

# HPC (U-3).

## Parallel Communication

Page No.

Date

### \* One to all broadcast & All to one Reduction -

- Processes  $\rightarrow$  exchange data  $\rightarrow$  other process  $\rightarrow$  parallel algorithms.
- change of data  $\rightarrow$  efficiency affected  $\rightarrow$  interaction delays  $\rightarrow$  execution.
- communication  $\rightarrow$  various parallel architecture  $\rightarrow$  improve performance & reduce development effort.

### \* Different architecture are-

- 1) One to all Broadcast (scatter).
- 2) All to one Reduction (gather).
- 3) All to All Broadcast & Reduction.
- 4) All-Reduce & prefix-sum.
- 5) Scatter & gather.
- 6) All to All personalized.
- 7) Circular shift.

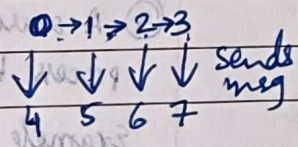
### Linear arrays, meshes & hypercubes.

$t_s$  = It is startup time for data transfer

$t_w$  = per word transfer time

It is same b/w all pair of nodes.

$M_1$  = communication cost in parallel machine.



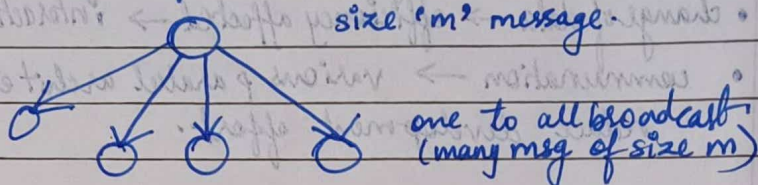
### \* One to One Broadcast with Recursive Doubling.

- source node broadcasts msg to all others efficiently.
- eight node ring with Node 0 as source.
- each step doubles distance b/w sender & receiver nodes, forms hierarchical comm<sup>n</sup> plan.
- recursive doubling technique ensures efficient one to all broadcast across eight node ring.
- Nodes with msg forward to opp neighbor (not already sent to)
- Repeats for  $\log_2(\text{no. of nodes})$  steps  $\log_2(8)$



## ① One to all Broadcast -

- ① Parallel algorithm often require a single process.
- ② It is used to send identical data to all processes or subset of them.
- ③ one to all Broadcast.
- ④ one processor has a piece of data size 'm' it needs to send everyone.

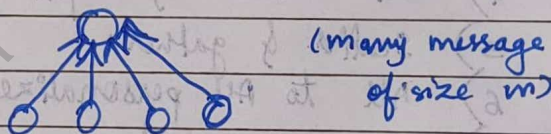


## ② All to one Reduction -

- ① Reduction can be performed by simply reversing direction.
- ② In all to one reduction each processor has m units of data.
- ③ these data items must be combined piece-wise using same associative operators.

0 1 2 3  
4 5 6 7

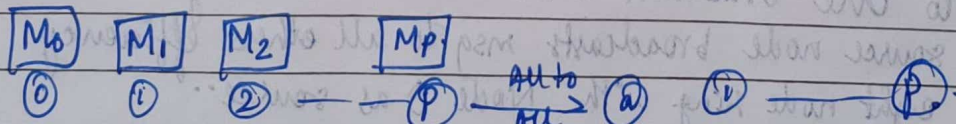
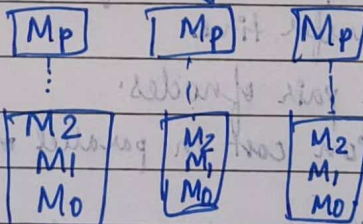
eg → addition or min.



## ③ All to all broadcast -

- ① All to all broadcast is generalisation of broadcast → processor is source.
- ② simultaneously every process out of p processes initiates broadcast.
- ③ A source process sends same m-word message to every other process but different processes may broadcast diff messages.

Example -



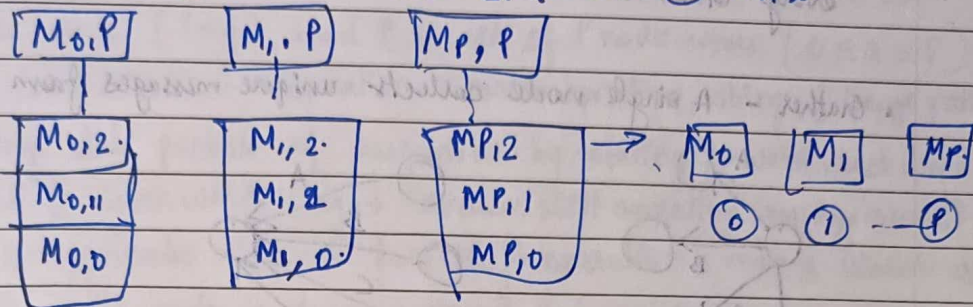
All to all Broadcast



#### ④ All to All Reduction -

All to All Reduction is reverse of all to all broadcast in every node. destination of all to one reduction.

$$M_x = M_{0,x} \oplus M_{1,x} \oplus M_{2,x} \dots \oplus M_{P,x}$$



#### \* All reduce & prefix sum operators -

• All reduce communication operation -

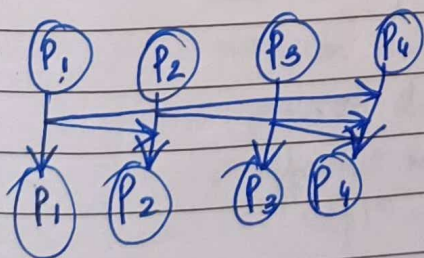
- ① Identical to all to one reduction followed one to all broadcast.
  - ② formulation is not efficient.
  - ③ faster way to perform all reduce by using pattern All to all broadcast.
  - ④ only diff is message size doesn't increase here.
- $$CA - T = (t_s + t_w M) \log P$$

#### \* Prefix - Scan Communication Operation -

• Scan operation.

example - Input - 1 2 3 4 operation - A.D.D. - Output - 1 3 6 10

- each output  $\rightarrow$  sum of all numbers  $\rightarrow$  input upto given point.
- histogram uses prefix scan. Also used for data compression.
- scan operation perform parallel prefix.

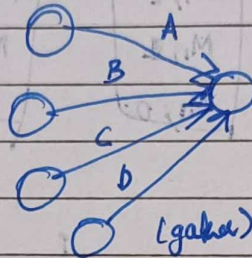
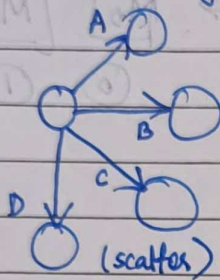


Scan communication.



## \* Scatter & Gather -

- Scatter - single node sends a unique message of size  $m$  to every other node. called as one to all personalized communication.
- Gather - A single node collects unique messages from each node.



C.A P Processes :- No. of steps  $\log P$ .

$$\text{Total time } T = t_s \log P + t_w (P-1)$$

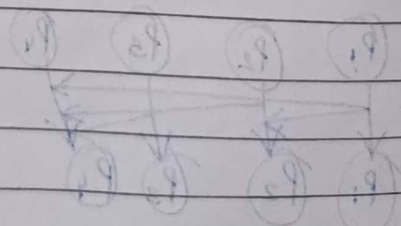
## \* All to All Personalised Communication -

- Communication  $\rightarrow$  Total exchange.
- each node  $\rightarrow$  sends same msg  $\rightarrow$  nodes in Broadcast.
- each node sends  $\rightarrow$  direct msg of size  $m \rightarrow$  every node
- parallel algorithm  $\rightarrow$  parallel sorting, fast fourier transform.

CA = processes =  $P-1$

Transmitted msg size =  $m(P-1)$ .

$$\text{Total time} = T = t_s (P-1) + \sum_{i=1}^{P-1} P t_w m = t_s (P-1) + (P-1) P t_w m$$





## \* Circular shift -

- particular permutation (redistribution in set)
- circular  $q$ -shift - special permutation in which node  $i$  sends a data packet to node  $(i+q) \bmod P$  in set of  $P$  nodes where  $(0 \leq q \leq P)$
- circular communication  $\rightarrow$  operation  $\rightarrow$  matrix operation, string & image pattern
- permuting data packets of each node by rotating forward, backward, upward & downwards. (eg  $\rightarrow$  circular shift operation on hypercube).
- develop a hypercube algo  $\rightarrow$  for shift operation, map a linear array with  $2d$  nodes onto a  $d$ -dimensional hypercube
- perform mapping  $\rightarrow$  assign node of linear array to node  $j$  of linear array of node  $i$  of any  $q$  in  $p$ -node hypercube is most  $2 \log(P-1)$
- $T = (t_s + t_{wm}) (2 \log p - 1)$

## \* Improve speed of commun<sup>n</sup> operation -

- need speed of communication  $\rightarrow$  3 operations

### 1) one to all broadcast -

- splitting & routing messages into parts scatter + All to all broadcast.
- improve speed of one to all, split message into  $P$  part & perform scatter operation.
- scatter - divide msg into  $n$  part & every node carry each part of msg.
- CA -  $T = 2 \times (t_s \log p + t_{wm})$

- shift repositions elements of array, by specified distance cyclically
- moves element from one end of array to other.
- data reordering & rotation in parallel algo
- shift can be left or right
- preserves data locality & reduce memory access impact pattern
- efficient manipulation of large dataset
- Appl<sup>n</sup>  $\rightarrow$  cryptography, signal processing & image manipulation



## 2) All to one Reduction -

- performed by gather operation. All to all reduction + gather.
- msg are reduced to one message of gather reduced msg.

## 3) All reduce -

- performed by performing all to one reduction followed by one to all broadcast.
- asymptotically optimal algorithms for two operation.
- construct a similar algorithm all reduce operation.
- gather & scatter operations - cancel each other require all to all reduction & broadcast.

• All to one reduction + one to all broadcast.

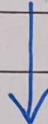


All to all reduce  
+ Gather



Scatter

All to all broadcast



(scatter & gather)  
(cancel operation)

All to all reduction  
+ All to all Broadcast

Fig. All reduce operation.



V-3

## \* Block communication using MPI in HPC -

- Initialization - MPI  $\rightarrow$  Message Passing Interface  
 $\rightarrow$  Is initialized to setup communication.
- Process Setup  $\rightarrow$  multiple processors are created each on separate node or processor.
- Blocking send  $\rightarrow$  sender process executes MPI-send, halting until msg is sent.
- Blocking Receive  $\rightarrow$  receiver process executes MPI-recv, halting until msg is received.
- completion - Both resume after msg trans<sup>n</sup>, continuing further processing.

Adv  $\rightarrow$  • Performance overhead  
 • ensure data exchange before continuing

Disadv  $\rightarrow$  • performance overhead  
 • scalability with large processes.

For use case  $\rightarrow$  • simple comm<sup>n</sup> patterns • small msg & low comm<sup>n</sup>  
 • critical synchronization needs

## \* Non-blocking using MPI -

- Initialization  $\rightarrow$  to set up comm<sup>n</sup> (MPI)
- Process Setup  $\rightarrow$  Multiple processors are created each on separate node.
- Non-blocking send  $\rightarrow$  sender process executes MPI-Isend, allows to continue execution immediately.
- Non-blocking receive  $\rightarrow$  receiver process executes MPI-Irecv, allows to continue execution immediately.
- completion check  $\rightarrow$  Both sender & receiver process check completing using MPI-test or similar function.

Adv  $\rightarrow$  • Improved performance  
 • Better Scalability  
 • Flexibility

Disadv  $\rightarrow$  • more complex programming  
 • Potential for deadlock

Use case  $\rightarrow$  • complex communication patterns • overlapping comm<sup>n</sup>  
 • large msg for long time



## \* Prefix Sum Operation

- also known as scan
- computes cumulative sum of elements in sequence
- crucial for various algo. like parallel reduction.
- sum of all elements up to certain index in array
- Result at each index is sum of all elements preceding it.
- application like  $\rightarrow$  computational geometry, image processing, cryptography
- requiring prefix info for each element.
- efficiency & applicability in various domains