

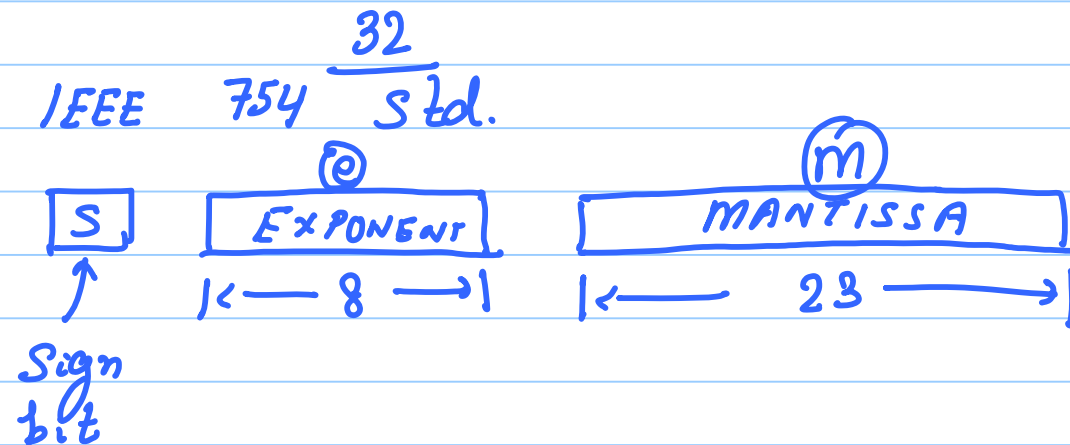
Sep-5

Note Title

05-09-2012

Represent floating point in binary

$[2.96]_d$ $[5 \times 10^{15}]$ $[-1.923 \times 10^{-27}]$



$$13 = 8 + 4 + 0 + 1$$

$$2^3 + 2^2 + 0 + 2^0$$

$$\boxed{1 \ 1 \ 0 \ 1}$$

$$1.75$$

$$1 + 0.5 + 0.25$$

$$2^0 + 2^{-1} + 2^{-2}$$

Subset of numbers of the form: $1.x$

$$1.x$$

$$1 + 2^{-1}(x_1) + 2^{-2}(x_2) + \dots$$

Any Number (N) (decimal)

$$N = (1.x) \times 2^{\text{exp}} \times \boxed{(\pm 1)}$$

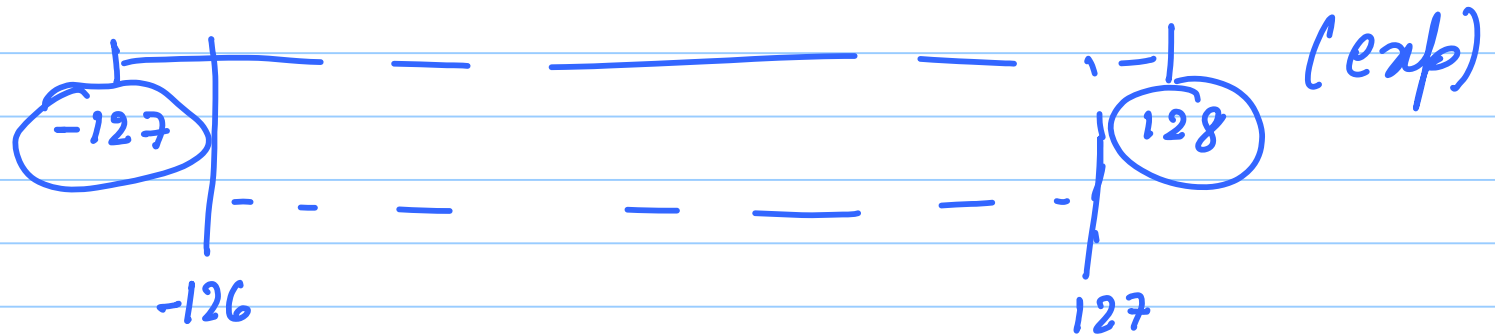
Sign.

Sign ± 1 (sign bit) $\textcircled{1}$

$e \rightarrow$ biased notation. $\boxed{8}$

$$\text{exp} = e - \text{BIAS}$$

(0, .. 255)
(127)



Mantissa {23 bits}

$$N = (-1)^s 1.m \times 2^{(e-127)}$$

Special case:

	$e=0, m=0$	N 0	
{	$s=0, e=255, m=0$	∞	
	$s=1, e=255, m=0$	$-\infty$	
{	$e=255, m \neq 0$	NAN	{ NAN+ ∞ = NAN NAN $\times 2$ = NAN }
	$e=0, m \neq 0$	Denormal Numbers.	

1) Floating point numbers are approximate.

2) There is some amount of error in representation.

$$x > 0, \quad x/2 > 0$$

Normal floating point numbers.

$$\text{Largest: } \left\{ \left[\overbrace{1.1 \dots 1}^{23} \right] \times 2^{127} \right. \\ \left. \begin{array}{l} \text{(+ve)} \\ \text{binary} \end{array} \right\}$$

(Normal)



$$(1 + 2^{-1} + \dots + 2^{-23})$$

$$(2 - 2^{-23}) \times 2^{127} \\ = 2^{128} - 2^{104}$$

Smallest:
(+ve)
(Normal)

$$1 \cdot \underbrace{0 \dots 0}_{23} \times 2^{-126} = 2^{-126}$$

```

prog. { x = 2-126
      if (x/2 > 0)
          print ("hi");
  
```

Denormal number -

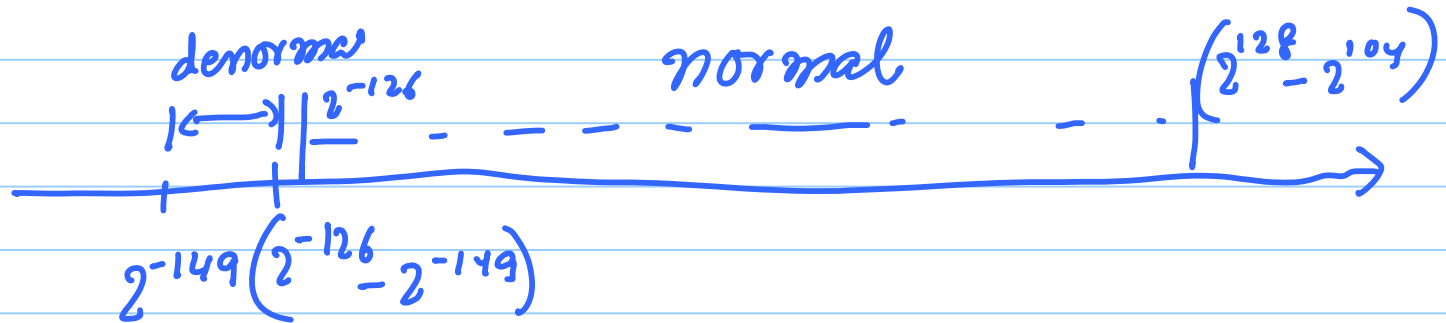
$$\boxed{e=0} \quad \boxed{m \neq 0}$$

$$\text{Form: } (-1)^s \cdot 0.M \times 2^{-126}$$

Largest denormal number:
(+ve)

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

Smallest (+ve) denormal number: 2^{-149}



Is, $(\Delta x > 0)$

$$\left[\begin{array}{l} x + \Delta x > x \\ \nearrow 2^{100} \quad \nearrow 2^{-100} \end{array} \right.$$

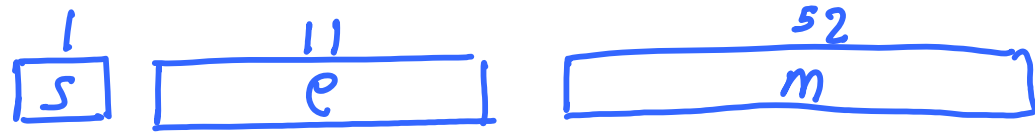
printf (x + Δx - x);

$x = 2^{-149}$
if (x / 2 > 0)
printf ("hi");

float → 32 bits

double → 64 bits

double.



$$\text{exp} = e - 1023$$