

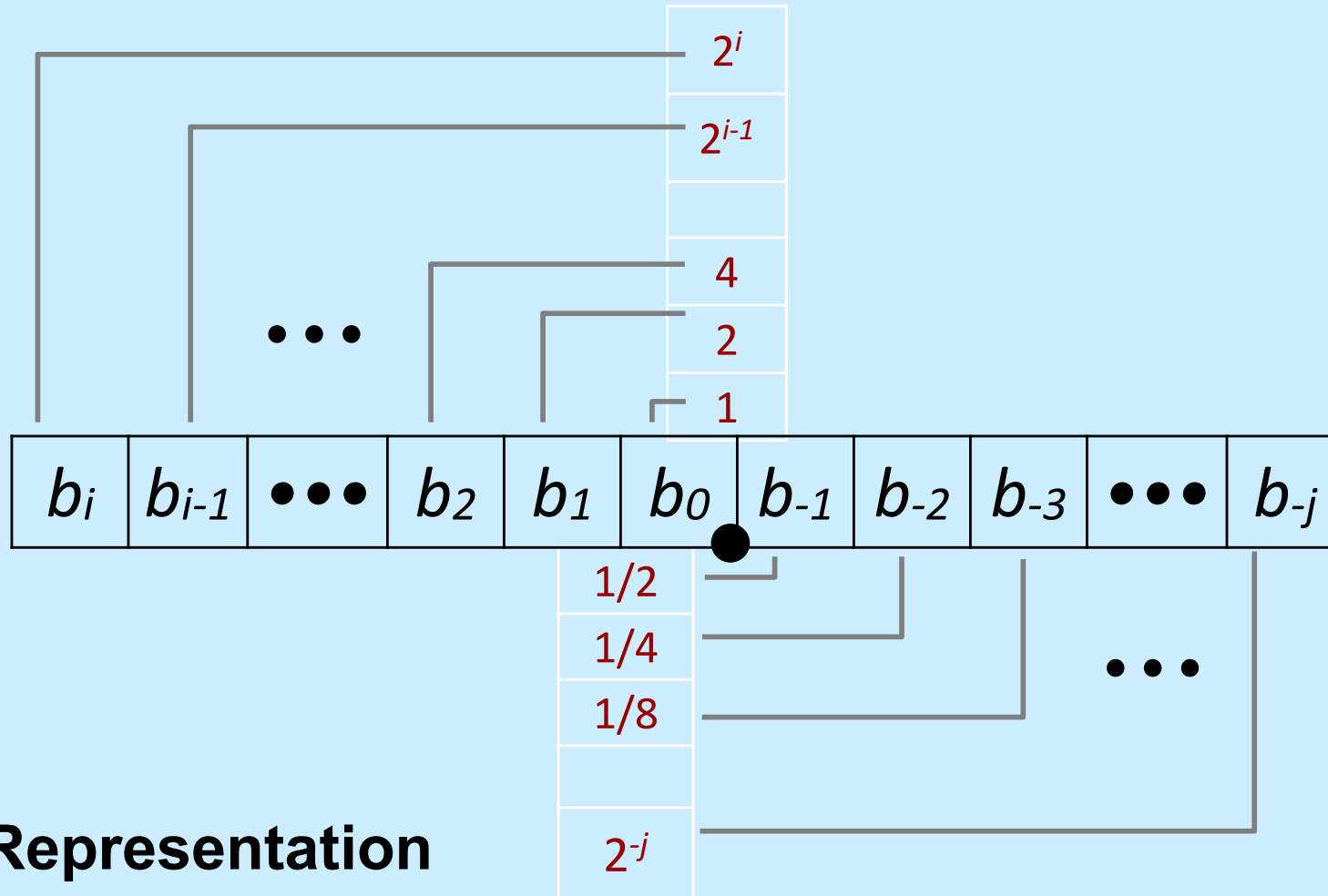
# CS 33

## Data Representation (Part 3)

# Fractional binary numbers

- What is  $1011.101_2$ ?

# Fractional Binary Numbers



- **Representation**

- bits to right of “binary point” represent fractional powers of 2
- represents rational number:

$$\sum_{k=-j}^i b_k \times 2^k$$

# Representable Numbers

- **Limitation #1**

- can exactly represent only numbers of the form  $n/2^k$

- » other rational numbers have repeating bit representations

- value      representation

- » 1/3      0.0101010101[01]...<sub>2</sub>

- » 1/5      0.001100110011[0011]...<sub>2</sub>

- » 1/10     0.0001100110011[0011]...<sub>2</sub>

- **Limitation #2**

- just one setting of decimal point within the  $w$  bits

- » limited range of numbers (very small values? very large?)

# IEEE Floating Point

- **IEEE Standard 754**
  - established in 1985 as uniform standard for floating point arithmetic
    - » before that, many idiosyncratic formats
  - supported by all major CPUs
- **Driven by numerical concerns**
  - nice standards for rounding, overflow, underflow
  - hard to make fast in hardware
    - » numerical analysts predominated over hardware designers in defining standard

# Floating-Point Representation

- Numerical Form:

$$(-1)^s M 2^E$$

- sign bit **s** determines whether number is negative or positive
- significand **M** normally a fractional value in range [1.0,2.0)
- exponent **E** weights value by power of two

- Encoding

- MSB **s** is sign bit **s**
- exp field encodes **E** (but is not equal to E)
- frac field encodes **M** (but is not equal to M)



# Precision options

- **Single precision: 32 bits**



- **Double precision: 64 bits**



- **Extended precision: 80 bits (Intel only)**



# “Normalized” Values

- When:  $\text{exp} \neq 000\dots 0$  and  $\text{exp} \neq 111\dots 1$
- Exponent coded as *biased* value:  $E = \text{Exp} - \text{Bias}$ 
  - *exp*: unsigned value  $\text{exp}$
  - $\text{bias} = 2^{k-1} - 1$ , where  $k$  is number of exponent bits
    - » single precision: 127 (Exp: 1...254, E: -126...127)
    - » double precision: 1023 (Exp: 1...2046, E: -1022...1023)
- Significand coded with implied leading 1:  $M = 1.\text{xxx}\dots\text{x}_2$ 
  - $\text{xxx}\dots\text{x}$ : bits of *frac*
  - minimum when  $\text{frac}=000\dots 0$  ( $M = 1.0$ )
  - maximum when  $\text{frac}=111\dots 1$  ( $M = 2.0 - \epsilon$ )
  - get extra leading bit for “free”



# Normalized Encoding Example

- **Value:** `float F = 15213.0;`

$$\begin{aligned} - 15213_{10} &= 11101101101101_2 \\ &= 1.1101101101101_2 \times 2^{13} \end{aligned}$$

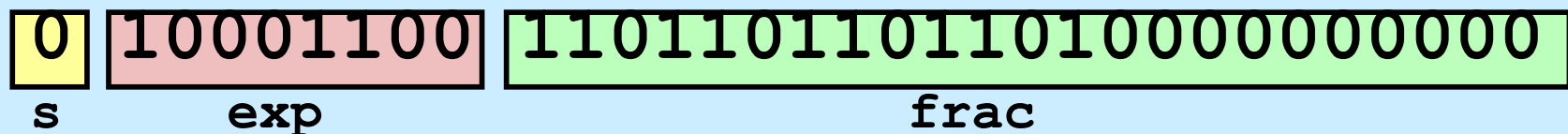
- **Significand**

$$\begin{aligned} M &= 1.\underline{1101101101101}_2 \\ \text{frac} &= \underline{1101101101101}0000000000_2 \end{aligned}$$

- **Exponent**

$$\begin{aligned} E &= 13 \\ \text{bias} &= 127 \\ \text{exp} &= 140 = 10001100_2 \end{aligned}$$

- **Result:**



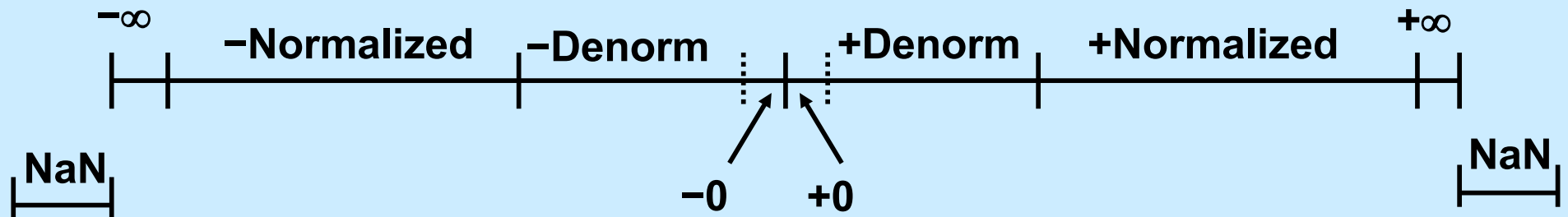
# Denormalized Values

- Condition:  $\text{exp} = 000\dots 0$
- Exponent value:  $E = -\text{Bias} + 1$  (instead of  $E = 0 - \text{Bias}$ )
- Significand coded with implied leading 0:  
 $M = 0.\text{xxx}\dots\text{x}_2$ 
  - $\text{xxx}\dots\text{x}$ : bits of  $\text{frac}$
- Cases
  - $\text{exp} = 000\dots 0, \text{frac} = 000\dots 0$ 
    - » represents zero value
    - » note distinct values:  $+0$  and  $-0$  (why?)
  - $\text{exp} = 000\dots 0, \text{frac} \neq 000\dots 0$ 
    - » numbers closest to  $0.0$
    - » equispaced

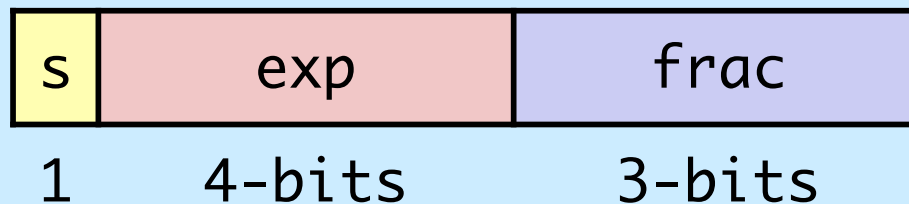
# Special Values

- **Condition:  $\text{exp} = 111\dots 1$**
  - **Case:  $\text{exp} = 111\dots 1$ ,  $\text{frac} = 000\dots 0$** 
    - represents value  $\infty$  (infinity)
    - operation that overflows
    - both positive and negative
    - e.g.,  $1.0/0.0 = -1.0/-0.0 = +\infty$ ,  $1.0/-0.0 = -\infty$
  - **Case:  $\text{exp} = 111\dots 1$ ,  $\text{frac} \neq 000\dots 0$** 
    - not-a-number (NaN)
    - represents case when no numeric value can be determined
    - e.g.,  $\text{sqrt}(-1)$ ,  $\infty - \infty$ ,  $\infty \times 0$
-

# Visualization: Floating-Point Encodings



# Tiny Floating-Point Example



- **8-bit Floating Point Representation**
  - the sign bit is in the most significant bit
  - the next four bits are the exponent, with a bias of 7
  - the last three bits are the *frac*
- **Same general form as IEEE Format**
  - normalized, denormalized
  - representation of 0, NaN, infinity

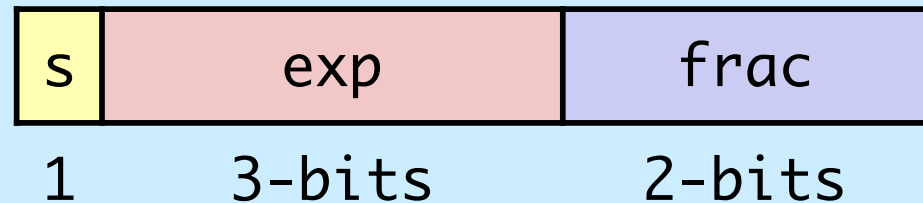
# Dynamic Range (Positive Only)

	s	exp	frac	E	Value	
Denormalized numbers	0	0000	000	-6	0	closest to zero
	0	0000	001	-6	$1/8 * 1/64 = 1/512$	
	0	0000	010	-6	$2/8 * 1/64 = 2/512$	
	...					
	0	0000	110	-6	$6/8 * 1/64 = 6/512$	largest denorm
	0	0000	111	-6	$7/8 * 1/64 = 7/512$	
Normalized numbers	0	0001	000	-6	$8/8 * 1/64 = 8/512$	smallest norm
	0	0001	001	-6	$9/8 * 1/64 = 9/512$	
	...					
	0	0110	110	-1	$14/8 * 1/2 = 14/16$	closest to 1 below
	0	0110	111	-1	$15/8 * 1/2 = 15/16$	
	0	0111	000	0	$8/8 * 1 = 1$	closest to 1 above
	0	0111	001	0	$9/8 * 1 = 9/8$	
	0	0111	010	0	$10/8 * 1 = 10/8$	
	...					
	0	1110	110	7	$14/8 * 128 = 224$	largest norm
	0	1110	111	7	$15/8 * 128 = 240$	
	0	1111	000	n/a	inf	

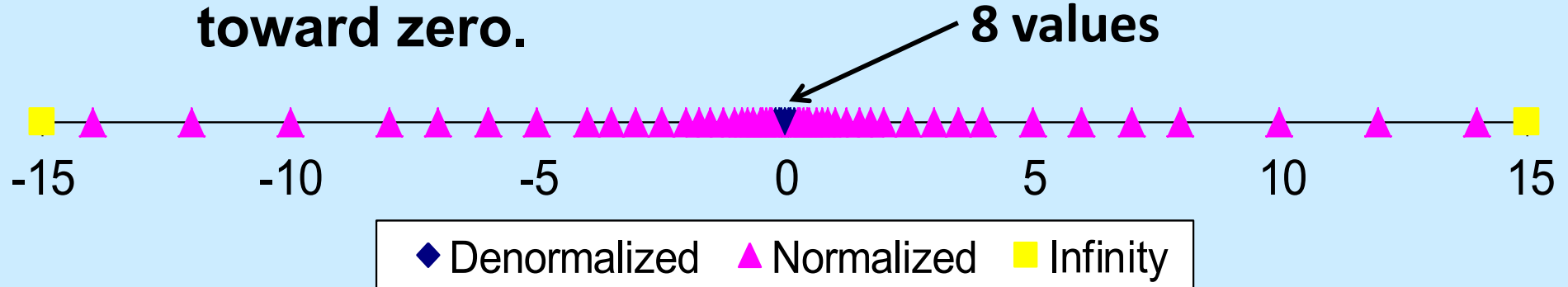
# Distribution of Values

- 6-bit IEEE-like format

- e = 3 exponent bits
- f = 2 fraction bits
- bias is  $2^{3-1}-1 = 3$



- Notice how the distribution gets denser toward zero.



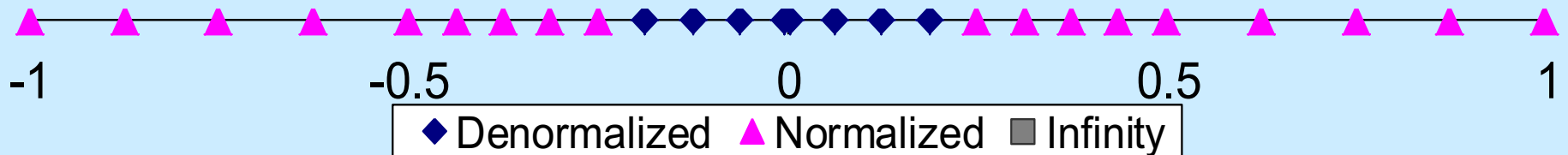
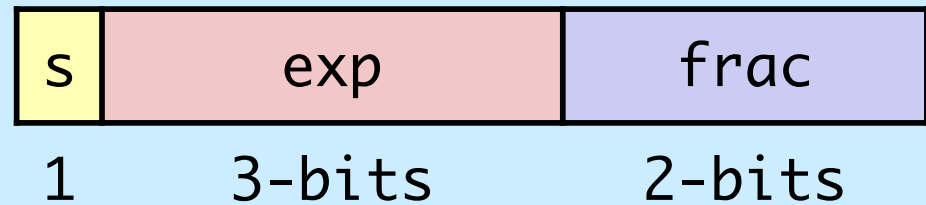
# Distribution of Values (close-up view)

- 6-bit IEEE-like format

- e = 3 exponent bits

- f = 2 fraction bits

- bias is 3

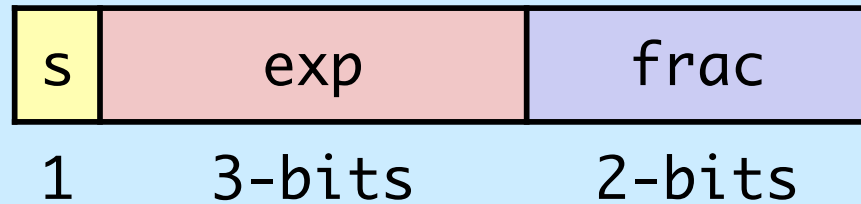




# Quiz 1

- **6-bit IEEE-like format**

- e = 3 exponent bits
- f = 2 fraction bits
- bias is 3



**What number is represented by 0 011 10?**

- a) 12
- b) 1.5
- c) .5
- d) none of the above

# Floating-Point Operations: Basic Idea

- $x +_f y = \text{Round}(x + y)$
- $x \times_f y = \text{Round}(x \times y)$
- Basic idea
  - first **compute exact result**
  - make it fit into desired precision
    - » possibly overflow if exponent too large
    - » possibly **round to fit into** frac

# Rounding

- Rounding modes (illustrated with \$ rounding)

	\$1.40	\$1.60	\$1.50	\$2.50	-\$1.50
towards zero	\$1	\$1	\$1	\$2	-\$1
round down ( $-\infty$ )	\$1	\$1	\$1	\$2	-\$2
round up ( $+\infty$ )	\$2	\$2	\$2	\$3	-\$1
nearest integer	\$1	\$2	?	?	?
nearest even (default)	\$1	\$2	\$2	\$2	-\$2

# Creating a Floating Point Number



- **Steps**

- normalize to have leading 1
- round to fit within fraction
- postnormalize to deal with effects of rounding

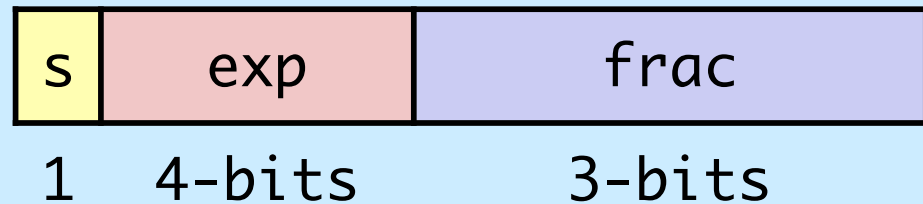
- **Case study**

- convert 8-bit unsigned numbers to tiny floating-point format

**example numbers**

128	10000000
13	00001101
33	00010001
35	00010011
138	10001010
63	00111111

# Normalize



- **Requirement**
  - set binary point so that numbers of form 1.xxxxx
  - adjust all to have leading one
    - » decrement exponent as shift left

<i>Value</i>	<i>Binary</i>	<i>Fraction</i>	<i>Exponent</i>
128	10000000	1.0000000	7
13	00001101	1.1010000	3
17	00010001	1.0001000	4
19	00010011	1.0011000	4
138	10001010	1.0001010	7
63	00111111	1.1111100	5

# Rounding

1.BBG**RXXX**

Guard bit: LSB of result

Round bit: 1<sup>st</sup> bit removed

Sticky bit: OR of remaining bits

- Round-up conditions

- round = 1, sticky = 1  $\Rightarrow > 0.5$

- guard = 1, round = 1, sticky = 0  $\Rightarrow$  round to even

<i>Value</i>	<i>Fraction</i>	<i>GRS</i>	<i>Incr?</i>	<i>Rounded</i>
128	1.000 <b>0000</b>	000	N	1.000
13	1.101 <b>0000</b>	100	N	1.101
17	1.000 <b>1000</b>	010	N	1.000
19	1.001 <b>1000</b>	110	Y	1.010
138	1.000 <b>1010</b>	011	Y	1.001
63	1.111 <b>1100</b>	111	Y	10.000

# Postnormalize

- **Issue**
  - rounding may have caused overflow
  - handle by shifting right once & incrementing exponent

<i>Value</i>	<i>Rounded</i>	<i>Exp</i>	<i>Adjusted</i>	<i>Result</i>
128	1.000	7		128
13	1.101	3		13
17	1.000	4		16
19	1.010	4		20
138	1.001	7		134
63	10.000	5	$1.000 * 2^6$	64

# Floating-Point Multiplication

- $(-1)^{s1} M1 2^{E1} \times (-1)^{s2} M2 2^{E2}$
- **Exact result:**  $(-1)^s M 2^E$ 
  - sign  $s$ :  $s1 \wedge s2$
  - significand  $M$ :  $M1 \times M2$
  - exponent  $E$ :  $E1 + E2$
- **Fixing**
  - if  $M \geq 2$ , shift  $M$  right, increment  $E$
  - if  $E$  out of range, overflow (or underflow)
  - round  $M$  to fit `frac` precision
- **Implementation**
  - biggest chore is multiplying significands



# Floating-Point Addition

- $(-1)^{s1} M1 2^{E1} + (-1)^{s2} M2 2^{E2}$

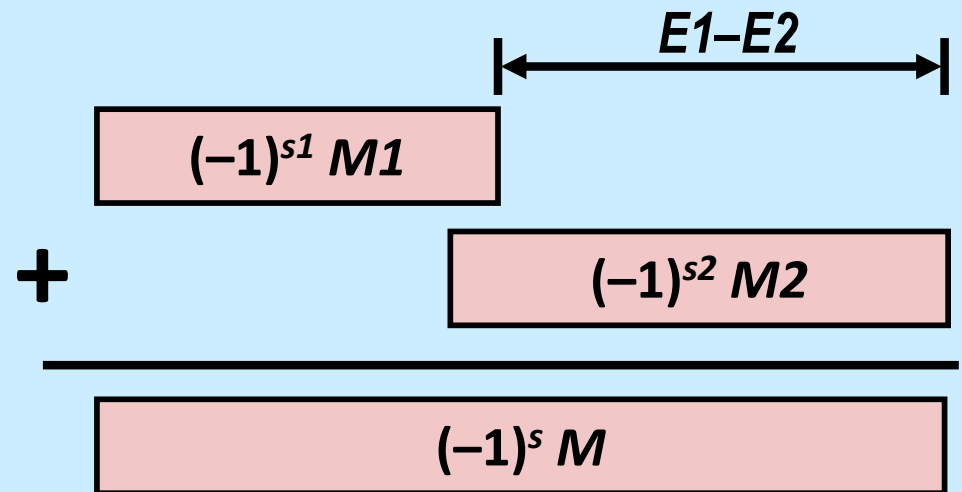
–assume  $E1 > E2$

- **Exact result:**  $(-1)^s M 2^E$

–sign  $s$ , significand  $M$ :

» result of signed align & add

–exponent  $E$ :  $E1$



- **Fixing**

–if  $M \geq 2$ , shift  $M$  right, increment  $E$

–if  $M < 1$ , shift  $M$  left  $k$  positions, decrement  $E$  by  $k$

–overflow if  $E$  out of range

–round  $M$  to fit **frac** precision

# Floating Point in C

- **C guarantees two levels**
    - `float`      single precision
    - `double`     double precision
  - **Conversions/casting**
    - casting between `int`, `float`, and `double` changes bit representation
    - `double/float`  $\rightarrow$  `int`
      - » truncates fractional part
      - » like rounding toward zero
      - » not defined when out of range or NaN: generally sets to TMin
    - `int`  $\rightarrow$  `double`
      - » exact conversion, as long as `int` has  $\leq 53$ -bit word size
    - `int`  $\rightarrow$  `float`
      - » will round according to rounding mode
-

# Quiz 2

Suppose  $f$ , declared to be a `float`, is assigned the largest possible floating-point positive value (other than  $+\infty$ ). What is the value of  $g = f + 1.0$ ?

- a)  $f$
- b)  $+\infty$
- c) NAN
- d) 0

# Float is not Rational ...

- **Floating addition**
  - **commutative:  $a +^f b = b +^f a$** 
    - » **yes!**
  - **associative:  $a +^f (b +^f c) = (a +^f b) +^f c$** 
    - » **no!**
      - **$2 +^f (1e20 +^f -1e20) = 2$**
      - **$(2 +^f 1e20) +^f -1e20 = 0$**

# Float is not Rational ...

- **Multiplication**

- **commutative:  $a *^f b = b *^f a$**

- » **yes!**

- **associative:  $a *^f (b *^f c) = (a *^f b) *^f c$**

- » **no!**

- **$1e20 *^f (1e20 *^f 1e-20) = 1e20$**

- **$(1e20 *^f 1e20) *^f 1e-20 = +\infty$**

# Float is not Rational ...

- **More ...**
  - **multiplication distributes over addition:**  
$$a *^f (b +^f c) = (a *^f b) +^f (a *^f c)$$
    - » no!
    - »  $1e20 *^f (1e20 +^f -1e20) = 0$
    - »  $(1e20 *^f 1e20) +^f (1e20 *^f -1e20) = \text{NaN}$
  - **loss of significance:**  
$$x = y + 1$$
$$z = 2 / (x - y)$$
$$z == 2?$$
    - » not necessarily!
      - consider  $y = 1e20$