

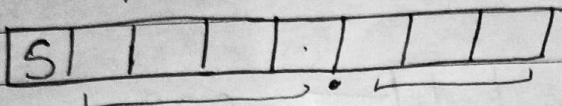
Floating point Representation

Lecture 12

Real No.

- Fixed point Notation
- Floating point Notation

* Fixed point Notation



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

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

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

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

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

$+ \text{Ve}$ $\bullet \text{Range} = 0 \ 1111 \cdot 111$ $+ (15 \cdot 875)_{10}$	$- \text{Ve}$ $1111 \cdot 111$ $- 15 \cdot 875$
-----------------------------------------------------------------------------------------	-------------------------------------------------------

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

<u>Stepup</u>	$\rightarrow 0.125$.
: 000	—	.0
: 001	—	.125
: 010	—	.250
: 011	—	.375
: 100	—	.500
,	,	,
,	,	,



$$2^8 = 256$$

\downarrow \downarrow

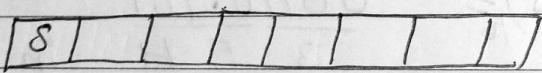
+ve	-ve
128	128
\downarrow	
0 to 15.875	-0 to -15.875

000 0. 000
1 1 1

1111 0 000
1 1 1

$$16 \times 8 = 128$$

Disadvantage \rightarrow Range will be very less



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

Floating Point Representation

Part 2

* Floating point notation

$$363.4829 \times 10^0$$

↓

$$36.34829 \times 10^1$$

$$3.634829 \times 10^2$$

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

$$363.4829$$

↓

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

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

Normalized notation

Scientific notation

$$a * b^e$$

\downarrow \downarrow

real no. base

\rightarrow exponent
must be integer

$$\begin{array}{c} + \leq \\ m * b^e \\ \downarrow \\ 0 < m < 1 \end{array} \rightarrow \begin{array}{l} \text{Integer} \\ \text{base} \end{array}$$

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

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

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

$$83.4 \times 10^{4.7} \rightarrow \text{Neither normalized nor scientific}$$

* Base 2

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

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

$$.1 \times 2^3 \rightarrow N, S$$

$$.01 \times 2^4 \rightarrow S$$

$$.001 \times 2^{4.2} \rightarrow X$$

* Normalized

Binary

• 1 → it must be '1'

in other base

$\downarrow \cdot \downarrow$
0 ≠ 0

Ex:

- $123 * 10^5 = \text{Normalize}$
- $93 * 10^4 = N \dots$
- $1 * 2^3 = N \dots$
- $101 * 2^5 = N \dots$
- $01 * 2^3 = X \quad \text{cancel}$

* Rule for Binary

at base

$\downarrow \rightarrow_2$

\downarrow \downarrow
Should not be 0 after point

Should be 0 before point

for mantissa

$\boxed{S \mid E \mid M}$

$$0 \leq |m| < 1$$

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

format

$\boxed{S \mid E \mid M}$
 $\boxed{S \mid M \mid E}$
 $\boxed{M \mid S \mid E}$

$\boxed{M \mid E \mid S}$
 $\boxed{E \mid M \mid S}$
 $\boxed{E \mid S \mid M}$

$\boxed{S \mid E \mid M} \Rightarrow 16 \text{ bit}$

$\downarrow \quad \downarrow \quad \downarrow$

1 8 7 = bits

\searrow

excess 128

Let number $\rightarrow (3 \cdot 5)_{10} \rightarrow (11 \cdot 1)_2$

Convert this into normalized form
 $(11 \cdot 1)_2 \Rightarrow \bullet 111 \times 2^2$

Let's store it according to sem

$\begin{array}{c} \boxed{0} \underline{10000010} \quad \boxed{0} \underline{1110000} \\ \downarrow \quad \downarrow \\ \text{Sign(s)} \quad 128+2 \\ \downarrow \quad \downarrow \\ \text{excess } 128 \quad 2 \text{ is exponent} \\ \bullet 111 \times 2^2 \end{array}$

$\# (-25)_{10} = (\bullet 01 \quad \bullet 1 \times 2^{-1})$
 $\bullet 0111111000000$
 $3FC0H \underline{\underline{AH}}$

Hex Representation

0100000101110000

4170H AH

Normalized

$-(6 \cdot 75)_{10} \rightarrow (110 \cdot 11)_2 = \bullet 11011 \times 2^3$

$\begin{array}{c} \boxed{1} \quad E \quad M \\ \boxed{1} \underline{100000011} \underline{1101100} \\ C \quad | \quad ECH \underline{\underline{AH}} \end{array}$

o $\underline{-8.6} \rightarrow 1000.\overline{100} \rightarrow \overline{10001001...}$
 Mental
 S E M
 $\underline{1}\underline{0000100}\underline{1000100}$

C244H Ans

o 3FCOH

0011111100 0000

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

\Downarrow

$(025)_{10}$ Ans

o 8CCOH

100011001000 0000
 S E M

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

\Downarrow

-2^{104} Ans

* [S | M | E]
 1 8 7 \rightarrow excess 64

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

0110000001000001
~~60444~~ 60444 A.

* $E | S | M \rightarrow 22 \text{ bit} = 32 \text{ bit}$
 ↓ ↓
 excess 127

• $(13.875)_{10} \rightarrow (1101.111)_2 \Rightarrow .1101111_2$

$\underline{1000001000}\underline{110111000\dots}$	$\frac{32}{4} = 8 \text{ digit}$
$82378000H$ $\downarrow +3 = 3 \text{ digit}$	<u>Ans</u>

* Floating point Representation

Part 3

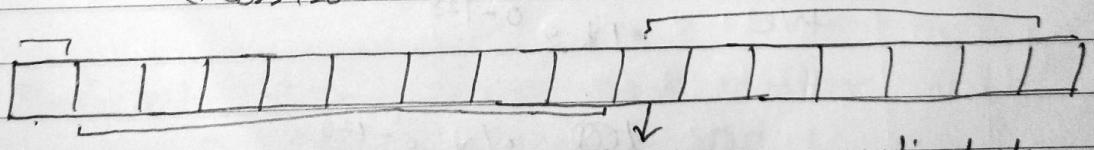
* Range

- ↳ Normalized $\rightarrow .1$
- ↳ De normalize $\rightarrow .0$

S	E	M
---	---	---

1	8	7
---	---	---

↳ excess 128



If then Denormalize else
normalize

or

*

Q. T.O

* Normalized

- largest +ve number

$$0 \underbrace{1111111}_{S} \cdot \underbrace{1111111}_{M}$$

$$\cdot 1111111 * 2^{255-128}$$

$$\cdot 1111111 * 2^{127} \Rightarrow (1-2^{-7}) * 2^{127}$$

$$1111111 * 120^{\underline{120}}$$

$$127 * 2^{\underline{120}} \text{ Ans}$$

- Smallest -ve number

$$\frac{1}{3} \underbrace{11111111}_{S} \cdot \underbrace{1111111}_{M}$$

$$- \cdot 1111111 * 2^{255-128}$$

$$- \cdot 1111111 * 2^{127} \Rightarrow (1-2^{-7}) * 2^{127}$$

$$- 127 * 120^{\underline{120}}$$

- Smallest +ve normalized number :-

$$0 \underbrace{00000000}_{\downarrow \text{+ve}} \cdot 1000000$$

$$01 * 2^{0-128}$$

~~$$01 * 2^{-128}$$~~

~~$$1 * 2^{-129}$$~~

$$2^{-129} \text{ Ans}$$

- largest -ve & normalized number

1 0000 0000 1000000

$$\begin{array}{r} \downarrow \\ - 0.1 * 2 \end{array} \quad 0-128$$

$$\begin{array}{r} - 0.1 * 9 \\ - 2 \end{array} \quad \begin{array}{l} -128 \\ -129 \end{array}$$

Aux

* Denormalized number

- largest +ve denormalized number

0 | / / / / / / / 0 | / / / / /

$$\begin{array}{l}
 \downarrow \\
 + \quad 011111 * 2^{125} \Rightarrow 111111 * 2^{126} \\
 + \quad 0111111 * 2^{127} \quad ((1 - 2^{-6}) * 2^{-128}) \\
 \hline
 \text{Ans}
 \end{array}$$

- Smallest -ve denormalized number

~~4001~~ - 63 * 2120 Ave

- Smallest +ve normalized number

$$\begin{array}{r} \underline{\textcircled{O}} \quad \underline{\textcircled{O} \textcircled{O} \textcircled{O} \textcircled{O} \textcircled{O} \textcircled{O}} \quad \underline{\textcircled{O} \textcircled{O} \textcircled{O} \textcircled{O} \textcircled{O}} \\ \textcircled{S} \qquad \textcircled{e} \qquad \textcircled{M} \\ \textcircled{+} \cdot \textcircled{O} \textcircled{O} \textcircled{O} \textcircled{O} \textcircled{O} \textcircled{O} \textcircled{O} \textcircled{1} * \textcircled{2}^{\textcircled{0}-\textcircled{128}} \end{array}$$

$$2^{-7} * 2^{-128}$$

$$\begin{array}{c} \cdot 1 \rightarrow 2^{-1} \\ \cdot 01 \rightarrow 2^{-2} \\ \hline \end{array}$$

- largest -ve denormalized number
 -2^{-135} Ans

Note • By default number will be in normalized form because there are very less chances of denormalized number

Ex $\boxed{81E1M1} \rightarrow$
 ↓ ↓
 1 excess 128

$(3.5)_{10} \rightarrow$ Convert into normalized form

1101 \Rightarrow Binary

$0.111 * 2^2 \Rightarrow$ Normalized number

Note Rule \Rightarrow After point number should be '0' in denormalized form

$0.111 * 2^2 \Rightarrow$ Normalized number

many form
 of denorm-
 lized
 number

$0.0111 * 2^3 \Rightarrow$ denormalized number $0.01110 * 2^4 \quad \underline{\hspace{2cm}} \quad 4$ $0.000111 * 2^5 \quad \underline{\hspace{2cm}} \quad 4$ $0.0000111 * 2^6 \quad \underline{\hspace{2cm}} \quad 4$

$0.111 * 2^3 \Rightarrow$ 0100 0011 0111 000
41B0H

S E M $01000001011000 \Rightarrow 41B8H$

Mantissa's first bit is '0' so it is denormalized number

$$\circ 0111 * 2^{128+3-128}$$

$$\circ \underline{0111} * 2^8 \Rightarrow (11.1)_2 = (3.5)_{10}$$

Normalized this no. to get the ans

$$\Rightarrow 0111 * 2^2 \Rightarrow \text{Normalized form}$$

$$\circ .00111 * 2^4$$

$$0111.00111 * 2^{128+4-128}$$

$$\circ \underline{00111} * 2^7 = (0011.1)_2 = (3.5)_{10}$$

Normalized this number to get ans

$$0111 * 2^2 \Rightarrow \text{Normalized form}$$

* [S | E | M] → 20bit
 \downarrow \downarrow
 16bit 11bit excess 1024

largest normalized number

$$\left(1 - \frac{1}{2}^{20}\right) * 2^{1023}$$

Smallest positive de normalized number

$$\left(1 - \frac{1}{2}^{127}\right) * 2^{-127} \Rightarrow \left(1 - \frac{1}{2}^{127}\right) * 2^{1022}$$

$$\text{Smallest} = 2^{-7} * 2^{-127} \Rightarrow 2^{-134}$$

- $\boxed{S|E|M} \rightarrow 52\text{bit}$
 16bit excess 3023
 10 bit, excess 511
 +ve normalized no.

$$\left(1 - 2^{-21}\right) * 2^{512}$$

$$\text{Smallest} + \text{ve} \cdot 1 * 2^{-511} \rightarrow 2^{-512} \Rightarrow 0.5 * 2^{-511}$$

Largest +ve denormalized number

$$1 - 2^{-20} * 2^{-511} \Rightarrow 2^{-532}$$

Smallest

$$2^{-21} * 2^{-511} \rightarrow 2^{-532}$$

- $\boxed{S|E|M} \rightarrow 52\text{bit}$
 16bit excess 3023
 (11bits)

* ~~Denormalized~~

E-18

$$(1 - 2^{-52}) * 2^{2047-1023}$$

$$0.1 * 2^{-1023} \Rightarrow 2^{-1024} \Rightarrow 0.5 * 2^{-1023}$$

Largest +ve denormalized no.

$$(1 - 2^{-51}) * 2^{1023}$$

$$\underline{\text{SPD}} \rightarrow 2^{-52} * 2^{-1023} = 2^{-1075} \quad \underline{\text{Ans}}$$

o

S	E	M
---	---	---

 \rightarrow 13 bits
 ↓ ↓
 1 bit ex 32(6 bits)

Largest +ve normalized no.

$$1 - 2^{-32} * 2^{31}$$

$$\underline{\text{SPN}} \rightarrow 0.1 * 2^{-31} \Rightarrow 0.1 * 2^{-32} \Rightarrow 0.5 * 2^{-31}$$

Smallest +ve denormalized number

$$2^{-13} * 2^{-32}$$

Largest:

$$(1 - 2^{-12}) * 2^{30} \quad \underline{\text{Ans}}$$