

Computer Arithmetic

①

Addition and Subtraction of signed numbers, Design of fast adders, Multiplication of positive numbers, signed-operand multiplication, Fast multiplication, Integer division, Floating point numbers and operations

A basic operation in all digital 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 arithmetic and logic unit (ALU) of the processor.

Addition and Subtraction of signed numbers :-

The following table shows the logic truth table for the sum and carry out functions for adding equally weighted bits x_i and y_i .

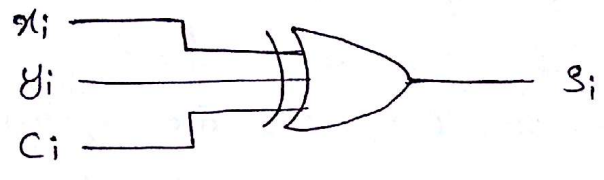
x_i	y_i	Carry in C_i	Sum S_i	Carry out C_{i+1}
0	0	0	0	0
0	0	1	1	0
0	1	0	1	0
0	1	1	0	1
1	0	0	1	0
1	0	1	0	1
1	1	0	0	1
1	1	1	1	1

here

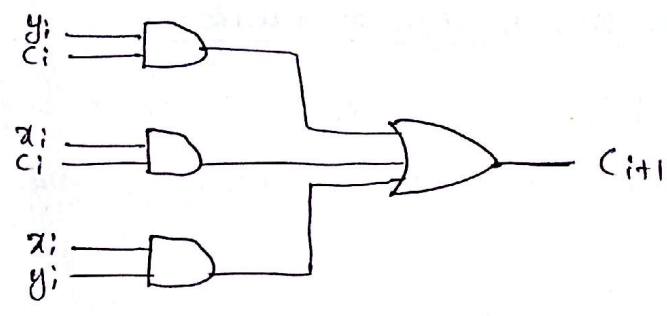
$$S_i = \bar{x}_i \bar{y}_i c_i + \bar{x}_i y_i \bar{c}_i + x_i \bar{y}_i \bar{c}_i + x_i y_i c_i = x_i \oplus y_i \oplus c_i$$

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

The sum expression can be implemented with a 3 input XOR gate like below.



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



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

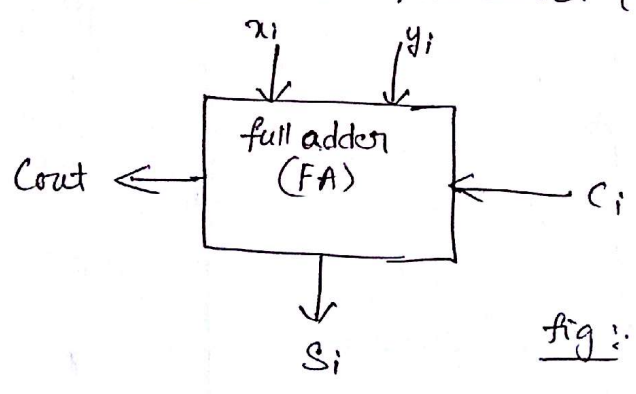


fig: logic for a single stage

A cascade connection of n full adder blocks 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.

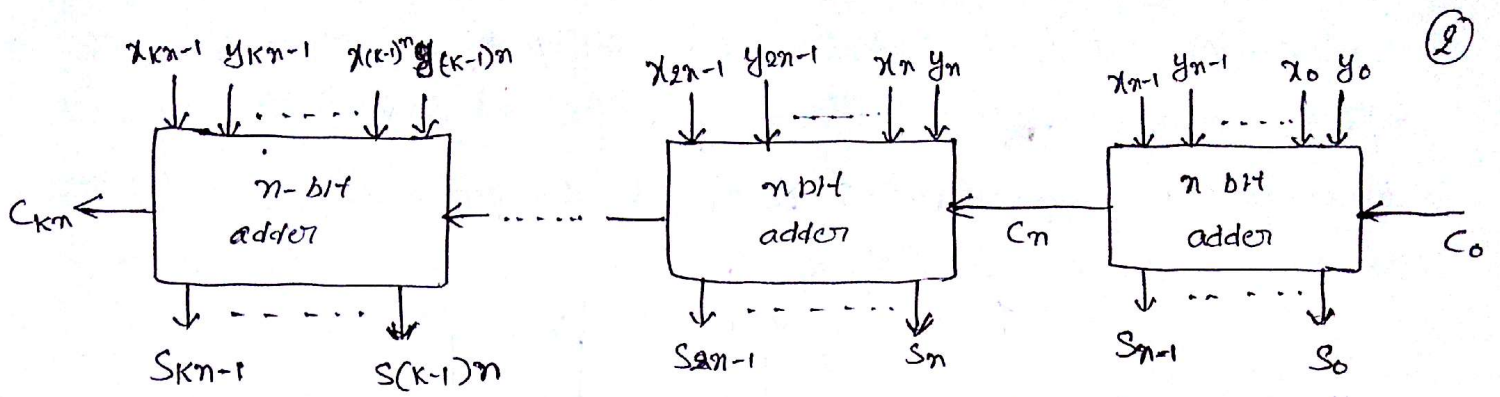


fig: cascade of k n-bit adder

Addition / Subtraction Logic unit :-

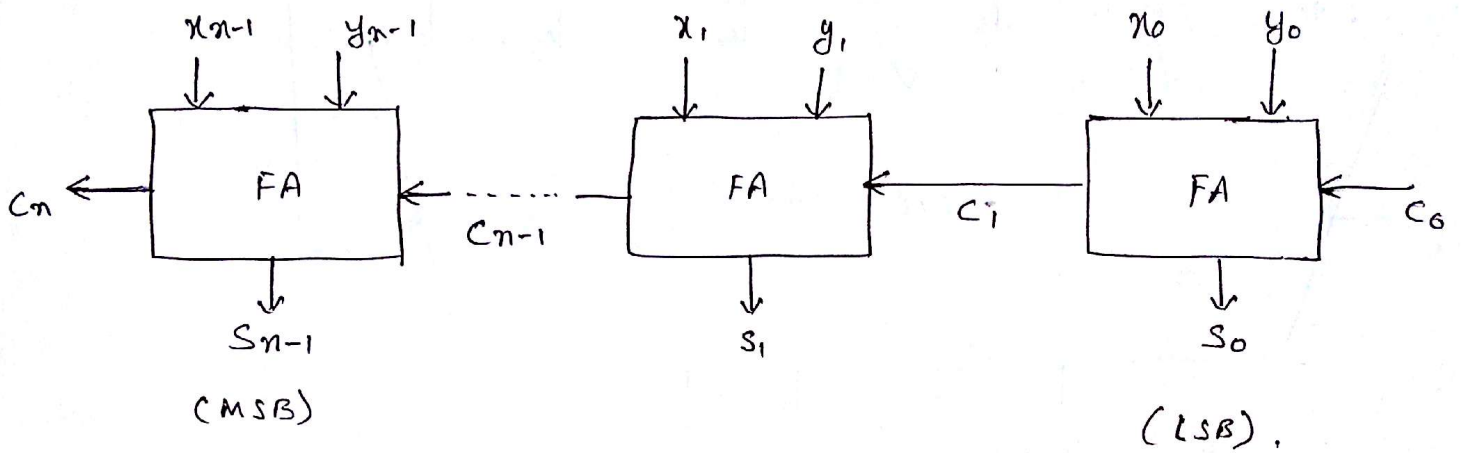


fig: An n-bit ripple - carry adder

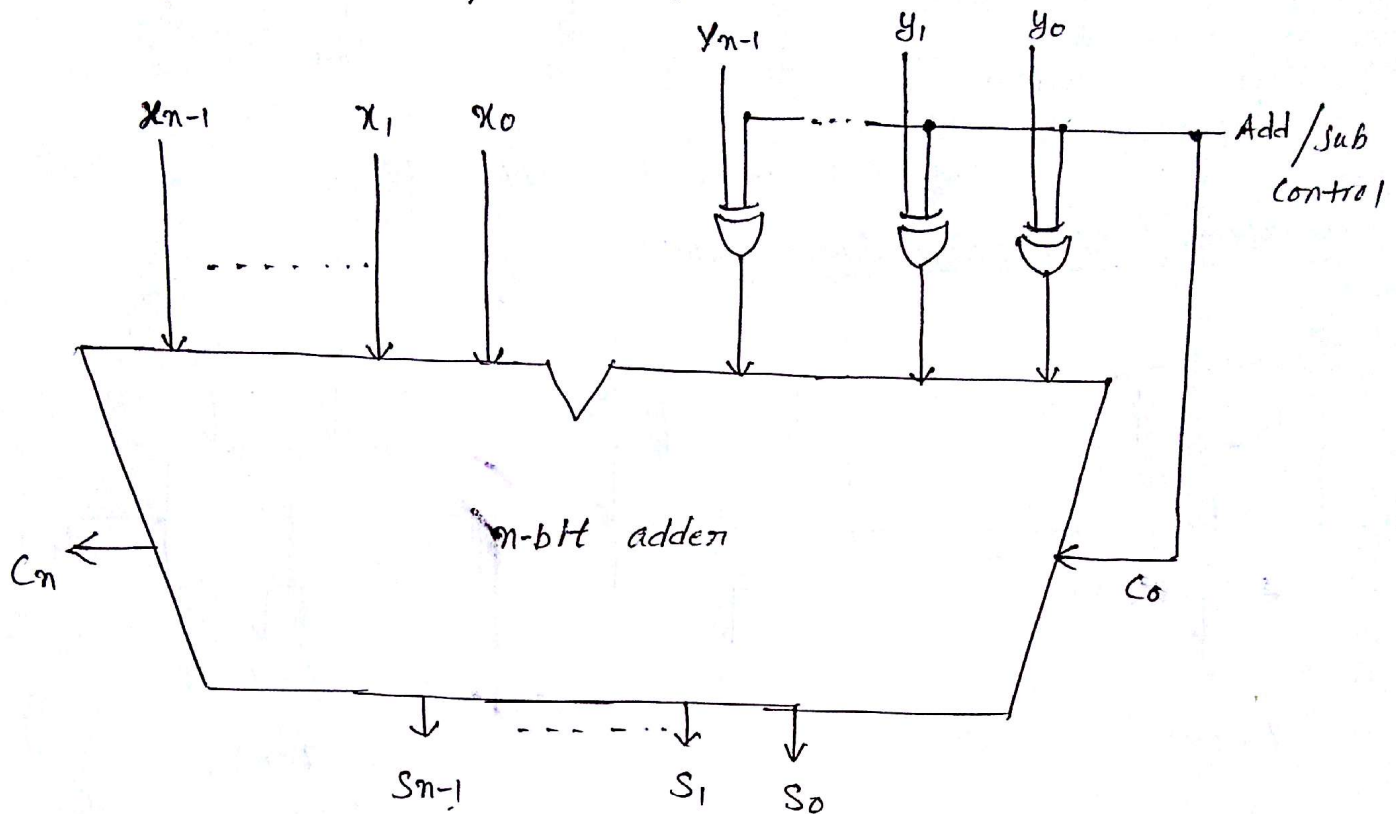
Addition / Subtraction Logic Unit :-

Overflow can only occur when the signs of the two operands are the same. An overflow may occur if the sign of the result is different. A circuit to detect overflow can be added to the n-bit adder by implementing the logic expression.

$$\text{Overflow} = x_{n-1}y_{n-1}\bar{s}_{n-1} + \bar{x}_{n-1}\bar{y}_{n-1}s_{n-1}$$

In order to perform subtraction operation $X - Y$, we take the 2's complement of Y and add it to X . The following circuit diagram performs the either addition or subtraction based on the value of add/sub input control line.

When the add/sub control line is set to 0 it will perform the addition operation. When the add/sub control line is set to 1 it will take the 1's complement of Y by the XOR gates and C_0 is set to 1 to get the 2's complement of Y .



DESIGN OF FAST ADDERS :-

If an n -bit adder ripple-carry adder is used in the addition/subtraction, it may take more time for producing outputs. Whether to accept this delay or not depends upon the context of the speed of other processor components and the data transfer times of registers and cache memories. The delay through a network of logic gates depends on the integrated circuit electronic technology.

Using the logic implementation in logic for a single stage C_{n-1} is available in $2(n-1)$ gate delays, and S_{n-1} is correct one XOR gate delay later. The final carry-out C_n is available after $2n$ gate delays.

Two approaches can be used to reduce delay in adders.

The first approach is to use the fastest possible electronic technology in implementing the ripple-carry logic design. The second approach is to use an augmented logic gate network structure that is larger than parallel adder which has seen previously.

Carry - Lookahead Addition:-

One method of speeding up the process by eliminating inter stage carry delay is called "lookahead-carry addition". It uses fast adder circuit to speed up. The logic expressions for S_i and C_{i+1} are

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

$$C_{i+1} = x_i y_i + x_i C_i + y_i C_i$$

$$C_{i+1} = x_i y_i + (x_i + y_i) C_i$$

We can write this equation as

$$C_{i+1} = G_i + P_i C_i$$

where $G_i = x_i y_i$ and $P_i = (x_i + y_i)$.

Here the expressions G_i and P_i are called "generate" and "propagate" of carry. The generate function produces on carry when both x_i and y_i are one, regardless of the input carry. Each stage (bit stage) contains an AND gate to form G_i , an OR gate to form P_i and three input XOR gate to form S_i .

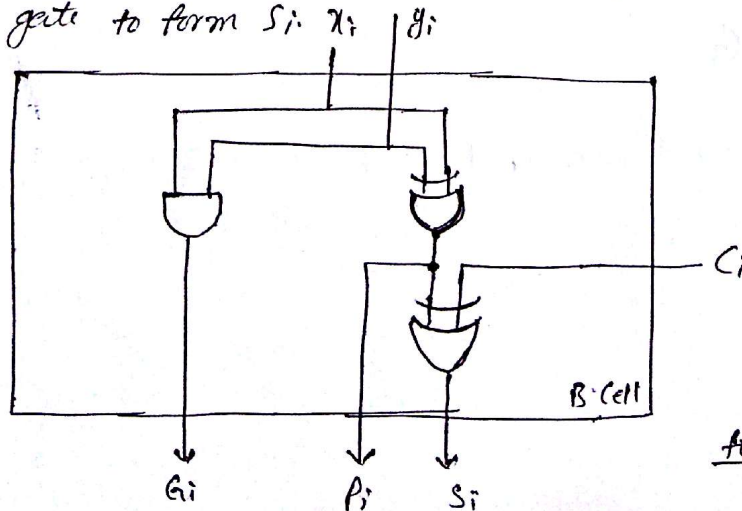


Fig: Bit-stage cell.

Now C_{i+1} can be expressed as a sum of products function of the P and G outputs of all the preceding stages.

$$C_1 = G_0 + P_0 C_m$$

$$C_2 = G_1 + P_1 C_1 = G_1 + P_1 G_0 + P_1 P_0 C_m$$

$$C_3 = G_2 + P_2 C_2 = G_2 + P_2 G_1 + P_2 P_1 G_0 + P_2 P_1 P_0 C_m$$

⋮

The following diagram shows the general form of a carry-lookahead adder circuit

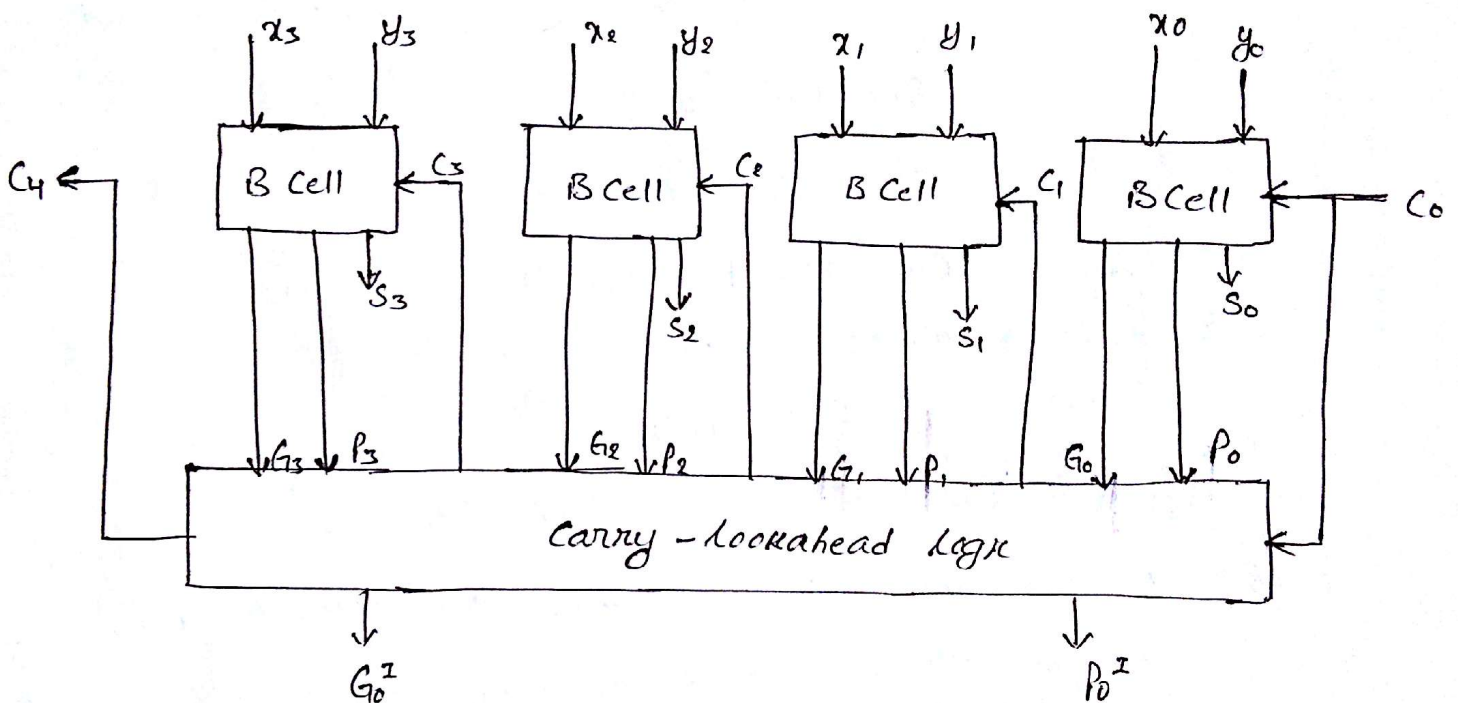


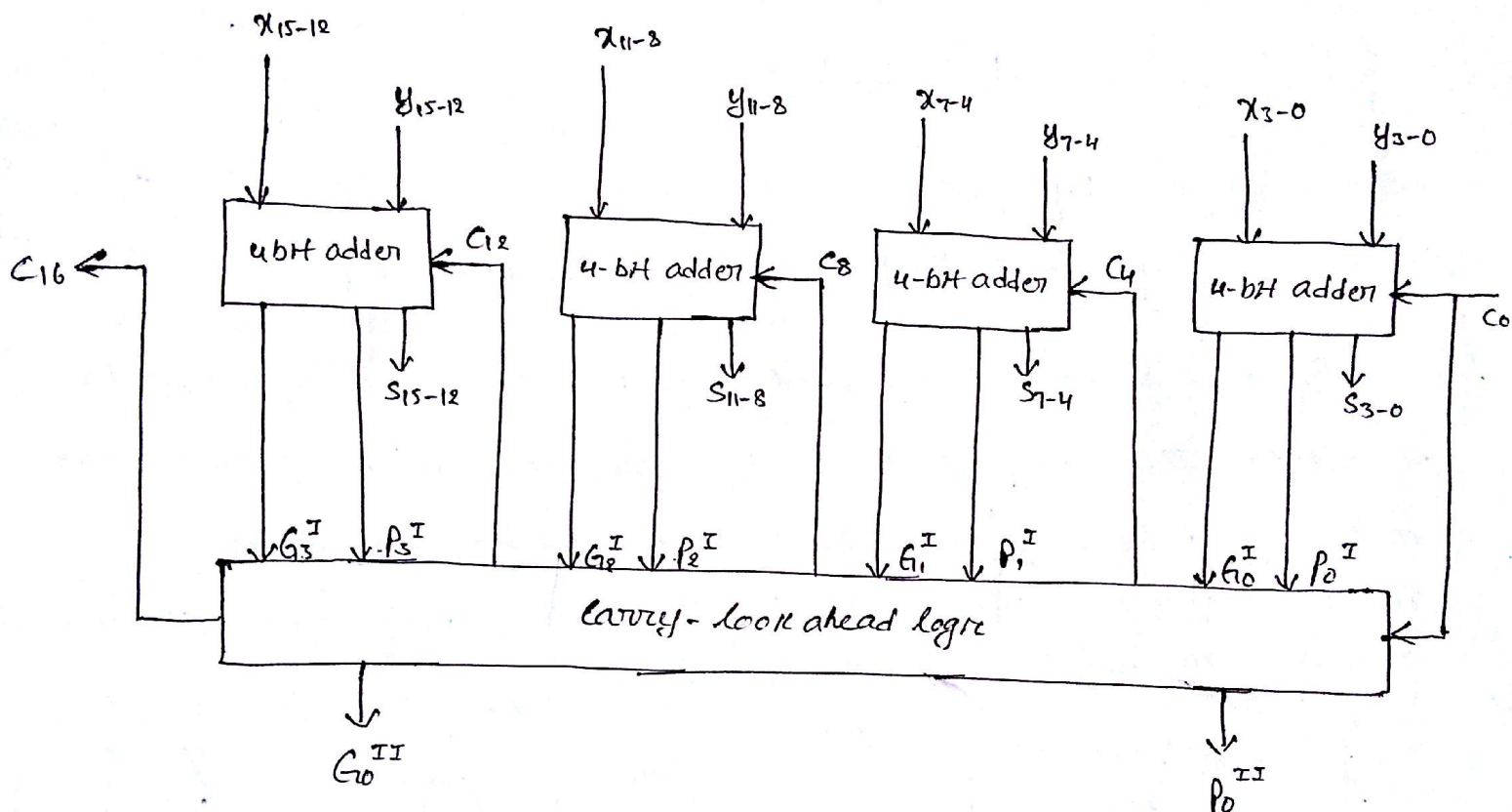
Fig: 4-bit carry-lookahead adder

$$C_4 = G_3 + P_3 G_3$$

$$= G_3 + P_3 G_2 + P_3 P_2 G_1 + P_3 P_2 P_1 G_0 + P_3 P_2 P_1 P_0 C_m$$

High level Generate and Propagate function:-

Previously we have seen only 1 bit adder block circuit. Same circuit we can implement for n-bit adder which is high level. By using high-level block generate and propagate functions, it is possible to use the lookahead approach to develop the carries C_4, C_8, C_{12}, \dots in parallel.



In first block $P_0^I = P_3 P_2 P_1 P_0$.

and $G_0^I = G_3 + P_3 G_2 + P_3 P_2 G_1 + P_3 P_2 P_1 G_0$

C_{16} will be formed like below by using carry-lookahead logic

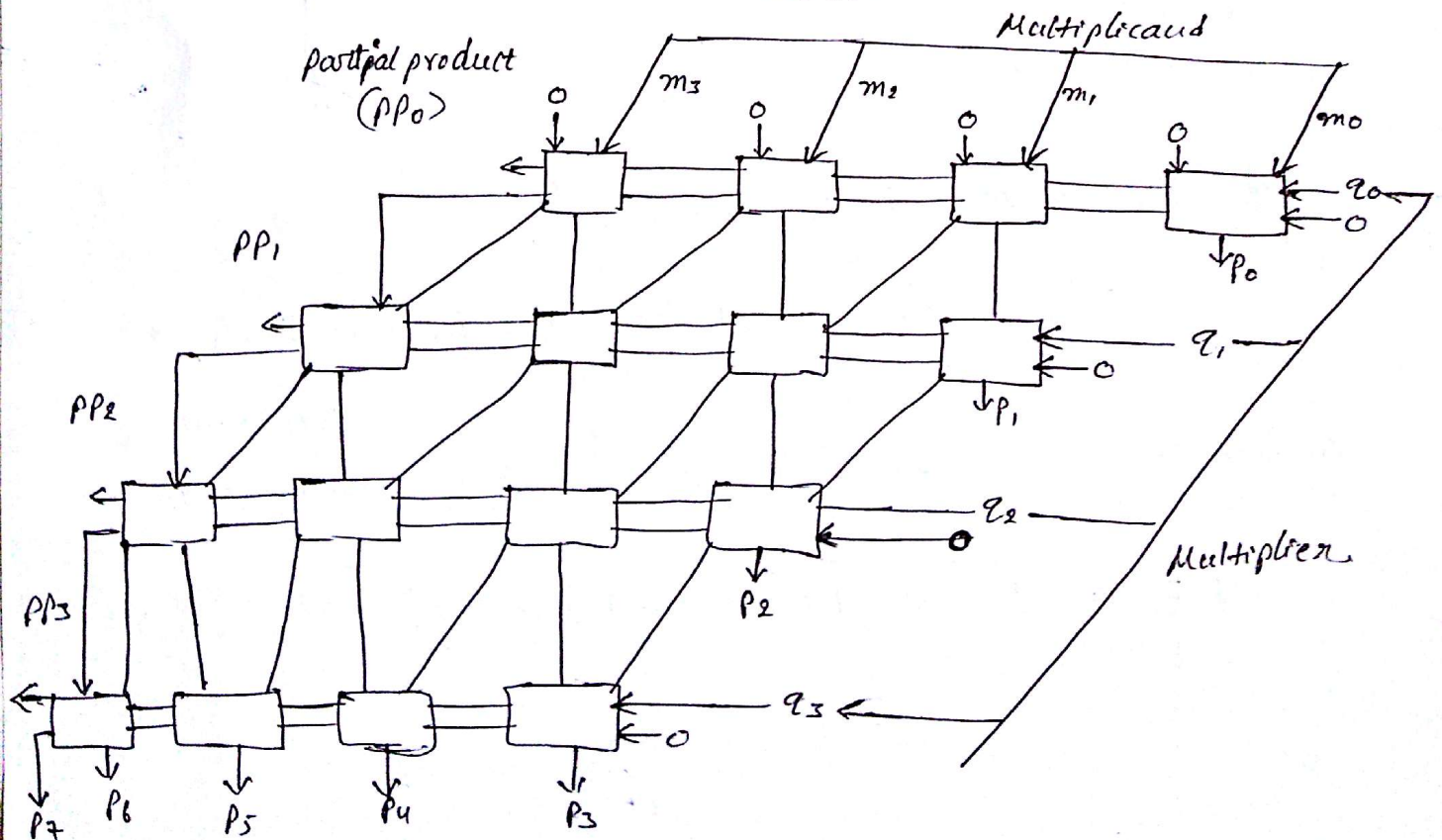
$$C_{16} = G_3^I + P_3^I G_2^I + P_3^I P_2^I G_1^I + P_3^I P_2^I P_1^I G_0^I + P_3^I P_2^I P_1^I P_0^I C_0$$

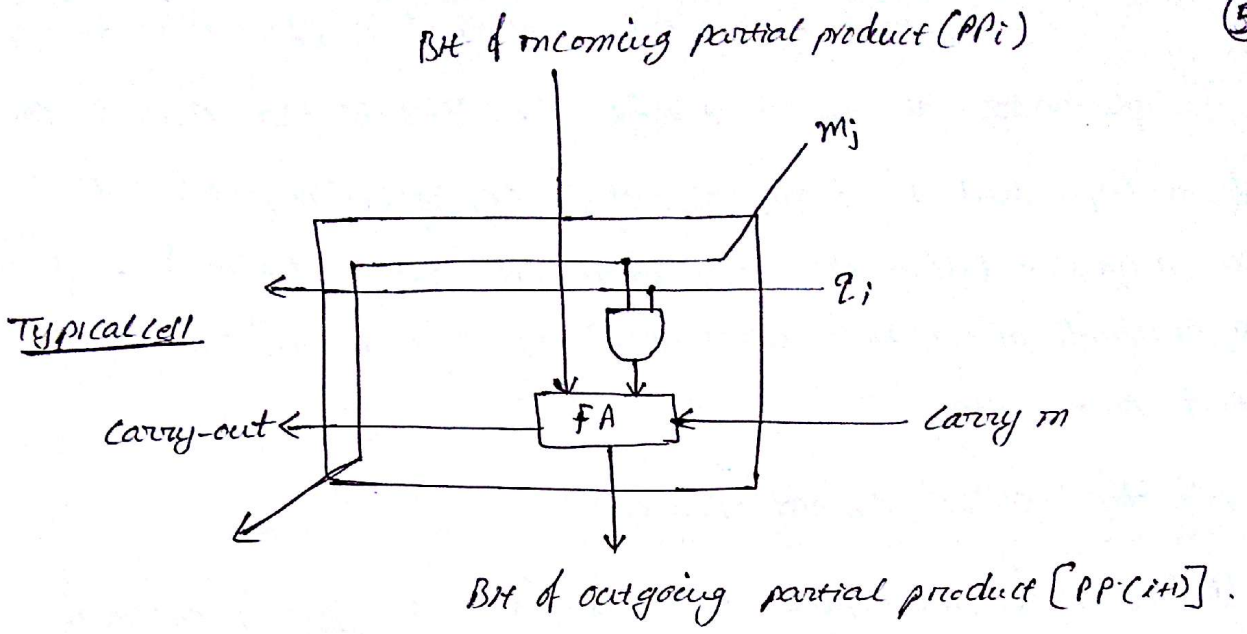
Multiplication of Positive Numbers:-

The usual algorithm for multiplying integers by hand is like below.

$$\begin{array}{r}
 1101 \quad (13) \text{ Multiplicand } M \\
 1011 \quad (11) \text{ Multiplier } Q \\
 \hline
 1101 \\
 1101 \\
 0000 \\
 1101 \\
 \hline
 10001111 \quad (143) \text{ Product } P.
 \end{array}$$

This algorithm can be applied for unsigned numbers and positive numbers. The product of two n -digit numbers can be accommodated in $2n$ digits. Binary multiplication of positive operands can be implemented in a combinational two-dimensional logic array. The main component in each cell is a full adder FA.





The simplest way to perform multiplication is to use the adder circuitry in the ALU for a number of sequential steps. The following diagram shows the hardware arrangement for sequential multiplication.

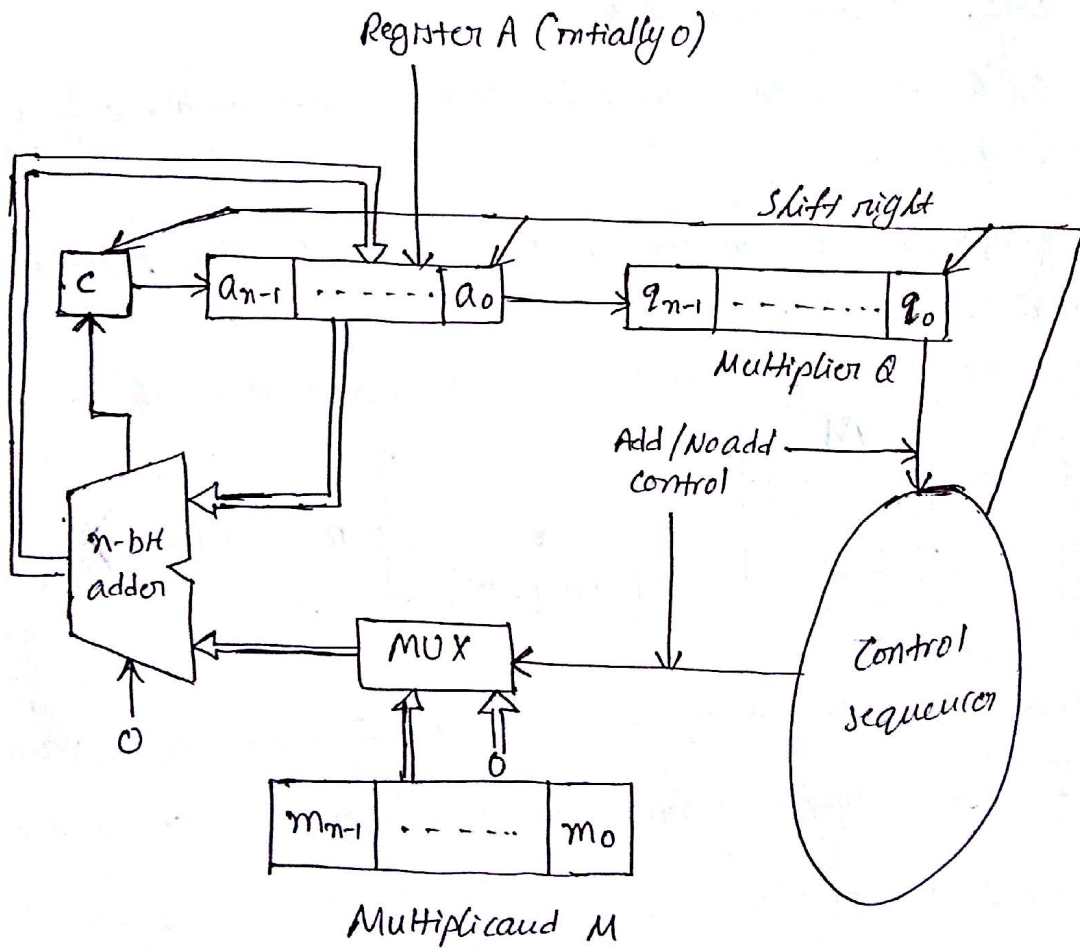


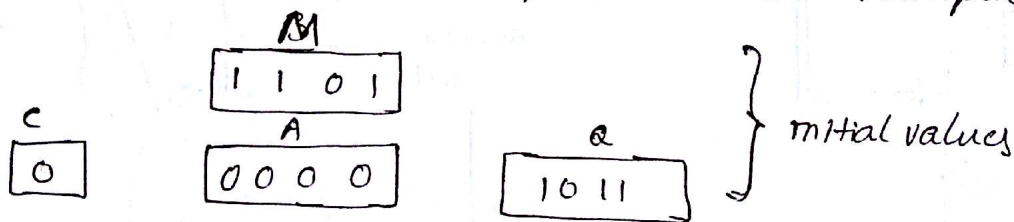
Fig: Register configuration

Registers A and Q combined hold PP; while multiplier Q₁ generates the signal Add/NoAdd. This signal controls the addition of multiplicand M to PP; to generate PP(i+1). The product is computed in n cycles. Using the sequential hardware structure it is clear that a multiplication instruction takes much more time to execute than an Add instruction.

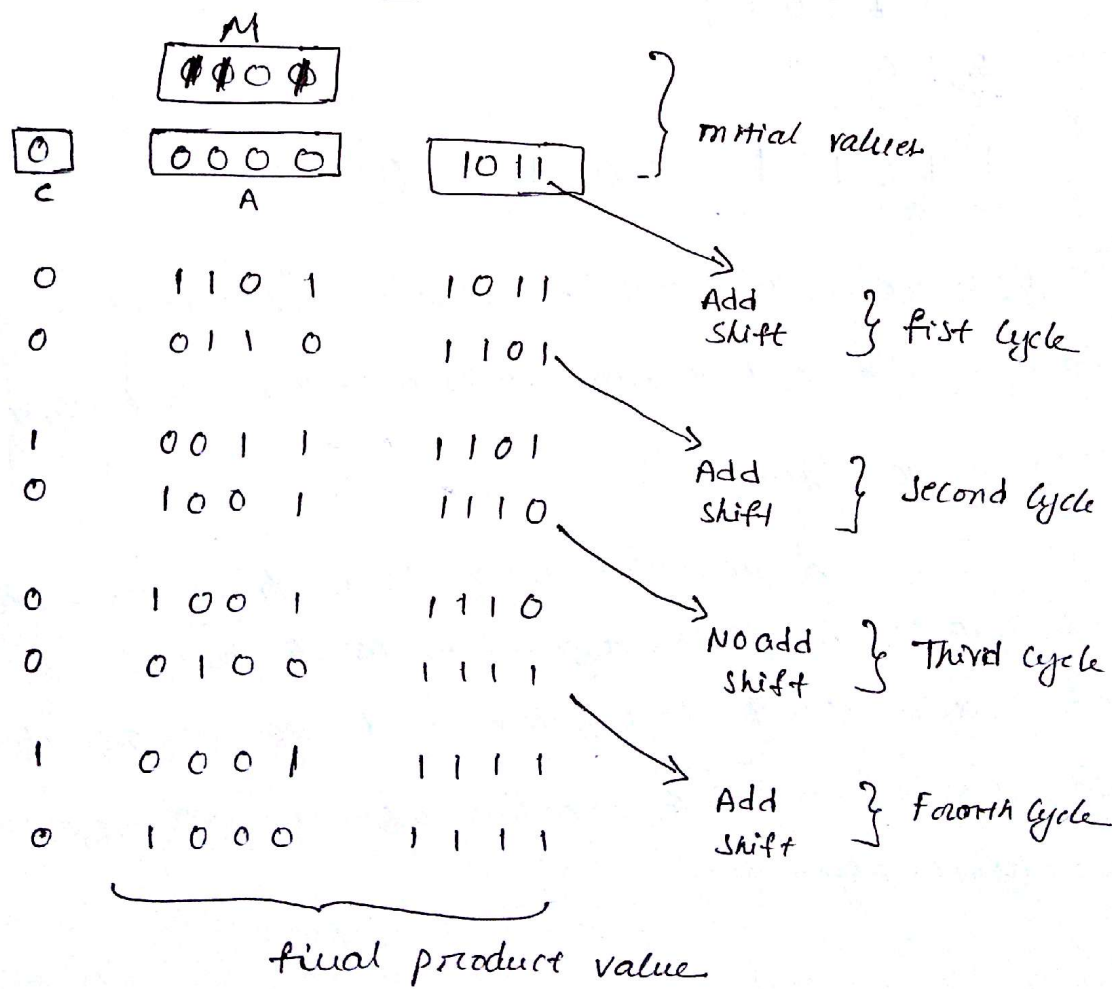
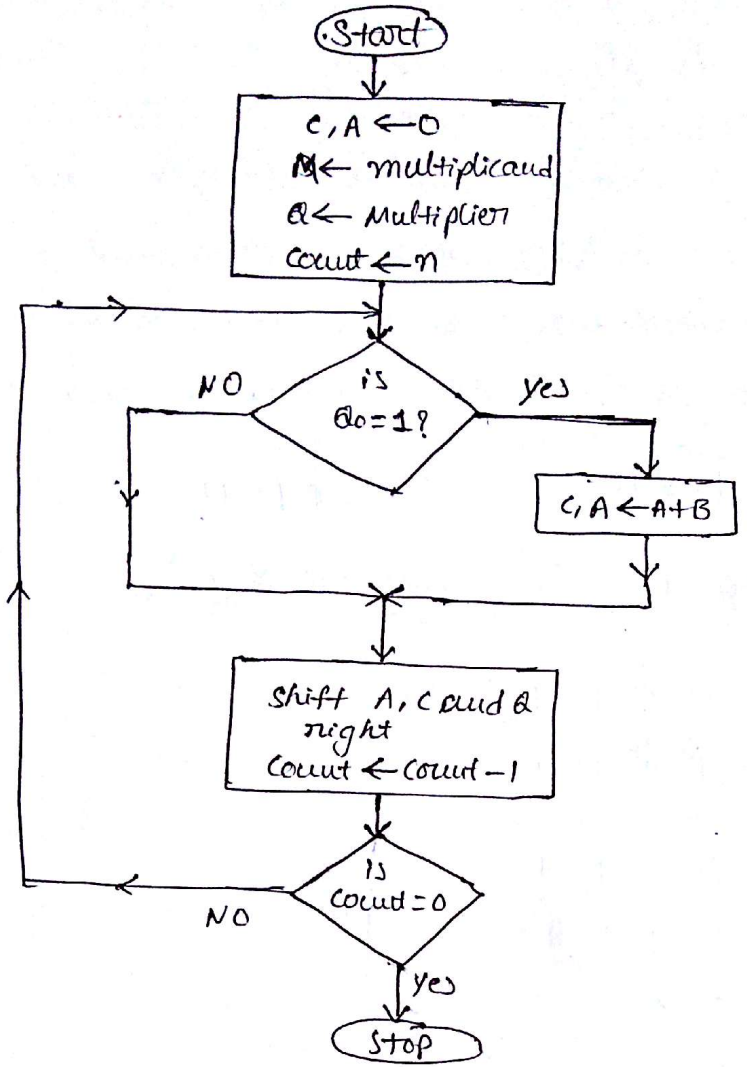
Multiplication operation steps :-

- ① Bit 0 of multiplier operand (Q₀ of Q register) checked.
- ② If bit 0 (Q₀) is 1 then multiplicand and partial product are added and all bits of C, A and Q register are shifted to the right one bit position. So the C bit goes into A_{n-1}, A₀ goes into Q_{n-1} and Q₀ is lost. C bit cleared to zero.
- ③ If bit 0 (Q₀) is 0, then no addition is performed, only shift operation is carried out.
- ④ Repeat step 1, 2, 3 n times to get the desired result in the A and Q registers.

⇒ Consider the example Multiplicand = 1101 Multiplier = 1011



Initially the A register and C (carry bit) values sets to zero. The following flowchart indicates the sequential multiplication operation



SIGNED - OPERAND MULTIPLICATION:-

Consider the case of a positive multiplier and a negative multiplicand. when we add a negative multiplicand to a particular partial product, we must extend the sign-bit value of the multiplicand to the left as far as the product will extend. we ~~now~~ consider the following example -13×11

$$-13 = 10011 \quad (\text{2's complement of } 13).$$

$$\begin{array}{r} 10011 \quad (-13) \\ 01011 \quad (11) \\ \hline 111110011 \\ 111110011 \\ 000000000 \\ 1100011 \\ 0000000 \\ \hline 1101110001 \quad (-145) \end{array}$$

Now we will discuss a technique which works equally well for both negative and positive multiplier called Booth algorithm

Booth Algorithm:-

A powerful algorithm for signed-number multiplication is Booth's algorithm, which generates a $2n$ -bit product and treats both positive and negative numbers uniformly. we can use a concept of recoding when we get negative numbers in Booth Algorithm.

the general Booth's algorithm recoding scheme can be given as

Multiplier -1 times the shifted multiplicand is selected when moving from 0 to 1, $+1$ times the shifted multiplicand is selected when moving from 1 to 0, and 0 times the shifted multiplicand is selected for none of the above case as multiplier is scanned from right to left.

Multiplier		version of multiplicand selected by b_i
b_i	b_{i-1}	
0	0	$0 \times M$
0	1	$+1 \times M$
1	0	$-1 \times M$
1	1	$0 \times M$

Fig: Booth multiplier recoding table

EX: ① Recode the multiplier 101100

Sol: 1 0 1 1 0 0 0 \leftarrow implied zero Multiplier

-1 +1 0 -1 0 0 Recoded Multiplier

\rightarrow 0 1 0 1 0 1 0 1 0 1 0 1 0 1
 +1 -1 +1 -1 +1 -1 +1 -1 +1 -1 +1 -1
 } worst-case multiplier

\rightarrow 1 1 0 0 0 1 0 1 1 0 1 1 1 0 0
 0 -1 0 0 +1 -1 +1 0 -1 +1 0 0 0 -1 0 0
 } ordinary multiplier

\rightarrow 0 0 0 0 1 1 1 1 0 0 0 0 1 1 1
 0 0 0 +1 0 0 0 0 -1 0 0 0 +1 0 0 -1
 } Good multiplier

⇒ Multiply 01110 (+14) and 11011 (-5).

Sol:

recode the multiplier -5

1 1 0 1 1

0 - 1 1 0 - 1 → Recoded multiplier.

						0	1	1	1	0
						0	-	1	0	-
1	1	1	1	1	1	0	0	1	0	
0	0	0	0	0	0	0	0	0	0	
0	0	0	0	1	1	1	0			
1	1	1	0	0	1	0				
0	0	0	0	0	0					
1	1	1	0	1	1	1	0	1	0	(-70)

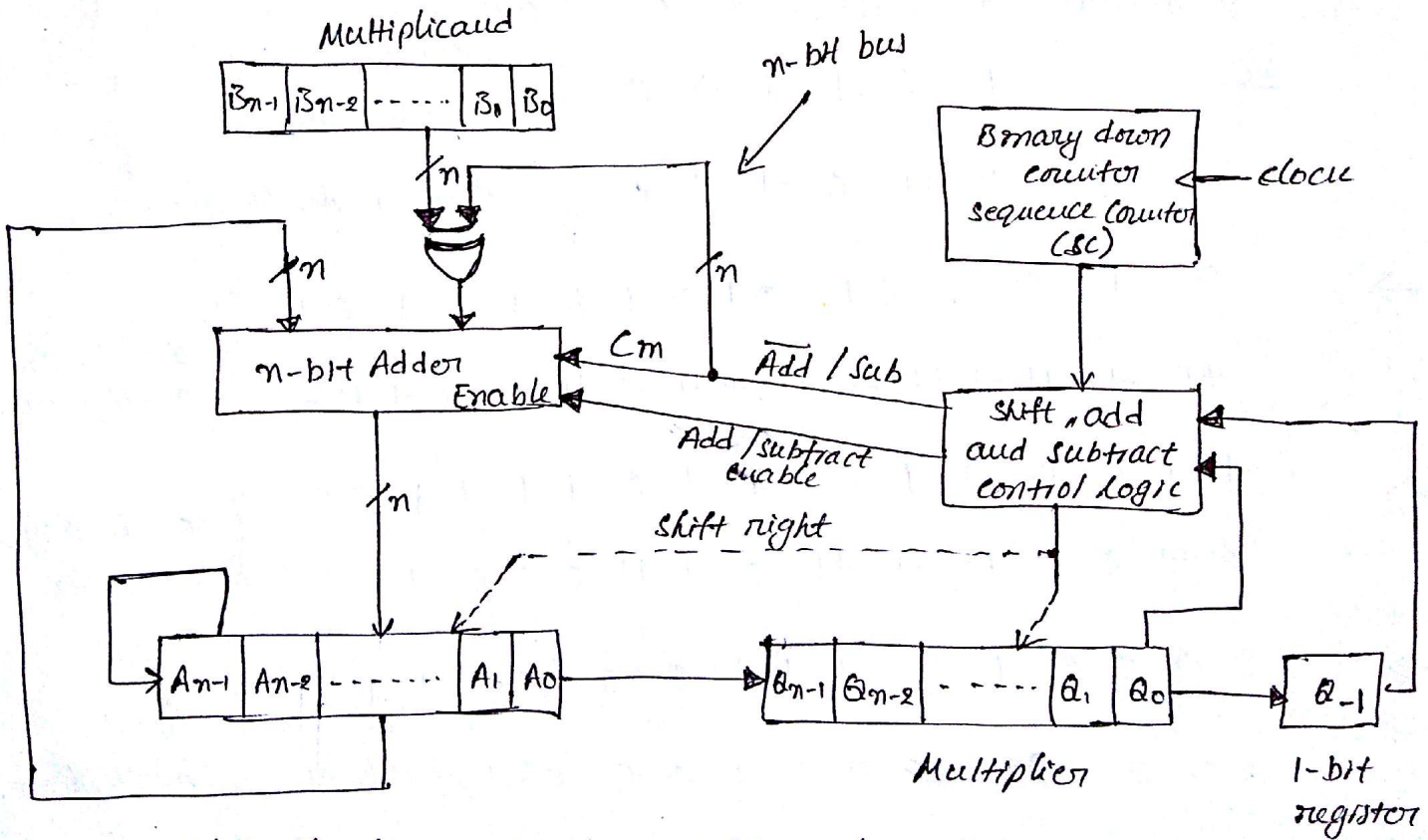


Fig: Hardware implementation of signed binary multiplication
 → initial settings: $A \leftarrow 0$ and $Q_{-1} \leftarrow 0$

Q_0	Q_{-1}	Add/sub	Add/subtract enable	shift
0	0	x	0	1
0	1	0	1	1
1	0	1	1	1
1	1	x	0	1

(8)

Fig: truth table for shift, add, and subtract control logic.

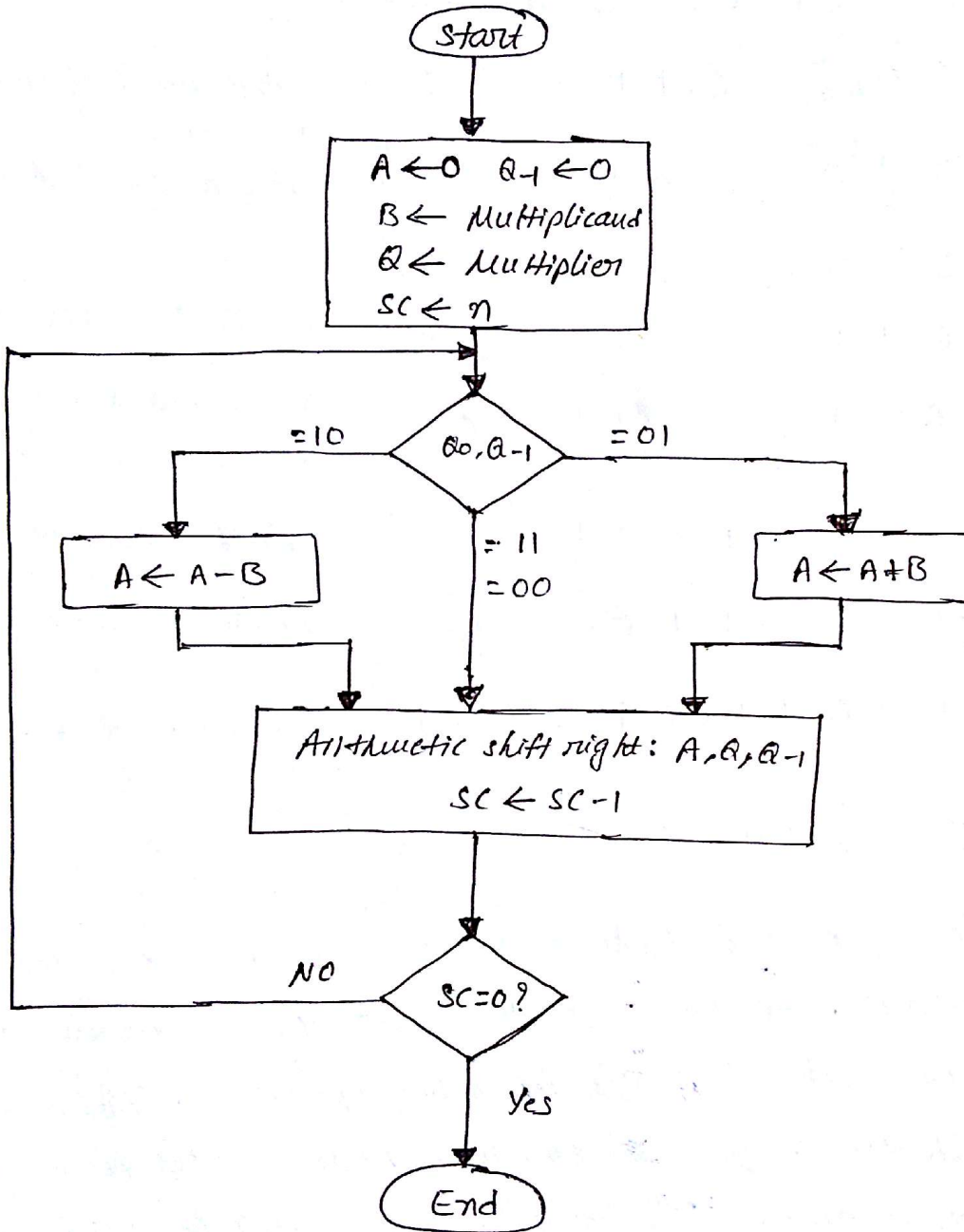


Fig: Booth's algorithm for signed multiplication.

Example! Multiply 5×-4

Multiplicand (B) = 0101 (5), Multiplier (Q) = 1100 (-4)

SC	A	Q	Q-1	
100	0000	1100	0	→ initial values
011	0000	0110	0	Arithmetic shift right
010	0000	0011	0	Arithmetic shift right
	0000			
	1011			$Q_0 Q_{-1} = 10$
	1011	0011	0	$\therefore A \leftarrow A - B$
001	1101	1001	1	Arithmetic shift right
000	1110	1100	1	Arithmetic shift right

Result: 111011001 = (-20) (2's complement of 20).

FAST MULTIPLICATION :-

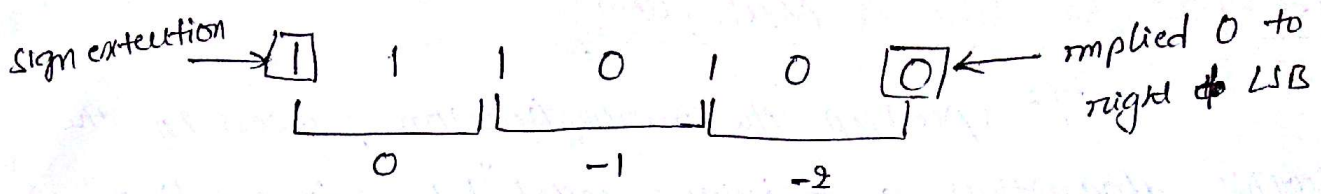
There are two techniques for speeding up the multiplication process. In the first technique the maximum number of summands are reduced to $n/2$ for n -bit operands. The second technique called the carry save addition reduces the time needed to add the summand. The first approach will be called as BH-pair recording of multipliers.

Bit-Pair Recoding of Multipliers:-

To speedup the multiplication process in the Booth's algorithm a technique called "bit-pair recoding" is used. It is also called "modified Booth's algorithm". It decreases the no. of summands. In this technique, the Booth-recoded multiplier bits are grouped in pairs. Then each pair is represented by its equivalent single bit multiplier reducing total no. of multiplier bits to half. In this technique for a n -bit multiplier we will get $n/2$ summands. The pair $(+1, -1)$ is equivalent to the pair $(0 +1)$, similarly $(+1, 0)$ is equivalent to $(0 +2)$. The following table indicates the bit-pair code for all possible multiplier bit options.

Multiplier bit-pair		Multiplier bit on the right	BH pair recoded multiplier bit at position i
$i+1$	i		
0	0	0	0
0	0	1	+1
0	1	0	+1
0	1	1	+2
1	0	0	-2
1	0	1	-1
1	1	0	-1
1	1	1	0

Ex: Find the bit-pair code for multiplier 11010



Ex: Multiply given signed 2's complement numbers using bit pair recoding

$$A = 110101$$

$$B = 011011$$

$$\begin{array}{r}
 110101 \\
 +2 \quad -1 \quad -1 \\
 \hline
 000000001011 \\
 00000001011 \\
 11101010 \\
 \hline
 111011010111 \quad (-297)
 \end{array}$$

Carry - Save Addition Of Summands :-

Multiplication requires the addition of several summands. A technique called carry-save addition (CSA) speeds up the addition process.

Example: Multiply 01101 (+13) and 11010 (-6) using Booth's recoding and bit-pair recoding techniques

Sol:- recoding of multiplier 11010 in Booth's technique is

$$0 -1 +1 -1 0$$

$$\begin{array}{cccccc}
 & & & 0 & 1 & 1 & 0 & 1 & & \\
 & & & 0 & -1 & +1 & -1 & 0 & & \\
 \hline
 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
 1 & 1 & 1 & 1 & 1 & 0 & 0 & 1 & 1 & \\
 0 & 0 & 0 & 0 & 1 & 1 & 0 & 1 & & \\
 1 & 1 & 1 & 0 & 0 & 1 & 1 & & & \\
 0 & 0 & 0 & 0 & 0 & 0 & & & & \\
 \hline
 1 & 1 & 1 & 0 & 1 & 1 & 0 & 0 & 1 & 0 \quad (-78)
 \end{array}$$

$$\begin{array}{cccccc}
 & & & 0 & 1 & 1 & 0 & 1 & & \\
 & & & 0 & & -1 & & -2 & & \\
 \hline
 1 & 1 & 1 & 1 & 1 & 0 & 0 & 1 & 1 & 0 \\
 1 & 1 & 1 & 1 & 0 & 0 & 1 & 1 & & \\
 0 & 0 & 0 & 0 & 0 & 0 & & & & \\
 \hline
 1 & 1 & 1 & 0 & 1 & 1 & 0 & 0 & 1 & 0 \quad (-78)
 \end{array}$$

while multiplying with -2, take the 2's complement for 26 as the multiplicand is 13

$$13 \times -2 = -26$$

Z.

The following diagrams are the circuitary diagrams of Ripple-carry and carry save arrays for the multiplication operation $M \times Q = P$ for 4-bit operands

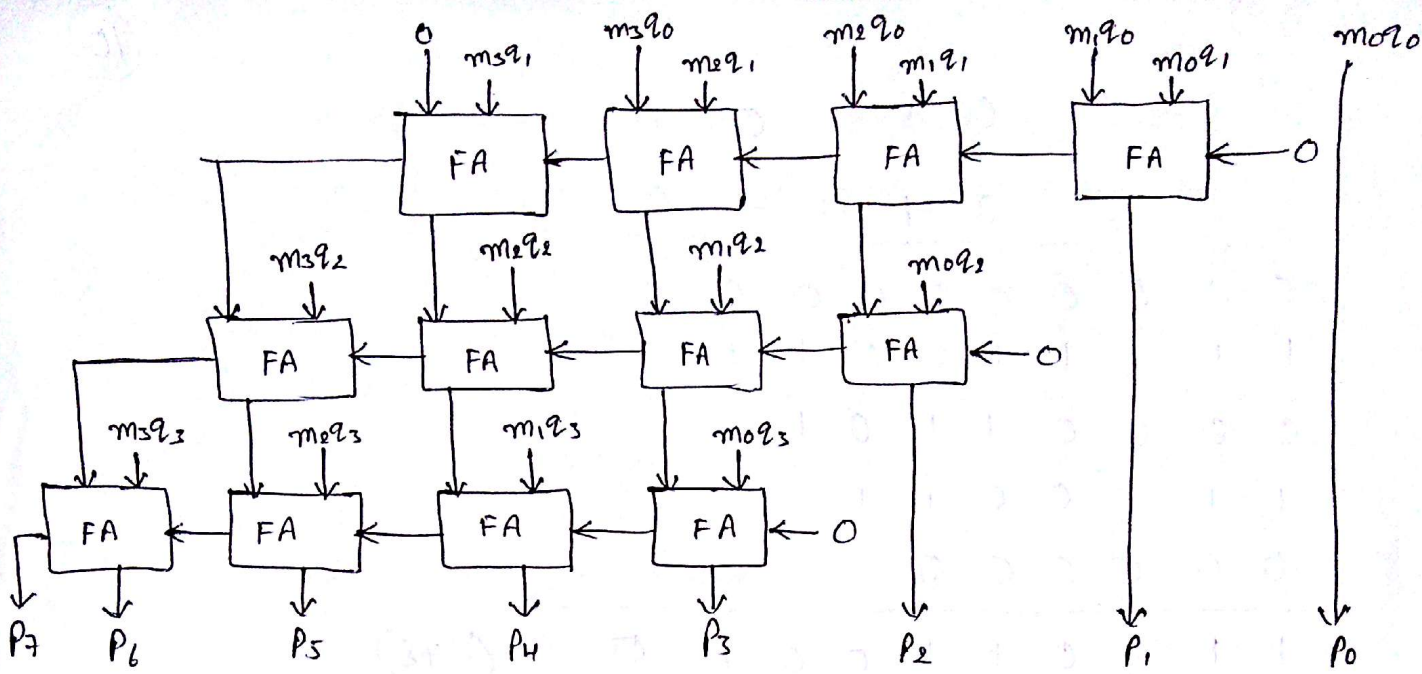


Fig: ripple carry array

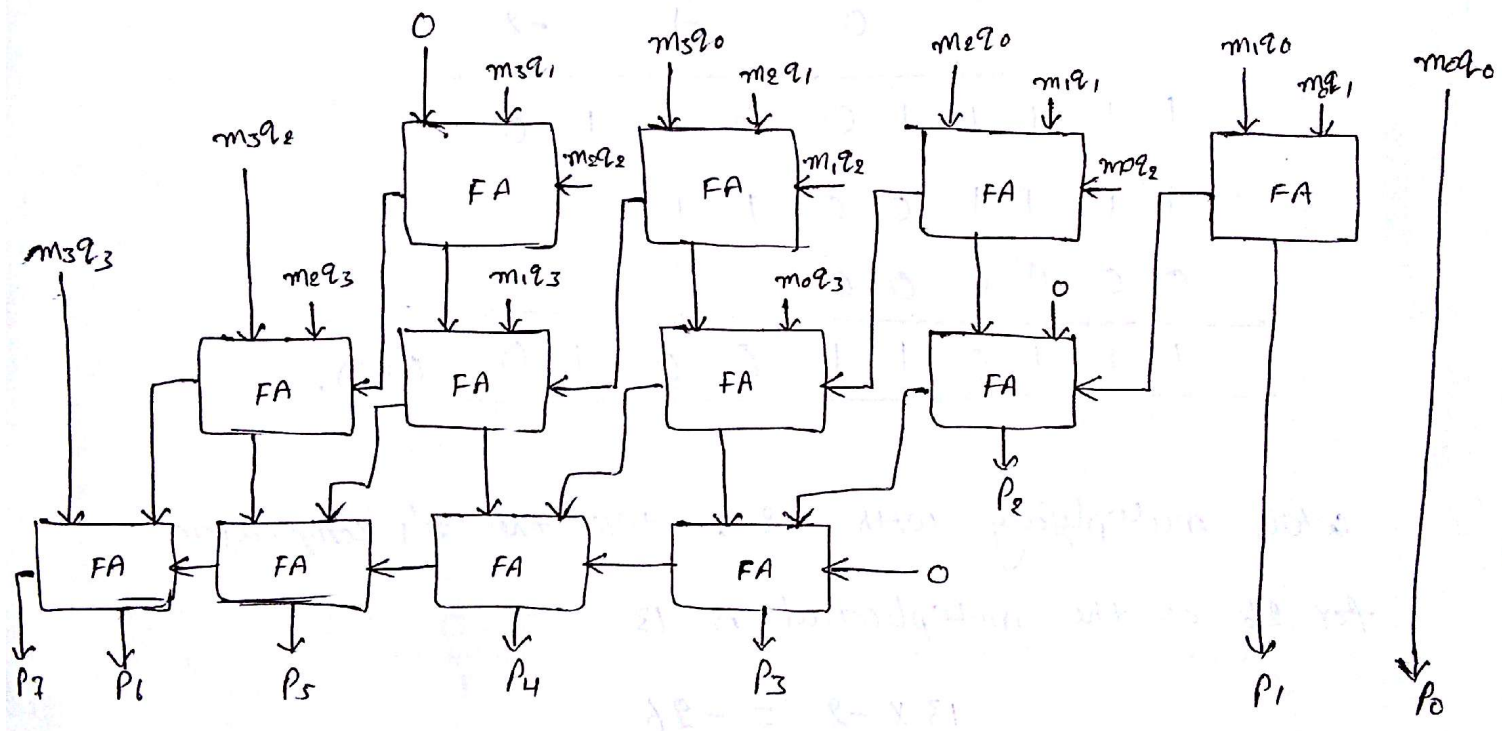


Fig: carry save array

Instead of letting the carries ripple along the rows, they can be saved and introduced into the next row, at the correct weighted positions. This frees up an input to three full adders in the first row. Delay through the carry ^{save} array is somewhat less than delay through the ripple carry array.

Consider a situation where addition of many summands required in the multiplication of longer operands. We can group the summands in threes and perform carry save addition on each of these groups in parallel to generate a set of S and C vectors in one full-adder delay. Next we group all of the S and C vectors into threes, and perform carry-save addition on them, generating a further set of S and C vectors in one more full-adder delay. We continue with this process until there are only two vectors remaining. Then they can be added in a ripple carry or a carry lookahead adder to produce the desired product.

	1	0	1	1	0	1	(15)	M
x	1	1	1	1	1	1	(63)	Q
<hr/>								
	1	0	1	1	0	1	A	
	10	1	1	0	1		B	
	101	1	0	1			C	
	1011	0	1				D	
	10110	1					E	
	101101						F	
<hr/>								
	1011000	1	00	11			(2835)	product.

Same multiplication can be done in carry-save addition as follows:

$$\begin{array}{r}
 101101 \quad M \\
 \underline{111111 \quad Q} \\
 \hline
 \end{array}$$

$$\begin{array}{r}
 101101 \quad A \\
 101101 \quad B \\
 \underline{101101 \quad C} \\
 \hline
 11000011 \quad S_1 \\
 00111100 \quad C_1
 \end{array}$$

$$\begin{array}{r}
 101101 \quad D \\
 101101 \quad E \\
 \underline{101101 \quad F} \\
 \hline
 11000011 \quad S_2 \\
 00111100 \quad C_2
 \end{array}$$

$$\begin{array}{r}
 11000011 \quad S_1 + C_1 + S_2 \\
 00111100 \\
 \underline{11100001} \\
 \hline
 \end{array}$$

$$\begin{array}{r}
 11010100011 \quad S_3 \\
 00001011000 \quad C_3 + C_2 \\
 00111100 \\
 \hline
 \end{array}$$

$$\begin{array}{r}
 010111010011 \quad S_4 \\
 01010100000 \quad C_4 \\
 \hline
 \end{array}$$

$$101100010011 \quad \text{product}$$

INTEGER DIVISION :-

The division operation is more complex than multiplication. For simplicity we will see division for positive numbers only. The division operation will be done manually like below.

$$\begin{array}{r}
 \text{divisor} \leftarrow 12 \leftarrow \text{dividend} \\
 12 \overline{) 169} \leftarrow \text{Quotient} \\
 \underline{12} \\
 49 \leftarrow \text{partial remainder} \\
 \underline{48} \\
 1 \leftarrow \text{remainder value}
 \end{array}$$

same in binary

$$\begin{array}{r}
 1100 \overline{) 10101001} \leftarrow \text{remainder} \\
 \underline{1100} \\
 10010 \\
 \underline{1100} \\
 01100 \\
 \underline{1100} \\
 00001
 \end{array}$$

The following diagram indicates the binary division.

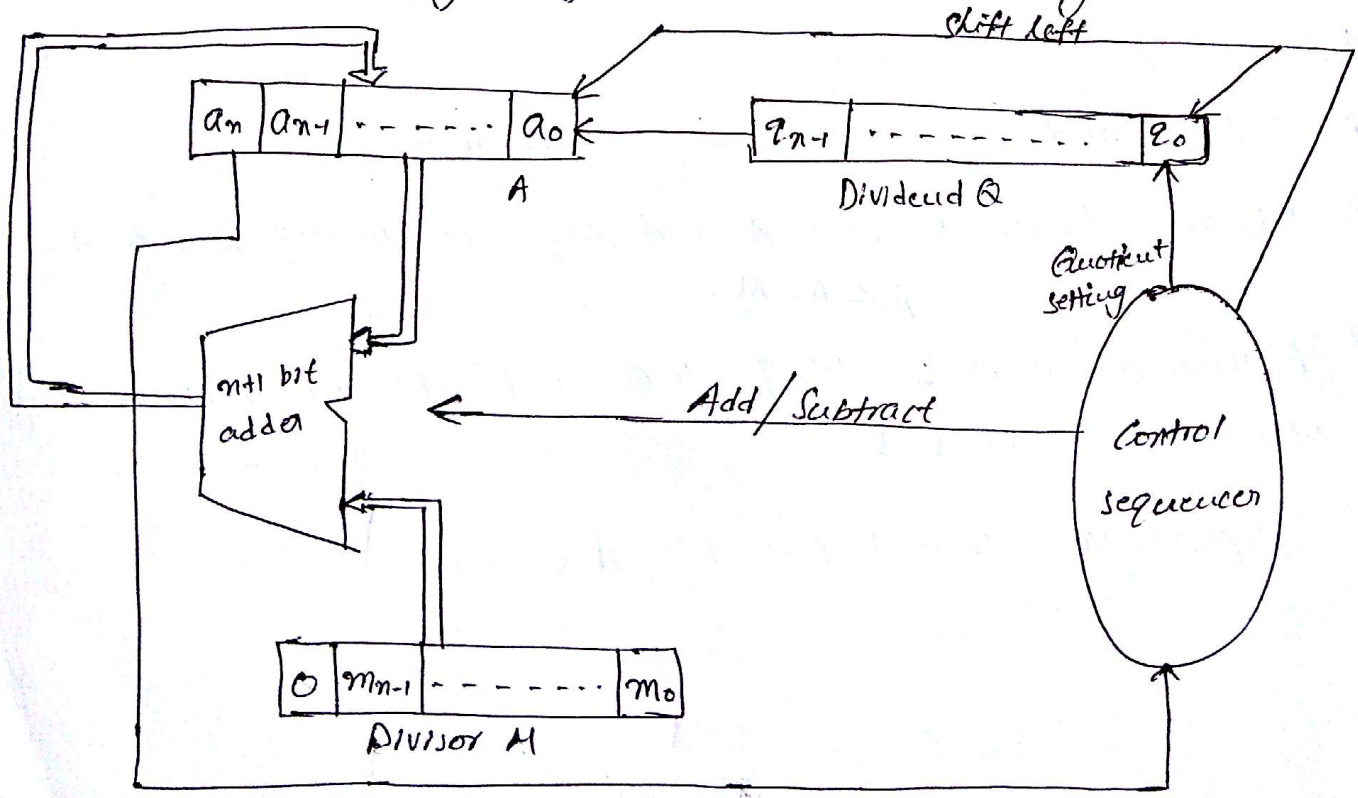


Fig: Circuit arrangement for binary division

In both the decisions, division process is same, only in binary division quotient bits are 0 and 1. We will now see binary division in detail. First the bits of the dividend are examined from left to right, whether it is equal or greater than divisor. Until the condition occurs 0's are replaced in the quotient from left to right. When the condition satisfied, a 1 is placed in the quotient and the divisor is subtracted from the partial product dividend. The process continues until all the bits of the dividend are brought down and result is still less than the divisor.

Restoring division:-

The divisor ~~and~~ which is n -bit positive number is loaded into register M and an n -bit positive dividend is loaded into register Q at the start of division operation. Register A is initially set to zero. The division operation is then carried out. After the division is complete, the n -bit quotient is in register Q and the remainder is in register A .

division operation steps:-

- ① shift A and Q left one binary position.
- ② Subtract divisor M from A and place the answer back in A .
 $A \leftarrow A - M$.
- ③ If the sign of A is 1, set Q_0 to 0 and add divisor M back to A .
otherwise set Q_0 to 1
- ④ Repeat the steps 1, 2 and 3 n times.

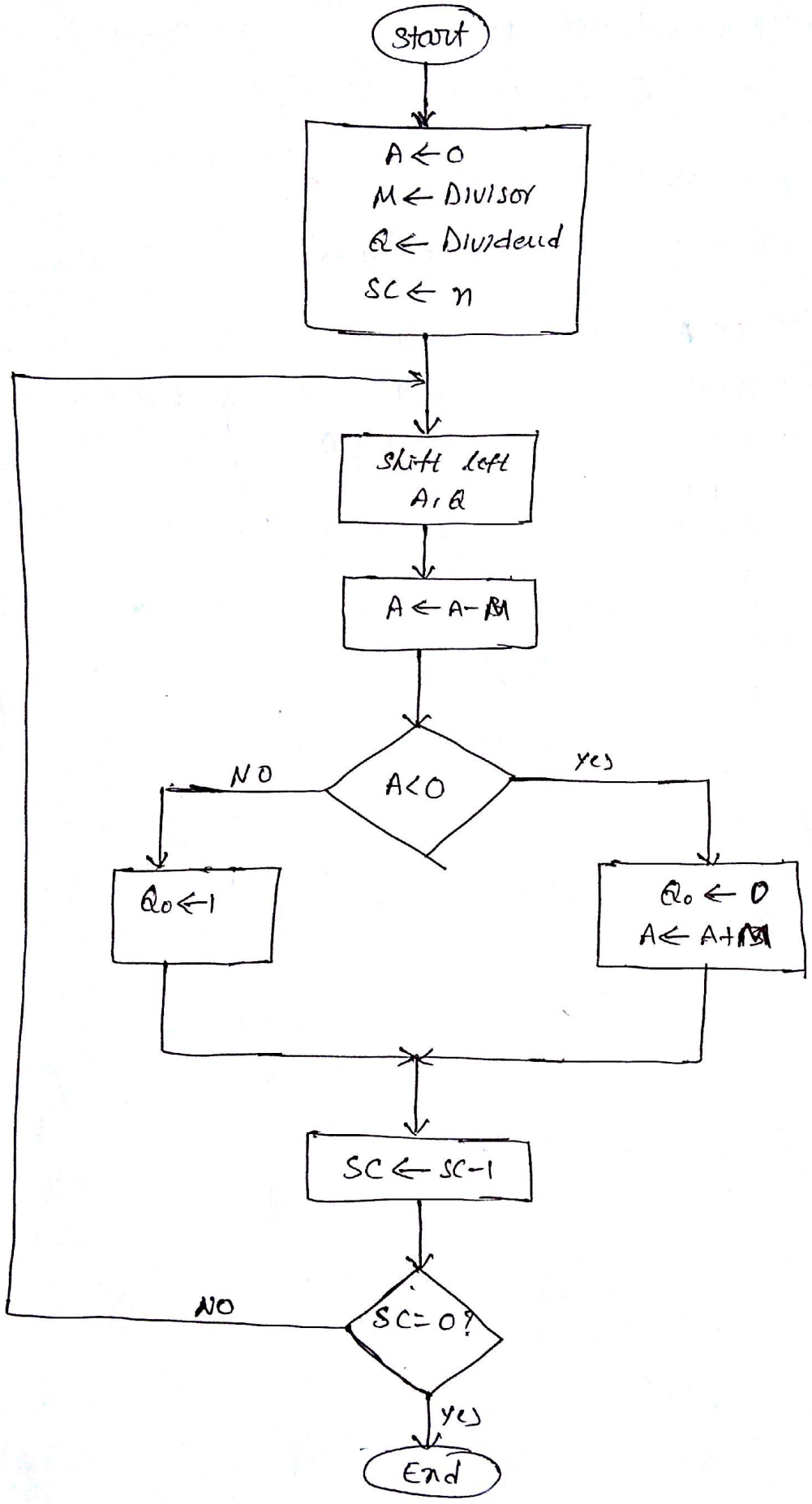


Fig: flow chart for restoring division operation

Example:- Consider 4-bit dividend and 2-bit divisor

Dividend = 1 0 0 0

Divisor = 0 0 1 1

	E	A register	Q register	
initially	0	0000	1000	
shift	0	0000	000	
subtract M	0	1110	000	} 1 st cycle
Set Q ₀ to 0	0	1110	0000	
Restore	0	0001	0000	
shift	0	0010	000	
subtract M	1	1110	000	

Example: Consider 4-bit dividend and 2-bit divisor.

Dividend = 1 0 0 0 @

Divisor = 0 0 1 1 M

		A		Q	
Initially	0	0 0 0 0		1 0 0 0	
Shift	0	0 0 0 1		0 0 0	
Subtract M	1	1 1 0 1			} 1st Cycle
Set Q ₀	1	1 1 1 0			
Restore	1	1 1 1 0		0 0 0 0	
		1 1		0 0 0 0	
Shift	0	0 0 1 0		0 0 0	
Subtract M	1	1 1 0 1			} 2nd Cycle
Set Q ₀	1	1 1 1 1			
Restore	1	1 1 1 1		0 0 0 0	
		1 1		0 0 0 0	
Shift	0	0 1 0 0		0 0 0	
Subtract M	1	1 1 0 1			} 3rd Cycle
Set Q ₀	0	0 0 0 1			
Restore	0	0 0 0 1		0 0 0 1	
		1 1		0 0 0 1	
Shift	0	0 0 1 0		0 0 1	
Subtract M	1	1 1 0 1			
Set Q ₀	1	1 1 1 1			
Restore	1	1 1 1 1		0 0 ¹ 0	
		1 1		0 0 ¹ 0	
Remainder	0	0 0 1 0		0 0 1 0	Quotient

The division algorithm just discussed needs restoring register A after each unsuccessful subtraction. Subtraction is said to be unsuccessful if the result is negative. Therefore it is referred to as restoring division algorithm. This algorithm is improved, giving non-restoring division algorithm.

Non-restoring algorithm :-

The steps for non-restoring algorithm are

Step 1: If the sign of A is 0, shift A and Q left one bit position and subtract divisor from A, otherwise shift A and Q left and add divisor to A. If the sign of A is 0, set Q_0 to 1 otherwise set Q_0 to 0.

Step 2: Repeat step 1 for n times.

Step 3: If the sign of A is \neq , add divisor to A.

Ex: perform division operation for the following using non-restoring algorithm.

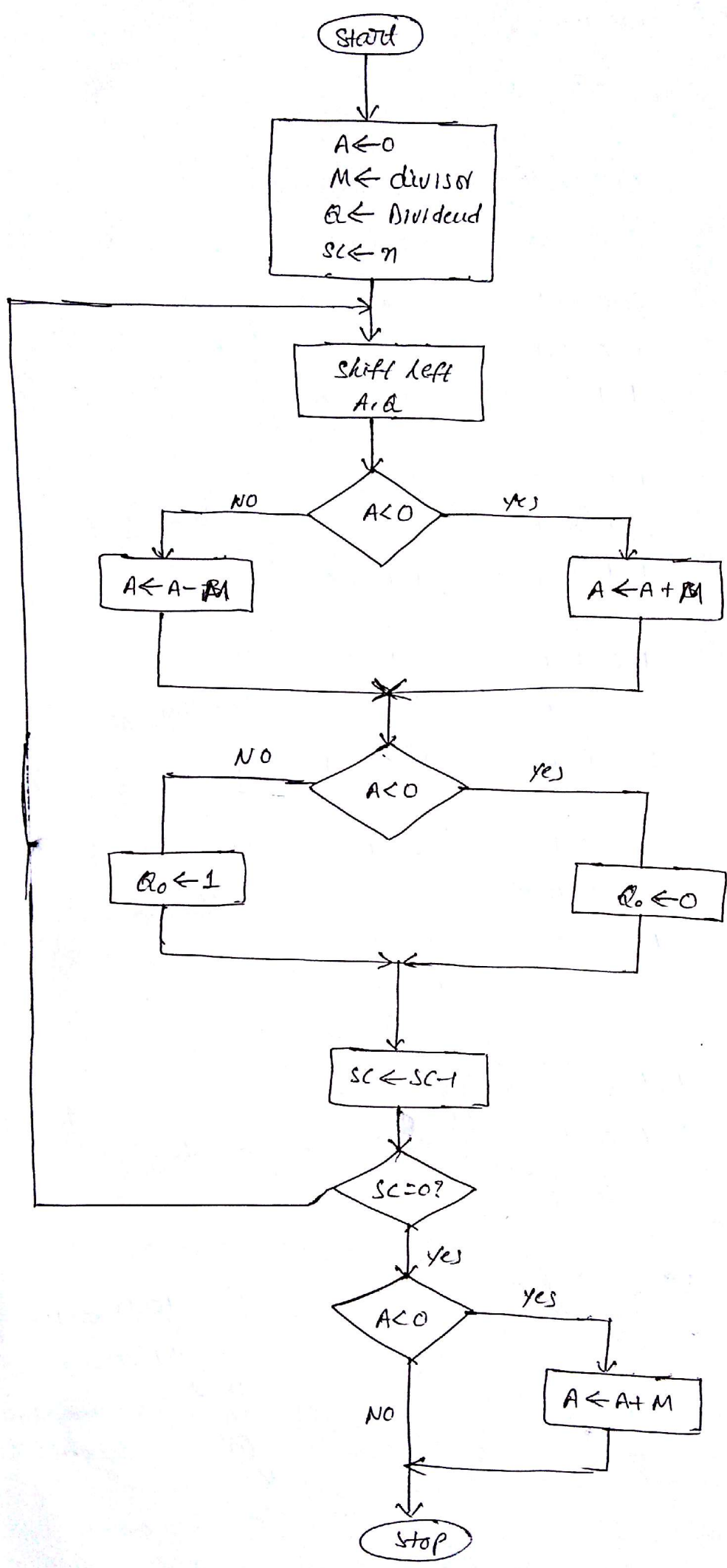
(i) Dividend = 1010

Divisor = 0011

(ii) Dividend = 1011

Divisor = 0101

The following diagram gives the flow of non-restoring algorithm.



Ex: Dividend = 1011 @
 Divisor = 0101 M

	E	A register	Q register	
initially	0	0000	1011	
Shift	0	0001	011	} First Cycle
Subtract	1	1011		
Set Qo	1	1100	0110	
Shift	1	1000	110	} Second Cycle
Add	0	0101		
Set Qo	1	1101	1100	
Shift	1	1011	100	} Third Cycle
Add	0	0101		
Set Qo	0	0000	1001	
Shift	0	0001	001	} Fourth Cycle
Subtract	1	1011		
Set Qo	1	1100	0010	

Quotient.

Add	1	1100	} Restore remainder
	0	0101	
	0	0001	

Remainder

1011 = 11
 0101 = 5
 5) 11 (2 → ~~Remainder~~ = 0010
 10

 11
 10

 11
 10

 11
 10

 11
 10

 11

Quotient
 Remainder = 0001

Floating-point numbers and operations:-

Till now we have seen fixed-point numbers and have considered them as integers. If B is a binary vector, then we have seen that it can be interpreted as an unsigned integer by:

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

This vector has an implicit binary point to its immediate right, then it can be represented as

$$b_{n-1} b_{n-2} b_{n-3} \dots b_1 b_0 \bullet \text{ implicit binary point}$$

Suppose if the binary vector is interpreted with the implicit binary point is just left of the sign bit:

$$\text{implicit binary point} \bullet b_{n-1} b_{n-2} b_{n-3} \dots b_1 b_0$$

Then the value of b is given by

$$V(b) = b_{n-1} \cdot 2^{-1} + b_{n-2} \cdot 2^{-2} + b_{n-3} \cdot 2^{-3} + \dots + b_1 \cdot 2^{-n-1} + b_0 \cdot 2^{-n}$$

While using floating point numbers we are going to have fraction numbers. Those fractions can be represented ~~the~~ by using above notations.

The value of the unsigned binary fraction is given by

$$V(B) = b_{31} \cdot 2^{-1} + b_{30} \cdot 2^{-2} + \dots + b_1 \cdot 2^{-31} + b_0 \cdot 2^{-32}$$

the range can be given by

$$0 \leq V(b) \leq 1 - 2^{-32} \approx 0.9999999998$$

In general for a n -bit binary fraction, the range value is given by

$$0 \leq V(b) \leq 1 - 2^{-n}$$

IEEE Standard Floating-Point Numbers :-

Previous representations have a fixed point. Either the point is to the immediate right or it is to the immediate left. This is called fixed point representation. Here there is a drawback that the representation can only represent a finite range of numbers.

A more convenient representation is the scientific representation, where the numbers are represented in the form:

$$x = m_1 \cdot m_2 m_3 \times b^{\pm e}$$

Components of these numbers are:

$m \rightarrow$ Mantissa

$b \rightarrow$ implied base

$e \rightarrow$ exponent.

A number such as the following is said to have 7 significant digits.

$$x = \pm 0. m_1 m_2 m_3 m_4 m_5 m_6 m_7 \times b^{\pm e}$$

Fractions in the range 0.0 to 0.9999999 need about 24 bits of precision. For example the binary fraction with 24 bits:

111111111111111111111111 = 0.9999999404.

Every real number between 0 and 0.9999999404 can be represented by a 24-bit fractional number.

The smallest non-zero number that can be represented is:

000000000000000000000001 = 5.96046 x 10^-8

Every other non-zero number is constructed in increments of this value.

A 24-bit mantissa can approximately represent a 7-digit decimal number, and an 8-bit exponent to an implied base of 2 provides a scale factor with a reasonable range. One bit is needed for the sign of the number. Since the leading nonzero bit of a normalized binary mantissa must be a 1, it does not have to be included explicitly in the representation. Therefore a total of 32 bits is needed.

The standard representation of floating-point number in 32-bits has been developed and specified in detail by the Institute of Electrical and Electronics Engineers (IEEE).

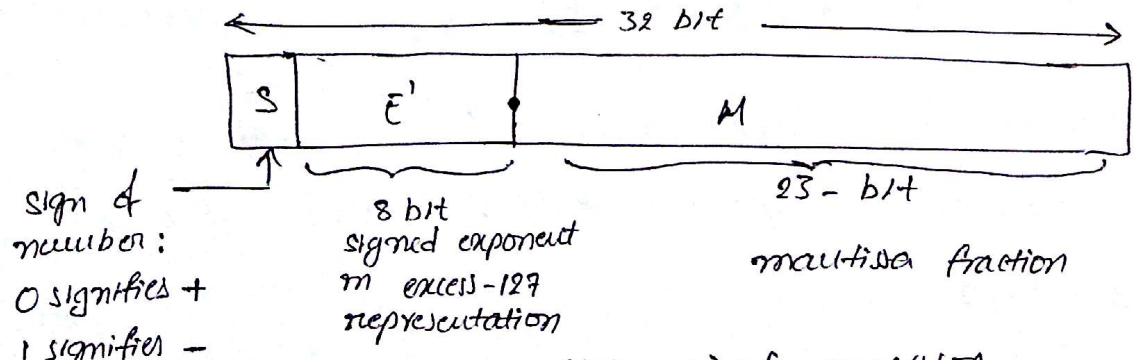
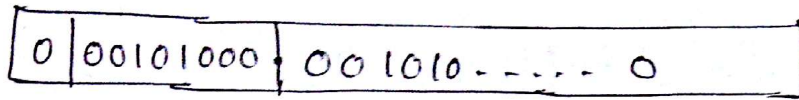


Fig: single precision

value represented in single precision is

Ex: $\pm 1.M \times 2^{E'-127}$



value represented = $1.001010 \dots 0 \times 2^{-87}$

NOTE!:- There are two representation for IEEE ~~float~~ Standard floating-point representation. They are

(i) single precision

(ii) double precision

→ single precision floating point represented in 32-bit with the exponent in excess of 127.

→ double precision floating point represented in 64-bit with the exponent in excess of 1023.

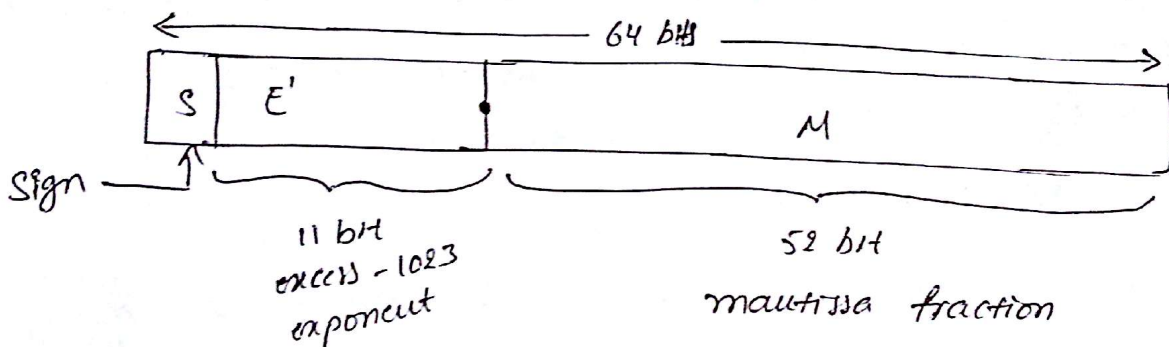


fig: Double precision.

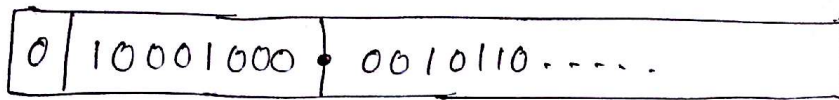
value represented as = $\pm 1.M \times 2^{E'-1023}$

Instead of the signed exponent E , the ~~value~~ value actually stored in the exponent field is an unsigned integer

$E' = E + 127.$

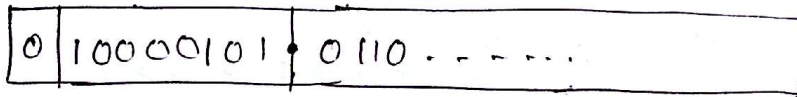
This is called the excess-127 format. Thus E' is in the range $0 \leq E' \leq 225$. The range for normal values is $1 \leq E' \leq 224$. This means that actual exponent E is in the range $-126 \leq E \leq 127$. We should consider two aspects of operating with floating-point numbers.

First if a number is not normalized, it can always be put in normalized form by shifting the fraction and adjusting the exponent.



value represented as = $+0.0010110..... \times 2^9$

Fig:- unnormalized value



value represented as = $+1.0110..... \times 2^6$

Fig:- Normalized value

Above example is for underflow.

→ underflow occurs when the result of a floating point is larger than the largest positive value or smaller than the smallest negative value

→ overflow occurs when the result of a floating point operation is smaller than the smallest positive number or larger than the largest negative value

Arithmetic Operations on Floating point Numbers :-

If the exponents are different, the mantissas of floating-point numbers must be shifted with respect to each other before they are added or subtracted. Consider the following decimal example in which we wish to ~~add~~ perform addition.

$$2.9400 \times 10^2 \quad \text{with} \quad 4.3100 \times 10^4$$

Here the exponents are different so we need to shift the floating point number. So we rewrite 2.9400×10^2 to 0.0294×10^4 .

Addition and Subtraction :-

Consider two floating point numbers:

$$X = m_1 \cdot \gamma^{e_1} \quad \text{and}$$

$$Y = m_2 \cdot \gamma^{e_2} \quad \text{Assume } e_1 \geq e_2.$$

Let us see the rules for addition and subtraction.

Step 1: Select the number with a smaller exponent and shift its mantissa right, a number of steps equal to the difference in exponents i.e. $|e_2 - e_1|$.

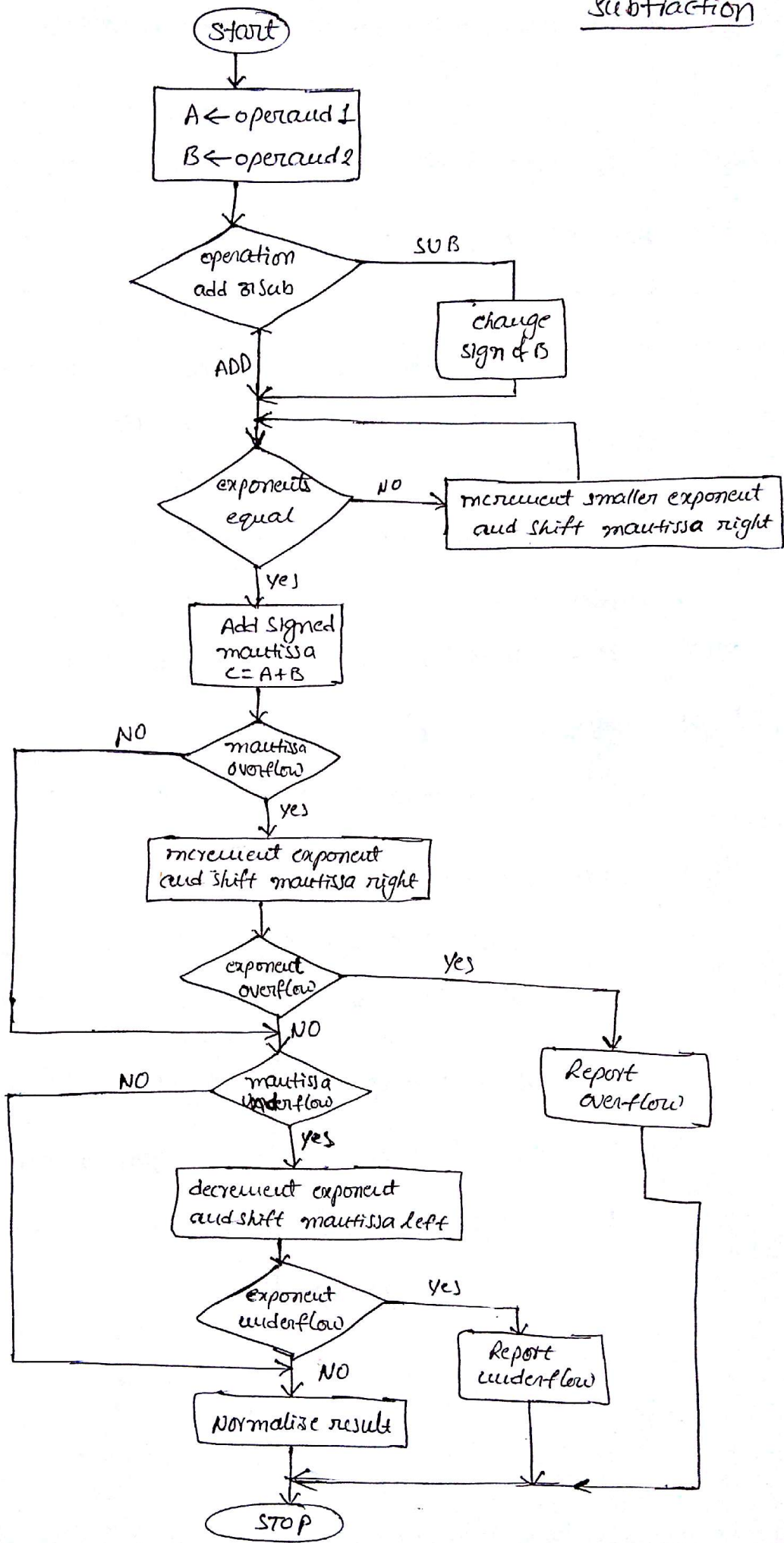
Step 2: Set the exponent of the result equal to the larger exponent.

Step 3: Perform addition/subtraction on the mantissa and determine the sign of the result.

Step 4: Normalize the result, if necessary.

flow chart for floating point addition and

Subtraction



Ex:- Add single precision floating point numbers A and B where $A = 44900000H$ and $B = 42A00000H$.

Sol:-

Step 1: Represent numbers in single precision format.

$$A = 0 \ 10001001 \ 0010000 \ \dots \ 0$$

$$B = 0 \ 10000101 \ 0100000 \ \dots \ 0$$

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

$$\text{Actual exponent} = 137 - 127 = 10$$

$$\text{Exponent for B} = 10000101 = 133$$

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

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

Step 2: Shift mantissa

$$\text{Shifted mantissa of B} = 00000100 \ \dots \ 0$$

Step 3: Add mantissa

$$\text{Mantissa of A} = 00100000 \ \dots \ 0$$

$$\text{Mantissa of B} = 00000100 \ \dots \ 0$$

$$\text{Mantissa of result} = 00100100 \ \dots \ 0$$

As both numbers are positive, sign of the result is positive.

$$\therefore \text{Result} = 0 \ 10001001 \ 00100100 \ \dots \ 0$$

$$= 44920000H.$$

Problems in Floating Point Arithmetic :-

- ① Mantissa overflow :- The addition of two mantissas of the same sign may result in a carryout of the most significant bit. If so, the mantissa is shifted right and the exponent is incremented.
- ② Mantissa underflow :- In 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.
- ③ 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 $-\infty$.
- ④ Exponent underflow :- Exponent underflow occurs when a negative exponent exceeds the maximum possible exponent value. In such cases, the number is designated as zero.

Multiplication rules :-

Floating point multiplication is somewhat easier than addition and subtraction because an alignment of mantissa is not required.

- ① Add the exponents and subtract bias. i.e 127 in case of single precision number and 1023 in case of double precision numbers.
- ② Multiply the mantissas and determine the sign of the result.
- ③ Normalize the resulting value, if necessary.

Division Rules:-

- ① Subtract the exponents and add bias (127 in case of single precision numbers and 1023 in case of double precision number).
- ② Divide the mantissas and determine the sign of the result.
- ③ Normalize the resulting value, if necessary.

Guard Bits and Truncation:-

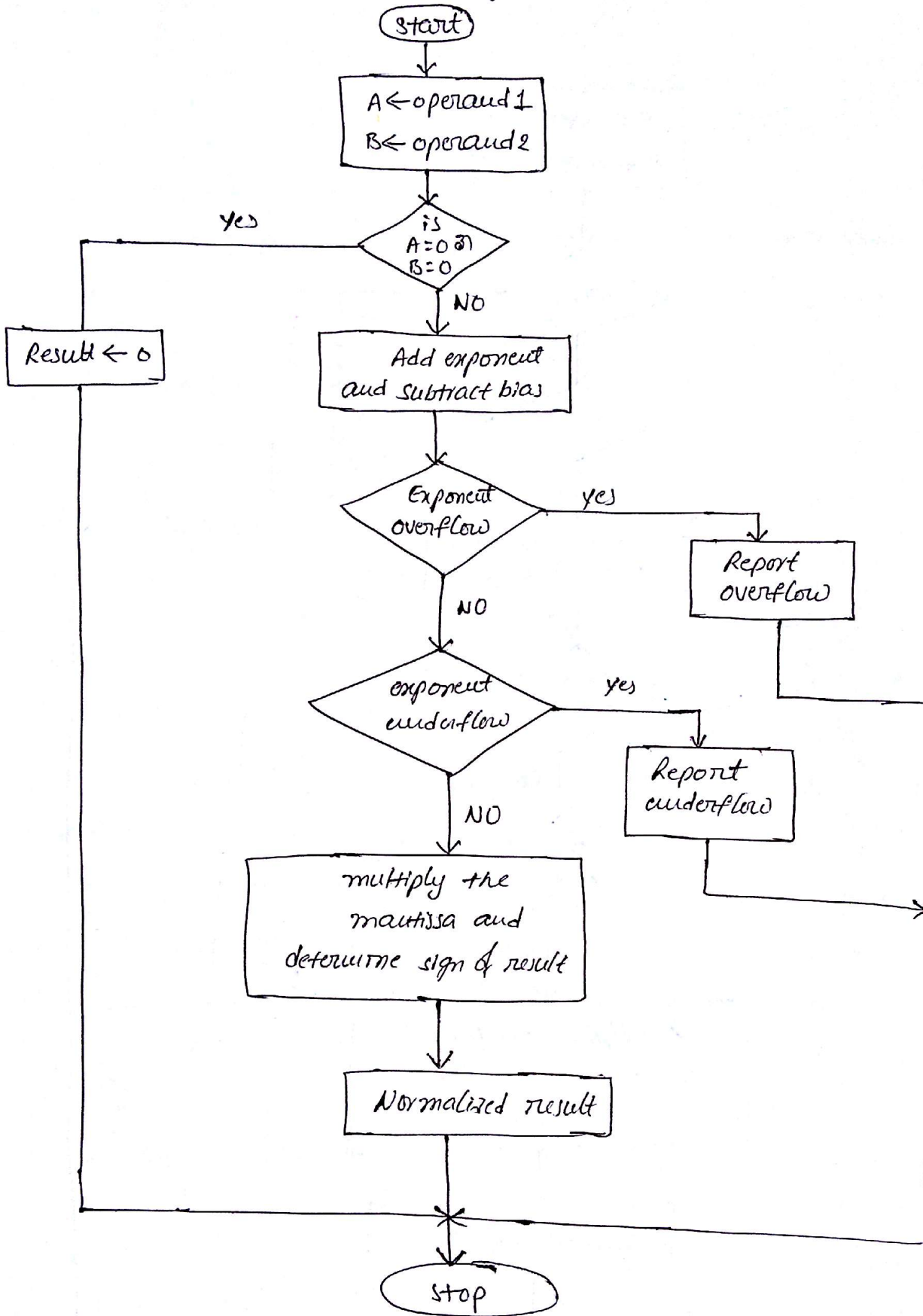
While adding two floating point numbers with 24-bit mantissas, we shift the mantissa of the number with the smaller exponent to the right until the two exponents are equalized. This implies that mantissa bits may be lost during the right shift. To prevent this floating point operations are implemented by keeping guard bits.

Removing guard bits in generating a final result requires that the extended mantissa be truncated to create 24-bit number that approximates the longer version.

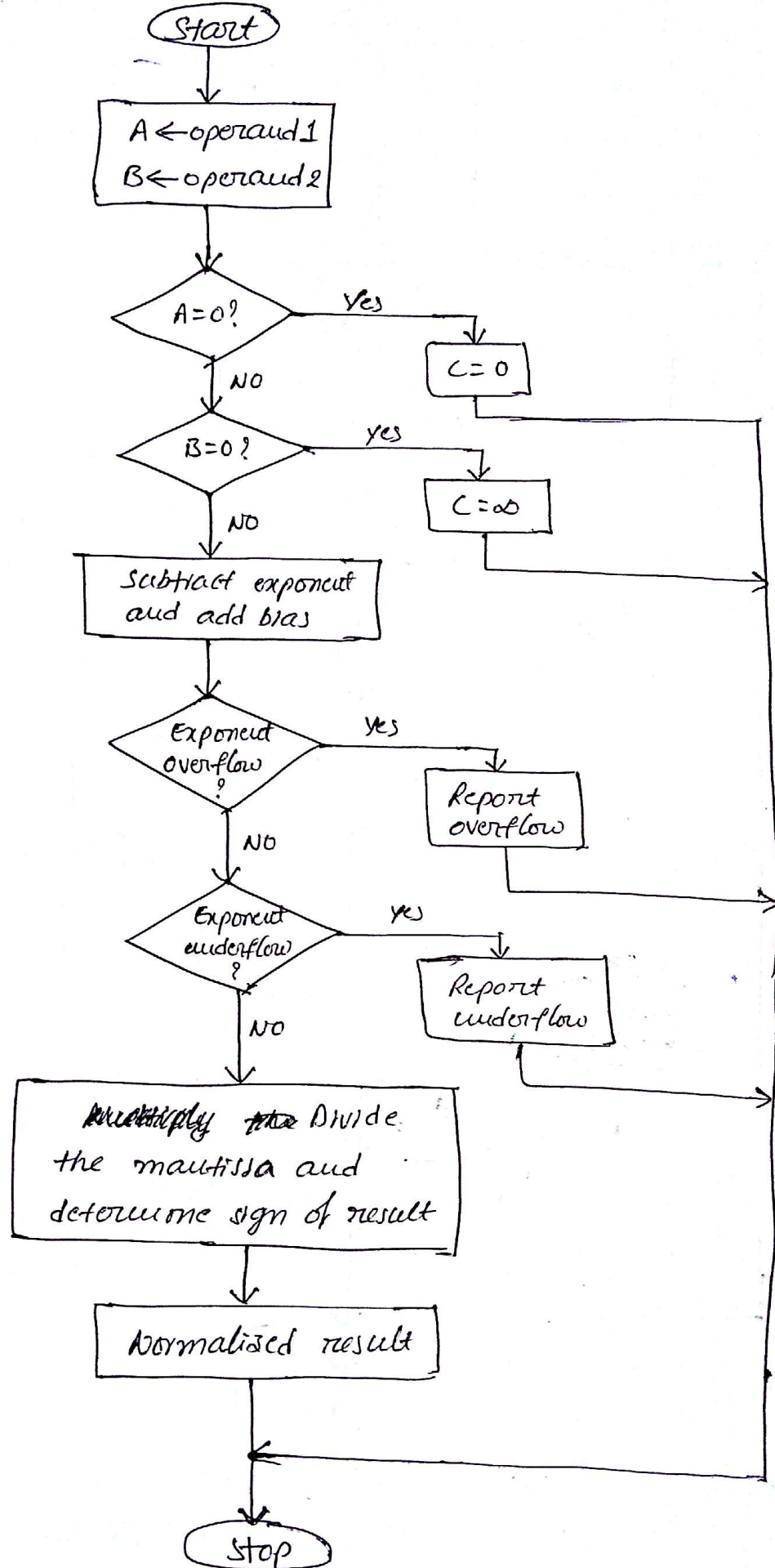
There are several ways to truncate. The simplest way is to remove the guard bits and make no changes to the retained bits. The next simplest method of truncation is Von Neuman rounding. If the guard bits are all 0, they are dropped. However if any bit of the guard bit is a 1, then the LSB of the retained bit is set to 1.

The third truncation method is a rounding procedure. If there is a 1 in the MSB of the guard bit then a 1 is added to the LSB of the retained bits.

flow chart for floating point multiplication.



flow chart for floating point Division



Implementing Floating-Point Operations :-

The hardware implementation of floating-point operations involves a considerable amount of logic circuitry. These operations can also be implemented by s/w routines.

In most general-purpose processors, floating-point operations are available at the machine-instruction level, implemented in hardware. The following diagram gives the implementation of floating-point operations.

In addition/subtraction operation, the first rule is to compare exponents to determine how far to shift the mantissa of the number with the smaller exponent. The shift-count value n , is determined by the 8-bit subtractor circuit in the upper left corner of the figure. The magnitude of the difference $E'_A - E'_B$, $\geq n$ is sent to the SHIFTER unit.

The ^{sign} difference from comparing exponents determines which mantissa is to be shifted. The sign is sent to the SWAP network. If the sign is 0 then $E'_A \geq E'_B$ and the mantissas M_A and M_B are sent straight through the SWAP network. This results in M_B being sent to the SHIFTER, to be shifted n positions to the right. The other mantissa, M_A is sent directly to the mantissa adder/subtractor. If the sign is 1, then $E'_A < E'_B$ and the mantissas are swapped before they are sent to the SHIFTER.

Step 2: The exponent of the result E' is determined as E'_A if $E'_A \geq E'_B$, $\geq E'_B$ if $E'_A < E'_B$ based on the sign of the difference resulting from comparing exponents in step 1. This is done by two-way multiplexer.

32-bit operands $\left\{ \begin{array}{l} A: S_A, E_A', M_A \\ B: S_B, E_B', M_B \end{array} \right\}$

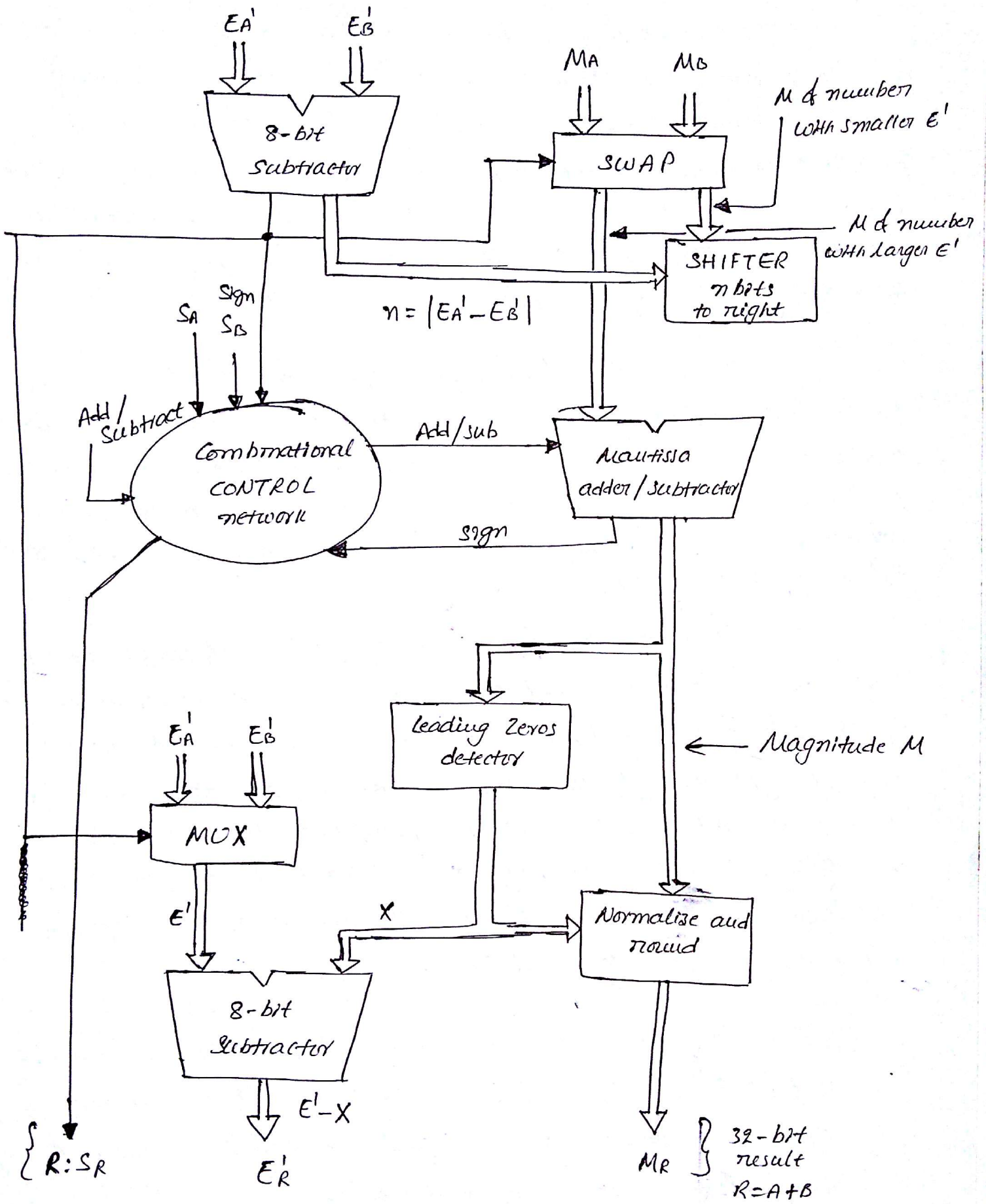


Fig: Floating-point addition-subtraction unit.

(23)

Step 3: involves the major component i.e. the mantissa adder / subtractor. The control logic determines whether the mantissas are to be added or subtracted. It is decided by the signs of operands i.e. S_A and S_B . The control logic also determines the sign of result S_R .

For example ~~the sign~~ A is negative ($S_A = 1$) and B is positive ($S_B = 0$) and the operation is $A - B$, then the mantissas are added and the sign of result is negative ($S_R = 1$).

On the other hand if A and B are positive and the operation is ~~positive~~ $A - B$, then the mantissas are subtracted. The sign of the result S_R now depends up on the mantissas subtraction.

- (i) if $E'_A > E'_B$ then $M_A - (\text{Shifted } M_B)$ is positive and the result is positive.
- (ii) if $E'_B > E'_A$ then $M_B - (\text{Shifted } M_A)$ is positive and the result is negative.

Step 4: Add / Subtraction operation result will be normalized if necessary

~~→~~ ~~→~~ N. Subrah