

Number Systems

Decimal \rightarrow base 10

\rightarrow digits $\{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$

Example of computation of decimal numbers :

524.503

$$= 5 \times 10^2 + 2 \times 10^1 + 4 \times 10^0 + 5 \times 10^{-1} + 0 \times 10^{-2} + 3 \times 10^{-3}$$

base 10

Binary \rightarrow base 2

\rightarrow digits $\{0, 1\}$

Example of binary computation :

101.011

$$= 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 + 0 \times 2^{-1} + 1 \times 2^{-2} + 1 \times 2^{-3}$$

base 2

Repeating numbers

decimal system examples : $\frac{1}{3} = 0.\overline{3333} = 0.\bar{3}$

$$\frac{1}{7} = 0.\overline{142857} \\ = 0.\bar{142857}$$

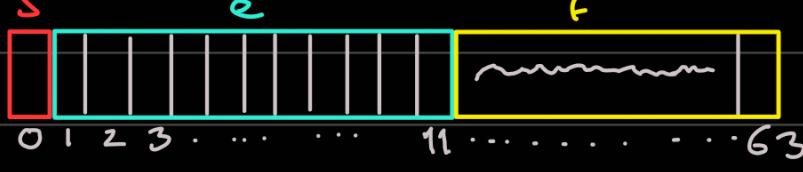
binary system examples : $\frac{1}{3} = 0.\overline{01}$

$$\frac{1}{7} = 0.\overline{001001\dots} \\ = 0.\bar{001}$$

These repeating numbers repeat indefinitely and are shortened with a bar over their heads

Double Precision

The IEEE arithmetic standard :



formula to calculate #s
 $= (-1)^S \times 2^{E-1023} \times F$

S : sign bit (either a 0 or 1)

e : biased exponent of 2 (contains 11 bits; first position bit, i.e. bit 1 gives $e = 2^{048}$)

f : significand (52 bits)

IEEE standard also has some special reserved numbers that are represented as :

(i) e is all 1's and

$f = 0$ we say the number is Inf (infinity)

or, $f \neq 0$ we say the number is NaN (not a number)

(ii) e is all 0's, then 1.f becomes 0.f (graceful underflow)

Calculating real max/real min

For real max e is almost all 1's (last e = 0 in bit position 11 to avoid infinity).

$$e \implies 1111111110$$

$$\therefore e = 1024 + 512 + 256 + 128 + 64 + 32 + 16 + 8 + 4 + 2 + 0$$

or, $e = 2^{11}-1-1$ [gives the sum of first 10 bits

subtracts the 1 from the last bit position since the bit is a 0]

and all the f's (all the 52 bits) are 1 \Rightarrow

$$f = \frac{1}{2} + \frac{1}{4} + \dots + \frac{1}{2^{51}} + \frac{1}{2^{52}} + 1?$$

or, $f = 1 + (1 - \frac{1}{2^{52}})$ sum upto bit position 52
f starts with bit position 0 as well?

the s=0 (+ve number)

$$\text{realmax} = (-1)^0 \times 2^{1023-1023} \times (2 - 2^{-52}) \\ = 1.7977e^{308}$$

Important

To understand why
 $f = \frac{1}{2} + \frac{1}{4} + \dots + \frac{1}{2^{51}} + \frac{1}{2^{52}}$
 is equal to $1 - \frac{1}{2^{52}}$
 assume we want
 first 3 decimal
 places of the binary
 system:
 $f = \frac{1}{2} + \frac{1}{4} + \frac{1}{8}$
 or, $f = 2^{-3}(4 + 2 + 1)$
 or, $f = 2^{-3}(2^3 - 1)$
 or, $f = 1 - \frac{1}{8} = 1 - \frac{1}{2^3}$

Double precision number formats

$$(1)_{10} = (-1)^0 \times 2^{1023-1023} \times 1.0 \quad (\text{both sides are in decimal})$$

binary for e $\rightarrow 1023 = 2^{10} - 1$ and should equal)

$$\text{or } (1023)_{10} = (10000000000)_2 - (00000000001)_2 \\ = (0111111111)_2$$

\therefore The binary/computer format of 1 is 0011111111 00...0
f's all 0

$$\left(\frac{1}{2}\right)_{10} = (-1)^0 \times 2^{1022-1023} \times 1.0$$

binary for e $\rightarrow 1022 = 2^9 + 2^8 + \dots + 2^1$
 $\therefore e = (011111110)_2$

computer representation of $\frac{1}{2} = (001111111100\dots0)_2$
f's all 0

$$\frac{1}{3} = \frac{1}{4} \left(1 + \frac{1}{3}\right)$$

$$\left(\frac{1}{3}\right)_{10} = (-1)^0 \times 2^{1021-1023} \times 1 \cdot \frac{1}{3}$$

$$\therefore e \rightarrow 1021 = 2^9 + 2^8 + \dots + 2^0 - 2 = (0111111101)_2$$

$$f \rightarrow \frac{1}{3} = 0.\overline{01}$$

$$\therefore \left(\frac{1}{3}\right)_{10} = (001111110101010101\dots01)_10$$

Signed vs. Unsigned

Unsigned → +ve numbers

$$\text{binary example : } (1000\ 1101)_2 = (141)_{10}$$

signed → accounting for -ve numbers

(0 for +ve / 1 for -ve)

$$\text{example : } 0100\ 0110 = (20)_{10} \text{ +ve}$$

$$1000\ 1101$$

$$= -1 \times 2^7 + 0 \times 2^6 + 0 \times 2^5 + 0 \times 2^4 + 1 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 1 \\ = (-115) \text{ -ve}$$



Word → 32 bits long

Half word → 16 bits long

Binary to unsigned

$$(-100)_{10}$$

$$\text{i) Convert } 100 \text{ to binary : } (0110\ 0100)_2$$

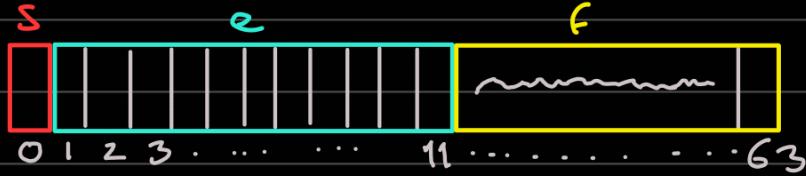
$$\begin{array}{l} \text{recall } \begin{cases} 1+0=1 \\ 1+1=10 \end{cases} \end{array}$$

$$\text{ii) Take 2's complement : } 1001\ 1011 \text{ flipping bits}$$

$$\begin{array}{r} + 1 \\ \hline \end{array} \text{ adding 1}$$

$$1001\ 1100 \rightarrow -100 \text{ in binary}$$

Eps



$$1 \text{ is represented as : } (-1)^0 \times 2^0 \times 1.0$$

$$\text{The next largest number is : } (-1)^0 \times 2^0 \times (1 + 2^{-52}) \quad [\text{last bit of } f \text{ is a } 1] \\ = (1 + 2^{-52})$$

$$\text{machine eps} = (1 + 2^{-52}) - 1 = 2^{-52}$$

In a similar fashion :

$$\text{next largest number after } 2 : (-1)^0 \times 2^1 \times (1 + 2^{-52}) = 2(1 + 2^{-52})$$

∴ difference between 2 and next largest number

$$= 2(1 + 2^{-52}) - 2 = 2 \times 2^{-52}$$