

Data Representation – Floating Point

CSCI 2400 / ECE 3217: Computer Architecture

Instructor:

David Ferry

*Slides adapted from Bryant & O'Hallaron's slides
via Jason Fritts*

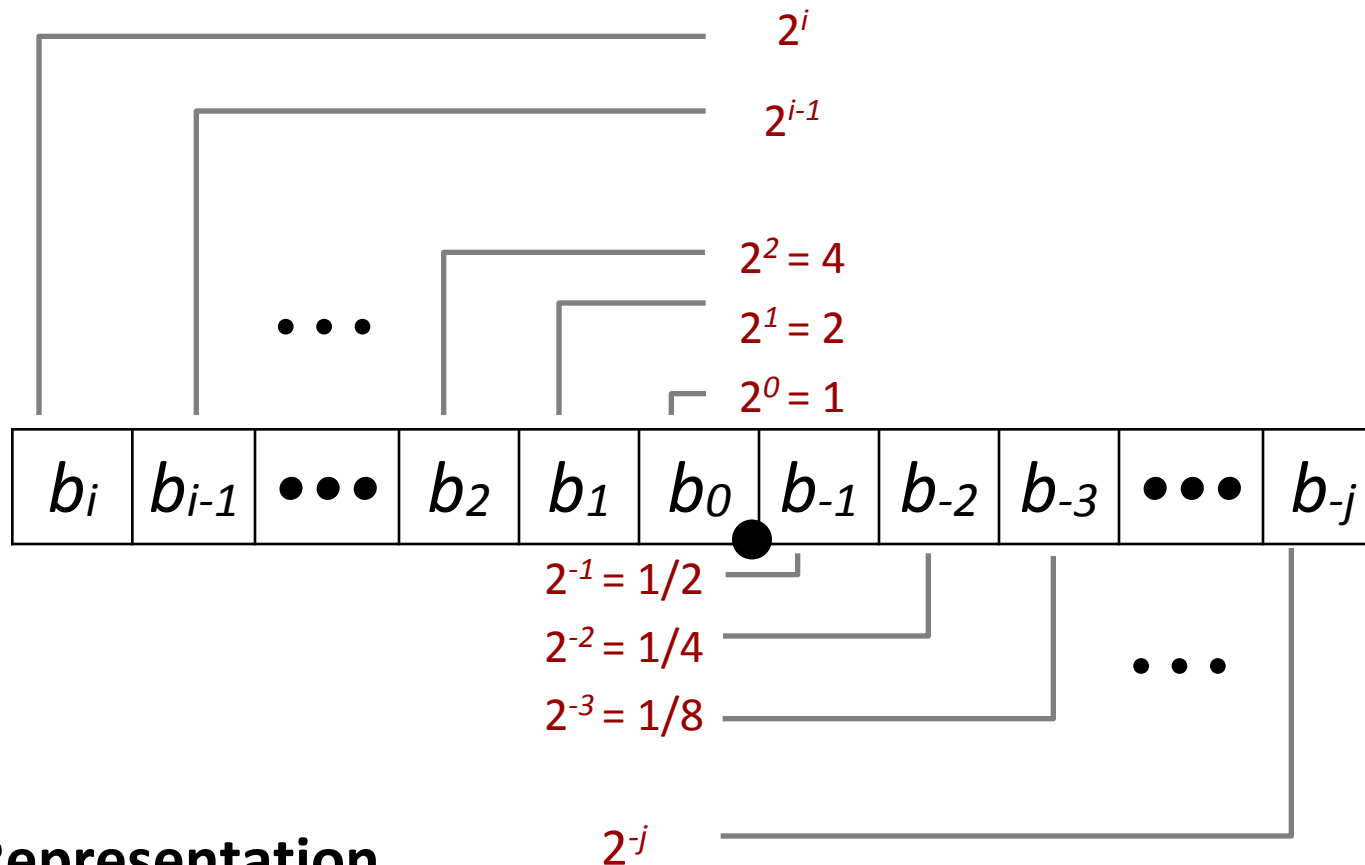
Today: Floating Point

- **Background: Fractional binary numbers**
- **Example and properties**
- **IEEE floating point standard: Definition**
- **Floating point in C**
- **Summary**

Fractional binary numbers

- What is 1011.101_2 ?
- How can we express fractions like $\frac{1}{4}$ in binary?

Place-Value Fractional Binary Numbers



■ Representation

- Bits to right of “binary point” represent fractional powers of 2
- Represents rational number:

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

Fractional Binary Numbers: Examples

Value	Representation	
$5 \frac{3}{4}$	101.11_2	$= 4 + 1 + \frac{1}{2} + \frac{1}{4} = 5 \frac{3}{4}$
$2 \frac{7}{8}$	10.111_2	$= 2 + \frac{1}{2} + \frac{1}{4} + \frac{1}{8} = 2 \frac{7}{8}$
$\frac{25}{64}$	0.011001_2	$= \frac{1}{4} + \frac{1}{8} + \frac{1}{64} = \frac{25}{64}$

Observations

- Divide by 2 by shifting right
- Multiply by 2 by shifting left

Limitations

- Can only exactly represent numbers of the form $x/2^k$
- Other rational numbers have repeating bit representations

<u>Value</u>	<u>Representation</u>
$1/3$	$0.0101010101[01]..._2$
$1/5$	$0.001100110011[0011]..._2$

- Limited range when used with “fixed point” representations

Quick Check

■ Convert:

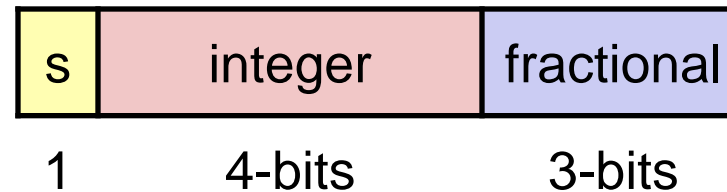
$255 \frac{9}{16}$ to binary

10101.10101_2 to decimal

Quick Check

Suppose an 8-bit *fixed-point* representation with:

- One sign bit
- Four integer bits
- Three fractional bits



Convert:

$12 \frac{1}{8}$ to binary

$-6 \frac{3}{8}$ to binary

11010110_2 to decimal

What bit pattern(s) have the largest positive value? What is it?

What bit pattern(s) have the value closest to zero?

What bit pattern(s) have the value of zero?

Today: Floating Point

- Background: Fractional binary numbers
- **Example and properties**
- IEEE floating point standard: Definition
- Floating point in C
- Summary

Floating Point Representation

- Similar to scientific notation

E.g.: $1.25 \times 10^3 = 1,250$

E.g.: $2.78 \times 10^{-2} = 0.0278$

- **FP is this concept but with an efficient binary format! But...**

- Uses base 2 instead of base 10
- Places restrictions on how certain values are represented
- Deals with finiteness of representation

Floating Point Representation

■ Numerical Form:

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

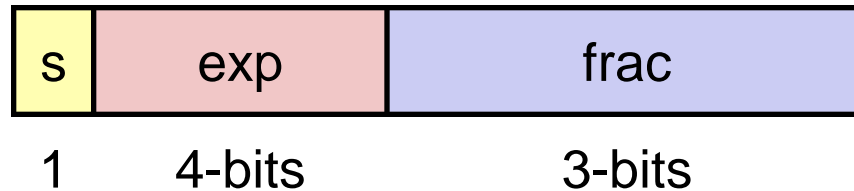
- **Sign bit** s determines whether number is negative or positive
- **Significand** (mantissa) M normally a fractional value in range $[1.0, 2.0)$
- **Exponent** E weights value by power of two

■ Encoding

- 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*)



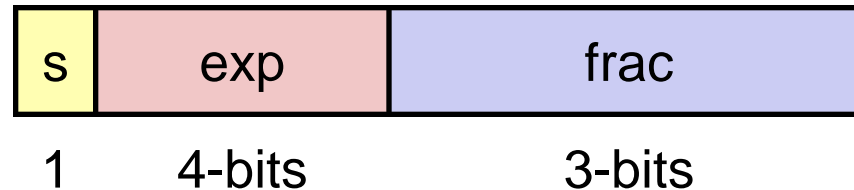
Tiny Floating Point Example



■ 8-bit Floating Point Representation

- the sign bit is in the most significant bit
- the next four bits are the exponent (**exp**)
 - **exp** (not E) encoded as a 4-bit unsigned integer
 - Uses a *bias* to represent negative exponents
- the last three bits are the fraction (**frac**)
 - encodes fractional part of a fractional binary number

Tiny Floating Point Example



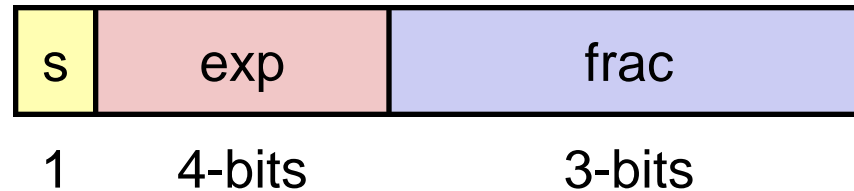
■ Fractional (Mantissa) Value

- Three bits- encoded left-to-right as $b_0b_1b_2$
- Place value of b_0 is $1/2$, of b_1 is $1/4$, and b_2 is $1/8$
- Value usually includes an implied 1:

$$M = 1 + \frac{1}{2}b_0 + \frac{1}{4}b_1 + \frac{1}{8}b_2$$

- If **exp** == 0000 then the implied 1 is not used

Tiny Floating Point Example



■ Exponent bias

- enable exponent to represent both positive and negative powers of 2
- use half of range for positive and half for negative power
- given k exponent bits, bias is then $2^{k-1} - 1$

■ Exponent Value

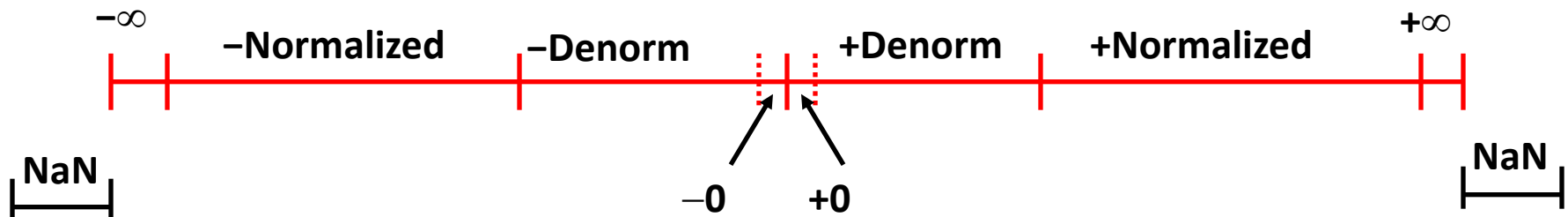
- is usually computed $E = \text{exp} - \text{bias}$
- if $\text{exp} == 0000$ then $E = 1 - \text{bias}$

Floating Point Encodings and Visualization

■ Five encodings:

- Two general forms: normalized, denormalized
- Three special values: zero, infinity, NaN (*not a number*)

<u>Name</u>	<u>Exponent (exp)</u>	<u>Fraction (frac)</u>
zero	exp == 0000	frac == 000
denormalized	exp == 0000	frac != 000
normalized	0000 < exp < 1111	frac != 000
infinity	exp == 1111	frac == 000
NaN	exp == 1111	frac != 000



Dynamic Range (Positives)

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

norm: $E = \text{Exp} - \text{Bias}$
denorm: $E = 1 - \text{Bias}$

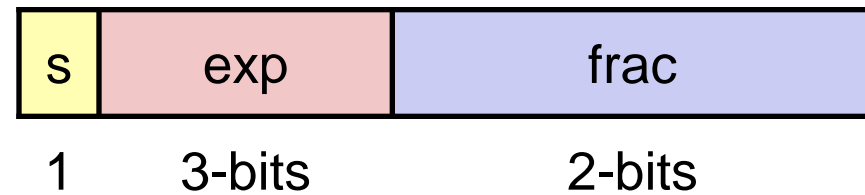
	s	exp	frac	E	Value	
Denormalized numbers	0	0000	000	-6	0	
	0	0000	001	-6	$1/8 * 1/64 = 1/512$	closest to zero
	0	0000	010	-6	$2/8 * 1/64 = 2/512$	
	...					
	0	0000	110	-6	$6/8 * 1/64 = 6/512$	
	0	0000	111	-6	$7/8 * 1/64 = 7/512$	largest denorm
Normalized numbers	0	0001	000	-6	$8/8 * 1/64 = 8/512$	smallest norm
	0	0001	001	-6	$9/8 * 1/64 = 9/512$	
	...					
	0	0110	110	-1	$14/8 * 1/2 = 14/16$	
	0	0110	111	-1	$15/8 * 1/2 = 15/16$	closest to 1 below
	0	0111	000	0	$8/8 * 1 = 1$	
	0	0111	001	0	$9/8 * 1 = 9/8$	closest to 1 above
	0	0111	010	0	$10/8 * 1 = 10/8$	
	...					
	0	1110	110	7	$14/8 * 128 = 224$	
	0	1110	111	7	$15/8 * 128 = 240$	largest norm
	0	1111	000	n/a	inf	infinity
	0	1111	xxx	n/a	NaN	NaN (not a number)

Distribution of Values

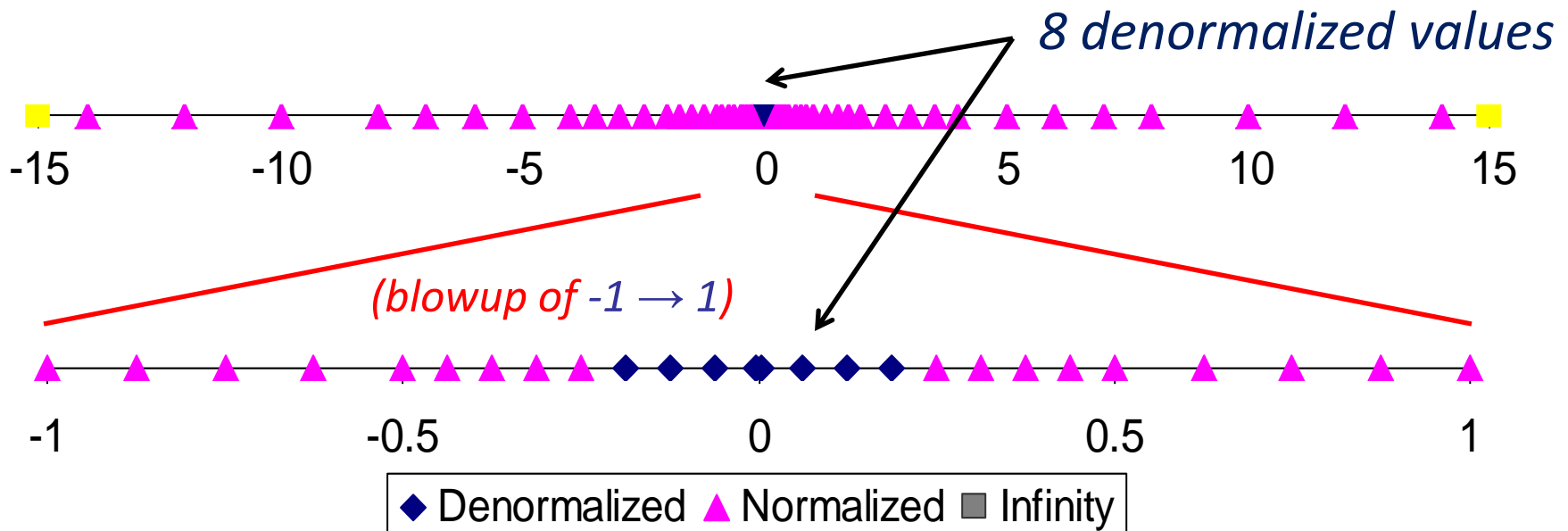
(reduced format from 8 bits to 6 bits for visualization)

■ 6-bit IEEE-like format

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



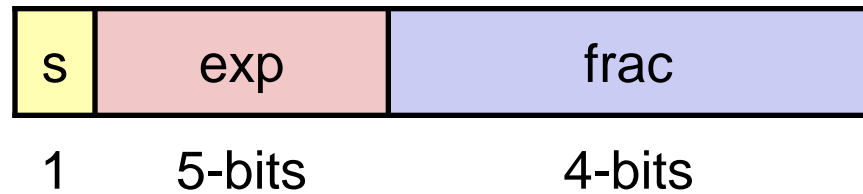
■ Notice how the distribution gets denser toward zero.



Quick Check

■ 10-bit IEEE-like format

- $e = 5$ exponent bits
- $f = 4$ fraction bits

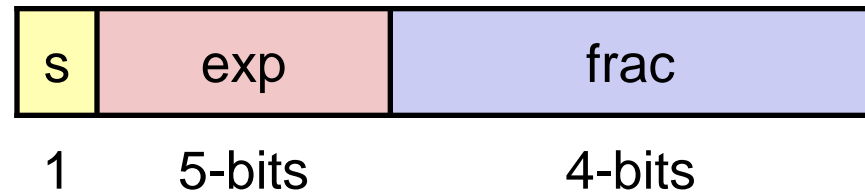


■ What is the exponent bias?

Quick Check

■ 10-bit IEEE-like format

- $e = 5$ exponent bits
- $f = 4$ fraction bits

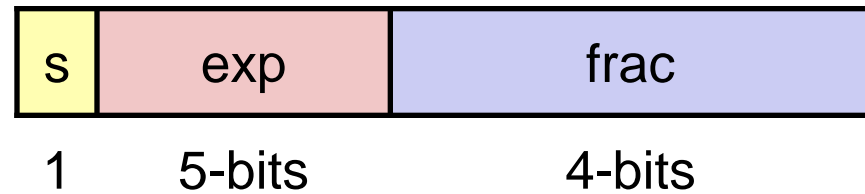


- What is the exponent bias? $2^{5-1} - 1 = 15$

Quick Check

■ 10-bit IEEE-like format

- $e = 5$ exponent bits
- $f = 4$ fraction bits

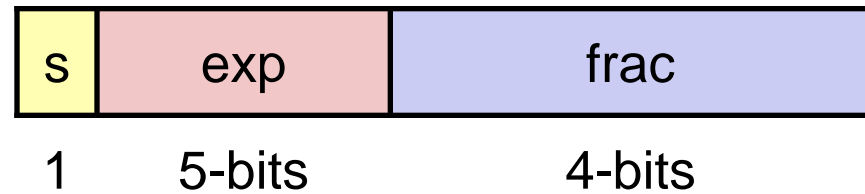


- What is the exponent bias? $2^{5-1} - 1 = 15$
- How many denormalized numbers are there?

Quick Check

■ 10-bit IEEE-like format

- $e = 5$ exponent bits
- $f = 4$ fraction bits

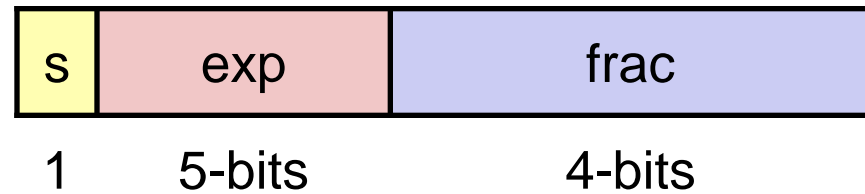


- What is the exponent bias? $2^{5-1} - 1 = 15$
- How many denormalized numbers are there?
 - Exponent = 00000, so 2^4 positive and 2^4 negative

Quick Check

■ 10-bit IEEE-like format

- $e = 5$ exponent bits
- $f = 4$ fraction bits

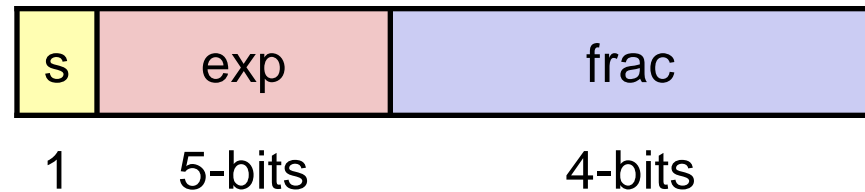


- What is the exponent bias? $2^{5-1} - 1 = 15$
- How many denormalized numbers are there?
 - Exponent = 00000, so 2^4 positive and 2^4 negative
- What is the bit pattern of the maximum value number?
- What is the bit pattern of the number closest to zero?

Quick Check

■ 10-bit IEEE-like format

- $e = 5$ exponent bits
- $f = 4$ fraction bits

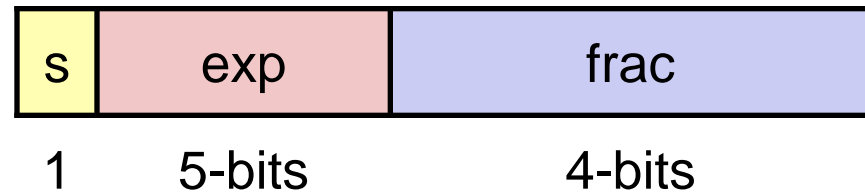


- What is the exponent bias? $2^{5-1} - 1 = 15$
- How many denormalized numbers are there?
 - Exponent = 00000, so 2^4 positive and 2^4 negative
- What is the bit pattern of the maximum value number?
 - Sign = 0, Exponent = 11110, frac=1111, so 0111101111
- What is the bit pattern of the number closest to zero?
 - Sign = ?, Exponent = 00000, frac=0001, so ?000000001

Quick Check

■ 10-bit IEEE-like format

- $e = 5$ exponent bits
- $f = 4$ fraction bits

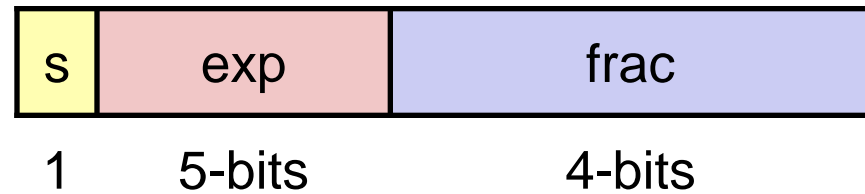


- ## ■ What is the bit pattern of the smallest positive normal number?

Quick Check

■ 10-bit IEEE-like format

- $e = 5$ exponent bits
- $f = 4$ fraction bits



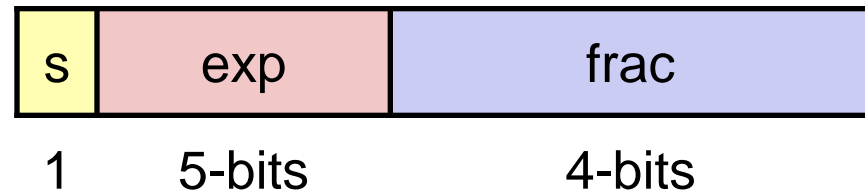
■ What is the bit pattern of the smallest positive normal number?

- Sign = 0, exp = 00001, frac = 0000; so 0000010000

Quick Check

■ 10-bit IEEE-like format

- $e = 5$ exponent bits
- $f = 4$ fraction bits



■ What is the bit pattern of the smallest positive normal number?

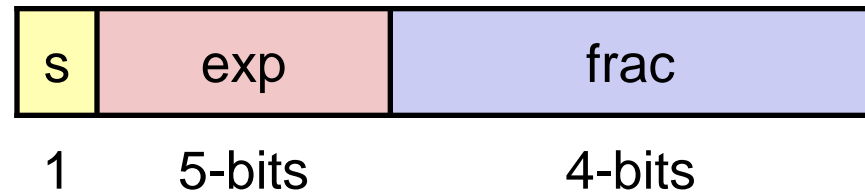
- Sign = 0, exp = 00001, frac = 0000; so 0000010000

■ What is the value of the smallest positive normal number?

Quick Check

■ 10-bit IEEE-like format

- $e = 5$ exponent bits
- $f = 4$ fraction bits



■ What is the bit pattern of the smallest positive normal number?

- Sign = 0, exp = 00001, frac = 0000; so 0000010000

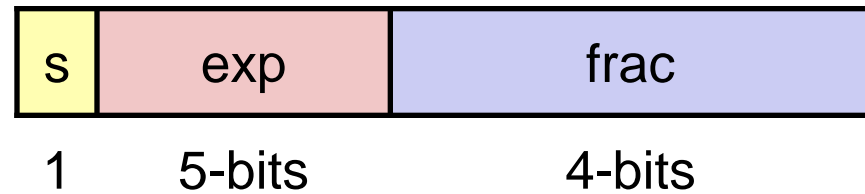
■ What is the value of the smallest positive normal number?

- Value = $(-1)^S \times M \times 2^E$
- $S = 0$
- Exponent bias = 15, so $E = 1 - 15 = -14$
- $M = 1 + 0 \times \frac{1}{2} + 0 \times \frac{1}{4} + 0 \times \frac{1}{8} + 0 \times \frac{1}{16} = 1$
- Value = $(-1)^0 \times 1 \times 2^{-14} = 1/2^{14} = 0.00006103515$

Quick Check

■ 10-bit IEEE-like format

- $e = 5$ exponent bits
- $f = 4$ fraction bits

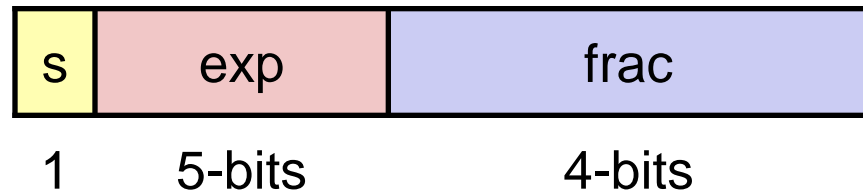


- Given a 32-bit floating point number, and a 32-bit integer, which can represent more discrete values?

Quick Check

■ 10-bit IEEE-like format

- $e = 5$ exponent bits
- $f = 4$ fraction bits



■ Given a 32-bit floating point number, and a 32-bit integer, which can represent more discrete values?

- Both can represent 2^{32} values, but some bit patterns duplicate values, e.g. $+0/-0$, $+\infty/-\infty$, and many NaNs (exponent = 11...1, frac \neq 00...0)

Today: Floating Point

- Background: Fractional binary numbers
- Example and properties
- **IEEE floating point standard: Definition**
- Floating point in C
- Summary

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** (mantissa) M normally a fractional value in range $[1.0, 2.0)$
- **Exponent** E weights value by power of two

■ Encoding

- 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*)



Standard Types

■ Single precision: 32 bits (c type: `float`)



■ Double precision: 64 bits (c type: `double`)



■ Extended precision: 80 bits (Intel only)



Normalized Values

- **Condition:** $exp \neq 000\dots 0$ and $exp \neq 111\dots 1$
- **Exponent coded as *biased* value:** $E = Exp - Bias$
 - Exp : unsigned value of exp field
 - $Bias = 2^{k-1} - 1$, where k is number of exponent bits
 - Single precision: **127** ($exp: 1\dots 254 \Rightarrow E: -126\dots 127$)
 - Double precision: **1023** ($exp: 1\dots 2046 \Rightarrow E: -1022\dots 1023$)
- **Significand coded with implied leading 1:** $M = \underline{1}.xxx\dots x_2$
 - $xxx\dots x$: bits of $frac$
- **Decimal value of normalized FP representations:**
 - Single-precision: $Value_{10} = (-1)^s \times 1.frac \times 2^{exp-127}$
 - Double-precision: $Value_{10} = (-1)^s \times 1.frac \times 2^{exp-1023}$

Normalized Encoding Example

■ Value: float $F = 15213.5$;

$$15213.5_{10} = 11101101101101.1_2$$

$$= 1.11011011011011_2 \times 2^{13}$$

shift binary point by K bits so that only one leading 1 bit remains on the left side of the binary point (here, shifted right by 13 bits, so $K = 13$), then multiply by 2^K (here, 2^{13})

■ Significand

$$M = 1.11011011011011_2$$

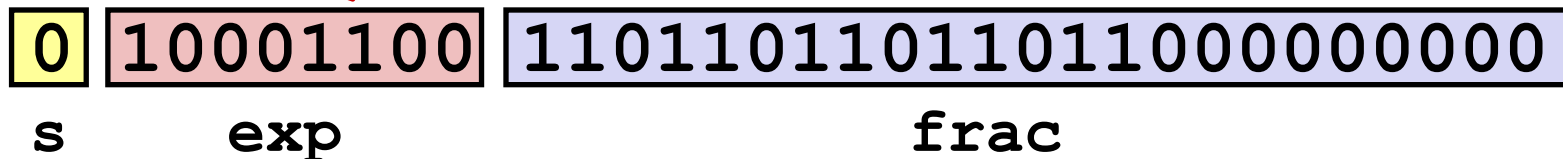
$$\text{frac} = 110110110110110000000000_2$$

■ Exponent ($E = \text{Exp} - \text{Bias}$)

$$E = 13$$

$$\text{Bias} = 127$$

$$\text{Exp} = E + \text{Bias} = 140 = 10001100_2$$



Denormalized Values

- **Condition:** $\text{exp} = 000\dots 0$
- **Exponent value:** $E = 1 - \text{Bias}$ (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 **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 very close to 0.0
 - Lose precision as get smaller
 - Equispaced

Special Values

- **Special condition: $\text{exp} = 111\dots 1$**

- **Case: $\text{exp} = 111\dots 1$, $\text{frac} = 000\dots 0$**
 - Represents value ∞ (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., $\text{sqrt}(-1)$, $\infty - \infty$, $\infty \times 0$

Interesting Numbers

{single, double}

<i>Description</i>	<i>exp</i>	<i>frac</i>	<i>Numeric Value</i>
■ Zero	00...00	00...00	0.0
■ Smallest Pos. Denorm.	00...00	00...01	$2^{-\{23,52\}} \times 2^{-\{126,1022\}}$
<ul style="list-style-type: none"> ■ Single $\approx 1.4 \times 10^{-45}$ ■ Double $\approx 4.9 \times 10^{-324}$ 			
■ Largest Denormalized	00...00	11...11	$(1.0 - \epsilon) \times 2^{-\{126,1022\}}$
<ul style="list-style-type: none"> ■ Single $\approx 1.18 \times 10^{-38}$ ■ Double $\approx 2.2 \times 10^{-308}$ 			
■ Smallest Pos. Normalized	00...01	00...00	$1.0 \times 2^{-\{126,1022\}}$
<ul style="list-style-type: none"> ■ Just larger than largest denormalized 			
■ One	01...11	00...00	1.0
■ Largest Normalized	11...10	11...11	$(2.0 - \epsilon) \times 2^{\{127,1023\}}$
<ul style="list-style-type: none"> ■ Single $\approx 3.4 \times 10^{38}$ ■ Double $\approx 1.8 \times 10^{308}$ 			

Precision

- How *accurate* are floating point numbers? Consider again:

$$15213_{10} = 11101101101101.1_2$$



110110110110110000000000

frac

- Size of **frac** limits how many decimal places we can represent
- E.g. single precision has 23 explicit bits + 1 implicit bit

$$\begin{array}{lcl}
 123456789_{10} & = & 111010110111100110100010101_2 \\
 M & = & 1.\underline{11010110111100110100010} \mathbf{101}_2 \times 2^{26} \\
 \text{frac} = & & \underline{11010110111100110100011}_2 \text{ (round up)}
 \end{array}$$

Loss of Precision Example

■ **Value:** float F = 123456789;

■ $123456789_{10} = 111010110111100110100010101_2$



$$= 1.11010110111100110100010101_2 \times 2^{26}$$



■ **Significand**

$$\begin{array}{lcl} M & = & 1.\underline{11010110111100110100010101}_2 \\ \text{frac} & = & \underline{11010110111100110100011}_2 \end{array}$$

■ **Reconstructed Value**

$$(-1)^S * M * 2^E$$

$$M = 1 + 0.5 + 0.25 + \dots = 1.83964955806732177734375$$

$$V = 1.83964955806732177734375 * 2^{26}$$

$$V = 123456792$$

Today: Floating Point

- Background: Fractional binary numbers
- Example and properties
- IEEE floating point standard: Definition
- **Floating point in C**
- Summary

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** → **int**
 - Truncates fractional part
 - Like rounding toward zero
 - Not defined when out of range or NaN: Generally sets to TMin
- **int** → **double**
 - Exact conversion, as long as **int** has ≤ 53 bit word size
- **int** → **float**
 - Will round according to rounding mode

Today: Floating Point

- Background: Fractional binary numbers
- Example and properties
- IEEE floating point standard
- Floating point in C
- **Rounding, addition, multiplication**
- Summary

Floating Point Operations: Basic Idea

$$\blacksquare \mathbf{x} +_{\mathbf{f}} \mathbf{y} = \mathbf{Round}(\mathbf{x} + \mathbf{y})$$

$$\blacksquare \mathbf{x} \times_{\mathbf{f}} \mathbf{y} = \mathbf{Round}(\mathbf{x} \times \mathbf{y})$$

■ 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 (illustrate 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 Even (default)	\$1	\$2	\$2	\$2	-\$2

Closer Look at Round-To-Even

■ Default Rounding Mode

- Hard to get any other kind without dropping into assembly
- All others are statistically biased
 - Sum of set of positive numbers will consistently be over- or underestimated

■ Applying to Other Decimal Places / Bit Positions

- When exactly halfway between two possible values
 - Round so that least significant digit is even
- E.g., round to nearest hundredth

7.8949999	7.89	(Less than half way)
7.8950001	7.90	(Greater than half way)
7.8950000	7.90	(Half way—round up)
7.8850000	7.88	(Half way—round down)

Rounding Binary Numbers

■ Binary Fractional Numbers

- “Even” when least significant bit is 0
- “Half way” when bits to right of rounding position = $100..._2$

■ Examples

- Round to nearest $1/4$ (2 bits right of binary point)

Value	Binary	Rounded	Action	Rounded Value
$2 \frac{3}{32}$	$10.00\textcolor{red}{011}_2$	10.00_2	($<1/2$ —down)	2
$2 \frac{3}{16}$	$10.00\textcolor{red}{110}_2$	10.01_2	($>1/2$ —up)	$2 \frac{1}{4}$
$2 \frac{7}{8}$	$10.11\textcolor{red}{100}_2$	11.00_2	($=1/2$ —up)	3
$2 \frac{5}{8}$	$10.10\textcolor{red}{100}_2$	10.10_2	($=1/2$ —down)	$2 \frac{1}{2}$

Scientific Notation Multiplication

■ $(2.5 \times 10^3) \times (3.0 \times 10^2) = ?$

■ **Compute result by pieces:**

- Sign : $sign_{left} * sign_{right} = 1 * 1 = 1$
- Significand : $M_{left} * M_{right} = 2.5 * 3.0 = 7.5$
- Exponent : $E_{left} + E_{right} = 3 + 2 = 5$

Result: 7.5×10^5

FP 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
 - Round M to fit **frac** precision
- **Implementation**
 - Biggest chore is multiplying significands

Scientific Notation Addition

■ $(2.5 \times 10^3) + (3.0 \times 10^2) = ?$

■ Assume E_{left} is larger than E_{right}

■ Align by decimal point:

■ Significand :

$$\begin{array}{r} 2.5 \\ + .3 \\ \hline 2.8 \end{array}$$

■ Exponent : $E = E_{\text{left}} = 3$

Result: 2.8×10^3

Floating Point Addition

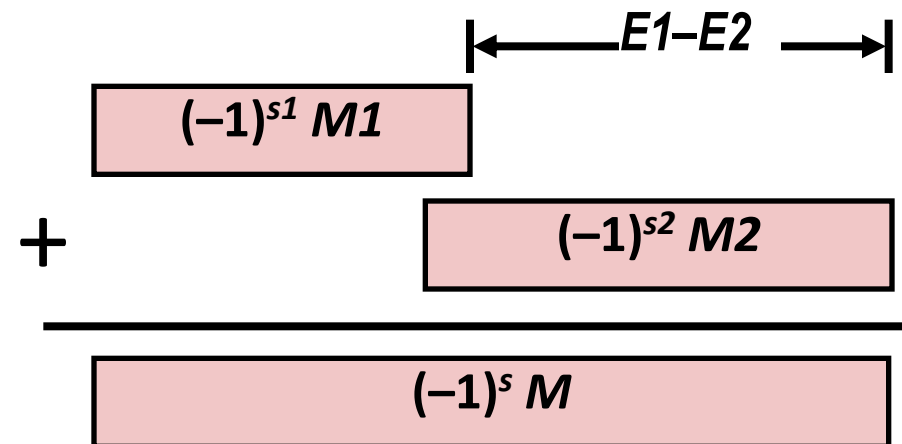
$$\blacksquare (-1)^{s1} M1 2^{E1} + (-1)^{s2} M2 2^{E2}$$

- Assume $E1 > E2$

$$\blacksquare \text{Exact Result: } (-1)^s M 2^E$$

- Sign s , significand M :
 - Result of signed align & add
- Exponent E : $E1$

Get binary points lined up



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

Mathematical Properties of FP Add

■ Compare to those of Abelian Group

- Commutative?

Yes

- $(a + b) = (b + a)$

- Associative?

No

- Overflow and inexactness of rounding

- $(3.14 + 1e10) - 1e10 = 0, \quad 3.14 + (1e10 - 1e10) = 3.14$

Mathematical Properties of FP Mult

■ Compare to Commutative Ring

- Multiplication Commutative?

Yes

- Ex: $(1e20 * 1e-20) = (1e-20 * 1e20)$

- Multiplication is Associative?

No

- Possibility of overflow, inexactness of rounding

- Ex: $(1e20 * 1e20) * 1e-20 = \text{inf}$, $1e20 * (1e20 * 1e-20) = 1e20$

- Multiplication distributes over addition?

No

- Possibility of overflow, inexactness of rounding

- $1e20 * (1e20 - 1e20) = 0.0$, $1e20 * 1e20 - 1e20 * 1e20 = \text{NaN}$

Summary

- **Represents numbers of form $M \times 2^E$**
- **One can reason about operations independent of implementation**
 - As if computed with perfect precision and then rounded
- **Not the same as real arithmetic**
 - Violates associativity/distributivity
 - Makes life difficult for compilers & serious numerical applications programmers