

High Performance Computing

Performance →

Algorithms	→ Computer
Deterministic	Non deterministic
Abstract	Not abstract
Always Terminates	Might not terminate.

Unit of Performance - Flops / Sec

$$= \text{Frequency} \times \text{Number of Cores} \times \text{Flops per cycle}$$

$$= \frac{2.5 \times 10^9}{\text{cycle sec}} \times 4 \text{ Flops / Cycle}$$

$$= 10 \text{ GFLOPS / Sec} \quad \leftarrow \text{Serial Programming}$$

output that
reaches
(Theoretical Limit)

Performance is determined by :-

- 1) Architecture → Cache
- 2) Algorithm
- 3) Software Environment

- Scalability - Able to utilise higher numbers of cores

- Portability - Must run on other architectures too.

For calculation of performance of a code, it must take into account the memory access time required (As we have a lot of data which needs to be loaded)

- Systems Perspective - Architecture, Algorithm and Software Environment

Non deterministic

Deterministic

Software Environment

Algorithm

and Software Environment

Get

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

- A node can have one or more sockets

Each socket has a processor

Each processor can have multiple cores

Each core performs Flops

Several such nodes

= Supercomputer

④ Properties of Memory - Latency and Bandwidth → (MB/s)
 Increases in Order - Secondary, RAM, Cache, Registers to Processor

- Serial - Only one instruction at any moment in time.
 Improvement Latency (How fast you can compute 1 instr.)
- Parallel - Solving concurrently
 Throughput is good (How fast n instructions can be computed)
- Speedup = $\frac{\text{Serial Time}}{\text{Parallel Time}}$

* Multiplicity :- There are multiple memory

- Multiple data paths
- Multiple processor

- Types of parallelism - Node Level
 Thread Level

Instruction level → Pipeline



④ A

$$\text{Throughput} = \frac{1}{\text{Latency}} \text{ instructions / cycle}$$

- 1 result / cycle after wind up phase
- Depth of pipeline = m stages
- Total independent operations = N
- Total time of the pipeline = $N + m - 1$
 (Without pipeline) - Serial = mN

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

If $N \gg m$ then Speedup = m

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

when $N \gg m$, Throughput $\approx \frac{1}{N}$

Thus for ~~for~~ smaller problem size, throughput won't be high hence there would be lower performance levels.

- If P_1 requires 40 bytes of data; P_2 requires 1 KB, P_3 requires 1 MB then all of them will take the same access time.

Latency = Time required to access 0 bytes of data.

- Von Neumann Bottleneck :- Not able to provide data to the process at the desired rate ie The rate at which processor processes the instruction.

- Vector Triad = 1 Multiplication + 2 addition - $a + b * c$

- Plot the graph of Performance vs N

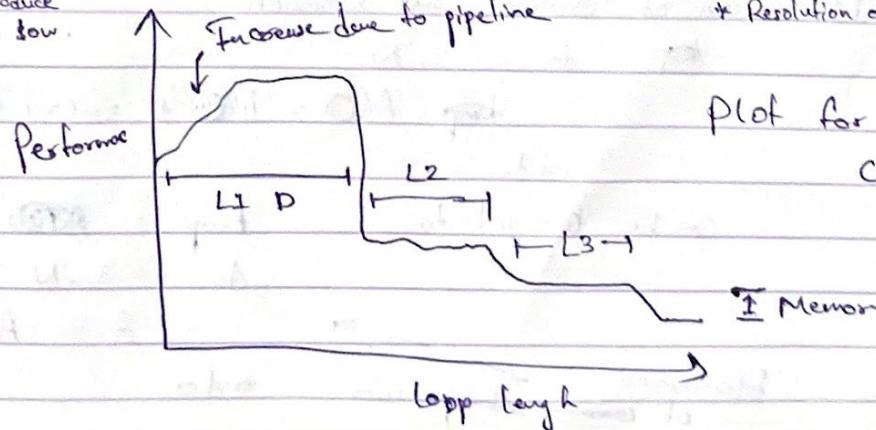
pipeline cannot produce

of 1. Thus it is slow.

low throughput to 1.

L2 is low

L2 is higher
in bus channel



(Measure Time
Timer function for elapsed
* Resolution of times.

Plot for single
core

Memory

Addition and Multiplication have very short pipeline. Thus both

Machine Balance	Code Balance
- Machine Balance = $\frac{\text{Memory Bandwidth}}{\text{Peak performance}}$	Code Balance of loop = $\frac{\text{Data traffic}}{\text{Number of operations}}$
Unit = GB / FLOPs	$= \frac{4(A+B+C)}{2}$
- Difference for cache, RAM, disk etc.	
- Historically decreasing trend, will decrease in future too due to increase in peak performance (More cores)	$= 2$
- Code Balance of Vectorized Reciprocal of code balance = Computational Cost	
- Lesser the value of code balance, more efficient will be the code.	
Lecture Miss on 15 th Jan	Granularity = Ratio of computation to communication
<u>Simple Optimization of Serial Code</u>	Coarse \rightarrow Large computation between stages Fine \rightarrow Small computation between stages

- Using Exit \rightarrow Lesser computations of non required things
- ~~A = A + B * 2.0~~ \times \rightarrow ~~A + B * p * (2 ln B)~~
- ~~A = A + B * 2.0~~ \times \rightarrow ~~A + B * p * (2 ln B)~~
~~A = A + B * B~~ \checkmark Better because pipeline of multiplication is lower. (Lower stages in pipeline)
- Declaration in tmp of things which is computed over and over
- Elimination of common subexpression

~~tmp~~ do i = 1, N
~~tmp~~ A(~~i~~) = A(~~i~~) + sin(x) ~~* i~~
end do

can be changed to

tmp = A(~~i~~) sinx + A
do i = 1, N

A = A + tmp * i

Elimination of common subexpression

end do

- Single precision vs Double precision

Single is better as size of single is less, more things can be stored in cache \Rightarrow Higher Cache hit ratio.

Q. Which is better $B^k B + B$ 40 times or B^{10}

- Avoid Conditional Branching \rightarrow Reason - Branch Miss can happen.

The compiler predicts that the result of \oplus if conditionals based on statistical algorithms of compiler. It will start calculation based on its assumption. But if the condition fails then it has to empty its pipeline and re-start computation.

- do i = 1, N, 2

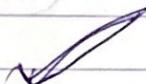
do j = 1, N

$$c(i) = a(i,j) + b(j)$$

$$c(i+1) = a(i+1,j) * b(j)$$

end do

end do



do i = 1, N

do j = 1, N

$$c(i) = a(i,j) + b(j)$$

end do

end do



First is better because of Loop Unrolling. The c array will be loaded in the memory. \Rightarrow Higher cache hit.

If $c(i)$ is loaded it will load $c(i), c(i+1) \dots (i+n-1)$ if $a(j,i)$ and $b(j)$ loads up the cache then $c(i)$ will be needed to be loaded again in the cache.

- Organisation of Parallel Platforms



Logical

↓
Programming view

Physical

↓
Architecture

Distributed
(Shared / Master)

Expressing parallel task

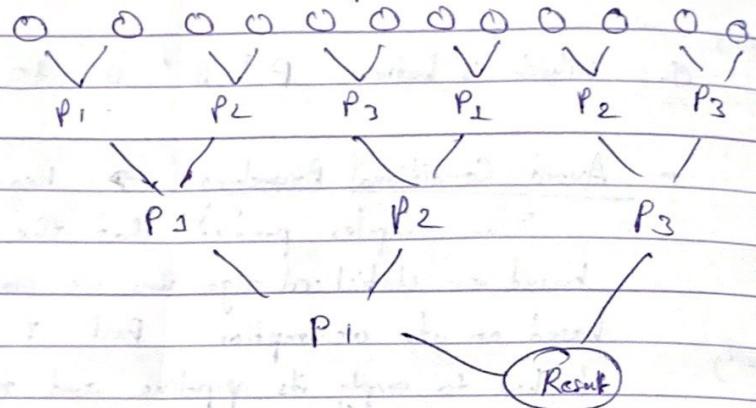
→ Mechanism

for communication

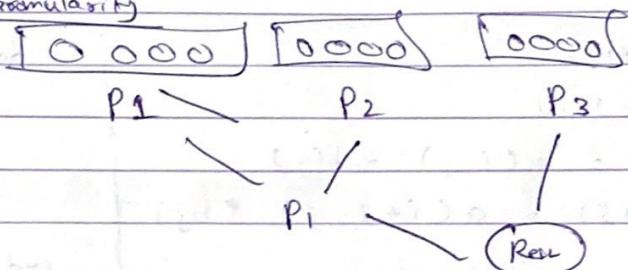
↳ Several models
(Control Structure)

- Shared Memory Parallelism

- Fine Granularity for N addition (3 processor)



- Coarse Granularity



For matrix multiplication of $n \times n$ matrix. There are 4 processes

Finest granularity = n^2 (Each processor doing 1 multiplication of cell)

1 row 1 column by a processor \rightarrow Granularity level - n
 $n/4$ rows ~~columns~~^{wh} by a processor \rightarrow Granularity $\rightarrow n/4$

OpenMP \rightarrow MP stands for multi processing

Open MPI \rightarrow Message Passing \rightarrow For Distributed system

Q) 64 bit DDRAM interface, 8 channels, 1 GHz clock speed.

Access latency of 200 cycles - ?

Ans Access Throughput = $\frac{8 \text{ bytes}}{\text{bytes/sec}} \times \frac{8 \text{ channels}}{\text{transfers}} \times \frac{2 \text{ transfers}}{\text{clock}} \times 10^9 \text{ clocks/sec}$

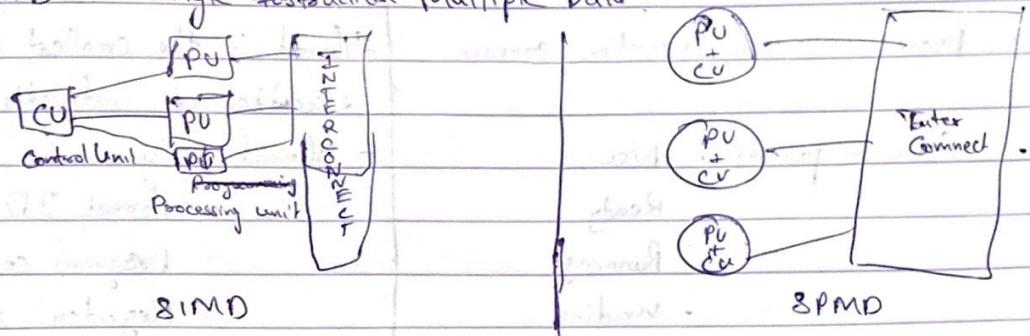
$$= 128 \text{ GBps}$$

Double Data

~~200~~

- SIMD → Single Instruction Multiple Data

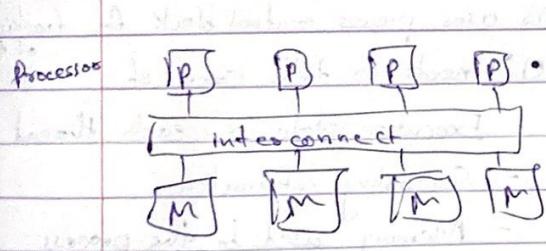
SPMD - Single Instruction Multiple Data



Q. Which is faster → Ans - Depends on the program.

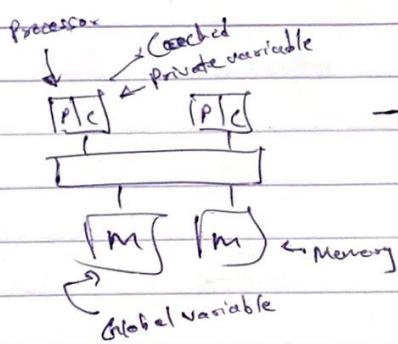
- Global Memory → Easy to write / Difficult to write.

- UMA → Uniform Memory Access | NUMA → Non-uniform Memory Access



Each processor will take
some time for process

(Shared)

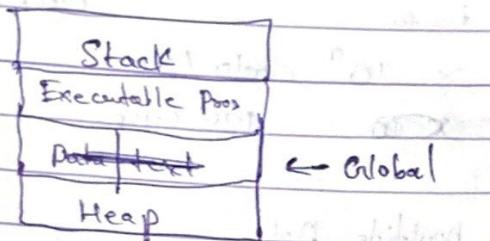


(Distributed)

(Memory can be still
accessed through interface
but will require lot of time)

~~Stack is faster than heap~~

- Process :- Running Program
 - These will be context switch
 - All process has own id. - pid
 - Status of process - Waiting, running etc



- Process is an execution stream

- Status of process -
 - New
 - Ready
 - Running
 - Waiting
 - Terminated

Threads :- Lightweight process

They run concurrently
Thread too has id.



- Thread is the smallest unit of sequential execution of instructions

- Thread includes
 - Thread ID
 - Program counter.
 - Register state (pointer)
 - Stack (pointer)

- Shared Resources

- Text section

- Data section

- Other OS resources

OS uses process control block for tracking of threads

OS needs to take care of

- Execution state of each thread
- Scheduling information
- Memory used by the process
- Information about the open file.

- Scheduling of SIMD is easier than SPMD because in SPMD there are multiple control unit which needs to be scheduled accordingly.

* Dependencies

\\$

for $i = 1, n$
 for $j = 2, m$
~~Data Dependencies~~ $b[i,j] = \dots + b[i+j-1]$ // j is dependent.

- Types of dependencies - Flow, Anti, Input & Output

~~Flow/True Statement~~

Statement ' i ' precedes ' j ' and computes the values that ' j ' uses

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

Number of flow dependencies - $\textcircled{1} \rightarrow \textcircled{2}$
 $\textcircled{2} \rightarrow \textcircled{4}$

~~Anti~~

- Statement ' i ' precedes ' j ' and ^{uses} computes the values that ' j ' computes

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

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

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

~~Output~~

- . Statement ' i ' precedes ' j ' and ' i ' computes a value that ' j ' also computes

$$\textcircled{1} \rightarrow \textcircled{3} \quad | \quad \textcircled{1} \rightarrow \textcircled{4} \quad | \quad \textcircled{3} \rightarrow \textcircled{4}$$

~~Input~~

- . Statement ' i ' precedes ' j ' and ' i ' ^{uses} computes a value that ' j ' also uses.

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

- If dependence flows from ' i ' to ' j '

Source \downarrow Sink

Flow / True dependency is difficult to remove

- In data dependency graph, nodes are the statements, whereas edges are the dependency relations

Eg - $a[i] = b[i] + c[i]$
 $d[i] = a[i]$

There is flow dependency but no loop carried dependency

Dependence distance = 0

- $a[i] = b[i] + c[i]$

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

This has loop carried dependency \rightarrow Dependence distance = 1.

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

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

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

$c[j] = a[j]$

Can be parallelized

Locality is worse

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

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

$\& c[j] = a[j]$



Cannot be parallelized

Locality is better

We use Loop Fusion so as to get better utilisation of spatial locality

We use Loop Distribution so that we can use parallelism.

- Directives : - Special preprocessor instructions

- For shared systems - All processes share common memory
 - Easy to code.
 - Bus contention leads to poor scalability

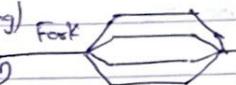
- In distributed, each processor has private memory
 - Scales well
 - Difficult memory management

- Fork - Join

Fork is called by thread (parent) to create new thread (child).

Fork helps in concurrency (increase concurrency)

Join removes concurrency (decrease concurrency)



All must unite
before continuing.
(Explicit Barrier)
(Implicit Barrier)

Parent must join with child

- A structure block :- There is one point of entry, one point of exit
Inside structure block & if condition is not allowed.

- In general branching is not allowed
- Check if exit() is allowed or not.

A structured block can be executed in parallel. (mostly)

- Scope → Access (Who can access it)
 - how long is it visible.

- Performance of Shared Memory depends on

- Load Balancing (mapping work loads with thread scheduling)

- Caches

- Locality

- How locality affects scheduling algorithm

↳ 1003 to 4 threads
250 each.

The rest 3 can be given to the one who finds the 250 tasks first.

- Work Sharing directives just shares the work, does not create threads

pragma omp parallel ← thread creates

pragma omp for ← work sharing

if the variable is in for loop then it is converted to private.

pragma omp for

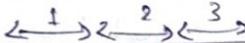
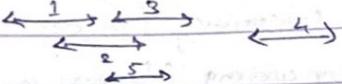
for (i=0 ; i < N ; i++)

Since work sharing is used.

The i will be taken as private to the thread

②

* Index variable of a worksharing loop is private by default.

- Synchronous -  (Block Mechanism)
- Asynchronous -  (Non-Blocking)

- Work sharing - for directive \rightarrow shares iterations

section directive \rightarrow

critical directive \rightarrow serializes section of code.

\rightarrow For work sharing directives:

- No implicit barrier at entry. (No need to wait for others)
- Explicit barrier at exit. (Every thread waits for process to complete)
- No wait clause is used to remove this barrier.

③ Scheduling Mechanism for For \rightarrow static, dynamic, guided.

Default size is 1. chunk.



↳ Barrier
Threads can know down only after

- Thread Synchronisation :-

Implicit Event Synchronization

Explicit Synchronisation - critical, master directives in OpenMP

- critical

- atomic, single, master, barrier, flush,
omp_set_lock().

- OpenMP - Reduction.

- It is binary operators. | Operator

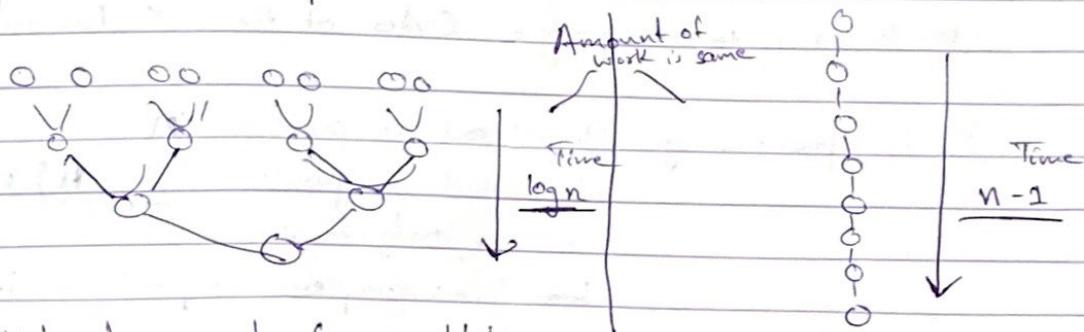
(Floating point operation
 \neq Associative)

- Identity element Operation must be associative

- Identity element must be present
- In reduction if it is shared as well as private.
→ Opt Best performance will be obtained.

* Module - 3 - Parallel Algorithm Design

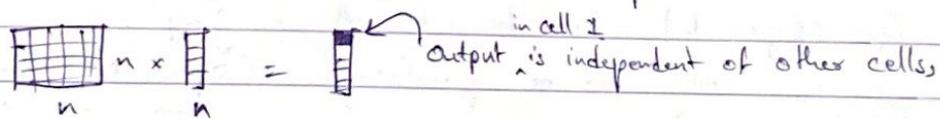
- Dividing program problem into task that can be executed concurrently / in parallel. → Decomposition



Task dependency graph for n additions in parallel

Task dependency graph for n addition in series.

General Rule :- Map the output to the processes.



- For good design, control dependencies must be removed

Number of tasks that can be executed in parallel → degree of concurrency
It increases if granularity is low and vice-versa.

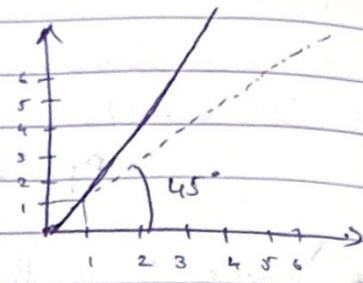
- Each node is the task, each edge is a dependency.
Nodes should be more, edges should be less.

Note: Task dependency graph deals with control dependency.

- Superlinear Speedup

Not possible in theory.

But it is possible because as we introduce more processors, we introduce more amount of cache, thus it is possible to have superlinear speedup.



Serial cost / time complexity = Order of time

Parallel cost / time complexity = Order of time \times No of processors

In the previous eg., Work cost in serial = n

Work cost in parallel = $(\log_2 n) n$

Time complexity in serial = n

Time complexity in parallel = $\log_2 n$.

- Work = Total no of nodes / break / vertices if all tasks do same amt of work, else summation of work done by each.

- Span / Depth (D) : - longest path in graph.

$\frac{W}{D}$ = Avg. work of parallelism. (Optimal number of processors to keep all of them busy on average)

- $T_p(n) \geq D(n)$ || p is number of processes. || D is depth of graph
 T_p = Execution time.

The minimum execution time is the depth of the graph

* Work Span Law : $T_p(n) \geq \max \left\{ D_n, \left\lceil \frac{W(n)}{p} \right\rceil \right\}$

- Ratio of total work to critical path length. = ~~No. degree of concurrence~~
 $= ?$

IP Address of HPC Cluster
→ 10.100.71.130

- Critical path length (E_J) $a \rightarrow 27, b \rightarrow$
- Each task → 10 time units.

$$W = 40 + 15 + 8 = \frac{63}{\begin{array}{c} 7 \\ 27 \\ 3 \end{array}} \approx 3 \text{ processors.}$$

Types of Dependencies

* Task Interaction Graph

- Data
- Loop
- Control
- Structural

If the entire data is not replicated across all the tasks, they will have to communicate elements of the vector.

$$A \begin{bmatrix} & & \end{bmatrix} B \begin{bmatrix} & & \end{bmatrix} = \begin{bmatrix} & & \end{bmatrix}$$

Eg. Multiplication of sparse matrix with vector

- Only non zero elements of A takes place in computation
- Matrix B is divided into partitions. Each stored distributedly

Task 0 has 1st row of A and B[0][j] and need to compute. Thus it will require B[1][j], B[4][j], B[8][j] as A[0][1][j], A[0][4][j] and A[0][8][j] is non zero.

Note Task Interaction Graph will be used to deal with data dependencies

- Helps in Load Balancing

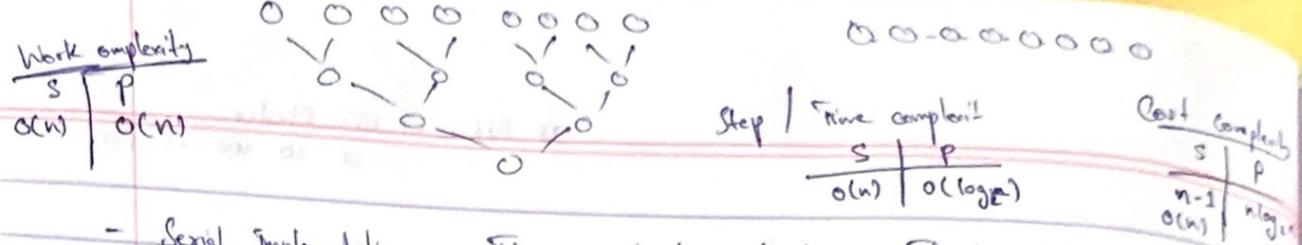
Note As granularity increases (finer) the ^{cost of} overhead computation increases.

- Mappings are determined by both task dependency and task interaction graph
- Task Dependency graph → Load Balancing.
→ Optimise the use of processors

Cab Access Cluster through

SSH ID @ 10.100.71.130.

SSH ~~GICS@~~ GICS@



- Serial Implementation - Time Complexity of algo , If dependency or not.
- Implementation of approach - Can add mapping graph
- Complexity of parallel .
- Cost of parallel .
- Theoretical Speedup = $\frac{T_{\text{serial}}}{T_{\text{parallel}}} = \frac{n}{\log n}$
- Estimated Serial Fraction - \times
- Upper bound on Amdahl - \times
- Number of memory Access - Inside code , put a counter . , Find code bkg
- Number of computation - Can be found by counters .
- Time curve (Constant problem size , very ~~number~~ of process)
- Time curve (Constant processor , very ~~problem size~~) -
- Speed up analysis .
- Efficiency curve \rightarrow Speedup / Number of processor (If actually den)
- Kaop Flatt metric $\rightarrow \times$
- Further Detailed $\rightarrow \times$ (First - 3 sessions)

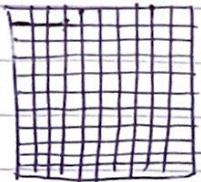
Decomposition Techniques

- Recursive
- Divide and Conquer \rightarrow Decomposed to set of sub prob.
- Imp
 - Data :
 - Exploratory
 - Speculative

- \rightarrow In recursive , the granularity increases
- In reduction , granularity decreases

- Data Decomposition :- Partition the data on which computation are performed , across various tasks

- Input Decomposition \rightarrow Output not known.
- \rightarrow Output Decomposition \rightarrow Each element of output can be calculated independently of other task.
- Intermediate Decomposition -



\rightarrow For it to be multiplied.

1st row with 1st column, 2nd column.

But there is a possibility that cache will be filled and 1st few elements will be removed.

Thus if we have divided it then cache ~~with~~ hit will be more.
Also work complexity remains the same. The size of block will be based on cache size.

- Work inefficient if work complexity (parallel) $>$ work complexity (serial)
This can be valid even when step complexity (parallel) $<$ step complexity (serial).

- Intermediate Decomposition of matrix multiplication

Task dependency \rightarrow ① ⑤ ③ ② ⑦ ⑥ ④ ⑧ ⑨ ⑩ ⑪ ⑫
 ② ⑥ ⑪ ⑫

* Owner Computer Rule.

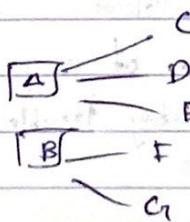
- Process assigned a particular data is responsible for all computation associated with it

* Exploratory Decomposition

- Decomposition goes hand-in-hand with its execution.

Eg - Game playing etc \rightarrow Non deterministic systems.

* Speculative Decomposition :- Assumes that one of them would be ^{coarse}

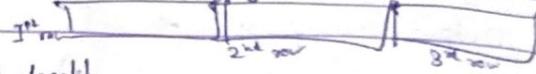


Starts processing C

as if speculates that C will be executed

- For 2-D array - It is stored in memory as flat.

Thus there is problem of spatial locality.



In Lab → Check the schedule directives

- Scan Prefix sum

Given array - $A = \{a_0, a_1, \dots, a_n\}$

$$\text{Scan}(A) = [I, a_0, a_0 + a_1, a_0 + a_1 + a_2]$$

This is difficult to parallelise → Each element depends on its previous.

Scan does $\log n$ parallel iterations.

$$(n-1) + (n-2) \dots \log n \text{ terms}$$

$$\Rightarrow n \log n - (n-1) \approx O(n \log n) \text{ Complexity}$$

Thus work complexity is more than serial ⇒ Work inefficient.

Work efficient → Build balanced binary tree → Blelloch

→ Problem if distributed



These depends on the previous, becomes serial!

Scan can be used in compaction.

1 2 3 4 5 6 7

0 1 0 0 1 1 0

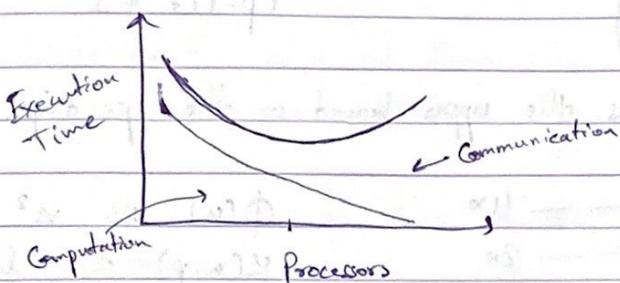
Scan 0 1 1 1 2 3 3

2 | 5 | 6 | 1 | 3 | 4 | 7 | 0.

Lab

- 1) Amdahl's Law → possibility of parallelization
- 2) Gustafsson-Baeris Law - Evaluation of a parallel program
- 3) Kepp-Flatt Metric - Which is creating trouble? Serial or Overhead
- 4) Isoefficiency Metric.

- Sequential Time - $\sigma(n)$ $n \rightarrow$ problem size
- Parallel Time - $\phi(n)$ $p \rightarrow$ processor
- Overhead $\rightarrow K(n, p)$
- Speedup ~~(n, p)~~ $\rightarrow \Psi(n, p) = \frac{T_p}{T_s} \leq \frac{\sigma(n) + \phi(n)}{\sigma(n) + \frac{\phi(n) + K(n, p)}{p}}$
- If p increases Communication time increases
Computation time decreases.



$$\text{Efficiency} = \left(\frac{T_p}{T_s} \right) \times \frac{1}{p}$$

Amount of processor utilization

$$0 \leq e \leq 1$$

$$E(n, p) = \frac{1}{p} \left(\frac{\sigma(n) + \phi(n)}{\sigma(n) + \frac{\phi(n) + K(n, p)}{p}} \right)$$

Thus efficiency $e \leq \frac{\sigma(n) + \phi(n)}{p(\sigma(n) + K(n, p)) + \phi(n)}$

Thus if $\sigma(n) = 0$ and $K(n, p) = 0$ Then $e \approx 1$.

if $\phi(n) = 0$ then $e = \frac{\sigma(n)}{p(\sigma(n) + K(n, p))}$

Thus $e \rightarrow 0$ when $p \rightarrow \infty$

• Amdahl's Law : - $\Psi(n, p) \leq \frac{\sigma(n) + \phi(n)}{\sigma(n) + \frac{\phi(n) + K(n, p)}{p}}$

Ignoring overheads

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

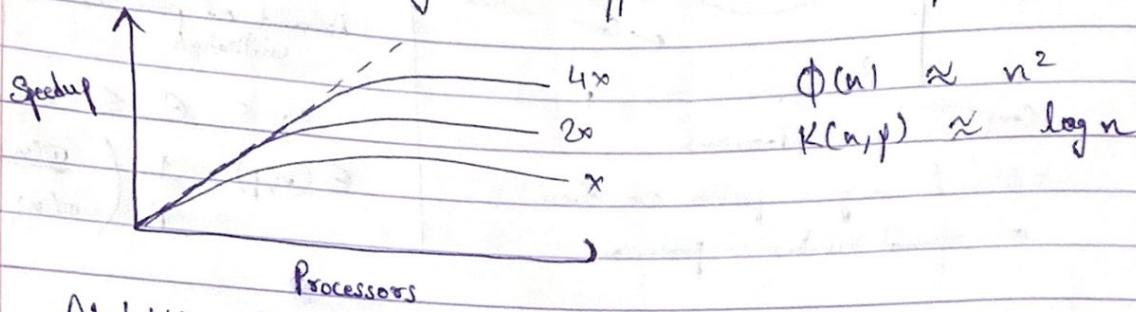
Let $f = \frac{\sigma(n)}{\sigma(n) + \phi(n)} \rightarrow$ Fraction of serial work

$$\therefore \psi(n, p) = \frac{\sigma(n) + \phi(n)}{\sigma(n) + \phi(n) + \frac{p}{p}} = \frac{\sigma(n) + \phi(n)}{\sigma(n) + \phi(n) + \frac{1}{p}}$$

$$\frac{f}{1-f} = \frac{\sigma(n)}{\phi(n)} \quad \therefore \psi(n, p) = \frac{1 + \left(\frac{1-f}{f}\right)}{1 + \left(\frac{1-f}{pf}\right)} = \frac{1}{f} \left| \frac{pf + 1 - f}{pf} \right.$$

$$\therefore \psi(n, p) = \frac{1}{\frac{f}{1-f} + \left(1 - \frac{1}{pf}\right)} = \frac{p}{(p-1)f + 1}$$

Thus Amdahl's law gives the upper bound on the speedup:



Amdahl's effect :- For fixed number of processors, the speedup increases as the problem size increases.

- Gustafsson - Barcis law :- Time is fixed

Let s = Fraction of time spent in the parallel code performing serial work

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

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

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

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

$$\therefore \Psi(n, p) \leq s + p(1-s)$$

$$\leq p + (1-p)s$$

- Overheads in parallel code are due to the following
 - Creation of threads → Process Setup
 - Communication
 - Extra Computation.
 - Coordination and Synchronization

- Experimentally determined Serial fraction

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

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

Experimentally determined Serial fraction

$$e = \frac{\text{idle time} + \text{overhead}}{(p-1)\sigma(n) + pK(n, p)}$$

$$\begin{aligned} \text{idle time} &= (p-1)\sigma(n) \\ \text{tot_overhead} &= pK(n, p) \end{aligned}$$

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

$$= \frac{pT(n, p) - \phi(n) - \sigma(n)}{(p-1)T(n, p)}$$

$$= \frac{pT(n, p) - T(n, 1)}{(p-1)T(n, p)}$$

$$\Rightarrow pT(n, p) = e \{ (p-1)T(n, 1) \} + T(n, 1)$$

$$\Rightarrow pT(n, p) = epT(n, 1) + T(n, 1) - eT(n, 1)$$

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

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

$$\text{As } \Psi = \frac{T(n, 1)}{T(n, p)}$$

$$1 = \Psi e + \Psi \left(\frac{1-e}{p} \right)$$

$$\Rightarrow \frac{1}{\Psi} = e + \frac{1-e}{p} \Rightarrow \frac{1}{\Psi} = (p-1)e + 1$$

$$\Rightarrow \left(\frac{p-\Psi}{\Psi} \right) \frac{1}{p-1} = e$$

$$\boxed{e = \frac{\frac{1}{\Psi} - \frac{1}{p}}{1 - \frac{1}{p-1}}}$$

p	2	3	4	5	8
Ψ	1.82	2.5	3.08	3.57	4.71
e	0.1	0.1	0.1	0.1	0.1

Since the value of e is constant, there is no overhead

\Rightarrow Problem with serial \rightarrow Maybe large computation need to be done

p	2	3	4	5	8
Ψ	1.87	2.61	3.2	3.73	4.71
e	0.07	0.075	0.08	0.085	0.1

Since the value of e increases, there is large overhead. Thus there is scope of improvement.

- For a given fixed problem size, the efficiency decreases as you increase the number of processors

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

Measure of the ability to increase the performance as the number of processors increases.

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

$\text{Overhead} = T_o(n, p) = (p-1)\sigma(n) + pK(n, p)$

Thus efficiency $E(n, p) = \frac{\sigma(n) + \phi(n)}{\sigma(n) + \phi(n) + T_o(n, p)}$

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

$$\Rightarrow \frac{T_o(n, p)}{T(n, 1)} \leq \frac{1 - E(n, p)}{E(n, p)}$$

$$\Rightarrow \left[\frac{T(n, 1)}{T(n, p)} \geq \left(\frac{E(n, p)}{1 - E(n, p)} \right) T_o(n, p) \right]$$

It tells you the rate at which the problem size must increase such that efficiency remains constant.

Block Matrix Multiplication

(Can use part to calculate cache hits).

Let cache hit time = 1 cycle

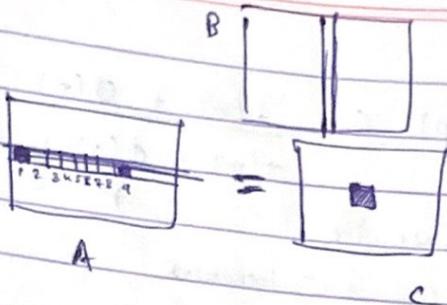
Miss penalty ≈ 100 cycles

Average Access time for 99% = $\frac{99}{100} \times 1 + \frac{1}{100} \times 100 = 1.99$ cycles

Avg access time for 97% = $\frac{97}{100} \times 1 + \frac{3}{100} \times 100 = 3.97$ cycles

\Rightarrow Takes almost twice the value.

\Rightarrow Avg access time for 80% = $\frac{80}{100} + \frac{20}{100} \times 100 = 20.8$ cycles



$L_1 = 32 \text{ Kb}$

$L_2 = 256 \text{ Kb}$

$L_3 = \cancel{64 \text{ Mb}} 6-8 \text{ Mb}$

Matrix elements = 8 bytes

Total operations = $O(N^2)$

Cache block size = 64 bytes

$$\text{Cache Miss} = \frac{N}{8} + N \xrightarrow{\text{For column,}} \frac{9N}{8}$$

(For row)

This is for 1 output \rightarrow Thus Cache Miss $O\left(N^2 \cdot \frac{9N}{8}\right)$

④

- Using Blocks - Block size = B^2
The cache can hold 3 blocks ($3B^2 < C$)
- Cache Miss for each block = $(B^2/8) \cdot \frac{2N}{B}$

$$\text{Total Cache Miss} = \frac{2N}{B} \times \frac{B^2}{8} = \frac{nB}{4}$$

$$\text{Total Blocks} = \left(\frac{N}{B}\right)^2$$

$$\therefore \text{Total Cache miss} = \left(\frac{n}{B}\right)^2 \times \left(\frac{nB}{4}\right) = \frac{n^3}{4B}$$

- General observation - Keep your working data set very small.
Keep stride sizes small. (jumps in memory)