

NO:

DATE:

Gamble, Risk

17/02/24 Decomposition Techniques.

#### 4) Speculative Decomposition :-

6/6/63

⇒ Branch prediction.

#### 5) Hybrid Decomposition :-

→ we combine decomposition technique. (data, recursive, ....)

#### Mapping Schemes :-

⇒ # of tasks

} → apriori

⇒ size of each task

to know  
something  
before time

⇒ if known, then

Static Mapping → decides which task  
at which process

⇒ unknown,

Dynamic

Mapping

Static

causes load imbalance,

know about size or #

at time of actual  
execution.

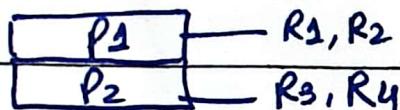
↓  
decides  
during  
execution.

# ↑ Static Mapping

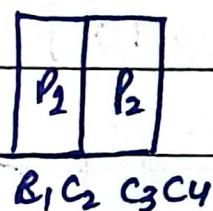
NO: \_\_\_\_\_  
DATE: \_\_\_\_\_

## ↓ Array Distribution Schemes :-

Block distribution → 1D → (divide rows or columns)  
row-wise dist



column-wise → n/p rows as block



2D :-

$$P = p_1 * p_2$$

4x4 grid  
2x8

$$\text{chunk size} = n/p_1 * n/p_2$$

Cyclic :-

↳ no of chunks > # of p's

4x4 grid, p=3

make sure no

0 1 2

load imbalance

0 1 2

Block-Cyclic dist :- (issue)

NO:

DATE:

Sparse Matrix  
→ Task 9 inter & aut

issue → some P's  
idle, some  
have more  
load.

Sol:- randomize the  
mapping

$V = \{0, 1, 2, 3, \dots, 11\}$  elements.

$\text{random}(V) = \{8, 2, 6, \dots, 10\}$

mapping =  $\underbrace{8 \ 2 \ 6}_{P_0} \ \dots \ \dots \ \dots$ ,

Threads:-

- ↓ ↳ inter-process comm → more expensive
- ↓ ↳ inter-thread comm → less "

Pthread → POSIX

independent stream of ins.

→ Strcpy → don't do bound

checks, so not a

recommended fun → buffer

• overflow.

↗ integer overflow:-

## OpenMP:-

override default value

NO:

DATE: 19th Feb 25

omp\_set\_num\_threads(4);

openMP: → by default 4 threads.

- ↳ we can set environment variable e.g it can be 2,  
`omp_get_thread_num()`
- ↳ return thread id of threads.

$\Rightarrow$  must include <omp.h>

In visual  
Studio, enable

#pragma omp parallel

$\{$  // Run in parallel

3

OpenMp.

Output: hello(0) world(0)

hello (1) world (1)

hello (2) world (2)

hello (3)    world (3)

```
#pragma omp parallel num-threads(4) ...
```

directives.

get # of threads ← open\_get\_num\_threads()  
actually created

NO: \_\_\_\_\_

DATE: \_\_\_\_\_

## Loop work sharing construct:-

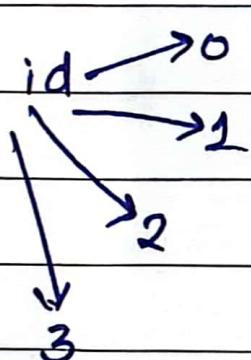
Sequential code add 1D array

open mp parallel  
region

master thread

first thread always  
has 1D value = 0

→ give array  
parts to threads.



istart = 0 iend = 25

istart = 25 iend = 50

istart = 50 iend = 75

istart = 75 iend = 100

In-built construct for FOR loop of  
Open mp

# pragma omp parallel      can write here  
# pragma omp for      no need for  
                              separate line  
                              →omp assign  
for(i=0 ; i < N ; i++) {  
    a[i] = a[i] + b[i]; }  
                              itself to  
                              every threads  
                              in proper  
                              way.

NO:

DATE:

## For R-join model :-

Compile:

Linux:

g++ -fopenmp hello-world.cpp  
open mp  
enabled.

Pi :-

$$\int_0^1 \frac{4 \cdot 0}{(1+x^2)} dx = \pi$$

⇒ break the curve into very small curves or rectangles and cal area of  $\square$ s.

we can do by openmp.

$$\sum_{i=0}^N (F(x_i) \Delta x) \approx \pi$$

middle of  
 each interval

width  
 height

Step = the width of each  $\square$

$$= \frac{1}{100,000} \rightarrow x = 0.5$$

$$i = 0 \quad i = 1 \quad 100,000$$

$$i = 2 \quad x = \frac{0.5}{100,000}$$

$$x = \frac{1.5}{100,000}$$

NO: \_\_\_\_\_

DATE: \_\_\_\_\_

$\Rightarrow \text{nthreds} = \text{omp-get-num-threads}$  (good practice) use it instead of just using `#define Num-threads 4`.

### Synchronization :-

mutual exclusion :- only one thread at a time can enter critical section.

`# pragma omp critical [name]`

{

code-block

}

- [name] optional name that identifies critical directive.

Previously sum array, but now

Sum in critical region so no issues

due to sum  $\rightarrow$  only one  $\frac{1}{3}$  can update <sup>(P1)</sup> at a time

also its initialized

inside parallel so they are private to each thread.

separate copy  
↑ of every thread.

OpenMP

```

graph TD
    A[OpenMP] --> B[env variable]
    A --> C[directives]
    A --> D[library functions]
  
```

→ openMP for constructs.

→ openMP parallel section

→ openMP critical section.

## Static Scheduling

# pragma omp for schedule (static, chunk-size)

if not

Specified, it  
will split  
into equal  
chunks e.g

i=100  
chunk-size=25

e.g  
i=100

chunk size =  
10

means 10  
iterations in  
one thread.

(total e.g +  
threads)

10	20	30	40
T1	T2	T3	T4

50	60	70	80
----	----	----	----

90	100
----	-----

NO:

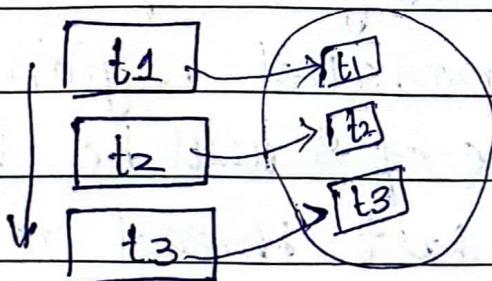
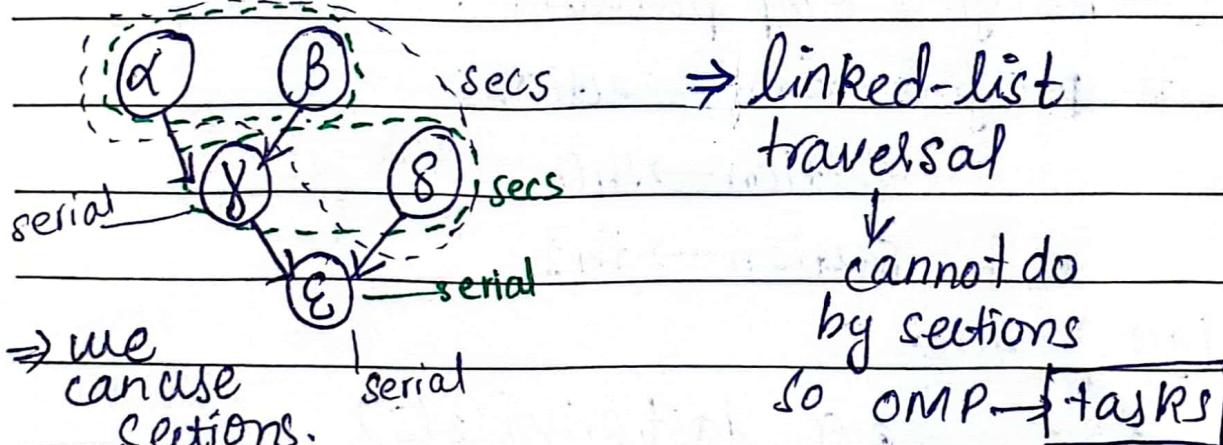
DATE:

⇒ if last private not set, then S value is value that set before # pragma.

## ⇒ Conditional parallelism :-

if (`{scalar expression}`)

## ⇒ OPEN MP task construct :-



only one thread will do it

→ tasks forming by while loop

→ made by one thread.

→ put in (add on) pool

→ their execution by dynamically allocated threads.

## firstprivate

e.g. `firstvariable(s)`

master thread set value will be assigned to s (private variable) of threads initially.

# pragma omp parallel section

Section → th0

Section → th1

# pragma omp parallel

# pragma omp sections

section → th0

section → th1

## lastprivate

e.g. `lastprivate(s)`

→ out of parallel region then

value of s will be last update.

(e.g. last thread done work on s)

NO:

DATE:

## Private

# pragma omp parallel for private(j)

⇒ i will be shared every thread will have own copy of j  
has thread ka private ban jayega

[before] → (shared)

var, [parallel]

then by default

[after] then by default (private)

default (private) ← iteration  
kay liye use nota ha

[default (none)]

in force programmers to tell about each var either its private or shared, good prog practice.

\* const var, don't write priv or shared

⇒ [default (shared)]

↳ if not mention, then by default shared

NO:

DATE:

5th March '25

## Scheduling loops

### 1) Static:

[ ] (optional)

Static, 10 → 10 iterations  
thread.

schedule ( static [ , chunk-size] )

if not mention, it  
divides itself  
equally to # of  
threads.

### 2) Dynamic:

↓ at run time

if NO chunk size, then

default size = [ 1 ]

schedule ( dynamic [ , chunk-size] )

### 3) Guided: It decides chunk size

heuristically.

schedule ( guided )

if no chunk size = then min [ 1 ]

schedule ( guided , C )

↳ chunksize

schedule ( runtime )

[ OMP\_SCHEDULE ]

value of variable  
either  
static or dynamic

→ environment variable

NO:

DATE:

omp-set-num-threads(2);

#pragma omp parallel default(none)

#pragma omp single

#pragma omp task

printf(stderr, "A\n");

#pragma omp task

printf(stderr, "B\n");

}

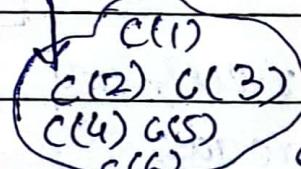
into pool,

execution separate from

creation of tasks.

task pool

thread 0: a, c(1), c(7),



c(5),

thread 1: [c(2)]

thread 2: [c(3)] [c(6)]

out of  
parallel

thread 3: [c(4)]

region so  
executed by  
master  
thread.

# Distributed Systems

NO:

12 March 2025

DATE:

based on → Hadoop → big data  
map-reduce model : Apache  
Spark

# Google PageRank algo

IOT (Internet of things)

## Cloud computing

## Cluster Computing :-

Set of comps, near to each other, connected, performing like they are single machine.

Gridcomputing: → interconnected but owned by multiple organizations

↓ mostly for scientific applications

## On-prem

"on-premises"

Cloud computing ← switch to

## ایضاً دفتر

due to (i) cost is less (a lot cheaper)

means ← ~~once gone  
can't be  
recovered~~ Sunk cost (initially spent  
to start business)  
in On premises

NO: \_\_\_\_\_

DATE: \_\_\_\_\_

## Monte Carlo Method :-

↳ randomized algo

↓ generate many points in square

$$\Rightarrow \left[ (x-a)^2 + (y-a)^2 < r^2 \right] \quad \begin{array}{l} \text{in circle or} \\ \text{out of circle} \end{array}$$

↓  
(inside square)  
must

① calculate # of points in circle

② divide points in circle and  $\times 4$  ③

Point in  $\square$  → total points  
as # of random value actually.

Increase,  $\pi$  value becomes more accurate.

parallelize → divide iterations.

thread ①

# pragma omp master

↓ master thread  
will execute it  
(like single)

TRUE

Nested parallelism

reasonable # of threads on upper layer are enough.

NO:

DATE:

**OMP\_NESTED**

False (default)

True

# pragma omp parallel num-threads(4)

{ :

| to t<sub>1</sub> t<sub>2</sub> t<sub>3</sub>

# pragma omp parallel num-threads(2)

{ → this

If FALSE

decides if new team of threads created or not or original used.

to → [t<sub>0,0</sub> t<sub>0,1</sub>]

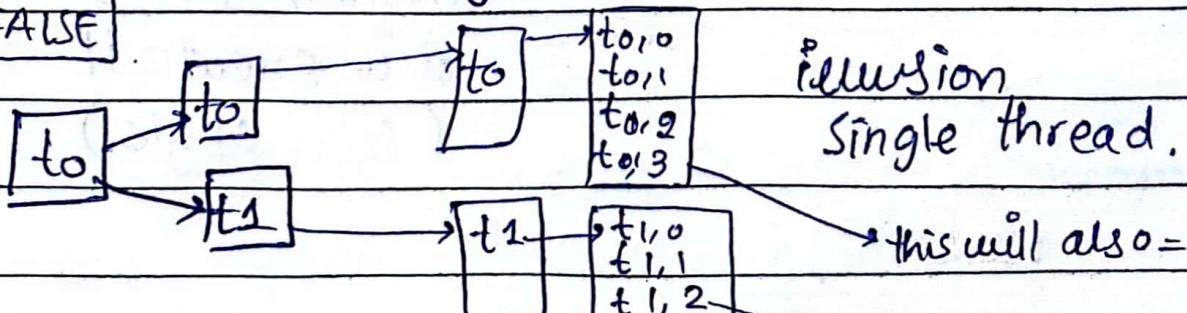
not creating new threads,

image of multithreading but not really multithreading

t<sub>1</sub> → [t<sub>1,0</sub> t<sub>1,1</sub>] → actually → executes 1 thread. Serially.

TRUE → enables nested parallelism, makes actually a new team of threads.

If FALSE



If TRUE

then

threads = 4

even if single (nested)

# Environment variables

NO:

DATE:

10 March 2025

at runtime



void omp\_set\_dynamic (int flag)  
↳ called from outside of  
parallel region

omp\_get\_max\_threads()

omp\_get\_num\_procs()

↳ # of processors in system.

⇒ omp\_set\_dynamic(1) → it will dynamically  
decide how many threads

omp single, nowait,

↳ remove implicit barriers

not necessary

8 threads,

can be minimum

or more (not so)

much

(means other threads

thread will

execute this

line.

not going to

wait for

it)

## OMP\_SCHEDULE

→ static, dynamic,  
guided.

↳ chunk size

default = 1

windows: \$env: OMP\_SCHEDULE = "guided, 4"

ubuntu: export OMP\_SCHEDULE = " "

NO:

DATE:

## Message passing:-

↓ request / response

→ synchronous (indication from server if it's ready)

→ asynchronous (no need to wait)

buffers

→ faster plus less reliable.

## Communication Model :-

$$\frac{d}{s} = \frac{d_{prop}}{d_{proc}} + \frac{d_{trans}}{d_{proc}} \quad \text{th} \rightarrow \text{time for each hop}$$
$$\frac{L}{R} = \frac{n}{d_{prop}}$$
$$n = \# \text{ of hops}$$

## Architectures:-

e.g.

↓ P2P

→ Napster, Gnutella,

↓ Client / Server

↓ to download  
music (mp3)  
file sharing

## Issues:

→ Programming complexity

→ Scarce Robustness

$n_1 \rightarrow n_2$   
 $n_1 \rightarrow n_5$  } same type

NO:

DATE:

17 March 2025

## Basic Communication Operations

### Assumptions:-

1) Support cut-through routing



message should be able to transfer from router even if one bit of a packet is received

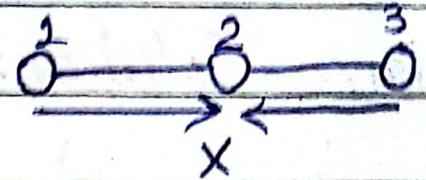
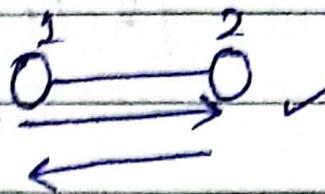
2) Links are bi-directional

data sending can be on both ways at one time or at one single clock cycle

3) Single-port

In one clockcycle

one place receive or send at single time.



$m$  = size of message     $t_w$  = trans. time for one word

$t_s$  = startup time

$$\frac{m t_s}{100MB * 1} = \boxed{10s}$$

10Mbps

NO:

DATE: 19th March 2025

circular [5-shift]

# of nodes  
 $P = 16$

$$q \bmod \sqrt{P}$$

# of shifts = 5

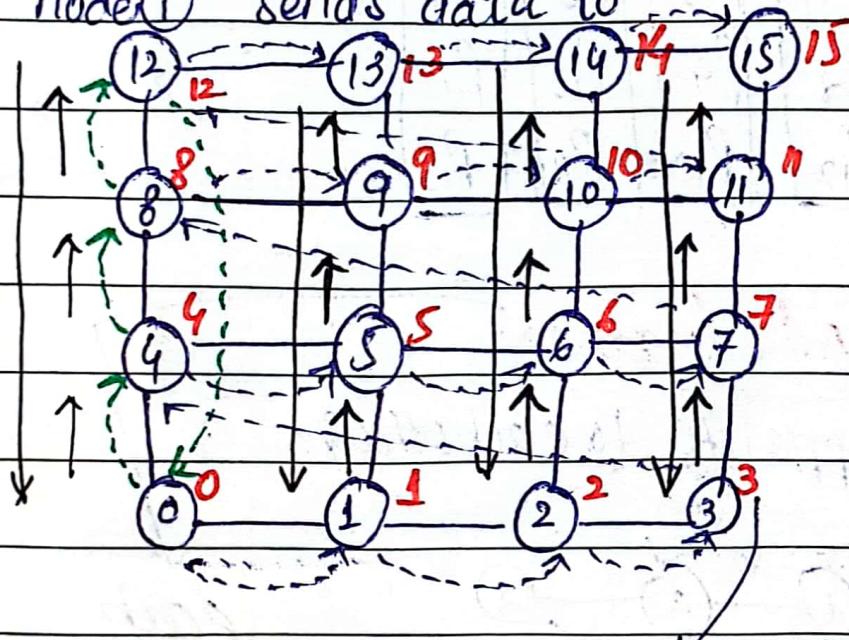
# of rows &  
cols =  $\sqrt{P}$   
= 4

4x4 mesh.

$$\sqrt{P} * \sqrt{P}$$

circular shift → naive

node 0 sends data to



+ all  
links are  
bidirectional

(t1)

① right shift  
all

(t2)

② one shift  
up of  
left most  
column

manage mod  
stuff.

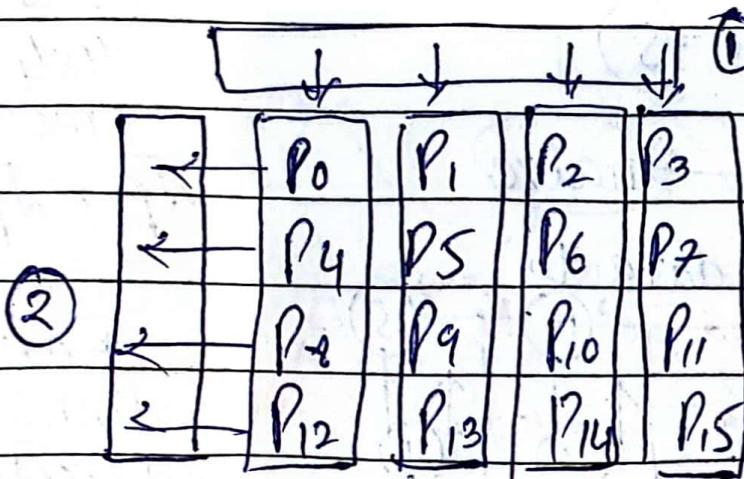
③ each one  
shift up  
(t3)

These  
values  
shift.

## → Matrix Vector Multiplication :-

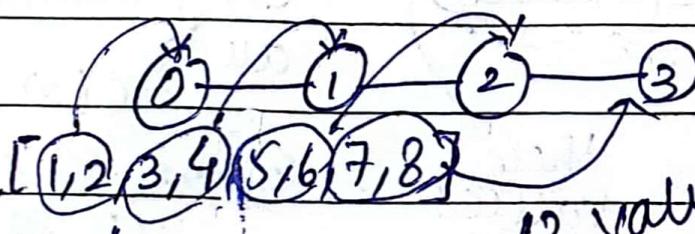
- 1) One to All broadcast
- 2) All to One reduction.

log P



## Scatter and Gather

opposite to each other.



Send  
mutually  
executive

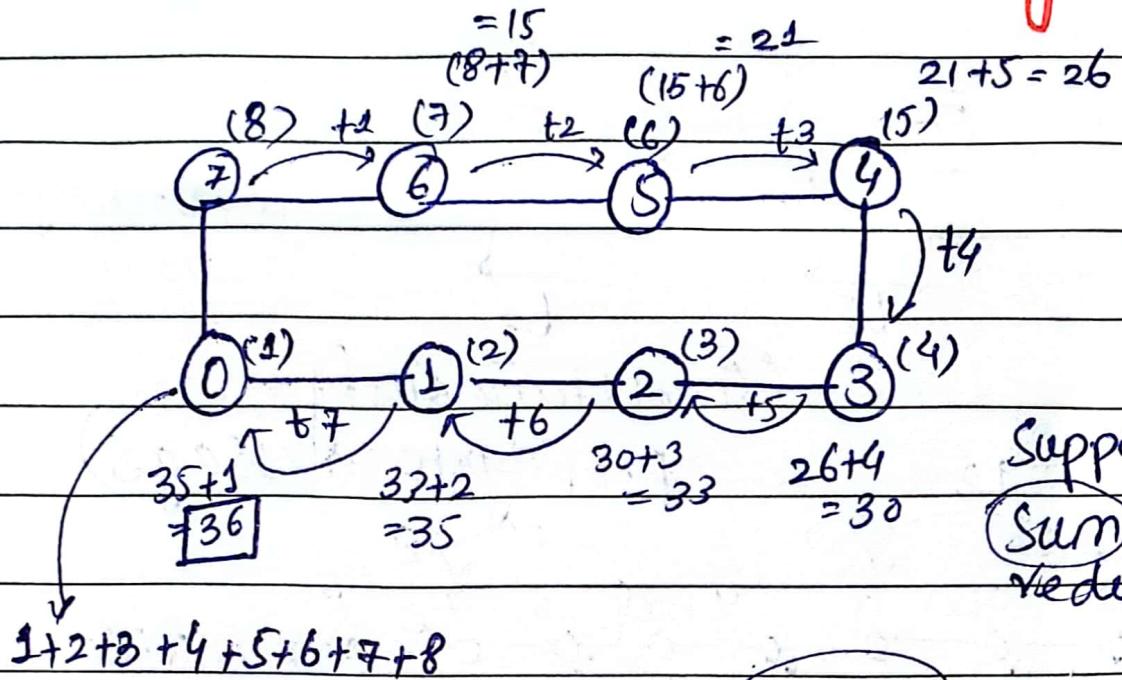
Values to  
Separate nodes.

(12 values)  
each  
node.

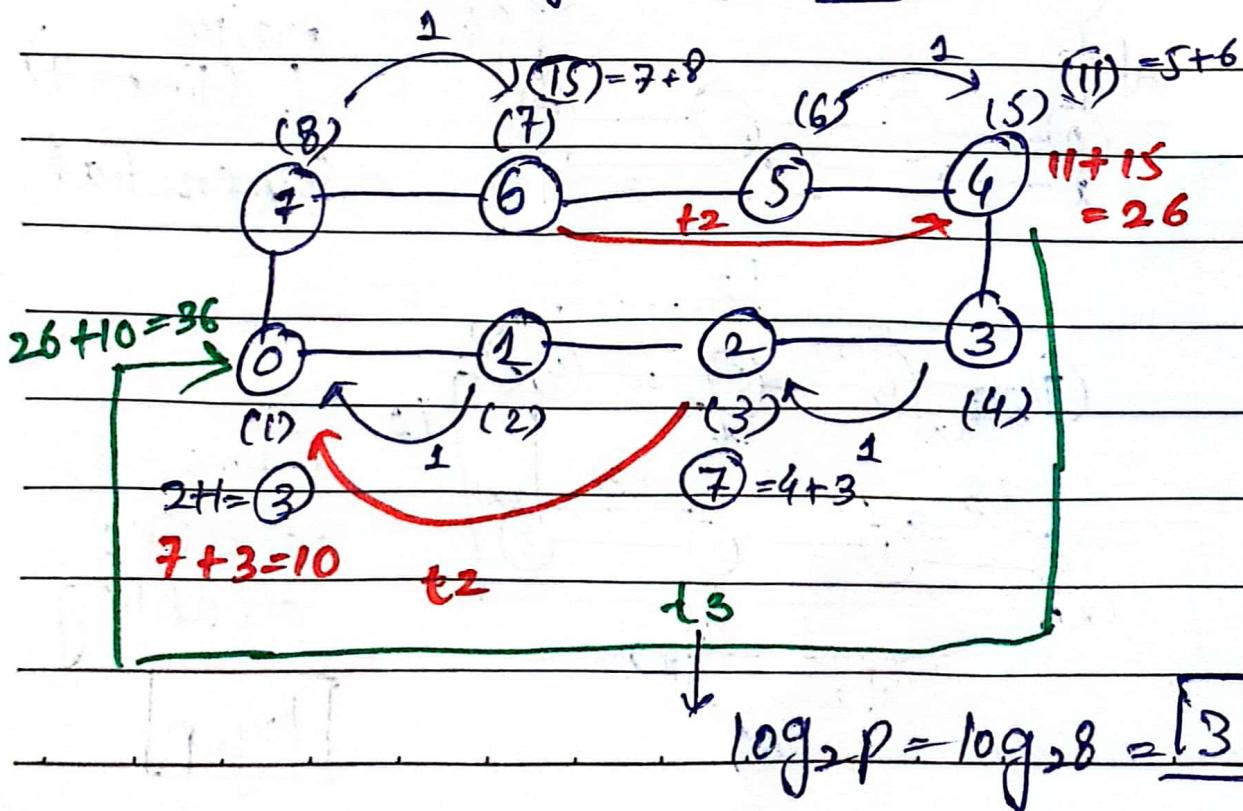
one  
node  
may have  
8 values.

NO:

DATE:

• naive ( $P-1$ )Reduction: • recursive doubling ( $\log_2 P$ )7 clockcycles  $\rightarrow$  naive :-

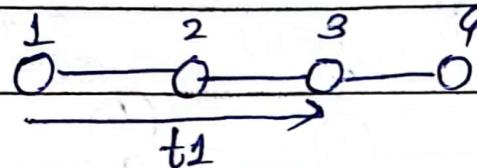
opposite of broadcast  
 recursive doublings  $\rightarrow$  all to one  
 same time.



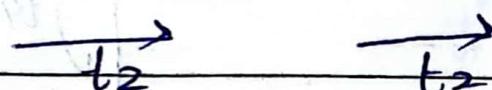
NO:

DATE:

## Recursive doubling :-



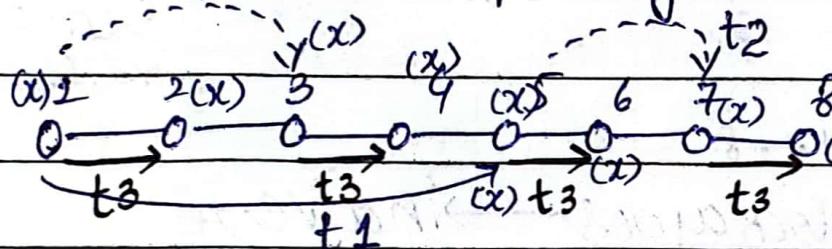
2 clockcycles



now.

(Simultaneously

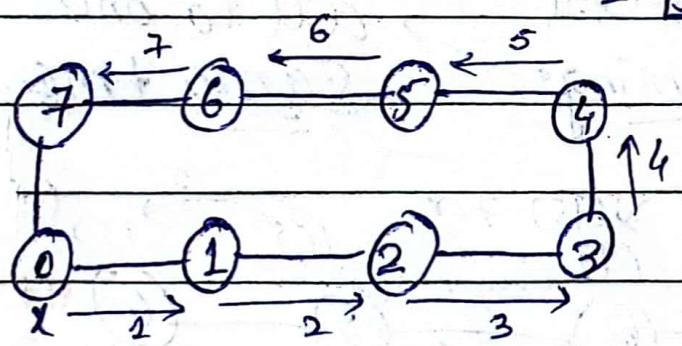
propagate message)



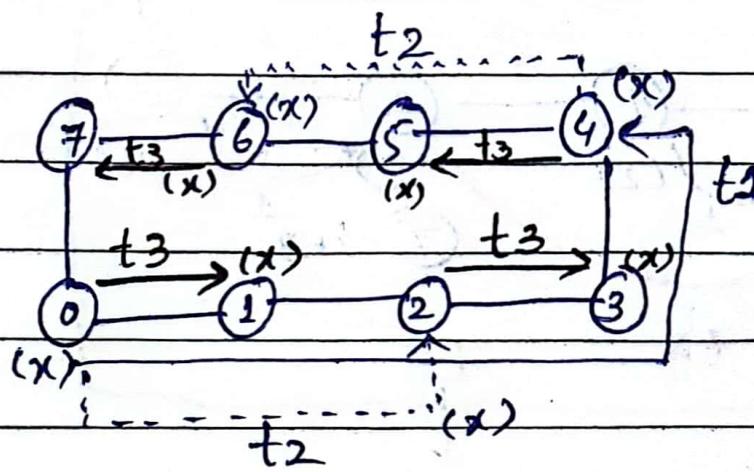
↳ 7  
clock  
cycles.

$$\Rightarrow \log_2 P = \log 8$$

= 3 clockcycles.



↓  
naive  
Sequential



$$\log_2 P$$

NO: \_\_\_\_\_

DATE: \_\_\_\_\_

## One to All Broadcast :-

→ single process sends identical data to all other processes.

each have copy of

size m

All to One → opposite of broadcast

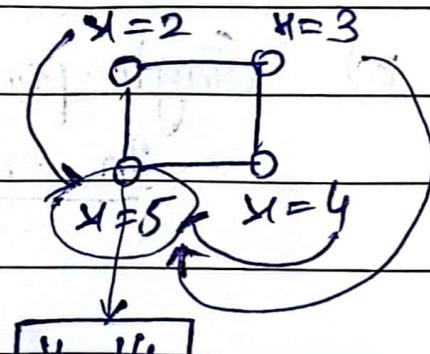
## Reduction :-

All processes data are accumulated at single destination.

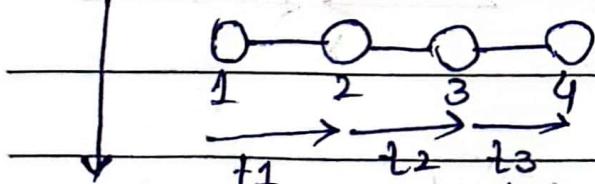
into one buffer of size m.

through associative operator  
→ +, -, \*, ...

## Linear Array or Ring:



## Naive Solution:-



sequential message transfer

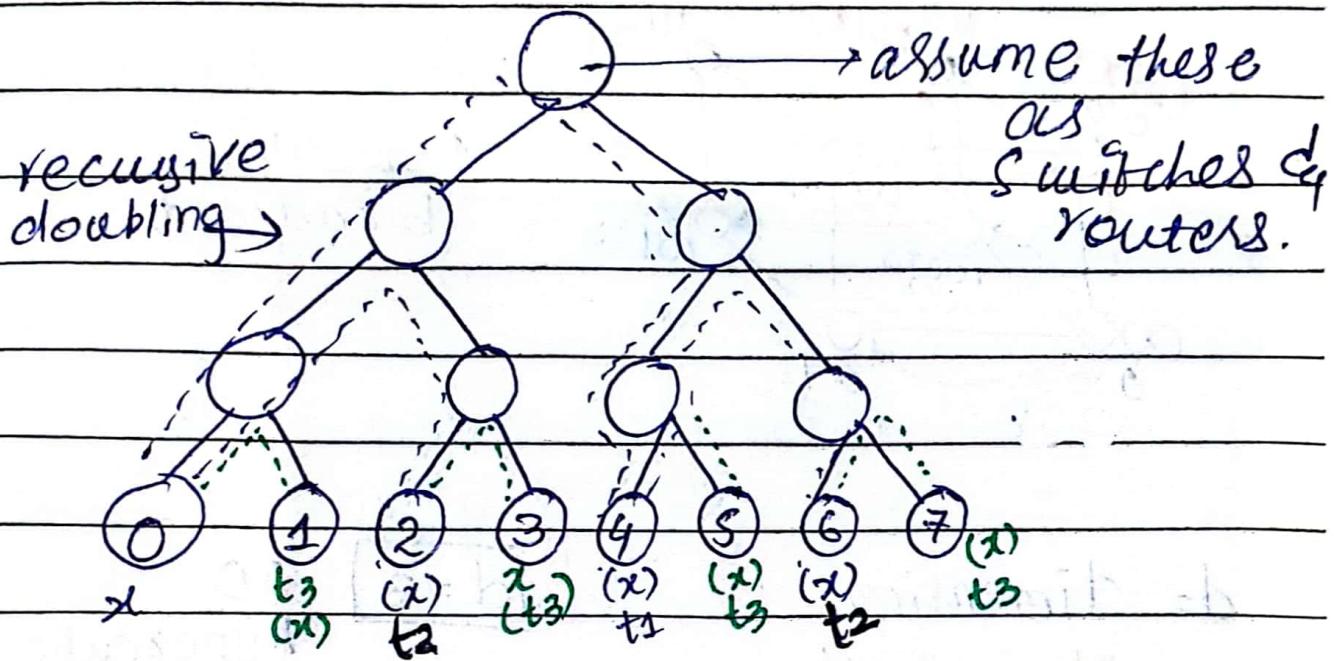
$[P - 1]$  message

clockcycle for P processes.

NO:

DATE:

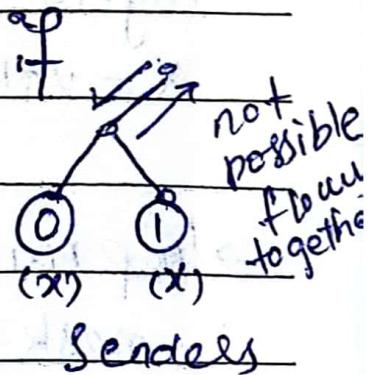
# Binary tree One to All Broadcast



③ 3 clockcycles

$$\log_2 8 = 3$$

log P



## All to One Reduction

opposite

① every leaf send to its left neighbour

(1, 3, 5, 7)

② (6, 4, 3, 0) (odd values)

③ 4, 0

[7, 6, 5, 4]

[7, 6, 5]

[7, 6]

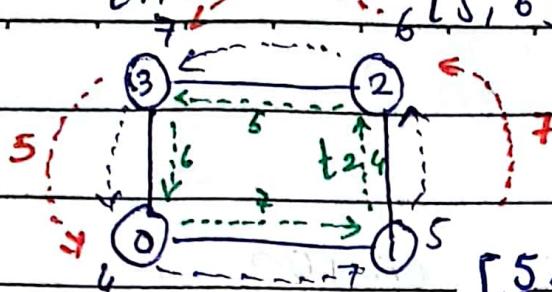
[6, 5, 4, 7]

[6, 5, 4]

[5, 6] t1

NO:

DATE:



⇒ move

anti-clock  
wise.

[4, 7]

[4, 7, 6]

[4, 7, 6, 5]

[5, 4, 7]

[5, 4, 7, 6]

associative operation

All to All Reduction: → every node is applying all to all reduction.

[7+6+5+4]

goal

(3)

(2)

(0)

(1)

[6+4+5+7]

[4+6+7+5]

[7+4+5+6]

[6+7+4+5]

⇒ move

clock  
wise

Naive

[7+4]

[6+7]

[6+7+4]

t3

t1

t2

[4+5+6+7]

[4+5]

[4+5+6]

[5+6+7]

[5+6+7+4]

NO:

DATE:

Id ual  $\rightarrow$  opposite.

Scatter: 1 node has several values  
So, send some of these  
values to some nodes

Gather: It is opp of gather. But  
we are not applying any  
associate operations

all to all broadcast :-

$\rightarrow$  Recursive

doubling better  
then  
naive

but in this  
case Naive

modification is  
better

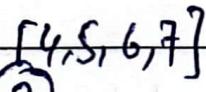
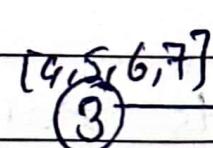
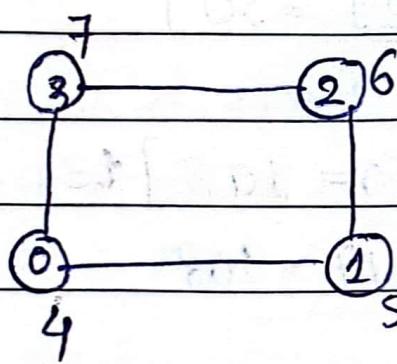
goal:

{4, 5, 6, 7}

{4, 5, 6, 7}

[4, 5, 6, 7]

[4, 5, 6, 7]



[4, 5, 6, 7]

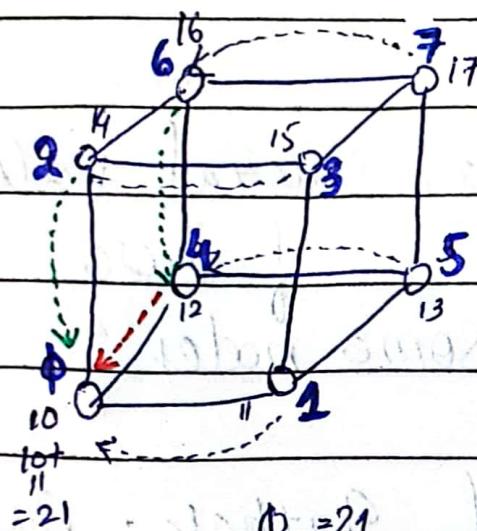
[4, 5, 6, 7]

NO:

DATE:

24 March 2025

Topology :- linear array,  
2-D ring, Hypercube,  
2D-mesh.



all to 1  
reduction  
(associative operation)

$$\phi = 21$$

] t1

$$4 = 13 + 2 = 25$$

$$2 = 15 + 14 = 29$$

$$6 = 17 + 16 = 33$$

$$4 = 33 + 25 = 58 ] t2$$

$$\phi = 29 + 21 = 50 ]$$

$$\phi = 58 + 50 = 108 ] t3$$

$$10 + 11 + 12 + 13 + 14 + 15 + 16 + 17 = 108$$

↓ to check correction.

$m \cdot t_w$  = transmission delay.

$t_s$  = start time

$\downarrow$   
size of message → time to transmit one word.

loading time /

ending time

$[t_s + m \cdot t_w]$  → time for each of transfers.

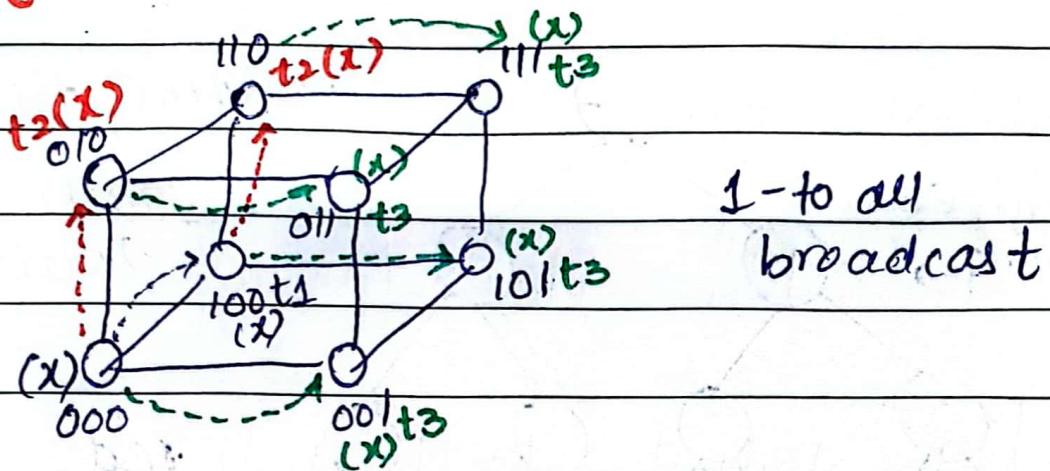
$$T = (t_s + m \cdot t_w) \log p \rightarrow \boxed{\text{Cost}}$$

NO:

DATE:

## Hypercube:-

P = 8 nodes



d = dimension

my-id = node-id

x = value

$\boxed{d=3}$  3D

hypercube

$$\text{mask} = 2^3 - 1 = \boxed{7}$$

$$(7)_10 = (111)_2$$

→ loop runs d times

$$\boxed{i=2 = d-1}$$

$$\text{mask} = \text{mask XOR } 2^i$$

$$(111) \text{ XOR } 2^2 (100)$$

$$\text{mask} = \boxed{1011}$$

to decide

sends or

receives.

$$\text{if } (000 \text{ AND } 1011) = 000$$

$$\text{if } (\text{my id and } 2^2)$$

$$000 \text{ AND } 100 = 000$$

1 node All-to-1 Reduce

↑  
1 to All Broadcast

NO:

DATE:

All - Reduce:



MPI All Reduce.

All to All Reduction:

→ Same like

Strategies:

M#1) Use all to one reduction followed by one to all broadcast

M#2) step by step sharing of msg in dimension (1, 2, 3)  
replace union with associative operator.

Prefix-Sum: → like Cumulative Sum

$$\begin{array}{|c|c|c|c|c|c|} \hline 1 & 2 & 3 & 4 & 5 & 6 \\ \hline \end{array} \quad O(n)$$

↓↓↓↓↓

$$\begin{array}{|c|c|c|c|c|c|} \hline 1 & 3 & 6 & 10 & 15 & 21 \\ \hline \end{array}$$

NO:

DATE:

[13-16]

[9-12]

[5-8]

[1-4]

in column  
fault array is  
Shared rather  
the value.

add (1) & (2)

$$\text{Total timer} = \boxed{2ts(\sqrt{p}-1) + mt_w(p-1)}$$

all to all broadcast      clockcycle  $\rightarrow \lceil \log_2 p \rceil$

hypercube  $\rightarrow$  dimension wise

cost est.

$$1^{\text{st}} \text{ step} = (ts + mt_w) \xrightarrow{\substack{1 \text{ value} \\ \text{swap (shared)}}}$$

$$2^{\text{nd}} \text{ step} = (ts + 2mt_w) \xrightarrow{\substack{2 \text{ values} \\ \text{swap}}}$$

$$i^{\text{th}} \text{ step} = (ts + 2^{i-1}mt_w)$$

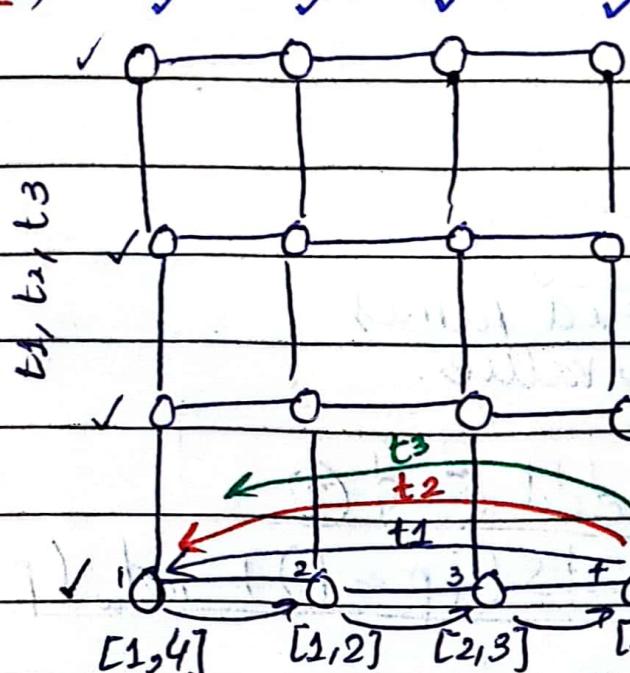
① Total cost  $= \boxed{\sum_{i=1}^{\log(p)} (ts + 2^{i-1}mt_w)}$

Simplify  $\rightarrow T = \boxed{ts \log p + mt_w(p-1)}$

costest  $\rightarrow$  all to all reduction

NO:

DATE:

(Mesh) $t_4 \ t_5 \ t_6$ 

rowwise ①

$[1,4] \rightarrow [1,2] \rightarrow [2,3] \rightarrow [3,4]$

$[1,4,3] \rightarrow [2,1,4] \rightarrow [3,2,1] \rightarrow [4,3,2,1]$

$[1,4,3,2] \rightarrow [2,1,4,3] \rightarrow [3,2,1,4]$

② columnwise

$$p = 16 \quad \sqrt{p} = 4 \quad \# \text{ of rows}$$

$\# \text{ of cols.}$

$$\underline{T(\text{first phase})} = (ts + mtw)(\sqrt{p}-1) \quad \text{--- } ①$$

↓                                  ↓  
rowwise                            # of cols.

$$\underline{T(\text{second phase})} = (ts + (\sqrt{p})mtw)(Tp-1)$$

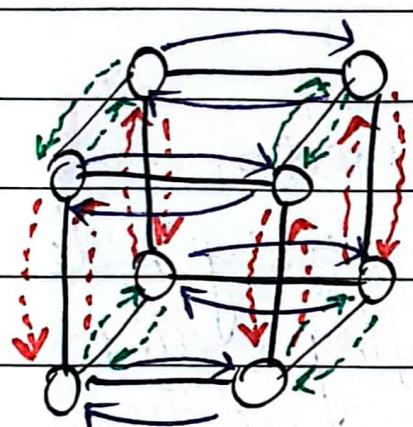
columnwise                          ↓                          ↓  
every node                            share an array  
                                        rather than numbers

②

NO: \_\_\_\_\_

DATE: \_\_\_\_\_

⇒ all to all broadcast on  
2D mesh

Hypercube:  dimension  
cuise

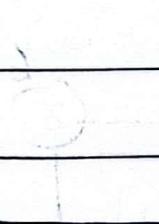
+1 +2 +3

$\lceil \log_2 P \rceil$  clockcycle

cost

26 March 2025

costest

↳ all to all broadcast  (Linear ring)

$$T = (t_s + m t_w) (p-1)$$

↓  
Startup  
time

↓  
#of  
processing  
nodes.

time to  
transfer  
a word.

theoretical  $\rightarrow$  no parallel overhead  
experimental  $\rightarrow$  ✓

NO:

DATE:

cost analysis:  $\sqrt{p}$  processors.

$$(t_s + \underbrace{(t_w m P_2)}_{\substack{\downarrow \\ \text{transmission} \\ \text{time for word.}}})(\sqrt{p} - 1)$$

$\Rightarrow$  Hypercube:

OpenMP

NO:

DATE:

$m(p-1)$  } size of message      # of clockcycles  
↓  
# of unique messages.

$p-1$

Auto All  
personalized [Ring]:

size  
decrease  
at each t

Cost Analysis:

$$T = \sum_{i=1}^{(p-1)} (ts + (p-i)m tw)$$

total  
time,  
communication  
cost

= Simplify

$$= ts + \left(\frac{1}{2}\right) pm tw (p-1)$$

$\Rightarrow$  mesh

$0 \rightarrow 1 \rightarrow 2$  row level done

② column level propagate

NO:

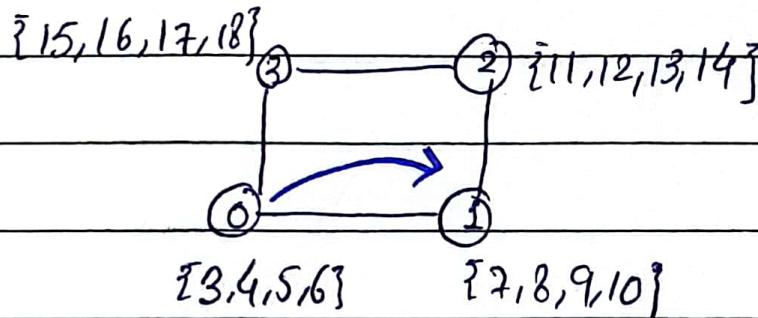
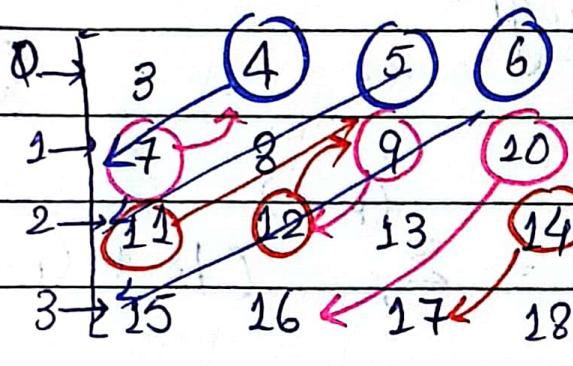
DATE:

7th April 2025

All to All Personalized  $\rightarrow$  all nodes

$M(0, p-1) \rightarrow$  destination  
source

do Scatter and  
it's happening  
in parallel



$\{4\{0,1\}, 5\{0,2\}, 6\{0,3\}\}$  Simultaneously  
 $\downarrow \downarrow$  in parallel

$p-1 \rightarrow p-2 \rightarrow p-3 \rightarrow p-4$   
at every clock cycle,

$\{7\{1,0\}, 9\{1,2\}, 10\{1,3\}\}$

Size of message decrease  
by 1

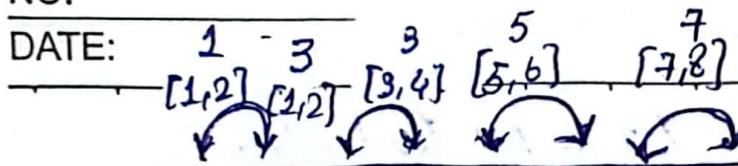
$\{12\{2,0\}, 12\{2,1\},$

$m(p-1) \rightarrow m(3)$   
 $\downarrow$  without message

$14\{2,3\}\}$

NO:

DATE:



Prefix Sum

$d=3$	1	2	3	4	5	6	7	8
	000	001	010	011	100	101	110	111

hypercube.

t1

$$\begin{matrix} [3,4] \\ 7 \end{matrix} \quad \begin{matrix} [5,6] \\ 11 \end{matrix} \quad \begin{matrix} [7,8] \\ 15 \end{matrix}$$

11

$$[5,6][7,8]$$

$$[7,8][5,6]$$

$$15+5+6$$

$$= 26$$

t2

1	2	3	4	5	6	7	8
$[1,2][3,4]$ $[1\dots 4]$	$[1\dots 4]$	$[1\dots 4]$	$[1,2]$	$[5,6][7,8]$	$[7,8][5,6]$		
1	3	$3+1+2$	$7+1+2$	5	$7+5+6$		

⇒ increase

distance

of neighbours

in powers of

2 at every  
clockcycle

t3

1	2	3	4	5	6	7	8
---	---	---	---	---	---	---	---

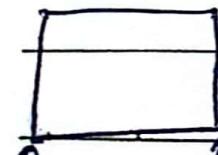
inverse gather

$$3 \quad 6 \quad 10 \quad 15 \quad 21 \quad 28 \quad 26+4+3+2+1 = 36$$

Scatter:

$$\begin{matrix} \{2,1\} \\ 3 \\ 2 \end{matrix} \quad \begin{matrix} \{4,3\} \\ 2 \end{matrix} \quad \{8,7,6,5,4,3,2,1\}$$

$\underbrace{0}_{0} \quad \underbrace{1}_{1} \quad \underbrace{2}_{2} \quad \underbrace{3}_{3}$



{8,7}

{6,5}