

2.5 Floating point division

$$y_1 = M_1 \times 2^{E_1}$$

$$y_2 = M_2 \times 2^{E_2}$$

$$z = (M_1 \div M_2) \times 2^{(E_1 - E_2)}$$

1. Normalise y_1 & y_2
2. Unsign if necessary
3. Divide mantissas.
4. Subtract exponents.
5. Resign (if necessary)
6. Normalise

Example

$$1.0100110e0010$$

$$\div 0.0010110e0100$$

Unsign: $1.0100110e0010 \Rightarrow 0.1011010e0010$ (Normalised)

Normalise #2: $0.1011000e0010$

Example continued:

$$M_1 = 0.1011010$$

$$M_2 = 0.1011000$$

00000001	000001	
01011000	01011010	00000000
01011010		↓ ↓ ↓ ↓ ↓
-01011000		00000000
00000000		0110000
		01011000
		01011000

discard

error in division

$$M_1 \div M_2 = 0010000001 \text{ e } 0000$$

Reformatting into Q0.7

$$\therefore M_1 \div M_2 = 0.10000001 \text{ e } 0001$$

Subtract exponents:

$$E_1 - E_2$$

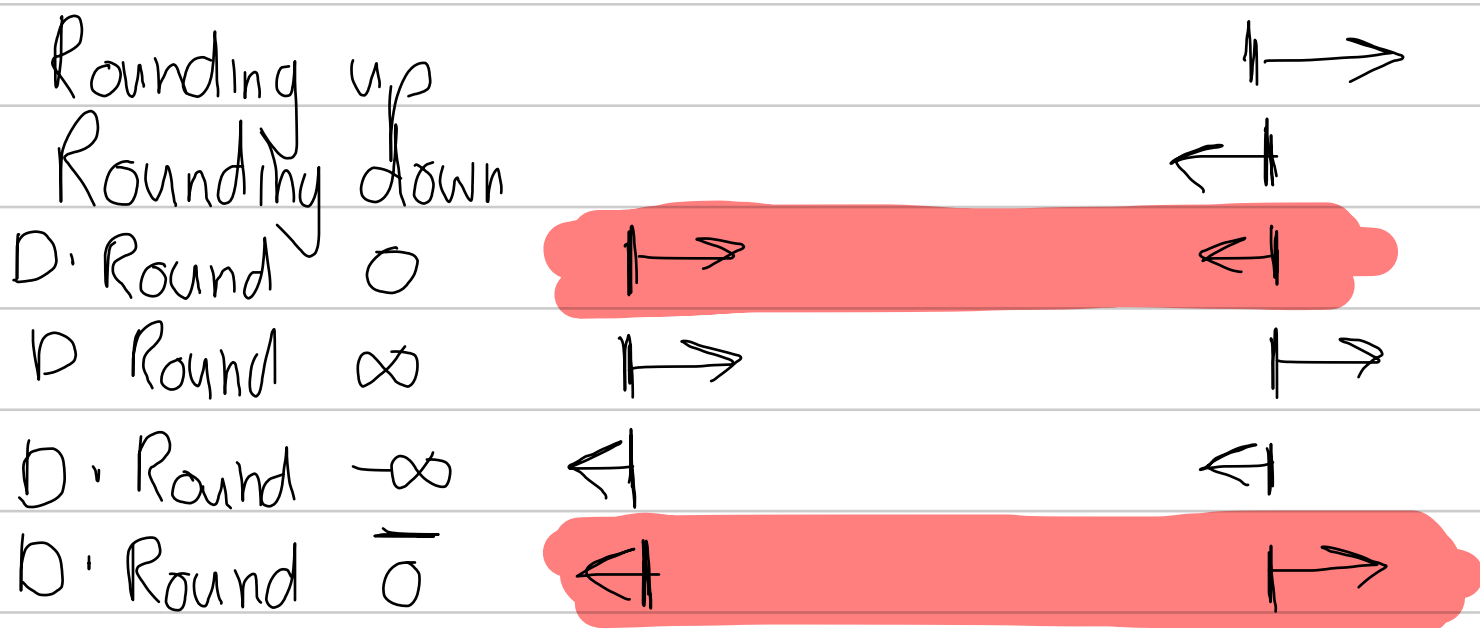
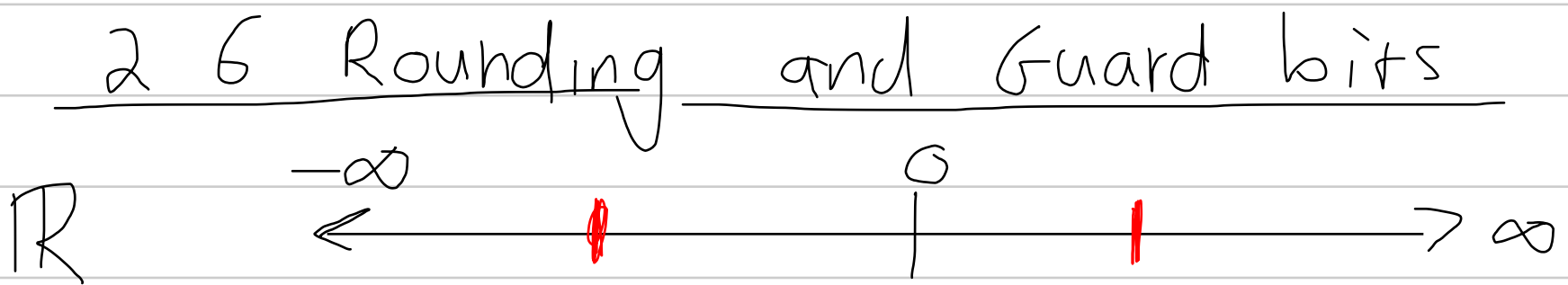
$$\begin{array}{r} 0010 \\ -0010 \\ \hline 0000 \end{array}$$

Example continued

$$0.1000001e(0001+0000)$$

Resign since y_1 was negative

$$z = 1.0111111e0001$$



Guard digits

In $Q \cdot m$ subtraction/addition adding a guard bit after the LSB improves precision.

Add guard digit before normalisation step
Round guard digit towards 0.

* $0.101e001 - 0.011e000$

Without guard digit

$$\begin{array}{r} 0.101e001 \\ - 0.001e001 \\ \hline 0.100e001 \end{array}$$

With guard digit

$$\begin{array}{r} 0.1010e001 \\ - 0.0011e001 \\ \hline 0.1010e001 \end{array}$$

$$\begin{aligned} 0.101e001 &= 1.25_{10} \\ 0.011e000 &= 0.375_{10} \end{aligned}$$

$$1.25_{10} - 0.375 = 0.875$$

$$0.100e001 = 1_{10}$$

$$0.110e000 = 0.875$$

2.7 IEEE floating point format

First proposed in 1985. Now present in every mobile device, PC etc... Is the format used to store data on fixed/permanent storage.

S(±)	8 exp bits	23 fraction bits
------	------------	------------------

1 Sign bit

8 Exponent bits Q7.0 signed

23 Fractional bits Q0.23 **unsigned**

Exponent has a 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 interpreted differently.

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

$+1_{10} \rightarrow 3F80\ 0000$ (IEEE FP)