

Chapter



Data Types and Computer Arithmetic

Syllabus

- ❖ Scalar Data Types
- ❖ Fixed and Floating point numbers
- ❖ Signed numbers
- ❖ Integer Arithmetic
- ❖ 2's Complement multiplication
- ❖ Booth's Algorithm
- ❖ Hardware Implementation
- ❖ Division
- ❖ Restoring and Non Restoring algorithms
- ❖ Floating point representations
- ❖ IEEE standards
- ❖ Floating point arithmetic

Information is represented in computer in the form of Binary digit, called bit. We have to learn how arithmetic calculations are performed through these bits.

2.1 Number Systems :

A number system of base r is a system with r distinct symbols for r digits.

System	Base	Symbols (digits)
Decimal	10	0, 1, 2, 3, 4, 5, 6, 7, 8, 9
2	0, 1	
Binary	8	0, 1, 2, 3, 4, 5, 6, 7
Octal	16	0, 1, 2, 3, 4, 5, 6, 7, 8, 9, A(10), B(11), C(12), D(13), E(14), F(15)
Hexadecimal		

To determine the quantity that the number represents we multiply the digit by an integer power of r (base) depending on the place it is located and then find the sum of weighted digits.

Example 1 : Let us consider a number ' $d_4 d_3 d_2 d_1 d_0$ ' with base ' r ' then

$$\text{the value of the number} = d_4 \times r^4 + d_3 r^3 + d_2 r^2 + d_1 r^1 + d_0.$$

Solution :

A number $d_4 d_3 d_2 d_1 d_0$ in base r is written as,

$$(d_4 d_3 d_2 d_1 d_0)_r$$

2.1.1 Decimal Numbers :

Decimal number system has **10 digits** represented by 0, 1, 2, 3, 4, 5, 6, 7, 8, 9. Any decimal number can be represented as a string of these digits.

A decimal number 7452.694 should be written as,

$$7 \times 10^3 + 4 \times 10^2 + 5 \times 10^1 + 2 \times 10^0 + 6 \times 10^{-1} + 9 \times 10^{-2} + 4 \times 10^{-3}$$

2.1.2 Binary Number :

Binary number system has **2 digits** represented by 0, 1. Any binary number can be represented as a string of these two digits called bits. The base of binary number system is 2.

Decimal equivalent :

$$\begin{aligned}
 (101101.101)_2 &= 1 \times 2^5 + 0 \times 2^4 + 1 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 \\
 &\quad + 1 \times 2^{-1} + 0 \times 2^{-2} + 1 \times 2^{-3} \\
 &= 2^5 + 2^3 + 2^2 + 1 + \frac{1}{2} + \frac{1}{2^3} = 45.625
 \end{aligned}$$

Thus, $(101101.101)_2 = (45.625)_{10}$

2.1.3 Octal Numbers :

As octal system has **eight digits** represented as 0, 1, 2, 3, 4, 5, 6, 7. Any octal number can be represented as a string of these digits. The base of octal number system is 8.

Decimal equivalent :

$$\begin{aligned}
 (45.6)_8 &= 4 \times 8^1 + 5 \times 8^0 + 6 \times 8^{-1} \\
 &= 37 + \frac{6}{8} = 37.75 \\
 \therefore (45.6)_8 &= (37.75)_{10}
 \end{aligned}$$

2.1.4 Hexadecimal Numbers :

The hexadecimal system has **16 digits**, which are represented as 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E, F. The base of hexadecimal number system is 16.

Decimal Equivalent:

$$\begin{aligned}
 (F3.4)_{16} &= F \times 16^1 + 3 \times 16^0 + 4 \times 16^{-1} \\
 &= 15 \times 16 + 3 + \frac{4}{16} \\
 &= 243.25 \\
 (F3.4)_{16} &= (243.25)_{10}
 \end{aligned}$$

2.1.5 BCD (Binary Coded Decimal) Numbers :

In BCD system, each digit of decimal number is represented using a 4 bit binary number.

Decimal digit	Primary representation
0	0000
1	0001
2	0010
3	0011
4	0100
5	0101
6	0110
7	0111
8	1000
9	1001

BCD equivalent of a decimal number :

$(5943)_{10} = 0101 \ 1001 \ 0100 \ 0011$ in BCD

In the above representation, both sign and magnitude of a number are represented separately.



Fig. 2.2 : A $(n + 1)$ bit number

The sign of the number is represented using the left most bit. The left most bit is 0 for positive numbers and 1 for negative numbers. Thus, a number of n bits can be represented as $n + 1$ bit number, where $(n + 1)^{\text{th}}$ bit is the sign bit and rest n bits represent its magnitude.

Example 2 :

Represent 12 and -12 using sign magnitude representation. Assume size of register = 8 bits.

Solution :

Number	Representation
+12	<u>0_0001100</u> Sign Magnitude
-12	<u>1_0001100</u> Sign Magnitude

Disadvantages of sign magnitude representation :

- There are two representation for 0.
- +0 is represented as 0 000 0000
- 0 is represented as 1 000 0000

To do the arithmetic addition with one positive number and one negative number, we have to check the magnitude of numbers. The number having smaller magnitude is then subtracted from the bigger number and the sign of bigger number is selected. The implementation of such a scheme in digital hardware is complicated.

1's complement representation :

In 1's complement representation, negative values are obtained by complementing each bit of the corresponding positive number.

Example 3 :

Represent 12 and -12 using 1's complement representation. Assume size of register = 8 bits.

2.2 Representation of Data in Computer :

The binary number system is most nature for computer. Computer stores information in flip-flops, which are two state devices. Three systems are widely used for representing both positive and negative numbers :

- Sign magnitude representation
- 1's - complement representation
- 2's - complement representation

In the fixed point numbers we assume that the position of the binary point is at the end. It implies that all integers are represented using **fixed point representation**. Real numbers are represented using **floating point representation**.

2.2.1 Fixed Point Representation :

Solution :	Number	Representation
+ 12	$\begin{array}{ c } \hline 0 \\ \hline 00001100 \\ \hline \end{array}$	Sign bit 0 indicates it is a positive number
- 12	$\begin{array}{ c } \hline 1 \\ \hline 110011 \\ \hline \end{array}$	Sign bit 1 indicates it is a negative number

Disadvantages :

- There are two representations for 0.
- + 0 is represented as 00000000
- 0 is represented as 11111111

2's complement representation :

- In 2's complement representation, negative number is obtained by adding 1 to 1's complement representation.
- Example 4 : Represent 12 and - 12 using 2's complement representation. Assume size of register = 8 bits.

Number	Representation
+ 12	0000 1100
- 12	1's complement of (0000 1100) + 1 = 1111 0011 + 1 = 11101000

Advantages :

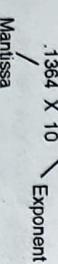
- Only one representation for 0.
- Subtraction is carried out like addition.

2.2.2 Floating Point Representation :

Floating point number representation consists of two parts :

- Mantissa
- Exponent

A number 13.64 can be written as 1.364×10^2 .



Sign of number
Mantissa
Exponent

Exponent indicates the correct decimal location. An exponent of + 2, indicates that actual position of decimal point is two places to the right of the assumed position, while - 2 indicates that the assumed position of the point is two places towards the left of assumed position. Decimal point is assumed and it is never represented physically.

Floating point numbers are often represented in normalized form. A floating point number whose mantissa does not contain zero as the most significant digit of the number is considered to be in normalized form.

Binary number**Normal form****Representation**

0100	.100 $\times 2^{-1}$	Sign
$\begin{array}{ c } \hline 0 \\ \hline 1000000000 \\ \hline \end{array}$	$\begin{array}{ c } \hline 100 \\ \hline 0000000000 \\ \hline \end{array}$	Sign Mantissa Exponent

11.0101	.110101 $\times 2^2$	Sign
$\begin{array}{ c } \hline 1101010000 \\ \hline 010 \\ \hline \end{array}$	$\begin{array}{ c } \hline 1101010000 \\ \hline 010 \\ \hline \end{array}$	Sign Mantissa Exponent

Register size is assumed to be of 16 bits. 12 bits are for mantissa and four bits are for exponent.

2.3 Computer Arithmetic :**2.3.1 Addition of Positive Numbers :**

Binary addition is performed in a manner similar to that of decimal addition. Corresponding bits are added, and if a carry 1 is produced, it is added to the binary digit at the left.

$$\begin{array}{r}
 & \overset{0}{+} & \overset{0}{+} & \overset{1}{+} & \overset{1}{+} \\
 & 0 & 0 & 1 & 1 \\
 \text{Carry} & \swarrow & \swarrow & \swarrow & \swarrow \\
 \text{Sum} & 0 & 0 & 0 & 1 \\
 \text{Carry} & \swarrow & \swarrow & \swarrow & \swarrow \\
 \text{Sum} & 0 & 0 & 0 & 1 \\
 \text{Carry} & \swarrow & \swarrow & \swarrow & \swarrow \\
 \text{Sum} & 0 & 0 & 0 & 1 \\
 \text{Carry} & \swarrow & \swarrow & \swarrow & \swarrow \\
 \text{Sum} & 0 & 0 & 0 & 1
 \end{array}$$

Fig. 2.3(a) : Addition of 1 bit numbers

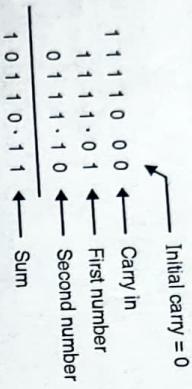
Fig. 2.3(b) : Addition of 1 bit numbers with previous carry

x _i	y _i	Carry-in c _{i-1}	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	1	0
1	0	0	1	0
1	0	1	0	1
1	1	0	0	1
1	1	1	0	1
1	1	1	1	1

Example 5 : Find the sum of two binary numbers 1111.01 and 0111.10.

Solution :

$$\begin{array}{r} 1111.01 = (15.25)_{10} \\ + 0111.10 = (7.50)_{10} \\ \hline 10110.11 = (22.75)_{10} \end{array}$$



Example 6 : Find the expression for sum s_i and carry-out c_{i+1} , in terms of x_i, y_i and c_i .

Solution :

From the truth table given in Fig. 2.4, we have,

$$\begin{aligned} s_i &= \bar{x}_i \bar{y}_i c_i + \bar{x}_i y_i \bar{c}_i + x_i \bar{y}_i \bar{c}_i + x_i y_i c_i \\ &= x_i \oplus y_i \oplus c_i \\ &\quad \text{exclusive or} \\ c_{i+1} &= \bar{x}_i y_i c_i + x_i \bar{y}_i \bar{c}_i + x_i y_i \bar{c}_i + x_i \bar{y}_i c_i \\ &= x_i y_i + x_i c_i + y_i c_i \end{aligned}$$

The implementation of full adder follows directly from the logic expression of truth table in Fig. 2.4.

Example 7 :

Design an adder to add two 4-bit numbers.

Solution :

$$\begin{array}{ll} \text{First number} & X = x_3 x_2 x_1 x_0 \\ \text{Second number} & Y = y_3 y_2 y_1 y_0 \end{array}$$

Here x_i and y_i ($i = 0$ to 3) represent a bit.

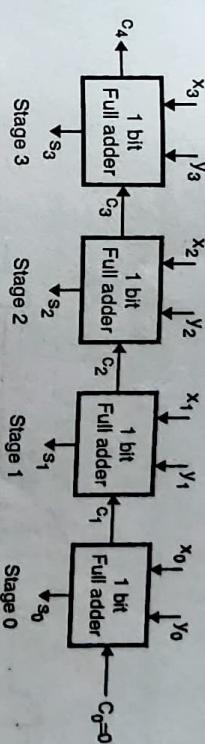


Fig. 2.5 : Logic circuit for addition of two 4-bit binary numbers using carry propagation

s_3, s_2, s_1, s_0 represents the overall sum and overall carry is c_4 , from the third stage. The main feature of this adder is that the carry of each lower bit is fed to the next higher bit addition stage, it implies that addition of the next higher bit has to wait for the previous stage addition. Such an adder is called ripple carry adder. The ripple carry addition becomes time consuming when we add two long binary numbers.

- To overcome the problem of excessive delay in ripple carry adder, a fast adder, carry look ahead adder can be designed.
- In carry look-ahead adders, the carry-in for various stages can be generated directly by the logic expressions.

Example 8 :

Design a 4-bit carry look ahead adder.

Solution :

We know that,

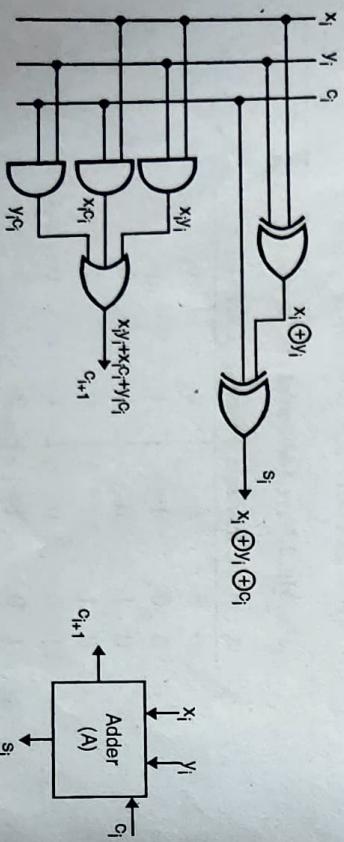
$$c_{i+1} = x_i y_i + x_i c_i + y_i c_i$$

and $s_i = c_i \oplus x_i \oplus y_i$

Let us introduce two more variables g_i and p_i or generate and propagate respectively. g_i and p_i are defined by the following logic equations,

$$g_i = x_i y_i$$

$$p_i = x_i + y_i$$



(a) Logic for one bit adder (full adder) (b) Schematic symbol for 1 bit full adder
Fig. 2.4 : One bit full adder

Continuing in this way, c_{i+1} can be expressed as a sum of products functions of the p and g outputs of all the preceding stages. The carries in a four-stage carry-look ahead are defined as follows :

$$\begin{aligned} c_1 &= g_0 + p_0 c_0 \\ &= g_1 + p_1 c_1 \\ &= \dots (1) \end{aligned}$$

Similarly, c_i can be expressed in terms of g_{i-1} , p_{i-1} and c_{i-1} .

$$c_i = g_{i-1} + p_{i-1} c_{i-1}$$

on substituting Equation (2) into Equation (1) we obtain,

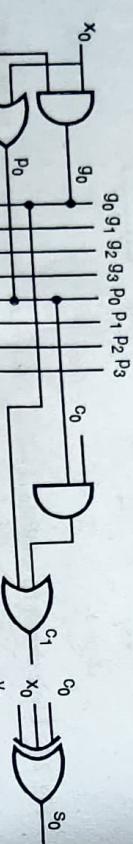
$$c_{i+1} = g_i + p_i (g_{i-1} + p_{i-1} c_{i-1})$$

Or,

$$c_{i+1} = g_i + p_i g_{i-1} + p_i p_{i-1} c_{i-1}$$

2's complement representation is the best choice for addition of two signed numbers or for finding difference of two numbers.

If x and y are two positive numbers then :



Example 9 :

Find the value of the following expressions using 2's complement arithmetic

- (1) 55 + 27
- (2) -55 + 27
- (3) 55 - 27
- (4) -55 - 27

Assume register size to be of 8 bits.

Solution :

Number	Binary equivalent	2's complement
55	00110111	11001001
27	00011011	11100101

(1) 55 + 27 :

$$\begin{array}{r}
 55 \quad 00110111 \\
 + 27 \quad 00011011 \\
 \hline
 82 \quad 01010010
 \end{array}$$

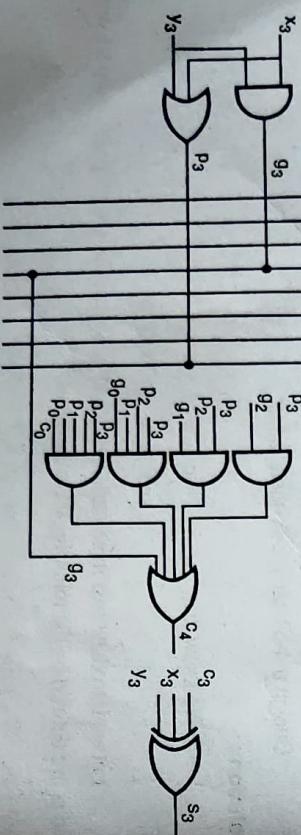


Fig. 2.6 : 4 bit carry look ahead adder

(2) $-55 + 27 :$

$$\begin{array}{r} -55 \\ +27 \\ \hline -28 \end{array}$$

1100100
00011011
1100100

We can verify that -28 in 2's complement representation is 1100100, by taking 2's complement of 11100100. It should yield 28(0001100).

2's complement of 11100100 = $1 + 00011011$

$$= 00011100$$

(3) $55 - 27 :$

$$\begin{array}{r} 55 \\ -27 \\ \hline -28 \end{array}$$

00110111
11100101
00011100

Carry out from the most significant bit should be discarded
[2's complement of 27]

(4) $-55 - 27 :$

$$\begin{array}{r} -55 \\ -27 \\ \hline -82 \end{array}$$

11001001
11100101
10101110

Carry out should be discarded

Since, the sign of the result has changed, an underflow has occurred.

2.3.2.1 Overflow/Underflow in 2's Complement Arithmetic :

Overflow can occur when two positive numbers are added. An overflow is said to have occurred if the sign of the resultant sum of two numbers is different.

Example 10 :

Add two numbers 65 and 75 using 2's complement arithmetic. Assume register size = 8 bits.

Solution :

$$\begin{array}{r} 65 \\ 75 \\ \hline 140 \end{array}$$

01000001
01001011
10001100
[Carry out from the sign bit is discarded]

Since, the sign of the result has changed, an overflow has occurred.

Underflow can occur when two negative numbers are added. An underflow is said to have occurred if the sign of the resultant sum of two numbers is different.

Solution :

Add two numbers -65 and -75 using 2's complement arithmetic. Assume register size = 8 bits.

Number Binary equivalent 2's complement

Number	Binary equivalent	2's complement
65	01000001	10111110
75	01001011	10110101

$$\begin{array}{r} -65 \\ -75 \\ \hline 10111110 \\ 10110101 \\ \hline 01110011 \end{array}$$

[Carry out from the sign bit is discarded]

Sign of the result is different. It has become a positive number.

Example 12 :
Design a 4-bit subtractor using four full adders.

Solution :

$$X = x_3 x_2 x_1 x_0$$

$$Y = y_3 y_2 y_1 y_0$$

$$\therefore X - Y = X + 2^4 \text{ complement of } Y$$

$$= x_3 x_2 x_1 x_0 + \overline{y_3} y_2 y_1 y_0 + 1$$

Above expression can be realized by adding X and complement of Y . Initial carry c_0 should be set to 1.

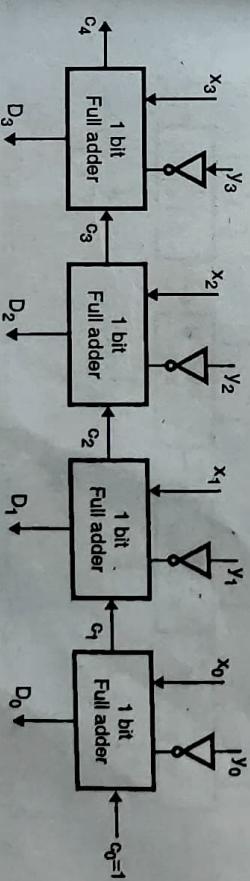


Fig. 2.7 : Logic circuit for subtraction of two 4-bit numbers

Example 13: A 4-bit addition/subtraction circuit if the external control $ctrl = 0$, it should add and if $ctrl = 1$, it should perform subtraction.

Design a 4-bit adder to perform addition of the two 4-bit numbers and

Solution :

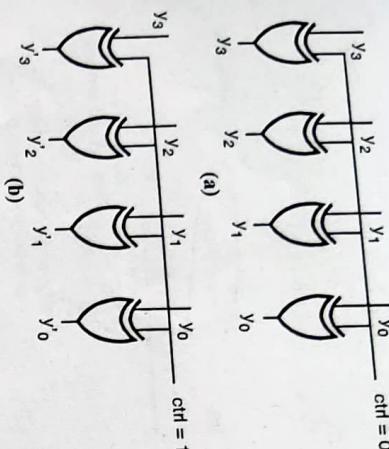


Fig. 2.8

When the signal $ctrl = 0$, the binary number passed through the network of exclusive or " " as an output without any change. We know,

$x \oplus 0 = x$ [signal passes without any change]
 $x \oplus 1 = x'$ [signal is inverted]

When the signal $ctrl = 1$, the binary number passes through the network of exclusive OR gates will appear inverted.

The total product is produced by summing the partial products.

Algorithmic approach

Let us take a variable x initialized with 0, i.e.,

卷之三

Step 1: Since the b_0 bit of multiplier is

Step 2: Right shift x

Step 3. Since the bits of $b_3 b_2 b_1 b_0$ are 0 1 0 1, we have

Step 3 : Since, the b_1 bit of multiplier is 0

$$x = x + 0000 = 0101 \mid 1 + 0000 = 0101 \mid 1$$

When $\text{ctrl} = 0$, circuit will produce sum of the input numbers

$$X + Y = X \bar{X} \bar{Y} \bar{Z} + \bar{Y} X \bar{Z} \bar{X} Z X = X + Y$$

when $\text{ctrl} = 1$, circuit will produce difference of input numbers

$$X - Y = X - \left(X + \frac{Y}{2} \right) = \frac{X - Y}{2}$$

2.3.3 Binary Multiplication (Unsigned Numbers)

The usual algorithm for multiplying two integers is shown in the Fig. 2.10. This algorithm can be used for multiplication of two unsigned or positive numbers. The product of two n-bit numbers can be accommodated in $2n$ bits.

Multiplicand M = 11

$b_3 \ b_2 \ b_1 \ b_0 \ - \ \text{bit}$
position

$$\left\{ \begin{array}{cccccc} & 1 & 0 & 1 & 1 & - & p_0 \\ & 0 & 0 & 0 & 0 & - & p_1 \\ \text{Partial products} & 1 & 0 & 1 & 1 & - & p_2 \\ 1 & 0 & 1 & 1 & - & - & p_3 \end{array} \right\} \quad \begin{matrix} p_1 \text{ is added to } p_0 \text{ after} \\ \text{right shifting } p_0 \end{matrix}$$

Example 13 : Design a 4-bit addition/subtraction circuit if the external control $\text{ctrl} = 0$ it should perform addition of the two 4-bit numbers and if $\text{ctrl} = 1$, it should perform subtraction.

Solution :

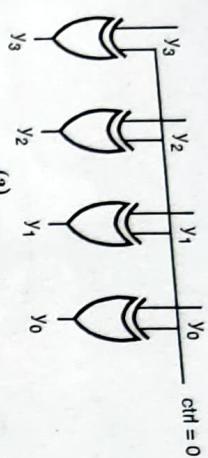


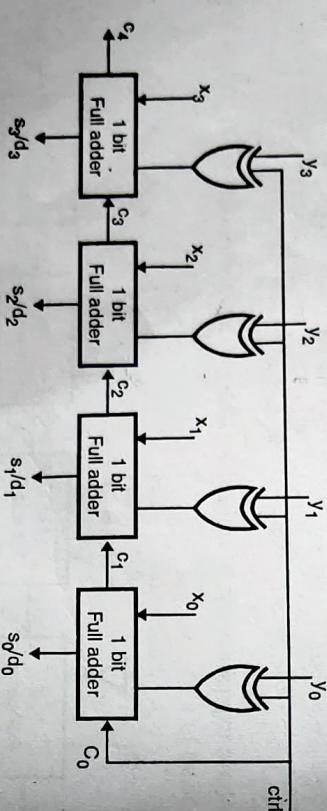
Fig. 2.8

- When the signal $\text{ctrl} = 0$, the binary number passed through the network of exclusive OR gates will appear on output without any change. We know,

$$x \oplus 0 = x \text{ [signal passes without any change]}$$

and $x \oplus 1 = x'$ [signal is inverted]

- When the signal $\text{ctrl} = 1$, the binary number passes through the network of exclusive OR gates will appear inverted.



When $\text{ctrl} = 0$, circuit will produce sum of the input numbers

$$X + Y = x_3 x_2 x_1 x_0 + y_3 y_2 y_1 y_0 + 0$$

when $\text{ctrl} = 1$, circuit will produce difference of input numbers

$$X - Y = x_3 x_2 x_1 x_0 + y_3 y_2 y_1 y_0 + 1$$

2.3.3 Binary Multiplication (Unsigned Numbers): The usual algorithm for multiplying two integers is shown in the Fig. 2.10. This algorithm can be used for multiplication of two unsigned or positive numbers. The product of two n-bit numbers can be accommodated in $2n$ bits.

Multiplicand $M = 11$

Multipplier $Q = 13$

$$\begin{array}{r} b_3 \ b_2 \ b_1 \ b_0 = \\ 1 \ 0 \ 1 \ 1 \ 0 \ 1 \ 1 \ 0 \ 1 \ 1 \ 0 \ 1 \ 1 \ 0 \ 1 \ 1 \end{array} \quad \begin{array}{l} \text{bit} \\ \text{position} \end{array}$$

$$\begin{array}{r} \hline 1 \ 0 \ 1 \ 1 \ 0 \ 1 \ 1 \ 0 \ 1 \ 1 \ 0 \ 1 \ 1 \ 0 \ 1 \ 1 \end{array} \quad \begin{array}{l} \hline p_0 \\ p_1 \\ p_2 \\ p_3 \end{array}$$

$$\begin{array}{c} \text{Partial products} \\ \left\{ \begin{array}{r} 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \\ 1 \ 0 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \\ \hline 1 \ 0 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \end{array} \right\} p_i \text{ is added to } p_0 \text{ after right shifting } p_0. \end{array}$$

$$\begin{array}{r} \hline 1 \ 0 \ 0 \ 0 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \end{array} \quad \begin{array}{l} \hline p_1 \\ p_2 \\ p_3 \end{array}$$

The total product is produced by summing the partial products.

Algorithmic approach

Let us take a variable x initialized with 0, i.e.,

$$x = 0000 \text{ [a four bit binary number]}$$

Step 1 : Since the b_0 bit of multiplier is 1

$$x = x + 1011 = 0000 + 1011 = 1011$$

Step 2 : Right shift x

$$\begin{array}{r} b_3 \ b_2 \ b_1 \ b_0 \\ 1 \ 0 \ 1 \ 0 \ 1 \ 1 \ 1 \ 1 \end{array}$$

Step 3 : Since, the b_1 bit of multiplier is 0

$$x = x + 0000 = 0101 \mid 1 + 0000 = 0101 \mid 1$$

Initial value of registers :

$$\begin{array}{r} \text{i.e.} \\ \begin{array}{r} 0010 \quad 11 \\ + \quad 1011 \\ \hline 1101 \quad 11 \end{array} \end{array}$$

Step 6 : Right shift the result

$$x = 0110 \mid 111$$

Step 7 : Since, the b_3 bit of multiplier is 1

$$x = x + 1011$$

$$\begin{array}{r} 0110 \quad 111 \\ + \quad 1011 \\ \hline 10001 \quad 111 \end{array}$$

Step 8 : Right shift the result

$$x = 1000 \mid 1111$$

Example 14: Show multiplication of two numbers 11 and 13 using 4-bit registers.

Solution :

$$(11)_{10} = (1011)_2$$

$$(13)_{10} = (1101)_2$$

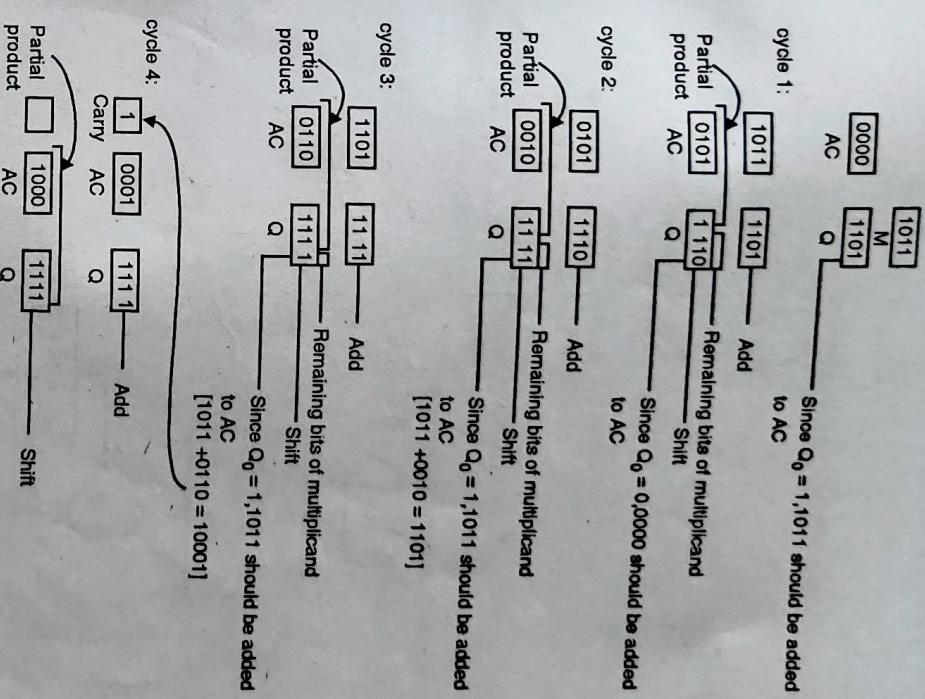
Let us assume existence of the following 4-bit registers :

AC → Accumulator

M → Multiplicand

Q → Multiplier

Algorithm discussed above can be implemented with the help of three registers AC, M and Q.



Two 4-bit numbers are multiplied in four cycles.

If Q_0 is 1 then M is added to AC and combined AC, Q register is shifted right by 1 bit.

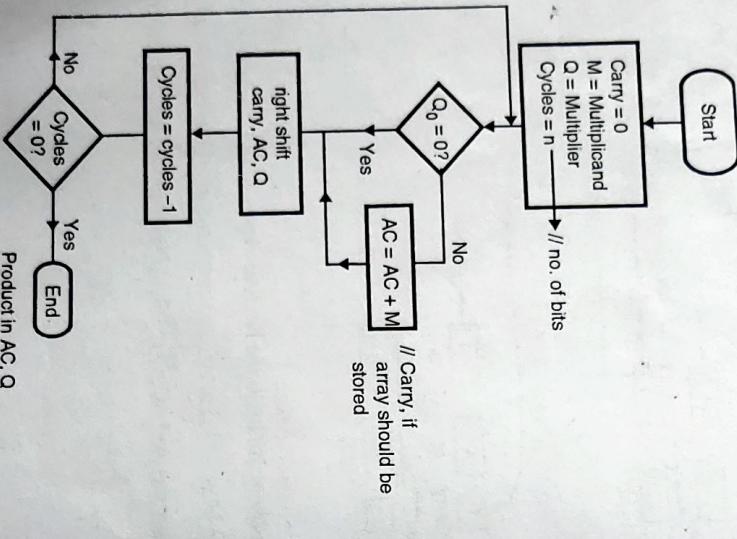
If Q_0 is 0 then combined AC, Q register is shifted right by 1 bit (Addition not required).

Example 15 :
Show multiplication of two numbers 13 and 14 using 4-bit registers.

Solution :

Carry	AC	Q	
0	0000	1110	Initial values
0	0000	0111	shift] 1 st cycle since Q ₀ =0, no addition is required
0	1101	0111	Add(M)] 2nd cycle since Q ₀ =1, at the end of first cycle, M is added
0	0110	1011	Add(M)] 3rd cycle since Q ₀ =1, at the end of 2nd cycle, M is added
1	0011	1011	AC = AC+M // no. of bits
0	0110	1101	AC = AC+M
1	0110	1101	AC = AC+M Add(M)] 4th cycle shift since Q ₀ =1, at the end of 3rd cycle, M is added
0	1011	0110	

Flowchart for unsigned binary multiplication :



Hardware implementation of unsigned binary multiplication :

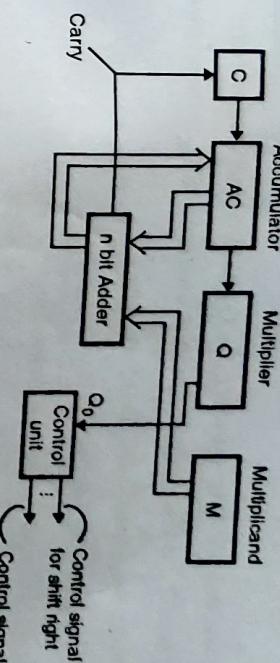


Fig. 2.10 (b)

- Control unit checks the bits Q₀. If the Q₀ bit is found to be 1 then : control unit generates the required control signals to perform AC = AC + M and subsequently, C, AC, Q registers are right shifted by 1 bit.
- If the Q₀ bit is found to be 0 then : control unit generates the required control signals to right shift C, AC, Q registers.

2.3.4 Binary Multiplication (Signed Number) using Booth's Algorithm :

Straight forward multiplication does not work with negative numbers. The same can be understood through an example :

$$\begin{array}{r}
 1001 \quad (9) \\
 \times 0101 \quad (5) \\
 \hline
 00001001 \quad (1001) \times 2^0 \\
 00100100 \quad (1001) \times 2^1 \\
 \hline
 00101101 \quad (45)
 \end{array}$$

$$\begin{array}{r}
 0111 \quad (-9) \\
 \times 0101 \quad (-5) \\
 \hline
 11110111 \quad (-9) \times 2^0 \\
 11011100 \quad (-9) \times 2^1 \\
 \hline
 11010011 \quad (-45)
 \end{array}$$

(a) Unsigned integers
(b) 2's complement integers

Fig. 2.11 : Comparison of multiplication of unsigned and 2's complement integers

Multiplication can not be used directly in the manner it is used for unsigned numbers. One of the most common algorithms for multiplication 2's complement numbers is Booth's algorithm.

Fig. 2.10 (a)

Product in AC, Q

2.3.4.1 Booth's Algorithm :

May 2005, May 2007

- Multiplicand and multiplier are placed in the Q and M registers respectively.
- A 1-bit register Q_{-1} is placed right of the least significant bit Q_0 of the register Q.
- The final result will appear in AC and Q registers.
- AC and Q_{-1} registers are initialized to 0.
- Multiplication of number is done in n cycles.

In each cycle Q_0 and Q_{-1} bits are examined:

- If Q_0 and Q_{-1} are (1 - 1 or 0 - 0) then all the bits of AC, Q and Q_{-1} registers are shifted to the right by 1-bit.
- If Q_0 and Q_{-1} are 01 then multiplicand is added to AC. After addition AC, Q and Q_{-1} registers are shifted to the right by 1 bit.
- If Q_0 and Q_{-1} are 10 then multiplicand is subtracted from AC. After subtraction AC, Q and Q_{-1} registers are shifted to the right by 1-bit.

Booth's algorithm uses the concept of arithmetic right shift. In arithmetic right shift, the left most bit of AC is not only shifted right by 1 bit but it also remains in the original position.

Example 16 :

Right shift contents of AC, Q, Q_{-1} .

Solution :

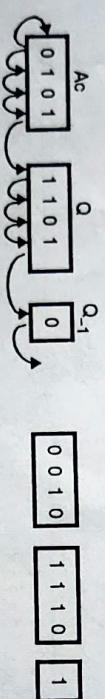


Fig. (a) Positive number

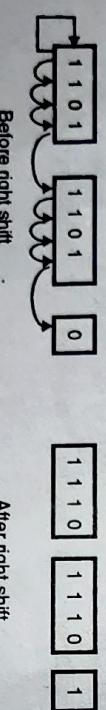


Fig. (b) Negative number

Fig. 2.12 : Arithmetic right shift operation on signed numbers

2.3.4.2 Flowchart of Booth's Algorithm :

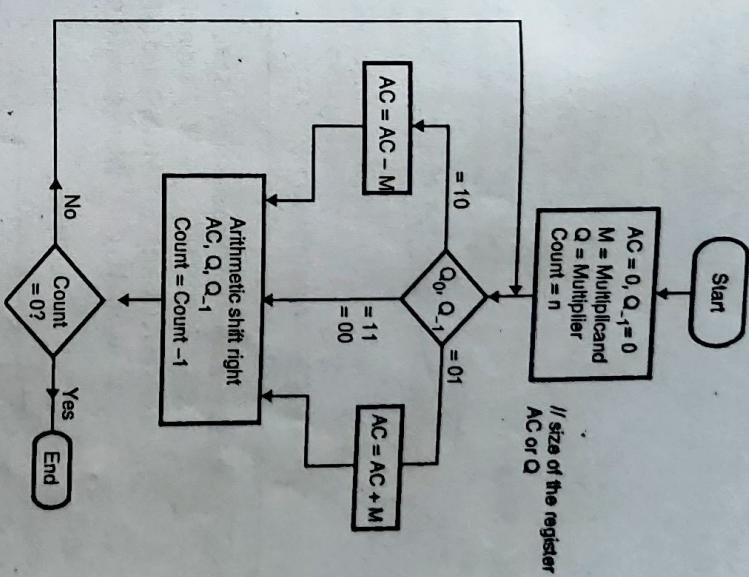


Fig. 2.13 : Flowchart of Booth's algorithm for 2's complement multiplication

Example 17 :

Multiply 7 and 3 using Booth's algorithm. Register size = 4-bits.

Solution :

$$(7)_{10} = (0111)_2 \quad \text{Multiplicand (M)}$$

$$(-7)_{10} = (1001)_2 \quad 2\text{'s complement of } (7)_{10} (-M)$$

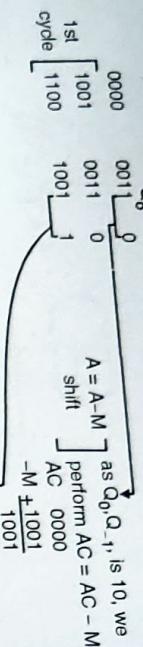
$$(3)_{10} = (0011)_2 \quad \text{Multiplier (Q)}$$

Multiplication will require 4 cycles as the register size n = 4-bit.

AC-M is same as AC + 2's complement of M.

Example 19 : Multiply -7 and -3 using Booth's algorithm. Register size = 4 bits.

Solution :



$$\begin{aligned} (-7)_{10} &= (0111)_2 && \text{(-M) 2's complement of M} \\ (-7)_{10} &= (1100)_2 && \text{Multiplicand (M)} \\ (3) &= (0011)_2 && \text{Multiplier (Q)} \end{aligned}$$

$M = 0111$

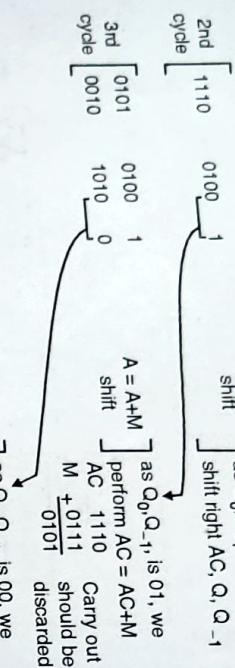
$Q = 1001$

$AC = 0000$

$a = 0$

$a_{-1} = 0$

$initial$



Example 18 : Multiply -7 and 3 using Booth's algorithm. Register size = 5 bits.

Solution : Multiply -7 and 3 using Booth's algorithm. Register size = 4 bits.

$$\begin{aligned} (-7)_{10} &= (0111)_2 && \text{(-M) 2's complement of M} \\ (-7)_{10} &= (1100)_2 && \text{Multiplicand (M)} \\ (3) &= (0011)_2 && \text{Multiplier (Q)} \end{aligned}$$

$M = 0111$

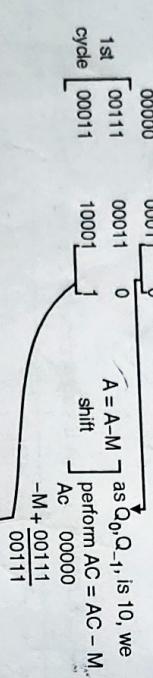
$Q = 1001$

$AC = 0000$

$a = 0$

$a_{-1} = 0$

$initial$



Example 20 :

Explain Booth's algorithm to multiply the following pair of signed 2's complement numbers:

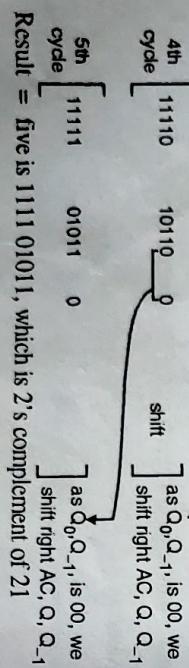
$$\begin{aligned} A &= 110011 \text{ multiplicand} \\ B &= 101100 \text{ multiplier} \end{aligned}$$

Solution :

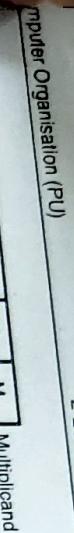
$$\text{Multiplicand (M)} = 110011$$

$$2\text{'s complement of multiplicand} = 001101 = (-M)$$

$$\text{Multiplier } Q = 101100$$



Multiplication will require 6 cycles as the register size $n = 6$ -bits.



Initial setting : $AC \leftarrow 0$ and $Q_{-1} \leftarrow 0$
 $Q \leftarrow$ Multiplexer
 $M \leftarrow$ Multiplicand

Fig. 2.14 : Hardware implementation of Booth's algorithm

- In case of subtraction, Add / sub line is 1, therefore $C_0 = 1$ and multiplicand is complemented and then applied to the n-bit adder. As a result, the 2's complement of multiplicand (M) is added to the register AC.

2.3.5 Bit-pair Recording of Multipliers (A Fast Multiplication Method) :

- Above technique halves the maximum number of summands. Bit-pair recording technique is derived directly from the Booth's algorithm. In Booth's algorithm, multiplier is examined two bits at a time, starting from the right. Table given below gives operation based on pair of bits.

$Q_0 \quad Q_1$	Operation	Notation
0 0	No add/subtract	0
1 1	No add/subtract	0
0 1	Add	+ 1
1 0	Subtract	- 1

Fig. 2.15 : Operations as per Booth's algorithm

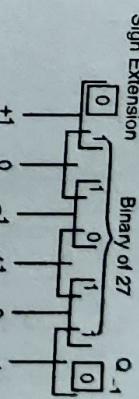
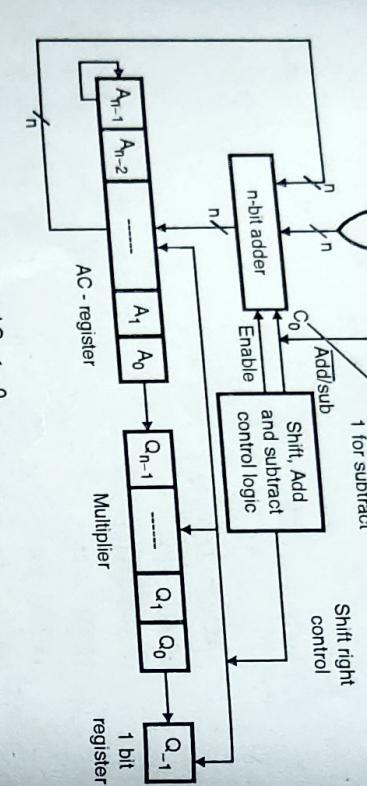


Fig. 2.16 : Operations for the multiplier = 27, using 0, +1 and -1 notation

In bit-pair recording, two adjacent operations of Booth's algorithm are combined together.



Let us try to understand operations $+2, -2, +1, -1, 0$ on the data $M = 9$ and register size $n = 5$ bits.

$$M = 9 = (01001)_2 = \underbrace{00000 \ 01001}_{\text{10 bit representation of '9'}} \quad \text{as the register size = 5}$$

$-9 = 2$'s complement of 9 = 11111 10111

Operation Value (9) Multiplicand

Operation	Value (9)	Multiplicand
0	00000 00000	
+ 1	00000 01001	
- 1	11111 10111	
+ 2	00000 10010	\leftarrow data is left shifted
- 2	11111 01110	$\leftarrow -9$ is left shifted

'0' stands for no operation. + 1 stands for addition and - 1 stands for subtraction.

$$\text{Multiplicand } M = \underbrace{11111}_{\sim} 110011$$

As $n = 6$, number is represented in 12 bits

As the number given is a negative number, it is extended by 1's.

$\neg M \equiv 2\$$ complement of $M = 00000001101$

Operation	Value (M)
0	000000 000000
+1	111111 100111
-1	000000 001101
+2	111111 100110
-2	000000 011010

)
Results
2's complement
of 54

A truth table for a 4-to-1 multiplexer. The inputs are labeled Q = [1 0 1 1] and S = [0 0 0 1]. The outputs are labeled Y = [0 1 1 0]. The table shows how the selected input line (S) connects to the output Y based on the address inputs.

S	Y
0	0
1	1
2	1
3	0

$$\begin{array}{r} \text{↑} \\ - \\ \hline 0 & 1 & -1 & 0 \\ 0 & -1 & +1 & 0 \\ \hline 1 & 0 & 0 & 1 \end{array} \quad \text{↑}$$

Pairing of operations

	0	1	0	0	1	1	Multiplicand
	0	-1	-2				
	0	1	0	0	1		
	0	-1	-2				
	0	0	0	0	0	0	
	0	0	0	0	0	1	
	0	0	0	0	0	1	
	0	0	0	0	0	0	
	+	0	0	0	1	1	
	0	0	0	1	0	0	
	0	0	0	1	0	1	
	0	0	0	1	0	0	
	+	0	0	0	1	0	
	0	0	0	1	0	0	
	0	0	0	1	0	0	
	=	(-54)					

Example 23:

Multiply the following numbers using bit-pair recording method.

Exodus 3

Implement multiplication of the following pair of signed 2's complement numbers using bit-pair recording

$A = 110011$ multiplicand

B = 101100 millijar

(May 2005)

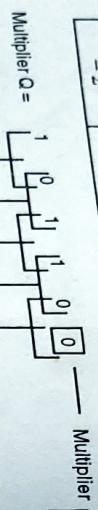
Solution:

Multiplicand	01111 (15)
Multiplier	10110 (- 10)

(May 2007)

Computer Organisation (PU)

Operation	Value (M)
0	00000 00000
+1	00000 01111
-1	11111 10001
+2	00000 11110
-2	11111 00010



Single value can be taken as it is → -1
+1 0 -1 0
-1 +2 -2

0 1 1 1 – Multiplicand

x
1 1 1 1 0 0 1 0
0 0 0 1 1 1 0
+ 1 1 0 0 0 1
—————

-1 +2 -2

M is left shifted ← M is left shifted

← - M is left shifted

- Bits of the dividend are examined from left to right, until the set of bits examined represents a number greater than or equal to divisor. Until this event occurs, 0's are placed in the quotient from left to right.
- When the divisor becomes greater than or equal to the group of bits of dividend under examination, a 1 is placed in the quotient and the divisor is subtracted from the partial dividend.
- The result so produced is referred to as a partial remainder.
- Additional bits from the dividend are appended to the partial remainder until the result is greater than or equal to the divisor.
- Now, the divisor is subtracted from this number to produce a new partial remainder.
- The process continues until all the bits of the dividend are exhausted.

2.3.6.1 Restoring Division :

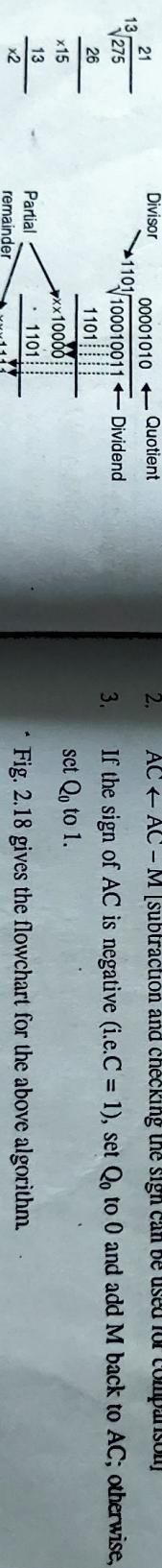
Division algorithm can be implemented using a machine.

Algorithm :

Register, M ← Divisor Register, Q ← Dividend Register, AC ← 0 Register, C ← 0 [1-bit carry register] Count ← n [size of the register]	Initial values
---	----------------

2.3.6 Integer Division :

Division is carried out on the same principle as integer. But it is some what more complex. The basis for the algorithm is the paper and pencil approach, and the operation involves repetitive shifting and subtraction. Fig. 2.17 shows an example of division of unsigned binary integers.



(a) Decimal numbers (b) Binary numbers

Fig. 2.17 : Paper and pencil approach for division

Dec. 2005, Dec. 2006, May 2007

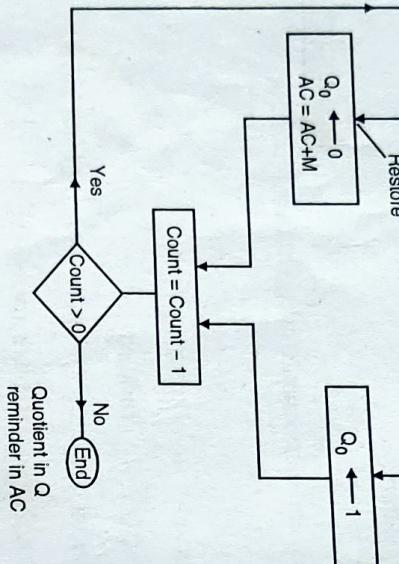
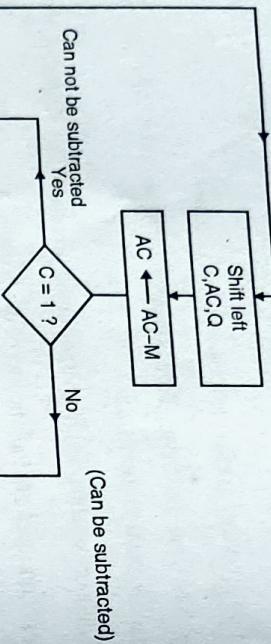
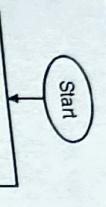


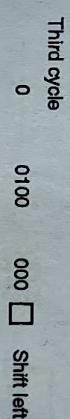
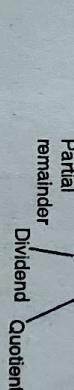
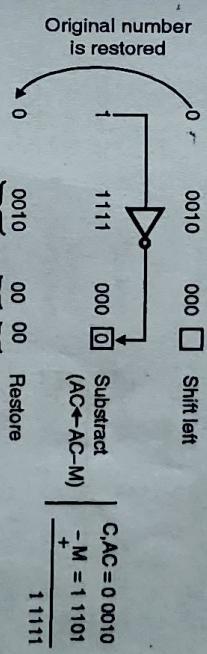
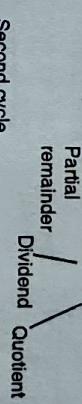
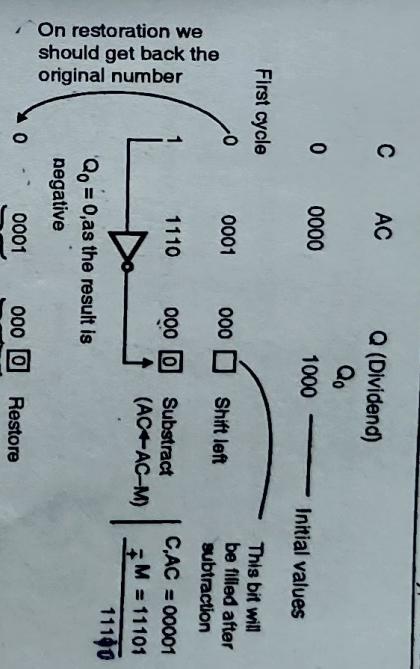
Fig. 2.18 : Flowchart for unsigned binary division (Restoring division)

Perform division of the following positive numbers using restoring division.
 $1000 \div 11$, register size, $n = 4$ bits.

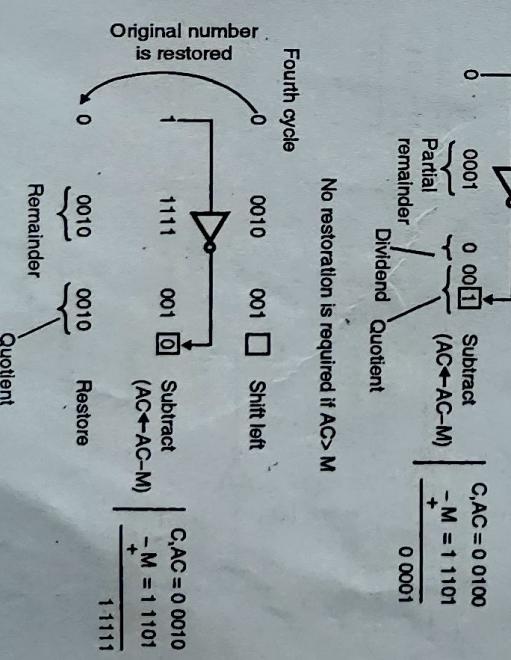
Solution :

$$\text{Divisor } M = \underbrace{0001}_{n+1 \text{ bits}} \quad \text{and} \quad -M = \underbrace{11101}_{n+1 \text{ bits}}$$

One additional bit has been taken to handle borrow. Borrow (C) will be set whenever $M > AC$ and we perform $AC = AC - M$.



No restoration is required if $AC > M$



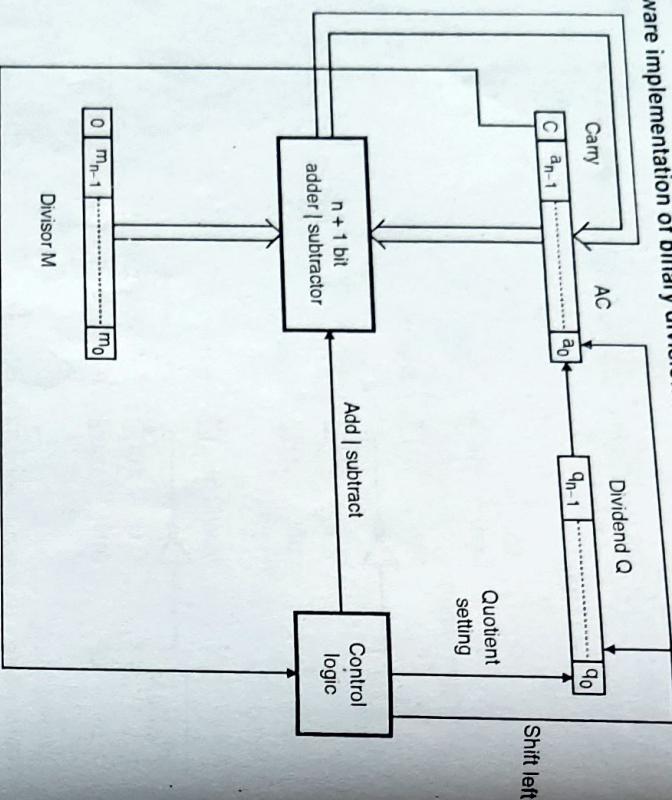
Hardware implementation of binary division :

Fig. 2.19 : Circuit for binary division [Both restoring and non-restoring]

- An n-bit positive divisor is loaded into register M .
- An n-bit positive dividend is loaded into register Q .
- Register $AC(A)$ is set to 0.
- Initial carry C is set to 0.
- After the division is complete, the n-bit quotient is in register Q and the remainder is in register AC .

2.3.6.2 Non-restoring Division Technique :

May-2006, Dec-2006

The algorithm of restoring division can be improved by avoiding restoring after an unsuccessful subtraction. Subtraction is said to be unsuccessful if the result is negative.

Algorithm for non-restoring division :

Step 1 : Do the following n times :

- If the sign of AC is positive ($C = 0$), then shift C , AC and Q left one bit position and subtract M from AC .
- else
 - Shift C , AC and Q left one bit position and add M to AC .
 - If the sign of AC is positive then
 - Set q_0 to 1
 - else
 - Set q_0 to 0

Step 2 : If the sign of AC is negative then add M to AC .
Step 2 is needed to leave the proper positive remainder in AC at the end of n cycles.

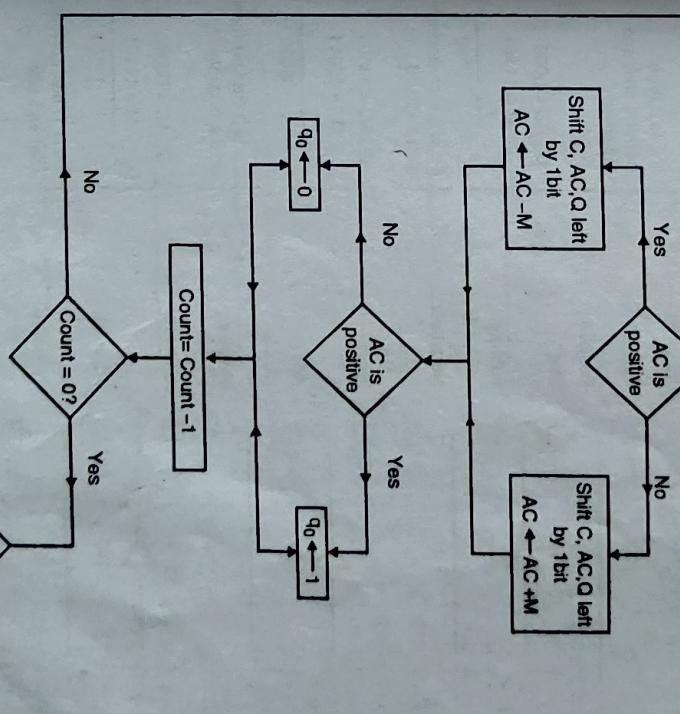
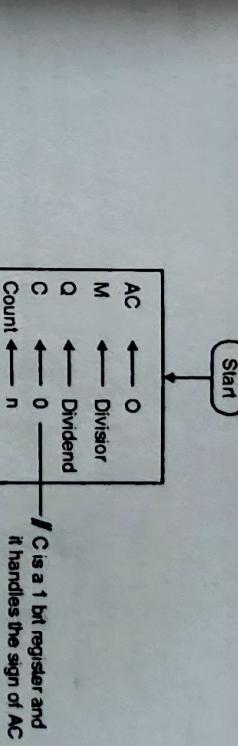
Flowchart for non-restoring division operation :

Fig. 2.20 : Flowchart for non-restoring division

Flowchart for non-restoring division is given in the Fig. 2.20. Logic circuit for both restoring and non-restoring division can be handled on the circuit given in Fig. 2.19.

Example 29 :

Perform division of the following numbers using non-restoring division.

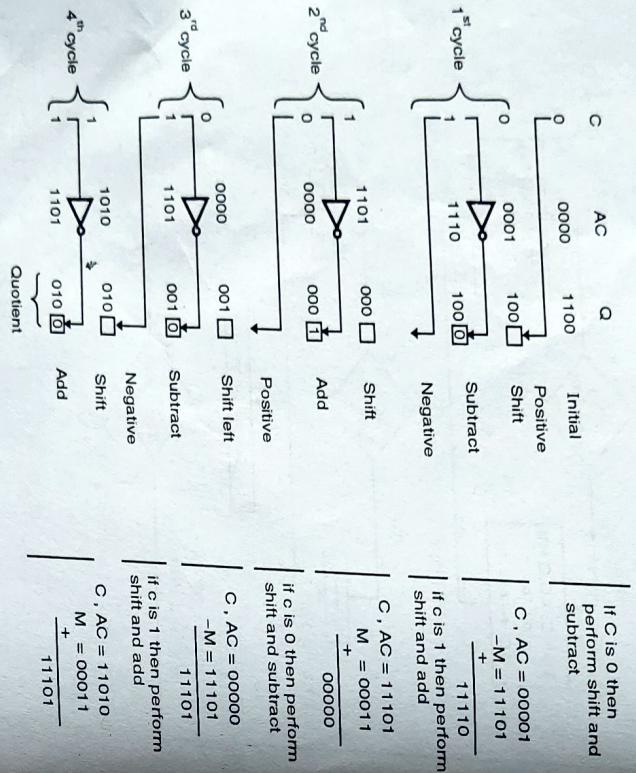
$$1100 \div 11$$

Solution : Dividend Q = 1100

$$\text{Divisor } M = \underbrace{0\ 0011}_n 1,$$

n + 1 bits – one additional bit to handle carry/borrow/sign

$$-M = 11101$$



Quotient = 0011 = 3
Remainder = 0010 = 2

2.4 Floating Point Numbers :

The range of numbers that can be represented by a fixed-point number is insufficient for many applications. In scientific applications, very large and very small numbers are encountered. Scientific notation permits us to represent such numbers using relatively few digits.

For example,

$$2.5 \times 10^{10}$$

represents a fixed point integer 2500000000.

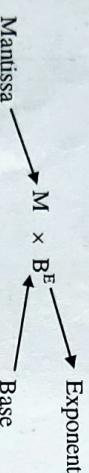
The floating-point codes used in computers are represented in binary.

Format of representation :

Three numbers are associated with a floating point number :

- (1) A mantissa (M)
- (2) An exponent (E)
- (3) Base (B)

The mantissa M is also referred to as the significant or fraction. These three components together represent the real number.



For example, in 2.5×10^{10}

Mantissa = 2.5

Exponent = 10

Base = 10

0	1	8	9	31
Sign	Biased exponent = 8 bit	Significant = 23 bits		

Fig. 2.21 : Floating point number representation (Binary)

In a typical representation of a floating point number :

- (1) Exponent is biased.
- (2) Mantissa or significand is normalized.

Biased exponent :
Exponent is represented using an excess-128 code. An 8 bit binary number represents values from 0 to 255. However, as we are adding 128 in the biased exponent, the actual exponent values represented will be -128 to 127.

Exponent value	Exponent value after biasing (+128)
-128	0 (000 0000) ₂
-127	1 (0000 0001) ₂
0	128 (1000 0000) ₂
127	255 (1111 1111) ₂

Advantages of biasing :

Very large or small exponents can easily be represented.

For example,

If the exponents are in the -1050 to -900, we can select the biasing factor as 1050.

With 1050 as biasing,
-1050 will be represented as -1050 + 1050 = 0
- 900 will be represented as -900 + 1050 = 150

Normalized mantissa :

A binary floating point number is represented in a normalized form, that is, the number is of the form ± 0 . (Significand starting with non-zero bit) $\times 2^{\pm} (\text{exponent value})$.

Example,

Binary number	Its normal form
11.101	.11101 $\times 2^2$
.00101	.101 $\times 2^{-2}$
1.01×2^5	.101 $\times 2^6$

Advantages of normalization :

As for a normalized mantissa, the left most bit can not be zero, therefore, it has to be 1. Thus, it is not necessary to store this first bit and it is assumed implicitly for the number. Therefore, a 23 bit mantissa can represent $23 + 1 = 24$ bit significand.

Example 31 :

Represent $(13.54)_{10}$ in 32 bit register with
Exponent = 8 bits

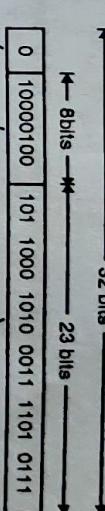
Mantissa = 23 bits

Exponent is biased with biasing of 128 and Mantissa is normalized.

Solution :

Step 1 : Converting 13.54 to its equivalent binary form.

$$\begin{aligned}13 &= 1101 \\.54 \times 2 &= 1.08 \\.08 \times 2 &= 0.16 \\.16 \times 2 &= 0.32 \\.32 \times 2 &= 0.64 \\.64 \times 2 &= 1.28 \\.28 \times 2 &= 0.56 \\.56 \times 2 &= 1.12 \\.12 \times 2 &= 0.24 \\.24 \times 2 &= 0.48 \\.48 \times 2 &= 0.96 \\.96 \times 2 &= 1.92 \\.92 \times 2 &= 1.84 \\.84 \times 2 &= 1.68 \\.68 \times 2 &= 1.36 \\.36 \times 2 &= 0.72 \\.72 \times 2 &= 1.44 \\.44 \times 2 &= 0.88 \\.88 \times 2 &= 1.76 \\.76 \times 2 &= 1.52 \\.52 \times 2 &= 1.04 \\.04 \times 2 &= 0.08 \\.08 \times 2 &= 0.16 \\.16 \times 2 &= 0.32 \\.32 \times 2 &= 0.64\end{aligned}$$



Representation of $(13.54)_{10}$ in floating point format.

Disadvantages of normalisation :

Minimum value of significand :

The implicit first bit as 1 followed by 23 0's.

$1000\ 0000\ 0000\ 0000\ 0000\ 0000$ (1 is hidden)

Decimal equivalent = $1 \times 2^{-1} = 0.5$

Maximum value of significand :

The implicit first bit 1 followed by 23 1's.

$1111\ 1111\ 1111\ 1111\ 1111$

Decimal equivalent :

$$\begin{array}{r} \text{Binary :} & 0.1111\ 1111\ 1111\ 1111\ 1111 \\ & + 0.0000\ 0000\ 0000\ 0000\ 0001 = 2^{-24} \\ \hline & 1.0000\ 0000\ 0000\ 0000\ 0000\ 0000 \end{array}$$

\therefore Maximum value of significand = $1 - 2^{-24}$

Therefore, in normalized mantissa and biased exponent form, the format of Fig. 2.21 can represent binary floating point number in the range.

Lowest negative number : Maximum significand and maximum exponent

Step 2 : Biasing

$$\therefore 13.54 = 1101.100010100011110101110000$$

$$= .1101100010100011110101110000 \times 2^4$$

Exponent = 4

Exponent after biasing = $128 + 4$

$$= 132$$

$$(132)_{10} = (10000100)_2$$

Highest negative numbers = Minimum significand and minimum exponent

$$= -0.5 \times 2^{-128}$$

Lowest positive number : 0.5×10^{-128}

Highest positive number : $(1 - 2^{-24}) \times 10^{127}$

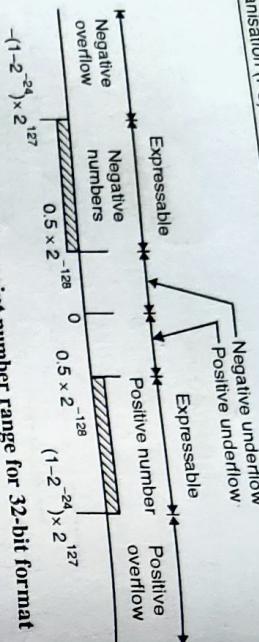


Fig. 2.22 : Binary floating point number range for 32-bit format

From Fig. 222 it is clear that 0 can not be represented.

- (1) From the Fig. 2.2, it can be seen that the range of floating point numbers is $-0.5 \times 2^{-128} <\!\!> +0.5 \times 2^{128}$. It can not represent values between -0.5×2^{-128} and $+0.5 \times 2^{128}$.

(2) A number in the range $0.5 \times 2^{-128} <\!\!> -0.5 \times 2^{128}$ can cause an underflow.

(3) ∞ can not be represented.

2.4.1 IEEE Standards

May 2003, Dec. 2003, May 2007

- IEEE floating point standards addresses a number of such problems.
 - Zero has definite representation in IEEE format.
 - $\pm \infty$ has been represented in IEEE format. A $+\infty$ indicated that the result of an arithmetic expression is too large to be stored.
 - If an underflow occurs, implying that a result is too small to be represented as a normalized number, it is encoded in a denormalized scale.
 - Fig. 2.23 gives the representation of floating point numbers.

Fig. 2.23 gives the representation of floating point numbers.

178 = 10111100

Exponent (E)	Significand (N)	Value/Comments
255	Not equal to 0	Does not represent a number
255	0	$-\infty$ or $+\infty$ depending on sign bit
Normalized scale	$0 < E < 255$	$\pm(1.N) \cdot 2^{E-127}$
Denormalized scale	0	$\pm(0.N) \cdot 2^{-128}$
0	0	± 0 depending on sign bit

Double precision (64 bits) :

Exponent (E)	Significand (N)	Value/comments
2047	Not equal to 0	Does not represent a number
2047	0	$-\infty$ or $+\infty$ depending on sign bit
Normalized scale	$0 < E < 2047$	$\pm(1.N) \cdot 2^{E-1023}$
Denormalized scale	0	$\pm(0.N) \cdot 2^{-1022}$
0	0	± 0 depending on sign bit

Fig. 2.24 : Values of floating point numbers as per IEEE format

Represent $(178.1875)_{10}$ in single and double precision floating point format. (Dec. 2005)

Convert given decimal number into its equivalent binary

Convert given decimal number into its equivalent binary

178 = 1011100

- base = 2
 - Significand is in normalized form i.e the first bit is 1 and it is hidden.
 - S is sign bit.
$$.3750 \times 2 = 0.7500$$

$$.7500 \times 2 = 1.500$$

$$.500 \times 2 = 1.000$$

0 1
11
12

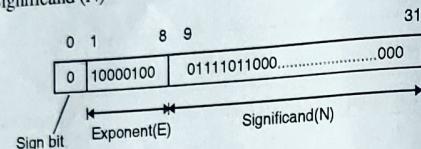
Sign bit

Fig. 2.23 : IEEE standard format

$$101111.011 = 1.01111011 \times 2^5$$

(Decimal is shifted left by five places)

$$= 1.0\underbrace{1111011}_{\text{Significand (N)}} \times 2^{\underbrace{132 - 127}_{\text{exponent (E)}}}$$

**(b) Double precision format :**

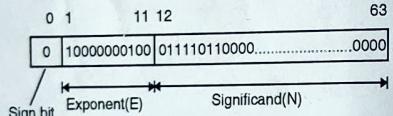
In IEEE double precision format, the value of a number for given exponent (E) and significand (N) is given by $(1.N) 2^{E-1023}$.

In order to represent $(10111100.0011)_2$, we must convert it into the form $(1.N) 2^{E-1023}$.

$$101111.011 = 1.01111011 \times 2^5 = 1.01111011 \times 2^{\underbrace{1028 - 1023}_{\text{Exponent (E)}}}$$

Significand (N) Exponent (E)

$$(1028)_{10} = (1000000100)_2$$

**Example 33 :**

Represent $(309.1875)_{10}$ in single precision and double precision format. (May 2006)

Solution : Convert given decimal number into its equivalent binary.

$$(309)_{10} = (100110101)_2$$

$$.1875 \times 2 = 0.3750$$

$$.3750 \times 2 = 0.7500$$

$$.7500 \times 2 = 1.5000$$

$$.5000 \times 2 = 1.0000$$

$$\therefore 309.1875 = 100110101.0011$$

(a) Single precision format :

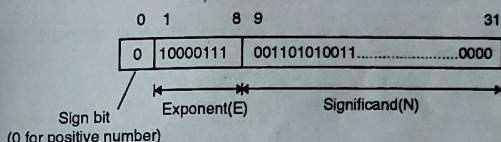
In IEEE single precision format, the value of a number for given exponent (E) and significand (N) is given by $(1.N) \times 2^{E-127}$.

In order to represent 100110101.0011 , we must convert it into the form $(1.N) \times 2^{E-127}$.

$$100110101.0011 = 1.001101010011 \times 2^8$$

$$= 1.0\underbrace{01101010011}_{\text{Significand (N)}} \times 2^{\underbrace{135 - 127}_{\text{Exponent (E)}}}$$

$$(135)_{10} = (10000111)_2$$

**(b) Double precision format :**

In IEEE double precision format, the value of a number for given exponent (E) and significand (N) is given by,

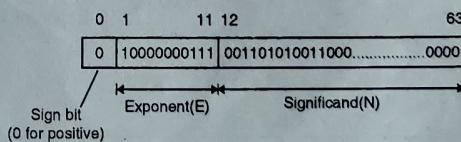
$$(1.N) \times 2^{E-1023}$$

In order to represent 100110101.0011 , we must convert it into the form $(1.N) \times 2^{E-1023}$.

$$100110101.0011 = 1.001101010011 \times 2^8$$

$$= 1.0\underbrace{01101010011}_{\text{Significand (N)}} \times 2^{\underbrace{1031 - 1023}_{\text{Exponent (E)}}}$$

$$(1031)_{10} = (10000000111)_2$$



3.4.2 Floating Point Arithmetic :

A number of operations are performed on floating point numbers.

- A number of operations are possible:

 - (1) Addition
 - (2) Subtraction
 - (3) Multiplication
 - (4) Division

(1) Addition (2) Subtraction
 (3) Multiplication (4) Division

For addition and subtraction, it is necessary that both the numbers must have the same exponent. This may require shifting the decimal point of one of the number to achieve alignment. Consider a decimal example in which we wish to add 2.4500×10^2 to 4.5100×10^5 .

We rewrite 2.4500×10^2 as

$.0024500 \times 10^5$ and then perform the addition.

$$\begin{array}{r} 4.5100 \times 10^5 \\ + 0.0245 \times 10^5 \\ \hline 4.6145 \times 10^5 \end{array}$$

Multiplication and division are straight forward. Problems may arise as the result of these operations. These are :

- Exponent overflow
 - Exponent underflow
 - Significand underflow
 - Significand overflow

2.4.2.1 Addition and Subtraction :

May 2005, May 2006, Dec. 2006

In floating point arithmetic, addition and subtraction are more complex than multiplication and division. Addition and subtraction operations are carried out in four basic phases :

- (1) Check for zeros. (2) Align the significands
 (3) Add or subtract the significands (4) Normalize the result

First number X = M × B^x

Second number $Y = N \times B^y$

Arithmetic operations :

$$\left. \begin{array}{l} X + Y = (M \times B^{x-y} + N) B^y \\ X - Y = (M \times B^{x-y} - N) B^y \end{array} \right\} \text{Exponent } x \leq \text{exponent } y$$

$$X * Y = (M * N) B^{x+y}$$

$$X \div Y = (M \div N) \times B^{X-Y}$$

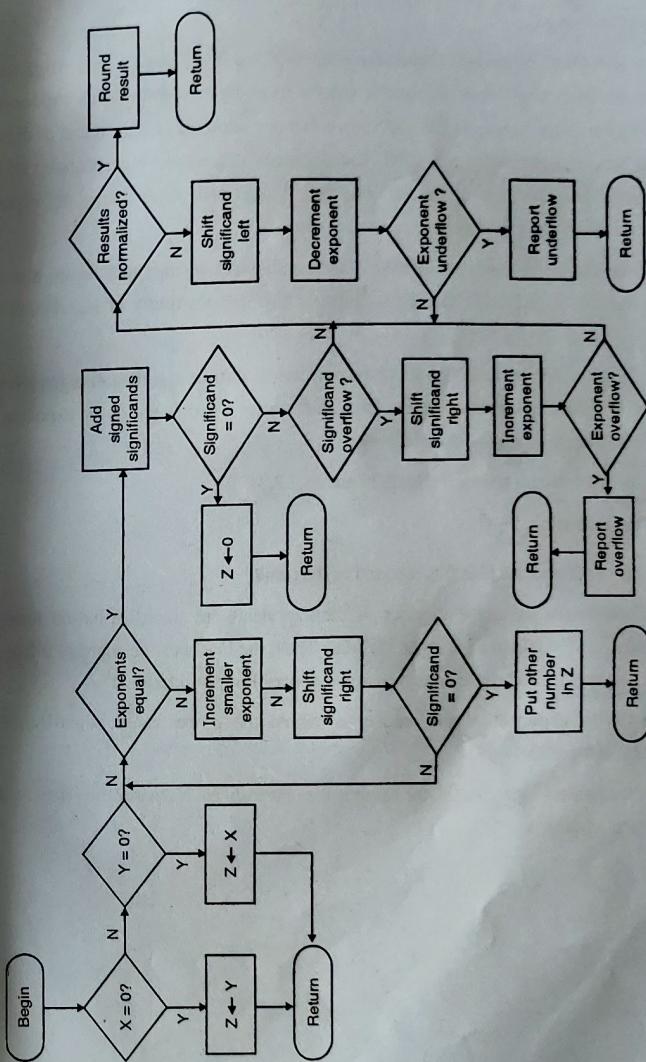


Fig. 2.25 : Floating point numbers and arithmetic operations

- Addition and subtraction are almost similar operations. In subtraction, sign of the second number is changed.
- In the next phase, exponents of the two numbers X and Y are made equal. Alignment is achieved by shifting either the smaller number to its right (increasing its exponent) or shifting the larger number to the left. Since either operation may result in loss of digits, it is the smaller number that is shifted. The alignment is achieved by repeatedly shifting the magnitude portion of the significand right 1 digit and incrementing the exponent until the two digits exponents are equal.
- Next, the two significands are added together, taking into account their signs. Since the sign may differ, the result may be 0. There is also the possibility of significand overflow.
- Next, result is normalized. Normalization consists of shifting significand digits left until the most significant digit is nonzero. Each shift causes a decrement of the exponent and thus could cause an exponent underflow.
- Detailed flowchart is shown in Fig. 2.26 and Fig. 2.27.

2.4.2.2 Multiplication :

- If either of the operand is 0, 0 is reported as the result.
- Next, exponents are added together. If the exponents are stored in biased form, the exponent sum would have doubled the bias. Thus, the bias must be subtracted from the sum. The result could cause either an exponent overflow or underflow.
- Next, if the exponent of the product is within the proper range, significands are multiplied together.
- After the product is calculated, the result is then normalized. Normalization may result in exponent underflow

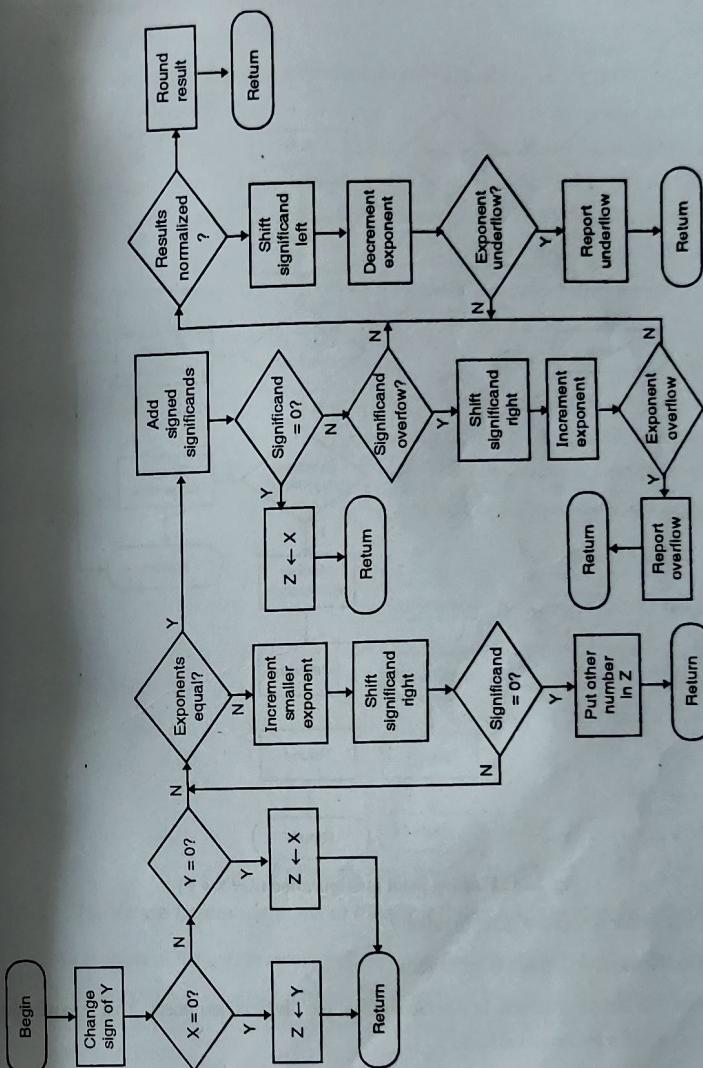
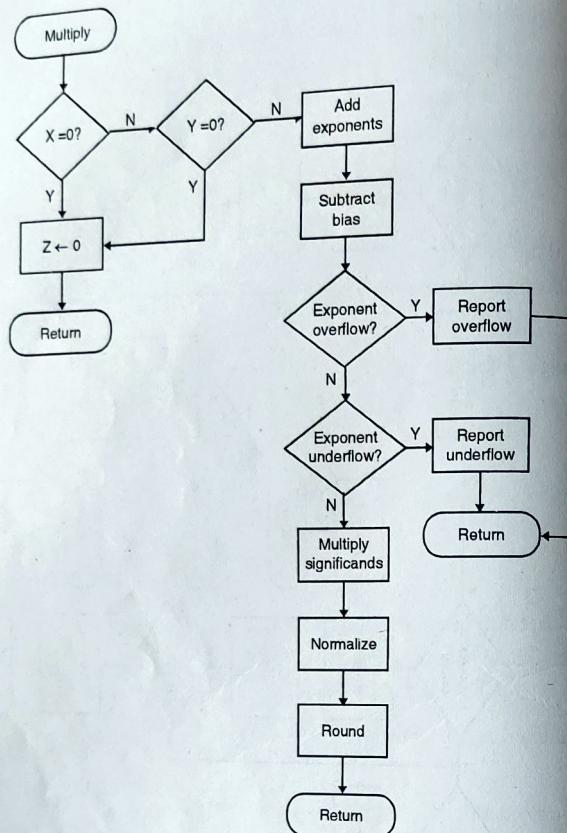


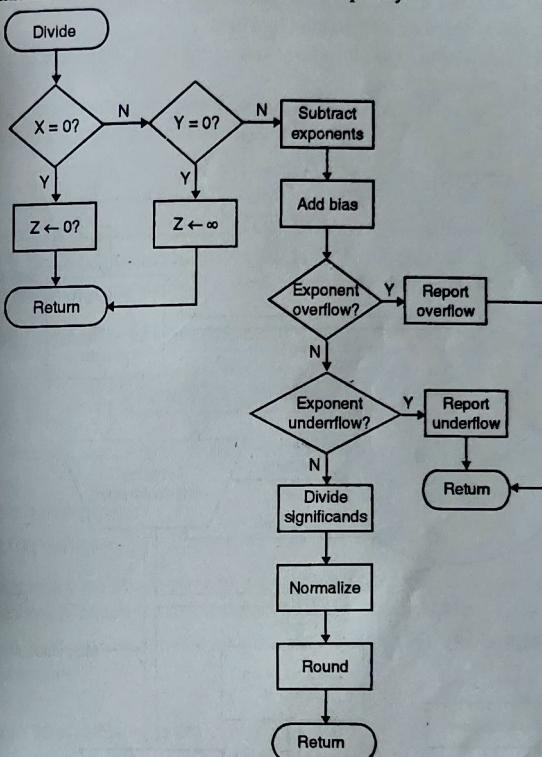
Fig. 2.27 : Floating point subtraction ($Z = X - Y$)

2.4.2.3 Division :

Fig. 2.28 : Floating point multiplication ($Z = X * Y$)

- If the divisor is 0, result is set to infinity.
- If the dividend is 0, result is set to zero.
- Next, the divisor exponent is subtracted from the dividend exponent. This removes the bias, which must be added back.

- Exponent is checked for underflow or overflow.
- Next, significands are divided.
- Normalization and rounding are carried out subsequently

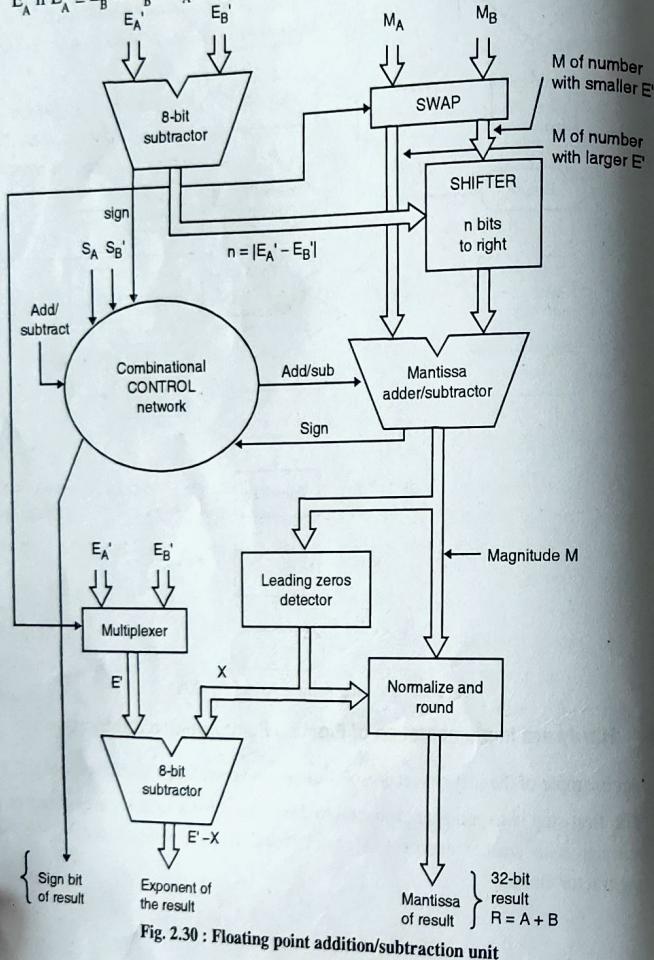
Fig. 2.29 : Floating point division ($Z = X/Y$)

2.4.3 Hardware Implementation of Floating Point Addition/Subtraction :

An example of the implementation of floating point operation is shown in the Fig. 2.30.

- The first step is to compare exponents to determine how far to shift the mantissa of the number with smaller exponent. The shift count value n is determined by the 8-bit subtractor circuit.

- The magnitude of the difference $E'_A - E'_B$ sent to the SHIFTER unit.
- The sign of the difference that results from comparing exponents determines which mantissa is to be shifted. Therefore, the sign is sent to the SWAP unit.
- One of the mantissas is sent directly to the adder/subtractor whereas other mantissa is sent to the adder/subtractor after shifting by n bits.
- Multiplexer unit is used to select the exponent of the result. Exponent of the result E' is E'_A if $E'_A \geq E'_B$ or E'_B if $E'_A < E'_B$.



- A : S_A = Sign bit of the number A
 E'_A = Exponent of number A
 M_A = Mantissa of number A
- B : S_B = Sign bit of the number B
 E'_B = Exponent of number B
 M_B = Mantissa of number B

- The control logic determines whether the mantissas are to be added or subtracted. This is decided by the signs of the operands (S_A and S_B) and the operation (Add or subtract) that is to be performed on the operands.
- After Add/subtract, normalization is carried out on the mantissa M. The number of leading zeroes in M determines the number leading zeroes in M determine, the number of bit shifts, X, to be applied to M. The value of X is also subtracted from the tentative result exponent E' to generate the true result exponent, E'_R .

2.5 University Questions and Answers :

May 2005 : Total Marks 32

- Q. 1** Explain Booth's algorithm to multiply the following pair of signed two's complement numbers :

$$A = 110011 \text{ multiplicand}$$

$$B = 101100 \text{ multiplier}$$

Also, implement the above using bit-pair recording and explain how it achieves faster multiplication. (Section 2.3.4.1, Examples 20 and Example 22) (12 Marks)

- Q. 2** Compare restoring and non-restoring division algorithm with the help of the following binary numbers

$$1100 \div 11$$

(Example 26 and Example 29) (12 Marks)

- Q. 3** Explain IEEE floating point number formats. (Section 2.4.1) (4 Marks)

- Q. 4** Draw the flowchart for floating point addition. (Section 2.4.2.1) (4 Marks)

Dec. 2005 : Total Marks 26

- Q. 5** Draw the flowchart for restoring division algorithm and solve the following using above algorithm.

$$\text{Dividend} = 17$$

$$\text{Divisor} = 3$$

(Section 2.3.6.1 and Example 25) (10 Marks)

Q. 6 Represent $(178.1875)_{10}$ in single and double precision floating point format.
(Example 32)

Q. 7 Using Booth's algorithm multiply the following :

Multiplicand = + 15

Multiplier = - 6

(Example 21)

May 2006 : Total Marks 32

Q. 8 Draw the hardware implementation of Booth's algorithm and explain the same.
(Section 2.3.4.2)

Q. 9 Represent $(309.1875)_{10}$ in single precision and double precision format.
(Example 33)

Q. 10 Draw the flowchart for non-restoring division algorithm and perform division of the following numbers using non-restoring division algorithm :

Dividend = 1011

Divisor = 0011

(Section 2.3.6.2 and Example 30)

Q. 11 Draw the flowchart for floating point addition and explain. **(Section 2.4.2.1)** (6 Marks)

Dec. 2006 : Total Marks 36

Q. 12 Compare restoring and non-restoring division algorithm. Perform division of the following numbers using restoring division algorithm :

Dividend = 1101

Divisor = 11

(Sections 2.3.6.1, 2.3.6.2 and Example 26)

(12 Marks)

Q. 13 Describe the IEEE standards for single and double precision floating point numbers. **(Section 2.4.1)** (6 Marks)

Q. 14 Multiply the following pair of signed two's complement numbers using Booth's algorithm :

A : 110011

B : 101100

(Example 20)

Q. 15 Explain and draw the flowchart for floating point subtraction. **(Section 2.4.2.1)** (12 Marks)

(6 Marks)

May 2007 : Total Marks 24

Q. 16 Draw the flowchart for Booth's algorithm and solve the following using bit-pair recording method :

Multiplicand = 01111

Multiplier = 10110

(Section 2.3.4.1 and Example 23)

(8 Marks)

Q. 17 Draw flowchart for restoring division algorithm and perform division of the following numbers using restoring division algorithm.

A : 1100 B : 0100

(Section 2.3.6.1 and Example 27)

(10 Marks)

Q. 18 Write short note on IEEE standards. **(Section 2.4.1)**

(6 Marks)

□□□