

Floating Point

Math 426

University of Alaska Fairbanks

September 21, 2020

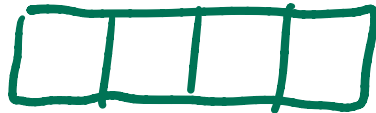
Signed Integers (Offset Binary)



E.g., 4-bit representation of signed integers.



Unished:



↑
0,1

2^4 possibilities.

$$0 - 15 = 2^4 - 1$$

0 0 0 0 → 0

0 0 0 1 → 1

0 0 1 0 → 2

⋮

1 1 1 1 → $8 + 4 + 2 + 1 = 15$

Signed Integers (Offset Binary)

E.g., 4-bit representation of signed integers.

Offset binary subtract off an offset (bias)

0 0 0 0 \rightarrow -7

0 0 0 1 \rightarrow -6

0 0 1 0 \rightarrow -5

⋮

1 0 0 0 \rightarrow 1

⋮

1 1 1 1 \rightarrow 8
15

offset for 4 bits

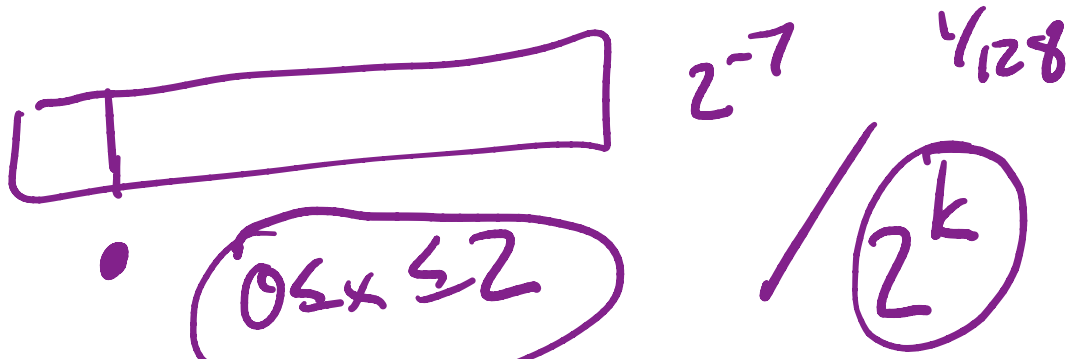
is $1 - 2^3$

Signed Integers (Offset Binary)

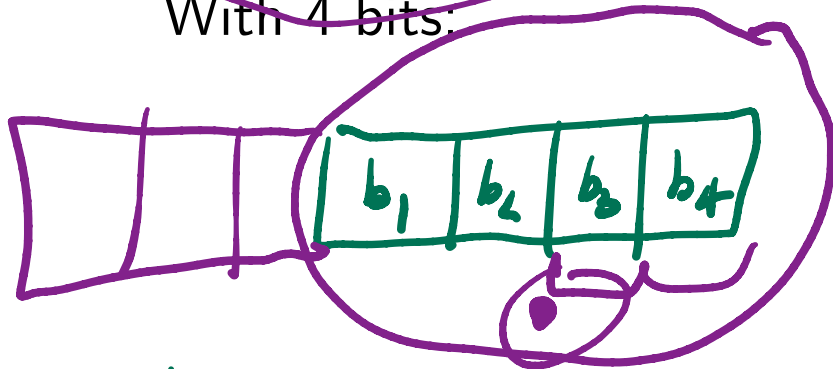
E.g., 4-bit representation of signed integers.

In most other applications, two's complement is used!

Unsigned Fixed Point Numbers



With 4 bits:

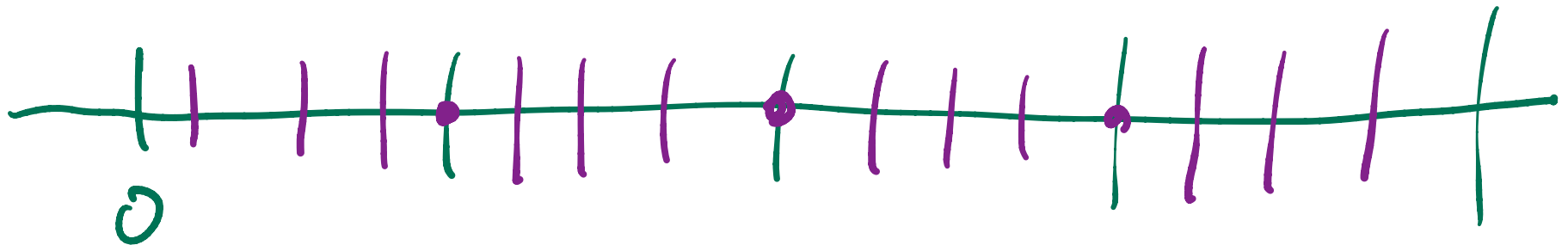


$$(b_1 b_2 . b_3 b_4)_2$$

smallest $(00.01)_2 = \frac{1}{4}$

biggest $(11.11)_2 = 3\frac{3}{4}$

end
↓ 4



Floating Point Numbers

in between we use
 m .

$$(-1)^s$$

$$s \cdot m \cdot 2^E$$

- ▶ s , sign, ± 1
- ▶ m , mantissa, $1 \leq m < 2$ → fixed point number
- ▶ E , an integer, exponent

↳ offset binary (signed number)

$$1: s = 0, m = 1, E = 0 \quad (-1)^0 \cdot 1 \cdot 2^0 = 1$$

$$2: s = 0, m = 1, E = 1$$

Fictional 8-bit format



- ▶ sign: 1 bit $s=0 \rightarrow +, s=1 \rightarrow - \quad (-1)^s$
- ▶ exponent: 3 bits
- ▶ mantissa: 4 bits



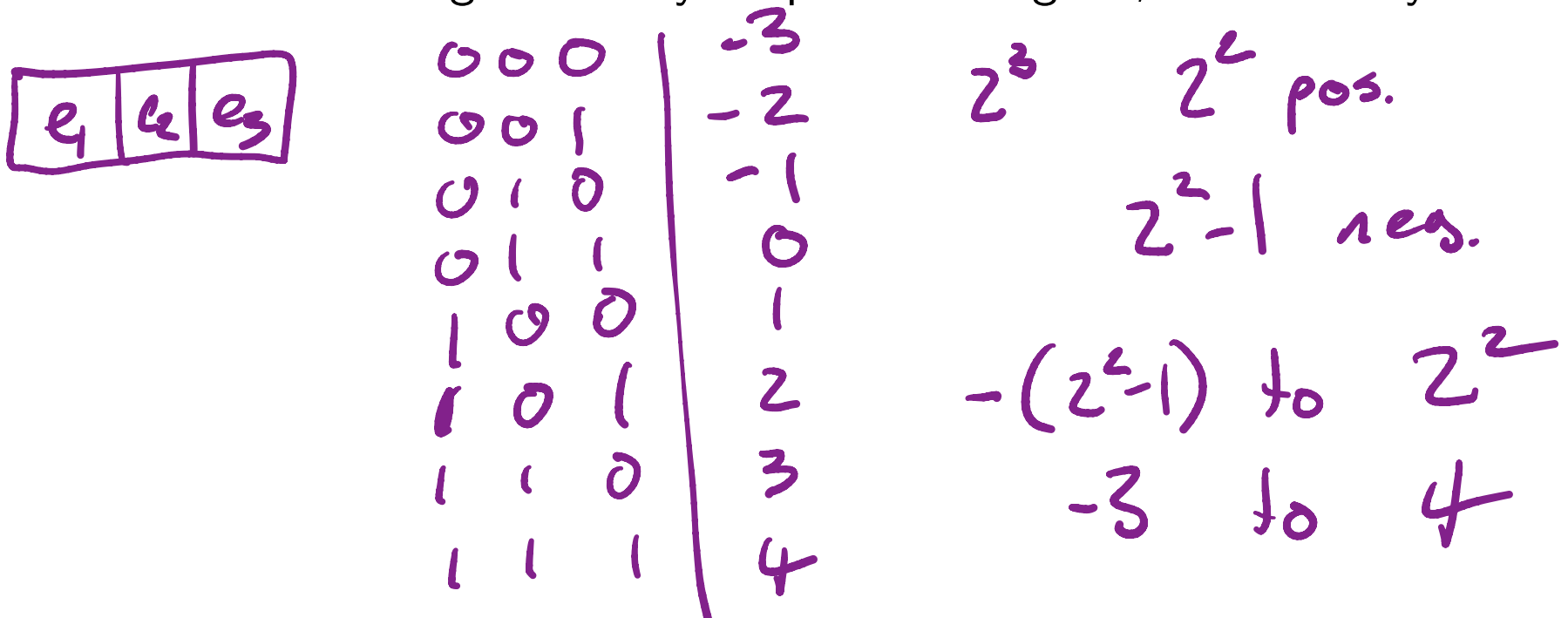
$1 \leq m < 2$	
$1.000)_2$	1
1.001_2	$1 + \frac{1}{8}$
\vdots	\vdots
1.111_2	$1 + \frac{7}{8}$

hidden bit format.
 The 1 is implicit.
 $(1.b_1b_2b_3b_4)_2$
 $1.0000_2 = 1$
 $1.0001_2 = 1 + \frac{1}{16}$
 \vdots
 $1.1111_2 = 1 + \frac{15}{16}$

Fictional 8-bit format

- ▶ sign: 1 bit
- ▶ exponent: 3 bits
- ▶ mantissa: 4 bits

Mantissa is in unsigned binary. Exponent is signed, offset binary.



Hidden bit representation

All numbers between 1 and 2 in base two start with a 1, so we can save a bit and gain precision by making the 1 implicit.

Hidden bit representation

All numbers between 1 and 2 in base two start with a 1, so we can save a bit and gain precision by making the 1 implicit. All available numbers:

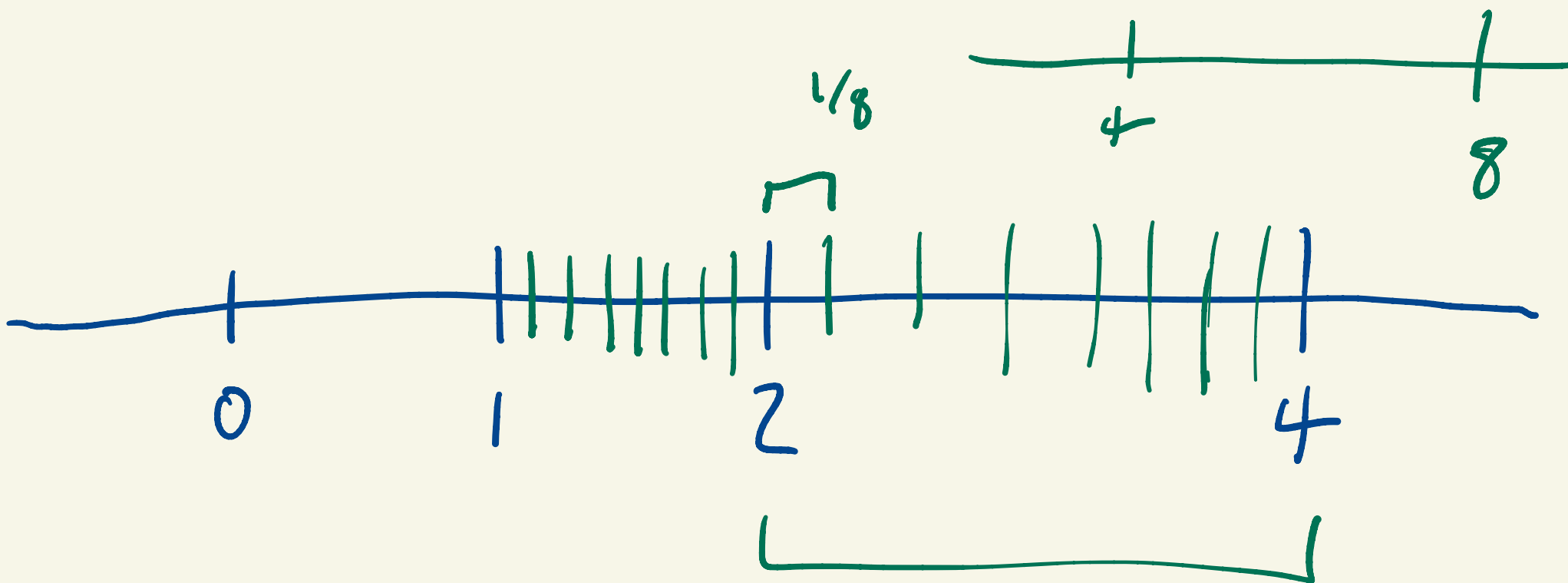
Available numbers:

$$\pm \left[1, \frac{1}{16}, \dots, \frac{15}{16} \right] \times \left[\begin{array}{c} [-3, -2, \dots, 4] \\ 2 \\ \hline \frac{1}{8}, \frac{1}{4}, \dots, 8, 16 \end{array} \right]$$

Available numbers:

$$\pm \left[1, \frac{1}{16}, \dots, \frac{15}{16} \right] \times \left[-3, -2, \dots, 4 \right]$$

$\frac{1}{8}, \frac{1}{4}, \dots, 8, 16$

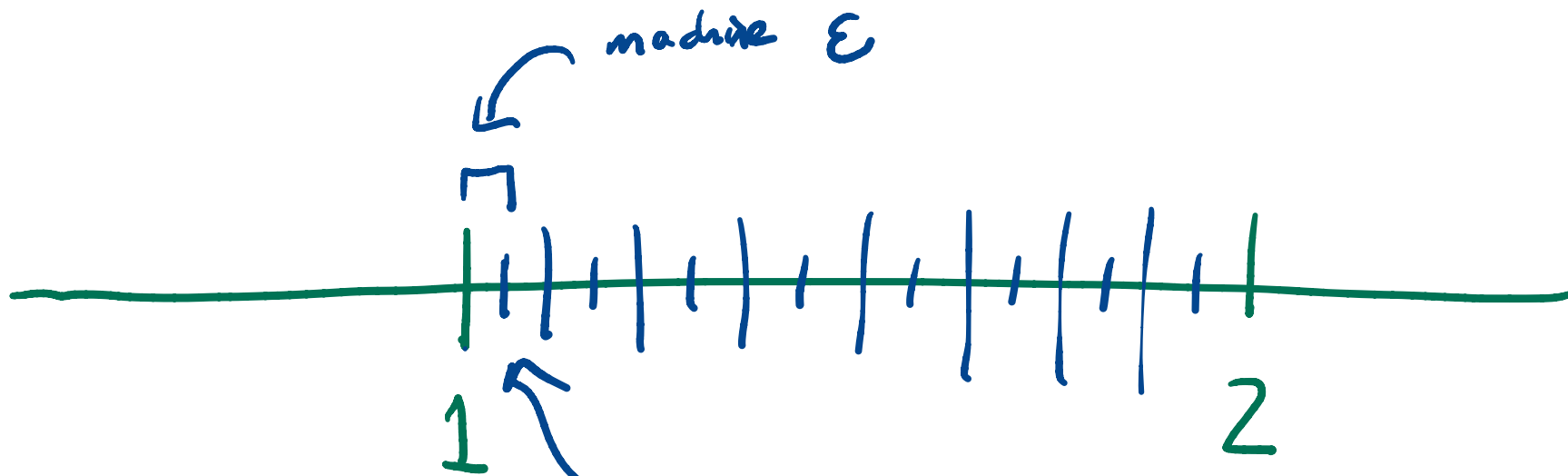


$$s=0$$

$$E=0 \quad 2^0$$

$$E=1, 2^1$$

Precision and Machine Epsilon

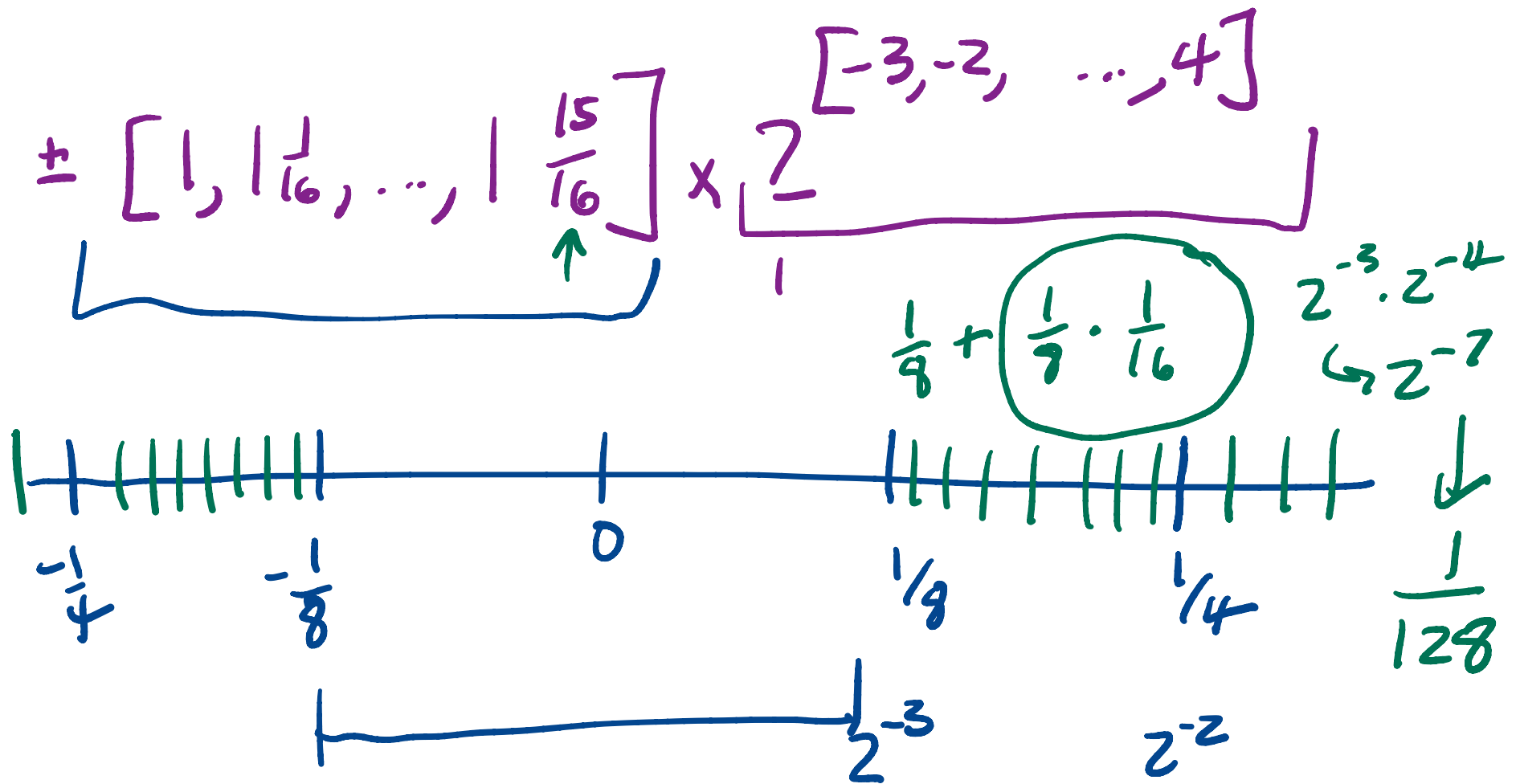


8 bit format $\epsilon = 1/16$

- $1 + \epsilon$
 - $1 + 2\epsilon$
 - \vdots
 - $2 - \epsilon$
-

The Monster Gap Around Zero

Smallest positive number



Subnormal Numbers (And two zeros!)

Give up some exponents:

$$000 \rightarrow -3 \quad 1-2^2$$

$$001 \rightarrow -2$$

⋮

⋮

⋮

⋮

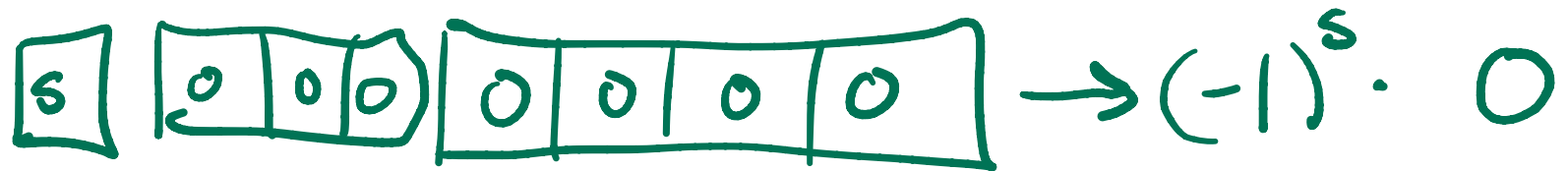
$$110 \rightarrow 3$$

$$111 \rightarrow 4 \quad 2^2$$

$$\begin{array}{c} 2-2^2 \\ \vdots \\ 2^2-1 \end{array}$$

Subnormal Numbers (And two zeros!)

000 \longmapsto very small
111 \longmapsto weird



$$s = 0 \rightarrow +0$$

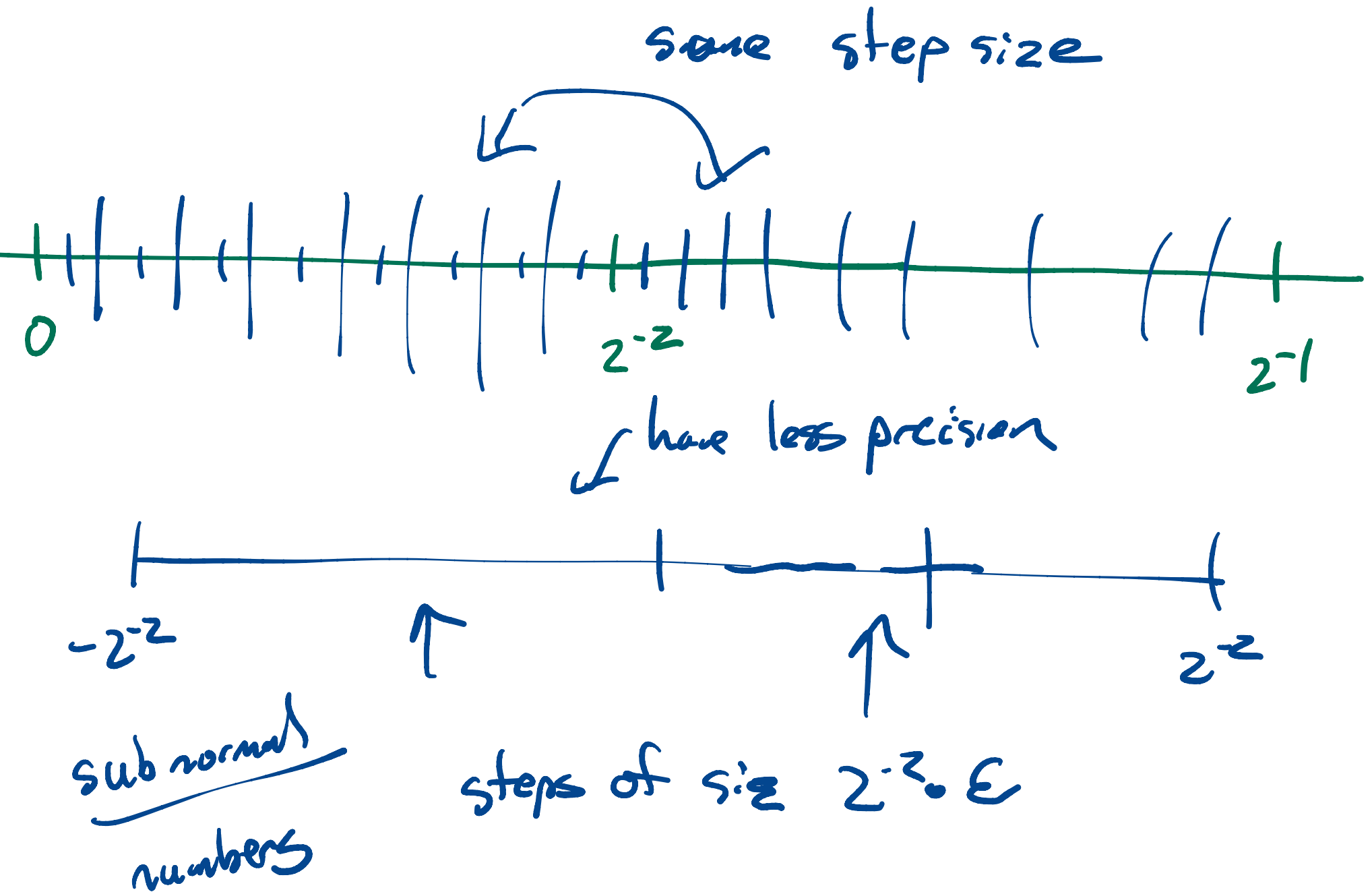
$$s = 1 \rightarrow -0$$

s 000 b_1 b_2 b_3 b_4

$$(-1)^s (0.b_1 b_2 b_3 b_4)_2 \cdot 2^{-2}$$

$\frac{1}{16}$

Subnormal Numbers (And two zeros!)



Infinity

$e: 111 \rightarrow \cancel{4}$

Inf
- Inf

Exponent is all 1's. Mantissa is all 0's.

$\rightarrow \pm \infty$

Positive infinity: 0 111 0000

Negative infinity: 1 111 0000

$x + y$
 \uparrow

$\infty + ?$

$\infty + 5 = \infty$

$7/\infty \rightarrow 0$

$-7/\infty \rightarrow -0$

Infinity

Exponent is all 1's. Mantissa is all 0's.

Positive infinity: 0 111 0000

Negative infinity: 1 111 0000

Infinity

Exponent is all 1's. Mantissa is all 0's.

Positive infinity: 0 111 0000

Negative infinity: 1 111 0000

$$\frac{1}{+0} \rightarrow +\text{Inf}$$
$$\frac{1}{-0} \rightarrow -\text{Inf}$$

Any other pattern s 111 $b_1b_2b_3b_4$ is Not a Number (NaN).

$$x = \text{NaN}$$

$$\frac{0}{0} \rightarrow \text{NaN}$$

IEEE 754

Single precision: 32 bits.

1. sign: 1 bit
2. exponent: 8 bits
3. mantissa: 23 bits

Machine epsilon: $2^{-23} \approx 2.2 \times 10^{-7}$.

Smallest (normal) number: $2^{-126} \approx$