

## 2.0 Floating point numbers

$$y = M \times x^E \leftarrow \text{exponent. } \in \mathbb{Z}$$

↑                      ↖  
mantissa              base  
(fixed point)

$$M \equiv Q_{0.m} \text{ signed}$$

$$E \equiv Q_{n.o} \text{ signed}$$

$$1.38 \times 10^{-4}$$

$$2.71 \times 10^3$$

$$1.38 E - 4$$

$$2.71 E 3$$

### Example

$$\begin{aligned} 1011e011 &= (-1 + 0.25 + 0.125) \times 2^3 \\ &= -0.625 \times 2^3 \\ &= -5.00_{10} \end{aligned}$$

$$0.101e011$$

$$\begin{aligned} 0.101_2 \times 2_{10}^3 &= 01.01_2 \times 2_{10}^2 = 010.10_2 \times 2_{10}^1 = \\ &= 0101 \end{aligned}$$

$$\begin{aligned} 1.50 \times 10^3 &\Rightarrow 150 \times 10^3 &= & 15.0 \times 10^2 \\ + 1.01 \times 10^2 &+ 0.10 \times 10^3 &= & + 1.01 \times 10^2 \end{aligned}$$

$$\begin{aligned} 0.733 \times 10^5 &= 7.330 \times 10^4 \\ + 0.015 \times 10^7 &= 1.5 \times 10^5 \end{aligned}$$

# Binary examples

## Normalising:

$$\begin{aligned} & 111101e00100 = (-1 + 0.5 + 0.25 + 0.125 + 0.03125) \\ & = 111010e00011 \quad \times 2^4 = -1.5_{10} \\ & = 110100e00010 \\ & = 101000e00001 = (-1 + 0.25) \times 2^1 = -0.75 \times 2 = -1.5 \end{aligned}$$

## 2.1 Normalisation

Act of maximising precision.

left shift mantissa and subtract exponent.  
UNTIL  $MSB = \overline{(MSB_{-1})}$ .

## 2.2 Justification

Act of matching exponents

Arithmetic right shift and increase exponent  
UNTIL exponent matches target exponent

## Justification example

$$\begin{aligned} & 101101100e0000 && \text{to} && e0111 \\ = & 110110110e0001 \\ = & 111011011e0010 \\ = & 111101101e0011 \\ = & 111110110e0100 \\ = & 111111011e0101 \\ = & 111111101e0110 \\ = & 111111110e0111 \end{aligned}$$

$$000101e0010 + 011001e0011$$

$$\begin{aligned} & 0.10100e0000 \\ + & 0.11001e0011 \\ = & 0.06010e0011 \\ + & 0.11001e0011 \\ \hline & 0.11011e0011 \end{aligned}$$

## Overflow example

$$\begin{aligned} & 0.101e000 \\ + & 0.101e000 \end{aligned}$$

$$\begin{aligned} & 0.101e000 \\ + & 0.101e000 \\ \hline \text{overflow } & 1.00e000 \\ & 0.101e001 \end{aligned}$$

## 2.3 Floating point addition and subtraction

1. Normalise.
2. Justify on max exponent.
3. Add/Subtract.
4. Deal with carry/overflow.
5. Normalise answer.

## 2.4 FP Multiplication

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

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

$$\begin{aligned} z &= y_1 \times y_2 && E_1 + E_2 \\ &= (M_1 \times M_2) \times x \end{aligned}$$

1. Normalise
2. Fixed point multiplication of mantissas.
3. Fixed point addition of exponents.
4. Normalise

# Example

$$101 \ e 0110$$

$$111 \ e 0100 = 100 \ e 0010$$

$$M_1 \times M_2$$

$$Q0.2 \quad 011$$

$$Q1.2 \times 0100$$

$$\underline{\quad 000}$$

$$0000$$

$$\sum 01100$$

$$\underline{000000}$$

$$00.1100$$

$$Q1.4$$

(Added additional bit to be sign aware)

$$0.1100 \quad Q0.4$$

Add exponents

$$\cancel{0}0110$$

$$\cancel{0}0010$$

$$\underline{\cancel{0}1000}$$

$$= 01000$$

overflow, adding bit to preserve sign

$$z = 01100 \ e 01000$$