

# Representing and Manipulating Information: Floating point denormalized numbers and mastery

Yipeng Huang

Rutgers University

February 29, 2024

# Table of contents

## Announcements

Programming assignment 3

## Floats: Overview

### Floats: Normalized numbers

Normalized: exp field

Normalized: frac field

Normalized: example

### Floats: Denormalized numbers

Denormalized: exp field

Denormalized: frac field

Denormalized: examples

### Floats: Special cases

### Floats: Summary

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

Deep understanding 2: Why have denormalized numbers?

Deep understanding 3: Why is bias chosen to be  $2^{k-1} - 1$ ?

# Programming assignment 3

## Programming assignment 3

1. Due Friday 3/8.
2. Get started early! Plenty of background already for Parts 1 through 3.

# Table of contents

## Announcements

Programming assignment 3

## Floats: Overview

### Floats: Normalized numbers

Normalized: exp field

Normalized: frac field

Normalized: example

### Floats: Denormalized numbers

Denormalized: exp field

Denormalized: frac field

Denormalized: examples

### Floats: Special cases

### Floats: Summary

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

Deep understanding 2: Why have denormalized numbers?

Deep understanding 3: Why is bias chosen to be  $2^{k-1} - 1$ ?

# Floating point numbers

## Avogadro's number

$$+6.02214 \times 10^{23} \text{ mol}^{-1}$$

## Scientific notation

- ▶ sign
- ▶ mantissa or significand
- ▶ exponent

# Floating point numbers

## Before 1985

1. Many floating point systems.
2. Specialized machines such as Cray supercomputers.
3. Some machines with specialized floating point have had to be kept alive to support legacy software.

## After 1985

1. IEEE Standard 754.
2. A floating point standard designed for good numerical properties.
3. Found in almost every computer today, except for tiniest microcontrollers.

## Recent

1. Need for both lower precision and higher range floating point numbers.
2. Machine learning / neural networks. Low-precision tensor network processors.



# 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



# The IEEE 754 number line

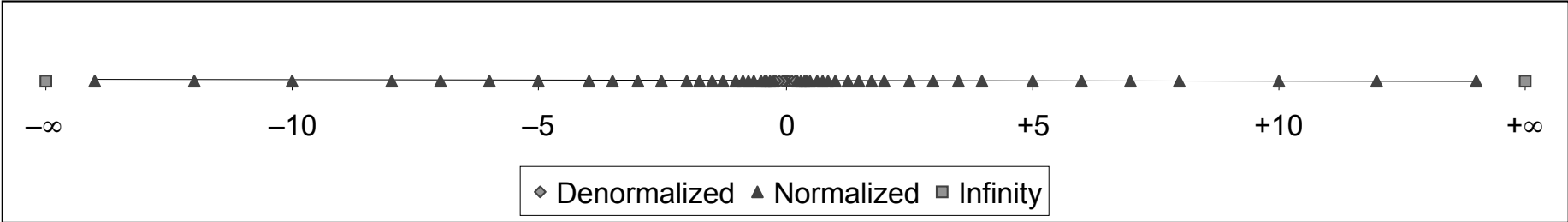


Figure: Full picture of number line for floating point values. Image credit CS:APP

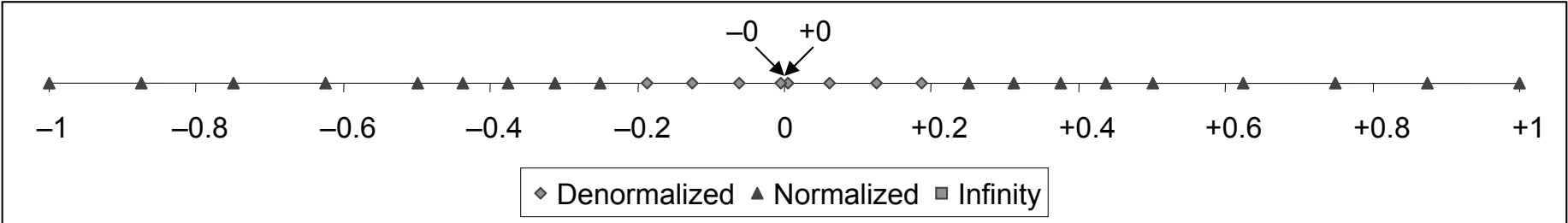


Figure: Zoomed in number line for floating point values. Image credit CS:APP

# 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

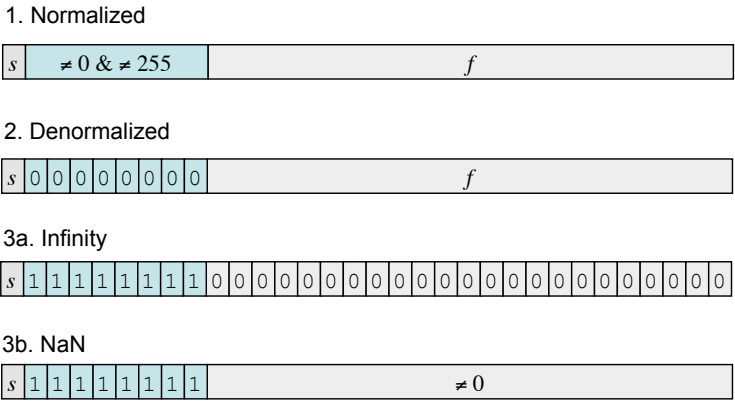


Figure: Different cases within a floating point format. Image credit CS:APP

## Normalized and denormalized numbers

Two different cases we need to consider for the encoding of  $E, M$

# Table of contents

## Announcements

Programming assignment 3

## Floats: Overview

### Floats: Normalized numbers

Normalized: exp field

Normalized: frac field

Normalized: example

### Floats: Denormalized numbers

Denormalized: exp field

Denormalized: frac field

Denormalized: examples

### Floats: Special cases

### Floats: Summary

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

Deep understanding 2: Why have denormalized numbers?

Deep understanding 3: Why is bias chosen to be  $2^{k-1} - 1$ ?



# Normalized: frac field

$$M = 1.\text{frac}$$

# Normalized: example

- ▶ 12.375 to single-precision floating point
- ▶ sign is positive so  $s=0$
- ▶ binary is  $1100.011_2$
- ▶ in other words it is  $1.100011_2 \times 2^3$
- ▶  $\text{exp} = E + \text{bias} = 3 + 127 = 130 = 1000\_0010_2$
- ▶  $M = 1.100011_2 = 1.\text{frac}$
- ▶  $\text{frac} = 100011$

# Table of contents

## Announcements

Programming assignment 3

## Floats: Overview

### Floats: Normalized numbers

Normalized: exp field

Normalized: frac field

Normalized: example

### Floats: Denormalized numbers

Denormalized: exp field

Denormalized: frac field

Denormalized: examples

### Floats: Special cases

### Floats: Summary

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

Deep understanding 2: Why have denormalized numbers?

Deep understanding 3: Why is bias chosen to be  $2^{k-1} - 1$ ?

$$\begin{array}{l} 6.02 \cdot 10^{23} \\ 6.03 \cdot 10^{22} \\ 0. \end{array} \quad \begin{array}{l} \updownarrow \\ \updownarrow \end{array} \quad 10^{21}$$



# The IEEE 754 number line

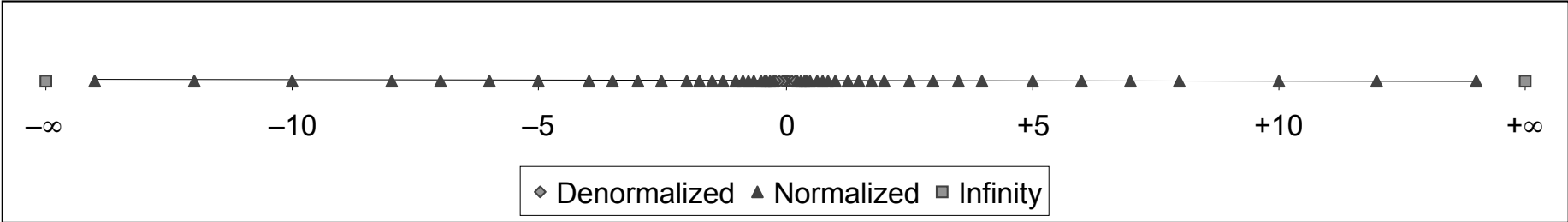


Figure: Full picture of number line for floating point values. Image credit CS:APP

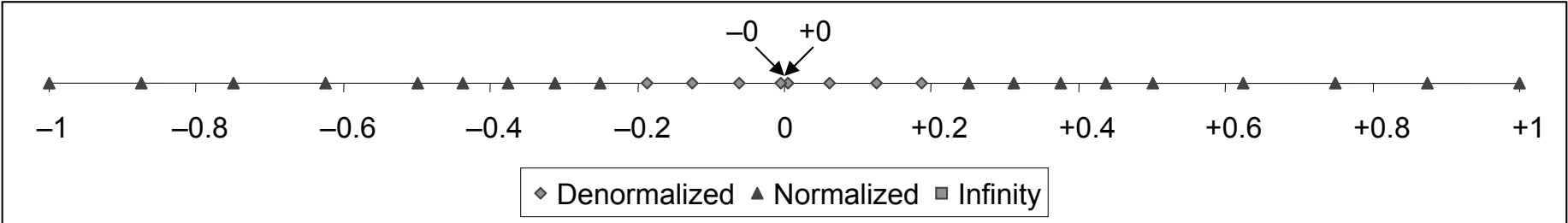


Figure: Zoomed in number line for floating point values. Image credit CS:APP



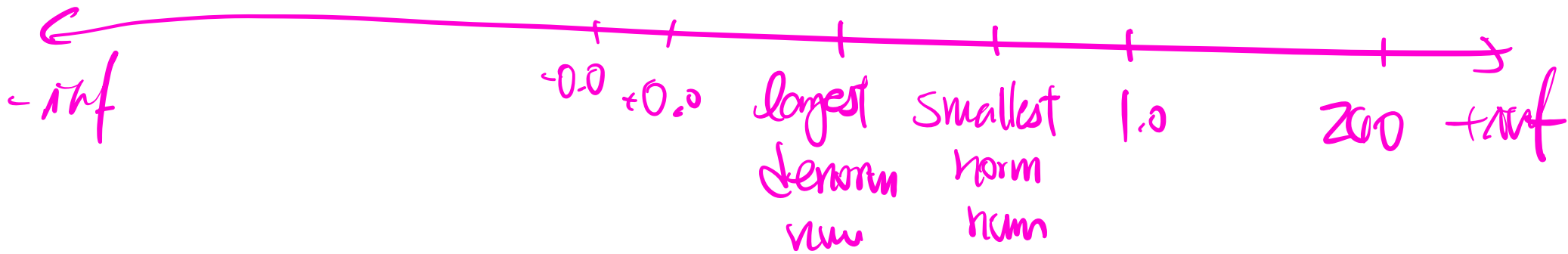
# Denormalized: frac field

$M = 0.\text{frac}$

value represented leading with 0

# Denormalized: examples

8-bit, 1 bit sign, 4-bit exp field  
3 bit frac.



$$(-1)^S \cdot M \cdot 2^E$$

+inf. 0\_1111\_000

largest  
pos  
number

$$0_1110_111 \Rightarrow (-1)^0 \cdot M \cdot 2^E = +1.875 \cdot 2^7 = 240$$

$$E = \text{exp. bias}$$

$$= 14 - (2^{k-1} - 1)$$

$$= 14 - (2^3 - 1)$$

$$= 7$$

$$M = 1.111_2$$

$$= 1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{8}$$

$$= 1.75 + .125$$

$$= 1.875$$

$$+1.0 = (-1)^S \cdot M \cdot 2^E \\ = (-1)^0 \cdot 1.0 \cdot 2^0$$

$$1.000 = 1.\text{frac}$$

$$0_0111_000$$

$$E = \text{exp. bias}$$

$$0 = 7 - 7$$

smallest positive  
norm number

$$0\_0001\_000$$

$$(-1)^0 \cdot M \cdot 2^E$$

$$= +1.0 \cdot 2^E$$

$$= +1.0 \cdot 2^{-6}$$

$$= \frac{1}{64}$$

$$E = \text{exp-bias}$$

$$= 1 - 7$$

$$= -6$$

largest denorm.  
pos number

$$0\_0000\_111$$

$$\Rightarrow (-1)^0 \cdot M \cdot 2^E = 0.875 \cdot 2^{-6}$$

$$M = 0.111 = \frac{1}{2} + \frac{1}{4} + \frac{1}{8} = 0.875$$

$$E = 1 - \text{bias} = -6$$

smallest pos  
num (denorm)

$$0\_0000\_001$$

$$M = 0.001_2 = \frac{1}{8}$$

$$E = -6$$

$$\Rightarrow 0.125 \cdot 2^{-6} = \frac{1}{8} \cdot \frac{1}{64} = \frac{1}{512}$$

+0.0

$$0\_0000\_000$$

$$\frac{+0.0}{\text{inf}} = +0.0$$

-0.0

$$1\_0000\_000$$

$$\frac{-1.0}{\text{inf}} = -0.0$$

Commutative :

$$A+B = B+A$$

$$238 + 239 = 239 + 238$$

$$A+(-B) = (-B)+A$$

$$A \cdot B = B \cdot A$$

< In the tiny float  
number system >

Associative

$$(A+B)+C = A+(B+C)$$

$$\begin{array}{l} (240+1)+(-239) \\ = \infty + (-239) \\ = \infty \end{array} \left\{ \begin{array}{l} 240+(1+(-239)) \\ = 240+(-238) \\ = +2 \end{array} \right.$$

Distributive

$$A(B+C) = AB+AC$$

$$\begin{array}{l} \frac{1}{2}(240+1) \\ = \frac{1}{2} \cdot \infty \\ = \infty \end{array} \quad \begin{array}{l} \frac{1}{2} \cdot 240 + \frac{1}{2} \cdot 1 \\ = 120 + \frac{1}{2} \\ = 120.5 \end{array}$$

# Table of contents

## Announcements

Programming assignment 3

## Floats: Overview

### Floats: Normalized numbers

Normalized: exp field

Normalized: frac field

Normalized: example

### Floats: Denormalized numbers

Denormalized: exp field

Denormalized: frac field

Denormalized: examples

### Floats: Special cases

### Floats: Summary

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

Deep understanding 2: Why have denormalized numbers?

Deep understanding 3: Why is bias chosen to be  $2^{k-1} - 1$ ?



What's the best that  
 a double precision float can  
 do to represent  $1 + 10^{12} = (10^3)^4 + 1$

$$(-1)^S \cdot M \cdot E$$

$$\approx 1024^4 + 1$$

$$= (2^{10})^4 + 1$$

$$S = 0$$

$$E \approx 40$$

$E = \text{exp-bias}$

$$\text{exp} = E + \text{bias}$$

$$= 40 + (2^{(11-1)} - 1)$$

$$= 40 + (2^{10} - 1)$$

$$= 40 + 1023$$

$$= 1063$$

52 bits

$$M = 1.\underbrace{000 \dots 000}_{52 \text{ bits}}$$

$$1.\underbrace{000 \dots 001}_{52 \text{ bit frac}}$$

$$\underline{52} > \underline{40} ?$$



# Table of contents

## Announcements

Programming assignment 3

## Floats: Overview

### Floats: Normalized numbers

Normalized: exp field

Normalized: frac field

Normalized: example

### Floats: Denormalized numbers

Denormalized: exp field

Denormalized: frac field

Denormalized: examples

### Floats: Special cases

### Floats: Summary

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

Deep understanding 2: Why have denormalized numbers?

Deep understanding 3: Why is bias chosen to be  $2^{k-1} - 1$ ?

# 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 has implied leading 1	M = 0.frac M has leading 0
	greater range large magnitude numbers denser near origin	greater precision small magnitude numbers evenly spaced

**Table:** Summary of normalized and denormalized numbers

# Table of contents

## Announcements

Programming assignment 3

## Floats: Overview

### Floats: Normalized numbers

Normalized: exp field

Normalized: frac field

Normalized: example

### Floats: Denormalized numbers

Denormalized: exp field

Denormalized: frac field

Denormalized: examples

### Floats: Special cases

### Floats: Summary

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

Deep understanding 2: Why have denormalized numbers?

Deep understanding 3: Why is bias chosen to be  $2^{k-1} - 1$ ?

# 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

## Announcements

Programming assignment 3

## Floats: Overview

### Floats: Normalized numbers

Normalized: exp field

Normalized: frac field

Normalized: example

### Floats: Denormalized numbers

Denormalized: exp field

Denormalized: frac field

Denormalized: examples

### Floats: Special cases

### Floats: Summary

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

Deep understanding 2: Why have denormalized numbers?

Deep understanding 3: Why is bias chosen to be  $2^{k-1} - 1$ ?

## 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`?

$$1.125 \times 2^{-7}$$

What is the decimal value of `0b0_0000_111`?

$$1.875 \times 2^{-7}$$

What is the decimal value of `0b0_0001_000`?

$$2.000 \times 2^{-7}$$

## 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`?

$$0.125 \times 2^{-6}$$

What is the decimal value of `0b0_0000_111`?

$$0.875 \times 2^{-6}$$

What is the decimal value of `0b0_0001_000`?

$$1.000 \times 2^{-6}$$

# Table of contents

## Announcements

Programming assignment 3

## Floats: Overview

### Floats: Normalized numbers

Normalized: exp field

Normalized: frac field

Normalized: example

### Floats: Denormalized numbers

Denormalized: exp field

Denormalized: frac field

Denormalized: examples

### Floats: Special cases

### Floats: Summary

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

Deep understanding 2: Why have denormalized numbers?

Deep understanding 3: Why is bias chosen to be  $2^{k-1} - 1$ ?