

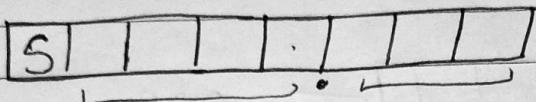
Floating point Representation

Lecture 12

Real No.

- Fixed point Notation
- Floating point Notation

* Fixed point Notation



$$+ (4.5)_{10} =$$

$$\begin{array}{r} \underline{00100} \cdot \underline{100} \\ 2 \quad 4H \end{array} \rightarrow \text{in fixed point notation}$$

in hex

$$\bullet (6.25)_{10} = \begin{array}{r} \underline{00110} \cdot \underline{010} \\ 3 \quad 2H \end{array} = \text{in hex}$$

$$\bullet (11.75)_{10} = \begin{array}{r} \underline{11011} \cdot \underline{110} \\ D \quad EH \end{array}$$

$$\bullet \text{Range } (6.625)_{10} = \begin{array}{r} \underline{00110} \cdot \underline{101} \\ 3 \quad 5H \end{array} = \text{in hex}$$

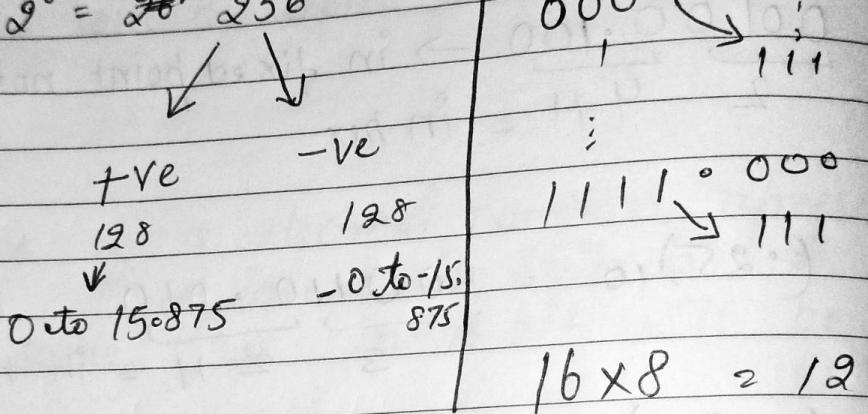
| + Ve | - Ve |
|-------------------------|-------------------|
| $0\ 1111 \cdot 111$ | $11111 \cdot 111$ |
| $+ (15 \cdot 875)_{10}$ | $- 15 \cdot 875$ |

We representing all real no. from -15.875 to 15.875] X

| | | |
|---------------|---------------------|------|
| <u>Stepup</u> | $\rightarrow 0.125$ | .0 |
| .000 | — | .125 |
| .001 | — | .250 |
| .010 | — | .325 |
| .011 | — | .500 |
| .100 | — | ! |
| | ! | ! |

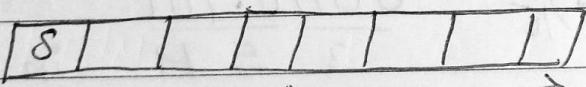


$$2^3 = \cancel{2^4} 256$$



$$16 \times 8 = 128$$

Disadvantage \rightarrow Range will be very less



if we shift point this side range will decrease
and accuracy will increase

Floating Point RepresentationPart 2

* Floating point notation

$$363.4829 \times 10^0$$

↓

$$363.4829$$

↓

$$36.34829 \times 10^1$$

$$3634.829 \times 10^{-1}$$

$$3.634829 \times 10^2$$

$$36348.29 \times 10^{-2}$$

$$\underline{3.634829 \times 10^3}$$

Normalized notation

Scientific notation

$$a * b^e$$

↙ ↓ → exponent
real no. base must be integer

$$m * b^e$$

↑ ↓ → base
 $1 \leq |m| < 1$

$$3.6 \times 10^{-5} \rightarrow \text{Scientific}$$

$$0.36 \times 10^{-4} \rightarrow \text{Normalized, Scientific}$$

$$4.98 \times 10^{-6} \rightarrow \text{Scientific}$$

$$3.04 \times 10^{4.7} \rightarrow \text{Neither normalized nor scientific}$$

• Base 2

$$101.11 \times 2^3 \rightarrow \text{Scientific}$$

$$1.01 \times 2^4 \rightarrow S$$

$$0.1 \times 2^3 \rightarrow N, S$$

$$0.01 \times 2^4 \rightarrow S$$

$$0.001 \times 2^{4.2} \rightarrow X$$

* Normalized

Binary

• 1 → it must be '1'

in other base

$\downarrow \quad \downarrow$ after point it not = to '0'
 $0 \neq 0$

Ex:

• $123 * 10^5$ = Normalize

• $93 * 10^4$ = N — "

• $1 * 2^3$ = N — "

• $101 * 2^5$ = N — "

• $01 * 2^3$ = X ~~Normalize~~

* Rule for Binary

$a * b^e$
 $\downarrow \quad \rightarrow_2$

\downarrow \downarrow
 Should not be 0 after point

Should be 0 before point

for mantissa

$\boxed{S | E | M}$

• $| \leq |m| < 1$

while storing we'll not store the base
 because : is obvious i.e. 2

format

$\boxed{S | E | M}$

$\boxed{S | M | E}$

$\boxed{M | S | E}$

$\boxed{M | E | S}$

$\boxed{E | M | S}$

$\boxed{E | S | M}$

S E M \Rightarrow 16 bit
 ↓ ↓ ↓
 1 8 7 = bits

\searrow
excess 128

Let number $\rightarrow (3.5)_{10} \rightarrow (11.1)_2$

Convert this into normalized form
 $(11.1)_2 \Rightarrow 1.11 \times 2^2$

Let's store it according to 8em

Sign(s) 0 10000010 1110000 M
 ↓ ↓ ↓
 128+2 2 is exponent
 excess 128 $\cdot 111 \times 2^2$

$$+(\cdot 25)_{10} = (\cdot 01)_2$$

$$\cdot 1 \times 2^{-1}$$

Q 01111110000000

3FC0H Ans

Hex Representation

0100000101110000

4170H Ans

Normalized

$$-(6.75)_{10} \rightarrow (110.11)_2 \Rightarrow \cdot 11011 \times 2^3$$

E M
 110000111101100
 C I E CH Ans

- $-8.6 \rightarrow 1000.100 \rightarrow 1000100\ldots$ Mantissa

| | | |
|---|---------|---------|
| S | e | M |
| 1 | 1000100 | 1000100 |

C244H Ans

- 3FC0H

0011111110000000

$$+ \cdot 1 * 2^{127-128} \Rightarrow \cdot 1 * 2^{-1} = (01)_2$$

(025)₁₀ Ans

- 8CC0H

1000110010000000
S E M

$$-ve \cdot 1 * 2^{25-128} \Rightarrow -\cdot 1 * 2^{-103} = -1 * 2^{-104}$$

-2^{104} Ans

*

| | | |
|---|---|---|
| S | M | E |
| 1 | 8 | 7 |

 \rightarrow excess 64

- $(1.5)_{10} \rightarrow (1.1)_2 \rightarrow \cdot 11 * 2^1$

0110000601000001

~~6041H~~ Ans

* E | S | M → 22 bit = 32 bit
 ↓ ↓
 9 excess ~~12~~ 56

• ~~Ans~~ $(13.875)_{10} \rightarrow (1101.111)_2 \Rightarrow .1101111 \times 2^4$

| | |
|---------------------------------|----------------------------------|
| <u>1000001000110111000...</u> | $\frac{32}{4} = 8 \text{ digit}$ |
| <u>82378000H</u> <u>Ans</u> | |
| \downarrow 5 + 3 = 8 digit | |