



Base 2 : used by computer to represent integers & real numbers

i.e. 2 digits $\rightarrow 0 \text{ \& } 1$ [powers of 2]

Positive Integers : represented directly by converting decimal value into power of 2

Ex: decimal number 5 \rightarrow $\begin{matrix} 8 & 4 & 2 & 1 \\ 0 & 1 & 0 & 1 \end{matrix}$

Revised

$$0101 \rightarrow 0 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0$$

$$\rightarrow 4 + 1$$

$$= 5$$

Negative Integers : typically represented using a technique called 2's complement

Ex: -5 in 8-bit system

1 byte = 8 bits $-128 \text{ to } +127$

1 1 1 1 1 0 1 1

2's complement

Step 1

5 in 8-bit system = 00000101

128	64	32	16	8	4	2	1
0	0	0	0	0	1	0	1

Step 2

1's complement [flip all the bits i.e. $1 \rightarrow 0$
 $0 \rightarrow 1$]

MSB

1 --- ve
 0 --- +ve

left most bit
 MSB \rightarrow

0	0	0	0	0	1	0	1
1	1	1	1	1	0	1	0

Step 3 Add 1 bit to the 1's complement

$$\begin{array}{r} 1 \ 1 \ 1 \ 1 \ 0 \ 1 \ 0 \\ \hline 1 \ 1 \ 1 \ 1 \ 0 \ 1 \ 1 \end{array}$$

Verification (Reverse)

c1 $11111011 \rightarrow -5$

c2 Subtract 1 bit: 11111011

left most bit
m.s.B \rightarrow $\boxed{1}1111010$

c3 flip bits: $\begin{array}{r} 11111010 \\ \hline 00000101 \end{array}$

$$\begin{array}{r} 8421 \\ 0101 \\ \hline 5 \end{array}$$

Real Numbers

5.625 \rightarrow IEEE754 format (binary representation)

\uparrow

used for floating point arithmetic representation

for both single precision (32 bit) & double precision (64 bits)

Step 1:

Convert integer part to binary

$$\begin{array}{r} 8421 \\ 101 \end{array}$$

$$5_{10} = 101_2$$

Step 2:

Convert fractional part to binary

$$5.625 \rightarrow 0.625$$

0.625 to binary, repeatedly multiply by 2 & note the integer part

$$\begin{array}{rcl}
 0.625 \times 2 & = & 1.25 \text{ --- (1)} \\
 0.25 \times 2 & = & 0.5 \text{ --- 0} \\
 0.5 \times 2 & = & 1.0 \text{ --- (1)} \\
 0 \times 2 & = & 0 \text{ --- Stop here}
 \end{array}$$

$$\therefore 0.625 \rightarrow 0.101$$

Step 3: Combine both integer & fractional part

$$5.625_{10} = 101.101_2$$

Step 4: Normalize binary number in scientific notation $\rightarrow 1.01101 \times 2^2$ (exponent)

$$101.101 \rightarrow 1.01101 \times 2^2$$

- 1) move toward fraction & keep only 1 digit before point
- 2) the exponent will be 2 cause we moved 2 bit

Step 5: Now represent in IEEE 754 format (Single precision)

float f = (float) 5.625 32 bit
double d = 5.625 64 bit

$$\text{Sign} \cdot \text{exponent} \cdot \text{mantissa}$$

bits of IEEE 754

→ Sign bit : 0

→ exponent = 2 + 127 = 129 $\leftarrow (10000001)_2$

Single precision : 127

double precision : 1023
& 11 bits

→ Mantissa = The fractional part from step 4

$$i.e. \quad 0.10100000 \dots 0 \quad \text{clay padded} \quad 23 \text{ bits}$$

128	64	32	16	8	4	2	1
1	0	0	0	0	1	0	1

5.625

32 bit

Sign

exponent

fraction

23

0

10000001

01101000000000000000000

18

Mantissa

32 (single prec) : 23

64 (double) : 52
prec

4-bit 5.25 — Binary format

$$0 \quad (00000) \quad 011010000 - \dots - 0$$

4)

Reverse

$$0 \quad (00000) \quad 011010000 - \dots - 0$$

5 18
23

1 bit 8 bit 2

Convert back to real number

$$\text{Exponent} = 128 - 127 = \textcircled{2} \quad (1.01101) \times 2^2$$

Normalized mantissa: 1.01101×2^n

$$\frac{101.101}{2} = ()_{10}$$

binary fraction 0.101_2

→

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

$$= 1 \times \frac{1}{2} + 0 \times \frac{1}{4} + 1 \times \frac{1}{8}$$

$$= 0.5 + 0 + 0.125$$

$$= \underline{\underline{0.625}}$$

$$0101.101_2 = \underline{\underline{5.625}}$$

$$\begin{array}{r} 8421 \\ \hline 5 \end{array}$$

Bit wise Not (!)

1) Flip 1st bit

2) change the sign

3) Interpret 2's complement

→ flip the bit (1's complement)

2) add extra bit

$$(12)_{10} = -$$

→ 8-bit system

Step 1:

128	64	32	16	8	4	2	1
0	0	0	0	1	1	0	0

Step 2:

most -ve

1	1	1	1	0	0	1	1
---	---	---	---	---	---	---	---

Step 3: 2's complement

1 1 1 1 0 0 1 1

(1) 1's complement [flip bits]

1	1	1	1	0	0	1	1
<hr/>							
0	0	0	0	1	1	0	0

(2) Add 1 bit

0	0	0	0	1	1	0	0
							1
<hr/>							
0	0	0	0	1	1	0	1
							←
							8 4 2 1
 (-13)							