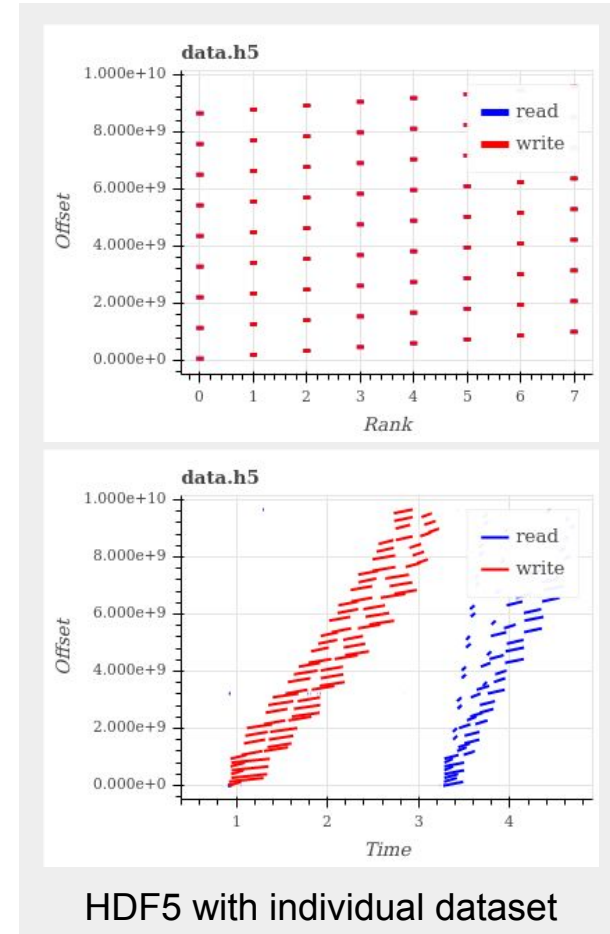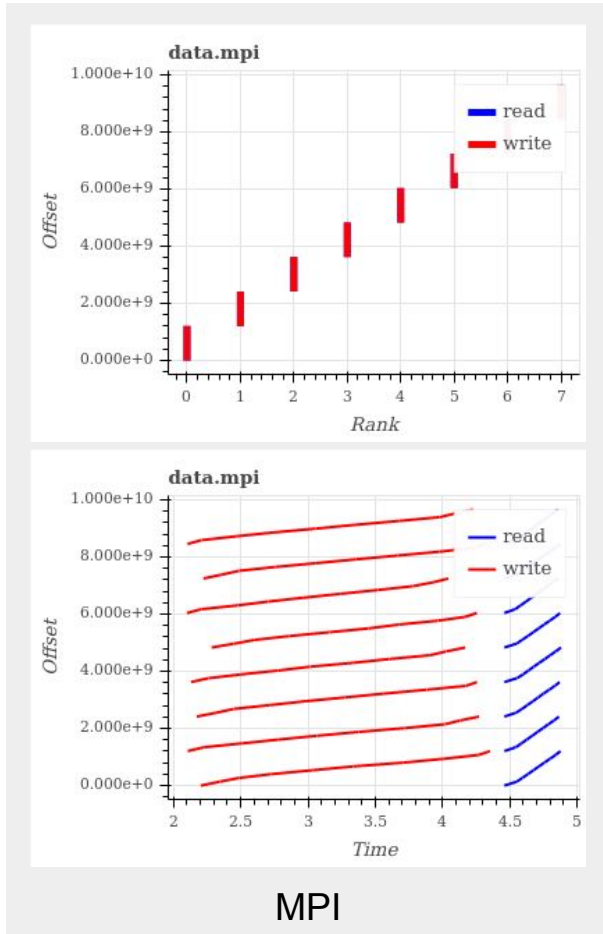# An I/O Study of ECP Applications

Setup

8 Processes: 4 nodes x 2 ranks/node
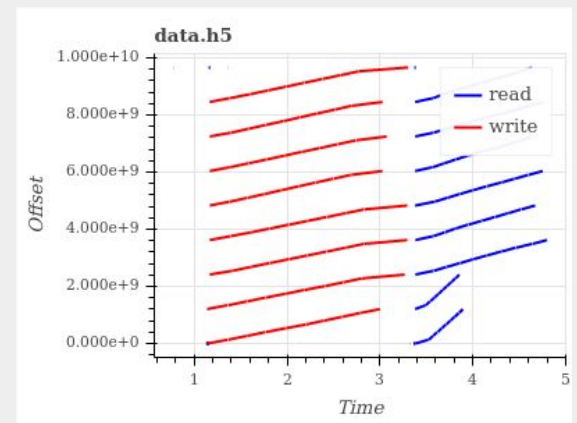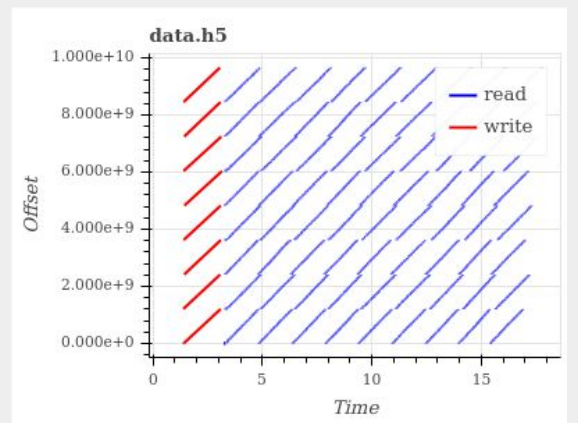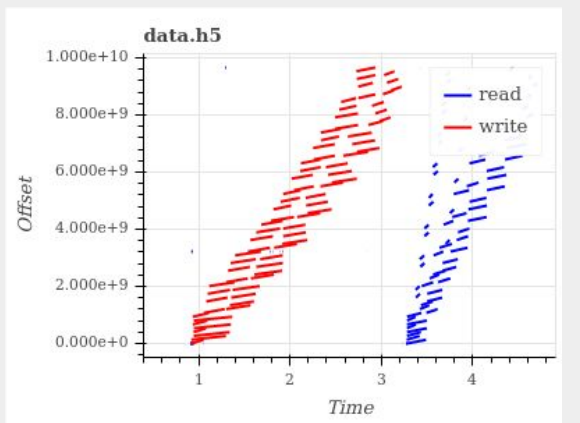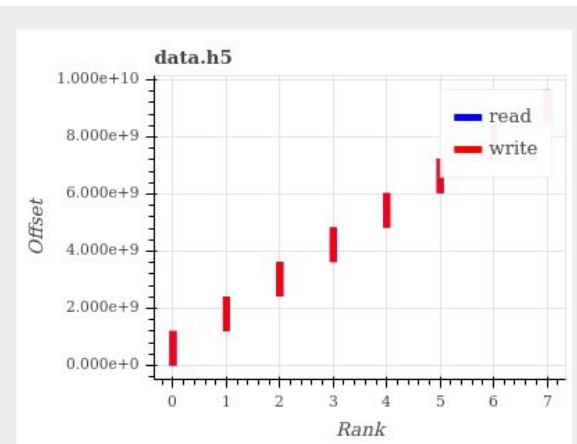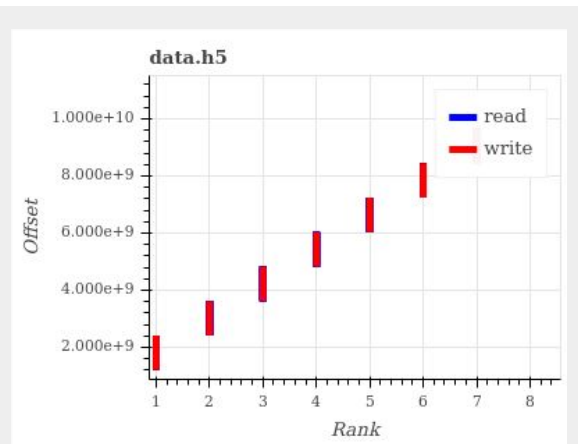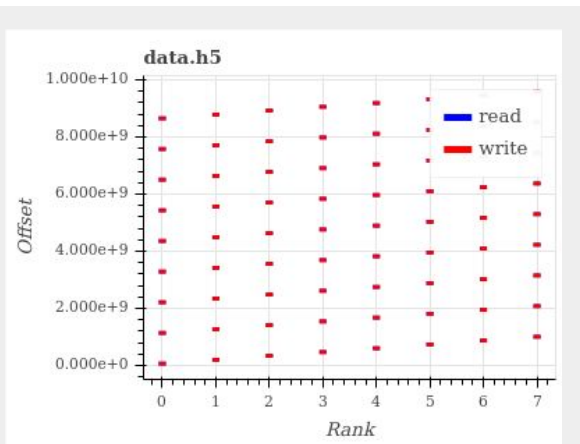
1GB per variable, 9GB in total

1. Individual dataset (9 dataset, one for each variable)
2. Compound data type
3. With new HDF5 API - AllocaMulti

# HACC-IO: MPI vs HDF5, why HDF5 is slow



MPI

HDF5 with individual dataset

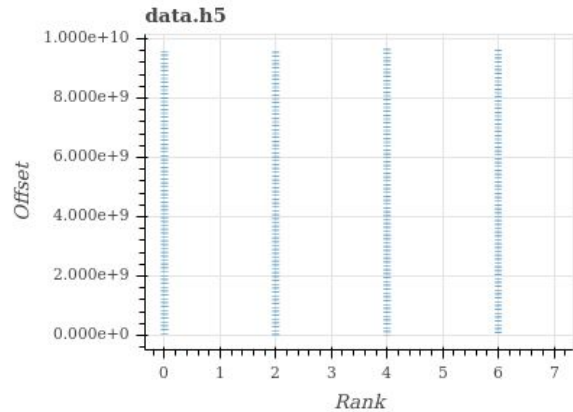# HACC-IO: HDF5 using different I/O methods



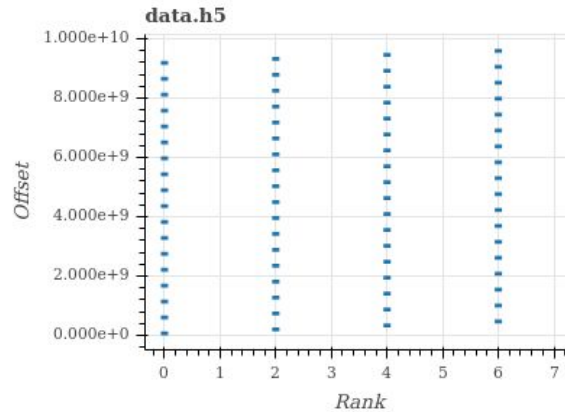Individual Dataset

Compound Datatype

Alloca Multi

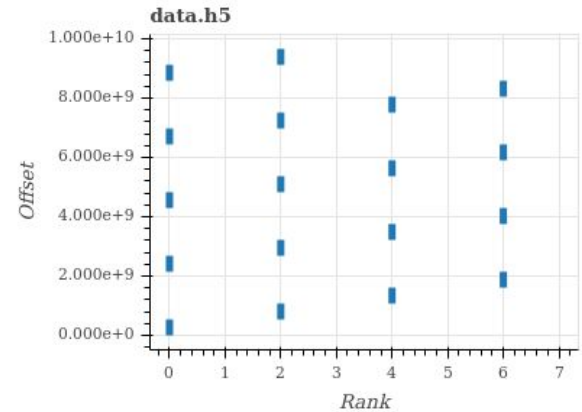# HACC-IO: will collective I/O help with individual dataset?

HDF5 with individual dataset but with collective I/O enabled
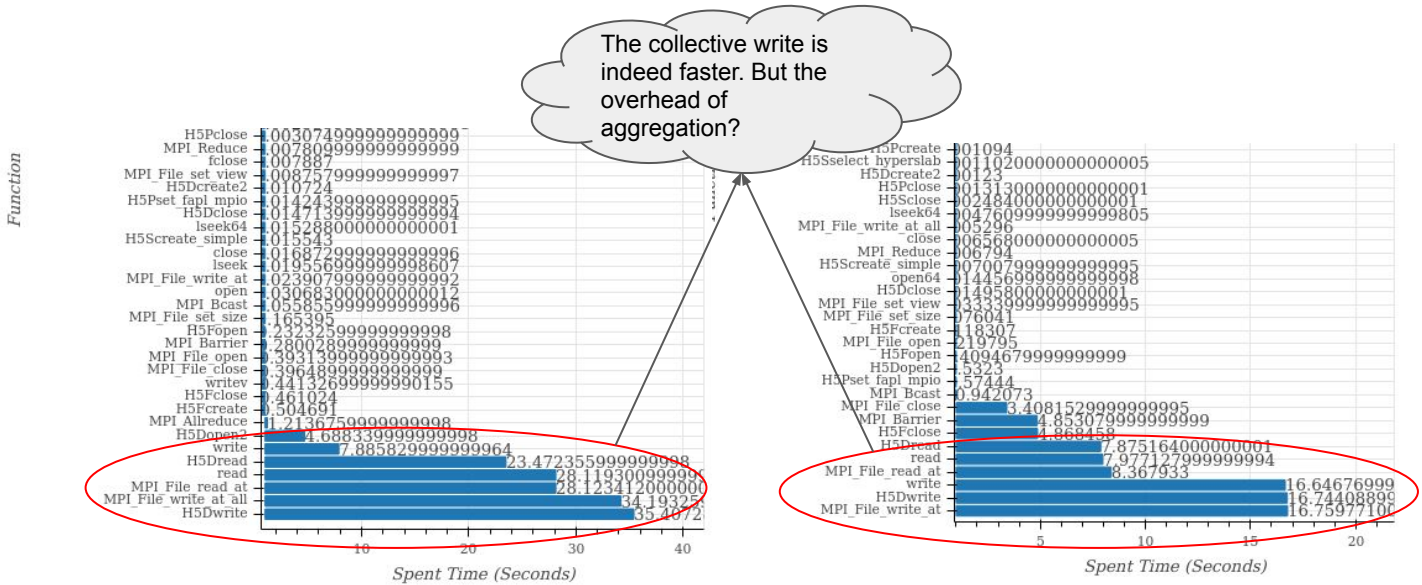


Stripe Size: 1M,
Stripe Count: 4

Stripe Size 128M
Stripe Count: 4

Stripe Size 512M
Stripe Count: 4

# HACC-IO: Is collective I/O always better?

- When the request size is big, the collective communication overhead increases and the benefits from collective I/O becomes limited.



Collective write
Stripe size:128M

Independent write

# Recorder

- Multi-level I/O tracing library that captures function calls from  HDF5, MPI and POSIX.
- No aggregation, it keeps every function and its parameters. Useful to exam access patterns.
- Built-in visualizations for access patterns, function counters, I/O sizes, etc.
- Also reports I/O conflicts such as write-after-write, write-after-read, etc. Useful for consistency semantics check (File systems with weaker consistency semantics).
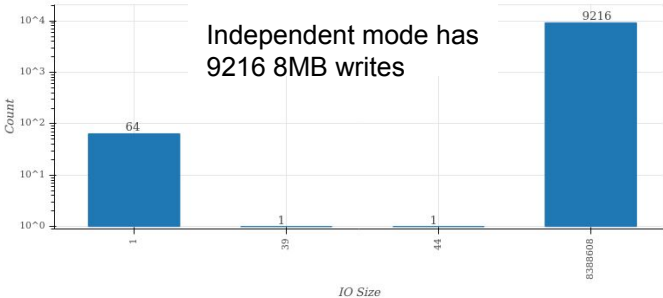
Wang, Chen, Jinghan Sun, Marc Snir, Kathryn Mohror, and Elsa Gonsiorowski. "Recorder 2.0: Efficient Parallel I/O Tracing and Analysis." In IEEE International Workshop on High-Performance Storage (HPS), 2020.
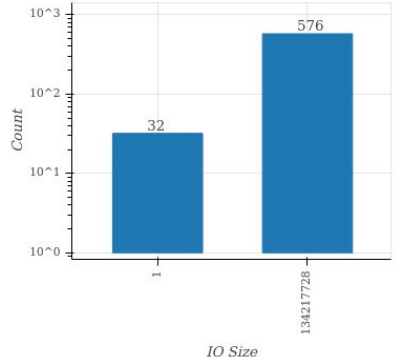https://github.com/uiuc-hpc/Recorder

Please ignore the rest slides.

# HACC-IO: Is collective I/O always better?

- Strip size of 128M



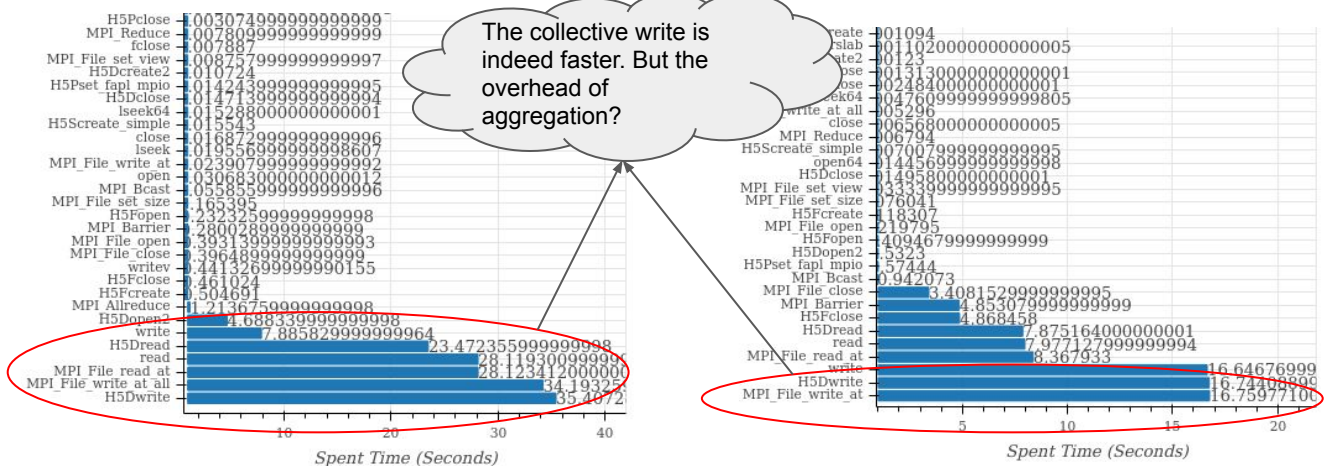Independent mode has 9216 8MB writes



Collective mode has 576 128M writes

Collective

Collective

The collective write is indeed faster. But the overhead of aggregation?

Setup

1024 Processes: 32 nodes x 32 ranks/node

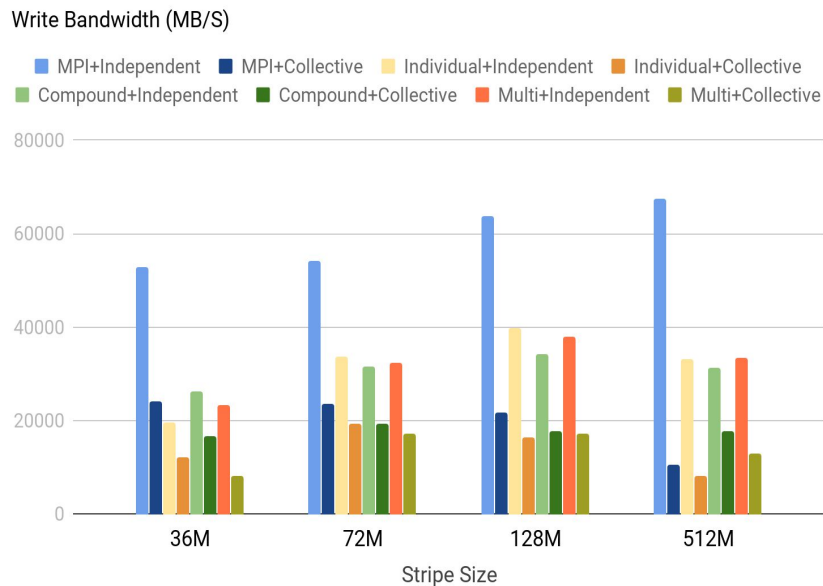HACC-I/O: write only. HDF5 collective metadata is **disabled**.

1.   Compound data type
2.   Individual dataset (9 dataset, one for each variable)
3.   With new HDF5 API Multi Allocation

Each of the above three methods are run with both independent and collective MPI-IO mode. So total of 6 runs.
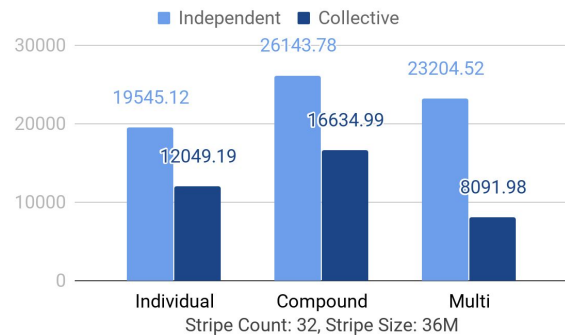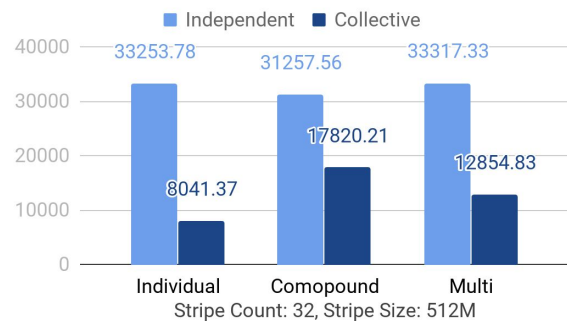
# Aggregated write bandwidth with different stripe size

1024 Processes: 32x32

8GB per variable, 72GB total.

# Aggregated write bandwidth with different stripe size

1024 Processes: 32x32
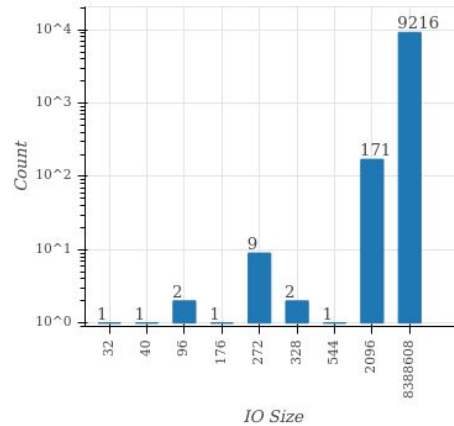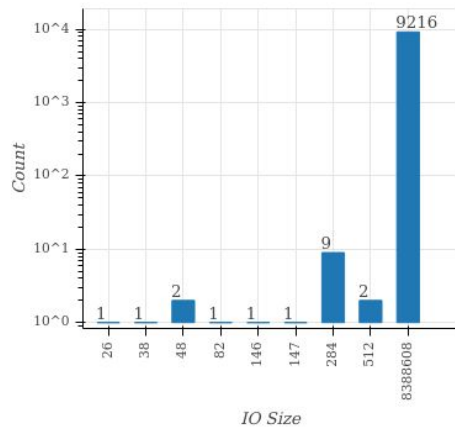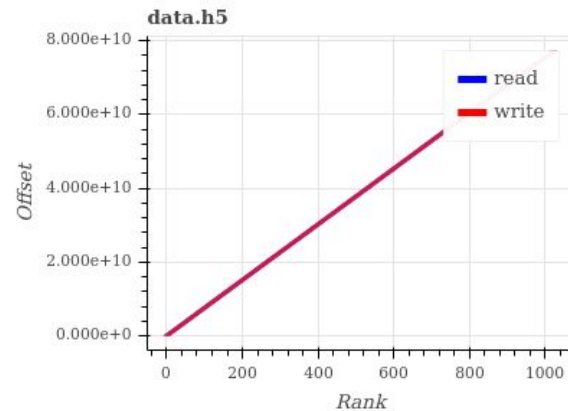
8GB per variable, 72GB total.



Write Bandwidth (MB/S)

Independent    Collective

26143.78

23204.52

19545.12

16634.99

12049.19

8091.98
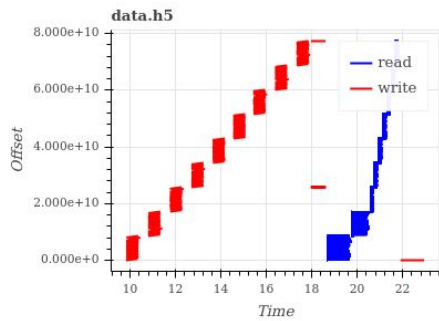
Individual    Compound    Multi

Stripe Count: 32, Stripe Size: 36M



Write Bandwidth (MB/s)

Independent    Collective

33253.78    31257.56    33317.33

17820.21

12854.83

8041.37

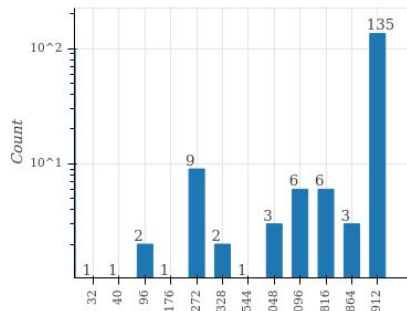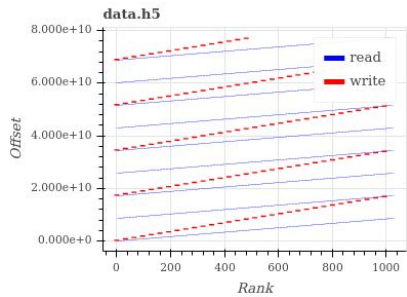Individual    Comopound    Multi
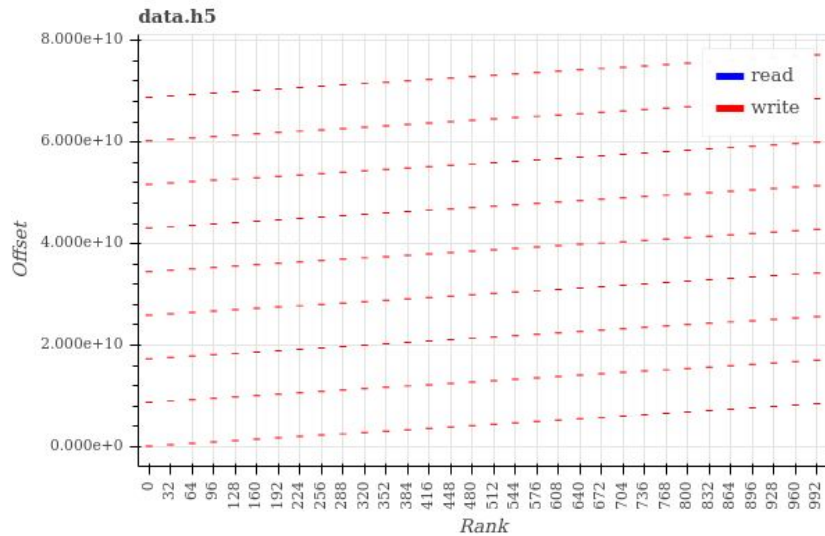
Stripe Count: 32, Stripe Size: 512M

# Independent

# Collective

HDF5 Individual Dataset with
Collective MPI-IO