

ACM/CS 114

Parallel algorithms for scientific applications

Michael A. G. Aïvázis

California Institute of Technology

Winter 2010

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     MPI_Comm workers;
6     MPI_Group world_grp, workers_grp;
7
8     /* initialize MPI; for brevity all status checks are omitted */
9     MPI_Init(&argc, &argv);
10
11     /* get the world communicator to build its group */
12     MPI_Comm_group(MPI_COMM_WORLD, &world_grp);
13
14     /* build another group by excluding a process */
15     MPI_Group_excl(world_grp, 1, 0, &workers_grp);
16
17     /* now build a communicator out of the processes in workers_grp */
18     MPI_Comm_create(MPI_COMM_WORLD, worker_grp, &workers);
19
20     /* etc... */
21
22     /* shut down MPI */
23     MPI_Finalize();
24
25     return 0;
26 }
```

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);
```

- ▶ the function

```
1 double MPI_Wtime();
```

returns the time in seconds from some arbitrary time in the past

- ▶ guaranteed not to change only for the duration of the process
- ▶ you can compute the elapsed time for any program segment by making calls at the beginning and the end and computing the difference
- ▶ no guarantees about synchronized clocks among different processes
- ▶ you can compute the clock resolution by using

```
1 double MPI_Wtick();
```

Other collective operations

- ▶ `MPI_Scan` computes partial reductions: the p^{th} process receives the result from processes 0 through $p - 1$

```
1 int MPI_Scan(  
2     void* send_buffer, void* recv_buffer,  
3     int count, MPI_Datatype datatype, MPI_Op operation,  
4     MPI_Comm communicator  
5 );
```

- ▶ `MPI_Reduce` collects the result at only the given process root

```
1 int MPI_Reduce(  
2     void* send_buffer, void* recv_buffer,  
3     int count, MPI_Datatype datatype, MPI_Op operation,  
4     int root, MPI_Comm communicator  
5 );
```

- ▶ synchronization is also a global operation:

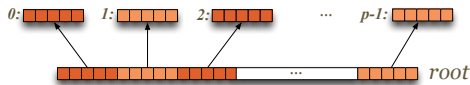
```
1 int MPI_Barrier(MPI_Comm communicator);
```

participating processes block at a barrier until they have all reached it

Scatter

- ▶ `MPI_Scatter` sends data from `root` to all processes

```
1 int MPI_Scatter(  
2     void* send_buffer, int send_count, MPI_Datatype send_datatype,  
3     void* recv_buffer, int recv_count, MPI_Datatype recv_datatype,  
4     int root, MPI_Comm communicator  
5 );
```

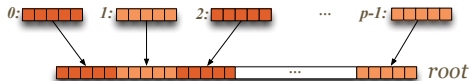


- ▶ it is as if the data in `send_buffer` were split in p segments, and the i^{th} process receives the i^{th} segment
- ▶ the `send_xxx` arguments are only meaningful for `root`; they are ignored for other processes
- ▶ the arguments `root` and `communicator` must be passed identical values by all processes

Gather

- ▶ the converse is `MPI_Gather` with `root` receiving data from all processes

```
1 int MPI_Gather(  
2     void* send_buffer, int send_count, MPI_Datatype send_datatype,  
3     void* recv_buffer, int recv_count, MPI_Datatype recv_datatype,  
4     int root, MPI_Comm communicator  
5 );
```



- ▶ it is as if p messages, one from each processes, were concatenated in rank order and placed at `recv_buffer`
- ▶ the `recv_xxx` arguments are only meaningful for `root`; they are ignored for other processes
- ▶ the arguments `root` and `communicator` must be passed identical values by all processes

Broadcasting operations

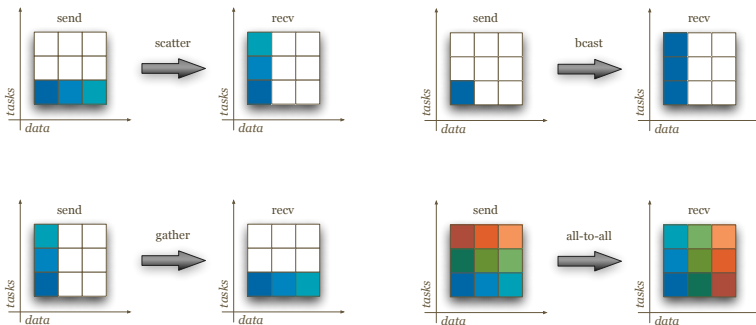
- ▶ `MPI_Alltoall` sends data from all processes to all processes in a global scatter/gather

```
1 int MPI_Alltoall(  
2     void* send_buffer, int send_count, MPI_Datatype send_datatype,  
3     void* recv_buffer, int recv_count, MPI_Datatype recv_datatype,  
4     MPI_Comm communicator  
5 );
```

- ▶ use `MPI_Bcast` to send the contents of a buffer from `root` to all processes in a communicator

```
1 int MPI_Bcast(  
2     void* buffer, int count, MPI_Datatype datatype,  
3     int root, MPI_Comm communicator  
4 );
```

Data movement patterns for the collective operations



Virtual topologies

