

Introduction to MPI

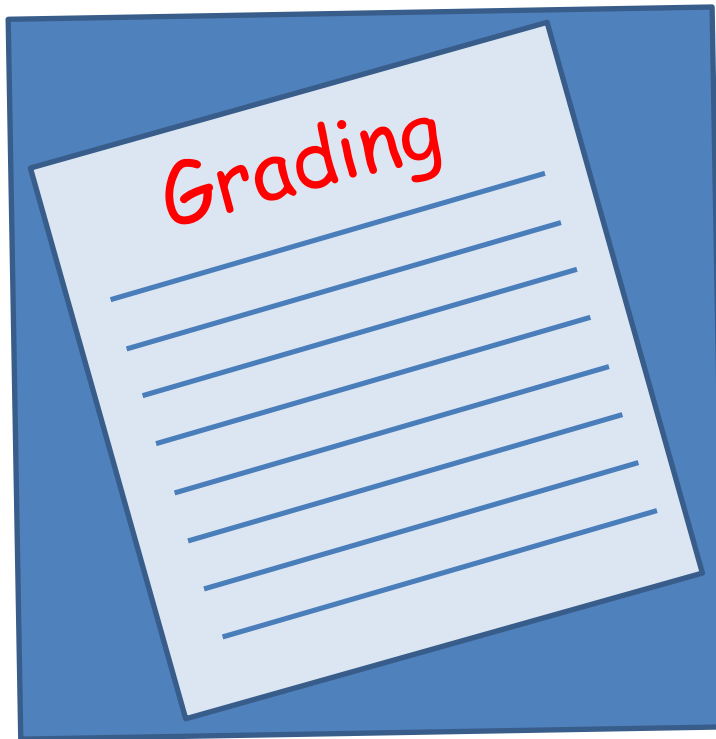
Preeti Malakar

pmalakar@cse.iitk.ac.in

An Introductory Course on High-Performance
Computing in Engineering

28th September 2019

Parallelism



Instructor

60 hours



TAs

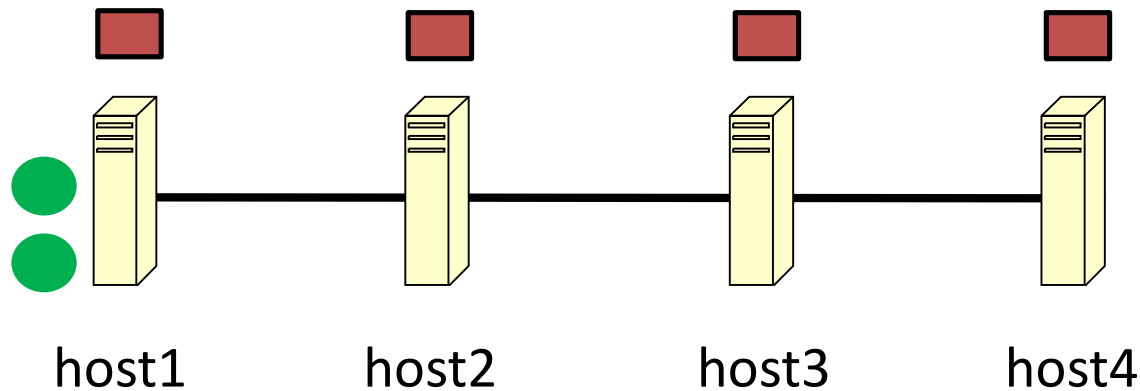
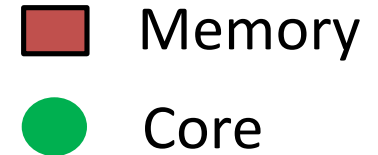
4 hours due to limitations

Total parallel time = $20 + 4 = 24$, Speedup = $60/24$
May not always achieve ideal speedup

This Talk

- Shared memory
 - OpenMP, Pthreads, ...
- Distributed memory
 - MPI, UPC, ...
- Hybrid
 - MPI + OpenMP

System Model



- Interconnected systems
- Distributed memory
- NO centralized server/master

Parallel Code – Getting started

Q: How should I write a parallel code to add up a million numbers using 4 processes (on 4 nodes)?

- Distribute the numbers to the 4 processes
- Collect the result back at one of the processes for further processing

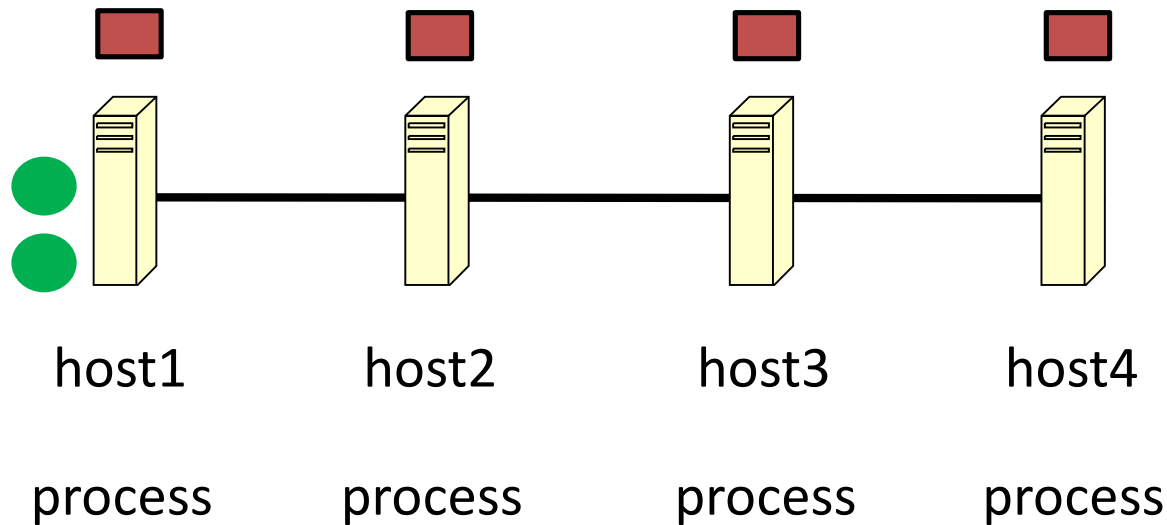
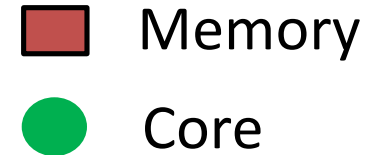
Message Passing Interface (MPI)

- Standard for message passing in a distributed memory environment
- Efforts began in 1991 by Jack Dongarra, Tony Hey, and David W. Walker
- MPI Forum
 - Version 1.0: 1994
 - Version 2.0: 1997
 - Version 3.0: 2012

MPI Implementations

- MPICH (ANL)
- MVAPICH (OSU)
- Intel MPI
- OpenMPI

Parallel Processes



Q: How do the processes communicate among each other?

Communication Channels



- Sockets for network communication
- MPI handles communications, progress etc.

Reference: Design and Evaluation of Nemesis, a Scalable, Low-Latency, Message-Passing Communication Subsystem by Buntinas et al.

Message Passing Paradigm

- Message sends and receives
- Explicit communication

Communication types

- Blocking
- Non-blocking

Getting Started

```
#include <mpi.h>
#include <stdio.h>

int main(int argc, char** argv) {

    // Initialize the MPI environment
    MPI_Init(NULL, NULL);

    // Get the number of processes
    int size;
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    // Get the rank of the process
    int rank;
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);

    // Get the name of the processor
    char processor_name[MPI_MAX_PROCESSOR_NAME];
    int name_len;
    MPI_Get_processor_name(processor_name, &name_len);

    // Print off a hello world message
    printf("Hello I am rank %d out of %d processes\n", rank, size);

    // Finalize the MPI environment.
    MPI_Finalize();
}
```

Function names:
MPI_*

How many
outputs?

Initialization and Finalization

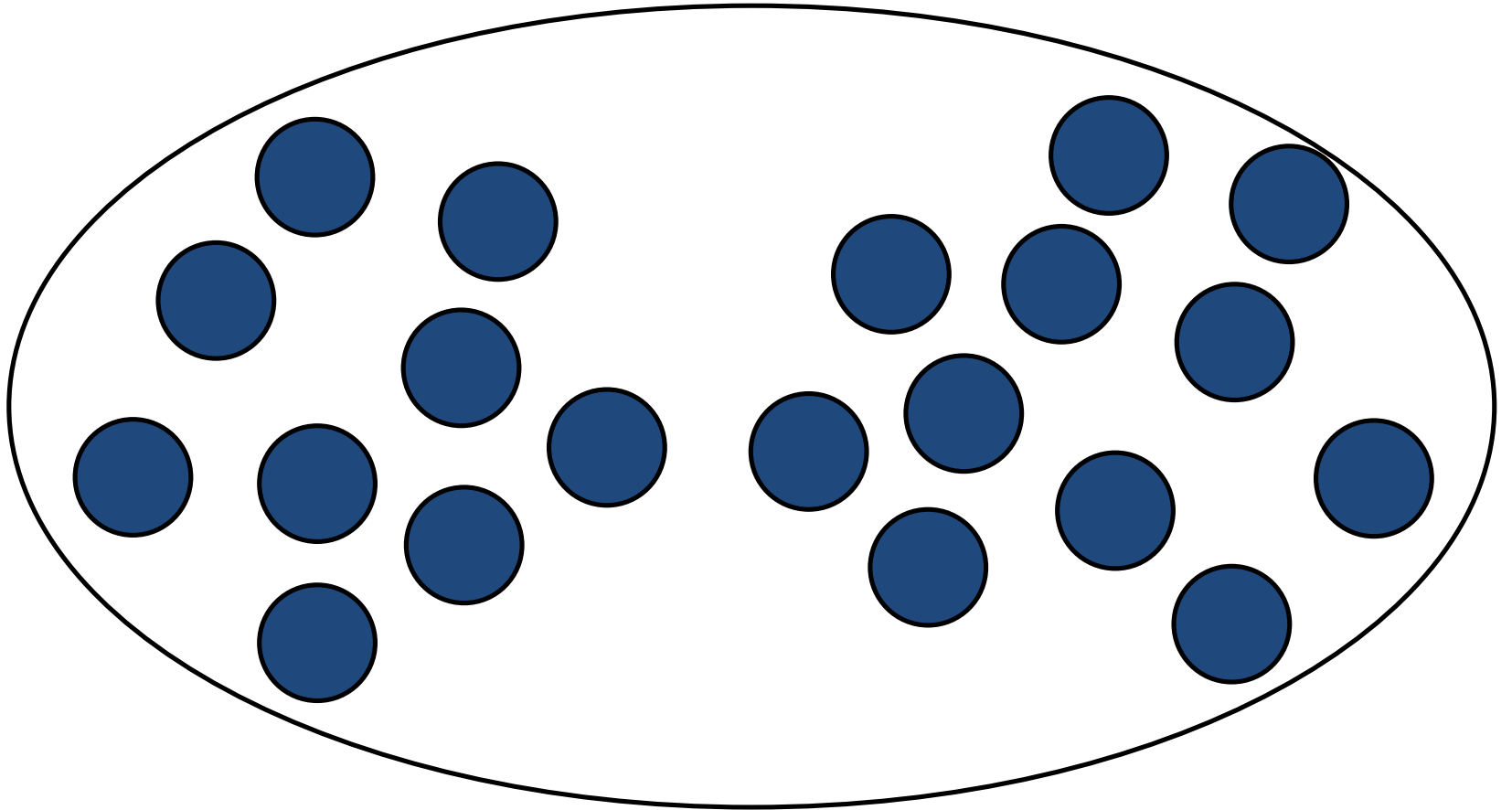
MPI_Init

- gather information about the parallel job
- set up internal library state
- prepare for communication

MPI_Finalize

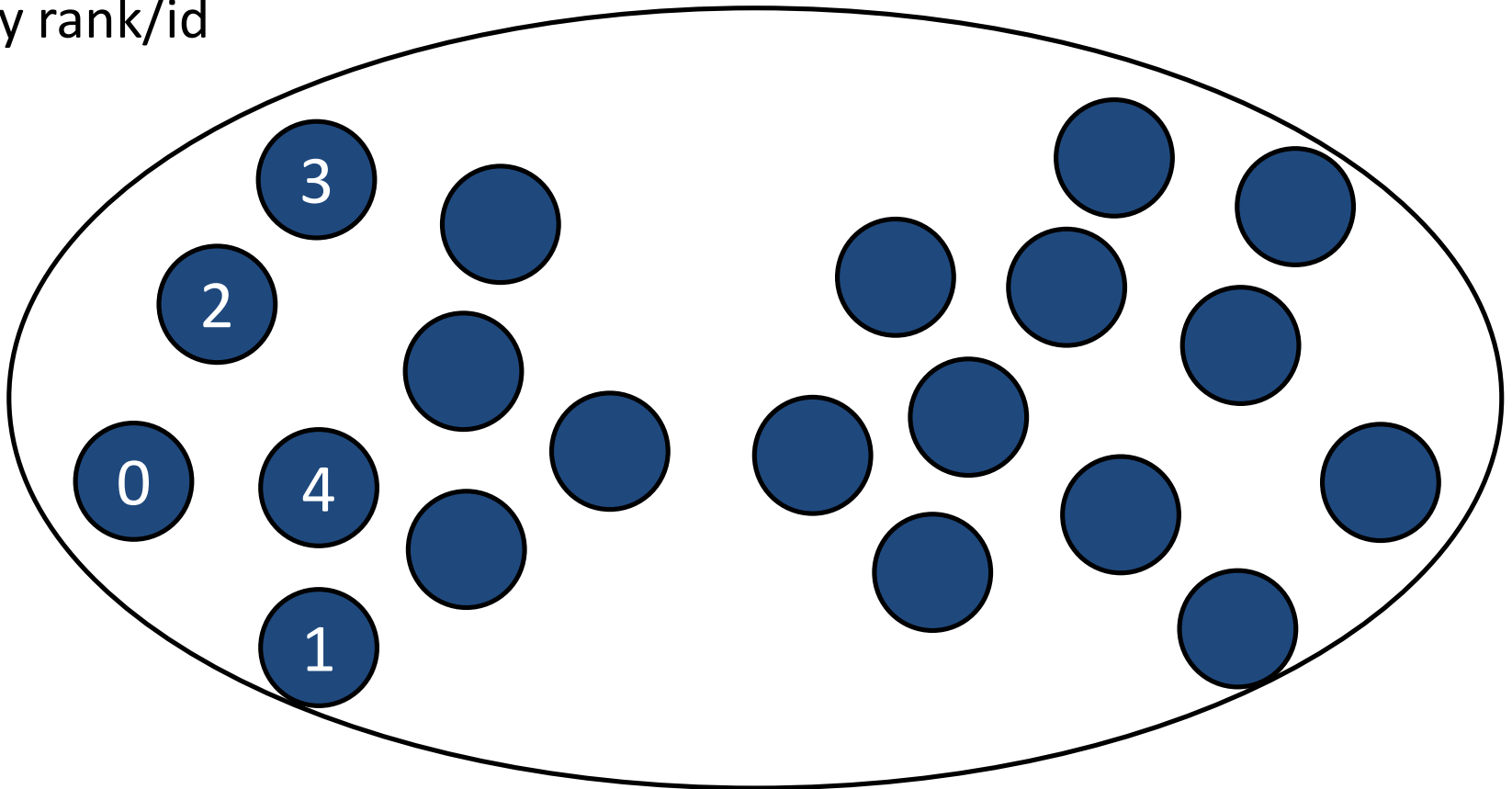
- cleanup

MPI_COMM_WORLD



MPI_COMM_WORLD

Process identified
by rank/id



Communication Scope

Communicator (communication handle)

- Defines the scope
- Specifies communication context

Process

- Belongs to a group
- Identified by a rank within a group

Identification

- `MPI_Comm_size` – total number of processes in communicator
- `MPI_Comm_rank` – rank in the communicator

Getting Started

```
#include <mpi.h>
#include <stdio.h>

int main(int argc, char** argv) {

    // Initialize the MPI environment
    MPI_Init(NULL, NULL);

    // Get the number of processes
    int size;
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    // Get the rank of the process
    int rank;
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);

    // Get the name of the processor
    char processor_name[MPI_MAX_PROCESSOR_NAME];
    int name_len;
    MPI_Get_processor_name(processor_name, &name_len);

    // Print off a hello world message
    printf("Hello I am rank %d out of %d processes\n", rank, size);

    // Finalize the MPI environment.
    MPI_Finalize();
}
```

Total
number of
processes

Rank of a
process

Executing MPI codes

```
mpicc -o executable program.c
```

```
mpirun -np 4 ./executable
```

MPI Message

- Data and header/envelope
- Typically, MPI communications send/receive messages

Message Envelope

Source: Origin of message

Destination: Receiver of message

Communicator

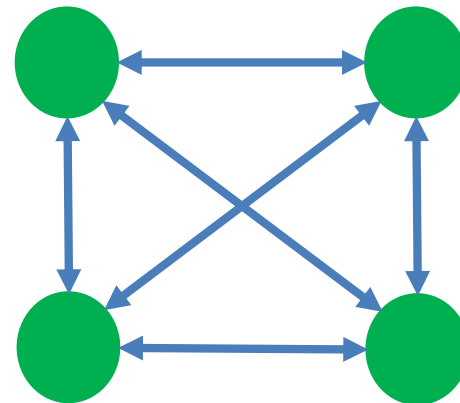
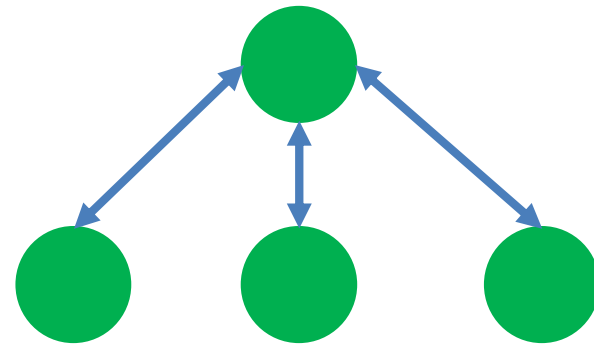
Tag (0:MPI_TAG_UB)

MPI Communication Types

Point-to-point

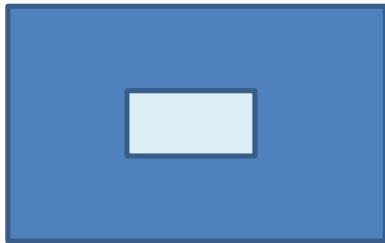


Collective



Point-to-point Communication

- MPI_Send



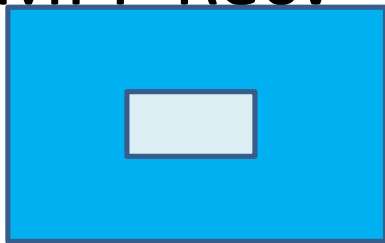
SENDER

Blocking send and receive

```
int MPI_Send (const void *buf, int count,  
MPI_Datatype datatype, int dest, int tag,  
MPI_Comm comm)
```

Tags should match

- MPI_Recv



RECEIVER

```
int MPI_Recv (void *buf, int count,  
MPI_Datatype datatype, int source, int tag,  
MPI_Comm comm, MPI_Status *status)
```

MPI_Datatype

- MPI_BYTE
- MPI_CHAR
- MPI_INT
- MPI_FLOAT
- MPI_DOUBLE

Example 1

```
MPI_Comm_rank (MPI_COMM_WORLD, &myrank);

// Sender process
if (myrank == 0)      /* code for process 0 */
{
    strcpy (message, "Hello, there");
    MPI_Send (message, strlen(message)+1, MPI_CHAR, 1, 99,
MPI_COMM_WORLD);
}

// Receiver process
else if (myrank == 1) /* code for process 1 */
{
    MPI_Recv (message, 20, MPI_CHAR, 0, 99, MPI_COMM_WORLD, &status);
    printf ("received :%s\n", message);
}
```



Message
tag



Message
tag

MPI_Status

- Source rank
- Message tag
- Message length
 - MPI_Get_count

MPI_ANY_*

- MPI_ANY_SOURCE
 - Receiver may specify wildcard value for source
- MPI_ANY_TAG
 - Receiver may specify wildcard value for tag

Example 2

```
MPI_Comm_rank (MPI_COMM_WORLD, &myrank);

// Sender process
if (myrank == 0 || myrank == 2) /* process 0 and 2 */
{
    sprintf (message, "Hello, there from %d", myrank);
    MPI_Send (message, strlen(message)+1, MPI_CHAR, 1, 99,
MPI_COMM_WORLD);
}

// Receiver process
else if (myrank == 1) /* process 1 */
{
    MPI_Recv (message, 40, MPI_CHAR, MPI_ANY_SOURCE, 99,
MPI_COMM_WORLD, &status);
    printf ("received :%s\n", message);
}
```



Bug?

Example 2 (correct)

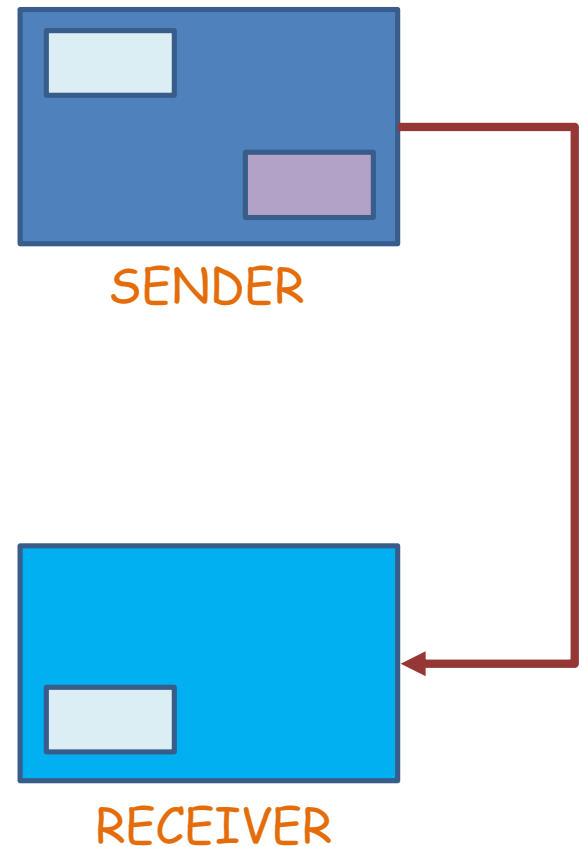
```
MPI_Comm_rank (MPI_COMM_WORLD, &myrank);

// Sender process
if (myrank == 0 || myrank == 2) /* process 0 and 2 */
{
    sprintf (message, "Hello, there from %d", myrank);
    MPI_Send (message, strlen(message)+1, MPI_CHAR, 1, 99,
MPI_COMM_WORLD);
}

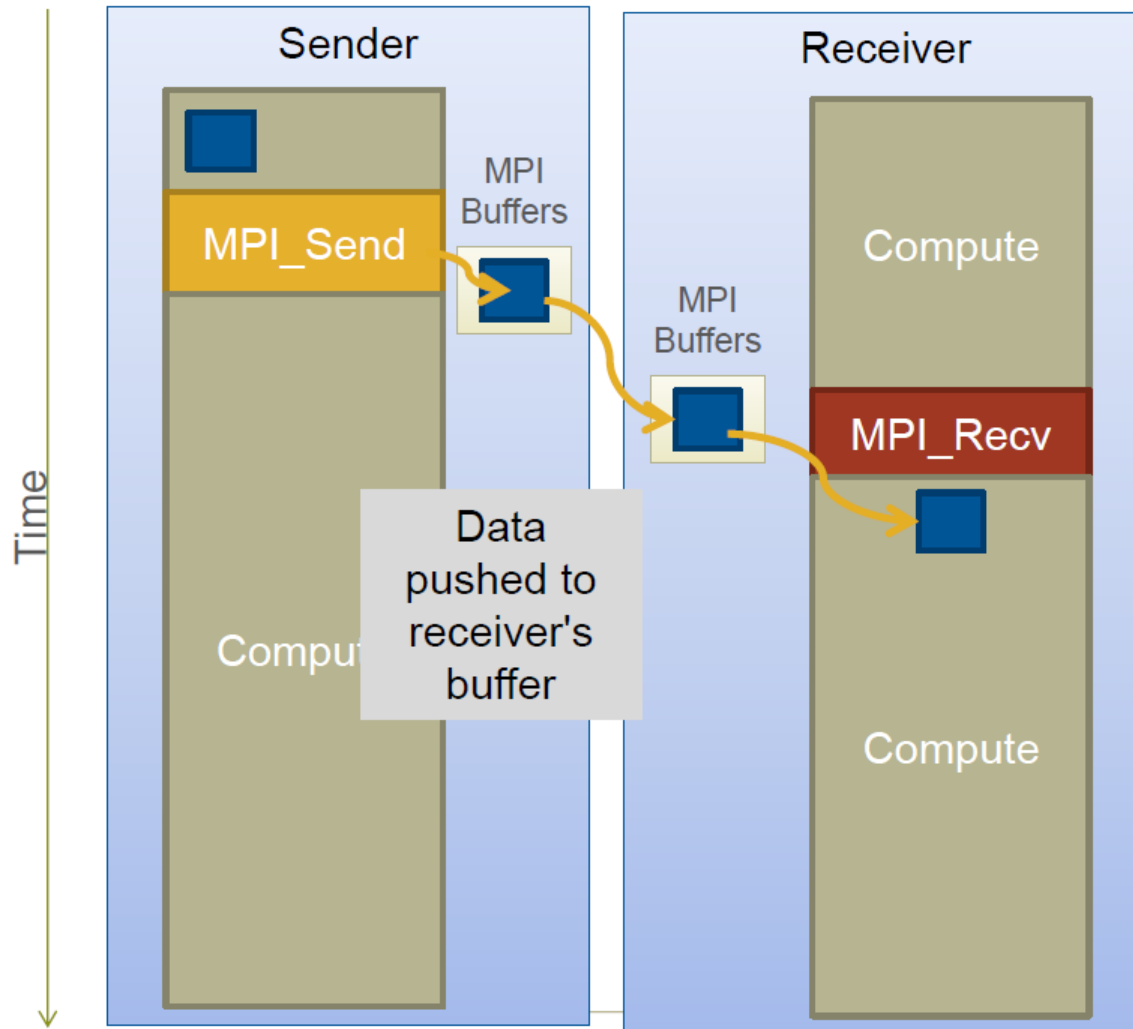
// Receiver process
else if (myrank == 1) /* process 1 */
{
    MPI_Recv (message, 20, MPI_CHAR, MPI_ANY_SOURCE, 99,
MPI_COMM_WORLD, &status);
    MPI_Recv (message, 20, MPI_CHAR, MPI_ANY_SOURCE, 99,
MPI_COMM_WORLD, &status);
}
```

MPI_Send (Blocking)

- Does not return until buffer can be reused
- Message buffering
- Implementation-dependent
- Standard communication mode



Buffering



Message Protocols

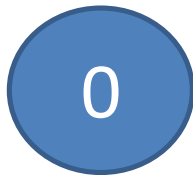
- Short
 - Message sent with envelope/header
- Eager
 - Send completes without acknowledgement from destination
 - Small messages – typically 128 KB (at least in MPICH)
 - `MPIR_CVAR_CH3_EAGER_MAX_MSG_SIZE` (check mpivars)
- Rendezvous
 - Requires an acknowledgement from a matching receive
 - Large messages

Other Send Modes

- MPI_Bsend Buffered
 - May complete before matching receive is posted
- MPI_Ssend Synchronous
 - Completes only if a matching receive is posted
- MPI_Rsend Ready
 - Started only if a matching receive is posted

Non-blocking Point-to-Point

- MPI_Isend (buf, count, datatype, dest, tag, comm, request)
- MPI_Irecv (buf, count, datatype, source, tag, comm, request)
- MPI_Wait (request, status)



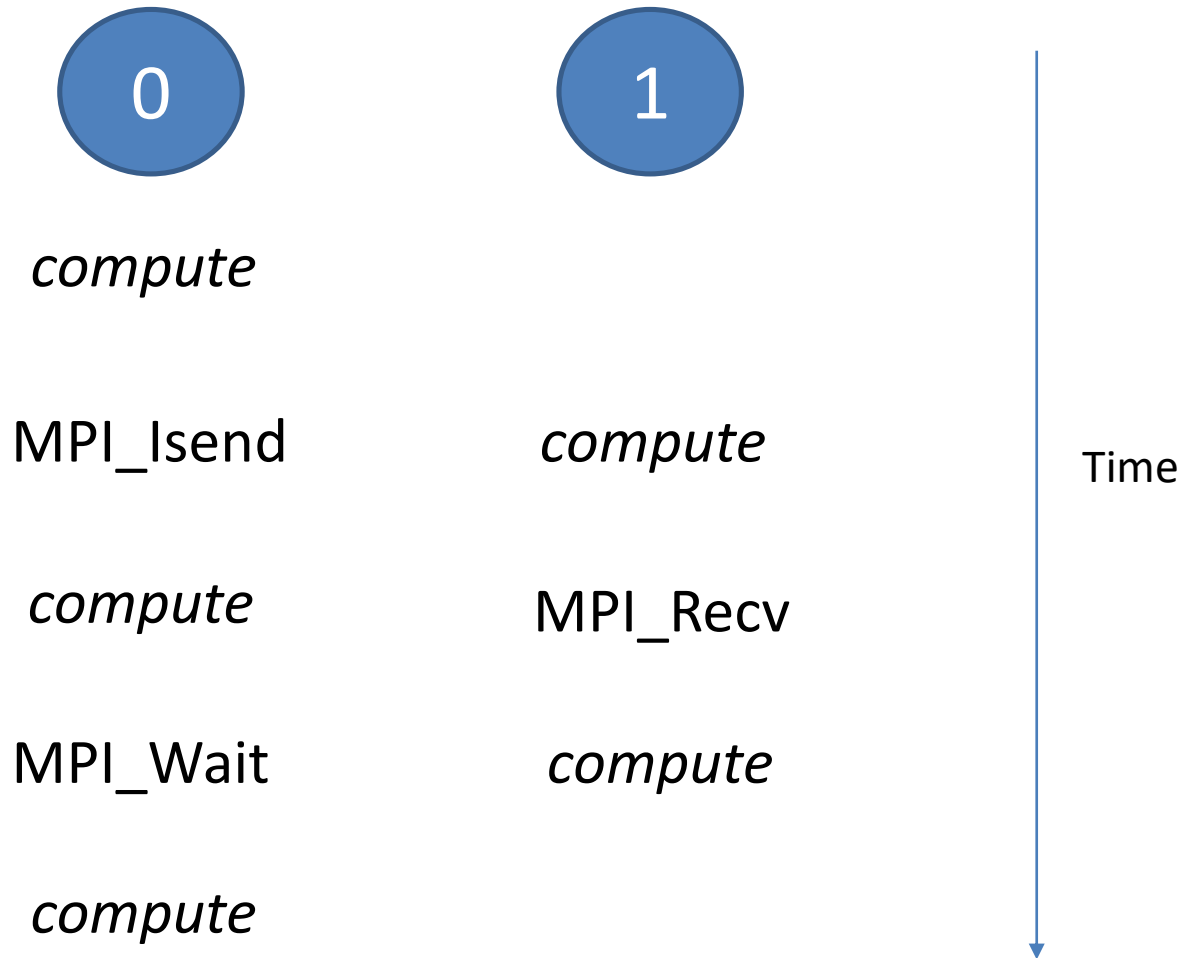
MPI_Isend
MPI_Irecv



MPI_Isend
MPI_Irecv

Safe

Computation Communication Overlap



Collective Communications

- Must be called by all processes that are part of the communicator

Types

- Synchronization (MPI_Barrier)
- Global communication (MPI_Bcast, MPI_Gather, ...)
- Global reduction (MPI_Reduce, ...)

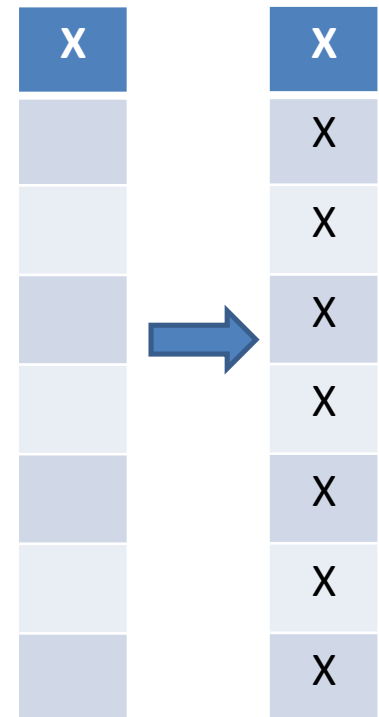
Barrier

- Synchronization across all group members
- Collective call
- Blocks until all processes have entered the call
- `MPI_Barrier (comm)`

Broadcast

- Root process sends message to all processes
- Any process can be root process but has to be the same in all processes
- `int MPI_Bcast (buffer, count, datatype, root, comm)`
- Number of elements in buffer – count

Q: Can you use point-to-point communication for the same?



Example 3

```
int rank, size, color;
MPI_Status status;

MPI_Init (&argc, &argv);
MPI_Comm_rank (MPI_COMM_WORLD, &rank);
MPI_Comm_size (MPI_COMM_WORLD, &size);

color = rank + 2;
int oldcolor = color;
MPI_Bcast (&color, 1, MPI_INT, 0, MPI_COMM_WORLD);

printf ("%d: %d color changed to %d\n", rank, oldcolor, color);
```

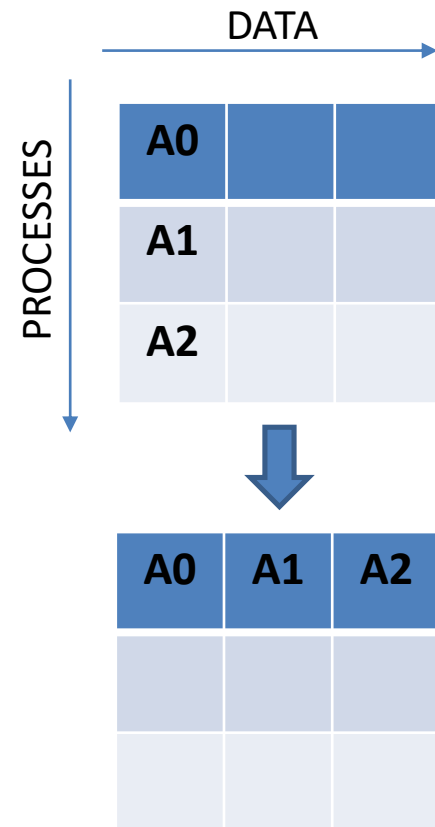
0: 2 color changed to 2

1: 3 color changed to 2

2: 4 color changed to 2

Gather

- Gathers values from all processes to a root process
- `int MPI_Gather (sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, root, comm)`
- Arguments `recv*` not relevant on non-root processes



Example 4

```
MPI_Comm_rank (MPI_COMM_WORLD, &rank);
MPI_Comm_size (MPI_COMM_WORLD, &size);

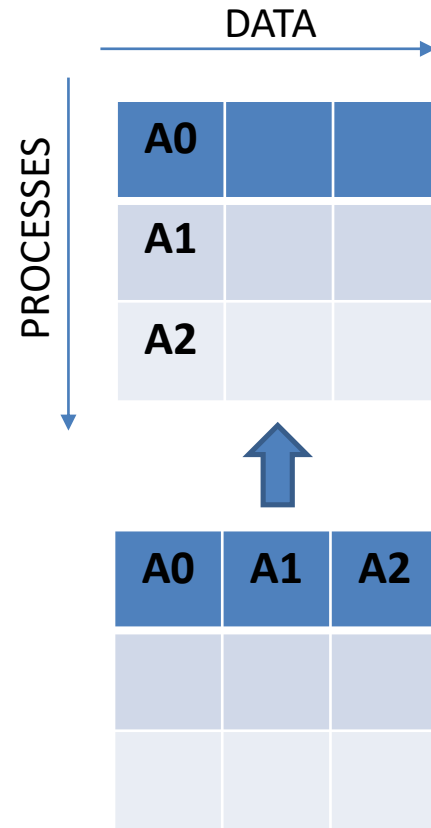
color = rank + 2;

int colors[size];
MPI_Gather (&color, 1, MPI_INT, colors, 1, MPI_INT, 0,
MPI_COMM_WORLD);

if (rank == 0)
    for (i=0; i<size; i++)
        printf ("color from %d = %d\n", i, colors[i]);
```

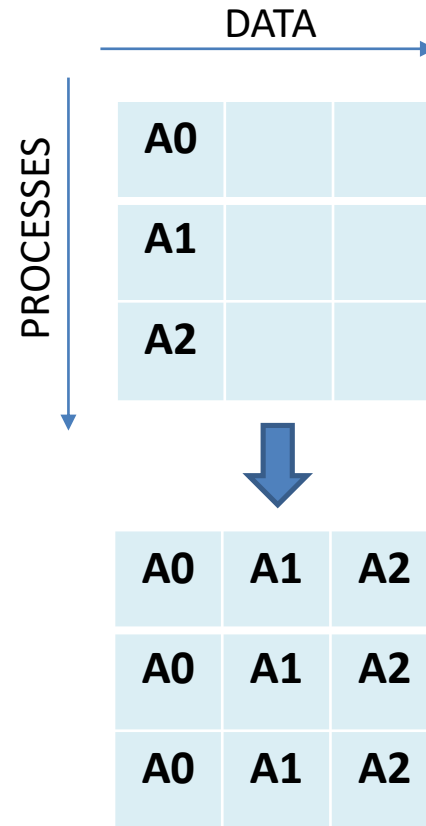
Scatter

- Scatters values to all processes from a root process
- `int MPI_Scatter (sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, root, comm)`
- Arguments `send*` not relevant on non-root processes
- Output parameter – `recvbuf`



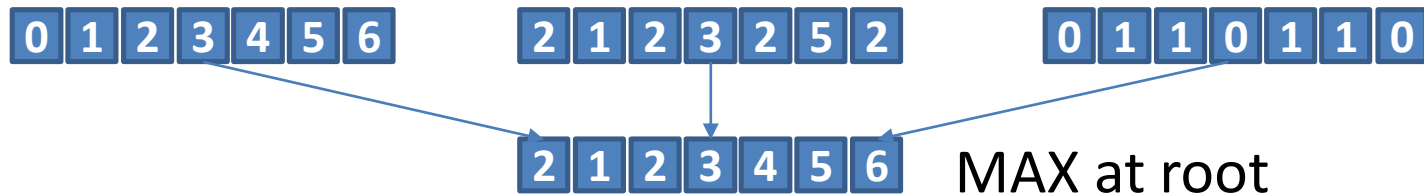
Allgather

- All processes gather values from all processes
- `int MPI_Allgather`
(`sendbuf`, `sendcount`,
`sendtype`, `recvbuf`,
`recvcount`, `recvtype`,
`comm`)



Reduce

- MPI_Reduce (inbuf, outbuf, count, datatype, op, root, comm)
- Combines element in inbuf of each process
- Combined value in outbuf of root
- op: MIN, MAX, SUM, PROD, ...



MPI_Reduce Example

```
int local_marks[number_of_TAs];  
  
// compute local marks in parallel  
  
MPI_Reduce (&local_marks,  
&total_marks, 1, MPI_FLOAT,  
MPI_SUM, 0, MPI_COMM_WORLD);
```



Instructor

60 hours



TAs

4 hours due to limitations

Allreduce

- MPI_Allreduce (inbuf, outbuf, count, datatype, op, comm)
- op: MIN, MAX, SUM, PROD, ...
- Combines element in inbuf of each process
- Combined value in outbuf of each process

0 1 2 3 4 5 6

2 1 2 3 2 5 2

0 1 1 0 1 1 0

2 1 2 3 4 5 6

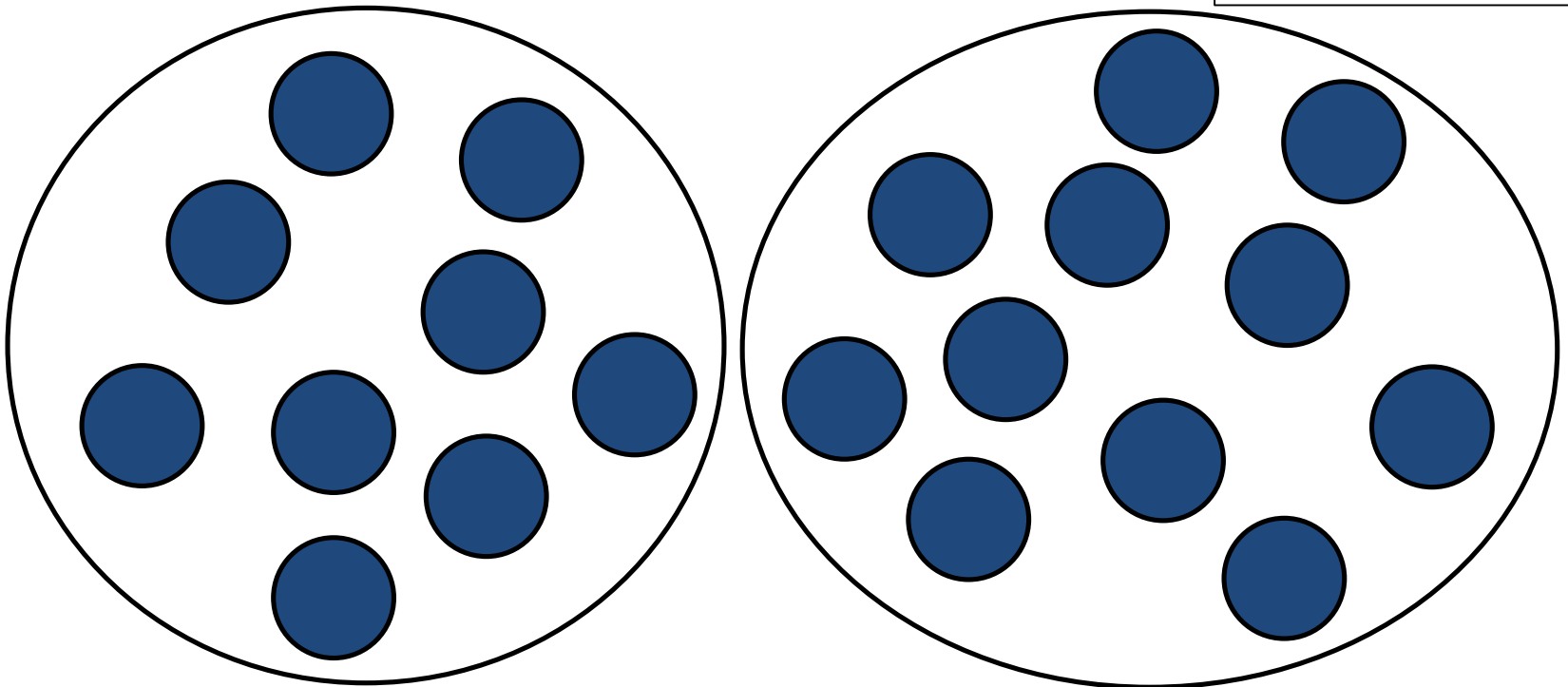
2 1 2 3 4 5 6

2 1 2 3 4 5 6

MAX

Sub-communicator

- Logical subset
- Different contexts

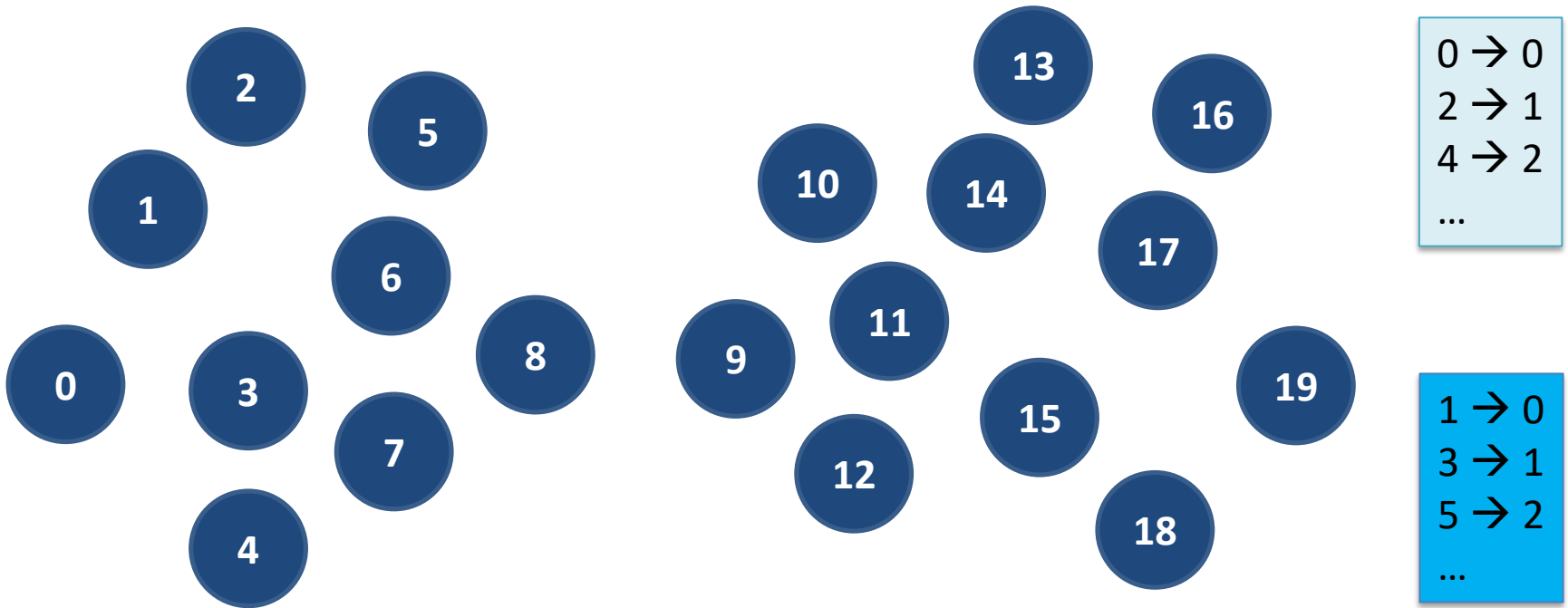


MPI_COMM_SPLIT

MPI_Comm_split (MPI_Comm oldcomm, int color, int key, MPI_Comm *newcomm)

- Collective call
- Logically divides based on *color*
 - Same color processes form a group
 - Some processes may not be part of newcomm (MPI_UNDEFINED)
- Rank assignment based on *key*

Logical subsets of processes



How do you assign one color to odd processes and another color to even processes ?
color = rank % 2

How to run an MPI program on a cluster?



MPI process
MPI process

MPI process
MPI process

MPI process
MPI process

MPI process
MPI process

4 nodes, ppn=2

`mpiexec -n <number of processes> -f <hostfile> ./exe`

```
<hostfile>
host1:2
host2:2
host3:2
...
```

How to run an MPI program on a managed cluster/supercomputer?



MPI process
MPI process

MPI process
MPI process

MPI process
MPI process

MPI process
MPI process

4 nodes, ppn=2

Execution on HPC2010: `qsub sub.sh`

How to run an MPI program on your lab cluster?

Install MPICH or MVAPICH (open source)

Where to install?

- Shared file system
- Mount on other systems

```
<hostfile>
```

```
host1:2
```

```
host2:2
```

```
host3:2
```

```
...
```

Reference Material

- Marc Snir, Steve W. Otto, Steven Huss-Lederman, David W. Walker and Jack Dongarra, MPI - The Complete Reference, Second Edition, Volume 1, The MPI Core.
- William Gropp, Ewing Lusk, Anthony Skjellum, Using MPI : portable parallel programming with the message-passing interface, 3rd Ed., Cambridge MIT Press, 2014.

Thank You

pmalakar@cse.iitk.ac.in