

FLOATING POINT NUMBERS:-

In a fixed point arithmetic, only small numbers can be handled. For example, in a 6-bit number, the maximum numbers represented are

011111 $\rightarrow (+ve)$

111111 $\rightarrow (-ve)$

If the value of the number is to be increased, then the no of bits is to be increased. This is impossible in computer because it is expensive.

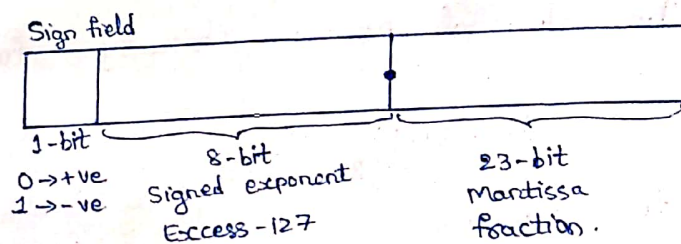
To represent very large nos like Avogadro's number (6.024×10^{23}) or very small no like Planck's constant (6.625×10^{-34}), the floating point representation becomes essential.

The general format of floating point no is

$$\pm 1.M \times 2^E$$

where M is the mantissa which must be a fraction
E is the exponent can be +ve or -ve integer.

32 bit number in floating point representation
(IEEE standard) :-



The binary point is always assumed at the left most position of mantissa.

The 8-bit exponent part should represent both negative as well as positive exponent.

The highest no that can be represented by 8-bit is 256. This is divided into 2 portions. One part is used to represent positive exponent, & the other to represent negative exponent.

The last 23 bits represents mantissa. Since, binary

normalization is used, the most significant bit of the mantissa is always '1'. This bit is not represented & it is assumed to be to the immediate left of the binary point.

Hence, the 23 bits stored in the Mantissa-field actually represents the fractional part of the mantissa, i.e., the bits which are present to the right of binary point.

The range of E' ($E' = E + 127$) is $0 \leq E' \leq 255$.

The end values of this range, 0 & 255 are used to represent special values. Therefore, the range of E' for normal values is

$$1 \leq E' \leq 254.$$

This means that the actual exponent (E) is in the range

$$-126 \leq E \leq 127.$$

The 32-bit standard representation is called Single Precision Representation because, it occupies a single 32-bit word. The Scale-Factor has a range of 2^{-126} to 2^{+127} .

- i) Represent $(+1259.125)_{10}$ in a single Precision format.
Converting the decimal no in binary format

$$\equiv +(01.0011101011 \cdot 001)_2$$

Normalizing the number

$$\equiv +1.0011101011001 \times 2^{10}$$

From the given number

$$S = 0; E = 10; E' = 137_{(10)}; M = 0011101011001$$

$$\therefore E' = 89_H = 10001001$$

0	10001001	0011101011001000....0
---	----------	-----------------------

$$\equiv 449D6400_{(H)}$$

$$\begin{array}{r} 16 \overline{) 1259} \\ 16 \overline{) 78} - B \\ 4 - E \\ 4EB_{(H)} \end{array}$$

$$\begin{array}{r} 16 \overline{) 125} \\ 16 \overline{) 78} - B \\ 4 - E \\ 4EB_{(H)} \end{array}$$

$$\begin{array}{r} 0.125 \times 2 \\ 0 \leftarrow 0.2500 \\ 0 \leftarrow 0.5000 \\ 1 \leftarrow 1.0 \end{array}$$

$$\begin{array}{r} 16 \overline{) 137} \\ 8 - 9 \end{array}$$

$$(-307.1875)_{10}$$

$$\equiv -000100110011.0011_{(2)}$$

~~5248~~

$$\equiv -1.001100110011 \times 2^8$$

$$S = 1 ; E = 8_{(10)} = 1000_{(2)} ; E' = 8 + 127 = 135_{(10)}$$

$$E' = 135_{(10)}$$

$$= 87_{(H)} = 10000111$$

$$M = 001100110011$$

1	10000111	0011001100110....0
---	----------	--------------------

$$\equiv \underline{\underline{C3999800}}_{(H)}$$

$$(0.0625)_{10}$$

$$\equiv 0.0001_{(2)}$$

$$\equiv 1.0000... \times 2^{-4}$$

$$\therefore S = 0 ; E = -4 ; E' = 123 = 7B_{(H)} = 01111011$$

$$M = 00000$$

0	01111011	0000.....0
---	----------	------------

$$\equiv 3D800000_{(H)}$$

$$(+1.725)_{10}$$

$$\equiv 1.10111001100110011001100 \times 2^0$$

$$S = 0 ; E = 0 ; E' = 127 = 7F = 01111111$$

$$M = 10111001100110011001100$$

0	01111111	10111001100110011001100
---	----------	-------------------------

$$\equiv \underline{\underline{3FBCCCCC}}_{(H)}$$

$$\begin{array}{r} 16 \overline{) 307} \\ 16 \overline{) 19} - 3 \\ \hline 1 - 3 \\ \hline 133_{(H)} \end{array}$$

$$\begin{array}{r} 0.1875 \times 2 \\ \hline 0.375 \times 2 \\ \hline 0.75 \times 2 \\ \hline 1.5 \times 2 \\ \hline 1 \leftarrow 1 \\ \hline 0.011_{(2)} \end{array}$$

$$\begin{array}{r} 16 \overline{) 135} \\ \hline 8 - 7 \end{array}$$

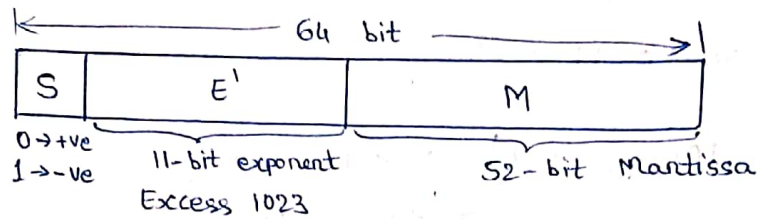
$$\begin{array}{r} 0.0625 \times 2 \\ \hline 0.125 \times 2 \\ \hline 0.25 \times 2 \\ \hline 0.5 \times 2 \\ \hline 1.0 \end{array}$$

$$\begin{array}{r} 16 \overline{) 123} \\ \hline 7 - B \end{array}$$

$$\begin{array}{r} 0.725 \times 2 \\ \hline 1.450 \times 2 \\ \hline 0.9 \times 2 \\ \hline 1.8 \times 2 \\ \hline 1.6 \times 2 \\ \hline 1.2 \times 2 \\ \hline 0.4 \times 2 \\ \hline 0.8 \times 2 \\ \hline 1.6 \times 2 \\ \hline 1.2 \times 2 \\ \hline 0.4 \times 2 \\ \hline 0.8 \end{array}$$

$$\begin{array}{r} 16 \overline{) 127} \\ \hline 7 - F \end{array}$$

* To provide more precision & range for floating point numbers, the IEEE standard also specifies a DOUBLE PRECISION Format. The Double Precision format has an increased exponent & mantissa range.



" $\rightarrow 2048$
0 $\rightarrow 2047$

The range of E' is $0 \leq E' \leq 2047$. The end values of this range, 0 & 2047 are used to represent special values. Therefore, the range of E' for normal values is $1 \leq E' \leq 2046$. This means that the actual exponent ($E = E' - 1023$) is in the range $-1022 \leq E \leq 1023$.

i) Represent $+1259.125$ in double precision format

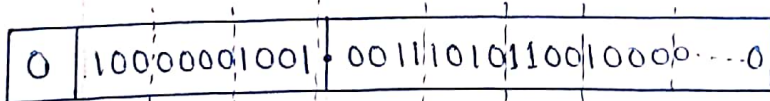
$$1259.125 \equiv +(010011101011.001)_2$$

$$\equiv +1.0011101011001 \times 2^{10}$$

$$\therefore S = 0; M = 0011101011001000 \dots 0$$

$$E = 10 \Rightarrow E' = 10 + 1023 = 1033_{(10)} \\ = 409_{(H)} = 010000001001_{(2)}$$

$$\begin{array}{r} 16 \overline{) 1033} \\ 16 \overline{) 64} - 9 \\ \hline 4 - 0 \end{array}$$



$$\equiv 4093AC8000000000_{(H)}$$

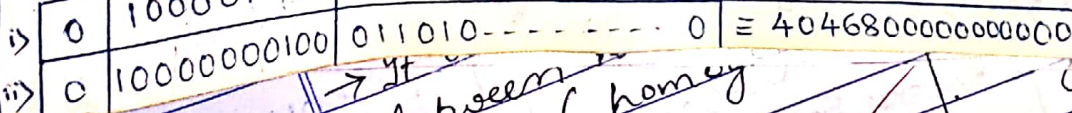
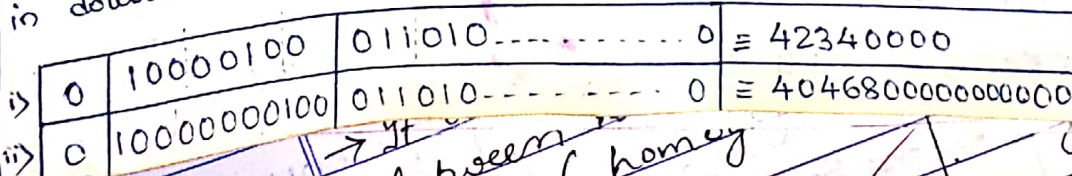
ii) $+45$ in single precision & double precision format

$$45 \equiv 00101101 \times 2^0 \equiv 1.01101 \times 2^5$$

$$S = 0; M = 01101; E = 5$$

$$\text{in single precision; } E' = E + 127 = 132 = 84_{(H)} \\ = 10000100$$

$$\text{in double precision; } E' = E + 1023 = 1028 = 404_{(H)} \\ = 10000000100$$



$$\begin{array}{r} 16 \overline{) 45} \\ 2 - 0 \\ \hline = 00101101 \\ 16 \overline{) 132} \\ 8 - 4 \\ \hline 16 \overline{) 1028} \\ 16 \overline{) 64} - 4 \\ \hline 4 - 0 \end{array}$$

The Hexadecimal value of π is 3.243F6A8885A308D3...
 work out the IEEE standard representation π in single & double precision format.

$$\pi \approx 0011.001001000011111011010100010001000010110100011000100011010011001$$

$$\approx 1.1001001000011111011010100010001000010110100011000100011010011 \times 2^1$$

a) Single precision format

$$S = 0; E = 1; M = 1001001000011111011010$$

$$E' = 1 + 127 = 128 = 80_{(CH)} = 1000\ 0000$$

0	10000000	1001001000011111011010
---	----------	------------------------

$$\approx 40490FDA$$

b) Double precision format:

$$S = 0; E = 1; M = 1001001000011111011010100010001000010110100011000$$

$$E' = 1024 = 400_{(CH)}$$

0	100000000000	1001001000011111011010100010001000010110100011000
---	--------------	---

$$\approx 400921FB54442D18$$

$$\left. \begin{array}{l} -25.125 \\ -0.08125 \end{array} \right\} \text{S.P.F \& D.P.F}$$

$$\approx 00011001.001 \approx 1.001001 \times 2^4; E' = 131 = 83_{(CH)}$$

$$E' = 1027 = 403_{(CH)}$$

1	10000011	100100100.....0	$\approx C1C90000$
1	10000000011	100100100.....0	$\approx C039200000000000$

$$\begin{array}{r} 19 \quad 0.125 \times 2 \\ 0 \leftarrow 0.25 \times 2 \\ 0 \leftarrow 0.5 \times 2 \\ 1 \end{array}$$

$$-0.08125$$

$$\approx -0.000101001100110011..... \times 2^{-4}; E' = 123 = 7B_{(CH)}$$

$$\approx -1.01001100110011..... \times 2^{-4}; E' = 1019 = 3FB_{(CH)}$$

1	01111011	010011001100110011.....	$\approx BPA66666$
1	0111111011	010011001100110011.....	$\approx BFB4CCCCCCCC$

$$\begin{array}{r} 0.08125 \times 2 \\ 0 \leftarrow 0.1625 \times 2 \\ 0 \leftarrow 0.325 \times 2 \\ 0 \leftarrow 0.65 \times 2 \\ 1 \leftarrow 1.3 \times 2 \\ 0 \leftarrow 0.6 \times 2 \\ 1 \leftarrow 1.2 \times 2 \\ 0 \leftarrow 0.4 \times 2 \\ 0 \leftarrow 0.8 \times 2 \\ 1 \leftarrow 1.6 \times 2 \\ 1 \leftarrow 1.2 \times 2 \\ 0 \leftarrow 0.4 \times 2 \end{array}$$

* SPECIAL VALUES:-

The end values 0 & 255 of the excess-127 E' are used to represent special values.

- 1) When $E' = 0$ & $M = 0$, the value exact 0 is represented.
- 2) When $E' = 0$ & $M \neq 0$, a denormal number is represented. Their value is $\pm 0.0 \times 2^{-126}$ i.e., their nos are smaller than smallest normal numbers.
- 3) When $E' = 255$ & $M = 0$, the value infinity is represented, where infinity is the result of dividing a normal number by 0.
- 4) When $E' = 255$ & $M \neq 0$, the value represented is called "NOT A NUMBER" (NaN). NaN is the result of invalid operation such as $\frac{0}{0}$, $\sqrt{-1}$.

* EXCEPTIONS:-

According to IEEE standards, the processor sets the Exception flags if Underflow, Overflow, \div by zero, Inexact or Invalid condition occurs during the program Execution.

- > UNDERFLOW: In a single precision, if the number requires an exponent less than (-126) or in a double precision, if the number requires an exponent less than (-1022) to represent its normalized form, UNDERFLOW occurs.
- > OVERFLOW: In a single precision, if the number requires an exponent greater than $(+127)$ or in a double precision, if the number requires an exponent greater than $(+1023)$ to represent its normalized form, OVERFLOW occurs.
- > DIVIDE BY ZERO: Divide by zero exception occurs when any number is divided by zero.
- > INEXACT: Inexact is the name for a result that requires rounding in order to be represented in one of the normal format.
- > INVALID: An Invalid exception occurs when an operation such as $\frac{0}{0}$ or $\sqrt{-1}$ are attempted.

When exception occurs, the result are set to special values. Systems & User-defined routines are used to handle such exceptions.

CHOPPING-OFF (OR) TRUNCATING:

If an 8-bit number is to be reduced to a 3-bit number, simply chop-off the high order bits as follows:

Suppose the binary number is (0.10101001) . This can be reduced to 3-bit by neglecting all the bits after the third, so that the approximation is (0.101) .

This process is called TRUNCATING or CHOPPING-OFF.

ROUNDING:

Instead of truncating, if the closed available value is chosen, it is called ROUNDING. In the above case, the rounded value is (0.11) . In computers, after any arithmetic operation, it is always required to perform either Rounding or Truncating to get a representable value.

ARITHMETIC OPERATION USING FLOATING POINT NUMBER:

ADDITION & SUBTRACTION:

Choose the number with smaller exponent & shift its mantissa right number of times equal to difference in exponent.

Set the exponent of the result equal to larger exponent.

Perform Addition / Subtraction on the Mantissa & determine the sign of the result.

Normalize the resulting value.

Check for over-flow / Under-flow & take appropriate action.

$$\text{Let } A = M_x 2^{E_x} \text{ and } B = M_y 2^{E_y}$$

$$\text{Then } Z = A \pm B \text{ [Make exponents equal]}$$

Add single precision floating point numbers A & B
 where $A = 4490000$ & $B = 42A00000$

I Step:

$$A \equiv \begin{array}{c|c|c} S & E & M \\ \hline 0 & 1000100 & 10010000000000000000 \end{array}$$

$$B \equiv \begin{array}{c|c|c} S & E & M \\ \hline 0 & 1000010 & 10100000000000000000 \end{array}$$

\therefore A & B are positive

$$E'_A = 89_{(10)} = 137_{10} \Rightarrow E_A = E'_A - 127 = 10$$

$$M_A = 001000 \dots$$

$$E'_B = 85_{(10)} = 133_{10} \Rightarrow E_B = E'_B - 127 = 6$$

$$M_B = 01000 \dots$$

$$\therefore A \equiv 1.001 \times 2^{10} \quad B \equiv 1.010 \times 2^6$$

Adjusting value of exponent of B is 10

\therefore B has smaller exponent with difference
 \therefore Its mantissa is shifted right by 4-bits as shown

$$\text{Shifted } M_B = 000101000 \dots$$

II Step:

Addition of Mantissa

$$\begin{array}{r} M_A \equiv 00100000 \dots \\ + M_B \equiv 00010100 \dots \\ \hline \text{Result} \equiv 00110100 \dots \end{array}$$

$$\begin{aligned} \therefore A + B &\equiv 0 \mid 10001001 \mid 001101000000000000000000 \\ &\equiv 449A0000 \\ &= 1.001101 \times 2^{10} \end{aligned}$$

Subtraction of Mantissa

$$\begin{array}{r} M_A = 00100000 \dots \\ - M_B = 00010100 \dots \\ \hline \text{Result} = 00001100 \end{array}$$

$$A - B = 0 \mid 10001001 \mid 000011000 \dots = 44860000$$

* FLOATING POINT MULTIPLICATION :

Let $X = M_x \times 2^{E_x}$ & $Y = M_y \times 2^{E_y}$ be the 2 floating point numbers. Then $Z = X * Y$

$$= (M_x * M_y) \times 2^{(E_x + E_y)}$$

i) Add the exponent & subtract the bias (127 in case single precision numbers & 1023 in case of double precision numbers)

ii) Multiply the mantissa & determine the sign of the result.

iii) Normalize the result.

iv) Check for overflow/underflow & take the appropriate action.

Ex: Perform the multiplication of

$$A = 36100000$$

$$B = D4100000$$

$$A \equiv 010110110000100$$

$$B \equiv 110101000001000$$

$$A; S=0; E' = 6C = 108_{(10)}; E = E' - 127 = 108 - 127 = -19$$

$$\therefore A \equiv +1.001 \times 2^{-19} \equiv +1001 \times 2^{-22}$$

$$B; S=1; E' = A8 = 168_{(10)}; E = E' - 127 = 168 - 127 = 41$$

$$\therefore B \equiv -1.001 \times 2^{41} \equiv -1001 \times 2^{38}$$

$$\begin{aligned} * \frac{1001}{1001} &\equiv * \frac{1001}{0111} \equiv * \frac{1001}{2^{-1}} \\ &\quad \begin{array}{r} 0000111 \\ 10010 \\ \hline \end{array} \\ \text{Result} &\equiv 1001111 \end{aligned}$$

$$\begin{array}{r} 1001 \\ 01110 \\ +1001 \\ \hline 2^{-1} \end{array}$$

$$\begin{aligned} \therefore \text{Result} &\equiv 1001111 \times 2^{38-22} \\ &\equiv 1001111 \times 2^{16} \equiv 1.001111 \times 2^{22} \end{aligned}$$

$$22 \times 127$$

$$\begin{aligned}\text{Result} &= (-1.001 \times 1.001) \times 2^{41-19} \\ &= -1.010001 \times 2^{22}\end{aligned}$$

$$\begin{aligned}E &= 22 \\ E' &= 149 \\ 16 \overline{) 149} \\ &\quad 9-5\end{aligned}$$

1 | 10010101 | 010001000000.....

≡ CAA20000

* FLOATING POINT DIVISION:-

Let $X = M_x \times 2^{E_x}$ & $Y = M_y \times 2^{E_y}$

$$Z = \frac{X}{Y} = \frac{M_x}{M_y} \cdot 2^{E_x - E_y}$$

- i) Subtract the exponent. Then the bias being vanished, add the bias (127 or 1023). Without the addition of bias, it will be zero-exponent.
- ii) Divide the mantissa & determine the sign of the result.
- iii) Normalize the resulting value if necessary.
- iv) Check for overflow / underflow & take suitable action.

Ex Divide $A \equiv 36100000$ by $B \equiv D4100000$

$$A \equiv 1.001 \times 2^{-19}$$

$$B \equiv -1.001 \times 2^{41}$$

$$\text{Result} = \frac{1.001}{-1.001} \times 2^{-19-41} = -1.0000 \times 2^{-60}$$

$$E = -60 ; E' = 127 - 60 = 67$$

$$\begin{array}{r} 16 \overline{) 67} \\ \underline{4-3} \end{array}$$

$$\begin{aligned}\therefore \text{Result} &\equiv 1 | 01000011 | 00000000..... \\ &= \underline{\underline{A1800000}}\end{aligned}$$

* MEMORY LOCATION & ADDRESSES:
Refex Text