

## ① Representações de números no computador

- \* bit :  $\rightarrow$  0 ou 1 (com ou sem corrente)
  - $\rightarrow$  menor unidade de informação
  - $\rightarrow$  números binários

\* Representações dos números:

$$\begin{aligned}(1328)_{10} &= 1 \times 10^3 + 3 \times 10^2 + 2 \times 10^1 + 8 \times 10^0 \\ &= 1000 + 300 + 20 + 8\end{aligned}$$

$\Rightarrow$  em qualquer base:

$$(a_3 a_2 a_1 a_0)_\beta = (a_3 \beta^3 + a_2 \beta^2 + a_1 \beta^1 + a_0 \beta^0)_{10}$$

onde  $0 \leq a_i \leq \beta - 1$

$\Rightarrow$  base 2:  $(10010)_2 = (1 \times 2^4 + 0 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 0 \times 2^0)_{10}$

$$= (16 + 2)_{10} = (18)_{10}$$

•  $(0)_2 = (0 \times 2^0)_{10} = (0)_{10}$

•  $(1)_2 = (1 \times 2^0)_{10} = (1)_{10}$

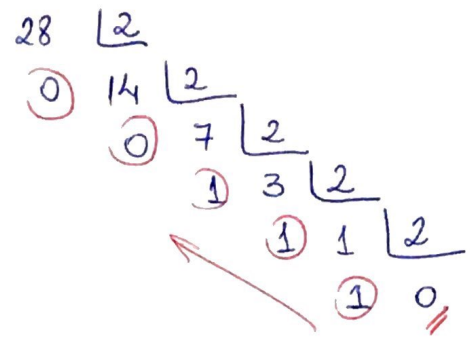
•  $(10)_2 = (1 \times 2^1 + 0 \times 2^0)_{10} = (2)_{10}$

•  $(11)_2 = (1 \times 2^1 + 1 \times 2^0)_{10} = (3)_{10}$

•  $(100)_2 = (1 \times 2^2 + 0 \times 2^1 + 0 \times 2^0)_{10} = (4)_{10}$

•  $(28)_{10} = ?$

$= (11100)_2$



•  $x = 2 \rightarrow (10)_2 \rightarrow 2$  bits

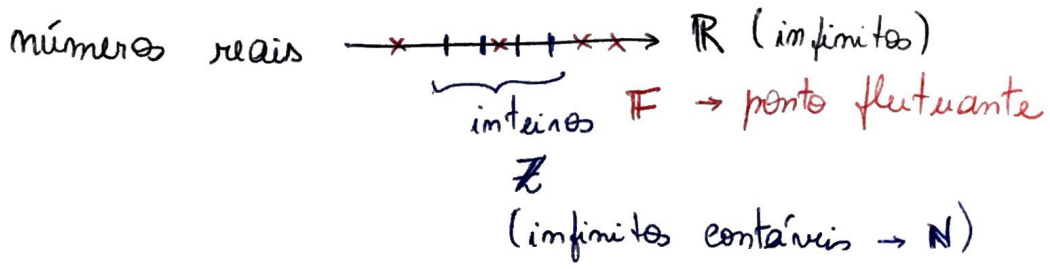
•  $x = 0,25$  ?

•  $x = 1/3 = 0,333$  ?

•  $x = \sqrt{2}$  ?

} não tem soluções exatas no computador

↳ irracional,  $\neq \frac{A}{B}$



$\Rightarrow \mathbb{F} \rightarrow$  ponto flutuante  $\rightarrow$  número finito de casas decimais  
 ↳ arredondamento

$x = 1/3 = 0,3333$

$x = 2/3 = 0,6667$

$\Rightarrow$  posição do número decimal não é fixa

$x = (-1)^s (0.\underbrace{a_1 a_2 \dots a_t}_m) \beta^e = (-1)^s m \beta^{e-t}$

•  $\beta =$  base  $\geq 2$  (inteiro)

•  $m =$  mantissa (inteiro)  $0 \leq a_i \leq \beta - 1$

•  $t =$  nº de dígitos da mantissa

•  $e =$  expoente (inteiro),  $s =$  sinal  $\begin{cases} 0 & \oplus \\ 1 & \ominus \end{cases}$

ex) . 250 = (-1)^0 0,25 x 10^3 = (-1)^0 25 x 10^1

. - 0,33 = (-1)^1 0,33 x 10^0 = (-1)^1 33 x 10^-2

. 0,001 = (-1)^0 0,1 x 10^-2 = (-1)^0 1 x 10^-3

obs: se a1 ≠ 0, representação é única.

\* Padrão IEEE 754

(-1)^s 1 . m x 2^{e-b} E = e - b

- byte = 8 bits
half = 16 bits (2 bytes)
single = 32 bits (4 bytes)
double = 64 bits (8 bytes)

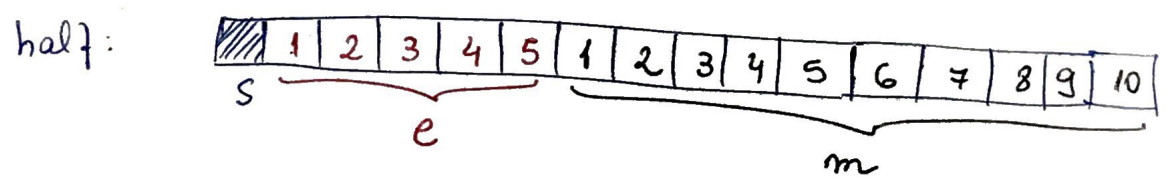


Table of special cases: 0 000 00 000 000 0000 -> (0)10, 0 111 11 '' -> +inf, 1 111 11 '' -> -inf

obs: e = { 0 0 0 0 0, 1 1 1 1 1 } são casos especiais

(1)10 = (1)2 -> 0 01111 000 000 0000 = 1 x 2^0 = (-1)^0 1.0 x 2^0 s=0, m=0, E=0

⇒ half: 16 bits: 1(s), 5(e), 10(m)

$F(2, 11, -14, 15)$

•  $e_{max} = (11110)_2 = (30)_{10}$

•  $e_{min} = (00001)_2 = (1)_{10}$

•  $e = [1, 30] \Rightarrow b = 15 \Rightarrow E = [-14, 15]$

números pequenos  
↓

↑  
números grandes

•  $|X|_{max} = 0\ 11110\ 1111111111$

$= (-1)^0 (1, 1111111111)_2 \times 2^{(11110)_2 - (15)_{10}}$

$= [1 \times 2^0 + (1 - 2^{-10})] \times 2^{15}$

$= (2 - 2^{-10}) 2^{15} = 65504 = 6,55 \times 10^4$

•  $|X|_{min} = 0\ 00001\ 000\ 0000000$

$= (-1)^0 1,0 \times 2^{(00001)_2 - (15)_{10}}$

$= 1 \times 2^{-14} = 1/16384 = 6,10 \times 10^{-5}$

• precisão: 10 (mantissa) + 1 = 11 bits

$\varphi_{10} = 3$

⇒ single: 32 bits: 1(s), 8(e), 23(m)

(5)

$\mathbb{F}(2, 24, -126, 127)$

•  $e_{\max} = (11111110)_2 = (254)_{10}$

•  $e_{\min} = (00000001)_2 = (1)_{10}$

•  $e = [1, 254] \Rightarrow b = 127 \Rightarrow E = [-126, 127]$

•  $|x|_{\max} = (2 - 2^{-23}) \times 2^{127} = 3,4028235 \times 10^{38}$

•  $|x|_{\min} = 1 \times 2^{-126} = 1,1754944 \times 10^{-38}$

•  $p_{10} = 7$  (24 bits em  $\beta=10$ )

⇒ double: 64 bits: 1(s), 11(e), 52(m)

$\mathbb{F}(2, 53, -1022, 1023)$

•  $e_{\max} = (1111111110)_2 = (2046)_{10}$

•  $e_{\min} = (0000000001)_2 = (1)_{10}$

•  $e = [1, 2046] \rightarrow b = 1023 \rightarrow E = [-1022, 1023]$

•  $|x|_{\max} = (2 - 2^{-52}) 2^{1023} = 1,797693 \dots \times 10^{308}$

•  $|x|_{\min} = 1 \times 2^{-1022} = 2,2250 \dots \times 10^{-308}$

•  $p_{10} = 15$  (53 bits em  $\beta=10$ )

$\mathbb{F}(\beta, t, L, U)$

$L \leq E \leq U$

obs: Quaternion:  $\mathbb{F}(2, 53, -1021, 1024) \Rightarrow (-1)^s 0.1_m 2^{e-b}$



① Representação de números no computador

⇒ ponto flutuante:  $F(\beta, t, \underbrace{L, U}_{\text{decimal}})$ ,  $L \leq E \leq U$

\* Padrão IEEE 754:  $x = (-1)^s \underbrace{1}_{\text{implícito}} m \cdot 2^{e-b}$  ← bias

$$\left\{ \begin{array}{l} \beta = \text{base} = 2 \\ t = \text{precisão binária} = m + 1 \quad (m = \text{mantissa}) \\ E = \text{expoente} = e - b \end{array} \right.$$

• half:  $F(2, 11, -14, 15)$   $\left. \begin{array}{l} 16 \text{ bits: } 1(s), 5(e), 10(m) \\ b = 15 \end{array} \right\}$

• single:  $F(2, 24, -126, 127)$   $\left. \begin{array}{l} 32 \text{ bits: } 1(s), 8(e), 23(m) \\ b = 127 \end{array} \right\}$

• double:  $F(2, 53, -1022, 1023)$  64 bits: 1(s), 11(e), 52(m)

ex:  $\left\{ \begin{array}{l} x = (1)_{10} \\ x = (1)_2 \end{array} \right. \quad x = (-1)^0 1,0 \cdot 2^0 \quad s = 0, E = 0, m = 0$   
 $e = E + b = 0 + 15 = (15)_{10}$   
 $e = 1111$

$x_{\text{half}} = 0 \ 01111 \ 0000000000$

obs: ~~single~~ <sup>double</sup> Quaterni:  $F(2, 53, -1021, 1024)$

pois  $x = (-1)^s 0,1 m \cdot 2^{e-b}$

$|x|_{\min} = 2^L$ ,  $|x|_{\max} = (2 - 2^{-\overset{t-1}{m}}) 2^U$

\* precisão em decimal:

• half: 11 bits em  $\beta=2 \Rightarrow (11111111111)_2 = (2047)_{10}$   
 $p_{10} = 3$  (até 999).

• single: 24 bits em  $\beta=2 \Rightarrow (111\dots 11)_2 = 16.777.215$   
 $p_{10} = 7$  (até 9.999.999)

• double: 53 bits em  $\beta=2 \Rightarrow (111\dots 11)_2 = ?$   
 $p_{10} = 15$

- 
- \*  $|x| > |x|_{max} \Rightarrow \text{overflow } (\pm\infty)$
  - \*  $|x| < |x|_{min} \Rightarrow \text{underflow (zero)}$
- 

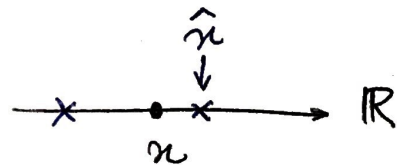
\* exemplo memória: 1 milhão de dados

- half:  $10^6 \times 2 \text{ bytes} \approx 2 \text{ Mb}$
- single:  $10^6 \times 4 \text{ bytes} \approx 4 \text{ Mb}$
- double:  $10^6 \times 8 \text{ bytes} \approx 8 \text{ Mb}$

↳ usado pela RAM ou armazenamento binário

↳ ASCII → caracteres que representam números →  
~ 3x arquivo binário.

1.1 Erro em ponto flutuante:



•  $EA_n = |x - \hat{n}| \rightarrow$  erro absoluto, minimizado na escolha de  $\hat{n}$  (arredondamento)

•  $ER_n = \frac{|x - \hat{n}|}{|x|} \rightarrow$  erro relativo, ~ constante em ponto flutuante

ex)  $\beta = 10, p_{10} = 4$

⇒ ponto fixo:

$$\begin{array}{l}
 \text{a) } x = 3507,6 \\
 \hat{x} = 3508
 \end{array}
 \left. \vphantom{\begin{array}{l} \text{a) } x = 3507,6 \\ \hat{x} = 3508 \end{array}} \right\}
 \begin{array}{l}
 EA_x = 0,4 \\
 ER_x = \frac{0,4}{3507,6} \approx 1,1 \times 10^{-4}
 \end{array}$$

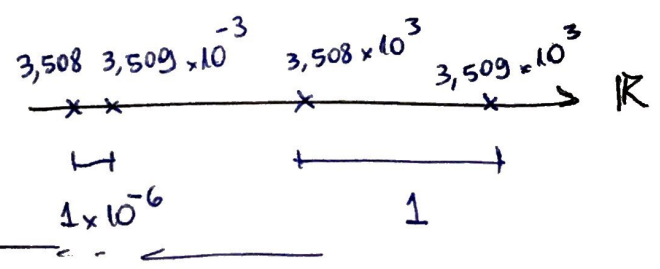
$$\begin{array}{l}
 \text{b) } x = 0,0035076 \\
 \hat{x} = 0,004
 \end{array}
 \left. \vphantom{\begin{array}{l} \text{b) } x = 0,0035076 \\ \hat{x} = 0,004 \end{array}} \right\}
 \begin{array}{l}
 EA_x = 0,0004924 \\
 ER_x \approx 0,14
 \end{array}$$

⇒ ponto flutuante:

$$\begin{array}{l}
 \text{a) } x = 3507,6 = 3,5076 \times 10^3 \\
 \hat{x} = 3,508 \times 10^3
 \end{array}
 \left. \vphantom{\begin{array}{l} \text{a) } x = 3507,6 = 3,5076 \times 10^3 \\ \hat{x} = 3,508 \times 10^3 \end{array}} \right\}
 \begin{array}{l}
 EA_x = 0,4 \\
 ER_x \approx 1,1 \times 10^{-4}
 \end{array}$$

$$\begin{array}{l}
 \text{b) } x = 0,0035076 = 3,5076 \times 10^{-3} \\
 \hat{x} = 3,508 \times 10^{-3}
 \end{array}
 \left. \vphantom{\begin{array}{l} \text{b) } x = 0,0035076 = 3,5076 \times 10^{-3} \\ \hat{x} = 3,508 \times 10^{-3} \end{array}} \right\}
 \begin{array}{l}
 EA_x = 0,0004 \times 10^{-3} \\
 ER_x = \frac{0,0004 \times 10^{-3}}{3,5076 \times 10^{-3}} \approx \frac{1,1 \times 10^{-4}}{1}
 \end{array}$$

obs: quanto maior o valor absoluto, maior o erro absoluto e o espaçamento na reta F



\* operações aritméticas com ponto flutuante

a) associativa:  $(2 + 3) + 4 = 2 + (3 + 4)$   
 $(2 \times 3) \times 4 = 2 \times (3 \times 4)$

x F



b) comutativa:  $x + y = y + x$

$xy = yx$

✓ F

c) distributiva:  $2(1+3) = 2 \cdot 1 + 2 \cdot 3$  X F

erros de arredondamento a cada operação.

$$\text{ex) } \left\{ \begin{array}{l} (23,4 + 5,18) + 3,05 = 31,7 \\ 23,4 + (5,18 + 3,05) = 31,6 \end{array} \right.$$

$$\text{ex) } \left\{ \begin{array}{l} 3,18 (5,05 + 11,4) = 52,5 \\ 3,18 \cdot 5,05 + 3,18 \cdot 11,4 = 52,4 \end{array} \right.$$

$$\text{ex) } \frac{(1+x)^{-1} - 1}{x} = 1$$

se  $x = 1 \times 10^{-15}$  (próximo da precisão de máquina)

$$\frac{(1+x)^{-1} - 1}{x} = 1,1102... \quad \text{erro de 11\% !}$$

obs: - trabalhe com números na ordem de 1. !

- Cap 1 Quaterni: