



Using MPI @ FASRC

Francesco Pontiggia, PhD
Sr. Research Facilitator

March 2020
CS205



Objectives

- To introduce the basic concepts of MPI, and give you the minimum knowledge to write simple parallel MPI programs
- To provide the basic knowledge required for running your parallel MPI applications efficiently on the FAS-RC cluster



Outline

- Introduction and MPI Basics
- MPI Exercises
- Summary



Introduction and MPI Basics

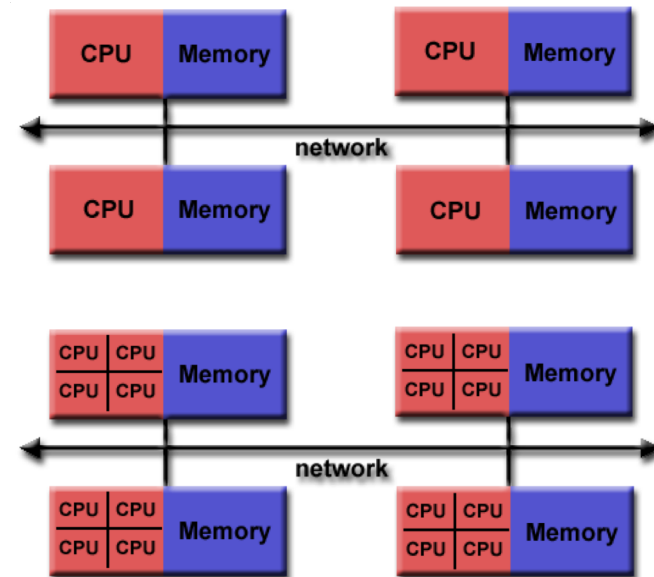


What is MPI?

- **M P I** = **M**essage **P**assing **I**nterface
- MPI is a **specification** for the developers and users of message passing libraries. By itself, it is **NOT** a library
- MPI primarily addresses the message-passing parallel programming model: data is moved from the address space of one process to that of another process through cooperative operations on each process
- Most recent version is MPI-3.1
- Actual MPI library implementations differ in which version and features of the MPI standard they support

MPI Programming Model

- ❑ Originally MPI was designed for distributed memory architectures
- ❑ As architectures evolved, MPI implementations adapted their libraries to handle shared, distributed, and hybrid architectures
- ❑ Today, MPI runs on virtually any hardware platform
 - Shared Memory
 - Distributed Memory
 - Hybrid
- ❑ Programming model remains clearly distributed memory model, regardless of the underlying physical architecture of the machine
- ❑ Explicit parallelism – programmer is responsible for correct implementation of MPI





Reasons for using MPI

- ❑ **Standardization** - MPI is the only message passing specification which can be considered a standard. It is supported on virtually all HPC platforms
- ❑ **Portability** - There is little or no need to modify your source code when you port your application to a different platform that supports (and is compliant with) the MPI standard
- ❑ **Performance Opportunities** - Vendor implementations should be able to exploit native hardware features to optimize performance. Any implementation is free to develop optimized algorithms
- ❑ **Functionality** - There are over 430 routines defined in the most recent MPI
- ❑ **Availability** - A variety of implementations are available, both vendor and public domain



MPI Language Interfaces

- C/C++
- Fortran
- Java
- Python (pyMPI, mpi4py, pypar, MYMPI)
- R (Rmpi)
- Perl (Parallel::MPI)
- MATLAB (DCS)
- Others



Compiling MPI Programs on Odyssey

MPI Implementation	Compiler	Flag
OpenMPI MVAPICH2 Intel MPI	mpicc mpicxx mpif90	None

Intel:

```
module load intel/19.0.5-fasrc01
module load openmpi/4.0.2-fasrc01
mpicxx -o mpi_test.x mpi_test.cpp
```

GNU:

```
module load gcc/9.2.0-fasrc01
module load openmpi/4.0.2-fasrc01
mpicxx -o mpi_test.x mpi_test.cpp
```

<https://docs.rc.fas.harvard.edu/kb/mpi-software-on-odyssey/>



Running MPI Programs on FAS-RC cluster (1)

Interactive test jobs:

(1) Start an interactive bash shell

```
> srun -p test -n 4 --pty --mem=4G -t 0-06:00 /bin/bash
```

(2) Load required modules, e.g.,

```
> module load gcc/9.2.0-fasrc01 openmpi/4.0.2-fasrc01
```

(3) Compile your code (or use a Makefile)

```
> mpicxx -o hello_mpi.x hello_mpi.cpp
```

(4) Run the code

```
> mpirun -np 4 ./hello_mpi.x
```

```
Hello world from process 0 out of 4  
Hello world from process 1 out of 4  
Hello world from process 2 out of 4  
Hello world from process 3 out of 4
```



Running MPI Programs on FAS-RC cluster (2)

Batch jobs:

(1) Compile your code, e.g.,

```
> module load gcc/9.2.0-fasrc01 openmpi/4.0.2-fasrc01
> mpicxx -o mpi_hello.x mpi_hello.cpp
```

(2) Prepare a batch-job submission script

```
#!/bin/bash
#SBATCH -J mpi_job                # Job name
#SBATCH -o slurm.out              # STD output
#SBATCH -e slurm.err              # STD error
#SBATCH -p shared                  # Queue / partition
#SBATCH -t 0-00:30                # Time (D-HH:MM)
#SBATCH --mem-per-cpu=4000        # Memory per MPI task
#SBATCH -n 8                       # Number of MPI tasks
module load gcc/9.2.0-fasrc01 openmpi/4.0.2-fasrc01 # Load required modules
srun -n $SLURM_NTASKS --mpi=pmix ./hello_mpi.x
```

(3) Submit the job to the queue

```
> sbatch mpi_test.run
```



Running MPI Programs on FAS-RC cluster (3)

- Sometimes programs can be picky about having MPI available on all the nodes it runs on so it could be useful to have MPI module loads in your `.bashrc` file
- Some codes are topology sensitive thus the following slurm options can be helpful
 - `--contiguous` # Contiguous set of nodes
 - `--ntasks-per-node` # Number of tasks per node
 - `--hint` # Bind tasks according to hints
 - `--distribution, -m` # Specify distribution method for tasks
- For hybrid mode jobs you would set both `-c` and `-n`

<https://slurm.schedmd.com/sbatch.html>

https://slurm.schedmd.com/mc_support.html

<https://docs.rc.fas.harvard.edu/kb/hybrid-mpiopenmp-codes-on-odyssey/>



MPI Exercises



Exercises - Setup

- Login to the cluster

- Make a directory for this session, e.g.,

```
> mkdir ~/MPI
```

- Get a copy of the MPI examples. These are hosted at Github at https://github.com/fasrc/User_Codes/tree/master/Courses/CS205/MPI_2020

```
> cd ~/MPI
```

```
> git clone https://github.com/fasrc/User_Codes.git
```

```
> cd User_Codes/Courses/CS205/MPI_2020
```

- Load compiler and MPI library software modules

```
> module load gcc/8.2.0-fasrc01
```

```
> module load openmpi/4.0.1-fasrc01
```

(using gcc-8 instead of gcc-9 for fortran compatibility)



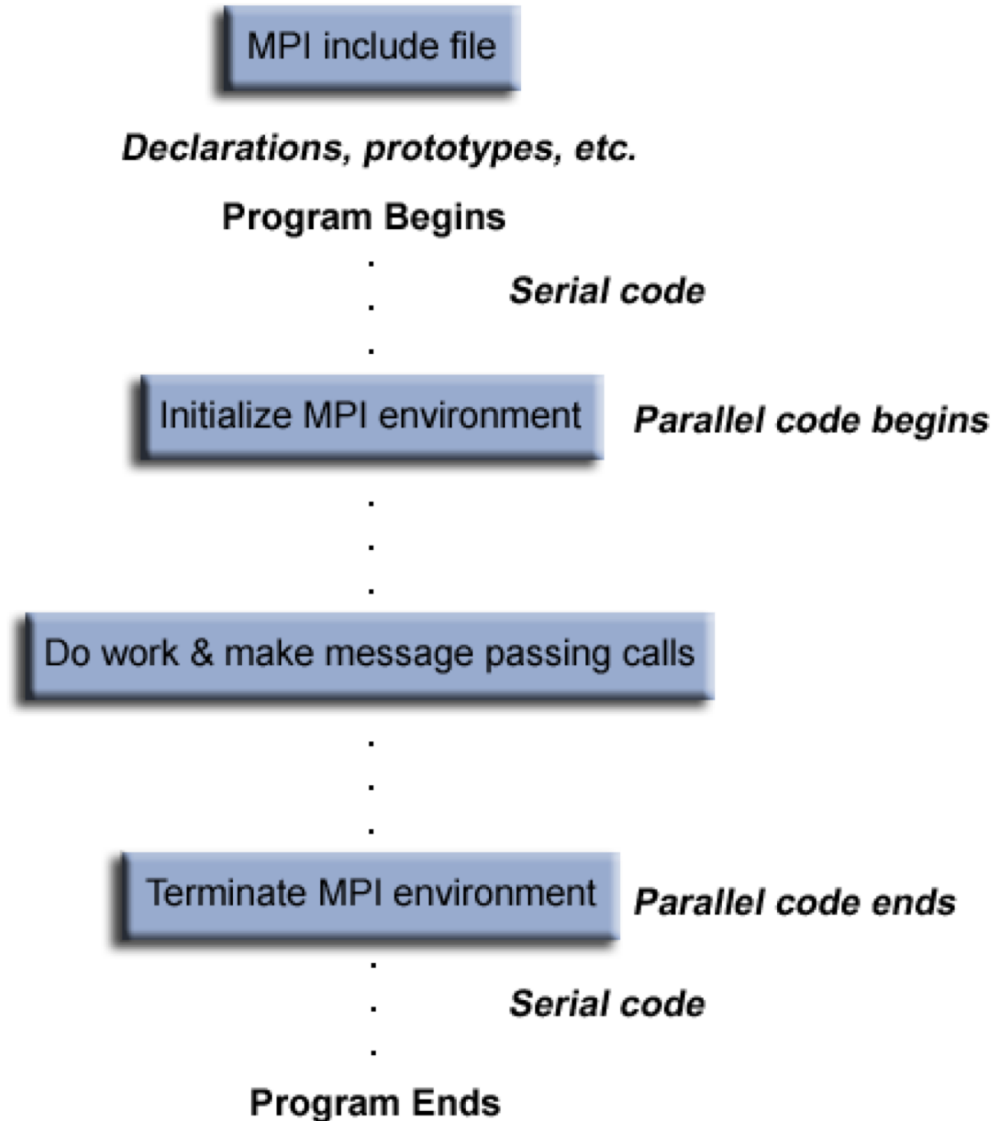
MPI Exercises - Overview

1. MPI Hello World program
2. Parallel FOR loops in MPI – dot product
3. Scaling – speedup and efficiency
4. Parallel Matrix-Matrix multiplication
5. Parallel Lanczos algorithm

https://github.com/fasrc/User_Codes/tree/master/Courses/CS205/MPI_2020



MPI Program Structure



MPI Communicators

- MPI uses objects named communicators to define which processes can communicate with each other
- Most MPI routines require you to specify a communicator as an argument
- `MPI_COMM_WORLD` is a predefined communicator including all MPI processes
- Within a communicator, each process is identified by its rank – a unique integer identifier. Ranks are contiguous and start at 0

MPI_COMM_WORLD





MPI Header Files

- Required for all MPI programs that make MPI library calls

C/C++ Include File	Fortran Include File
#include "mpi.h"	include 'mpif.h'

- With MPI-3 Fortran, the **use mpi_f08** module is preferred over using the include file **mpif.h**

Format of MPI Calls

C/C++ Bindings	
Format	<code>rc = MPI_Xxxx(parameter, ...)</code>
Example	<code>rc = MPI_Bsend(&buf,count,type,dest,tag,comm)</code>
Error Code	Returned as "rc"
Fortran Bindings	
Format	<code>CALL MPI_XXXX(parameter, ...)</code> <code>call mpi_xxxx(parameter, ...)</code>
Example	<code>CALL MPI_BSEND(buf,count,type,dest,tag,comm,ierr)</code>
Error Code	Returned as "ierr"



Exercise 1: MPI Hello World

```
#include <iostream>
#include <mpi.h>
using namespace std;
// Main program.....
int main(int argc, char** argv){
    int i;
    int iproc;
    int nproc;
// Initialize MPI.....
    MPI_Init(&argc, &argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &iproc);
    MPI_Comm_size(MPI_COMM_WORLD, &nproc);

    for ( i = 0; i < nproc; i++ ){
        MPI_Barrier(MPI_COMM_WORLD);
        if ( i == iproc ){
            cout << "Hello world from process " << iproc
                 << " out of " << nproc << endl;
        }
    }
// Shut down MPI.....
    MPI_Finalize();
    return 0;
}
```



Exercise 1: MPI Hello World

(1) Description – a simple parallel “Hello World” program printing out the number of MPI parallel processes and process IDs

(2) Compile the program

```
> cd ~/MPI/User_Codes/Courses/CS205/MPI_2019/Example1  
> make
```

(3) Run the program (the default is setup to 4 MPI tasks)

```
> sbatch sbatch.run
```

(4) Explore the output (the “omp_hello.dat” file), e.g.,

```
> cat mpi_hello.dat  
Hello world from process 0 out of 4  
Hello world from process 1 out of 4  
Hello world from process 2 out of 4  
Hello world from process 3 out of 4
```

(5) Run the program with a different MPI process number – e.g., 2, 4, 8



Parallelizing DO / FOR Loops

In almost all scientific and technical applications, the hot spots are likely to be found in DO / FOR loops.

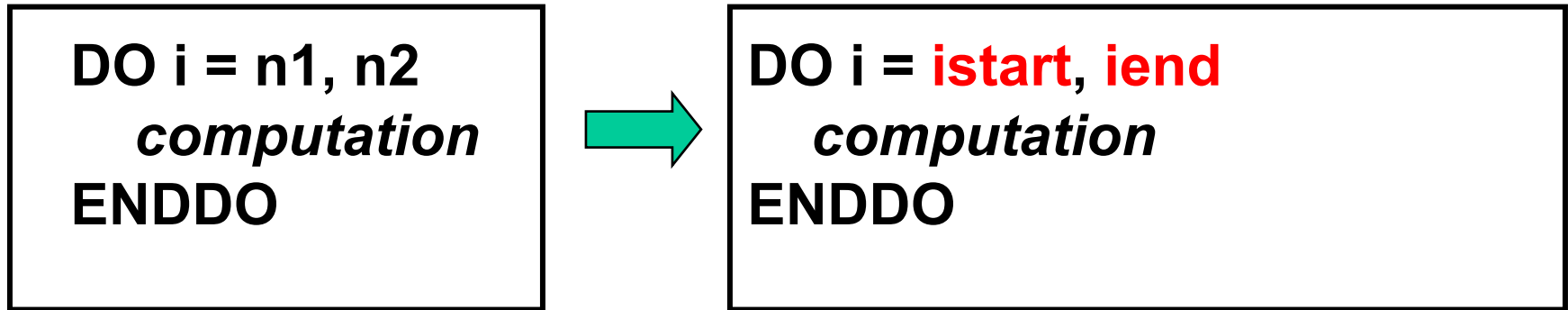
Thus parallelizing DO / FOR loops is one of the most important tasks when you parallelize your program.

The basic technique of parallelizing DO / FOR loops is to distribute iterations among MPI processes and to let each process do its portion in parallel.

Usually, the computations within a DO / FOR loop involve arrays whose indices are associated with the loop variable. Therefore distributing iterations can often be regarded as dividing arrays and assigning chunks (and computations associated with them) to MPI processes.

Block Distribution

In block distribution, iterations are divided into p parts, where p is the number of MPI processes to be executed in parallel.



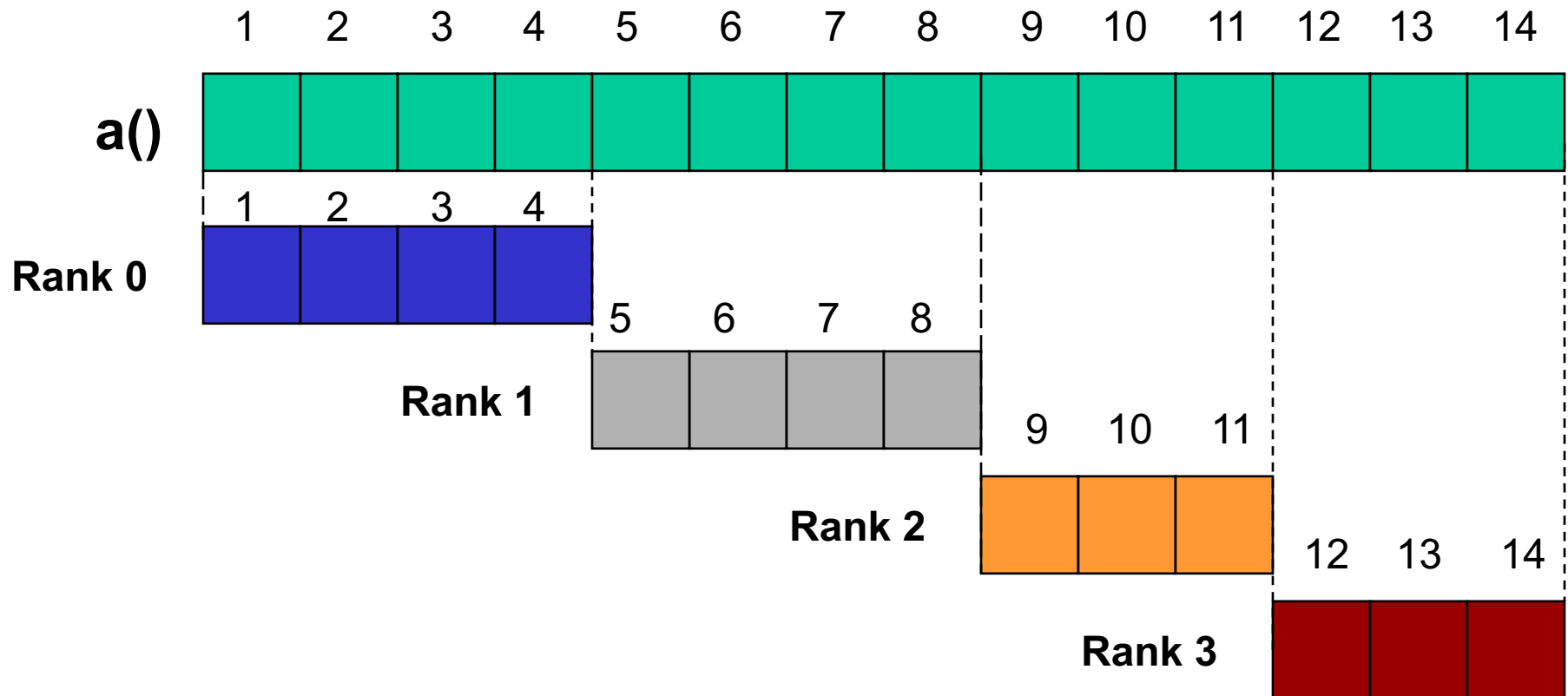
Example: Distributing **14** iterations over **4** MPI tasks

		istart					iend							
Iteration	1	2	3	4	5	6	7	8	9	10	11	12	13	14
Rank	0	0	0	0	1	1	1	1	2	2	2	3	3	3

Shrinking Arrays

Extremely important for efficient memory management!

Block distribution of 14 iterations over 4 cores. Each MPI process needs only part of the array a()





Shrinking Arrays, cont'd

Fortran Example

```
integer(4), allocatable :: a(:)
...
call para_range(1, n, nproc, iproc, istart, iend)
if ( .not. allocated(a) ) allocate( a(istart:iend) )
sum = 0.0
do i = istart, iend
    sum = sum + a(i)
end do
...
if ( allocated(a) ) deallocate(a)
```

The `para_range` subroutine

Computes the iteration range for each MPI process

FORTRAN implementation

```
subroutine para_range(n1, n2, nprocs, irank, istart, iend)
  integer(4) :: n1           !Lower limit of iteration variable
  integer(4) :: n2           !Upper limit of iteration variable
  integer(4) :: nprocs       !Number of MPI ranks
  integer(4) :: irank        !MPI rank ID
  integer(4) :: istart       !Start of iterations for rank iproc
  integer(4) :: iend         !End of iterations for rank iproc
  iwork1 = ( n2 - n1 + 1 ) / nprocs
  iwork2 = MOD(n2 - n1 + 1, nprocs)
  istart = irank * iwork1 + n1 + MIN(irank, iwork2)
  iend = istart + iwork1 - 1
  if ( iwork2 > irank ) iend = iend + 1
  return
end subroutine para_range
```



The `para_range` subroutine, cont'd

Computes the iteration range for each MPI process

C / C++ implementation

```
void para_range(int n1, int n2, int &nprocs, int &irank, int &istart, int &iend){  
    int iwork1;  
    int iwork2;  
    iwork1 = ( n2 - n1 + 1 ) / nprocs;  
    iwork2 = ( ( n2 - n1 + 1 ) % nprocs );  
    istart = irank * iwork1 + n1 + min(irank, iwork2);  
    iend = istart + iwork1 - 1;  
    if ( iwork2 > irank ) iend = iend + 1;  
}
```



Exercise 2: Parallel MPI for / do loops

(1) Description – Program performs a dot product of 2 vectors in parallel

(2) Review the source code and compile the program

```
> cd ~/MPI/User_Codes/Courses/CS205/MPI_2019/Example2  
> make
```

(3) Run the program (the default is setup to 4 MPI tasks)

```
> sbatch sbatch.run
```

(4) Explore the output (the “mpi_dot.dat” file), e.g.,

```
> cat mpi_dot.dat  
Global dot product: 676700  
Local dot product for MPI process 0: 11050  
Local dot product for MPI process 1: 74800  
Local dot product for MPI process 2: 201050  
Local dot product for MPI process 3: 389800
```

(5) Run the program with a different MPI process number – e.g., 2, 4, 8



Exercise 2: Parallel MPI for / do loops

```
...
// Initialize MPI.....
MPI_Init(&argc,&argv);
MPI_Comm_rank(MPI_COMM_WORLD,&iproc);
MPI_Comm_size(MPI_COMM_WORLD,&nproc);

// Call "para_range" to compute lowest and highest iteration ranges for each MPI task
para_range( 0, N-1, nproc, iproc, ista, iend );
...

// Calculate local vector dimension and allocate memory.....
loc_dim = iend - ista + 1; // Local DIM
a = new float[loc_dim];
b = new float[loc_dim];

// Calculate local dot product.....
pdot = 0.0;
for ( i = 0; i < loc_dim; i++ ) {
    d1 = a[i];
    d2 = b[i];
    pdot = pdot + ( d1 * d2 );
}

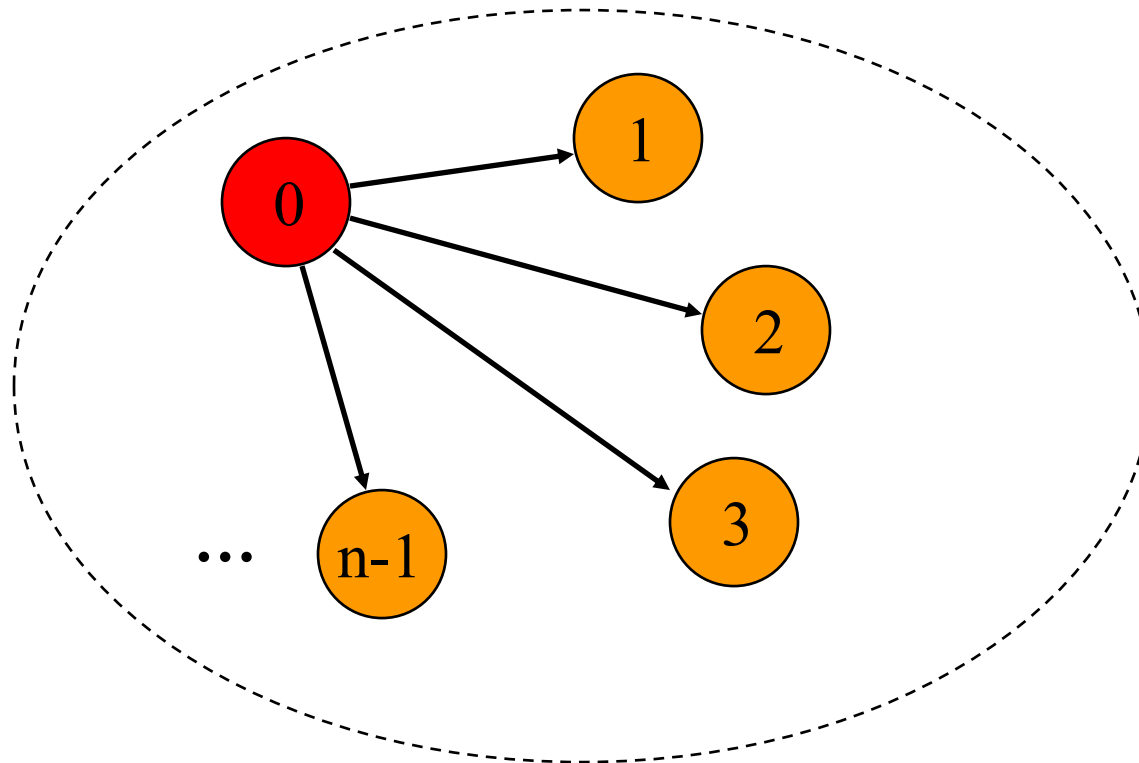
// Get global dot product.....
MPI_Reduce(&pdot, &ddot, 1, MPI_REAL, MPI_SUM, 0, MPI_COMM_WORLD);
...

```

Collective Communication

Collective communication allows you to exchange data among a group of processes. The communicator argument in the collective communication subroutine calls specifies which processes are involved in the communication.

MPI_COMM_WORLD



MPI_Reduce

Usage:

Fortran: CALL MPI_REDUCE(sendbuf, recvbuf, count, datatype, op, root, comm, ierror)

C / C++: MPI_Reduce(&sendbuf, &recvbuf, count, datatype, op, root, comm)

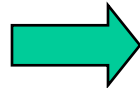


Operation	Data Type
MPI_SUM, MPI_PROD	MPI_INTEGER, MPI_REAL, MPI_DOUBLE_PRECISION, MPI_COMPLEX
MPI_MAX, MPI_MIN	MPI_INTEGER, MPI_REAL, MPI_DOUBLE_PRECISION
MPI_MAXLOC, MPI_MINLOC	MPI_2INTEGER, MPI_2REAL, MPI_2DOUBLE_PRECISION
MPI_LAND, MPI_LOR, MPI_LXOR	MPI_LOGICAL
MPI_BAND, MPI_BOR, MPI_BXOR	MPI_INTEGER, MPI_BYTE

Cyclic Distribution

In cyclic distribution, the iterations are assigned to processes in a round-robin fashion.

```
DO i = n1, n2
  computation
ENDDO
```



```
DO i = n1+iprocc, n2, iprocc
  computation
ENDDO
```

Example: Distributing **14** iterations over **4** cores in round-robin fashion

Iteration	1	2	3	4	5	6	7	8	9	10	11	12	13	14
Rank	0	1	2	3	0	1	2	3	0	1	2	3	0	1

Scaling and Efficiency

How much faster will the program run?

Speedup:

$$S(n) = \frac{T(1)}{T(n)}$$

Time to complete the computation
on **one** MPI process

Time to complete the computation
on **n** MPI processes

Efficiency:

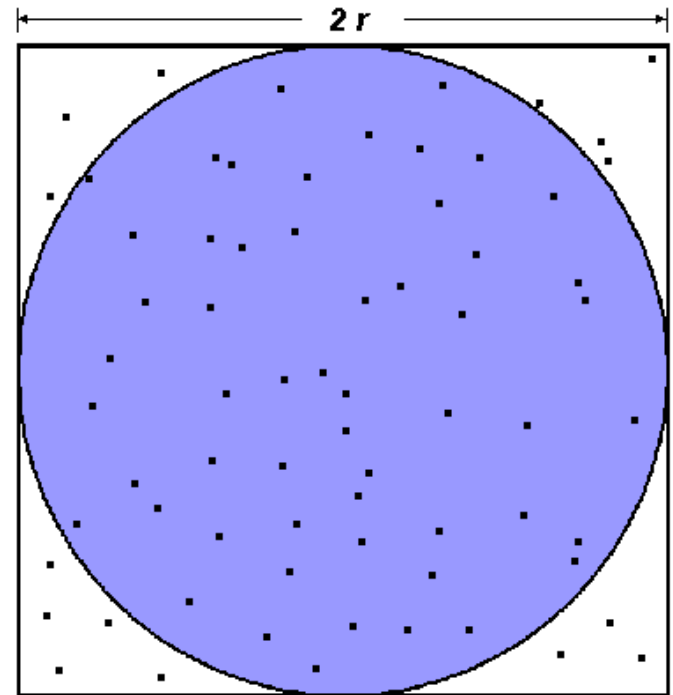
$$E(n) = \frac{S(n)}{n}$$

Tells you how efficiently you parallelized
your code

Example 3: Scaling – Compute PI in Parallel

Monte-Carlo Approximation of PI:

1. Inscribe a circle in a square
2. Randomly generate points in the square
3. Determine the number of points in the square that are also in the circle
4. Let r be the number of points in the circle divided by the number of points in the square
5. $PI \sim 4 r$
6. Note that the more points generated, the better the approximation



$$A_S = (2r)^2 = 4r^2$$

$$A_C = \pi r^2$$

$$\pi = 4 \times \frac{A_C}{A_S}$$



Example 3: Scaling – Compute PI in Parallel

(1) Description – Program performs parallel Monte-Carlo approximation of PI

(2) Review the source code and compile the program

```
> cd ~/MPI/User_Codes/Courses/CS205/MPI_2020/Example3
> make
```

(3) Run the program

```
> sbatch sbatch.run
```

(4) Explore the output (the “mpi_dot.dat” file), e.g.,

```
> cat mpi_pi.dat
...
Elapsed time = 3.429223 seconds
Pi is approximately 3.1416008000000000, Error is 0.0000081464102069
Exact value of PI is 3.1415926535897931
...
```

(5) Run the program with a different number of MPI processes – 1, 2, 4, 8, 16 – and record the run times for each case. This will be needed to compute the speedup and efficiency (the default is set to run on 1, 2, 4, 8, and 16 MPI tasks)

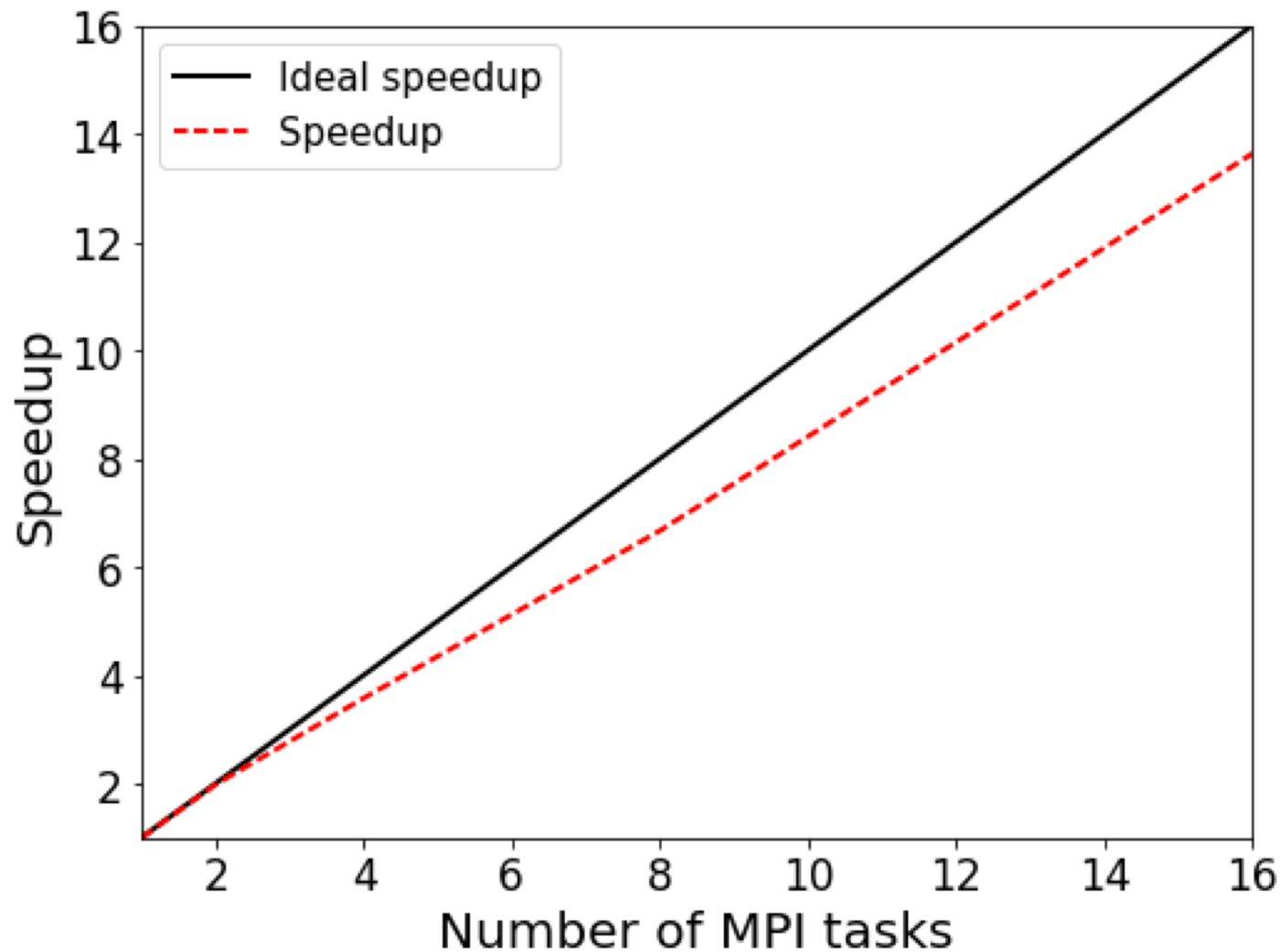


Example 3: Scaling – Compute PI in Parallel

You may use the `speedup.py` Python code to calculate the speedup and efficiency. It generates the below table plus a speedup figure

# MPI tasks	Walltime	Speedup	Efficiency (%)
1	12.27	1.00	100.00
2	6.19	1.98	99.11
4	3.43	3.58	89.43
8	1.84	6.67	83.36
16	0.90	13.63	85.21

Example 3: Scaling – Compute PI in Parallel





Exercise 4: Matrix Multiplication

- A standard problem in computing is matrix multiplication:

$$\mathbf{C} = \mathbf{A} \times \mathbf{B}$$

- In this example we take a naive approach to parallelizing matrix multiplication
- The matrix A is divided up based on its rows by the number of ranks. Each sub array is sent to its relevant ranks by using the `MPI_Scatter` command
- Matrix B is simply sent to all processors using the `MPI_Bcast` command.
- Each rank computes its subset of C based on the part of A it received
- The full solution for C is brought together using `MPI_Gather`



Exercise 4: Matrix Multiplication

(1) Description – A simple algorithm for matrix multiplication demonstrating the use of `MPI_Bcast`, `MPI_Scatter`, and `MPI_Gather`

(2) Compile the program

```
> cd ~/MPI/User_Codes/Courses/CS205/MPI_2020/Example4  
> make
```

(3) Run the program (the default is setup to 4 ranks)

```
> sbatch sbatch.run
```

(4) Explore the output (the “`mmult.dat`” file).

(5) Run the program with different MPI process number – e.g., 1, 2, 4, 8. See how the run time varies depending on number of ranks (**HINT**: use `sacct`, `time`, or `MPI_Wtime` to get duration). Try varying the size of the matrix allowed to see how long it takes and then do scaling tests to see how well this code is parallelized.



MPI Collective Communication Subroutines

Category	Subroutines
1. One buffer	MPI_BCAST
2. One send buffer and one receive buffer	MPI_GATHER , MPI_SCATTER, MPI_ALLGATHER , MPI_ALLTOALL, MPI_GATHERV, MPI_SCATTERV, MPI_ALLGATHERV, MPI_ALLTOALLV
3. Reduction	MPI_REDUCE , MPI_ALLREDUCE , MPI_SCAN, MPI_REDUCE_SCATTER
4. Others	MPI_BARRIER, MPI_OP_CREATE, MPI_OP_FREE

* Subroutines printed in boldface are used most frequently



Exercise 5: Parallel Lanczos diagonalization

(1) Description – Program performs parallel Lanczos diagonalization of a random symmetric matrix of dimension 100 x 100

(2) Review the source code and compile the program

```
> cd ~/MPI/User_Codes/Courses/CS209/MPI_2020/Example5  
> make
```

(3) Run the program

```
> sbatch sbatch.run
```

(4) Explore the output (the “planczos.dat” file), e.g.,

```
> cat planczos.dat
```

```
...  
iteration:           50  
      1  50.010946087355691      50.010946087355670  
      2  5.1987505309251540      5.1987505309260547  
      3  5.1856783411618199      5.1856783411660166  
      4  5.0014143249256930      5.0014143250821714  
      5  4.8032829796748748      4.8032831650372225  
Lanczos iterations finished...
```

(5) Run the program with a different number of MPI processes – 1, 2, 4, 8.



Summary and hints for efficient parallelization

- ❑ Is it even worth parallelizing my code?
 - Does your code take an intractably long amount of time to complete?
 - Do you run a single large model or do statistics on multiple small runs?
 - Would the amount of time it take to parallelize your code be worth the gain in speed?

- ❑ Parallelizing established code vs. starting from scratch
 - Established code: Maybe easier / faster to parallelize, but may not give good performance or scaling
 - Start from scratch: Takes longer, but will give better performance, accuracy, and gives the opportunity to turn a “black box” into a code you understand



Summary and hints for efficient parallelization

- Increase the fraction of your program that can be parallelized. Identify the most time consuming parts of your program and parallelize them. This could require modifying your intrinsic algorithm and code's organization
- Balance parallel workload
- Minimize time spent in communication
- Use simple arrays instead of user defined derived types
- Partition data. Distribute arrays and matrices – allocate specific memory for each MPI process
- For I/O intensive applications implement parallel I/O in conjunction with a high-performance parallel filesystem, e.g., Lustre



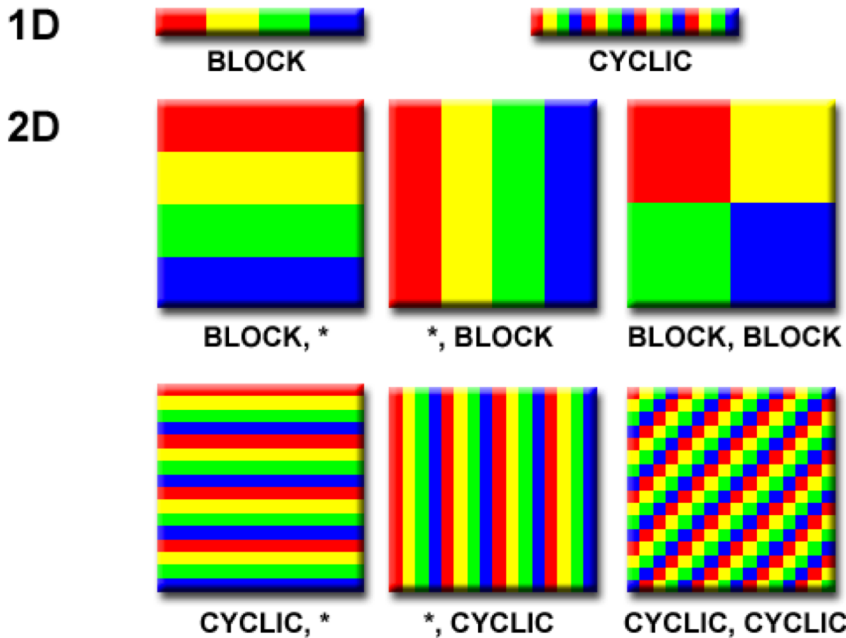
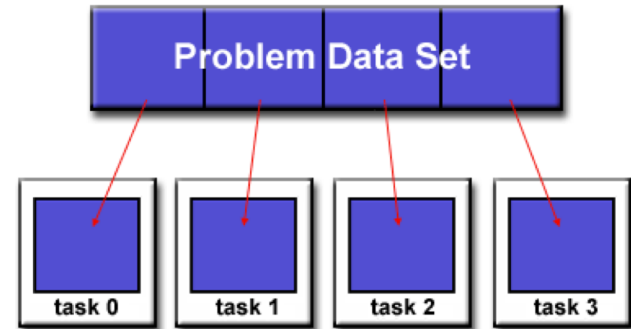
Extra Slides

Designing parallel programs – partitioning

One of the first steps in designing a parallel program is to break the problem into discrete “chunks” that can be distributed to multiple parallel tasks.

Domain Decomposition:

Data associated with a problem is partitioned – each parallel task works on a portion of the data



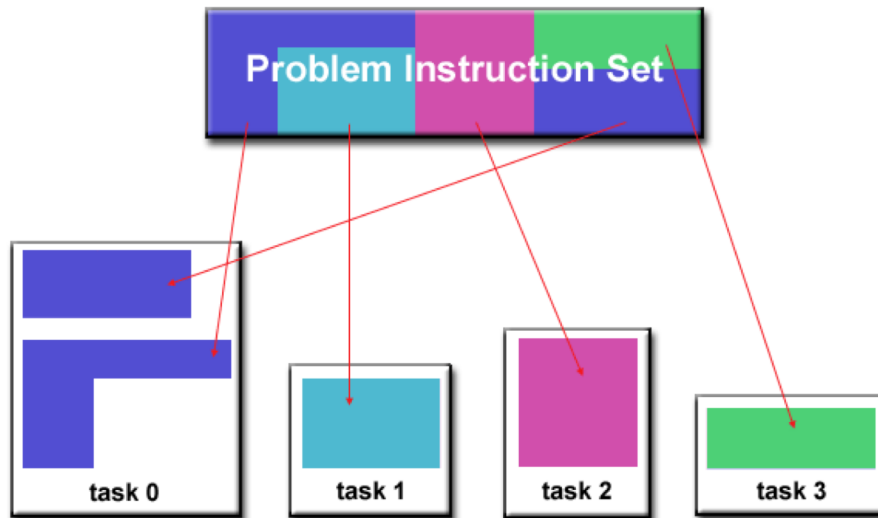
There are different ways to partition the data

Designing parallel programs – partitioning

One of the first steps in designing a parallel program is to break the problem into discrete “chunks” that can be distributed to multiple parallel tasks.

Functional Decomposition:

Problem is decomposed according to the work that must be done. Each parallel task performs a fraction of the total computation.





Designing parallel programs – communication

Most parallel applications require tasks to share data with each other.

Cost of communication: Computational resources are used to package and transmit data. Requires frequently synchronization – some tasks will wait instead of doing work. Could saturate network bandwidth.

Latency vs. Bandwidth: Latency is the time it takes to send a minimal message between two tasks. Bandwidth is the amount of data that can be communicated per unit of time. Sending many small messages can cause latency to dominate communication overhead.

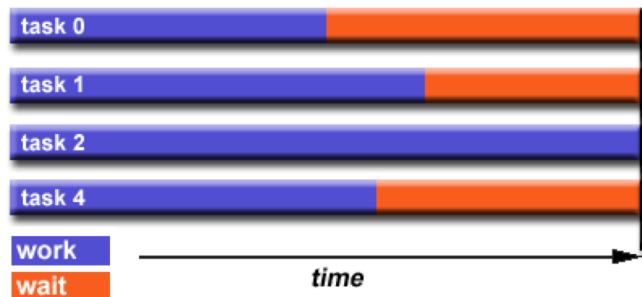
Synchronous vs. Asynchronous communication: **Synchronous** communication is referred to as **blocking** communication – other work stops until the communication is completed. **Asynchronous** communication is referred to as **non-blocking** since other work can be done while communication is taking place.

Scope of communication: Point-to-point communication – data transmission between tasks. Collective communication – involves all tasks (in a communication group)

This is only partial list of things to consider!

Designing parallel programs – loadbalancing

Load balancing is the practice of distributing approximately equal amount of work so that all tasks are kept busy all the time.



How to Achieve Load Balance?

Equally partition the work given to each task: For array/matrix operations equally distribute the data set among parallel tasks. For loop iterations where the work done for each iteration is equal, evenly distribute iterations among tasks.

Use dynamic work assignment: Certain class problems result in load imbalance even if data is distributed evenly among tasks (sparse matrices, adaptive grid methods, many body simulations, etc.). Use scheduler – task pool approach. As each task finishes, it queues to get a new piece of work. Modify your algorithm to handle imbalances dynamically.

Designing parallel programs – I/O

The Bad News:

- I/O operations are inhibitors of parallelism
- I/O operations are orders of magnitude slower than memory operations
- Parallel file systems may be immature or not available on all systems
- I/O that must be conducted over network can cause severe bottlenecks

The Good News:

- Parallel file systems are available (e.g., Lustre)
- MPI parallel I/O interface has been available since 1996 as a part of MPI-2

I/O Tips:

- Reduce overall I/O as much as possible
- If you have access to parallel file system, use it
- Writing large chunks of data rather than small ones is significantly more efficient
- Fewer, larger files perform much better than many small files
- Have a subset of parallel tasks to perform the I/O instead of using all tasks, or
- Confine I/O to a single tasks and then broadcast (gather) data to (from) other tasks