

* IEEE 754 Notations

Lecture 13

Part 1

Single Precision
32 bit (4 byte)

double precision
64 bit (8 byte)

* format

$\boxed{S \mid E \mid M} \rightarrow 32\text{bit}$
 S → 1-bit
 E → 8-bit
 M → 23-bit
 ex - 127

* format

$\boxed{S \mid E \mid M} \rightarrow 52\text{bit}$
 S → 1-bit
 E → 11-bit
 M → 32-bit
 ex - 1023

* IEEE Normalized

do not store
(implicit)
 $1.0 \times \dots$

Ex : $(3.5)_{10} \rightarrow$ Convert it into IEEE single precision normalized number

Step 1 \Rightarrow Binary = $(11.1)_2$

Step 2 \Rightarrow Normalized = $(11.1)_2 \Rightarrow 1.11 \times 2^1$

sign \swarrow 01000000 \searrow exponent 11000
 40600000H \swarrow mantissa Any

* $(0.5)_{10} \Rightarrow (0.1)_2 \Rightarrow 1.0 \times 2^{-1}$

00111110.0...

3F000000H Ans

* $(6.25)_{10} \Rightarrow (110.01)_2 \Rightarrow 1.1001 \times 2^3$

0100000011001...

40C8000H Ans

* 3FC0000H \rightarrow IEEE Single Precision, decimal?

001111110000000...
S E M

$+ 1.0 \times 2^{127-127} \Rightarrow 1.0 \times 2^0 \Rightarrow (1.0)_2$

$(1.5)_{10}$ Ans

* 40C8000H \rightarrow IEEE Single Precision, decimal?

0100000011001000...
E M

$+ 1.1001 \times 2^{129-127} \Rightarrow 1.1001 \times 2^2 \Rightarrow (110.01)_2$

$(6.25)_{10}$ Ans

* Denormals \Rightarrow exponent all 0's

if exponent is all 0 means
it is normalized

"1" is not implicit

* Ex: $s \ e \ m$
0 0000000 0000 - - -

• all exponent 0 means it is
denormalized and "1" is not
implicit.

0 0000000 000 - - -

$$+ 0 * 2^{0-127} \Rightarrow 0$$

• 1 0000000 000 - - -

$$- 0 * 2^{0-127} = -0 \quad \underline{\text{Ans}}$$

$$\begin{aligned} 0 & \quad 00000000H \Rightarrow +0 \\ 80000000H & \Rightarrow -0 \end{aligned}$$

* These two are special number

S E M
0 111111 000 -

$+ \infty$

$7F800000H = +\infty \Rightarrow +ve \text{ infinity}$

S E M
1 111111 000 - - -

$FF800000H \Rightarrow -\infty \Rightarrow -ve \text{ infinity}$

both are special number

* $0111111 \neq 0$ - - Not a number

$1111111 \neq 0$ - - Not a number

S.P.N	S	E	M
$+0$	0	all 0's	all 0's
-0	1	all 0's	all 0's
$+\infty$	0	all 1's	all 0's
$-\infty$	1	all 1's	all 0's
NAN	0/1	all 1's	$\neq 0$
Denormals	0/1	all 0's	$\neq 0$

Part 2

Largest * normalized number

- Largest +ve

S E M
0 / 1111111 / 011111 - - -

$$7 \text{ F}7\text{FFFFF} \text{ H} \quad \underline{\text{Ans}}$$

$$1.0 / 111... \times 2^{254-127} \Rightarrow 2-2^{-2^3} \times 2^{127}$$

- Smallest +ve

S E M
0 00000000 0000 - - -

0 08000000 H Ans

$$1.0 \times 2^{1-127} \Rightarrow 1.0 \times 2^{-126} \Rightarrow 2^{-126} \quad \underline{\text{Ans}}$$

- Smallest -ve

- $(2-2^{-2^3}) \times 2^{127}$
FF7FFFFF H Ans

- Largest -ve

- $1.0 \times 2^{-126} \Rightarrow$
80800000 H Ans

* Denormalized number

- Largest +ve

s E
0 00000000 1111--1

007FFFFFH M exception

$$+ \quad 1 - 2^{-23} * 2^{1-127} \Rightarrow 1 - 2^{-23} * 2^{-127}$$

o .

Exception: when all the exponent are 0's
 we'll consider it as "1"

- Smallest ~~+ve~~ -ve

~~0~~100000000 1111--.

807FFFFFH

- Smallest +ve

s E M
0 00000000 0000--1

$$00000001H \underline{M} \Rightarrow 2^{-23} * 2^{1-127}$$

- Largest -ve

1000000000--1

80000001H

$$-2^{-23} * 2^{-126} \Rightarrow -2^{-149} \underline{M}$$

$$\Rightarrow 2^{-23} * 2^{-126}$$

$$\Rightarrow 2^{-149} \underline{M}$$

Background :-

- * Find the smallest diff b/w two consecutive normalized +ve number

SPN

$$0 \ 00000001 \ 0000 - - - \\ 00800000H \Rightarrow 1.0 * 2^{1-127} = 2^{-126}$$

\downarrow Next consecutive number

$$0 \ 00000001 \ 0000 - - - 1 \\ 1 + 2^{-23} * 2^{-126}$$

gap in consecutive no. $\Rightarrow (1 + 2^{-23}) * 2^{-126} - 1 * 2^{-126}$

$$2^{-126} + 2^{-149} - 2^{-126}$$

* Minimum gap : 2^{-149}

* Max gap :

$$0 \underline{11111110} 11111 - - - 1 \\ 7F7FFFH \Rightarrow 2^{-2-23} * 2^{127}$$

$$0 \underline{11111110} 0000 - - - 0 \\ 1.0 * 2^{254-127} = 2^{127}$$

$$0 \underline{11111110} 0000 - - - 1 \\ (1 + 2^{-23}) * 2^{127} \Rightarrow 2^{104}$$

$$\underline{\text{gap}} \Rightarrow 2^{127} + 2^{104} - 127 \Rightarrow 2^{104} \text{ Ans}$$

$2^{-149} \rightarrow$ Min gap b/w normalized number

$2^{104} \rightarrow$ Max — a — 1

IEEE \rightarrow non uniform distribution

* gap in denormals

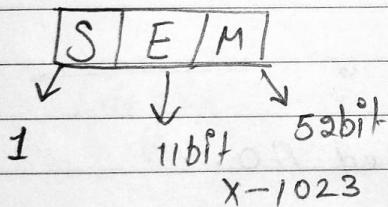
$0\ 0000\ 0000\ 000\dots1$ $\xrightarrow{\text{to maintain the minimum}}$ gap

*

IEEE - Notations

Part 3

o Double Precision



$$(1.5)_{10} \Rightarrow (1.1)_2$$

$\underbrace{0\ 0}_{3}\ \underbrace{11111111}_{F}\ \underbrace{0\dots0}_{8\ 00000000H}$ Any

$$(.25)_{10} \Rightarrow (.01)_2 \Rightarrow 1.0 \times 2^{-2}$$

$\underbrace{0\ 011111101000\dots0}_{3\ F\ D\ 0\ 000000000000H}$

(-40.1) \rightarrow find single and double precision

$(101000.00011)_2 \Rightarrow 1.010000001111$
Normalized form

Single Precision

11000100010000000110011001100110
C 9206666H

* double precision

1100000001000100000011001100
C 0440CCCCCCCCCCCCH A4

Largest Positive Normalized No.

011111110 111 -- 1

7FFFFFFFAFFFFFH

$$(2 - 2^{-52}) * 2^{2046-1023}$$

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

Smallest positive -ve normalized no.

1 11111110 1111 ~ ~ 1

f. F E F F F F F F F F F F F F F F H

$$-(2 - 2^{-52}) \times 2^{2046 - 1023} \Rightarrow$$

$$-2^{-2^{-52}} \times 2^{1023}$$

• Smallest positive normalized no.

0 0000000000010000 ~ ~ .

$$1.0 \times 2^{-1023} \Rightarrow 2^{-1022} \text{ Ans}$$

• Largest -ve no.

$$-2^{-1022} \text{ Ans}$$

* Denormals

smallest

• Largest +ve denormalized no.

0 0000000000000000 ~ ~ 1

00000000000000001 H

$$2^{-52} \times 2^{1-1023} \Rightarrow 2^{-1074}$$

Ans

- Largest -ve denormalized no.

$$-2^{-1074} \text{ Ans}$$

- Largest +ve normalized no.

0 0000000000 1111 ... 1

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

- Smallest -ve denormalized no.

$$\{-1 - 2^{-52}\} * 2^{-1022} \text{ Ans}$$

- * Max gap b/w two consecutive normalized numbers

0 111111110 0000 ... 0

$$1.0 \times 2^{2046 - 1023} \Rightarrow 2^{1023}$$

Next consecutive no.

0 1111111110 000 ... 1

$$1.0 * 2^{-52} * 2^{2046 - 1023}$$

$$(1 + 2^{-52}) 2^{1023}$$

$$2^{1023} + 2^{971}$$

Max gap: 2^{971}

$$2^{971} = 10^n$$

$$971 = n * \log_2 10$$

$$n = 320$$

$$10^{320} \text{ Approx}$$

Min gap:

$$0.0000000001000\ldots$$

$$1.0 * 2^{1-1023} \Rightarrow 2^{-1022}$$

Next Consecutive no

$$0.0000000001000\ldots 1$$

$$(1.0 + 2^{-52}) * 2^{1-1023}$$

$$2^{1022} + 2^{-1074}$$

$$\text{Min gap} \Rightarrow 2^{-1074}$$

* Denormals \leftarrow min gap, same
max

0 000000000000 0000--1

$$2^{-52} * 2^{1-1023} \Rightarrow 2^{-1074} \quad \underline{\text{Ans}}$$

Next consecutive no.

0 000000000000 0000--10

$$2^{-51} * 2^{1-1023} = 2^{-1073}$$

gap $2^{-1073} - 2^{-1074}$

$$2^{-1073} - 2^{-1} * 2^{-1073}$$

$$\frac{2^{-1073}}{2} \Rightarrow 2^{-1074}$$

Representation of number system
is completed