

Date _____

Day _____

CHAPTER # 02

▷ VON NEUMANN ARCHITECTURE :

→ The classical von neuman architecture consist of

① Main memory ② CPU ③ Processor / Core ④ Interconnect^{on} _{btw mem & CPU.}

↳ Main Memory consist of collection of locations used for storing instruction & data

↳ CPU is divided in control unit & ALU.

CU → decides which instr. to execute thru a program counter registers.

ALU → responsible for executing actual instr.

CPU contains register which contains state of a program

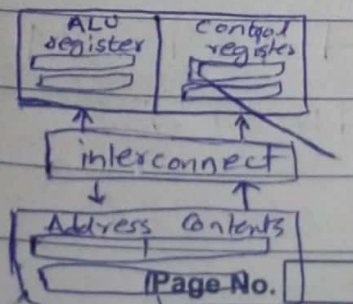
↳ Bus acts as an interconnection btw mem & CPU

→ A von neuman machine executes a single instr. at a time & each instr operates on only a few pieces of data.

→ Transfer of data from mem to CPU is called fetch or read.

→ Transfer of data from CPU to memory is called written to mem or stored.

CPU



Date _____

Day _____

→ The separation of mem and CPU is often called von Neuman bottleneck, since the interconnect determines the rate at which instr & data can be accessed. Therefore the rate of data transfer through interconnected bus is way slower than the CPU's processing speed. CPU

Solution:

- ① Using cache memory which are fast & close to CPU.
- ② Pipelining i.e. to break down inst into stages.
- ③ Parallelism i.e. to use multiple processing

"THE BASICS OF CACHING"

Cache is a collection of memory locations that can be accessed in less time than some other memory locations.

↳ The principle that an access of one location is followed by an access of a nearby loc. is called locality.

↳ Nearby location → Spatial locality

↳ Location to be used in near future → Temporal locality

Cache Block / Lines: Memory access will effectively operate on blocks of data and instr instead of single data item / inst.

Levels of Cache: Cache is usually divided into levels.

① L1, L2, L3,

L1 being the smallest & the fastest.

Date _____

Day _____

Multilevel caches don't duplicate information that's available in another level. When data is needed by CPU, it searches from top to bottom & then in main mem.

→ If the data is found in the cache it's called cache hit or L1 hit or hit.

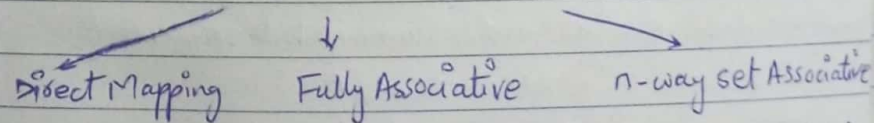
→ If not found then cache miss / L1 miss / miss.

→ If a value in cache is updated, the val in cache & main mem are inconsistent. Two approaches to deal with this:

↳ Write-through caches, the line is written to main mem when it is written to the cache.

↳ Write-back caches, the data isn't written immediately. Rather, the updated data is marked (dirty) and when the cache line is replaced by a new cache line from memory, the dirty line is written to memory.

"CACHE MAPPINGS"



① Direct Mapping: Is mein no. of hits karm or miss zyada hojate hai aur k sabse loc fixed hai to empty space zyada hota hai

→ Is mein hum har ek block of main memory ko ek fixed cache line assign kr dete hai. Assign kaise kr liye we use $K \bmod n$, where K = Block No. & n = line no.

Date _____

Day

Cache

Main Memory

↑ keel de hi category.

L_0	$B_0 \ B_4 \ B_8 \ B_{12} \ -B_{23}$	$W_0 \ W_1 \ W_2 \ W_3$	B_0
L_1	$B_1 \ B_5 \ B_9 \ B_{13} \ -B_{24}$	$W_4 \ W_5 \ W_6 \ W_7$	B_1
L_2	$B_2 \ B_6 \ B_{10} \ B_{14} \ -B_{30}$	$W_8 \ W_9 \ W_{10} \ W_{11}$	B_2
L_3	$B_3 \ B_7 \ B_{11} \ B_{15} \ -B_{31}$	$W_{12} \ W_{13} \ W_{14} \ W_{15}$	B_3
		\vdots	\vdots
		$W_{124} \ W_{125} \ W_{126} \ W_{127}$	B_{31}

16 words

line size = block size

16 words

Line Size = Block Size

$$\frac{16}{4} = 4$$

128 Words

Each block have 4 words

so total no. of blocks = $\frac{128}{4} = 32$

$$B_0 \rightarrow 0 \pmod{4} = 0$$

$$b_2 \rightarrow 2 \pmod{4} = 2$$

$$B_{31} \rightarrow 31 \pmod{4} = 3$$

- Ab jab hi B₁ cache mein jayega to wo L1 mein jayega. Sab blocks ki location fixed krdi hai.

- Finding Physical Address in Main Mem:

→ As we have total 128 words so no. of bits required to store the address of 128 words = 2^7

Block Number Block offset \Rightarrow Physical Add.

5 bits 2 bits

7 bits

0001010

2 2 (block size = 2)

Block no 2

$\frac{w}{scl}$ 42nd Indet

Page No.

Date _____

Day

- Finding Address in Cache

7 bit

Tag	Line No	Block offset
3 bits	2 bits	2 bits

ek line mei total 3 blocks
mein se koi bhi aasakta hai
To humen ki 3 bits with help use to
find block no.

find the

Tag line Block Off

001 01 00 $\rightarrow \sqrt{120}$

001	01	01	w_{21}
-----	----	----	----------

001	01	10	w_{22}
-----	----	----	----------

001 01 11 11

ye sub words
B5 mein hai to
iska mtlb cache ki
L1 mein B5 para hai

② Fully Associative Mapping:

→ Ismein koi bhi block kahi par bhi aa sakta hai.

→ Ismein hits to zyada hojate hai lekin no. of comparisons to check that whether a block is present in cache is equal to total no. of lines.

③ N-1klay Associative Mapping:

→ N ki value humein given hogi like 2-way, 4-way, n -way associative mapping.

→ Hum open cache ko n sets mein divide kr denge

→ Main memory se cache mein jo block/foame ayega wo $K \bmod S$ (set no) ki basis par ayega. Ab ek set mein 2, 3, 4, -- ki number bhi lines hosakti hai. ab agar koi block So mein store hua hai to wo So ki kisi bhi block mein hosakta hai.

7

Tag	Set No	B offset
4 bits	1 bits	2 bits

Page No.

Page No.

Date _____

Day _____

→ LRU (Least Recently Used) Algo is used to replace the cache lines.

! Revise Virtual memory, pages, TLB from OS Notes.

"INSTRUCTION LEVEL PARALLELISM"

① PIPELINING:

In CPU it allows multiple functional units such as fetching, decoding, execute, store to work on diff stages of multiple inst at a same time.

It speeds up the execution by overlapping the execution of inst.

- ↳ Pipelining operates at speed of slowest stage
- ↳ Delay causes the pipeline to stall.

② Multiple Issue:

It enhances processor performance by allowing multiple inst to be issued and executed simultaneously in a single clock cycle.

This is achieved by replicating functional units within the processor enabling it to handle more than one operation at a time.

If a processor has 2 adders then it can execute two addition operations at a time.

Date _____

Day _____

Types of Multiple Issue:

↳ Static Multiple Issue

- Inst are scheduled at compile time
- The compiler determines which inst can be executed in // and schedules them accordingly.

↳ Dynamic Multiple Issue

- Inst are scheduled at runtime.
- The processor itself decides which instruction can be executed simultaneously, based on the current state of execution.
- Processors using dynamic MI are called superscalar.

► Speculation in Multiple Issue:

This includes making educated guess about certain outcomes and execute them in // with other inst.

↳ Compiler based Speculation: The compiler inserts extra code to check whether the speculation was correct and if not it performs necessary action.

↳ Hardware based Speculation: The hardware stores speculative results in a buffer and if its correct it's committed to register or memory. If the speculation is incorrect the results are discarded and re-executed.

Date _____

Day _____

"HARDWARE MULTITHREADING"

It is used to improve the performance of processors by allowing them to execute multiple threads of a program simultaneously.

Granularity:

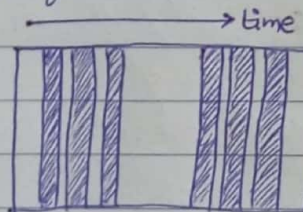
It is the ratio of computation to communication. Periods of computation are typically separated from periods of communication by synchronization events.

• **Fine-grain Parallelism:** Relatively small amounts of computational work are done b/w communication events.

↳ Tasks are broken down into very small units of computation.
 ↳ Since the tasks are small there is frequent communication and syn b/w processors, leading to higher overhead in terms of managing the //ism.

↳ Eg: Pipelining within CPUs, multithreading within loops.

↳ Facilitates Load balancing as it focuses more on distributing the task among all cores, it gives less opportunity for performance enhancement.



▨ Communication
 □ Computation

Instruction level //ism exploits this by executing multiple insts per clk cycle using techniques like superscalar execution or pipelining

Page No.

Date _____

Day _____

• **Coarse-Grain Parallelism:** Relatively larger amount of computational work is done b/w communication & sync events.

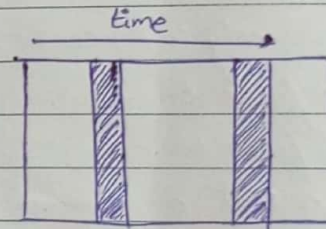
↳ Tasks are divided into larger units, such as whole function, process or independent computations.

↳ Larger independent chunks of computation being processed in //.

↳ Lower overhead compared to fine-grain bcz syn b/w task is less frequent.

↳ Programs run for longer period of time.

↳ Harder to load balance efficiently.



▨ Communication
 □ Computation

Multiprocessing where different process are assigned to separate processors, distributed computing

Thread-level //ism:

Attempts to provide //ism thru simultaneous execution of diff threads, it provides a coarser-grained //ism than ILP. threads are larger or coarser than the the fine-grained units.

Page No.

Date _____

Day _____

▷ Hardware Multithreading:

Provides a means for system to continue doing useful work when the task being currently executed has stalled—for example, if the current task has to wait for data to be loaded from memory.

For this system must support very rapid switching b/w threads.

↳ Fine-grained Multithreading—The processor switches b/w threads after each instr., skipping threads that are stalled. Longer threads may have to wait for longer period of time.

↳ Coarse-grained Multithreading—only switches the threads that are stalled waiting for a time-consuming operation to complete (eg: load from main memory)

↳ Simultaneous Multithreading—or SMT is a variation on fine-grained multi—g. It attempts to exploit superscalar processors by allowing multiple threads to make use of multiple functional units.

"Preferred threads" → threads that have many inst ready to execute.