

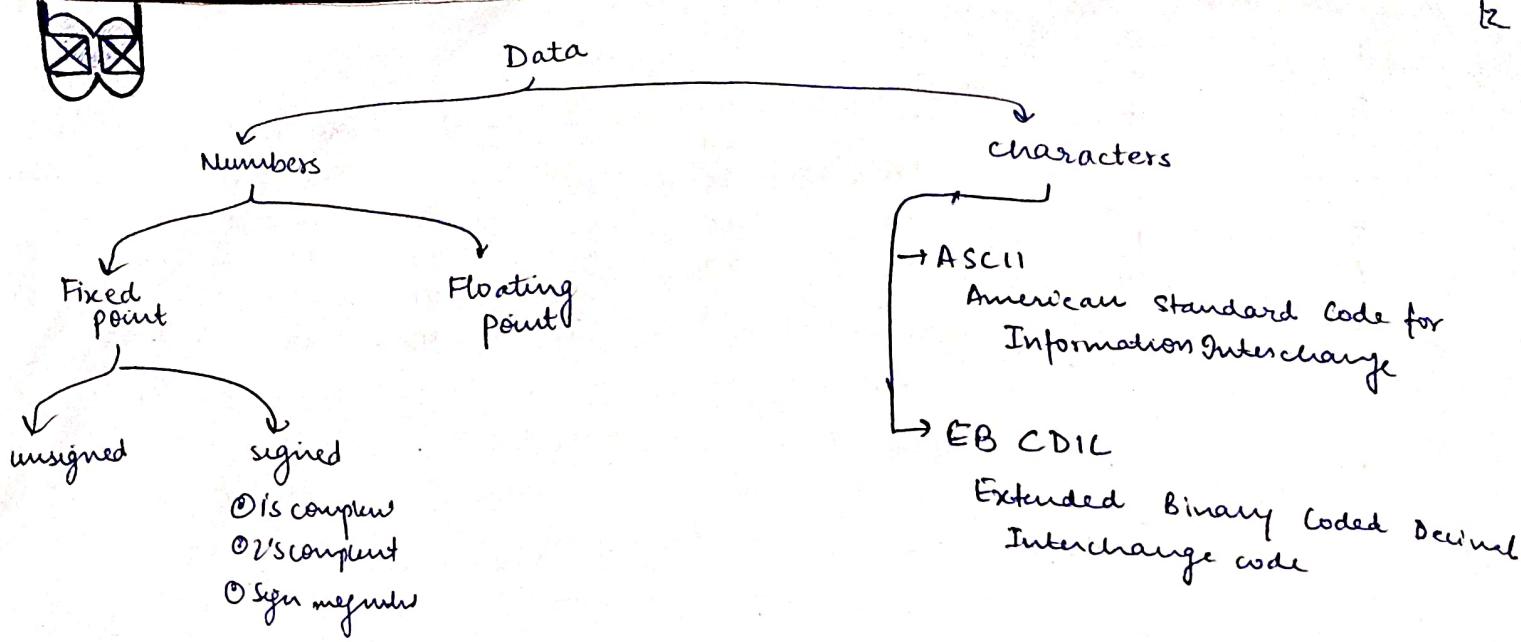
1	8	15	22	29	36	43	50	57	64
2	9	16	23	30	37	44	51	58	65
3	10	17	24	31	38	45	52	59	66
4	11	18	25	32	39	46	53	60	67
5	12	19	26	33	40	47	54	61	68
6	13	20	27	34	41	48	55	62	69
7	14	21	28	35	42	49	56	63	70

COMPUTER ORGANIZATION AND ARCHITECTURE

- [Architecture] 1. Introduction
- 2. Floating Point Representation
- 3. Basics of Computer System & Microoperations
- 4. Instruction and Addressing Modes
- 5. CPU, Data Path & Control Unit.
- [Organization] 6. IO Organization
- 7. Memory & cache
- 8. Pipelining.

ARCHITECTURE —
conceptual design & fundamental operational structure.

CPU design, Instructions, Addressing modes,
Data Format.



Stored Programs Architecture.

The program currently executing should be stored in memory (RAM)

2 architectures —

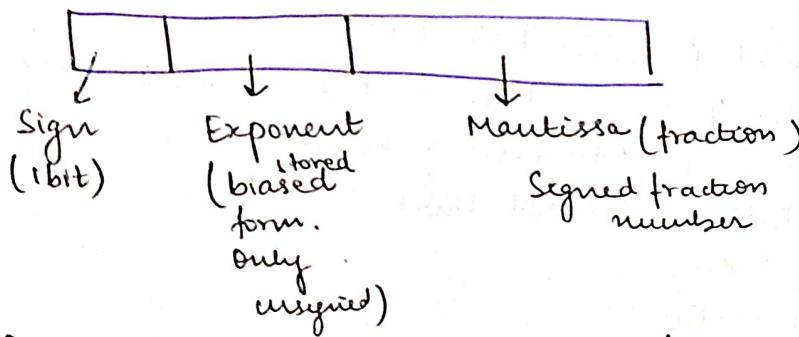
- Von-neumann architecture
- Harvard Architecture

VNA

Not required for GATE

Floating Point Representation

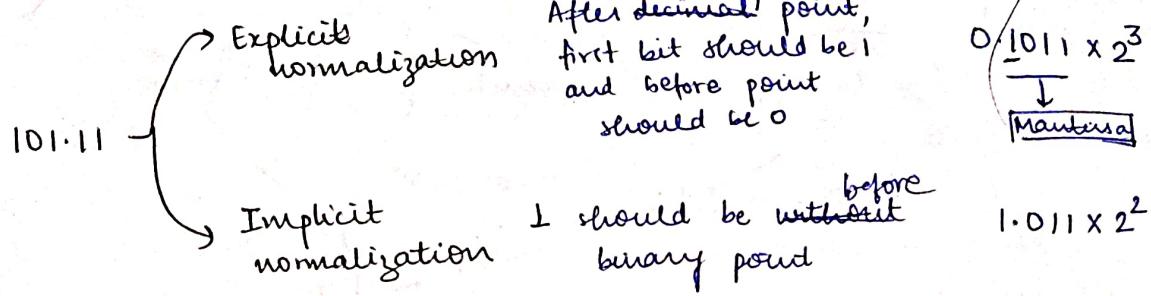
provides larger range of numbers with limited number of bits.



Bias value is added in exponent.

Bias value = $2^n - 1$ where n = no. of bits in exponent.
also known as excess notation

Mantissa is normalized.



Explicit normalization is default (if not given)

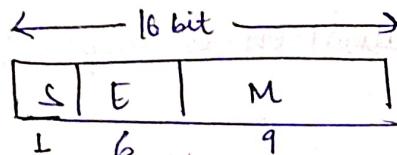
- ① Binary point is not stored
- ② Digits before binary point are not stored.
- ③ Base of binary number is not stored

$$\text{Value in implicit normalization} = (-1)^s \times 0.M \times 2^{E-\text{bias}}$$

$$\text{Value in external normalization} = (-1)^s \times 0.M \times 2^{E-\text{bias}}$$

Question: Consider a 16 bit register used to store floating point number. The mantissa is normalized signed fractional number. Exponent is represented in excess 32 form. What is the 16 bit value for $+19.25_{10}$ in this register?

52 1001011
1 000101
k=5



bias = +32

$$\therefore \text{No. of bits in } E = 5$$

$$2^{k-1} = 32 \Rightarrow k-1 = 5$$

$$k = 6$$

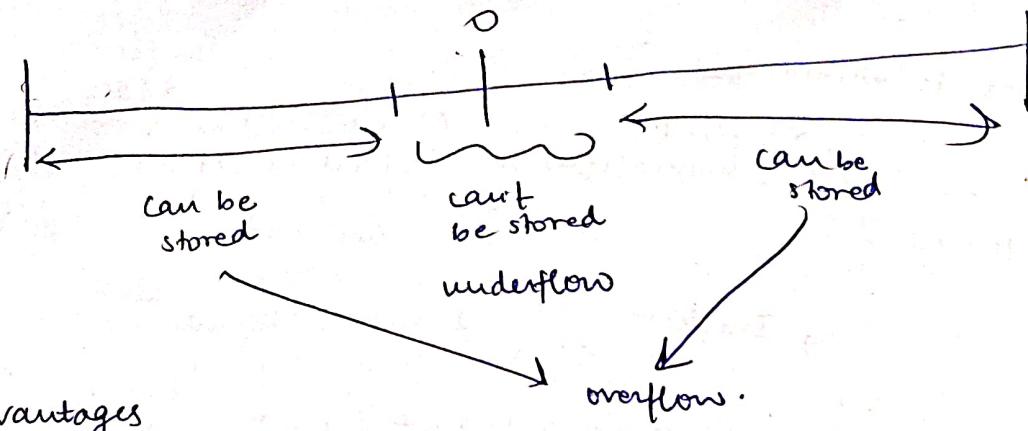
$$19.25 = 10011.01$$

not mentioned : explicit normalization

$$\therefore 19.25 = 10011.01 \times 2^5$$

$$E = 5 + 32 = 37 = 100101$$

$$0\ 100101\ 00110\ \underline{10100}$$



Disadvantages

- can't store 0
- can't store values very close to 0

IEEE - 754 Floating point Representation

① Standard.

Single Precision
(32 bits)

S	E	M
1	11 1023	23

Bias = 127

Double Precision
(64 bits)

S	E	M
1	11 1023	52

Bias = 1023

For IEEE, default normalization is implicit.

Special numbers

S	E	M	Number
0	00.....0	0....0	+0
1	00....0	0...0	-0
0	11....1	0...0	+∞
1	11....1	0....0	-∞
0/1	11....1	M ≠ 0...0	Not a number (NaN)
0/1	00....0	M ≠ 0...0	Denormalized number (Fraction no.)
0/1	E ≠ 00..0 and E ≠ 100..0	M = xx...x	Normal no. (Implicit normalize)

$$\text{Value (implicit)} = (-1)^s \times 1.M \times 2^{E-\text{bias}}$$

$$\text{Value (denormalized)} = (-1)^s \times 0.M \times 2^{-126 \text{ or } -1022}$$

very very small number that cannot be normalized

$$\text{Smallest possible value of } E = (000\cdots)_{10} = 1$$

$$\text{hence, } e = 1-127 = -126$$

$$\text{If number} = 0.\underbrace{00000}_{1} \cdots 011 \Rightarrow \text{Implicit norm} \\ M = 1.1 \\ e = -126$$

\therefore the number is stored in denormalized form.

$$0.\underbrace{0000}_{128 \text{ times}} \cdots 011 = 0.0011 \times 2^{-126}$$

$$e = -129 + 127 = -2$$

\therefore not possible to normalize

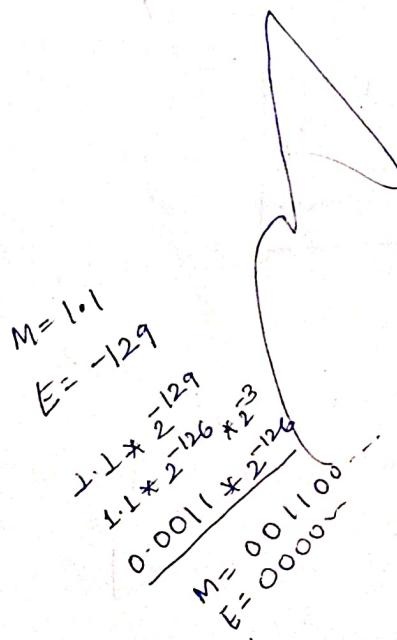
$$M = 0011$$

$$E = 000000\cdots$$

For denormalized no.,

$$e = -126$$

always



COMPONENTS OF COMPUTER

- ① CPU
 - Control unit (gives signal to different components)
 - ALU (computation)
- ② Memory
 - Primary (RAM, ROM)
 - Secondary (HDD)
- ③ I/O Devices
 - Input
 - Output.

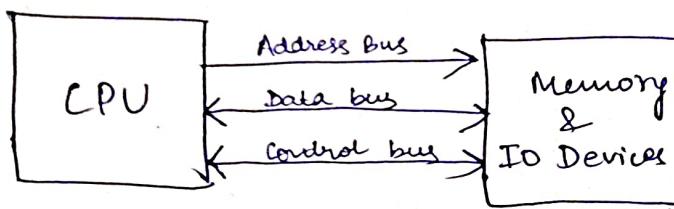
Other components

- ④ System Bus
- ⑤ Registers

System Bus

- ⑥ Collection of communication lines between components of computer
- ⑦ 3 Types

- Address Bus
 - Data Bus
 - Control Bus
- based on component content carried.



Operations that can be performed on memory -

- Read [content goes from memory to CPU]
- Write [data goes from CPU to memory]

If a memory has 2^n no. of addresses,
Address length (No. of bits in the address bus) = n bits.

17

Memory Access : Read

① CPU sends address to memory through Address Bus.

② CPU sends control signal to memory.

③ Memory performs read operation and sends ^{data} through data bus.

Memory Access : Write

① CPU sends address to memory through address bus

② CPU sends data to memory through data bus

③ CPU sends control signal to memory via control bus

④ Memory performs write operation.

CPU Registers

① small storage locations within CPU.

② 2 types

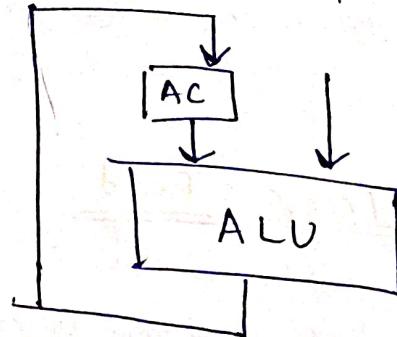
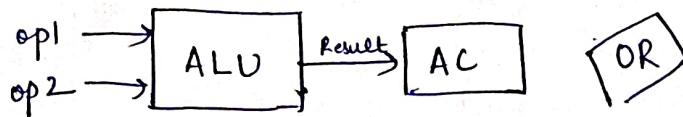
→ General Purpose Registers (GPR) R₀, R₁, R₂...

→ General & Special Purpose Registers (SPR)

- ① → Accumulator
- ② → Program Counter
- ③ → Instruction Register
- ④ → Stack Pointer
- ⑤ → Flag Register / Program Status Word (PSW)
- ⑥ → Address / Memory address Register (MAR)
- ⑦ → Data / Memory Data Register (MDR / MBR)

① Accumulator

used to store the result of ALU and sometimes one of the operands for ALU too.



2 different architecture based on system design

② Program Counter

used to store address of next instruction to be executed.

③ Instruction Register

used to store currently executing instruction.

④ Stack Pointer

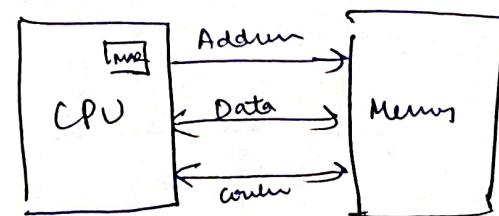
used to store address of the top of stack

⑤ Flag Register | Status Register | Program Status Word

used to store status of ALU & result

Z	S	ca
↓	↓	↓	
Is result 0?	sign of result	carry of result	
			(condition code)

(condition check.
also known as
condition code)



⑥ Memory Address Register

used to send address to memory

⑦ Memory Data Register

used to send data from CPU to memory or memory to CPU

Micro-Operations

The operations executed on values stored in registers

Symbolic Notation to describe microoperations = Register Transfer Language

Register Transfer

$$R1 \leftarrow R2$$

Two mutually exclusive operations can be performed simultaneously

$$R1 \leftarrow R2, \quad PC \leftarrow PC + 1$$

Memory Transfer

Memory Read

$$R \leftarrow M[Add]$$

Memory write

$$M[Add] \leftarrow R2$$

Q:- If the content of register R1 is 5 and content of memory address 1000 is 20. Then, the content of register R2 after following code execution is

$$R1 \leftarrow R1 + 1$$

$$R1 = 6$$

$$R2 \leftarrow R1 + M[1000]$$

$$R2 = 6 + 20 = \underline{\underline{26}}$$

Instruction Length

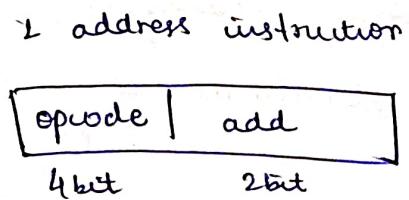
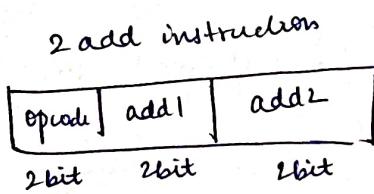
Fixed length instruction

fixed variable length opcode.

Variable length instruction

fixed length opcode.

Q: Consider a system which supports 2 address instructions and 1 address instructions both. The system has 6 bit instructions and 2 bits addresses. If there are three 2-address instruction in the system, then, the 1-address instructions the system can support?



If CPU receives an instruction, it needs to decide if it is 2 address instruction or 1 address instruction.

$$\text{No. of bits in } 2 \text{ add ins} = 2 \quad [\because 2 \text{ bits are checked}]$$

Only 3 2 address instructions are possible
 \therefore unused combination of bits = 1

\therefore First 2 bits of opcode of 1 address instruction are fixed.

$$\therefore \text{Max no. of 1 address instruction} = 2^2 = 4$$

$$\text{Min no. of 1 address instruction} = \underline{\underline{1}} \quad (\text{minimum is always 1})$$

No. of 2 address instructions supported	Unused opcode	No. of 1 address instruction supported
3	1	$1 \times 2^1 = 2$
2	2	$2 \times 2^2 = 8$
1	3	$3 \times 2^3 = 12$

$$2 \text{ bcoz } \frac{(\text{bits opcode})_{\text{1addr}} - (\text{bits opcode})_{\text{2addr}}}{2} = 2$$

Q:- Consider a system with 24 bit instruction and 9 bit addresses. If there are 60 2-address instructions, then, maximum how many 1 address instructions can be formulated in the system?

2 address instruction

opcode	add1	add2
6 bit	9 bit	9 bit

1 address instruction

opcode	add
15 bit	9 bit

$$\text{No. of } 2\text{-address instructions} = 60$$

$$\text{No. of 2 address instructions possible} = 2^6 = 64$$

$$\text{No. of unused combinations} = 4$$

First 6 bits in opcode of 1 address instruction can be any of the four combinations.

$$\therefore \text{Max. No. of 1 address instructions that can be formulated} = 4 \times 2^9 = 2^2 \times 2^9 = 2^{11} = \underline{\underline{2048}}$$

Q:- Consider a system with 32 bit instructions and 12 bit addresses. If there are 254 2-address instructions and 8000 1-address instructions then, maximum how many 0-address instructions can be formulated?

2 address ins

opcode	add1	add2
8 bit	12 bit	12 bit

1 address ins

opcode	add
20 bit	12 bit

0 address

opcode
32 bit.

$$(2^8 - 254) \times 2^{12} = \text{Max. no. of 1 address instructions}$$

$$2 \times 2^{12} = 2^{13}$$

$$\text{Max. no. of 0 address instructions} = (2^{13} - 8000) \times 2^{12}$$

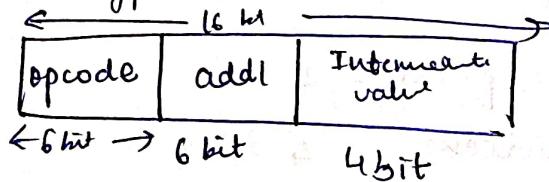
$$= (8 \times 1024 - 8000) \times 2^{12}$$

$$= (8192 - 8000) \times 2^{12} = \underline{\underline{192 \times 2^{12}}}$$

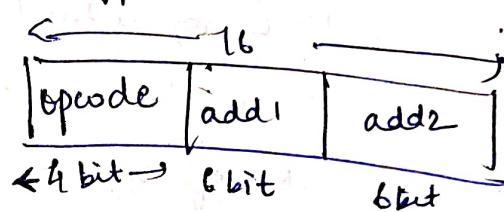
~~Q2) GATE 2020~~

A processor has 64 registers and ^{uses} 16 bit instruction format.
 It has 2 types of instructions i-type & R-type.
 Each I-type instruction contains an opcode, a register name & a 4 bit intermediate value.
 Each R-type instruction contains an opcode and 2 register names.
 If there are 8 distinct I-type opcodes, the max. no. of distinct R-type opcodes is

I-type instruction



R-type instruction



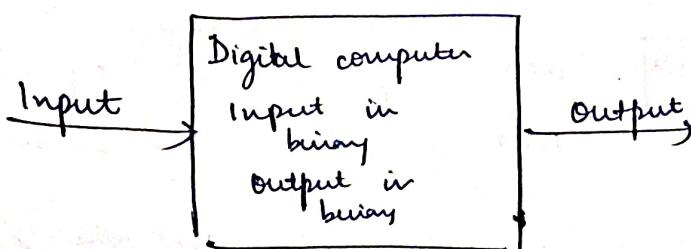
Let max. no. of R-type instructions = x

$$(2^4 - x) \times 2^2 = 8$$

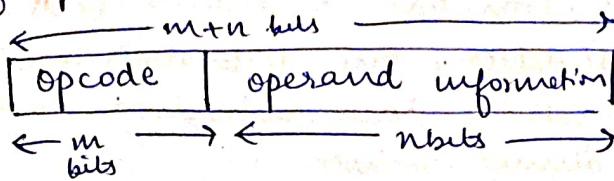
$$2^4 - x = 2 \quad , \quad x = 2^4 - 2 = 16 - 2 = \underline{\underline{14}}$$

INSTRUCTION

- ① The commands to CPU for performing operations
- ② Group of bits which tell computer system to perform an operation



The type of operation to be performed is determined by opcode



2^m types of instructions supported.

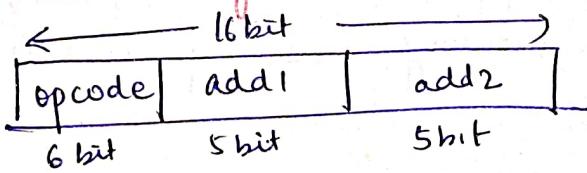
Instruction Set Architecture (ISA)

Collection of all instructions supported by CPU

Types of instruction based on operands.

- 3 Address Instruction (3 addresses used)
 - 2 Address Instructions (2 addresses used)
 - 1 Address Instructions (1 address used)
 - 0 Address Instructions (0 address used) → 2nd operand is accumulator & result is stored in accm
- supported in
accumulator based
architecture stack
based
architecture

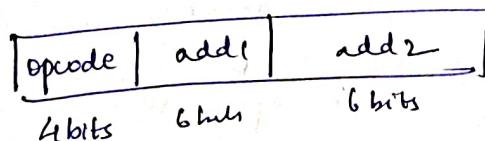
Q:- Consider a digital computer which supports only 2 address instructions each with 16 bits. If address length is 5 bits then maximum & minimum how many instructions the system can support?



$$\text{Max. no. of instructions} = 2^6 = 64$$

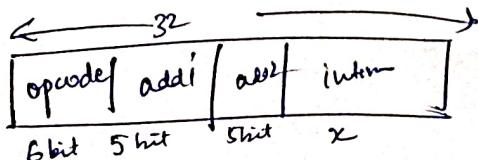
$$\text{Min.} = 1$$

Q:- Consider a digital computer which supports 16 2 address instructions if address length is 6 bits then the bytes of one instruction is —



Size of instruction
= 16 bits

Q:- A processor has 60 distinct instructions and 24 general purpose registers. A 32-bit instruction word has an opcode, 2 register operands and an intermediate operand. The no. of bits available for intermediate operand is —



$$x = 32 - (6 + 5 + 5)$$

$$= 32 - 16 = 16 \text{ bits}$$

Gate 2016
 Consider a processor with 64 registers and an instruction set of size 12. Each instruction has 5 distinct fields — opcode, 2 source register identifiers, one destination register identifier and a 12 bit immediate value. Each instruction must be stored in memory in byte aligned fashion.

If 100 instructions in program, the amount of memory (in bytes) consumed by the program text is _____?

opcode	add1	add2	add3	immediate
4 bit	6 bit	6 bit	6 bit	12 bit

length of instruction = $4 + 6 \times 3 + 12 = 34$

Since memory is byte aligned,

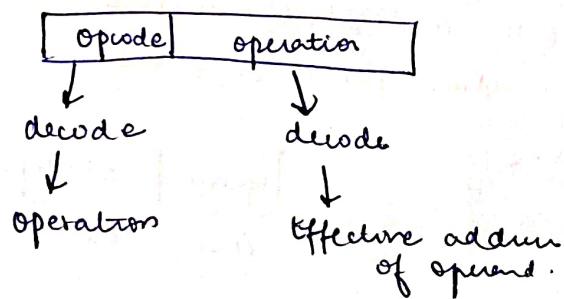
∴ 5 bytes are required to store 1 ins.

Memory req. to store 100 instructions = $100 \times 5 = 500$ bytes

Effective Address: address in memory where operand is stored.

How instruction is executed?

- ① Instruction Fetch
- ② Instruction Decode
- ③ Effective Address Calculation
- ④ Operand Fetch
- ⑤ Execute
- ⑥ Write Back Result



Instruction cycle

All phases are not required for all types of instructions.

115

Fetch
Decode
Execute

Mandatory
in all
instructions

Effective Address Calculation

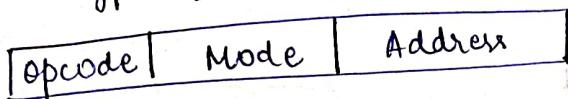
Operand fetch

Write Back Result

Optional
in some
instructions

Addressing modes

Tells the type of address in the instruction.



It specifies how and from where the operands are obtained ~~from~~ for the instruction.

Non
computable
model
i.e. no
computation
required.

1. Implied mode
2. Immediate mode
3. Direct Addressing mode
4. Indirect Addressing mode
5. Register mode.

6. Register Indirect mode
7. Autoincrement & Autodecrement mode.
8. Index Register mode
9. PC Relative mode
10. Base Register mode

1. Implied mode

The opcode definition itself defines the operand.

Example: INCA (Increment accumulator)

0 address instruction

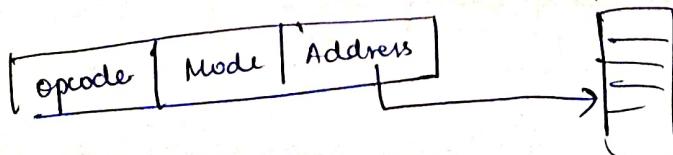
2. Immediate mode

Address field of instruction gives operand value.

→ used to initialize registers with constant values

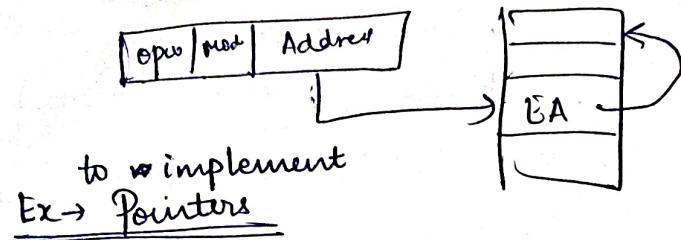
3. Direct mode

Address field specifies the effective address (address of memory).



4. Indirect mode

Address field of instruction specifies the address of effective address.

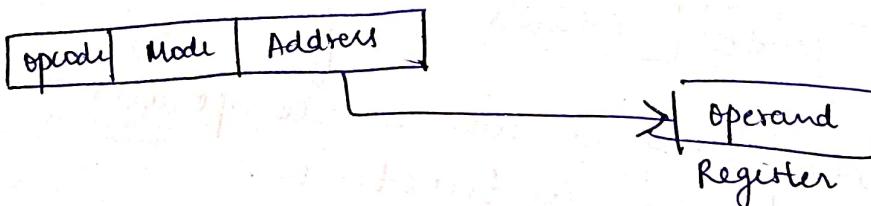


No. of memory accesses = 2.
get operand

to implement
Ex → Pointers

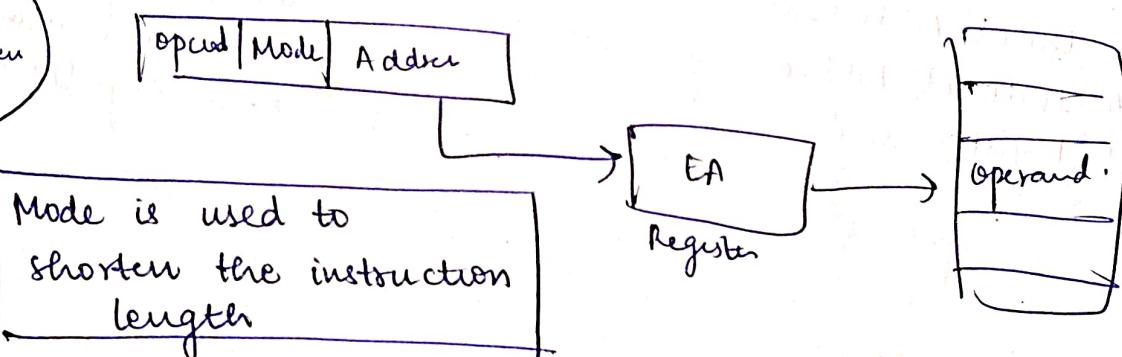
5. Register mode

Address field specifies the register which holds the operand.

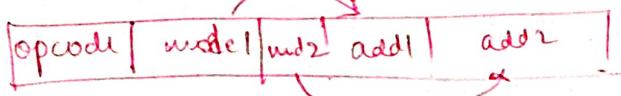


6. Register Indirect mode

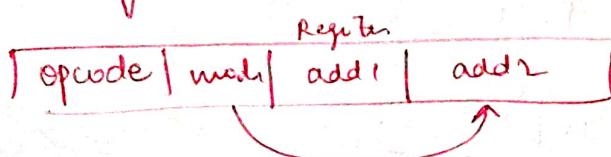
Address field specifies the address of register which holds effective address.



2 address instruction format



Register memory architecture → one address is fixed on register.



⑦ Autoincrement / Autodecrement mode

- ① Variant of register indirect mode.
- ② Content of register is automatically incremented / decremented to access the sequential data.

Auto increment \rightarrow Post increment \rightarrow perform operation then increment value of register.

Autodecrement \rightarrow Pre decrement \rightarrow first decrement then use the value of register.

⑧ Indexed / Index Register mode.

- ① used to access a single element of array.

^{Effective}
② Address of operand = Value in address field + value in index register.

⑨ PC - Related mode / Position independent mode.

- ① used for branch instruction [intra segment branching]

Effective address = Value in address field + Value in program counter

⑩ Base Register mode

- ① used for intersegment branching.

② Base Register value is added to address field of instruction to get effective address.

Effective address

\leftarrow address

~~next~~

instruction dt

Example

index	200	Memory
	opcode mode	
201	Address = 500	
202	Next instruction	
399	450	
400	700	
500	800	
600	900	
702	Target instruction	
800	300	

After instruction fetch, PC=202

202	PC = 200
	RI = 400
index	XR = 100
	AC

Mode	Effective address	Operand
Immediate mode	201	500
Direct mode	500	800
Indirect mode	800	300
Register mode	400	400
Register Indirect mode	700	700
Autodecrement mode	399	450
Indexed mode	600	900
PC-Relative mode	702	400 Target

Q:- An instruction is stored at location 300 and its address field at location 301. The address field has the value 250. A processor register R1 contains number 200. Evaluate effective address if addressing mode is

1) Direct

$$EA = 250$$

2) Relative (PC relative)

$$EA = 302 + 250 = 552$$

3) Register indirect

$$EA = 200$$

Q:- A relative branch mode type instruction is stored in memory at address 300. The branch is made to an address 450.

1. What should be the value of relative adder field of the instruction?
2. Determine the value of PC before instruction fetch, after the fetch and after execution phase.

1. 149

2. Before instruction fetch = 300

After fetch = 301

After execute = 450

Q:- Consider a RISC machine where each instruction is exactly 4 bytes long. Conditional and unconditional branch instruction use PC-relative addressing mode with offset specified in bytes to the target location of the branch instruction. Consider the following instruction sequence

Ins No.

i

Ins
add R2 R3 R4

i+1

sub R5 R6 R7

i+2

cmp R1 R9 R10

i+3

beq R1, offset

Branch ins.

PC value = i+4

$$\text{Target}^4 = \text{PC} + \text{offset}$$

$$\text{offset} = (\text{target} - \text{PC}) \times 4$$

$$= \cancel{10000}$$

$$(i+4 - i) \times 4$$

$$= \underline{-16}$$

Central Processing Unit (CPU)

It is the brain of a computer, containing all the circuitry needed to process input, store data and give output results.

CPU cycle time = Time CPU takes to perform a micro operation

$$\boxed{\text{CPU cycle time} = \frac{1}{\text{clock rate}} \rightarrow \text{given in GHz.}}$$

CPI = Cycles Per Instruction

$$\text{Time taken to execute instruction} = \text{CPI} * \text{CPU cycle time} = \frac{\text{CPI}}{\text{clock rate}} \rightarrow \text{avg}$$

Q: Consider computing the overall CPI for a machine A for which the following performing measures were recorded when executing a set of benchmark programs. Assume that the clock rate = 200MHz.

Instruction Category	% of occurrence	No. of cycles per ins.
ALU	48	1
Load & Store	10	3
Branch	39	4
Other	3	5

$$\begin{aligned}
 \text{CPI}_{\text{avg}} &= 1 \times 0.48 + 3 \times 0.1 + 4 \times 0.39 + 5 \times 0.03 \\
 &= 0.48 + 0.3 + 1.56 + 0.15 \\
 &= \underline{\underline{2.49}}
 \end{aligned}$$

0.48
0.30
1.56
0.15
2.49

(21)

MIPS (Million Instructions per second)
Used to give CPU's performance

Gate 2014

Consider two processors P₁ and P₂ executing the same instruction set. Assume that under identical conditions, for the same input, a program running on P₂ takes 25% less time but incurs 20% more CPI as compared to the program running on P₁. If the clock frequency of P₁ is 1GHz then the clock frequency of P₂ (in GHz) is _____?

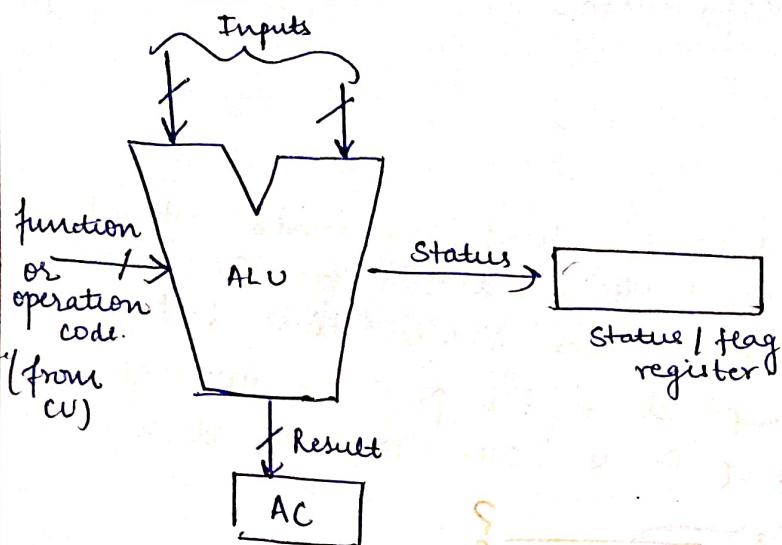
No. of instructions in P₁ and P₂ are same.

	P ₁	P ₂
Execution time	t ₁	0.75 t ₁
CPI	c ₁	1.2 c ₁
Freq.	1GHz	?

$$\text{No. of instructions executed, } n = \frac{\text{Execution time} \times \text{clock rate}}{\text{CPI}}$$

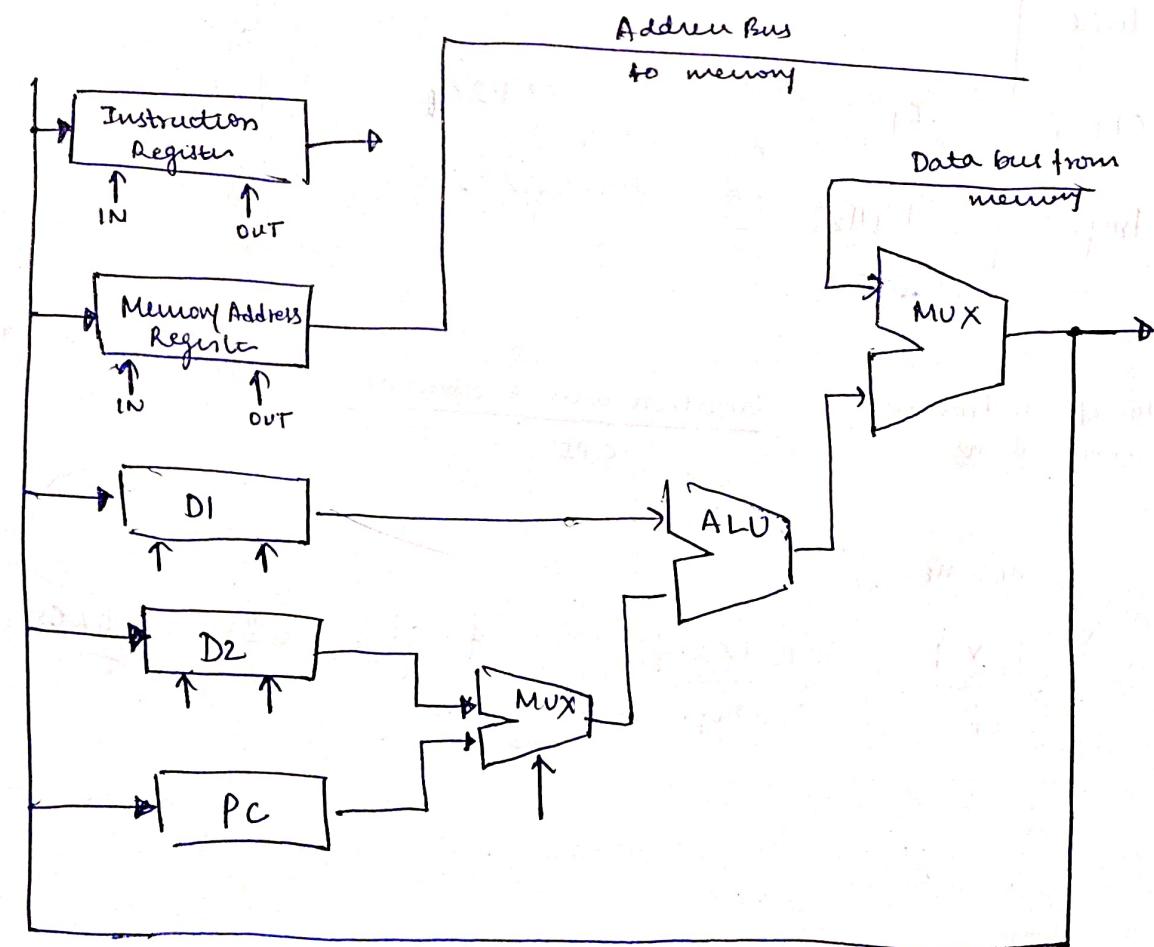
$$n_1 = n_2 \\ \Rightarrow \frac{t_1 \times 1}{c_1} = \frac{0.75 t_1 \times f}{1.2 c_1} \Rightarrow f = \frac{1.2}{0.75} = \frac{1.2 \times 4}{3} = \underline{\underline{1.6 \text{ GHz}}}$$

ALU



Datapath (CPU Design)

- ① Collection of functional units such as ALU and multiplexers.
- ② Perform data processing operations.



Instruction Fetch -

$$IR \leftarrow M[PC]$$

$$1. AR \leftarrow PC$$

$$2. IR \leftarrow M[AR] \quad \{ \text{mutually exclusive}$$

$$\textcircled{B} \quad PC \leftarrow PC + 1 \quad \} \text{ in same step.}$$

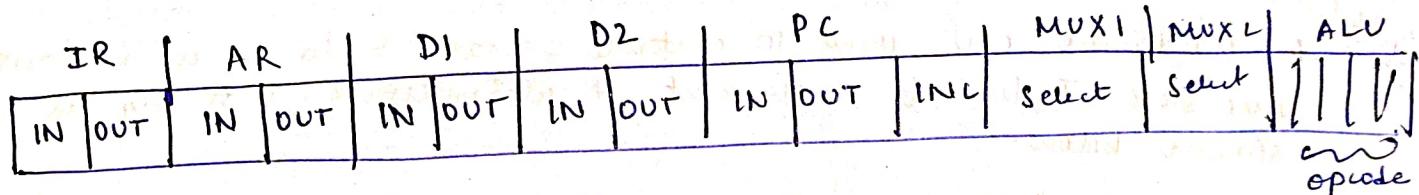
(23)

Control unit

generates control signals and sends to various components of computer. All components operate according to these control signals.

control word - ^{variable} name of control signal

control word - collection of signals generated by control unit at once.



Control word

$$1. AR \leftarrow PC$$

Pout

① Control unit generates control word.

② control word signals are sent to different components.

③ Components perform microoperations based on control signal

Hardwired Control Unit

Control logic is implemented with gates, flip flops, decoders & other digital circuits.

Advantage

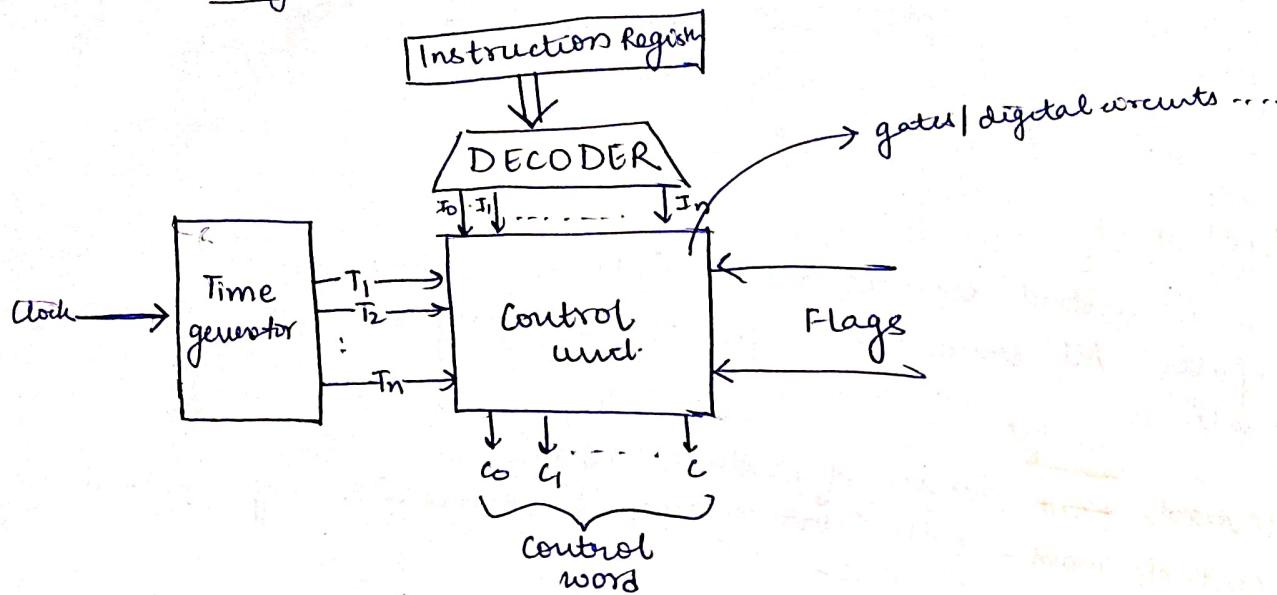
can be optimized to produce a faster mode of operation.

Disadvantage

Rearranging the wires among various components is difficult.

Updation in control logic is difficult.

Design of Hardwired Control unit



Gate 2005

A hardwired CPU uses 10 control signals S_1 to S_{10} in various time steps T_1 to T_5 to implement 4 instructions I_1 to I_4 as shown below—

	T_1	T_2	T_3	T_4	T_5
I_1	$S_1 S_3 S_5$	$S_2 S_4 S_6$	$S_1 S_7$	S_{10}	$S_3 S_8$
I_2	$S_1 S_3 S_5$	$S_8 S_9 S_{10}$	$S_5 S_6 S_7$	S_6	S_{10}
I_3	$S_1 S_3 S_5$	$S_7 S_8 S_{10}$	$S_2 S_6 S_9$	S_{10}	$S_1 S_3$
I_4	$S_1 S_3 S_5$	$S_2 S_6 S_7$	$S_5 S_{10}$	$S_6 S_9$	S_{10}

Which of the following pairs of expressions represent the circuit for generating control signals S_5 and S_{10} respectively?

- A. $S_5 = T_1 + I_2 \cdot T_3$ and
 $S_{10} = (I_1 + I_3) \cdot T_4 + (I_2 + I_4) \cdot T_5$
- B. $S_5 = T_1 + (I_2 + I_4) \cdot T_3$ and
 $S_{10} = (I_1 + I_3) \cdot T_4 + (I_2 + I_4) \cdot T_5$
- C. $S_5 = T_1 + (I_2 + I_4) \cdot T_3$ and
 $S_{10} = (I_2 + I_3 + I_4) \cdot T_2 + (I_1 + I_3) \cdot T_4 + (I_2 + I_4) \cdot T_5$
- D. $S_5 = T_1 + (I_2 + I_4) \cdot T_3$ and
 $S_{10} = (I_2 + I_3) \cdot T_2 + I_4 \cdot T_3 + (I_1 + I_3) \cdot T_4 + (I_2 + I_4) \cdot T_5$

$$S_5 = T_1 + T_3 \cdot I_2 + T_3 \cdot I_4 = T_1 + (I_2 + I_4) T_3$$

$$S_{10} = (I_2 + I_3) \cdot T_2 + I_4 \cdot T_3 + (I_1 + I_3) T_4 + (I_2 + I_4) T_5$$

Microprogrammed Control Unit

- ① Control logic is implemented with micro-programme.
- ② All possible control words are stored in memory & based on requirement word is sent to different components.

Advantage

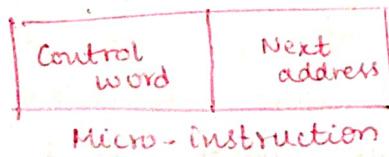
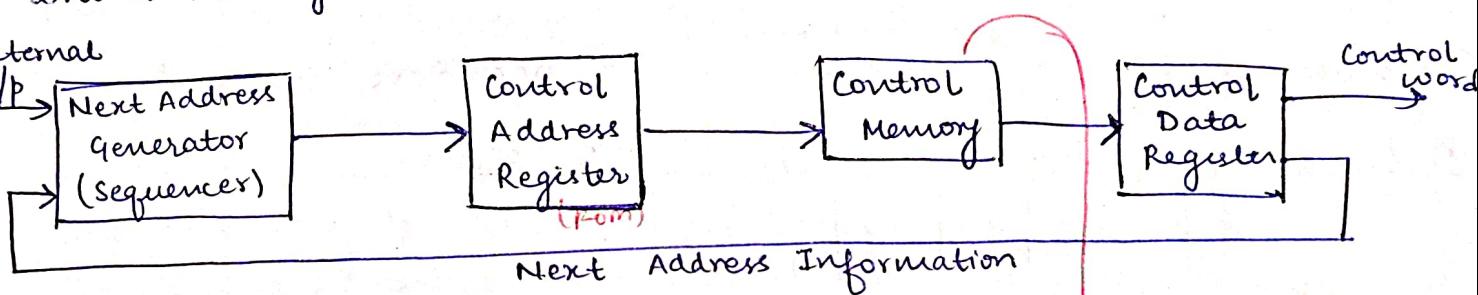
Updating the control logic is easy

Disadvantage

slower mechanism (because memory is involved)

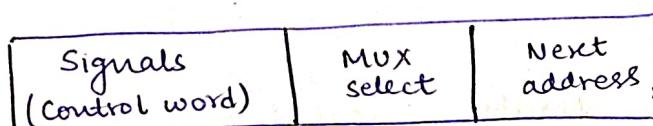
control memory — control words are stored.

External i/p



Each row stores the control word + next word information.

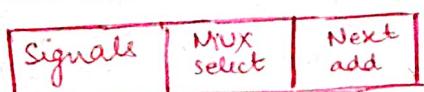
Standard micro instruction format



Control word sequencing.

Q: Consider a CPU

1. 120 control signal bits
2. Microprogrammed CU micro-instruction



3. 16 inputs MUX
4. 16 distinct MUX supported by CPU

64 microoperations per instruction

No. of bits in signal = 120

No. of bits in MUX select = 4

$$\begin{aligned} \text{Total no. of micro-instructions in control mem} &= 16 \times 64 \\ &= 2^4 \times 2^6 = 2^{10} \end{aligned}$$

∴ Next add → 10 bits

$$\therefore \text{Size of microin} = 120 + 4 + 10 = 134$$

$$\begin{aligned} \therefore \text{Size of memory} &= 134 \times 2^{10} \text{ bytes} \\ &= 134 \text{ kB} \end{aligned}$$

Type of microprogrammed

↓
Horizontal

- ① One bit for each control signal

↑
Vertical

signals are divided into multiple groups such that at a time only 1 signal can be active.

Each group's info is stored or encoded from

- ② Larger control word.

- ③ Smaller control word.

- ④ No decoder

- ⑤ Decoders are used.

- ⑥ faster

- ⑦ slower.

Gate 2002

Horizontal microprogramming:

- (A) Does not require use of signal decoders.
- (B) Results in larger sized microinstructions than vertical.
- (C) uses one bit for each control signal
- (D) all of the above

Gate 1999

Decreasing order of operating speeds

Vertical



Horizontal



Hardwired

A control unit generates 120 control signals, which are divided into 6 groups of mutually exclusive signals as below:

Group 1: 30

5x32 Decoder

No. of bits req. in

Group 2: 13

4x16 Decoder

vertical = 5+4+4 + 2+5+6 = 26 bits

Group 3: 12

4x16 Decoder

No. of bits req. in

Group 4: 3

2x4 Decoder

horizontal = 120 bits

Group 5: 27

5x32 Decoder

∴ No. of bits saved = 120 - 26 = 94 bits

Group 6: 35

6x64 Decoder

How many bits are saved by using vertical microprogrammed control and as compared to horizontal?

Types of CPU Design

27

RISC (Reduced Instruction Set Computer)

Less number of instructions.

Fixed length instructions
(\therefore variable ^{width} opcode)

Simple instructions

Limited addressing modes.

Hardwired CU

Only Register-Register arithmetic operations

More no. of registers

<Preferable for pipelined CPU>

CISC (Complex Instruction Set Computer)

More number of instructions.

Variable length instruction
(\therefore fixed length opcode)

Complex instructions

More & complex addressing modes.

Microprogrammed CU

Register-memory & memory to operations

Less no. of registers

Date 2018

Consider the following processor design characteristics -

I. Register to register arithmetic operations only.

II. Fixed length instruction format.

III. Hardwired control unit.

Which of the following characteristics above are used in design of RISC processor.

A) I & II only

(C) I and III only

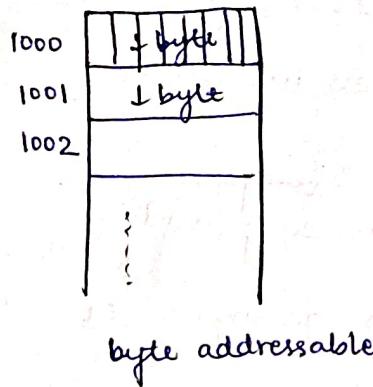
(B) II and III only

(D) I, II and III

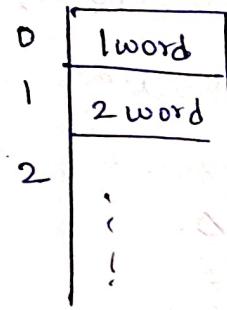
Byte v/s word addressable memory

Addressable meaning -

which unit का एक unique address दिया जाया है,
ज्यादा वो byte है? या वो word है? ~~word~~
जानने के लिए क्या रीहर हमारे साथ



byte addressable



word addressable

where word size = 32 bit

Byte addressable

Address of first word = 0

Address of second word = 4

Address of third word = 8

:

Word addressable

Address of first word = 0

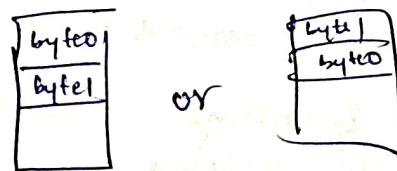
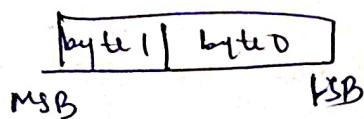
Address of 2nd word = 1

Address of 3rd word = 2

Byte Ordering

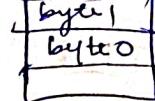
How multi byte data is stored in memory?

from LSB to MSB or from MSB to LSB.



This ordering is decided by byte order.

Big Endian MSB to bytes are stored first



Little Endian LSB bytes are stored first



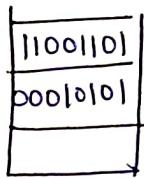
Q:- 16 bit word -

Assume

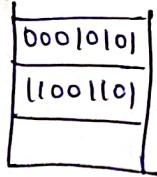
1 word = 2 bytes
memory is
Byte addressable

11 00 11 01 000 10101

Big Endian



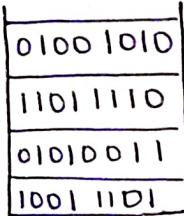
Little Endian



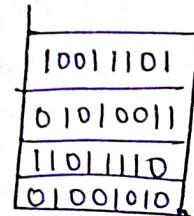
Q:- Assume 1 word = 4 bytes

word = 4ACD539D

Big Endian



little Endian



ARCHITECTURE L&H !!



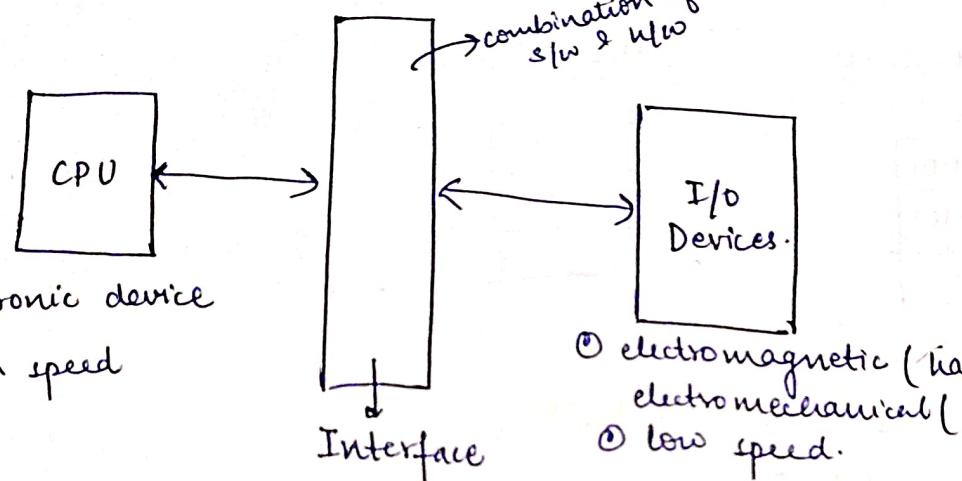
COMPUTER ORGANIZATION

IO Organization

Peripheral device: Devices connected to processor externally.



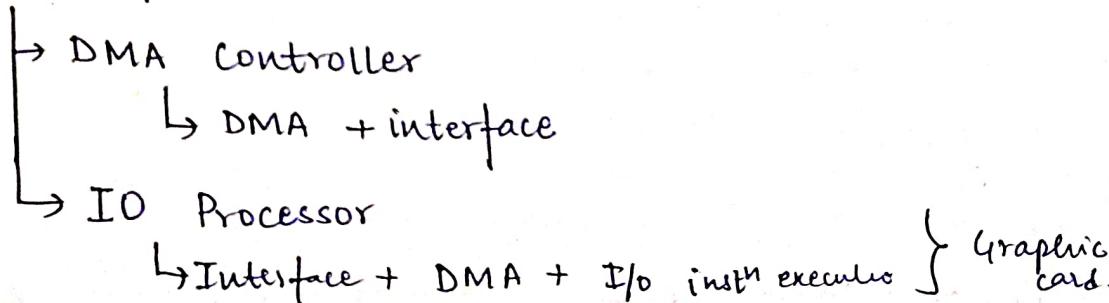
- ① electronic device
- ② high speed



- ① electromagnetic (harddisk)
- electromechanical (printer) etc.
- ② low speed.

- ③ converts signal
- ④ speed management synchronization
- ⑤ conversion of data codes & data formats.

IO Interface



IO v/s memory bus

(3)

- ① 3 ways to connect memory and I/O with CPU through buses.

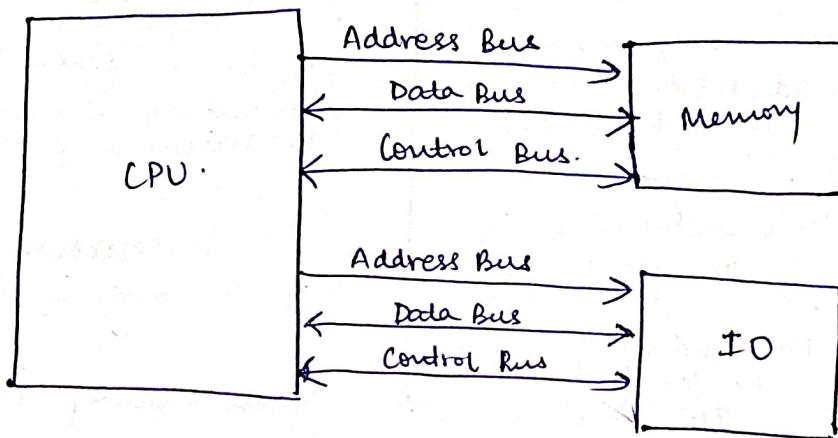
1) Separate Buses for Both

Disadvantage 😞

- ↳ costly
(bcz more buses are required)

Advantage 😊

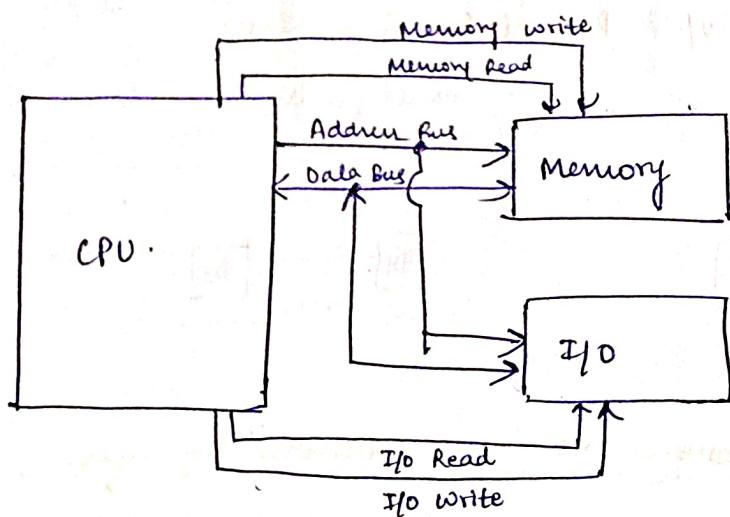
- ↳ easy to design



2) Common Data, Address Bus

Advantage 😊

- ↳ no memory wastage



CPU distinguishes b/w memory & I/O using control signal.

I/O Devices have their own addresses

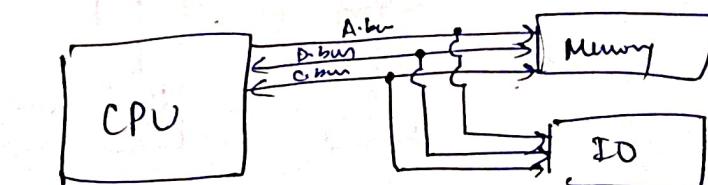
- ∴ I/O mapped I/O
- ∴ Isolated I/O

3) Common Address, Data & Control Bus

Disadvantage 😞

- ↳ memory wastage

Advantage 😊



X
I/O addresses do not have their own address.
↳ known as Memory mapped I/O

↳ memory addresses are assigned to I/O device.
whenever CPU addresses I/O,
it goes directly to I/O

correctness

Some locations in memory are used to store address of I/O.

So, that CPU can access I/O devices from those assigned addresses.

Memory mapped I/O

Memory wastage

No separate address space for IO

All memory access instructions are used by IO also

More instructions for IO access

more addressing modes for IO access

IO Mapped IO

No memory wastage

Separate address spaces for IO

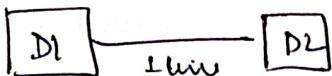
IO access instructions & memory access instructions are different.

less instructions

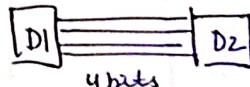
less addressing modes

Serial v/s Parallel Transfer:

1 bit at a time
slow

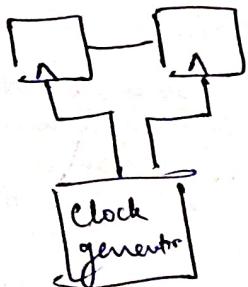


multiple bits at a time
fast
costly

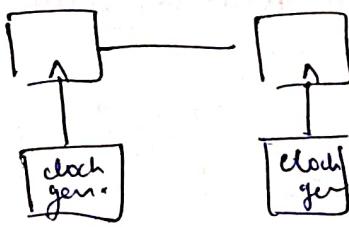


Synchronous v/s asynchronous transfer:

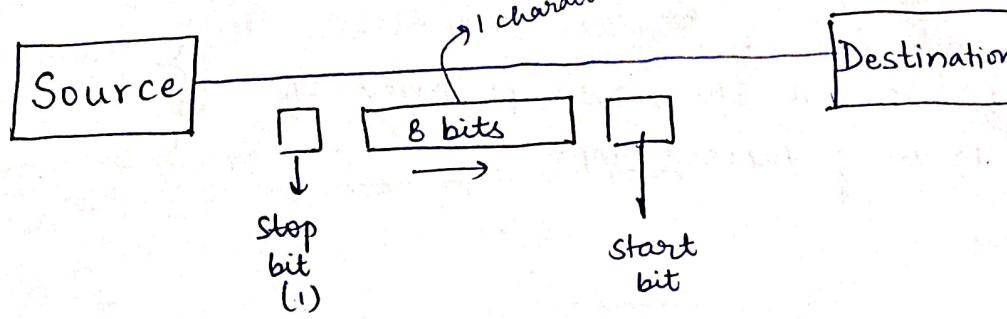
all devices operate on single clock



separate each device has its own clock.



Serial Asynchronous Transfer



- ① No transfer means 1 state.
- ② 0 (start bit) detected means transmission started
- ③ character bits follow start bit.
- ④ After last bit of character '1' (stop bit) will be sent.

$$\text{Efficiency of transmission line} = \frac{8}{10} = 0.8 = \underline{\underline{80\%}}$$

Q. How many 8 bit characters can be transmitted per second over 9600 bits/second serial communication link using a parity synchronous mode of transmission with 1 start bit, 8 data bits, 2 stop bits and 1 parity bit.

To send one 8 bit character, $1 + 8 + 2 + 1 = 12 \text{ bits}$ are required.

9600 bits can be transferred in 1 sec.

$$\therefore \text{No of characters} = \frac{9600}{12} = \underline{\underline{800 \text{ characters/sec.}}}$$

Q. An asynchronous serial communication is employing 8 character bits, 1 parity bit, 2 start bits & 1 stop bit. To maintain rate of 400 char/sec the minimum transfer rate required is _____?

$$1 \text{ char} \rightarrow 1 + 2 + 1 + 8 = 12 \text{ bits}$$

400 char/sec

$$400 \times 38.33 = 8k \text{ char}$$

$$= 400 \times 12 \text{ bits/sec}$$

$$= \underline{\underline{4800 \text{ bits/sec}}}$$

Modes of data transfer

- Programmed I/O or program controlled I/O } Data transfer b/w CPU & I/O
- Interrupt initiated I/O or Interrupt Driven I/O } Data transfer b/w I/O & memory.
- Direct Memory Access (DMA)

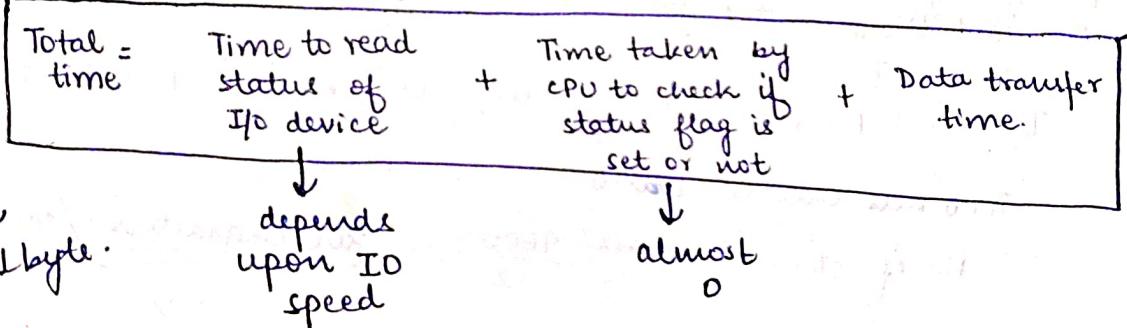
Programmed I/O

- ① No provision through which I/O can inform to CPU about data transfer.
- ② I/O sets its own status bit and waits.
- ③ CPU runs a program periodically and checks the status of each device one by one.
- ④ If any device has its status ~~is~~ set, then, CPU performs data transfer.



→ time wastage of CPU

Time required in Programmed I/O



By default,
status flag = 1 byte.

Ques

Consider a device which operates with 20 MBPS operating speed. The device is operating on program controlled mode of I/O & it has to transfer 20 B of data from it. Data is transferred byte wise.

Size of status register = 2 bytes

Total time needed to perform data transfer in microseconds?

$$\begin{aligned} \text{Total bytes} &= 2 + 20 \\ &= 22 \text{ bytes} \end{aligned}$$

$$\begin{aligned} \text{Total time} &= \frac{22}{20} \mu\text{s} \\ &= 1.1 \mu\text{s} \end{aligned}$$

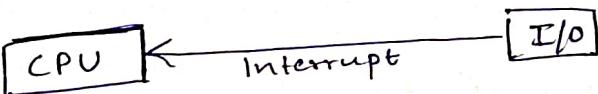
X To transfer 1 byte, 3 bytes are needed [bcoz byte wise transfer]
status = 2 bytes

To transfer 20 B, 60 bytes are needed

$$\therefore \text{Time req.} = \frac{60}{20 \times 10^6} \text{ s} = \underline{\underline{3 \mu\text{s}}}$$

Interrupt IO

- ① IO device has a provision (Interrupt signal) to inform to CPU about communication.
- ② When CPU receives interrupt:
 - completes execution of current instruction.
 - saves the status (PC, PSW etc.) of current process onto the stack.
 - Branches to service the interrupt.
 - Resumes the previous process by taking out the value from stack.



ISR = Interrupt Service Routine [CPU executes ISR when it receives interrupt]

Vector Address → address or location of ISR.
(reference)

Interrupt

Vectored interrupt

Device sends interrupt and vector.

[vector address is sent only after interrupt is acknowledged by CPU]

Data bus is used to send vector.

Interrupt → control signal.

Non-vectored interrupt (scalar)

Device sends only interrupt

→ CPU executes default service routine
CPU obtains location of actual ISR.

Maskable interrupt

CPU can accept or reject the interrupt.

Non Maskable interrupt

CPU always accepts the interrupt.

Interrupt

External interrupt

I/O Device generates the interrupt.

Internal interrupt

Interrupt due to unexpected error during instruction execution.

(Ex → page fault)
System call

Time Required in Interrupt IO

$$= \text{Interrupt overhead time} + \text{Service time.}$$

Simultaneous Interrupt

- ① 2 device generate interrupt simultaneously.
- ② The highest priority device's interrupt should be serviced first

Priority Based Interrupt Handling

Software solution



run a program
resolve the
priority

Hardware solution

Serial

Daisy chaining

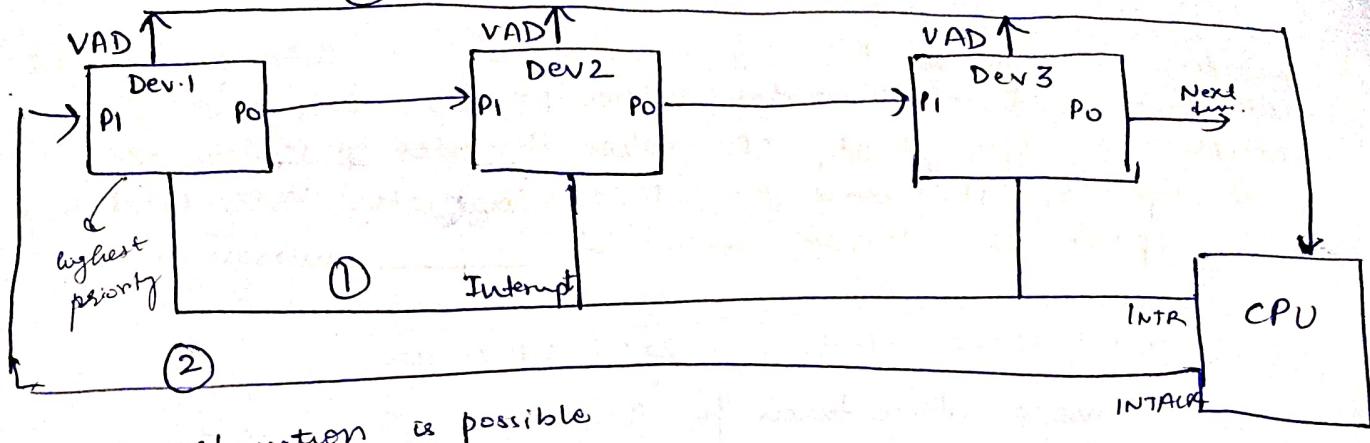
Parallel



VAD = Vector Address

Daisy chaining

135



starvation is possible

Gate 2005

A device with data transfer rate $10 \times 10^3 \text{ kB/sec}$ is connected to a CPU. Data is transferred byte wise. Let the interrupt overhead be 4 microsecond. The byte transfer time b/w 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?

$$\text{Performance gain} = \frac{\text{performance of interrupt}}{\text{performance of programmed I/O}} = \frac{\text{prog I/O time}}{\text{Interrupt I/O time}}$$

or
speed up.

For programmed I/O -

Status register = 1 byte (default)

$$\text{Time to read status register} = \frac{1}{10 \times 10^3} = 10^{-4} \text{ sec} = 100 \mu\text{s}$$

Data transfer rate = 0

$$\therefore \text{Total programmed I/O time} = 100 \mu\text{s}$$

For interrupt I/O -

$$\text{Total time} = \text{overhead} = 4 \mu\text{s}$$

$$\therefore \text{Performance gain} = \frac{100}{4} = \underline{\underline{25}}$$

Ques

Consider a CPU which takes 2 microseconds as interrupt overhead time when a device generates interrupt for CPU, and CPU accepts it. After that, CPU takes 9 cycles to service the interrupt. If CPU runs on 2MHz clock rate, then total time CPU spends for interrupt service is _____ microseconds.

$$\text{clock rate} = 2\text{MHz} = 2 \times 10^6 \text{ cycles/sec.}$$

$$\text{Service Time} = \text{time taken for 9 cycles} = \frac{9}{2 \times 10^6} = 4.5 \mu\text{s}$$

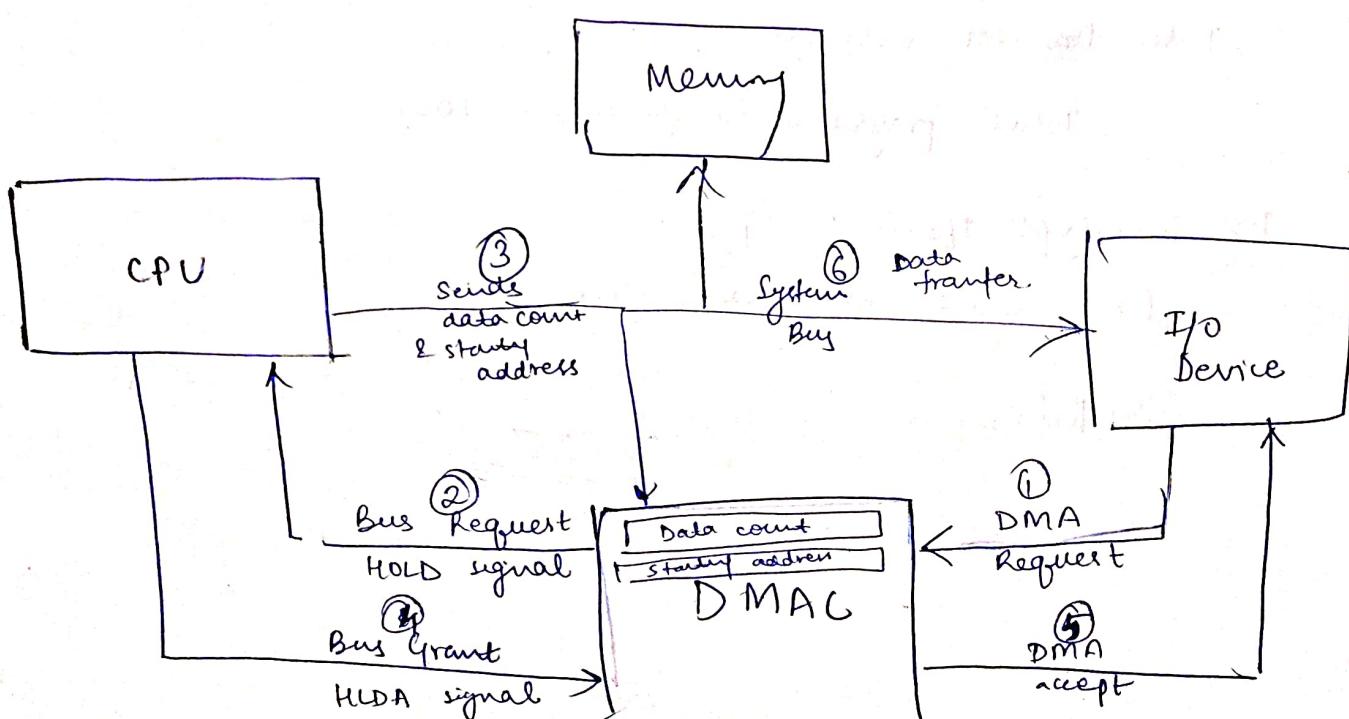
$$\text{Interrupt overhead} = 2 \mu\text{s}$$

$$\therefore \text{Total time} = 2 + 4.5 = \underline{\underline{6.5 \mu\text{s}}}$$

DIRECT MEMORY ACCESS (DMA)

Enables data transfer operation b/w I/O device and memory without CPU intervention.

Need of hardware — DMA Controller



Starting address = memory address starting from where the data transfer should be started.

Data Count = No. of bytes or words to be transferred.
depends upon memory byte addressable / word addressable.

DMA During DMA transfer, CPU can perform only those operations which do not require system bus
→ means CPU will be mostly blocked.

time is saved bcoz there is no need to store present state of CPU in stack it is a time consuming process means a lot!!

Date 20/6

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 main memory. Memory is byte addressable. The minimum number of times the DMA controller needs to get the control of system bus from the processor to transfer the file from the disk to main memory is

Max. no. of bytes that can be transferred if DMA controller has got the control of system bus

$$= \frac{2^{16}}{2^{16}-1} \text{ bytes}$$

$\approx 2^{16}-1$ bytes
 -1 is necessary

Max. no. that can be represented by 16 bits
 $= 2^{16}-1$

$$\therefore \text{No. of times control of system bus needed} = \frac{29154}{2^{16}-1} = \frac{29154}{65535} = 455.53 = 456$$

Modes of DMA Transfer

- ① Entire data is not transferred at once from I/O to memory.
- ② For sometime, DMAC has control of bus then CPU takes back the control again and gives it back to DMAC.

- Burst Mode
- cycle Stealing mode
- Interleaving mode.

to enhance user experience
Screen freeze
no jags if
CPU is blocked

BURST MODE

A burst of data (1KB - 2KB generally) is transferred in one time before CPU takes back the control of the buses.

CYCLE STEALING MODE

- slow I/O device takes some time to prepare a byte or word.
- During this time CPU keeps the control of the buses.
- whenever the word/byte is ready, then, CPU gives the control of the buses to DMAC for 1 memory cycle.
- During this cycle, the word/byte is transferred to memory.

INTERLEAVING MODE

CPU will never be in waiting mode.

Whenever CPU performs internal operations and does not require the buses then DMAC will be given control of the buses.

% of time CPU is blocked

Time required to prepare the data in I/O = t_x

Time required to transfer the data to memory = t_y

$$\boxed{\% \text{ of time CPU is blocked due to DMA} = \frac{t_y}{t_x + t_y}}$$

Burst mode

$$\boxed{\% \text{ of time CPU is blocked} = \frac{t_y}{t_x}}$$

Cycle stealing mode.

$$\boxed{\% \text{ of time CPU is blocked} = 0}$$

Interleaving mode

t_x is dependent on the speed of I/O device. (given or given)
 t_y can be given in question or depends upon cycle time.

Ques

Consider a device operating on 2Mbps speed & transferring data to memory using cycle stealing mode of DMA. If it takes $2\mu s$ to transfer 16B of data to memory. What is the ready prepared time?

$\% \text{ of time CPU is blocked due to DMA} = ?$

Suppose x B of data is to be transferred.

$$t_x = \frac{x}{16} \times 2\mu s$$

$$t_y = \frac{x}{2} \mu s$$

$$\therefore \% = \frac{t_y}{t_x} = \frac{x \times 16}{2 \times x} = 4\%$$

16 bytes of data is to be transferred.

t_y = time taken to transfer = $2\mu s$

$$t_x = \frac{16}{2} \mu s = 8\mu s$$

$$\therefore \% \text{ of time CPU is blocked} = \frac{t_y}{t_x} = \frac{2}{8} = \frac{1}{4} = \underline{\underline{25\%}}$$

Ques A device is constantly transferring the data with 40 kBps speed to memory using DMA. Assume that the transmission time to memory from IO once the data (4 bytes) is ready is 10 μs.

1. % of time CPU is blocked in burst mode.

2. % of time CPU is blocked in cycle stealing mode

Data size = 4 bytes

$$t_y = \frac{4 \text{ bytes}}{10 \mu\text{s}} \quad t_x = \frac{4}{40 \times 10^3 \text{ s}} = \frac{1}{10^4} \text{ s} = 100 \mu\text{s}$$

1) % of time CPU is blocked in burst mode = $\frac{t_y}{t_x + t_y} = \frac{10}{10 + 100} = \frac{10}{110} = \frac{1}{11} \approx 9.09\%$

2) % of time CPU is blocked in cycle stealing

$$\text{mode} = \frac{t_y}{t_x} = \frac{10}{100} = \underline{10\%}$$

GATE 2021

consider a ~~DMA support~~ computer system with DMA support.

The DMA module is transferring one 8 bit character in one CPU cycle from a device to memory through cycle stealing at regular intervals.

Consider a 2 MHz processor.

If 0.5 cycles are used for DMA, the data transfer rate of the device is _____ bits per second.

8 bit character in 1 CPU cycle

2 MHz processor $\Rightarrow 2 \times 10^6$ cycles per second

$\frac{0.5}{100} \times 2 \times 10^6$ cycles are used for DMA $\approx 10^4$.

$\Rightarrow 10^4$ cycles are used for DMA per second

1 cycle \rightarrow 8 bit data

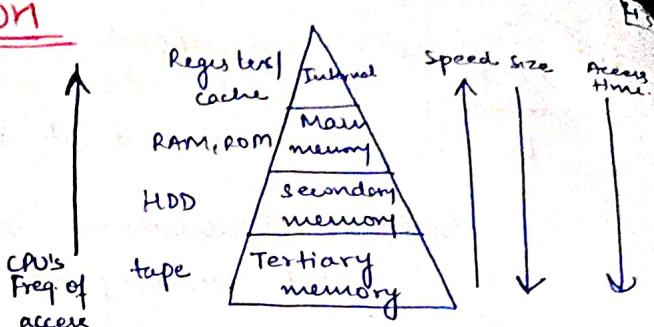
$\therefore 10^4$ cycles $\rightarrow 8 \times 10^4$ bit data

\therefore Data transfer rate = 80000 bits/sec

Memory Organization

Memory hierarchy is used

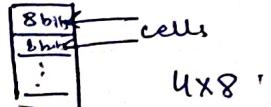
- To maximize access speed.
- To minimize per bit storage cost.



Memory Representation

$$= \text{No. of cells} \times \text{Each cell capacity}$$

$$= \text{No. of memory locations} \times \text{Bits per location}$$



Default unit of storage = Bits (if not given or question)

Ques
Consider a CPU with 13 bits address bus and 16 bits in data bus. Maximum memory size which can be supported is 16 kbytes

$$\begin{aligned} & 2^{13} \times 2^{16} = 2^{29} \text{ bits} \\ & = 2^{14} \text{ bytes} = \underline{\underline{16 \text{ KB}}} \end{aligned}$$

Ques
A 32 bit wide memory has 24 bits address to be accessed. Maximum memory capacity =

$$\text{Size of each memory location} = 32 \text{ bit}$$

$$\text{No. of memory locations} = 2^{24}$$

$$\therefore \text{Maximum memory capacity} = 2^{24} \times 32 = \underline{\underline{2^{29} \text{ bits}}}$$

Memory cycle time = time req. to perform read / write operation in 1 location.

Ques
The memory cycle time of a memory is 200 nsec. The maximum rate with which the memory can accessed bytes/sec.

Note: byte addressable

200 nsec \rightarrow 1 byte data is accessed

$$\therefore \text{Rate} = \frac{1}{200 \text{ nsec}} \times 10^9 \text{ bytes/sec} \quad \frac{10^9}{200} \text{ bytes/sec} = \underline{\underline{5 \times 10^6 \text{ bytes/sec}}}$$

gate 2016

A processor can support a maximum memory of 4GB where memory is word addressable (2 bytes word). The size of address bus of the processor is at least 31 bits.

Size of 1 memory location = 2 bytes

Total memory size = 4GB = 4×2^{30} bytes

$$\therefore \text{No. of locations} = \frac{4 \times 2^{30}}{2} = 2 \times 2^{29} = 2^{31}$$

\therefore No. of bits in address bus = $\log(2^{31})$ 31 bits

Main memory

Used to store currently running programs and their data.

RAM Random Access Memory (volatile)

ROM Read Only Memory (non-volatile)

P.O.S.T (Power on self test) hardware connected to CPU

Bootstrapping (loading of OS)

Bootstrap program = Program responsible for booting.

ROM is accessed only when power is turned off.

Once OS is loaded, ROM is not accessed again until next start.

Types of RAM

145

Static

- ① Mode = flip flops
- ② No refresh req.
- ③ Faster
- ④ Expensive
- ⑤ used for cache
- ⑥ High Power Consumption

Dynamic (DRAM)

- ① Capacitors.
- ② Storage in form of electric charge
- ③ Periodic refresh required.
- ④ Slower
- ⑤ less expensive.
- ⑥ used for main memory.
- ⑦ low operational power consumption

Ques The amount of ROM needed to store the table of multiplication of two 4 bit unsigned integers is?

- a) 64 bits
- b) 128 bits
- c) 1K bits
- d) 2K bits.

4 bit unsigned integers
No. of bits in product = $4+4=8$ bits.

No. of combinations of product = $2^4 \times 2^4 = 2^8$

$$\therefore \text{Size of ROM} = 2^8 \times 8 = 2^8 \text{ bits}$$

= 2K bits

Ques A ROM is used to store the table of multiplication of two 8 bit unsigned integers. The size of ROM required is

- a) 256×16
- b) $64K \times 8$
- c) $4K \times 16$
- d) $64K \times 16$

8 bit

Size of each cell = No. of bits in product = $8+8=16$ bits.

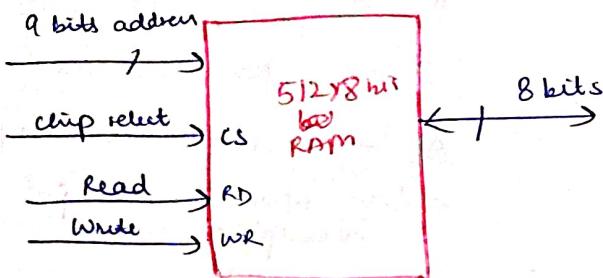
No. of combinations of product = $2^8 \times 2^8 = 2^{16}$

$$\therefore \text{Size of ROM} = 2^{16} \times 16$$

= 64K \times 16

Decoder required to access memory
of size $128K \times 16$ bits
 $2^7 \times 2^10 = 2^17$ = 17×2^{17} decoder

RAM chip

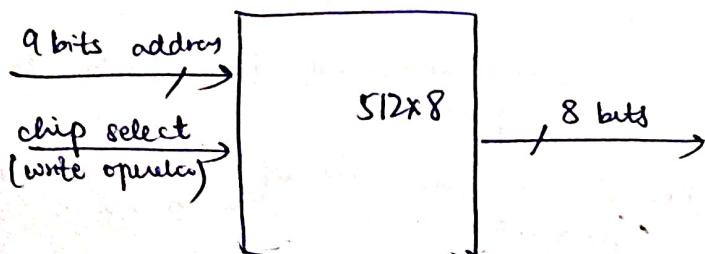


$$\begin{aligned} \text{No. of pins required in RAM chip} &= 9 + 8 + 1 + 1 + 1 \\ &= \underline{\underline{20}} \end{aligned}$$

CS	RD	WR	Operation
0	X	X	No operation
1	0	0	No operation
1	0	1	Write
1	1	X	Read

even if both RD & WR are 1, a read operation takes place.

ROM chip



$$\begin{aligned} \text{No. of pins required} &= 8 + 9 + 1 \\ &= \underline{\underline{18}} \end{aligned}$$

CS	Operation
0	No operation
1	Write read

Ques A ROM chip of 2048×16 bits has four select inputs and operates from a 5 volt power supply. How many pins are needed for the IC package?

No. of pins for address line = 11

No. of pins for data line = 16

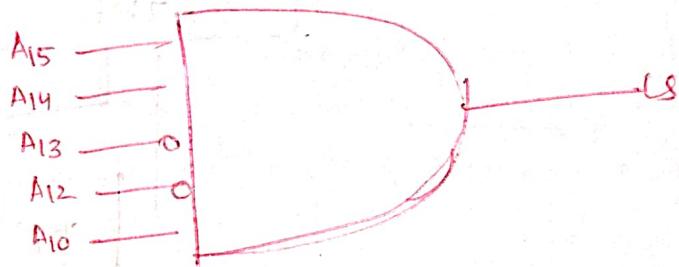
No. of pins for chip select = 4

2 connections are required →
→ power - 1 bit
→ ground - 1 bit

$$\therefore \text{No. of pins} = 11 + 16 + 4 + 2 = \underline{\underline{33}}$$

Gate 2019

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



(A) C800 to CFFF

(C) C800 to C8FF

- (B) CADD to CAFF
(D) DA60 to DFFF

11001000000000000000 → Beginning

1100111111111111 → Ending

$\therefore \text{C800 to CFFF}$

Multiple chips in single system

Total capacity of system = no of chips * capacity of 1 chip.

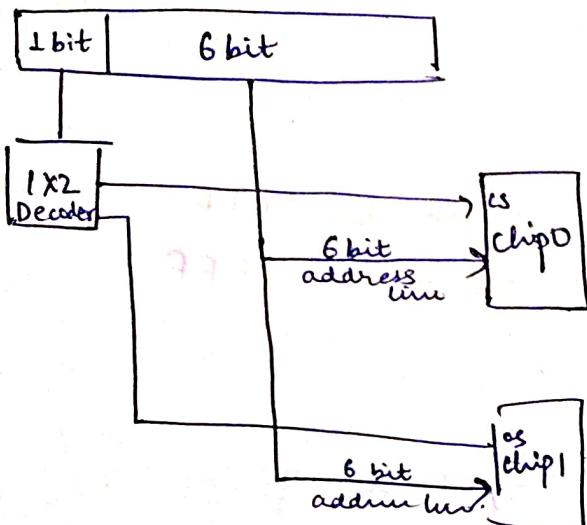
Vertical Arrangement

used when number of addresses required in memory system is more than individual chips.

Ex → chip0 — RAM 64×8 bit
chip1 — RAM 64×8 bit

$$\text{Total capacity} = 128 \times 8 \text{ bit}$$

7 bit address generated by CPU.



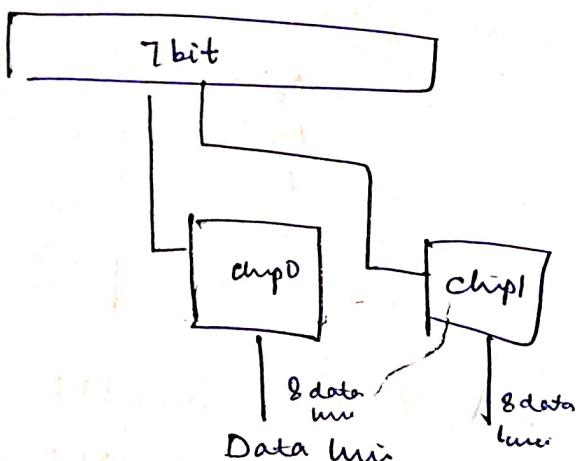
Horizontal Arrangement

No of data lines required are greater than number of data lines in each chip.

CPU generates 7 bit address.

RAM chip — 128×8 bit

$$\text{Total capacity} = 128 \times 16 \text{ bit}$$

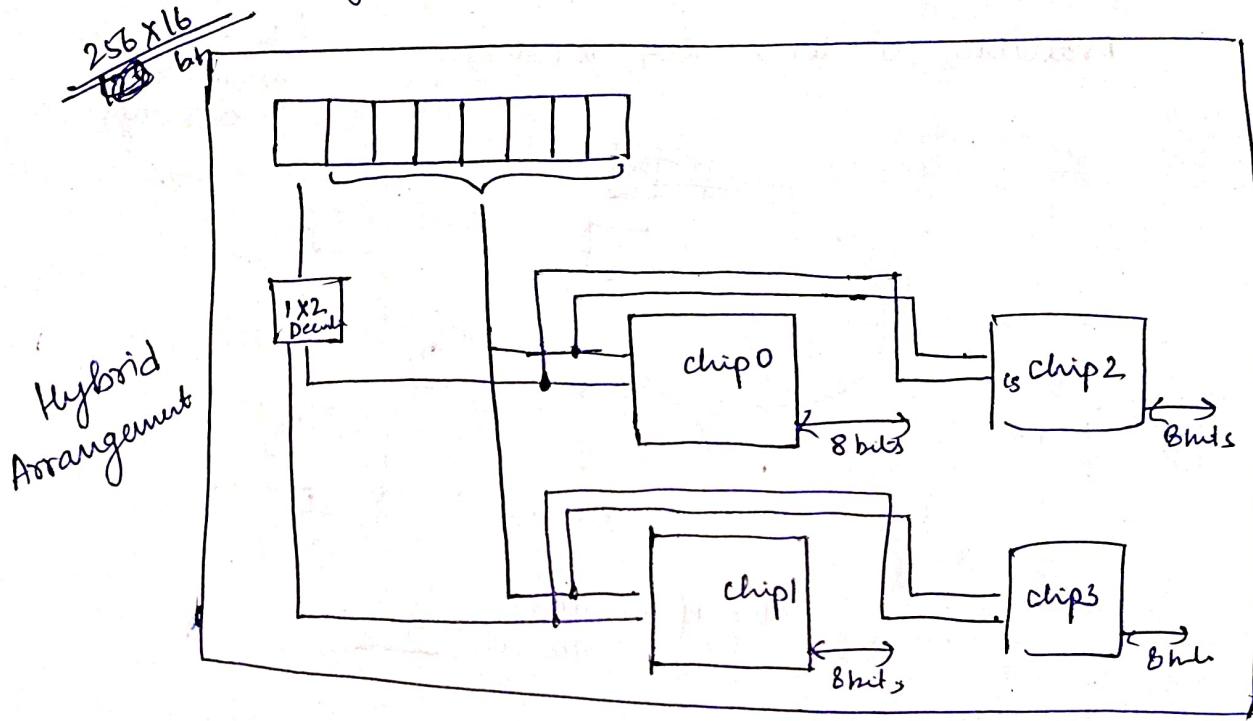


Used when data per location required is more than one chip locations

Ques How many 128×8 bit RAM chips are needed to provide a memory capacity of 256×16 bits.

$$\text{No. of chips required} = \frac{256 \times 16}{128 \times 8} = 4 \text{ RAM chips.}$$

CPU generates 8 bit address.



Gate 2009

How many $32K \times 1$ RAM chips are needed to provide a memory capacity of $256 K$ Bytes?

- (A) 8
- (B) 32
- (C) ~~64~~
- (D) 128

$$\text{No. of RAM chips required} = \frac{256K \times 8 \text{ bits}}{32K \times 1 \text{ bits}}$$

$$= \frac{2^8 \times 2^3}{2^5} = 2^6$$

$$= \underline{\underline{64}}$$

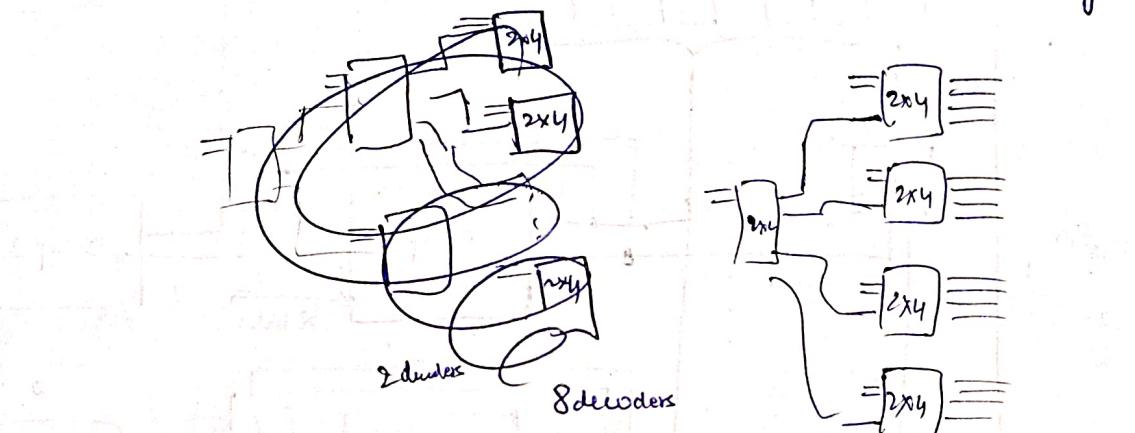
Gate 2013

A RAM chip has a capacity of 1024 words of 8 bits each ($1K \times 8$).
The number of 2×4 decoders with enable line needed to construct a $16K \times 16$ RAM from $1K \times 8$ RAM is $\leftarrow ?$

$$\text{No. of RAM chips required} = \frac{16K \times 16}{1K \times 8} = 16 \times 2 = 32 \text{ chips}$$

$$\text{Number of } 2 \times 4 \text{ Decoder required} = \frac{16 \times 2}{5 \times 32} \text{ Decoder}$$

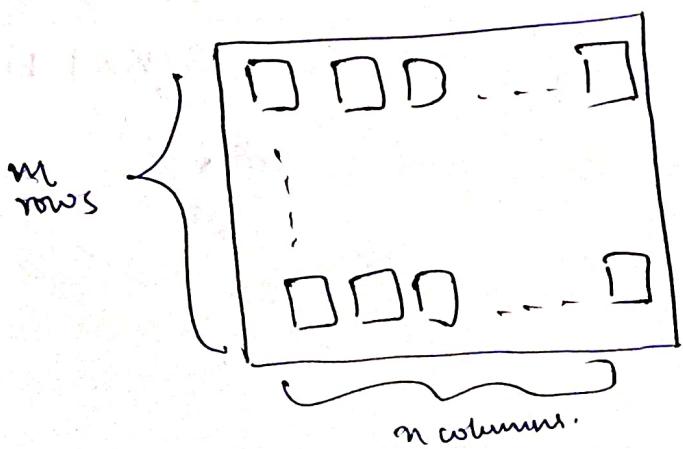
(bcz 16 chips arranged vertically)



$$\text{No. of decoders required} = 5$$

DRAM Refresh

DRAM needs to be periodically refreshed because it is made up of capacitors which store data through electric charge.



In one refresh operation, entire row gets refreshed.

\therefore Total time required for refresh = $\frac{\text{No. of rows}}{1 \text{ chip}} \times \text{Time taken in 1 refresh operation}$

* Time required to refresh n chips = $\frac{1}{n}$ chip refresh-time

Ques A DRAM chip of $512K \times 8$ bits has 1K rows of cells with 512 cells in each row. If DRAM takes 20ns for 1 refresh, then, total refresh time of the DRAM is _____ microseconds.

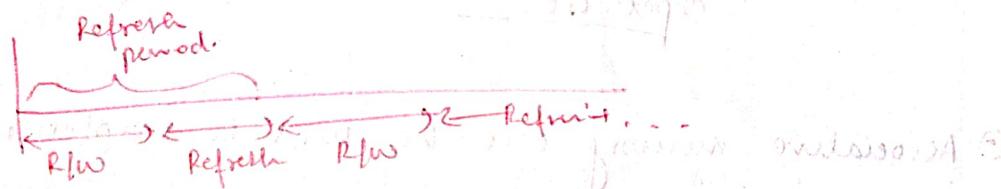
$$\text{No. of rows in DRAM chip} = 1K = 2^{10}$$

$$\text{Time for 1 refresh} = 20\text{ns}$$

$$\text{Total refresh time} = \frac{1K}{2^{10}} \times 20\text{ns} = \frac{2^{10}}{2^{10}} \times 20 \times 10^{-9}\text{s} = 20\mu\text{s}$$

Refresh period

Time period in which the DRAM chip should be refreshed atleast once.



Gate 2018

A 32 bit wide memory unit with a capacity of 1GB & built using $256M \times 4$ bit RAM chips. The number of rows of memory cells is 2^{14} . Time taken to perform 1 refresh operation is 50 ns. The refresh period is 2 milliseconds.

The percentage of time available for performing other memory read/write operations is _____?

$$\text{No. of rows} = 2^{14}$$

$$\text{Total refresh time} = 2^{14} \times 50\text{ns} = 0.819200\text{ms} = 819.2\text{ns}$$

$$\text{R/w in refresh period} = 2\text{ms} - 0.8192 = 1.18\text{ms}$$

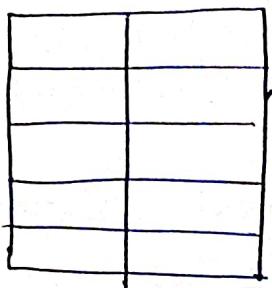
$$\times \text{ of time for r/w} = \frac{1.18}{2} \times 100 = 0.59 \times 100 = \underline{\underline{59\%}}$$

For solving question, $2^{10} = 1K \approx 1000$

Associative memory

① Known as content addressable memory.

② Cells do not addresses.



every cell is divided into 2 parts.

1st part is like address.

If CPU sends a content, it is compared against first part and if that's matched, 2nd part of the cell is sent.

③ The searching takes place parallelly in all the cells.

∴ Associative memory is very fast. (faster than SRAM)
∴ expensive

④ Associative memory can be used to implement TLB (Translation lookAhead buffer.)

Ques

The memory cycle time of a memory is 500 nsec. The maximum rate with which memory can be accessed?

byte addressable

a) 500 Bytes/sec

b) 2000 bytes/sec

c) 2 Mbytes/sec

d) 2 GBytes/sec

Memory cycle time

= time to read/write operation in memory in one cell.

500 nsec → 1 byte can be addressed

$$\therefore \text{Rate} = \frac{1}{500 \times 10^{-9}} = \frac{10^9}{500}$$

$$= \frac{2 \times 10^6}{500} = \underline{\underline{4000}}$$

2 Gbytes/sec

Ques Consider two 16 bits unsigned values A and B. What will be the maximum size of result for

1. Addition of A and B $16+1 = \underline{17 \text{ bits}}$

2. Multiplication of A and B $16 \times 16 = \underline{32 \text{ bits}}$

Size of addition table = $2^{32} \times 17 = \underline{\underline{2^{32} \times 17 \text{ bits}}}$

Size of multiplication table = $2^{32} \times 32 = \underline{\underline{2^{32} \times 32 \text{ bits}}}$

16

Ques A memory M is defined as number of words multiplied by the number of bits per word. The number of address and data lines needed for the memory $M = 2K \times 32 \text{ bits}$ are

No. of address lines needed = $\log_2(2K) = \log_2 2^11 = \underline{11 \text{ lines}}$

No. of data lines needed = $\underline{32}$

Ques

Consider a memory of size $2K \times 8$ bits. What is the size of decoder needed to access the cells of the memory completely?

No. of address lines = $\log_2 2K = 11$ address lines

\therefore Size of decoder needed = $\underline{11 \times 2^11}$

Ques

Statement 1: Static RAM memory devices retain data for as long as power is supplied.

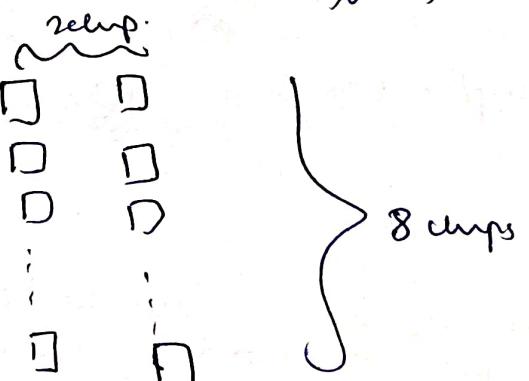
X Statement 2: Static RAM is used when the size of read/write memory required is large.

TOC DBMS OS
 COA CD
 DL CN

Ques How many $8K \times 8$ bits RAM chips are needed to provide a memory capacity of $64K \times 16$ bits? And how the chips are arranged?

$$\text{No. of RAM chips req.} = \frac{64K \times 16}{8K \times 8} = \underline{\underline{16 \text{ RAM chips}}}$$

~~8x2 hybrid~~



Ques A DRAM chip of $512K \times 8$ bits has 256 rows of cells with $2K$ cells in each row. If DRAM takes 20ns for 1 refresh then, total refresh time of the DRAM is _____ microseconds.

$$\text{Refresh time} = \frac{\text{time of refresh}}{1 \text{ row}} = 20\text{ns}$$

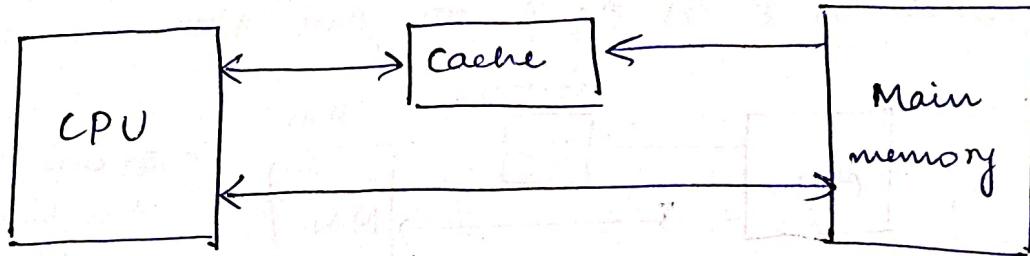
$$\begin{aligned} \text{Total Time to refresh DRAM} &= 20\text{ns} * 256 \\ &= 5120\text{ ns} = \underline{\underline{5.12\mu\text{s}}} \end{aligned}$$

Locality of Reference

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

Temporal locality of reference = CPU address same location again after sometime (based on time)

Spatial locality of reference = CPU addresses nearby location (based on space)



Cache Hit

CPU's demanded content is present in cache.

Cache Miss

CPU's demanded content is not present in cache.

$$\text{Cache Hit Ratio (H)} = \frac{\text{Cache to No. of hits}}{\text{Total memory references.}}$$

$$\text{Miss Ratio} = 1 - H$$

In case of miss, CPU goes to main memory accesses the required content and then brings the block (which contains demanded content) to cache from main memory.

$$\text{Average Memory Access Time} = H * \text{Memory access time for each hit} + (1 - H) \text{Memory access time for each miss.}$$

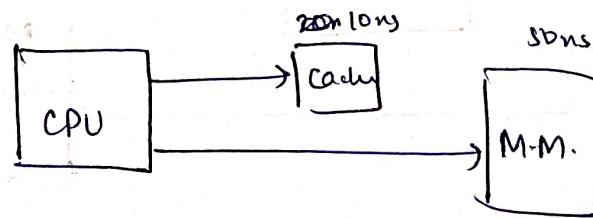
Types of Cache Accesses

1. Simultaneous access

Signal is sent to both cache and main memory at the same time requesting for an address.

अगर data cache में मिल गया तो उसे use कर दिया जाएगा और memory की search करना बहुत कर दिया जाएगा

अगर cache miss होता है तो memory में पढ़ते ही search होते ही होता है। तो avg. time लगता



In case of miss,
time taken = 50ns

$$T_{avg} = H * t_{CM} + (1-H) * t_{MM}$$

2. Hierarchical access

Memory की search तो ही start होती जाती है।
cache miss के बारे में बाकी की बातें

$$T_{avg} = H * t_{CM} + (1-H) (t_{CM} + t_{MM})$$

$$= H * t_{CM} + t_{CM} + t_{MM} - H t_{CM} - H t_{MM}$$

$$T_{avg} = t_{CM} + (1-H) t_{MM}$$

T_{avg} = Avg. memory access time

H = cache hit ratio.

t_{CM} = cache memory access time

t_{MM} = main memory access time

If t_{CM} , t_{MM} are not given explicitly,

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

else

if hierarchical/ level/ serial word is mentioned in question

$$T_{avg} = t_{CM} + (1-H)t_{MM}$$

if simultaneous/ parallel word is mentioned in question

$$T_{avg} = H * t_{CM} + (1-H)t_{MM}$$

Cache look up time
is zero

default
if not found

T_{avg} when Locality of Reference is not included.

$$\text{Simultaneous} \rightarrow H * t_{CM} + (1-H)t_{Block}$$

$$\text{Hierarchical} \rightarrow t_{CM} + (1-H)(t_{Block} + t_{CM})$$

$$t_{Block} = \frac{\text{time required to access block}}{\text{size}} = \text{size} * t_{MM}$$

T_{avg} when Block Transfer is included

$$\text{Simultaneous} = H * t_{CM} + (1-H)(t_{Block} + t_{CM})$$

$$\text{Hierarchical} = t_{CM} + (1-H)(t_{Block} + t_{CM} + t_{Block})$$

Note:- In one access of cache, entire block can be accessed

cache memory
is designed
in such a way.

Ques

Assume that for a certain processor, a read request takes 50 ns on a cache miss and 10 ns on cache hit. Suppose while running a program, it is observed that 90% of the processor's read requests result in a cache hit. The avg. read access time in nanoseconds is _____?

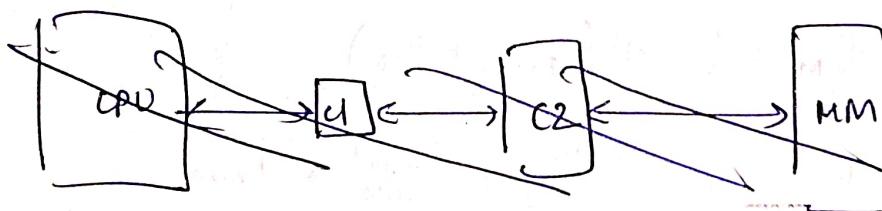
$$\begin{aligned} T_{avg.} &= H * T_{cache\ hit} + (1-H) * T_{cache\ miss} \\ &= 0.9 * 10 + 0.1 * 50 \\ &= 9 + 5 = \underline{\underline{14\ ns}} \end{aligned}$$

Ques

In a 2-level hierarchy, if the top level has an access time of 30 ns and the bottom level has an access time of 150 ns, what is the hit rate on the top level required to give an average access time of 45 ns.

~~$$T_{avg.} = t_{cM} + (1-H)t_{MM}$$~~

Top level memory = closest to CPU



For C2,

~~$$T_{avg.} = t_{cM} + (1-H)t_{MM}$$~~

$$\Rightarrow 150 =$$

~~$$t_{cM} = 30\text{ ns}$$~~

~~$$t_{MM} = 150\text{ ns}$$~~

~~$$T_{avg.} = t_{cM} + (1-H)t_{MM}$$~~

$$\Rightarrow 45 = 30 + (1-H)150$$

$$\Rightarrow \frac{1}{10} = 1-H \Rightarrow H = \frac{9}{10} = \underline{\underline{0.9}}$$

Cache Write

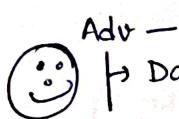
189

Updating value in cache.

Cache Write problem :- If CPU updates the content in cache then the original content in main memory also should be updated.
Also known as write propagation

Write Through

Parallel update of content in both cache memory as well as main memory is done by the CPU.
Simultaneous write operation.



Adv -

→ Data consistency [cache & m.m. will ^{always} have same value of a content]



Disadv -

→ Time consuming [for write operation, mm is accessed regardless of hit or miss]

Simultaneous

Hierarchical

$$t_{\text{Avg Read}} = H \cdot t_{\text{cm}} + (1-H) t_{\text{mm}} \quad H \cdot t_{\text{cm}} + (1-H)(t_{\text{cm}} + t_{\text{mm}})$$

$$t_{\text{Avg Write}} = \max(t_{\text{cm}}, t_{\text{mm}}) = t_{\text{mm}}$$

$$\therefore t_{\text{Avg}} = \frac{\text{fraction of read}}{\text{read}} \cdot t_{\text{Avg Read}} + \frac{\text{fraction of write}}{\text{write}} \cdot t_{\text{Avg Write}}$$

Ques

A system has write through cache with access time 100ns and hit ratio of 90%. The main memory access time is 500ns. The 10% of memory oper references are read.

$$1. t_{\text{Avg. Read}} = H \cdot t_{\text{cm}} + (1-H) t_{\text{mm}} = 0.9 \times 100 + 0.1 \times 500 = 90 + 50 = \underline{140 \text{ ns}}$$

$$2. t_{\text{Avg. write}} = \max(t_{\text{cm}}, t_{\text{mm}}) = 500 \text{ ns}$$

$$3. t_{\text{Avg}} = 0.7 \times 140 + 0.3 \times 500 = 98 + 150 = \underline{248 \text{ ns}}$$

4. Effective hit ratio
let H_e be effective hit ratio

$$H_e \cdot 100 + (1-H_e) \cdot 500 = 248$$

$$\Rightarrow 100H_e + 500 - 500H_e = 248 \Rightarrow 400H_e = 252 \Rightarrow H_e = \frac{252}{400} = \underline{63\%}$$

6

CPU updates only cache. • whenever a block is replaced from cache, it is updated in memory.

Write Back

When a block is to be removed from cache memory, if upgradation is done in the block, it is transferred to main memory.



Adv

→ time saving



Disadv

→ Data inconsistency

$$t_{\text{Avg read or write}} = H * t_{\text{CM}} + (1-H) (t_{\text{MM}} + t_{\text{write back}})$$

[Simultaneous]

$$H * t_{\text{CM}} + (1-H) (t_{\text{MM}} + t_{\text{CM}} + t_{\text{write back}})$$

[Hierarchical]

Write Allocate & No Write Allocate

For read operation

→ for a cache miss, a block is brought from main memory to cache memory.

For write operation

→ Write allocate

For a write miss, missed block should be copied to cache memory from main memory.

→ No write allocate

For write misses, upgradation is done in memory directly without bringing it to cache.

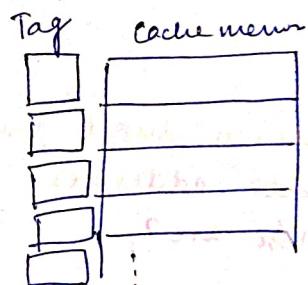
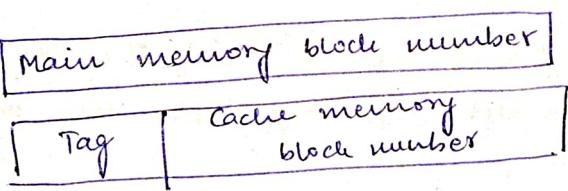
For better performance,
Write through \Rightarrow used no allocate write allocate.
Write back \Rightarrow uses write allocate

Cache Mapping

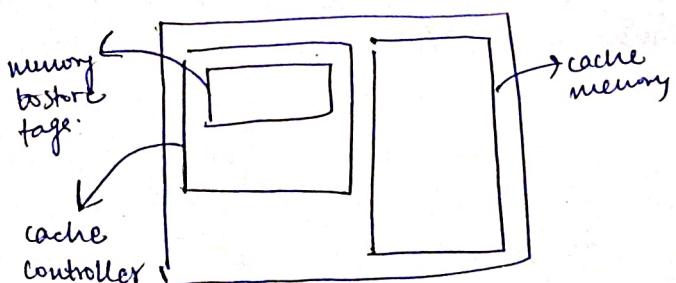
- ① CPU generates main memory address.
- ② Conversion of main memory address to the main memory address is called cache mapping.
- ③ Types -
 - Direct Mapping
 - Set Associative mapping
 - Fully Associative mapping.

Direct Mapping -

- ① Cache memory block = Main memory number $\%.$ No. of blocks in cache memory.
- ② Tag identifies which block of main memory (out of all possible) is present in cache at a given time.



- ③ No. of tag fields required = No. of blocks in cache (line)



Main memory address =

Block number	Block offset
--------------	--------------

Cache memory address =

Tag	Block number	Block offset
-----	--------------	--------------

deciding whether hit/miss

Block offset is not used to tell if it is a miss or hit.

Ques

Block size = 16 bytes

Size of cache memory = 64KB

Size of main memory address = 32 bits

Direct mapping

$$1) \text{ Bits in byte offset} = \log_2 [\text{Block size}] = \log_2 16 = \underline{\underline{4}}$$

$$2) \text{ Bits in cache block number } \log_2 \left(\frac{64 \text{ KB}}{16 \text{ B}} \right) = \log_2 4 \text{ K} = \underline{\underline{12 \text{ bits}}}$$

$$3) \text{ Bits in tag} = \log_2 32 - (4+12) = \underline{\underline{16 \text{ bits}}}$$

$$4) \text{ Tag Directory size} = 16 \times 4 \text{ K} = 2^4 \times 2^{12} = 2^{16} = \underline{\underline{64 \text{ K bits}}} \\ / \text{size of metadata} \\ = \underline{\underline{8 \text{ KB}}} \quad \underline{\underline{8 \text{ kilobytes}}}$$

Ques

Consider a direct mapped cache of size 128KB. The CPU generates 32 bits addresses. The number of tag bits in main memory address are?

Size of cache memory = 128KB

$$\text{No. of tags} = \frac{2^{32}}{128 \text{ KB}} = \frac{2^{32}}{2^7 \times 2^{10} \times 2^3} = \frac{2^{32}}{2^{17}} = \underline{\underline{2^5}}$$

$$\text{No. of tag bits} = \underline{\underline{15}}$$

$$\text{No. of tag bits} = 32 - \log_2 (\text{Cache memory size})$$

Index bit = No. of bits in cache block number

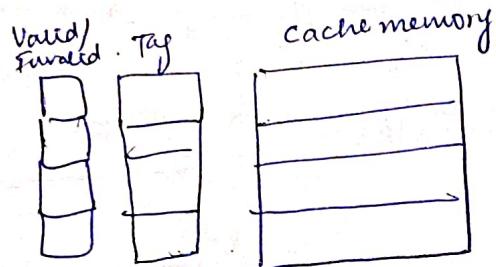
163

Cache Initialization

When computer restarts, garbage value is present in cache.

Possible that CPU generated tag bit matches the tag & garbage tag of block.
So, garbage value is accessed.

To solve this problem, an extra bit called valid/invalid bit is present with each block.



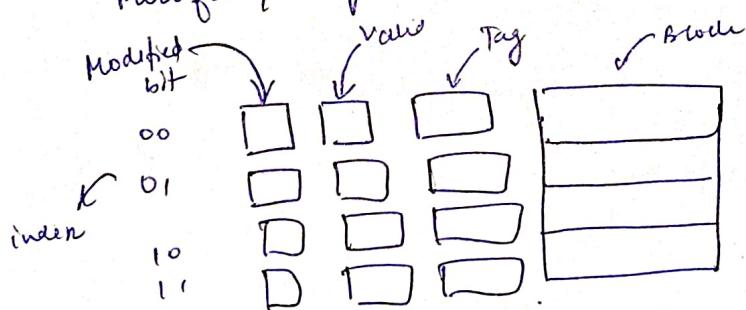
0 → Invalid
1 → Valid

When CPU starts, all the valid/invalid bits are set to 0.

Write Back Cache Performance

Performance can be improved if we can keep track of each cache block which is updated or not.

So, one more bit is present for each block in cache called Modified/Dirty bit.



Dirty bit
0 → only read
1 → write

$$\text{Tag Directory size} = \frac{\text{No. of blocks in cache num}}{\text{No. of blocks in cache num}} \times \left[\begin{array}{l} \text{Tag bits} \\ + \text{Extra bits} \end{array} \right]$$

Ques A 16KB direct mapped write back cache is organized into 512 blocks, each of size 32 bytes. The processor generates 30 bit address. The cache controller maintains free tag information for each cache block comprising of the following —

What is the total size of memory needed at the cache controller to store meta tags for ten cache.

$$\text{No. of blocks} = \frac{16\text{ kB}}{32\text{ B}} = \frac{16}{2^5} = 2^4 \times 2^{10} = 2^9$$

store meta tags for few cache.

$$\text{Size of metadata} = 2^9 [16 + 1 + 1] = 2^9 \times 18 = 2^9 \times 2^4 = 2^{13} \text{ bytes}$$

$$\text{in main memory} = \frac{2^{30}}{2^5} = 2^{25} = 2^9 \times 2 \times 9 = 9 \text{ kbytes}$$

$$\therefore \text{Tag bits} = \log_2 \left(\frac{2^8}{2^9} \right) = 0 \text{ log}_2 16$$

= 16 bits

Targ for write back cache -

→ Simultaneous

65

Consider a machine with a byte addressable main memory of 2^{20} bytes, block size of 16 bytes and a direct mapped cache having 2^{12} cache lines.

Let the addresses of 2 consecutive bytes in memory be $(E201F)_{16}$ and $E202(E2020)_{16}$. What are the cache and tag line address (in hex) for main memory address $(E201F)_{16}$

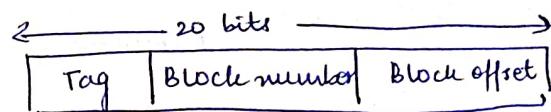
$$\text{Block size} = 16 \text{ bytes}$$

$$\text{No. of cache lines} = 2^{12}$$

$$\text{Main memory address} =$$

$$\text{Line size} = 16 \text{ bytes}$$

$$\text{Block offset} = 4 \text{ bits}$$



$E201F =$
 A binary string $1110\ 0010\ 0000\ 0001\ 1111$ is shown. Arrows point from the first four bits ('1110') to the label 'Tag bits' and from the next four bits ('0010') to the label 'cache line'.

$E201F =$
 A binary string $1110\ 0010\ 0000\ 0001\ 1111$ is shown. Arrows point from the first four bits ('1110') to the label 'Tag', from the next four bits ('0010') to the label '(201)', and from the last four bits ('0001') to the label '(F)'.

(A) E, 201

(B) F, 201

(C) E, E20

(D) 2010F

Ques

Consider a 32kB cache with 32 bytes of blocks. The main memory is of 16MB. The main memory address is 0x67AB1E has tag, block and offset values as.

$$\text{Size of cache} = 32 \text{ kB} = 2^5 \times 2^{10} \text{ bytes} = 2^{15} \text{ bytes}$$

$$\text{Size of 1 block} = 32 \text{ bytes}$$

$$\therefore \text{No. of blocks} = \frac{2^{15}}{2^5} = 2^{10} \text{ blocks}$$

$$\boxed{\text{No. of bits in block offset} = \log_2 32 = 5}$$

$$\text{No. of tag bits} = 32 - (10 + 5) = 17$$

$0x67AB1E =$
 A binary string $0110\ 0111\ 1010\ 1011\ 0001\ 1110$ is shown. Arrows point from the first five bits ('0110 0111') to the label 'Tag', from the next four bits ('1010') to the label 'block no.', and from the last four bits ('1011 0001') to the label 'offset'. Below the string, the values are given as $(0CF)_{16}$, $(158)_{16}$, and $(1E)_{16}$ respectively.

Problem with Direct Mapping
Blocks in main memory have fixed position in cache. Preq. access of blocks occupying same line.

Set Associative mapping

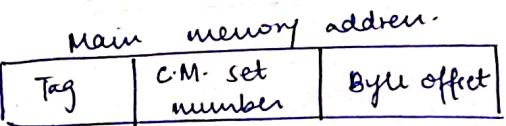
① Multiple blocks at an index.

② n-way set associative mapping \Rightarrow n blocks at one index = $\boxed{8x}$

$$\text{Associativity} = n$$

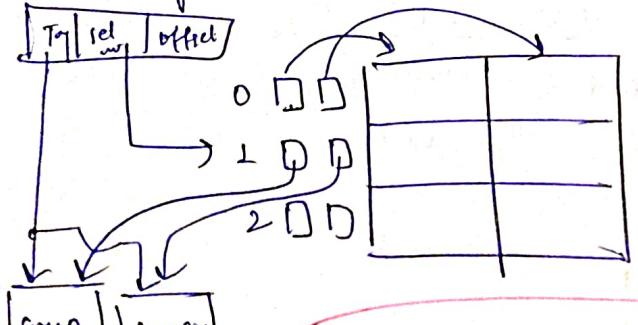
$$\text{No. of indices/sets} = \frac{\text{No. of blocks in cache}}{\text{Associativity}}$$

$$\text{Cache memory set no.} = \frac{\text{Main memory block}}{\text{No. of sets}}$$



$$\text{Tag directory size} = \frac{\text{No. of blocks in cache}}{\text{Associativity}} + [\text{Extra bits} + \text{Tag}]$$

Tag bit is present with each block of a particular set.



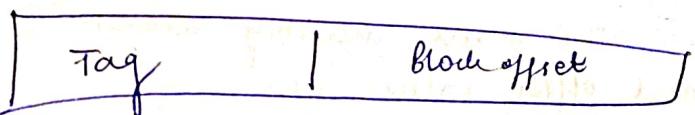
Tag bits are compared with tags of each block in a set.

If matches, it is a hit, if not, then miss.

Fully Associative set mapping

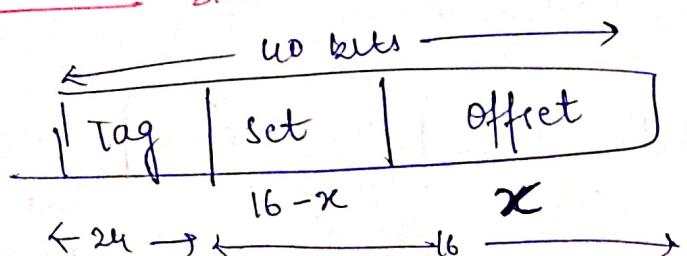
Only 1 set [Similar to set associative mapping]

Main memory address



~~Y'all 1016~~

The width of the physical address on a machine is 40 bits. The width of tag field in a 512 kB 8-way set associative cache is _____ bits



$$\text{No. of bits in tag field} = 24$$

$$\text{Set block size} = 2^x$$

$$\text{No. of blocks} = \frac{2^{19}}{2^x} = 2^{19-x}$$

$$\text{No. of sets} = \frac{2^{19-x}}{8} = 2^{16-x}$$

Hardware implementation of Direct mapping -

167

No. of comparators required = 1

Size of comparator = No. of tag bits.

No. of multiplexers required = No. of tag bits

Size of multiplexers = No. of blocks \times 1.

Hit Latency Time

$$= \underbrace{1 \text{ MUX delay}}_{\text{for tag selection}} + \underbrace{1 \text{ comparator delay}}_{\text{for tag comparison}}$$

Hardware implementation of k-way set associative mapping

No. of comparators required = k

Size of comparator = No. of tag bits

No. of multiplexers req. = $k \times$ No. of tag bits

Size of multiplexers = No. of sets in cache \otimes 1

k-input OR gate = 1

Hit Latency time

$$= \underbrace{1 \text{ MUX delay}}_{\text{for tag selection}} + \underbrace{1 \text{ comparator delay}}_{\text{for tag comparison}} + \underbrace{\text{OR gate delay}}_{\text{delays}}$$

Hardware implementation of fully associative set mapping

No. of comparators required = No. of blocks in ~~main~~ cache memory.

Size of comparator = No. of tag bits

No. of multiplexers = 0

1 OR gate — n input OR gate

Hit Latency time

$$= \underbrace{1 \text{ comparator delay}}_{\text{delays}} + \underbrace{\text{OR gate delay}}_{\text{delays}}$$

Cache Miss Penalty

Time Required to bring a missed block from main memory to cache.

Ques 2019

A certain processor deploys a single level cache. The cache block size is 8 words and word size is 4 bytes. The memory system uses a 60MHz clock. To service cache miss, the memory controller first takes 1 cycle to accept the starting address of the block, and finally transfers the words of the requested block at the rate of 1 word/cycle. The maximum bandwidth for the memory system when the program running on the processor issues a series of read operations is $\underline{\quad \times 10^6 \text{ bytes/sec.}}$

$$\text{Block size} = 8 \times 4 = 32 \text{ bytes}$$

$$\text{Clock rate} = 60 \text{ MHz}$$

$$\text{Time to transmit 1 word} = 1 \text{ cycle}$$

$$8 \text{ words} = 8 \text{ cycles}$$

$$\text{Total cycle in case of miss} = \underline{\text{cycles}}$$

$$= 8 + 1 + 3$$

$$= 12 \text{ cycles}$$

32 bytes

To transfer 32 bytes, $\frac{12}{32}$ cycles are req.

To transfer 1 byte

$$\text{Bandwidth} = \frac{32}{12} \text{ bytes/cycle}$$

$$= \frac{32}{12} \times 60 \times 10^6 \text{ bytes/s.}$$

$$= 160 \times 10^6 \text{ bytes/s.}$$

Types of cache miss

1) Cold/Compulsory miss

First time access of a block will always cause miss.

Solution: increase no. of blocks the block size of cache/memory.

2) Capacity miss

If it is not cold miss and miss occurs when the cache is full.

Solution to reduce - increase the size of cache.

3) Conflict miss

If it's not a cold & capacity miss but due to the set is full and miss occurs.

Solution to reduce → Increase associativity.

Example

- ① No of blocks in cache = 4
- ② 2 way set associative cache
- ③ Cache is initially empty
- ④ LRU Replacement policy
- ⑤ CPU Requests for main memory

{ No of sets in cache = 2

blocks - 0, 4, 0, 8, 0, 4, 1, 3, 1, 5, 1, 3

0	8 4	4 0
1	4 3	1 1

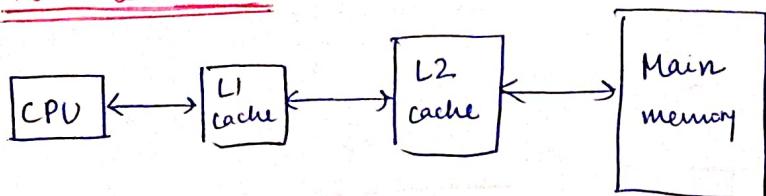
Concept
Timeit
Answers &
JCTd
etc!!

Goals of using cache memory

1. Minimize Access Time → use small cache
2. Maximize Hit Rate → use large cache
3. Minimize Miss Penalty.

} multilevel cache

Multilevel cache



Access time	t_1	t_2	t_{MM}
Hit ratios	H_1	H_2	

L1

Simultaneous access -

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

169

Hierarchical access -

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

OR

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

Policy

→ Inclusion: Blocks in L1 are also present in L2 [Block is copied from L2 to L1]

→ Exclusion: Blocks in L1 need not be present in L2 [Block is moved from L2 to L1]

L2 is filled with only replaced (victim blocks) of L1

∴ it is known as victim cache.

Ques

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

$$t_{CM} = 25\text{ns} \quad t_{MM} = 300\text{ns} \quad t_{block} = 300 \times 4 = 1200\text{ns}$$

$$\text{twrite-back} = \text{% of dirty blocks} \times t_{block} = 0.1 \times 1200 = 120\text{ns}$$

$$T_{avg} = H \cdot t_{CM} + (1-H) [t_{CM} + t_{block} + \text{twrite-back}]$$

$$= 0.9 \times 25 + (1-0.9) (25 + 1200 + 120)$$

$$= 25 + 0.1 (1200 + 120)$$

$$= 25 + 0.1 \times 1440 = 25 + 144 = \underline{\underline{169\text{ ns}}}$$

Magnetic Disk

$$\text{Disk Access Time} = \text{Seek Time} + \text{Rotational Latency} + \text{Transfer Time}$$

Data transfer
dist. etc off
dist., rotation
speed of
disk same &
 t_{eff}

Time required
to position read/
write head on
a particular track
(will be mentioned
directly in
ques]

Time required
to position read
write head on a
particular sector
in track

$$RL = \frac{\text{1 rotation time}}{2}$$

Time
required
to
transfer data
from sector.

$$TT = \frac{\text{1 rotation time}}{\text{no. of sectors per track}}$$

* At any given time, only one read-write head can work

Ques

Consider a disk with 16 platters, 2 surfaces per platter, 1K tracks per surface, 1K sectors per track and 2048 bytes per sector. Disk rotates with 3000 rpm. Seek time is 10ms.

- 1) Capacity of disk
- 2) No. of bits req. for addressing
- 3) Disk access time
- 4) Disk Transfer rate.

3000 rotations — 60 sec

$$1 \text{ rotation} = \frac{60}{3000} \text{ sec} = \frac{1}{50} \text{ sec} = \frac{1000}{50} \text{ ms} = 20 \text{ ms.}$$

$$\text{Rotational latency} = \frac{20}{2} = 10 \text{ ms}$$

$$\text{Capacity of disk} = 16 \times 2 \times 2^{10} \times 2^{10} \times 2048 = 2^4 \times 2 \times 2^{20} \times 2^{10} = \underline{\underline{2^{\underline{\underline{36}}}}}$$

$$\text{No. of bits in addr} = \underline{\underline{2^{\underline{\underline{36}}}}} \quad \text{No. of sectors} = 16 \times 2 \times 1K \times 1K = \underline{\underline{2^{\underline{\underline{25}}}}} \therefore \underline{\underline{25 \text{ bits}}}$$

$$\text{Disk access time} = 10 \text{ ms} + 10 \text{ ms} + \frac{20 \text{ ms}}{2^1 \times 1000} = 20 + 0.02 = \underline{\underline{20.02 \text{ ms}}}$$

~~Rotational latency~~ ~~20 ms~~ bytes/sec. ~~0.02~~ 0.02 ms \rightarrow 1 sector data = 2048 bytes $= 2 \text{ KB}$

$$\therefore 1 \text{ s} \rightarrow \frac{2048 \times 1000}{0.02} \text{ bytes} = 2^{10} \times \underline{\underline{100 \text{ MBPS}}}$$

Addressing is done

↳ track wise

↳ cylinder wise.

Cylinder no., surface no., sector no. \rightarrow

No. of sectors per cylinder =

No. of sectors
in 1 track * No. of surfaces.

Address $\langle C, h, s \rangle$

No. of sectors per track = n_t

No. of sectors per cylinder = n_c

$$\text{Sector no.} = C * n_c + h * n_t + s$$

$$C = \left\lfloor \frac{\text{sector no.}}{n_c} \right\rfloor$$

$$h = \left\lfloor \frac{(\text{sector no.}) / n_c}{n_t} \right\rfloor$$

$$s = (\text{sector no.} \% n_c) \% n_t$$

Ques Consider a disk with 2K sectors per track. Disk rotates with 6000 rpm. Seek time is 20ms. Calculate disk access time.

Seek time = 20ms

$$6000 \text{ rotations} = 60 \text{ s} \\ 1 \text{ rotation} = \frac{60}{6000} \text{ s} = 0.01 \text{ s} = 10 \text{ ms}$$

$$\text{Rotational latency} = \frac{10}{2} = 5 \text{ ms.}$$

$$\text{Data transfer rate} = \frac{10}{2000} = 0.5 \times 0.01 \text{ ms} = 0.005 \text{ ms.}$$

$$\text{Disk access time} = 20 \text{ ms} + 5 \text{ ms} + 0.005 \text{ ms} = \underline{25.005 \text{ ms}}$$

Constant linear Velocity — Variable sector capacity
Constant angular velocity — Constant sector capacity

* Ques

$$\text{Ans} \quad \frac{60}{6000} = ?$$

Flynn's classification of Computer

Pipeline system	SISD	Single Instruction stream	single Data stream.
	SIMD	single Instruction stream	Multiple Data stream
Multiple Pipeline	MISD	Multiple Instruction stream	single Data stream (theoretical)
	MIMD	Multiple Instruction stream	Multiple Data stream

PIPELINING

- ① Pipelining is useful when same processing is applied over multiple inputs.
- ② Pipelining is a technique of decomposing a sequential process into suboperations.
Suboperations are performed in segments.
All segments perform their respective suboperations parallelly on different inputs.

Pipeline Cycle Time

Minimum time in which all segments can perform their respective suboperation.

$$\text{Max}[t_{\text{seg}1}, t_{\text{seg}2}, t_{\text{seg}3}, \dots]$$

General Consideration

Consider a k segment pipeline.

With clock cycle time = t_p

to perform n tasks.

Time required to perform 1st task = $k \times t_p$

Time required to perform remaining $n-1$ tasks = $(n-1) t_p$

∴ Total time required for all n tasks = $(k+n-1)t_p$

Consider a non-pipeline system that takes t_n time to perform a task

Time required to perform n tasks = $n \times t_n$.

Performance of pipeline = Speedup (S) = $\frac{\text{Nonpipeline time}}{\text{pipeline time}}$

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

When $n \gg k$,

$$\frac{n t_{p1}}{n t_p} = \frac{t_n}{t_p}$$

more/ideal speedup
value of n is not given.
so

When time required to perform 1 task is equal a pipeline and non-pipeline

$$t_n = k t_p \Rightarrow s = \frac{t_n}{t_p} = k$$

Speed up is always less than k

Pipeline hazard.

Situations that prevent the next instruction from being executed during its designated clock cycle.

- Structural Hazard / Resource Conflict
- Data Hazard / Data Dependency
- Control Hazard / Branch Difficulty

Structural Hazard

2 different segments try to use same resource at a same time.

Solution — → Incur stall cycles

→ Increase number of resources

Data Hazard / Dependency

Result of a segment is input to the next segment.

$$\begin{aligned} i: R1 &\leftarrow R2 + R3 \\ it: R4 &\leftarrow R1 + R5 \end{aligned} \quad \left\{ \text{dependency.} \right.$$

IF DA OF EX WB

IF DA

DA

EX WB

DA

WB

A general instruction pipeline cannot detect data dependency.

The solution of data dependency is provided by compiler.

$$R1 \leftarrow R2 + R3$$

====

$$R5 \leftarrow R1 + R4$$

} independent (no operation instructions)

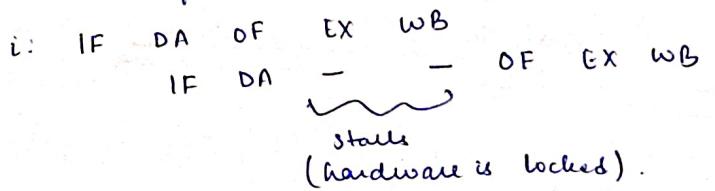
Delayed Load.

Software Solution.

③ Hardware solutions are used when pipeline can detect and solve data dependency.

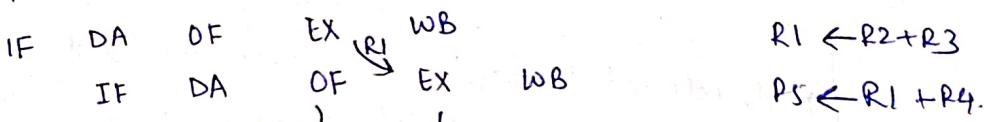
Hardware interlock.

stalls in case of data dependency.



Operand Forwarding

The result generated by i th instruction is transferred to EX phase of $(i+1)$ th instruction from the accumulators.



wrong value of R1 is fetched
 R1 is also fetched.

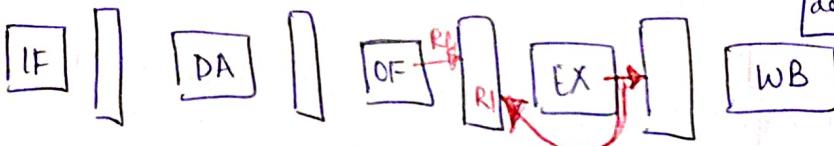
value of R1 from accumulator is used.

$$R1 \leftarrow R2 + R3$$

$$PS \leftarrow R1 + R4.$$

no stall

Operand forwarding can be used in this example because ALU operation in both the instructions.



$$i\text{th} \rightarrow R1 \leftarrow M[\text{Add}]$$

$$(i+1) \rightarrow R2 \leftarrow R1 + R5$$

} operand forwarding cannot be used.

No stalls for ALU to ALU dependency.
 but for memory to ALU or ALU to memory dependency, stalls required.

CONTROL HAZARD / BRANCH DEPENDENCY

for branch instruction, the CPU has to wait till branch outcome is available.

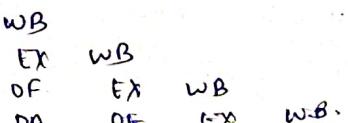
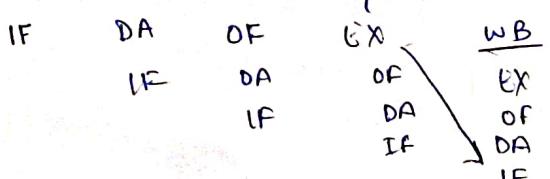
A general pipeline cannot detect branch dependency

↓
 compiler provides solution — Delayed branch.

inserts 3-4 no operation instructions below branch instruction.

→ after execute phase, target instruction is found out.

Branch



→ Target inf.

Hardware solution.

→ Branch Prediction

Predicts whether the current branch instruction will take jump or not.

If prediction is correct — no problem.

If prediction is not correct — Roll back to correct branch instruction & execute target.

Classification of Data Hazard / Dependency

True dependency:

→ RAW (Read After Write)	i: $R1 \leftarrow R2 + R3$ j: $R5 \leftarrow R1 + R4$	[j reads before written by i. Hence j incorrectly gets old value]
→ WAW (Write After Write)	i: $R1 \leftarrow R2 + R3$ j: $R1 \leftarrow R4 * R5$	[j writes in destination before i writes it]
→ WAR (Write After Read)	i: $R1 \leftarrow R2 + R3$ j: $R2 \leftarrow R9 + R7$	[if j modifies R2 before i uses its value i incorrectly reads new value]

Anti dependency

false dependency

Solution :-

Register Renaming

$$\begin{aligned} i: R1 &\leftarrow R2 + R3 \\ j: R6 &\leftarrow R9 + R7 \end{aligned}$$

(WAR)

$$\begin{aligned} i: R6 &\leftarrow R2 + R3 \\ j: R4 &\leftarrow R4 * R5 \end{aligned}$$

(WAW)

CPI (Cycles Per Instructions)

No. of cycles needed to execute 1 instruction.

For pipelined system,

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

In ideal case,

$$CPI = 1$$

If because of hazards, total extra (stall) cycles = x

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

In ideal case,

$$CPI = \frac{n+x}{n} = 1 + \frac{x}{n}$$

With hazard,
one instruction execution cycle = CPI_{avg} + fp.

Pipelining Hazards

- ① Structural Hazards → [diff. segments same resource required] → incur stalls → increase no. of stalls.
- ② Data Hazards → RAW (Read after write) → True dependency
WAW → False dependency
WAR → anti dependency.
→ operand forwarding → stalls.
- ③ Control Hazard → predict branch → static (branch history table)
→ stalls

Amdahl's law

Madeleine risk dimensions

used to calculate speedup in multiprocessor system.

Program Status word)₃₂ = std::vector<int> { 0, 1, 2, 3, 4, 5, 6, 7 };

Program Status Word (PSW) is a collection of data that operating system maintains to keep track of system's current state.

8 bit
Register

Four flags are called conditional flags -

- CY (Carry)
 - AC (Auxiliary Carry)
 - Parity (P)
 - OV (Overflow)

RS0 and RS1 are register selection bit

The diagram illustrates the 8-bit register structure with fields labeled from 7 to 0 (LSB) at the top. The fields are:

- CY**: Carry Flag
- AC**: Auxiliary carry flag
- RSI**: Register bank status
- RSO**: Register bank status
- OV**: Overflow flag
- P**: Parity flag

Annotations below the register fields include:
 - A bracket under the first two fields (CY and AC) is labeled "Auxiliary carry".
 - A bracket under the last three fields (RSI, RSO, OV) is labeled "Register bank status".
 - A red arrow points to the RSI field with the label "carry".
 - A red arrow points to the RSO field with the label "status".
 - A red arrow points to the OV field with the label "Overflow flag".
 - A red arrow points to the P field with the label "Parity flag".
 - A red box encloses the RSI, RSO, and OV fields with the label "carry status".
 - A red box encloses the RSI, RSO, and OV fields with the label "status".
 - A red arrow points to the RSI field with the label "carry".
 - A red arrow points to the RSO field with the label "status".
 - A red arrow points to the OV field with the label "Overflow flag".
 - A red arrow points to the P field with the label "Parity flag".

In 16 bit PSW —

→ 7 bits are unused

→ 9 bits are used

→ 3 control flags

→ 6 status flags

Overflow — 11th bit ~~last digit~~

Control Flags -

12-15, 1, 3, 5 - unrued

CE \rightarrow 0th

$$p_1, p_2 \rightarrow 2$$

Rainy

C - zero flag

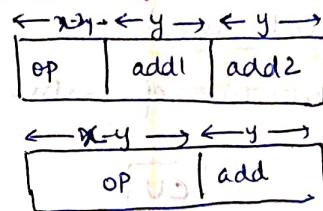
→ DF (Direction Flag) — 10th bit

→ IF (Interrupt Flag) — 9th bit

→ TF (Trap Flag) — 8th bit

System supporting variable length opcode

length instruction
2 address instruction:



1 address instructions:

m 2 address instructions supported
n 1 address instructions supported

$$\text{Max. } n = \left(2^{x-2y} - m\right) \cdot 2^{(x-y) - (x-2y)}$$

Min n = 1

$$\left(2^{x-2y} - m\right) \cdot 2^{\text{No. of unused combinations}}$$

opcode of address instruc.
address inc -
opcode length of 2 address instructions.

IEEE - 754 Floating Point Representation

Bias = 127

Single precision

$$(-1)^S \times 1.M \times 2^{E-127}$$

32 bits
1 sign bit (s)
8 exponent (E)
23 mantissa (M)

For denormalized no., E = -126 always

10

11 10 10 Number

S E M

0 000....0 00...0 +0

1 000...-0 00...0 -0

0 111...-1 00...0 +∞

1 111...-1 00...0 -∞

0/1 111...-0 M ≠ 0 Not a number

0/1 000...-0 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0 Denormalized no.

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0 Implicit Normalization

0/1 E ≠ 11...-1 M ≠ 0 (Normal no.)

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

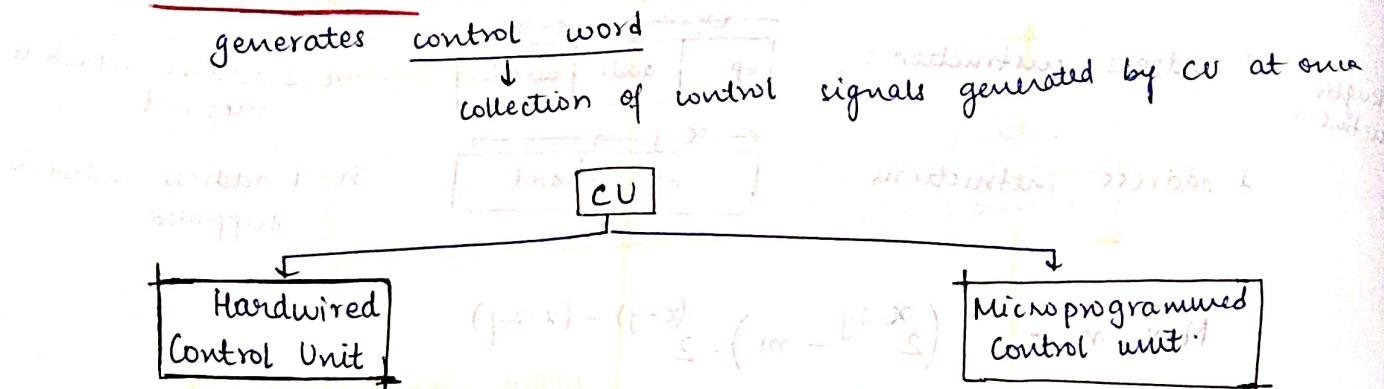
0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...-1 M ≠ 0

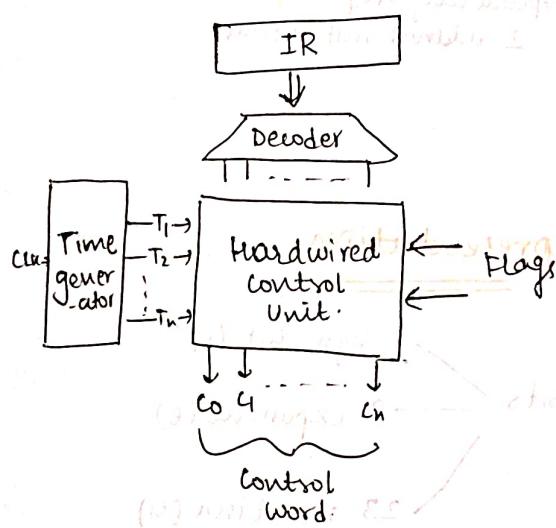
0/1 E ≠ 00...-0 M ≠ 0

0/1 E ≠ 11...

Control unit



- ① faster
- ② difficult to rearrange / change control logic



- ① slower
- ② updation of control logic is easy.

Vertical Microprogramming

① signals divided into groups such that at a time, only 1 signal can be active

- ② smaller control word
- ③ slower
- ④ decoders used

Horizontal Microprogramming

① one bit for every control signal.

- ② larger control word
- ③ faster
- ④ decoders not used.

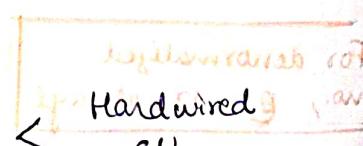
On the basis of operating speed -

Vertical Microprogrammed CU

Slowest

Horizontal Microprogrammed CU

Fastest



RISC

- ① Less no. of instructions
- ② Fixed length ins (variable opcode)
- ③ Less addressing modes
- ④ Hardwired CU
- ⑤ More no. of registers
- ⑥ only R-R operation

CISC

- ① More no. of instructions
- ② variable length ins (fixed opcode)
- ③ More addressing modes
- ④ μprogrammed CU
- ⑤ less no. of registers
- ⑥ Reg. to Memory operation also present.

Preferrable for pipelined CPU

IO Mapped IO / Isolated IO

common data bus and address bus.

different control bus

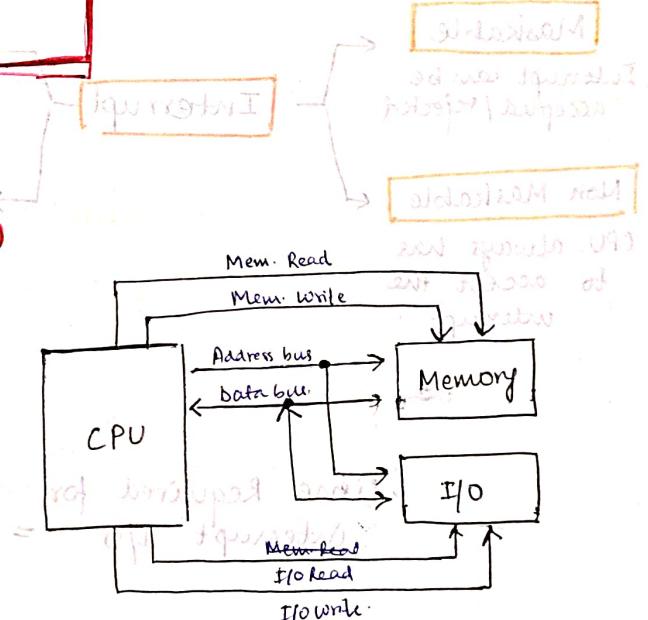
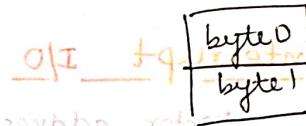
CPU distinguishes b/w memory & I/O using control signal.



Big endian -



little endian -

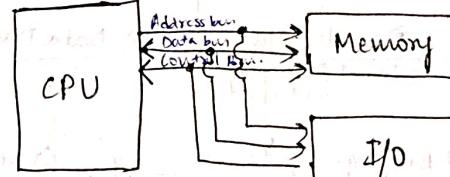


→ no wastage of memory.

Memory Mapped IO

common address bus, data bus, control bus.

Some addresses in memory are reserved for I/O.



→ memory wastage

< Data bus is 2 sided
Address bus is 1 sided >

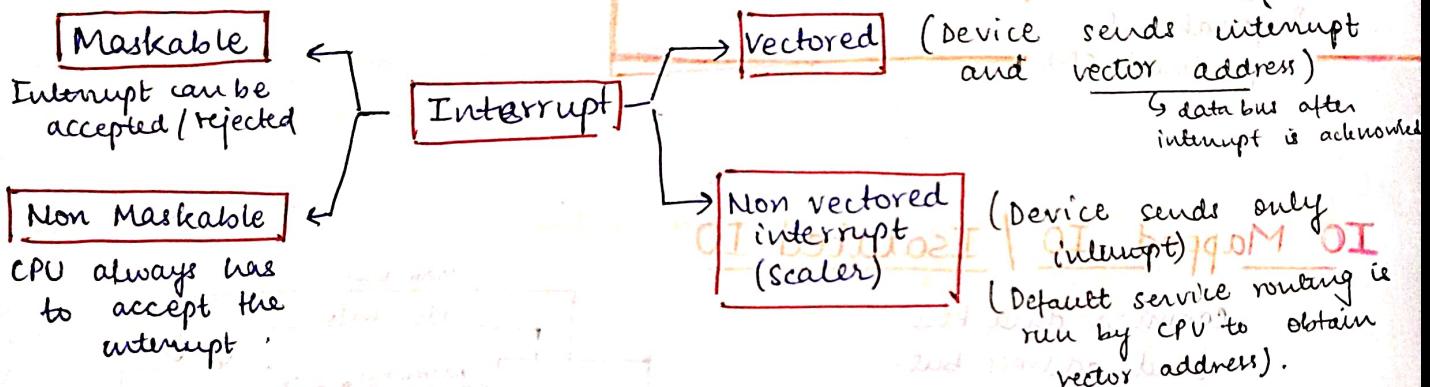
Modes of IO Transfer

1. Programmed I/O

$$\text{Total Time} = \frac{\text{Time to read status of I/O device}}{\text{Instruction execution time}} + \frac{\text{Time to check if status flag is set or not}}{\text{Instruction execution time}} + \frac{\text{Data transfer time}}{\text{Data transfer rate}}$$

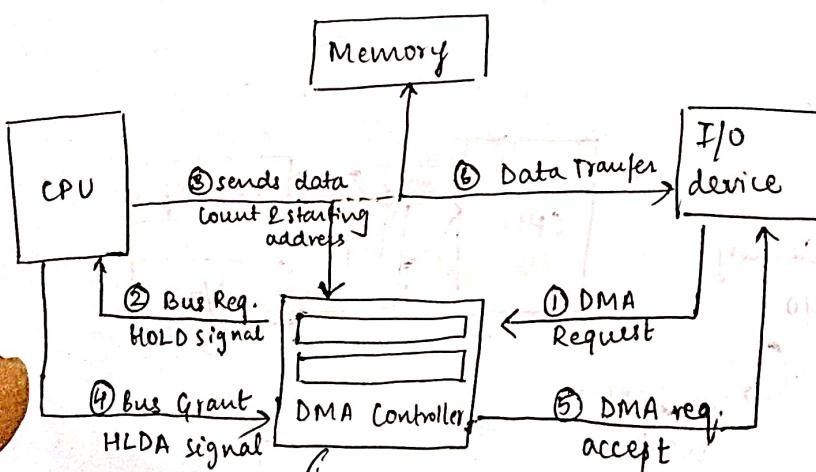
2. Interrupt I/O

Vector address = address / location of interrupt service Routine.



$$\text{Time Required for interrupt I/O} = \frac{\text{Interrupt overhead time}}{\text{Instruction execution time}} + \frac{\text{Service time}}{\text{Instruction execution time}}$$

3. Direct Memory Access



Data count - No. of bytes/words to be transferred

Starting memory address - Starting from where data is to be transferred.

Time is saved in this method as CPU does not have to store the present state of the system in stack.

Time taken to prepare the data in I/O = t_x will be given in question

Time taken to transfer data to memory = t_y depends upon clock freq./ cycle time.

Imp.

GATE
2022

% of Time CPU is blocked —

$$1. \text{Burst mode} = \frac{t_y}{t_x + t_y}$$

$$2. \text{Cycle stealing mode} = \frac{t_y}{t_x}$$

$$3. \text{Interleaving mode} = 0$$

Size of ROM required to store table —

n bit A \times n bit B.

Size of multiplication table = $2^{2n} \times 2^n$

Size of addition table = $2^{2n} \times (n+1)$

RAM
(Random Access
Memory)

- Static
 - flip flops
 - no refresh
 - faster
 - expensive
 - ★ → higher power consumption.
- Dynamic
 - data is stored in form of electric charge.
 - periodic refresh.
 - slower
 - less expensive.
 - ★ → low power consumption.

- ① In one refresh operation, entire row gets refreshed.
- ② Time to refresh n chips = Time to refresh 1 chip.

$$\boxed{\text{Refresh period} = \frac{\text{Read/write time}}{\text{Total refresh time}} + \text{Time}}$$

① Associative memory is used to implement

TLB (Translation Look Aside Buffer)

↳ stores frequently used pages.

TLB Hit —

Page is retrieved from TLB

$$\text{Time} = T_{TLB}$$

TLB Miss —

Page is retrieved from page table & then memory.

$$\rightarrow \text{Time} = T_{TLB} + T_{MM} + T_{MM} = T_{TLB} + 2T_{MM}$$

coz page table is stored in main memory.

$$\therefore T_{avg} = H \times T_{TLB} + (1-H)(T_{TLB} + 2T_{MM})$$

Cache write

→ Write Through (simultaneous update in M.M.)

→ Write Back (cache block transferred to M.M. when replaced)

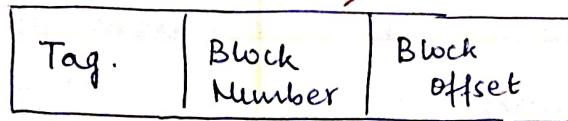
→ Write allocate — write miss → missed block should be copied from M.M. to cache memory

→ No write allocate — write miss → update is done in M.M. directly

Cache Mapping

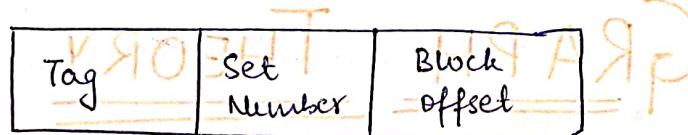
→ Direct Mapping.

decide if it is a cache hit or cache miss.



$$\text{Cache memory block number} = \frac{\text{Main memory block number}}{\text{No. of blocks in cache memory}}$$

→ Set Associative Mapping.

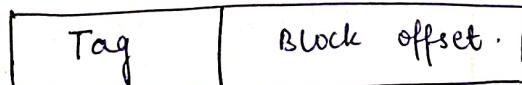


$$\text{Cache memory set no.} = \frac{\text{Main memory block number}}{\text{No. of sets in cache memory}}$$

→ Fully Associative mapping

only 1 set.

→ similar to associative memory.



$$\text{Size of metadata/Tag directory size} = \left[\frac{\text{No. of Tag bits}}{\text{Extra bits}} \right] * \text{No. of lines in cache memory.}$$

Types of cache miss.

1) cold / Compulsory miss

Soln → increase block size of cache memory.

2) Con Capacity miss

Soln → increase size of cache.

3) Conflict miss

Soln → increase associativity of cache.