Data representation: Fixed Point, Floating Point Normalized Numbers

Yipeng Huang

Rutgers University

October 14, 2025

Table of contents

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

Unsigned fixed-point binary for fractions

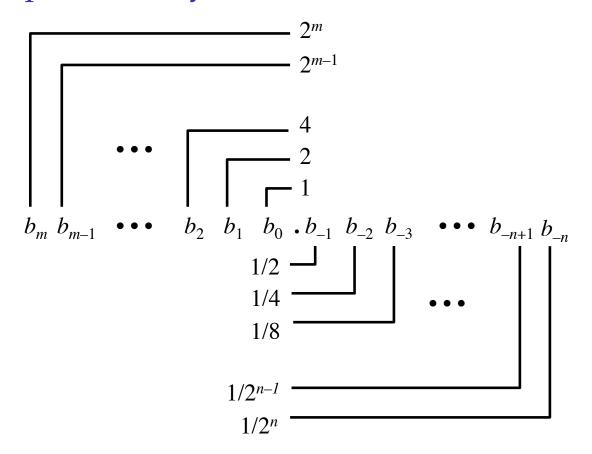


Figure: Fractional binary. Image credit CS:APP

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 $.625 = .5 + .125 = 0000.1010_2$
- ightharpoonup 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 r
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

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

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

-76 = -3.[4159 -... -8.00000° 3-1409 3-1409 2 1,14159 0.125

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

<< *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
- $ightharpoonup 16 >> 3 = 0001_0000_2 >> 3 = 0000_0010_2 = 2 = 16/2^3$
- $-16 >> 3 = 1111_0000_2 >> 3 = 1111_1110_2 = -2 = -16/2^3$

Table of contents

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

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

Table of contents

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

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

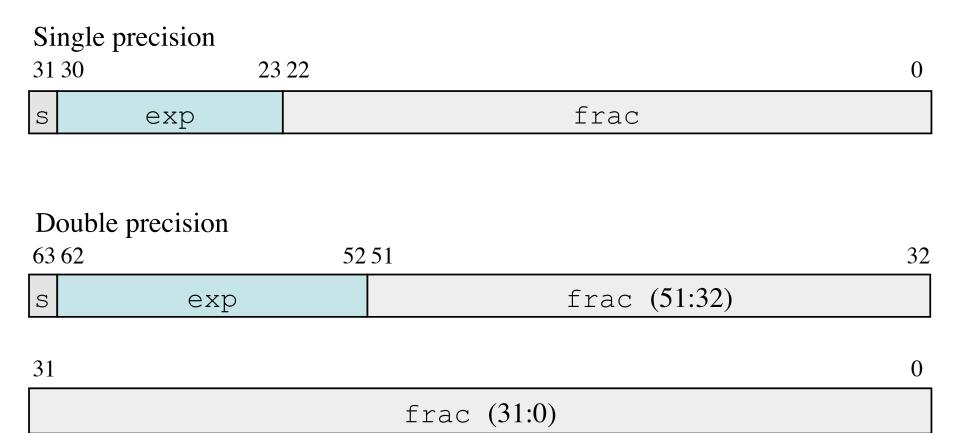


Figure: The two standard formats for floating point data types. Image credit CS:APP

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

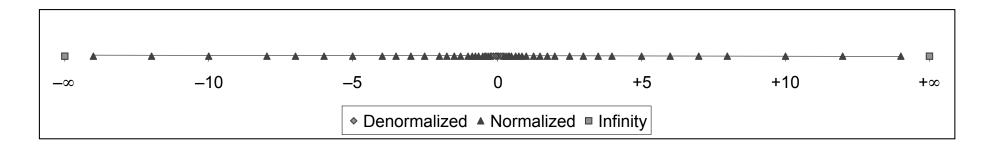


Figure: Full picture of number line for floating point values. Image credit CS:APP

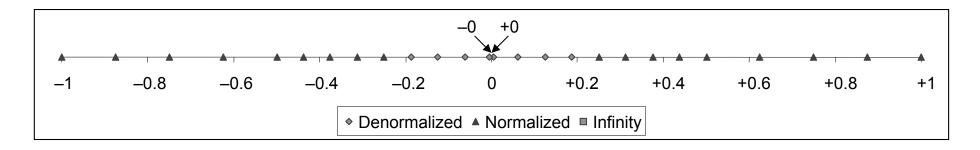


Figure: Zoomed in number line for floating point values. Image credit CS:APP

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

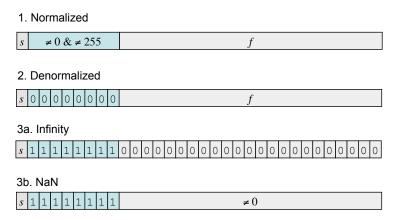


Figure: Different cases within a floating point format. Image credit CS:APP

Normalized and denormalized numbers

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

Table of contents

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: 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
- ightharpoonup 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

- ► 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