

Optimising I/O on the Cray XE6

Cray XE6 Performance Workshop
11-12 July

Primary File Systems on HECToR

- There are two primary types of filesystem on HECToR

Home Space

- Mounted as /home (NFS)
- Smaller size (70TB total)
- Fully and regularly backed up
- Only available on the login nodes

Designed and optimised for compilation, editing, long term storage of critical files.

Work Space

- Mounted as /work (Lustre)
- Large size (1.02PB total)
- **IT IS NOT BACKED UP**
- Available on the compute and login nodes

Designed and optimised for scratch large files and high bandwidth transfers (e.g. scientific output, restart files)

There are no local disks on the compute nodes

HECToR Archives

- **HECToR Archive**
 - **Tape archive system available to users**
 - **1.02 PB of Tape storage available**
- **Scheduled archiving occurs at 22:00 each evening**
 - Files in the ARCH and ARCH2 directories are written to tape
 - Once successfully copied they are then removed.
 - Files in ARCH have a single copy written to tape
 - Files in ARCH2 have two copies written to different sets of tapes
 - Using ARCH2 counts double towards quotas
 - We recommend using ARCH2 to avoid the risk that tapes cannot be restored.
- **Files can be actively restored or archived at other times from the command line.**
- **Users wishing to store critical large files should use the archive!**

A Very Short Overview of the Archiver CLI - 1

- **Quotas - cresquota**

```
ted@hector-xe6-4:~ > cresquota
```

User Name	Group Name	User (GB)	Group (GB)	Quota (GB)
ted	y02	248.9	389.0	6000
ted	y05	0.0	0.0 *	NONE

- **Archiving - cresarc**

```
ted@hector-xe6-4:~WORK > cresarc -1 example-file
```

```
WARNING: Path "example-file" is actually "/esfs2/y02/y02/ted/example-file".
```

```
Successful "/esfs2/y02/y02/ted/example-file"
```

```
Archive successful (this Archive was NOT duplicated and is at a higher risk of loss).
```

A Very Short Overview of the Archiver CLI - 2

- **Listing - creslist**

```
ted@hector-xe6-4:~WORK > creslist
Archive Time      Keyword  Mod Time      Path
-----
07/03/2012 11:39          Jul 03 11:38:27 /esfs2/y02/y02/ted/example-file
```

- **Restoring - cresrestore**

```
ted@hector-xe6-4:~WORK > cresrestore /esfs2/y02/y02/ted/example-file
Restore from backup performed on "07/03/2012 11:39:01" successful
```

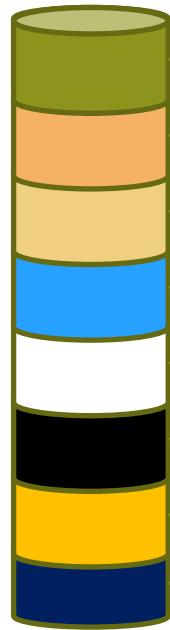
Understanding Parallel Filesystems

Concepts for reading or writing files to lustre

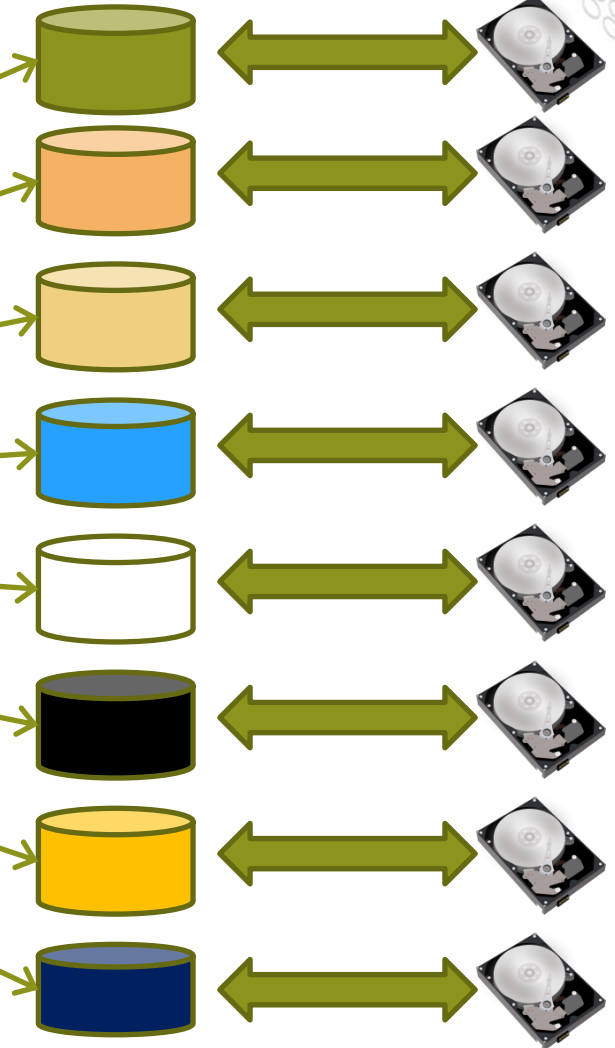
File System Fundamentals



Single Logical File
e.g. /work/example

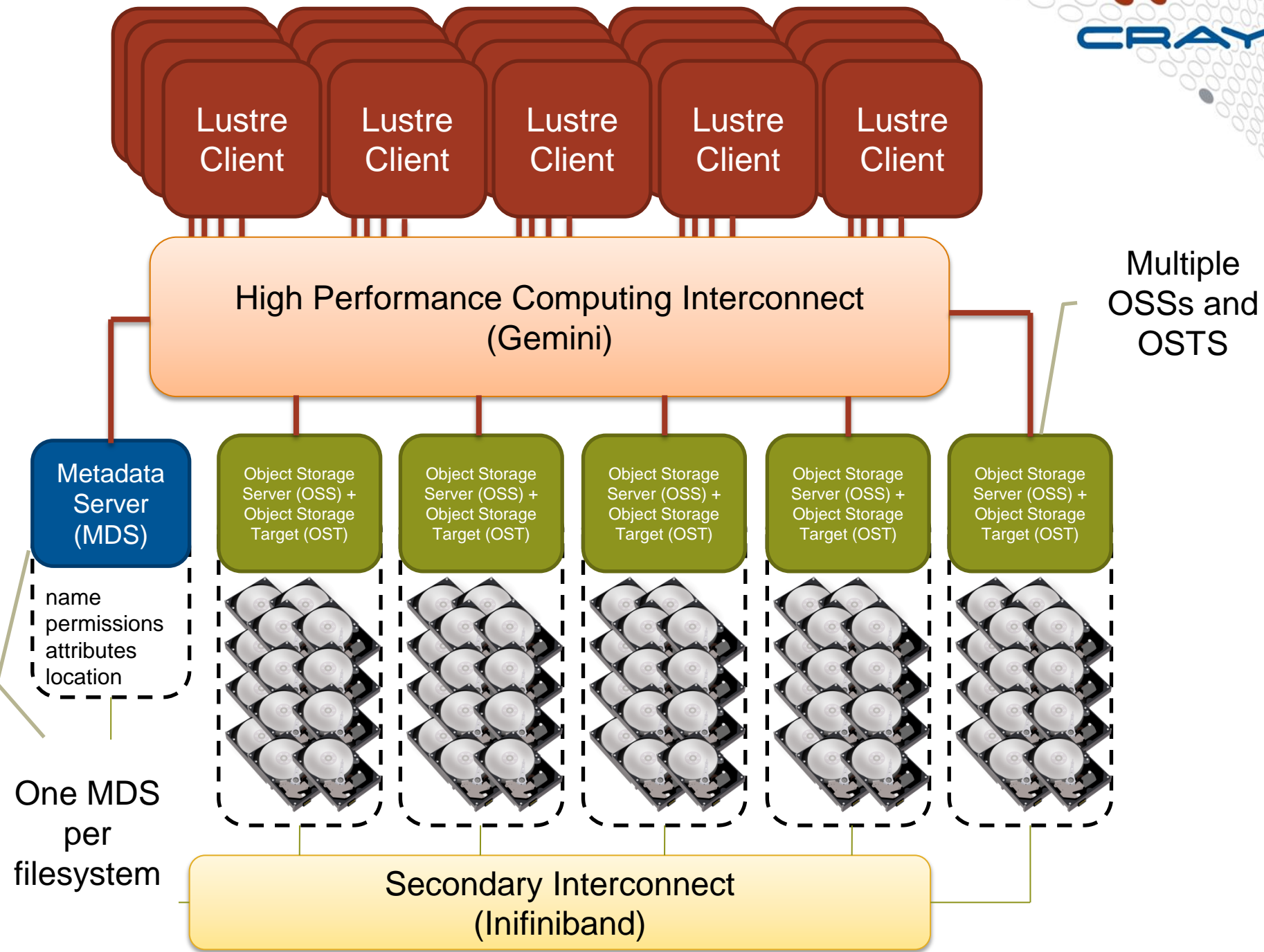


File automatically
divided into stripes

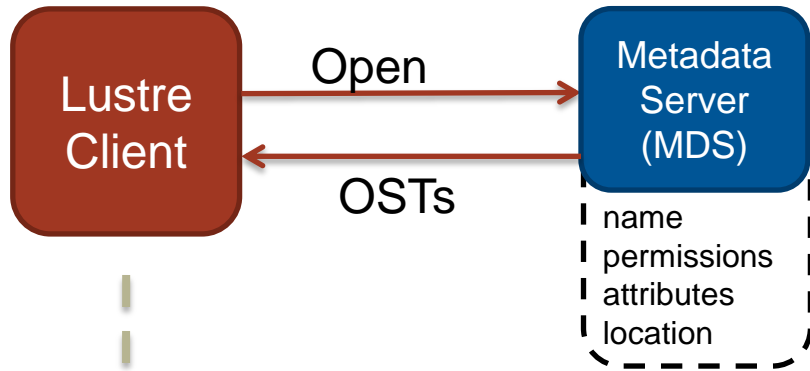


Stripes are written/read
from across multiple drives

- **A scalable cluster file system for Linux**
 - Developed by Cluster File Systems -> Sun -> Oracle.
 - Name derives from “Linux Cluster”
 - The Lustre file system consists of software subsystems, storage, and an associated network
- **MDS – metadata server**
 - Handles information about files and directories
- **OSS – Object Storage Server**
 - **The hardware entity**
 - The server node
 - Support multiple OSTs
- **OST – Object Storage Target**
 - **The software entity**
 - This is the software interface to the backend volume



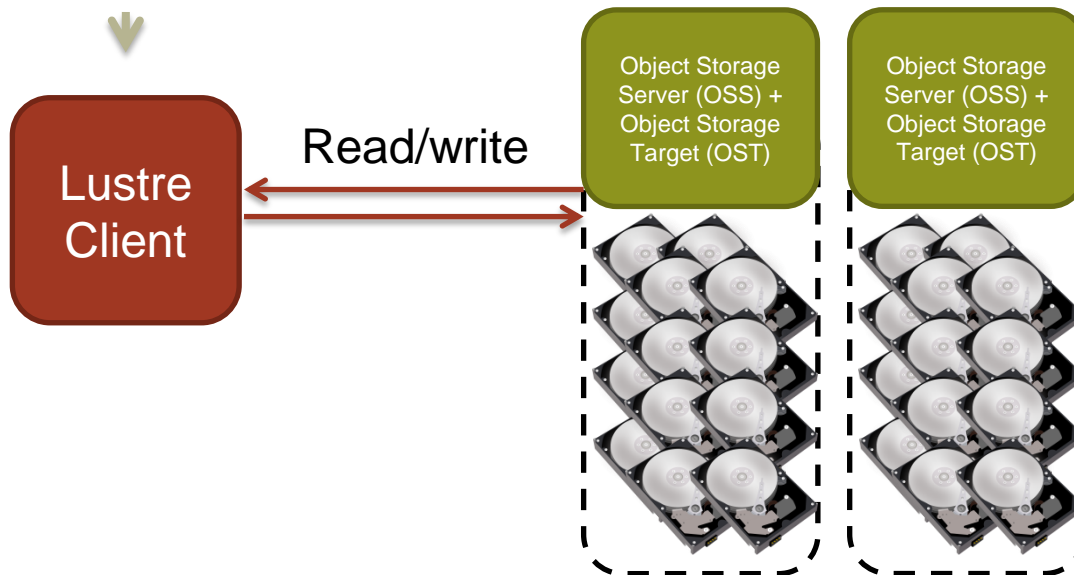
Opening a file



The client sends a request to the MDS to opening/acquiring information about the file

The MDS then passes back a list of OSTs

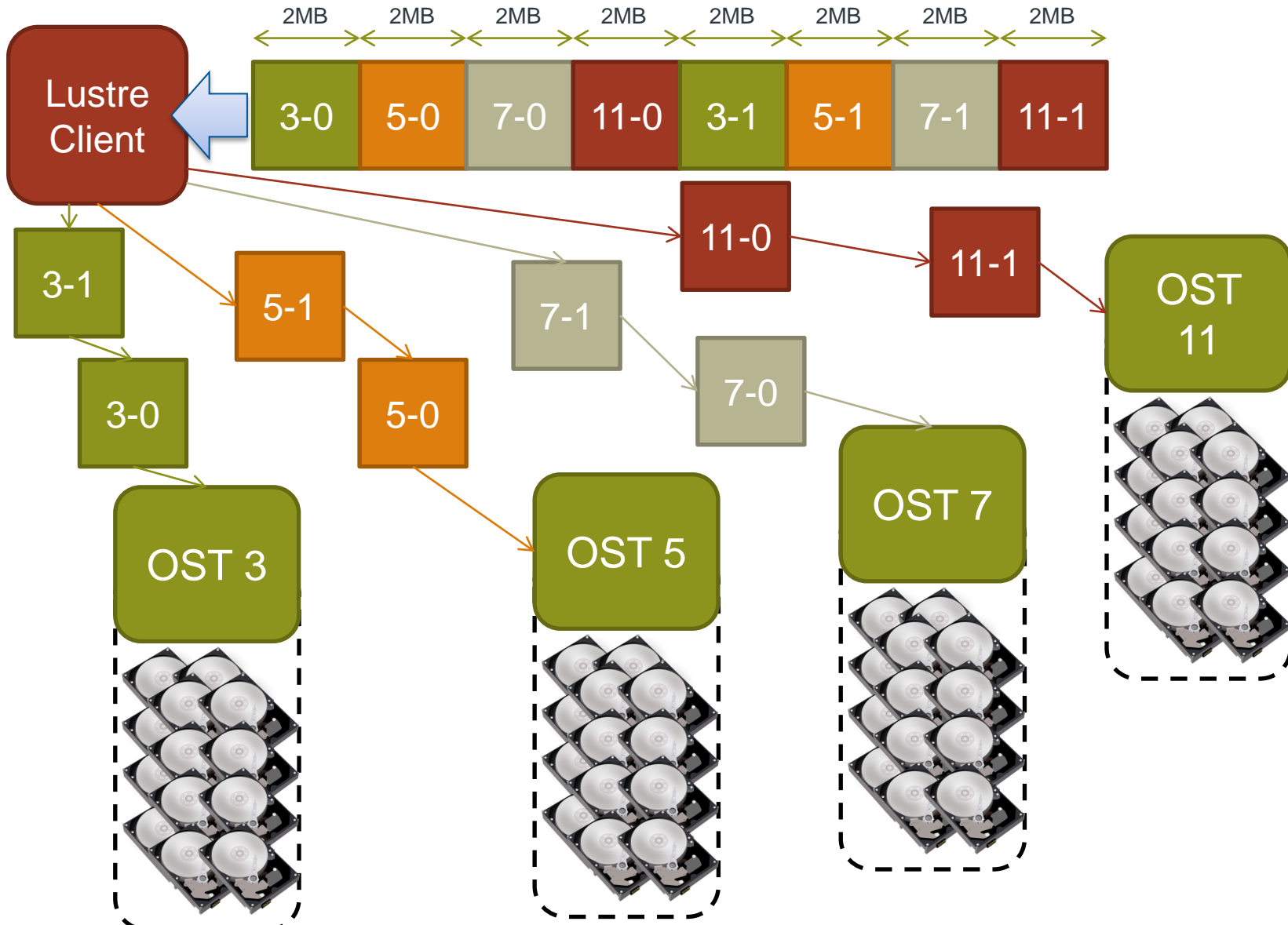
- For an existing file, these contain the data stripes
- For a new files, these typically contain a randomly assigned list of OSTs where data is to be stored



Once a file has been opened no further communication is required between the client and the MDS

All transfer is directly between the assigned OSTs and the client

File decomposition – 2 Megabyte Stripes



Controlling Lustre Striping

- **lfs** - the lustre utility for setting the stripe properties of new files, or displaying the striping patterns of existing.
- The most used options are :
 - setstripe – Set striping properties of a directory or new file
 - getstripe – Return information on current striping settings
 - osts – List the number of OSTs associated with this file system
 - df – Show disk usage of this file system
- For help execute lfs without any arguments

```
$ lfs
lfs > help
Available commands are:
    setstripe
    find
    getstripe
    check
.....
```

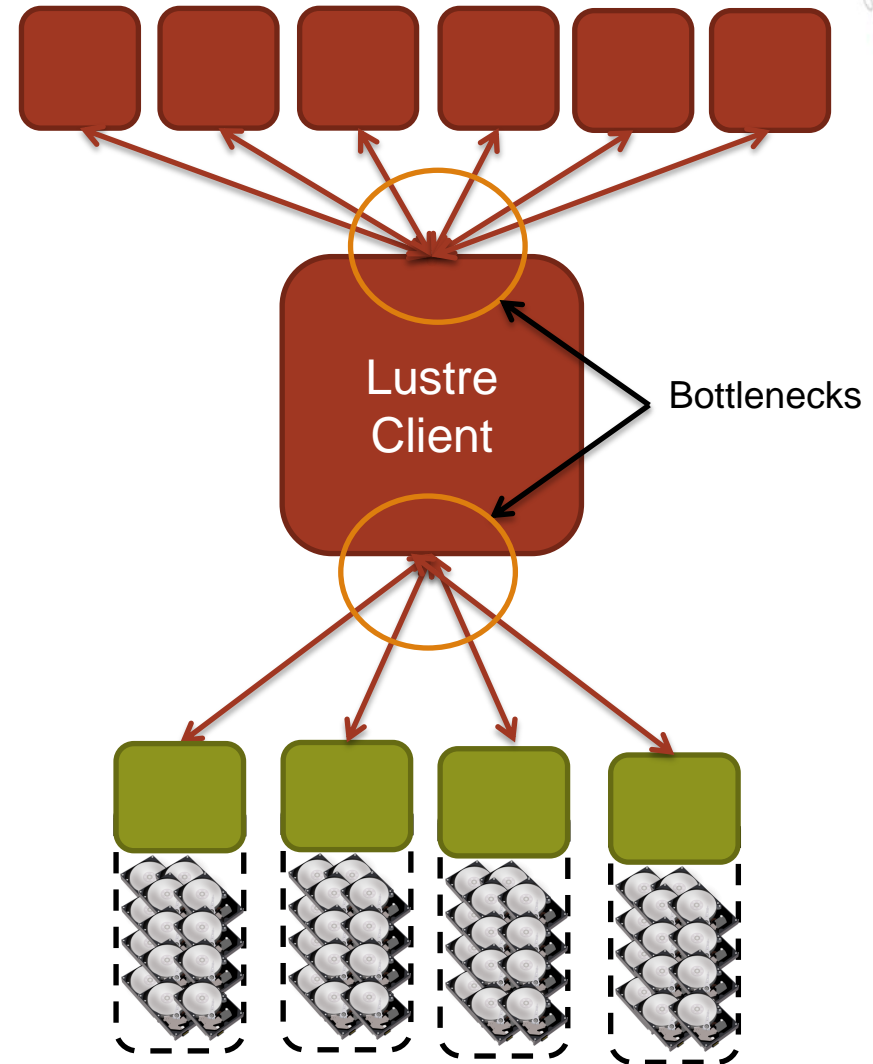
lfs setstripe

- **Sets the stripe for a file or a directory**
- `lfs setstripe <file|dir> <-s size> <-i start> <-c count>`
 - size: Number of bytes on each OST (0 filesystem default)
 - start: OST index of first stripe (-1 filesystem default)
 - count: Number of OSTs to stripe over (0 default, -1 all)
- **Comments**
 - Can use lfs to create an empty file with the stripes you want (like the touch command)
 - Can apply striping settings to a directory, any children will inherit parent's stripe settings on creation.
 - The stripes of a file is given when the file is created. It is not possible to change it afterwards.

Common I/O Paradigms

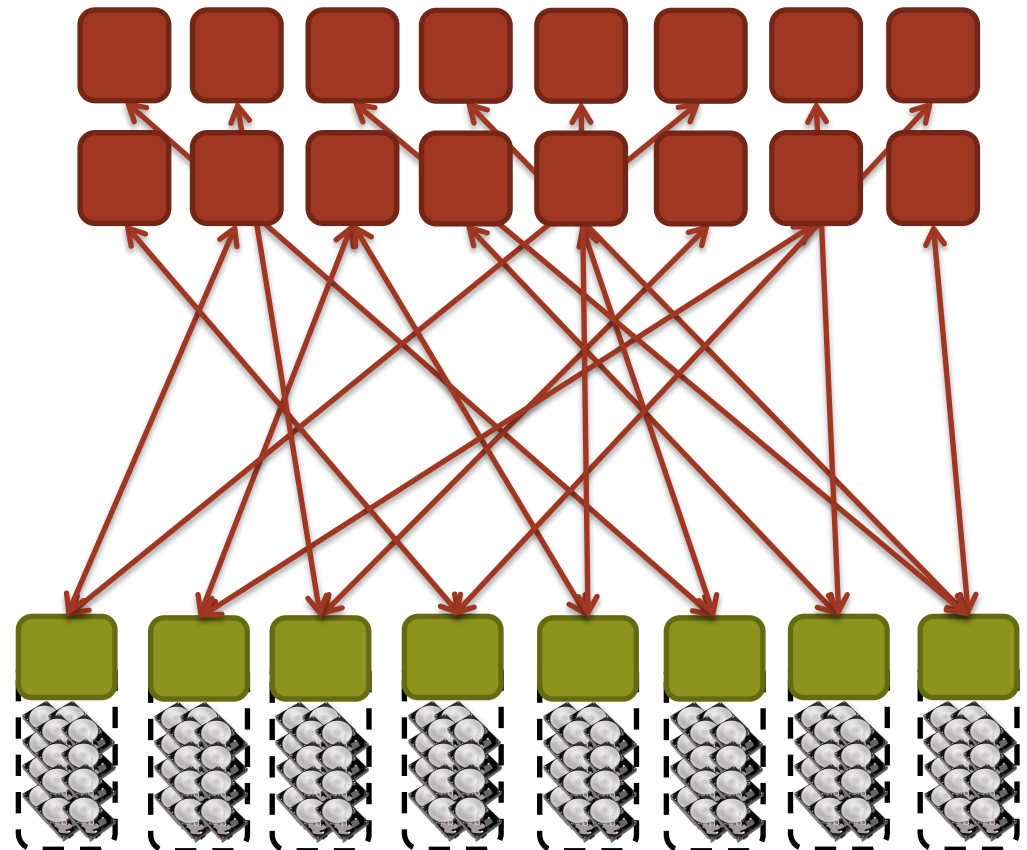
Spokesperson

- **One process performs I/O.**
 - Data Aggregation or Duplication
 - Limited by single I/O process.
- **Easy to program**
- **Pattern does not scale.**
 - Time increases linearly with amount of data.
 - Time increases with number of processes.
- **Care has to be taken when doing the „all to one“-kind of communication at scale**
- **Can be used for a dedicated IO Server (not easy to program)**



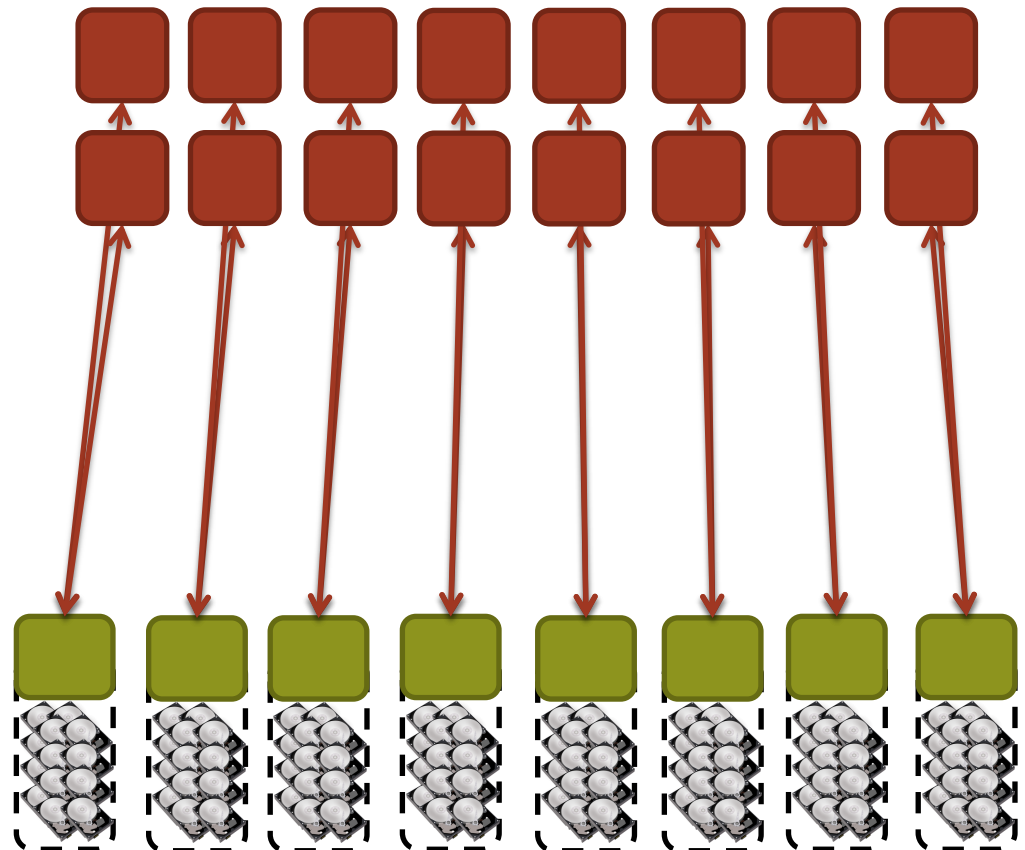
Multiple Writers – Multiple Files

- **All processes perform I/O to individual files.**
 - Limited by file system.
- **Easy to program**
 - Requires job to always run on the same number of cores
- **Pattern does not scale at large process counts.**
 - Number of files creates bottleneck with metadata operations.
 - Number of simultaneous disk accesses creates contention for file system resources.



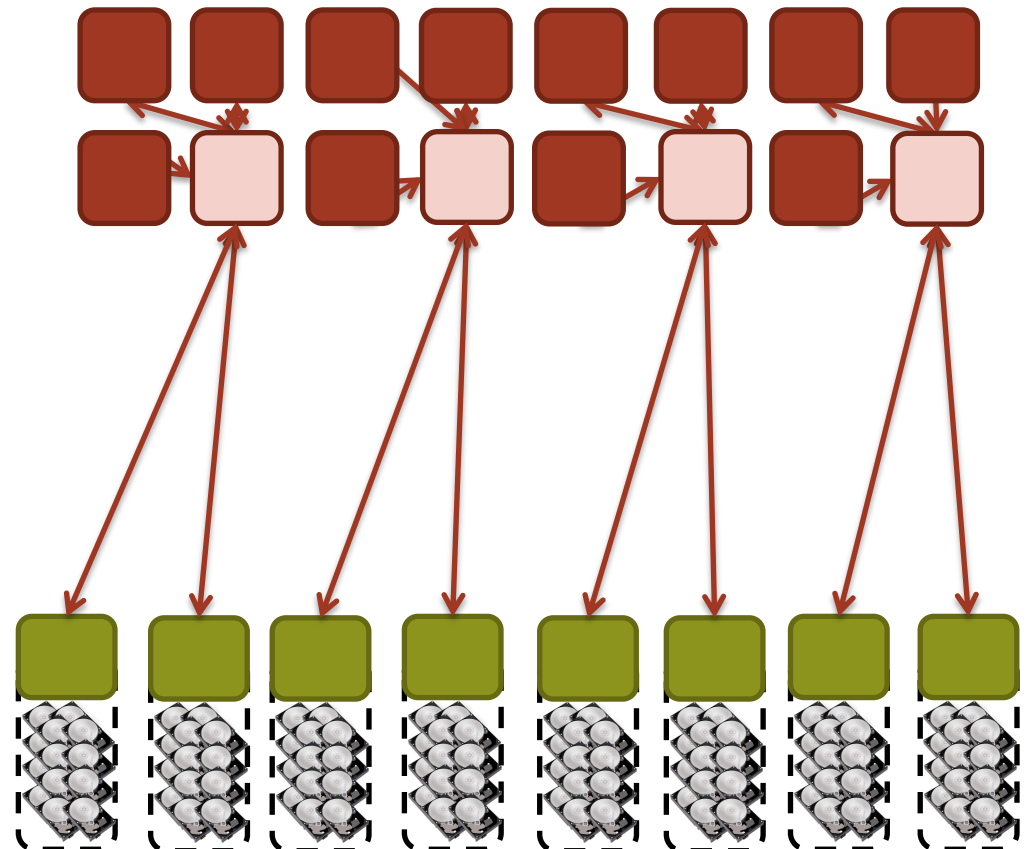
Multiple Writers – Single File

- Each process performs I/O to a single file which is shared.
- Performance
 - Data layout within the shared file is very important.
 - At large process counts contention can build for file system resources.
- Not all programming languages support it
 - C/C++ can work with fseek
 - No real Fortran standard



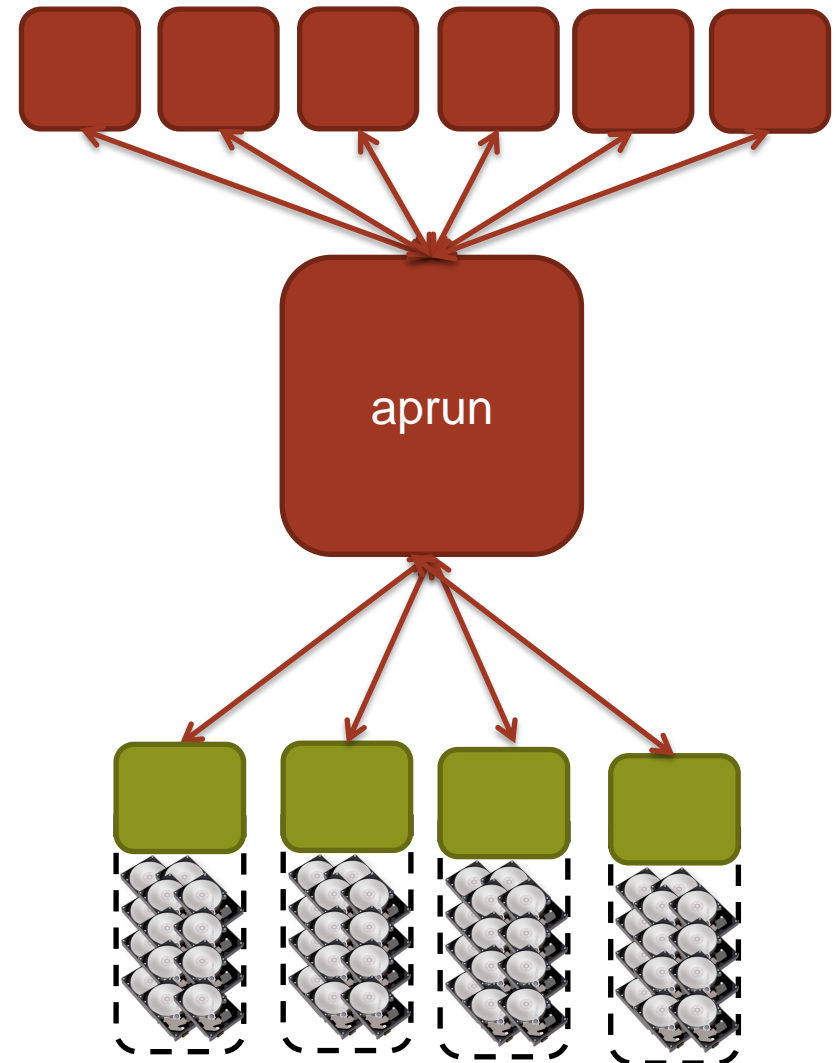
Collective IO to single or multiple files

- **Aggregation to a processor in a group which processes the data.**
 - Serializes I/O in group.
- **I/O process may access independent files.**
 - Limits the number of files accessed.
- **Group of processes perform parallel I/O to a shared file.**
 - Increases the number of shares to increase file system usage.
 - Decreases number of processes which access a shared file to decrease file system contention.



Special Case : Standard Output and Error

- Standard Output and Error streams are effectively serial I/O.
- All STDIN, STDOUT, and STDERR I/O serialize through aprun
- Disable debugging messages when running in production mode.
 - “Hello, I’m task 32,000!”
 - “Task 64,000, made it through loop.”
 - ...



CRAY IO Software stack

Application

HDF5

NETCDF

MPI-IO

POSIX I/O

Lustre File System

Optimising IO in Applications



1. Tune your parameters

Any easy and non-invasive approach

Select best striping values

- **Selecting the striping values will have an impact on the I/O performance of your application**
- **Rule of thumb :**
 1. $\# \text{ files} > \# \text{ OSTs} \Rightarrow$ Set `stripe_count=1`
You will reduce the lustre contention and OST file locking this way and gain performance
 2. $\# \text{ files} == 1 \Rightarrow$ Set `stripe_count=#OSTs`
Assuming you have more than 1 I/O client
 3. $\# \text{ files} < \# \text{ OSTs} \Rightarrow$ Select `stripe_count` so that you use all OSTs
Example : You have 8 OSTs and write 4 files at the same time, then select `stripe_count=2`

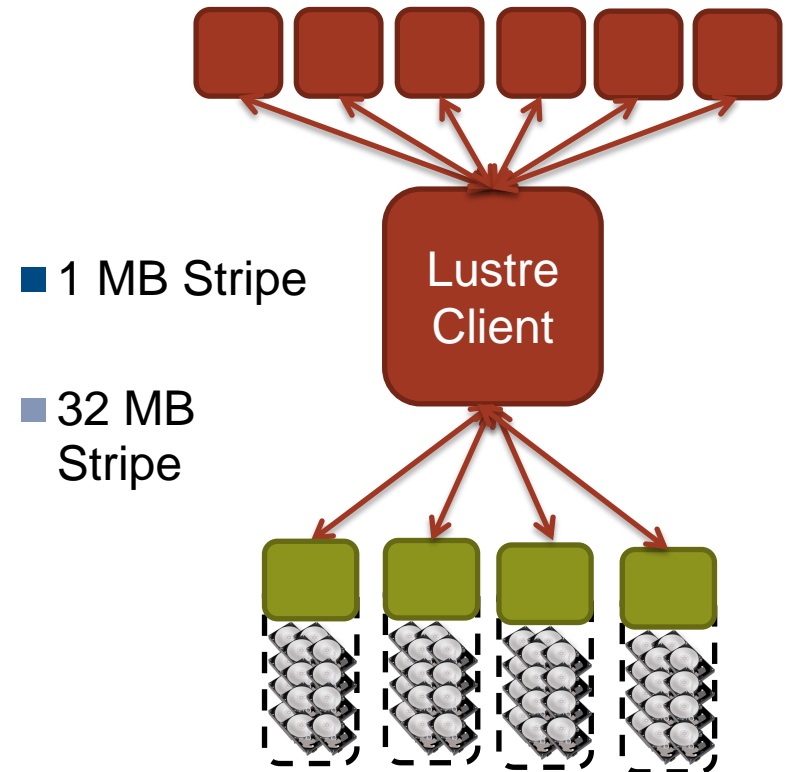
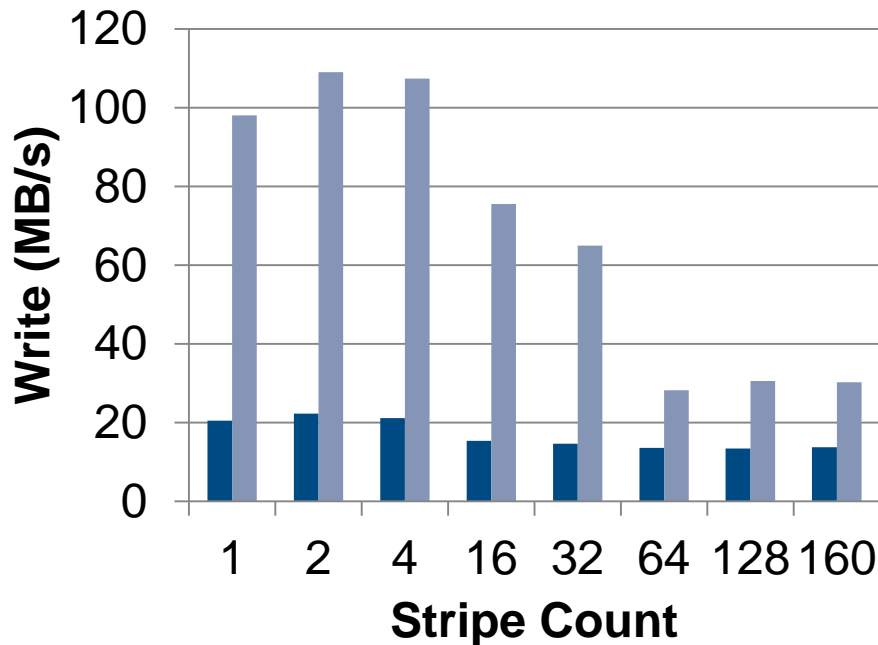
Always allow the system to choose OSTs at random!

(Ensures even loading of the OSTs and prevents accidental contention)

Case Study 1 : Spokesman

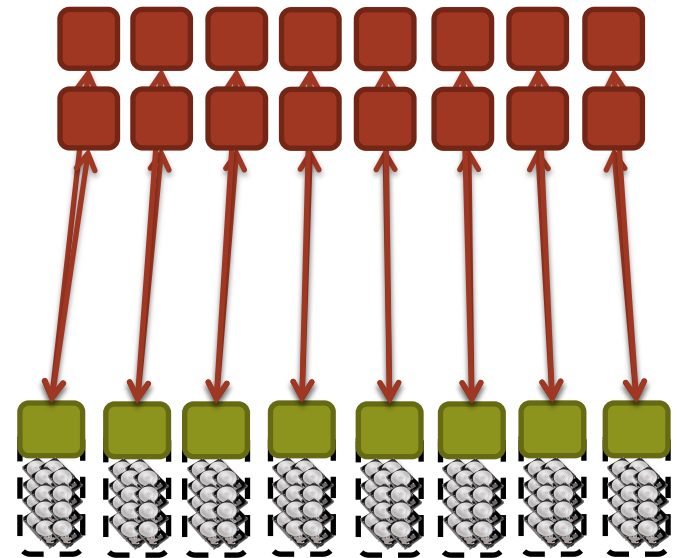
- **32 MB per OST (32 MB – 5 GB) and 32 MB Transfer Size**
 - Unable to take advantage of file system parallelism
 - Access to multiple disks adds overhead which hurts performance
 - Note : XE6 numbers might be better

Single Writer Write Performance



Case Study 2 : Parallel I/O into a single file

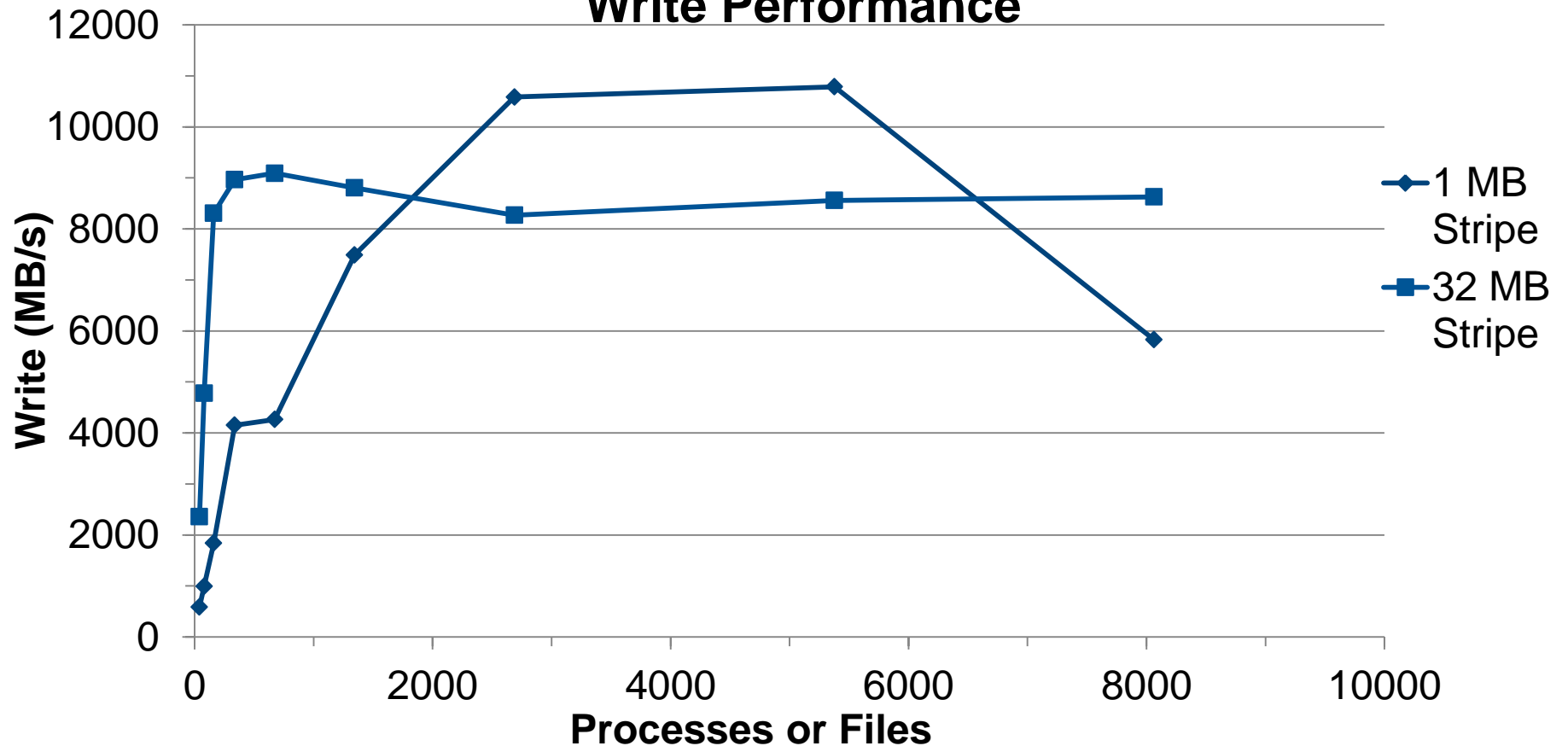
- A particular code both reads and writes a 377 GB file.
Runs on 6000 cores.
 - Total I/O volume (reads and writes) is 850 GB.
 - Utilizes parallel HDF5
- **Default Stripe settings: count =4, size=1M, index =-1.**
 - 1800 s run time (~ 30 minutes)
- **Stripe settings: count=-1, size=1M, index =-1.**
 - 625 s run time (~ 10 minutes)
- **Results**
 - 66% decrease in run time.



Case Study 3 : Single File Per Process

- 128 MB per file and a 32 MB Transfer size, each file has a stripe_count of 1

**File Per Process
Write Performance**



2. Try to hide the I/O

Asynchronous I/O

Majority of data is output, allow computation to overlap
Double buffer arrays to allow computation to continue while data flushed to disk

1. Use asynchronous POSIX calls

- Only covers the I/O call itself, any packing/gathering/encoding still has to be done by the compute processors
- Not currently supported by Lustre (calls become synchronous)

2. Use 3rd party libraries

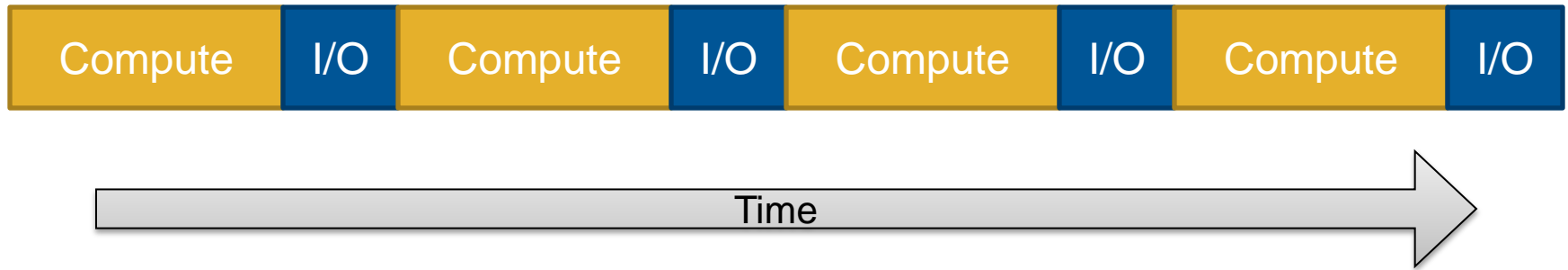
- Typical examples are MPI-IO (see later)
- Again, packing/gathering/encoding still done by compute processors

3. Add I/O Servers to the application

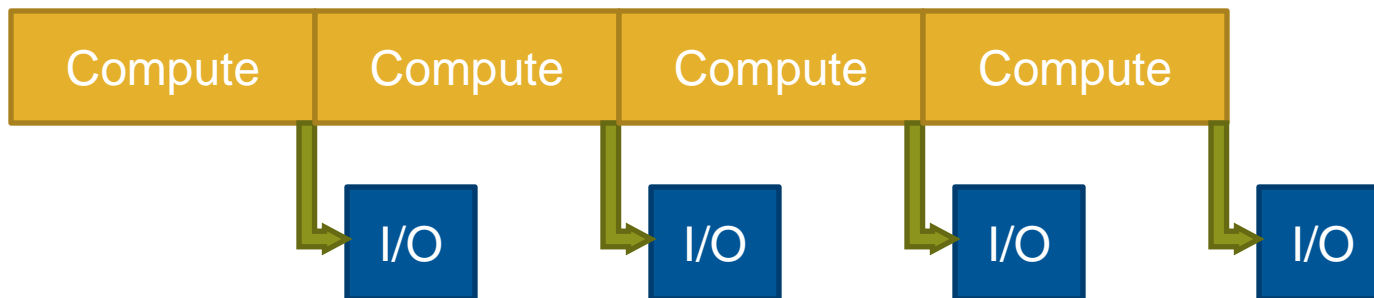
- Add processors dedicated to performing time consuming operations
- More complicated to implement than other solutions
- Portable across platforms (works on any parallel platform)

Asynchronous I/O

Standard Sequential I/O



Asynchronous I/O



Naive IO Server Pseudo Code

User more nodes to act as I/O Servers

Compute Node

```
do i=1,time_steps
  compute(j)
  checkpoint(data)
end do

subroutine checkpoint(data)
  MPI_Wait(send_req)
  buffer = data
  MPI_Isend(IO_SERVER, buffer)
end subroutine
```

I/O Server

```
do i=1,time_steps
  do j=1,compute_nodes
    MPI_Recv(j, buffer)
    write(buffer)
  end do
end do
```

IO Servers

- **Successful strategy deployed in multiple codes.**
- **Strategy has become more successful as number of nodes has increased.**
 - Addition of extra nodes only cost 1-2% in resources
- **Requires additional development that can pay off for codes that generate large files.**
- **Typically still only one or a small number of writers performing I/O operations (not necessarily reaching optimum bandwidth).**

I/O Performance : To keep in mind

- There is no “One Size Fits All” solution to the I/O problem.
- Many I/O patterns work well for some range of parameters.
- Bottlenecks in performance can occur in many locations. (Application and/or File system)
- Going to extremes with an I/O pattern will typically lead to problems.
- I/O is a shared resource. Expect timing variation

3. Parallelise e.g. MPI-IO

Change how the application handles I/O

A simple MPI-IO program in C

```
MPI_File fh;
MPI_Status status;

MPI_Comm_rank(MPI_COMM_WORLD, &rank);
MPI_Comm_size(MPI_COMM_WORLD, &nprocs);
bufsize = FILESIZE/nprocs;
nints = bufsize/sizeof(int);

MPI_File_open(MPI_COMM_WORLD, 'FILE',
              MPI_MODE_RDONLY, MPI_INFO_NULL, &fh);
MPI_File_seek(fh, rank * bufsize, MPI_SEEK_SET);
MPI_File_read(fh, buf, nints, MPI_INT, &status);
MPI_File_close(&fh);
```

And now in Fortran using explicit offsets

```
use mpi ! or include 'mpif.h'
integer status(MPI_STATUS_SIZE)
integer (kind=MPI_OFFSET_KIND) offset ! Note : might be integer*8

call MPI_FILE_OPEN(MPI_COMM_WORLD, 'FILE', &
    MPI_MODE_RDONLY, MPI_INFO_NULL, fh, ierr)
nints = FILESIZE / (nprocs*INTSIZE)
offset = rank * nints * INTSIZE
call MPI_FILE_READ_AT(fh, offset, buf, nints, MPI_INTEGER, status,
    ierr)
call MPI_GET_COUNT(status, MPI_INTEGER, count, ierr)
print *, 'process ', rank, 'read ', count, 'integers'
call MPI_FILE_CLOSE(fh, ierr)
```

- **The *_AT routines are thread safe (seek+IO operation in one call)**

Write instead of Read

- Use `MPI_File_write` or `MPI_File_write_at`
- Use `MPI_MODE_WRONLY` or `MPI_MODE_RDWR` as the flags to `MPI_File_open`
- If the file doesn't exist previously, the flag `MPI_MODE_CREATE` must be passed to `MPI_File_open`
- We can pass multiple flags by using bitwise-or `|` in C, or addition `+` or `IOR` in Fortran
- If not writing to a file, using `MPI_MODE_RDONLY` might have a performance benefit. Try it.

MPI_File_set_view

- **MPI_File_set_view** assigns regions of the file to separate processes
- Specified by a triplet (*displacement, etype, and filetype*) **passed to MPI_File_set_view**
 - *displacement* = number of bytes to be skipped from the start of the file
 - *etype* = basic unit of data access (can be any basic or derived datatype)
 - *filetype* = specifies which portion of the file is visible to the process

- **Example :**

```
MPI_File fh;
for (i=0; i<BUFSIZE; i++) buf[i] = myrank * BUFSIZE + i;
MPI_File_open(MPI_COMM_WORLD, "testfile", MPI_MODE_CREATE |
    MPI_MODE_WRONLY, MPI_INFO_NULL, &fh);
MPI_File_set_view(fh, myrank * BUFSIZE * sizeof(int), MPI_INT,
    MPI_INT, 'native', MPI_INFO_NULL);
MPI_File_write(fh, buf, BUFSIZE, MPI_INT, MPI_STATUS_IGNORE);
MPI_File_close(&fh);
```

MPI_File_set_view (Syntax)

- Describes that part of the file accessed by a single MPI process.
- Arguments to MPI_File_set_view:
 - MPI_File file
 - MPI_Offset disp
 - MPI_Datatype etype
 - MPI_Datatype filetype
 - char *datarep
 - MPI_Info info

Collective I/O with MPI-IO

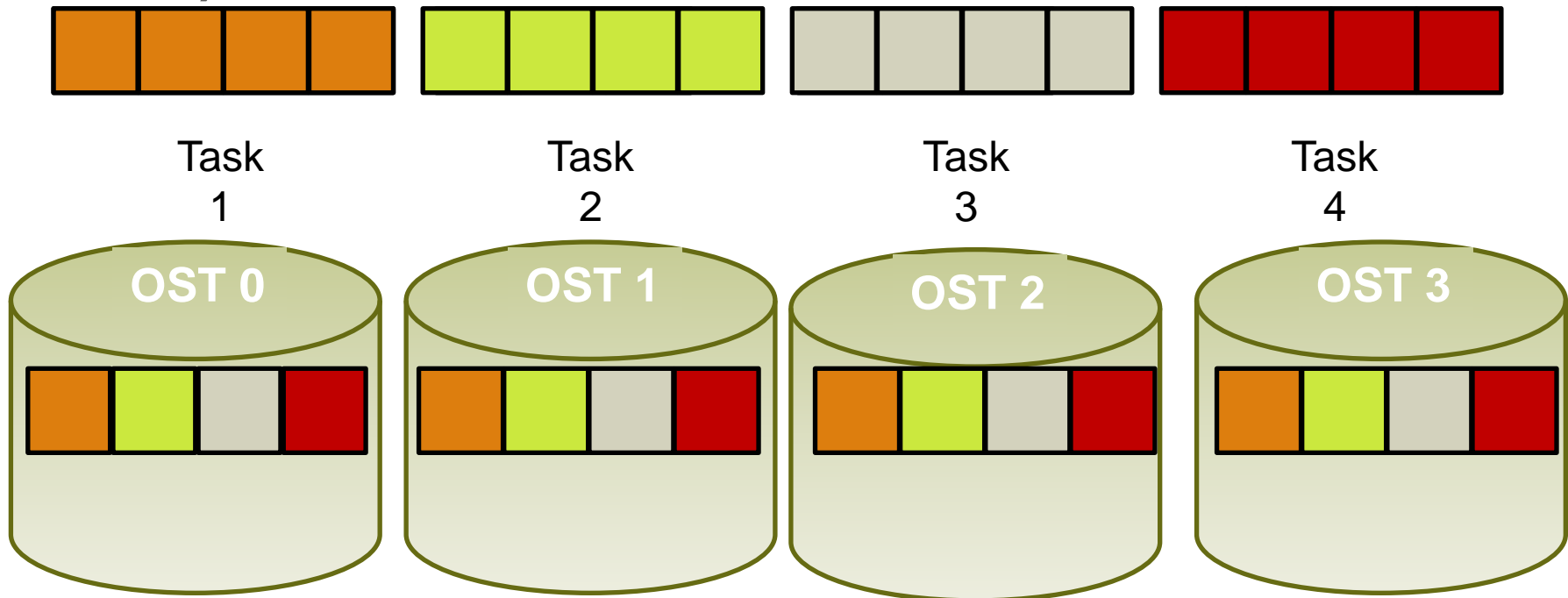
- **MPI_File_read_all, MPI_File_read_at_all, ...**
- **_all indicates that all processes in the group specified by the communicator passed to MPI_File_open will call this function**
- **Each process specifies only its own access information – the argument list is the same as for the non-collective functions**
- **MPI-IO library is given a lot of information in this case:**
 - Collection of processes reading or writing data
 - Structured description of the regions
- **The library has some options for how to use this data**
 - Noncontiguous data access optimizations
 - Collective I/O optimizations

2 Techniques : Sieving and Aggregation

- **Data sieving is used to combine lots of small accesses into a single larger one**
 - Reducing # of operations important (latency)
 - A system buffer/cache is one example
- **Aggregation/Collective Buffering refers to the concept of moving data through intermediate nodes**
 - Different numbers of nodes performing I/O (transparent to the user)
- **Both techniques are used by MPI-IO and triggered with HINTS**

Lustre problem : „OST Sharing“

- A file is written by several tasks :
- The file is stored like this (one single stripe per OST for all tasks) :



- => Performance Problem (like ‚False Sharing‘ in thread programming)
- Flock mount option needed. Only 1 task can write to an OST any time

MPI-IO Interaction with Lustre

- **Included in the Cray MPT library.**
- **Environmental variable used to help MPI-IO optimize I/O performance.**
 - MPICH_MPIIO_CB_ALIGN Environmental Variable. (Default 2)
 - MPICH_MPIIO_HINTS Environmental Variable
 - Can set `striping_factor` and `striping_unit` for files created with MPI-IO.
 - If writes and/or reads utilize collective calls, collective buffering can be utilized (`romio_cb_read/write`) to approximately stripe align I/O within Lustre.
- **HDF5 and NETCDF are both implemented on top of MPI-IO and thus also uses the MPI-IO env. Variables.**

MPICH_MPIIO_CB_ALIGN

- If set to 2, an algorithm is used to divide the I/O workload into Lustre stripe-sized pieces and assigns them to collective buffering nodes (aggregators), so that each aggregator always accesses the same set of stripes and no other aggregator accesses those stripes.
If the overhead associated with dividing the I/O workload can in some cases exceed the time otherwise saved by using this method.
- If set to 1, an algorithm is used that takes into account physical I/O boundaries and the size of I/O requests in order to determine how to divide the I/O workload when collective buffering is enabled.
However, unlike mode 2, there is no fixed association between file stripe and aggregator from one call to the next.
- If set to zero or defined but not assigned a value, an algorithm is used to divide the I/O workload equally amongst all aggregators without regard to physical I/O boundaries or Lustre stripes.

MPI-IO Hints (part 1)

- **MPICH_MPIIO_HINTS_DISPLAY** – Rank 0 displays the name and values of the MPI-IO hints
- **MPICH_MPIO_HINTS** – Sets the MPI-IO hints for files opened with the `MPI_File_Open` routine
 - Overrides any values set in the application by the `MPI_Info_set` routine
 - Following hints supported:

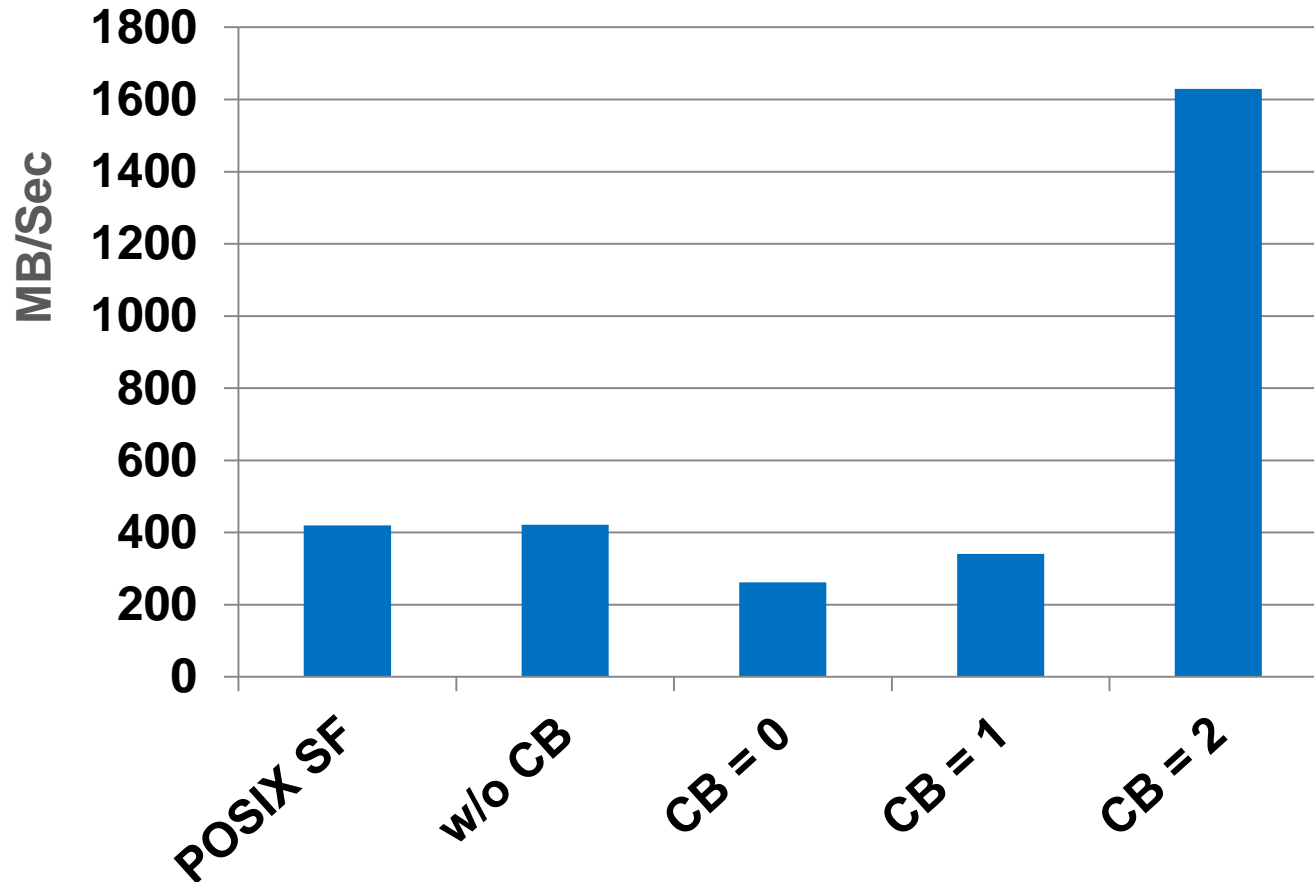
direct_io	cb_nodes	romio_ds_write
romio_cb_read	cb_config_list	ind_rd_buffer_size
romio_cb_write	romio_no_indep_rw	Ind_wr_buffer_size
cb_buffer_size	romio_ds_read	striping_factor
		striping_unit

Env. Variable MPICH_MPIO_HINTS (part 2)

- If set, override the default value of one or more MPI I/O hints. This also overrides any values that were set by using calls to `MPI_Info_set` in the application code. The new values apply to the file the next time it is opened using a `MPI_File_open()` call.
- After the `MPI_File_open()` call, subsequent `MPI_Info_set` calls can be used to pass new MPI I/O hints that take precedence over some of the environment variable values.
Other MPI I/O hints such as `striping_factor`, `striping_unit`, `cb_nodes`, and `cb_config_list` cannot be changed after the `MPI_File_open()` call, as these are evaluated and applied only during the file open process.
- The syntax for this environment variable is a comma-separated list of specifications. Each individual specification is a `pathname_pattern` followed by a colon-separated list of one or more `key=value` pairs. In each `key=value` pair, the key is the MPI-IO hint name, and the value is its value as it would be coded for an `MPI_Info_set` library call.
- Example:
`MPICH_MPIO_HINTS=file1:direct_io=true,file2:romio_ds_write=disable,/scratch/user/me/dump.*:romio_cb_write=enable:cb_nodes=8`

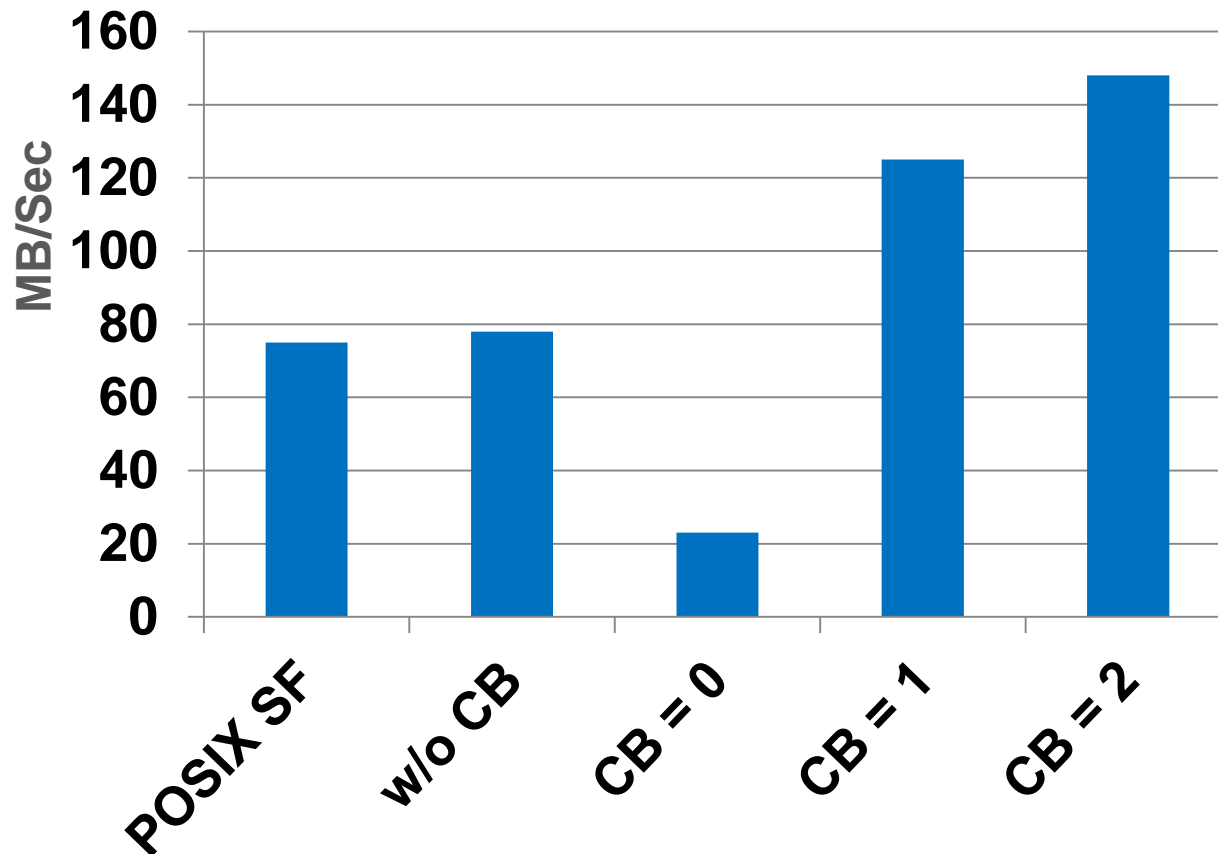
IOR benchmark 1,000,000 bytes

MPI-IO API , non-power-of-2 blocks and transfers, in this case blocks and transfers both of 1M bytes and a strided access pattern. Tested on an XT5 with 32 PEs, 8 cores/node, 16 stripes, 16 aggregators, 3220 segments, 96 GB file



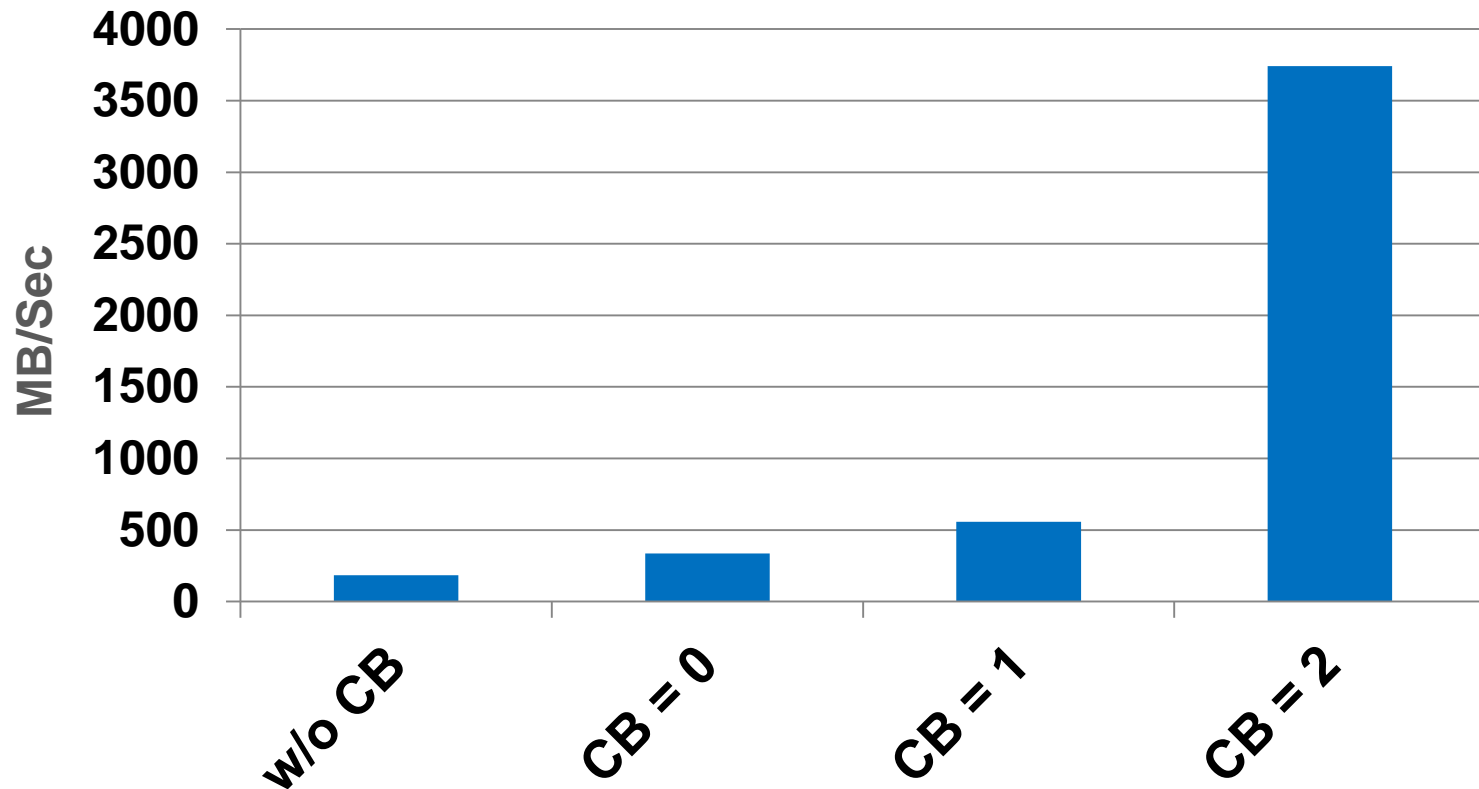
IOR benchmark 10,000 bytes

MPI-IO API , non-power-of-2 blocks and transfers, in this case blocks and transfers both of 10K bytes and a strided access pattern. Tested on an XT5 with 32 PEs, 8 cores/node, 16 stripes, 16 aggregators, 3220 segments, 96 GB file



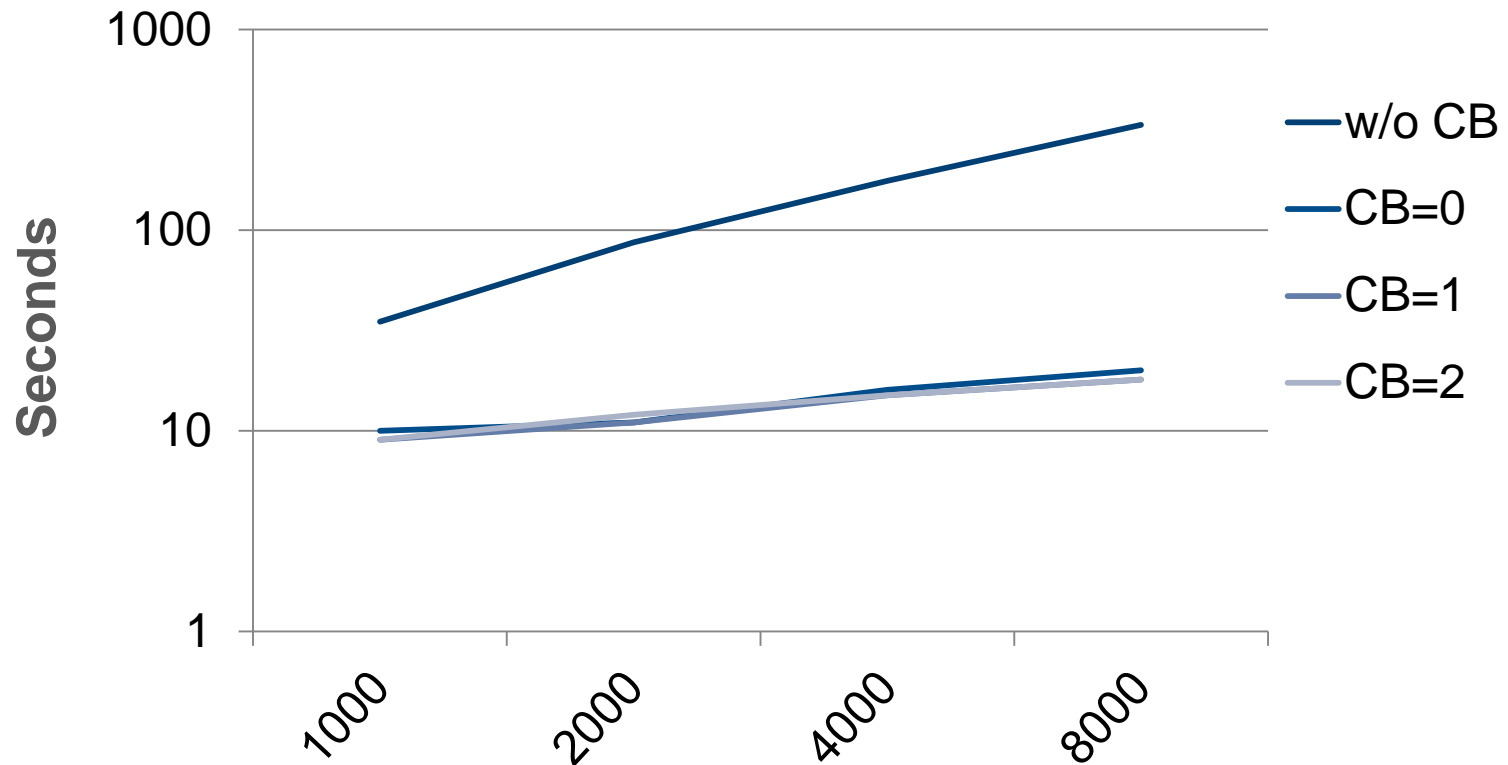
HYCOM MPI-2 I/O

On 5107 PEs, and by application design, a subset of the PEs(88), do the writes. With collective buffering, this is further reduced to 22 aggregators (cb_nodes) writing to 22 stripes. Tested on an XT5 with 5107 PEs, 8 cores/node



HDF5 format dump file from all PEs

Total file size 6.4 GiB. Mesh of 64M bytes 32M elements, with work divided amongst all PEs. Original problem was very poor scaling. For example, without collective buffering, 8000 PEs take over 5 minutes to dump. Note that disabling data sieving was necessary. Tested on an XT5, 8 stripes, 8 cb_nodes

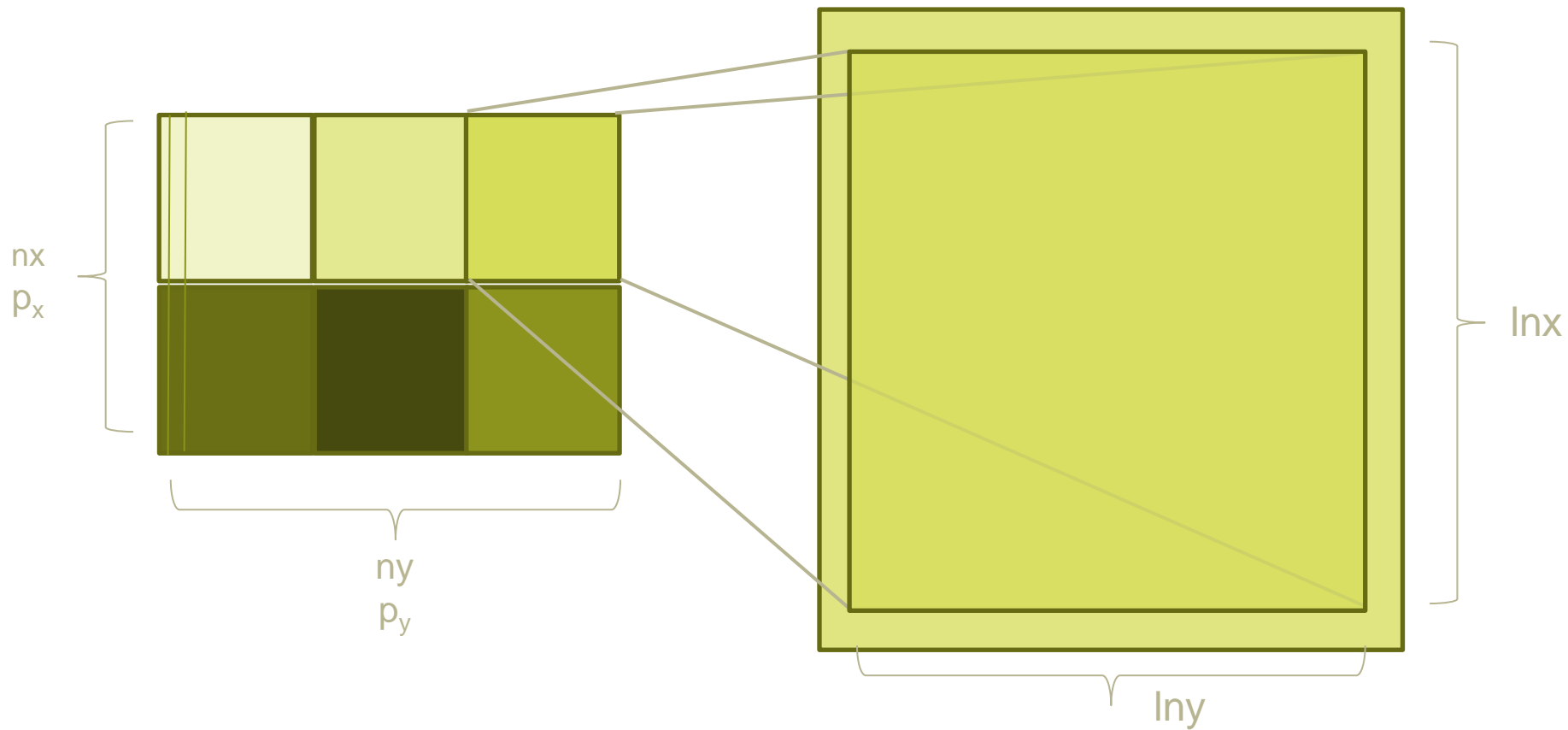


MPI-IO Example

Storing a distributed Domain into a single
File

Problem we want to solve

- We have 2 dim domain on a 2 dimensional processor grid
- Each local subdomain has a halo (ghost cells).
- The data (without halo) is going to be stored in a single file, which can be re-read by any processor count
- Here an example with 2x3 processor grid :



Approach for writing the file

- **First step is to create the MPI 2 dimensional processor grid**
- **Second step is to describe the local data layout using a MPI datatype**
- **Then we create a “global MPI datatype” describing how the data should be stored**
- **Finally we do the I/O**

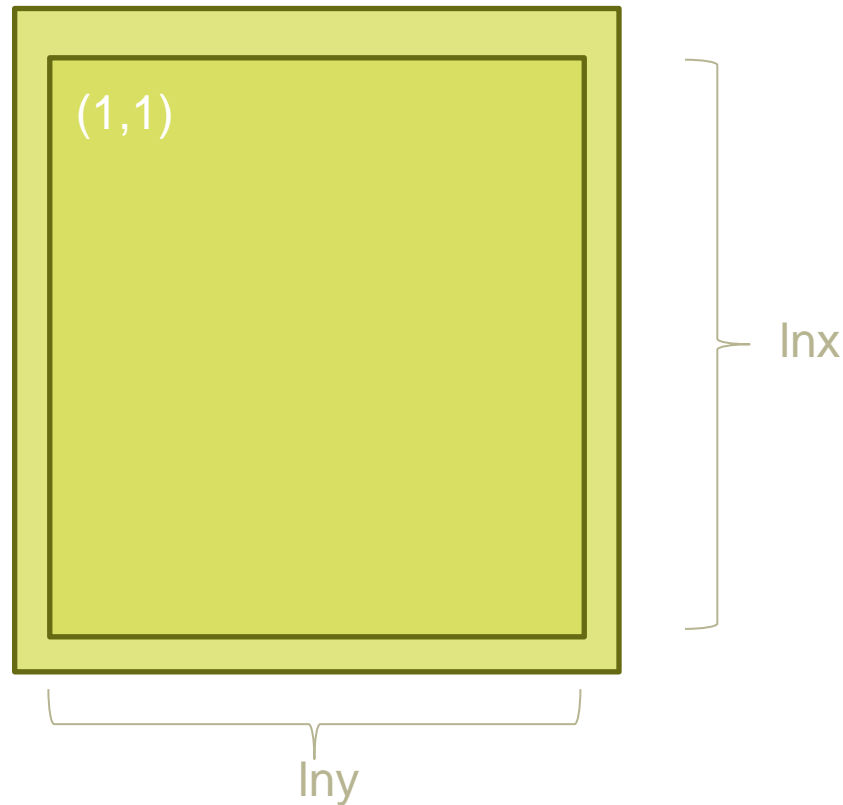
Basic MPI setup

```
nx=512; ny=512 ! Global Domain Size
call MPI_Init(mpierr)
call MPI_Comm_size(MPI_COMM_WORLD, mysize, mpierr)
call MPI_Comm_rank(MPI_COMM_WORLD, myrank, mpierr)

dom_size(1)=2; dom_size(2)=mysize/dom_size(1)
lnx=nx/dom_size(1) ; lny=ny/dom_size(2) ! Local Domain size
periods=.false. ; reorder=.false.
call MPI_Cart_create(MPI_COMM_WORLD, dim, dom_size, periods,
reorder, comm_cart, mpierr)
call MPI_Cart_coords(comm_cart, myrank, dim, my_coords, mpierr)

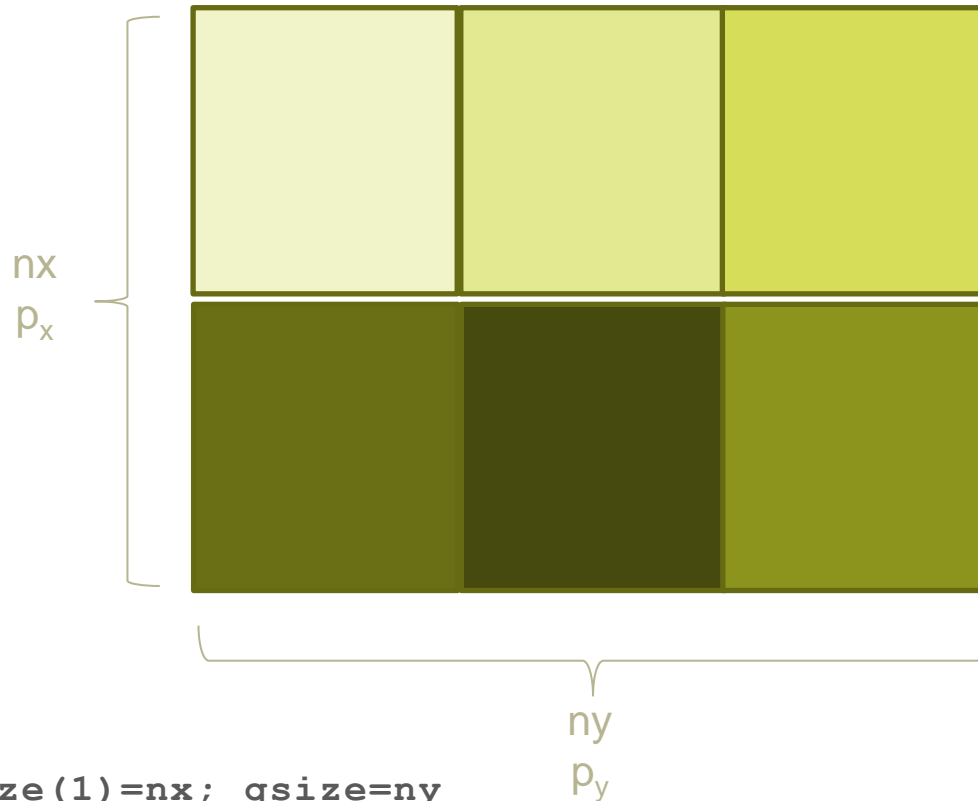
halo=1
allocate (domain(0:lnx+halo, 0:lny+halo))
```

Creating the local data type



```
gsize(1)=lnx+2; gsize(2)=lny+2
lsize(1)=lnx; lsize(2)=lny
start(1)=1; start(2)=1
call MPI_Type_create_subarray(dim, gsize, lsize, start,
    MPI_ORDER_FORTRAN, MPI_INTEGER, type_local, mpierr)
call MPI_Type_commit(type_local, mpierr)
```

And now the global datatype



```

gsize(1)=nx; gsize=ny
lsize(1)=lnx; lsize(2)=lny
start(1)=lnx*my_coords(1); start(2)=lny*my_coords(2)
call MPI_Type_create_subarray(dim, gsize, lsize, start,
    MPI_ORDER_FORTRAN, MPI_INTEGER, type_domain, mpierr)
call MPI_Type_commit(type_domain, mpierr)
  
```

Now we have all together

```
call MPI_Info_create(fileinfo, mpierr)
call MPI_File_delete('FILE', MPI_INFO_NULL, mpierr)
call MPI_File_open(MPI_COMM_WORLD, 'FILE',
    IOR(MPI_MODE_RDWR, MPI_MODE_CREATE), fileinfo, fh, mpierr)

disp=0 ! Note : INTEGER(kind=MPI_OFFSET_KIND) :: disp
call MPI_File_set_view(fh, disp, MPI_INTEGER, type_domain, 'native',
    fileinfo, mpierr)
call MPI_File_write_all(fh, domain, 1, type_local, status, mpierr)
call MPI_File_close(fh, mpierr)
```


I/O Performance Summary

- **Buy sufficient I/O hardware for the machine**
 - As your job grows, so does your need for I/O bandwidth
 - You might have to change your I/O implementation when scaling
- **Lustre**
 - Minimize contention for file system resources.
 - A single process should not access more than 4 OSTs, less might be better
- **Performance**
 - Performance is limited for single process I/O.
 - Parallel I/O utilizing a file-per-process or a single shared file is limited at large scales.
 - Potential solution is to utilize multiple shared file or a subset of processes which perform I/O.
 - A dedicated I/O Server process (or more) might also help
 - Did not really talk about the MDS

And there is more

- <http://docs.cray.com>
 - Search for MPI-IO : „Getting started with MPI I/O“, „Optimizing MPI-IO for Applications on CRAY XT Systems“
 - Search for lustre (a lot for admins but not only)
 - Message Passing Toolkit
- **Man pages (man mpi, man <mpi_routine>, ...)**
- **mpich2 standard :**
<http://www.mcs.anl.gov/research/projects/mpich2/>

Cray Scientific Libraries (Libsci)

- **Traditional model**
 - Tuned general purpose codes
 - Only good for dense
 - Not problem sensitive
 - Not architecture sensitive

Cray Scientific Libraries Today

- **Goal of scientific libraries**
 - Improve Productivity at optimal performance**
- **Cray use four concentrations to achieve this**
 - **Standardization**
 - Use standard or “de facto” standard interfaces whenever available
 - **Hand tuning**
 - Use extensive knowledge of target processor and network to optimize common code patterns
 - **Auto-tuning**
 - Automate code generation and a huge number of empirical performance evaluations to configure software to the target platforms
 - **Adaptive Libraries**
 - Make runtime decisions to choose the best kernel/library/routine

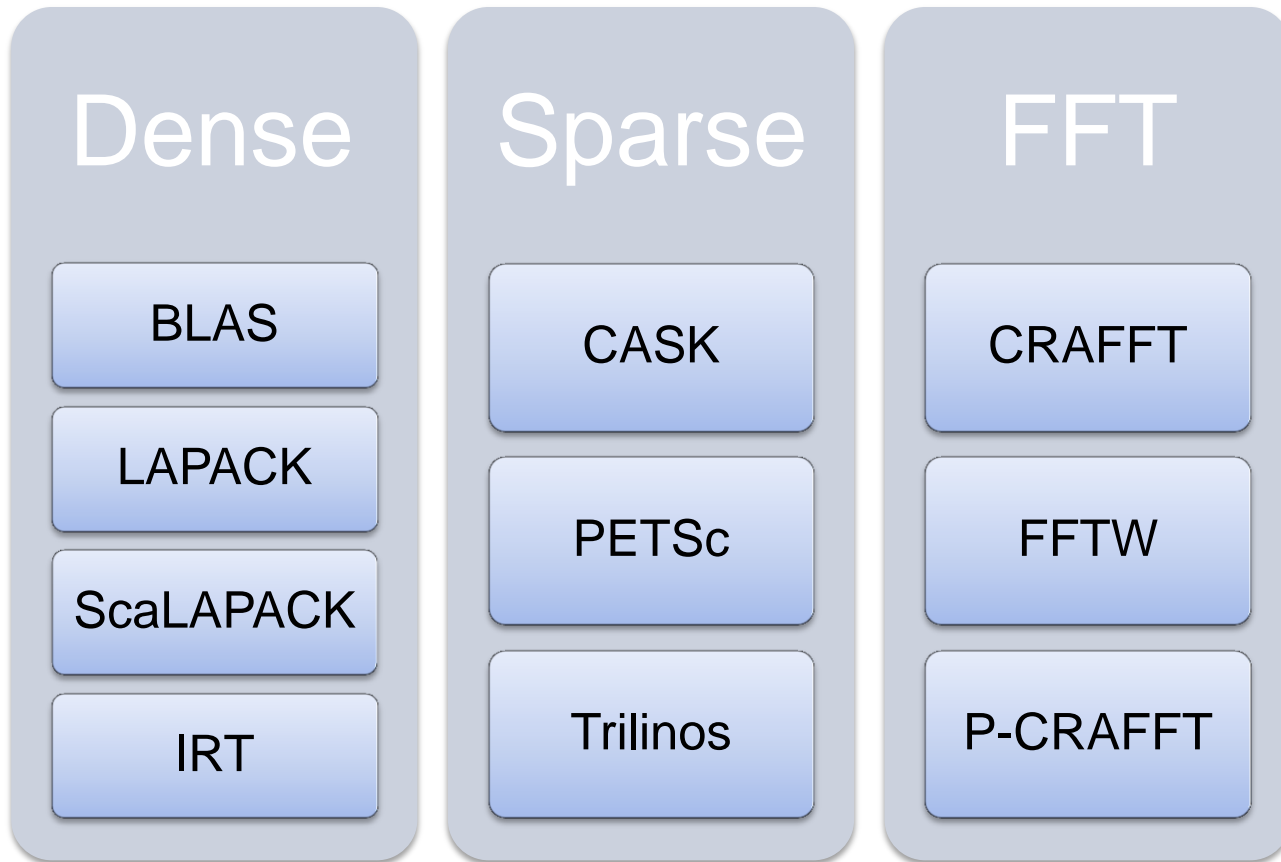
Standardization

- **Three separate classes of standardization, each with a corresponding definition of productivity**
 1. Standard interfaces (e.g., dense linear algebra)
 - Bend over backwards to keep everything the same despite increases in machine complexity. Innovate 'behind-the-scenes'
 - Productivity -> innovation to keep things simple
 2. Adoption of near-standard interfaces (e.g., sparse kernels)
 - Assume near-standards and promote those. Out-mode alternatives. Innovate 'behind-the-scenes'
 - Productivity -> innovation in the simplest areas
 - (requires the same innovation as #1 also)
 3. Simplification of non-standard interfaces (e.g., FFT)
 - Productivity -> innovation to make things simpler than they are

Hand Tuning

- **Algorithmic tuning**
 - Increased performance by exploiting algorithmic improvements
 - Sub-blocking, new algorithms
 - **LAPACK, ScaLAPACK**
- **Kernel tuning**
 - Improve the numerical kernel performance in assembly language
 - **BLAS, FFT**
- **Parallel tuning**
 - Exploit Cray's custom network interfaces and MPT
 - **ScaLAPACK, P-CRAFFT**

Cray Scientific/Math Libraries



IRT – Iterative Refinement Toolkit

CASK – Cray Adaptive Sparse Kernels

CRAFFT – Cray Adaptive FFT

Threaded BLAS Library

- **Libsci includes standard BLAS1,2,3**
- **Most BLAS in libsci are highly tuned and threaded**
- **The emphasis is on the routines that are the most important to users – feedback always welcome**
- **There are single and multi-threaded libraries on the system**
- **The multi-thread library is linked by default**
- **The single-thread library is there for specialist use and debugging –no real reason to try it**
- **Usage – just as standard BLAS**

Threading

- **LibSci is compatible with OpenMP**
- **Control the number of threads to be used in the your program with OMP_NUM_THREADS**
 - E.g. in job script
 - `export OMP_NUM_THREADS=16`
- **What behaviour you get from the library depends upon your code**
 1. No threading in code
 - The BLAS Call will use OMP_NUM_THREADS threads
 2. Threaded code, outside a parallel region
 - The BLAS call will use OMP_NUM_THREADS threads
 3. Threaded code, inside a parallel region
 - The BLAS call will use a single thread

PETSc (Portable, Extensible Toolkit for Scientific Computation)

- **Serial and Parallel versions of sparse iterative linear solvers**
 - Suites of iterative solvers
 - CG, GMRES, BiCG, QMR, etc.
 - Suites of preconditioning methods
 - IC, ILU, diagonal block (ILU/IC), Additive Schwartz, Jacobi, SOR
 - Support block sparse matrix data format for better performance
 - Interface to external packages (ScaLAPACK, SuperLU_DIST)
 - Fortran and C support
 - Newton-type nonlinear solvers
- **Large user community**
 - DoE Labs, PSC, CSCS, CSC, ERDC, AWE and more.
- **<http://www-unix.mcs.anl.gov/petsc/petsc-as>**

Usage and External Packages

- **To use Cray-PETSc, load the appropriate module :**

```
module load petsc
```

```
module load petsc-complex
```

(no need to load a compiler specific module)

- **Treat the Cray distribution as your local PETSc installation**

Third Party Scientific Libraries

- **Cray provides state-of-the art scientific computing packages to strengthen the capability of PETSc**
 - **Hypre 2.7.0b**: scalable parallel pre-conditioners
 - AMG (Very scalable and efficient for specific class of problems)
 - 2 different ILU (General purpose)
 - Sparse Approximate Inverse (General purpose)
 - **ParMetis 3.2.0**: parallel graph partitioning package
 - **MUMPS 4.9.2**: parallel multifrontal sparse direct solver
 - **SuperLU_dist 2.5**: Parallel sparse, direct linear-system solvers
 - **SuperLU 4.1**: sequential version of SuperLU_dist
 - **SUNDIALS 2.4.0**: Suite of Nonlinear and Differential/Algebraic equation solvers
 - **Scotch 5.1.1**: Sequential and parallel graph partitioning library
- **Load with**
 - `module load tps1`

Trilinos

- **The Trilinos Project <http://trilinos.sandia.gov/>**
“an effort to develop algorithms and enabling technologies within an object-oriented software framework for the solution of large-scale, complex multi-physics engineering and scientific problems”
- **A unique design feature of Trilinos is its focus on packages.**
- **Very large user-base and growing rapidly. Important to DOE.**
- **Cray’s optimized Trilinos released on January 21**
 - Includes 50+ trilinos packages
 - Optimized via CASK
 - Any code that uses Epetra objects can access the optimizations
- **Usage :**
`module load trilinos`

Cray Adaptive Sparse Kernel (CASK)

- **CASK is a product developed at Cray using the Cray Auto-tuning Framework (Cray ATF)**
- **The CASK Concept :**
 - Analyze matrix at minimal cost
 - Categorize matrix against internal classes
 - Based on offline experience, find best CASK code for particular matrix
 - Previously assign “best” compiler flags to CASK code
 - Assign best CASK kernel and perform Ax
- **CASK silently sits beneath PETSc and Trilinos on Cray systems**



Support Model



Large-scale application

- Highly portable
- User controlled

PETSc / Trilinos / Hypre

- Highly portable
- User controlled

All systems

CASK

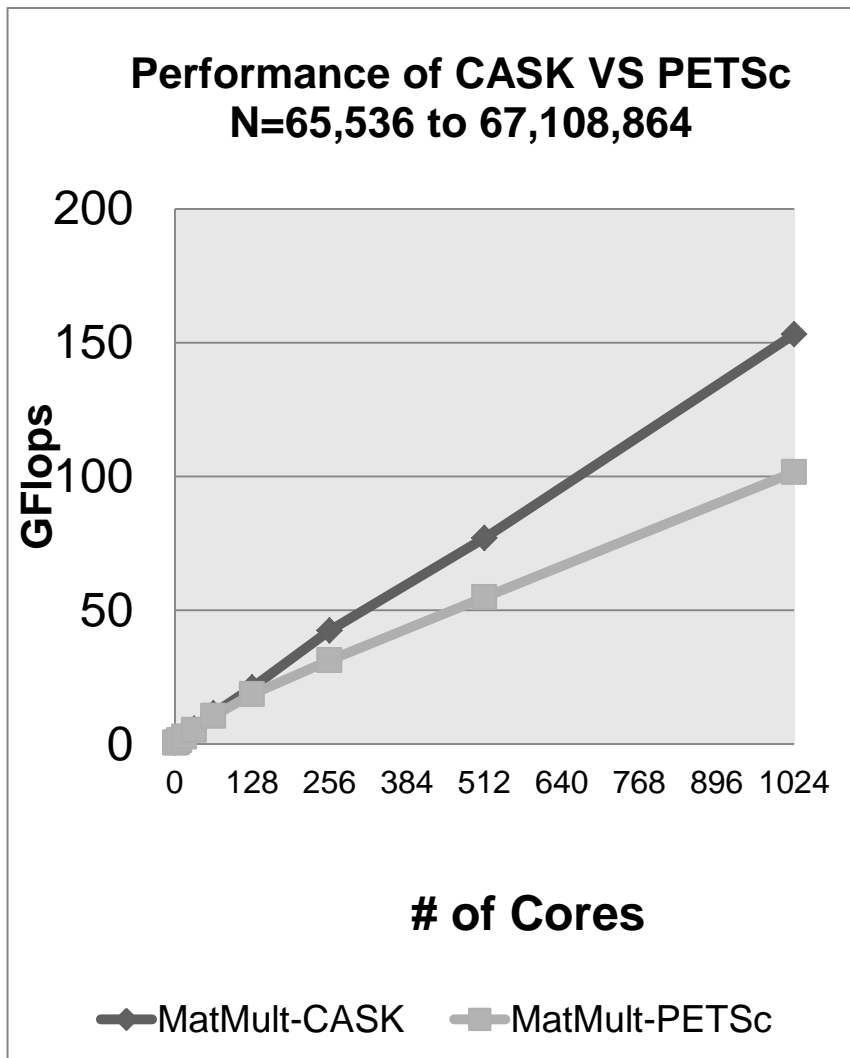
- Cray specific / tuned
- Invisible to User

Cray only

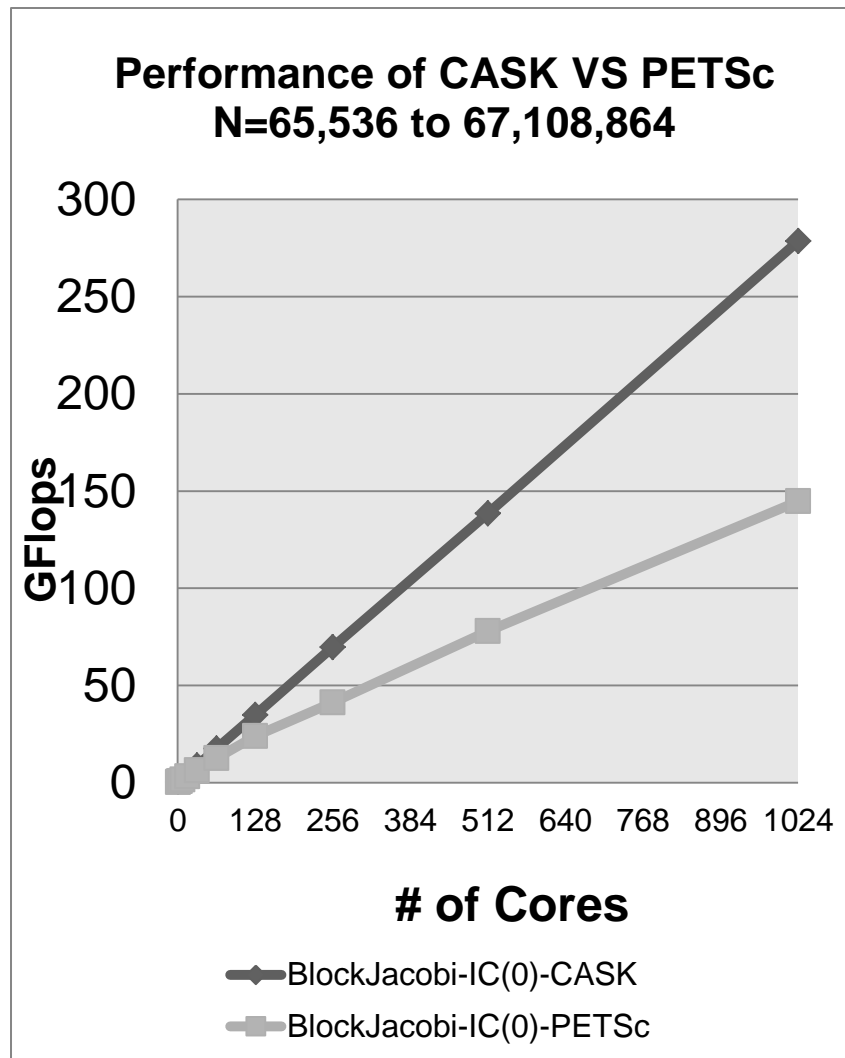
CASK + PETSc Scalability (XT4)



SpMV



Block Jacobi Preconditioning



Cray Adaptive FFT (CRAFFT)

- **In FFTs, the problems are**
 - Which library choice to use?
 - How to use complicated interfaces (e.g., FFTW)
- **Standard FFT practice**
 - Do a plan stage
 - Deduced machine and system information and run micro-kernels
 - Select best FFT strategy
 - Do an execute

Our system knowledge can remove some of this cost!

CRAFFT library

- **CRAFFT is designed with simple-to-use interfaces**
 - Planning and execution stage can be combined into one function call
 - Underneath the interfaces, CRAFFT calls the appropriate FFT kernel
- **CRAFFT provides both offline and online tuning**
 - Offline tuning
 - Which FFT kernel to use
 - Pre-computed PLANs for common-sized FFT
 - No expensive plan stages
 - Online tuning is performed as necessary at runtime as well
- **At runtime, CRAFFT will **adaptively select the best** FFT kernel to use based on both offline and online testing (e.g. FFTW, Custom FFT)**
- **man intro_crafft**

CRAFFT usage

1. Load module fftw/3.2.0 or higher.
2. Add a Fortran statement “use crafft”
3. call `crafft_init()`
4. Call `crafft transform` using `none`, `some` or `all` optional arguments (as shown in red)

In-place, implicit memory management :

```
call crafft_z2z3d(n1,n2,n3,input,ld_in,ld_in2,isign)
```

in-place, explicit memory management

```
call crafft_z2z3d(n1,n2,n3,input,ld_in,ld_in2,isign,work)
```

out-of-place, explicit memory management :

```
crafft_z2z3d(n1,n2,n3,input,ld_in,ld_in2,output,ld_out,ld_out2,isign,work)
```

Note : the user can also control the planning strategy of CRAFFT using the `CRAFFT_PLANNING` environment variable and the `do_exe` optional argument, please see the `intro_crafft` man page.

Parallel CRAFFT

- **As of December 2009, CRAFFT includes distributed parallel transforms**
- **Uses the CRAFFT interface prefixed by “p”, with optional arguments**
- **Can provide performance improvement over FFTW 2.1.5**
- **Currently implemented**
 - complex-complex
 - Real-complex and complex-real
 - 3-d and 2-d
 - In-place and out-of-place
- **Upcoming**
 - C language support for serial and parallel

parallel CRAFFT usage

1. Add “use crafft” to Fortran code
2. Initialize CRAFFT using `crafft_init`
3. Assume MPI initialized and data distributed (see manpage)
4. Call `crafft`, e.g. (optional arguments in red)

2-d complex-complex, in-place, internal mem management :

```
call crafft_pz2z2d(n1,n2,input,isign,flag,comm)
```

2-d complex-complex, in-place with no internal memory :

```
call crafft_pz2z2d(n1,n2,input,isign,flag,comm,work)
```

2-d complex-complex, out-of-place, internal mem manager :

```
call crafft_pz2z2d(n1,n2,input,output,isign,flag,comm)
```

2-d complex-complex, out-of-place, no internal memory :

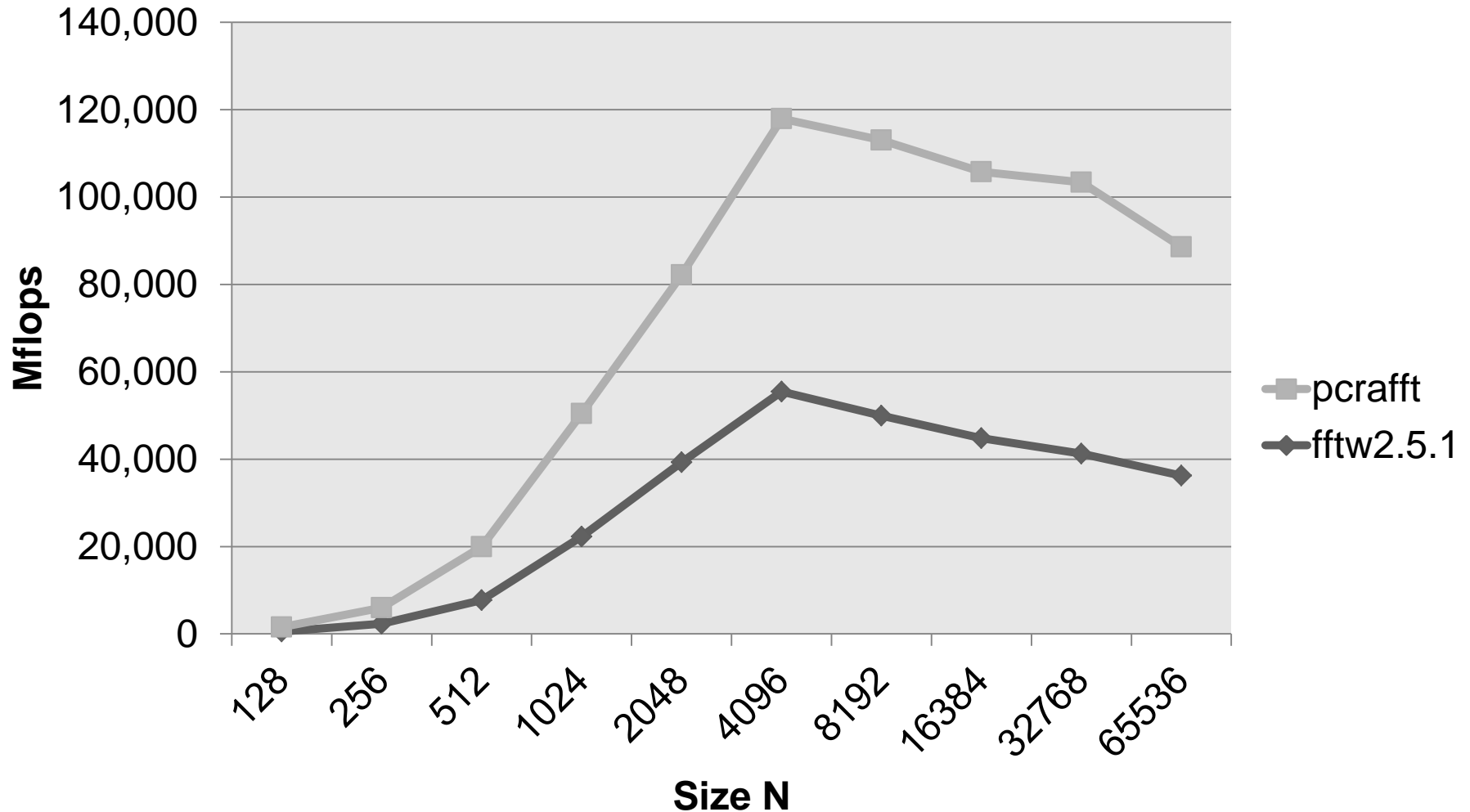
```
crafft_pz2z2d(n1,n2,input,output,isign,flag,comm,work)
```

Each routine above has manpage. Also see 3d equivalent :

```
man crafft_pz2z3d
```

Parallel CRAFFT performance

2D FFT (N x N, transposed), 128 cores

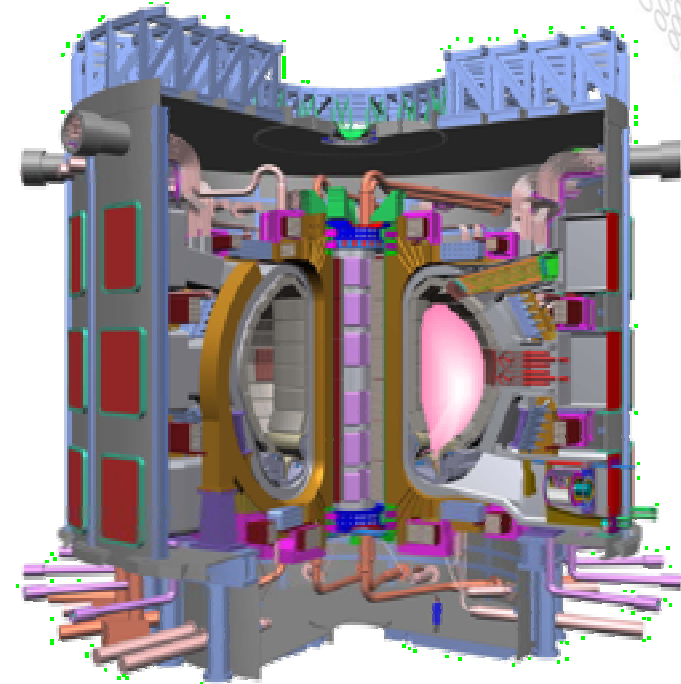


Iterative Refinement Toolkit

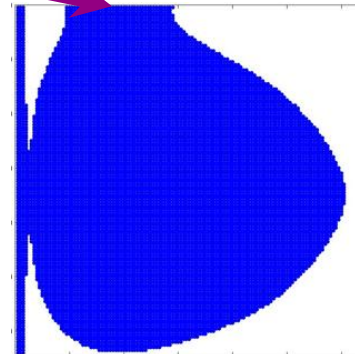
- Solves linear systems in single precision
- Obtaining solutions accurate to double precision
 - For well conditioned problems
- Serial and Parallel versions of LU, Cholesky, and QR
- 2 usage methods
 - **IRT Benchmark routines**
 - Uses IRT 'under-the-covers' without changing your code
 - Simply set an environment variable
 - Useful when you cannot alter source code
 - **Advanced IRT API**
 - If greater control of the iterative refinement process is required
 - Allows
 - condition number estimation
 - error bounds return
 - minimization of either forward or backward error
 - 'fall back' to full precision if the condition number is too high
 - max number of iterations can be altered by users

Example: AORSA Fusion Energy

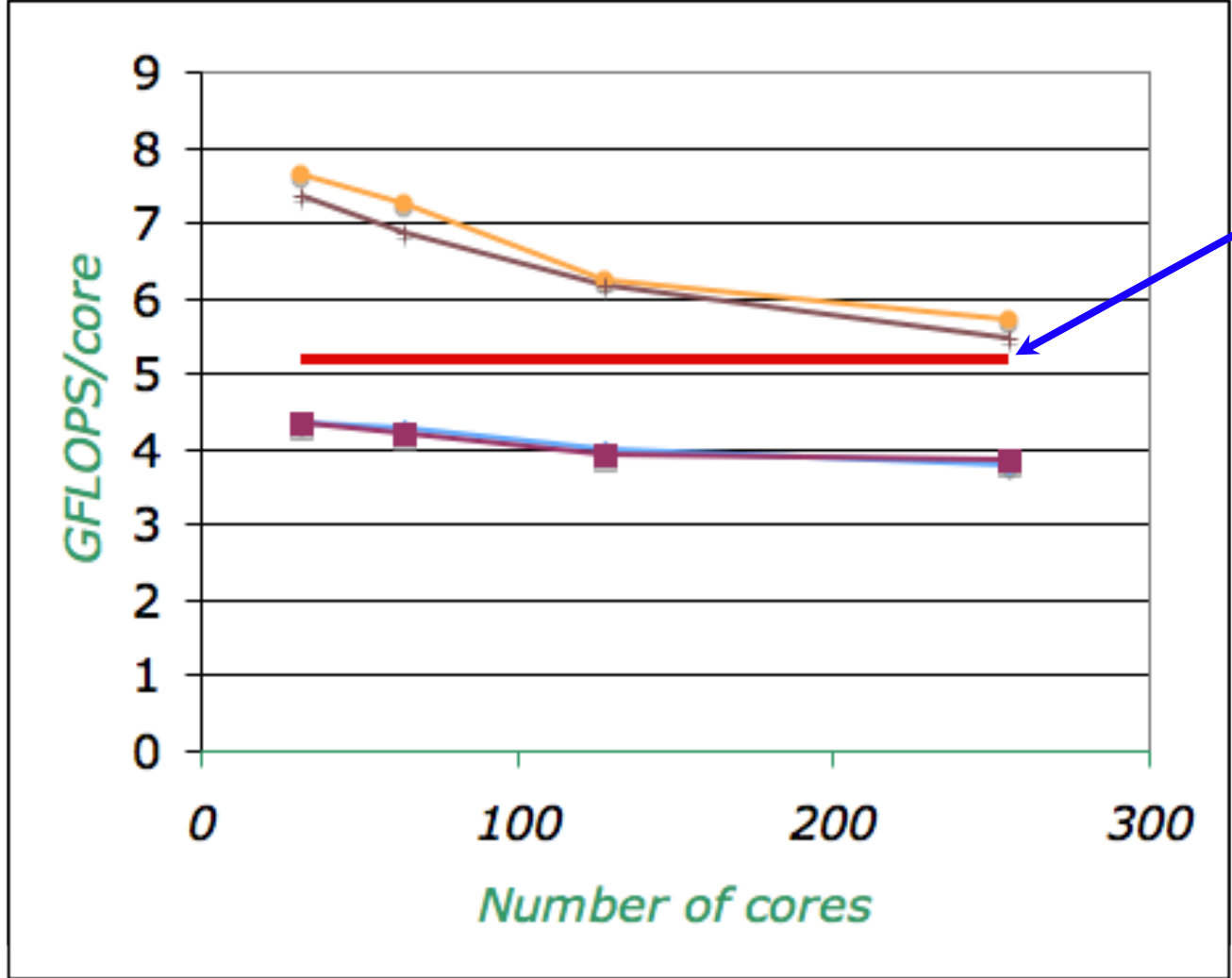
- “High Power Electromagnetic Wave Heating in the ITER Burning Plasma”
- RF heating in Tokamak
- Maxwell-Boltzmann Eqns
- FFT
- Dense linear system
- Calc Quasi-linear op



ITER-FEAT



Example: AORSA Fusion – solver performance on 128x128 grid



Theoretical Peak

IRT usage

Decide if you want to use advanced API or benchmark API

benchmark API :

```
setenv IRT_USE_SOLVERS 1
```

Advanced API :

1. locate the factor and solve in your code (LAPACK or ScaLAPACK)
2. Replace factor and solve with a call to IRT routine
 - e.g. dgesv -> irt_lu_real_serial
 - e.g. pzgesv -> irt_lu_complex_parallel
 - e.g. pzposv -> irt_po_complex_parallel
3. **Set advanced arguments**
 - Forward error convergence for most accurate solution
 - Condition number estimate
 - “fall-back” to full precision if condition number too high

man intro_irt

Recent Updates

- **LibSci 11.0.06 March 15th 2012**
 - Updates for the recent versions from Netlib
 - LAPACK 3.4.0
 - ScaLAPACK 2.0.1
 - GotoBLAS 1.13
- **LibSci 11.1.0 June 16th 2012 (Available soon)**
 - First version of a new BLAS library, CrayBLAS
 - Replaces the current gotoBLAS implementation
 - CrayBLAS optimised extensively using autotuning and runtime adaption.
 - Provides a tuned routine for the specific BLAS call.
 - Can be switched off using `CRAYBLAS_AUTOTUNING_OFF=1`