

in $\bar{y}z + \bar{y}\bar{z}$

* By using carry bits let Cin be the carry into the MSB & can't be carry out from MSB

$$\text{Cin} @ \text{Cout} = \begin{cases} 0 & \text{No overflow} \\ 1 & \text{Overflow} \end{cases}$$

When overflow occurs, no. of bits is increased by 1 (i.e. increasing the range).

a) Floating & fixed point representation

$$\rightarrow (54.3)_{10} \quad (71.25)_8 \quad (1110.001)_2$$

↓ ↓ ↓
Decimal point Octal point Binary point-

$$(x_3 x_2 x_1 \cdot x_{-1} x_{-2})_8$$

↓
Radix Point-

→ Fixed point representation : To represent any number, bit selection is pre-decided (number is represented with a fixed no. of digits before & after the decimal point).

Ex: 8 bit no.

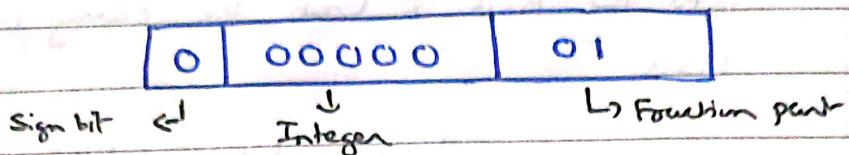
1 bit	5 bits	2 bits
-------	--------	--------

↓ ↓ ↓
Sign bit Integer Part Fractional Part-

00000 -> 0	00 -> 0
01 -> 0.25	
10 -> 0.5	
11 -> 0.75	

Eg Represent $(0.3125)_{10}$ in the above fixed point representation.

M $(0.3125)_{10} = (0.0101)_2$



→ We could have easily represented 0.0101 using 8 bits but we have to take the approximate value of $(0.3125)_{10}$

{Disadvantage of FPR}

→ If we assign more bits for the integer part, then possible range increases but if we assign more bits for the fraction part, then precision increases.

→ Fixing the no. of bits for integer part & fraction part is the disadvantage of FPR (fixed point representation), therefore we move into floating point representation

→ Floating point representation

Better Precision

$$(5.625)_{10} \rightarrow (101.101)_2$$

↓ Earlier, this was stored as follows :-

[101101, 3]

↓ Radix point is present after the 3rd bit from MSB

$$(5.625)_{10} \rightarrow 0.\underbrace{5625}_{\text{Mantissa}} \times 10^{\underbrace{1}_{\text{Exponent}}}$$

$$(101.101)_2 \rightarrow 0.\underbrace{101101}_{\text{Mantissa}} \times 2^{\underbrace{3}_{\text{Exponent}}}$$

S	Exponent-	Mantissa
---	-----------	----------

* Need of normalization

$$(101.101)_2 \rightarrow \begin{cases} 0.101101 * 2^3 \\ 1.01101 * 2^2 \\ 0.0101101 * 2^4 \\ 101101 * 2^{-3} \end{cases}$$

Different Representations

* Normalization :- Basically are two types :-

a) Explicit normalization :- Move the radix point to the LHS of most significant '1' in the bit sequence

$$(101.101)_2 \rightarrow 0.101101 * 2^3$$

Better
Precision

b) Implicit normalization :- Move the radix point to the RHS of most significant '1' in the bit sequence.

$$(101.101)_2 \rightarrow 1.01101 * 2^2$$

S	Exponent-	Mantissa
---	-----------	----------

↓
Exponent can be +ve or -ve

Exponent-

-8 -7 -6 -5 -4 -3 -2 -1 0 1 2 3 4 5 6 7

$\downarrow +8$ (Biasing)

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

Exponent is called excess 8

$$(0.0101)_2 \rightarrow 0.101 * 2^{-1}$$

S | Exponent | Mantissa

Explicit Normalization

0	0	1	1	1	0	1	0	0

$$(1.01101)_2 \rightarrow 1.01101 * 2^{\text{bias}}$$

S | Exponent | Mantissa

Implicit Normalization

0	1	0	1	0	0	1	1	0	1

Formula

$$\text{i) Explicit Normalization} : (-1)^S * 0.M * 2^{E-\text{Bias}}$$

$$\text{ii) Implicit Normalization} : (-1)^S * 1.M * 2^{E-\text{Bias}}$$

Example 8

$$(5.625)_{10} \rightarrow (101.101)_2$$

Let us represent this using both explicit & implicit normalization

Explicit $\rightarrow (101.101)_2 \rightarrow 0.101101 * 2^3$

Normalization | Exponent | Mantissa

0	1	0	1	1	1	0	1	1	0

$(-1)^s * 0.M * 2^{E-Bias}$

$$(101.101)_2 \rightarrow (5.5)_{10}$$

More Precision

$$(101.101)_2 \rightarrow 1.01101 * 2^2$$

Implicit \rightarrow | Exponent | Mantissa

0	1	0	+ 0	0	1	1	0	1

$(-1)^s * 1.M * 2^{E-Bias}$

$$(101.101)_2 \rightarrow (5.625)_{10}$$

Although, we have come across explicit & implicit normalization, however the required bits to store a particular floating point number is yet to be fixed. Therefore IEEE has come across with a standard

→ IEEE 754 standard (1985)

↳ Responsible for standardization of floating point numbers.

↳ IEEE - 1985 (First)

IEEE - 2008

IEEE - 2019

(Revised)

Name	Common Name	Precision	Exponent Bits		Significand Bits		Exponent Bias		Organization of Bits	
			S	E (bias)	M (sign)	E (bias)	M (sign)	S	E (bias)	M (sign)
binary 16	Half Precision	11	1	5	24	8	127	15	1023	16383
binary 32	Single Precision	24	1	8	53	—	1023	15	1023	16383
binary 64	Double Precision	53	1	11	113	—	1023	19	1023	262143
binary 128	Quadruple Precision	113	1	19	237	—	1023	19	1023	262143
binary 256	Octuple Precision	237	1	19	—	—	1023	19	1023	262143

IEEE 754

→ IEEE 754 - Single Precision

s(1 bit) E(8 bits) M(23 bits) Representation

0001 0000 0000 0000.....0 ±
 (E=0) (M=0)

$$0\text{ or }1 \quad \begin{array}{|c|c|} \hline & \text{1111} \\ \hline \end{array} \quad \begin{array}{|c|c|} \hline & \text{0000} \dots \dots \text{0} \\ \hline \end{array} \quad \pm\infty$$

$(E=255)$ $(M=0)$

0 or 1 $1 \leq E \leq 254$ $M = x \times x \times \dots \times$ Implicit Normalized Form

Case I	$E = 0$	$M \neq 0$	Fochical Form
	$E = 255^\circ$	$M \neq 0$	NAN

→ IEEE 754 – Double Precision

S(8 bits) E(11 bits) M(52 bits) Representation

$$\text{Orbits} \quad 000\ 0000\ 0000 \quad 000\dots 0 \quad \pm 0$$

$(E=0)$ $(H=0)$

$$\text{Octal} \quad \begin{array}{cccccc} | & | & | & | & | & | \\ \text{III} & \text{III} & \text{I} & \text{II} & \text{I} & \text{II} \end{array} \quad \begin{array}{c} \text{000...0} \\ (\text{H=0}) \end{array} \quad \pm \infty$$

(E=2047)

0 or 1 $1 \leq E \leq 2046$ $M = x_0x_1\dots x_n$ Implicit Normalized Form

$E=0$ $M \neq 0$ Fractional Form
 $E=2047$ $M \neq 0$ NRM

Ex 1 Consider a 32 bit register which stores the Floating Point numbers in IEEE single precision format,

- i) Find the decimal value of the following 32 bits

0	1000 0011	1100 ... 0
---	-----------	------------

- ii) Find the decimal value of the following 32 bits

0	1111 1111	1100 ... 01
---	-----------	-------------

Sol i) Biased exponent : 131 (1000 0011)

Bias (Weight) : 127

$$\text{Value} = (-1)^S \times 1.M \times 2^{E-\text{Bias}}$$

$$=(-1)^0 \times 1.11 \times 2^{131-127} = (11100)_2 = (28)_{10}$$

ii) Biased exponent : 1255 \$ M \neq 0\$

Represent NAN

Ex 2 Consider a 64 bit register which stores the floating point numbers in IEEE double precision format,

- i) Find the decimal value of the following 64 bits :-

1	100 0000 0011	1101100 ... 0
---	---------------	---------------

- ii) Find the decimal value of the following 64 bits :-

0	111 1111 1111	0000 ... 0
---	---------------	------------

Sol i) Biased Exponent $\therefore 1027$ (10000000011)

$$\text{Bias (Weight)} = 1023$$

$$\text{True exponent} = 1027 - 1023 = 4$$

$$\text{Value} : (-1)^s \times 1.M \times 2^{E-\text{Bias}}$$

$$(-1)^1 \times 1.11011 \times 2^4 = (-11101.1)_2$$

\uparrow
 $(-29.5)_{10}$

ii) Biased exponent $\therefore 2047$ & $M=0$ ($s=0$) $\rightarrow \infty$

Ex 3 Represent $(116)_{10}$ in IEEE 754 single precision format.

$$\text{Sol} \quad (116)_{10} \rightarrow (1110100)_2$$

$$\text{Implicit Normalized Form} : 1.110100 * 2^6$$

\rightarrow Single precision \therefore Biased exponent $= 6 + 127 = 133$

0	1000 0101	110100 ... 0
---	-----------	--------------

\rightarrow Double precision \therefore Biased exponent $= 6 + 123 = 1029$

0	100 0000 0101	110100 ... 0
---	---------------	--------------

Ex 4 Determine the range of (positive) floating point numbers represented using fractional form of IEEE 754 single precision format.

Sol Fractional form $\therefore (-1)^s \times 0.M \times 2^{-126}$

$$\text{Smallest fraction} : (-1)^0 \times 0.000...01 \times 2^{-126} = \frac{1}{2^{126}} = \frac{1}{2^{23}} = 2^{-149}$$

Largest fraction: $(-1)^0 * 0.111\dots11 * 2^{-126}$

$$\frac{1}{2^1} \frac{1}{2^2} \frac{1}{2^3} \dots \frac{1}{2^{23}} \frac{1}{2^{-22}}$$

$$= \left(\frac{1}{2} + \frac{1}{2^2} + \dots + \frac{1}{2^{23}} \right) * 2^{-126}$$

$$= \left[\frac{\frac{1}{2} \{ (\frac{1}{2})^{23} - 1 \}}{\frac{1}{2} - 1} \right] * 2^{-126} = (1 - 2^{-23}) * 2^{-126}$$

$$\text{Range: } 2^{-149} \longleftrightarrow (1 - 2^{-23}) * 2^{-126}$$

Ex5 Add the following numbers in 8-bit register using signed 2's complement notation:

- a) 50 \$ -5 b) 45 \$ -65 c) 75 \$ 85

$Cout = Cin \Rightarrow \text{No overflow}$

Carry out

Carry in

Sign Bit

Sign Bit

Overflow

Underflow

Sol a) 2's complement

$$\begin{array}{r} 50 \\ -5 \\ \hline \end{array}$$

$$\begin{array}{r} 00101011 \\ 10111101 \\ \hline 10010101 \end{array}$$

$$-5 - 011111011$$

Answer: 45

Carry in = Carry out \Rightarrow No underflow/overflow

Answer is positive (Neglect +)

b) $45 \div -65$

$$\begin{array}{r} 0\ 0\ 1\ 0\ 1\ 1\ 0\ 1 \\ -65 \\ \hline 1\ 1\ 1\ 0\ 1\ 1\ 0\ 0 \end{array}$$

\times $\underbrace{\quad}_{\substack{0\ 0\ 0\ 1\ 0\ 0\ 1\ 1}} + \underbrace{\quad}_{\substack{1\ 1\ 1\ 0\ 1\ 1\ 0\ 0}}$

Carry in = Carry out = 0 \therefore No overflow / underflow

No Carr \therefore Take 2's complement of
↓
final answer

Answer is negative (-2^0)

c) $75 \div 85$

$$\begin{array}{r} 0\ 1\ 0\ 0\ 1\ 0\ 1\ 1 \\ 85 \\ \hline 1\ 0\ 1\ 0\ 0\ 0\ 0\ 0 \end{array}$$

Carry out \neq Carry in \therefore Overflow

Normalization of floating point numbers \therefore A floating point number is said to be normalized if the most significant position of the mantissa is a non-zero.

$$(37.25)_{10} = (100101.01)_2 = 1.0010101 * 2^5$$

$$(7.625)_{10} = (111.101)_2 = 1.11101 * 2^2$$

$$(0.3125)_{10} = (0.0101)_2 = 1.01 * 2^{-2}$$