# ADMIRE
malleable data solutions for HPC

admire-eurohpc.eu

# Adaptive multi-tier intelligent data manager for Exascale

# Current progress in ad-hoc storage systems within ADMIRE project

Javier Garcia-Blas

uc3m | Universidad Carlos III de Madrid

Marc-André Vef

JOHANNES GUTENBERG UNIVERSITÄT MAINZ     JG|U

**HPC-IODC: HPC I/O in the Data Center Workshop**

EuroHPC
Joint Undertaking

GOBIERNO DE ESPAÑA     MINISTERIO DE CIENCIA E INNOVACIÓN

# The EuroHPC ADMIRE project

- € 7.9 million overall budget (2021-2024)

- 19 partners in 6 EU countries

- Key points:
  - Adaptive multi-tier data management
  - Computational and I/O malleability
  - Continuous monitoring
  - Focus on ad-hoc storage systems
  - Transparent data staging
  - Use cases from industry and academics
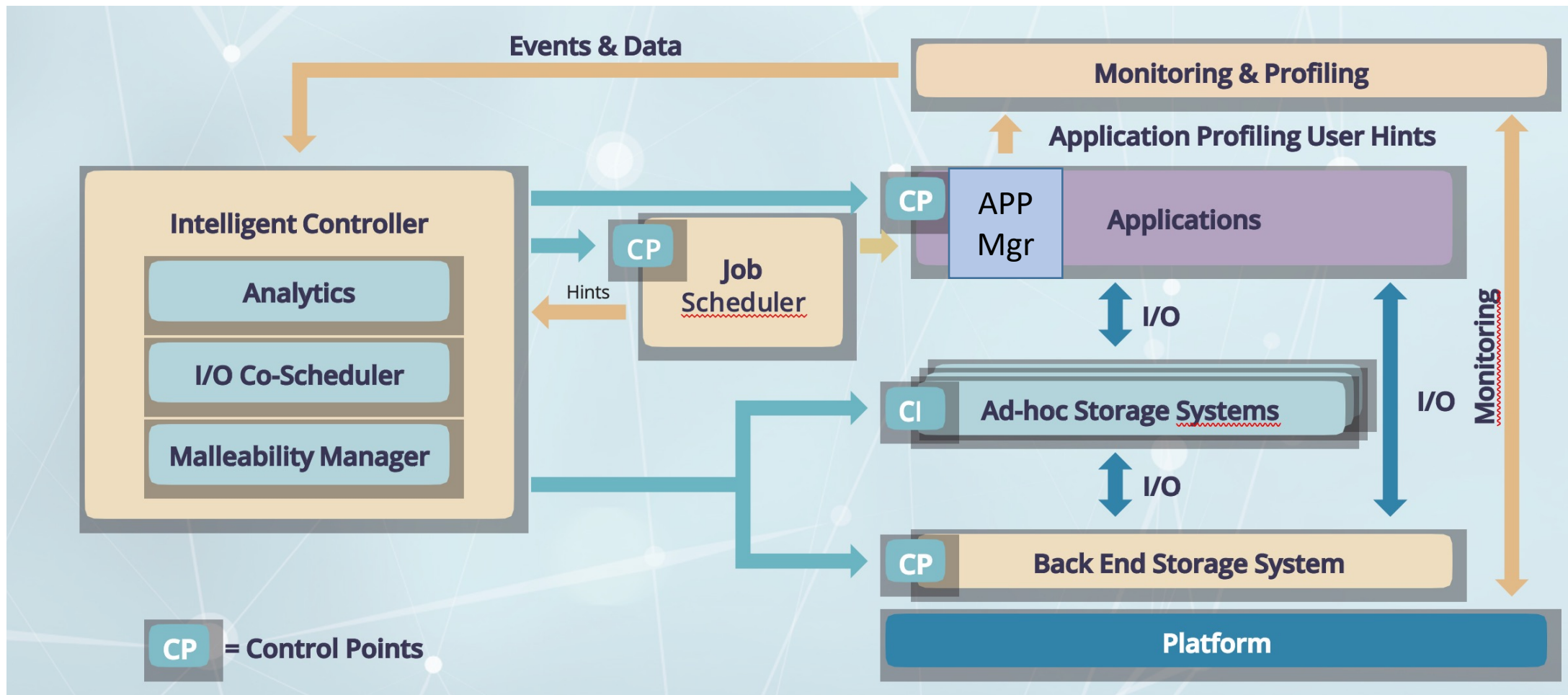


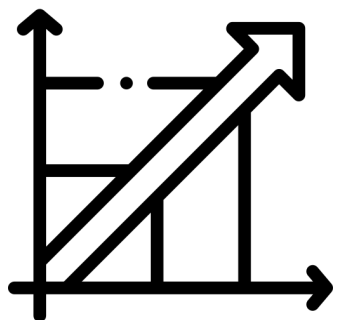ADMIRE project architecture @ https://admire-eurohpc.eu

# Motivation

- I/O-intensive HPC-based applications have been primarily based on distributed object-based file systems
  - **Separate data** from **metadata** management
  - Enable each client to **communicate in parallel** with multiple storage servers.
- Exascale I/O raises the throughput and storage capacity requirements by several orders of magnitude.
- Current challenges
  - Systems already developed for data analytics are not directly applicable to HPC due to the **fine-granularity** involved in scientific applications.
  - Semantic gap between the application requests and the way they are managed by the storage back-end at the block level.

# The framework

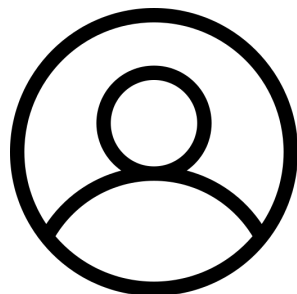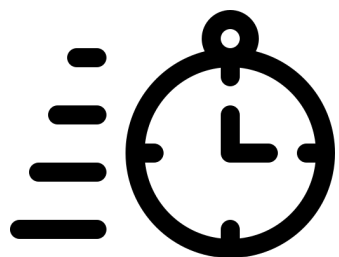# Adaptive I/O. Ad-hoc storage systems

- Deployed for an application or workflow
  - Implement temporary ad hoc storage systems with resilience features
  - Support fast node-local storage technologies (e.g., RAM, NVMe, SSD, ..)

- Four types of ad hoc storage systems
  - Distributed file system based on the **GekkoFS** burst buffer file system
  - Distributed object-oriented data store based on **dataClay**
  - **Hercules In-Memory Storage System** and
  - **Expand,** ad-hoc DFS based on MPI with external connectors**.**

- All ad-hoc storage systems support the same APIs
  - POSIX -> LD_PRELOAD
  - Custom I/O interfaces possible, but may require application changes

- General data staging approach in ADMIRE via I/O scheduler
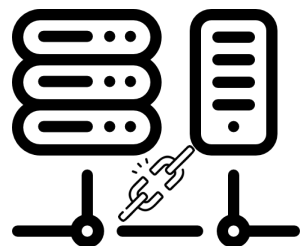  - Allows backend PFS independence

# GekkoFS



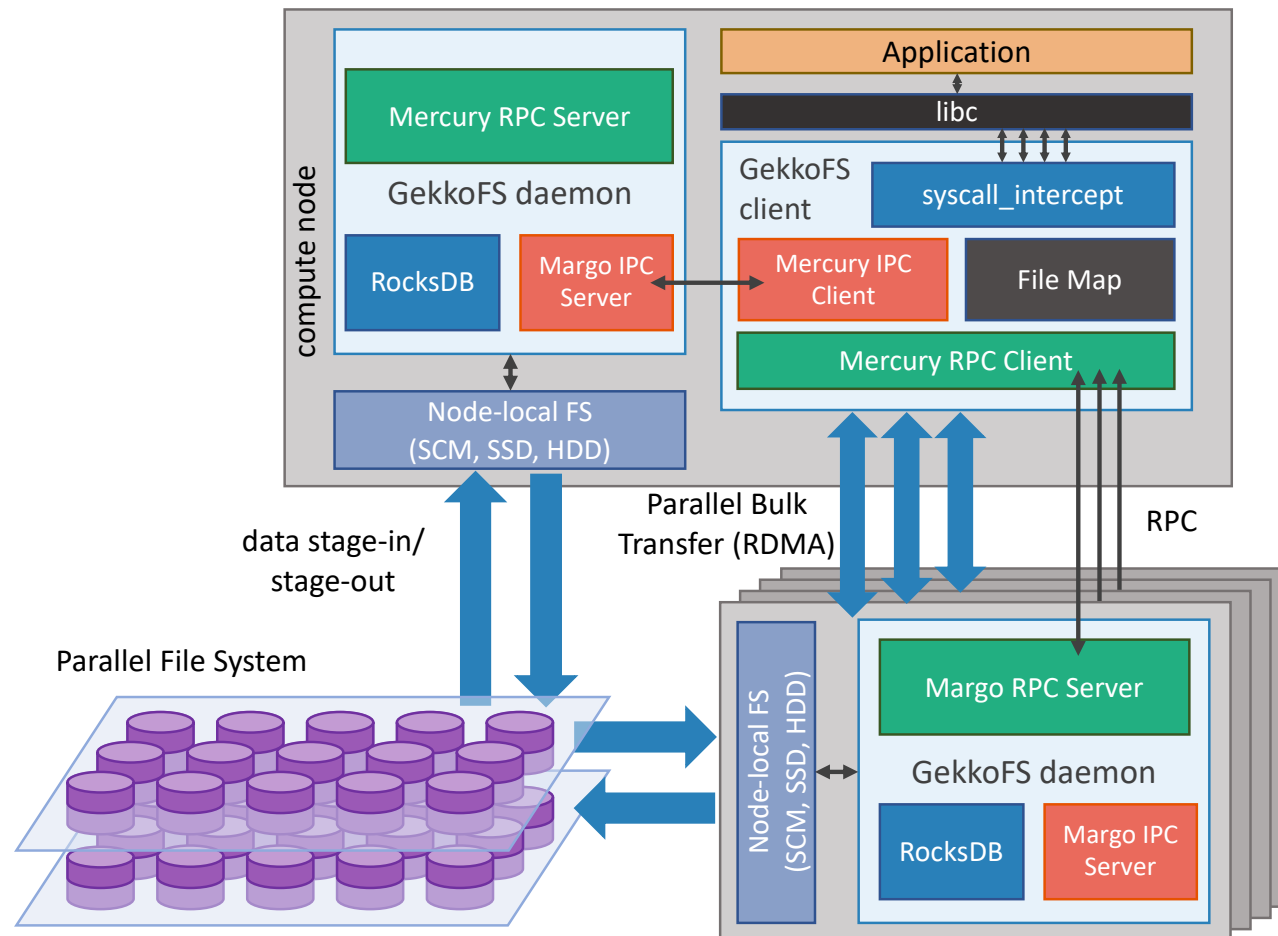**1. Scalability**

**3. User space**
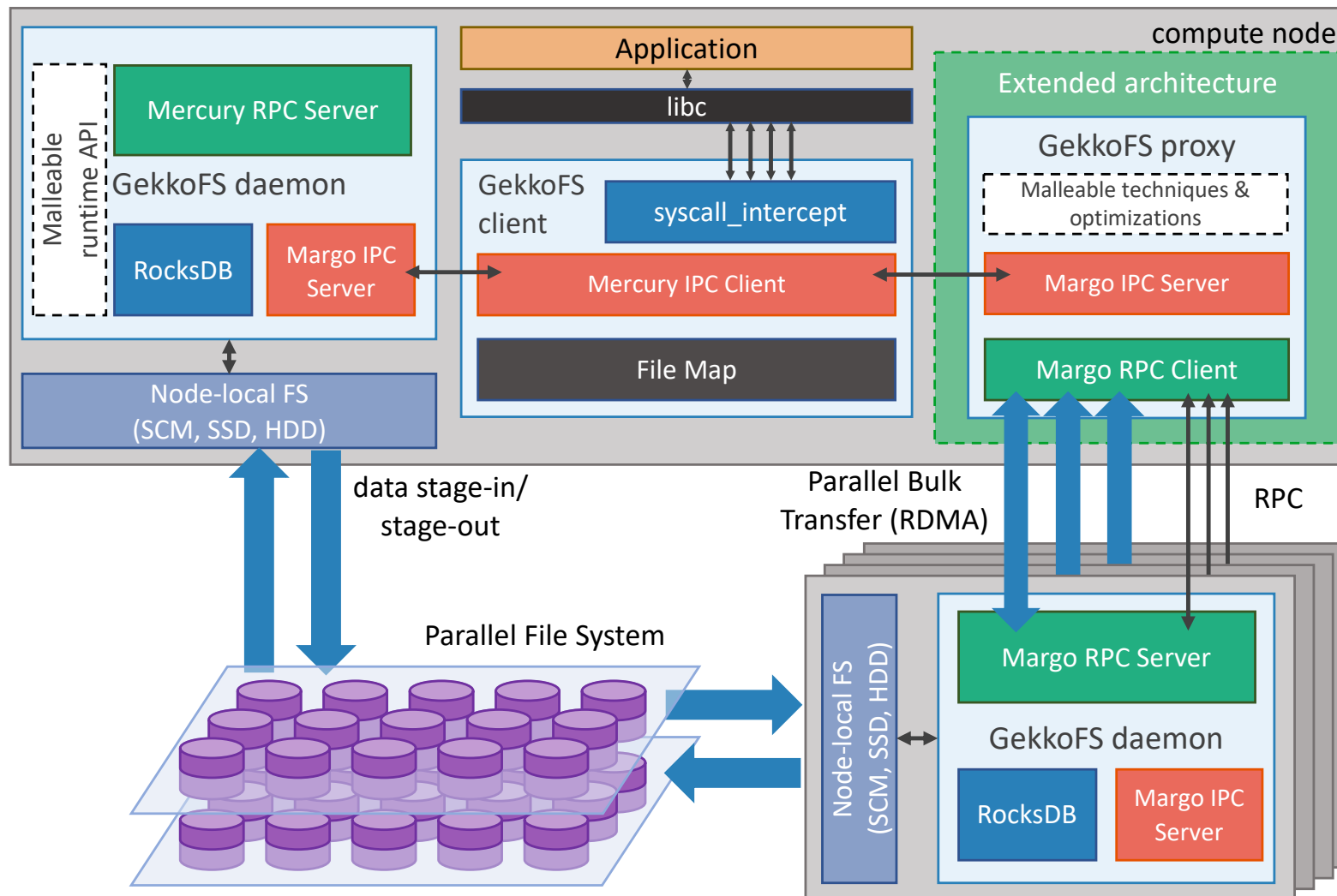
**2. Fast deployment**

**4. Hardware independence**

M.-A. Vef, N. Moti, T. Süß, M. Tacke, T. Tocci, R. Nou, A. Miranda, T. Cortes, A. Brinkmann. GekkoFS – A Temporary Burst Buffer File System for HPC Applications. In Journal of Computer Science and Technology (JCST), 2020

GekkoFS is open source:
https://storage.bsc.es/gitlab/hpc/gekkofs/

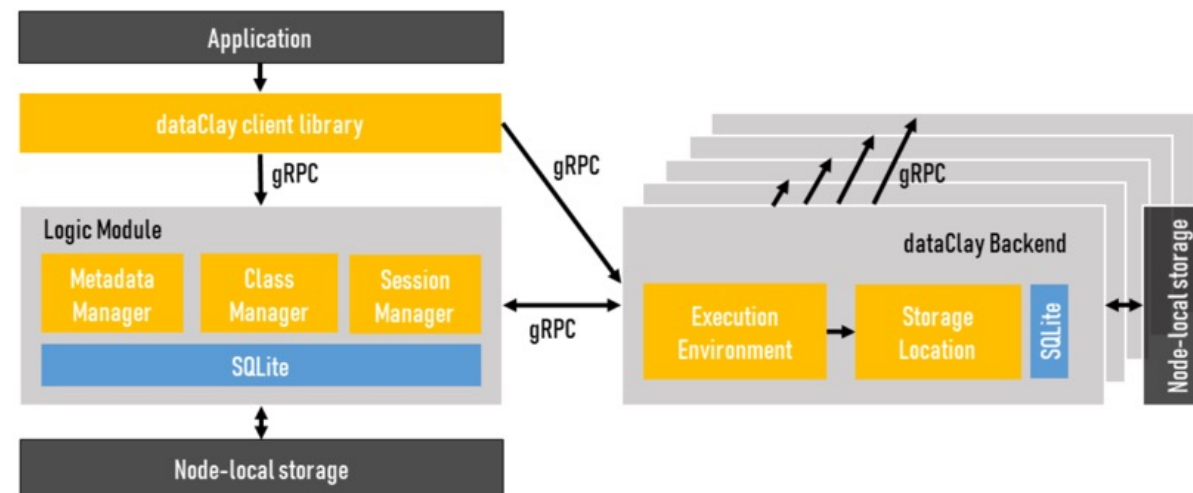# Challenges of LD_PRELOAD-based FS libraries

- Intercepts any function
  - Supporting all libc I/O functions is a daunting task
  - Other I/O interfaces, e.g., MPI I/O need special treatment
  - **Syscall_intercept helps to an extend**
- No valid mount point requires path checking
  - Detangle relative paths
  - **Separate file descriptor management necessary**
- Threads and forks don't mix
  - I/O library has no control over application or 3rd party library threads
  - `fork()` can lead to dead-locks
  - **Minimal interception library required**
- Malleability techniques limited
  - Any technique is confined to one process that can stop at any time
  - **Node-wide client solution necessary**

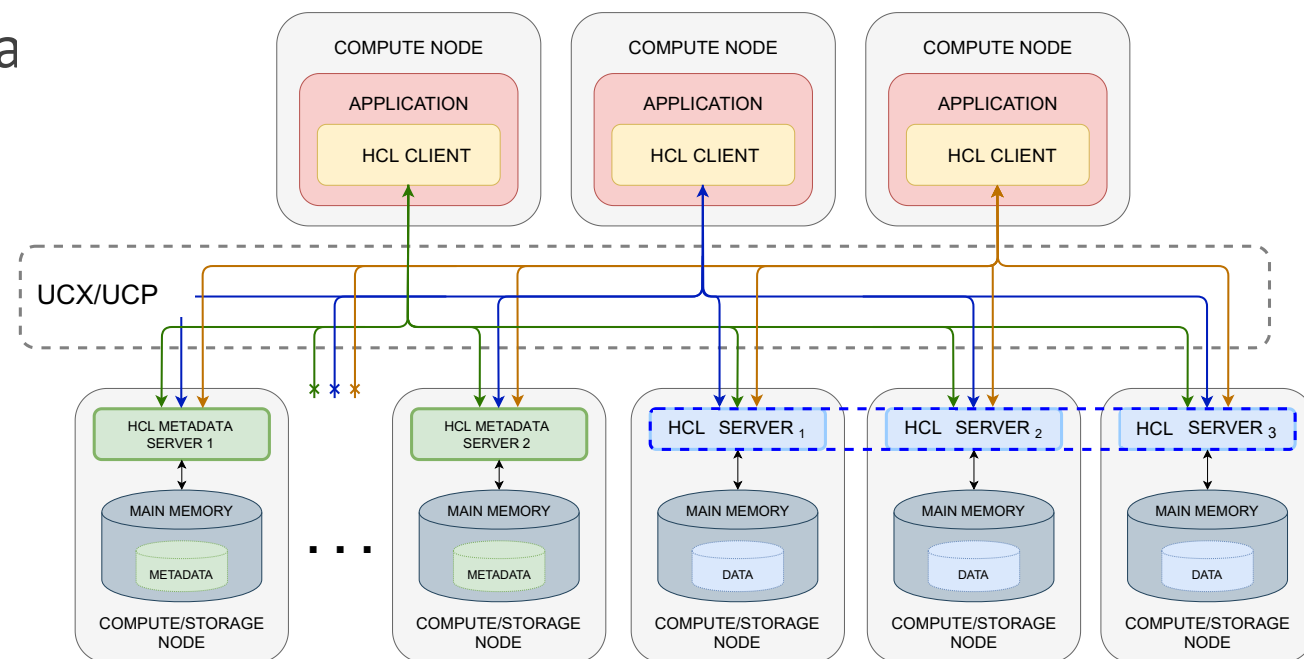HPC-IODC: HPC I/O in the Data Center Workshop. 2023.

# dataClay

- An active object store for HPC, Big Data, and Edge-to-Cloud applications

- Runs on heterogeneous infrastructures

- Brings computation to the data, enabling execution of custom user code
  - Objects = data + methods
  - Reduces communication costs and improves application performance

- Manages objects in-memory during execution
  - Avoids transformations and reduces disk accesses

- APIs
  - Object-oriented: transparent persistence (Python and Java bindings)
  - Object store: Get/Put operations

- In-memory IO accelerator for volatile data

- Use cases:
  - HPC workflows.
  - Checkpointing engines.
  - Distributed locking.

- POSIX-based support for portability.

- Main features
  - Supports attached/detached storage servers.
  - Supports several data distribution policies.
  - Efficient storage backend based on GLIB.
  - Write on close policy.
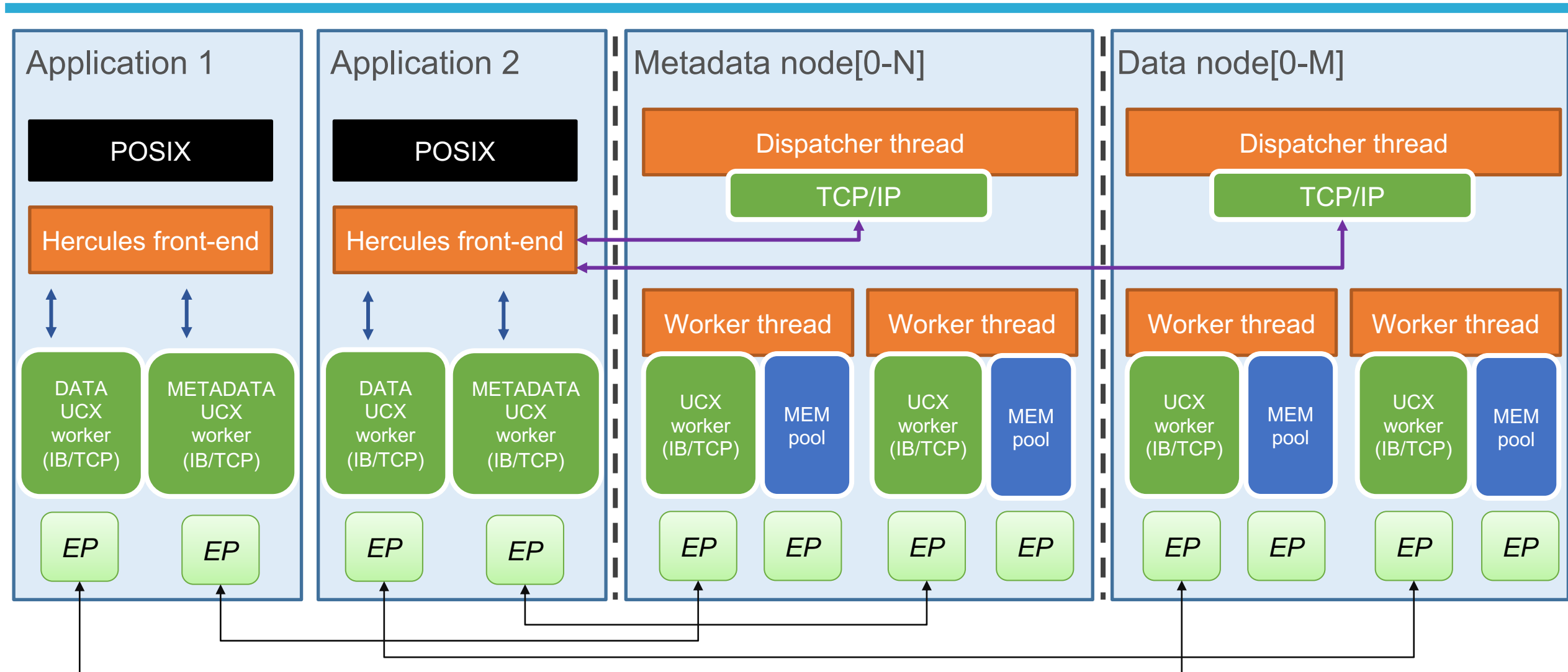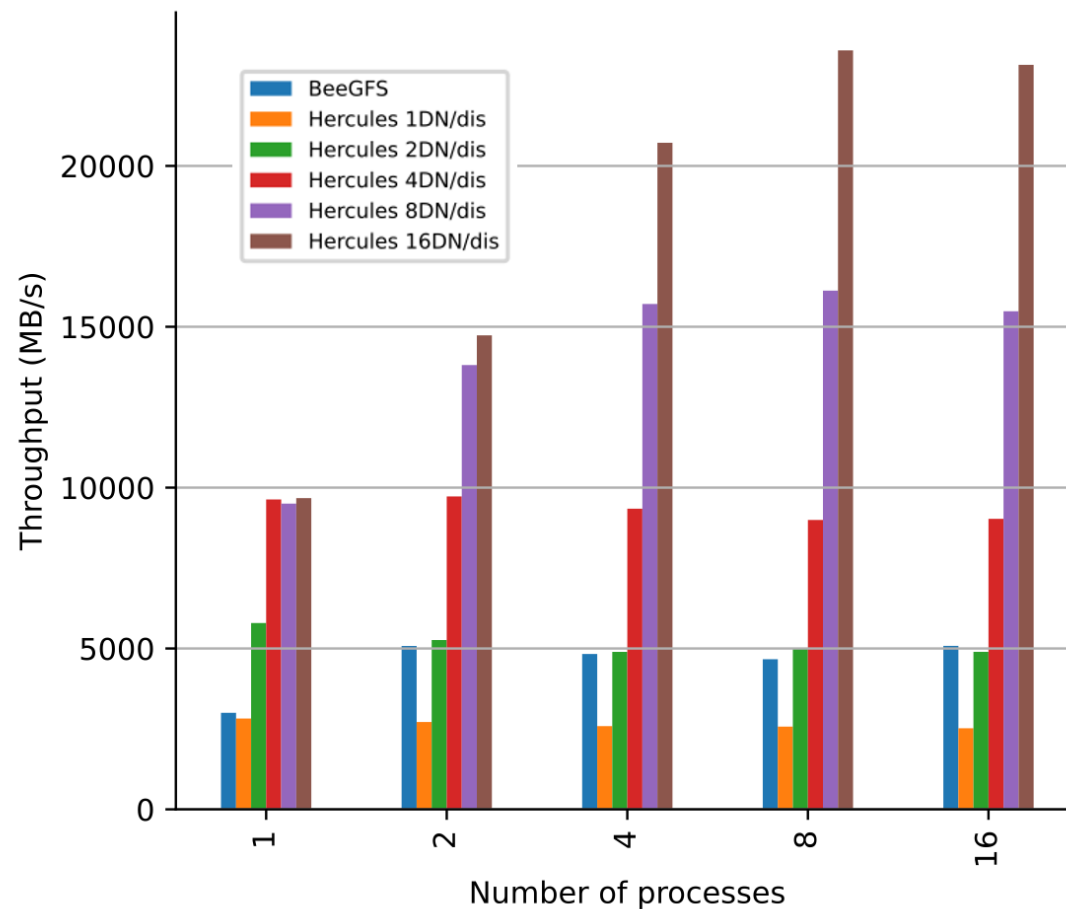  - Integrated with SLURM.
  - Data replication.



*Javier Garcia-Blas, Genaro Sanchez-Gallegos, Cosmin Petre, Alberto Riccardo Martinelli, Marco Aldinuchi and Jesus Carretero. Hercules: scalable and network portable in-memory ad-hoc file system for data-centric and high-performance applications. EuroPAR conference. 2023.*
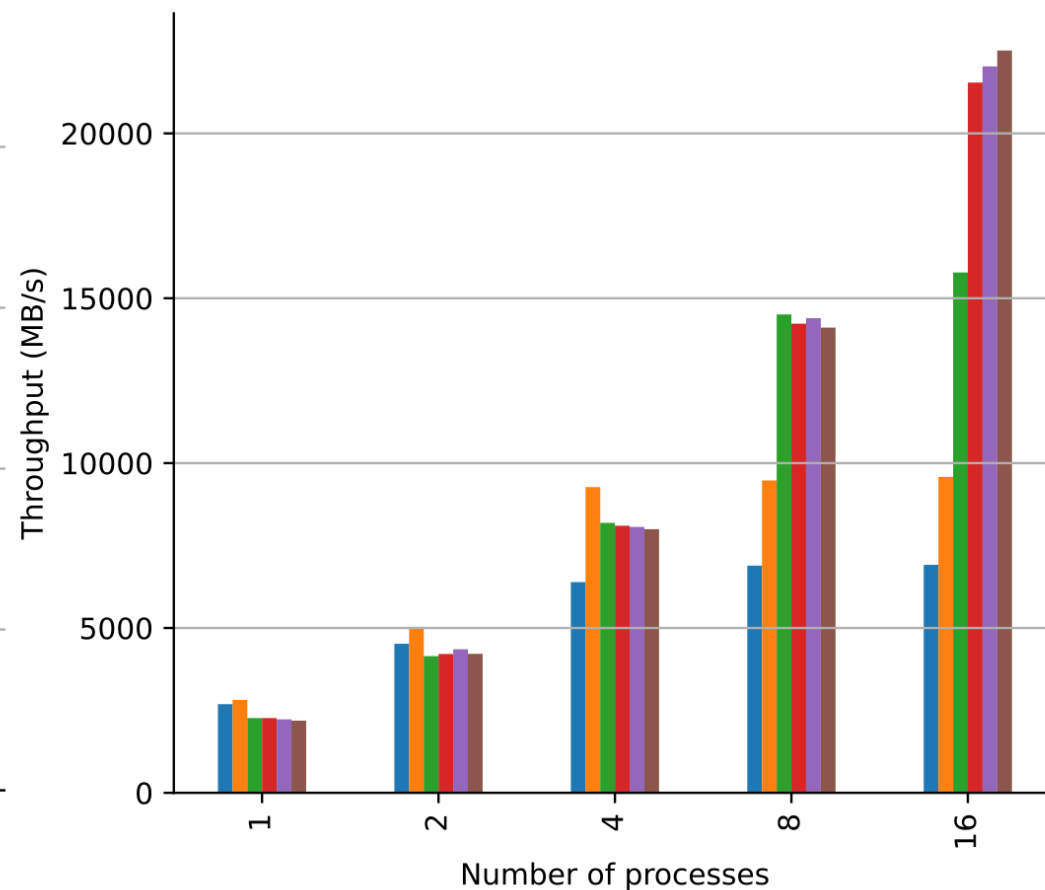
https://gitlab.arcos.inf.uc3m.es/admire/imss

# Portable communication layer

❑ Based on UCX

  o Unified Communication - X Framework (UCX) is a new acceleration library, integrated into the OpenMPI, OpenSHMEM, Dask, Charm++, …

❑ Key aspects

  o Multiple network interfaces/protocols available (TCP/IP, Omnipath, Infiniband supported).

  o Zero-copy message transfers of large data packages (>= 1 Mbytes).

  o Eliminated internal copies from application to network layer.

  o Asynchronous communication between peers.

  o RDMA QoS isolation.

  o End-point/two-sided-based communication.

# UCX communication layer developed

❏ Non-blocking/tag-based communication (MPI style)

➢ **Asynchronous** communication between peers.

❏ Low-level communication schema (in contrast to Margo RPC)

❏ Frontend

➢ Data and metadata UCX's workers enables **communication overlap**.

➢ QoS

❖ Interfaces and protocols can be enabled/disabled to adapt **network requirements**.

❖ Communication can be upgraded/downgraded (Infiniband/Omnipath to TCP).

➢ Communication parameters configured by using environment variables/config file.

❏ Backend

➢ One single listener per worker thread.

➢ Stores a pool of active end-points (two-sided communication).

HPC-IODC: HPC I/O in the Data Center Workshop. 2023.

Write Throughput

Read Throughput

# Expand Parallel File System
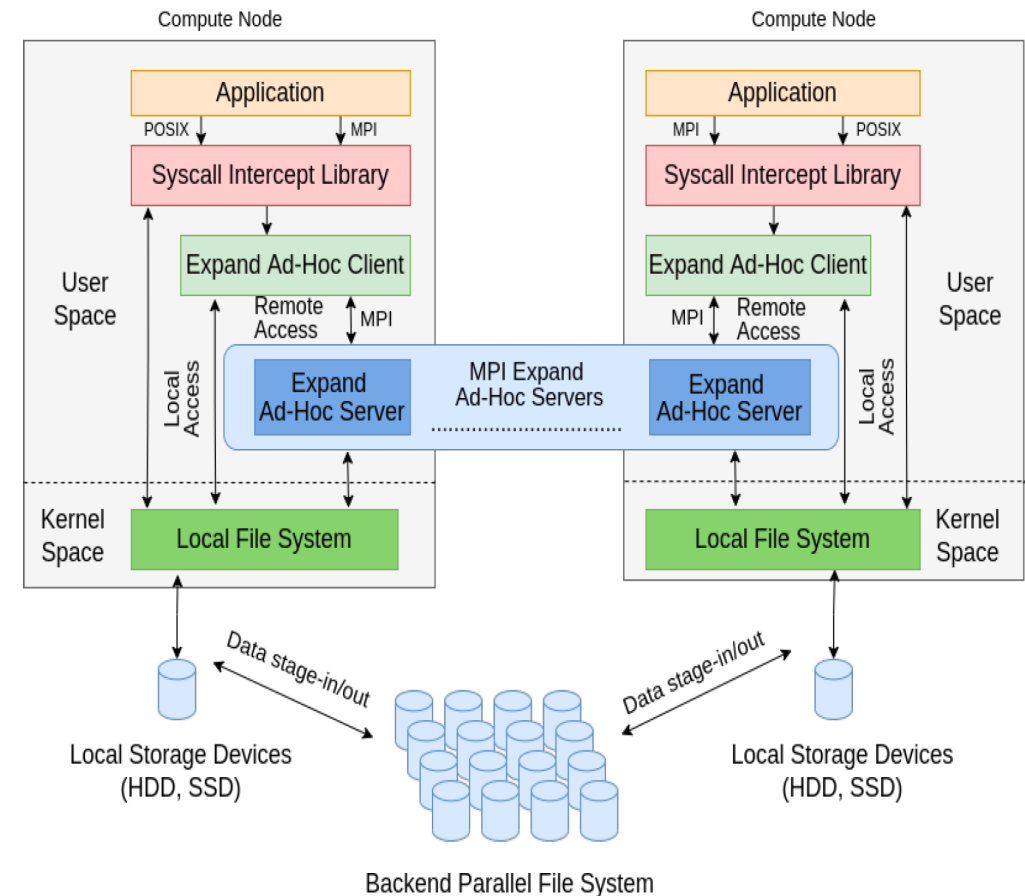
- **Expand** ad-hoc file system
  - Parallel, distributed file system based on the use of existing standard storage servers.
  - MPI-based
  - Distribution based on subfile/server and Round Robin
  - Fully distributed metadata (B0 of each file) by hash
  - Aggressive parallel I/O ops
  - Client-servers may be co-located

*F. García, A. Calderón, J. Carretero, J. Fernández, J. M. Pérez. The Design of the Expand Parallel File System. International Journal of High Performance Computing Applications. Vol. 17. Nº 1, Spring 2003. Pags. 21-37*

https://arcos-xpn.github.io
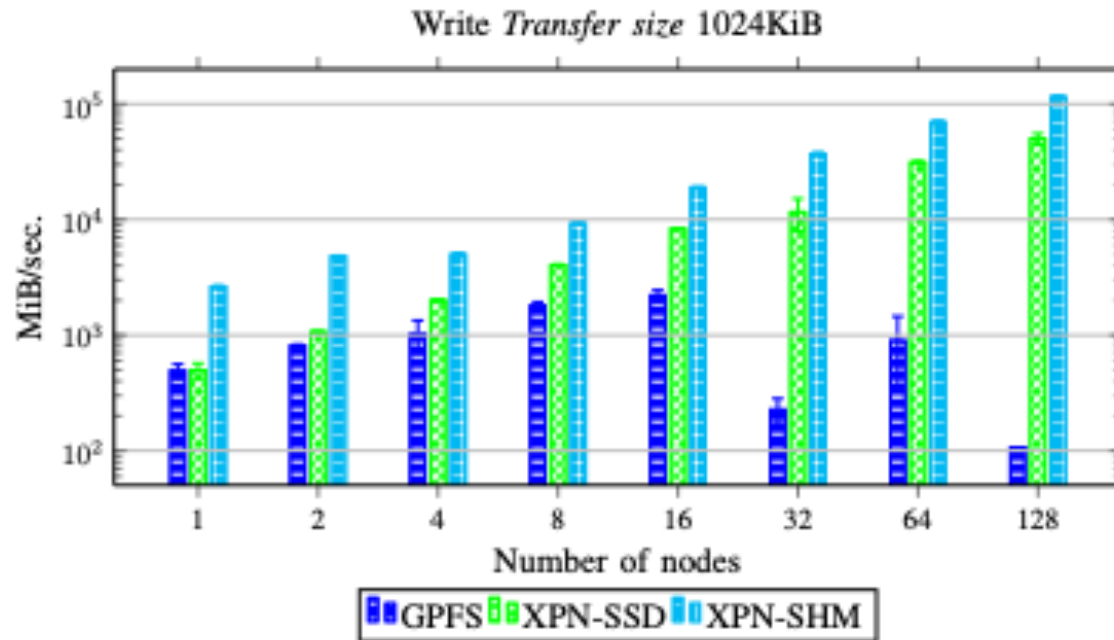
❑ IOR in MareNostrum



Fig. 4. GPFS vs. Expand Ad-Hoc. Bandwidth (MiB/sec.) writing with different transfer sizes (64KiB, 512KiB, and 1MiB), compute nodes (1, 2, 4, 8, 16, 32, 64, and 128), with 8 client process per node and shared file. Results in logarithmic scale
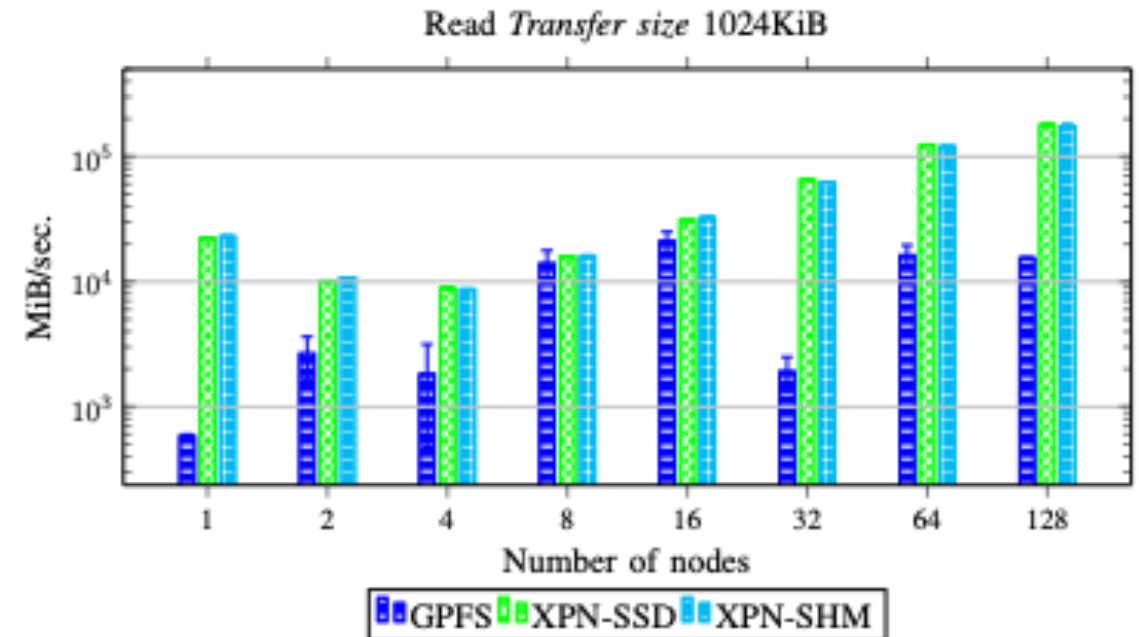
Fig. 5. GPFS vs. Expand Ad-Hoc. Bandwidth (MiB/sec.) reading data with different transfer sizes (64KiB, 512KiB, and 1MiB), compute nodes (1, 2, 4, 8, 16, 32, 64, and 128), with 8 client process per node and shared file. Results in logarithmic scale

# Current efforts

❑ Malleability

➢ Expanding/shrinking the number of data nodes.

❑ Monitoring

➢ Integration with existing monitoring tools.

❑ Replication / Mirroring

❑ Error correction schemes

➢ ECC/parity.

➢ To implement distributed parity calculations on advanced data placement schemes.