

# Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat University of Leeds edu\_assist\_pr

Lecture 9: Point-to-point com

## Previous lectures

# Assignment Project Exam Help

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



- <https://eduassistpro.github.io>

- Standard API for low-level programming

Passing Interface.

- Processing units are **processes** rather

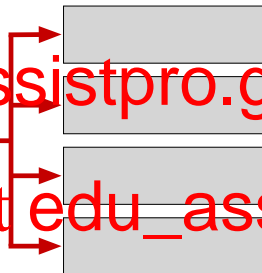
- Saw a 'Hello World' program for MPI.

Add WeChat edu\_assist\_pro

## mpiexec or mpirun

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

```
mpiexec -n 4 ./helloWorld
```



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

# Assignment Project Exam Help

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



<https://eduassistpro.github.io>



**Vector addition**, the same problem we had in memory systems in Lecture 3.



How exceeding the **buffer** size for some patterns can lead to **deadlock**.

## Vector addition

# Assignment Project Exam Help

Recall, vector addition can be written mathematically as

and in C

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

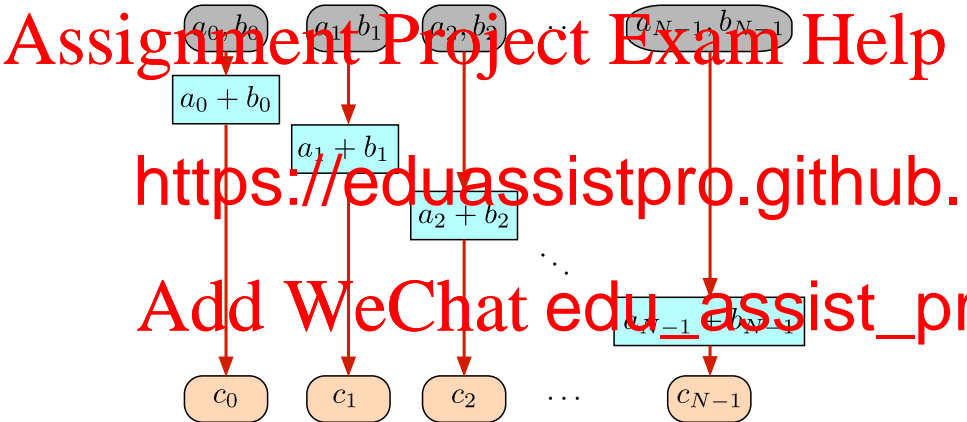
where vectors **a**, **b** and **c** all have  $N$  el

This is a **data parallel** problem, also kno

---

<sup>1</sup>By convention, indexing starts from 1 for mathematics notation but 0 in code.

## Vector addition as a map<sup>1</sup>



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

## Vector addition on a distributed memory system

Code on Minerva: `vectorAddition.c`

# Assignment Project Exam Help

Suppose vectors **a** and **b** initially lie in the memory space of **one** **proc**

- <https://eduassistpro.github.io>

Therefore must:

- 1 **Distribute** vectors **a** and **b** across **ch**
- 2 Perform the calculations in parallel working on a different **segment** of the arrays.
- 3 **Gather** the segments together on a single process.

## Point-to-point communication

# Assignment Project Exam Help

The simplest way to send data from one process to another is to use MPI\_Send() and MPI\_Recv():



<https://eduassistpro.github.io>

Recall from last lecture that after initialising MPI the total number of processes numProcs and the 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     fo
4
5
6     // The s
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:

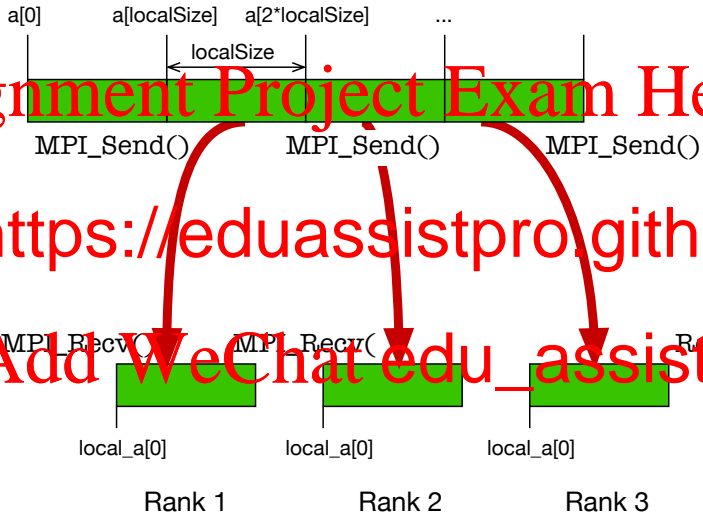
```
1 MPI_Status status;  
2 ...  
3 if( rank  
4 {  
5     MPI_Recv(  
6         localSize,           // Pointer to  
7         MPI_FLOAT,           // The size being sent  
8         0,                   // The data type  
9         0,                   // Source process rank  
10        0,       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

And similarly for local\_b.

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pro



## Completing the calculation

# Assignment Project Exam Help

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

- <https://eduassistpro.github.io>

Note that in the code, `local` arrays are not the full arrays.

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

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

- 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?

**Assignment Project Exam Help**  
The MPI standard does not specify how the communication is actually performed.



<https://eduassistpro.github.io>

- For HPC machines (where nodes do not have local disks) could use **Link** layer protocols or bsp

**Add WeChat edu\_assist\_pro**

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

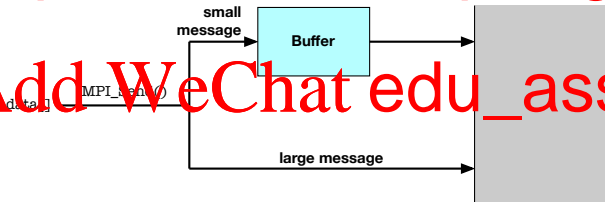
# Assignment Project Exam Help

- Each data message has a **header** containing information such as the source and destination ranks<sup>1</sup>.

•

- <https://eduassistpro.github.io>

Add WeChat edu\_assist\_pro



<sup>1</sup>Maximum header size is `MPI_BSEND_OVERHEAD`, defined in `mpi.h`.

## Blocking communication

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

Blo

<https://eduassistpro.github.io>

This **blocking** communication affects 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.

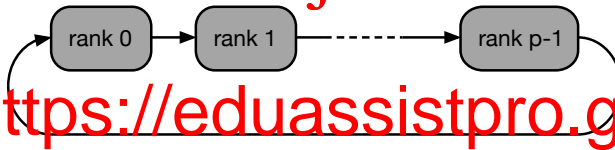
- We will cover non-blocking communication in Lecture 12.



## Cyclic communication

Code on Minerva: `cyclicSendAndReceive.c`

Consider a problem where the communication pattern is cyclic



Encode this concisely using the ternary operator '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

# Assignment Project Exam Help

If the data is small enough to fit on the buffer

- 1 Each process calls `MPI_Send()` to send data 'to its right.'
- 2 \_\_\_\_\_ turns.

3 <https://eduassistpro.github.io>

If the data is too large for the buffer, the application

- 1 `MPI_Send()` does **not** return until the receiver 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

Assignment Project Exam Help

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

Dea

E

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr

In this case the 'synchronisation event' is the block receive that required the destination process to re

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

## Resolving communication deadlocks

# Assignment Project Exam Help

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



- <https://eduassistpro.github.io>

There are various ways to resolve this deadlock problem



Change the program log c [here]



Use **non-blocking communication**



Allocate your own memory for a buffer and use **buffered send** `MPI_Bsend()`.

## Staggering the send and receives

Assignment Project Exam Help

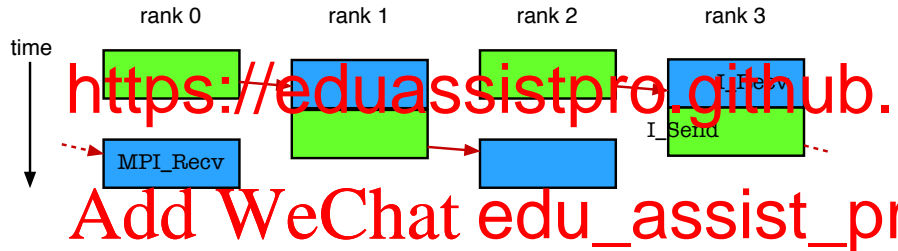
For this example, it is easiest to change the program logic to use **staggered** sends and receives<sup>1</sup>:

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

---

<sup>1</sup>Recall  $i\%2==0$  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

# Assignment Project Exam Help

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

- <https://eduassistpro.github.io> using
- These routines are **blocking**, a similar to **synchronous communication**.
- Exceeding the buffer can lead to

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