

OpenMP: Shared Memory  
Distributed

C, OpenMP, OpenMPI → Technical

[www.letshpc.org](http://www.letshpc.org)

Performance → operations / computations per unit time

Algorithms → Computer

- |                  |                   |
|------------------|-------------------|
| 1) Deterministic | Non deterministic |
| 2) Abstract      |                   |
| 3) Terminating   | Non terminating   |

System / Architecture make deterministic → non deterministic

Agenda → Accuracy + Efficiency

CPU  
core

Performance  
↓ unit

flops (floating point operations processed)

flops / sec.

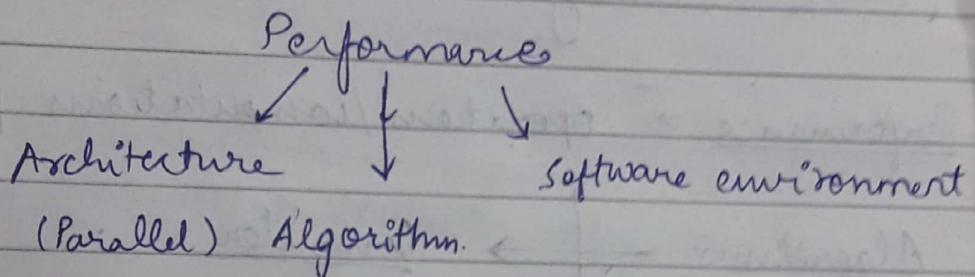
- Frequency of the processor (cycles/s)
- Number of cores
- Flops / cycle

Rotary  
N.Y. Dist. 3050

Clocks/s ↓	Flops/cycle
$2.5 \times 10^9$ / second	$\times 4$ Flops
	$\times$
$= 10 \text{ G Flops / sec}$	

Multiply it with the number of cores.

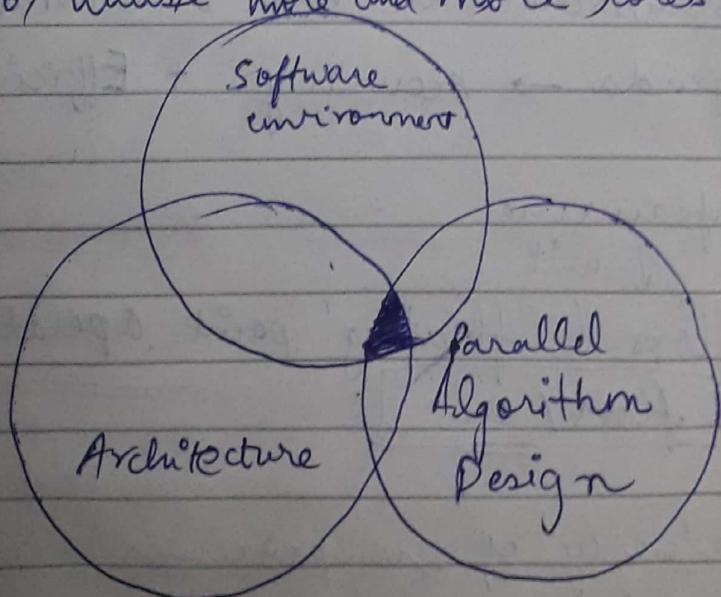
Serial programming  $\rightarrow$  10G Flops/sec.



### Cache memory

- ① Portability
- ② Scalability

able to utilise more and more cores,



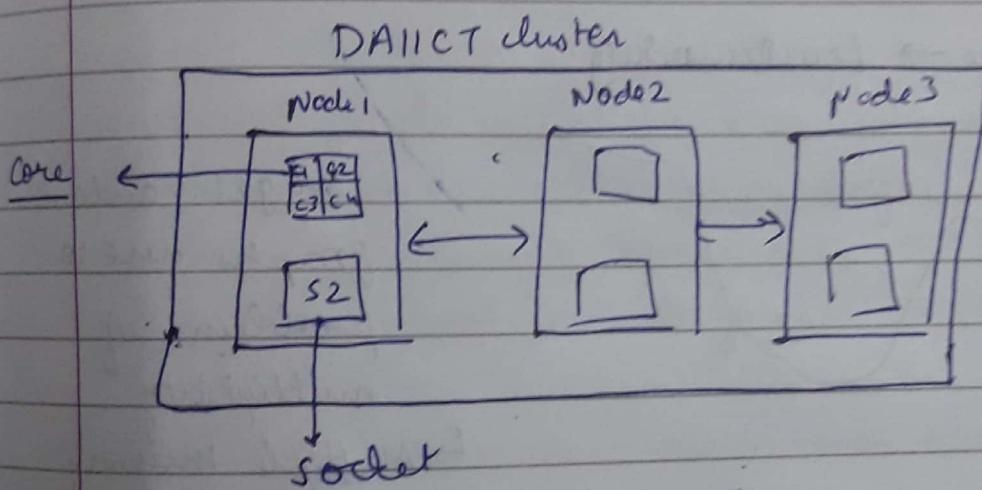
computation  $\rightarrow$  operation + data

Architecture  $\rightarrow$  ND (Non Deterministic)

Algorithm Design  $\rightarrow$  D

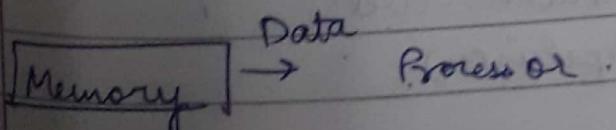
Software environment  $\rightarrow$  ND

Max. possible performance out of a particular system for  
a particular problem



Supercomputer = Several nodes  $\rightarrow$  HPC

Node performance in gflops =  $(\text{CPU speed in GHz}) \times (\text{number of cores})$   
 $\times (\text{CPU instruction per cycle}) \times$   
 $(\text{number of CPUs per node})$



① Bandwidth  $\rightarrow$  MB/s

② Latency  $\rightarrow$  seconds

von Neumann bottleneck

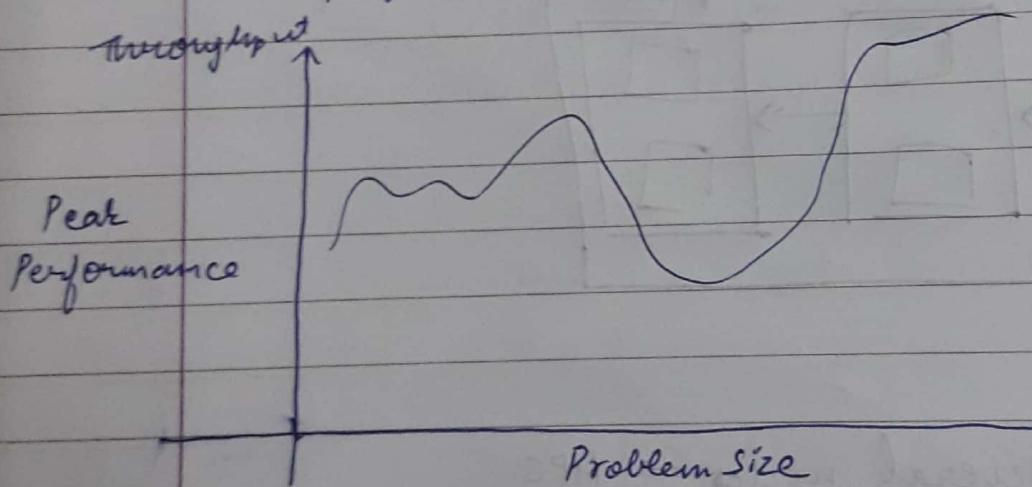
SISD ( single instruction Single Data )  
SIMD MISD MIMD

- Instruction Stream
- Data Stream

Task.

$$\text{Speedup} = \frac{\text{Serial Time}}{\text{Parallel Time}}$$

Peak performance  $\rightarrow$  benchmarking



We get random graphs due to problem of multiplicity:

- ↳ Multiple memory
- ↳ Multiple data paths
- ↳ Multiple processes

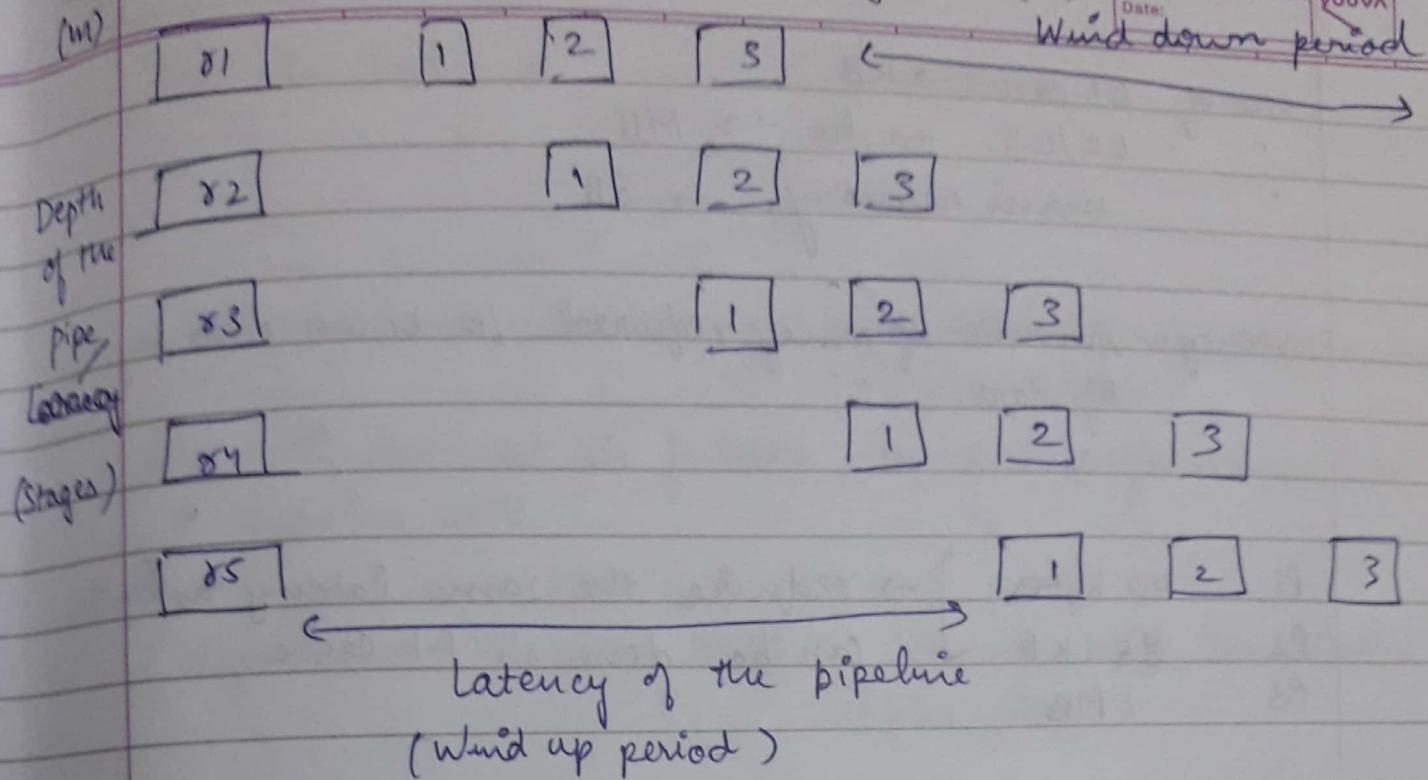
Our aim is to get best performance

Parallelism possible in various level :

\* Node \* Thread \* Instruction Level  
( Pipelines )

# Pipelines

→ cycles (n)



- 1 result / cycle after wind up phase
- Depth of the pipeline =  $m$  stages.
- Total Independent operations  $\rightarrow N$   
Time of pipeline to finish  $\rightarrow (N+m-1)$

serial (no pipeline) time  $\rightarrow Nm$

$$\boxed{\text{Speedup} = \frac{Nm}{N+m-1}}$$

$$N \gg m$$

we get much higher speedup.

$$\boxed{\text{Throughput of the pipeline} = \frac{N}{N+m-1}}$$

For higher values of  $N \Rightarrow \boxed{\text{Throughput} \approx 1}$

Size of L1 cache  $\rightarrow$  kB

" " L2/L3 cache  $\rightarrow$  MB

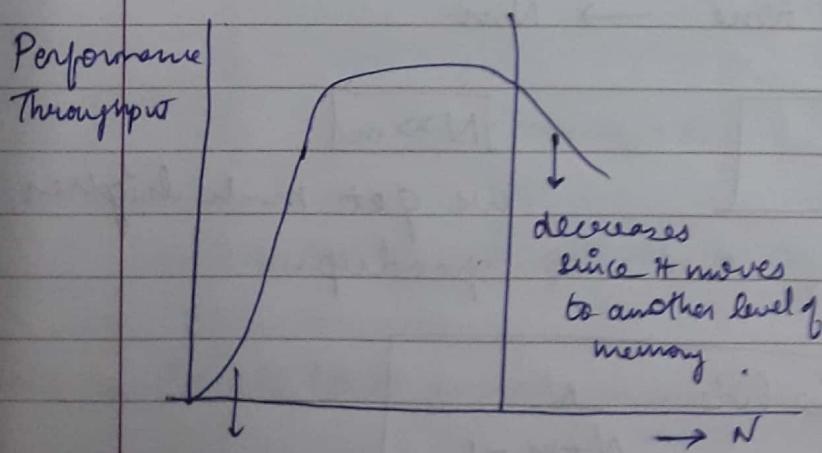
" " main memory  $\rightarrow$  GB

Latency: Amount of time required to access one byte of data.

P1 40 bytes ]  $\rightarrow$  require the same latency as both  
P2 8B [ KB can be done in L1 cache.  
P3 1 MB

Performance of CPU is increasing at a much greater rate than performance of memory (DRAM) [von Neumann bottleneck].

→ Parallel vector triad benchmark.



increases because of instruction level parallelism (pipelines)

Machine balance

Code balance

Machine balance = memory bandwidth / peak performance  
(bytes/flops)

It will decrease in future as peak performance ↑ at a greater rate

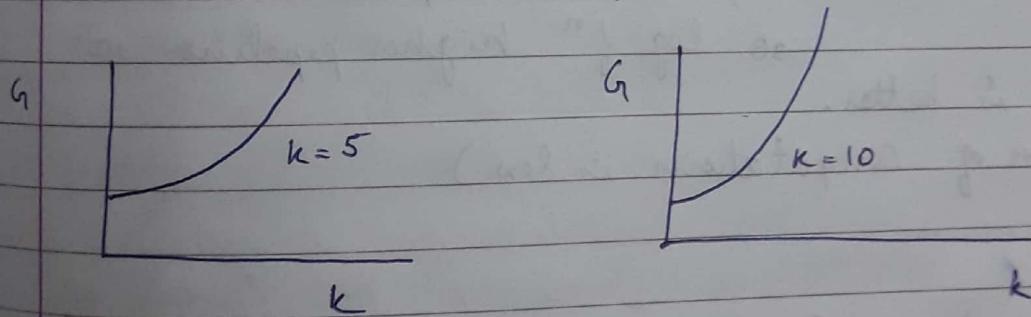
Code balance = data traffic / floating point operations  
(bytes/flops)

Code balance of vector triad =  $4/2 = 2 \text{ bytes/flops}$ .

Reciprocal of code balance  $\rightarrow$  computational intensity

$\therefore$  Lower code balance is desirable

$$\text{Grain } G = \frac{k}{\cancel{x} + (1-x)k} \quad 1 \leq G \leq k$$



2 factors : speed and size

$\rightarrow$  Cache line length = 16  $\Rightarrow$  Hit ratio =  $\frac{15}{16}$   
(spatial locality)

(1<sup>st</sup> A, cannot hit, not present on cache,  
initially A [1:n] bring 16 to cache, I can't hit  
 $\frac{15}{16}$  hit ratio)

Spatial + Temporal locality  
→ (2 terms related to)

Page No.:

Date:

Yours

Cache is  $K$  times faster than RAM  
↳ (latency + Bandwidth)

Cache reuse ratio is  $\gamma$  (fraction)  
Access to main memory ( $T_m$ )

$$\text{Cache access time} = \frac{T_m}{K} = T_c$$

$$\begin{aligned}\text{Average access } T_{av} &= \frac{\gamma \times T_m + (1-\gamma)T_m}{K} \\ &= \gamma T_c + (1-\gamma)T_m\end{aligned}$$

$$\text{Performance gain} = \frac{T_{av}}{T_m} = \frac{\gamma T_c + (1-\gamma)T_m}{\gamma T_c + (1-\gamma)K T_c} = \frac{1}{K}$$

Case A  $A = A + B * \pi^{2.0}$  (case A)  $A, B$  are float

Case B  $A = A + B * B$  (case B)

B is better as  $(B * \pi^{2.0}) \leftrightarrow \exp(2 \ln B)$

so  $\log f^n$  higher pipeline  $\Rightarrow$

(case B) is better.

(Number of computations is less)

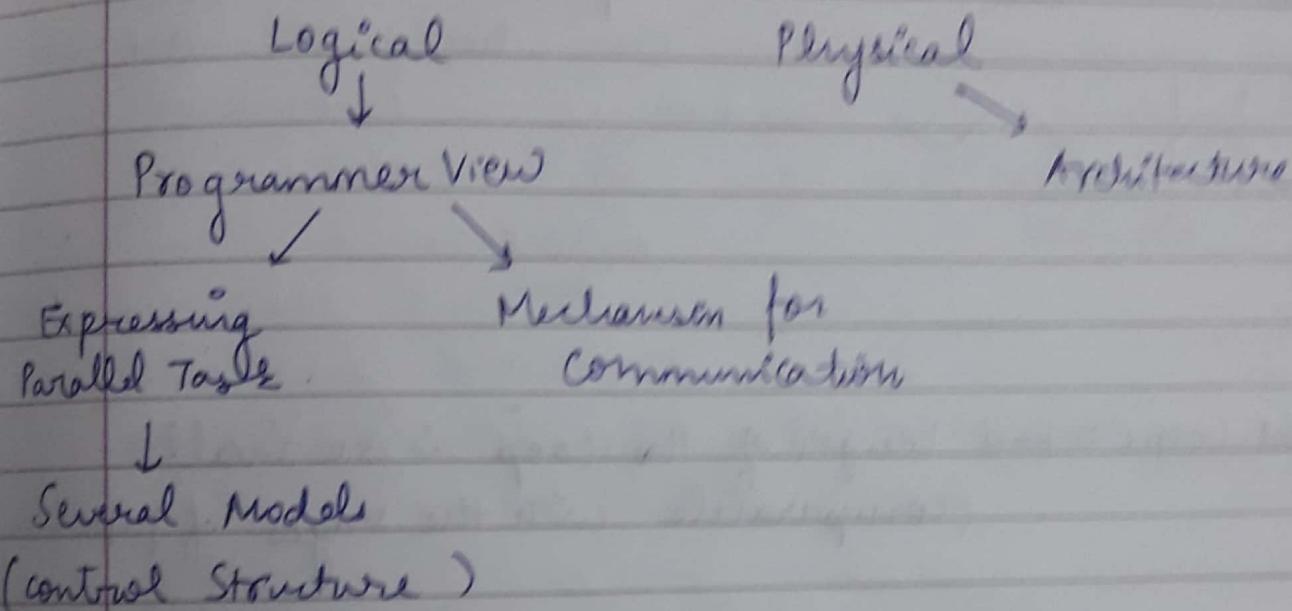
Avoid conditional branching:

- since in conditional branching → there is a possibility of branch miss
- compiler ~~too~~ makes a guess and if it is proved wrong he has to flush the pipeline

## Loop Unrolling

Chances of cache hit increases.

## Organisation of Parallel Platforms



Granularity: Ratio of computation to communication.

Coarse: large computational work b/w communication

Fine: small

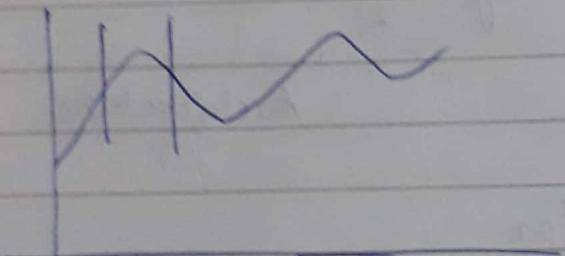
## Matrix Multiplication

To finest level of granularity =  $O(N^3)$   $N^2$   
Coarse " " " "  $\frac{N}{4}$  meets

OpenMP (Multi Processing  $\rightarrow$  MP) (Shared systems)  
Open MPI (Message passing) (Distributed systems)

*start*

- \* not necessary to compile using -fopenmp  
(How to measure one clock time?)
- \* Resolution of the clock  
time (We need to take multiple measurements)



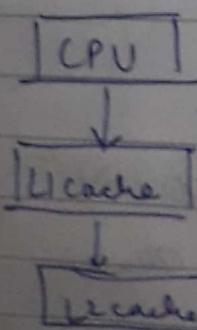
Tight loops  $\approx$  too length of the loop is so small  
(comparable with the stages of pipeline)

Loop Size  $\rightarrow 10^3$

4 arrays  $\rightarrow 4 \times 10^3$

(Double precision)  $\rightarrow 32 \times 10^3 = 32 \text{ K}$  (Maximum

$\rightarrow$  size of L1  
cache)



register

Latency of L2 cache more than  
L1 cache.

Memory system

64 bit DDR AM interface  
8 channels  
1 GHz clock frequency

(DDR)  $\rightarrow$  Double Data Rate

Access latency of 200 cycles

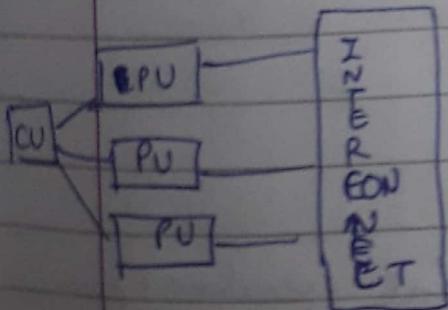
$$\rightarrow \text{Access Throughput} = \frac{8 \text{ bytes}}{\text{transfer}} \times 8 \text{ channels} \times 2 (\text{DDR}) \times \text{channel} \times (\text{clock} / \text{cycle}) \times 10^9 \frac{\text{clock}}{\text{second}}$$

$$= 128 \frac{\text{Gbps}}{\text{second}}$$

L3 cache is a shared cache (multiple threads can access)

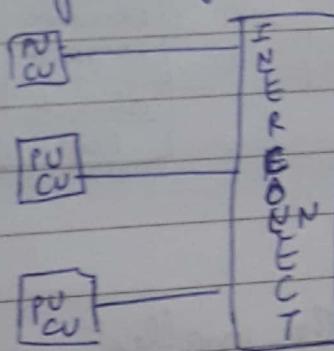
## SIMD

(Single Instruction  
Multiple Data)



## SPMD

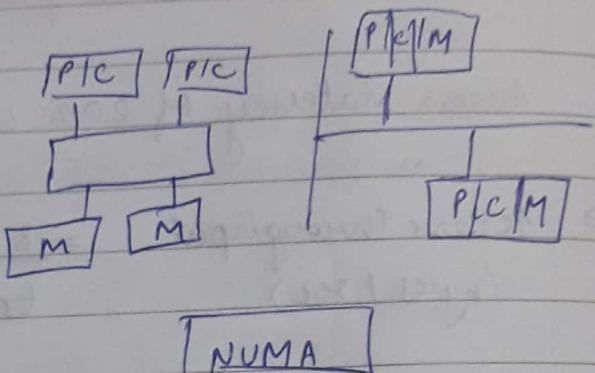
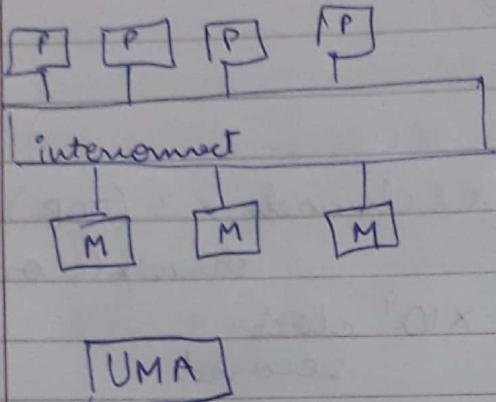
(Single Program Multiple Data)



Memory :   
  $\rightarrow$  Primary Private  
 $\rightarrow$  Shared

UMA: Uniform Memory Access

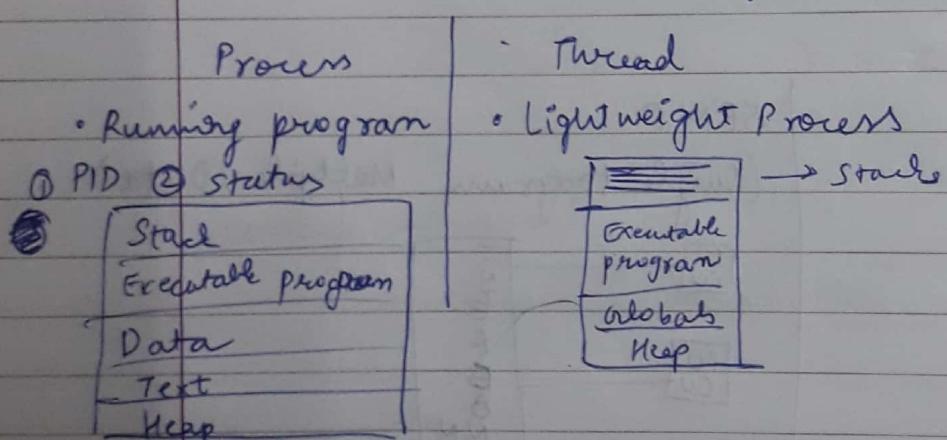
NUMA: Non uniform Memory Access -



NUMA

( Each process can access its own cache very fast however cannot access other cache quickly )

- context switching (diff. from parallel programming)



\* Thread is the smallest unit of sequential execution of instructions.

\* Process is an execution stream

Status of the process :

- New • Running • Ready • Waiting • Terminated
- Threads have independent stacks but share heap and global variables.

- \* Thread includes (individual)
- \* Thread ID \* Program Counter \* Register Stack \* Stack (pointer) (pointer)
- (shared)
- \* Text section \* Data section \* other OS resources

OS uses process control block

- \* execution state for each thread
- \* scheduling information
- \* Memory used by the process
- \* Information about open files.

~~(SIMD takes more time in scheduling as it has more only)~~

~~(SIMD takes less time in scheduling as it has only a single CU whereas SPMD takes more time in scheduling)~~

~~(Every thread can read at the same time whereas while writing only one thread can write at a time)~~

### Dependencies:

for  $i = 1 \dots n$

for  $j = 1 \dots m$

$b[i, j] = \dots + b[i, j-1]$  → Dependency on  $j$

4 kinds of dependences:  
① Flow, ② Anti ③ Input ④ Output

→ Flow (true dependence)

Statement 'i' precedes 'j'  
and 'i' computes a value that 'j' uses

$$\textcircled{1} \quad n=1 \quad \textcircled{2} \quad y=n+2 \quad \textcircled{3} \quad n=z-w \quad \textcircled{4} \quad n=y/2$$

$\textcircled{1} \rightarrow \textcircled{2}$  [These are the dependencies]  
 $\textcircled{2} \rightarrow \textcircled{4}$

→ Anti dependence

Statement 'i' precedes 'j'  
and 'i' uses a value that 'j' computes

$$\textcircled{1} \quad n=1 \quad \textcircled{2} \quad y=n+2 \quad \textcircled{3} \quad n=z-w \quad \textcircled{4} \quad x=y/w$$

$$\textcircled{2} \rightarrow \textcircled{3}$$

$$\textcircled{2} \rightarrow \textcircled{4}$$

→ Output dependence

Statement 'i' precedes 'j'

and 'i' computes a value that 'j' also computes

$$\textcircled{1} \quad n=1 \quad \textcircled{2} \quad y=n+2 \quad \textcircled{3} \quad n=z-w \quad \textcircled{4} \quad n=y/2$$

$$\textcircled{1} \rightarrow \textcircled{4}$$

$$\textcircled{3} \rightarrow \textcircled{4}$$

→ Input dependence

Statement 'i' precedes 'j'

and 'i' uses a value that 'j' also uses

$$\textcircled{1} \quad x=1 \quad \textcircled{2} \quad y=n+2 \quad \textcircled{3} \quad n=z-w \quad \textcircled{4} \quad n=y/z$$

③ → ④

- \* Difficult to remove flow dependence.
- Dependence flows from  
 $i \xrightarrow{\quad \text{to} \quad} j$   
 ↓           ↓  
 source      sink

Two important dependencies:

- Data dependency → Task dependency  
 (represented using DAGs)

### Data dependency graph

- Nodes are the statement
- Edges are the dependencies relations

\*\* Loop carried  
dependence

$$a[i] = b[i] + c[i] \rightarrow (\text{Flow dependence})$$

$$d[i] = a[i] \quad [\text{loop independent}]$$

loop independent  
depend

loop distribution

$$a[i] = b[i] + c[i] \rightarrow (\text{loop carried dependence})$$

$$d[i] = a[i-1]$$

for ( $i=0; i < N; i++$ )

$$a[i-1] = b[i]$$

for ( $j=0; j < r; j++$ )

$$c[j] = a[j]$$

can be parallelized  
locality is worse  
(more overhead due  
to parallel processing)

for ( $i=0; i < N; i++$ )

$$a[i-1] = b[i]$$

$$c[i] = a[i]$$

cannot be parallelised.  
locality is better  
(spatial locality)

- Loop fusion done for better utilisation of spatial locality
- Loop distribution to bring parallelism.

directives: special preprocessor instructions

Shared: all processors have a shared memory  
 → bus contention leads to program's scalability

Distributed: all processors have distributed memory  
 → memory management is difficult

### Fork-Join concept

- Fork creates concurrency
- Join removes concurrency
- Fork children can join with the parent in any order.

OpenMP → structured blocks: one point of entry and exit

(Inside structure blocks if condition is not allowed)

Scope → Access  
 Scope → P Time it is valid

- Global variables shared among threads
- Private variables:
  - exist only within the new scope; they are uninitialized and undefined outside the scope
  - loop variables.

### Create Parallel Regions

C

```
# pragma omp parallel
{
```

code to be executed by each thread

}

- `foromp`.

- `OMP_NUM_THREADS`

```
export OMP_NUM_THREADS=4
setenv "        " 4
```

// determines how many threads are utilized

~~#include "omp.h"~~  
~~void main()~~

### SPMD

↗ (decreases concurrency)

```
# pragma omp critical // only one thread can
global_result += my_result; execute
```

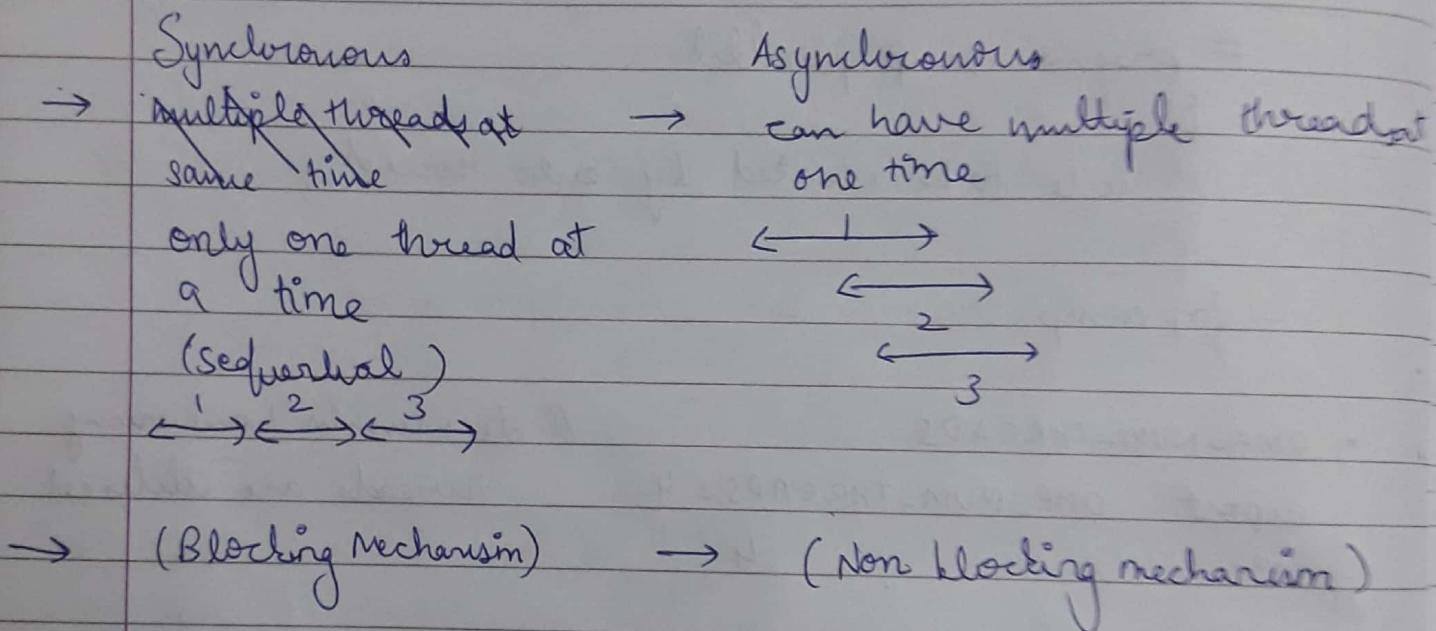
# pragma omp parallel → creates threads

# pragma omp for → no directive does work sharing

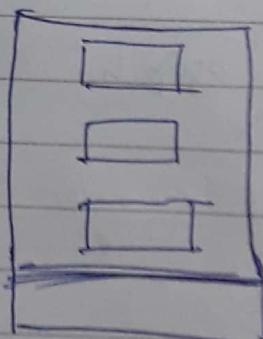
Three ways:

- ① Define a variable  $i$  explicitly inside parallel
- ② # pragma omp for (use inside this)
- ③

Index variable of a worksharing loop is private by default.



Data Scope  
Coordination  
Synchronisation  
Load balancing  
work sharing  
Scheduling



Barrier

- (a) event synchronisation  
(b) explicit synchronisation



\* By defining number of threads we control granularity

Floating point operations are not associative

(due to round-off error)

While using reduction there should be an identity element.

f

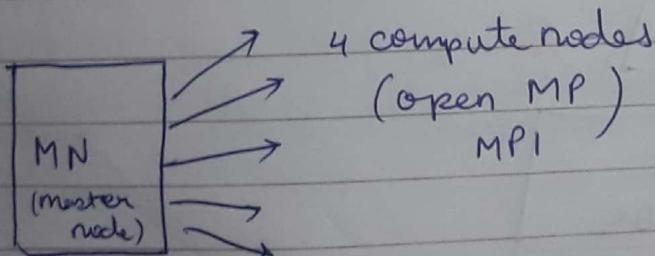
Types of dependencies we have seen :

- control
- data
- loop
- structured

Task dependency graph → control dependency

Task interaction graph → data dependency

IP of HPC cluster: 10.100.71.130



Any directory created on any node will be created in the rest of the nodes

SSH ID @ 10.100.71.130

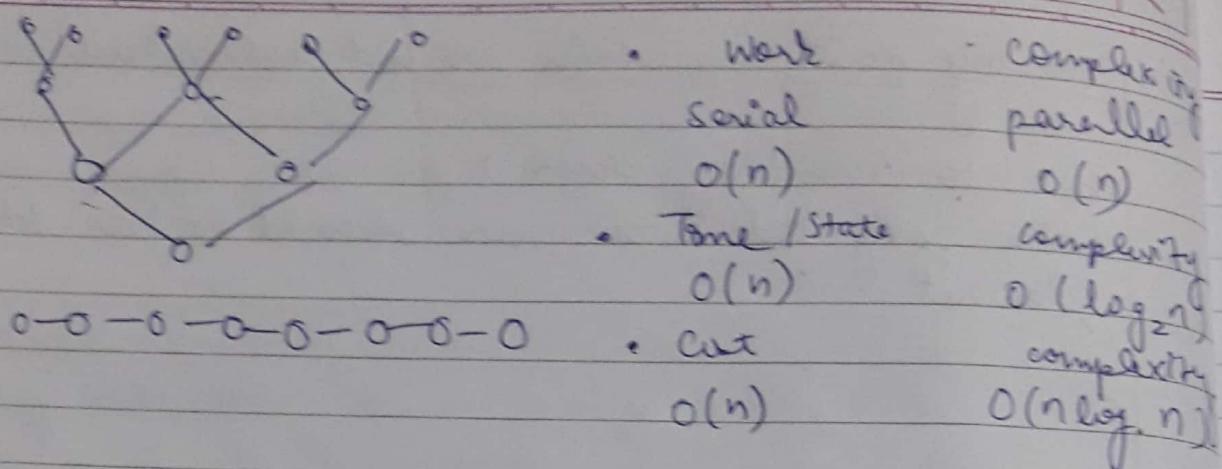
Top

SSH GICS0

GICS1

GICS2

- Write all entries in an editor first and then post on Lets HPC
- Serial : Time implementation, whether dependency or not.



$$\text{Theoretical speedup} = \frac{T_s}{T_p} = \frac{n}{\log_2 n}$$

→ Ignore estimated serial fraction and upper bound and Karp-Flatt metrix analysis.

Number of memory access and computation.

Speedup ( $n, p_j$ )

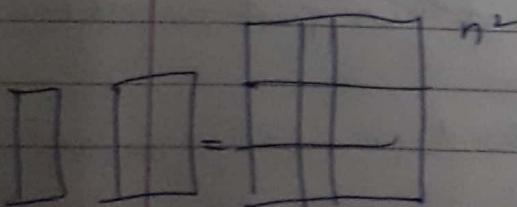
Problem size      Number of processor.

$$\text{Efficiency} = \frac{\text{Speed up}}{\text{Processor}}$$

→ Recursive Decomposition  $\rightarrow$  Quicksort  
(granularity is increasing)  $\rightarrow$

→ Data Decomposition

Output data decomposition:



$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix} \rightarrow \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix}$$

$$C_{11} = A_{11}B_{11} + A_{12}B_{21} \\ C_{12} = A_{11}B_{12} + A_{12}B_{22}$$

We can further divide the task further but the concurrency will be equal.

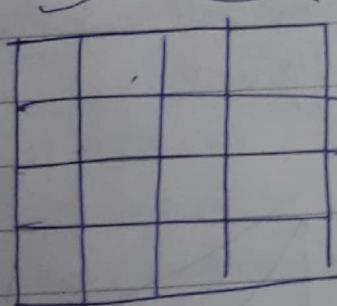
More concurrency will be achieved but more memory will be required.

### Block Matrix Multiplication

$n = 10000$

$$10000 \times 8 = 80,000 \text{ bytes} = 80 \text{ KB}$$

$\frac{8}{4}$



Work in efficient parallel algorithm:  
if word complexity ( $P$ ) > val complexity ( $S$ )  
but step complexity ( $P$ ) < step complexity ( $S$ )  
can

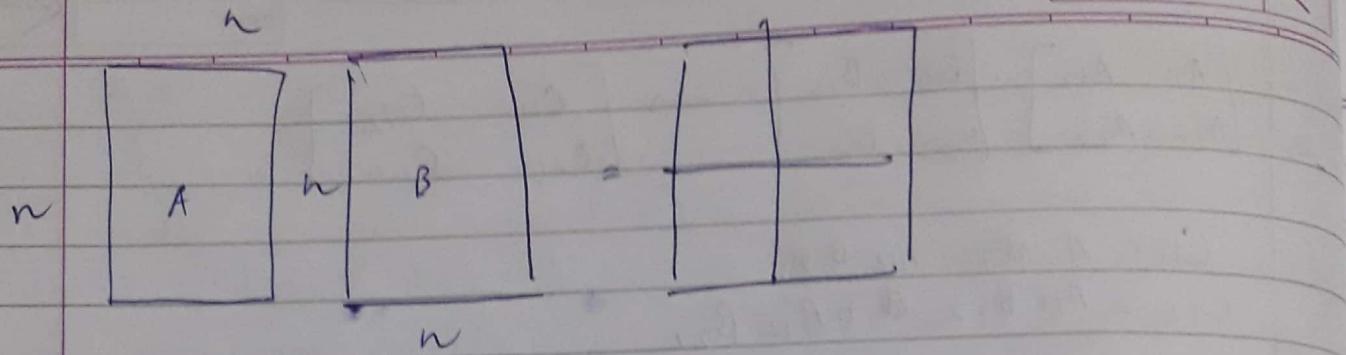
- Decide the size of the block based on your cache size

### Intermediate data decomposition

Divide the task into 8 parts and combine it again

$$D_{11} = A_{11}B_1$$

$$D_{12} = A_{12}B_{21}$$



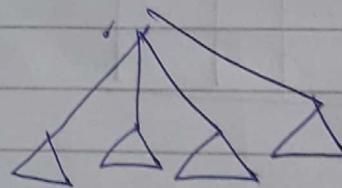
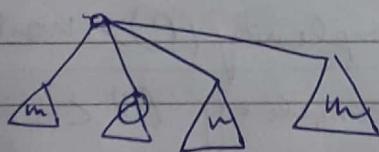
If one processor assigned a data then associated with all the computations associated with it.

Sorting an array; input data decomposition.

Exploratory decompositions:

Based on Since the program is non deterministic work complexity will be different.

Anomaly



Serial work:  $2m+1$

Parallel work: 1

Serial work:  $m$

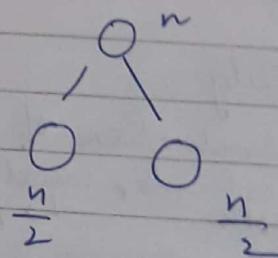
Parallel work:  $4m$

Speculative Decomposition: Don't know where to go next

Hybrid decomposition: mix of data and recursive decomposition.

Commands

screen top ps &amp; ssh scp



$$\begin{cases} \frac{n}{b^n} = 1 \\ l=2 \end{cases}$$

$$\log n = n \log b$$

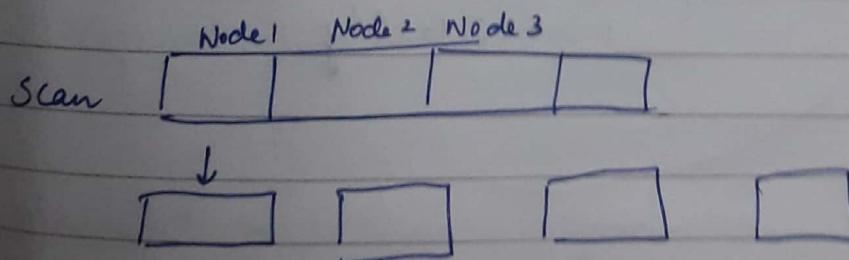
$$n = \frac{\log n}{\log b} = \log_b n$$

Work complexity = sum of all steps.

Step complexity :  $\log_b n$

### Scan Algorithm

$$\begin{aligned}
 & n + (n-1) + (n-2) + \dots + \\
 & n \log n - (1+2+4+\dots) \\
 & n \log n - (2n-1) \\
 & = O(n \log n)
 \end{aligned}$$



## Performance Analysis

- 1) Andahl's Law  $\rightarrow$  possibility
- 2) Gustafson-Barsis Law  $\rightarrow$  evaluation of a parallel program
- 3) Karp-Flatt Metric  $\rightarrow$  Overheads, serial, parallel  $\rightarrow$  problem size  
 ↳ both      ↳ problem size
- 4) I/O - efficiency analysis.  
 ↳ scalability

Sequential Time  $\rightarrow \sigma(n)$

Problem size  $\rightarrow n$

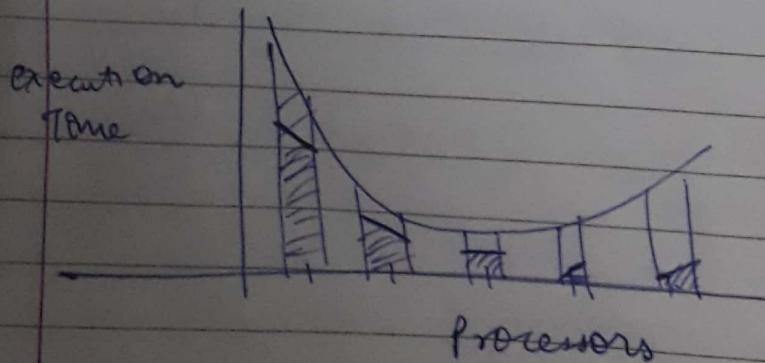
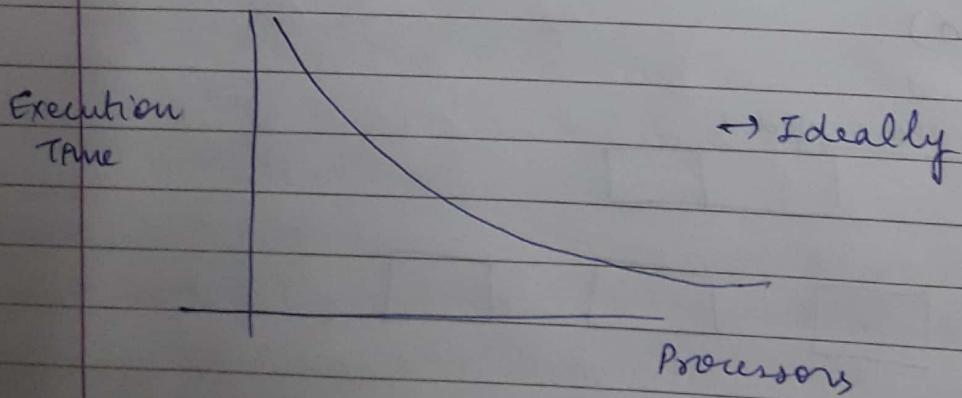
P  $\rightarrow$  processor

Parallel Time  $\rightarrow \phi(n)$  [Total time]

Overhead  $\rightarrow K(n, p)$

$$\text{Speedup } \Psi(n, p) = \frac{T_s}{T_p} \leq \frac{\sigma(n) + \psi(n)}{\sigma(n) + \psi(n) + K(n, p)}$$

$p \uparrow$  = communication time  $\uparrow$  + computation time  $\downarrow$



■  $\rightarrow$  computation  
 □  $\rightarrow$  communication  
 fixed problem size

For a given problem size there is an optimal number of processors that will give you.

$$\text{Efficiency} = \left( \frac{T_s}{T_p} \right) \perp$$

$$\boxed{0 \leq E \leq 1 \\ (n, p)}$$

Amount of processor utilization.

$$\begin{aligned} E(n, p) &\leq \left( \frac{\sigma(n) + \varnothing(n)}{\sigma(n) + \varnothing(n) + \kappa(n, p)} \right) \perp \\ &= \frac{\sigma(n) + \varnothing(n)}{p\sigma(n) + \varnothing(n) + p\kappa(n, p)} \end{aligned}$$

In an ideal code,  $\sigma \rightarrow 0, \kappa \rightarrow 0$

$$E(n, p) \leq \underline{\underline{1}}$$

In worst case

$$\geq \underline{\underline{0}} \quad \frac{\sigma(n)}{p(\sigma(n) + \kappa(n, p))} = \underline{\underline{0}}$$

$$\boxed{\therefore 0 \leq E(n, p) \leq 1}$$

Amdahl's Law:

$$\Psi_{(n, p)} \leq \frac{\sigma(n) + \varnothing(n)}{\sigma(n) + \frac{\varnothing(n)}{p} + \kappa(n, p)}$$

$$\Psi_{(n, p)} \leq \frac{\sigma(n) + \varnothing(n)}{\sigma(n) + \varnothing(n)/p}$$

Analyze Amdahl's Law as [not very practical]  
as not considering overhead

\* Fixed problem size

Page No.:

Date:

Yours

$$f = \frac{\sigma(n)}{\sigma(n) + \Omega(n)} \rightarrow \text{serial part}$$

$$\Psi(n, p) \leq \frac{\sigma(n) + \Omega(n)}{\sigma(n) + \frac{\Omega(n)}{p}} = \frac{\sigma(n)}{\sigma(n)} f = \sigma(n) f$$

$$\Omega(n) = \frac{\sigma(n) - f\sigma(n)}{f} = \frac{\sigma(n)(1-f)}{f} = \sigma(n)(\frac{1}{f} - 1)$$

$$\Psi(n, p) \leq \frac{\sigma(n)/f}{\sigma(n) + \frac{\sigma(n)(1/f - 1)}{p}} = \frac{1/f}{1 + \frac{1-f}{p}}$$

$$\boxed{\Psi(n, p) \leq \frac{1/f}{1 + \frac{1-f}{p}}}$$

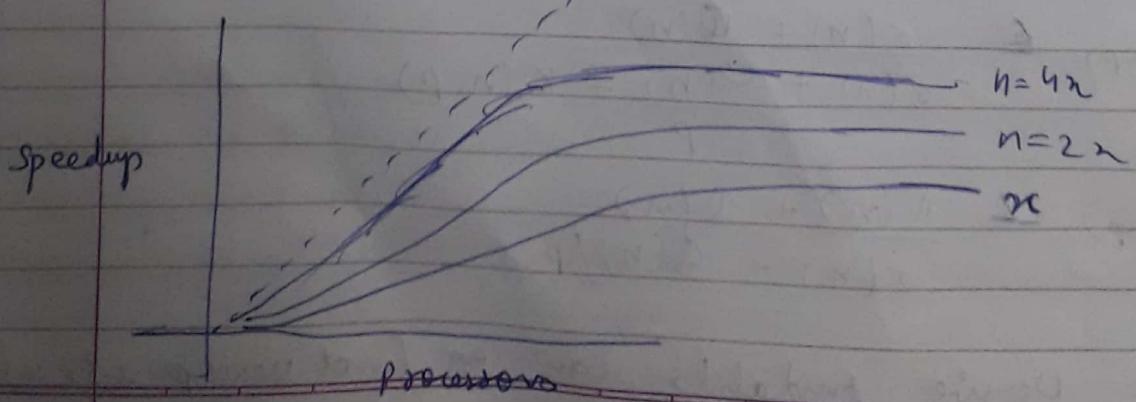
If  $f \uparrow$  then  
poor speedup

If  $f \downarrow$  then good  
speedup

\* However bad analysis since not considered overhead

- Amdahl's law gives us an upper bound on speedup with a fixed problem size and neglecting overhead.

$\Psi(n, p)$  grows slower than  $\Omega(n)$  for the problem size



For a fixed number of processors speedup increases as problem size increases (Amdahl's effect)

→ Do serial matrix v/s block matrix w.

2 ways to do BMM :

$$\psi(n, p) \leq \sigma(n) + \frac{\phi(n)}{p} + \kappa(n, p)$$

Gustafsson - Barsis Law : // Here too we neglect the overhead

→ fraction of time spent in the serial parallel code performing serial work.

$$s = \frac{\sigma(n)}{\sigma(n) + \frac{\phi(n)}{p}}$$

$$1-s = \frac{\phi(n)/p}{\sigma(n) + \phi(n)/p}$$

$$\sigma(n) = s\sigma(n) + \frac{s\phi(n)}{p}$$

$$\phi(n) = (\sigma(n) + \frac{\phi(n)}{p})(1-s)p$$

$$\boxed{\sigma(n) = \frac{s\phi(n)}{p(1-s)}}$$

$$\psi(n, p) \leq \left( \sigma(n) + \frac{\phi(n)}{p} \right) (s + (1-s)p)$$

$$\sigma(n) + \frac{\phi(n)}{p}$$

$$\boxed{\psi(n, p) \leq p + (1-p)s}$$

performance

// The speedup we can get by putting p number of processors.

## Karp Flatt Metric

overhead:

- ↳ creation of threads. → process set up.
- ↳ communication
- ↳ extra computation
- ↳ synchronization.

Experimentally determined serial fraction.

$$\rightarrow T(n, p) = \sigma(n) + \frac{\phi(n)}{p} + K(n, p)$$

$$\rightarrow T(n, 1) = \sigma(n) + \phi(n) \rightarrow \text{serial}$$

~~e~~ Total amount of ideal  
 ↓ idle      ↓ overhead

$$e = \frac{(p-1)\sigma(n) + K(n, p).p}{(p-1)T(n, 1)}$$

$$p\sigma(n) - \sigma(n) + pK(n, p)$$

$$= pT(n, p) - e(n) - \sigma(n) = .pT(n, p) - T(n, 1)$$

$$pT(n, p) = e(p-1)T(n, 1) + T(n, 1)$$

$$pT(n, p) = T(n, 1)eP - T(n, 1)e + eT(n, 1)$$

$$T(n, p) = T(n, 1)e + T(n, 1)\left(1 - \frac{e}{p}\right)$$

$$T(n, p) = T(n, p)e + T(n, p)\left(\frac{1-e}{p}\right)$$

$$\frac{\frac{1}{3} - \frac{2}{5}}{\frac{2}{3}} =$$

$$\left\{ \Psi = \frac{T(n, 1)}{T(n, P)} \right.$$

4 → 0.99

Page No.:  
Date:

YOUVA

$$1 = \Psi e + \Psi(1-e)/P$$

$$\frac{1}{\Psi} = e + (1-e)/P$$

$$\Rightarrow \frac{1}{\Psi} = e + \frac{1}{P}$$

$$\left[ e = \frac{1/4 - 1/P}{1 - 1/P} \right]$$

P	2	3	4	5	...	8
$\Psi$	1.82	2.5	3.08	3.57		4.71
e	0.1	0.1				0.1

→ Ideal time remains same when e is constant  
 Overhead is zero  
 And we cannot improve at our code.  
 // problem in the serial part

P	2	3	4	5	...	8
$\Psi$	1.87	2.61	3.2	3.73		4.71
e	0.07	0.075	0.08	0.085		0.07

When the overhead is very high

## Isoefficiency Metric

$$V(n, p) \leq \frac{\sigma(n) + \phi(n)}{\sigma(n) + \phi(n) + K(n, p)}$$

Able to increase performance as number of processors increase.

$$\leq \frac{p(\sigma(n) + \phi(n))}{\sigma(n) + \phi(n) + (p-1)\sigma(n) + K(n, p)}$$

overhead.  $T_0(n, p) = (p-1)\sigma(n) + pK(n, p)$

$$E(n, p) \leq \frac{\sigma(n) + \phi(n)}{\sigma(n) + \phi(n) + T_0(n, p)}$$

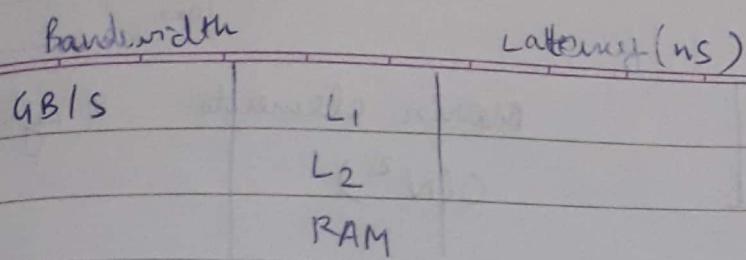
$$E(n, p) \leq \frac{1}{1 + \frac{T_0(n, p)}{\sigma(n) + \phi(n)}}$$

$$\frac{T_0(n, p)}{T_0(n, 1)} \leq \frac{1 - E(n, p)}{E(n, p)}$$

$$T(n, 1) \geq \frac{E(n, p)}{1 - E(n, p)} T_0(n, p)$$

Rate at which the problem size has to be increased so that efficiency remains constant.

Slope of LHS much greater than RHS



Latency increases by a factor of 10 at every step

Cache hit rate time = 1 cycle      99% hit rate

miss penalty = 100 cycle

$$\text{Average access time} = \frac{99}{100} \times 1 + \frac{1}{100} \times 100 = \underline{1.99} \approx 2$$

~~80% miss~~  $\frac{97 \times 1}{100} + \frac{3 \times 100}{100} = \frac{397}{100} = 3.97 \approx 4$       97% hit rate

(2 times slower)

than 99% hit rate

80% hit rate

$$\frac{80 \times 1}{100} + \frac{20 \times 100}{100} = 20.8$$

Miss rate = 1 - cache hit

\*\* Miss penalty is nearly around 100 times than a hit

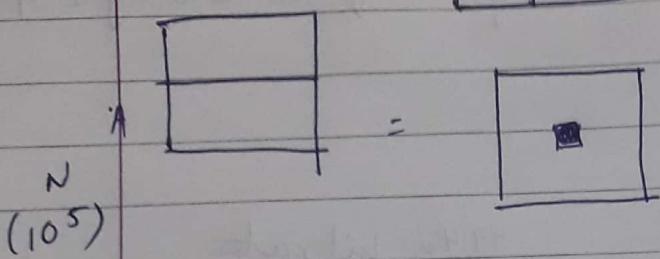
L1 - 32K

L2 - 256K

L3 - 20 MB (cluster)  
6 MB (lab PC)

B

Matrix elements = 8 bytes  
 $O(N^2)$



Cache block size ( $c$ ) = 64 bytes

Cache hit miss

$$= \frac{\Delta}{8} + N = \frac{9N}{8} = O\left(\frac{N^2 \times 9N}{8}\right) = O\left(\frac{9N^3}{8}\right) \approx O(N^3)$$

(Row)

Cache miss for serial

Block size  $\sim B^2 = (B \times B)$

Cache can hold around 3 such blocks.

~~$3B^2 \leq C$~~ 

$$3B^2 \leq C$$

Cache miss for each block.

$$\sim \frac{2N}{B} \times \frac{B^2}{8} = \frac{NB}{4}$$

$$\text{Number of blocks} = \left(\frac{N}{B}\right)^2$$

$$\text{Cache miss for each blo} = \frac{N^3}{4B} = O(N^3)$$

Larger the block size quicker the processing.

$$\text{Faster Speedup} = \frac{9N^3}{8 \times B^3} \times 9B = \underline{\underline{4.5B}}$$

Lessons:

- 1) keep your working dataset very very small
- 2) keep your stride size small.

Problem:

$T_p$

$n =$  processor takes  $(\log n)^2$  times  $\Rightarrow$  parallel  
serial time is  $\Rightarrow n \log n$   $\Rightarrow$  serial

- ① speedup
- ② efficiency

$$\text{Speed} = \frac{n \log n}{(\log n)^2} = \frac{n}{\log n}$$

Speedup  $\uparrow$  problem size  $\uparrow$

$$\text{Efficiency} = \frac{n}{\log n} \times \frac{1}{n} = \frac{1}{\log n}$$

Efficiency  $\downarrow$  problem size  $\uparrow$

Cost complexity  $\Rightarrow n(\log n)^2 \Rightarrow$  Not cost optimal by a factor of  $\log 2$

Let make  $p \leq n$

Now the parallel time will be  $\frac{p}{n} (\log n)^2$

$$\text{Speed up} = \frac{p}{\log n}$$

speedup decrease as problem size increased since it is not cost optimal