

# 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?

$$x = \pm 1.\textcolor{red}{f} \times 2^{\textcolor{red}{m}}$$

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.

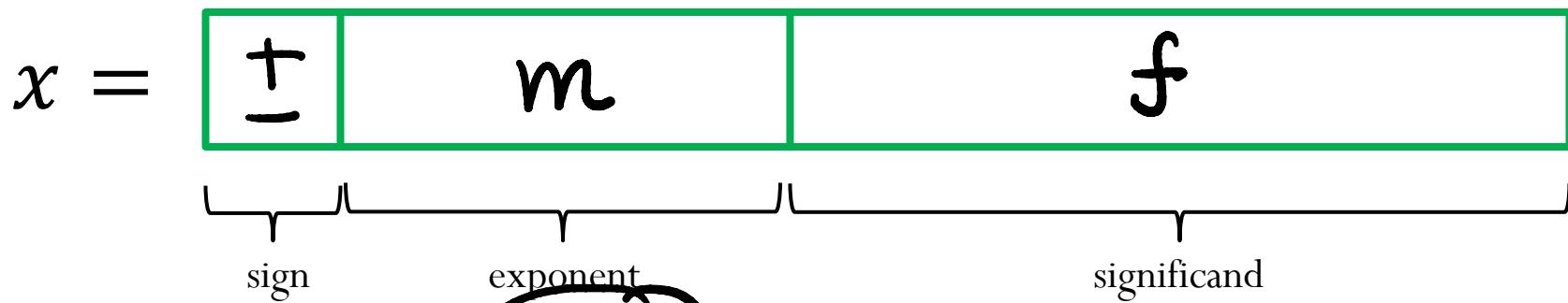
# Floating-point number representation

Numerical form:

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

Representation in memory:

$$\begin{aligned} m &\in [L, U] \\ m &\in [-4, 4] \end{aligned}$$



$c = m + \text{shift}$

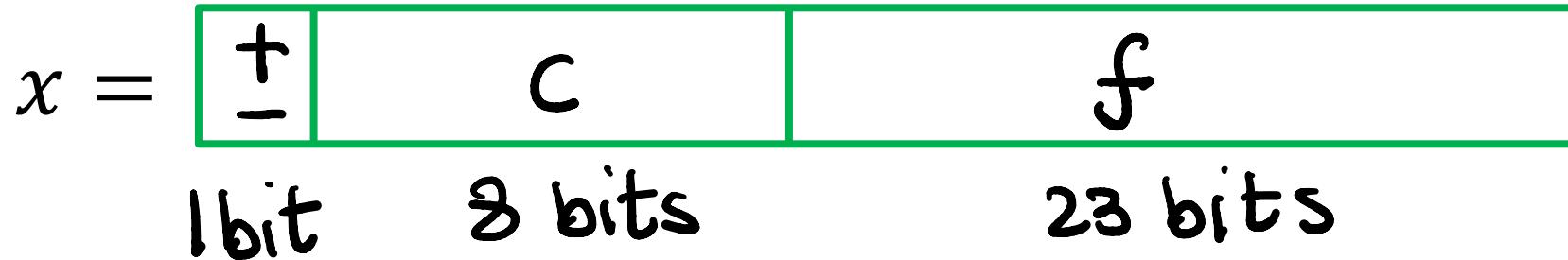
↑  
unsigned int      signed int

# Precisions:

Finite representation: not all numbers can be represented exactly!

$$x = \pm 1.f \times 2^{c-shift}$$

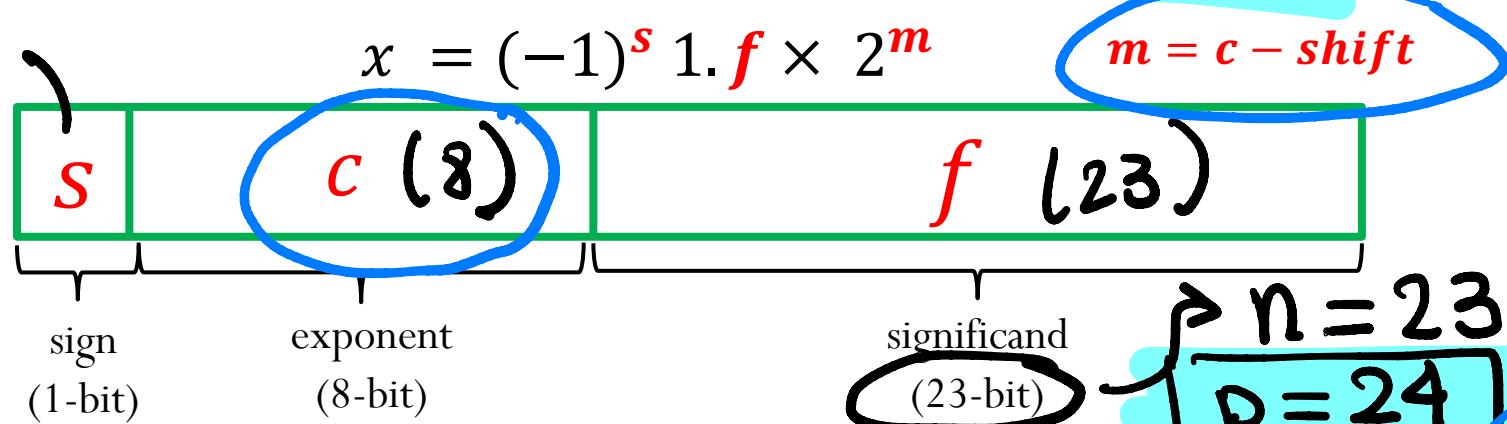
IEEE-754 Single precision (32 bits):



IEEE-754 Double precision (64 bits):



# IEEE-754 Single Precision (32-bit)



$$(00000000)_2 = (0)_{10}$$

$$(11111111)_2 = (255)_{10}$$

Reserve 0, 255 → special cases

$$\begin{aligned} & 1 \leq c \leq 254 \rightarrow 1 \leq m + \text{shift} \leq 254 \\ & \text{Set } [\text{shift} = 127] \rightarrow [-126 \leq m \leq 127] \quad m \in [-126, 127] \end{aligned}$$

# IEEE-754 Single Precision (32-bit)

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

$s \begin{cases} 0 \rightarrow (-1)^0 \Rightarrow \text{Positive} \\ 1 \rightarrow (-1)^1 \Rightarrow \text{Negative} \end{cases}$

Example: Represent the number  $x = -67.125$  using IEEE Single-Precision Standard

$$(67.125)_{10} = (1000011.001)_2 = (1\underbrace{000011001}_2 \times 2^6)_{10}$$

$m=6$   
 $c = m+127$   
 $c = (133)_{10}$

$$s = 0$$

$$f = \underbrace{000011001000 \dots 0}_{23 \text{ bits}}$$

$$c = (\underbrace{10000101}_2)_2$$

$\underbrace{\hspace{1cm}}_{3 \text{ bits}}$

$$0 \underbrace{1000101}_{1 \text{ bit}} \underbrace{0001101 \dots 0}_{23 \text{ bits}}$$

$\underbrace{\hspace{1cm}}_{8 \text{ bits}}$

# IEEE-754 Single Precision (32-bit)

$P=24$

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

8      23

- Machine epsilon ( $\epsilon_m$ ): is defined as the distance (gap) between 1 and the next larger floating point number.

$$(1)_{10} = \underbrace{1.000\dots00}_{23 \text{ bits}} \times 2^0$$

( - )

$$= \underbrace{1.000\dots01}_{23 \text{ bits}} \times 2^0$$

$0.000\dots01 \times 2^0$

$\epsilon_m = 2^{-23}$

$g_m = 2^{-n}$

- Smallest positive normalized FP number:

$$UFL = 2^L = 2^{-126} \approx 10^{-38}$$

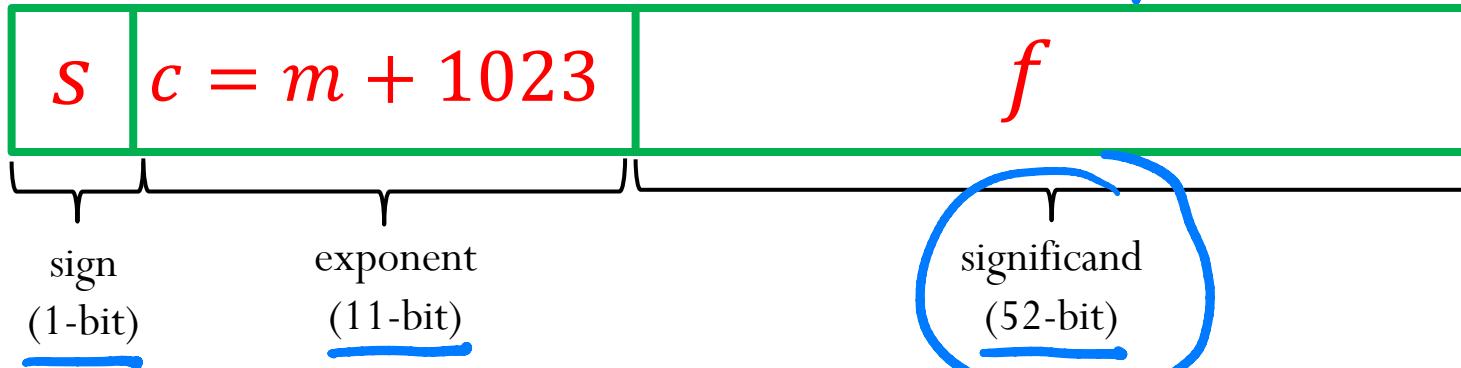
- Largest positive normalized FP number:

$$OFL = 2^{U+1} (1 - 2^{-P}) = 2^{128} (1 - 2^{-24}) \approx 10^{38}$$

# IEEE-754 Double Precision (64-bit)

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

$p=53$  ( $n+1$ )



$s = 0$ : positive sign,  $s = 1$ : negative sign

Reserved exponent number for special cases:

$$\begin{cases} c = (0000000000)_2 = 0 \\ c = (1111111111)_2 = 2047 \end{cases}$$

$$c = m + \text{shift}$$

Therefore  $1 \leq c \leq 2046$

$$1 \leq m + \text{shift} \leq 2046$$

$$\text{shift} = 1023$$

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

# IEEE-754 Double Precision (64-bit)

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

- **Machine epsilon ( $\epsilon_m$ )**: is defined as the distance (gap) between 1 and the next larger floating point number.

$$(1)_{10} = \boxed{0 \quad 0111 \dots 111 \quad 000000000000 \dots 000000000}$$

$$\epsilon_m = 2^{-n}$$
$$n = 52$$

$$(1)_{10} + \epsilon_m = \boxed{0 \quad 0111 \dots 111 \quad 000000000000 \dots 000000001}$$

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

- Smallest positive normalized FP number:

$$m \in [-1022, 1023] \text{ UFL} = 2^L = 2^{-1022} \approx 2.2 \times 10^{-308}$$

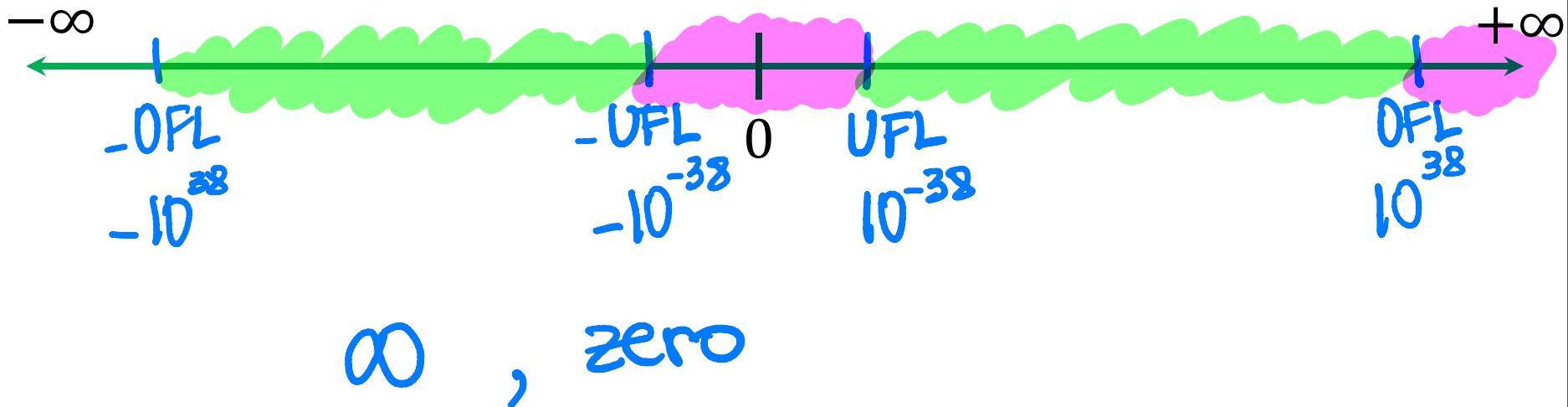
$$p=52$$

- Largest positive normalized FP number:

$$\text{OFL} = 2^{U+1}(1 - 2^{-p}) = 2^{1024}(1 - 2^{-53}) \approx 1.8 \times 10^{308}$$

$$U = 1023$$

# Normalized floating point number scale (single precision)



# Special Values:

$$c = (0000 \dots 0)$$

$$c = (111 \dots 1)$$

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

1) Zero:

$$x = \boxed{s} \boxed{000 \dots 000} \boxed{0000 \dots 0000}$$

$8, 11$        $23, 52$

2) Infinity:  $+\infty$  ( $s = 0$ ) and  $-\infty$  ( $s = 1$ )

$$x = \boxed{s} \boxed{111 \dots 111} \boxed{0000 \dots 0000}$$

$8, 11$        $23, 52$

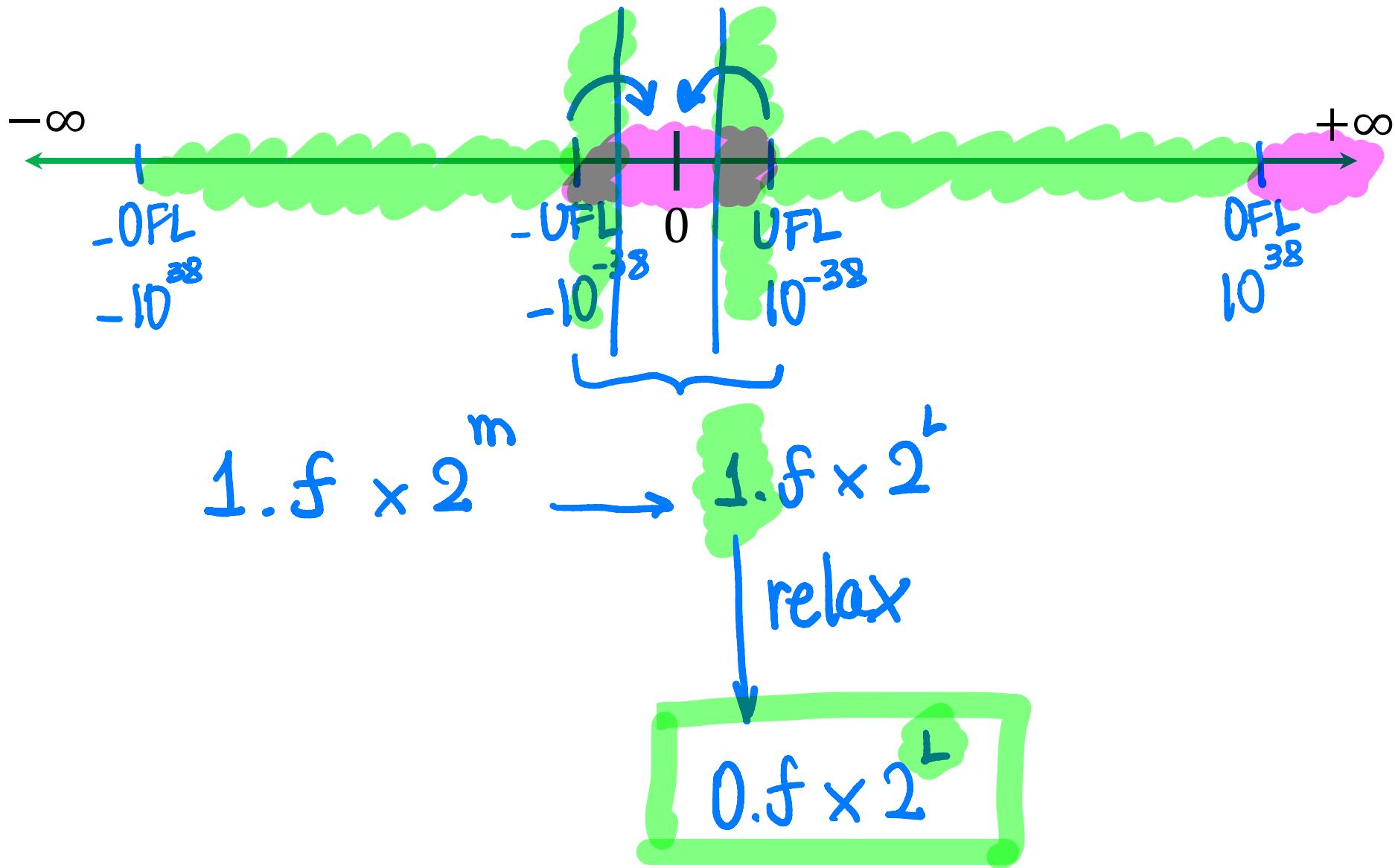
3) NaN: (results from operations with undefined results)

$$x = \boxed{s} \boxed{111 \dots 111} \boxed{\text{anything } \neq 00 \dots 00}$$

$8, 11$        $(100 \dots 010)$

4)  $c = (000 \dots 0)$      $f = (\text{anything}) \rightarrow \text{subnormal}$

# Normalized floating point number scale (single precision)



# Subnormal (or denormalized) numbers

- Noticeable gap around zero, present in any floating system, due to normalization
  - ✓ The smallest possible significand is 1.00
  - ✓ The smallest possible exponent is  $L$
- Relax the requirement of normalization, and allow the leading digit to be zero, only when the exponent is at its minimum ( $m = L$ )

$$x = (-1)^s \cdot f \times 2^L$$

Handwritten annotations:

- A green bracket encloses the entire expression  $x = (-1)^s \cdot f \times 2^L$ .
- A blue star is placed above the expression.
- A blue arrow points from the handwritten text "m = C-shift" to the exponent  $L$  in the formula.
- A blue circle contains the handwritten text "m = L".
- Below the formula, a blue arrow points from the handwritten text "C = (0000...) → subnormal" to the fraction part  $f$  of the formula.

# Subnormal (or denormalized) numbers

IEEE-754 Single precision (32 bits):

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

Exponent set to  $m = -126$

$$0.f \times 2^{-126}$$

Smallest positive subnormal FP number:

$$0.0000\dots 01 \times 2^{-126} = 2^{-23} \times 2^{-126} \approx 1.4 \times 10^{-45}$$

IEEE-754 Double precision (64 bits):

$$c = (0000000000)_2 = 0$$

Exponent set to  $m = -1022$

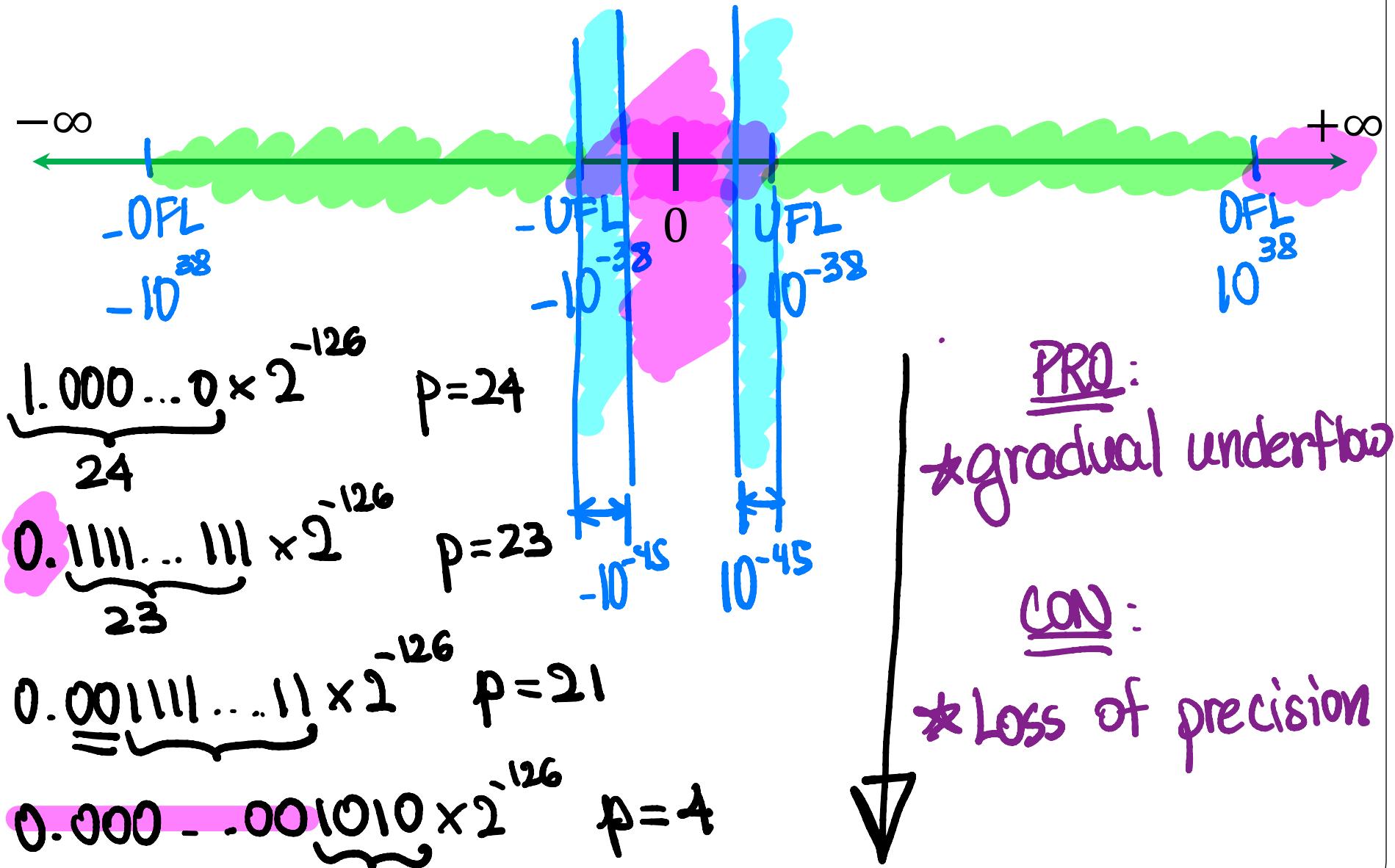
$$0.f \times 2^{-1022}$$

Smallest positive subnormal FP number:

$$0.000\dots 001 \times 2^{-1022} = 2^{-52} \times 2^{-1022} \approx 10^{-324}$$

52

# Normalized floating point number scale (single precision)



# Subnormal (or denormalized) numbers

Another special case:

$$x = \boxed{s \quad c = 000 \dots 000} \quad f$$

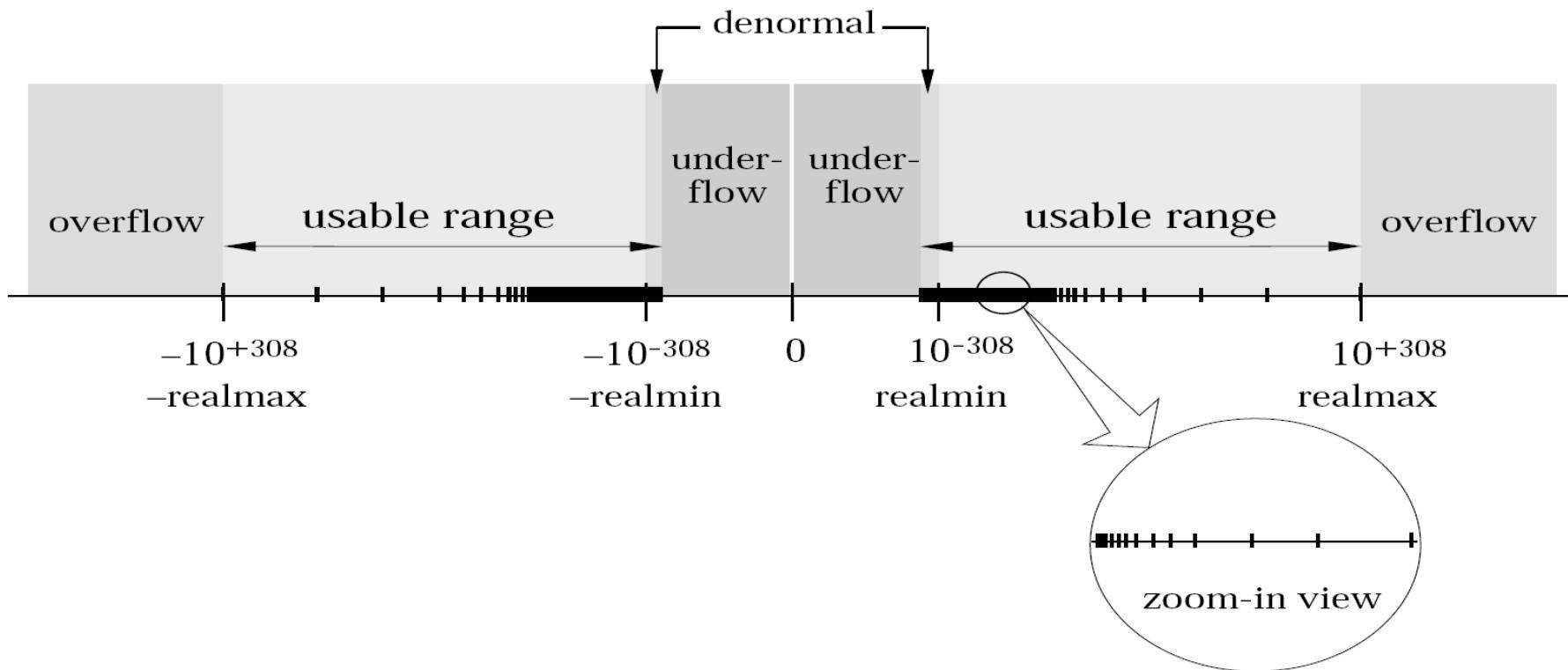
$$x = (-1)^s 0.\textcolor{red}{f} \times 2^{\textcolor{red}{L}}$$

Note that this is a special case, and the exponent  $\textcolor{red}{m}$  is **not** evaluated as  $\textcolor{red}{m} = \textcolor{red}{c} - \text{shift} = -\text{shift}$ .

Instead, the exponent is set to the lower bound,  $\textcolor{red}{m} = \textcolor{red}{L}$

- PROS: More gradual underflow to zero
- CONS: - Computations with subnormal numbers are often slow;
  - Loss of precision

# IEEE-754 Double Precision



# Summary for Single Precision

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

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