

Floating Point

CSC3501 Computer Organization & Design

Instructors:

Hao Wang

Today: Floating Point

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

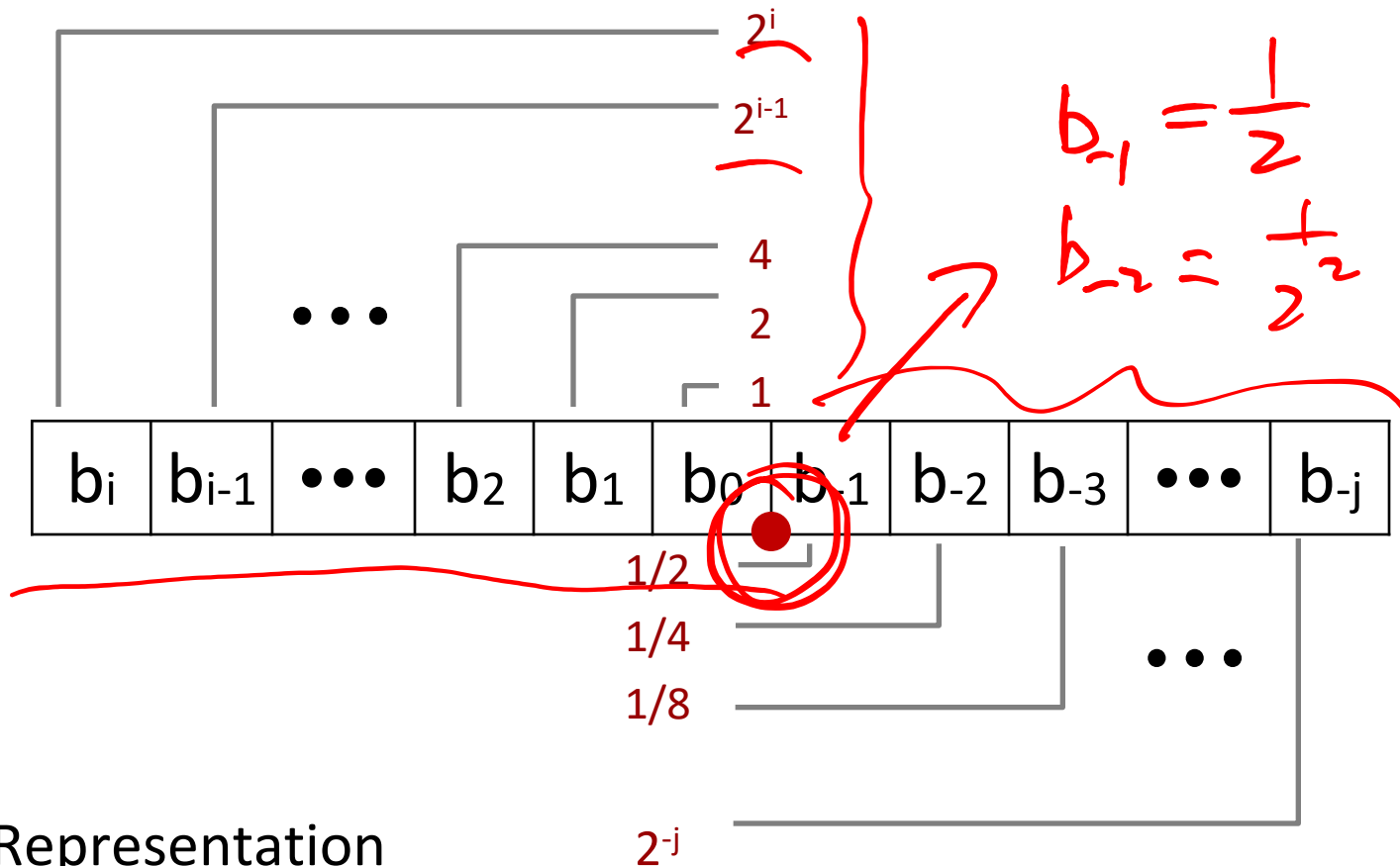
Fractional binary numbers

- What is 1011.101_2 ?

Handwritten work for converting 1011.101_2 to decimal:

$$\begin{array}{r}
 2021 \quad \frac{1}{2} \quad \frac{1}{8} \\
 \uparrow \\
 1011.101_2 \\
 \hline
 (11 + \frac{5}{8})_{10} \\
 \hline
 (11.625)_{10}
 \end{array}$$

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

5 3/4

2 7/8

1 7/16

Representation

5 101.11₂

10.111₂

1.0111₂

$$\frac{1}{2} + \frac{1}{4} = \frac{3}{4}$$

$$\rightarrow 5 + \frac{3}{4}$$

Observations

- Divide by 2 by shifting right (unsigned)
- Multiply by 2 by shifting left
- Numbers of form 0.111111...₂ are just below 1.0
 - $1/2 + 1/4 + 1/8 + \dots + 1/2^i + \dots \rightarrow 1.0$
 - Use notation $1.0 - \epsilon$

$$\cdot \frac{1}{2} \quad \frac{1}{2} \quad \dots$$

Representable Numbers

■ Limitation #1

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

Value	Representation
■ $1/3$	$0.0101010101 [01] \dots_2$
■ $1/5$	$0.001100110011 [0011] \dots_2$
■ $1/10$	$0.0001100110011 [0011] \dots_2$



■ Limitation #2

- Just one setting of binary point within the w bits
 - Limited range of numbers (very small values? very large?)

Today: Floating Point

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

IEEE 802.x
WiFi

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

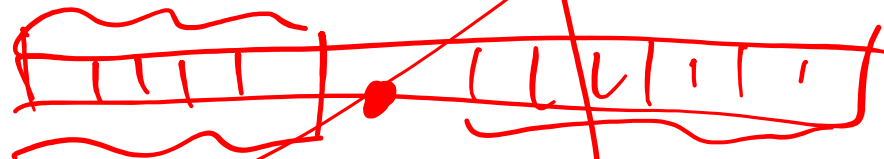
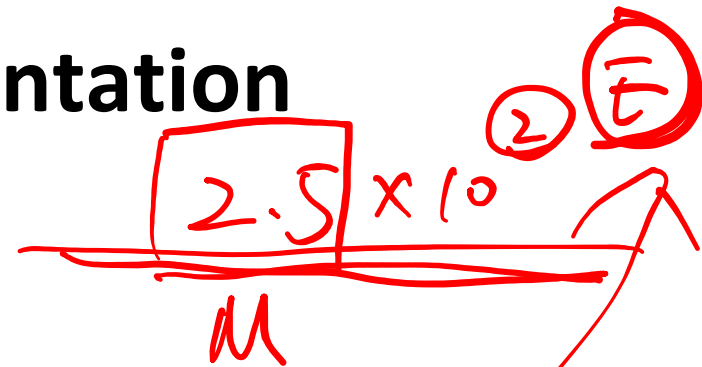
Numerical Form:

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

- 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

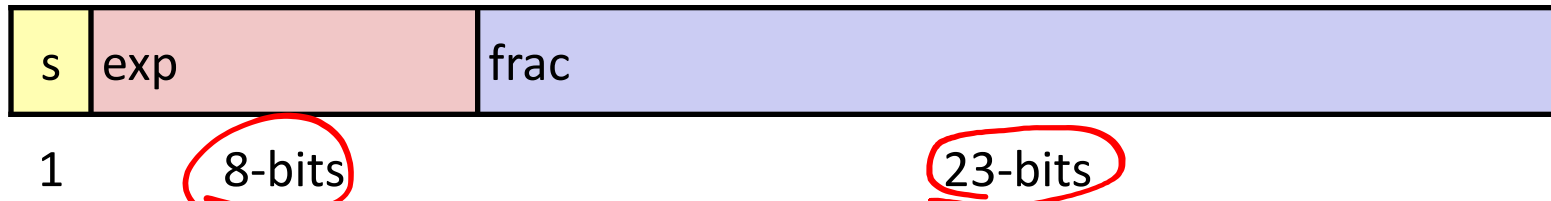
Encoding

- 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)

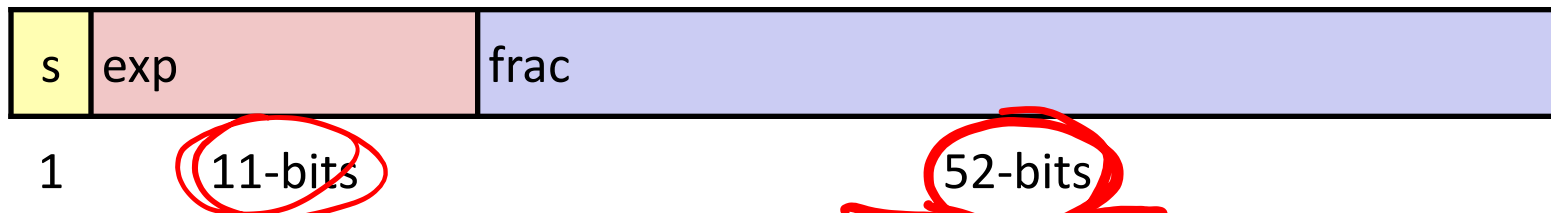


Precision options

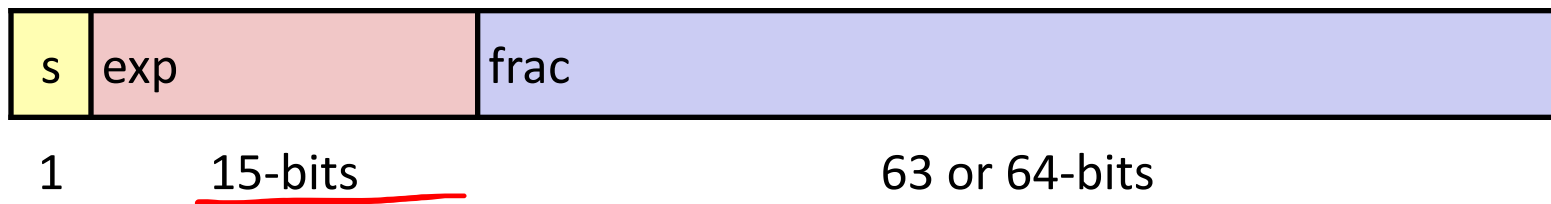
- Single precision: 32 bits



- Double precision: 64 bits



- Extended precision: 80 bits (Intel only)



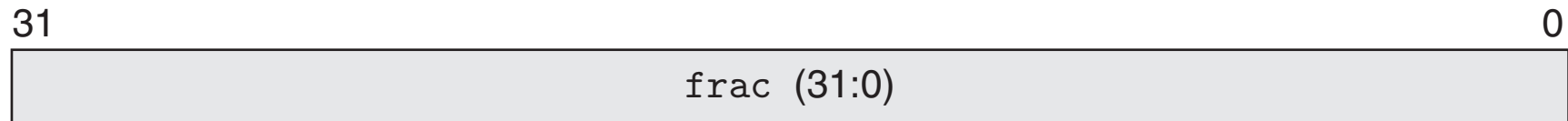
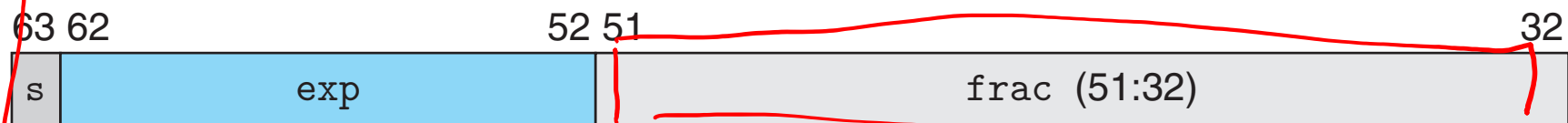
Aligned Memory View

Single precision



word.

Double precision

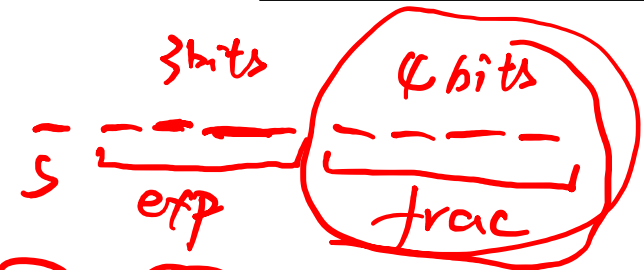


word.

“Normalized” Values

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

- When: $\text{exp} \neq 000\dots 0$ and $\text{exp} \neq 111\dots 1$



- Exponent coded as a biased value: $E = \text{Exp} - \text{Bias}$

- Exp: **unsigned** value of exp field
- 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)

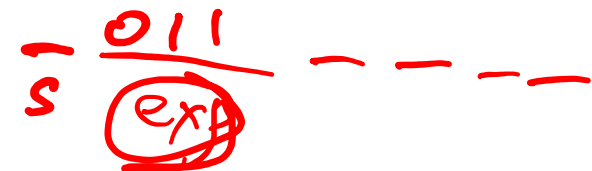
$$E = \text{Exp} - \text{Bias} = 2^{3-1} - 1 = 3$$

0 1 1

$$M = 1 + 0.\text{frac}$$

- Significand coded with **implied leading 1**: $M = 1.\text{xxx}\dots\text{x}_2$

- xxx...x: bits of frac field
- Minimum when frac=000...0 ($M = 1.0$)
- Maximum when frac=111...1 ($M = 2.0 - \epsilon$)
- Get extra leading bit for “free”



Normalized Encoding Example

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

$$E = \text{Exp} - \text{Bias}$$

Value: float $F = 15213.0$;

$$15213_{10} = 11101101101101_2 \rightarrow 13 \text{ bits}$$

$$= 1.1101101101101_2 \times 2^{13}$$

Significand

$$M = 1.1101101101101_2$$

$$\text{frac} = 1011011011010000000000_2$$

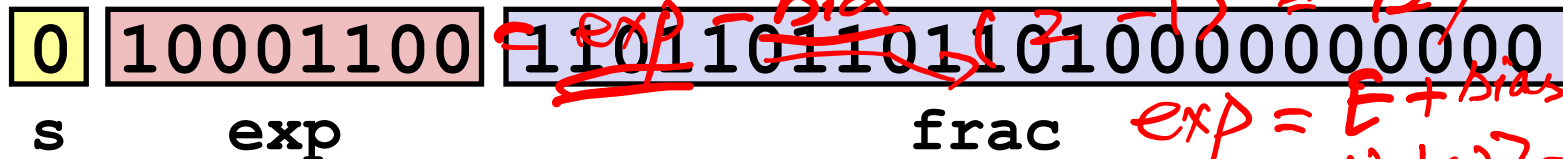
Exponent

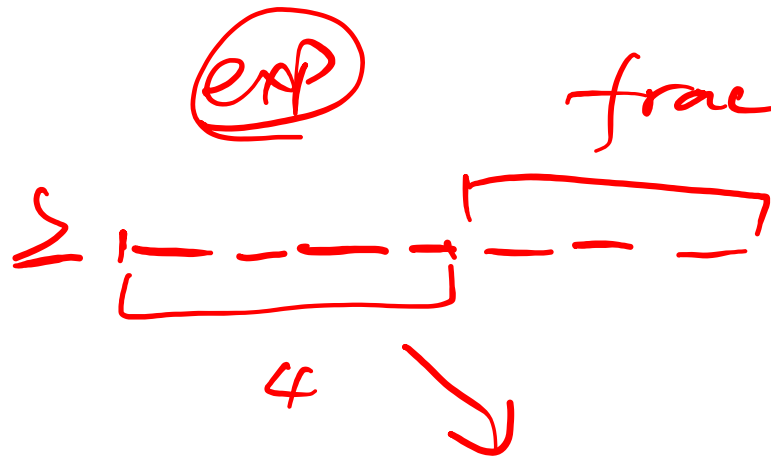
$$E = 13$$

$$\text{Bias} = 127$$

$$\text{Exp} = 140 = 10001100_2$$

Result:




 2^E

$$\bar{E} = \text{exp} - \text{bias}$$

$$\frac{4-1}{2-1} = 7$$

$$\frac{(-1)^s M \times 2^E}{M = 1 + \text{frac}}$$

$$M = 1 + \text{frac}$$

$$\begin{array}{r} 0000 \rightarrow 0 \\ 1111 \rightarrow 15 \\ \hline 8421 \end{array}$$

$$[7, 8]$$

$$(-1)^s M \times \left(\frac{1}{2} \right) \approx 2$$

$$[1, 2)$$

$$[1, 2) \cup (2 \times 2^{15})$$

$$0.1111$$

Denormalized Values

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

$$E = 1 - \text{Bias}$$

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

Special Values

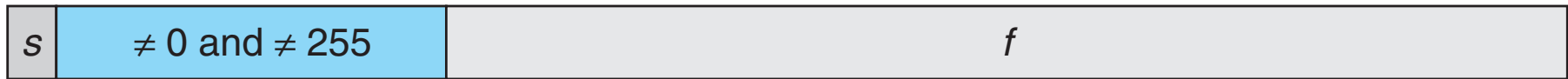
- Condition: **exp = 111...1**

- Case: **exp = 111...1, frac = 000...0**
 - Represents value ∞ (infinity)
 - Operation that overflows
 - Both positive and negative
 - E.g., $1.0/0.0 = -1.0/-0.0 = +\infty$, $1.0/-0.0 = -\infty$

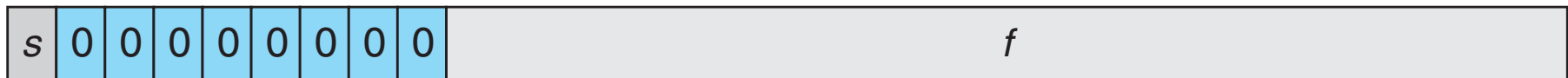
- Case: **exp = 111...1, frac \neq 000...0**
 - Not-a-Number (NaN)
 - Represents case when no numeric value can be determined
 - E.g., $\text{sqrt}(-1)$, $\infty - \infty$, $\infty \times 0$

The Three Cases

1. Normalized



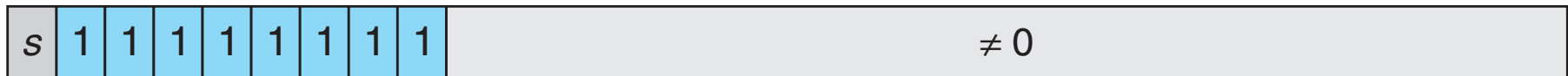
2. Denormalized



3a. Infinity



3b. NaN



Visualization: Floating Point Encodings

