

Computer Architecture & Organization: Implementation of Computer Architecture

1980s machine

1990s - more complex & time sharing.

2000s - 1990s computer

2000s - Cell phone

had need for embedded systems & nuclear powered

	Desktop	Server	Embedded system
Price	\$500 - \$1000	\$1000 - \$10000	\$10 - \$100
Power	100W - 1000W	100W - 1000W	10W - 100W
Reliability	High	High	High
Availability	High	High	High
Performance	High	High	High
Flexibility	High	High	High
Scalability	High	High	High

Measuring Relative Summarizing the performance of computer

→ Program time & the time for start & completion of particular event

→ Throughput → Amount of work done per unit time

→ Execution time → Rate of completion of particular task

Example: X is faster than Y

It is execution time of Y

Execution time of X

not discussed for system

Since $Performance = \frac{1}{\text{exec time}}$

$$n = \frac{\text{performance of } X}{\text{performance of } Y}$$

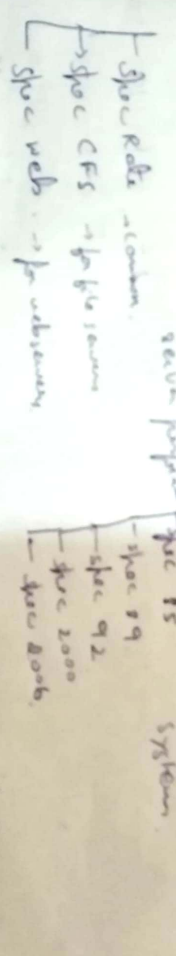
Ex TS

Benchmark

- 1) Real program : by user
- 2) Kernel \rightarrow Key piece from real application
- 3) Toy Benchmark \rightarrow as realistic, assigned program (random)
- 4) Synthetic Benchmark \rightarrow for program with known data, the load of program, the number of real application

1) Desktop Benchmark \rightarrow to measure performance of desktop

2) Server Benchmark



3) TP Benchmarks \rightarrow TPC
Computing system
Transaction Processing

- TPC
- \rightarrow TPC
- \rightarrow TPC IoT \rightarrow Big data

Standard for performance evaluation comparison (spec)

The couple A is 10 times faster than B for program 1.
 B is 10 times faster than A for program 2.
 A is 20 times faster than C for program 1.
 C is 50 times faster than A for program 2.
 B is 2 times faster than C for program 1.
 Identify the execution time for each of the system for
 from 142 & also identify the total execution time for
 system A B & C

→ p1

$A = 10B \quad B = B$

$A = 20C \quad C = C$

$B = 2C \quad C = C$

p2

p1	A B C		
	10	1	10/20
p2	1	10	50

p1	A B C		
	1	10	20
10	$\frac{20(50)}{10}$	$\frac{20(10)}{10}$	$\frac{10}{50}$
11	10	10	20
p2	1000	100	20
Total time	1001	110	4

Summarizing Performance

eg. If spec ratio is 50 and implies computer under test is 50 times better than reference computer

$$\text{Spec Ratio} = \frac{\text{Execution time of reference computer}}{\text{Execution time on computer under test}}$$

Spec Ratio of Computer A on Benchmark was 1.25 times higher than Computer B - so represent this defn in terms of performance.

$$\left[\begin{array}{l} \frac{\text{perf of A}}{\text{perf of B}} = \frac{1}{1.25} = \frac{1}{1.25} = \frac{1}{1 + \frac{1}{4}} = \frac{4}{5} \\ \text{Spec Ratio of A} \\ \text{Spec Ratio of B} = 1.25 \end{array} \right] \Rightarrow \frac{\text{Spec Ratio of A}}{\text{Spec Ratio of B}} = 1.25$$

$$\frac{\text{Spec Ratio of A}}{\text{Spec Ratio of B}} = 1.25 = \frac{\frac{\text{exe of Ref}}{\text{exe of A}}}{\frac{\text{exe of Ref}}{\text{exe of B}}} = \frac{\text{exe of B}}{\text{exe of A}} = \frac{\text{perf of A}}{\text{perf of B}} = 1.25$$

$$\boxed{\frac{\text{Spec of A}}{\text{Spec of B}} = \frac{\text{perf of A}}{\text{perf of B}} = \frac{\text{exe of B}}{\text{exe of A}}}$$

$$\text{Geometric Mean} = \sqrt[n]{\prod_{i=1}^n \text{sample}_i} \quad \begin{array}{l} 2.4 \\ \sqrt{2(4)} \\ \sqrt{8} = 2\sqrt{2} \end{array}$$

Q2) Show that the Ratio G.M is equal to G.M of performance Ratio works computer A & B

⇒

$$\frac{A \cdot B}{AB} = \frac{B}{A}$$

Geometric Mean of A

$$\sqrt[n]{\prod_{i=1}^n \text{Spec Rate } A_i}$$

$$= \sqrt[n]{\frac{\sum_{i=1}^n \text{Spec Rate } A_i}{n}}$$

$$= \sqrt[n]{\frac{\sum_{i=1}^n \text{excellent Ref}}{n}}$$

$$= \sqrt[n]{\frac{\sum_{i=1}^n \text{excellent Ref } B_i}{n}}$$

$$= \sqrt[n]{\frac{\sum_{i=1}^n \text{performance } A_i}{n}}$$

↑
geometric mean of performance

Quantitative principles of computer design

3 principles

- 1) Take the advantage of parallelism
- 2) Principle of locality → shared resource
- 3) Focus on core sense of can

* Rules of thumb

90% general
only 10% specific

(Optimise a can not be properly used)

Andels law

→ performance improved to be gained from using some faster mode execution is limited by fraction of the time the faster mode can be used.

$$\text{speedup} = \frac{\text{Performance of entire task using enhancement}}{\text{Performance of entire task without using enhancement}}$$

$$= \frac{1}{\text{Execution time using enhancement}}$$

Execution time of entire task without using enhancement.

$$= \frac{\text{Execution time of entire task without using enhancement}}{\text{Execution of entire task using enhancement}}$$

$$\text{speedup overall} = \frac{\text{Execution time old}}{\text{Execution time new}}$$

$$= \frac{1}{\left(1 - \text{Fraction enhanced}\right) + \frac{\text{Fraction enhanced}}{\text{speedup enhanced}}}$$

Qualitative principles of Design

2 factors Andels. 1. given a 2 ways for some enhancement.

→ the fraction of the computation time in the original code that can be converted.

Fraction enhanced is less than 1 //

speed of enhanced. greater than 1
2nd factor...

$$Execution\ Time_{New} (ET) = ET_{old} \left((1 - F_{io/enhanced}) + \frac{F_{enhanced}}{Speedup_{enhanced}} \right)$$

Overall Speedup Ratio

$$Speedup_{overall} = \frac{ET_{old}}{ET_{new}} = \frac{1}{(1 - F_{enhanced}) + \frac{F_{enhanced}}{Speedup_{enhanced}}}$$

Q. Although
ref. factors
→ means

a) Suppose that we want to enhance the processor used for web serving. the ~~new~~ processor is 10 times faster.

Assuming original process is busy with computation - 40% of the time & is waiting for I/O 60% of time. what is the overall speedup gained by enhancing the processor.

⇒ processor enhanced → amount of time computer is doing computation.

$$\frac{1}{(1 - 10) + \frac{10}{10}}$$

$$\frac{1}{1 - 0.4 + \frac{0.4}{10}}$$

$$= \frac{1}{\left(1 - \frac{40}{100}\right) + \frac{\frac{40}{100}}{10}}$$

$$= \frac{1}{(1 - 0.4) + \frac{0.4}{10}}$$

$$= \frac{1}{0.6 + 0.04}$$

$$Speedup_{overall} = \frac{1}{0.64} = \frac{100}{64} = \frac{25}{16} = 1.56$$

b) ~~exer~~ 2

A. Conner transform a record as graph, processor & sensor root. Implementation of floating point (FP) is more root vary significantly in ~~for~~ performance, especially across processor designed for graphics. Unlike FP speed = $\frac{100}{20} = 5$ sensor root (FPS) is replaced for 20% of execution of critical graphed bandwidth. A processor is about 20% better had a speed up factor of 10.

$$= \frac{1}{(1-0.2)} + \frac{0.2}{10}$$

$$= \frac{1}{0.8 + 0.02} = \frac{1}{0.82} = 1.22$$

2nd a. $\text{factor} = \frac{1}{2}$ of execution time = 0.5 speed up = 1.6

$$= \frac{1}{(1-0.5)} + \frac{0.5}{1.6}$$

$$= \frac{1}{1-0.5} + \frac{0.5}{1.6} = \frac{1}{0.5 + 0.3125} = 1.23$$

↑ better

Means, improving the performance of the FP operation overall slightly better because of higher freq.

Q) Derive the CPU performance equation.

→ we do it in terms of clock cycle

$$\text{CPU time} = \text{CPU clock cycles for program} \times \text{clock cycle time} \quad \text{--- (1)}$$

$$\text{CPU time} = \text{CPU clock cycles for program} \times \frac{1}{\text{clock rate}}$$

Number of instructions executed in the instruction path
 length = instruction count
 clock cycle per instruction CPE = $\frac{\text{clock cycle}}{\text{instruction count}}$
 CPE, CPU clock cycle for a program — (2)

from (1) & (2)

$$\text{CPI} \times \text{Instruction count} = \frac{\text{CPU time}}{\text{CPU clock cycle time}}$$

$$\text{CPU time} = \text{CPE} \times \text{IC} \times \text{CPU clock cycle time}$$

where

$$\begin{aligned} \text{CPU time} &= \frac{\text{clock cycle}}{\text{Instruction}} \times \frac{\text{Instruction}}{\text{program}} \times \frac{\text{seconds}}{\text{clock cycle}} \\ &= \frac{\text{seconds}}{\text{program}} = \text{CPU time} \end{aligned}$$

$$\text{CPU clock cycles} = \sum_{i=1}^n \text{IC}_i \times \text{CPI}_i \quad \text{--- (3)}$$

where

IC_i represent number of times instruction i is executed in a program

CPI_i represents the average number of clock per instruction for instruction i

(3) RD

$$\text{CPU time} = \left(\sum_{i=1}^n \text{IC}_i \times \text{CPI}_i \right) \times \text{clock cycle time}$$

The overall CPI can be computed now using eqn.

(3) + (2)
over all CPI

$$CPI = \frac{\sum_{i=1}^n IC_i \times CPI_i}{\text{Instruction count}}$$

$$CPI = \sum_{i=1}^n \frac{IC_i}{\text{Instruction count}} \times CPI_i$$

no of times the instruction is fetched

Example 3.

→ suppose we have made the following measurement

freq of FP operation is 25%.

Avg CPI of FP operation is 4.

Avg CPI of our instr is 1.33.

Freq of FPSQR = 2%.

CPI of FPSQR is 20.

→ Assume that 2 design alternatives are given CPI of FPSQR to 2 or decrease the avg CPI of all FP op to 2.5. Compare the 2 design alternatives using the performance eqn.

$$\begin{aligned} CPI_{\text{overall}} &= \sum CPI_i \times \frac{IC_i}{\text{Instruction count}} \\ &= \text{FP} \quad \text{other} \\ &= 4 \times 25\% + 1.33 \times 75\% \\ &= \frac{4 \times 25}{100} + \frac{1.33 \times 75}{100} \\ &= 1 + \frac{3.99}{4} = 0.99 \\ &= \underline{\underline{2.0}} \end{aligned}$$

$$\begin{aligned} CPI &= \frac{25 \times 4}{100} + \frac{2 \times 20}{100} \\ &= 1 + 0.4 \\ &= \underline{\underline{1.4}} \\ &= \frac{25(2.5)}{100} + \frac{2(40)}{100} \end{aligned}$$

25/100

Compute the CPI for the enhanced FP50K by substituting the cycles saved for the original CPI.

$$\begin{aligned} \text{CPI}_{\text{with new FP50K}} &= \text{CPI}_{\text{original}} - 2\% \cdot (\text{CPI}_{\text{old FP50K}} - \text{CPI}_{\text{new FP}}) \\ &= 2 - 0.02(20 - 2) \\ &= 2 - 0.02(18) \\ &= 2 - 0.36 \\ &= \underline{\underline{1.64}} \end{aligned}$$

$$\begin{aligned} \text{CPI}_{\text{new FP}} &= \text{CPI}_{\text{original}} - 25\% \cdot (\text{CPI}_{\text{old FP}} - \text{CPI}_{\text{new FP}}) \\ &= 2 - 25\% \cdot (4 - 2.5) \\ &= 2 - \frac{25}{100} (1.5) \\ &= 2 - \frac{1}{4} (1.5) \quad 30 \quad 20 \\ &= 2 - 0.375 \\ &= \underline{\underline{1.625}} \end{aligned}$$

Now if CPI is low the machine is better.

FP Performance is marginally better with FP enhancement, since CPI is slightly lower.

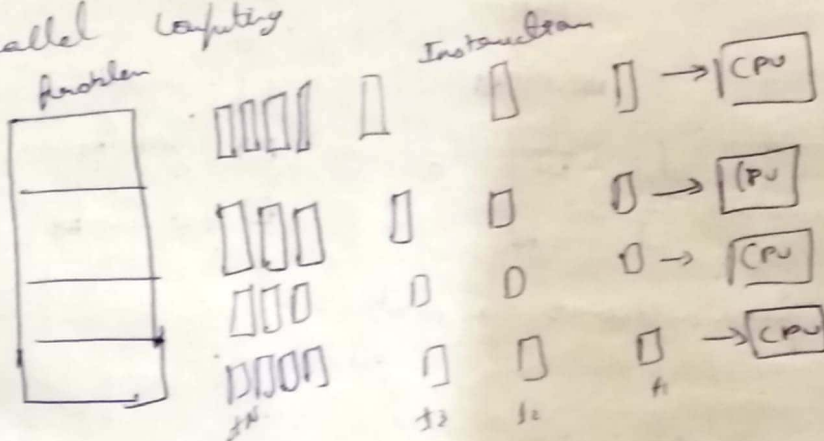
$$\begin{aligned} \text{Speedup w.r. to new FP} &= \frac{\text{CPU time original}}{\text{CPU new FP}} \\ &= \frac{2}{1.625} = \underline{\underline{1.23}} \end{aligned}$$

Introduction to parallel processing

→ simultaneous execution multiple instructions of a computational task

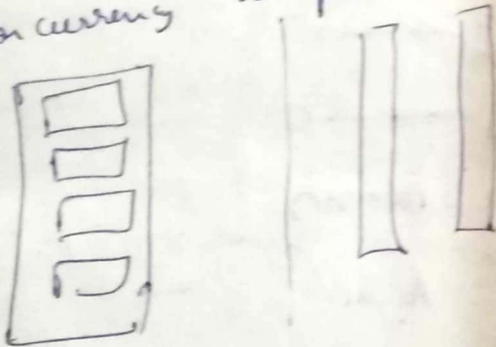
style → thread level → multiprocessor → concurrent execution
Instruction level → pipelined → one after other

Parallel computing



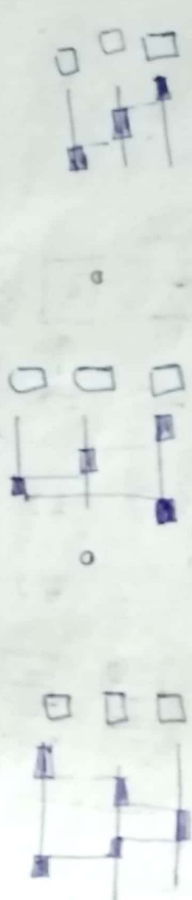
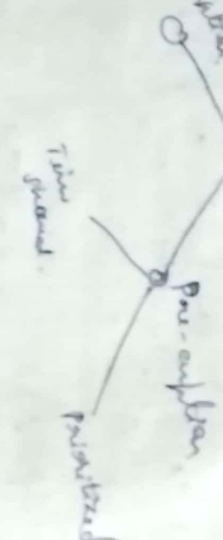
data level parallelism
 procedure level parallelism

concurrency vs parallelism



the multiple client & 1 server
 → to apply all client we have to apply Pre-emption rule.

Pre-emption
 Non pre-emption

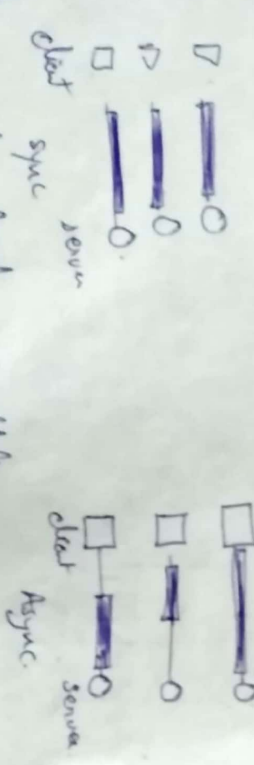


In most cases pre-emption is a must.

Parallel execution

N client N server model
 synchronous or Asynchronous

Start of
 same time



* Types and level of parallelism
 available → in program → make waiting lot we need to do it
 utilized → during execution → parallel & order-level

Types of available parallelism

- Data parallelism
→ same for all the data
 $a(1) + b(1) + c(1)$
- Functional parallelism
→ separates tasks into different operations to different levels

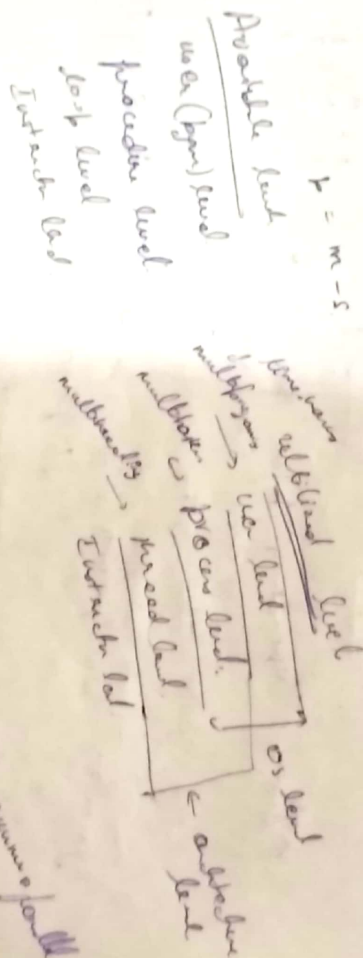
$$a \leftarrow 2$$

$$b \leftarrow 3$$

$$m = (a+b)/c$$

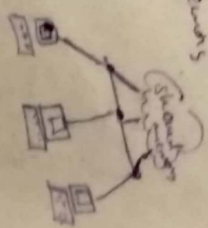
$$s = (a^2 + b^2)^{1/2}$$

$$p = m - s$$

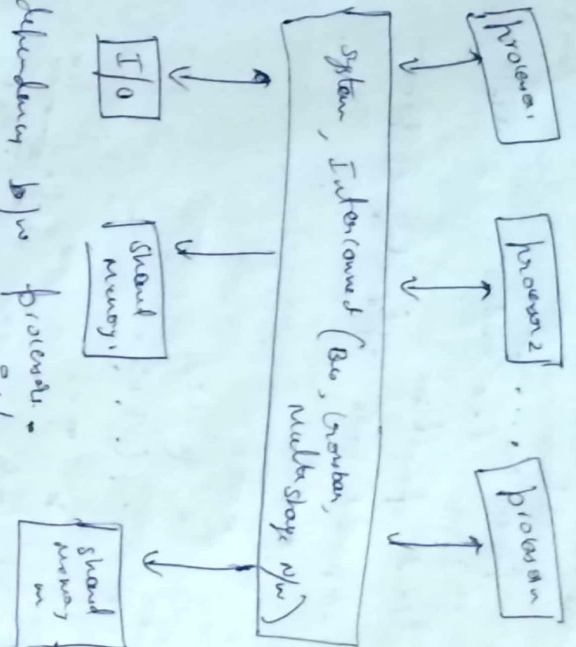


Shared Memory Multiprocessors

- Uniform memory access (UMA)
- Non-uniform memory access (NUMA)
- Cache-only memory architecture (COMA)



Highly coupled process dependencies



There is dependency b/w processes -
to access same variable

→ tightly coupled system

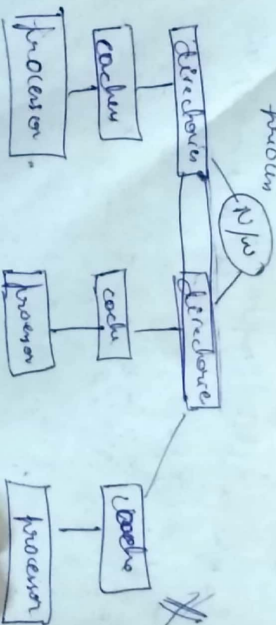
symmetric multiprocessing → process have equal access to all hardware

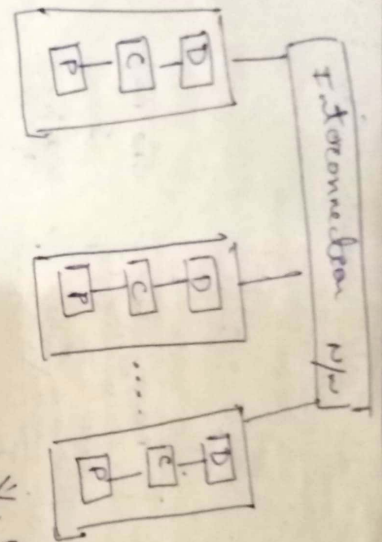
Asymmetric multiprocessing → Only 1 or subset processes can have I/O or executive capability
of master slave

Remaining process (slave) don't have I/O capability
(Attached processor)

NUMA → shared memory is divided into small memory or local processor distributed into clusters.

COMP

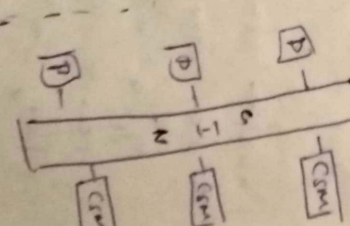
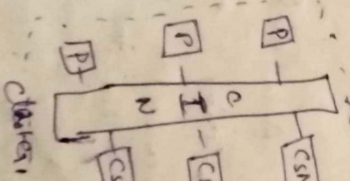
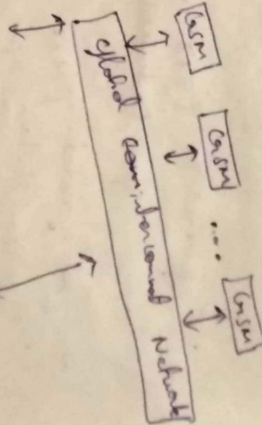
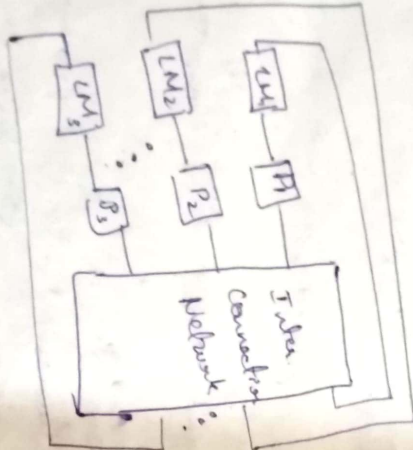




D → data
C → cache
P → processor

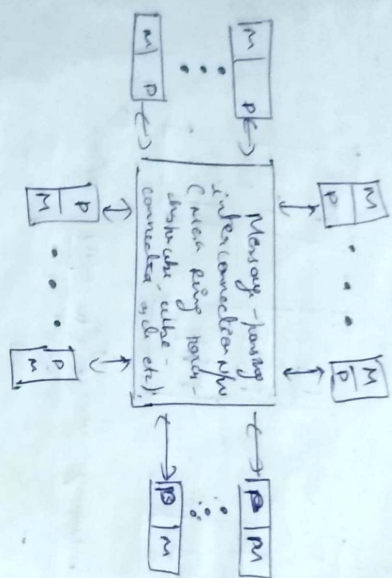
Non uniform memory access (NUMA)

~ 100%



CEM → cluster shared memory
NEM → global network memory
CEM → cluster shared memory
P → processor

distributed Memory Matrix Computer (Grand alex. 2)



1st generation (1933-1987) as one hand on fluoroscopy technology

→ 2nd generation (1988-1992) were influenced with much + computer
 architecture, w/ more memory, growing
 hand on for grain multi
 3rd gen (1993-1997) were

3rd year (1993-1997)

Problem lead a
Andell, low
a enhancement for the known

Problem load a more
→ we can consider an enhancement to the program
of a web search. The new CPU is 20 times faster in
search queries than the old program. The old program
busy with search queries 70% of the time, and is
idle spending time waiting for input/output.

$$sk_{old} = \frac{1 - 0.7}{0.7}$$

$$\frac{1}{(0.3) + 0.035} = \frac{1}{0.335} = 2.985$$