

COMPUTER ARCHITECTURE & ORGANIZATION - I

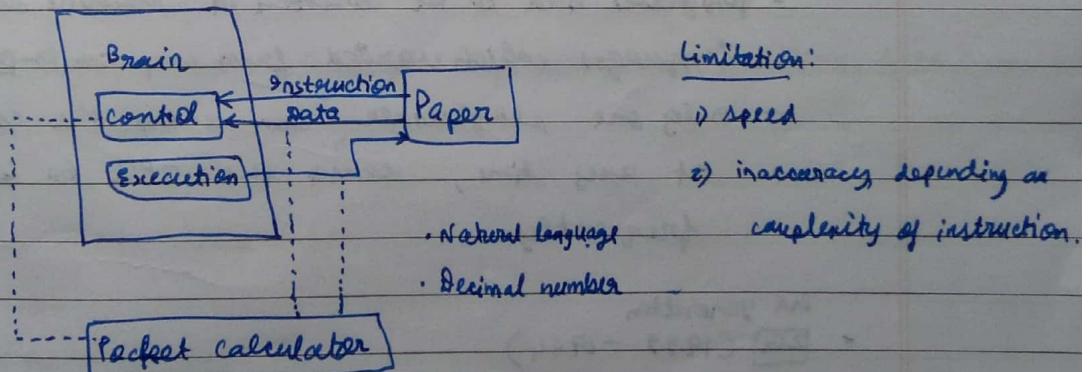
Computer Architecture - interface b/w software and hardware

- Abstract model & programmers view
- in terms of
 - instructions
 - addressing modes
 - registers
- describes what computer does
- deals with high level design issues
 - ↳ for eg. "is there a multiplication instruction?"

Computer Organization - deals with component of connection in a system

- expresses the realization of architecture.
- describes how computer does a task
- organization is done on basis of architecture
- deals with low level design issues
 - ↳ e.g. "is there a multiplication unit or is it done by repeated addition?"

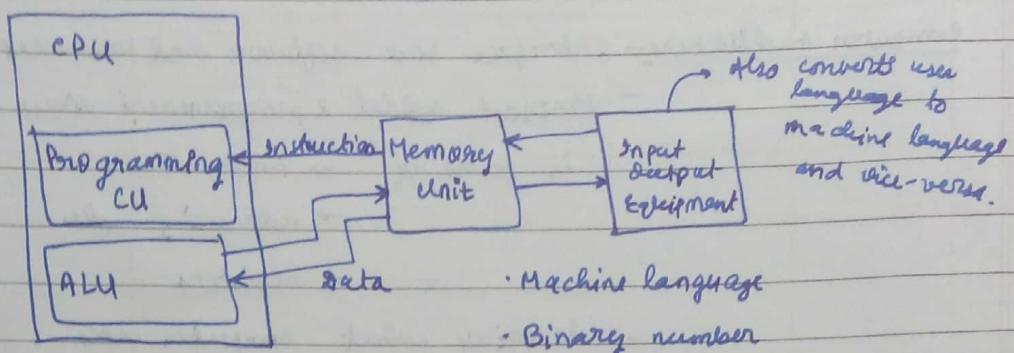
Human computer



- control interprets the instruction and ensures that they are executed in proper sequence.

- Paper contains instructions, data or operand and intermediate results.

Machine Computer



History of computer

1st generation

- **1G** (late 1940's to early 1950's)

- control of computer centralized to single cell.
- all operations in the system, like transfer of a single word of information b/w I/O device and memory, required direct intervention by CPU.

Notes

- very little system software compared to modern machines.
- programs had to be written in machine or assembly languages which varied from computer to computer.
- only one programmer can access the computer at any time, causes the CPU to be idle frequently.

2nd generation

- **2G** (1955 - 1964)

- Transistors which had been invented in 1948 at AT&T Bell laboratories gradually replaced vacuum tubes in the design of switching circuit.
- Cathode ray tube memories and delay line memories were replaced by ferrite core and magnetic drums as the technologies in main memory.
- use of floating point arithmetic hardware.

- machine independent high level programming languages such as FORTRAN, COBOL were introduced to simplify programming.
- special processor (for I/O) was introduced to support I/O operations.
- computer manufacturers began to design ^{system} software.

• **[3G] (1964 - 1971)**

- integrated circuits began to replace the discrete transistor circuits used in 2G, resulting in substantial reduction of cost & size.
- semiconductor memory replaced ferrite cores.
- use micro programming to simplify the design of CPU.
- pipelining and multi-processing were introduced for concurrent or parallel processing to increase the effective speed of execution.
- efficient method for automatic sharing of facilities / resources of the system.
- user interacts through keyboards and monitors interface with an operating system.

• **[4G] (1971 - Present)**

- uses VLSI technology
- compact, reliable, chip
- gave rise to PC revolution
- time-sharing, real time networks, distributed OS introduced
- high level languages like C, C++, etc., were used.
- development of GUI, mouse, handheld device
- great developments in the field of networks.
- concept of internet was introduced
- more powerful computer could be linked together to form networks.

EG (in development)

- multiprocessor based system
- use of AI
- use of optical circuits.
fibres
- Magnetic enabled chip
- Huge development of storage.

Instruction

Binary string

0011 1011 0000 0000 0111
 ← 20 bits →
 = size of instruction in 1G

Add content of memory location
 with accumulator (CPU register)

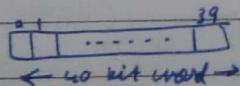
Symbolic programming

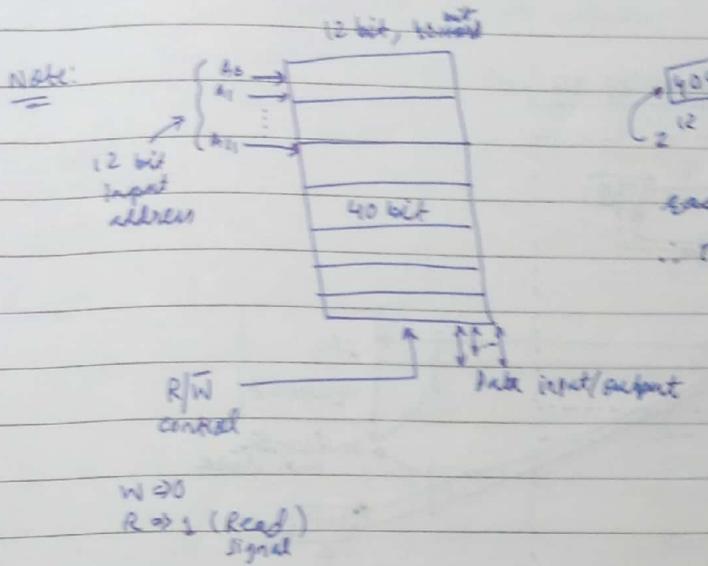
ADD X1

↳ add content of memory location X1
 with accumulator.

1G

- 40 bit word ⇒ it is possible to access 40 bit data at a time from the memory.
- 12 bit address





4096

 2^{12} locations possible in memory

each 12, 40 bit

 \therefore total = $2^{12} \times 40$ bit memory.

40 bit → can be either → data
or → instruction

Consider for the moment → data

Fixed point number:

Put a decimal separator

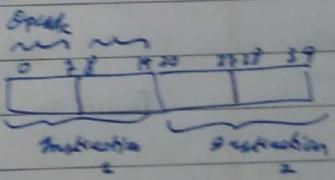
Floating point number:

Binary position has position 0 at by the 40 slots available.

now, consider - instruction

$$40 = (0-19) + (20-39)$$

↓ ↓
instruction instruction



8 bit ← (0-7) = opcode

$(20-27) = 8$ bits

12 bit ← (8-19) = Address

$(27-39) = 12$ bits

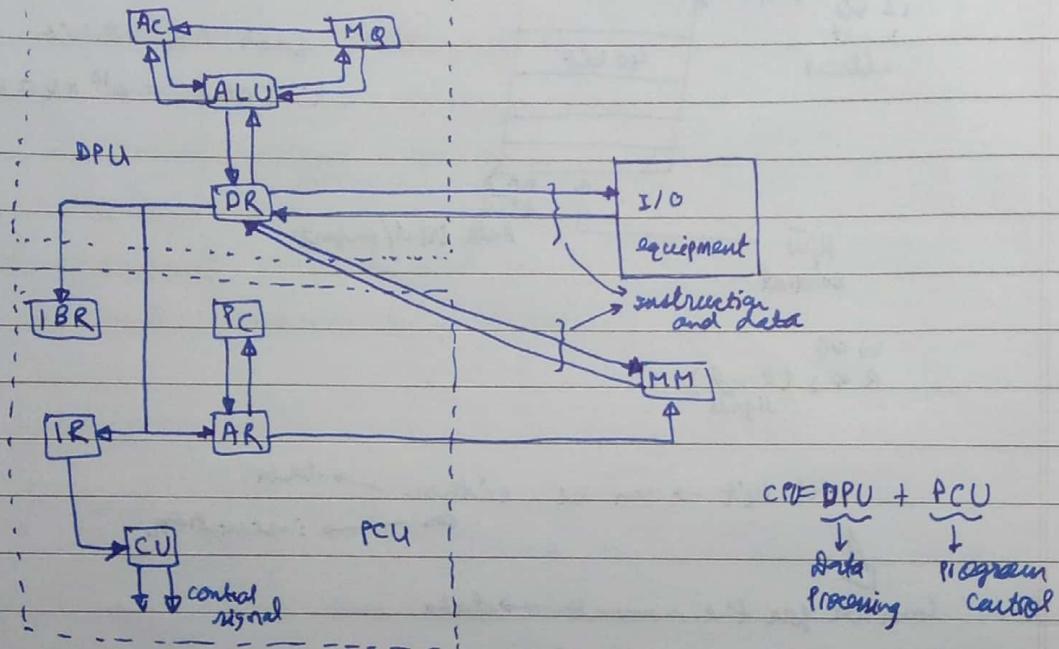
(∴ 19 bit address)

opcode indicates operation

Address indication location of available operand

size of opcode = n bits $\rightarrow 2^n$ instructions can be executed
in various generations by that processor

System Organization of 1G



AC = Accumulator

MQ = Multiplier Quotient

DR = Data Register

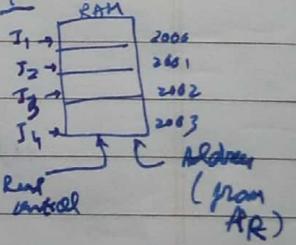
IBR = Instruction Buffer Register

PC = Program Counter

IR = Instruction Register

AR = Address register

programmer has written
program with 4 instructions
→ executions & below



Function of PC: Hold the address of the instruction to be executed.

PC sends address of instruction to AR.

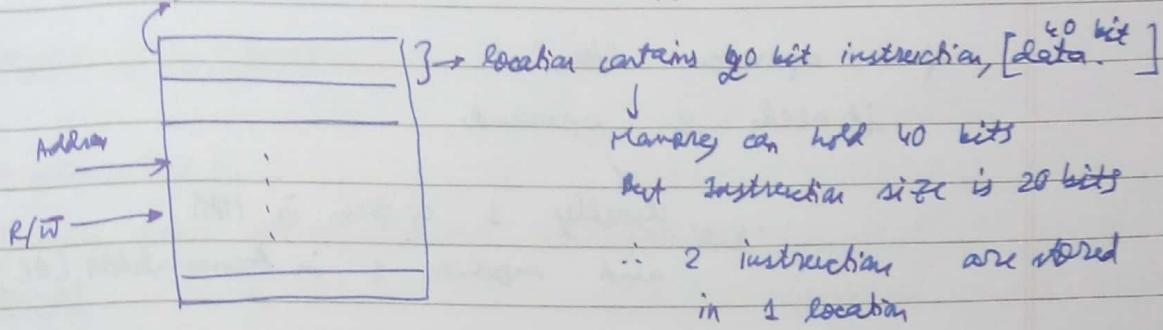
AR sends it to main memory.

CU sends control signal to MM to read instruction, from the address sent to it by the AR.

4. for 1G, 32 bit address for I₁ is 2000.

∴ MM will be signalled to process 2000 address which contains instruction I₁.

each location has 12 bit address.



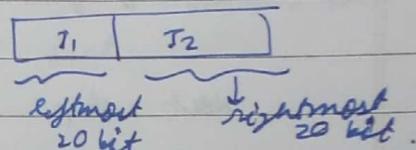
∴ DR size ← 2 instructions come out
is also 40 bit and are stored in the
data register.

But processor can execute only
1 instruction at a time

i.e. 1 instruction enters

Buffer Register (BR)

Let's assume I₂ enters BR.



Now for I₁ → operation to be performed

| | |
|--------|---------|
| opcode | Address |
| 8 bit | 12 bit |

From 8 bit opcode, operation to be performed needs to be understood.
This is called decoding of opcode.

∴ opcode part of the instruction enters IR.
From IR, it goes to CU

IR decodes opcode.

CU generates corresponding signal from this decoded opcode.

From 12 bit address, the address of the operands in main memory
can be known. ∴ the address field of the instruction goes to
AR. From there it goes to MM. Data comes out and enters DR.
∴ Data is now available at DR.

Next step is execution.

If an operator like '+'
it needs 2 operand.

usually 1 of them in MM
and another 1 in Accumulator (AC)

∴ DR sends 1 operand → sent to ALU
AC sends 1 operand ←
CU has already sent control signal from previous cycle

ALU does operation on the operand
to give output result.

Next step is to store result

either in internal register or in main memory
ALU → DR → MM
Basically, ALU sends it to AC

To execute instruction I_2 ,
it is simply retrieved by the process I_2
from IBR, does not need to read from MM
∴ saves time.

MB is required during multiplication / division

Floating Point Numbers

→ why needed?

→ when fixed point numbers not sufficient?

Answer:

- Range of numbers represented by fixed point number is insufficient for many applications particularly scientific computations where very large and very small numbers are frequently encountered.

Scientific notation permits such numbers to be represented using relatively few digits.

e.g. 1.0×10^{18} → scientific notation → floating point

1.000...0 → fixed point
↓
38 zeros

Floating point format:

$$M \times B^E$$

M: Mantissa

B: Base

E: Exponent

→ No. of bits in Mantissa indicate accuracy/precision

e.g. $1010 \cdot \underbrace{11010010}_{3 \text{ bits}} \times 2^2$ [we don't know size of floating point register.
↓
Suppose it can store only 5 bits after decimal
∴ loss of 3 bits of data.]

∴ More bits in Mantissa to improve accuracy

→ exponent field indicates range of the number

e.g. For base 2, $E = -8 \Rightarrow (2^{-8} \rightarrow 2^0)$ is the range.
 $= 2^{55}$ numbers

$$\left. \begin{array}{l} 1.0 \times 10^{18} \\ 0.1 \times 10^{19} \\ 1000000 \times 10^{12} \\ 0.000001 \times 10^{24} \end{array} \right\} \begin{array}{l} \text{same number} \\ \therefore \text{unique representation not} \\ \text{possible.} \end{array}$$

But for computers, we would like a unique representation.

\therefore Normalization is required to specify a unique normal form for floating point numbers in computer implementation.

The Mantissa part is said to be normalized if there is no leading zero in the magnitude part of the number, that means the bit after the binary point must be 1.

In the above number, only 0.1×10^{19} has normalized mantissa.

\therefore It is the unique representation of 1.0×10^{18} .

Clearly, for other 3 forms can be normalized by shifting mantissa towards right or left.

Encoded or Biased exponent

For very small numbers, we need negative exponent.

∴ to represent floating point number, we need signed exponent.

Using 8-bits, list of unsigned exponents is as follows.

| | | |
|------------------|----------------------|----------|
| $255 = 11111111$ | Exponent Bit pattern | 155 |
| $254 = 11111110$ | | 254 |
| $253 = 11111101$ | | 253 |
| \vdots | | \vdots |
| $127 = 01000000$ | | 127 |
| $128 = 00000001$ | | 128 |
| $129 = 00000000$ | | 129 |
| $0 = 00000000$ | | 0 |

$$\text{Biased exponent} = \text{Desired exponent} + \underbrace{k}_{\text{Bias}}$$

$$\text{Bias for 8-bit exponent} = \underbrace{127 \text{ or } 128}$$

↳ consider this by default, if not specified.

| Unsigned exponent | signed exponent | |
|-------------------|------------------|------------------|
| | Bias = 127 | Bias = 128 |
| 255 | 127 | 127 |
| 254 | 127 | 126 |
| 253 | 126 | 125 |
| \vdots | \vdots | \vdots |
| 129 | 2 | 1 |
| 128 | 1 | 0 |
| 127 | 0 | -1 |
| \vdots | \vdots | \vdots |
| $\frac{1}{2}$ | $-\frac{125}{2}$ | $-\frac{126}{2}$ |
| $\frac{1}{4}$ | $-\frac{126}{4}$ | $-\frac{127}{4}$ |
| 0 | $-\frac{127}{1}$ | $-\frac{128}{1}$ |

$$\text{unsigned} = \text{signed exponent} + \text{Bias}$$

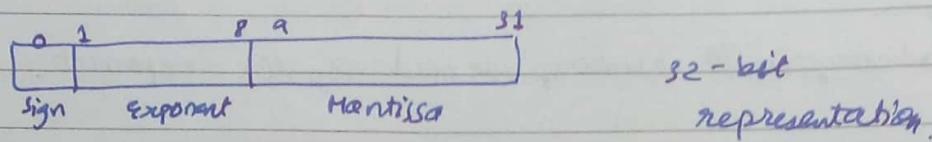
$$\text{Bias} = 127$$

$$(-127 \rightarrow 128) \rightarrow -2^{n-1} + 1 \rightarrow 2^{n-1}$$

where n = no. of bits
in the exponent field

$$(-128 \rightarrow 127) \rightarrow -2^{n-1} \rightarrow 2^{n-1} - 1$$

IEEE 754 standard Floating point Number Format:



- ve number \Rightarrow sign-bit = 1
- Precision is determined by mantissa
- Mantissa is normalized always.
- Mantissa = 23 bits

\therefore Numbers are like

$$0.\underline{1}xyzabc\dots$$

\downarrow always 1

\therefore don't waste space to store it

Do $\rightarrow 1.\underline{xyzabc\dots}$

\downarrow store this in the mantissa
(23 bits)

Initially
adjust the value in exponent value.

$$1.111\dots1 = M_{\text{Max}}$$

$$1.\underbrace{000\dots0}_{23 \text{ bits}} = M_{\text{Min}}$$

$$N = (-1)^S \cdot (2)^{(E-127)} \cdot (1.M)$$

$$0 < E < 255$$

e.g.

| | | |
|---|----------|----------------------------------|
| 1 | 10000001 | 10110100000000000000000000000000 |
|---|----------|----------------------------------|

$$N = (-1)^1 \cdot (2)^{(129-127)} \cdot (1.10110100000000000000000000000000)$$

$$\text{Ex. } N = -1.5$$

Represent it using 32 bit IEEE floating point format.

$$N = (-1.5)_{10} = (-1.1)_2 = (-1.1) \times 2^0$$

$E = \text{desired exponent} + \text{bias}$

$$\therefore E = 0 + 127 = 127 \rightarrow (10000000_0)$$

Max value of $E = 254$ ($\because E \in (0, 255)$)

$\therefore \text{Max number} \rightarrow [0|1111110|111111\dots1]$

$$N = (-1)^0 \cdot (2)^{254-127} \cdot (1.1111\dots1)$$

$$\begin{aligned} 1.11111111 &= 1 + \underbrace{0.1111\dots1}_\text{less num of fp} \\ &\quad \times 2^{-1} \frac{(1+2^{-1})^{23}}{1-2^{-1}} = 1 - 2^{-23} \\ &= 1 + 1 - 2^{-23} \\ &= 2 - 2^{-23} \end{aligned}$$

$$\therefore N_{\max} = (2^{127})(2 - 2^{-23})$$

Min number $\rightarrow [0|00000001|00\dots0]$

$$\begin{aligned} N_{\min} &= (-1)^0 \cdot 2^{(1-127)} \cdot (1.00\dots0) \\ &= 2^{-126} \end{aligned}$$

Similarly, -ve N_{\max} and -ve N_{\min} can be found.

Ex. How is $-6\frac{5}{8}$ represented in 32 bit format

$$-6\frac{5}{8} = -\left(4 + 2 + \frac{1}{2} + \frac{1}{8}\right) = -(110.101)_2 = -(1.10101)_2 \times 2^2$$

$$E = 127 + 2 = 129 = 10000001$$

$$M = \underbrace{10101}_{5 \text{ bits}} \underbrace{0000\dots0}_{23 \text{ bits}}$$

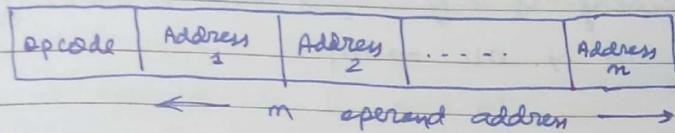
Note: 0 \rightarrow can't be normalized

\rightarrow represented as $0.\underbrace{0\dots0}_{23 \text{ zeros}}$

0×2^E \rightarrow consider only
value for $E \in (0, 255)$

- Result of a floating point operation is not a valid floating point number. Special code referred to as NAN is used
e.g. $\frac{0}{0}$, $\sqrt{-ve}$
- Overflow \rightarrow Magnitude of the number is too big to represent in normal format.
Result is referred to as inf.
- Underflow \rightarrow A result is non-zero, but too small to represent as a normalized number. Such a number is encoded in de-normalized form characterized by $E=0$ and it is represented as $0.M$.

- Instruction format



Processor that allows m operand addresses in the address field of the instruction, is called m -address machine.

Problem: Opcode $\rightarrow 1$ bit

Address $\rightarrow 10$ bits

$$\therefore \text{Total instruction size} = 8 + 10m$$

↓

Which is large \rightarrow lot of memory needed to store instruction.

\therefore we need to reduce the size of instruction

↓

reduce opcode size

- Instead of fixed size,
we use variable size
depending on how
frequently those
instructions are needed.

↓

reduce address part of the instruction

Method 1:

Variable opcode size (variable opcode) in brackets for fixed sized opcode,
in word how needed 3 bits per instruction

| Fixed length | I | Probability of occurrence P_i | 0.5 (1) | 0.5 (1) | 0.5 (1) | 0.5 (1) | → (1) |
|--------------|----------------|------------------------------------|------------|------------|------------|------------|-------|
| 000 | I ₁ | 0.5 (1) | 0.3 (01) | 0.3 (01) | 0.3 (01) | 0.3 (01) | → (0) |
| 001 | I ₂ | 0.3 (01) | 0.08 (000) | 0.12 (001) | 0.12 (001) | 0.12 (001) | → (0) |
| 010 | I ₃ | 0.08 (000) | 0.06 (001) | 0.08 (000) | 0.08 (000) | 0.08 (000) | → (0) |
| 011 | I ₄ | 0.06 (001) | 0.06 (001) | 0.06 (001) | 0.06 (001) | 0.06 (001) | → (0) |
| 100 | I ₅ | 0.06 (001) | | | | | → (0) |

gut regular
Huffman coding.

Average no. of bits per instruction

$$= \sum_{i=1}^5 P_i c_i = 0.5 + 0.6 + \dots = 1.82$$

$$\therefore \text{Savings} = \frac{3 - 1.82}{3} \times 100$$

Advantage of fixed length of code:

- Decoding is easy. just read first bit, bit by bit thereafter.
e.g. 011 \Rightarrow I₄

Disadvantage of variable length opcode.

$$I_1 \rightarrow 1$$

$$I_2 \rightarrow 01$$

$$I_3 \rightarrow 000$$

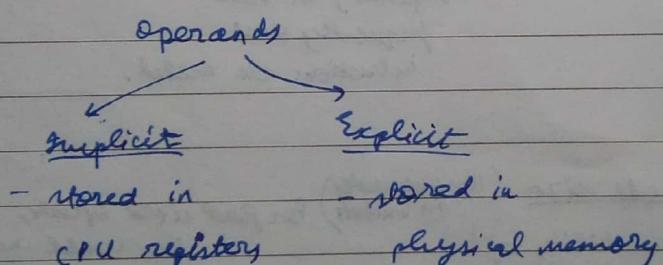
$$I_4 \rightarrow 0011$$

$$I_5 \rightarrow 0010$$

Method 2:

Address has to be specified in the instruction if the operands are in memory.

But if operands are in internal register of CPU, we don't have to specify the address.



e.g. ADD R₁ ^{register}

\Downarrow \Rightarrow 1 operand in AC and 1 in Register R₁

$$AC = AC + R_1$$

Now, if $AC = 03$ (previously)
 $R_1 = 04$ (from mem R₁)

then $AC = 07$ (now)

\therefore instruction contains only opcode. By decoding opcode, processor understands that

Opcode: ADD R₁

[opcode]

Decoded by processor to understand that:

→ Operation = Addition

→ Operands in R₁ and AC

→ Destination of result = AC

Not only size of instruction ↓

also speed ↑

(∴ we are using only registers)

ADD M(X)

[opcode | X]

⇒ AC = AC + content of memory location
where address is X.

3 address M/C

MULTIPLY T, A, B

(T = A × B)

2 address M/C

MULTIPLY T, B

(T = T × B)

1 address M/C

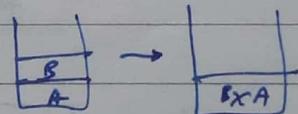
MULTIPLY B

(AC = AC × B)

0 - address M/C

MULTIPLY

Stack
↓



Addressing Modes

| | |
|--------|---|
| opcode | x |
|--------|---|

To execute instruction, processor requires the current value of x, which can be specified in several ways, known as addressing modes.

→ Immediate Mode → Operand itself specified in address field of instruction.

| | |
|--------|----|
| opcode | 07 |
|--------|----|

e.g. ADD x is constant
Value of x is 07.

ADD 07 in immediate mode

would mean ACC \leftarrow AC + 07
(Accumulator)

- 1. instruction is in memory
- 2. memory access by processor
- 1 to read op code
- 1 to read operand

→ Direct/Absolute Mode

| | |
|--------|---------|
| opcode | Address |
|--------|---------|

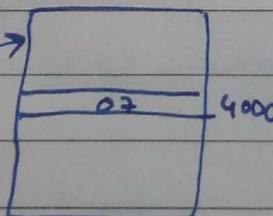
Address of the operand

e.g.

| | |
|--------|------|
| opcode | 4000 |
|--------|------|

4000 applied to address bus

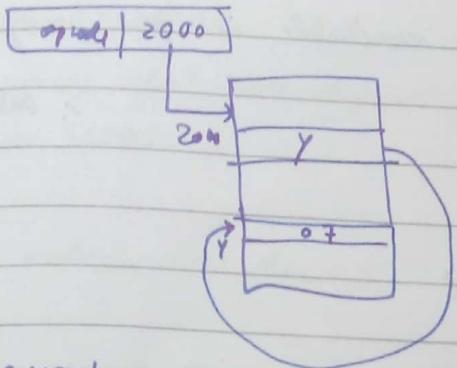
07 came out of data bus



3. memory access by processor:

- 1 to read op code
- 1 to read operand
- 1 to read operand

→ Indirect Addressing mode



Address field contains a pointer
that points to a memory location
where the operand address is available.

In memory access is made
by processor.

Advantages:

Y can be easily changed to Z
for using other operands.

→ Register Addressing Mode

ADD B

$\Rightarrow AC \leftarrow AC + B$

B: Register inside the processor.

∴ It is used, when operand is
available inside register

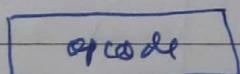
Processor understands

that operand is in register B

and not register C

∴ opcode → ADD C } are different.
 + opcode → ADD B

Instruction format:



• Only 0 memory access.

∴ Very less time needed.

∴ Size of instruction = size of opcode.

→ Register indirect

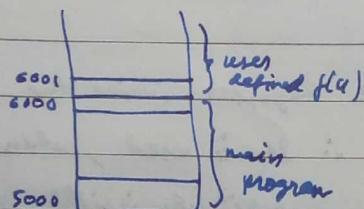
Address of operand is available in the register.

∴ only 1 memory access would be needed

→ Relative Addressing Mode

In the address field, we specify the displacement of the operand, not the address of the operand.

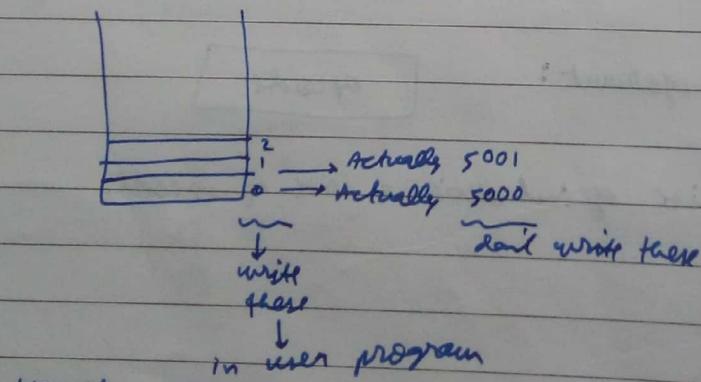
- what is displacement of operand?



After modification of main program,

it may not fit within 5000 - 6000)

∴ Take \sim Register $R = 5000$ ∵ user-defined function may be overwritten.



$$5000 + 0 = 5000$$

$$5000 + 1 = 5001$$

Base address Displacement Effective Address Applied to address base of memory to get data.

To avoid the problem of over-writing user-defined function memory,

Take $R = 8000$

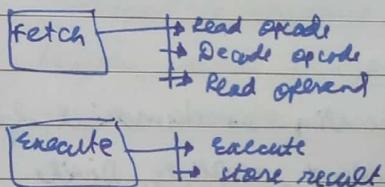
assuming that user-defined function doesn't cross 8000 starting from 6000.

\therefore problem resolved.

To represent 5000 $\rightarrow 2^{13}$ (bits) needed \therefore 16-bit register

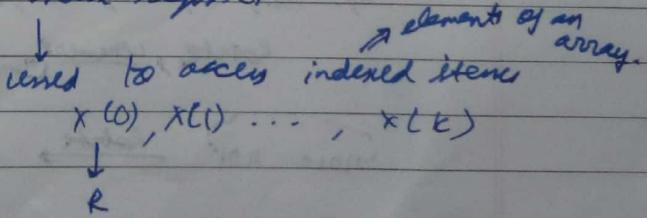
Displacement = distance from base address.

ALU is the one which
calculate the (+) operation in
base (+) Disp. = effective



In some processors,
fetch and execute operations go on simultaneously.
Then, if they have only 1 ALU,
 \rightarrow above problem happen.

R = Base register / PC / index register



→ Auto increment / decrement addressing mode

Indexed items are frequently addressed sequentially so that a reference to $x(k)$ at memory location A is immediately followed by a reference to $x(k+1)$ or $x(k-1)$ at location $A+1$ or $A-1$ respectively.

| | |
|-----|----------|
| A+1 | $x(k+1)$ |
| A | $x(k)$ |
| A-1 | $x(k-1)$ |

∴ a sequence of items can be accessed by auto incrementing or decrementing an address.

Instruction implementation

Instructions are divided into 5 major groups:

(1) Arithmetic instruction

- performs operations on Numerical data

e.g. Add, Subtract, Multiply, Divide,

Absolute, Negate, Increment, Decrement

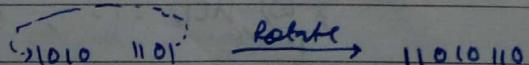
↓ ↓
replace change sign
operand by of operand
if absolute
value

(2) Logical instruction

- Boolean/ non-numerical operation

e.g. AND, OR, NOR, XOR, XNOR, shift,

rotate, convert → change the format.



(3) Data transfer instruction

- from source to destination

e.g. Move, Store, Load, Exchange, Clear, Set

swap source
 & destination
 content

Transfer all
 0's to
 destination

Source - Memory location } Load
 Destination - Register

Source - Register } Store
 Destination - Memory location

(4) Program control instruction

- If not used, execution of instructions in program is always in sequence by default

- Program control instruction changes sequence of program execution.

e.g. Jump, jump conditional, jump to subroutine, Return, execute skip, Test, compare

2001 : I_1

2002 : I_2

2003 : JUMP X $\rightarrow I_1, I_2$ (Go to instruction at address X) I_4, I_5
 end execute it

2004 : I_4

2005 : I_5

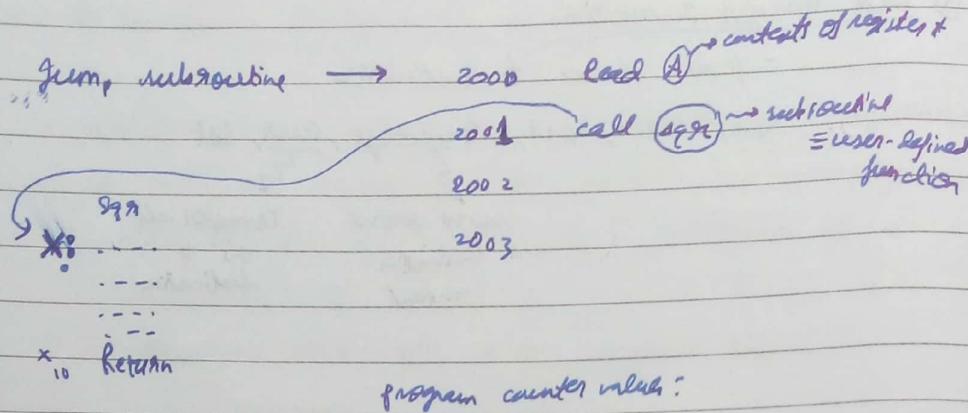
X : ~

Jump conditional \rightarrow if ($R_1 > 0$)
 execute I_2

else

execute I_3

Program → contains address of counter instruction to be executed.



2000
 2001
 2002 → 2002 calls the subtracting
 X₀ actually.
 2002 stored in stack.
 X₁₀ control goes to X
 2002 control returned
 when "Return" to
 2002 was popped from
 stack/
 stored in
 program counter.

Skip → increment program counter
to skip next instruction

Test → specified condition based on outcome.

(5) I/O instruction

- cause information to be transferred between processor or main memory and external I/O devices

e.g. Input, Output, Start I/O, Test I/O, Host I/O.

for 8085 Microprocessor

Input → read data from I/O device

IN < Port No > → address of input device
 IN F4H → hex value

read input data and load into accumulator

OUT F4H
 To data from accumulator
 goes to I/O

Start IO → CPU reads instruction
 Transfer it to I/O processor (see History of computers)
 to initiate I/O.

Test IO → transfer status information operation from
 I/O system to specified destination.

Halt IO → to terminate I/O operation.

* Arithmetic instruction implementation

→ Hardware requirement in CPU to execute such instructions.
 Algorithms used for implementation

→ Addition and subtraction with signed magnitude data

Addition:

Operation

Add magnitude

subtract magnitude

$A > B$ $A < B$ $A = B$

result = 0
 $0 = +0$
 $0 = -0$

$$(+A) + (+B) \quad + (A+B)$$

$$+ (A-B) \quad - (B-A) \quad + (A-B)$$

$$(+A) + (-B)$$

$$- (A-B) \quad + (B-A) \quad + (A-B)$$

$$(-A) + (+B)$$

$$- (A+B)$$

$$- (A-B) \quad + (B-A) \quad + (A-B)$$

Subtraction:

$$(+A) - (+B)$$

$$+ (A-B) \quad - (B-A) \quad + (A-B)$$

$$(+A) - (-B)$$

$$+ (A+B)$$

$$(-A) - (+B)$$

$$- (A+B)$$

$$- (A-B) \quad + (B-A) \quad + (A-B)$$

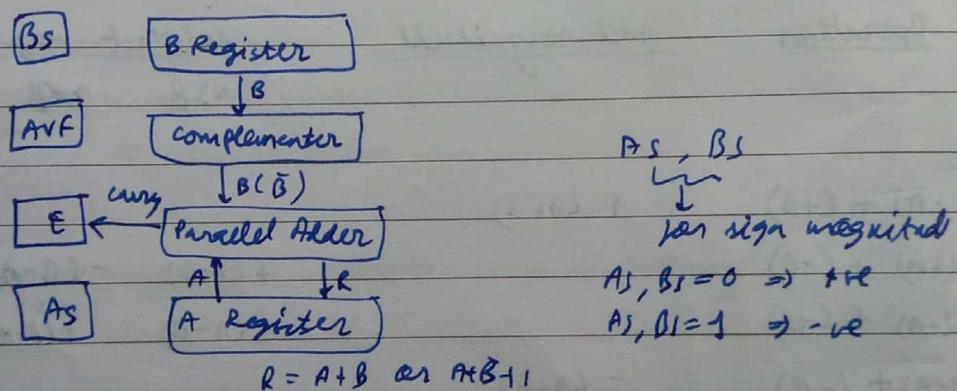
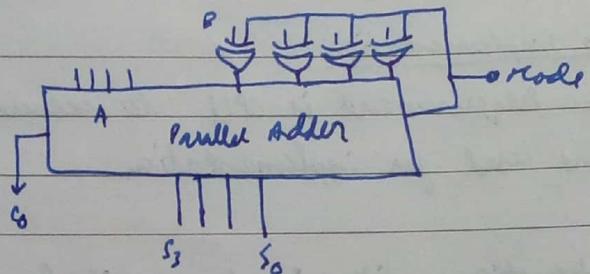
$$(-A) - (-B)$$

A + B

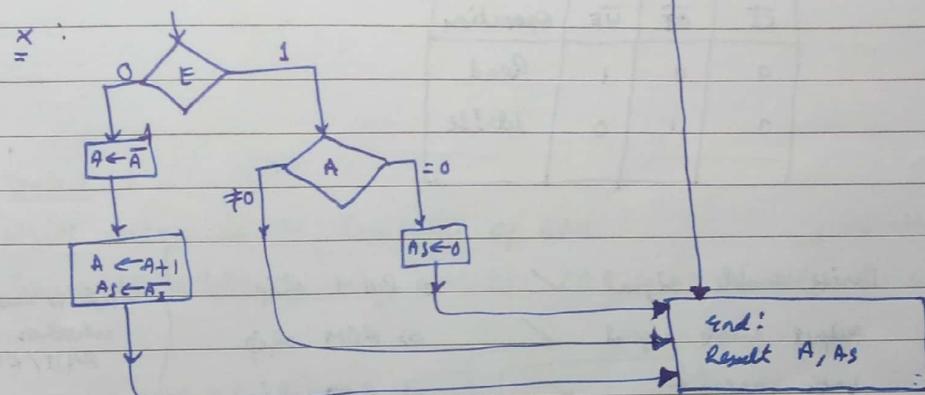
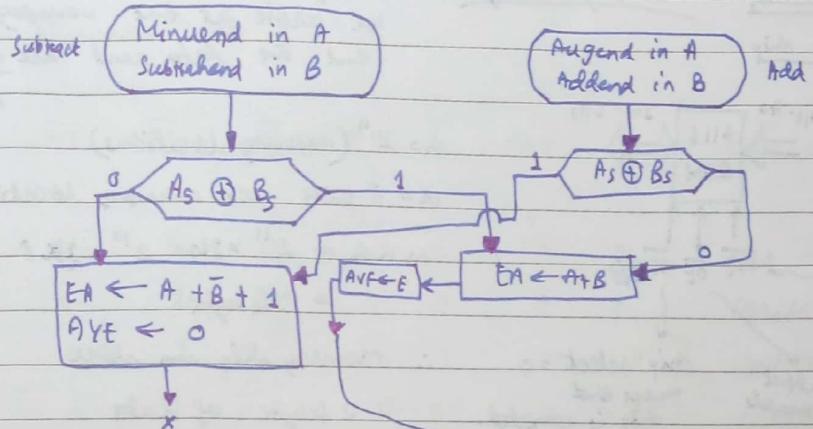
- Two register \rightarrow initially stores A, B; finally one stores result.
- One adder
- One FF (to store carry)

$$\begin{array}{r} A - B \\ \hline \hookrightarrow A - B = A + \bar{B} + 1 \end{array}$$

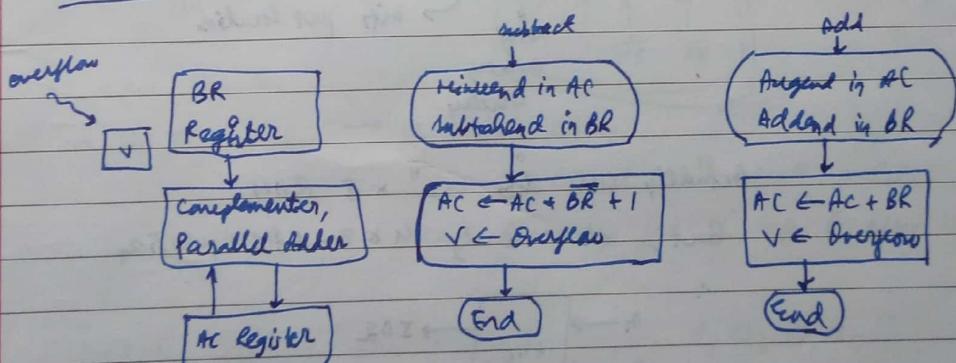
↑ 2's complement of B
↓ 1's complement of B



AVF: overflow (\Rightarrow carry bit generated)
flag



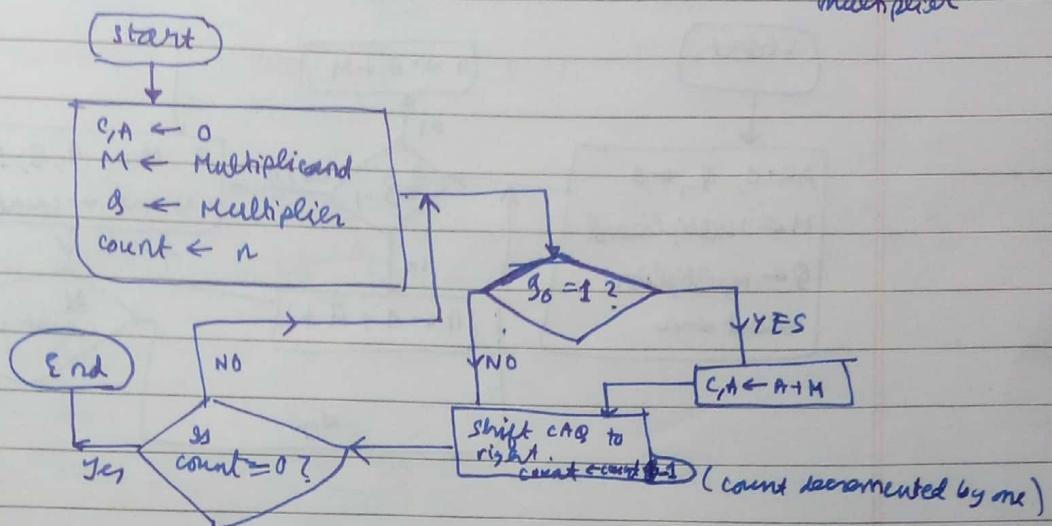
* Addition and subtraction using 2's complement



MULTIPLICATION

Unsigned Integer

count: No. of bits in multiplier

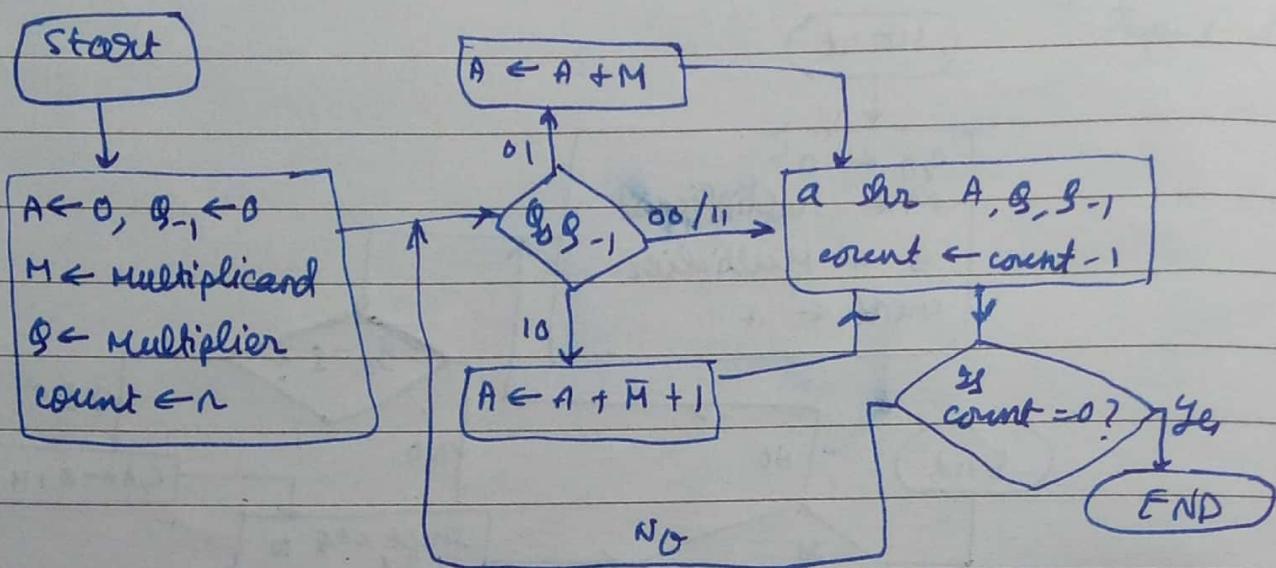


Ex.

$$\begin{array}{r}
 110 \\
 \times 101 \\
 \hline
 110 \\
 000 \\
 \hline
 110
 \end{array}
 \quad \text{At first, } B = 1 \\
 C, A \leftarrow 0 + M \\
 \uparrow \uparrow \\
 000 \quad 110$$

| Step | Operations | C | A | B | M |
|------|-----------------------|---|------|------|------|
| 4 | Add M | 0 | 0001 | 1101 | 1011 |
| 3 | shift CAG | 0 | 1011 | 1101 | |
| 2 | b=0 shift CAG | 0 | 0101 | 1110 | |
| 1 | shift CAG | 0 | 0010 | 1111 | |
| 0 | b=1, Add M shift B | 0 | 1011 | 1111 | |

$B = 11, M = 13$
 $B \times M = 143$

Signed integer multiplication (unigned else)Booth's algo

ashra : arithmetic
shift right (Keeps most significant
bit same)

$\begin{smallmatrix} 1 & 0 & 1 & 1 \\ \downarrow & & & \\ 1 & 1 & 0 & 1 \end{smallmatrix}$

shri : logical shift
ca@

$\begin{smallmatrix} 1 & 0 & 1 & 1 \\ \downarrow & & & \\ 0 & 1 & 1 & 0 \end{smallmatrix}$

$$\begin{array}{r} x \times y = a \\ \hline (n) & (m) \end{array}$$

example:

0000000000000000

| Operation | (Partial product) A | B_n | B_{n-1} | (Multiplied) M | count |
|-----------|------------------------|-------|-----------|-------------------|-------|
| | 000 | 11010 | | 0111 | 100 |

$$\begin{array}{r} 00B_{-1}=10 \\ \text{sub } M \\ \hline 100 \end{array}$$

$$\begin{array}{r} \text{ashr} \rightarrow 1100 \\ 1101 \\ 1110 \\ \hline 011 \end{array}$$

$$\begin{array}{r} 00B_{-1}=0 \\ \text{add } M \\ \hline 0011 \end{array}$$

$$\begin{array}{r} \text{ashr} \rightarrow 0001 \\ 1110 \\ 1110 \\ \hline 010 \end{array}$$

$$\begin{array}{r} 00B_{-1}=10 \\ \text{sub } M \\ \hline 1010 \end{array}$$

$$\begin{array}{r} \text{ashr} \rightarrow 1101 \\ 0111 \\ 0111 \\ \hline 001 \end{array}$$

$$\begin{array}{r} \text{ashr} \rightarrow 1110 \\ 1011 \\ 1011 \\ \hline 000 \end{array}$$

2's complement of 2,
 $= -2$

Assignment - 1:

$$1) (-7) \times 3$$

$$2) (-7) \times (-3)$$

Worst case: 10101010

Best case: 11111111

} Possible multipliers for
 Booth's algo.

Booth's encoding:

$$\left\{ +100 -100 00 \right.$$

$$10110110$$

+ 1+10 -1+10 -10 } encoded multiplier

- + → Add 1's of A
- → Sub 1's of A

Bit pair multiplication

- Booth's algo has 2 attractive features
 - Handles both +ve and -ve multipliers
 - Efficiency in terms of speed which depends on multiplier
- Bit pair multiplication reduces maximum number of operations and hence improves speed.

| $i+1$ $\xleftarrow{\text{multiplicand}} \xrightarrow{i}$ bit pair | $i-1$ Multiplicand bit at right $B-1$ | Booths | Bit pair |
|---|--|---------------------------------|----------|
| 0 0 | 0 | 0 0 $\rightarrow 0 \times M$ | |
| 0 0 | 1 | 0 +1 $\rightarrow +1 \times M$ | |
| 0 1 | 0 | +1 -1 $\rightarrow +1 \times M$ | |
| 0 1 | 1 | +1 0 $\rightarrow +2 \times M$ | |
| 1 0 | 0 | -1 0 $\rightarrow -2 \times M$ | |
| 1 0 | 1 | -1 +1 $\rightarrow -1 \times M$ | |
| 1 1 | 0 | 0 -1 $\rightarrow -1 \times M$ | |
| 1 1 | 1 | 0 0 $\rightarrow 0 \times M$ | |

$$+100 -10060$$

$$\hookrightarrow -1 \times 2^4 + 1 \times 2^7$$

$$+1 -1$$

$$\hookrightarrow +1 \times 2^1 - 1 \times 2^0 = 2 - 1 = +1$$

$$\text{Ex. } 13 \times (-6)$$

Multiplicand = -6

$$6 = 00110$$

$$-6 = 11001$$

$$\begin{array}{r} +1 \\ \hline 11010 \end{array}$$

we add this bit to make no. of bits in multiplication even.

$$\begin{array}{r} 110100 \\ \uparrow \\ i-1 \end{array}$$

Processor adds this

to left blank bits are filled with 0. if -ve, the left blank bits are filled with 1

$$11010 \text{ becomes}$$

$$\begin{array}{r} 111010 \\ 00-1+1-10 \end{array}$$

Bit pair

Booth's

$$\begin{array}{r} 13 = 01101 \\ -6 = 0-1+1-10 \\ \hline 0 \rightarrow 0000000000 \\ -13 \rightarrow 1111100111 \\ +13 \rightarrow 00001101 \\ -13 \rightarrow 11100111 \\ 0 \rightarrow 0000000 \end{array}$$

$$\boxed{1110110010}$$

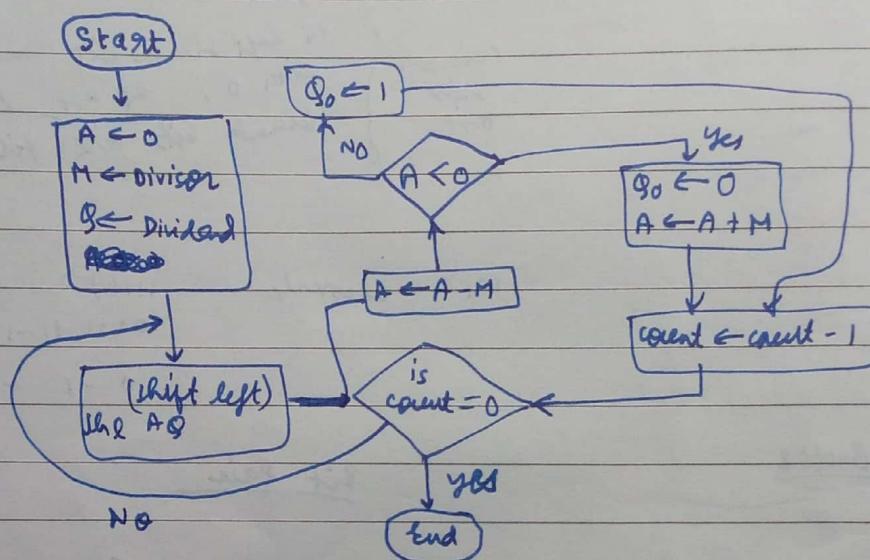
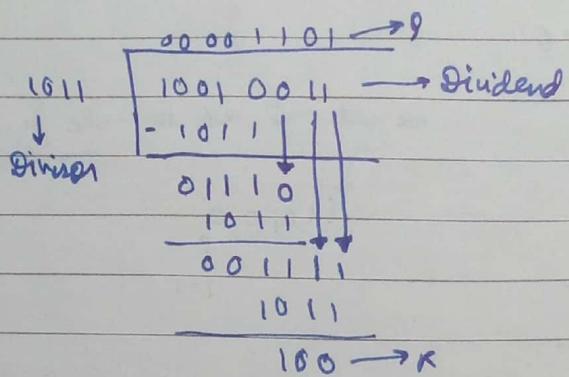
\downarrow sign bit 10 bits = 78

\downarrow sign bit = (+ve) $: Ans. = -78$

Bit Pair

$$\begin{array}{r} 13 = 01101 \\ -6 = 0-1-2 \\ \hline 1111100110 \rightarrow -28 \\ 11110011 \rightarrow -13 \\ 000000 \rightarrow 0 \end{array}$$

$$1110110010$$

Division

Flowchart

↑
only for unsigned integer division.

$$\begin{array}{r} \text{by, } 7 \div 3 \\ \hline \end{array}$$

| | A | B | M | Covert |
|---------------------|---|------|-------|--------|
| Initial value | 0000 | 0111 | 0011 | 100 |
| Shl A B Sub M | <div style="display: flex; align-items: center;"> 0000 1110 </div> <div style="display: flex; align-items: center;"> 0011 1110 </div> 1101 1110 | | | |
| A < 0 ∴ add M | <div style="display: flex; align-items: center;"> 0011 1110 </div> 0000 1100 | | | 011 |
| Shl A B | 0001 | | | |
| Sub M | 0011 | | | |
| Add M, $g_0 = 0$ | <div style="display: flex; align-items: center;"> 1110 1100 </div> 0011 1100 | | | 010 |
| Shift A B Sub M | <div style="display: flex; align-items: center;"> 0011 1000 </div> <div style="display: flex; align-items: center;"> 0000 1001 </div> 0001 0010 | | ----- | ---00 |
| Shift A B Sub M | <div style="display: flex; align-items: center;"> 0011 0010 </div> 1110 | | | |
| $g_0 \leftarrow 1$ | | | | |
| Shift A B Sub M | <div style="display: flex; align-items: center;"> 0011 0010 </div> 0011 | | | |
| $g_0 \leftarrow 0$ | | | | |
| Add M | | | | |

$$\begin{array}{r} 1101 \\ \downarrow \\ M \end{array} \quad \begin{array}{r} 1001 \\ \downarrow \\ A \end{array}$$

$$A < M \Rightarrow A - M < 0$$

\therefore unsuccessful
 \therefore restore

Division of
signed
integers
↓

- ① Division in register M
- ② Dividend in register A
- ③ Init A = 0
- ④ If M, A have same sign, perform $A \leftarrow A - M$
else $A \leftarrow A + M$

⑤ Preceding operation is successful if sign of A is same before and after operation

⑥ If operation is successful,

$$Q_0 \leftarrow 1$$

If operation is unsuccessful

set $Q_0 \leftarrow 0$ & return previous value of A

⑦ repeat till quotient = 0.

⑧ If sign of divisor & dividend are same,
quotient is in register Q,
otherwise correct quotient is 2's complement of Q

$$\text{Ex. } (-7) \div 3$$

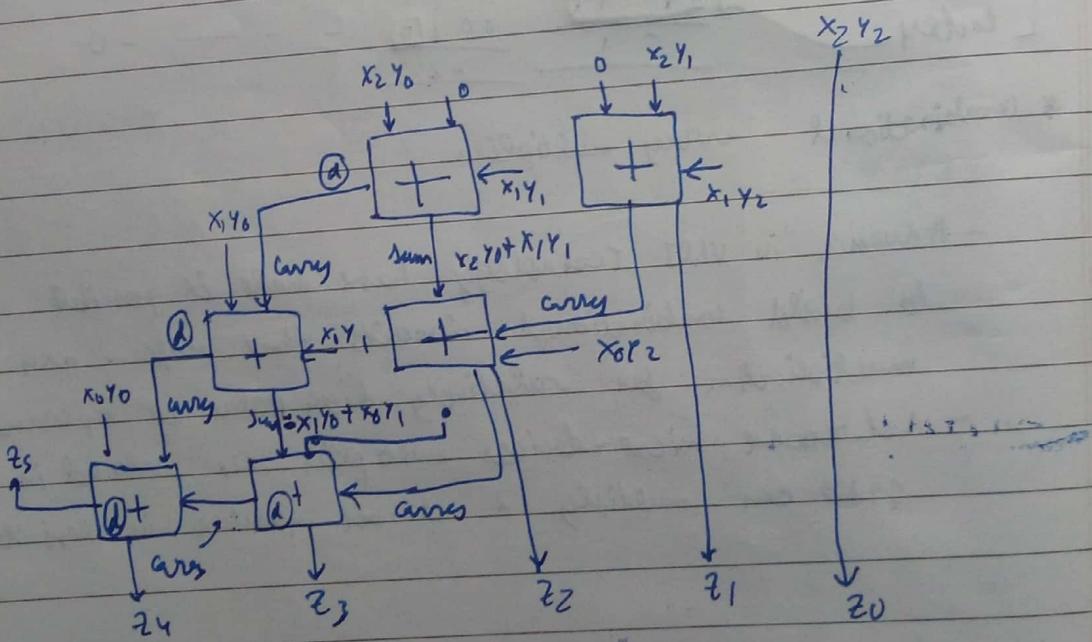
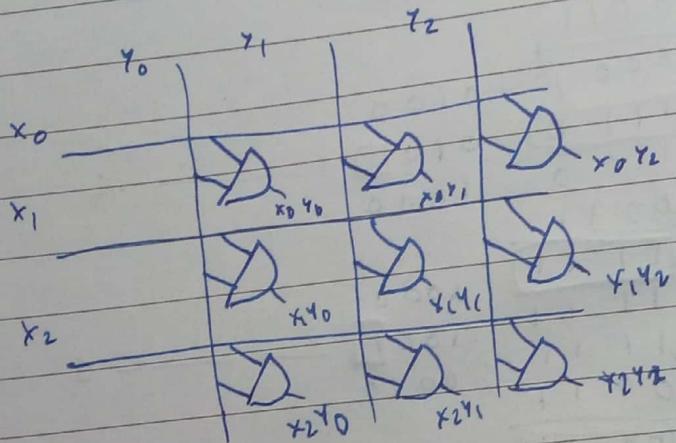
$B = 2^1$'s complement of 0010
 $R = 2^1$'s complement of 1111

| Operation | (dividend) | (divisor) | M. | Count |
|-----------|--|--|------|-------|
| | A | B | M. | |
| Sub A/B | $\begin{array}{r} 1111 \\ - 0011 \\ \hline 0010 \end{array}$ | $\begin{array}{r} 1001 \\ (-7) \end{array}$ | 0011 | 4 |
| Add M | | | (3) | |
| Restore | $\begin{array}{r} 1111 \\ - 1110 \\ \hline 0011 \end{array}$ | $\begin{array}{r} 0010 \\ 0100 \\ \hline \end{array}$ | | 3 |
| The A/B | | | | |
| Sub A/B | $\begin{array}{r} 1111 \\ - 0000 \\ \hline 1111 \end{array}$ | $\begin{array}{r} 0100 \\ 1000 \\ \hline 0011 \end{array}$ | | 2 |
| Add M | | | | |
| Restore | $\begin{array}{r} 1111 \\ - 1100 \\ \hline 0011 \end{array}$ | $\begin{array}{r} 1000 \\ 1001 \\ \hline \end{array}$ | | 1 |
| The A/B | | | | |
| Sub A/B | $\begin{array}{r} 1111 \\ - 0010 \\ \hline 1111 \end{array}$ | $\begin{array}{r} 0010 \\ 1010 \\ \hline \end{array}$ | | 0 |
| Add M | | | | |
| Restore | | | | |

* Combinational array multiplier.

- Advances in VLSI technology have made it possible to build combinational circuits that perform $n \times n$ bit multiplication for relatively high value of operand.
- Advanced micro-devices multipliers chip produced in 1985 can multiply 2 32-bit number in less than 80ns.

$$\begin{array}{c}
 \begin{array}{ccc} Y_0 & Y_1 & Y_2 \end{array} \\
 \begin{array}{c} X_0 \quad X_1 \quad X_2 \\ \hline X_2Y_0 \quad X_2Y_1 \quad X_2Y_2 \end{array} \rightarrow 2^0 X_2 Y \\
 \begin{array}{c} X_1Y_0 \quad X_1Y_1 \quad X_1Y_2 \\ \hline X_0Y_0 \quad X_0Y_1 \quad X_0Y_2 \end{array} \rightarrow 2^1 X_1 Y \\
 \begin{array}{c} X_0Y_0 \quad X_0Y_1 \quad X_0Y_2 \\ \hline Z_0 \quad Z_1 \quad Z_2 \end{array} \rightarrow 2^2 X_0 Y \\
 \hline
 \end{array} \rightarrow \sum_{j=0}^{2^n-1} 2^{n-1-j} X_j Y$$



[FA. has delay = d] \rightarrow digital logic knowledge

$$\text{Total delay of circuit above} = d + d + d + d = 4d$$

In terms of n ,

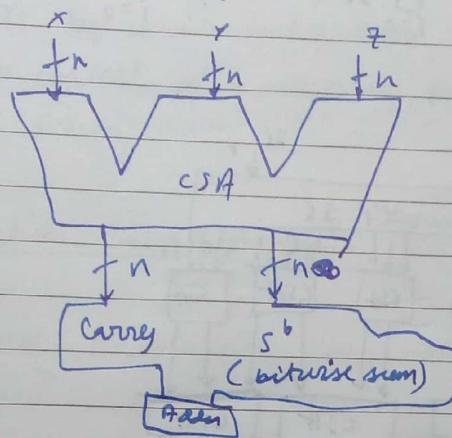
$$\text{Total delay} = 2(n-1)d \rightarrow \text{delay of 1 FA.}$$

\downarrow
n bits
in multiplier

No. of bits in multiplier = n

\Rightarrow No. of partial products = n

Carry Save Adder (CSA)



$$S^b = (s_0, s_{n-1}, \dots, s_1, s_0)$$

↳ leading bit always 0.

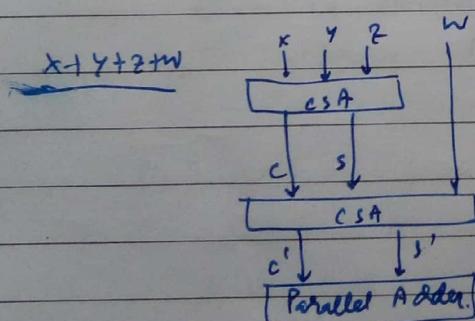
$$C = (c_n, c_{n-1}, \dots, c_1, 0)$$

↳ always 0

$$S_i^* = x_i \oplus y_i \oplus z_i$$

$$c_{i+1} = (x_i y_i) \text{ or } (y_i z_i) \text{ or } (z_i x_i)$$

$$\begin{array}{r} s_0 = 1 \\ \hline x = 001011 & i=0 \\ y = 010101 \\ z = 111101 \\ \hline s^b = 010001 \\ \hline c = 011101 \\ \hline 1011101 \end{array}$$

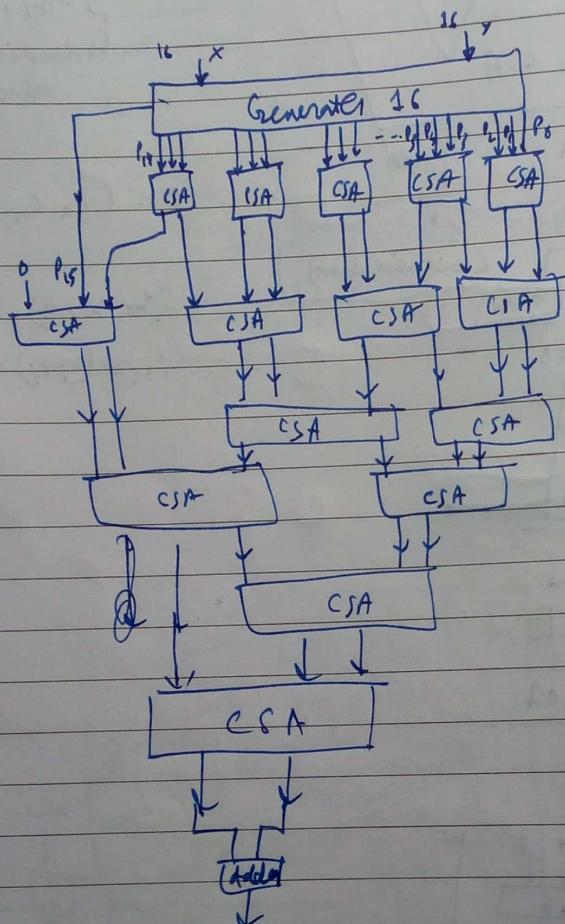


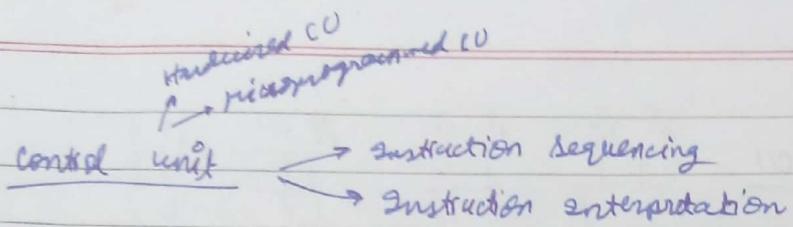
Multiplication using CSA

n-bit M in Y

n-bit Q in X

$$\text{Product } P = \sum_{i=0}^{n-1} x^i z^i y = \sum_{i=0}^{n-1} p_i \quad \rightarrow \text{Partial products}$$



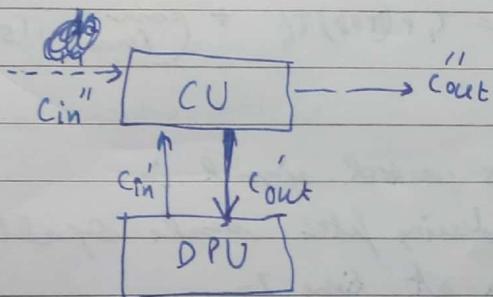


→ Instruction sequencing

Methods by which instructions are selected for execution, or how control of processor transfers from one instruction to other.

→ Instruction interpretation

CU interprets instruction → determine control signal to be issued.



$Cout$ directly control operation of DPU.

Cin' indicate occurrence of unusual condition such as errors in DPU.

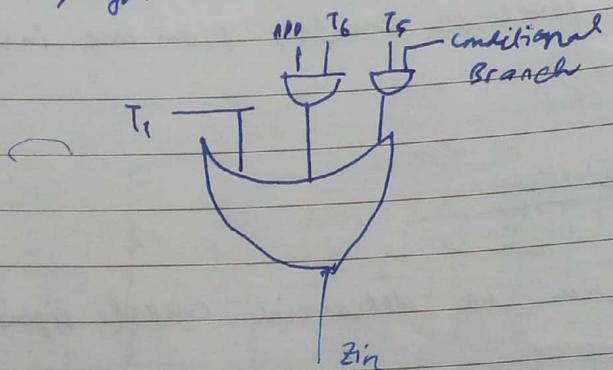
Cin' allows data-dependent decisions to be made.

Cin'' receives from other CU (for eg. from a supervisor or controller). Typically include START, STOP, timing info.

$Cout''$ transmitted to other CU may indicate status condition such as busy, operation condition, etc.

Hardwired CU

- Ex. let control signal Z_1 to be activated at T_1 for all instance.
 T_6 for ADD instruction
 T_5 for condition branch instruction.

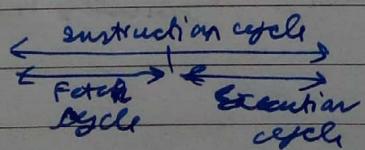


$$Z_{in} = T_1 + (\text{ADD})T_6 + (\text{cond branch})T_5$$

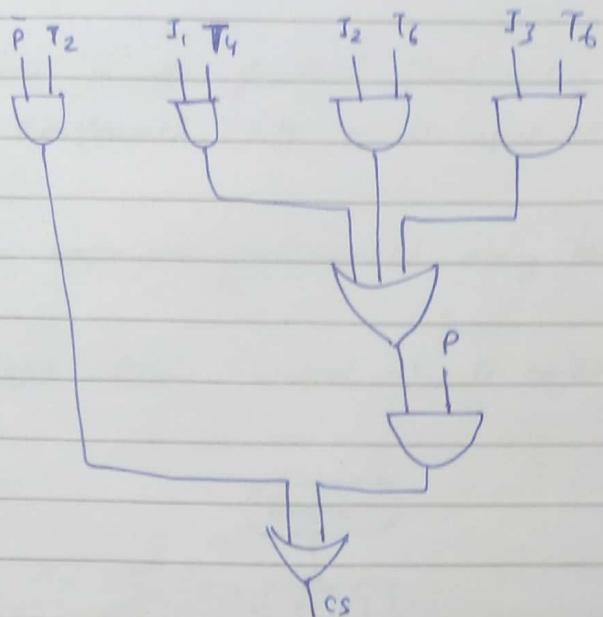
Ex. Suppose C_5 control signal is required during fetch cycle of all instruction at time T_2 .

[During execution of I_1 at T_4
 ... I_2 at ~~T_5~~ T_6
 ... I_3 at ~~T_5~~ T_6 .]

$p=0$ for fetch
 $p=1$ for execute



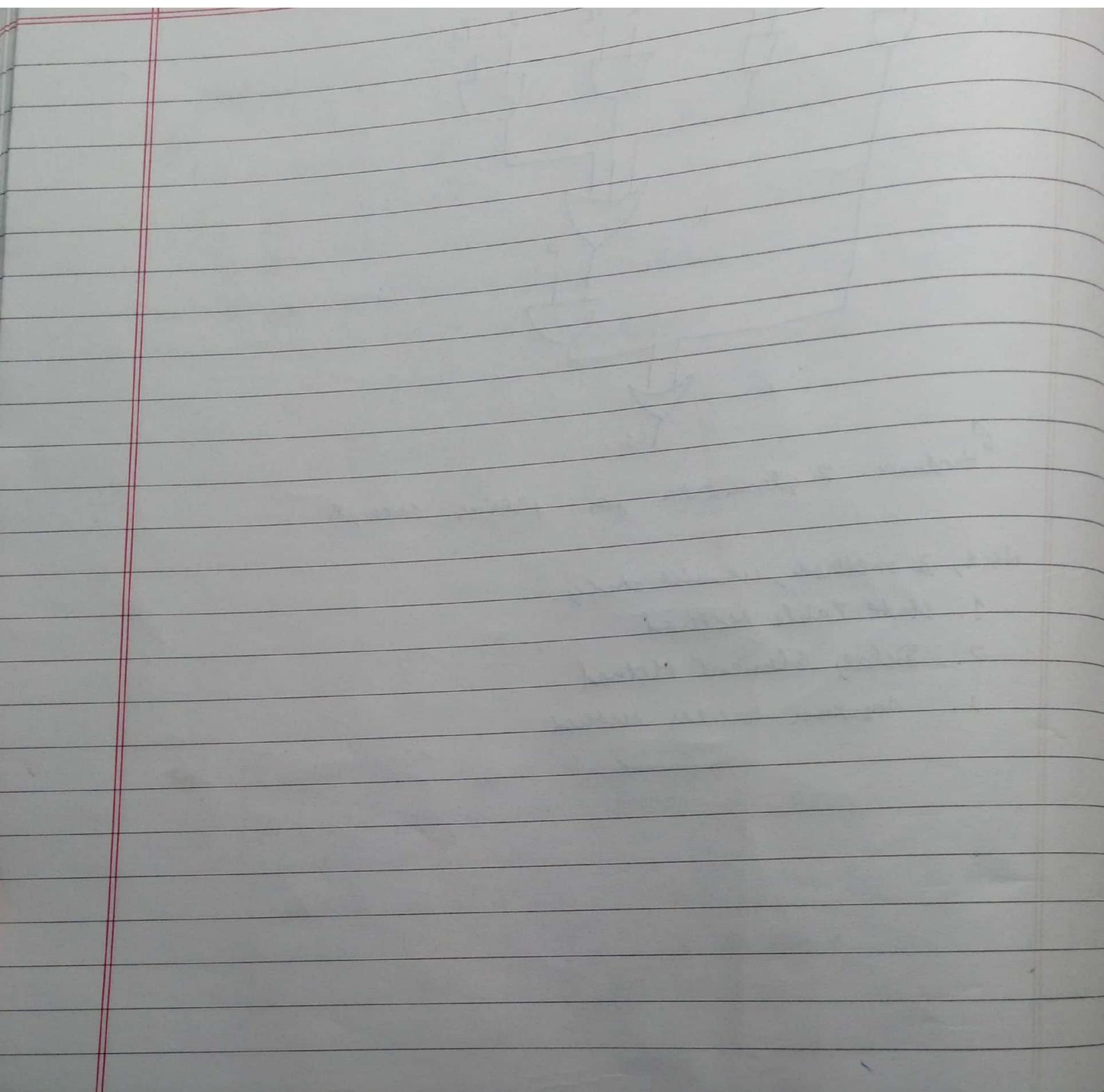
$$C_5 = \bar{p} T_2 + p (I_1 T_4 + I_2 T_6 + I_3 T_6)$$



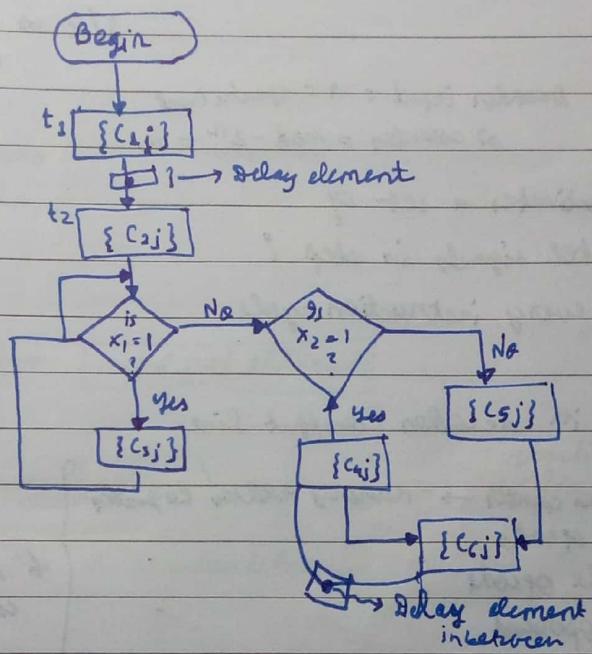
Hardware Implementation for previous example

Next, 3 methods we will study.

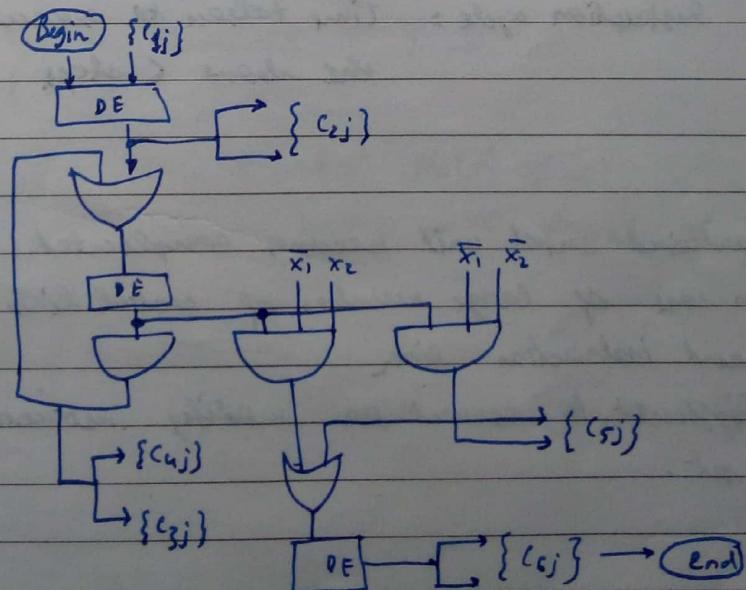
1. state Table Method
2. Delay Element Method
3. Sequence Counter Method



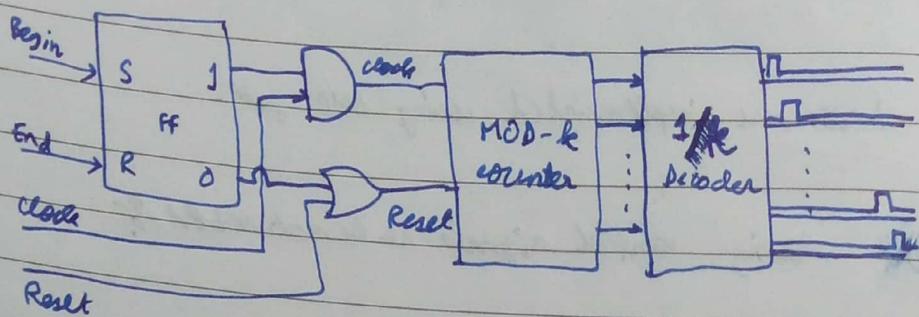
- Rules:*
- K lines in the precedent that merge into a common line are transformed into a K-input OR gate.
 - Decision boxes are implemented using AND gate.
 - Flowchart showing control signals to be activated :-



- Circuit diagram using Rules(to convert the flowchart above) :



Sequence counter Method:



$1/k \Rightarrow 1$ output

and 1 input

only are
high at any
time.

Decoder input = $n = \text{counter input}$
 $\Rightarrow \text{counter} = \text{mod } 2^n = k$

Φ_i activates a set of
control signals in step i
of every instruction cycle.

Φ_i : i^{th} decoder output line

- | | | |
|--------------------|--|---|
| {c _{ij} } | 1. Program counter \rightarrow Memory Address Register | } |
| {c _{ij} } | 2. Read opcode | |
| {c _{ij} } | 3. Decode opcode | |
| {c _{ij} } | 4. Read operand | |
| {c _{ij} } | 5. Execute I | |
| {c _{ij} } | 6. Store result | |

6 steps to
execute 1
instruction

Instruction cycle: Time taken to execute all
the above 6 steps.

• Hardwired control unit becomes complicated
in case of large numbers of control lines
and instruction lines

• Difficult to correct or modify instruction
set.

Microprogrammed CU

I Micro-program (cp)

$$I_1 \quad M_P \rightarrow \begin{bmatrix} uI_1 \\ uI_2 \\ \vdots \\ uI_m \end{bmatrix}$$

$$I_2 \quad M_P$$

:

$$I_n \quad M_P$$

microprogram micro-instruction

microprogram

$$M_P \xrightarrow{x} uI_1 \rightarrow \boxed{\text{control field}} \quad \boxed{\text{Address field}}$$

$$I_1 \rightarrow \boxed{\text{opcode}} \quad \boxed{\text{addr}}$$

Now this x is calculated

x : starting address of M_P

(for the previous part example, $uI_1, uI_2, uI_3, \dots, uI_6$)

$$\begin{array}{l} x(uI_1) \quad \boxed{\text{control}} \quad \boxed{x+1} \\ x+1(uI_2) \quad \boxed{\text{control}} \quad \boxed{y} \\ y(uI_3) \end{array}$$

→ control field of uI_1
contains control signals of $\{c_{ij}\}_{ij}$ group

Instruction = operation

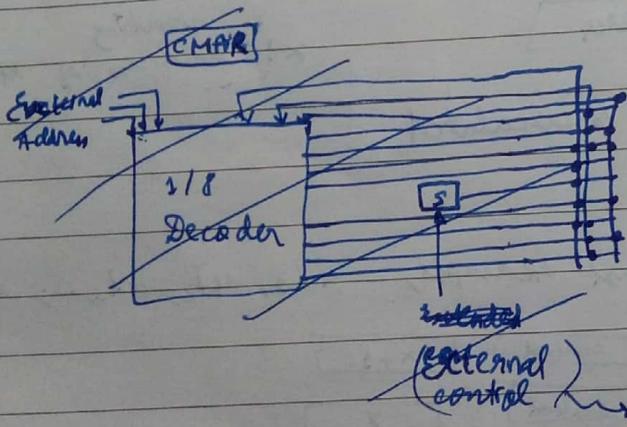
↳ can be divided into a few u operations

Load X
 Load accumulator by content of
 memory location X.
 $AC \leftarrow M[X]$

Instruction: Load R

micro-operation: $R \leftarrow R_1$
 $R \leftarrow R_2$
 $R \leftarrow R_3$
 $R \leftarrow R_4$

Willems Design



3 bits for address
 8 micro-instructions

$$2^3 = 8$$

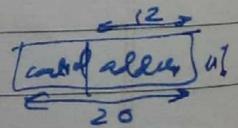
∴ these 2 are related

e.g. $(S_3 AC = 0) ?$

Redraw at
down

$S \rightarrow$ switch

with 12 bits,
 No. of add = 2^{12}



∴ size of control memory = $2^{12} \times 2^0$

→ Micro-programmed control unit is slower than hardwired control.

uProgram counter goes from 001 to 001---

Address field of uInstruction points to next one (\therefore can be solved by memory are to iAddress)

Then, don't need to provide

address field of u-Instruction.

Just send external address from uProgram counter

Instructions in RAM

uProgram in control memory

} \rightarrow 2 different memories involved

Read opcode of I₁,

Decide Opcode of I₁,

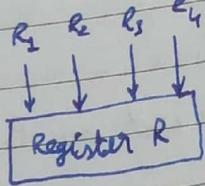
From Opcode, you get starting address of MIs

uP₁ had 10 MI.

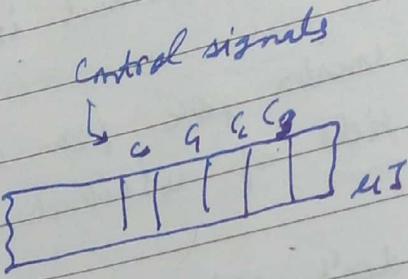
\therefore Time for I₁ = total time for fetch and execute
of all the MI for uP₁.

\therefore This is why uP control unit is slower.

Horizontal UI format



$* R \leftarrow F_1$
 $R \leftarrow F_2$
 $R \leftarrow R_3$
 $R \leftarrow R_4$



e.g. $C_0 \rightarrow 1000$
 $C_1 \rightarrow 0100$
 $C_2 \rightarrow 0010$
 $C_3 \rightarrow 0001$
 NO operation $\rightarrow 0000$

Horizontal = long format

n-control signals needed

\Rightarrow size of control field = n

Vertical UI

| $K_1 K_2 K_3$ | |
|---------------|--------------------------------|
| 0 0 0 | \Rightarrow NOP |
| 0 0 1 | $\Rightarrow R \leftarrow R_1$ |
| 0 1 0 | $\Rightarrow R \leftarrow R_2$ |
| 0 1 1 | $\Rightarrow R \leftarrow R_3$ |
| 1 0 0 | $\Rightarrow R \leftarrow R_4$ |

We can have at least 1 bit per UI.

By Reading we get $K_1 K_2 K_3$

\therefore less bits than horizontal.

But to decode $K_1 K_2 K_3$ into $C_0 C_1 C_2 C_3$
we need extra decoder in case of vertical.

height = No. of instructions

width = size of instruction

} attempts have been made to
reduce this?

↳ reduce no. of
instructions

reduces
height

size of control memory = width * height

- Examples of independent & dependent -

$AC \leftarrow 0$

$AC \leftarrow AC - DR$

$PC \leftarrow AR$

$AC \leftarrow AC \text{ OR } DR$

$PC \leftarrow RC + 1$

In micro-programming level, independent MI are available & executed in parallel. Hence one MI can specify signals that can execute multiple micro-operations.

CLOCK

Parallelism in MI

Q) No. of micro-operation = 46

Degree of parallelism = 2

Size of control-field = 9

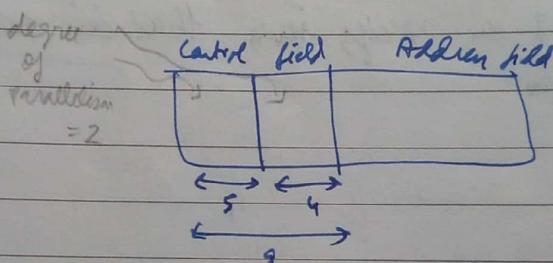
$$\begin{aligned} q &= 7+2 \\ &= 1+1 \\ &= 6+3 \\ &= 5+4 \end{aligned} \quad \left. \right\} \text{which one?}$$

$$q = 5+4$$

$$\text{Because } 2^5 + 2^4 = 48$$

48 is closest to 46
and 2 operations are no operation C must be

$$\therefore 46+2 = 48 = 2^5 + 2^4$$



Micro-operation timing

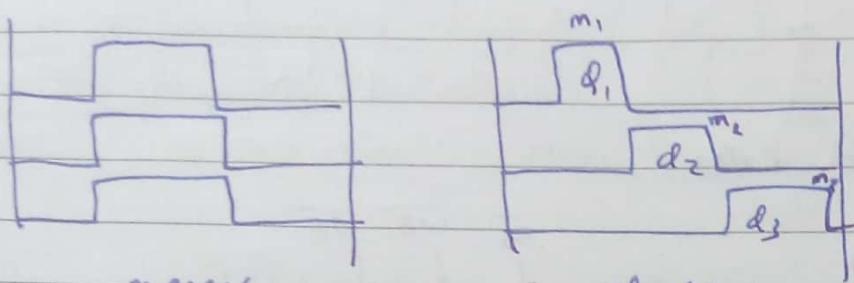
$$m_1: R_2 \leftarrow f(R_1) \rightarrow M_1$$

$$m_2: R_3 \leftarrow f(R_2) \rightarrow M_2$$

$$m_3: R_4 \leftarrow f(R_3) \rightarrow M_3$$



Monophase: Control signals required to execute micro-operation remain active throughout MI cycle.



monophase

polyphase

no. of u.I.
reduces
from 3 to 1

Though u-op are dependent, but
we can execute them in
3 - different phases in the same
u.I.

Back to micro-operation timing. (polyphase)

(Micro-code compaction) (minimizing H)

| <u>NAME</u> | <u>u-Operation</u> | <u>Phase used</u> | <u>Resource used</u> |
|-------------|------------------------------------|-------------------|----------------------|
| m_1 | $RLEFT \leftarrow AC$ | Φ_1 | |
| m_2 | $RRIGHT \leftarrow 1$ | Φ_1 | |
| m_3 | $ROUT \leftarrow RLEFT + RRIGHT$ | Φ_2 | Adder |
| m_4 | $MAR \leftarrow RRIGHTSHIFT(ROUT)$ | Φ_3 | shifter |
| m_5 | $M(MAR) \leftarrow MDR$ | Φ_1 | Bus. |
| m_6 | $RLEFT \leftarrow 1$ | Φ_1 | |
| m_7 | $RRIGHT \leftarrow R_0$ | Φ_1 | |

$$MI_1 = \{m_1, m_2, m_3, m_4\}$$

$$MI_2 = \{m_5, m_6, m_7\}$$

\therefore no. of micro-instructions = 2 for polyphase

\therefore height \downarrow [H=7 \rightarrow H=2]

| <u>② NAME</u> | <u>u-Operation</u> | <u>Phase</u> | <u>Resource</u> |
|---------------|--------------------------------------|--------------|-----------------|
| m_1 | $A \leftarrow R_1$ | Φ_1 | |
| m_2 | $C \leftarrow A + R_2$ | Φ_2 | Adder |
| m_3 | $R_3 \leftarrow D$ | Φ_2 | |
| m_4 | $R_6 \leftarrow R_6 + 1$ | Φ_2 | Adder |
| m_5 | $R_7 \leftarrow R_3 \text{ or } R_6$ | Φ_2 | ALU |

$$MI_1 = \{m_1, m_2, m_3\}$$

$$MI_2 = \{m_4\}$$

$$MI_3 = \{m_5\}$$

$$MI_1 = \{m_1, m_2, m_3\}$$

$$MI_2 = \{m_4\}$$

$$MI_3 = \{m_5\}$$

H=3

further reduce \Rightarrow

H=2

Minimizing w

Objective is to derive a format of the micro-instruction control field such that the total number of bits in the control field is minimum.

∴ We want a control field encoding method that uses as few bits as possible

Compatibility class

c_i = control signal

I_j = Micro-instruction

$c_i \in I_j \rightarrow c_i$ is activated by I_j

control signals c_1 & c_2 are compatible
if $c_1 \in I_j$ but $c_2 \notin I_j$

Maximal compatibility class

Defined as compatibility classes to which no control signal can be added without introducing a pair of incompatible control signals

e.g. UI control signals

| | |
|-------|------------|
| I_1 | a, b, e, g |
| I_2 | a, c, e, h |
| I_3 | a, d, f |
| I_4 | b, e, f |

extra reqd. 1, 2, 4
- 8 control signals
target "g" to reduce "g"
using more no. of bits

Algo is enacted step-by-step next

Step-1: Determine the set of MCCs.

| | | | | |
|-------|--|--|--|--|
| S_1 | (a) bcd & ghi | a, b a, c a, d a, e a, f a, g a, h | b, g b, e b, d b, c b, f b, g b, h | ca cb cd ce cf cg ch |
| S_2 | bcd, bde, bgh, (cd), deg, dgf, dhf ef, fg, fg, gh, gh | | | |
| S_3 | (bde) (bgh) (deg) (dgh) (efg) (fg) | | | |

new and no
intervall for
other different
examples $\rightarrow \{ \}$

Process terminates when no new compatibility classes can be formed.

Step-2: Determination of all minimal set of MCCs that include each control signal.
Each of these set is minimal MCC cover.

COVER TABLE
 \hookrightarrow

Each row is a compatibility class obtained above.

Each column is a control signal.

| $c_i \setminus c_j$ | a | b | c | d | e | f | g | h |
|---------------------|---|---|---|---|---|---|---|---|
| $c_1 = a$ | x | | | | | | | |
| $c_2 = cd$ | | | x | x | | | | |
| $c_3 = bde$ | x | | x | x | | | | |
| $c_4 = bdh$ | x | | x | | | x | | |
| $c_5 = deg$ | | | x | x | x | | | |
| $c_6 = dgh$ | | | x | | x | x | | |
| $c_7 = efg$ | | | x | x | x | | | |
| $c_8 = fgh$ | | | x | x | x | | | |

Certain rows and columns can be deleted from cover table to simplify determination of minimal MCC cover.

\rightarrow column c_j contains only 1 x occurring in row c_i , then c_i is said to be essential MCC and they must appear in every minimal MCC cover since it only covers c_j . All rows corresponding to essential MCCs can be

c_1, c_2
essential

columns
a, c
deleted

deleted from coverable and all columns with \times mark in essential rows can be deleted.

columns
d, g
deleted

→ column c_i is said to dominate column c_j if c_i contains a \times marks in every row where c_j contains a \times mark and c_i has more \times .

Dominating column can be deleted as every MC that covers c_j automatically covers c_i .

Rows
aeg, dgh
deleted

→ Row c_i dominates c_j if c_i contains a \times in every column where c_j contains \times and c_i containing more \times than c_j .

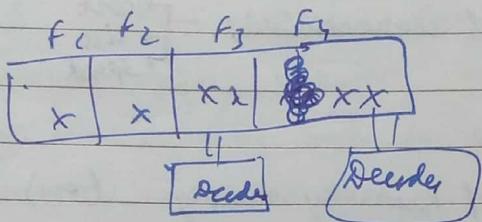
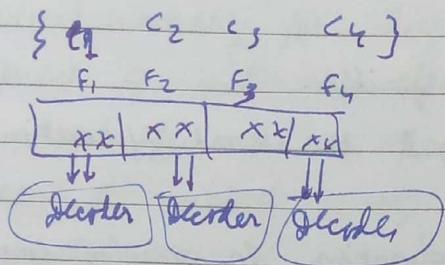
Dominated row c_i can be deleted since c_i covers all control signals that are covered by c_j .

Minimal cover for this table is

$$*\{c_1, c_2, c_4, c_7\} \quad \begin{matrix} a \\ a \\ c \\ bdm \\ ebg \end{matrix} \quad \rightarrow \text{these are the 2 choices possible.}$$

$$*\{c_1, c_2, c_3, c_8\} \quad \begin{matrix} a \\ a \\ c \\ bde \\ fgh \end{matrix}$$

$$\begin{aligned} R_1, R_2, R_3, R_4 &\rightarrow c_1 = a \\ c_2 - c_1 &= cd \\ c_3 - bde &= be \\ c_4 - bgf &= fg \end{aligned}$$



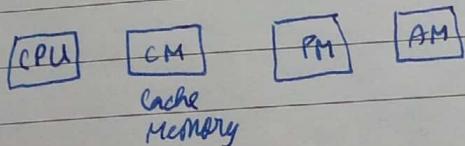
MEMORY SUBSYSTEM

= storage devices + Algorithm
to control or manage the stored information.

- distributes stored information in complex fashion over a variety of different memory units with different physical characteristics.

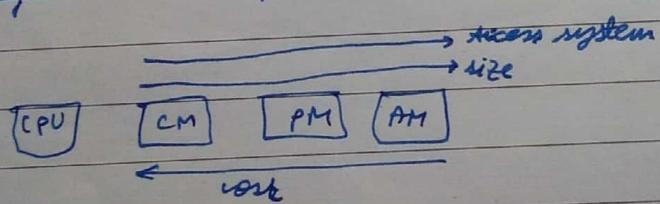
→ cost
→ size
→ speed

- Eg. → Register
- physical / primary memory (PM)
- secondary / auxiliary memory (AM)



Logical position of CM is
inbetween CPU and PM.

Objective of memory system design.



- Memory system provides adequate storage capacity with an acceptable level of performance at a reasonable cost.

- Methods to fulfill objective:

1. Use of a number of different memory devices with different cost / performance ratios organized to achieve high average performance at a low cost per bit.

2. individual memory for a hierarchy of storage devices
3. development of automatic storage allocation method to make most efficient use of most available memory space.
4. Development of virtual memory/ auxiliary memory.
 - development of virtual memory concept to make ordinary user free from memory management and make programs largely independent of memory configuration.
5. Design of communication links to the memory system so that all processors connected to it can operate at or near their maximum rate.

Memory device characteristics

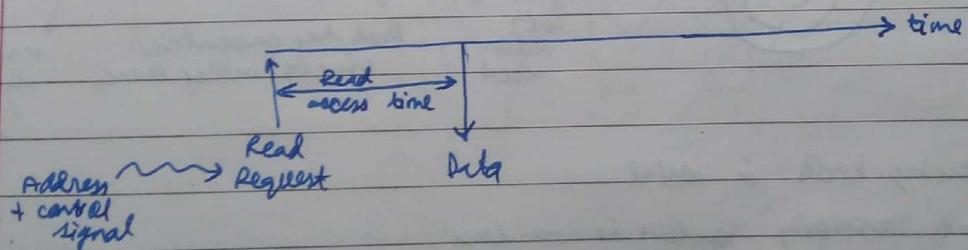
- cost → cost of information storage cells and access circuitry.

$$C = \text{cost in Rs.}$$

$$C = \text{cost per bit} = \frac{C}{S}$$

$$S = \text{storage capacity in bits}$$

- Access time (t_A)



(Similarly for write access time)

- t_A should be low.

Access mode

- Random access mode
- Serial access mode

- Random access mode

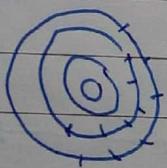
- = Each storage location can be accessed independently of other locations
- = costly, as every location has its own access mechanism

e.g. ferrite-core

- Serial access mode

- = Access mode shared among memory locations.
- = high TA, low cost, TA depends on location to be accessed
- e.g. magnetic tape, audio cassette

Magnetic disks



access [(side), (track no.), (sector no.)]
 ↓ ↓ ↓
 both sides used Disk has concentric tracks
 circles called tracks divided
 into sectors

Accessing track is serial

but accessing sector is random

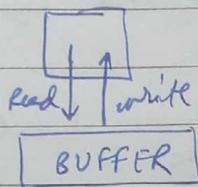
iff every track has its individual read-write head

Performance of storage

- Destructive Read Out (DRO)
- Dynamic Memory
- volatile

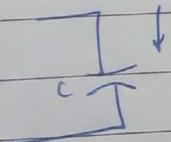
→ DRO

% Method of reading memory destroys stored information, then we store data in a buffer, and follow each read operation by a write operation.

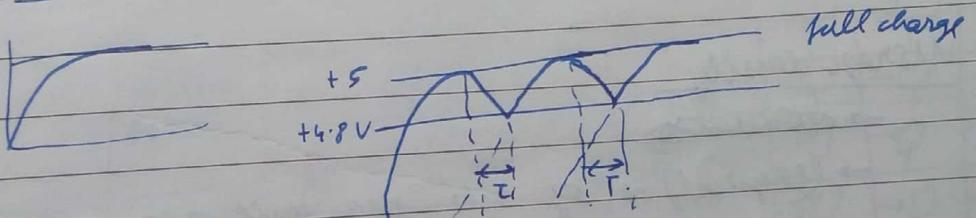


→ Dynamic Memory

Data determined by capacitor



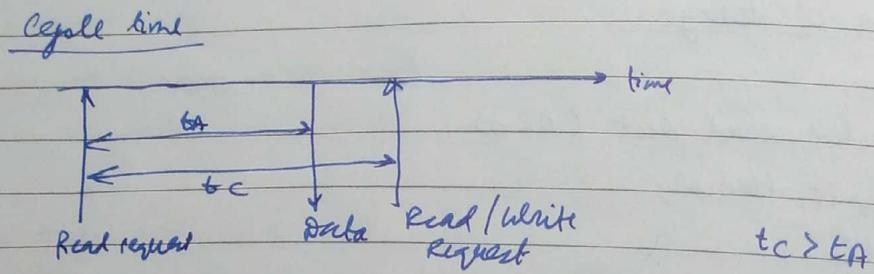
when C is fully charged,
content of memory locations is logic 1



every T seconds, "refresh" operation
occurs to keep C fully charged

→ volatile

Semiconductor memories are volatile
magnetic memories are not volatile



t_c : cycle time

for DRAM, one read cycle includes
one read followed by one write.

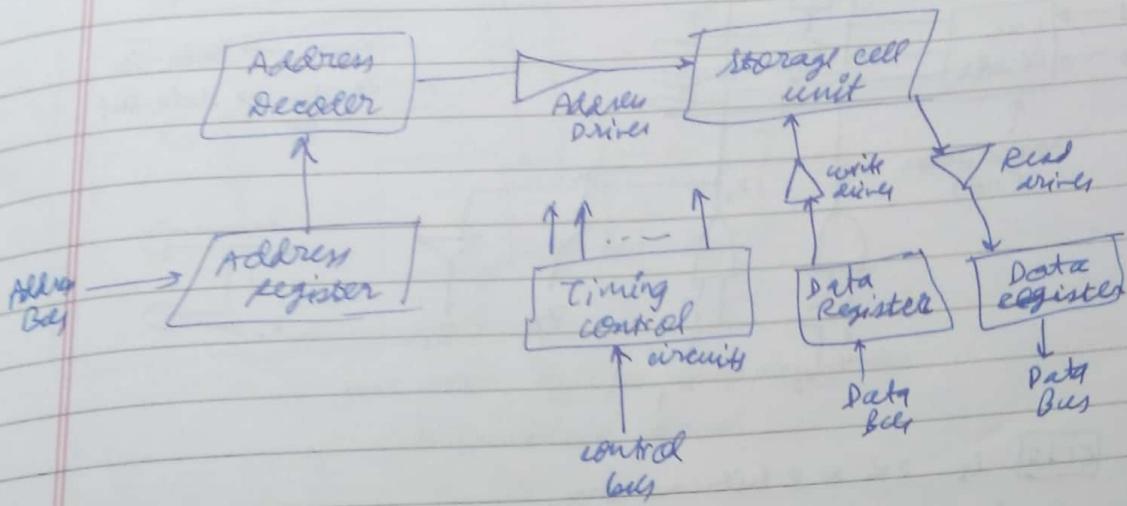
$$\frac{\text{Data transfer rate}}{(\text{bm})} = \frac{1}{t_c} \quad (\text{bits/sec or words/sec})$$

$\frac{\text{bus width}}{(w)} = \text{No. of bits that can be transferred simultaneously for the memory bus.}$

$$bm = \frac{w}{t_c}$$

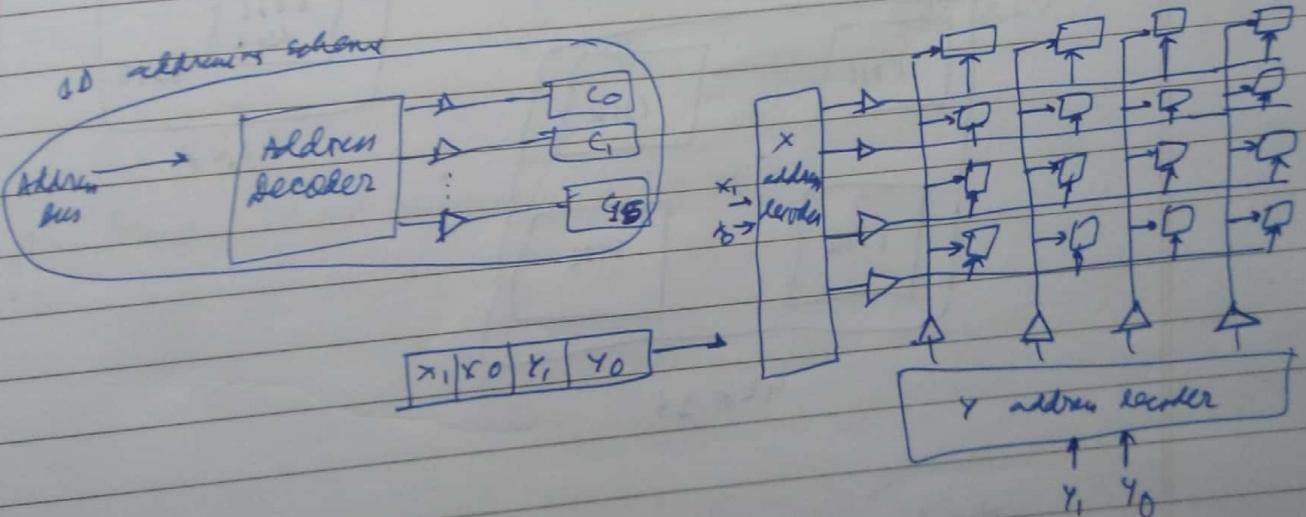
Storage density

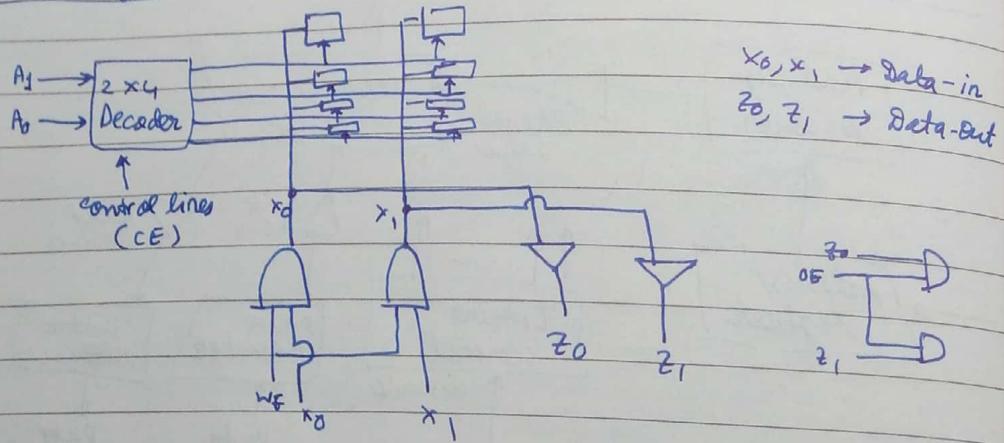
- ~~Reliability~~
- (Physical) no. of bits per unit area or volume
- Determines portability of memory.
- Large energy consumption combined with high storage density requires expensive cooling equipment.
- Reliability
 - less reliable if parts are moving
 - less reliable if storage density is very high or data transfer rate is very high.

Random access memory

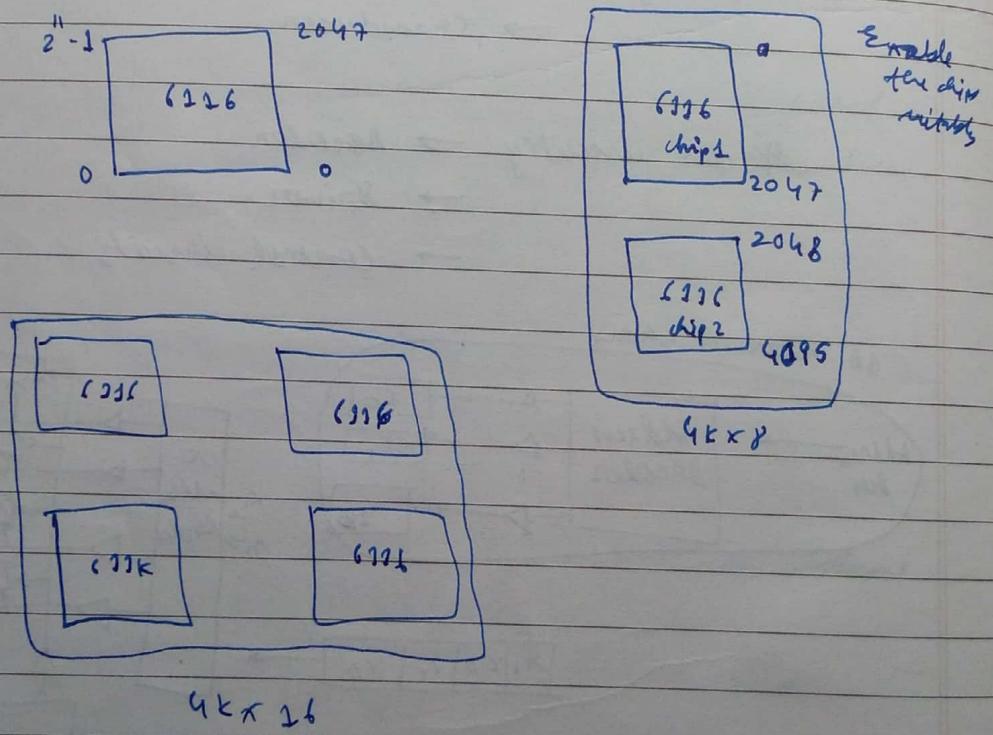
Drives maybe → Amplifier
→ transducer

Access circuitry → decoder
→ driver
→ control circuits



RAM design

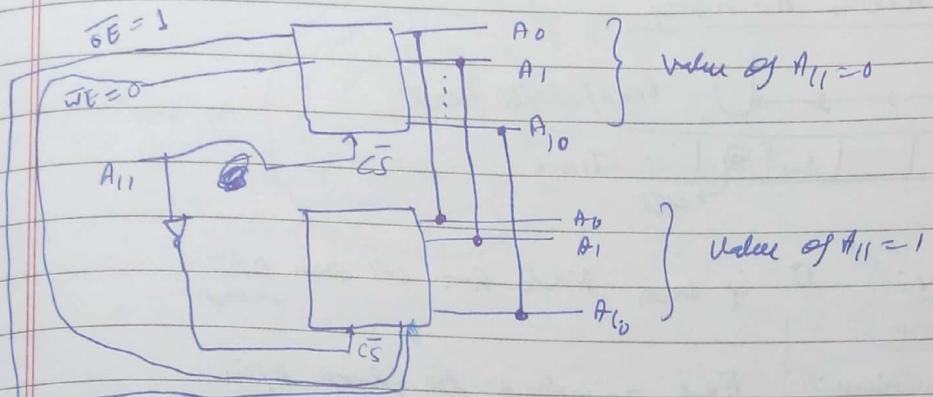
Q) 6116 is $2K \times 8$ bit memory.
Design $4K \times 8$ bit memory using 6116.



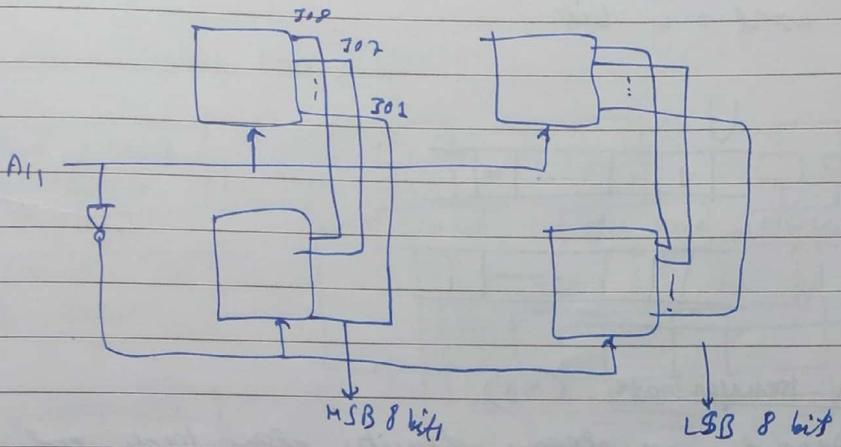
Number of rows in these designs

$$= \frac{\text{No. of desired memory location}}{\text{No. of available memory location}} \quad (y = \frac{4K}{2K})$$

$$\text{Number of columns} = \frac{\text{No. of desired bits per location}}{\text{No. of available bits per location}} \quad (y. \frac{8}{8} = 1)$$

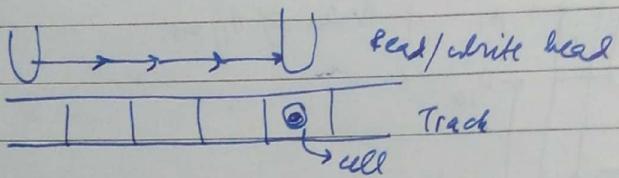


We use A_{11} for supplying the chip-select signal



H.W. Q) Design 32x8 memory using 16x4.

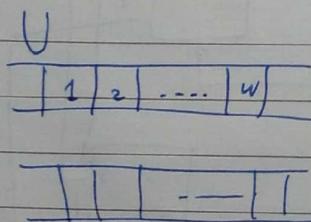
Serial Access Memory



Seek time: O if each track has its own R/W head

Latency time: time required to place R/W head over the required cell in a track

1 word = w bits



Data transfer rate (t_B)

Depends upon storage density along track and speed at which stored information is moving.

T = storage density (bits/cm)

V = velocity (cm/sec)

TV = Data transfer rate (bits/sec)

t_B = Time required to access a block of information.

N = track capacity (words)

r = Rate of revolution (revs/sec)

n = Number of words per block

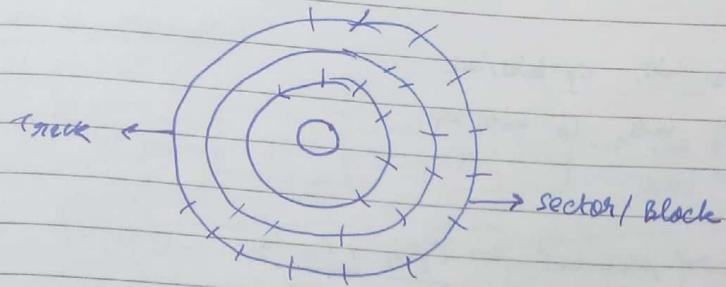
rN = data rate of memory (words/sec)

$\frac{n}{rN}$ = time required to transfer n words

t_s = seek time

Average latency time = $\frac{1}{2r}$

1 rev = 1 completed track

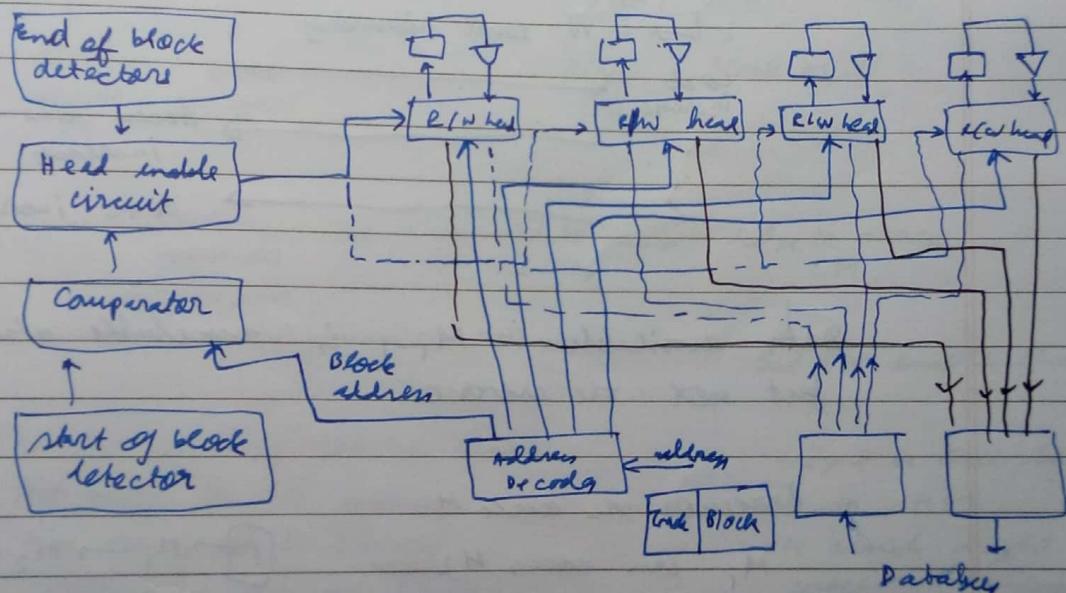


$$t_B = t_s + \frac{1}{2\pi} + \frac{n}{\pi N}$$

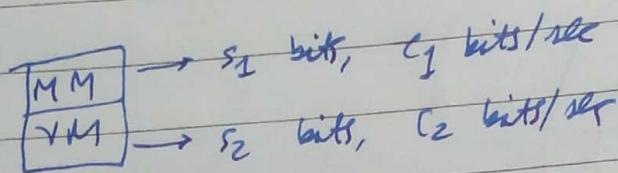
↓ ↓ ↓
seek latency data transfer
time time time

NOTE:

Data transfer + Data transfer time

serial Access Memory unitReasons of using virtual memory:

- To free programmer from the need to carry out storage allocation and to permit efficient sharing of memory space among different users.
- To make programs independent of the configuration and capacity of memory systems, used during their execution.
- High access rate and low cost can be achieved with a memory hierarchy.



$$\text{Average cost per bit} = \frac{c_1 s_1 + c_2 s_2}{s_1 + s_2}$$

If $s_2 \gg s_1$, average cost per bit $\approx c_2$

Memory hierarchy

$M_1, M_2, M_3, \dots, M_{N-1}, M_N$

↳ N -level hierarchy

cost increase



Access time increase



Size increase

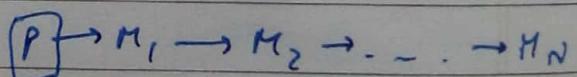
Data available in M_{i-1} is available also in M_i , but not the reverse.

Processor can access M_1 ,

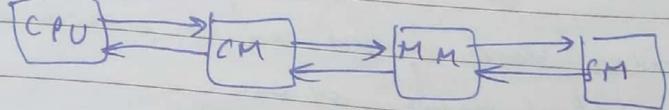
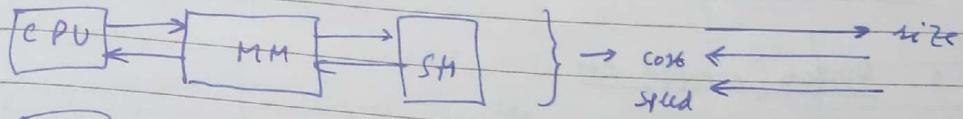
M_1 can access M_2

:

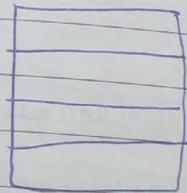
not on



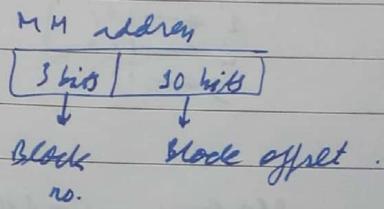
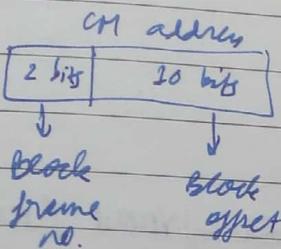
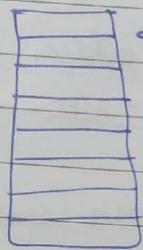
Cached Memory



CM having
4 block
frame



MM having
8 blocks



Processor generates MM address → converted to CM address → helps to access CM.

But CM smaller than MM

Processor generates 100 MM address.

Hit occurs for 80 accesses.

Miss occurs for 20 accesses.

$$H = \text{Hit ratio} = \frac{80}{100} \quad \text{Miss ratio} = \frac{20}{100}$$

No. of times CM accessed = 100

t_C = access time of CM

t_M = access time of MM

t_A = average access time

$$\text{Access time} = 100 t_C + 20 t_M$$

$$t_A = t_C + \frac{20 t_M}{5}$$

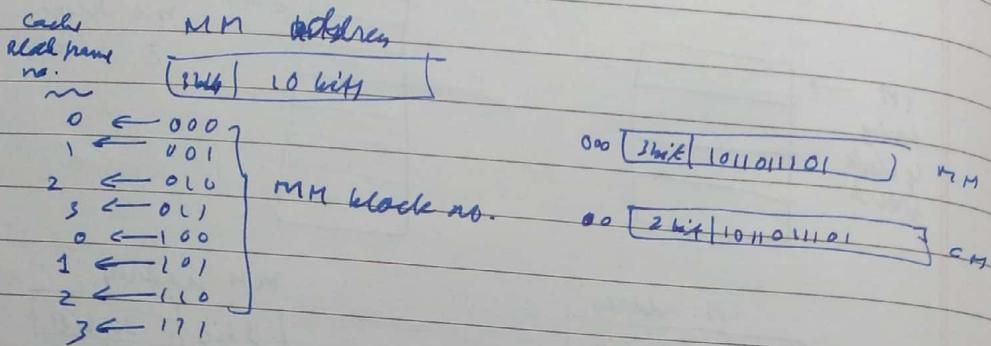
$$\text{Access time } t_A = \left(\frac{\text{No. of times CM accessed}}{\text{No. of times MM accessed}} \right) \times t_C + \left(\frac{\text{No. of times MM accessed}}{\text{No. of times MM accessed}} \right) \times t_M$$

$$t_A = \frac{\text{Access time}}{\left(\frac{\text{No. of times CM accessed}}{\text{No. of times MM accessed}} \right)}$$

$$\therefore t_A = t_C + (1-H) t_M$$

Address Mapping1. Direct mapped CM

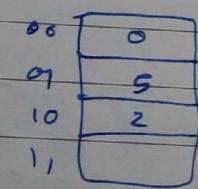
- one particular block of MM can be placed in one particular block frame of CM.



Block frame address = (Block address) Mod (No. of blocks in CM).

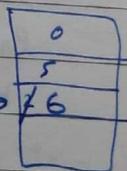
CPU generates request for the following blocks:

0, 2, 5, 6, 7, 1, 4



↑
Block Frame Address

← But 6 has to come here also. ∵ doing →



Initially CM empty

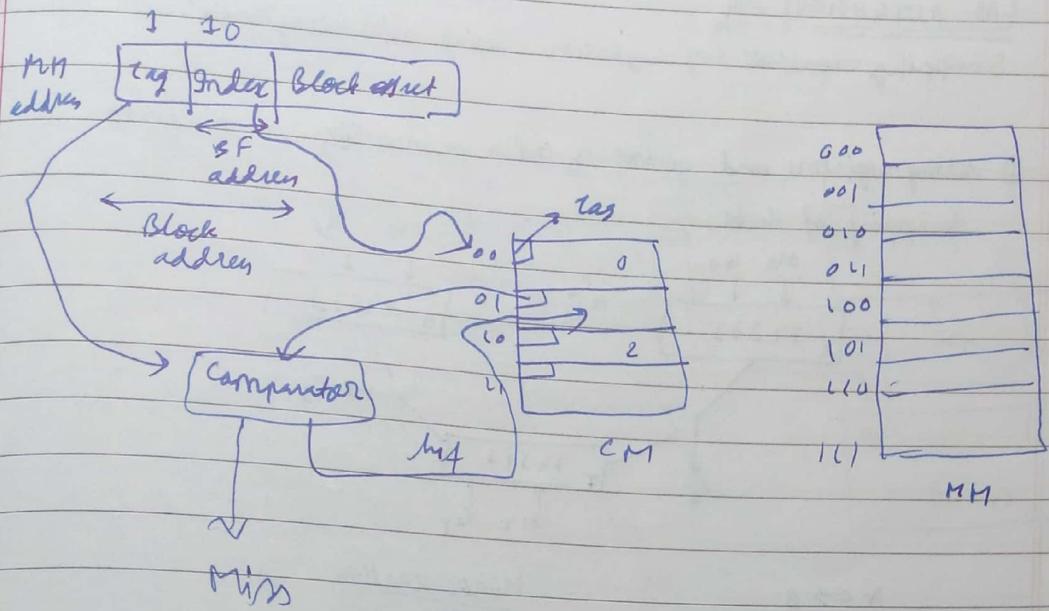
∴ we get direct miss and conflicting miss.

∴ Direct mapping is bad.

∴ 4 misses are compulsory. (for each BF address)

Then 3 conflict miss also occur.

∴ Total = 4 + 3 = 7 misses



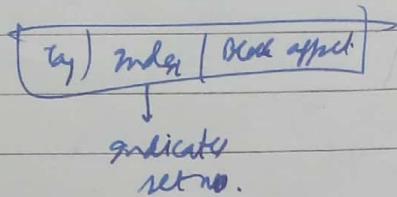
2 Fully modifiable mapping

3. Set Associative Mapping (SAM)

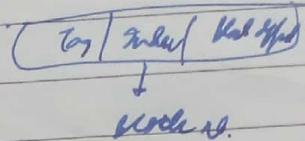
Associativity = No. of block frames in a set
 [for direct mapping, associativity = 1]

k -way set associative : A particular block can be placed in one of the k -block frames in a set.

For set associative mapping



For direct mapping



$$(\text{Set No.}) = (\text{Block no.}) \bmod (\text{No. of sets in memory})$$

Eg. sequence of memory block requests

i, 4, 3, 25, 8, 19, 5, 25, 8, 16, 35, 45, 22, 83,
 16, 25, 7

CM contains 8 block frames.

For DM

| | |
|---|----|
| 0 | 16 |
| 1 | 25 |
| 2 | |
| 3 | 25 |
| 4 | |
| 5 | 45 |
| 6 | 22 |
| 7 | |

19 → conflict miss
 25 → hit
 8 → hit
 5 → conflict miss
 35 → conflict miss

| | | | |
|---|----|----|----|
| 0 | 16 | 8 | 26 |
| 1 | 25 | | |
| 2 | | | |
| 3 | 25 | 35 | 3 |
| 4 | 4 | | |
| 5 | 45 | | |
| 6 | 22 | | |
| 7 | 7 | | |

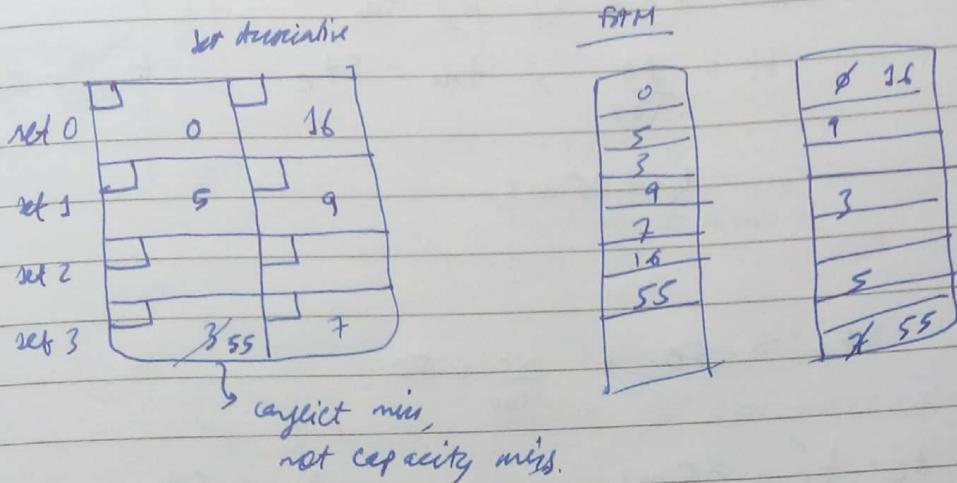
For FAM

| | | |
|---|----|----|
| 0 | 4 | 45 |
| 1 | 3 | 22 |
| 2 | 25 | |
| 3 | 8 | |
| 4 | 19 | 3 |
| 5 | 6 | 7 |
| 6 | 16 | |
| 7 | 35 | |

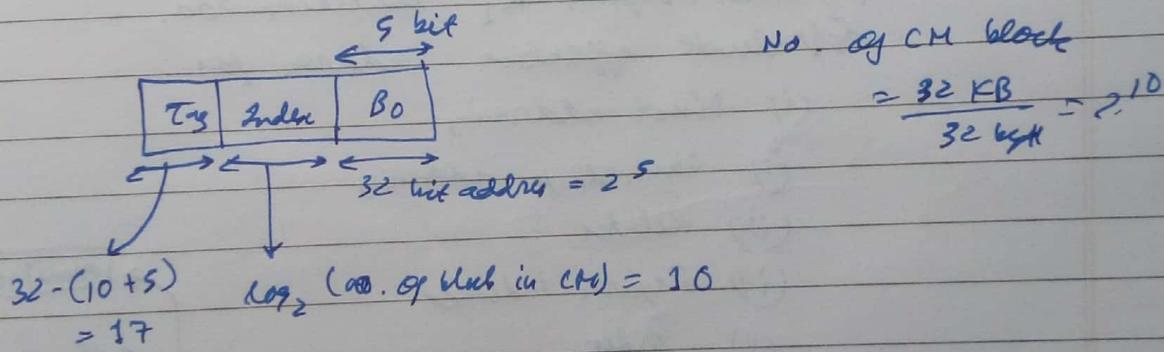
4 is LRU ~ capacity miss
 ∴ Remove 4, put 45.

LRU: least recently used.

Ex. 4-way set-associative memory with 4 sets and 8 cache block frames. 0-317 (128 blocks) in MM.
 LRU used to replace cache block frames.
 sequence of memory block requests
 $\rightarrow 0, 5, 3, 9, 7, 0, 16, 55$

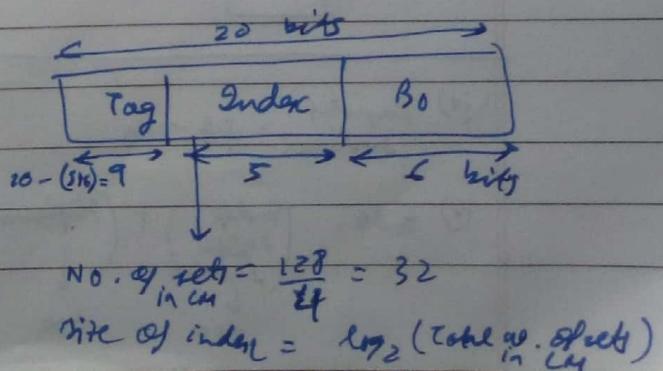


Cy. DM Cache of 127 32 KB.
 Block size = 32 bytes.
 CPU generates 32 bit address.
 Find no. of bits needed for cache indexing and tag bits



Ex. 4-way set associative, consisting of 128 blocks, cache
 CPU generates 22 bit address. Block size = 64 words,
 find no. of bits in tag, index and offset field.

No. of CM blocks = 128
 CM block size = 64 words
 MM address = 20 bits
 4 way SMT.



$$\text{No. of sets} = \frac{128}{4} = 32$$

$$\text{Size of index} = \log_2(\text{Total no. of sets})$$

Ex. A 2-level memory consisting of MM and CM.

CM is 5 times faster than MM. Hit ratio = 20%.

If average access time is increased by 20%.

From 50ns, what is approx. change in hit ratio?

$$t_A = t_C + (1 - H) t_M$$

$$H = \frac{20}{100}, t_M = 5 t_C, t_A = 50$$

$$\therefore t_C = 25 \text{ ms}$$

$$t_A' = 50 + \frac{20}{100} \times 50$$

$$t_C = t_C' = 25 \text{ ms}$$

$$t_M = 5 t_C$$

$$H = ? \text{ (idk)}$$

Ex. An instruction is stored at location 300

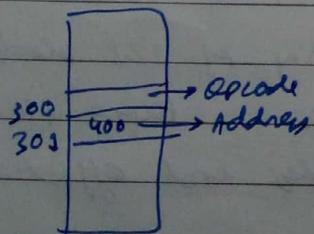
with its address field at location 301.

Address field value is 400.

Processor register R1 contains 200.

Evaluate effective address for

- (i) Direct addressing
- (ii) Immediate addressing
- (iii) Relative
- (iv) Register indirect
- (v) Index



(iv) Register indirect: EA = 200

(v) Index: (Index register content) + (Displacement) EA

(i) Direct: EA = 400

(ii) Immediate: EA = 301 \rightarrow [opcode | 400]

(iii) Relative:

Read from 300 \rightarrow PC = 300

Read from 301 \rightarrow PC = 301

\therefore Now \rightarrow PC = 302

Displacement = 400

EA = 302 + 400 = 702

Displacem

ent w.r.t base

base regis

ter (if)

not specifi

ed, we add

program

counter

Displacement

w.r.t base

base regis

ter (if)

not specifi

ed, we add

program

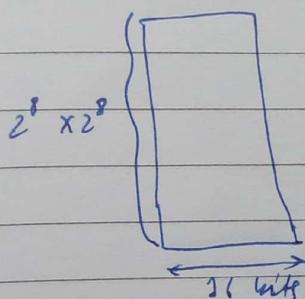
counter

Q. What should be size of ROM to store the multiplication table of 2 digit - unsigned integer?

| | | | |
|----|-------------|-----|---|
| | 1 2 3 4 ... | 10 | 2 ⁸ terms |
| 1 | 1 2 3 4 ... | 30 | |
| 2 | 2 4 6 8 ... | 20 | |
| 3 | | | (10x10 is just an example value to demonstrate it.) |
| 4 | | | |
| | 1 | | |
| | 1 | | |
| | 1 | | |
| | 1 | | |
| 10 | | 100 | product terms |

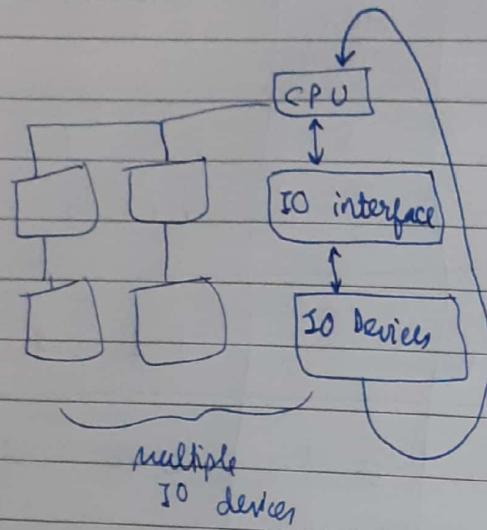
$$\text{Here, } \# \text{ of product terms} = 2^8 \times 2^8$$

$$\text{No. of memory locations in ROM needed} = 2^8 \times 2^8$$

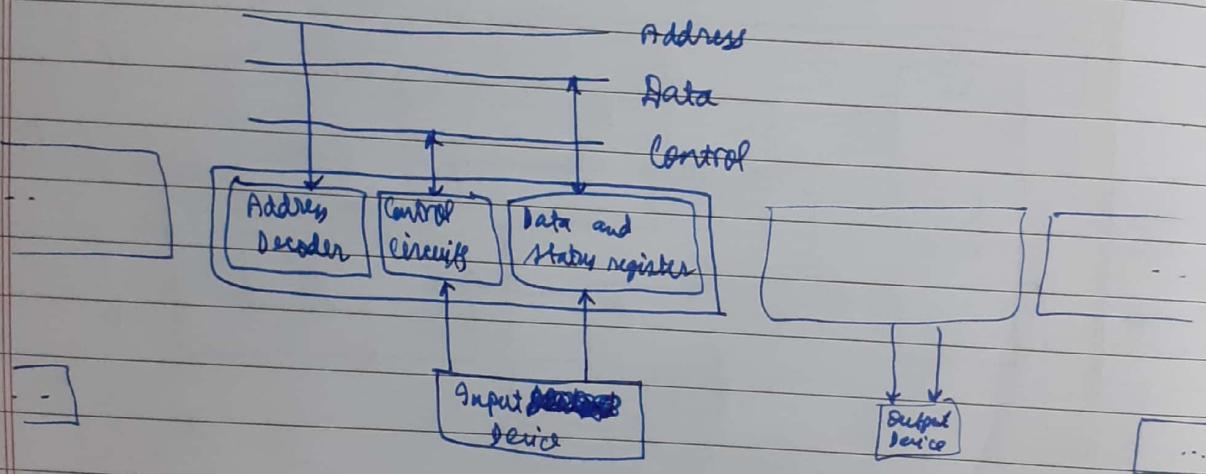


$$\therefore \text{ROM size} = 2^8 \times 2^8 \times 16$$

Peripherals (I/O devices)



each I/O device has its own address.



Status register

↓
Status → whether transfer of data b/w I/O device and processor is successful or not.

e.g. of control signals.

I/O R

I/O W

ice has its

Bus* Bus architecture

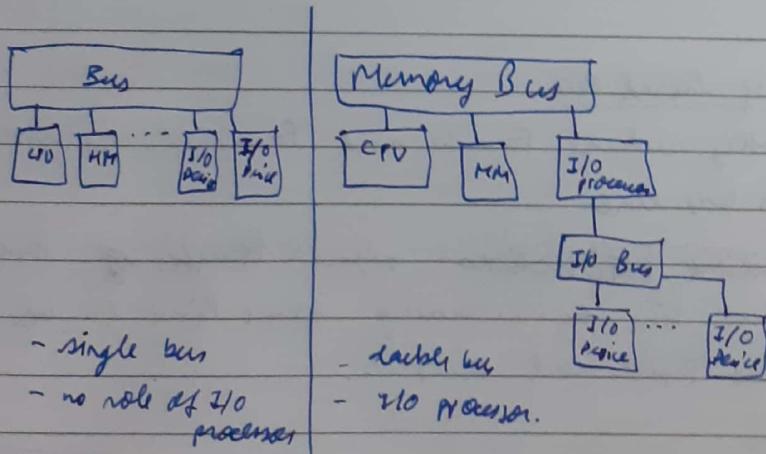
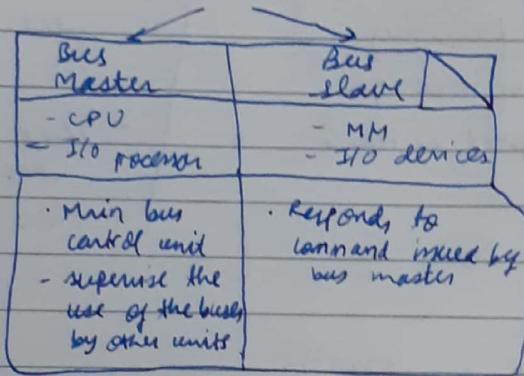
Bus interconnects several processor-level components.
(MM, I/O processor, CPU, I/O devices)

* Bus control

controls 2 things

- time of transferring data via bus
- selection process by which a device can gain control of the bus.

→ Devices are divided into 2 categories:



* Bus communication

Asynchronous

- Sending and receiving ends work at some speed.
This is possible if sender and receiver are close-by, as they can then use same clock pulse.
- For larger separation, different clocks have to work at same frequency
- speed of operation
 - = speed of clock

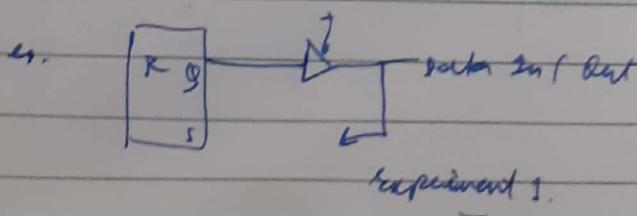
Synchronous

- "Handshaking signals"
- analog-digital converter
- sending and receiving end exchange handshaking signals, before starting exact communication.
- speed of operation = speed of ^{transf. device} clock
- Used for long-distance communication
- More complex bus control circuit

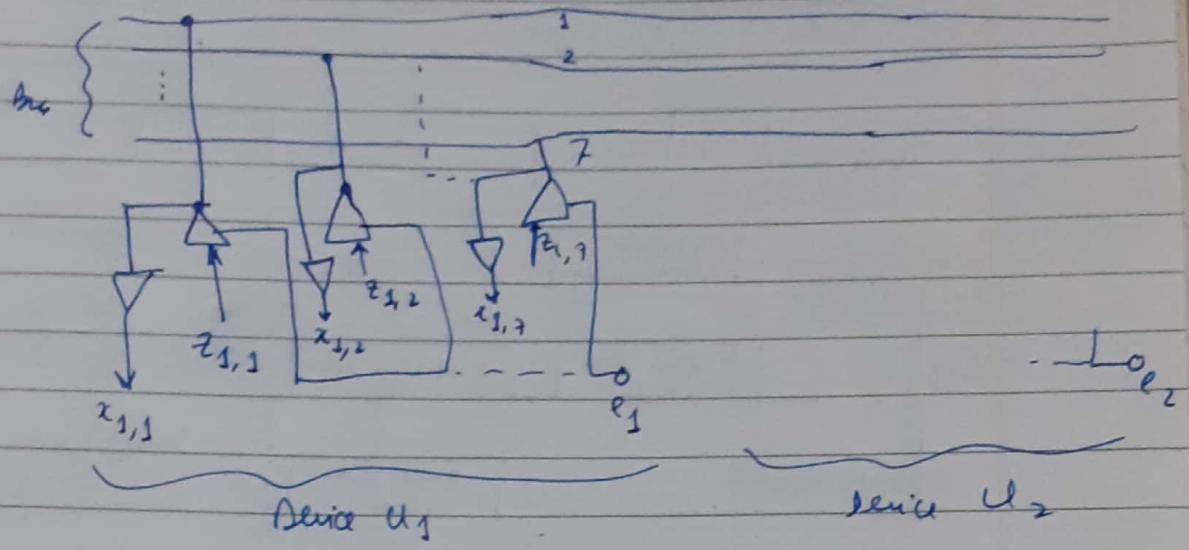
* Shared Bus

Advantage of shared bus

- Greatly reduce fan-in and fan-out constraints on bus lines
- Facilitate bidirectional signal transfer over a bus line by allowing same line to be used as input and output.



Tristate used here provides similar functions as shared bus.



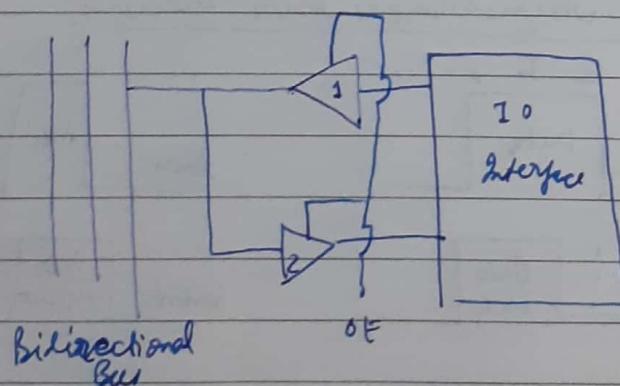
$[e_1 = 1, e_2 = 0]$: $U_1 \rightarrow \text{write}$
 $U_2 \rightarrow \text{read from buffer}$

$[e_1 = 0 \text{ (read)}, e_2 = 1 \text{ (write)}]$

$[e_1 = 1, e_2 = 1]$: not allowed \rightarrow can't write simultaneously

$[e_1 = 0, e_2 = 0]$: allowed \rightarrow can read simultaneously

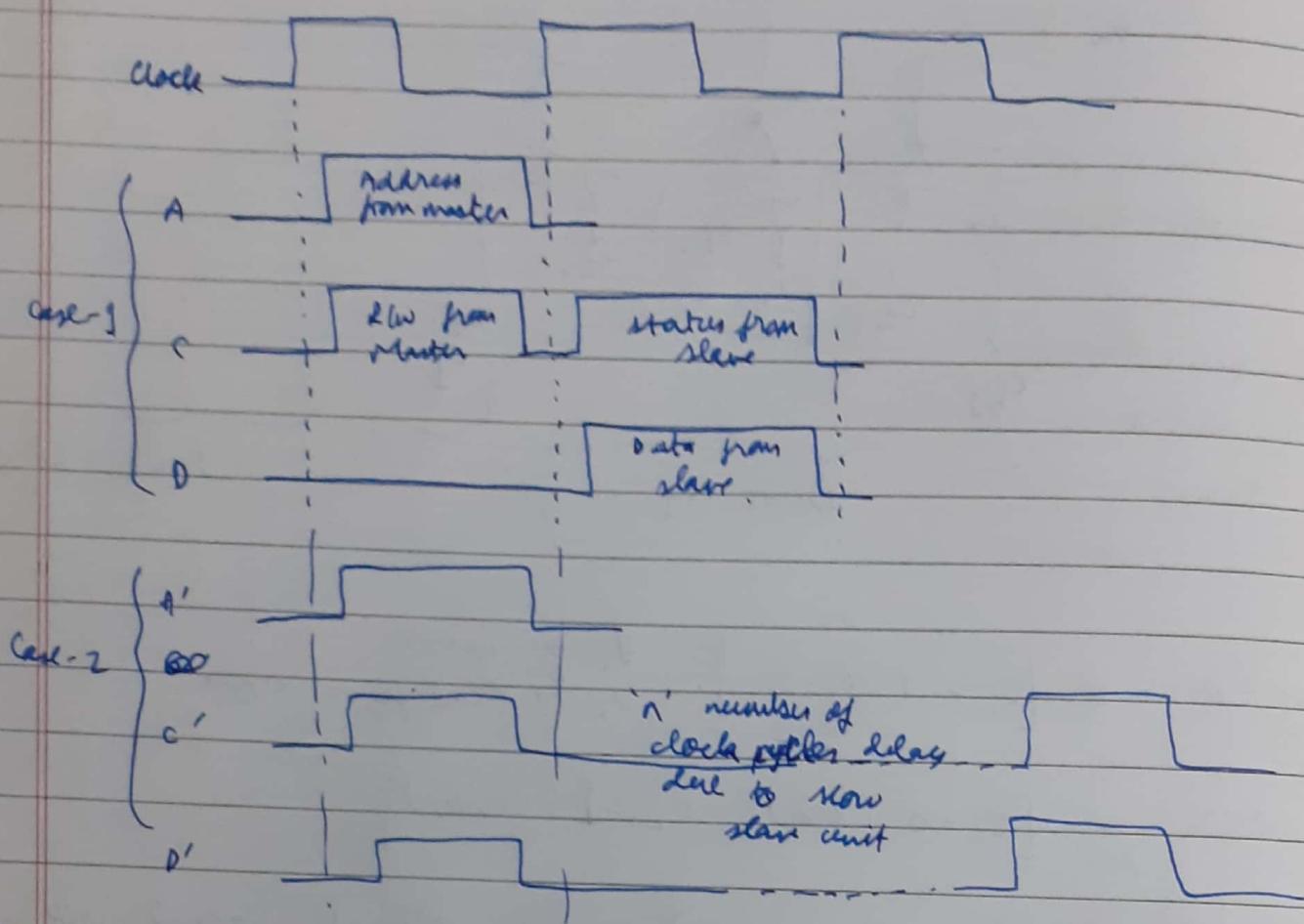
* Transceiver



1 = Driver ($OE = 1$)
2 = Receiver ($OE = 0$)

Driver drives the bus while maintaining its output level within specified limit

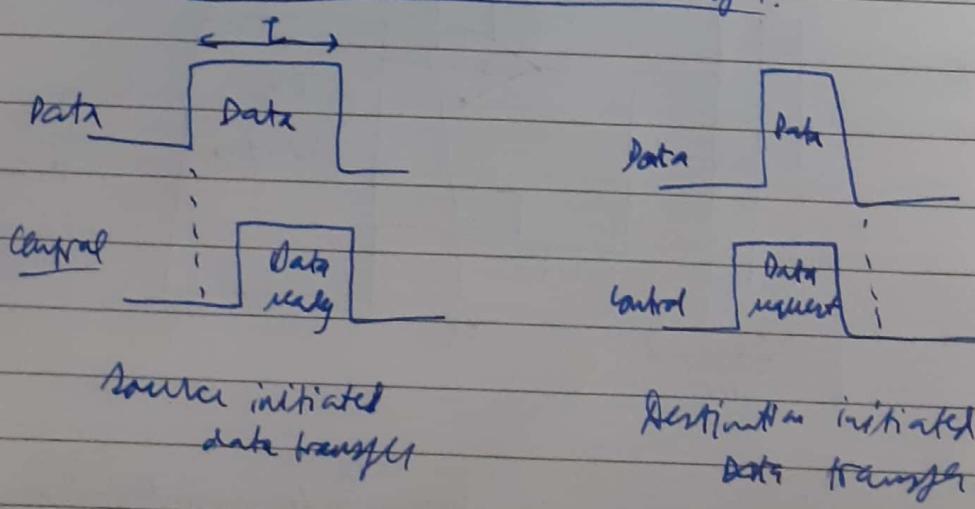
* Timing diagram for synchronous data transfer



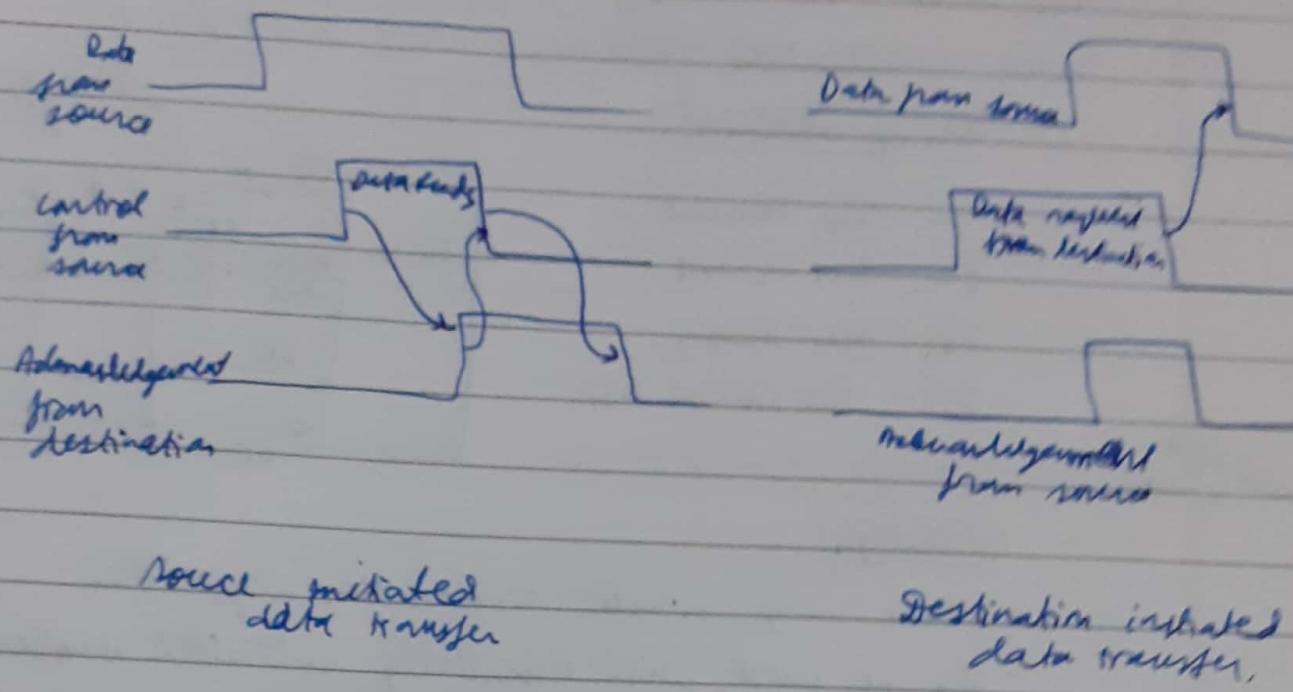
* Asynchronous data transfer

- Used by device having widely different data rates.

→ One-way asynchronous data transfer.



→ 2-way synchronous data transfer



Bus Arbitration

Arbitration is the process where

Next device to become bus master is selected

Three arbitration schemes

1) Daisy chaining

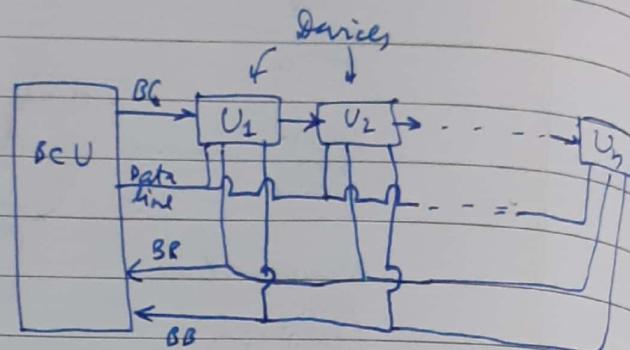
BR: Bus Request

BG: Bus Grant

BB: Bus Busy

BCU: Bus Control Unit

* Control lines = 2ⁿ⁺¹



Advantages:

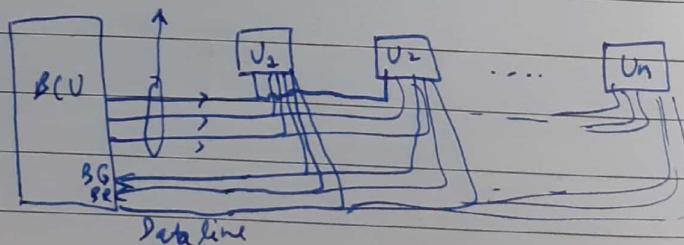
- Needs very simple arbitration algorithm
- Can be used with essentially unlimited number of bus units.
- It has 2 control lines (other than BB)

Difficulty:

- Failure of 1 bus unit affects the entire scheme
- Priority is decided in terms of location of device w.r.t BCU.
- Now. (BG for U_n has to pass through $U_1 \rightarrow U_2 \rightarrow U_3 \rightarrow \dots \rightarrow U_{n-1}$)

2) Polling

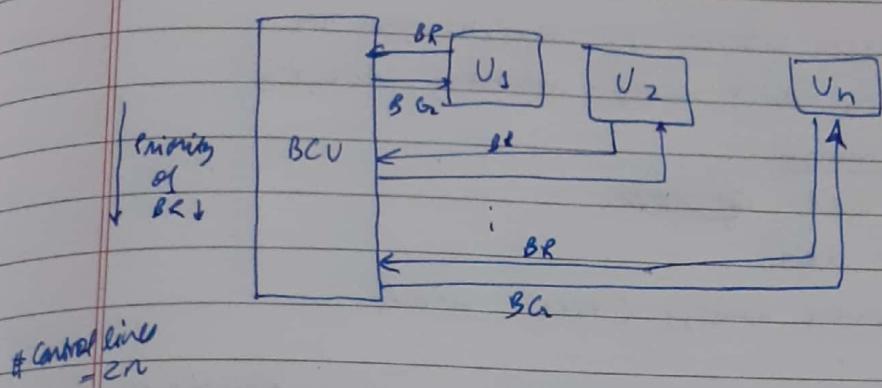
3 poll control lines



"Polling sequence" \rightarrow e.g. 5, 4, 3, 1, 2, 6

'from most priority'

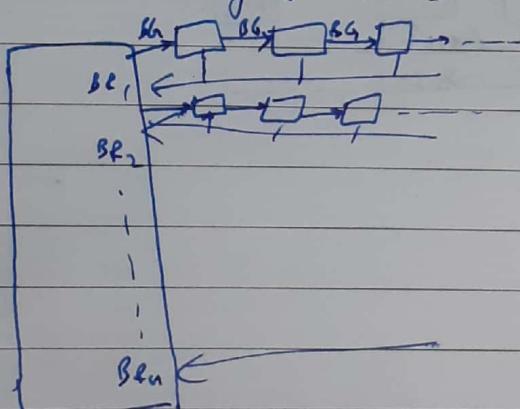
3) Independent Requesting



$$\text{No. of control lines} = 2 \times (\text{No. of devices})$$

4) Mixing system

- combination of 1 and 3



m - BR lines

n - devices

$$m < n$$

priority: $BR_1 > BR_2 > BR_3 > \dots$

\therefore higher priority devices
should be connected
with higher priority
BR. lines.

BG links are
arranged as in
daisy chaining

I/O operation

Data transfer between I/O devices and RM.

→ I/O device can't access memory directly

Mode of data transfer:

→ Programmed mode ~~(DMA)~~

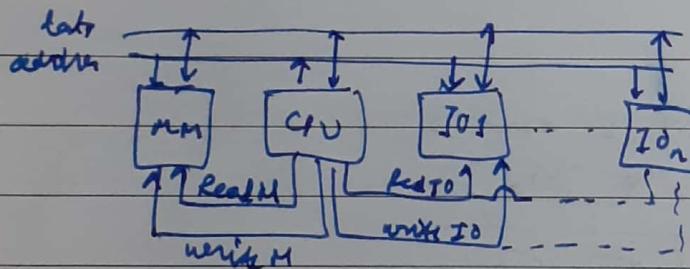
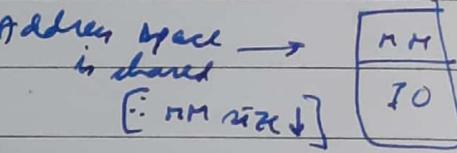
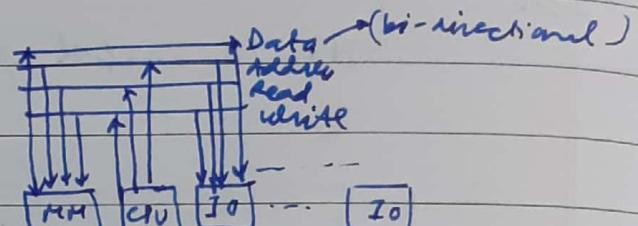
→ Interrupt mode

→ Direct memory access (DMA) mode

Programmed Mode

→ memory mapped I/O

→ I/O mapped I/O



- If address from RM part, memory is referenced.
Otherwise, I/O is referenced.

After this, absent for 1 class (6 Pages)

classmate

Date _____

Page _____

RAID

- Redundant array of inexpensive disks
 - ↳ storage system based on multiple disks
 - 7 different configurations (or RAID levels)
 - ↳ RAID-0 to RAID-6
 - ↳ improving speed.

RAID-1 (Mirroring)

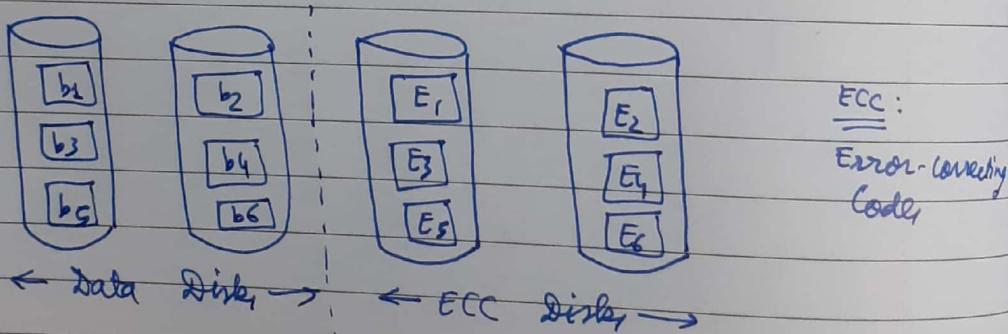
N = number of data disks

N = number of redundant disks

$$\boxed{2N}$$

RAID-2

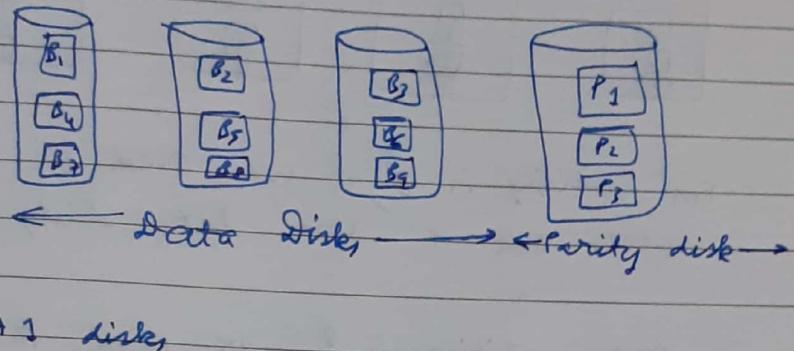
- strips the bits of a file across the disks.



- uses lot of disk
- need to maintain synchronization of all these disks, which is expensive.
- random reading is fair, good for sequential reading and fair for sequential writing.
- single-bit error can be detected and corrected.
- multiple errors can be detected, but cannot be corrected.

RAID-3

- strips the bytes of a file across the disk.

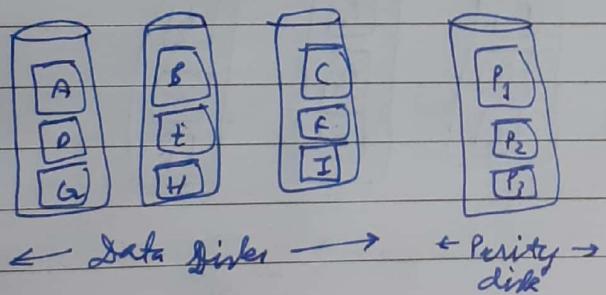


$N+1$ disks,

- high read speed
- low write speed
- If 2 disks are damaged, data cannot be retrieved simultaneously

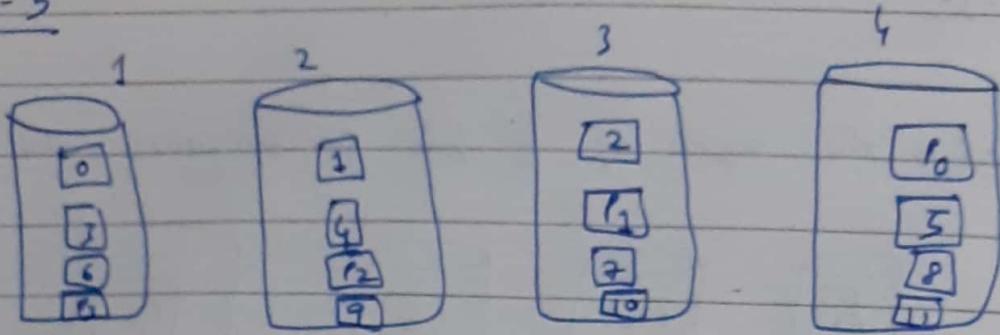
RAID-4

- uses Block-level stripping



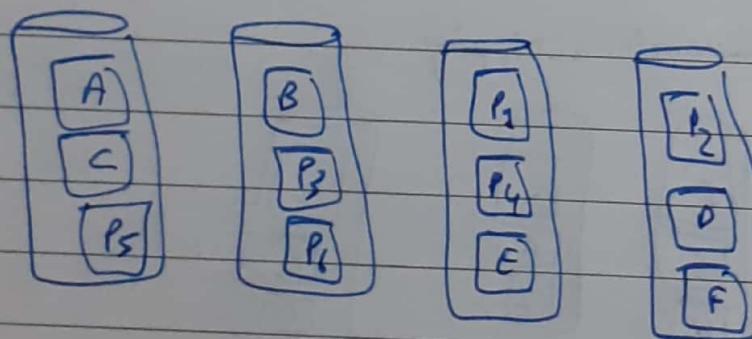
- needs hardware support for parity calculation
- write operation is slower due to parity calculation
- good for sequential read/write data.
- Good random read, bad random write.
- For each I/O write, calculation of parity block is needed.

RAID - 5



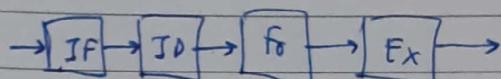
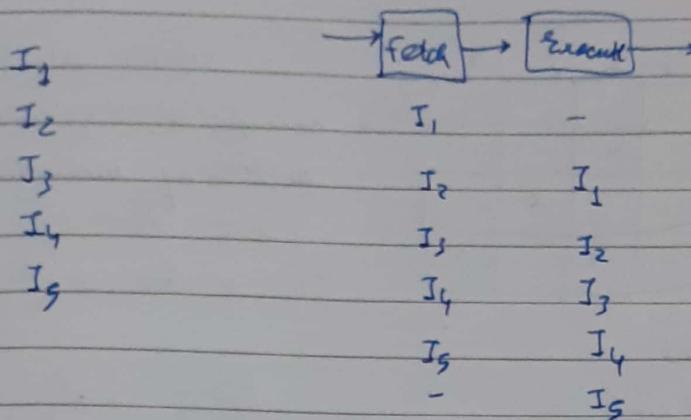
- For writing in 3, we need to access 1st & 3rd data disk,
- For writing in 8, we need to access 4 and 2.
- All the disks can be accessed simultaneously.
- But block 1 and 3 cannot be accessed simultaneously.

RAID - 6



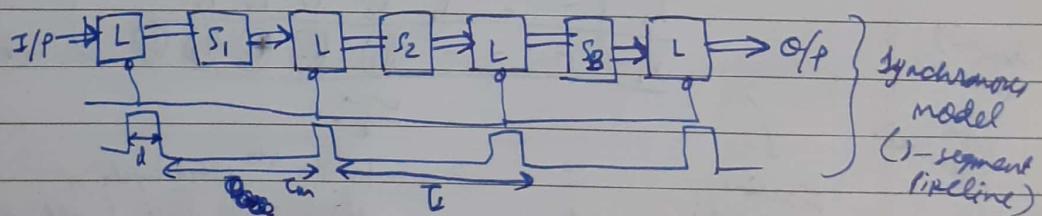
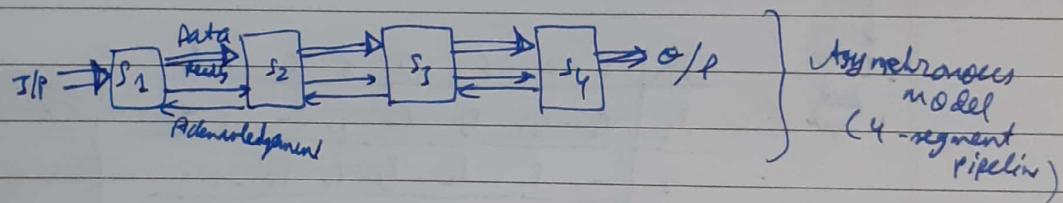
-
-

Pipelining



* Linear pipeline processor

- cascade of stages which are linearly connected to perform a fixed function over a string of data moving from one end to the other.



$L \rightarrow$ clocked master-slave
flip-flop.

$S_1, S_2, S_3 \rightarrow$ combinational
circuits

$$f = \text{clock frequency} = \frac{1}{T}$$

$d = \text{latency}$

$T = \text{clock time}$

$$T = d + C_m$$

$$C_m = \max \left\{ \frac{t_1}{t_2}, \frac{t_2}{t_3} \right\}$$

t_1 : time taken by S_1
for its task

t_2 : time taken by S_2

t_3 : time taken by S_3

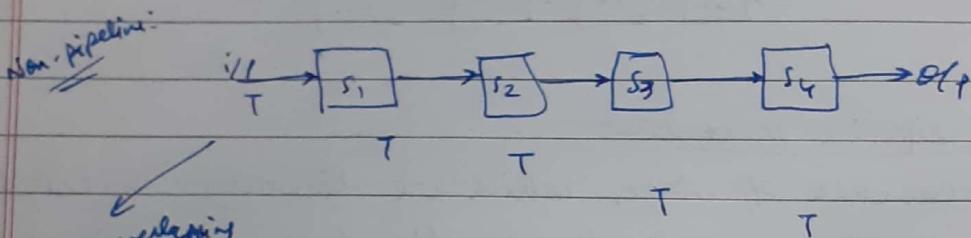
* Reservation table for a 4-segment pipeline:

| | 1 | 2 | 3 | 4 |
|-------|---|---|---|---|
| S_1 | X | | | |
| S_2 | | X | | |
| S_3 | | | X | |
| S_4 | | | | X |

→ specifies the utilization pattern of successive stages

S_j = segment

cc = clock cycle



∴ 4 clock cycles needed for task T.

↓ Task, 4 task cycle

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|-------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|-----|
| S_1 | T ₁ | T ₂ | T ₃ | T ₄ | T ₅ | - | - | - |
| S_2 | - | T ₁ | T ₂ | T ₃ | T ₄ | T ₅ | - | - |
| S_3 | - | - | T ₁ | T ₂ | T ₃ | T ₄ | T ₅ | - |
| S_4 | - | - | - | (1) | (2) | (3) | (4) | (5) |

⇒ Scan-Time Diagram

→ When a task reaches segment S_4 , it is completed

in the 4th clock cycle,
pipeline is full (all segments have some task)

↓
pipeline concept:
5 tasks, 4 segment pipeline,

1 clock cycle

$$\text{No. of clock cycles needed} = k + (n-1)$$

↓ ↓

no. of segments no. of tasks

$$4 + (5-1) = 8$$

$$4 + 1 + 1 + 1 + 1 = 8$$

* Speedup of a k -segment pipeline (S_k)

• Non-pipeline

processor :



1 task $\rightarrow kT$

n tasks $\rightarrow n k T$

clock-time = T

• Pipeline-processor



n tasks $\rightarrow \{k + (n-1) \bullet\} T$

$$\cdot \boxed{S_k = \frac{[n k]T}{[k + (n-1)]T}} \quad \begin{array}{l} \xrightarrow{\quad} S_{k \max} = k, \text{ for } n \rightarrow \infty \\ \xrightarrow{\quad} S_{k \min} = 1, \text{ for } n=1 \end{array}$$

* Efficiency of a k -segment pipeline (E_k)

$$\cdot \boxed{E_k = \frac{S_k}{k} = \frac{1}{k+n-1}} \quad \begin{array}{l} \xrightarrow{\quad} E_{k \max} = 1, \text{ when } n \rightarrow \infty \\ \xrightarrow{\quad} E_{k \min} = \frac{1}{k}, \text{ when } n=1 \end{array}$$

* Throughput of a k -segment pipeline (H_k)

→ output per unit time

• $[k + (n-1) \bullet]T$ time is required for evaluating n -tasks

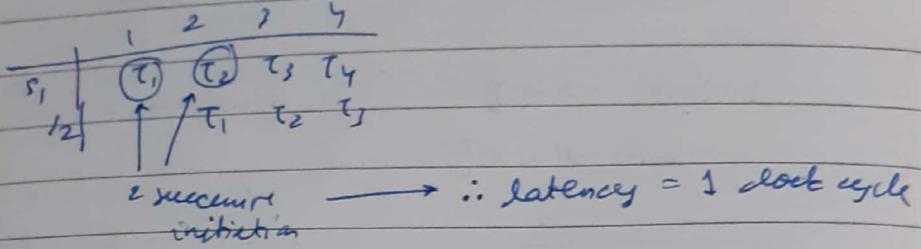
$$\cdot \boxed{H_k = \frac{n}{[k + (n-1)]T} = \frac{E_k}{T} = E_k f}$$

Relation b/w H_k , E_k and S_k :

$$\boxed{H_k = E_k f = \frac{S_k}{kT}}$$

* Pipeline latency

↳ delay that occurs w/o 2 successive initiation
in a pipeline



→ for static pipeline, pipeline latency = 1

* Pipeline Performance depends on k.

- If we increase the number of segments,
we feel that time needed \downarrow
- But k can't be increased too much.

- t = total time required for a nonpipeline sequential program.
to execute the same program on a k-stage pipeline,
with an equal flow of delay d .

$$\text{clock period} = \left(\frac{t}{k} \right) + d$$

\downarrow time taken to execute a task per segment

↳ latch delay.

$$\text{Throughput} = \frac{1}{\left(\frac{t}{k} + d \right)}$$

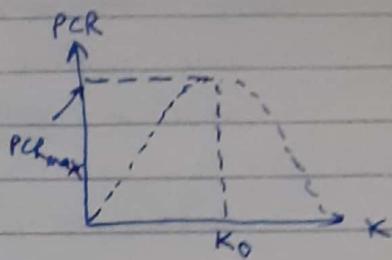
$$\text{Pipeline cost} = c + kh$$

\downarrow cost of each latch

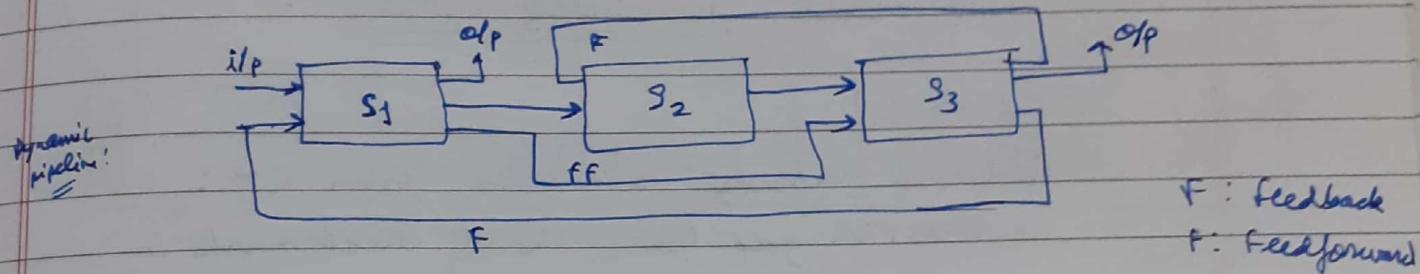
cost of logic stages

$$\text{performance cost ratio} = \text{PCR} = \frac{\text{throughput}}{\text{cost}} = \frac{\frac{1}{(\frac{t}{K} + d)(r + kh)}}{}$$

$$K_0 = \text{optimal cost} = \sqrt{\frac{tc}{dh}}$$



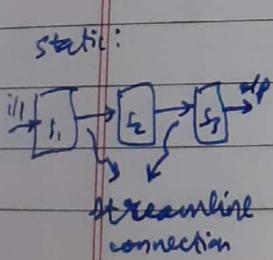
* Nonlinear pipeline



F: feedback
F: feedforward

| s^{cc} | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|----------|---|---|---|---|---|---|---|---|
| s_1 | x | | | | x | | x | |
| s_2 | | x | x | . | | | | |
| s_3 | | | x | x | x | | | |

RT for function x



| s^{cc} | 1 | 2 | 3 | 4 | 5 | 6 |
|----------|---|---|---|---|---|---|
| s_1 | y | | | | y | |
| s_2 | | y | y | y | y | y |

RT for function y

- * static pipeline can evaluate a single function, whereas a dynamic pipeline can evaluate multiple functions using the same pipeline

- In a static pipeline, it is easy to partition a given function into a sequence of linearly ordered sub-functions which is not possible in dynamic pipeline due to presence of F and FF connections

* Latency Analysis:

- any attempt by 2 or more initiations to use the same pipeline stage at the same time would cause a collision.
- latency for which collision occurs is called forbidden latency
- latency for which there is no collision is called permissible latency.
- forbidden latency calculation from RT:
 → Time gap b/w the 2 checkmarks in the same row in the reservation table

$$\begin{aligned} FL_x &= 5, 2, 7, 4 \\ PL_x &= 1, 3, 6, 8 \end{aligned} \quad \left. \begin{array}{l} \text{for function } x \\ (\text{see previous page}) \end{array} \right\}$$

$$FL_x = 2, 4$$

Diagram to show that if latency = 2, collision would occur for function X.

| s ^{cc} | 2 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|----------------------------|----------------|----------------|----------------|-------------------------------|-------------------------------|-------------------------------|---|----------------|----------------|----|
| s ₁ | x ₁ | | x ₂ | x ₃ | x ₁ | | x ₁ x ₂ | | x ₂ | |
| s ₂ | | x ₁ | | x ₁ x ₂ | | x ₂ x ₃ | | | | |
| s ₃ | | | x ₁ | | x ₁ x ₂ | | x ₁ x ₂ x ₃ | x ₂ | | |

* Latency sequence:

- it sequence of permissible non-forbidden latencies, b/w successive task initiation

* Latency cycle:

- latency sequence which repeats the same subsequence (cycle) indefinitely.

for function x

| Latency cycle (1, 2) $\rightarrow 1, 2, 1, 2, 1, 2, \dots$ | | | | | | | | | | | | | | | | | |
|--|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 |
| S ₁ | x ₁ | x ₂ | | | | x ₁ | x ₂ | x ₁ | x ₂ | x ₃ | x ₄ | | | x ₃ | x ₄ | x ₃ | x ₄ |
| S ₂ | | x ₁ | x ₂ | x ₁ | x ₂ | | | | | x ₃ | x ₄ | x ₃ | x ₄ | | | | |
| S ₃ | | | x ₁ | x ₂ | x ₁ | x ₂ | x ₁ | x ₂ | | x ₃ | x ₄ | x ₃ | x ₄ | x ₃ | x ₄ | | |

← cycle (x₁, x₂) → ← cycle (x₃, x₄) →

$$\text{Average latency} = \frac{1+2}{2} = 4.5$$

Latency cycle (3) $\rightarrow 3, 3, 3, \dots$

| cc | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | | | | |
|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|--|
| S ₁ | x ₁ | | | x ₂ | | | x ₁ | x ₂ | x ₁ | x ₂ | x ₃ | | | x ₂ | x ₃ | x ₂ | x ₃ | |
| S ₂ | | x ₁ | | | x ₁ | x ₂ | | x ₂ | x ₁ | x ₂ | | x ₁ | | | | | | |
| S ₃ | | | x ₁ | | x ₁ | x ₂ | x ₁ | x ₂ | x ₃ | x ₂ | x ₃ | | x ₂ | x ₃ | | | | |

← cycle ← cycle →

$$\text{constant latency} = 3$$

Collision-free scheduling

When scheduling events in a pipeline, the main objective is to obtain the shortest average latency b/w initiations without causing collisions.

Collision-vector

For a reservation table having i columns

max forbidden latency = m

$$m \leq n-1$$

e.g. for function X , $n=8$, $m=7$

Permissible latency (p) should be as small as possible
 $1 \leq p \leq m-1$

e.g. for function X , $1 \leq p \leq 6$

$p=1 \rightarrow$ static pipeline

collision vector (CV) = $c_m c_{m-1} \dots c_1$

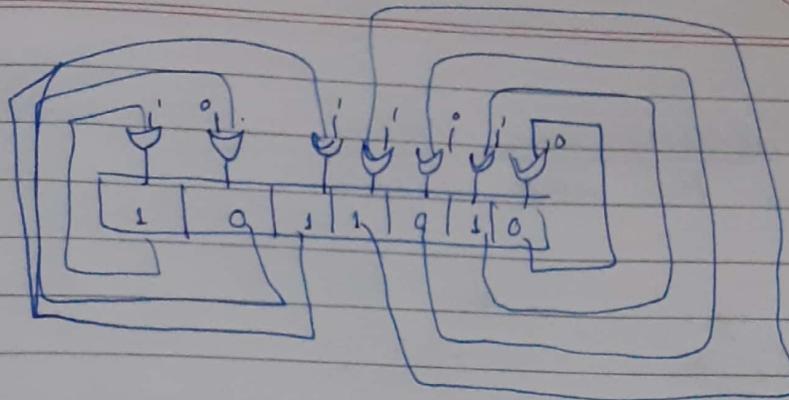
c_i = i th bit in CV

$c_i = 1$ if i is a FL

$c_i = 0$ if i is not a FL.

e.g., for X , $c_7 c_6 c_5 c_4 c_3 c_2 c_1$ \oplus
 $1 \quad 0 \quad 1 \quad 1 \quad 0 \quad 1 \quad 0$

For X , $CV = 1010$



1 unsafe
0 safe

0 emerges for PL
1 safe

1 emerges for PL
0 unsafe -

[101101]

$\rightarrow 011101 \rightarrow 0$ } 1st shift
 $\rightarrow 0010110 \rightarrow 1$ } 2nd shift

1011010

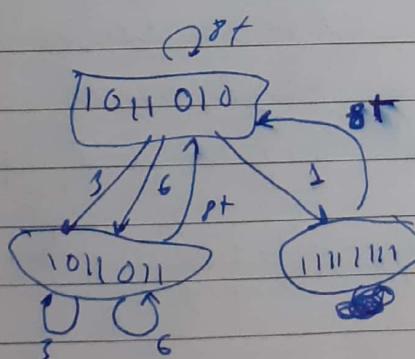
$$\begin{array}{r} + 0101101 \\ \hline 1111111 \end{array}$$

op:

→ in register
[1011010]
shift 3 times to set
nxt

0001011

$$\begin{array}{r} 1011010 \\ \rightarrow 0001011 \\ \hline 1011011 \end{array}$$



state diagram
for function X

were all points
make transition that
avoid collision.

min latency = 3.
(average)

are infinitely many latency cycles, one can trace from the state diagram.

For eg.:

(1, 8)

(1, 8, 6, 8)

(1, 3, 8)

(1, 8, 6)

(3, 8)

(7, 6, 3), ...

Simple cycle is a latency cycle in which ~~some~~ each state appears only once.

↳ ex. $\underbrace{(1, 8)}_{2,5}, \underbrace{(3)}_3, \underbrace{(6)}_6, \underbrace{(3, 8)}_{5,5}, \underbrace{(8)}_8, \underbrace{(6, 8)}_7$

Not simple cycle

↳ ex. (1, 8, 6, 8) → traverses 10, 1, 8, 10 twice

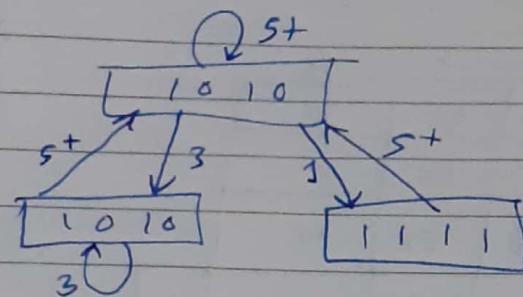
Simple cycle → Average latency (eg. 4.5, 5.5, 7)
→ constant latency (eg., 3, 6, 1)

Greedy cycle = simple cycle with average latency lower than other simple cycles.

Ex. $\underbrace{(1, 8)}_{4,5}$ and $\underbrace{(3)}_3$ are greedy
↳ least avg. ↳ least constant

minimum average latency = 3

4. State diagram $D_1 \gamma$.



$\{1, s\}$ and $\{3\}$ are
greedy sets.

CPI

↳ clock per instruction

MIPS

↳ million of instruction per second
→ helps to find out about speed of processor.

CPI

I_c = Number of instructions in a program.

T = CPU time to execute I_c no. of instruction

J = clock time

$$T = I_c \times CPI \times J$$

MIPS

c = total no. of clock cycles required to execute
 I_c no. of instructions

$$CPI = \frac{c}{I_c}$$

I_c : total no. of instructions

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

T : Time to execute I_c no. of instructions.

$$= \frac{I_c}{T \times CPI \times 10^6}$$

$$\text{Instructions per second} = \frac{I_c}{T}$$

$$= \frac{4}{CPI \times T \times 10^6}$$

$$= \frac{6}{CPI \times 10^6}$$

$$= \frac{6 I_c}{C \times 10^6}$$

Pipeline in superscalar processor

$IF =$ instruction fetch
 $ID =$ instruction decode
 $EX =$ execute
 $W =$ register write-back

| clock cycle | 1 | 2 | 3 | 4 |
|-------------|----|----|----|----|
| | IF | ID | EX | W |
| | IF | ID | EX | W |
| | IF | ID | EX | W |
| | | IF | ID | EX |
| | | IF | ID | EX |
| | | IF | ID | EX |
| | | IG | ID | EX |
| | | IG | ID | EX |

- uses multiple independent instruction pipeline having multiple stages.
- fetches multiple instructions at a time. Then attempts to find ~~nearby~~ independent instructions, that are independent of one another and can therefore be executed in parallel.

Superscalar processor of degree m

$m =$ no. of pipelines in the processor

Here $m = 3$.

Consider $m = 2, n = m$. of tasks = 4

Assume:

All tasks are independent (for this semester)

| cc | 1 | 2 | 3 | 4 | 5 | 6 |
|-------|-------|-------|-------|-------|-------|-------------------------|
| S_1 | T_1 | T_3 | | | | 5 clock cycles required |
| S_1 | T_2 | T_4 | | | | |
| S_2 | | T_1 | T_3 | | | |
| S_2 | | T_2 | T_4 | | | |
| S_3 | | | T_1 | T_3 | | |
| S_3 | | | T_2 | T_4 | | |
| S_4 | | | | T_1 | T_3 | |
| S_4 | | | | T_2 | T_4 | |

S_i : segment of pipeline

for base processor, $m = 2$, 4 tasks

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|-------|----------|----------|----------|----------|----------|----------|----------|
| s_1 | τ_1 | τ_2 | τ_3 | τ_4 | | | |
| s_2 | | τ_1 | τ_2 | τ_3 | τ_4 | | |
| s_3 | | | τ_1 | τ_2 | τ_3 | τ_4 | |
| s_4 | | | | τ_1 | τ_2 | τ_3 | τ_4 |

$$\text{No. of clock cycles} = 7$$

If $k = \text{no. of segments}$

$N = \text{no. of tasks}$

$m = \text{degree of processor}$.

No. of clock cycles

For base processor $\rightarrow k + 1(N - 1)$

For superscalar processor $\rightarrow k + \frac{1}{m}(N - m)$

Speedup

$T(1, 1) \rightarrow 1 \text{ instruction in 1 clock cycle}$

$$\rightarrow OPT = 1$$

$$T(m, 1) = k + 1(N - m)$$

$$T(m, 1) = k + \frac{1}{m}(N - m)$$

$T(m, 1) \rightarrow m \text{ instructions in } 2 \text{ clock cycles}$

$$\rightarrow OPT = \frac{1}{m}$$

~~Time per task~~

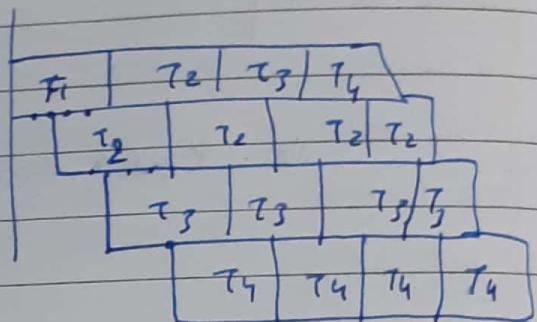
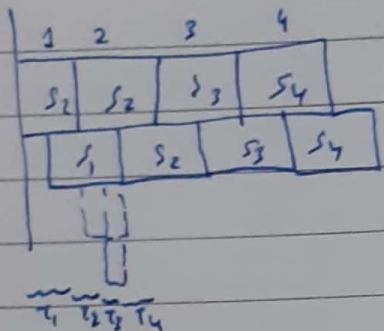
$$S(m, 1) = \frac{T(1, 1)}{T(m, 1)}$$

$$S(m, 1) = \frac{k + N - 1}{k + \frac{1}{m}(N - m)} = \frac{m(k + N - 1)}{mk + N - m}$$

$$N \rightarrow \infty, S(m, 1) \rightarrow m$$

Pipeline in superscalar processor

Exploits the fact that many pipeline stage perform tasks that require less than half a clock cycle.



$$CIS = \frac{1}{m}$$

degree = m

$\Rightarrow \left(\frac{1}{m}\right)$ parts of a clock cycle
required for completing
a task

$\tau(s, m) \rightarrow$ instruction is executed in
 $\frac{1}{m}$ clock cycle for superscalar processor
of degree m .

$$\tau(s, m) = k + \frac{1}{m}(N-1)$$

$$S(s, m) = \frac{k + N - 1}{\frac{k+1}{m}(N-1)} = \frac{m(k+N-1)}{mk+N-1}$$

RISC

- Relatively few instructions & few addressing modes.
- Memory access related to load and store instruction.
- All operations done within CPU register.
- Fixed length easily decoded instruction format.
- Hardwired control (not microprogrammed).

CISC

- Large no. of instructions typically from 100 to 250.
- Some instructions that perform specialized tasks are used infrequently.
- Large variety of addressing modes.
- Variable length instruction format.

Q) How many 28x1 RAM chips are needed to design a memory of capacity 2048 bytes? What is the size of the address bus to access 2048 byte of memory? How many of these lines will be common to all chips? How many lines must be decoded for chip select? Specify size of decoder!

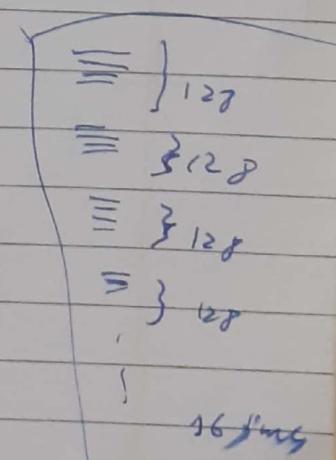
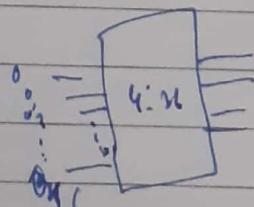
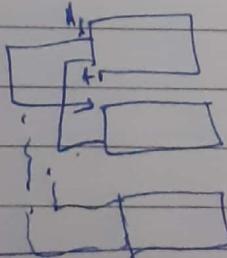
$$\text{Required memory} \rightarrow 2048 \text{ bytes} = 2048 \times 1$$

$$\text{Available memory} \rightarrow 128 \text{ bytes} = 128 \times 1$$

$$\therefore \text{No. of chips} = 16$$

$$\therefore \text{No. of lines in address bus} = 11$$

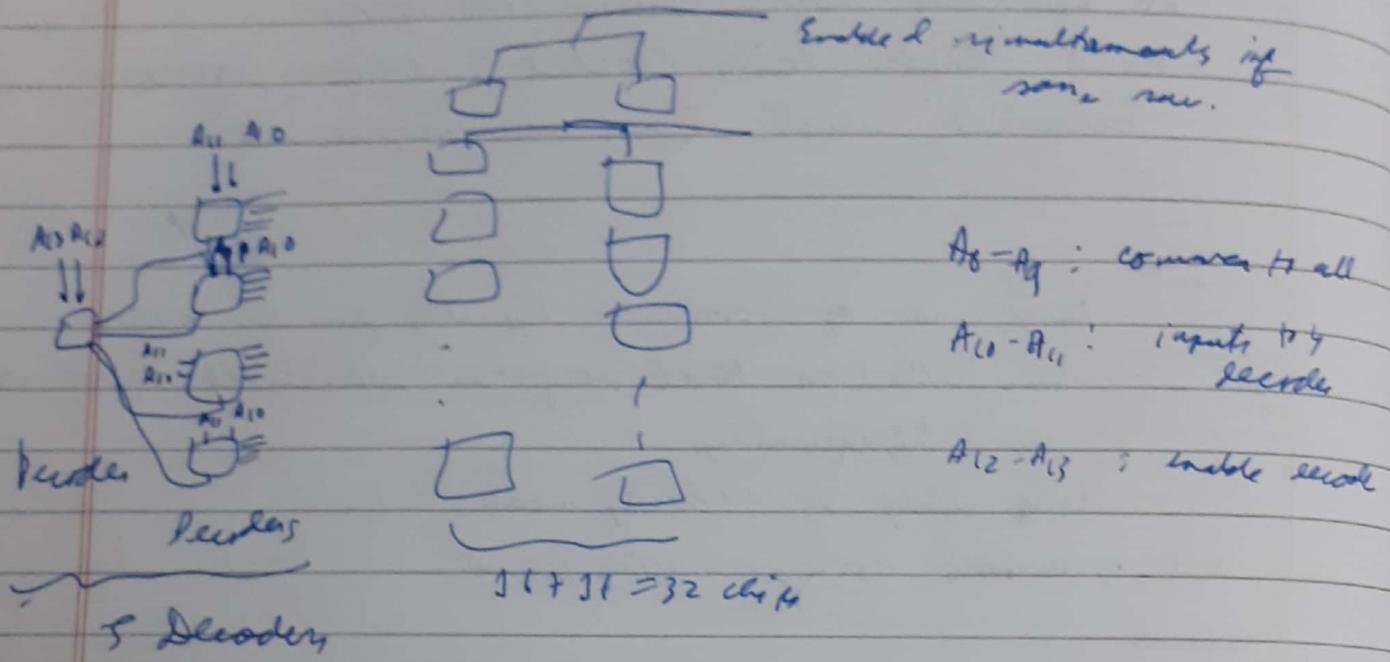
Decoder has chip-select role



(new
diagram
(assignment))

4) ton ship capacity 5F58

How many 2-to-4 decoders are required to generate enable for combinability
 $g(k \times 36)$ ROM from $3k \times 1$)



$$A_{13} \ A_{12} \ A_{11} \ A_{10} + \gamma \ldots - A_0$$

0 0 0 0 -

o o o o

0 0 0 0

8° 8

0 0 0 1 1

Digitized by srujanika@gmail.com

1024

34
dinner for
wife

7

302
address 10

chip 16

- 8) If magnetic disk has the following specification.
- No. of tracks = 1024
 - No. of sectors per track = 512
 - No. of bytes per sector = 512
 - Disk rotational speed = 7200 rpm
 - Seek time = 9 ms

Find data transfer rate and average access time

$$7200 \rightarrow \text{1 rotation}$$

$$1 \text{ rotation} \rightarrow \frac{1}{120} \text{ sec}$$

$$1 \text{ rotation} = 1 \text{ track access}$$

$$(512 \times 512) \text{ bytes} = 1 \text{ track size.}$$

$$\therefore 512 \times 512 \text{ bytes can be accessed in } \frac{1}{120} \text{ sec}$$

$$\text{No. of bytes accessed/sec} = 512 \times 512 \times 120$$

$$\text{Avg access time} = \text{seek time} + \text{latency time} + \text{data transfer time}$$

$$= (9 \times 10^{-3}) \text{ s} + \left(\frac{1}{2} \cdot \frac{1}{120} \right) + \left(\begin{array}{l} \text{Time to access} \\ \text{data from 1 sector} \\ 1 \text{ sector} = 512 \text{ bytes} \end{array} \right)$$

$$\text{Time required to access 512 bytes} = \frac{512}{512 \times 512 \times 120}$$

$$\therefore \text{Avg access time} = \frac{9}{1000} + \frac{1}{240} + \frac{1}{512 \times 120}$$

- g) Pipeline \rightarrow speedup factor = 50
 \rightarrow operating efficiency = 80%
 No. of stages in the pipeline = ?

$$E_K = \frac{\sum K}{K}$$

$$\frac{80}{100} = \frac{10}{K} \Rightarrow K = \lceil \frac{100}{8} \rceil = 13$$

- g) A non-pipeline system takes 50 ns to process a task.
 Same task can be processed in a 6-segment pipeline with clock cycle = 10 ns.
 Find speed up ratio of pipeline for 100 tasks.

$$\text{Speed up} = \frac{50 \text{ ns}}{\left[8 + (100-1) \right] \times 10 \text{ ns}} = \frac{5}{105}$$

- g) A processor takes 12 cycles to complete an instruction execution. The corresponding pipeline processor uses 6 stages with the execution times of 3, 2, 5, 4, 6 and 2 cycles respectively. What is the speedup that a very large no. of instructions are to be executed?

$$\text{Speed up} = \frac{12}{\max\{2, 3, 4, 5, 6\}} = 2$$

- g) Pipeline with 2 stages - each stage takes 10 ns. (multiplication, addition). How much time is required to complete following code segment?

for i = 1 to 100 do $A[i] = B[i] * (C[i] + D[i])$

$$\underbrace{\{K + (n-1)\}T}_{\text{No. of clock cycles required}} = \{2 + (100-1)\}10 = 1010 \text{ ns.}$$

Q) A pipeline has 4 stages with time delay
 $(50\text{ ns}, 60\text{ ns}, 90\text{ ns}, 80\text{ ns})$

Interface latch delay = 10 ns

Find clock frequency of the pipeline?

$$\text{Cycle time} = [\max \{50, 60, 90, 80\} + 10] \text{ ns}$$

$$= 100 \text{ ns}$$

$$\text{Clock frequency} = \frac{1}{100} \text{ GHz} = 10 \text{ MHz}$$

Q) A 4 stage pipeline has stage delay $150, 120, 160, 140\text{ ns}$ respectively. Registers that are used b/w the stages have a delay of 5 ns each.
 find total time taken to process 1000 data items on this pipeline.

on 1st data item

$$150 + 5 + 120 + 5 + 160 + 5 = 585 \text{ ns}$$

\sum^5

$$T_m + R = 165 \text{ ns}$$

\sum

$$\max \{150, 120, 160, 140\}$$

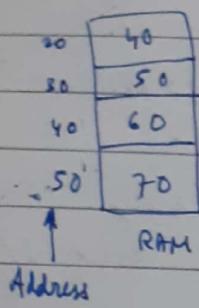
for all the other items,
 165 ns is time required.

$$585 + (1000 - 1) \times 165 \text{ ns}$$

Every operand field of an instruction is associated with some piece of data. In order to execute the instruction, the processor needs the current value of data, which can be specified in several ways.

Addressing Mode

Q)



Instruction: ~~Load AC~~

- (i) Load immediate 20
 - (ii) Load direct 20
 - (iii) Load indirect 20
 - (iv) Load immediate 30
 - (v) Load direct 30
 - (vi) Load indirect 30
- instruction address
already given
note

Find AC content
after execution of
instruction

- (i) $AC = 20 \rightarrow 2$ memory access \rightarrow read address 20 \rightarrow read data at the address (40)
- (ii) $AC = 40$
- (iii) $AC = 60$

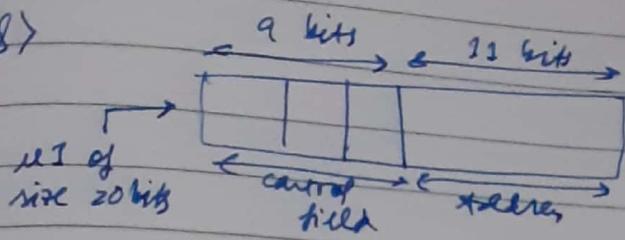
Q) In asynchronous serial communication controller uses a start-stop scheme for controlling the serial I/O of a system, it is programmed for a string of length 7 bits, 1 odd parity bit and 1 stop bit.

Transmission rate = 1200 bits/sec

- (i) What is the complete bitstream that is transmitted for the string 0110101 ?
- (ii) How many such strings can be transmitted per second?

- (i) $\xrightarrow{\text{parity}} 0110101 \xrightarrow{\text{stop (would be last bit)}} \rightarrow 10 \text{ bits long string}$
- (ii) $\frac{1200}{60} = 20$

Q)



$$\text{No. of M-op} = 2^3$$

Degree of parallelism = 3

Find the size of each control field and size of control ~~field~~ memory

$$2^3 + 2^3 + 2^3 - 3 = 21$$

$$\therefore \text{size of each control field} = 3$$

$$\text{size of each location} = 20 \text{ bit}$$

$$\text{no. of locations} = 2^3$$

$$\text{control memory size} = 2^3 \times 20$$