Data representation: Integers, Fixed Point, Floating Point

Yipeng Huang

Rutgers University

October 9, 2025

Table of contents

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

Ways to represent	negative numbers (Mpanlan	Addition
1. Sign magnitude		unconventional
2. 1s' complement	redundant Zeno	Sound, remember to add
3. 2's complement	UMque Zeno	Sound Corry

- 1s' complement (-
- → Flip all bits
 - ► Addition in 1s' complement is sound 🗸
 - ► In this encoding there are 2 encodings for 0
 - ► -0: 0b1111
 - ► +0: 0b0000

most regaline

number

Ob1000-2000

Ob1((1_(11))

-(270)

06/1111_1110 One -100 One

2's complement

power-of-two

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

- 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

most veg numbre
0 6 000 - 0000
= - 1 2 f 60

most positue number 060111_1111 = +(270

negone 06(1112-1111)

one 060000-000

 $+123_{00} - 52_{00}$ $1 \times 69 + 1 \times 32 + 1 \times 16 + 1 \times 9 + 0 \times 9 + 1 \times 21 \times 16$ $0 \times 60 + 1 \times 10 + 1 \times 11$

00/100-1100

170

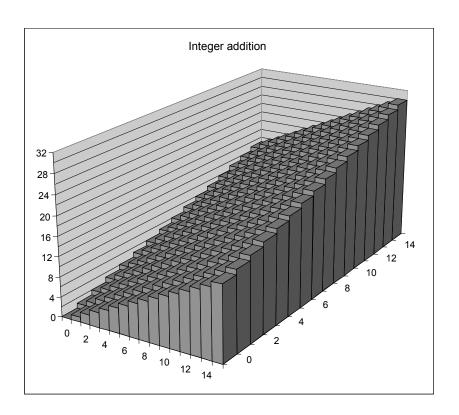
2's complement

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

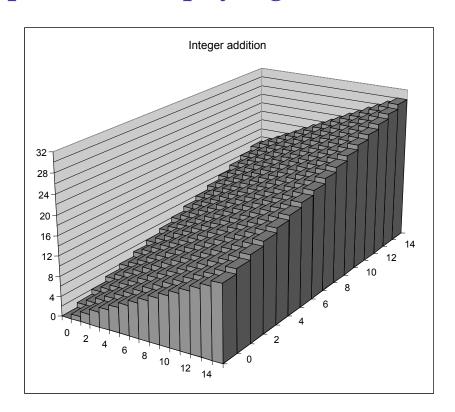


Unsigned addition (4-bit word) Overflow Normal

Figure: Image credit: CS:APP

Figure: Image credit: CS:APP

Importance of paying attention to limits of encoding



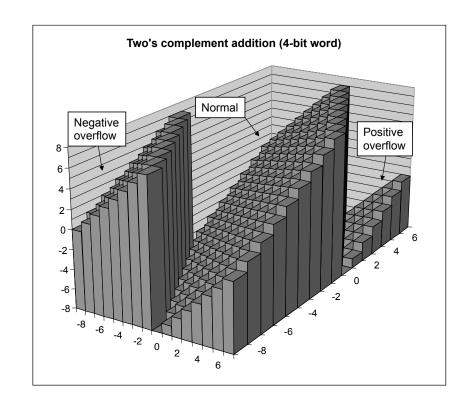


Figure: Image credit: CS:APP

Figure: Image credit: CS:APP

https://www.theatlantic.com/technology/archive/2014/12/ how-gangnam-style-broke-youtube/383389/ signed int- (32-bie)

Most pos number 060[11_----__1111

$$Z^{(6)} = Z^{(10+6)} = Z^{(2)} = Z^{(2)} = [024 \cdot Z^{(6)}]$$

$$= Z^{(3+10+2)} = Z^{(3+10+2)} = Z^{(3+10+2)} = Z^{(2)} = Z^{(2$$

Table of contents

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

Unsigned fixed-point binary for fractions

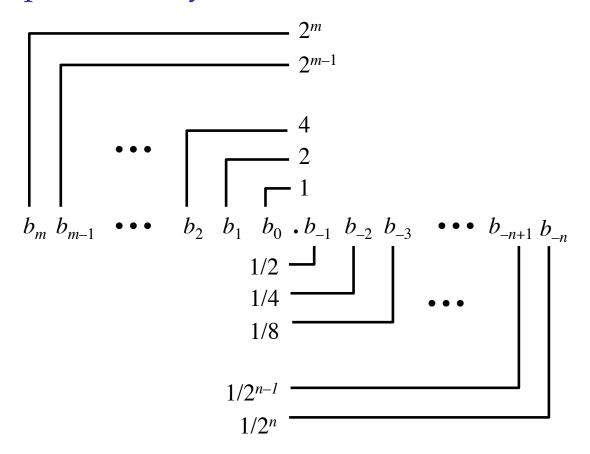


Figure: Fractional binary. Image credit CS:APP

Unsigned fixed-point binary for fractions

unsigned fixed-point char example	weight in decimal
1000.0000	8
0100.0000	4
0010.0000	2
0001.0000	1
0000.1000	0.5
0000.0100	0.25
0000.0010	0.125
0000.0001	0.0625

Table: Weight of each bit in an example fixed-point binary number

- ightharpoonup $.625 = .5 + .125 = 0000.1010_2$
- \triangleright 1001.1000₂ = 9 + .5 = 9.5

Signed fixed-point binary for fractions

signed fixed-point char example	weight in decimal
1000.0000	-8
0100.0000	4
0010.0000	2
0001.0000	1
0000.1000	0.5
0000.0100	0.25
0000.0010	0.125
0000.0001	0.0625

Table: Weight of each bit in an example fixed-point binary number

$$-.625 = -8 + 4 + 2 + 1 + 0 + .25 + .125 = 1111.0110_2$$

$$ightharpoonup 1001.1000_2 = -8 + 1 + .5 = -6.5$$

Limitations of fixed-point

- \triangleright 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

<< *N* Left shift by N bits

- ightharpoonup multiplies by 2^N
- \triangleright 2 << 3 = 0000_0010₂ << 3 = 0001_0000₂ = 16 = 2 * 2³
- $-2 << 3 = 1111_1110_2 << 3 = 1111_0000_2 = -16 = -2 * 2^3$

>> *N* Right shift by N bits

- ightharpoonup divides by 2^N
- $ightharpoonup 16 >> 3 = 0001_0000_2 >> 3 = 0000_0010_2 = 2 = 16/2^3$
- $-16 >> 3 = 1111_0000_2 >> 3 = 1111_1110_2 = -2 = -16/2^3$

123₁₀ 060111_1011

signed char number= 123; pointf (god", number >> 0); // (Z] printf ("%), numbers 1): 1/6/: 123/2 060///_(0//) 050011 - 1101; 3216+8+9+1 326 8471 = 61 prints (%) number >>6; " (number >>6) & Obl 00014+ 1811 8 0h0000-000

Table of contents

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

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

Table of contents

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

Floating point numbers

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

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

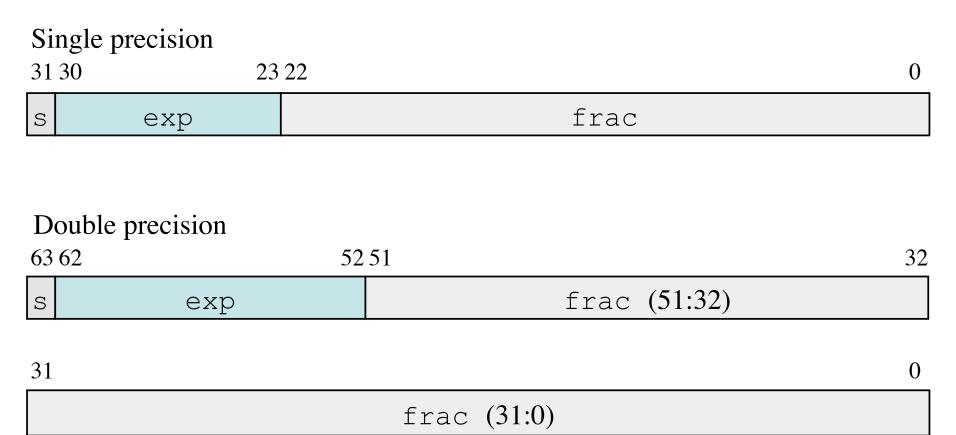


Figure: The two standard formats for floating point data types. Image credit CS:APP

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

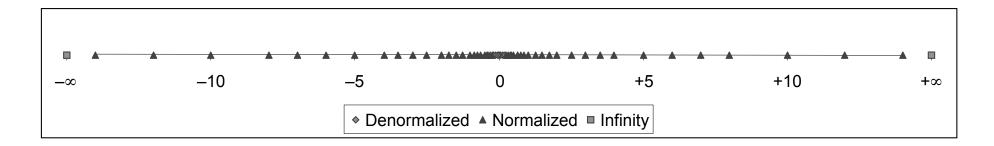


Figure: Full picture of number line for floating point values. Image credit CS:APP

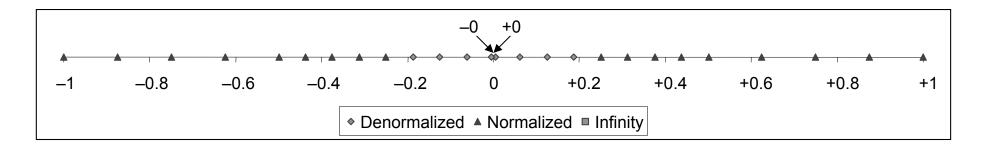


Figure: Zoomed in number line for floating point values. Image credit CS:APP

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

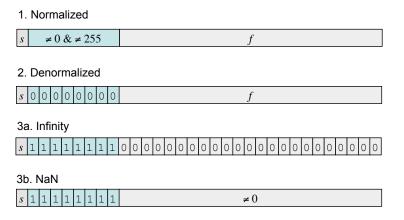


Figure: Different cases within a floating point format. Image credit CS:APP

Normalized and denormalized numbers

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

Table of contents

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

Normalized: exp field

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

exp is a k-bit unsigned integer

Bias

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

For normalized numbers, E = exp-bias

In other words, exp = E + 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.frac

Normalized: example

- ► 12.375 to single-precision floating point
- \triangleright sign is positive so s=0
- ▶ binary is 1100.011₂
- \triangleright in other words it is 1.100011₂ \times 2³
- ightharpoonup exp = $E + \text{bias} = 3 + 127 = 130 = 1000_0010_2$
- $M = 1.100011_2 = 1.frac$
- ightharpoonup frac = 100011