

5-6 (5 bit)

$$\begin{array}{r}
 00101 \quad (5) \\
 + 11010 \quad (-6) \\
 \hline
 11111 \quad (-1)
 \end{array}$$

$$\begin{array}{r}
 6 \\
 00110 \\
 11001 \\
 +1 \\
 \hline
 (-6) 11010
 \end{array}$$

$$\begin{array}{r}
 00101 \quad 5 \\
 + 11111 \quad -1 \\
 \hline
 \text{carry } \times 00100 \quad 4
 \end{array}$$

$$\begin{array}{r}
 01100 \quad 12 \\
 00110 \quad +6 \\
 \hline
 10010 \quad -14 \\
 \text{overflow}
 \end{array}$$

$$\begin{array}{r}
 \textcircled{-14} \\
 10010 \\
 11001 \\
 +1 \\
 \hline
 01110 \\
 14
 \end{array}$$

Real numbers

5 bits fixed point

3 bits integer
2 bits fraction

$2^2 \ 2^1 \ 2^0 \ 2^{-1} \ 2^{-2}$
1 0 1 . 0 1
5.25

$(0, .25, .5, .75)$
 $1 \times 2^2 = 4$
 $0 \times 2^1 = 1$
 $1 \times 2^0 = 1$
 $0 \times 2^{-1} =$
 $1 \times 2^{-2} =$
 $\frac{1}{2^2} = \frac{1}{4}$

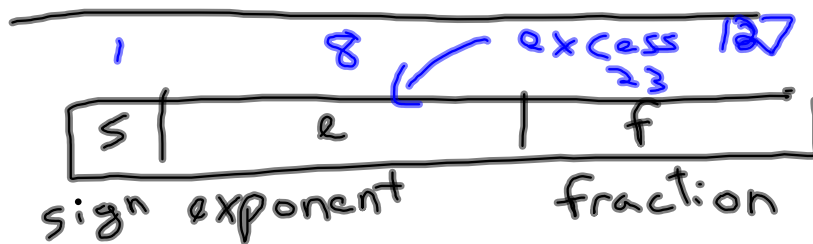
Floating Point

0.003

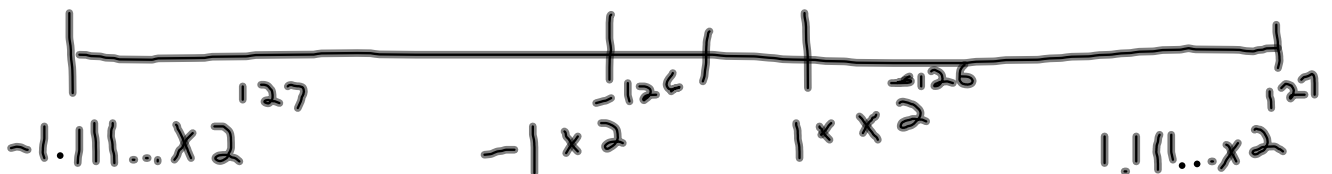
3.0×10^{-3}

50

5.0×10^1



$$v = (-1)^s \times f \times 2^e$$



Format

Normalized

\pm	$0 < e < \text{max}$	f
-------	----------------------	-----

 $\pm 1.f \times 2^e$

Zero

\pm	00...0	000...	0
-------	--------	--------	---

Denormalized

\pm	00...0	f : any non-zero
-------	--------	--------------------

 $0.f \times 2^{-127}$

Infinity

\pm	111...	1	000... 0
-------	--------	---	----------

 $\pm \infty$

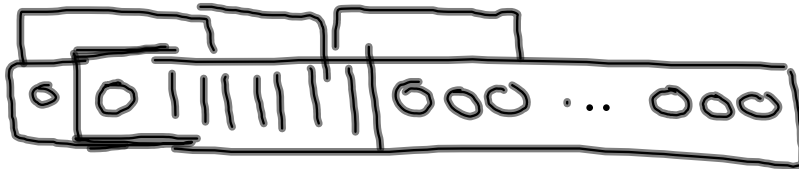
Not a number

\pm	111...	1	non-zero
-------	--------	---	----------

NaN

	float	double
sign	1	1
exp	8	11
frac	23	52
exp system	Excess 127	Excess 1023
Decimal Range	$\sim 10^{-38}$ to 10^{38}	$\sim 10^{-308}$ to 10^{308}
Smallest denorm	10^{-45}	10^{-324}

Decimal 1



$$1_{10} = +1.0 \times 2^0$$

$$0x3F800000$$

$$0x42E48000$$



$$e = 133 - 127 = 6$$

$$+ 1.11001001 \times 2^6$$

$$1110010.0100000000..$$

$$\begin{array}{r} 1 \\ 64 \\ + 32 \\ + 16 \\ + 2 \\ \hline 114 \end{array}$$

$$+ 114.25$$

0x80400000

1000 0000 0100 0000 0000...

-127

base 2

$$-0.1 \times 2^{-127}$$

$$= -0.5 \times 2^{-127}$$

Denormalized \uparrow

0x758E1000

0111 0101 1000 1110 0001 0000...

X (nevermind)

0x7F8E1000

0111 1111 1000 1110 0001 0000...

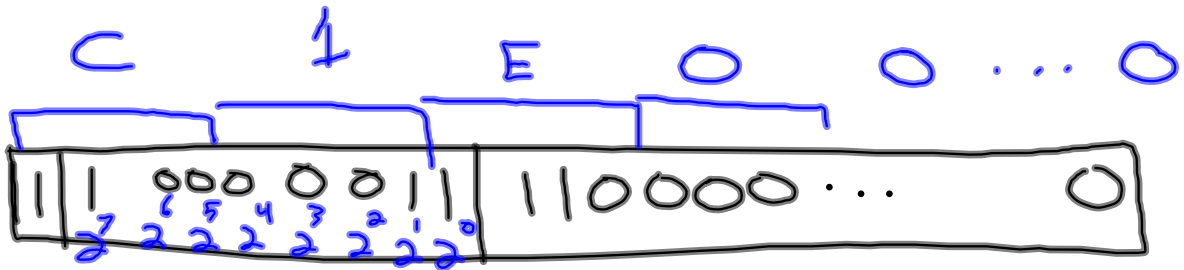
+

e

f

NAN

$$-3.5 \times 2^3$$



exp $4 + 127 = \underline{131}$ f - $3.5_{10} = 11.1_2$
convert unsigned $(1.11 \times 2^1) \times 2^3$

0xC1E00000

1.11×2^4