# CC755: Distributed and Parallel Systems

**Lecture 10: Programming Using the Message Passing Paradigm II**

**Slides by: Ananth Grama, Anshul Gupta, George Karypis, and Vipin Kumar**

**To accompany the text "Introduction to Parallel Computing", Addison Wesley, 2003.**
**With some edits from other sources**

## Dr. Manal Helal, Spring 2016

moodle.manalhelal.com

# Topic Overview

- Topologies and Embedding
- Overlapping Communication with Computation
- Collective Communication and Computation Operations
- Groups and Communicators

# Communicators

- A communicator defines a *communication domain* - a set of processes that are allowed to communicate with each other.

- Information about communication domains is stored in variables of type `MPI_Comm`.

- Communicators are used as arguments to all message transfer MPI routines.

- A process can belong to many different (possibly overlapping) communication domains.

- MPI defines a default communicator called `MPI_COMM_WORLD` which includes all the processes.

# Avoiding Deadlocks

## Consider:

```
int a[10], b[10], myrank;
MPI_Status status;
...
MPI_Comm_rank(MPI_COMM_WORLD, &myrank);
if (myrank == 0) {
    MPI_Send(a, 10, MPI_INT, 1, 1, MPI_COMM_WORLD);
    MPI_Send(b, 10, MPI_INT, 1, 2, MPI_COMM_WORLD);
}
else if (myrank == 1) {
    MPI_Recv(b, 10, MPI_INT, 0, 2, MPI_COMM_WORLD);
    MPI_Recv(a, 10, MPI_INT, 0, 1, MPI_COMM_WORLD);
}
...
```

If MPI_Send is blocking, there is a deadlock.

# Avoiding Deadlocks

Consider the following piece of code, in which process i sends a message to process $i$ + 1 (modulo the number of processes) and receives a message from process $i$ - 1 (module the number of processes).

```
int a[10], b[10], npes, myrank;
MPI_Status status;
...
MPI_Comm_size(MPI_COMM_WORLD, &npes);
MPI_Comm_rank(MPI_COMM_WORLD, &myrank);
MPI_Send(a, 10, MPI_INT, (myrank+1)%npes, 1,
      MPI_COMM_WORLD);
MPI_Recv(b, 10, MPI_INT, (myrank-1+npes)%npes, 1,
MPI_COMM_WORLD);
...
```

Once again, we have a deadlock if MPI_Send is blocking.

# Avoiding Deadlocks

We can break the circular wait to avoid deadlocks as follows:

```
int a[10], b[10], npes, myrank;
MPI_Status status;
...
MPI_Comm_size(MPI_COMM_WORLD, &npes);
MPI_Comm_rank(MPI_COMM_WORLD, &myrank);
if (myrank%2 == 1) {
        MPI_Send(a, 10, MPI_INT, (myrank+1)%npes, 1,
                MPI_COMM_WORLD);
        MPI_Recv(b, 10, MPI_INT, (myrank-1+npes)%npes, 1,
                MPI_COMM_WORLD);
}
else {
        MPI_Recv(b, 10, MPI_INT, (myrank-1+npes)%npes, 1,
                MPI_COMM_WORLD);
        MPI_Send(a, 10, MPI_INT, (myrank+1)%npes, 1,
                MPI_COMM_WORLD);
}
...
```

# Sending and Receiving Messages Simultaneously

To exchange messages, MPI provides the following function:

```
int MPI_Sendrecv(void *sendbuf, int sendcount,
    MPI_Datatype senddatatype, int dest, int
    sendtag, void *recvbuf, int recvcount,
    MPI_Datatype recvdatatype, int source, int recvtag,
    MPI_Comm comm, MPI_Status *status)
```

The arguments include arguments to the send and receive functions. If we wish to use the same buffer for both send and receive, we can use:

```
int MPI_Sendrecv_replace(void *buf, int count,
    MPI_Datatype datatype, int dest, int sendtag,
    int source, int recvtag, MPI_Comm comm,
    MPI_Status *status)
```

# Topologies and Embeddings

- MPI allows a programmer to organise processors into logical $k$-d meshes.

- The processor ids in `MPI_COMM_WORLD` can be mapped to other communicators (corresponding to higher-dimensional meshes) in many ways.

- The goodness of any such mapping is determined by the interaction pattern of the underlying program and the topology of the machine.

- MPI does not provide the programmer any control over these mappings.

# Topologies and Embeddings



(a) Row–major mapping  (b) Column–major mapping  (c) Space–filling curve mapping  (d) Hypercube mapping

Different ways to map a set of processes to a two-dimensional grid. (a) and (b) show a row- and column-wise mapping of these processes, (c) shows a mapping that follows a space-filling curve (dotted line), and (d) shows a mapping in which neighbouring processes are directly connected in a hypercube.

# Creating and Using Cartesian Topologies

- We can create cartesian topologies using the function:

```
int MPI_Cart_create(MPI_Comm comm_old, int ndims,
                                int *dims, int *periods, int
    reorder,                                MPI_Comm *comm_cart)
```

  This function takes the processes in the old communicator and creates a new communicator with dims dimensions.

- Each processor can now be identified in this new cartesian topology by a vector of dimension dims.

# Creating and Using Cartesian Topologies

- Since sending and receiving messages still require (one-dimensional) ranks, MPI provides routines to convert ranks to cartesian coordinates and vice-versa.

```
int MPI_Cart_coord(MPI_Comm comm_cart, int rank, int maxdims,
                   int *coords)

int MPI_Cart_rank(MPI_Comm comm_cart, int *coords, int *rank)
```

- The most common operation on cartesian topologies is a shift. To determine the rank of source and destination of such shifts, MPI provides the following function:

```
int MPI_Cart_shift(MPI_Comm comm_cart, int dir, int s_step,
                   int *rank_source, int *rank_dest)
```

# Overlapping Communication with Computation

- In order to overlap communication with computation, MPI provides a pair of functions for performing non-blocking send and receive operations.

```
int MPI_Isend(void *buf, int count, MPI_Datatype datatype,
                 int dest, int tag, MPI_Comm comm,
                 MPI_Request *request)
int MPI_Irecv(void *buf, int count, MPI_Datatype datatype,
                 int source, int tag, MPI_Comm comm,
                 MPI_Request *request)
```

- These operations return before the operations have been completed. Function `MPI_Test` tests whether or not the non-blocking send or receive operation identified by its request has finished.

```
int MPI_Test(MPI_Request *request, int *flag,
         MPI_Status *status)
```

- `MPI_Wait` waits for the operation to complete.

```
int MPI_Wait(MPI_Request *request, MPI_Status *status)
```

# Avoiding Deadlocks

Using non-blocking operations remove most deadlocks. Consider:

```
int a[10], b[10], myrank;
MPI_Status status;
...
MPI_Comm_rank(MPI_COMM_WORLD, &myrank);
if (myrank == 0) {
   MPI_Send(a, 10, MPI_INT, 1, 1, MPI_COMM_WORLD);
   MPI_Send(b, 10, MPI_INT, 1, 2, MPI_COMM_WORLD);
}
else if (myrank == 1) {
   MPI_Recv(b, 10, MPI_INT, 0, 2, &status, MPI_COMM_WORLD);
   MPI_Recv(a, 10, MPI_INT, 0, 1, &status, MPI_COMM_WORLD);
}
...
```

Replacing either the send or the receive operations with non-blocking counterparts fixes this deadlock.

# Collective Communication and Computation Operations

- MPI provides an extensive set of functions for performing common collective communication operations.

- Each of these operations is defined over a group corresponding to the communicator.

- All processors in a communicator must call these operations.

# Collective Communication Operations

- The barrier synchronisation operation is performed in MPI using:

    ```
    int MPI_Barrier(MPI_Comm comm)
    ```

    The one-to-all broadcast operation is:

    ```
    int MPI_Bcast(void *buf, int count, MPI_Datatype datatype,
            int source, MPI_Comm comm)
    ```

- The all-to-one reduction operation is:

    ```
    int MPI_Reduce(void *sendbuf, void *recvbuf, int count,
            MPI_Datatype datatype, MPI_Op op, int target,
            MPI_Comm comm)
    ```

# Predefined Reduction Operations

| Operation | Meaning | Datatypes |
|---|---|---|
| MPI_MAX | Maximum | C integers and floating point |
| MPI_MIN | Minimum | C integers and floating point |
| MPI_SUM | Sum | C integers and floating point |
| MPI_PROD | Product | C integers and floating point |
| MPI_LAND | Logical AND | C integers |
| MPI_BAND | Bit-wise AND | C integers and byte |
| MPI_LOR | Logical OR | C integers |
| MPI_BOR | Bit-wise OR | C integers and byte |
| MPI_LXOR | Logical XOR | C integers |
| MPI_BXOR | Bit-wise XOR | C integers and byte |
| MPI_MAXLOC | max-min value-location | Data-pairs |
| MPI_MINLOC | min-min value-location | Data-pairs |

# Collective Communication Operations

- The operation `MPI_MAXLOC` combines pairs of values $(v_i, l_i)$ and returns the pair $(v, l)$ such that $v$ is the maximum among all $v_i$'s and $l$ is the corresponding $l_i$ (if there are more than one, it is the smallest among all these $l_i$'s).

- `MPI_MINLOC` does the same, except for minimum value of $v_i$.

| Value | 15 | 17 | 11 | 12 | 17 | 11 |
|---|---|---|---|---|---|---|
| Process | 0 | 1 | 2 | 3 | 4 | 5 |

```
MinLoc(Value, Process) = (11, 2)
MaxLoc(Value, Process) = (17, 1)
```

An example use of the `MPI_MINLOC` and `MPI_MAXLOC` operators.

# Collective Communication Operations

MPI datatypes for data-pairs used with the `MPI_MAXLOC` and `MPI_MINLOC` reduction operations.

| MPI Datatype | C Datatype |
|---|---|
| `MPI_2INT` | pair of ints |
| `MPI_SHORT_INT` | short and int |
| `MPI_LONG_INT` | long and int |
| `MPI_LONG_DOUBLE_INT` | long double and int |
| `MPI_FLOAT_INT` | float and int |
| `MPI_DOUBLE_INT` | double and int |

# Collective Communication Operations

- If the result of the reduction operation is needed by all processes, MPI provides:

```
int MPI_Allreduce(void *sendbuf, void *recvbuf,
                  int count, MPI_Datatype datatype,
                  MPI_Op op, MPI_Comm comm)
```

- To compute prefix-sums, MPI provides:

```
int MPI_Scan(void *sendbuf, void *recvbuf, int
             count,MPI_Datatype datatype, MPI_Op op,
             MPI_Comm comm)
```

# Collective Communication Operations

- The gather operation is performed in MPI using:

```
int MPI_Gather(void *sendbuf, int sendcount,
               MPI_Datatype senddatatype, void *recvbuf,
               int recvcount, MPI_Datatype recvdatatype,
               int target, MPI_Comm comm)
```

- MPI also provides the MPI_Allgather function in which the data are gathered at all the processes.

```
int MPI_Allgather(void *sendbuf, int sendcount,
          MPI_Datatype senddatatype, void *recvbuf,
               int recvcount, MPI_Datatype recvdatatype,
               MPI_Comm comm)
```

- The corresponding scatter operation is:

```
int MPI_Scatter(void *sendbuf, int sendcount,
               MPI_Datatype senddatatype, void *recvbuf,
               int recvcount, MPI_Datatype recvdatatype,
               int source, MPI_Comm comm)
```

# Collective Communication Operations

- The all-to-all personalized communication operation is performed by:

```
int MPI_Alltoall(void *sendbuf, int sendcount,
    MPI_Datatype senddatatype, void *recvbuf, int
    recvcount, MPI_Datatype recvdatatype,
    MPI_Comm comm)
```
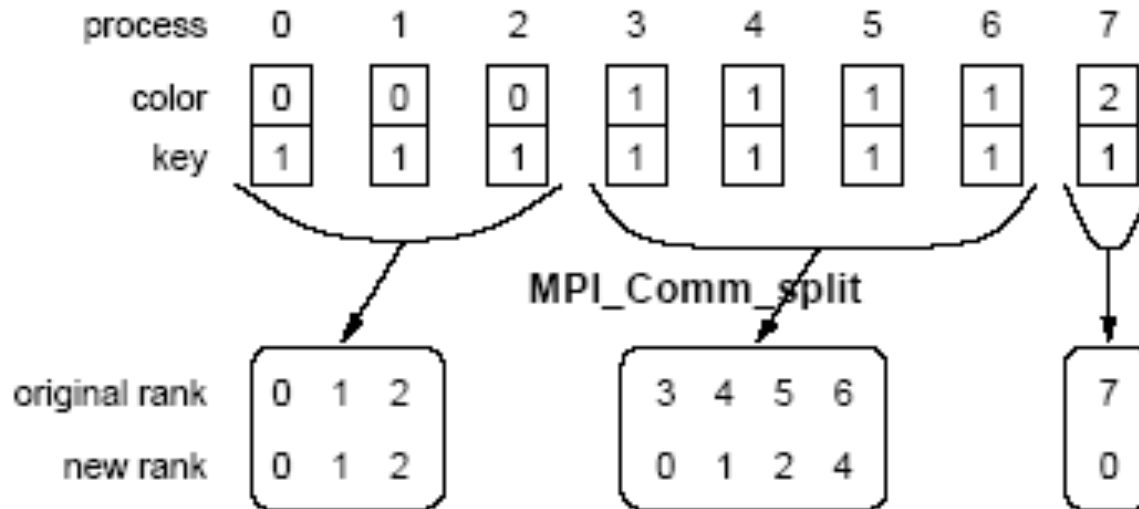
- Using this core set of collective operations, a number of programs can be greatly simplified.

# Groups and Communicators

- In many parallel algorithms, communication operations need to be restricted to certain subsets of processes.

- MPI provides mechanisms for partitioning the group of processes that belong to a communicator into subgroups each corresponding to a different communicator.

- The simplest such mechanism is:

```
int MPI_Comm_split(MPI_Comm comm, int color, int
                   key, MPI_Comm *newcomm)
```

- This operation groups processors by color and sorts resulting groups on the key.

# Groups and Communicators



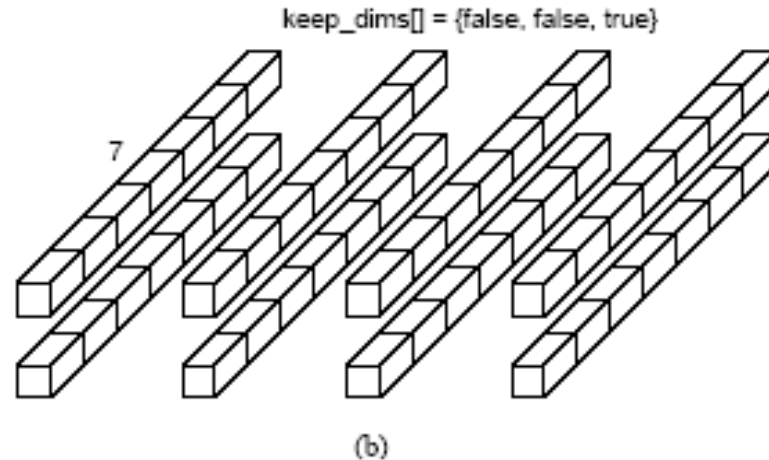Using `MPI_Comm_split` to split a group of processes in a communicator into subgroups.
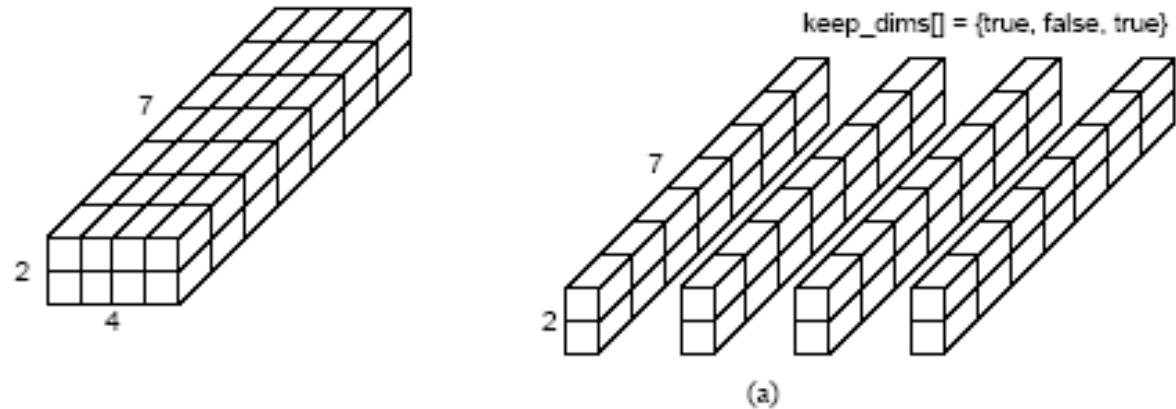
# Groups and Communicators

- In many parallel algorithms, processes are arranged in a virtual grid, and in different steps of the algorithm, communication needs to be restricted to a different subset of the grid.

- MPI provides a convenient way to partition a Cartesian topology to form lower-dimensional grids:

```
int MPI_Cart_sub(MPI_Comm comm_cart, int *keep_dims,
                 MPI_Comm *comm_subcart)
```

- If `keep_dims[i]` is true (non-zero value in C) then the `i`th dimension is retained in the new sub-topology.

- The coordinate of a process in a sub-topology created by `MPI_Cart_sub` can be obtained from its coordinate in the original topology by disregarding the coordinates that correspond to the dimensions that were not retained.

# Groups and Communicators



keep_dims[] = {true, false, true}

(a)

keep_dims[] = {false, false, true}

(b)

Splitting a Cartesian topology of size 2 x 4 x 7 into (a) four subgroups of size 2 x 1 x 7, and (b) eight subgroups of size 1 x 1 x 7.