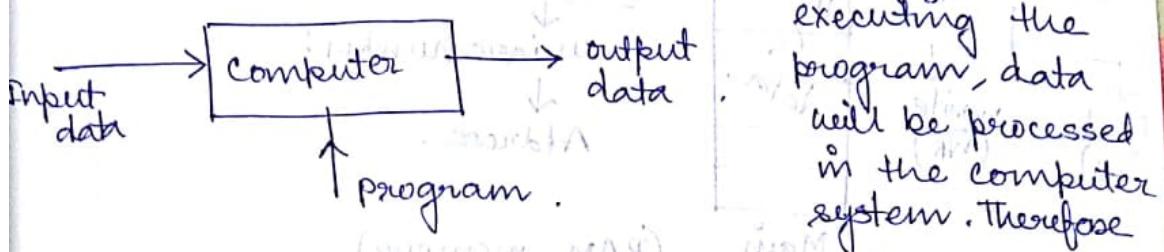


Computer fundamentals

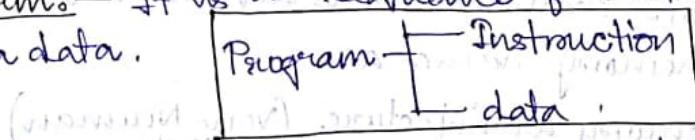
Computer :- It is a computational device used to process the data under the control of a program. By



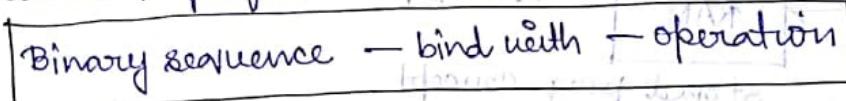
executing the program, data will be processed in the computer system. Therefore

computer functionality is program execution.

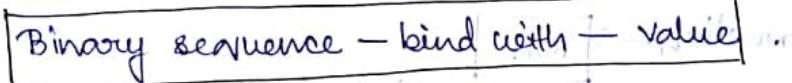
Program :- It is a sequence of instructions along with a data.



Instruction :- It is a binary sequence designed inside the processor to perform some task.



Data :- It is a binary sequence which is associated with value based on the data format.



functionality of computer :-

computer functionality is program execution. Let us consider computer components as a reference to justify the computer functionality.

① CPU (Processing)

② Memory (Storage)

③ I/O (External communication)

Memory organisation :-

→ Memory is a storage element in the computer, so program is always stored in the memory.

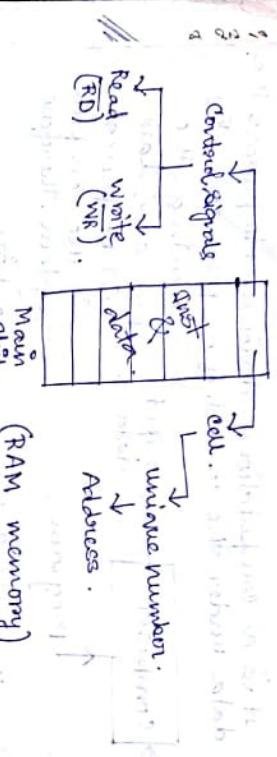
→ Memory chip is divided into equal parts called as cells. Each cell is identified with unique number called as address.

→ Memory chip recognises the control signals to indicate the type of operation.

→ Address register → Processor → Control Unit
data present in data bus

$(0000\ 0000\ 00000000)_2$ address space.

Each cell is not given 16 bit address line. Design is implemented to address each cell logically by implementing address decoding mechanism.



Note :- Processor architecture is divided into 2 types based on the memory name as:

(i) Single memory architecture (Non-Neumann)

Ex :- MP (Memory Processor).
Storage prog concept.

(ii) Multi-memory architecture :- [Harvard arch]

Instruction \rightarrow ROM
Data \rightarrow RAM. Ex :- microcontroller.

\rightarrow CPU generates the memory request to access the program from memory chip i.e., memory cell.

Memory request \rightarrow .

Address | Control signal
↓
indicates type of memory cell operation.
Memory cell.

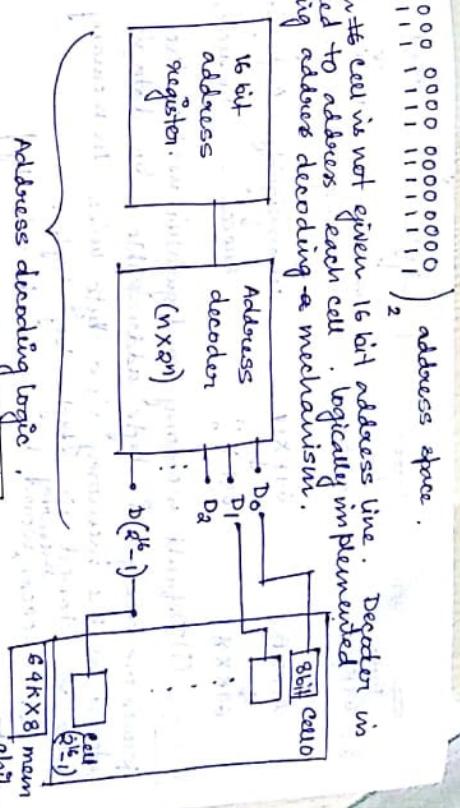
Memory chip configuration :-
Number of address bits required to access the memory always depends on memory capacity. Memory chip config is represented as $64K \times 8$, $64K \times 16$, $64K \times 32$, $128K \times 16$, $256K \times 8$ memory chip. Total cell size (No. of bits stored in cell) is $2^6 \times 8 = 2^6 \times 16$ cells.

cells in memory chip \times cell size = $2^{10} \times 2^6 = 2^{16}$ cells.
16 bit address. Address range / address space mem map of a memory chip.

Address decoding logic.

| MP | Bit | Nibble | Byte | Word |
|-------|-----|--------|------|------|
| 8bit | 4 | 4 | 8 | 8 |
| 16bit | 1 | 4 | 8 | 16 |
| 32bit | 1 | 4 | 8 | 32 |

Byte is always having unique meaning. word is having multiple meaning hence ambiguous & unsuccessful.



Byte addressable vs word addressable memories :-
Based on the cell size, memory chip configuration is divided into two types.

i) Byte addressable memory :-
ii) Word addressable memory :-

\rightarrow When the memory cell size is 8 bit, then the corresponding address space is called as byte addressable.

Ex :- 4x8, 2 bit address cell holds 8 ".

Ex :- 8x8 : 8 bit address cell holds 8 ".

Ex :- 16x8 : 4 bit address cell holds 16 ".

Ex :- 32x8 : 5 bit address cell holds 32 ".

Ex :- 64x8 : 6 bit address cell holds 64 ".

Ex :- 128x8 : 7 bit address cell holds 128 ".

Ex :- 256x8 : 8 bit address cell holds 256 ".

Ex :- 512x8 : 9 bit address cell holds 512 ".

Ex :- 1024x8 : 10 bit address cell holds 1024 ".

Ex :- 2048x8 : 11 bit address cell holds 2048 ".

Ex :- 4096x8 : 12 bit address cell holds 4096 ".

Ex :- 8192x8 : 13 bit address cell holds 8192 ".

Ex :- 16384x8 : 14 bit address cell holds 16384 ".

Ex :- 32768x8 : 15 bit address cell holds 32768 ".

Ex :- 65536x8 : 16 bit address cell holds 65536 ".

Ex :- 131072x8 : 17 bit address cell holds 131072 ".

Ex :- 262144x8 : 18 bit address cell holds 262144 ".

Ex :- 524288x8 : 19 bit address cell holds 524288 ".

Ex :- 1048576x8 : 20 bit address cell holds 1048576 ".

eg. Byte Bit

N° 1000
64 KW

Byte Word

0000

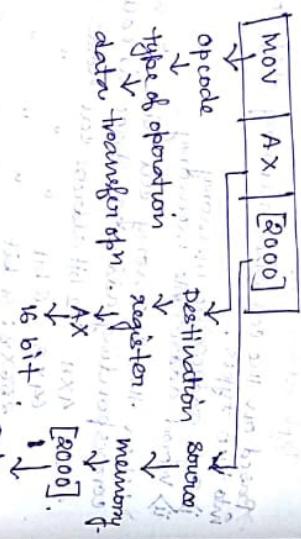
① TB [2000] contain HB when
AX: 1Q23H

② The [0001] contain HB

Note:- (i) Default memory configuration used in the processor design is byte addressable. So data is stored in the memory in a byte wise. The memory data is stored in the low byte.

(ii) In the processors operations are performed on a word format. So when the word length of a processor is > 8 bits then multiple cell accessing is required in parallel to access the data from memory in a word wise.

In 8086 up
Instruction:



→ To handle the above problem, memory address interpretation mechanism (endian -ness mechanism) is introduced. These mechanisms shows the order of one word. These mechanisms shows the order of data storage in the memory. There are 2 kinds.

↑ Little endian

Little Big endian mechanism shows the order of data storage. Little endian mechanism shows the order of data storage. Address contains lower byte & higher addresses contain higher byte.

| | | | | | |
|---|---|---|---|---|---|
| 5 | 4 | 3 | 2 | 1 | 0 |
| 5 | 4 | 3 | 2 | 1 | 0 |

: Address : Byte storage

Little endian technique lower address contains higher byte.

→ In big endian technique lower address constraint
higher byte and vice versa : so address of byte

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

 : Byte storage
Note: : one interpretation mechanism in the instruction

Default address interpretation is little endian so the above processes design is little endian too. The above generates the following op.

$$AX = 1223 H$$

In most registers contain 16 bit data, but memory addressing is byte address type. Hence 2 memory locations will be accessed parallelly. Interfacing in done to access 2 8 bit cell parallelly to get 16 bit word.

→ During the execution of some outcomes are generated due to lack of order of a data storage in the memory i.e

Pins structure :-

Processor consists of a set of hardware pins, need to perform the various operations on the externally connected components i.e. memory & I/O.

Hardware pins are categorized into 3 groups :- active low, active high, time multiplexed pin.

Active low pin :-

This pin is enabled when clock pulse is in low state. Denoted as pin name e.g. RD, WR, INTA etc.

Active high pin :-

This pin is enabled when the clock pulse is in high state. Denoted as pin name e.g. ALE, INTB, HOLD, HLDA etc.

Time multiplexed pin :-

This pin is used to carry the multiple meanings but not at the same time. In the processor design data pins are time multiplexed with address pins so advantage is reduces the number of hardware pins.

e.g. In 8085 AD₀ - AD₇ bits.

In 8086 AD₀ - AD₁₅ bits.

ALE pin is used to differentiate the meaning of a time multiplexed pin.

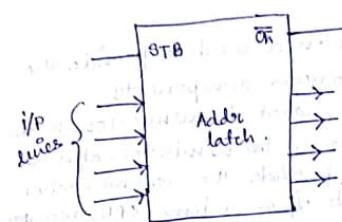
ALE pin is enable during burst or single memory access.

2 time multiplexed pins carries data. 2 time multiplexed pin carries the address.

Address latch :- It is a temporary storage comp used to hold the address. It consists of 2 control pins (STB & \overline{OE}) used to control the input & output lines of the latch respectively.

STB (strobe pin) → 0 : Input lines are disabled
→ 1 : Input lines are enabled.

\overline{OE} (output enable pin) → 0 : o/p lines are enabled
→ 1 : o/p lines are disabled.

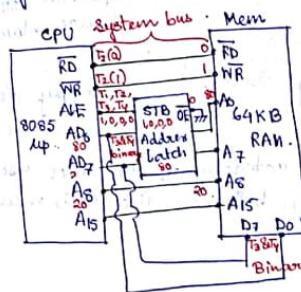


7.4.16

Memory Interfacing

This concept is used to integrate the CPU and memory into one component by providing the mapping between CPU pins & memory pins. Let us consider 8085 processor as a reference CPU to analyse memory interfacing.

It is a 8 bit microprocessor supports 64KB physical memory.



STB line 1
i/P lines of latch
enabled.

Machine instru: LDA[8080]

→ Address of Acc ← M[8080].

CPU generates the mem request

8080 H RD

① Send the address

T₁: ALE = 1

80 → AD₇ - AD₀

80 → A₁₅ - A₈

② Send the 'CS' & Mem

T₂: ALE = 0

RD = 0

WR = 1

③ Read the data

T₃ & T₄: ALE = 0

(Mem) D₇ → D₀ → AD₇ - AD₀.

Delay

T₅ & T₆: ALE = 0

D₇ → D₀ → AD₇ - AD₀

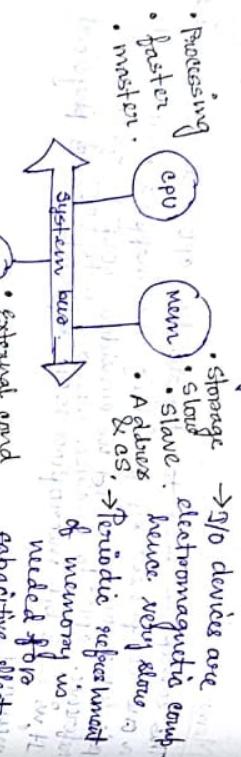
Delay

T₇ & T₈: ALE = 0

D₇ → D₀ → AD₇ - AD₀

BUS

- (i) Bus is a communication channel used to provide the communication between various components.
- characteristic of a bus is shared transmission medium.
- limitation of bus is only one transmission at a time.
- A bus which is used to provide the communication between major components of a system (CPU, memory and I/O) is called as system bus.



To devices are not always slave. Power supply via CS > control signal.

→ System bus consists of 3 category of the lines used to perform the communication between the CPU, memory and I/O, name as

- Address lines
- Data lines
- Control lines

Address lines: These lines are used to carry the address to memory & I/O. So address lines are unidirectional. Based on the address lines, we can determine the capacity of memory system.

Ex:- In 8085 AD₀-AD₇, A₈-A₁₅. Total 16 address lines we have. 2¹⁶ cells of memory = 64K memory.

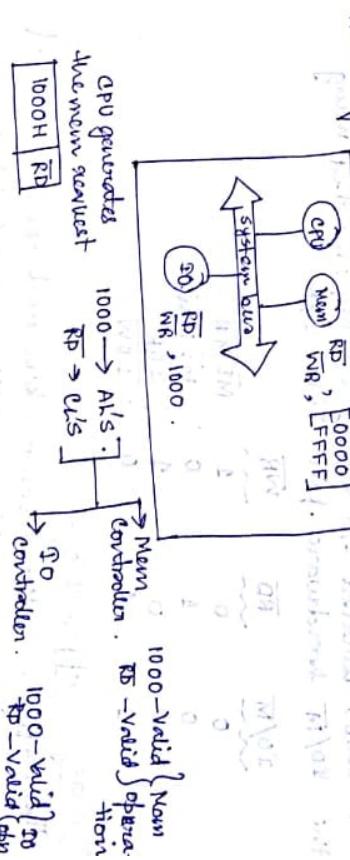
Ex:- In 8088 AD₀-AD₁₅, A₈-A₁₉. 2¹⁶ address lines → 1K × 1K cells → 1M cells → 1MB.

Data lines :- These lines are used to carry the binary sequence between the CPU, memory & I/O. So data lines are bi-directional.

Based on the data lines we can determine the word length of a processor. Based on word length we can determine the performance of a processor. ex:- In 8085 AD₀-AD₇ so word length is 8 bit. So operations are performed on a 8-bit data format. Similarly in 8086 AD₀-AD₁₅ so word length is equal to 16 bit so operations are performed on 16 bit data format.

Control lines :- These lines are used to carry the control signals & timing signals. These are unidirectional control signals indicates the type of operation memory and signals are used to synchronise the memory clock.

Note when there is a common bus, common control signal and common address space is used between memory and I/O. Now there is a possibility of ambiguity.



→ To handle the above problem bus configurations are used in the computer design.

BUS configuration

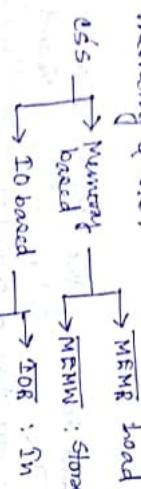
| COMMON BUS | TOP | ISOLATED I/O | M>M<map |
|--------------------------|-----|--------------|---------|
| COMMON as Control signal | X | X | ✓ |
| COMMON Address space | ✓ | ✓ | ✓ |
| | | X | X |

→ There are 3 different bus config used in the system design name as i) TOP (input-output-processors)

ii) Isolated - I/O [isolated so] iii) Memory mapped I/O.

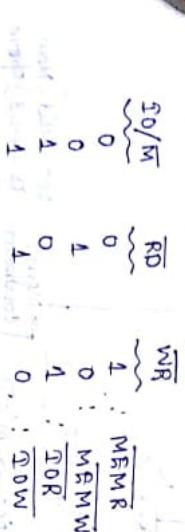
10b - This configuration uses the common control signals and common address space but different busses for both memory and I/O. Additional control mechanism is required to manage I/O bus and memory bus as it is expensive. Therefore it is suitable in the high performance processor design.

Isolated I/O - This configuration uses the common bus and common address space but different control signals for both memory & I/O.



When the processor supports I/O out instructions then we can conclude that isolated I/O config used in the system design. 8085 and 8086 both use the procedures we use the isolated I/O config.

→ In 8085, isolated I/O config was implemented using the I/O/M hardware pin. i.e. $\overline{I/O}$

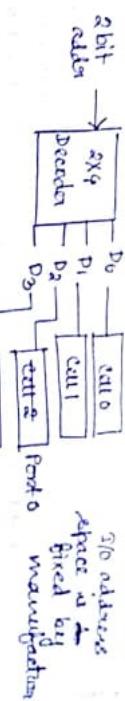
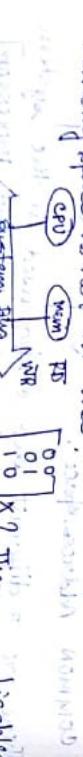


Memory mapped I/O

This config uses the common bus and common control bus but unique address space. Unique address space is maintained by taking some of the address from memory space and assign them to the I/O ports.

Disable the same address in the memory space.

In this config I/O address become the part of a memory address space so limitation is the complete memory space is not in use.



(5) Consider a hypothetical processor which supports 256M word memory. Processor uses the memory mapped I/O configuration. In which when the 3 MSB bits of address is 4 then assign them to I/O ports. How many are I/O addresses and memory address are possible in the processor?

$$\text{# I/O addresses} = 2^{25} \text{ cells} = 2^{28} \text{ cells}$$

Address = $A_0 A_1 A_2 A_3 \dots A_{23}$ → 28 bits.

| | | | | | | | | | | | | | | | | | | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| A_0 | A_1 | A_2 | A_3 | A_4 | A_5 | A_6 | A_7 | A_8 | A_9 | A_{10} | A_{11} | A_{12} | A_{13} | A_{14} | A_{15} | A_{16} | A_{17} | A_{18} | A_{19} | A_{20} | A_{21} | A_{22} | A_{23} |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

$$\# \text{ I/O addresses} = 4 \times 2^{25} = 16 \text{ bits remaining} = 16(600 \dots) 25 \text{ zeros}$$

32 bits → 32 bit → 25 bit.

Memory address = 7×2^{25} .

Total space = I/O space + Mem space

$$= (1 \times 2^{25}) + (7 \times 2^{25})$$

$$= (1+7) \times 2^{25}$$

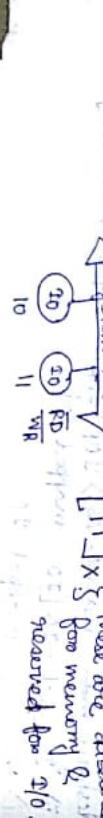
$$= 8 \times 2^{25}$$

$$= 2^8.$$

Instruction cycle

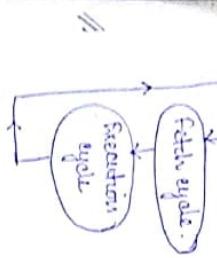
This concept describes the execution sequence of a program. It consists of 2 sub-phases namely no field cycle and execution cycle.

→ Sequence diagram is



Halt instruction

Note → Halt instruction is used as a machine control instruction to indicate the end point of a program.



→ This instruction invokes the unconditional jump with effective address as target address.

PC : 1000 T1 → Halt instruction working

nothing is interpreted in & perspectives, User-point of view.

Halt indicates end of the program execution. CPU point

of view Halt means execute the same instruction upto an infinite number of times.

PC : 1000 T1 → executed

1001 T2 → executed

1002 T3 → execution = JMP 1003

1003 T4 → unconditional jump

1003 T5 → Halt → JMP 1003

1004 T6 → PC : 1004 1003.

fetch cycle.

→ objective of a fetch cycle is instruction fetch, In this process CPU reads the instruction from the memory based on a PC.

→ PC is a mandatory register in the CPU design used to hold the instruction address. So CPU executes the program based on a program counter.

→ Starting address of a program is supplied by the programmer along with an executable command. and the next instruction address is calculated by incrementing the PC with a word length of the processor. Finally at the end of fetch cycle, PC holds the address of next instruction during the execution of

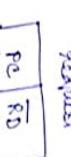
current instruction.

Trick

- + = branching or jumping ; executable statements

t → trace memory
execute pending instruction
store in memory created
from previous created
frame
update
halt

PC → CPU generates memory request



stack address.

PC ← PC + step size.

Di → data line

Ah → Address line

Ch → Control line

(b) Consider a 32-bit hypothetical processor used to execute the following programme.

(a)

. Assume that program size (words) is stored in the memory with a starting address of 1000 decimal onwards.

153:1012-1015 will be the value present in PC during execution of a 16 instruction.

2 1016-1023

1024-1027 from

1028-1031

1032-1035

1036-1039

1040-1043

1044-1047

1048-1051

1052-1055

1056-1059

1060-1063

1064-1067

1068-1071

1072-1075

1076-1079

PC contains 1032.

valid PC.

(b) assume memory is word addressable and the program is stored in memory with starting address of 732 decimal onwards. What could be the value present in PC during execution of T6.

means each cell contain 1 word ie 32 bits.

∴ PC holds T11.

PC

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

- Q) A CPU has 24 bit instructions. Program is loaded in the memory with a starting address of 300 decimal onwards. Which one of the following is a valid program counter?

a) 400 b) 500 c) 600 d) 700

Only 4 memory register address is needed in acc. computation, to manipulate the data.

- execution cycle
- opcode
- instruction format
- layout design

$$24 \text{ bit} = 3 \text{ bytes}$$

Every instruction needs 3 cells space.

Byte Addressable memory by default

| | | |
|-----|---|-----|
| 1 | ° | 300 |
| 1 | ° | 302 |
| 303 | — | 305 |
| 306 | — | 308 |
| 309 | — | 311 |

valid per multiplicative factors of 3.

↓
Transfer of
control op.
Target address
memory address

Execution cycle :-
Objective of an execution cycle is processing of a ready fetched instruction. Operate is required to process the instruction.
→ Operate is a part of instruction specified by the instruction.

→ Instruction format describes the layout design (internal structure of an instruction)

Note:— There is no common instruction format used in the CPU design so let us consider accumulators CPU as a reference model to take the instruction format.

In this CPU, first operand is always sequenced in accumulators & end operand is present either in register or in memory. After the data manipulation result is always placed in accumulators.

→ In the processor design only 1 accumulator is present

so it become a default location so this address is not specified in the instruction.
→ compatible instruction format of all accumulators
CPU via

| | |
|--------------|---------|
| Op Code | Address |
| long word | 5 bits |

register van instruction adres
pe geleverde instructie bevat:
IR geeft data address
hence op code taken

\rightarrow = starting address of : executable statement

The diagram illustrates the flow of memory access:

- PC → CPU generates mem req.**
- PC ← PC + step size.**
- Memory** (represented by a box) receives the address from the PC and returns data to the CPU.
- CPU** (represented by a box) receives the data from memory and stores it in its registers.

② Design time binding

binomial negative binomial

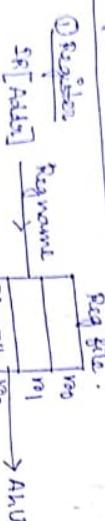
卷之三

EQUATIONS

ACC
ACM
ACM/ACM
ACM/ACM

③ Operand Address (Op)

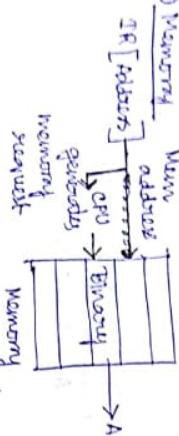
Reg file.



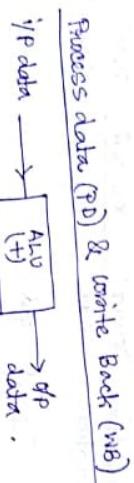
→ In the system design, extra storage is provided to hold the status of an instruction during the execution called as PSW or flag register.

PSW, Addressing mode, Instruction set] → RISC processor

④ Memory Address



⑤



Instruction

⑥ ex:- 8-bit ALU
Instruction : ADD 100

Acc: 10001100
no : 11001110
0110111010

Acc → Acc + no.

extra register
(PSW).

→ CPU reads the instruction from memory based on a PC called as instruction fetch.

→ IR is a mandatory register in the CPU designed to hold the currently fetched instruction to decode because instruction format is predefined in the instruction register.

→ During the decoding process of an instruction the respective hardware is enabled in the system to perform the operation.

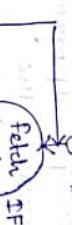
→ Based on the type of a CPU and the associated addressing modes of an instruction, data will be accessed either from register or from memory called as operand fetch.

→ Input data is processed based on the enabled hardware called as process data. When result is placed into a destination called as write back.

Sequence diagram is
① Start.

INT

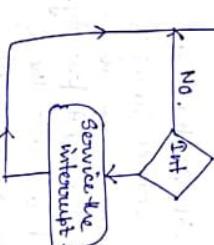
→ sequence diagram states that, CPU will respond to interrupt only after the completion of a current instruction execution.



→ During the program execution, after completion of every instruction, CPU

reads the status of interrupt pins to detect the interrupt.

→ If INT is present → go to the next "instruction fetch" from the main program.



→ When the CPU detects the interrupt, then it saves the user prog context (PC, PSW &

some important registers) into a stack. Then through the control to interrupt vector table (IVT).

→ IVT is a part of the memory where two different interrupt subprograms' information is present.

→ Interrupt subprograms are subprogrammed.

→ Subroutine is performed with vector address & end with RET / RETI instruction.

→ During execution of interrupt subprograms, when the CPU executes the RET instruction, then it invokes POP operation to restore the program context into the registers, therefore, after servicing the interrupt, control will be transferred to a user program.

Note: Objective of an interrupt cycle is interrupt subprogram initialisation. In this process, PC value is saved into a stack, later vector address is loaded into a PC.

→ There are 2 approaches present to service the interrupt.

- (i) Single interrupt processing approach
- (ii) Nested interrupt

In the first approach, CPU disables the further interrupts during the execution of a current interrupt subprogram.

In the second approach CPU will be servicing the

interrupts based on the priority assignment i.e. when the higher priority interrupt is occurred during the execution of current interrupt then it saves current interrupt status into stack and therefore the control

of high priority interrupt.

(b) Consider a hypothetical processor which supports 4 interrupts T_1, T_2, T_3, T_4 with respective service times of 20, 30, 15, 10 us. Response time of an interrupt is 5 us

T_1 has the highest priority & T_4 has the least priority. When in the range of time required to service the T_4 interrupt when the interrupt may or may not occur simultaneously.

Delay between request line and acknowledgement line is response time.

Q. minimum time required from T_4 (without simultaneous occurrence)

Only one INT at a time

The interrupt goes to interrupting processor

↓
Performance + Response time

max time time required for T_4 : (T_1, T_2, T_3, T_4) .
with simultaneous occurrences

5 to 15 ns to 45 ns.

: 15 ns. need to execute the following program segment. Program is

RISC (vs) CISC :-

characteristic of RISC :- (Reduced Instruction Set Comp).

- ① It supports more number of registers.
- ② It has fixed length instructions.
- ③ It has one instruction per cycle ($CPI=1$).
- ④ It has successful pipeline ($CPI=1$).
- ⑤ It has smaller instruction set.
- ⑥ It is a RISC computer.
- ⑦ It is used in real time application.
- ⑧ It is a high expensive processor.
- ⑨ It has dedicated control unit.

example:— Motorola processor
Power PC processor ARM (Advanced RISC machine)

Characteristics of CISC (Complex Instruction Set Comp)

- ① It supports less no of registers.
- ② " " more no of addressing mode.
- ③ Variable length instruction.
- ④ " " no one instruction per cycle (CPI ≠ 1).
- ⑤ " " unsuccessful pipeline (CPI ≠ 1).
- ⑥ " " longer instruction set.
- ⑦ It is a general purpose computer.
- ⑧ " " is used in general purpose applications.
- ⑨ It is a less expensive processor.
- ⑩ It uses the microprogrammed control unit ex:- pentium processor.

Components of Computer system:

Computer system consists of 3 fundamental comp viz

CPU, memory, I/O.

CPU organisation: CPU consists of 3 internal components used to execute the instructions name as

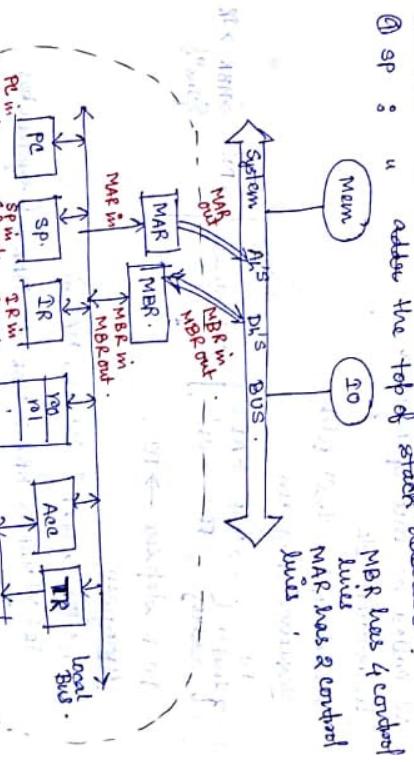
① Register ② ALU ③ Control unit.

Register: It is a temporary storage component used to hold the frequently referred data.

→ In the computer system, program is always present in a memory so frequent communication takes place between CPU & memory during program execution.

→ Memory is a slow component so processor performance degrades if there is no proper balancing. Therefore memory hierarchy design is used as a synchronization mechanism in computer design to minimize the speed gap between CPU & memory.

→ This design organises system & supported memory into 4 levels i.e.



→ In the computer system, in instruction is always present in the memory. During its execution, data becomes most important element therefore keep the smallest storage component (register) in the CPU need to hold the frequently referred data. Therefore every processor contains a minimum set of mandatory registers. They are

- ① PC : holds the instruction address fetched instruction
- ② IR : holds the currently operand & ALU op.
- ③ ACC : holds the 1st ALU operand
- ④ TR : holds the 2nd ALU operand
- ⑤ MAR : " " main address connected to address bus of system bus.
- ⑥ MBR : holds u memory content connected to MBR. data lines of system bus.
- ⑦ GPR : general purpose register, holds the frequently referred data.
- ⑧ PSW : holds status of an instn
- ⑨ SP : " " adder the top of stack address

Micro-operation :-

→ Micro-operation is a elementary operation in the microprocessor.

→ Register to register transfer operation is one kind of micro-operation.
→ Micro operation or average taking 1 cycle to complete the execution (when the microoperation contain memory reference then it takes more than 1 cycle to complete the execution).
→ Control signals are required to execute the micro operation.

Machine Inst: MOV R0,R1

Reg transfer:

$R0 \leftarrow R1$

Hardware design:

Memory

ALU

Control

Register

PC

Memory

| ④ | Medien/Gerüste | ↑ 1 | ↑ 2 | ↑ 3 |
|----------------|----------------|-------|-------|-----|
| T ₁ | S1/S4 | S1/S5 | S1/S3 | |
| T ₂ | S1/S2 | S2/S4 | S2/S4 | |
| T ₃ | S0,S5 | S3,S5 | S3,S4 | |
| T ₄ | S0,S5 | S2,S5 | S2,S4 | |

control signal data base is implemented into a control unit by using the following approaches.

- ↳ handwired approach
- ↳ micro programmed approach.

Handwriting practice unit

kinetics design. Control signal is expressed in a sum of product expression formed. This control expression is directly realised by the independent hardware.

→ It is a fast method used in real time applications
→ It uses model modification techniques re-design and
→ Fewer number of modifications required

→ → RISE another unit as a harmonized vertical unit.

generated a major winner. 4 control signals (S_0, S_1, S_2, S_3) and 3 instructions (I_1, I_2, I_3) to each instruction takes 4 microseconds (T_0, T_1, T_2, T_3)

complete the execution. Following table shows the control signals generated for each operation, from each instruction. Get the control expression for s_0, s_1, s_2, s_3 signals.

| $\frac{S_1}{T_1}$ | $\frac{S_2}{T_2}$ | $\frac{S_3}{T_3}$ |
|-------------------|-------------------|-------------------|
| S_0, S_1, S_2 | S_1, S_2, S_3 | S_2, S_3 |

for Ω_2 w^o even in only $m = T_1$ $T_1 \rightarrow 1$
 but $f_{\Omega_2} T_1, T_3, T_4$. No operation no ($T_2 \Omega_2$) or ($T_3 \Omega_2$). one
 \therefore op of Ω_2 state is only $m = T_1$ map only.

For SO signal \rightarrow output of OR gate is SO signal.
 $T_1 \oplus T_2 \oplus T_3 \oplus T_4$ in all T_1, T_2, T_3, T_4 .
 $T_1 \oplus T_2 \oplus T_3 \oplus T_4 = 1$ so I_1, I_2, I_3, I_4 \rightarrow OR operation is 1 in all
 T_1, T_2, T_3, T_4 .

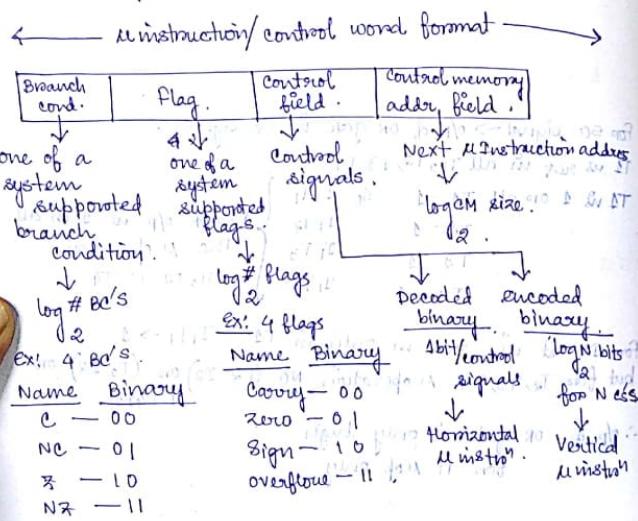
$$S_0 : \pi(\mathfrak{I}_1 + \mathfrak{I}_2) + T_2 \mathfrak{I}_1 + T_3 (\mathfrak{I}_1 + \mathfrak{I}_3) + T_4 (\mathfrak{I}_1 + \mathfrak{I}_3) \\ S_1 : T_2 (\mathfrak{I}_2 + \mathfrak{I}_3) + T_3 (\mathfrak{I}_1 + \mathfrak{I}_2 + \mathfrak{I}_3) + T_4 (\mathfrak{I}_1 + \mathfrak{I}_2)$$

$T_5 = T_3 \oplus T_4$
 $T_4 = T_1 \oplus T_2$

Micro programmed CU :-

Micro programmed CPU :- ROM contains microprogram. User program is available in RAM.

- In this design, control unit is programmed with the micro programs.
 - This control unit contains the control memory used to store the microprograms.
 - Control memory is a permanent memory i.e ROM.
 - In this design micro sequencer unit is present to generate the next micro instruction address.
 - CAR & CDR registers are associated with a control memory used to hold the control memory address & control memory content respectively.
 - In the control memory, micro-instruction is stored in a following format.



SOP, decoded encoded formatted signals \leftarrow control signals.
 Based on the way of representing control signals,
 micro instruction is divided into 2 kinds —
 i) Horizontal microinstruction [Decoded form of CS]
 ii) Vertical micro instruction [Encoded form of CS].

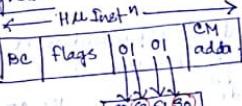
Based on the kind of micro instruction programmed in the control memory, control unit is divided into two types.

- control memory, etc.

 - ④ Horizontal microprogrammed control unit.
 - ⑤ Vertical microprogrammed control unit.

Horizontal programming: $\{ s_0, s_1, s_2, s_3 \}$

Horizontal minstruction
design with 5 so. say
signals.



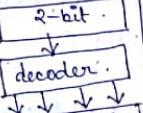
Vertical Programming:

- # CS's in the H/W : {SO, S1, S2, S3} . Vertical 4 Instⁿ design with
 encoded form of a control signal : \log_2 bit.
 i.e. \log_2 bits : \Rightarrow 2 bit. {function codes are used }
 $\frac{1}{2}$ {in the control field to }
 $\frac{1}{2}$ {map memory }

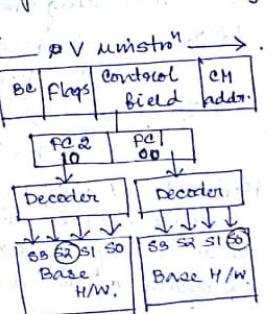
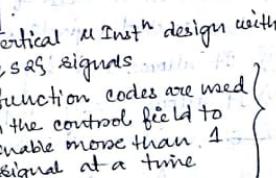
| |
|---------|
| 00 - S3 |
|---------|



Control field



~~s₃ s₂ s₁ s₀~~
Base



Horizontal microprogrammed CU :-

- (i) In this unit, control signal is expressed in a decoded binary format i.e. 1bit/cs.
- Ex:- when the hardware contain 100 control lines then 100 bits are required to represent 100 signals.
- (ii) It supports longer control word.
- (iii) there is no need of external decoders to generate the control signals hence it is faster than vertical.
- (iv) It allows high degree of parallelism i.e. none or more than 1.
- (v) It is bit flexible compared with hardware control unit.

Vertical microprogrammed CU :-

- (i) In this design, control signal is expressed in a encoded binary format i.e. $\log N$ bit format.
- Ex:- when the CPU hardware contain 100 control lines then 7-bit code is required to represent the control signal.
- (ii) It supports shorter control word.
- (iii) There is a need of external decoders to generate the control signals so it is slower than horizontal.
- (iv) It allows low degree of parallelism i.e. none or 1.
- (v) It allows easy implementation of new instructions.
- (vi) so it is more flexible.

Note :- Ascending order of a control unit design

(a) In terms of speed is,

Vertical < Horizontal < hardwired.

(b) In terms of flexibility is

hardwired < horizontal < vertical.

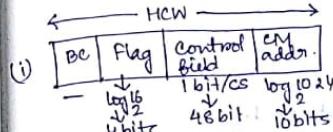
Horizontal microprogrammed control unit is not in the practice because directly working with decoding signal is very complex task therefore default microprogrammed CU is a vertical CU.

3/4/16

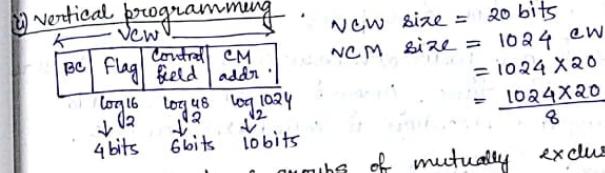
HCW size = 62 bits

HCM size = 1024 CW

$$= \frac{1024 \times 62}{8} \text{ bytes}$$



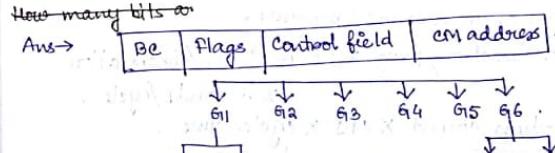
Vertical programming



6. A CU supports 6 groups of mutually exclusive control signals are

| Group | G ₁ | G ₂ | G ₃ | G ₄ | G ₅ | G ₆ |
|-------|----------------|----------------|----------------|----------------|----------------|----------------|
| CS | 1 | 3 | 5 | 8 | 20 | 30 |

How many bits are saved using vertical over horizontal programming.



$$\text{Max bits in VCW} = 1 + 3 + 5 + 8 + 20 + 30 = 67 \text{ bits}$$

(Horizontal)

$$\text{Max bits in VCW} = \log_2 1 + \log_2 3 + \log_2 5 + \log_2 8 + \log_2 20 + \log_2 30$$

$$= 1 + 2 + 3 + 3 + 5 + 5 = 19 \text{ bits}$$

$$\therefore \text{No of bits saved} = 67 - 19 = 48 \text{ bits}$$

(g) A microprogrammed CU supports 256 control signals. Micro instruction is designed such that 3 signals are active. What is the size of control field in micro instruction?

Ans :- default microprogrammed CU = vertical.

$$\# CSs > 256$$

$$\# \text{bits per CS} = \log_2 256 = 8 \text{ bits}$$

Horizontal microprogrammed CU :-

- (i) In this unit, control signal is expressed in a denary binary format i.e. 1bit/cs.
 - ex:- when the hardware contain 100 control lines then 100 bits are required to represent 100 signals.
 - (ii) It supports longer control word.
 - (iii) There is no need of external decoder to generate the control signals because it is faster than vertical.
 - (iv) It allows high degree of parallelism i.e. none or more than 1.
 - (v) It is less flexible compared with hardware control unit.
- Vertical microprogrammed CU :-
- (i) In this design, control signal is expressed in a encoded binary format i.e. log₂ bit format.
 - ex:- when the CPU hardware contain 100 control lines then 7-bit code is required to represent the control signal.
 - (ii) It supports shorter control word.
 - (iii) There is a need of external decoder to generate the control signals so it is slower than horizontal.
 - (iv) It allows less degree of parallelism i.e. never or 1.
 - (v) It allows easy implementation of new instructions so it is more flexible.
- Note :-** Ascending order of a control unit design
- (a) Intensity of speed \Rightarrow
 - (b) Vertical \Rightarrow Horizontal \Rightarrow Hardwired.
 - (c) In terms of flexibility \Rightarrow
- Horizontal \Rightarrow Hardwired \Rightarrow Vertical.
- Horizontal microprogrammed control unit is not in the practice because directly working with decoding signal is very complex task therefore default program is in a vertical CU.

2/1/16

| Horizontal programming | | | | | | new size = 20 bits |
|------------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|---------------------|
| PC | Flag | Control field | CM | CS | CR | new size = 1024 c/w |
| log ₂ 1024 | log ₂ 1024 | log ₂ 1024 | log ₂ 1024 | log ₂ 1024 | log ₂ 1024 | 1024 X 20 bytes |
| 4 bits | 4 bits | 10 bits | 10 bits | 10 bits | 10 bits | |

- b. A CU supports 6 groups of mutually exclusive control signals \Rightarrow
- | Group | G ₁ | G ₂ | G ₃ | G ₄ | G ₅ | G ₆ |
|-------|----------------|----------------|----------------|----------------|----------------|----------------|
| cs | 4 | 3 | 5 | 8 | 20 | 30 |
- Max many bits are saved using vertical over horizontal programming.

| Ans \Rightarrow | PC | Flags | control field | cm address |
|-------------------|----------------|----------------|----------------|----------------|
| | G ₁ | G ₂ | G ₃ | G ₄ |



$$\text{Max bits req} = 1 + 3 + 5 + 8 + 20 + 30 = 67 \text{ bits}$$

(Horizontal)

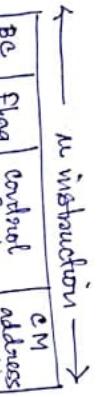
$$\text{bits per cs} = \log_2 1 + \log_2 3 + \log_2 5 + \log_2 8 + \log_2 20 + \log_2 30$$

$$= 1 + 2 + 3 + 3 + 5 + 5 = 19 \text{ bits}$$

$$\therefore \text{No of bits saved} = 67 - 19 = 48 \text{ bits}$$

- (d) A microprogrammed CU supports 256 control signals. Micro instruction is designed such that 3 signals are active. What is the size of control field in micro instruction?
- Ans:- default. A programmed CU \equiv vertical.

$$\# \text{ bits per CS} = \log_2 256 = 8 \text{ bits}$$



Performance evaluation of a processor.

Performance is an indirect measurement. So it is based on execution time. Amount of time required to complete program execution is called execution time.

System X

System Y

Task 1
ET_X = 5 ns
ET_Y = 10 ns.

Execution time is calculated based on processor clock that operated with a const freq.

CPU time = no. of seconds / program
= no. of instruction / program × no. of cycles/instruction

= instruction count × CPS × cycle time

CPU time = $\sum (I_i \times CPT_i) \text{ cycle time}$

→ Program is a combination of data transfers, data manipulation and transfer of control. Instruction will different instr. consumes different cycles to complete execution. So long execution time.

CPU time = $\sum (I_i \times CPT_i) \text{ cycle time}$

→ Speed up (S) factors is used to calculate the performance gain. It is a comparison factor used to compare diff CPU performance level to support the performance gain.

$$S = \frac{\text{Performance}_X}{\text{Performance}_Y} = \frac{\frac{1}{ET_X}}{\frac{1}{ET_Y}} = \frac{ET_Y}{ET_X}$$

$$S = n \rightarrow \text{same const.}$$

→ System X will run n times faster than system Y.

(8) 200 MHz 2.3 GHz processor used to execute the following program sequence.

| Op Type | Inst Count | CPI | |
|------------|------------|-----|---|
| Load | 900 | 11 | a) Avg. Instn execution time of a program |
| Store | 200 | 9 | b) MIPS rate |
| Arithmetic | 300 | 8 | c) Avg. ET |
| Logical | 150 | 6 | |
| Branch | 50 | 4 | |

a) Avg. instn ET = $\frac{\text{Total no. of cycles / program}}{\text{Total no. of instructions / program}}$

$$= \frac{(300 \times 11) + (900 \times 9) + (300 \times 8) + (150 \times 6) + (50 \times 4)}{300 + 200 + 300 + 150 + 50}.$$

$$= \frac{8600}{1000} = 8.6 \text{ cycles.}$$

$$= 8.6 \times \frac{1}{2.3 \text{ ns}} \text{ ns.}$$

(b) MIPS rate
(million instruction per second)

$$1 \text{ instn} \xrightarrow{3.74 \text{ ns}} 1 \text{ sec} \Rightarrow \frac{1 \text{ instn}}{3.74 \text{ ns}}$$

$$\text{With } ? \text{ no. of executing instns} = 3.74 \text{ ns}$$

$$= \frac{10^9}{3.74 \text{ sec}} \text{ instruction/sec}$$

$$1 \text{ million} = 10^6$$

$$= 0.2674 \times 10^9$$

(c) Program execution time

$$= \frac{\text{No. of instns in the program}}{\text{Avg. instn execution time}} = \frac{a}{2}$$

$$= (\text{No. of instns in the program} \times \text{avg. instn ET})$$

$$= (1000 \times 3.74 \text{ ns})$$

$$= 10^3 \times 3.74 \times 10^{-9} \text{ sec} = 3.74 \text{ nsec.}$$

(8) A 3.2 ns clock cycle processor consumes 8 cycles/branch & store instrn, 6 cycles from ALU instrn and 4 cycles from branch instrn. Relative freq of these instrn are 40%, 40%, 80% respectively. Processor is enhanced then the clock cycle time is changed to 3.2 ns with an avg CPS of 4. How much performance

is improved in enhancement?

Execution time of old system

$$FT_{old} = \sum (2e_i \times cpi) \text{ cycle time}$$

$$= [(0.4 \times 1) + (0.4 \times 6) + (0.2 \times 4)] \times 3 \text{ ns.}$$

$$= 14.7 \text{ ns.}$$

$FT_{new} = \sum (2e_i \times cpi) \text{ cycle time}$

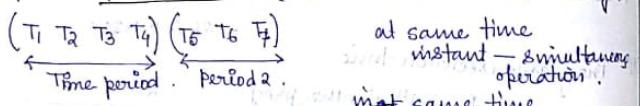
$$= [(0.4 \times 1) + (0.4 \times 1) + (0.8 \times 1)] \times 3.2 \text{ ns.}$$

$$= 3.2 \text{ ns.}$$

$$S_1 = \frac{FT_{old}}{FT_{new}} = \frac{14.7}{3.2} = 4.6 \text{ times faster}$$

New system is running 4.6 times faster than old system
High performance CPU design

High performance processors exploits the concurrency. Concurrency means 2 or more events execution. Event may be a program, or subprogram or instruction or stage of an instruction. Based on the event, the corresponding concurrency is named as program level or subprogram level or inter-instruction level or intra-instruction level of concurrency. Concurrency implies to parallelism or simultaneous or pipelining.



- Parallelism means two or more events execution in the same time period.
- Simultaneous means two or more events execution in the same time instance.
- Pipelining means two or more events execution in an overlapped time span.

According to Flynn's classification, processor architecture is divided into 4 types

1) SISD - Single instruction stream and single data

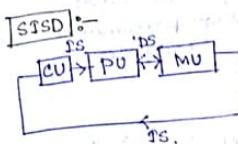
Here relative data given hence adding the values gives the average.

2) SIMD - Single instruction stream and multiple data stream
MISD - Multiple instruction stream and single data stream
MIMD - Multiple instruction stream and multiple data stream

Short-cut

CPU { CU: Control unit
PU: Processing unit (ALU)
MU: Memory unit

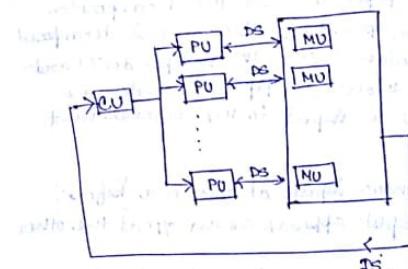
IS: Instruction Stream
DS: Data Stream



In this architecture, instruction & data both are present in the same memory. After the operation result is also present in the same memory called as stored program concept. Von Neumann machine was designed based on the stored prog concept. This architecture is implemented in a uniprocessor system design.

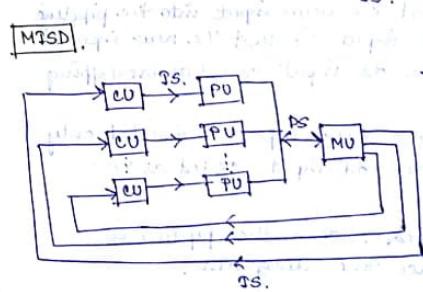
3) SIMD. Single control unit, multiple processing unit.

This architecture is implemented in a vector processor system design.
ex - Stream Processor P.P.P.

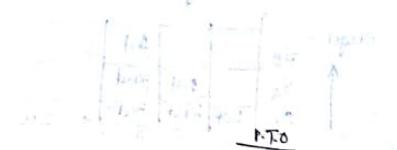


(concurrent operation can be performed. High performance computer) as different PU

$a1 \times b1$ can be $a1_1 \times b1_1$ done simultaneously



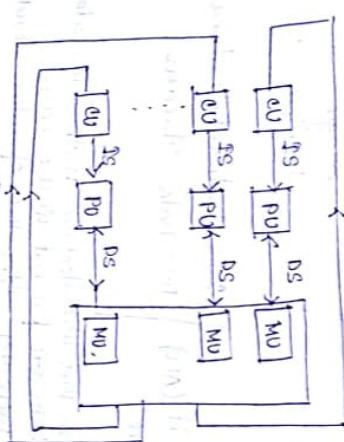
→ In this architecture, physically multiprocessors are present but logically only one processor is in the use therefore this architecture is not yet implemented.



Implemented in a multi processor system design.

or "Cray processor design process".

→ System processing concurrency possible. Large high performance computer.



→ Two performance processors.

Note: Pipelining is an enhancement technique in the processor design used to improve the performance of a uni processor system.

Pipelining

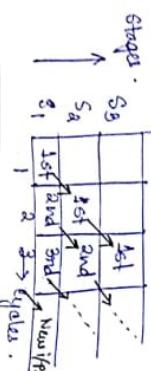
→ Function of a pipeline. → Pipeline was the decomposition technique. In this process, problem statement is decomposed into subproblems and assign them to independent hardwiring. Later connect them in pipeline order i.e first unit of connected as a input to the second unit and so on.

Definition: Accepting the new inputs at one end before the previously accepted input appears as an off at the other end.

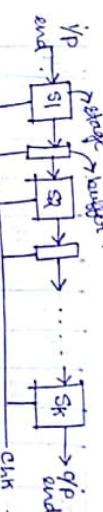
Definition states that insert the new inputs into the pipeline before completion of old input. So that the new input is executing along with our old input called as overlapping execution.

In the new pipeline process new input is inserted only after the completion of our old input called as non overlapping execution.

→ Overlapping execution sequence in the pipeline is described using the space time diagram.

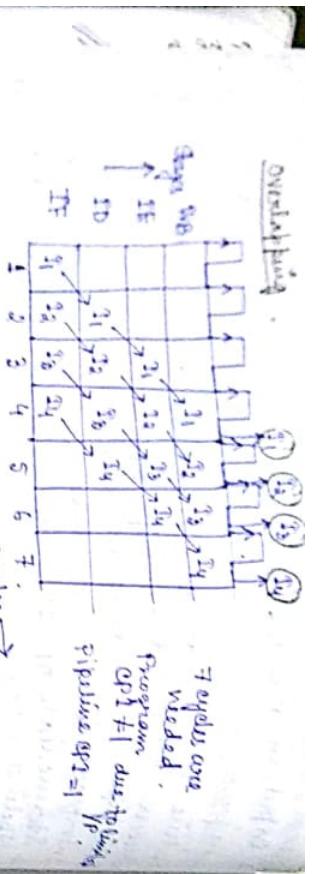


- Successful characteristics of pipelining → in every new cycle new input is inserted into a pipeline therefore $CPI=1$.
- Pipeline has a ends name as input end and output end.
- Between these ends, multiple pipes are interconnected. Each pipe in the pipeline is called as stage or segment.
- Intermediate register is used between the stages to hold the intermediate op. It is also named as buffer or latch or pipeline register.
- All the stages in the pipeline along with a interface register is controlled by the common clock no clock adjustment is very important factor.
- Note: → When all the stages are perfectly balanced (uniform delay) then $\text{cycle time} = \text{stage delay}$
- ② When pipeline stages are not perfectly balanced (non uniform delay) then $\text{cycle time} = \max(\text{stage delay})$
- ③ If buffer delay is included then $\text{cycle time} = \max(\text{stage delay} + \text{buffer delay})$
- Instruction pipeline:
- The pipeline which is used to process the instruction is called as instruction pipeline.
- Instruction execution process is decomposed into 4 sub operations
- ① Instruction fetch (IF) — memory
 - ② Instruction decode (ID) — Decoder unit
 - ③ Execute (EX) — ALU unit
 - ④ Write back (WB) — Register unit
- Sub operations are assigned to an independent hardware when connect them in a pipeline order becomes 4 stage instruction pipeline i.e IF



Consider 4 instructions (I_1, I_2, I_3, I_4) used to execute on a above pipeline. Overlapping execution sequence is

Overlapping

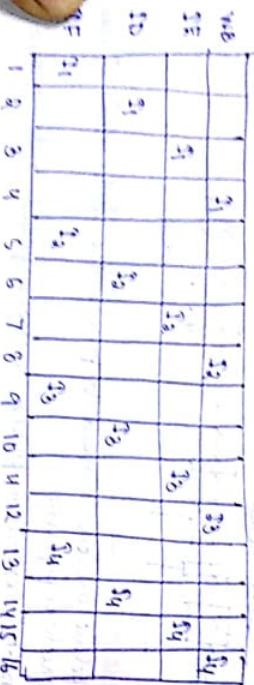


Using 4 clock cycles, 7 cycles are needed for T2. Similarly T3 and T4 are also taking only 4 extra cycles each.

$$\text{extra total cycles} = 4 + (4 - i) = 7$$

In cycles one needs from each into n , and in may case more than $n+n-1$ cycles will be needed.

Non-overlapping



Performance evaluation of pipeline

Consider K sequential pipeline with cycle time t_p used to execute n tasks. The very first task is executed in the pipeline in a non-overlapping method in K cycles.

The remaining $(n-1)$ tasks are emerged from the pipe at the rate of 1 task per cycle, so $(n-1)$ cycles are required to complete the $(n-1)$ tasks. Therefore total time required to complete n tasks in K sequential pipeline is

$$E\Gamma_{\text{pipe}} = ((K+n-1))t_{\text{cycle}} + (K+n-1)t_p$$

\rightarrow Consider non-pipeline processor used to execute n tasks in which each task takes t_n time to complete.

Therefore total time required to complete n tasks in a non-pipeline processor is

$$E\Gamma_{\text{non-pipe}} = n \cdot t_n$$

$$\rightarrow \text{Performance gain of pipeline is } S = \frac{\text{Performance pipe}}{\text{Performance non-pipeline}}$$

$$S = \frac{t_n}{(K+n-1)t_p} = \frac{E\Gamma_{\text{non-pipe}}}{E\Gamma_{\text{pipe}}} \quad S = \frac{n \cdot t_n}{(K+n-1)t_p}$$

\rightarrow Number of pipeline stages increases from n becomes K stages. \rightarrow Using the number of tasks increases from n to $(K+n-1)$ tasks.

$$(K+n-1) \rightarrow n \quad \text{when this condition is true}$$

$$S = \frac{t_n}{t_p} \quad t_n: \text{one task execution time in non-pipeline.} \quad t_p: \text{cycle time of pipeline.}$$

$t_n = t_p$ One task execution time in non-pipeline.

\rightarrow When the pipeline stages are perfectly balanced then one task $E\Gamma$ in the non-pipeline is also equal to the number of stages in the pipeline i.e. $t_n = K \cdot t_p$.

$$\rightarrow S = \frac{K \cdot t_p}{t_p} = K \rightarrow \text{number of stages in pipeline} = \text{pipeline depth.}$$

\rightarrow When the system is operating with 100% efficiency, then maximum build up is possible that is also equal to a pipeline depth.

$$100\% \rightarrow S_{\max} \cdot \frac{n}{m} = \frac{S}{t_p} = \frac{S}{t_{\max}}$$

\rightarrow When the system is operating with 100% efficiency, then maximum build up is possible that is also equal to a pipeline depth.

$$S = \frac{100}{100} \cdot \frac{n}{m} = S$$

$(\text{Throughput}) \text{TP}_{\text{pipe}} = \frac{\text{total time taken}}{40 \text{ pipeline tasks}}$

$$(\text{Throughput}) \text{TP}_{\text{pipe}} = \frac{\text{# tasks processed}}{\text{total time taken}} / \text{to process the task.}$$

(b) consider an instruction pipeline which has a speedup factor 10 while operating with 70% efficiency. What could be the no. of stages in this pipeline?

$$S=10 \quad \eta=70\% \quad \eta = \frac{S}{K} \Rightarrow 70\% = \frac{10}{K} \quad K = \frac{10}{0.7} = \frac{100}{7} = 14.28$$

- (b) Consider a pipeline A & B where A is having 8 stages of uniform delay of 2 ns while B is having 5 stages with a respective delay of 2ns, 6ns, 4ns, 3ns & 1ns. How much time is saved when 100 tasks are pipelined using A instead of B?

$$(E)A = (k+n-1) \times t_p = (8+100-1) \times 2 \text{ ns} = (8+100-1) \times 2 \text{ ns} = (8+100-1) \times 6 \text{ ns} = (107 \times 2) \text{ ns} = 214 \text{ ns}$$

$$(E)B = (k+n-1) \times t_p = (5+100-1) \times 1 \text{ ns} = (5+100-1) \times 1 \text{ ns} = 105 \text{ ns}$$

$$\therefore (E)B - (E)A = (105 - 214) \text{ ns} = 498 \text{ ns}$$

$$= (214 - 214) \text{ ns} = 0 \text{ ns}$$

$$= (214 - 214) \text{ ns} = 0 \text{ ns}$$

$$= (214 - 214) \text{ ns} = 0 \text{ ns}$$

- (c) Consider a 4 stage pipeline with a respective delay of 20 ns, 30 ns, 10 ns, 15 ns. What is the approx speedup and eff. of pipeline when very large no. of tasks are executed?

$$S = \frac{tn}{tp} = \frac{20+30+10+15}{15} = 2.5$$

$$\eta = \frac{S}{k} = \frac{2.5}{4} = 62.5\%$$

If all the stages have same delay then $S=k \rightarrow 100\%$

- (d) Consider 4 stage pipeline with respective delay of 20ns, 40ns, 25ns, 10ns. But overall delay will be 90ns. Interface delay is 10ns. But in non-pipeline we need between the stages have delay of 5ns. What is performance gain of pipeline?

$$S = \frac{tn}{tp}; n \rightarrow \infty$$

$$tp = \max(\text{stage delay} + \text{buffer delay}) = \max(35, 45, 30, 20) = 45 \text{ ns}$$

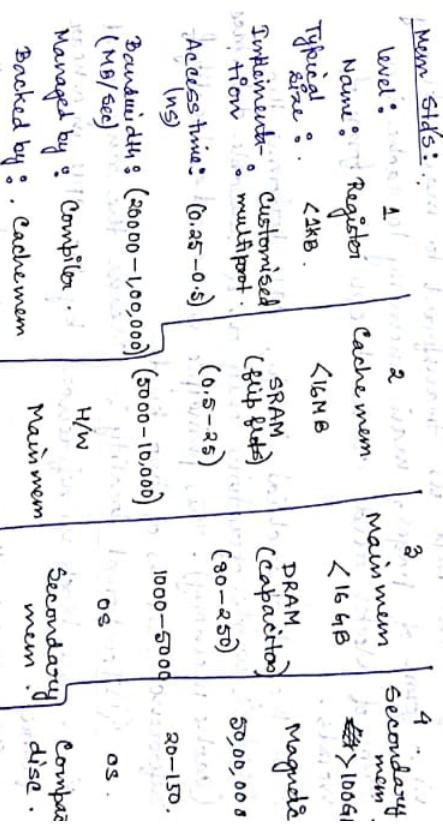
last stage cannot be taken.

$$tp = \frac{tn}{S} = \frac{100}{4} = 25 \text{ ns}$$

total time

$$tp = 75 \times 4 = 300 \text{ ns}$$

With reference to hierarchical design, decreasing sequence of a system supported memories is as follows, so



According to a memory hierarchy design system supported memory standard is as follows.

Memory Organisation: Hierarchy design system according to a memory hierarchy design system supported memory standard is as follows.

Buffers not required in non-pipeline. There is no need to store intermediate off as after completion of one stage the other stage starts. Note: Interface register is not required in non-pipeline because data dependency condition does not occur.

Last stage of doesn't need buffer.

$$tn = (S_1 + S_2 + S_3 + S_4) = 35 + 40 + 25 + 20 = 110 \text{ ns}$$

$$tn = 115 \text{ ns}$$

But the last stage of completion of one stage the other stage starts.

→ CPU generated memory request is initially referred to the cache memory. When the data is in cache, then the operation becomes hit so the hit data is transferred to CPU in a form of words. Otherwise (data is not in the cache), the operation is miss so the reference is forwarded to main memory.

→ When the operation is hit in main memory, then the data is transferred to cache in a form of blocks. So operation is always hit in the secondary mem. Therefore hit data is transferred to main memory in a form of pages, main memory to cache in a form of blocks and cache to CPU in a form of words.

→ In the hierarchical design, CPU performs the read & write operations only on a reusable data store (cache memory). This principle is named as locality of reference. It is of two kinds.

(i) Spatial Locality:— Means the same word in a same block is referenced by the CPU in a near future.

(ii) Temporal Locality:— Means adjacent words in the same block is referenced by the CPU in a sequence like words copied from main mem to cache.

→ In hierarchical design, lower level memory design is always become a subset of higher level memory data called as inclusion.

→ Amount of time needed to access one more data from memory is called as average access time (Tau_{avg}) i.e.

$$\text{Tau}_{\text{avg}} = \text{Hc} \cdot \text{Te} + (1 - \text{Hc}) \cdot \text{Hm} (\text{Tm} + \text{Tr}) + \text{Mc} \cdot \text{Mm} \cdot \frac{\text{Hs}(\text{Tf} + \text{Tr})}{(1 - \text{Hm})}$$

→ Hc → Hit cache
Te → Access time of data from cache.
Mc → Miss in cache = $(1 - \text{Hc})$
Hm → Hit into main mem.
Tm → Access time from data from cache.

Main memory to cache memory organisation:

→ In the hierarchical memory organisation, data is transferred from main memory to cache mem in form of blocks. So both of the memories are organised into blocks based on block size.

$$\text{Number of blocks} = \frac{\text{CM size}}{\text{block size}}$$

$$\text{Number of blocks} = \frac{\text{MM size}}{\text{block size}}$$

$$\text{Number of cache memory blocks} = \frac{\text{Data moved b/w memories}}{\text{block size}}$$

$$= \frac{32 \text{ KB}}{32 \text{ B}} = 1 \text{ K} = 2^{10}$$

$$\text{Number of memory blocks in main memory} = \frac{32 \text{ KB}}{32 \text{ B}} = 2^{27}$$

Cache memory contains 2^{10} block space. Main mem contains 2^{27} information i.e. blocks. Block size is same.

→ Secondary memory to main memory organisation:— In the hierarchical design data is transferred from SM to MM in a form of pages. So both of the memories are organised in form of pages based on page size.

$$\# \text{ of pages (frames)} = \frac{\text{Main mem size}}{\text{page size}}$$

$$\# \text{ of pages in SM} = \frac{\text{SM size}}{\text{page size}}$$

$$\text{Physical address} = \frac{\text{Physical address}}{128} = 2^7$$

$$\text{Physical address} = \frac{\text{Physical address}}{40 \text{ bit logical address}} = 2^{32}$$

$$\rightarrow \text{No. of main MM blocks} = \frac{2^{32}}{2^{27}} = 2^5 \text{ frames}$$

Note :- Replacement algorithms are needed in the memory design used to replace the data when the memory is full. They are

- ① FIFO →
- ② LRU → Least Recently used.
- ③ Optimal replacement.

→ In FIFO, replace the block or page which is having the longest time stamp.

→ In LRU, replace the block or page which is least in the memory longest time without reference.

→ In optimal replacement, replace the block or page which is not referenced in future upto a long time.

This algorithm is not in practice because references are not known in advance.

Types of memory organisation :-

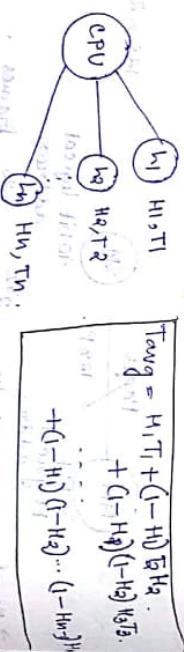
Based on the way of accessing the system supports MM, memory organisation is divided into 2 types

- ▷ Simultaneous access
- ▷ Hierarchical access

MM organ.

In this design CPU is directly connected to all the levels of memories.

Accessing sequence of this design is, when there is a miss in level 1 then CPU directly access the data from level 2 without copying into level 4. If not in level 2 then CPU directly access the data from level 3 without copying into level 2 & level 1 & go on required until level 4, memory no locality of reference is not present.



$$T_{avg} = H_1 T_1 + (1-H_1) H_2 T_2 \\ + (1-H_1)(1-H_2) H_3 T_3 \\ + (1-H_1)(1-H_2) \dots (1-H_4) H_4 T_4$$

$$H_4 = 1$$

$$T_{avg} = 0.7 \times 20 + (0.3) \times 0.8 (20) + (0.3) (0.2) (10) \\ = 14 + 4.8 + 1.2 = 20 \text{ ns}$$

(b) Give a 2 level memory organisation whose access time is 8 times faster than level 2 & its access time is 40ns less than the avg access time. What is the hit ratio?

If 'S' contain hierarchical word then was hierarchical organisation, else was simultaneous organisation. This belongs to simultaneous.

$$H_1 = T_1, H_2 = T_2 \rightarrow T_2 > T_1 \\ S = \frac{T_2}{T_1} \Rightarrow 8 = \frac{T_2}{T_1} \Rightarrow T_2 = 8T_1$$

$$T_1 = T_{avg} - 40 \\ \boxed{T_{avg} = T_1 + 40}$$

$$T_{avg} = H_1 T_1 + (1-H_1) H_2 T_2 \\ = H_1 \times 20 + (1-H_1) \times 1 \times 160 \\ \Rightarrow H_1 =$$

$$\text{Let } T_1 = 20 \text{ ns} \\ T_2 = 160 \text{ ns} \\ T_{avg} = (20+40) \text{ ns} = 60 \text{ ns}$$

(c) 3 level memory organisation has the following specification

| Level | Access time(ns) | Block size in words | Hit ratio |
|-------|-----------------|---------------------|-----------|
| 1 | 20 | 4 | 0.7 |
| 2 | 200 | 8 | 0.8 |
| 3 | 1000 | 16 | 1. |

If the required block is not in L1 then transfer the data from L2 to L1. If not in L2 then transfer the data from L3 to L2 & L2 to L1. How long time does it take to access the data?

Hierarchical.

* access time/block

$$T_1 = 20 \text{ ns}$$

$$T_2 = 200 \text{ ns}$$

$$T_3 = 1000 \text{ ns}$$

$$T_{avg} = H_1 T_1 + (1-H_1) H_2 (T_2 + T_1) + (1-H_1)(1-H_2) H_3 (T_3 + T_2 + T_1) \\ = 0.7 \times 20 + (0.3) \times 0.8 (200) + (0.3) (0.2) (1000)$$

$$= 14 + 4.8 + 12 = 30.8 \text{ ns}$$

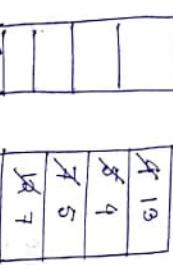
Q) Consider 4 block cache (empty) with the following main memory block references. 4, 5, 7, 12, 4, 5, 13, 4, 5, 7.

Identify the hit ratio using

a) FIFO b) LRU c) optimal replacement.

$$\text{Hit ratio} = \frac{\text{Number of hits}}{\text{Total no. of accesses}}$$

FIFO.



$$H = \frac{8}{10}$$

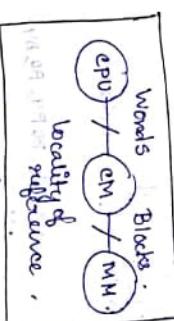
$$H = 0.8$$

Initially cache empty.

first in first out and
first when miss - off
top element will be
replaced then
second and
so on till = 10

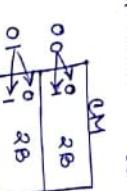
$$H = \frac{8}{10} = 0.8$$

(main memory is slow)
Memory organisation:-
In the hierarchical system, data is
transferred from main memory
to cache memory in a form of blocks
based on the
no. memories are organised into blocks
size.



- ① Memory organisation
- ② Mapping techniques
- ③ Replacement algorithms
- ④ Writing policies

$$\# \text{ lines} = \frac{8}{2} = 4$$



3 bit memory address is coming from B0. 1 bit is set address (line 0). Remaining 2 rows will be stored in cache control reg and will be treated as tag.

No. of sets = 2
Hence for B0: $0 \bmod 2 = 0 \Rightarrow$ means set 0

$B_0 \rightarrow 4 \bmod 2 = 0 \Rightarrow u = 0$

we can store B0 and B1 in a sequence when B2 comes $2 \bmod 2 = 0$ but set 0 is full hence replacement needed. That will be done depending on 3 replacement tech (FIFO/ LRU/ ...)

→ In this cache design explicit replacement tech are used to replace the data when the set is full.

(b) Consider a 64 KB direct mapped cache with a block size of 32 B. CPU generates the 40 bit physical address. How many tag bits are required in the cache line and what is the tag memory size in the cache control?

$$CM size = 64KB$$

$$\text{Block size} = 32B$$

$$\text{address} = 40B$$

mapping tech = direct map

$$\# lines(N) = \frac{64K}{32} = \frac{2^{16}}{2^5} = 2^{11}$$

addresses format of direct map

| Tag | Line offset | Word offset |
|--------|-------------|-------------|
| 10 bit | 10 bit | 10 bit |

| Tag | Line offset | Word offset |
|--------|--|--|
| 24 bit | 24 bit | 24 bit |
| 24 bit | long of block | long of block |
| 24 bit | $\frac{\text{long of block}}{2^5} = 5$ | $\frac{\text{long of block}}{2^5} = 5$ |

∴ 10 bit tag + 10 bit line offset + 5 bit word offset = 24 tag bit.

(b) Repeat the question with 8-way and associative cache design.

Set associative mapping

$$\begin{aligned} \text{Set} &= 1 \\ \text{CM size} &= 64KB \\ \text{Block} &= 32B \\ \text{Add} &= 40bit \end{aligned}$$

$$\text{No. of sets} = \frac{N}{\text{P-way}} = \frac{2^{11}}{8} = 2^8.$$

27 bits in tag.

| Tag | Set offset | WO |
|---------|------------|----------|
| 27 bits | length 8 | length 2 |

$$= 8.8 = 51$$

$$\text{Tag memory size} = \# \text{sets} * \# \text{lines} * \# \text{bits in line}$$

$$\begin{aligned} \# \text{sets} &= 2^8 \\ \# \text{lines} &= 2^4 \\ \# \text{bits in line} &= 2^8 * 8 * 27 \text{ bits} \\ &= 54K \text{ bits.} \end{aligned}$$

(3) Replacement algos → refer the previous section

(4) Updating techniques:

→ CPU performs the read & write operations only on a cache memory or in the hierarchical memory system. Then only the cache memory is updated. The same update is not performed in the main memory so the same cache block is updated. The same update is not performed in the main memory at different places. This kind of inconsistency prob in memory is called cache coherence.

→ Cache coherence causes the data loss. To handle this problem writing policies are used in the memory design. There are two types:

1) Write through:

→ In the write-through protocol CPU performs the main simultaneous write op in both cache mem & main mem. So no coherence prob.

→ In the write-back protocol each cache line maintains one extra bit name as update bit. This bit is set when the CPU updates the cache data. Therefore based on the status of this bit, CPU writes back the data into main memory before replacement. Therefore data is saved.

Secondary storage:

Classification of a secondary storage components is as follows.

Secondary storage

- Direct access storage devices

- Magnetic disks

- Hard disk (high capacity, low cost/bi)

- Floppy disk (low capacity, slow, cheap)

- Optical discs

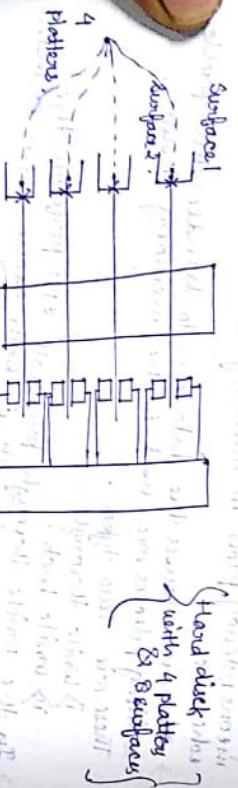
- CD-ROM [compact disc - ROM]

- Serial Access storage devices

- Magnetic tape (high speed sequential access)

Hard disk

Hard disk contain a bunch of magnetic coated platters used to store the data. Each platter contains two surfaces. Every surface contain main track read write head.



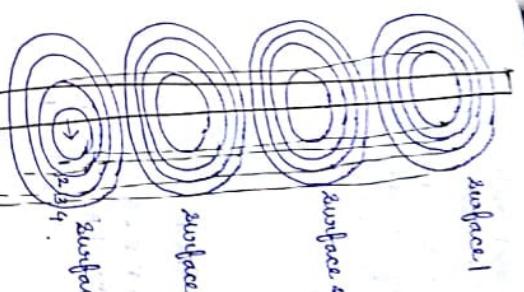
→ Seek time is amount of time required to move the read-write head to the desired track in the surface.
→ Access time of a data from hard disk is calculated using the following formula i.e

$$\text{Total time} = \text{Seek time} + \frac{\text{Average rotational latency}}{\text{Rotational speed}}$$

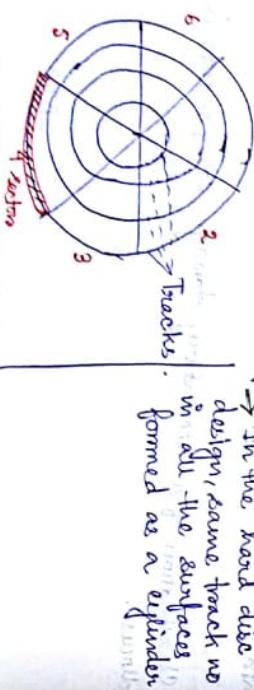
① Disk capacity = # cylinders * cylinder capacity

$$\text{② Track capacity} = \frac{\text{# sectors}}{\text{track}} \times \text{bytes per sector}$$

$$\text{③ Cylinder capacity} = \# \text{ surfaces} \times \text{tracks in the cylinder} \times \text{bytes per sector}$$



→ In the hard disk design, same track no. in all the surfaces formed as a cylinder. Overhead means additional delays in the hard-disk operation.



(Q) Consider a hard-disk with a rotational speed of 7200 rpm. Seek time of a disk is 8.3 ms. Disk contains 8 platters. Each surface contains 32 tracks. Each track contains 8 sectors. Sectors holds the 512 bytes of data. Record size is 4KB. How much time is required to access the second from the hard-disk?

$$\textcircled{1} \# cylinders = 32.$$

$$\textcircled{2} \text{ track capacity} = 64 \times 512 \text{ B.}$$

$$\textcircled{3} \text{ Disk capacity} = 32 \times 16 \times 64 \times 512 \text{ B}$$

$$= 2^{25} \times 2^4 \times 2^6 \times 2^9 \text{ B} = 2^{24} \text{ B}$$

$$= 16 \text{ MB.}$$

Alternative

$$\Rightarrow 16 \times 32 \times 64 \times 512 \text{ B.}$$

$$\Rightarrow 2^4 \times 2^5 \times 2^6 \times 2^9 \text{ B.}$$

$$\Rightarrow 16 \text{ MB.}$$

- $\textcircled{1}$ Seek time = 8.3 ms.
 $\textcircled{2}$ Rotational time: Depends on the rotational speed.

$$7200 \text{ revn} \xrightarrow{\text{1 revn}} 60 \text{ sec}$$

$$4 \text{ secn} \xrightarrow{\text{1 secn}} \frac{60}{7200} \text{ sec.} = \frac{1}{120} \text{ sec.} = 8.33 \text{ ms.}$$

$$\text{Average rotational latency} = \frac{1}{2} \text{ revolution time or}$$

$$= \frac{1}{2} \cdot \frac{60}{7200} \text{ secn} = 0.833 \text{ ms.}$$

- $\textcircled{3}$ Transfer time: Depends on rotational speed.

$$\text{Record size} = 4 \text{ KB} = 4 \times 1024 \text{ B.}$$

$$\text{Rotational speed} = 7200 \text{ rpm.}$$

$$\frac{1}{2} \text{ secn} = \frac{60}{7200} \text{ sec.} = \frac{1}{120} \text{ sec.} = 8.33 \text{ ms.}$$

$$\text{Record size} = \frac{4 \text{ KB}}{512 \text{ B}} = \frac{4 \times 1024}{512} = 8 \text{ sectors}$$

$$\text{Transfer time} = \frac{8 \text{ secn}}{4 \text{ tracks}} = 2 \text{ secn} = 20 \text{ ms.}$$

$$\text{Total time} = 8.33 + 20 = 28.33 \text{ ms.}$$

$$\text{Total time} = 28.33 \text{ ms.}$$

$$\text{Time} \xrightarrow{\text{?}} 8 \text{ sectors}$$

$$\Rightarrow \frac{7.69 \text{ ms}}{64} \times 8$$

$$4.96 \text{ ms} = 1 \text{ track data} = 64 \text{ sectors} = 8 \text{ sectors.}$$

$$4.96 \text{ ms} = \frac{60}{7200} \text{ sec} = 0.00833 \text{ sec.}$$

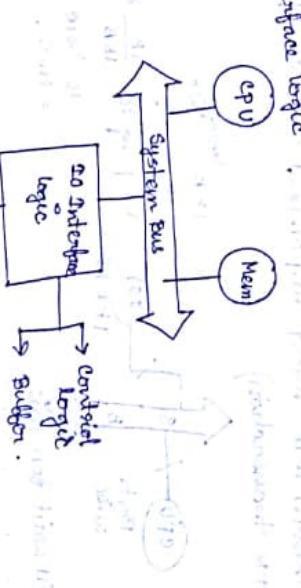
$$4.96 \text{ ms} = \frac{60}{87800} \text{ sec} = 0.00068 \text{ sec.}$$

$$\text{Total time} = 28.33 + 0.00068 = 28.33068 \text{ ms.}$$

$$\text{Total time} = (8.33 + 0.00068) \text{ ms} = 8.33068 \text{ ms.}$$

I/O Organisation

Two devices are electro-magnetic components, a CPU is electronic. So there is a difference in terms of operating modes, comp. So there is a difference in terms of operating modes, word formats, data transfer rate, processor speed, To synchronise the I/O speed with a processor speed, high-speed interface logic (I/O interface or I/O module) is used as a synchronisation technique. I/O interface takes the responsibility of I/O. So processor time is utilised for other useful task. In the system design all the I/O devices are connected to a system bus via I/O interface logic.



→ CPU initializes the I/O interface along with a I/O command → CPU initiates the I/O interface along with other useful task. → I/O interface sends the I/O command to the I/O interface logic. I/O interface logic interacts the I/O command and enables the respective port from the operation. → Based on the speed of a device consumes the time to complete the operation. When status of an operation is transferred to buffer, the status then again I/O interface generates the interrupt signal to CPU. → When the buffer contains the status then CPU waits for acknowledgement value.

→ After receiving the acknowledgement signal the buffer content will be transferred to a CPU. → In this entire process CPU time is wasted by taking the data from the buffer. So buffer delay is very less. Speed of I/O is minimised.

- Based on the speed of a device consumes the time to prepare the data, later enables the DMA R&B signal.
- When the DMA module receives this signal, then it enables the hold signal to gain the control of a system bus and waiting for HDMA signal.
- After receiving the HDMA Signal, DMA module enables the DMA ACK line. When the device receives this signal, then it transfers the data to main memory via DMA module until the count becomes zero. After DMA operation, bus connection is re-established to CPU.
- In the DMA operation, CPU is present in 2 states:
 - Busy state
 - Blocked state / HOLD state.
- CPU is busy until the data is prepared by the device.
- CPU is blocked until the data is transferred to main memory.
- Let x is a preparation time, y is a transfer time.
Then % time CPU is busy = $\frac{x}{x+y} \times 100\%$.
- % time CPU is blocked = $\frac{y}{x+y} \times 100\%$.
- DMA module is operating in 3 modes:
 - i) Burst mode
 - ii) Cycle stealing mode
 - iii) Block mode.

→ In burst mode of DMA, bulk amount of data is transferred to main memory based on the HOLD and HDMA signals.

→ In a cycle stealing mode of DMA, few bytes of data is transferred to main memory by the possible suspension of a CPU operation.

→ In a block mode of DMA, data is transferred to main memory in a block wise based on the HOLD & HDMA signals.