

Overview of Parallel HDF5 and Performance Tuning in HDF5 Library

HDF and HDF-EOS Workshop VI
Elena Pourmal, Albert Cheng

Outline

- Overview of Parallel HDF5 design
- Setting up parallel environment
- Programming model for
 - _ Creating and accessing a File
 - _ Creating and accessing a Dataset
 - _ Writing and reading Hyperslabs
- Performance tuning in HDF5
- Parallel tutorial available at
 - _ <http://hdf.ncsa.uiuc.edu/HDF5/doc/Tutor>

PHDF5 Initial Target

- Support for MPI programming
- Not for shared memory programming
 - _Threads
 - _OpenMP
- Has some experiments with
 - _Thread-safe support for Pthreads
 - _OpenMP if called “correctly”

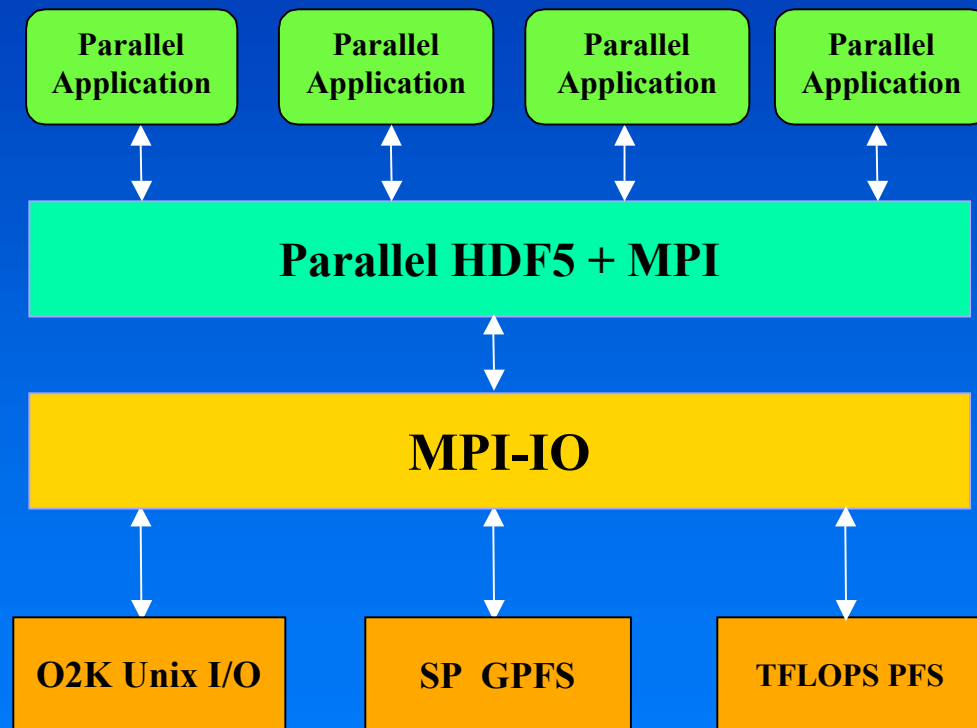
PHDF5 Requirements

- PHDF5 files compatible with serial HDF5 files
 - _ Shareable between different serial or parallel platforms
- Single file image to all processes
 - _ One file per process design is undesirable
 - Expensive post processing
 - Not useable by different number of processes
- Standard parallel I/O interface
 - _ Must be portable to different platforms

Implementation Requirements

- No use of Threads
 - _ Not commonly supported (1998)
- No reserved process
 - _ May interfere with parallel algorithms
- No spawn process
 - _ Not commonly supported even now

PHDF5 Implementation Layers



User Applications

HDF library

Parallel I/O layer

Parallel File systems

Parallel Environment Requirements

- MPI with MPI-IO
 - _Argonne ROMIO
 - _Vendor' s MPI-IO
- Parallel file system
 - _IBM GPFS
 - _PVFS

How to Compile PHDF5

- h5cc _ HDF5 compiler command
_ Similar to mpicc
- To compile:
% h5cc h5prog.c
- Show the compiler commands without executing them (i.e., dryrun):
% h5cc _show h5prog.c

Collective vs. Independent Calls

- MPI definition of collective call
 - _ All processes of the communicator must participate in the right order
- Independent means not collective
- Collective is not necessarily synchronous

Programming Restrictions

- Most PHDF5 APIs are collective
- PHDF5 opens a parallel file with a communicator
 - _ Returns a file-handle
 - _ Future access to the file via the file-handle
 - _ All processes must participate in collective PHDF5 APIs
 - _ Different files can be opened via different communicators

Examples of PHDF5 API

- Examples of PHDF5 collective API
 - _ File operations: H5Fcreate, H5Fopen, H5Fclose
 - _ Objects creation: H5Dcreate, H5Dopen, H5Dclose
 - _ Objects structure: H5Dextend (increase dimension sizes)
- Array data transfer can be collective or independent
 - _ Dataset operations: H5Dwrite, H5Dread

What Does PHDF5 Support ?

- After a file is opened by the processes of a communicator
 - _ All parts of file are accessible by all processes
 - _ All objects in the file are accessible by all processes
 - _ Multiple processes write to the same data array
 - _ Each process writes to individual data array

PHDF5 API Languages

- C and F90 language interfaces
- Platforms supported:
 - _ IBM SP2 and SP3
 - _ Intel TFLOPS
 - _ SGI Origin 2000
 - _ HP-UX 11.00 System V
 - _ Alpha Compaq Clusters
 - _ Linux clusters
 - _ SUN clusters
 - _ Cray T3E

Creating and Accessing a File Programming model

- HDF5 uses access template object to control the file access mechanism
- General model to access HDF5 file in parallel:
 - _ Setup MPI-IO access template
 - _ Open File
 - _ Close File

Setup access template

Each process of the MPI communicator creates an access template and sets it up with MPI parallel access information

C:

```
herr_t H5Pset_fapl_mpio(hid_t plist_id,  
                      MPI_Comm comm, MPI_Info info);
```

F90:

```
h5pset_fapl_mpio_f(plist_id, comm, info);  
integer(hid_t) :: plist_id  
integer        :: comm, info
```

plist_id is a file access property list identifier

C Example

Parallel File Create

```
23     comm = MPI_COMM_WORLD;
24     info = MPI_INFO_NULL;
26     /*
27      * Initialize MPI
28      */
29     MPI_Init(&argc, &argv);
33     /*
34      * Set up file access property list for MPI-IO access
35      */
36     plist_id = H5Pcreate(H5P_FILE_ACCESS);
37     H5Pset_fapl_mpio(plist_id, comm, info);
38
42     file_id = H5Fcreate(H5FILE_NAME, H5F_ACC_TRUNC,
43                        H5P_DEFAULT, plist_id);
49     /*
50      * Close the file.
51      */
52     H5Fclose(file_id);
54     MPI_Finalize();
```


F90 Example

Parallel File Create

```
23 comm = MPI_COMM_WORLD
24 info = MPI_INFO_NULL
26 CALL MPI_INIT(mpierror)
29 !
30 ! Initialize FORTRAN predefined datatypes
32 CALL h5open_f(error)
34 !
35 ! Setup file access property list for MPI-IO access.
37 CALL h5pcreate_f(H5P_FILE_ACCESS_F, plist_id, error)
38 CALL h5pset_fapl_mpio_f(plist_id, comm, info, error)
40 !
41 ! Create the file collectively.
43 CALL h5fcreate_f(filename, H5F_ACC_TRUNC_F, file_id,
    error, access_prp = plist_id)
45 !
46 ! Close the file.
49 CALL h5fclose_f(file_id, error)
51 !
52 ! Close FORTRAN interface
54 CALL h5close_f(error)
56 CALL MPI_FINALIZE(mpierror)
```

Creating and Opening Dataset

- All processes of the MPI communicator open/close a dataset by a collective call
 - _ C: H5Dcreate or H5Dopen; H5Dclose
 - _ F90: h5dcreate_f or h5dopen_f; h5dclose_f
- All processes of the MPI communicator extend dataset with unlimited dimensions before writing to it
 - _ C: H5Dextend
 - _ F90: h5dextend_f

C Example

Parallel Dataset Create

```
56 file_id = H5Fcreate(...);
57 /*
58  * Create the dataspace for the dataset.
59  */
60 dimsf[0] = NX;
61 dimsf[1] = NY;
62 filespace = H5Screate_simple(RANK, dimsf, NULL);
63
64 /*
65  * Create the dataset with default properties collective.
66  */
67 dset_id = H5Dcreate(file_id, "dataset1", H5T_NATIVE_INT,
68                   filespace, H5P_DEFAULT);

70 H5Dclose(dset_id);
71 /*
72  * Close the file.
73  */
74 H5Fclose(file_id);
```

F90 Example

Parallel Dataset Create

```
43 CALL h5fcreate_f(filename, H5F_ACC_TRUNC_F, file_id,  
    error, access_prp = plist_id)  
73 CALL h5screate_simple_f(rank, dimsf, filespace, error)  
76 !  
77 ! Create the dataset with default properties.  
78 !  
79 CALL h5dcreate_f(file_id, "dataset1", H5T_NATIVE_INTEGER,  
    filespace, dset_id, error)  
  
90 !  
91 ! Close the dataset.  
92 CALL h5dclose_f(dset_id, error)  
93 !  
94 ! Close the file.  
95 CALL h5fclose_f(file_id, error)
```

Accessing a Dataset

- All processes that have opened dataset may do collective I/O
- Each process may do independent and arbitrary number of data I/O access calls
 - _C: H5Dwrite and H5Dread
 - _F90: h5dwrite_f and h5dread_f

Accessing a Dataset Programming model

- Create and set dataset transfer property
 - _C: H5Pset_dxpl_mpio**
 - _H5FD_MPIO_COLLECTIVE
 - _H5FD_MPIO_INDEPENDENT (default)
 - _F90: h5pset_dxpl_mpio_f**
 - _H5FD_MPIO_COLLECTIVE_F
 - _H5FD_MPIO_INDEPENDENT_F (default)
- Access dataset with the defined transfer property

C Example: Collective write

```
95  /*
96   * Create property list for collective dataset write.
97   */
98  plist_id = H5Pcreate(H5P_DATASET_XFER);
99  H5Pset_dxpl_mpio(plist_id, H5FD_MPIO_COLLECTIVE);
100
101  status = H5Dwrite(dset_id, H5T_NATIVE_INT,
102                  memspace, filespace, plist_id, data);
```

F90 Example: Collective write

```
88  ! Create property list for collective dataset write
89  !
90  CALL h5pcreate_f(H5P_DATASET_XFER_F, plist_id, error)
91  CALL h5pset_dxpl_mpio_f(plist_id, &
                           H5FD_MPIO_COLLECTIVE_F, error)
92
93  !
94  ! Write the dataset collectively.
95  !
96  CALL h5dwrite_f(dset_id, H5T_NATIVE_INTEGER, data, &
                  error, &
                  file_space_id = filespace, &
                  mem_space_id = memspace, &
                  xfer_prp = plist_id)
```


Writing and Reading Hyperslabs

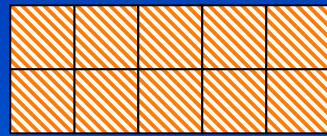
Programming model

- Distributed memory model: data is split among processes
- PHDF5 uses hyperslab model
- Each process defines memory and file hyperslabs
- Each process executes partial write/read call
 - _ Collective calls
 - _ Independent calls

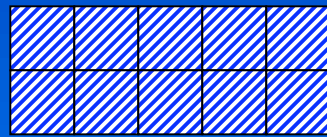
Hyperslab Example 1

Writing dataset by rows

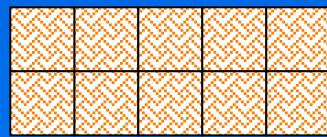
P0



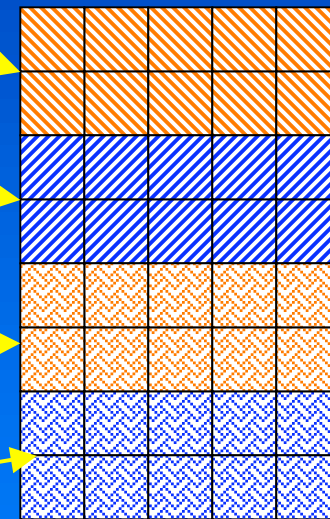
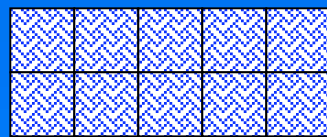
P1



P2



P3



File

Writing by rows

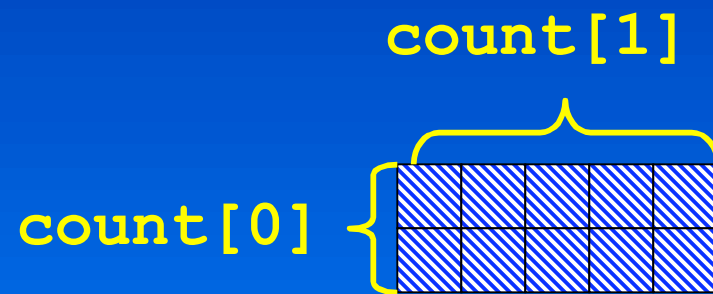
Output of h5dump utility

```
HDF5 "SDS_row.h5" {
GROUP "/" {
  DATASET "IntArray" {
    DATATYPE  H5T_STD_I32BE
    DATASPACE SIMPLE { ( 8, 5 ) / ( 8, 5 ) }
    DATA {
      10, 10, 10, 10, 10,
      10, 10, 10, 10, 10,
      11, 11, 11, 11, 11,
      11, 11, 11, 11, 11,
      12, 12, 12, 12, 12,
      12, 12, 12, 12, 12,
      13, 13, 13, 13, 13,
      13, 13, 13, 13, 13
    }
  }
}
}
```

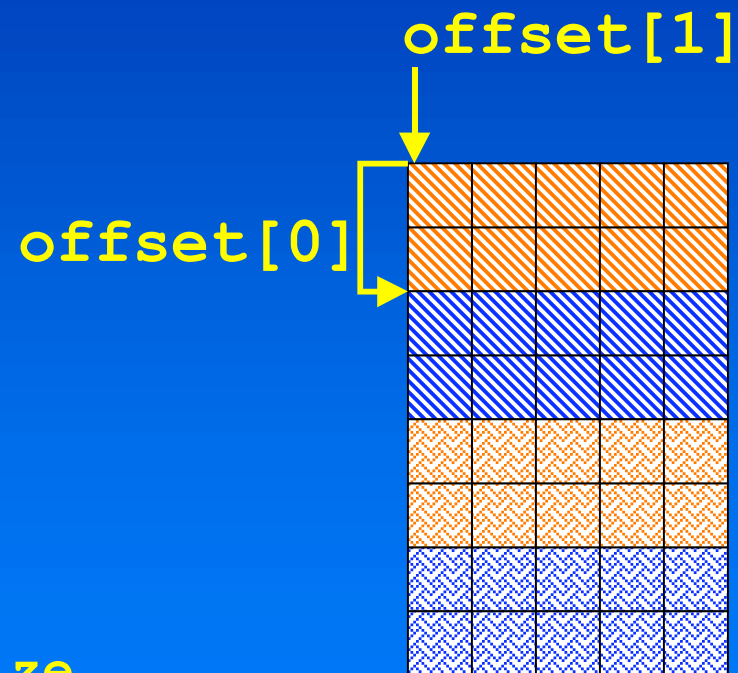
Example 1

Writing dataset by rows

P1 (memory space)



File



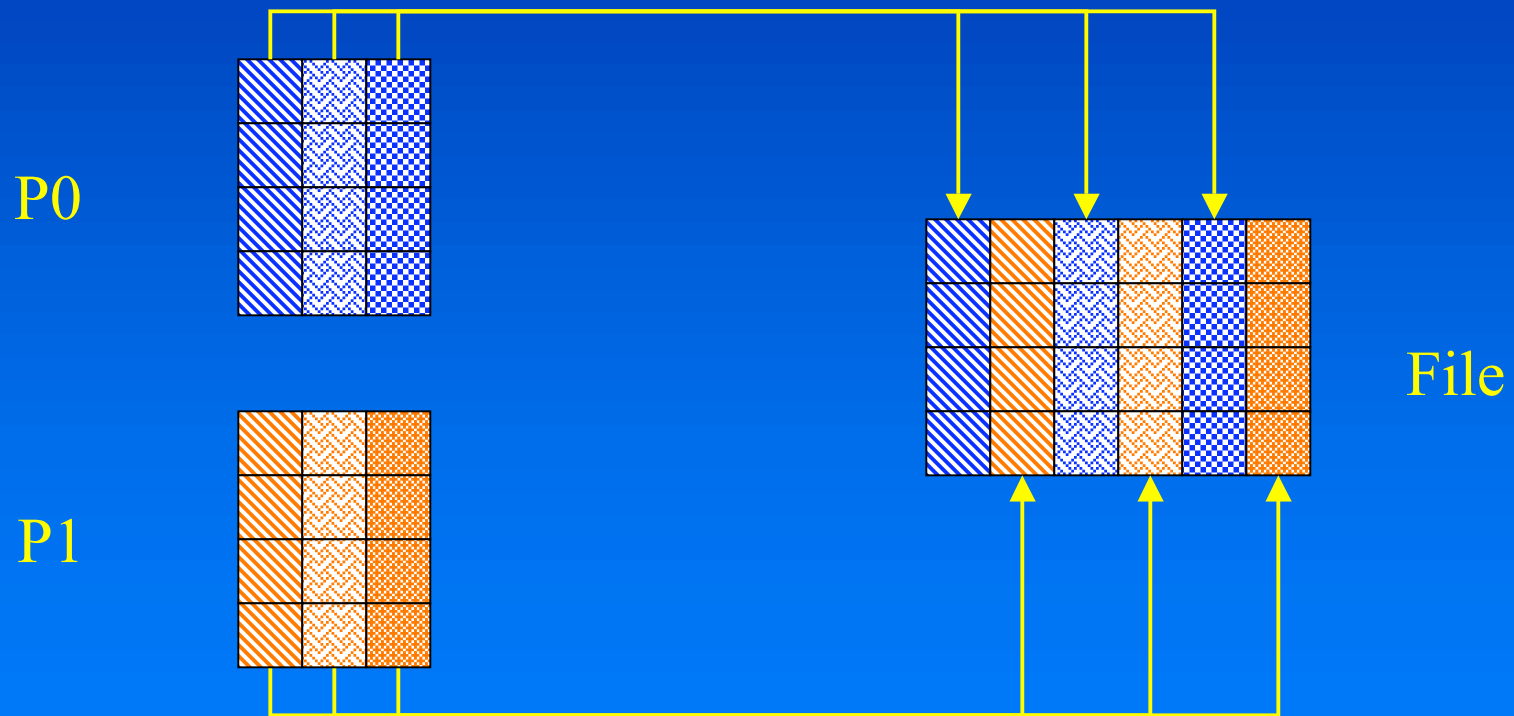
```
count[0] = dimsf[0]/mpi_size  
count[1] = dimsf[1];  
offset[0] = mpi_rank * count[0]; /* = 2 */  
offset[1] = 0;
```

C Example 1

```
71  /*
72   * Each process defines dataset in memory and
73   * writes it to the hyperslab
74   * in the file.
75   */
76   count[0] = dimsf[0]/mpi_size;
77   count[1] = dimsf[1];
78   offset[0] = mpi_rank * count[0];
79   offset[1] = 0;
80   memspace = H5Screate_simple(RANK, count, NULL);
81  /*
82   * Select hyperslab in the file.
83   */
84   filespace = H5Dget_space(dset_id);
85   H5Sselect_hyperslab(filespace,
86                      H5S_SELECT_SET, offset, NULL, count, NULL);
```

Hyperslab Example 2

Writing dataset by columns



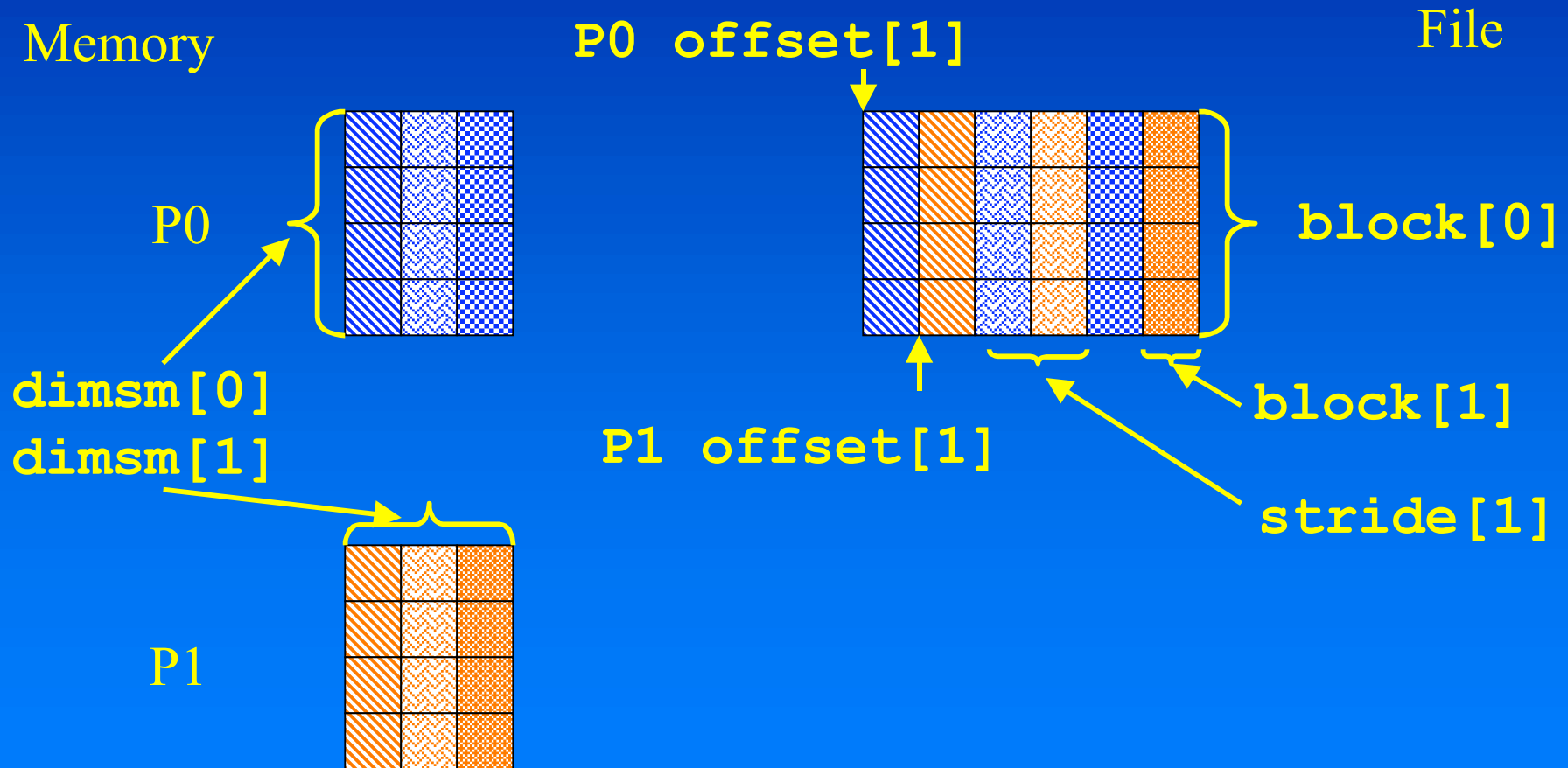
Writing by columns

Output of h5dump utility

```
HDF5 "SDS_col.h5" {
GROUP "/" {
  DATASET "IntArray" {
    DATATYPE  H5T_STD_I32BE
    DATASPACE SIMPLE { ( 8, 6 ) / ( 8, 6 ) }
    DATA {
      1, 2, 10, 20, 100, 200,
      1, 2, 10, 20, 100, 200,
      1, 2, 10, 20, 100, 200,
      1, 2, 10, 20, 100, 200,
      1, 2, 10, 20, 100, 200,
      1, 2, 10, 20, 100, 200,
      1, 2, 10, 20, 100, 200,
      1, 2, 10, 20, 100, 200,
      1, 2, 10, 20, 100, 200
    }
  }
}
}
```

Example 2

Writing Dataset by Column



C Example 2

```
85      /*
86      * Each process defines hyperslab in
87      * the file
88      */
89      count[0] = 1;
90      count[1] = dimsm[1];
91      offset[0] = 0;
92      offset[1] = mpi_rank;
93      stride[0] = 1;
94      stride[1] = 2;
95      block[0] = dimsf[0];
96      block[1] = 1;
97
98      /*
99      * Each process selects hyperslab.
100     */
101     filespace = H5Dget_space(dset_id);
102     H5Sselect_hyperslab(filespace,
        H5S_SELECT_SET, offset, stride,
        count, block);
```

Hyperslab Example 3

Writing dataset by pattern

P0



P1



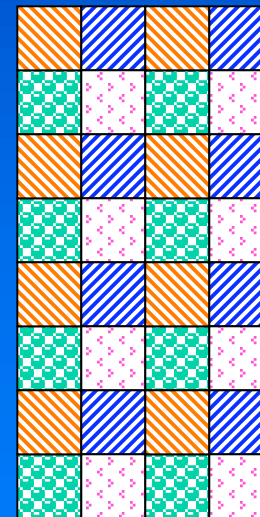
P2



P3



File



Writing by Pattern

Output of h5dump utility

```
HDF5 "SDS_pat.h5" {
GROUP "/" {
  DATASET "IntArray" {
    DATATYPE  H5T_STD_I32BE
    DATASPACE  SIMPLE { ( 8, 4 ) / ( 8, 4 ) }
    DATA {
      1, 3, 1, 3,
      2, 4, 2, 4,
      1, 3, 1, 3,
      2, 4, 2, 4,
      1, 3, 1, 3,
      2, 4, 2, 4,
      1, 3, 1, 3,
      2, 4, 2, 4
    }
  }
}
}
```

Example 3

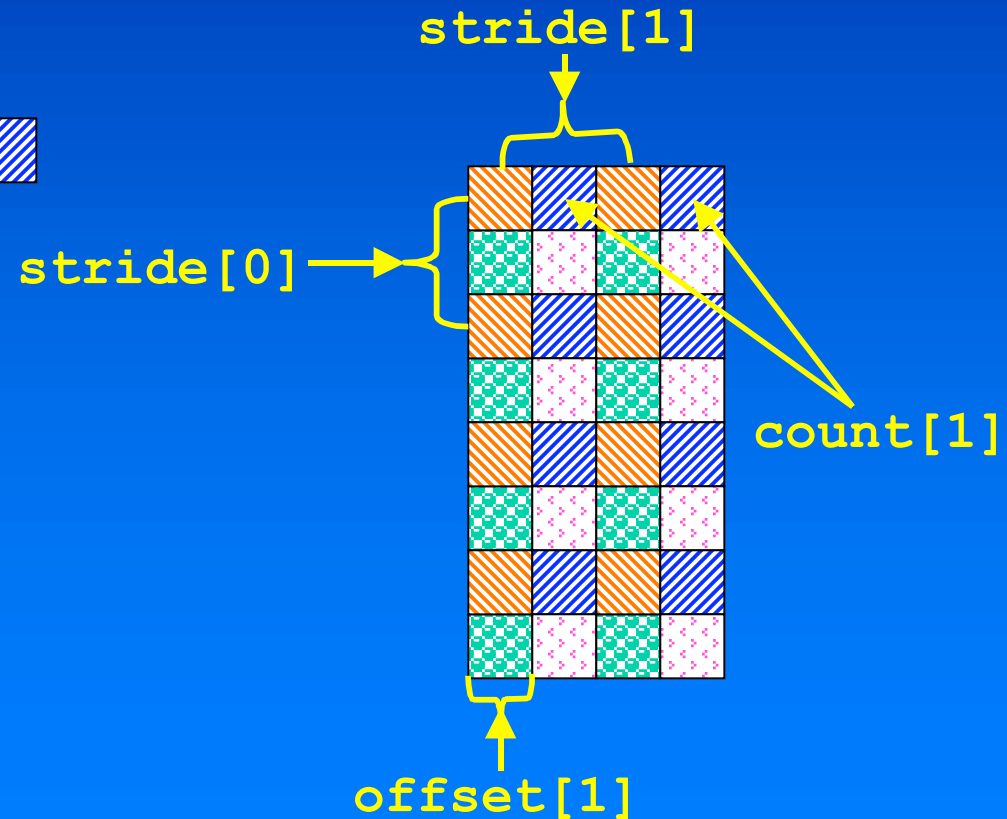
Writing dataset by pattern

Memory

P2



File



```
offset[0] = 0;  
offset[1] = 1;  
count[0] = 4;  
count[1] = 2;  
stride[0] = 2;  
stride[1] = 2;
```

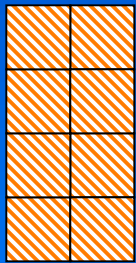
C Example 3: Writing by pattern

```
90     /* Each process defines dataset in memory and
91     * writes it to the hyperslab
92     * in the file.
93     */
94     count[0] = 4;
95     count[1] = 2;
96     stride[0] = 2;
97     stride[1] = 2;
98     if(mpi_rank == 0) {
99         offset[0] = 0;
100        offset[1] = 0;
101    }
102    if(mpi_rank == 1) {
103        offset[0] = 1;
104        offset[1] = 0;
105    }
106    if(mpi_rank == 2) {
107        offset[0] = 0;
108        offset[1] = 1;
109    }
110    if(mpi_rank == 3) {
111        offset[0] = 1;
112        offset[1] = 1;
113    }
```

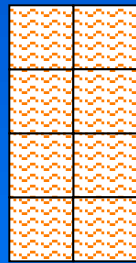
Hyperslab Example 4

Writing dataset by chunks

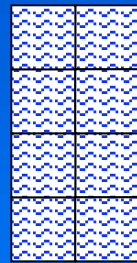
P0



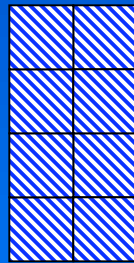
P1



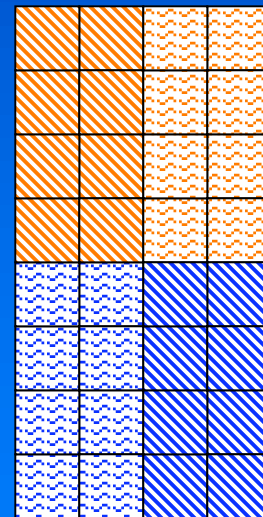
P2



P3



File



Writing by Chunks

Output of h5dump utility

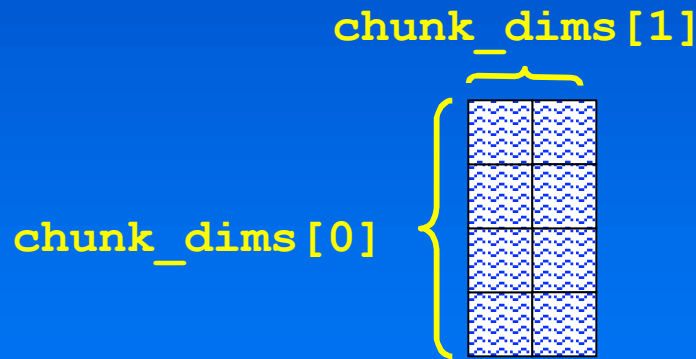
```
HDF5 "SDS_chnk.h5" {
GROUP "/" {
  DATASET "IntArray" {
    DATATYPE  H5T_STD_I32BE
    DATASPACE SIMPLE { ( 8, 4 ) / ( 8, 4 ) }
    DATA {
      1, 1, 2, 2,
      1, 1, 2, 2,
      1, 1, 2, 2,
      1, 1, 2, 2,
      3, 3, 4, 4,
      3, 3, 4, 4,
      3, 3, 4, 4,
      3, 3, 4, 4
    }
  }
}
}
```

Example 4

Writing dataset by chunks

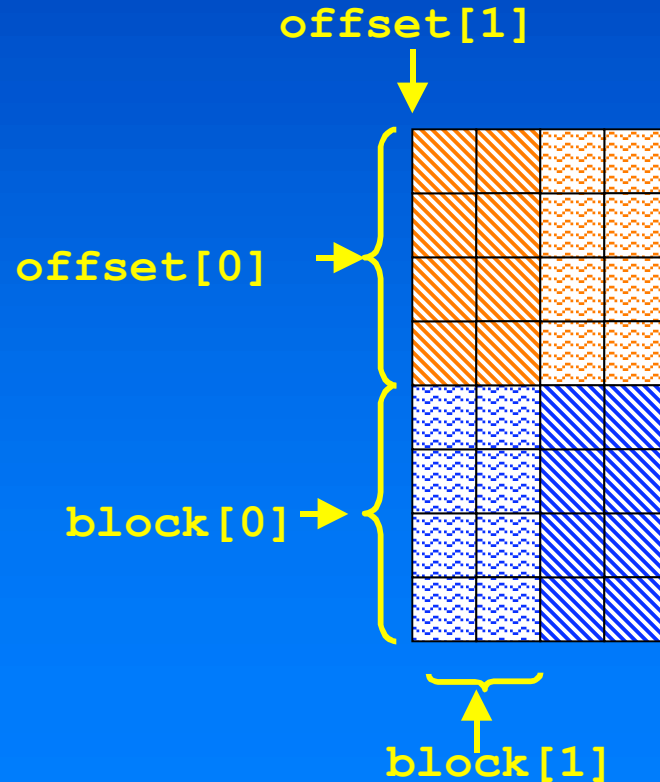
Memory

P2



```
block[0] = chunk_dims[0];  
block[1] = chunk_dims[1];  
offset[0] = chunk_dims[0];  
offset[1] = 0;
```

File



C Example 4

Writing by chunks

```
97     count[0] = 1;
98     count[1] = 1 ;
99     stride[0] = 1;
100    stride[1] = 1;
101    block[0] = chunk_dims[0];
102    block[1] = chunk_dims[1];
103    if(mpi_rank == 0) {
104        offset[0] = 0;
105        offset[1] = 0;
106    }
107    if(mpi_rank == 1) {
108        offset[0] = 0;
109        offset[1] = chunk_dims[1];
110    }
111    if(mpi_rank == 2) {
112        offset[0] = chunk_dims[0];
113        offset[1] = 0;
114    }
115    if(mpi_rank == 3) {
116        offset[0] = chunk_dims[0];
117        offset[1] = chunk_dims[1];
118    }
```

Performance Tuning in HDF5

Two Sets of Tuning Knobs

- File level knobs
 - _ Apply to the entire file
- Data transfer level knobs
 - _ Apply to individual dataset read or write

File Level Knobs

- H5Pset_meta_block_size
- H5Pset_alignment
- H5Pset_fapl_split
- H5Pset_cache
- H5Pset_fapl_mpio

H5Pset_meta_block_size

- Sets the minimum metadata block size allocated for metadata aggregation.
- Aggregated block is usually written in a single write action
- Default is 2KB
- *Pro:*
 - _ Larger block size reduces I/O requests
- *Con:*
 - _ Could create “holes” in the file and make file bigger

H5Pset_meta_block_size

- When to use:
- File is open for a long time and
 - _ A lot of objects created
 - _ A lot of operations on the objects performed
 - _ As a result metadata is interleaved with raw data
 - _ A lot of new metadata (attributes)

H5Pset_alignment

- Sets two parameters
 - _ Threshold
 - Minimum size of object for alignment to take effect
 - Default 1 byte
 - _ Alignment
 - Allocate object at the next multiple of alignment
 - Default 1 byte
- Example: (threshold, alignment) = (1024, 4K)
 - _ All objects of 1024 or more bytes starts at the boundary of 4KB

H5Pset_alignment Benefits

- In general, the default (no alignment) is good for single process serial access since the OS already manages buffering.
- For some parallel file systems such as GPFS, an alignment of the disk block size improves I/O speeds.
- *Con: File may be bigger*

H5Pset_fapl_split

- HDF5 splits to two files
 - _Metadata file for metadata
 - _Rawdata file for raw data (array data)
 - _Two files represent one logical HDF5 file
- *Pro:* Significant I/O improvement if
 - _metadata file is stored in Unix file systems (good for many small I/O)
 - _raw data file is stored in Parallel file systems (good for large I/O).

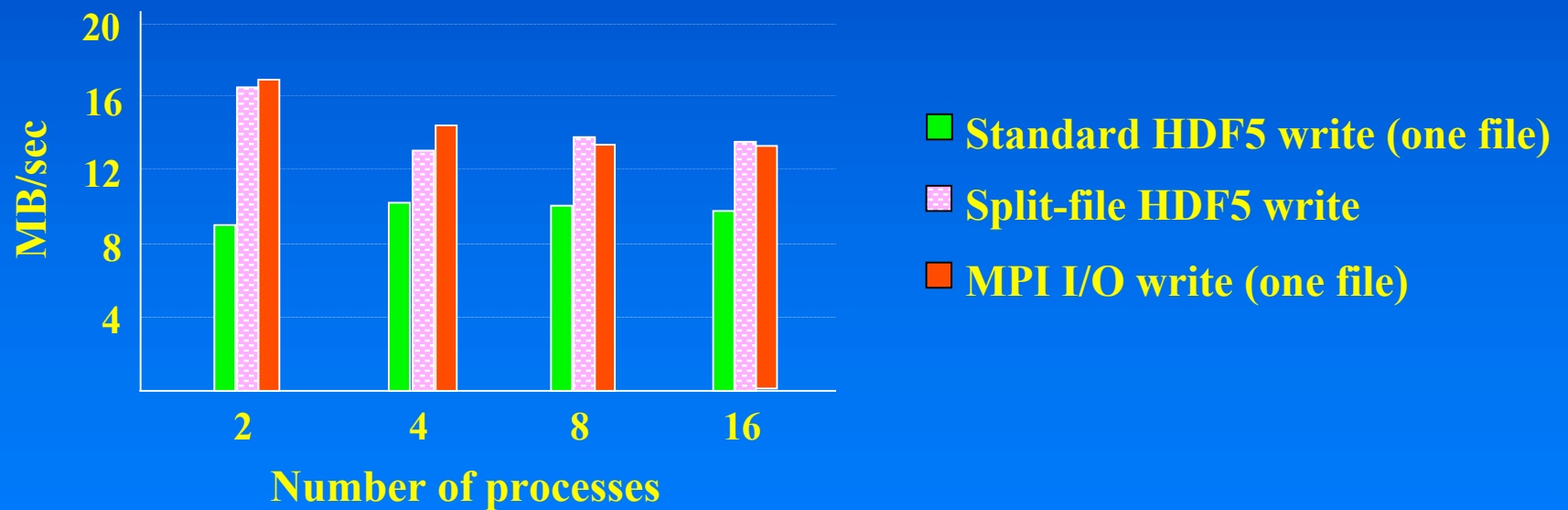
H5Pset_fapl_split

- *Con:*
 - _ Both files should be “kept together” for integrity of the HDF5 file
 - _ Can be a potential problem when files are moved to another platform or file system

Write speeds of Standard vs. Split-file HDF5 vs. MPI-IO

Results for ASCI Red machine at Sandia National Laboratory

• Each process writes 10MB of array data



H5Pset_cache

- Sets:
 - _ The number of elements (objects) in the meta data cache
 - _ The number of elements, the total number of bytes, and the preemption policy value (default is 0.75) in the raw data chunk cache

H5Pset_cache (cont.)

- Preemption policy:
 - _ Chunks are stored in the list with the most recently accessed chunk at the end
 - _ Least recently accessed chunks are at the beginning of the list
 - _ $X \times 100\%$ of the list is searched for the fully read/written chunk; X is called preemption value, where X is between 0 and 1
 - _ If chunk is found then it is deleted from cache, if not then first chunk in the list is deleted

H5Pset_cache (cont.)

- The right values of N
 - _ May improve I/O performance by controlling preemption policy
 - _ 0 value forces to delete the “oldest” chunk from cache
 - _ 1 value forces to search all list for the chunk that will be unlikely accessed
 - _ Depends on application access pattern

Chunk Cache Effect by H5Pset_cache

- Write one integer dataset
256x256x1024 (256MB)
- Using chunks of 256x16x1024
(16MB)
- Two tests of
 - _Default chunk cache size (1MB)
 - _Set chunk cache size 16MB

Chunk Cache Time Definitions

- Total
 - _ Time to open file, write dataset, close dataset and close file
- Dataset write
 - _ Time to write the whole dataset
- Chunk write
 - _ Time to write a chunk
- User time/System time
 - _ Total Unix user/system time of test

Chunk Cache Size Results

Cache buffer size (MB)	Chunk write time (sec)	Dataset write time (sec)	Total time (sec)	User time (sec)	System time (sec)
1	132.58	2450.25	2453.09	14.00	2200.10
16	0.376	7.83	8.27	6.21	3.45

Chunk Cache Size Summary

- Big chunk cache size improves performance
- Poor performance mostly due to increased system time
 - _ Many more I/O requests
 - _ Smaller I/O requests

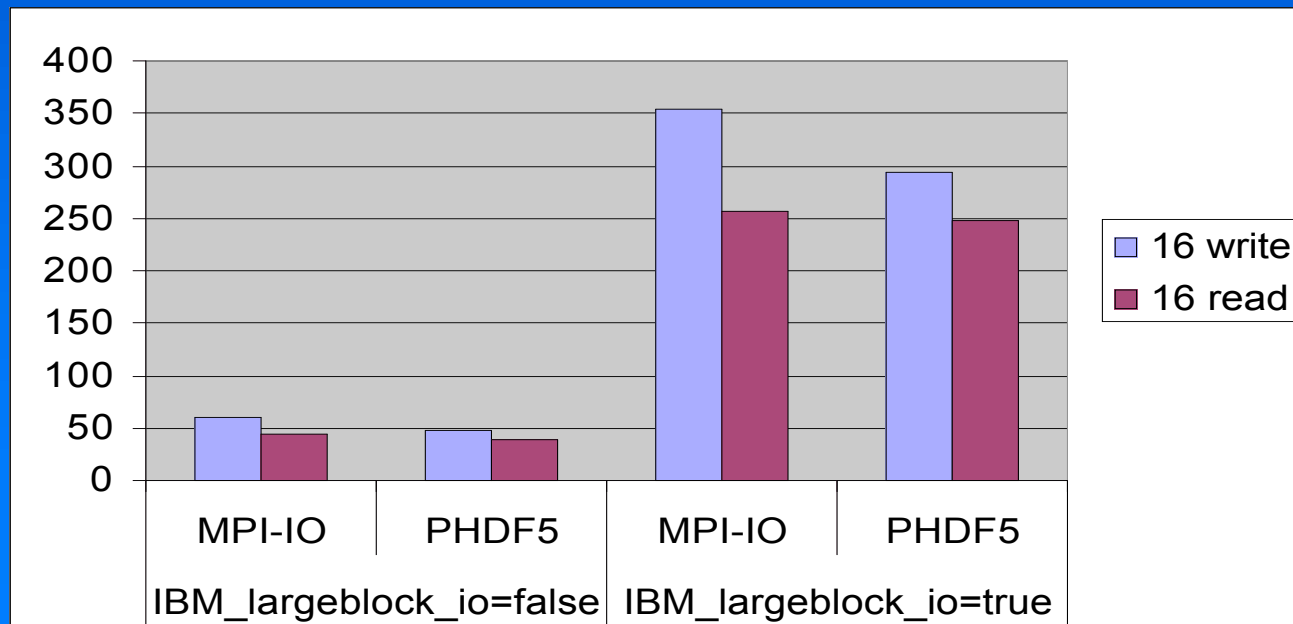
I/O Hints via H5Pset_fapl_mpio

- MPI-IO hints can be passed to the MPI-IO layer via the Info parameter of H5Pset_fapl_mpio
- Examples
 - _ Telling Romio to use 2-phases I/O speeds up collective I/O in the ASCI Red machine
 - _ Setting IBM_largeblock_io=true speeds up GPFS write speeds

Effects of I/O Hints

IBM_largeblock_io

- GPFS at Livermore National Laboratory
ASCI Blue machine
 - _ 4 nodes, 16 tasks
 - _ Total data size 1024MB
 - _ I/O buffer size 1MB



Effects of I/O Hints IBM_largeblock_io

- GPFS at LLNL Blue
 - _ 4 nodes, 16 tasks
 - _ Total data size 1024MB
 - _ I/O buffer size 1MB

		IBM_largeblock_io=false		IBM_largeblock_io=true	
Tasks		MPI-IO	PHDF5	MPI-IO	PHDF5
16	write	60	48	354	294
16	read	44	39	256	248

Data Transfer Level Knobs

- H5Pset_buffer
- H5Pset_sieve_buf_size

H5Pset_buffer

- Sets size of the internal buffers used during data transfer
- Default is 1 MB
- Pro:
 - _ Bigger size improves performance
- Con:
 - _ Library uses more memory

H5Pset_buffer

- When should be used:
 - _Datatype conversion
 - _Data gathering-scattering (e.g. checkerboard dataspace selection)

H5Pset_sieve_buf_size

- Sets the size of the data sieve buffer
- Default is 64KB
- Sieve buffer is a buffer in memory that holds part of the dataset raw data
- During I/O operations data is replaced in the buffer first, then one big I/O request occurs

H5Pset_sieve_buf_size

- Pro:
 - _ Bigger size reduces I/O requests issued for raw data access
- Con:
 - _ Library uses more memory
- When to use:
 - _ Data scattering-gathering (e.g. checker board)
 - _ Interleaved hyperslabs

Parallel I/O Benchmark Tool

- h5perf
 - _ Benchmark test I/O performance
- Four kinds of API
 - _ Parallel HDF5
 - _ MPI-IO
 - _ Native parallel (e.g., gpfs, pvfs)
 - _ POSIX (open, close, lseek, read, write)

Useful Parallel HDF Links

- Parallel HDF information site
_ http://hdf.ncsa.uiuc.edu/Parallel_HDF/
- Parallel HDF mailing list
_ hdfparallel@ncsa.uiuc.edu
- Parallel HDF5 tutorial available at
_ <http://hdf.ncsa.uiuc.edu/HDF5/doc/Tutor>