Data representation: Floating Point Denormalized Numbers and Mastery

Yipeng Huang

Rutgers University

October 21, 2025

Table of contents

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: Understanding its design

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: Properties

Floating point multiplication

Floats: Summary

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

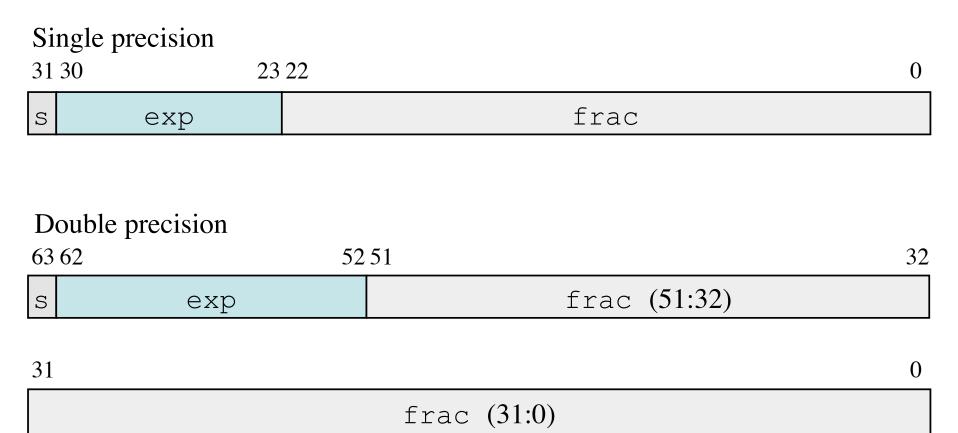


Figure: The two standard formats for floating point data types. Image credit CS:APP

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

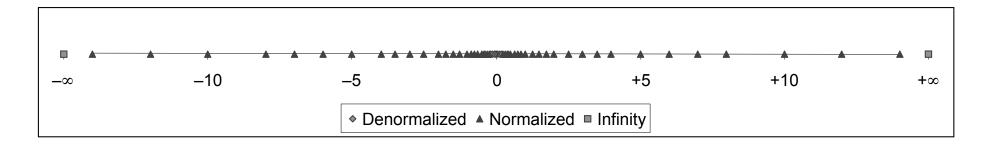


Figure: Full picture of number line for floating point values. Image credit CS:APP

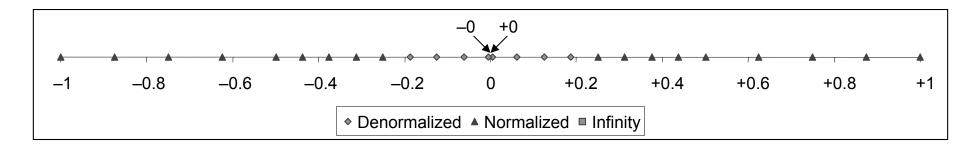


Figure: Zoomed in number line for floating point values. Image credit CS:APP

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

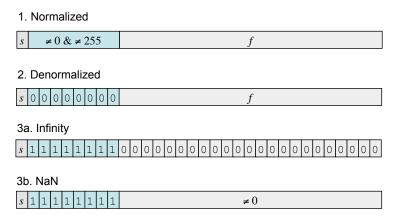


Figure: Different cases within a floating point format. Image credit CS:APP

Normalized and denormalized numbers

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

Table of contents

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: Understanding its design

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: Properties

Floating point multiplication

Floats: Summary

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

positive infinity

largest normalizes number

epsilon
$$0 = 0(1 = 0.01)$$

= $(-1)^{5} - (1.00(2) \cdot 7^{-7})$
= $(1 + \frac{1}{6}) \cdot 2^{5} = \frac{9}{6}$

Opsilon: F

one

0-0111-000

Smallest normalized number

0.0001.000
$$= (-1)^{S} \cdot (monthsa) \cdot Z^{E}$$

$$= (+1) \cdot (1.000_{2}) \cdot Z^{(exp-bias)}$$

$$= (+1) \cdot ((-0) \cdot Z^{(1-7)})$$

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

bias =
$$(2^{(4-1)})$$

= $2^{(4-1)}$
= 7

Table of contents

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: Understanding its design

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: Properties

Floating point multiplication

Floats: Summary

The IEEE 754 number line

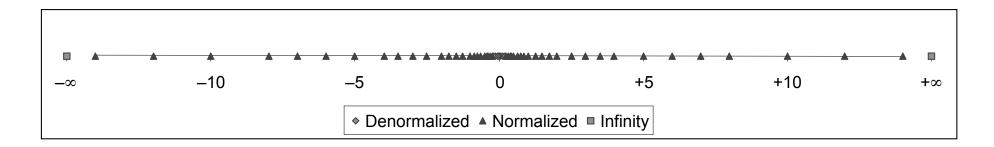


Figure: Full picture of number line for floating point values. Image credit CS:APP

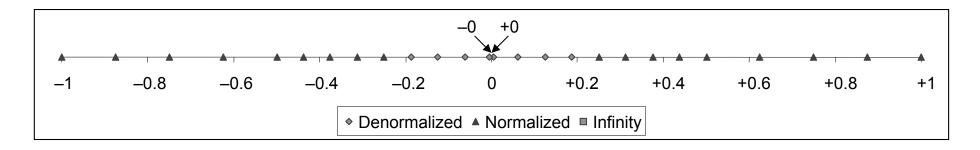


Figure: Zoomed in number line for floating point values. Image credit CS:APP

Denormalized: exp field

For denormalized numbers, exp = 0

Bias

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

For denormalized numbers, E = 1-bias

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

Table: Summary of denormalized exp field

Denormalized: frac field

M = 0.frac value represented leading with 0 Denormalized: examples

largest denormalizes number

Without special treatment and how it works)

with spewal treatment.

$$0_{-}0000_{-}111$$
= $(-1)^{s} \cdot (0_{-}111_{z}) \cdot 7$
= $(-1)^{s} \cdot (0_{-}$

Smallest denormalized number

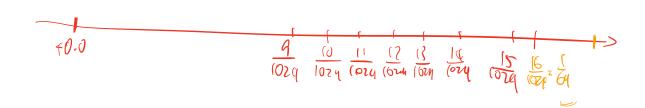
Without speared treatment < not how it works

$$0.0000.00[$$

$$=)(-1)^{3}-(1.001)\cdot Z$$

$$=([+\frac{1}{5})\cdot Z^{-7}]$$

$$=\frac{9}{5}\cdot Z^{-7}=\frac{9}{7^{2}\cdot 2^{-7}}=\frac{9}{(0.29)}$$



With special treatment

$$0-0000-00|$$

$$= -1 \cdot (0.00) \cdot 7$$

Table of contents

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: Understanding its design

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: Properties

Floating point multiplication

Floats: Summary

Floats: Special cases

number class	when it arises	exp field	frac field
+0 / -0		0	0
<pre>+infinity / -infinity</pre>	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

Why have both +0.0 and -0.0?

$$\frac{1}{\sqrt{\infty}} = +\infty$$

Table of contents

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: Understanding its design

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: Properties

Floating point multiplication

Floats: Summary

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?

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.

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?

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 $0b1_0110_1111?$ -1.875×2^{-1}

What is the decimal value of $0b1_0111_000?$ -2.000×2^{-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 × 2⁻⁷

What is the decimal value of $0b0_0000_111?$ 1.875×2^{-7}

What is the decimal value of $0b0_0001_000?$ 2.000×2^{-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 $0b0_0000_001$? 0.125×2^{-6}

What is the decimal value of $0b0_0000_111$? 0.875×2^{-6}

What is the decimal value of $0b0_0001_000?$ 1.000×2^{-6}

Table of contents

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: Understanding its design

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: Properties

Floating point multiplication

Floats: Summary

How to multiply scientific notation?

Recall: $log(x \times y) = log(x) + log(y)$

Floating point multiplication

Carnegie Mellon

FP Multiplication

- $-(-1)^{s1} M1 2^{E1} x (-1)^{s2} M2 2^{E2}$
- Exact Result: (-1)^s M 2^E
 - Sign s: s1 ^ s2
 - Significand M: M1 x M2
 - Exponent E: E1 + E2
- Fixing
 - If M ≥ 2, shift M right, increment E
 - If E out of range, overflow
 - Round M to fit frac precision
- Implementation
 - Biggest chore is multiplying significands

Table of contents

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: Understanding its design

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: Properties

Floating point multiplication

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
$\dot{\mathrm{M}}$	M = 1.frac	M = 0.frac
	M has implied leading 1	M has leading 0
	greater range	greater precision
	large magnitude numbers denser near origin	small magnitude numbers evenly spaced

Table: Summary of normalized and denormalized numbers

Connecting to actual number ranges for 32-bit floor and 64-bit double: