

Representing and Manipulating Information: Fixed point, floating point normalized and denormalized numbers

Yipeng Huang

Rutgers University

February 22, 2024

Table of contents

Announcements

Programming assignment 2

Integers and basic arithmetic

Representing negative and signed integers

Fractions and fixed point representation

`monteCarloPi.c` Using floating point and random numbers to estimate PI

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

Programming assignment 2

Programming assignment 2

1. Due Friday 2/23.
2. Graph algorithms and hash table.

Table of contents

Announcements

Programming assignment 2

Integers and basic arithmetic

Representing negative and signed integers

Fractions and fixed point representation

`monteCarloPi.c` Using floating point and random numbers to estimate PI

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

Representing negative and signed integers

Ways to represent negative numbers

1. Sign magnitude
2. 1s' complement
3. 2's complement

Representing negative and signed integers

2's complement

1) flip all bits
+1

2) 

signed char	weight in decimal
00000001 ↙	1
00000010	2
00000100	4
00001000	8
00010000	16
00100000	32
01000000	64
10000000	-128

Table: Weight of each bit in a signed char type

- ▶ what is the most positive value you can represent? 127
- ▶ what is the most negative value you can represent? -128
- ▶ how to represent -1? 11111111
- ▶ how to represent -2? 11111110

Representing negative and signed integers

2's complement

signed char	weight in decimal
00000001	1
00000010	2
00000100	4
00001000	8
00010000	16
00100000	32
01000000	64
10000000	-128

Table: Weight of each bit in a signed char type

- ▶ MSB: 1 for negative
- ▶ To make a number negative: flip all bits and add 1.
- ▶ Addition in 2's complement is sound

Importance of paying attention to limits of encoding

$$15 + 15 = 30$$

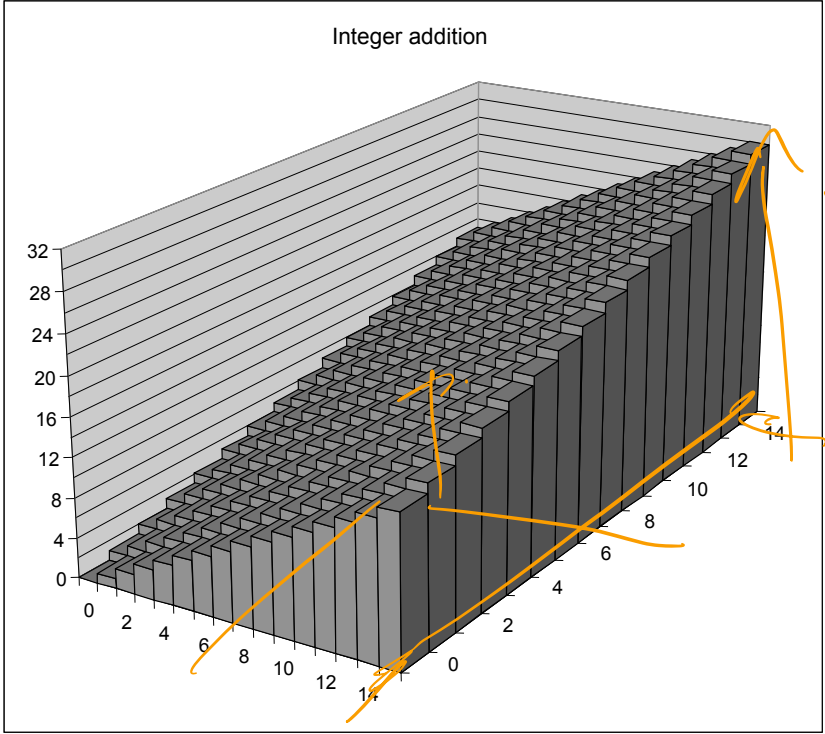


Figure: Image credit: CS:APP

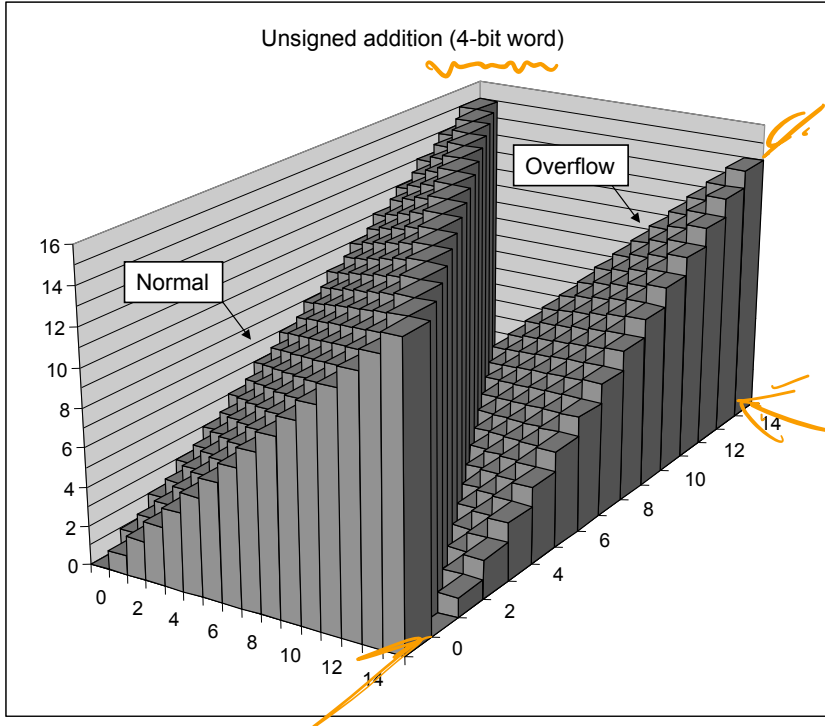


Figure: Image credit: CS:APP

Importance of paying attention to limits of encoding

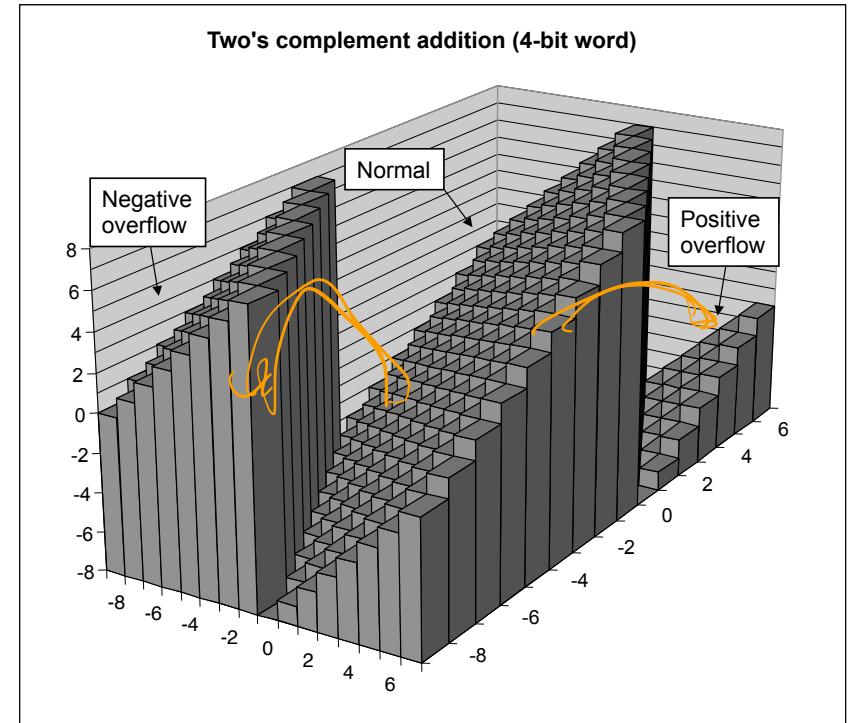
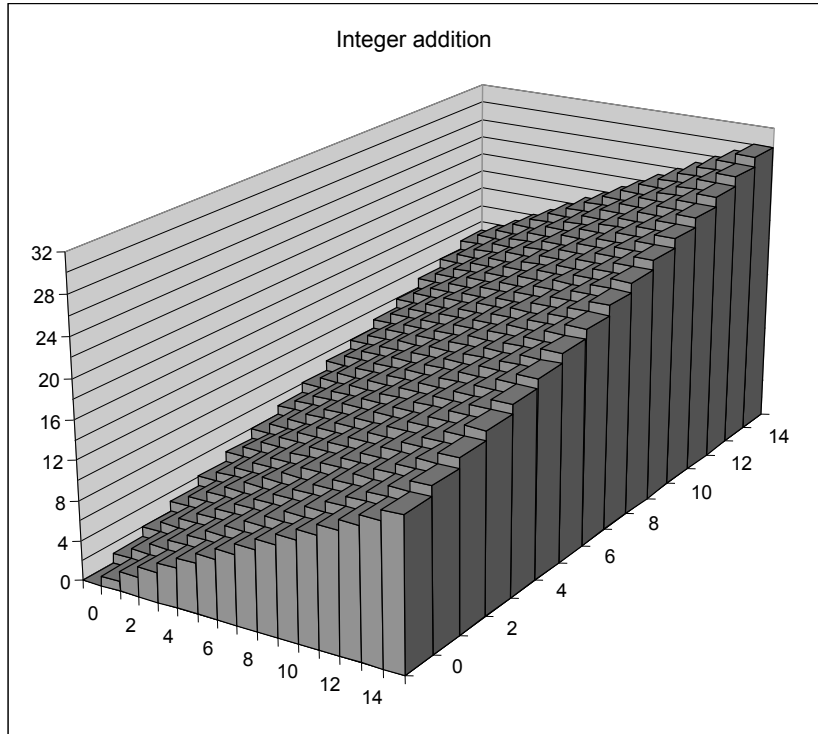


Figure: Image credit: CS:APP

Figure: Image credit: CS:APP

<https://www.theatlantic.com/technology/archive/2014/12/how-gangnam-style-broke-youtube/383389/>

$$32 = 2^{30+2} = 2^{30} \times 2^2 = 2^2 \times 2^{30}$$

$$= 4 \times 2^{(10+10+10)} = 4 \times (2^{10})^3$$

$$\approx 4 \times (1000)^3 = 4 \times 10^9 = 4 \text{ Billion}$$

GB = gigabyte = 10^9 bytes

GiB = gibibyte = $(2^{10})^3$

Table of contents

Announcements

Programming assignment 2

Integers and basic arithmetic

Representing negative and signed integers

Fractions and fixed point representation

`monteCarloPi.c` Using floating point and random numbers to estimate PI

Floats: Overview

Floats: Normalized numbers

Normalized: exp field

Normalized: frac field

Normalized: example

Floats: Denormalized numbers

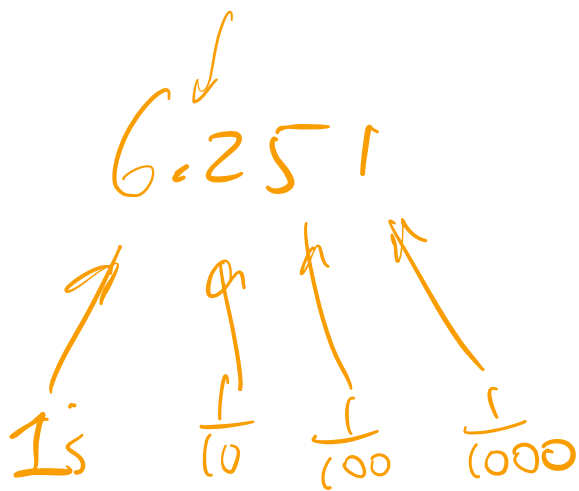
Denormalized: exp field

Denormalized: frac field

Denormalized: examples

Floats: Special cases

Floats: Summary



06.251

6.2510000.

2^{''} ←

2.000000'' ←

Unsigned fixed-point binary for fractions

fractional
decimal
scientific notation

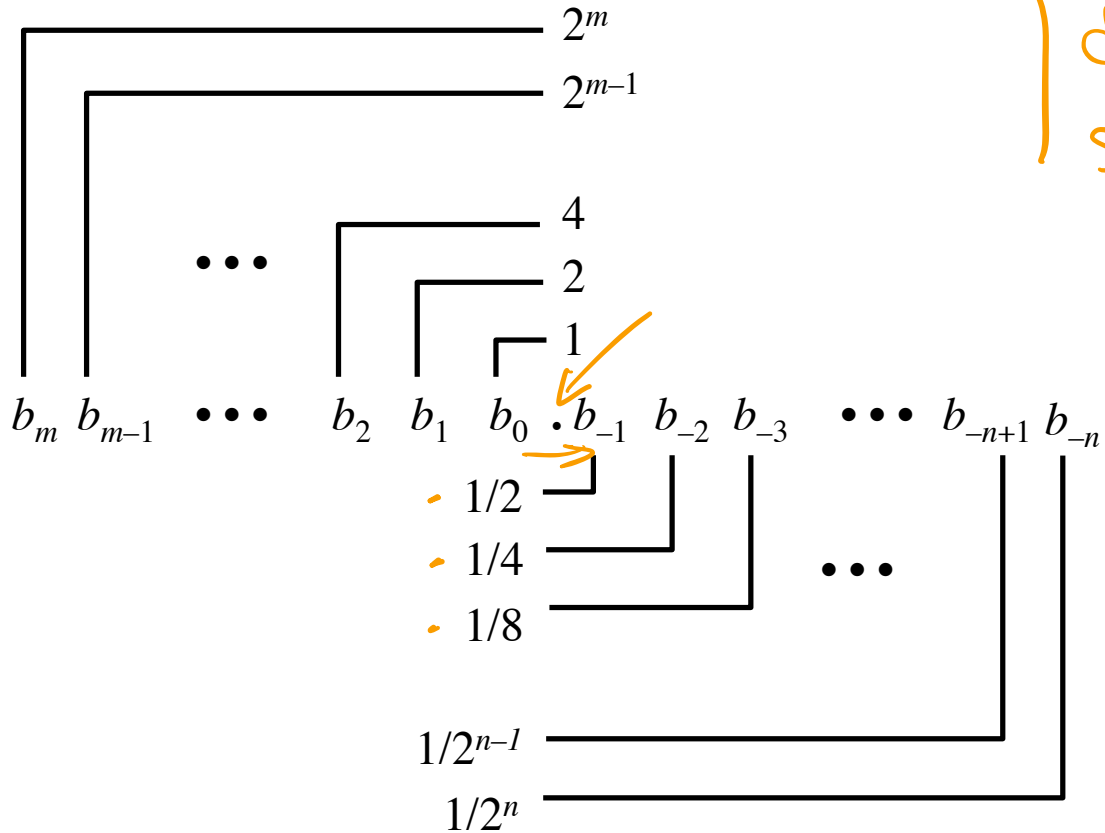


Figure: Fractional binary. Image credit CS:APP

$$3.14 \xrightarrow{\text{shift left by 1}} 31.4$$

"

$$3.14 \times 10^1$$

$$3.14_{10}$$

8-bit unsigned fixed point w/
binary point at 4 places
from MSB.

$$\begin{array}{cccc|cccc} \hline / & / & / & / & - & / & / & / & / \\ 8 & 4 & 2 & 1 & & \frac{1}{2} & \frac{1}{4} & \frac{1}{8} & \frac{1}{16} \\ \hline \end{array}$$

$$0011.0010 \rightarrow 2 + 1 + \frac{1}{8} = 3.125$$

$$\begin{array}{r} 0.14_{10} \\ 0.125_{10} \\ \hline 0.015_{10} \\ 0.0625 \end{array}$$

$$\frac{1}{3} = 0.\overline{3}_{10}$$

$$\begin{aligned} 0.\overline{1}_2 &= 0.1\overline{1}_2 = 0.11\overline{1}_2 \\ &= \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \frac{1}{16} + \frac{1}{32} + \frac{1}{64} \\ &\rightarrow 1 \end{aligned}$$

Limitations of fixed-point

- ▶ Can only represent numbers of the form $x/2^k$
- ▶ Cannot represent numbers with very large magnitude (great range) or very small magnitude (great precision)

Bit shifting

$\ll N$ Left shift by N bits

- ▶ multiplies by 2^N
- ▶ $2 \ll 3 = 0000_0010_2 \ll 3 = 0001_0000_2 = 16 = 2 * 2^3$
- ▶ $-2 \ll 3 = 1111_1110_2 \ll 3 = 1111_0000_2 = -16 = -2 * 2^3$

$\gg N$ Right shift by N bits


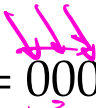


- ▶ divides by 2^N
- ▶ $16 \gg 3 = 0001_0000_2 \gg 3 = 0000_0010_2 = 2 = 16/2^3$

- ▶ -16 \gg 3 = 1111_0000_2 \gg 3 = 1111_1110_2 = -2 = $-16/2^3$


Table of contents

Announcements

Programming assignment 2

Integers and basic arithmetic

Representing negative and signed integers

Fractions and fixed point representation

`monteCarloPi.c` Using floating point and random numbers to estimate PI

Floats: Overview

Floats: Normalized numbers

Normalized: exp field

Normalized: frac field

Normalized: example

Floats: Denormalized numbers

Denormalized: exp field

Denormalized: frac field

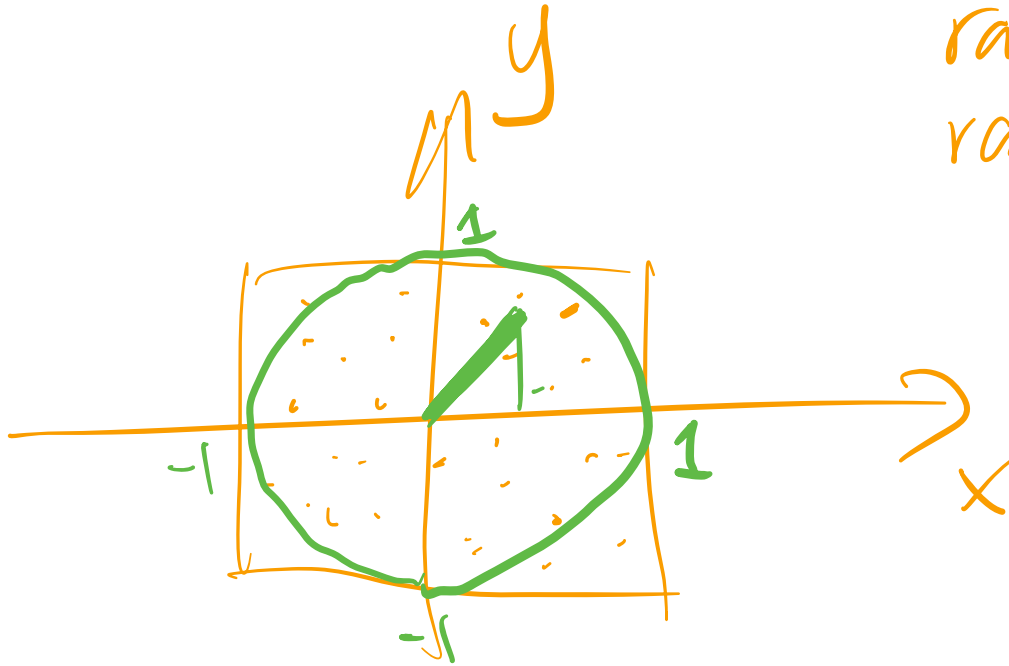
Denormalized: examples

Floats: Special cases

Floats: Summary

monteCarloPi.c Using floating point and random numbers to estimate PI

π



rand x

rand y

```
for (1 to 100000) {  
  if ((rand x + rand y) < 1)
```

add to
total insides

add total

```
}  
inside  
total
```

$$1 \cdot 1 \cdot \pi = \text{ratio} \cdot 1 \cdot 1 \cdot 4$$

$$\pi = 4 \cdot \text{ratio}$$

Table of contents

Announcements

Programming assignment 2

Integers and basic arithmetic

Representing negative and signed integers

Fractions and fixed point representation

`monteCarloPi.c` Using floating point and random numbers to estimate PI

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

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

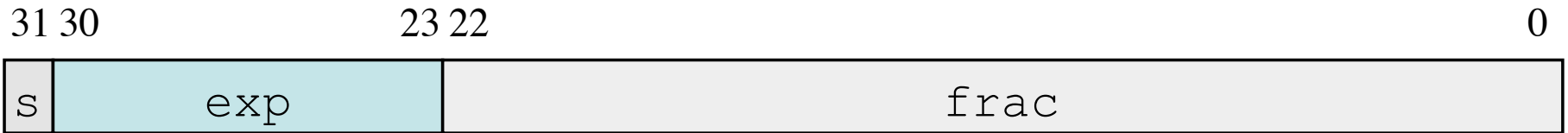
1. IEEE Standard 754.
2. A floating point standard designed for good numerical properties.
3. Found in almost every computer today, except for tiniest microcontrollers.

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

Single precision



Double precision

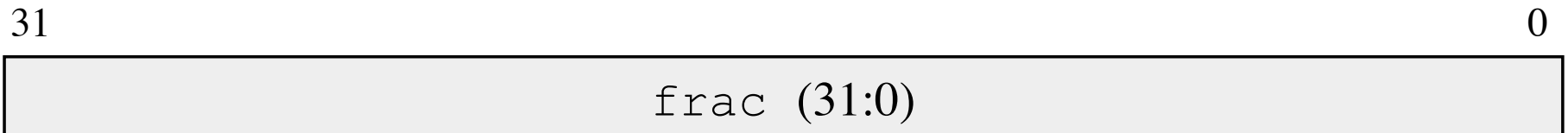
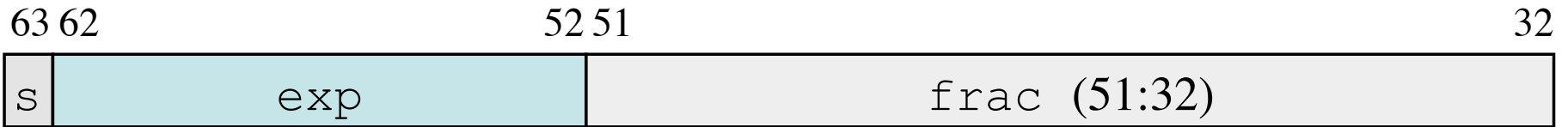


Figure: The two standard formats for floating point data types. Image credit CS:APP

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

The IEEE 754 number line

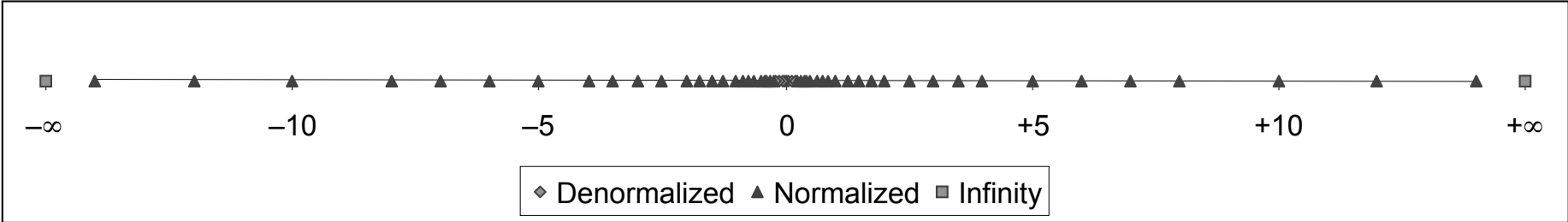


Figure: Full picture of number line for floating point values. Image credit CS:APP

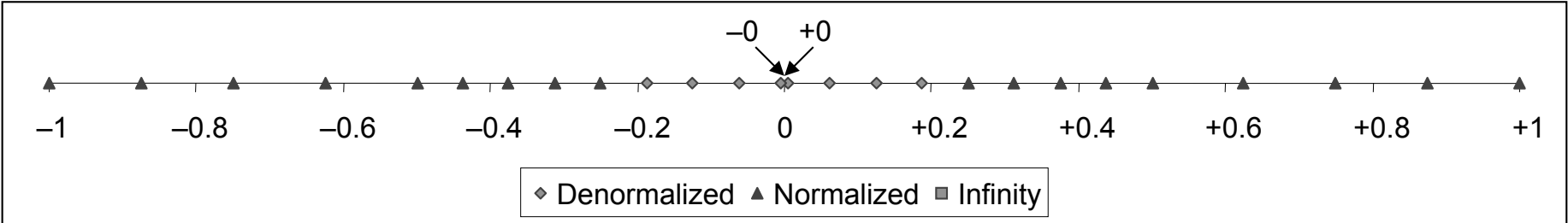


Figure: Zoomed in number line for floating point values. Image credit CS:APP

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

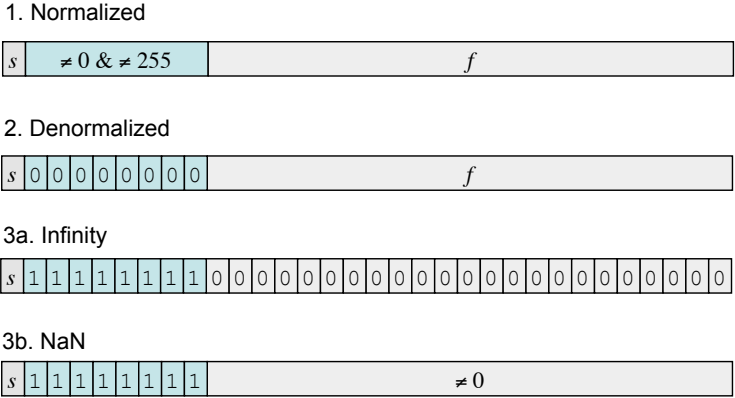


Figure: Different cases within a floating point format. Image credit CS:APP

Normalized and denormalized numbers

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

Table of contents

Announcements

Programming assignment 2

Integers and basic arithmetic

Representing negative and signed integers

Fractions and fixed point representation

`monteCarloPi.c` Using floating point and random numbers to estimate PI

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

Normalized: exp field

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 = $2^{k-1} - 1$
- ▶ 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

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 2

Integers and basic arithmetic

Representing negative and signed integers

Fractions and fixed point representation

`monteCarloPi.c` Using floating point and random numbers to estimate PI

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

The IEEE 754 number line

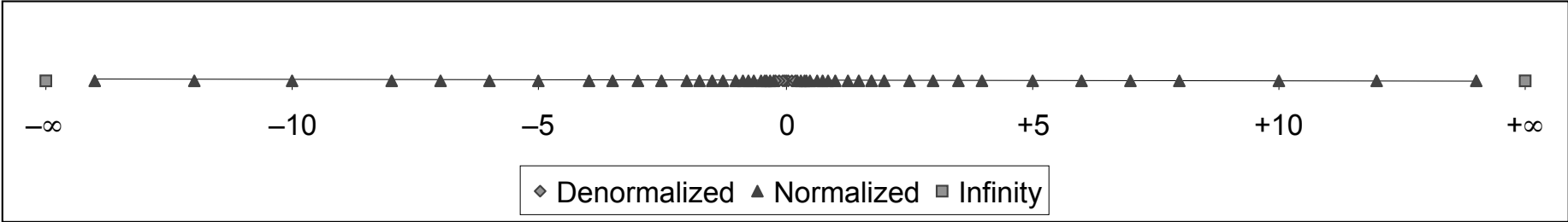


Figure: Full picture of number line for floating point values. Image credit CS:APP

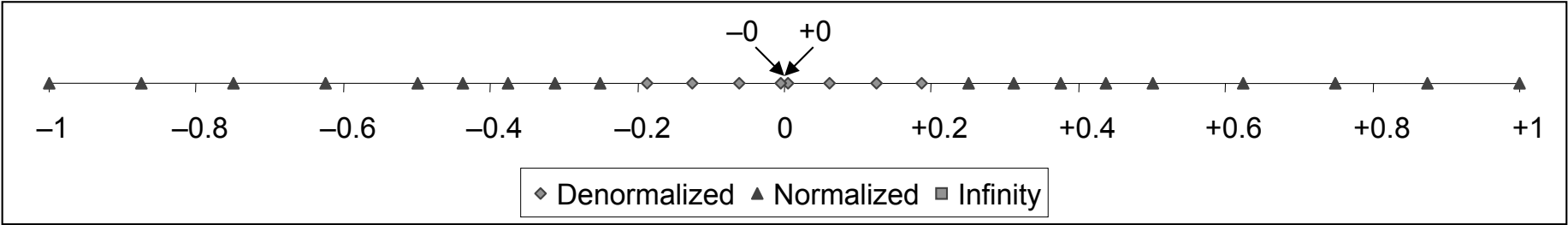


Figure: Zoomed in number line for floating point values. Image credit CS:APP

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 2

Integers and basic arithmetic

Representing negative and signed integers

Fractions and fixed point representation

`monteCarloPi.c` Using floating point and random numbers to estimate PI

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

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 2

Integers and basic arithmetic

Representing negative and signed integers

Fractions and fixed point representation

`monteCarloPi.c` Using floating point and random numbers to estimate PI

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

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