

AC1 - Guia Prático 9 (Vírgula Flutuante)

$N = (+/-) 1.f \times 2^{Exp}$

Diagram illustrating the IEEE 754 floating-point format structure:

- N : mantissa
- $(+/-)$: sign
- $1.f$: significando
- 2^{Exp} : expoente

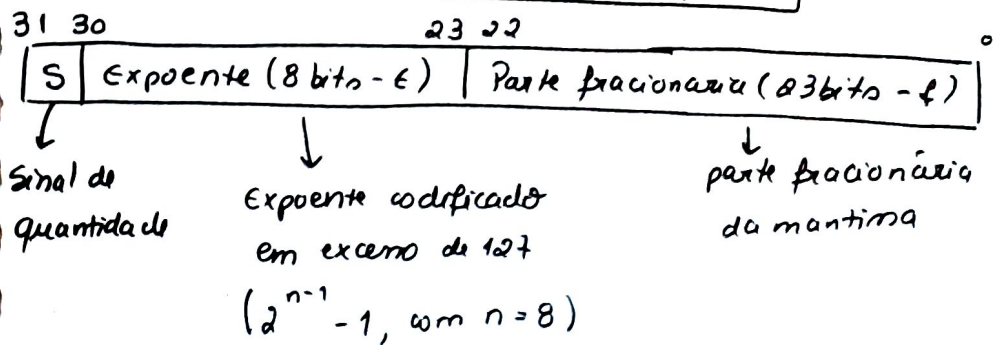
- O problema da divisão do espaço de armazenamento coloca-se também neste caso, mas agora na determinação do número de bits ocupados pela parte fracionária e pelo expoente.

- Esta divindade é um compromisso entre gama de representação e precisão:

\Rightarrow Aumento do n° de bits da parte fracionária
↓
maior precisão na representação

⇒ Aumento do n° de bits do expoente
↓
maior gama de representação.

Norma IEEE 754 (precisão simples)



→ A representação é sempre normalizada: o bit à esquerda do ponto binário é sempre 1.
Como é sempre 1, esse bit não é explicitamente representado (hidden bit)

- Example:

⑦ Qual o valor, em decimal, representado em
ox 4158 0000?

$0 \quad 10000010 \quad 101100000000000000000000$
 \downarrow
 Signo positivo
 $2^7 + 2^1 = 128 + 2 = 130$
 $\text{Exponente} = 130 - \text{offset}$
 $= 130 - 127$
 $= 3$
 \downarrow
 \in

$\text{Mantissa} = (1 + f) =$
 $= 1 + .1011 = 1.1011$
 A quantidade representada será
 $+1.1011 \times 2^3$

$$R = +1.1011 \times 2^3 = (1 + 1 \times 2^{-1} + 0 \times 2^{-2} + 1 \times 2^{-3} + 1 \times 2^{-4}) \times 2^3$$
$$= +1.6875 \times 8 = \underline{\underline{+13.5}}$$

② codificar no formato vírgula flutuante IEEE 754
precisão simples, o valor $-12593.75_{10} \times 10^{-3} =$

$$= -12.59375$$

Parte inteira = $12_{10} = 1100_2$

Parte fracionária = $0.59375_{10} = 0.10011_2$

$$12.59375_{10} = 1100.10011_2 \times 2^0$$

Normalização: $1100.10011_2 \times 2^0 =$
 $= 1.10010011_2 \times 2^3$

Exponente codificado: $+3 + 127 =$
 $= 130_{10} = 10000010_2$

1 10000010 100100110000000000000000
C 1 4 9 8 0 0 0 0

0xC1498000

$$\begin{array}{r} 0.59375 \\ \times \quad 2 \\ \hline \end{array}$$

$$\begin{array}{r} ①. 18750 \\ \times \quad 2 \\ \hline \end{array}$$

$$\begin{array}{r} 0.18750 \\ \times \quad 2 \\ \hline \end{array}$$

$$\begin{array}{r} ②. 37500 \\ \times \quad 2 \\ \hline \end{array}$$

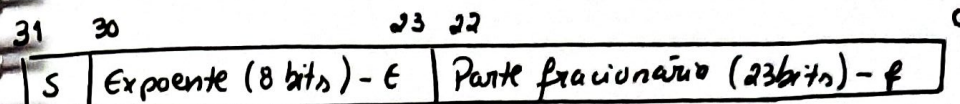
$$\begin{array}{r} ③. 75000 \\ \times \quad 2 \\ \hline \end{array}$$

$$\begin{array}{r} ④. 50000 \\ \times \quad 2 \\ \hline \end{array}$$

$$\begin{array}{r} ⑤. 50000 \\ \times \quad 2 \\ \hline \end{array}$$

$$\begin{array}{r} ⑥. 00000 \\ \times \quad 2 \\ \hline \end{array}$$

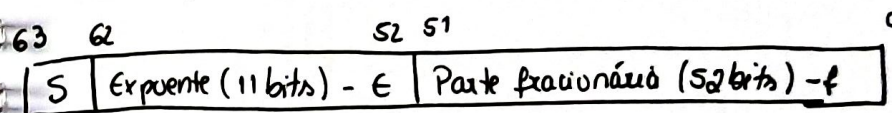
• Representação de quantidades em precisão simples (32 bits)



$$N = (-1)^S \cdot 1.f \times 2^{(E-127)}$$

Precisão simples \rightarrow tipo float

• Representação de quantidades em precisão dupla (64 bits)



$$N = (-1)^S \cdot 1.f \times 2^{(E-1023)}$$

Precisão dupla - tipo double

\hookrightarrow A quantidade zero é representável no formato descrito.

• Técnicas de arredondamento do resultado

① Truncatura (exemplo com $d = 2$ bits)
 \hookrightarrow parte frac.

val	Trunc(val)	Erro
x.00	x	0
x.01	x	$-\frac{1}{4}$
x.10	x	$-\frac{1}{2}$
x.11	x	$-\frac{3}{4}$

$$\begin{aligned} \text{Erro médio} &= \frac{0 - \frac{1}{4} - \frac{1}{2} - \frac{3}{4}}{4} = \\ &= \frac{-\frac{1}{4} - \frac{2}{4} - \frac{3}{4}}{4} = \\ &= \frac{-\frac{6}{4}}{4} = -\frac{6}{16} = -\frac{3}{8} \end{aligned}$$

② Arredondamento simples ($d = 2$ bits)

val	Arred(val)	Erro
x.00	x	0
x.01	x	$x - x.25 = -1/4$
x.10	x + 1	$(x+1) - x.5 = +1/2$
x.11	x + 1	$(x+1) - x.75 = +1/4$

$$\text{Erro médio} = \frac{0 - 1/4 + 1/2 + 1/4}{4} = 1/8 //$$

→ Soma-se 1 ao 1º bit à direita do ponto binário e truncamos o resultado

$$\Leftrightarrow \boxed{\text{arred}(\text{val}) = \text{trun}(\text{val} + 0.5)}$$

$$\begin{array}{r} x.00 \\ +0.1 \\ \hline x.10 \end{array} \quad \begin{array}{r} x.01 \\ +0.1 \\ \hline x.11 \end{array} \quad \begin{array}{r} x.10 \\ +0.1 \\ \hline x+1.00 \end{array} \quad \begin{array}{r} x.11 \\ +0.1 \\ \hline x+1.01 \end{array}$$

→ O erro médio é mais próximo de zero do que no caso da truncatura, mas ligeiramente polarizado do lado positivo

③ Arredondamento para o par mais próximo

val	Arred(val)	Erro	val	Arred(val)	Erro
x0.00	x0	0	x1.00	x1	0
x0.01	x0	$-1/4$	x1.01	x1	$-1/4$
x0.10	x0	$-1/2$	x1.10	x1+1	$+1/2$
x0.11	x1	$+1/4$	x1.11	x1+1	$+1/4$

$$\begin{array}{r} x0.10 \\ +0.0 \\ \hline x0.10 \end{array} \quad \begin{array}{r} x1.10 \\ +0.1 \\ \hline x1+1.00 \end{array}$$

$$\begin{aligned} \text{Erro médio} &= \frac{(0 - 1/4 - 1/2 + 1/4) + (0 - 1/4 + 1/2 + 1/4)}{4} \\ &= -1/8 + 1/8 = 0 // \end{aligned}$$

Norma IEEE 754 - arredondamentos



G - Guard Bit

R - Round Bit

↓ (S) - Sticky Bit

resultado da soma lógica de todos os bits à direita do bit R
(p.e. se houver à direita de R pelo menos um bit a '1',
então $S = '1'$).