

# Representing real numbers in the computer

floating point

$$\underbrace{3.271}_{\text{mantissa}} \times \underbrace{10^{-4}}_{\text{base}}$$

Signed integer exponent

$$\left. \begin{array}{l} 32.71 \times 10^{-5} \\ 0.327 \times 10^{-3} \end{array} \right\} \text{point can "float"}$$

Scientific notation.

Usually, we write  $3.271 \times 10^{-4}$

only one digit to the left of the decimal point.

If exponent is 0, we don't usually write  $3.271 \times 10^0$  but instead write just 3.271

# Scientific Notation vs Fixed point

$$5.781 \times 10^{-4} \quad \text{vs} \quad 0.0005781$$

Which can store more information in least amount of space?

$-4$  one exponent.

Floating Point

$$\underbrace{5.781}_{4 \text{ digits}} \times 10$$

Fixed point.

$$\underbrace{0.0005781}_{7 \text{ digits}}$$

hard to store really small numbers or really large numbers.

Format statements

Matlab

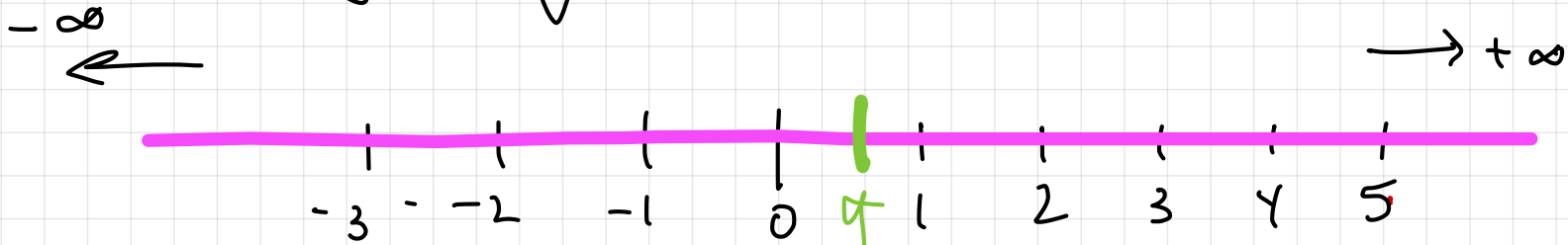
```
fprintf('% .4e', 5.781e-4)
```

```
fprintf('% .8f', 5.781e-4)
```

```
>> x = 5.781e-4;  
>> fprintf('% .4e\n', x)  
5.7810e-04  
>> fprintf('% .7f\n', x)  
0.0005781
```

# Digital representation of real numbers

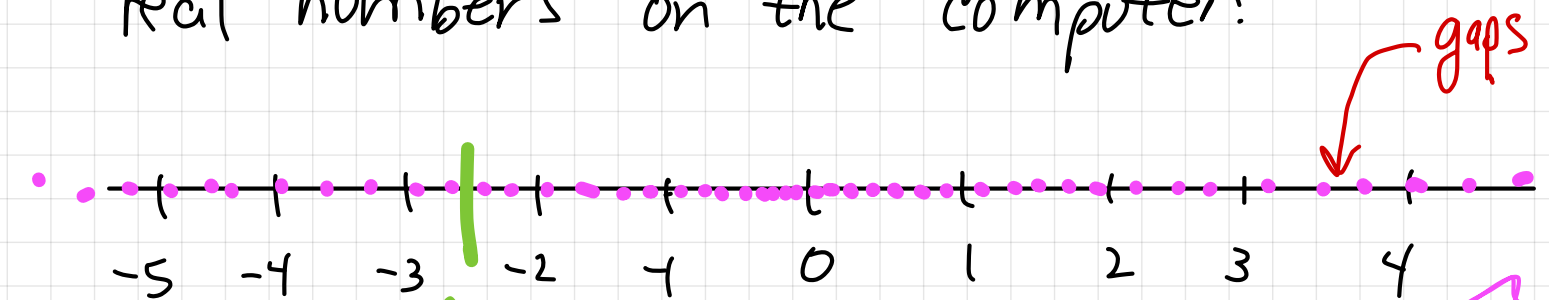
① infinity many real numbers



"infinite precision"

always hit a red number

② We can only represent a finite set of real numbers on the computer.



might not hit exact representation.

numbers are more densely packed near 0

less dense for large positive/negative numbers.

What numbers can we represent on the computer?

Numbers that can be represented in binary as:

$$\pm (1 + b_1 2^{-1} + b_2 2^{-2} + b_3 2^{-3} + \dots) \times 2^e$$

where  $b_i = 0$  or  $1$ ,  $i = 1, 2, 3, \dots$  and

where  $e$  is an integer.  $e = \pm 1, \pm 2, \pm 3, \dots$

Example:  $x = 0.5$  (in decimal)

$$x = 1.0 \times 2^{-1} \quad (\text{base } 2) \quad b_1 = b_2 = b_3 = \dots = 0$$
$$e = -1$$

Example:  $x = 2$

$$x = 1.0 \times 2^1 \quad b_1 = b_2 = b_3 = \dots = 0$$
$$e = 1$$

$$13 = 1 + 4 + 8 = 1 + 2^2 + 2^3$$

Example:

$$.25 = 2^{-2} :$$

$$13 \cdot 25 = 2^2 + 1 + 2^2 + 2^3$$

$$X = 13 \cdot 25 =$$

$$= 2^3 (2^{-5} + 2^{-3} + 2^{-1} + 1)$$

$$= (1 + 2^{-1} + 2^{-3} + 2^{-5}) \times 2^3$$

← factor out  $2^3$  out

$$b_1 = 1 \quad b_2 = 0 \quad b_3 = 1 \quad b_4 = 0 \quad b_5 = 1$$

$$b_6 = b_7 = \dots = 0$$

$$e = 3$$

Example:  $X = 0.1$

$X = 0.1 = 1 \times 10^{-1} \Rightarrow$  cannot be represented exactly!

$$X = (1 + 2^{-1} + 2^{-4} + 2^{-5} + 2^{-8} + 2^{-9} + 2^{-12} + 2^{-13} + 2^{-16} + 2^{-17} + 2^{-20} + 2^{-21} + \dots) \times 2^{-4}$$

$$= 0.09999 \dots = 0.0\bar{9}$$

$$= 0.1$$

continues for ever!

Example:  $X = 0$

$$(1 + \dots) \times 2^0$$

do we have to have this  
1 here? 😞

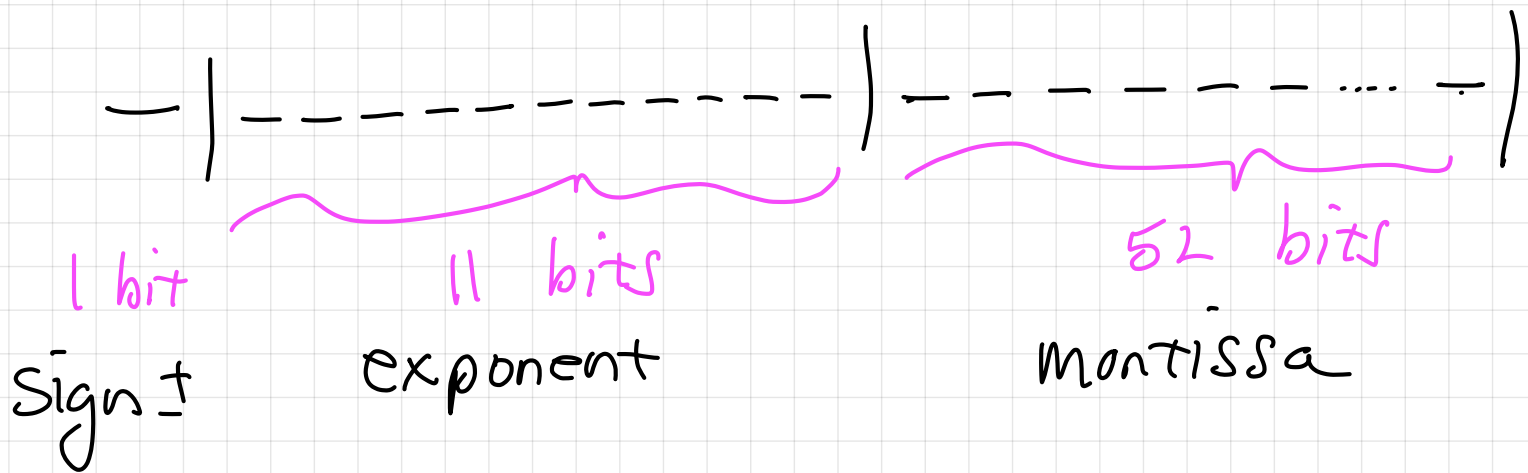
Example:  $X = -5$

How do we get a negative value?

How do we encode the binary representation on the computer?

- Need a "signed exponent"
- Need a mantissa
- Need a  $\pm$  sign
- Need a 0

- Encode as a string of 0s and 1s
- IEEE** Floating point standard. 754 standard.



"64-bit precision" = "double precision"

Sign bit: 1 bit

0 : positive  
1 : negative

exponent: 11 bits

"(biased)"

to get negative exponents

mantissa:  $b_1 b_2 b_3 \dots b_{52}$

## more on the exponent

11 bits allows us to represent

$$E = e_0 2^0 + e_1 2^1 + e_2 2^2 + \dots + e_{10} 2^{10}$$

Example:  $E = 6$   $e_i = 0 \text{ or } 1$

$\underline{e_{10}} \dots \underline{e_3} \underline{e_2} \underline{e_1} \underline{e_0}$

$$E = 2 + 4 = 2^1 + 2^2$$

$$e_0 = 0 \quad e_1 = e_2 = 1 \quad e_3 = e_4 = \dots = e_{10} = 0$$

Smallest  $E$ :  $e_0 = e_1 = \dots = e_{10} = 0$   $E = 0$

Largest  $E$ :  $e_0 = e_1 = \dots = e_{10} = 1$

$$E = 2^{11} - 1 = 2047$$

But how do we get a negative exponent? Shift the exponent

First half represent negative numbers  
second half represents positive numbers



0 1 2 3 4 ..... 2047

2048 possible values  
for E.

Designate 0, 1, ..., 1022 as "negative"

1023 as 0

1024, 1025, ..., 2047 as  
"positive"

"bias the exponent"  $e = E - 1023$

To represent  $X = 0.5$

$$E = e + 1023$$

$$X = 1 \times 2^{-1} \quad e = -1$$

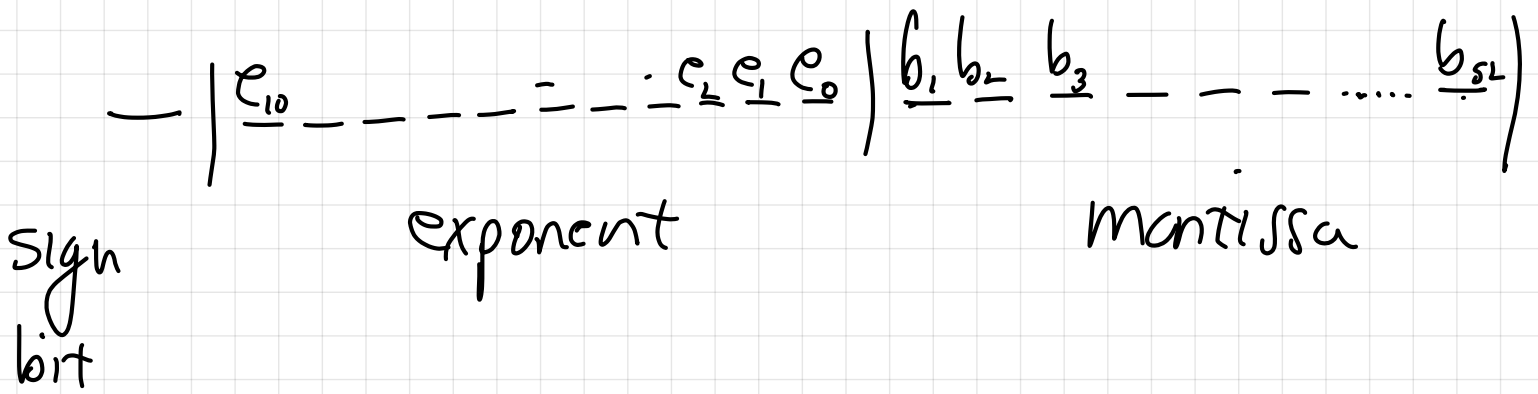
$$E - 1023 = -1 \Rightarrow E = 1022$$

$$1022 = 2^1 + 2^2 + 2^3 + \dots + 2^9 \quad \begin{matrix} 1024 - 2 \\ 2^{10} - 2 \end{matrix}$$

$$E = \underline{e_{10}} \underline{e_9} \underline{e_8} \underline{e_7} \underline{e_6} \underline{e_5} \underline{e_4} \underline{e_3} \underline{e_2} \underline{e_1} \underline{e_0}$$

$$= (0111111110)_b$$

## Full representation



"normalized mantissa"

$$1, b_1, b_2, b_3, \dots, b_{s_2}$$

↖ "1" is assumed

Example:  $X = 0.5$

$$X = 1.0 \times 2^{-1}$$

$$= (0 \mid 01111111110 \mid 000000 \dots 0)_b$$

biased exponent = 1022

52 0s

Example

$$X = 2$$

$$10^3 - 1$$

$$\frac{1000}{0999}$$

$$X = 1.0 \times 2^1$$

$$\text{exponent: } 1 = 1 \times 2^0$$

$$e = E - 1023 = 0$$

$$E = 1023 = 2^{10} - 1$$

0 1 1 1 1 1 1 1 1 1 1 1

(0 | 0 1 1 1 1 1 1 1 1 1 | 0 0 0 0 ... 0)<sub>2</sub>

+1

$$E = e - 1023$$

mantissa = 52 0s

$$1 = 2 - 1 = 2^1 - 1 = (01)_2$$

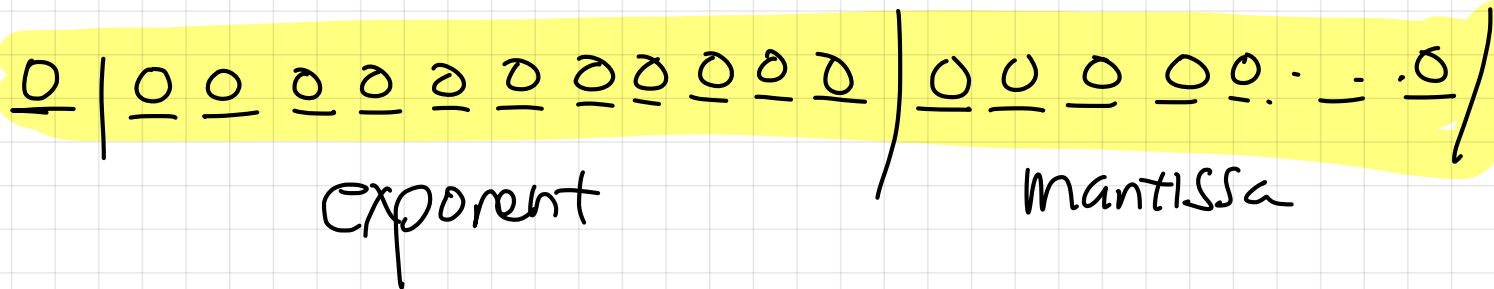
$$3 = 4 - 1 = 2^2 - 1 = (011)_2$$

$$7 = 8 - 1 = 2^3 - 1 = (0111)_2$$

$$15 = 16 - 1 = 2^4 - 1 = (01111)_2$$

$$1023 = 1024 - 1 = 2^{10} - 1 = (0111111111)_2$$

How do we represent 0?

$$X =$$


means:  $X = 0$

Not  $\Gamma = 0 \Rightarrow e = E - 1023$   
 $= -1023$

mantissa 1.0 ?

$X = 1.0 \times 2^{-1023}$  NO!

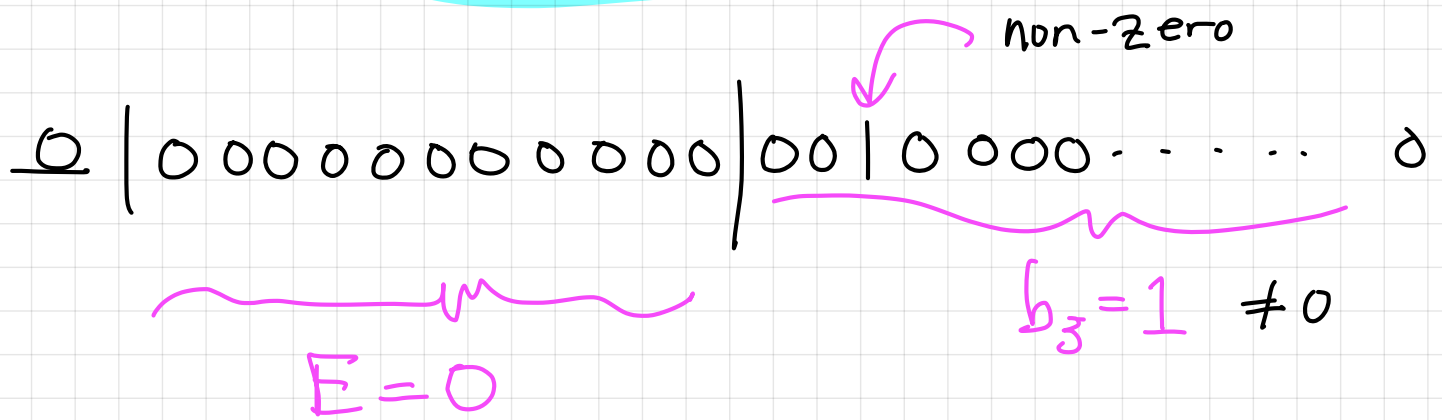
The  $\mathbb{I}\mathbb{E}\mathbb{E}\mathbb{E}$  standard tells us how

0 is represented  $\Rightarrow$  all bits set to 0.

The mantissa is "denormalized"

If all bits in the exponent are 0, then we drop the "1" on the mantissa. <sup>not 1023</sup> instead of  $1.0 \times 2^{-1023} \Rightarrow 0.0 \times 2^{-1022}$

# Denormalized numbers



If  $E = 0 \Rightarrow$  exponent is  $-1022$  (not  $-1023$ )

Drop implied "1" on mantissa  
 $\Rightarrow$  Get a "denormalized number"

denormal  
number

$$0.125 \times 2^{-1022} \approx 10^{-308}$$

Denormalized numbers are less precise, because they will have fewer digits in the mantissa.

Something like:  $0.00000034$   
vs  $3.4713671 \times 10^{-7}$

these extra digits  
are lost.

# Other special numbers in the IEEE standard

- $\pm 0$
- all bits set to 0
  - possible to have  $+0$   $-0$  depending on how sign bit is set.

- $\pm \infty$
- all exponent bits set to 1
  - all mantissa bits set to 0
  - sign bit can be 1 or 0

- NAN
- all exponent bits set to 1
  - mantissa bits can be set to any non-zero value
  - "quiet Nans" - propagate through calculations
  - "signaling Nans" - signals a floating point exception
- "not a number"
- the end 