# Representing and Manipulating Information: Floating point mastery

Yipeng Huang

Rutgers University

March 3, 2022

# Table of contents

# Quizzes and programming assignments

### Short quiz 4

- Due tonight, just before midnight. All about integers.

### Short quiz 5

- We will have a short quiz next week Tuesday to Thursday, all about floating point, to help you with PA3.

### Programming assignment 3

- Has been out, due next week, Thursday before spring break.

# Reading and class session plan

## Reading: CS:APP Chapter 3

- Chapter 3: Machine-level representation of programs
- Read Chapter 3.1 through 3.5 for now.

## Class session plan

- Today: finish up deep topics in floating point.
- Next Tuesday: new chapter on assembly.

# Table of contents

# Floating point numbers

### Avogadro's number
$+6.02214 \times 10^{23} \, mol^{-1}$

### Scientific notation

- ▶ sign
- ▶ mantissa or significand
- ▶ exponent

# Floats and doubles

Single precision

| 31 | 30 | 23 | 22 | 0 |
|---|---|---|---|---|

| s | exp | frac |
|---|---|---|

Double precision

| 63 | 62 | 52 | 51 | 32 |
|---|---|---|---|---|

| s | exp | frac (51:32) |
|---|---|---|

| 31 | 0 |
|---|---|

| frac (31:0) |
|---|

Figure: The two standard formats for floating point data types. Image credit CS:APP

# Floats and doubles

| property | half* | float | double |
|---|---|---|---|
| total bits | 16 | 32 | 64 |
| s bit | 1 | 1 | 1 |
| exp bits | 5 | 8 | 11 |
| frac bits | 10 | 23 | 52 |
| C printf() format specifier | None | "%f" | "%lf" |

Table: Properties of floats and doubles

# Different cases for floating point numbers

**Value of the floating point number = $(-1)^s \times M \times 2^E$**

- ▶ $E$ is encoded the exp field
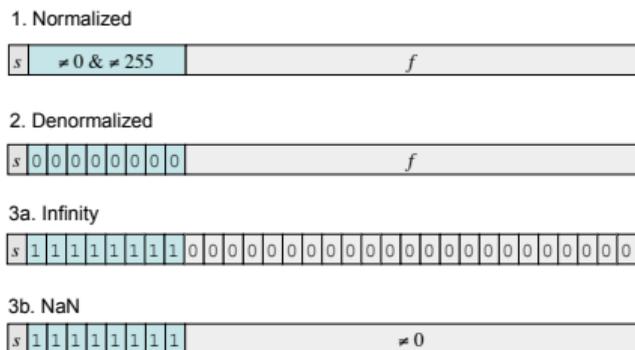- ▶ $M$ is encoded the frac field



Figure: Different cases within a floating point format. Image credit CS:APP

## Normalized and denormalized numbers

Two different cases we need to consider for the encoding of E, M
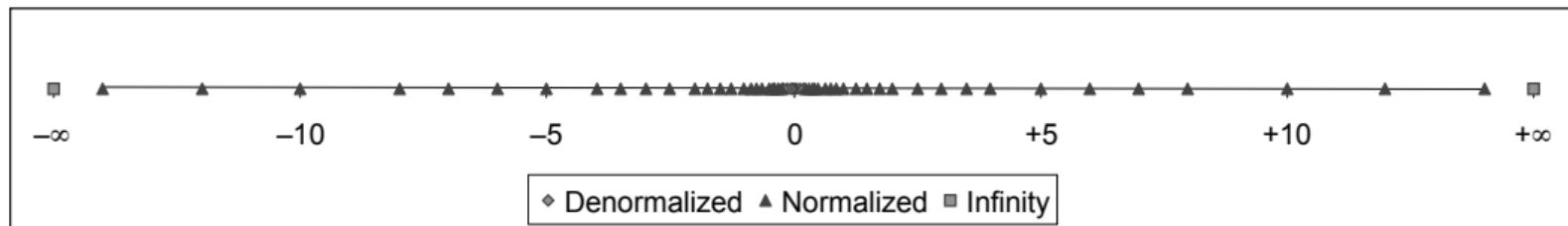
# The IEEE 754 number line



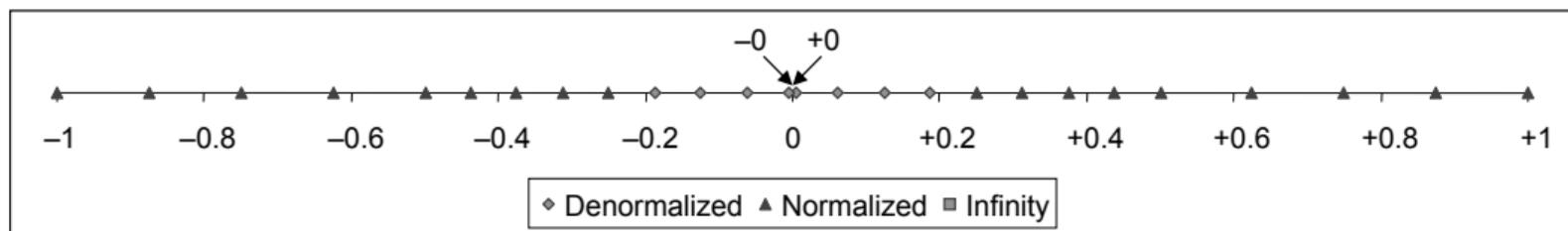Figure: Full picture of number line for floating point values. Image credit CS:APP



Figure: Zoomed in number line for floating point values. Image credit CS:APP

# Floats: Summary

|  | normalized | denormalized |
|---|---|---|
| value of number | $(-1)^s \times M \times 2^E$ | $(-1)^s \times M \times 2^E$ |
| E | E = exp-bias | E = -bias + 1 |
| bias | $2^{k-1} - 1$ | $2^{k-1} - 1$ |
| exp | $0 < exp < (2^k - 1)$ | $exp = 0$ |
| M | M = 1.frac | M = 0.frac |
|  | M has implied leading 1 | M has leading 0 |
|  | greater range | greater precision |
|  | large magnitude numbers | small magnitude numbers |
|  | denser near origin | evenly spaced |

Table: Summary of normalized and denormalized numbers

# Table of contents

# Deep understanding 1: Why is exp field encoded using bias?

exp field needs to encode both positive and negative exponents.

Why not just use one of the signed integer formats? 2's complement, 1s' complement, signed magnitude?

# Deep understanding 1: Why is exp field encoded using bias?

exp field needs to encode both positive and negative exponents.

Why not just use one of the signed integer formats? 2's complement, 1s' complement, signed magnitude?

Answer: allows easy comparison of magnitudes by simply comparing bits.

# Deep understanding 1: Why is exp field encoded using bias?

exp field needs to encode both positive and negative exponents.

Why not just use one of the signed integer formats? 2's complement, 1s' complement, signed magnitude?

Answer: allows easy comparison of magnitudes by simply comparing bits.

Consider hypothetical 8-bit floating point format (from the textbook)

1-bit sign, $k = 4$-bit exp, 3-bit frac.

What is the decimal value of 0b1_0110_111?

What is the decimal value of 0b1_0111_000?

# Deep understanding 1: Why is exp field encoded using bias?

exp field needs to encode both positive and negative exponents.

Why not just use one of the signed integer formats? 2's complement, 1s' complement, signed magnitude?

Answer: allows easy comparison of magnitudes by simply comparing bits.

Consider hypothetical 8-bit floating point format (from the textbook)

1-bit sign, $k = 4$-bit exp, 3-bit frac.

What is the decimal value of
0b1_0110_111?
$-1.875 \times 2^{-1}$

What is the decimal value of
0b1_0111_000?
$-2.000 \times 2^{-1}$

# Table of contents

# Deep understanding 2: Why have denormalized numbers?

Why not just continue normalized number scheme down to smallest numbers around zero?
Answer: makes sure that smallest increments available are maintained around zero.

Suppose denormalized numbers NOT used.

| What is the decimal value of 0b0_0000_001? | What is the decimal value of 0b0_0000_111? | What is the decimal value of 0b0_0001_000? |
|---|---|---|
| $1.125 \times 2^{-7}$ | $1.875 \times 2^{-7}$ | $2.000 \times 2^{-7}$ |

# Deep understanding 2: Why have denormalized numbers?

Why not just continue normalized number scheme down to smallest numbers around zero?

Answer: makes sure that smallest increments available are maintained around zero.

Suppose denormalized numbers ARE used.

| | | |
|---|---|---|
| What is the decimal value of 0b0_0000_001? | What is the decimal value of 0b0_0000_111? | What is the decimal value of 0b0_0001_000? |
| $0.125 \times 2^{-6}$ | $0.875 \times 2^{-6}$ | $1.000 \times 2^{-6}$ |

# Table of contents

# Floats: Special cases

| number class | when it arises | exp field | frac field |
|---:|:---:|:---|:---|
| +0 / -0 | | 0 | 0 |
| +infinity / -infinity | overflow or division by 0 | $2^k - 1$ | 0 |
| NaN not-a-number | illegal ops. such as $\sqrt{-1}$, inf-inf, inf*0 | $2^k - 1$ | non-0 |

Table: Summary of special cases

# Table of contents

# How to multiply scientific notation?

Recall: $log(x \times y) = log(x) + log(y)$

# Floating point multiplication

## FP Multiplication

- $(-1)^{s1} M1 \; 2^{E1} \; \times \; (-1)^{s2} M2 \; 2^{E2}$
- Exact Result: $(-1)^{s} M \; 2^{E}$
  - Sign s:             s1 ^ s2
  - Significand M:      M1 x  M2
  - Exponent E:         E1 + E2

- Fixing
  - If M ≥ 2, shift M right, increment E
  - If E out of range, overflow
  - Round M to fit `frac` precision

- Implementation
  - Biggest chore is multiplying significands

# Properties of floating point

## Mathematical Properties of FP Add

- Compare to those of Abelian Group
  - Closed under addition?                              Yes
    - But may generate infinity or NaN
  - Commutative?                                        Yes
  - Associative?                                        No
    - Overflow and inexactness of rounding
    - `(3.14+1e10)−1e10 = 0, 3.14+(1e10−1e10) = 3.14`
  - 0 is additive identity?
  - Every element has additive inverse?                 Yes
    - Yes, except for infinities & NaNs                 Almost
- Monotonicity
  - a ≥ b ⇒ a+c ≥ b+c?                                  Almost
    - Except for infinities & NaNs

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

29

# Properties of floating point

## Mathematical Properties of FP Mult

- Compare to Commutative Ring
  - Closed under multiplication?                     Yes
    - But may generate infinity or NaN
  - Multiplication Commutative?                       Yes
  - Multiplication is Associative?                     No
    - Possibility of overflow, inexactness of rounding
    - Ex: `(1e20*1e20)*1e-20= inf, 1e20*(1e20*1e-20)=1e20`
  - 1 is multiplicative identity?                      Yes
  - Multiplication distributes over addition?           No
    - Possibility of overflow, inexactness of rounding
    - `1e20*(1e20-1e20)=0.0, 1e20*1e20 - 1e20*1e20 =NaN`
- Monotonicity
  - $a \geq b$ & $c \geq 0 \Rightarrow a * c \geq b * c$?                  Almost
    - Except for infinities & NaNs

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

30