

# HPC

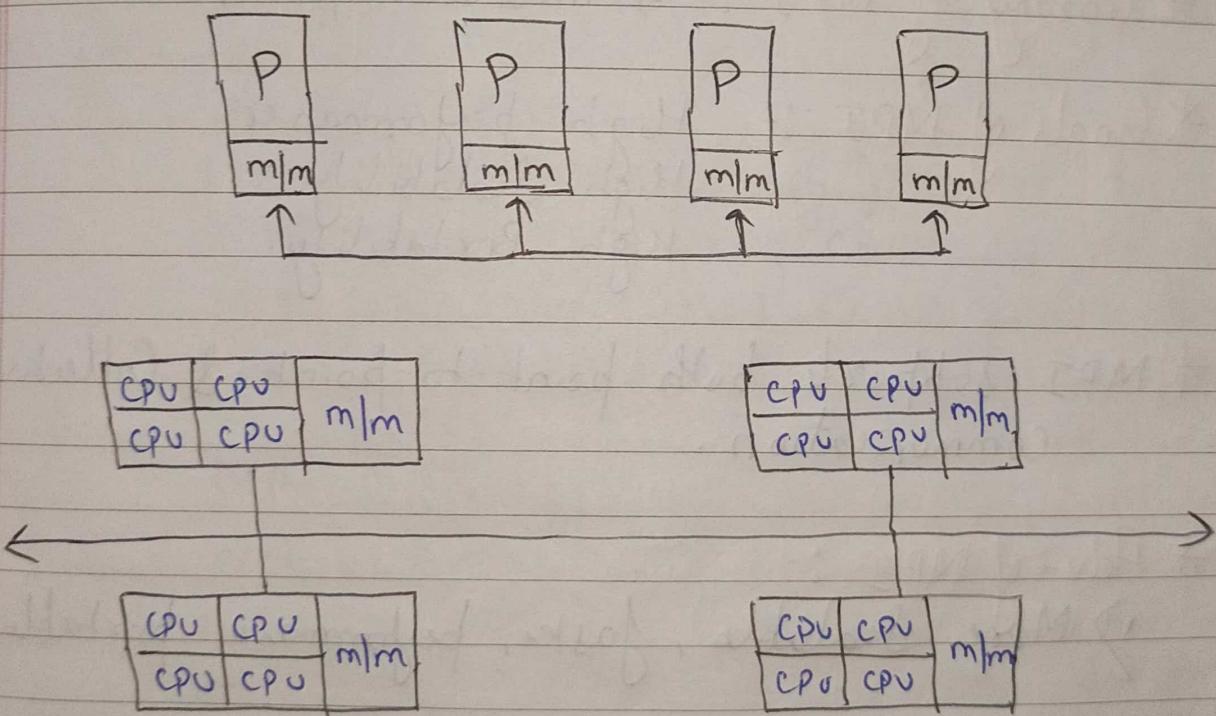
## Unit - 5

MPI :- (Message Passing Interface)

→ MPI concept came in 1990s.

→ It is used for communication between two or more processors.

\* MPI is a standardize means of exchanging message between multiple computers running parallelly across distributed memory.



**Note :-** MPI is needed when there is connection between distributed memory system. We don't need for shared memory system.

(i.e required when there is no way for one processor to directly access the address of another).

However, it can be regarded as programming model

and used on shared-memory or hybrid system.

**Need of MPI?** → It provides mechanism to create processes to execute process on different processors.  
Also provide mechanism to send and receive message.

\* In earlier days, we were using MPMD model which had certain limitations but nowadays we use SPMD which have no restrictions.  
MPMD uses PVM (Parallel virtual Machine) as library.

\* Library of MPI :- `#include <mpi.h>`

\* Goals of MPI :- High performance  
High Scalability  
High Portability.

\* MPI supports both point-to-point & collective communication.

\* Adv. of MPI :-

→ More Scalable, faster, performance, portable

History :-

MPI1 (Version 1) → 1994

- most application runs on this
- used to create parallel SPMD code
- available on almost all language C, C++, Fortran.

MP1 2 (version 2) → 1996

→ MPI + One-sided communication

dynamic process control

Parallel I/O (MPI-I/O).

↑ can change according to  
nature of model.

### MPI with a Code :-

```
#include <stdio.h>
#include <mpi.h>
int main() {
    int rank, size;
    MPI_Init(&argc, &argv);
    MPI_Comm_size(MPI_COMM_WORLD, &size);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    printf("Hello World", I am %d of %d", rank, size);
    MPI_Finalize();
    return 0;
}
```

→ MPI\_Init → Initialize parallel environment

→ MPI\_COMM\_WORLD → Default commutator of MPI

→ MPI\_COMM\_SIZE → No. of processes

→ MPI\_COMM\_RANK → Rank of calling process.

↓  
(unique identifier which identify each process).

→ MPI\_Finalize → End parallel environment

Communication → point-to-point  
Communication → Collective

Point-to-Point → two processes in the same computer are going to communicate each other.

Collective → All processes communicate each other.

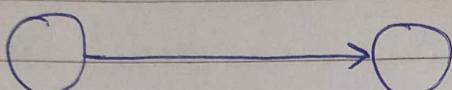


Fig. point-to-point

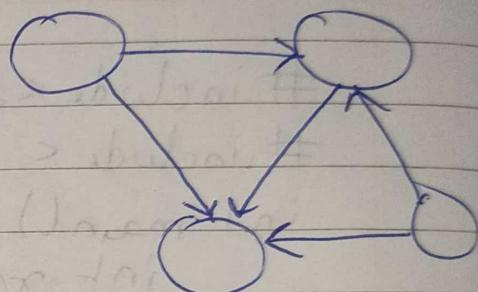
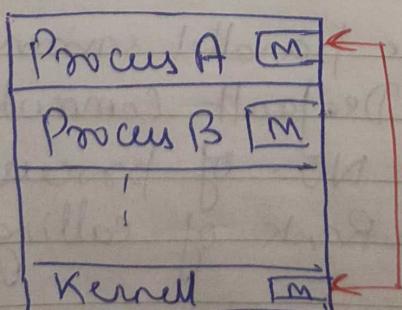


Fig. Collective

\* We need 2 operations to perform MP.  
In distributed environment, computer processes may reside in different computers connected by a network. This kind of scenario is called as MP System (MPS). There is no need of shared m/m.



\* MPI provides 2 type of operations :-  
① MPI\_send ()    ② MPI\_receive ()

\* Sender \* Sending and receiving are two fundamental concept of MPI. Almost every single function in MPI can be implemented with basic send & receive calls.

\* Usage - send by a process can be either fixed or variable.

fixed → The system level implementation is straight forward as compared to variable. But it creates problem and makes the task of programmer more difficult.

Variable → require more complex system but the program task becomes easy or simpler.

### Prototype of MPI\_Send()

MPI\_Send( )

void \*data,

int count,

MPI\_Datatype datatype

int destination,

int tag,

MPI\_Comm communicator)

### Prototype of MPI\_Recv()

MPI\_Recv( )

void \*data,

int count

MPI\_Datatype datatype,

int source,

int tag,

MPI\_Comm communicator),

MPI\_Status \* status)



For communication between A → B, a communication link must exist between two processes. This link can be implemented in variety of ways:-

① Direct or Indirect

② Synchronous & Asynchronous

③ Automatic or explicit buffering.

Direct → Each process that want to communicate, must explicitly name the communication.

Send(P, message) receive(q, message)

Indirect → not explicitly called.  
→ use mailbox

Send a message to mailbox A  
receive a message from mailbox A.

Synchronous → Communicate after scheduling.  
Ex → video call  
→ require immediate action.

Asynchronous → there is no need of scheduling.  
→ No immediate action is required.

Automatic → automatically allocates space  
→ provide Queue with indefinite block  
→ no need to wait

Explicit buffering → need to wait until the block becomes available  
→ Sender cannot send any message till then.

Communication → Blocking  
Non-Blocking

Blocking :- Simultaneous send & receive function → MPI\_SEND(), MPI\_RECV()

Non-Blocking :- Immediate execution.  
WGL queue to store msg.  
function → MPI\_ISEND(), MPI\_IRECV().

Collective Communication :-

→ One to many

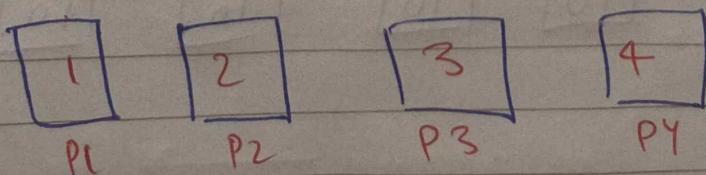
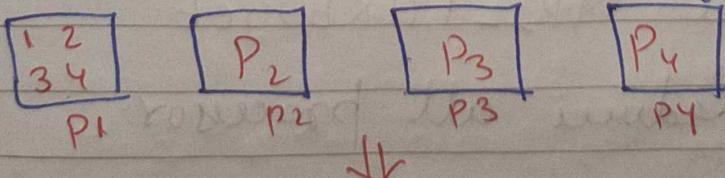
3 types :-

1) Barrier Synchronization :- MPI\_BARRIER().  
↳ Block until all process have reached a synchronization point.

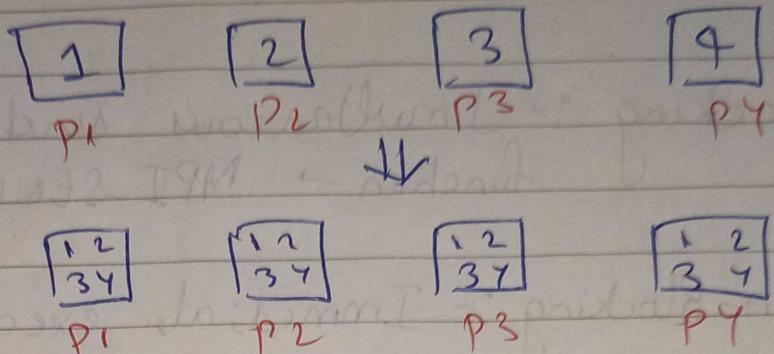
2) Data Movement (Global Communication)

↳ Broadcast MPI\_BCAST → (to all)

↳ Scatter MPI\_SCATTER → Send data from 1 task to all other task in a group.



→ Gather : MPI\_GATHER → (opposite to Scatter).  
 → All-Gather : MPI\_ALLGATHER

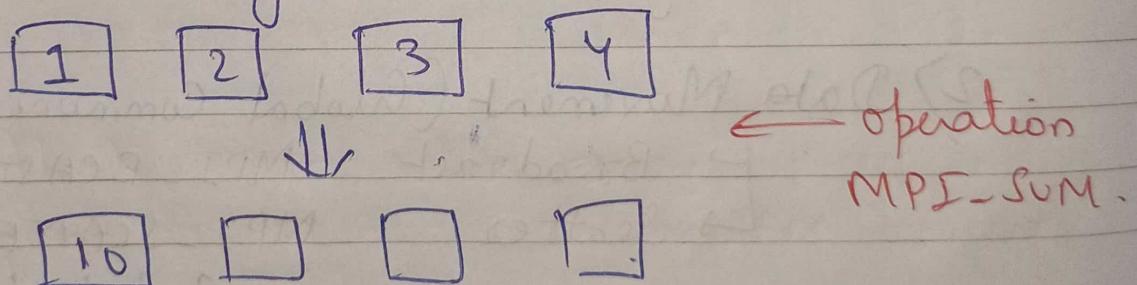


### Collective operation (Global Reduction) :-

→ One process from Communicator collects data from each process & performs an operation on that data to compute a result.

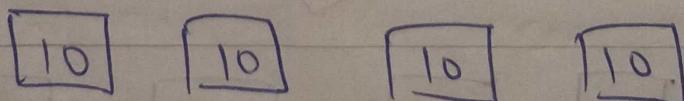
#### 1) \* Reduce → MPI\_REDUCE

↳ It reduces the value from all process to a single value



#### \* All Reduce MPI\_ALLREDUCE

→ Reduces all processor.



## \* Reduce Scatter MPI\_REDUCE\_SCATTER

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

↓↓

$$\begin{bmatrix} 4 \end{bmatrix}$$

$$\begin{bmatrix} 8 \end{bmatrix}$$

$$\begin{bmatrix} 12 \end{bmatrix}$$

$$\begin{bmatrix} 16 \end{bmatrix}$$

## MPI Scan :-

commutative result

$$\begin{bmatrix} 1 \end{bmatrix}$$

$$\begin{bmatrix} 2 \end{bmatrix}$$

$$\begin{bmatrix} 3 \end{bmatrix}$$

$$\begin{bmatrix} 4 \end{bmatrix}$$

↓

$$\begin{bmatrix} 1 \end{bmatrix}$$

$$\begin{bmatrix} 3 \end{bmatrix}$$

$$\begin{bmatrix} 6 \end{bmatrix}$$

$$\begin{bmatrix} 10 \end{bmatrix}$$

# Virtual topologies

A virtual topology represent a way that MPI processes communicate.

MPI provides 2 types of topologies:-

- (i) Cartesian grid based topology
- (ii) Graphs Based topology.

(i) Cartesian grid based topology :-

It is a regular grid like structure that can be visualized as 1D, 2D and 3D array of process. To create this we use `mpi_cart_create()` function.

Grid structure are defined by a number of dimensions and no. of processes in each coordinate direction. It is mapped using a row-major numbering system. It is machine dependent.

0 (0,0)	1 (0,1)	2 (0,2)	3 (0,3)
(1,0) 4	(1,1) 5	(1,2) 6	(1,3) 7
(2,0) 8	(2,1) 9	(2,2) 10	(2,3) 11

function :- 1) `MPI_Cart_create (MPI_Comm comm-old, int ndims, int *dims, int *periods, int reorder, MPI_Comm *comm-cart)`.

2) `int MPI_Dims_create (int nodes, int ndims, int *dims)`

2) Graph based topology :-

- It can only be used in Intra Communicators
  - No. of graph nodes must not be greater than No. of processes.
  - It is a more general structure that can represent arbitrary connections between processes.
  - It uses Nodes and edges to represent the process and connection link respectively.
- MPI\_Graph\_Create() function is used to create this.

functions :- 1) MPI\_Graph\_Create (MPI\_Comm comm\_old, int nnodes, int \*index, int \*edges, int reorder, MPI\_Comm \*comm\_graph).

2) MPI\_Topo\_test (MPI\_Comm comm, int \*top\_type)  
↳ determine type of topology associated.

3) MPI\_Cartdim\_get (MPI\_Comm comm, int \*ndims)  
↳ retrieve Cartesian topology information.

4) MPI\_Cart\_Rank (MPI\_Comm comm, int \*coords, int \*rank)  
↳ determine process rank in communicator given Cartesian location.

5) MPI\_Graph\_neighbour\_Count()

↳ return count of neighbours

6) MPI\_Graph\_neighbour()

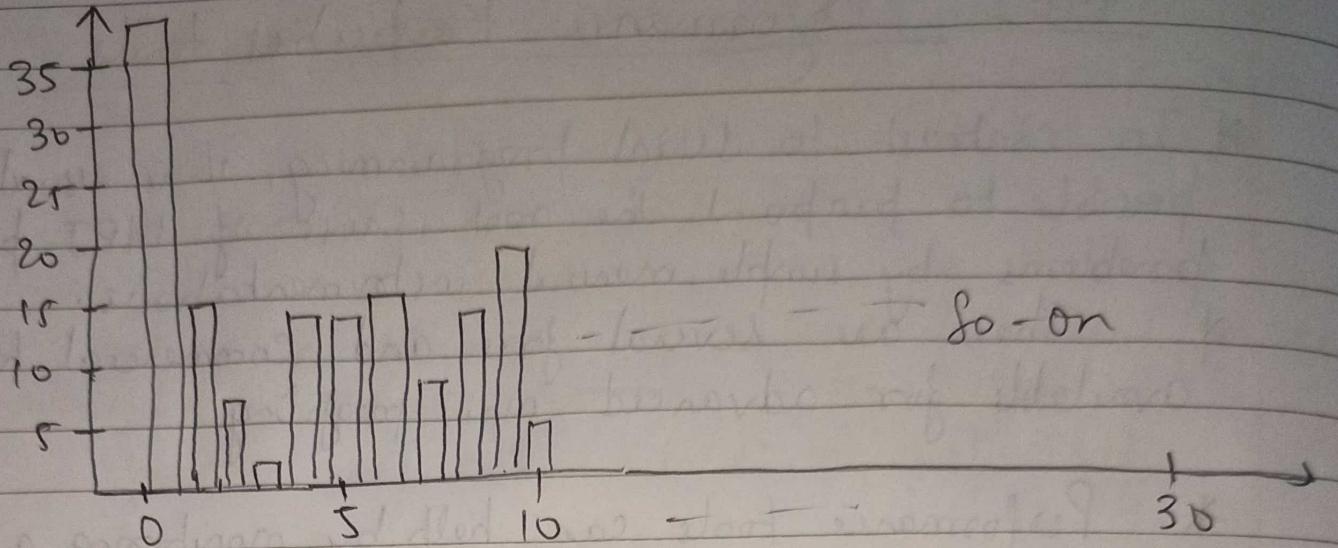
↳ return neighbours

7) MPI\_Cart\_Displace()

↳ for displacement

## MPI Performance Profiling tools

- \* In contrast to serial programming, it is usually not possible to pinpoint the root cause of MPI performance problems by simple manual instrumentations.
- \* There are several free and commercial tools available for advanced mpi profiling.
- \* Performance tools can help by monitoring a program execution and reducing performance data that can be analysed to locate and understand the area of poor performance.
- \* IPM is a simple and low-overhead tool that is able to retrieve this information.
- \* IPM is a portable profile infrastructure for parallel codes. It provides a low overhead performance. It is scalable and easy to use.
- \* Like most MPI profilers, IPM uses the MPI profiling interface, which is part of standard.
- \* Each MPI function is wrapped around the actual function whose name starts with 'PMPI-'.
- \* In case of IPM, it is sufficient to preload a dynamic library and run the application. Information about dynamic data volumes (per process and per process pair), time spent in MPI calls, load imbalance etc is then accumulated over the application's runtime and can be viewed in graphical form.



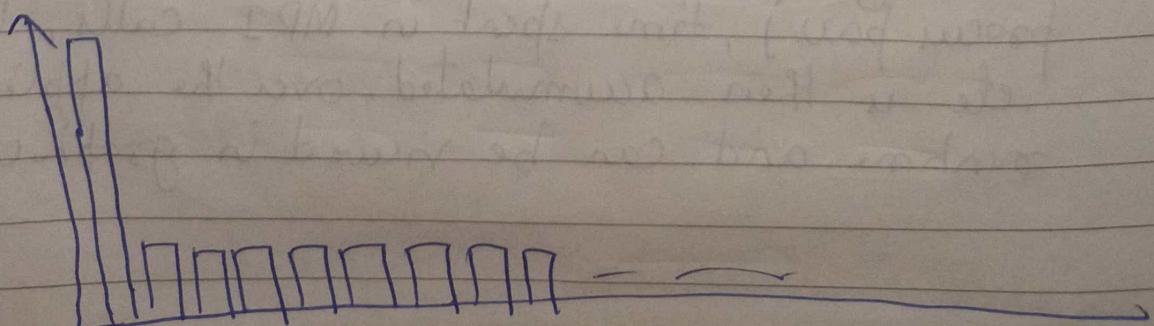
→ IPM "Communication balance" of a master work style parallel application.

Each bar represent MPI rank and shows how much ~~time~~ the process spend in different MPI functions.

In this, the runtime was 38 seconds.  
Rank 0 (the master) distributes the work among the workers, so it spend most of its runtime in MPI\_recv() waiting for result.

The workers are obviously quite load imbalanced and between 5 and 50% of their time is wasted waiting at barrier.

A small change in parameter (reducing the size of work packages) was able to correct this problem. Overall 25% reduction in runtime is seen.



## Problems in MPI :-

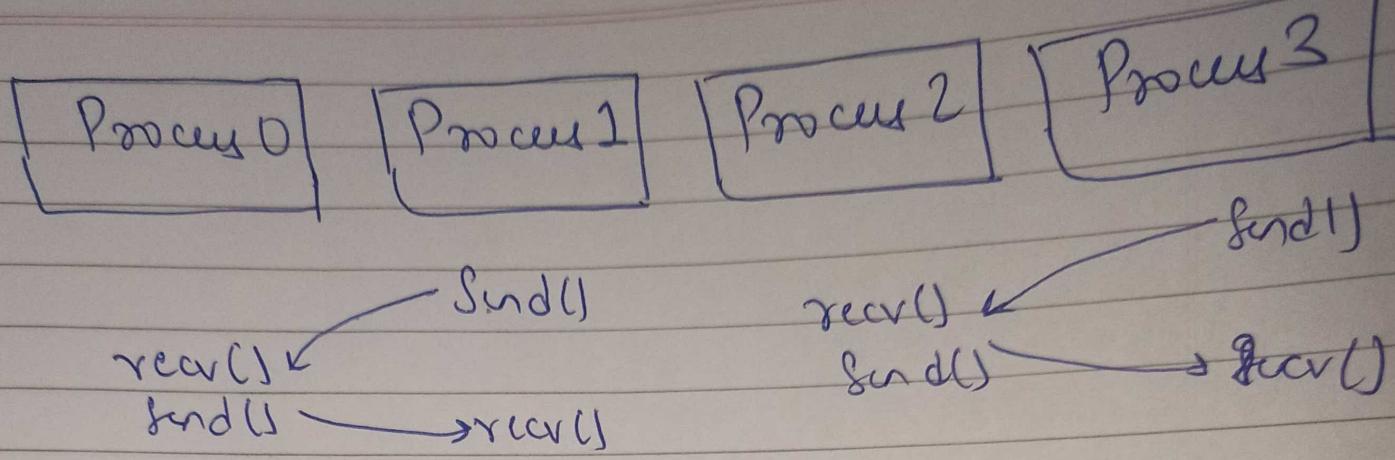
- order of message is not assured.
- MPI does not preserve the temporal order of messages unless they are sent between same sender and receiver (or with same tag).
- Rank 0 does not call MPI\_recv() before returning from its own execution.
- A common source of error in program that uses point-to-point communication is called deadlock.
- Consider a case where 2 process exchange data, each process perform both send and receive. The idea is that Process 0 will send <sup>data</sup> to P1 who will receive it from P0 and at same time P1 sends to P0.
- Deadlock occur when every process is forced to wait on an action of another process to continue executing.

How to overcome it?

- 1) One possible solution of fixing a deadlock problem or to avoid it is simply reversing the order of one of the receive or send pairs.

Send/receive → Receive/send

2)



- odd processes are send/recv() pair
- even processes are recv()/send() pair

## HPC UNIT – 5

### Non-blocking vs asynchronous communication in MPI

In the context of MPI (Message Passing Interface), the terms "non-blocking" and "asynchronous" are often used interchangeably, but they have slightly different meanings. Let's explore the difference between the two:

**1. Non-blocking communication:** Non-blocking communication refers to the ability to initiate a communication operation and continue executing other tasks without waiting for the operation to complete. In MPI, non-blocking communication is achieved using functions like `MPI\_Isend` and `MPI\_Irecv`. When you use these non-blocking communication functions, the MPI library returns immediately, allowing your program to perform other computations while the communication operation progresses in the background. However, you still need to ensure that you don't access the data involved in the communication until the operation is completed, to avoid race conditions.

**2. Asynchronous communication:** Asynchronous communication, in the context of MPI, implies that multiple communication operations can be in progress simultaneously, and they can progress independently of each other. In other words, you can initiate multiple communication operations and not necessarily wait for each of them to complete before initiating others. This level of asynchrony allows for potential overlap of communication and computation, leading to increased performance. Non-blocking communication is a key technique used to achieve asynchronous communication in MPI.

In summary, non-blocking communication is a mechanism that allows you to initiate a communication operation and continue with other tasks, while asynchronous communication is a broader concept that encompasses the ability to have multiple communication operations progressing simultaneously. Non-blocking communication is often used to enable asynchronous communication in MPI, but it's important to note that not all non-blocking communication is necessarily asynchronous if only one communication operation is being used.

## **Explain 1. serialization 2. Synchronization 3. Contention 4. Implicit serialization in MPI.**

1. **Serialization:** Serialization in MPI refers to the process of converting data from a structured format in memory to a linear format suitable for transmission over the network or storage. When sending or receiving data using MPI communication routines, the data needs to be serialized before transmission and deserialized upon reception. Serialization ensures that the data can be transmitted as a sequence of bytes, independent of the memory layout or data structures used by the sender and receiver processes. MPI provides functions like `MPI\_Pack` and `MPI\_Unpack` to facilitate serialization and deserialization of complex data types.
2. **Synchronization:** Synchronization in MPI involves coordinating the execution of multiple processes to ensure that they reach specific points in their execution in a coordinated manner. Synchronization is essential to enforce dependencies, order of operations, and consistency in parallel programs. In MPI, synchronization can be achieved using various mechanisms, such as collective communication operations like `MPI\_Barrier`, where processes wait until all processes in a communicator reach the barrier before continuing execution. Synchronization is crucial in parallel algorithms to avoid race conditions and maintain correctness.
3. **Contention:** Contention refers to a situation in which multiple processes or threads are competing for the same resource or piece of data, leading to potential conflicts and decreased performance. In the context of MPI, contention can arise when multiple processes simultaneously access shared resources, such as memory locations or I/O devices, leading to contention for access to those resources. Contention can result in increased communication delays, decreased scalability, and even deadlocks. Proper design and synchronization techniques are necessary to minimize contention and ensure efficient execution in parallel programs.
4. **Implicit serialization:** Implicit serialization in MPI refers to the automatic handling of data serialization and deserialization by MPI library functions. When using point-to-point communication routines like `MPI\_Send` and `MPI\_Recv`, MPI automatically serializes the data being sent and deserializes it upon reception, abstracting away the serialization process from the programmer. This implicit serialization simplifies the programming model and allows the developer to focus on the application logic rather than low-level data serialization details. However, it's important to note that implicit serialization might not be as efficient as explicit serialization using `MPI\_Pack` and `MPI\_Unpack` for certain complex data types, where explicit control over the serialization process may be required for performance optimization or compatibility with external systems.