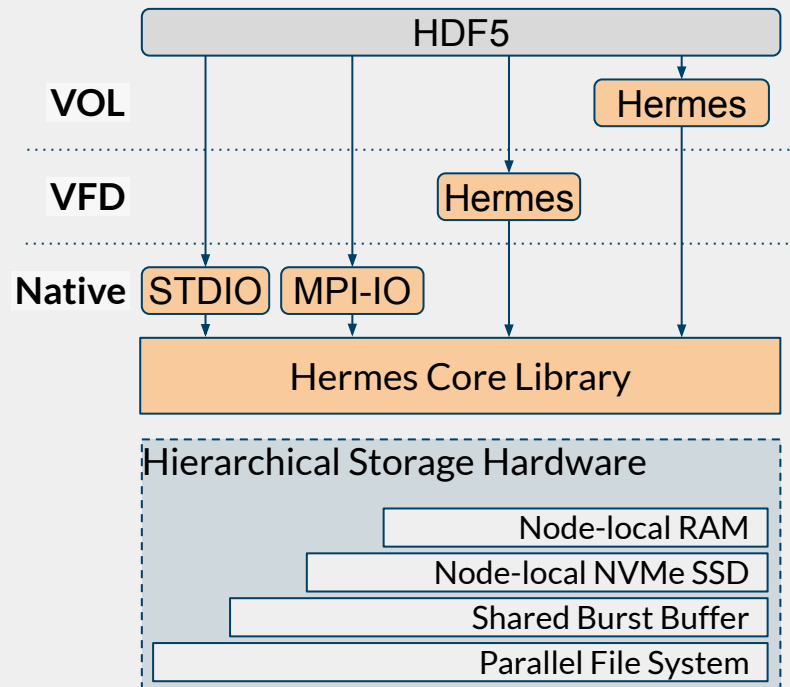


# Hermes Adapter Layer

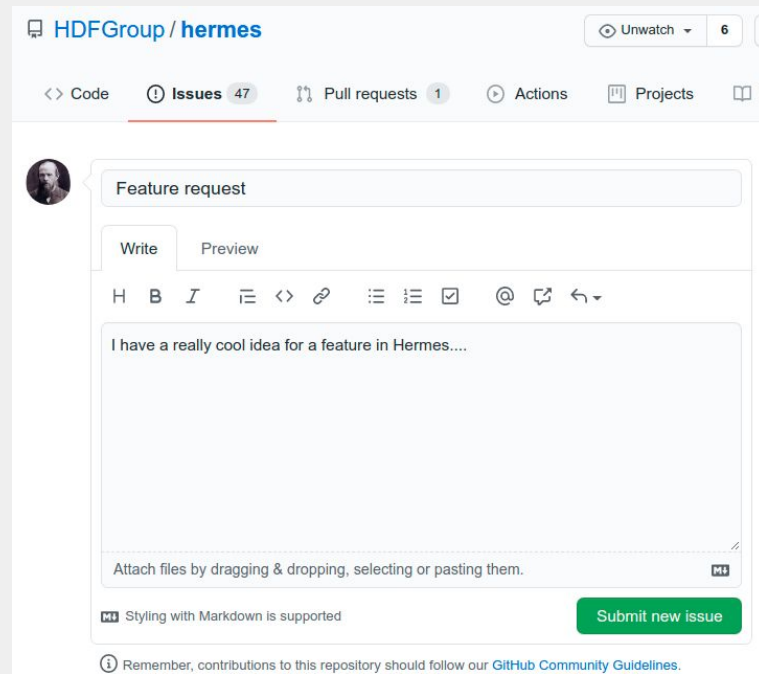
11

Applications can natively interact with Hermes using existing I/O Interfaces

- Standard Interceptors
  - STDIO
  - POSIX
  - MPI-IO
- HDF5 Level
  - Hermes VFD
  - Hermes VOL



- Github repo:  
<https://github.com/HDFGroup/hermes>
- Create an issue to submit feedback, use cases, or feature requests.
- Note: Hermes is still under active development.
- Join us in the Hermes Forum:
  - <https://forum.hdfgroup.org/c/hermes>










HDFGroup / hermes

<> Code Issues 47 Pull requests 1 Actions Projects V

Feature request

Write Preview

H B I  <>     @  

I have a really cool idea for a feature in Hermes....

Attach files by dragging & dropping, selecting or pasting them.

Styling with Markdown is supported

Submit new issue

Remember, contributions to this repository should follow our [GitHub Community Guidelines](#).



Demos

- Beta release
  - Feedback appreciated!
  - Monthly release schedule
  - [Wiki](#) overview
  - [Doxygen](#) overview
- Deployment and configuration guide
- Running IOR with Hermes



FAQ

## How does Hermes integrate into modern HPC environments?

- As of now, applications link to Hermes (re-compile or dynamic linking). We plan to integrate Hermes with the system scheduler so to provision buffering resources.
- Hermes can also be run as a service by spawning the Hermes daemons before the application starts. This way, workflows are supported where one job (typically a simulation) produces data and another job (typically analytics) consumes the data directly from the buffering system.
- Lastly, for HPC environments that support containers, we have bundled Hermes into Docker images. Singularity containers may also be supported if requested.



## What are Hermes build dependencies?

- A C++ compiler that supports C++ 17.
- [Thallium](#) - RPC library for HPC.
- [GLOG](#) - The Google logging library (v0.4.0).
- Google [ORTOOLS](#) for constraint optimization.
- MPI (tested with MPICH 3.3.2 and OpenMPI 4.0.3).
- The [Catch2](#) testing framework (only required if built with -DBUILD\_TESTING=ON).
- [Hermes README](#)

## Does running Hermes require any elevated privileges?

- No. Hermes does not require administrative privileges. Hermes is designed to run in user space as an application extension.

## How to get started with Hermes? Do I need to make changes to my application?

- Start with provided adapters, currently, POSIX, STDIO, MPI-IO
  - No code changes, just `LD_PRELOAD`
- Use the HDF5 Hermes VFD; required code change: a new HDF5 file access property or no code change: default HDF5 file access property
- Use Hermes API directly; mainly for new development
  - Benefits: maximum flexibility and greatest potential performance gains; the Hermes native API is simple and minimalistic.
    - Buckets and Blobs
    - Virtual Buckets and Traits

## How scalable is Hermes?

- On paper, Hermes should be as scalable as the RPC components of Mochi.
- Currently, our development testbed consists of only 64 nodes, but Exascale is the target, clearly.
- Scalability testing will be an important activity over the next year.
- Help us to explore more :)

## Which MPI implementations are supported?

- For the native Hermes API, we support OpenMPI, MPICH and Intel MPI.
- Currently, at the MPI-IO adapter supports only MPICH. The following updates to MPI\_Status happen in 3 places in mpiio.cc, and won't compile with OpenMPI.

```
mpi_status->count_hi_and_cancelled = 0;  
mpi_status->count_lo = ret;
```

- We intend to remove the specifics of the MPICH implementation and support more MPI implementations.

## What data placement strategies are available in Hermes? Can I bring my own?

- Data placement in the hierarchical space is configurable and extensible.
- We currently provide three strategies: Random, Round-Robin and MinimizeTime.
  - **Random**, each blob is distributed randomly between buffering targets favoring load balancing.
  - **Round-Robin**, each blob is placed to the next buffering target in order (mimics PFSs)
  - **Linear Optimization**, each blob is decomposed and distributed to a subset of the best available buffering targets. This is the default engine and is customizable through the set of constraints the algorithms receives. For example, the user can request targets with limited remaining capacity to not be considered for placement.
- We will provide an interface for user-defined data placement strategies.

## How do I prevent specific paths from being intercepted by an adapter?

- If you want to exclude a Hermes adapter from intercepting certain paths, you currently add the directory to `hermes::adapter::kPathExclusions` (in [interceptor.h](#)) and recompile.
- We will provide an environment variable to facilitate that.

## Are there examples of how to use the Hermes native API?

- Several examples are available [here](#).
- Start with [end\\_to\\_end\\_test](#).



## How much RAM should I set aside for Hermes buffering space?

- Configurable by the user.
- 10-20% of the available RAM is recommended.

## What is the RAM trade-off between applications and Hermes?

- More RAM for Hermes can lead to higher performance (e.g., data-intensive science).
- No RAM for Hermes means skipping the DRAM tier entirely (e.g., older servers)

## What is the RAM footprint of Hermes metadata management?

- Total bytes for metadata =  
280 +  
 $(192 * \text{num\_targets}) +$   
 $(288 * \text{max\_buckets\_per\_node}) +$   
 $(352 * \text{max\_vbuckets\_per\_node}) +$   
 $(64 * \text{max\_blobs\_per\_node}) +$   
 $(64 * \text{max\_blobs\_per\_node} * \text{avg\_buffers\_per\_blob})$
- We plan to implement a technique where unused metadata information is swapped into NVMe to relieve DRAM pressure.

## Is there any interference between Hermes and the OS page cache?

- There is interference.
- It's a long story.

## Is prefetching currently supported in Hermes?

- We are going to support it, just not in this release.
  - We have an initial implementation and test results, but more polishing is needed.
- The prefetching decision is based upon the importance score of a blob.
  - If a blob is deemed important by either the user or by the system (using collected statistics reflecting the access characteristics), the prefetcher will kick in and elevate the blob to upper tiers of the hierarchy.
    - Less important blobs will move in the opposite direction and start moving down in the hierarchy.

## How is Hermes configured?

- Currently, through a configuration file; better tooling will be forthcoming.
- Configuration examples can be found [here](#).
- Adjustments for system *and* application might be needed. See [this Hermes Wiki entry](#).



Multi-Tiered  
Distributed I/O  
Buffering System

Thank you.

Contact us

[akougkas@iit.edu](mailto:akougkas@iit.edu)

[gheber@hdfgroup.org](mailto:gheber@hdfgroup.org)

Learn more

[tinyurl.com/hermes-buffering](http://tinyurl.com/hermes-buffering)

<https://github.com/HDFGroup/hermes>

Feel free to join us on a roundtable  
conversation at 12:20pm CST

<https://meet.google.com/wfr-ivsr-zuz>

Or dial: (US) +1 234-985-0092

PIN: 717 562 972#

The  Group

ILLINOIS INSTITUTE  
OF TECHNOLOGY

