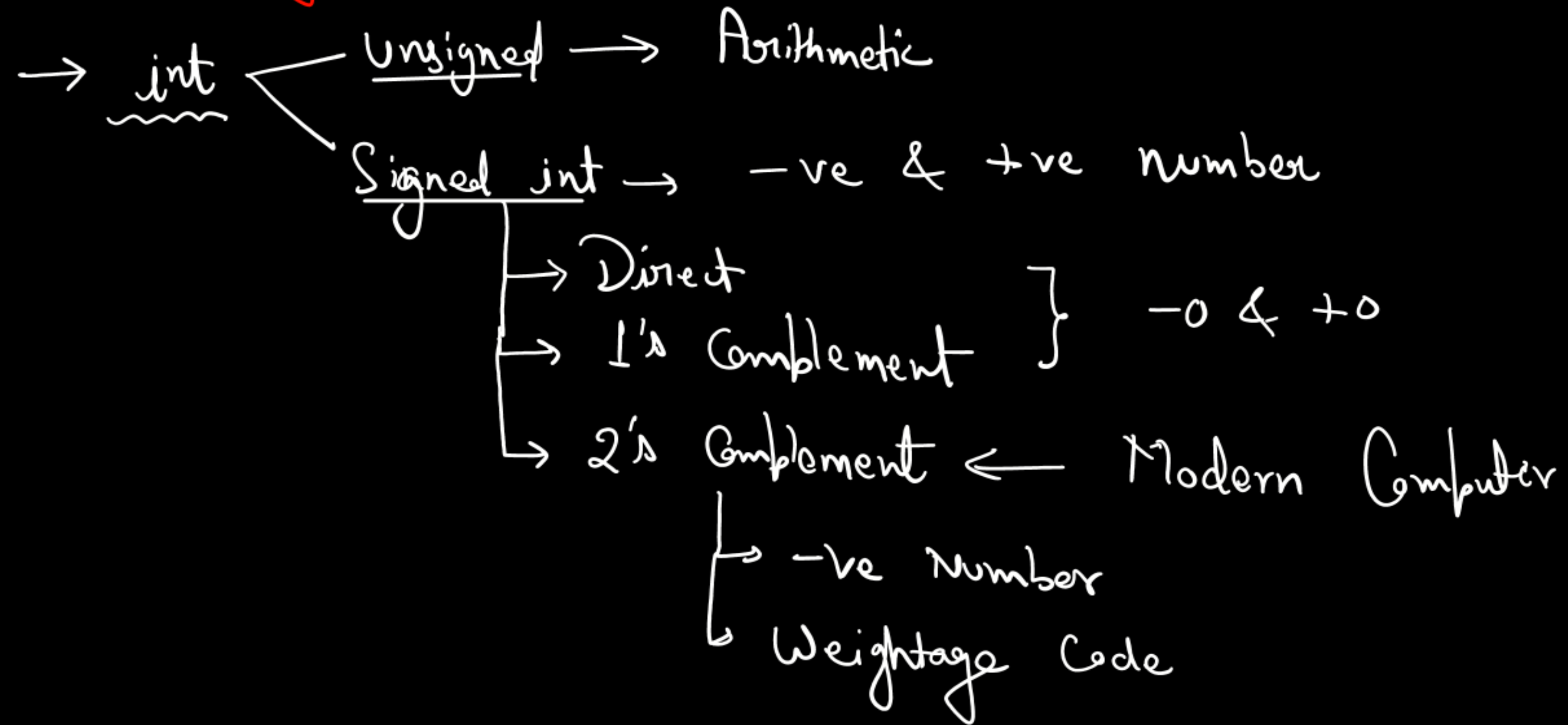


Storing Real Numbers



→ float Numbers: The number that consists point

Ex 3.14159, 2.718, 6.022×10^{-23}

float

- fixed point → The location of point is fixed Ex 1.758, 37.95
 - floating point → The location of point floats from left to right or vice-versa
↳ opposite
- Ex 6.022×10^{-23} , 1.75×10^4 ← Scientific Notation
 60.375×10^{-5}

(A) Fixed Point Representation

If fixed point representation of a number is $\overbrace{IIII}^{\text{Integral Part}} . \overbrace{FFF}^{\text{Fractional Part}}$

So we can store minimum value 0000.0001 & Maximum Value
will be 9999.9999

Three Parts

(a) Sign bit (Sign field)

b) Integral part (integer field)

c) Fractional part (fraction)

Signed Number

* Bits are fixed to store fixed Point Number

→ integral part
→ fractional part
→ sign

Ex: $-45.5 \rightarrow 0.5$



Sign bit

Integral Part

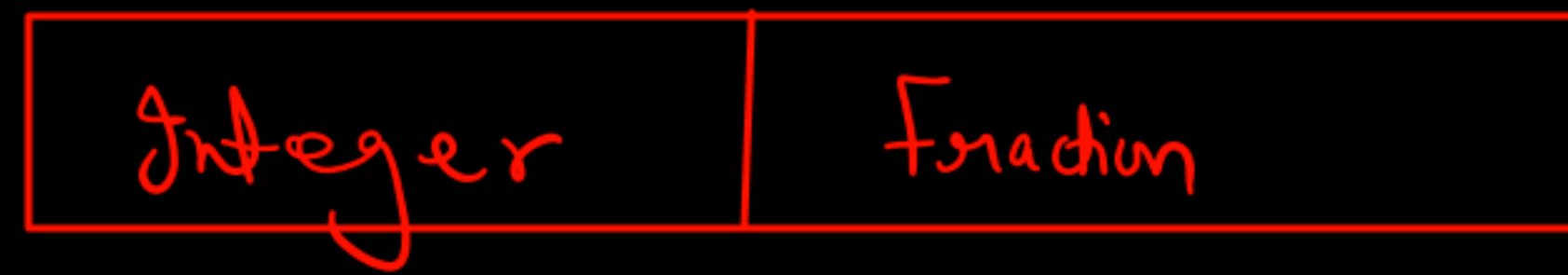
Fractional Part

1 -ve
0 +ve

12 bits fixed Point Number

→ 1 bit ← sign
→ 6 bits ← integral
→ 5 bits ← fractional

Unsigned Number Representation \leftarrow No Sign bit



Range in fixed point representation

for 'k' bits :

Signed Representation : $-(2^{k-1} - 1)$ to $(2^{k-1} - 1)$ \leftarrow Do Not use

Unsigned Representation : $(0 \text{ to } 2^k - 1)$

Signed 2's Complement Representation :

(-2^{k-1}) to $(2^{k-1} - 1)$

(B) Floating Point Representation

$$\Rightarrow (5.625)_{10} \longrightarrow (101.101)_2$$

$$\begin{array}{c} \downarrow \\ 0.5625 \times 10^1 \end{array}$$

Mantissa Radix

Exponent

$$\begin{array}{c} \downarrow \\ 0.101101 \times 2^3 \end{array}$$

Mantissa Radix

Exponent

Representation \Rightarrow

Sign		
S	Exponent	Mantissa
0	011	101101

$$\Rightarrow (0011101101)_2$$

$$\begin{array}{l} (101.101)_2 \\ \left\{ \begin{array}{l} 1.01101 \times 2^2 \\ 0.101101 \times 2^3 \end{array} \right\} \\ \Downarrow \\ \text{Normalization} \end{array}$$

Normalization : आसान तरीका

a) Explicit Normalization : Move the radix point to LHS of the Most Significant '1' in the Sequence

$$(101.101)_2 \rightarrow 0.101101 \times 2^3$$

b) Implicit Normalization : Move the radix point to the RHS of the Most Significant '1' in the Sequence

$$(101.101)_2 \rightarrow 1.01101 \times 2^{(2)} \leftarrow \text{Exponent}$$

Mantissa

Biasing: जितने Bits का Exponent होता है उसकी Range के सबसे बड़े Number को हमें Add करना होता है Exponent के साथ ताकि Exponent +ve हो जाए,

$$(000101)_2 \rightarrow 1.01 \times 2^{(-3)} \leftarrow \text{Negative}$$

2's Complement form \rightarrow Range $\Rightarrow (-2^{k-1})$ to $(2^{k-1}-1)$
 Suppose we have 4 bit Exponent

4 bits = 2^4
 -2^{4-1} to $2^{4-1}-1$
 $= -2^3$ to 2^3-1
 $= [-8 \text{ to } +7]$

largest value (ignore sign)

-8	-7	-6	-5	-4	-3	-2	-1	0	1	2	3	4	5	6	7	} 16 combinations
+8	+8	+8	+8	+8	+8	+8	+8	+8	+8	+8	+8	+8	+8	+8	+8	
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	

Without Bias = 1.01×2^{-3}

With Bias = $1.01 \times 2^{-3+8}$
 $= 1.01 \times 2^5 \leftarrow +ve$

Converting the Number into actual form:

a) Explicit Normalization : $(-1)^{\delta} \times 0.M \times 2^{\text{Exponent} - \text{Bias}}$, $\delta \in \{0, 1\}$

b) Implicit Normalization : $(-1)^{\delta} \times 1.M \times 2^{\text{Exponent} - \text{Bias}}$

$$-1^0 = 1$$

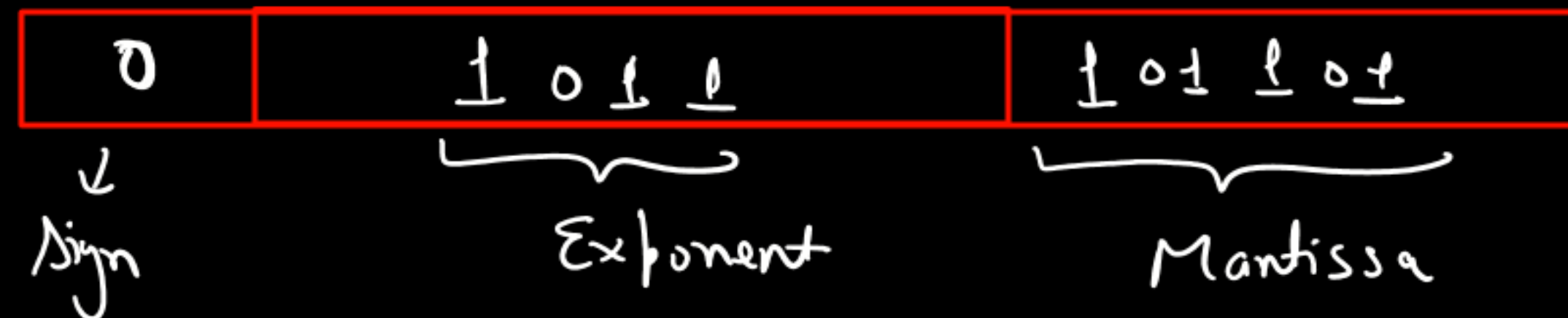
$$-1^1 = -1$$

Ex $(101.101)_2 \rightarrow$ Explicit : 0.101101×2^3

$$\begin{aligned} &\rightarrow 0.101101 \times 2^{3+8} \\ &= 0.\underline{101101} \times 2^{11} \end{aligned}$$

Ex Exponent = 4 bits

4 bits = 16
 \downarrow
 -2^{k-1} to $2^{k-1}-1$
 $= -8$ to 7



$\Rightarrow (0 \ 1011 \ 101101)_2$

$(11)_2 \rightarrow 1011$

Implicit Mode: $(\underline{101}, 101) \rightarrow 1.01101 \times 2^2$

Adding Bias = $1.01101 \times 2^{2+8}$
 $= 1.01101 \times 2^{10} \leftarrow \text{Biased Exponent}$

Exponent = 4 bits
 $\hookrightarrow 16 \text{ Comb.}$
 $2^k \text{ Comb.} = -2^{k-1} \text{ to } 2^{k-1}-1$
 $\hookrightarrow -8 \text{ to } +7$

$(10)_{10} \rightarrow (1010)_2$



\downarrow
 Sign
 \downarrow
 1 bit

$\underbrace{\hspace{2cm}}$
 Exponent
 \downarrow
 4 bit

\downarrow
 5 bits

$\Rightarrow \text{total} = 10 \text{ bits}$

o r o

Suppose \rightarrow Exponent $\rightarrow 10 \text{ bits} \rightarrow \underbrace{000000}_{\text{Bit Extension}} 1010$

Mantissa $\rightarrow 10 \text{ bits} \rightarrow 01101 \underbrace{00000}_{\text{bit Extension}}$

bit Extension

IEEE Standard:

↳ Precision

→ Half Precision: (16 Bit)

↳ 1 bit → Sign

5 bits → Exponent

10 bits → Mantissa

→ Single Precision (32 bit)

↳ 1 bit → Sign

8 bit → Exponent

23 bits → Mantissa

→ Double Precision (64-bit)

↳ 1 bit = Sign

11 bits = Exponent

52 bits = Mantissa

→ Quadruple Precision (128 bit)

↳ 1 bit → Sign

15 bit ⇒ Exponent

112 bit ⇒ Mantissa