

Date \_\_\_\_\_

Day \_\_\_\_\_

## "PARALLEL SOFTWARE"

### ► EAVESDAYS :-

#### Single Program Multiple Data :-

↳ Same program working on different parts of data.

```
    { if (I am thread / process 0)
        Operate on first half of array;
    else
        Operate on second half;
```

```
    { if (I am thread / process 0)
        do this ;
    else
        do this ;
```

→ Divide the work among process / threads in such a way that each process / thread gets roughly same amount of work, is called load balancing.

→ Programs that can be used by simply dividing the work among process / threads are called embarrassingly parallel.

Date \_\_\_\_\_

Day \_\_\_\_\_

→ SHARED MEMORY &

→ Communication is done thru Shared & private variable

→ Dynamic & Static Threads:

↳ In Dynamic threads program, there is a master thread and at any given instant a (possible empty) collection of worker threads.

When a request arrive for a work to be done master forks the workers thread, which after completion of task joins the master thread.

↳ In Static Thread paradigm, all of the threads are forked and all the threads run until all the work is completed. After completing work they join master thread and then they are freed.

If a thread is idle, its resource can't be freed.

→ Non determinism:

↳ A computation is non-deterministic if an input can result in diff outputs.

`printf("Thread %d, my_val = %d\n", my_rank, my_val);`  
The order in which threads execute is diff in each run.

But in shared memory System this problem can lead to race condition

Page No.

Date \_\_\_\_\_

Day \_\_\_\_\_

↳ To avoid threads to simultaneously access critical section we use mutual exclusion lock or mutex or lock.

↳ The use of mutex enforces serialization of critical section.

↳ Buggy waiting is an alternative to the mutex. It wastes system resources (cpu cycles).

```
my_val = Compute_val (my_rank);
if (my_rank == 1)
    while (!OK_for_1);
```

```
n += my_val;
if (my_rank == 0)
    OK_for_1 = true;
```

↳ Semaphores are similar to mutexes

↳ A monitor provides mutual exclusion at higher level. It is an object whose method can only be executed by one thread at a time.

↳ Transactional Memory: In DBMS, a transaction is an access to database that the system treats as a single unit. It either completes successfully or if it encounters any error during transaction it rolls back erasing all the data.

Page No.

Date \_\_\_\_\_

Day \_\_\_\_\_

## → Thread Safety :

Example: The C library function strtok splits an input string into subsequent strings. When it's first called it's passed a string, and on subsequent calls it returns successive substrings. If thread 0 makes its 1st call to strtok, and then thread 1 makes its 1st call to it before thread 0 has completed splitting its string, then t0's string will be lost or overwritten, and on subsequent calls it may get substring of t1's strings. So strtok is not thread safe.

## ▷ DISTRIBUTED - MEMORY :

## → Message Passing :

```
char message[100];
```

```
my_rank = Get_rank();
if (my_rank == 1) {
    Sprintf(message, "Greetings from process 1");
    Send(message, MSG_CHAR, 100, 0);
}
else if (locrank == 0) {
    Recieve(message, MSG_CHAR, 100, 1);
    printf("Process 0 Received %s\n", message);
}
```

Type of element      ↑  
 no. of rank of  
 elements receiving  
 ↓                      process

Date \_\_\_\_\_

Day \_\_\_\_\_

## ▷ INPUT AND OUTPUT :

→ In DMS programs, only process 0 will access stdin. In SMS programs, only master thread or thread 0 will access stdin.

→ In both DMS & SMS programs, all the processes/threads can access stdout & stderr.

→ Only a single process/thread will attempt to access any single file other than stdin, stdout or stderr. For eg: each process can open its own private file for reading or writing, but no two processes/thread will open same file.

→ Debug output should always include rank or id of process/thread.

## ▷ PERFORMANCE - Speedups &amp; Efficiency:

$$\text{Speedup} = \frac{T_{\text{serial}}}{T_{\text{parallel}}} = S$$

Let  $P = \text{no. of cores}$   
 Then,  $T_{\text{parallel}} = \frac{T_{\text{serial}}}{P}$

Then we say our program has linear speedup.

Page No.

Date \_\_\_\_\_ Day \_\_\_\_\_

$$\text{Q for linear speedup, } S = \frac{T_{\text{serial}}}{T_{\text{parallel}}} = P$$

The value  $\frac{S}{P}$  is called efficiency

$$E = \frac{S}{P} = \frac{T_{\text{serial}}}{P - T_{\text{parallel}}}$$

$E/P$  will get smaller & smaller as  $P$  increases

If we include the necessary "parallel overhead" cause due to parallelizing, because of communication and mutual exclusion then,

$$T_{\text{parallel}} = T_{\text{serial}} / P + T_{\text{overhead}}$$

### " Amdahl's Law "

If we are able to utilize 90% of a serial program.

If  $T_{\text{serial}} = 20 \text{ sec}$ , Then,

$$T_{\text{parallel}} = 0.9 \times \frac{T_{\text{serial}}}{P} = \frac{18}{P}$$

& For unparallel part,  $0.1 \times T_{\text{serial}} = 2$

$$T_{\text{parallel}} = \frac{0.9 \times 18 + 2}{P}$$

Page No.

Formula:  $S(N) = \frac{1}{\alpha + \frac{(1-\alpha)}{N}}$

Where,  
 $\alpha$  = serial fraction  
 $N$  = no. of processor

Date \_\_\_\_\_ Day \_\_\_\_\_

$$\text{Speedup} = S = \frac{T_{\text{serial}}}{\frac{18}{P} + 2} = \frac{20}{\frac{18}{P} + 2}$$

Now if 'P' gets larger than the value  $18/P$  comes closer to zero. This means how regardless of how much we incres.  $P$  (no. of cores) the speedup will remain less than  $20/2$  i.e. 10. Even if we use 1000 cores we can't get speedup better than 10.

It means if we have unparallelized part  $\gamma$ , then we can't get speedup better than  $\gamma$ . Here  $\gamma = (1 - 0.9) = 1/10$

### GUSTAFSON'S LAW :

This law considers the workload's scalability and disregards the idea of having a fixed portion of the workload.

As the computing resources (typically the no. of processors) increase, the problem size & workload can be scaled up to utilize the available resources effectively.

→ Adv. Ldisadv from PDF. Week #05 & 06

Formula:  $S = 1 + (N-1) \times P$

$S = N + (1-N) \times S$  As we incr. the no. of processor ( $N$ ) the workload or problem size also increases.

$N$  = No. of processor Page No.

$S$  = fraction of time for Serial part

Date \_\_\_\_\_

Day \_\_\_\_\_

→ Strongly Scalable : When we incr the no. of process we can keep efficiency fixed without incr the problem size , the program is said to be strongly scalable.

→ If we can keep the efficiency fixed by incr the problem size at the same rate as we incr the no. of process , then the program is said to be weakly scalable.

Page No. 

Date \_\_\_\_\_

Day \_\_\_\_\_

**CHAPTER # 03****Message Passing Interface "**

→ Command to run MPI programs

`mpiexec -n <no of process> ./mpi-hello`

`mpiexec -n (4) ./mpi-hello`

(This program will be executed by 4 threads)

▷ MPI\_Init & MPI\_Finalize :

`MPI_Init` initializes the setup of message passing like allocating storage for message buffers and might decide which process gets which rank.

`MPI_Init (&argc, &argv)`

optional parameters

`MPI_Finalize();`

(tells that we are done with MPI)

▷ Communicator, MPI\_Comm\_Size And MPI\_Comm\_Rank :

`MPI Communicator` is a collection of process that can send messages to each other. This communicator is called `MPI_COMM_WORLD`. Initialize by `MPI_Init()`:

```
int MPI_Comm_Size (MPI_COMM_WORLD, no_of_processes);
```

```
int MPI_Comm_Rank (           ,           , curr_process_rank);
```

Page No.

Date \_\_\_\_\_ Day \_\_\_\_\_

Date \_\_\_\_\_

→ We use the approach, single program multiple data (SPMD), in which we use same program and create diff branches acc to rank.

#### ▷ MPI\_Send():

`MPI_Send(msg_buffer, msg_size, msg_type, dest, tag, communicator);`

tag → tag is a non-negative int. It is used to distinguish messages. Suppose P1 is sending floats to P0. Some of floats should be printed while others should be used in computation. So process 1 can use tag of 0 for msgs that should be printed & tag 1 for computational msgs.

Two process having diff communicator can not send & rec msgs.

#### ▷ MPI\_Recv():

If received msg is larger than buffer then an overflow err will occur & return MPI\_ERR\_TRUNCATE.

`MPI_Recv(msg_buffer, msg_size, msg_type, source, tag, communicator, status_p);`

Q: This can happen sometimes that we don't know in which order the msgs will be received. So instead

Page No. Page No. 

Date \_\_\_\_\_ Day \_\_\_\_\_

Date \_\_\_\_\_ Day \_\_\_\_\_

of waiting source we can use MPI\_ANY\_SOURCE same for tag MPI\_ANY\_TAG.

↳ But these are not valid for MPI\_Send();

#### ▷ status-p Argument:

→ MPI\_Status is a struct with at least three members MPI\_SOURCE, MPI\_TAG and MPI\_ERROR.

→ We have to make a variable MPI\_Status status.

→ This is passed as &status as last argument of MPI\_Recv.

→ After that we can access the source thru status.MPI\_SOURCE and tag thru status.MPI\_TAG.

→ To get the amount of data received we have to pass this status variable in

`MPI_Get_Count(&status, rec_type, &count);`

↳ Now the count variable has the amount of data received.

Page No. Page No.

Date \_\_\_\_\_

Day \_\_\_\_\_

### ▷ Structure of Message-Passing Programming :

#### ① Asynchronous Paradigm :

All concurrent tasks execute asynchronously. This makes it possible to implement any parallel algo. These types of programs can have non-deterministic behaviors due to race condition.

#### ② Loosely Synchronous :

In such programs, tasks or subset of tasks synchronize to perform interactions.

#### ▷ Direct Memory Access :

DMA allows copying of data from one memory location to another (e.g.: communication buffer) without CPU support.

#### ▷ Blocking Non-Buffered Send & Receive :

##### ↳ Blocking Non-Buffered Send :

- The operation is considered blocking when the sending process waits until the message has been fully transmitted to the receiver or until the receiver is ready to receive the message.
- Non-buffered means the sys doesn't use intermediate

Page No. 

Date \_\_\_\_\_

Day \_\_\_\_\_

memory (buffers) to temporarily store the message

- The send oper only completes when the receiver has started to receive the data, ensuring that the msg is immediately consumed.

##### ↳ Blocking Non-Buffered Receive :

- The receive operation is considered blocking when the receiving process waits until the msg arrives.
- It will only proceed once the sender has transmitted the data & receiver has successfully obtained it.
- Non-buffered means the system does not allocate a temporary buffer; the receiver waits for the data to be sent directly.

- Both the sending & receiving process must be ready because if one isn't ready the other would be in idle (waiting for the other process).

##### ↳ Deadlocks in Blocking Non-Buffered Operations

P<sub>0</sub>

Send (&amp;a, 1, 1);

Receive (&amp;b, 1, 1);

P<sub>1</sub>

{Send (&amp;a, 1, 0);

Receive (&amp;b, 1, 0);

Done ek doosre ke data send karne hai or done wait karne hai k data receive ho doosre process ki tar se.

Page No.

22K-4187 Samiye Ali

## ASSIGNMENT #02

Date \_\_\_\_\_

Day \_\_\_\_\_

This results in indefinite waiting, which results in deadlock.

### QUESTION #01 (a)

#### ► Blocking Buffered Send / Receive :

##### ↳ Blocking Buffered Send :

- The sending process uses a buffer to temporarily store the message, allowing the sender to proceed with waiting for the receiver to be ready to accept the msg.
- Once the msg is copied to buffer, the sender can continue with its execution.
- The process is only blocking until the data is copied to buffer.

##### ↳ Blocking Buffered Receive :

- The receiver waits for the message to be available, blocking until the msg arrives.
- The msg is transferred from the buffer to receiver's memory when the receiver calls the receive operation.

`MPI_Buffer_attach(buffer, BUFFER_SIZE);`

`if (rank == 0) {`

`data = 100;`

`MPI_BSend(&data, 1, MPI_INT, 1, 0, MPI_COMM_WORLD);`

`printf("Process 0 sent data %d", data);`

Page No.

Date \_\_\_\_\_

Day \_\_\_\_\_

`else if (rank == 1) {`

`MPI_Recv(&data, 1, MPI_INT, 0, 0, MPI_COMM_WORLD, MPI_STATUS_IGNORE);`

`printf("Process 1 received data %d\n", data);`

`}`

`MPI_Buffer_detach(buffer, &BUFFER_SIZE);`

`MPI_Finalize();`

9

- Buffer protocols alleviate idling overheads at the cost of adding buffer management.
- If a program is highly sync, then non-buffered may perform better.

#### ↳ Impact of Finite buffers in Message Passing :

P<sub>0</sub>

P<sub>1</sub>

`for (i=0; i<1000; i++) {`

`foo(i);`

`produce_data(&a);`

`receive(&a, 1, 0);`

`send(&a, 1, 1);`

`consume_data(&a);`

`};`

Q: P<sub>1</sub> was slow to getting loop, the P<sub>0</sub> might have sent all of its data - If there is enough buffer space

Page No.

Date \_\_\_\_\_

Day \_\_\_\_\_

space, then both process can proceed; however if the buffer is not sufficient (i.e. buffer overflow), the sender would have to be blocked until some of the corresponding receive operations had been posted, thus freeing up buffer space.

#### ↳ Deadlock in Buffered Operation:

P <sub>0</sub>	P <sub>1</sub>
rec(&a, 1, 1)	rec(&a, 1, 0);
send(&b, 1, 1)	send(&b, 1, 0);

#### QUESTION # 2 (a)

#### ► Non-Blocking Message Passing :

##### ↳ Non-blocking Send: (Non-buffered)

- The sending process initiates the transmission of a msg & then immediately proceeds with its computation. The send oper does not wait for the receiver to be ready.
- The data might not have been completely transferred when the send operation returns, so the user must ensure that the data buffer is not modified until the communication completes.

Page No. 

Date \_\_\_\_\_

Day \_\_\_\_\_

↳ Non-blocking receive:

- The receiving process initiates a request to receive data but continues its exec without waiting for data to arrive.
- The process can check whether the msg has been received at a later point or can proceed with other work.

#### MPI\_Request request;

```
if (rank == 0) {
    data = 100;
```

```
MPI_Isend (&data, 1, MPI_INT, 1, 0, MPI_COMM_WORLD,
&request);
```

```
printf ("Process 0 initiated send \n");
```

```
// Do some other work
```

```
// Encuse send oper completes
```

```
MPI_Wait (&request, &status); // P0 completed
```

```
} else if (rank == 1) {
```

```
MPI_Irecv (&data, 1, MPI_INT, 0, 0, MPI_COMM_WORLD,
&request); // P1 initiated
```

```
// Do some other work
```

```
MPI_Wait (&request, &status); // P1 completed
```

Page No.

Date \_\_\_\_\_

Day \_\_\_\_\_

↳ Buffered :

The msg is temp stored in buffer, allowing the sending process to continue immediately.

" MPI - LIBRARY "

// To get name of the processor.

```
char processor_name[MPI_MAX_PROCESSOR_NAME];
int name_len;
```

MPI\_Get\_Processor\_name(processor\_name, &name\_len);

// MPI\_Send Recv Simultaneously.

```
int a[10], b[10], npes, myrank;
MPI_Status status;
```

MPI\_Comm\_Size(MPI\_COMM\_WORLD, &npes);

MPI\_Comm\_Rank( , &myrank);

```
MPI_SendRecv(a, 10, MPI_INT, (myrank+1)%npes, 1,
             b, 10, MPI_INT, (myrank-1+npes)%npes,
             MPI_COMM_WORLD, &status);
```

↳ avoids circular deadlock problem

Page No. 

Date \_\_\_\_\_

Day \_\_\_\_\_

// MPI Send & Recieve for with Single buffer

```
MPI_Sendrecv_replace(void *buff, int count,
                     MPI_Datatype datatype, int dest,
                     MPI_Status send_tag, int source,
                     MPI_Status recv_tag, MPI_COMM_WORLD,
                     &status);
```

Received data replaces the data that was sent out of the buffer.

// Checking whether a non-blocking oper completed?

MPI\_Test(MPI\_Request \*req, int \*flag, &status);  
 Git tests whether or not the operation send/recv identified by its req has finished. And returns true in flag otherwise false.

If non-blocking operation has finished, the request object is deallocated & req is set to MPI\_Request\_NULL

// For explicitly deallocating req object.

MPI\_Request\_free(&request);

Page No.

Date \_\_\_\_\_

Day \_\_\_\_\_

- MPI\_BARRIER() enforces sync among all processes in a communication world by making them wait until every process reaches the barrier, ensuring coordination at specific program points.

`MPI_BARRIER(MPI_COMM_WORLD);`

- MPI\_BCAST() → A broadcast func is one of the std collective communication technique. One process sends the same data to all process in a communicator.

The process designated as "root" sends the msg.

`MPI_BCAST(&data, 1, MPI_INT, 0, MPI_COMM_WORLD)`

is data available means bini  
data boga wo subke mil jaga bank of process that  
Sends the data

↳ Useful when one process has data that needs to be distributed (shared) with all other processes, such as input parameters or configuration parameters.

- MPI\_SCATTER() → Similar to broadcast but used to distribute different chunks of data from root process to all other process in a communicator, including itself.

Page No.

Page No.

Date \_\_\_\_\_

Day \_\_\_\_\_

`MPI_Scatter(&sendbuffer, &int sendCount, MPI_Datatype send datatype, void *recvbuffer, int recvCount, MPI_Datatype recvType, int root, MPI_Comm MPI_COMM_WORLD)`

- if we have an array → {10, 20, 30, 40, 50, 60}'s then root P0 will receive 10, P1 receive 20, P2 30, P3 40 & so on. (P0 sendCount = 1)
- if sendCount > 2 each process receives 2 elements.
- Usually sendCount = no of element / no of process.

- MPI\_GATHER() → It is inverse of scatter. It takes elements from many process and gathers them to one single process.

It takes elements from each process & gathers them to the root process.

Same Syntax as scatter.

↳ These only root process needs to have a valid buffer receive buffer, other calling process can pass NULL for rec-data.

↳ Here recvCount is the count of elements received per process, not the total summation of counts from all process.

Day

Date

- MPI\_Allgather() → Use for many-to-many communication  
Given set of elements distributed across all processes  
MPI Allgather will gather all of the elements to all the processes.

`MPI_Allgather(void *sendbuf, int sendCount,  
MPI_Datatype sendType, void *recvbuf,  
int recvCount, MPI_Datatype recvType,  
MPI_Comm world);`

- MPI\_Reduce() → Works similar to gather, the only difference is that it reduces the result to the root process by doing some operation on data.

`MPI_Reduce(void *send_data, void *recv_data, int count,  
MPI_Datatype datatype, MPI_Op op,  
int root, MPI_Comm world);`

Some predefined operation can be applied.

MPI\_MAX

MPI\_MIN

MPI\_SUM

MPI\_PROD

MPI\_BAND (Logical AND)

MPI\_LOR

MPI\_BAND (bitwise AND)

MPI\_BOR

MPI\_EXOR

MPI\_BXOR

MPI\_MAXLOC (return max value & rank of process that owns it)

Page No.  ⑥ 18

• MPI Timer

Use to calculate time taken by MPI Programs

double d1

d1 = MPI\_

!

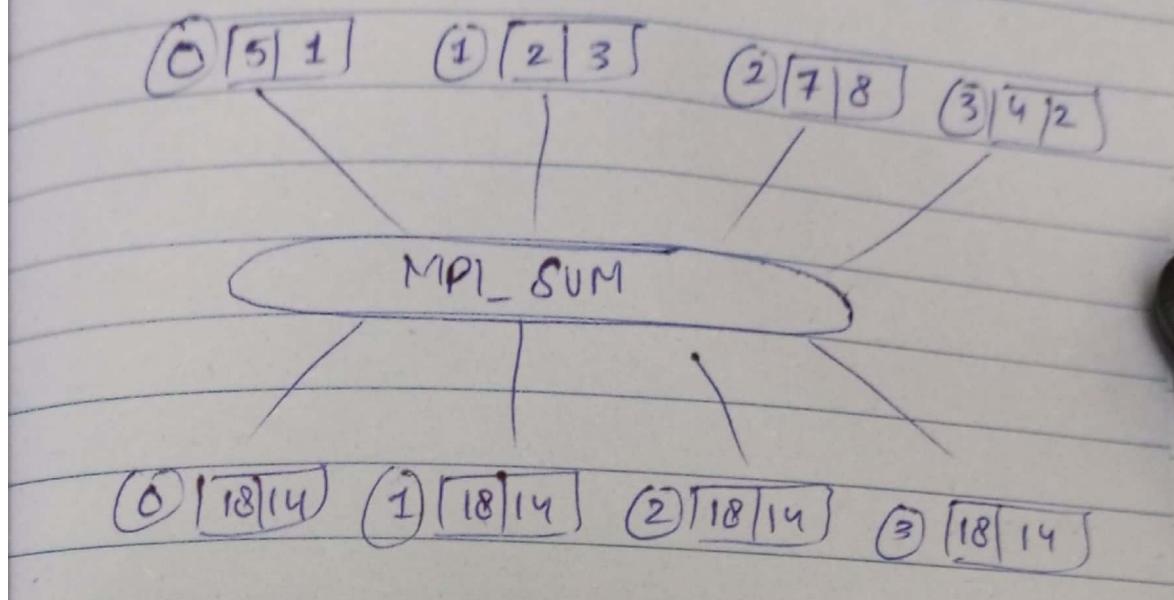
d2 = MPI\_

printf("%T,

Date

Day

• MPI\_Allreduce :



MPI\_Allreduce (void \*sendbuf, void \*recvbuf, int count,  
MPI\_Datatype datatype, MPI\_Op op,  
MPI\_Comm world);

• MPI\_Timers

Use to calculate elapsed time b/w 2 points in an MPI Program.

double d1, d2;

d1 = MPI\_Wtime();

{

d2 = MPI\_Wtime();

printf ("Time : %f", d2 - d1);