


Aula AOC

24/08

- **Aritmética de ponto flutuante:**
  - ↳ Notação de ponto fixo (= ponto da notação decimal)
  - ↳ Notação de ponto flutuante
- O problema do ponto fixo é que o tamanho da parte inteira e da fracionária fica fixo com relação a seu armazenamento em memória.
- O ponto flutuante se move dinamicamente para uma posição conveniente usando a notação científica.



- ↳ A mantissa é armazenada em complemento-dois
- ↳ O expoente é representado na notação **excess** ou **biased (polarizado)**

Exemplo:

- ✗ Excess (bias) 128 significa
- ✗ Campo expoente de 8 bit
- ✗ Faixa de valores puros 0-255
- ✗ Subtraia 128 para obter o valor correto
- ✗ Nova faixa -128 to +127

⋮

## Números Fracionários

$$0.81 \rightarrow 9 - 1001.100\dots$$

$$\begin{array}{l}
 0.81 \xrightarrow{\times 2} 1,62 \rightarrow 1 \quad \left| \begin{array}{l} \text{se } e \geq 1 \\ \text{else } = 0 \end{array} \right. \\
 \overbrace{^1}^{0.62} \rightarrow 1,24 \rightarrow 1 \\
 0.24 \rightarrow 0,48 \rightarrow 0 \\
 0,48 \rightarrow 0,96 \rightarrow 0
 \end{array}$$

⋮

|                         | Sign   | Biased exponent | Fraction                 | Value               |
|-------------------------|--------|-----------------|--------------------------|---------------------|
| positive zero           | 0      | 0               | 0                        | 0                   |
| negative zero           | 1      | 0               | 0                        | -0                  |
| plus infinity           | 0      | all 1s          | 0                        | $\infty$            |
| minus infinity          | 1      | all 1s          | 0                        | $-\infty$           |
| quiet NaN               | 0 or 1 | all 1s          | $\neq 0$ ; first bit = 1 | qNaN                |
| signaling NaN           | 0 or 1 | all 1s          | $\neq 0$ ; first bit = 0 | sNaN                |
| positive normal nonzero | 0      | all 1s          | f                        | $2_{e-16383}(1.f)$  |
| negative normal nonzero | 1      | all 1s          | f                        | $-2_{e-16383}(1.f)$ |
| positive subnormal      | 0      | 0               | $f \neq 0$               | $2_{e-16383}(0.f)$  |
| negative subnormal      | 1      | 0               | $f \neq 0$               | $-2_{e-16383}(0.f)$ |

Continuação 25/08

## • Operações aritméticas c/ ponto flutuante

$$X = 0,3 \cdot 10^2 = 30$$

$$Y = 0,2 \cdot 10^3 = 200$$

$$X + Y = (0,3 \cdot 10^{2-3} + 0,2) \cdot 10^3 = 0,23 \cdot 10^3 = 230$$

$$X - Y = (0,3 \cdot 10^{2-3} - 0,2) \cdot 10^3 = -0,17 \cdot 10^3 = -170$$

$$X \cdot Y = (0,3 \cdot 0,2) \cdot 10^{2+3} = 0,06 \cdot 10^5 = 6000$$

$$X \div Y = (0,3 \div 0,2) \cdot 10^{2-3} = 1,5 \cdot 10^{-1} = 0,15$$

### • Fases importantes:

- Verificar se é zero
- Ajustar os expoentes
- Realizar a operação
- Normalizar o resultado

## • Bits de GUARDA

- Bits dentro das arquiteturas que têm como função principal melhorar as aproximações do processo aritmético de ponto flutuante.

## • Arredondamento

- Round to nearest
- Round towards  $+\infty / -\infty$
- Round towards 0

## • Infinito

• O infinito é tratado como o caso limite da aritmética real.

## • NaN

- Quiet - Erro sem aviso
- Signaling - Erro c/ aviso

## Casos QUIET:

Soma e subtrações c/  $\infty$

Multiplicar 0 ·  $\infty$

Dividir  $\frac{0}{0}$  ou  $\frac{\infty}{\infty}$

Raiz de números negativos