



| S. No | Topic  | Reference                          | Page/Section                           |
|-------|--|------------------------------------|--|
| 1.    | Introduction to parallel programming.  | Introduction to Parallel Computing | 1.1 to 1.2.4                           |
| 2.    | Programming models and parallel architectures, superscalar   | Introduction to Parallel Computing | 2.1                                    |
| 3.    | Why pipelining? Performance analysis using pipelining on a single processor.                               | Introduction to Parallel Computing | 2.1                                    |
| 4.    | Impact of memory latency on performance. How to increase memory bandwidth.                                 | Introduction to Parallel Computing | 2.2, 2.2.1                             |
| 5.    | Cache and memory bandwidth. Performance considerations using the cache.                                    | Introduction to Parallel Computing | 2.2.2                                  |
| 6.    | Shared and distributed memory architectures. Interconnection networks in distributed memory architectures. | Introduction to Parallel Computing | 2.3, 2.3.1, 2.3.2, 2.4.2, 2.4.3, 2.4.4 |

|     |  |                                    |        |
|-----|--|------------------------------------|--------|
| 7.  | What is OpenMP? Basic program using OpenMP.                                  | Introduction to Parallel Computing | 7.10   |
| 8.  | Programming using single thread.   | Introduction to Parallel Computing | 7.10.1 |
| 9.  | Parallel programming using multiple threads.                                 | IIT-Delhi Book                     | 7.10.1 |
| 10. | Context switching threads, thread functions.                                 | IIT-Delhi Book                     |        |
| 11. | Sequential consistency model   | IIT-Delhi Book                     |        |
| 12. | Critical section, race conditions  | IIT-Delhi Book                     |        |
| 13. | Handling of race conditions using openmp                                     | IIT-Delhi Book                     |        |
| 14. | OpenMP scoping of variables  | IIT-Delhi Book                     |        |
| 15. | Thread private variables   | IIT-Delhi Book                     |        |
| 16. | Manual distributions of work and critical section                            | IIT-Delhi Book                     |        |
| 17. | Loops and reductions   | IIT-Delhi Book                     |        |
| 18. | Vector-Vector operations, matrix-vector operations, matrix-matrix operations | IIT-Delhi Book                     |        |
| 19. | Tasks  | IIT-Delhi Book                     |        |
| 20. | Variable accessing in task   | IIT-Delhi Book                     |        |
| 21. | Completion of tasks and scoping variables in tasks                           | IIT-Delhi Book                     |        |
| 22. | Recursive task spawning  | IIT-Delhi Book                     |        |
| 23. | Message passing interface  | From class discussion              |        |
| 24. | Basic MPI calls  | From class discussion              |        |
| 25. | Broadcast implementation using MPI   | From class discussion              |        |
| 26. | MPI Non blocking calls   | From class discussion              |        |
| 27. | MPI Collectives and MPI broadcast  | From class discussion              |        |
| 28. | Characterization of interconnects  | From class discussion              |        |
| 29. | Hockney model  | From class discussion              |        |
| 30. | Broadcast, Reduce, Scatter, Gather, Reduce Scatter, AlltoAll                 | From class discussion              |        |
| 31. | Introduction to GPU architectures  | Slides and class discussion        |        |
| 32. | Data parallel algorithms, SIMD, SIMT,  | Slides and class discussion        |        |
| 33. | Introduction to CUDA programming   | Slides and class discussion        |        |
| 34. | Programme implementation using CUDA, Kernel launch                           | Slides and class discussion        |        |
| 35. | Grid, blocks, CUDA functions   | Slides and class discussion        |        |
| 36. | Multidimensional mapping of data space synchronization                       |                                    |        |
| 37. | Histogram  |                                    |        |

# OpenMP



→ multiprocessing

1. Hello world in C

multiple threads

```
#include <omp.h>
```

```
#include <stdio.h>
```

```
int main( int argc, char *argv[] ) {
```

/\* whenever we need a block to be  
executed in parallel we have  
to enclose it in #pragma omp

compiler directive

header file  
(  
containing  
prototypes  
for various  
OpenMP func)

```
#pragma omp parallel {
```

```
    printf("Hello world");
```

```
}
```

```
return 0;
```

```
}
```

To compile & run →

```
$ gcc -fopenmp HelloWorld.c
```

```
$ ./a.out
```

To define no of threads

```
$ export OMP_NUM_THREADS = 4
```

\$ ./a.out

env variable

# \* Threads

single stream of control in the flow of program

Note - Understanding memory of a Program

int i = 4;

int j;

fun main {

    int n;

    160 → g(n);

101 →

}

120 → g(int n) {

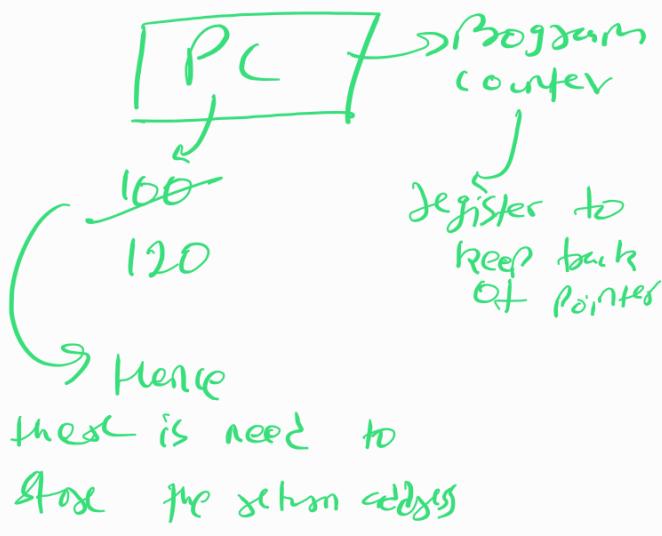
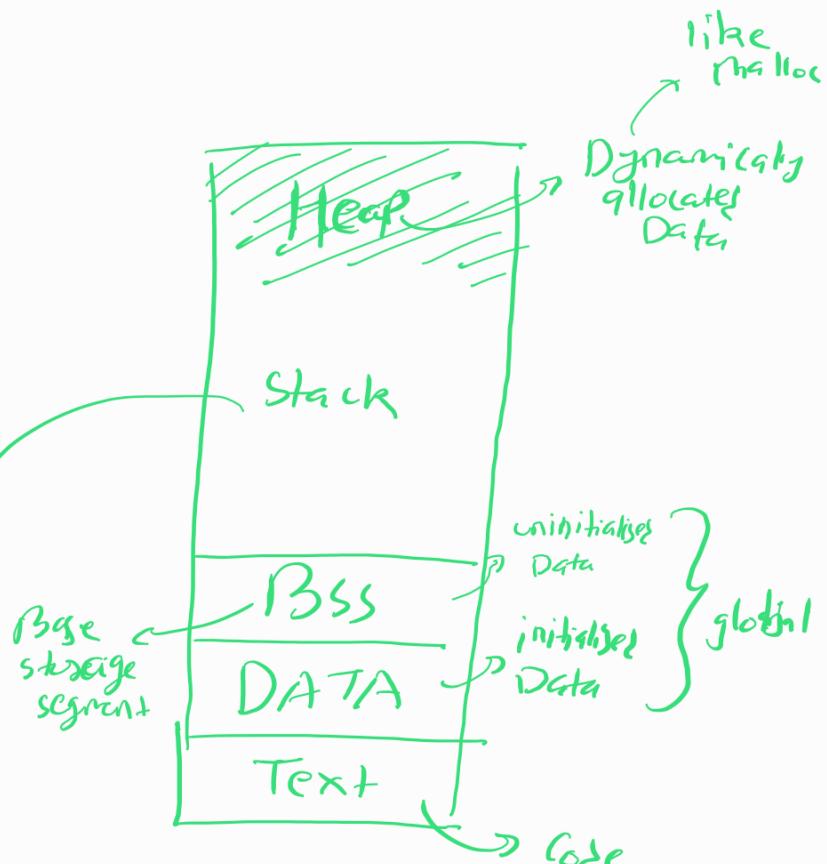
    int y;

    :

}

lifetime of  
y is =  
lifetime of  
g(n)

at every call of  
g another frame  
is created  
independent  
from previous



Now for AL Program

↳ we can run multiple  
strings of code flows  
in

multiple threads

for doing so the main issue is  
memory of program will be stack as  
each thread will have its own control of  
flows



Here each thread has its  
own stack

most devices are multitasking



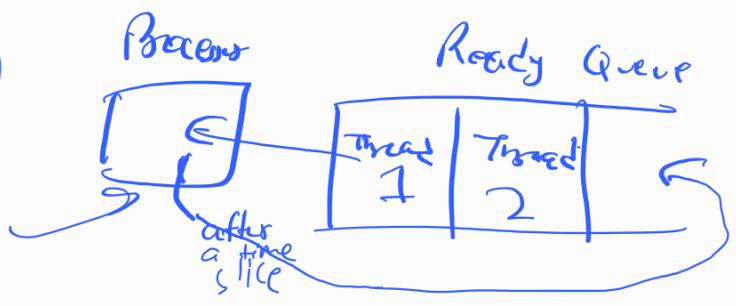
Thread 1 is scheduled

for a time slice

then Thread 2, doesn't

need to complete Thread 1

The scheduling  
is done by  
Ready Queue



## Context Switching



Process by which an OS saves the current state ("context") of the executing thread or Process & loads the previously saved state of another

whenever a context is switched a copy is created of the registers and is stored in thread local space

On single core  $\rightarrow$  context switching creates an illusion of concurrency

On multi-core  $\rightarrow$  true parallelism is possible - multiple threads can run truly each on its own core, context switching can still happen if needed

# \* Open MP basic thread func

```
#include <omp.h>
#include <stdio.h>
int main( int *argc, char **argv[ ] ) {
    #Pragma omp parallel {
        /* get the total number of OMP threads
        int numt = omp_get_num_threads(); */
        /* get the thread id of this thread.
        int tid = omp_get_thread_num(); */
        printf("HelloWorld from thread %d of %d, tid,
        numt);
    }
    return 0;
}
```

Note - the order isn't same at every run in the output, but execution takes place sequentially to make it atomic we have to take help of locks & mutual exclusion.

Note - if numt & tid is defined outside the `for` block.

```
int main (int *argc, char *argv[]) {
```

```
    int tid, numt;
```

```
#pragma omp parallel {
```

```
    numt = omp_get_num_threads();
```

```
    tid = omp_get_thread_num();
```

```
    printf("HelloWorld from thread %d of %d, tid, numt);
```

```
}
```

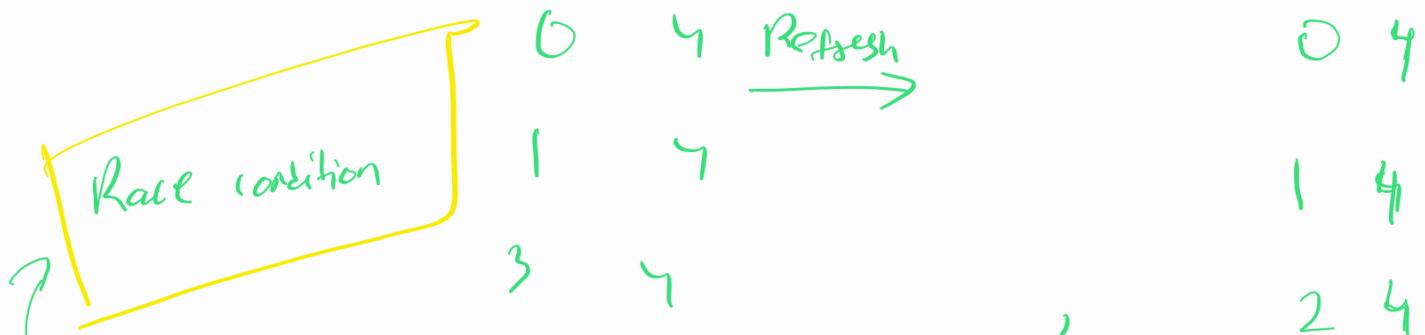
```
dtyn 0;
```

```
}
```

Output →

Hello world from thread 2 of 4

Hello world from thread 1 or 4



(change in output thread 3)  
didn't print anything, what

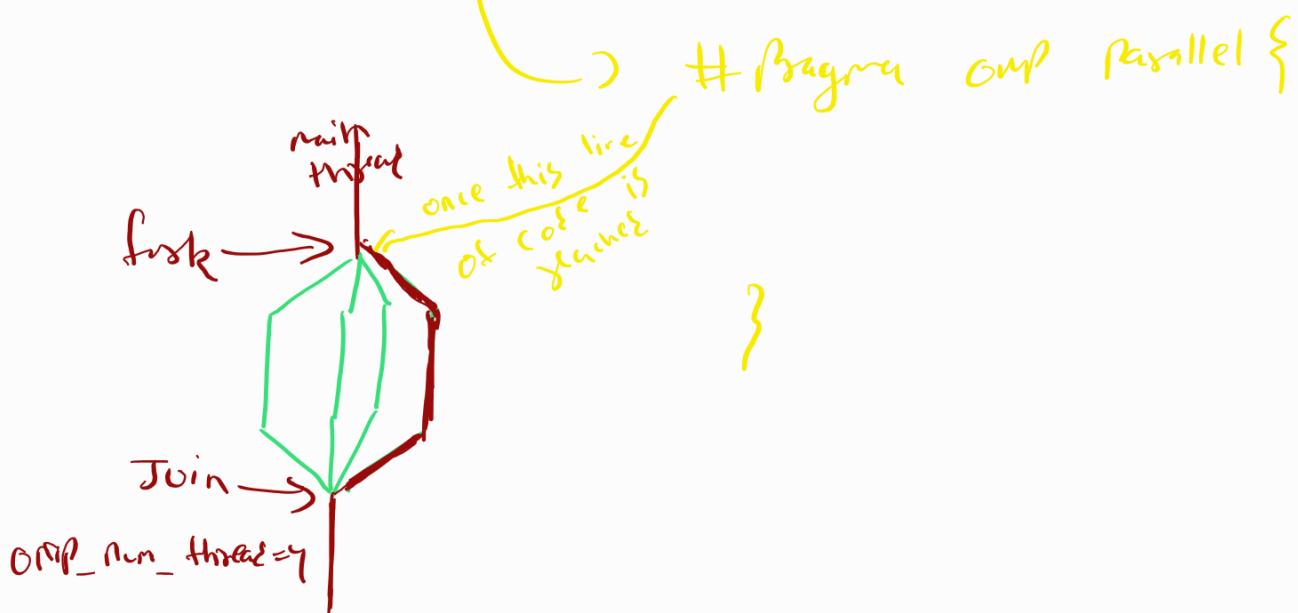
actually happened is thread 3's output was overwritten by thread 1

# \* About OpenMP

- API for shared memory model Programming
  - compiler directive (e.g. `#pragma omp parallel`)
  - Runtime Library routines (e.g. `omp_get_num_threads`)
  - Env variables (e.g. `OMP_NUM_THREADS`)
- Adv -
  - Standardization & Portability
  - Ease of use

- Thread Based
  - threads share memory space except stacks

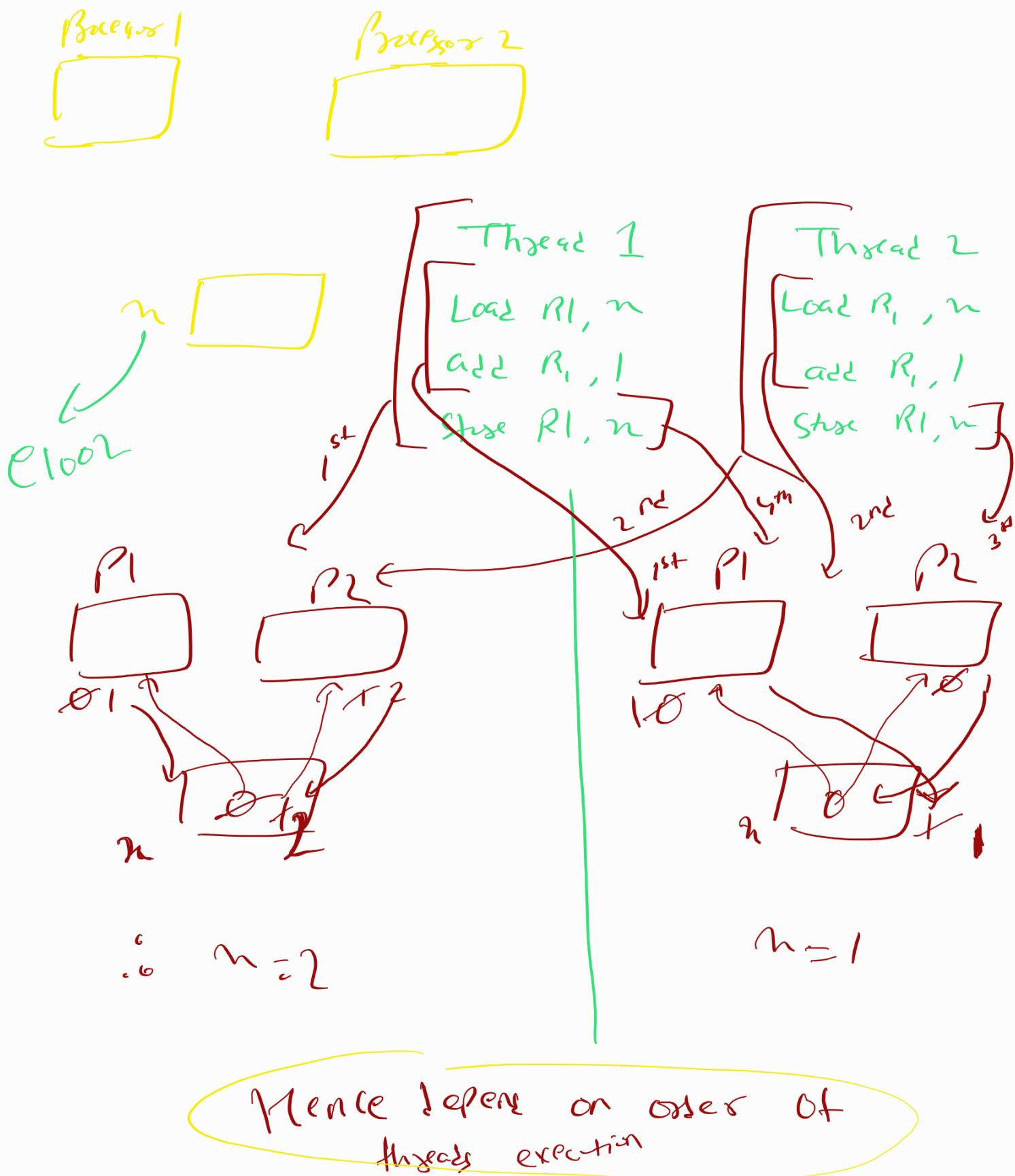
- Fork-Join model



## \* Race condition -

If correctness of the PL program depends on the order in which threads are scheduled.

e.g.



To fix Race condition -

```
eg- int main(---){  
    int numt, tid;  
    /*#pragma omp parallel {  
        numt = omp_get_num_threads();  
        tid = omp_get_thread_num();  
        printf("Hello --- , tid, numt);  
    }  
    return 0;  
}
```

int main(---){  
 int numt, tid;  
 /\* specify Private & Shared variable  
 #pragma omp parallel  
 Private(tid) Shared(numt){  
 numt = --- ,  
 tid = --- ;  
 printf(" --- , tid, numt);  
 }  
 return 0;  
}

note- default behavior for variable in omp is shared.

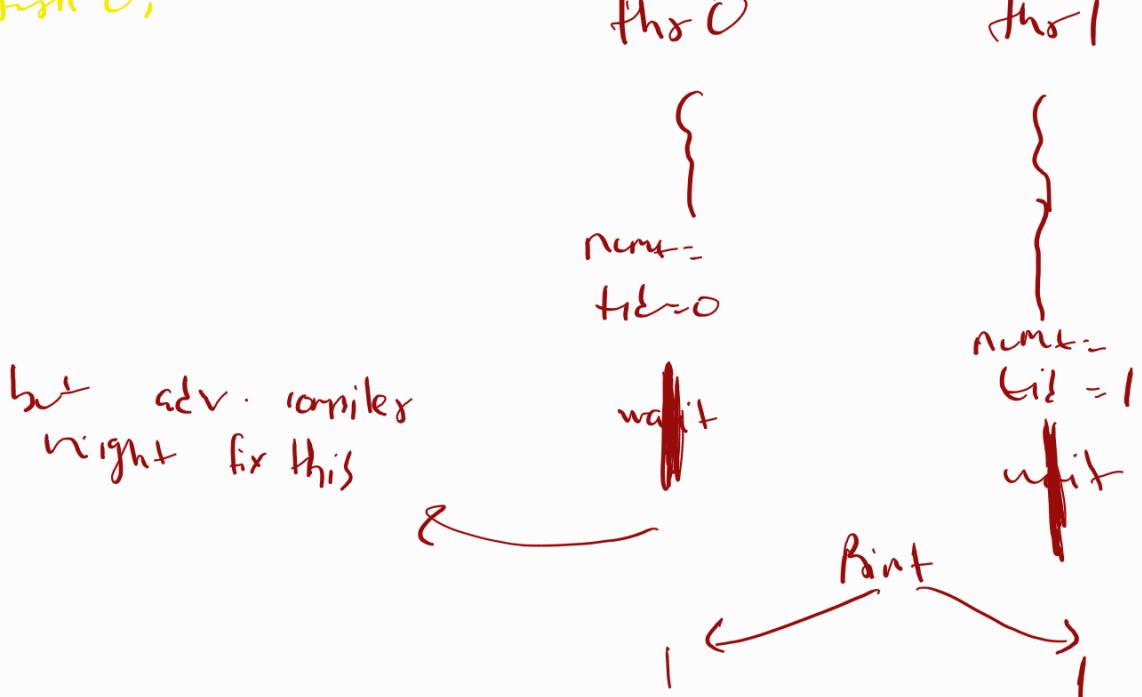
means every thread has its own tid  
variable is global & everyone access safe location

Note - still how can we make sure that  
the bug is fixed for sure.



By running the code with shared flag fail  
(almost everytime),

```
int main ( - - - - ) {  
    int numt, tid;  
    #pragma omp parallel {  
        numt = omp_get_num_threads();  
        tid = omp_get_thread_num();  
        /* intialize wait so other thread get scheduled .  
           for (j=0->100000000; j++) ;  
        Pint(Hello - - - - , tid, numt);  
    }  
    return 0;  
}
```



but adv. compiler  
might fix this

so now for fixing + for size check

```
int main ( - - - ) {  
    int nunt, tid;  
    #pragma omp parallel private(tid) shared(nunt){  
        nunt = omp_get_num_threads();  
        tid = omp_get_thread_num();  
        /* inform each so other thread get scheduled.  
         * for (j=0 → 1000000000; j++) ;  
         * print(Hello - - - , tid, nunt);  
    }  
    return 0;  
}
```

Now as we can observe that `nunt` is a shared resource hence, it's getting called multiple times, so let's try to define it globally

```
int ran ( - - - ) {  
    int tid;  
    int numt = omp_get_num_threads();  
    #pragma omp parallel private(tid)  
    tid = omp_get_thread_num();
```

```
    Pitt(Hello - - - , tid, nunt);  
}  
return 0;  
}
```

Output →

Hello - - . Thread 2 of 3  
3 of 1  
0 of 1  
1 of 1

bcz when defined globally the no of thread was only 1

Now to avoid calling `numt` multiple times  
and at the same time not defining it  
globally below is the soln -

↳ Only one thread can call `numt`

```
int main ( - - - ) {  
    int tid; numt
```

# pragma omp parallel default (shared)

tid = omp\_get\_thread\_num();

/\* Initialise numt if tid == 0

if (tid == 0) {

numt = omp\_get\_num\_threads();

}

printf(Hello - - - , tid, numt);

}

return 0;

}

Output →

HelloWorld from thread 0 of 7

1 of 4 <sup>done again</sup> --- of 7

2 of 4 . . . of 7

3 of 4 - - - of 7

garbage value

- - - 0 of 7

--- 1 of 7

. . . 2 of 7

- - - 3 of 7

to make it fail everytime  
(almost)

```
int main ( - - - ) {  
    int tid; numt
```

# pragma omp parallel default (shared)

```
    tid = omp_get_thread_num();
```

/\* Initialise numt if tid == 0

```
    if (tid == 0) {
```

/\* introduce wait so other threads get scheduled

```
        for (j=0 → 100000; j++) {
```

```
            numt = omp_get_num_threads();
```

```
        }
```

```
        Pmt(Hello - - - - , tid, numt);
```

```
}
```

```
return 0;
```

```
}
```

so now for fixing + size check

```
int main( -- -- ) {
```

```
    int numt;
```

```
#pragma omp parallel default (shared) {
```

```
    tid = omp_get_thread_num();
```

```
    if (tid == 0) {
```

```
        for (j = 0 → 1000000; j++) {
```

```
            numt = omp_get_num_threads();
```

```
}
```

```
}
```

```
#pragma omp parallel default (shared) {
```

```
    int tid = omp_get_thread_num();
```

```
    cout << "Hello --- "
```

```
}
```

```
return 0;
```

```
}
```

↳ not elegant

her calls break  
able to getue calls.

On other way to solve multiple calls



Thread Private Variable



Making TID persistent  
across all sections.

directive used to give  
thread its own instance  
of a global variable &  
keep that instance persistent  
across multiple TL regions

Middle ground b/w Shared & Private

- each thread has its own copy of variable
- But the value persists across different parallel sections unlike private where it gets reinitialized

int tid;

#pragma omp threadprivate(tid)

int main( int \*argc, char \*argv[] ) {

    int numL;

    # pragma omp parallel default (shared) {

        tid = OMP\_get\_thread\_num();

        if ( tid == 0 ) {

            for ( j=0; j<10000000; j++ ) {

                numL = OMP\_get\_num\_threads();

            }

        }

    # pragma omp parallel default (shared) {

        printf ("Hello world %d.%d.%d,%d\n", i,

        j, numL,

        numL );

    }

Note - we usually avoid  
using implicit barriers

using #pragma omp parallel  
as barriers

```
int tid;
#pragma omp threadprivate(tid)
int main( int *argc, char *argv[] ) {
    int numL;
    #pragma omp parallel default (shared)
        tid = OMP_get_thread_num();
        if ( tid == 0 ) {
            for ( j=0; j<10000000; j++ ) {
                numL = OMP_get_num_threads();
            }
        }
    # pragma omp barrier
```

printf ("Hello world %d.%d.%d,%d\n", i,

    j, numL,

Here, there's  
a compiler directive  
for safe use



printf ("Hello world %d.%d.%d,%d\n", i,

    j, numL,

if we want only one thread to perform any activity in a synchronised manner

#pragma omp single

```
int main ( int * argc , int * argv [ ] ) {
```

```
    int numt;
```

#pragma omp parallel default (shared)

```
    int j , tid;
```

#pragma omp single {

```
    for (j = 0 → 100000000) {
```

```
        numt = omp_get_num_threads();  
    }
```

```
    tid = omp_get_thread_num();
```

```
    printf (" How %d of %d , tid , numt );
```

```
}
```

if

other threads  
can go forward

#pragma omp single nowait {

```
    for (j = 0 → 100000000) {
```

```
        numt = omp_get_num_threads();  
    }
```

if

#pragma omp single master {

```
    for (j = 0 → 100000000) {
```

```
        numt = omp_get_num_threads();  
    }
```

only master  
thread is  
executed  
tid = 0

# \* Shared memory consistency models

memory consistency rules → set of rules which defines the order in which reads & writes

hardware & compiler might reorder the code for optimization

conflict b/w programmers & system

## 1. Sequential consistency model -

- operation of all processor appears in some global order that is consistent with each processor's local order of operation.
- very simple to reason
- expensive to implement
- Processors & compiler can't reorder memory ops freely.

## 2. TSO (Total store order)

- 3. release consistency
- 4. weak consistency model

Note -

omp\_get\_wtime() = time in second  
omp\_get\_wtick() = resolution of timer

## I. Computing Sum

### Sequential code

```
#include <omp.h>
#include <stdio.h>

#define ARR_SIZE =
    1000000000

int a[ARR_SIZE];

int main( int *argc, char *argv[] ) {
    double t1, t2;
    sum = 0;
    for( int i=0 → ARR_SIZE ) {
        arr[i] = 1;
    }
    t1 = omp_get_wtime();
    for( i=0 → ARR_SIZE ) {
        sum += arr[i];
    }
    t2 = omp_get_wtime();
    printf ("sum = %d : time = %.9f", sum, t2-t1)
    return 0;
}
```

## Parallel code -

```
#include <omp.h>
#include <stdio.h>

#define ARR_SIZE=
    1000000000
int a[ARR_SIZE];
int main( int *argc, char *argv[] ) {
    int sum = 0; int numt, tid;
    double t1, t2;
    for( int i=0 → ARR_SIZE ) {
        a[i] = 1;
    }
    t1 = omp_get_wtime();
}
```

```
# Bragzzi omp parallel private(tid) {
```

```
    tid = omp_get_thread_num();
    numt = Omp_get_num_threads();
    for ( int i=0 → ARR_SIZE ) {
        sum += a[i];
    }
}
```

```
t2 = omp_get_wtime();
```

```
printf("sum = %d · time = %g", sum, t2-t1);
```

```
return 0;
```

```
}
```

asking each  
thread to do  
some task

## Problems -

- Fundamental —  
not truly Parallel

- Race Condition  
in sum

→ sum > 1 billion

•  $t_0 > t_1$

→ cache  
coherency

expected sum =

4 billion

$4 \times 1 \text{ billion}$

default  
no. of threads

→ shared variable

Output →

sum = 1373461697 · time = 31.1403

# Fixing #1 code

```
int tid, numt, sum=0;  
double t1, t2;  
for(i=0→ARR_MAX) {  
    arr[i]=i;  
    }  
    t1 = omp_get_wtime;
```

```
#pragma omp parallel private(tid) {
```

```
int from, to;
```

```
tid = omp_get_thread_num();
```

```
numt = omp_get_num_threads();
```

```
from = (ARR_SIZE / numt) * tid;
```

```
to = (ARR_SIZE / numt) * (tid + 1) - 1;
```

```
i++ (tid == numt - 1) {
```

```
to = ARR_SIZE - 1;
```

```
}
```

```
for (i = from → to) {
```

```
#pragma omp critical {  
    sum += a[i];
```

```
}
```

```
printf(“ - - - . ”)
```

```
dmn 0;
```

```
}
```

fetch the  
value & update

so test & set(x)  
is used

• Fundamental  
problem solved  
now distribution of  
work ;

•  $t_p < t_s$

very important to  
ensure section is atomic  
so at one time only  
one thread can  
access.

a lock is  
maintained to do  
so

lock itself run  
take a while condition  
sometimes

Even after #pragma omp critical the  
time will be significantly hence  
we will introduce Psum (local variable)

```
#include<omp.h>
```

```
#include <stdio.h>
```

```
#Define ARR_SIZE=1000000
```

```
int a[ARR_SIZE];
```

```
int main(int argc, char *argv[]) {
```

```
    int tid, numt, sum=0;
```

```
    for (int i = 0 → ARR_SIZE) {
```

```
        a[i] = i;
```

```
}
```

```
# pragma omp parallel default(shared) private(tid)
```

```
    int Psum=0, from, to;
```

```
    tid = omp_get_thread_num();
```

```
    numt = omp_get_num_threads();
```

```
    from = (ARR_SIZE / numt) * tid;
```

```
    to = ((ARR_SIZE / numt) * (tid + 1)) - 1;
```

```
    for (i = from → to) {
```

```
        Psum += a[i];
```

```
}
```

```
# pragma omp critical {
```

```
    sum += Psum;
```

```
}
```

```

#include <stdio.h>
#include <omp.h>
#define ARR_SIZE = 1000000000
int a[ARR_SIZE];
int main(int argc, char *argv[]) {
    int , sum=0;
    for (int i=0 → ARR_SIZE) {
        a[i] = 1;
    }
}

```

```

#pragma omp parallel default(shared) {
    int psum=0;
    #pragma omp for
    for (i=0 → ARR_SIZE) {
        psum += a[i];
    }
}

```

```

#pragma omp critical {
    sum += psum;
}

```

↓ to reduce the overhead of  $psum$

$\uparrow$    
 $\#$  **Pragmas**  $\#$  **omp parallel default (shared)** Reduction ( $+=$  sum)

```

#pragma omp for
for (i=0 → ARR_SIZE) {
    sum += a[i];
}

```

gives each thread its private copy of sum

row to even slice the Barray one by

# Barray one parallel for default(shape)  
reduction(+:sum){}

```
for(i=0 → ARR_SIZE){  
    sum += a[i];  
}
```

## \* Linear Algebra

BLAS (Basic Linear Algebra Subroutines)

- i) BLAS-1      vector, vector - vector
- ii) BLAS-2      vector - matrix
- iii) BLAS-3      matrix - matrix

### I. Dot Product ( $y = A \cdot B$ )

sequential →

```
int n;  
dot = 0  
for (i=0 → n){  
    dot += A[i] * B[i];  
}
```

parallel →

```
int N;  
dot = 0;
```

# Barray one parallel by  
for (i=0 → N){  
 dot += A[i] + B[i];  
}

*fixed iteration time*

*slice of N iterations & hand them to threads in team*

*equal size blocks*

schedule(static/dynamic, chunk\_size)

## II. Matrix X Vector

sequential  $\rightarrow$

```

for (i=0 → n) {
    for (j=0 → n) {
        y[i] = A[i,j] * x[j];
    }
}

```

$$n \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}_A \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}_X^n = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}_y$$

$\mathcal{V} \rightarrow$

Note - C stores matrix in Row major order

Hence in C column is not contiguous

U  
false sharing situation

((cache line size is too large, after not storing it create illusion like sharing))

Sol<sup>A</sup>  $\rightarrow$  store in Multiple Variables

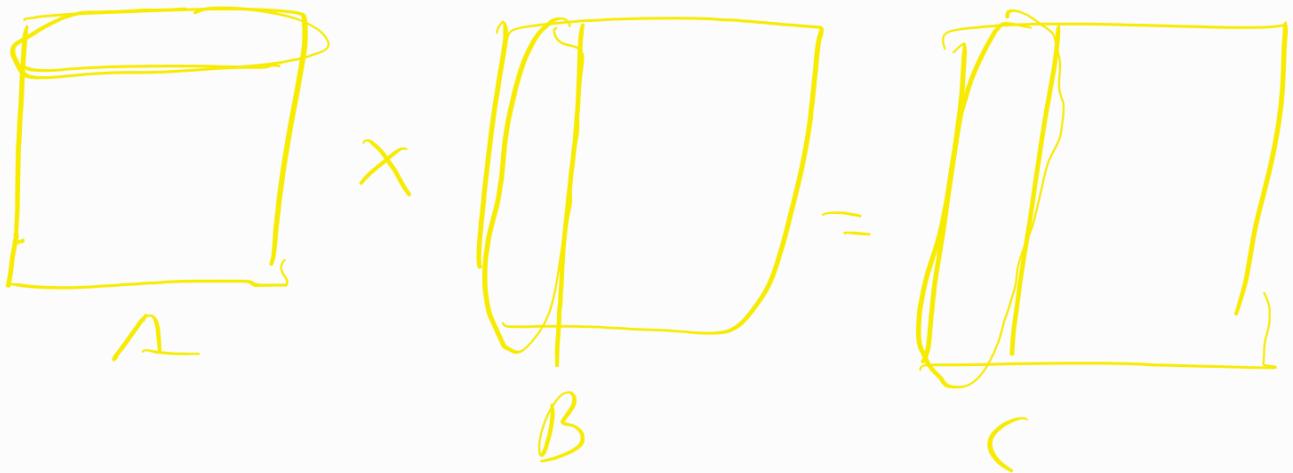
# Pdagra and Parallel for Schedulip (dynamic, chunk size)

```

for (i=0 → n) {
    for (j=0 → n) {
        y[i] += A[i,j] * x[j];
    }
}

```

### III Matrix - Matrix multiply



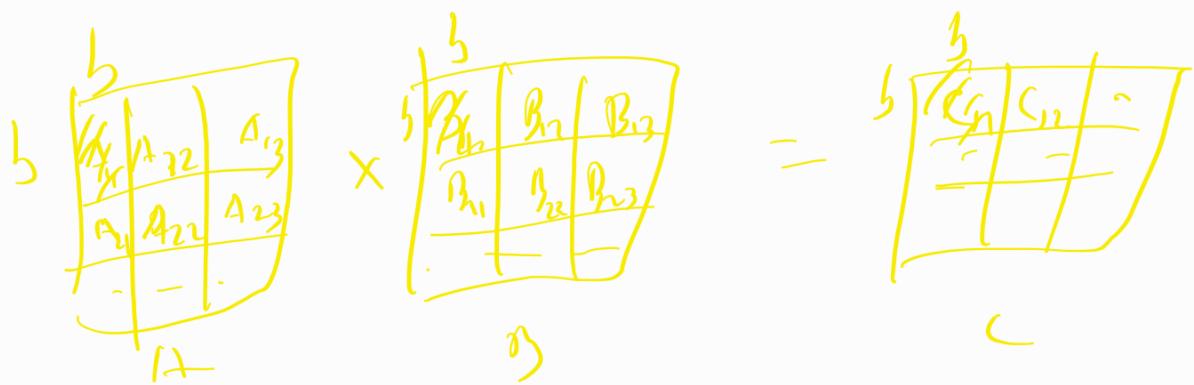
the column major in C won't be much of an issue but the column major in B is an issue and here we can't get benefit of data locality

#### Theoretical

$$\frac{\# \text{ compute ops}}{\# \text{ memory ops}} \approx \frac{2n^3}{3n^2} = O(n)$$

we want to be compute bound

so what we can do is divide the matrices in  $b \times b$  blocks



$$\frac{\# \text{ compute ops}}{\# \text{ memory ops}} \approx b$$

$$C_{11} = A_{11} \times B_{11} + A_{12} \times B_{21} + \dots$$

## \* Tasks

→ Piece of code that needs to be executed at some point of time.

e.g - counting array sum using task

```
sum = 0  
ARR_SIZE = 600
```

```
STEP_SIZE = 100
```

```
int a[ARR_SIZE];
```

#pragma omp parallel {

```
#pragma omp for { }  
for (i=0 → ARR_SIZE) { }  
start = i, end = i + STEP_SIZE - 1;  
printf( " - - - - " )
```

#pragma omp task {

```
int psum = 0
```

```
printf( " - - - . " )
```

```
for ( i=start → end ) { }
```

```
psum += a[i];
```

```
}
```

#pragma omp critical {

```
    sum += psum;
```

```
}
```

```
}  
printf( " ... " )
```

## Task Queue

Executed using Linked List

```
struct node {
    int data;
    Node *next;
}

// Program One Parallel {
    // Program one single {
        for (ptr= head; ptr != null; ptr = ptr->next)
            // Program one task {
                complete-inverse (ptr->matrix);
            }
    }
}
```

A single type  
out of the  
for loop type  
Program serial

Arte - To access variable of this task ~~the~~ here  
it's data env. associated to task major  
copy of all variables of thread or state.

Note - # pragma on parallel first private( )

(create copy of variable  
& initiate two variable  
just before this point)

Note - # pragma on for nowait()

generally, most of the func has  
an implicit barriers at the  
end of their code, nowait  
will & voids that barrier &  
start executing next lines.

Note - # pragma on barrier

for directives which doesn't have  
a implicit barrier & have  
nowait - to Put a Barrier  
Barriers.

Note - # pragma on task wait

it's not a  
barrier

wait for all tasks spawned by  
this task to complete

## Using task for recursive splitting

```
int do_sum(int start, int end){  
    int mid, x, y, res;  
    if (end == start) {  
        res = a[start];  
    } else {  
        mid = (start + end) / 2;  
    }
```

```
    printf (" [TID: %d] sum (%d, %d) =  
            sum (%d) + sum (%d, %d),  
            tid, start, end, start, mid,  
            mid+1, end);
```

# Pthreads omp task shared(x)  
 $x = \text{do\_sum}(\text{start}, \text{mid});$

# Pthreads omp task shared(y)  
 $y = \text{do\_sum}(\text{mid}+1, \text{end});$

# Pthreads omp taskwait  
}  $res = x + y;$

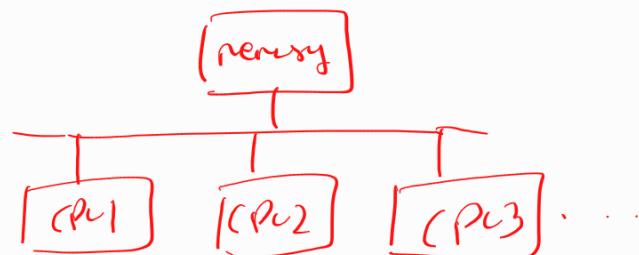
```
# define ARR_SIZE 3  
# define STEP_SIZE 1  
int a[ARR_SIZE];  
int main( int argc, int argv[]){  
    int i, sum=0;  
    for(i=0; i<ARR_SIZE; i++)  
        a[i]=i;  
    # Pthreads omp parallel  
    # Pthreads omp single  
    sum = do_sum(0, ARR_SIZE-1);  
    printf ("sum = %d", sum);  
    return 0;  
}
```

not a barrier just  
wait for tasks spawned  
by task to complete

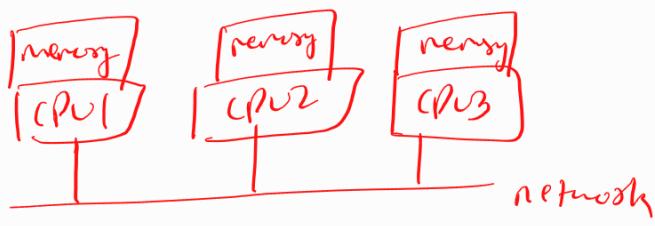
# Open MPI

## Parallel Architecture

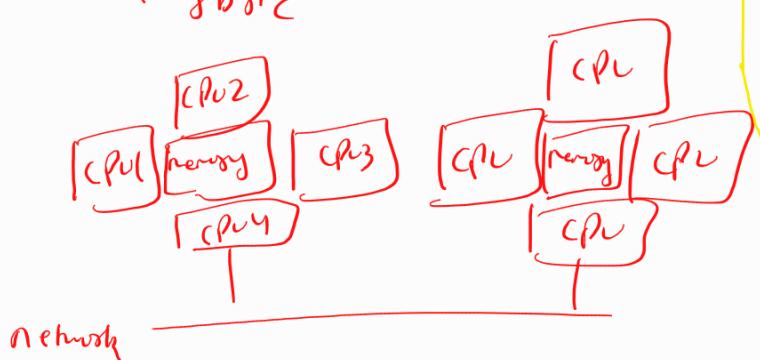
- Shared memory



- Distributed memory

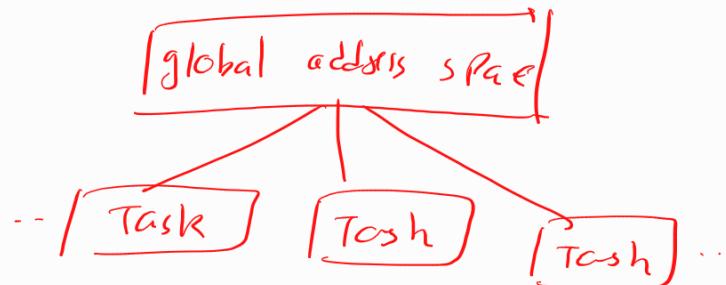


- Hybrid

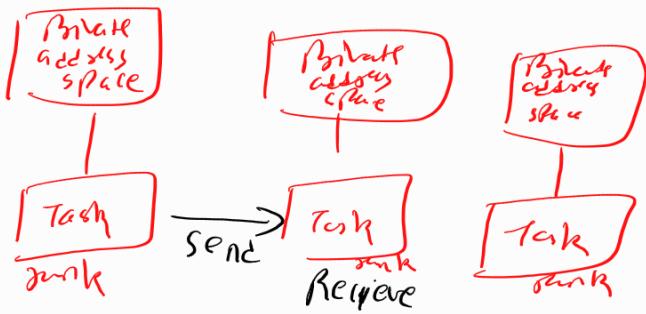


## & Programming models

- Shared memory  
Open MP



- Message Passing  
MPI



- message passing +

OpenMP + MPI



## Structure of MPI code

```
#include "mpi.h"
:
main() {
    :
    MPI_Init(...),
    :
    work → using MPI calls for
    : message passing
    :
    MPI_Finalize(...),
}
```

Note - → All MPI  
all processes will  
execute same  
code,  
if you want one  
process to do  
something diff. then  
we take use of  
rank to  
differentiate

Few important MPI calls → → MPI\_Comm\_World  
→ communicator

1. MPI\_Comm\_size (comm, &size);  
↓  
sets total no. of MPI Process

2. MPI\_Comm\_rank (comm, &rank);  
↓  
unique identifier for  
process

## • Implementing Broadcast

→ Some rank has data and it wants to pass that data to other ranks

#include "mpi.h"

;

main( . . . ) {

int bvalue, num\_ranks, rank;

MPI\_Init( . . . );

MPI\_Comm\_size(MPI\_COMM\_WORLD, &num\_ranks);

MPI\_Comm\_rank(MPI\_COMM\_WORLD, &rank);

if (rank == 0) {

bvalue = 100;

if (rank == 0) {

MPI\_Send {  
Pointers to buffer,  
l = Count,  
Element type,  
Destination,  
t = tag,  
cwid = Communicator,

MPI\_Send(&bvalue, 1, MPI\_INT, 1, 1, MPI\_COMM\_WORLD),

}

MPI\_Receive {

Pointers to buffer,

l = Count,

Element type,

datatype,

t = tag,

cwid = Communicator,

&status,

}

if (rank < numrank - 1) {

    MPI\_Receive(&bvalue, 1, MPI\_INT, rank - 1, l, MPI\_COMM\_WORLD, &status);

    MPI\_Send(&bvalue, 1, MPI\_INT, rank, l, MPI\_COMM\_WORLD),

}

else {

    MPI\_Receive(&bvalue, 1, MPI\_INT, rank - 1, l, MPI\_COMM\_WORLD, &status);

}

Note - Blocking calls - doesn't return until the MPI library has completed the necessary work  
 e.g. MPI\_Send(...), MPI\_Recv(...)

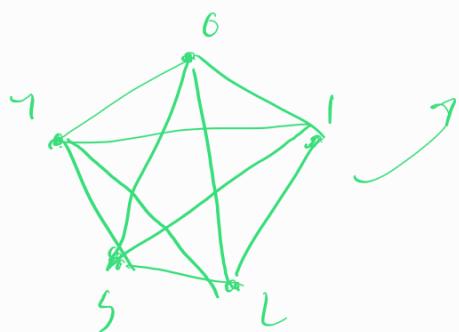
Non-Blocking calls -

Returns immediately, handing you a MPI\_Request while the actual data transfer proceeds in background  
 e.g. MPI\_Isend(...), MPI\_IRecv(...)

If we think this way of broadcasting might be slow, if per sent call cost = t then,  $(n-1)t$   
 $\text{cost}$

but can we use a better strategy for more efficient calls making

i) clique (complete graph) :



many calls have to be made from each node  
 not scalable

ii) binary tree (Recursive Doubling) :

1<sup>st</sup>: 0 → 1  
 2<sup>nd</sup>: 0 → 2, 1 → 3  
 3<sup>rd</sup>: 0 → 4, 1 → 5, 2 → 6, 3 → 7 } not reliable will keep track of doubling behaviour

senting data ← Receiving

1<sup>st</sup>: 0  
 2<sup>nd</sup>: 0, 1  
 3<sup>rd</sup>: 0, 1, 2, 3

1<sup>st</sup>: 1  
 2<sup>nd</sup>: 2, 3  
 3<sup>rd</sup>: 4, 5, 6, 7

2<sup>nd</sup> iteration n'

Algo to implement it

```
for (ctr=1; ctr < numrank; ctr++) {  
    if (rank < ctr) { /* i have data */  
        if (rank+ctr < numrank)  
            MPI_Send(&value, --, rank+ctr, --);  
  
    } else if (rank > 2*ctr) {  
        MPI_Recv(&value, --, rank-ctr, --);  
    }  
}
```

Note - MPI Send/Recv (~~&buff, count, type, source/destination,~~  
~~tag, community and &status);~~)

To identify which  
MPI\_Status. MPI\_Source  
MPI\_Status. MPI\_Tag

MPI\_Status  
↓  
MPI\_get\_Status (&status, datatype, &buff)

## \* MPI collectives

→ operations that are performed by all ranks together

- i) synchronisation operation
- ii) Data movement (Distribution & collection)
- iii) collective computation → like Reduce operation

### 1. Synchronisation -

a) MPI\_BARRIER (comm) → mpi-comm-world

main( ) {

; compute  
MPI\_BARRIER( --- ),  
;

every process  
barrier has to  
wait for all  
the processes to  
complete.

}

b) `MPI_Bcast` (& buffer, count, datatype, root, comm);

`Broadcast`

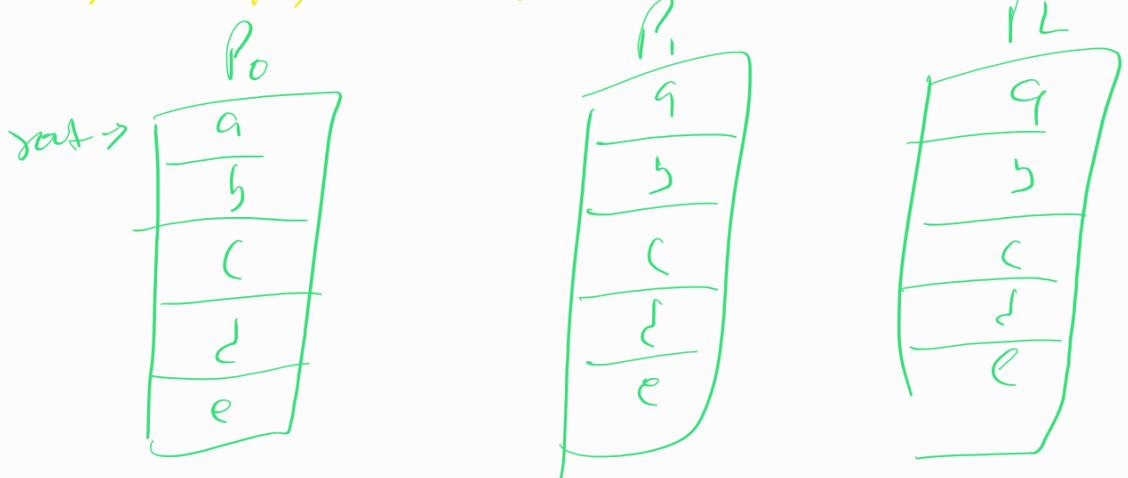
points to the start of data

rank which is doing the Broadcast

Note - `MPI_Bcast` is more efficient than `mpi_send`

Send → sends data node to Node

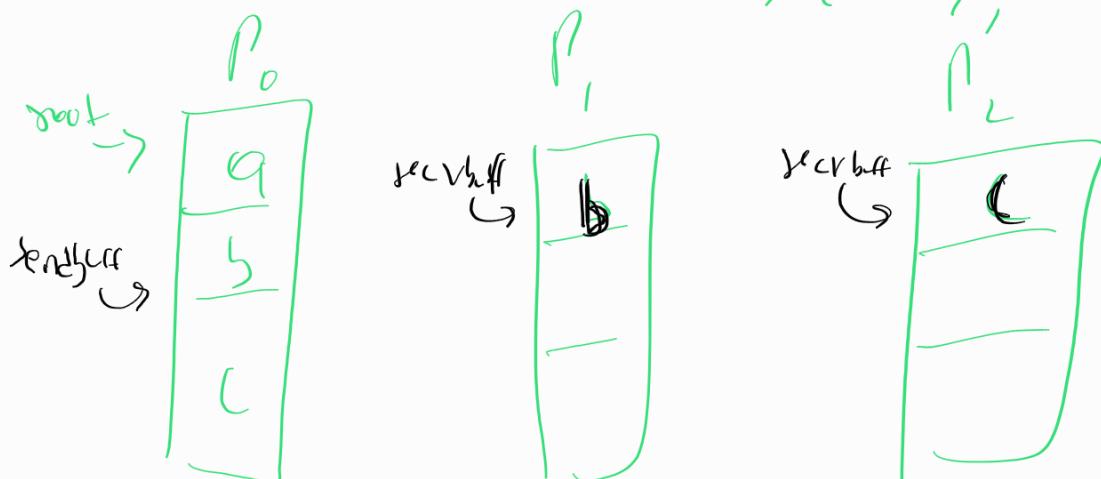
Bcast → makes it available for all on network



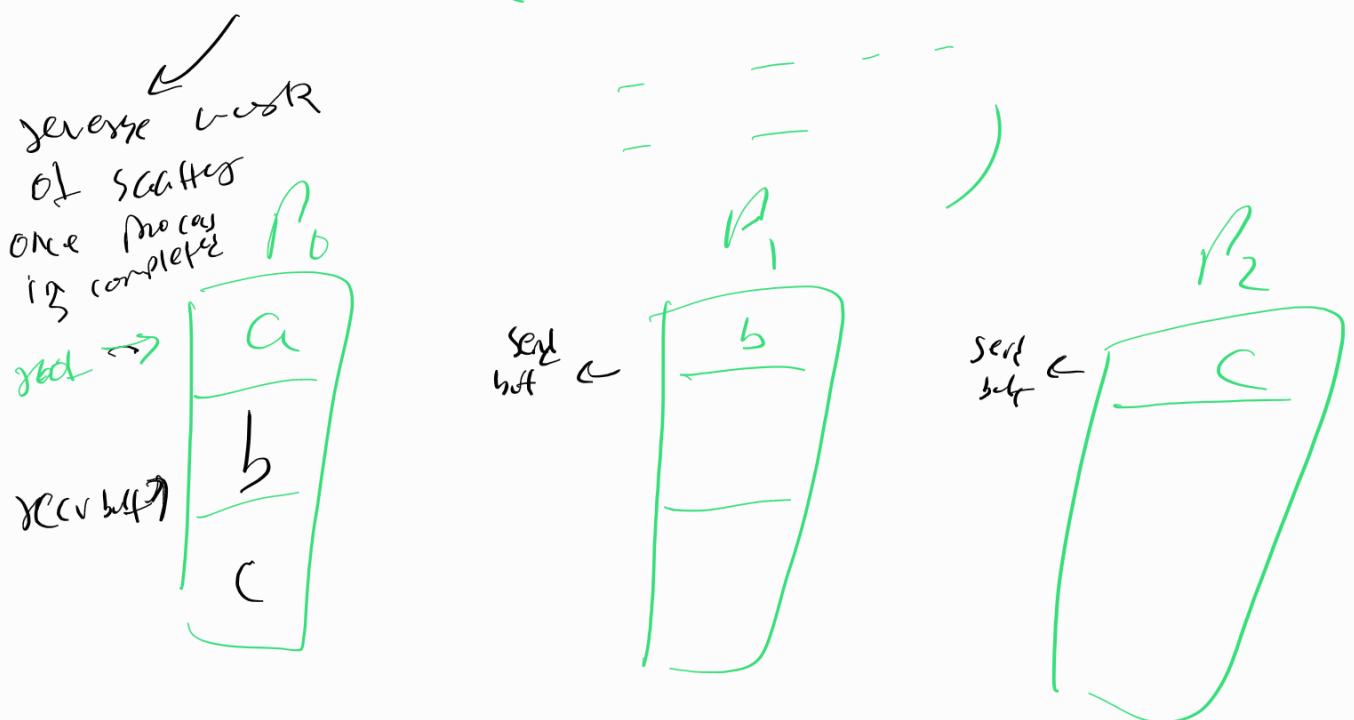
## 2. Data movement

In a master-slave type setting  
scatter helps in distributing data to other processors.

a) `MPI_Scatter` (& sendbuff, sendcnt, sendtype, & recvbuff, recvcnt, recvtype, root, comm);



b) MPI\_gather ( - - - )



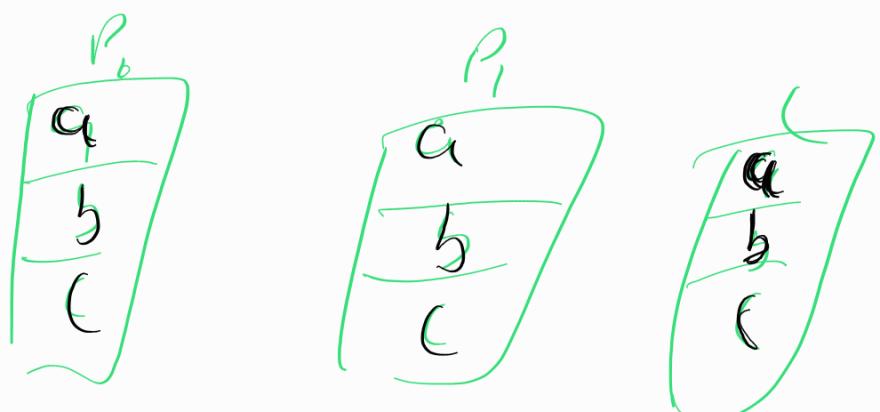
c) MPI\_gatherv (& sendbuf, sendcnt, sendtag,  
& recvbuf, recvcnt, recvtag,  
(communications))

→ no work

recv buf  $\Rightarrow$



recv buf  $\Rightarrow$



note - Blocking calls  $\rightarrow$  buffer isn't reusable

- Non-Blocking calls  $\rightarrow$  for overlaid communication with computation

a) MPI\_Irecv(...)

$\hookrightarrow$  returns immediately, even if data hasn't been fully copied out

$\hookrightarrow$  buffer becomes

not reusable till it has completely been copied

$\text{MPI_Irecv}(\&\text{buf}, \text{cnt}, \text{datatype}, \text{destination}, \text{tag},$   
 $\text{comm}, \&\text{request})$

$\hookrightarrow \text{MPI_Request}$

b) MPI\_IRecv(...)

$\hookrightarrow$  when we don't know when we'll be receiving message back

$\text{MPI_IRecv}(\&\text{buf}, \text{cnt}, \text{datatype}, \text{source}, \text{tag},$   
 $\text{comm}, \&\text{request})$

Note -

$\text{MPI\_Test}(\&\text{request}, \&\text{flag}, \&\text{status})$

if completed  $\Rightarrow \text{flag} = 1$   
if not completed  $\Rightarrow \text{flag} = 0$

int  
 $\text{MPI\_status}$

help us know  
the progress  
of Isend

$\text{MPI\_TestAny}(\text{count}, \&\text{array\_of\_requests}, \&\text{index},$   
 $\&\text{flag}, \&\text{status});$

Checks if  
any of the  
request is  
completed

$\text{MPI\_TestAll}$

Note -  $\text{MPI\_wait}(\&\text{request}, \&\text{status})$

→ A Blocking call  
which waits for all  
non-blocking calls to  
be completed.

$\text{MPI\_WaitAll}(\text{count}, \&\text{array\_of\_requests},$   
 $\&\text{array\_of\_status});$

$\text{MPI\_WaitAny}$

- Application : Distributed Histogram updation  
→ Large histogram

### 3 Collective computation

a) MPI\_Reduce(& sendbuff, & Recvbuff, cnt, datatype,

op, root, comm.);

Operation to  
be performed

the sink in with result  
is expected

MPI\_Sum  
MPI\_Max

MPI\_Min

:

e.g. MPI\_Sum

Send buff →

0  
3

1  
5

2  
7

Recv buff →

15

MPI\_AllReduce (---)

(---, root, ---)

no sent

e.g. MPI\_Sum

Send buff →

0  
3

1  
5

2  
7

Recv buff →

15

15

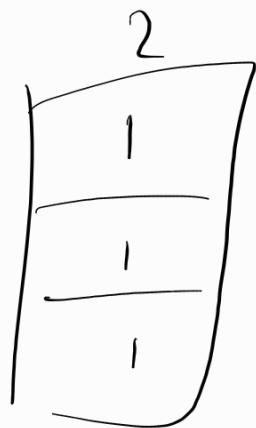
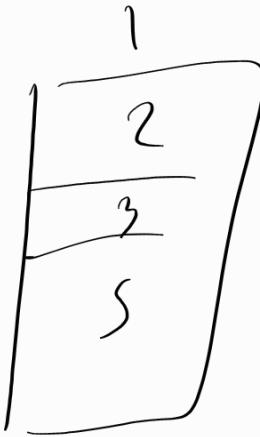
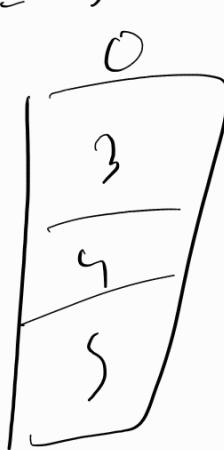
15

Sum available to all

↳ MPI\_Reduce\_Scatter (& sendbuf, &recvbuf, &count,  
datatype, op, comm);

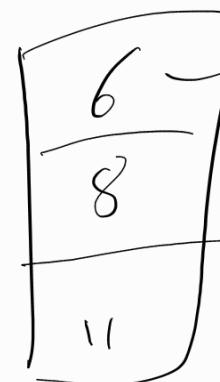
e.g. MPI\_Sum

Send buf →



Sum-1

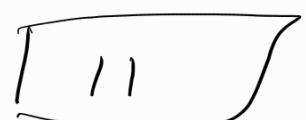
⇒ reduce ⇒



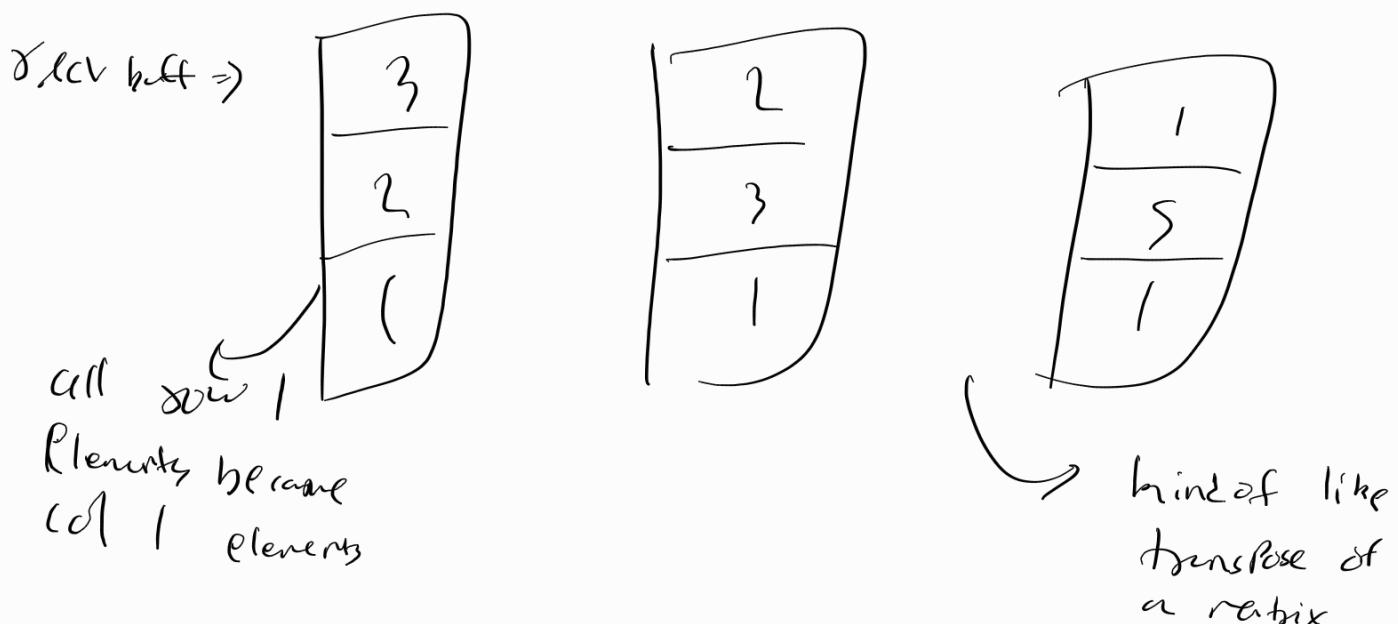
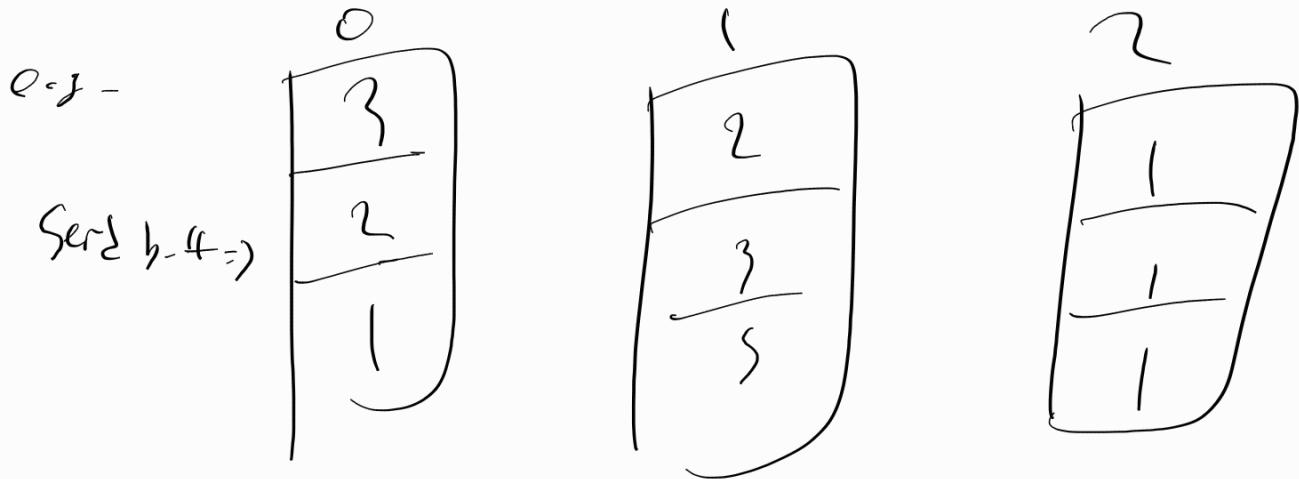
sum of  
first two  
elements

Sum-2 → Scatter ⇒

Recv buf →



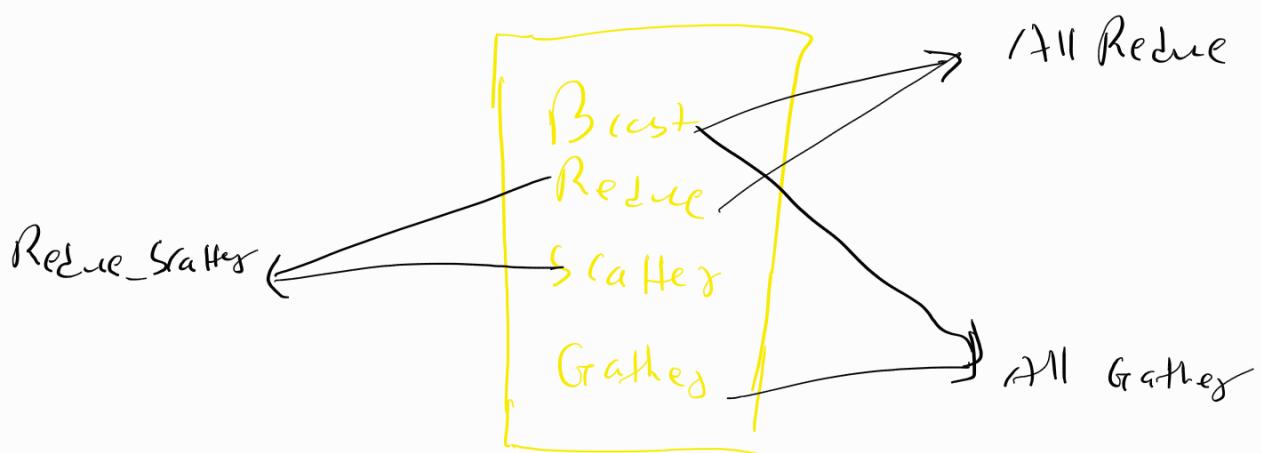
c) MPI- AlltoAll (& Sendbuff, sendcount, sendtype  
 & Recvbuff, recvcount, recvtype,  
 (omm));



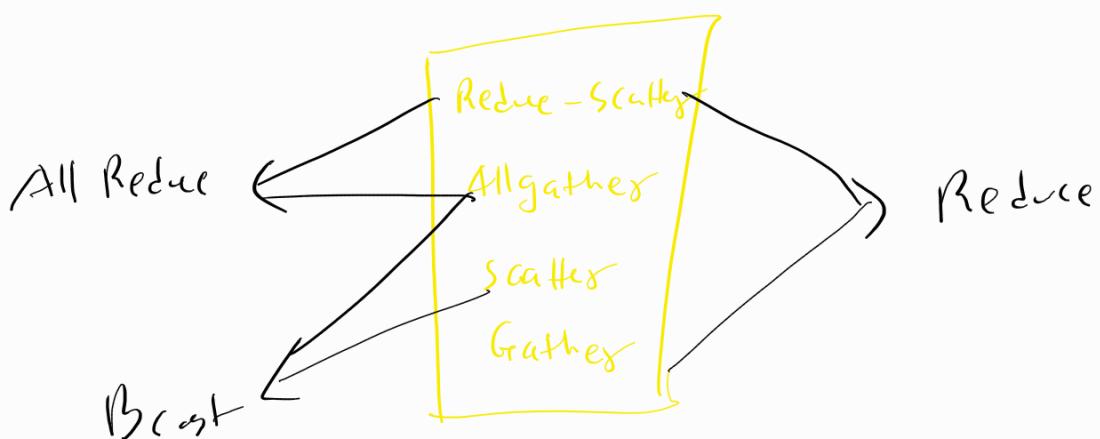
Note - To implement these calls

↳ usually we use few building blocks using combination of which we can implement other calls

i) Spanning tree Algo



ii) Bucket Algo





## Interconnects

Metrics to be considered for a suitable interconnection

- i) # links
- ii) diameter
- iii) degree (max)
- iv) Bisection width

} ↓ - metrics

BWP → ↑ connectivity

min no of nodes

to separate, that cause  
the network to be  
separated in two equal  
halves

Also degrds  
the lower bound

min value of communication  
allowed b/w two halves of  
a network

min no of links required  
to separate to create two halves

Decides cost of network

|                 | Linear | 1D torus        | 2D mesh     | 2D torus    | 3D torus                     | Mesherule      | Trees     | fat tree   | compl. cost |
|-----------------|--------|-----------------|-------------|-------------|------------------------------|----------------|-----------|------------|-------------|
| # links         | $P-1$  | $P$             | $2P$        | $2P$        | $2P$                         | $P \log P / 2$ | $2P$      | $P \log P$ | $P^2 / 2$   |
| Degree          | 2      | 2               | 4           | 4           | 4                            | $\log P$       | 2         | 4          | $P$         |
| Diameter        | $P-1$  | $\frac{P-1}{2}$ | $2\sqrt{P}$ | $\sqrt{P}$  | $\frac{P^{1/2} \times d}{2}$ | $\log P$       | $2\log P$ | $2\log P$  | 1           |
| Bisection width | 1      | 2               | $\sqrt{P}$  | $2\sqrt{P}$ | $2 \cdot P^{\frac{d-1}{2}}$  | $P/2$          | 1         | $P/2$      | $P^2 / 4$   |

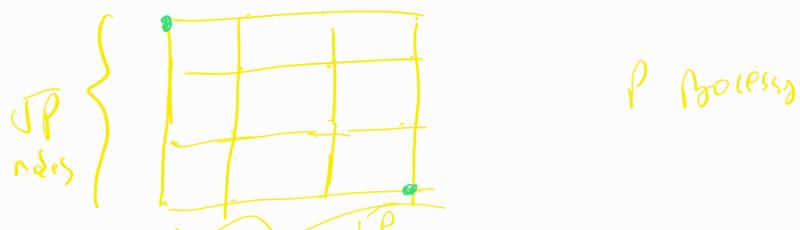
i) Linear



ii) 1-D torus

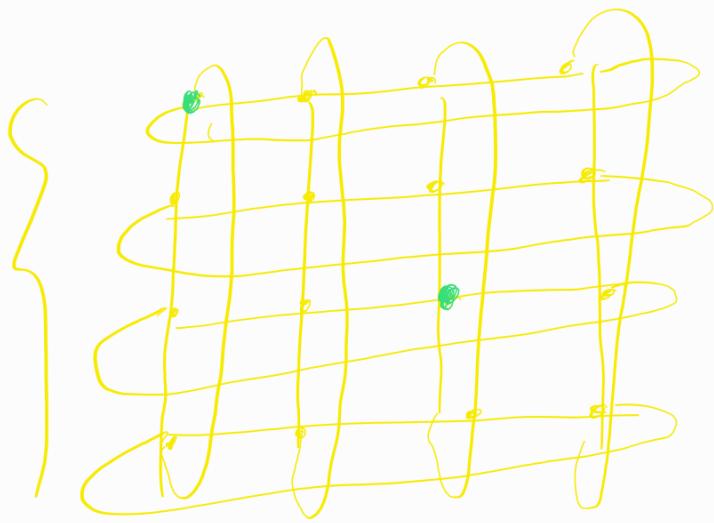


iii) 2-D mesh



iV)

JP



v) General formula for d-dimensional Trees  
p-roles

$$\text{degree} = 2d$$

$$\text{links} = \frac{2d \times p}{2} \sim pd$$

$$\text{diameter} = \underbrace{p^{1/2}}_2 \times d$$

$$\text{fisction width} = 2(p^{1/2})^{d-1}$$

## vi) Hypercube

↳ rank nodes in a dimension =  $L$

$$2 \times 2 \times 2 \times \dots \times 2 \quad P \text{ nodes}$$



$$\therefore 2^d = P$$

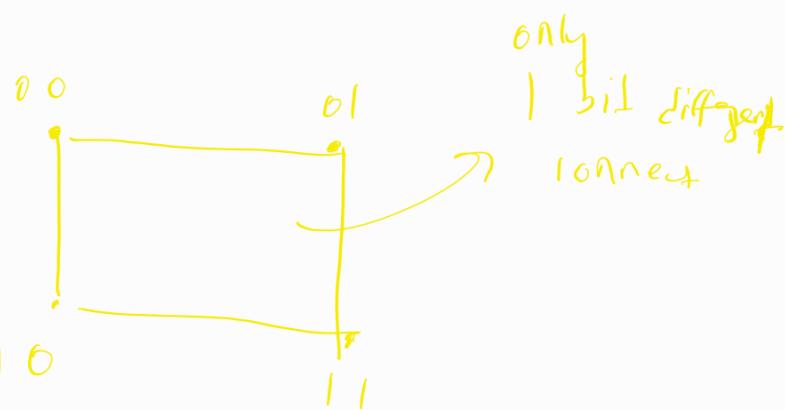
$$\therefore d = \log P \Rightarrow \text{hypercube}$$

0-D Hypercube :

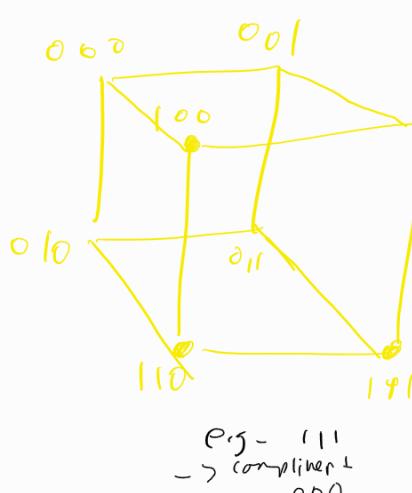
1-D Hypercube :



2-D Hypercube :



3-D Hypercube :



degree :  $\log P$

# links =  $P \log P / 2$

diameter :  $\log P$

bisection width =  $P/2$

↳ two heat furthest  
rule form a bit is  
is

To know if any random node is connected to which node

e.g. 100

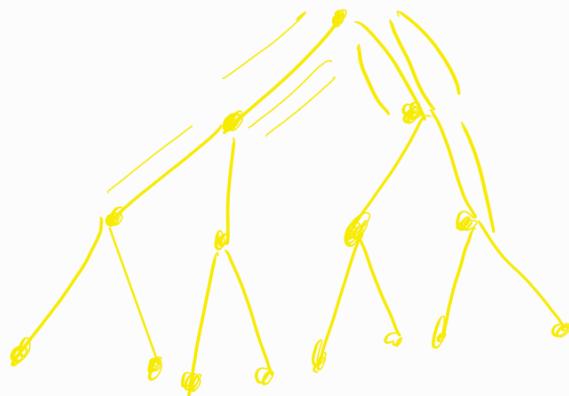
→ check all  $P^n$  by changing 1 bit  
000, 110, 101

### vii) Trees



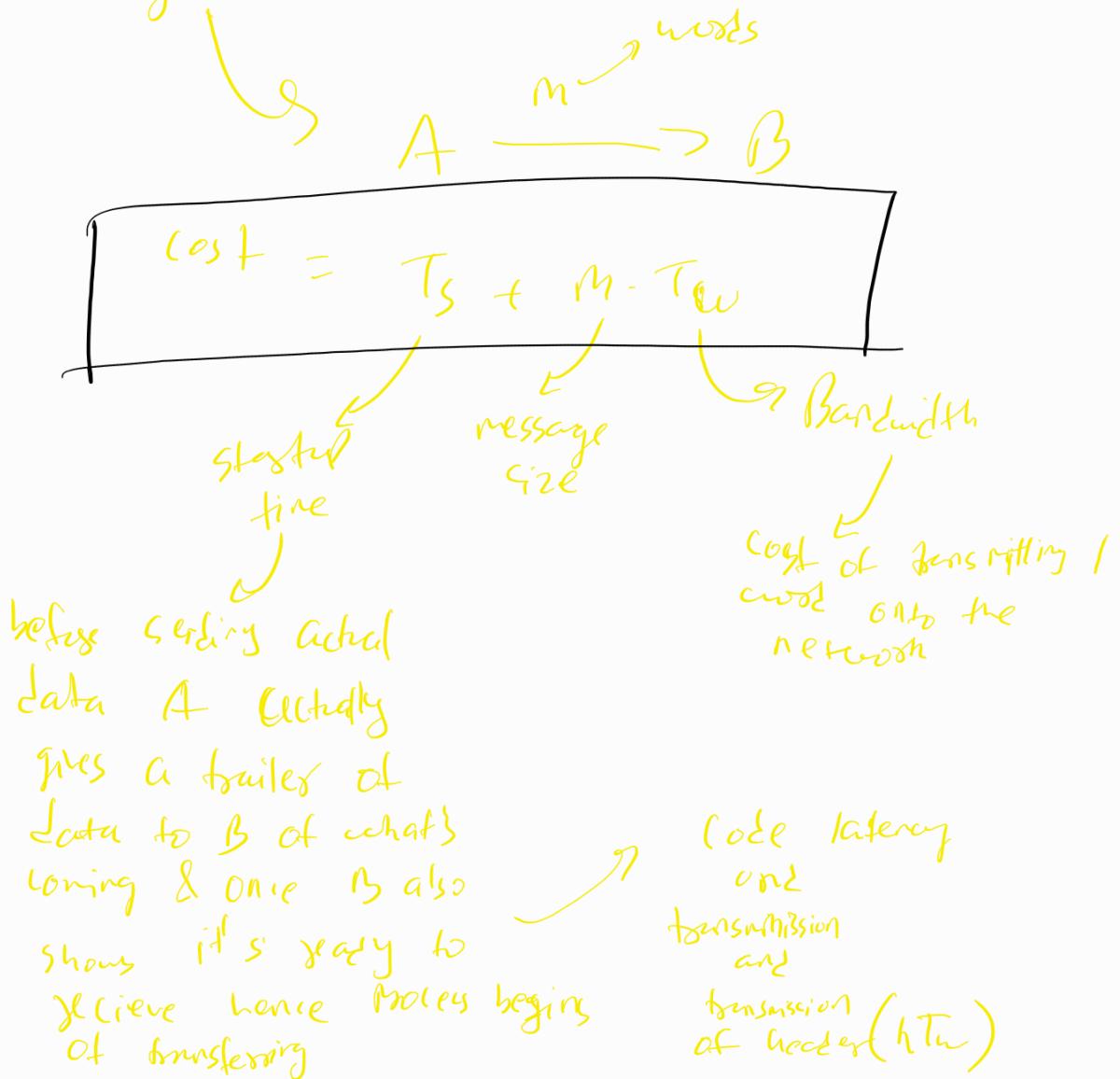
### viii) Fat Tree

↳ as we go up we keep doubling the links



# \* Evaluation of Communications in II Programs

## 1. Hockney models



- `v can do many I send & I recv simultaneously`
- `communication with any rate hQB since  $T_s$`
- `Duplex: all links are bidirectional`

now we'll try to analyse cost of few collectives



II

Diagram illustrating Broadcast: A single source node labeled 0 is connected to all other nodes in a ring of  $P-1$  nodes labeled 1 through  $P-1$ . The formula for cost is given as  $\text{cost} = (P-1)(T_s + m \cdot T_u)$ , with a note below it stating "very expensive".

Hence, choose Doubling / Halving

Diagram illustrating the recursive doubling process: A single source node labeled 0 is connected to  $\frac{P}{2}$  nodes. These  $\frac{P}{2}$  nodes are further connected to  $\frac{P}{4}$  nodes, and so on, forming a tree structure.

$$0 \rightarrow \frac{P}{2}, \quad \frac{P}{2} \rightarrow \frac{3P}{4}$$

$P = \text{no. of nodes}$

$\therefore \log P = \text{no. of bits}$

e.g. if  $P=16 \Rightarrow \text{no. of bits} = 4$

0000

$\dots y_1 \ y_2 \ y_3 \ \dots y_{\log P}$

2<sup>nd</sup> bit from  
left flipped

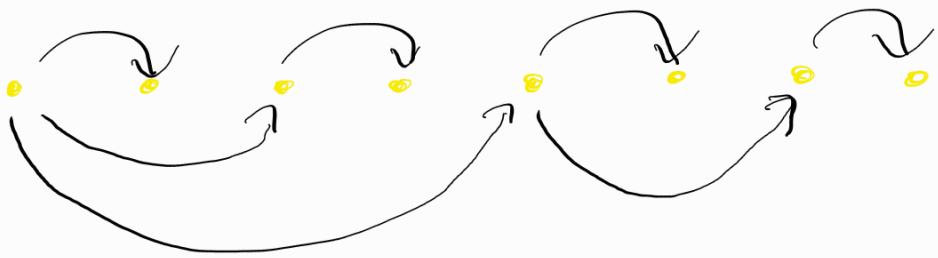
0000  $\rightarrow$  1000

0000  $\rightarrow$  0100

1000  $\rightarrow$  1100

here going in opp.  
dir as it's  
convenient to only  
flip one bit

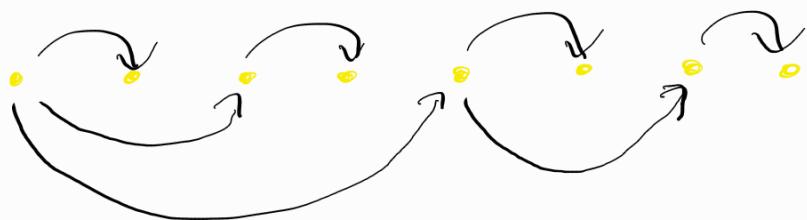
Diagram illustrating the recursive halving process: A group of  $P$  nodes is shown. The cost of the collective is given as  $\text{cost} = \log P (T_s + m \cdot T_u)$ , where the term  $\log P$  is annotated as "by Recursive Doubling". The cost is also noted to be proportional to "no. of rounds".



$$\therefore \text{cost} = \log P(\tau_s + m \cdot \tau_w)$$

(let's) check on different architectures

a) linear

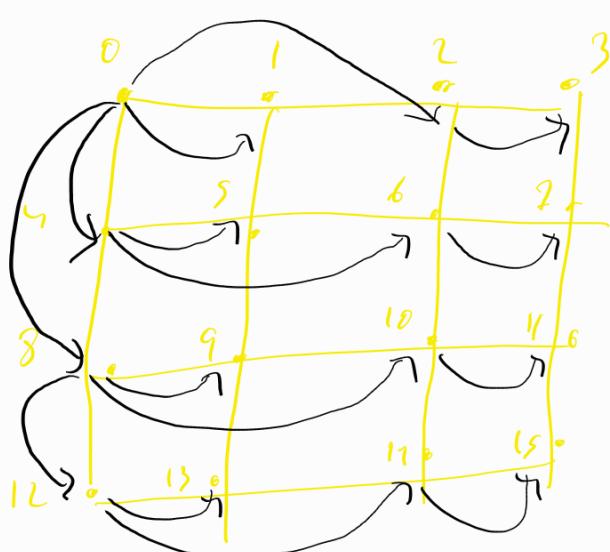


$$\therefore \text{cost} = \log P(\tau_s + m \cdot \tau_w)$$

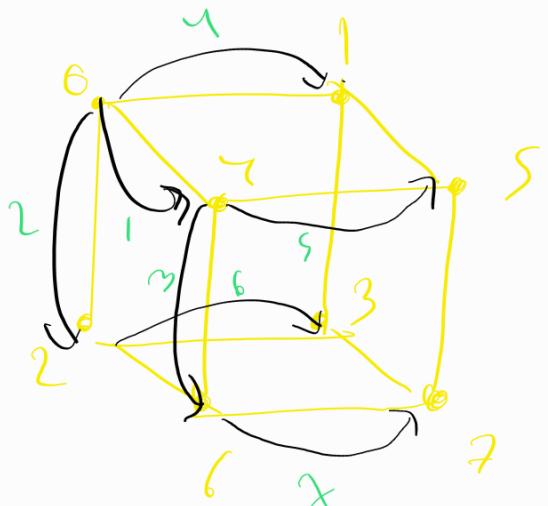
b) 2D mesh ( $\sqrt{P} \times \sqrt{P}$ )

$$6 \rightarrow 16/L$$

$$0 \rightarrow$$



(c) Hyper cube



(d) tree



Conclusion - cost same  
Same for all  
Architecture.

Note - The cost will be same of Reduce  
as cost as its just reverse  
of Broadcast, hence deal of Broadcast

$$\text{Cost of Reduce} = \log P(T_S + m \cdot T_B)$$

ii) Scatter

{  
gather}



$$\text{cost} = T_s + \left( \frac{mP}{2} \right) \cdot T_w$$

$$+ \\ T_s + \left( \frac{mP}{3} \right) \cdot T_w \\ \vdots$$

$$\text{cost} = T_s \log P + m \left( \frac{P}{2} + \frac{P}{3} + \dots + 1 \right) T_w$$

$$\boxed{\text{cost} = T_s \log P + (P-1)m T_w}$$

same for all  
architectures

iii) Reduce Scatter

&

All gather

$$: \text{cost}(\text{Reduce}) + \text{cost}(\text{scatter})$$

$$\log p (\tau_s + m \tau_w) + \tau_s \log p + (p-1)m \tau_w$$

$$\therefore \boxed{\text{cost} = 2\tau_s \log p + 2(p-1)m \tau_w}$$

(can be optimized further)

by recursive doubling

↓↓

$$\boxed{\text{cost} = (p-1)(\tau_s + m(p-1)\tau_w)}$$

Note - which approach is better to calc cost  
depends on  $\tau_s$  &  $\tau_w$  aka the size  
of the msg, if msg size big

↓

Recursive doubling  
approach

if msg size are small, latency term dominate,

↓

Reduce + scatter

## role - different sizes of nodes

case1: size of the nodes grain same, then ( $m$ )

(case2: sizes of some/all nodes are  $mp$ , then ( $mp$ )

? Broadcast:  $m$

All reduce:  $m$

Scatter:  $mp \rightarrow m \Rightarrow mp$

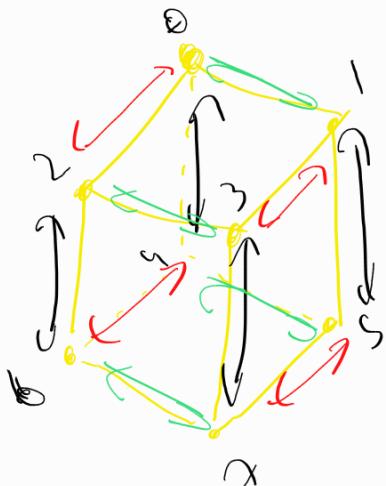
Reduce Scatter:  $mp \rightarrow m \rightarrow mp$

Gather:  $m \rightarrow mp \Rightarrow mp$

## GN mesh architecture

$$\text{Cost} = 2(\sqrt{P}-1)T_S + m(P-1)T_W$$

## ON Hypercube architecture



| 0: | 0[0] | 1[0] | 2[0] | ... | 3[0] |
|----|------|------|------|-----|------|
| 1: | 0[1] | 1[1] | 2[1] | ... | 3[1] |
| 2: | 0[2] | 1[2] | 2[2] | ... | 3[2] |
| 3: | .    | .    | .    | .   | .    |
| 4: | .    | .    | .    | .   | .    |
| 5: | .    | .    | .    | .   | .    |
| 6: | .    | .    | .    | .   | .    |
| 7: | 0[7] | 1[7] | 2[7] | ... | 3[7] |

Iteration 1 : 0: 0[0+4], 1[0+7], 2[0+7], 3[0+7]  
 4: 4[0+4], 5[0+4], 6[0+7], 7[0+4]  
 1: 0[1+s], 1[1+s], 2[1+s], 3[1+s]  
 5: 4[1+s], 5[1+s], 6[1+s], 7[1+s]  
 2: 0[2+6], 1[2+6], 2[2+6], 3[2+6]  
 6: 4[2+6], 5[2+6], 6[2+6], 7[2+6]  
 3: 6[3+2], 1[3+7], 2[3+7], 3[3+7]  
 7: 4[3+2], 5[3+7], 6[3+7], 7[3+2]

Iteration 2 : 0: 0[0+4+2+6], 1[0+4+2+6]  
 2: 2[0+4+2+6], 6[0+4+2+6]  
 1: 0[1+3+5+7], 1[1+3+5+7]  
 3: 2[1+3+5+7], 3[1+3+5+7]  
 5: 4[1+5+3+7], 5[1+5+3+7]  
 7: 6[1+5+3+7], 7[1+5+3+7]  
 4: 4[0+4+2+6], 5[0+4+2+6]  
 6: 6[0+4+2+6], 7[0+4+2+6]

Iteration 3 : 0: 0[1+2+3+4+5+6+7]  
 1: 1[1+2+3+4+5+6+7]  
 2:  
 3:  
 4:  
 5:  
 6:  
 7: 7[1+2+3+4+5+6+7]

cost  
 on  
 hypercube =  $t_s \cdot \log p + t_w \cdot m \cdot (p-1)$   
 for  
 Reduce-  
 Scatters

iv) All Reduce :

using hypercube architecture -

use same iterations just the message  
size changes.

cost  
 on  
 hypercube =  $(t_s + m t_w) \log p$   
 for  
 ReduceAll

v) All to All  $\rightarrow$  kind of like  
Transpose of  
matrix

By again Analyzing function of Hypercube  
with some iterations we can analyse cost  
just the size of message changes

∴

$$\text{cost}_{\text{cn}} \text{ (Hypercube for All to All)} = \left( t_s + \frac{m}{2} t_u \right) \log P$$

Note - Conclusion :  $P \text{ Broadcast} \rightarrow (2P-2)t_s + 2m\left(1-\frac{1}{P}\right)t_w$

↑  
to optimize  
we can  
pipeline

Broadcast :  $(t_s + m t_w) \log P [R \cdot D]$

lower bound  
" mtw

Scatter :  $t_s \log P + m(p-1)t_w$

longer bound  
" mtw

Reduce - scatter :  $2t_s \log P + mt_w(P \log P + P-1) [R \cdot D]$

Ring :  $(t_s + t_w \cdot m)(P-1)$

lower bound

mesh :  $2t_s(\sqrt{P}-1) + t_w \cdot m(P-1)$

mtw  $\{\sqrt{P}-1\}$

Hypercube :  $t_s \log P + t_w \cdot m(P-1)$

All-Reduce :  $2(t_s + m t_w) \cdot \log P [R \cdot D]$

Hypercube :  $(t_s + m t_w) \log P$

lower bound  
mtw

Optimal  $\rightarrow$  Chain :  $2t_s(P-1) + 2t_w \cdot m\left(1-\frac{1}{P}\right)$

All to All : Hypercube :  $(t_s + m t_w) \log P$

E-cube routing :  $(P-1)\left(t_s + \frac{m}{P}t_w\right)$

lower bound =  
mtw

Note Again to compare with approach is better if  
depends on size of message

short message  $\rightarrow$  Recursive Doubling

large message  $\rightarrow$  Other Alg.

Note - The lower bound for 'All to All' was  $m(t_w)$  but from Hypercube architecture we're getting a lower bound of  $m(t_w) \log p$  which isn't correct so a new approach needs to be found



## E-cube Routing



→ simplest, deadlock-free way to send a package from any source  $s$  to any destination

▷ in an  $n$ -dimensional hypercube works

by fixing one different bit at a time

for All-to-All

$$(P-1) \left( t_{sf} + \frac{m}{P} t_w \right)$$

$n$  cube  $\Rightarrow 2^n$  nodes  $\Rightarrow n$  bit binary string

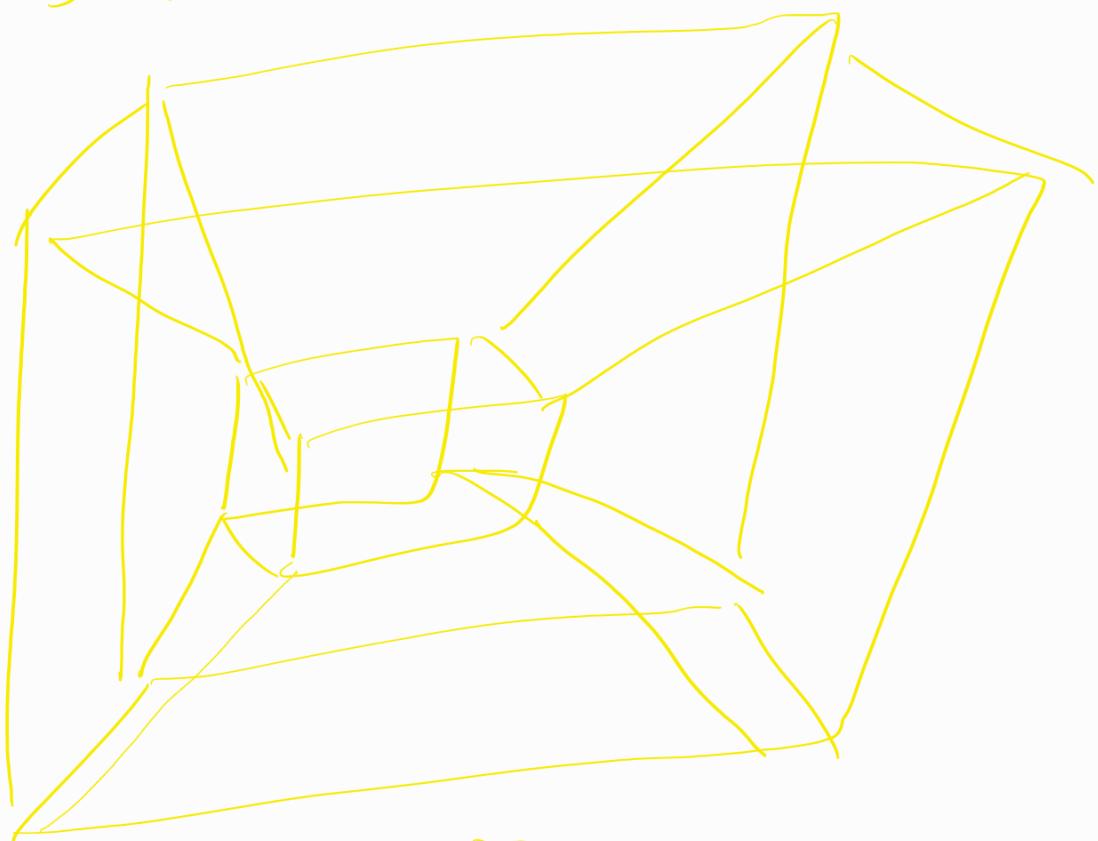
e.g. if 3-D Hypercube, iteration = 5, rep the route  $\rightarrow 2^3 = 8$  nodes

$$5 \Rightarrow 101$$



e.g. Find the Route Q.

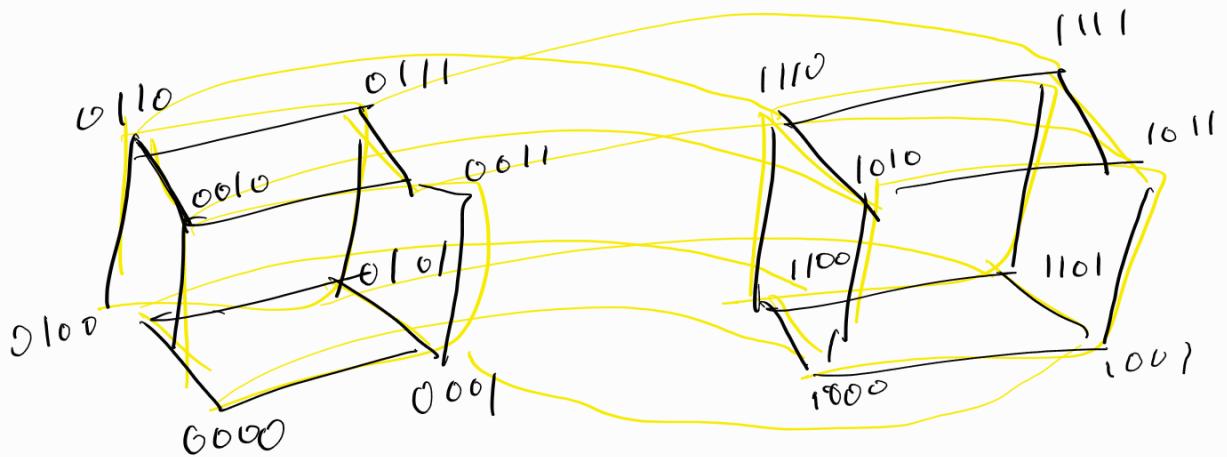
4-D Hypercube



4-bit binary

0 8

$$S = 0110 \\ D = 1101$$



→

$$\text{Step-1} \Rightarrow S \oplus D \Rightarrow 0110 \oplus 1101$$

0..

$$0110 \rightarrow 0111 \rightarrow 0101 \rightarrow 1101$$

Path

1011

3<sup>rd</sup> bit from USB  
won't flip

