# Normalized floating-point numbers

$$x = \pm\, q \times 2^m = \pm\, 1.b_1 b_2 b_3 \ldots b_n \times 2^m = \pm\, 1.f \times 2^m$$

- **Exponent range**: $[L, U]$

- **Precision**: $p = n + 1$

- **Smallest positive normalized FP number:**

$$\text{UFL} = 2^L$$

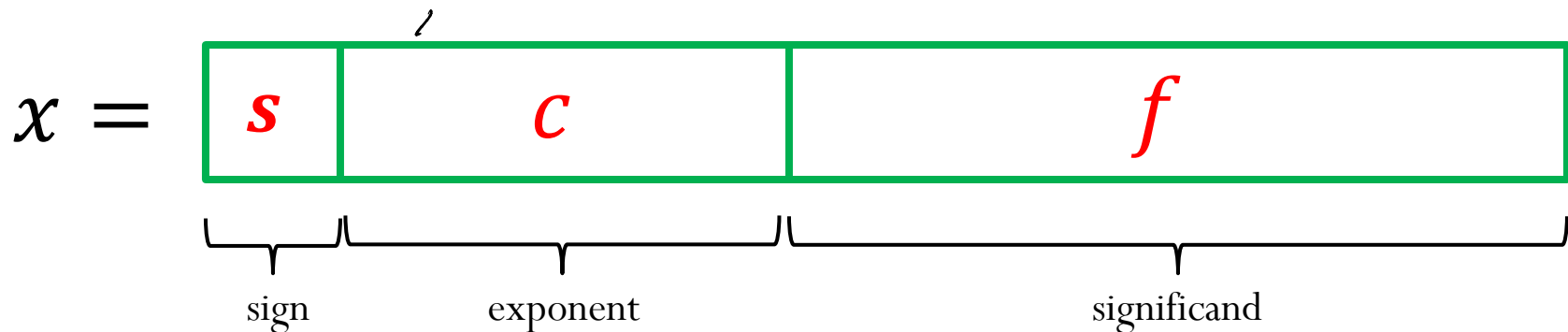- **Largest positive normalized FP number:**

$$\text{OFL} = 2^{U+1}(1 - 2^{-p})$$

# Floating-point number representation

**Numerical form:**

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

**Representation in memory:**

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

$\underbrace{\phantom{xxx}}_{\text{sign}} \quad \underbrace{\phantom{xxxxxx}}_{\text{exponent}} \quad \underbrace{\phantom{xxxxxxxxxxxx}}_{\text{significand}}$

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

$$m = c - shift$$

# Single Precision

$c \Rightarrow 8$ bits

$(000 \ldots 00)_2 \rightarrow (0)_{10}$

$(111 \ldots 11)_2 \rightarrow (255)_{10}$

$f \Rightarrow 23$ bits

$p = 24$

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

OFL : $2^{128}(1 - 2^{-24}) \approx 10^{38}$

$1 \leq c \leq 264$

$m = c - \text{shift}$

$\text{shift} = 127$

$-126 \leq m \leq 127$

$L$ 　 $U$

$\epsilon_m = 2^{-23}$

# Double Precision

$c = 11$ bits

$(000 \ldots 00)_2 \rightarrow (0)_{10}$

$(111 \ldots 11)_2 \rightarrow (2047)_{10}$

$f = 52$ bits

$p = n+1 = 53$

UFL : $2^{-1022} \approx 10^{-308}$

OFL : $2^{1024}(1 - 2^{-53}) \approx 10^{308}$

$1 \leq c \leq 2046$

$\text{shift} = 1023$

$-1022 \leq m \leq 1023$

$L$ 　 $U$

$\epsilon_m = 2^{-52}$

# Special Cases

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

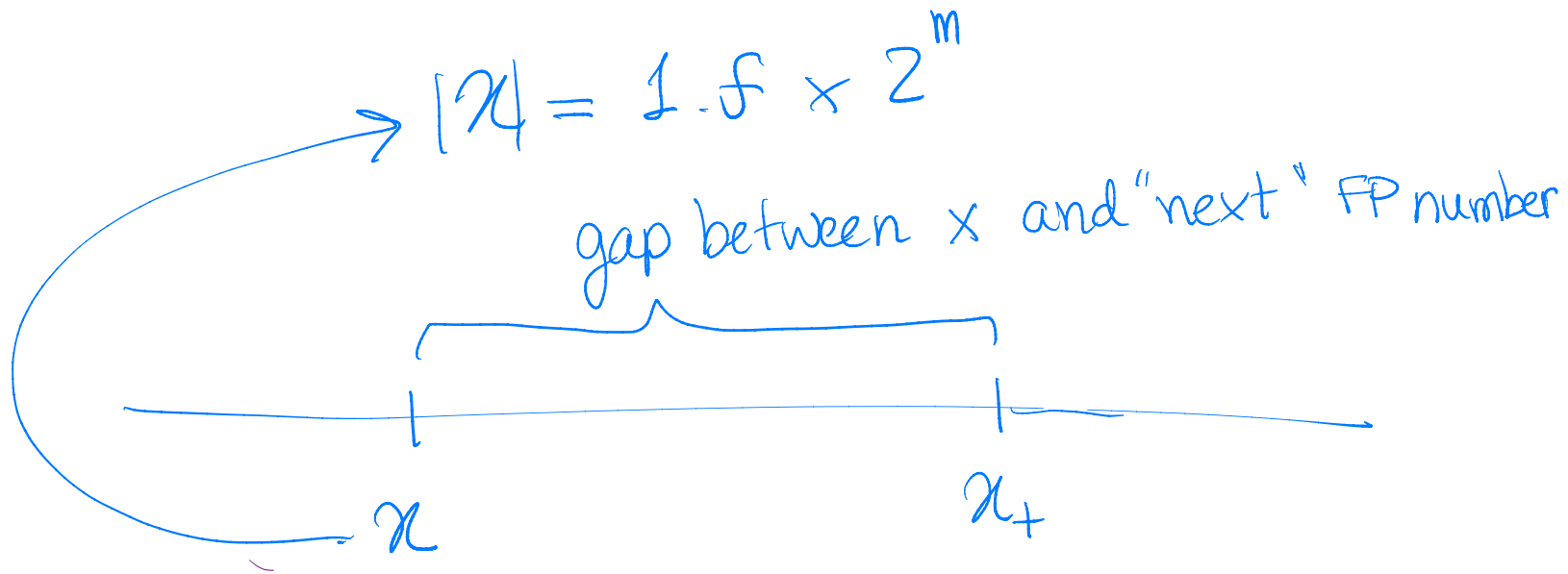| c | f | |
|---|---|---|
| All zeros | All zeros | $\pm 0$ |
| All zeros | $f \neq 0$ | Subnormal $x = \pm 0.f \times 2^L$ |
| All ones | All zeros | $\pm \infty$ |
| All ones | $f \neq 0$ | NaN |

# Normalized floating point number scale

$$0.00000 \sim \cdots \sim 001 \times 2^L$$

subnormal numbers

$0.f \times 2^L$

$2^{-n}$

smallest subnormal # is $2^{-n} \times 2^L$

$-\infty$             $+\infty$

$-$OFL

$-$UFL

0

UFL

OFL

overflow

overflow

$1 \cdot f' \times 2^L$

underflow

$$\rightarrow |x| = 1.f \times 2^m$$

gap between x and "next" FP number

$$x \qquad x_+$$

$$\text{gap} = |x - x_+| = \epsilon_m \times 2^m = e_a$$

Example: $2^4 \times 2^{-52} = 2^{-48}$ $(x = 2^4)$