

MPI deadlocks and blocking

Dr. Christopher Marcotte — *Durham University*

MPI and deadlocks

In this exercise we will modify a simple code that sends a message between two processes.

In this directory you will find `ptp-deadlock.c`, and the relevant bit of code is this piece, which sends and receives a message between processes:

`ptp-deadlock.c`

```
26  if (rank == 0) {  
27      MPI_Send(sendbuf, nentries, MPI_INT, 1, 0, MPI_COMM_WORLD);  
28      MPI_Recv(recvbuf, nentries, MPI_INT, 1, 0, MPI_COMM_WORLD, MPI_STATUS_IGNORE);  
29  } else if (rank == 1) {  
30      MPI_Send(sendbuf, nentries, MPI_INT, 0, 0, MPI_COMM_WORLD);  
31      MPI_Recv(recvbuf, nentries, MPI_INT, 0, 0, MPI_COMM_WORLD, MPI_STATUS_IGNORE);  
32  }
```

Recall that `MPI_Send(...)` behaves like `MPI_Bsend(...)` when there is available buffer space, and behaves like `MPI_Ssend(...)` when there is not, blocking for the corresponding `MPI_Recv(...)`.

Exercise

Compile the program `ptp-deadlock.c` using

```
mpicc -o ptp-deadlock ptp-deadlock.c
```

and run the program, supplying `nentries` as an integer command line argument:

```
mpirun -n 2 ./ptp-deadlock <nentries>
```

Working from 1 in powers of 2, how large can you make `<nentries>` before the program deadlocks?

Hint

You will need `Control-c` (or `scancel <jobid>` on the batch system) to end the job once it deadlocks.

Solution

This depends on your machine. On my laptop it is 1024 entries.

Exercise

Modify the code to use `MPI_Ssend` instead of `MPI_Send`. How large can you make `nentries` now?

Solution

The calling syntax is the same, so just add that extra `s` in the function name:

ptp-deadlock-solution.c

```
34  if (rank == 0) {
35      MPI_Ssend(sendbuf, nentries, MPI_INT, 1, 0, MPI_COMM_WORLD);
36      MPI_Recv(recvbuf, nentries, MPI_INT, 1, 0, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
37  } else if (rank == 1) {
38      MPI_Ssend(sendbuf, nentries, MPI_INT, 0, 0, MPI_COMM_WORLD);
39      MPI_Recv(recvbuf, nentries, MPI_INT, 0, 0, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
40  }
```

And, of course, it deadlocks for any size message because both processes are waiting inside `MPI_Ssend(...)` for `MPI_Recv(...)`s that will never come.

Exercise

Modify the code to use `MPI_Sendrecv(...)` instead of `MPI_Ssend(...)` and `MPI_Recv(...)`.

Solution**ptp-deadlock-solution.c**

```
42  if (rank == 0) {
43      MPI_Sendrecv( sendbuf, nentries, MPI_INT, 1, 0, /* Send parameters */
44                  recvbuf, nentries, MPI_INT, 1, 0, /* Recv parameters */
45                  MPI_COMM_WORLD, MPI_STATUS_IGNORE);
46  } else if (rank == 1) {
47      MPI_Sendrecv( sendbuf, nentries, MPI_INT, 0, 0, /* Send parameters */
48                  recvbuf, nentries, MPI_INT, 0, 0, /* Recv parameters */
49                  MPI_COMM_WORLD, MPI_STATUS_IGNORE);
50  }
```

This code completes for any size message (up to some very unrealistically large sizes); `nentries=1073741824 (230)` took a few seconds but completed.

Challenge

Instead, use `MPI_Bsend` and set up the necessary buffers. Does this deadlock?

Global Gather

In this exercise we'll implement a rudimentary collective operation, where we gather some information from each process onto rank 0, and fill a buffer. This involves a lot of communication, and so gives an opportunity to test latency and throughput of some blocking and non-blocking MPI calls.

Exercise

Complete an MPI code in which rank-0 gathers a message from every process and places it in an array at a position corresponding to the rank of the sender. Implement two functions, `gather_blocking(...)` and `gather_nonblocking(...)`, so you can test their performance.

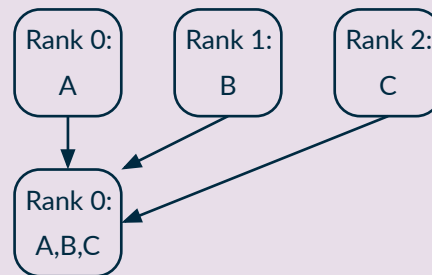


Figure 1:

See the `int main(...)` function in `gather-from-all.c`:

`gather-from-all.c`

```

39 int main(int argc, char **argv){
40     MPI_Init(&argc, &argv);
41
42     MPI_Comm comm = MPI_COMM_WORLD;
43     int rank, size;
44     MPI_Comm_rank(comm, &rank);
45     MPI_Comm_size(comm, &size);
46
47     int local;
48     local = rank * rank;
49
50     int *blocking = NULL;
51     int *nonblocking = NULL;
52     if (rank == 0) {
53         /* Allocate space for output arrays -- 1 int per process. */
54         blocking = malloc(size*sizeof(*blocking));
55         nonblocking = malloc(size*sizeof(*nonblocking));
56     }
57     double start, end;
58     start = MPI_Wtime();
59     for (int i = 0; i < 100; i++) {
60         gather_blocking(&local, blocking, comm);
61     }
62     end = MPI_Wtime();
63     if (rank == 0) {
64         printf("Blocking gather takes %.3g s\n", (end - start)/100);
65     }
66     start = MPI_Wtime();
67     for (int i = 0; i < 100; i++) {
68         gather_nonblocking(&local, nonblocking, comm);
69     }
70     end = MPI_Wtime();
71     if (rank == 0) {
72         printf("Non-blocking gather takes %.3g s\n", (end - start)/100);
73     }
74     free(blocking);
75     free(nonblocking);
76     MPI_Finalize();
77     return 0;
78 }

```

c

We allocate space for P processes outputs (a single `int` each) on rank 0, and then time how long it takes to call the two gathering functions a fixed number of times (in this case, 100, arbitrarily). Compare the timing of the two versions. Which performs better? Does it depend on the total number of messages, P ?

Hint

Process 0 uses a blocking `MPI_Recv(...)` for all receives in `gather_blocking(...)`, and uses a non-blocking `MPI_Irecv(...)` followed by `MPI_Waitall(...)` in `gather_nonblocking(...)`.

Solution

In both functions, we switch to different behavior for rank 0, first inserting the rank 0 value into index 0, and either creating the requests, immediate-receiving, and waiting – or using a blocking receive. Note that in both cases, we use `&(recvbuf[i])` for the receive buffer so the result is offset into the allocated array.

gather-from-all-solution.c

```
5 void gather_nonblocking(const int *send, int *recvbuf, MPI_Comm comm){c
6     int size, rank;
7     MPI_Comm_size(comm, &size);
8     MPI_Comm_rank(comm, &rank);
9
10    if (rank == 0) {
11        MPI_Request *requests;
12        requests = malloc((size-1) * sizeof(*requests));
13        recvbuf[0] = send[0];
14        for (int i = 1; i < size; i++) {
15            MPI_Irecv(&(recvbuf[i]), 1, MPI_INT, i, 0, comm, &(requests[i-1]));
16        }
17        MPI_Waitall(size-1, requests, MPI_STATUSES_IGNORE);
18        free(requests);
19    } else {
20        MPI_Send(send, 1, MPI_INT, 0, 0, comm);
21    }
22 }
23
24 void gather_blocking(const int *send, int *recvbuf, MPI_Comm comm){
25     int size, rank;
26     MPI_Comm_size(comm, &size);
27     MPI_Comm_rank(comm, &rank);
28
29     if (rank == 0) {
30         /* Put my value in the first spot */
31         recvbuf[0] = *send;
32         for (int i = 1; i < size; i++) {
33             /* Receive from each process in turn */
34             MPI_Recv(&(recvbuf[i]), 1, MPI_INT, i, 0, comm, MPI_STATUS_IGNORE);
35         }
36     } else {
37         /* Send to rank 0 */
38         MPI_Send(send, 1, MPI_INT, 0, 0, comm);
39     }
40 }
```

On my laptop, I get roughly a factor of $10\times$ between the versions:

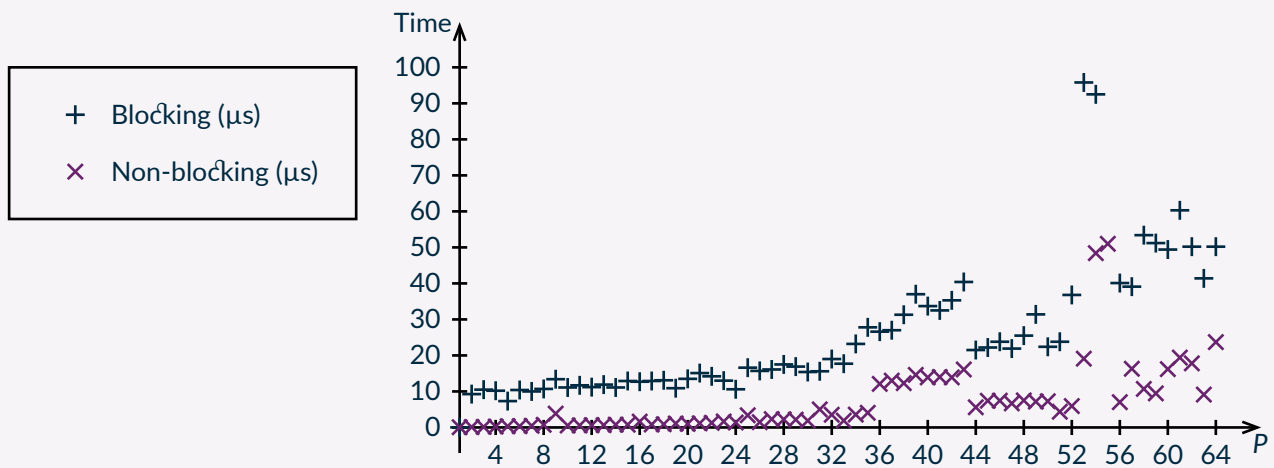


Figure 2: Scaling of ring allreduce as a function of P .

Challenge

Modify the code to compute the average time for these exchanges, over multiple instances – this will give more reliable results for such short times.

Aims

- Introduction to MPI deadlocks and how they impact program structure
- Understanding basics of blocking communications
- Understanding the primary failure mode of MPI