# Computer representation of numbers (Overton, Sec 3 & 4)

Binary & decimal representation:

$$(71)_{10} = 7 \times 10^1 + 1 \times 1$$

$$(1000111)_2 = 1 \times 1 + 1 \times 2 + 1 \times 4 + 0 \times 8 + 0 \times 16 +$$
$$0 \times 32 + 1 \times 64$$

↗ $i$th position corresponds to $2^{i-1}$

$$(5.5)_{10} = 5 \times 1 + 5 \times \frac{1}{10}$$
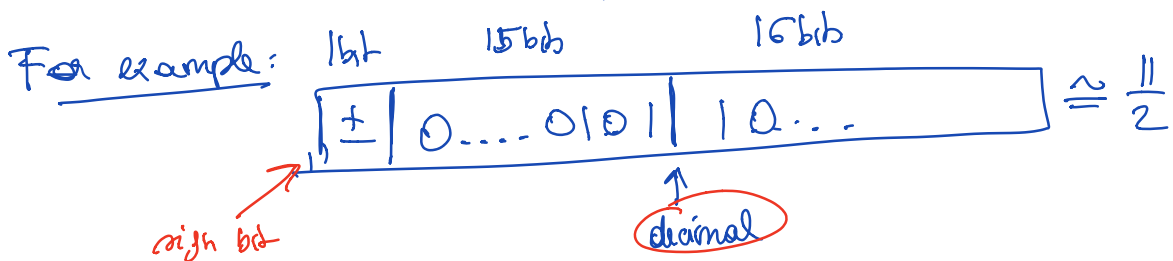$$(101.1)_2 = 1 \times 4 + 0 \times 2 + 1 \times 1 + 1 \times \frac{1}{2}$$

Some numbers have a finite representation with respect to one basis, but not another one:

$$\frac{1}{10} = (0.1)_{10} = (0.00011001100\ldots)_2$$

↗ finite     ↗ infinite

## Fixed point representation (base 2, binary)

Storage   32 bit  ↗ 4 byte   single precission
          64 bit  ↘ 8 byte   double precission   "default"

For example:   1 bit   15 bits        16 bits

$$\boxed{\pm \mid 0 \ldots 0101 \mid 10\ldots} \cong \frac{11}{2}$$

↗ sign bit           ↑ decimal

disadvantage : Some numbers cannot be represented at all if
they are too small or large.

## Floating point representation

base 10: $\qquad x = \pm S \times 10^{E} \qquad 1 \leq S < 10$

← integer, "exponent"

↗ mantissa

$$2.3456 \times 10^{9}$$

base 2: $\qquad x = \pm S \times 2^{E} \qquad 1 \leq S < 2 \ (\rightarrow S = 1...)$

sign bit

8 bits for $E$

23 bits for mantissa

$$S = (1. b_1 b_2 ..... )$$

IEEE standardized representation: don't store 1, "hidden bit"
in '70-80 s

round-off, dealing with $Inf, -Inf,$ NaN

# Single format

$$\boxed{\pm \mid a_1 \ldots a_8 \mid b_1 \quad \ldots \quad b_{23}}$$

biased representation
stored is $E + 127$

$$\left( \sum_{i=1}^{23} b_i 2^{-i} + 1 \right) \times 2^{\left[ \sum_{i=1}^{\infty} a_i 2^{8-i} - 127 \right]}$$

## smallest number:

$$\boxed{\pm \mid 0 \ldots 01 \mid 0 \ldots \ldots 0}$$

$$\simeq 1 \times 2^{-126} \approx 1.2 \times 10^{-38} \quad \text{realmin (single)}$$

## large number:

$$\boxed{\pm \mid 1 \ldots 10 \mid 1 \ldots \ldots 1}$$

$$\simeq 10^{38} \quad \text{realmax (single)}$$

## special:   $a_1 \ldots a_8 = 1$

$$\boxed{\pm \mid 1 \ldots 1 \mid 0 \ldots \ldots 0} \longrightarrow \pm \infty$$

$$\boxed{\pm \mid 1 \ldots 1 \mid 0 \ldots 1 \ldots 0} \longrightarrow \text{NaN}$$

$$\boxed{\pm \mid 0 \ldots 0 \mid 0 \ldots 0} \longrightarrow 0$$

Table 4.1: IEEE Single Format

$$\pm \mid a_1 a_2 a_3 \ldots a_8 \mid b_1 b_2 b_3 \ldots b_{23}$$

| If exponent bitstring $a_1 \ldots a_8$ is | Then numerical value represented is |
|---|---|
| $(00000000)_2 = (0)_{10}$ | $\pm(0.b_1 b_2 b_3 \ldots b_{23})_2 \times 2^{-126}$ |
| $(00000001)_2 = (1)_{10}$ | $\pm(1.b_1 b_2 b_3 \ldots b_{23})_2 \times 2^{-126}$ |
| $(00000010)_2 = (2)_{10}$ | $\pm(1.b_1 b_2 b_3 \ldots b_{23})_2 \times 2^{-125}$ |
| $(00000011)_2 = (3)_{10}$ | $\pm(1.b_1 b_2 b_3 \ldots b_{23})_2 \times 2^{-124}$ |
| $\downarrow$ | $\downarrow$ |
| $(01111111)_2 = (127)_{10}$ | $\pm(1.b_1 b_2 b_3 \ldots b_{23})_2 \times 2^{0}$ |
| $(10000000)_2 = (128)_{10}$ | $\pm(1.b_1 b_2 b_3 \ldots b_{23})_2 \times 2^{1}$ |
| $\downarrow$ | $\downarrow$ |
| $(11111100)_2 = (252)_{10}$ | $\pm(1.b_1 b_2 b_3 \ldots b_{23})_2 \times 2^{125}$ |
| $(11111101)_2 = (253)_{10}$ | $\pm(1.b_1 b_2 b_3 \ldots b_{23})_2 \times 2^{126}$ |
| $(11111110)_2 = (254)_{10}$ | $\pm(1.b_1 b_2 b_3 \ldots b_{23})_2 \times 2^{127}$ |
| $(11111111)_2 = (255)_{10}$ | $\pm\infty$ if $b_1 = \cdots = b_{23} = 0$, NaN otherwise |

Subnormals:

$$\boxed{\pm \mid 0 \cdots 0 \mid 1\, 0 \cdots \cdots \cdots 0}$$

$$\simeq (0.1)_2 \times 2^{-126} \simeq 2^{-127}$$

$$\boxed{\pm \mid 0 \cdots 0 \mid 0\, 0 \cdots \cdots \cdots 01}$$

$$= 2^{-149}$$

Double precission:   64 bit

$$\boxed{\pm \mid a_1 \cdots \quad a_{11} \mid b_1 \qquad \qquad b_{52}}$$

11 bit                      52 bit

<mark>Machine epsilon</mark>  $\big(eps(1)\big)$ : gap between 1 and
next largest number

$$\boxed{\pm \mid 0\, 1 \cdots 1 \mid 0 \cdots \cdots \cdots 0}$$

$$1.0 \cdots 0 \times 2^0$$

$$\longrightarrow \; eps = \underline{2^{-52} \approx 10^{-16}} \quad (double)$$

$$\approx 10^{-7} \qquad (single)$$