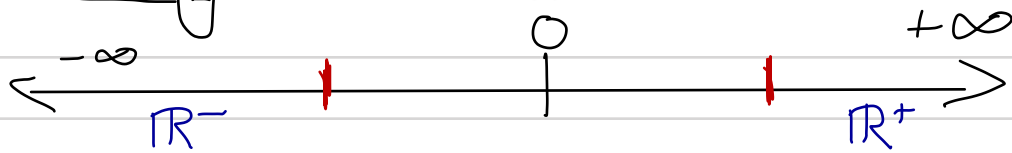


# Rounding and Guard bits



Rounding up:

Rounding down:

D. round 0:

D. Round  $+\infty$ :

D. round  $-\infty$ :

D. round  $\bar{0}$ :



In  $Q.O.m$  subtraction and addition adding a guard bit after the LSB ( $Q.O.(m+1)$ ) improves precision.

Add guard bit before normalisation step. Round to  $Q.O.m$  format towards zero.

$$0.101e001 - 0.011e000$$

Without guard bit:

$$\begin{array}{r} 0.101e001 \\ - 0.001e001 \\ \hline 0.100e001 \end{array} \leftarrow \text{information lost during justification}$$

With guard bit:

$$\begin{array}{r} 0.\overset{1}{\cancel{0}}\overset{1}{\cancel{0}}e001 \\ - 0.0011e001 \\ \hline 0.0111e001 \end{array}$$

→ 1.00

→ 0.75

→ 0.875 ✓

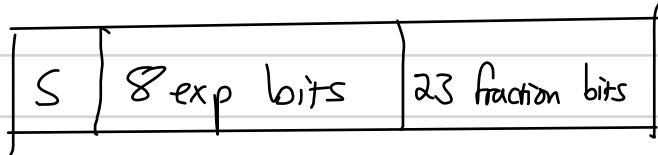
= 0.011e001

= 0.111e000

# IEEE Floating point format

First proposed in 1985: now present in every mobile device regardless of architecture.

THE format used to store data on fixed or removable media.



1 sign bit

8 exponent bits (Q7.0 signed)

23 fractional bits Q0.23 (**unsigned**)

Exponent has bias of -127

$$E_{min} = 1 - 127 = -126$$

$$E_{max} = 254 - 127 = +127$$

(offset of 127 subtracted from stored value)

Stored exponent of 0x00 and 0xFF are interpreted differently.

| Exponent    | Mantissa 0 | Mantissa 1 | Equation                                  |
|-------------|------------|------------|---|
| 0x00        | +0/-0      | sub-normal | $(-1)^{s_b} \times 2^{-126} \times 0.FB$  |
| 0x01 → 0xFE | Normalised | Normalised | $(-1)^{s_b} \times 2^{E-127} \times 1.FB$ |
| 0xFF        | +∞/-∞      | NaN        | Quiet signalling                          |

$$+1_{10} \Rightarrow \underbrace{0}_{(s_b)} \ 0 \ 111111 \ 0000 \dots 000$$

$$\Rightarrow 3F80 \ 0000$$

$$-2 \Rightarrow (000 \ 0000$$

$$1100000 \ 00 \dots 000$$

DEADBEEF = 1101 1110 1010 1101  
IEEE 1011 1110 1110 1111

$$SB = 1$$

$$Exp = 1011101 = 189 - 127 = 62$$

$$1.0101101101111011101111 = 1.357893309$$

$$DEADBEEF = -6.259853398707798016 \times 10^{18}$$

$$y = (-1)^{sb} \times 2^{E-127} \times 1.FB$$

$$DEADBEEF = (-1)^1 \times 2^{189-127} \times 1.3578933091758433$$