

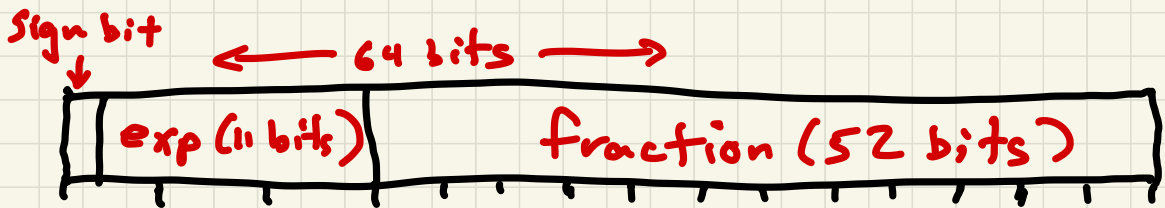
# Properties of floating-point

CS 111

Nov. 10, 2020



# IEEE Standard Float64



REPRESENTS THE NUMBER:

$$\text{If } 1 \leq \text{exp} \leq 2046, \\ \pm (1.\text{frac}) \times 2^{\text{exp} - 1023}$$

If  $\text{exp} = 0$  and  $\text{frac} = 0$ ,

0

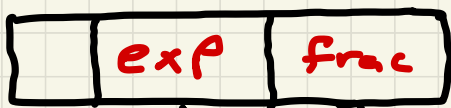
If  $\text{exp} = 2047$  ( $\equiv 7fff$ ) and  $\text{frac} = 0$ ,

Inf

If  $\text{exp} = 2047$  and  $\text{frac} \neq 0$ ,

NaN (not-a-number)

# Toy 5-bit floating-point



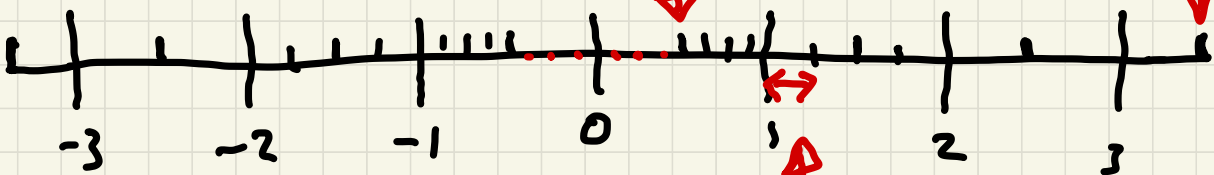
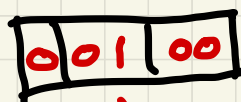
sign 9 2 bits 2 bits

$$1 \leq \text{exp} \leq 3: \pm(1.\text{frac}) \times 2^{\text{exp}-2}$$

max:  $1.11 \times 2^1 = 11.1 = 3\frac{1}{2}$



smallest normalized:  $1.00 \times 2^{-1} = \frac{1}{2}$



machine epsilon =

distance from 1  
to the next larger

Floating-point number =  $\frac{1}{4}$

# Catastrophic Cancellation Example

$$\begin{aligned}(x-1)^7 &= x^7 + 7x^6 \\ &\quad + 21x^5 - 35x^4 \\ &\quad + 35x^3 - 21x^2 \\ &\quad + 7x - 1.\end{aligned}$$