

July 27,

# HIGH PERFORMANCE COMPUTING.

WEDNESDAY

- Optimize run-time but also consider hardware.
- Asymptotic analysis.
- Not faith-based (No established laws)  
There are thumb rules and guidelines but no absolute truths.
- There are parallel computing maps that are formed.
- Processors are called socket. Multi-core sockets.
- Number of transistors used decides performance  
HPC cluster  $\equiv$  Nodes  $\in$  Sockets.  $\in$  Core
- Based on the algorithm used, the code can be parallelized in different ways.
- Speed up vs. number of cores  
Speed up vs. problem size } Plots

July 29, FLOPS = freq.  $\times$  no of op.  $\times$  no. of cores.  
Clock cycle

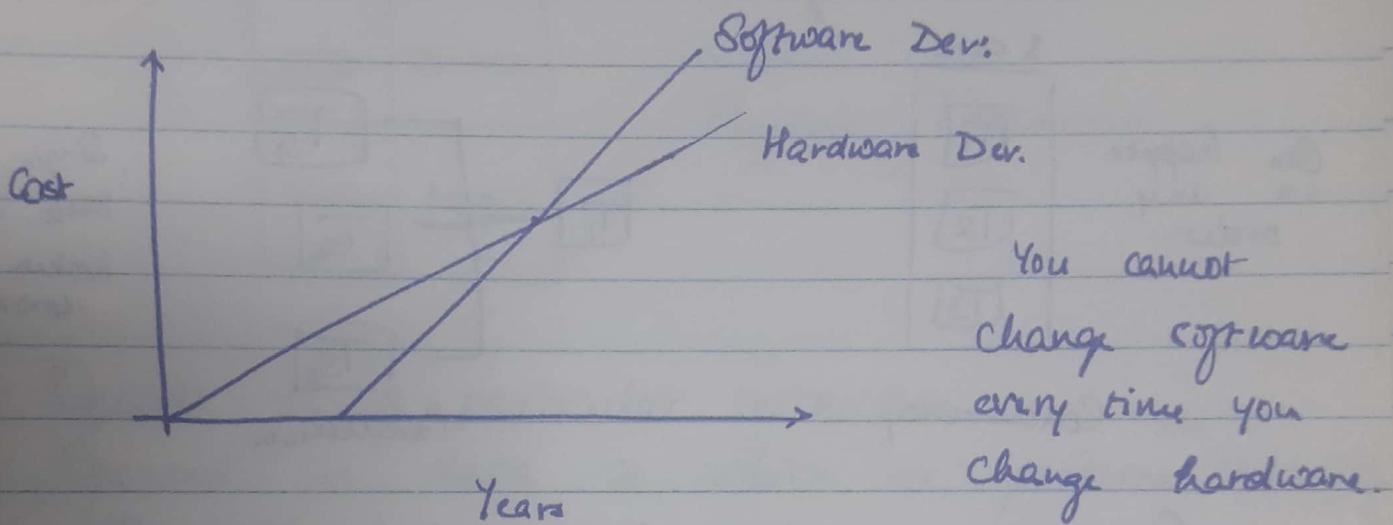
FRIDAY

This is how performance is measured.

- Initially, the number of cores was not too important. The frequency could also not be increased much beyond a point. (Due to heat)
- Then, there were two new design philosophies:
  - i) Latency minimization (unit: sec.)
  - ii) Throughput maximization
- Latency: Time required to complete a task  
All CPU's are based on this
- Multi-core system. (4 - 20 cores)
- Throughput: No. of operations per cycle.  
How much total work you can do in a given time. Does not matter how much time each operation takes as long as a lot of operations are performed.
  - Many-core (100 - 1000 cores)
- Practically, the design is a combination of both philosophies.
- For a many-core, you need a smarter/bigger control unit. Thus, in general, multi-core is used.
- GPU uses throughput max.

## → Factors affecting performance.

The code must be scalable and portable.



Portable: Can work on any hardware / platform

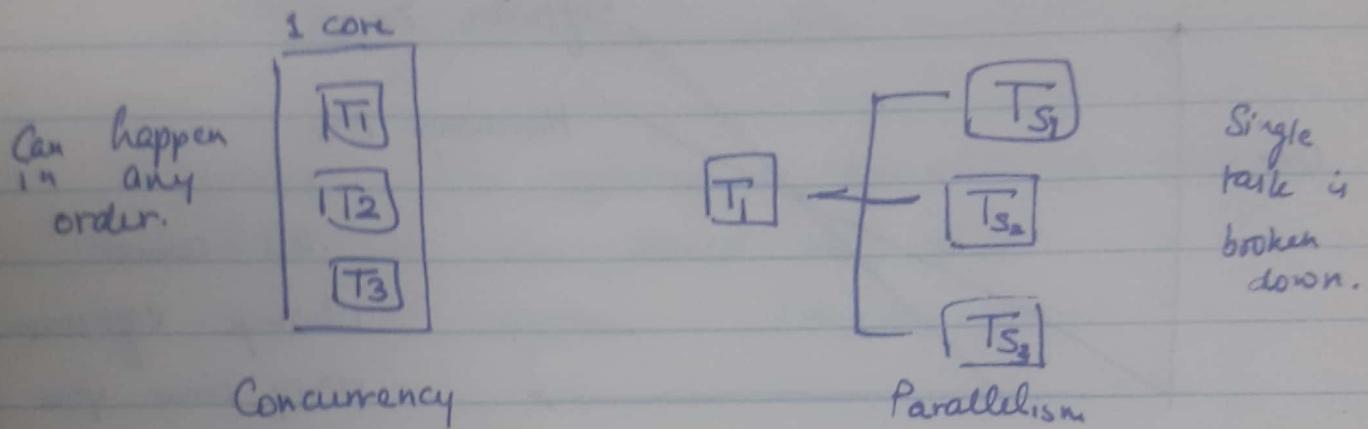
Scalable: Can work on a higher number  
of cores as well.

- Memory access time decreases by shared memory.
- Bandwidth tells us about throughput.  
bytes/sec.
- The constraints ~~are~~<sup>is</sup> computer architecture. So, focus on parallelism

## → Multi-tasking vs Parallelism

Time slicing and multi-scheduling leads to multi-tasking. Switching between

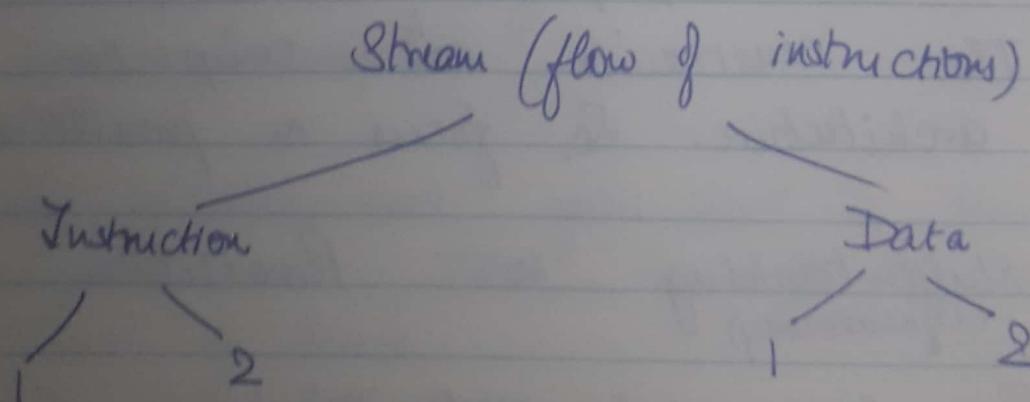
different processes very fast isn't parallelism  
You are only sharing resources. Both  
tasks run and neither is completed  
In parallelism, both tasks run  
simultaneously.



Concurrency control is used for communication

- For more performance, increase power

August 2, Tuesday. There is instruction and data level parallelism.



SIS - Single Instruction Single  
SIM - Single Instruction Multiple

These are serial instructions.

|      |      |
|------|------|
| SISD | SIMD |
| MISD | MIMD |

Flynn's Taxonomy.

- Opcode is a particular instruction.
- Operand is data on which we perform <sup>the</sup> instruction.  
Thus, there is always an opcode & operand.

- Pipelining: Some kind of assembly code
- Shared memory: Memory shared by diff. processes
- Distributed memory: Cannot access directly but through some network

Shared memory has very fast access.

- Granularity: Computation  
Communication

The ratio is important because they happen at different time scales. Computation is much higher. Thus, we want a ratio  $< 1$ .  
[Ideal. Not possible]

$> 1$

Coarse granularity

$< 1$

Fine granularity

- Observed Speedup: No unit. Ratio.

$$S = \frac{\text{Time on serial computer}}{\text{Time on parallel computer}}$$

Processors must be kept the same.

Single core vs. Multi core

Other hardware must be same to allow comparisons.

$$S = \frac{T_s}{T_{s/p}} = P \rightarrow \text{no } \uparrow \text{ processing}$$

Ideal. Doesn't usually happen.

If this relation is satisfied, linear speedup

There is also superlinear speedup ie higher than  $P$ .

If the problem is coarsely granular, you can increase efficiency by / to get (?) linear speedup

- Synchronization: Even after division, it might be required that some processor finish a task before some others

- Parallel Overhead:  
Launch threads, kill threads.  
Explicit synchronization statements
- Massively Parallel: Huge number of cores
- Embarrassingly Parallel: High Granularity.  
There are the best algorithms.
- Scalability

e.g. SIMD

$\text{Sum} = 0$       Loop:  $i \rightarrow n$       Add 1

Divide the array among different processes. Compute individually, store in private and then merge.

Efficient parallelization is not parallelization of the same algorithm but generating a new algorithm altogether.

- Why parallelism? Problems:
  - Communication among cores
  - Load balancing
  - Synchronization.
- Complexity does not necessarily remain the same. May/may not be better than serial

- Data parallelism is Task based.  
↳ set of instructions

→ Amdahl's Law:

Maximum speedup is governed by

$$\text{Speedup} = \frac{1}{1-P}$$

$P$  = fraction that can be parallelised

$S$  = serial part

$$S+P=1$$

e.g. If  $P=0.95$ , speedup = 20.

$P=0$  Speedup = 1.

Theoretical maximum  $\Rightarrow$  can be  $\infty$ .

$$P=1$$

Speedup depends on  $S$ .

$$\text{Speedup} = \frac{1}{\frac{P}{N} + S}$$

↓  
 no. of  
 processors

If  $N >> P$ ,  
 $\frac{P}{N} \approx 0$   
 $\therefore \text{Speedup} = 1/S$

This is just theoretical. Does not account for parallel overheads.

- Efficiency: Every time you increase the number of processes, your efficiency decreases.

Increasing the size of the problem decreases communication time (or no. of communication)

e.g. Size =  $n$       3 processes

Ideally, if Size =  $2n$ , we would use 6 processes and thus 6 communications.

But, using 3 processes limits communication to 3.

This increases granularity.

- We measure time of each part of the code and not the whole code. This is because entire code time might also include wait time. Thus, put some timing functions in the portions of code. This is called elapsed time. Resolution is the order of time that a function can entertain. Use appropriate functions (timer functions) with resolution for computing elapsed time.

August 3, Wednesday. Linux Commands: lscpu, /proc/cpuinfo  
 Speed up =  $T_m$  for hardware details

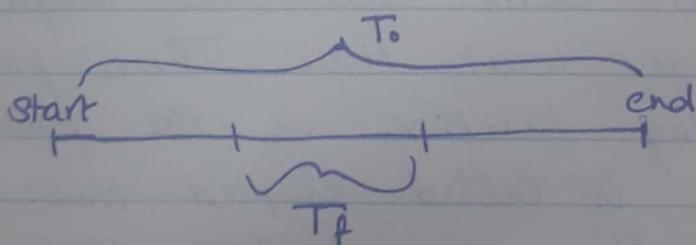
Cache memory for instruction differs from cache instruction for data memory

- a. Change problem size } Look at the Speedup
- b. Change core size      (Not necessarily linear)

→ Optimizing the code.

To understand relation between code and hardware (time it takes for execution of code)

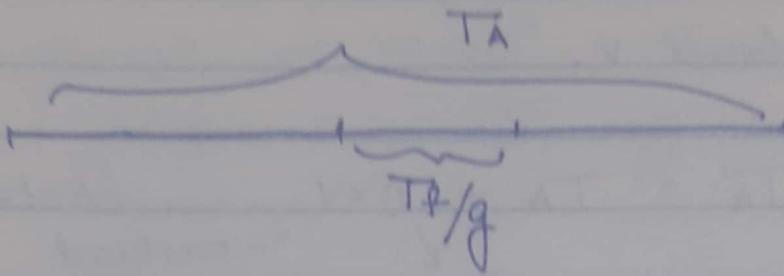
- Amdahl's Law.



T<sub>o</sub>: Time taken for serial code

T<sub>f</sub>: Part of time taken by code that can be serialized/parallelized

g: Peak performance gain for accelerated portion



$$S = T_0 / T_A$$

(Spud up)

$$f = T_F / T_0$$

f: Fraction which can be accelerated.

$$T_A = (1-f) T_0 + \left| \frac{f}{g} \right| T_0$$

$$\left( S = \frac{T_0}{T_A} = \frac{1}{(1-f) + \left| \frac{f}{g} \right|} \right)$$

$$T_A = T_0 - T_F + \frac{T_F}{g}$$

$$T_A = T_0 \left( 1 - \frac{T_F}{T_0} \right) + \left( \frac{T_F}{T_0} \right) \left( \frac{1}{g} \right) T_0$$

Theoretical

If we increase the problem size, the amount of overhead does not catch up with the increase in computation?

So, increasing problem size increases granularity

- Overhead v.

$$T_A' = T_A + n \times v \quad (\text{Actual})$$

↗ overhead  
 Due to  
 each of the section  
 where we did the  
 speed up.

- Measure overhead within the parallel computation
- Different Time functions have different precision (resolution)
- Use time function with much higher precision. (of us)

$$N = 10^7$$

| Processor cores | 1 | 2             | 3           | 4            |
|-----------------|---|---------------|-------------|--------------|
| Speed up        | 1 | $\approx 1.9$ | $\approx 2$ | $\times 2.5$ |

linear ~~scale~~ up

$$3 = \frac{3}{2}(1.9)$$

$$4 = 2(1.9) \approx 3.8$$

but we  
got 3.2

Efficiency ↓ decreases.

As  $N=10^7 \rightarrow N=10^9$  Relation Overhead decreases. So, efficiency relatively increases.

- Scalability

→ Strongly scalable  
→ Weakly scalable

If we increase the number of cores and increase the problem size in such a way that the efficiency remains constant. Then, the problem is known as Scalable.

$$T_{\text{serial}} = n \text{ sec.}$$

$$T_{\text{parallel}} = \left( \frac{n+1}{P} \right) \text{ sec.} \quad \text{linear scale-up}$$

$$\text{problem size} = n$$

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

$$\frac{n}{P} \gg 1$$

$$\text{Problem size} \gg P$$

$$\frac{n}{P} \gg 1$$

$$\text{So, } T_{\text{parallel}} = \frac{n}{P} + 1 \\ \approx \frac{n}{P}$$

$E = \frac{\text{Speed up}}{\text{Processors}}$

$$E = \frac{n}{\left(\frac{n}{P} + 1\right)} \times \frac{1}{P} = \frac{n}{n+P}$$

If we increase processors  $k$  times,  
then we will have to increase  
problem size. ' $x$ ' times, such that  
efficiency remains same

$$E = \frac{xn}{xn+kP}$$

$x = k$  for efficiency to remain same.

If we are not able to find an  
' $x$ ' such that efficiency remains  
same then it is 'weakly scalable'

- Profiling: Information about program's behaviour

Run-time  $\rightarrow$  "hot-spots"  $\rightarrow$  performance  
bottle neck  
are areas which take more time and we should focus on those.

Optimization.

• Function profiling:

If short functions, then we look at time taken by the function.

In operation → how we are accessing data. Essential to look at this.

e.g. gprof from GNU.

i. Flat function profile

ii. Call graph profile (butterfly graph)

gprof: Compile with ~~f95 -pg~~  
~~f95 -m pg~~ file.f

i. Flat profiling:

• la.out Then we get gmon.out

The information profiling is sorted by 3<sup>rd</sup> column.

The 1<sup>st</sup> column shows percentage contribution of each method.

August 5,  
Wednesday  
FRIDAY

Problem size =  $n$ ,

$T_{\text{Serial}} = n$

$$T_{\text{Parallel}} = \frac{n}{P} + 1$$

If we increase the processes, then problem size has to be increased.

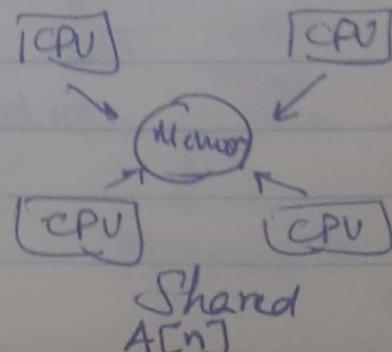
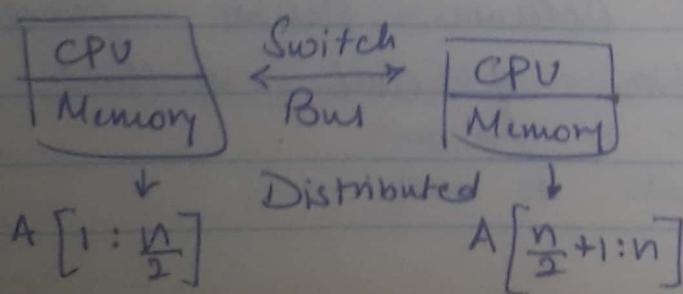
- **Strongly scalable:** If we increase number of processes without increasing problem size and efficiency remains same.
- **Weakly scalable:** If we increase number of processes and we have to increase the problem size so that efficiency remains same.

→ Shared Memory Parallelism:

Open MP - Multi processing.

# pragma omp { . . . }

This can be used when all cores have shared memory.



In order to access the other memory, we require to communicate with it using switch/bus  
(Explicitly pass messages)

We have to ensure that there is no collision (synchronous access) to the same memory.

• Fork: brings concurrency.

Creates child threads which run concurrently with parent

• Join: Removes concurrency.

Child joins with parent process. Parent waits for child to exit.

C:      # pragma omp parallel

{      // code to be executed by each thread.

}      (Each thread is launched on different cores)

• Specify Threads: Inside the code,

In code, setenv OMP\_NUM\_THREADS=4

In bash, export OMP\_NUM\_THREADS=4.

Compiler Directive

Control constructs  
(Launching threads)

Data Scoping  
(Race conditioning,  
Shared access)

Synchronization

August 9. • Structured block is a part of  
TUESDAY. code that has one entry  
point and one exit point.  
No branching.

These are easiest to parallelize.  
With branching via "if" statement,  
cannot determine exact flow.

- Take a serial code and make minimum changes. This should get maximum change in speed-up.
- Trying to discretize a continuous function leads to computational inaccuracy.

e.g. Integration.

```
h = (b-a)/n;  
approx = (f(a) + f(b))/20;  
for (i=1; i <= n-1; i++) {  
    xi = a + i*h;  
    approx += f(xi);  
}  
approx = h*approx;
```

[To compile: gcc -fopenmp hello.c]

To parallelize, we can assign some trapezoids to some cores, thus dividing the work. Each core individually computes a ~~the~~ private result and all these are added at the end, to get the final result.

global-result += my-result.

Mutual exclusion      Race Condition

Problems of concurrency will arise when multiple threads try to update it. Thus, add mutex (equivalent critical directive)  $\leftarrow$  #pragma omp critical.

- If you make my-result a function, we are actually serializing it because that function will run <sup>one</sup> at a time.
- If my-result is an array, when trying to access [0][1] simultaneously, it tries to access cache lines ~~a~~ which creates a problem (TBD)?
- \* Default scope of variables declared before a parallel block is shared.

→ Reduction:

Reduce to a single reduced variable.  
? Similar to calling piece of code critical.

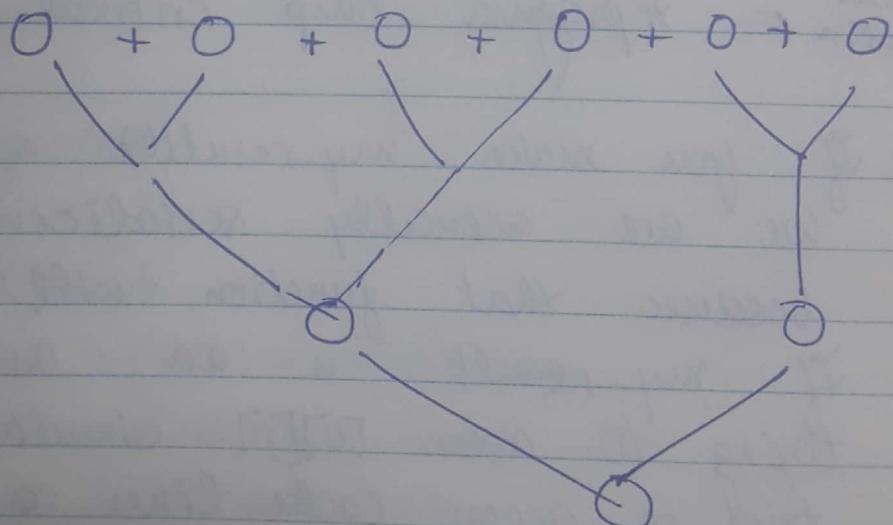
- Floating point in binary is not associative.

reduce(<operator>: <variable-list>)

- Exactly like critical but the compiler does it for you

Binary operators.

- Reduction algorithms.



Thus, still some threads can work in parallel while updating.

More efficient, requires less changes, than 'critical'.

• Dependency algorithms: Programs where the next step cannot be computed before the previous one.

e.g. Fibonacci.

You must remove all <sup>loop</sup> dependencies.

This is loop dependency.

Data dependencies may be fine.

#pragma omp for

This specifies sections where the code has to work. Improves efficiency.

August 10,

Wednesday.

Intel profiling ← Commercial profiling  
Use timers instead of gprof.

→ Optimum Techniques:

e.g. if do( )

if( $a > 2$ )

A) b exit

endif

do( )

B) if ( $a > 2$ )

endif

Use A, because ends <sup>break</sup> once we of loop is over.

eg.  $A = B^{**} 2.0$  → Exponent is actually calculated by  
 $A = B * B$   
 $\checkmark$   
 More primitive  
 so faster

$$x^y = e^{y \ln x}$$

Thus, more demanding.

eg. ~~for i=1,10~~

$$A=0$$

$$i=1, 180$$

$$A=A+h$$

$$\cos(A)$$

To compute cos() function faster,  
 store the values in a look up table. Thus,  
 we can have a

~~The value table for all values in the cache and access is direct.~~

• Cache exploits locality → {Temporal  
 Look-up tables are generally stored in  
 1<sup>st</sup> / 2<sup>nd</sup> level of cache.

Check: How many clock cycles does division require vs multiplication? 4?

Prediction

eg. Branching is costlier and takes more time for predictions.

eg. do ( )

do ( )

if ( )

else if ( ) } Branch

else

- Each processor has its own cache. Thus, when you optimize for one processor, optimize for all processors.  
This leads to superlinearity.  
(Speed-up becomes more than theoretical as theory does not account for multiple cache)
- In different languages we different defaults in terms of storage.  
C is row-major  
Fortran is column-major. <sup>spatial</sup>  
Thus, this would change the locality in memory.

$$\begin{bmatrix} M_{00} & M_{01} & M_{02} \\ M_{10} & M_{11} & M_{12} \\ M_{20} & & \\ M_{30} & & \end{bmatrix} \xrightarrow{\text{Row}} \boxed{M_{00} \quad M_{01} \quad M_{02}} \rightarrow \overbrace{\boxed{M_{10} \quad M_{11} \quad M_{12}}}^{\text{Column}}$$

### → Scheduling:

How many jobs are assigned to each thread? How this is broken affects.

#### • Types:

- i. Static: Simply equally divide. If not exactly divisible, give extras from the beginning threads.

This would be load imbalance.

Scheduling done at compile time  
May be possible that slower thread gets more work.

- We are not aware of how the hardware exactly works. Thus, some threads might work faster than others if the hardware of that core is better.

i) Static - 3. : Specify chunk sizes. ?

Extra: That much allocated to each

ii) Dynamic: Based on <sup>availabilities</sup> abilities of that core, jobs are assigned to the threads.  
*Not work at run-time*

iii) Guided: Looks at the work and divides and assigns. Then, remaining work is ~~assigned~~ divided again and assigned.  
Keep dividing.

→ Work Sharing:

$$x = 50$$

\* pragma omp for private(x)

{      i = 1, 20

$$x = i$$

x specifies that work on which is in stack of thread.

? print x.

$$\underline{\text{Ans} = 50}$$

`lastprivate()` → Stores the value  
of the last division.  
This doesn't mean that the last  
thread is supposed to be the last  
thread computed. But, value stored  
is the last one.  
`Ans = 20.`

August 16, TUESDAY. "Computational science is all about linear algebra"

Three combinations: Vector - Vector  
Vector - Matrix  
Matrix - Matrix

→ Vector Dot Product:

Multiplication of two vector followed  
by summation

Private variables are stored in stack.  
Shared are in heap.

This is the same as local and global.

# pragma omp parallel for default(shared)  
  private(i)

Is `private(i)` necessary?  
Find difference between global and shared  
variables.

Makes collapse(n) For nested loops.  
single loop of  $N \times M$  and parallelizes.

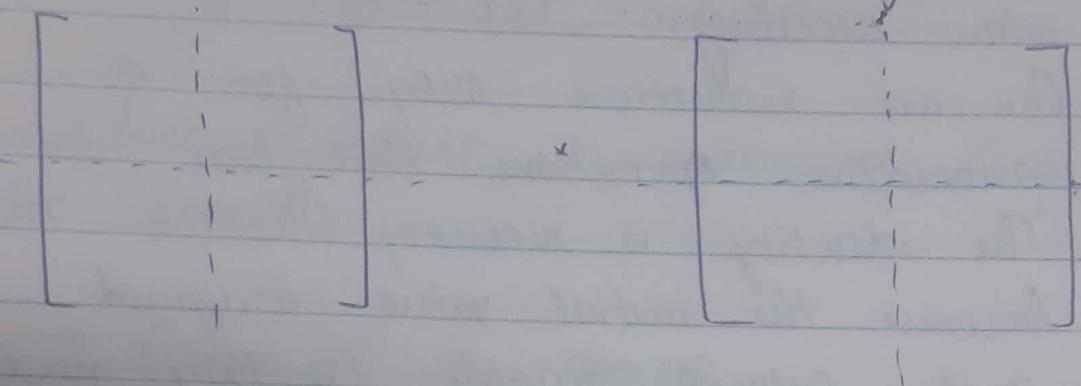
No. of columns<sup>of second</sup> must be equal to number  
of rows in first

Sequential code:

Each statement does 2 memory access  
and 2 operations. All of this  
affects the performance.  
 $O(n^3)$

- While parallelizing, the main issue is data access. This is what will adversely affect the speedup.  
Try and see the pattern of data access and if any part is reusable.  
Once read, see if the value can be reused!
- Instead of reading whole rows and columns, divide the data into blocks and then compute in partitions.  
This is called Block Matrix Multiplication.  
This reduces data access time.
- Cache is used to reduce data access lag.
- When an array goes into cache, locality is also very important. Thus, when you partition the matrix into blocks, the <sup>temporal</sup> locality of each block greatly increases, thus improving speedup.

Also have to think about which of the three loops you will parallelize. For this, look at how C or traverse arranged matrices in column row major.



- Coalescing also takes place. This changes depending on the language you use.
- If C is row major. So, if you have  $A[0][0]$ ,  $A[0][1]$ ,  $A[1][0]$ ,  $A[0][0]$  and  $A[0][1]$  will coalesce.

Coalescing means combining adjacent blocks.

- Thus, how what blocks are made affects the time taken for computation.

$$\text{Eqn. 2. } \pi = 4 \sum_{k=0}^{\infty} \frac{(-1)^k}{2k+1}$$

$\text{Sum} = 0 \rightarrow \text{Shared.}$

This needs to be mutually exclusive.  
Three ways to do this:

Critical

Reduction

Atomic (?)

reduction (+; sum)

When using reduction, sum becomes a private variable. It is private when being computed but global when accessing.

Can use reduction only for associative operations that have identity. The identity is necessary because that becomes the initial value assigned to the reduced variable. Do they need to be commutative? No.

Sum=0;

reduction (+; sum);

for i=0, N

{

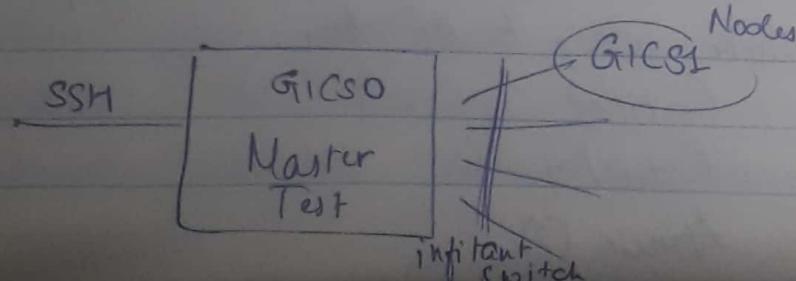
$$\text{Sum} += \frac{(-1)^k}{2k+1}$$

}

How do you use the k? It is dependent on previous but won't be valid during threads.

Can use  $i \% 2$ .

August 19, ssh <ID>@10.100.71.130.  
FRIDAY.



Classify problems into various categories  
and approach is different for each.

(i) Divide and Conquer

e.g. Trapezoidal integration

- Most important to understand Cache and locality.

→ Components of OpenMP

i. Directives → work-sharing functionality.

ii. Runtime libraries → no. of threads, thread ID, nested parallelism, timer.

iii. Environment Variables → setting no. of threads, scheduling

OMP\_DYNAMIC → Breaks the code into chunk sizes.

Scheduling → Chunk size  
→ Static/dynamic

There is always an optimum chunk size.

• Implicit barrier while exiting a parallel section i.e. cannot exit unless all threads have been computed.

If do not, no-wait()

\* pragma omp master

Only the master thread is allowed to access the share variable.

### → Embarrassingly Parallel Computations:

Can be divided into completely independent parts. The threads do not need to communicate

e.g. Matrix addition

Random number generator.

Image processing (Geometric transformation of each pixel)

Most important is how you divide.

At times, may get worse results than serial if not divided correctly.

Like a 2D-array. 2 loops required to traverse. Each pixel has 3 data points (R, G, B). This is converted to a black and white picture where each pixel only carries the intensity.  
Data access has to be considered.

$$x' = (x - c_x) \cos\theta + (y - c_y) \sin\theta + c_x$$

$$y' = -(x - c_x) \sin\theta + (y - c_y) \cos\theta + c_y$$

Colour scaling is more efficient than transformation. Why?

? The location of the write pixel changes and so data access takes more time.

When we run the program, data is shifted from main memory to RAM.

There Two types of registers:

i) General purpose

Data storage. Operands

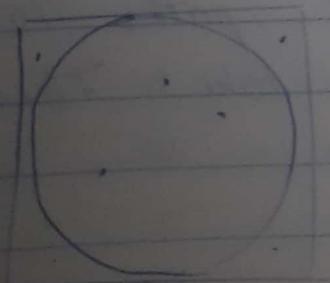
Can it have instructions?

ii) Special purpose

PC, stack pointer, instruction register  
current instruction

What happens during la.out? What goes into RAM?  
exec file?

e.g. Monte Carlo method for  $\pi$



Generating coordinates randomly.

Finally, look at probability of the point lying inside the circle.

- Uniformly distributed: Equally dispersed over the number-line. No bias.
- Statistically independent: Future value is not dependent on previous values. This makes it much easier to parallelize.
- Can use histograms to check uniform distribution. Can also check average of all. How to use random number generator for parallel programs? Does it matter?  
Go to same random?

August 23, Memory access is very important. Have to look at latency (time taken for one access)

- If you request 0 bytes, bandwidth (byte/sec) is 0. However, latency is not zero. Some fixed minimum time will always be used.
- Memory hierarchy.

The actual distance from the CPU is looked at. L1 closer than L2. This is what determines the latency.

- Main memory: DRAM (Dynamic RAM)
- There are many different paths to reach that the program can take. This decides performance.

- General purpose      Special purpose registers.
    - ~ A program can access only them.
  - Two important instructions : STORE, LOAD.
  - To probe the memory hierarchy and see when saturation happens (proportional to problem size)  
(Peak performance which in turn is measured through speedup)

→ Microbenchmarking:

Difference between processes: instruction sets  
Basic form of parallelism: System level parallelism.

(iv) Instruction level parallelism (ILP): pipeline

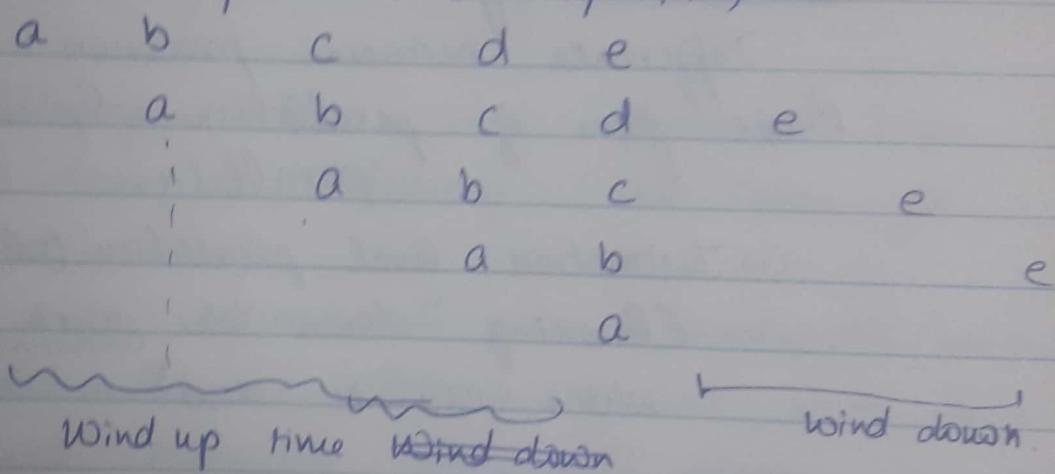
## Breaking down the work in steps

which can be done in Constant time.

- X
- Upper triangle of pipeline matrix is called Wind down. This is the part where instructions are leaving. Overhead.
  - First half is wind up, since no new instructions enter. This is also Overhead. ?

*problem size* • Longer pipelines imply greater overhead.

•  $n=1000$  loops better than  $n=100$   
(nested loop  $\uparrow$  wind up  $\uparrow$ )



Wind up is the time taken for an instruction to be executed.

Wind down is the time the last instruction takes to be executed.

$$\text{Overhead} = \text{Wind up} + \text{Wind down}$$

- After wind up, we will obtain one instruction per clock cycle. All stages of the pipeline get filled.

Wind-up and Wind-down take much more time if the pipeline has more stages. During parallelization, they can be considered as overhead.

$N +, -, *$  → Diff. operators. Look at time of execution.

$$T_{\text{exec}} = E - S$$

$\frac{T_{\text{exec}}}{N}$  = Performance of diff. operators

August 24, Hardware:

WEDNESDAY

Size of cache memory, bandwidth for cache access, bus speed, processor architecture...

- Difference between low-end and high-end (GISC) machines.
- Cache size is larger for high-end machines (Xeon) making parallel code performance much higher.
- Lifetime of a thread stack in the lifetime of the thread.  
(variables are lost after)

- Stacks are faster than heaps. But stacks have limited size (They can also be extended?)

$i = 5, N$

$$C[i] = A[i] + B[i]$$

|             |   |   |   |   |   |                          |
|-------------|---|---|---|---|---|--------------------------|
| 5           | 1 | 2 | 3 | 4 | 5 | *                        |
| 2           |   | 1 | 2 | 3 | 4 |                          |
| 3           |   |   | 1 | 2 | 3 | Max. Throughput          |
| 4           |   |   |   | 1 | 2 | = 1 instruction / cycle. |
| 5           |   |   |   |   | 1 |                          |
| Latency     |   |   |   |   |   |                          |
| Wind Up = 4 |   |   |   |   |   |                          |

- We can get maximum throughput after wind-up
- For  $n$ -depth pipeline, we obtain max throughput after  $(n-1)$  cycles.
- Pipeline for different functions is different.  
(? What happens to wind-up & down)

e.g.  $n = 5$

Depth =  $n$

Instruction =  $N$

JLP, cache, superscalar architecture = Throughput more than 1

$$\text{Speed up} = \frac{nN}{N+(n-1)} \rightarrow \begin{array}{l} N \text{ instructions, each takes } \\ n \text{ cycles,} \\ \text{Serial so no pipeline} \\ nN \end{array}$$

Pipeline is the parallelization here JLP.  $T_{\text{Parallel}} = N(n-1)$

Wind-up cycles (overhead)

For complicated instructions, pipelines are much longer so latency is much higher.

$$\text{Speed up} = \frac{n}{1 + \frac{(n-1)}{N}} \quad (n-1) \ll N$$

In this scenario,  $(n-1) \ll N$  speedup =  $n$   
(speed up  $\propto$  Depth of pipeline)

$$\left[ \frac{n-1}{N} \rightarrow 0 \right]$$

$$\text{Speed-up} = \frac{nN}{N+(n-1)} = \frac{\cancel{(n-1)}}{\cancel{(n-1)}} + \underbrace{\left( \frac{1+N}{n-1} \right)}_{\frac{N}{(n-1)} + 1}$$

$$[S \rightarrow 1.]$$

For example of 'N' addition,

$$\begin{aligned} \text{Throughput} &= \text{No. of operations/cycle} \\ &= \frac{N}{N+(n-1)} \xrightarrow{\substack{\text{operations} \\ \text{pipeline depth}}} \end{aligned}$$

For parallel code,

$$\text{No. of cycles} = N+(n-1)$$

$$\text{No. of operations} = N \text{ additions.}$$

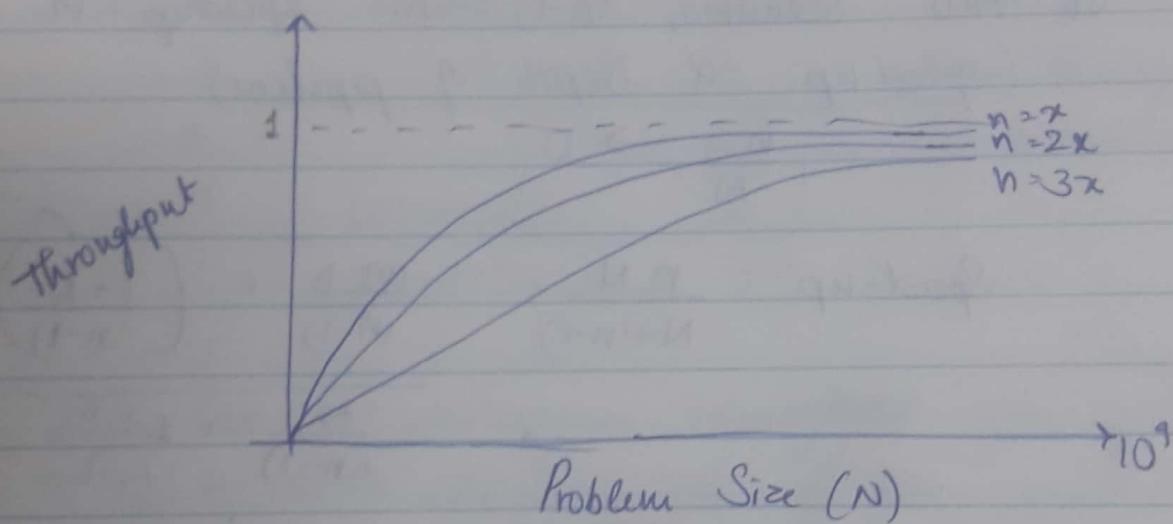
$$\text{Speed-up} = n \times \text{Throughput.}$$

$$\boxed{\text{Speed-up} = \text{Depth} \times \text{Throughput}}$$

$$\text{Throughput} = \frac{N}{N+(n-1)} = \frac{1}{1 + \frac{(n-1)}{N}} \rightarrow \text{no. operations}$$

Addition  $\rightarrow$  Multiplication  $\rightarrow$  Division.

• Depth(n) of the pipeline changes.



Smaller depth will saturate to maximum throughput faster than operations with larger depth. (As has lesser overhead)  
So, we should try and use simplified operations with smaller depth.

$\rightarrow$  Loop Dependency.

$$\text{eg. } A[i] = A[i-n]$$

$$A[1] = A[0] \\ A[2] = A[1]$$

$\text{MOV } A[i-1], A[i-2]$      $\text{MOV } A[i], A[i+1]$   
 $\text{MOV } A[i], A[i-1]$      $\text{MOV } A[i+1], A[i+2]$

Because second instruction, problem in pipeline.  
 Have to wait for first inst to be complete before executing

Since first, doesn't matter.  
 Don't wait, can directly execute without problem.  
 $A[i+1]$  is updated only after.

T/

$n = -1$

$n = +1$

→ Real dependency

Pseudo dependency.

$A[i] = A[i-1]$  Until  $A[i-1]$  is not completed,  $A[i]$  won't be either.

$A[i] = A[i+1]$  → Faster. Don't need to flush.

$A[i+1]$  is already present which  
 can be used to directly update  $A[i]$   
 Value is available so update is done  
 without waiting.

Mem Read Fetch Decode Execution Write

$A[i] = A[i-1]$

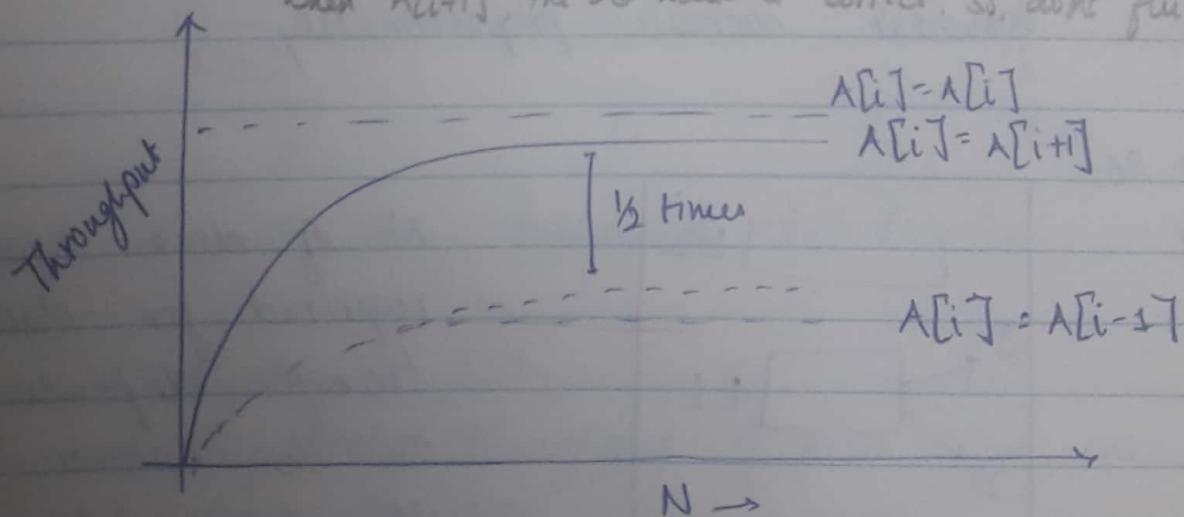
$A[i-1] = A[i-2]$

Talking about pipeline flush.

When  $A[i] = A[i-1]$ , it is possible that old value of  $A[i-1]$  is stored in  $A[i]$ . execution stage.

This needs to flush pipeline to get updated value. Then execute So, more time.

When  $A[i+1]$ , the old value is correct. So, don't flush.



CISC is always complex as compared to RISC  
 But for larger throughput, we use  
 CISC? Shorter depth of pipeline required  
 So, wind-up and flush quidar

Super Speedup  $\rightarrow$  more than one instruction architecture per cycle.

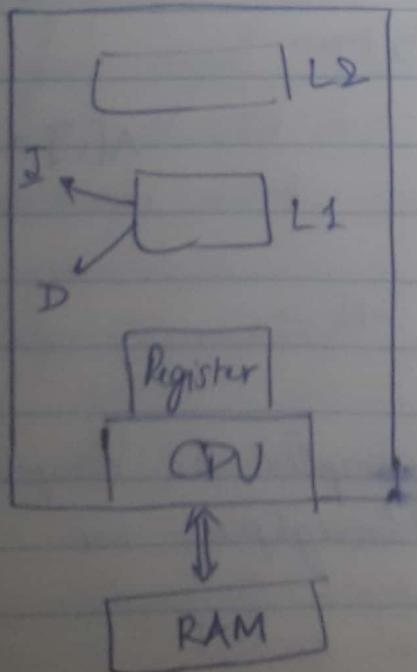
Not able to give data fast enough but processor speed keeps increasing DRAM GAP.

- Von-Neumann Bottleneck

Limitation on throughput caused by the standard personal computer architecture caused by the standard personal computer architecture, like slower memory access.

Faster a processor is, more time it has to spend idle as the memory access time wouldn't decrease proportionately.

August 26,  
FRIDAY.

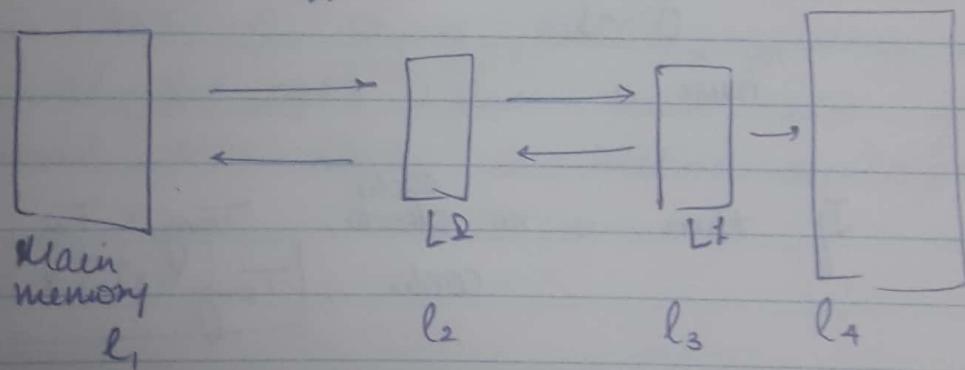


No. of clock cycles required for data transfer depends on the distance of the memory from the CPU.

Suppose bus has a clock rate of 600MHz (BW)  
Width = 32 bits ?

- Latency : Time for a single operation  
Memory latency as well.  
Time delay to transfer 0 bytes of data  
(which depends on different devices  
(memory units))

• Delay before actual data is put on link.  
Also, different distances of RAM from CPU  
result into different latency cache.



latency:  $l_1 > l_2 > l_3 > l_4$

1. Pipeline : ILP
2. Superscalar: Throughput of more than one
3. RISC to CISC
4. Out of order execution - GPU (Compiler takes care)
5. SIMD

- Performance of CPU is increasing at a much higher rate than memory access rate
  - Processor-Memory Gap
- Temporal locality - Reuse of same cache data over time
  - e.g. repeated access  $a[i]$
  - Temporal → spatial locality.

$$\begin{array}{l} \text{Memory} \rightarrow T_m \\ \text{Cache} \rightarrow T_c = \frac{T_m}{r} \end{array}$$

If reuse ratio =  $\frac{1}{2}$   
 Cache  $\approx T_m$   
 Cache takes double the time than it would require.

Reuse ratio  $\rightarrow r$

No reuse  $\leftarrow 0 \rightarrow 1$   
 Complete reuse (Theoretical)

If there is no cache,  $T_{avg} = T_m$   
 cache,  $T_{avg} = rT_c + (1-r)T_m$

- Cache line: Whenever we are accessing from main memory, the cache also copies the adjacent locations to the value we are searching. Typically, 64-128 bytes.  
 For matrices, rows in C.

e.g. Data = 128 bytes

Bandwidth = 50 Gb/s

Latency = 50 ns

$$= \frac{128}{50} \times 10^{-9} \approx 2.4 \text{ ns}$$

Here, latency > Latency due to bandwidth.

Thus, it becomes a more critical issue.

Cache length = 16 bytes

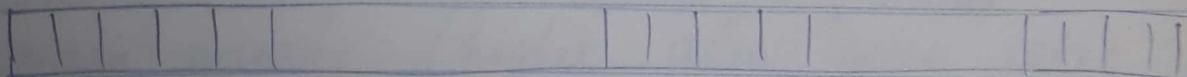
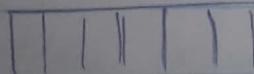
Hit ratio = 0.5

→ We can find the data only half the time.

So, actual utility =  $16 \times \frac{1}{2} = 8 \text{ bytes}$

(50% of spatial locality is being used)

- Look at hit ratio and reuse ratio for different size blocks in matrix multiplication



Prefetch: To predict which data/instruction would be more useful.

Done by compiler. It will bring that even before requiring data/instruction.

Cache bandwidth, cache speed, cache

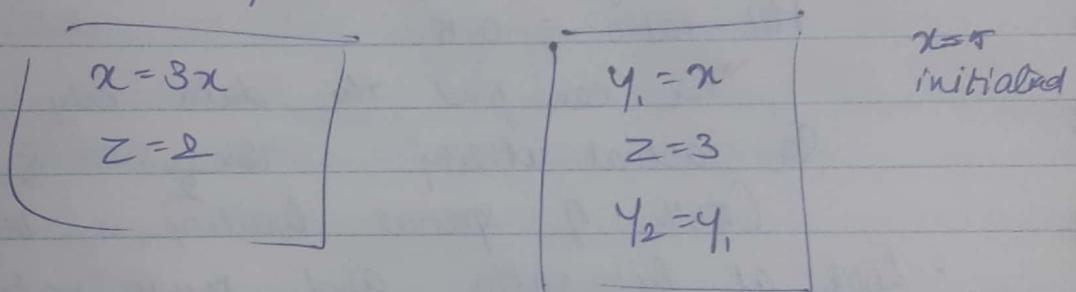
throughput

- Why use smaller?

Even if locality is not required in the program, we might end up blocking the whole bus.

- Cache Coherency:

Changes in each cache should be reflected in memory.



In order to avoid race conditions,  
mutex lock.

barrier, atomic no wait

master single collapse

sections first private thread private

August 30, Data Access: Change the datapaths  
TUESDAY based on which we can change performance.

→ Performance Modelling:

i) Latency Delay

process, as well  
to memory

ii) Bandwidth of path

• How do we reduce latency?

i) Instruction level parallelism

ii) CISC → RISC (smaller instructions so quicker)  
(No. of operations reduced)

• Main memory: Bandwidth =  $10^{10}$  bytes/sec.

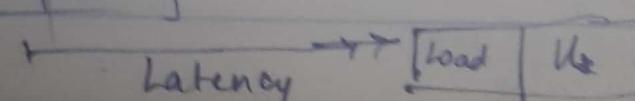
• Bandwidth of cache is higher than main memory (higher data access rate)

→ Overlapping Latency:

Clock cycle: 1 | 2 | 3 | 4 | 5

| Load | Use  |      |      |     |  |
|------|------|------|------|-----|--|
| LD   | Use  |      |      |     |  |
|      | Load | Use  |      |     |  |
|      |      | Load | Use  |     |  |
|      |      |      | Load | Use |  |

Suppose cache has 4 contiguous memory locations



- If there is a cache miss, the data will be obtained after a latency (due to access from memory)
- In order to reduce latency, compiler prefetches the instruction. (Compiler predicts the possibility of usage of the cache miss element and thus starts prefetching and so in turn reduces latency)

→ Rule of 72.

Growth rate =  $x\%$ .

$\frac{72}{x}$  years will be required to double performance

DRAM :  $\frac{72}{9} = 8$  years       $9\%/\text{year}$

CPU processor :  $\frac{72}{60} \approx 1.2$  years       $60\%/\text{year}$ .

Thus, Processor-Memory Gap.

→ Parallel Computing: Limitations

a) Von Bottleneck.

Hardware → Constraints in Software

(i) Moore's law:

More transistors → More Speed  
↓  
More power  
More heat → Less Speed  
(CPU starts giving less performance)

(ii) Microbenchmarking:

- a. Data Access
- b. Processor's capability
- c. Understanding hardware
- d. Programming capability

→ Balance Analysis:

- i. Machine Balance
- ii. Code Balance

iii. Machine Balance: It is intrinsic.

= Memory Bandwidth (B/s)  
Peak performance (FLOP<sub>s</sub>/s)

=  $\frac{\text{Bytes}}{\text{No. of operations}}$

FLOPs =  $\frac{\text{No. of operations}}{\text{cycles}}$   
 $\times \frac{\text{No. of cycles}}{\text{second}}$

- Depends on which memory you access  
(Bandwidth is different)

$$\text{Balance of PDM} = \frac{10^{10} \text{ B/s}}{3.3 \text{ GHz} \times 4 \times 3} \xrightarrow{\substack{\text{No. of cycles} \\ \text{second}}} \text{No. of cores}$$

$\xrightarrow{\substack{\text{No. of operations} \\ \text{cycle}}}$

Machine balance is going to decrease in the future as DRAM gap increases.

(ii) Code Balance e.g.  $A + B \times C$

$$\text{Data Traffic} = 3$$

$$\text{No. of operations} = 2$$

(iii) Code Balance:

$$= \frac{\text{Data Traffic (B/s)}}{\text{Floating point operations}}$$

$$\frac{1}{\text{Code Balance}} = \frac{\text{Computational Intensity of the code.}}{}$$

We always want to do more floating point operations with the data. Hence, for better performance, code balance should decrease.

- Machine and code balance have the same units
  - Data traffic: Performance limiting data path  
(hardware dependent)
- Data  $\leftarrow$  Load/Store

\* Peak performance of a code is the expected maximum fraction of peak performance of a code with Balance Code  $B_c$  and Machine Balance  $B_m$

$$\text{Light speed} = \min \left( 1, \frac{B_m}{B_c} \right)$$

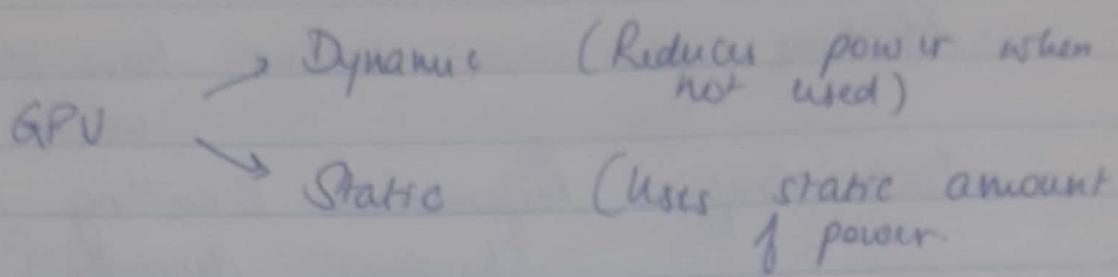
If Expected maximum fraction of peak performance of code = 1.

$B_m$  is constant.

If  $B_c \downarrow$ ,  $\frac{B_m}{B_c} \uparrow$

If the DRAM-gap increases, light-speed decreases.

- |                              |   |  |
|------------------------------|---|--|
| i, Data access               | } | Our codes generally underutilize hardware. $\leftarrow$ serial code. |
| ii, Processor                |   |  |
| iii, Performance of Algo.    |   |  |
| iv, Pipeline                 |   |  |
| v, Peak until latency hiding |   | Parallel tries to overcome this.                                     |



September 6, 2018 Parallel Vector Triad benchmark

TUESDAY

```

do j=1, NITER
  do i=1, N
    A(i) = B(i) + C(i) + D(i);
  enddo
  if i>C then
    C
  enddo

```

Code Balance =  $\frac{4}{2}$  → No. of operations  
 (Assignment is not an operation because it is a store operation under data access)  
 $B_C = 2$ . We want  $B_C$  to be as low as possible.

[Loop fusion, jamming, unloading]

```

eg. for i=1, N {
  A[i] = B[i] + C[i]
}

```

for  $i=1, N \{$

$$D[i] = B[i] + E[i]$$

}

Increase speedup by loop fusion

Lightspeed  $\ell = \left(1, \frac{B_m}{B_c}\right)$

$\frac{B_m}{B_c}$  shows what percentage of  $B_m$  we are able to utilize.  
peak performance.

Your machine can be more efficient, but you are only able to use  $B_c$  portion of it.

[Look at this during the assignment/project]  
Only way to improve after a point is by improving data access.

e.g. A present  $[N]$

A apart  $[N]$

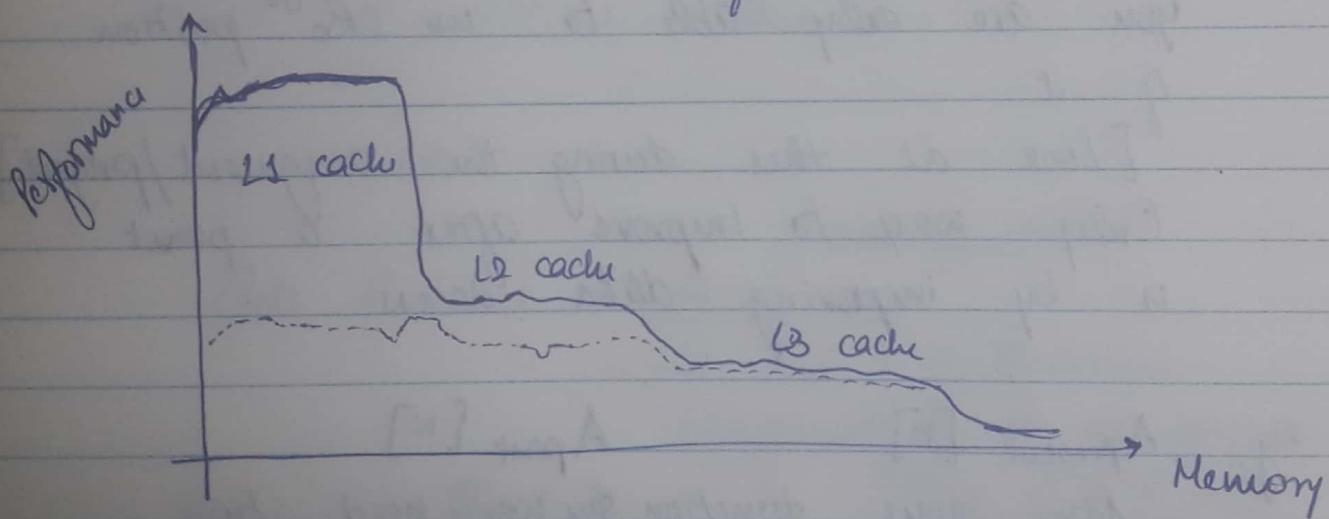
After some function, say we need two array information to be stored.

Rather than making them two separate arrays, it is much better to define a single two-dimensional array.

This has a greater chance of cache hits.

$A_{pp} [N, 2]$

- In the parallel vector benchmark, the outer loop has been added just to increase the code size so that looking at performance is easier. The bottom if is necessary because without that, the compiler is smart enough to ignore the outer loop.
- Vector processing: The processor takes chunks of data instead of single points. This makes a huge difference



There is a small dip in the initial part. This is due to latency and bandwidth. Bandwidth becomes a problem only for larger problems. However, the latency exists even for smaller problems. Always have to wait.

- There is a step drop between the cache proceedings. This is because L2 has to interact with both RAM and L1 and so on for L3. (? also slower in general)
- The initial drop is due to the wind up time in pipelines. The performance will increase only once the size of the problem is greater than the depth of the pipeline. Prefetching, etc are problems.

Ideally, we would have thought that the graph should have started off at a higher performance and then decrease <sup>saturate</sup>. because the whole program would fit in the cache and so less misses.

However pipeline plays a more important role.

- \* Look at performance.  
 $B_m/B_c$ , latency, etc.

Some algorithms don't have a scope of optimizing beyond a point. So, you may have to change the algo as well.

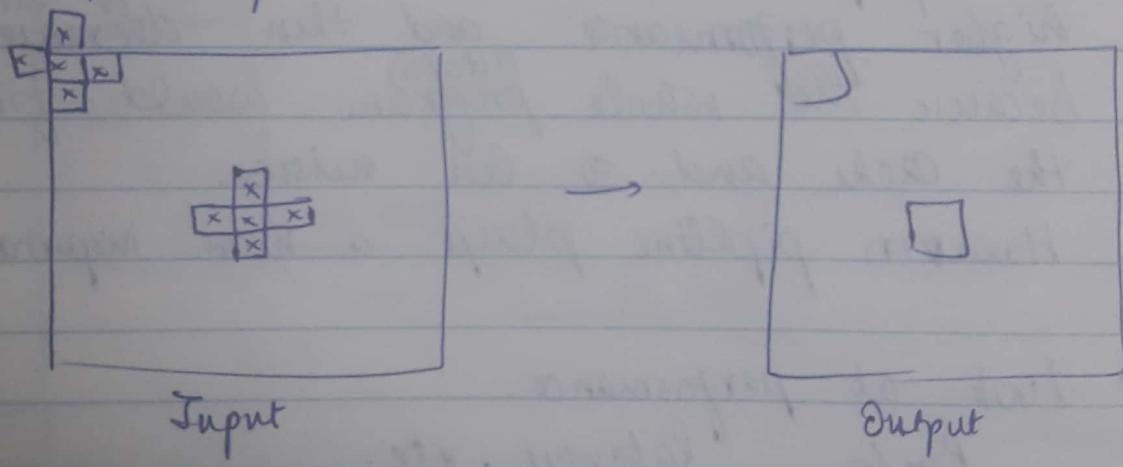
e.g. Image processing

$$\boxed{1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7 \ 8} \text{ Input}[1D]$$

$$\boxed{2 \ 3 \ 4 \ 3 \ 2} \text{ Mask}$$

$$\boxed{1 \ 21 \ 5970 \ 841 \ 1} \text{ Output}[1D]$$

Masks and other things have odd sizes in general. This is because it has symmetry



The mask acts in the form of a convolution & Multiply each of the mask to each of the input and add. (Five terms at a time)

If you have to exceed the input, padding based on how long the mask is.  
(ghost elements. Generally assumed to be 0)

$$\text{Mask length} = (2N+1)^2$$

- The filter/mask is constant. It won't be changed every time. It is already loaded into cache. Don't have to re-fetch. Thus, when you increase the size of the filter, you are increasing the number of operations per cycle without increasing data access.

Increasing the image size changes data access time.

- Arrays are generally not stored in registers. They are kept in cache.
  - $A[i] + \text{temp}$  ~ in register. Faster
  - $A[i] + B[i]$  No one in register/cache
- 'Can be considered to be one load/store as access time for temp is very less in comparison (Low latency)
- Latency is an more important factor as compared to bandwidth.

3D convolution - used in MRI scans, etc.

# September 7, Performance Modeling.

WEDNESDAY.

- 1, Amdahl's Law → possibility of parallelization
- 2, Gustafson's Barsis → evaluation of problem
- 3, Karp flat Metric → Investigate
- 4, Isoefficiency Analysis → Scalability

Restriction on performance :

- Parallel Overhead      • Serial part

Part

$$\text{Parallel Time} = \phi(n, p)$$

$$\text{Serial Time}_{\text{part}} = \sigma(n)$$

$$\text{Problem size} = n$$

$$\text{Processors} = p$$

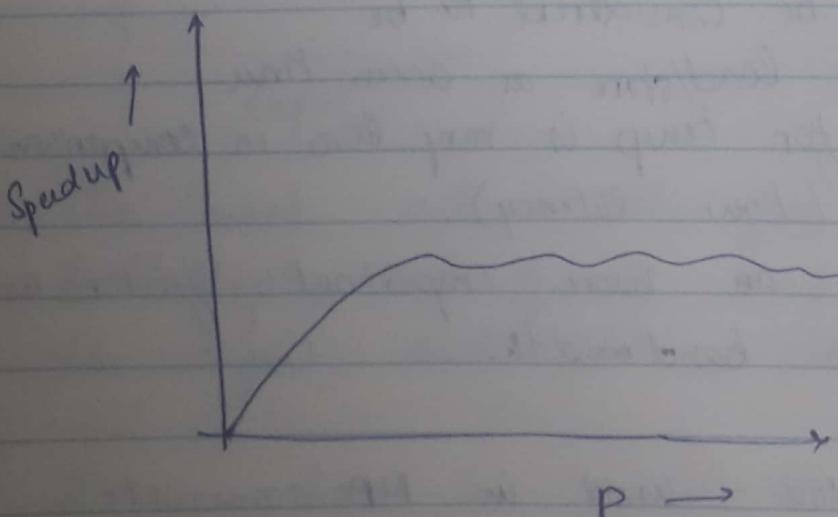
$$\text{Overhead} = k(n, p)$$

Speedup

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

$$P$$

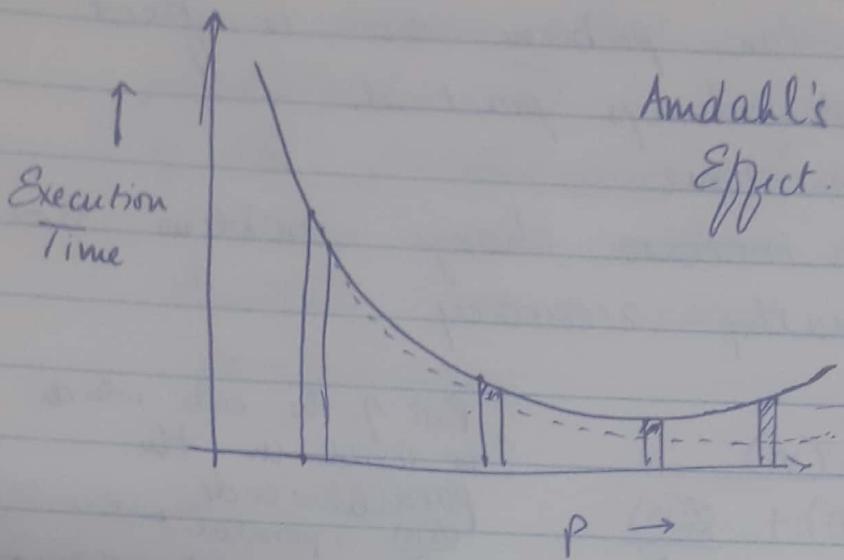
$$\text{Parallel Time} = \phi(n, p)$$



In terms of  
Amdahl's  
law, speedup  
saturates.

as it  
does not  
take

communication into account.



Amdahl's Effect.

This is the time required for communication. The other part is actual processing time which eventually saturates.

- Thus, there is always an optimum number of processes for any given problem

$$\text{Efficiency} = \frac{\text{Speedup}}{\text{No. of processes}}$$

$$\leq \frac{\sigma(n) + \varphi(n)}{P(\sigma(n) + \frac{\varphi(n)}{P} + K(n, P))} \leq 1$$

The denominator has a term of  $P \cdot \sigma(n)$ .  $\sigma(n)$  proportionately decreases as we increase problem size. Thus, proportion of serial part decreases.

$$\begin{aligned} \text{Speedup} &\leq \frac{\sigma(n) + \varphi(n)}{\sigma(n) + \frac{\varphi(n)}{P} + K(n, P)} \\ &\leq \frac{\sigma(n) + \varphi(n)}{\sigma(n) + \frac{\varphi(n)}{P}} \end{aligned}$$

• For Amdahl's law, problem size is fixed

This is not always practical.

Instead, fix time.

This would increase change problem size and eventually accuracy.

$$S_1 = \frac{\sigma(n)}{\sigma(n) + \frac{\phi(n)}{P}}$$

Part of the code which  
is serial in the  
parallelized code.  
 $\frac{\phi(n)}{P}$  is parallel part  
 $\sigma(n) + \frac{\phi(n)}{P} = 1$

$$1 - S_1 = \frac{\phi(n)/P}{\sigma(n) + \frac{\phi(n)}{P}}$$

$$\text{Speedup} \sim \Psi(n, p) \leq \frac{\sigma(n) + \frac{\phi(n)}{P}}{\sigma(n) + \frac{\phi(n)}{P}} = S + (1-S)\frac{p}{P}$$

$$= p + S + P - Sp$$

$$= p + (1-p)S$$

$p$  = no. of processors  
 $S$  = serial fraction

$$\boxed{\Psi(n, p) = p + (1-p)S}$$

Gustafson's Law.

• Only the initial assumptions are different.

Time taken for execution here is fixed

i.e.  $\sigma + \frac{\phi}{P}$ . Problem size can vary

This assumes that the code can be parallelized. It is used for evaluation of the parallelization, how accurately it was done within a certain time. Amdahl talked about possibility of parallelization.

September 13,

- TUESDAY.
  - Rotation of images.

Basically, taking one vector, & applying rotational matrix and getting new vector.

- Scaling of images

$$A\vec{x} = \lambda \vec{x}$$

Equivalent to computing eigen values.

- Rotation of images is similar to convolution masking which can be used in machine learning for feature extraction.

- Random variables must be

i. Completely unrelated independent

ii. Equally probable.

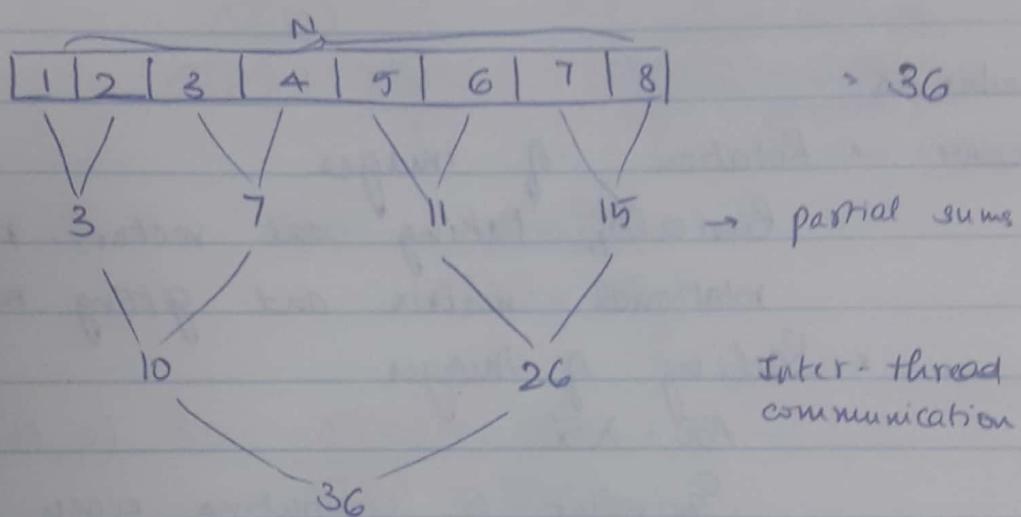
PJ computing for using random numbers uses probabilistic methods to get a deterministic answer.

Look at how rand() is applied in C

- Cycle of Random number generator. The

numbers will begin repeating after a point.  
 So, cannot give the same seed to each thread as they would begin generating the same random numbers.

- Reduction addition



There is some concept known as step complexity. Even if you have infinite processor, cannot compute in a single step.

$$\text{Here, Steps} = \log_2 N \quad N = \text{size of array.}$$

$$\begin{aligned}
 \text{Operations at each step} &= \frac{N}{2} + \frac{N}{4} + \dots + 1 \\
 &= \left( \frac{N-1}{2} \right) \cdot 1 + 2 + \dots + 2^{\lfloor \log_2 N \rfloor - 1} \\
 &= \frac{1}{2} \left( 1 + 2^{\lfloor \log_2 N \rfloor - 1} \right) - N - 1
 \end{aligned}$$

### Amdahl's Law

$$f = \frac{\sigma(n)}{\sigma(n) + \theta(n)}$$

- Serial in serial code
- Problem size is fixed
- Scaled speed up

$$\Psi \leq \frac{1}{P + \frac{(1-P)}{P}}$$

Does not consider overhead

### Gustafson's Law

$$S = \frac{\sigma(n)}{\sigma(n) + \frac{\theta(n)}{P}}$$

- Serial in parallel code
- Execution time is fixed
- Scaled speed up

$$\Phi \leq P + (1-P)S$$

Does not consider overhead

### Karp Flatt Metric

|        |   |   |   |   |
|--------|---|---|---|---|
| P      | 1 | 2 | 3 | 4 |
| $\Psi$ |   |   |   |   |

Include in report.

$$\text{Efficiency } e = \frac{(p-1)\sigma(n)}{(p-1)T(n,1)} + PK(n, p)$$

The wasted time is scaled by this as is divided among processor. Only one processor will work on it and rest will be idle. Thus, time wasted

$$= (p-1)\sigma(n)$$

$\sigma(n)$  = serial time

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

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

$$= pT(n, p) - \phi(n) - \sigma(n)$$

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

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

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

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

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

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

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

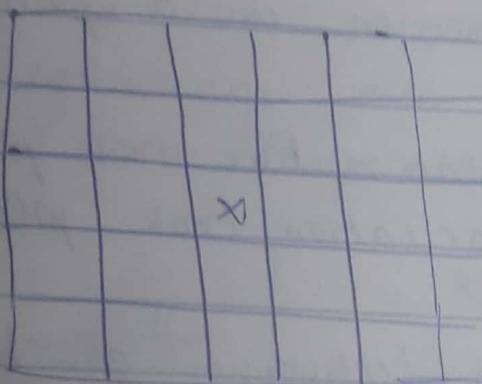
$$e = \frac{\gamma_\psi - \gamma_p}{1 - \frac{1}{p}}$$

September 20, eg. Cellular automata → logic → physics  
TUESDAY

Conway's Game of life → evolution  
Society

Engineering Problems of

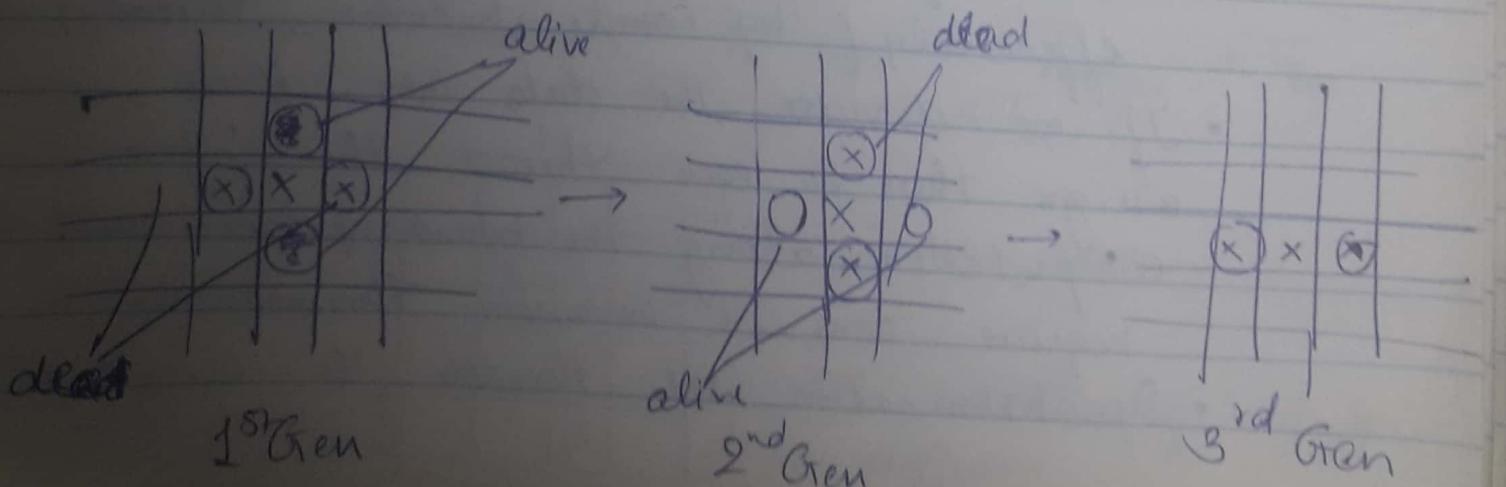
- Any social problem is dynamic in nature in both space and time
    - 2-D domain
    - States - sit, stand, lie
- States may have rules/ may evolve



Size of problem can be controlled by no. of generations we. with the problem to continue.

Size of problem does not remain in the 2-D domain.

States of the problem → live, die.



This process will continue in an oscillatory system.

- Application:

- i) Studying pattern

- ii) Evolution of society - Society may have started at some location & then shifted to some other location (also became a bigger comm.)

- iii) How to spread rumours in a particular scenario?

- iv) Location in a place - Fire took place.  
How will evacuation took place?

- Which parallelization techniques can be used?
  - Every step (in the next generation) depends on the previous step.
  - Spatial locality will be used because we are looking at neighbours.
  - Data from adjacent cells can be used effectively (by comm. between threads)
  - If we divide the data into blocks, assign these to each thread, then for computation along the boundary we will require comm. between threads
  - Synchronization bet. threads is essential

before moving to the next generation.

(Load balance / load imbalance)

- Load Balance - Same no. of FLOPs are done by each

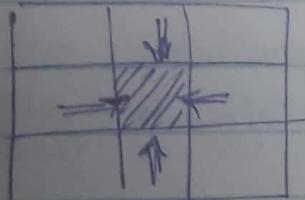
e.g. Simulating Heat transfer.

Differential equation

Application: Heat flow in processes

Heat flow depends on the gradient of temperature.

If we take a cell and look at the neighbouring cells if they are at higher or lower temperature.

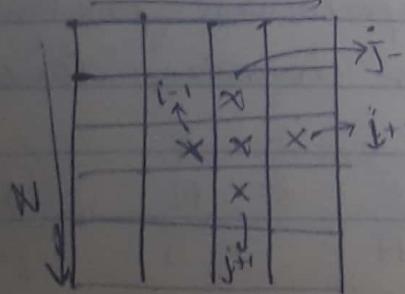


We should stop when there is no heat flow. Steady state

$$T_{\text{new}} = T_{\text{old}} + k \left\{ (T_{\text{top}} - T_{\text{old}}) + (T_{\text{bottom}} - T_{\text{old}}) + (T_{\text{left}} - T_{\text{old}}) + (T_{\text{right}} - T_{\text{old}}) \right\}$$

$$T_{\text{new}} = T_{\text{old}} + k \left\{ \frac{T_{\text{top}} + T_{\text{bottom}} + T_{\text{left}} + T_{\text{right}}}{4} - T_{\text{old}} \right\}$$

September 21,  
WEDNESDAY.



$$\Delta Q = \frac{\partial Q}{\partial t}$$

The behaviour of water/gas diffusion

The problem will try to reach the steady state (based on law of nature). This depends on the neighbourhood.

$$\text{(time)} \frac{\partial Q}{\partial t} = \frac{\partial^2 Q}{\partial x^2} + \frac{\partial^2 Q}{\partial y^2} = \Delta Q \text{ (space)}$$

Try to vectorize the equation

$$Ax = B$$

$$= Q(i+1, j) + Q(i-1, j) - 2Q(i, j) \\ + Q(i, j+1) + Q(i, j-1) - 2Q(i, j)$$

System of linear eqn → Put it in a matrix →  $x = A^{-1}B$  → Solve!  
use matrix multiplication

Engineering & mathematical problems  
diagonal / tri-diagonal elements → sparse matrix.  
are present

$$\begin{bmatrix} & & 0 \\ & & \\ 0 & & \end{bmatrix}$$

Matrix multiplication

Sparse matrix using data compression.

Data compression will increase the possibility of the data in cache.

$$N \begin{bmatrix} & & 0 & 0 \\ & & 0 & 0 \\ & & 0 & 0 \\ & & & \bullet \end{bmatrix}$$

Such sparse matrix can be compressed. Computational problem generally have boundary condition.

$$\frac{N^2}{5N} \Rightarrow \frac{N}{5}$$

Reduce by a factor of  $\approx N$

→ Compression of Matrix:

$$\begin{bmatrix} 0 & 0 & 4 & 5 \\ 0 & 0 & 0 & 0 \\ 2 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \xrightarrow{\text{4 rows}}$$

Compress Sparse Row (CSR):

$$\text{elements} \rightarrow [4 \ 5 \ 2 \ 1 \ 1] \quad \text{row\_no.} \times \text{row\_length}$$

$$\text{Column} \rightarrow [2 \ 3 \ 0 \ 1 \ 2] \quad + \text{col\_index}$$

$$\text{row} \rightarrow [0 \ 2 \ 2 \ 4 \ 5] \quad \begin{array}{l} \text{show} \\ \text{the amt of} \\ \text{memory} \end{array} \quad \text{for the first two arrays}$$

$$\begin{array}{cccc} (2-0) & (2-2) & (4-2) & (5-4) \\ =2 & =0 & =2 & =1 \\ (\text{elements} \quad \text{in 1st row}) & (\text{end row}) & (\text{3rd row}) & (\text{4th row}) \end{array}$$

But that cannot be used  
as the matrix is  
sparse. All elements-

The last index-element

shows no. of elements (because row length is not fixed)  
present in the compressed matrix.

$$\text{length of element array} = \text{no. of elements} = \frac{\text{length of column array}}{\text{row}}$$

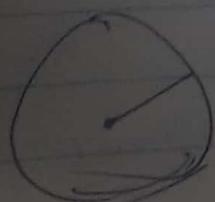
$$\begin{bmatrix} 6 & 0 & 4 & 5 \\ 7 & 0 & 0 & 0 \\ 2 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \quad \text{row} \rightarrow [0, 3, 4, 6, 7]$$

$$\begin{array}{l} \text{length of row} \\ \text{array} = N+1 \end{array}$$

→ Tot. no. of  
elements in sparse  
matrix.

$$\rightarrow NxN$$

eg. Generating image → using complex numbers (Mandelbrot Set)  
Evolve complex numbers to generate images.



Size of image → radius changed.

September 23, Scan (Prefix Sum)

FRIDAY

$$A = [a_0, a_1, \dots, a_{n-1}]$$

$$\text{scan}(A) = [I, a_0, (a_0 + a_1), \dots, (a_0 + a_1 + \dots + a_n)]$$

$\top$  - identity for the operation.

eg  $3 \ 1 \ 7 \ 0 \ \top \ 1 \ 6 \ 3 \ \top \ \top$

Don't need to keep.  $0 \ 3 \ \top \ 11 \ 11 \ 15 \ 16 \ 20$

In sparse matrix, the third operation is this.

A little like reduction but here, the dimensions of the result are the same as the input.

The final total gives you the length required to store all the values.

Don't need to

Can overwrite the output onto the same array.

Difficult to parallelize due to dependence on the previous value.

Inclusive Scan

3 4 11 0 15 16 22 25

Includes the last element.

Exclusive Scan

Include identity and not last

- Can go from inclusive to exclusive and vice versa?

→ Reduction:

- For reduction, number of step computations is still the same. No. of steps is much less.

Work complexity → work efficient

(Time) Step Complexity → step efficient.

- Follows the simple law of Divide and Conquer

No. of active cores halves at each step.

→ Scan.

→ Scan.

$$\text{No. of active cores} = \frac{N - \text{Stride size}}{2^{\text{Stride size}}}$$

Step Operations at each

$$\text{Step} = N - 2^i \quad \sum_{i=0}^{\log_2 N} (N - 2^i)$$

$$= \left[ \sum_{i=0}^{\log_2 N} (N - 2^i) \right] \log_2 N = N \log N - \frac{N}{2^{\log_2 N}}$$

$$= N \left( 1 \frac{\left( \left( \frac{N}{2} \right)^{\log_2 N} - 1 \right)}{1/2 - 1} \right)$$

$$= N \left( \frac{1}{2^0} + \frac{1}{2^1} + \frac{1}{2^2} + \dots + \frac{1}{2^{\log_2 N}} \right)$$

$$= N \left( \frac{1/2^{\log_2 N} + 1}{1/2} \right)$$

$$- \log N.$$

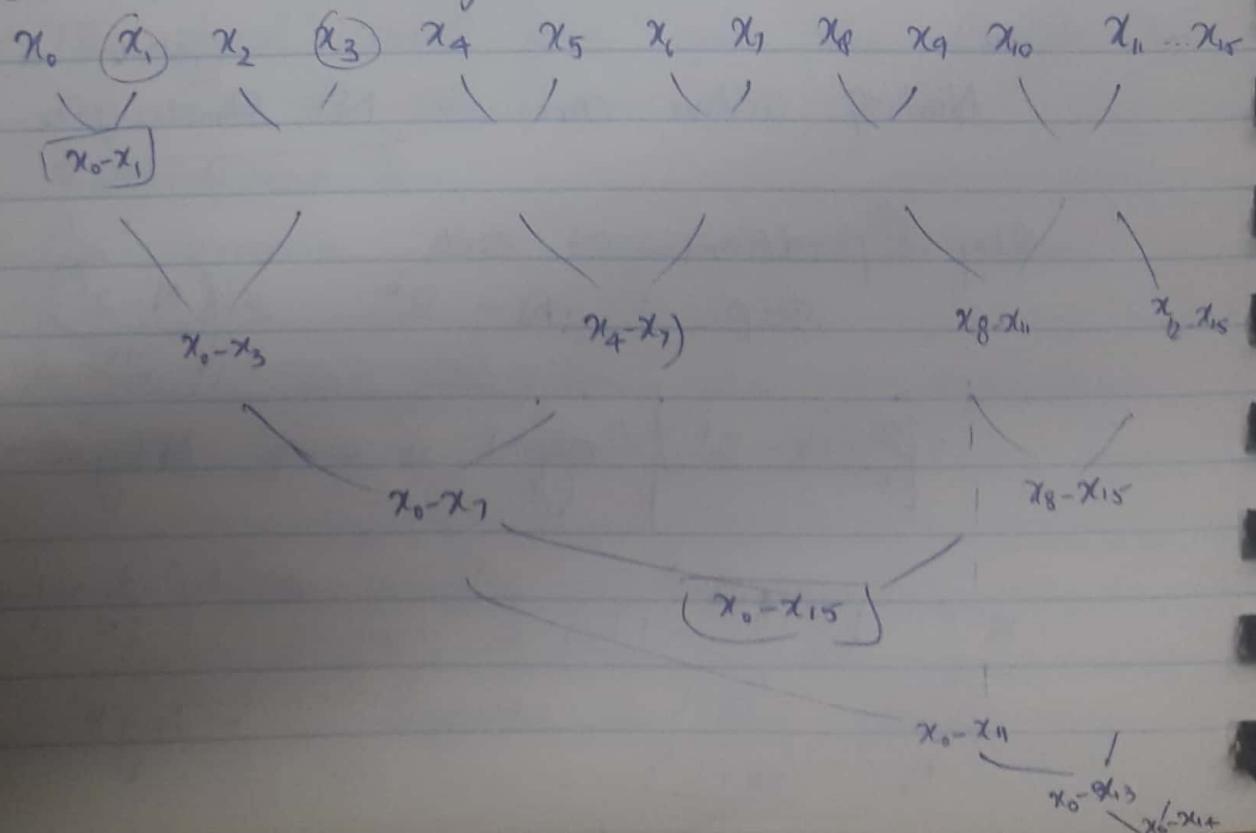
September 27, Reduction is a pattern, not necessarily  
TUESDAY always a sum. "Reduction tree"

### Reverse Reduction:

This is used to parallelize scan.  
After using reduction once, you have a lot of partial unused sums stored in the array.

Through this, by creating (almost) a reverse reduction tree i.e repeated same steps in opposite direct order, you can compute all the scan intermediate values.

[Store sums at higher order index]



## Scope of optimization:

- (i) Processor
- (ii) Memory
- (iii) Data path

### • Computer To Memory Ratio

eg From simple to block matrix multiplication  
we increased this ratio.

Basically, try and reuse memory.

## Three types of problems:

### (i) Throughput

Not necessarily much of data access but  
lots of computation eg. image generation

### (ii) Server / Bandwidth

More interested of bandwidth is high  
communicating with other processor

### (iii) Computational Science

Both.

In all three, we are basically hiding  
latency.

### • Analytical speedup of Reduction

$$\left(\frac{n}{\log n}\right)$$

Serial  
Parallel

$$\text{Overhead} = \text{Cost of Parallel} - \text{Cost of Serial}$$

$$T_0 = P T_p - T_s$$

NOTE: Calculate efficiency as  $\frac{T_s}{P T_p} = \frac{1}{1 + \frac{T_0}{T_s}}$  in assignments

$$\text{Speedup of Scan} = \frac{n}{n \log n}$$

So, no speedup.

### Reduction

$$\text{Step} \rightarrow \log n$$

$$\text{Work} \rightarrow n$$

$$\text{Cost} \rightarrow n \log n$$

$$\frac{n}{\log n}$$

$$\frac{n}{n}$$

$$\frac{n}{n \log n}$$

Not cost optimal. Work / Step optimal.

September  
October 28,  
WEDNESDAY

$$\text{Cost of parallel code} = (\text{No. of processors}) * (\text{No. of steps for parallel code})$$

We assume that  $n=p$ . But not always so.

So, look at efficiency

$$\text{Efficiency} = \frac{\text{Speedup}}{\text{no. of processors}}$$

$$= \frac{T_s / \text{cost}}{\text{cost}} = \frac{S}{P}$$

$$\text{Cost} = \frac{T_s}{\text{Efficiency}} = \frac{T_s}{S/P} = \frac{P T_s}{S}$$

Cost optimal if  $\text{Eff.} > 1$

$$\text{e.g. Speedup} = \frac{T_s}{T_p} = \frac{n}{\log n}$$

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

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

$$\frac{T_{\text{Serial}}}{T_{\text{Parallel}}} = \frac{n \log n}{(n \log n)^2} = \frac{1}{n \log n}$$

$$\frac{T_{\text{Serial}}}{T_{\text{Parallel}}} = \frac{1}{n \log n}$$

$\rightarrow$  Not cost optimal by a factor of  $\log n$

Now, for  $P$  processors,

$$T_p = \frac{n(\log n)^2}{P}$$

Time taken by  $P$  processors:

Each processor takes  $T_p = \frac{n(\log n)^2}{P}$  time thus total time for  $n = P n \log n$ .

$$\text{Speedup} = \frac{P}{\log n}$$

This shows that speedup decreases if problem size increases. Hence, not cost optimal.

Scaling Down: Assume  $n$  processors <sup>for n problem size</sup>  
 total Then reality is  $P$  processors, but still assume total  $n$ . (Pad virtual processors which are sitting ideal)  
 So, each processor has to do more computations by factor of  $\frac{n}{P}$

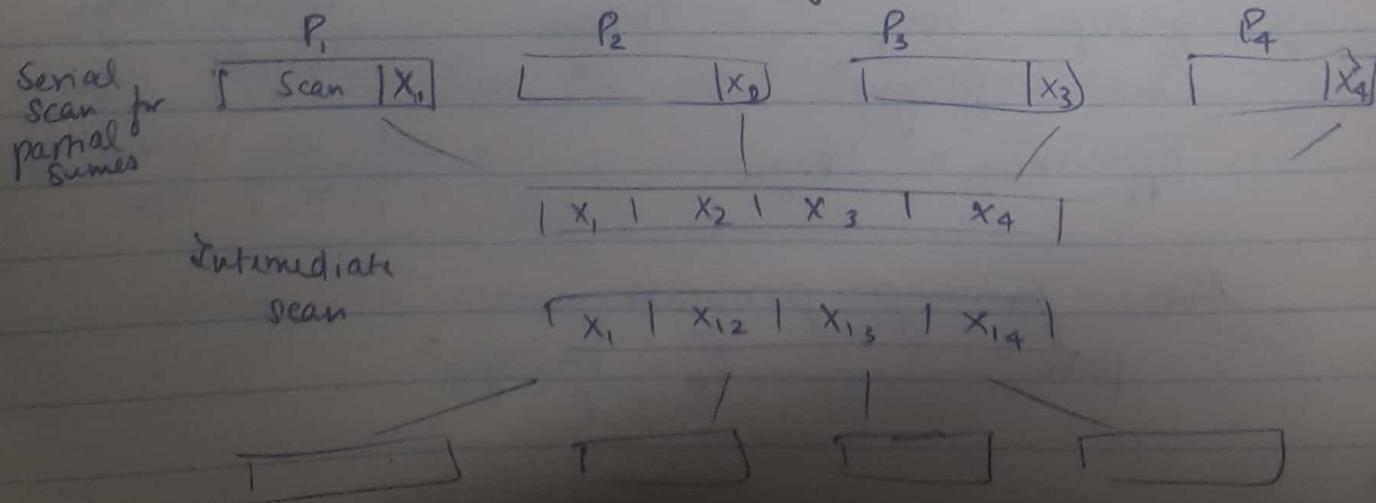
For reduction,  $P < n$

$$O\left(\frac{n}{P} \log P\right) \sim O(n \log P)$$

Step 1: Serial Time =  $n P$

Step 2: Adding partial sum =  $\log P$

$$\therefore T_P = O\left(\frac{n}{P} + \log P\right) \sim O(n) \quad P < n$$



September 30,  
FRIDAY.

$$C = \frac{T_s}{P T_p}$$

$$C = \frac{1}{1 + \frac{T_o}{T_s}}$$

|                            |   |   |   |    |    |    |     |
|----------------------------|---|---|---|----|----|----|-----|
| Processing<br>Elements (P) | 2 | 4 | 8 | 16 | 32 | 64 | 128 |
|----------------------------|---|---|---|----|----|----|-----|

|                           |   |   |   |    |     |     |     |
|---------------------------|---|---|---|----|-----|-----|-----|
| Speed up                  | 2 | 4 | 8 | 16 | 32  | 64  | 128 |
| Processors<br>size<br>inc | 2 | 3 | 6 | 8  | ... | ... | 80  |

Overhead / communications increase with  
increase in problem size.

$\sigma(n)$  is the part that cannot be parallelized.  
So, this is done only by one processor  
and rest are idle then.

$K(n, p)$  overhead for each processor.

$$\therefore \text{Idle time} = (p-1)\sigma(n) + pK(n, p)$$

This is a major reason why speedups  
don't come up to  $p$  theoretical  
speedups.

$$e = \frac{1/s - 1/p}{1 - 1/p}$$

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

$e \begin{cases} \text{constant} \\ \text{increasing.} \end{cases}$

- If  $e$  is constant, with inc. in processors, it shows that the  $\sigma(n)$  factor is the restricting factor.  
This is because  $K(n,p)$  depends on  $p$  and but its effect (though it increases) is not reflected in  $e$ .  
 $\sigma(n)$  only depends on  $n$ .
- If  $e$  is increasing,  $K(n,p)$  affects.

→ Scalable Parallel Systems: Isoefficiency:

Cost optimum algorithms have efficiency of  $O(1)$

{ No. of processors  
Problem size

If you reach a part where  $e$  becomes constant, you cannot optimize further because it means that you are going into serial part and that cannot be touched.

$$\begin{aligned} \text{Speedup } \Psi(n, p) &\leq \frac{\sigma(n) + \phi(n)}{\sigma(n) + \frac{\phi(n)}{p} + k(n, p)} \\ &= \frac{p(\sigma(n) + \phi(n))}{\sigma(n) + \phi(n) + \underbrace{(p-1)\sigma(n) + pk(n, p)}_{\text{Tot. overhead}}}, \end{aligned}$$

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

$$= \frac{1}{1 + \frac{T_0(n, p)}{\sigma(n) + \phi(n)}} \geq \frac{\text{parallel overhead}}{\text{serial cost}}$$

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

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

$$T_o(n, i) \geq \left( \frac{E(n, p)}{1 - E(n, p)} \right) T_o(n, p)$$

We want to keep this ratio constant. ie we want efficiency to remain constant.

$$T(n, i) = \text{constant } T(n, p)$$

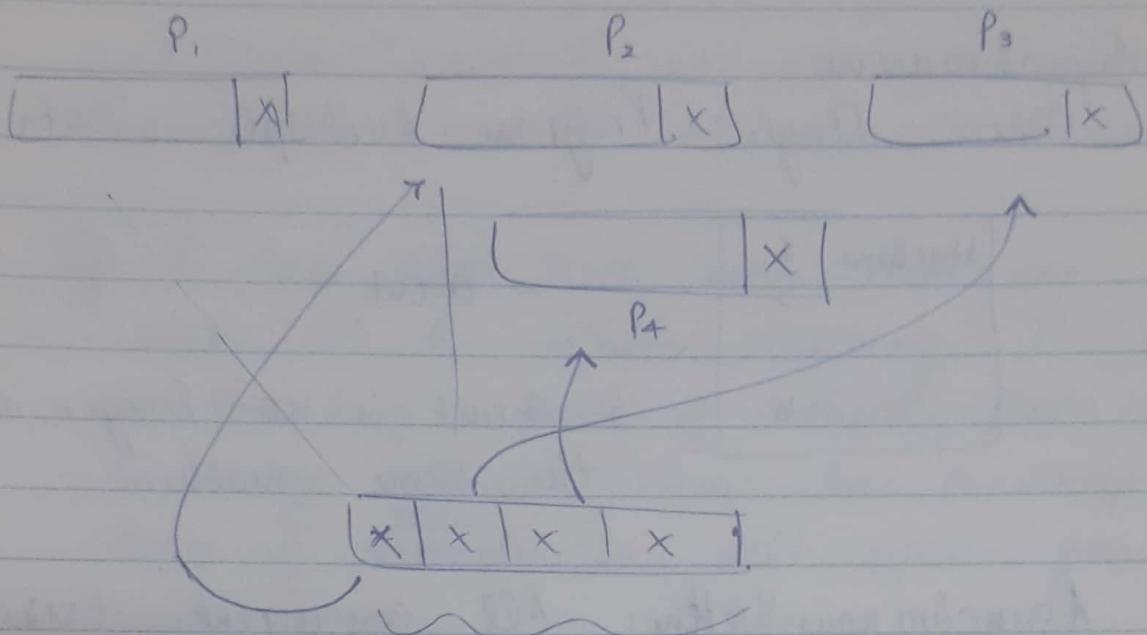
To keep it constant, we have must increase the problem size along with the number of processors.

With only increase in processors, efficiency would decrease.

There are other factors that affect efficiency as well. Memory access, Cache organization, etc.

Thus, Scan and Reduction are designed only for equal size-processors.

If size is more, idle time increases. Thus, reduce problem size to no. of processors first and then apply reduction on that.



Apply Scan here. on last elements.  
 Now, add each of the equivalent  
 elements to corresponding threads  
 subarrays to get the complete  
 scan. Can used exclusive scan!

October 4, (i) Shared Memory Parallel Computer  
 TUESDAY (ii) Distributed - Memory " "

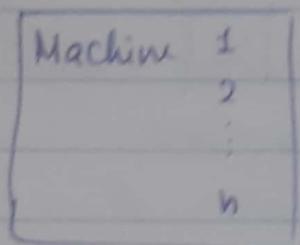
OpenMP cannot be used or distributed.  
 So, we MPI (

In shared, you create a.out and  
 then the OS directly takes data.

In distributed, you create a.out and  
 then explicitly send it to each system.

Asynchronous

SPMD: Single Program Multiple Data



a-out

a-out is no longer on  
the same machine.

Asynchronous System: All concurrent tasks  
execute asynchronously.

P<sub>0</sub>

a=100;  
send (2a, 1, 1);  
a=0;

P<sub>1</sub>

receive (2a, 1, 0);  
printf("%d\n", a);

We want P<sub>1</sub> to receive a = 100 and not  
a=0. So,

- One method is to wait for P<sub>1</sub> to receive a before P<sub>0</sub> advances to a=0 command.
  - ↳ This causes lots of overhead.
- We can use ~~to~~ buffered space instead,  
Copy a into another variable b and  
send b. Thus, P<sub>0</sub> can manipulate b & a  
with any issues.

no. of  
elements size

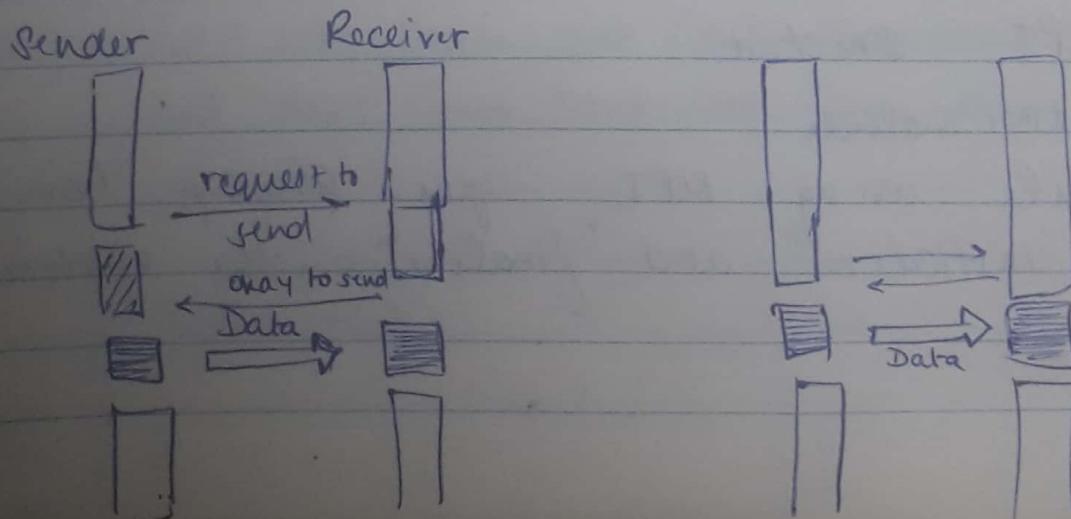
```

send(Void *sendbuf, int nelems, int dest)
receive(Void *recvbuf, int nelems, int source)

```

- If  $P_0$  can also send more than one data to  $P_i$ . While receiving,  $P_i$  has to keep a track of which data is being received. Thus, has to keep a track of <sup>order</sup> variables as well apart from just source.
- In non-buffered blocking, idle time increases and deadlocks are a big issue.
- In buffered, the sender copies the data into the designated buffer and returns after the copy operation has been completed. The receiver can also use the buffer and will use the data only when required by it.

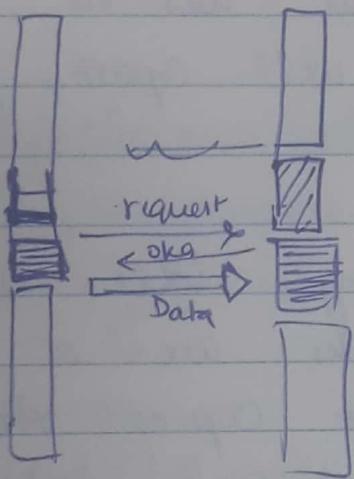
Telling overheads decrease but copying overheads increase.



Second case has minimum idling time.

However, this is very unlikely to happen because the two processes are independent and run asynchronously.

In the first case, sender has to wait for the receiver to be ready.



Here, receiver needs some information but the sender cannot send yet. This leads to idling time for receiver.

→ Buffered Blocking Message Passing Operations:

Whether buffering works depends on hardware as well

→ MPI structure.

MPI includes

- While using MPI, you always have to "initialize" and "finalize" the environment.

`MPI_Init(&argc, &argv)` ↗  
command line  
arguments  
Same as main()  
`MPI_Finalize()`

→ No other MPI call  
can be made after this.

You have several processors for each process  
but each process can have multiple threads

→ Communicators and Process Groups:

- Communicators are groups of processes that can communicate with one another.  
eg Student can be part of HPC group and DBMS, etc. However, does not have to be part of all.
- Defines a communication domain. Can only communicate in this domain.
- One process can belong to many different communication domains.

int MPI\_Comm\_size (MPI\_Comm comm, int \*size)

int MPI\_Comm\_rank (MPI\_Comm comm, int \*rank)

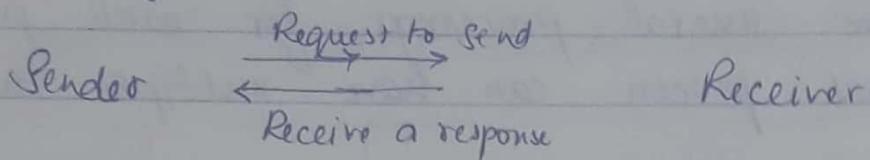
no. of  
CPUs

↓  
processes  
within that

October 7, • Save time in communication

FRIDAY - Correct information is obtained at the receiver's end.

- Communication is enabled (done correctly) using hardware.



- Buffer can reduce idle time.

Blocking Mechanism - guarantee that data will be transferred correctly

- Overlap communication with computation.  
But it is essential that we take guarantee that data required to be transferred should not be overwritten.

$P_0 \rightarrow P_1 \leftarrow P_2$  • Send data to a process anonymously is possible in MPI.

• So, we will not be able

to find if the data is correct or not (and

1 element (can be specified from which processor) whether int, float)

$P_0$   
Send(&A, 1, 1)

tags 1

Receive(&B, 1, 1)

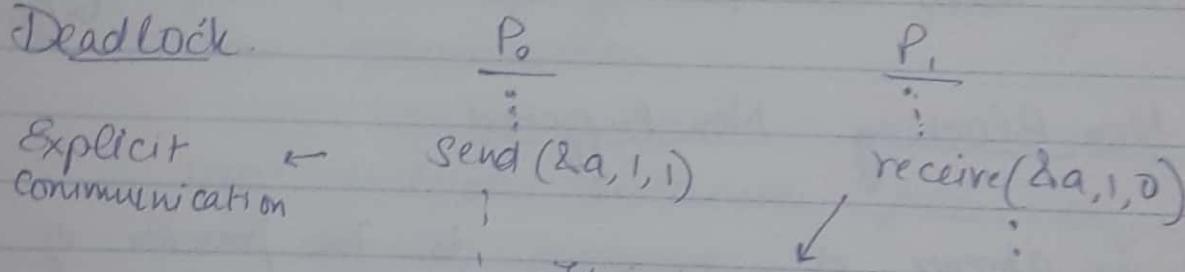
$P_1$   
Send(&A, 1, 0)

tags

receive(&B, 1, 0)

- If we use non-blocking then it will lead to dead-locking. Because the 'send' from  $P_0$  should be received by  $P_1$ .
- But  $P_1$  rather than receiving, sends data to  $P_0$ . So, now  $P_1$  is blocked as  $P_0$  is waiting until data is received from  $P_1$  (which is not sent).
- If we use non-blocking send-receive method then the code will function.

### Deadlock



If instead, the function is  
to deadlock in this blocking state  
always

- Non-blocking mechanism
  - Have to ensure that the update is done correctly
  - Use buffered space. So, fast & without incorrect code.
- CPU doesn't necessarily have to take part in the communication.  
Sender  $\rightarrow$  sending process waits till a copy of data and stores it in the buffer.
- Talking time is completely gone by using buffer.
- Place it in the buffer of receiver
- Placed in the buffer of hardware. Communication buffer

### i) Blocking buffered transfer.

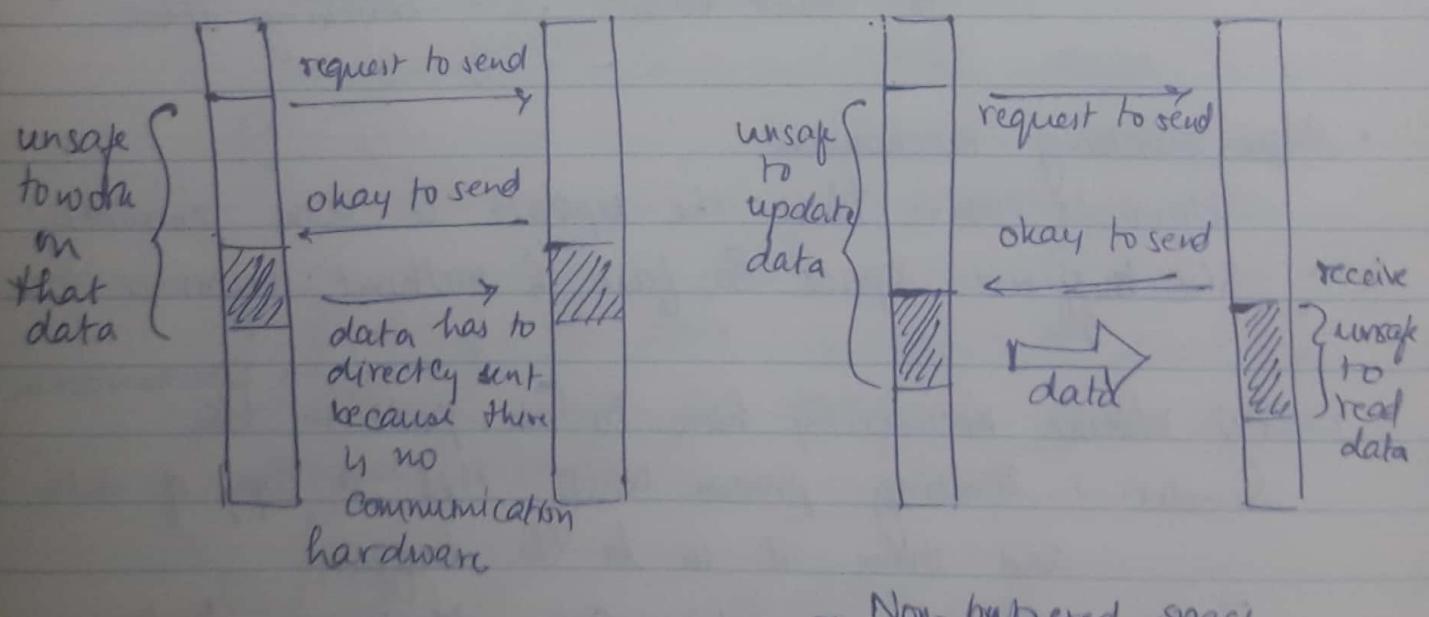
- In the presence of communication hardware with buffer at the receiver and sender's end.
- Buffer is supported by hardware
- Sender only copies data into its buffer and the receiver can access data at anytime
- Buffer is not supported by hardware
- Sender obstructs the receiver from continuing & copies data in its buffer.

### Non-Blocking

### Non-Buffered

a, In absence of communicator hardware

b, In presence of hardware



Communication hardware

→ Non-buffered space where it stores data

Simply ensures the CPU does not need to overlook communication. Assigns to hardware & continues its computation.

October 12,  
WEDNESDAY

Several processes run each process. Thus, each has to initialize and finalize MPI environment respectively. So, everything comes within that.

**MPI\_COMM\_WORLD**: All processes are under this communicator.

You can communicate between two processes only if they are part of the same communicator.

**MIND**: Multiple executable files on different processes.

However, we will mostly be working with <sup>mostly</sup> SIMD, ie one shared file

A process can be a part of several communication domains.

You can also have unicast or multicast.

- Buffered - non-buffered
- Send-receive / Sca, etc
- Type of communication.

→ Communicators:

From no of cpus, no of processors in domain and current cpu id, it can get what the part of the code that cpu is supposed to do.

e.g. `Npz_Send ( buffer, lengthsize of buffer, datatype,`  
`only length  
of data to be  
comm.  
rank of comm. proc, tag, comm. domain )`  
`to identify  
send-receive pair`

`Npz_Recv ( . . . , &status)`  
Always blocking!

Helps with error handling  
like size/length mismatch

Non-blocking receive is unsafe.

There is no 'status' argument on the send side. Other functions like `Npz_Success` which shows this.

~~mp~~ MPI\_send ("a")  
a[0] = 1.

↙

This statement will not execute unless  
the entire data of 'a' is sent.  
Thus, we can check status.

### Point-to-point communication

Collective " " : Faster than  
point-to-point. However, can be  
used in only some cases.

'Status' is an object (like a stack)  
rather than a simple variable.  
Thus, much data can be stored there.  
Used as arguments for a lot of  
other functions.

### → Collective Calls:

Synchronization is done through  
~~variables~~ 'barrier'?

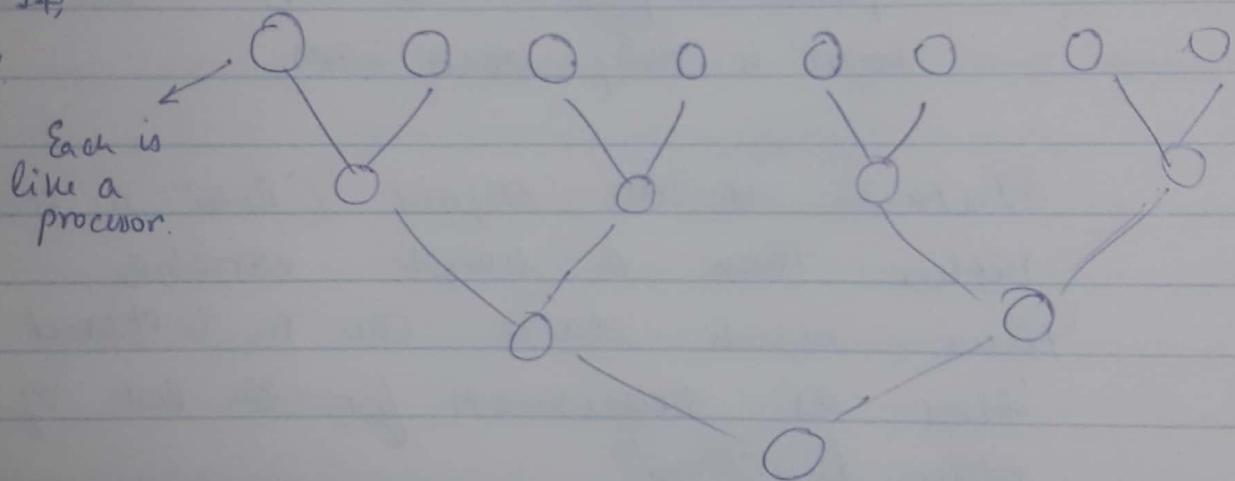
MPJ - Barrier()

Explicitly synchronize  
among processors  
within that comm. domain

- MPJ-Bcast: Broadcast to all other processes. All processes call  $\$$  Bcast with same root and communicator arguments.
- \* Look at all the function, return types and see if they can be used as arguments for elsewhere.

October 14,

FRIDAY



Since there are dependencies on the previous levels, regardless of how many processors you have, you cannot complete it in one step. This is called Span

Here, the degree of concurrency is 8. This is because we treat each node as an independent processor.

We have 8 nodes (or 8 data points) so

each of them could have executed independently.  
That would lead to 8 processes.

Don't get stuck with reduction. The tree signifies any parallel computation.  
[Reversing the tree makes it easier to see the processes. 1 process is eventually divided into 8]

Now, on reduction.

Threads =  $p$

i) Work =  $T$  (no. of computations)

Work Law :  $\boxed{pT_p \geq T_1}$

ii) Cost =  $(\text{num\_t}) \log(\text{num\_t})$  Steps =  $\log(\text{num\_t})$

iii) Span / Depth / Steps =  $\log_2^n$  if  $T_\infty = \infty$   
OR

Critical Path length (Shortest Time)

iv) Speedup =  $\frac{T_1}{pT_p}$

Step Law :  $\boxed{T_p \geq T_\infty}$

If we have  $P$  processors, no. of steps taken

would be atleast steps taken by  $\infty$  processes

If  $n > p$ ,

$$T_p = \frac{n}{p} + \log_2 p$$

$$T_{\infty} = \log_2 n$$

$$\frac{T_1}{p} \leq T_p \leq \frac{T_1}{p} + \log_2 n$$

Brent's  
Graham's  
Law:

$$\left| \frac{T_1}{p} \leq T_p \leq \frac{T_1}{p} + T_{\infty} \right)$$

$\frac{T_1}{p}$  signifies when there is a single step of computation.

RHS ie the upper bound is taken span into consideration. That assumes a unit time of computation at each step and so, adding.

y, Iso efficiency.

At what rate should you increase your problem size ( $w$ ) so that the efficiency remains constant.

$$w = k T_0(p, w)$$

$k \propto E$

For reduction,  $T_0 = p \log p$   $P \geq w$   
 [Overhead at each step for each processor]

If  $P \rightarrow P_i$

$$\frac{W_1}{w} = \frac{p_i \log p_i}{P \log P}$$

$$W_1 = \left( \frac{p_i \log p_i}{P \log P} \right) w$$

$$\frac{W_1}{w} = \frac{p^2 + p \log p + w^2 p^2}{P^2 + P \log P + w^2 P^2}$$

$$W_1 = p^2 + p \log p + w^2 p^2$$

Not always possible to write the iso-efficiency equation for all problems. May be that problem size is not exactly scalable in that way.

e.g. Matrix addition :  $w^2$  problem size.

$$\text{e.g. } W = K(p^2 + p \log p + w^2 p^2)$$

Problem size should be increased based on the asymptotic upperbound of overhead.

NOTE: If strictly talking theory, you can assume no. of processors to be equal to problem size.

$$\text{eg. } T_S = n \log n$$

$C_1, C_2, C_3, C_4 \rightarrow +$  algorithms

|            | $C_1$ | $C_2$    | $C_3$      | $C_4$             |
|------------|-------|----------|------------|-------------------|
| processors | $n^2$ | $\log n$ | $n$        | $\sqrt{n}$        |
| $T_P$      | 1     | $n$      | $\sqrt{n}$ | $\sqrt{n} \log n$ |

|                                  | $n \log n$ | $\log n$ | $\sqrt{n} \log n$ | $\sqrt{n}$ |
|----------------------------------|------------|----------|-------------------|------------|
| [ $C_1$ is most optimum by this] |            |          |                   |            |

|  | $\frac{\log n}{n}$ | 1 | $\frac{\log n}{\sqrt{n}}$ | 1 |
|--|--------------------|---|---------------------------|---|
| [ $C_2$ & $C_4$ are most optimum here] |                    |   |                           |   |

|                                    | $n^2$ | $n \log n$ | $n \sqrt{n}$ | $n \log n$ |
|------------------------------------|-------|------------|--------------|------------|
| [ $C_2$ & $C_4$ most optimum here] |       |            |              |            |

Overall,  $C_2$  is the most optimum.

Q?

How do you get check optimum minimum time with minimum no. of processors?  
ie minimum time and cost optimal.

Ans:

$$T_P = \frac{n}{P} + \log_2 P$$

$$\frac{d T_P}{P} = -\frac{n}{P^2} + \frac{1}{P} = 0$$

$$\therefore P = n$$

Thus,  $P=n$  gives minimum time.

$$\text{Cost} = P \times \log p$$

$T_p = \log p$   
Nb. of steps = time taken in parallel

$$E = \frac{S}{P} \quad \therefore P = \frac{S}{E}$$

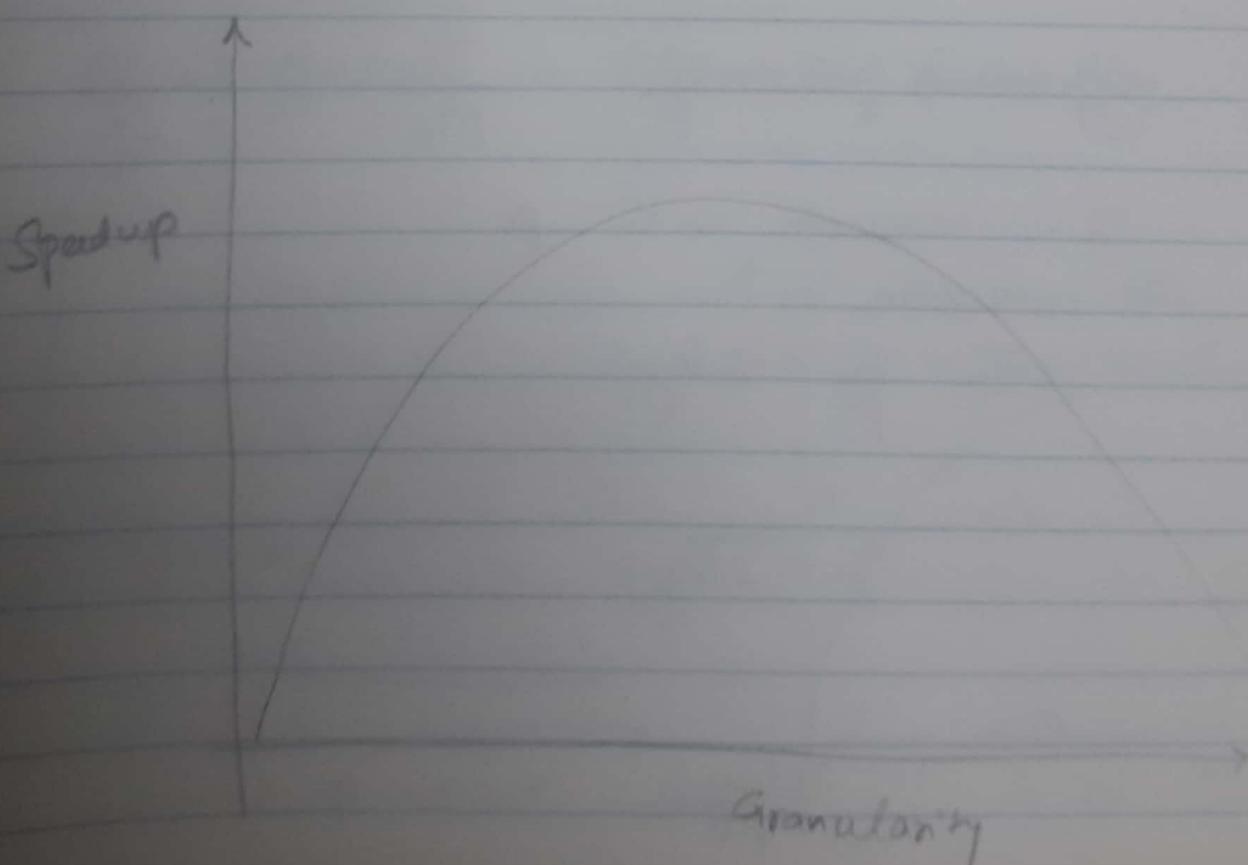
$$C = \frac{S}{E} \times \log p$$

$$= \frac{\frac{T_S}{T_P}}{E} \times T_P$$

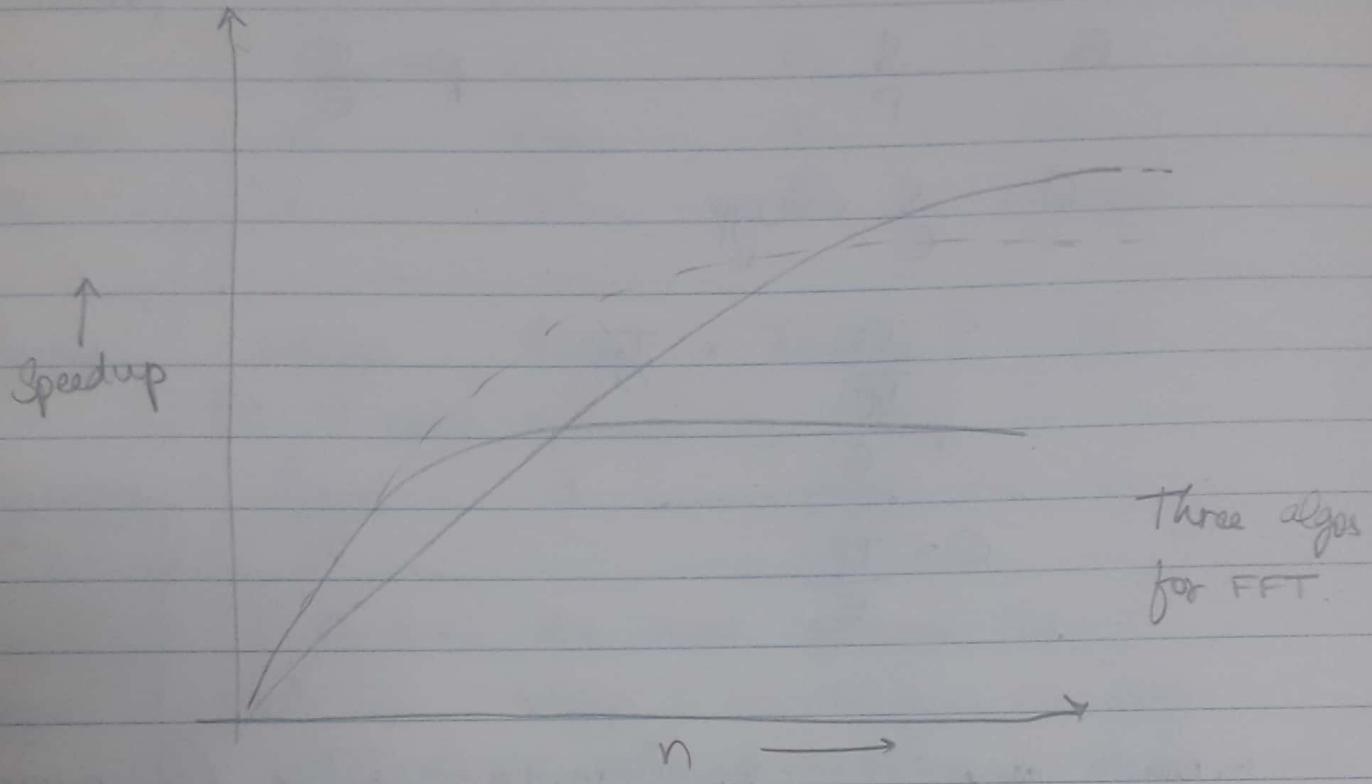
$$C = \frac{T_S}{E}$$

September 30  
FRIDAY

## Effect of Granularity on Performance



# Scalability of Parallel Systems



Increasing cost with increasing  $n$ .

$$\text{Efficiency } E = \frac{S}{P} = \frac{T_s}{P T_p} = \frac{1}{1 + \frac{T_o}{T_s}}$$

$T_o$  increases with  $P$ .

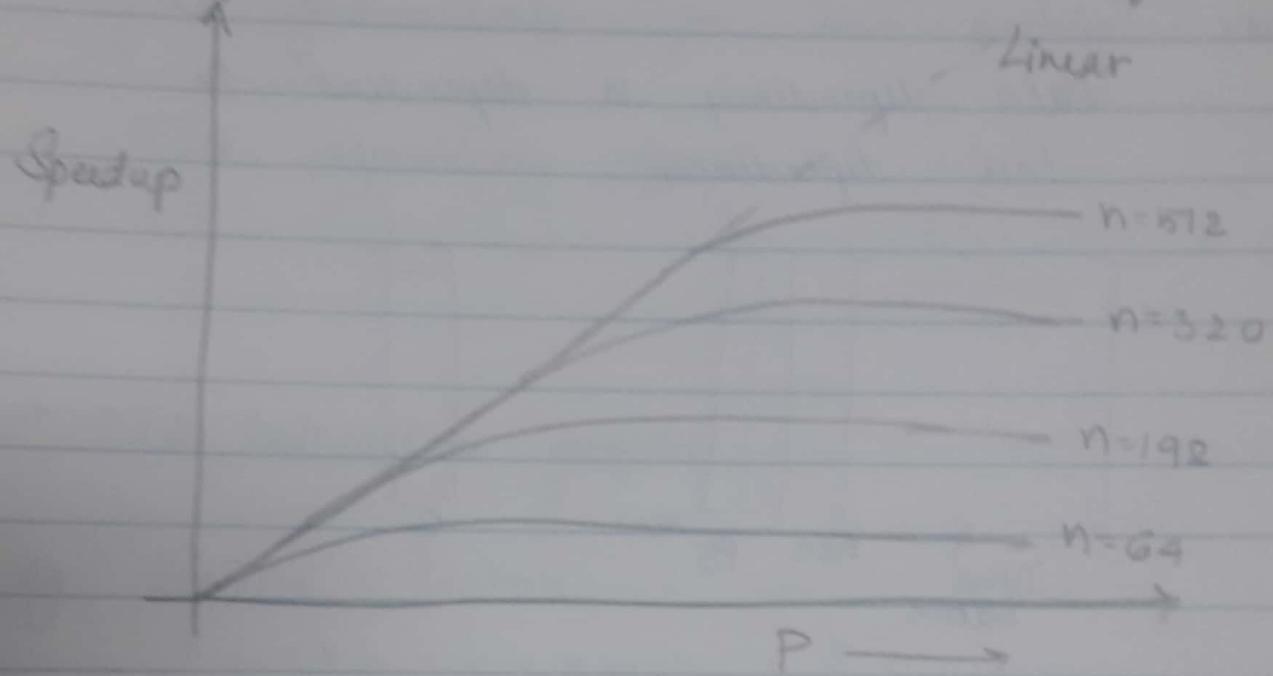
For reduction,

$$T_p = \frac{n}{P} + 2 \log p$$

$$S = \frac{n}{\frac{n}{P} + 2 \log p}$$

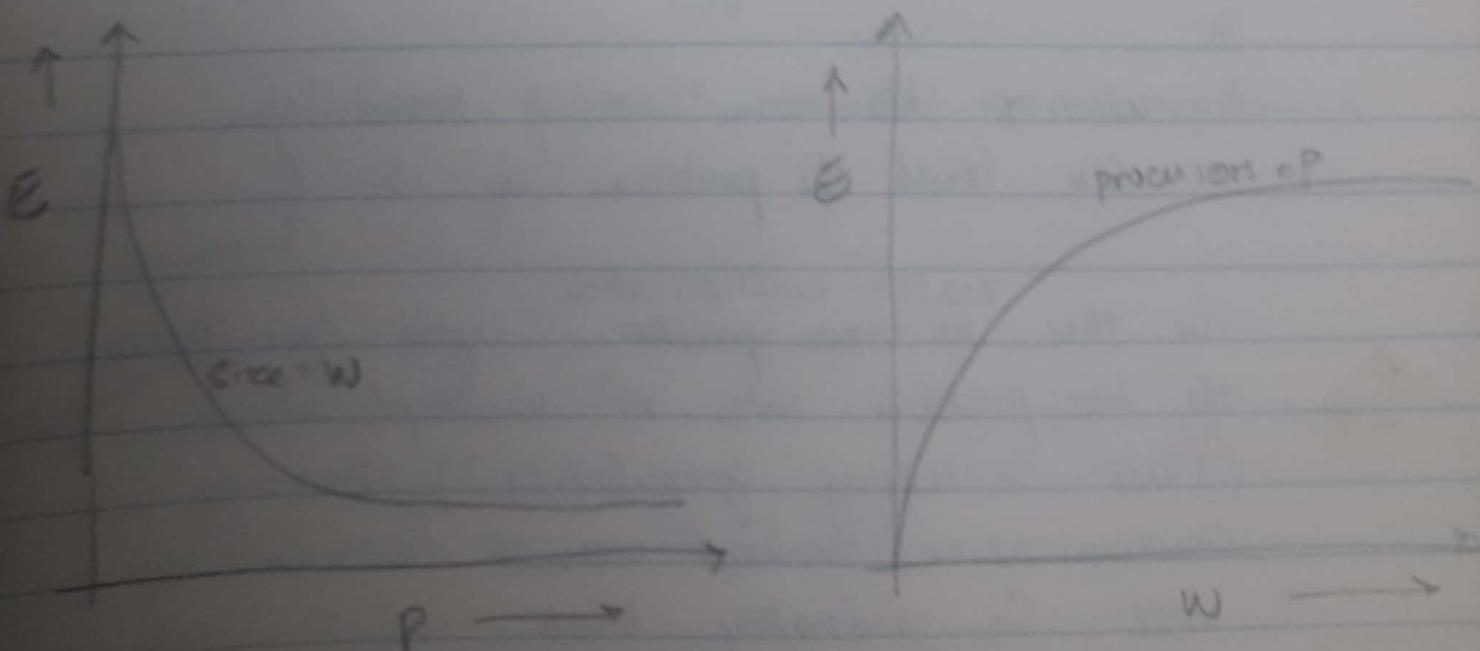
$$E = \frac{1}{1 + \frac{2 \log p}{n}}$$

## Scaling Characteristics of Parallel Program



When assumed  $p=n$ , instead of linear curve, we get saturating curve.

## Efficiency - Metric of Scalability



October 21, Decomposition, Tasks and Dependency graphs:  
FRIDAY.

Data dependency is dependent on task dependency.

$$\begin{bmatrix} \text{---} & | & | & | & | \\ & * & | & | & | \\ \text{A} & n \times n & \text{B} & n \times 1 & = \\ \text{Dense Matrix.} & & & & \text{C} \end{bmatrix}$$

Each row of A will require entire B. There is no task dependency since each element of C can be calculated individually.

However, data dependency exists.

- You will want to increase granularity to make it more efficient.

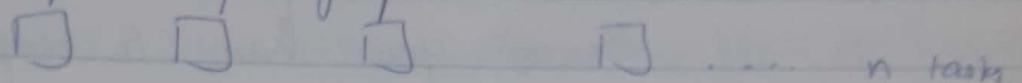
¶

Granularity - Maximum no. of tasks you can break the problem in.

For this matrix multiplication, you can break the total problem into  $n^2$  multiplications. (Thus increasing granularity?) However, communication overhead would increase as well. So, generally break in n tasks.

Thus, there is an inherent bound on how fine granularity for each problem.

Task dependency graph.



→ Task interaction graph:

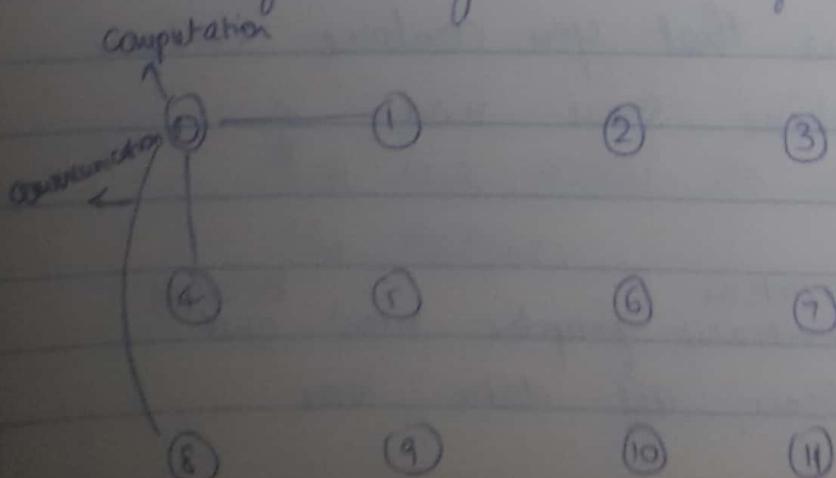
e.g. Sparse Matrix

|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|---|----|----|
| 1 | x | x |   |   | x | . | . | x |   |   |    |    |
| 2 | x |   |   |   |   |   |   |   |   |   |    |    |
| 4 | x |   |   |   |   |   |   |   |   |   |    |    |
| 8 | x |   |   |   |   |   |   |   |   |   |    |    |

b

|      |        |   |
|------|--------|---|
| b[0] | node 0 | - |
| b[1] | node 1 | - |
| -    | -      | - |
| -    | -      | - |

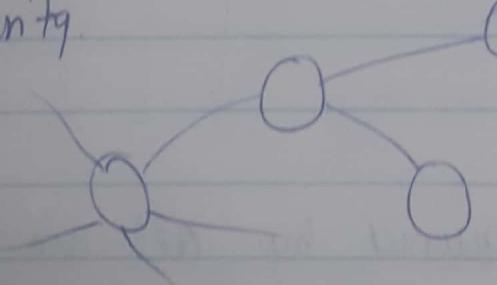
$b[0]$  is required by all non-zero elements of the first column of A.



- These graphs are independent non-directional. Only shows the interaction between two tasks. Who sends or receives doesn't matter.
- [In distributed memory, b isn't available to all. Hence, you have to explicitly send it to other processes through MPI]

Task dependency graphs are inherently directed. Each level is computed only when the previous has been computed.

Now, once the entire graph is formed, you see that each task has to send data to multiple other no. tasks. This increases communication overheads and thus reduces granularity.



Multiple nodes interact with the same node. Can combine.

One solution is that you combine certain tasks. i.e. combine some nodes of the graph.

(Talk about interaction dependency graphs now because on shared systems, all data was)

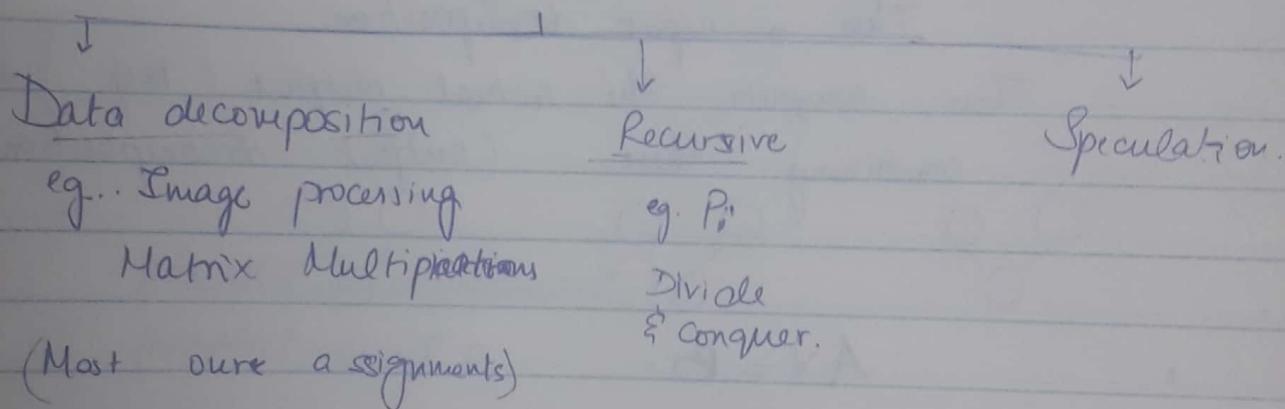
available to everyone. Not true in MPI)

## → Processors and Mappings:

- Determined by both dependency and interaction graphs.

Task dependency → for load balance.

Task interaction → communications reqd



## Output Data Decomposition

eg. Matrix multiplication.

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

Each task works on a its part of the output.

Input Data Decomposition  
eg. Simple filtering.

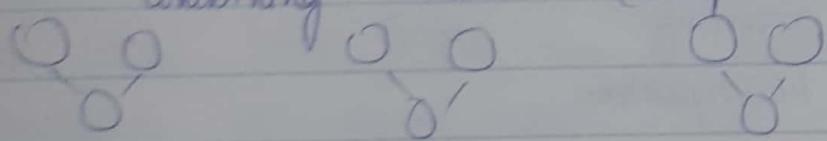
Each task only works on its part  
of the input.

Intermediate

eg. Block multiplication.

Compute partial sums from each block.  
This is input decomposition.

Then, compute the actual output by  
combining these. (output decomposition)



$$A \cdot X = B$$

Very memory intensive.

Avoid taking  $A^{-1}$  because if  $A$  was a sparse  
matrix, the inverse would become dense.

October 25, IN-SEM 2.

TUESDAY, 13, Called loop unloading.

do i = 1, n-1 , i+2

$$b(i) = a(i) + S^* a(i+1)$$

$$b(i+1) = a(i+1) + S^* a(i+2)$$

end do

3 accesses for 4 comp instead of 2 for 2.

$$\begin{array}{r} 0 \\ 0 \quad 1250000 \\ \hline (1025000+1) + 25000 \\ 25000 \end{array}$$

SUM

#pragma  
for (1, 10<sup>7</sup>):

sum += f()

$$4 \left( 1 - \frac{1}{3} + \frac{1}{5} - \dots \right)$$

parallel-measurement.c  
parallel-script.c