

Sec 6

General

$$\begin{cases} 0.125 = 1.25 \times 10^{-1} \\ 6.000 = 6 \times 10^0 \\ 5\,000\,000 = 5 \times 10^6 \end{cases}$$

Scientific notation

Significant mantissa

$(+)$ 1.25 $\times 10^1$

bits bits

How Computers

Sign exponent mantissa/significand

no of bits

Size of mantissa determines — Precision ✓
 Size of exponent determines — Range

Floating point

Significant

Example

decimal

$32 = 2^5 = 1.0 \times 2^5$ or 0.1×2^6

Exp Mantissa / Significand

$0.1 = 0.5$

$32 = 2^5$

$= 1.0 \times 2^5$

$= 0.1 \times 2^6$

Significand Exponent

We prefer this notation where the decimal point is implied.

Significand = 0.1

0.1

$110 = 6$

Mp. Significand.

$0.01 = 2^{-2}$

2^{-3}

S E Significand

$0.01 = 2^{-2}$

2^{-3}

$0.1 = 2^{-1}$

To resolve:

representation).

- First digit of the significand must be 1.
(normalization)

- $0.1xxxxx \dots$

$$4.5_{10} = 100.1 \times 2^0$$

$$= 10.01 \times 2^1$$

$$= 1.001 \times 2^2$$

$$= 0.1001 \times 2^3 \leftarrow \text{normalization}$$

- For negative exponents we use a "biased" exponent
Bias - number that is approximately
on the range of values

0 15

5 bit exponent = 32 bits
mid value = 16 (bias)

Excess 16 bit representation

In this model - exponent values less than 16 are negative, fractional nos.

$$32 = 16 \times 2^5 = 0.1 \times 2^6$$

we excess 16 biased exponent.

$$6 \text{ will be represented as } = 6 + 16 = 22$$

so we have Sign Exponent Mantissa

0	10110	10000000
---	-------	----------

bias on excess 16 reprt. $2^{22-16} = 2^6$ $0.1 = \pm 0.1 \times 2^6$

Ex. $0.0625 = 1.0 \times 2^{-4} = 0.1 \times 2^{-3}$

now in terms of excess 16 biased exponent = $16 - 3 = 13$

sign

bit Exponent

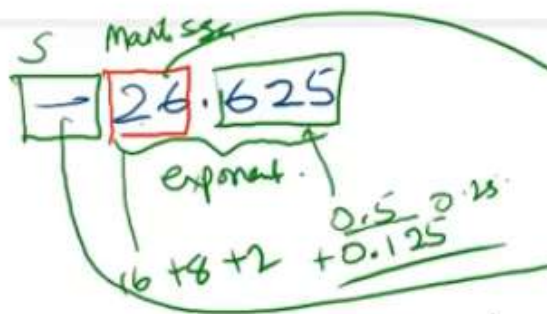
Significand

0	01101	10000000
---	-------	----------

⊕ The first digit of the significand should be / must be 1.
 with no ones to the left of the radix point.
 - ^{this} process ~~is~~ is called normalization

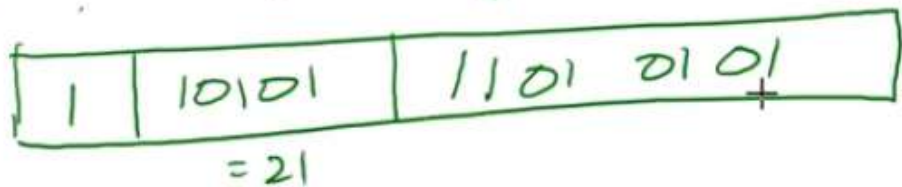
0.1xxxxx. all significands or
mantissa ✓

$$4.5 = \text{~~100~~ } 100.1 \times 2^0 = 1.001 \times 2^2 \text{ or } = \boxed{\underline{0.1001 \times 2^3}}$$



$$\begin{aligned}
 & \underline{11010.101} \\
 &= 0.\underline{11010101} \times \frac{5}{2} \\
 & \quad \downarrow \text{significand} \quad \text{exponent} \\
 & \quad \quad \quad 16 + 5 = 21
 \end{aligned}$$

-ve sign



Standard

Single precision standard is called IEEE 754. double precision

IEEE 754 Single precision floating point standard uses.
(32 bits) 8 bit exponent, 23 bit for significand
Bias = 127

IEEE 754 double precision floating point standard uses.
(64 bits) 11 bit exponent, 52 bit significand.
Bias = 1023.

In both — the significand has an implied 1 to the left of the radix point.
format is 1.XXXX instead of 0.1xxx

The largest no representable by 7 bit is
1111111 = 127

$$4.5 = 0.1001 \times 2^3$$

$$\text{IEEE Std format} = \underbrace{(1)}_{\text{implicit}} \underbrace{001}_{\text{mantissa}} \times 2^2$$

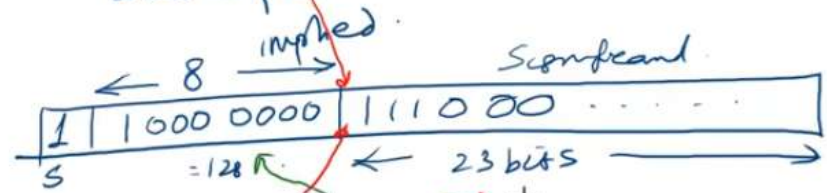


$$S = \underline{\underline{011}}$$

After the
0.5
+ 0.25

$$\begin{aligned} -3.75 & \quad \text{IEEE single precision.} \\ -3.75 &= -11.11 = -1.111 \times 2^1 \quad (\text{one shift}) \\ &= \underline{\underline{1}}.111 \times 2^1 \end{aligned}$$

$$\text{Bias is } 127 \therefore \text{no is } = 127 + 1 = 128$$



$$\begin{aligned} & \text{value} \quad \text{Bias} \\ & - (1).111 \times 2^1 + (128 - 127) = -1.111_2 \times 2^1 = -11.11 = -3.75 \end{aligned}$$

$\boxed{\text{Range}}$ of FP number
 for a 32 bit - 8 bit exponent
 $\pm 2^{256} \approx 1.5 \times 10^{77}$ } largest no

$\boxed{\text{Accuracy}} \rightarrow 23 \text{ bit significand} \rightarrow 2^{-23} \approx 1.2 \times 10^{-7}$
 About 6 decimal places.

$\boxed{64 \text{ bits}}$ Range $\pm 2^{2048}$
 Accuracy $\frac{1}{2^{-52}}$

$$2^{10} = 1024 \sim 1K$$

bits

Decimal value of an IEEE no

$$(1-2s) * (1.f) * 2^{e-bias}$$

$s, f, e =$ decimal.

$(1-2s) = 1$ or -1 Whether $s = 0$ or 1
fractional field; we add the implicit 1.

$bias = 127$ or 1023
SBP. DBP

$$\begin{array}{lcl}
 \begin{array}{c} 1 \\ \hline s \end{array} & \begin{array}{c} 0111100 \\ \hline e \end{array} & \begin{array}{c} 11000 \dots 0 \\ \hline 0.5 + 0.25 \\ 0.75 \end{array} \\
 s = 1, (-) & 124. & \\
 (1-2) * (1+0.75) * 2^{124-127} & & \\
 = -1 * 1.75 * 2^{-3} = -\frac{1.75}{8} = -0.21875
 \end{array}$$

Exo

347.625
↓
Binary

101

0.5 x 1
+ 0.25 x 0
+ 0.125 x 1

Binary to binary
Decimal to FP-

101011011.101 x 2⁰

1.0101101101 x 2⁸

S = 0, exp = 8 + 127 = 135, 0101101101...

+

Special values:

E	F	meaning	Explan:
0000 0000 0000 0000	000...0 XXXXX	0 Valid no.	+0.0 or -0.0. (-1) ^S × 2 ⁻¹²⁶ (0xf) (unnormalized)
1111 1111 1111 1111	000000... XXXXXXXX	∞ NaN. (Not a number)	

Find the range, smallest nos, largest nos, Precision -
 NaN, ∞ , overflow, underflow

(Not a number)

E	Real Exp.	F	Value
0000 0000	Reserved	0000 xxx...	0
0000 0001	-126	{	
1111 1111	Reserved	000... xxx...	$(-1)^S 2^E \cdot (L.F)$ Normalized ∞ NaN

bonnet
a cross
venn

Smallest and highest number in 32 bit floating point

int $x = 33554431 \leftarrow$ In float it can't be represented,
float $y = -$ (causes issue)
 $0.10 = 0.0001100110011 \dots$



This 0.10 is not
represented in floating
point representation

HW floating point x multiplication & Addition:-

