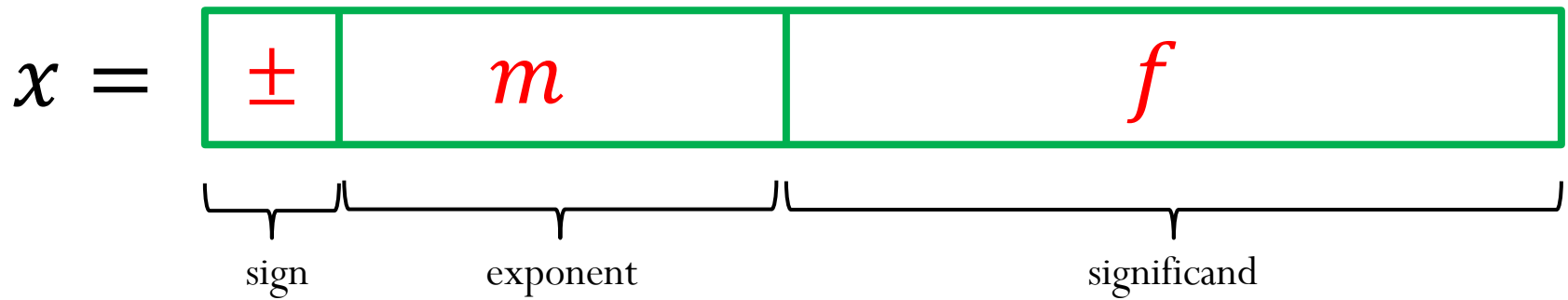


Machine numbers: how floating point numbers are stored?

Floating-point number representation

What do we need to store when representing floating point numbers in a computer?

$$x = \pm 1.f \times 2^m$$



Initially, different floating-point representations were used in computers, generating inconsistent program behavior across different machines.

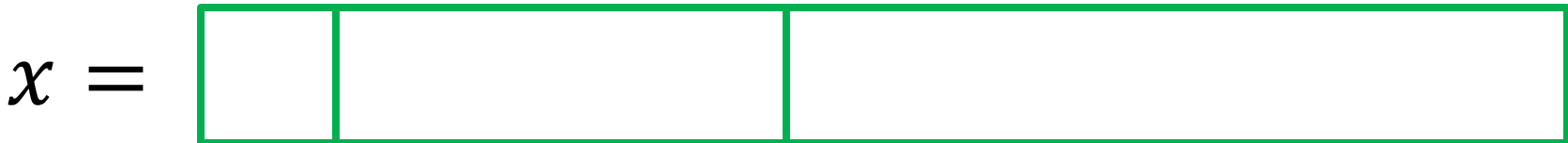
Around 1980s, computer manufacturers started adopting a standard representation for floating-point number: IEEE (Institute of Electrical and Electronics Engineers) 754 Standard.

Floating-point number representation

Numerical form:

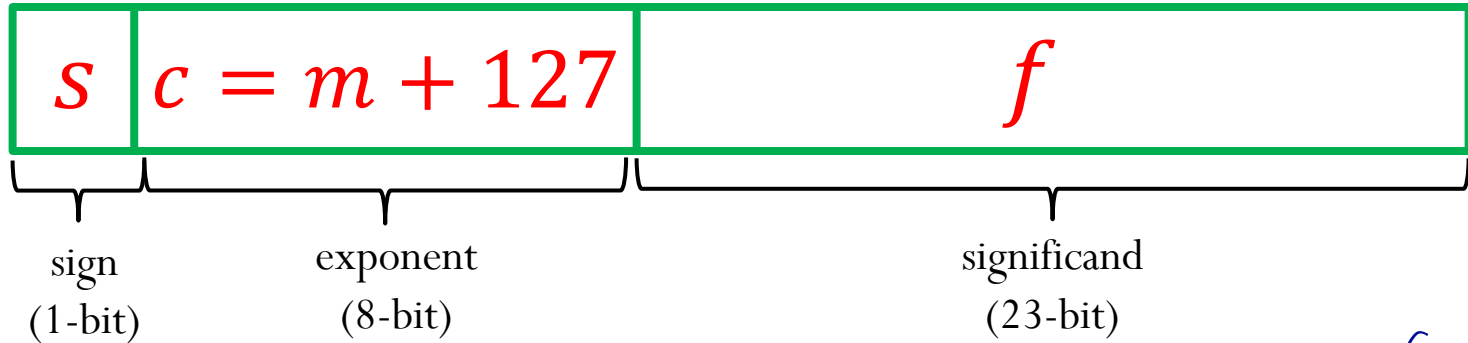
$$x = \pm 1.f \times 2^m$$

Representation in memory:



IEEE-754 Single Precision (32-bit)

$$x = (-1)^s 1.f \times 2^m$$



$C = (00 \dots 0)_2 \rightarrow$ reserved $\rightarrow C = (0)_{10}$
 $C = (11 \dots 1)_2 \rightarrow$ NaN or $\infty \rightarrow C = (255)_{10}$

$$1 \leq C \leq 254 \rightarrow 1 \leq m + 127 \leq 254$$

$$-126 \leq m \leq 127$$

range of
exponent \Rightarrow

$$m \in [-126, 127]$$

$$n = 23 \Rightarrow p = 24$$

$$m \in [-126, 127]$$

IEEE-754 Single Precision (32-bit)

$$x = (-1)^s 1.f \times 2^m = \boxed{s \quad c \quad f} \quad c = m + 127$$

- Machine epsilon** (ϵ_m): is defined as the distance (gap) between 1 and the next largest floating point number.

$$1.\underbrace{000 \dots 00}_{23 \text{ bits}} \times 2^0$$

$$1.0000 \dots 01 \times 2^0$$

$$\frac{0.\underbrace{000}_{2^{-1}2^{-2}2^{-3}} \dots \underbrace{01}_{2^{-23}} \times 2^0}{}$$

$$\epsilon_m = 2^{-23} \approx 1.2 \times 10^{-7}$$

$$\epsilon_m = 2^{-n}$$

- Smallest positive normalized FP number:**

$$UFL = 1.00 \dots 0 \times 2^{-126}$$

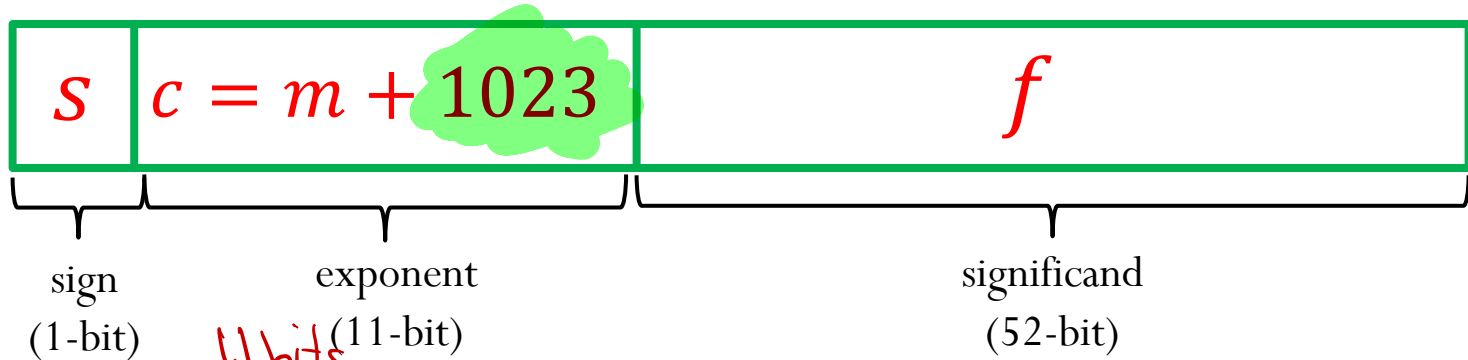
$$UFL = 2^4 \approx 10^{-38}$$

- Largest positive normalized FP number:**

$$OFL = 2^{127} (1 - 2^{-24}) = 2^{127} (1 - 2^{-24}) \approx 10^{38}$$

IEEE-754 Double Precision (64-bit)

$$x = (-1)^s 1.f \times 2^m$$



$$c = (\underbrace{000 \dots 00}_{11 \text{ bits}})_2 = (0)_{10} \rightarrow \text{reserved for zero}$$

$$c = (111 \dots 11)_2 = (2047)_{10} \rightarrow \text{NaN and } \infty$$

$$1 \leq c \leq 2046 \rightarrow 1 \leq m + 1023 \leq 2046$$

$$m \in [-1022, 1023]$$

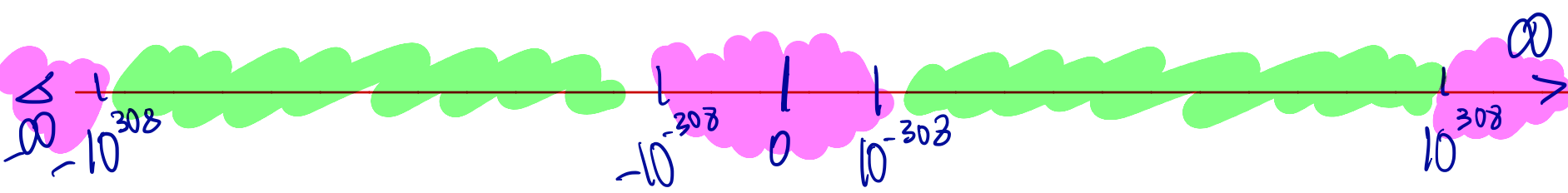
$$-1022 \leq m \leq 1023$$

$m \in [-1022, 1023]$ $m=52 \rightarrow p=53$
 IEEE-754 Double Precision (64-bit)

$$x = (-1)^s 1.f \times 2^m = \boxed{s \quad c \quad f} \quad c = m + 1023$$

- **Machine epsilon** (ϵ_m): is defined as the distance (gap) between 1 and the next largest floating point number.

$$\epsilon_m = 2^{-n} \rightarrow \epsilon_m = 2^{-52} \approx 2.2 \times 10^{-16}$$



- **Smallest positive normalized FP number:**

$$UFL = 2^L \rightarrow UFL = 2^{-1022} \approx 2.2 \times 10^{-308}$$

- **Largest positive normalized FP number:**

$$OFL = 2^{U+1} (1 - 2^{-p}) = 2^{1024} (1 - 2^{-53}) \approx 10^{308}$$

Can we represent # smaller than UFL

$$x = L.f \times 2^m$$

$$x = 0.f \times 2^L$$

denormalized / subnormal

Let's make

$$c = (00 \dots 00)_2 \quad \text{all zeros}$$

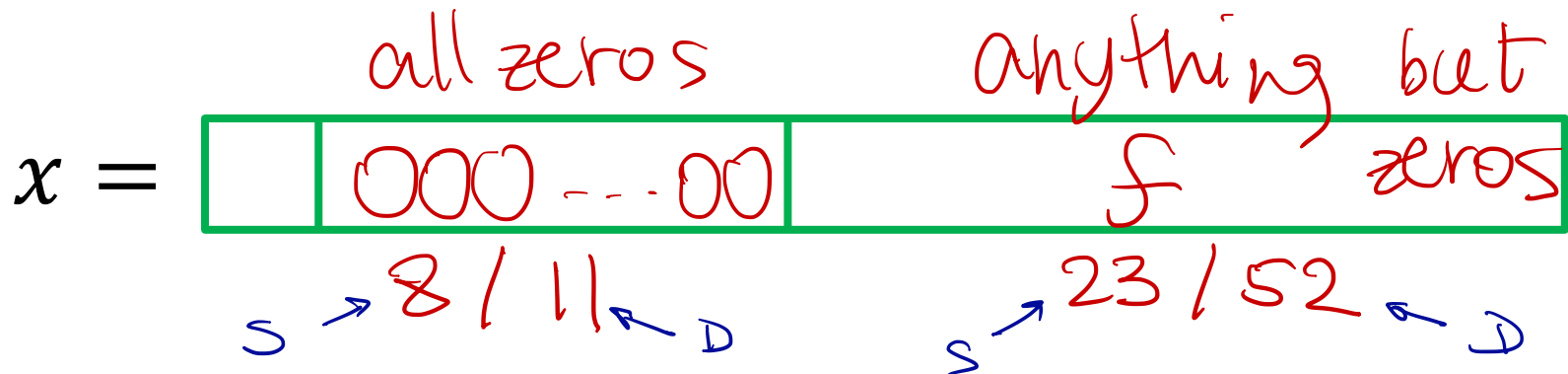
$$f = \text{anything but zeros}$$

used to "indicate" that exponent is $m = L$ (and NOT to evaluate m)

Subnormal (or denormalized) numbers

- Noticeable gap around zero, present in any floating system, due to normalization
- Relax the requirement of normalization, and allow the leading digit to be zero, only when the exponent is at its minimum ($m = L$)
- Computations with subnormal numbers are often slow.

Representation in memory (another special case):



Numerical value:

$$X = \pm 0. \mathbf{f} \times 2^L \quad \left| \begin{array}{l} m \in [-126, 127] \text{ single} \\ m \in [-1022, 1023] \text{ double} \end{array} \right.$$

Subnormal (or denormalized) numbers

IEEE-754 Single precision (32 bits):

$$c = (00000000)_2 = 0$$

smallest subnormal : $(0.\underset{2^{-1}}{0}\underset{2^{-2}}{0}\underset{2^{-3}}{0}\dots\underset{2^{-23}}{001})_2 \times 2^{-126} = 2^{-23} \times 2^{-126} \approx 1.4 \times 10^{-45}$

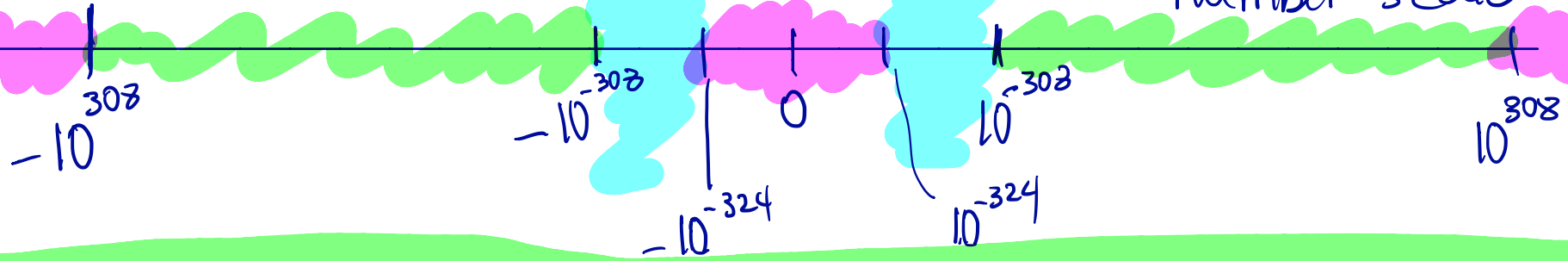
IEEE-754 Double precision (64 bits):

$$c = (000000000000)_2 = 0$$

smallest subnormal : $(0.\underset{2^{-1}}{0}\underset{2^{-2}}{0}\dots\underset{2^{-52}}{001})_2 \times 2^{-1022} = 2^{-52} \times 2^{-1022} \approx 4.9 \times 10^{-324}$

Subnormal (or denormalized) numbers

double precision number scale



$1.\underbrace{000 \dots 00}_{23 \text{ bits}} \times 2^{-126} \rightarrow \text{UFL} \quad p = 24$

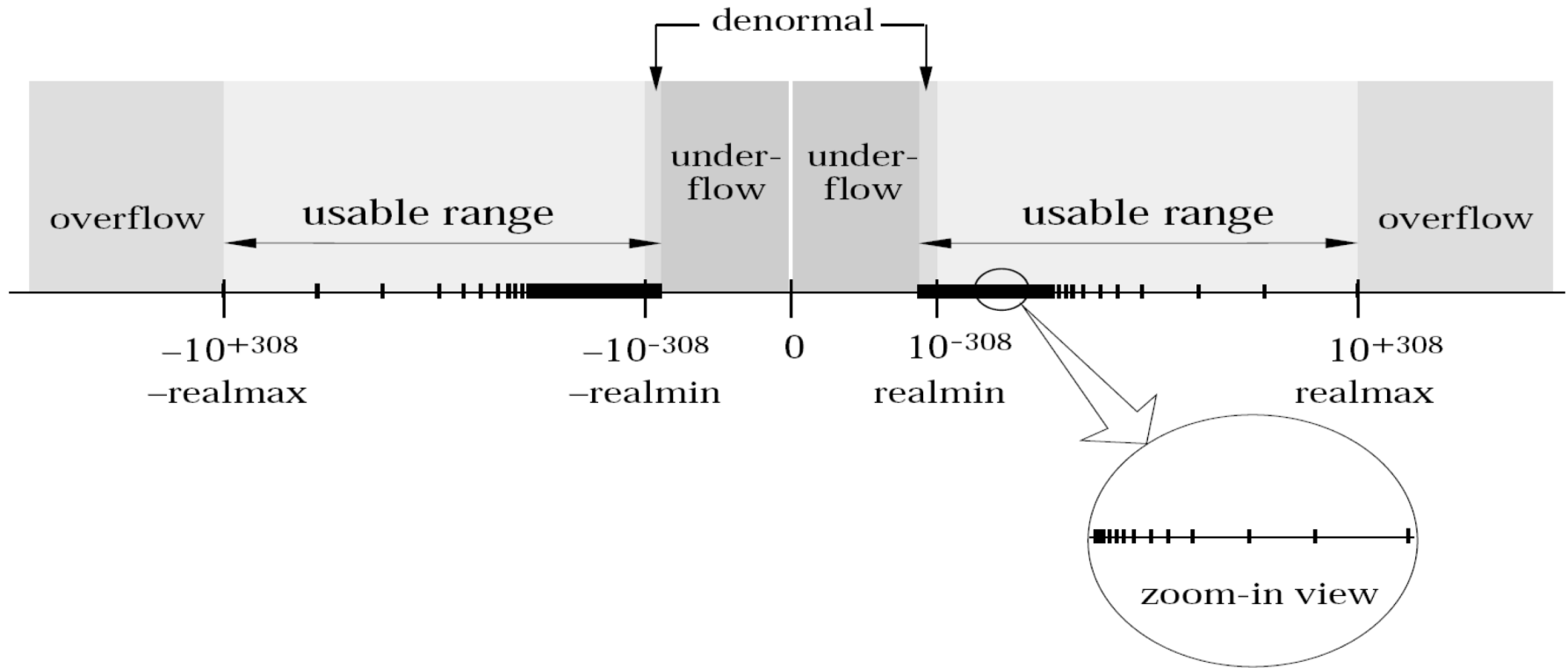
$0.\underbrace{1111 \dots 11}_{23 \text{ bits}} \times 2^{-126} \rightarrow p = 23$

$0.\underbrace{000 \dots 00111}_{23 \text{ bits}} \times 2^{-126} \rightarrow p = 3$

→ not significant bits
Loss of precision !!

Example using single precision

IEEE-754 Double Precision



Summary for Single Precision

$$x = (-1)^s 1.f \times 2^m = \boxed{\begin{array}{|c|c|c|} \hline s & c & f \\ \hline \end{array}} \quad m = c - 127$$

Stored binary exponent (c)	Significand fraction (f)	value
00000000	0000...0000	zero
00000000	$any\ f \neq 0$	$(-1)^s 0.f \times 2^{-126}$
00000001	$any\ f$	$(-1)^s 1.f \times 2^{-126}$
⋮	⋮	⋮
11111110	$any\ f$	$(-1)^s 1.f \times 2^{127}$
11111111	$any\ f \neq 0$	NaN
11111111	0000...0000	infinity

What is the equivalent decimal number?

0 00000000 000000000000000000000000

1 11111111 000000000000000000000000

0 11111111 1111111111000011111111

0 00000000 111100000000000000000000

0 01111111 000000000000000000000000

Clicker question

$$p = n + 1 = 4$$

A number system can be represented as $x = \pm 1.b_1b_2b_3 \times 2^m$
for $m \in [-5, 5]$ and $b_i \in \{0, 1\}$.

$$1.000 \times 2^{-5} = 2^{-5}$$

1) What is the smallest positive normalized FP number:

- a) 0.0625 b) 0.09375 **c) 0.03125** d) 0.046875 e) 0.125

2) What is the largest positive normalized FP number:

- a) 28 **b) 60** c) 56 d) 32 1.111×2^5 or $2^{u+1}(1-2^{-p}) = 2^6(1-2^{-4})$

3) How many additional numbers (positive and negative) can be represented when using subnormal representation?

- a) 7 **b) 14** c) 3 d) 6 e) 16

$$0.001 \times 2^{-5}$$

$$0.100$$

$$.010$$

$$.011$$

$$.101$$

$$.110$$

$$.111$$

4) What is the smallest positive subnormal number?

- a) 0.00390625** b) 0.00195313 c) 0.03125 d) 0.0136719

5) Determine machine epsilon

- a) 0.0625 b) 0.00390625 c) 0.0117188 **d) 0.125**

$$\epsilon_m = 2^{-n} = 2^{-3}$$

A number system can be represented as $x = \pm 1.b_1b_2b_3b_4 \times 2^m$
 for $m \in [-6,6]$ and $b_i \in \{0,1\}$.

$n=4$ $p=5$

1) Let's say you want to represent the decimal number 19.625 using the binary number system above. Can you represent this number exactly?

$$(19.625)_{10} = (10011.101)_2 = 1.0011101 \times 2^4$$

$$1.0011 \times 2^4$$

integer range \rightarrow until 2^p
 double precision $\rightarrow 2^{53}$

2) What is the range of integer numbers that you can represent exactly using this binary system?

$$(1)_{10} = (1)_2 = 1.00000 \times 2^0$$

$$(2)_{10} = (10)_2 = 1.00000 \times 2^1$$

$$(3)_{10} = (11)_2 = 1.1000 \times 2^0$$

$$(15)_{10} = (1111.0)_2 = 1.1110 \times 2^3$$

$$(11111)_2 = (31)_{10} = 1.1111 \times 2^4$$

$$(32)_{10} = 1.00000 \times 2^5$$

$$1.00001 \times 2^5 = 34$$

cannot represent $(33)_{10} !!$

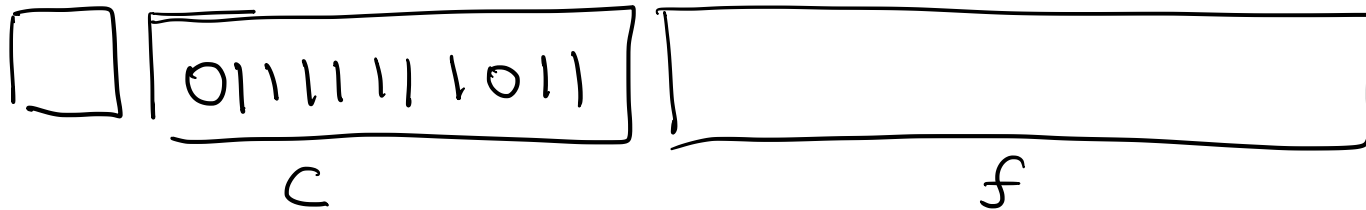
Rounding errors

Example

Show demo: "Waiting for 1".

Determine the double-precision machine representation for 0.1

$$0.1 = (0.000110011 \overline{0011} \dots)_2 = (1.100110011 \dots)_2 \times 2^{-4}$$



$$m = -4 \longrightarrow c = m + 1023 \longrightarrow c = (1019)_{10}$$

$$f = 1001 \ 1001 \ 1001 \ \dots \ 1001 \ 1001 \ \dots$$

Diagram illustrating the fraction field 'f' with a red bracket under the first 52 bits labeled '52'. The next 4 bits are labeled '1010' and the remaining bits are labeled 'IGNORE' with an arrow pointing to the right.

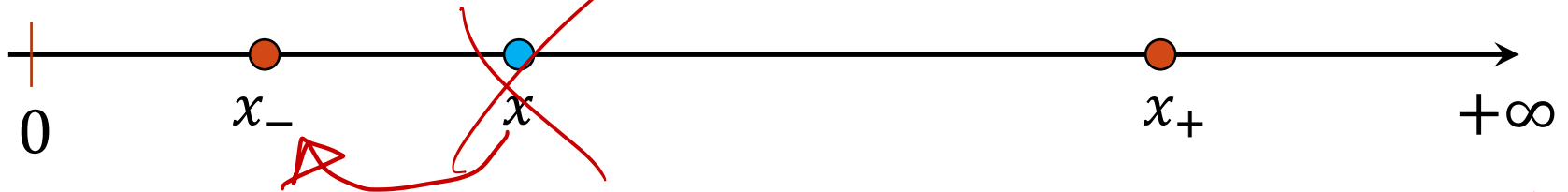
Machine floating point number

- Not all real numbers can be exactly represented as a machine floating-point number.

- Consider a real number in the normalized floating-point form:

$$x = \pm 1.b_1b_2b_3 \dots b_n \dots \times 2^m$$

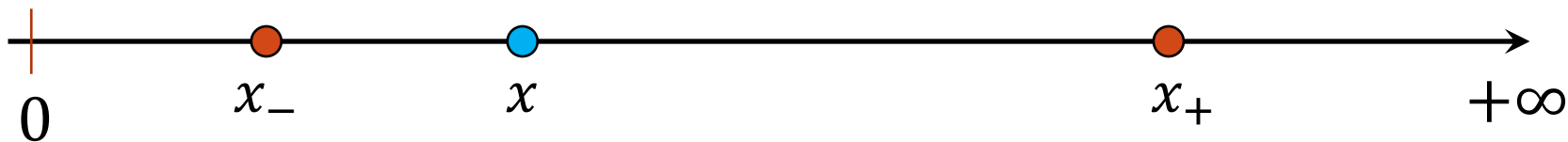
- The real number x will be approximated by either x_- or x_+ , the nearest two machine floating point numbers.



$$x_- = 1.b_1b_2b_3 \dots b_n \times 2^m$$

$$x_+ = \underbrace{1.b_1b_2b_3 \dots b_n \times 2^m}_{x_-} + \underbrace{0.000 \dots 01 \times 2^m}_n$$

$$x_+ - x_- = 2^{-n} \times 2^m$$



Exact number: $x = 1.b_1b_2b_3 \dots b_n \dots \times 2^m$

$$x_- = 1.b_1b_2b_3 \dots b_n \times 2^m$$

$$x_+ = 1.b_1b_2b_3 \dots b_n \times 2^m + \underbrace{0.000 \dots 01}_{\epsilon_m} \times 2^m$$

Gap between x_+ and x_- : $|x_+ - x_-| = \epsilon_m \times 2^m$

Examples for single precision:

x_+ and x_- of the form $q \times 2^{-10}$

x_+ and x_- of the form $q \times 2^4$:

x_+ and x_- of the form $q \times 2^{20}$:

x_+ and x_- of the form $q \times 2^{60}$:

Handwritten calculations in red:

$$2^{-23} \times 2^{-10} = 2^{-33} \sim 10^{-10}$$

$$\rightarrow 2^4 \times 2^{-23} \sim 10^{-6}$$

$$\rightarrow 2^{20} \times 2^{-23} \sim 0.125$$

$$\rightarrow 2^{60} \times 2^{-23} = 2^{37} \sim 10^{11}$$

The interval between successive floating point numbers is not uniform: the interval is smaller as the magnitude of the numbers themselves is smaller, and it is bigger as the numbers get bigger.

Gap between two successive machine floating point numbers

A "toy" number system can be represented as $x = \pm 1.b_1b_2 \times 2^m$
for $m \in [-4,4]$ and $b_i \in \{0,1\}$.

$(1.00)_2 \times 2^0 = 1$	$(1.00)_2 \times 2^1 = 2$	$(1.00)_2 \times 2^2 = 4.0$
$(1.01)_2 \times 2^0 = 1.25$	$(1.01)_2 \times 2^1 = 2.5$	$(1.01)_2 \times 2^2 = 5.0$
$(1.10)_2 \times 2^0 = 1.5$	$(1.10)_2 \times 2^1 = 3.0$	$(1.10)_2 \times 2^2 = 6.0$
$(1.11)_2 \times 2^0 = 1.75$	$(1.11)_2 \times 2^1 = 3.5$	$(1.11)_2 \times 2^2 = 7.0$

$(1.00)_2 \times 2^3 = 8.0$	$(1.00)_2 \times 2^4 = 16.0$	$(1.00)_2 \times 2^{-1} = 0.5$
$(1.01)_2 \times 2^3 = 10.0$	$(1.01)_2 \times 2^4 = 20.0$	$(1.01)_2 \times 2^{-1} = 0.625$
$(1.10)_2 \times 2^3 = 12.0$	$(1.10)_2 \times 2^4 = 24.0$	$(1.10)_2 \times 2^{-1} = 0.75$
$(1.11)_2 \times 2^3 = 14.0$	$(1.11)_2 \times 2^4 = 28.0$	$(1.11)_2 \times 2^{-1} = 0.875$

$(1.00)_2 \times 2^{-2} = 0.25$	$(1.00)_2 \times 2^{-3} = 0.125$	$(1.00)_2 \times 2^{-4} = 0.0625$
$(1.01)_2 \times 2^{-2} = 0.3125$	$(1.01)_2 \times 2^{-3} = 0.15625$	$(1.01)_2 \times 2^{-4} = 0.078125$
$(1.10)_2 \times 2^{-2} = 0.375$	$(1.10)_2 \times 2^{-3} = 0.1875$	$(1.10)_2 \times 2^{-4} = 0.09375$
$(1.11)_2 \times 2^{-2} = 0.4375$	$(1.11)_2 \times 2^{-3} = 0.21875$	$(1.11)_2 \times 2^{-4} = 0.109375$