



International  
Centre for  
Radio  
Astronomy  
Research

# MPI

Research Associate Professor  
Kevin Vinsen



Curtin University



THE UNIVERSITY OF  
WESTERN  
AUSTRALIA



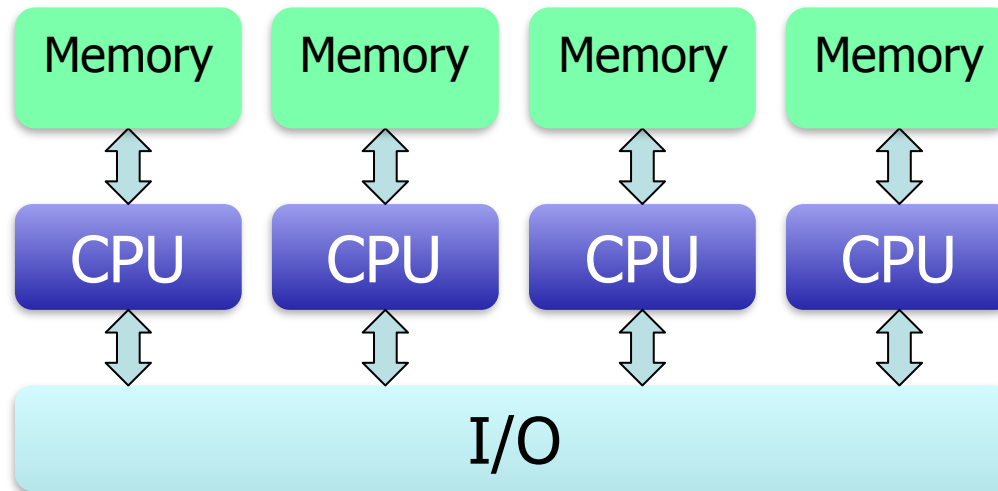
Government of Western Australia  
Department of the Premier and Cabinet  
Office of Science



# Distributed Memory HPC

Communicate via messages

Message Passing Interface (MPI)





# MPI - Message Passing Interface

---

- **MPI is a standard, not a library!**
  - It defines the functions, constants, and behaviours for a set of commands to exchange information between computers.
- **Many libraries implement the MPI standard. ie:**
  - mpich2
  - lam-mpi
  - OpenMPI
  - HP MPI
  - Intel MPI



# MPI binding

---

**MPI is not a programming language**

- **MPI libraries support**

- C
- C++
- Fortran
- Python
- Java
- many others!



# What MPI can Not

---

**MPI cannot magically transform your program into an efficient parallel program**

- 1. The algorithm must be suitable for parallelization**
- 2. Communication between processors takes time**
- 3. Remember the model**



# When Not to Parallelise Code

---

- **Code will only be used once (or infrequently)**  
Efficient parallel code takes time to develop!
- **Current performance is acceptable and execution time is short**
- **Frequent & significant code changes**
- **Some algorithms simply do not parallelize**



# When To Parallelise Code

---

**Code is physically incapable of running on one computer**

- memory requirements are too great
- run time would be months

**Code will be reused frequently**

- Parallelization is an investment

**Data structures are simple, calculations are local**

- Easy to communicate and synchronize between processors



# MPI Programming Model

---

MPI employs MIMD in SPMD framework.

1. The user issues a directive to the operating system that has the effect of placing a copy of the executable program on each processor.
2. Each processor begins execution of its copy of the executable.
3. Different processes can execute different statements by branching within the program based on their process identity (rank)

```
    if (my_rank != 0)
```

```
        .  
        .  
        .
```

```
    else
```

```
        .  
        .  
        .
```





# MPI Basics

---

## 1. Starting and Finishing

<i>MPI_Init</i>	initialise MPI
<i>MPI_Finalize</i>	terminate computation

## 2. Identifying yourself

<i>MPI_Comm_size</i>	number of processes
<i>MPI_Comm_rank</i>	my process identifier

## 3. Sending and Receiving messages

<i>MPI_Send</i>	send a message
<i>MPI_Recv</i>	receive a message



# MPI\_Init

---

C: *int MPI\_Init(int \*argc, /\* in/out \*/  
char \*\*argv) /\* in/out \*/*

Fortran: **call MPI\_Init(ierr)**

## Initialises MPI subroutines

- Connects to other processes
- Processes any MPI specific command line arguments
- Must be called before any other MPI function is used
- Must be called before your program examines any command lines arguments



# MPI\_Finalize

---

C: *MPI\_Finalize(void)*

Fortran: **call MPI\_Finalize(ierr)**

## Cleans up MPI sessions

- Must be called by all processes before exiting
- After being called, MPI functions cannot be used
- If forgotten, program will sometimes (but not always) crash at the end with strange errors



# MPI Communicator

---

## Communicator

- Collection of processes
- Determines scope to which messages are related
- Identifier of process (rank) is relative to communicator
- Defines the scope of communications (broadcast, etc.)



# MPI\_Comm\_size

---

**C:** *MPI\_Comm\_size(MPI\_Comm comm, /\* in \*/  
int \*size) /\* out \*/*

**Fortran:** *call MPI\_Comm\_size(MPI\_Comm, integer, integer ierr)*

**Retrieves the number of processes in the communicator**

The *MPI\_Comm comm* parameter specifies the “communicator”, or communication group, to perform the function on  
In general, use *MPI\_COMM\_WORLD*, which means all processes



# MPI\_Comm\_rank

---

**C:** *MPI\_Comm\_rank(MPI\_Comm comm, /\* in \*/  
int \*rank) /\* out \*/*

**Fortran:** *call MPI\_Comm\_rank(MPI\_Comm, integer, integer ierr)*

**Retrieves the rank of the process**

- Rank is 0-based. (ie. the first process has rank 0, the last has a rank of nProcesses-1)
- Use this value to decide which work this process should do



# MPI Messages

---

**Message content, a sequence of bytes**

**Message needs wrapper (envelope for a letter)**

- Destination
- Source
- Message type
- Size (count)
- Communicator
- Broadcast



# Point to Point Blocking

---

**MPI\_Send**

**send a message**

**MPI\_Recv**

**receive a message**





# MPI message protocol

---

Communicator (sixth parameter in `MPI_Send` and `MPI_Recv`) determines the context for destination and source ranks  
`MPI_COMM_WORLD` is automatically supplied communicator, which includes all processes created at start-up  
Other communicators can be defined by user to group processes and to create virtual topologies



# MPI\_Send

---

```
MPI_Send(void *buf,          /* in */
         int count,          /* in */
         MPI_Datatype datatype, /* in */
         int dest,           /* in */
         int tag,            /* in */
         MPI_Comm comm)      /* in */
```

Fortran: **INTEREG err**

*buf* is a pointer to the data you want to send

*count* is the number of data elements to send

*datatype* is the type of the data

MPI\_INT, MPI\_FLOAT, MPI\_DOUBLE, etc

*dest* is the *rank* of the process that should receive data

*tag* is a unique message identifier

*comm* is generally MPI\_COMM\_WORLD



# MPI\_Recv

---

```
MPI_Recv(void *buf,          /* out */
         int count,          /* in */
         MPI_Datatype datatype, /* in */
         int src,            /* in */
         int tag,            /* in */
         MPI_Comm comm,      /* in */
         MPI_Status *status) /* out */
```

Fortran: **INTEREG err**

*buf* is a pointer to the location the data should be stored

*count* is the number of data elements to receive

*datatype* is the type of the data

MPI\_INT, MPI\_FLOAT, MPI\_DOUBLE, etc

*src* is the *rank* of the process that sends the data (or MPI\_ANY\_SOURCE)

*tag* is a unique message identifier (or MPI\_ANY\_TAG)

*comm* is generally MPI\_COMM\_WORLD

*status* contains information about the message received



# MPI message protocol

---

Status of message received by MPI\_Recv is returned in the *status* parameter

Number of items actually received can be determined from status by using function MPI\_Get\_count

```
int MPI_Get_count( MPI_Status *status, MPI_Datatype  
datatype, int *count )
```



# First MPI Program (in C)

---

```
#include <mpi.h>
```

```
#include <string.h>
```

```
#include <stdio.h>
```

```
int main(int argc, char* argv[])
```

```
{
```

```
    int my_rank;                // rank of process
```

```
    int p;                     // number of processes
```

```
    int source;                // rank of sender
```

```
    int dest;                  // rank of receiver
```

```
    int tag=0;                 // tag for messages
```

```
    char message[100];         // storage for messages
```

```
    MPI_Status status;         // return status for receive
```

```
    ...
```



# First MPI Program (cont.)

---

*/\* Start up MPI \*/*

```
MPI_Init(&argc, &argv);
```

*/\* Find out process rank \*/*

```
MPI_Comm_rank(MPI_COMM_WORLD, &my_rank);
```

*/\* Find out number of processes \*/*

```
MPI_Comm_size(MPI_COMM_WORLD, &p);
```



# First MPI Program (cont.)

---

```
if(my_rank != 0) {                                // for slave processes
    sprintf(message, "Greetings from process %d!", my_rank);
    dest = 0;

    /* lets send message */
    MPI_Send(message, strlen(message)+1, MPI_CHAR, dest,
              tag, MPI_COMM_WORLD);
}
```



# First MPI Program (cont.)

---

```
else {                                     // for master process
    printf("I am the Master! My rank is %d.\n ", my_rank);
    for(source = 1; source < p; source++) {
        /* lets receive message from source */
        MPI_Recv(message, 100, MPI_CHAR, source, tag,
MPI_COMM_WORLD, &status);
        printf("%s\n", message);
    }
}
```





# First MPI Program (cont.)

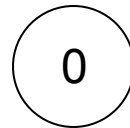
---

```
/* Shut down MPI */  
    MPI_Finalize();  
}
```



# Blocking Type of Communications in Trapezoid Rule Program

---



Active

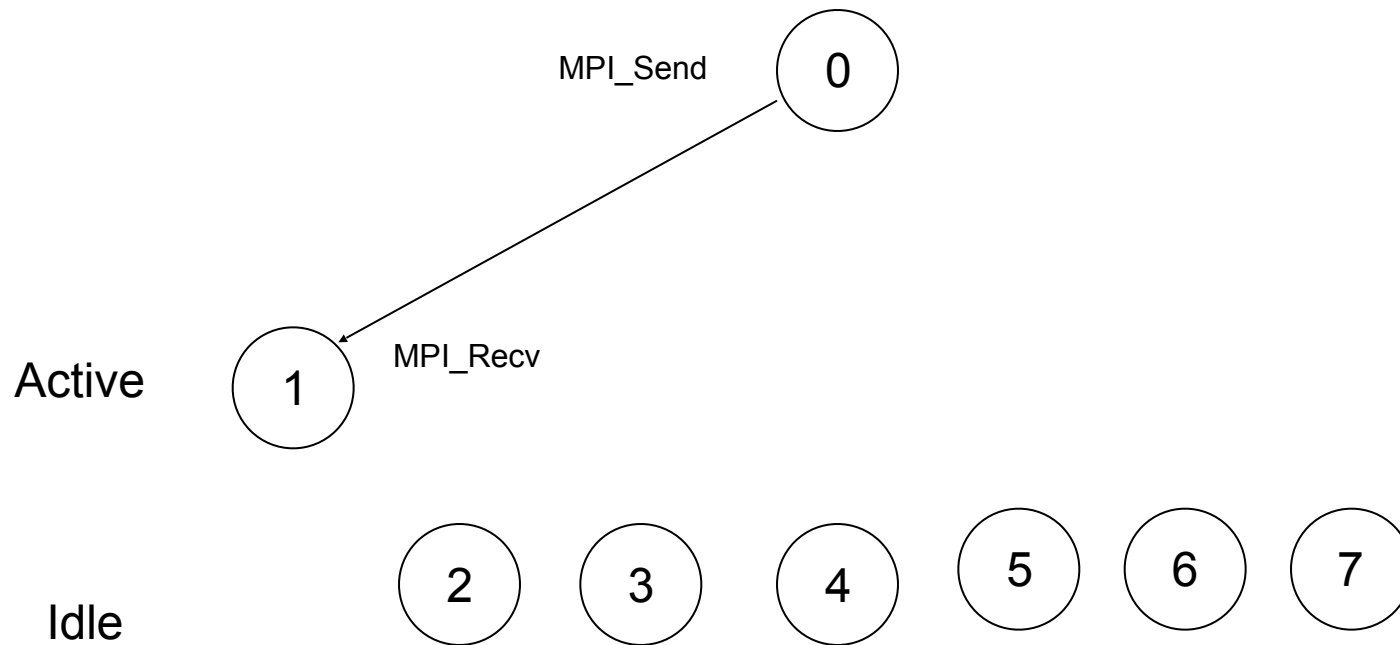
Idle





# Blocking Type of Communications in Trapezoid Rule Program

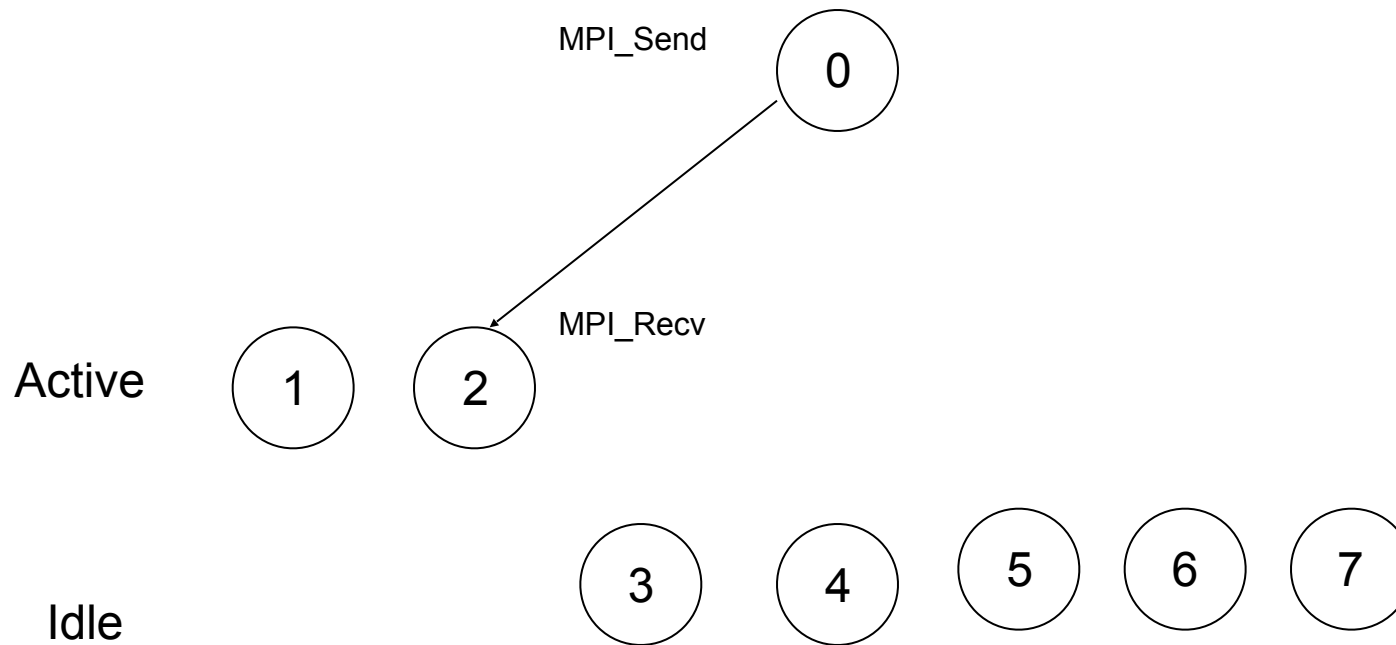
---





# Blocking Type of Communications in Trapezoid Rule Program

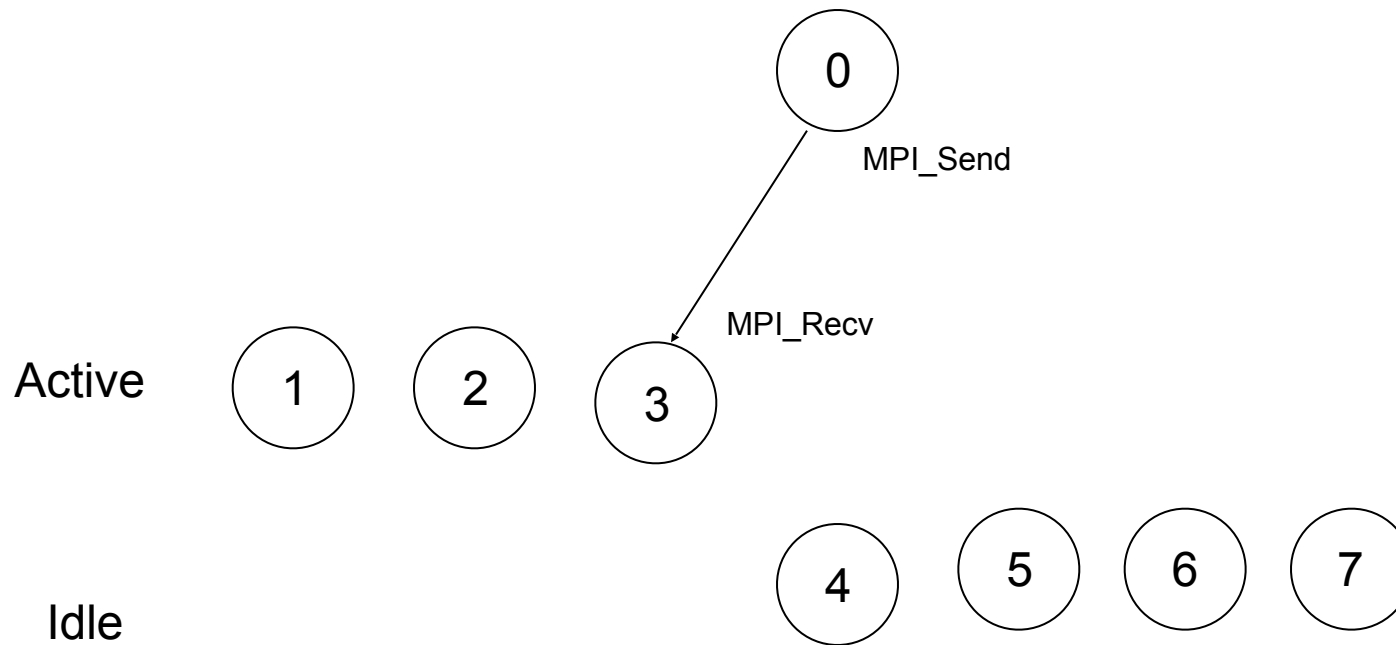
---





# Blocking Type of Communications in Trapezoid Rule Program

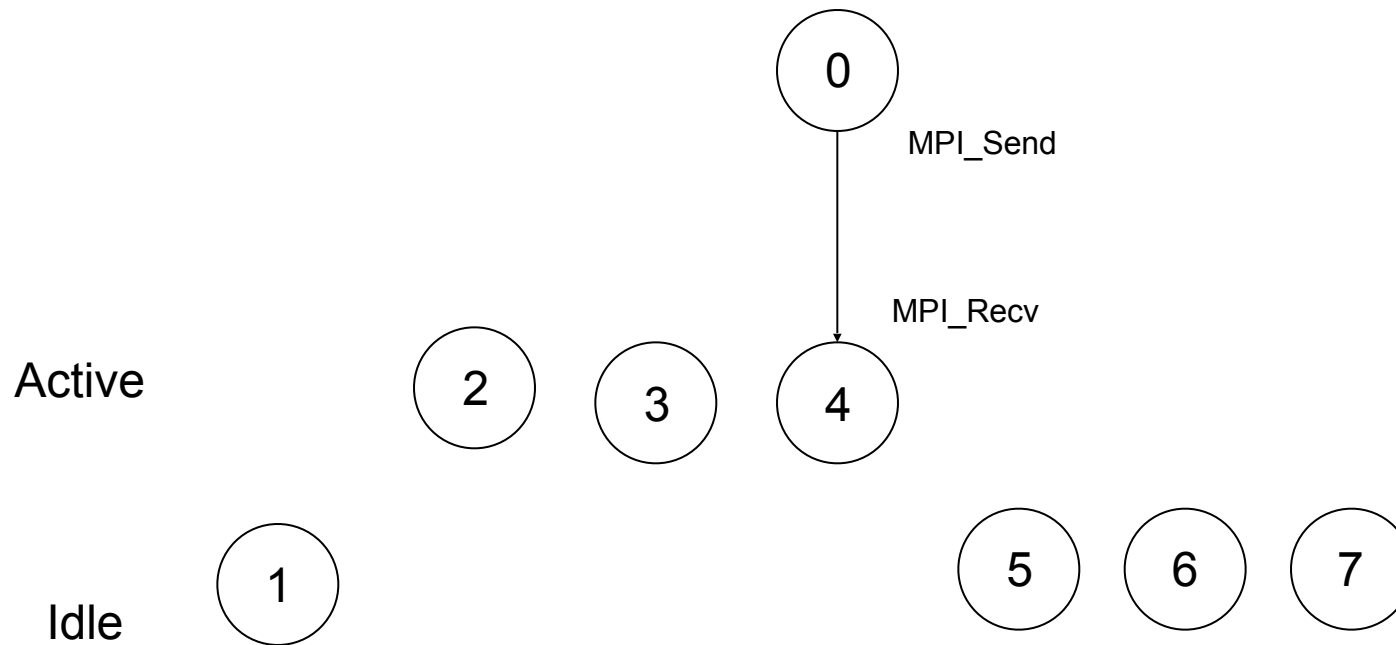
---





# Blocking Type of Communications in Trapezoid Rule Program

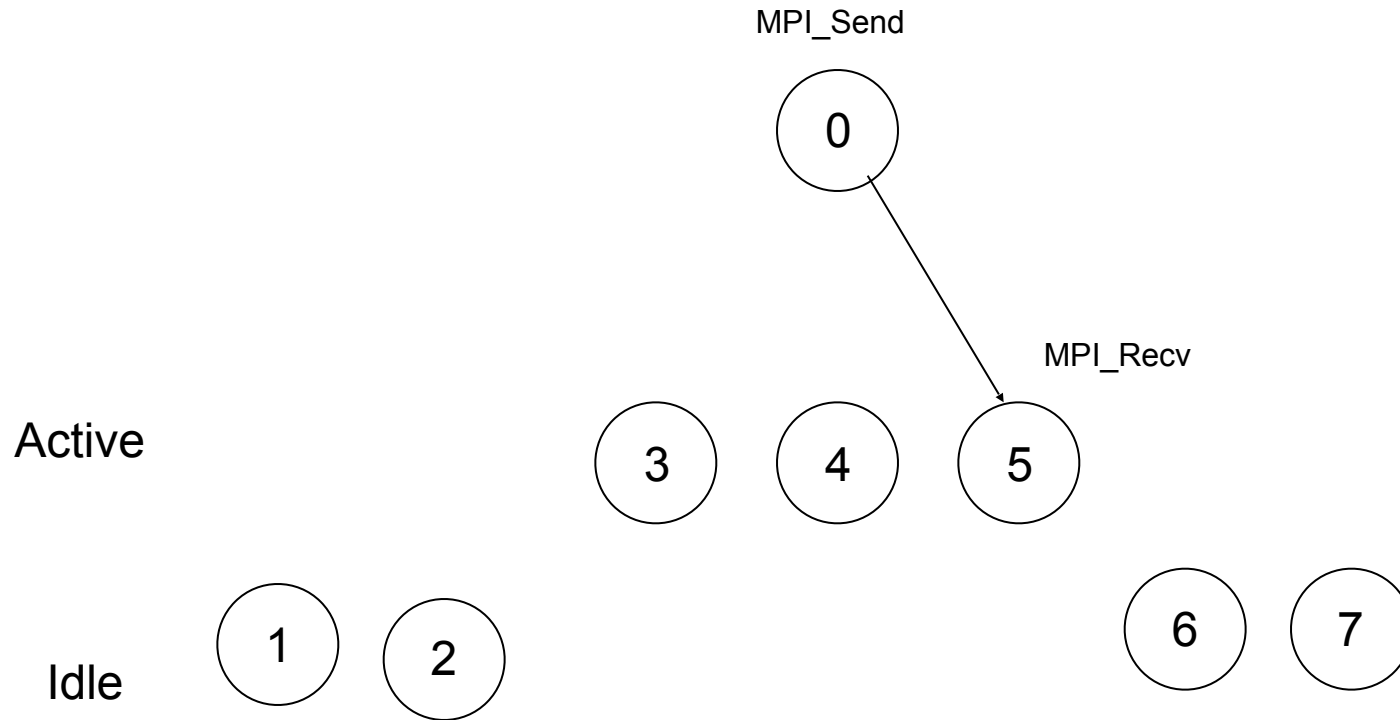
---





# Blocking Type of Communications in Trapezoid Rule Program

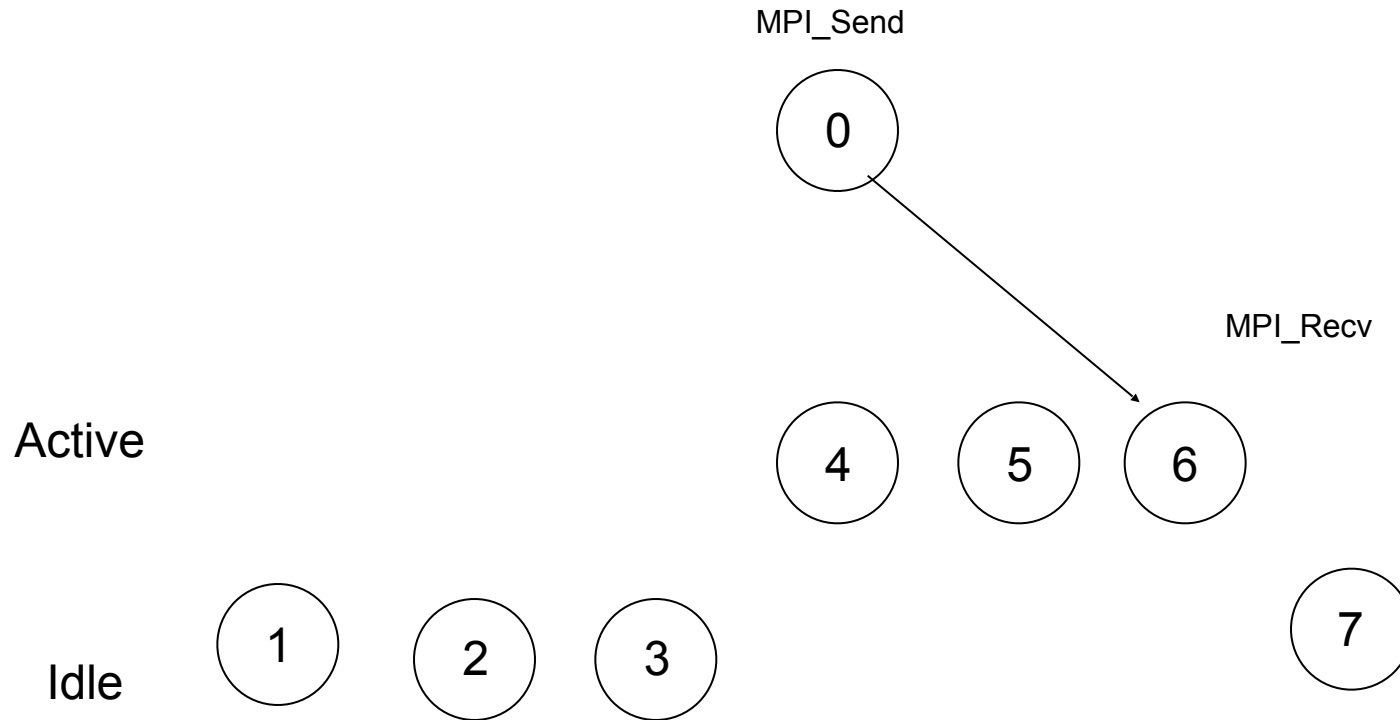
---





# Blocking Type of Communications in Trapezoid Rule Program

---

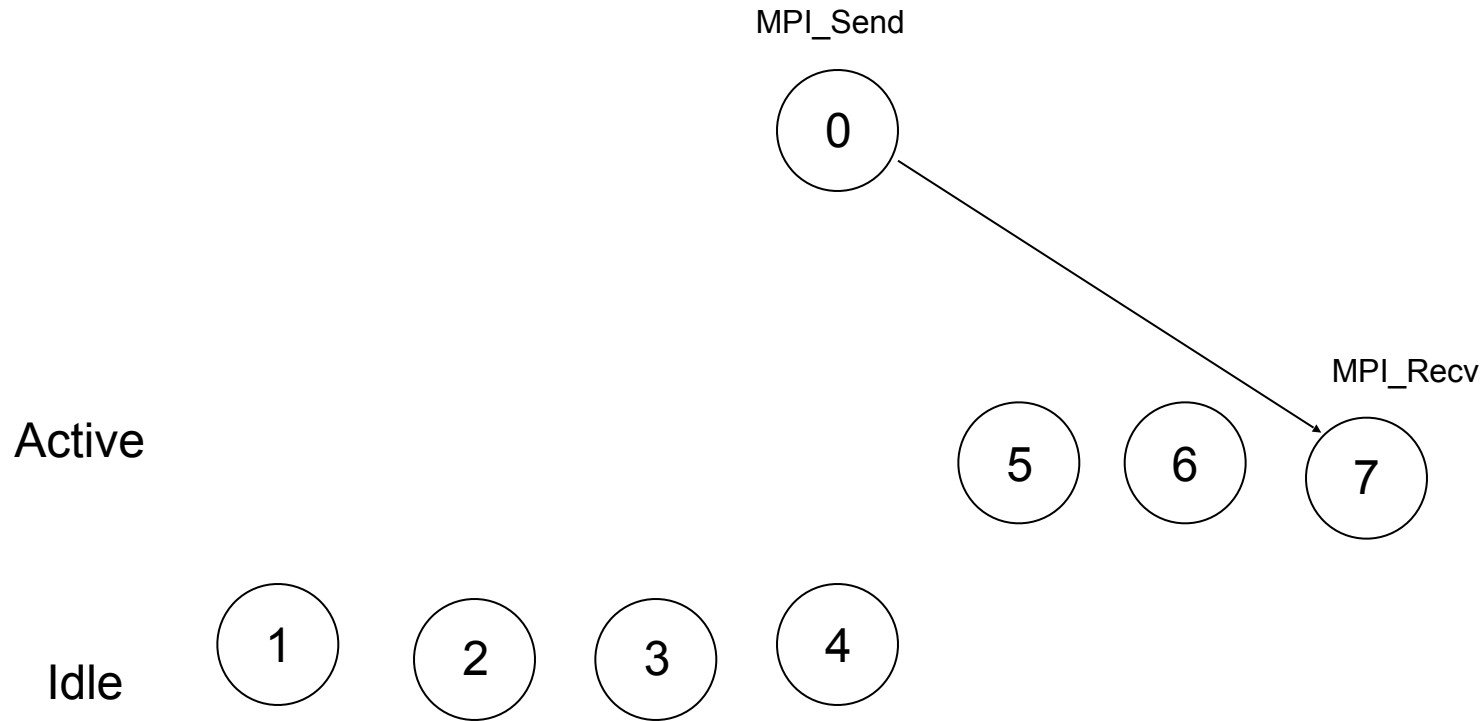






# Blocking Type of Communications in Trapezoid Rule Program

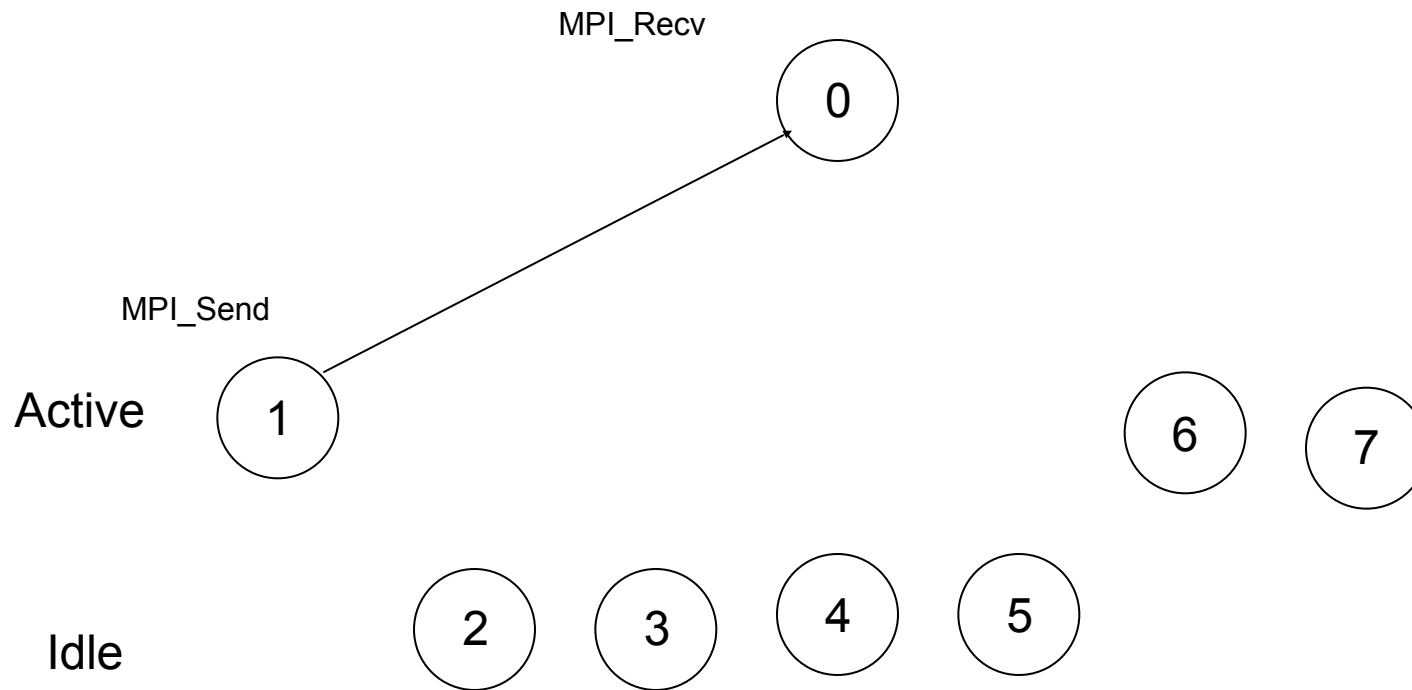
---





# Blocking Type of Communications in Trapezoid Rule Program

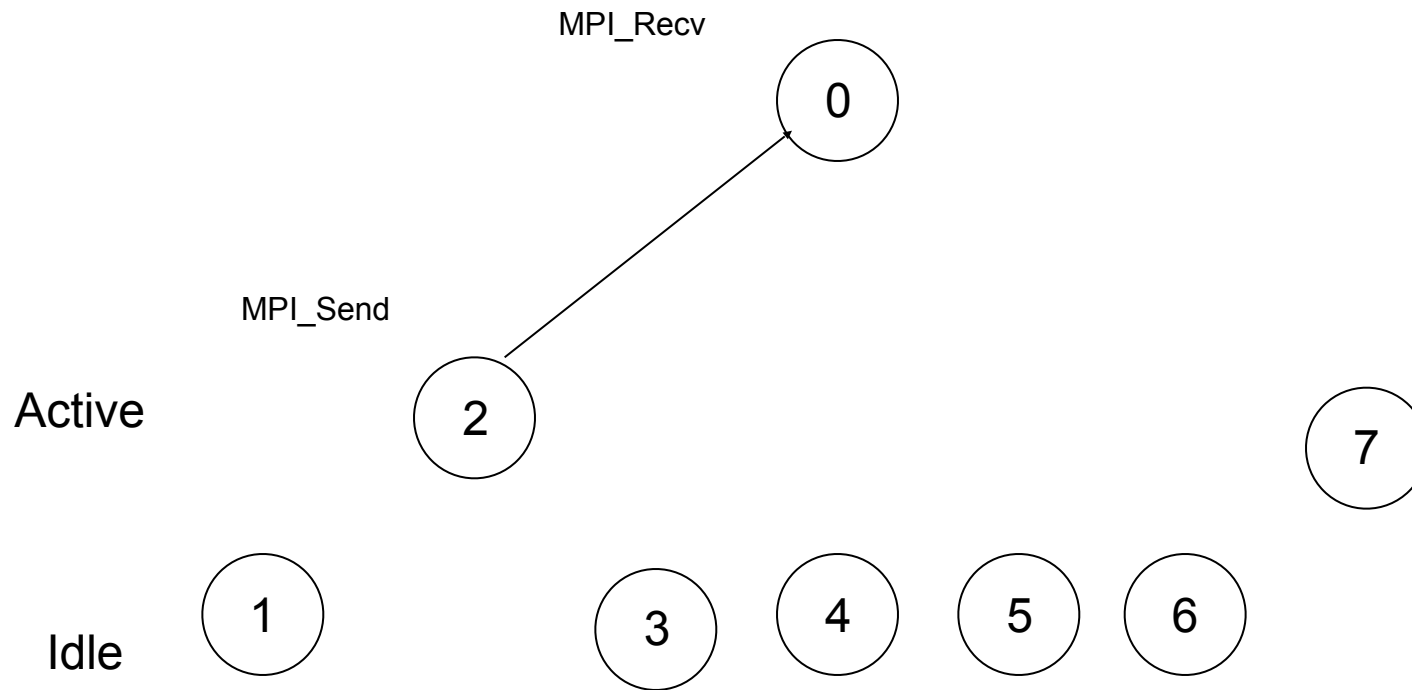
---





# Blocking Type of Communications in Trapezoid Rule Program

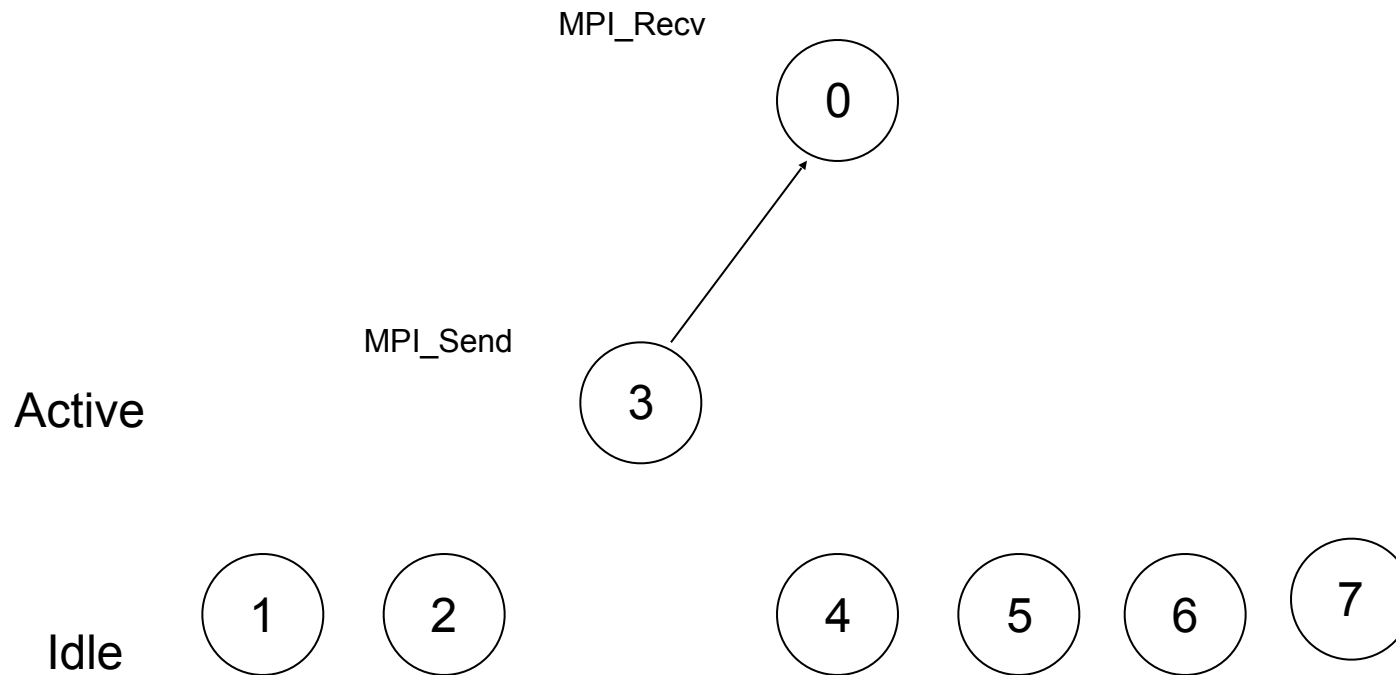
---





# Blocking Type of Communications in Trapezoid Rule Program

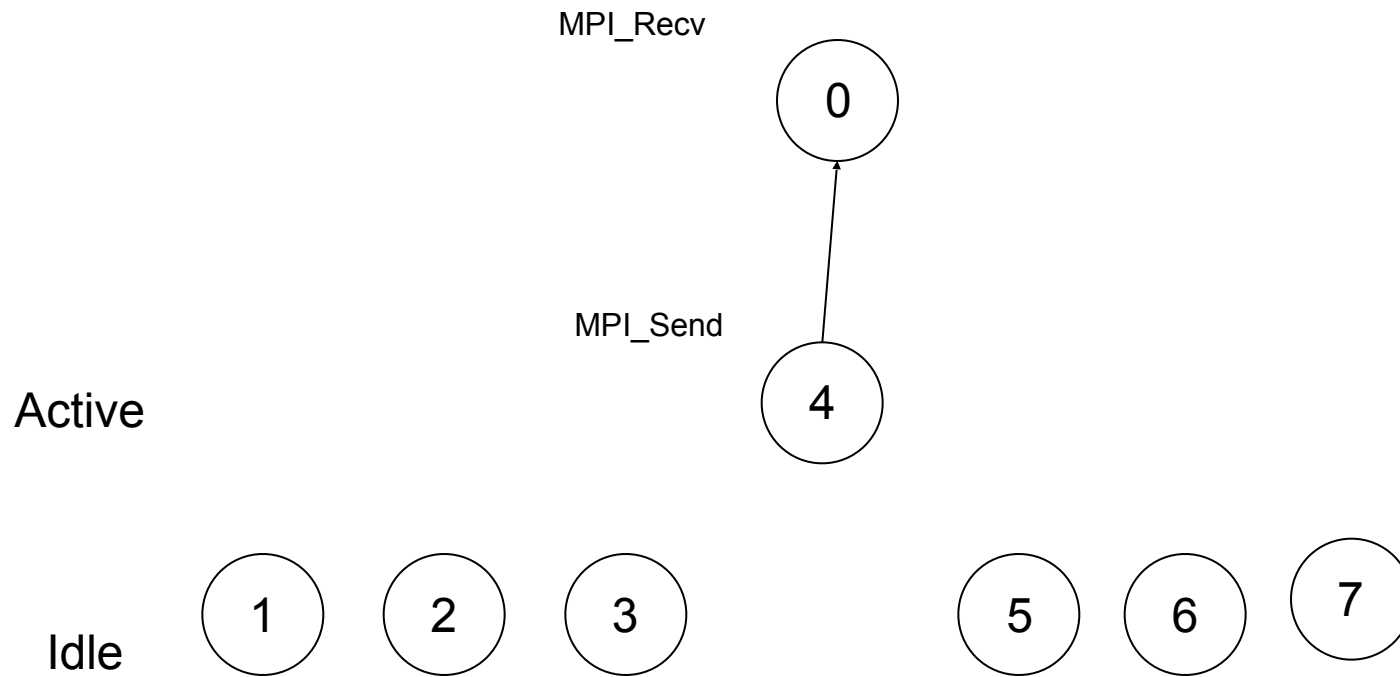
---





# Blocking Type of Communications in Trapezoid Rule Program

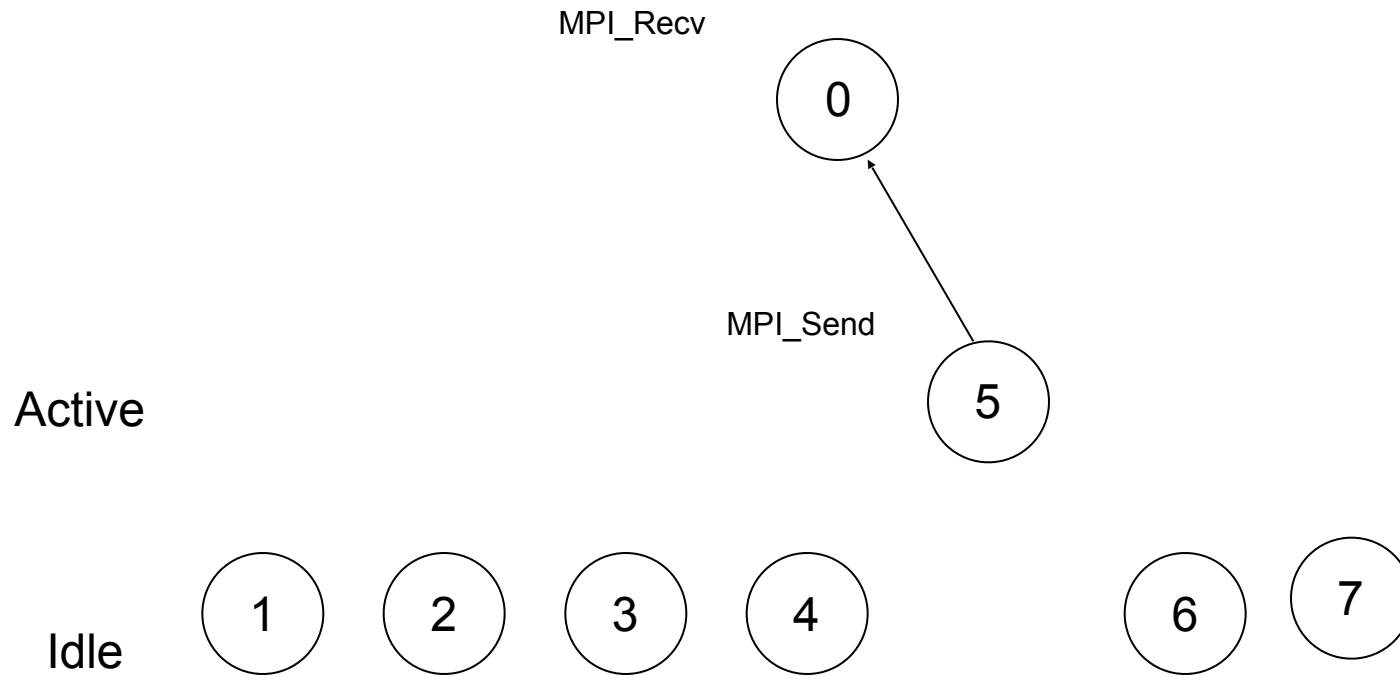
---





# Blocking Type of Communications in Trapezoid Rule Program

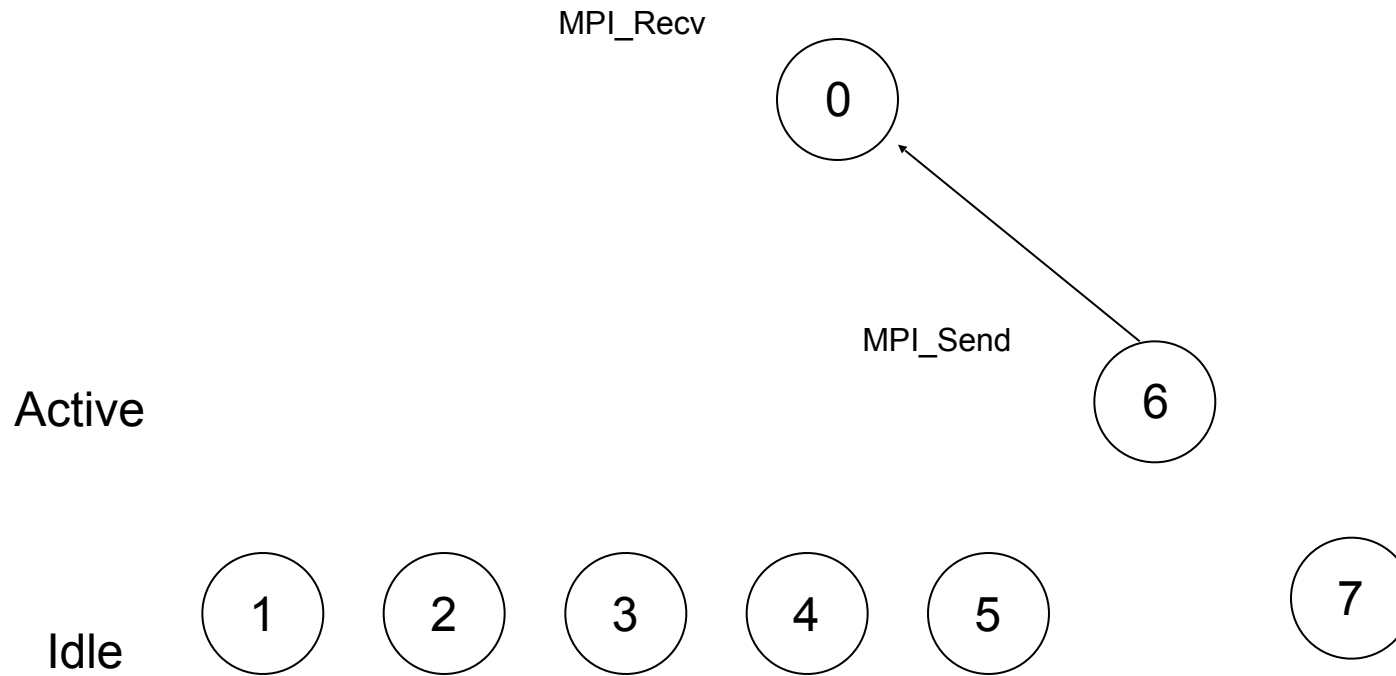
---





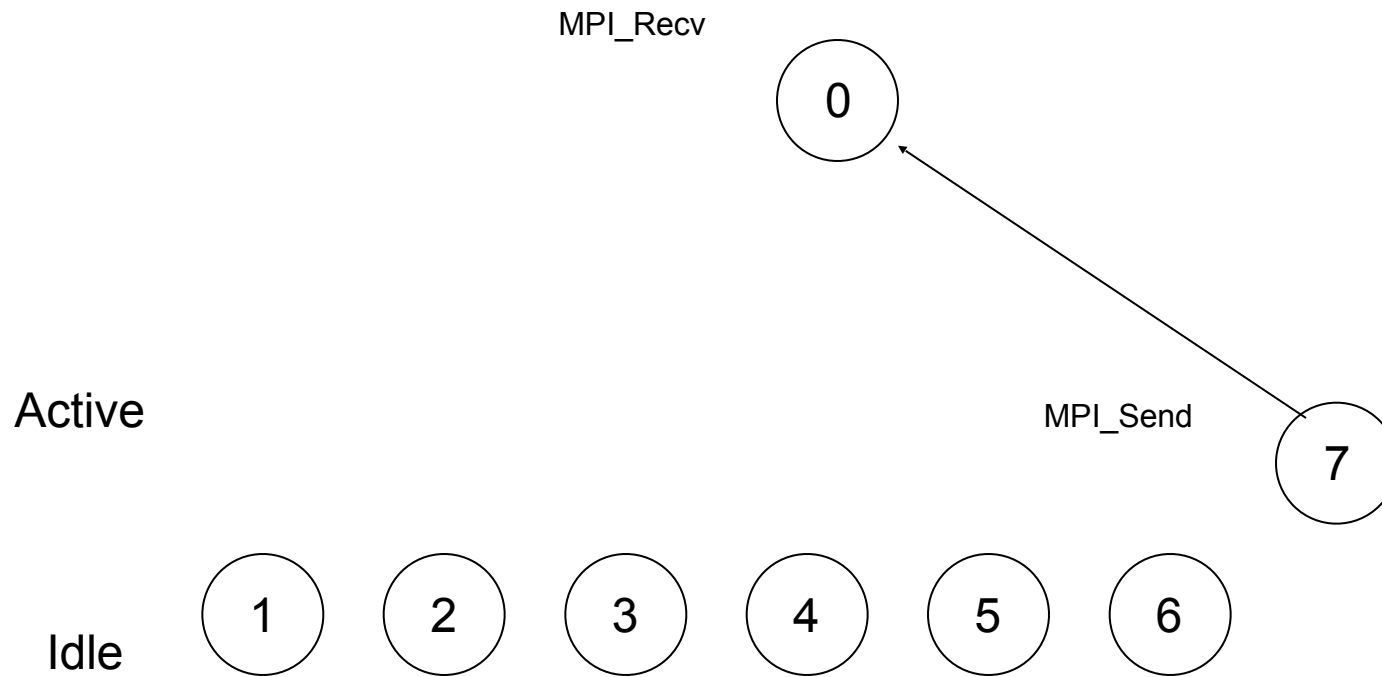
# Blocking Type of Communications in Trapezoid Rule Program

---





# Blocking Type of Communications in Trapezoid Rule Program

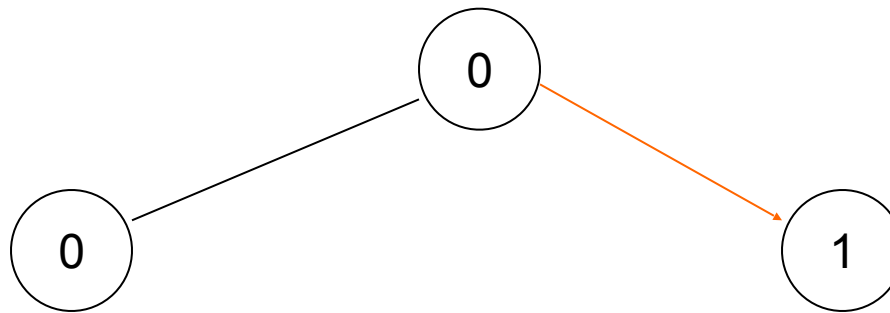






# Processed Configured as a Tree

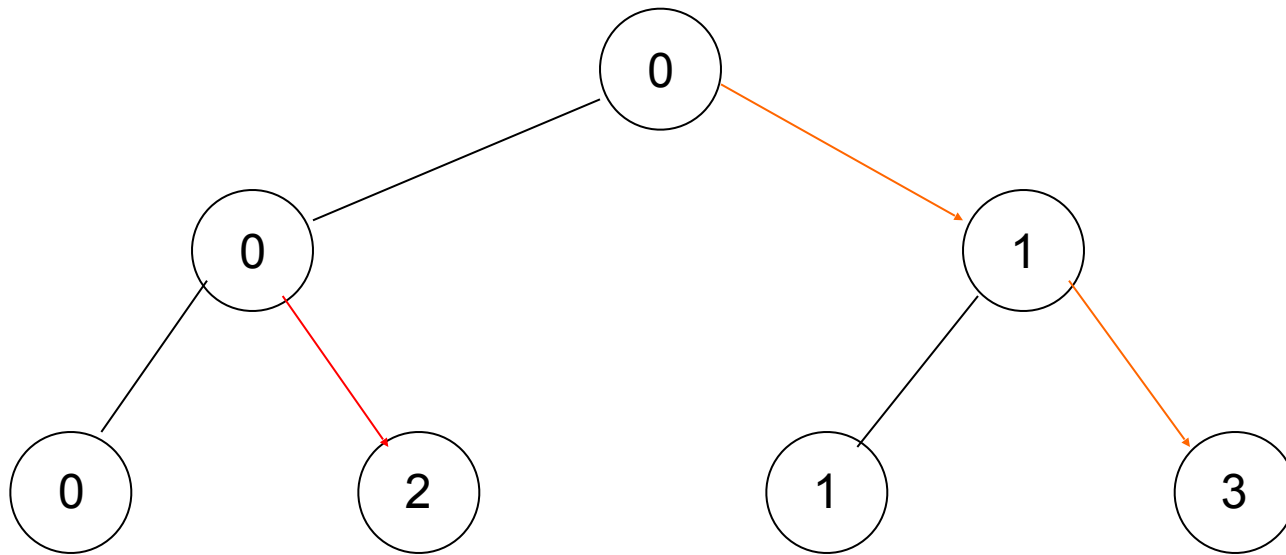
---





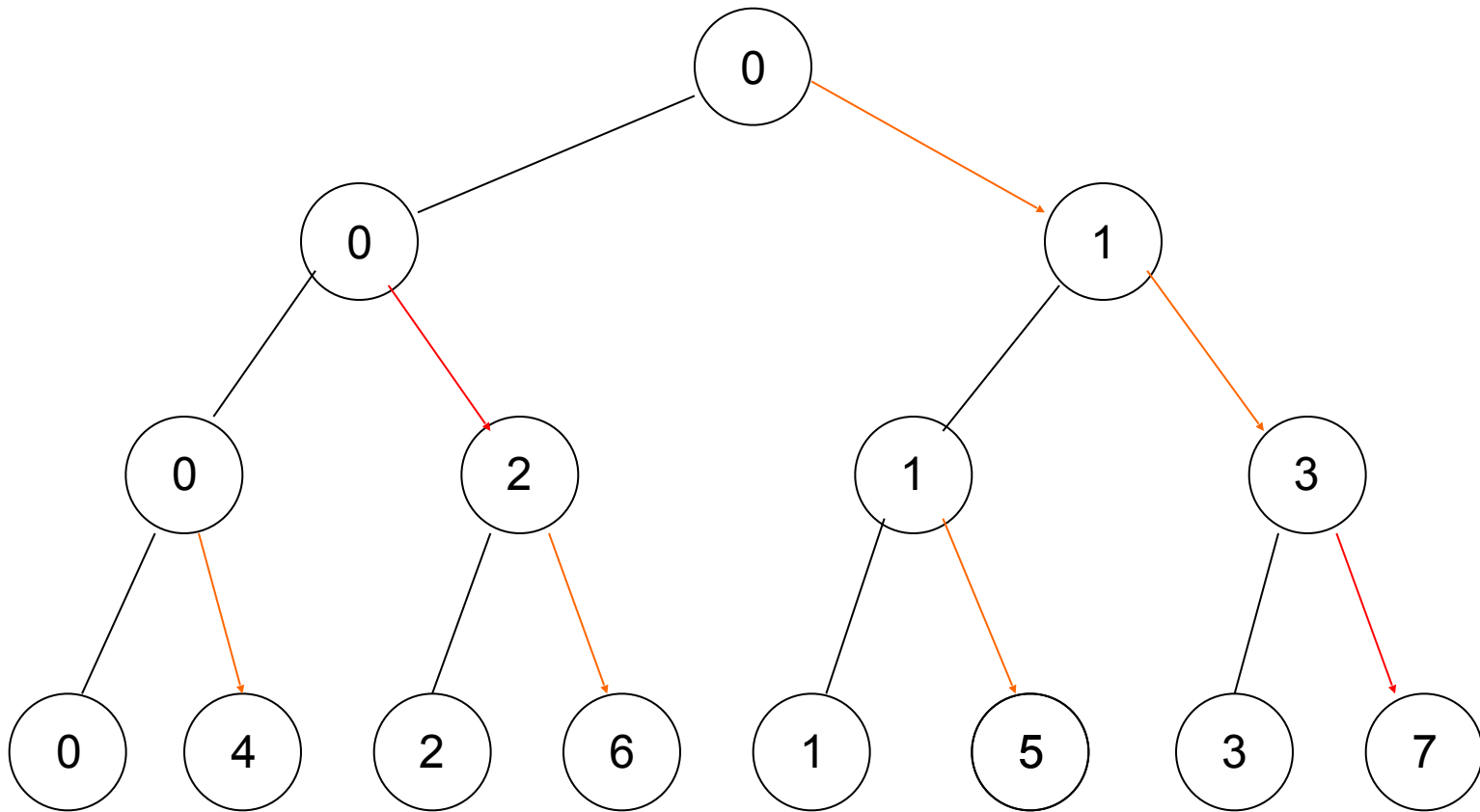
# Processed Configured as a Tree

---





# Processed Configured as a Tree





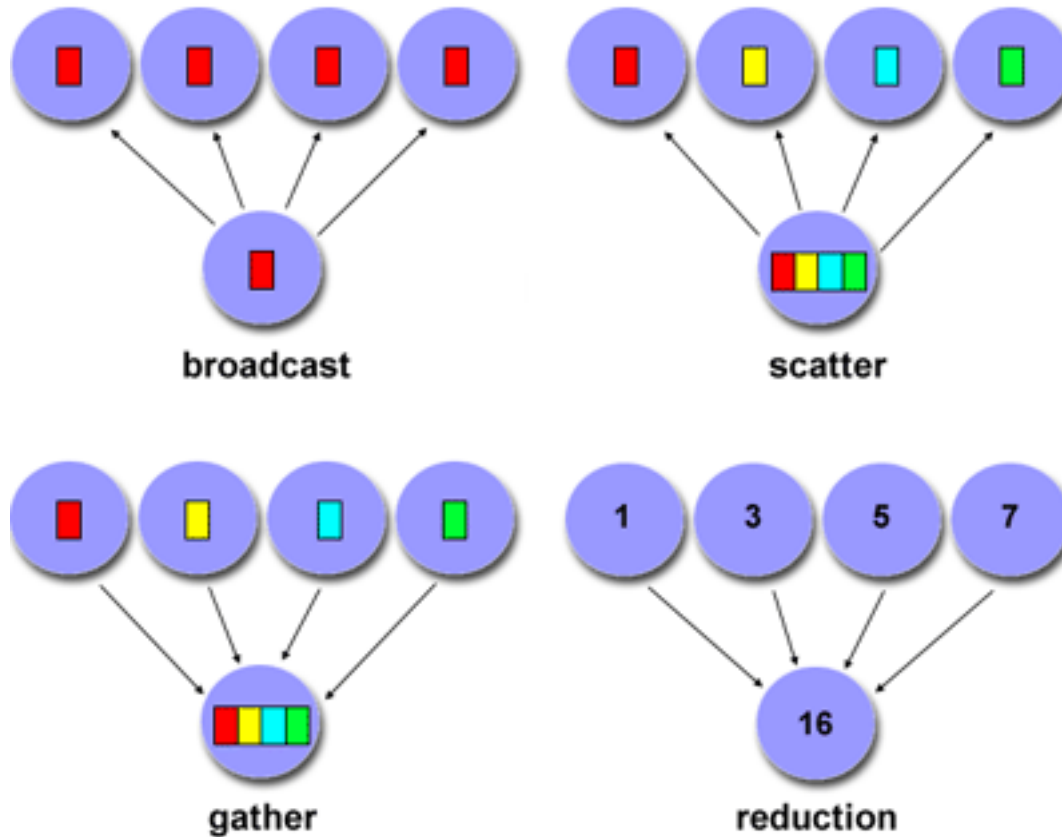
# Collective Communication

---

A communication pattern that involves all the processes in a communicator is a Collective Communication.



# Interactions





# Broadcast

---

**A Broadcast is a collective communication in which a single process sends data to every process in communicator.**



# Reducing

---

**Reducing is the act of receiving data from all processes onto one process and performing a simple action on it to get a final result.**



# Broadcasting and Reducing in MPI

---

**These are very common MPI operations and they have special functions**

MPI\_Bcast

MPI\_Reduce

**These can simplify things a lot!**





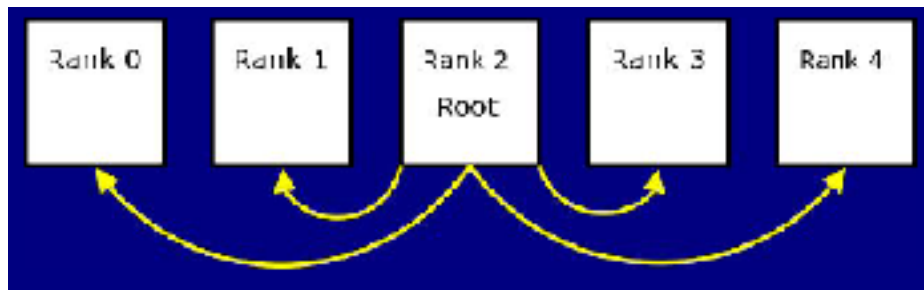
# MPI\_Bcast

MPI\_Bcast sends data from *root* process to all others

On root rank, *message* is send location

On other ranks, *message* is receive location

```
int MPI_Bcast(  
    int count,                /* in */  
    void *message,            /* in/out */  
    MPI_Datatype datatype,     /* in */  
    int root,                  /* in */  
    MPI_Comm comm              /* in */) 
```

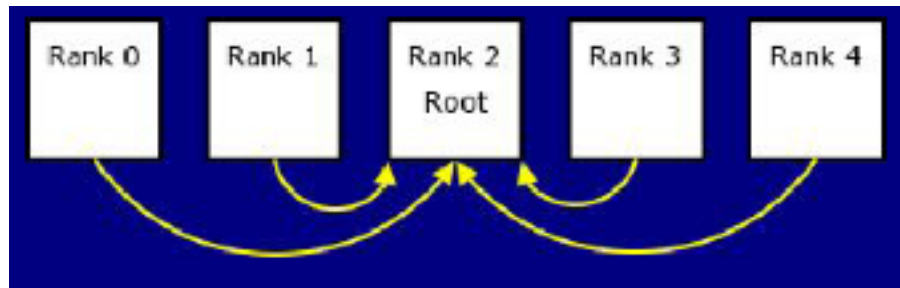




# MPI\_Reduce

MPI\_Reduce combines the operands stored in the memory referenced by *operand* using operation *operator* and stores the result in *\*result* on process *root*.

```
int MPI_Reduce(  
    void*      operand,    /* in */  
    void*      result,     /* out */  
    int        count,      /* in */  
    MPI_Datatype datatype, /* in */  
    MPI_Op      operator, /* in */  
    int        root,       /* in */  
    MPI_Comm    comm       /* in */) 
```





# MPI\_Reduce

---

Both *operand* and *result* refer to *count* memory locations with type *datatype*.

MPI\_Reduce must be called by all processes in the communication *comm*

*count*, *datatype*, *operator*, and *root* must be the same on each process



# MPI\_Reduce

---

*operator* can take on one of the followed predefined values

<i>Operation name</i>	<i>Meaning</i>
<b>MPI_MAX</b>	<b>Maximum</b>
<b>MPI_MIN</b>	<b>Minimum</b>
<b>MPI_SUM</b>	<b>Sum</b>
<b>MPI_PROD</b>	<b>Product</b>
<b>MPI_LAND</b>	<b>Logical and</b>
<b>MPI_BAND</b>	<b>Bitwise and</b>
<b>MPI_LOR</b>	<b>Logical or</b>
<b>MPI_BOR</b>	<b>Bitwise or</b>
<b>MPI_LXOR</b>	<b>Logical exclusive or</b>
<b>MPI_BXOR</b>	<b>Bitwise exclusive or</b>
<b>MPI_MAXLOC</b>	<b>Maximum and location of maximum</b>
<b>MPI_MINLOC</b>	<b>Minimum and location of minimum</b>



# MPI\_Reduce

---

Example for Trapezoid Rule program

```
MPI_Reduce(&integral, &total, 1, MPI_DOUBLE,  
           MPI_SUM, 0, MPI_COMM_WORLD);
```



# Sources of Overheads in Parallel Programs

---



# Sources of Overheads in Parallel Programs

---

1. **Inter-process interactions:** Processors working on any non-trivial parallel problem will need to talk to each other.



# Sources of Overheads in Parallel Programs

---

1. **Inter-process interactions:** Processors working on any non-trivial parallel problem will need to talk to each other.
2. **Idling:** Processes may idle because of load imbalance, synchronization, or serial components.





# Sources of Overheads in Parallel Programs

---

1. **Inter-process interactions:** Processors working on any non-trivial parallel problem will need to talk to each other.
2. **Idling:** Processes may idle because of load imbalance, synchronization, or serial components.
3. **Excess Computation:** This is computation not performed by the serial version. This might be because the serial algorithm is difficult to parallelize, or that some computations are repeated across processors to minimize communication.



# Synchronization in MPI

---

```
int MPI_Barrier( MPI_Comm comm /* in */) 
```

This function causes each process in comm to block until every process in comm has called it.



# Timing in MPI

---

*double MPI\_Wtime(void)*

Returns an elapsed time on the calling processor in seconds since an arbitrary time in the past. If a process is interrupted by the system, the time it spends idle will be added into the elapsed time.

**This is a function, declared as DOUBLE PRECISION MPI\_WTIME() in Fortran.**