

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

Yipeng Huang

Rutgers University

February 27, 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

6.02214 0000 ...

$6.02214 \times 10^{23}$   
 $6.02215 \times 10^{23}$  ↑

↕

$10^{18}$  ↓

# 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

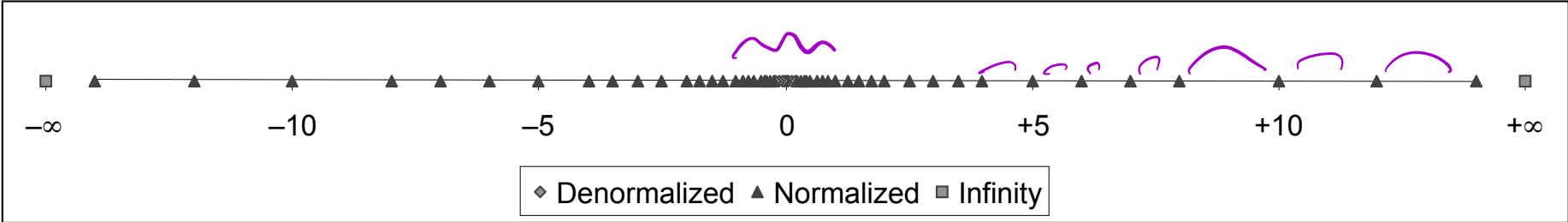


Figure: Full picture of number line for floating point values. Image credit CS:APP

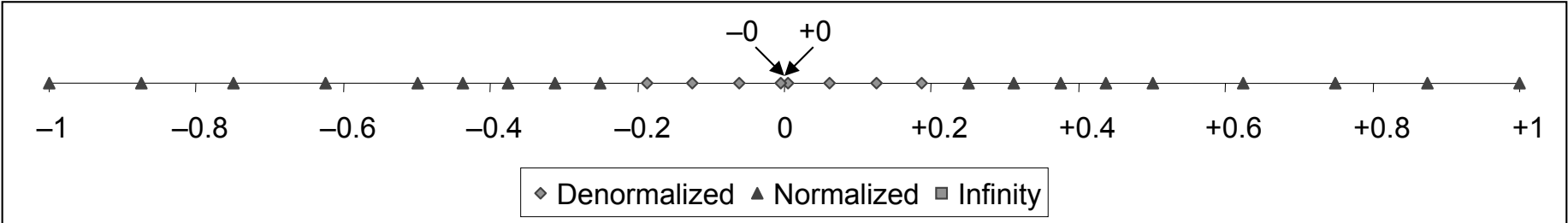


Figure: Zoomed in number line for floating point values. Image credit CS:APP

# 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

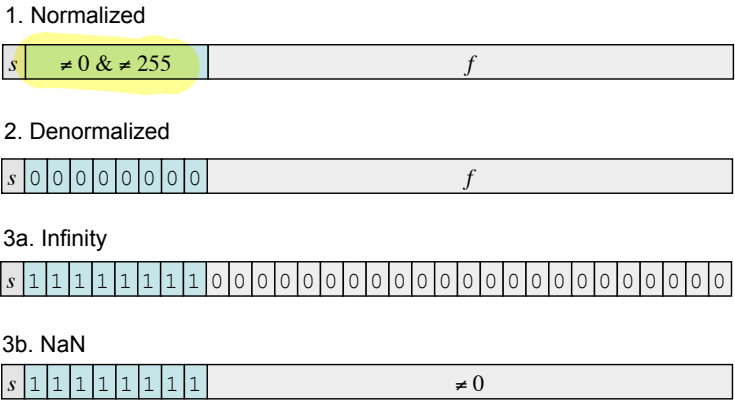


Figure: Different cases within a floating point format. Image credit CS:APP

## 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?

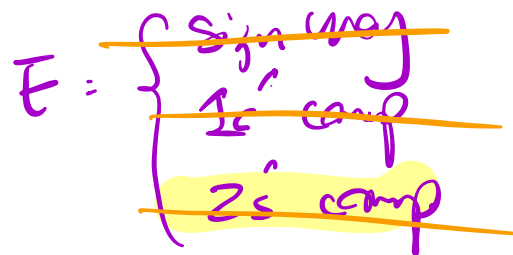
Deep understanding 3: Why is bias chosen to be  $2^{k-1} - 1$ ?

# Normalized: exp field

$$(-1)^s \times M \times 2^E$$

For normalized numbers,  
 $0 < \text{exp} < 2^k - 1$

- ▶ exp is a  $k$ -bit unsigned integer



## Bias

- ▶ need a bias to represent negative exponents

▶ bias =  $2^{k-1} - 1 = 2^{(k-1)} - 1 = 1024 - 1 = 1023$

- ▶ bias is the  $k$ -bit unsigned integer:

011..111

For normalized numbers,

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

In other words,  $\text{exp} = E + \text{bias}$

	property	float	double
	k	8	11
	bias	127	1023
	smallest E (greatest precision)	-126	-1022
	largest E (greatest range)	127	1023

Table: Summary of normalized exp field

$$\begin{aligned} E &= 0000\text{..}0001 - \text{bias} \\ &= 1_{10} - 127_{10} \\ &= -126 \end{aligned}$$

$$\begin{aligned} E &= 1111\text{..}1110 - \text{bias} \\ &= (2^k - 1 - 1) - \text{bias} \\ &= 2023 - 1 - 1 - 1023 = 1023 \end{aligned}$$

Normalized: frac field

$$\text{Number} = (-1)^S \times M \times 2^E$$

↓

$$1.01110$$

↑

$M = 1.\text{frac}$

# Normalized: example

$$12.375_{10} = (-1)^s \times M \times 2^E$$

$$= (-1)^0 \times \left( 8_{10} + 4_{10} + \frac{1}{4} + \frac{1}{8} \right)$$

- ▶ 12.375<sub>10</sub> to single-precision floating point
- ▶ sign is positive so s=0
- ▶ binary is 1100.011<sub>2</sub>
- ▶ in other words it is 1.100011<sub>2</sub> × 2<sup>3</sup>
- ▶ exp = E + bias = 3 + 127 = 130 = 1000\_0010<sub>2</sub>
- ▶ M = 1.100011<sub>2</sub> = 1.frac
- ▶ frac = 100011

$$1 \times 8 + 1 \times 4 + 0 \times 2 + 0 \times 1 + 0 \times \frac{1}{2} + 1 \times \frac{1}{4} + 1 \times \frac{1}{8}$$

$$= 8 + 4 + 0.25 + .125$$

$$= 12.375_{10}$$

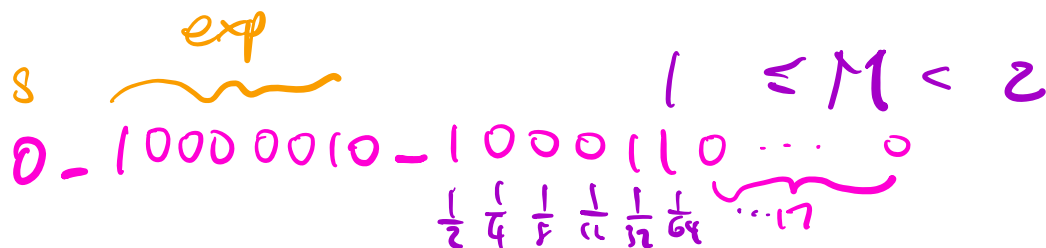
$$= 1100.011_2$$

$$= 1.100011_2 \times 2^3$$

$$60.22 \times 10^{22}$$

$$6.022 \times 10^{23}$$

$$0.6022 \times 10^{24}$$



16-bit float.

1-bit sign.

5-bit exp field

10-bit frac.

1\_01100\_0010100000

$$\begin{aligned} \text{Number} &= (-1)^S \times M \times 2^E \\ &= (-1)^1 \times 1.15625 \times 2^{-3} \end{aligned}$$

$$\begin{aligned} M = 1.\text{frac.} &= 1.0010100000 \\ &= 1. + 0 \times \frac{1}{2} + 0 \times \frac{1}{4} + 1 \times \frac{1}{8} + 0 \times \frac{1}{16} + 1 \times \frac{1}{32} \\ &= 1 + 0.125 + 0.03125 \\ &= 1.15625 \end{aligned}$$

E: exp-bits

$$= 01100_2 - 01111_2$$

$$= (8+4) - (2^4 - 1)$$

$$= 12 - 15$$

$$= -3$$

# 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  $2^{k-1} - 1$ ?



# The IEEE 754 number line

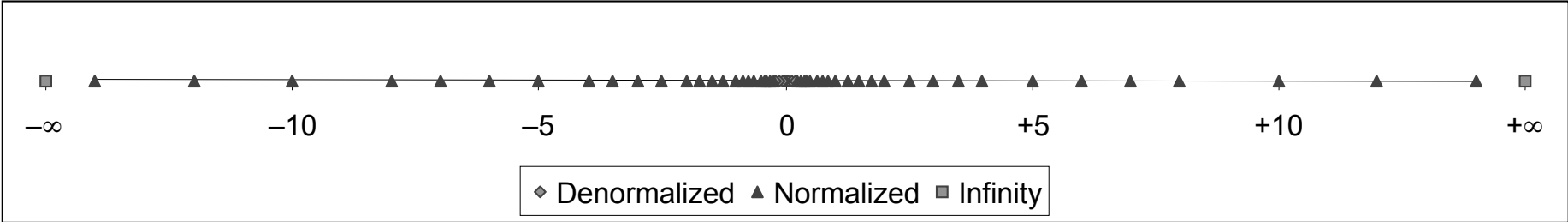


Figure: Full picture of number line for floating point values. Image credit CS:APP

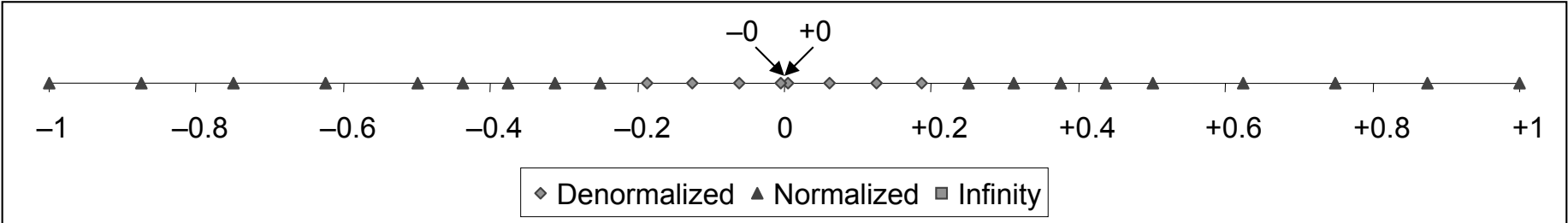


Figure: Zoomed in number line for floating point values. Image credit CS:APP

# Denormalized: exp field

For denormalized numbers,  $\text{exp} = 0$

## Bias

- ▶ need a bias to represent negative exponents
- ▶  $\text{bias} = 2^{k-1} - 1$
- ▶ bias is the  $k$ -bit unsigned integer:  
011..111

For denormalized numbers,  
 $E = 1\text{-bias}$

property	float	double
k	8	11
bias	127	1023
E	-126	-1022

Table: Summary of denormalized exp field

# Denormalized: frac field

$M = 0.\text{frac}$

value represented leading with 0

# Denormalized: examples

# 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: 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

## 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  $2^{k-1} - 1$ ?