

## Parallel Programming and MPI- Lecture 3

Abhik Roychoudhury  
CS 3211  
National University of Singapore

Sample material: Parallel Programming by Lin and Snyder, Chapter 7.

1

CS3211 2009-10 by Abhik Roychoudhury

## Summary of previous lectures

- ▶ MPI as a programming interface
- ▶ Message passing communication
  - ▶ Communicating sequential processes
- ▶ Entering and Exiting MPI
  - ▶ MPI\_Init, MPI\_Finalize
- ▶ Point-to-point communication
  - ▶ Blocking & Non-blocking
  - ▶ MPI\_Send, MPI\_Recv, MPI\_Isend, MPI\_Irecv
  - ▶ Wait and test operations to complete communication.
- ▶ Collective communication

▶ 2

CS3211 2009-10 by Abhik Roychoudhury

## In today's discussion

- ▶ Managing communicators in MPI
  - ▶ Defines a communication domain.
  - ▶ Used implicitly several times in our discussion in defining the communication primitives.

▶ 3

CS3211 2009-10 by Abhik Roychoudhury

## Size and rank

- ▶ `int MPI_Comm_size(comm, &size);`
  - ▶ # of processes in the communicator
- ▶ `int MPI_Comm_rank(comm, &rank);`
  - ▶ Rank of the process that calls it
  - ▶ In the range  $0 \dots \text{size}-1$
- ▶ There is a pre-defined communicator
  - ▶ `MPI_COMM_WORLD`

▶ 4

CS3211 2009-10 by Abhik Roychoudhury

## So, what is a communicator?

- ▶ A group is an ordered set of processes.
- ▶ A communicator is a handle to a group of processes.
- ▶ A communicator thus defines a **communication domain**.
- ▶ Even for the same group of processes  $\langle p_1, \dots, p_N \rangle$ , it might be convenient to describe disparate communication domains containing the same group of processes.
  - ▶ Why?
    - ▶ To separate library code execution from user code execution.
    - ▶ A send in library may be received by a receive in user code.
    - ▶ This can be prevented by making the library and user code operate in different communication domains!

▶ 5

CS3211 2009-10 by Abhik Roychoudhury

## Intra- and Inter-communicators

Intra-communicator

For communication within a group of processes.

Inter-communicator

For point-to-point communication between disjoint groups of processes.

▶ 6

CS3211 2009-10 by Abhik Roychoudhury

## Can we ignore communicators?

- ▶ There is a single global communicator
  - ▶ MPI\_COMM\_WORLD
  - ▶ Contains all processes.
  - ▶ We can only work with this one.
- ▶ However, it may be advantageous to separate out certain communications, to prevent executions with arbitrary send-receive matching!

▶ 7

CS3211 2009-10 by Abhik Roychoudhury

## Creating communicators

- ▶ `int MPI_Comm_dup( comm, newcomm)`
  - ▶ MPI\_Comm comm
  - ▶ MPI\_Comm \*newcomm
  - ▶ Creates a new communicator with the same group of processes.
- ▶ `int MPI_Comm_create(comm, group, newcomm)`
  - ▶ MPI\_Comm comm
  - ▶ MPI\_Group group
  - ▶ MPI\_Comm \*newcomm
  - ▶ The argument group must be a subset of the group of comm
  - ▶ Always possible to use, with MPI\_COMM\_WORLD

▶ 8

CS3211 2009-10 by Abhik Roychoudhury

## Exercise

- ▶ We are trying to define a parallel library which does multi-cast (a variant of MPI\_Bcast)
  - ▶ Differences between MPI\_Bcast and our library
    - ▶ Instead of the root process in MPI\_Bcast, the function takes a flag which is true if the calling process is root, and false otherwise.
- ▶ Signature of MPI\_Bcast
  - ▶ `int MPI_Bcast(buffer, count, datatype, root, comm)`
    - ▶ Starting address of buffer
    - ▶ # of entries in buffer
    - ▶ Data type of buffer
    - ▶ Rank of the broadcasting process
    - ▶ The communicator capturing the group of processes.

▶ 9

CS3211 2009-10 by Abhik Roychoudhury

## Exercise

- ▶ Signature of mcast
  - ▶ `Mcast(buf, count, type, isroot, comm)`
    - ▶ Output buffer at root, input buffer at other processes
    - ▶ Number of items to be broadcast
    - ▶ Type of items to be broadcast
    - ▶ Flag saying whether the process is a root
    - ▶ Communicator.

▶ 10

CS3211 2009-10 by Abhik Roychoudhury

## Example Multi-cast library

```
void mcast(void *buff, int count, MPI_Datatype type,
           int isroot, MPI_Comm comm)
{
    MPI_Comm_size(comm, &size);
    MPI_Comm_rank(comm, &rank);

    int numleaves, /*number of leaves in broadcast tree */
        childleaves, /*number of leaves in child's broadcast tree */
        child; /* rank of child */
    if (isroot){
        numleaves = size - 1;
    }
    else{ ...
```

▶ 11

CS3211 2009-10 by Abhik Roychoudhury

## Example Multi-cast library

```
else { /* not a root process */
    MPI_Recv(&numleaves, 1, MPI_INT, MPI_ANY_SOURCE, 0, comm);
    MPI_Recv(buff, count, type, MPI_ANY_SOURCE, 0, comm);
}
while (numleaves > 0){
    /* pick child in the middle of current leaf processes */
    child = (rank + (numleaves + 1) / 2) % size;
    childleaves = numleaves / 2;
    /* send leaf count and message to child */
    MPI_Send(&childleaves, 1, MPI_INT, child, 0, comm);
    MPI_Send(buff, count, type, child, 0);
    numleaves -= (childleaves + 1); /* remaining number of leaves */
}
}
```

▶ 12

CS3211 2009-10 by Abhik Roychoudhury

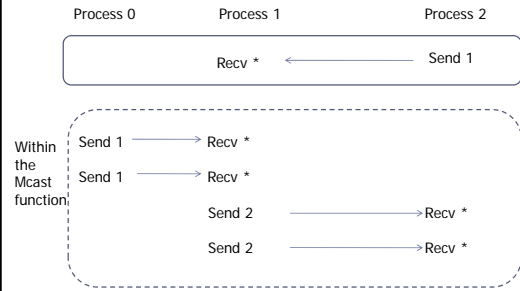
Now, consider the following code

```

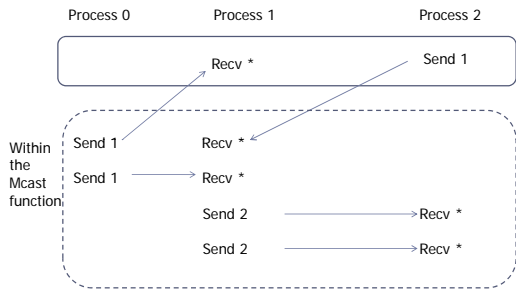
> MPI_Comm_rank( comm, &myrank );
> if (myrank == 2){
>   MPI_Send( ..., 1, MPI_INT, 1, 0, comm);
> } else if (myrank == 1){
>   MPI_Recv(..., 1, MPI_INT, MPI_ANY_SOURCE, 0, comm)
> }
> mcast(..., 1, MPI_INT, (myrank == 0), comm);

```

“Expected” Behavior



Possible & “Unexpected” Behavior



How can this happen?

- ▶ Process 0 starts executing multi-cast earlier than other processes.
  - ▶ The processes are executing on different processors after all.
  - ▶ Different processors run at different speeds!
- ▶ Process 1 executes the MPI\_Recv in the caller code
  - ▶ This matches with the first MPI\_Send of process 0 executed inside the mcast library!
- ▶ This is why separate communicators are needed!

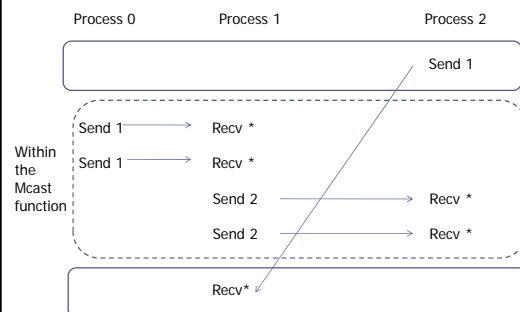
A more complex scenario

```

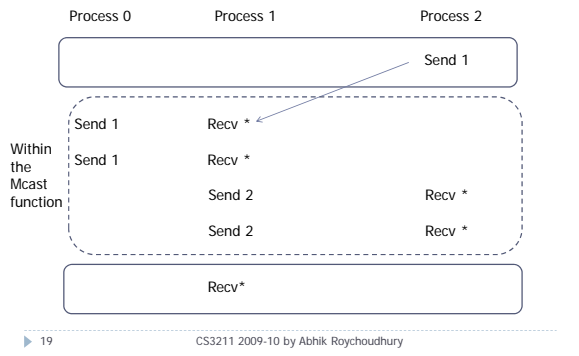
> MPI_Comm_rank( comm, &myrank );
> if (myrank == 2){
>   MPI_Send( ..., 1, MPI_INT, 1, 0, comm);
> }
> mcast(..., 1, MPI_INT, (myrank == 0), comm);
> if (myrank == 1){
>   MPI_Recv(..., 1, MPI_INT, MPI_ANY_SOURCE, 0, comm)
> }

```

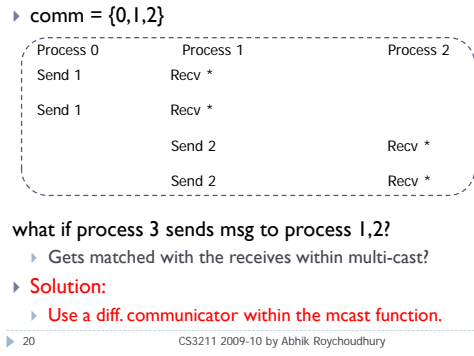
The “expected” behavior



## The “unexpected” behavior



## Another possibility



## Using different communicators

```
void mcast(void *buf, int count, MPI_Datatype type,
           int isroot, MPI_Comm comm)
{
    int size, rank, numleaves, child, childleaves;
    MPI_Status status;
    MPI_Comm pcomm; /* private communicator */

    MPI_Comm_dup(comm, &pcomm);
    MPI_Comm_size(pcomm, &size);
    MPI_Comm_rank(pcomm, &rank);

    /* dynamically build up a broadcast tree now */
```

▶ 21      CS3211 2009-10 by Abhik Roychoudhury

## Using different communicators

```
▶ if (isroot){
    numleaves = size - 1;
} else { /* receive from parent */
    MPI_Recv(&numleaves, 1, MPI_INT, MPI_ANY_SOURCE, 0, pcomm, &status);
    MPI_Recv(buf, count, type, MPI_ANY_SOURCE, 0, pcomm, &status);
}
while (numleaves > 0){
    child = (rank + (numleaves + 1) / 2) % size;
    childleaves = numleaves / 2; /* send to child in the next 2 lines */
    MPI_Send(&childleaves, 1, MPI_INT, child, 0, pcomm);
    MPI_Send(buf, count, type, child, 0, pcomm);
    numleaves -= (childleaves + 1); /* compute remaining number of leaves */
}
MPI_Comm_free(&pcomm);
}
```

▶ 22      CS3211 2009-10 by Abhik Roychoudhury

## Exercise

- ▶ Can there be other solutions which avoid the additional communicator allocation (pcomm)?
  - ▶ How about inserting a barrier at the beginning and at the end of the mcast function?
  - ▶ Can this solution be consistently employed for any parallel library?
  - ▶ What are the implications on
    - ▶ Performance?
    - ▶ Correctness?
      - Try out the 2 communication scenarios we discussed earlier.

▶ 23      CS3211 2009-10 by Abhik Roychoudhury

## Wrapping up

- ▶ MPI programming
  - ▶ Explicit message passing, as opposed to shared memory.
- ▶ Important concepts
  - ▶ Point to point communication
    - ▶ Blocking send receives --- MPI\_Send, MPI\_Recv
    - ▶ Non-blocking send receives ---- MPI\_Isend, MPI\_Irecv
  - ▶ Collective communication
    - ▶ Scatter, Gather
    - ▶ MPI\_Reduce
  - ▶ Communicators
    - ▶ The default communicator is MPI\_COMM\_WORLD

▶ 24      CS3211 2009-10 by Abhik Roychoudhury