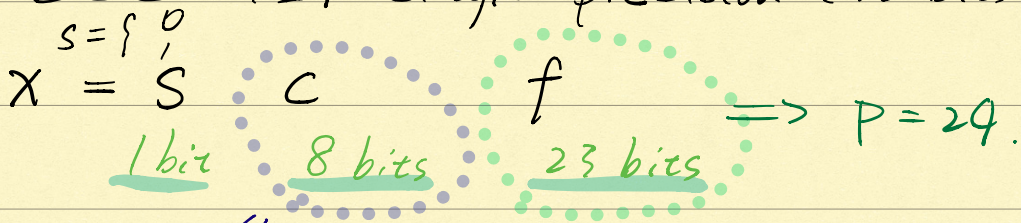Numerical Form: $x = (-1)^s \cdot 1.f \times 2^m$

$m = c - \text{shift}.$

Representation in memory: $x = \underbrace{\pm}_{\text{Sign}} \underbrace{m}_{\text{exponent}} \underbrace{f}_{\text{Significand}}$

$\underset{\substack{\text{unsigned int} \quad \text{Signed int.}}}{C = m + \text{shift}}$

# IEEE-754 Single precision (32 bits):

$s = \begin{cases} 0 \\ 1 \end{cases}$

$x = \underset{1 \text{ bit}}{S} \quad \underset{8 \text{ bits}}{c} \quad \underset{23 \text{ bits}}{f} \quad \Rightarrow P = 24.$

$(0\,0\,0\,0\,0\,0\,0\,0)_2 = (0)_{10}$

$(1\,1\,1\,1\,1\,1\,1\,1)_2 = (255)_{10}$

$C \in [0, 255].$

reserve $\underline{0}.\underline{255}$ for special cases.

$\Rightarrow C \in [1, 254]. \Rightarrow 1 \le m + \text{shift} \le 254.$

Set shift $= 127 \longrightarrow -126 \le m \le 127.$

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

Example: $x = -67.125$

$(67.125)_{10} = (1000011.001)_2 = (1.0000110001)_2 \times 2^6$

$S = 0, \quad f = \underbrace{0000110010000\cdots0}_{23 \text{ bits}}$

$m = 6, \quad C = m + \text{shift} = 6 + 127 = (133)_{10}$

$$C = (\underbrace{1\,0\,0\,0\,0\,1\,0\,1}_{8\text{ bits}})_2$$

$$\underset{\underset{1\text{ bit}}{\underbrace{\phantom{0}}}}{0}\ \ \underset{8\text{ bits}}{\underbrace{1\,0\,0\,0\,0\,1\,0\,1}}\ \ \underset{23\text{ bits}}{\underbrace{0\,0\,0\,0\,1\,1\,0\,0\,1\,0\,0\,0\cdots 0}}$$

# Machine Epsilon $(\varepsilon_m)$ : $\varepsilon_m = 2^{-23}$

UFL
Smallest positive normalized FP number : $2^{-126} \approx 10^{-38}$

Largest positive normalized FP number: $2^{127+1}\cdot(1-2^{-24})$
OFL
$$= 2^{128} - 2^{104}$$
$$\approx 10^{38}$$

# IEEE - 754 Double precision (64 bits) :

$$x = \pm \quad c \quad\quad f$$
$$\underset{1\text{ bit}}{\phantom{x}}\ \ \underset{11\text{ bits}}{\phantom{x}}\ \ \underset{52\text{ bits}}{\phantom{x}} \Rightarrow p = 53.$$

Set shift $= 1023 \Rightarrow c = m + 1023$.

$(0\,0\,0\,0\,0\,0\,0\,0\,0\,0\,0)_2 = (0)_{10}$

$(1\,1\,1\,1\,1\,1\,1\,1\,1\,1\,1)_2 = (2047)_{10}$

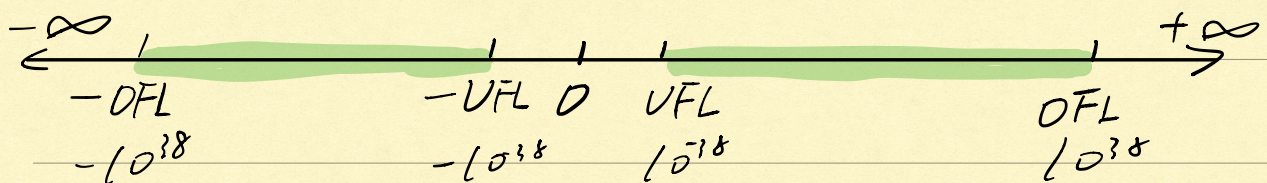$C \in [1, 2046]$.

$\Rightarrow \quad m \in [-1022, 1023].$

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

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

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

## Special Cases

Single precision



① represent $\underline{0}$ : $x = S \underbrace{00\cdots000}_{8 / 11 \text{ bits}} \underbrace{000 \text{---} 0000}_{23 / 52 \text{ bits.}}$

② represent $+\infty$ (S=0) and $-\infty$ (S=1)

$x = S \underbrace{111\cdots111}_{8/11 \text{ bits}} \underbrace{000 \text{---} 0000}_{23/52 \text{ bits.}}$

③ represent $\underline{NaN}$

$x = S \underbrace{111\cdots111}_{8/11 \text{ bits}} \underbrace{\text{anything} \neq 00\ldots00}_{23/52 \text{ bits.}}$

④ represent <u>subnormal numbers</u>

$$x = S \; \underline{000 \text{---} 000} \quad \underline{\text{anything} \neq 00 \text{---} 00}$$

8 / 11 bits     23 / 52 bits.

Smaller than UFL.

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

<u>Subnormal (or denormalized) numbers.</u>

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

IEEE-754 Single precision (32-bits)

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

$$m = -126$$

Smallest positive subnormal FP number:

$$0.000 \text{---} 1 \times 2^{-126} = 2^{-23} \times 2^{-126} = 2^{-149} \approx 1.4 \times 10^{-45}$$
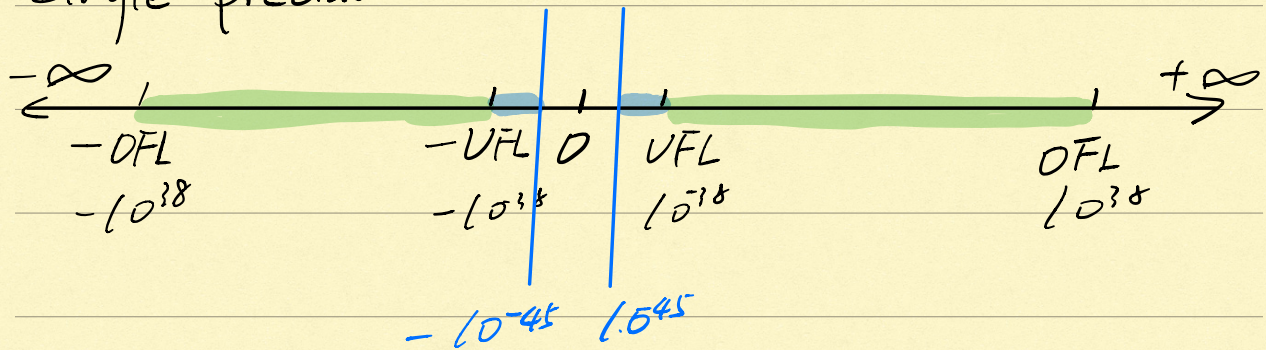
IEEE-754 double precision (64-bits)

$$c = (00000000000)_2 = (0)_{10}$$

$$m = -1022$$

Smallest positive subnormal FP number:

$$2^{-52} \times 2^{-1022} = 2^{-1074} \approx 10^{-324}$$

Single precision.



$$\underbrace{1.000 - - 0}_{24} \times 2^{-126} \qquad p = 24$$

$$0.\underbrace{111 - - 1}_{23} \times 2^{-126} \qquad p = 23$$

$$0.00 \cdot - - 0 \underbrace{1010}_{4} \times 2^{-126} \qquad p = 4$$

Overflow: $X : |X| > OFL \longrightarrow X = \pm OFL$

Gradual Underflow: $X : |X| < UFL \longrightarrow X = 0, \pm UFL$,

subnormal
number