# COMP201
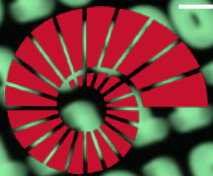
# Computer Systems & Programming

Lecture #05 –Floating Point

KOÇ UNIVERSITY

Aykut Erdem // Koç University // Fall 2020

# Recap: Bitwise Operators

- You're already familiar with many operators in C:
  - **Arithmetic operators:** +, -, *, /, %
  - **Comparison operators:** ==, !=, <, >, <=, >=
  - **Logical Operators:** &&, ||, !


- **Bitwise operators:**
  - **Logical operators:**  &, |, ~, ^,
  - **Bit shift operators:** <<, >>

# COMP201 Topic 2: How can a computer represent real numbers in addition to integer numbers?

# Learning Goals

Understand the design and compromises of the floating point representation, including:

- Fixed point vs. floating point

- How a floating point number is represented in binary

- Issues with floating point imprecision

- Other potential pitfalls using floating point numbers in programs

# Plan For Today

- Representing real numbers
- Fixed Point
- Floating Point

**Disclaimer:** Slides for this lecture were borrowed from

—Nick Troccoli's Stanford CS107 class

# Lecture Plan

- Representing real numbers
- Fixed Point
- Floating Point

# Real Numbers

- We previously discussed representing integer numbers using two's complement.

- However, this system does not represent real numbers such as 3/5 or 0.25.

- How can we design a representation for real numbers?

# Real Numbers

**Problem**: unlike with the integer number line, where there are a finite number of values between two numbers, there are an *infinite* number of real number values between two numbers!

**Integers between 0 and 2:** 1

**Real Numbers Between 0 and 2:** 0.1, 0.01, 0.001, 0.0001, 0.00001,…

We need a fixed-width representation for real numbers. Therefore, by definition, *we will not be able to represent all numbers*.

# Real Numbers

**Problem**: every number base has un-representable real numbers.

**Base 10:** $1/6_{10} = 0.16666666\ldots_{10}$

**Base 2:** $1/10_{10} = 0.00011001100110011001\ldots_2$

Therefore, by representing in base 2, *we will not be able to represent all numbers*, even those we can exactly represent in base 10.

# Fixed Point

- **Idea:** Like in base 10, let's add binary decimal places to our existing number representation.

$$5\ 9\ 3\ 4\ .\ 2\ 1\ 6$$

$10^3 \quad 10^2 \quad 10^1 \quad 10^0 \qquad 10^{-1} \quad 10^{-2} \quad 10^{-3}$

$$1\ 0\ 1\ 1\ .\ 0\ 1\ 1$$

$2^3 \quad 2^2 \quad 2^1 \quad 2^0 \qquad 2^{-1} \quad 2^{-2} \quad 2^{-3}$

# Lecture Plan

- Representing real numbers
- Fixed Point
- Floating Point

# Fixed Point

- **Idea:** Like in base 10, let's add binary decimal places to our existing number representation.

$$1\ 0\ 1\ 1\ .\ 0\ 1\ 1$$

8s    4s    2s    1s        1/2s  1/4s  1/8s

- **Pros:** arithmetic is easy!  And we know exactly how much precision we have.

# Fixed Point

- **Problem:** we have to fix where the decimal point is in our representation. What should we pick?  This also fixes us to 1 place per bit.

. 0 1 1 0 0 1 1

<span style="color:red">1/2s    1/4s    1/8s       …</span>

1 0 1 1 0 . 1 1

<span style="color:red">16s     8s     4s     2s     1s              1/2s   1/4s</span>

# Fixed Point

- **Problem**: we have to fix where the decimal point is in our representation. What should we pick?  This also fixes us to 1 place per bit.

Base 10

Base 2

$$5.07E30 = 10\underbrace{\phantom{.....}}_{\text{100 zeros}}0.1$$

$$9.86E\text{-}32 = 0.0\underbrace{\phantom{.....}}_{\text{100 zeros}}01$$

To be able to store both these numbers using the same fixed point representation, the bitwidth of the type would need to be at least 207 bits wide!

# Let's Get Real

What would be nice to have in a real number representation?

• Represent widest range of numbers possible

• Flexible "floating" decimal point

• Represent scientific notation numbers, e.g. $1.2 \times 10^6$

• Still be able to compare quickly

• Have more predictable over/under-flow behavior

# Lecture Plan

- Representing real numbers
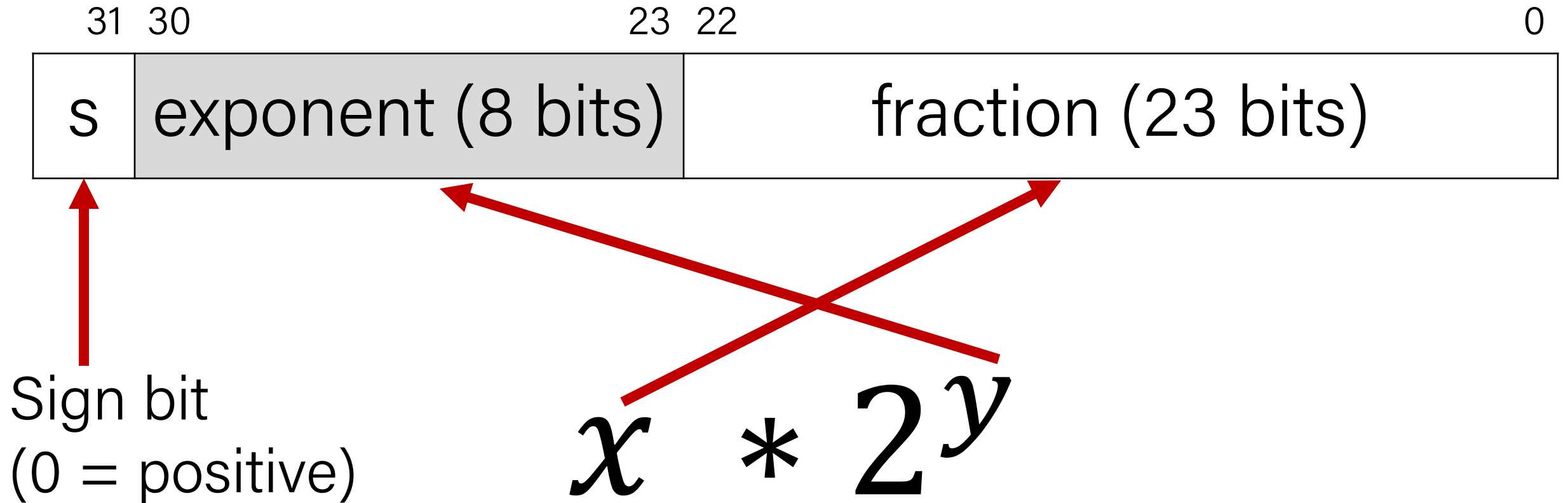- Fixed Point
- Floating Point

# IEEE Floating Point

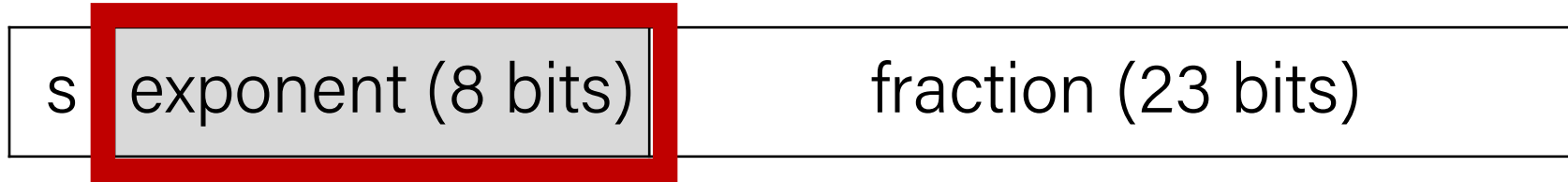Let's aim to represent numbers of the following scientific-notation-like format:

$$x * 2^y$$

With this format, 32-bit floats represent numbers in the range ~1.2 x$10^{-38}$ to ~3.4 x$10^{38}$!  Is every number between those representable?  **No.**

# IEEE Single Precision Floating Point

# Exponent

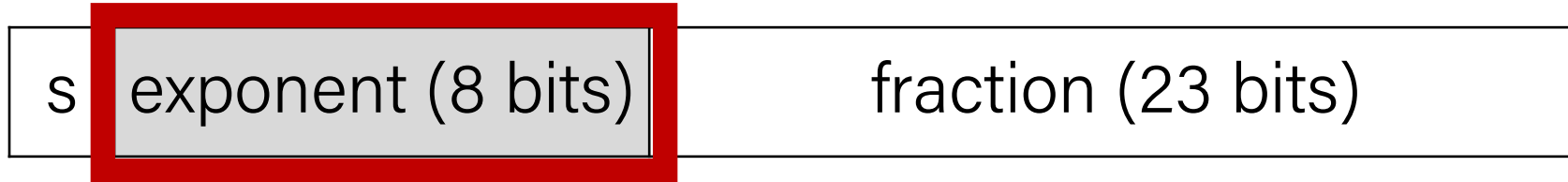| s | exponent (8 bits) | fraction (23 bits) |

| Exponent (Binary) | Exponent (Base 10) |
|:---:|:---:|
| 11111111 | ? |
| 11111110 | ? |
| 11111101 | ? |
| 11111100 | ? |
| … | ? |
| 00000011 | ? |
| 00000010 | ? |
| 00000001 | ? |
| 00000000 | ? |

# Exponent

| s | exponent (8 bits) | fraction (23 bits) |

| Exponent (Binary) | Exponent (Base 10) |
|---|---|
| 11111111 | RESERVED |
| 11111110 | ? |
| 11111101 | ? |
| 11111100 | ? |
| ... | ? |
| 00000011 | ? |
| 00000010 | ? |
| 00000001 | ? |
| 00000000 | RESERVED |

# Exponent

| s | exponent (8 bits) | fraction (23 bits) |

| Exponent (Binary) | Exponent (Base 10) |
|:---:|:---:|
| 11111111 | RESERVED |
| 11111110 | 127 |
| 11111101 | 126 |
| 11111100 | 125 |
| … | … |
| 00000011 | -124 |
| 00000010 | -125 |
| 00000001 | -126 |
| 00000000 | RESERVED |

# Exponent

| s | exponent (8 bits) | fraction (23 bits) |
|---|---|---|

- The exponent is **not** represented in two's complement.

- Instead, exponents are sequentially represented starting from 000...1 (most negative) to 111...10 (most positive).  This makes bit-level comparison fast.

- **Actual value = binary value – 127 ("bias")**

| | |
|---|---|
| 11111110 | 254 – 127 = 127 |
| 11111101 | 253 – 127 = 126 |
| … | … |
| 00000010 | 2 – 127 = -125 |
| 00000001 | 1 – 127 = -126 |

# Fraction

| s | exponent (8 bits) | fraction (23 bits) |
|---|---|---|

$$x * 2^y$$

- We could just encode whatever x is in the fraction field.  But there's a trick we can use to make the most out of the bits we have.

# An Interesting Observation

**In Base 10:**

$42.4 \times 10^5 = 4.24 \times 10^6$

$324.5 \times 10^5 = 3.245 \times 10^7$

$0.624 \times 10^5 = 6.24 \times 10^4$

We tend to adjust the exponent until we get down to one place to the left of the decimal point.

**In Base 2:**

$10.1 \times 2^5 = 1.01 \times 2^6$

$1011.1 \times 2^5 = 1.0111 \times 2^8$

$0.110 \times 2^5 = 1.10 \times 2^4$

**Observation:** in base 2, this means there is always a 1 to the left of the decimal point!

# Fraction

| s | exponent (8 bits) | fraction (23 bits) |
|---|---|---|

$$x * 2^y$$

- We can adjust this value to fit the format described previously. Then, x will always be in the format **1.XXXXXXXX...**

- Therefore, in the fraction portion, we can encode just what is *to the right* of the decimal point! This means we get one more digit for precision.

**Value encoded = 1._[FRACTION BINARY DIGITS]_**

# Practice

| Sign | Exponent | | | | | Fraction | | | |
|------|----|------|---|---|---|----|---|---|------|
| 0 | 0 | … | 0 | 0 | 0 | 1 | 0 | 1 | 0 | … |

Is this number:

**A) Greater than 0?**

**B) Less than 0?**

Is this number:

**A) Less than -1?**

**B) Between -1 and 1?**

**C) Greater than 1?**

# Skipping Numbers

- We said that it's not possible to represent *all* real numbers using a fixed-width representation. What does this look like?

Float Converter

- https://www.h-schmidt.net/FloatConverter/IEEE754.html

Floats and Graphics

- https://www.shadertoy.com/view/4tVyDK

# Let's Get Real

What would be nice to have in a real number representation?

- Represent widest range of numbers possible

- Flexible "floating" decimal point

- Represent scientific notation numbers, e.g. $1.2 \times 10^6$

- Still be able to compare quickly

- Have more predictable over/under-flow behavior

# Representing Zero

The float representation of zero is all zeros (with any value for the sign bit)

| Sign | Exponent | Fraction |
|------|----------|----------|
| any | All zeros | All zeros |

• This means there are two representations for zero! ☹

# Representing Small Numbers

If the exponent is all zeros, we switch into "denormalized" mode.

| Sign | Exponent | Fraction |
|---|---|---|
| any | All zeros | Any |

- We now treat the exponent as -126, and the fraction as *without* the leading 1.
- This allows us to represent the smallest numbers as precisely as possible.

# Representing Exceptional Values

If the exponent is all ones, and the fraction is all zeros, we have +- infinity.

| Sign | Exponent | Fraction |
|------|----------|----------|
| any  | All ones | All zeros |

- The sign bit indicates whether it is positive or negative infinity.
- Floats have built-in handling of over/underflow!
  - Infinity + anything = infinity
  - Negative infinity + negative anything = negative infinity
  - Etc.

# Representing Exceptional Values

If the exponent is all ones, and the fraction is nonzero, we have
**Not a Number (NaN)**

| Sign | Exponent | | | | | | Fraction |
|------|---|---|---|---|---|---|----------|
| any | 1 | … | … | … | … | 1 | Any nonzero |

- NaN results from computations that produce an invalid mathematical result.
  - Sqrt(negative)
  - Infinity / infinity
  - Infinity + -infinity
  - Etc.

# Number Ranges

- 32-bit integer (type **int**):
  - › -2,147,483,648 to 2147483647
  - › Every integer in that range can be represented

- 64-bit integer (type **long**):
  - › −9,223,372,036,854,775,808 to 9,223,372,036,854,775,807

- 32-bit floating point (type **float**):
  - – ~1.2 x10$^{-38}$ to ~3.4 x10$^{38}$
  - – Not all numbers in the range can be represented (not even all integers in the range can be represented!)
  - – Gaps can get quite large! (larger the exponent, larger the gap between successive fraction values)

- 64-bit floating point (type **double**):
  - – ~2.2 x10$^{-308}$ to ~1.8 x10$^{308}$
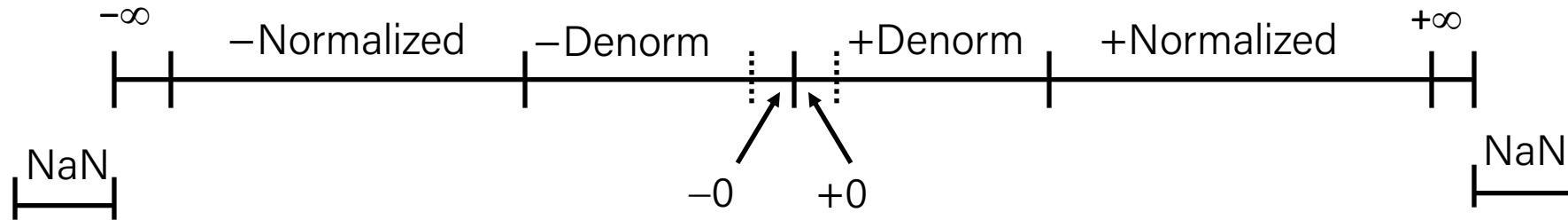
# Precision options

- Single precision: 32 bits

| s | exp | frac |
|---|-----|------|
| 1 | 8-bits | 23-bits |

- Double precision: 64 bits

| s | exp | frac |
|---|-----|------|
| 1 | 11-bits | 52-bits |

- Extended precision: 80 bits (Intel only)

| s | exp | frac |
|---|-----|------|
| 1 | 15-bits | 63 or 64-bits |

# Visualization: Floating Point Encodings

# Additional Reading



**What Every Computer Scientist Should Know About Floating-Point Arithmetic**

DAVID GOLDBERG

*Xerox Palo Alto Research Center, 3333 Coyote Hill Road, Palo Alto, California 94304*

Floating-point arithmetic is considered an esoteric subject by many people. This is rather surprising, because floating-point is ubiquitous in computer systems: Almost every language has a floating-point datatype; computers from PCs to supercomputers have floating-point accelerators; most compilers will be called upon to compile floating-point algorithms from time to time; and virtually every operating system must respond to floating-point exceptions such as overflow. This paper presents a tutorial on the aspects of floating-point that have a direct impact on designers of computer systems. It begins with background on floating-point representation and rounding error, continues with a discussion of the IEEE floating-point standard, and concludes with examples of how computer system builders can better support floating point.

Categories and Subject Descriptors: (Primary) C.0 [**Computer Systems Organization**]: General—*instruction set design*; D.3.4 [**Programming Languages**]: Processors—*compilers, optimization*; G.1.0 [**Numerical Analysis**]: General—*computer arithmetic, error analysis, numerical algorithms* (Secondary) D.2.1 [Software Engineering]: Requirements/Specifications—*languages*; D.3.1 [**Programming Languages**]: Formal Definitions and Theory—*semantics* D.4.1 [**Operating Systems**]: Process Management—*synchronization*

General Terms: Algorithms, Design, Languages

Additional Key Words and Phrases: denormalized number, exception, floating-point, floating-point standard, gradual underflow, guard digit, NaN, overflow, relative error, rounding error, rounding mode, ulp, underflow

[What Every Computer Scientist Should Know About Floating-Point Arithmetic](),
David Goldberg, ACM Computing Surveys, 23(1), 1991

# Recap

- Representing real numbers
- Fixed Point
- Floating Point

**Next time:** *More on floating points. How can a computer perform arithmetic operating on floating points?*