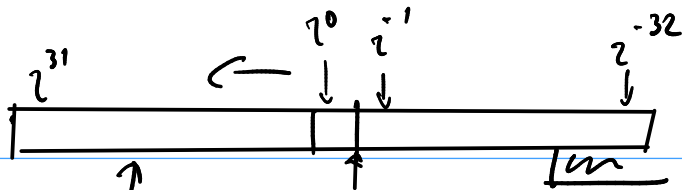


Overview

- Floating point

→ Example 2
→ IETF



Fixed point

Flaw: Uneven amount of accuracy
between small and large
numbers

$$(Number) = 1 \in (\leq 2) \cdot 2^{\text{exponent}}$$

$\underbrace{1.111111}_{\text{significant}}$

$$2 \cdot 2^{\text{exponent}}$$

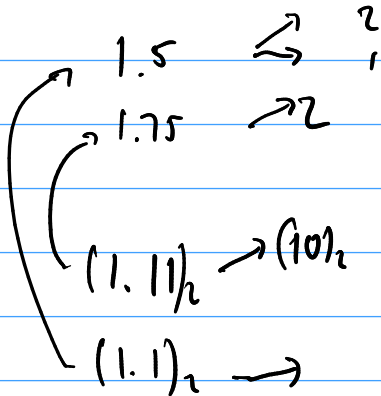
$$= 1 \cdot 2^{\text{exponent} + 1}$$

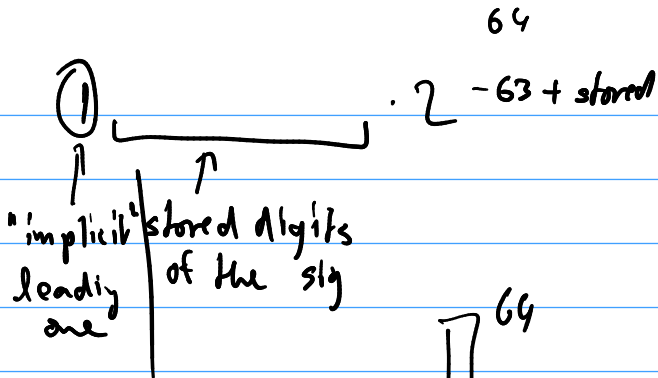
$$123.125 = (1111011001)_2$$

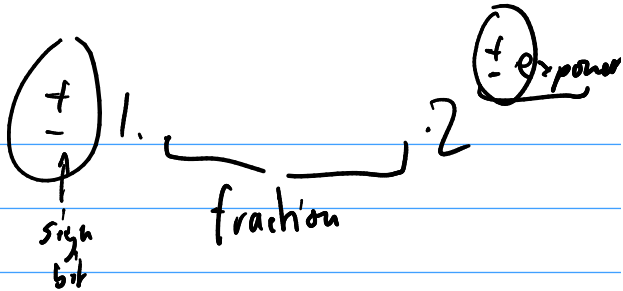
$2^6 = 64$

$\begin{matrix} \uparrow & \uparrow & \uparrow & \uparrow & \uparrow \\ 2^6 & 2^5 & 2^4 & 2^3 & 2^2 \end{matrix}$

$$= \left(\underbrace{1111011}_{\text{int}} \underbrace{001}_{\text{rounded}} \right) \cdot 2^6$$





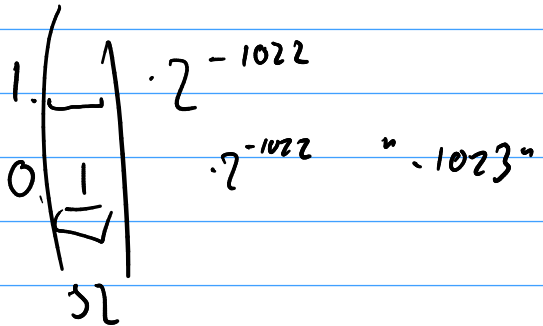


$$2^{-1022} = (1. \text{---})_2 \cdot 2^{-1022}$$

$$2^{-1023} = (0. \underline{10000})_2 \cdot 2^{-1022}$$

↑ subnormal numbers

1.000000001



Floating Point numbers

Convert $13 = (1101)_2$ into floating point representation.

What pieces do you need to store an FP number?

In-class activity: Floating Point

Unrepresentable numbers?

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

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

Demo: Picking apart a floating point number

Subnormal Numbers

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

Subnormal Numbers (II)

Why would you want to know about subnormals? Because computing with them is often slow, because it is implemented using 'FP assist', i.e. not in actual hardware. Many C compilers support options to 'flush subnormals to zero'.

Demo: Density of Floating Point Numbers

Demo: Floating Point vs. Program Logic

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.

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

Can we give that number a name?