Data representation: Floating Point Normalized Numbers and Denormalized Numbers

Yipeng Huang

Rutgers University

October 16, 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: Summary

Floating point numbers

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

Scientific notation

- sign 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

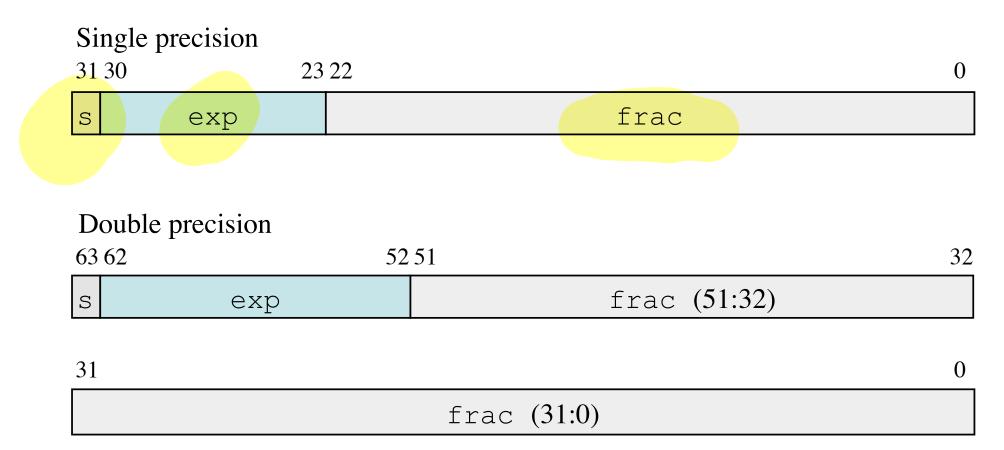


Figure: The two standard formats for floating point data types. Image credit CS:APP

Floats and doubles

text book
tiny float

8
1
4
3

property	half*	float	double	quad
total bits	16	32	64	(28)
s bit	1	1	1	1
exp bits frac bits	5	8	11	15
frac bits	10	23	52	112
C printf() format specifier	None	''%f''	''%lf''	

Table: Properties of floats and doubles

The IEEE 754 number line

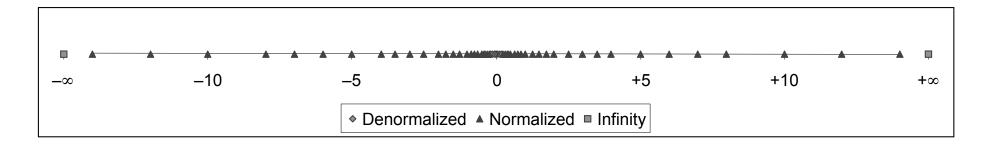


Figure: Full picture of number line for floating point values. Image credit CS:APP

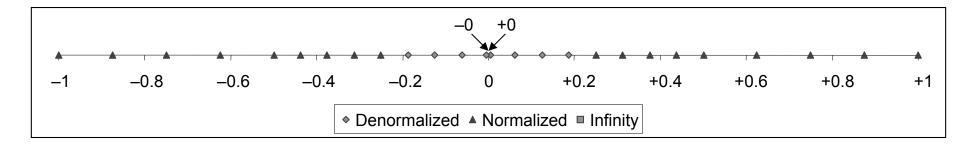


Figure: Zoomed in number line for floating point values. Image credit CS:APP

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

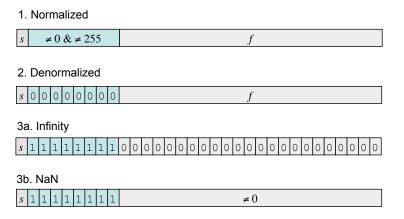


Figure: Different cases within a floating point format. Image credit CS:APP

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: 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

inflimity

0-1111-000

largest (normalized) number

Mantitsa: 1. frac
$$= 1.[1]_{b} = [+2+4+5] = [5]$$

$$= 1.55$$

$$\frac{239}{239} \rightarrow 0_{110} = (1.10_{1}) \cdot (1.4600) \cdot 2^{2}$$

$$= (1.110_{1}) \cdot 12f$$

$$= (1+1+2) \cdot 12f = 1.12f:32-120$$

epsilon

ON

$$[0.0] = (0.0) - Z_0$$

Smallest (normalized) number < non Zero>

exp frac
0.0001 000

 $[-1] \cdot [1-frac] \cdot Z$ $= (+1) \cdot (1-000) \cdot Z$ $= [-2] \cdot Z$ $= [-3] \cdot Z$ $= [-3] \cdot Z$ $= [-6] \cdot Z$

Zen

s exp frac. 0_0000_000 Sentity

0

1

Comnutative

AB=BAV

associationery

 $(AB)C = A(BC) \times$

740+1-64

[240+1] -64 = 00-69 =00

Z40+ ((-64) = Z40-63 = [7]

dytributative

240 (69-64) = 240×0 = 0

Z90.69 - Z90.69

= 00 - 10

= NaN

Inverse

Yes, negation always constr

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: Summary

The IEEE 754 number line

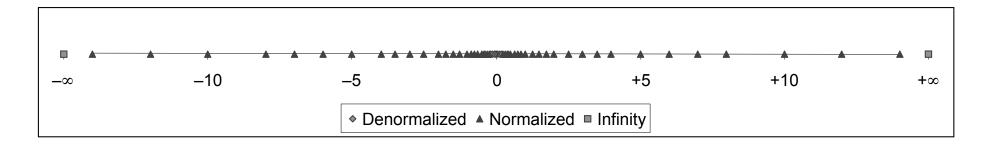


Figure: Full picture of number line for floating point values. Image credit CS:APP

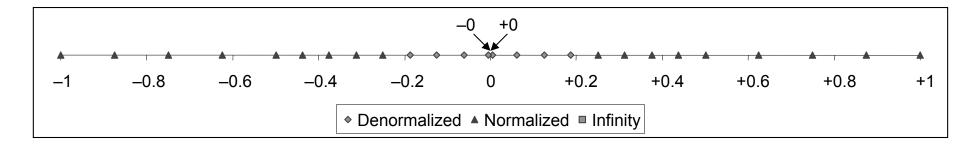


Figure: Zoomed in number line for floating point values. Image credit CS:APP

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 denormatized number

smallest denormed of number

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: 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

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: Summary

Floats: Summary

	normalized	denormalized
value of number	$(-1)^s \times M \times 2^E$	$(-1)^s \times M \times 2^E$
	$E = \exp$ -bias	E = -bias + 1
bias	$2^{k-1}-1$	$2^{k-1}-1$
exp	$0 < exp < (2^k - 1)$	exp = 0
$\dot{ ext{M}}$	M = 1.frac	M = 0.frac
	M has implied leading 1	M has leading 0
	greater range	greater precision
	large magnitude numbers	small magnitude numbers
	denser near origin	evenly spaced

Table: Summary of normalized and denormalized numbers

I. Why use bias encoding?

I. Why have Lenormalized numbers?

PROPERTIES.