

1/2/2020  
Saturday.

## unit - 2

### Arithmetic

#### 1. Introduction:-

- A basic operation in all computers is the addition or subtraction of two numbers.
- Arithmetic operations occur at the machine instruction level. They are implemented along with basic logic functions such as AND, OR, NOT, and EXCLUSIVE-OR (XOR), in the ALU subsystem of the processor.
- But in this chapter, Logic circuits used to implement arithmetic operations.
- Multiply & divide operations require more complex circuitry than either addition & subtraction operations, and also affect performance.
- Compared with arithmetic operations, logic operations are simple to implement combinational logic. They require only independent Boolean operations whereas carry/borrow lateral signals are required in arithmetic operations.
- To perform addition and subtraction operations 2's complement is the best representation.
- A logic circuit designed to add unsigned binary numbers can also be used to add signed numbers in 2's complement.

## Difference between Signed and Unsigned Numbers

| <u>unsigned</u>                            | <u>Signed</u>  |
|--|--|
| → valid only for positive nos              | → valid for both +ve & -ve nos   |
| → NO sign bit concept.                     | → sign bit concept is used.  |
| → It is the straight binary Representation | → MSB of the entire number is known as sign bit(s).  |
| → NO further categorization                | $S=0 \Rightarrow$ Positive Number<br>$S=1 \Rightarrow$ Negative No.  |
|  | → Here there are 3 types of categorization.<br>1. Signed Magnitude Rep.<br>2. 1's Complement<br>3. 2's Complement. |

## TOPIC 1: Addition and Subtraction of Signed Number

| <u><math>x_i</math></u> | <u><math>y_i</math></u> | <u>carry-in (<math>c_{i-1}</math>)</u> | <u>Sum (<math>s_i</math>)</u> | <u>Carry-out (<math>c_i</math>)</u> |
|-------------------------|-------------------------|--|-------------------------------|-------------------------------------|
| 0                       | 0                       | 0                                      | 0                             | 0                                   |
| 0                       | 0                       | 1                                      | 1                             | 0                                   |
| 0                       | 1                       | 0                                      | 1                             | 0                                   |
| 0                       | 1                       | 1                                      | 0                             | 1                                   |
| 1                       | 0                       | 0                                      | 1                             | 0                                   |
| 1                       | 0                       | 1                                      | 0                             | 1                                   |
| 0                       | 1                       | 1                                      | 1                             | 1                                   |
|                         |                         |  |                               |                                     |

$$s_i = \bar{x}_i \bar{y}_i c_i + \bar{x}_i y_i \bar{c}_i + x_i \bar{y}_i \bar{c}_i + x_i y_i c_i$$

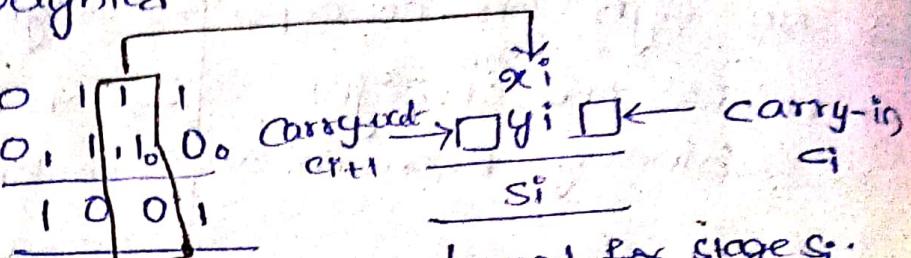
$$= x_i \oplus y_i \oplus c_i$$

$$c_{i+1} = y_i c_i + x_i c_i + x_i y_i$$

→ The above figure shows the logic truth table for the sum ( $s_i$ ) and carry out ( $c_{i+1}$ ) functions for adding equally weighted bits  $x_i$  &  $y_i$  in two nos.

Ex:

$$\begin{array}{r} x \\ + y \\ \hline z \end{array} = \begin{array}{r} 7 \\ + 6 \\ \hline 13 \end{array}$$



Legend for stage  $s_i$ :

fig: Logic specification for a stage of binary addition

→ The logic expression for  $s_i$  can be implemented with a 3-input XOR gate.

sum



fig: Logic for a single stage of binary addition.

→ The carry-out function,  $c_{i+1}$ , is implemented with a two-level AND-OR logic circuit.

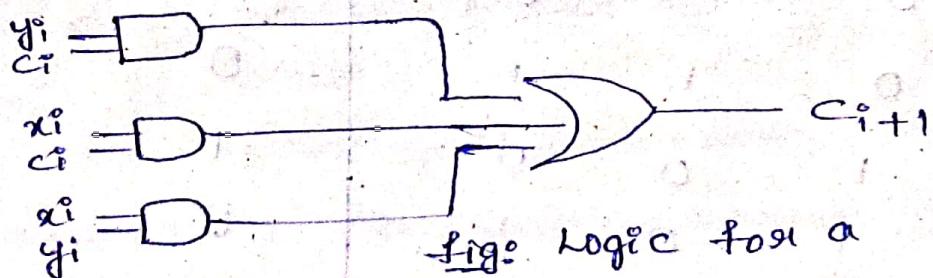
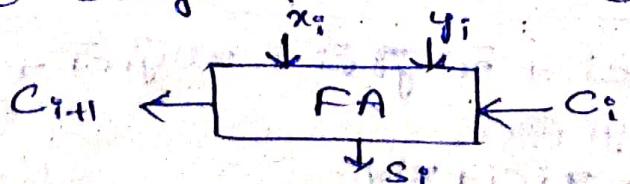
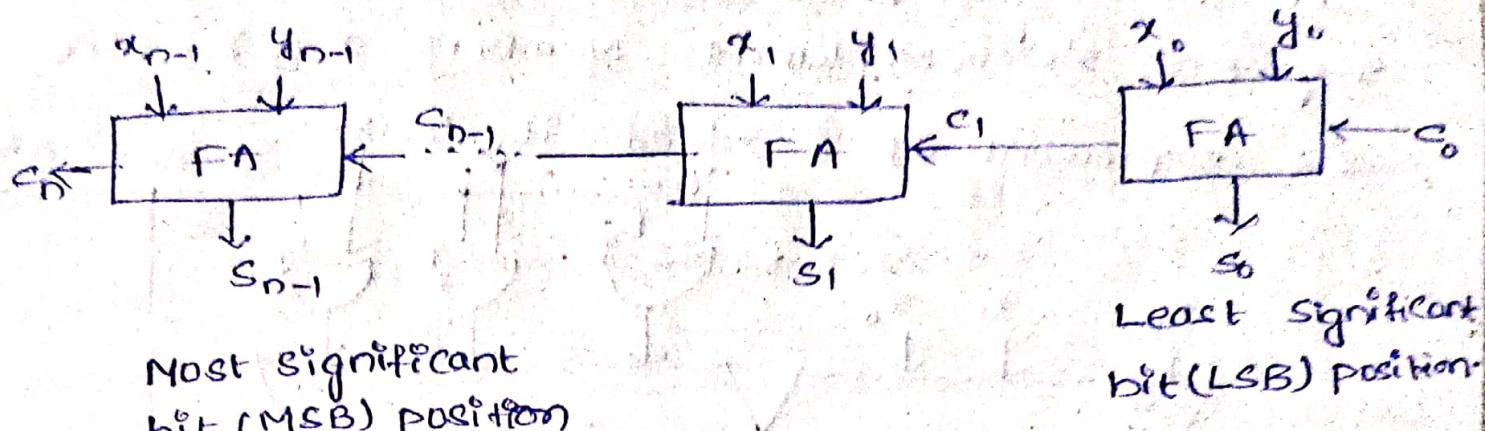


fig: Logic for a single stage

→ A convenient symbol for the complete circuit for a single stage of addition called a full adder (FA).



→ The cascaded connection of n-bit adder blocks.



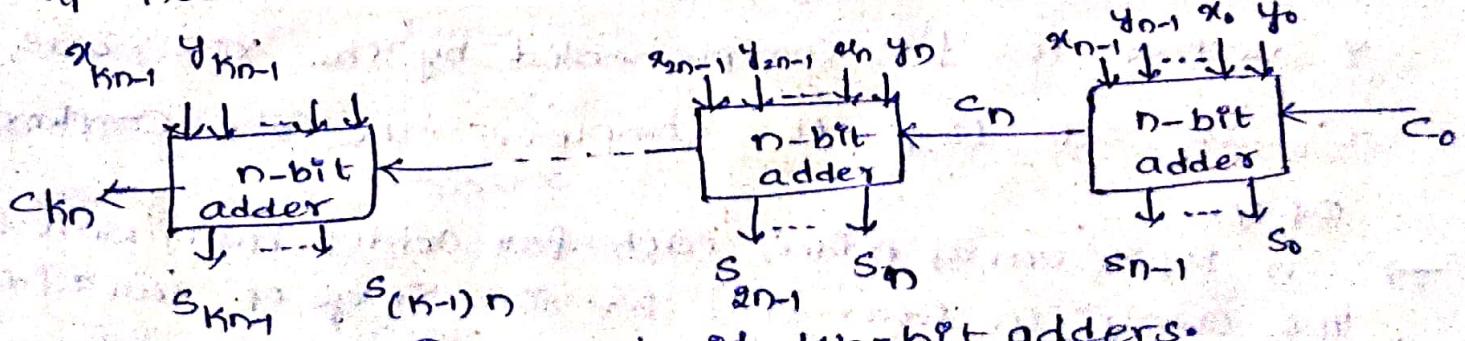
(b) fig: An n-bit ripple-carry adder.

→ It can be used to add two n-bit numbers since the carries must propagate (or) ripple through this cascade, the configuration is called an n-bit ripple-carry adder.

→ The carry-in ( $c_o$ ) into the least-significant-bit (LSB) position provides a convenient means of adding 1 to a number.

→ for instance, forming the 2's complement of a no involves adding 1 to the 1's complement of the number.

→ The carry-signals are also useful for interconnecting n-adders to form an adder capable of handling 16 bits that are  $16n$  bits long, as shown in fig.



(c) cascade of  $k$  n-bit adders.

All (a) to (c) fig name for Logic for Addition of binary vectors.

## Addition | subtraction Logic unit :-

→ In order to perform the subtraction operation  $x-y$  on 2's complement numbers  $x$  &  $y$ , the 2's complement of  $y$  added to  $x$ .

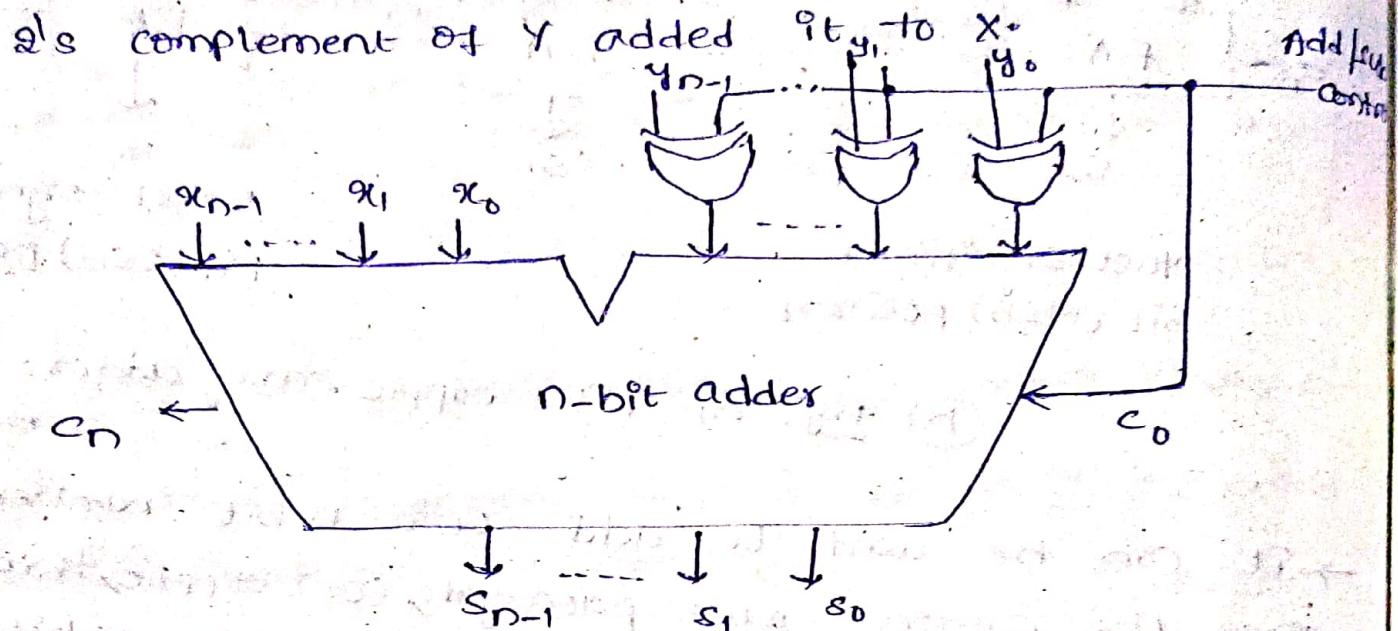


fig: Binary addition-subtraction logic network.

→ This add|sub logic unit can be used to perform either addition (or) subtraction based on the value given to the Add|sub control line.

→ This line, i.e., set to 0 for addition, applying the  $y$  vector unchanged to one of the adder inputs along with a carry-in signal,  $c_0$  of 0.

→ When the Add|sub control line is set to 1, the  $y$  vector is 2's complemented by the XOR gates &  $c_0$  is set to 1 to complete the 2's complementation of  $y$ .

→ In Ripple carry Adder, each full adder has to wait for its carry in from its previous stage of full adder.

→ This causes a delay & makes ripple carry adder extremely slow to overcome this carry look-ahead adder.

## TOPIC 2: Design of Fast Adders

- If an n-bit ripple-carry adder is used in the addition/subtraction unit, it may have too much delay in developing its outputs.
- In order to reduce the delay the approach is to use the fastest possible electronic technology in implementing the Ripple-carry-logic design (or) carry-looks ahead addition.

Carry-Look-Ahead Addition :-  
 → Carry-Look-Ahead adder is an improved version of the Ripple-carry adder.  
 → The logic expressions for full adder sum ( $s_i$ ) & carry ( $c_{i+1}$ )

Carry-out's given as,

$$s_i = x_i \oplus y_i \oplus c_{i-1} \quad \rightarrow (1)$$

$$c_{i+1} = x_i y_i + x_i c_{i-1} + y_i c_{i-1} \quad \rightarrow (2)$$

factoring the eqn (2) we get,

$$c_{i+1} = x_i y_i + @ (x_i + y_i) c_{i-1}$$

we can write

$$c_{i+1} = G_i + P_i c_i$$

$$G_i = x_i y_i \quad P_i = (x_i + y_i)$$

→ The  $G_i$  and  $P_i$  are called the Generate & Propagate functions.

→ The Propagate function means that an input carry will produce an output carry when either  $x_i$  is 1 (or)  $y_i$  is 1.

→ Each bit stage contains an AND gate to form  $G_i$  an OR gate to form  $P_i$  & a three input XOR gate to form  $s_i$ .

→ The adequate propagate functions can be visualized as,  $P_i = x_i \oplus y_i$  which differs from  $P_i = x_i + y_i$  only when  $x_i = y_i = 1$ .

$$\begin{aligned} \rightarrow C_{i+1} &= x_i y_i P + x_i^c C_i + y_i^c C_i \\ &= x_i^c y_i^c + (x_i + y_i) C_i \\ C_{i+1} &= G_i + P_i C_i \end{aligned}$$

when  $i=0$ ,

$$C_{0+1} = G_0 + P_0 C_0 \Rightarrow C_1 = G_0 + P_0 C_0$$

$$\begin{aligned} i=1, C_{1+1} &= G_1 + P_1 C_1 \Rightarrow C_2 = G_1 + P_1 C_1 \\ &= G_1 + P_1 (G_0 + P_0 C_0) \\ C_2 &= G_1 + P_1 G_0 + P_1 P_0 C_0 \end{aligned}$$

$$i=2, C_{2+1} = G_2 + P_2 C_2 \\ C_3 = G_2 + P_2 [G_1 + P_1 G_0 + P_1 P_0 C_0]$$

$$C_3 = G_2 + P_2 G_1 + P_2 P_1 G_0 + P_2 P_1 P_0 C_0$$

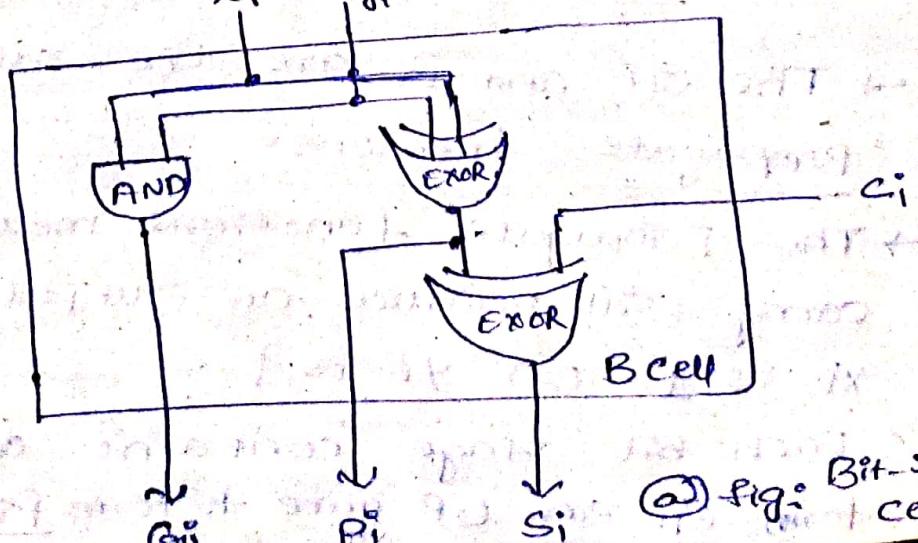
$$\begin{aligned} i=3, C_{3+1} &= G_3 + P_3 C_3 \\ C_4 &= G_3 + P_3 [G_2 + P_2 G_1 + P_2 P_1 G_0 + P_2 P_1 P_0 C_0] \\ C_4 &= G_3 + P_3 [G_2 + P_2 G_1 + P_3 P_2 G_0 + P_3 P_2 P_1 G_0 + P_3 P_2 P_1 P_0 C_0] \end{aligned}$$

$$\text{Here } C_{i+1} = G_i + P_i C_i$$

$$G_i^P = x_i y_i$$

$$P_i^c = x_i^c y_i^c$$

$$S_i = x_i^c y_i^c \oplus C_i$$



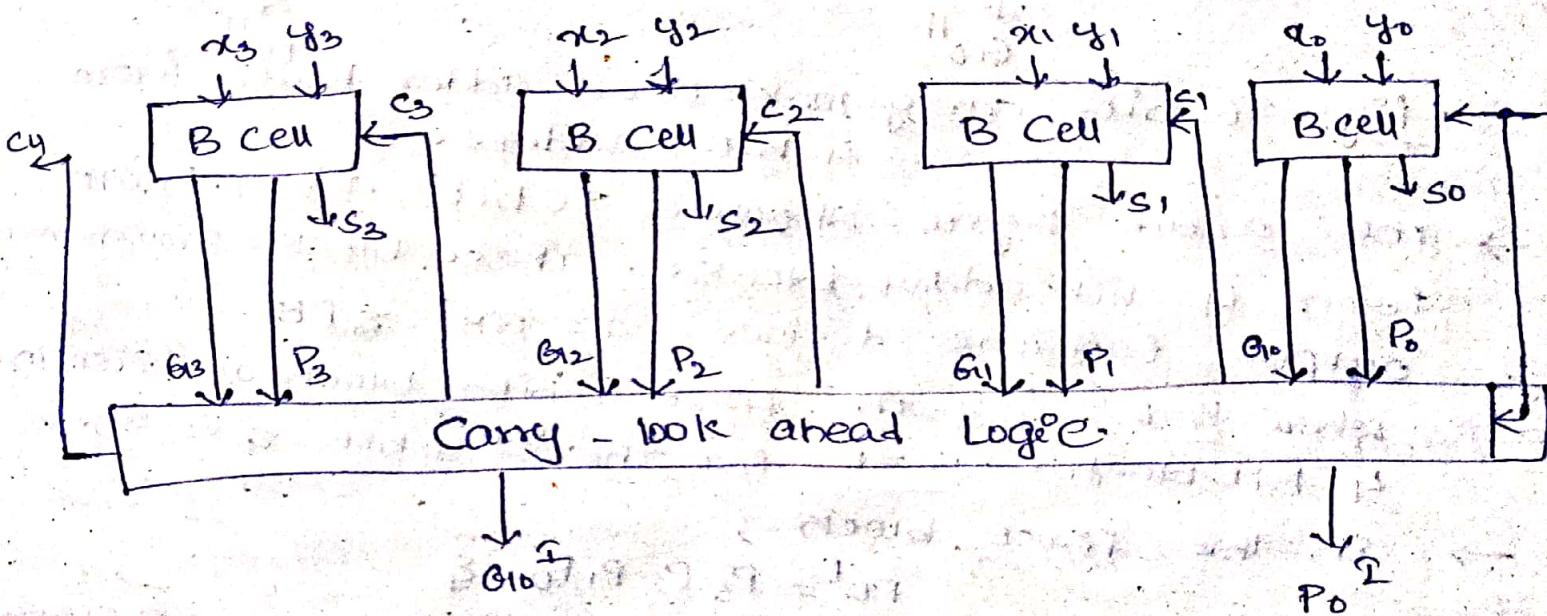
### → Advantages :-

1. Circularity Reduce that's why it is better.
2. Only two Pin Gates instead of 3 pin gate. only  $(x_1, y_1)$
3. Two inputs are allowed.
4. It reduces the propagation delay.

### Drawbacks :-

4-Gate Delay. [One AND Gate Delay & 1 or 2 Gate Delay &  
( $c_i, g_i$  & EXOR)]

→ It gets more complicated as the no. of bits increase.



(b) 4-bit Adder

fig: 4-bit carry-lookahead Adder.

→ The complete 4-bit adder, here the carries are implemented in the carry look ahead logic. An adder implemented in this form is called a carry-lookahead adder.

### Drawbacks:-

→ 3 gate Delay for carry bit & 4 gate delay for all sum bits.

Note: 4-bit Ripple-Carry Adder requires 7 gate delays for  $S_3$  & 8 gate delays for  $C_4$ .

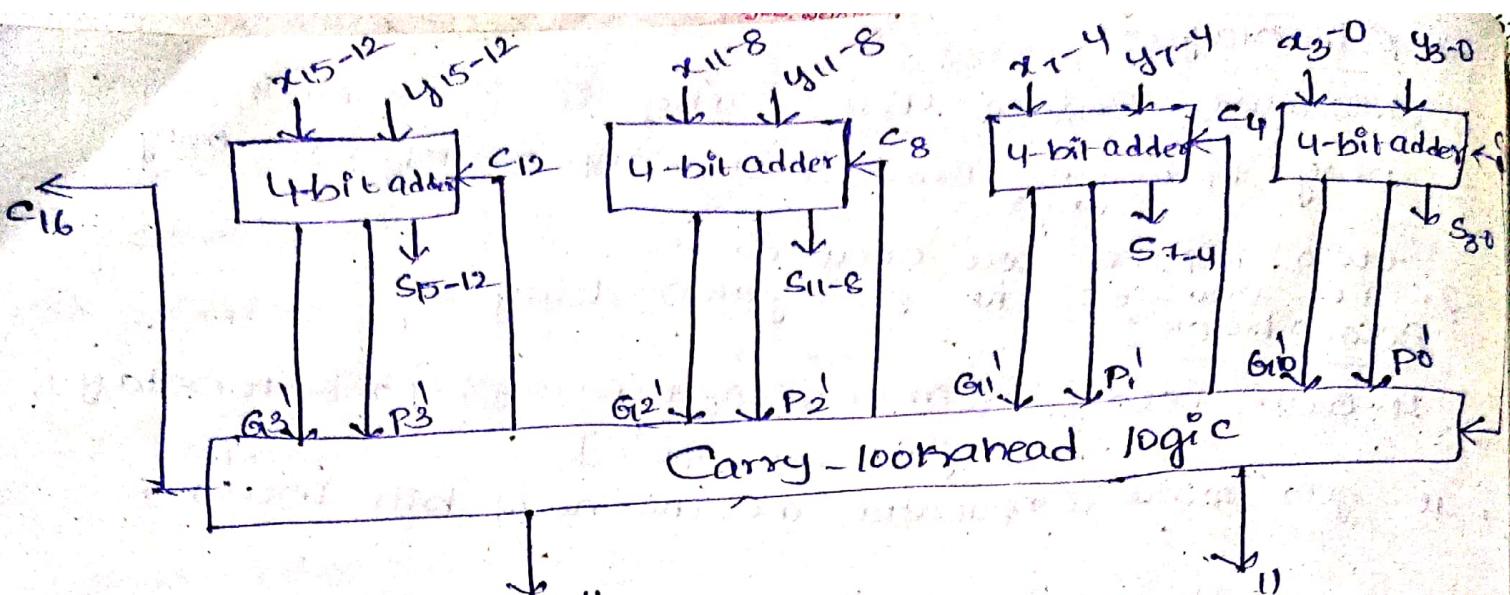


fig: 16-bit carry-lookahead adder built from 4-bit adders.

- The above figure shows a 16-bit adder built from 4-bit adder blocks. These blocks provide new output functions defined as  $G_i^1$  &  $P_i^1$ .
- where  $i=0$  for the first 4-bit block as shown in 4-bit blocks.  $i=1$  for the 2nd 4-bit & so-on.
- In the first block,  

$$P_0^1 = P_3 P_2 P_1 P_0 \text{ &}$$
  

$$G_0^1 = G_3 + P_3 G_2 + P_3 G_1 2 G_1 1 + P_3 G_1 2 G_1 0$$
- In the first level  $G_i^1$  &  $P_i^1$  functions determine whether bit stage  $i$  generates (or) propagates a carry. And the 2nd level  $G_i^1$  &  $P_i^1$  functions determine whether block  $i$  generates (or) propagates a carry.
- The carry  $C_{16}$  ie formed by one of the carry-lookahead circuits and given as,

$$C_{16} = G_3^1 + P_3^1 G_2^1 + P_3^1 P_2^1 G_1^1 + P_3^1 P_2^1 P_1^1 G_0^1 + P_3^1 P_2^1 P_1^1 P_0^1$$

### TOPIC 3: Multiplication of Positive Numbers:

→ The multiplication is a complex operation than addition and subtraction. It can be performed in  $4 \times 4$  (or)  $5 \times 5$ .

$$\begin{array}{r}
 & 1 \ 1 \ 0 \ 1 \\
 & 1 \ 0 \ 0 \ 1 \\
 \hline
 & 1 \ 1 \ 0 \ 1 \\
 & 0 \ 0 \ 0 \ 0 \\
 & 0 \ 0 \ 0 \ 0 \\
 \hline
 & 1 \ 1 \ 1 \ 0 \ 0 \ 1
 \end{array}$$

(13) multiplicand (M)

(9) multiplier (A)

} Partial product.

} Final Product.

- Multiplication process involves generation of partial products, one for each digit in the multiplier.
- In the binary system the partial products are easily defined when the multiplier bit is 0, the partial product is multiplicand.
- The partial product is shifted one position to the left.

Register Configuration (or) How implementation of unsigned Binary Multiplication:-

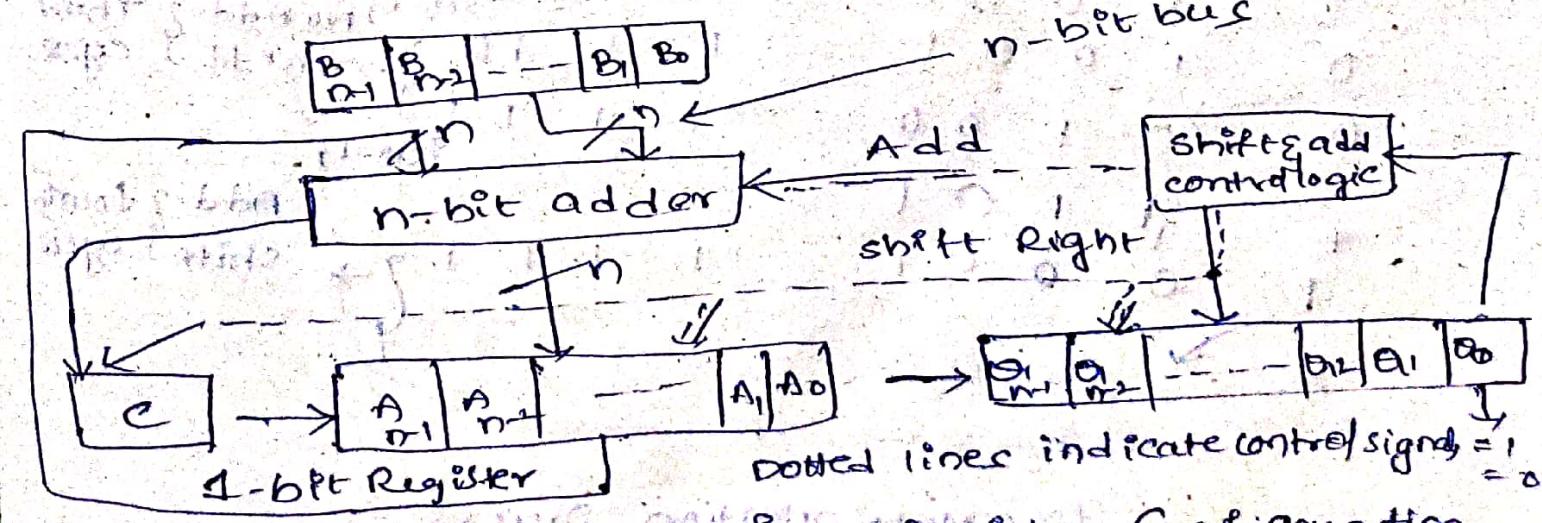


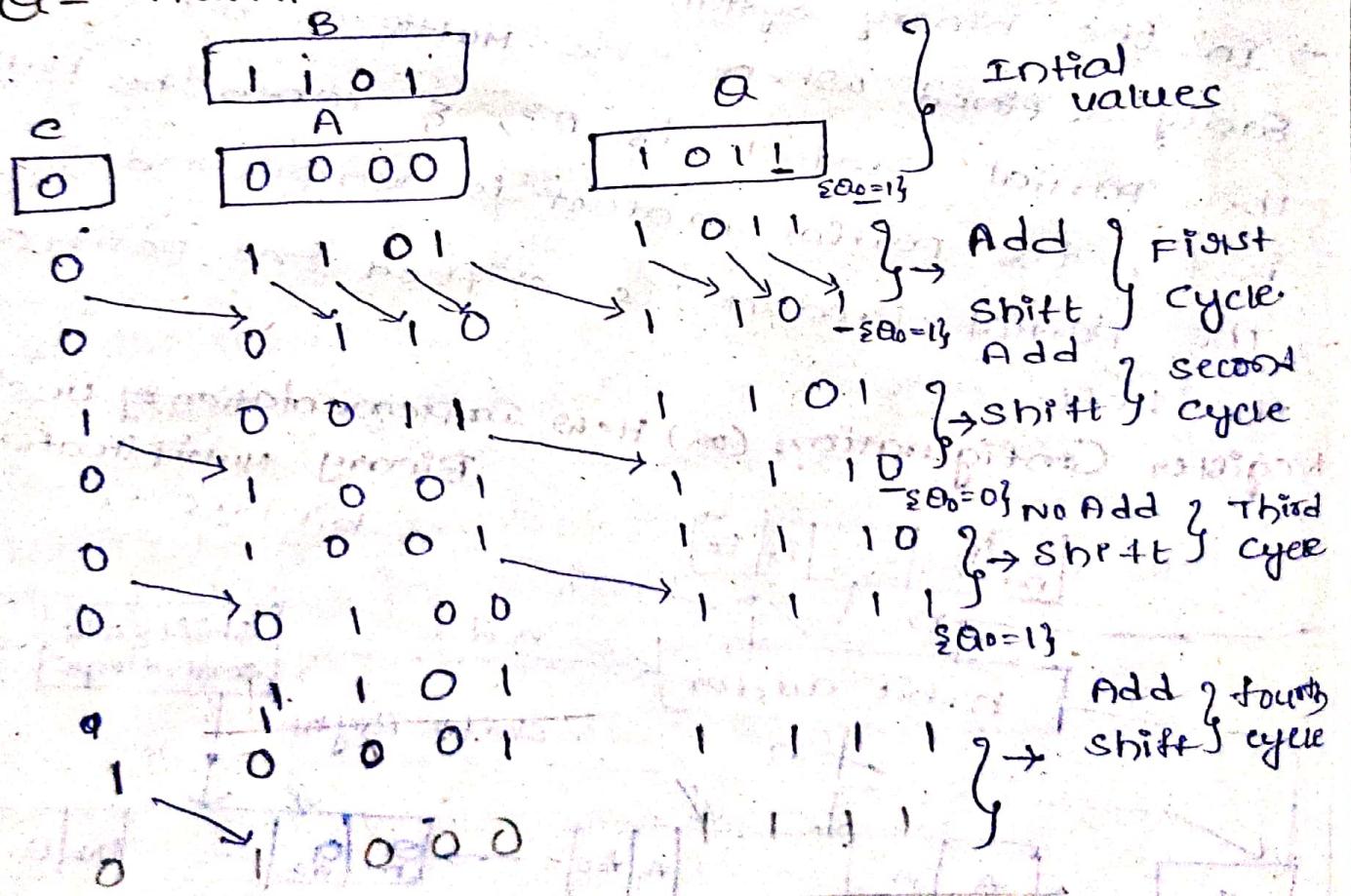
Fig: Register Configuration

- Bit 0 of multiplier operand  $A_0$  of  $A$  is checked.
- If  $A_0$  is one, then multiplicand and partial product are added & all bits C, A and D Registers are shifted to the right one bit.
- If  $A_0$  is zero, then no addition is performed only shift operation is carried out.
- These steps are repeated until to get the desired results in the A and D Registers.

Eg: consider 4-bit Multiplier & Multiplicand

$$B = \text{Multiplicand} = 1101$$

$$A = \text{Multiplier} = 1011$$



Eg: Multiplication Process.

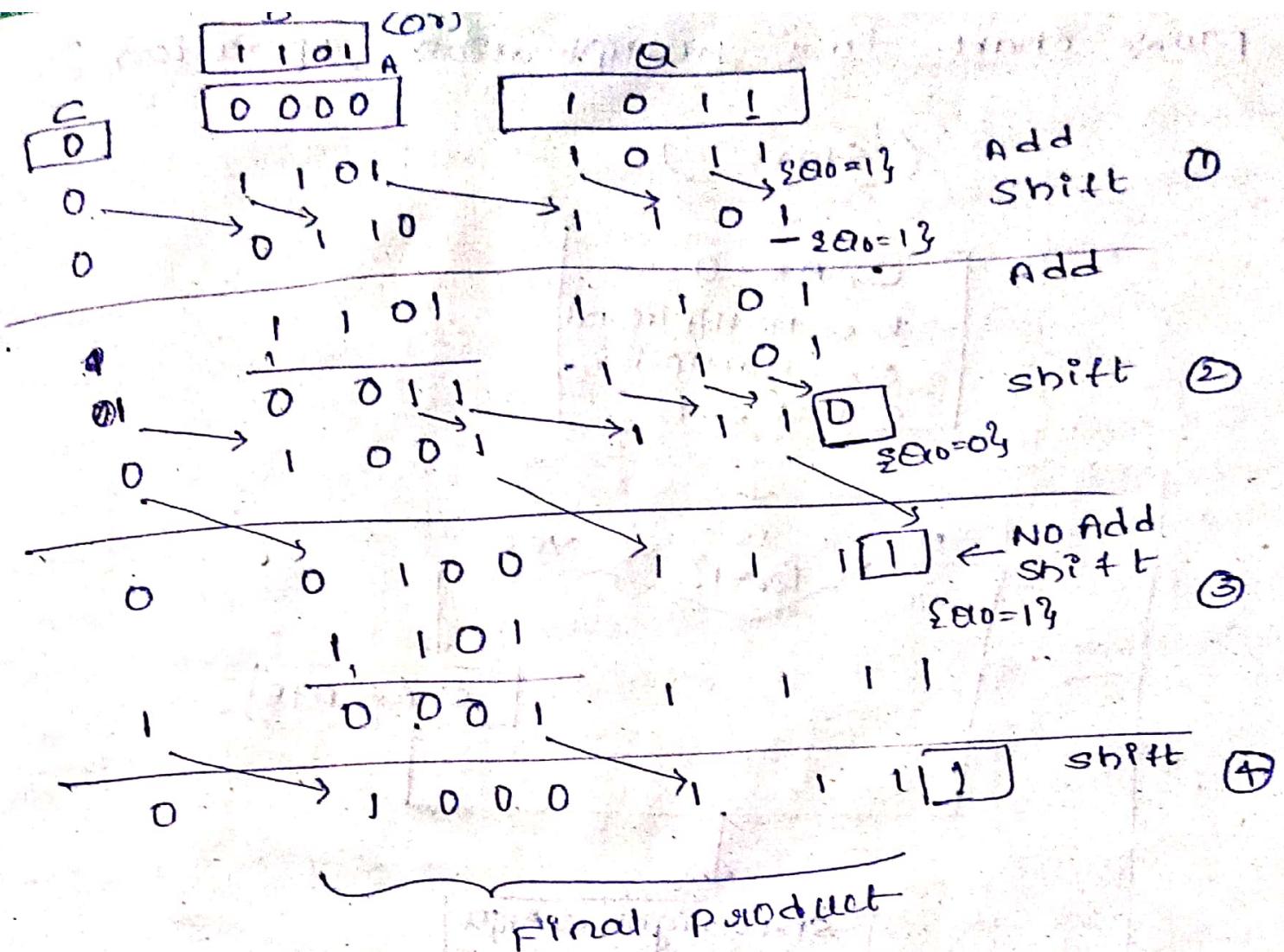


fig: Multiplication process

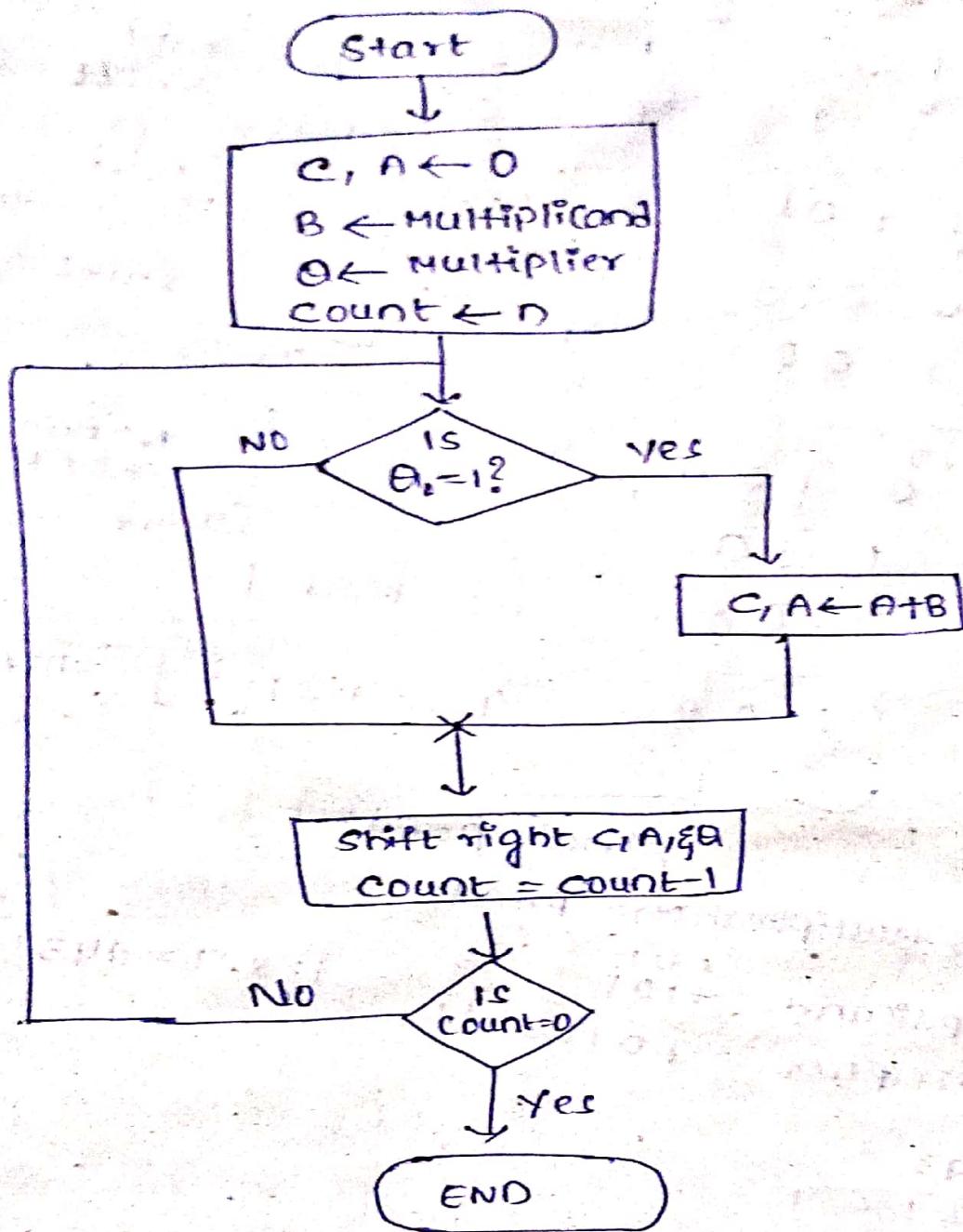
$$\text{B - multiplicand} = \frac{8421}{1101} = 13 \quad 13 \times 11 = 143$$

$$\text{A - multiplier} = \frac{-1011}{-1011} = 11$$

$$\begin{array}{r}
 143 \\
 \times 11 \\
 \hline
 143 \\
 143 \\
 \hline
 154
 \end{array}$$

$$\begin{array}{r}
 1000 \quad 1111 \\
 128 \ 64 \ 32 \ 16 \quad 84 \ 21 \\
 \hline
 128 + 15 = 143
 \end{array}$$

Flow chart for a multiplication operation :-



## TOPIC 4: Multiplication of Signed Numbers (or)

### Signed Operand Multiplication

- In signed operand multiplication, first consider the case of a positive multiplier & negative multiplicand.
- When we add a negative multiplicand to a partial product, we must extend the sign-bit value of the multiplicand.
- For, a negative multiplier is to form the 2's complement of both the multiplier & the multiplicand. And proceed as in the case of a positive multiplier. A technique that works equally well for both negative & positive multiplier called the Booth alg.

$$\begin{array}{r}
 \text{8} \quad \text{4} \quad \text{2} \quad \text{1} \\
 \hline
 M \rightarrow -13 \rightarrow 1 \ 1 \ 0 \ 1 \text{, comp } \rightarrow 0 \ 0 \ 1 \ 0 \\
 D \rightarrow +11 \rightarrow 1 \ 0 \ 1 \ 1
 \end{array}$$

0 0 1 1

As it is -ve number so place '1' in front of it.

$$\begin{array}{r}
 C(-13) \rightarrow 1 \ 0 \ 0 \ 1 \ 1 \\
 (+11) \rightarrow \times 0 \ 1 \ 0 \ 0 \ 1 \\
 \hline
 \end{array}$$

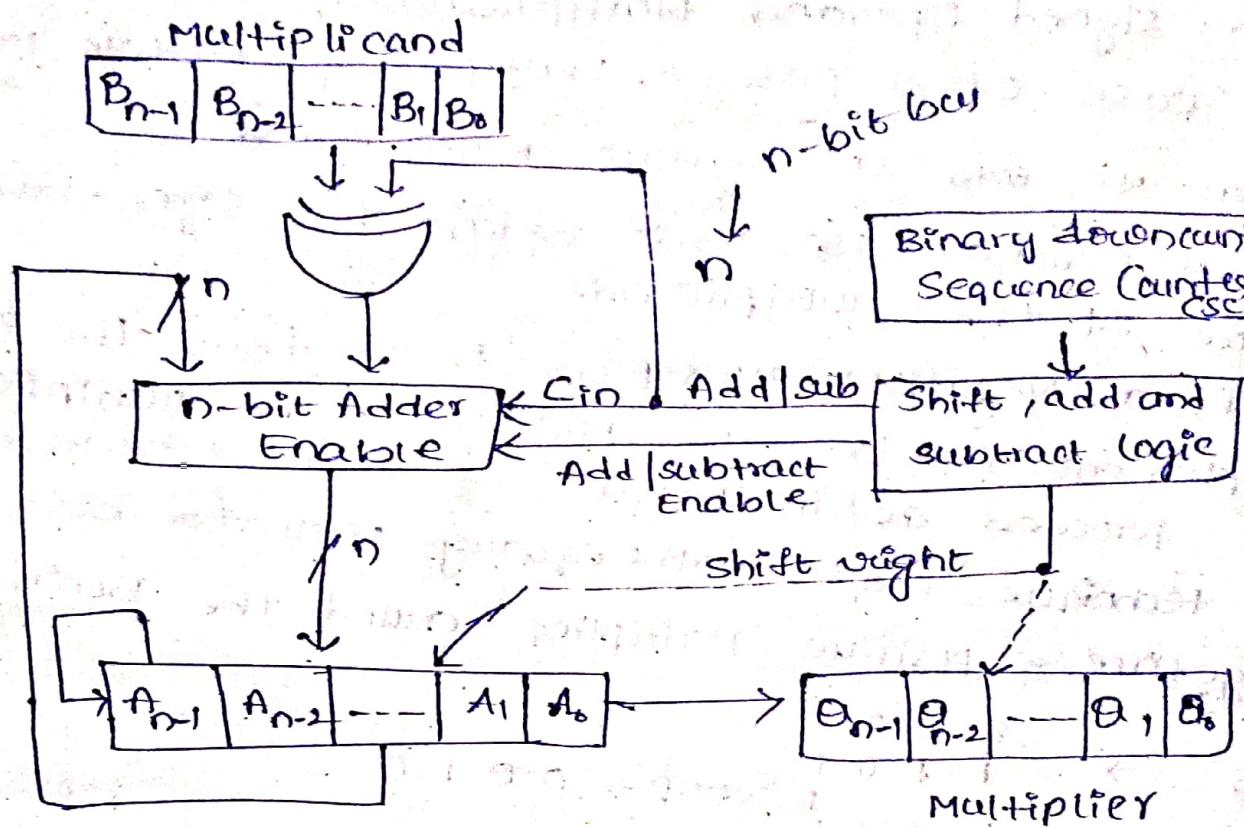
(+ve no. so placed)

$$\begin{array}{r}
 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 0 \ 0 \ 1 \ 1 \\
 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 0 \ 0 \ 1 \ 1 \\
 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \\
 \hline
 1 \ 1 \ 1 \ 0 \ 1 \ 1 \ 1 \ 1 \ 0 \ 0 \ 1 \ 1 \\
 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 1 \\
 \hline
 1 \ 1 \ 1 \ 0 \ 1 \ 1 \ 1 \ 1 \ 0 \ 0 \ 1 \ 1 \ (-143)
 \end{array}$$

sign extension of negative multiplicand.

## Sign Extension Of Negative Multiplicand (Cor)

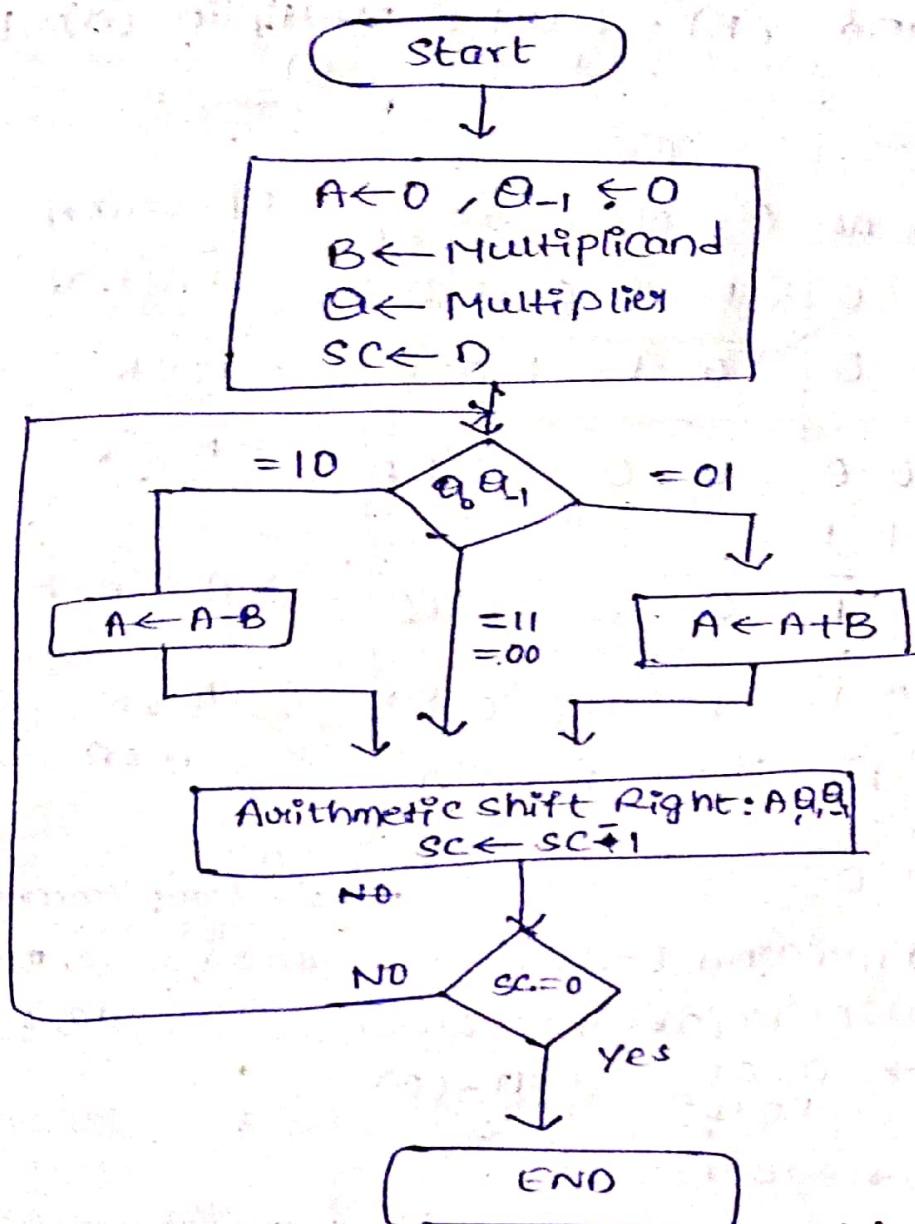
How implementation of signed binary multiplication



$$\text{Initially : } A \leftarrow 0 \text{ & } A_{-1} = 0$$

- The Booth's algorithm can be implemented as shown in above figure.
- The circuit is similar to Circuit for positive no multiplication.
- It consists of n-bit adder, shift, add subtract control Logic & four registers. A, B, Q & Q-1
- The multiplier & multiplicand are loaded into register Q and register B.
- The register A and Q-1 are initially set to 0.
- The sequence counter, SC is set to a number n equal to the number of bits in the multiplier.

## Booth algorithm for signed multiplication :-



case i) : Both positive multiplicand ( $B = 010$ ) [5] & multiplier = 0100

| SC  | A   | Q               | operation              |
|---|---|-----------------|------------------------|
| A <sub>3</sub> A <sub>2</sub> A <sub>1</sub> A <sub>0</sub> | Q <sub>3</sub> Q <sub>2</sub> Q <sub>1</sub> Q <sub>0</sub> | θ <sub>-1</sub> |                        |
| 1 0 0   | 0 0 0 0   | 0 1 0 0         | Initial                |
| 0 1 1   | 0 0 0 0   | 0 0 1 0         | Arithmetic shift Right |
| 0 1 0   | 0 0 0 0   | 0 0 0 1         | A SR                   |
| 1 0 1 1   | 0 . 0 0 1   | 0               |                        |
| 1 0 1 1   | 0 . 0 0 1   | 0               | A ← A-B                |
| 0 0 1   | 1 0 0 0   | 1               | A SR                   |
| 0 0 1   | 1 0 0 0   | 1               | A ← A+B                |
| 0 0 0 0   | 1 0 0 0   | 1               | Res: 0001 0100         |

2. Negative Multiplier ( $5 \times -4$ )

Multiplicand (B) = 0101 & Multiplier (A) = 1100  
 $\begin{array}{r} 1100 \\ \times 0101 \\ \hline 1010 \\ +1010 \\ \hline 1011(-B) \end{array}$

| SC   | A   | Q   | Operation            |
|------|---|---|----------------------|
| 100  | A <sub>3</sub> A <sub>2</sub> A <sub>1</sub> A <sub>0</sub> | Q <sub>3</sub> Q <sub>2</sub> Q <sub>1</sub> Q <sub>0</sub> | Initial              |
| 011  | 0 0 0 0   | 1 1 0 0   | ASR                  |
| 010  | 0 0 0 0   | 0 1 1 0   | A SR                 |
| (+1) | 0 1 1 1   | 0 0 1 1 0   | A $\leftarrow$ A - B |
| 001  | 1 0 1 1   | 0 1 0 1 0   | A SR                 |
| 000  | 1 1 1 0   | 1 1 0 0 1   | A SR                 |

Result = 1 1 1 0 1 1 0 0  $\rightarrow$   $\frac{-20}{(2^8 \text{ complement of } 20)}$

3. Negative Multiplicand (-5)  $\rightarrow$  B

Multiplier (A) = 4  $\Rightarrow$  (0100)

$$\begin{array}{r} \text{Multiplicand: } B \rightarrow 0101 \\ \text{Multiplier: } A \rightarrow 0100 \\ \text{Subtraction: } A - (B) \\ -B \rightarrow 1011 \end{array}$$

$$\begin{array}{r} 1100 \\ 0100 \\ 0100 \\ \hline 0000 \\ 1100 \\ \hline 0100 \\ 0100 \\ \hline 0000 \end{array}$$

| SC   | A   | Q   | Operation            |
|------|---|---|----------------------|
| 100  | A <sub>3</sub> A <sub>2</sub> A <sub>1</sub> A <sub>0</sub> | Q <sub>3</sub> Q <sub>2</sub> Q <sub>1</sub> Q <sub>0</sub> | Initial              |
| 011  | 0 0 0 0   | 0 1 0 0   | ASR                  |
| 010  | 0 0 0 0   | 0 0 1 0   | ASR                  |
| (+1) | 0 0 1 0   | 0 0 0 1 0   | A $\leftarrow$ A - B |
| 001  | 0 1 0 1   | 0 0 0 0 1 0   | ASR                  |
| 000  | 1 0 1 1   | 1 0 0 0 0 1   | A $\leftarrow$ A + B |

4. Both Negative ( $-5 \times -4$ )

$$\text{Multiplicand (B)} = 1011 (-5)$$

$$\text{Multiplier (A)} = 1100 (-4)$$

| sc                               | A   | multiplier A  | operations           |
|----------------------------------|---|---|----------------------|
|                                  | A <sub>3</sub> A <sub>2</sub> A <sub>1</sub> A <sub>0</sub> | A <sub>3</sub> A <sub>2</sub> A <sub>1</sub> A <sub>0</sub> |                      |
| 1 0 0                            | 0 0 0 0   | 1 1 0 0   | Initial              |
|                                  | ↓ 0 0 0 0   | 0 1 1 0   | ASR                  |
| 0 1 1                            | 0 0 0 0   | 0 0 1 1   | ASR                  |
| 0 1 0                            | 0 1 0 1   | 0 0 1 1   | $A \leftarrow A - B$ |
|                                  | ↓ 0 1 0 1   | 1 0 0 1   | ASR                  |
| 0 0 1                            | 0 0 1 0   | 0 1 0 0   | ASR                  |
|                                  | ↓ 0 0 1 0   | 0 1 0 0   |                      |
|                                  |   | 1   |                      |
| Result = 0 0 0 1 . 0 1 0 0 = +20 |   |   |                      |

### Topic 5: Fast multiplication.

- There are two techniques for speeding up multiplication process.
- In first technique guarantees that the maximum number of summands [version of the multiplicand].
- are must be added is  $n/2$  for  $n$ -bit operands.

- The second technique reduces the time needed to add the summands.

1. Bit-Pair Recording of Multipliers

2. Carry-Save Addition of summands.

# 1. Bit-pair recoding of Multipliers.

- 1. for fast multiplication & process.
- 2. Speed the multiplication process.

## Bit-pair recoding Table

| Multiplexer bit pair<br>i+1 i | Multiplexer bit at the right<br>i-1 | Multiplicand selected<br>at position i |
|-------------------------------|-------------------------------------|--|
| 0 0                           | 0                                   | 0 × M                                  |
| 0 0                           | 1                                   | +1 × M                                 |
| 0 1                           | 0                                   | +1 × M                                 |
| 0 1                           | 1                                   | +2 × M                                 |
| 1 0                           | 0                                   | -2 × M                                 |
| 1 0                           | 1                                   | -1 × M                                 |
| 1 1                           | 0                                   | -1 × M                                 |
| 1 1                           | 1                                   | 0 × M                                  |

## (b) Table of multiplicand selection decision

$$\begin{array}{l} \text{Eg :- } B = M = (13) = \underline{\quad 0 \quad 1 \quad 1 \quad 0 \quad 1} \text{ (positive so first 0)} \\ \text{Multiplexer - A = } (-6) = \underline{\quad 1 \quad 1 \quad 0 \quad 1 \quad 0} \text{ [use } \begin{array}{l} 0110 \\ \text{so } 1001 \\ \text{1st] } 11 \\ \hline 1010 \end{array} \text{ for addition]} \end{array}$$

Take multiplier value  $\underline{\quad 1 \quad 1 \quad 1 \quad 0 \quad 1 \quad 0 \quad 0}$  [0]

MSB Position

Place 1

0 -1 -2

[Add '0' in LEB Position]

$\begin{array}{r} 0 \quad -1 \quad -2 \\ \hline 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 1 \quad 1 \quad 0 \end{array}$  C2's complement of Multiplier and add

0 1 1 0 1 → multiplicand

0 -1 -2

$$\begin{array}{r}
 & + 0 0 1 \\
 & 1 1 1 1 0 0 1 1 0 \rightarrow [2^1 \text{ complement of multiplicand}] \\
 & 1 1 1 1 0 0 1 1 \rightarrow \text{and multiply } 2^{10} \\
 & 0 0 0 0 0 0 0 1 0 \\
 & \hline
 & 1 1 0 1 1 0 0 1 0
 \end{array}$$

$\rightarrow$   $(2^1)$  and multiply  $2^{10}$   
 $\downarrow$   
 $\text{Mul } 2^1 \text{ comp} \rightarrow 1 0 0 1 0$

1 0 0 1 1

Ignore carry. Cross verify.  $13 \times -6 = (-78)$

$$\begin{array}{r}
 64.32 16 8 4 2 1 \\
 \overline{+ 1 0 0 1 1 1 0} \\
 \overline{0 1 1 0 0 1} \\
 \overline{0 1 1 0 0 1 0}
 \end{array}$$

→ add 3 more

fig: Multiplication requiring  
only  $n/2$  summands.

Eg: Multiply  $B$  (multiplicand)  $\times A = -11 = 2^1 \text{ comp of } 11 \Rightarrow (1) 1 0 1 1$ ,  
 $\text{Multiplicand} \rightarrow AB = 27 = 0 1 1 0 1 1$ .  $\frac{1}{1} \overbrace{1 0 1 0 1}$

Take Multiplier value for Recoding

$0 1 1 0 1 1$   $\boxed{0}$

$+2 -1 -1$

Multiplicand value is  $\rightarrow 1 1 0 1 0 1 \times$

$+2 -1 -1$  (leave the bit)

$$\begin{array}{r}
 \text{(extend 6 morebit)} \quad 0 0 0 0 0 0, 0 0 0 1 0 1 0 \quad (2^1 \text{ comp of } B) \\
 \text{(last bit is } 1 \text{ so extending)} \quad \overline{1 1 1 0 1 0 1 0 1 1} \quad B \times C(10)
 \end{array}$$

## 2. Carry - save Addition of summands

- Multiplication requires the addition of several summands.
- A technique called carry-save addition (CSA)
- speeds up the addition process.
- consider the array for  $4 \times 4$  multiplication shown in fig. with the first row consisting of just the AND gates that implement the bit products  $m_3a_0, m_2a_0, m_1a_0$ , and  $m_0a_0$ .

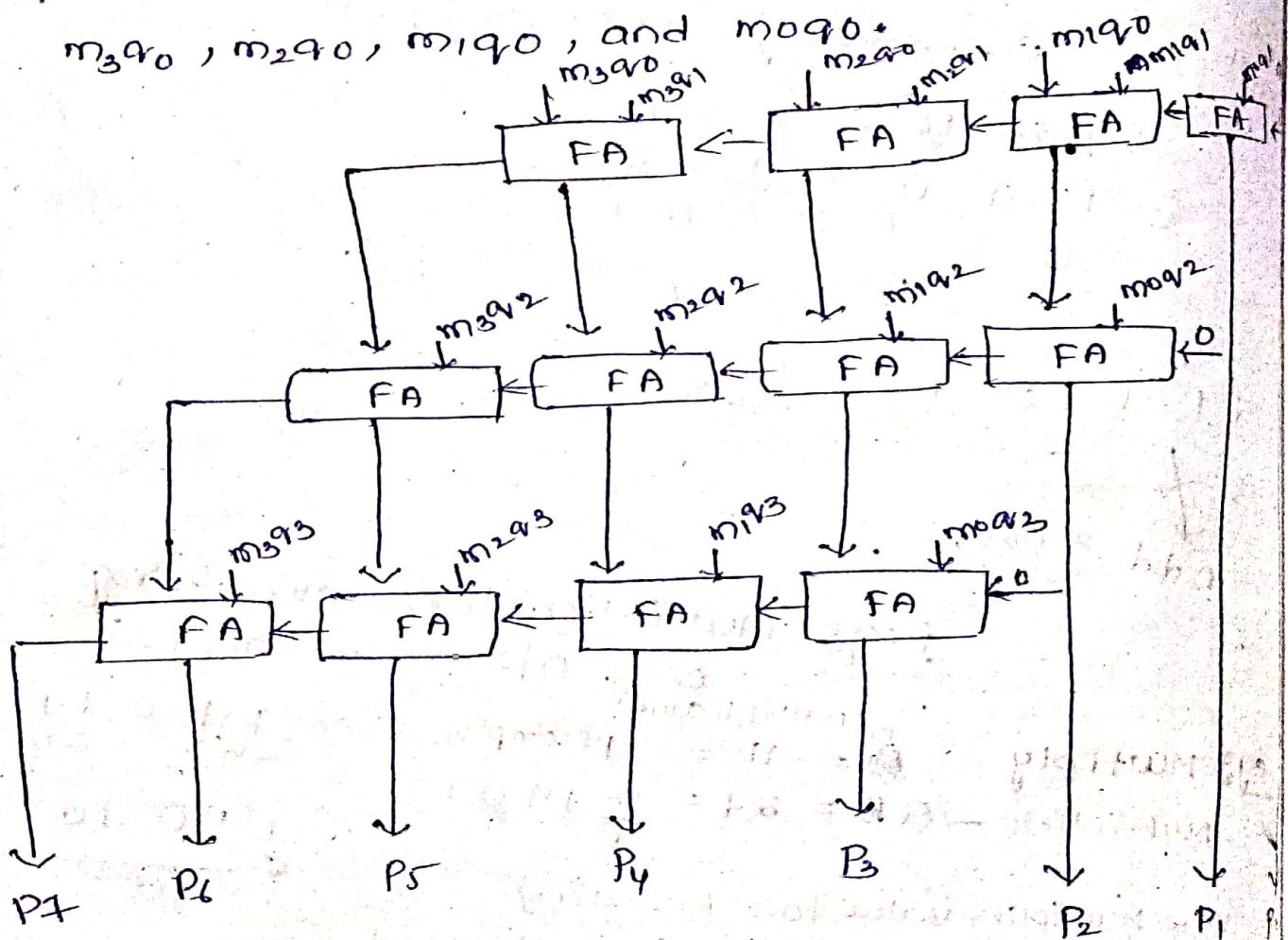


fig.(a) Ripple carry array.

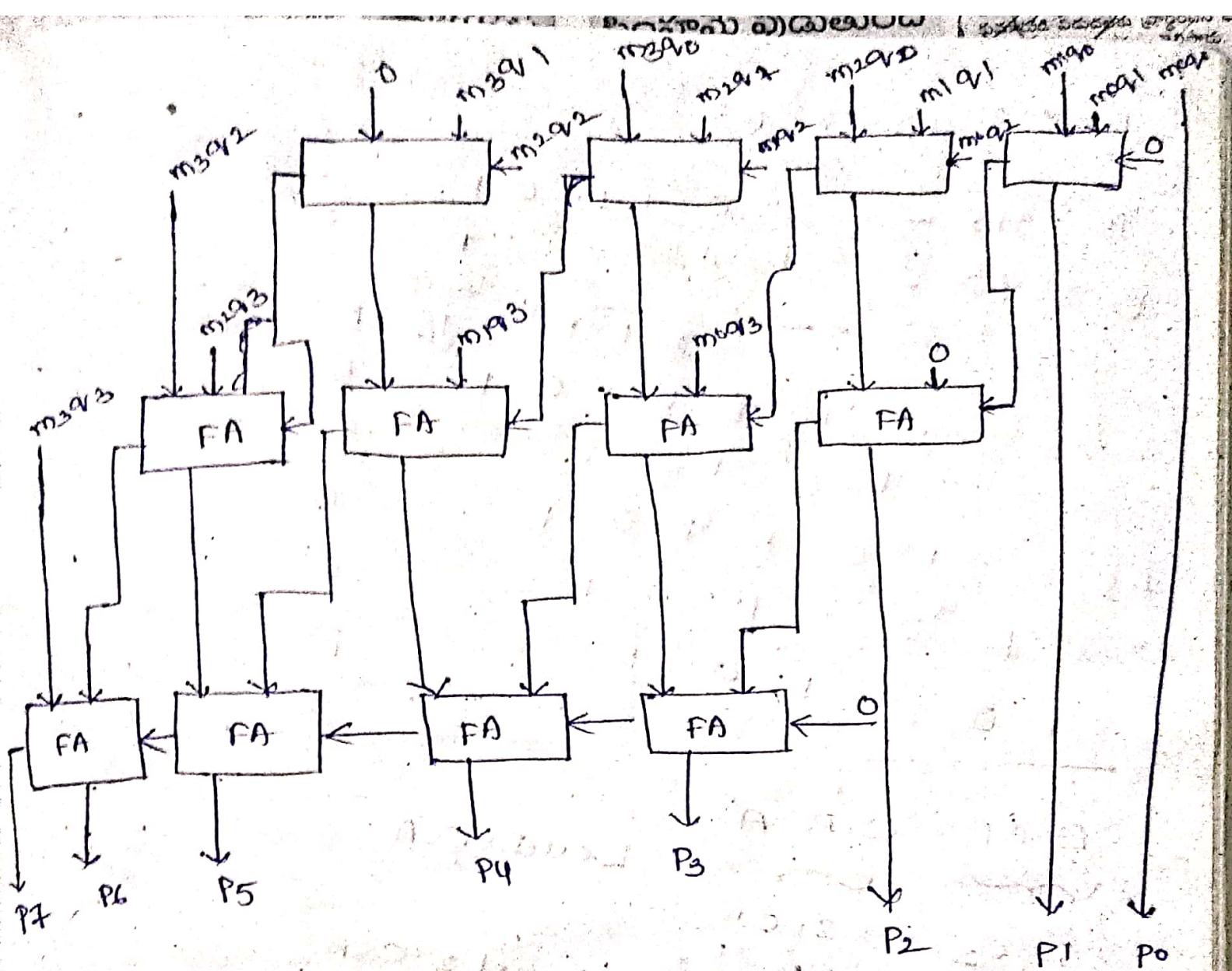


fig : (b) carry - save array

The above fig (a) And (b) are Ripple-Carry and Carry-save arrays for the multiplication operators,  $M \times N = P$  for 4-bit operands.

$$\text{Ex: } 45 \times 63$$

$$\begin{array}{r}
 B \quad 45 = 1 \ 0 \ 1 \ 1 \ 0 \ 1 \\
 O \quad 63 = 1 \ 1 \ 1 \ 1 \ 1 \ 1 \\
 \hline
 & 1 \ 0 \ 1 \ 1 \ 0 \ 1 & A \\
 & 1 \ 0 \ 1 \ 1 \ 0 \ 1 & B \\
 & 1 \ 0 \ 1 \ 1 \ 0 \ 1 & C \\
 & 1 \ 0 \ 1 \ 1 \ 0 \ 1 & D \\
 & 1 \ 0 \ 1 \ 1 \ 0 \ 1 & E \\
 & 1 \ 0 \ 1 \ 1 \ 0 \ 1 & F \\
 \hline
 & 1 \ 0 \ 1 \ 1 \ 0 \ 1 &
 \end{array}$$

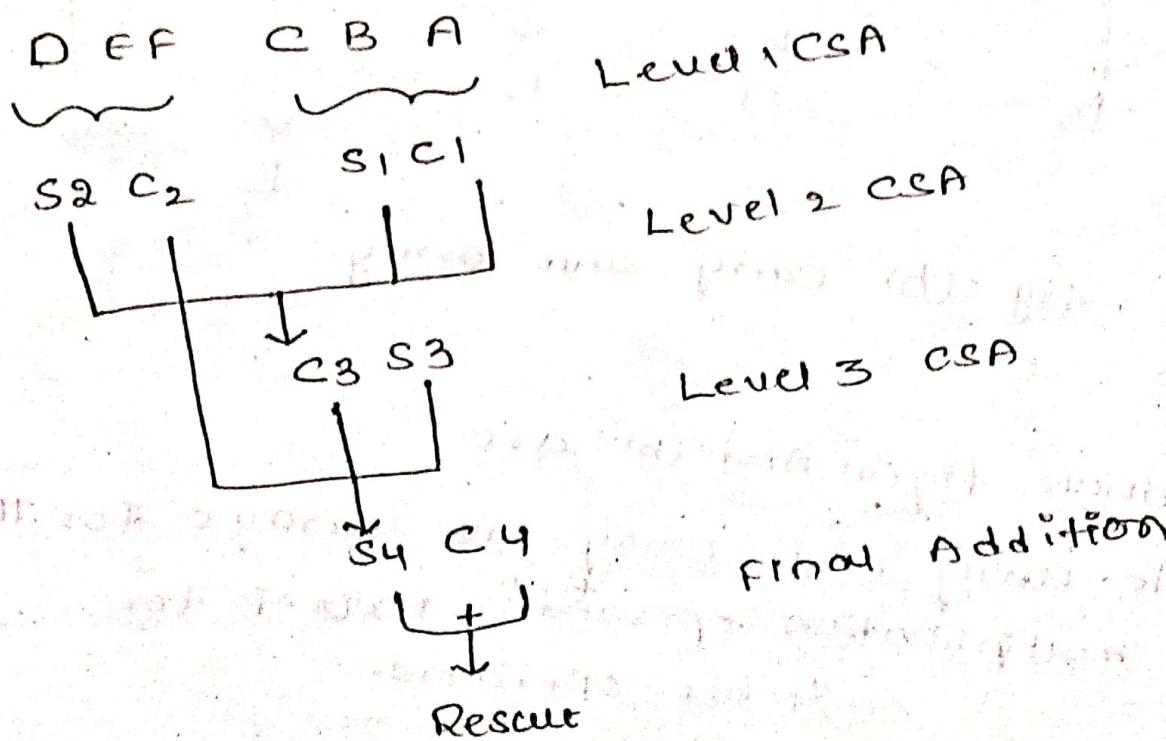


fig: schematic representation of the carry  
save addition.

$$\begin{array}{r}
 101101 \text{ A} \\
 101101 \text{ B} \\
 101101 \text{ C} \\
 \hline
 110000 \text{ F. S1} \\
 0111100 \times \text{C1} \\
 \hline
 110000011 \text{ S1} \\
 0111100 \times \text{C1} \\
 \hline
 110000011000 \text{ S2} \\
 \hline
 110110100011 \text{ S3} \\
 00010110000 \times \text{C3} \\
 \hline
 1101101100011 \text{ S3} \\
 00010110000 \times \text{C2} \\
 \hline
 1011110100011 \text{ S4} \\
 10110100000 \times \text{C4} \\
 \hline
 1011110100011 \text{ S}
 \end{array}$$

↓  
283  
45×63

fig. The multiplication example for ~~45~~ 45×63  
 performed, using carry save addition

## Topic : 6 :- Integer Division

The division is more complex than multiplication.

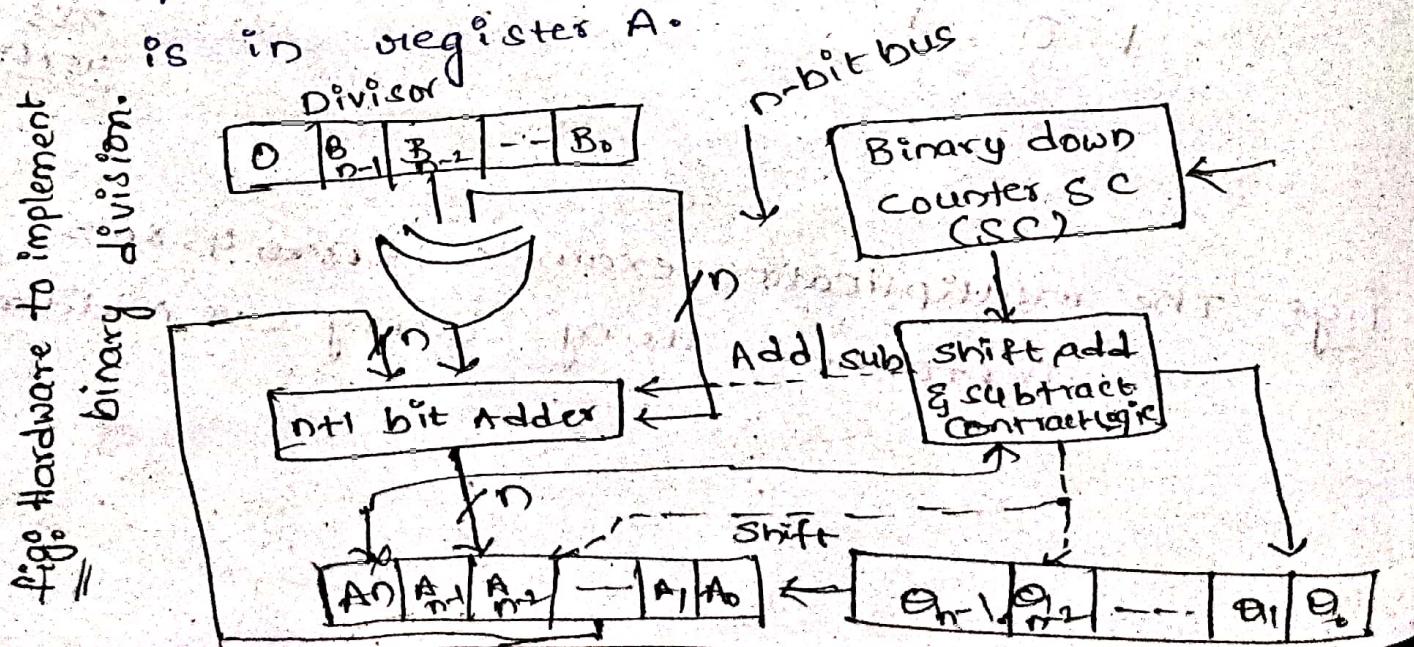
$$\begin{array}{r}
 \text{→} \\
 \text{Divisor} \\
 \overline{)12) \cdot 169 \quad (14 \leftarrow \text{Quotient}} \\
 12 \\
 \hline
 49 \quad \leftarrow \text{Partial remainder} \\
 48 \\
 \hline
 1 \quad \leftarrow \text{Remainder}
 \end{array}$$

## Restoring Division :-

- It consists of n-bit binary adder, shift, add, and sub control logic and registers A, B and Q.

→ Divisor and dividend are loaded into register B and register Q respectively.

→ Register A is initially set to zero. After the division process completes, the n-bit quotient is in register Q and the remainder is in register A.



## Division Operation

Steps 8- (Restoring division):

1. shift A and Q left one binary position.
2. subtract divisor from A & place answer back.
3. If the sign bit of A is '1', set Qd to 0 and add divisor back to A. Otherwise, set Qd to 1.
4. Repeat steps 1, 2, 3 in 4 times.

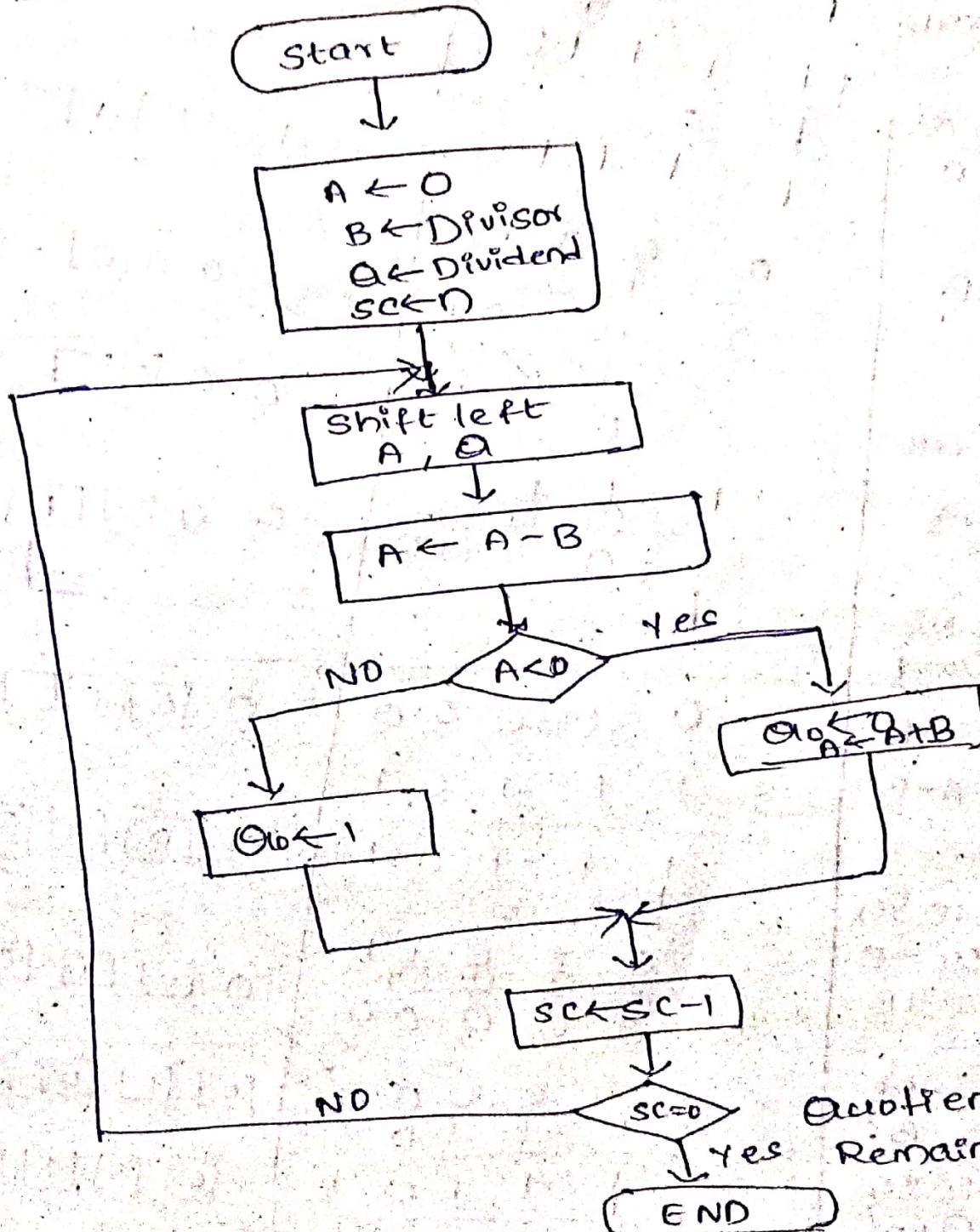


Fig:- flowchart for restoring division operation-

e.g.: Dividend = 1100 (a) (12)  
 Divisor = 0011 (B) (3)

| Steps  | A  | Q (dividend) | Divisor   |
|--|--|--------------|-----------|
| Initial  | 0 0 0 0 0  | 1 1 0 0      | 0 0 1 1   |
| Shift Left<br>A $\leftarrow A - B$ (sub)<br>AD < 0 Yes<br>so Q <sub>0</sub> = 0                    | 0 0 0 0 1<br>1 1 1 0 1<br>-----<br>1 1 1 1 0       | 1 0 0 0 0    | 1 0 0 1 1 |
| A $\leftarrow A + B$   | 0 0 0 0 1<br>0 0 0 1 1<br>-----<br>0 0 0 0 0       | 1 0 0 0 0    | 0 0 0 1 1 |
| Shift Left<br>A $\leftarrow A - B$<br>AD < 0 NO<br>so Q <sub>1</sub> = 1                           | 0 0 0 0 0<br>0 0 0 1 1<br>1 0 1 1 0 1<br>0 0 0 0 0 | 0 0 0 0 0    | 0 0 0 1 1 |
| Shift Left<br>A $\leftarrow A - B$<br>AD < 0 YES<br>set Q <sub>2</sub> = 0<br>A $\leftarrow A + B$ | 0 0 0 0 0<br>1 1 1 0 1<br>-----<br>1 1 1 0 1       | 0 0 1 0 0    | 0 0 1 0 0 |
| Shift Left<br>A $\leftarrow A - B$<br>AD < 0 YES<br>so Q <sub>3</sub> = 0<br>A $\leftarrow A + B$  | 0 0 0 0 0<br>1 1 1 0 1<br>-----<br>0 0 0 0 0       | 0 1 0 0 0    | 0 1 0 0 0 |
| Remainder Q in A.<br>A Restoring Division example.   | 0 0 0 0 0<br>1 1 1 0 1<br>-----<br>0 0 0 0 0       | 0 1 0 0 0    | 0 1 0 0 0 |

## NDD - Restoring Division or

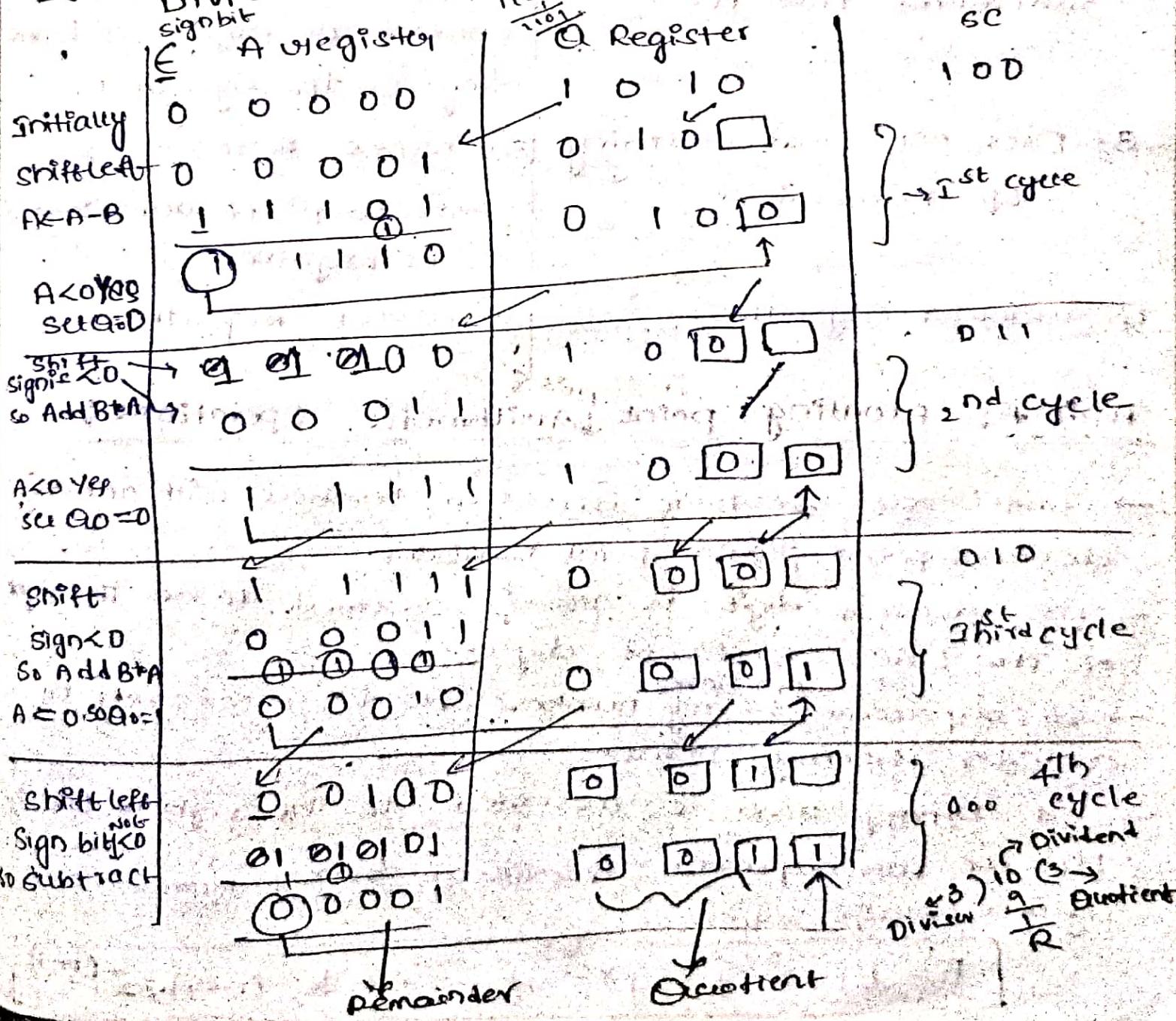
1. If the sign of A is 0, shift A, and Q.
2. If the left bit position & subtract divisor from A, otherwise shift A & Q Left and add divisor to A, if the sign of A is 0, set Q<sub>0</sub> to 1; otherwise set Q<sub>0</sub> to 0.
3. Repeat steps 1 & 2 for 'n' times.
4. If the sign of A is 1, add divisor to A.

Non - Restoring division example :-

$$\text{Dividend} = 1010 \quad (10) \rightarrow (0)$$

$$\text{Divisor} = 0011 \quad (3) \rightarrow (B)$$

signbit  
eg: A register



కడవ మాట ఎంచు అవినార్థక  
ప్రాణ కుటుంబమే

## Differences b/w Restoring & Non-Restoring Division

### Restoring Division

1. Needs restoring of register A if the result of the subtraction is negative.
2. In each cycle, content of Register A is first shifted left and then divisor is subtracted from it.
3. Does not need restoring of remainder.
4. Slower Algorithm.

### Non-Restoring Division

1. Does not need restoring.
2. In each cycle content of Register A is first shifted left & then divisor is added (or) subtracted with the content of register A depending on the sign of A.
3. Needs Restoring of remainder if remainder is negative.
4. Faster Algorithm.

## Topic 7: Floating point Arithmetic Operations

- Until now, we have discussed numbers with only decimal point, fixed point numbers.
- The decimal digit is always assumed to be the right of the least significant digit.
- To represent binary numbers it is necessary to consider binary point.
- If binary point is assumed to the right of the sign bit we can represent fractional binary number as given below.

$$B = b_0 \times 2^0 + b_1 \times 2^{-1} + \dots + b_{-(n-1)} \times 2^{-(n-1)}$$

→ with this fractional number system we can represent the fractional numbers in the following range :  $2^{-(n-1)} \leq F \leq 1 - 2^{-(n-1)}$

→ If  $n=32$ , then the value range is approximately

$$2^{(-31)} \leq F \leq 1 - 2^{-(31)}$$

$$0 \left[ \frac{-31}{2} \right] = -2 \cdot 3285 \times 10^{-31}$$

→ But this range is not sufficient to represent fractional numbers.

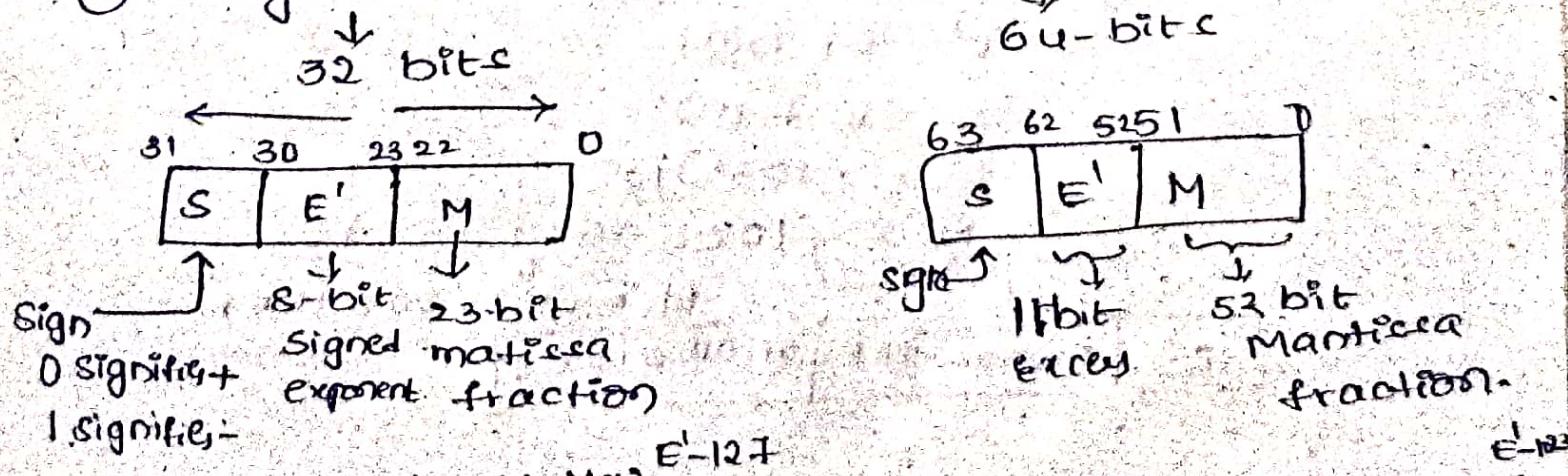
→ To accommodate very large integers and very small fractions, a computer must be able to represent numbers and operate on them in such a way that the position of the binary point is variable and is automatically adjusted.

→ In this case, the binary point is said to float, and the numbers are called "floating-point numbers".

### Floating point Representation:-

Floating point can be represented in two forms

(a) single-precision      (b) Double precision.



$$\text{Value Represented} = \pm 1 \cdot M \times 2^{E-127}$$

$$\text{Value Represented} = \pm 1 \cdot M \times 2^{E-1023}$$

రద్దు మూలిక ఎండ్రు అవస్థలో

TOP 76  
BOTTOM -001

Eg:  $(1460 \cdot 125)_{10} \rightarrow$

$$(1460 \cdot 125)_{10} = (10110110100 \cdot 001)_2$$

↓  
Normalize

$$1 \cdot 0110110100001 \times 2^{10} \rightarrow \text{Normalized Result}$$

Single precision :-

$s=0$  for positive

$E \rightarrow$  Exponent = 10

$M \rightarrow$  Mantissa (after decimal)

$M = 0110110100001$

$$E' = E + bias$$

$$E' = 10 + 1023 = 137$$

$$\boxed{E' = 137}$$

$$E' = (137)_{10} = (10001001)_2$$

| S | E'       | M             |
|---|----------|---------------|
| 0 | 10001001 | 0110110100001 |

Double precision :-

$s=0$

$$E' = E + bias$$

$$= 10 + 1023$$

$$E' = (1023)_{10}$$

$$M = (100000001001)_2$$

| S | E'           | M             |
|---|--------------|---------------|
| 0 | 100000001001 | 0110110100001 |

## Floating point Addition and subtraction:-

Consider two floating point numbers:

$$x = m_1 \cdot r^{e_1} \text{ and}$$

$$y = m_2 \cdot r^{e_2} \text{ Assume } e_1 > e_2.$$

### Rules for Addition and subtraction :-

1. Select the number with a smaller exponent and shift its mantissa right, a no. of steps equal to the difference in exponents  $|e_2 - e_1|$ .
2. set the exponent of the result equal to the larger exponent.
3. perform addition/subtraction on the mantissa and determine the sign of the result.
4. Normalize the result, if necessary.

Eg: Add single precision floating point numbers

A and B, where  $A = 44900000H$  and

$$B = 42A00000H$$

→ In the first step, convert these Decimal numbers to binary.

Sol: Represent numbers in single precision format

$$\begin{array}{r} M \\ \times \quad E \\ \hline \end{array}$$

$$A = 0 \quad 10001001001000000000000000000000$$

$$B = 0 \quad 10000101101000000000000000000000$$

$$\text{Exponent for } A = 10001001 \Rightarrow 137$$

$$\text{Actual exponent} = 137 - 127 \text{ (Bias)}$$

$$e_1 = 137 - 127 \\ = 10$$

$$\text{Exponent for } B = 10000101$$

$$= 133$$

$$\text{Actual exponent} = 133 - 127 \text{ (Bias)} = 6$$

∴ Number B has smaller exponent with difference  
4. Hence its mantissa is shifted right by  
4 bits as shown below.

Step 2: Shift Mantissa

$$\text{Mantissa of } B = 0 \ 1 \ 0 \ 0 . 0 \ 0 \ 0 \ 0 \ 0 \ 0 \cdots 0$$

$$\text{Ist shift Mantissa of } B = \begin{array}{ccccccccc} \downarrow & \searrow & \downarrow & \searrow & \downarrow & \searrow & \downarrow & \searrow & \downarrow \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{array} \cdots 0$$

$$\text{IInd " " } B = \begin{array}{ccccccccc} \downarrow & \searrow & \downarrow & \searrow & \downarrow & \searrow & \downarrow & \searrow & \downarrow \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \end{array} \cdots 0$$

$$\text{IIIrd shift Right " " } = \begin{array}{ccccccccc} \downarrow & \searrow & \downarrow & \searrow & \downarrow & \searrow & \downarrow & \searrow & \downarrow \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \end{array} \cdots 0$$

$$\text{IVth shift Right " " } = \begin{array}{ccccccccc} \downarrow & \searrow & \downarrow & \searrow & \downarrow & \searrow & \downarrow & \searrow & \downarrow \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \end{array} \cdots 0$$

$$\text{shifted Mantissa of } B = 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \cdots 0$$

Step 3: Add Mantissas

$$\text{Mantissa of } A = 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \cdots 0$$

$$\text{Mantissa of } B = 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \cdots 0$$

$$\text{Mantissa of Recut} = \underline{0 \ 0 \ 1 \ 0 \ 0 \ 1 \ 0 \ 0 \cdots 0}$$

As both numbers are positive, sign of the result

$$\therefore \text{Result} = 0 \ 1 \ 0 \ 0 \ 1 \ 0 \ 0 \ 1 \ 0 \ 0 \cdots 0$$

$$= 44920000H$$

→ we have seen addition process for floating

point numbers. Similar process is used for

Subtraction of floating point numbers.

→ In subtraction, two mantissas are subtracted instead of addition and the sign of greater mantissa is assigned to the result.

Eg: subtract single point precision & floating point numbers A and B ( $A-B$ ), where  $A=44900000H$ ,  $B=42A00000H$ .

Sol:

Step 1: Represent numbers in single precision format

$$A = 0 \ 100010010010000 \dots 0$$

$$B = 0 \ 10000101010000 \dots 0$$

$$\begin{aligned} \text{Exponent for } A &= 10001001 \\ &= 137 \end{aligned}$$

$$\begin{aligned} \text{Actual exponent} &= 137 - 127 \text{ (Bias)} \\ &= 10 \end{aligned}$$

$$\begin{aligned} \text{Exponent for } B &= 10000101 \\ &= 133 \end{aligned}$$

$$\therefore \text{Actual exponent} = 133 - 127 \text{ (Bias)}$$

$$= 6$$

$\therefore$  Number B has smaller exponent with difference

$\therefore$  Hence its mantissa is shifted right by 4 bits as shown below.

Step 2: Shift Mantissa

shifted Mantissa of B = 00000100-0

Step 3: subtract Mantissa

$$\begin{array}{r} \text{Mantissa of } A = 00100000 \dots 0 \\ \text{Mantissa of } B = 00000100 \dots 0 \end{array}$$

$$\begin{array}{r} \text{Mantissa of } B = 00000100 \dots 0 \\ \text{Mantissa of } A = 00100000 \dots 0 \end{array}$$

$$\therefore -B \rightarrow \underline{\underline{11111011}}$$

$$\underline{\underline{11111011}} \quad \underline{\underline{0001001100}} \quad \dots \dots 0$$

$$\therefore B \rightarrow \underline{\underline{11111100}}$$

$$\underline{\underline{11111100}} \quad \underline{\underline{0001001100}} \quad \dots \dots 0$$

Mantissa for A is greater than mantissa for B.

Therefore sign of result is sign of A.

$$\text{Result} = 0 \ 1000100100011100 \dots 0$$

$$= 448E0000H$$

## Problems in Floating Point Arithmetic &

Mantissa Overflow:- The addition of two mantissas of the same sign may result in a carry-out of the MSB bit. If so, the mantissa is shifted right and the exponent is incremented. [ $M = M_{MAX}$ ]

Mantissa Underflow :- If the process of aligning mantissas, digits may flow off the right end of the mantissa. In such case truncation methods such as chopping, rounding are used. [ $M > M_{MIN}$ ]

Exponent Overflow :- Exponent overflow occurs when a positive exponent exceeds the maximum possible exponent value. In some systems this may be designed as + $\infty$  or - $\infty$ . [ $E = E_{MAX}$ ]

Exponent Underflow :- Exponent underflow occurs when a negative exponent exceeds the maximum possible exponent value. In such cases, the number is designed as zero. [ $E > E_{MIN}$ ]

Flowchart & Algorithm for floating point Addition &

Algorithm :- To write algorithm for floating point subtraction, addition as follows.

Declare registers : AM [0: M-1] ; For mantissa of A.  
BM [0: M-1] ; For mantissa of B.  
AE [0: E-1] ; For exponent of A.  
BE [0: E-1] ; For exponent of B.  
E [0: E-1] ; For temporary storage of exponent.

Step 1: Read Numbers:  $AM \leftarrow$  Mantissa of A  
 $BM \leftarrow$  Mantissa of B  
 $AE \leftarrow$  exponent of A  
 $BE \leftarrow$  exponent of B

Step 2: Compare the two exponents:  $E = AE - BE$

Step 3: Equalize

If  $E < 0$  then

right shift  $AM$ ;

$E = E+1$  and go to step 3;

If  $E > 0$  then

right shift  $BM$ ;

$E = E-1$  go to step 3;

Step 4: Add Mantissas:  $AM = AM + BM$  and  
 $E = \text{Max}(AE, BE)$

Step 5: Check exponent overflow

If overflow occurs during mantissa addition

then if  $E = E_{\text{max}}$  then report overflow

else

right shift  $E$  &  $AM$

$E = E+1$  go to step 3.

Step 6: Adjust for zero result

If  $AM = 0$  then  $E = 0$ ; go to step 3.

Step 7: Normalize result

If  $AM$  normalized then go to step 3.

Step 8: Check exponent underflow

If  $E > E_{\text{min}}$  then

left shift  $AM$ ;

$E = E-1$  go to step 7.

Step 9: Stop.

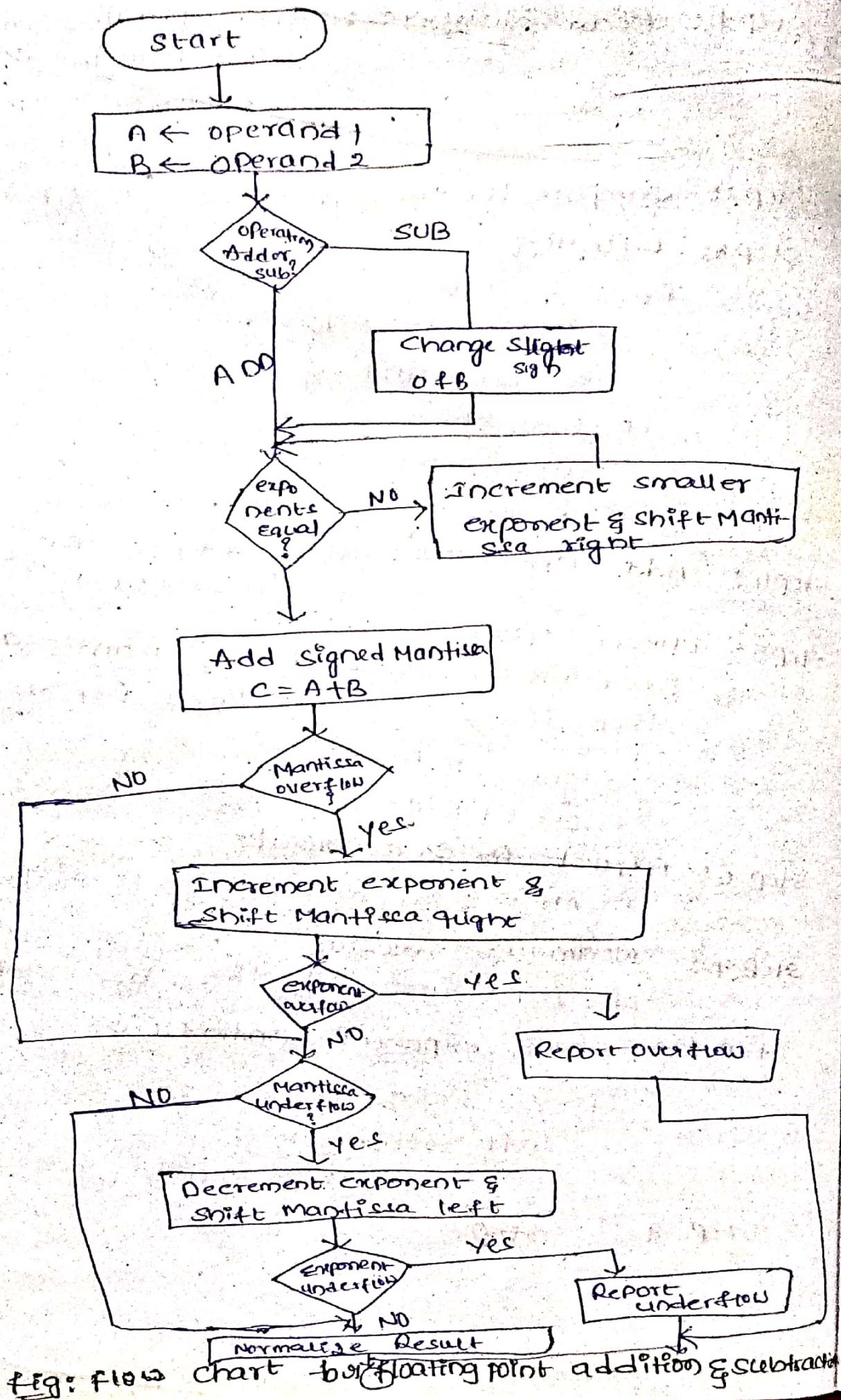
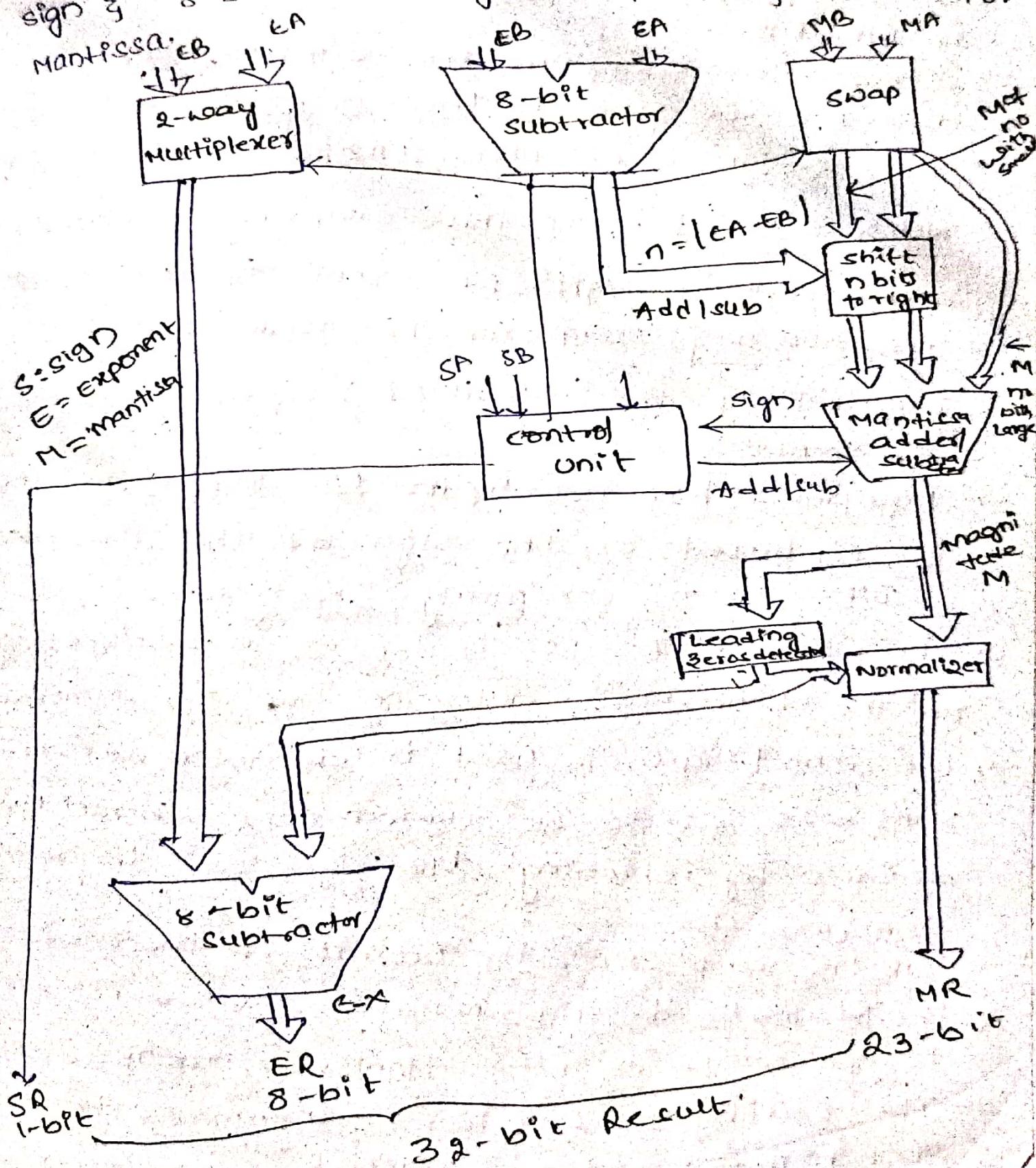


Fig: flow chart for floating point addition & subtraction

# Implementing Floating Point Operations (Add & sub) :-

- It shows the hardware implementation for the addition & subtraction of 32-bit floating point operations.
- It has the single precision format i.e. 1 bit for sign & 8-bits for signed exponent & 23-bit for mantissa.



- To compare the exponents of Numbers to determine a number with a smaller exponent & then determine to shift the Mantissa of that number with, so that its exponent matches with the other Number.
- The steps for the 8-bit subtractor are Exponent values of two number & the operation performed is  $E_A - E_B$ .
- The subtraction gives the difference  $D$ . The difference is sent to the SHIFTER unit.
- If the sign  $D \neq 0$ , then  $E_A \geq E_B$  and if to swap Networks per 1. In this swapping is enabled and Mantissa  $M_A$  is sent to the SHIFTER.
- The two way Multiplexer is used to set the exponent of the result ( $E$ ) equal to the larger exponent.
- Multiplexer has two inputs  $E_A$  and  $E_B$  and the output is based on the sign of the difference result from comparing exponents.
- The next step is to perform addition & subtraction on the Mantissas & determine the sign of the result.
- The control logic is used to determine whether the Mantissas are to be added (or) subtracted, it decides by checking the signs of the operands (SA and SB).
- In the last step, the result of the Mantissa ( $M$ ) is normalized. The normalized value is truncated to generate the 23-bit Mantissa, MR of the result.
- The Leading zeros detector determines the number of bit shifts,  $x$  to be applied to the Mantissa ( $M$ ).