

System: Scientific simulation using
gcc, g++ compiler.

OpenMP - shared memory system.

www.letsthpcc.org → Refer this website to generate reports using this.

Performance = ?

- | | |
|--|---|
| <u>algorithms</u> <ul style="list-style-type: none"> - they are deterministic - they are abstract - terminates - we cannot isolate system and learn parallel programming - our agenda is to solve a problem correctly. - In HPC, we have a change in agenda which is efficiency. | <u>computer hardware</u> <ul style="list-style-type: none"> - Performance is non-deterministic - not abstract - not terminating - no parallel programming - agenda is accuracy - efficiency |
|--|---|

unit of performance: computers per second. /

floating point operations per second?

flops per second.

Performance will be written in terms of

architecture: ~~in terms of architecture~~

$$\text{Performance} = \text{frequency of processor} \times \text{No. of cores} \times \text{flops}$$

$$= \frac{\text{cycles}}{\text{second}} \times \frac{\text{flops}}{\text{cycle}} \times \text{No. of cores}$$

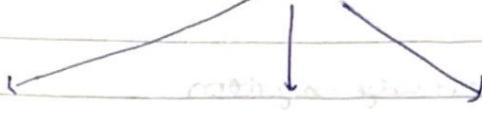
$$2.5 \times 10^9 \text{ cycles/second} \times 4 \text{ flops/cycle}$$

$$\text{Typical performance} = \frac{2.5 \times 10^9 \text{ cycles}}{\text{second}} \times \frac{4 \text{ flops}}{\text{cycle}}$$

$$\text{Typical processor speed} = 10 \text{ flops/second} = 10^{10} \text{ flops/seconds}$$

$$\text{Typical ability} = 40 \text{ flops/second.}$$

Performance is decided by 3 important factors.



architecteur algorithm software
parallel environment

- This is all what defines HPC.
 - Designing the algorithm is critical, architecture needed to optimize the problem, it will give performance to the processor.
 - Main problem is to design architecture that fits well with system such as cache.
 - Code must be scalable and portable.

1. Options, Futures and Derivatives : J C Hull and Basu : 9th edition (2014) - Amazon
2. Simulations, By Ross. - Amazon - Amazon

For N elements, we have $O(N)$

Computation = Operation + Data.

$\mu = 10^9$ N m^{-2} $\text{mole}^{-1} \text{dm}^{-3} \text{ s}^{-1}$ mol^{-1}

$O(N)$ → computation complexity is only taken into consideration, the memory complexity is not taken into consideration.

If $N = 10^{20}$, everything does not go into RAM, accessibility has to be taken into consideration.

- unit of operation : pccm / second.

$$\text{rate} = 10^9 \text{ additions} = 1 \text{ g/sec} / \text{second}$$

10 seconds. 1.28 minutes long.

* Theoretical performance = 10 Gbytes/second,

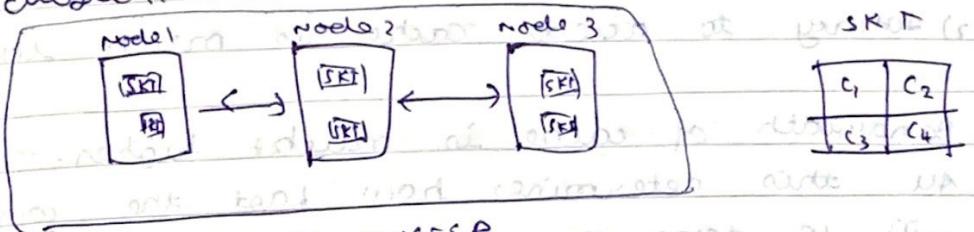
so try to find a way to reach 10 Gbps/sec.

list of things in system architecture:-

- 1) architecture - non deterministic
- 2) algorithm design - Deterministic
- 3) software environment - non deterministic

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

- nodes are connected with each other in a cluster.



- Inside a node, there is a socket.
- the cores perform 2 hops in the cluster.
- Inside a socket there is a core

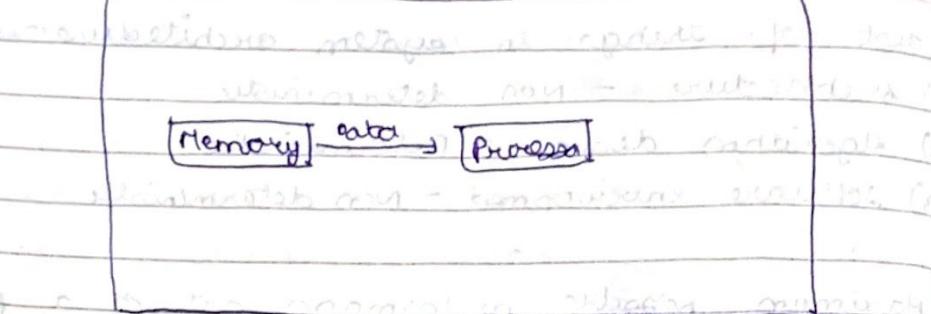
supercomputer = several nodes

node performance = CPU speed * no. of CPU cores

CPU instruction per cycle * no. of CPUs per node

- Report about hardware, match performance

- from where performance will come in node
- device technology in cache sits on processor, RAM sits outside processor, so cache is much faster than RAM



Important factors to consider while memory access

- 1) Bandwidth. \rightarrow MB/sec, GB/sec.
- If bandwidth is not enough, then use the high MB/sec is useless.
- 2) Latency to access cache is much lower.
- Bandwidth of cache is much higher.
- All this determines how fast the comparison will be done.
- Write algorithm in such a way that memory complexity is reduced.

Serial:

- Break a problem into a discrete form of instructions.
- Executes instructions one by one on a single processor.

Parallel:

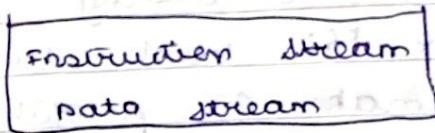
- Break problem into discrete problems.
- Each part broken into a series of instructions.
- Each instruction executes simultaneously on several processors.
- Coordination mechanism needed.
- Cannot continue to scale processor frequencies.

(no 10 GHz chips), because more power will burn processor.

- we cannot increase transistor density

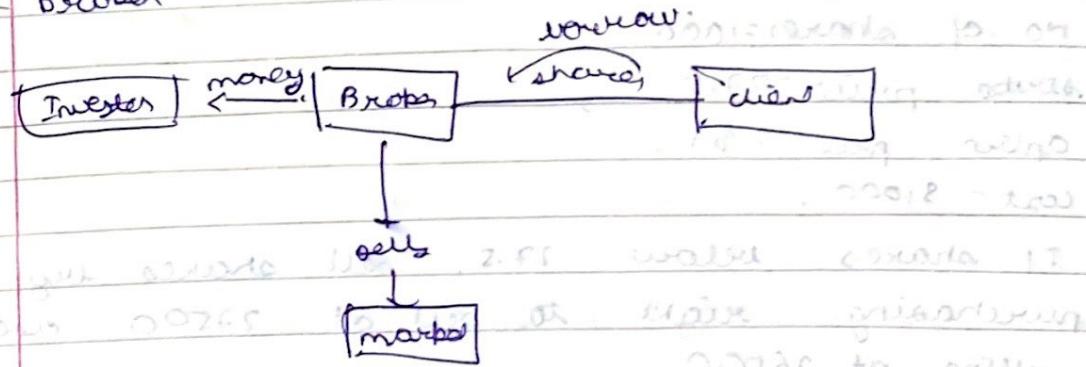
this is why parallel algorithms are needed.

The landscape from Parallel computing Research:
A view from Berkeley



$$\text{observed speedup} = \frac{\text{serial time}}{\text{parallel time}}$$

- stock options
- call option: right to buy.
- put option: right to sell.
- buying 100 shares of stock and selling it in London gives \$300 profit.
- short selling: involving an asset not owned.
- Broker



- investor buys shares, pay an income like interest to broker and then to bank.
- when investor earns enough money, the money is returned with interest to client.

now investor is owner of the shares.

Investor Broker Cient

- Investor borrows shares from client
- Broker sells and gives money to investor
- Investor buys same volume of shares
- Give proceedings to client via broker

Now, buying / selling of stocks:-

- Buy at \$60 after borrowing from bank
- After 1 year, give \$63 back to bank and sell at \$66 to make \$3 profit
- Now also, sell stock at \$60, give in bank buy at \$58, earn \$63, make \$5 profit
- forward agreement: agreement to buy or sell an asset at a certain future price at a certain price

call option: Right to buy at a fixed price

put option: Right to sell at a fixed price

- Price of 1 share = \$28.

No. of shares = 1000.

strike price = \$27.5

option price = \$1.

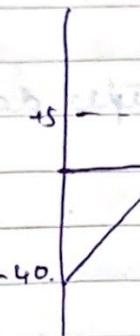
cost = \$1000.

If shares below 27.5, sell shares by purchasing right to sell at 27500 and selling at 26500.

- If share above 27.5, no need for put option

27) $k=40$, $x_0 = 5$. option cost = 5.

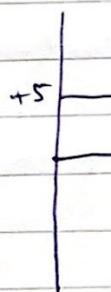
selling means short on market from



$$\text{Loss} = 40$$

$$\text{Profit} = 5.$$

If we have a call option:



$$\text{Max profit} = 5$$

$$\text{Max loss} = \infty.$$

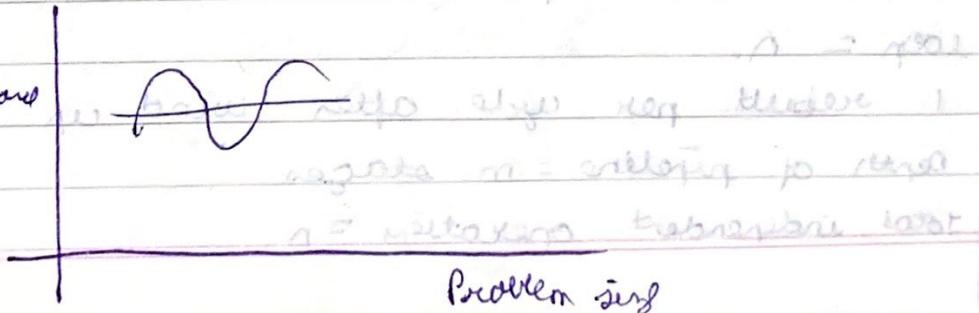
- 28) HPC - ground of 2. mod. 2 - volume too much

peak performance \rightarrow for benchmarking, setting some standard

- write a code in such a way so we know about its performance.

in order to measure it, $n = m = 1000$

peak
Performance

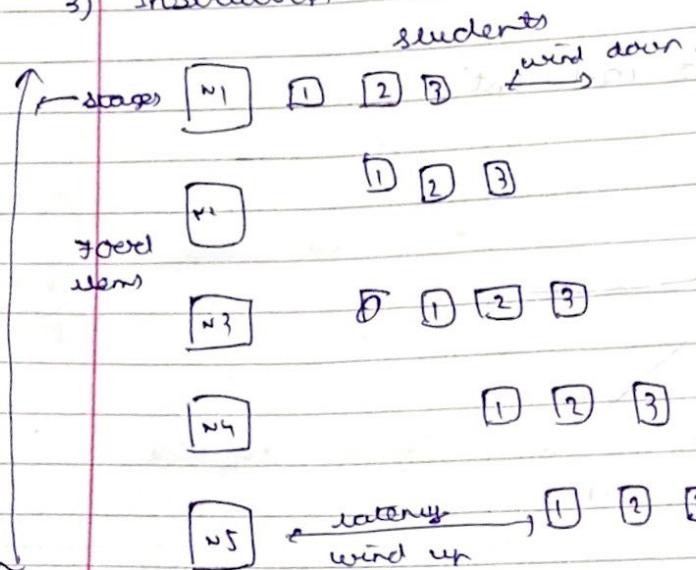


low level benchmarks - program to understand
chief performance characteristics of a system.

+ Multiplicity:-

- o we have multiple memory, multiple database
and multiple processors.

- 1) node
- 2) threads run on cores.
- 3) instruction



- when we complete 5 components, we do an operation
- ensure that all 5 stages are efficient, take same amount of time.
- one operation per cycle.
- Depth or latency of the pipeline.
- Depth = m, no. of students are also in a loop = n.
- 1 result per cycle after wind up phase.
- Depth of pipeline = m stages
- Total independent operations = n.

Time of pipeline to finish = $m+n-1$

serial time = mn .

$$\text{speedup} = \frac{\text{serial time}}{\text{parallel time}}$$

$$= \frac{nm}{n+m-1}$$

conclusion: If $n > m$, we will get good speedup only if loop is large compared to no. of stages. (loop \rightarrow total no. of instructions).

- we have an optimal no. of stages which has to be figured out.

$$\text{throughput of pipeline} = \frac{n}{n+m-1}$$

$$\text{Dividing by } n, \text{ we get: } \frac{1}{1+\frac{m-1}{n}}$$

The maximum throughput is 1 when $n < m$.

- von neumann architecture \rightarrow main memory, CPU, interconnects
- most important limitation is data access.

L1 cache = in kilobytes.

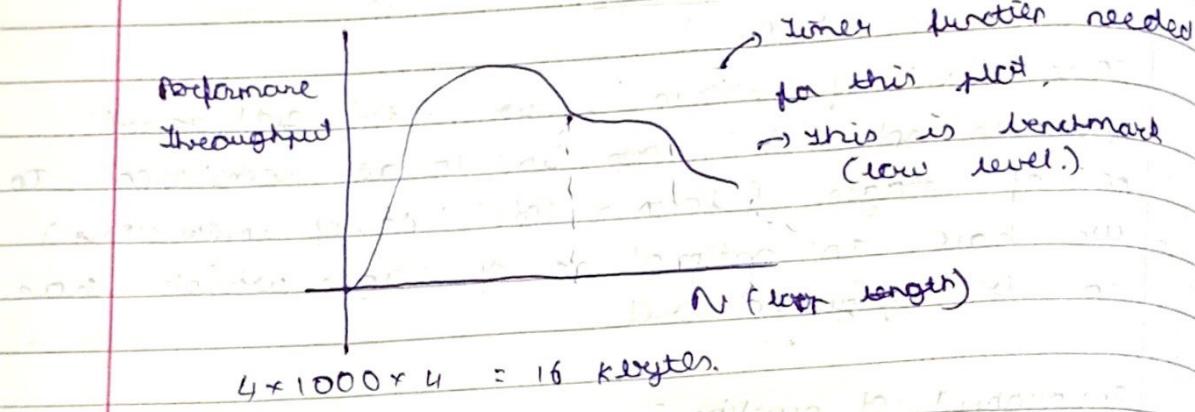
L2/L3/main memory = in megabytes.

latency: - amount of time required to access 0 bytes of data.

- von neumann bottleneck

Microbenchmarking:-

- separate influences, for
- determine specific machine capabilities.



Machine balance = Memory bandwidth GB/sec.
Peak Performance.

Balance analysis: Theoretical performance of even used codes.

Machine balance

code balance.

Machine balance will decrease in future because no. of cores will increase and core performance will increase.

Code balance of even = data traffic

$\frac{\text{no. of floating point operations}}{\text{no. of variables used}}$

Data traffic = no. of variables used.

Computational intensity = $\frac{\text{work}}{\text{code balance}}$

Data → load, store, execution of operation.

Data traffic → performance limiting data path.

- took up: wall clock timer, $a = \text{gettime_of_day}()$
- $\text{get time of day}()$ written, $b = \text{get_time_of_day}()$
- Each clock has a resolution, $c = a - b / \text{resolution}$
- Be careful about error, $c = a - b$
- Plot performance vs no. of loops
- the graph is used as a low level benchmark

- principle of locality - spatial and temporal locality.
- Memory access will be on blocks of items instead of individual items, these are called data blocks.
 - Cache is k times faster than RAM. (latency + bandwidth).
 - Cache reuse ratio = r . (fraction).
 - When we are working on data, can we use the data such that we have a, b, c ?
 - Every time we access something else, understand cache hits or miss may take longer.
 - $a[i] = b[i] + c[i]$, if next row, write down we need to reuse $b[i], c[i], a[i]$, we may have to reuse cache before throwing it away.
 - fraction of load and store reused from cache.
 - Access to main memory time = T_m .
 - Cache access time = $T_c = \frac{T_m}{k}$.
 - Average access time = $T_{av} = rT_c + (1-r)T_m$.

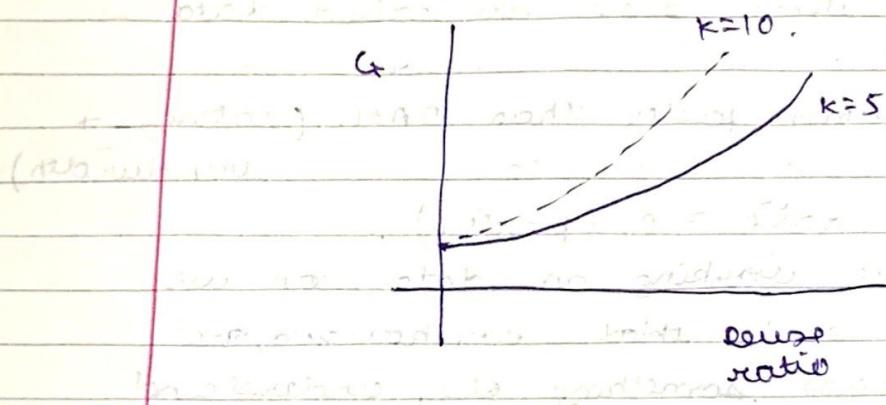
$$\text{Performance gain, } G = \frac{T_m}{T_{av}}$$

$$\text{given, } G = \frac{kT_c}{rT_c + (1-r)T_m}$$

$$G = \frac{K}{r + (1-r)k} \rightarrow \text{formula for gain.}$$

reduces miss rate \rightarrow more reuse \rightarrow less misses \rightarrow less time

- If there is huge reuse, i.e. in the case of a fast cache, the gain G will be much faster.
- RAM will sit outside the processor.
- If $r=0$, no gain.
- write code in such a way, the reuse ratio is greater.



- We called data from RAM to cache, the next time, we keep on ~~do~~ taking data from the cache.
- If no temporal locality, we use spatial locality we do not bring 1 elements, we bring a whole lot of elements together.

$A[1] \rightarrow$ not only used, use all elements.

$A[1] \dots A[16]$

1 --- 16.

if \rightarrow bring everything
at once is better

latency :- the time for 0 bytes of data, latency best is same, so latency per element increases

- we will increase the cache hit ratio.
- whenever there is spatial locality, the problem of latency is solved.

cache line length = 16.

spatial locality fixes the hit ratio at?

hit ratio = $\frac{15}{16}$ (in the first time, you have to use latency, because nothing very good ↪ was in cache, but for the rest, hit ratio. there is a hit in the cache.)

- parallel algorithm :- no one unique way of doing something.
- design is on how the data is accessed in
 - cache
- efficient parallelization is not parallelizing each step, but to design a new algorithm.
- many numbers added simultaneously but in parallel's case: parallelism with dependency

$$\text{speedup} = \frac{1}{1-p}$$

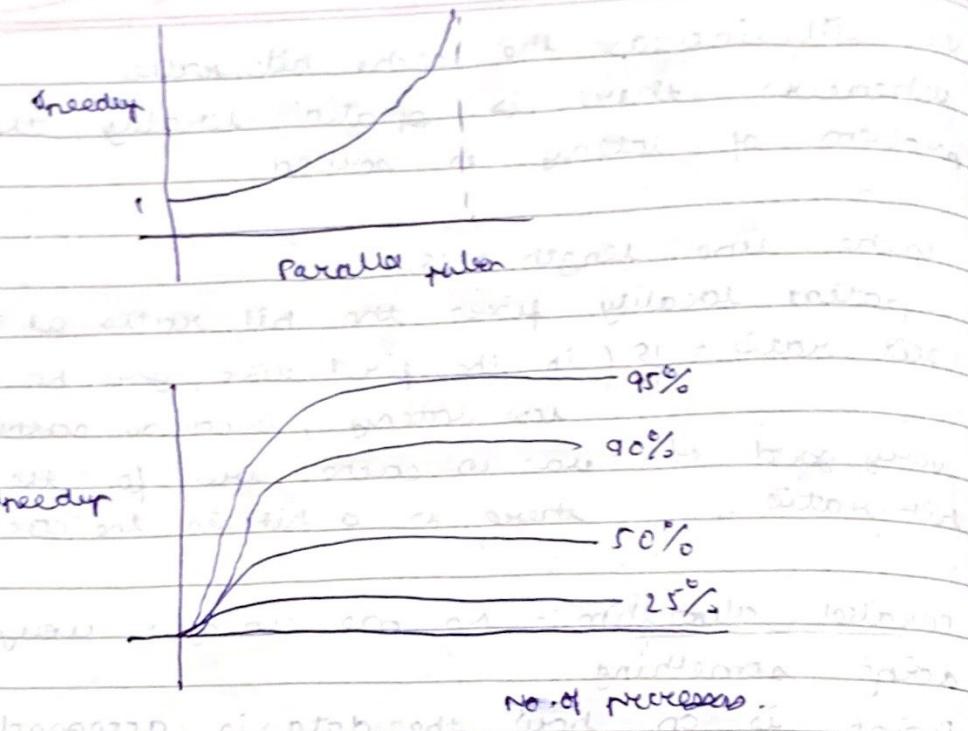
p = fraction of code that can be parallelized.

$p=0 \rightarrow$ no speedup.

If there is ∞ dependency, everything cannot be parallel.

$$p=0.5, \text{ speedup} = \frac{1}{1-0.5} = \frac{1}{0.5} = 2.1$$

Theoretical maximum speedup = ∞ .



- optimizing code's Performance modelling
- Profiling - serial code is not the same as parallel code
- Find subroutine taking maximum amount of time and optimize that subroutine (that particular part of the code).

compute : f95 -pg file.f77

file : gmon.out has to be saved by

command : gprof gmon.out > gmon.out

or gprof gmon.out > gmon.out

$a = a + b * b$. following all terms

$a = a + b * b$. equation 2.0

both will give the same answer.

$$P + 2.0 \rightarrow e^{2lnB}$$

This optimization technique is called strength reduction.

- A lookup table can be used.

case A:-

do $i=1, N$

$$A(i) = A(i) + s + r * \sin(x)$$

end do

case B:-

$$tmp = s + r + \sin(x)$$

do $i=1, N$

$$A(i) = A(i) + tmp$$

end do

Due to cache hit rate, one process will be much faster than the other.
 $\beta = 2.0 \rightarrow \exp(2 \ln \beta) \rightarrow$ very heavy function as
 more no. of pipeline stages

This is called local iteration reorganization to make it faster.

- whenever we have a conditional branching, we should try to avoid it.
- the problem is branch misses (like cache miss), if branch miss happens, more time is taken.
- If there is something false, the pipeline has to be flushed, this will take more time, based on a probability model.

as many as
branching

following branches
exist

clear branches
(.moreso & less)

a) do $i=1, N$
 do $j=1, N$
 $c(i) = c(i) + a(i,j) * b(j)$
 $c(i+1) = c(i+1) + a(j,i+1) * b(j)$
 enddo
 enddo.

b) do $i=1, N$
 do $j=1, N$
 $c(i) = c(i) + a(j,i) * b(j)$
 enddo
 enddo.

In the first case cache hit is taking place, but for the second case, cache hit is not taking place, this is very unrolling.

Parallel Programming: scope of parallelism, granularity, locality, load balance, coordination and synchronization.

- Parallel programming is harder than sequential programming.

Organization of parallel problems.

→ Using \downarrow logical \downarrow physical
 programmer view. \downarrow architecture

expressing parallel task. \downarrow Mechanism for communication.

several models.
 (context switcher.)

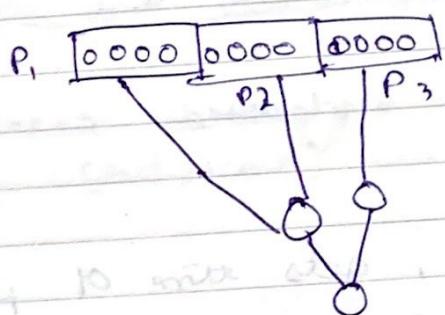
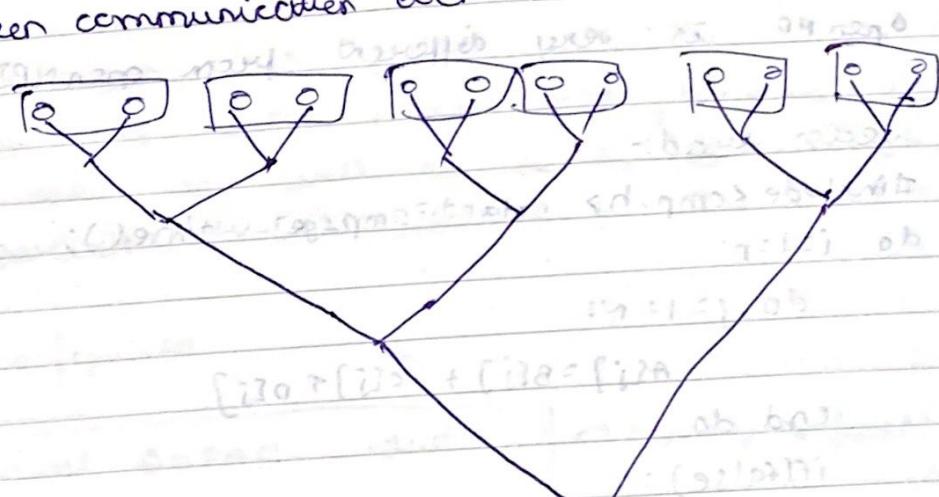
* shared memory parallelism:

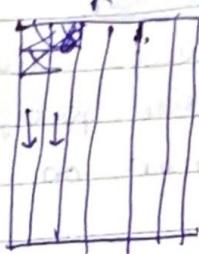
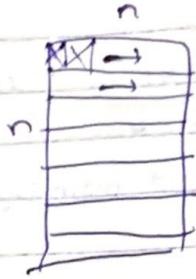
- A set of compiler devices and library routines for parallel application programmers.
- Performance improvement on shared memory machines.
- use OpenMP.
- Multi vendor support in OpenMP.
- It is portable.

granularity: Ratio of computation to communication.

coarse - large amount of computational work between communication events.

Fine - small amounts of computational work between communication events.





we have n processes. processes

4 processes.

$n > 4$.

Finest level of granularity = n^2

Highest level of granularity = $n \times n$

Logical level will decide level of granularity.

OpenMP is very different from openMPI.

vector triad:-

```
#include <omp.h>. start = omp_get_wtime();
```

```
do i=1:r;
```

```
do j=1:n!
```

$$A[ij] = B[ij] + C[ij] * D[ij].$$

```
end do.
```

```
if(false):
```

$$\text{call dummy}(A,B).$$

```
end do.
```

```
end do.
```

```
total = omp_get_wtime
```

wtime(): world time, gets time of process,
not the system.

- 1) compile: gcc filename.c -fopenmp -o filename-pg.

2) Run: ./filename.out

3) gprof -b filename.out gmon.out

timer function

elapsed time

resolution of clock.

Right now: stages are small, comparable to stages of pipeline.

[CPU]

Registers

L1 cache

L2 cache

- when L2 cache comes into play, L1 has to take as well as send, not only send, so performance will reduce.

2) cat /proc/cpuinfo

64 bit DDR4 RAM

8 channels

16Hz clock frequency

Memory system

Access throughput = 8 bytes \times 8 channels
(in bytes) second

transfer

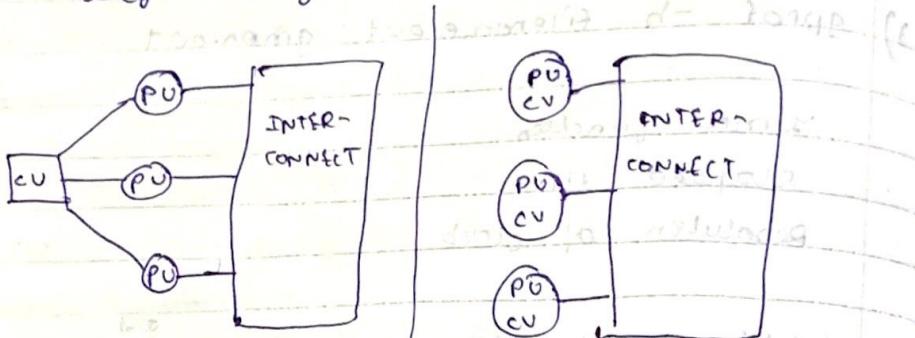
$$\times \frac{2 \text{ transfers}}{\text{clock cycle}} \times 10^9 \text{ clock cycles} \text{ per second}$$

$$= 128 \text{ GB/second}$$

- we already can access shared cache, but not private cache

SIMD - single instruction multiple data

SPMD - single program multiple data



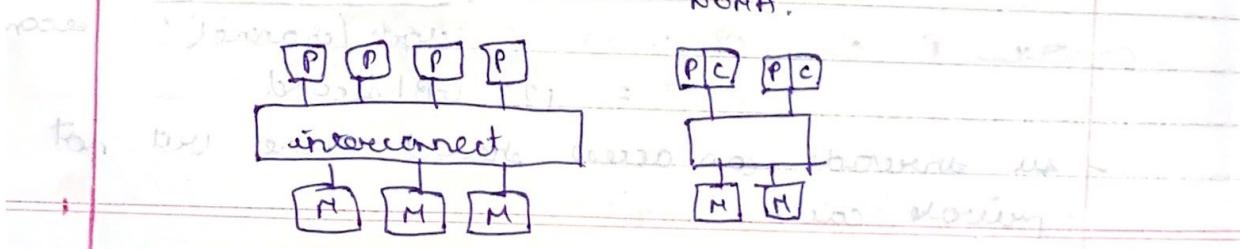
Relatives between SIMD and SPMD to work

- Control unit gives a
- single instruction, executed on multiple processes.
- Speed will depend on the program.
- For every clock cycle, one instruction is executed.
- When we have a global memory, it is very easy to read, but very difficult to write, so there comes the concept of private memory.
- In parallel programming, we deal with global and private memory.

UMA - Uniform Memory Access

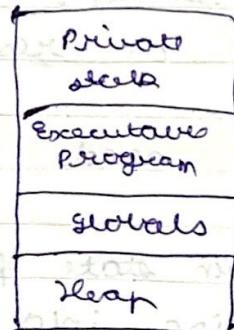
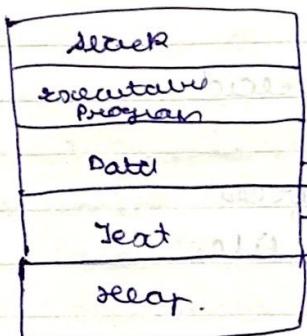
NUMA - Non-Uniform Memory Access

2.2.2. ~~Architecture~~ (X86)



- Each process will take same amount of time, in first diagram, but in other diagrams, different amount of time is taken.
- global / private variables

<u>process</u>	<u>thread</u>
- It is a running program.	- It is a lightweight process.
- All processes have a process ID.	
- All processes have a status.	



several stacks

- Thread is the smallest unit of execution of instructions.
 - Process is an execution stream.
- Ready good for new system
 Running good for real time
 Waiting good for multitasking
 Terminated good for time keeping
- swapping is possible between the intermediate stages.
 - moment to moment threads are related, so can be assigned to them.

threads at once and yet priority

Thread includes :-

- 1) Thread ID
- 2) PC
- 3) Register set, set (pointer)
- 4) Stack. (pointer)

Shared resources:-

- 1) Stack session
- 2) Data session.
- 3) Other OS resources.

The OS needs to keep track of the :-

- 1) Execution state for each thread.
- 2) Scheduling information
- 3) Memory used by the process.
- 4) Information about open files

- The OS will use a process control block.
- Scheduling has overhead in form of a few clock cycles elapsed.

SIMD - one control unit gives instruction to multiple data, very easy.

SPMD - scheduling time is more as many control units give multiple instructions.

Mutual exclusion, when we write to a file, we will not allow others to write.

- writing by multiple users is sequential,

reading can be done in parallel.

Dependencies:-

- for $i=1, n$
- for $j=2, m$
- $b[i, j] = \dots + b[i, j-1]$
- for different values of i 's we can launch threads, but this is not possible for different values of j as there is a dependency on value of j .

4 types of dependencies:-

flow, anti, input, output

flow dependency

- statement ' i ' precedes ' j ' and ' i ' captures a value that ' j ' uses.

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

Dependencies:-

$$\textcircled{1} \rightarrow \textcircled{2}$$

$$\textcircled{1} \rightarrow \textcircled{3}$$

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

- anti dependency: after statement ' i ' precedes ' j ' and ' i ' captures a value that ' j ' computes.

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

Dependencies:-

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

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

3) Output dependency:

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

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

Dependencies:

$$\textcircled{1} \rightarrow \textcircled{3}$$

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

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

4) Input dependency:

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

Dependencies:

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

Dependencies: (Input dependency not a problem)

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

- the most challenging problem is how to remove the dependencies.

- It is very difficult to remove true dependency.

Dependency flows from 'i' to 'j'

↳ sink

source ↳ sink

- we can have data and task dependency
use DAG, directed graph.