

Qn

Program Counter [PC]: Contains starting address of the next instruction to be executed.

Memory Address Register: Contains memory address for either read/write operations.

DR/MDR/MBR: holds the instructions or data.

Instruction Register: Contains instruction which is currently executed by the CPU.

Accumulator: It contains temporary result of all operations.

GPR: for processing the data.

PSW: stores the status of all result.

Stack pointer: Contains top of stack address.

Steps in Instructions Cycle

1. Instructions address calculation
 2. Instructions fetch
 3. Decoding
 4. Operand address calculation
 5. Operand fetch
 6. Data processing
 7. Result Storage.
- Fetch Cycle Execute Cycle

- * If interrupt occurs, then after completing current instruction execution, the interrupt will be serviced.
- * When interrupt occurs it pushes the PC value into stack as a return address and control transferred to "Interrupt Sub Routine".

System Bus

1. Address line/bus: carry address towards memory and I/O. (Unidirectional)
2. Data line/bus: carry data (bidirectional)
3. Control lines/bus: carry control signals (individually unidirectional and collectively bidirectional).

→ Memory is represented as 2^m . $n = \# \text{ address lines}$
 $m = \# \text{ data lines}$.

- * Address line specifies the capacity of memory.
- * Data line specifies the capacity of data (cell size).
- * n bit address line can represent 2^n memory cells, Range $0 \text{ to } 2^n - 1$
- * Default configuration in memory is byte addressable, and default in the CPU is word.

Pins: Processor contains set of hardware pins to perform operations.

1. Active low pins: Enabled when ip is 0 or clock pulse is low state.
2. Active high pins: Enabled when ip is 1 or clock pulse is high state.
3. Dual pins:
 1 :- then memory operations.
 0 :- then I/O operations.
4. Time multiplexed pins: this pin carries multiple meaning, but one at a time. Address pins are time multiplexed with data pins to carry the data and address (one at a time)

Opcode: Operational code., n bit opcode can perform 2^n operations.



→ N Operations, then

$$\text{Opcode} = \lceil \log_2 N \rceil \text{ bits.}$$

* If memory capacity is M , then address field = $\lceil \log_2 M \rceil$ bits

Elements of Machine Instructions

- * Operations code
- * Source operand reference
- * Result operand reference
- * New Instructions reference

Storage.

→ Register (Fastest)

→ Memory
 → Cache memory
 → Main memory
 → Secondary memory

Faster: Registers \rightarrow Cache \rightarrow Main memory \rightarrow Secondary memory.

Size: Secondary memory \rightarrow Main memory \rightarrow Cache memory \rightarrow Registers.

* Register \rightarrow Reg.-A.F.

Immediate field: n bit.

* Memory \rightarrow Memory-A.F

Unsigned range = 0 to $2^n - 1$.

Signed range = -2^{n-1} to $+ (2^{n-1})$

Instruction Set Architecture

1. Stack organisation (0AF)
2. Single accumulator (1AF)
3. General register organisation.
 - Register memory reference (2AF)
 - Register - Register reference (3AF)

Expand Opcodes Technique

1. Primitive instructions
(lowest bit in opcode)
2. Derived instructions
(higher opcode bit)
3. Further derived instructions
(higher opcode bit)

Expand Opcode Steps

1. Identify the primitive instructions in the Cpu
2. Calculate the total no of possible operations
3. Identify the free opcodes after allocating existing instructions
4. Calculate the no of derived instructions possible by multiplying
 $\Rightarrow \text{Free opcodes} \times 2^{\text{Increment bit in opcode}}$

* Addressing is a technique used to calculate the effective address.

Addressing Modes:

- | | |
|--------------------------|---------------------------|
| 1. Immediate Ans | 7. Base Register Ans |
| 2. Direct/Absolute Ans | 8. Indexed Register Ans |
| 3. Memory indirect Ans | 9. Implied / Implicit Ans |
| 4. Register direct Ans | 10. Auto Decrement Ans |
| 5. Register Indirect Ans | 11. Auto increment Ans |
| 6. PC-relative Ans. | |

- ① Immediate Am: Operand are present in Address field of instruction.
- ② Direct/Absolute Am: Operand are present in memory and Address field of the instructions contain effective address.
- ③ Memory Indirect Am: Operand and effective address present in memory, instructions contains address of effective address.
- ④ Register Direct Am: Am operand are present in registers.
- ⑤ Register Indirect Am: Am operand present in memory and Effective address is stored in register.
- ⑥ Pc relative Am: $\left\{ EA = \text{Current} + \frac{\text{AF (offset)}}{\text{Pc value}} \right\}$
- ⑦ Base register Am: $\left\{ EA = \text{Base register} + \frac{\text{AF (offset)}}{\text{Value}} \right\}$
- ⑧ Index register Am: $\left\{ EA = \text{Index reg val} + \frac{\text{AF (offset)}}{} \right\}$
- ⑨, ⑩ Auto Increment and Decrement Am: Similar to register indirect in which register value increments or decrements.
- ⑪ Implied/Implicit Am: Am Operand are present in the opcode itself.

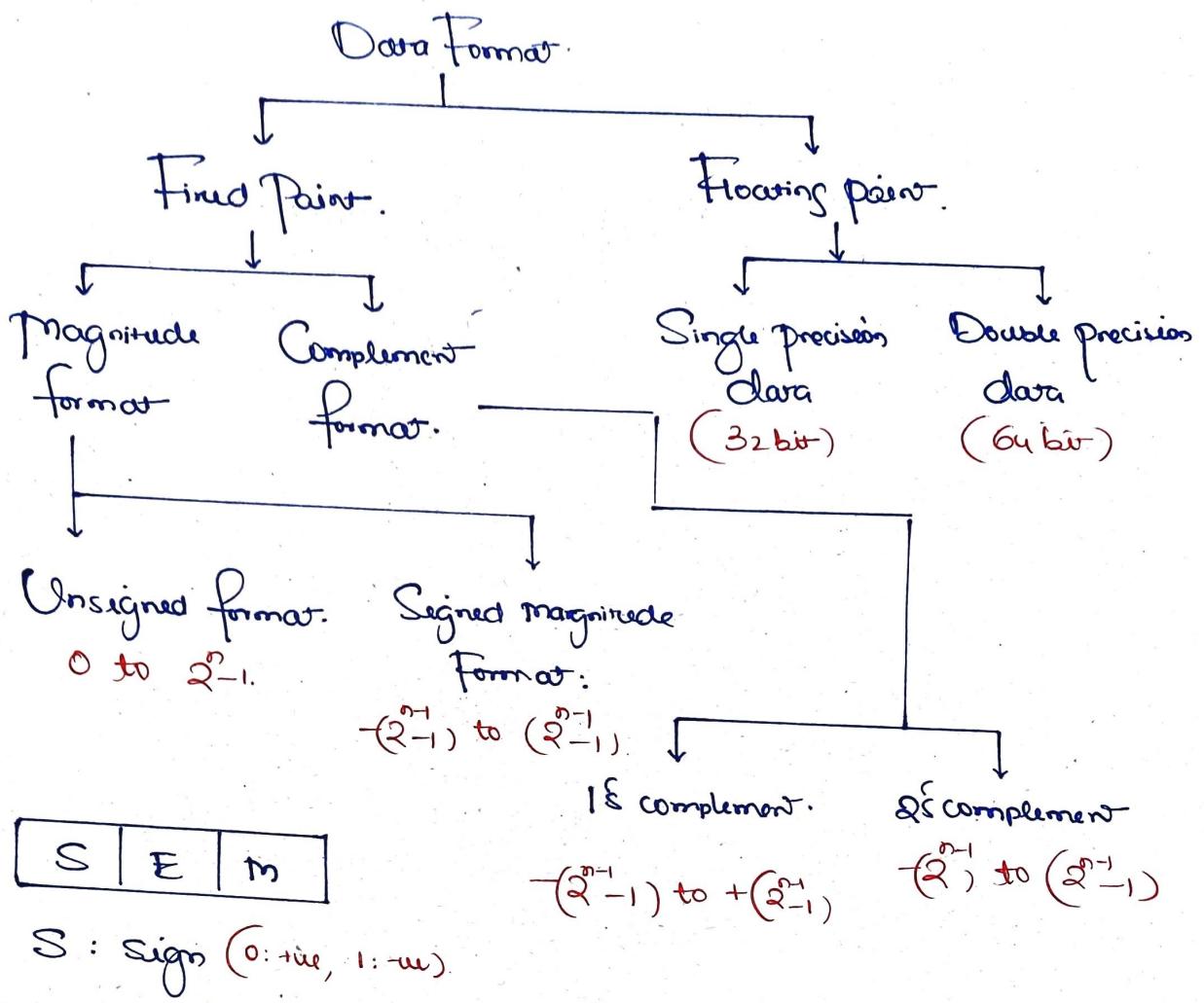
Constant \rightarrow Immediate Am

Variable \rightarrow Direct Am

Pointer \rightarrow Indirect Am

Array \rightarrow Index Am

Reallocations \rightarrow Base Reg Am.



S : Sign (0: +ve, 1: -ve)

E : Exponent $\rightarrow E = e + b_{\text{bias}}$, if exponent is 'k' bit bias = $\underline{\underline{2^{k-1}}}$

m : Mantissa.

Mantissa Alignment.

+ shift.

* 2^+ shift. Power (+) = Right alignment \rightarrow increment the exponent.

* 2^- shift.

Power (-) = Left alignment \rightarrow decrease the exponent.

Explicit Normalised:

$$0.\underbrace{1.m}_{\text{Mantissa}} \times 2^e \Rightarrow (-1)^s \times 0.m \times 2^{E-\text{bias}}$$

Implicit Normalised: $1.\underbrace{m}_{\text{Mantissa}} \times 2^e \Rightarrow (-1)^s \times 1.m \times 2^{E-\text{bias}}$

IEEE 754. Floating Point Representation,

\rightarrow Single precision (32 bit), Ex: 1.27

\rightarrow Double precision (64 bit), Ex: 1.023

S	E	m
1 bit	8 bit	23 bit

S	E	m
1 bit	11 bit	52 bit

Single Precision (32 bits), *double precision is similar with bit changed

Sign (1bit)	E (8bits)	m (23 bits)	Value
0 or 1	$E = 0$	$m = 0$	± 0
0 or 1	11111111	00000...0	$\pm \infty$
0 or 1	$1 \leq E \leq 254$ 254.	$m = \dots$ any no.	Implicit normalized $(-1)^E \times 1.m \times 2^{E-127}$
0 or 1	$E = 0$	$m \neq 0$.	Denormalized number $(-1)^E \times 0.m \times 2^{E-bias}$
0 or 1	$E = 255$	$m \neq 0$	Not a Number

Little Endian: lower address contains lower byte and higher address contains higher byte.

Big Endian: Lower address contains higher byte and higher address contains lower byte.

ALU: "Arithmetic and logic unit", performs arithmetic and logical operations, Comparisons and Conditions Checking.

→ No. of multiplexer = no. of bits in the register.

→ Size of multiplexer = no. of registers.

Micro-program: Sequence of micro Operations to perform some operations in hardware level.

Control Signals can be implemented in the Control Unit by

1. Hardwired CU design: (Faster CU design, eg: RISC)

2. Micro programmed CU design.

Horizontal Programming

- * Control Signals expressed in a decoded format.
- * $1 \text{ bit} \rightarrow 1 \text{ CS}$
- $N \text{ bit} \rightarrow N \text{ CS}$
- * Supports longer control words.
- * No external decoder req - wired.
- * Supports high degree of parallelism

Vertical Programming

- * Control Signals are expressed in encoded format.
- * $n \text{ bit} \rightarrow 2^n \text{ CS.}$
- $N \cdot \text{CS} \rightarrow \lceil \log_2 N \rceil \text{ bits}$
- * Supports shorter control words.
- * External decoder required.
- * Supports low degree of parallelism

Speed : Hardwired \succ Horizontal \succ Vertical

Time Consumption : Vertical \succ Horizontal \succ Hardwired

Program Execution Time is based on:

1. Cycle

2. Cycle time, cycle time of $\frac{1}{\text{Clock freq}}$

Program ET = $I_c \times CPI$

\times Cycle time.

$I_c = \text{instr count}$

$CPI = \text{cycle per instr}$

$$\text{Program ET} = \left[\sum_i (I_{c_i} \times CPI_i) \right] \text{cycle time.}$$

$$CPI = \sum_{i=1}^n (CPI_i \times I_i) / I_c$$

$$\text{MIPS rate} = \frac{I_c}{T \times 10^6} = \frac{f}{CPI \times 10^6} \text{ MIPS.}$$

$$T = I_c \times CPI \times \tau.$$

$I_c = \text{Instruction Count}$

$$\tau = 1/f.$$

Performance gain

$$= \frac{\text{Perf. of new}}{\text{Perf. of old.}}$$

Pipelining: Stalls are executed in overlapping manner.

$$E.T. \text{ pipeline} = [k + (n-1)] \times t_p \quad k = \# \text{ stages (segments)}$$

$n = \text{No. of instructions}$

$$E.T. \text{ non pipeline} = n \times t_n \quad t_n = \text{each instruction execution time}$$

$$\text{Speedup factor } [S] = \frac{E.T. \text{ non pipe}}{E.T. \text{ pipe}}$$

$$S = \frac{t_n}{t_p}$$

- * In ideal case where pipeline stages are perfectly balanced the one task execution in E.T. pipeline = $E.T. \text{ non pipe} = k \times t_p$.

Only when very large no. of tasks executed / or no of tasks not given.

$$\text{Throughput} = \frac{n \times t_n}{[k + (n-1)] t_p} \quad \begin{matrix} \text{(tasks completed} \\ \text{per unit time)} \end{matrix}$$

$$\text{Throughput} = \frac{1}{t_p} \quad \begin{matrix} \text{- when no of tasks} \\ \text{not given} \end{matrix}$$

In uniform delay: 1 Task execution in pipeline = 1 Task E.in non pipeline.

$$t_p = \frac{\text{Stage delay} + \text{buffer delay}}{\text{Pipeline}}$$

In non-uniform delay: 1 Task E_{pipe} ≥ 1 Task E_{non pipe}.

$$t_p = \max_{\text{stage}} (\text{stage delay} + \text{buffer delay})$$

$$\text{Efficiency} = \eta = \frac{S}{K} \quad \begin{matrix} \text{--- speedup} \\ \text{--- no. of stages} \end{matrix}$$

RISC Pipeline Stages

1. Instructions fetch (IF)
2. Instructions decode (ID)
3. Execute (EX)
4. Memory Access (MA)
5. Write Back (WB)

- * Every dependency may/may not create stall.
- * ANTI and OUTPUT dependency never creates stall.
- * Data dependency may/may not create stall.

* If the target address of the branch instruction is available at

k^{th} stage: # stalls from branch instructions = $k-1$.

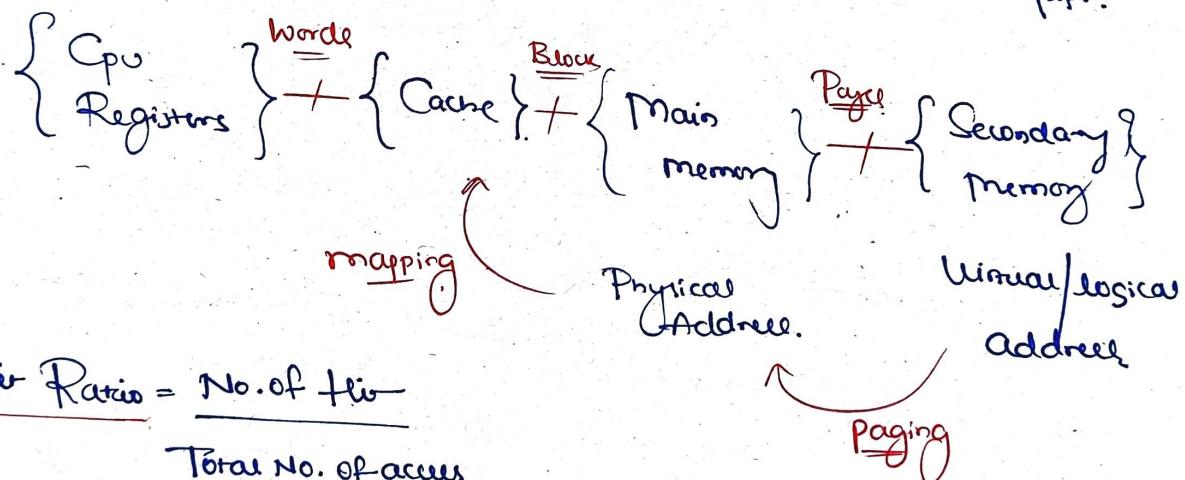
* Average Instructions ET = $1 + \left(\frac{\text{Branch instruction frequency}}{\text{(due to control hazard), CPI} \neq 1} \times \frac{\# \text{ stalls due to branch instructions}}{\text{Instructions}} \right)$

$$\rightarrow \text{CPI} = \left(1 + \frac{\# \text{ stalls/instructions}}{\text{Instructions}} \right)$$

$\rightarrow \# \text{ stalls per instruction} = \frac{\text{Branch instruction frequency}}{\text{Instructions}} \times \text{Branch penalty.}$

$\rightarrow \text{Avg Instructions ET} = \left(1 + \frac{\# \text{ stalls}}{\text{Instructions}} \right) \times \text{Cycle time}$

$\rightarrow \text{Speedup factor using Stalls: } S = \frac{\text{Avg Instructions ET}_{\text{non pipe}}}{\text{Avg Instructions ET}_{\text{pipe}}} = \frac{\text{Avg Instructions ET}_{\text{non pipe}}}{\text{Avg Instructions ET}_{\text{pipe}}}$



$$\underline{\text{Hit Ratio}} = \frac{\text{No. of hit}}{\text{Total No. of access.}}$$

Avg. memory Access Time (Tau)

$$\underline{\text{Tau}} = \text{hit} \times \frac{\text{Time taken by memory during hit}}{\text{hit}} + (1-\text{hit}) \left(\frac{\text{Time taken by memory during miss}}{\text{miss}} \right)$$

$$\underline{\text{Simultaneous Access: }} \text{Tau} = h \times t_c + (1-h) t_m.$$

h : hit ratio

t_c : Cache access time

t_m : main memory access time

$$\rightarrow T_{avg} = b_1 t_1 + (1-b_1) b_2 t_2 + (1-b_1)(1-b_2) b_3 t_3 \dots (1-b_{n-1}) b_n t_n$$

Hierarchical Access.

$$\rightarrow T_{avg} = b_1 t_1 + (1-b_1) b_2 (t_2 + t_1) + (1-b_1)(1-b_2) b_3 (t_3 + t_2 + t_1) + \dots + (1-b_{n-1}) b_n (t_n + t_{n-1} + \dots + t_2 + t_1)$$

$$\rightarrow T_{avg} = b_c t_c + (1-b_c) [t_m + t_c]$$

\rightarrow 1 word access time = T_{avg} .

$$\# \text{words/sec} = \frac{1}{T_{avg}} \quad \begin{cases} \text{data transfer} \\ \text{rate} \end{cases}$$

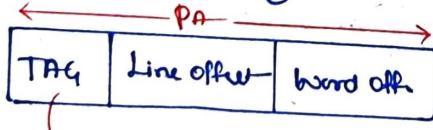
$$\rightarrow \text{No. of mm blocks} = \frac{\text{mm size}}{\text{Block size}}$$

$$\rightarrow \text{No. of cm blocks} = \frac{\text{cm size}}{\text{Block size}}$$

$$\text{Line offset} = \lceil \log_2 \# \text{lines} \rceil$$

$$\text{Word offset} = \lceil \log_2 \text{Block size} \rceil$$

Direct Mapping.



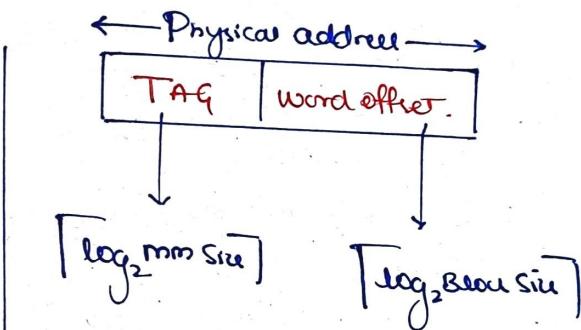
$$\text{TAG} = \frac{\text{mm size}}{\text{cm size}}, \quad \text{Tag memory: } \Rightarrow \# \text{cm lines} \times \text{Tag bits}$$

Mapping function: $k \bmod N = i$

k : mm block no.

N = no of cm lines

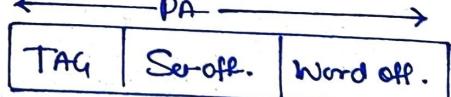
i = cache line no.



$$\text{TAG memory size} = \# \text{lines} \times \text{Tag bits}$$

(depends on mapping technique)

Set Associative Mapping



$$\# \text{sets} = \frac{\# \text{lines}}{N\text{-way.}} \quad \begin{matrix} \text{cm size} \\ \text{Block size} \end{matrix}$$

$$\text{Set offset} = \lceil \log_2 \text{sets} \rceil$$

$$\text{TAG} = \text{PA} - (\text{so} + \text{w.o})$$

$$\text{Mapping: } k \bmod S = i$$

Cause set Cache set no.

- * In direct mapping: $\text{hit latency} = \text{latency of tag comparator}$.
- * In set associative mapping: $\text{hit latency} = \text{latency of tag comparator} + \text{latency of mux.}$

Direct Mapping.

- * No. of comparator = 1.
- * Size of Comparator = # bits in Tag.

Set Association

- * No. of comparator = N-way
- * Size of Comparator = # bits in Tag.

$$\text{Tag bits in N-way} = \text{Tag bits in direct mapping} + \log_2 N \text{ association.}$$

$$\text{Local Miss Rate} = \frac{\# \text{ misses in Cache}}{\# \text{ access to that Cache}}$$

$$\text{Global Miss rate} = \frac{\# \text{ misses in the cache}}{\text{Total } \# \text{ CPU generated reference.}}$$

Access Time in Multilevel Cache

$$\text{Level 1 : } T_{avg} = b_1 t_1 + (1-b_1) (t_m + t_1)$$

or

$$T_{avg} = t_1 + (1-b_1) t_m$$

$$\text{Level 2 : } T_{avg} = b_1 t_1 + (1-b_1) b_2 (t_2 + t_1) + (1-b_1) (1-b_2) (t_m + t_2 + t_1)$$

or

$$T_{avg} = t_1 + (1-b_1) [t_2 + (1-b_2) t_m]$$

WriteThrough

- * Word updating time "tw" = max (word update time in cache, word update time in mm)

$$\rightarrow T_{avg}(\text{read}) = b_{1r} t_c + (1-b_{1r}) (t_m + t_c) \quad \begin{matrix} \text{read now} \\ \text{read data} \\ \text{read now} \end{matrix} \quad \begin{matrix} \text{word update time in cache} \\ \text{read data from Cn.} \end{matrix} \quad \begin{matrix} \text{word update time in mm} \end{matrix}$$

$$\rightarrow T_{avg}(\text{write}) = b_{1w} t_w + (1-b_{1w}) (t_m + t_c) \quad \begin{matrix} \text{Write hit} \\ \text{Word update} \end{matrix} \quad \begin{matrix} \text{Write miss} \\ \text{Word update} \end{matrix} \quad \begin{matrix} \text{Write update} \\ \text{Word update} \end{matrix}$$

$$\rightarrow Q_{\text{av}} = \frac{1}{\tan \theta T}$$

Simultaneous

Acurve : Tangent = $b r t c + (1-tn) t m.$

Write Block Updating Technique

$$\rightarrow \text{Tang(read)} = b \cdot tc + (1 - b) \left[\begin{array}{l} \text{1. clean} \\ \text{modified} \end{array} \right] (tc + tm + tc) + \begin{array}{l} \text{1. clean} \\ \text{biv} \end{array} (tm + tc)$$

$$\rightarrow T_{avg}(\text{write}) = h_{wtrc} + (1-h_{wtr}) \left[\begin{array}{l} \text{if data modified} \\ (t_w + t_m + t_c) + \text{if clean bit} \end{array} \right] \\ \left[(t_m + t_c) \right].$$

$$\rightarrow T_{avg}^{(ws)} = Fr \times T_{avg}^{(read)} + Fw \times T_{avg}^{(write)}$$

↓ of
 read operations ↓
 frequency (↑) of write operations

$$\rightarrow \quad \eta_{WB} = \frac{1}{\text{Targ}_{WB}}$$

Rotational latency: Amount of time taken to rotate the track when the read/write head comes at to exact position.

\Rightarrow Rotational latency = $\frac{1}{2}$ Rotations time

$$\underline{\text{Disk Access Time}} = \underbrace{\text{ST} + \text{Rotational latency} + \text{DTT. + Overhead}}_{\substack{\text{Seek} \\ \text{time}}} + \underbrace{\text{idle transfer time}}_{\substack{\text{if any}}}$$

$$\rightarrow \text{Cylinder no } 'C' = \left[\frac{\text{Outer Sector}}{\# \text{sector/cylinder}} \right]$$

$$\rightarrow \text{Surface no. 'b' } = \left\lfloor \frac{\text{Given Sector} \cdot \# \text{ sectors/cylinder}}{\text{Total Sector per track}} \right\rfloor$$

$$\rightarrow \text{Sector Number} = S + S_r * b + S_t * T_c * C.$$

Sector Sector Surface Cylinder no.
 per track no. no. no.

I-O model

1. Programmed I/O. (CPU time = Speed of I/O device)
2. Interrupt-driven I/O (CPU time = I/O interface time (latency))
3. DMA

* In programmed I/O and Interrupt-driven I/O we cannot access the memory directly, we need help of CPU.

Data Count: The no. of byte/word transfer from I/O to memory until the count value becomes '0'.

$\Rightarrow x$: preparation time, y = transfer time

* % time CPU busy = $(\frac{x}{x+y}) \times 100$

* % time CPU blocked = $(\frac{y}{x+y}) \times 100$.