

Parallel Computing - MPI

# Message Passing Interface



Shweta Das  
[shwetad@cdac.in](mailto:shwetad@cdac.in)  
HPC - Tech, CDAC Pune

Centre for Development of Advanced Computing



# MPI - Message Passing Interface

MPI is built on 'Routines'

The basic MPI Routines :-

- ☐ MPI\_Init () ;
- ☐ MPI\_Comm\_rank () ;
- ☐ MPI\_Comm\_size () ;
- ☐ MPI\_Send () ;
- ☐ MPI\_Recv () ;
- ☐ MPI\_Finalize () ;
- ☐ -----

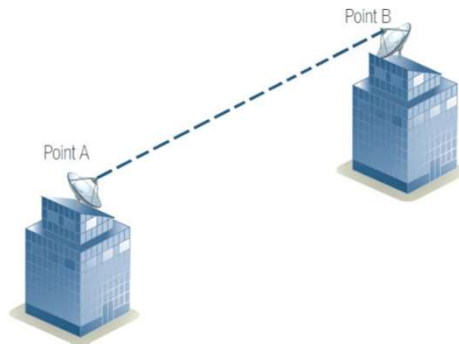
# MPI - Communication

# MPI - Communication

**Point to Point Comm<sup>n</sup>**

# MPI - Communication

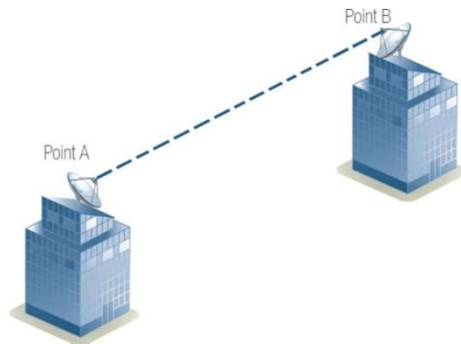
**Point to Point Comm<sup>n</sup>**



# MPI - Communication

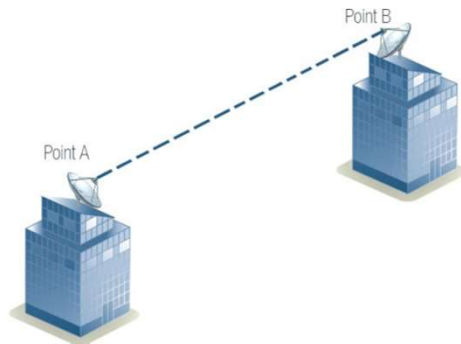
**Point to Point Comm<sup>n</sup>**

**Collective Comm<sup>n</sup>**

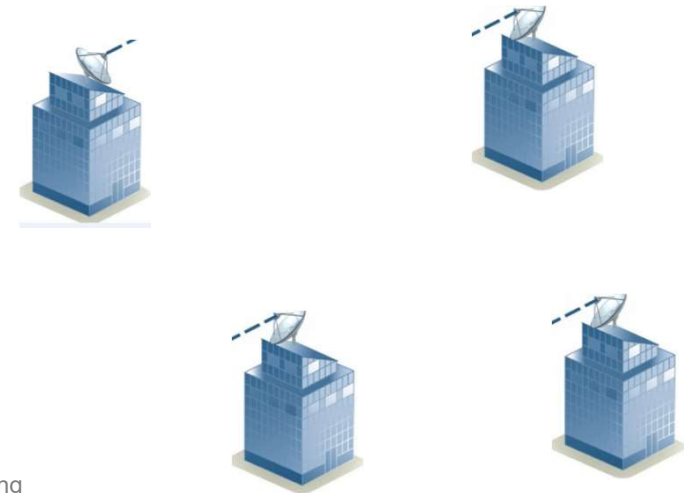


# MPI - Communication

**Point to Point Comm<sup>n</sup>**

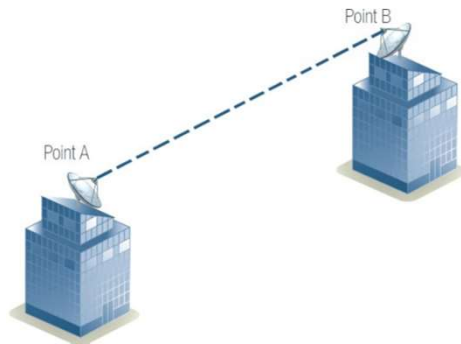


**Collective Comm<sup>n</sup>**

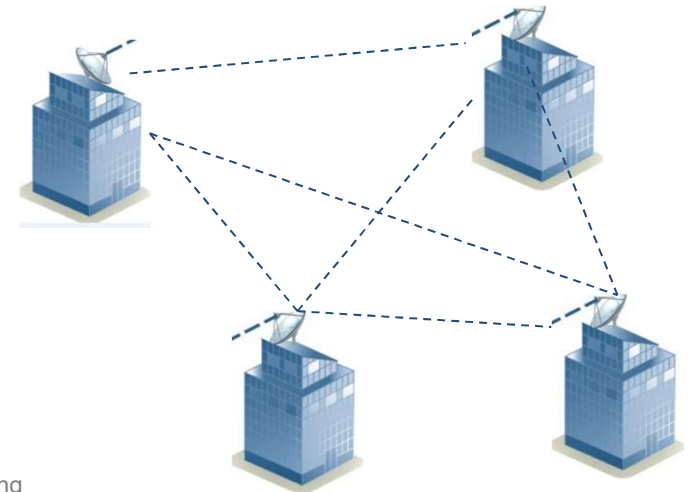


# MPI - Communication

**Point to Point Comm<sup>n</sup>**



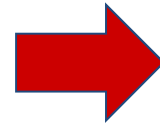
**Collective Comm<sup>n</sup>**



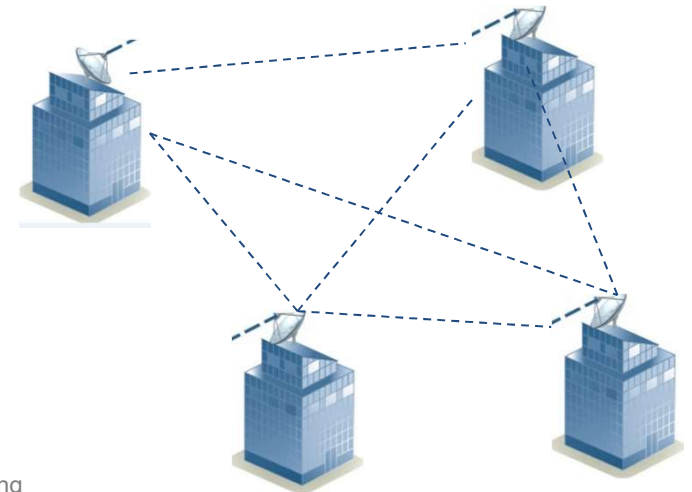
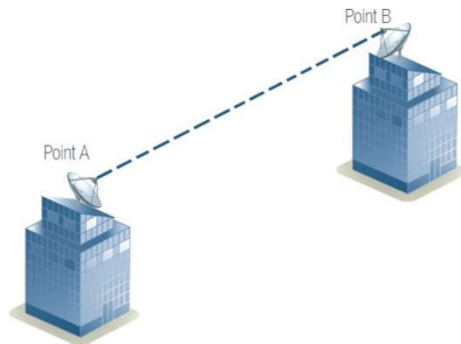


# MPI - Communication

Point to Point Comm<sup>n</sup>



Collective Comm<sup>n</sup>



# MPI - Collective Communication

- **Collective communication must involve all processes in the scope of a communicator.**
- **Involve coordinated communication within a group of processes identified by an MPI communicator.**

# Types of Collective Operations

- **Synchronization** - Processes wait until all members of the group have reached the synchronization point.
- **Data Movement** - broadcast, scatter/gather, all to all
- **Collective Computation (reductions)** - one member of the group collects data from the other members and performs an operation (min,max, add, multiply, etc.) on that data.

# Basic Collective Communication Routines

- **MPI\_Bcast( ) - Broadcast (one to all)**
- **MPI\_Scatter( ) - Scatter (one to all)**
- **MPI\_Gather( ) - Gather (all to one)**
- **MPI\_Reduce( ) - Reduce (all to one)**
- **MPI\_Allgather( ) - (all to all)**
- **MPI\_Allreduce( ) - (all to all)**

# MPI - Broadcast

## Syntax :

➔ `MPI_Bcast ( void* data , Int count , MPI_Datatype datatype , Int source_process , MPI_Comm comm ) ;`

- One process sends the same data to all processes in a communicator.

P0



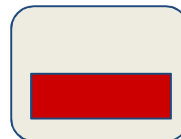
# MPI - Broadcast

## Syntax :

➔ `MPI_Bcast ( void* data , Int count , MPI_Datatype datatype , Int source_process , MPI_Comm comm ) ;`

- One process sends the same data to all processes in a communicator.

P0

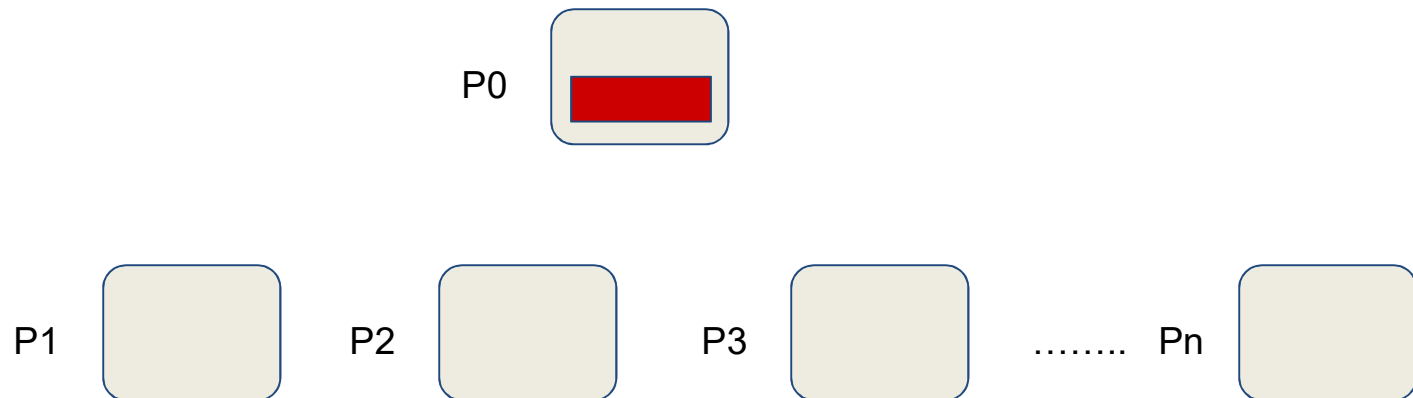


# MPI - Broadcast

## Syntax :

➔ `MPI_Bcast ( void* data , Int count , MPI_Datatype datatype , Int source_process , MPI_Comm comm );`

- One process sends the same data to all processes in a communicator.

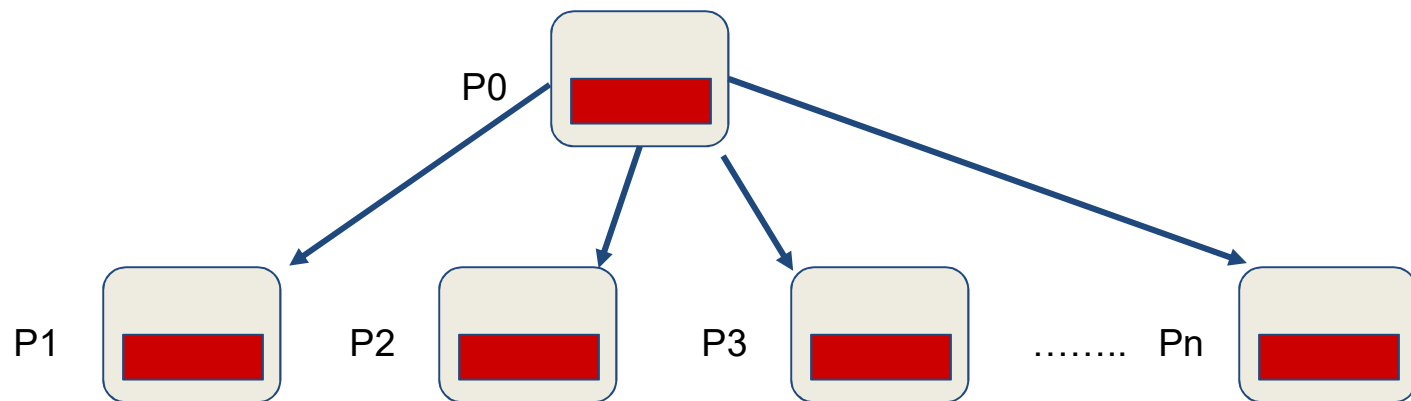


# MPI - Broadcast

## Syntax :

➔ `MPI_Bcast ( void* data , Int count , MPI_Datatype datatype , Int source_process , MPI_Comm comm ) ;`

- One process sends the same data to all processes in a communicator.





## MPI - Broadcast : Example

```
void Get input(int my rank ,Int comm_sz , double a_p , double b_p , int* n_p )
{
    if (my rank == 0)
    {
        printf("Enter a, b, and n \n");
        scanf("%lf %lf %d", a_p, b_p, n_p);
    }

    MPI Bcast(a_p, 1, MPI_DOUBLE, 0, MPI_COMM_WORLD);

    MPI Bcast(b_p, 1, MPI_DOUBLE, 0, MPI_COMM_WORLD);

    MPI Bcast(n_p, 1, MPI_INT, 0, MPI_COMM_WORLD);
}
```

# MPI - Reduce

## Syntax :

➔ `MPI_Reduce (void* input_data , void* output_data , Int count ,  
MPI_Datatype datatype , MPI_Op operator , Int  
Dest_process, MPI_Comm comm) ;`

# MPI - Reduce

## Syntax :

➔ `MPI_Reduce (void* input_data , void* output_data , Int count ,  
MPI_Datatype datatype , MPI_Op operator , Int  
Dest_process, MPI_Comm comm) ;`

```
MPI_MAX  
MPI_MIN  
MPI_SUM  
MPI_PROD  
MPI_LAND
```

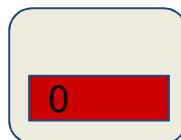
```
:  
:  
:  
:
```

# MPI - Reduce

## Syntax :

➔ `MPI_Reduce (void* input_data , void* output_data , Int count ,  
MPI_Datatype datatype , MPI_Op operator , Int  
Dest_process, MPI_Comm comm) ;`

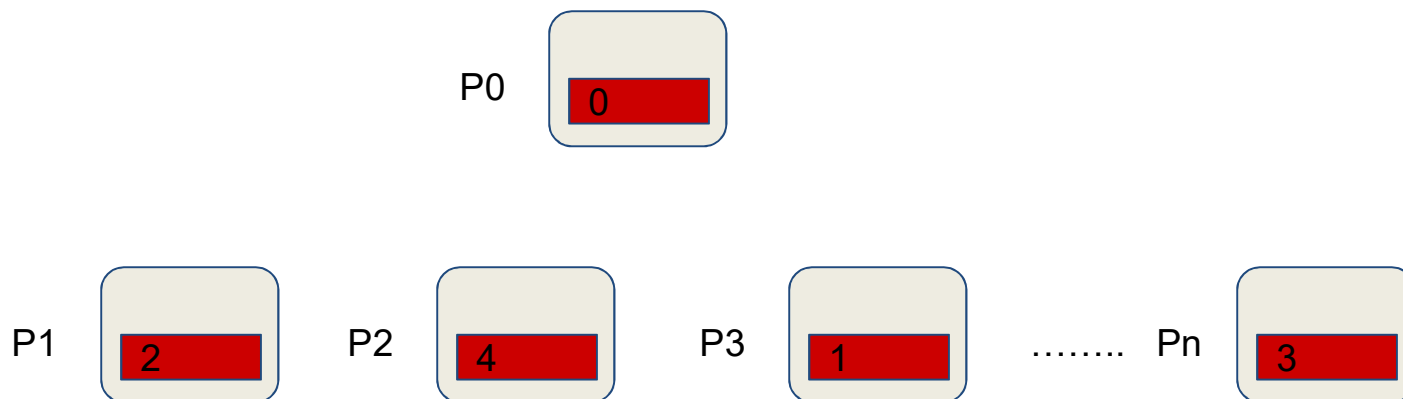
P0



# MPI - Reduce

## Syntax :

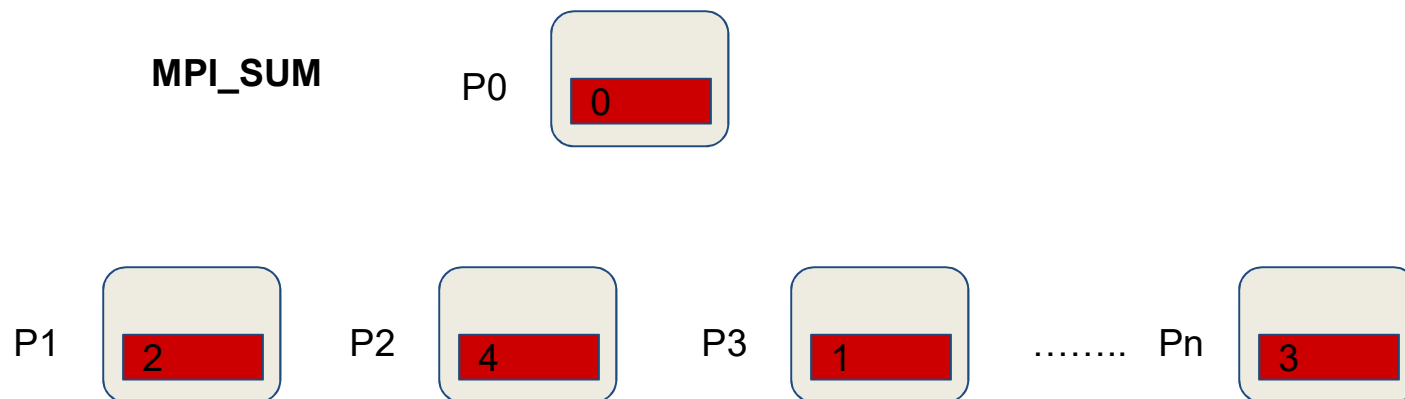
➔ `MPI_Reduce (void* input_data , void* output_data , Int count ,  
MPI_Datatype datatype , MPI_Op operator , Int  
Dest_process, MPI_Comm comm) ;`



# MPI - Reduce

## Syntax :

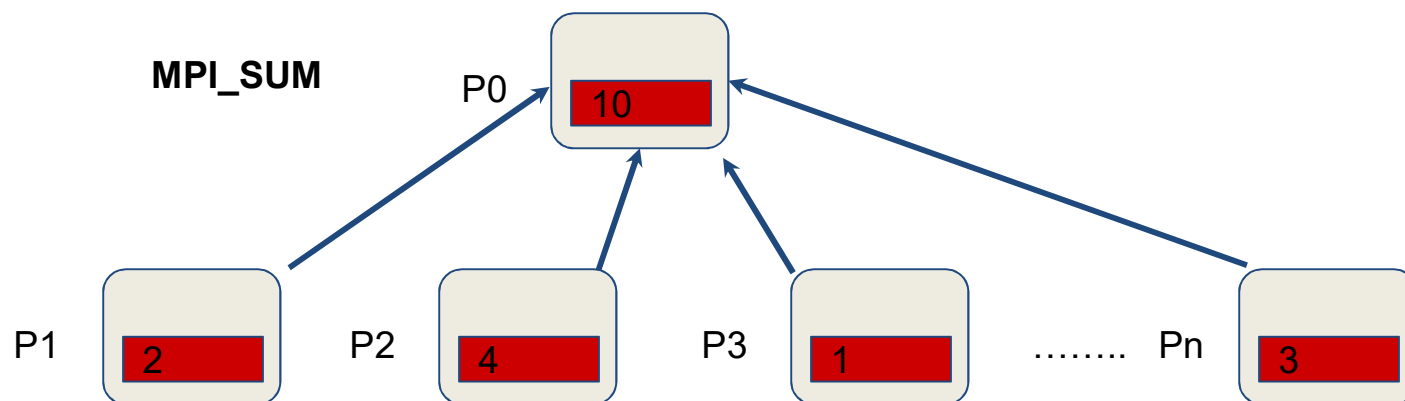
➔ `MPI_Reduce (void* input_data , void* output_data , Int count ,  
MPI_Datatype datatype , MPI_Op operator , Int  
Dest_process, MPI_Comm comm) ;`



# MPI - Reduce

## Syntax :

➔ `MPI_Reduce (void* input_data , void* output_data , Int count ,  
MPI_Datatype datatype , MPI_Op operator , Int  
Dest_process, MPI_Comm comm) ;`



# MPI - Reduce

## Syntax :

➔ `MPI_Reduce (void* input_data , void* output_data , Int count ,  
MPI_Datatype datatype , MPI_Op operator , Int  
Dest_process, MPI_Comm comm) ;`

Example : Many lines in Trap. example programs are replaced by this single line ...



# MPI - Reduce

## Syntax :

➔ `MPI_Reduce (void* input_data , void* output_data , Int count ,  
MPI_Datatype datatype , MPI_Op operator , Int  
Dest_process, MPI_Comm comm) ;`

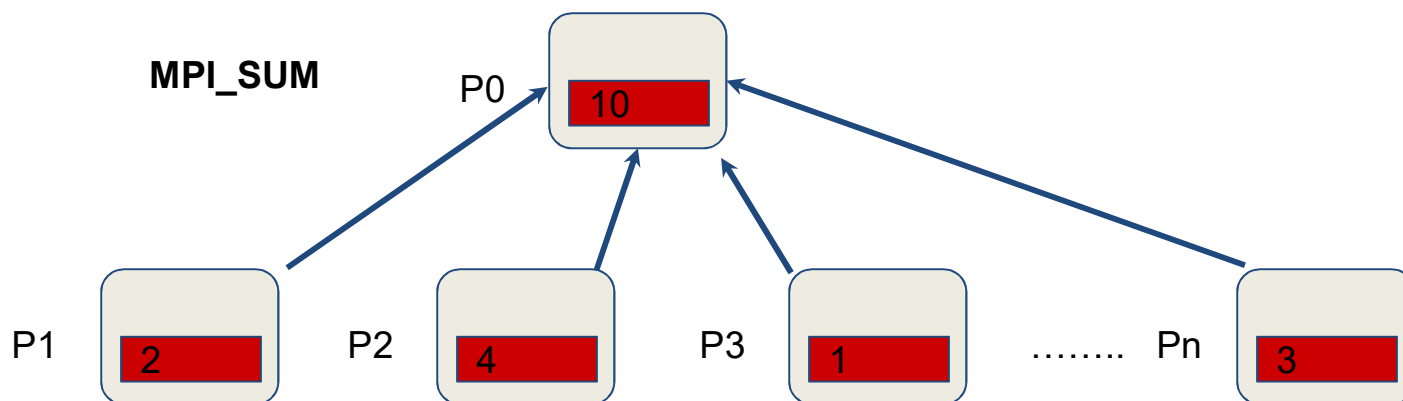
Example : Many lines in Trap. example programs are replaced by this single line ...

➔ `MPI_Reduce(&local_int, &total_int, 1, MPI_DOUBLE, MPI_SUM, 0,  
MPI_COMM_WORLD) ;`

# MPI - Allreduce

## Syntax :

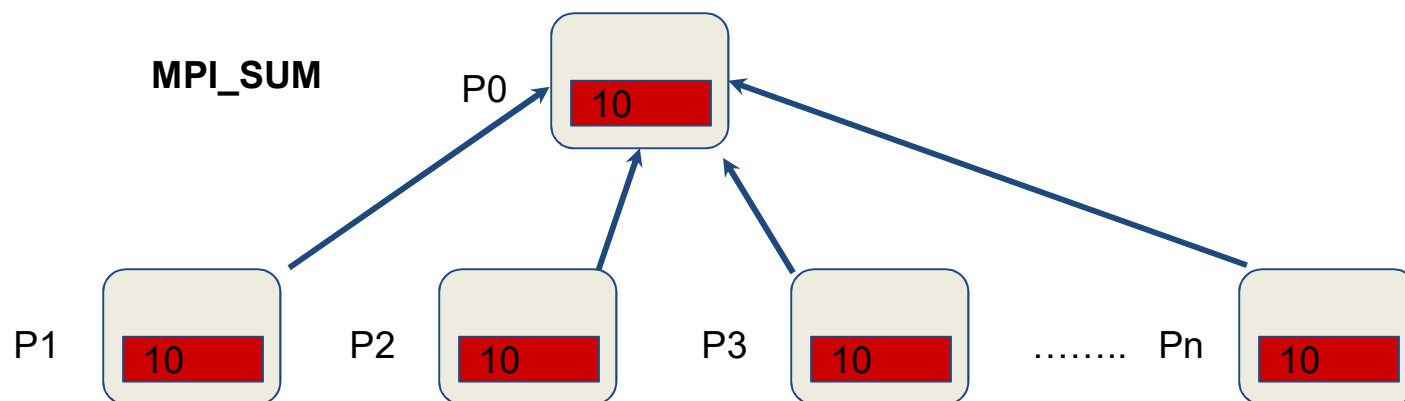
➔ `MPI_Reduce (void* input_data , void* output_data , Int count ,  
MPI_Datatype datatype , MPI_Op operator , MPI_Comm  
comm) ;`



# MPI - Allreduce

## Syntax :

➔ `MPI_Reduce (void* input_data , void* output_data , Int count ,  
MPI_Datatype datatype , MPI_Op operator , MPI_Comm  
comm) ;`

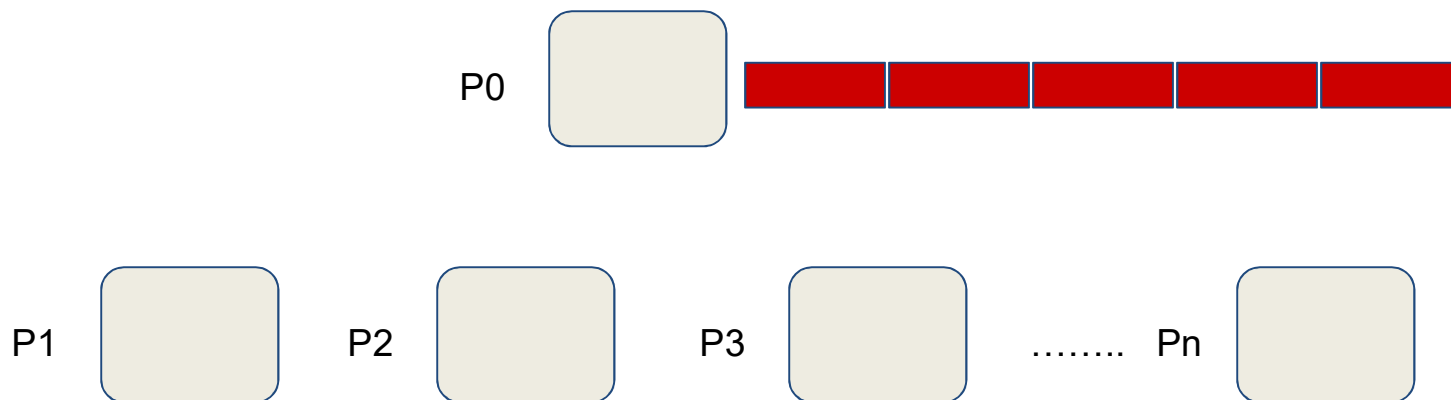


# MPI - Scatter

## Syntax :

➔ `MPI_Scatter (void* send_buffer , Int send_count , MPI_Datatype  
send_datatype , void* recv_buffer , Int recv_count ,  
MPI_Datatype recv_datatype , Int source_process ,  
MPI_Comm comm ) ;`

➤ MPI\_Scatter sends chunks of data to different processes..

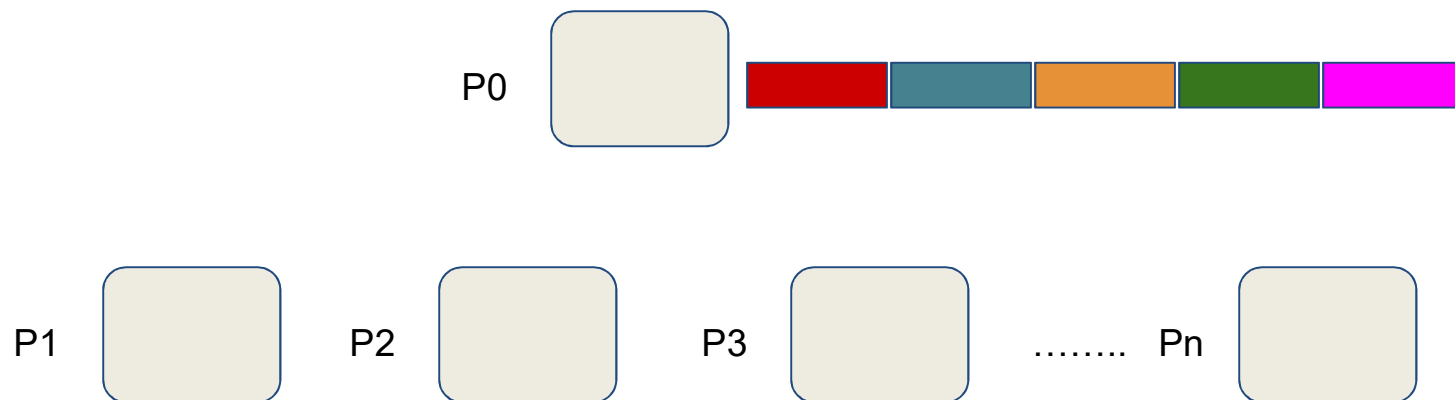


# MPI - Scatter

## Syntax :

➔ `MPI_Scatter (void* send_buffer , Int send_count , MPI_Datatype  
send_datatype , void* recv_buffer , Int recv_count ,  
MPI_Datatype recv_datatype , Int source_process ,  
MPI_Comm comm ) ;`

➤ MPI\_Scatter sends chunks of data to different processes..

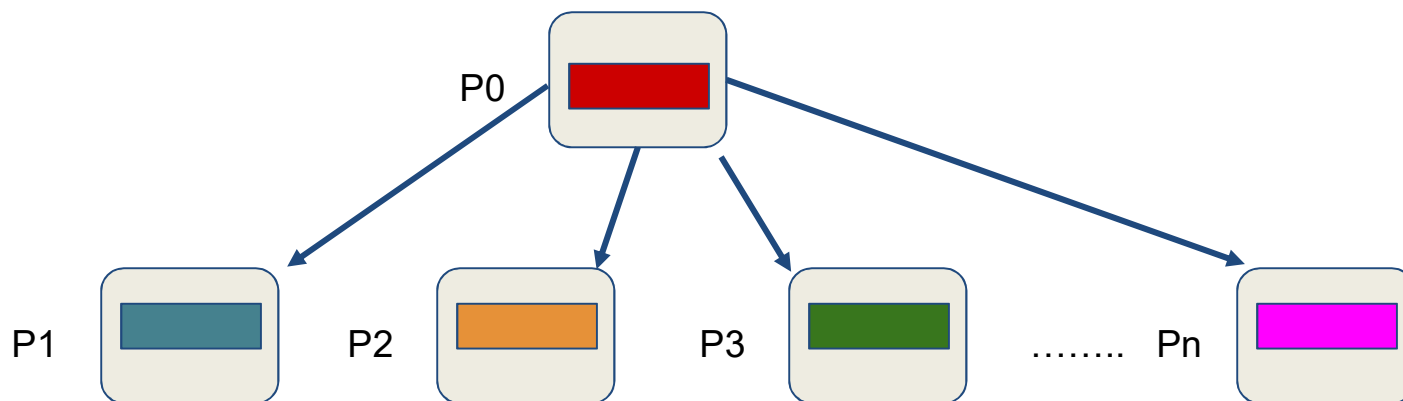


# MPI - Scatter

## Syntax :

➔ `MPI_Scatter (void* send_buffer , Int send_count , MPI_Datatype  
send_datatype , void* recv_buffer , Int recv_count ,  
MPI_Datatype recv_datatype , Int source_process ,  
MPI_Comm comm ) ;`

➤ MPI\_Scatter sends chunks of data to different processes..



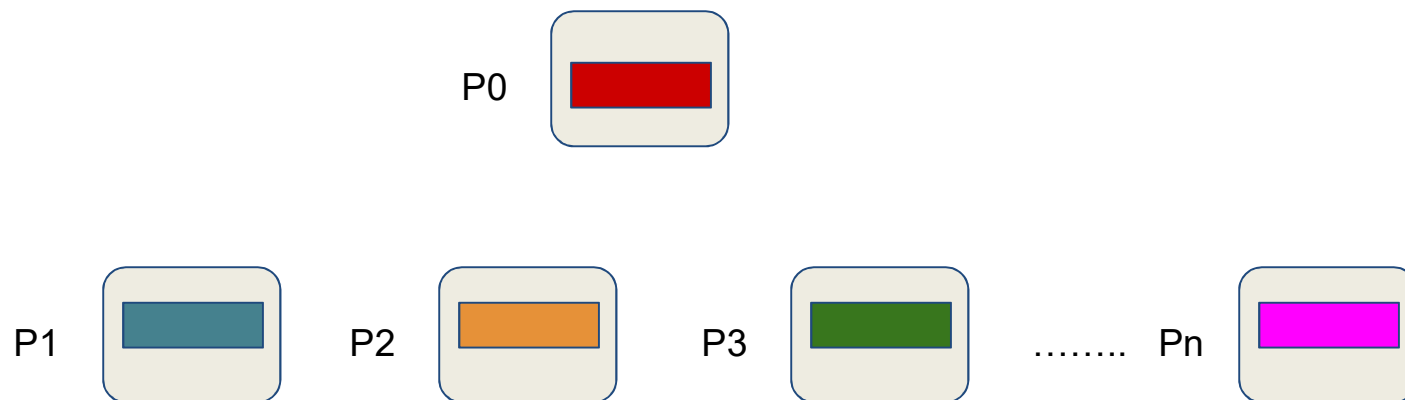
# MPI - Gather

## Syntax :

➔ **MPI\_Gather** (void\* send\_buffer , Int send\_count , MPI\_Datatype  
send\_datatype , void\* recv\_buffer , Int recv\_count ,  
MPI\_Datatype recv\_datatype , Int destination\_process

,  
MPI\_Comm comm ) ;

➤ MPI\_Gather collects chunks of data from different processes..



# MPI - Gather

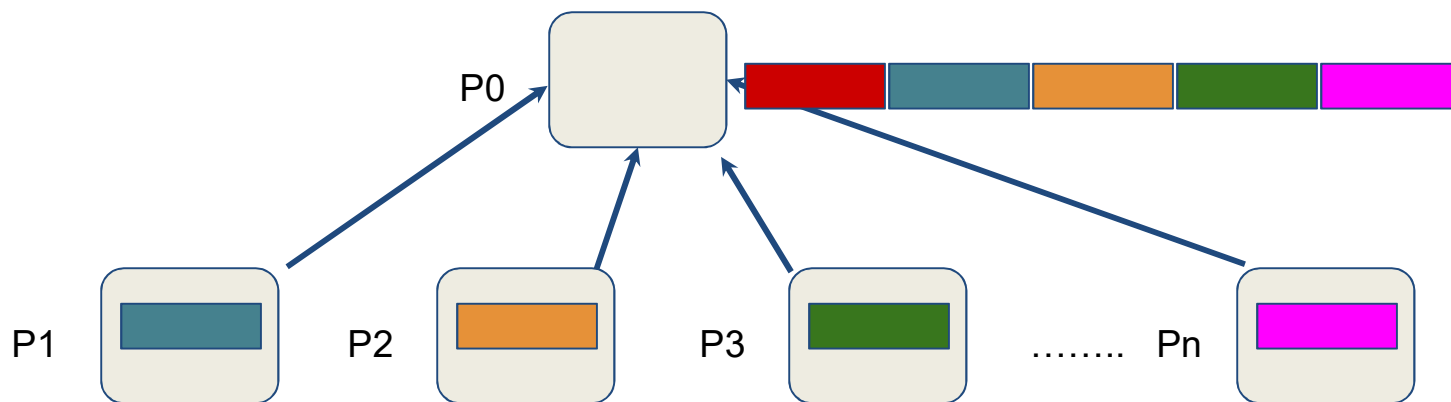
## Syntax :

➔ **MPI\_Gather** (void\* send\_buffer , Int send\_count , MPI\_Datatype  
send\_datatype , void\* recv\_buffer , Int recv\_count ,  
MPI\_Datatype recv\_datatype , Int destination\_process

,

**MPI\_Comm comm ) ;**

➤ MPI\_Gather collects chunks of data from different processes..





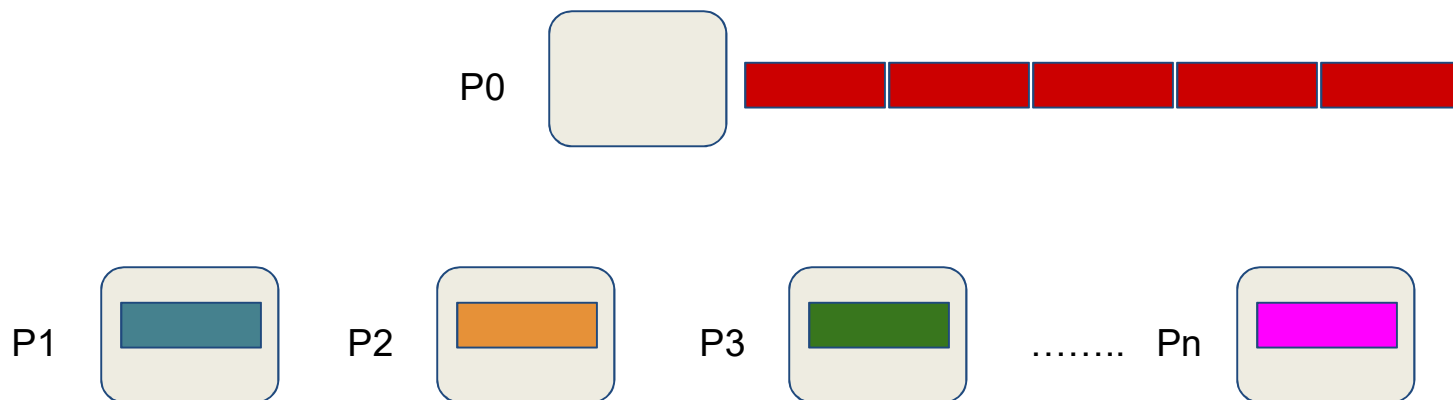
# MPI - Gather

## Syntax :

➔ `MPI_Gather (void* send_buffer , Int send_count , MPI_Datatype  
send_datatype , void* recv_buffer , Int recv_count ,  
MPI_Datatype recv_datatype , Int destination_process`

`,  
MPI_Comm comm ) ;`

➤ MPI\_Gather collects chunks of data from different processes..



# MPI - Allgather

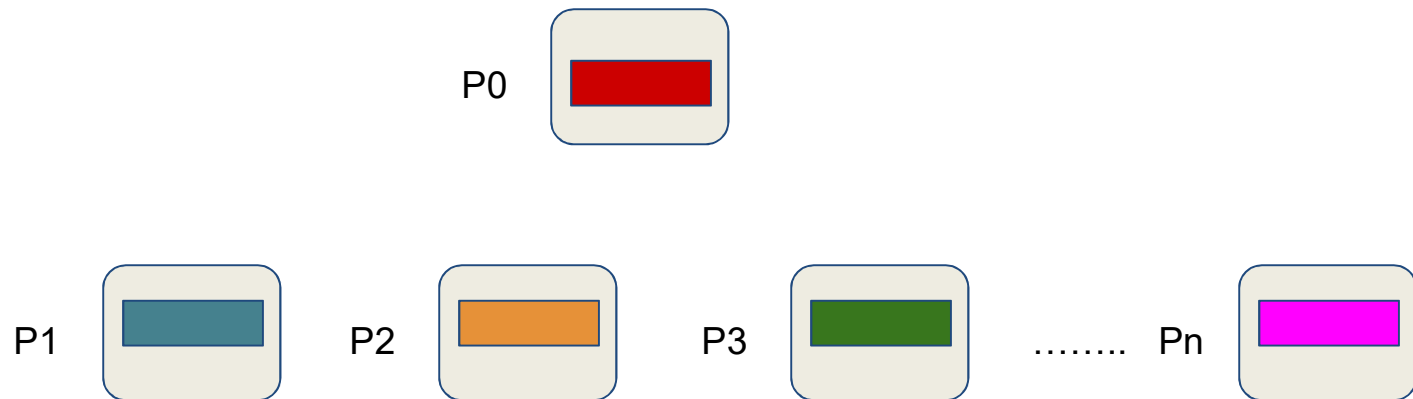
## Syntax :

➔ **MPI\_Gather (void\* send\_buffer , Int send\_count , MPI\_Datatype  
send\_datatype , void\* recv\_buffer , Int recv\_count ,  
MPI\_Datatype recv\_datatype ,  
MPI\_Comm comm ) ;**

# MPI - Allgather

## Syntax :

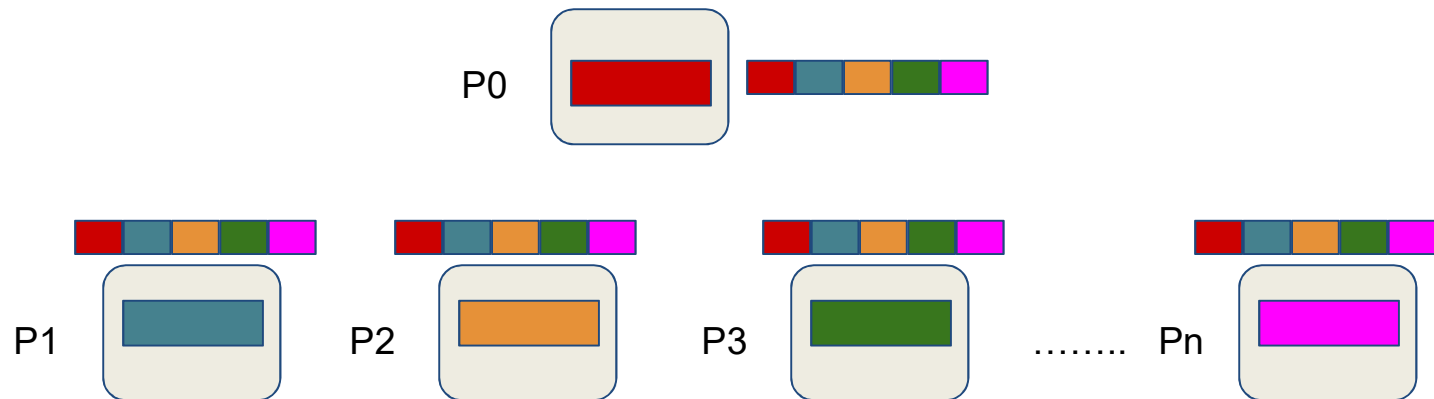
➔ `MPI_Gather (void* send_buffer , Int send_count , MPI_Datatype  
send_datatype , void* recv_buffer , Int recv_count ,  
MPI_Datatype recv_datatype ,  
MPI_Comm comm ) ;`



# MPI - Allgather

## Syntax :

➔ `MPI_Gather (void* send_buffer , Int send_count , MPI_Datatype  
send_datatype , void* recv_buffer , Int recv_count ,  
MPI_Datatype recv_datatype ,  
MPI_Comm comm ) ;`



# MPI - Synchronization

# MPI - Barrier

**Syntax :**

➔ `MPI_Barrier (MPI_Comm communicator) ;`

# MPI - Barrier

## Syntax :

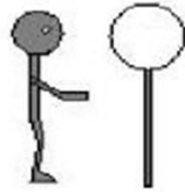
➔ **MPI\_Barrier (MPI\_Comm communicator) ;**

- Used to block the calling process until all processes have entered the function. The call will return at any process only after all the processes or group members have entered the call
- The MPI\_BARRIER routine blocks the calling process until all group processes have called the function. When MPI\_BARRIER returns, all processes are synchronized at the barrier

# MPI - Barrier

**Syntax :**

➔ `MPI_Barrier (MPI_Comm communicator) ;`

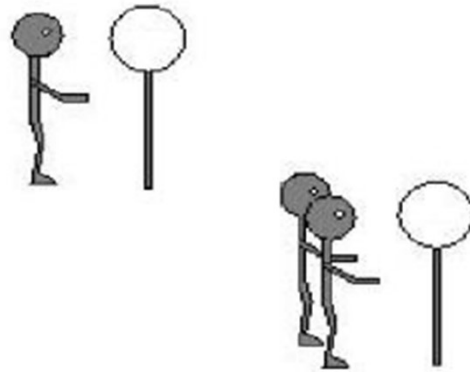




# MPI - Barrier

**Syntax :**

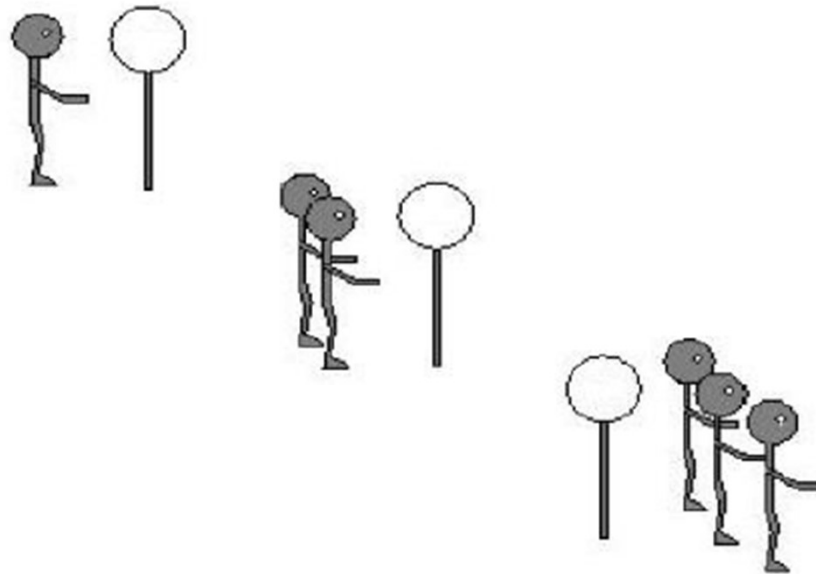
➔ `MPI_Barrier (MPI_Comm communicator) ;`



# MPI - Barrier

**Syntax :**

➔ `MPI_Barrier (MPI_Comm communicator) ;`



## Recap :

- Point to Point Vs Collective communication -
- MPI\_Broadcast(...)
- MPI\_Scatter(...)
- MPI\_Reduce(...)
- MPI\_Allreduce(...)
- MPI\_Gather(...)
- MPI\_Allgather(...)
- Miss MPI routines !
- ..... ..

A decorative graphic on the left side of the slide, featuring a dark green curved shape and a stylized map of India in green and yellow.

## References :

- [1] Barker, Brandon. "Message passing interface (mpi)." *Workshop: High Performance Computing on Stampede*. Vol. 262. 2015.
- [2] Yuan, Chung-Tsz, and Shenjian Chen. "Message Passing Interface (MPI)." (1996).
- [3] <https://computing.llnl.gov/tutorials/mpi/>

Thank  
You

