October 7-9, 2025



# ALCF Hands-on HPC Workshop



# I/O libraries for Parallel Perf Part 1: MPI-IO Using and tuning MPI-IO and HDF5

Rob Latham (robl@mcs.anl.gov)
Math and Computer Science
Argonne National Laboratory

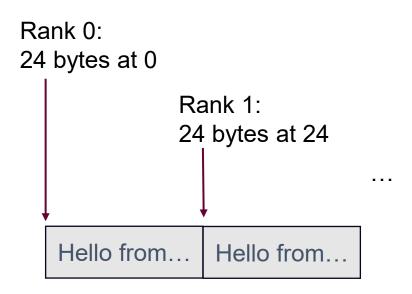
#### **MPI-IO**

- I/O interface specification for use in MPI apps
- Data model is same as POSIX: stream of bytes in a file
- Like classic POSIX in some ways...
  - Open() → MPI File open()
  - Pwrite() → MPI File write()
  - Close() → MPI File close()
- Features many improvements over POSIX:
  - Collective I/O
  - Noncontiguous I/O with MPI datatypes and file views
  - Nonblocking I/O
  - Fortran bindings (and additional languages)
- Implementations available on most (all?) platforms



#### "Hello World" MPI-IO style: contiguous

```
/* an "Info object": these store key-value strings for tuning the
* underlying MPI-IO implementation */
MPI Info create(&info);
snprintf(buf, BUFSIZE, "Hello from rank %d of %d\n", rank, nprocs);
len = strlen(buf);
/* We're working with strings here but this approach works well
* whenever amounts of data vary from process to process. */
MPI Exscan(&len, &offset, 1, MPI OFFSET, MPI SUM, MPI COMM WORLD);
MPI CHECK(MPI File open(MPI COMM WORLD, argv[1],
           MPI MODE CREATE | MPI MODE WRONLY, info, &fh));
/* all means collective. Even if we had no data to write, we would
* still have to make this call. In exchange for this coordination,
* the underlyng library might be able to greatly optimize the I/O */
MPI CHECK(MPI File write at all(fh, offset, buf, len, MPI CHAR,
           &status));
MPI CHECK(MPI File close(&fh));
```



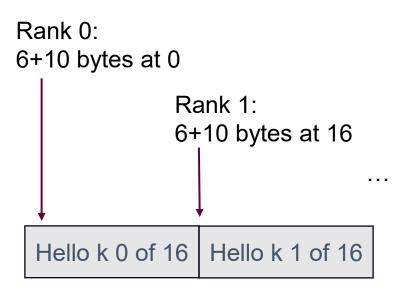


#### "Hello World" MPI-IO style: non-contiguous in memory

```
MPI Datatype memtype;
MPI Count memtype size;
/* sample string:
 * Hello from rank 8 of 16
 * the '-' indicates which elements an indexed type with
   lengths 6 and 10 at displacemnts 0 and
 * "10 from end of string" would select: */
int lengths[2] = {6, 10};
int displacements[2] = {0, len-10};
MPI_Type_indexed(2, lengths, displacements, MPI CHAR, &memtvpe);
MPI Type commit(&memtype);
MPI Type size x(memtype, &memtype size);
MPI CHECK(MPI File write at all(fh, offset, buf, 1, memtype,
            &status));
```

Hello from rank 1 of 16

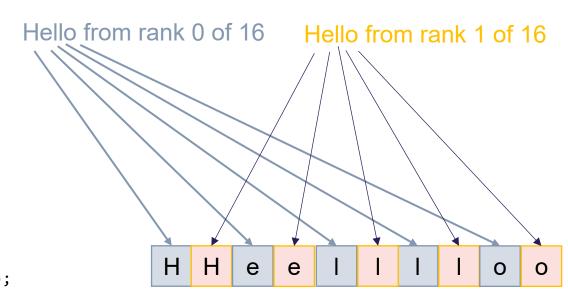
'lengths" and "displacements": each rank sends first six and last ten characters to file





#### "Hello World" MPI-IO style: non-contiguous in file

```
/* noncontiguous in file requres a "file view*/
MPI Datatype viewtype;
int *displacements;
displacements = malloc(len*sizeof(*displacements));
/* each process will write to its own "view" of the file:
 * Rank 0:
 *Hello from ...
 * Rank 1:
 * Hello from ...
for (int i=0; i< len; i++)</pre>
   displacements[i] = rank+(i*nprocs);
MPI Type create indexed block(len, 1, displacements, MPI CHAR, &viewtype);
MPI Type commit(&viewtype);
free(displacements);
MPI CHECK(MPI File open(MPI COMM WORLD, argv[1],
           MPI MODE CREATE | MPI MODE WRONLY, info, &fh));
MPI CHECK(MPI File set view(fh, 0, MPI CHAR, viewtype, "native", info));
MPI CHECK(MPI File write at all(fh, offset, buf, len, MPI CHAR,
           &status));
```



While this access describes lots of small regions, the library sees it as one single access and can optimize.



#### RUNNING

- Submit to the 'alcf training' queue and use the 'alcf training' project (aurora)
- I've provided a 'hello-aurora.sh' shell script
  - qsub -A alcf training -q alcf training ./hello-aurora.sh
- We'll use the DAOS file system
  - ALCF has made a "alcf training" pool on the "daos user" service
  - Job script will create your own container inside that pool
    - daos container create --type POSIX \$DAOS POOL \$DAOS CONT
  - DAOS is always running, but we have to "launch" the legacy file system view of it
    - launch-dfuse perf.sh \${DAOS POOL}:\${DAOS CONT}
    - Now shell tools can operate on /tmp/\${DAOS\_POOL}/\${DAOS\_CONT}
- There's a special "cpu binding" to place processes such that they use all 8 Aurora network cards.



#### **Output on Aurora**

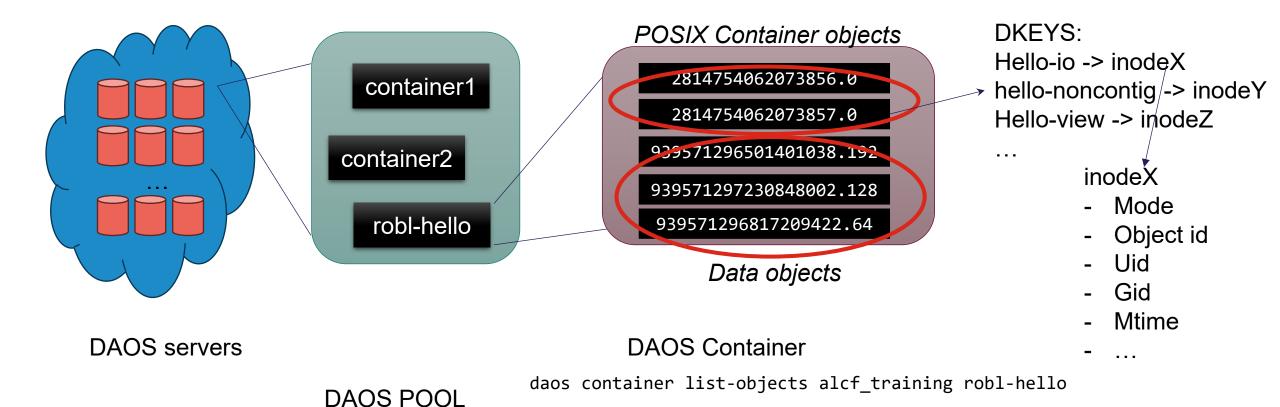
```
==== contiguous in memory and file
cat /tmp/ATPESC2025 0/robl-hello/hello.out
Hello from rank 0 of 16
Hello from rank 1 of 16
Hello from rank 15 of 16
==== noncontiguous in memory
cat /tmp/ATPESC2025_0/robl-hello/hello-
noncontig.out
Hello k 0 of 16
Hello k 1 of 16
Hello 15 of 16
```

#### Output of our hello programs



#### Under the hood: DAOS (essentially)

daos pool list-containers alcf training





#### Key takeaways

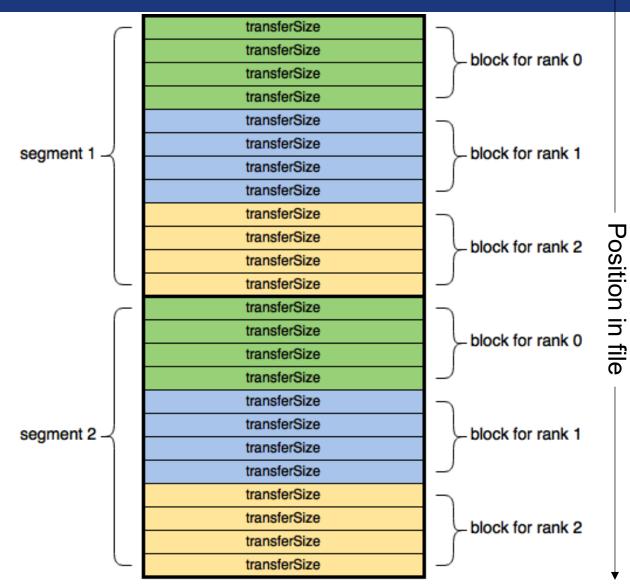
- Simple example but still captures important concepts
  - Info objects: tuning parameters:
    - enable/disable optimizations
    - Adjust buffer sizes
    - Select alternate strategies
  - Data placement in file specified by user
    - "shared file pointer" possible but not optimized
  - Collective vs independent I/O
  - Error checking!!!



#### The IOR benchmark

- MPI application benchmark
  - reads and writes data in configurable ways
  - I/O pattern can be <u>i</u>nterleaved <u>or random</u>
- Input:
  - transfer size, block size, segment count
  - interleaved or random
- Output: Bandwidth and IOPS
- Configurable backends
  - POSIX, STDIO, MPI-IO
  - HDF5, PnetCDF, S3, rados

https://github.com/hpc/ior



#### Hands-on: IOR and stripe size

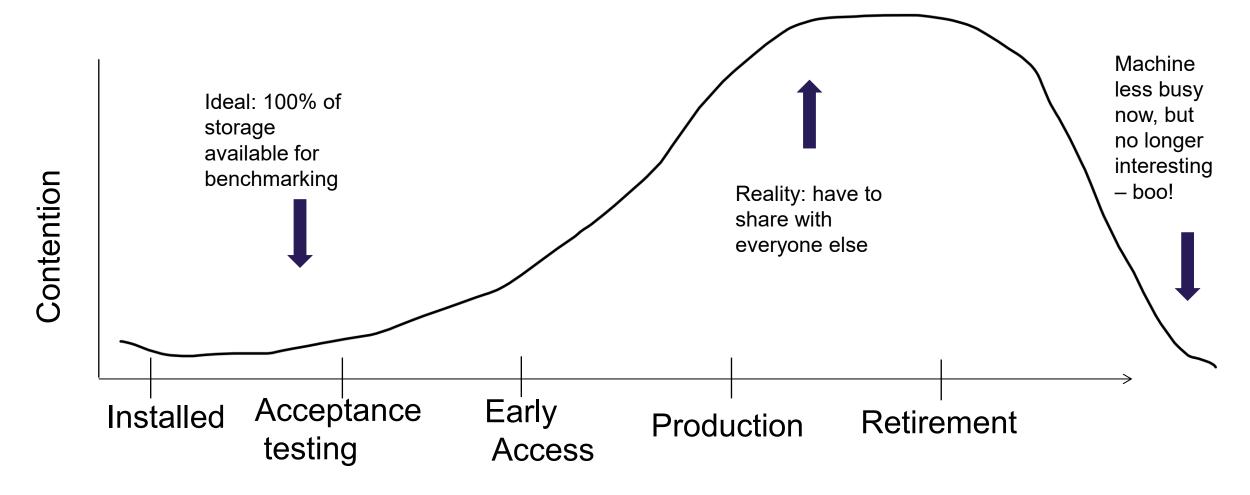
- For a fixed number of nodes, MPI processes, block size, and transfer size...
- Vary the stripe count
  - IOR environment variables
  - MPICH config file

```
$stripe=1
rm -f ${OUTPUT}/ior-stripe-$stripe.out
export IOR HINT MPI striping factor=$stripe
 # -a MPIIO: using MPI-IO so we can pass the "striping factor" hint
            : fsync after each write phase: push out dirty data to storage
            : reorder ranks: read from a different rank than the one that wrote
           : segments: each client will write to eight regions
           : repeat experiment five times: lots of variability in I/O
 # -i
           : transfer size: how big each request will be
 # -t
            : block size: how big each region will be in the file (needs to
  # -b
                    be a multiple of transfer size).
mpiexec -n ${NTOTRANKS} --ppn ${NRANKS PER NODE} \
       ior --mpiio.showHints -a MPIIO \
          -e -C -s 8 -i 5 \
         -t 1MiB -b 64MiB -o ${OUTPUT}/ior-stripe-$stripe.out
```

```
00000 11111 22222 ··· NNNN
```

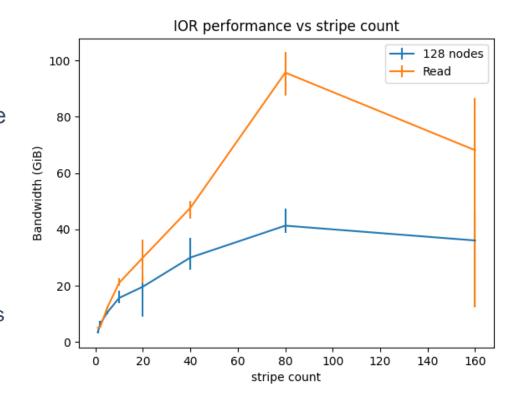


### Contention in benchmarkig



#### Hands on: IOR and stripe count

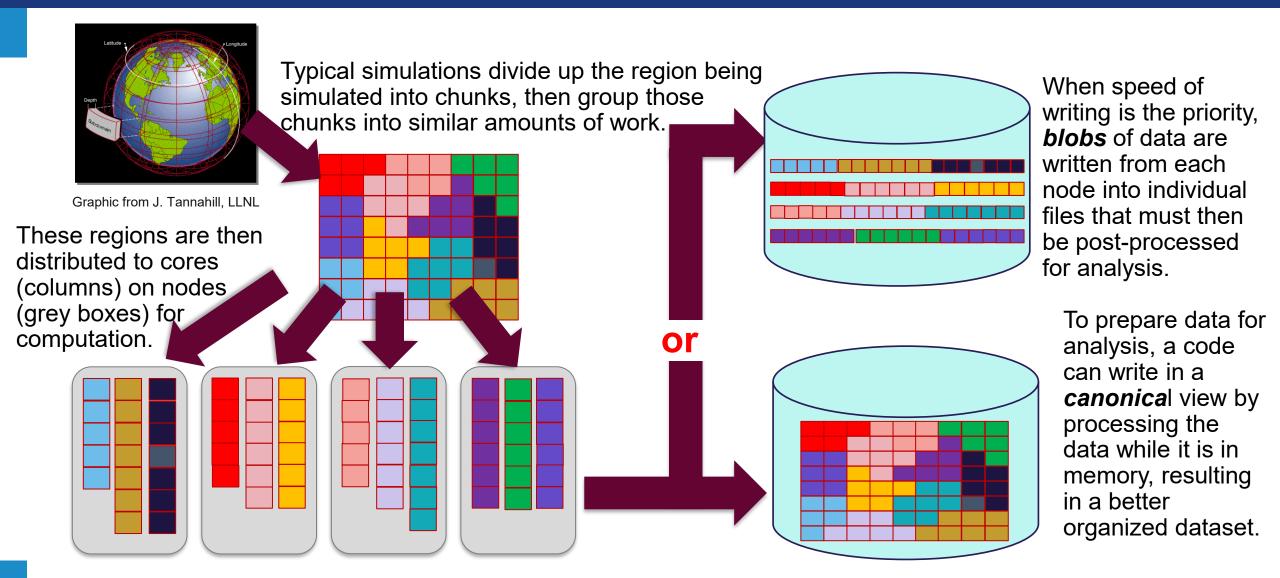
- Default stripe size is 1
  - Why? Most files small: optimizing for common case
- "All the servers" doesn't seem to hurt performance here
  - Ifs setstripe -1 /path/to/file
- Could go further with "overstriping"
  - Didn't work on Polaris: investigating
- "Where's my bandwidth?"
  - 128 nodes (network links) here
  - Shared file (so I can experiment with stripe count) means lustre locking overhead/coordination
- Graph at right from February 2023 any changes today?



visualization\_io/mpiio-hdf5/io-sleuthing/examples/striping

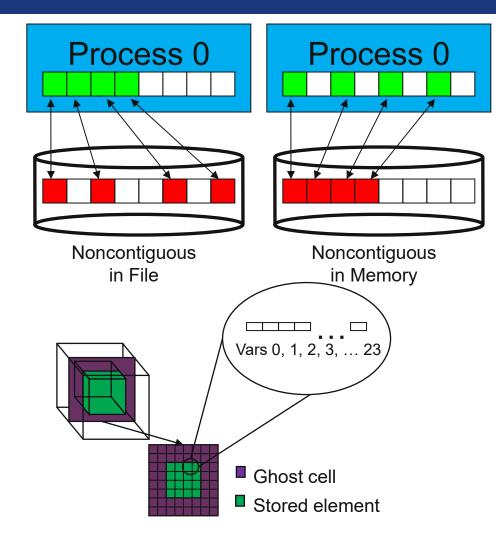


#### **Decomposition**



#### Contiguous and Noncontiguous I/O

- Contiguous I/O moves data from a single memory block into a single file region
- Noncontiguous I/O has three forms:
  - Noncontiguous in memory
  - Noncontiguous in file
  - Noncontiguous in both
- Structured data leads naturally to noncontiguous I/O (e.g., block decomposition)
- Describing noncontiguous accesses with a single operation passes more knowledge to I/O system



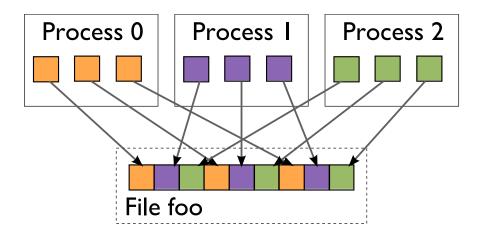
Extracting variables from a block and skipping ghost cells will code etc: <a href="https://github.com/argonne-lcf/ALCF">https://github.com/argonne-lcf/ALCF</a> Hands on HPC Workshop



#### I/O Transformations

Software between the application and the PFS performs transformations, primarily to improve performance

- Goals of transformations:
  - Reduce number of I/O operations to PFS (avoid latency, improve bandwidth)
  - Avoid lock contention (eliminate serialization)
  - Hide huge number of clients from PFS servers
- "Transparent" transformations don't change the final file layout
  - File system is still aware of the actual data organization
  - File can be later manipulated using serial POSIX I/O



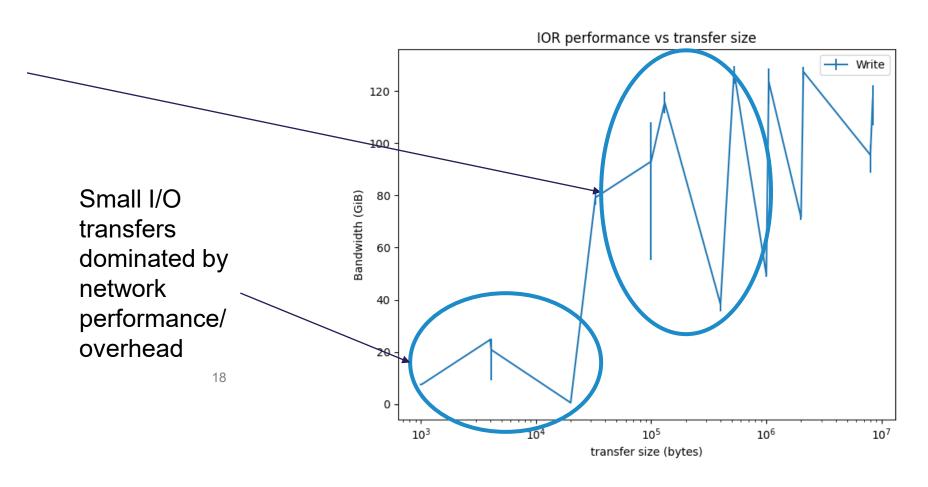
When we think about I/O transformations, we consider the mapping of data between application processes and locations in file



#### Request Size and I/O Rate

Sawtooth due to "power of 10" vs "power of 2" differences

In general, larger requests better.

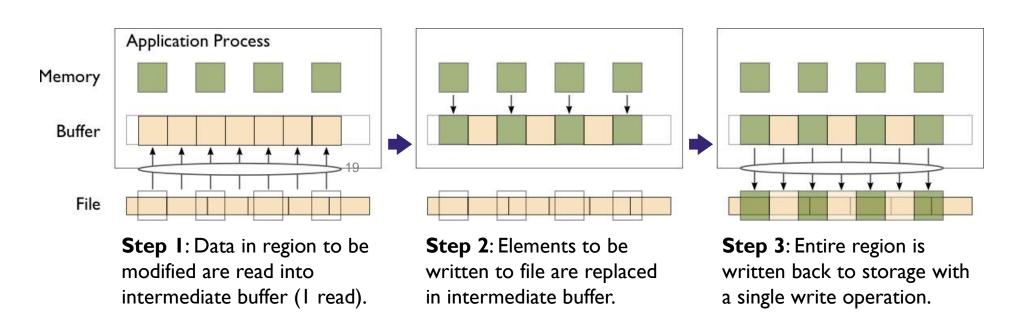


Tests run on 128 nodes, 16 process per node (2048 processes total) of HPE/Cray/Intel Auroa at Argonne, writing to DAOS



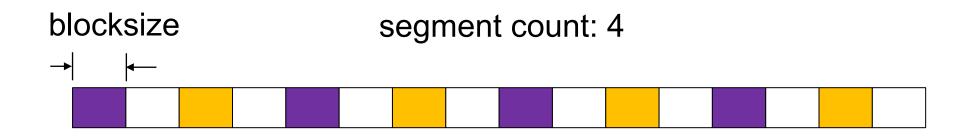
#### Reducing Number, Increasing Size of Operations

- Because most operations go over the network, I/O to a PFS incurs more latency than with a local FS
- *Data sieving* is a technique to address I/O latency by combining operations:
  - When reading, application process reads a large region holding all needed data and pulls out what is needed
  - When writing, three steps required (below)



# Noncontig with IOR

- IOR can describe access with an MPI datatype
  - --mpiio.useStridedDatatype -b ... -s ...
- (buggy in recent versions: use 4.0rc1 or newer)



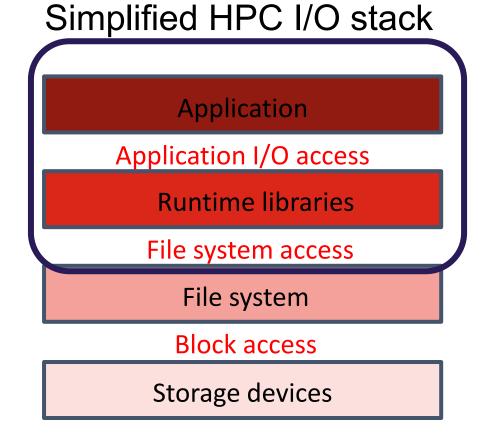


#### Darshan: Characterizing Application I/O

How is an application using the I/O system? How successful is it at attaining high performance?

Strategy: observe I/O behavior at the application and library level

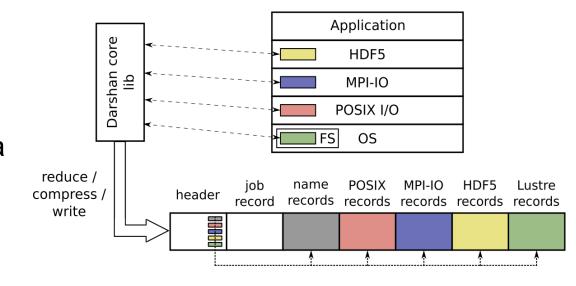
- What did the application intend to do?
- How much time did it take to do it?
- What can be done to tune and improve?





#### **How does Darshan work?**

- Darshan records file access statistics independently on each process
- At app shutdown, collect, aggregate, compress, and write log data
- After job completes, analyze Darshan log data
  - darshan-parser provides complete text-format dump of all counters in a log file
  - PyDarshan Python analysis module for Darshan logs, including a summary tool for creating HTML reports



- Originally designed for MPI applications, but in recent Darshan versions (3.2+) any dynamically-linked executable can be instrumented
  - In MPI mode, a log is generated for each app
  - In non-MPI mode, a log is generated for each *process*
- More information: <a href="https://docs.alcf.anl.gov/theta/performance-tools/darshan/">https://docs.alcf.anl.gov/theta/performance-tools/darshan/</a> or Shane's (concurrent) session



#### Data Sieving in Practice

Not always a win, particularly for writing:

IOR benchmark, fixed file size, increasing segments

Enabling data sieving instead made writes slower: why?

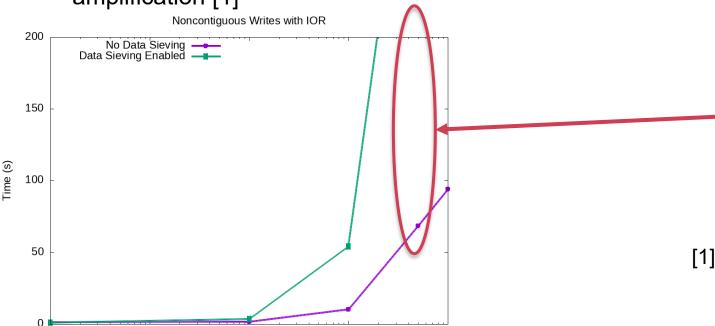
Locking to prevent false sharing (not needed for reads)

Multiple processes per node writing simultaneously

Internal ROMIO buffer too small, resulting in write

amplification [1]

10



1000

10000

	Naiive	Data Sieving
MPI-IO writes	960	960
MPI-IO Reads	0	0
Posix Writes	4 800 000	4 800 000
Posix Reads	0	4 800 784
MPI-IO bytes written	8.9 GiB	8.9 GiB
MPI-IO bytes read	0	0
Posix bytes read	0	2334 GiB
Posix bytes written	8.9 GiB	2343 GiB
Runtime (sec)	68.8	404.2

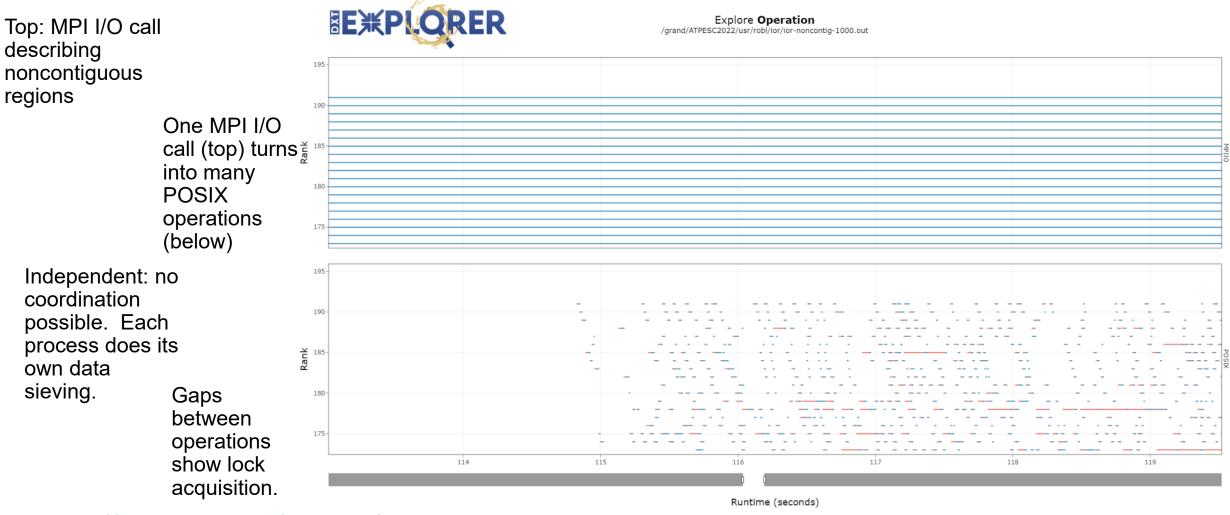
Selected Darshan statistics for 5000 segments



100

**Pieces** 

### Data Sieving: time line

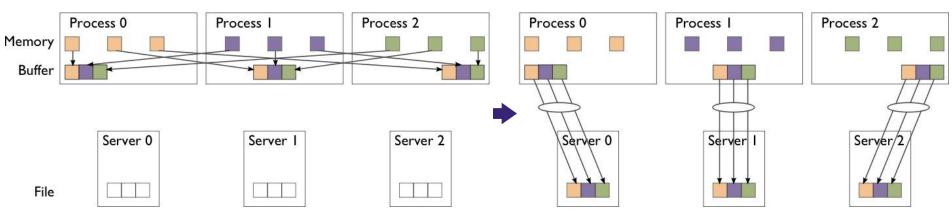


https://github.com/hpc-io/dxt-explorer Interactive log analysis tool by Jean Luca Bez



#### **Avoiding Lock Contention**

- To avoid lock contention when writing to a shared file, we can reorganize data between processes
- Two-phase I/O splits I/O into a data reorganization phase and an interaction with the storage system (two-phase write depicted):
  - Data exchanged between processes to match file layout
  - 0<sup>th</sup> phase determines exchange schedule (not shown)



**Phase I**: Data are exchanged between processes based on organization of data in file.

**Phase 2**: Data are written to file (storage servers) with large writes, no contention.



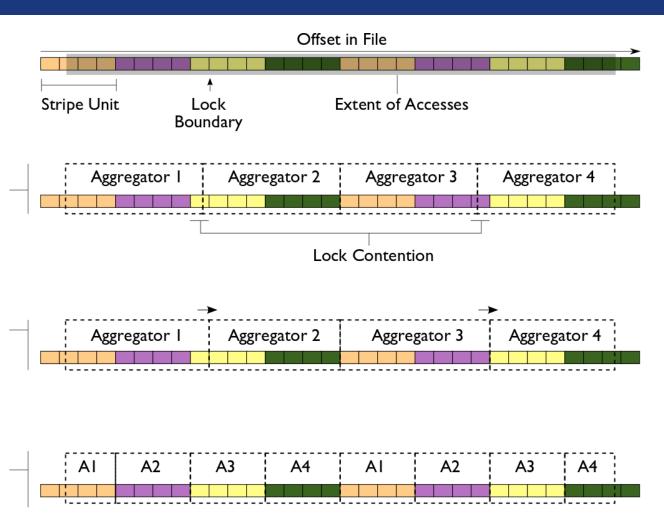
### Two-Phase I/O Algorithms

Imagine a collective I/O access using four aggregators to a file striped over four file servers (indicated by colors):

One approach is to evenly divide the region accessed across aggregators.

Aligning regions with lock boundaries eliminates lock contention.

Mapping aggregators to servers reduces the number of concurrent operations on a single server and can be helpful when locks are handed out on a per-server basis (e.g., Lustre).



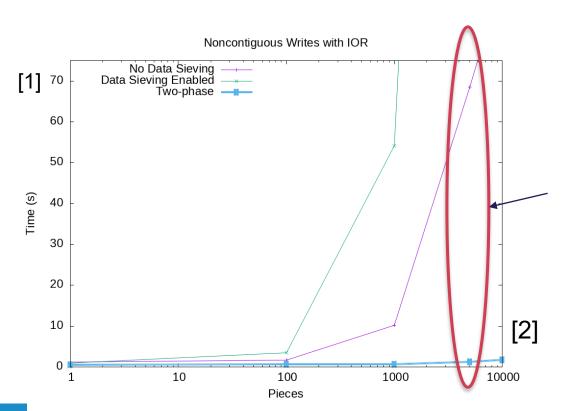
For more information, see W.K. Liao and A. Choudhary, "Dynamically Adapting File Domain Partitioning Methods for Collective I/O Based on Underlying Parallel File System Locking Protocols," SC2008, November 2008.



#### **Two-phase I/O in Practice**

- Consistent performance independent of access pattern
  - Note re-scaled y axis [1]
- No write amplification, no read-modify-write
- Some network communication but networks are fast

Requires "temporal locality" -- not great if writes "skewed", imbalanced, or some process enter collective late.

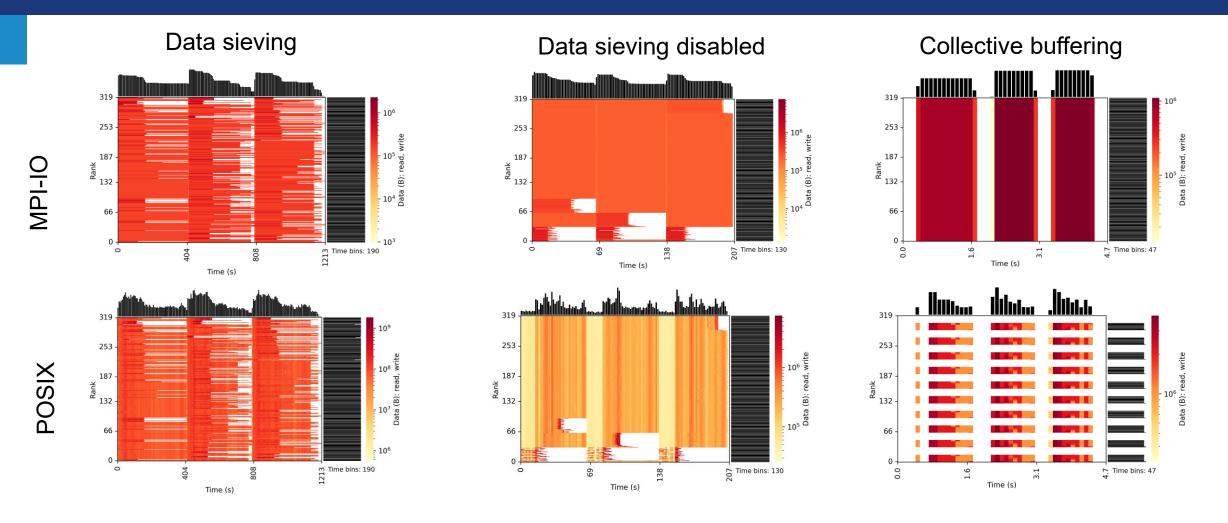


novou , irribularioou, or	Naiive	Data Sieving	Two-phase
MPI-IO writes	960	960	960
MPI-IO Reads	0	0	0
Posix Writes	4 800 000	4 800 000	9156
Posix Reads	0	4 800 784	0
MPI-IO bytes written	8.9 GiB	8.9 GiB	8.9 GiB
MPI-IO bytes read	0	0	0
Posix bytes read	0	2334 GiB	0
Posix bytes written	8.9 GiB	2343 GiB	8.9 GiB
Runtime (sec)	68.8	404.2	1.56

Selected Darshan statistics, 5000 segments



#### More investigation: Darshan heatmaps (Polaris, Lustre)

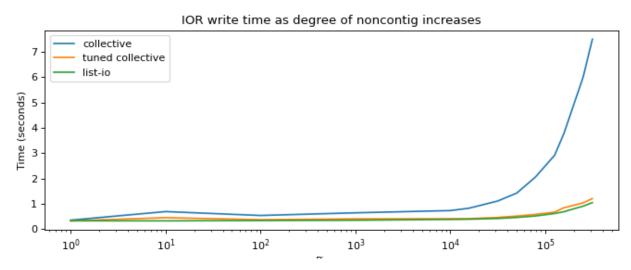


Effect of ROMIO optimizations on IOR benchmark: 5000 non-contiguous segments, three iterations. Note the x axis



#### DAOS: Collective I/O vs scatter-gather I/O

- Same IOR experiment but on Aurora this time
  - 2 nodes, 96 processes per node
- List-IO lets us avoid two sources of overhead
  - "rounds" of I/O no buffering at intermediate aggregator
  - No network exchange of data
- tuned: asking for more aggregators per node lets us use all 8 network cards
- Since List-IO does not aggregate, could be a problem at larger scale (evaluation "on my list")
  - Obviously, combining both approaches would be great (that's "on my list" now too...)



	Two- phase	Tuned Two- phase	List-IO
MPI-IO writes	1152	1152	1152
MPI-IO Reads	0	0	0
DAOS Writes	696	768	1152
DAOS Reads	0	0	0
MPI-IO bytes written	10.7 GiB	10.7 GiB	10.7 GiB
MPI-IO bytes read	0	0	0
DAOS bytes read	0	0	0
DAOS bytes written	10.7 GiB	10.7 GiB	10.7 GiB
Max MPI-IO write time	1.335 sec	0.35 sec	0.22 sec
Max DAOS write time	3.10 msec	3.485 msec	0.22 sec

Selected Darshan statistics, 5000 segments



#### **Tuning MPI-IO: info objects**

- You will likely never need these, but can help in specific situations:
- Both keys and values are strings
- Applicable to all ROMIO-based MPI-IO libraries

Hint	Default Value	effect
cb_buffer_size	16777216	An internal buffer for "two phase i/o". Bigger value takes away application memory, but results in fewer rounds of I/O
romio_cb_read romio_cb_write	Enable (on cray) automatic (ROMIO)	Turn on/off collective i/o: code will fall through to independent case
romio_no_indep_rw cb_config_list	True "*:*" (on Cray) or "*:1" elsewhere	"deferred open" – only i/o aggregators open the file. Open time not usually dominant factor unless total I/O moved per file fairly small
Cb_config_list	Default is "*:1" but should be "*:8" or higher	Aurora has eight network cards and needs 8 or more processes to obtain highest bandwidth



#### **Tuning MPI-IO: cray-specific hints**

- Hints that only work on Cray systems
- Perfectly fine to pass these (or anything) to any MPI library: libraries will ignore hints they don't recognize.
- More cray tuning at <a href="https://cpe.ext.hpe.com/docs/mpt/mpich/intro\_mpi.html#mpi-io-environment-">https://cpe.ext.hpe.com/docs/mpt/mpich/intro\_mpi.html#mpi-io-environment-</a> variables

Info key	Default value	effect
cray_cb_write_lock_mode	0	Set to "2" to try out "lock ahead": should allow greater concurrency
cray_cb_nodes_multiplier	1	Depending on stripe size and number of nodes, "2" or more might improve performance



#### **Data Model Libraries**

- Scientific applications work with structured data and desire more self-describing file formats
- PnetCDF and HDF5 are two popular "higher level" I/O libraries
  - Abstract away details of file layout
  - Provide standard, portable file formats
  - Include metadata describing contents
- For parallel machines, these use MPI and probably MPI-IO
  - MPI-IO implementations are sometimes poor on specific platforms, in which case libraries might directly call POSIX calls instead



#### The Parallel netCDF Interface and File Format

- Thanks to Wei-Keng Liao, Alok Choudhary, and Kaiyuan Hou (NWU) for their help in the development of PnetCDF.
- https://parallel-netcdf.github.io/



### Parallel NetCDF (PnetCDF)

- Based on original "Network Common Data Format" (netCDF) work from Unidata
  - Derived from their source code
- Data Model:
  - Collection of variables in single file
  - Typed, multidimensional array variables
  - Attributes on file and variables
- Features:
  - C, Fortran, and F90 interfaces (no python)
  - Portable data format (identical to netCDF)
  - Noncontiguous I/O in memory using MPI datatypes
  - Noncontiguous I/O in file using sub-arrays
  - Collective I/O
  - Non-blocking I/O
- Unrelated to netCDF-4 work
- Parallel-NetCDF tutorial:
  - https://parallel-netcdf.github.io/wiki/QuickTutorial.html
- Interface guide:
  - http://cucis.ece.northwestern.edu/projects/PnetCDF/doc/pnetcdf-c/index.html
  - 'man pnetcdf' on polaris (after loading module)



# Parallel netCDF (PnetCDF)

- (Serial) netCDF
  - API for accessing multi-dimensional data sets
  - Portable file format
  - Popular in both fusion and climate communities
- Parallel netCDF
  - Very similar API to netCDF
  - Tuned for better performance in today's computing environments
  - Retains the file format so netCDF and PnetCDF applications can share files
  - PnetCDF builds on top of any MPI-IO implementation

Cluster

PnetCDF

**ROMIO** 

Lustre

IBM AC922 (Summit)

**PnetCDF** 

Spectrum-MPI

**GPFS** 



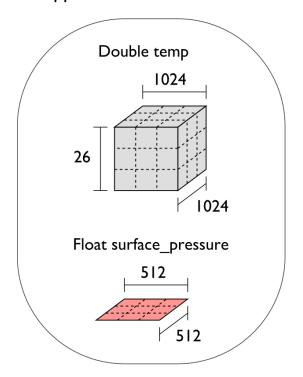
#### netCDF Data Model

• The netCDF model provides a means for storing multiple, multi-dimensional arrays in a single file.

Offset

Ξ.

#### **Application Data Structures**



#### netCDF File "checkpoint07.nc"

```
Variable "temp" {
  type = NC_DOUBLE,
  dims = {1024, 1024, 26},
  start offset = 65536,
  attributes = {"Units" = "K"}}

Variable "surface_pressure" {
  type = NC_FLOAT,
  dims = {512, 512},
  start offset = 218103808,
  attributes = {"Units" = "Pa"}}

< Data for "temp" >

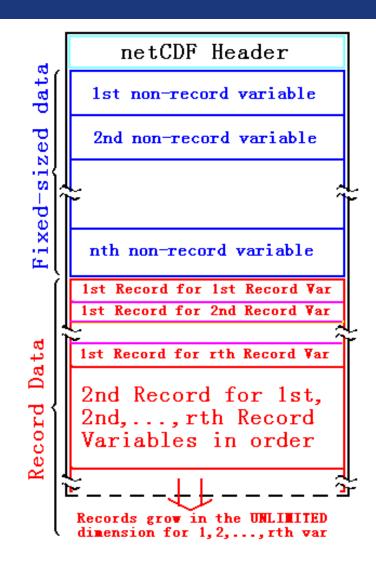
< Data for "surface_pressure" >
```

netCDF header describes the contents of the file: typed, multi-dimensional variables and attributes on variables or the dataset itself.

Data for variables is stored in contiguous blocks, encoded in a portable binary format according to the variable's type.

## Record Variables in netCDF

- Record variables are defined to have a single "unlimited" dimension
  - Convenient when a dimension size is unknown at time of variable creation
- Record variables are stored after all the other variables in an interleaved format
  - Using more than one in a file is likely to result in poor performance due to number of noncontiguous accesses





# Pre-declaring I/O

- netCDF / Parallel-NetCDF: bimodal write interface
  - Define mode: "here are my dimensions, variables, and attributes"
  - Data mode: "now I'm writing out those values"
- Decoupling of description and execution shows up several places
  - MPI non-blocking communication
  - Parallel-NetCDF "write combining" (talk more in a few slides)
  - MPI datatypes to a collective routines (if you squint really hard)



## "Hello world" Parallel-NetCDF style



## **Running on Polaris**

```
#!/bin/bash -1
#PBS -A ATPESC2024
#PBS -1 walltime=00:10:00
#PBS -l select=1
#PBS -1 place=scatter
#PBS -1 filesystems=home:eagle
#PBS -q debug
#PBS -N hello-io
#PBS -V
OUTPUT=/eagle/radix-io/${USER}/hello
mkdir -p ${OUTPUT}
NNODES=$(wc -1 < $PBS_NODEFILE)</pre>
NRANKS PER NODE=32
NTOTRANKS=$(( NNODES * NRANKS PER NODE ))
cd $PBS O WORKDIR
mpiexec -n $NTOTRANKS -ppn $NRANKS_PER_NODE \
        ./hello-pnetcdf ${OUTPUT}/hello-pnetcdf.nc
```

```
% ncmpidump /eagle/radix-io/${USER}/hello/hello-pnetcdf.nc
netcdf hello-pnetcdf {
// file format: CDF-2 (large file)
dimensions:
        d1 = 790;
variables:
        char v1(d1);
data:
 v1 = "Hello from rank 0 of 32\n",
    "Hello from rank 1 of 32\n",
    "Hello from rank 2 of 32\n",
    [...]
    "Hello from rank 27 of 32\n",
    "Hello from rank 28 of 32\n",
    "Hello from rank 29 of 32\n",
    "Hello from rank 30 of 32\n",
    "Hello from rank 31 of 32\n",
```

Job submission script

Output of "hello-pnetcdf"



# HANDS-ON: writing with Parallel-NetCDF

- 2-D array in file, each rank writes 'YDIM' (1) rows
- Many details managed by pnetcdf library
  - MPI-IO File views
  - offsets
- Be mindful of define/data mode: call ncmpi enddef()
- Library will take care of header i/o for you
- Define two dimensions
  - ncmpi def dim()
- Define one variable
  - ncmpi\_def\_var()
- Collectively put variable
  - ncmpi put vara int all()
  - 'start' and 'count' arrays: each process selects different regions
- Check your work with 'ncdump <filename>'
  - Hey look at that: serial tool reading parallel-written data: interoperability at work



# Solution fragments for Hands-on

#### Defining dimension: give name, size; get ID

```
/* row-major ordering */
NC_CHECK(ncmpi_def_dim(ncfile, "rows", YDIM*nprocs, &(dims[0])) );
NC_CHECK(ncmpi_def_dim(ncfile, "elements", XDIM, &(dims[1])) );
```

# Defining variable: give name, "rank" and dimensions (id); get ID Attributes: can be placed globally, on variables, dimensions

#### I/O: 'start' and 'count' give location, shape of subarray. 'All' means collective

```
start[0] = rank*YDIM; start[1] = 0;
count[0] = YDIM; count[1] = XDIM;
NC_CHECK(ncmpi_put_vara_int_all(ncfile, varid_array, start, count, values) );
```



0	1	2	3
10	11	12	13
20	21	22	23
30	31	32	33
240	74.5	2.2	

Full example in visualization\_io/mpiio-hdf5/hands-on/array



## **Inside PnetCDF Define Mode**

- In define mode (collective)
  - Use MPI\_File\_open to create file at create time
  - Set hints as appropriate (more later)
  - Locally cache header information in memory
    - All changes are made to local copies at each process
- At ncmpi\_enddef
  - Process 0 writes header with MPI\_File\_write\_at
  - MPI\_Bcast result to others
  - Everyone has header data in memory, understands placement of all variables
    - No need for any additional header I/O during data mode!



## **Inside PnetCDF Data Mode**

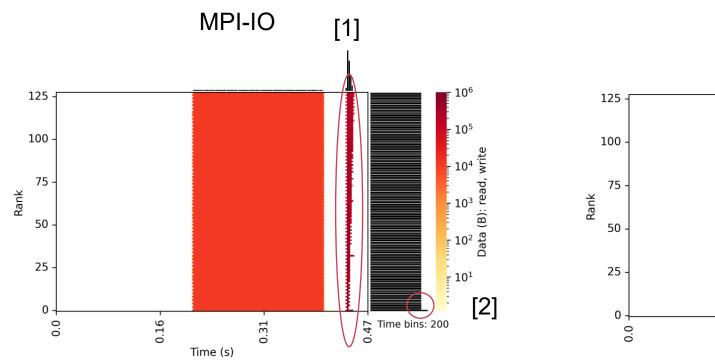
- Inside ncmpi\_put\_vara\_all (once per variable)
  - Each process performs data conversion into internal buffer
  - Uses MPI\_File\_set\_view to define file region
  - MPI\_File\_write\_all collectively writes data
- At ncmpi\_close
  - MPI\_File\_close ensures data is written to storage

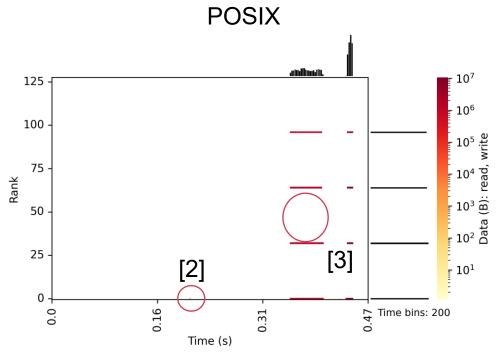
- MPI-IO performs optimizations
  - Two-phase possibly applied when writing variables
- MPI-IO makes PFS calls
  - PFS client code communicates with servers and stores data



## Inside PnetCDF: Darshan heatmap analysis

IOR writing Parallel-NetCDF (see visualization\_io/mpiio-hdf5/hands-on/ior/polaris/ior-pnetcdf.sh)



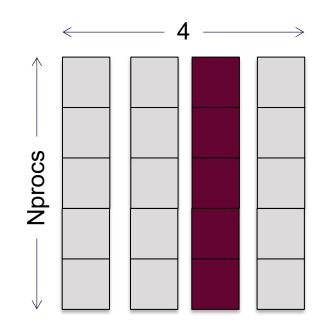


- [1]: all processes call MPI write and read re-reading going to be fast (cached)
- [2]: one process wrote header -- small: just one pixel in POSIX
- [3]: what you don't see only "aggregators" actually do I/O



# **HANDS-ON:** reading with pnetcdf

- Similar to MPI-IO reader: just read one row
- Operate on netcdf arrays, not MPI datatypes
- Shortcut: can rely on "convention"
  - One could know nothing about file as in previous slide
  - In our case we know there's a variable called "array" (id of 0) and an attribute called "iteration"
- Routines you'll need:
  - ncmpi\_inq\_dim to turn dimension id to dimension length
  - ncmpi\_get\_att\_int to read "iteration" attribute
  - ncmpi get vara int all to read column of array





# Solution fragments: reading with pnetcdf

#### Making **inq**uiry about variable, dimensions

```
NC CHECK (ncmpi inq var (ncfile, 0, varname, &vartype, &nr dims,
     dim ids,&nr attrs));
NC CHECK(ncmpi inq dim(ncfile, dim ids[0], NULL, &(dim lens[0])) );
NC CHECK(ncmpi inq dim(ncfile, dim ids[1], NULL, &(dim lens[1])) );
```

#### The "Iteration" attribute

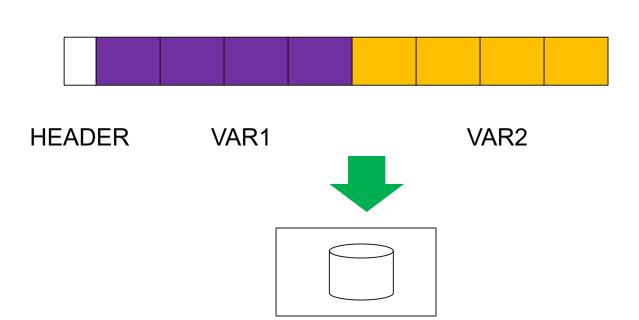
```
NC CHECK(ncmpi get att int(ncfile, 0, "iteration", &iterations));
```

No file views or datatypes: just a starting coordinate and size – everyone reads same slice in this case

```
count[0] = dim lens[0]; count[1] = 1;
starts[0] = 0; starts[1] = XDIM/2;
NC CHECK(ncmpi get vara int all(ncfile, 0, starts, count, read buf));
```



# Parallel-NetCDF write-combining optimization

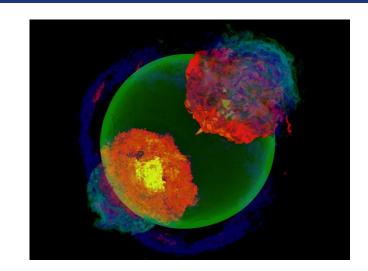


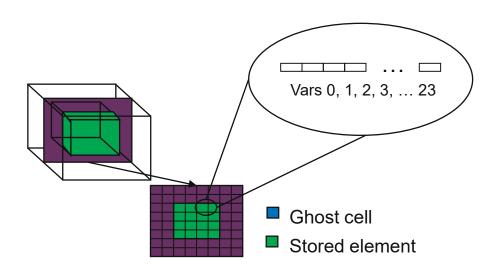
- netCDF variables laid out contiguously
- Applications typically store data in separate variables
  - temperature(lat, long, elevation)
  - Velocity\_x(x, y, z, timestep)
- Operations posted independently, completed collectively
  - Defer, coalesce synchronization
  - Increase average request size



## **Example: FLASH Astrophysics**

- FLASH is an astrophysics code for studying events such as supernovae
  - Adaptive-mesh hydrodynamics
  - Scales to 1000s of processors
  - MPI for communication
- Frequently checkpoints:
  - Large blocks of typed variables from all processes
  - Portable format
  - Canonical ordering (different than in memory)
  - Skipping ghost cells

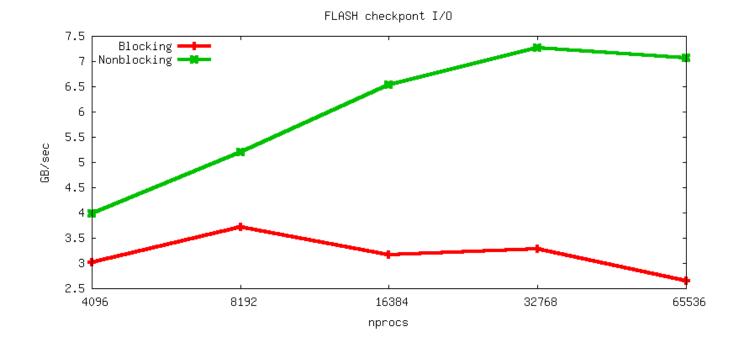






## FLASH Astrophysics and the write-combining optimization

- FLASH writes one variable at a time
- Could combine all 4D variables (temperature, pressure, etc) into one 5D variable
  - Altered file format (conventions) requires updating entire analysis toolchain
- Write-combining provides improved performance with same file conventions
  - Larger requests, less synchronization.





# **HANDS-ON:** pnetcdf write-combining

- 1. Define a second variable, changing only the name
- 2. Write this second variable to the netcdf file
- 3. Convert to the non-blocking interface (ncmpi\_iput\_vara\_int)
  - not collective "collectiveness" happens in ncmpi\_wait\_all
  - takes an additional 'request' argument
- 4. Wait (collectively) for completion



# Solution fragments for write-combining

#### Defining a second variable

#### The non-blocking interface: looks a lot like MPI

#### Waiting for I/O to complete

```
/* all the I/O actually happens here */
NC_CHECK(ncmpi_wait_all(ncfile, 2, reqs, status));
```



## **Hands-on continued**

- · Look at the darshan output. Compare to darshan output for single-variable writing or reading
  - Results on polaris surprised me: vendor might know something I don't
    - Maybe some kind of small-io optimization?



# PnetCDF Wrap-Up

- PnetCDF gives us
  - Simple, portable, self-describing container for data
  - Collective I/O
  - Data structures closely mapping to the variables described
- If PnetCDF meets application needs, it is likely to give good performance
  - Type conversion to portable format does add overhead
- Some limits on (old, common CDF-2) file format:
  - Fixed-size variable: < 4 GiB</li>
  - Per-record size of record variable: < 4 GiB</li>
  - 2<sup>32</sup> -1 records
  - Contributed extended file format to relax these limits (CDF-5, released in pnetcdf-1.1.0, November 2009, integrated in Unidata NetCDF-4.4)



# The HDF5 Interface and **File Format**





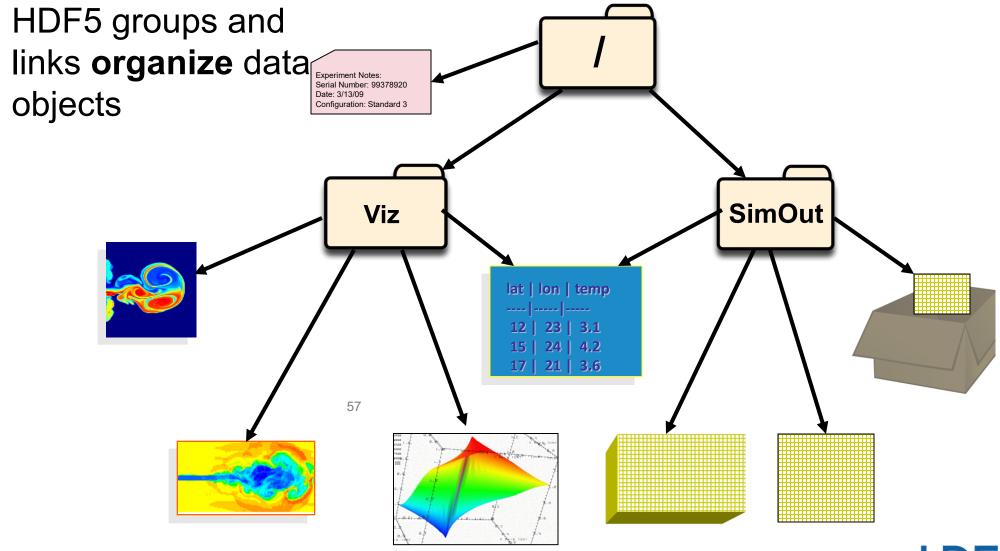
## HDF5

- Hierarchical Data Format, from The HDF Group (formerly of NCSA)
  - https://www.hdfgroup.org/
- Data Model:
  - Hierarchical data organization in single file
  - Typed, multidimensional array storage
  - Attributes on any HDF5 "object" (dataset, data, groups)
- Features:
  - C, C++, Fortran, Java (JNI) interfaces
    - Community-supported Python, Lua, R
  - Portable data format
  - Optional compression (even in parallel I/O mode)
  - Chunking: efficient row or column oriented access
  - Noncontiguous I/O (memory and file) with hyperslabs
- Parallel HDF5 tutorial:
  - https://portal.hdfgroup.org/display/HDF5/Introduction+to+Parallel+HDF5





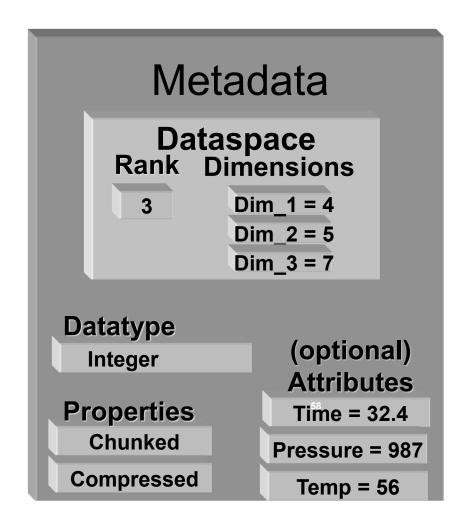
# **HDF5 Groups and Links**

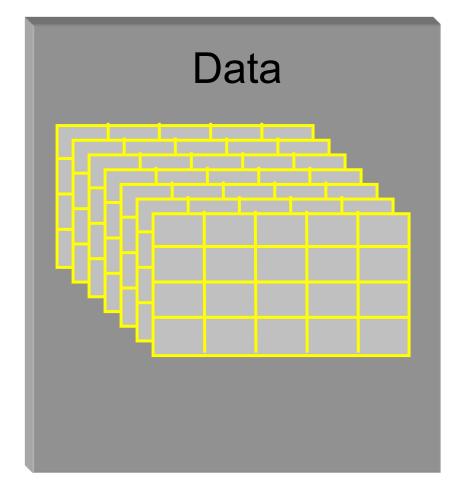






## **HDF5 Dataset**

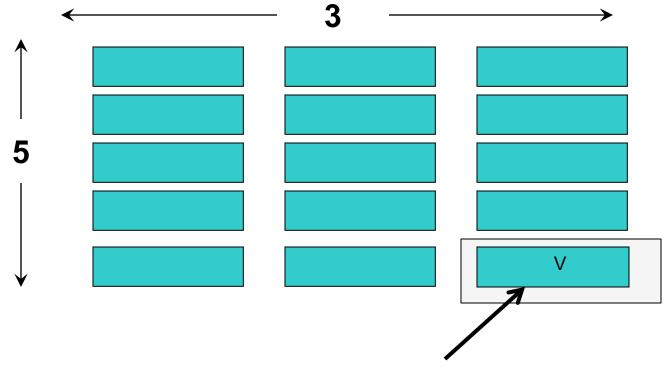








## **HDF5 Dataset**



Datatype: 16-byte integer

**Dataspace:** Rank = 2

Dimensions =  $5 \times 3$ 



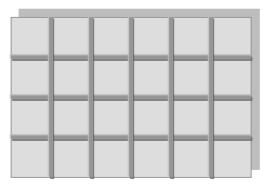


# **HDF5 Dataspaces**

#### Two roles:

Dataspace contains spatial information (logical layout) about a dataset stored in a file

- Rank and dimensions
- Permanent part of dataset definition



Rank = 2

Dimensions = 4x6
Subsets: Dataspace describes application's data buffer and data elements participating in I/O



Rank = 1

Dimension = 10

## **Basic Functions**

H5Fcreate (H5Fopen) create (open) File

H5Screate simple/H5Screate create dataspace

H5Dcreate (H5Dopen) create (open) Dataset

H5Sselect\_hyperslab select subsections of data

access Dataset H5Dread, H5Dwrite

H5Dclose close Dataset

close dataSpace H5Sclose

H5Fclose close File

NOTE: Order not strictly specified





## "Hello World" HDF5 style

Cannot fit all in one slide: here are some highlights (see 'hello-hdf5.c' for full example)

```
file = H5Fcreate(argv[1], H5F_ACC_TRUNC, H5P_DEFAULT,
file access property list);
```

- "property lists" used a lot in HDF5 (see next slide)
- Serial interface came first, with parallel features added later

```
/* in this simple example everyone writes their string to a
   1-d dataset; HDF5 supports variable length arrays ("ragged
   arrays") but these datatypes have odd interactions with parallel
i/o */

/* like writing to a plain file, we'll create one big variable
and everyone can write their string to the right (non-
overlapping) place in the file */
hid_t dataset, datatype, file_space;
hsize_t size=varlen;

file_space = H5Screate_simple(1, &size, NULL);
/* remember we got 'offset' from the MPI_Exscan above */
hsize_t start=offset, count=len;
status = H5Sselect_hyperslab(file_space, H5S_SELECT_SET,
   &start, NULL, &count, NULL);
```

- Lots of flexibility in how memory, file regions described
- Lots more we could say about "hyperslab"



## HDF5 example: opening with MPI-IO

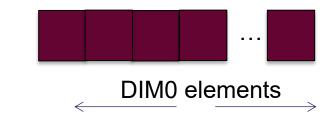
```
/* Initialize MPI */
MPI_Init(&argc, &argv);
/* Create an HDF5 file access property list */
fapl id = H5Pcreate (H5P FILE ACCESS);
/* Set file access property list to use the MPI-IO file driver */
ret = H5Pset_fapl_mpio(fapl_id, MPI_COMM_WORLD, MPI_INFO_NULL);
/* Create the file collectively */
file id = H5Fcreate(argv[1], H5F ACC TRUNC, H5P DEFAULT, fapl id);
/* Release file access property list */
ret = H5Pclose(fapl id);
```



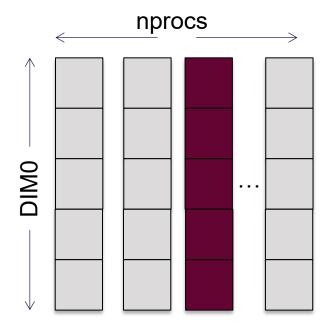
## HDF5 example: setting up data transfer

```
/* Select column of elements in the file dataset */
file_start[0] = 0; file_start[1] = mpi_rank;
file_count[0] = DIM0; file_count[1] = 1;
ret = H5Sselect_hyperslab(file_space_id, H5S_SELECT_SET,
       file start, NULL, file count, NULL);
ret = H5Sselect_hyperslab(mem_space_id, H5S_SELECT_SET,
       mem_start, NULL, mem_count, NULL);
/* Set up the collective transfer properties list */
dxpl id = H5Pcreate(H5P DATASET XFER);
ret = H5Pset dxpl mpio(dxpl id, H5FD MPIO COLLECTIVE);
/* Write data (one column of doubles) collectively */
ret = H5Dwrite(dset_id, H5T_NATIVE_DOUBLE, mem_space_id,
       file space id, dxpl id, write buf);
```

#### **MEMORY**



#### FILE





## **Effect of HDF5 Tuning**

- HDF5 property lists can have big impact on internal operations
- Collective I/O vs. Independent I/O
  - Huge reduction in operation count
  - Implies all processes hit I/O at same time
- Collective metadata (new in 1.10.2)
  - Further reduction in op count, especially reads (reading HDF5 internal layout information)
  - Big implications for performance at scale

Operation counts	Independent	Coll. I/O	Coll. MD
POSIX Write	3680007	9	9
MPI-IO Indep write	3680007	7	0
MPI IO Collective Write	0	16	48
POSIX Read	3680113	115	10
MPI-IO indep read	3680113	113	8
MPI-IO collective read	0	16	16

Selected Darshan statistics for 16 MPI processes writing 230 K doubles to HDF dataset, reading back same.

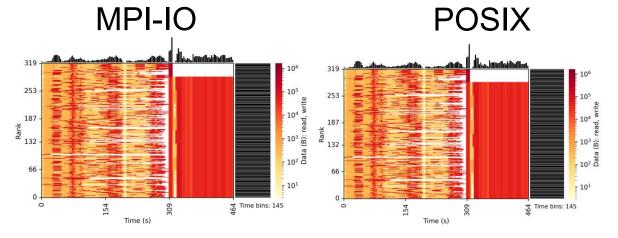
visualization\_io/mpiio-hdf5/hands-on/hdf5/h5par-comparison.c

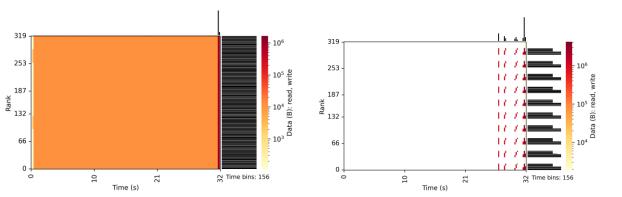




## **Effect of HDF5 Tuning**

- HDF5 property lists can have big impact on internal operations
- Collective I/O vs. Independent I/O
  - Huge reduction in operation count
  - Implies all processes hit I/O at same time
- Collective metadata (new in 1.10.2)
  - Further reduction in op count, especially reads (reading HDF5 internal layout information)
  - Big implications for performance at scale





visualization\_io/mpiio-hdf5/io-sleuthing/examples/hdf5





## **HDF5** in other languages

- Python:
  - H5py: <a href="http://www.h5py.org/">http://www.h5py.org/</a>
    - closely coupled with mpi4py and numpy;
    - some collective tuning not exposed at python level
- C++:
  - Highfive: <a href="https://github.com/BlueBrain/HighFive">https://github.com/BlueBrain/HighFive</a>
    - header-only interface to HDF5 C API



## **New HDF5 features:**

- New in HDF5-1.14.0
  - Async operations
    - Potential for background progress
  - Multi-dataset I/O
    - Similar to pnetcdf "operation combining"



## **Data Model I/O libraries**

- Parallel-NetCDF: <a href="http://www.mcs.anl.gov/pnetcdf">http://www.mcs.anl.gov/pnetcdf</a>
- HDF5: <a href="http://www.hdfgroup.org/HDF5/">http://www.hdfgroup.org/HDF5/</a>
- NetCDF-4: http://www.unidata.ucar.edu/software/netcdf/netcdf-4/
  - netCDF API with HDF5 back-end
- ADIOS: <a href="http://adiosapi.org">http://adiosapi.org</a>
  - Configurable (xml) I/O approaches
- SILO: https://wci.llnl.gov/codes/silo/
  - A mesh and field library on top of HDF5 (and others)
- H5part: http://vis.lbl.gov/Research/AcceleratorSAPP/
  - simplified HDF5 API for particle simulations
- GIO: <a href="https://svn.pnl.gov/gcrm">https://svn.pnl.gov/gcrm</a>
  - Targeting geodesic grids as part of GCRM
- PIO:
  - climate-oriented I/O library; supports raw binary, parallel-netcdf, or serial-netcdf (from master)
- Many more: consider existing libs before deciding to make your own.
- Note absence of a "machine learning" library research opportunity for someone!



# Wrap-up

- Lots of activity, history making I/O better... Still a lot to do!
  - Workflow, task-oriented, AI/ML
- ALCF consultants, research community eager to help

