# ExaHDF5: Delivering Efficient Parallel I/O on Exascale Computing Systems
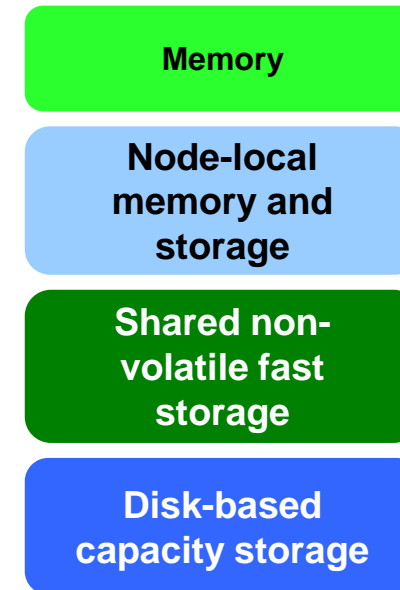
Suren Byna (Lawrence Berkeley Lab - LBL)

Scot Breitenfeld (The HDF Group - THG), Bin Dong (LBL), Quincey Koziol (LBL), Elena Pourmal (THG), Dana Robinson (THG), Jerome Soumagne (THG), Houjun Tang (LBL), Venkatram Vishwanath (Argonne National Lab), and Richard Warren (THG)

# Exascale I/O architectures and software

- Exascale storage hardware
  - Deepening hierarchy with:
    - Fast node-local storage and storage-class memory
    - Shared SSD-based storage layer
    - Disk-based capacity storage

- I/O software
  - High-level self-describing I/O libraries (HDF5, etc.)
  - Middleware (MPI-IO) and optimization layers
  - File systems (Lustre and Spectrum Scale/GPFS) and object storage (Intel DAOS)

- Challenges
  - Heterogeneity of storage devices and distributed across nodes
  - Disparity of I/O software stack (different tuning parameters)
  - Overheads of managing metadata in self-describing formats
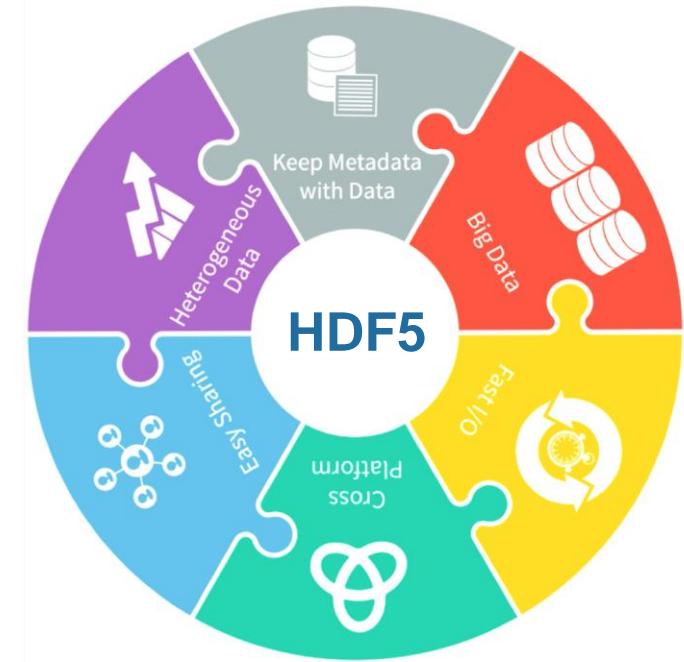  - Obtaining sustained I/O performance on exascale storage

**Exascale Storage Hardware**

| Memory |
|---|
| Node-local memory and storage |
| Shared non-volatile fast storage |
| Disk-based capacity storage |

**Exascale I/O Software**

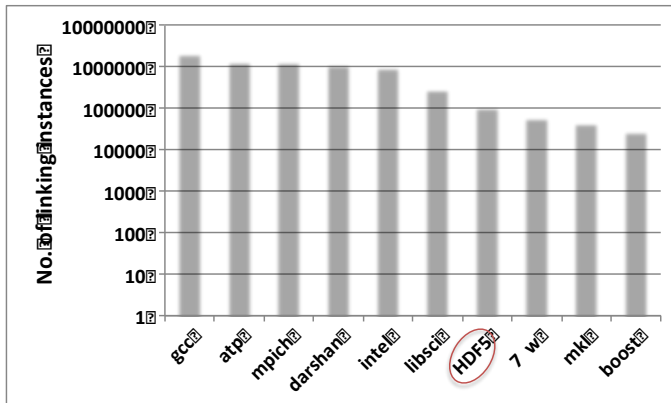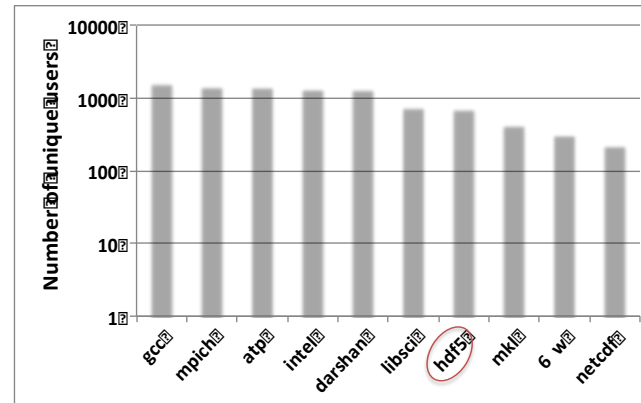| High-level lib (HDF5, etc.) |
|---|
| IO middleware (MPI-IO) |
| IO forwarding |
| Parallel file systems |

# ExaIO Project Products – HDF5

- HDF5 is designed to organize, store, discover, access, analyze, share, and preserve diverse, complex data in continuously evolving heterogeneous computing and storage environments.
  - Maintained by The HDF Group (THG)

- NASA/NOAA satellite data (Aura, JPSS-1, etc.)
  - Highest Technology Readiness Level (TRL 9) - "Flight proven" through successful mission operations

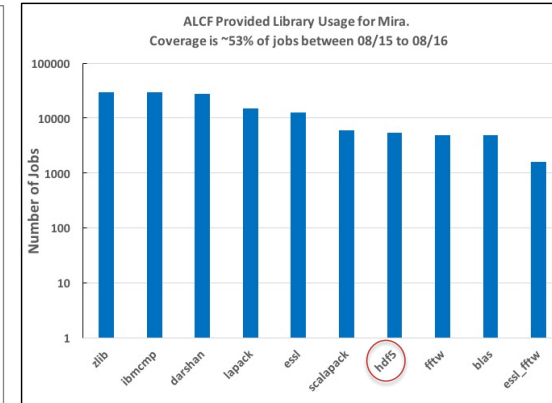- Heavily used on DOE supercomputing systems
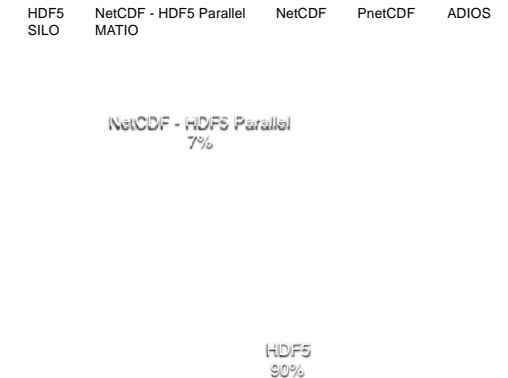
**HDF5: 2002 R&D 100 Award Winner**

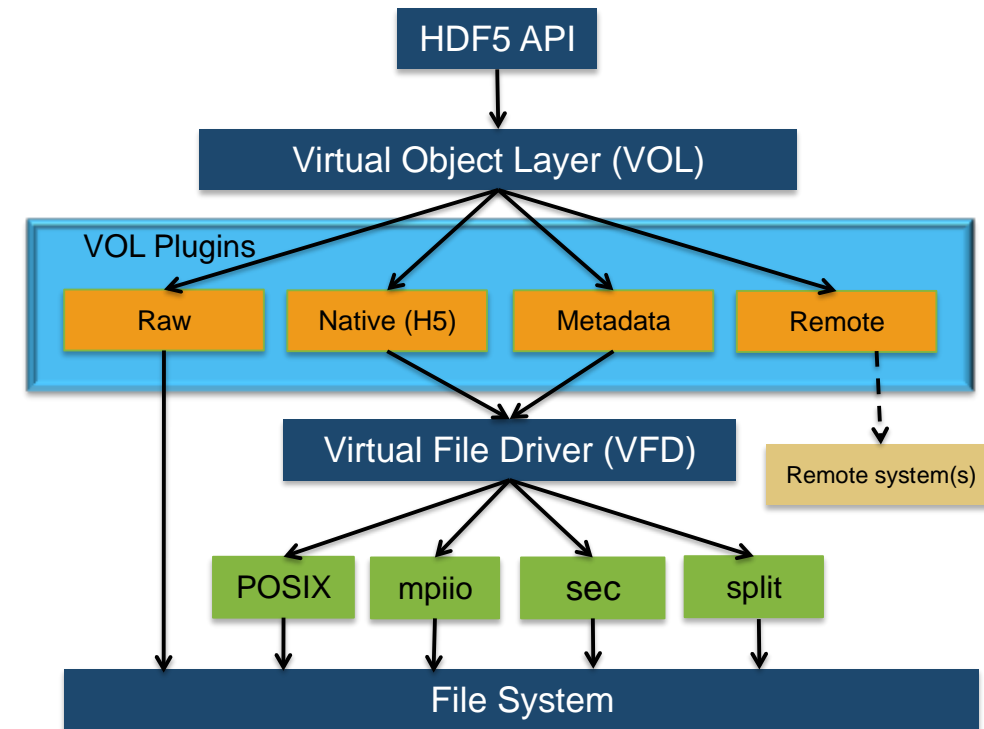a. Number of linking instances on Edison (NERSC)

b. Number of unique users on Edison (NERSC)

c. Number of linking instances on Mira (ALCF)

# Virtual Object Layer (VOL)

- Virtual Object Layer (VOL) provides an application with the HDF5 data model and API, but allow different underlying storage mechanisms

- Enables developers to use HDF5 on novel current and future storage systems easily
  - Prototype VOL connectors for using burst buffer storage transparently and for accessing DAOS are available
  - Developed VOL connectors for reading PnetCDF and ADIOS-BP data

- Integrated into the HDF5 trunk (will be released in 1.10.12 later this year)
  https://bitbucket.hdfgroup.org/projects/HDFFV/repos/hdf5/

- Allows ADIOS and other libraries to use HDF5 API

# HDF5 Data Elevator

- Data Elevator VOL connector
  - Transparent data movement in storage hierarchy - writes and reads
  - Intercepts file opens, write, read, and close function calls and places data in burst buffers temporarily; DE moves data asynchronously
  - Prefetches predicted chunks of data to burst buffer or memory
  - In situ data analysis capability using burst buffers
  - Phase 2 plan includes extending capabilities of Data Elevator for node-local storage



Memory
Node-local storage
Shared burst buffer
Parallel file system

Simulatioin processes    TEDM processes
HDF5/Others API    MPI IO API    Computing Node
IOCI
Redirected I/O    Async Data Movement
Append
Burst Buffer
f.h5.temp
DEMT
f.h5, f.h5.temp, ...
f.h5    PFS

Data Analysis Applications
Parallel Read
Insert, Update, Query , ...
Metadata Table

ARCHIE
-- metadata manager
-- consistency manager
-- prediction algorithm
-- parallel reader/writer
-- garbage collection
-- fault tolerance manager
-- ... ...

Hits
Miss
...
SSD File System
Cached chunks
...
Disk File System
Array to be analyzed
Parallel chunk prefetching

Computing    Writing Data    Moving Data from BB to PFS

VPIC + Lustre

VPIC + Burst Buffer DataWarp

VPIC + Burst Buffer DataWarp API

VPIC + Burst Data Elevator

VPIC multi-timestep

Time (s)
0    100    200    300

Stage (only for DataWarp)    Read #1    Read #2
Read #3    Read #4    Read #5
Read #6    Read #7    Read #8

Time (s)
80
60
40
20
0

0.8X
3.4X
5.8X

Lustre(Disk)    DataWarp(Disk+SSD)    ARCHIE(Disk)    ARCHIE(Disk+SSD)

https://bitbucket.hdfgroup.org/projects/HDF5VOL/repos/dataelevator/

# Asynchronous I/O with HDF5

- Asynchronous I/O allows an application to overlap I/O with other operations
- The asynchronous I/O feature has been implemented as a VOL (Virtual Object Layer) connector, <u>without</u> requiring major change the HDF5 library



Write VPIC data (256MB per process timestep, 5 timesteps)



Read VPIC data (256MB per process timestep, 5 timesteps)

https://bitbucket.hdfgroup.org/projects/HDF5VOL/repos/async/