

Computer Organization & Architecture

Computer Architecture: Deals with conceptual design & fundamental operational structure

Computer Organization: Deals with physical devices and their interconnection with a perspective of improving the performance.

| Computer Architecture | Computer Organization |
|--|---|
| <ul style="list-style-type: none"> → CPU Design → Instructions → Addressing modes → Data formats | <ul style="list-style-type: none"> → I/O organization → Memory organization → pipelining |

System Bus:

It has 3 things:

i) Address bus (unidirectional)

ii) Data bus (Bidirectional)

iii) Control bus (every line is individually unidirectional)

Memory may also send controls to CPU.

Those signals are wait, ready.

CPU Registers:

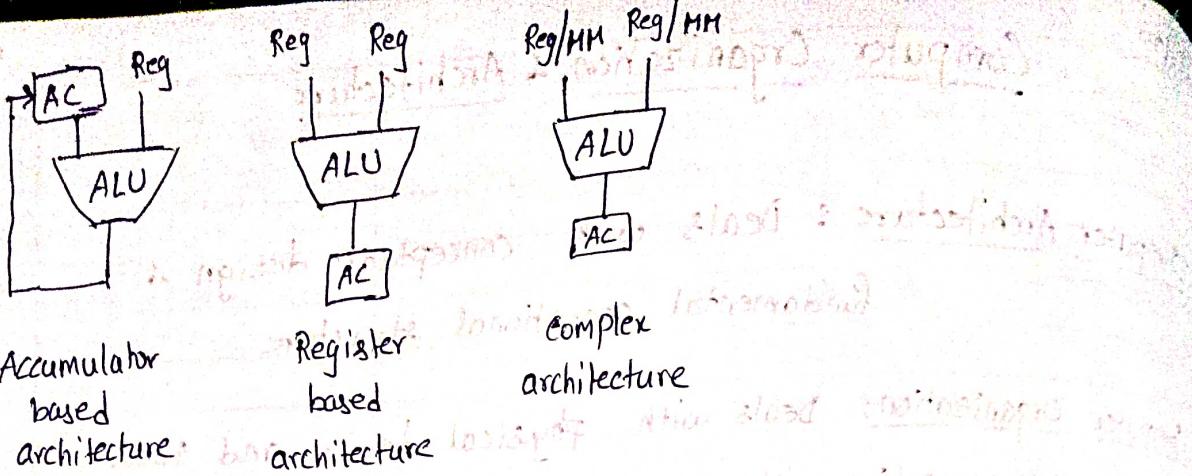
i) General Purpose Registers (store any general content)

ii) Special Purpose registers (stores special values)

Special Purpose registers

i) Accumulator: It is used to store result of ALU operations.

Sometime it may even store one i/p to ALU.



→ If inputs to ALU are given from stack then it is called stack based architecture.

32-bit Architecture:

- size of each i/p given to ALU is 32 bit.
- word size = 32 bits

64-bit Architecture

- size of each i/p given to ALU is 64 bit.
- word size = 64 bits.

(i) Program Counter:

- It stores the address of next instruction to be executed.

(ii) Instruction Register:

- It used to store current instruction which CPU is executing
(At the time of instruction fetch, the instruction is directly fetched into IR)

(iii) Stack Pointer:

- It used to store the address of top of the stack.

(iv) Flag Register / Program Status Word (PSW) / Status Register

Each bit of flag register has a ~~no purpose~~ purpose

(or)

Flag register stores status of ALU result

Z (zero flag) \leftarrow 1 (ALU's result = 0)

0 (ALU's result ≠ 0)

S (sign flag) \leftarrow 0 (+ve)

(-ve)

Logics for sign bit selection 0/1
and work given 0/1 or 1/0
(-ve) & +ve logic for sign bit selection 0/1

→ In some CPUs PSW is same as flag register.

In some CPUs PSW is a flag register + accumulator.

(vi) Memory Address Register:

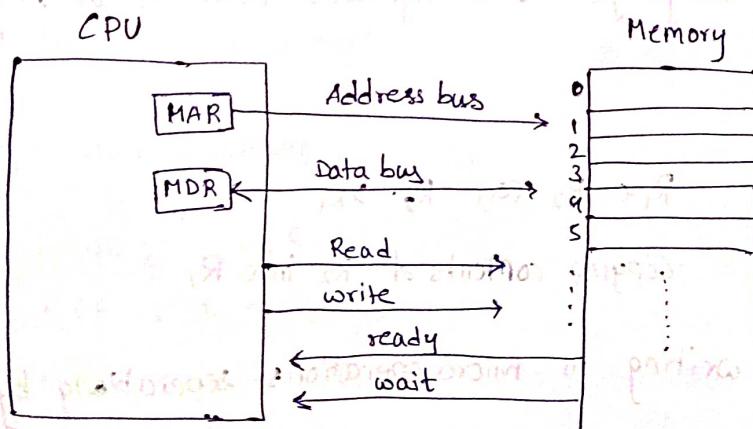
It is used to send address to memory.

MAP

(vii) Memory Data Register:

It is used to send data to memory (memory write) or

to receive data from memory (memory read).



Note: $2^{\text{size of PC}} \geq 2^{\text{size of MAR}}$ $\Rightarrow \text{size of PC} \geq \text{size of MAR}$

→ size of PC = size of MAR = Physical memory address length.

→ size of MBR = size of acc = length of each memory location.

If virtual memory is used $\text{size of PC} \geq \log_2(\text{VAS})$

$\text{size of MA} = \text{PA}$

Memory Read:

- i) CPU sends address to memory using address bus.
- ii) CPU enables read control signal.
- iii) Memory performs read operation & sends data to CPU using data bus.

Memory write:

- i) CPU sends address to memory using address bus.
- ii) CPU sends data using data bus. (from MBR)
- iii) CPU sends enabled write signal.
- iv) Memory performs write operation.

Micro Operation:

- The smallest indivisible operation performed by CPU is called micro-operation. Every μ-op is performed in 1 cycle.
- These are the operations performed on values stored in registers.
- Register Transfer Language is a symbolic notation of μ-ops.

RTL Notations:

Register transfer: $R_1 \leftarrow R_2$ (or) $R_2 \rightarrow R_1$

Copying contents of R_2 into R_1

Comma: writing to micro operations separating by comma means that they are performed parallelly.

Eg: $R_1 \leftarrow R_2 + R_3, PC \leftarrow PC + 1$

Memory transfer

Memory Read: $R_1 \leftarrow M[\text{address}]$

two ways

$\xrightarrow{\text{Directly specify address}} R_1 \leftarrow M[2000]$

$\xrightarrow{\text{Indirectly specify address}} R_1 \leftarrow M[\text{MAR}]$

Types

3 types

- i)
- ii)
- iii)

i) Arithmetic

- Addition
- Comparison
- 2's Complement
- Adder
- Incrementer

ii) Logic

- AND
- OR
- XOR
- X-NOR

iii) Shift

3 types

L

Memory write:

$$M[\text{Address}] \leftarrow R_1$$
$$M[2000] \leftarrow R_1 \quad M[\text{MAR}] \leftarrow R_1$$

Types of Micro operations:

3 types:

i) Arithmetic micro operations

ii) Logic Micro Operations

iii) Shift Micro Operations

i) Arithmetic micro operations:

- Addition, subtraction,

- Complement : $R_1 \leftarrow \overline{R_2}$

- 2's Complement : $R_1 \leftarrow \overline{R_2} + 1$

- Adding with 2's complement : $R_1 \leftarrow \overline{R_2} + \overline{R_3} + 1 \cong R_1 \leftarrow R_2 - R_3$

- Increment, Decrement

ii) Logic MicroOperations:

- AND : $R_1 \leftarrow R_2 \wedge R_3$

- OR : $R_1 \leftarrow R_2 \vee R_3$

- XOR : $R_1 \leftarrow R_2 \oplus R_3$

- X-NOR : $R_1 \leftarrow R_2 \odot R_3$

iii) Shift Micro Operations:

3 types:

- Logical shift

- Left shift

- Right shift

left shift : $1011 \Rightarrow 0110$

right shift : $1011 \Rightarrow 00101$

Every time a new bit is added, the new bit is 0.

2. Circular shift : (rotation)

i) left circular shift:

$1011 \Rightarrow 0111$

ii) right circular shift

$1011 \Rightarrow 1101$

3. Arithmetic shift:

* → Arithmetic shift can be performed only on signed numbers.

* → After performing arithmetic shift the sign must be retained.

Right shift:

~~1011~~ \Rightarrow ~~0011~~ \Rightarrow ~~0011~~ \Rightarrow ~~0011~~ \Rightarrow ~~0011~~

Left shift:

It is same as logical left shift but it is allowed by CPU only when sign is not going to change.

Eg: 1101
 $\swarrow \swarrow \swarrow$
 1010

Eg: 0101

Here apply left shift will result in change of sign

∴ This operation will not be allowed by CPU.
CPU generates Arithmetic left shift overflow error.

me a new bit
ded, the new bit

Word addressable & Byte addressable memory

word addressable:

Every word has an address. Bytes within a word doesn't have any address.

Byte addressable:

Every byte has an address.

Instructions:

Instruction :

| | |
|--------|----------|
| opcode | operands |
|--------|----------|

ISA (Instruction set Architecture):

It is collection of all those instructions supported by CPU

Types of instruction based on operands:



i) 4-address instruction:

within each instruction a maximum of 4 addresses can be specified

| | | | | |
|--------|------|------|------|------|
| opcode | add1 | add2 | add3 | add4 |
|--------|------|------|------|------|

→ next instruction's address.

→ The CPUs supporting 4-address instructions don't have

Program Counter.

→ Also here inst's need not to be stored contiguously

Disadvantage:

→ Large size instructions

∴ large program

→ Instruction fetch takes more time.

→ ~~Relocation~~ Relocation is not easy.

∴ we need to modify add4 of every instruction.
instruction has done aligned form

(ii) 3-Address Instruction:

→ Here we use PC instead of 4th address field.

→ So a maximum of 3 operands can be given

| | | | |
|--------|------|------|------|
| opcode | Add1 | Add2 | Add3 |
|--------|------|------|------|

→ All disadvantages of 4-address instruction are eliminated here.

(iii) 2-Address Instructions:

→ Maximum of 2 addresses can be specified within an instruction.

| | | |
|--------|------|------|
| opcode | add1 | add2 |
|--------|------|------|

→ Here one of the two operands acts as both source and destination.

Disadvantage:

→ More number of instruction in program compared to 3-address format.

→ So it require more memory for program

Eg: Consider $x = (a+b) * (c+d)$

3-address: 2-address

$$R_1 \leftarrow a+b$$

$$R_2 \leftarrow c+d$$

$$x \leftarrow R_1 * R_2$$

$$R_1 \leftarrow a$$

$$R_1 \leftarrow R_1 + b$$

$$R_2 \leftarrow c$$

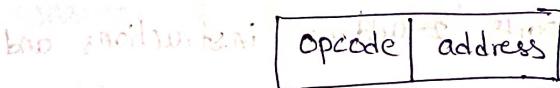
$$R_2 \leftarrow R_2 + d$$

$$R_1 \leftarrow R_1 * R_2$$

$$x \leftarrow R_1$$

(iv) 1-address instruction:

→ Maximum of 1 address is given in 1-address instruction.



→ The 2nd operand is accumulator (implicitly)

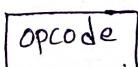
→ This accumulator is based on architecture.

Disadvantage:

→ More no of instructions in program

(v) 0-address instruction:

→ No address is mentioned in instruction



→ This is used in stack based architecture.

The operands are taken from the stack.

11/11/20

- Q1 Consider a digital computer which supports only 2-address instructions each with 16-bits. If address length is 6-bits then maximum and minimum how many instructions the system can support.

Sol:

| | | |
|---|--------|-------|
| 4 | 16 | 6 |
| ↓ | 16 → 4 | 6 → 6 |

$$\therefore \text{max} = 2^4 = 16$$

$$\text{min} = 1 \text{ instruction}$$

Multiple instructions supported: if 16 words of memory

- Q2 Consider a system which supports 2-address instructions and 1-address instructions both. The system has both 1-address and 2-address instructions. The system has 8-bit instructions and 2-bit addresses. If there are three ~~two~~ 2-address instructions in the system then maximum and minimum no of 1-address instructions supported is?

Sol:

2-add:

| | | |
|-----|----|----|
| OPC | A1 | A1 |
| 2 | 2 | 2 |

1-add:

| | |
|--------|---|
| Opcode | A |
| 4 | 2 |

$$\text{max: } \therefore 4 \text{ 2-add ins} \\ (00, 01, 10, 11)$$

but given only 3

2-add instructions

assume 00, 01, 10 are used.

Let 11 be unused for

2-address instruction

Now every instruction that starts with 11 has 1-address instruction.

| | |
|---|---|
| 4 | 2 |
|---|---|

opcode

1100

1101

1110

1111

1110 = A + C

1111 = D + E

∴ The system can have a maximum of 4 1-address instructions.

Minimum 1-address instructions = 1

- (Q3) Consider a system which supports 3-address and 2-address instructions. It has 30-bit instructions with 8-bit addresses. If there are ' x ' 3-address instructions then find the no of 2-address instructions possible.

Sol:

3-add:

| | | | |
|---|---|----|---|
| 6 | 8 | 18 | 8 |
|---|---|----|---|

↳

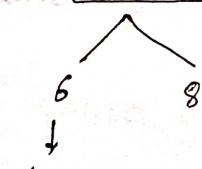
∴ 2^6 3-address instructions are possible if x

but only x are used

$2^6 - x$ are unused.

2-add:

| | | |
|----|---|---|
| 14 | 8 | 8 |
|----|---|---|



∴ $(2^6 - x)^8$ are max no of 2-address instructions possible.

(Q4) In the previous question, if there are 512 2-address instructions then find no of 3-address instructions.

Sol:

$$(2^6 - x)2^8 = 512$$

$$\frac{2^6 - x}{2^8} = \frac{512}{2^8}$$

$$2^6 - x = 2^8$$

$$x = 64 - 256 = 62$$

Addressing Modes

Effective address:

→ Address of operand in a computation type instruction or target address in a branch type instruction.

Instruction cycle:

There are several phases in execution of an instruction.

These phases are called instruction cycle.

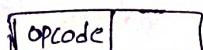
→ Generally there are 6 phases.

i) Instruction fetch:

→ Instruction is brought to IR from memory.
→ Now PC is incremented by the size of instruction that is fetched.

instruction address

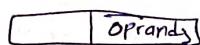
(ii) Instruction decode:



↓
decode

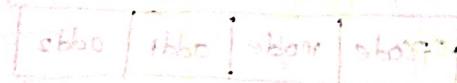
operation.

(iii) Effective address calculation:



↓
decode

↓
effective address



(iv) Operand fetch

From their effective address computed operands are brought to the CPU.

(v) Execution

The operation is performed and result is stored in the accumulator.

(vi) Copy back result:

Now the result is stored in memory.

This process is called instruction cycle.

* These phases are divided into 2 cycles

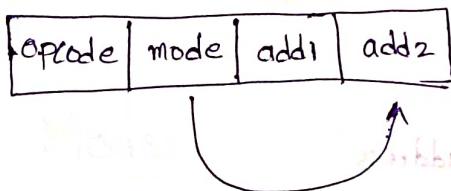
i) fetch cycle: Instruction fetch

ii) execution cycle: remaining 5 phases

→ Every instruction need not require all 6 phases of the instruction cycle.

Defn of addressing mode:

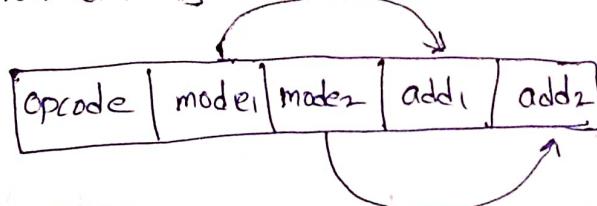
Addressing mode specifies how and from where the operand is obtained using address field value of instruction.



Here mode specifies how add2 has to be interpreted.

add1 is always taken from register.

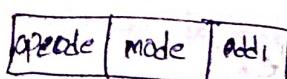
Another variant is



Different addressing modes are

i) Implied mode:

→ Here one of the operands is understood by the opcode itself.



Eg: INCA

This instruction is increment accumulator.

Here accumulator is implied from the opcode.

(iii) Immediate mode:

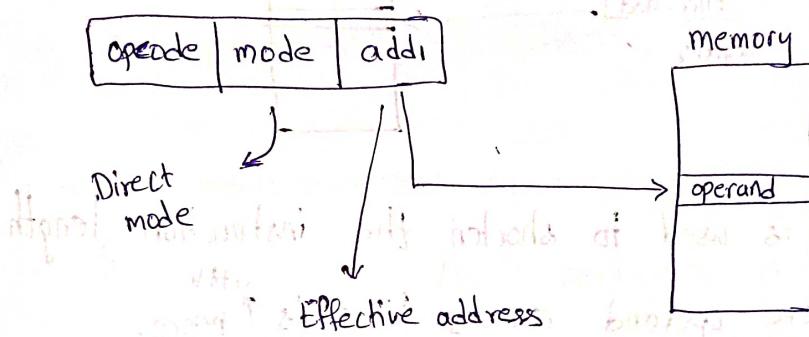
whatever is present in operand field that itself is considered operand.

generally operand range is small in immediate mode

→ It is used to initialize registers with constants.

(iv) Direct mode / Absolute mode:

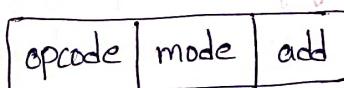
address field has address of operand.



→ Suitable for accessing static variables.

(iv) Indirect mode:

Address field has address of main memory in which effective address is present.

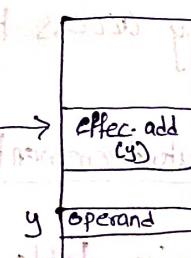


Denoted as

(loc)

Eg: Add R1, (1000)

This is not much suitable for pipeline execution as it requires 2 mem access for 1 operand fetch



→ This mode is used for implementing pointers.

Used in parameter passing.

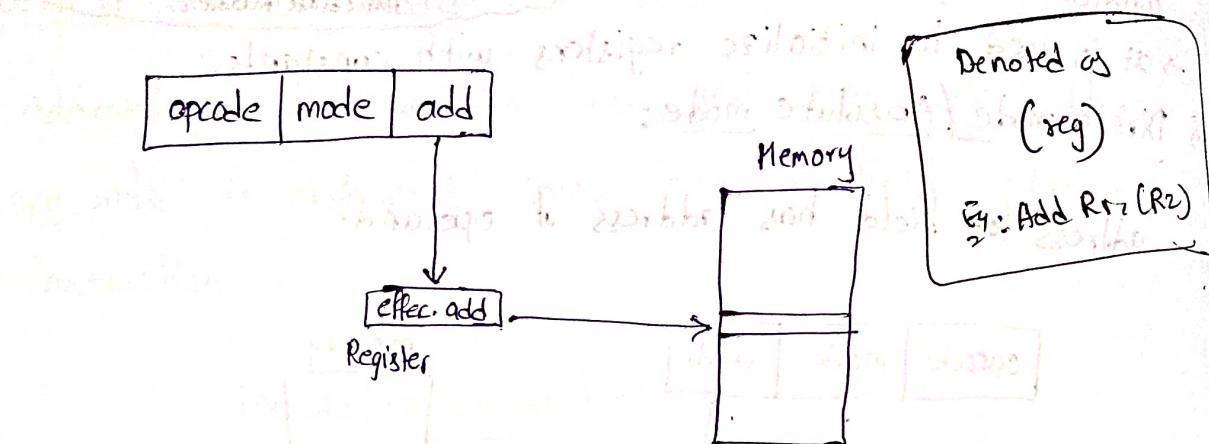
(v) Register mode / Register direct mode:

Address part of instruction specifies a register which hold the operands.

So indirection is generally used through registers

(vi) Register indirect mode:

- The address field specifies the address of register which contains effective address.



- This mode is used to shorten the instruction length.
- However here operand access time is more.

Direct mode → 1 memory access time.

Reg. indirect mode → 1 memory access time

+ 1 register access time

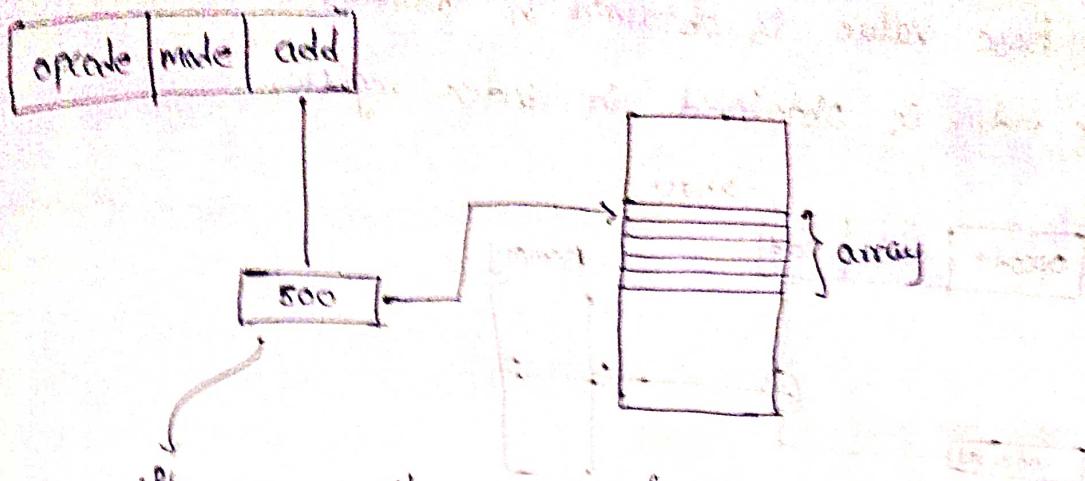
But register access time negligible compared to memory access time.

(vii) Auto increment / Auto decrement mode:

→ Used when accessing table of content (array) sequentially.

→ It is a variant of register indirect mode.

→ Content of register (effective address) is automatically incremented or decremented.



After 1 access, the content of register is incremented / decremented automatically.

* Auto increment mode \Rightarrow post increment

i.e., first memory access is performed, then value is incremented. Denoted as $(R_i) +$

* Auto decrement mode \Rightarrow pre decrement

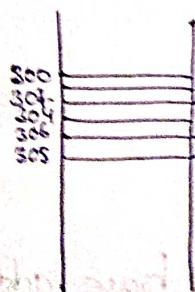
i.e., first decrement content of register then perform memory access.

The amount to be inc/dec is implicitly determined by the size of operand accessed

Denoted as $- (R_i)$

(iii) Indexed mode / Index register mode:

\rightarrow This method is used to access an array element.



$$\text{address of } A[3] = 300 + 3 * w$$

\rightarrow index value

$$= 300 + 6$$

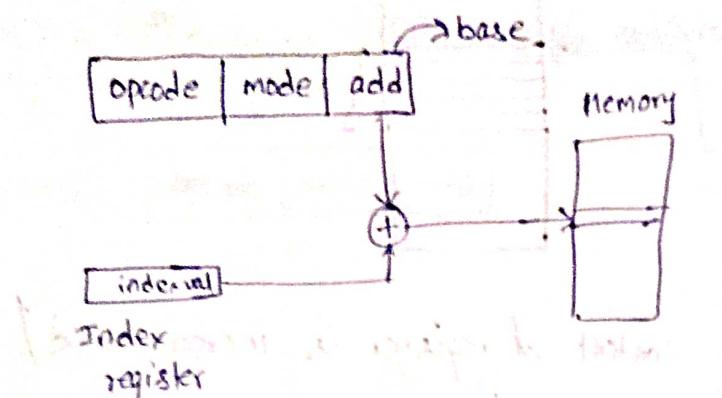
$$= 306$$

\therefore effective address = base + index value

\rightarrow It is denoted as $X(R)$

e.g.: Add $X(R_1), R_2$

- The base value is obtained in address field
- The index is obtained in index register.



Generally there will be ~~be~~ 2 instructions for accessing $A[i]$

1. Index register $\leftarrow w + i$

2. operand from base + index value.

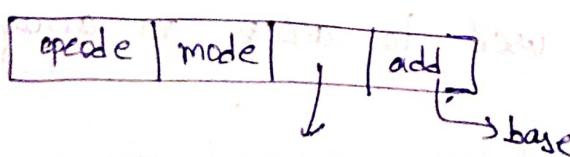
There are two ways for implementing this mode.

i) Index register is special purpose register:

Here CPU knows where index register is.

i.e., address of index register is implied.

ii) Index register is not a special purpose register:



Disadvantage:

- If array has to be relocated, then the base address which is stored in every instruction (array access instruction) has to be changed. This is a costly operation and hence this mode is not preferable. (updation of instruction is costly)

(ix) Index, Base register mode:

This is same as Indexed mode but the difference is that the base value is also put in a register.

| | | |
|--------|------|-----|
| opcode | mode | add |
|--------|------|-----|

→ This field has no use



→ Here the disadvantage of indexed mode is overcome.

→ Here also 2 implementations are possible.

(i) Index, base registers are special purpose registers:

| | | |
|--------|------|-----|
| opcode | mode | add |
|--------|------|-----|

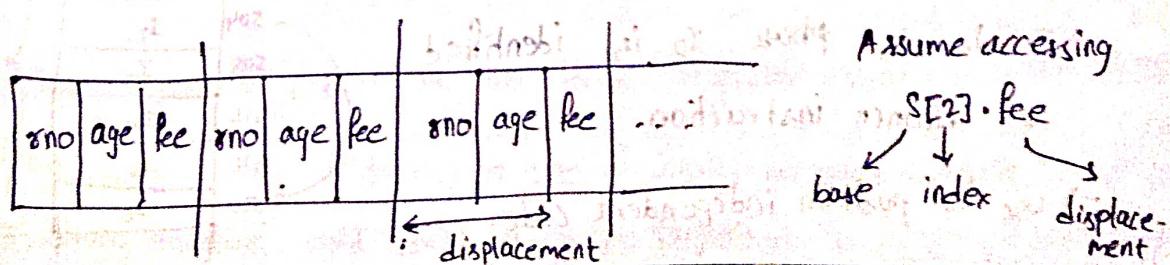
(ii) Index, base registers are special purpose registers with bits:

| | | | | |
|--------|------|-----------|----------|-----|
| opcode | mode | index reg | base reg | add |
|--------|------|-----------|----------|-----|

(x) Index, base register + Displacement mode:

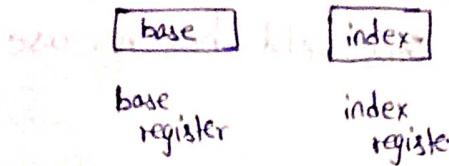
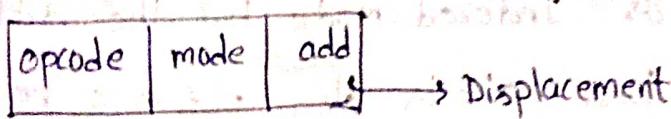
→ This method is used for accessing members within array of structures (array of structures).

Eg: student(rno,age,fee)



base.address
ray access
costly operation

\therefore Effective address = base + index + displacement



(xi) Scaled mode:

→ It is upgraded version of indexed register mode.

→ It is used for pointer calculations.

Effective address = address field value of instruction + (index register value * scaling factor)

taken implicitly.

→ In index register mode, we need to execute two instructions.

→ In scaled mode, we ~~add~~ just compute effective address in a different way (only with one instruction) and thus it is fast.

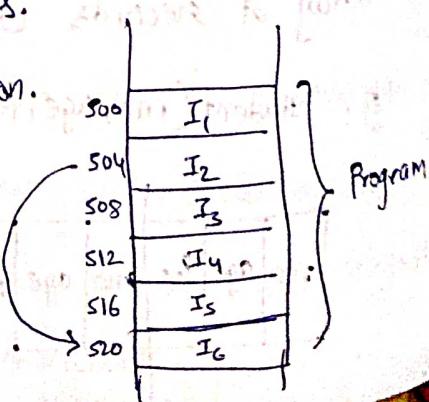
(xii) PC - Relative mode / Position independent mode:

→ This used for branch instructions.

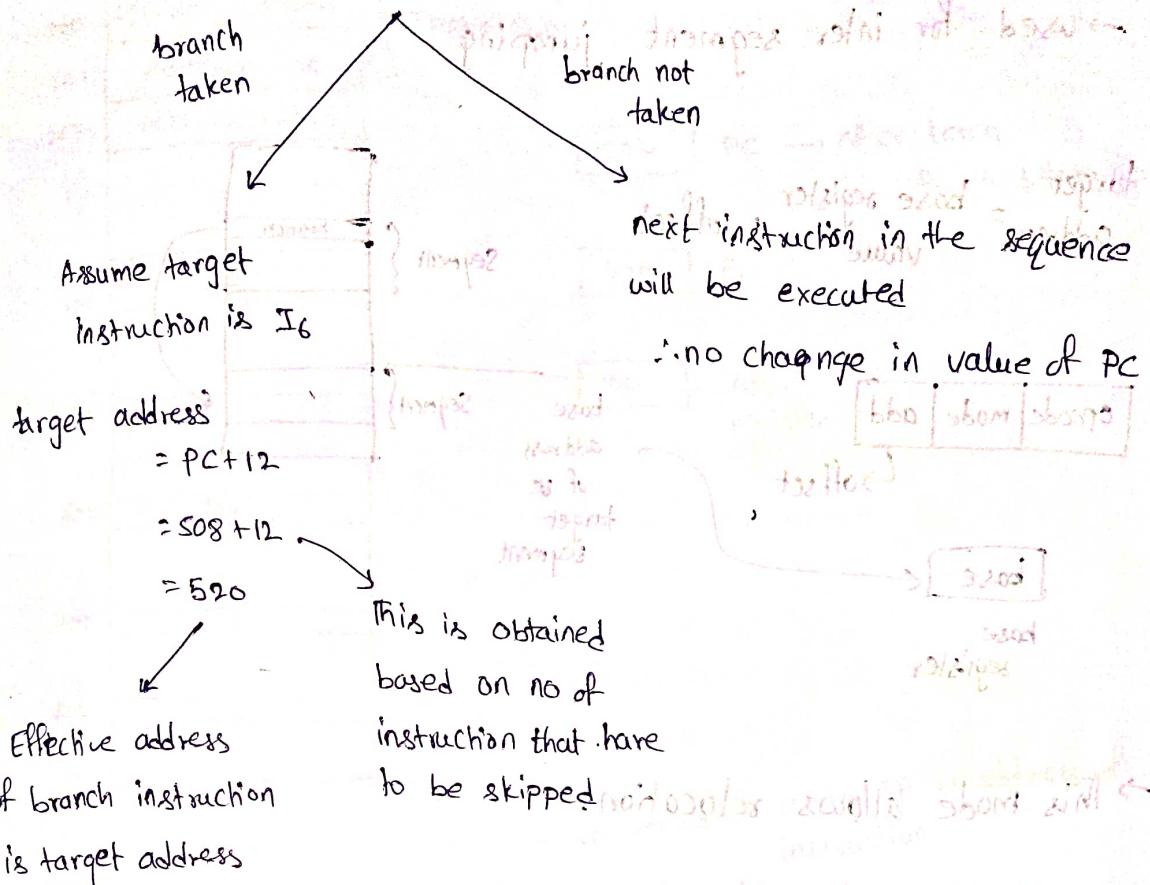
→ Assume, I_2 is currently under execution.
 $\Rightarrow PC = 508$

→ In decode phase I_2 is identified as branch instruction.

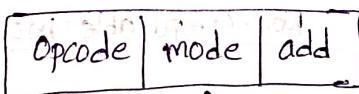
→ Use for position independent code.



now there will be 2 possibilities:



$$EA = PC + \text{no of locations to be skipped}$$



branching instruction
 branch to relative address
 instruction

no of locations to be skipped

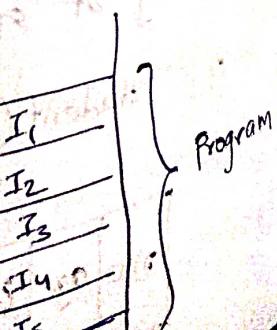
EA = PC + add field

offset

→ This mode is used for intra segment jump i.e., jump within same segment

For inter segment jump we use base register mode.

→ PC relative mode also reduces size of instruction because without address being relative, add field has to be larger.



(xiii) Base Register mode:

→ used for inter segment jumping

target address = base register value + offset

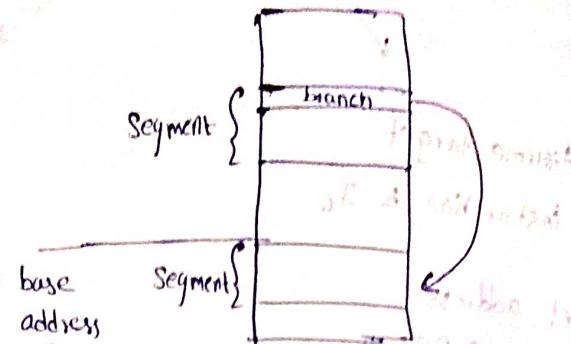
↳ for calculating target address

| | | |
|--------|------|-----|
| opcode | mode | add |
|--------|------|-----|

↳ offset

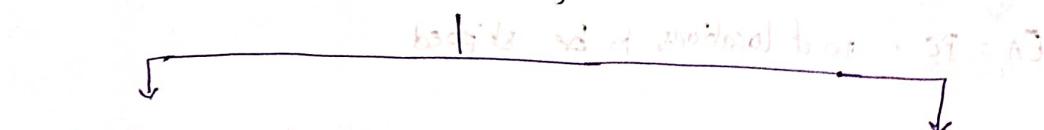
base

base register



→ This mode allows relocation.

Classification of addressing modes:



Computable modes

→ Computation is req. for obtaining EA

→ The modes under this category are:

(i) Auto inc/Auto dec.

(ii) Indexed mode

(iii) Index, base register mode

(iv) Index, base reg + displacement

(v) Scaled mode

(vi) PC-relative

(vii) base register mode

→ Remaining all are used for operands.

Non-computable mode

→ No computation required.

→ The modes in this category are:

(i) Implied

(ii) Immediate

(iii) Direct

(iv) Indirect

(v) Register direct

(vi) Register indirect

18/11/20

Eg:

| | |
|-----|------------------|
| 200 | opcode mode |
| 201 | Address = 500 |
| 202 | next instruction |
| 203 | : |
| 204 | : |
| 205 | : |
| 206 | : |
| 207 | : |
| 208 | : |
| 209 | 450 |
| 210 | 300 |
| 211 | : |
| 212 | 800 |
| 213 | 900 |
| 214 | : |
| 215 | 325 |
| 216 | : |
| 217 | 300 |
| 218 | : |

Instruction 200 has address 500, additional to 500 there is one more word length of instruction which is 201. So PC is set to 202 after fetch.

200 → PC → After fetch
PC is set to 202

400 → R1

100 → index register

(450) base ← target instruction

300 ← target instruction

In immediate mode EA is address of instruction

| Mode | Effective address | Operand |
|----------------------|-------------------|---------|
| 1. Immediate mode | 201 | 500 |
| 2. Direct mode | 500 | 800 |
| 3. Indirect mode | 800 | 300 |
| 4. Register mode | - | 400 |
| 5. Register indirect | 400 | 700 |
| 6. Auto decrement | 399 | 450 |
| 7. Indexed mode | 600 (500 + 100) | 900 |
| 8. PC relative | 202 + 500 = 702 | 300 |

EA is given only MM references but not register references
(Assume 500 is referring register)

Predecrement

000 200 bba

001 201 bba

002 202 bba

003 203 bba

Assuming that 500 is no of locations to be skipped

(Q5) An instruction is stored at location 300 with its address field at location 301. The address field has value 250. A processor register contains the number 200. Evaluate the effective address for each of the below addressing modes.

1. Direct $\rightarrow 250$
2. Immediate $\rightarrow 301$ (Address of instn)
3. Relative $\rightarrow 302 + 250 \therefore$ instruction at 552
4. Register indirect \rightarrow content of [200] 200
5. Auto increment $\rightarrow 200$ (EA)

Now register value is set to 201

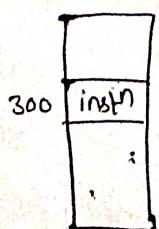
\rightarrow In PC-relative jumping can be either

i) forward jumping: jumping to next address
i.e., offset is +ve

ii) backward jumping: jumping to previous address
i.e., offset is -ve

(Q6) A relative branch mode type instruction is stored in memory at address 300. The branch is made to an address 450.

a) What should be the value of relative address field of the instruction.



PC value during execution is 301

$$\begin{aligned} 450 &= \text{PC} + \text{address field content} \\ &= 301 + \text{address field} \\ \therefore 149 & \end{aligned}$$

b) Determine and after

(Q7) A relative on address offset value

(Q8)
For comp
each add
following

250. \rightarrow Kef
the effective
modes.

instn)

instruction at
200 552

is set to 1201

address

ve

from statement

address

-ve.

from 1201

stored in

ade to an

address field

execution is 301

field content

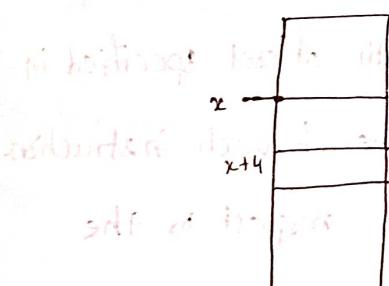
b) Determine the value of PC before instn fetch, after fetch
and after the execution phase

before fetch \rightarrow 300

after fetch \rightarrow 301

after execution \rightarrow 450

(Q7) A relative branch mode type instruction has taken branch to
an address 750. The ^{branch} instruction is 4 bytes and it has
offset value 400. This starting address of instn in memory is



PC $x+4 \rightarrow$ during execution

$$PC + offset = 750$$

$$(x+4) + 400 = 750$$

$$\therefore x = 346$$

$\therefore 346$ is the starting address.

(Q8) For computers based on three-address instruction formats,
each address field can be used to specify which of the
following

s1: A memory operand

s2: A processor register

s3: An implied accumulator register.

- a) either s_1 or s_2
- b) either s_2 or s_3
- c) only s_2 and s_3
- d) All of s_1, s_2, s_3

Sol:

An address field can be reference to memory or to a register but implied is not possible.

(References to implied are implied within operation)

\therefore opt @

Q9
9-17

Consider a RISC machine where each instruction is exactly 4 bytes long. Conditional and unconditional branch instructions use PC-relative addressing mode with offset specified in bytes to the target location of the branch instruction. Further the offset is always with respect to the address of next instruction in the program sequence.

Consider the following instruction sequence

instruction No:

Instruction

i:

add R₂, R₃, R₄

i+1:

sub R₅, R₆, R₇

i+2:

CMP R₁, R₉, R₁₀

i+3:

beq R₁, offset

If target of branch instruction is i, then the decimal value of the offset is _____

Sol:

while i+3 is under execution, PC will have i+4

\therefore i = PC + offset

$$i = (i+4) + \text{offset}$$

$$\text{offset} = -4 \text{ (instructions)}$$

$$= -16 \text{ (4*4)}$$

CPU:

→ It is main component of computer systems which performs operations and controls the system.

CPU Cycle:

It is amount of time required for CPU to perform one micro operation.

It is time \uparrow of clock.

$$\text{clock rate} = \frac{1}{\text{CPU cycle time}} = \frac{1}{\text{clock time period}}$$

CPI (Cycles per instruction):

No of CPU cycles required to execute an instruction.

Execution time:

time for execution of 1 instruction

i.e., $\text{CPI}_{\text{avg}} * \text{cycle time}$

∴ execution time for n instructions = $n * \text{CPI}_{\text{avg}} * \text{cycle time}$

→ CPI need not to be same for every instruction

decimal
arithmetic

$i+4$

MIPS : (Million Instructions Per Second)

CPU performance is given by MIPS

→ If CPU takes t -seconds for n instruction then

$$\text{no of instruction executed in one second} = \frac{n}{t}$$

$$\text{no of instructions executed in 1 sec} = \frac{\text{no of instruction}}{\text{total execution time}}$$

$$= \frac{\text{no of}}{t}$$

$$\therefore \text{performance in MIPS} = \frac{\text{no of instruction}}{(\text{execution time}) \times 10^6} \text{ MIPS}$$

$$= \frac{n}{(n \times \text{CPI}_{\text{avg}} \times \text{cycle time}) \times 10^6} \text{ MIPS}$$

$$\boxed{\text{MIPS} = \frac{1}{\text{CPI}_{\text{avg}} \times \text{cycle time} \times 10^6}}$$

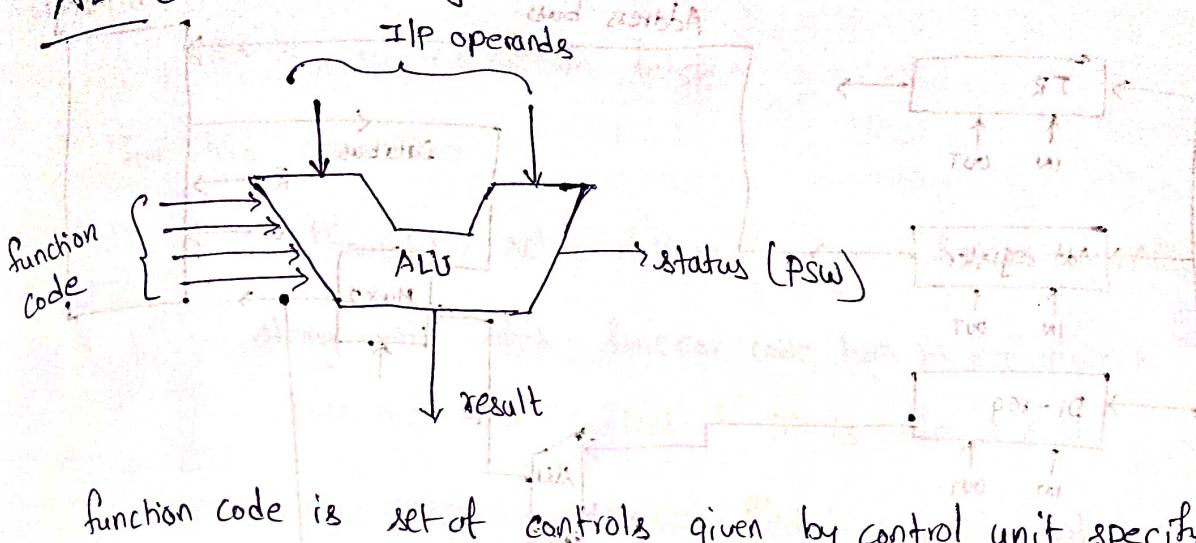
$$\boxed{\text{MIPS} = \frac{\text{Clock rate (in Hz)}}{\text{CPI}_{\text{avg}} \times 10^6}}$$

→ MIPS may provide false result for comparison of 2 processes.

because it doesn't consider complexity of instructions.

→ To overcome this we use FLOPS (Floating point operations per second)

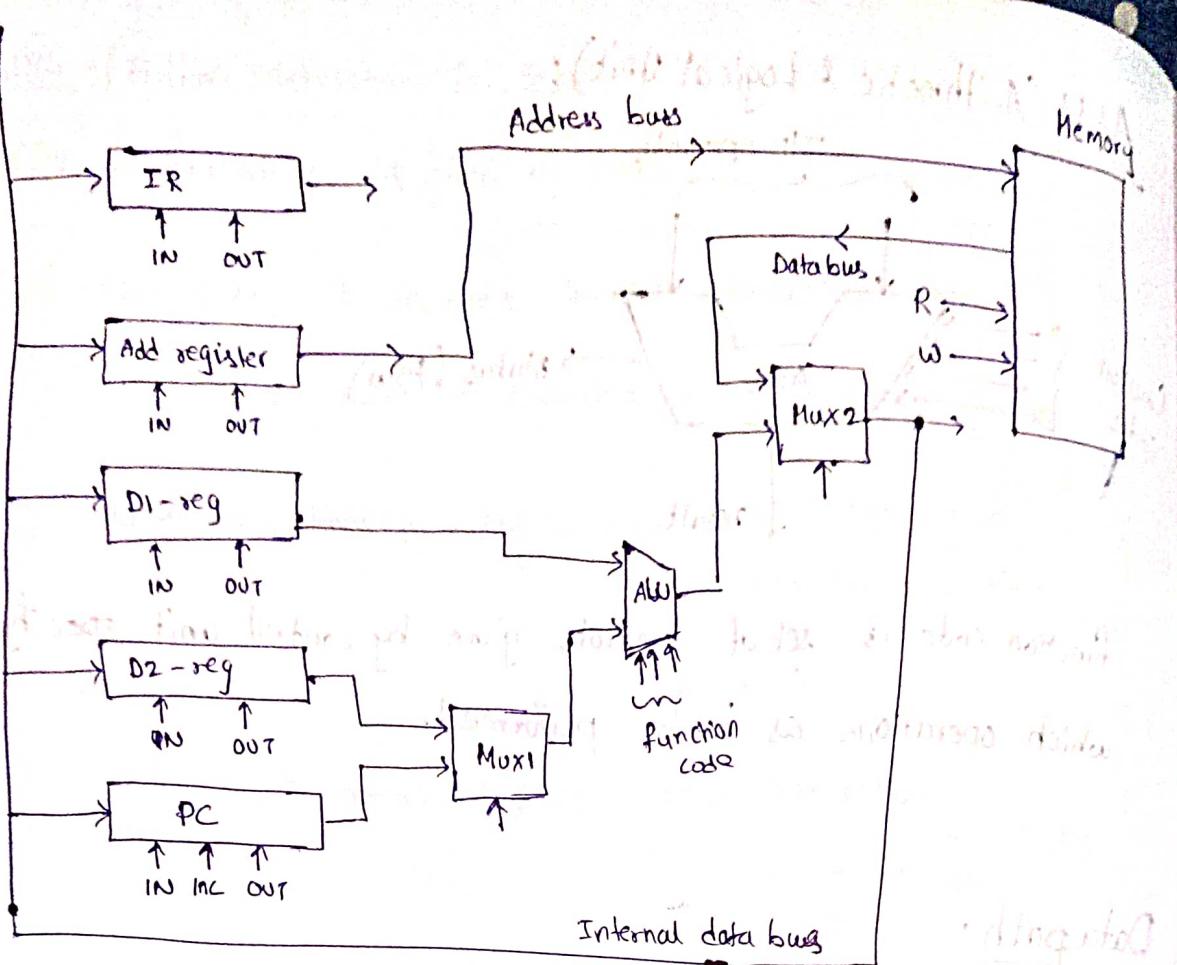
ALU (Arithmetic & Logical Unit):



function code is set of controls given by control unit specifying which operations as to be performed.

Data path:

- Collection of functional units such as arithmetic logic units or multipliers.
- Performs data processing operation.
- For any data processing operations, data follows a specific path. This path is called datapath.



Instruction fetch:

- $AR \leftarrow PC$
- $IR \leftarrow M[AR]$, $PC \leftarrow PC + 1$

Control Unit:

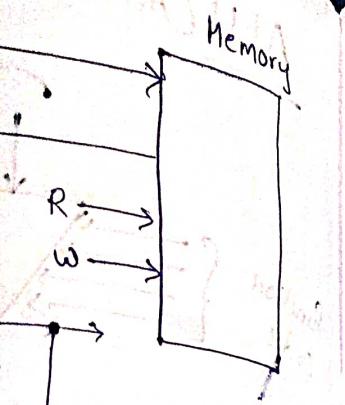
It generates control signals, sends those signal to different components and components work accordingly.

Control variable:

It is a name given to control signal

Control word:

Collection of all control signals. Each bit of control word corresponds to one signal. If the bit is 1, the corresponding control signal is sent.



Consider

$AR \leftarrow PC$ in instruction fetch.

for this operation

$PC_{out} = 1$; select of Mux1 = 1; select of Mux2 = 1; $AR_{in} = 1$;

along with this function code has to be given to ALU such that content of PC is O/P of ALU.

All these signals are specified in control word.

Consider

$IR \leftarrow MEAR$, $PC \leftarrow PC + 1$

$IR_{in} = 1$; $AR_{out} = 1$; read = 1; $PC_{inc} = 1$;
select of MUX2 = 0;

14/10/20

Control unit organization:

CU can be organized in 2 ways:

(i) Hardwired CU.

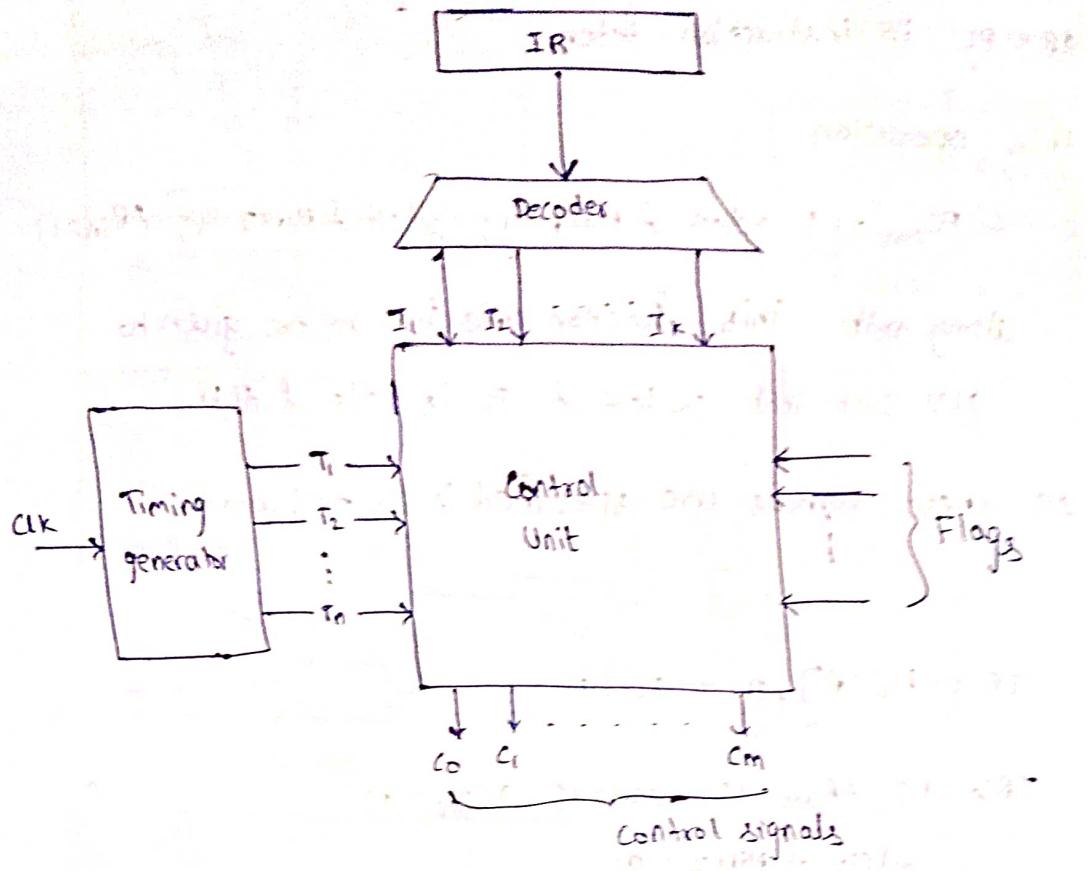
(ii) Microprogrammed CU.

(i) Hardwired CU:

Control logic is implemented with gates, FFs, decoders and other digital circuits.

Adv: Can be optimized to produce a faster mode of operation

Disadv: Updating is difficult cuz it requires rearranging wires.



We ~~can~~ create a table as shown below

| inst timer type | I ₁ | I ₂ | I ₃ | I ₄ | ... |
|--------------------|---|----------------|----------------|----------------|-----|
| T ₁ | C ₁ , C ₂ CS, CQ | -- | -- | -- | -- |
| T ₂ | CS, K ₆ | | | | |
| T ₃ | C ₁ , C ₃ | | | | |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |

no of rows = no of cycle in that type of instruction

Based on this table we build the hardware by obtaining logical function for each control signal.

→ Timing generator is reset after execution of every instruction.

Q10 Solve gate 2005 & on analysis of existing slide on Hardwired CU design

(ii) Microprogrammed Control Unit

Here all the possible control words are stored in memory.

This memory is known as control memory.

Based on req. specified control word is fetched.

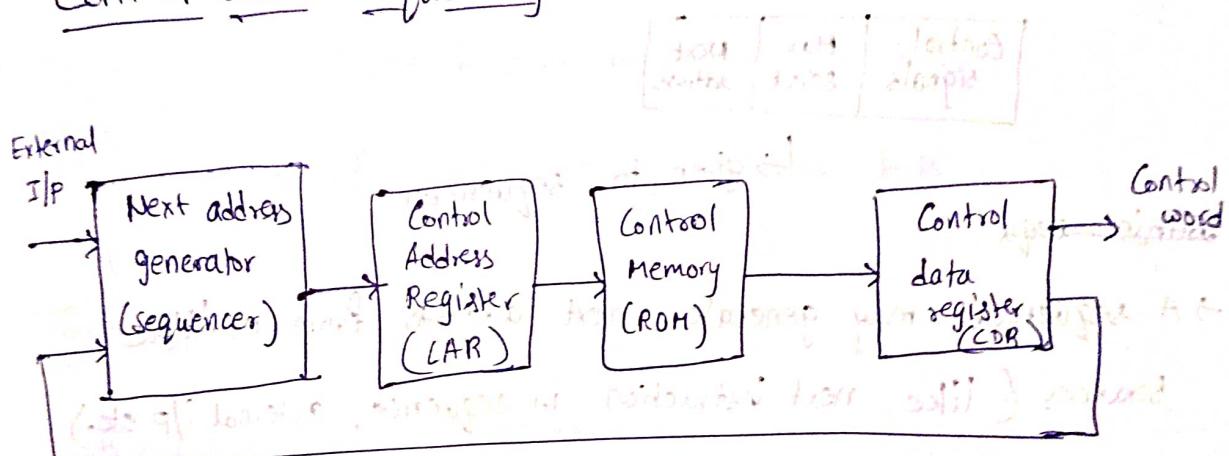
→ control logic is implemented with micro-programs.

Adv: Updating the control logic is easy.

Disadv: Slower than hardwired CU.

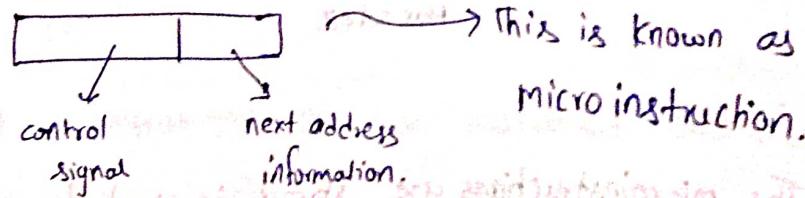
→ Microprogrammed CU is mainly used in complex highend computers.

Control word sequencing:

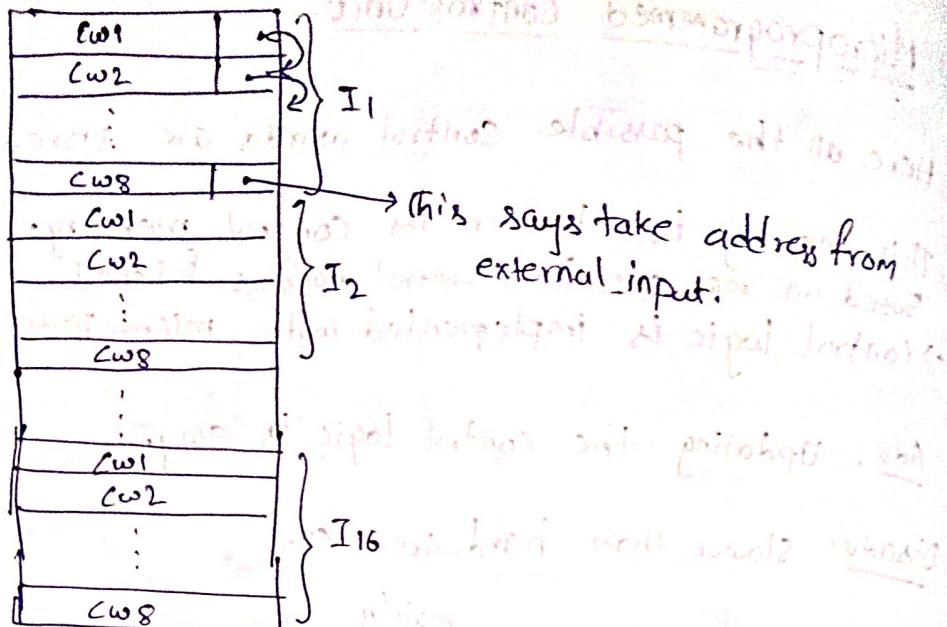


Each location

Each location of control memory is as shown below



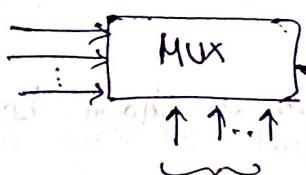
→ Assume a CPU has 16 different instructions, and each instruction has 8 μ-ops. The μ-ops are stored in control memory as shown



Standard format of micro-instruction:

| Control signals | Mux select | Next address |
|-----------------|------------|--------------|
|-----------------|------------|--------------|

- given to sequencer
- within seqn
- A sequencer may generate next address from multiple sources (like next instruction in sequence, external i/p etc.)
- Thus sequencer will have MUX logic



Q11
a-04

The microinstructions are stored in control memory of a processor have a width - - - - -

(Q2) A CPU has 64 distinct instructions and each instruction takes 8 micro-operation. Every micro instruction has 3 fields: control signals, mux select, address field. Let mux has 16 input. Assume there are 128 control signals. What is the size of control memory required?

Sol: no. of microinstructions to be stored = $64 \times 8 = 2^9$

| | | |
|-----------------|------------|----------|
| 128 | 4 | 9 |
| Control signals | Mux Select | next add |

\therefore size of each micro-instⁿ = 141 bits

\therefore size of control memory = $2^9 \times 141$ bits.

Types of micro-programmed Control Unit:

i) Horizontal micro programmed CU

ii) Vertical micro programmed CU

Horizontal:

- The CU that is discussed so far is called horizontal microprogrammed CU.
- In horizontal μ-programmed CU, we maintained one bit for each control signal.
- If there a large no of control signals and if a μ-inst needs only few signals then storing all the bits is waste of space.

Vertical μ -programmed CU:

→ Here we divide control signals into set of groups.

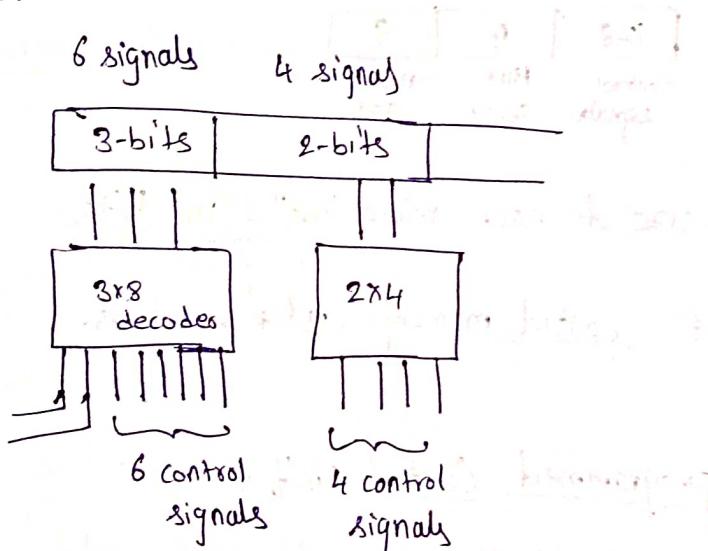
The groups are formed such that, for any micro-instruction

requires only one control signal from ~~that~~ group enabled

→ Group sizes may be unequal.

→ These bits of each group are encoded and thus size of control unit is reduced.

→ Consider,



These are 2 types of vertical μ -programmed CU

(By default we consider this)

→ Always exactly

one signal is

enabled from each group

→ Almost 1 signal.

is enabled from each

→ Here for a 4 signal

group we encode into
in 2 bits

→ Here for a 4 signal

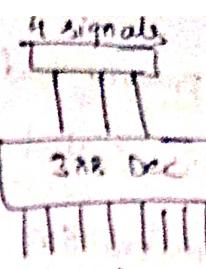
group we encode into

3 bits

(∴ we have 5 combinations)

groups.
micro-instr
group enabled

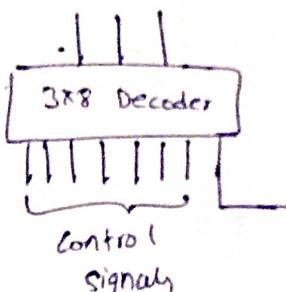
thus size



Control signals

For 100 to 111 all control signals will be disabled.

→ If we have 3 signals
storage as 8x8 memory square then we need to
represent 8 combinations
∴ 3 bits required



Control
signals

For 111 all control signals are
disabled.

∴ 2³ programming words

→ Vertical μ -programmed CU is slower (due to decoders)

Compared to horizontal μ -programmed CU.

Notes:

→ Sometimes there may be some signals which cannot be grouped. So these signals will be stored in horizontal manner.

1 signal
from each
group

on 4 signal
encode into
5 combinations

(Q12) A microprogrammed CU is required to generate a total of 25 control signals. Assume that during any microinstruction, at most 2 control signals are active. Minimum no. of bits required in the control word to generate the required control signals will be?

Ans: 5 bits

Here we don't know how to group.

So we create 22 groups each of size 25 signals.

No. of combination for each group = $25+1=26$

∴ it requires 5 bits

| control word | SS signal | SC signal |
|----------------------------------|-----------|-----------|
| 11111111111111111111111111111111 | G_1 | G_2 |
| | 5 | 5 |

∴ Minimum bits req = 10

Nano-programmed CU:

→ Some instructions may have common microinst.

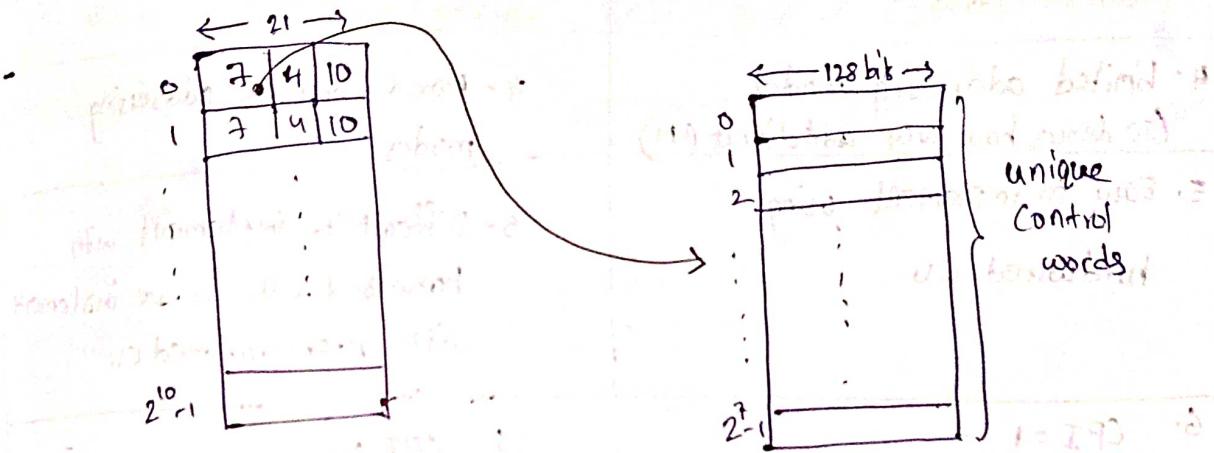
However we can still store them separately in microprogram control unit.

→ In nano programmed CU we store all unique microinst. separately in memory.

any i total
achieve.
words to generate.

→ Assume we have 128 unique micro instⁿ and 128 control signals. Also assume we have 16 instⁿ where each instⁿ has 64 M.op. Then organization of nano programmed CU is as shown below:

no of micro instructions = $16 \times 64 = 2^{10}$



∴ size: $2^7 \times 2^7 = 2^{14}$ bits

Thus in this a way memory required is less and flexibility is also good.

(Q13) Consider a μ-programmed CU which has to support 32 no of instⁿs - For each instⁿ execution, CU generates a sequence of 64 control words. Each microinstruction contains 3 fields: 118 control signals to support horizontal CU, a MUX select field to select one of 18 inputs and a next address field.

The size of control memory used is?

Sol:

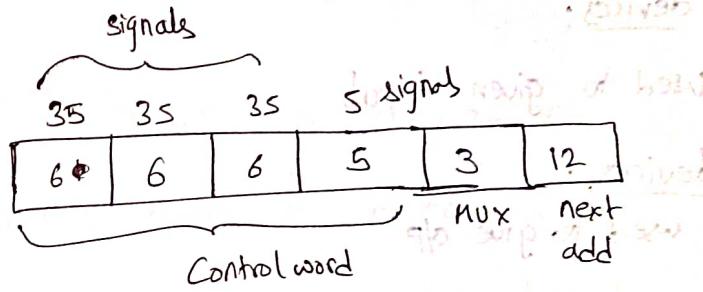
$2^{11} \cdot 2^5 + 2^6 (118 + 3 + 11) = 2^{11} * (132)$ bits

RISC vs CISC

| RISC | CISC |
|---|---|
| 1. Less no of instns | 1. More no of instns |
| 2. Fixed length instructions | 2. Variable length instructions |
| 3. Simple instns | 3. Complex instns |
| 4. Limited addressing modes (IO devices have only direct AM) | 4. More & complex addressing modes |
| 5. Easy to implement using hardwired CU | 5. Difficult to implement with hardwired CU. So we implement with programmed CU |
| 6. CPI = 1 | 6. CPI > 1 |
| 7. Register to register arithmetic operations only | 7. reg to mem or mem to reg arithmetic operations possible |
| 8. More no of registers are available | 8. Less no of registers are available. |
| 9. Provides register window (Efficient parameter passing) | 9. No register window |
| 10. Compiler design is easy Programmer overhead is more | 10. Compiler design is complex Programmer overhead is less |
| 11. More compatible for instruction pipeline | 11. Less compatible for instruction pipeline. |

(Q14) Design of a vertical microprogrammed CU, requiring to generate 40 signals. Out of first 35, only 3 signals can be active at a time and for remaining 5 anyone can be active any time. The micro instruction of the CU stores control signal information along with 3-bit mux select and 12 bit address field. size of control memory required is _____

Sol:



\therefore size of each micro instr = 38 bits

no of micro instrs = 2^{12}

\therefore size of control memory = $2^{12} \times 38$ bits



Peripheral Devices

IO Organization

Peripheral Devices:

→ Devices that are connected to the processor externally are known as peripherals (apart from main memory).

There are 3 types peripheral devices

i) Input devices:

Used to give input

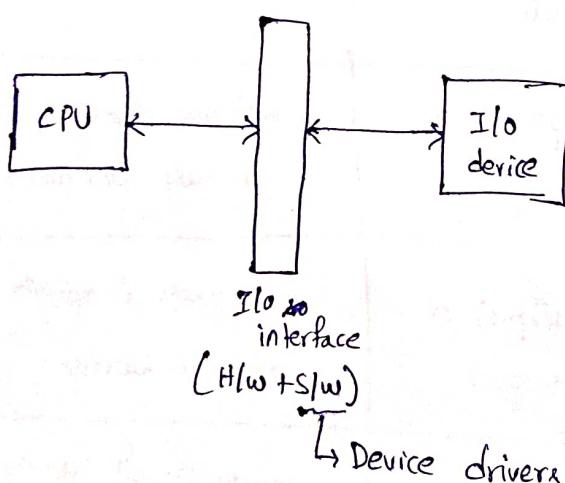
ii) O/P devices:

Used to give o/p

iii) Storage devices:

Used for storage purposes

→ CPU cannot access IO directly. It connected through I/O interface.



→ Interface is needed for signal conversion.

(i) providing synchronization b/w CPU & I/O

(ii) Data format conversion

(iii) Management & Control (so that a peripheral doesn't disturb operation of others)

IO vs Memory Buses

- Every IO device has an address so that CPU uniquely accesses an IO device.
- ∴ ~~we need address bus~~
- for IO bus also we have address bus, data bus, control bus.

→ ~~IO & memory have separate buses.~~

~~But they have common address~~

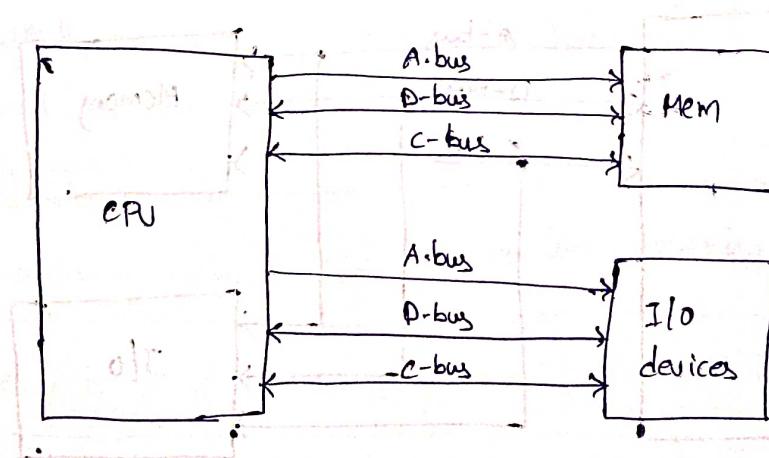
There are 3 types of bus configurations

(i) Separate buses for memory & IO

(ii) Common address and data bus but different control buses.

(iii) All are common buses

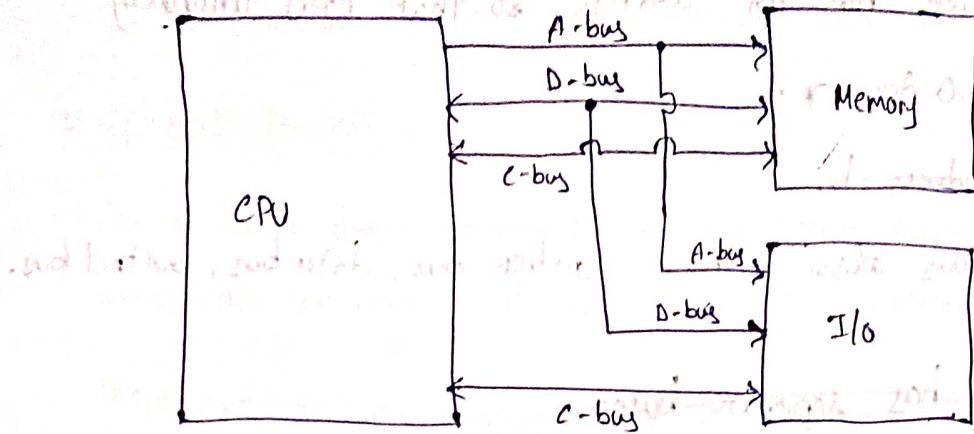
i) Separate buses for memory & IO:



→ This organization is costly

→ This is generally used in computers that provide separate I/O processor. Here memory communicates with both CPU and IOP using memory bus.

(ii) Common address & Data buses:

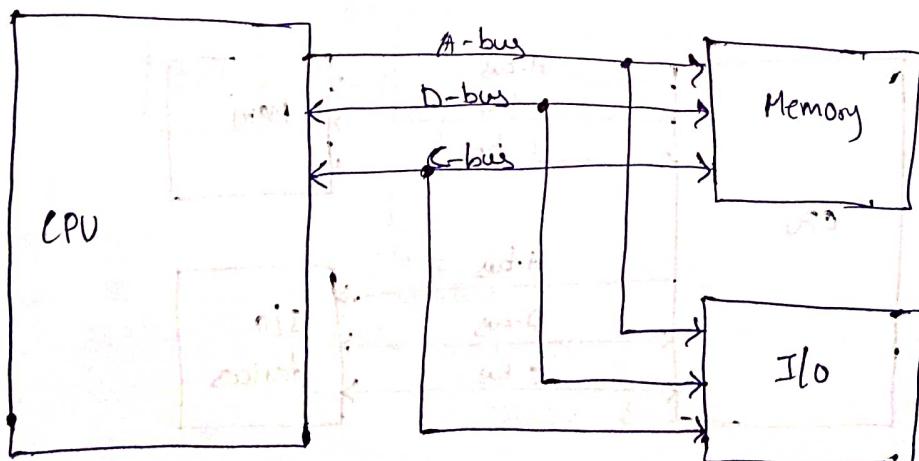


→ Since Address and data buses are same, to resolve the conflict regarding which device CPU wants access, the control signals (read & write) are used.

→ This configuration is called I/O mapped I/O or port-mapped I/O / Isolated I/O

→ This is mostly used

(iii) Common address, data, control buses:



→ Assume -

→ In this configuration I/O adds device a given address from memory.

→ Here computers have only one set of read & write signals.

Assume a system with memory 8 bytes & 2 I/O devices d_1, d_2



Assume if I/O devices are ~~allocated~~ given address as

$d_1 \rightarrow 010$ & $d_2 \rightarrow 101$

Here interface register are mapped to memory

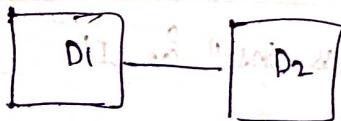
This configuration is called memory mapped I/O

→ here interface register data is stored in memory. So load & store can be used to access ~~I/O~~ I/O devices also.

| Memory mapped I/O | I/O mapped I/O |
|--|---|
| (i) I/O devices don't have separate address space | (i) I/O have own address space |
| (ii) Some memory is wasted | (ii) No wastage |
| (iii) All the memory access inst ⁿ s can be used to access I/O devices also | (iii) I/O access inst ⁿ s and memory access inst ⁿ s are different. |
| (iv) More inst ⁿ for I/O access (\because it has all possible memory access inst ⁿ s) | (iv) Less inst ⁿ for I/O |
| (v) \therefore more addressing modes for I/O | (v) less addressing modes for I/O |
| (vi) Can have more no. of I/O devices | (vi) It has less no. of I/O devices |

Asynchronous Data Transfer (serial)

- If two devices involved in data transfer operate under same clock then such data transfer is called synchronous data transfer.
- In synchronous data transfer, the clock rate is set based on slower devices and due to this performance falls.
- Thus we use asynchronous data transfer in which we have separate clocks for devices.
- Since there is no synchronization provided by clock, we provide extra external synchronization.
This synchronization is provided by I/O interface.
- There are two types of data transfers
 - (i) Serial : single line for data transfer.
 - (ii) Parallel : multiple lines are used



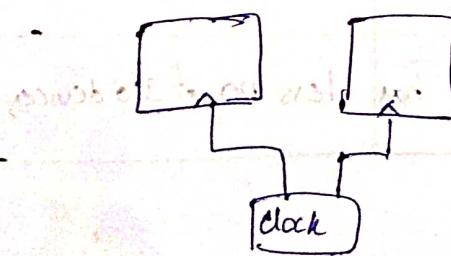
Serial

(less cost)

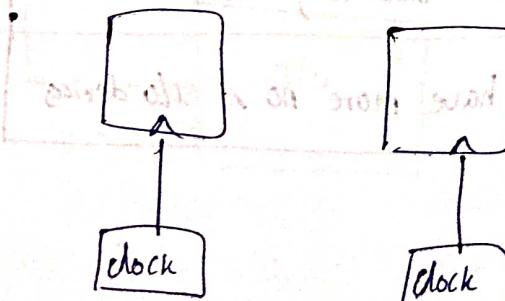


Parallel

(costly)



Synchronous



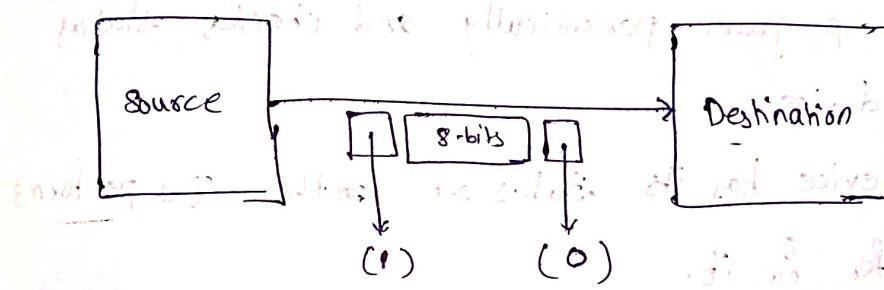
Asynchronous.

→ Let us discuss asynchronous serial data transfer

Note:

* whenever data is sent on a line, the last bit of data sent will remain on communication line until new data is sent.

→ In this transfer we provide synchronization by adding a start bit & a stop bit for an 8-bit data.



Here as long as 1 is present on data line, it means no data is being sent.

Once a zero is seen, destination understands that data is being sent.

→ Here for every 8-bits, we send 10 bits.

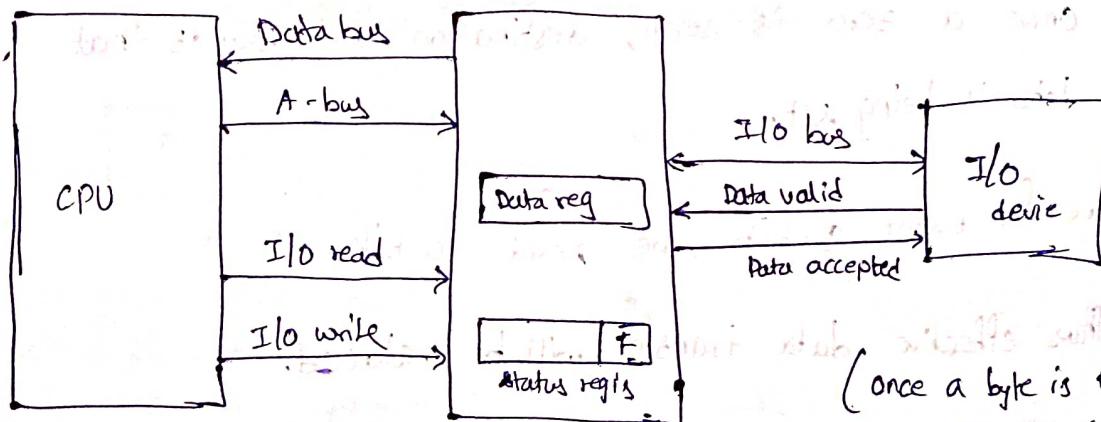
Thus effective data transfer will be reduced.

Q15
Ans
Gray

How many 8-bit characters

Mode of transfer:

- It tells about the way in which data can be transferred from an I/O device.
- (i) Programmed I/O:
- There is no provision through which I/O can inform CPU about data transfer.
- I/O sets its status & waits.
- CPU runs a program periodically and checks status of each device.
- If any device has its status set the CPU performs data transfer for it.
- This works on the approach of polling.

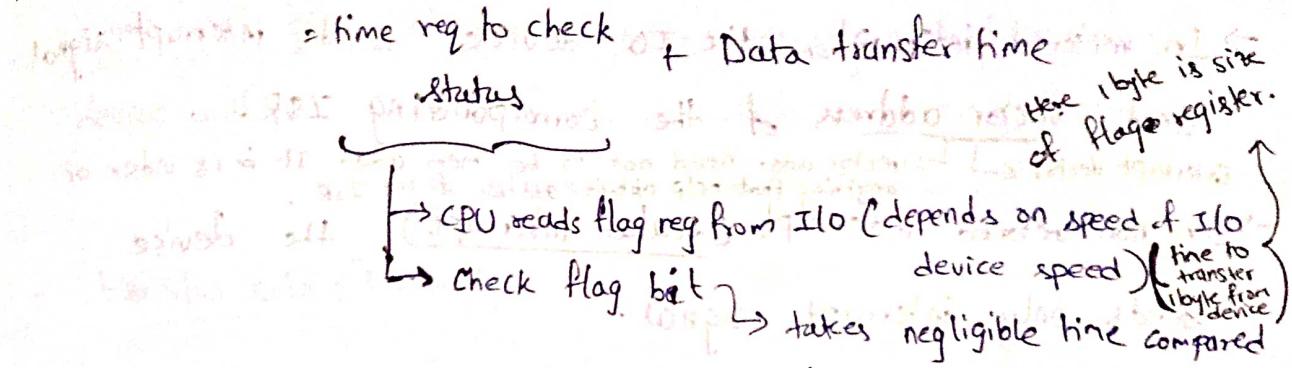


(once a byte is transferred from I/O to Interface, data accepted is enabled by Interface. This is not disabled until the byte is read by CPU)

- I/O device writes data into data register of the interface and sets its flag saying it is ready for I/O.

** time req in programmed I/O

**



(iii) Interrupt Initiated I/O

→ I/O device has provision to interrupt CPU abt communication.

→ When CPU receives interrupt: (from I/O device)

(i) It completes execution of current inst

(ii) Saves the status (PC, PSW etc) of current process onto the stack

(iii) Branches to service the interrupt

(iv) Resumes the previous process by taking out the values from stack.

→ Memory never generates interrupt.

ISR:

→ ISR is set insts executed to service an ~~int~~ interrupt.

→ All the ISRs are part of device drivers.

Vectored vs Non-Vectored interrupt:

- In vectored interrupt, the IO device sends interrupt signal and vector address of the corresponding ISR.
- Interrupt Vector → vector add. need not to be mem add. It is index or anything that help obtain address of the ISR.
- In non-vectored interrupt (scalar interrupt), the device sends only interrupt signal.

Here CPU runs a default service routine through which it obtains the location of the corresponding ISR.

Maskable interrupt vs Non-Maskable interrupt:

- Maskable interrupt: interrupt that can be rejected if required.
- Non-maskable interrupt: interrupt that must be serviced.
- When an interrupt is masked, CPU may completely reject it or can keep it pending.

Internal interrupt vs External interrupt:

External interrupt: If any device generates intr. then it is called external interrupt.

Internal interrupt: Sometimes CPU may encounter an unexpected error while execution. This is known as internal interrupt. Eg: page fault.

Show
systems
calls

During this kind of interrupt, CPU stops execution of current instn, solves the problem & restarts the instn.

Simultaneous Interrupts:

interrupt signal
is index or device

→ when CPU gets interrupts simultaneously, higher priority device will be serviced first.
i.e., priority based interrupt handling

HW soln (Polling)

→ A SW runs first before than ISR or default service routine. Using this CPU finds the higher priority device.

→ This is based on polling.

CPU checks higher priority device first if it has generated interrupt. Later it checks next higher

priority device and so on...

→ Disadv: If there are more ~~interrupts~~ devices, time req to poll could exceed the time

(ii) H/W solns: available to service I/O device.

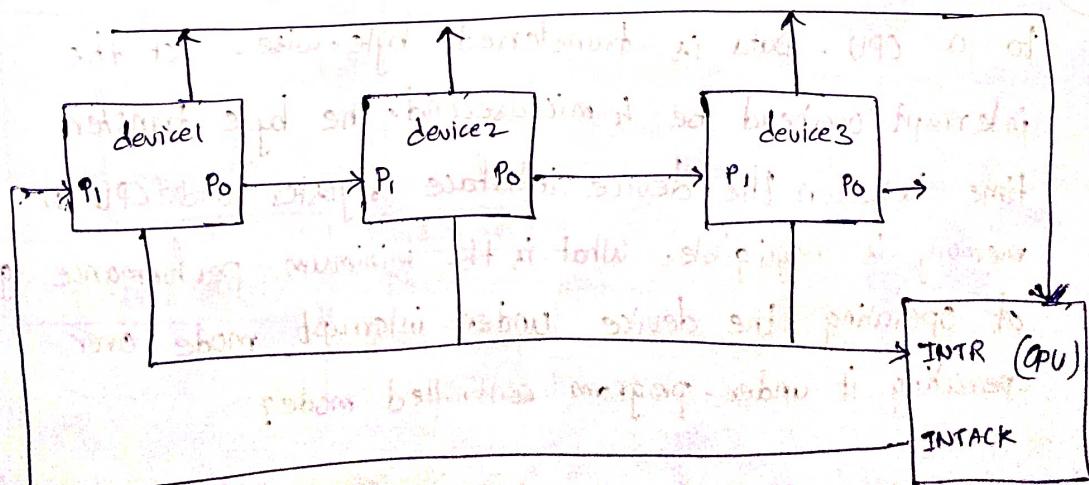
Serial. soln (Daisy Chaining)

parallel soln

High priority are assigned to interrupts which, if not serviced, would lead to serious consequences.

Also higher priority is given to faster data transfer (Eg: hard disk) and low priority for slow devices (Eg: keyboard)

Daisy Chaining:



→ CPU receives intr and understands interrupt occurs.
→ Now it sent intack to device first.

Device 1 will accept ack if it has interrupted and if sends vector address. and it sends '0' to device 2.

Otherwise

Otherwise Device 1 will send 1 to device 2; and this process continues until ~~device is found~~ until INTR is set to 0.

After this CPU will disable INTACK.

→ ~~These devices~~ is selected such that

→ In daisy chaining the device that is electrically nearest gets highest priority.

→ Time required in interrupt IO:

interrupt overhead + service routine

- execution of current instruction
- pushing data onto stack
- default service routine etc.

Data transfer rate refers to internal speed of device. The rate at which data is produced or processed.

(Q16) A device with data transfer rate 10 KB/sec is connected to a CPU. Data is transferred byte-wise. Let the interrupt overhead be 4 microsecond. The byte transfer time between the device interface register and CPU or memory is negligible. What is the minimum performance gain of operating the device under interrupt mode over operating it under program controlled mode?

80) If average processing time for each I/O is 100 μs
 time req in programmed I/O is
 = overhead + actual byte data read + transfer time
 = $100 \mu\text{s} + 0$
 $\approx 100 \mu\text{s}$
 time req in interrupt I/O
 = overhead + service time
 = 4 to
 ≈ 4
 performance gain $\approx \frac{100}{4} \approx 25$

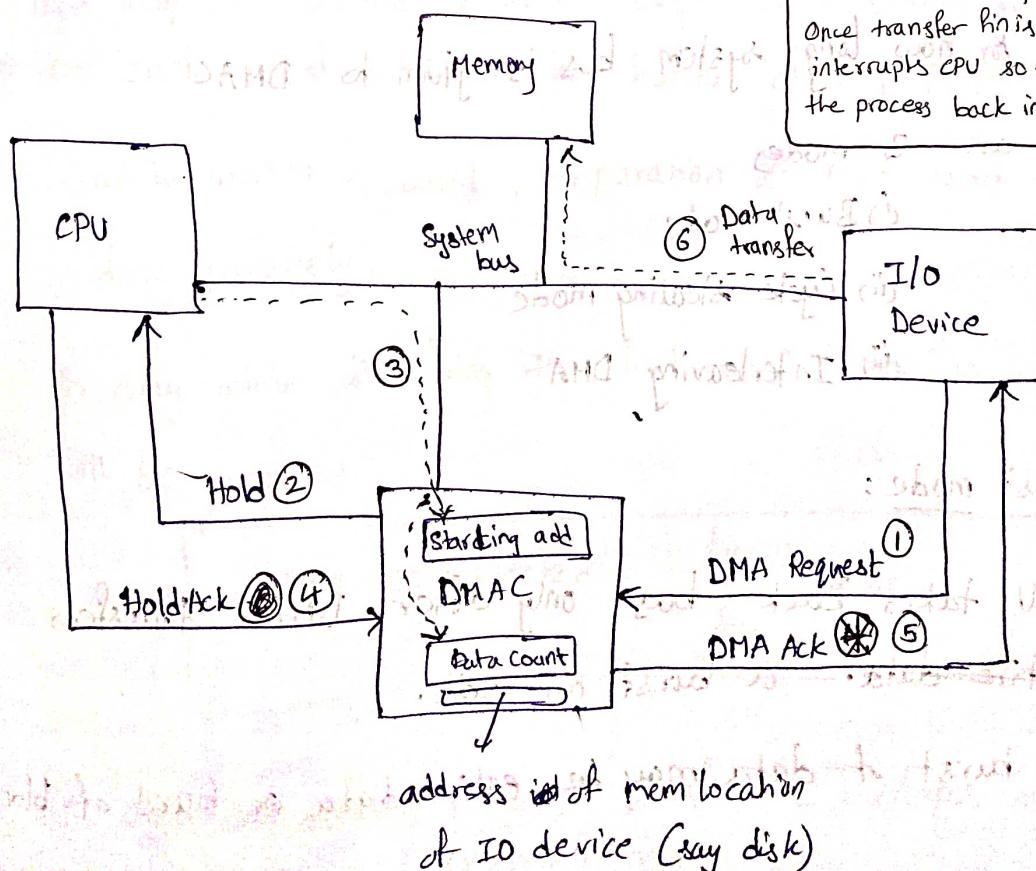
(Ans consider 100 μs. SAMD action of each multiple interrupt for
 interrupt handling)

Direct Memory Access

→ DMA is a technique that enables data transfer b/w I/O and memory without CPU intervention.

→ H/w required for DMA is DMAC.

When DMA transfer is initiated
 OS puts the corresponding process in block state and may schedule another process.
 Once transfer finishes, DMAC interrupts CPU so that OS put the process back in ready state



- Starting address is address from which memory access has to be done
- Data count is no. of memory locations to be accessed from starting address. (words/bytes)
- Those two information is sent by CPU to DMAC
- After transfer of each byte/word data count value is reduced by 1. and starting address is incremented by 1.
- The data transfer stops once data count reaches 0.
- When system bus is with DMAC, CPU performs only those operations that doesn't require system bus. This means mostly CPU will be blocked.
(Cuz even instruction fetch requires memory)

Modes of DMA Transfer:

→ Tells for how long system bus is given to DMAC

→ There are - 3 modes

(i) Burst mode

(ii) Cycle stealing mode

(iii) Interleaving DMA

i) Burst mode:

→ CPU takes back bus only after DMAC transfers ~~entire data~~. a burst of data.

→ A burst of data may be entire data or burst of blocks.

i) Cycle Stealing mode:

- Preferable when IO-device is slow in preparing data
- So, here the idea is to give system bus to DMA for short time when data is already with IO device.
- Here system bus is given to IO for DMA ≈ 1 cycle in which the prepared word is transferred to memory.
- During preparation of word CPU is not blocked.
- CPU is blocked during data transfer only.
- Let time required to prepare data (t_p) = t_x
- Time required to transfer data = t_y
- $\therefore \% \text{ time CPU blocked} = \frac{t_y}{t_x + t_y} * 100$ (Burst mode)

ii) Cycle Stealing Mode

- There may be second possible case (Cycle stealing mode)
- Here each word is stored in buffer after preparation.
- ∴ while writing a word, preparation of next word is done parallelly.
- So only while writing the last word, ~~no word~~ will be prepared.

$$\% \text{ time CPU blocked} = \frac{t_y}{t_x} * 100$$

In cycle stealing mode, transfer time overlaps with preparation time

(iii) Interleaving DMA:

Whenever CPU does not require system bus (doing internal work) then only control of the buses will be given to DMAC.
→ CPU will not be blocked in this.

Note:

- Among all modes, data transfer will be fast in burst mode, second fastest is cycle stealing and finally interleaving DMA
- During DMA transfer, DMA sends read/write signals to memory.
- DMAC is a special purpose processor which controls data transfer b/w memory and I/O.
- ~~For~~ DMA HLDA may be enabled even between the execution of instn.

∴ DMA service can't be done during execution of instn.

Interrupt service is done only after end of instn.

→ DMA & CPU can't work parallelly. Then how DMA improves performance over interrupt I/O.

In interrupt I/O CPU state has to be switched.

In DMA no need to switch the state.

∴ DMA gives better performance for ^{data} transfer b/w I/O & memory.

(doing internal
given to DMAc.)

burst mode,
interleaving DMA

signals to

Controls

even between

instruction.

f instn.

DMA

switched.

transfer b/w

- However for ~~short~~ data transfer b/w CPU & I/O interrupt is better.
- Also we use DMA for large ~~data~~ transfers and interrupt I/O for short data transfers.

16/11/20

- (Q1) Consider a device operating on 2Mbps speed and transferring the data to memory using cycle stealing mode of DMA. If it takes 2 μs to transfer 16 bytes data to memory when it is ready/prepared. Then % time CPU blocked due to DMA is

Sol:

$$2 \text{ MBPS} \rightarrow 2 \times 10^6 \text{ bytes/sec}$$

$$16 \text{ bytes} \rightarrow \frac{1}{2 \times 10^6} \times 16 \text{ sec}$$

(Ans. 3.2 ms and 16 bytes $\rightarrow 2 \times 10^6$ bytes/second)

So transfer time = $16 \times 10^{-6} \text{ sec} = 16 \mu\text{sec}$

Here transfer time overlaps with preparation time

$$\therefore \% \text{ time CPU blocked} = \frac{16}{8} \times 100 = 25\%$$

Q18
G-16

The size of data count register of a DMA controller is 16 bits. The processor needs to transfer a file of 29154 kB from disk to MM. The memory is byte addressable. The min no of times DMA needs to get control of system bus from processor for this transfer.

Sol:

DC \rightarrow 16 bits

Memory is byte addressable.

so max of $2^{16}-1$ bytes can be transferred at once

$$\text{no of times control required} = \frac{29154 \text{ kB}}{2^{16}-1} = 455.55$$

so, which has wrong because no fifferent control is required for each of 16 bits & above value is 456 times

so, total number of bytes transferred is $456 \times 16 = 7296$

Add Additional Info:

In asynchronous data transfer, synchronization can be provided in 3 ways

(i) strobe signal (generally used b/w CPU & I/O)

(ii) Handshaking signal (generally used b/w I/O & interface unit)

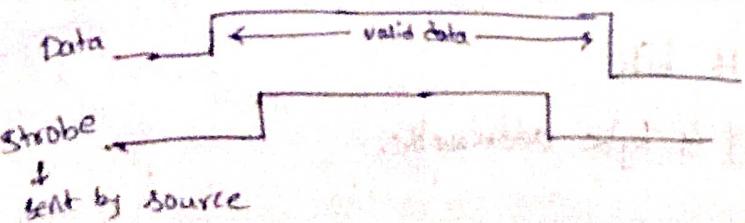
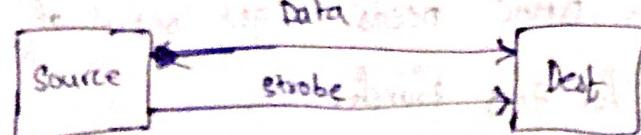
(iii) asynchronous serial data transfer

(i) strobe signal:

→ Here we maintain a control line to time each transfer

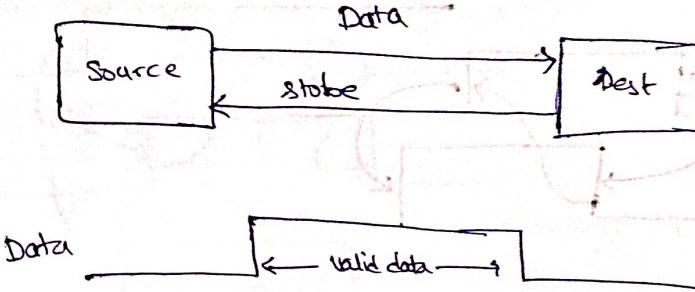
through which source sends strobe signal to destination.

Source initiated transfer:



The source place strobe signal after placing data.

Destination initiated transfer:



In both the cases, destination often uses (not necessary) falling edge of strobe to trigger destination registers.

- whenever strobe signal is disabled implies that data on the bus is not valid.

(ii) Handshaking:

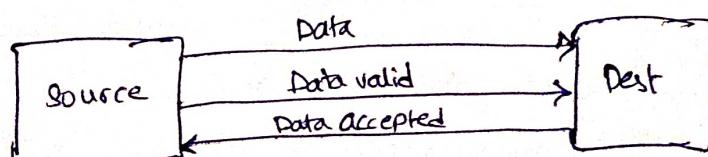
- The problem with strobe signal is that source never knows whether destination has received the data or not. There is no acknowledgement.

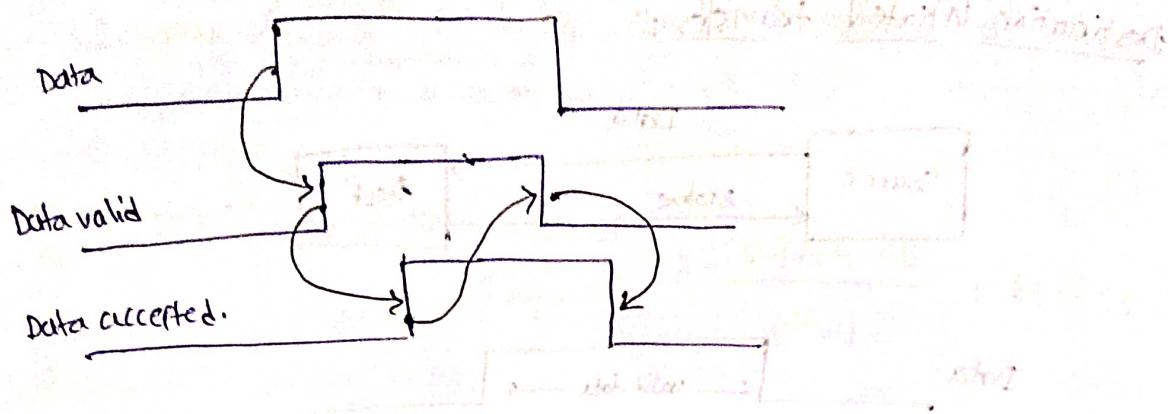
- Here we use two control signals

Data valid : says data on bus is valid (sent by source)

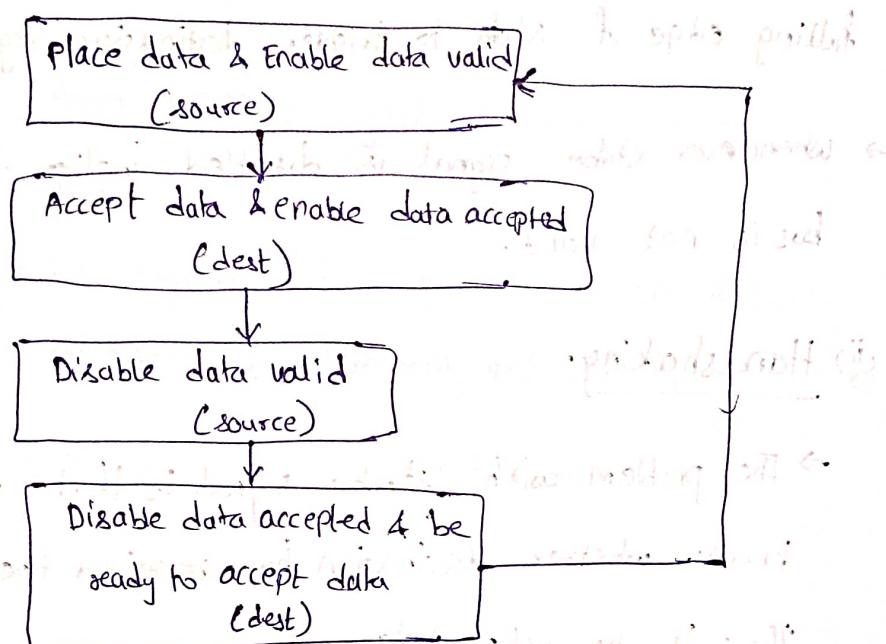
Data accepted : says data is received (sent by destination) and turned off only after it is ready for next data

Source initiated transfer:





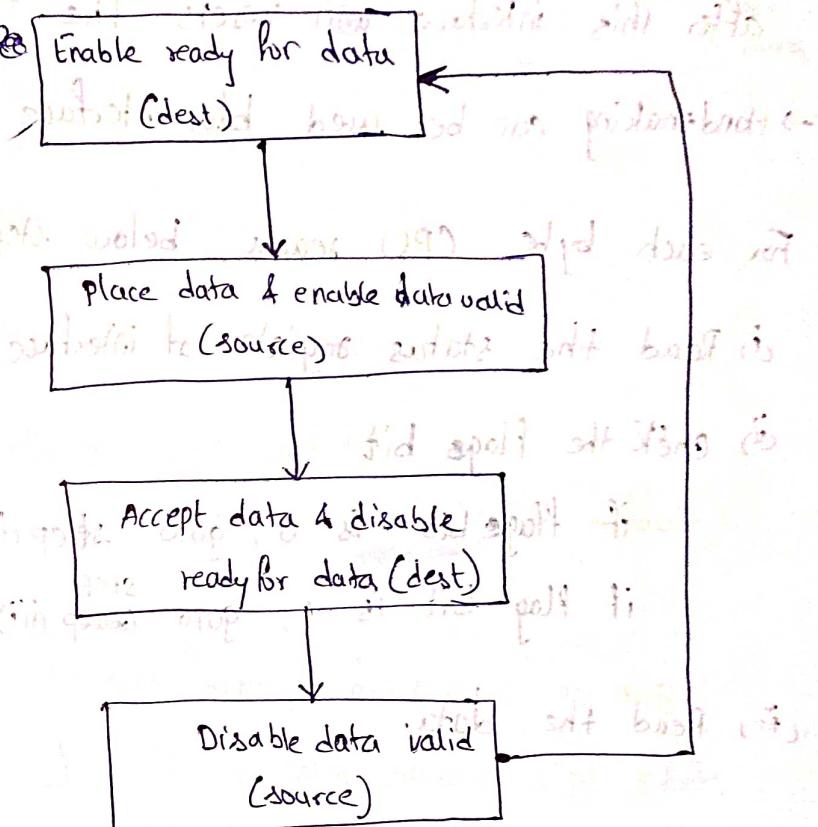
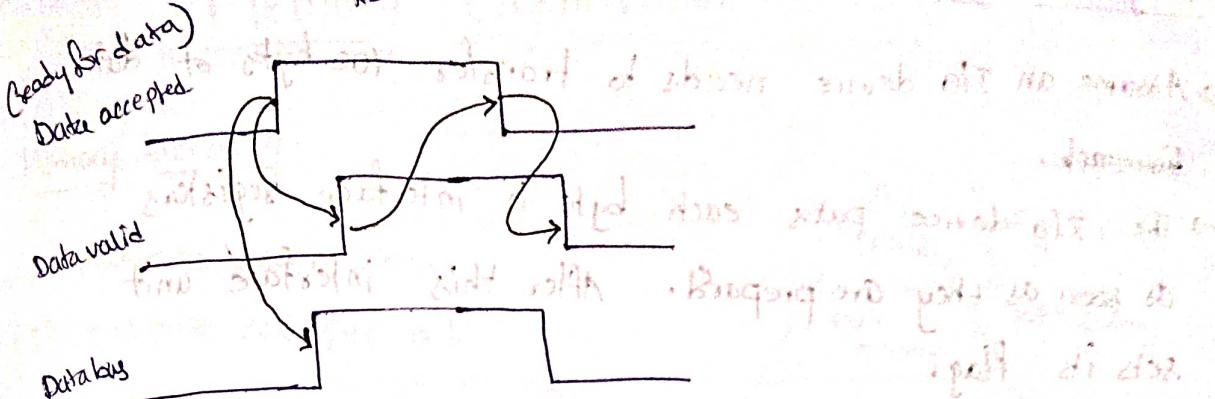
→ Source sends next data only after destination disables data accepted. (By disabling data accepted destination says it is willing to receive data).



Destination initiated transfer:

→ Here we see 'data accepted' line as 'ready for data'.

→ So destination enables 'ready for data' (data accepted) to start the data transfer.



→ If one of the two devices is faulty, data transfer will

not be completed. This kind of error can be detected by means of a timeout mechanism.

Programmed I/O

→ Assume an I/O device needs to transfer 100 bytes of data.

~~for each~~

→ The I/O device puts each byte in interface registers as soon as they are prepared. After this interface unit sets its flag.

→ The CPU sees the flag and reads the data and after this interface unit resets the flag.

→ Handshaking can be used b/w interface & I/O device.

For each byte CPU reads, below steps are performed

(i) Read the status register of interface unit

(ii) Check the flag bit

 if flag-bit is 0, goto step (i)

 if flag-bit is 1, goto step (iii)

(iii) Read the data.

→ Here it is clearly understood that, CPU loops as long

a data is being prepared and made available in interface register.

→ So the time depends on data transfer rate of device.

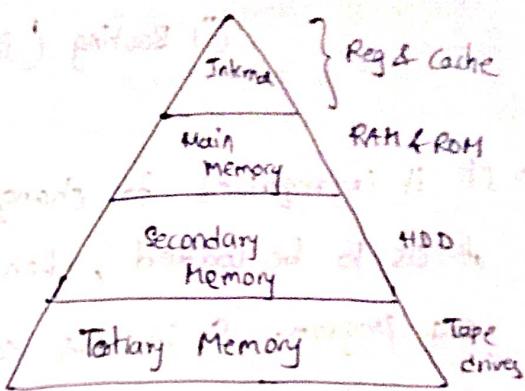
If device needs 't' time to produce data and make it available in registers, then CPU would also take around 't' time to finish the I/O operation.

17/11/20

Memory Organization

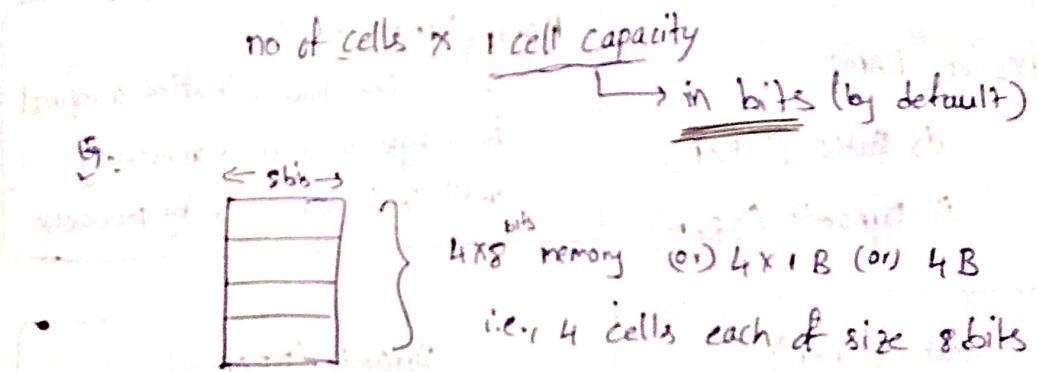
Memory hierarchy:

- maximize access speed
- minimize cost per bit.



Memory Representation:

- Memory is represented by grouping of memory cells into data blocks.

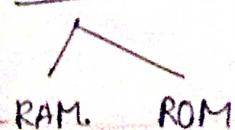


If it is given like 4KB memory and need to find its size
split 4 into 2 and 2 into 1024

Main Memory:

Defn: It is memory that is used to store currently running programs & their data.

types:



(ROM: Non-volatile)
RAM: Volatile

- ROM:
- Initially when CPU is turned on it gets initial program from ROM. This program does 2 things:
 - i) Hardware check (POST (Power-on Self-test))
 - ii) Booting (Bootstrap program)

- If it is required to change the initial set of programs of OS to be loaded, then we need to change bootstrap program. So in this case we store bootstrap program on disk and store bootstrap loader in the ROM.

RAM:

- stores current running program & their data along with the OS.

types of RAM:

(i) Static (SRAM)

(ii) Dynamic (DRAM)

Memory Latency:

It is time b/w initiating a request for a byte or a word in memory until it is retrieved by processor

| Static (SRAM) | Dynamic (DRAM) |
|--------------------------|---|
| → made up of Flip-Flops | → made up of capacitors |
| → No refresh is required | • Data is stored in the form of electric charges (Charges tend to discharge and data may lose) |
| → faster & costlier | → Periodic refresh is required |
| → Eg: Cache memory | → slower & less costlier Eg: Main memory |

Static

- idle power consumption is very less
- Operational power consumption is high comparatively

Dynamic

- idle power consumption is more
- operational power consumption is less

programs

bootstrap

trap program
& ROM.

with the

bring a request
in memory
by processor

it's
form

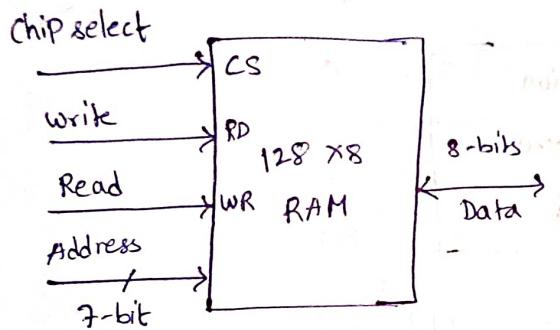
charge and
change

required

18/11/20

RAM Chip: RAM is made from dynamic memory chip

along with it 2 address lines are also



Every computer has

at least 2 chips

(i) RAM (ii) ROM

CS is used to select
that chip that we need.

| CS | Read | Write | Operation |
|----|------|-------|--------------|
| 0 | X | X | No operation |
| 1 | 0 | 0 | No operation |
| 1 | 0 | 1 | write |
| 1 | 1 | X | Read |

→ no of pins needed for this RAM chip

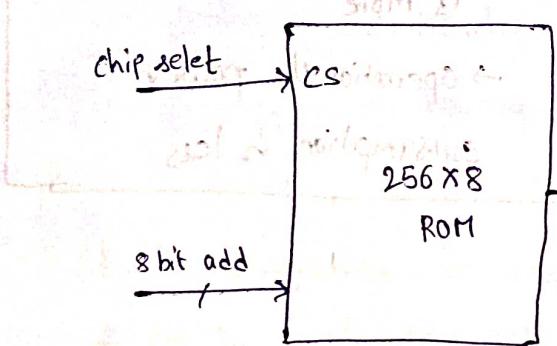
$$(Add) + (\text{data}) + (\text{CS}) + (\text{RD}) + (\text{WR}) = 18$$

data bus is not 8 bits

data bus is 8 bits

according to what we

ROM chip:



The only operation performed on ROM is read.

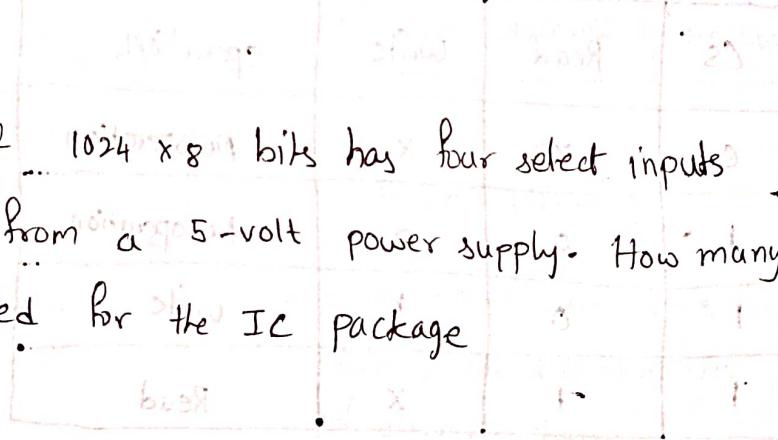
So we don't have RD & WR signals.

| CS signal | operation |
|-----------|--------------|
| 0 | no operation |
| 1 | read |

$$\rightarrow \text{No of pins} = 8 \text{ (add)} + 8 \text{ (data)} + 1 \text{ (CS)} = 17$$

(Q19) A ROM chip of 1024×8 bits has four select inputs and operates from a 5-volt power supply. How many pins are needed for the IC package

Sol:



Add $\rightarrow 10$

Data $\rightarrow 8$

CS $\rightarrow 4$

Power $\rightarrow 1$

whenever there is power (we must consider a chip for ground also)

\therefore ground $\rightarrow 1$

\Rightarrow total pins = 24

Multiple Chip Support:

Total memory capacity = no of chips * capacity of each chip

→ Assume we have two 128x8 bit chips;

from this we can build 256x8 bit RAM
or 128x16 bits RAM

add size of each chip = 7 bits

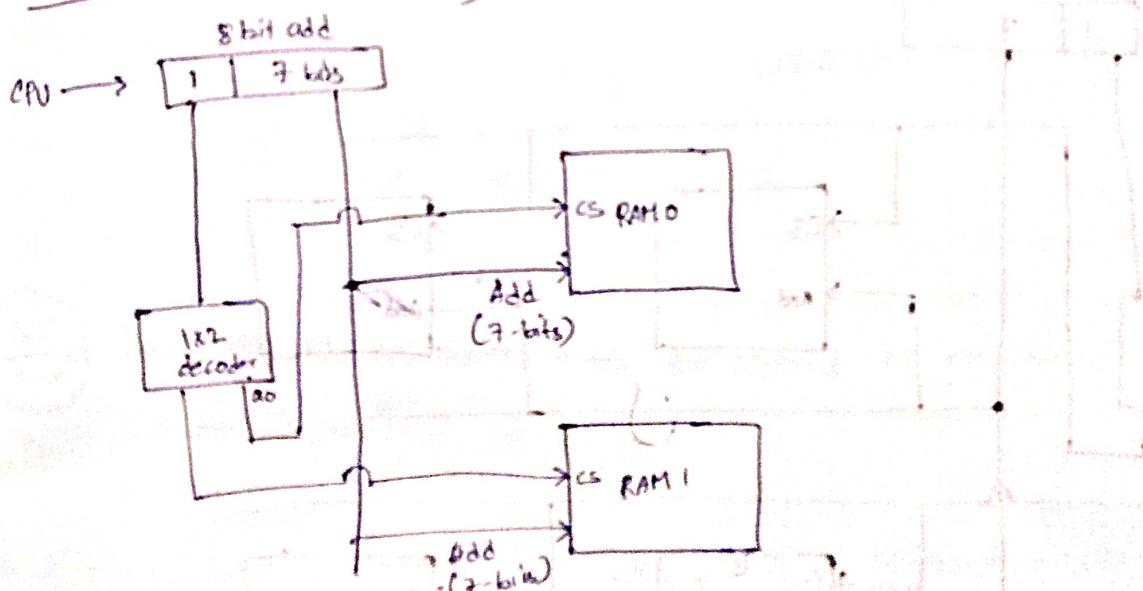
but CPU sees MM of capacity 256x8 bits

∴ CPU generates 8 bit address

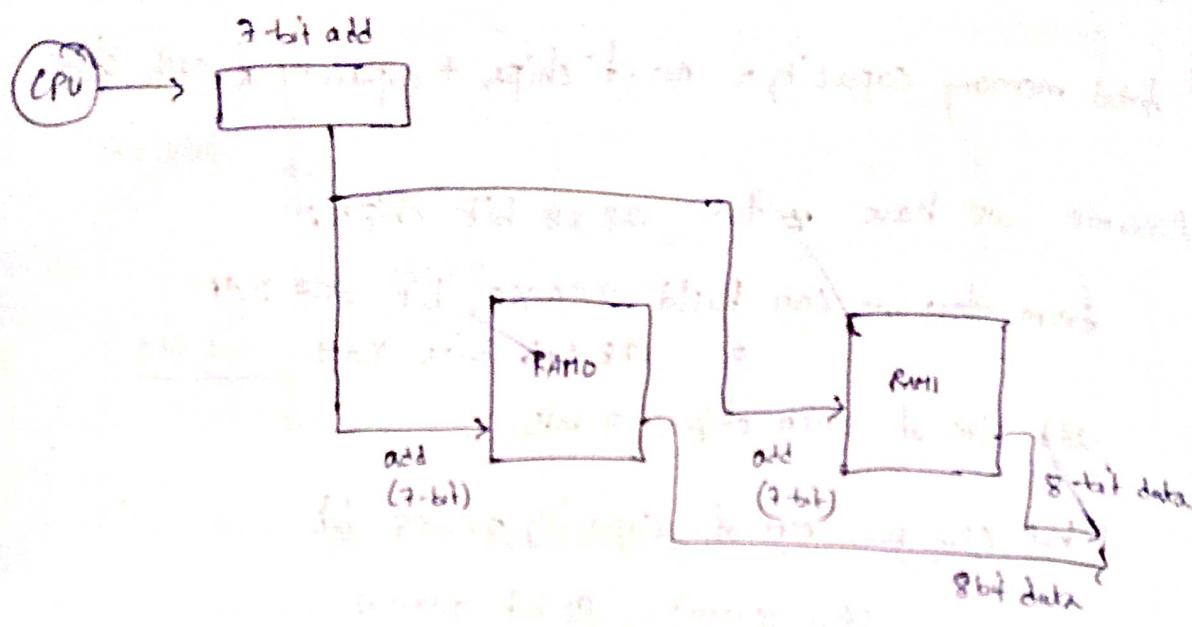
| | | |
|-----|-----------------|------------|
| 0 | 0 0 0 0 0 0 0 0 | chip RAM 0 |
| 1 | 1 1 1 1 1 1 1 1 | |
| 127 | 0 1 1 1 1 1 1 1 | chip RAM 1 |
| 128 | 1 0 0 0 0 0 0 0 | |
| 255 | 1 1 1 1 1 1 1 1 | |
| 256 | 0 0 0 0 0 0 0 0 | |

Deciding which RAM to access is done using MSB bit.

256x8 ; (vertical arrangement)

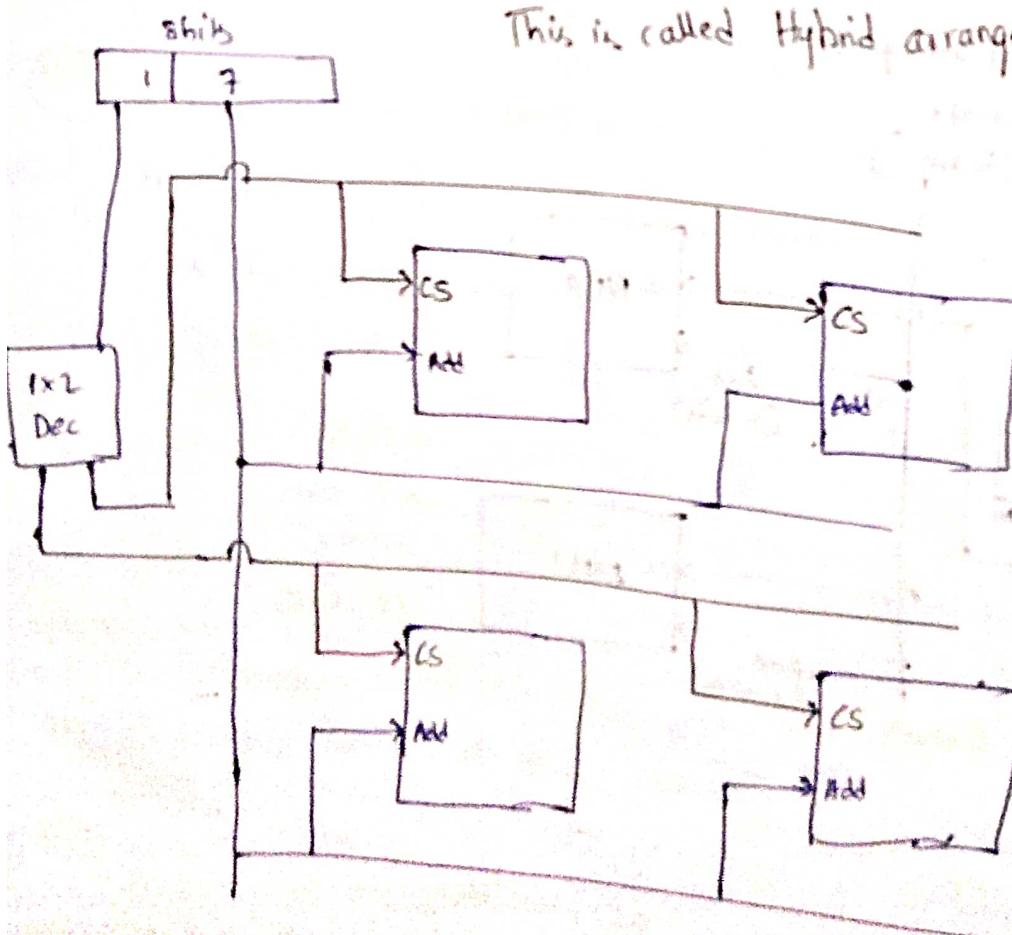


128x16: (Horizontal arrangement)



- If required addresses are more, then we use vertical arrangement.
- If data required at each memory location is more then we use horizontal arrangement.
- Now consider building 256×16 from 8×8 chips;

This is called Hybrid arrangement



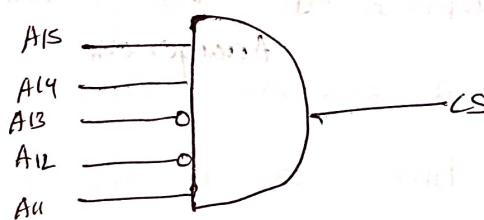
Q.20
6-19

The chip select logic for a certain DRAM chip in a memory system design is shown below. Assume that

The memory system has 16 address lines denoted by A₁₅ to A₀

What is range of address (in hexadecimal) of the memory system that can get enabled by the chip select (CS) signal

- a) E800 to CFFF
- b) CA00 to CAFF
- c) C800 to C8FF
- d) DA00 to DFFF



So, Address bits A₁₅ A₁₄ A₁₃ A₁₂ A₁₁
1 0 1 0 0 1 = all the 1's bits

∴ Address range is 11001000 0000 0000

(when address 1 is min) (when address 0 is max)
i.e., C800 to CFFF

Q.21
6-09

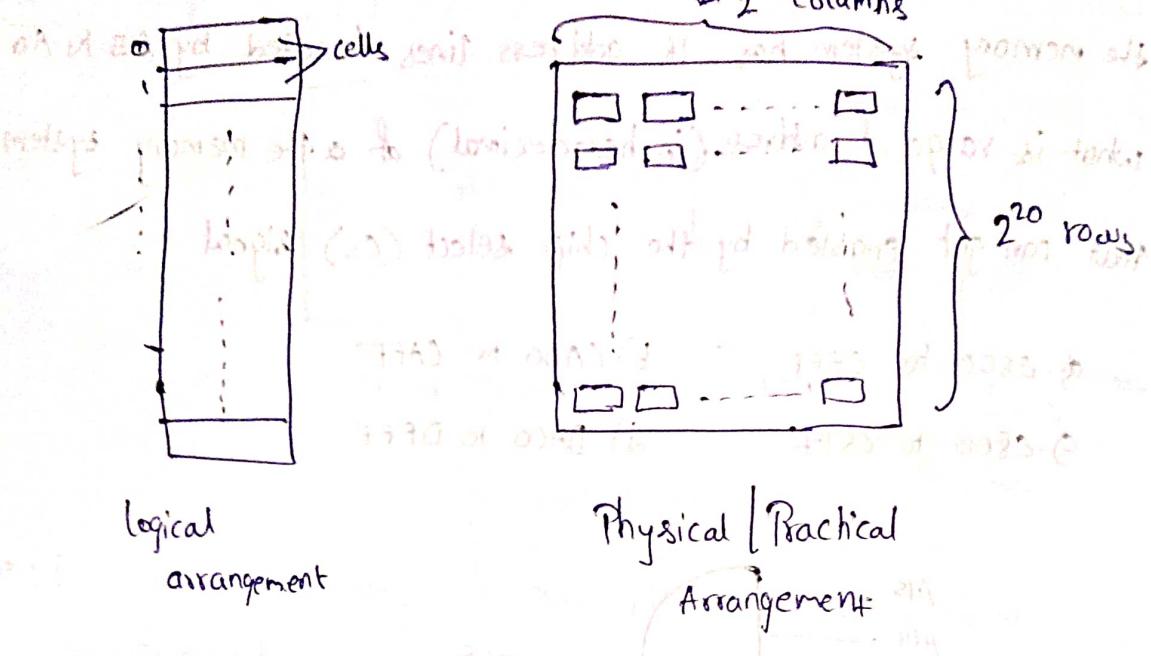
How many 82k x 1 RAM chips are needed to provide memory capacity of 256 kBytes.

So, given off one smit

$$\frac{256 \text{ k} \times 8}{32 \text{ k} \times 1} = \frac{2^8 \times 8}{2^{15} \times 1} = 2^3 \times 8 = 64$$

DRAM chip cell Arrangement

Consider $1G \times 1B$ memory $\Rightarrow 2^{30}$ cells



DRAM Chip Refresh

- * In one refresh operation, 1 row of cells can be refreshed.

\therefore The above chip needs 2^{20} refreshes.

$$\therefore \text{total chip refresh time} = \left(\frac{\text{no of rows}}{\text{in the chip}} \right) * \left(\frac{1 \text{ refresh operation}}{\text{time}} \right)$$

- \rightarrow If a system has n chips in memory

total refresh time of the memory system

$$= \frac{n \text{ chip refresh time}}{(n - \text{chip refresh time})}$$

* \because Each chip can be refreshed parallelly

(Q22) A DRAM chip of $512k \times 8$ bits has 256 rows of cells with 2k cells in each row. If DRAM takes 20ns hr 1 refresh then total refresh time of DRAM is _____

Sol:

we need 256 refreshes

$$\therefore 256 \times 20 \text{ ns}$$

$$= 5120 \text{ ns}$$

(Q23) G-10 A main memory unit with a capacity of 4 MB is built using $1M \times 1$ -bit DRAM chips. Each DRAM chip has 1k rows of cells with 1k cells in each row. The time taken for a single refresh operation is 100ns. The time required to perform one refresh operation on all the cells in the memory unit is

- a) 100 ns
- b) $100 \times 2^{10} \text{ ns}$
- c) $100 \times 2^{20} \text{ ns}$
- d) $3200 \times 2^{20} \text{ ns}$

Sol:

time req. of entire mem refresh = refresh time of single chip

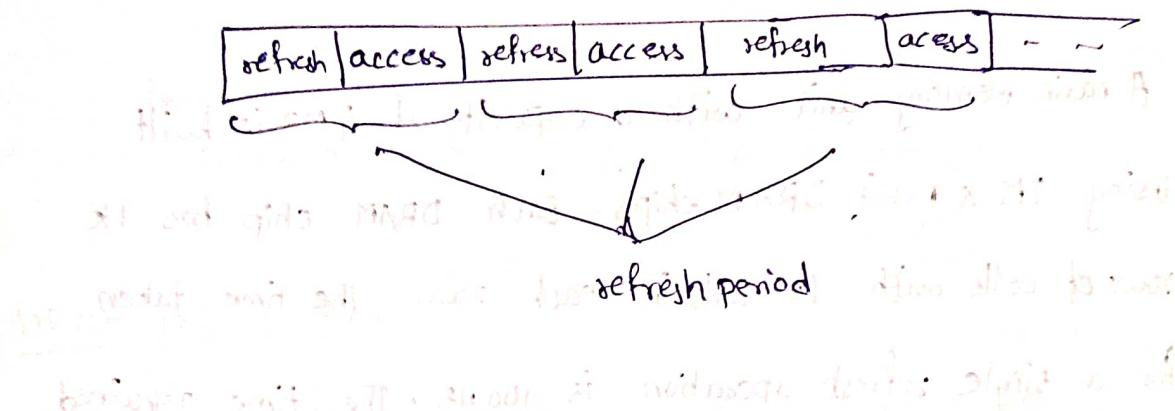
i.e., no of rows \times 1 refresh time

$$2^{10} \times 100 \text{ ns}$$

\therefore opt (b)

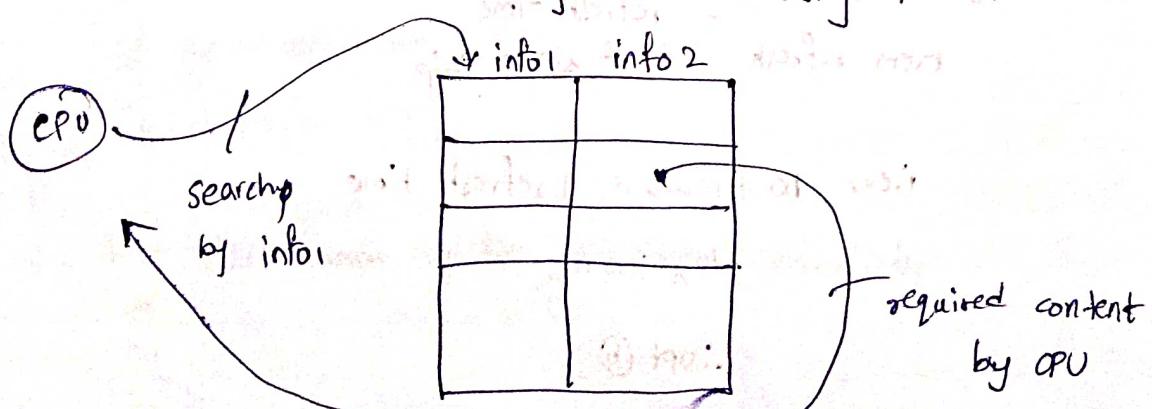
Refresh period

- Refresh period is amount of time after which refresh is performed once.
- Refresh is performed regularly.
- After performing the refresh some part time is used for memory access.



Associative Memory:

- This is also known as content addressable memory.
- The memory cells of associative memory don't have addresses.
- Every cell contains 2 information - out of which one is already known by CPU and other is the one required by CPU.
- So here searching is done using 1st information.



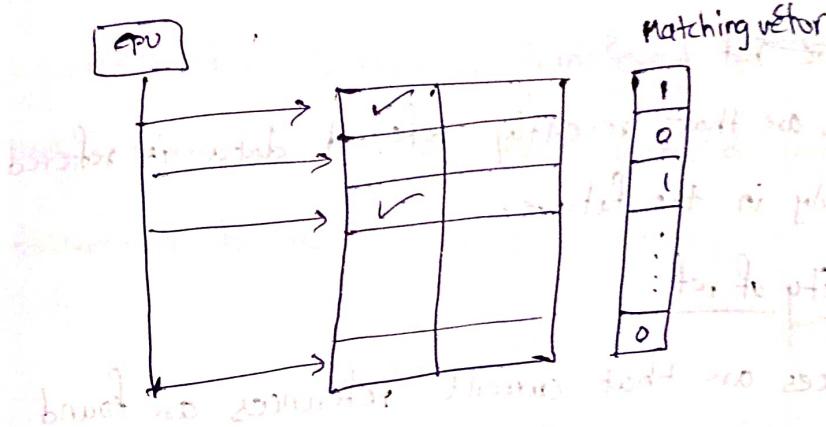
for which
is used for

- Comparison of each cell is done parallelly using separate matching logic for each cell.
- ∴ searching is very fast (faster than SRAM also).

Also cost is very expensive

- In associative memory, we maintain a matching vector which contains one bit corresponding to each cell.

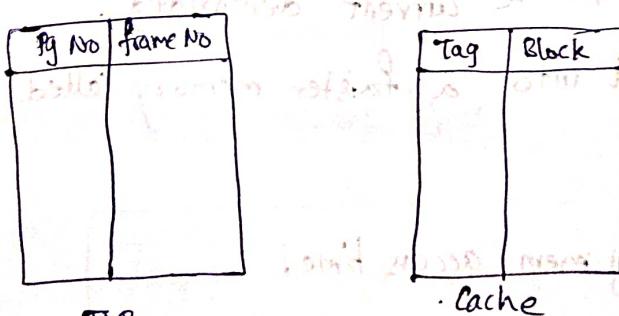
If match occurs the bit is set.



Applications:

- i) for cache & TLB implementation

↳ we use when fully associative mapping is required.



(Search by Pg.No)
(Search by tag information)

content

Locality of Reference:

If CPU has requested one address for memory access, then that particular address or nearby address will be accessed soon.

types :

i) Temporal locality of ref.

ii) Spatial locality of ref.

(i) Temporal locality of ref

The future ref

Chances are that recently referred data is referred most probably in the future.

(ii) Spatial locality of ref

Chances are that current references are found around the previous references.

Cache Memory:

→ Based on locality of ref, the current demanded localities are brought into a faster memory called cache.

→ Using cache reduces avg mem. access time.

Cache hit: If CPU's demanded content is present in cache then we say it is cache hit.

Cache miss: If CPU's demanded content is not present in cache, then we say it is cache miss.

→ Performance of cache is given by cache hit ratio

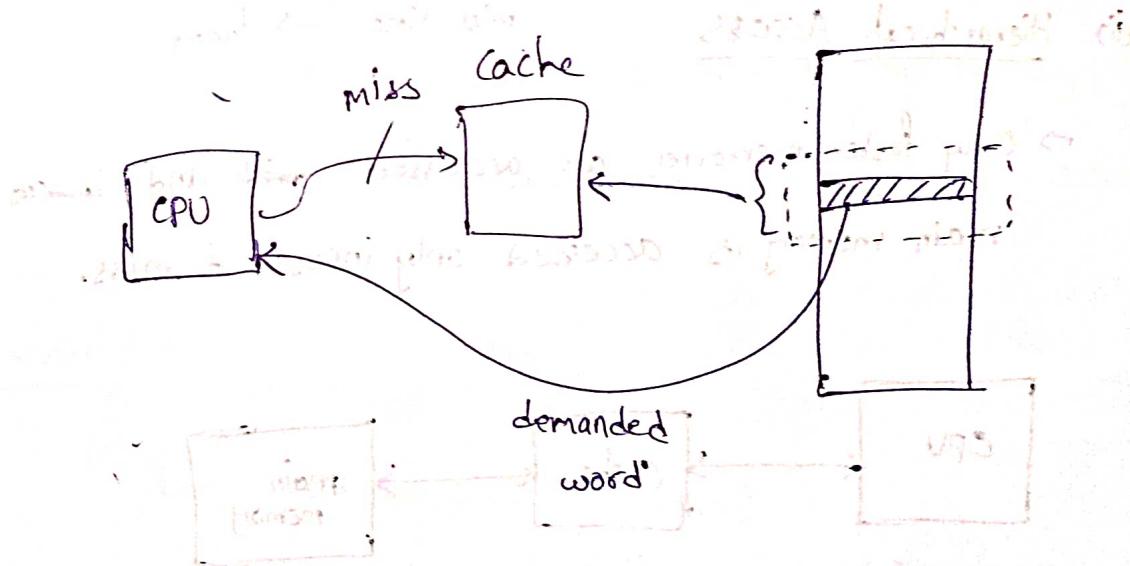
$$\text{hit ratio} = \frac{\text{no of hits}}{\text{total no of references}}$$

$$\Rightarrow \text{miss ratio} = 1 - \text{hit ratio}$$

~~hit ratio = ?~~

$$\text{hit rate} = (\text{hit ratio}) * 100$$

→ When a cache miss occurs the block in which required word is present is transferred to the cache and at the same time the requested word is directly transferred to the CPU from memory



Average Memory access time:

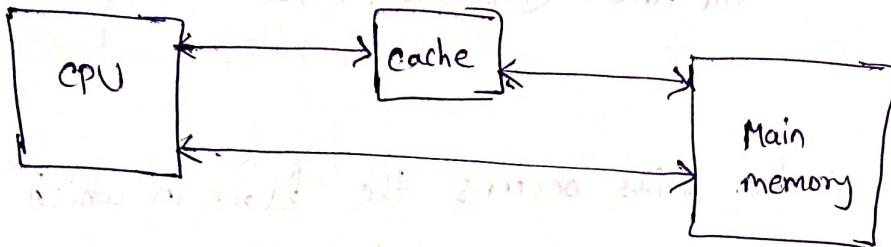
$$T_{avg} = h * (\text{hit time}) + (1-h) * \text{miss time}$$

where $h \rightarrow \text{hit ratio}$

Types of cache accesses:

i. Simultaneous Access:

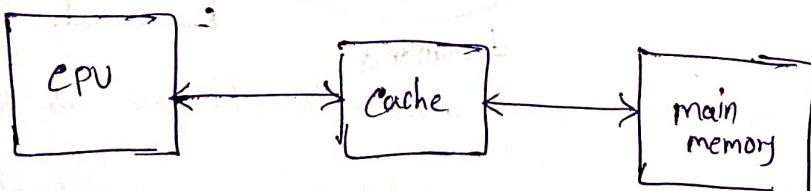
Here request for cache & main memory are generated simultaneously.



$$T_{avg} = (h * t_{cache}) + (1-h) t_{mm}$$

ii. Hierarchical Access

Only faster memories are accessed and in case of miss, main memory is accessed only.



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

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

hit time $\rightarrow t_{cm}$

miss time $\rightarrow t_{cm} + t_{mm}$

Note:

→ If words like level or hierarchy are used then we go for hierarchical access.
otherwise simultaneous access.

→ Also use simultaneous access if stmt like "ignore initial time req; to check if miss or hit occurred"

19/11/20

Tavg when locality of reference is added, included:

Simultaneous:

$$T_{avg} = h * t_{cm} + (1-h) t_{block}$$

t_{block} → block access time for main memory

*t_{block} = block size * t_{mm} [memory is fast]*

Hierarchical:

$$T_{avg} = h * t_{cm} + (1-h) (t_{block})$$

Tavg when block transfer included:

Simultaneous:

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

Note: In one cache access 1 cache block can be accessed

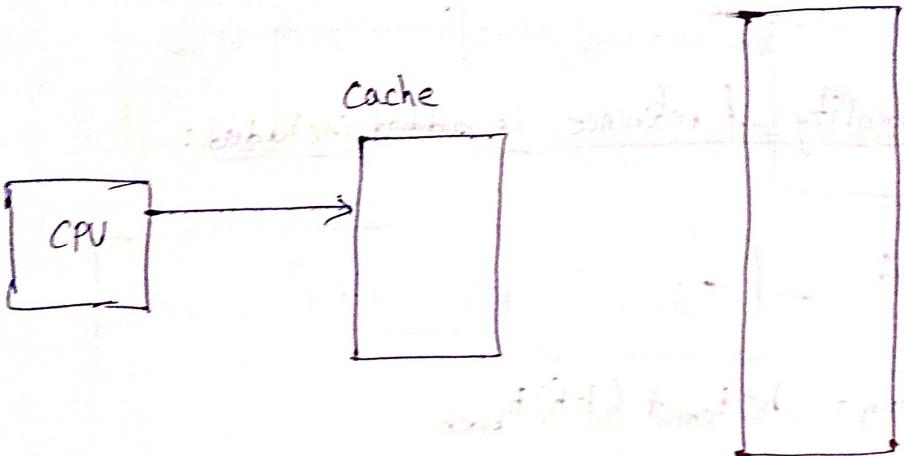
Hierarchical:

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

irresp
for
→ T_{avg}
→ $(T_{avg})_{avg}$

Cache Write/Write Propagation:

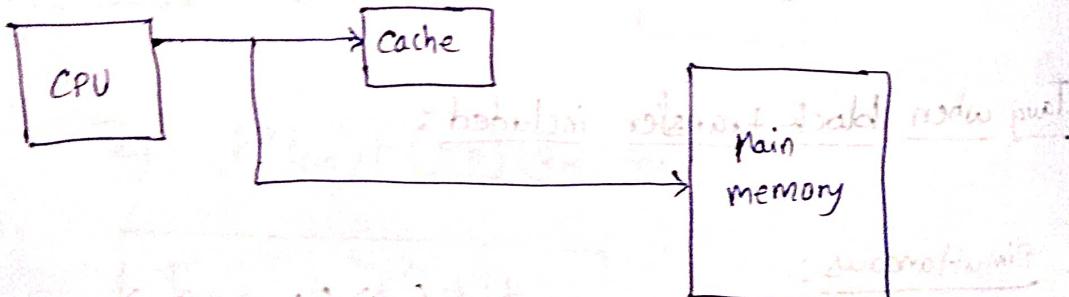
number of memory blocks in cache



→ If CPU update cache memory content, then the original content in memory should also be updated.

This is done in two ways:

i) Write through:



As soon as CPU updates cache memory contents, simultaneously main memory content is also updated.

Effective hit

$$h_{eff} =$$

Q24 Consider

If,

are for w

So,

90%

(only CM access)

i.e., irrespective of hit or miss, CPU will access MM.

→ Tang for read remains same.

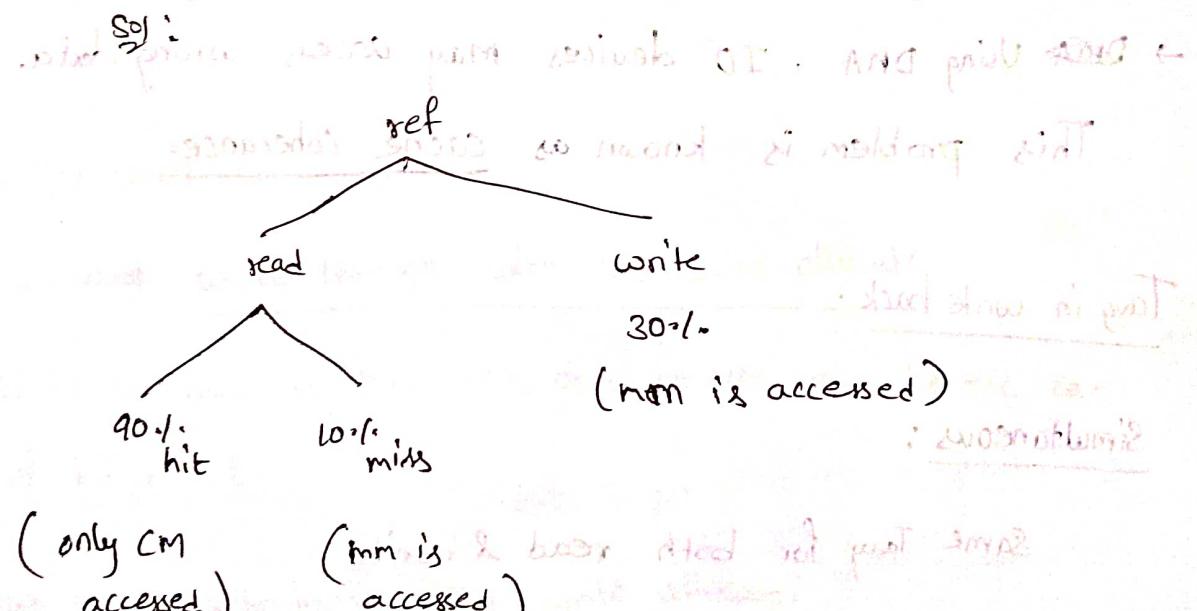
$$\rightarrow (T_{avg})_{write} = \max \{t_{cm}, t_{mm}\} = t_{mm}$$

Effective hit ratio: → Effective hit ratio = $\frac{\text{no of ref in which only Cache is accessed}}{\text{total no of ref (both read & write)}}$

$$h_{eff} = \frac{\text{no of ref in which only Cache is accessed}}{\text{total no of ref (both read & write)}}$$

(Q24) Consider a write through cache with 90% hit rate.

If, 70% of references are for read and 30% are for write, then find effective hit ratio.



$$\therefore \text{effective hit ratio} = 0.7 * 0.9 = 0.63$$

i) Write Back:

→ CPU updates only in cache & when a block is replaced it is copied into MM.

Adv of write through:

→ MM always contains correct information (Consistency)

Disadv of write through:

→ Takes more time.

Adv of write back:

→ Takes less time

Disadv of write back:

→ Inconsistency b/w values in cache & main memory

→ Using DMA, IO devices may access wrong data.

This problem is known as cache coherence.

Tang in write back:

Simultaneous:

Same Tang for both read & write

$$(Tang)_{read} = (Tang)_{write} = h \cdot t_{cm} + (1-h)(t_{block} + t_{write_back})$$

Hierarchical:

$$(Tang)_{read} = (Tang)_{write} = h \cdot t_{cm} + (1-h)(t_{cm} + t_{block} + t_{write_back})$$

$$twink_block = x * t_{block}$$

$x \rightarrow$ fraction of dirty block

20/11/20

Block size of addressable part

Block size of addressable part

Block size of addressable part

- * Based on the action we perform for write miss we have 2 policies

(i) Write Allocate:

The block in MM is loaded into cache on write miss,

followed by write-hit action

(ii) NO write allocate:

The block is modified in the MM itself and is not loaded into the cache.

Note:

* For better performance, we use Write Through with no write allocate.

* We used write back with write allocate.

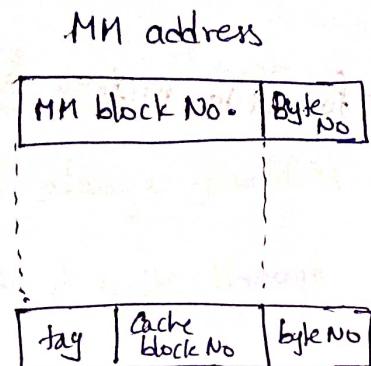
It is because any how CPU goes for MM even in the case

of hit or not.

* We used write back with write allocate.

Cache Mapping

Direct Mapping:



(tag information is also called metadata)

$$\text{Cache block num} = \frac{\text{MM block no.}}{\text{no of cache blocks}}$$

$$\text{no of bits in byte offset} = \log_2(\text{block size})$$

$$\text{no of bits in cache block} = \log_2(\text{no of cache blocks})$$

$$\text{no of tag bits} = \log_2\left(\frac{\text{size of MM}}{\text{size of cache}}\right) = \log_2\left(\frac{\text{no of MM blocks}}{\text{no of cache blocks}}\right)$$

$$\text{Tag directory size} = (\text{no of blocks}) * (\text{no of tag bits})$$

(when no info abt v/I bit and other is given)

Cache Initialization

→ when system is turned on, the content of tag is invalid.

To maintain information about validity of content of a cache block we maintain a valid bit per each block.

→ when system is turned on, all valid bits are set to 0.

→ so whenever we index into a cache block, we also check if the valid bit is set or not.

Modified bit / Dirty bit

- In write back cache we don't need to perform write back if the cache is ~~not~~ block content is not modified.
- So we maintain this information using dirty bit

dirty bit = 1 ⇒ dirty block (modified)

dirty bit = 0 ⇒ clean block (unmodified)

Note:

- Cache Controller is a hardware that ~~controls~~ controls cache accessing and has memory to store meta data (tags, valid bit, modified bit etc.)

~~22/11/20~~
22/11/20

~~page fault~~ ~~with invalid~~, ~~l1/l2~~

Set Associative Mapping

$$\boxed{\text{set No.} = (\text{MM block No.}) \div (\text{no of sets})}$$

- no of sets in cache = $\frac{\text{no of blocks in cache}}{\text{associativity}}$

→ used for cache indexing.

| | | |
|----------|---------|-------------|
| tag bits | set No. | byte offset |
|----------|---------|-------------|

$$\rightarrow \text{no of tag bits} = \log_2 \left(\frac{\text{no of MM blocks}}{\text{no of sets}} \right) = \log_2 \left(\frac{\text{MM size}}{\text{cache size}} \right) + \log_2 (\text{associativity})$$

Ex: MM add: 32 bits

Cache size: 64 kB

block size: 32 B

4-way set associative

$$\rightarrow \text{no of blocks} = \frac{2^{16}}{2^5} = 2^{11}$$

$$\text{no of sets} = \frac{2^{11}}{4} = 2^9$$

$$\text{no of tag bits} = \log_2 \left(\frac{\text{MM size}}{\text{cache size}} \right) + \log_2 (\text{associativity})$$

$$= \log_2 (2^{32-16}) + \log_2 (4)$$

$$= 16 + 2 = 18$$

address consisting full combination of (tag bits) + (set no) + (byte offset)

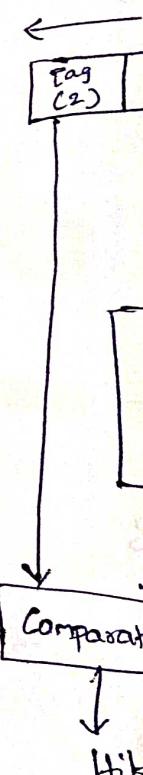
| tag bits | set no | byte off |
|----------|--------|----------|
| 18 | 9 | 5 |

→ 32 → 18 + 9 + 5

Fully Associative Mapping

| Tag bits | byte offset |
|----------|-------------|
|----------|-------------|

$$\rightarrow \text{no of tag bits} = \log_2 (\text{no. of MM blocks}) = \log_2 (2^{11})$$



→ n-bits

| Set no | Block | Address |
|--------|-------|--------------------------|
| 0 | 0 | 000000000000000000000000 |

→ Hit latency

Hardware implementation of Mapping

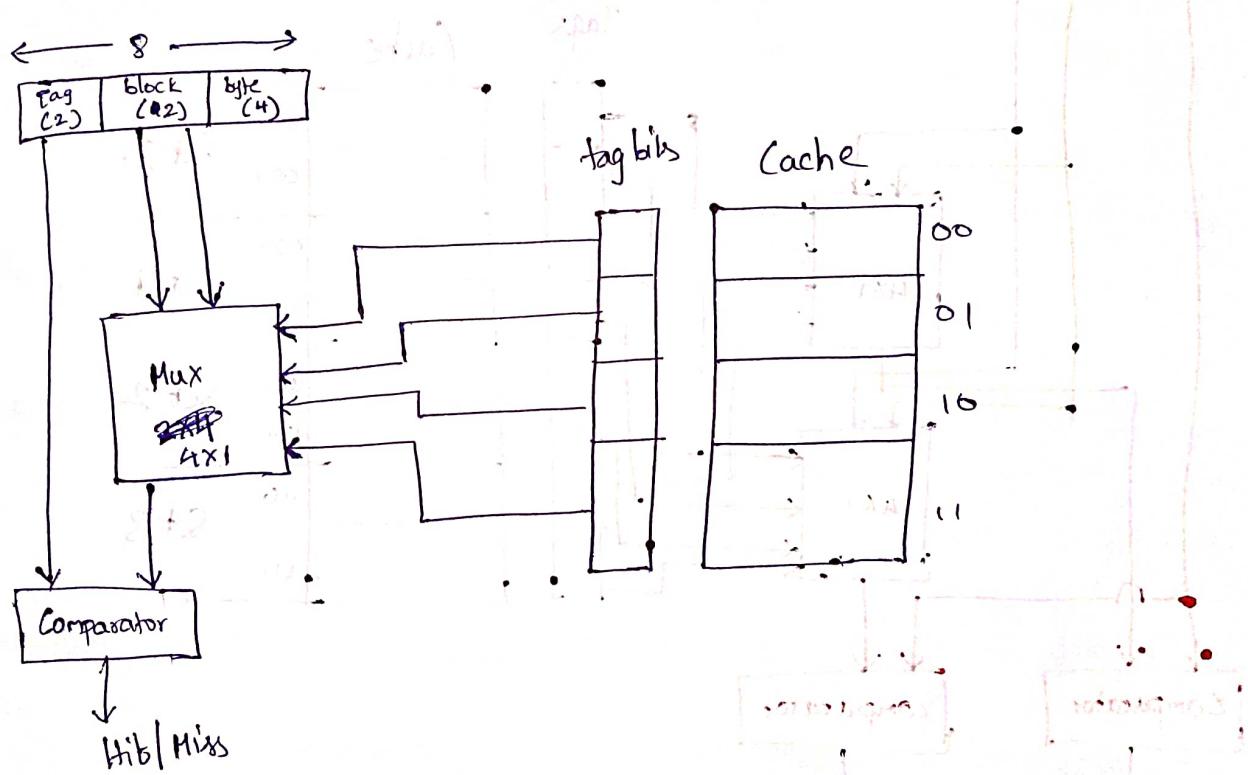
Direct Mapping:

Consider MM size: 2^8 B

Cache size: 2^6 B

Block size: 2^4 B

$$\Rightarrow \text{no of cache blocks} = 2^6/2^4 = 2^2 = 4 \text{ blocks}$$



\rightarrow no of Muxes need = no of tag bits

\rightarrow size of Mux = $2^x \times 1$

~~no of tag bits~~ $2^x \rightarrow$ no of blocks

\rightarrow n-bit comparator

$n \rightarrow$ no of tag bits.

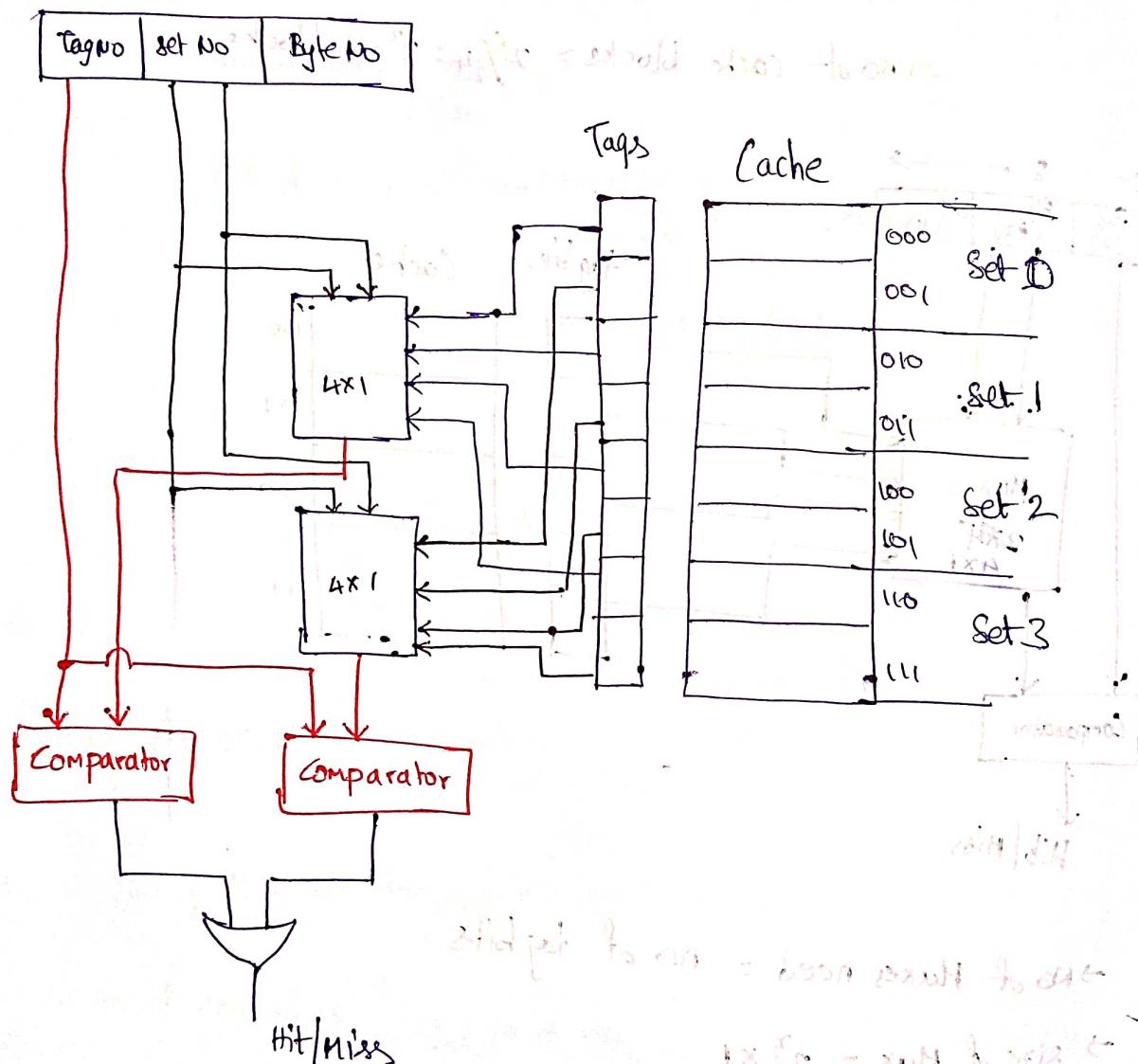
Hit latency: (Time taken to get data) = $100\text{ns} \times 10 = 1\mu\text{s}$

Amount of time required for cache controller to determine whether hit or miss occurred.

→ In direct mapping all the Muxes operate in parallel.

$$\therefore \text{Hit latency in direct mapping} = \text{Mux delay} + \text{Comparator}$$

Set Associative Mapping:



$$\rightarrow \text{Size of Mux} = 2^x \times 1$$

↳ no of sets

$$\rightarrow \text{no of Mux} = (\text{no of tag bits}) * (\text{associativity})$$

→ n-bit comparator ; $n \rightarrow \text{no of tag bits}$

→ no of comparators = Associativity

→ k-input OR gate ; $k \rightarrow \text{Associativity}$

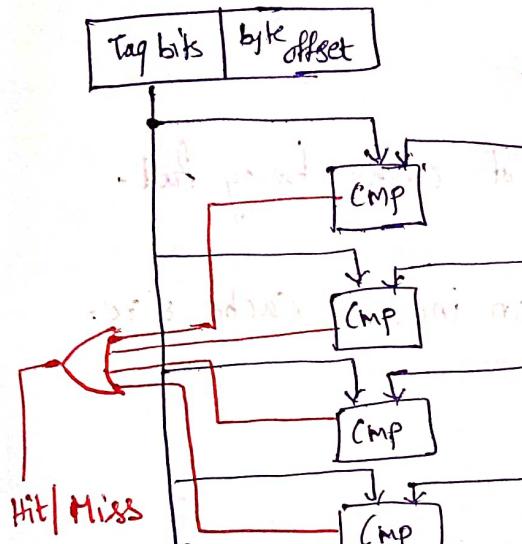
hit latency

$$\text{in-set associative} = \text{Mux delay} + \text{Comparator delay} + \text{OR gate delay}$$

Mapping delay \leq set associativity delay \leq hit latency \leq block size \times block delay

Set associativity \rightarrow number of bits in tag \leq number of sets in cache

Fully Associative Mapping:



Tag Cache

Cache [lineup] (b)

\rightarrow size of comparator = no of tag bits = $\log_2(\text{no of MM blocks})$

\rightarrow no of comparators req = no of cache blocks

\rightarrow size of OR gate = no of comparators = no of cache blocks

Block Replacement:

\rightarrow Direct Mapping doesn't req any replacement policy.

\rightarrow For set associative & fully associative we use below replacement policies

(i) FIFO

(ii) Optimal

(iii) LRU

Types of Cache Misses:

(i) Cold / Compulsory miss:

This is miss due to first time reference of a cache block.

→ Compulsory miss can be reduced by increasing block size.

(ii) Capacity Miss:

→ It is cache miss occurring because of cache being full.

→ To reduce capacity miss, we can increase Cache size.

→ It is not cold miss.

(iii) Conflict Miss:

If cache set is full (to which MM block maps) and hence miss occurs due to tag mismatch.

→ Conflict miss is not cold miss & not capacity miss.

→ Conflict miss can be reduced by increasing associativity.

Note:

→ When asked to check what kind of miss a given miss is, check only in the below order.

i) Cold miss

ii) Capacity miss

iii) Conflict miss

Note: Conflict miss is never possible in fully associative mapping.

Eg: Consider

No of blocks in cache = 4

Let cache be 2-way set associative
and let replacement policy be LRU.

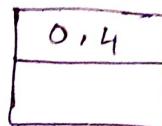
Consider ref string: 0, 4, 0, 8, 0, 4, 1, 3, 1, 5, 1, 3

0: ~~Conflict~~

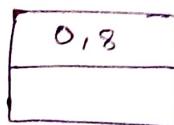
Compulsory miss (bit of address)

1: Compulsory miss

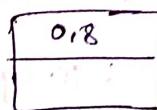
2: Hit



3: Compulsory Miss

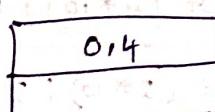


4: Hit



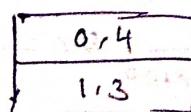
* 5: Here cache is not full.
but set to which map is full.

∴ Conflict miss

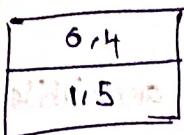


1: Cold miss

2: hit



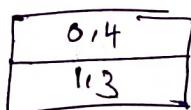
5: Compulsory Miss



! : hit

* 3: Here cache is full & set is also full
still we say it is capacity miss

∴ Capacity miss



Cache miss penalty

It is time required to bring a missed block from main memory to cache.

Assume

→ Cycles required to send address to MM : 1 cycle

Cycles req to access 1 MM cell (word) : 10 cycles

Cycles req to transfer 1 cell data to cache : 1 cycle

| Cache block size | MM cell size | Miss penalty |
|------------------|--------------|---|
| 4 bytes | 1 byte | 1 + (4 * 10) + 4 = 45 cycles → Address is sent only once. |
| 4 bytes | 2 bytes | 1 + (2 * 10) + 2 = 23 cycles from then accessing of entire block is done. |
| 4 bytes | 4 bytes | 1 + (1 * 10) + 1 = 12 cycles |

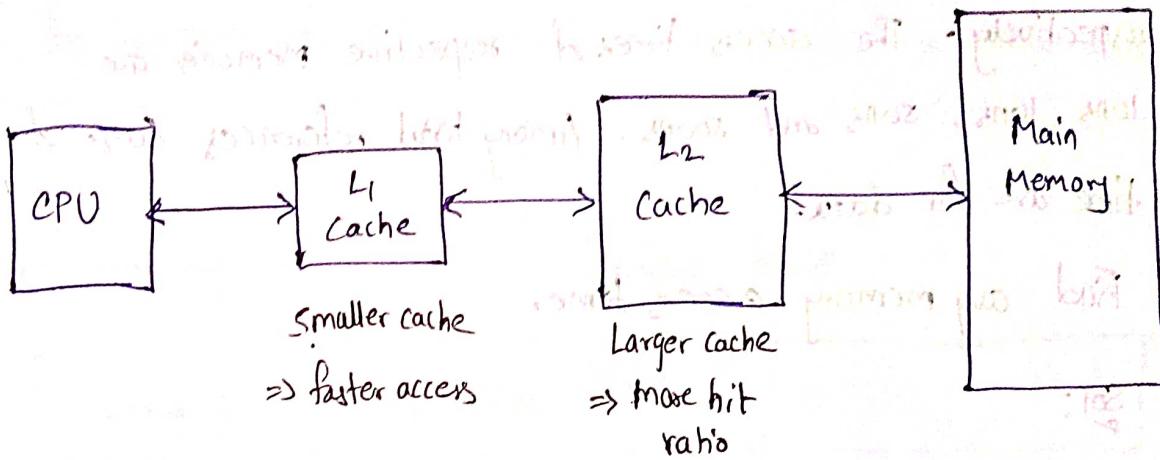
→ If CPU runs on 500 MHz clock then

$$45 \text{ cycles} = 45 \times \frac{1}{500 \times 10^6 \text{ sec}} = 9 \times 10^{-8} = 90 \text{ ns}$$

Goals of using cache memory:

- Minimize access time. → (smaller cache) } Multilevel
- Maximize hit rate. → (larger cache) } cache
- Minimize miss penalty.

Multilevel Cache:



Tang in Multilevel cache:

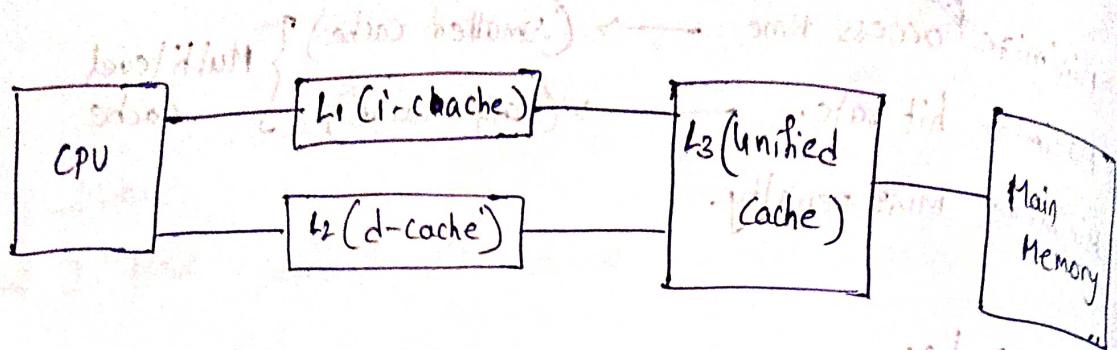
i) Simultaneous access:

$$T_{avg} = h_1 * t_1 + (1-h_1) [h_2 * t_2 + (1-h_2) t_{mm}]$$

ii) Hierarchical access:

$$\begin{aligned} T_{avg} &= h_1 * t_1 + (1-h_1) [h_2 * (t_1 + t_2) + (1-h_2) (t_1 + t_2 + t_3)] \\ &\approx t_1 + (1-h_1) t_2 + (1-h_2) t_3 \end{aligned}$$

Eg: The multilevel memory hierarchy is given:



The hit ratios of L₁, L₂, L₃ and MM are 0.8, 0.9, 0.95 and 1.0 respectively. The access times of respective memories are 10ns, 10ns, 50ns and 500ns. Among total references 60% of time are for data.

Find avg memory access time.

Sol:

Avg access time for instruction

$$\begin{aligned}
 &= t_1 + (1-h_1) t_2 + (1-h_2) t_3 \\
 &= 10 + (1-0.8)(50) + (1-0.8)(1-0.95)(500) \\
 &= 10 + 10 + (0.2)(0.05)(500) \\
 &= 10 + 10 + 5 \\
 &= 25
 \end{aligned}$$

Avg access time for data

$$\begin{aligned}
 &= 10 + (1-0.9)(50) + (1-0.9)(1-0.95)(500) \\
 &= 10 + 5 + 2.5 \\
 &= 17.5
 \end{aligned}$$

Avg memory access time = $0.6(17.5) + 0.4(25)$

Cache Inclusion Policies:

There are two types of policy

(i) inclusion policy

(ii) Exclusion policy

(i) Inclusion policy

The content of L₁ cache is present in L₂ cache also.

i.e., L₂ is inclusive of L₁.

For read operations, consider below table

| Scenario | Operation |
|---|---|
| Hit in L ₁ | CPU gets content from L ₁ |
| Miss in L ₁ Hit in L ₂ | <p>Copy a block from L₂ to L₁</p> <p>A replaced block is copied only to memory and no participation of L₂</p> |
| Miss in L ₁ & Miss in L ₂ & Hit in MM | <p>Copy the block from MM to L₁ & L₂</p> <p>both</p> <p>→ If a block is replaced from L₂ then send back invalidation to L₁ (i.e., set valid bit in L₁)</p> |

(ii) Exclusion Policy:

Content of L₁ cache need not to be present in L₂ cache.

→ For read operations, consider below table:

| Scenario | Operation |
|--|--|
| Hit in L ₁ | CPU gets contents from L ₁ |
| Miss Hit in L ₁ & Hit in L ₂ | Move a block from L ₂ to L ₁ , the replaced block of L ₁ is moved to L ₂ . |
| Miss in L ₁ & Miss in L ₂ | Copy the block from MM to L ₁ only. The replaced block (if any) is moved to L ₂ cache. |

→ In exclusion policy, L₂ is called victim cache, since it is being filled by victim blocks (replaced blocks) of L₁.

(iii) Value inclusion policy:

→ In inclusion policy, ~~all~~ all the blocks in L₁ are present in L₂ also, but it is not needed that values are same in L₁ & L₂.

→ Value inclusion policy ensures that value of blocks in L₁ is same as value of blocks in L₂.

24/Intro

- (Q25) what hit rate is required to reduce the effective memory access time from 150ns to 42ns, if cache access time is 30ns.

Sol: ~~if there is no level or hierarchy~~ \Rightarrow hit rate = 0.9

These is no word like level or hierarchy.

so go for simultaneous access

$$42 = h(30) + (1-h)(150)$$

$$\Rightarrow h = 0.9 \text{ id. of been hit } 42$$

$$\therefore \text{hit rate} = 90\%$$

$$(access time)_{cache} + (0.1)_{miss} = 42$$

- (Q26) A computer system contains a cache. Uncached memory access takes a times longer than access to cache.

If cache has hit ratio 0.8, The ratio cache memory access time to uncached memory access time is _____

Sol: ~~if there is no level or hierarchy~~ \Rightarrow hit rate = 0.9

let uncached access time = m ~~in case of mem per bit~~

\therefore mem access time

~~any access time~~ \Rightarrow cache access time = $\frac{m}{9}$

$$\therefore \frac{0.8(m/a) + 0.2(m)}{m} = \frac{\frac{8}{9} + 2}{10} = \frac{26}{90} = \frac{13}{45}$$

(Q27) In a two-level hierarchy, the top level has an access time of 20 ns and the bottom level has an access time of 80 ns. The hit rate on the top level is 90%. If the block size of cache is 8 bytes then avg mem access time is _____.

(Note: Consider the system uses locality of ref)

Sol:

$$(Cost)(H \cdot 1) + (Cost)d = ?$$

Here we need to use hierarchical access

$$T_{avg} = 0.9(20) + 0.1(20 + 8 * 80)$$

block access time

$$= 18 + 0.1(660)$$

$$= 18 + 66 = 84 \text{ ns}$$

(Q28) 2-level hierarchy; Cache access time = 15 ns; MM access = 110 ns
hit rate is 90%. Block size = 16 bytes;

Find avg mem access time including miss penalty

Sol:

Miss penalty is time to bring block from MM to Cache

$$T_{avg} = 0.9(15) + 0.1(15 + 110 * 16 + 15)$$

$\underbrace{\hspace{1cm}}$ block access \downarrow transferring block to Cache

$$= 13.5 + 179 = 192.5 \text{ ns}$$

(Q29) Consider a memory hierarchy with a write back cache. The cache has an access time of 25 ns - 90% of all the memory accesses are found in the cache itself. The main memory access time is 300 ns. The 10% of the cache blocks are dirty. The cache block size is 4 bytes. Find avg mem access time.

Sol:

$$t_{avg} = h * t_{cm} + (1-h) (t_{cm} + t_{block} + t_{write-back})$$

$$= 0.9(25) + 0.1(25 + 4*300 + 0.2*4*300)$$

$$= 22.5 + 0.1(25 + 1200 + 240)$$

$$= 169 \text{ ns}$$

(Q30) Consider direct mapped cache with

32-bit architecture

MM size = 4 GB

Cache size = 512 kB

block size = 16 words

Find

address division -

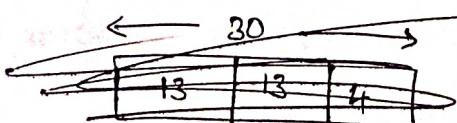
Sol:

~~word size = 4 words~~

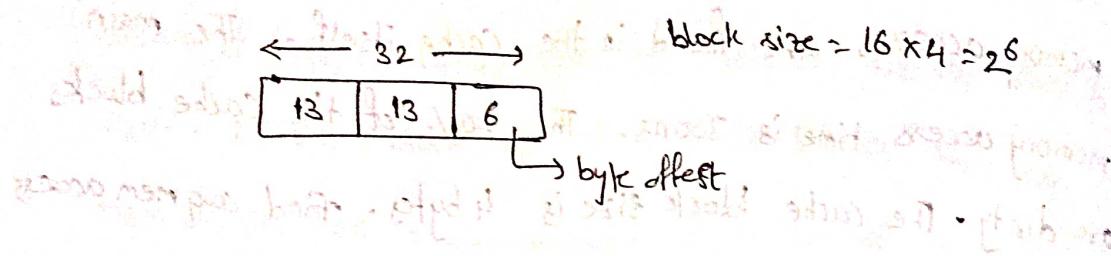
~~MM size = 2^{32} bytes = 2^{30} words~~

~~Cache size = 2^{19} bytes = 2^{17} words~~

~~no. of MM blocks = 2^{26} ; no. of cache blocks = 2^{13}~~



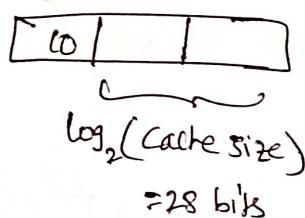
Since no information is given we consider memory is by te addressable.



Q31 Consider a direct mapped cache of size 256 MB.

No of tag bits = 10. Find max size of main memory supported in the system is _____

Sol:



∴ size of PA = 38 bits

∴ size of MM = 256 GB

Q32 MM size: 2^{36} bytes

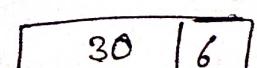
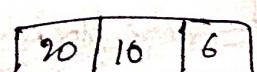
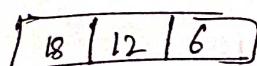
Cache size: 2^{18} bytes

block size: 2^6 bytes

Find no of tag bits

Find add. division in direct, 4-way set ass., full associative

Sol:



27/11/20

Q33 Consider a direct mapped cache.

① 32 bit address of main memory. How many bits?

Block size: 32 bytes. No. of blocks?

No. of cache blocks: 256. Why?

In which block of cache should we look for each of the following addresses?

1. 1A2BC012

2. FFFF00FF

3. 12345678

4. C109 D532

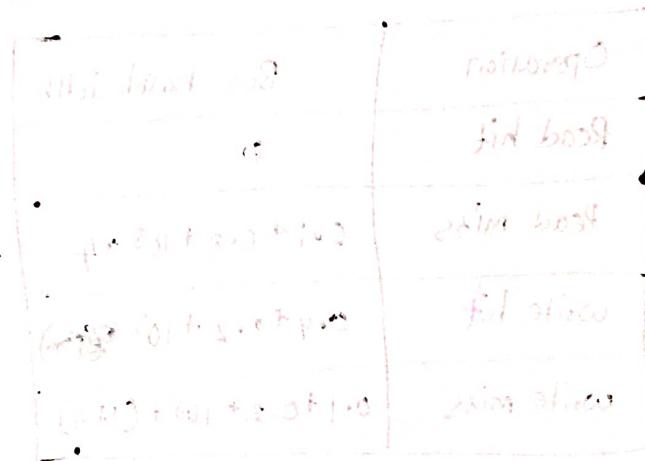
Sol:

Ans: 1 - 0

2 - 7

3 - 179

4 - 169



Note:

CPU requires data bus only when it needs to access main memory. To access cache memory CPU uses internal bus.

**

Q34

Cache hit ratio = 0.9

Block size: 4 words & use locality of reference.

CPU sends reference to the cache at a rate of 10⁸ words per second.

20% of the above ref are write.

The bus reads or writes a single word at a time.

Assume at any one time, 20% of the block frames in the cache have been modified.

The cache uses write allocate on a write miss.

The percentage of the bus bandwidth used on average in

1. write through Cache: _____

2. Write back cache: _____

Sol:

1. write through Cache:

| Operation | Bus bandwidth |
|------------|----------------------------|
| Read hit | 0 |
| Read miss | $0.1 * 0.8 * 10^8 * 4$ |
| write hit | $0.9 * 0.2 * 10^8 * 4$ |
| write miss | $0.1 * 0.2 * 10^8 * (1+4)$ |

$$\text{total} = (0.32 + 0.18 + 0.1) * 10^8$$

$$= 0.51 * 10^8$$

$$\therefore \text{bandwidth used} = 51\%$$

2. Write back Cache:

| | |
|------------|------------------------------------|
| Read hit | 0 |
| Read miss | $0.1 * 0.8 * 10^8 * (4 + 0.2 * 4)$ |
| write hit | 0 |
| write miss | $0.1 * 0.2 * 10^8 * (4 + 0.2 * 4)$ |

$$\text{total} = (0.384 + 0.096) * 10^8$$

$$= 0.48 \times 10^8$$

∴ bandwidth used = 48%.

Q35

Consider a cache with 2^{13} blocks of size 32 bytes each.

The CPU generates 32-bit address.

Cache controller has 1 valid bit, 1 modified bit and tag bits.

size of cache controller memory is 18kbytes.

Cache is organized as k-way set associative.

Find max value of k to utilize the controller memory maximum is _____?

Ans: 8

Q36

A computer system uses 16 bit memory addresses. It has 2k bytes cache. organized in direct mapped manner.

Block size: 64 bytes.

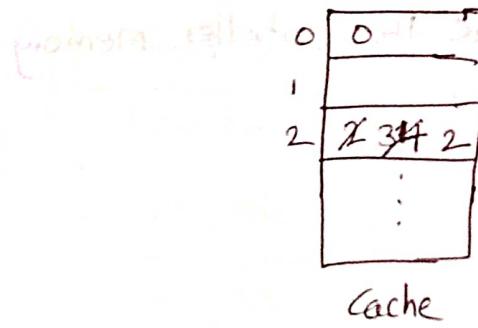
CPU generates below references (decimal values)

28, 144, 2176, 2180, 128, 2176.

Assuming cache is initially empty. Find no of misses.

- a) 3 b) 4 c) 5 d) 9

| MM address | MM block (-16k) | Cache block (-1-32) | Hit / Miss |
|------------|--------------------|------------------------|------------|
| 28 | 0 | 0 | Miss |
| 144 | 2 | 2 | Miss |
| 2176 | 34 | 2 | Miss |
| 2180 | 34 | 2 | Hit |
| 128 | 2 | 2 | Miss |
| 2176 | 34 | 2 | Miss |



D
D

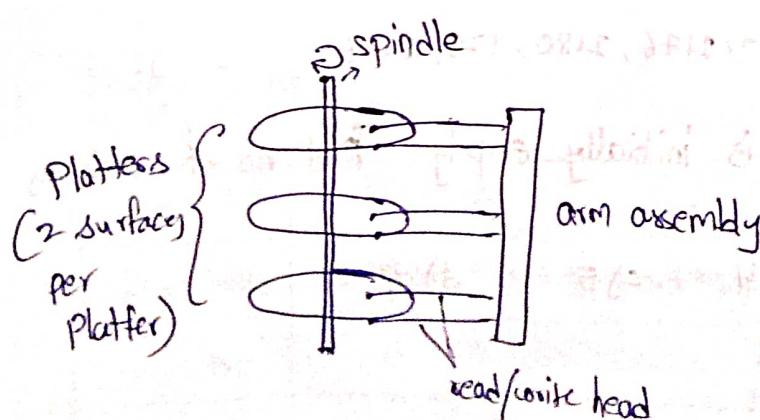
Conc

→ Al

→ Thu

28/11/20 more important topics in beginning, then end of

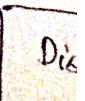
Magnetic Disk

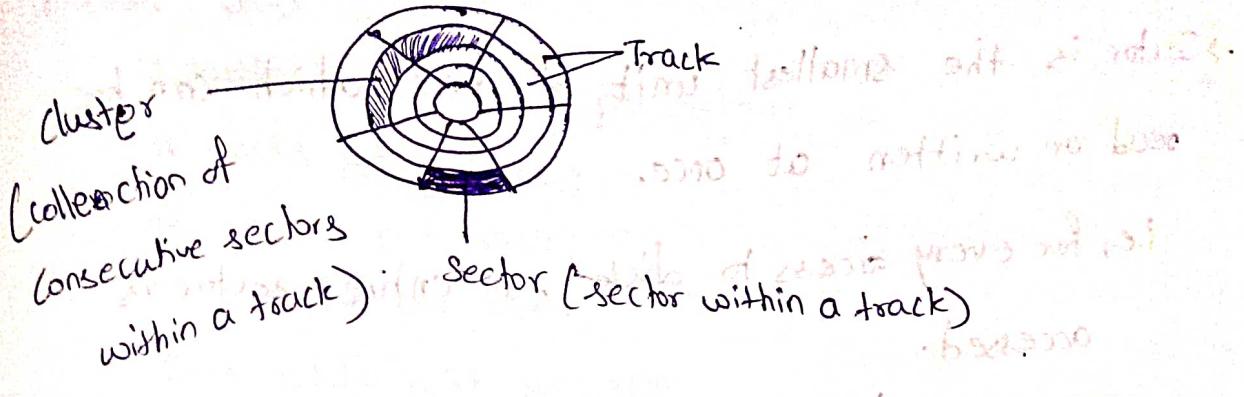


→ Here
Vaniak

→ Sto

→ Thu





Disk Capacity

Data is stored on disk in 2 ways:

- i) Constant sector capacity
- ii) Variable sector capacity

Constant Sector Capacity: (Default consideration)

- All sectors store same amount of data.
- Thus density of storage is lesser as we move to the outer tracks.

$$\text{Disk capacity} = \text{no of platters} * \text{no of tracks}$$

$$* \text{no of sectors per track} * \text{sector capacity}$$

- Here angular velocity is constant.

Variable Sector Capacity:

- Storage density is constant in all the tracks.
- Thus sector capacity varies from track to track.

$$\text{Disk capacity} = \text{no of platters} * \text{surface capacity.}$$

(These linear velocity is constant)

Note:

→ Sector is the smallest unit of disk which can be read or written at once.

i.e., for every access to disk, an entire sector is accessed.

→ Sector is the ^{least} addressable unit of disk.

Disk Access Time:

→ Disk access time is time required for accessing 1 sector.

$$\text{Disk access time} = \text{seek time} + \text{Rotational latency} + \text{transfer time}$$

(for 1 sector)

where transfer time is additional delay given in question.

Seek time: time req. to position the arm over the desired track

Rotational Latency: time req. to rotate desired sector under R/W head

Transfer time: time req. to read or write 1 sector.

Note:

→ In one rotation time, 1 track can be transferred.

$$\therefore \text{transfer time of } 1 \text{ sector} = \frac{1 \text{ rotation time}}{\text{no of sectors per track}}$$

Q: Consider disk with

16 platters; 2 surfaces per platter;

1k tracks per surface;

2k sectors per track;

2048 bytes per sector;

Disk rotates with 3000 rpm

seek time = 10 ms

a) Find capacity of disk

$$2 \times 16 \times 2^{10} \times 2^{11} \times 2048 \text{ B}$$

$$= 2^{37} \text{ Bytes}$$

$$= 128 \text{ GB}$$

b) no of bits required for addressing the disk?

$$= \log_2 (\text{no of sectors})$$

$$= \log_2 \left(\frac{2^{37}}{2^{11}} \right)$$

$$= 26 \text{ bits}$$

c) Find disk access time.

seek time + rot. + transfer time

latency

$$60 \text{ sec} \rightarrow 3000 \text{ rot}$$

$$1 \text{ rot} \rightarrow \frac{1}{3000} \text{ sec} = 20 \text{ ms}$$

transfer time of
1 sector

$$= \frac{1}{60 \text{ sec}} = \frac{20}{2^{11}} \text{ ms}$$

$$= 10 + \frac{1}{120}$$

$$= 10 + \frac{20}{2} + \frac{20}{2^{11}}$$

$$= 10 + 10 + 0.01$$

$$= 20.01 \text{ ms}$$

d) Find disk transfer rate.

in 1 rotation, 1 track is transferred

$$20 \text{ ms} \rightarrow 2^{11} * 2^{11} \text{ bytes}$$

$$20 \text{ ms} \rightarrow 2^{22} \text{ bytes}$$

$$1 \text{ sec} \rightarrow \frac{2^{22}}{20 \times 10^{-3}}$$

$$= \frac{2^{22} \times 10^2}{2}$$

$$= 2^{21} \times 10^2$$

$$= 200 \text{ MBPS}$$

= 200 Mega Bytes Per Second.

e) If the disk is used in cycle stealing mode of DMA such that whenever 64-bit word is available, it will be transferred in 16 ns. Find % of time CPU is blocked?

$$\% \text{ time CPU blocked} = \frac{\text{transfer time}}{\text{Preparation time}}$$

1 sector prep time = 0.01 ms

$$1 \text{ byte} = \frac{0.01}{2048} \text{ ms} = \frac{10^4}{2 \times 10^3} \text{ ns} = 5 \text{ ns}$$

$$8 \text{ bytes} = 40 \text{ ns}$$

$$\% \text{ time CPU blocked} = \frac{16}{40} \times 100 = 40\%$$

Multiple Sectors Access time:

(i) Sequential:

Here sectors are stored on the same track consecutively.

Transfer time of n sectors = seek time + Rotational + n * transfer latency (for 1 sector)

(ii) Random:

Here each sector is stored in a random location.

Transfer time of

$$n \text{ sectors} = n * (\text{seek time} + \text{rot. latency} + \text{transfer time for 1 sector})$$

Cylinder: Collection of tracks of same radius from all surfaces.

→ Data is stored cylinder wise.

→ Using the concept of cylinder we reduce seek time and thus improve speed.

→ No of cylinders = no of tracks per surface.

Note:

A sector is addressed as $\langle c, h, s \rangle$.

c → cylinder number

h → surface number

s → sector number in track.

⇒ Sector number = $(c * \text{no of surfaces} * \text{sectors/track}) + (h * \text{sectors/track}) + s$

(Assuming starting address is $\langle 0, 0, 0 \rangle$)

→ So for given sector number 'n' we can find corresponding $\langle c, h, s \rangle$

$$c = n // (\text{sectors/cylinder}) \rightarrow \text{integer division}$$

$$h = [n \% (\text{sectors/cylinder})] / (\text{sectors/track})$$

$$s = [n \% (\text{sectors/cylinder})] \% (\text{sectors/track})$$

Parallel Processing:

→ It is simultaneous processing of data.

8 types of parallel processing

i) Vector processing

ii) Array processing

iii) Pipelining

} not there in the syllabus.

Flynn's Classification of Computers:

i) SISD (Single Instruction stream, Single data stream):

At a time one instⁿ is fetched and executed.

↳ e.g.: Von neuman computers

ii) SIMD (Single Instⁿ stream, Multiple data stream):

One instⁿ is fetched at a time, but multiple instⁿs are executed at a time.

↳ e.g.: Pipelined processor.

iii) MISD (Multiple Instⁿ stream, Single Data stream)

More instⁿs are fetched at a time but only one instⁿ is executed at a time.

→ This is not practical as fetch multiple instⁿs but executing only one at a time is of no use.

(iv) MIMD (Multiple Instⁿ stream, Multiple data Stream)

- Here more than one instⁿ is fetched at a time & more than one instⁿ is executed at a time.
- This type systems have multiple pipelines
- Eg: Super scalar Computers.
- This systems have Instⁿ level parallelism (ILP)

Pipelining:

- Pipelining is useful when same processing is applied over multiple inputs.
- Pipelining decomposes a sequential process into sub-operations.
- Suboperations are performed in stage (segment).
- Each segment works parallelly

Eg:

Consider $A_i * B_i + C_i$ for $i = 1 \text{ to } 5$

Sub operations:

$$\text{Seg1: } R_1 \leftarrow A_i, R_2 \leftarrow B_i$$

$$\text{Seg2: } R_3 \leftarrow R_1 * R_2, R_4 \leftarrow C_i$$

$$\text{Seg3: } R_5 \leftarrow R_3 + R_4$$

at a time
a time.

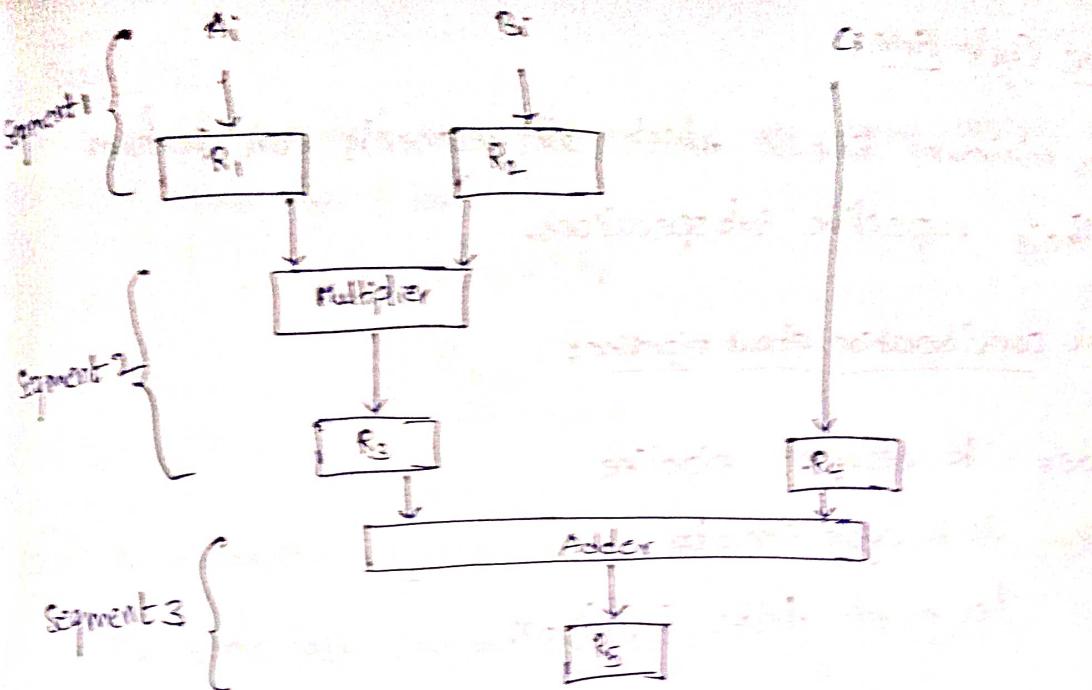
es

multiple? A_1, A_2, A_3

(GLP)

is applied

into sub-operating
ment)



| Clock Pulse | Segment 1 | | Segment 2 | | Segment 3 | |
|-------------|-----------|-------|-------------|-------|-------------------|-------------------|
| | R1 | R2 | R3 | R4 | R5 | R6 |
| 1 | A_1 | B_1 | $A_1 * B_1$ | | | |
| 2 | A_2 | B_2 | $A_1 * B_1$ | C_1 | | |
| 3 | A_3 | B_3 | $A_2 * B_2$ | C_2 | $A_1 * B_1 + C_1$ | |
| 4 | A_4 | B_4 | $A_3 * B_3$ | C_3 | $A_2 * B_2 + C_2$ | |
| 5 | A_5 | B_5 | $A_4 * B_4$ | C_4 | $A_3 * B_3 + C_3$ | |
| 6 | | | $A_5 * B_5$ | C_5 | $A_4 * B_4 + C_4$ | |
| 7 | | | | | | $A_5 * B_5 + C_5$ |

no. of cycles req with pipelining = 7

no. of cycles req without pipelining
(sequential) = $5 * 3 = 15$

Pipeline Cycle time:

It is minimum time in which all segments can perform all their respective sub operations.

General consideration about pipeline:

Consider k segment pipeline

$$\Leftrightarrow \text{clock cycle time} = t_p$$

$$\text{let no of tasks (insts)} = n$$

$$\text{time req to perform 1 task} = k * t_p$$

(i.e., output of 1st task)

$$\text{time req to perform next } n-1 \text{ tasks} = (n-1) t_p$$

$$\therefore \text{total time req for } n \text{ tasks} = (k + n - 1) t_p$$

~~Total time is~~

Note:

performance of pipeline is given by speed up ratio

$$\text{speed up ratio} = \frac{\text{non-pipeline time}}{\text{Pipeline time}}$$

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

$t_n \rightarrow$ time req for each task in sequential processing.

if ~~n~~ is large.

$$S_{\text{ideal}} = \frac{t_n}{t_p}$$

\rightarrow This is ideal case.

Special case:

If pipeline & non-pipeline req same amount of time for 1 task
i.e., $t_n = k * t_p$

$$S_{ideal} = k$$

- (Q37) A non-pipeline system takes 50 ns to process a task. The same task can be processed in a 6-segment pipeline with a clock cycle of 10 ns. Determine the speedup ratio of the pipeline for 100 tasks. What is the max speed up that can be achieved?

Sol:

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

$$= \frac{100 * 50}{(6 + 99) * 10} = 4.76$$

$$S_{max} = \frac{t_n}{t_p} = \frac{50}{10} = 5$$

29/11/20

29/11/20

Synchronous Pipeline:

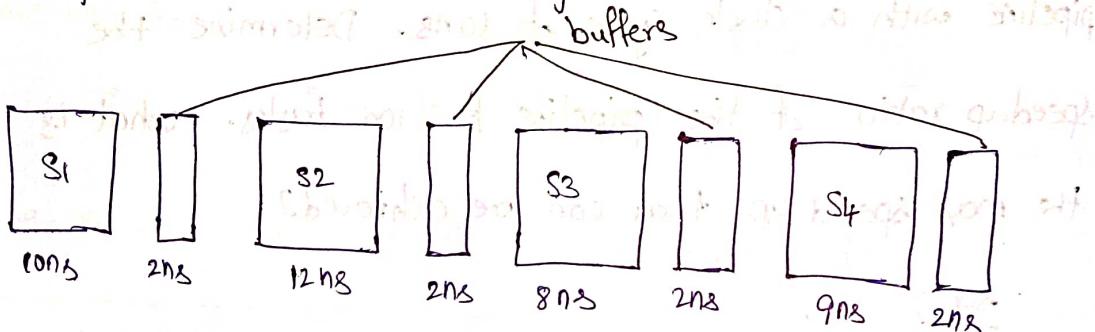
→ Each segment / stage may have different delays.

So we need to synchronize them.

Cycle time is chosen as $\max(\text{stage delay})$

→ OLP of one stage is ILP to next stage.

Since delays are different, we need intermediate buffers b/w stage to store OLP of a stage.



$$\text{Clock period} = \max(\text{seg. delay}) + \text{reg. delay}$$

$$= 12 + 2 = 14 \text{ ns}$$

But in a ~~non-pipeline~~ non-pipeline system,

$$\begin{aligned} \text{time for 1 task} &= 10 + 12 + 8 + 9 \quad (\text{seg. processing doesn't} \\ &\qquad\qquad\qquad \text{need reg delays}) \\ &= 39 \text{ ns} \end{aligned}$$

Latency:

It is amount of time after which machine takes next input.

In above example,

$$\text{latency in non-pipeline} = t_n = 39 \text{ ns}$$

$$\text{latency in pipeline} = \text{cycle time}, t_p = 14 \text{ ns}$$

Throughput :

It is no of ips processed per unit time.

In pipelined system,

$$\text{time for } n \text{ tasks} = (k+n-1)t_p$$

$$\Rightarrow \text{throughput} = \frac{n}{(k+n-1)t_p}$$

for ideal case

$$\text{throughput} = \frac{n}{n*t_p} = \frac{1}{t_p}$$

- (Q38) 5-stage pipeline; cycle time = 15 ns; find processing time for 500 tasks.

Sol:

$$\text{time} = (5 + (500-1))15 = 504 * 15 = 7560 \text{ ns}$$

- (Q39) 6-segment pipeline;

stage delays $\rightarrow 20, 26, 21, 21, 24, 28$ (in ns)

find processing time for 1000 tasks & speed up

Sol:

$$\text{cycle time} = \max(\text{delays}) = 28$$

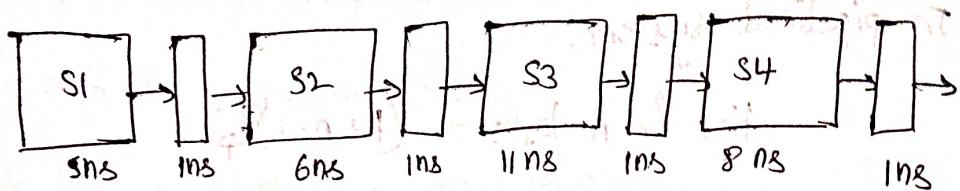
$$\Rightarrow \text{processing time} = (6 + 1000 - 1) * 28$$

$$= 1005 * 28 \\ = 28,140$$

$$\text{speed up} = \frac{(20+26+21+21+24+28)(1000)}{28140} = \frac{140000}{28140} = 4.97$$

$$\text{Max speed up} = \frac{t_n}{t_p} = \frac{20+26+21+21+24+28}{28} = \frac{160}{28} \approx 5$$

Q40
a-ii



Find approximate speed up under ideal conditions.

Sol:

$$\text{pipeline cycle time} = \max(5, 6, 11, 8) + 1$$

$$\frac{1}{\frac{1}{5} + \frac{1}{6} + \frac{1}{11} + \frac{1}{8}} = t_p = 12$$

execution time of 1 task in non-pipeline

$$\text{system} = 5 + 6 + 11 + 8 = 30$$

$$\text{ideal speed up} = \frac{t_n}{t_p} = \frac{30}{12} = 2.5$$

Note:

for n operations,

$$\text{time under pipeline system} = (k+n-1)t_p$$

$$\Rightarrow \text{avg time for 1st inst} = \frac{(k+n-1)t_p}{n}$$

$$\Rightarrow CPI = \frac{k+n-1}{n}$$

$$\text{under ideal conditions, } CPI_{\text{ideal}} = 1$$

$$\text{time per operation} = t_p$$

- (Q4) Stage delays of a 4-stage pipeline are 800, 500, 400 and 300.
 If it is upgraded to 5-stage pipeline by replacing 1st stage by 2 stages of delays 600 and 350 then ~~the~~ the throughput increase of the pipeline is _____ percent.

Sol:

4-stage:

$$t_p = 800$$

$$\text{Time for 1 inst}^n = 800 \text{ ps} \Rightarrow \text{Throughput} = \frac{1}{800}$$

5-stage:

$$t_p = 600$$

$$\text{Time for 1 inst}^n = 600 \text{ ps} \Rightarrow \text{Throughput} = \frac{1}{600}$$

$$\Rightarrow \% \text{ Increase} = \frac{800 - 600}{600} \times 100 = \frac{200}{600} \times 100 = 33.33\%$$

$$\text{Throughput increase} = \frac{1}{600} - \frac{1}{800}$$

$$\frac{1}{800}$$

$\times 100$

$$= \frac{800 - 600}{800 \times 600}$$

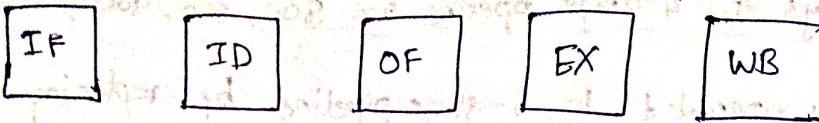
$$\times 100 = \frac{200}{600} \times 100 = \frac{1}{3} \times 100 = 33.33\%$$

Instruction Pipeline:

→ It is pipeline implemented on instruction cycle.

Assume we have 5 stages.

- Instⁿ Fetch (IF)
- Instⁿ decode & Address Calculation (ID)
- Operand fetch (OF)
- Execution (EX)
- Write back (WB)



| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|-----------------|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| I ₁ | IF | ID | OF | EX | WB | | | | | | | | | | |
| I ₂ | | IF | ID | OF | EX | WB | | | | | | | | | |
| I ₃ | | | IF | ID | OF | EX | WB | | | | | | | | |
| I ₄ | | | | IF | ID | OF | EX | WB | | | | | | | |
| I ₅ | | | | | IF | ID | OF | EX | WB | | | | | | |
| I ₆ | | | | | | IF | ID | OF | EX | WB | | | | | |
| I ₇ | | | | | | | IF | - | - | | | | | | |
| I ₈ | | | | | | | | | | | | | | | |
| I ₉ | | | | | | | | | | IF | ID | OF | EX | WB | |
| I ₁₀ | | | | | | | | | | | IF | ID | OF | EX | WB |

stall cycles/bubbles

Assume in cycle 7,

I₆ is decoded and found to be branch instⁿ.

→ whenever branch instⁿ is found, ~~next execution~~ pipeline is employed - i.e., no other instⁿ's next to branch are executed.

→ Assume I₆

Note:

→ The target instⁿ of branch is found after execution stage.

Here in this example, target instⁿ is found at the end of 9th cycle.

Let target instⁿ be I₉, so IF of I₉ is done in 10th cycle.

Note:

→ In the previous example, even if branch condition of I6 evaluated to false, we still need to fetch I7 again in 10th cycle.

→ In this example, no of instⁿs executed = 8.

$$\text{no of cycle req} = 15$$

$$\text{But no of cycles that actually needed} = 5 + 8 - 1 = 12$$

Here we req 3 extra cycles

$$\hookrightarrow (15 - 12)$$

i.e., 3 stall cycles

→ Here result of target instⁿ is found at the end of 4th stage

$$\therefore \text{no of stall cycles} = 4 - 1 = 3$$

i.e., If branch to target is evaluated in i^{th} stage, then

$$\text{no of stall cycles} = i - 1$$

Q42

Consider an instⁿ pipeline

5-13

Sol:

$$t_p = 10 + 1 = 11 \text{ ns}$$

For each branch, no of stall cycles = 4 - 1 = 3

$$\text{no of instⁿ executed} = 12 - 4 = 8$$

$$\text{no of cycle} = 5 + 8 - 1 = 12$$

1 branch inst^h \Rightarrow 3 stall cycles

$$\therefore \text{total cycles} = 12 + 3 = 15$$

$$\therefore \text{time} = 15 \times 11 = 165 \text{ ns}$$

Pipeline Hazards:

→ Hazard is a situation that prevents the next inst^h from being executed in the expected cycle.

\therefore Hazards lead to stall cycles.

3 types of hazards:

i) Structural Hazard / Resource Conflict

ii) Data Hazard / Data Dependency

iii) Control Hazard / Branch Difficulty

i) Structural Hazard / Resource Conflict:

→ This is situation when 2 different inputs (inst^h) try to use same resource at same time.

Eg: Assume MUL req 3 cycle for execution and ADD req 2 cycles for execution

MUL: IF ID OF EX EX EX WB

ADD: IF ID OF - - EX WB

Assuming that MUL & ADD using same operands

→ Also consider below case

I₁: IF ID OF EX WB

I₂: IF ID OF EX

I₃: IF ID OF

Both operand & instⁿ
needs to be fetched
from memory.

so in this case we use two cache memories

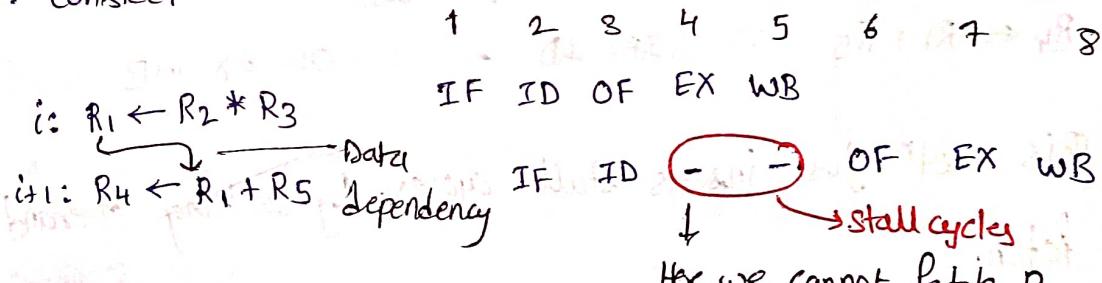
i.e., instⁿ cache & data cache

and thus avoid the conflict.

j) Data Hazard | Data Dependency

This is situation where result of an instⁿ is used as input by next instⁿ.

Eg: Consider



→ general pipeline hardware cannot detect data dependency

s/w solⁿ for data dependency

→ This solⁿ is
given by compiler

H/w solⁿ

H/w interlock
Operand Forwarding

→ In slow solⁿ, compiler generates instⁿ such that there is no data dependencies.

→ This solⁿ is known as Delayed Load.

Delayed Load:

$$R_1 \leftarrow R_2 * R_3$$

$$R_5 \leftarrow R_1 + R_4$$

Here we have seen that there are 2 stall cycles.

So compiler inserts 2 instⁿ b/w the two instⁿ.

→ So we insert either independent instⁿ or NOP instⁿ.

(This insertion is done by compiler)

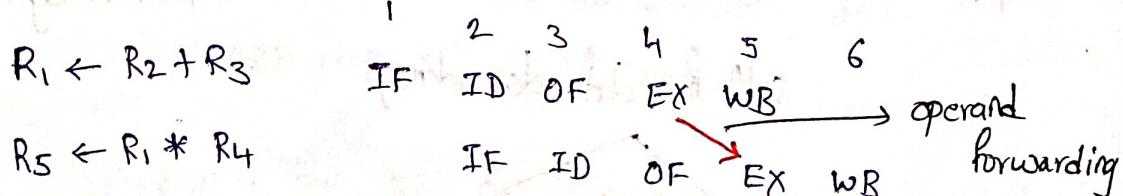
H/W Interlock:

$$R_1 \leftarrow R_2 * R_3 \quad \text{IF ID OF EX WB}$$

$$R_4 \leftarrow R_1 + R_5 \quad \text{IF ID - - OF EX WB}$$

→ This solⁿ just inserts stall cycles by locking operand fetch.

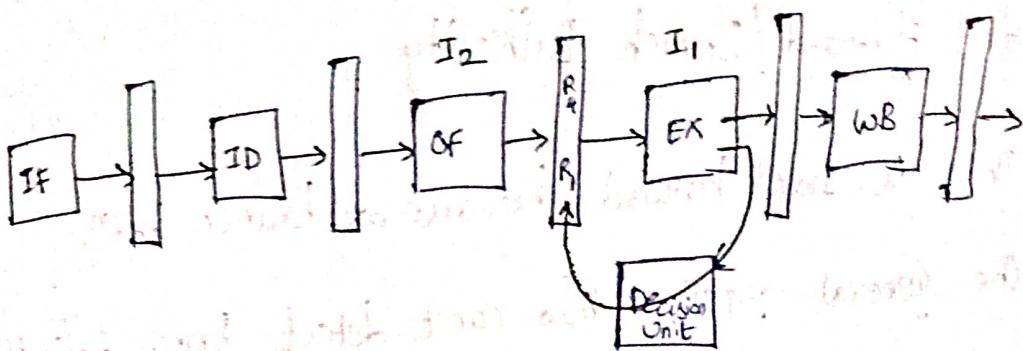
Operand Forwarding / By-passing:



At end of 3rd cycle

I₂ fetch R₁ & R₄ (However R₁'s value is wrong)

But at the end of 4th cycle, I₁'s execution phase finishes
and it ~~will~~ will forward the operand to execution phase
of I₁.



- However this solⁿ is ~~too~~ expensive
- ~~Decision unit~~ performs operand forwarding only if there is data dependency.

Note:

→ Operand Forwarding can be used only if there is ALU to ALU dependency.

i.e. ~~for other~~

→ For other dependencies operand forwarding doesn't work

i.e. for example

$$\left. \begin{array}{l} R_1 \leftarrow M[\text{add}] \\ R_4 \leftarrow R_1 * R_5 \end{array} \right\} \begin{array}{l} \text{IF ID - OF WB} \\ \text{IF ID - OF WB} \end{array}$$

Here dependency is not from ALU to ALU possible when loading can be done in execution phase
∴ we cannot use operand forwarding.

Similarly for below ~~scenario~~ scenario also operand forwarding doesn't work.

$$R_4 \leftarrow R_1 * R_5 \quad \text{IF ID OF EX WB}$$

$$M[\text{add}] \leftarrow R_4 \quad \text{IF to ID - OF EX WB}$$

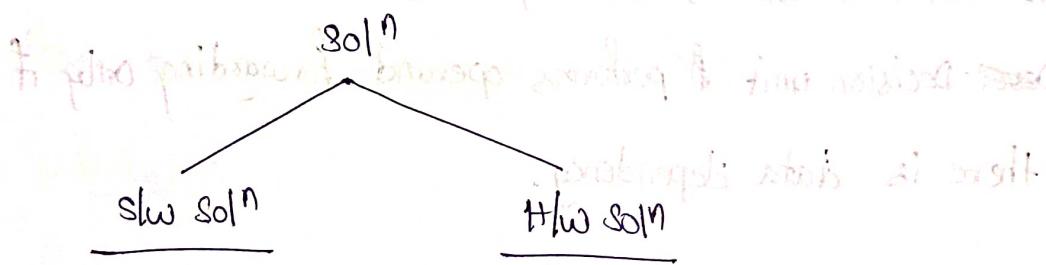
So in these cases stall cycles must exist.

2

Control Hazard / Branch Difficulty

→ This is hazard because of branch instrn.

→ The general pipeline H/w can't detect branch difficulty.



→ Delayed Branch
(provided by compiler)

→ Prefetch target instrn

→ Branch Prediction

→ Loop Buffer

→ Branch Target Buffer

Delayed Branch:

→ Based on no of stall cycles required compiler ~~inserts~~ inserts those many no of NOP instrn.

→ Some independent instrn is put in delayed slot by compiler, if no such instrn is found NOP is placed.

Prefetch Target Instruction

→ Here when branch instⁿ is detected, two pipelines will be initiated.

One pipeline is started assuming branch not taken

Other pipeline is started assuming branch is taken

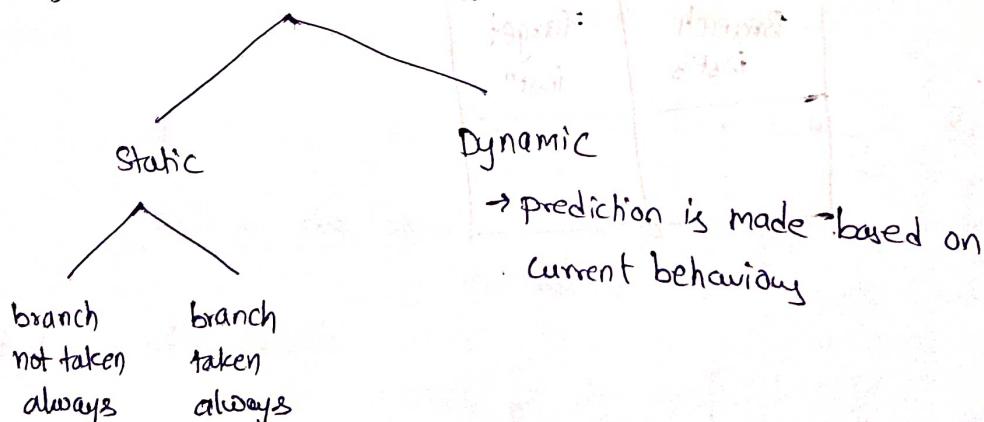
→ After finding the correct branch target, other pipeline will be discarded.

Branch Prediction:

→ Here H/w makes prediction whether branch occurs or doesn't and thus starts appropriate instⁿ.

→ If the prediction is found to be wrong, then the instⁿ is rolled back.

2 types of prediction



→ prediction is made based on current behaviour

Loop Buffer:

Consider a loop with 2 instructions that executes 100 times.

loop:
I₁:
I₂:
and so on

H/w detects that loop needs to execute 100 times.
After this the two inst's are stored in buffer and ^{loop} they are executed 100 times.

Branch Target Buffer:

- This buffer is used with branch prediction H/w.
- This buffer contains all branch inst's and their corresponding previous target inst's.
- Based on the content of this buffer prediction is made.

| Branch inst's | Target inst |
|---------------|-------------|
| : | Shifting |
| : | Shifting |

Data Hazard Classification

→ Assume there are 2 instⁿ i & j such that i executes before j.

RAW (Read After Write)

→ j reads source value before i writes it. So j gets incorrect (old) value

$$i: R_1 \leftarrow R_2 + R_3 \quad \} \text{Data dependency } i \rightarrow j$$

$$j: R_4 \leftarrow R_1 * R_5 \quad \}$$

Here j reads R_1 before i saves result in R_1 .

~~solutions discussed under data hazard~~

WAW (Write After Write)

→ j tries to write an operand before it is written by i.

$$i: R_1 \leftarrow R_2 * R_3 \quad \} \text{Write dependency } j \rightarrow i$$

$$j: R_1 \leftarrow R_4 + R_5 \quad }$$

If j writes R_1 first then final value in R_1

will be the value written by i.

WAR (Write after Read)

→ j writes to a destination before i reads it. So j incorrectly gets new value.

$$i: R_1 \leftarrow R_2 + R_3 \quad \} \text{Anti-Dependency}$$

$$j: R_2 \leftarrow R_4 * R_5 \quad }$$

This problem occurs if j writes R_2 before i reads it.

→ write dependency & anti-dependency are false

dependencies.

The sol'n for these two dependencies is register renaming.

Register renaming:

→ Here we use another register instead of common register that is causing problem.

i.e., i: $R_1 \leftarrow R_2 * R_3$

j: $R_6 \leftarrow R_4 + R_5$

(solving write dependency)

i: $R_1 \leftarrow R_2 + R_3$

j: $R_6 \leftarrow R_4 * R_5$

(solving anti-dependency)

→ Register Renaming is H/w soln.

i.e., H/w has to detect these dependencies and has to rename the registers.

Q43
G-OB7

Consider a pipelined processor with following four stage....

sol:

1 2 3 4 5 6 7 8

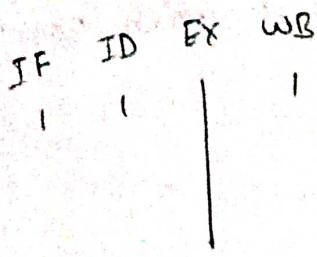
ADD IF ID EX WB

MUL IF ID EX EX EX WB

SUB IF ID -- EX WB

∴ 8 cycles

shortcut:



| | | |
|-----|---|-----|
| ADD | 1 | - 0 |
| SUB | 1 | - 0 |
| MUL | 3 | - 2 |

$$\text{no of cycles req (without stalls)} = 4 + 3 - 1 = 6$$

$$\text{no of mul operations} = 1$$

$$\Rightarrow \text{no of stalls} = 1 * 2 = 2$$

$$\therefore \text{total cycles req} = 6 + 2 = 8$$

(This shortcut works only if one phase takes extra cycles)

Pipeline Efficiency:

→ Because (stall ~~sup~~ cycles) efficiency of pipeline reduces.

$$\text{efficiency of pipeline} = \frac{s'}{s} \times 100$$

$s' \rightarrow$ speed up with hazards

$s \rightarrow$ speed up without hazards

Q44

Consider a 5-stage pipeline;

Consider execution of 1000 inst's of which 200 inst's cause 2 stall cycles each.

1. Calculate CPI of pipeline.

Sol:

No of cycles for execution

$$= (k+n-1) + \text{extra cycles}$$

$$= (5+1000-1) + (200 \times 2)$$

$$= 1404 \text{ cycles}$$

$$\Rightarrow CPI = \frac{1000}{1404} = 1.404$$

2. If pipeline cycle time is 3ns, what is avg inst's execution time.

Sol:

$$CPI * t_p = 1.404 * 3 = 4.212$$

3. Find CPI under ideal conditions.

~~CPI_{ideal} = (for any pipelined system)~~

$$CPI_{\text{ideal}} = (n + \text{extra cycles}) / 1000$$

$$= (1000 + 200 \times 2) / 1000$$

$$= \frac{1400}{1000}$$

$$= 1.4$$

Note:

$$CPI_{\text{ideal}} = 1 + (\text{stall freq}) * (\text{no of stall cycles})$$

Consider a 6-stage pipeline . . .

Sol:

$$CPI_{\text{(pipeline)}} = 1 + (0.25)(2) = 1.5$$

$$\text{1 instn execution time} = 1.5 t_p$$

$$\text{1 instn execution time in non-pipeline} = 6 * t_p$$

$$\Rightarrow \text{speed up} = \frac{6 t_p}{1.5 t_p} = 4$$

$$\text{Efficiency of pipeline} = \frac{s^1}{s_{\max}} \times 100 = \frac{4}{6} \times 100 = 66.66\%$$

instn

Q46

Consider below stage delay for instn I_1, I_2, I_3, I_4 for 4-stage pipeline.

Find no of cycle for below program

for($i=1$; $i \leq 2$; $i++$)

I_1, I_2, I_3, I_4);

| | s_1 | s_2 | s_3 | s_4 |
|-------|-------|-------|-------|-------|
| I_1 | 2 | 1 | 1 | 2 |
| I_2 | 1 | 3 | 1 | 3 |
| I_3 | 2 | 2 | 2 | 1 |
| I_4 | 1 | 3 | 1 | 2 |



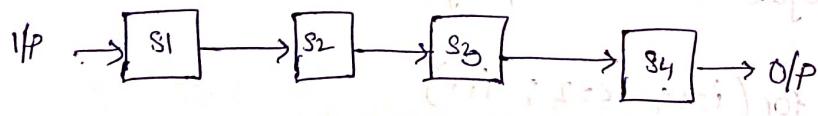
Q2:

| | S1 | S2 | S3 | S4 |
|----|----|----|----|----|
| I1 | 2 | 3 | 4 | 6 |
| I2 | 3 | 6 | 7 | 10 |
| I3 | 5 | 8 | 10 | 11 |
| I4 | 6 | 11 | 12 | 14 |
| I5 | 8 | 12 | 13 | 16 |
| I6 | 9 | 15 | 16 | 19 |
| I7 | 11 | 17 | 19 | 20 |
| I8 | 12 | 20 | 21 | 23 |

solution

Non-Linear Pipeline (chain of tasks with tasks interleaved)

Linear pipeline:

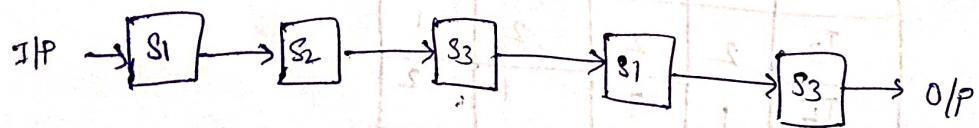


A task goes through each segment in order, exactly once.

Non-linear pipeline:

A segment may be used multiple times during execution of one task.

Eg:



Non-linear pipeline is represented using reservation table

This table shows which segment is used at which cycle.

Consider 3-stage non-linear pipeline

| clock cycle | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|-------------|---|---|---|---|---|---|---|---|
| S1 | X | | | | | X | X | |
| S2 | | X | X | | | | | |
| S3 | | | X | X | X | | | |

Latency: no of time units (clock cycles) b/w two initiations of a pipeline providing ip in pipeline.
pipeline is called latency. Latency is always non-negative.

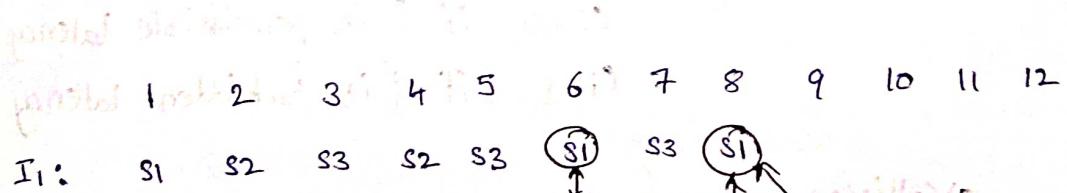
~~Collision~~: without any hazard, latency in linear pipeline is 1 cycle.

Collision: When two or more initiations are done at same pipeline stage at same time, then we say a collision has occurred.
This is resource conflict.

Permissible Latency: Latency which does not cause a collision

Forbidden Latency: Latency which causes collision.

Eg: for above reservation table, consider



(latency=5) I₂:

(latency=2) I₃:

(latency=7) I_n:

So here for S_1 ,

forbidden latencies are 2, 5, 7

an easy way for determining these latencies is

S_1 is used in 1, 6, 8 cycles

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|
| 1 | X | | | | | | | | | | | | | | |
| 2 | | X | | | | | | | | | | | | | |
| 3 | | | X | | | | | | | | | | | | |

$$\therefore \text{forbidden latencies } 6-1=5$$

$$8-6=2$$

$$8-1=7$$

Similarly we calculate for other stages also

$$\text{For } S_2, 4-2=2$$

$$5-3=2$$

$$7-5=2$$

$$7-3=4$$

$$\therefore \text{Forbidden latencies are } 2, 4, 5, 7$$

$$\Rightarrow \text{permissible latencies are } 1, 3, 6, 8$$

Collision Vector:

Collision vector is an n-bit vector provided : ~~permitted~~ ~~forbidden~~

c_n, \dots, c_2, c_1 such that

$c_i = 0$, if i is permissible latency

$c_i = 1$, if i is forbidden latency

\Rightarrow Collision vector for this example is

| | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 |
|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | |

State Diagram:

→ Take a state S .

→ For every 0 at position i

* shift collision vector of state S to the right by i positions.

* Take bitwise OR with original collision vector.

~~Ex:~~ State 01011010:

(i) 0 at position 1:

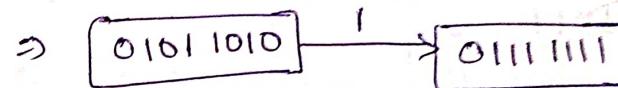
: right shift by 1

$\Rightarrow 00101101$

bitwise OR

$$\begin{array}{r} 01011010 \\ 00101101 \\ \hline 01111111 \end{array}$$

↳ This is a new state.



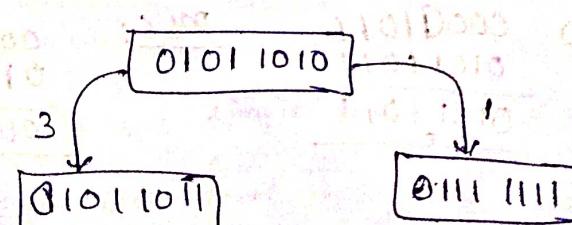
(ii) 0 at position 3:

: right shift by 3

original state from 01011010
when possible?

$$\begin{array}{r} 00000000 \\ 000001011 \\ 01011010 \\ \hline 01011011 \end{array}$$

↳ This new state

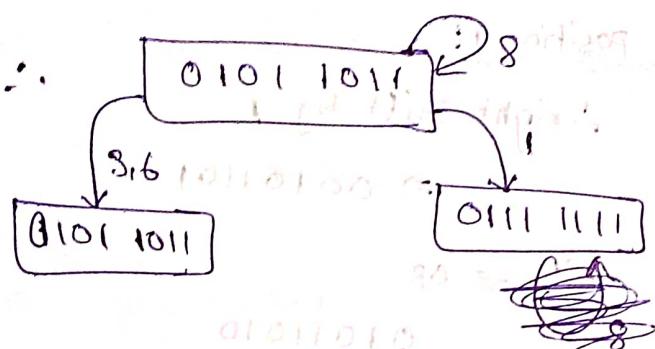


iii) 0 at position 6:

$$\text{R-shift} \Rightarrow \begin{array}{r} 0000\ 0001 \\ 0101\ 1010 \\ \hline 0101\ 1011 \end{array}$$

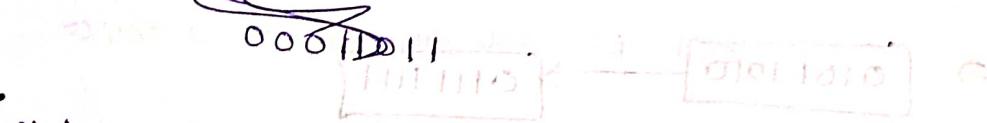
padding 3rd triplet added & start ~~new state~~ This is same state reached by B.

$$\text{R-shift} \Rightarrow \begin{array}{r} 0000\ 0000 \\ 0101\ 1011 \\ \hline 0101\ 1011 \end{array}$$



~~State 01011011:~~

0 at pos 3:



state 01011011 :

$$\text{R-shift} \Rightarrow \begin{array}{r} 0000\ 0000 \\ 0101\ 1011 \\ \hline 0101\ 1011 \end{array}$$

$$\begin{array}{r} 0101\ 1011 \\ \hline 0101\ 1011 \end{array}$$

:8 padding to 0's

Here we must take original
Collision vector.

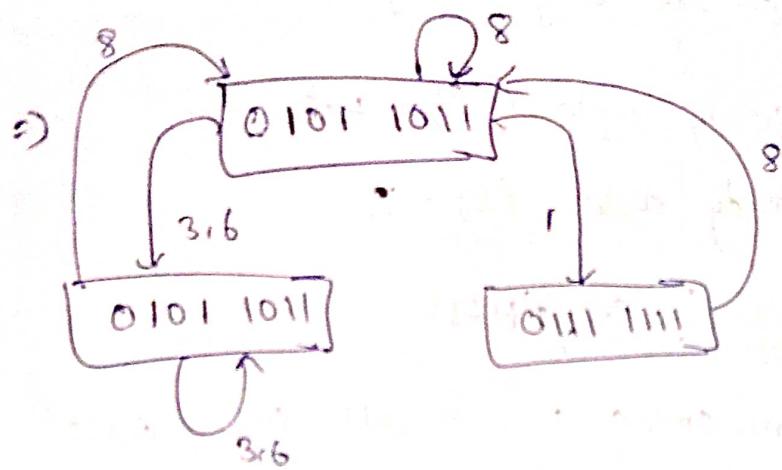
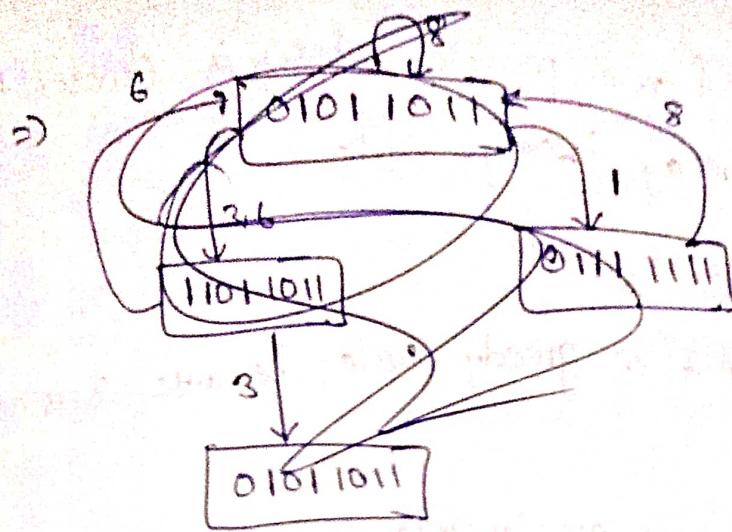
state 01011011 :

pos 3: R-shift \Rightarrow

$$\begin{array}{r} 0000\ 1011 \\ 0101\ 1011 \\ \hline 0101\ 1011 \end{array}$$

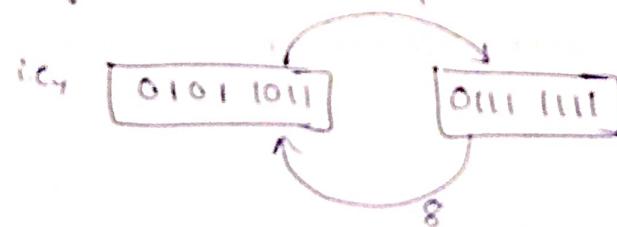
pos 6:

$$\begin{array}{r} 000000\ 01 \\ 01011011 \\ \hline 01011011 \end{array}$$



Simple Cycle:

A simple cycle is a latency cycle in which each state appears only once.



$\therefore (1, 8)$ is a simple cycle.

slly (8) is a simple cycle

slly (3, 8) (6, 8), (3), (6) are also simple cycle.

Greedy Cycle:

A greedy cycle is a simple cycle whose edges are all made with minimum latencies from their respective starting states.

0000001
011011
011011

i.e., At each state of min. possible latency, is chosen then it is called greedy, latency cycle.

$\Rightarrow (1, 8)$ is a greedy cycle.

But $(3, 8)$ is not a greedy cycle because 3 is not minimum.

Similarly (3) is also a greedy cycle.

avg of greedy cycle $(1, 8) = 4.5$

avg of greedy cycle $(3) = 3$

Minimum Average Latency (MAL)

It is minimum among avg of all greedy cycles.

i.e., here min avg latency = 3

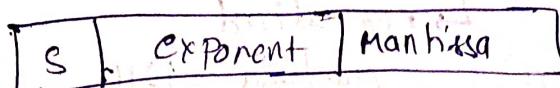
Note: The cycle $(1, 8)$ means provide 1st ip with latency 1, then with latency 8 then with latency 1 and so on..

~~Notes~~

Floating point Representation

→ Floating point representation can represent a larger range of numbers compared to fixed point number representation given same no. of bit.

→ General representation of floating point representation is



1-bit sign bit (positive or negative)

- Mantissa is signed normalized (implicit/explicit) fraction number.
- Exponent is stored in biased form.

Biased Exponent:

→ Biased exponent says that exponent should be only unsigned

Q. By consider 4-bit exponent

range with 4-bits : -8 to 7

now we need to change it to positive range

⇒ req. range : 0 to 15

i.e., we add a bias 8

$$\Rightarrow E = e + 8$$

↓ ↓ ↗
 biased original bias
 exp exp
 (or)
 stored exp

→ If exponent is of k-bits then bias = 2^{k-1}

i.e., exponent is stored as excess 2^{k-1} code.

→ Mantissa is number after decimal point.

Normalized Mantissa:

Explicit Normalization (Default consideration)

- In explicit normalization a number is represented with $0 \cdot \underline{\quad}$ and the digit follows.

$(101.11)_2 \xrightarrow{\text{explicit norm.}} 0 \cdot 10111 \times 2^3 \rightarrow \text{original exponent}$

\downarrow normalized
mantissa = 10111

$\Rightarrow e = 3$ (original exponent)

Implicit normalization / Implicit preceding:

Here a number is represented as $1 \cdot \underline{\quad}$ and before immediate digit following ~~the~~ can be any.

$(101.11)_2 \xrightarrow{\text{implicit norm}} 1 \cdot 0111 \times 2^2$

$\Rightarrow M = 0111$

$E = 2 + \text{bias}$

Value formula:

| | | |
|---|---|---|
| S | E | M |
|---|---|---|

Explicit normalization:

$$\text{value} = (-1)^S \times 0.M \times 2^{E-\text{bias}}$$

Implicit normalization

$$\text{value} = (-1)^S \times 1.M \times 2^{E-\text{bias}}$$

Q) Consider a 16-bit register used to store floating point numbers. The mantissa is normalized signed fraction number. Exponent is represented in excess-32 form. what is 16-bit value for $(13.5)_{10}$ in this register.

Sol:

$$\text{bias} = 32$$

\Rightarrow 6 bits of exponent

$$(13.5)_{10} = (1101.1)_2$$

explicit normalization: 0.11011×2^4

$$E = e + 32 = 36$$

$$E = 100100$$

$$S = 0$$

| | | | |
|---|---|---|---------------|
| 1 | 6 | 9 | |
| S | E | M | \Rightarrow |

| | | |
|---|--------|-----------|
| 0 | 100100 | 110110000 |
|---|--------|-----------|

Q) Consider

| | | |
|---|---|---|
| S | E | M |
| 1 | 6 | 9 |

i) Find max value that can be represented by above register.

$$(-1)^S \times 0.M \times 2^{E-\text{bias}}$$

for val to be max

$$M = 1111\ldots1111$$

$$E = 1111\ldots1111$$

$$(-1)^0 \times 0.1111\ldots1111 \times 2^{63-32}$$

$$0.1111\ldots1111 \times 2^{31} = 1111\ldots1111 \times 2^{22}$$

$$= 551 \times 2^{22}$$

$$= (2^9 - 1) \times 2^{22}$$

$$= 2^{31} - 2^{22}$$

(i) Find min value?

Sol:

Here every thing remain same except sign

$$\therefore -(2^{31}-2^{22})$$

* *
* (ii) Find smallest positive value?

0 00000 | 0000 0000

$$val = (-1)^0 \times 2^{-32} \times 0.100000000$$

$$= 0.1 \times 2^{-32}$$

$$= 0.1 \times 2^{-33}$$

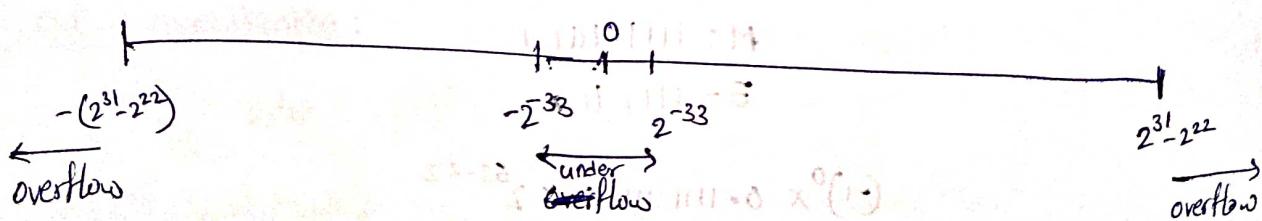
In explicit normalized form, mantissa can never be all zeroes since after decimal point there must be one.

(iv) Find largest negative value

i.e., -2^{-33}

Note:

In above question, the range of number than can be represented by above format is



→ If fixed point, with 16 bits we, max value is $(2^{15}-1)$
but here max value is $2^{31}-2^{22}$

i.e., range of number with floating point represn is more.

Note:

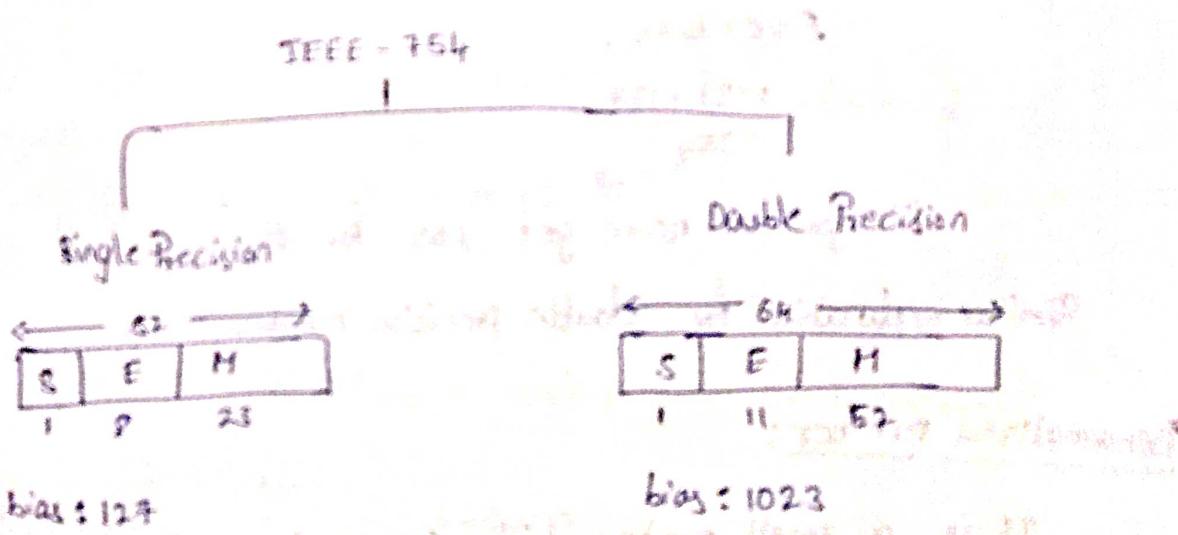
- If more number of bits are used in mantissa, then precision (accuracy) is more.
- If more number of bits are used in exponent, range is more.

Disadvantages of conventional floating point representation:

- It cannot represent 0.
- It cannot represent infinity.
- It cannot store a number which cannot be normalized.

IEEE - 754 Floating point Representation:

This, overcoming the disadvantages of conventional representation.



→ The bias is different for different levels of precision.

E = all 0's
or
E = all 1's

special Numbers

| S | E | M | number |
|-----|---------|---------|-------------------------|
| 0/1 | 00...0 | 0...00 | ± 0 (zero) got zero |
| 0/1 | 11...1 | 11...11 | $\pm \infty$ |
| 0/1 | 00....0 | M#0 | Denormalized |
| 0/1 | 11....1 | M#0 | fraction |
| 0/1 | rest | rest | NAN (Not a number) |

Explanation of terms: M is mantissa, E is exponent.

NAN:

Consider single precision, ~~normalizing float, half, part - 2^127~~

interior digits if $e=127$ (i.e., max possible exponent)

$$E = e + \text{bias}$$

$$\begin{aligned} &= 127 + 127 \\ &= 254 \end{aligned}$$

i.e., you can never get 255 for E .

Similar explanation for double precision number.

Denormalized number:

It is a small number which cannot be normalized.

Ex: Consider below example for single precision

$$0.00000\dots11$$

\downarrow implicit normalization

$$1 \cdot 1 \times 2^{-130}$$

$$\text{i.e., } M=1; e=-130; E=-130+127$$

$$= -3$$

but stored exponent (biased) can have only ~~non-negative values~~ non-negative values.

In conventional representation, storing this kind of numbers would be possible.

↑ But in IEEE-754,

we have to normalize the number upto 23 bits

Consider we get 0.011×2^{126}

In this case we store it below

$$E = 0000\ 0000$$

$$M = 011$$

Value Formula:

Implicit normalization:

$$\text{value} = (-1)^E \times 1.M \times 2^{E - \text{bias}}$$

bias = 127 or 1023

Denormalized:

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

$$(-1)^E \times 0.M \times 2^{-1022}$$

(Q51) Represent $(-14.25)_{10}$ in single precision (in hexadecimal)

Ans: $C1640000$

(Q52) Find value of below single precision representation

01000000111000...00

Ans: 28

(Q53) Find value of below single precision representation

0000000001100...00

Ans? (Q53) $E \rightarrow 0$ all 0's

$$\text{val} = (-1)^0 \times 0.11 \times 2^{-126} = 3 \times 2^{-128}$$

Q54 Represent $(1 \cdot 0)_{10}$ in single precision

Ans:

| | | |
|---|-----------|----------|
| 0 | 0111 1111 | 000...00 |
|---|-----------|----------|

Q55 Find max value possible in IEEE-754 (single precision)

Ans:

S E M

| | | |
|---|-----------|---------|
| 0 | 0111 1110 | 111...1 |
|---|-----------|---------|

$$\Rightarrow \text{value} = 1 \cdot 111\ldots 1 \times 2^{2^{\text{E}}-1-2^{\text{M}}}$$

$$= \underbrace{111\ldots 1}_{2^{\text{E}}} \times 2^{104}$$

$$= (2^{\text{E}}-1)(2^{104}) = (2^{128}-2^{104})$$

Q56 Fin min value possible in single precision

Ans: $-(2^{128}-2^{104})$

Q57 Find min positive val for IEEE-754 single precision.

Sol:

This is possible under denormalized form

| | | |
|---|-----------|-----------|
| 0 | 0000 0000 | 000...001 |
|---|-----------|-----------|

$$\text{value} = (-1)^0 \times 0.000\ldots 001 \times 2^{-126}$$

$$= 1 \cdot 0 \times 2^{-149}$$

$$= 2^{-149}$$

Practice Questions:

(Q58) Consider a digital computer which supports 2^{32} address instⁿs. Size of each instⁿ is 38 bits. The system supports a word addressable memory with word size of 32 bits. Then maximum allowable size of memory is _____ kb?

Ans: 8

For same question if memory is byte addressable then

Ans: 2

(Q59) Consider a digital computer which supports max of 32KB word addressable with word size 32-bits. The Computer supports 40 distinct instⁿs each of which has 4 fields

i) Opcode

ii) Mode field to specify one of 7 add. modes

iii) Reg field field to specify one of 24 processor reg.

iv) Mem add field.

The size of instⁿ is _____ bits

Ans: 27

(Q60) Consider a hypothetical processor with largest instⁿ length being 32 bits and total registers are 16. R0-R15. Processor supports following instructions:

ADD R_i, R_j

SUB R_i, R_j

AND R_i, R_j

NOT R_i

MOV R_i, R_j

LOAD Address

Store Address

JUMP address

J2 Address

The max no of address lines on this processor is 2²⁸

Sol:

If inst^h size is fixed \Rightarrow opcode length varies;

If inst^h size is variable \Rightarrow opcode size is fixed

Here we have variable length inst^h

for inst^h 1, 2, 3, 5, the format is

| | | |
|--------|------|------|
| opcode | add1 | add2 |
|--------|------|------|

Instruction format 4 + 4 + 4

changes according to inst^h size = 12 bits know after observing

for inst^h 4 the format is

| | |
|--------|------|
| opcode | add1 |
|--------|------|

4 + 4

inst^h size = 8

for inst^h's 6, 7, 8, 9

| | |
|--------|--------|
| opcode | add |
| 4 | 32 - 4 |
| | = 28 |

\therefore Max no of address lines = 2²⁸

Inst^h = 26 bits, add = 2²⁸ bits, 2²⁸ address lines total bno of add

1-add & 0-add inst^hs are supported.

Find range of 1-add & 0-add inst^hs possible

1-add insts their word length will be $0\text{-add} + 2^7 = 255$

Q17) If digit of sum of 8 bit number is 0 or 1, then 0-add
will be possible
↓
it is given to implement in 8 bit only ~~if only 1-add~~ 1-add is given
we can't use all 512
present word length is 8 bits

Combinations of 0-add
are also supported

$$\Rightarrow \max \text{0-add} = 511 \times 2^7$$

Range: 1 to 511

Range: 1 to 511×2^7

whenever it is given that a certain k-add insts are supported
then there must be at least 2^{k+1} insts of that type
to implement word length 8 bits
which is given after the formula to go the next. Then if 2-add, 1-add, 0-add, insts are supported.

Q62) 2-add, 1-add, 0-add, insts are supported.

inst - 20 bits; add - 7 bits

find max no of 0-add insts possible.

Sol:

Ans:

$$8063 * 2^7$$

Q63)

inst - i bits; add - a bits

0-add, 1-add, 2-add insts are supported.

If no of 2-add insts = x

" " 1-add " = y

" " 0-add " = ?

Sol:

total no of combinations = 2^i

0-add:

| | | | |
|---|---|---|---|
| i | 2 | a | a |
|---|---|---|---|

1-add:

| | | |
|---|---|---|
| i | a | a |
|---|---|---|

$$x \times 2^{2a} + y \times 2^a + z = 2^i \quad 0\text{-add: } i$$

$$z = 2^i - x \cdot 2^{2a} - y \cdot 2^a$$

Q64) Consider a PC-relative mode-type branch inst^h which takes branch on address 740 in memory. The inst^h has offset value 16_0 . What is the address of this inst^h in memory if inst^h is stored on 4 memory locations.

Ans: 576

Q65) Consider a 6-word inst^h

| | | | | |
|--------|-------|-------|------|------|
| opcode | mode1 | mode2 | add1 | add2 |
|--------|-------|-------|------|------|

The first operand (destination) uses register indirect mode and second operand uses indirect mode. Assume each operand is of size 2 words, each add is of 2 words and main memory takes 50 ns for 1 word.

Total memory access time required is

1. Fetch cycle of inst^h
2. Execution cycle of inst^h
3. Inst^h cycle of inst^h

sol:

$$1. 6 * 50 = 300 \text{ ns}$$

$$2. \text{Reading: } 2 * 50 + 4 * 50$$

$$\text{Writing: } \frac{2 * 50}{4} = 25 \text{ ns}$$

$$400$$

$$3. \text{total: } 700 \text{ ns}$$