

# 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

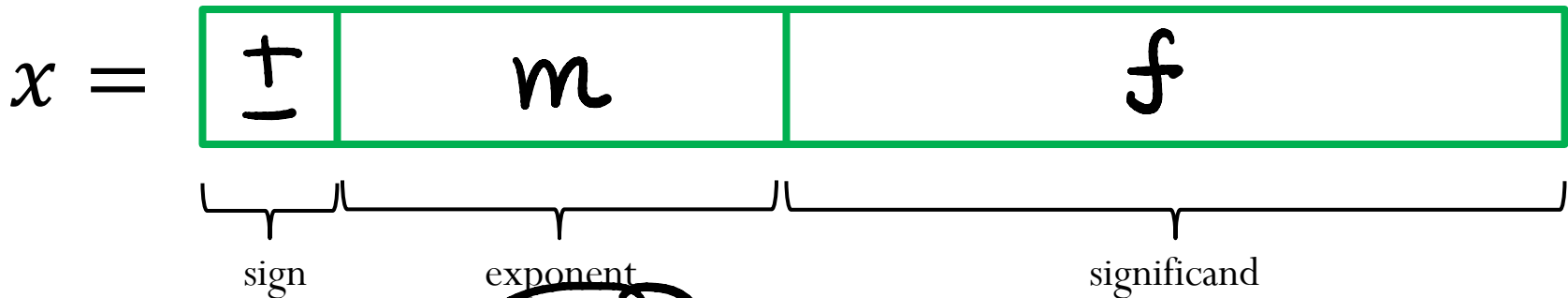
Numerical form:

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

Representation in memory:

$$m \in [L, U]$$

$$m \in [-4, 4]$$



$c = m + \text{shift}$

↓ unsigned int

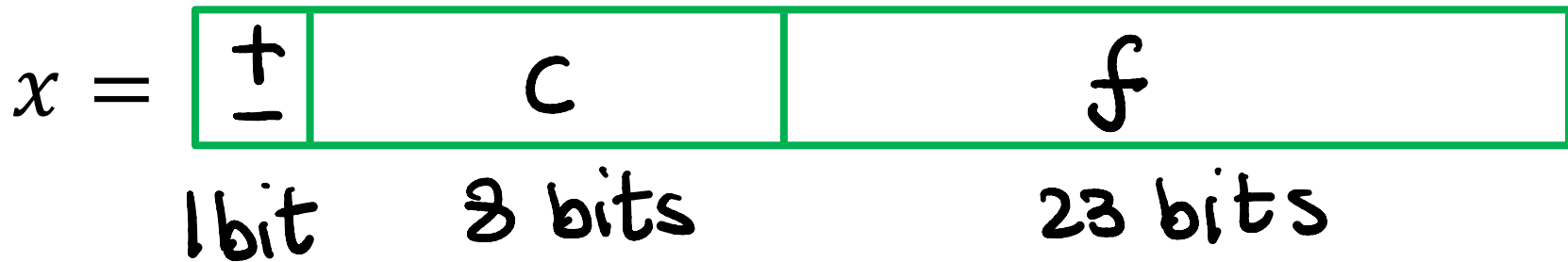
signed int

# Precisions:

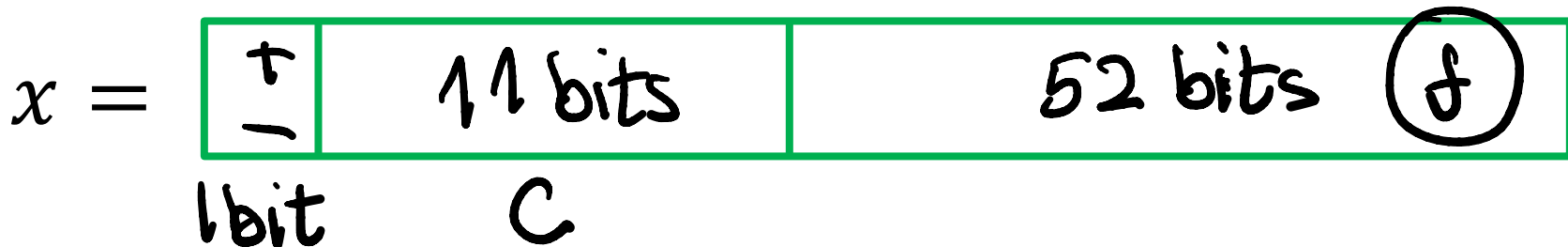
Finite representation: not all numbers can be represented exactly!

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

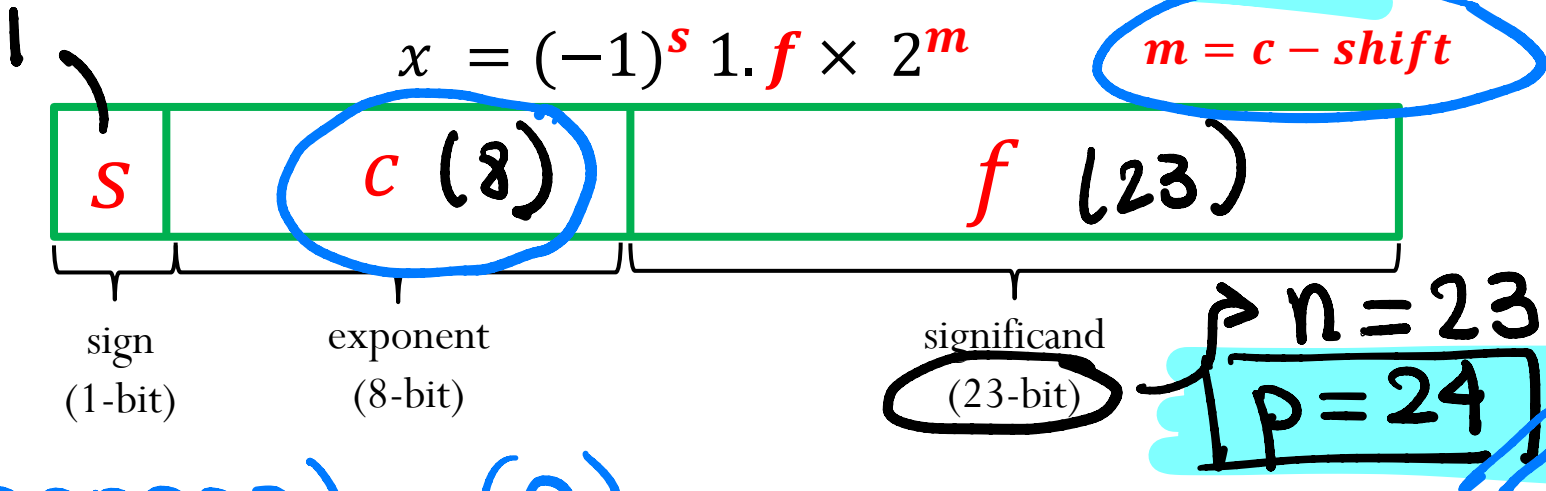
IEEE-754 Single precision (32 bits):  $C = m + \underline{\underline{\text{shift}}}$



IEEE-754 Double precision (64 bits):



# IEEE-754 Single Precision (32-bit)



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

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

$0 \leq c \leq 255$

reserve 0, 255 → special cases

$1 \leq c \leq 254 \rightarrow 1 \leq m + \text{shift} \leq 254$

$\underline{\underline{\text{set}}}$   $\text{shift} = 127 \rightarrow -126 \leq m \leq 127$ 
 $m \in [-126, 127]$

# IEEE-754 Single Precision (32-bit)

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

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

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

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

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

$$s = 0$$

$$f = \underbrace{000011001000\dots0}_{23 \text{ bits}}$$

$$C = (\underbrace{10000101}_2) \underbrace{01000101}_{8 \text{ bits}} \underbrace{0001101\dots0}_{23 \text{ bits}}$$

# IEEE-754 Single Precision (32-bit)

$p=24$

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

$\swarrow$  8
 $\swarrow$  23

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

$$\begin{aligned}
 (1)_{10} &= 1.\underbrace{000 \dots 00}_{23 \text{ bits}} \times 2^0 && 0.000 \dots 01 \times 2^0 \\
 & && \downarrow \\
 & && 2^{-23} \\
 & && \downarrow \\
 & && \epsilon_m = 2^{-23} \\
 & && \downarrow \\
 & && \epsilon_m = 2^{-n}
 \end{aligned}$$

$\ominus$

- Smallest positive normalized FP number:

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

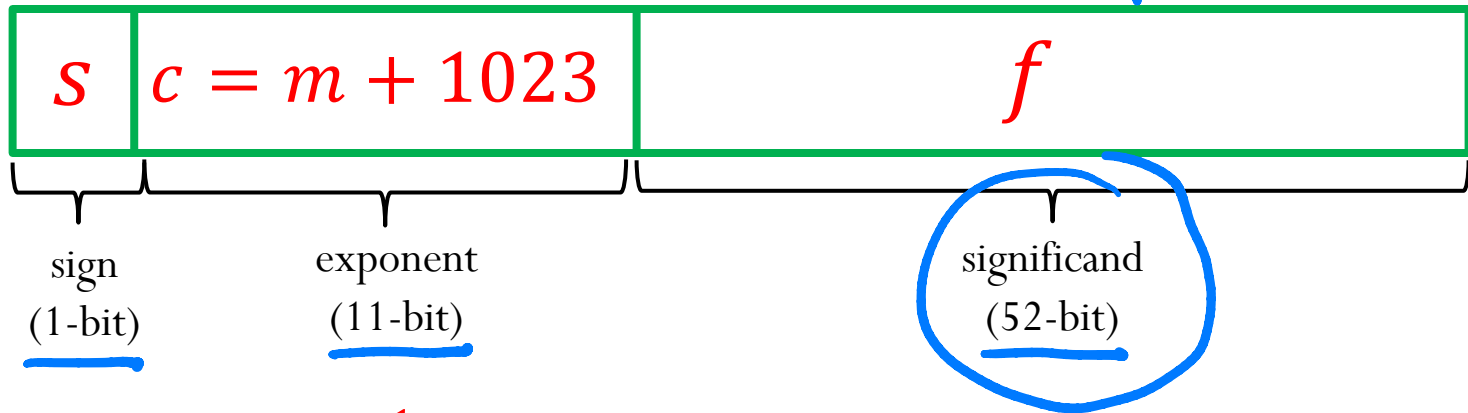
- Largest positive normalized FP number:

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

# IEEE-754 Double Precision (64-bit)

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

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



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

Reserved exponent number for special cases:

$$\begin{cases} c = (00000000000)_2 = 0 \\ c = (11111111111)_2 = 2047 \end{cases}$$

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

Therefore  $1 \leq c \leq 2046$

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

$$\text{shift} = 1023$$

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

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



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

$$(1)_{10} = \boxed{0 \quad 0111 \dots 111 \quad 000000000000 \dots 0000000000}$$

$$(1)_{10} + \epsilon_m = \boxed{0 \quad 0111 \dots 111 \quad 000000000000 \dots 0000000001}$$

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

$$n = 52$$

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

- **Smallest positive normalized FP number:**

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

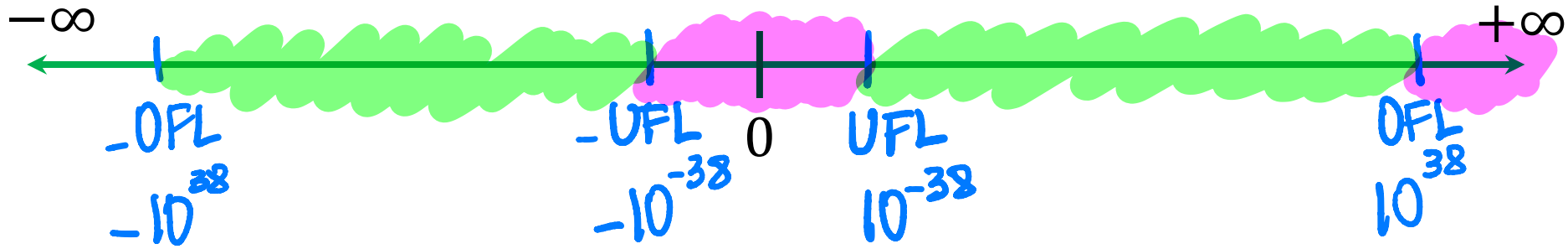
- **Largest positive normalized FP number:**

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

$$p = 52$$

$$U = 1023$$

# Normalized floating point number scale (single precision)



$\infty$  , zero

# Special Values:

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

$$c = (111 \dots 11)$$

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

1) Zero:

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

$8, 11$ 
 $23, 52$

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

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

$8, 11$ 
 $23, 52$

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

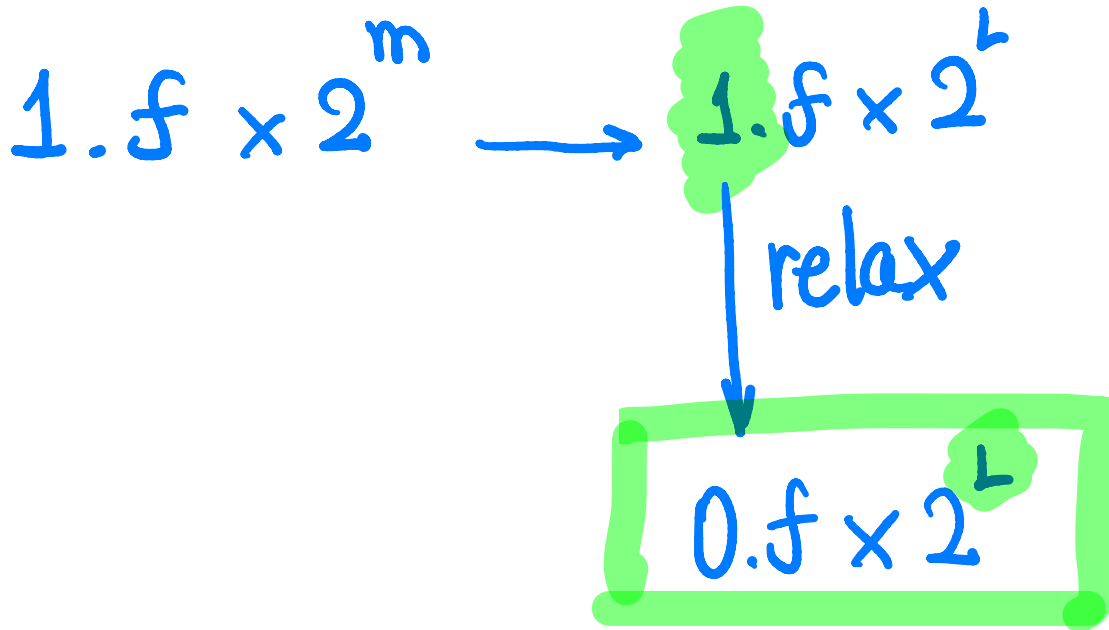
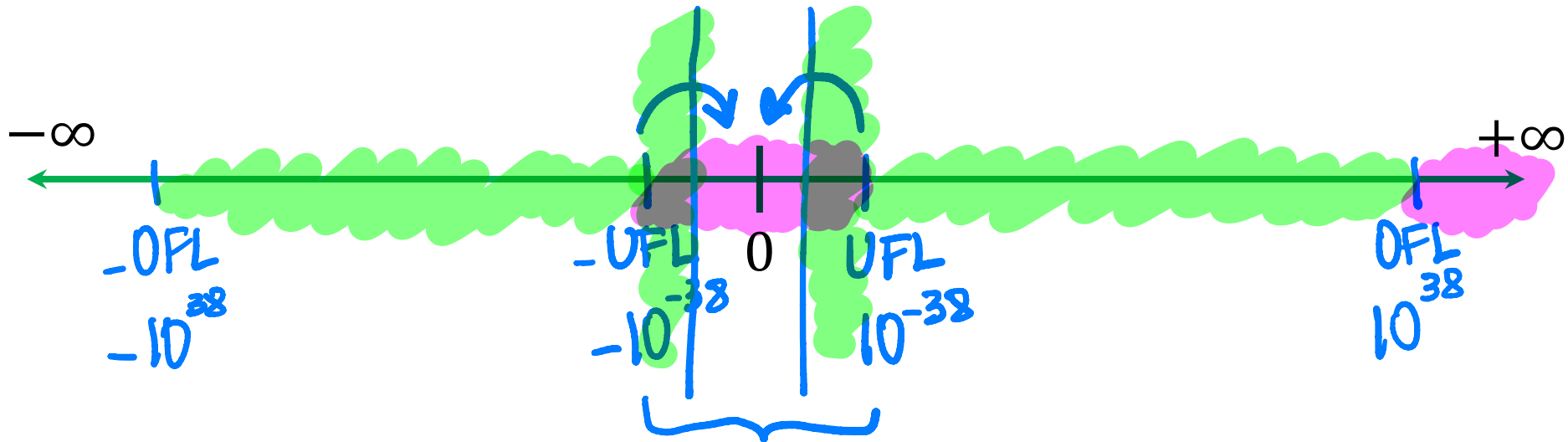
$$x = \boxed{s} \quad \boxed{111 \dots 111} \quad \boxed{\text{anything} \neq 00 \dots 00}$$

$8, 11$ 
 $(100 \dots 010)$

4)  $c = (000 \dots 0)$

$f = (\text{anything}) \rightarrow$  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 ( $m = L$ )

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

★

$m = C$  - shift

$m = L$

$C = (0000\dots) \rightarrow$  subnormal

# Subnormal (or denormalized) numbers

## IEEE-754 Single precision (32 bits):

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

Exponent set to  $m = -126$

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

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):

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

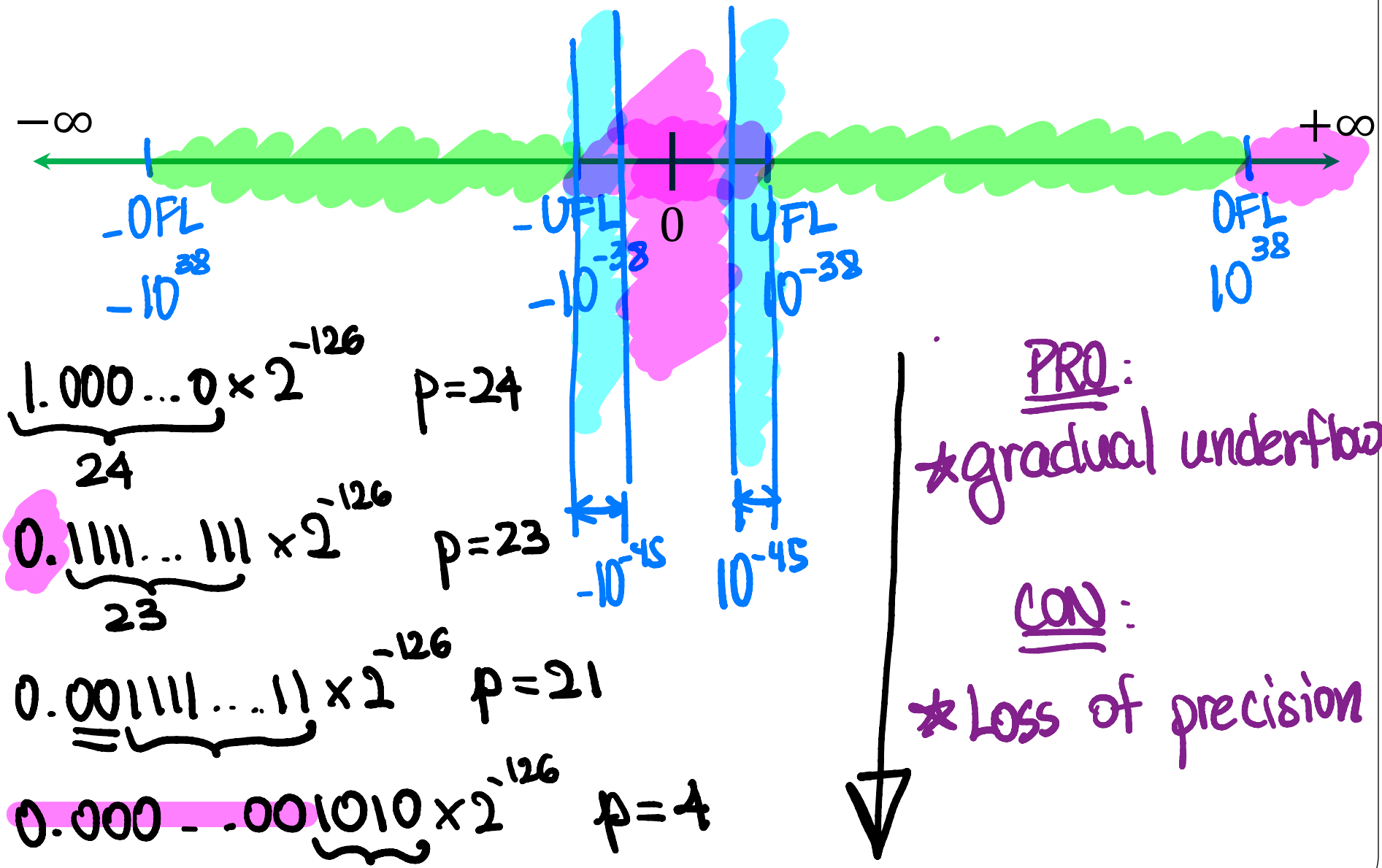
Exponent set to  $m = -1022$

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

Smallest positive subnormal FP number:

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

# Normalized floating point number scale (single precision)



# Subnormal (or denormalized) numbers

Another special case:

$$x = \boxed{\begin{array}{|c|c|c|} \hline s & c = 000 \dots 000 & f \\ \hline \end{array}}$$

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

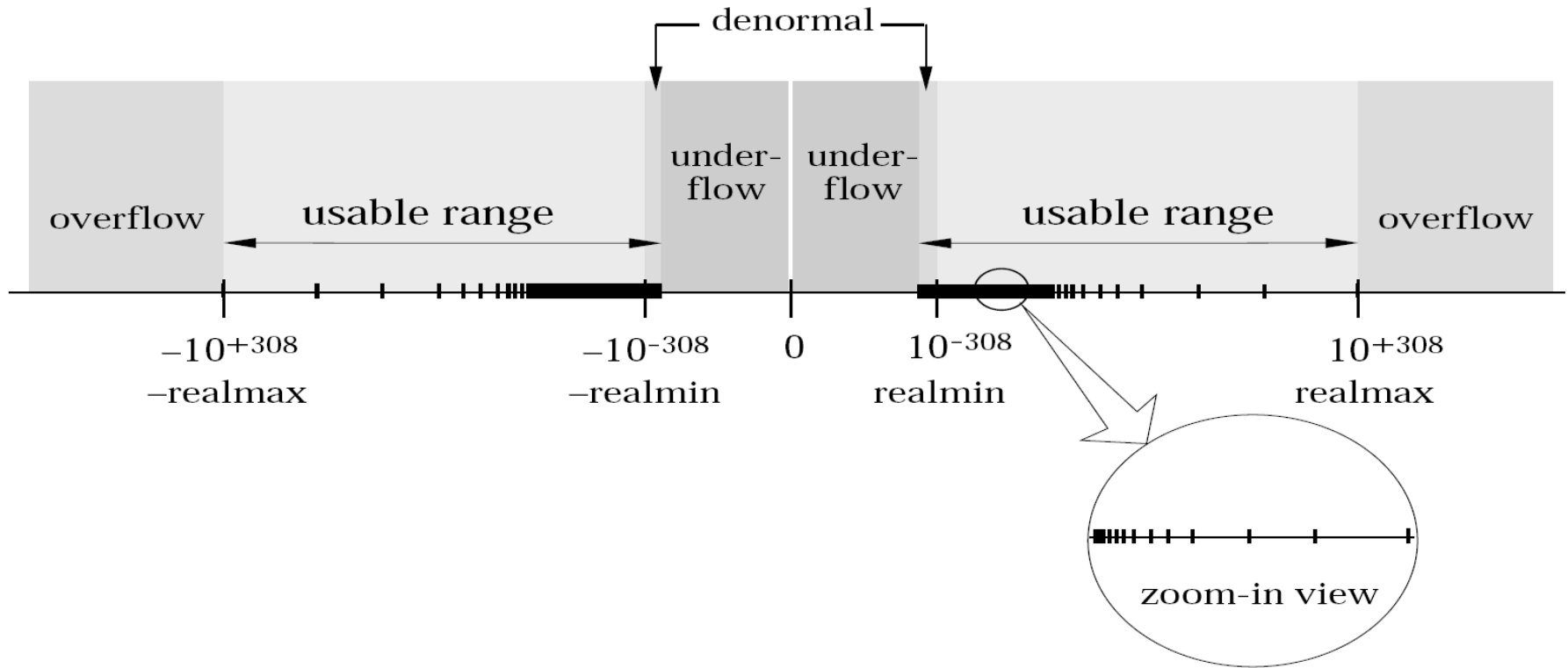
Note that this is a special case, and the exponent  $m$  is not evaluated as  $m = c - \text{shift} = -\text{shift}$ .

Instead, the exponent is set to the lower bound,  $m = L$

- 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

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

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