

Single precision IEEE floating point (32 bits)

- First bit is the sign: 1 is negative, 0 is non-negative
- Next 8 bits represent the exponent in excess 127 notation (stored value is 127 larger than the exponent)
- Remaining 23 bits represent the fraction (mantissa), excluding the 1 to the left of the radix point
- Positive or negative infinity: exponent is all 1's and the fraction is 0
- Not-a-number (NaN): exponent is all 1's and fraction is anything non-zero

s - sign bit

f - fraction representation

e - exponent representation

$$1.f \times 2^{e-127}$$

negative if s is 1

Decimal

$$\begin{array}{r} 0.13526 \\ \text{place} \quad \frac{1}{10} \quad \frac{1}{100} \quad \frac{1}{1000} \\ 10^{-1} \quad 10^{-2} \quad 10^{-3} \quad \dots \end{array}$$

Binary

$$\begin{array}{r} 1.01101001 \\ \text{---} \quad \uparrow \uparrow \uparrow \\ \frac{1}{2} \quad \frac{1}{4} \quad \frac{1}{8} \\ 2^{-1} \quad 2^{-2} \quad 2^{-3} \end{array}$$

- 62.125 to IEEE single-precision

$$62_{10} = 111110_2$$

$$0.125_{10} = 0.001_2$$

11110.00 1

Shift radix point 5 places $\rightarrow 1.111000_2 \times 2^{5_{10}}$

Exponent is 5, add 127 → 132

$$132_{10} = 10000100_2$$

Fraction, drop the 1 to the left of radix Point, pad with zeros on the right

$$0 \times \frac{110000100111}{\begin{array}{r} 6 \\ - 3 \\ \hline 3 \end{array}} \quad \frac{1000}{\begin{array}{r} 8 \\ - 8 \\ \hline 0 \end{array}} \quad \frac{1000000000000000}{\begin{array}{r} 0 \\ - 0 \\ \hline 0 \end{array}}$$

0.1

$$2^{-4} + 2^{-5} + 2^{-8} + 2^{-9} + 2^{-12} + 2^{-13} + \dots$$

(0.0625)

0.0001100110011001100 ...

Exponent is -4, add 127 \rightarrow 123

1111011

Fraction

| 0011001100110011001100