

# IEEE-754

---

Junseo Shin

April 11, 2024

- Introduction
- Floating Point Arithmetic
- Simulations and Examples

# Introduction

---

# Should We Trust Computers?

- The error in floating point as “noise” in the data.
- The computer automatically “preprocesses” the data in a way that is not always what you want.

## How does a computer operate on numbers?

178956970 – 178957034

# How does a computer operate on numbers?

178956970 – 178957034



# How does a computer operate on numbers?

178956970 – 178957034



= -64

## Adding/Subtracting two 32 bit numbers

$$\begin{array}{r} (0)000\ 1010\ 1010\ 1010\ 1010\ 1010\ 1010\ 1010 \\ - (0)000\ 1010\ 1010\ 1010\ 1010\ 1010\ 1110\ 1010 \\ \hline (1)000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0100\ 0000 \end{array}$$



## Adding/Subtracting two 32 bit numbers

$$\begin{array}{r} (0)000\ 1010\ 1010\ 1010\ 1010\ 1010\ 1010\ 1010 \\ - (0)000\ 1010\ 1010\ 1010\ 1010\ 1010\ 1110\ 1010 \\ \hline (1)000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0100\ 0000 \end{array}$$

$$= -2^6 = -64$$

# Multiplying two binary numbers

$$170 \times 170 = 28900$$

# Multiplying two binary numbers

$$170 \times 170 = 28900$$

									1	0	1	0	1	0	1	0	
									×	1	0	1	0	1	0	1	0
										0	0	0	0	0	0	0	0
									1	0	1	0	1	0	1	0	
							0		0	0	0	0	0	0	0		
						1	0		1	0	1	0	1	0			
					0	0	0		0	0	0	0	0				
				1	0	1	0		1	0	1	0					
			0	0	0	0	0		0	0	0	0					
				1	0	1	0		1	0	1	0					
			0	0	0	0	0		0	0	0						
		1	0	1	0	1	0		1	0	1	0					
		0	0	0	0	0	0		0	0	0						
	1	0	1	0	1	0	1		0								
1	1	1	0	0	0	0	0		1	1	1	0	0	1	0	0	

# Dividing Two Binary numbers

$$\begin{array}{r} \phantom{10101010} 10101010 \\ 10101010 \overline{) 111000011100100} \\ \underline{10101010} \phantom{00} \\ 0011011111 \\ \phantom{00} \underline{10101010} \phantom{00} \\ 0011010100 \\ \phantom{00} \phantom{00} \underline{10101010} \phantom{00} \\ 0010101010 \\ \phantom{00} \phantom{00} \phantom{00} \underline{10101010} \phantom{00} \\ 00000000 \end{array}$$

# Problems with Binary Integers

No way to represent rational numbers

# Problems with Binary Integers

No way to represent rational numbers

- Division is not accurate

$$3/2 = 1.5 \rightarrow 101/10 = 1.\overset{\cdot}{1} = 1$$

# Problems with Binary Integers

No way to represent rational numbers

- Division is not accurate

$$3/2 = 1.5 \rightarrow 101/10 = 1.\overset{.}{1} = 1$$

- Real-Valued Functions (e.g.  $\sin(x)$ ,  $e^x$ ,  $\sqrt{x}$ )

# IEEE-754: Floating Point Arithmetic

---



# From Decimal to Floating Point

Scientific Notation:  $n \times 10^m$

- $1234 = 1.234 \times 10^3$
- $0.01234 = 1.234 \times 10^{-2}$

# From Decimal to Floating Point

Scientific Notation:  $n \times 10^m$

- $1234 = 1.234 \times 10^3$
- $0.01234 = 1.234 \times 10^{-2}$

Floating Point: (sign, significand, exponent)

$$(-1)^s \times m \times 2^e$$

# Converting Decimal to Floating Point

Decimal  $\rightarrow$  Binary  $\rightarrow$  Floating Point

$$7.45 \rightarrow 7 + 0.45$$

$$\rightarrow 111.01110011001100110011001100110011 \dots$$

1 bit for sign, 8 bits for exponent, 23 bits for significand

# Converting Decimal to Floating Point

Decimal  $\rightarrow$  Binary  $\rightarrow$  Floating Point

$$7.45 \rightarrow 7 + 0.45$$

$$\rightarrow 111.01110011001100110011001100110011 \dots$$

$$\begin{array}{c} \text{sign} \\ \boxed{+} \\ \text{significand} \end{array} \underbrace{1.1101 \dots}_{\text{significand}} \times 2^{\underbrace{2}_{\text{exponent}}}$$

1 bit for sign, 8 bits for exponent, 23 bits for significand

$$\boxed{0 \mid 1000 \ 0001 \mid 1101 \ 1100 \ 1100 \ 1100 \ 1100 \ 110} 0110 \dots$$

$$(-1)^0 \times 2^{129-127} \times 1.1101 \ 1100 \ 1100 \ 1100 \ 1100 \ 110$$

# Error in Floating Point

Machine Epsilon

$$\epsilon = 2^{-(p-1)}$$

- Single Precision:  $p = 24$ :  $\epsilon = 2^{-23} \approx 1.19 \times 10^{-7}$

$$1 + \epsilon = 1.0000\ 0000\ 0000\ 0000\ 0000\ 001$$

- Double Precision:  $p = 53$ :  $\epsilon = 2^{-52} \approx 2.22 \times 10^{-16}$

# Simulations and Examples

---

```
python3
```

```
>>> 0.1 + 0.2
```

```
0.30000000000000004
```

```
>>> 0.3 / 0.10
```

```
2.9999999999999996
```

## GNU Multiple Precision Arithmetic Library (GNU MPFR)

- Used in Mathematica
- Follows IEEE-754 standard but with arbitrary precision
- Does this solve the problem?



# Stirling's Approximation vs Factorial

$$n! \approx \sqrt{2\pi n} \left(\frac{n}{e}\right)^n$$

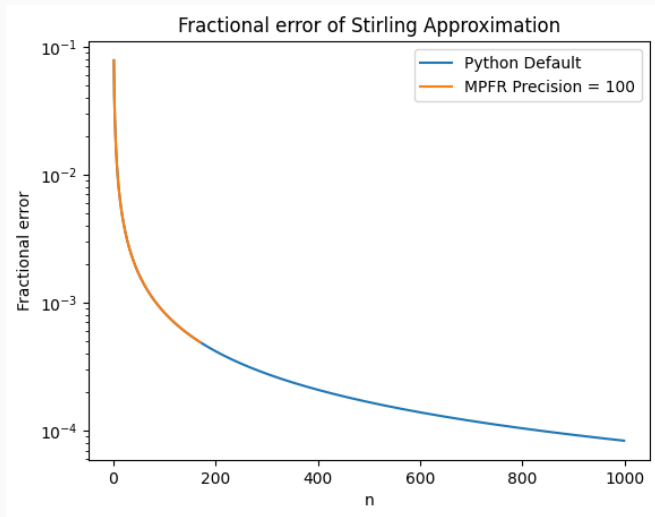
# Stirling's Approximation vs Factorial

$$n! \approx \sqrt{2\pi n} \left(\frac{n}{e}\right)^n$$

Python Default (Float64,  $p = 53$ ) vs MPFR ( $p = 100$ )

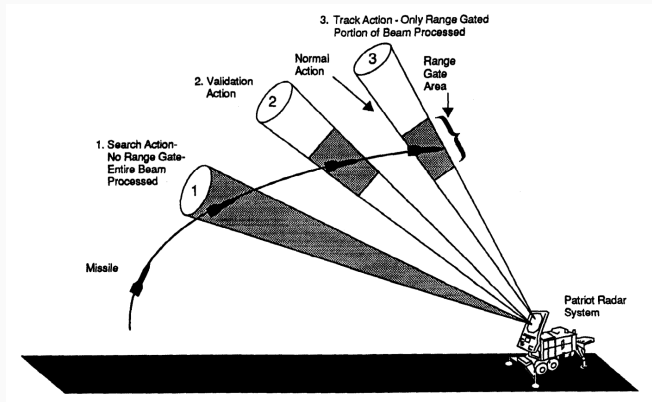
- $n = 171 \rightarrow$  Overflow in Float64
- At  $n = 100$ , Fractional Error of 0.000833 for both!
- Fractional Error of  $10^{-6}$  at  $n = 83334$

# Float64 vs MPFR



# Patriot Defense Missile Failure

- February 25, 1991: Patriot missile defense system failed to intercept SCUD missile
- 28 soldiers died



## Converting Integer to 24 bit Floating Point

- Integer to 24 bit floating point: 1 sign bit, 8 exponent bits, 15 + 1 significand bits
- 100 hours or 360,000 seconds  $\rightarrow$  0 10010001 010111111001000

$$(-1)^0 \times 2^{145-127} \times 1.01111111001000 = 360,000$$

- 360,010 seconds  $\rightarrow$  360,000 seconds

# Simulating Precision Loss

- Clock Speed: 10 ticks per second

$$0.1 \rightarrow 0.0999755859375$$

Calculated time: 359912.109375

- Machine Epsilon  $\epsilon = 2^{-20}$ ; Propogated Error =  $\epsilon * \text{seconds}$

# Government Report

Hours	Seconds	Calculated Time (Seconds)	Inaccuracy (Seconds)	Approximate Shift In Range Gate (Meters)
0	0	0	0	0
1	3600	3599.9966	.0034	7
8	28800	28799.9725	.0275	55
20 <sup>a</sup>	72000	71999.9313	.0687	137
48	172800	172799.8352	.1648	330
72	259200	259199.7528	.2472	494
100 <sup>b</sup>	360000	359999.6667	.3433	687

<sup>a</sup>Continuous operation exceeding about 20 hours—target outside range gate

<sup>b</sup>Alpha Battery ran continuously for about 100 hours

[3]

- You don't need to open a black box to understand where it fails
- Think about the limitations of a system. Just because it works in one instance doesn't mean it will work in another.
- Measure noise before you measure signal





IEEE standard for floating-point arithmetic.

**IEEE Std 754-2019 (Revision of IEEE 754-2008), pages 1–84, 2019.**



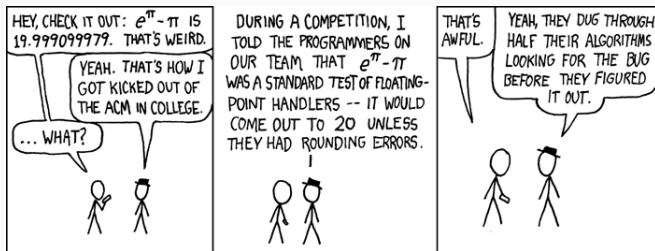
David Goldberg.

**What every computer scientist should know about floating-point arithmetic.**

*ACM Computing Surveys*, 23:5–48, Mar 1991.

[2] [1]

- <https://floating-point-gui.de/>
- Patriot Missile Defense Software Problem Led to System Failure at Dhahran, Saudi Arabia [3]
- Floating Point Converter



## 24 Bit Floating Point

$$e^{\pi} - \pi = 19.9970703125$$

<https://xkcd.com/217/>