

# Floating point

Recap:

$$217.3 = \underbrace{(1.1011)}_{(1 \leq \cdot \leq 2)} \cdot 2^{\underbrace{7}}_{\lceil \cdot \rceil}$$

↓  
(1 - - - - .)

## Unrepresentable numbers?

- Can you think of a somewhat central number that we cannot represent as

$$x = (1.\text{-----})_2 \cdot 2^{-p}?$$

0  $\rightarrow$  Positive 0: significand: all zeros + special exp. -1024

Exp. range -1023... 1024

$\uparrow$   
-1024

$\uparrow$   
special meaning: "turn off 'the leading 1'"

## Demo: Picking apart a floating point number

$$\text{Exponent} = \overset{\sim}{-1023} \rightarrow \text{stored} \uparrow$$

↑  
implicit  
(not stored)

$$(3)_{10} = (11)_2 = (1.1) \cdot 2^1 = (1.1) \cdot 2^{\sim{1023} + 1024}$$

---

$$\left. \begin{aligned} 2^{-1022} &= (1) \cdot 2^{-1022} \\ 2^{-1023} &= (0.\overbrace{100000}^{\text{mantissa}}) \cdot 2^{-1023} \\ 0 &= (0.\underline{0}\underline{0}\underline{0}\underline{0}\underline{0}\underline{0}) \cdot 2^{-1023} \end{aligned} \right\} \begin{array}{l} \text{normal} \\ \text{Subnormal} \end{array}$$

## Subnormal Numbers ✓ normal

- What is the smallest representable number in an FP system with 4 stored bits in the significand and an exponent range of  $[-7, 7]$ ?

$$\begin{array}{l}
 \begin{array}{c|cccc}
 (1 & \underline{0} & \underline{0} & \underline{0} & \underline{0})_2 \cdot 2^{\overbrace{-7}} & = 2^{-7} \\
 (0 & \underline{1} & \underline{0} & \underline{0} & \underline{0})_2 \cdot 2^{\underbrace{-8}} & = 2^{-8} \\
 \text{not stored} & & \text{stored} & & & \\
 \hline
 ( & \downarrow & 0 & 0 & 0 & 1 )_2 \cdot 2^{-8} & = 2^{-11}
 \end{array}
 \end{array}$$

"FP assist"  $\rightarrow$  working w/ subnormal numbers  
 $\rightarrow$  super slow

**Demo:** Density of Floating Point Numbers

**Demo:** Floating Point vs. Program Logic

"Underflow"

"Subnormals" → "gradual underflow"

## Floating Point and Rounding Error

What is the relative error produced by working with floating point numbers?

- What is smallest floating point number  $> 1$ ? Assume 4 stored bits in the significand.

$$(1.\underbrace{\dots}_{2^{-4}}1)_2 \approx 1 + \underbrace{2^{-4}}_{\text{machine eps}}$$

- What's the smallest FP number  $> 1024$  in that same system?

U LP : unit in the last place

- Can we give that number a name?

- What does this say about the relative error incurred in floating point calculations?

Rel. error introduced in every fp. op. is  
 $\sim$  machine eps

- What's that same number for double-precision floating point? (52 bits in the significand)

$$2^{-52}$$

## Implementing Arithmetic

- How is floating point addition implemented?

Consider adding  $a = (1.101)_2 \cdot 2^1$  and  $b = (1.001)_2 \cdot 2^{-1}$  in a system with three bits in the significand.

*shift onto same exponent*

$$\begin{array}{r} (1.101)_2 \cdot 2^1 \\ + (1.001)_2 \cdot 2^{-1} \\ \hline \end{array}$$

$$\begin{array}{r} \sim (1.101) \cdot 2^1 \\ \sim (0.01001) \cdot 2^1 \\ \hline (1.111) \cdot 2^1 \end{array}$$

$$(1.111)$$

**Demo:** Floating point and the harmonic series



## Problems with FP Addition

- What happens if you subtract two numbers of very similar magnitude?  
As an example, consider  $a = (1.1011)_2 \cdot 2^0$  and  $b = (1.1010)_2 \cdot 2^0$ .

**Demo:** Catastrophic Cancellation

**In-class activity:** Floating Point 2