# Floating Point Arithmetic

Want: Something like the real numbers... in a computer

Have: Integers, made of bits

$$23 = 0 \cdot 2^5 + 1 \cdot 2^4 + 0 \cdot 2^3 + 1 \cdot 2^2 + 1 \cdot 2^1 + 1 \cdot 2^0$$

$$(10111)_2 = (23)_{10}$$

How should we even represent fractions?

Idea: Keep going down past exponent zero

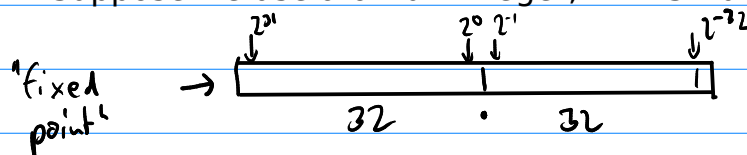$$23.625 = 0 \cdot 2^5 + 1 \cdot 2^4 + 0 \cdot 2^3 + 1 \cdot 2^2 + 1 \cdot 2^1 + 1 \cdot 2^0$$
$$+ 1 \cdot 2^{-1} + 0 \cdot 2^{-2} + 1 \cdot 2^{-3}$$
$$0.5 \qquad\qquad 0.25 \qquad\qquad 0.125$$

So:  Could store
  - a fixed number of bits with exponents >= zero
  - a fixed number of bits with exponents < zero

Suppose we use a 64-bit integer, with 32 bits >= 1 and 32 bits < 1.

"fixed point"  →

$$2^{31} \qquad 2^0 \; 2^{-1} \qquad 2^{-32}$$

32  •  32

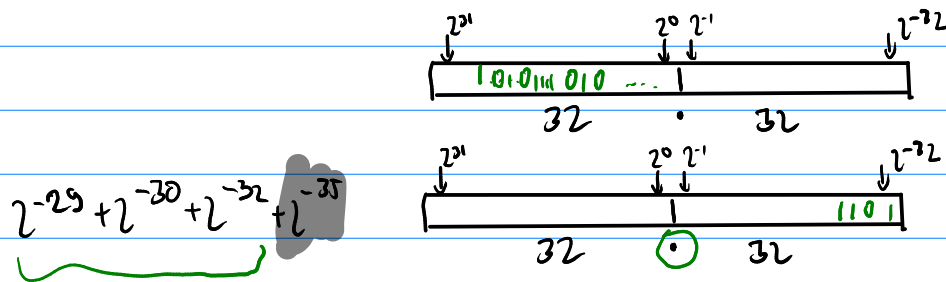What is the smallest number we can represent?

$$2^{-32} \approx 10^{-10}$$

What is the biggest number we can represent?

$$2^{31} + 2^{30} + 2^{29} + \cdots \approx 10^{9}$$

What's our range then?

$$10^{-10} \quad \cdots \quad 10^{-9}$$

This is called fixed-point arithmetic, and it's pretty bad.



$2^{-29} + 2^{-30} + 2^{-32} + 2^{-35}$

Should be able to do better.

Set a few bits aside to store the largest exponent. How?

$$2^{29} + 2^{27} \qquad = \left(2^0 + 2^{-2}\right) \cdot 2^{29} \qquad = (1.01)_2 \cdot 2^{29}$$

$$2^{-29} + 2^{-30} + 2^{-32} + 2^{-35} = \left(2^0 + 2^{-1} + 2^{-3} + 2^{-6}\right) \cdot 2^{-29} = \left(1.101001\right)_2 \cdot 2^{-29}$$

"Floating point"

"significand"          "exponent"