

Floating Point Addition (Binary)

1. Align binary points:

- Align binary point of the number with smaller exponent

2. Add significands:

3. Normalize result:

4. Round and renormalize if necessary:

NB. We assumed that significands can be only 4 bits or digits.

repeated multiplication by 2

$$\begin{array}{r} 0.5 \\ \hline 0.5 \times 2 = 1.0 \\ \hline 0.5 = 0.1_2 \end{array}$$

Problem: Perform binary floating-point addition and convert the resulting values to IEEE-754 single and IEEE-754 double precision format. Finally convert them to hexadecimal values. 0.5 + 0.4375

$$0.25 + 0.125 + 0.0625 = 0.4375$$

Binary point ↓

Weight	0.5	0.25	0.125
	1	0	0

$$0.5_{10} = 0.1_2$$

$$= 1.0 \times 2^{-1}$$

Normalized

Weight	0.5	0.25	0.125	0.0625
	0	1	1	1

$$0.4375_{10} = 0.0111_2$$

$$= 1.11 \times 2^{-2}$$

Normalized

1. Align binary points

$$0.5 = 1.000 \times 2^{-1}$$

$$0.4375 = 0.111 \times 2^{-1}$$

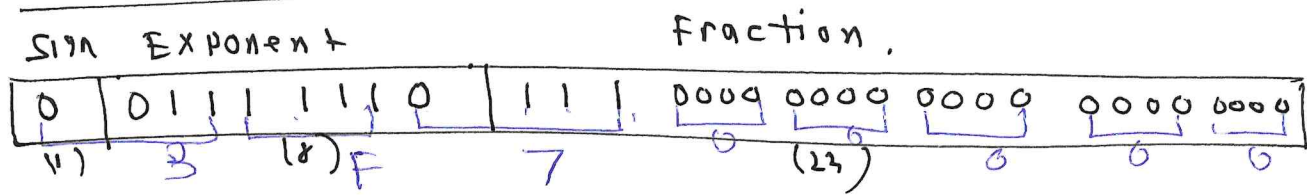
significands

2. Add significands

$$\begin{array}{r} 1.000 \\ 0.111 \\ \hline 1.111 \end{array}$$

3. Result = 1.111×2^{-1}

IEEE-754 single precision Format



$$S = 0$$

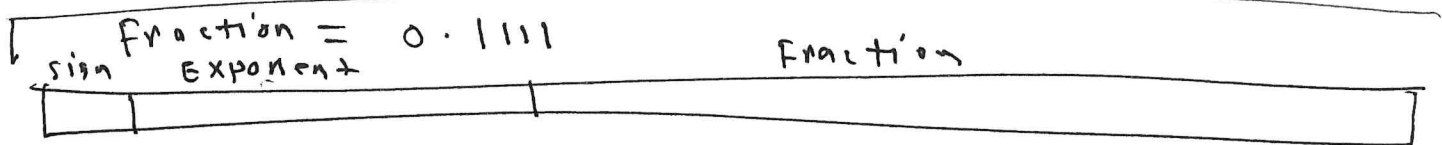
$$\text{exponent} = -1 + 127 = 126_{10} = 01111110_2$$

$$\text{Fraction} = 0.11100 \dots$$

Hex value in single precision = 0x3f700000

IEEE-754 Double precision Format

$$S = 0, \quad \text{exponent} = -1 + 1023 = 1022_{10} = 0111111110_2$$



Floating Point Subtraction (Binary)

1. Align binary points:

- Align binary point of the number with smaller exponent

2. Add significands:

3. Normalize result:

4. Round and renormalize if necessary:

NB. We assumed that significands can be only 4 bits or digits.

Problem: Perform binary floating-point subtraction and convert the resulting values to IEEE-754 single. Finally convert them to hexadecimal values.

$$0.5 - 0.4375$$

$$\begin{array}{r} 0.5 \\ \hline 1.0 \times 2^{-1} \end{array}$$

$$\begin{array}{r} -0.4375 \\ \hline -1.11 \times 2^{-2} \end{array}$$

$$\begin{array}{r} 8 \\ + (-3) \\ \hline 5 \end{array}$$

$$\begin{array}{r} 8 \\ - 3 \\ \hline 5 \end{array}$$

$$\begin{array}{r} 8 \quad 8 \\ 2's \text{ complement of } (+3) = -3 \\ \hline 5 \end{array}$$

$$+8 = 1000$$

$$+3 = 0011$$

$$+3 = 0011$$

$$1's \text{ complement} = 1100$$

$$+1$$

$$-3 = 2's \text{ complement} = 1101$$

$$\rightarrow (+3)$$

$$+8 = 1000$$

$$-3 = 1101$$

$$\begin{array}{r} \boxed{1} 0101 \\ \hline \end{array}$$

5

1. Align Binary points

$$0.5 = 1.000 \times 2^{-1}$$

$$-0.4375 = -\boxed{0.111} \times 2^{-1}$$

2. Add significands

$$2's \text{ complement of } 0.111 = 1.000$$

$$\begin{array}{r} 1.000 \\ + 0.111 \\ \hline 1.001 \end{array}$$

$$\begin{array}{r} 1.000 \\ 1.001 \\ \hline \cancel{1.001} \times 2^{-1} \end{array}$$

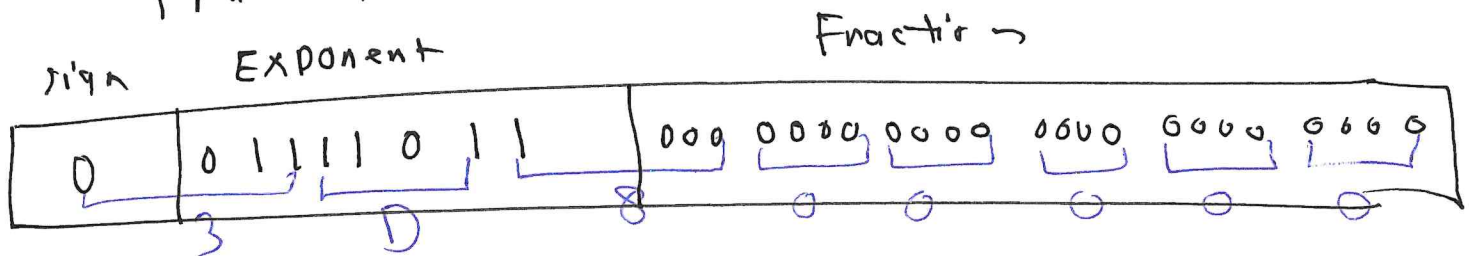
3. Normalized result = 1.000×2^{-4}

IEEE-754 single precision format

$$S = 0$$

$$\text{Exponent} = -4 + 127 = 123_{10} = 0111011_2$$

$$\text{Fraction} = 0.0$$



$$\text{Hex} = 0x3D800000$$

Floating Point Multiplication

1. Add exponents:
2. Multiply significands:
3. Normalize result:
4. Round and renormalize if necessary:
5. Determine the sign of result from signs of operands

NB. We assumed that significands can be only 4 digits of the significands and two digits of the exponents.

Floating Point Multiplication (Decimal)

Consider a 4-digit decimal example: $(1.110 \times 10^{10}) \times (9.2 \times 10^{-5})$

1. Add exponents

$$e_1 = 10 \quad e_2 = -5$$

$$e = e_1 + e_2 = 10 - 5 = 5.$$

2. Multiply significands = 1.110×9.2
= 10.212

$$\text{Result} = 10.212 \times 10^5.$$

3. Normalize result = $1.0212 \times 10^6.$

Solution

$$0.5_{10} = 0.1$$

$$\text{Normalized} = 1.000 \times 2^{-1}$$

1. Add exponent)

$$e_1 = -1$$

$$e_2 = -2$$

$$e = e_1 + e_2 = -3$$

[resulting exponent]

2. Multiply significand)

$$1.000 \times 1.110 = 1.110$$

$$\text{result} = 1.11 \times 2^{-3}$$

3. Normalize result = 1.11×2^{-3}

Problem: Perform binary floating-point multiplication and convert the resulting values to IEEE-754 single precision format. Finally convert them to hexadecimal values.

$$0.5 * -0.4375$$

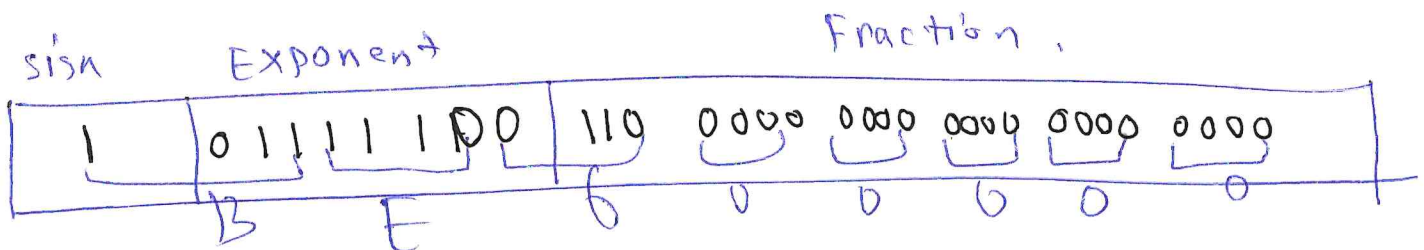
$$4. \text{ Sign} = (+) * (-) = -$$

$$\text{Sign} = 1$$

$$\text{Exponent} = -3 + 127 = 124_{10} = 01111100$$

$$\text{Fraction} = 0.11$$

$$\begin{array}{r} 1.000 \\ \times 1.110 \\ \hline 0000 \\ 1000 \\ 1000 \\ \hline 10000000 \end{array}$$



0xBE600000