

Prerequisite:

1) Rounding off

13.46

$\hookrightarrow 0 \leq n \leq 4$ same

$\hookrightarrow 5 \leq n \leq 9$ +1 previous digit

2) Unit Time Bit/Byte/Address

Kilo (K)	10^3	2^{10}
Mega (M)	10^6	2^{20}
Giga (G)	10^9	2^{30}
Tera (T)	10^{12}	2^{40}
milli (m)	10^{-3}	2^{-10}
micro (μ)	10^{-6}	2^{-20}
nano (n)	10^{-9}	2^{-30}
Pico (p)	10^{-12}	2^{-40}

3) $2^0, 2^1, 2^2, 2^3, 2^4, 2^5, 2^6, 2^7, 2^8, 2^9, 2^{10}$

1, 2, 4, 8, 16, 32, 64, 128, 256, 512, 1024

Basics:

1) Architecture

• CPU design \rightarrow instruction set
 \hookrightarrow Addressing modes

• Data format

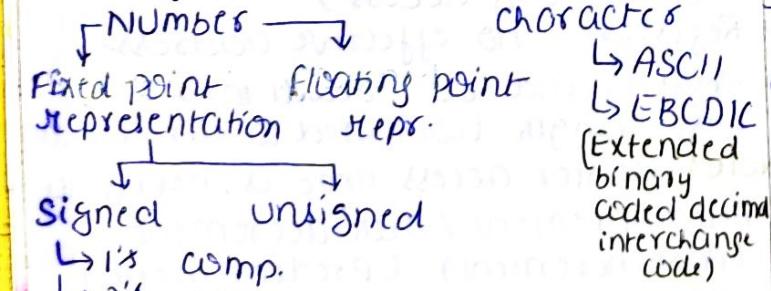
2) Organization

• IO organization

• Memory

• Performance improvement
 \hookrightarrow Pipelining

3) Data (in binary) \rightarrow



NOTE: 1 Byte = 8 bit

B \rightarrow Byte; b \rightarrow bit

* Everything in memory is in binary. CPU need to know format to understand it.

4) Buses

- Address bus (unidirectional) $\xrightarrow{\text{CPU to memory/IO}}$
- Data bus (Bidirectional) $\xleftrightarrow{\quad}$
- Control bus

5) Control signal send by ... to CPU?

• IO devices: interrupt signal

• Memory: Wait, Ready signal

6) There is only 1 each GPR, in system

7) 32 bit architecture means

AC is of 32 bits i.e. max. word size to ALU is 32 bits. & MDR size = 32 bits.

Word size of CPU = 32 bits

8) Word size of CPU & memory are different thing.

If memory is given word addressable

& its word size not given then consider it equal to word size of CPU.

9) Flag register is updated after each ALU operation.

(10) Program status word (PSW) =

PAC + Flag

11) Address register & Data register are used to access memory.

Micro Operation:

1) Smallest operation that CPU can perform in single step.

2) Denoted by Register transfer lang.

3) Time taken to execute 1 uop is called 1 CPU cycle time.

4) When 2 uops are independent, they can be performed simultaneously.

5) Types of uop:

• Arithmetic • Logical • Shift

6) Types of Shift uop:

• Logical shift } unsigned no
 • Circular } signed no

• Arithmetic } signed no

After shift sign should remain same

If sign changes it gives Arithmetic overflow error.

7) Memory \rightarrow Byte addressable (Default)

\hookrightarrow word \leftarrow word will be read or written at a time

NOTE: NO. of Inst from Ia to Ib = (b - a + 1)

- 8) If memory has n cells then
addr size (in bits) = $\lceil \log_2 n \rceil$
CPU of that system should have
min size of
 $AR, SP, PC = \lceil \log_2 n \rceil$ bits
Addr bus should have atleast $\log_2 n$ lines

Instruction:

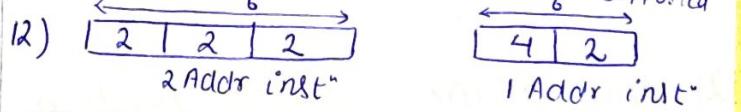
tells about operation **Opcode** | **operand info**

- i) If Opcode is of n bits, max no. of unique operations CPU can perform = Max no. of instⁿ supported by CPU = $2^{n - \text{size of ISA}}$ (Type)
- ii) Instⁿ set Architecture (ISA) is set of all instⁿ supported by CPU.
- iii) Change of instⁿ is costlier process, so avoid it.
- iv) Max 3 address can be specified within each instruction. (^{4 addr not used now})
- v) AC based archi supports only 1 address or 0 address inst. Stack based archi \rightarrow 0 "
- vi) If a system supports x address instⁿ then it can support all addr instⁿ below x .
- vii) If a system supports x address instⁿ then there must be at least 1 instⁿ of this type supported.
- viii) If there are 'n' no. of GPR's, no. of bits for register reference = $\log_2 n$

- ix) Multiple instⁿ supported: Always start solving with instⁿ having min no. of **opcode bits**.

- x) In n bit data

- Unsigned nos: max = $(2^n - 1)$
min = 0
- Signed nos: max = $+(2^{n-1} - 1)$
min = $-(2^{n-1} - 1)$ signed magnitude
 (-2^{n-1}) 2's comp

- xi) Fixed length instⁿ \Rightarrow variable length Opcode
xii) Variable len instⁿ \Rightarrow fixed len Opcode
size of Opcode = $\log_2 (\text{Total no. of inst}^n)$


- To get max 1 Addr instⁿ, use one 2 addⁿ
- To get min 1 Addr. instⁿ, use 3. 2 addⁿ instⁿ

Addressing Mode:

- 1) Effective address: operand addr or target address.
It's always memory addr not register reference.

- 2) 6 Phases of instⁿ cycle:

- Instⁿ fetch Fetch cycle
- Instⁿ decode
- Eff addr calculation
- Operand fetch Execution cycle
- Execution
- Write back result

- 3) In execution phase of branch instⁿ condition is checked and value of PC is updated accordingly jump

- 4) Special purpose registers (SPR's) are implicitly accessed by CPU, no need of explicit address.

- 5) Types of addressing modes:

- Implied (INCA)
- Immediate - to initialize registers with constant
- Direct / Absolute
- Indirect - used to implement pointer (2 memory access)
- Register - no effective address
- Register indirect - used to shorten instⁿ length (compared to direct mode)

NOTE: Register access time is negligible.

- Autoincrement / Autosubtract (Post increment / Predecrement)
Used to access array / loop
- Indexed / Index register causes extra
eff. addr = Base addr + index value (instⁿ) (index register)
- Used to access an array element.
- Time memory access.
- 2 instⁿ occurred, 1 to set index register value and 1 to fetch memory.

- PC Relative mode: **intrasegment**
Eff. addr = PC + Offset
→ present in inst.
- Base Register mode - **intersegment**
Eff. addr = Base addr + offset
→ Base register & inst.

* Forward jumping → +ve offset
Backward " → -ve "

6) Auto incr, decr, indexed, scaled
relative, Base Register modes are computable, rest are non computable

7) PC relative & Base Register mode supports reallocation.

8) When instr is fetched, PC is incremented by size of current instr.

NOTE: If instr size not given, take by default 1.

#CPU: 1 instr → 6 phases → each phase has many **mop**

1) CPU cycle: time taken to perform 1 **mop**.

2) Clock rate / clock freq = $\frac{1}{\text{no. of } \mu\text{op performed in sec}}$ unit (Hz)

3) Avg instr execution time = CPI avg * 1 CPU cycle

For n instr, = $n * \text{CPI avg} * 1 \text{ cycle}$ time

4) MIPS (Million instr per sec)

MIPS = $\frac{\text{no. of instr}}{\text{Total exec time}} * 10^6 = \frac{10^6}{\text{CPI avg * cycle}}$
= $\frac{\text{Clock freq} * 10^6}{\text{CPI avg}}$

• MIPS is not a good parameter.

5) FLOPS - floating point operations per sec. (good parameter)

6) ALU i/p: Operands, function code

O/P: Result, status → control signal

NOTE: If constant is present in instr no need to put it in register to give to ALU. (directly give)

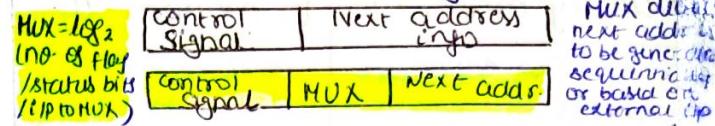
7) Control word: set of control signals to perform 1 **mop**.

8) Control word length is fixed for all Operations.

9) Types of control units:

- Hardwired - faster, not flexible
- MicroProgrammed - flexible
Control words are stored in control memory sequentially.

10) For each **mop**, 1 **control word** or control word is stored in control memory.



11) Set of control word in control memory is called **microprogram**.

12) Types of microprogrammed CU:

- | | |
|------------------------------|-------------------------|
| Horizontal | Vertical |
| • 1 bit for each ctrl signal | • Control word is small |
| • Control word large | • Decoder needed |
| • No decoder needed, faster | • Slower for each grp |
| • Degree of 11 is more | • Less |

13) Signals which can't be grouped, & stored in **Horizontal** manner.

14) Speed Hardwired > Horizontal > vertical

15) If there are 10 signals in microprog CU & at a time at most 2 can be active make 2 grp each having (10 signals + 1) (1 for case when no signal is active)

16) Types of processors:

- | | |
|--|--|
| RISC | CISC |
| • Reduced instr set computer | • Complex instr set computer |
| • Less no. of instr | • More |
| • Fixed length instr | • Variable length |
| • Simple instr | • Complex instr |
| • Limited addressing modes | • More & complex addressing modes |
| • 1 cycle per instr | • Multiple cycles |
| • More no. of registers | • Less no. of registers |
| • Register to register arithmetic opr only | • Register to memory, memory to register |
| • Easy to implement using hardwired CU | • Difficult using hardwired CU |

NOTE: Stack grows upward means address increases.

NOTE: ROM is a control memory used in microprogrammed CU to store **instr**.

* Expanding opcode is common in RISC to have fixed length instr.

9) Organization:

- 1) RAM is not IO device.
- 2) MM is directly connected to CPU & CPU connected to IO device via interface.
- 3) Speed: CPU > Main memory > IO devices

4) Memory mapped IO

- don't have separate address.
 - Memory wastage
 - distinguish memory & IO via address.
 - More inst & addressing modes for IO also
 - inst & addr modes of memory & used for IO
 - more IO devices
 - ALU opr can be directly performed
- 5) Baud \rightarrow Bits/Sec
- 6) Efficiency of transmission = $\frac{\text{no. of bits in data}}{\text{total bits transferred (including synchronization bit)}} * 100\%$

$$7) \text{Effective transfer rate} = \text{Efficiency} * \text{Transfer rate}$$

8) In asynchronous data transfer, eff data transfer rate decreases, time increases.

9) Modes of Transfer:

- (i) Programmed IO (Polling) (IO to CPU)
- Time = time taken by IO to share its status register (based on ID speed) + Data transfer time (depends on IO speed)

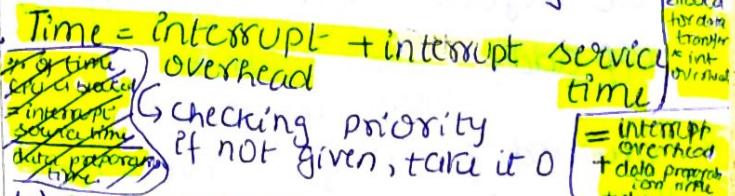
NOTE: By default status register = 1 Byte

- (ii) Interrupt initiated (IO to CPU)
- ISR - interrupt service routine (Present in main memory)
 - part of driver of IO device.
 - Internal interrupt - S/W interrupt
 - External " " - H/W "
 - CPU don't know about external interrupt while executing inst".
 - vector address is send via databus

NOTE: Among IO devices, Harddisk has highest priority by default.

- Priority based interrupt handling
- software sol. \downarrow Hardware sol. \downarrow

Serial Daisy chaining Parallel



$$10) \text{Performance} = \frac{\text{Older tech time}}{\text{newer (faster) tech time}}$$

$$= \frac{\text{Performance of new}}{\text{Per of Old}}$$

If speed up > 1 \Rightarrow New one better

Speed up < 1 \Rightarrow Old one better

11) Programmed IO & interrupt IO may have data transfer b/w memory & IO via CPU.

11) DMA: can generate address & send control signal to memory.

• Direct data transfer from IO to memory

• In DMA, CPU is used only once at starting to initialize DMA controller with starting address & data count.

• Modes of DMA transfer:

\rightarrow Constantly transferring data when gets control

\rightarrow Burst Mode: no overlapping

\rightarrow Cycle Stealing Mode:

• DMA gets bus control for 1 CPU cycle

• Overlapping of time

• Data is prepared with IO speed, transferred with memory speed

Time = Preparation + transfer

\rightarrow Interleaving mode: CPU gives bus control only when its free.

12) % of time CPU is blocked due to DMA

$t_x \rightarrow$ time taken by IO to prepare

$t_y \rightarrow$ time taken by memory to transfer data

• In burst mode

$$= \frac{t_y}{t_x+t_y} * 100\%$$

• In cycle stealing mode

$$= \frac{t_y}{t_x} * 100\%$$

• In interleaving mode = 0%

13) During DMA transfer, CPU state (register values) is intact.

14) Speed: Programmed < Interrupt < DMA

For very small data, DMA & interrupt

15) In DMA, CPU can send HLDA even in blw of inst execution.

16) In quer device speed given, use it to find preparation time. Transfer time directly given or in terms of memory cycle.

17) During inst^x execution, DMA transfer can be done but not interrupt transfer
18) CPU sends 2 info to DMAC: Data count & starting address

19) Max data transferred using DMA without CPU intervention = $2^x - 1$
 $x \rightarrow$ no. of bits in data count

20) In RAM/ROM 1 pin is required for each line.

	chip select	Read	Write	operation
0	X	X		No Opr.
1	0	0		No Opr.
1		1	X	Read
1		0	1	Write

21) If there's power pin, ground pin is also needed.

22) Total capacity = no. of chips * chip capacity

23) Default storage unit is bits.

$$\text{Ex:- } 32\text{K} \times 1 \text{ RAM} = 32\text{K bit RAM}$$

24) If require address more

→ Vertical arrangement

If require data more on each cell
→ Horizontal arrangement

If require address & data more
→ Hybrid arrangement

25) In vertical arrangement, MSB bit of address is used as chip select.

26) CPU can initiate memory request only when memory is ready.

27) In a refresh operation, one complete row of cells is refreshed.

28) 1 chip refresh time = no. of rows * refresh opr^y time
 \downarrow
= n chip refresh time

All chips are refreshed parallelly.

29) Associative memory (content addressable memory) - faster than SRAM & expensive too.

- used for TLB, cache
- cells don't have address.

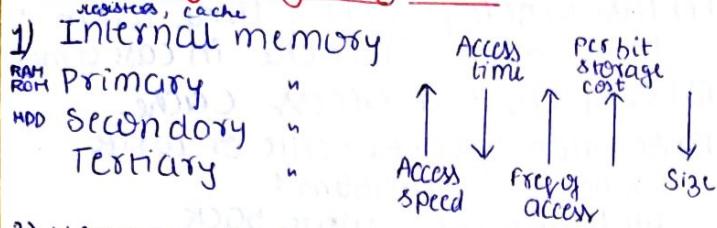
30) Speed:

Associative memory > SRAM > DRAM

31) Spatial Locality of reference - accessing near by addresses again
Temporal LOR - Accessing same address again.

32) If ram is of $64\text{M} \times 16$ bits then 64M cells are arranged in rows & columns.

Memory Organization



2) Memory representation:

$$\text{No. of cells} \times 1 \text{ cell size.}$$

3) By default - Byte addressable

4) Memory cycle time - time to access 1 memory location

5) Decoder used to access a specific cell using address.

6) RAM - volatile, ROM - non volatile

7) Memory access rate = $\frac{1}{\text{Memory cycle time}}$

8) Memory access decoder = $a \times b$

$$a = \text{address size}, b = \text{no. of cells} = 2^a$$

9) When CPU is turned on, program in ROM is runned for POST & Booting.

10) Bootstrap program's execution does booting.

11) Types of RAM:

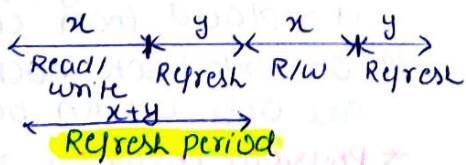
Static (SRAM)

- Made of flip-flops
- Faster & costlier
- No refresh required
- Idle power consumption less
- Operational power consumption more
- used for cache memory

Dynamic (DRAM)

- Made of capacitors stores data in form of electric charge.
- Slower & cheaper
- Refresh needed periodic
- more
- less
- used for MM.

12) In DRAM



$$\% \text{ of time used for refresh} = \frac{y}{x+y} \times 100\%. \quad \% \text{ of time used for R/W} = \frac{x}{x+y} \times 100\%$$

13) Multiplication table for two n-bit unsigned numbers = $2^{2n} + 2n$ bits

14) Addition table for two n-bit unsigned numbers = $2^{2n} + (n+1)$ bits

Cache Memory → TOP level memory

1) Hit Ratio = $\frac{\text{No. of hits in cache}}{\text{Total memory reference}}$ (Hit + Miss)

2) Miss Ratio = $\frac{\text{No. of miss}}{\text{(Hit + Miss)}}$

3) Hit Ratio + Miss Ratio = 1

4) General formula: Avg memory access

$$T_{avg} = H * \text{time required in case of hit} + (1-H) * \text{time required in case of miss}$$

5) Types of cache access:

Simultaneous & Hierarchical access

* $H \rightarrow$ Hit ratio

$t_{cm} \rightarrow$ Cache memory access time

$t_{mm} \rightarrow$ main memory access time

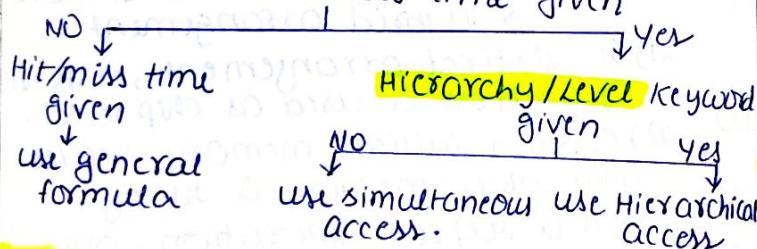
6) Simultaneous access →

$$T_{avg} = H * t_{cm} + (1-H) t_{mm}$$

7) Hierarchical access → $T_{avg} = t_{cm} + (1-H) t_{mm}$

8) Which formula to use:

if CM or MM access time given



9) If given, search time for Hit/Miss in cache is negligible → use simultaneous access.

10) Access time of memory does not change. T_{avg} changes with hit ratio.

11) $t_{cm} \leq T_{avg} \leq t_{mm}$

↳ 100% H

↳ 0% Hit rate

NOTE: Cache → only cache

Cached memory → Cache + MM

Uncached → Only MM

12) Cached memory access time = T_{avg}

13) T_{avg} when LDR is used:

Simultaneous: $T_{avg} = H * t_{cm} + (1-H) T_{block}$

Hierarchical: $T_{avg} = t_{cm} + (1-H) T_{block}$

↳ block access time from MM

14) $T_{block} = \text{block size} * T_{MM}$

15) In 1 T_{MM} , 1 MM cell is accessed.

16) If $H = 100\%$, $\Rightarrow T_{avg} = t_{cm}$

17) Miss penalty & extra time taken to bring block in cache in case of miss

18) Only CPU can access cache.

19) Solution of cache write or write propagation problem:

Write through, Write back

20) T_{avg} in write through cache:

• $T_{avg\ read} = H * t_{cm} + (1-H) t_{mm}$ (Simultaneous)

• $T_{avg\ read} = t_{cm} + (1-H) t_{mm}$ (Hierarchical)

• $T_{avg\ write} = t_{mm}$ or $\max(t_{cm}, t_{mm})$

↳ always simultaneous access

• $T_{avg} = \frac{\% \text{ of read}}{\text{Opr.}} * T_{avg\ read} + \frac{\% \text{ of write}}{\text{Opr.}} * T_{avg\ write}$

• Eff. Hit = $\frac{\% \text{ of read}}{\text{both read & write opr.}} * \text{Hit ratio for read opr.}$ for only read opr.

• For write, Hit ratio = 0

21) Every time there is read miss, block is brought from MM to CM.

22) Every time there is write miss, block is not brought from MM to CM (no write allocate)

23) For better performance use

• Write through with no write allocate

• Write back with write allocate

24) For no write allocate, miss penalty = 0.

25) Cache line → Cache block

26) Addressing latency / latency time of memory system is time taken by decoder to select a cell for operation

27) In write through cache, block is replaced from cache directly.

28) In write back cache, dirty blocks are only written back to MM.

* Physical memory → MM

* If using cache doubles speed of the memory, it means $T_{avg} = \frac{T_{MM}}{2}$

* Cache is accessed using internal bus, system bus is needed to access MM.

Cache Mapping:

- 1) CPU always generates **MM address**. ^{By default}
- 2) Cache mapping ^{> Set Associativity}
- 3) Direct set associative fully assoc.
- 4) Direct Mapping: (1-way set assoc mapping)
- CM block no = MM block no. / no. of blocks in CM
- 5) Tag identifier among all MM blocks which maps to one index, which one is present in cache
- 6) NO. of blocks = size of cache / block size
- 7) MM address in direct mapping
- | | | |
|-----|-------------|-------------|
| Tag | CM block no | Byte offset |
|-----|-------------|-------------|
- MM block no. $\rightarrow \log_2(\text{CM size})$
- 8) NO. of bits for Byte offset = $\log_2(\text{Block size})$ ^{not dependent on block size}
- 9) NO. of bits for CM = $\log_2(\text{no. of blocks in CM})$
- 10) Block size of CM & MM is same.
- 11) Index in direct mapping \rightarrow CM block number
- 12) Tag in direct mapping = MM address size - $\log_2(\text{CM size})$
- 13) Byte offset is NOT used to check hit/miss in any of the mapping
- 14) Cache controller = HW to check hit/miss + Metadata storage (like tag bits)
- 15) Metadata size or tag directory size = NO. of blocks in CM * (Tag + extra bits) (for all mapping)
- 16) Dirty bit - to know block is updated or not
 $0 \rightarrow$ Not modified $1 \rightarrow$ modified
- 17) CM is **not addressable**, it is indexable.
- 18) For a given cache size, block size & MM size, tag and CM block no. is same for byte and word addressable MM. Only MM address & byte offset size changes.
- 19) If MM address is given in decimal
 • MM block no. = $\frac{\text{Address}}{\text{block size}}$
- CM block no. = MM block no. % no. of blocks in CM
- 20) Set Associative Mapping:
- CM set no. = MM block no. / no. of sets in CM
- 21) n-way set associative means n blocks at each set in CM.
- 22) NO. of sets = NO. of blocks in CM / Associativity
- 23) Index in set associative \rightarrow set no.
- 24) MM address in set associative mapping:
- | | | |
|-----|-----------|-------------|
| Tag | Set (no.) | Byte offset |
|-----|-----------|-------------|
- MM block no. $\rightarrow \log_2(\text{CM size}) - \log_2(\text{Associativity})$
- 25) Tag bits are maintained for each block in cache.
- 26) Tag bit size in set assoc mapping = $\text{MM addr} - (\log_2(\text{size}) - \log_2 k)$
- * k \rightarrow Associativity
- 27) If associativity is doubled, set no. size reduces by 1 & tag size increases by 1.
- 28) Fully Associative Mapping:
- | | |
|-----|-------------|
| Tag | Byte offset |
|-----|-------------|
- 29) Index bit = 0 bits
- 30) Mapping: $\frac{1}{\text{NO. of sets}} = 0 \leftarrow$ single set
 $\frac{\text{MM block no.}}{\text{NO. of sets}} \leftarrow$ NO. of set
- 31) Tag bit = MM address size - $\log_2(\text{Block size})$
- 32) Tag = MM block no.
- 33) Size of tag is max in fully assoc & min in direct mapping.
- 34) Size of index is max in direct map. & min in fully assoc.
- 35) Max difference in tags (of directly mapping & fully assoc) = $\log_2(\text{NO. of blocks in CM})$
- NOTE:** If given CPU is 32 bit arch, it never means memory is word addressable. By default memory is byte addressable.

BLOCK Replacement :

- 1) Direct mapping don't need replacement policy. (below block packing at 1 index)
- 2) Replacement policies:
 - FIFO
 - Optimal
 - LRU
- 3) In case of miss penalty, block is transferred to CM byte by byte if MM is byte addr. Otherwise word by word. (MM & CM accessed multiple times) to transfer 1 block from MM to CM
- 4) Transfer rate = Bandwidth
- 5) Finding transfer rate from miss penalty:
Find cycles required to transfer block
→ Time → transfer rate
- 6) Types of cache misses:
 - cold or compulsory miss
 - capacity miss → block got replaced below CM got full
 - conflict miss → replaced due to set got full
- 7) To reduce cold miss - increase block size.
- 8) At a time max cold miss possible = no. of blocks in MM
- 9) To reduce capacity miss - increase CM size
- 10) To reduce conflict miss - increase associativity.
- 11) There is no conflict miss in fully associative cache.

Direct Mapping hardware:

- 1) NO. of MUX for tag selection = tag bit size
- 2) Size of MUX = NO. of blocks in CM : L
- 3) NO. of comparator = 1
- 4) size of comparator = tag bit size
- 5) Hit latency = 1 MUX delay + 1 comparator delay
- 6) time to declare hit → All MUX works parallelly

K-way set associative mapping HWs

- 1) Num. of MUX for tag select = k bunch = k * tag bit size
- 2) Size of MUX = NO. of sets : 1
- 3) NO. of comparator = K
- 4) size of comparator = tag bit size
- 5) OR gate = 1 (K input OR gate)
- 6) Hit latency = MUX delay + comparator delay + OR gate delay

NOTE: 4 i/p OR gate ≡ three 2 i/p OR gate
2 i/p OR gate ≡ 2:1 MUX

- ## # Fully associative mapping HWs:
- 1) MUX not required (no need to select set)
 - 2) NO. of comparator = NO. of blocks in CM
 - 3) size of comparator = tag bit size
 - 4) OR gate = 1 (i/p = no. of blocks in CM)
 - 5) Hit latency = comparator delay + OR gate delay

NOTE: If any of the delay is not given, take it zero.

- Associative memory can be used to implement fully associative cache.

Multilevel caches:

- 1) To minimize access time & maximize hit rate we use multilevel cache.
- | | | |
|--------------------------|-------------------------|----|
| (smaller)
L1
cache | (Bigger)
L2
cache | MM |
|--------------------------|-------------------------|----|
- access time: $t_1 < t_2 \leq t_3$
Hit ratio: $H_1 < H_2 \leq H_3$

2) Simultaneous access:

$$T_{avg} = H_1 t_1 + (1-H_1) H_2 t_2 + (1-H_1)(1-H_2) t_3$$

3) Hierarchical access:

$$T_{avg} = H_1 t_1 + (1-H_1)[t_2 + (1-H_2)t_3] \\ = H_1 t_1 + (1-H_1) H_2 (t_1 + t_2) + (1-H_1)(1-H_2) (t_1 + t_2 + t_3)$$

- 4) Prob of access of L1 = H_1
- " " " " L2 = $(1-H_1) H_2$
- " " " " MM = $(1-H_1)(1-H_2)$

- 5) Total access time = Number of memory reference * Tavg

Rel. b/w content of L1 & L2:

- 1) Inclusion Policy: Content of L1 should be present in L2, value inside block may be different.
- 2) Exclusion Policy: Content of L1 need not be present in L2.
L2 contains only blocks which are replaced from L1. So L2 is called victim cache
- * Smaller chip has lesser access time.
- * Smaller block has less miss penalty.



1) Tavg for instruction (Hierarchical)

$$T_{avg\ inst} = H_1 \cdot \text{cache}_1 \cdot T_{ic} + (1 - H_1) \cdot \text{cache}_1$$

$$[H_2 \cdot (T_{ic} + T_2) + (1 - H_2)(T_{ic} + T_2 + T_{MM})]$$

2) $T_{avg\ data} = H_{dc} T_{dc} + (1 - H_{dc}) [H_2(T_{dc} + T_2)$

$$+ (1 - H_2)(T_{dc} + T_2 + T_{MM})]$$

3) $T_{avg} = (\% \text{ of inst} * T_{avg\ inst}) + (\% \text{ of data} * T_{avg\ data})$

4) When block size of L1 & L2 is different
in case of miss in L1, block will
be brought from L2 to L1 whose
size is equal to **block size of L1**.

5) In one cache access, **one entire block** is accessed.

NOTE: $T_{avg} \propto$ Miss rate

$$T_{avg} \propto \frac{1}{\text{Hit Rate}}$$

X NOTE X

1) Pointer, recursion \rightarrow **indirect addr mode**

Array, structure, record \rightarrow **indexed addr mode**

2) To send 16 bit addr via 8 bit bus
first 8 bit is send then 8 bit.
8 MUX (2:1) required.

3) In Ques to find min register, if
 $inst^*$ is given $R_2 \leftarrow R_2 + R_3$

don't perform it as $R_2 \leftarrow R_3 + R_2$.
(unless given to optimize)

4) To have single data transfer using
DMA, system bus is accessed twice
one for data transfer from **IO to DMAC** & again from **DMAC to MM**.

5) When transfer is complete, **DMAC**
in DMA sends interrupt signal to CPU.

* By default array is row major
order.

* A sequence $\langle a_i \rangle$ is bounded above
by K iff all $a_i \leq K$.

* Bounded lower by K iff all $a_i \geq K$

* Jump stmt may jump to same inst.

Magnetic Disk

1) Each recording surface has 1
read/write head.

At a time only 1 is active.

2) NO. of surface = 2 * NO. of platters

3) Sector of a track \rightarrow disk sector

4) NO. of tracks = $2 * \text{NO. of platters} * \frac{\text{NO. of tracks per surface}}{\text{NO. of sectors per track}}$

5) NO. of sectors = $2 * \text{NO. of platters} * \frac{\text{NO. of tracks per surface} * \text{NO. of sectors per track}}{\text{NO. of sectors per track}}$

6) consecutive sectors collection

7) By default, each $\xrightarrow{\text{track}} \text{sector}$ has
constant capacity

• Sectors of inner track have **more density**.

• **constant angular velocity**.

8) If sectors have variable capacity then
has **constant density & constant linear velocity**. (Per unit block has constant capacity, so calculate by circumf.)

9) NO. of Bytes = NO. of sectors * 1 sector
on disk capacity

10) Smallest addressable **unit \rightarrow Sector**

11) 1 sector is accessed at a time.

12) NO. of bits in disk address = $\log_2 (\text{NO. of sectors in disk})$

13) **Disk access time** = seek time +
rotational latency + 1 sector transfer
time + other delays

↳ if not mentioned in ques take 0

14) Rotation of disk is only in 1 direction.

15) If current & target head are given, then find rotational latency using rotation speed.

• If not given then use average rotational latency.

$$\text{Avg rotational latency} = \frac{1}{2} \text{ rotation time}$$

16) 1 sector transfer = $\frac{1}{\text{rotation time}} \cdot \frac{\text{NO. of sectors per track}}{2}$

* Disk keeps rotating with some speed
whether read/write takes place or not.

17) Rotational latency = $\frac{1 \text{ rotation time}}{\text{NO. of sectors per track}}$

* NO. of sectors to rotate \times (first-instant)

18) To calculate rotation time from RPM, divide 60000 by RPM.
Answer will be in msec.

19) Disk transfer rate is used in
• DMA to find data preparation time
• Programmed IO for finding status check time.

20) Multiple sector access time :

- Sequentially stored (cluster)

$$t = 1 \text{ seek time} + 1 \text{ rotational latency} + n * 1 \text{ sector transfer time}$$

- Random

$$t = n [1 \text{ seek time} + 1 \text{ rotational latency} + 1 \text{ sector transfer time}]$$

21) NO. of cylinders = NO. of tracks per surface

22) Disk scheduling is done to have min seek time.

23) Addressing is done as cylinder \rightarrow surface \rightarrow sector

$\langle C, H, S \rangle$
cylinder no \nwarrow surface no \searrow sector no

24) Sector no. for given addr $\langle C, H, S \rangle$

$$= C * \text{NO. of sectors per cylinder} + H * \text{NO. of sectors per track} + S$$

25) If sector no. given S, find address

$$C = \lfloor S / \text{NO. of sectors per cylinder} \rfloor$$

$$H = \left\lfloor \frac{(S \% \text{NO. of sectors per cylinder})}{\text{sectors per track}} \right\rfloor$$

\hookrightarrow Remainder = S

26) If file is stored on n sectors, to find last sector no.

Find starting sector no. from addr, add $(n - 1)$

\hookrightarrow From sector no., find address of last sector.

27) In disk + DMA block

- Find transfer rate from track capacity & 1 rotation time.
- Find transfer time then solve for DMA queue.

28) If disk rotates with constant linear velocity, each track has different capacity. Find according to circumference of tracks.

Capacity/cm circumference is constant.

29) If disk rotates with constant angular velocity, capacity of each track is constant.

30) If no. of platters not given take it 1.

* We use concept of cylinder to store a file so that seek time required only once.

* Latency \rightarrow delay

* If seek time is not given calculate average seek time.

* Virtual memory size is limited by disk size & size of MAR (address register).

* Seek latency is not proportional to seek distance due to arm starting & stopping inertia.

* Larger disk block/sector size means less no. of blocks to fetch & hence better throughput.

* Disk requires device driver.

* RAID (Redundant array of independent disk) is used to provide fault tolerance & high speed (parallel accessing possible).

* Formatting of disk is done to store some metadata to make disk ready to use. So, formatted disk has slightly less capacity than unformatted disk.

\hookrightarrow In each sector some metadata is stored.

For eff data transfer rate consider sector capacity = $\frac{\text{Sector Capacity}}{\text{Metadata}}$

* Recording width = Outer \times inner \times if track width = 0.1 \Rightarrow No. of track $= \frac{1}{0.1} = 10$

11) Single Precision : (Denormalized no)

$$\text{Value} = (-1)^s \times 0.M \times 2^{-126}$$

Double Precision:

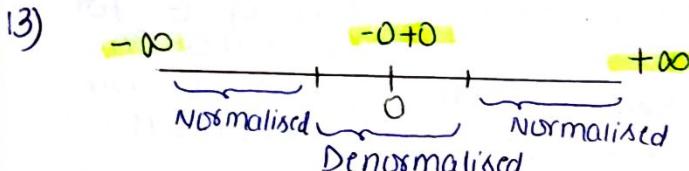
$$\text{Value} = (-1)^s \times 0.M \times 2^{-1022}$$

12) For denormalized no.

$$E = 00\cdots 0 = E_{\min} = E_{\max}$$

$$M_{\min} = 00\cdots 1$$

$$M_{\max} = 11\cdots 1$$



For max/min type ques remember this number line.

Pipelining: used when same process is done over multiple i/p

1) Flynn's classification of computers

◦ SISD (von Neumann)
single instruction stream
data stream execute

◦ SIMD (Single Instruction Multiple Data)
multiple data streams execute

◦ MISD (Multiple Instruction Stream)
multiple instructions streams execute

◦ MIMD (Multiple Instruction Multiple Data)
multiple instructions streams multiple data streams execute

◦ Pipelining (hypothetical)
multiple pipelines

2) If operation performed in all segment is said as 1 task completed.

3) Pipeline cycle time: min time in which all segment can perform their respective suboperation.

4) Cycle time $\rightarrow t_p$; n task, K segment pipeline

No. of cycles to complete 'n' task $= \frac{(K+n-1)}{K}$

Time needed $= (K+n-1) \times t_p$

5) Non pipeline system takes t_n time to perform a task.

Time (for 'n' task) $= n t_n$

6) Speed up $= \frac{\text{Slowest one time}}{\text{Faster one time}} = \frac{\text{Non Pipeline time}}{\text{Pipeline time}}$

$$S = \frac{n t_n}{(K+n-1) t_p}$$

$$S \propto n$$

7) If $n \gg (K-1)$, neglect $(K-1)$ to fill the pipe

$$\text{Speedup} = \frac{t_n}{t_p} \leftarrow \text{Max possible speedup}$$

8) Special case: If time taken to perform 1 task in Pipeline and non-pipeline is same ie delay of all segment same

$$t_n = K t_p$$

$$\text{then } S = K$$

$$9) S \leq K$$

10) If n not given or asked maxs or ideal S $\Rightarrow S = \frac{t_n}{t_p}$

$$\text{If } n \text{ given } \Rightarrow S = \frac{n t_n}{(K+n-1) t_p}$$

$$11) \text{Clock rate or Frequency} = \frac{1}{\text{cycle time}}$$

12) In ideal condition, pipeline gives 1 output per cycle.

13) If delay of all segments are different, we use buffer for synchronization.

$$14) t_p = \max(\text{segment delay}) + \frac{\text{Buffer delay}}{\text{Registers delay}}$$

$$15) t_n = \sum \text{ of all segments delay}$$

Registers not required

16) If all segments are having same delay & there is no buffer

$$S_{ideal} = K$$

17) To have better performance try to take segments with almost equal delay.

18) Latency: After how much time next i/p is given

Time taken by a task to complete

Non pipelined $\Rightarrow t_n$

Pipelined $\Rightarrow t_p$

19) Throughput: No. of operations per unit time

$$\text{Non pipeline} \Rightarrow \frac{1}{t_n}$$

$$\text{Pipeline} \Rightarrow \frac{n}{(K+n-1) t_p}$$

Latency $\Rightarrow \frac{1}{t_p}$

20) CPI of pipeline in ideal case without RAZORD = 1.

NOTE: Segment = Stage both r same chg

Instruction Pipeline :

1) When branch instⁱ is detected in decode phase, all inst^j which are fetched after branch instⁱ are removed from pipeline.

2) All branch instⁱ have stall cycles, whether jump is taken or not.

3) Stall cycles / Bubbles : Extra cycles

4) If after ith stage of branch instⁱ, execution outcome is available then no. of stall cycles = $i-1$ (due to branch) = no. of stages from decode to execution

5) If there are x branch instⁱ, total stalls = $x * \text{no. of stalls for branch inst}^i$

6) Stalls = Segment no. - segment per dependency (Ca) of write back - no. of OF Operand fetch

7) Total stalls = no. of dependent instⁱ * stalls per dependency

8) No. of stalls = Ca - No. of instⁱ b/w dependent instⁱ (Not consecutive)

If it's $\leq 0 \Rightarrow$ no stall

9) Loadⁱ instⁱ is performed completely in execution phase. (No need of write back phase). (In operand forwarding)

10) Operand forwarding can't be used if operand is not coming from ALU (load & store)

11) Solⁱ of Control Hazard:

Hardware Solⁱ

• Branch prediction

Stalls will be considered if prediction is wrong.

Software Solⁱ

• Delayed branch (done by compiler)

12) Data Hazard classification:

(i) RAW : True dependency

• H/W Solⁱ : Operand forwarding

• May have stalls due to memory access. Inorder dependency

(ii) WAW : False dependency (Output data dependency)

• H/W Solⁱ : Register renaming

• no stalls • Out of order dependency

(iii) WAR : False depc. (Antidependency)

• H/W Solⁱ : Register renaming

• no stalls • Out of order dependency

(13) RAR don't have any conflict.

(14) No. of stalls due to structural hazard = extra cycles/b/w each instruction is taking

(15) CPI = $\frac{\text{Total no. of cycles}}{\text{Total no. of inst}}$

NOTE: I₁: $R_1 \leftarrow R_2 + R_3$ gts only 1 dependency
 I₂: $R_1 \xrightarrow{\text{RAW}} \leftarrow R_1 + R_2$ I₁ $\xleftarrow{\text{I2}} I_2$
 I₃: $R_3 \xrightarrow{\text{I2}} \leftarrow R_1 + R_2$ as dep. b/w inst. not reg.
 (16) Asynchronous Pipeline: in a stage, different instⁱ take different no. of cycles.

Pipeline Hazards :

1) Situation that prevent next instⁱ to execute & leads to hazard, stalls.

2) Type:

- Structural hazard / Resource conflict
- Data hazard / Data dependency
- Control hazard / Branch difficulty

3) Structural hazard:

- When 2 segments try to use same resource at a time.
- no solution

4) Data hazard:

- Result of an instⁱ is used as i/p in next inst^j.

5) Solⁱ of Data dependency:

Software Solⁱ

(for pipeline which can't detect data dependency)

(Done by compiler)

• Delayed load

• Instⁱ reordering

• Using no opr. instⁱ

(delay dependent instⁱ)

Hardware Solⁱ

(for pipeline which can detect data dependency)

• H/W interlock

↳ Default

↳ Have stall cycles

• Operand forwarding

↳ or bypassing

↳ Operand directly forwarded from AC to ALU

↳ no stalls for WAW dependency

* CPI :

• Pipeline without hazard, $CPI = \frac{k+n-1}{n}$

Ideal condition, $CPI = 1$

• Pipeline with hazard, $CPI = \frac{(k+n-1)+x}{n}$

$x \rightarrow$ no. of stall cycles
 $n \rightarrow$ total no. of inst.

Ideal case (ignore $k-1$), $CPI = 1 + \frac{x}{n}$

18) Avg inst. Exec time = CPI * cycle time

$$19) CPI_{ideal} = 1 + \frac{x}{n} = 1 + \left(\frac{\text{fraction of inst causing stall}}{f} \times \frac{\text{stall per inst}}{s} \right)$$

$$= (1-f) + f \times (1+s)$$

20) In ideal case, with hazard

$$S = \frac{tn}{CPI_{avg} * tp}$$

$\leftarrow 1 \text{ inst exec time in non pipeline}$
 $\leftarrow 1 \text{ inst exec time in pipeline}$

21) classic RISC Pipeline :

- IF
- ID - Inst decode & Register Read
- Branch inst :- decode, target addr calculation, condition evaluation, PC updation.
- EX - ALU operation executed
- MEM - memory accessed for load/store
- WB - write back result of ALU store ops or load operation in registers

22) Non linear pipeline: segments are not used sequentially. Ex: S1 S2 S3 S2 S3

23) Non linear pipeline :

- find forbidden latency (which may cause collision)

- Find Permissible latency

- Find collision vector

- Draw state diagram

For each state, if 0 is at i th place, right shift by i & take OR with original collision vector.

- Find greedy cycle

- Find avg " ", its min is ans.

24) For a single task, non pipelined may work better than pipeline

Time non pipelined \leq Time pipelined $\boxed{t_n=18}$ $\boxed{tp=21}$

25) Register renaming can eliminate all WAW/WAR hazards. \rightarrow True

26) In ideal case without hazard Efficiency of pipeline = 100%.

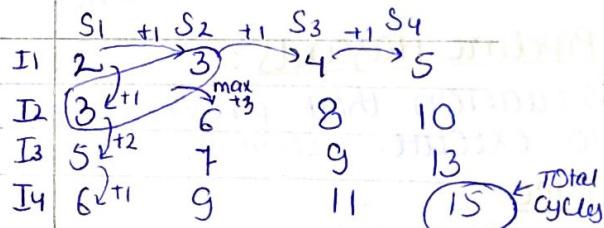
• with hazard

$$\text{Efficiency} = \frac{S}{S_{max}} = \frac{S}{K}$$

$S \leftarrow$ speedup
 $S_{max} \leftarrow$ given speedup
 $K \leftarrow$ no. of segments
 \leftarrow max speedup

* Asynchronous Pipeline ex:-

No. of cycles	S1	S2	S3	S4	\leftarrow Stages
I1	2	1	1	1	
I2	1	3	2	2	
I3	2	1	1	3	
I4	1	2	2	2	



* Pipeline with single segment cannot provide parallel processing.

* For single step, non pipeline can perform better than pipeline.

* All stages are perfectly balanced \rightarrow means all stages have equal delay

$$* CPI_{avg} = 0.0x(1) + 0.0y \times \text{oval}(1+2) + 0.0y \times 0.0b(1)$$

↓
No stall ↓
at b1.
at b2.
↓
stall ↓
No stall

$$* \text{Total exec time} = \frac{\text{no. of inst}}{\text{Avg inst exec time}}$$

* General pipeline: IF, ID, OF, EX, WB

Ex: branch inst condition check, load/store inst completes on 7 stage pipeline

* Latency of an inst is always less than on 5 stage pipeline. False

* 7 stage pipeline is always faster than 5 stage. False

$$* R_1 \leftarrow M[1000] \quad (\text{IF ID OF EX WB})$$

$$R_2 \leftarrow R_1 + R_3$$

- with operand forwarding, stalls = EX-OF = 1

- w/o " ", stalls = WB-DF = 2

* w/o operand forwarding, always fetch operand after write back

* Pipeline depth = no. of segments = k
* Efficiency: % of time every stage of the pipeline is being used