

288/FT

UNIT-III

COMPUTER ARITHMETIC

Addition & Subtraction:

Addition & subtraction with signed-Magnitude

Data:

Operation	Add	Subtract Magnitudes		
	Magnitudes	When $A > B$	$A < B$	$A = B$
$(+A) + (+B)$ Augend ↓ Addend	$+ (A+B)$			
$(+A) + (-B)$		$+ (A-B)$	$-(B-A)$	$+ (A-B) \Rightarrow 0$
$(-A) + (+B)$		$- (A-B)$	$+ (B-A)$	0 $+ (A-B)$
$(-A) + (-B)$ $(-A) - (+B)$ Minuend ↓ Subtrahend	$- (A+B)$			
$(+A) - (+B)$ Minuend ↓ Subtrahend		$+ (-A+B)$	$-(B-A)$	$+ (A-B)$
$(+A) - (-B)$ ($-B$)	$+ (A+B)$			
$(-A) - (+B)$	$- (A+B)$			
$(-A) - (-B)$		$- (A+B)$	$+ (B-A)$	$+ (A-B)$

→ For add operations, we add the magnitudes

when the signs are same.

- For subtract operations, we add Magnitude when the signs are different
- When the magnitudes are added, the sign of the result will be same as that of A.
- For Add, we subtract magnitudes if signs are different.
- For subtract, we subtract magnitudes if signs are same
- For $A > B$, condition, the sign of result is same as sign of A.
- In $A < B$ condition, sign of result is complement of sign of A.
- In $A = B$ condition, the sign of result is always +ve.

HARDWARE SYSTEM (for addition & subtraction)

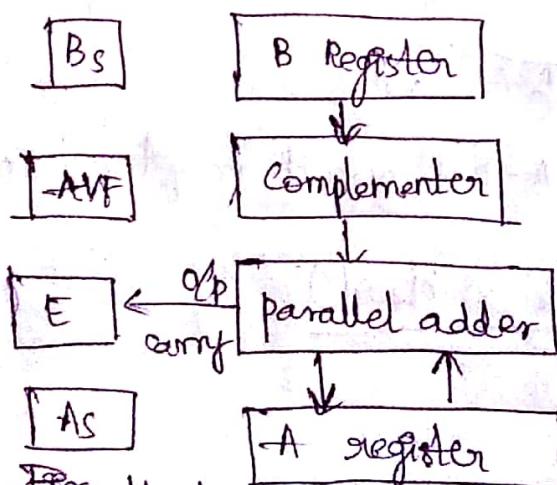


Fig Hardware for signed magnitude addition & subtraction

Hardware system is used for Addition and subtraction.

→ For addition, - A register & B registers contains operands and A_S & B_S flipflops indicates sign.

→ For subtraction, - A register - minuend, A_S - sign of minuend, B borrow - subtracted substrahend, B_S - sign of substrahend.

→ To perform $A+B$, $M=0$

→ To perform $A-B$, $A+B+1$ is performed.

To perform $A+B+1$, $M=1$.

→ After adding A & B, carry goes to E-FlipFlop

and Result is stored in E & A register

If $E=1 \Rightarrow$ Add overflow(AVE)
↳ of carry

To subtract the magnitudes, perform $A+B+1$

When $M=1 \rightarrow$ If carry & B is complemented

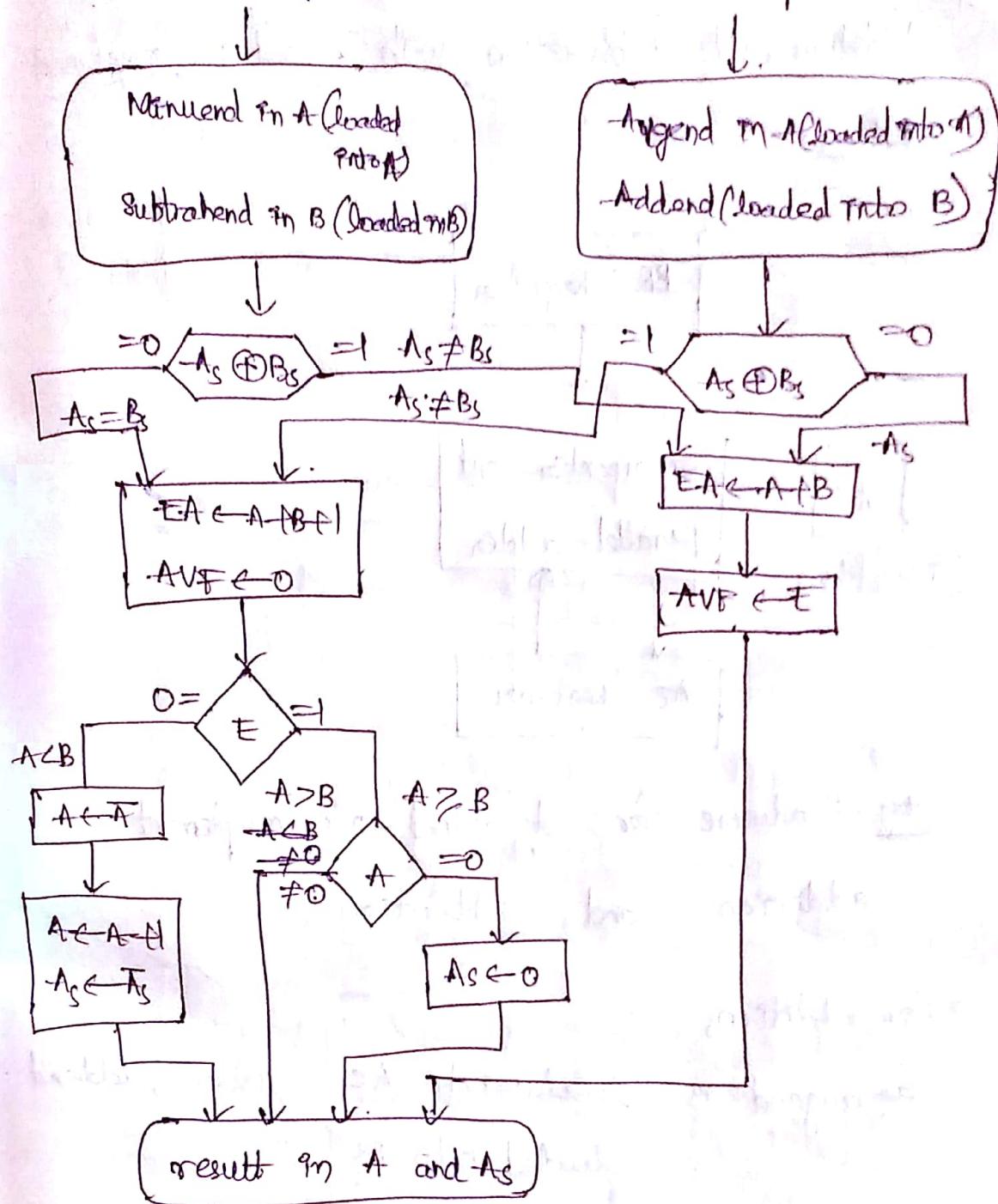
∴ of is $A+B+1 \Rightarrow A-B$

After doing $A+B+1$, result is stored in E

Flow chart (for add & subtract)

Subtract operation

Add operation



30/8/17

Addition and subtraction with signed-2's complement

Data :

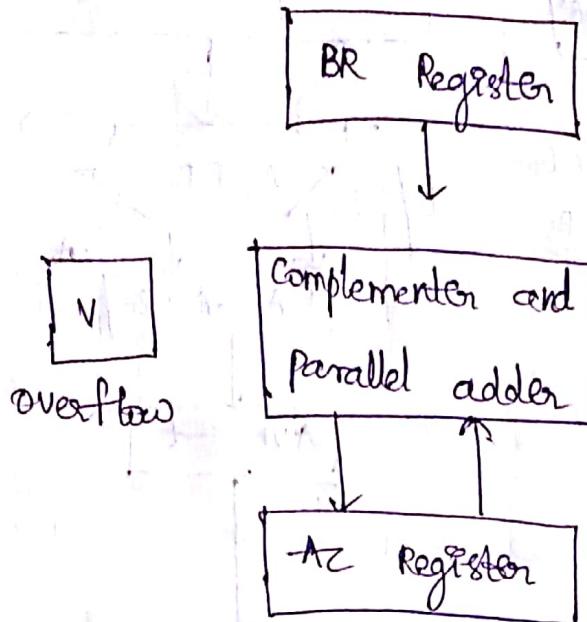


Fig - Hardware for signed-2's complement addition and subtraction

→ For addition,

augend is loaded into AC register, addend is loaded into BR "

Then AC and BR are added and the result

is stored in AC. If overflow occurs,

if overflow flip flop V is set.

→ For subtraction, minuend is loaded into AC and subtrahend is loaded into BR

BR is complemented (2's complement) and it is added to AC. The result $AC + BR + 1$ ($A_c = BR$) is stored in AC. V will be set if overflow occurs.

Flowchart for common hardware:

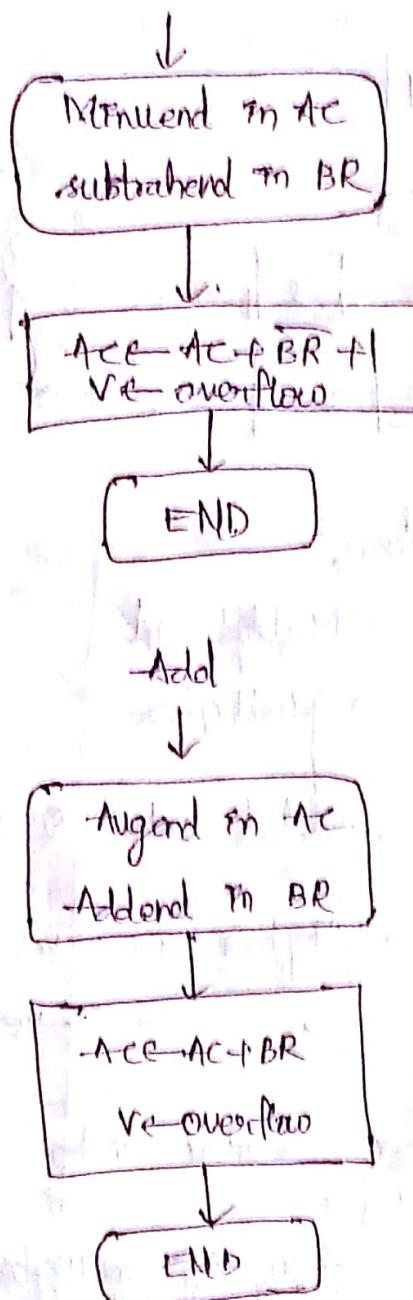
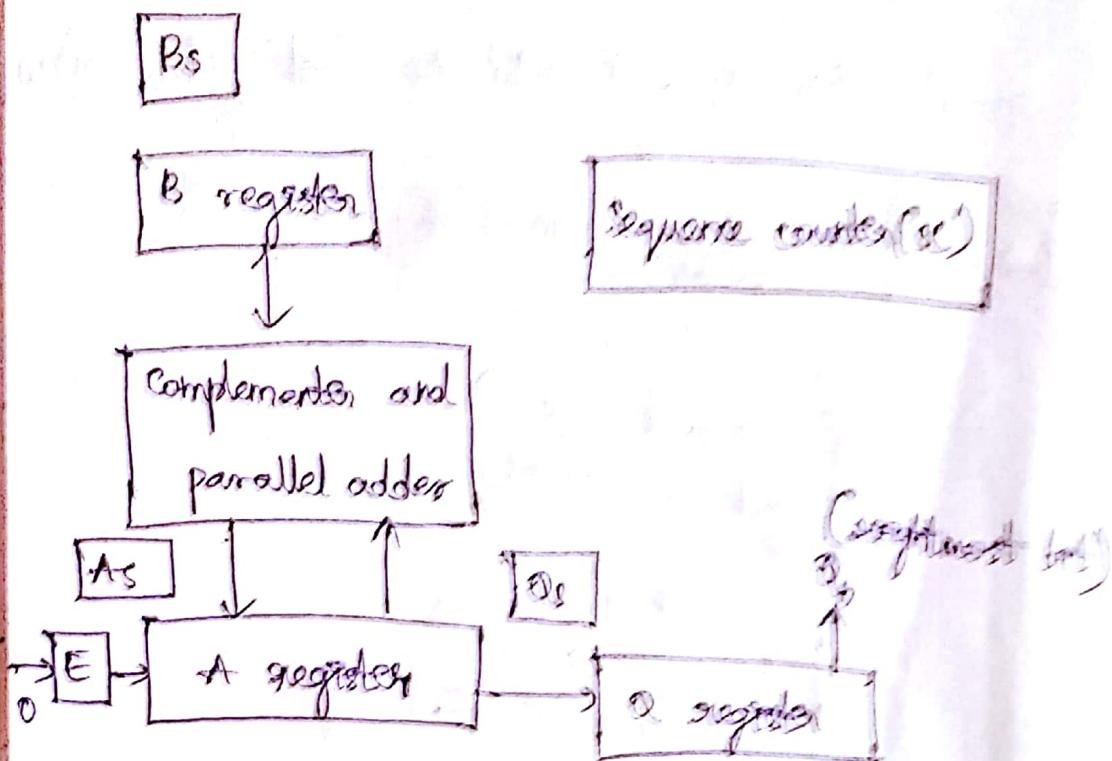


Fig: Algorithm for adding and subtracting numbers in signed 2's complement representation.

Multiplication - ALU Algorithm



In multiplication,

first operand \rightarrow Multiplicand

second " \rightarrow Multiplier

\rightarrow Multiplier is loaded into Q register

$Q_s = 0$ if multiplier is +ve

$Q_s = 1$ " " " " -ve

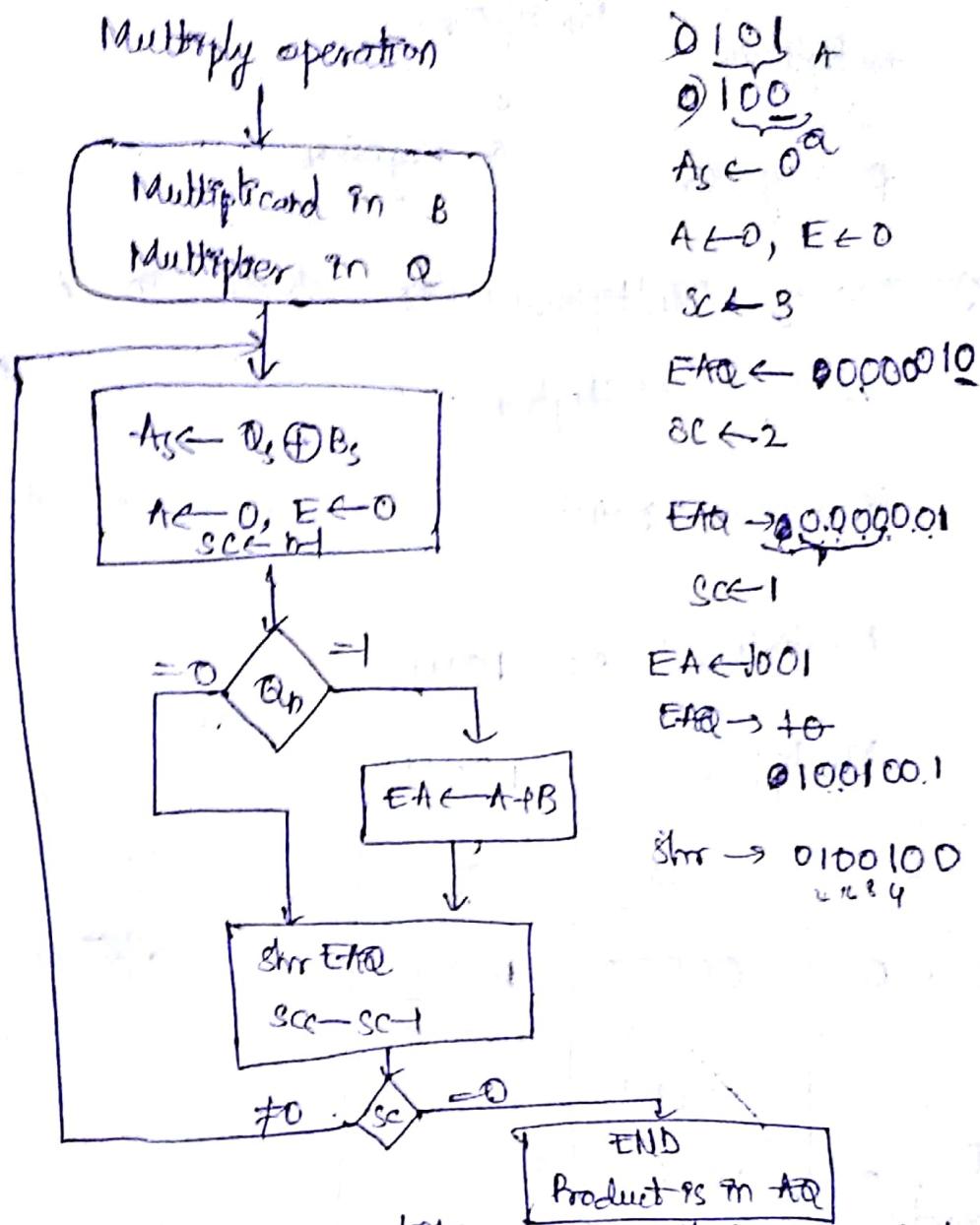
\rightarrow Multiplicand is loaded into register B

B_s is the sign of multiplicand

\rightarrow Result of multiplication is stored in A & Q registers, $A_s \rightarrow$ sign of result

\rightarrow sc contains no. of bits in multiplier excluding sign bit.

Algorithm:



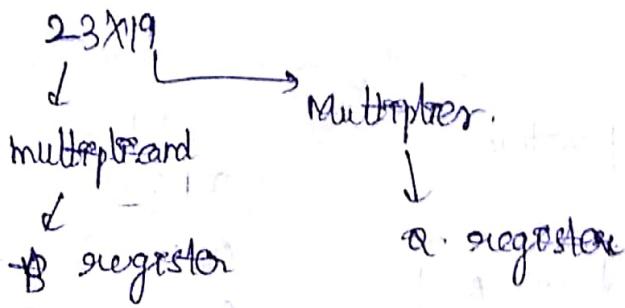
'n' indicates no.of bits in multiplier excluding sign bit

After adding A and B, if carry is generated, it is stored in E.

Shr → logic shift right (Ashr → Arithmetic shift right)

~~E ← 10~~ shr EAQ → indicates E is loaded into A and LSB of E is loaded into M&B of Q

~~Eq.~~



→ Sign of Multiplicand is stored in A₅ A₄
 " " Multiplicator "

Sign of result → A₅

Multiplicand B = 10111

Multiplicator Q = 10011

E	A	Q	SC
Initial	0	00000	101
Q _n = 1; add B First product	0	$\begin{array}{r} 10111 \\ + 10111 \\ \hline \end{array}$	
Shrt EAQ	0	$\begin{array}{r} 01011 \\ 10111 \\ \hline \end{array}$	100
Q _n = 1; add B Second partial Product Shrt EAQ	0	$\begin{array}{r} 00010 \\ + 0111 \\ \hline \end{array}$	011
Q _n = 0 Shrt EAQ	0	$\begin{array}{r} 01000 \\ + 0110 \\ \hline \end{array}$	010
Q _n = 0; Shrt EAQ	0	$\begin{array}{r} 00100 \\ + 1011 \\ \hline \end{array}$	001
Q _n = 1; add B Shrt EAQ	0	$\begin{array}{r} 11011 \\ + 01011 \\ \hline \end{array}$	000

$$437 = 01101101$$

The value EA before performing the operation
is called partial product.

Booth Multiplication Algorithm

The algorithm used for performing multiplication
of - signed 2's complement numbers is called
Booth multiplication algorithm.

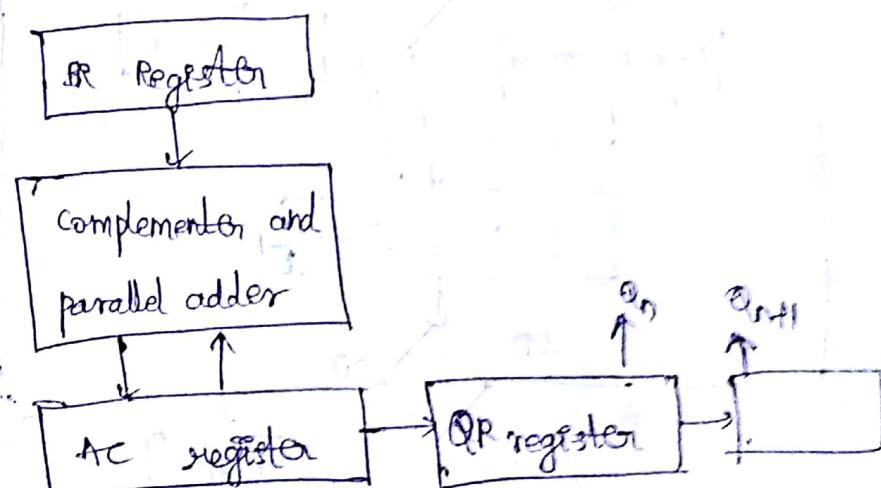
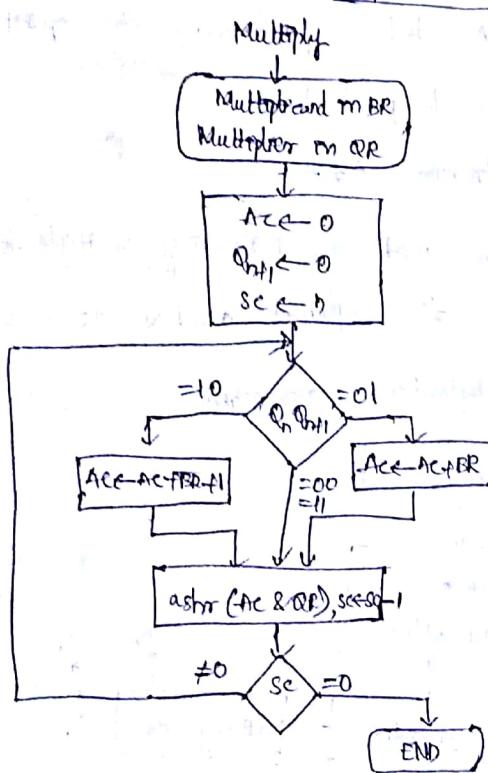


Fig. Hardware for Booth algorithm

QR contains the multiplier in 2's complement form
BR " multiplicand "

- SC contains the no. of bits in QR register.
- QR & PR contains the sign bits also.
- After multiplication, result is given - AC & QR
- The bit coming out from QR is goes into Q_{n+1}

Flowchart for Booth Multiplication:



Example of multiplication with Booth algorithms:

$$f(x)(-13) = 117$$

↓
Multiplicand

$$BR = 10111$$

$$PR_{ft} = 01001$$

$$Q_n Q_{n+1} =$$

$$\begin{array}{r} -9 \rightarrow 01001 \rightarrow \\ 10110 \end{array}$$

$$10111$$

$$\begin{array}{r} -13 \rightarrow 01101 \\ 00010 \end{array}$$

$$BR = 10011$$

$Q_n Q_{n+1}$	-AC	QR	Q_{n+1}	SC
Initial	00000	10011	0	101
Subtract to to BR	01001	01001		
ashr	00100	11001	1	100
ashr	00010	01100	1	011
Add BR	10111	11001		
ashr	11100	10110	0	010
ashr	11110	01011	0	001
Subtract BR	01001	01011		
ashr	00011	10101	1	000

Result is stored in -AC & QR

Array Multiplier (Fast multiplier)

concept

Consider 2-bit by 2-bit multiplier

$$\begin{array}{r}
 & b_1 & b_0 \\
 & a_1 & a_0 \\
 \hline
 a_0 & b_1 & a_0 & b_0 \\
 a_1 b_1 & a_1 b_0 & & \\
 \hline
 c_3 & c_2 & c_1 & c_0
 \end{array}$$

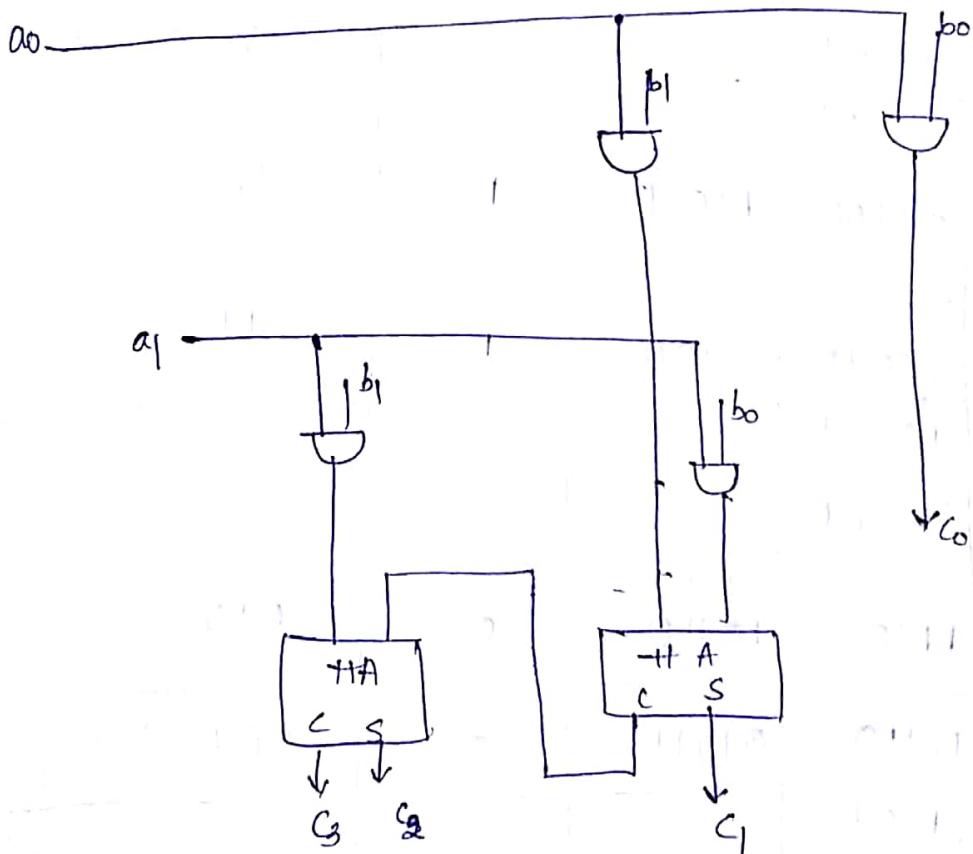


Fig: 2-bit by 2-bit Array Multiplier.

→ It is faster when compared to other multipliers

→ For j-multiplier bits and k-multiplicand

bits, we need $j \times k$ AND gates and

$j-1$, k -bit adders to produce a product of $j+k$ bits.

Consider $j=3, k=4 \rightarrow$ Multiplier and bits.
 \hookrightarrow Multiplier bits
 $j \times k = 12$ AND gates.
 $j-1 = 2$ four-bit adders

$k=4, j=3$ 4 bit by 3 bit multiplier

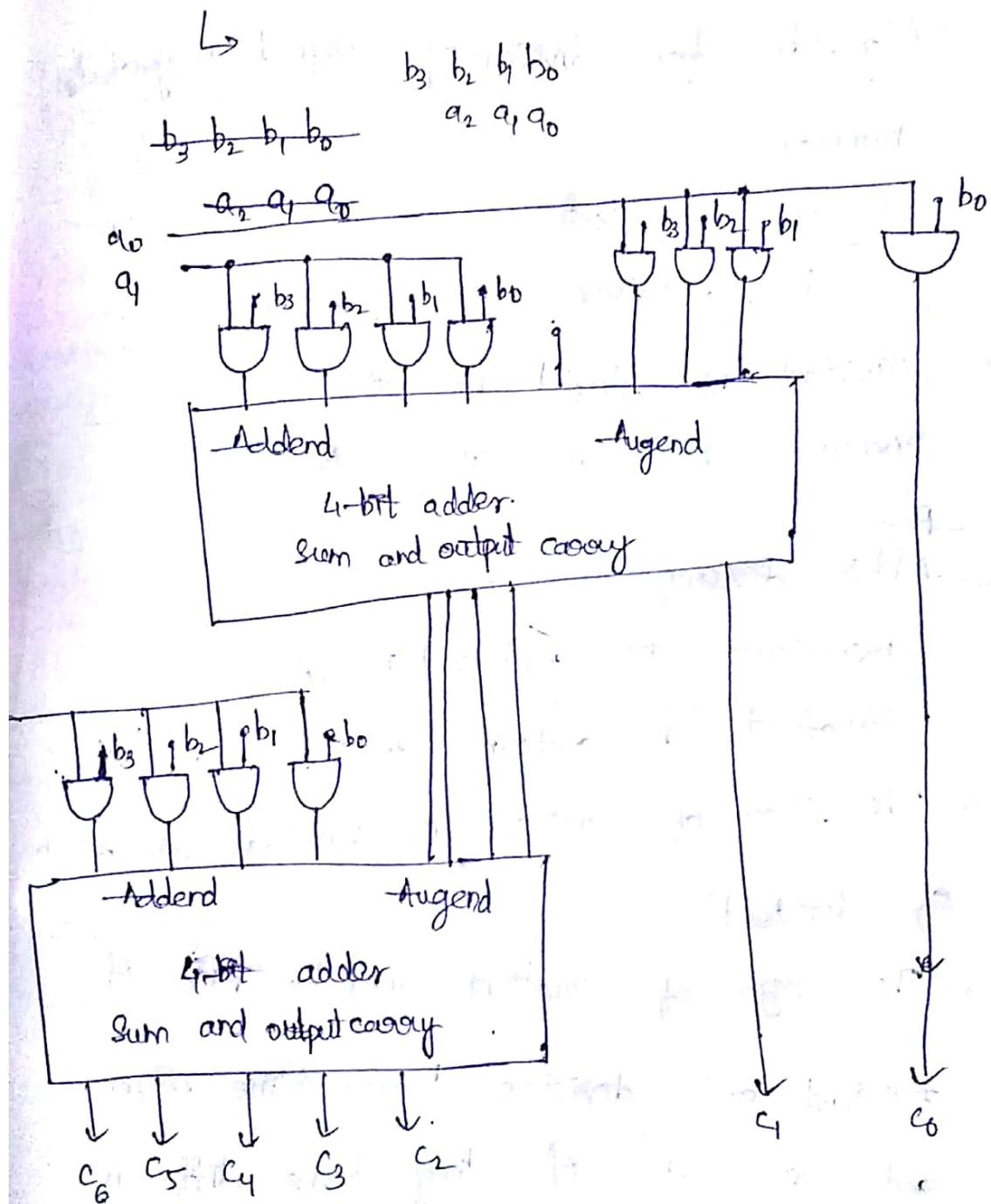


Fig: 4-bit by 3-bit array multiplier.

→ There is a 4x3 multiplier in which, a $(j-1)$ bit address to produce a product of $4+3=7$ bits output.

Division Algorithm:

→ Algorithm for division of signed magnitude numbers

Nr → Dividend

D → Divisor

→ Dividend is stored in A & Q

Divisor is in B

-Re-

→ After division,

Remainder is contained in A

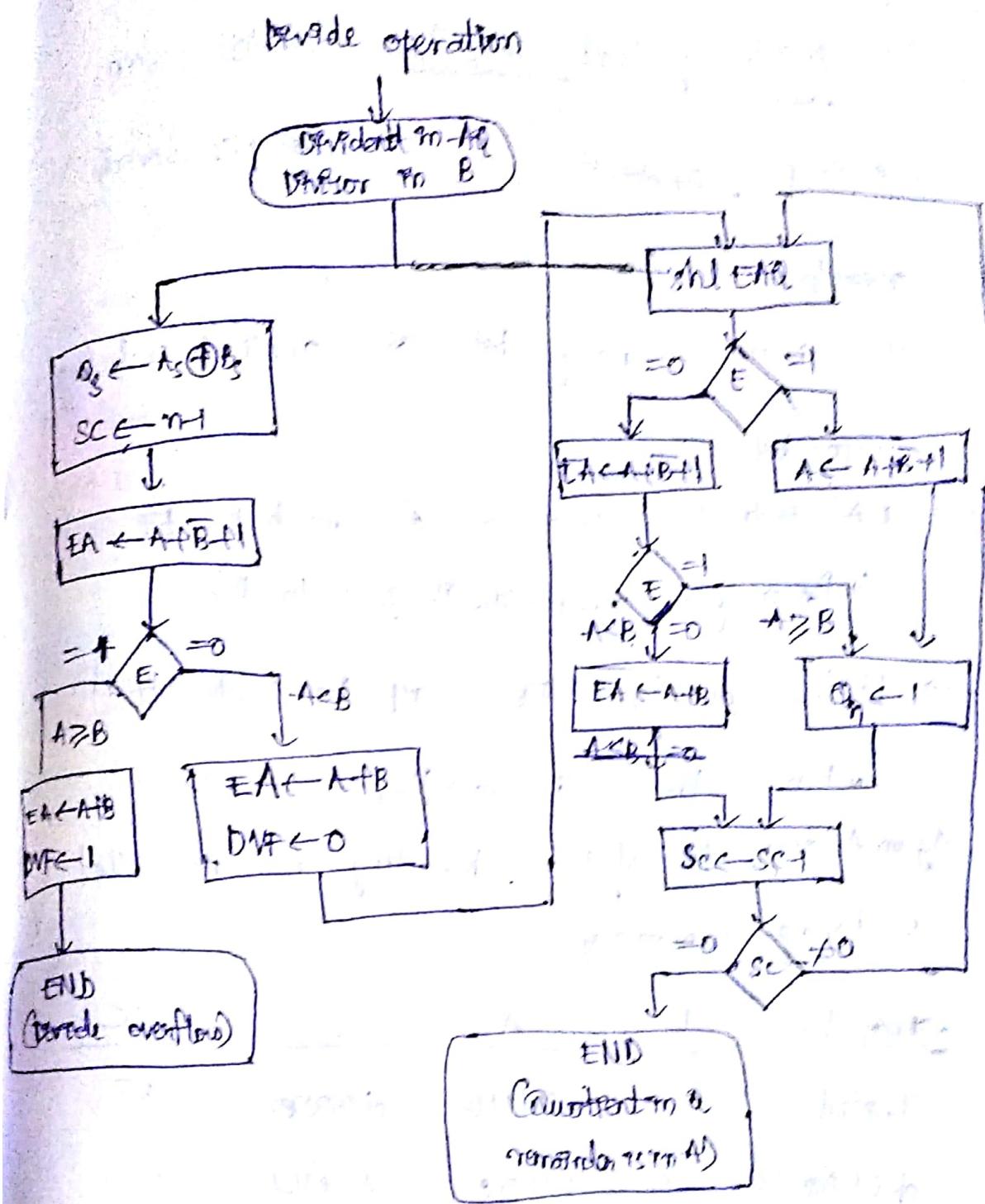
Quotient is contained in register Q

→ The sign of remainder is same as the sign of dividend

→ The sign of quotient will be +ve if dividend and divisor have same sign and -ve if they have diff. n

Sign of divisor is in B_s
a u - dividend = a A_s

Flowchart



Ex: $A = 448$

$B = 17$

$B = 10001$

$B+1 = 01111$

- The divide / quotient overflow occurs when $A \geq B$. quotient indicates divide overflow.
- 'm' indicates no. of bits in quotient including sign bit.
- EA indicates register A with E=FF.
(If carry is generated, it goes to E)
- Divide overflow flag will be set to 1 when there is overflow.
- Again Adding B after subtracting B is called restoring remainder.

<u>Initial</u>	<u>E</u>	<u>A</u>	<u>Q</u>	<u>SC</u>
Initial		01110	00000	5
shl EAQ	0	11100	00000	
add B+1		01111		
E=1	1	01011	00001	
Set Q _h =1	1	01011	00001	4
shl EAQ	0	10110	00010	
add B+1		01111		
E=E-1	1	00101		
Set Q _h =1	1	00101	00011	3
shl EAQ	0	01010	00110	
add B+1		01111		
E=0, leaves	0	11001		
Add B, leaves		10001		
Restoring remainder		01010		2

<u>E</u>	<u>A</u>	<u>Q</u>	<u>Sc</u>
0	10100 11111	01100	
Add B+1	00011		
E=1	00011	01101	1
set and	00010	11010	
shl E+1	01111		
Add B+1	10101		
E=0 leave Q=0	10001	H010	0
Add B-	00110		
	6	26	
	↓ : remainder	↓ quotient	

Floating-Point Arithmetic Operations:

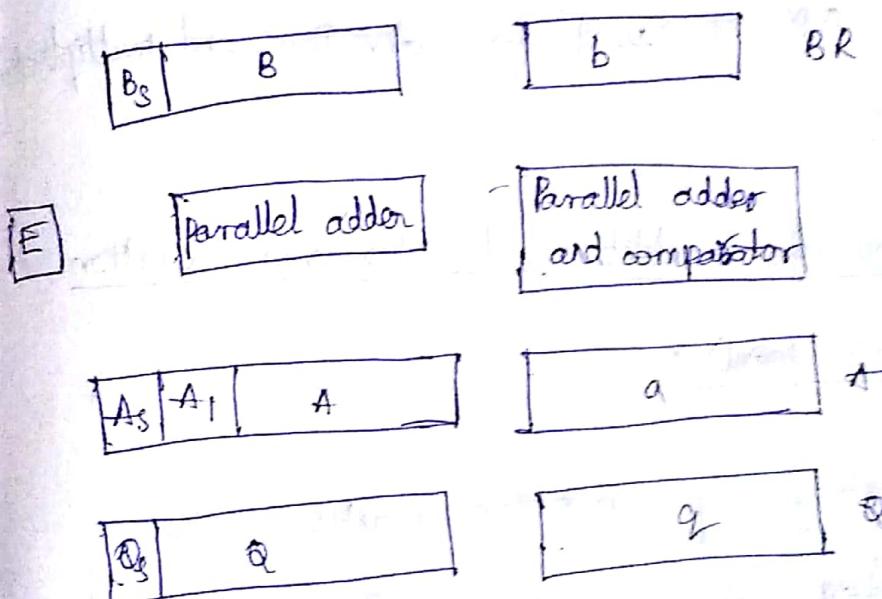


Fig: Register for floating-point arithmetic operations

B_S — sign of mantissa

B — magnitude of mantissa

b — exponent

Together it is called BR

Similarly A_S, A, A and 'a' indicate AC

A - MSB of the magnitude of mantissa.

If a_1, a_2, q indicates OR register

→ For addition operation

AC contains augend

ER " addend and result is stored in AC

→ For subtract operation

AC contains minuend

ER " ~~subtracted~~ ^{subtrahend} and result is stored in AC

→ Register OR is used for division and multiplication operation.

Algorithm for addition and subtraction operation

(floating point):

→ In floating representation, mantissa q_3

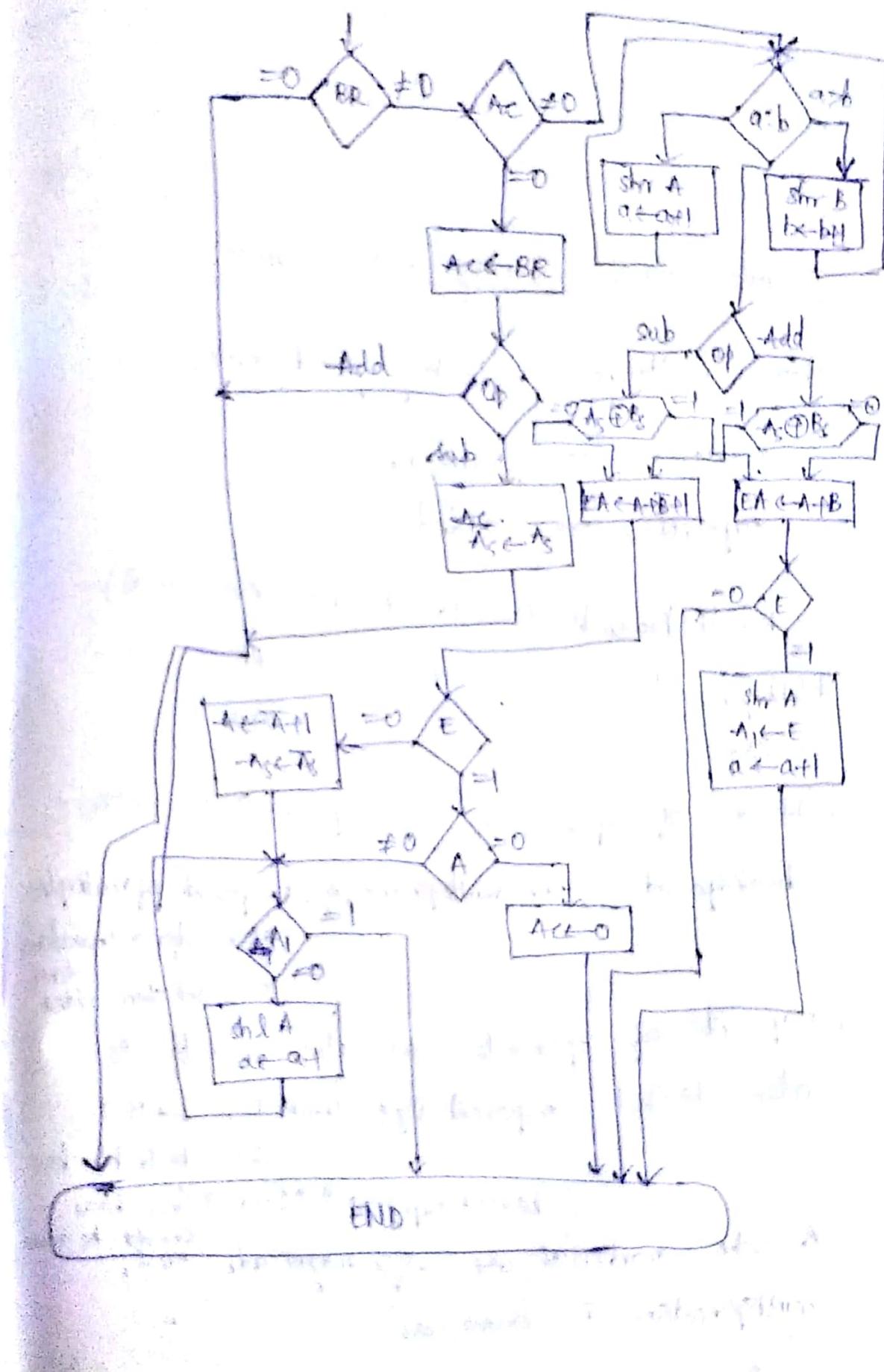
represented in sign-magnitude form and

exponent is represented in biased form

→ Biased exponent is always a fix value
which is given by adding the no. to
the actual exponent.

for addition and subtraction operations, exponent values of the operands should be made equal

Add or subtract



If $BP = 0$, no operation is performed ($AC \leftarrow AC + AC$
 $AC \leftarrow AC \times AC$)

When exponent is unmarked, mantissa is shifted right

→ if exponent is unmarked, $AC \leftarrow AC \times 10^0$
Left

Multiplication of floating point numbers

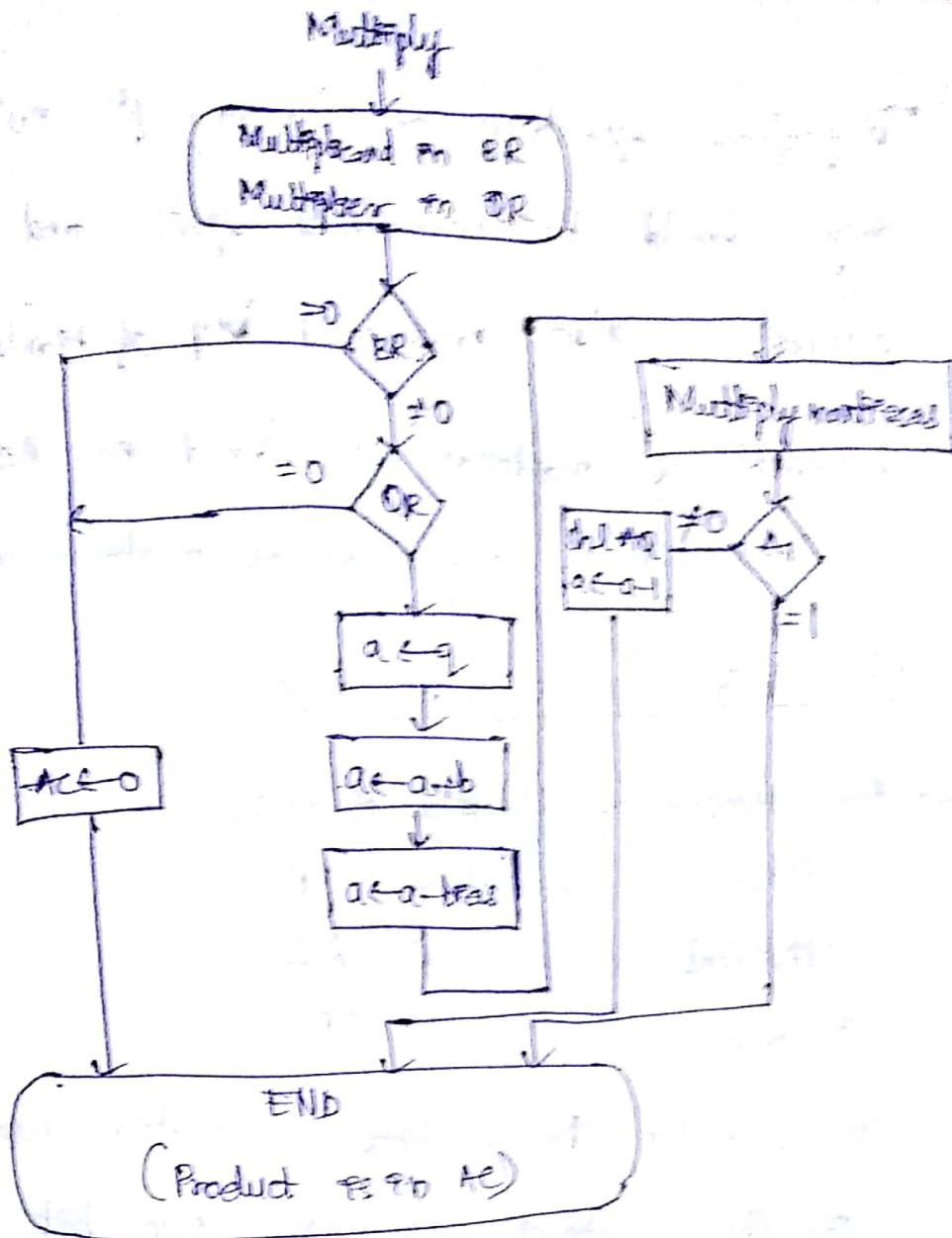
- After multiplying floating point no's,
mantissas are multiplied
exponents are added

Result (Product) is stored in AC (left)

Multiplicand " " BP
Multiplier " " OP

- If one of operand is zero, product is zero.
Exponent of BP multiplicand, q = exponent of multiplier
It is after transforming to d and then added.
- Exponents of operands and also result is also biased exponent. E.g. 4.11011 ; 2.11011

- AC → double bias
biased exp → 4.11011 (one bias)
The mantissas are sign magnitude no's,
multiplication is same as "



(It is shifted left only once $\therefore \text{e.g. } 0.110\bar{1} = 0.\underline{01}$)

To perform operation on floating pt. no's, they should be normalized first and result is also normalized (mag of mantissa)

product of mantissas is stored in AQ

$A \rightarrow$ MSB of magnitude of mantissa of result

Division of floating point numbers:

→ for division of floating pt no's,

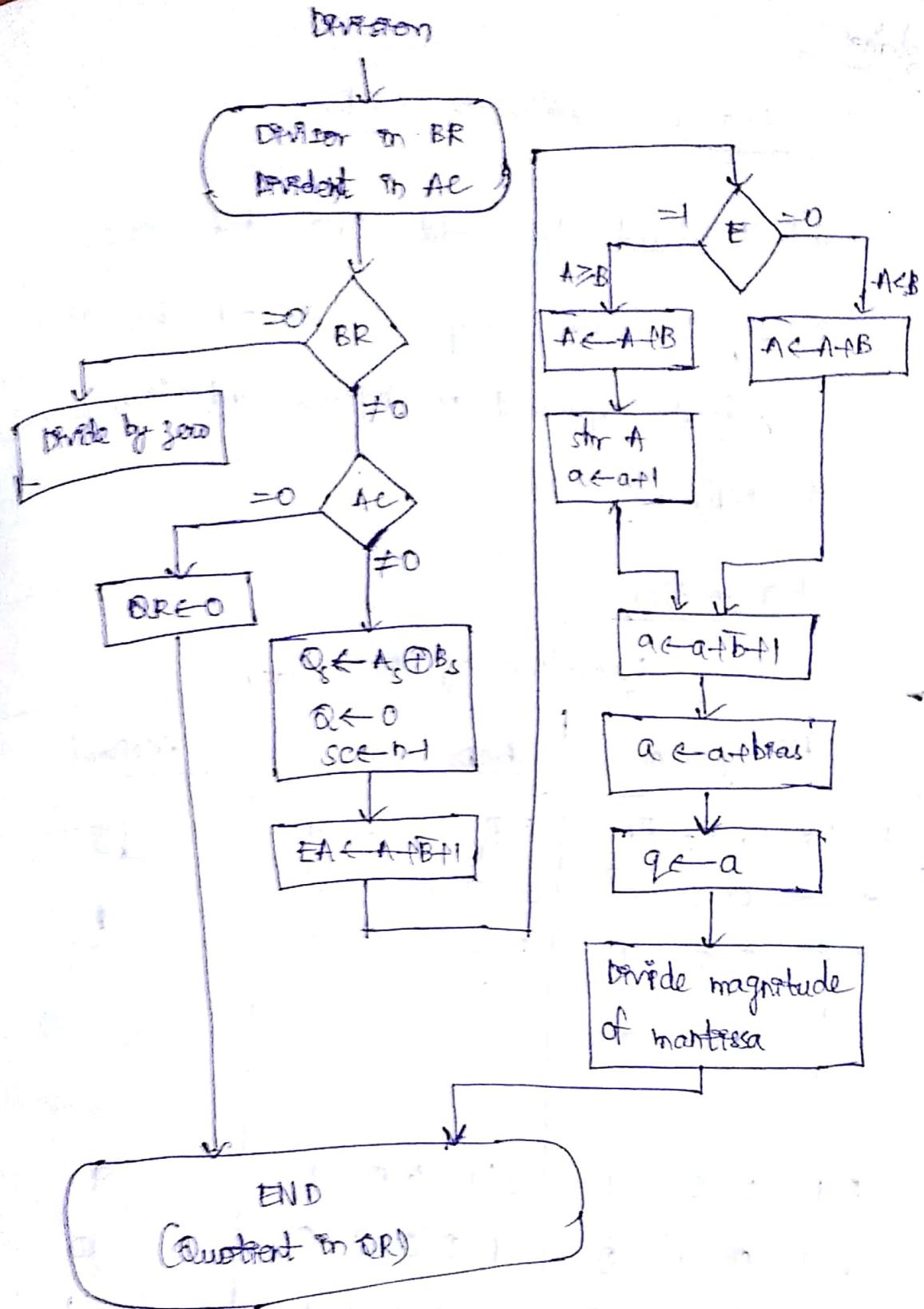
divisor is stored in BR

Dividend \rightarrow AC

Quotient \rightarrow DR

→ 'n' indicates no. of bits used for mantissa of the quotient excluding sign bit.

Overflow (divide) is avoided in floating point numbers by performing certain operations (ie. shift right)



~~gated~~ Decimal Methane, CH₄

It is used to add two BCD nos.

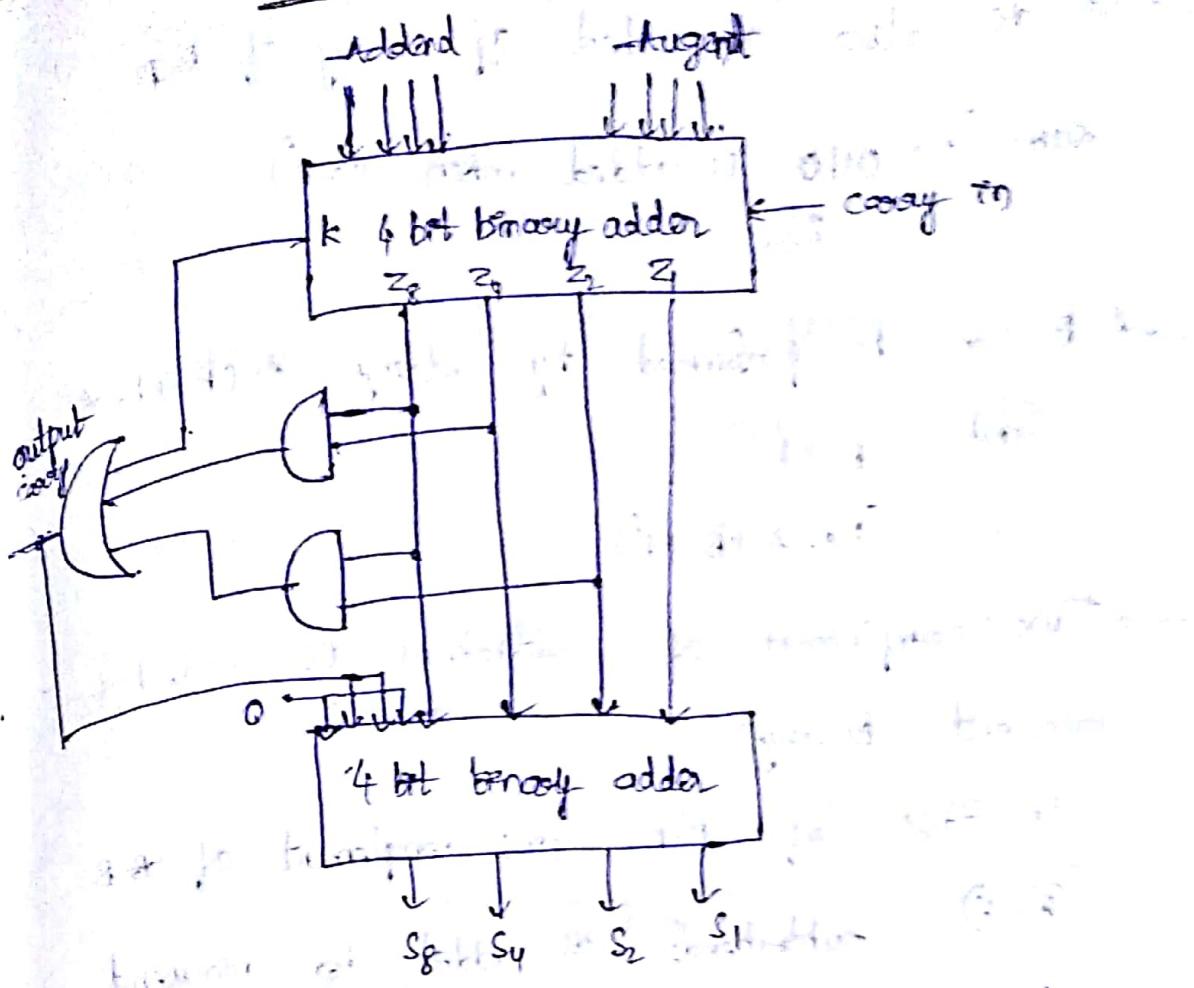
Count is also expressed in BCD. It is first expressed in binary and then converted to BCD.

BCD Adder

Implementation of BCD-Adder

Binary sum	BCD sum	Decimal
0 1 1 0 1 0 0 0	0 1 0 0 0 0 0 0	10
0 1 0 1 0 0 0 0	0 1 0 0 0 0 0 0	11
0 1 0 1 0 0 0 0	0 1 0 0 0 0 0 0	12
0 1 0 1 0 0 0 0	0 1 0 0 0 0 0 0	13
0 1 0 1 0 0 0 0	0 1 0 0 0 0 0 0	14
0 1 0 1 0 0 0 0	0 1 0 0 0 0 0 0	15
0 1 0 1 0 0 0 0	0 1 0 0 0 0 0 0	16
0 1 0 1 0 0 0 0	0 1 0 0 0 0 0 0	17
0 1 0 1 0 0 0 0	0 1 0 0 0 0 0 0	18
0 1 0 1 0 0 0 0	0 1 0 0 0 0 0 0	19

BCD Adder:



In addition, first operand is augend, second operand ~~is~~ is addend.

Binary and BCD codes are same from 0 to 9

Binary is converted into BCD when decimal > 9

i.e. (i) when $k=1$ (on) \rightarrow 1^{st} ip of OR gate

(ii) $z_3 z_4 = 1$ (on) \rightarrow 2^{nd} ip of OR gate

(iii) $z_3 z_2 = 1$ \rightarrow 3^{rd} " "

Binary is converted into BCD by adding 6 (0110)

$$C = k + z_3 z_4 + z_3 z_2$$

When $C=1$, 0110(6) is added to binary

- C is also called off carry of BCD sum (\because 010 is added when $C=1$)
 i.e., $A+B+C$
 → Subtraction can be performed by adding A with $\bar{B}+1$
 (i.e. $A+\bar{B}+1$).

→ Two's complement of subtractend is added to minuend (Borrow)
 In case of BCD, 10's complement of \bar{B} ($\bar{B}+1$) (subtractend) is added to minuend

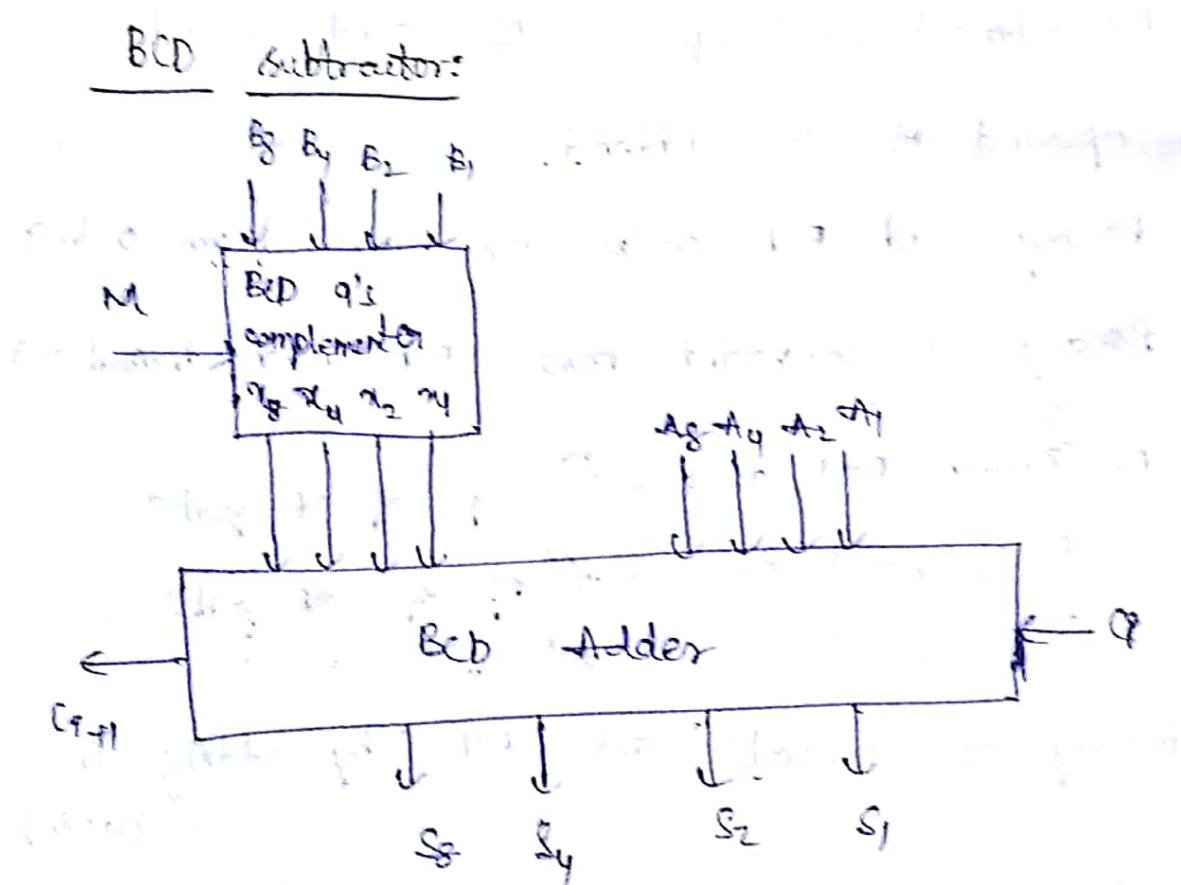


Fig: One stage of a decimal arithmetic unit or adder or part of adder.

n stages are required for n digit decimal no.

$M=0$ for BCD addition,

$M \neq 0$ for subtraction ($\bar{B}+1$) (\bar{B})

for addition, A and B are added

for subtraction, A and 9's complement of B is added.

C_0 is the 1's carry which is always '1' for first stage of BCD subtraction ($\bar{B}+1$)

$C_0 = 0$ for first stage of BCD addition

whereas for the other stages, it depends on the 1's carry of previous stages.

$$\underline{\underline{x_4}} = B_1 M + B_1' M$$

$$\underline{\underline{x_4}} = B_1 M + B_1' M$$

$$x_2 = B_2$$

$$x_3 = B_3 M' + (f_4' B_2 + B_1 B_2') M$$

$$x_4 = B_4 M' + B_3' B_2' B_1' M$$

$$M=0 \Rightarrow \underline{\underline{x_4}} = x_4 = B_1 \quad M=1 \Rightarrow x_4 = B_1'$$

$$x_2 = B_2$$

$$x_3 = B_3$$

$$x_4 = B_4$$

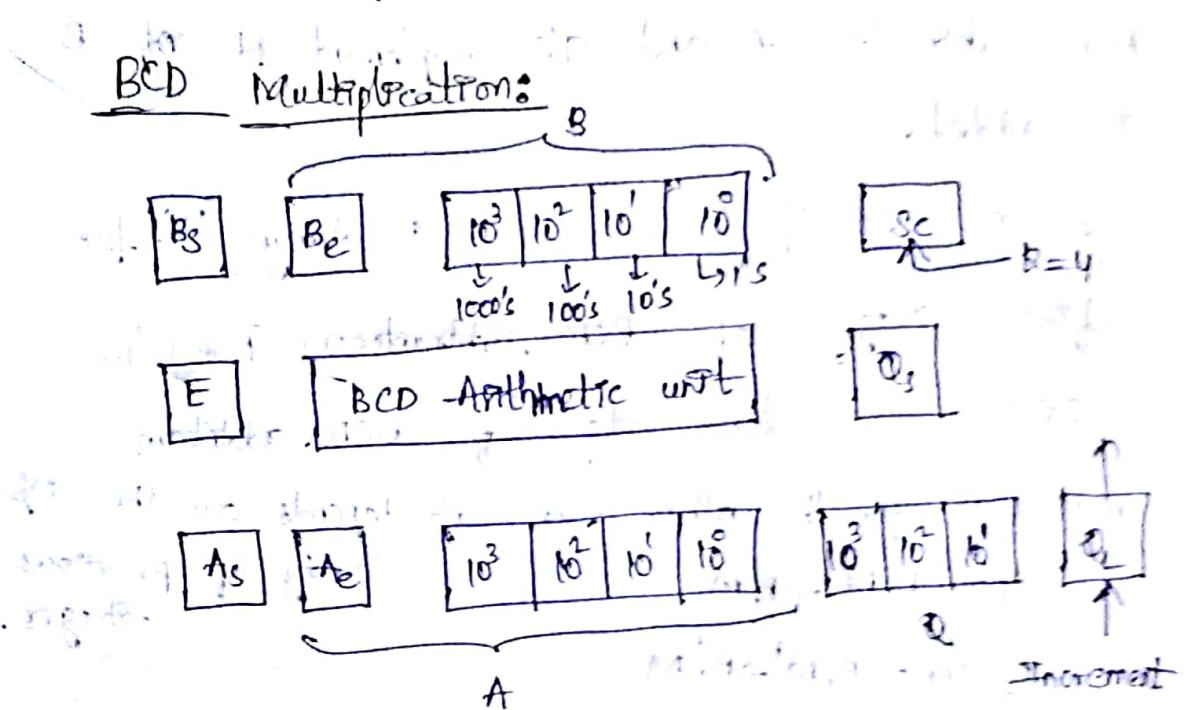
$$x_2 = B_2'$$

$$x_3 = B_3'$$

$$x_4 = B_4'$$

→ The algorithm applied for sign magnitude binary no's can also be applied for signed BCD no's.

→ Similarly, algorithm for addition and multiplication subtraction of 1's complement of binary no's can also be extended for $\frac{1}{2}$'s complement of BCD.



It is used for 4-digit BCD no's (16 bits)

Multiplicand - Each digit is represented in 4 bits.

Multiplicand is stored in B

Multiplexer u v w q

A_e, B_e are extension bits.

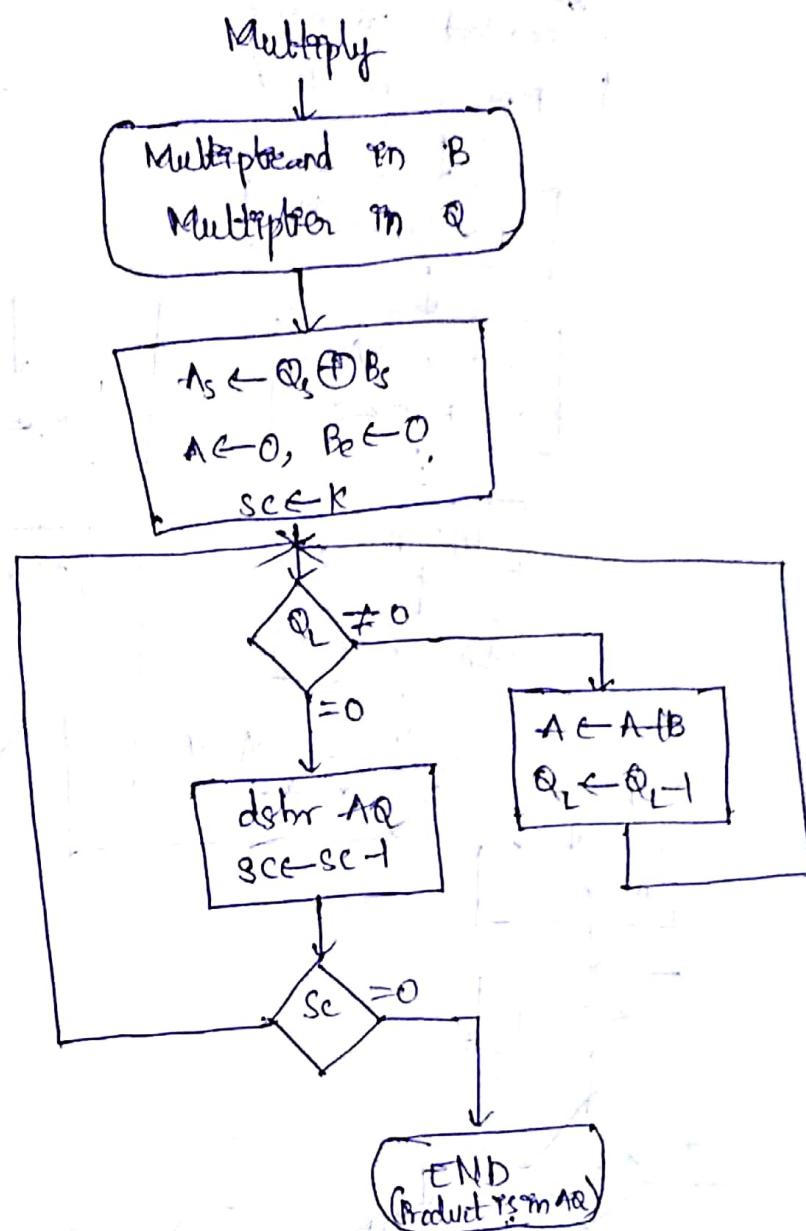
→ For division, dividend in ~~A_q~~ divisor in B

K indicates the no. of ~~bits~~ ^{digits} in multiplier ($k=4$)

For division, k - no of digits
in quotient

After multiplication, product is stored in AQ

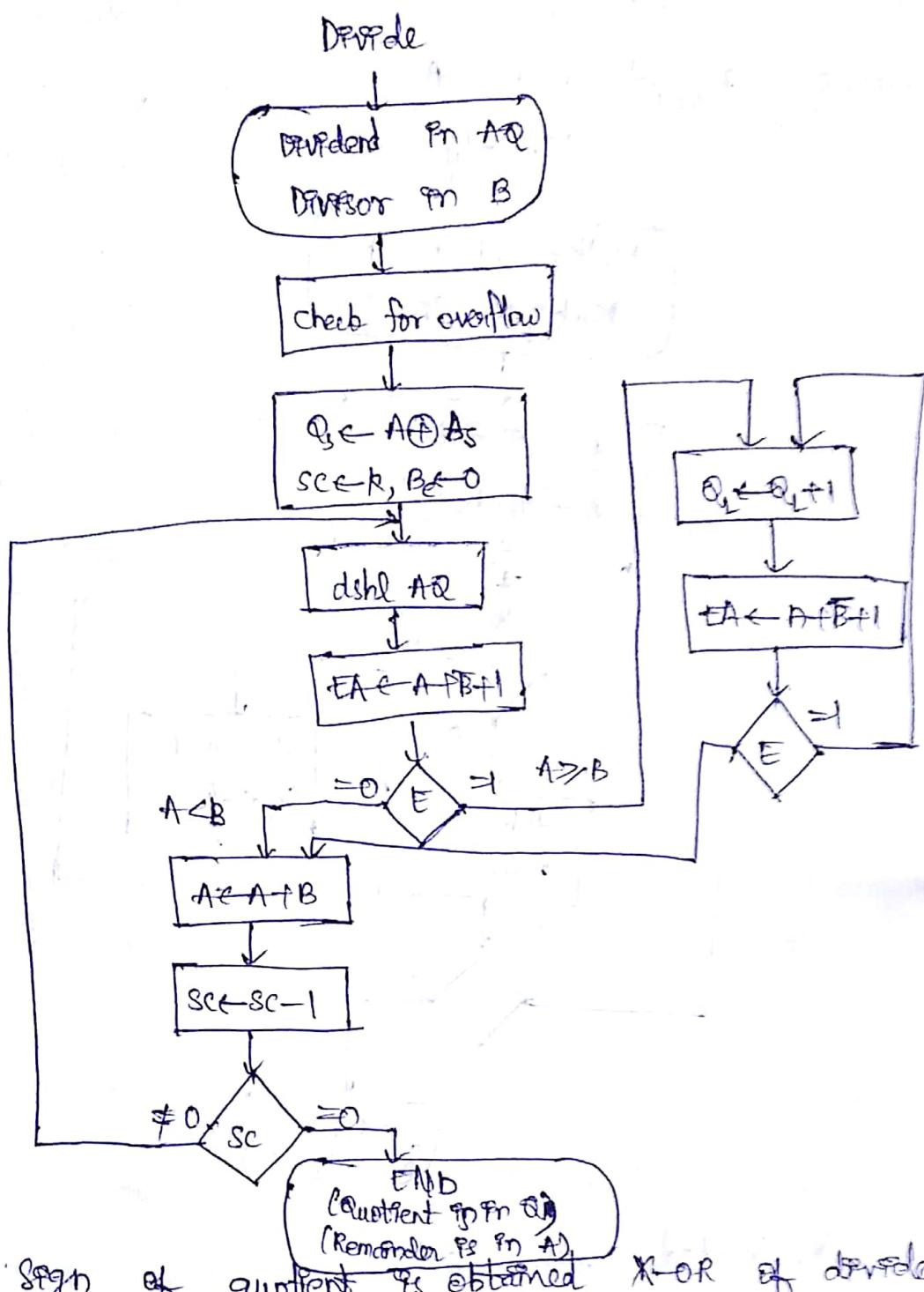
After division, quotient Q and
remainder is stored in A.



sc is loaded with 'k'

Q_L indicates the no of LSB of Q (multiplier or
dividend) dstr → decimal shift right

Division (Signed BCD)



Sign of quotient is obtained \times OR of dividend and divisor and sign of remainder is the sign of dividend.

k. indicates no. of digits in dividend quotient

$E=0 \Rightarrow A < B \Rightarrow A < B$ (\because carry not generated \Rightarrow --ve)

$E=1 \Rightarrow A \geq B$

for operations on
 algorithm applied float pt. binary no's can also
 " " " decimal no's (BRD) (+, -, *, /)
 be

THE MEMORY SYSTEM

Some basic concepts:

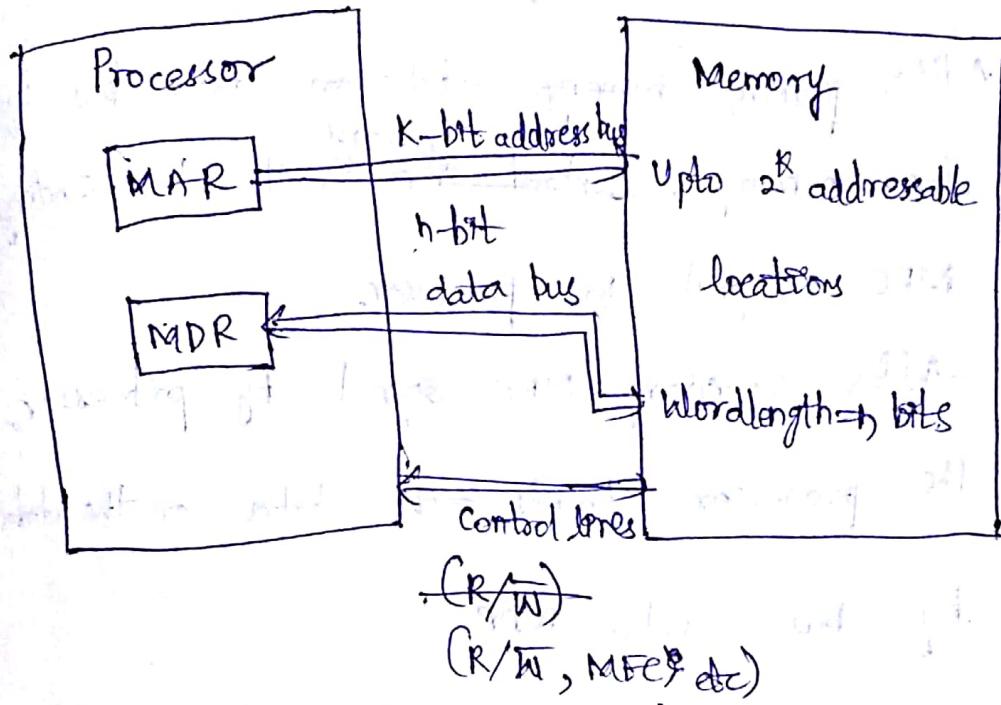


Fig. Connection of the memory to the processor

MAR (Memory Address Register)

MAR provides address of memory location to be accessed

MDR (Memory Data Register)

It contains the data that is to be read written into memory

Address bus is used to transfer address data

control lines are used to transfer control signals.

Eg: write/Read signal ($R/W = 0 \Rightarrow$ Read)

$R/W = 1 \Rightarrow$ Write)

MFC (Memory function completed)

After placing memory data on data bus by the memory, control lines will be sending MFC signal to processor.

After receiving MFC signal by processor, the processor loads the data on the data bus into MDR.

→ Size of MAR is k bits \Rightarrow it can provide 2^k no. of addresses from the memory.

→ Size of MDR is n -bit, we can transfer n -bits at a time from memory & processor.

$$\text{Size of MDR} = \text{word length} = n\text{-bits}$$

Word length \rightarrow no. of bits that can be transferred b/w processor & memory at a time.

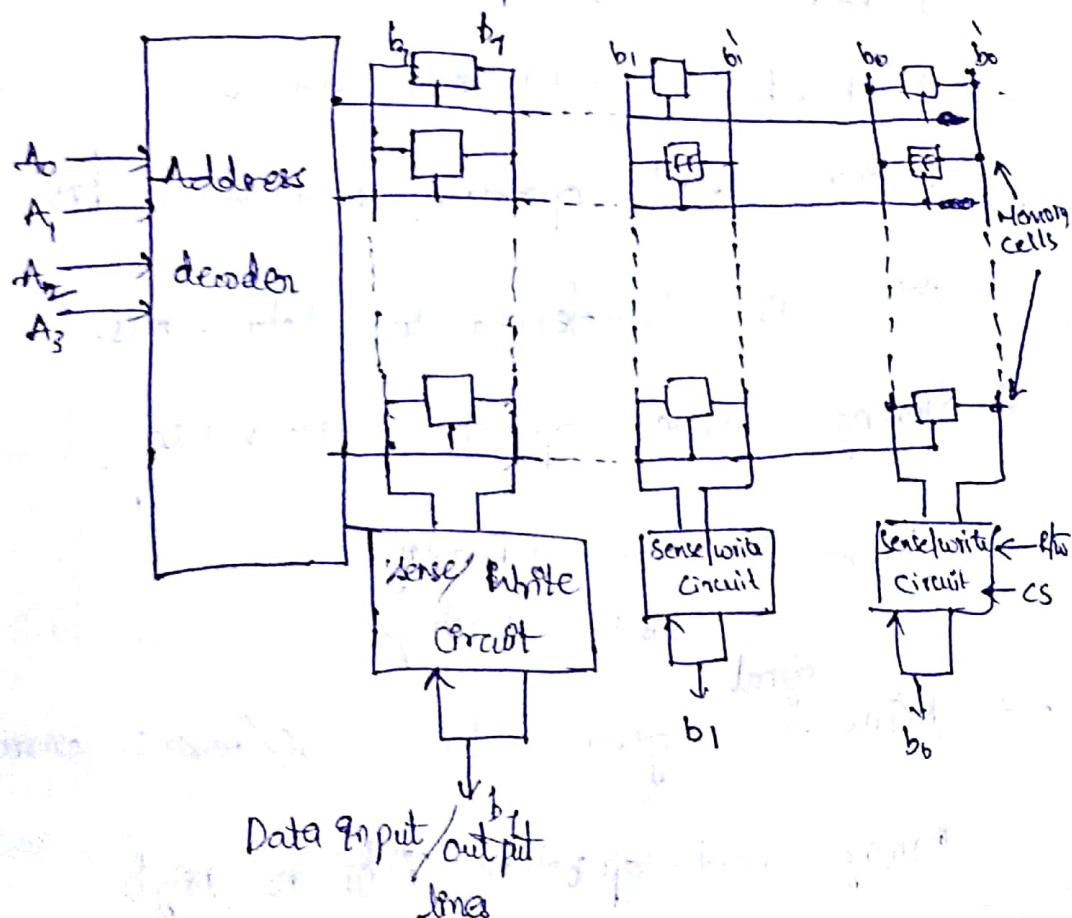
→ The no. of address lines depends on size of MAR.

→ no. of data lines in data bus
= no. of bits in MDR.

Semiconductor RAM Memories:

- if any location in a memory can be accessed in a fixed amount of time, it is called as Random Access Memory (RAM).
- If RAM is made up of semiconductor material, it is called as semiconductor RAM.

Internal organisation of memory chips:



→ For storing each bit, 1 memory cell is required

128 bits, 108 " "

- 1 memory cell is also called 1 flop flop
 - 16 bits are stored in 16 locations such that each location contains 8 bits.
 - 16 to 16 indicates 16 locations in the memory and each location contains 8 bits.
 - For accessing 16 locations, we require 4 bits for address ($2^4 = 16$)
 - After accessing some locations, other data is written into location by reading from data.
 - $b_7, b_6, b_5, b_4, \dots$ are called as bit lines.
 - $b_7, b_6, b_5, b_4, \dots$ → data lines (data bus)
 - During read operation, data from bit lines is transferred to data lines.
 - During write operation, vice versa.
- Sense → Sense → Reading
 Write → Writing
- P/T/R ^(signal) is given to sense/write circuit

During read operation, P/T/R is high

" write " , P/T/R is low.

Cs indicate chip select. When the memory

multiple bits, chop select pins used in order
to select a particular chop

Static Memories!

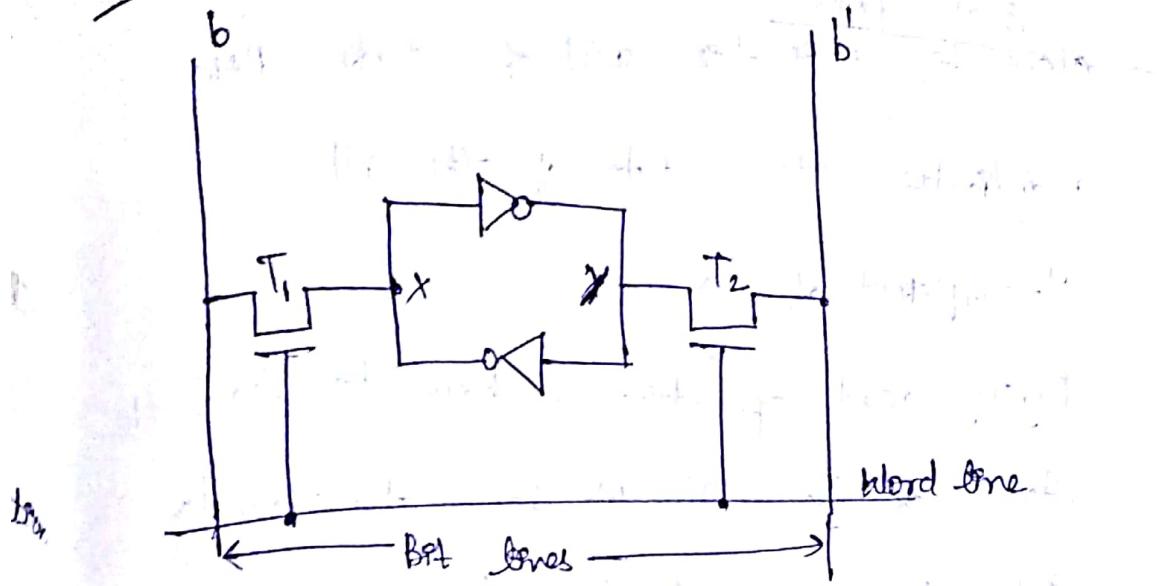


Fig. A static RAM cell.

→ In RAM, we can store information as long as power is on and it is lost when power is off.

→ It is of 2 types.

1) Static RAM

2) Dynamic RAM.

Static RAM!

(Cell is used to store one bit)

In this RAM, once information is stored, it

is returned as long as power is on.

~~Dynamically
selected~~

Dynamically PMA

In this PMA, cells are selected ~~postselectively~~

(presenting information) to produce information.

→ Above fig indicates cell of static PMA

X indicates the state of the cell.

Y-complement of X

During read operation, we have to search the state of the cell. We have to transfer information to data lines.

→ Word line is used to activate the word.

It indicates the amount of information stored in each block location function.

→ If the word contains 8 bits, 8 cells are required.

→ During read operation, word line is activated,

transistors T₁ and T₂ will be on and

b is transferred to b₁ and Y is transferred to b₁'. b & b₁' are called bit lines

which are connected to sense/write ch.

The purpose of sense/write cbt during read operation, is to transfer data from bit lines

to data lines and vice versa (during write)

During write operation, data is transferred to bit lines from data lines by activating word sense/write cbt and word line is activated and $b \rightarrow x$, and $b' \rightarrow y$.

→ ⁴ transistors for inverting and 2 for transferring are required. Totally, 6 transistors are required.

Advantage It is fast.

Disadvantage It is expensive.

If CMOS transistors for inverting, then power consumption will be less.

In dynamic RAM, there are of 2 types

1) Synchronous dynamic RAM

2) Asynchronous RAM

Difference between them is that

in synchronous and asynchronous

dynamical RAM, the clock is given to

synchronous dynamic RAM. In this RAM,

all operations are synchronised with clock.

Ast

Asynchronous Dynamic RAM (DRAM)

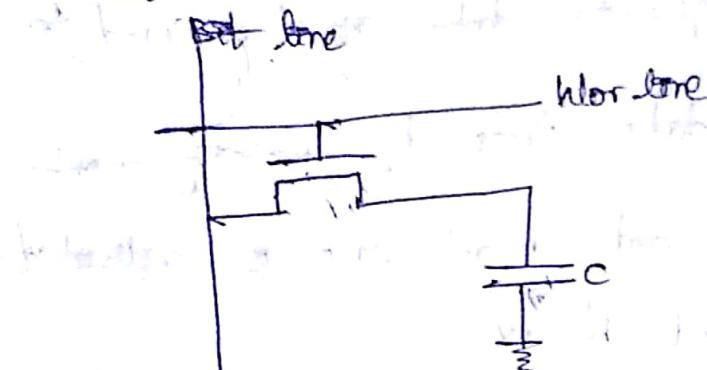


fig. A single ~~di~~ transistor dynamic

memory cell.

- It uses only 1 transistor
 - ∴ It is less expensive.
- When information is to be stored in cell,
transistor transistor is switched on and appropriate voltage is applied to bit line and capacitor is charged.
- When the transistor is off, capacitor discharges. ∴ Again capacitor should be recharged which is called to restore the value (to logic 1) which is called

periodical refreshing.

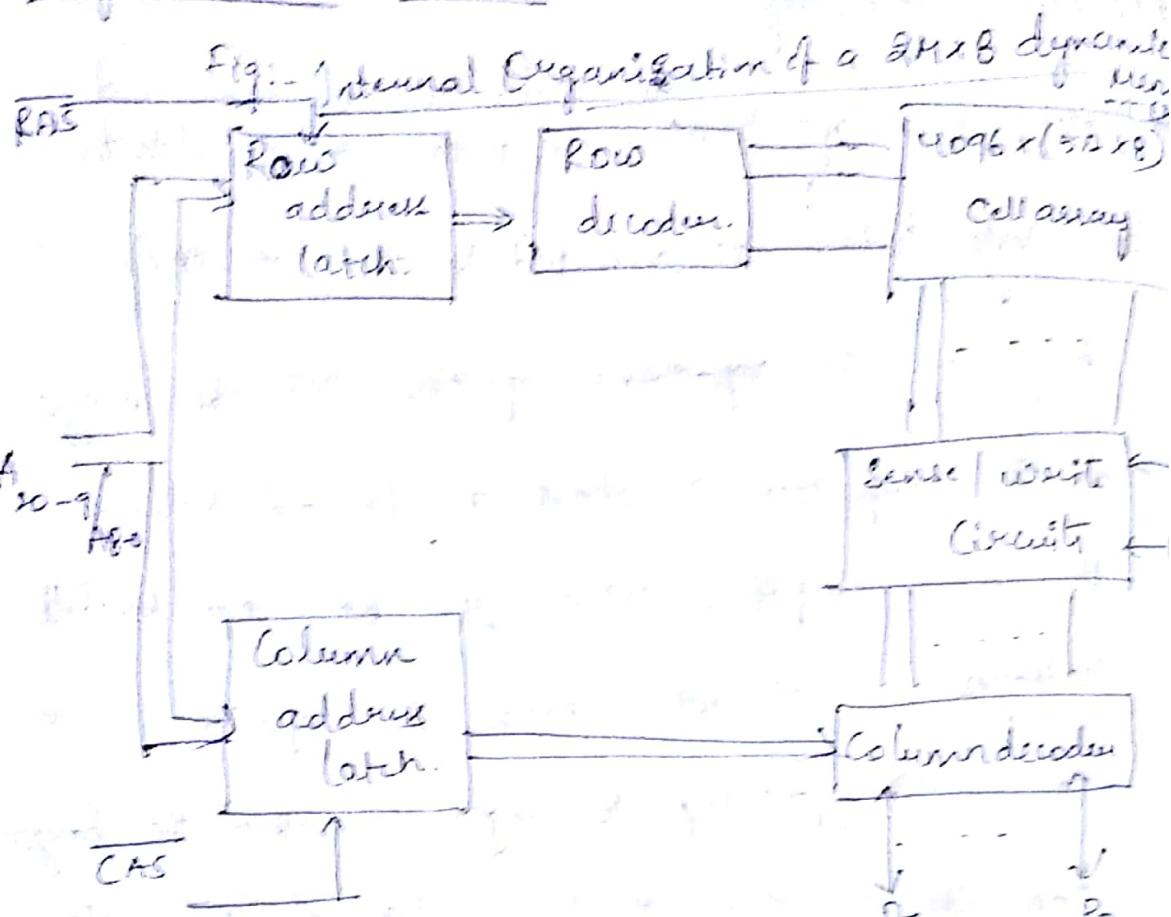
→ during read operation, sense/write circuit detects
the value which is connected to the bit line

The charge on the capacitor. If the charge
on the capacitor is above a threshold value,
the sense amplifier drives the bit lines to full
voltage \therefore it will charge the capacitor. " "
corresponding to logic 1. When the charge
on the capacitor is below threshold value, the
bit line is pulled to ground level which in turn
indicates no charge on capacitor. \Rightarrow logic 0
(Storing wrong information).

In order to avoid this, ^{cell} capacitor should
be periodically refreshed. It should be refreshed
before the charge on capacitor falls below
threshold values in order to retain the information
in the cell.

Asynchronous DRAMS

14/09/13



We have to design 16 Megabit dynamic Memory RAM
For this we need 16 Mega cells ($2M \times 8$). We are
Organizing 2 mega location which contains 8 cells
In each location total we get 16 cells. And we
need 2 mega addresses for 2 mega locations.

$$[1M = 2^{10} \times 2^{10}] \text{ For } 2M = 2 \times 2^{10} \times 2^{10}$$

- i) For addressing 2M locations we require $10 + 10 + 1 = 21$ address bits
 - (12 - Selecting row)
 - (9 - Selecting column)
- ii) Total 16M cells are arranged in a ~~2 rows~~ 4096 rows
and each row contains 4096 cells.

$$(4096 \times 4096 = 16M)$$

$$(512 \times 8)$$

- iii) Each ~~column~~ row we have 512 groups with 8 cells

$$(512 \times 8 = 4096 \text{ cells})$$

8 cells - 1 bit
for 512 bits \times 8 cells.
In each row we have 512 addresses.

- i. We have 512 columns in each row for selecting
ii) For a particular column we need 9 bits ($A_8 - 0$).
iii) For particular cell we should select the row first
and then column.
e.g. - 2 bit of 1 row 2nd column.
iv) $A_{20} - 19 \rightarrow$ indicates 12 address bits.

v) First row address loaded into Row address latch
by enabling RAS \rightarrow (Row signal which is in active low)
and row address decoded by the row address
decoder.

e.g. - If the Row address indicates address of 2nd row
the 2nd output of decoder is activated and then particular
row is selected.

vi) If we are performing read operation data $\xrightarrow{\text{in the row}}$
read by the Sense write (in read operation data
is transferred from Memory to data lines.)

In read operation \rightarrow acts like sense circuit

In write operation \rightarrow acts like write circuit

On sense of write select 4096 bits and the

Column decoder as 8 bits ($D_7 - D_0$) it selects

8 bits from 4096 bits (512×8)

\hookrightarrow (Column) for selecting
column we need 9 bits

vii) For selecting column we need 9 bits ($A_8 - 0$), that is
sent to Column address latch depending upon

the column address one of the column address is selected which is connected to column decoder.

During write operation data is transferred from data lines to Memory.

8 bits is transferred ($D_7 - D_0$) to 4096 bit pip to sense/write depending upon the column address.

From sense/write the circuit it is transferred to one of the row depending on row address.

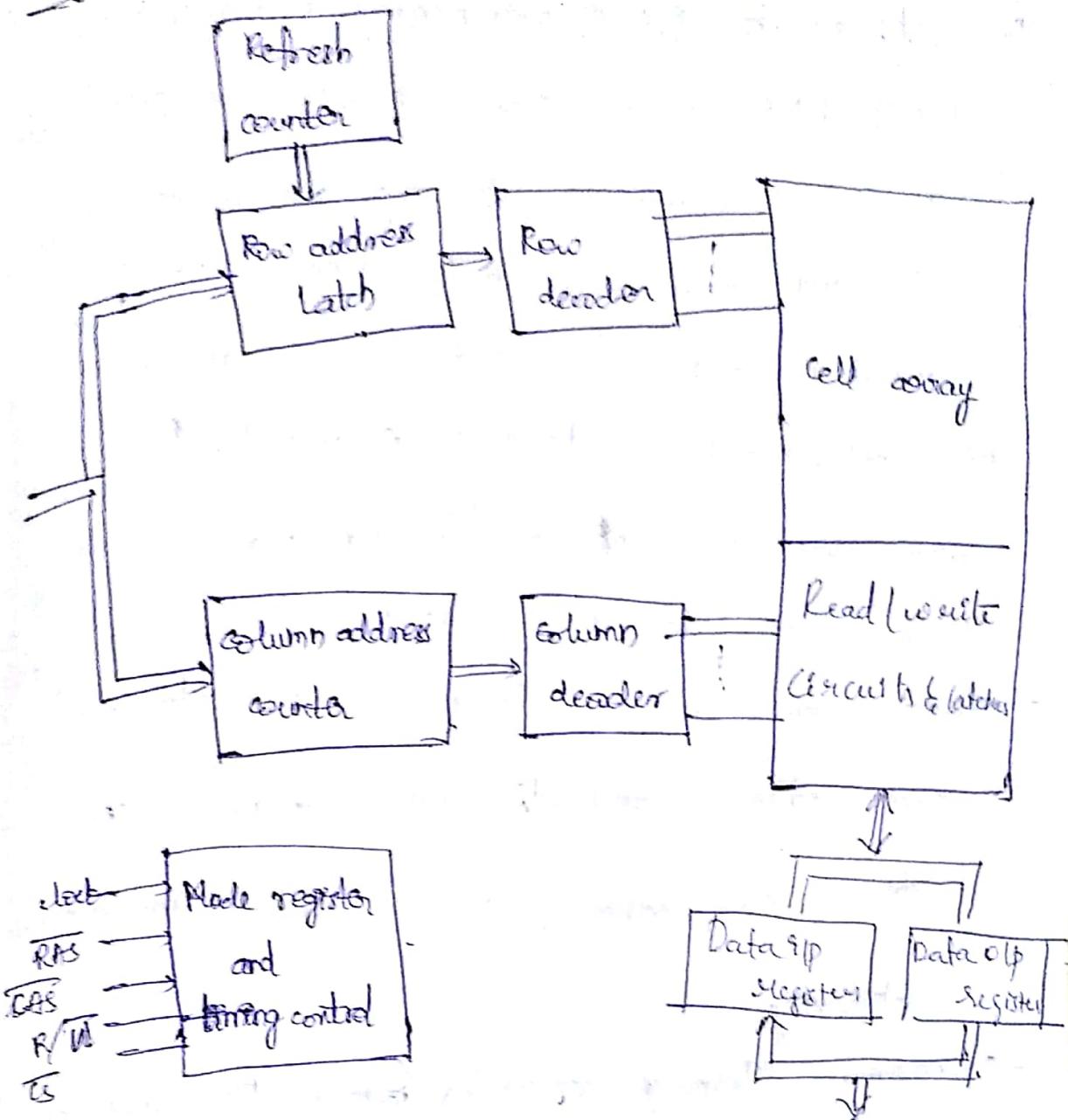
* When we perform block transfer (transferring multiple bytes) we have to transfer multiple bytes which indicates we have to transfer set of consecutive bytes (Consecutive bytes = bytes which have in consecutive address) for consecutive address row addresses same and column address to varied.

The bytes is in ~~regular~~ ^{particular} row indicates consecutive bytes with same row address and column address is varying in order to perform block transfer we have to transfer consecutive bytes.

For that purpose, we use latches at the output of sense/write circuit when we are using latches at the sense/write data is stored into latches.

When we are using latches at the output of sense/write circuit by providing only column address, we use consecutive bytes. This is called as fast page mode.

Synchronous DRAMs:



All operations are synchronised with clock in synchronous DRAM.

2Mx8 indicates it contains 2 Mega location and each location " " 8 cells.

Dynamiic RAM is used for main memory → more expensive
static RAM → cache → less expensive
→ more expensive (fast)

Speed of RAM is less as less as DRAM is used. As the size of main memory is more, more address bits are required. To decrease the no.

" ", we are dividing into two rows and, columns and it contains more than 1 word in each row. And we require row and column address to access a word. ~~first~~, ^{then} row address is sent and ^{then} column address. Hence, no. of address bits is reduced.

→ When same address lines are used to send row address and the column address one after the other, it is called ~~row~~ address multiplexing.

→ ~~Memory~~ Memory controller (among the memory and processor) divides 24 address bits in 12-row addresses and 9-column addresses.

→ For accessing two consecutive words in a row, column is incremented. Column address counter is used (per row). → Mode register is used to select the mode of operation of DRAM.

Block transfer indicates multiple word transfer.
During block transfer, mode register indicates
the no. of successive bits to be transferred (during
block transfer).

→ speed of transfer b/w main memory and processor
depends on speed of main memory and speed of bus
speed of main memory (DRATA) is expressed in terms

of latency (time taken to transfer one word
from/into memory). → Band width is another quantity
used to measure the performance of main memory,
which is used to measure speed of
during block transfer.

It indicates no. of bits or bytes transferred per
second.

Cache

Cache Memories:

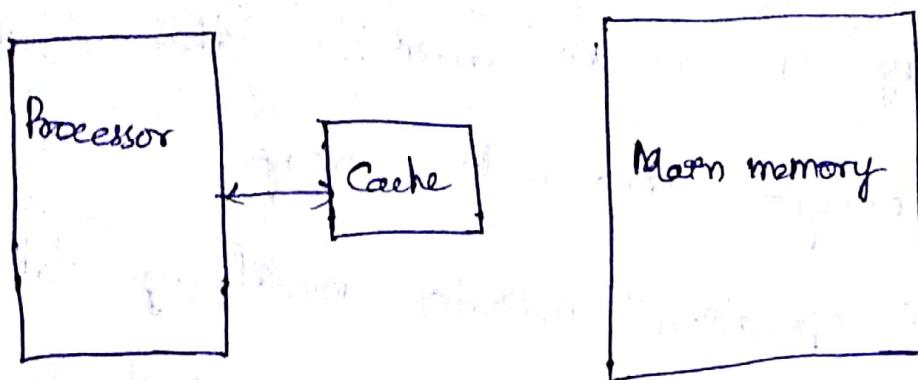


Fig. Use of a cache memory.

Main memory and cache indicates RAM (Read or write)

Main memory is DRAM

Cache memory is static RAM.

Main memory contains programs and data. As the size of main memory is large, to decrease the cost, DRAM is used. It is slow (DRAM) and performance is decreased (depends upon time of execution). To overcome this, first

time to read/write a memory, it is brought into cache memory (High speed static RAM).

When the word is read second time, it

directly read from cache memory (already

available).

During read operation, it is read from cache memory if it is available otherwise, it is brought into cache memory.

Write operation updates modifying data in some location. We have 2 protocols to update the data.

Write through protocol.

In this protocol, if the word which is to be modified in main memory is also available in cache memory, it is modified in both cache and main memory simultaneously.

If it is not available, it is modified only in main memory.

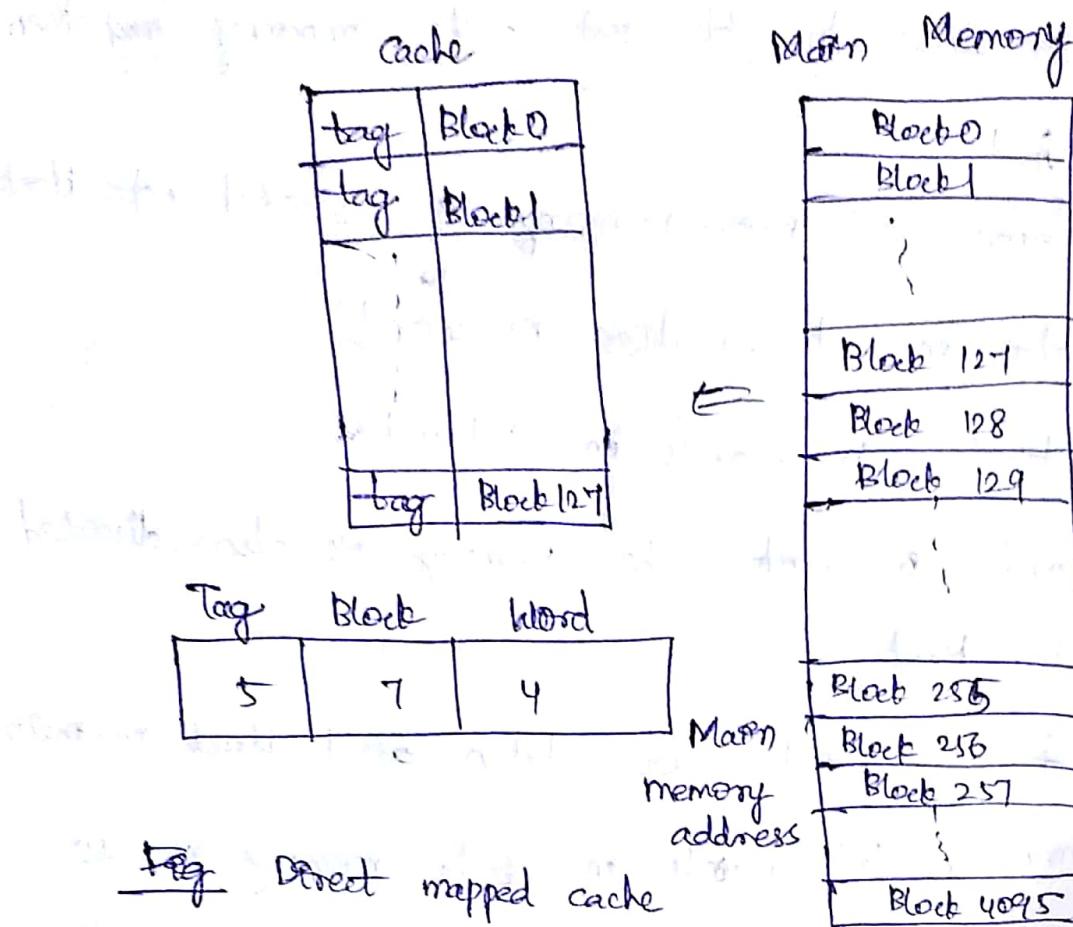
Write back protocol:

In this protocol, if the word to be modified is available in cache memory, it is modified only in cache memory initially.

- If the cache is full, word is removed from cache and then modified in main memory.
- If the word is not present in cache memory, it is first brought into cache memory and then modified.
- Words in main memory is divided into blocks
(Data in each location is word)
 - Eg: If there are 16 words in 1 block
- Words in cache memory is also divided into blocks.
- If 'n' words are taken as 1 block in main memory, 'n' words in cache memory is also taken as a block.
- If some word of main memory is to be transferred into some block of cache memory from main memory, the entire block is transferred.
- Mapping function decides where the block of main memory is to be placed in which block of cache memory.

Mapping functions:

Direct Mapping:



Tag Direct mapped cache

Here, each block contains 16 words.

$$\text{Total no. of words} = 4096 \times 16$$

Dimension of each block = 64 K locations

$$= 64 \times 1024 \\ = 2^6 \times 2^{10}$$

∴ 16 address bits are required to access 64 K locations.

Cache memory contains 128 blocks

$$\therefore \text{No. of words} = 128 \times 16$$

$$= 2048 \text{ words} \\ = 2^10 \text{ blocks}$$

Block 1 of main memory is placed in Block "modulo 3" of cache memory.

Block 0, 128, 256 of main memory is placed in Block 0 of cache

if, Block 1, 129, 257 ... in Block 1 ...

→ 16 bits in the figure indicates address of word of main memory.

→ 7 bits contains the block of cache memory where the word may be present.

→ Tag bit indicates which block of main memory is contained in block of cache (0, 128, 256)

→ If the tag bit in main memory matches with " " in cache memory, then it indicates that the required word is contained in cache memory (It is brought into cache if it is not present).

Associative mapping:

Any block of main memory can be placed in " " cache " "

Main memory address indicates the word of " " which is to be accessed.

Word is the data stored in a location

→ Tag bits of block of cache memory indicate which block of main memory is contained in that block.

→ If the tag bits of main memory matches with the tag bits of cache memory, it indicates that the word is present in cache memory.

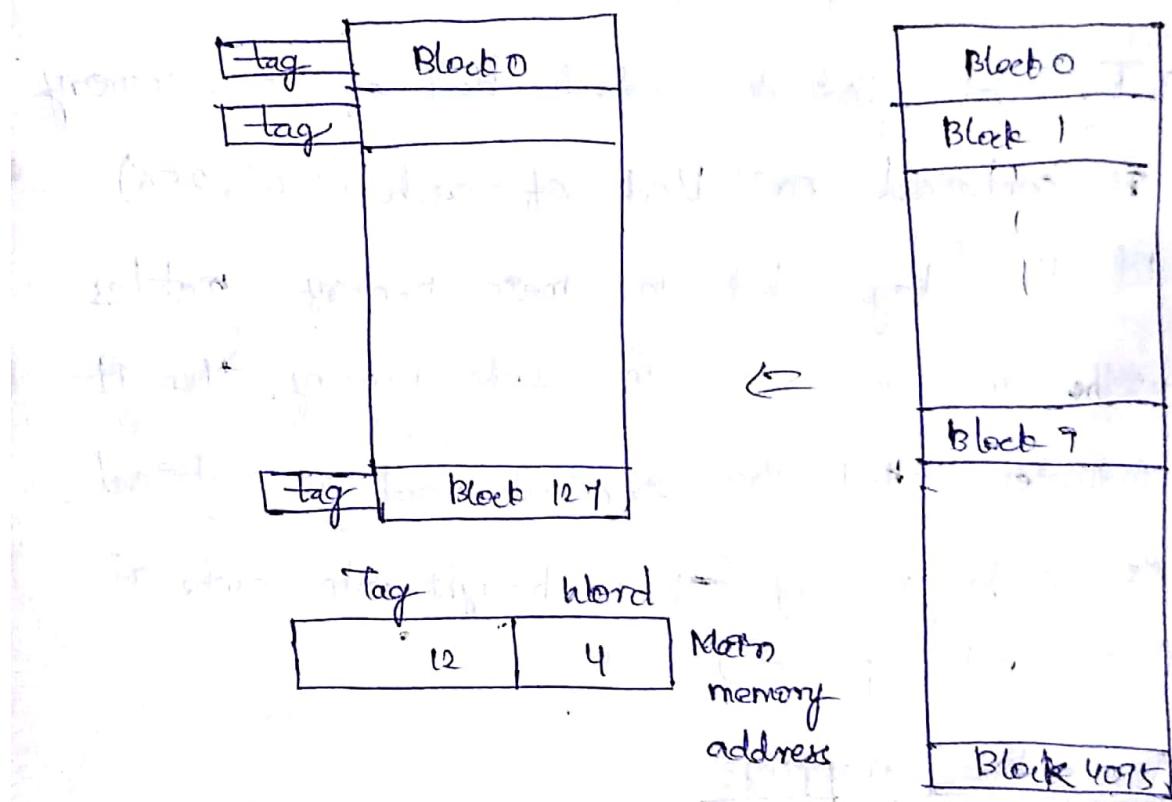


Fig. Associative mapped cache.

→ When new block is to be brought into cache memory (from main memory), one block of cache memory should be removed when the

cache is full.

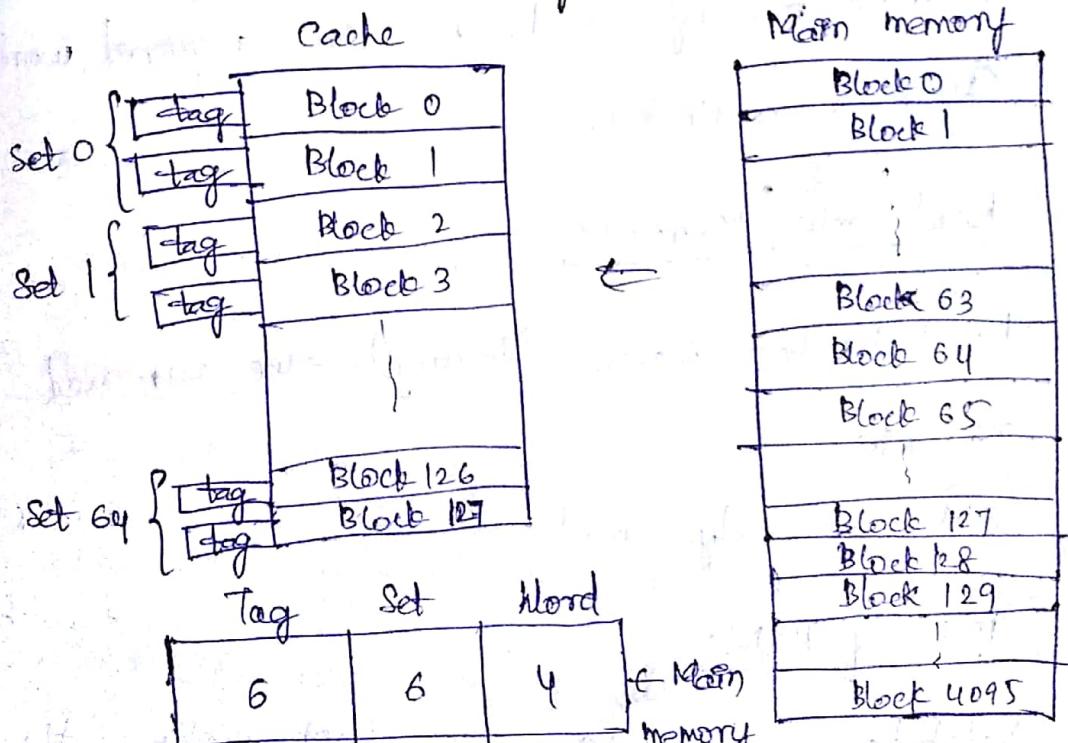
→ Algorithms which decides which of block of cache memory is to be replaced when the cache is full are called replacement algorithms.

→ One of the algorithms is LRU (Least Recently Used) block (which is accessed long back.)

→ Another algorithm is oldest algorithm (where the oldest block is removed.)

→ 4-bits in main memory address indicates which word in the block (16 words in a block) is to be accessed.

Set-Associative Mapping



Tag: Set associative mapped cache memory with two blocks per set.

→ In this mapping, in cache memory, some no. of blocks are taken as a set.

Here, each set contains 2 blocks

$$128 \text{ blocks} \Rightarrow 64 \text{ sets}$$

Block 0, 64, 128, ..., 4096 are placed in set 0 of cache memory (Block 0 or Block 1)

→ 16 bit main memory address is divided into

3 (parts) fields

6 bits ($2^6 = 64$) are used to indicate set no; set field.

→ In the corresponding set (e.g. set 0) the tag bits which matches indicates # with the tag bits in main memory add. is the required word (Block 0 or Block 1)

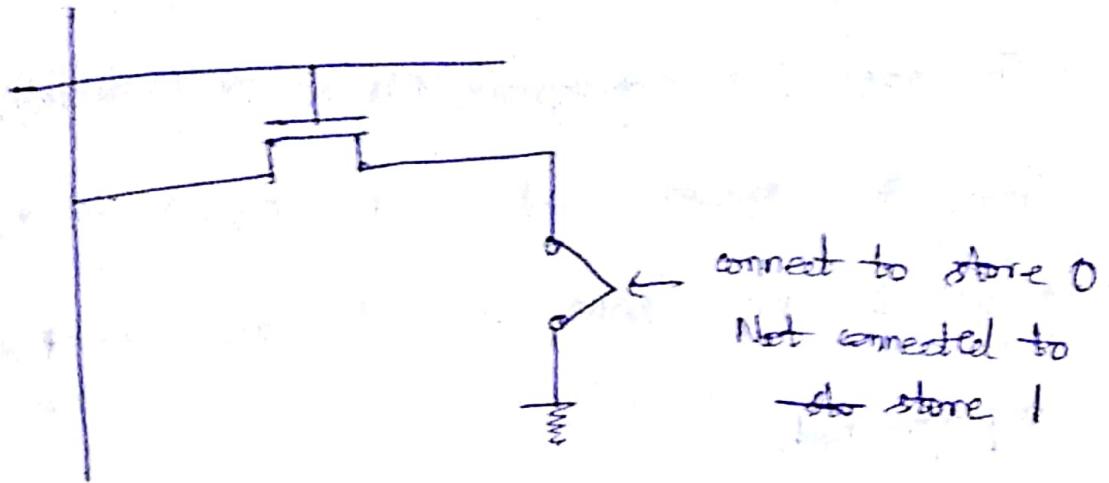
Read only Memory:

RAM (Random Access Memory) → we can read or write.

ROM can only read.

Types of ROM:

Volatile → information is lost when the power is off



~~RAM~~ → volatile

~~ROM~~ → non-volatile

To store bit 0, terminals at P are connected
to Vcc. To store bit 1, a fuse is not connected.

1) ~~ROM~~ → it is programmed by the manufacturer

2) ~~PRAM~~ → a user can write over it.

It is similar to ROM but all cells contain fuses b/w the terminals at point 'P'. Initially '0' is stored in PRAM. To store 1, fuses are burnt at point 'P' by the user by using high current pulse.

3) EROM

In this type of memory, we can erase the cells and reprogram the cell. Whereas, in ROM and ~~PRAM~~ PRAM, it is not erased once it is programmed.

To erase and reprogram, EPROM is removed from the circuit and it is exposed to UV rays. In this ROM, all cells are erased and reprogrammed.

4) EEPROM (Electrically Erasable PROM)

→ In order to erase and reprogram, it need not be removed from the circuit.

Select cells can be erased and reprogrammed.

Disadvantage: Diff. voltages are used to read, write and erase.

5) Flash memory:

It is similar to EEPROM.

We can read single cell but can write into blocks of cells.

Advantage: Low power consumption, uses only single power supply voltage, high density,

(density = no. of cells that can be accommodated in a given area) low cost per bit (cost of each cell is less)

Application: battery driven applications (due to low power consumption)

Virtual Memories

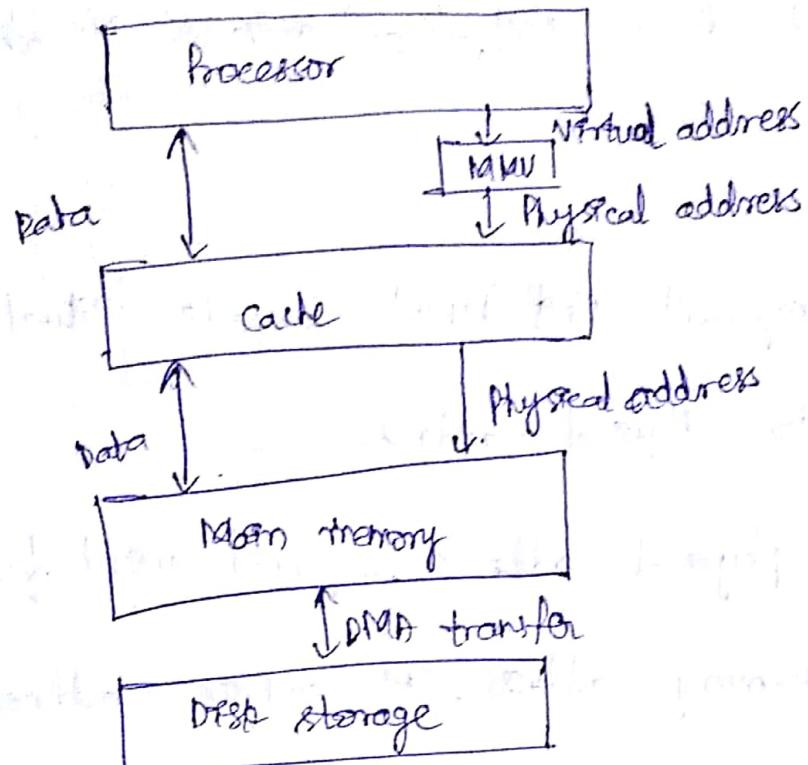


Fig Virtual memory organisation

- Initially, files will be stored in disk.
- To execute that file, it should be brought into main memory. For faster access, cache memory is used i.e. data from main memory is transferred to cache.
- Transfer b/w disk and main memory is done in terms of pages where each page contains 2K - 16K words.
- Initially, the address generated by processor is called virtual address. It contains page num

and offset. Page num indicates in which page required word is contained, offset indicates Pn where location is

→ Memory Management Unit (MMU) converts Virtual address into Physical address

→ Using physical address, required word from main memory address. It contains address of page in the main memory

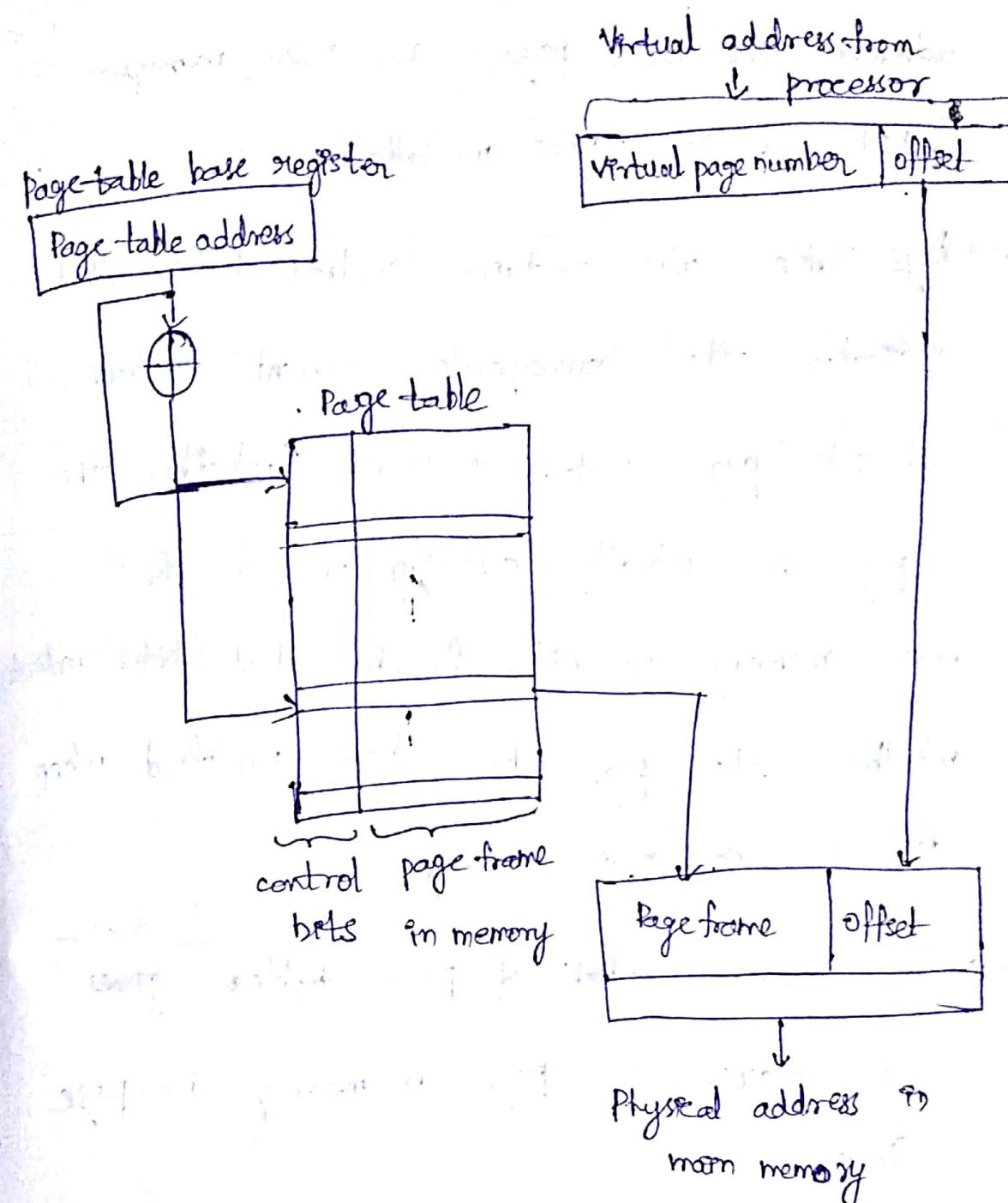
E.g. If virtual address gives page 0

location of page 0 is given by physical address.

→ Physical address contains page frame and offset. Page frame indicates the starting location of page in the main memory and offset indicates the location of word in page

Address translation:

(It is the process of converting virtual address into physical address)



→ Page table contains the address of each page in main memory.

→ Page table base register contains the starting address of the page table.

Ex: If page no. is 2, 2 is added to page table base register value and it gives the

address of the page) in main memory.

(which is in Page table)

→ Page Table also contains control bits which indicates the corresponding current status of the page. Status indicates whether the page is actually existing or in the main memory or not. By, the control bits indicate whether the page has been modified when it is in main memory.

→ The location of page tables gives the address of page in memory i.e. page frame.

→ Starting address of page in main memory is called page frame.

Eg: If page frame = 1000, offset = 4

word in 1004 is accessed.

→ Page table is present in main memory.

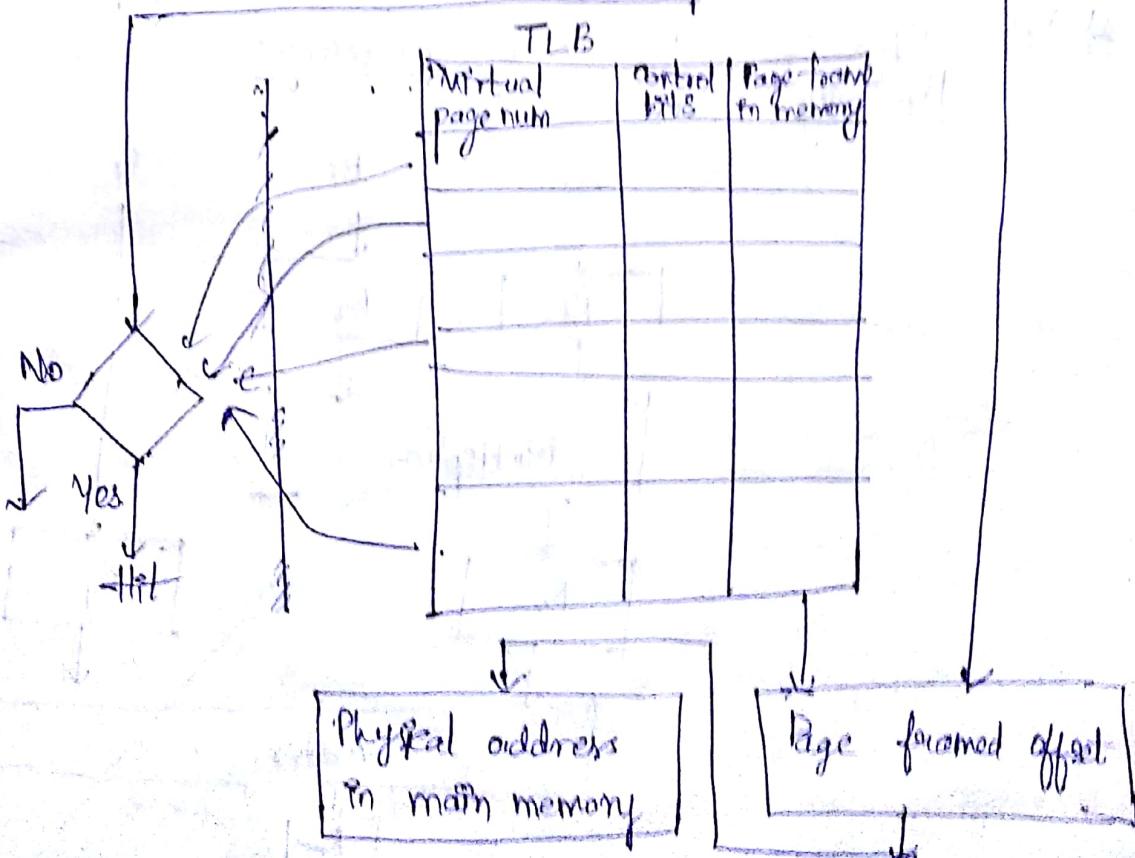
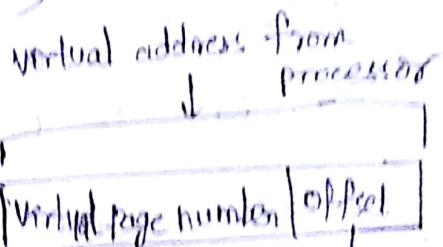
→ Conversion of virtual to physical address is done by MMU.

As page-table is contained in main memory, accessing-time of ^{converting} virtual address to physical address will be more.

In order to decrease the time for conversion from virtual to physical address, a portion of the page-table is placed in RAM in the cache memory. The portion of the page-table which is contained in RAM is called TLB.

(Translation Lookaside Buffer)

Virtual Memory:



TLB contains page table entries of recently accessed pages.

Eg: If TLB contains page 1, 4, 5, then virtual page num will be 1 or 4 or 5.

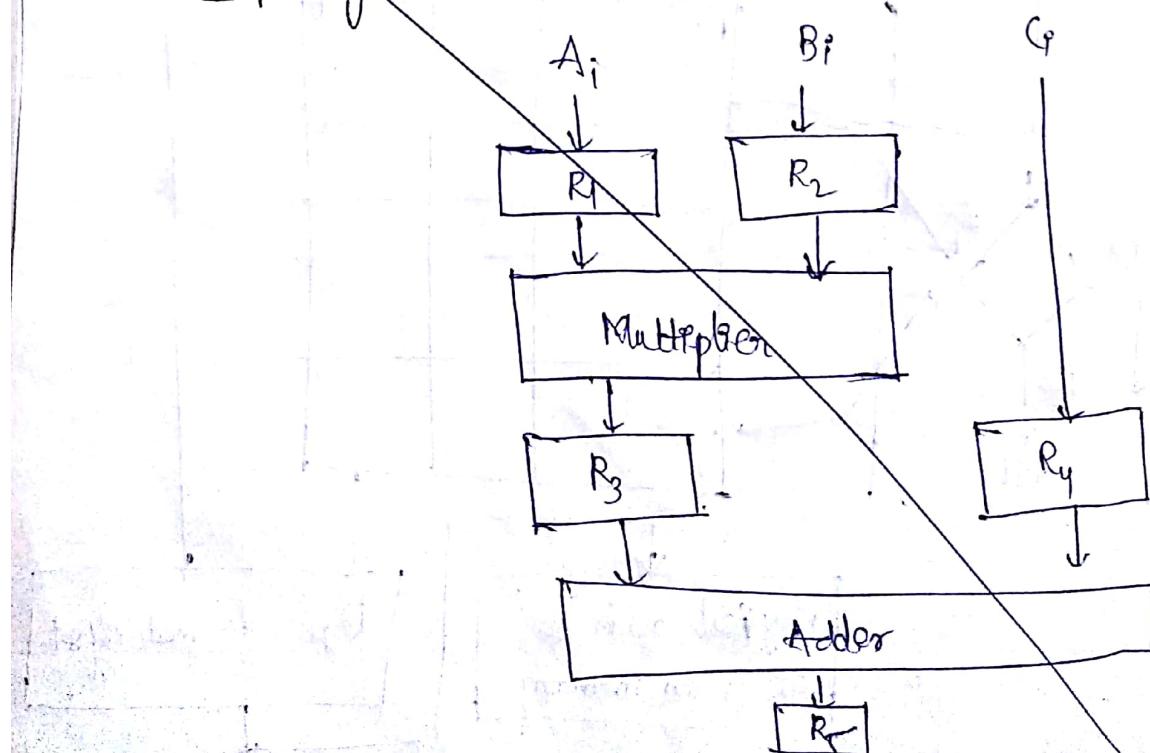
→ Virtual page num in virtual add. is compared with virtual page num in TLB.

If it matches with virtual page num, the required info. is contained in TLB.

Or, if it is not contained in TLB, it has to be taken from page table.

TLB contains only recently accessed pages.

Pipelining and Vector processing:



UNIT-IV

Pipelining and Vector Processing.

Pipelining

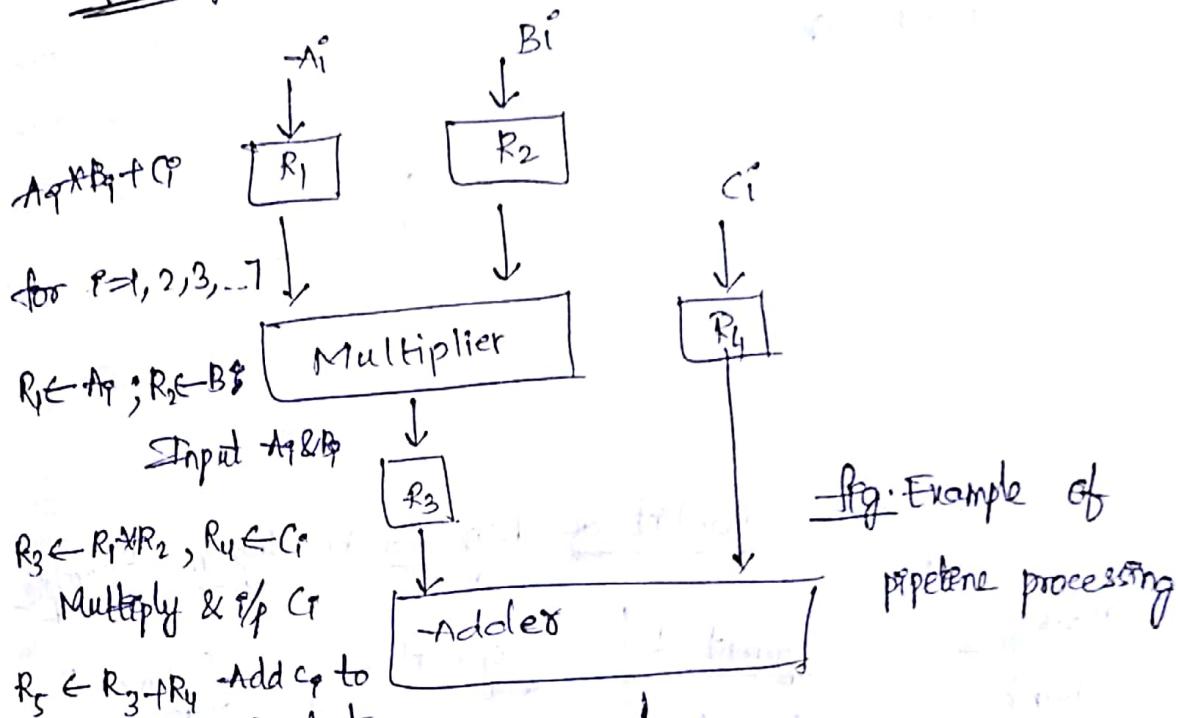


Fig. Example of pipelining processing

It performs operation on seven different data. Difference b/w pipelining & Non-pipelining ckt is non-pipelining ckt will process only one data at a time.

- After completion of processing of 1st data, processing of 2nd data is performed. In non-pipelining.
- In pipelining, we can process more than 1 data simultaneously.
- To perform $A_i * B_i + C_i$, this operation is divided into suboperations or segments

→ Here, it is divided into 3 segments.

→ During the first clock cycle, first suboperation is performed.

$$A_7 \# B_7 + C_7$$

Contents of Register in Pipeline Example

Clock pulse number	Segment 1		Segment 2		Segment 3	
	R ₁	R ₂	R ₃	R ₄	R ₅	
1	A ₁	B ₁				
2	A ₂	B ₂	A ₁ #B ₁	C ₁		
3	A ₃	B ₃	A ₂ #B ₂	C ₂	A ₁ #B ₁ +C ₁	
4	A ₄	B ₄	A ₃ #B ₃	C ₃	A ₂ #B ₂ +C ₂	
5	A ₅	B ₅	A ₄ #B ₄	C ₄	A ₃ #B ₃ +C ₃	
6	A ₆	B ₆	A ₅ #B ₅	C ₅	A ₄ #B ₄ +C ₄	
7	A ₇	B ₇	A ₆ #B ₆	C ₆	A ₅ #B ₅ +C ₅	
8			A ₇ #B ₇	C ₇	A ₆ #B ₆ +C ₆	
9					A ₇ #B ₇ +C ₇	

- In pipeline implemented, after performing $A \oplus B + C$, processing of data is performed one after the other (A_1 and B_1 are loaded).
- Whereas in ~~non-pipeline~~ pipeline, During 2nd clk pulse, 2nd suboperation on 1st data is performed simultaneously 1st suboperation on 2nd data is performed.

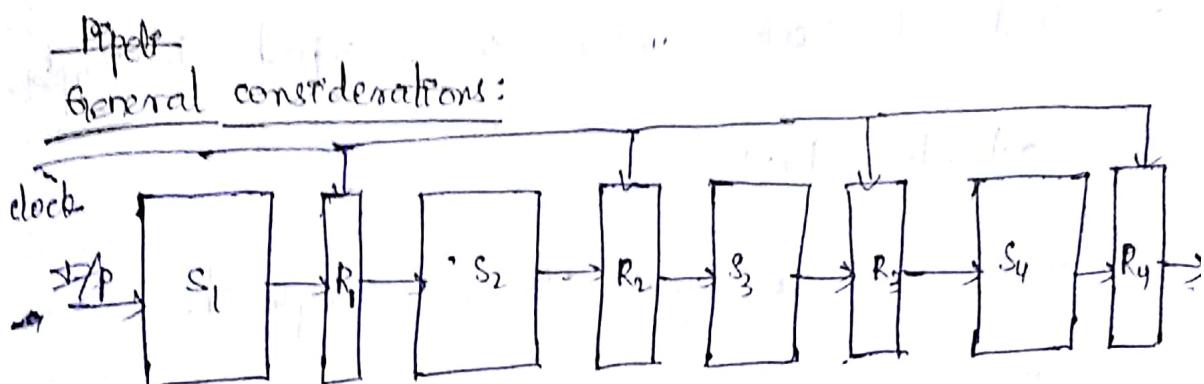


Fig. Four segment pipeline

- Four segment pipeline indicates 4 suboperations.
- In the previous eg, $n=7$ (7-tasks)

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

- Here, 6 tasks are performed $\Rightarrow n=6$
- Assume that
 - These tasks indicates the same operation $(A \oplus B + C)$ but on diff. data.
- 6 tasks are denoted as T_1, T_2, \dots, T_6
- All these 6 tasks should go through the 4 segments

Space-time diagram for pipeline:

	1	2	3	4	5	6	7	8	9	
1	T_1	T_2	T_3	T_4	T_5	T_6				→ clock cycle
2		T_1	T_2	T_3	T_4	T_5	T_6			
3			T_1	T_2	T_3	T_4	T_5	T_6		
4				T_1	T_2	T_3	T_4	T_5	T_6	

→ Total 9 clock cycles are required to complete all the tasks.

→ After 4 clock cycles, T_1 is completed,

5 u u T_2 u u

6 u u u u u

7 u u u u u

8 u u u u u

9 , T_6 u u

→ In 4+5 clock cycles, 6 tasks are completed. Here 4 indicates no. of segments.

→ $k+n$ clock cycles are required

to perform 'n' tasks, where k indicates no. of segments.

→ In nonpipelining, t_n if ' t_n ' indicates the time taken to complete one task,

the time taken to perform 'n' tasks is

t_{fn} .

's' indicates the speed up ratio which gives the diff in speed of pipelining and nonpipelining.

$$s = \frac{t_{\text{fn}}}{(k+n)t_p}$$

Here t_p indicates clock period.

→ If $s=1$, the time taken by nonpipelining is 4 times the " " by pipelining.

→ If n is large, $n > k+1$
 $\therefore K+1$ can be neglected.

$$s \approx \frac{n t_h}{n t_p}$$

$$s = \frac{t_h}{t_p}$$

$$\cancel{t_h} = t$$

The time taken by pipelining obt to complete
~~(first task)~~ single ~~t_h~~ = ~~$k + t_p$~~

k = no. of segments in pipelined pipelining

If we assume the time taken to

complete one task in pipe nonpipelining =
time taken to complete one task in Pipelinin

$$t_n = t_p$$

$$t_n = k t_p \quad | \quad t_p$$

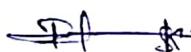
nonpipelining Pipelining (single task)

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

→ If S is the maximum theoretical speed that can be achieved using pipelining cbt.

Practically $S \leq k$ ($S \leq b$)

where k indicates no. of segments.



→ If $k=4$, speed of pipelining is 4 times greater than a nonpipelining

→ $S=k$ when t_p is large and $t_n = k t_p$

→ n' no. of nonpipelining cbts are required to achieve the speed of pipelining cbt

$S=4$ indicates 4-segment pipeline

→ for n' tasks, speed of pipelining $k \times$ speed of nonpipelining

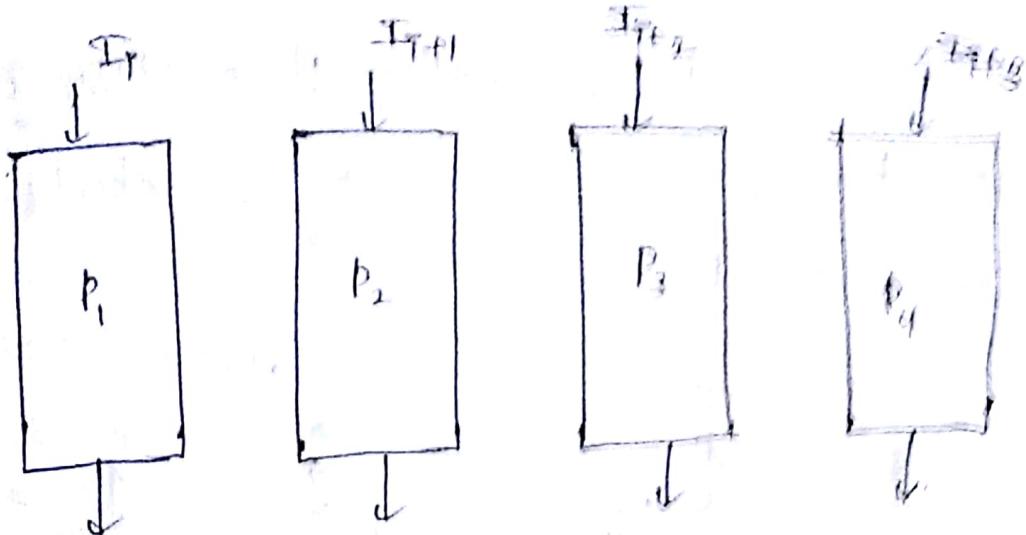


Fig Multiple functional units in UEL

(P_1, P_2, P_3, P_4)

Here 4 \downarrow nonpipelining dots are required to achieve the speed through pipelining.

-Arithmetic pipeline:

Here, adding floating pt. no's are added or subtracted using pipelining.

Adding / subtraction of float floating point no's is

divided into 4 sub operations.

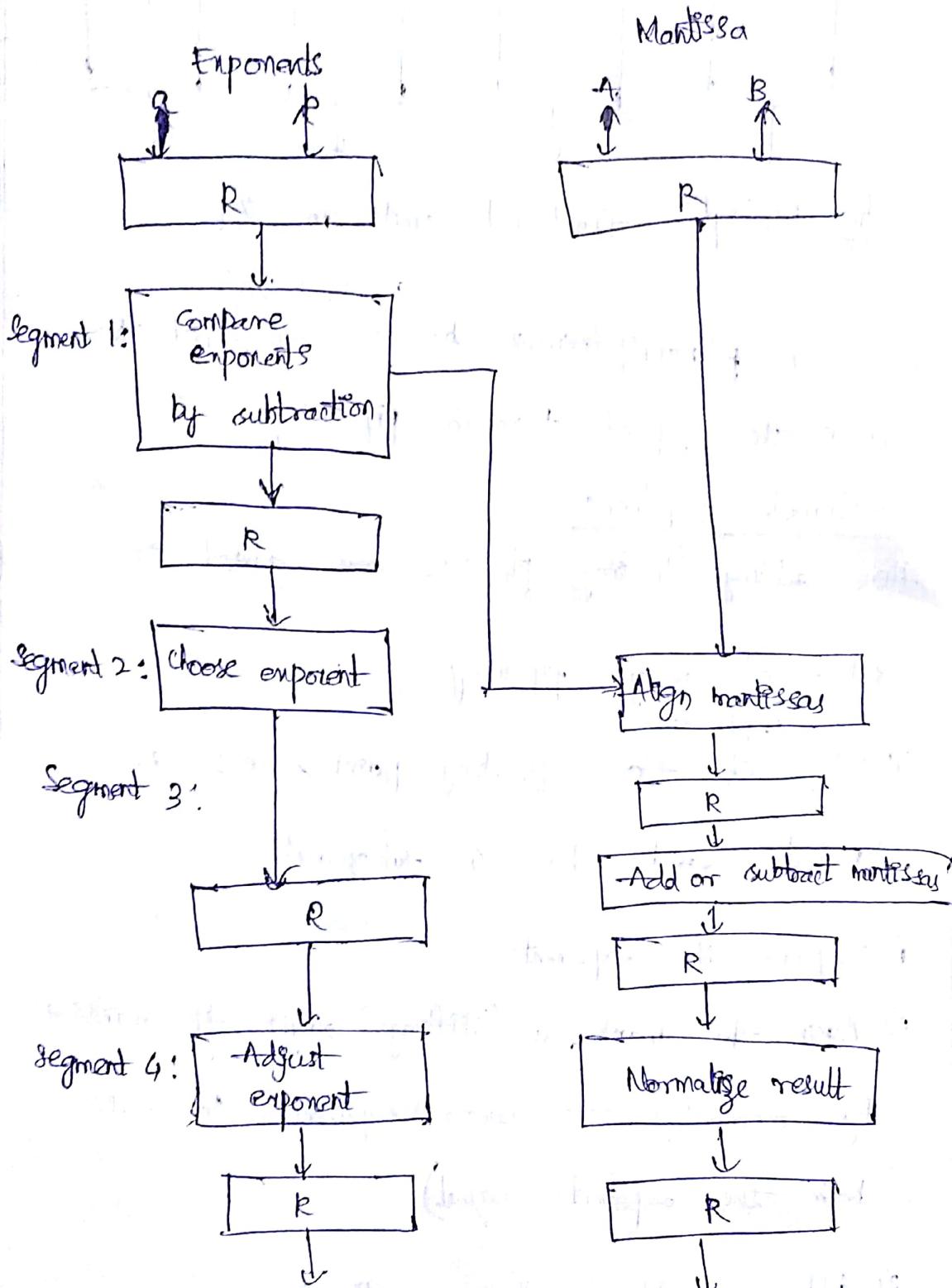
1) compare the exponents.

2) Align the mantissa (shifting right the mantissa by increasing the lower exponent to make both the exponents equal).

3) Add or subtract the mantissas

4) Normalize the result. (Most significant digit of floating pt. is non zero \rightarrow non normalized)

In case binary floating pt no, MSB bit of binary number should be equal to 1.



Instruction Pipeline

It is used to process multiple instructions simultaneously simultaneously. In order to process the instructions, following steps are to be followed are:

1) Fetch the instruction from memory

2) Decode the instruction

3) Calculate the effective address. (to find the location of operand.)

Eg: In case of branch instruction, eff. address indicates where to branch.

4) Fetch the operands from memory

5) Execute the instruction

Eg: Add, Mul, data transfer etc

6) Store the result in proper place

Eg: In case of addition, result is stored in memory of register

These steps are divided into 4 segments.

For condition and BR instruction, condition is

verified in segment 4. It will verify whether the condition is true. Processing takes place

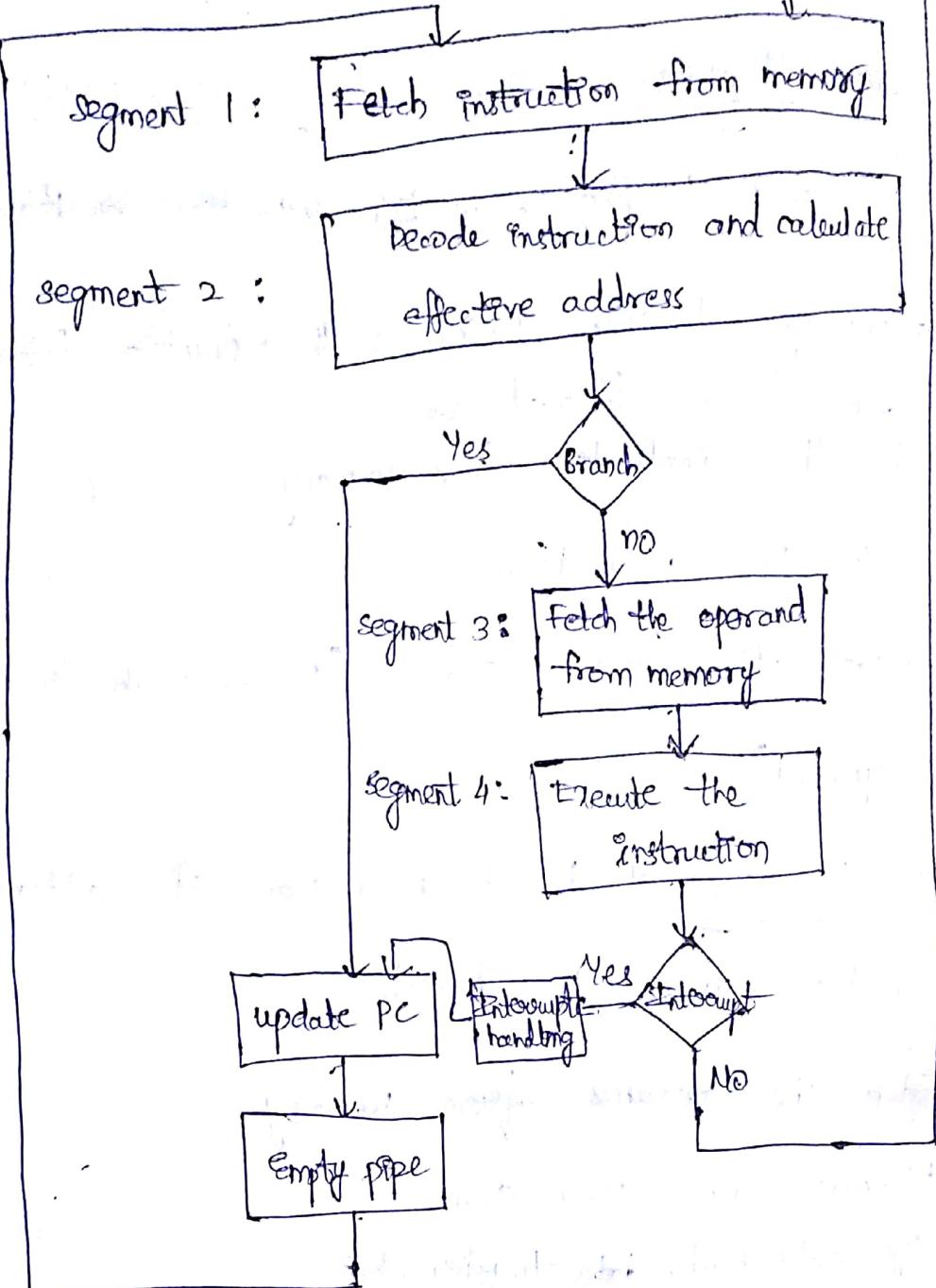


Fig Four-segment CPU pipeline

from branching address. After executing some instruction, before executing next instruction, it will check whether there is any interrupt. and if there is an interrupt,

ISR (Interrupt service Routine) is executed,

and PC is updated with starting address of PSR. For
 for unconditional BR, PC is updated with branching address, and pipe is emptied. Branching instruction is prefetched (before) before and condition is checked.

If the condition is false in case of cond. branching, the next instruction after BR instruct is executed

Step	1	2	3	4	5	6	7	8	9	10	11	12	13
Instruction 1	FI	DA	FO EX										
		FI	DA	FO EX									
(Branch)			FI DA	FO EX									
			FI	- - FI	DA	FO EX							
4				- - - FI	DA	FO EX							
5					- - - FI	DA	FO EX						
6						- - - FI	DA	FO EX					
7							- - - FI	DA	FO EX				

fig. Timing of instruction pipeline if the decoded instruction is BR instruction,

the next instruction is halted till the

execution of BR instruction. If the condition

is true, the first segment of instruction in branching add. is executed. Old, the fetched

~~instruction~~ is decoded. During decoding of first instruction after branching, second instruction after branching is decoded.

diff

Pipeline conflicts:

which arises

Problems in pipeline are called pipeline conflicts

1) Resource conflicts

→ When 2 segments of instruction pipeline

try to access the memory at the same time, it is called as resource conflicts.

(depends on instructions in pipeline)

2) Data dependency conflicts

→ When an instruction depends upon the result of previous instruction which is not yet available, data dependency conflict occurs.

Eg: The instruction is delayed till the result in previous is obtained.

→ In order to delay the instruction,

a hardware bit is used called as hardware interlocks.

Another way is to delay the instruct^h, i.e.,

by using compiler which inserts the No operation

instruction to delay the instruct^h. This is called

called as delayed load.

Instead of storing the result of previous

instruct^h in destination register from which

the next instruct^h has the access, it is directly forwarded by bypassing the register

forwarded to next instruct^h. The operation is called

as operand forwarding.

3) Branch

3) Branch difficulties:

During pipelining, operation of pipelining is

interrupted or halted which is called as branch

penalty. It degrades the performance of pipeline

In order to reduce minimize the degradation of

the performance of pipeline during branching

instruction, we will be prefetching the instruct^h's

- From branching address before the execution of instruction.

- the "n" instruction, to be done through hardware. It is very easy to handle.

Branch difficulties

→ easier to have a continuous flow of pipeline.

→ during branching instruction, compiler will

be either rearranging the instructions or

inserts the "No operation constraint", which

is done in RISC computer. This technique

is called as delayed branch

Examples for delayed load and delayed branch

RISC Pipeline:

Instruction pipeline corresponding to RISC computer is called

RISC pipeline. It is divided into 3 segments

segments

D : Instruction fetch

A : ALU operation → decoding and performing ALU operation.

E : Execute instruction

→ RISC computer consists of 3 types of

instructions.

1) Data manipulation instructions Ex: Add (+, -, *, /)

2) Data transfer Ex: load and store

3) Program control Ex: branch

→ During data manipulation instruction, we perform on ALU

data manipulation operation in 2nd segment and.

result of ALU is transferred to destination register in 3rd segment.

→ For data transfer instructions, (load or store)

in the 2nd segment, effective address is

calculated for load or store instruction

Ex: $M[x] \rightarrow R$

$\underbrace{\text{eff}}$
 Add

During 3rd segment, effective add. is transferred to the data memory for load or store operation

→ For program control instruction, during

2nd segment, branch address is calculated by ALU. During 3rd segment, branch add. is

transferred to PC.

Delayed load:

Eg:-

① LOAD : $R_1 \leftarrow M[\text{address } 1]$

② LOAD : $R_2 \leftarrow M[\text{address } 2]$

③ ADD : $R_3 \leftarrow R_1 + R_2$

④ STORE : $M[\text{address } 3] \leftarrow R_3$

Pipeline

clock cycles	1	2	3	4	5	6
1. Load R_1	S	A	E			
2. Load R_2		S	A	E		
3. Add $R_1 + R_2$			S	A	E	
4. Store R_3				S	A	E

a) Pipeline timing with data conflict
During 4th clock cycle, the value of

R_2 is not ready to be as R_2 is

loaded into destination register and data

dependency conflict occurs due to in

$R_1 + R_2$ and the 3rd instruction is delayed

until R_2 is loaded by compiler by first inserting
NO instruction b/w 2nd and 3rd instruction,
∴ In 5th clock-cycle, R_1+R_2 is performed.

clock cycle	1	2	3	4	5	6	7
1) Load R_1	S	A	E				
2) Load R_2		S	A	E			
3) No operation			S	A	E		
4) Add R_1+R_2				S	A	E	
5) Store R_3					S	A	E

→ This technique is called as delayed load.

7/10/17

Vector Processing:

→ Vector is a collection of data items
of same or diff type (array-name type)

→ Vector processing deals with operations
on vectors or arrays.

Vectors

→ Vectors

→ There are two types of processors to
perform any operation.

Scalar processor:

In scalar processor, we can process

only single data item at a time.

(No arithmetic pipeline.) ⇒ only single data item can be processed

2) Vector processor

ST processes multiple data items at a time.

It uses aerometric pipeline

→ speed of operation is more in vector

processor then scalar processor.

Eg: Vector instruction For adding vectors.

$$C(1:100) = A(1:100) + B(1:100)$$

A is vector containing 100 data items

10 11 12 13 14 15 16

44 45 46 47 48

vector processor uses vector instructions.

Scalar " a " scalar "

The above exp. indicates vector "u".

→ same operation can also be performed using

scalar instruction with more than one

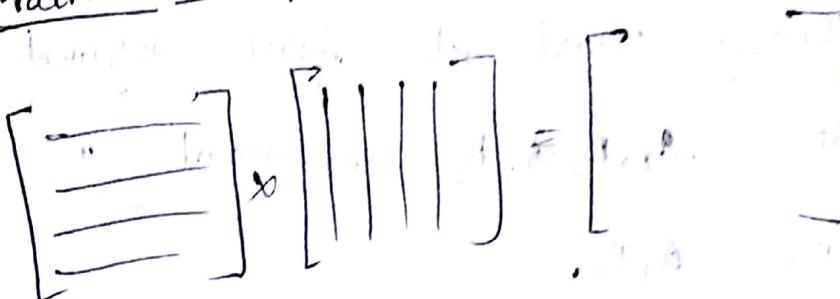
~~Instructions~~ using loops.

→ In scalar processor, single instruction will be operated on single data item at a time.

→ In vector processor, single instruction will be operated on multiple data items at a time due to arithmetic pipeline.

Eg:- Mat

Matrix Multiplication:



To prove:- To get inner product (element in product matrix),

To get inner product (element in product matrix),
matrices are multiplied.

$$C = A_1 B_1 + A_2 B_2 + A_3 B_3 + \dots + A_k B_k$$

A_1, A_2, \dots, A_k are row vectors

B_1, B_2, \dots, B_k are column vectors

C - inner product.

k - can be any value

→ thus inner product is obtained by

vector processor which uses pipelining.

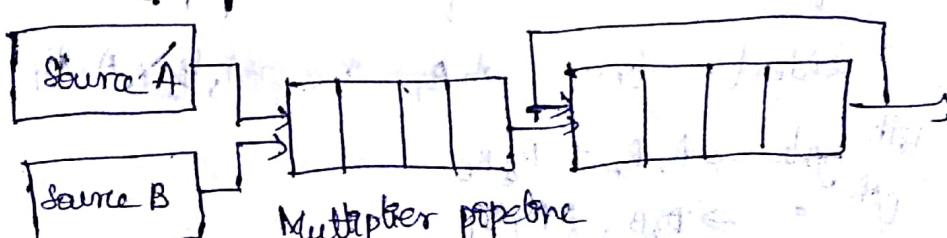


Fig. Pipeline for calculating an inner product

→ There are 4 segments in multiplier pipeline

and adder pipeline

→ During first cycle, first segment of multiplier pipeline contains A_7B_7 ,

→ Also during second cycle, first segment contains $A_2B_2 A_3B_3$ and second segment contains A_1B_1 .

→ During 4th cycle, the segments of multiplier pipeline contains $A_4B_4, A_3B_3, A_2B_2, A_1B_1$,

→ During 5th cycle, initially adder cell consists of 0. and 0 is added to A_1B_1 which is transferred to Adder pipeline

→ During 8 cycles, the

A_8B_8	A_7B_7	A_6B_6	A_5B_5
----------	----------	----------	----------

A_4B_4	A_3B_3	A_2B_2	A_1B_1
----------	----------	----------	----------

→ During 9th cycle,

first segment of adder pipeline is

added to A_1B_1 , i.e., $A_5B_5 + A_1B_1$

→ 10th cycle $\rightarrow A_2B_2 + A_8B_8$

11th $\rightarrow A_3B_3 + A_7B_7$

12th $\rightarrow A_4B_4 + A_8B_8$

→ Finally, partial products are obtained

$$C = A_1 B_1 + A_5 B_5 + A_9 B_9 + A_{13} B_{13} + \dots \quad \rightarrow ①$$

$$+ A_2 B_2 + A_6 B_6 + A_{10} B_{10} + A_{14} B_{14} + \dots \quad \rightarrow ②$$

$$+ A_3 B_3 + A_7 B_7 + A_{11} B_{11} + A_{15} B_{15} + \dots$$

$$+ A_4 B_4 + A_8 B_8 + A_{12} B_{12} + A_{16} B_{16}$$

Finally, 4 partial products in 4 segments
of adder pipeline are added.

Super scalar processor:

It contains multiple functional units or processor units such that multiple instructions are operated on multiple (separate) data items at the same time.

Super computer:

A commercial computer with vector instructions and pipelined floating point arithmetic operations is called a super computer.

→ In vector processor, data items are represented using floating point.

Array processors

- They are specially used to perform operations on arrays. (collection of data items of same type)
- There are 2 types of array processors.

1) Attached array processor

2) SIMD (Single instruction Multiple data) array processor

Attached Array Processor

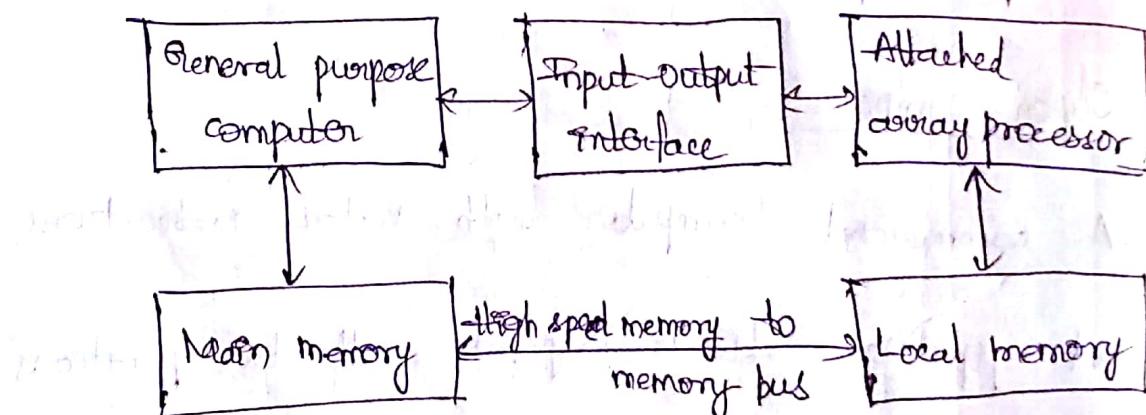


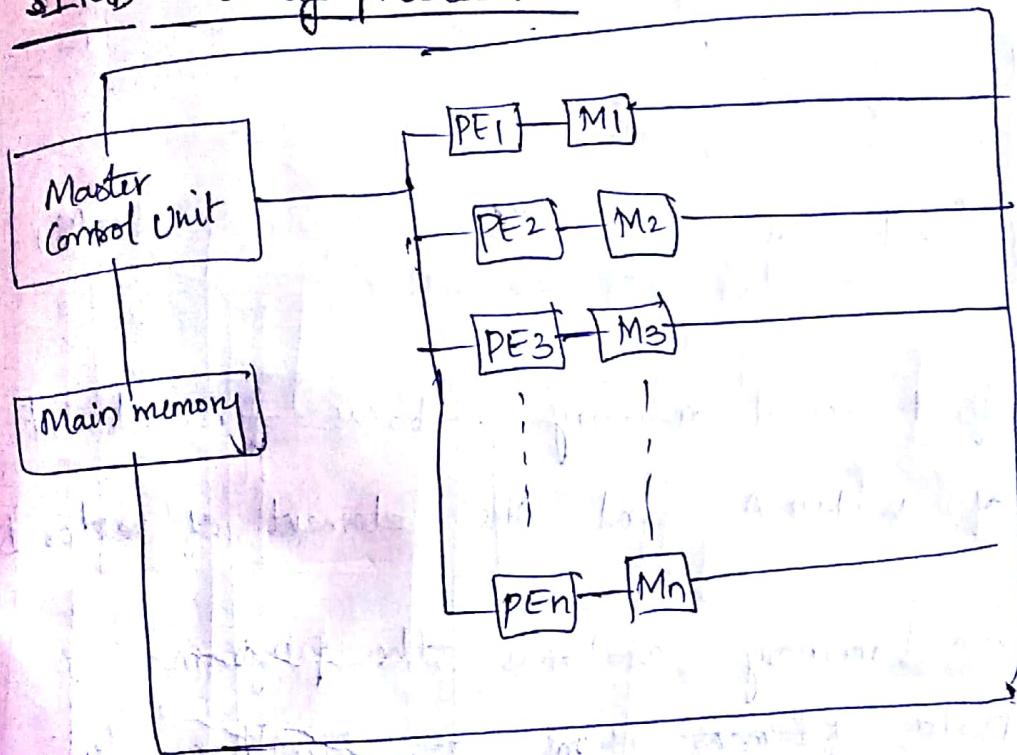
Fig Attached array processor with host

computer.

- It is attached to the general purpose computer. It acts as I/O device to the general purpose computer.

- It will access the data from local memory.
- The transfer of data from Main memory and local memory is done using high speed bus.
- It contains one or more pipeline floating pt. adders and multipliers.

SIMD Array Processor:



- In this processor, it contains multiple processing elements (PE) which contains ALU and registers.
(Floating Arithmetic unit)

→ M₁, M₂, ..., M_n indicates local memories

→ Master control unit decodes the instruction and

controls the operation of processing elements.

→ Scalar instructions and program control units, are executing using master control unit.

Vector instructions are executed by processing elements.

→ In case of vector instruction mentioned before

→ First element of vector A is given to M_1 ,

Second " " " " " " " M_2 ,

;

My first " " " " " " " M_1 ,

→ ∴ first local memory contains first element

of vector A and first element of vector B

→ Main memory contains the program

∴ Data elements items are simultaneously processed.

Multiprocessors:

- If a system contains multiprocessors then it is called multiprocessor system.
- In a microprocessor system different jobs are performed simultaneously (or) single job is partitioned into multiple tasks which are performed simultaneously.
- There are 2 types of multiprocessor systems
 - 1) Shared memory (or) tightly coupled multiprocessor system.
 - 2) Distributed memory (or) loosely coupled multiprocessor system.
- In this each processor uses its own local private memory.
- We have several structures for interconnecting the components of the multiprocessor system.

Interconnection structures.

1) Time shared common bus

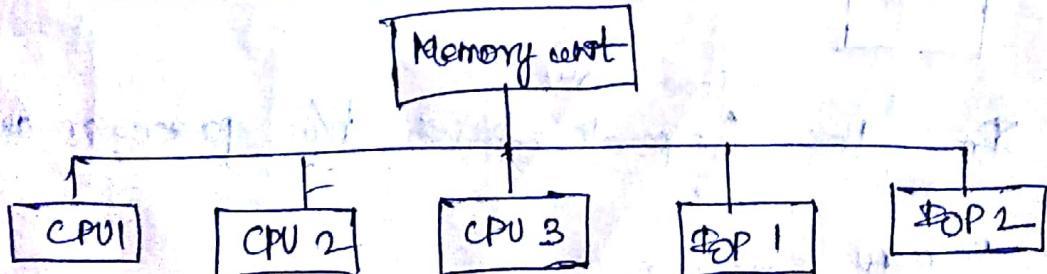
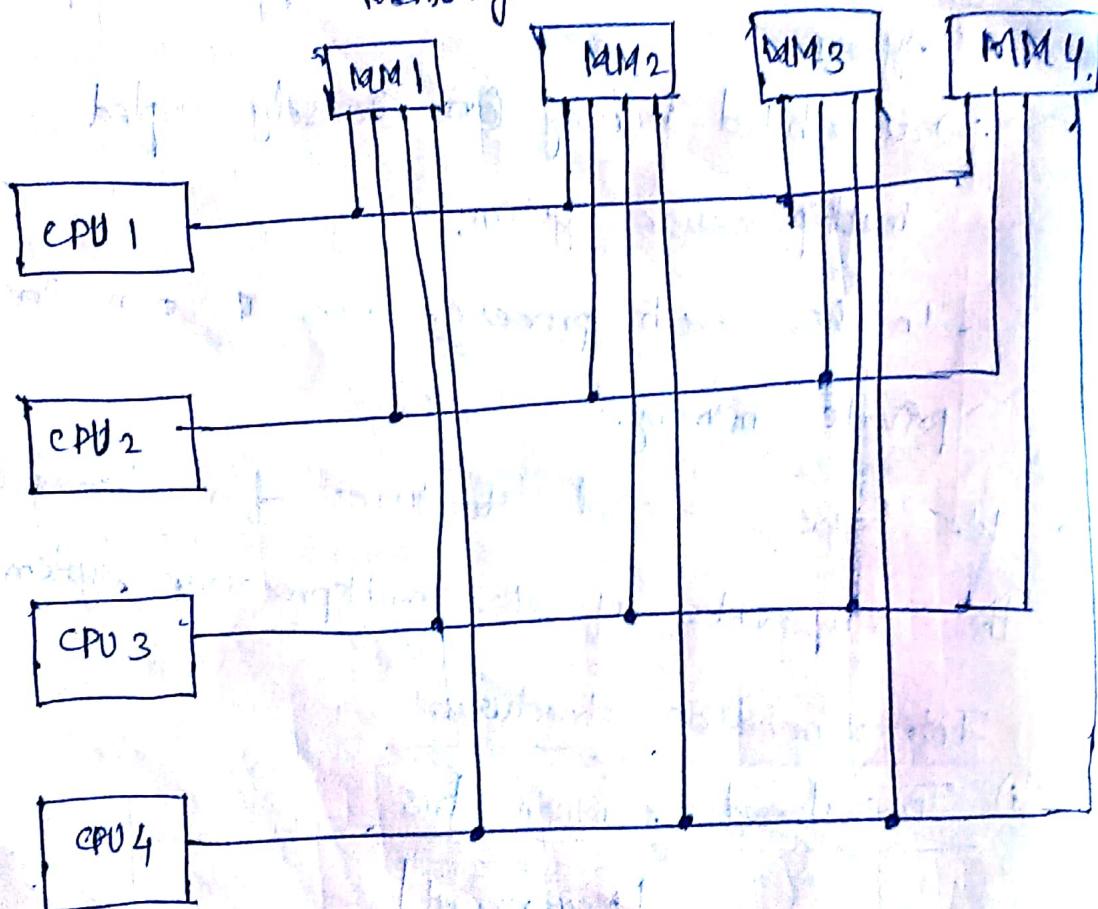


Fig: Time-shared common bus organization

- It indicates single bus (common) b/w the processors and memory
- When one processor uses memory, the other has to wait
- \Rightarrow If \rightarrow off processor which is used to perform i/p, o/p transfer

2) Multipoint memory

Memory modules.



\Rightarrow this, there are separate buses b/w processor and memory module.

Memory modules ^{are} parts of memory unit.

→ Memory module contains 4 ports. Each port is

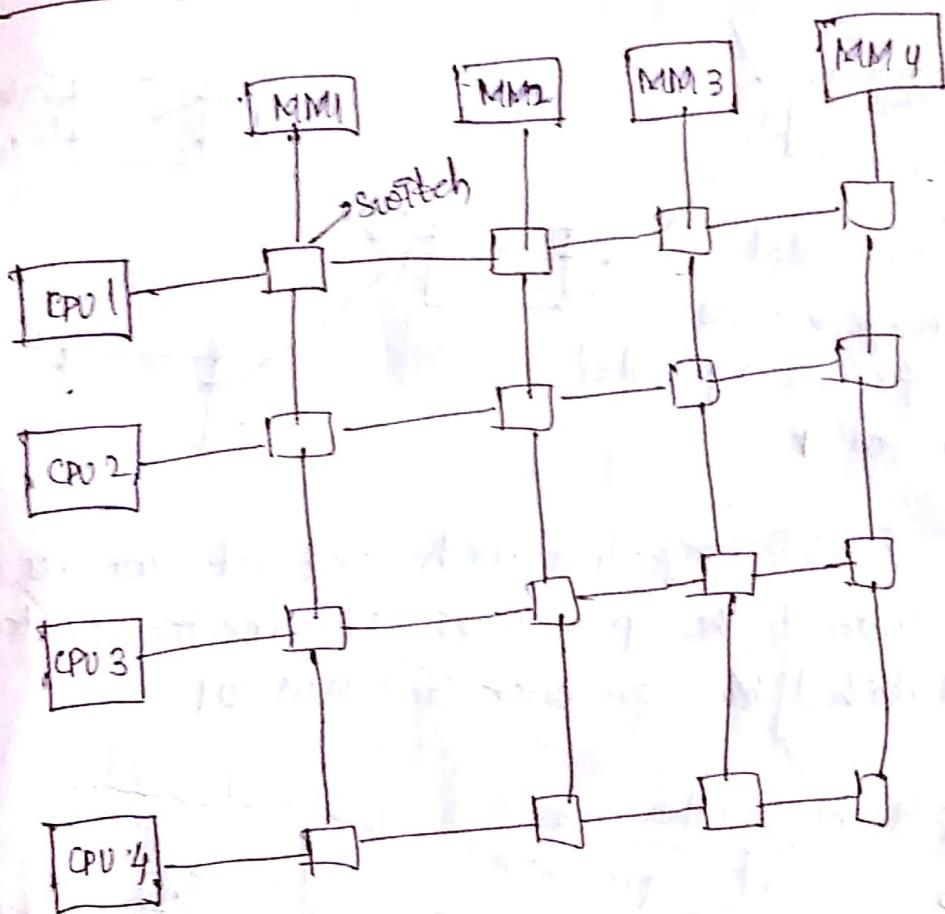
connected to 1 CPU. As it contains ~~multiple~~

multiple ports, it is called as multipoint memory.

It comes under shared memory (MM1, MM2...)

are ~~connected to~~ connected to CPU1, CPU2, ...

3) Crossbar switch:



Here, single port is used for memory module.

Here, switches are used to establish the

path b/w CPU and Memory module.

~~So~~ All the switches along the path should be ~~path~~ on for the transfer to take place.

④ Multistage Switching Networks:

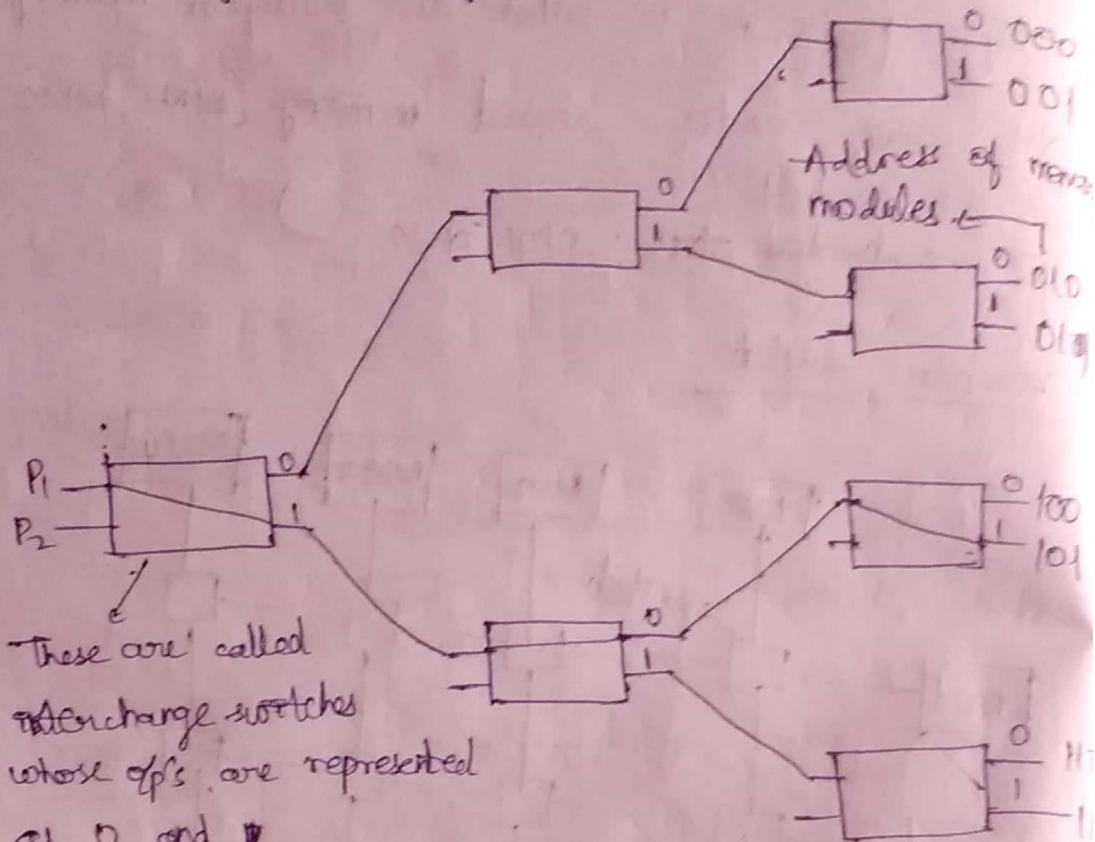
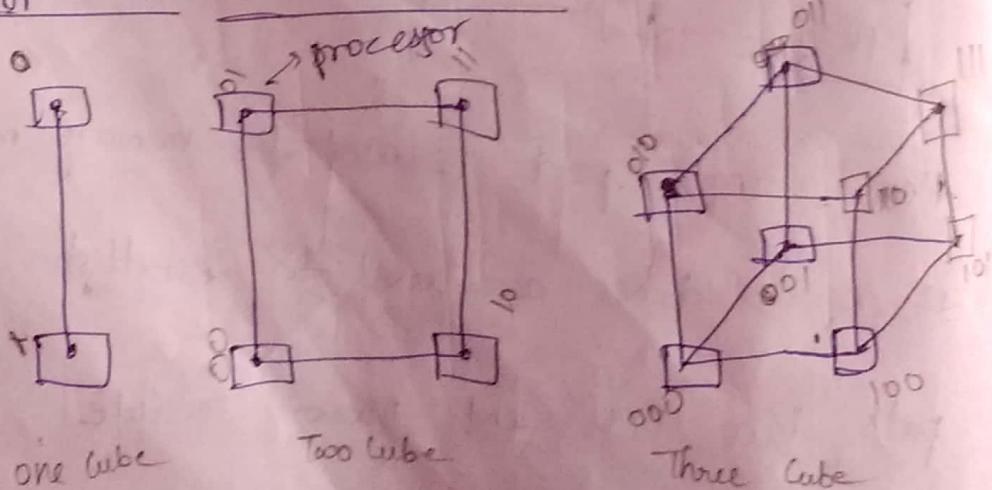


Fig: Binary tree with 2x2 st. switches.
P₁, P₂ indicates the processors. If connection is to be established b/w processor and M1 M1 101

⑤ Hypocube Interconnection



It indicates loosely coupled or distributed memory multiprocessor system

For 1 cube, 2^1 nodes.

{ For n cube, 2^n nodes

Each node indicates 1 processor

For n cube, we are interconnecting 2^n processors

00, 01, 10, 11 indicates address of nodes

Add. of successive nodes differ by only 1 bit



UNIT IV

INPUT-OUTPUT ORGANIZATION:

Input-output Interface:

- In order to compensate diff's b/w CPU and I/O devices, I/O interfaces are used.
- Eg:- CPU - high speed - electronic devices
I/O - slow - electro-mechanical or electromagnetic devices.
∴ signals values should be converted.

→ Another diff is data codes used by CPU and memory is diff from data code used by I/O devices. (e.g. ASCII or EBCDIC)

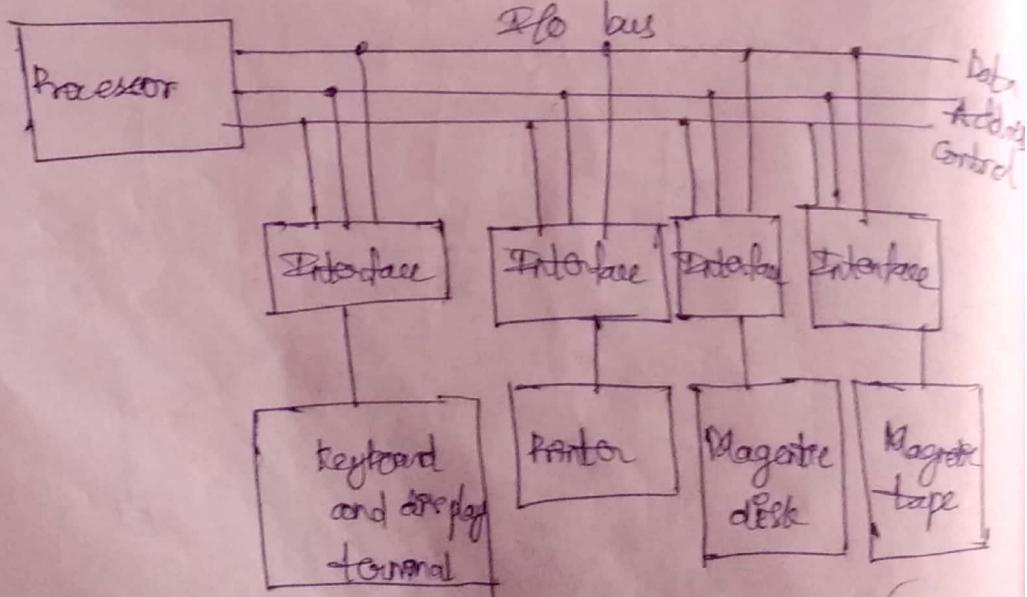


Fig Connection of I/O bus to Input-output device

Magnetic disk and Magnetic tape are used for storage devices

→ I/O devices are also called peripheral or external devices.

→ Bus contains data, address and control lines.

→ The address sent by the processor on the address lines ~~are~~ is used to select the I/O device.

→ Info transferred on the control lines is first function code which is referred as I/O command.

They are of 4 types.

I/O commands

1) Control command: It will inform the peripheral what to do.

2) Status command: It is used for diff. purposes eg. to know the status of peripheral re. whether it is ready to perform data transfer or not.

3) Whether there are any errors during the data transfer.

3) Data output command: It is used to transfer data on the data lines into the register of the interface.

4) Data Input command: It is used to transfer the data in the interface register into the data lines.

→ The bus which is used by processor and I/O devices is called ~~I/O~~ bus, which contains data, address and control lines. If by processor and memory → memory bus

→ There are 3 ways of using buses.

1) Using separate bus for I/O and memory

2) Using one common bus " " " " with separate control lines which are called as isolated I/O.

3) Using same common bus for I/O and memory with common control lines with common control which is called as memory mapped I/O.

Mem

Modes of Transfer

Either data is transferred from input device to memory or from memory to O/P device. This can be done through with or without CPU.

→ To transfer data from I/O device to memory, it is transferred from I/O to processor and from processor to memory.

If to transfer data from memory to I/O device, it is transferred from ^{memory} I/O device to processor and processor to ~~I/O~~ I/O device. It is done ~~with~~ through CPU.

There are one way to transfer

→ Second way is directly transferring data from memory and I/O device which is called as direct memory access.

3 Modes of transfer

1) Programmed I/O } through CPU

2) Interrupt Initiated I/O

3) Direct Memory access (DMA) → without CPU

In programmed I/O, as CPU is fast and I/O is slow, CPU has to continuously check whether the device is ready to perform data transfer. By this, CPU time is wasted.

which means that we are not efficiently using CPU.

2) To overcome this, interrupt initiated I/O is used.

In which we are using interrupt facility where the interface will be interrupting the CPU whenever the device is ready to perform data transfer. Meanwhile, CPU can be used to execute another program. In this CPU need not continuously check whether the device is ready or not.

2) DMA — We are directly transferring memory between memory and I/O.

12/10/17

Priority Interrupt:

To perform the data transfer.

1) Programmed

→ Programmed

→ Interrupted → Interrupt Initiate

→ Interrupt initiated I/O

→ In programmed I/O, CPU has to check whether I/O device is ready to transfer data or not. This is the drawback.

This can be overcome.

The time delay can be reduced.

because, CPU will get initiative from I/O device when I/O is ready to transfer.

The time in which CPU is waiting

- For I/O transfer we can perform another operations.
- For processor(CPU), we may have 2 or more I/O devices like keyboard, mouse etc. all these are interrupted from I/O.
- If all interrupts are occurring in the processor, processor responds to only one interrupt based on priority.
 - High speed I/O device - high priority
 - Low speed " " - lowest "
- Processor does servicing (or) responses to high priority when 2 or more devices are given can be done using 1) hardware
2) software.
- The task of servicing higher priority device by software is called polling.

Software \rightarrow by program (executing)

- In program, first we will be checking the 1st highest priority device is being ~~either~~ interrupted or not
- If interrupted, it will be serviced first
- If no first higher priority device, then 2nd

higher priority is selected and checked then

continued.

→ The task of sorting can also be done by 'hardware' i.e. 'Daisy chaining priority'.

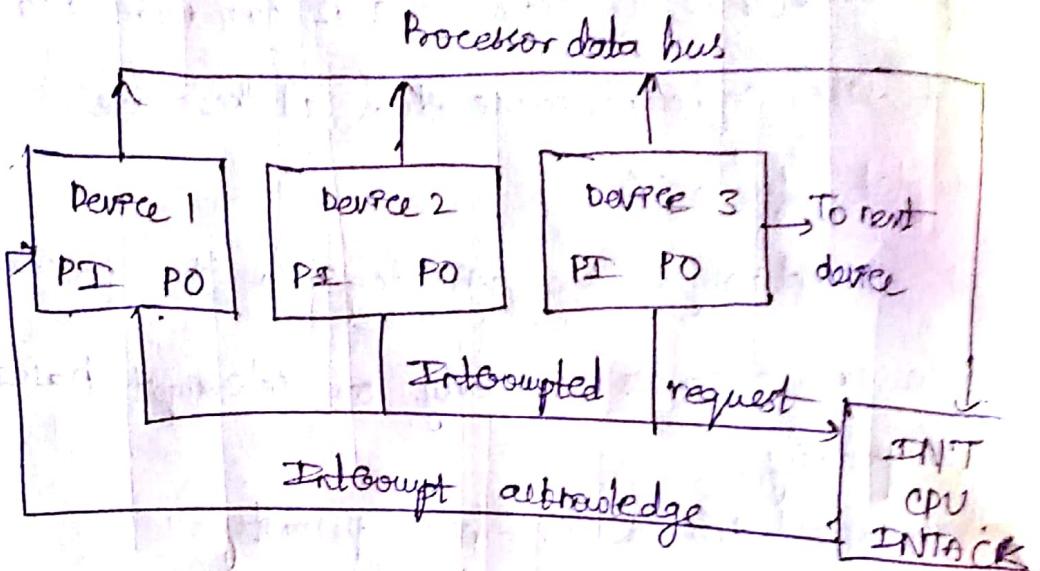


fig Daisy chaining priority Interrupt

→ In these, we are connecting devices serially, which 1st device is highest priority device & last device is lowest priority device.

INT → Interrupt Active low signal.

→ One (or more) no. of interrupt may be received to CPU.

→ In response of receiving interrupt, CPU generates INTACK signal which is active high signal.

INTACK - Interrupt acknowledge signal

→ If 1st interrupt acknowledgement signal is given by device with 1st highest priority device then if interrupt is present by device 1, IR will be '0'.

(Op of 1)

PO will be zero means, the acknowledgement signal is not passed to 2nd device.

→ If in case 2nd device is having to interrupt then PO(Op of 1st device) will be '1'.

That means, interrupt acknowledgement signal is passed to 2nd device and so no interrupt until the acknowledgement signal reaches the correct device.

→ The device which is having IR=1 & PO=0 is the device to be required by CPU.

→ During reading bit transfer data from device to CPU, that transfer is done by vector address (VAD).

→ Another hardware method is "parallel priority interrupt".

→ Here devices are connected parallelly whereas in previous hardware method, they are connected serially.

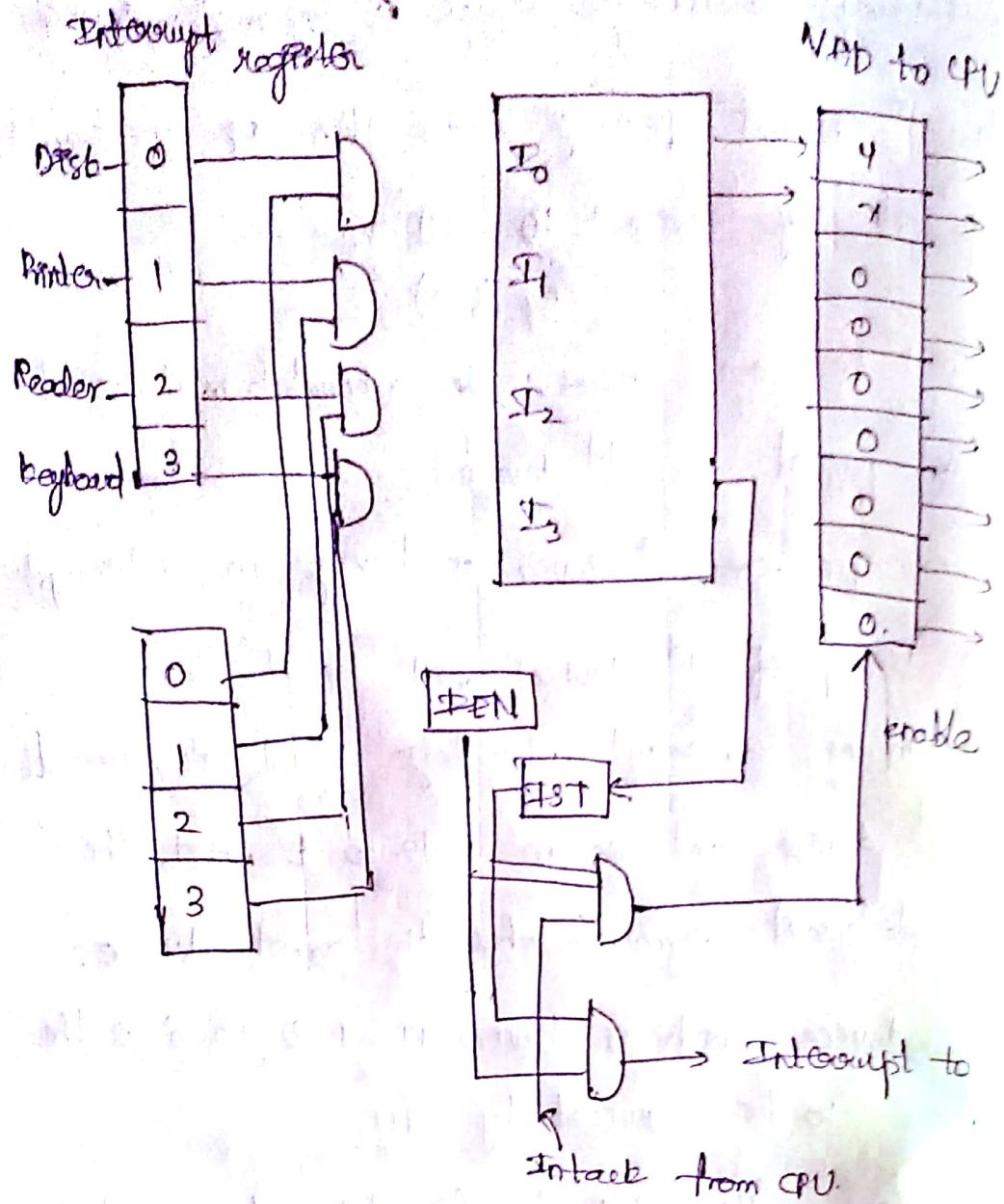


Fig Priority interrupt hardware

→ In interrupt register the devices are placed from top to bottom with high priority ~~1~~
first and lowest priority device at bottom.

→ The purpose of mask register is not to allow the lowest priority interrupt to go to CPU when CPU is dealing with highest

priority interrupt.

- Another purpose of Mask register is to allow the highest priority interrupt to CPU when lowest priority interrupt is being serviced.
- When there is an interrupt from interrupt register & it is not serviced, then op of AND gate is '1'
- Not masked means
 - If CPU is dealing with disk, then printer at same time getting interrupt, the Mask register should mask the printer device.
- If any of the op of AND gate is '1', IST flop is activated.
- The purpose of DEN - Interrupt enable signal
 - DEN = 1, the interrupt will be service &
 - DEN = 0, No " " taken
- After getting INTACK from CPU, we are enabling CPU containing VAD (Vector Address)
- VAD to CPU (register)

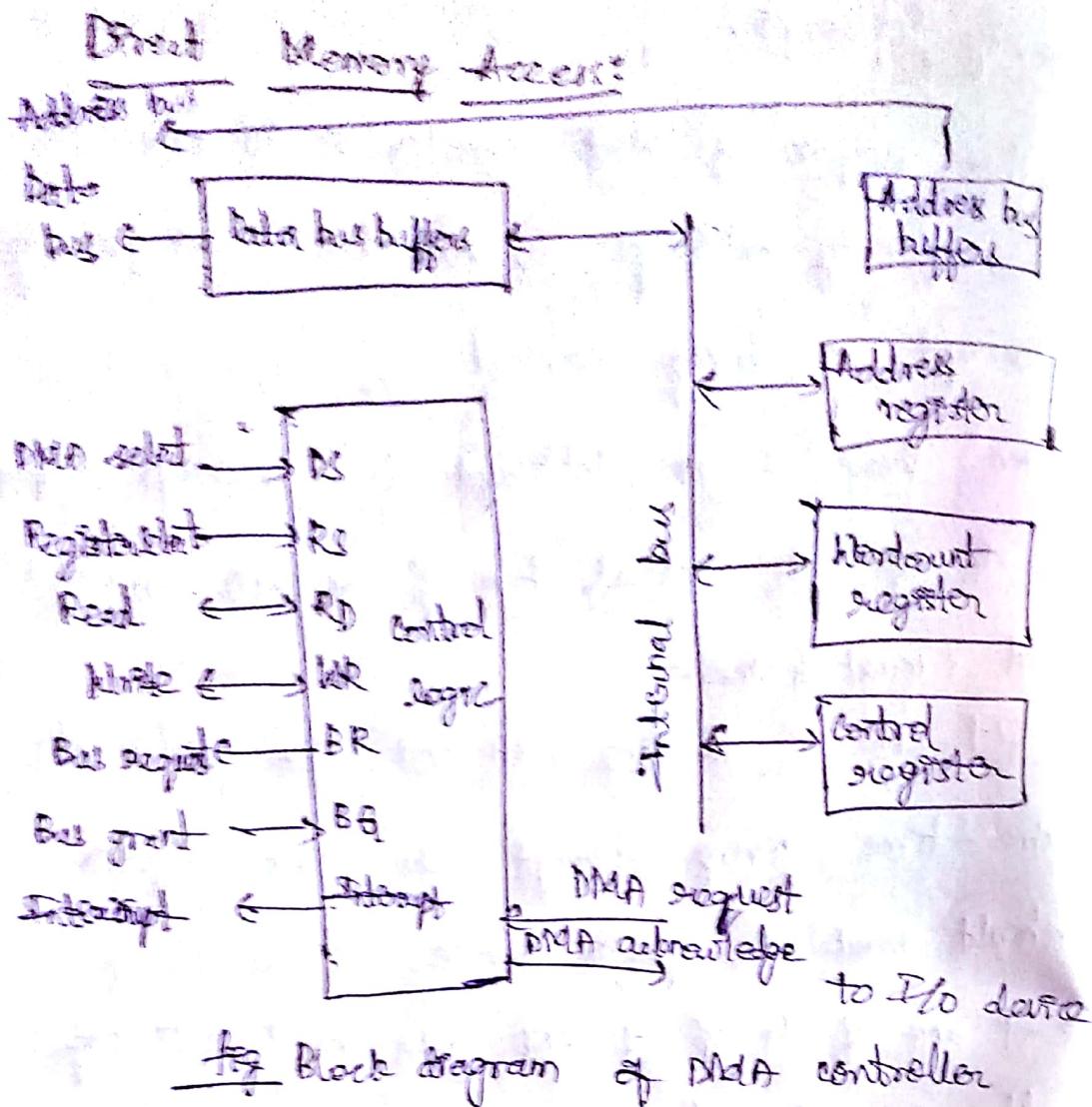
from device \rightarrow we have 00

" " 2 \rightarrow " " 01

" " 3 \rightarrow " " 10

" " 4 \rightarrow " " 11

IST - Interrupt status flop.



Block diagram of DMA controller

- It is one of modes of memory transfer b/w memory and I/O.
- In DMA, without the intervention of CPU, directly data is transferred b/w I/O and memory. For DMA directly direct access, DMA controller is used.
- There are 3 registers, which are initially
 - (before performing DMA transfer) initialised by CPU.

1) Address register

- To read (from memory \rightarrow I/O) and write (I/O to memory)
- It indicates the starting location of memory from which we have to read or write

2) Word count registers

- It indicates no. of words that have to be read/stored into memory during read/write operation.

3) Control register

- It indicates specifies the mode of transfer i.e. whether read or write. and also to start the DMA transfer

- When when
- When there is a request from I/O device to perform DMA transfer, (Initially memory buses are under control of CPU).

DMA controller should get the control of memory buses by sending bus request signal to requesting the the control of signal to the CPU.

- When bus request signal is sent, CPU terminates the execution of instruction (clock goes off that instant) and keeps the memory buses in high impedance which means that the memory buses are disconnected from CPU.
- Then, bus grant signal is sent by CPU (which means memory buses are under control of DMA controller)
- During ~~initial~~ initialization of the 3 DMA, data select pin and register select are enabled
- Read and write pins are bidirectional. They act as RP and CP. These pins act as RP pins to CPU and CP pins to memory. In order to load data into
- ~~When data is to be loaded into~~ 3 registers by CPU, Read/write pins are enabled.

- After each word transfer, word count becomes 300
 - When all words are transferred from wordcount register to I/O device & memory, interrupt is enabled for and wordcount will be 300.
 - In order to read the wordcount register, Read is enabled i.e. CPU sends the read signal.
 - DMA acknowledge will be generated to I/O device by cont DMA controller
- When I/O device receives DMA acknowledge, if it is an I/p device, it will place the data on the data bus. If, if it is a O/p device, it will take the data from data bus and word count register will be decremented. If the wordcount value ≠ 0, I/O device will send a DMA request and DMA controller check whether there is a DMA request. → If the I/O device is fast it will send a DMA request as soon as previous word is transferred and viceversa.

→ In case of slow I/O devices, DMA controller disables BR signal, and bus grant signal is disabled by CPU, i.e. memory buses are under control of CPU, and the CPU executes the program where it has stopped.

This is called as cycle ~~steering~~ ^{stealing}.

Handing over the memory buses to CPU after each coord transfer in case of slow I/O device is called cycle ~~steering~~ ^{stealing}.

It does not occur due to fast I/O devices. It will ~~do~~ handing over the control of memory buses ~~only~~ only after transferring all the coords. This is called burst transfer.

→ The 3 registers are ~~initialised~~ initialised by CPU through data bus. My word count register ^{value} is read by CPU through data bus.

After completion of all work, Interrupt is sent to CPU by M/H controller. CPU reads value of word count register \Rightarrow all words are transferred.

Input-output Processor

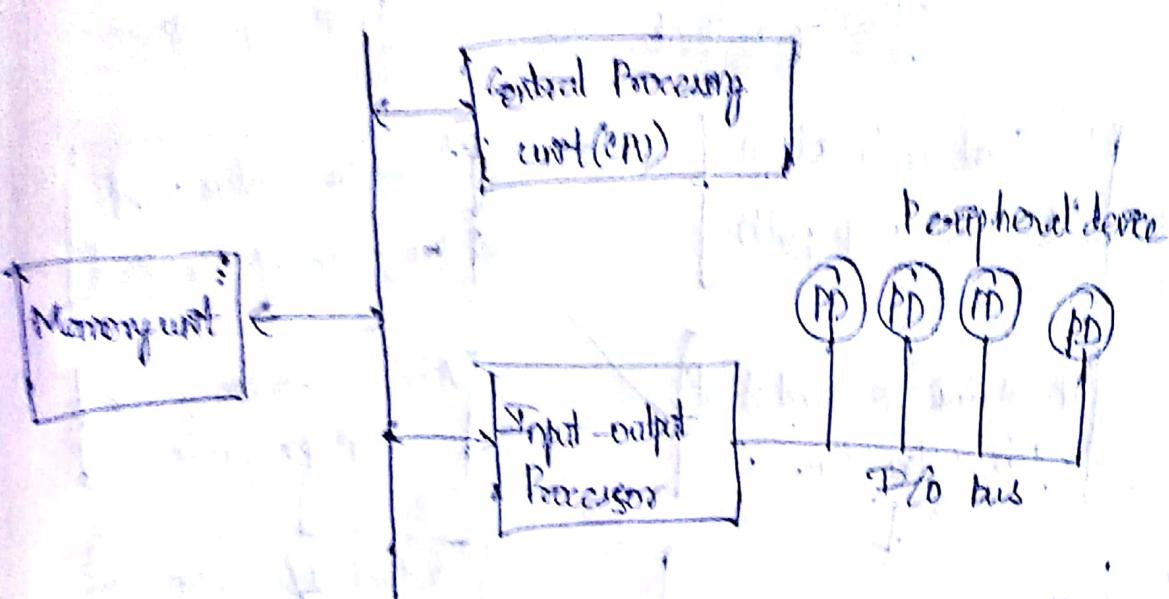


Fig.: Block diagram of a computer with I/O processor

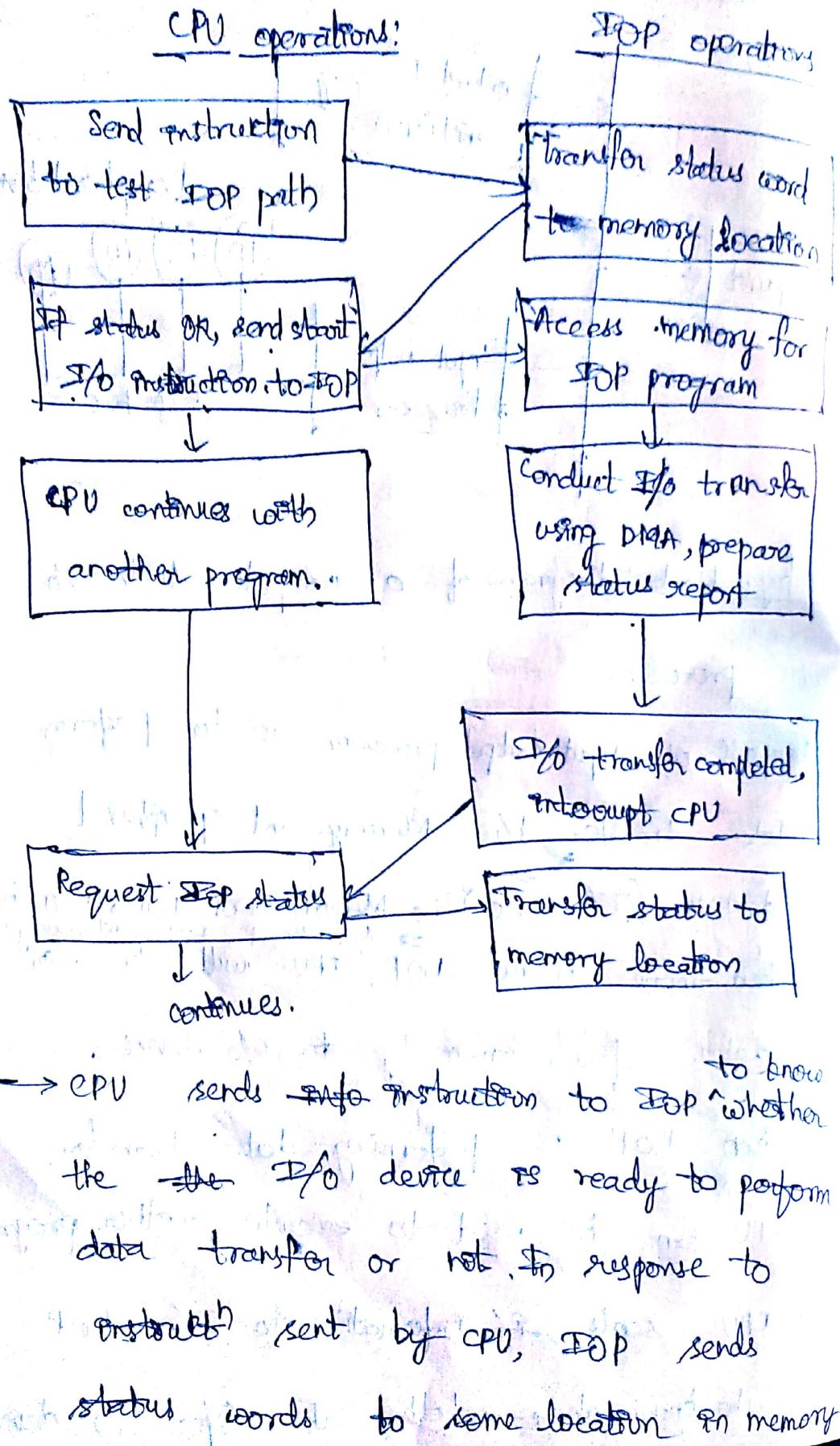
→ Purpose of Input-output processor is to perform data transfer b/w memory and peripheral devices (I/O devices) i.e. Memory \rightarrow DOP, DOP \rightarrow I/O (D)
or I/O \rightarrow DOP, DOP \rightarrow Memory unit
→ When there is no IOP, CPU will transfer data from memory to I/O device.

→ When DOP is performing data transfer, CPU can be used to execute another program

→ CPU sends ~~if~~ information to DOP. DOP to know where whether the ~~if~~ I/O device

PS: ready or not,

CPU-DOP communication



which indicates whether the I/O device is ready to perform data transfer.

- CPU reads the status words in memory and knows whether the device is ready to perform data transfer.
 - If status word is Ok → I/O device is ready
 - Then CPU sends start I/O instruction to SOP
 - When SOP receives start I/O instruction, SOP accesses SOP program from memory.
 - While SOP performs SOP program, CPU can be used to perform execute another program.
 - Status report indicates whether there is any error occurred during data transfer
 - After the completion of data transfer, SOP interrupts CPU. The purpose of interrupting CPU is to inform CPU that data transfer is completed. Then CPU requests the SOP status report from SOP and that status report is sent to memory location. From memory location, CPU reads the status word (report).

→ After reading status words, CPU knows whether any error has occurred or not.

Interrupt initiated I/O

→ In this I/O, I/O interrupts the CPU when the I/O device is ready to perform data transfer.

→ When CPU is interrupted, it executes interrupt service routine (ISR) which is nothing but I/O program.

→ To know the address of ISR, there are 2 ways.

1) Vectored interrupt

2) Nonvectored interrupt

→ In this case, the branch address (the address to which we have to jump in order to execute ISR) is contained in some fixed location in the memory.

∴ CPU has to take the branch address from memory.

→ In vectored interrupt, the branch address is provided by interrupt source. (One which is producing interrupt i.e. I/O or D/A device).

This branch interrupt information is called as interrupt vector.

∴ This interrupt vector indicates the starting address of the interrupt service routine i.e., D/A program which is used to perform I/O or D/A operation. (or) we get the address of address starting address of ISR in that address (direct or indirect.)