

Burnsley Lecture

Float | 1-bit sign
8-bit exponent
23-bit mantissa

$$\text{sign}(1 + \text{mantissa}) 2^{\text{exponent} - 127}$$

$$\text{Range: } -3.4 \times 10^{38} \text{ to } 3.4 \times 10^{38}$$

$$\text{Maximum: } \pm 3.4 \times 10^{38}$$

$$\text{Minimum: } \pm 1.7 \times 10^{-38}$$

Range

max 0
min 1

11111110

111...1

23, 15

$$\text{sign} \begin{cases} 0 \rightarrow + \\ 1 \rightarrow - \end{cases}$$

$$\text{Exponent} \mid 2^7 + 2^6 + \dots + 2 + 0 = 254$$

$$\text{Mantissa} \mid 2^{-1} + 2^{-2} + \dots + 2^{-23} = \sum_{i=1}^{23} \left(\frac{1}{2}\right)^i = 1 - \frac{1}{2^{23}}$$

$$\pm \left(1 + 1 - \frac{1}{2^{23}}\right) \cdot 2^{254 - 127} = \pm 3.4 \times 10^{38}$$

Minimum/smallest

complement of max range

0 0000001 000...0 23, 0s

$$\text{sign} \begin{cases} 0 \rightarrow + \\ 1 \rightarrow - \end{cases}$$

$$\text{Exponent} \mid 0 + \dots + 2^0 = 1$$

$$\text{Mantissa} \mid 0 + \dots + 0 = 0$$

$$\pm (1 + 0) 2^{1 - 127} = \pm 1.754944 \times 10^{-38}$$

Double | 1-bit sign
11-bit exponent
52-bit mantissa

$$\text{Range: } -1.7 \times 10^{308} \text{ to } 1.7 \times 10^{308}$$

$$\text{Greatest: } \pm 1.7 \times 10^{308}$$

$$\text{Smallest: } \pm 2.2 \times 10^{-308}$$

Range

max 0
min 1

1111111110

1...1 52, 15

$$\text{sign} \begin{cases} 0 \rightarrow + \\ 1 \rightarrow - \end{cases}$$

$$\text{Exponent} \mid 2^{10} + 2^9 + \dots + 2 + 0 = 2046$$

$$\text{Mantissa} \mid 2^{-1} + \dots + 2^{-52} = 1 - \frac{1}{2^{52}}$$

$$\pm \left(1 + 1 - \frac{1}{2^{52}}\right) \cdot 2^{2046 - 1023} = \pm 1.7 \times 10^{308}$$

Smallest

0 0000000000 0...0 52, 0s

$$\text{sign} \begin{cases} 0 \rightarrow + \\ 1 \rightarrow - \end{cases}$$

$$\text{Exponent} \mid 0 + \dots + 2^0 = 1$$

$$\text{Mantissa} \mid 0 + \dots + 0 = 0$$

$$\pm (1 + 0) \cdot 2^{1 - 1023} = \pm 2.22507 \times 10^{-308}$$