# Discussion 1:

## Floating point representation:

$$\left(a_n \cdots a_2 a_1 a_0 . b_1 b_2 \cdots\right)_\beta = \sum_{k=0}^{n} a_k \beta^k + \sum_{k=0}^{\infty} b_k \beta^{-k}$$

decimal : $\beta = 10$

binary : $\beta = 2$

octal : $\beta = 8$

hexadecimal: $\beta = 16$

$$0 \leq a_i < \beta$$

In binary, define a floating point number $X$ as

$$X = \pm q \cdot 2^m$$

$q$ : mantissa, $q = (1.t)_2$ for normal numbers

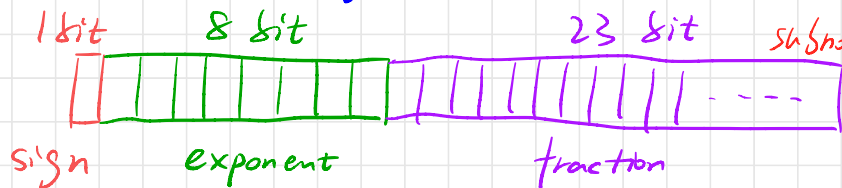$q = (0.t)_2$ for subnormal numbers
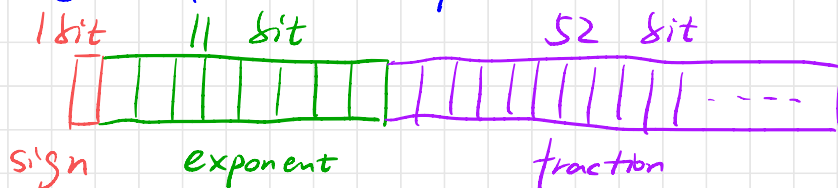
$m$ : exponent

eg. in decimal

normal : $6.5 \times 10^2$

$2.314 \times 10^{-5}$

subnormal : $0.065$

$0.54 \times 10^3$

## IEEE-754 Single precision:

1 bit      8 bit               23 bit

sign      exponent           fraction

## IEEE-754 double precision

1 bit      11 bit             52 bit

sign      exponent           fraction

Special cases:

Zero: when all exponent and fraction bits are 0. no restriction on the sign bit, so we have $+0$, $-0$

Infinity: all exponent bits are 1, all fraction bits are 0, $+\infty$ and $-\infty$ are distinguished from sign bit

NaN: "not a number" $(Inf-Inf, \frac{Inf}{Inf}, \frac{0}{0}, --)$. all exponent bits are 1, non-zero fraction bits

Subnormal numbers: all exponent bits are 0, non-zero fraction bits.

<span style="color:red">How do we define the exponent ?</span>

for float precision, what's the range can be represented by 8 bit?

smallest: $(00000000)_2 = 0$

largest: $(11111111)_2 = 2^0 + 2^1 + 2^2 + \cdots + 2^7 = 2^8 - 1 = 255$

the exponent is shifted by <span style="color:red">127</span> to avoid storing sign for exponent, as we saw above, $(00000000)_2$ and $(11111111)_2$ are reserved for special cases. So the actual range is $(1, 254)$, after shifting by 127, the exponent range is <span style="color:orange">$(-126, 127)$</span>

the largest normal number $\sim 2^{127} \sim 10^{38}$

-- smallest normal -- $\sim 2^{-126} \sim 10^{-38}$

for double precision, similar story:

   Smallest: $(00--0)_2 = 0$

           $\underline{11 \times 0}$

   largest: $(11---1)_2 = 2^{11} - 1 = 2047$

           $11 \times 1$

the actual useful range: $(1, 2046)$, shift by

$1023$ gives exponent range: $(-1022, 1023)$

   largest normal number $\sim 2^{1023} \sim 10^{307}$

   Smallest normal number $\sim 2^{-1022} \sim 10^{-308}$

Machine epsilon: the distance between 1 and the next largest floating point number.

Single precision:

$$1 = (+1) \times 2^0 \times (1.0000--0)_2$$

23th bit ↓ (pointing to last digit)

$$1 + \varepsilon_m = (+1) \times 2^0 \times (1.0000---1)_2$$

$$\text{So } \varepsilon_m = 2^{-23} \sim 10^{-7}$$

double precision:

$$1 = (+1) \times 2^0 \times (1.0000 \cdots 0)_2$$

$$1 + \varepsilon_m = (+1) \times 2^0 \times (1.0000 \cdots 1)_2$$

So   $\varepsilon_m = 2^{-52} \sim 10^{-16}$