

Assignment - 8

(1) Heat-Equation :-

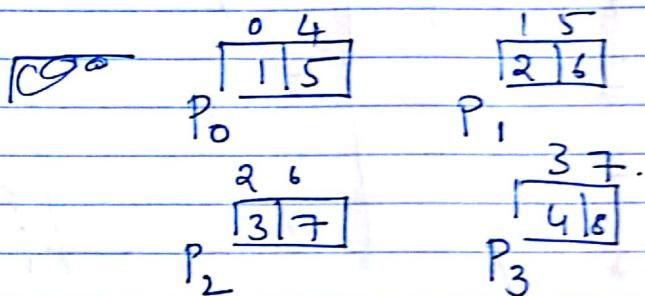
(1.1) Round-Robin Decomposition :-

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |

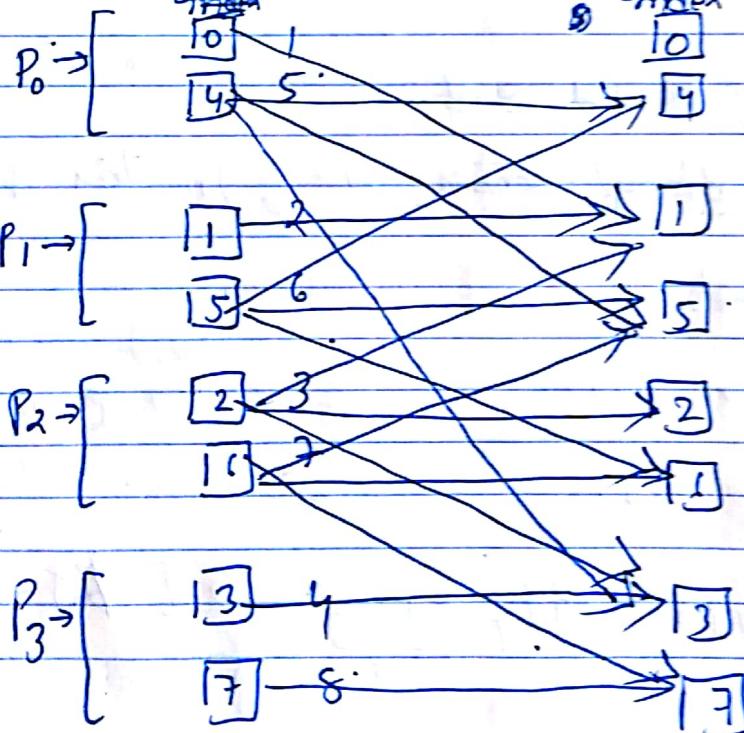
Let's say we have $P = 4$

P (no. of Processors)

$$N = 8$$



Every index has its global rank & local rank. (Assuming the Network is a clique)



{ g have a
gen function
to generate
data Locally }

// Pseudo mpi - program for Block wise decomposition

~~~~~~~~~ MPI\_Init(&argc, &argv);

int send\_value, recv\_value;

int my\_rank, total\_P;

~~int \*arr0 = new int[n/p]; int \*arr1 = new int[n/p];~~

MPI\_Comm\_rank(my\_rank, MPI\_COMM\_WORLD);

MPI\_Comm\_size(&total\_P, MPI\_COMM\_WORLD);

// arr0, arr1 have size( $n/\text{total\_P}$ )

// Every Processor Contains ( $N/p$  Size of array)

for (int k = 0; k < n; k++)

int P=0

for (int i = my\_rank; i < n; i = i + total\_P)

{

arr0[P] = ~~agen(i);~~

P++;

}

// arr0 denotes curr state

int send, recv; // arr1 denotes next state

for (int k = 0; k < n / total\_P; k++)

{

if (my\_rank == 0 && k == 0)

{

send = arr0[k];

MPI\_Send(&send, my\_rank + 1, MPI\_COMM\_WORLD);

MPI\_Recv(&recv, my\_rank + 1, MPI\_COMM\_WORLD);

arr1[K] = ~~2 \* arr0[K] + recv.~~

3.

}

$[P \rightarrow 0] \rightarrow$

else

~~if (my\_rank == total\_P - 1 && k == (n - 1))~~

{ // last element.

~~total\_P~~

else {

if (my\_rank ==  $\frac{total\_P-1}{total\_P} \cdot k$ ) for  $k = \left(\frac{n}{total\_P} - 1\right)$ .

// last element

MPI\_Send (&arr[K], my\_rank-1)

MPI\_Recv (&recv, my\_rank-1);

Compute arr[K]

arr[K] =  $\frac{arr[K] + recv}{2}$

};

else.

{ // here you will see. 2 sends & 2 recv for middle process

int send\_rank1 = my\_rank-1;

if (send\_rank1 < 0).

send\_rank1 = total\_P-1

int send\_rank2 = (my\_rank+1) % P

MPI\_Send (&arr[K], send\_rank1);

MPI\_Send (&arr[K], send\_rank2);

int recv1, recv2;

MPI\_Recv (&recv1, send\_rank1)

MPI\_Recv (&recv2, send\_rank2);

Compute arr[K]

arr[K] =  $\frac{arr[K] + recv1 + recv2}{3}$

3

## Communication in Round-Robin:

Except first & last Every process perform 2 sends & 2 Recv.

$$\text{Total Comm}^n = \text{first proc Comm}^n + \text{last proc Comm}^n + \text{rest all other proc}$$

$$= \left[ \left( \frac{N-1}{P} \right) * 2 + \frac{N-1}{P} \right] 2 + \frac{N-2}{P}$$

$$\text{total } P-2 * \left( \frac{N}{\text{total } P} \right) * 2 + 2 * \left[ \left( \frac{N-1}{\text{total } P} \right) * 2 + 1 \right]$$

$\underbrace{\hspace{10em}}$  ↓  $\underbrace{\hspace{10em}}$  ↓

other proc (middle proc)      first & last proc

so, here first & last proc also sends 2 sends + 2 Recv  
except for  $r_2(i=0)$  in  $\text{my\_rank}=0$  &  $(i=n)$  in  $\text{my\_rank}$

$$= \text{total } P-1$$

Per-link =  $\frac{N}{\text{total } P} * 2$  Everytime we are passing to send & 2 Recv

through the same Link.

Per-Node :-  $\frac{N}{\text{total } P} * 2$  for all-neighbors except first & last

for First  $\left[ \left( \frac{N-1}{\text{total } P} \right) * 2 + 1 \right] \rightarrow \text{send}$ .

For Last  $\left[ \left( \frac{N-1}{\text{total } P} \right) * 2 + 1 \right] \rightarrow \text{send}$ .

## 1.2 Block-wise decomposition

In Block-wise You need to store  $\frac{n}{p}$  elements in the each local-array of the process & output array to store  $\frac{n}{p}$  elements of next state.

```
int my_rank, P;
```

```
MPI_Initialized(&argc, &argv);
```

```
if (int K=0;
```

```
for (int i = my_rank * (N/P) ; i < (my_rank + 1) * N/P ; i++)
```

```
{
```

```
    arr0[K] = gen(i); // Block wise decomposition.
```

```
    K++;
```

```
? int send, recv;
```

```
if (my_rank == 0)
```

```
{
```

```
    for (int i=0; i < N/p; i++)
```

```
{
```

```
    if (i == N/p - 1)
```

```
{
```

```
    send = arr0[i];
```

```
MPI_Send(arr0[i], my_rank + 1);
```

```
MPI_Recv(recv, my_rank + 1);
```

```
Compute arr1[i] = arr0[i] + arr0[i-1] + recv;
```

```
?.
```

```
else { if (i == 0)
```

```
{ arr1[i] = 2 * arr0[i] + arr0[i+1];
```

```
?.
```

```
} else { arr1[i] = arr0[i] + arr0[i-1] + arr0[i+1];
```

else if ( $\text{my\_rank} == P - 1$ ).

{

for (int i = 0; i < N/p; i++)

{ if (i == 0).

{

$\text{MPI\_Send}(\text{arr}_0[i], \text{my\_rank} - 1);$

$\text{MPI\_Recv}(\text{arr}_1[0], \text{my\_rank} - 1);$

Compute heat eqn

$\text{arr}_1[i] = \frac{\text{recv} + \text{arr}_0[i] + \text{arr}_0[i+1]}{3}$

}.

else {

if (i == (N/p) - 1)

{  $\text{arr}_1[i] = \frac{2 * \text{arr}_0[i] + \text{arr}_0[i-1]}{3};$

}

else {

$\text{arr}_1[i] = \frac{\text{arr}_0[i-1] + \text{arr}_0[i] + \text{arr}_0[i+1]}{3},$

}

}

}.

else {

// All the process in the middle.

for (int i=0; i < N/p; i++)

{

if (i == 0)

{

MPI\_Send (arro[i], my\_rank - 1);

MPI\_Recv (~~arr~~recv, my\_rank - 1);

arr[i] = arro[i] + recv + arr[i+1];

3

}

else if (i == (N/p) - 1)

{

MPI\_Send (arro[i], my\_rank + 1);

MPI\_Recv (~~arr~~recv, my\_rank + 1);

arr[i] = arro[i] + arr[i-1] + recv;

3

}

else {

arr[i] = arro[i] + arr[i-1] + arr[i+1];

3

}

}

Communication in Block-wise decomposition:-

$$\begin{aligned}\text{Total-commr} &= \text{first\_procnt} + \text{middle\_procnt} + \text{last\_procn.} \\ &= 1 + (P-2)*2 + 1 \quad (\text{Per iteration}) \\ &= 2 + (P-2)*2 = 2 + 2P - 4 = 2P - 2.\end{aligned}$$

Per link = 2. (on an Average) } per iteration.

Per Node = Every node on average performs 2 seconds of  
2 recv/s whereas 1st & last broken make 1 Send & 1 Recv

(Q1) 1-D. partitioning (horizontal):-

(a) In horizontal partitioning all the matrix is going to be cut  $\frac{N}{P}$  pieces horizontally.

If you have a  $N \times N$  matrix we cut into:

$P$  pieces of  $\left[ \frac{N \times N}{P} \right]$  and vector of size  $\frac{N}{P}$  of  $P$ -parts

for (int i=0; i <  $\frac{N}{P}$ ; i++)  
 {  
 for (int j=0; j < N; j++)

$N \times N$  sized matrix

$N \rightarrow \dots$  vector

$P \rightarrow$  Total no of Proc.

for (int j=0; j < N; j++)

{

$a[i][j] = \text{Gen}((\text{my\_rank}) * \frac{N}{P} + i, j)$

}  
 }  
 // curr-state vector

for (int i=0; i <  $\frac{N}{P}$ ; i++)

{ curr-state[i] = Gen((my\_rank) \*  $\frac{N}{P}$  + i);

}

// so Every Processor has  $\frac{N}{P}$  part of Vector and it requires  $(P-1) * N$  data from other processors to compute the matrix - vector multiplication..

So, the current process should send  $\frac{N}{P}$  part of vector to all other processes (Broad-cast).

and Recv. the  $\frac{N}{P}$  th part of the data from all other processes.

Compute the matrix - Vector multiplication locally and save the result.

Copy the required  $\frac{N}{P}$  th <sup>curr state</sup> vector to next state vector &

Broad cast  $\frac{N}{P}$  th vector to all other in next iteration.

Continue for 10 iterations.

Total-Comm<sup>n</sup> :- (Horizontal)

Q) Memory each node consists are :-

$\left[\frac{N}{P} \times N\right]$  → To store matrix

$\left[\frac{N}{P}\right]$  → Size Array to store a Vector

$(P-1) * \frac{N}{P}$  → Size Array to store Recv-values

Total-Comm<sup>n</sup>:

Every process perform 1 Send of  $\frac{N}{P}$  th data &  $(P-1) * \frac{N}{P}$  recv-data. So,

$$\left[ \frac{N \times (P-1)}{P} + \frac{N \times (P-1) \times P}{P} \right] = \text{Total-Comm}^n$$

Horizontal-Cut:-

Total Communication per Link = 1 Send of  $\frac{N}{P}$  + 1 Recv of  $\frac{N}{P}$   
} Since N/w is clique }

Communication per Node =  $(P-1)$  Sends of  $\frac{N}{P}$  +  $(P-1)$  Recv  
of  $\frac{N}{P}$ .

## (2.2) Vertical - Decomposition :-

In vertical decomposition we have  $N \times N$  matrix cut.

to  $\frac{N}{P}$  pieces Vertically.

// So, ~~Every~~<sup>Total</sup> matrix of  $N \times N$  is ~~chopped~~ locally Generated  
in to each process. of size  $(\frac{N}{P} \times N)$  and matrix of size  $(N \times \frac{N}{P})$   
and output of size "N" to store the output.

// Once the data is Generated Locally we perform the matrix

vector b/w  $(N \times \frac{N}{P}) \times (\frac{N}{P} \times 1)$  ~~process~~, the results are.

stored as column vector (n vector that contains  $(\frac{N}{P})^{th}$  part).

-ginal answer. So, to get right answer we need to  
perform All reduce. So that every process has the  
right vector of 1<sup>st</sup> iteration.

All reduce = {All Broadcast + Reduce on Each  
                }  =  =   
(P-1) Sends + (P-1) Recv.

So, Now we take the required (my\_rank)<sup>th</sup>  
Part to curr-state & perform next iteration.

Continue until 10 iterations

$$N \times N = 3 \times 3 \quad P = 3.$$

$$\text{Ex:- } \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = \begin{bmatrix} 1 \\ 4 \\ 7 \end{bmatrix} [1] \otimes \begin{bmatrix} 2 \\ 5 \\ 8 \end{bmatrix} [2]$$

$$\begin{bmatrix} 3 \\ 6 \\ 9 \end{bmatrix} [3]$$

$P_1$

$P_2$

$P_3$ .

$$\begin{bmatrix} 3 \\ 3 \\ 3 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 1 \\ 4 \\ 7 \end{bmatrix} \begin{bmatrix} 4 \\ 10 \\ 16 \end{bmatrix} \begin{bmatrix} 9 \\ 18 \\ 27 \end{bmatrix}$$

After all-reduce Every process contains

$$\begin{bmatrix} 14 \\ 32 \\ 50 \end{bmatrix}$$

Total - Comm<sup>n</sup> :-

Memory :-  $(\frac{N \times N}{P})$  matrix.

$(\frac{N}{P})$  → size vector to store curr-state

$N$ -sized vector → Next state result  $\left( (\frac{N \times N}{P}) + (\frac{N}{P}) \right)$

$$O\left(\frac{N^2}{P} + \frac{N}{P} + N\right) = O\left(\frac{N^2}{P}\right)$$

Communication:  $P \times (P-1)$  sends of  $\frac{N}{P} + (P-1) \text{Recv } \frac{N}{P}$

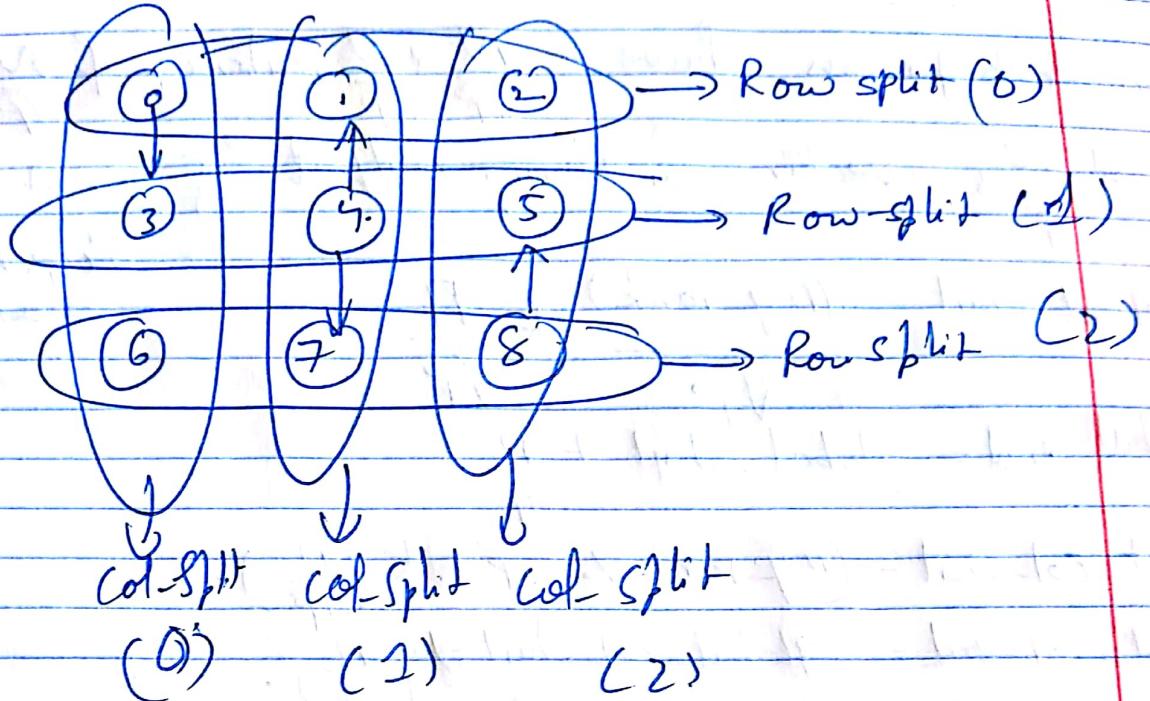
$$= P \times ((P-1)\frac{N}{P} + (P-1)\frac{N}{P}) = 2 \times (P-1) \times P \times \frac{N}{P}$$

$$= 2 \times (P-1) \times P \times \frac{N}{P}$$

Per link = 1 send  $\frac{N}{P}$  & 1 Recv  $\frac{N}{P}$ .

Per Node =  $(P-1)$  send  $\frac{N}{P}$  &  $(P-1)$  Recv  $\frac{N}{P}$ .

2.3 In Block wise decomposition.



① In Block-wise decomposition we store  $N \times N$  matrix in

different processor with a matrix of size  $\frac{N}{\sqrt{P}} \times \frac{N}{\sqrt{P}}$  and a vector of size  $\frac{N}{\sqrt{P}}$  & output vector of size  $\frac{N}{\sqrt{P}}$ .

② We compute local matrix vector multiplication of  $\frac{N}{\sqrt{P}}$ .

// Essentially we compute global row\_start & end, global start & col-end.

$$\text{global\_row\_st} = (\text{my\_rank}) / \text{processors in each colm}$$

## Block-decomposition:-

In each process have a  $\frac{N}{\sqrt{P}} \times \frac{N}{\sqrt{P}}$  matrix, &  $\frac{N}{\sqrt{P}}$  size vector and  $\frac{N}{\sqrt{P}}$  size vector for storing output. { genA is a function generated by user }

$$\text{global\_start} = (\text{my\_rank}) * \frac{N}{\sqrt{P}}$$

$$\text{global\_end} = \text{global\_start} + \frac{N}{\sqrt{P}}$$

$$\text{global\_col\_start} = (\text{my\_rank} \% \sqrt{P}) * \left(\frac{N}{\sqrt{P}}\right)$$

$$\text{global\_col\_end} = \text{global\_col\_start} + \left(\frac{N}{\sqrt{P}}\right)$$

1 Based on global-

for (int i = global\_start; i < global\_end; i++)

{ for (int j = global\_col\_start; j < global\_col\_end; j++)

{ grid [k] = genA (row, col); }

} }

by generate a vector.

Perform a row wise split based on  $(\text{my\_rank} / \frac{N}{\sqrt{P}})$

and col wise split based on  $(\text{my\_rank} \% \frac{N}{\sqrt{P}})$

so, perform All reduce in row\_split\_wise

so, for next iteration we do column wise Broadcast from the diagonal Processors for next iteration continue for 10 iterations

Total comm:

Memory:-  $\left(\frac{N}{\sqrt{P}} \times \frac{N}{\sqrt{P}}\right)$  matrix &  $\frac{N}{\sqrt{P}}$  vector for Giv & out put

$$= O\left(\frac{N^2}{P} + \frac{N}{\sqrt{P}} + \frac{N}{\sqrt{P}}\right) = O\left(\frac{N^2}{P}\right).$$

Total Comm<sup>n</sup> = row-wise + col-wise.

If Row-wise comm<sup>n</sup> performs full reduce }

$$\text{row-wis} = \left[ \underbrace{\left( \frac{N}{\sqrt{P}} \times (\sqrt{P}-1) \right) * \sqrt{P} + \left( \frac{N}{\sqrt{P}} \times (P-1) \right) * \sqrt{P}}_{\text{Sends.}} \times \sqrt{P} \right] \underbrace{\downarrow}_{\text{1 row}} \underbrace{\downarrow}_{\sqrt{P} \text{ row}}$$

Col-wise =

$$\text{Column takes} = \left( \underbrace{\frac{N}{\sqrt{P}} \times (\sqrt{P}-1)}_{\text{Sends.}} + \underbrace{\frac{N}{\sqrt{P}} * (\sqrt{P}-1)}_{\text{Recv.}} \right) \times \sqrt{P}$$

$$\therefore \text{Total Comm}^n = \text{row-wis} + \text{col-wis}.$$

$$= 2N * \sqrt{P}(\sqrt{P}-1) + 2N * (P-1)$$

$$= 2N(\sqrt{P}-1) \cdot (P+1) = 2N * (P-1)$$

On-link: 1 Send  $\frac{N}{\sqrt{P}}$  & 1 Recv  $\frac{N}{\sqrt{P}}$ .

most-loaded node: The central processor in above example at (4) in  $3 \times 3$  grid of proc.

Performs  $(\sqrt{P}-1) * \frac{N}{\sqrt{P}}$  recv/s of  $\frac{N}{\sqrt{P}}$  data

&  $(\sqrt{P}-1)$  sends of  $\frac{N}{\sqrt{P}}$  data

If we perform only diagonal reduce in processors

instead of all reduce we can

avoid lot of communication.

1 row can have  $\left( \frac{N}{\sqrt{P}} * (\sqrt{P}-1) \text{recv} + \frac{N}{\sqrt{P}} \text{send} \right)$

Total row Comm =  $\left( \frac{N}{\sqrt{P}} * (\sqrt{P}-1) + \frac{N}{\sqrt{P}} \right) \sqrt{P}$ .

Column =  $\left( \frac{N}{\sqrt{P}} * (\sqrt{P}-1) \text{send} + \frac{N}{\sqrt{P}} * (\sqrt{P}-1) \text{recv} \right) * \sqrt{P}$

$$= \frac{2N}{\sqrt{P}} * (\sqrt{P}-1) * \sqrt{P} + N(\sqrt{P}-1) + N.$$

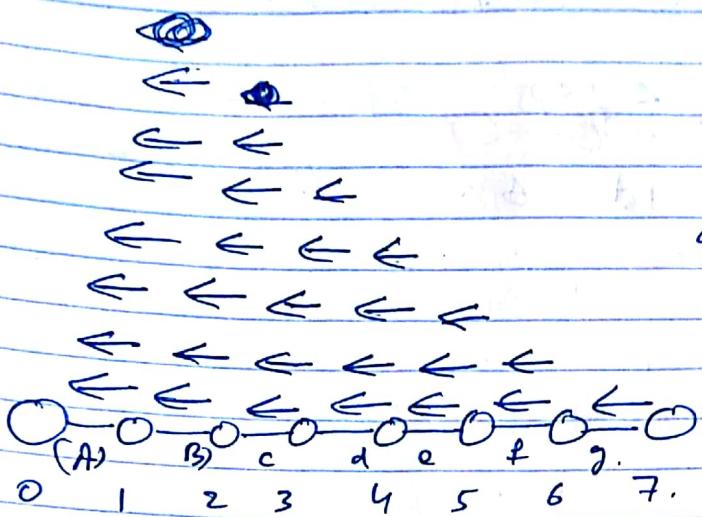
$$O(N\sqrt{P}) = 2N(\sqrt{P}-1) + N(\sqrt{P}-1) + N$$

$$= N(2\sqrt{P}-2 + \sqrt{P} - 1 + 1) = N(3\sqrt{P} - 2)$$

(3) Reduction:-

Reduce-Star on-chain :-

In Reduce Star all the results are send to zero.  
Let say  $P=8$ .



In the figure beside.

The most loaded link is  
(A) with 7 messages  
In terms of  $P$ ,  
 $(P-1) * 1$ .

The most loaded node.  
'1' ft is perform.

6 sends & 6 Recv

" $P-2$ " sends & " $P-2$ " Recv

$\Rightarrow 2P-4$  communication

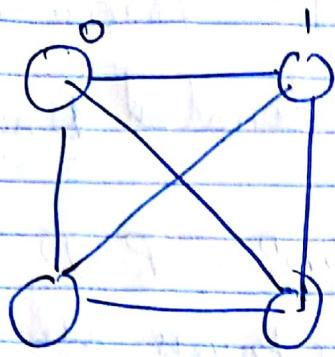
The longest chain of Comm happen from Processor 0 to proc 7

The " " " " " is  $\frac{(P-1)}{2}$  in above ex :-

$$\text{Total Comm} = 1 + 2 + 3 + \dots + \frac{P-1}{2} = \frac{P(P-1)}{2} = O(p^2)$$

Reduce-star on clique:-

Since the underlying N/w is a clique.



Every linknode has a link to Every other processor or.

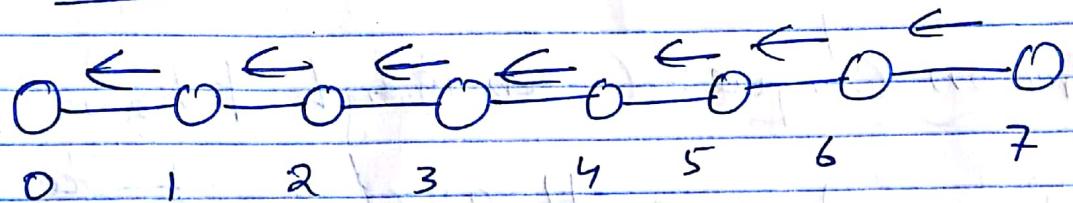
$$\text{Total Comm}^n = P - 1$$

bcz Every link has one value to send,

most loaded node is zero. because.

It is receiving  $(P-1)$  recv from all other processor  
Comm<sup>n</sup> is  $O(P)$       longest Comm<sup>n</sup> is 1

Reduce-chain:- {N/w is chain}



Every node performs Exactly one send & one recv. Except the last & first nodes. (because the last node did one send)

first node recv one value.

$$\text{Total Comm}^n = P - 1. = O(P).$$

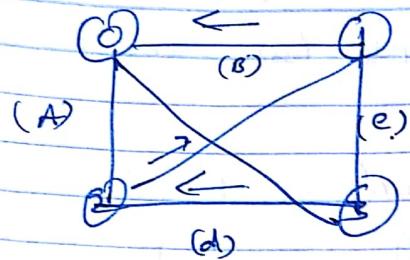
$$\text{longest chain of Comm}^n = P - 1 \text{ i.e } (0 - 7)$$

just look at the arrows.

data on the most loaded link = 1 data throughout the N/w on each Link.

data on most loaded Node =  $[P - 6]$   $\underbrace{[1 - P - 2]}_{\text{rank}}$   
Processor perform 1 Send & 1 Recv.

## Reduce-chain. (clique)-



here 1, 4 & 2 Process are  
the most loaded ones

$$\text{Total Comm}^n = (P-1) = O(P)$$

because Every process is sending 4 to  
data to its previous process.

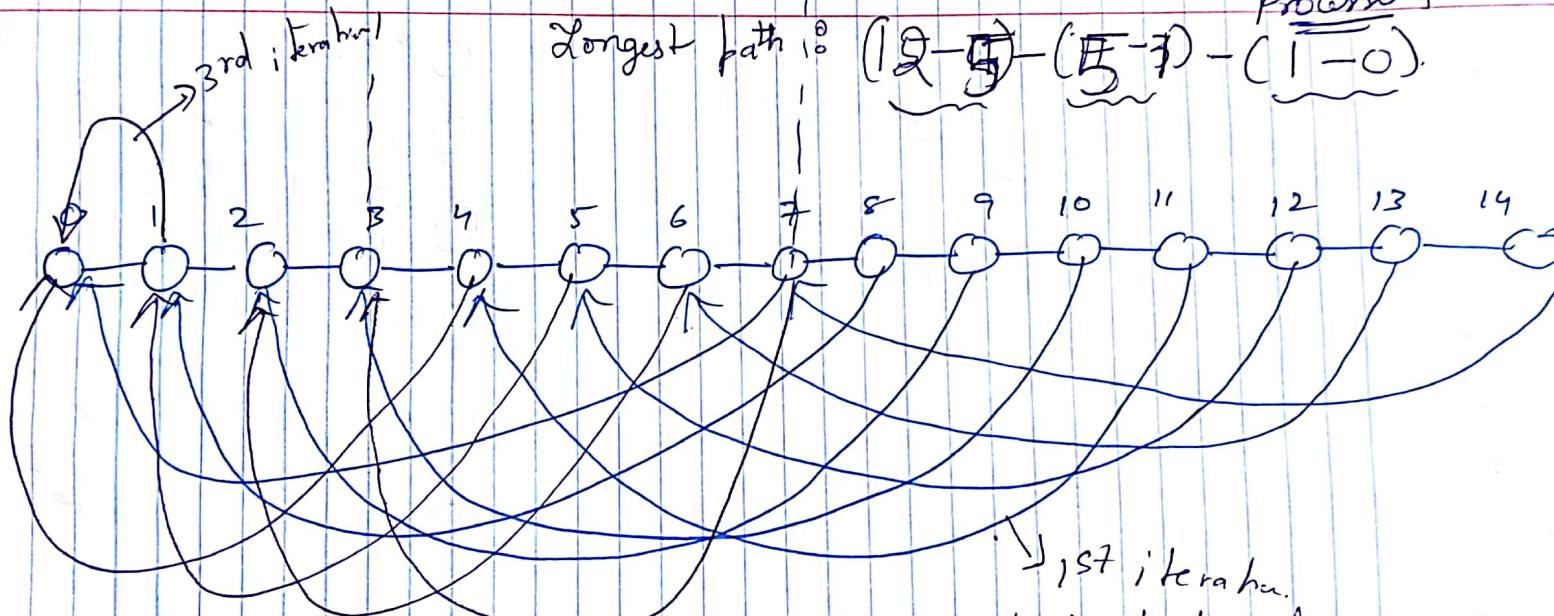
$$\text{longest Comm}^n = P-1 = O(P).$$

most loaded link. : 1 send.

most loaded nodes :- [1 - - - P-2] ranks performs 1 send  
of 1 Recv.

Reduce Tree on clique

(Even though  $\frac{N}{w}$  in clique), it's hard to express on 16 processors



Every process has 9 links to others

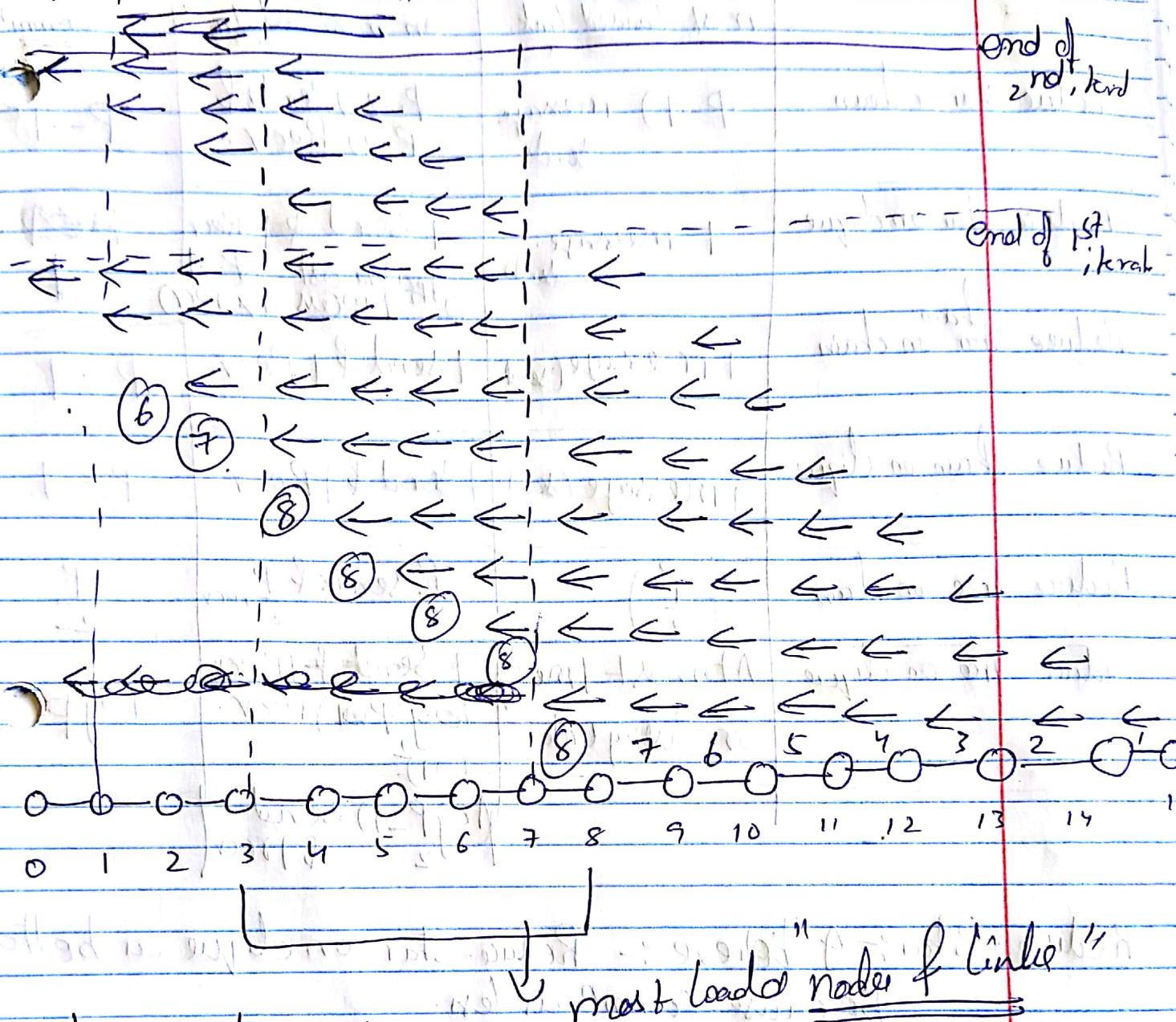
2nd iteration  
processors.

data on most loaded link is 1 send (one message).

data on most loaded node is  $\left[ \begin{smallmatrix} 1 \\ 0 \end{smallmatrix} \right]$  rank processor &  $\left[ \begin{smallmatrix} 0 \\ 1 \end{smallmatrix} \right]$  ranker

where  $n = \text{iterations}$  they perform  $2 \text{ recv } 2 \text{ & one send}$ .  
Critical path =  $O(\log P) =$  here it is  $3 = (12-5) - (5-1) (1-0)$

Reduce Tree on chain :-



whose rank are  $\frac{P}{2} - \frac{P}{2}$  in above Example

longest chain length is  $\left(\frac{P}{2}\right)$ .

data on most loaded node is  $\left(\frac{P}{2}\right)$ .

### Table:-

|                         | How much data on most loaded link. | Show much data on most loaded node                                    | Show long is commn |
|-------------------------|------------------------------------|-----------------------------------------------------------------------|--------------------|
| Reduce-star on chain    | (P-1) message sends                | (P-2) sends & (P-2) Recv.                                             | (P-1)              |
| Reduce-star on clique.  | 1 message send                     | 1 send & 1 Recv.<br><del>that most P-1 processors recd</del>          | <del>P-1</del>     |
| Reduce-chain on chain.  | 1 message (send)                   | 1 send & 1 Recv.                                                      | P-1                |
| Reduce chain on clique. | 1 message (send)                   | 1 send & 1 Recv.                                                      | P-1                |
| Reduce Tree on chain    | $\left(\frac{P}{2}\right)$         | $\frac{P}{2}$ sends & $\frac{P}{2}$ recv                              | $\frac{P}{2}$      |
| Reduce Tree on clique.  | Atmost 1 message on every link.    | + send & +recv.<br>$(\log_2 P+1)$ recv.<br>$(\log_2 P)$ send & 1 recv | $\log P$           |

Reduce-star :- I believe :- Reduce-star on clique is better because Commn is less.

Reduce-chain I believe : Reduce-chain on chain is better because clique is very large N/w in real. Time qts complex to model & handle.

Reduce-Tree: I believe :- Reduce-Tree on clique is better because The Commn is very less compare to Chain,