

→ Lec 2

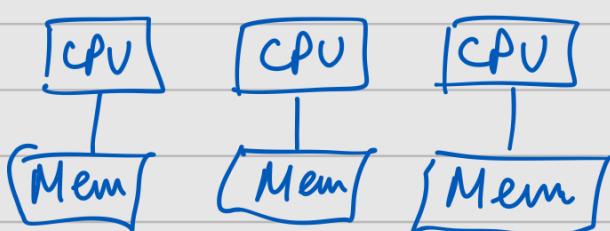
- Speedup = $\frac{\text{Time taken on 1 processor}}{\text{Time taken on } P \text{ processors}}$
- Efficiency = $\frac{\text{Speedup}}{P}$

→ Amdahl's Law

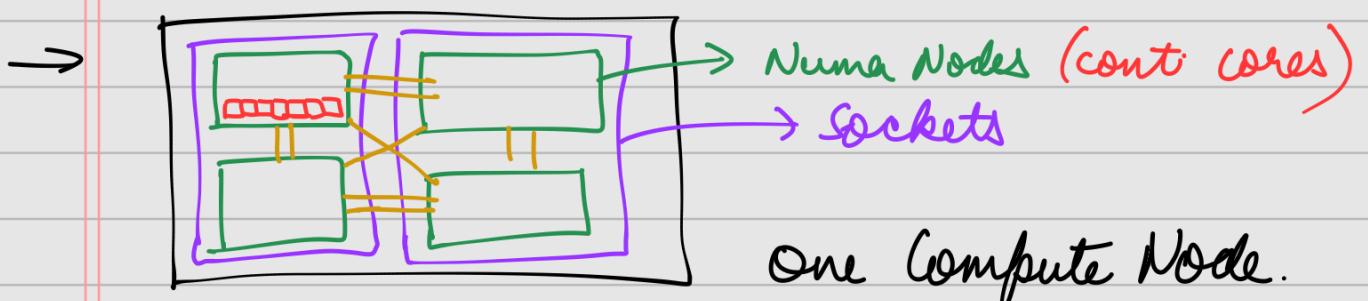
$$\text{Speedup } S = \frac{1}{(1-f) + f/P}$$

(f = frac of code ||sable)

→ NUMA



(Non uniform mem access)



One Compute Node.

→ MPI

Standard for message passing.
Distinct/Distributed memory environment
Explicit communication

→ Process - instance of prog execution

→ Execution context of a running prog - resources associated with a prog's exec"

→ PPN = process per node.

→ Round Robin

hostfile does not mention # processes on any node.

0 1 2 3
4 5 6 7

→ Sequential / Contiguous

hostfile mentions how many processes on any node.

0 2 4 6
1 3 5 7

- MPI_Comm_Size (MPI_COMM_WORLD, & size)

- MPI_Comm_Rank (MPI_COMM_WORLD, & rank)
global communicator

→ give values w/o communicator.

→ Process

belongs to a group.

identified by rank within the group.

→ Communicator

defines scope, specifies comm. context

→ MPI_Get_processor_name (hostname, & len)
: host name (cnews 3)

→ core Id = sched_getcpu()

: which core on that host that process.

→

Comm

P-2-P

client

→ Data types: BYTE, CHAR, INT, FLOAT, DOUBLE

→ MPI Send:

(buf, count, datatype, dest, tag, comm)

→ MPI Send:

(buf, count, datatype, src, tag, comm, states)
 $\text{size (B)} = \text{count} \times \text{sizeof (datatype)}$

→ Errors in Send - Recv

1) Recv → cts Send → cts
 $\text{cts} < \text{cts}$: runtime error

2) tag mismatch: no compile-time or run-time error. (Blocked)

- Strlen() : size except 10.
 Just string size.

3) MPI Send: blocking, standard mode.
 can transfer it to MPI buffer & move on
 • does not return till send buffer can be reused. message buffering affects this.

4) Eager v/s Rendezvous

Eager : buffering enabled.

— completes without ack from dest.
 small msgs, typically 128 KB.

Rendezvous

needs ack from matching receive.
large msgs.

→ MPI Status

MPI_SOURCE: source rank

MPI_TAG: message tag

count → **MPI_Get_Count**: msg len

* we take max of send & recv times.

** intra faster than inter

* distant node ⇒ time inc

→ Sources of variation

- network congestion / background traffic
- OS noise
- other appl.

* **MPI_Reduce**

applies a func to all values in arg0 of all ranks & stores result in rank 0 variable arg1.

* Need for diff recv [num ranks][data size]

only needed when many-to-one. Not needed for one-to-many.

* **MPI_Send**: blocking send
MPI_Recv: blocking recv

does not return till the msg data and envelope have been safely stored so that sender is free to modify the buffer

* MPI_ANY_SOURCE : recv specify wildcard for
 MPI_ANY_TAG : recv wildcard for tag.
 → only for recv

→ Many to One

recvarr [size][20]; buf [20];

MPI_Recv (buf, 20, __, MPI_ANY_SOURCE,
 MPI_ANY_TAG, __)

(copy) recvarr [status. MPI_SOURCE] = buf;

(we get status obj after Recv is done).

→ Over decomposition :

processes ↑ but time rem same.

→ Timing factors

(Same # processes, data size ↑ ⇒
 computation time ↑, comm time ↑)
 Same data size, # processes ↑ ⇒

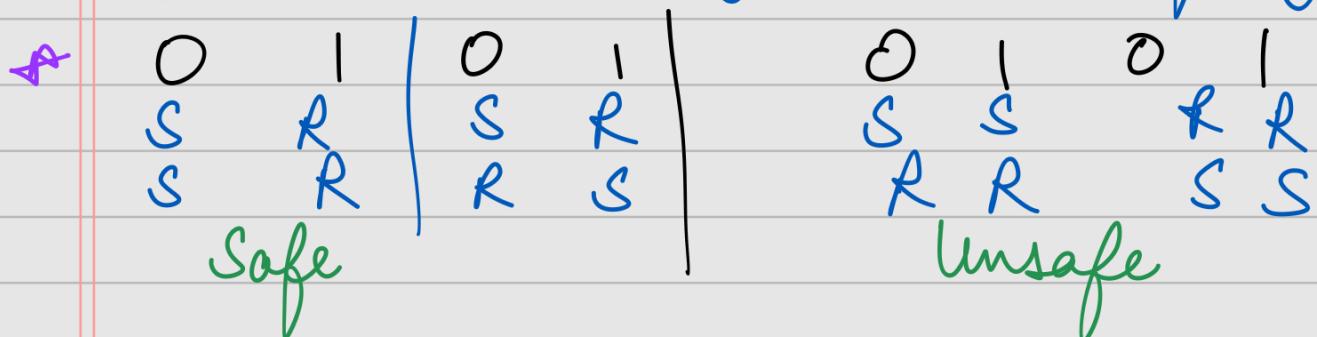
- i) Data size ↑ ⇒ comm time inc
- ii) Compute time scales almost.
- iii) No. of processes ↑ ⇒ comm. time inc

* Send (buf, count, INT, __)
 Recv (buffer, count, CHAR, __)

Runtime error: sent count × 4 bytes
 recv space allowed = count bytes

if msg size is smaller than recv count,
then only those locations corresponding
to message are modified

- * MPI-Send: non blocking till one point,
after that blocking. (no buffering)



→ Non-overtaking Messages

$$0 < S \left(\begin{pmatrix} 1, 1 \\ 1, 1 \end{pmatrix} \right) - R \left(\begin{pmatrix} 0, 1, 1 \\ 0, 1, 1 \end{pmatrix} \right) > 1$$

if src, tag match \Rightarrow first send matches
first recv only.

→ Diff Sends

- ① Send : blocking . Standard

 - started whether or not matching `recv` posted
 - may complete before matching `recv` posted
 - MPI may buffer messages .

- ② Send : buffered

- may complete before matching rec posted
 - does not depend on occurrence of matching rec

③ Ssend

Synchronous

- Completes only if a matching `recv` posted.

④ Rsend

Ready



- Started only if a matching `recv` posted.

- `MPI_BUFFER_ATTACH` (`buffer_addr, size`)
 - buffer in user's memory
 - used by msgs sent in buffered mode
 - only one buffer can be attached to a process at a time.

`MPI_BUFFER_DETACH` (`buffer_addr, size`)

Size = sum of sizes of all outstanding
Bsend + MPI_BSEND_OVERHEAD
for each Bsend.

*

O
ss
Ss

I
R
R

O
Ss
R

I
R
Ss

O
Ss
R
R

Safe

Unsafe/Deadlock

→ 2D NN with odd even

if $\gamma \% 2 == 0$ and $\gamma < p-1$

$S \left(\frac{_}{_}, \frac{\gamma+1}{\gamma+1}, \frac{\gamma+1}{\gamma+1}, \frac{_}{_} \right)$

$R \left(\frac{_}{_}, \frac{\gamma+1}{\gamma+1}, \frac{\gamma}{\gamma}, \frac{_}{_} \right)$

else if $\gamma \% 2 == 1$ and $\gamma > 0$

$R \left(\frac{\gamma-1}{\gamma-1}, \frac{\gamma}{\gamma}, \frac{_}{_} \right)$

$S(\tau^{-1}, \tau^+, \dots)$

if $(\tau / 2 == 0 \text{ and } \tau > 0)$

$$\begin{cases} S(\tau^{-1}, \tau^{-1}, \dots) \\ R(\tau^{-1}, \tau, \dots) \end{cases}$$

else if $(\tau / 2 == 1 \text{ and } \tau < P-1)$

$$\begin{cases} R(\tau+1, \tau, \dots) \\ S(\tau+1, \tau+1, \dots) \end{cases}$$


- NN-2 takes almost half time than NN-1 (both 2D). if scaled when internode. Intranode \rightarrow both same time.
- NN time more determined by no. of internode comm happening.

$\rightarrow \underline{\text{Isend}} / \underline{\text{Irecv}} : \text{Non Blocking p-to-p}$

MPI_Irecv (buf, count, datatype, dest, tag, comm, request)

MPI_Irecv (buf, count, datatype, src, tag, comm, request)

MPI_Wait (request, status)

MPI_Waitall (count, request, status)

of outstanding Irecv + Irecv called by that process

* Intra node v/s Inter node - no effect on computation time. Only comm time ↑

* Small messages are latency bound.
large messages bandwidth bound.

* Ring algo - data size remains same for each successive communication
Recursive Doubling - data size comm doubles

also # hops increases in rec doubling
while in ring we always comm. with neighbour