

## Lee 9

### "Variability"

- MPI\_Isend (buf, count, datatype, dest, tag, comm, req)
- MPI\_Irecv (buf, count, datatype, src, tag, comm, req)
- MPI\_Wait (req, status)
- MPI\_Waitall (count, req, status)
- MPI\_Test : non blocking, tests completion, starts progress.

22:40

\* MPI\_R\_CVAR\_ASYNC\_PROGRESS = 1

→ for same isend / irecv code, time decreases

→ Comm Graph / Matrix

Symmetric matrix

$$A[i][j] = A[j][i]$$

= Total no. of bytes comm b/w ranks i and j  
irresp of "dir" of comm.

→ Process Mapping / Alloc<sup>n</sup>

Given network topology and a given comm pattern (virtual topology)

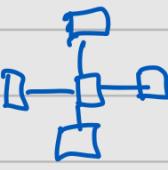
poor process mapping affects comm. time.

## Lec 10

### → Attributes of Interconnects

- Topology (connectivity) — structure of network
- Cost — no. of wires
- Diameter — max possible dist
- Bisection width — mincut

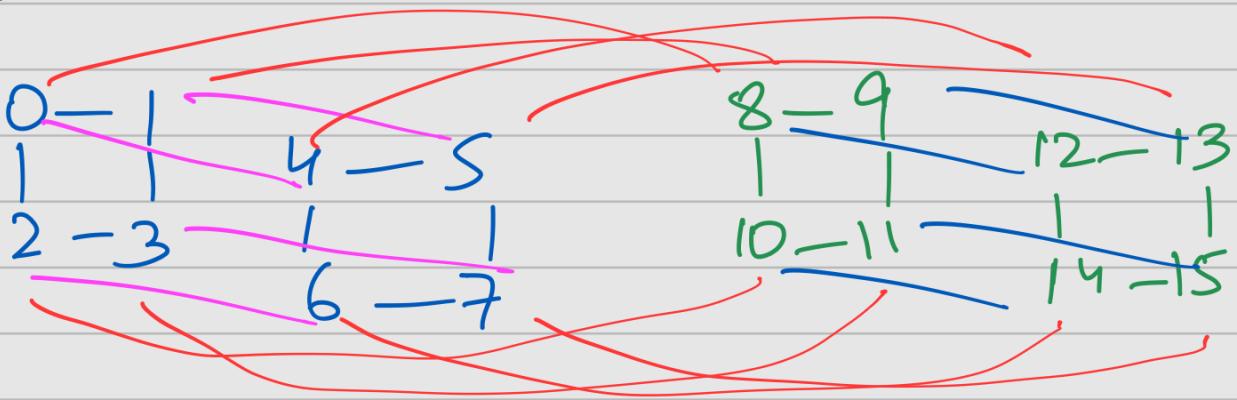
(# nodes =  $P$ )

| Topology   | Dia                                     | Cost                    | Bis Wid     |
|--|---|-------------------------|-------------|
|  Linear  | $P-1$                                   | $P-1$                   | 1           |
|  Ring  | $\text{floor}(P/2)$                     | $P$                     | 2           |
|  Star<br>(single pt of failure; cong. at central node) | 2                                       | $P-1$                   | —           |
| 2D Mesh  | $2(\sqrt{P}-1)$                         | $2\sqrt{P}(\sqrt{P}-1)$ | $5\sqrt{P}$ |
| 2D Torus   | $\sqrt{P}$ (even)<br>$\sqrt{P}-1$ (odd) | $2P$                    | $2\sqrt{P}$ |

$$3D \text{ Mesh} \quad 3(p^{y_{3-1}}) \quad 3(p^{2/3})(p^{y_{3-1}}) \quad p^{2/3}$$

## 3D Torus

Hypercube  $\log p$   $\frac{p \log p}{2}$   $p/2$



## Tree Networks

Prob: Comm bottleneck at higher levels

↳ Fat tree Networks

→ Comm. Cost

Bandwidth / Throughput =  $\frac{\# \text{ bits}}{\text{time}}$  bps

\* Latency bound region: Small msg size sat

\* Bandwidth bound region: Large msg size sat



$\rightarrow$  msg size

— for  $n$  words on  $l$  links

$$t_s + t_a(l) + \frac{n}{B}$$

cut thru

Startup time (at send/receive nodes)

Node latency (per hop time)

$B$  = bandwidth or  $t_w$ : per word transfer time

Store & fwd

$$t_s + t_a(l) + \frac{ln}{B}$$

## Hockney Model / Postal Model

$$T_{\text{comm}} = L + \frac{n}{B} \quad (L = t_s + l t_a)$$

- affect effective bandwidth  $\rightarrow$
- 1) Bulk comm
  - 2) fewer hops
  - 3) less volume comm
  - 4) network congestion
  - 5) network topology

Realistic model  $\Rightarrow$  overhead, contention

$$L' + \frac{n}{B'}$$

$\rightarrow$  MPI Groups

- ordered set of processes.
- contiguous ranks
- MPI\_Comm\_group (comm, \*group)
- MPI\_Group\_rank, MPI\_Group\_size
- Union & intersection of groups

### → MPI Communicators

- MPI object cont. a group of processes
- MPI\_COMM\_WORLD
- represents a comm. domain
- memory map to # processes in the group
- several commun. contexts may co-exist within a single communicator.

### MPI\_Comm\_Split

MPI\_Comm\_Split (oldcomm, color, key, \*newcomm)

↑  
inp

↑  
output

\* collective call — everyone in oldcomm has to do this.

\* color = division criteria

(Same col  $\Rightarrow$  same group)

- \* key - new rank assignment in inc order of key.
- MPI\_Group\_incl(g-group, N, ranks, & new-group)
- MPI\_Comm\_create\_group(MPI\_COMM\_WORLD, new-group, tag, & new-comm)
  - must be called by all processes in the comm,
  - tag must be same for all calls.

## Lec 12

### Collective Comm

- must be called by all processes that are part of the communicator

Barrier (sync)

Bcast      Reduce      Gather      Scatter

Allreduce      Allgather      Alltoall

- MPI\_BARRIER :

every rank needs to call this func  
and caller returns only after all processes  
have entered the call

- MPI - Gather:

Gathers values from all processes to root.

in comm, all processes ke sendbuf se sendcount elements hake root ke recvbuf me daal dega.

recvcount = size of any single receive

- Eq. P2P -

many to one

non-root send  $\rightarrow$  root recv

|    |  |  |
|----|--|--|
| A0 |  |  |
| A1 |  |  |
| A2 |  |  |

sendbuf



|    |    |    |
|----|----|----|
| A0 | A1 | A2 |
|    |    |    |
|    |    |    |

recvbuf

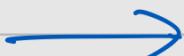
$$T_{\text{naive}} = (p-1) (L + \frac{n}{B})$$

$$T_{\text{opt}} = \log p L + \frac{n(p-1)}{B}$$

- MPI - Broadcast

comm ke sare processes ko root process buffer ke count elements bhejta, aur unke buffer me rakhega.

|   |   |   |   |
|---|---|---|---|
| X | Y | Z | W |
|   |   |   |   |
|   |   |   |   |

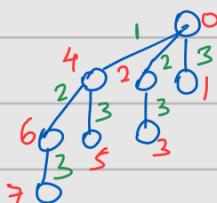


|   |   |   |   |
|---|---|---|---|
| X | Y | Z | W |
| X | Y | Z | W |
| X | Y | Z | W |

eq P2P — One to many. (Naive)  
 time =  $(P-1) \left( L + \frac{n}{B} \right)$  all 1 hop

optimised : "Binomial Algo" Hypercube  
 (dit half)

$$\text{Steps} = \log_2 P \quad T(P) = \log_2 P \left( L + \frac{n}{B} \right)$$



$P$  col = rank  $\times G/P$  or rank  $\% G$   
 Comm-Split  
 Comm-rank (news, &remain)  
 Comm-size()  
 int buf [2]  
 $P/G - 1$

```

if (news == news-1) {
    buf[0] = rank;
}
if (news == 0) {
    MPI_Send (&buf[1], news-1, news);
}
MPI_Bcast (buf, 2, INT, news-1, news);
    
```

### • MPI\_Scatter

sendcount — # elements sent to each process

root ke sendbuf se sendcount elements order wise comm ke ranks ke recvbuf me recvcount elements daalga.



naive — One to many.

$$T_{\text{naive}} = (P-1)L + \frac{n(P-1)}{B}$$

$$T_{\text{naive}} = (P-1)L + \frac{n(P-1)}{B}$$

$$T_{opt} = (\log_2 P)L + \left(\frac{n}{P}\right)\left(\frac{P-1}{B}\right)$$

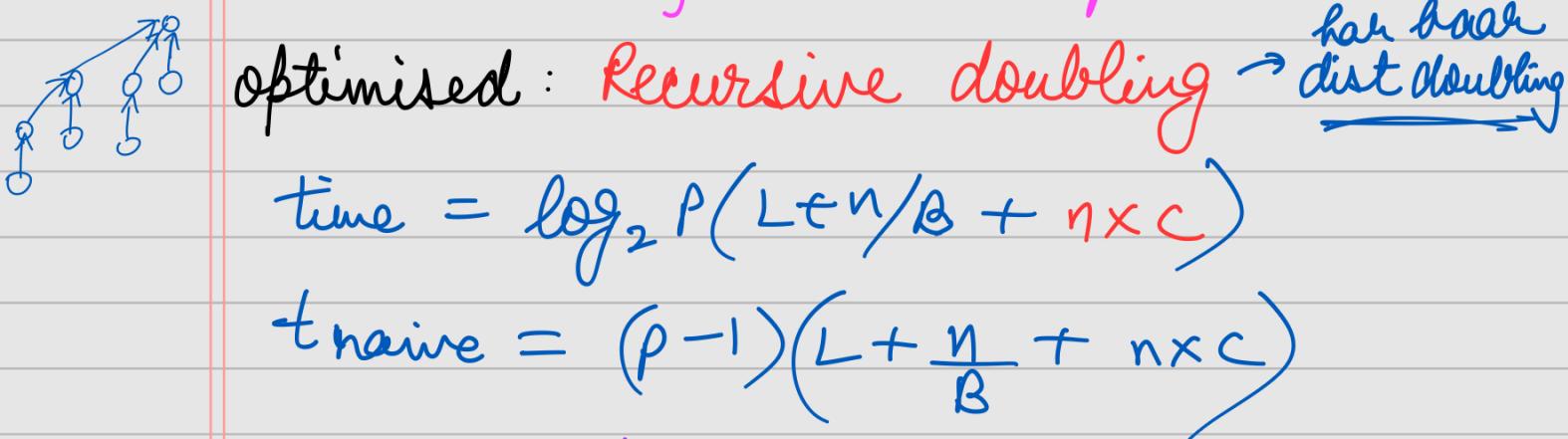
- MPI - Reduce

comm ke sare processes ke sendbuf se count elements ke operator lagake root ke outbuf me rakh dega.

MIN, Max, Sum, Prod

\* using Reduce instead of loop (manually Send - Recv) reduces time for same task.

Naive : many to one + computation



$$\text{time} = \log_2 P(L + n/B + nc)$$

$$+ t_{naive} = (P-1)\left(L + \frac{n}{B} + nc\right)$$

- MPI - Allgather

Sendbuf

|    |  |  |
|----|--|--|
| A0 |  |  |
| A1 |  |  |
| A2 |  |  |



|    |    |    |
|----|----|----|
| A0 | A1 | A2 |
| A0 | A1 | A2 |
| A0 | A1 | A2 |

recvbuf

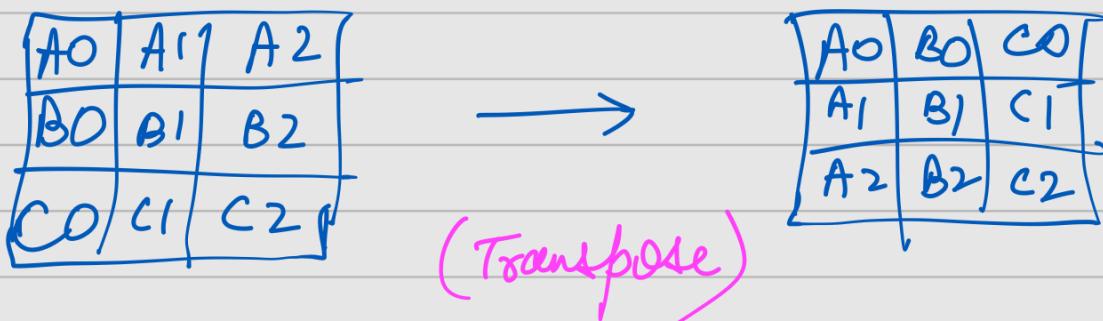
(gather + Broadcast)

comm ke sare processes ke sendbuf se sendcount elements behle ek ke recvbuf me rakhlo, fir uss recvbuf ko broadcast kro

- MPI - All reduce  
 $(\text{reduce} + \text{Bcast})$

- MPI - All to all

s b



each rank calls scatter  $\rightarrow & [\text{recvbuf} + \text{rank}]$   
each rank calls gather  
 $\hookrightarrow & [\text{sendbuf} + \text{rank}] \text{ se}$

- \* Binomial Algo for Broadcast

P processes

$K = \emptyset$

$(0 \rightarrow 15)$

$\log_2 P$   
steps

for ( $K = \emptyset$ ,  $K > 1$ ,  $K = K/2$ ) {

if ( $\text{rank} \% K == 0$ )

send ( , rank +  $K/2$ , rank )

else if ( $\text{rank} \% K == K/2$ )

recv ( , rank -  $K/2$ , rank -  $K/2$ , )

## → Domain Decomposition

$N$  grid points,  $P$  processes,  $\frac{N}{P}$  pts per process

|              | Comm                         | Comp          | Ratio                        |
|--------------|------------------------------|---------------|------------------------------|
| 1D domain    | 2                            | $N/P$         | $\frac{2P}{N}$               |
| 1D decompose |                              |               |                              |
| 2D domain    | $2\sqrt{N}$                  | $N/P$         | $\frac{2P}{\sqrt{N}}$        |
| 2D decompose |                              |               |                              |
| 2D domain    | $2\frac{\sqrt{N}}{\sqrt{P}}$ | $\frac{N}{P}$ | $2\frac{\sqrt{P}}{\sqrt{N}}$ |
| 5pt stencil  | $4\frac{\sqrt{N}}{\sqrt{P}}$ | $\frac{N}{P}$ | $4\frac{\sqrt{P}}{\sqrt{N}}$ |

MPI\_Pack (buf)  
MPI\_Send (buf)

MPI\_Recv (buf)  
MPI\_Unpack (buf)

(just)

"Gather the count first"

'one  
integer)

' for Gather variable size.