

1) IEEE 754 floating Representation

1) Single Precision - 32 bit

S E M  
1 8 23 b.i.h.

1) S (sign bit)

a)  $e$  (exponent)

3) M (mantissa (or) fractional part)

2) double Precision - 64 bit.

S E M

1 " 52 bits.

⇒ Single Precision to decimal conversion.

$$x_d = -1^8 x (1+M) \times 2^{E-127}$$

⇒ double precision to decimal conversion.

$$x_d = -15 \times (1 + 11) \times 2^{E-1025}$$

PRECISION: value (of) number

it represent the how much closer to the actual value.

Ex:  $1/6 = 0.166666666666 \dots$

## Single Precision.

$$\frac{s}{e}$$
$$\begin{array}{r} \epsilon \\ \hline 011100 \\ 124 \\ (-3) \end{array}$$
M

$(2A A A A A) = \text{AAAAA}$

$$\Rightarrow (-1)^0 \times 2^{-3} (1 + 0.33333334) = 0.166666675$$

Double      Precision:

S  
0

E  
3FCh  
(1020)

M  
5555555555555555h  
~~0000000000000000~~

$$= (-1)^0 \times 2^3 \times (1 + 0.3333333333333333)$$
$$= 0.1666666666666666$$

So we can see that there is no change in sign bit (0) exponent so only effect the precision ~~is~~ is fraction part (mantissa).