

Weighted vs nonweighted
fixed point codes

weighted

more weighted

High-magnitude \rightarrow integer \rightarrow fraction

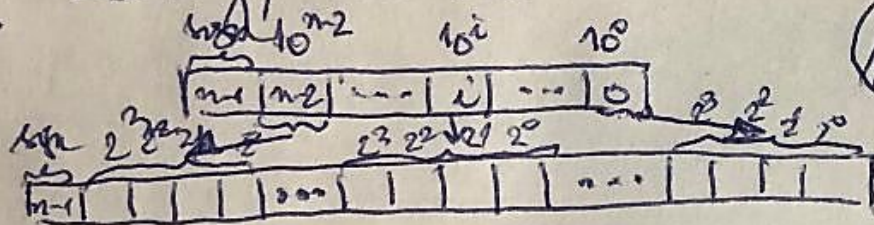
Diagram illustrating the structure of a floating-point number (IEEE 754 standard):

- Integer part (Left):** Consists of a sign bit (1 bit) and a magnitude field (n-1 bits). The magnitude is represented as 2^{n-1} .
- Fraction part (Right):** Consists of a sign bit (1 bit) and a magnitude field (n-1 bits). The magnitude is represented as 2^{1-n} .

True's complement
Two's complement

Three excess
Two - net of - Five

Binary to decimal


$$10^i \cdot 2^j$$

Binary vs Decimal

$$n=10 \rightarrow 2^{10} = 1024 \approx 1000 = 10^3$$

n bits $\rightarrow 2^n$ codes (words)

$$\frac{n}{4} \text{ decimal digits} \rightarrow 2^{0.83n}$$

$$\left. \begin{array}{l} 10 - x \\ x \end{array} \right\} \Rightarrow x = \frac{10n}{12} \approx 0,83n$$

* n bits $\rightarrow 2^n$ unsigned binary numbers

n bits $\rightarrow 10^{\frac{n}{5}}$ unsigned 2-out-of-5 decimal numbers

$$10^3 = 1000 \approx 1024 = 2^{10} \rightarrow \left. \begin{matrix} 3 \dots 10 \\ \frac{n}{5} \dots x \end{matrix} \right\} \Rightarrow x = \frac{\frac{n}{5} \times 10}{3} = \frac{2}{3}n \approx 0.66n$$

$\Rightarrow n$ bits $\rightarrow x \cdot 2^{0.66n}$ unsigned 2-out-of-5 decimal numbers

Binary floating point numbers

Number
notations

weighted
notation

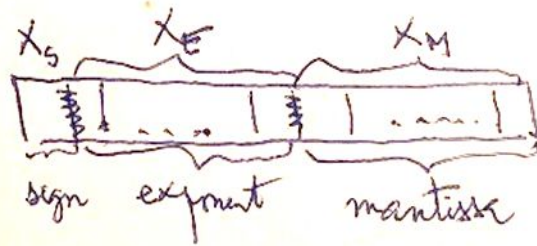
scientific
notation

$$N = \sum_{i=0}^{n-1} x_i \cdot r^i \text{ where } 0 \leq x_i < r$$

$$N = M \cdot B^E$$

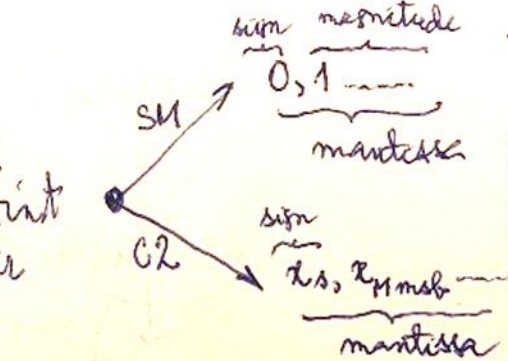
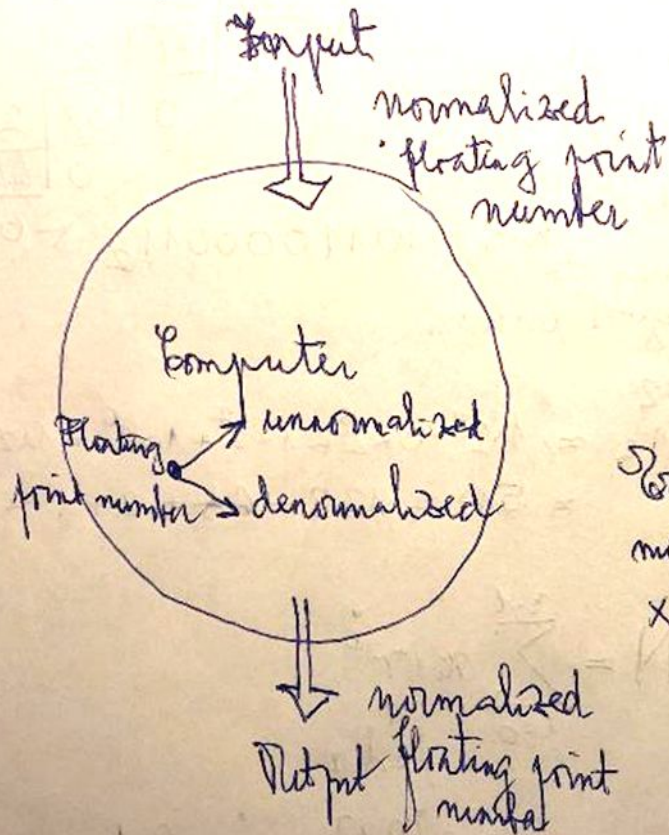
\uparrow mantissa

$E \leftarrow$ exponent
 $B \leftarrow$ base

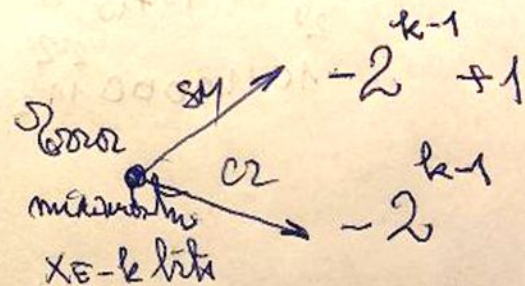


$$X = (-1)^{X_s} \cdot 2^{X_E} \cdot X_M$$

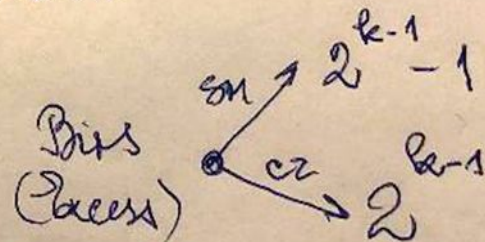
base mantissa



Mantissa problem



Exponent problem



$$X = 100.1110, 101 = 9,10.1110.101 \times 2^7$$

sign exp

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

$$Y = 0,001011 = 0,1011 \times 2^{-2}$$

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

Normalization

- Point leftshift \rightarrow Increasing exponent
- Point rightshift \rightarrow Decreasing exponent

$$X+Y = (X_M + Y_M^{Y_E - X_E}) \times 2^{X_E} \text{ where } X_E \geq Y_E$$

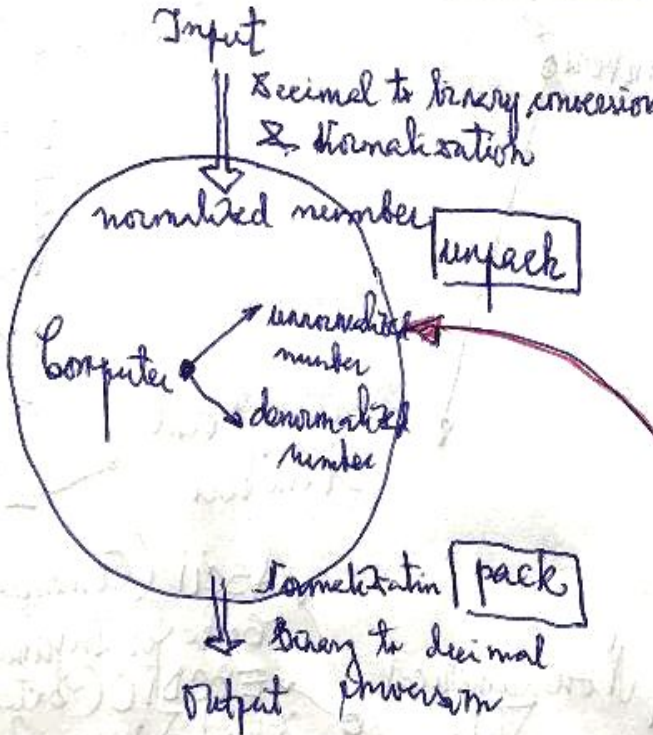
Floating Point Addition

1. Exponent comparison
ex 7 :: (-2) \Rightarrow
 $X_E - Y_E = 7 - (-2) = 9$
2. Rightshift for alignment with

$$\begin{array}{r} 0,1001110.101 + (X_E - Y_E) \text{ position} \\ 0,0000000.001011 \\ \hline 0,1001110.110011 \end{array}$$

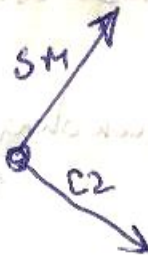
$$4. \text{ Normalization } 0,1001110.110011 \times 2^7$$

Mantissa Problem



Normalization Rules

Mantissa



* Observation:

- * Observation:
 - through normalization operation

$$\frac{1}{2} \leq |x_h| < 1$$

absolute value

$\frac{1}{2} \times 10 = 5$

for normalisation, point leftshift & increasing exponent or point rightshift & decreasing exponent

1/2	3/4	1
-----	-----	---

for normalization, point left shift
increasing exponent or point right shift
decreasing exponent, but
WARNING! for negative mantissa
right shift introducing 1s!!

Exponent Problem

$$X_{+0,2,10} = +0,0010001, \dots_2$$

$$\begin{array}{l} 0,2 \times 2 \\ 0,4 \times 2 \\ 0,8 \times 2 \\ 1,6 \times 2 \\ 0,2 \times 2 \\ 0,4 \times 2 \\ 0,8 \times 2 \\ 1,6 \times 2 \\ 0,2 \end{array}$$

$$X' = 2^{-3} + 2^{-6} = \frac{9}{64} = 0,140625_{10}$$

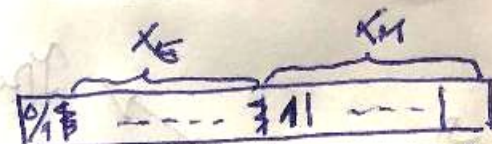
$$\Sigma_{\text{error}} = 2 - 0,1406 \pi = 0,059375_{10}$$

* from conversion algorithm \rightarrow truncation error

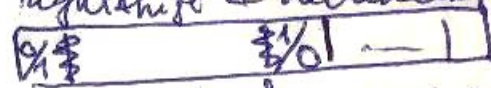
X'

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

Normalization
Rules



for normalization, point leftshift & increasing exponent or point rightshift & decreasing exponent



for normalization, point leftshift & increasing exponent or point rightshift & decreasing exponent, but

WARNING! for negative mantissa rightshift introducing 1s!!

* Observation:

- through normalization operation

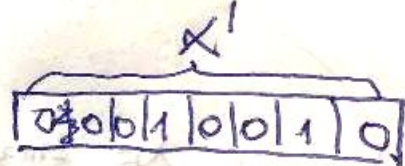
$$\frac{1}{2} \leq |X_M| < 1$$

absolute value

Exponent Problem

$$X = +0,2_{10} = +0,0010001, \dots_2$$

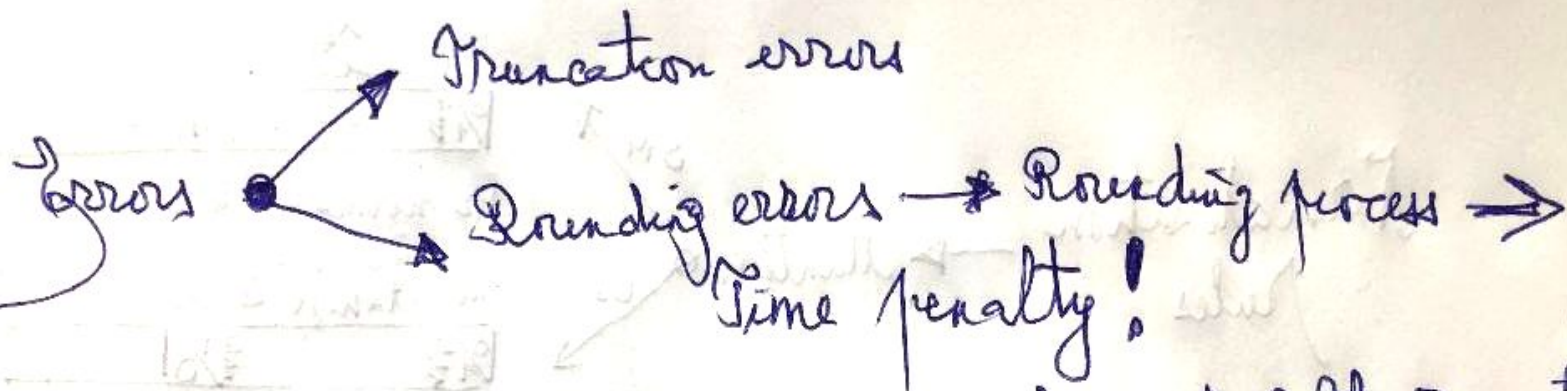
- 0,2 x 2
- 0,4 x 2
- 0,8 x 2
- 1,6 x 2
- 0,2 x 2
- 0,4 x 2
- 0,8 x 2
- 1,6 x 2
- 0,2



$$X' = 2^{-3} + 2^{-6} = \frac{9}{64} = 0,140625_{10}$$

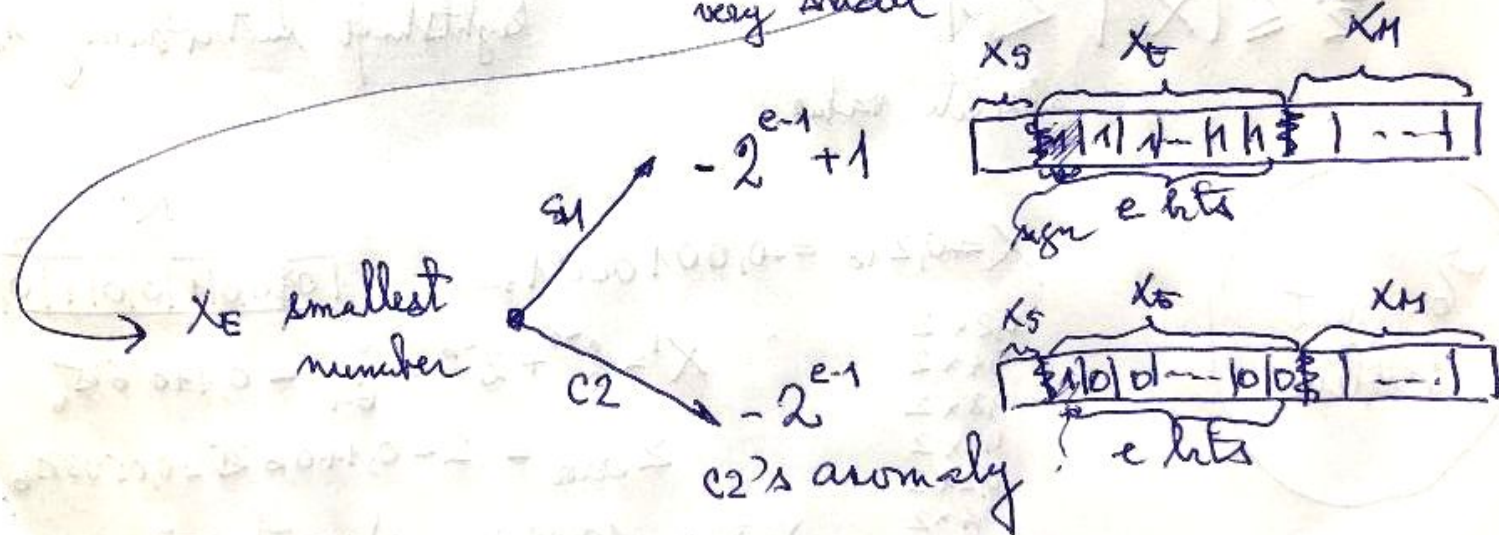
$$\Sigma \text{ error} = 2 - 0,140625 = 0,059375_{10}$$

* from conversion algorithm \rightarrow truncation error



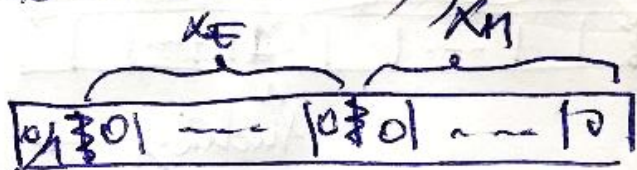
* By floating point operation → Partial result 0 mantissa large → significant error

Because errors → $Z = (-1)^{x_{n-1}} \cdot X_n \cdot 2^{\text{very small}}$

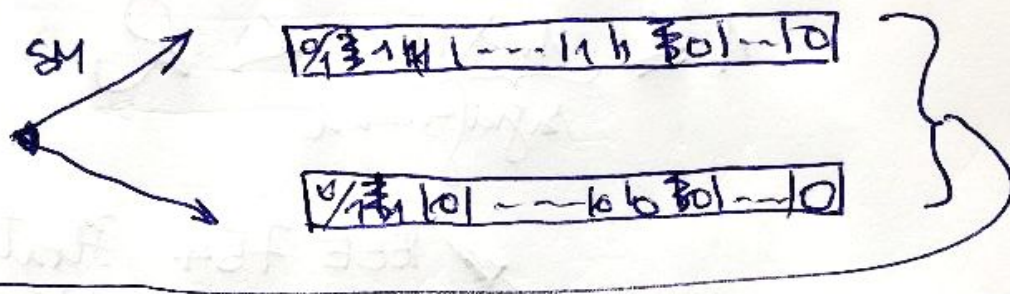


* But, for comparisons reasons (implementing of jumps or branch instructions), we want

0 representation as



and we have



Excess (Biased) Representation

SN \rightarrow $\text{Excess}(\text{Bias}) = 2^{e-1} = 1$

LS \rightarrow $\text{Excess}(\text{Bias}) = 2^{e-1}$

\rightarrow Exponent from signed binary integer number

through biased representation \rightarrow unsigned binary number