

ACM/CS 114

Parallel algorithms for scientific applications

Michael A. G. Aïvázis

California Institute of Technology

Winter 2012

Point to point communication

- ▶ to send a message

```
1 int MPI_Send(  
2     void* buffer, int count, MPI_Datatype datatype,  
3     int destination, int tag, MPI_Comm communicator  
4 );
```

- ▶ to receive a message

```
1 int MPI_Recv(  
2     void* buffer, int count, MPI_Datatype datatype,  
3     int source, int tag, MPI_Comm communicator  
4 );
```

- ▶ the tag enables choosing the order you may receive pending messages
- ▶ but for a given (source,tag,communicator) messages are received in the order they were sent
- ▶ receiving via wildcards: MPI_ANY_SOURCE and MPI_ANY_TAG
- ▶ in *standard* communication mode, sending and receiving messages are *blocking*, so the function does not return until you can safely access the buffer
 - ▶ to read, free, etc.

Communication modes

- ▶ in standard mode, the specification does not explicitly mention buffering strategy
 - ▶ buffering messages would remove some of the access constraints but it requires time and storage for the multiple copies
 - ▶ portability across implementations implies conservative assumptions about the order of initiation of sends and receives to avoid deadlock
- ▶ in *ready* mode, you must post a receive before the matching send can be initiated
 - ▶ `MPI_Rsend`, `MPI_Rrecv`
- ▶ in *buffered* mode, sends can be initiated, and may complete, regardless of when the matching receive is initiate
 - ▶ `MPI_Bsend`, `MPI_Brecv`
- ▶ in *synchronous* mode, sends can be initiated regardless of whether the matching receive has been initiated, but the send will not return until the message has been received
 - ▶ `MPI_Ssend`, `MPI_Srecv`

Asynchronous communication

- ▶ there are non-blocking versions of all these

```
1 int MPI_Isend(  
2     void* buffer, int count, MPI_Datatype datatype,  
3     int destination, int tag,  
4     MPI_Comm communicator, MPI_Request* request  
5     );
```

- ▶ faster, but you must take care to not access the message buffers until the messages have been delivered
 - ▶ more details later in the course, as needed
- ▶ for sends
 - ▶ standard mode: `MPI_Isend`
 - ▶ ready mode: `MPI_Irsend`
 - ▶ buffered mode: `MPI_Ibsend`
 - ▶ synchronous mode: `MPI_Issend`
- ▶ only one call for receives: `MPI_Irecv`
- ▶ extra `request` argument to check for completion of the request
 - ▶ `MPI_Test`, `MPI_Wait` and their relatives

Creating communicators and groups

- ▶ communicators and groups are intertwined
 - ▶ you cannot create a group without a communicator
 - ▶ you cannot create a communicator without a group
- ▶ the cycle is broken by `MP I_COMM_WORLD`

```
1  #include <mpi.h>
2
3  int main(int argc, char* argv[]) {
4      /* declare a communicator and a couple of groups */
5      int loner = 0;
6      MPI_Comm workers;
7      MPI_Group world_grp, workers_grp;
8
9      /* initialize MPI; for brevity all status checks are omitted */
10     MPI_Init(&argc, &argv);
11
12     /* get the world communicator to build its group */
13     MPI_Comm_group(MPI_COMM_WORLD, &world_grp);
14
15     /* build another group by excluding a process */
16     MPI_Group_excl(world_grp, 1, &loner, &workers_grp);
17
18     /* now build a communicator out of the processes in workers_grp */
19     MPI_Comm_create(MPI_COMM_WORLD, worker_grp, &workers);
20
21     /* etc.... */
22
23     /* shut down MPI */
24     MPI_Finalize();
25
26     return 0;
27 }
```

Manipulating communicators and groups

- ▶ releasing resources

```
1 int MPI_Group_free(MPI_Group* group);
2 int MPI_Comm_free(MPI_Comm* communicator);
3 int MPI_Comm_disconnect(MPI_Comm* communicator);
```

- ▶ you can make a new group by adding or removing processes from an existing one

```
1 int MPI_Group_incl(
2     MPI_Group grp, int n, int* ranks, MPI_Group* new_group);
3 int MPI_Group_excl(
4     MPI_Group grp, int n, int* ranks, MPI_Group* new_group);
```

- ▶ or by using set operations

```
1 int MPI_Group_union(
2     MPI_Group grp1, MPI_Group grp2, MPI_Group* new_group);
3 int MPI_Group_intersection(
4     MPI_Group grp1, MPI_Group grp2, MPI_Group* new_group);
5 int MPI_Group_difference(
6     MPI_Group grp1, MPI_Group grp2, MPI_Group* new_group);
```