

COMP3221 Parallel Computation

David Head

University of Leeds

Lecture 9: Point-to-point communication

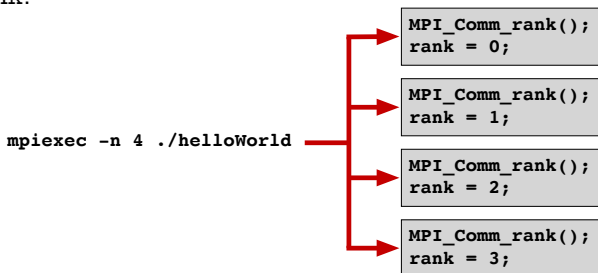
Previous lectures

Last lecture we started looking at **distributed memory systems**:

- Each **processing unit** can only see a fraction of memory.
- e.g. **clusters** of machines for HPC = High Performance Computing, where each node has its own memory.
- Standard API for low-level programming is MPI = Message Passing Interface.
- Processing units are **processes** rather than threads.
- Saw a 'Hello World' program for MPI.

mpiexec or mpirun

Launches **multiple executables** simultaneously, possibly on different machines/nodes, which are identical in every way **except** their rank:



All processes exist for the **duration of the program run**.

- Creation or destruction of **processes** is **expensive** (compared to *threads* in e.g. shared memory systems).

Today's lecture

Today we will start looking at using MPI to solve real problems.

- The simplest form of communication: **point-to-point**.
- Implemented in the MPI standard with `MPI_Send()` and `MPI_Recv()` (*plus variations*).
- **Vector addition**, the same problem we looked at for shared memory systems in Lecture 3.
- How exceeding the **buffer** size for some communication patterns can lead to **deadlock**.

Vector addition

Recall vector addition can be written mathematically as

$$\mathbf{c} = \mathbf{a} + \mathbf{b} \quad \text{or} \quad c_i = a_i + b_i, \quad i = 1 \dots N,$$

and in serial code as

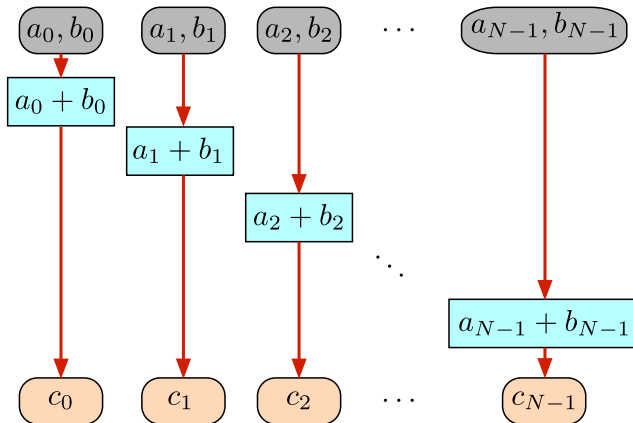
```
1 for( i=0; i<N; i++ )  
2   c[i] = a[i] + b[i];
```

where vectors **a**, **b** and **c** all have N elements¹.

This is a **data parallel** problem, also known as a **map** [*Lecture 3*].

¹By convention, indexing starts from 1 for mathematics notation but 0 in code.

Vector addition as a map¹



¹McCool et al., *Structured parallel programming* (Morgan-Kaufman, 2012).

Vector addition on a distributed memory system

Code on Minerva: `vectorAddition.c`

Suppose vectors **a** and **b** initially lie in the memory space of **one process only**, e.g. rank 0.

- Unlike in a shared memory system, the other processes **cannot see the data**.

Therefore must:

- ❶ **Distribute** vectors **a** and **b** across multiple processes.
- ❷ Perform the calculations in parallel **across all processes**, each working on a different **segment** of the arrays.
- ❸ **Gather** the segments together on a single process.

Point-to-point communication

The simplest way to send data from one process to another is to use `MPI_Send()` and `MPI_Recv()`:

- The process **sending** the data calls `MPI_Send()`.
- The process **receiving** the data calls `MPI_Recv()`.

Recall from last lecture that after initialising MPI, we determine the total number of processes `numProcs`, and the **rank** of 'this' process rank as follows:

```
1 int rank, numProcs;  
2  
3 MPI_Comm_size( MPI_COMM_WORLD, &numProcs );  
4 MPI_Comm_rank( MPI_COMM_WORLD, &rank );
```


MPI_Send()

For step 1, rank 0 **distributes** segments of the arrays a and b:

```
1 if( rank==0 )
2 {
3     for( p=1; p<numProcs; p++ )
4         MPI_Send(
5             &a[p*localSize],      // Pointer to the data
6             localSize,            // The size to be sent
7             MPI_FLOAT,            // The data type
8             p,                    // Destination process rank
9             0,                    // Tag; usually set to 0
10            MPI_COMM_WORLD        // Communicator
11        );
12 }
```

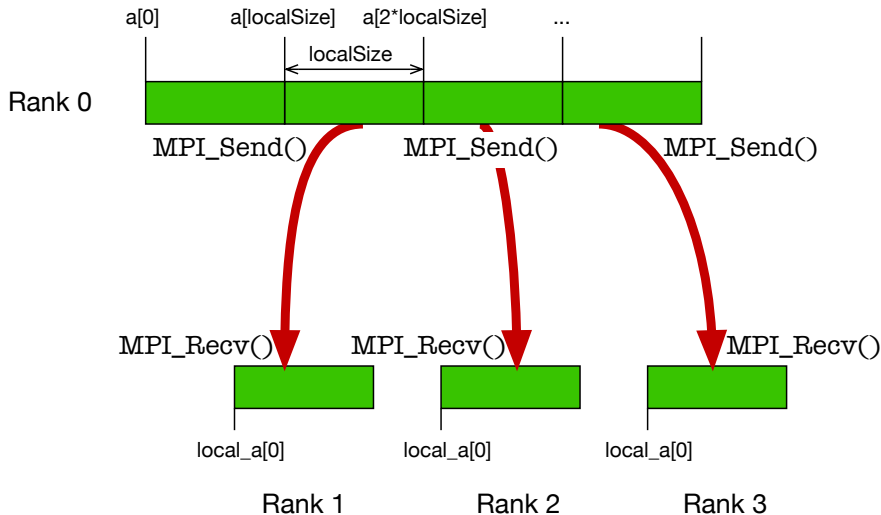
And similarly for b. Here, $\text{localSize} = N / \text{numProcs}$ is the **problem size per process**, *i.e.* the size of the **local** arrays / array segments.

MPI_Recv()

All ranks except 0 need to receive the data:

```
1 MPI_Status status;  
2 ...  
3 if( rank>0 )  
4 {  
5     MPI_Recv(  
6         local_a,           // Pointer to the data  
7         localSize,         // The size being sent  
8         MPI_FLOAT,         // The data type  
9         0,                 // Source process rank  
10        0,                 // Tag (can set to 0)  
11        MPI_COMM_WORLD,    // Communicator  
12        &status            // MPI_Status object  
13    );  
14 }
```

And similarly for local_b.



Completing the calculation

After rank 0 has distributed the full arrays `a` and `b` to the local arrays `local_a` and `local_b` on all other ranks:

- They all perform vector addition using their local arrays.
- The `local_c` arrays on each `rank>0` are sent to rank 0 using the same procedure as before.

Note that, in the code, **local** arrays are given separate names to the full arrays.

- e.g. `local_a` rather than `a`.

This is recommended (but not essential) to help keep track.

- The p-loop starts from 1, not zero. Sending 'to self' (e.g. from rank 0 to rank 0) is **undefined**.
 - Works in OpenMPI but *not* MPICH, so your code would not be portable.
- The **data type** is one of MPI_FLOAT, MPI_INT, MPI_DOUBLE, MPI_CHAR, ...
- &a[p*localSize] is a pointer to a **sub-array** that starts at element p*localSize of a.
- Most MPI calls return MPI_SUCCESS if successful, otherwise an error occurred.
- Can probe the status object to determine errors, rank of sending process *etc*.
- Can also replace &status with MPI_STATUS_IGNORE.

How is the communication performed?

The MPI standard does **not** specify **how** the communication is actually performed.

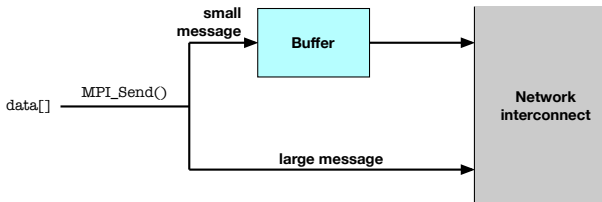
- If the nodes have IP addresses, could use standard internet protocols (e.g. sockets); the **Network** layer [*cf. COMP2221 Networks*].
- For HPC machines (where nodes do not have IP addresses), could use **Link** layer protocols or bespoke methods.

In this module we focus on **general** aspects of distributed system programming, not details of any MPI implementation.

- **Portable** code that should run on **any** implementation.

Common communication features

- Each data message has a **header** containing information such as the source and destination ranks¹.
- Message placed on a **buffer** ready to send.
- If it is too large, will send **direct** to the destination process.



¹Maximum header size is MPI_BSEND_OVERHEAD, defined in mpi.h.

Blocking communication

`MPI_Send()` and `MPI_Recv()` are examples of **blocking** routines.

Blocking routines do not return until all resources can be reused

This means e.g. values in the data array can be altered **without affecting the values sent**.

- Convenient from a programming perspective.

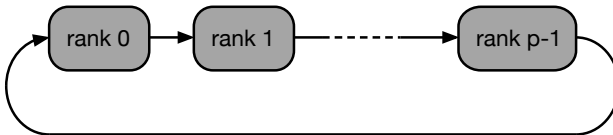
By contrast, **non-blocking** routines return 'immediately,' even though the data may still be being copied over.

- We will cover non-blocking communication in Lecture 12.

Cyclic communication

Code on Minerva: `cyclicSendAndReceive.c`

Consider a problem where the communication pattern is **cyclic**:



Encode this concisely using the ternary operator '`(a?b:c)`' to handle the wrap-around:

```
1 // Send data 'to the right'.
2 MPI_Send( sendData, N, MPI_INT,
3   ( rank==numProcs-1 ? 0 : rank+1 ), ... );
4
5 // Receive data 'from the left'.
6 MPI_Recv( recvData, N, MPI_INT,
7   ( rank==0 ? numProcs-1 : rank-1 ), ... );
```

Use of buffering

If the data is small enough to fit on the buffer:

- 1 Each process calls `MPI_Send()` to send data 'to its right.'
- 2 The data is **copied to the buffer** and `MPI_Send()` returns.
- 3 Each process calls `MPI_Recv()` and receives data from the process 'to its left.'

If the data is too large for the buffer, the application **hangs**:

- 1 `MPI_Send()` does **not** return until the destination process receives the data.
- 2 **All processes are in the same situation** - none of them reach their call to `MPI_Recv()`.
- 3 As no data is received, no process returns from `MPI_Send()`.

Deadlock

This is another example of **deadlock** that we first saw in Lecture 7:

Deadlock

Each process is waiting for a synchronisation event that never occurs.

In this case the ‘synchronisation event’ is the blocking send and receive that required the destination process to receive the data.

- Say more about the relationship between blocking and synchronisation in Lecture 12.

Resolving communication deadlocks

The buffer size is **not** specified by the MPI standard and varies between implementations.

- Even allowed to be zero size!
- Want to write code that works for **any** buffer size.

There are various ways to resolve this deadlock problem:

- ① Change the **program logic** *[here]*.
- ② Use **non-blocking communication** *[Lecture 12]*.
- ③ Allocate your own memory for a buffer and use **buffered send** `MPI_Bsend()`.

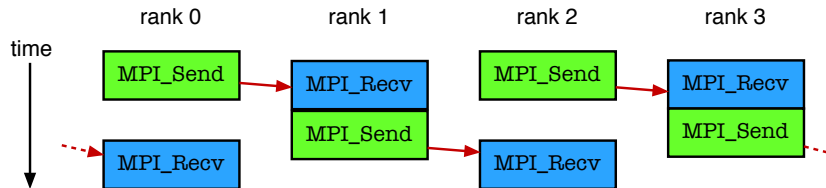
Staggering the send and receives

For this example, it is easiest to change the program logic to use **staggered** sends and receives¹:

```
1  if( rank%2 )
2  {
3      MPI_Send(sendData,N,MPI_INT,...);
4      MPI_Recv(recvData,N,MPI_INT,...);
5  }
6  else
7  {
8      MPI_Recv(recvData,N,MPI_INT,...);
9      MPI_Send(sendData,N,MPI_INT,...);
10 }
```

¹Recall $i\%2==0$ if i is even, and 1 if i is odd.

Processes with even-numbered ranks **receive** first **then** send, breaking the deadlock:



Note the arguments with each `MPI_Send()` and `MPI_Recv()`, including the source and destination ranks, **have not been altered**.

Summary and next lecture

Today we have looked at **point-to-point communication** in a distributed memory system:

- How to implement a **data parallel problem** (or a **map**) using `MPI_Send()` and `MPI_Recv()`.
- These routines are **blocking**, a similar concept to **synchronous communication**.
- Exceeding the buffer can lead to **deadlock**.

Next time we will look at some **performance considerations**, and how they can be improved by using **collective communication**.