

# CS 261

## Fall 2016

Mike Lam, Professor

# Floating-Point Numbers

# Floating-point

- Topics
  - Binary fractions
  - Floating-point representation
  - Conversions and rounding error

# Binary fractions

- Now we can store integers
  - But what about general real numbers?
- Extend binary integers to store fractions
  - Designate a certain number of bits for the fractional part
  - These bits represent negative powers of two
  - (Just like fractional digits in decimal fractions!)

**101.101**  
4    2    1    1/2   1/4   1/8

$$4 + 1 + 0.5 + 0.125 = \mathbf{5.625}$$

*(alternatively: 5 + 5/8)*

# Examples

Representation	Value	Decimal
$0.0_2$	$\frac{0}{2}$	$0.0_{10}$
$0.01_2$	$\frac{1}{4}$	$0.25_{10}$
$0.010_2$	$\frac{2}{8}$	$0.25_{10}$
$0.0011_2$	$\frac{3}{16}$	$0.1875_{10}$
$0.00110_2$	$\frac{6}{32}$	$0.1875_{10}$
$0.001101_2$	$\frac{13}{64}$	$0.203125_{10}$
$0.0011010_2$	$\frac{26}{128}$	$0.203125_{10}$
$0.00110011_2$	$\frac{51}{256}$	$0.19921875_{10}$

# Another problem

- For scientific applications, we want to be able to store a wide *range* of values
  - From the scale of galaxies down to the scale of atoms
- Doing this with fixed-precision numbers is difficult
  - Even signed 64-bit integers
    - Perhaps allocate half for whole number, half for fraction
    - Range:  $\sim 2 \times 10^{-9}$  through  $\sim 2 \times 10^9$

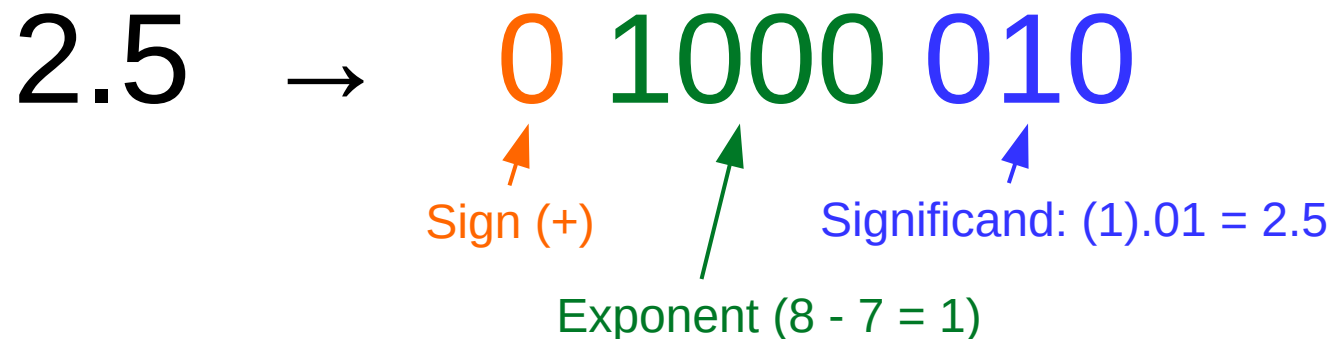
# Floating-point numbers

- Scientific notation to the rescue!
  - Traditionally, we write large (or small) numbers as  $x \cdot 10^e$
  - This is how **floating-point** representations work
    - Store **exponent** and fractional parts (the **significand**) separately
    - The decimal point “floats” on the number line
    - Position of point is based on the exponent

$$\begin{aligned} 1.23 &= 1.23 \times 10^0 = 0.123 \times 10^1 = 0.0123 \times 10^2 = \dots \\ &= 12.3 \times 10^{-1} = 123.0 \times 10^{-2} = \dots \end{aligned}$$

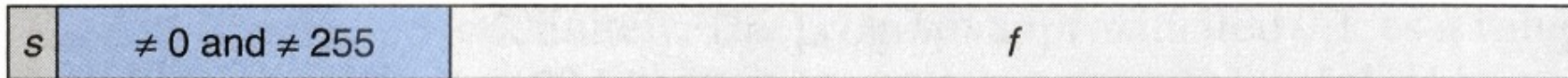
# Floating-point numbers

- However, computers use binary
  - So floating-point numbers use base 2 scientific notation ( $x \cdot 2^e$ )
- Fixed width field
  - Reserve one bit for the sign bit (0 is positive, 1 is negative)
  - Reserve n bits for **biased** exponent (bias is  $2^{n-1} - 1$ )
    - Avoids having to use two's complement
  - Use remaining bits for normalized fraction (implicit leading 1)
    - Exception: if the exponent is zero, don't normalize

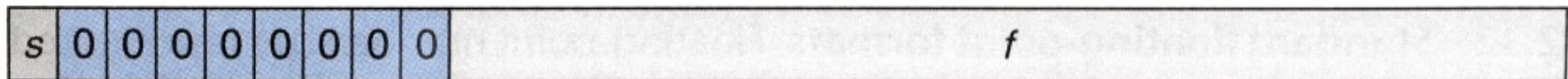


# Floating-point numbers

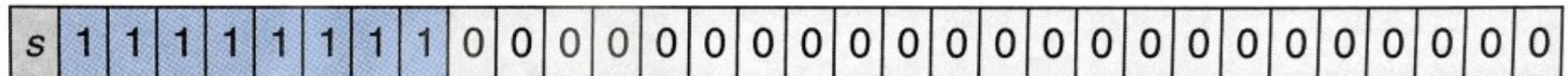
## 1. Normalized



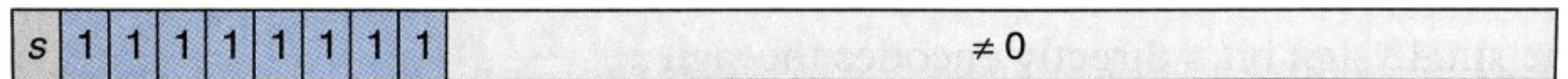
## 2. Denormalized



## 3a. Infinity



## 3b. NaN

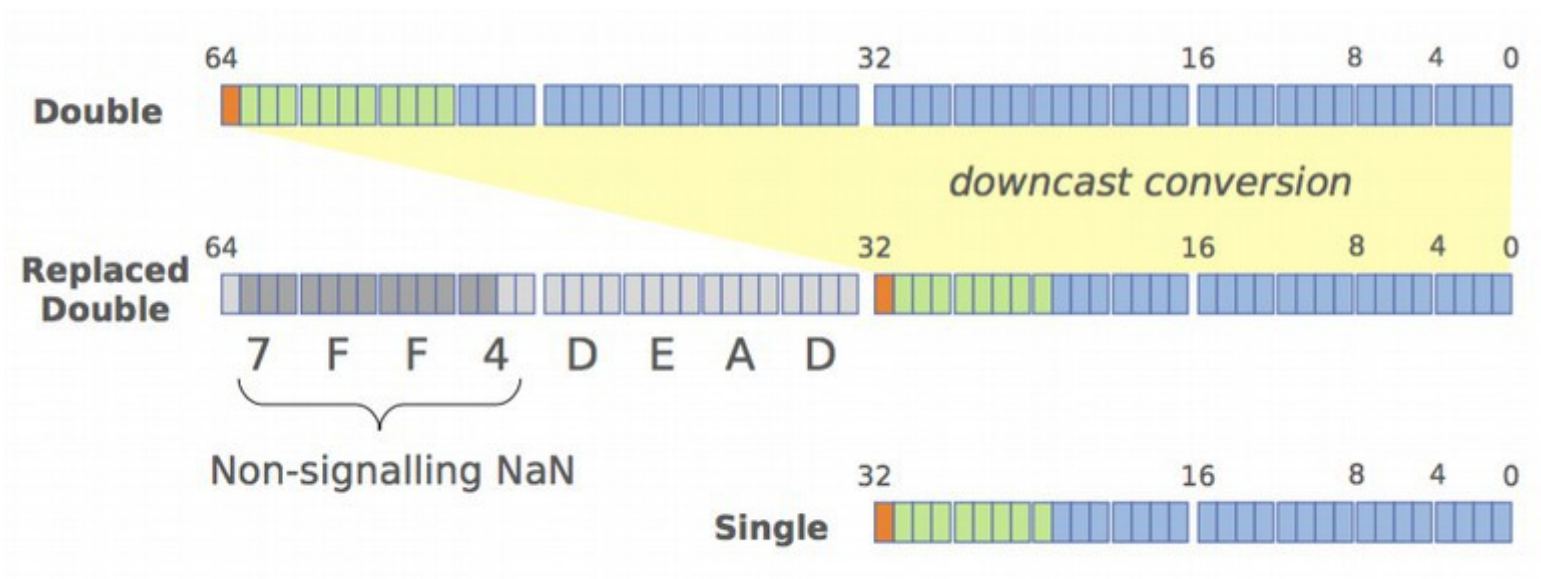


**Figure 2.33** Categories of single-precision floating-point values. The value of the exponent determines whether the number is (1) normalized, (2) denormalized, or (3) a special value.



# NaNs

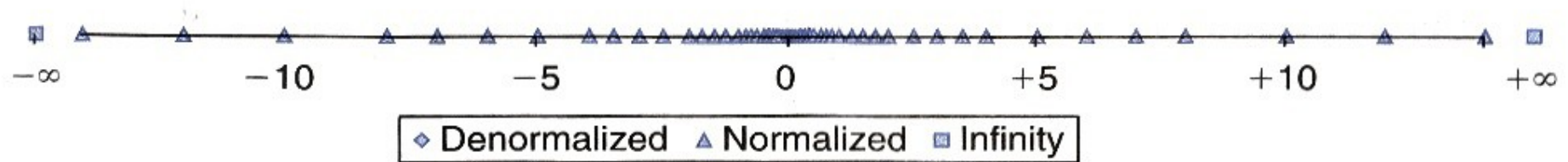
- NaN = “Not a Number”
  - Result of 0/0 and other undefined operations
  - Propagate to later calculations
  - Quiet and signaling variants (qNaN and sNaN)
  - Allowed a neat trick during my dissertation research:



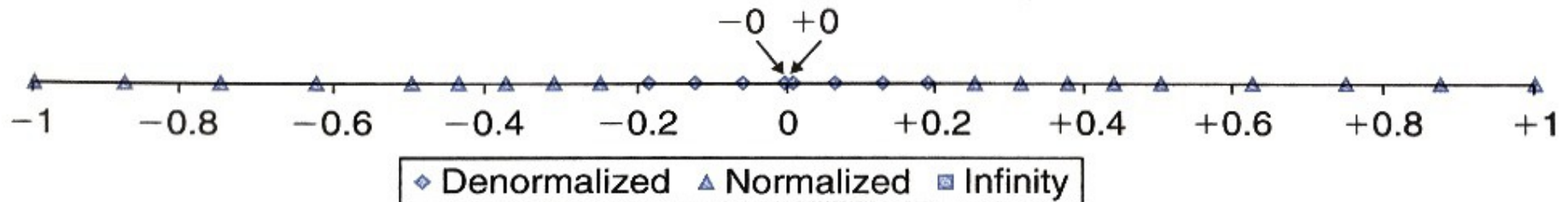
Description	Bit representation	Exponent			Fraction		Value		
		$e$	$E$	$2^E$	$f$	$M$	$2^E \times M$	$V$	Decimal
Zero	0 0000 000	0	-6	$\frac{1}{64}$	$\frac{0}{8}$	$\frac{0}{8}$	$\frac{0}{512}$	0	0.0
Smallest positive	0 0000 001	0	-6	$\frac{1}{64}$	$\frac{1}{8}$	$\frac{1}{8}$	$\frac{1}{512}$	$\frac{1}{512}$	0.001953
<i>gradual underflow near zero</i>	0 0000 010	0	-6	$\frac{1}{64}$	$\frac{2}{8}$	$\frac{2}{8}$	$\frac{2}{512}$	$\frac{1}{256}$	0.003906
	0 0000 011	0	-6	$\frac{1}{64}$	$\frac{3}{8}$	$\frac{3}{8}$	$\frac{3}{512}$	$\frac{3}{512}$	0.005859
	⋮								
Largest denormalized	0 0000 111	0	-6	$\frac{1}{64}$	$\frac{7}{8}$	$\frac{7}{8}$	$\frac{7}{512}$	$\frac{7}{512}$	0.013672
Smallest normalized	0 0001 000	1	-6	$\frac{1}{64}$	$\frac{0}{8}$	$\frac{8}{8}$	$\frac{8}{512}$	$\frac{1}{64}$	0.015625
<i>values &lt; 1</i>	0 0001 001	1	-6	$\frac{1}{64}$	$\frac{1}{8}$	$\frac{9}{8}$	$\frac{9}{512}$	$\frac{9}{512}$	0.017578
	⋮								
	0 0110 110	6	-1	$\frac{1}{2}$	$\frac{6}{8}$	$\frac{14}{8}$	$\frac{14}{16}$	$\frac{7}{8}$	0.875
	0 0110 111	6	-1	$\frac{1}{2}$	$\frac{7}{8}$	$\frac{15}{8}$	$\frac{15}{16}$	$\frac{15}{16}$	0.9375
One	0 0111 000	7	0	1	$\frac{0}{8}$	$\frac{8}{8}$	$\frac{8}{8}$	1	1.0
<i>values &gt; 1</i>	0 0111 001	7	0	1	$\frac{1}{8}$	$\frac{9}{8}$	$\frac{9}{8}$	$\frac{9}{8}$	1.125
	0 0111 010	7	0	1	$\frac{2}{8}$	$\frac{10}{8}$	$\frac{10}{8}$	$\frac{5}{4}$	1.25
	⋮								
	0 1110 110	14	7	128	$\frac{6}{8}$	$\frac{14}{8}$	$\frac{1792}{8}$	224	224.0
Largest normalized	0 1110 111	14	7	128	$\frac{7}{8}$	$\frac{15}{8}$	$\frac{1920}{8}$	240	240.0
Infinity	0 1111 000	—	—	—	—	—	—	$\infty$	—

**Figure 2.35** Example nonnegative values for 8-bit floating-point format. There are  $k = 4$  exponent bits and  $n = 3$  fraction bits. The bias is 7.

# Floating-point numbers



(a) Complete range



(b) Values between  $-1.0$  and  $+1.0$

**Figure 2.34** Representable values for 6-bit floating-point format. There are  $k = 3$  exponent bits and  $n = 2$  fraction bits. The bias is 3.

Not evenly spaced! (as integers are)

# Floating-point

- Some numbers cannot be represented exactly, regardless of how many bits are used!
  - E.g., 0.1 (dec)  $\rightarrow$  0.00011001100110011001100 ...
- This is no different than in base 10
  - E.g.,  $1/3 = 0.3333333333 \dots$

# Converting floating-point numbers

- Floating-point → decimal:

- 1) Sign bit: (value is “-1” if set, “1” if not)

- 2) Exponent:

- All zeroes: denormalized (exponent is 1-bias)

- All ones: NaN unless fraction is zero (which is infinity) – **DONE!**

- Otherwise: normalized (exponent is e-bias)

- 3) Fraction:

- If normalized:  $1 + f/2^{-m}$  (where m is the # of fraction bits)

- If denormalized:  $f/2^{-m}$  (where m is the # of fraction bits)

- Multiply *sign* x  $2^{exp}$  x *frac* to get final value

Note:

$$\text{bias} = 2^{n-1} - 1$$

(where n is the  
# of exp bits)

# Textbook's technique

$e$ : The value represented by considering the exponent field to be an unsigned integer

$E$ : The value of the exponent after biasing

$2^E$ : The numeric weight of the exponent

$f$ : The value of the fraction

$M$ : The value of the significand

$2^E \times M$ : The (unreduced) fractional value of the number

$V$ : The reduced fractional value of the number

Decimal: The decimal representation of the number

# Converting floating-point numbers

- Decimal → floating-point (normalized only)
  - Convert to fractional binary format
  - Normalize to 1.xxxxxx
    - Keep track of how many places you move the decimal and which direction
    - The “xxxxxx” bit string is the significand (pad with zeros or round if needed)
  - Encode resulting exponent
    - Add bias and convert to unsigned binary
    - If the exponent cannot be represented, result is zero or infinity

Note:  
bias =  $2^{n-1} - 1$   
(where n is the  
# of exp bits)

**Example (4-bit exp, 3-bit frac):**

2.75 (dec) → 10.11 (bin) → 1.011 x 2<sup>1</sup> (bin) → 0 1000 011

Bias =  $2^{4-1} - 1 = 7$

Exp: 1 + 7 = 8

# Example (textbook pg. 119)

$$12345_{10} \rightarrow 11000000111001_2$$

$$\rightarrow 1.1000000111001_2 \times 2^{13}$$

$$\text{exp} = 13 + 127 \text{ (bias)} = 140 = 10001100_2$$

$$\rightarrow 0 \ 10001100 \ 100000011100100000000000$$

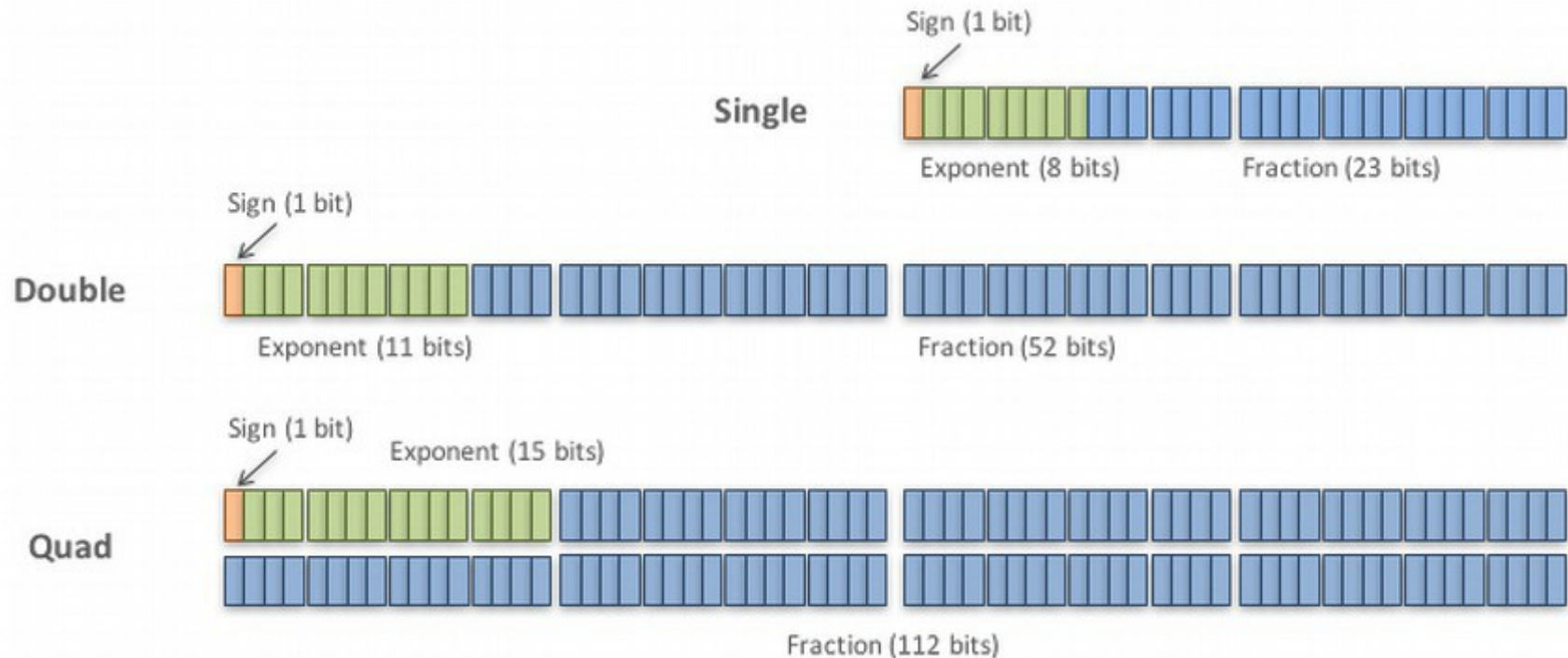
(note the shared bits that appear in all three representations)



# Floating-point numbers

## IEEE Floating-Point Numbers

Value is:  $(-1)^{\text{sign}} \times 1.\text{frac} \times 2^{\text{exp}}$



# Floating-point numbers

Name	Bits	Exp	Sig	Dec	M_Eps
IEEE half	16	5	10+1	3.311	9.77e-04
IEEE single	32	8	23+1	7.225	1.19e-07
IEEE double	64	11	52+1	15.955	2.22e-16
IEEE quad	128	15	112+1	34.016	1.93e-34

## NOTES:

- Sig is *<explicit>*[+*<implicit>*] bits
- Dec =  $\log_{10}(2^{\text{Sig}})$
- M\_Eps (machine epsilon) =  $b^{-(p-1)} = b^{(1-p)}$

# Conversion and rounding

To:

		Int32	Int64	Float	Double
<b>From:</b>	Int32	-	-	R	-
	Int64	O	-	R	R
	Float	OR	OR	-	-
	Double	OR	OR	OR	-

*O = overflow possible*  
*R = rounding possible*

*"-" is safe*

**10.00011** → 10.00

**10.00100** → 10.00

**10.00110** → 10.01

**10.10100** → 10.10

**10.01100** → 10.10

**10.11100** → 11.00

Mode	\$1.40	\$1.60	\$1.50	\$2.50	\$-1.50
Round-to-even	\$1	\$2	\$2	\$2	\$-2
Round-toward-zero	\$1	\$1	\$1	\$2	\$-1
Round-down	\$1	\$1	\$1	\$2	\$-2
Round-up	\$2	\$2	\$2	\$3	\$-1

Round-to-even: round to nearest,  
 on ties favor even numbers to  
 avoid statistical biases

**Figure 2.37** Illustration of rounding modes for dollar rounding. The first rounds to a nearest value, while the other three bound the result above or below.



# Floating-point issues

- Single vs. double precision choice
  - Theme: **system design involves tradeoffs**
  - Single precision arithmetic is faster
    - Especially on GPUs
  - Double precision is more accurate
    - More than twice as accurate!
  - Which do we use?
    - And how do we justify our choice?
    - Does the answer change for different regions of a program?
    - Does the answer change for different periods during execution?
    - **This is an open research question (talk to me if you're interested!)**

# Exercises

- What are the values of the following numbers, interpreted as floating-point numbers with a 3-bit exponent and 2-bit significand?
  - What about a 2-bit exponent and a 3-bit significand?

001100

011001

- Convert the following values to a floating-point value with a 4-bit exponent and a 3-bit significand. Write your answers in hex.

-3

0.125

120

$\infty$