

## Number Systems

A number system of base, or radix  $r$ , is a system that uses distinct symbols for  $r$  digits. Numbers are represented by a string of digit symbols. To determine the quantity that the number represents, it is necessary to multiply each digit by an integer power of  $r$  and then form the sum of all weighted digits.

For example, the decimal number system in everyday use employs the radix 10 system. The 10 symbols are 0, 1, 2, 3, 4, 5, 6, 7, 8 and 9. The string of digits 724.5 is interpreted to represent the quantity

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

ie., 7 hundreds plus 2 tens plus 4 units plus 5 tenths.

Binary number system uses the radix 2. The two digit symbols used are 0 and 1. The string of digits 101101 is interpreted to represent the quantity

$$1 \times 2^5 + 0 \times 2^4 + 1 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 = 15$$

To distinguish between different radix numbers, the number digits will be enclosed in parentheses and the radix of the number inserted as a subscript. For example, to show the equality between decimal and binary forty five will be written as  $(101101)_2 = (45)_{10}$ .

Octal and hexadecimal number system consists of radix 8 and 16 respectively. The eight symbols of the Octal system are 0, 1, 2, 3, 4, 5, 6, and 7. The 16 ~~num~~ symbols of the hexadecimal system are 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E and F. The symbols A, B, C, D, E, F correspond to the decimal numbers 10, 11, 12, 13, 14, 15 respectively.

## Conversion

The conversion of a decimal integer into a base  $r$  representation is done by successive divisions by  $r$  and accumulation of the remainders. The conversion of a decimal fraction to radix  $r$  representation is accomplished by successive multiplications by  $r$  and accumulation of the integer digits so obtained.

### Decimal to Binary

2 ways of converting from decimal to binary.

#### ① Sum-of-Weights Method

Determines the set of binary weight values whose sum is equal to the decimal number.

$$\text{eg. } 9 = 8 + 1 = 2^3 + 2^0.$$

By placing a 1 in the appropriate weight positions,  $2^3$  and  $2^0$  and a 0 in the other positions to get the binary number.

$$\begin{array}{cccc} 2^3 & 2^2 & 2^1 & 2^0 \\ (1 & 0 & 0 & 1)_2 \end{array}$$

$$12_{10} = 8 + 4 = 2^3 + 2^2 \rightarrow 1100_2$$

$$25_{10} = 16 + 8 + 1 = 2^4 + 2^3 + 2^0 \rightarrow 11001_2$$

$$58_{10} = 32 + 16 + 8 + 2 = 2^5 + 2^4 + 2^3 + 2^1 \rightarrow 111010_2$$

$$82_{10} = 64 + 16 + 2 = 2^6 + 2^4 + 2^1 \rightarrow 1010010_2$$

## ② Repeated Division - by -2 Method

The number is divided by 2 and then dividing each resulting quotient by 2 until there is a 0 quotient. The remainders generated by each ~~res~~ division form the binary number. The first remainder to be produced is the least significant bit (LSB) in the binary number.

$$\begin{array}{r}
 2 \overline{) 12} \\
 2 \overline{) 6} - 0 \text{ (LSB)} \\
 2 \overline{) 3} - 0 \\
 1 - 1 \uparrow \\
 \text{(MSB)}
 \end{array}$$

$(1100)_2$

$$\begin{array}{r}
 2 \overline{) 19} \\
 2 \overline{) 9} - 1 \\
 2 \overline{) 4} - 1 \\
 2 \overline{) 2} - 0 \\
 1 - 0
 \end{array}$$

$(10011)_2$

$$\begin{array}{r}
 2 \overline{) 45} \\
 2 \overline{) 22} - 1 \\
 2 \overline{) 11} - 0 \\
 2 \overline{) 5} - 1 \\
 2 \overline{) 2} - 1 \\
 1 - 0
 \end{array}$$

$(101101)_2$

## Converting Decimal Fractions to Binary

$$0.625_{10} = 0.5 + 0.125 = 2^{-1} + 2^{-3} = 0.101_2$$

Another method is repeated multiplication - by -2. Multiplying fractional number by 2 and then multiplying each resulting fractional part of the product by 2 until the fractional product is zero. The carries generated by each multiplication form the binary number.

$$0.3125 \times 2 = \overset{\text{carry}}{\downarrow} 0.625 \rightarrow 0 \text{ (MSB)}$$

$$0.625 \times 2 = 1.25 \rightarrow 1$$

$$0.25 \times 2 = 0.50 \rightarrow 0$$

$$0.50 \times 2 = 1.00 \rightarrow 1 \text{ (LSB)}$$

$$0.0101_2$$

Qn. Convert the following decimal number to binary using Sum of weights method and repeated division-by-2 method. (repeated multiplication-by-2 for fractions)

(a) 23 (b) 57 (c) 45.5 (d) 14 (e) 21 (f) 0.375

### Octal to Decimal

$$\begin{aligned}
 2374_8 &= 2 \times 8^3 + 3 \times 8^2 + 7 \times 8^1 + 4 \times 8^0 \\
 &= 2 \times 512 + 3 \times 64 + 7 \times 8 + 4 \times 1 \\
 &= 1024 + 192 + 56 + 4 = (1276)_{10}
 \end{aligned}$$

### Decimal to Octal

$$\begin{array}{r}
 8 \overline{) 359} \\
 \underline{8 \phantom{) 44} - 7} \text{ (LSD)} \\
 \phantom{8 \phantom{) 4} - 4} \\
 \phantom{8 \phantom{) 5} - 4} \\
 \phantom{8 \phantom{) 5} - 4} \text{ (MSD)}
 \end{array}$$

$(547)_8$

$$\begin{aligned}
 (0.325)_8 &= 3 \times 0.125 + 2 \times 0.015625 \\
 &\quad + 5 \times 0.001953 \\
 &= 0.375 + 0.03125 + 0.009765 \\
 &= \underline{\underline{(0.416015)_{10}}}
 \end{aligned}$$

## Octal to Binary

Each octal digit is represented by 3 bits as

|   |   |     |
|---|---|-----|
| 0 | - | 000 |
| 1 | - | 001 |
| 2 | - | 010 |
| 3 | - | 011 |
| 4 | - | 100 |
| 5 | - | 101 |
| 6 | - | 110 |
| 7 | - | 111 |

To convert an octal number to a binary number, simply replace each octal digit by the appropriate three bits.

$$\begin{array}{c} 13_8 \\ \swarrow \downarrow \\ (001011)_2 \end{array}$$

$$\begin{array}{c} 25_8 \\ / \quad \backslash \\ (010101)_2 \end{array}$$

$$47_8$$

$$170_8$$

$$75_8$$

$$5276_8$$

$$\begin{array}{c} 37.12_8 \\ / \quad | \quad \backslash \\ (011111.001010)_2 \end{array}$$

$$73_8$$

$$125_8$$

$$46_8$$

$$723_8$$

$$5624.37_8$$

## Binary to Octal

Break the binary number into groups of three bits and convert each group into the appropriate octal digit.

$$\begin{array}{c} 110 \quad 101 \\ \underline{\quad} \quad \underline{\quad} \\ (6 \quad 5)_8 \end{array}$$

$$\begin{array}{c} 101111001 \\ \underline{\quad} \quad \underline{\quad} \quad \underline{\quad} \\ 5 \quad 7 \quad 1_8 \end{array}$$

$$\begin{array}{c} 00101110010 \\ \underline{\quad} \quad \underline{\quad} \quad \underline{\quad} \quad \underline{\quad} \\ (1 \quad 3 \quad 4 \quad 6)_8 \end{array}$$

Qn. Convert the following.

Octal to decimal - (a)  $73_8$  (b)  $125_8$

Decimal to Octal - (a)  $98_{10}$  (b)  $163_{10}$

Binary to Octal - (a)  $001001101101101100_2$

(b)  $110101111_2$  (c)  $100110001.101111_2$

(d)  $10111111.0011_2$

### Hexadecimal

| <u>Decimal</u> | <u>Binary</u> | <u>Hexadecimal</u> |
|----------------|---------------|--------------------|
| 0              | 0000          | 0                  |
| 1              | 0001          | 1                  |
| 2              | 0010          | 2                  |
| 3              | 0011          | 3                  |
| 4              | 0100          | 4                  |
| 5              | 0101          | 5                  |
| 6              | 0110          | 6                  |
| 7              | 0111          | 7                  |
| 8              | 1000          | 8                  |
| 9              | 1001          | 9                  |
| 10             | 1010          | A                  |
| 11             | 1011          | B                  |
| 12             | 1100          | C                  |
| 13             | 1101          | D                  |
| 14             | 1110          | E                  |
| 15             | 1111          | F                  |

### Binary to Hexadecimal

Break the binary number into four bit groups starting at the binary point and replace each group with the equivalent hexadecimal symbol.

$$\begin{array}{cccc} \underline{1100} & \underline{1010} & \underline{0101} & \underline{0111} \\ (C & A & 5 & 7)_{16} \end{array}$$

$$\begin{array}{cccc} \underline{0011} & \underline{1111} & \underline{0001} & \underline{01101001} \\ (3 & F & 1 & 69)_{16} \end{array}$$

### Hexadecimal to Binary

Replace each hexadecimal symbol with the appropriate four bit binary code.

$$\begin{array}{c} 10AA \\ / \quad \backslash \\ (0001 \ 0000 \ 1010 \ 0100)_2 \end{array}$$

$$\begin{array}{c} (CF83)_{16} \\ / \quad \backslash \\ (1100 \ 1111 \ 1000 \ 0011)_2 \end{array}$$

Qn Convert the following Binary to hexadecimal

$$- (0011 \ 1001 \ 1000 \ . \ 1110)_2$$

Hex to Binary -  $(97A2)_{16}$

$$(D2E.8)_{16}$$

### Hexadecimal to Decimal.

First way is to convert the hexadecimal number to binary and then convert from binary to decimal.

$$\begin{array}{c} C \\ / \quad \backslash \\ (0001 \ 1100)_2 \end{array}$$

$$\begin{aligned} 00011100_2 &= 2^4 + 2^3 + 2^2 \\ &= 16 + 8 + 4 = 28_{10} \end{aligned}$$

$$\begin{array}{c} (A85)_{16} \\ / \quad \backslash \\ (1010 \ 1000 \ 0101)_2 \end{array}$$

$$\begin{aligned} 101010000101_2 &= 2^{11} + 2^9 + 2^7 + 2^2 + 2^0 \\ &= 2048 + 512 + 128 + 4 + 1 \\ &= (2693)_{10} \end{aligned}$$

Another way to convert is by multiplying each hexadecimal digit by its weight and then taking the sum of these products.

$$(E5)_{16} = E \times 16 + 5 \times 1 = 14 \times 16 + 5 \times 1 = 224 + 5 = 229_{10}$$

$$\begin{aligned} (B2F8)_{16} &= B \times 4096 + 2 \times 256 + F \times 16 + 8 \times 1 \\ &= 11 \times 4096 + 2 \times 256 + 15 \times 16 + 8 \times 1 = (45816)_{10} \end{aligned}$$

## Decimal to Hexadecimal.

$$16 \overline{) 650}$$

$$16 \overline{) 40} - 4$$

$$2 - 8$$

$$(650)_{10} = 28A_{16}$$

## Representation of Negative numbers

When an integer binary number is positive, the sign is represented by 0 and the magnitude by a positive binary number. When the number is negative, the sign is represented by 1 but the rest of the number may be represented in one of three possible ways:

- ① Signed-magnitude representation
- ② Signed-1's Complement representation
- ③ Signed-2's Complement representation.

The signed-magnitude representation of a negative number consists of the magnitude and a negative sign. In the other two representations, the negative number is represented in either the 1's or 2's Complement of its positive value. As an example, consider the signed number 14 stored in an 8-bit register. +14 is represented by a sign bit of 0 in the leftmost position followed by the binary equivalent of 14: 00001110. Each of the eight bits of the register must have a value and therefore 0's must be inserted in the most significant positions following the sign bit. Although there is only one way to represent +14, there are 3 different ways to represent -14 with eight bits.



- In Signed - magnitude representation - 1 0001110
- In Signed - 1's Complement representation - 1 1110001
- In Signed - 2's Complement representation - 1 1110010

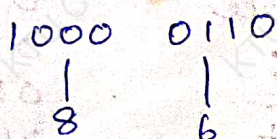
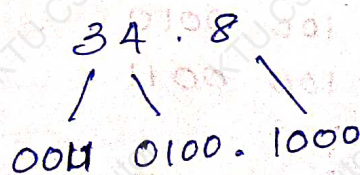
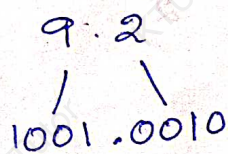
The signed-magnitude representation of -14 is obtained from +14 by complementing only the sign bit. The signed 1's complement representation of -14 is obtained by complementing all the bits of +14, including the sign bit. The signed 2's complement representation is obtained by taking the 2's complement of the positive number including its sign bit.

### Representation of BCD numbers

Binary Coded Decimal (BCD) means that each decimal digit is represented by a binary code of four bits. To express any decimal number in BCD, simply replace each decimal digit by the appropriate four bit code.

| Decimal | BCD  |
|---------|------|
| 0       | 0000 |
| 1       | 0001 |
| 2       | 0010 |
| 3       | 0011 |
| 4       | 0100 |
| 5       | 0101 |
| 6       | 0110 |
| 7       | 0111 |
| 8       | 1000 |
| 9       | 1001 |

The six four bit combinations that are not used are 1010, 1011, 1100, 1110 and 1111



Qn. Convert the following

Decimal to BCD - (a) 18 (b) 65 (c) 92.

(d) 150 (e) 321 (f) 1472

BCD to Decimal - (a) 00110001 (b) 01010011

(c) 100101110100 (d) 0001100001100000.0111

### Character Representation

Many applications of digital computers require the handling of data that consist not only of numbers, but also of the letters of the alphabet and certain special characters. Such a set contains between 32 and 64 elements (if only uppercase letters are included) or between 64 and 128 (if both uppercase and lowercase letters are included). In the first case, the binary code will require six bits and in the second case, seven bits.

### ASCII

American Standard Code for Information Interchange (ASCII) is the most widely used code to represent characters. It is a seven bit code in which the decimal the decimal digits are represented by the BCD code preceded by # 011.

| Character | ASCII code | Character | ASCII code |
|-----------|------------|-----------|------------|
| A         | 100 0001   | 0         | 011 0000   |
| B         | 100 0010   | 1         | 011 0001   |
| C         | 100 0011   | 2         | 011 0010   |
| ⋮         | ⋮          | ⋮         | ⋮          |
| P         | 101 0000   | 9         | 011 1001   |
| Q         | 101 0001   | Space     | 010 0000   |
| ⋮         | ⋮          | .         | 010 1110   |
|           |            | (         | 010 1000   |
|           |            | +         | 010 1011   |

## EBCDIC (Extended Binary Coded Decimal Interchange Code)

This is an eight bit code in which the decimal digits are represented by the BCD code preceded by 1111. Both lowercase and uppercase letters are represented in addition to numerous other symbols and commands. For example, uppercase A is 11000001, and lowercase a is 10000001.

ASCII and EBCDIC codes are commonly used in data transfer and computer interface applications.

### Tutorial - I

## Representation of floating point numbers

Two parts to represent a floating point number.

- ① Mantissa - a signed, fixed point number
- ② Exponent - designates the position of the decimal (or binary) point

The fixed point mantissa may be a fraction or an integer. For eg.  $+6132.789$  is represented as

| Fraction     | Exponent |
|--------------|----------|
| $+0.6132789$ | $+04$    |

The value of the exponent indicates that the actual position of the decimal point is four positions to the right of the indicated decimal point in the fraction. This representation is equivalent to the scientific notation  $+0.6132789 \times 10^{+4}$ .

A floating point binary number is represented in a similar manner except that it uses base 2 for the exponent. For eg,  $+1001.11$  is represented with an 8 bit fraction and 6 bit exponent as

| Fraction   | Exponent |
|------------|----------|
| $01001110$ | $000100$ |

The fraction has a 0 in the leftmost position to denote positive.

A floating point number is said to be normalized if the most significant digit of the mantissa is non-zero. For eg.  $350$  is normalized but  $00035$  is not.  $00011010$  is not normalized and is normalized by shifting three positions to the left and discarding the leading 0's to obtain  $11010000$ .

## Binary Arithmetic

Since digital systems do not process the decimal numbers, and they only process binary numbers, it is necessary to learn the binary arithmetic.

### Addition.

| A<br>(Augend) | B<br>(Addend) | S<br>(Sum) | C<br>(Carry) |
|---------------|---------------|------------|--------------|
| 0             | 0             | 0          | 0            |
| 0             | 1             | 1          | 0            |
| 1             | 0             | 1          | 0            |
| 1             | 1             | 0          | 1            |

①  $(1010)_2 + (0101)_2$

$$\begin{array}{r} 1010 \\ 0101 \\ \hline 1111 \end{array} = \text{decimal } 10$$

$$= \text{decimal } 5$$

$$\hline \text{decimal } 15$$

② Add  $(10)_{10} + (11)_{10}$  in binary

$$\begin{array}{r} 1010 + \\ 1011 \\ \hline (10101)_2 = (21)_{10} \end{array}$$

③ Add 26 and 13 in binary

$$26 = (11010)_2$$

$$13 = (01101)_2$$

$$\hline (100111)_2 = (39)_{10}$$

④ Add 1011.011 and 110.1

$$\begin{array}{r} 1011.011 \\ 110.1 \\ \hline 10001.111 \end{array}$$

$$\begin{array}{r} 11.375 \\ 6.5 \\ \hline 17.875 \end{array}$$

⑤ Add 101.11, 1101.01 and 10000.001

$$\begin{array}{r} 101.11 \\ 1101.01 \\ 10000.001 \\ \hline 100011.001 \end{array}$$

$$\begin{array}{r} 5.75 \\ 13.25 \\ 16.125 \\ \hline 35.125 \end{array}$$

### Subtraction

Binary Subtraction can be carried out in either one of two different ways.

- ① Direct Subtraction
- ② Complement Subtraction.

### Direct Subtraction

| A<br>(Minuend) | B<br>(Subtrahend) | D<br>(Difference) | B<br>(Borrow) |
|----------------|-------------------|-------------------|---------------|
| 0              | 0                 | 0                 | 0             |
| 0              | 1                 | 1                 | 1             |
| 1              | 0                 | 1                 | 0             |
| 1              | 1                 | 0                 | 0             |

① Subtract 10101 from 11011

$$\begin{array}{r} 11011 \\ - 10101 \\ \hline 00110 \end{array}$$

$$\begin{array}{r} 27 \\ - 21 \\ \hline 6 \end{array}$$



$$\begin{array}{r}
 \textcircled{3} \quad 30 - 25 \\
 11110 \quad 11001 \\
 \text{Minuend} \quad - \quad 11110 \\
 \text{1's Compl of Subtrahend} - 00110 \\
 \hline
 \textcircled{1} 00100 \quad + \\
 \quad \quad \quad \quad \quad \quad 1 \\
 \hline
 00101
 \end{array}$$

Add EAC

$$\begin{array}{r}
 \textcircled{4} \quad 25.5 = (11001.1)_2 \\
 (12.25)_{10} = (01100.011)_2 \\
 \text{1's Compl of } 12.25 = 10011.100 \\
 \begin{array}{r}
 11001.10 \quad + \\
 10011.100 \\
 \hline
 \textcircled{1} 01101.000 \quad + \\
 \quad \quad \quad \quad \quad \quad 1 \\
 \hline
 01101.001
 \end{array}
 \end{array}$$

Add EAC

$= (13.25)_{10}$

$$\begin{array}{r}
 \textcircled{5} \quad 10.625 = 1010.101 \\
 8.75 = 1000.110 \\
 \text{1's Compl of } 8.75 = 0111.001 \\
 \begin{array}{r}
 1010.101 \\
 0111.001 \\
 \hline
 \textcircled{1} 0001.110 \quad + \\
 \quad \quad \quad \quad \quad \quad 1 \\
 \hline
 0001.111
 \end{array}
 \end{array}$$

$= (1.875)_2$

Case (ii) - Subtraction of larger number from smaller number

Procedure

- ① Determine the 1's complement of the larger number (Subtrahend)
- ② Add 1's complement to the smaller number (minuend)
- ③ After addition, no carry will be generated but answer is in 1's complement form (negative number).



To get the answer in true form, take the 1's complement of it and assign negative sign to the answer.

①  $43 - 57$

$$\begin{array}{r} \text{Minuend} \quad 43 \quad - \quad 101011 \\ \text{1's Complement of} \\ \text{Subtrahend} \quad 57 \quad - \quad 000110 \\ \hline 110001 \end{array}$$

The result has no carry, so the answer is in 1's complement form.

1's Complement of 110001 is  $-(001110)$ .

Answer in true form =  $(-14)_{10}$ .

②  $8 - 10$

$$\begin{array}{r} 8 \quad - \quad 1000 \quad \text{Minuend} \quad - \quad 1000 \\ 10 \quad - \quad 1010 \quad \text{1's Complement of Subtrahend} \quad - \quad 0101 \\ \hline 1101 \end{array}$$

No carry, so the answer is in 1's complement form.

$\therefore$  True answer =  $-(0010)_2$

=  $(-2)_{10}$

③  $(8.75)_{10} = (1000.110)_2$

$(10.625)_{10} = (1010.101)_2$

$$\begin{array}{r} \text{Minuend} \quad - \quad 1000.110 \\ \text{1's Complement of} \\ \text{Subtrahend} \quad - \quad 0101.010 \\ \hline 1110.000 \end{array}$$

④  $\frac{5}{8} - \frac{7}{8}$  Minuend - 0.101

$$\frac{5}{8} = 0.101 \quad \text{1's Complement of Subtrahend} \quad - \quad 0.000 \\ \hline 0.101$$

$\frac{7}{8} = 0.111$

Answer is in 1's complement form.

$\therefore$  True answer is  $-(0.010)_2 = (-\frac{2}{8})_{10}$

$$\begin{aligned}
 & \times \textcircled{5} \quad 16.875 - 11.125 \\
 & 16.875 = 10000.111 \\
 & 11.125 = 01011.001
 \end{aligned}$$

$$\begin{array}{r}
 \text{Minuend} \quad - \quad \cancel{10000.00} \\
 \phantom{\text{Minuend}} \quad 10000.111 + \\
 \text{1's Complement of} \quad - \quad 10100.110 \\
 \text{Subtrahend} \quad \hline
 \phantom{\text{Minuend}} \quad \phantom{10000.111} \phantom{-} 00101.101
 \end{array}$$

### Advantages

- ① Since 1's Complement Subtraction can be accomplished with a binary adder, this method is useful in arithmetic logic circuits.
- ② It is very easy to find the 1's Complement of a number.

### Disadvantages

- ① Hardware implementation is difficult and it gives the concept of negative zero.

### 2's Complement Method

The 2's Complement of any binary number is determined by adding 1 to 1's Complement of that number. 2's complement form is used to represent negative numbers. 2's Complement of 1 is 1 and zero is 10.

Case (i) : Subtraction of a smaller number from larger number

### Procedure

- ① Determine the 2's Complement of the small number (Subtrahend)
- ② Add 2's Complement to the minuend.
- ③ Discard the carry generated.

$$\begin{array}{r} \textcircled{1} \quad (111.001)_2 - (101011)_2 \\ \text{Minuend} \quad - \quad 111001 \\ \text{2's complement of} \\ \text{Subtrahend} \quad - \quad 010101 \\ \hline \textcircled{1}001110 \end{array} \quad = \begin{array}{r} 57 \\ 43 \\ \hline 14 \end{array}$$

$$\begin{array}{r} \textcircled{2} \quad (100.5)_{10} - (50.75)_{10} \\ \text{Minuend} \quad - \quad 1100100.10 \\ \text{2's complement of} \\ \text{Subtrahend} \quad - \quad 1001101.01 \\ \hline \textcircled{1}0110001.11 \end{array} = (49.75)_{10}$$

$$\textcircled{3} \quad (1111)_2 - (1010)_2 \quad \textcircled{4} \quad (112)_{10} - (65)_{10}$$

$$\textcircled{5} \quad 22 - 7 \quad \textcircled{6} \quad (100 - 5)$$

Case (ii) : Subtraction of a larger number from smaller number

### Procedure

- ① Determine the 2's complement of the subtrahend.
- ② Add the 2's complement to minuend.
- ③ Answer is in 2's complement form. Take the 2's complement and assign negative sign to the answer.

No carry will be generated.

$$\begin{array}{r} \textcircled{1} \quad 7 - 22 \\ 00111 - 10110 \\ \text{Minuend} \quad - \quad 00111 \\ \text{2's compl of} \\ \text{Subtrahend} \quad - \quad 01010 \\ \hline \textcircled{1}0001 \end{array}$$

$$\text{Take 2's complement of } 1001 = 1111$$

$$= (-15)_{10}$$

②  $16.5 - 24.75$

$10000.10 - 11000.11$

Minuend =  $10000.10$

2's complement of Subtrahend =  $\frac{00111.01}{10111.11}$

2's complement of result =  $-(01000.01)_2 = -(8.25)_{10}$

Binary Multiplication

| Multiplicand<br>A | Multiplier<br>B | Product<br>P |
|-------------------|-----------------|--------------|
| 0                 | 0               | 0            |
| 0                 | 1               | 0            |
| 1                 | 0               | 0            |
| 1                 | 1               | 1            |

①  $7 \times 5$   
 $111 \times 101$

$$\begin{array}{r} 111 \times \\ 101 \\ \hline 111 \\ 000 \\ 111 \\ \hline 100011 = 35 \end{array}$$

②  $4.75 \times 3.625$

$100.110 \times 11.101$

$$\begin{array}{r} 100.110 \times \\ 11.101 \\ \hline 100110 \\ 000000 \\ 100110 \\ 100110 \\ \hline 1001001110 \end{array}$$

$= (17.21875)_{10}$

③  $22 \times 6$

④  $27 \times 21$

⑤  $\frac{3}{8} \times \frac{1}{4}$

## Binary Division

$$0 \div 1 = 0$$

$$1 \div 1 = 1$$

①  $50 \div 5$

$$110010 \div 101$$

$$= (1010)_2 = (10)_{10}$$

$$\begin{array}{r} 1010 \\ \hline 101 \overline{) 110010} \\ \underline{101} \phantom{0} \\ 0101 \phantom{0} \\ \underline{101} \phantom{0} \\ 0 \phantom{0} \end{array}$$

②

$$\begin{array}{r} 110.01 \\ \hline 100 \overline{) 11001} \\ \underline{100} \phantom{00} \\ 100 \phantom{00} \\ \underline{100} \phantom{00} \\ 0100 \phantom{0} \\ \underline{100} \phantom{0} \\ 0 \phantom{0} \end{array}$$

$$25 \div 4 = (6.25)_{10}$$

$$\begin{array}{r} 108421 \\ \hline 11001 \end{array}$$

$$\begin{array}{r} 6 \\ \hline 4 \overline{) 24} \\ \underline{24} \\ 0 \end{array}$$

③

$$\begin{array}{r} 10 \\ \hline 101.01 \overline{) 1010.10} \\ \underline{101} \phantom{00} \\ 00 \phantom{00} \end{array}$$

④

$$\begin{array}{r} 110.00 \\ \hline 10 \overline{) 1100.00} \\ \underline{10} \phantom{00} \\ 10 \phantom{00} \\ \underline{10} \phantom{00} \\ 00 \phantom{00} \\ \underline{00} \phantom{00} \\ 00 \phantom{00} \\ \underline{00} \phantom{00} \\ 00 \phantom{00} \end{array}$$

## Floating Point Arithmetic Operations

A floating point number in computer registers consists of two parts:

① a mantissa  $m$

② an exponent  $e$

The two parts represent a number obtained from multiplying  $m$  times a radix  $r$  raised to the value of  $e$ . thus

$$m \times r^e$$

The mantissa may be a fraction or an integer.  
eg. the decimal number 537.25 is represented  
in a register with  $m = 53725$  and  $e = 3$  and is  
interpreted to represent the floating point number

$$.53725 \times 10^3$$

A floating point number is normalized if the  
most significant digit of the mantissa is non-zero.  
A zero cannot be normalized because it does not  
have a non-zero digit. It is represented in floating  
point by all 0's in the mantissa and exponent.

Adding or subtracting two floating point numbers  
requires first an alignment of the radix point since  
the exponent parts must be made equal before

adding or subtracting the mantissas.

$$.5372400 \times 10^2 \quad +$$

$$.1580000 \times 10^{-1}$$

- Shift the first number three positions to the left,  
Or shift the second number three positions to the  
right.

- When the mantissas are stored in registers, shifting  
to the left causes a loss of most significant digits  
and shifting to the right causes a loss of least  
significant digits. So second one is preferable.

$$.5372400 \times 10^2 \quad +$$

$$.0001580 \times 10^2$$

---

$$.5373980 \times 10^2$$

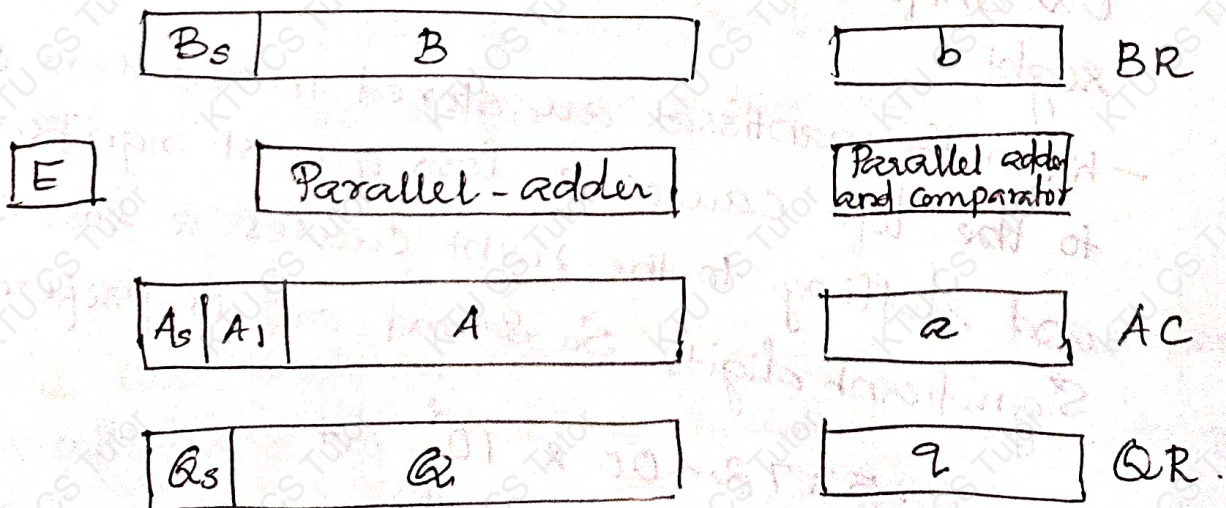
- When two normalized mantissas are added, the sum may contain an overflow digit.
- Overflow can be corrected by shifting the sum once to the right and incrementing the exponent.

$$\begin{array}{r}
 .56780 \times 10^5 \\
 .56430 \times 10^5 \\
 \hline
 .00350 \times 10^5
 \end{array}$$

To normalize, shift the mantissa to the left and decrement the exponent until a non-zero digit appears in the first position. ( $.3500 \times 10^3$ )

- For multiplication and division, alignment of mantissas are not required.
- Product can be formed by multiplying the two mantissas and adding the exponents.
- Division is accomplished by dividing the mantissas and subtracting the exponents.

### Register Configuration



- 3 registers BR, AC and CR.

each register is subdivided into two parts.

- ① Mantissa part - upper case letter symbol
- ② exponent part - lower case letter symbol.

$A_s, B_s, R_s$  represents the sign bit.

$A_1$  is the most significant bit of A. The bit in this position must be a 1 for the number to be normalized.

Parallel adder adds the two mantissas and transfers the sum into A and the carry to E. Separate parallel adder is used for the exponents.

The numbers in the registers are assumed to be initially normalized. After each arithmetic operation, the result will be normalized.

### Addition and Subtraction

- the two floating point operands are in AC and BR.
- Sum or difference formed in the AC.
- algorithm can be divided into 4 consecutive parts.

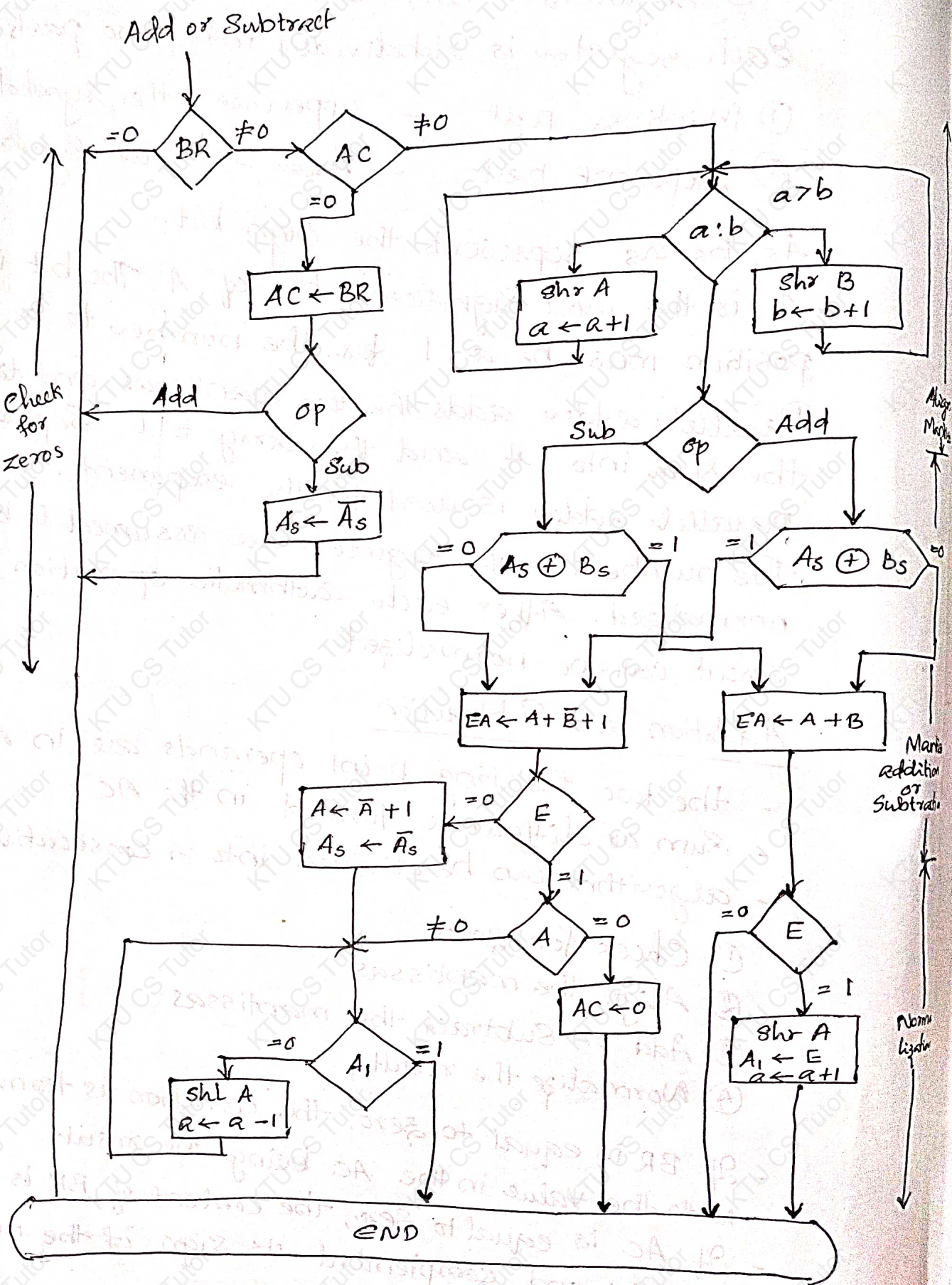
- ① Check for zeros
- ② Align the mantissas
- ③ Add or Subtract the mantissas
- ④ Normalize the result.

- If BR is equal to zero, the operation is terminated with the value in the AC being the result.

- If AC is equal to zero, the content of BR is transferred into AC and complement its sign if the numbers are to be subtracted.

If neither number is equal to zero, proceed to align the mantissas.





If the two exponents are not equal, the mantissa having smaller exponent is shifted to the right and its exponent incremented. This process is repeated until the two exponents are equal.

The magnitude part is added or subtracted depending on the operation and the signs of the two mantissas. If an overflow occurs when the magnitudes are added, it is transferred into flip flop  $E$ . If  $E$  is equal to 1, the bit is transferred into  $A$ , and all other bits of  $A$  are shifted right. The exponent must be incremented to maintain the correct number.

In the case of subtraction, the result may be zero or may have an underflow. If the mantissa is zero, the entire floating point number in the AC is made zero. Otherwise, the mantissa must have at least one bit that is equal to 1. The mantissa has an underflow if the most significant bit in position  $A_1$  is 0. In that case, the mantissa is shifted left and the exponent decremented. The bit in  $A_1$  is checked again and the process is repeated until it is equal to 1. When  $A_1 = 1$ , the mantissa is normalized and the operation is completed.



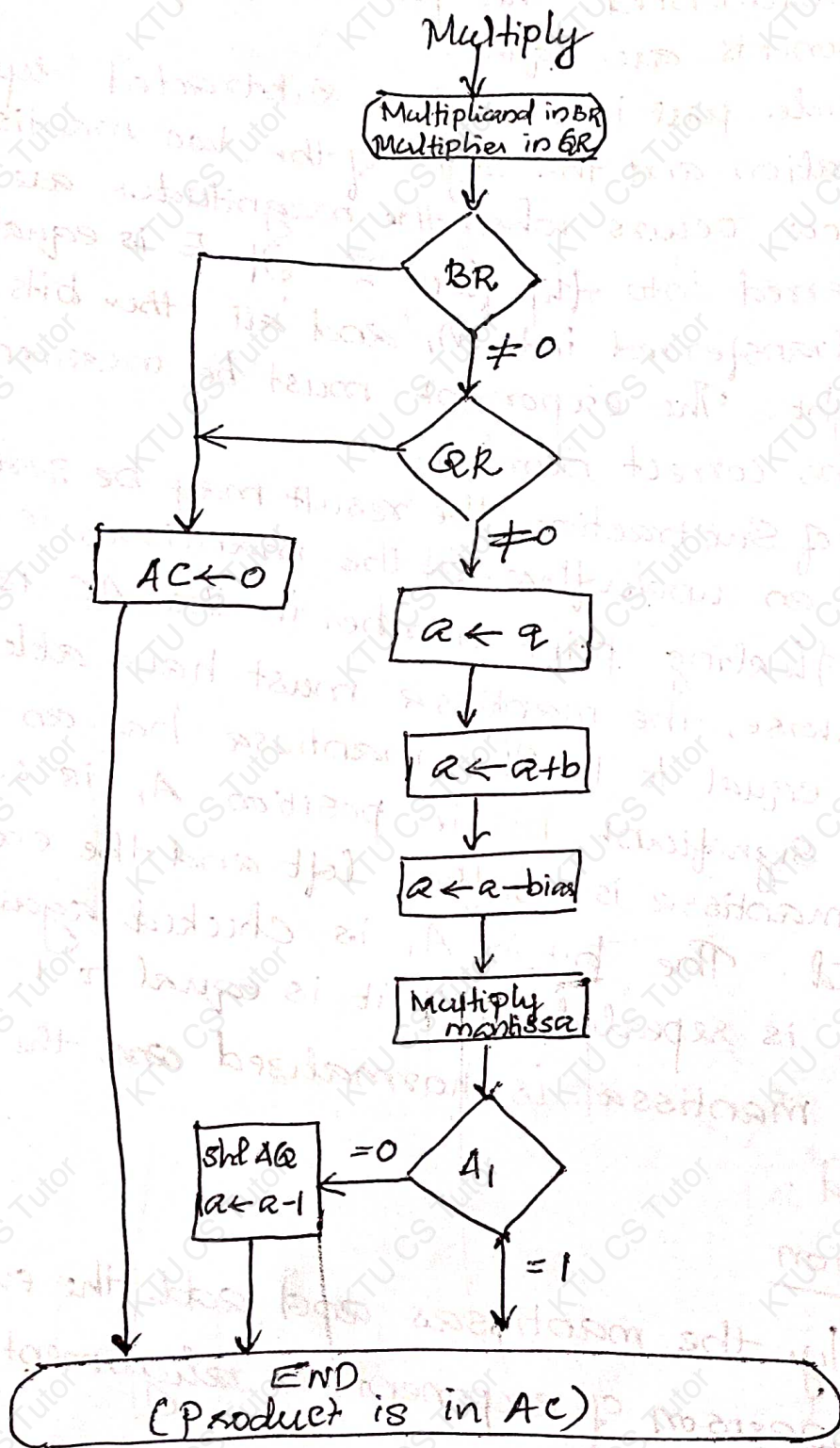
### Multiplication

- Multiply the mantissas and add the exponents.
- no comparison of exponents or alignment of mantissa is necessary.

- 4 parts:

- ① Check for zeros
- ② Add the exponents
- ③ Multiply the mantissas
- ④ Normalize the product

Steps 2 and 3 can be done simultaneously if separate addresses are available for the mantissas and exponents



The two operands are checked to determine if they contain a zero. If either operand is equal to zero, the product in the AC is set to zero and the operation is terminated.

If neither of the operands is equal to zero, the process continues with the exponent addition. The exponent of the multiplier is in  $q$  and the adder is between exponents  $a$  and  $b$ . The exponent from  $q$  should be transferred to  $a$  and add the two exponents.