

Parte 2

Representação de Números Reais

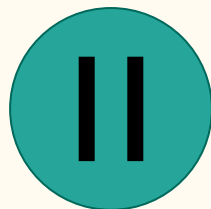
CI1164 - Introdução à Computação Científica
Profs. Armando Delgado e Guilherme Derenievicz
Departamento de Informática - UFPR

Sistemas de Numeração

Formas de representar quantidades ou números

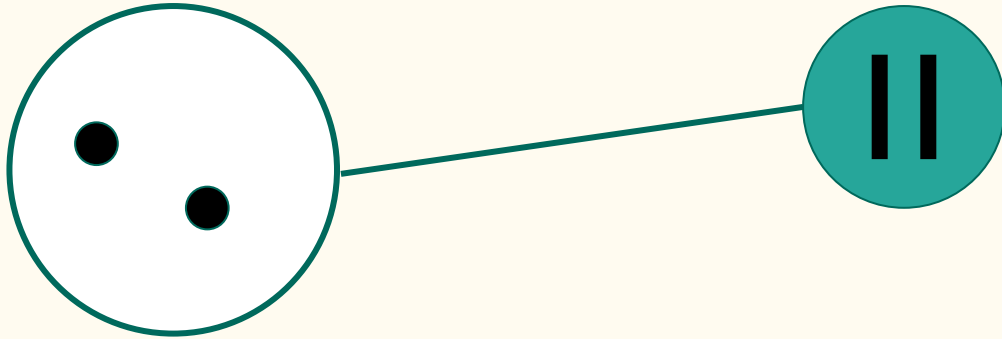
Sistemas de Numeração

Formas de representar quantidades ou números



Sistemas de Numeração

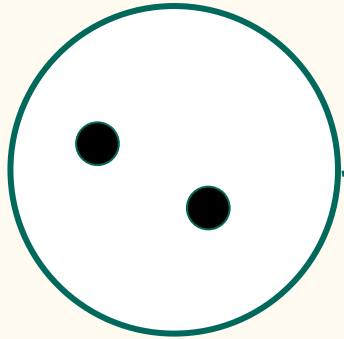
Formas de representar quantidades ou números



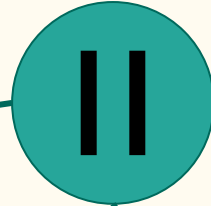
Sistema Numeral Romano

Sistemas de Numeração

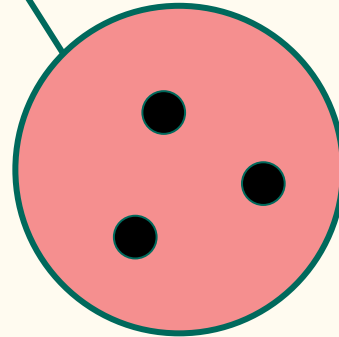
Formas de representar quantidades ou números



Sistema Numeral Romano

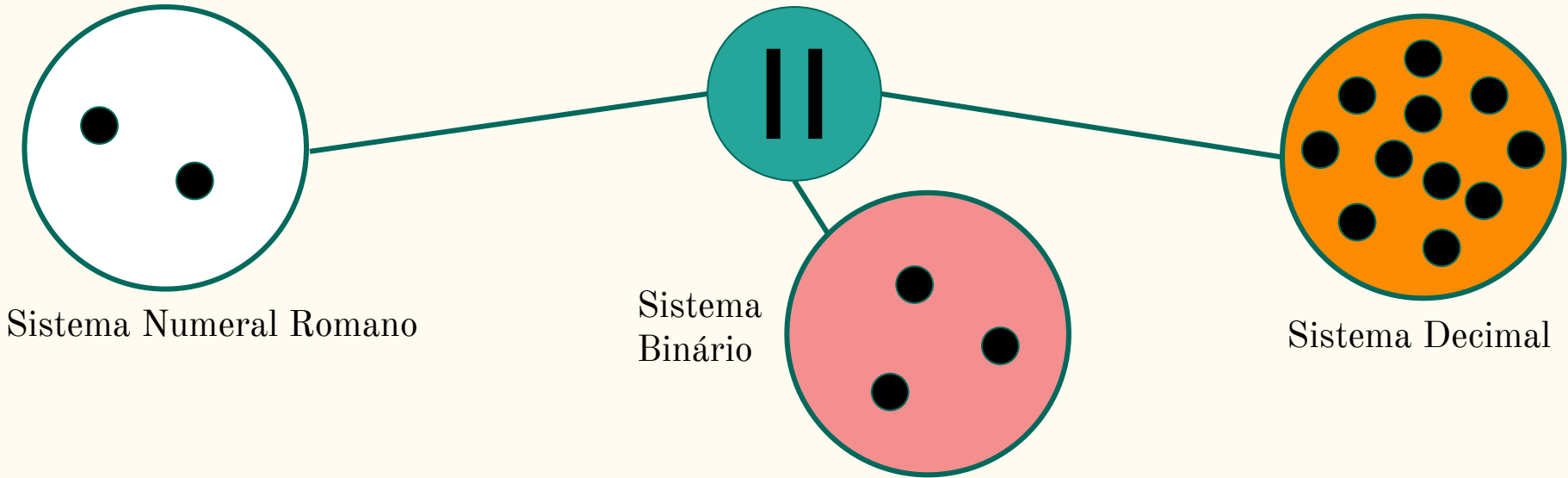


Sistema Binário



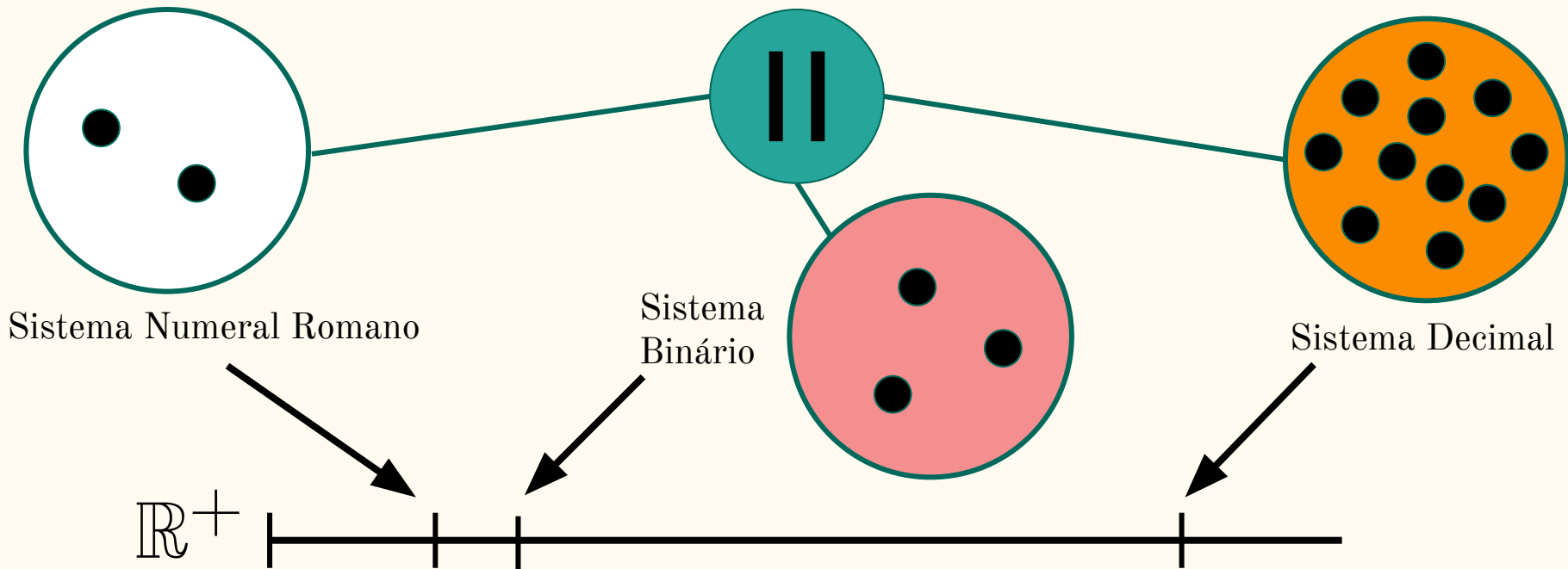
Sistemas de Numeração

Formas de representar quantidades ou números



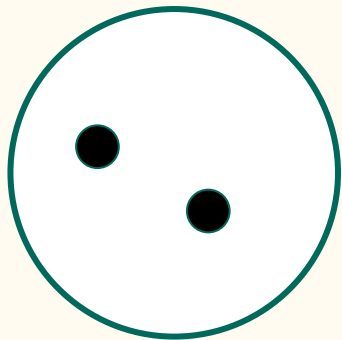
Sistemas de Numeração

Formas de representar quantidades ou números



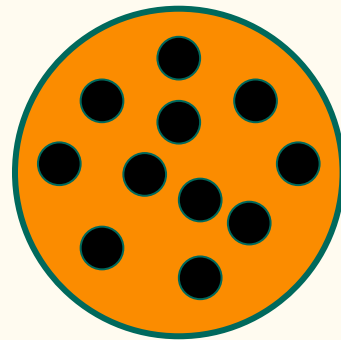
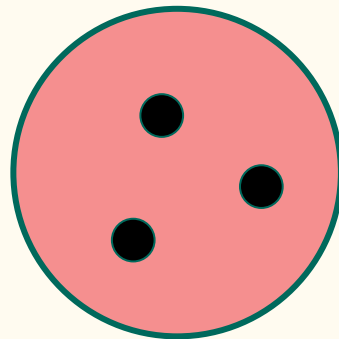
Sistemas de Numeração

Base: quantidade β de símbolos que podem ser utilizados para representar um número em dado sistema de numeração.



Sistema Numeral Romano

Sistema
Binário



Sistema Decimal



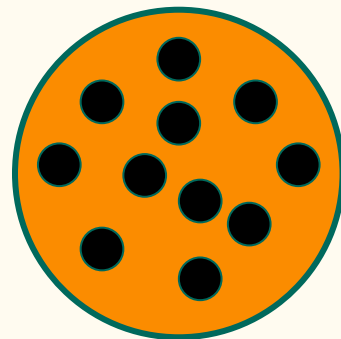
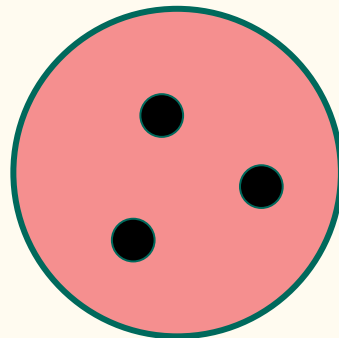
Sistemas de Numeração

Base: quantidade β de símbolos que podem ser utilizados para representar um número em dado sistema de numeração.

$\beta = 7$
{I, V, X, L, C, D, M}
Ex: VI

Sistema Numeral Romano

Sistema
Binário



Sistema Decimal



Sistemas de Numeração

Base: quantidade β de símbolos que podem ser utilizados para representar um número em dado sistema de numeração.

$$\beta = 7$$

{I, V, X, L, C, D, M}

Ex: VI

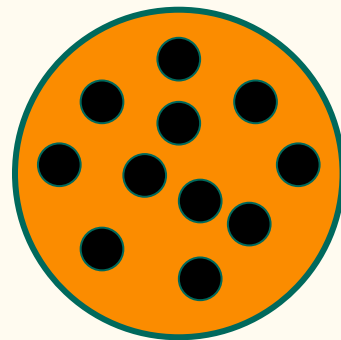
Sistema Numeral Romano

$$\beta = 2$$

{0, 1}

Ex: 110

Sistema
Binário



Sistema Decimal

\mathbb{R}^+



Sistemas de Numeração

Base: quantidade β de símbolos que podem ser utilizados para representar um número em dado sistema de numeração.

$$\beta = 7$$

{I, V, X, L, C, D, M}

Ex: VI

Sistema Numeral Romano

$$\beta = 10$$

{0, 1, 2, 3, 4, 5, 6, 7, 8, 9}

Ex: 6

Sistema Decimal

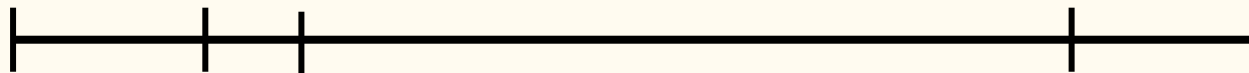
$$\beta = 2$$

{0, 1}

Ex: 110

Sistema
Binário

\mathbb{R}^+



Sistemas de Numeração

Base: quantidade β de símbolos que podem ser utilizados para representar um número em dado sistema de numeração.

$\beta = 7$
{I, V, X, L, C, D, M}
Ex: VI

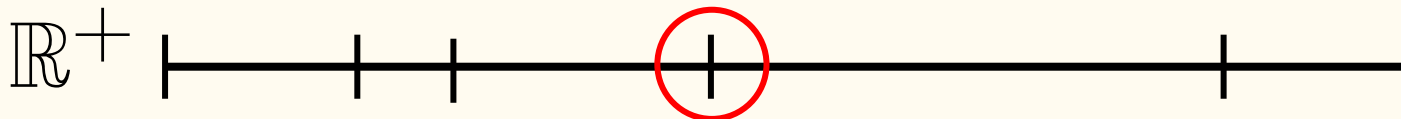
Sistema Numeral Romano

$\beta = 2$
{0, 1}
Ex: 110

Sistema Binário

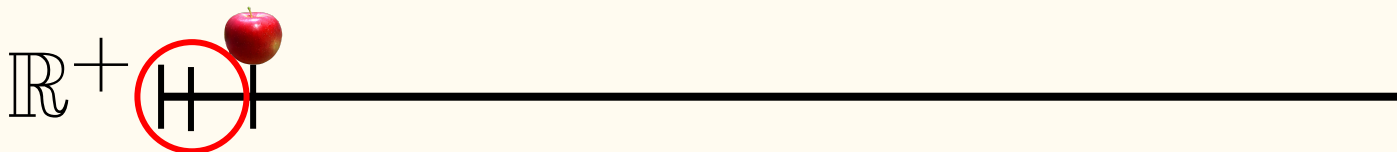
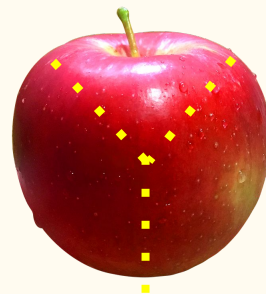
$\beta = 10$
{0, 1, 2, 3, 4, 5, 6, 7, 8, 9}
Ex: 6

Sistema Decimal



Sistemas de Numeração

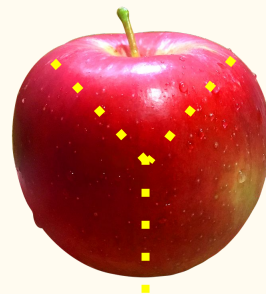
Base: quantidade β de símbolos que podem ser utilizados para representar um número em dado sistema de numeração.



Sistemas de Numeração

Base: quantidade β de símbolos que podem ser utilizados para representar um número em dado sistema de numeração.

Sistema Decimal: 0,33333...

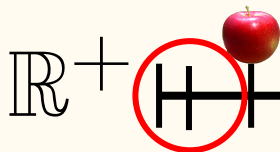
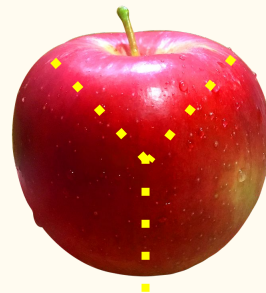


Sistemas de Numeração

Base: quantidade β de símbolos que podem ser utilizados para representar um número em dado sistema de numeração.

Sistema Decimal: $0,33333\dots$

$0,99999\dots = 1$

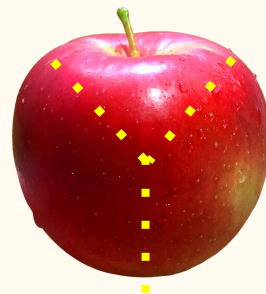


Sistemas de Numeração

Base: quantidade β de símbolos que podem ser utilizados para representar um número em dado sistema de numeração.

Sistema Decimal: 0,33333...

Sistema Numeral Romano: ::



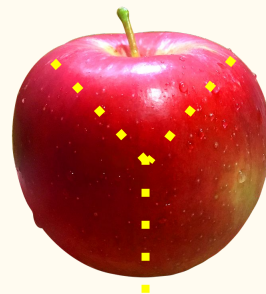
Sistemas de Numeração

Base: quantidade β de símbolos que podem ser utilizados para representar um número em dado sistema de numeração.

Sistema Decimal: 0,33333...

Sistema Numeral Romano: ::

Sistema Binário: 0,010101...



Sistemas de Numeração

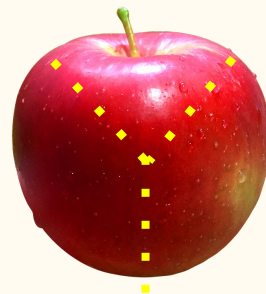
Base: quantidade β de símbolos que podem ser utilizados para representar um número em dado sistema de numeração.

Sistema Decimal: 0,33333...

Sistema Numeral Romano: ::

Sistema Binário: 0,010101...

Sistema Ternário: 0,1



Sistemas de Numeração

Base: quantidade β de símbolos que podem ser utilizados para representar um número em dado sistema de numeração.

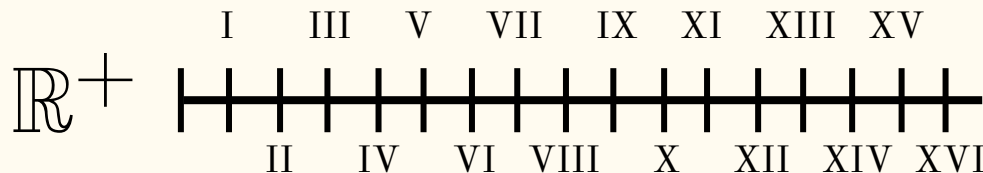
Sistema Decimal: 0,33333...

Sistema Numeral Romano: ::

Sistema Binário: 0,010101...

Sistema Ternário: 0,1

Sistema Aditivo



Sistemas de Numeração

Base: quantidade β de símbolos que podem ser utilizados para representar um número em dado sistema de numeração.

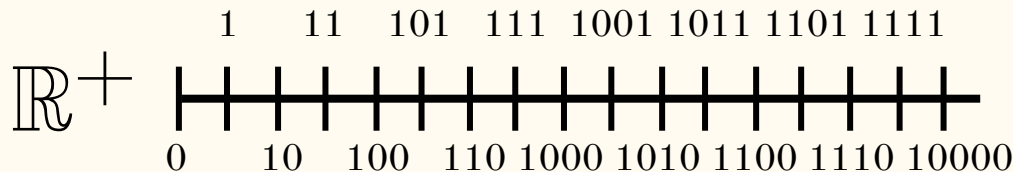
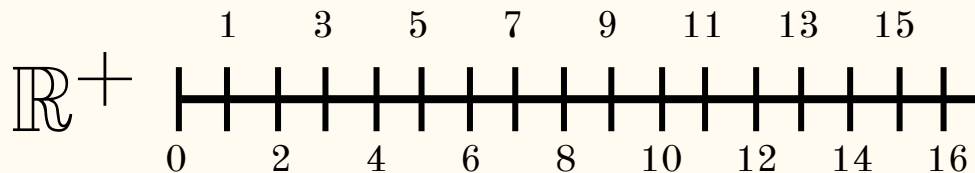
Sistema Decimal: 0,33333...

Sistema Numeral Romano: ::

Sistema Binário: 0,010101...

Sistema Ternário: 0,1

Sistema Posicional



Sistemas de Numeração

Base: quantidade β de símbolos que podem ser utilizados para representar um número em dado sistema de numeração.

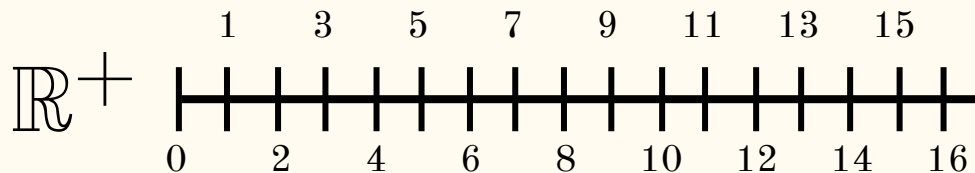
Sistema Decimal: 0,33333...

Sistema Numeral Romano: ::

Sistema Binário: 0,010101...

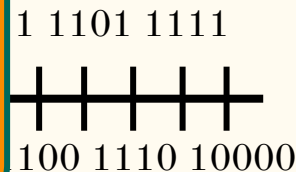
Sistema Ternário: 0,1

Sistema Posicional



$$n = a_2 a_1 a_0, b_1 b_2 b_3$$

$$n = a_2 * \beta^2 + a_1 * \beta^1 + a_0 * \beta^0 + b_1 * \beta^{-1} + b_2 * \beta^{-2} + b_3 * \beta^{-3}$$



Conversão de Base

Converter da base β para decimal

$$n = (a_2 a_1 a_0, b_1 b_2 b_3)_\beta$$

$$n = a_2 * \beta^2 + a_1 * \beta^1 + a_0 * \beta^0 + b_1 * \beta^{-1} + b_2 * \beta^{-2} + b_3 * \beta^{-3}$$

Conversão de Base

Converter da base β para decimal

$$n = (a_2 a_1 a_0, b_1 b_2 b_3)_\beta$$

$$n = a_2 * \beta^2 + a_1 * \beta^1 + a_0 * \beta^0 + b_1 * \beta^{-1} + b_2 * \beta^{-2} + b_3 * \beta^{-3}$$

Exemplo:

$$n = (101010,11)_2$$

$$\begin{aligned} n &= 1 * 2^5 + 0 * 2^4 + 1 * 2^3 + 0 * 2^2 + 1 * 2^1 + 0 * 2^0 + 1 * 2^{-1} + 1 * 2^{-2} \\ &= 32 + 8 + 2 + 0,5 + 0,25 \\ &= (42,75)_{10} \end{aligned}$$

Conversão de Base

Converter de decimal para a base β

$$n = (a_2 a_1 a_0, b_1 b_2 b_3)_{10}$$

Parte inteira:

$$\begin{array}{rcl} a_2 a_1 a_0 & \Big| & \beta \\ r_0 & & q_1 \\ & \Big| & \beta \\ & r_1 & q_2 \\ & & \Big| \beta \\ & & \dots \\ & & q_{n-1} \\ & & \Big| \beta \\ & r_{n-1} & 0 \end{array}$$

Conversão de Base

Converter de decimal para a base β

$$n = (a_2 a_1 a_0, b_1 b_2 b_3)_{10}$$

Parte inteira:

$$\begin{array}{r}
 a_2 a_1 a_0 \big| \beta \\
 \hline
 r_0 \quad q_1 \\
 q_1 \big| \beta \\
 \hline
 r_1 \quad q_2 \\
 q_2 \big| \beta \\
 \hline
 \dots \\
 q_{n-1} \big| \beta \\
 \hline
 r_{n-1} \quad 0
 \end{array}$$

Parte fracionária:

$$\begin{array}{l}
 0, b_1 b_2 b_3 \times \beta \\
 \hline
 d_1, b_1^1 b_2^1 b_3^1 \dots \rightarrow 0, b_1^1 b_2^1 b_3^1 \dots \times \beta \\
 \hline
 d_2, b_1^2 b_2^2 b_3^2 \dots \rightarrow 0, b_1^2 b_2^2 b_3^2 \dots \times \beta \\
 \hline
 d_3, b_1^3 b_2^3 b_3^3 \dots \\
 \hline
 \dots
 \end{array}$$

Conversão de Base

Converter de decimal para a base β

$$n = (a_2 a_1 a_0, b_1 b_2 b_3)_{10}$$

Parte inteira:

$$\begin{array}{r}
 a_2 a_1 a_0 \big| \beta \\
 \hline
 r_0 \quad q_1 \\
 q_1 \big| \beta \\
 \hline
 r_1 \quad q_2 \\
 q_2 \big| \beta \\
 \hline
 \dots \\
 q_{n-1} \big| \beta \\
 \hline
 r_{n-1} \quad 0
 \end{array}$$

Parte fracionária:

$$\begin{array}{l}
 0, b_1 b_2 b_3 \times \beta \\
 \hline
 d_1, b^1_1 b^1_2 b^1_3 \dots \rightarrow 0, b^1_1 b^1_2 b^1_3 \dots \times \beta \\
 \hline
 d_2, b^2_1 b^2_2 b^2_3 \dots \rightarrow 0, b^2_1 b^2_2 b^2_3 \dots \times \beta \\
 \hline
 d_3, b^3_1 b^3_2 b^3_3 \dots \\
 \hline
 \dots
 \end{array}$$

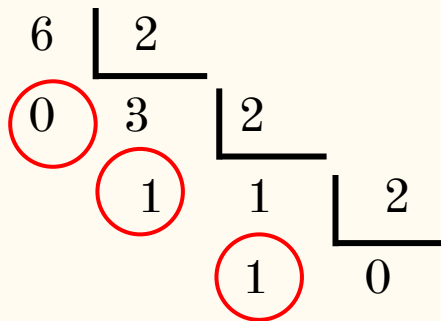
$$n = (r_{n-1} \dots r_1 r_0, d_1 d_2 d_3 \dots)_\beta$$

Conversão de Base

Converter de decimal para a base β

Exemplo: $n = (6,34)_{10}$

Parte inteira:



Conversão de Base

Converter de decimal para a base β

Exemplo: $n = (6,34)_{10}$

Parte inteira:

$$\begin{array}{r|l} 6 & 2 \\ \hline 0 & 3 \\ & \hline 1 & 2 \\ & \hline & 1 \\ & \hline & 1 \\ & \hline & 0 \end{array}$$

Parte fracionária:

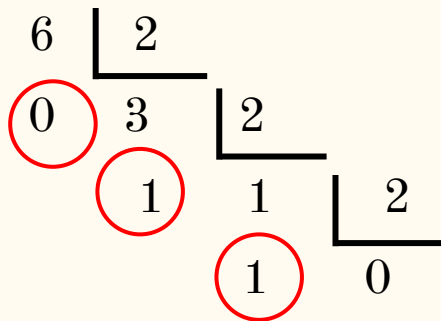
$$\begin{array}{l} 0,34 \times 2 \\ \hline 0,68 \rightarrow 0,68 \times 2 \\ \hline 1,36 \rightarrow 0,36 \times 2 \\ \hline 0,72 \rightarrow 0,72 \times 2 \\ \hline 1,44 \rightarrow 0,44 \times 2 \\ \hline 0,88 \rightarrow 0,88 \dots \end{array}$$

Conversão de Base

Converter de decimal para a base β

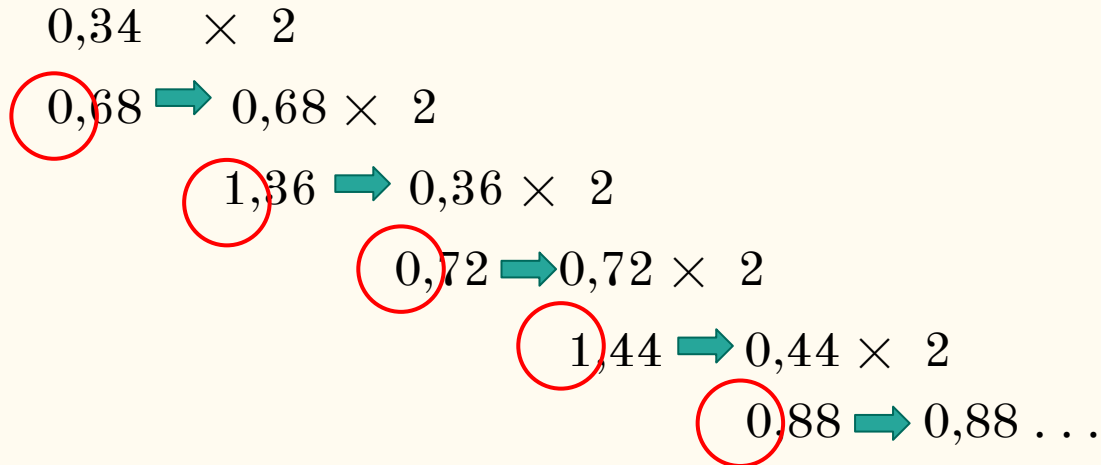
Exemplo: $n = (6,34)_{10}$

Parte inteira:



$$n = (110,01010\dots)_2$$

Parte fracionária:



Conversão de Base

Converter da base β_1 para a base β_2

1. Converter da base β_1 para decimal
2. Converter de decimal para a base β_2

Representação em Ponto Flutuante



$$\frac{1}{3} = 0,333333333333333333333333333333333333...$$

$$\sqrt{2} = 1,4142135623730950488016887242096...$$

$$\pi = 3,14159265358979323846264338327950...$$

Representação em Ponto Flutuante



Representação em Ponto Flutuante



000,00

000,01

000,02

000,03

...

999,99

Representação em Ponto Flutuante



000,00

000,01

000,02

000,03

...

999,99

\mathbb{R}^+



Representação em Ponto Flutuante



000,00

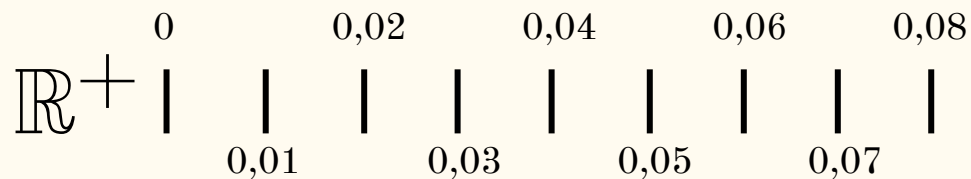
000,01

000,02

000,03

...

999,99



Representação em Ponto Flutuante



000,00

000,01

000,02

000,03

...

999,99

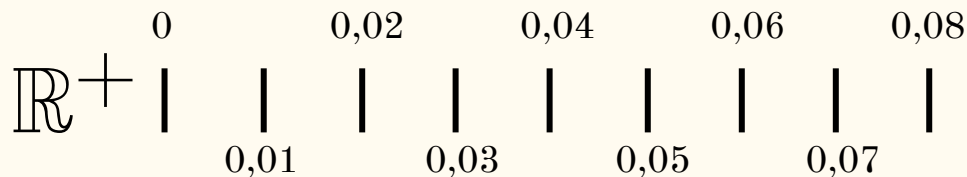
\mathbb{R}^+ ⁰ 999,99

Representação em Ponto Flutuante



000,00
000,01
000,02
000,03
...
999,99

$$0,08 \div 2 = 0,04$$

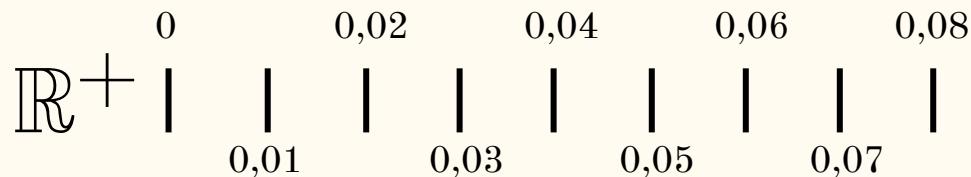
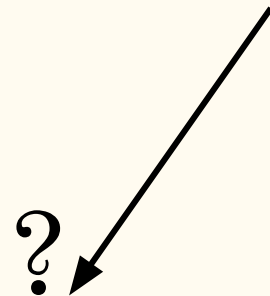


Representação em Ponto Flutuante



000,00
000,01
000,02
000,03
...
999,99

$$0,07 \div 2 = 0,035$$



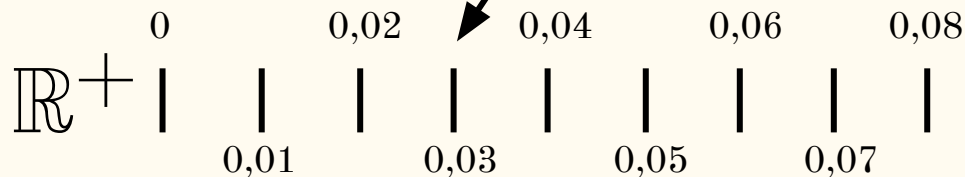
Representação em Ponto Flutuante



000,00
000,01
000,02
000,03
...
999,99

$$0,07 \div 2 = 0,03$$

truncamento



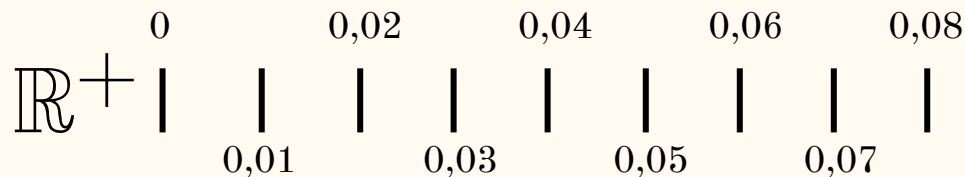
Representação em Ponto Flutuante



000,00
000,01
000,02
000,03
...
999,99

$$0,07 \div 2 = 0,04$$

arredondamento



Representação em Ponto Flutuante



000,00

000,01

000,02

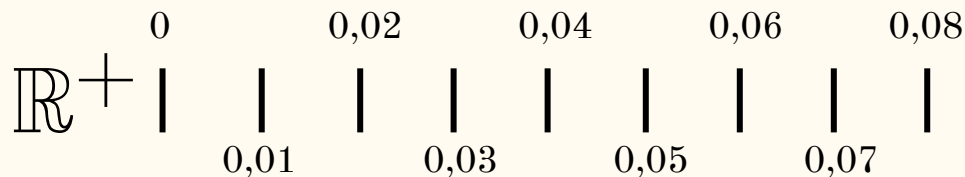
000,03

...

999,99

$$0,07 \div 2 = 0,04$$

$$EA \leq 0,01$$



Representação em Ponto Flutuante



000,00

000,01

000,02

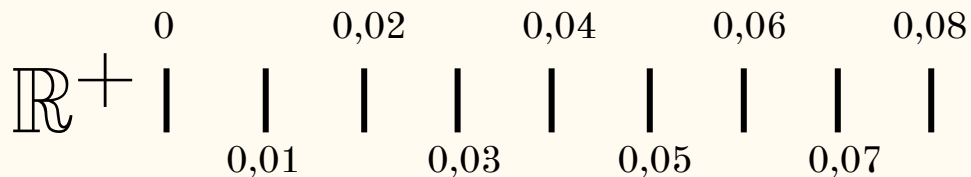
000,03

...

999,99

856,786 \rightarrow 856,79

$$\begin{aligned} ER &= 0,004 / 856,786 \\ &= 0,000005 \end{aligned}$$



Representação em Ponto Flutuante



000,00

000,01

000,02

000,03

...

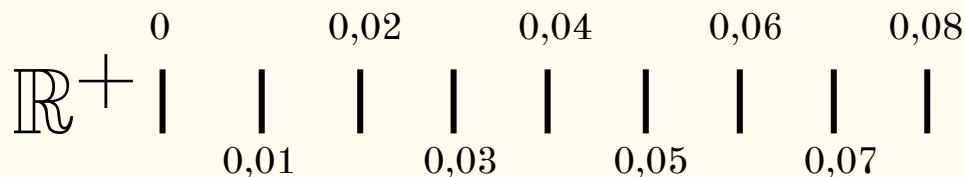
999,99

856,786 \rightarrow 856,79

$$\begin{aligned} ER &= 0,004 / 856,786 \\ &= 0,000005 \end{aligned}$$

23,546 \rightarrow 23,55

$$\begin{aligned} ER &= 0,004 / 23,546 \\ &= 0,0002 \end{aligned}$$



Representação em Ponto Flutuante



$$856,786 \rightarrow 8,56786 \times 10^2$$

$$23,546 \rightarrow 2,3546 \times 10^1$$

Representação em Ponto Flutuante



$$856,786 \rightarrow 8,56786 \times 10^2$$

$$23,546 \rightarrow 2,3546 \times 10^1$$

Representação em Ponto Flutuante



$$\begin{array}{l} 856,786 \rightarrow 8,56786 \times 10^2 \rightarrow 8,568 \times 10^2 \\ 23,546 \rightarrow 2,3546 \times 10^1 \rightarrow 2,355 \times 10^1 \end{array}$$

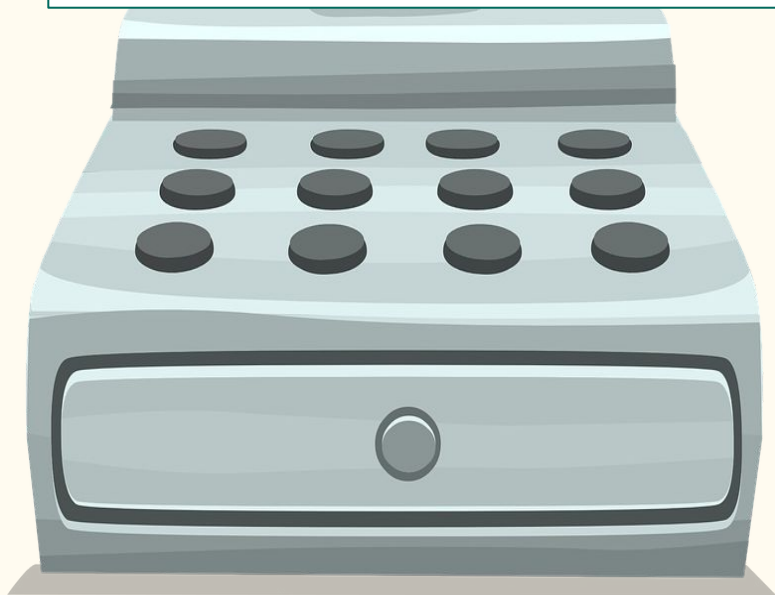
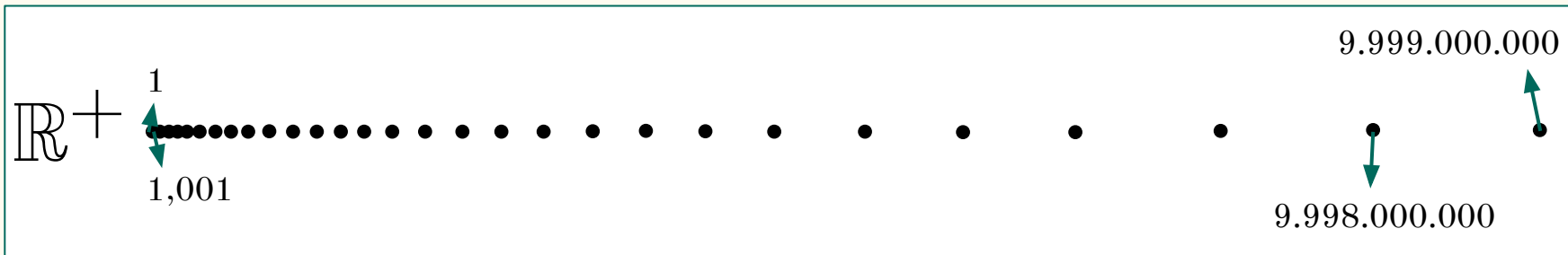
Representação em Ponto Flutuante



$856,786 \rightarrow 8,56786 \times 10^2 \rightarrow 8,568 \times 10^2$
 $23,546 \rightarrow 2,3546 \times 10^1 \rightarrow 2,355 \times 10^1$

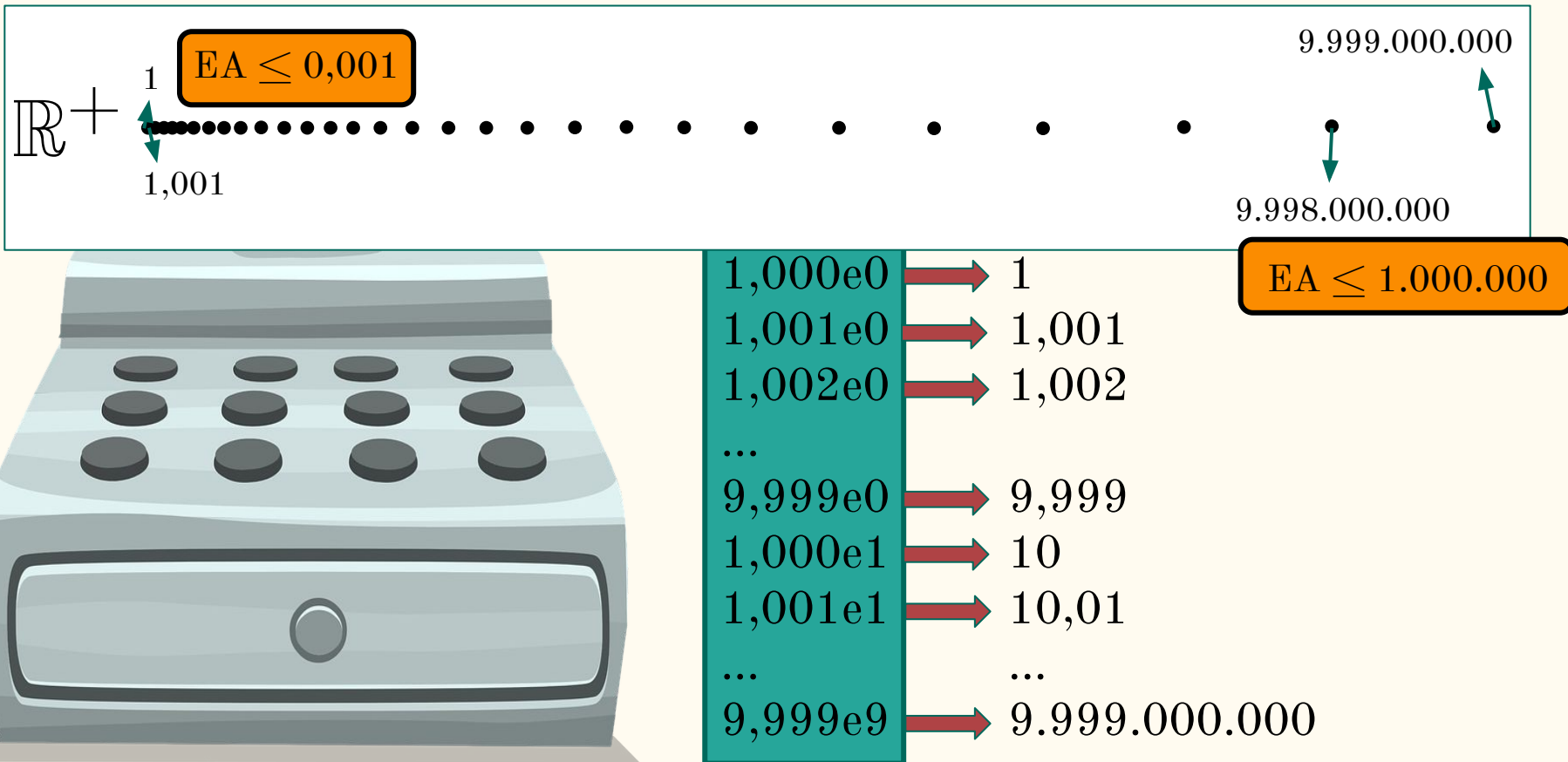
1,000e0	→	1
1,001e0	→	1,001
1,002e0	→	1,002
...		
9,999e0	→	9,999
1,000e1	→	10
1,001e1	→	10,01
...		...
9,999e9	→	9.999.000.000

Representação em Ponto Flutuante

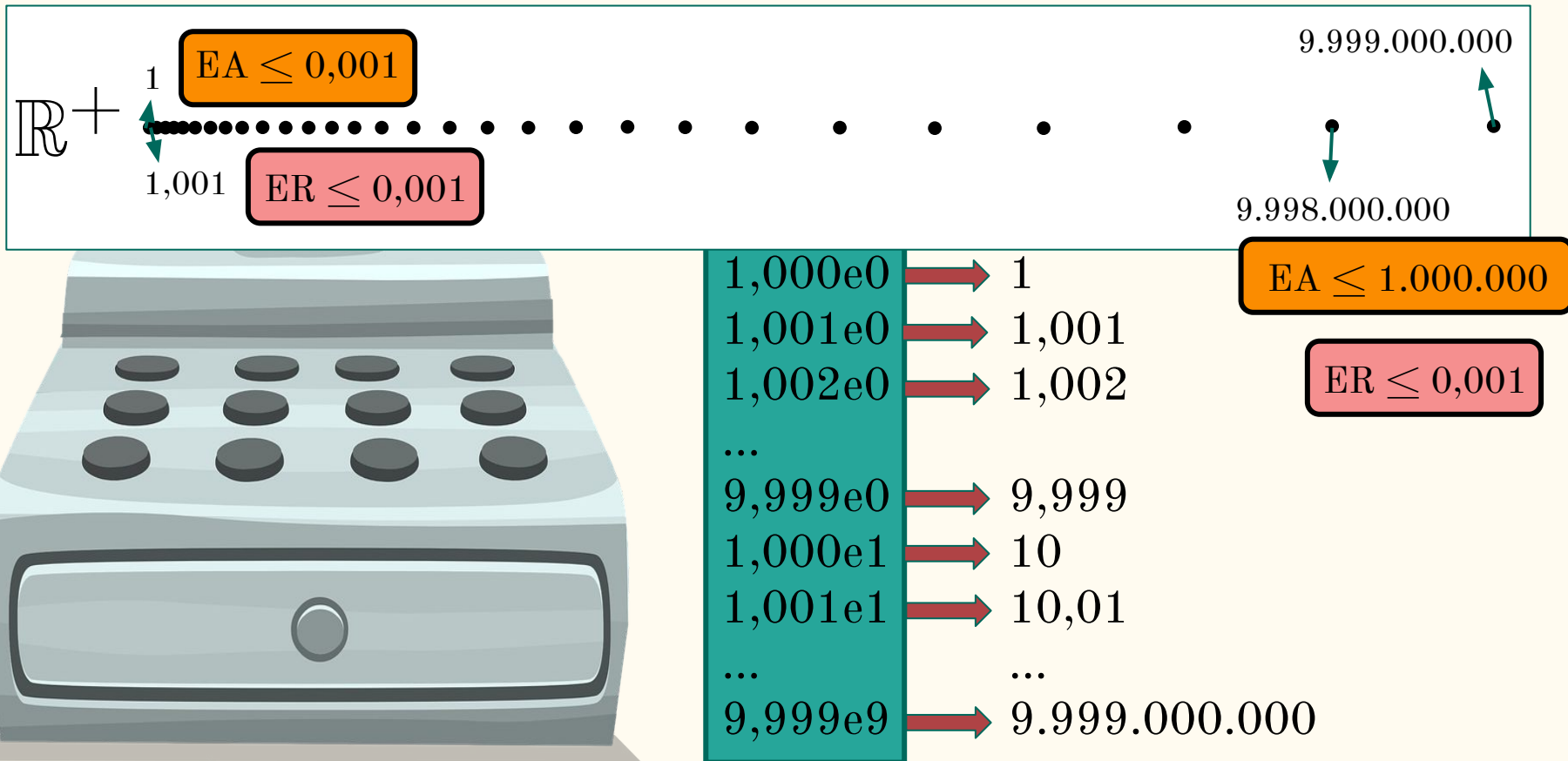


1,000e0	→	1
1,001e0	→	1,001
1,002e0	→	1,002
...		
9,999e0	→	9,999
1,000e1	→	10
1,001e1	→	10,01
...		...
9,999e9	→	9.999.000.000

Representação em Ponto Flutuante



Representação em Ponto Flutuante



Representação em Ponto Flutuante



$$8,56786 \times 10^2 \longrightarrow 8,568 \times 10^2$$

$$\begin{aligned} ER &= 0,014 / 856,786 \\ &= 0,00002 \end{aligned}$$

$$2,3546 \times 10^1 \longrightarrow 2,355 \times 10^1$$

$$\begin{aligned} ER &= 0,004 / 23,546 \\ &= 0,0002 \end{aligned}$$

Representação em Ponto Flutuante



$$d_1 d_2 d_3 \dots d_t \times 10^e$$

$(d_1 d_2 d_3 \dots d_t)$ é a mantissa com t dígitos

$0 \leq d_i \leq 9$, para todo $i = 1 \dots t$

$e \in [m, M]$ é o expoente (geralmente $m = -M$)

$d_1 > 0$ (números normalizados)

Representação em Ponto Flutuante



$$d_1, d_2 d_3 \dots d_t \times 10^e$$

$(d_1 d_2 d_3 \dots d_t)$ é a mantissa com t dígitos

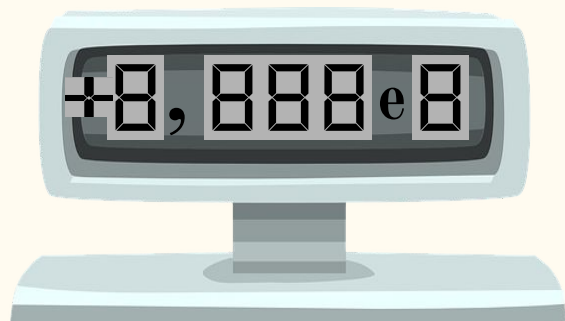
$0 \leq d_i \leq 9$, para todo $i = 1 \dots t$

$e \in [m, M]$ é o expoente (geralmente $m = -M$)

$d_1 > 0$ (números normalizados)

Como representar o zero?

Representação em Ponto Flutuante



- **Números negativos?**
“display” de sinal (\pm)

$$\pm(d_1d_2d_3\dots d_t) \times 10^e$$

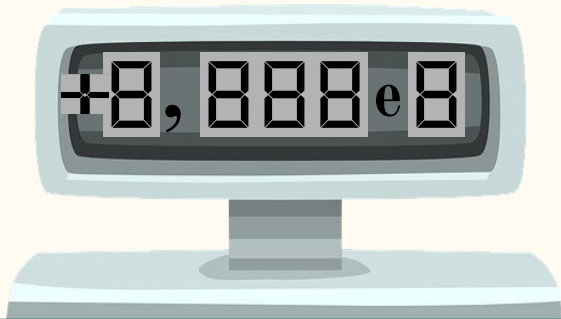
$(d_1d_2d_3\dots d_t)$ é a mantissa com t dígitos

$0 \leq d_i \leq 9$, para todo $i = 1\dots t$

$e \in [m, M]$ é o expoente (geralmente $m = -M$)

$d_1 > 0$ (números normalizados)

Representação em Ponto Flutuante



- **Números negativos?**
“display” de sinal (\pm)
- **Números menores que 1?**
expoente negativo: define-se um desvio Δ e subtrai do expoente

$$\pm(d_1d_2d_3\dots d_t) \times 10^{e'-\Delta}$$

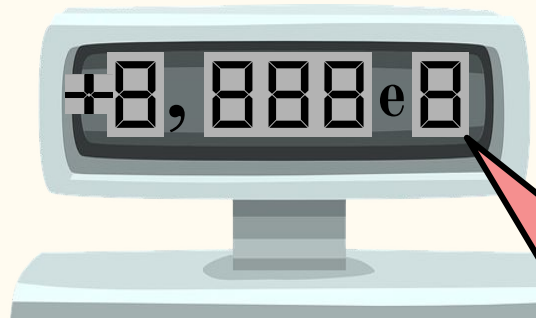
$(d_1d_2d_3\dots d_t)$ é a mantissa com t dígitos

$$0 \leq d_i \leq 9, \text{ para todo } i = 1\dots t$$

$e \in [m, M]$ é o expoente (geralmente $m = -M$)

$$d_1 > 0 \quad (\text{números } \underline{\text{normalizados}})$$

Representação em Ponto Flutuante



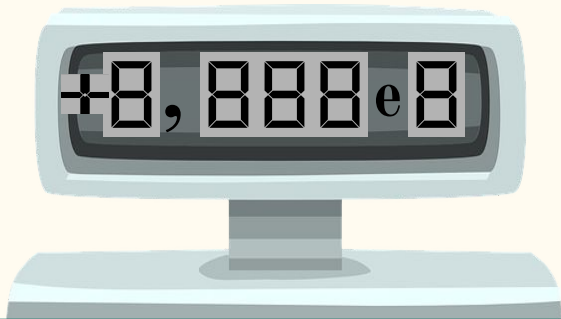
- **Números negativos?**
“display” de sinal (\pm)
- **Números menores que 1?**
expoente negativo: define-se um desvio Δ e subtrai do expoente

$$\pm(d_1d_2d_3\dots d_t) \times 10^{e'-\Delta}$$

$$\Delta = 4$$

$e' = 0$	10^{-4}	$e' = 5$	10^1
$e' = 1$	10^{-3}	$e' = 6$	10^2
$e' = 2$	10^{-2}	$e' = 7$	10^3
$e' = 3$	10^{-1}	$e' = 8$	10^4
$e' = 4$	10^0	$e' = 9$	10^5

Representação em Ponto Flutuante



- **Números negativos?**
“display” de sinal (\pm)
- **Números menores que 1?**
expoente negativo: define-se um desvio Δ e subtrai do expoente
- **Outras bases?**

$$\pm(d_1d_2d_3\dots d_t) \times \beta^{e'-\Delta}$$

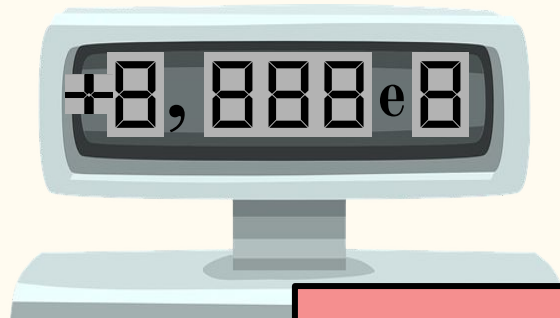
$(d_1d_2d_3\dots d_t)$ é a mantissa com t dígitos

$$0 \leq d_i \leq \beta-1, \text{ para todo } i = 1\dots t$$

$e \in [m, M]$ é o expoente (geralmente $m = -M$)

$$d_1 > 0 \quad (\text{números } \underline{\text{normalizados}})$$

Representação em Ponto Flutuante



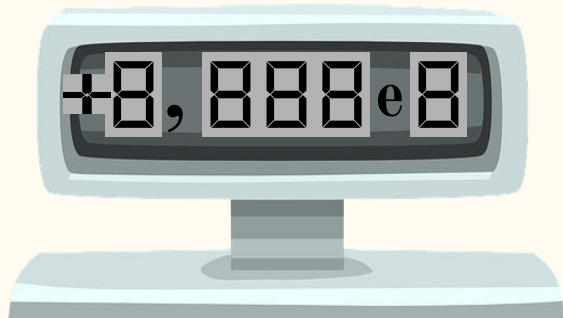
$$\pm(d_1d_2d_3\dots d_t) \times \beta^{e'-\Delta}$$

$(d_1d_2d_3\dots d_t)$ é a mantissa com t dígitos

- **Números negativos** $1101,011101 \times 2^0 = 11,01011101 \times 10^0$
 “display” de significando 2^2 geralmente $m =$
- **Números menores** 10^{-m}
 expoente negativo: define-se um desvio Δ e subtrai do expoente
- **Outras bases?**

$$d_1 > 0 \quad (\text{números } \underline{\text{normalizados}})$$

Representação em Ponto Flutuante



- **Números negativos?**
“display” de sinal (\pm)
- **Números menores que 1?**
expoente negativo: define-se um desvio Δ e subtrai do expoente
- **Outras bases?**

$$\pm(d_1 d_2 d_3 \dots d_t) \times \beta^{e' - \Delta}$$

$(d_1 d_2 d_3 \dots d_t)$ é a mantissa com t dígitos

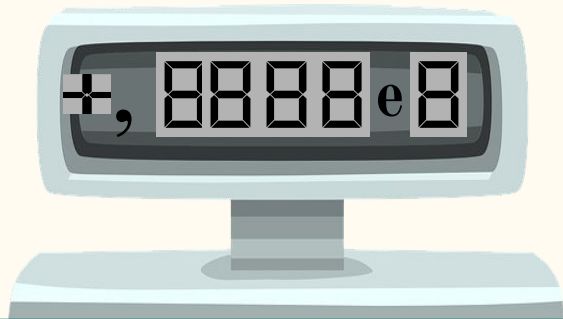
$0 \leq d_i \leq \beta - 1$, para todo $i = 1 \dots t$

$e \in [m, M]$ é o expoente (geralmente $m = -M$)

$d_1 > 0$ (números normalizados)

SPF(β, t, m, M)

Representação em Ponto Flutuante



- **Números negativos?**
“display” de sinal (\pm)
- **Números menores que 1?**
expoente negativo: define-se um desvio Δ e subtrai do expoente
- **Outras bases?**

$$\pm(0,d_1d_2d_3\dots d_t) \times \beta^{e'-\Delta}$$

$(d_1d_2d_3\dots d_t)$ é a mantissa com t dígitos

$0 \leq d_i \leq \beta-1$, para todo $i = 1\dots t$

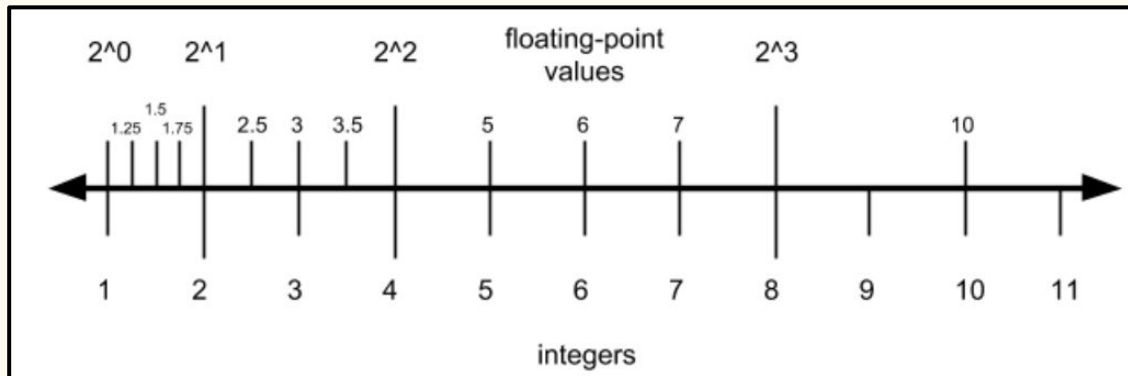
$e \in [m, M]$ é o expoente (geralmente $m = -M$)

$d_1 > 0$ (números normalizados)

SPF(β, t, m, M)

Exemplo

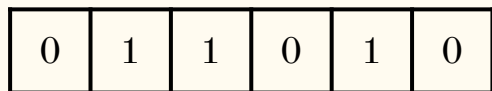
SPF($\beta=2$, $t=3$, $m=0$, $M=3$)



$e = 0$	$e = 1$	$e = 2$	$e = 3$
$1,00 \times 2^0 = 1_{10}$	$1,00 \times 2^1 = 2_{10}$	$1,00 \times 2^2 = 4_{10}$	$1,00 \times 2^3 = 8_{10}$
$1,01 \times 2^0 = 1,25_{10}$	$1,01 \times 2^1 = 2,5_{10}$	$1,01 \times 2^2 = 5_{10}$	$1,01 \times 2^3 = 10_{10}$
$1,10 \times 2^0 = 1,5_{10}$	$1,10 \times 2^1 = 3_{10}$	$1,10 \times 2^2 = 6_{10}$	$1,10 \times 2^3 = 12_{10}$
$1,11 \times 2^0 = 1,75_{10}$	$1,11 \times 2^1 = 3,5_{10}$	$1,11 \times 2^2 = 7_{10}$	$1,11 \times 2^3 = 14_{10}$

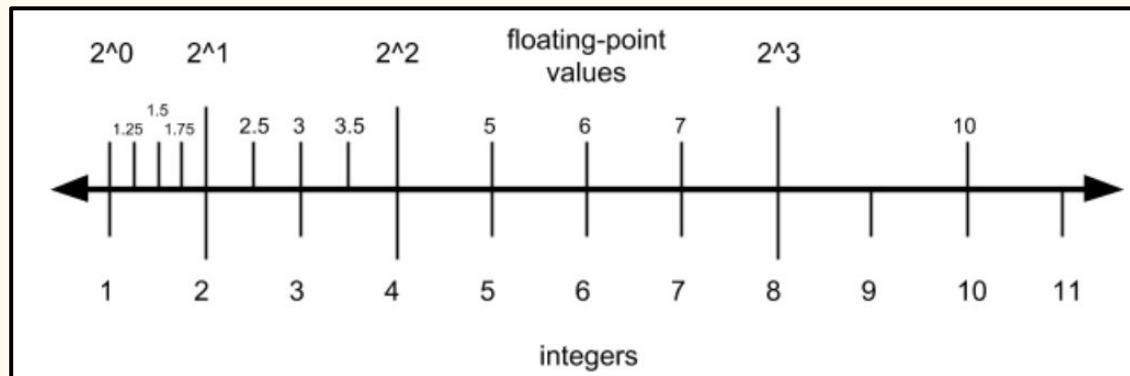
Exemplo

SPF($\beta=2$, $t=3$, $m=0$, $M=3$)



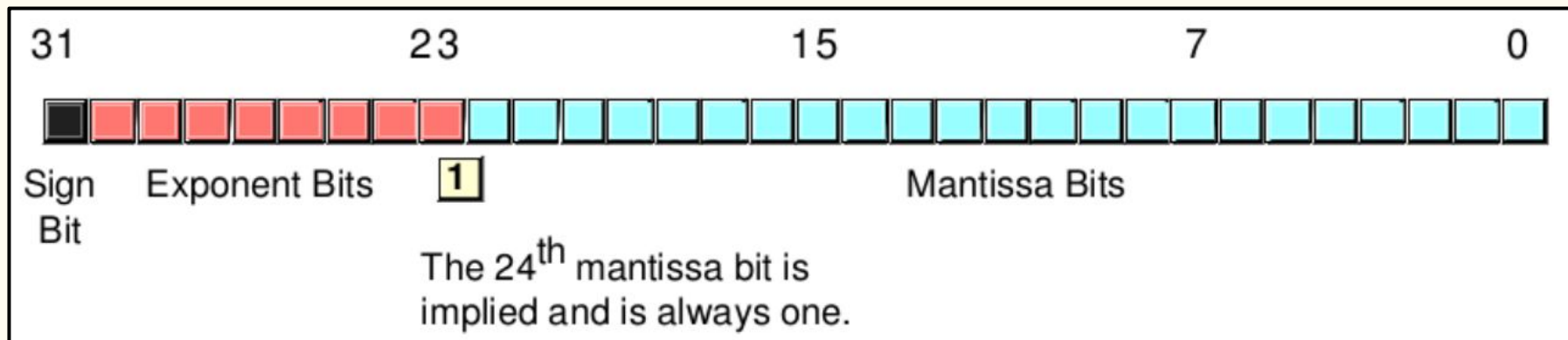
sinal
mantissa
expoente

$$(-1)^0 \times 1,10 \times 2^{(10 - \Delta)}$$



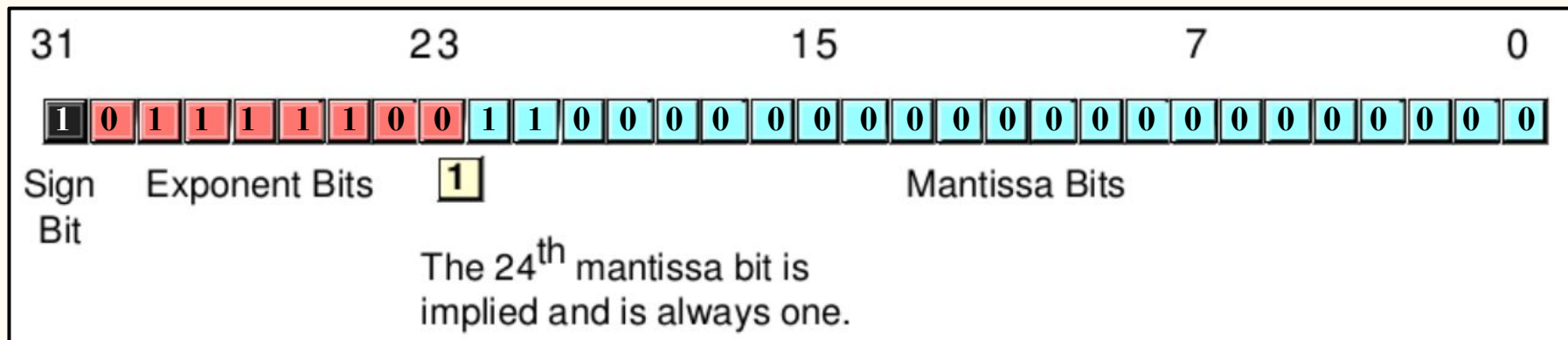
	e = 1	e = 2	e = 3
$1,00 \times 2^0 = 1_{10}$	$1,00 \times 2^1 = 2_{10}$	$1,00 \times 2^2 = 4_{10}$	$1,00 \times 2^3 = 8_{10}$
$1,01 \times 2^0 = 1,25_{10}$	$1,01 \times 2^1 = 2,5_{10}$	$1,01 \times 2^2 = 5_{10}$	$1,01 \times 2^3 = 10_{10}$
$1,10 \times 2^0 = 1,5_{10}$	$1,10 \times 2^1 = 3_{10}$	$1,10 \times 2^2 = 6_{10}$	$1,10 \times 2^3 = 12_{10}$
$1,11 \times 2^0 = 1,75_{10}$	$1,11 \times 2^1 = 3,5_{10}$	$1,11 \times 2^2 = 7_{10}$	$1,11 \times 2^3 = 14_{10}$

Formato IEEE 754



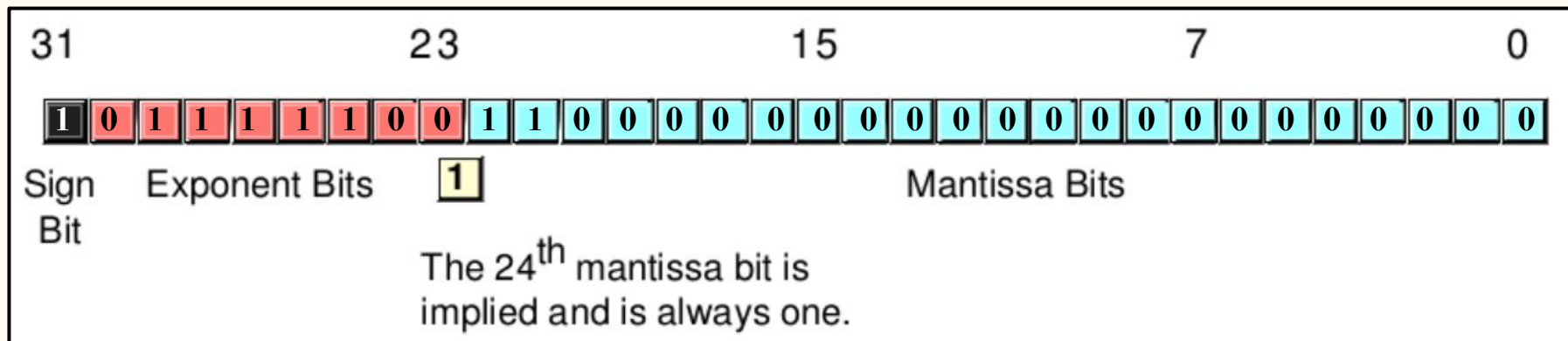
$$(-1)^{\text{Sign}} \times (1 + \text{Mantissa}) \times 2^{(\text{Exponent} - \Delta)}, \text{ onde } \Delta = 127 \text{ para float e } \Delta = 1023 \text{ para double}$$

Formato IEEE 754



$$(-1)^{\text{Sign}} \times (1 + \text{Mantissa}) \times 2^{(\text{Exponent} - \Delta)}, \text{ onde } \Delta = 127 \text{ para float e } \Delta = 1023 \text{ para double}$$

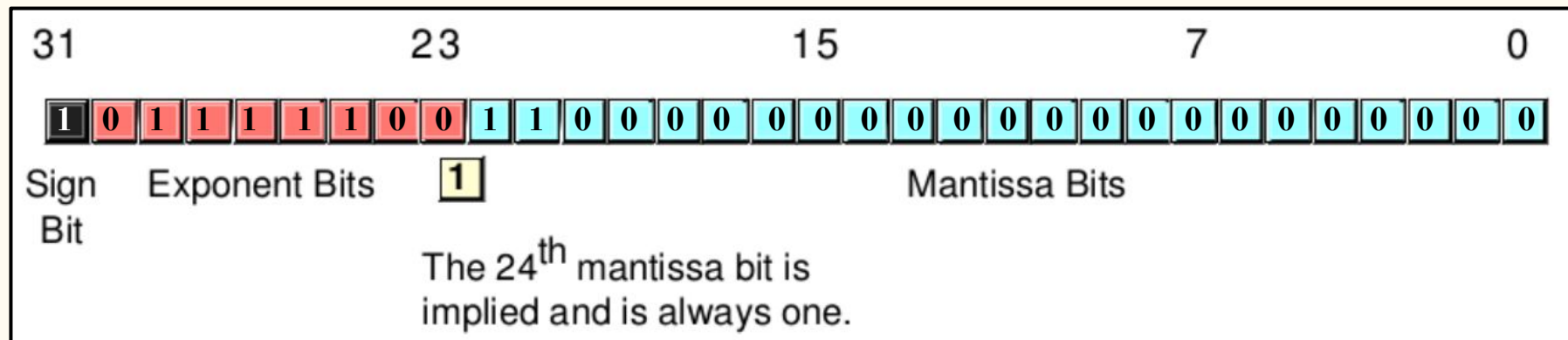
Formato IEEE 754



$$(-1)^{\text{Sign}} \times (1 + \text{Mantissa}) \times 2^{(\text{Exponent} - \Delta)}, \text{ onde } \Delta = 127 \text{ para float e } \Delta = 1023 \text{ para double}$$

Exemplo: $(-1)^1$

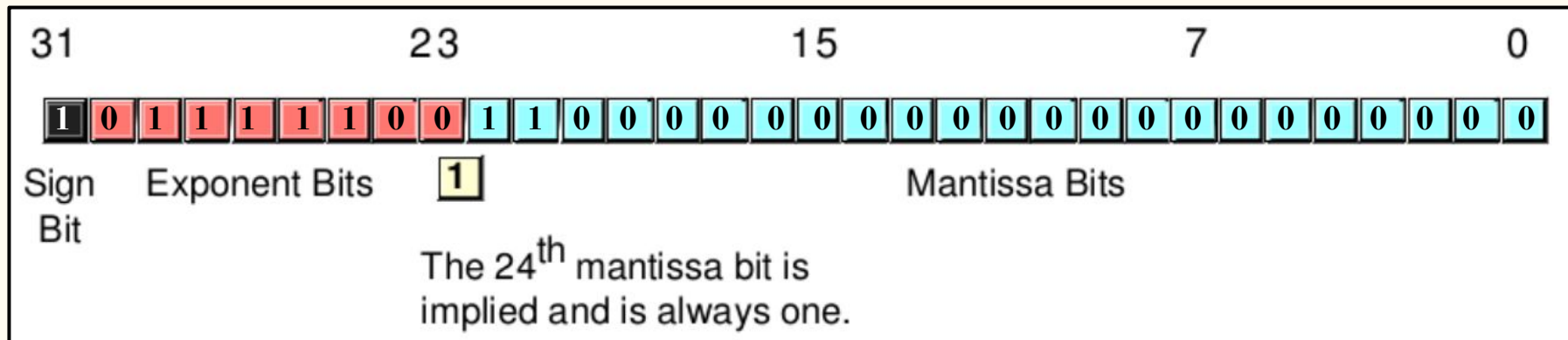
Formato IEEE 754



$$(-1)^{\text{Sign}} \times (1 + \text{Mantissa}) \times 2^{(\text{Exponent} - \Delta)}, \text{ onde } \Delta = 127 \text{ para float e } \Delta = 1023 \text{ para double}$$

Exemplo: $(-1)^1 \times (1 + 0,75)$

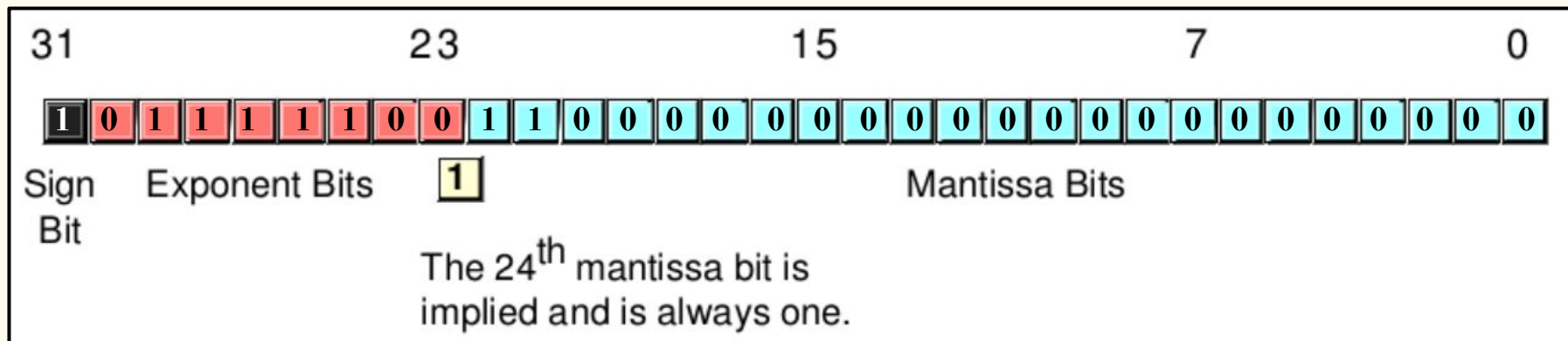
Formato IEEE 754



$$(-1)^{\text{Sign}} \times (1 + \text{Mantissa}) \times 2^{(\text{Exponent} - \Delta)}, \text{ onde } \Delta = 127 \text{ para float e } \Delta = 1023 \text{ para double}$$

Exemplo: $(-1)^1 \times (1 + 0,75) \times 2^{(124 - 127)}$

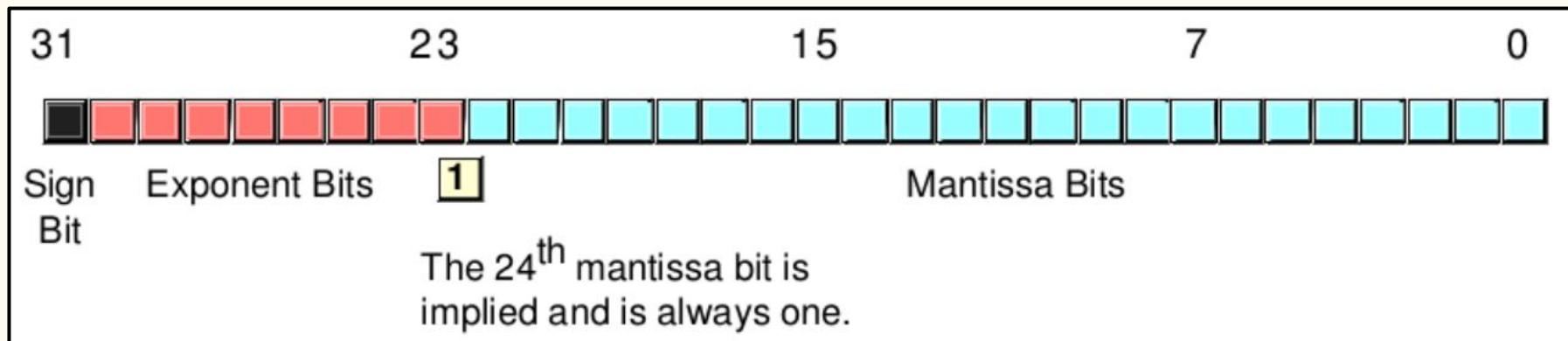
Formato IEEE 754



$$(-1)^{\text{Sign}} \times (1 + \text{Mantissa}) \times 2^{(\text{Exponent} - \Delta)}, \text{ onde } \Delta = 127 \text{ para float e } \Delta = 1023 \text{ para double}$$

Exemplo: $(-1)^1 \times (1 + 0,75) \times 2^{(124 - 127)} = -1,75 \times 2^{-3} = -0,21875$

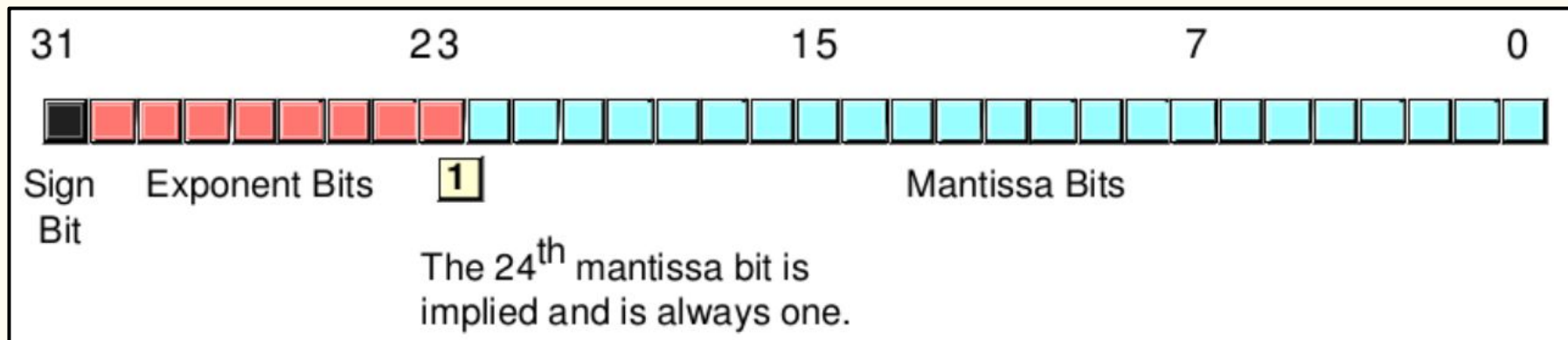
Formato IEEE 754



$$(-1)^{\text{Sign}} \times (1 + \text{Mantissa}) \times 2^{(\text{Exponent} - \Delta)}, \text{ onde } \Delta = 127 \text{ para float e } \Delta = 1023 \text{ para double}$$

Exemplo: $639,6875 = 1001111111,1011_2 = 1,0011111111011 \times 2^9$

Formato IEEE 754

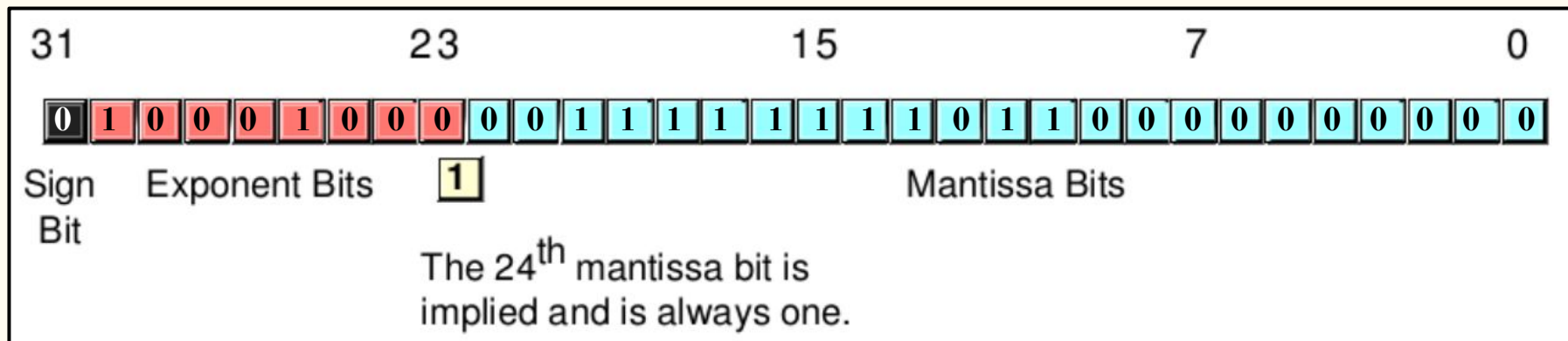


$$(-1)^{\text{Sign}} \times (1 + \text{Mantissa}) \times 2^{(\text{Exponent} - \Delta)}, \text{ onde } \Delta = 127 \text{ para float e } \Delta = 1023 \text{ para double}$$

Exemplo: $639,6875 = 1001111111,1011_2 = 1,0011111111011 \times 2^9$

Sign = 0 Exponent = $9 + 127 = 136 = 10001000_2$ Mantissa =
0011111111011

Formato IEEE 754

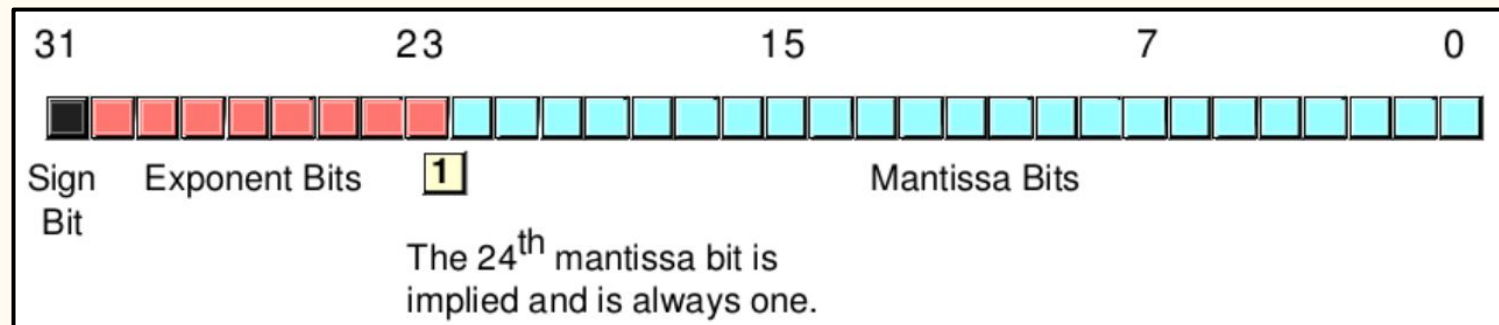


$$(-1)^{\text{Sign}} \times (1 + \text{Mantissa}) \times 2^{(\text{Exponent} - \Delta)}, \text{ onde } \Delta = 127 \text{ para float e } \Delta = 1023 \text{ para double}$$

Exemplo: $639,6875 = 1001111111,1011_2 = 1,0011111111011 \times 2^9$

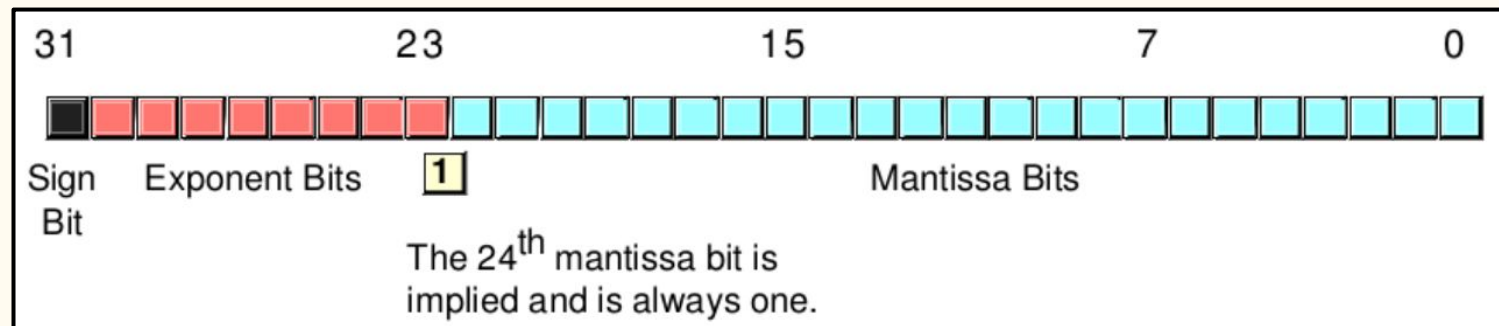
Sign = 0 Exponent = $9 + 127 = 136 = 10001000_2$ Mantissa =
0011111111011

Formato IEEE 754



	float	double
Bits de precisão (mantissa)	23+1	52+1
Precisão decimal (dígitos decimais)	6,5	14,5
Bits do expoente	8	11
Expoente máximo	127	1023
Expoente mínimo	-126	-1022
Deslocamento do expoente	127	1023
Maior/menor número	$10^{\pm 38}$	$10^{\pm 308}$

Formato IEEE 754



	float	double
Bits de precisão (mantissa)	23+1	52+1
Precisão decimal (dígitos decimais)	6,5	14,5
Bits do expoente	8	11
Expoente máximo	127	1023
Expoente mínimo	-126	-1022
Deslocamento do expoente	127	1023
Maior/menor número	$10^{\pm 38}$	$10^{\pm 308}$

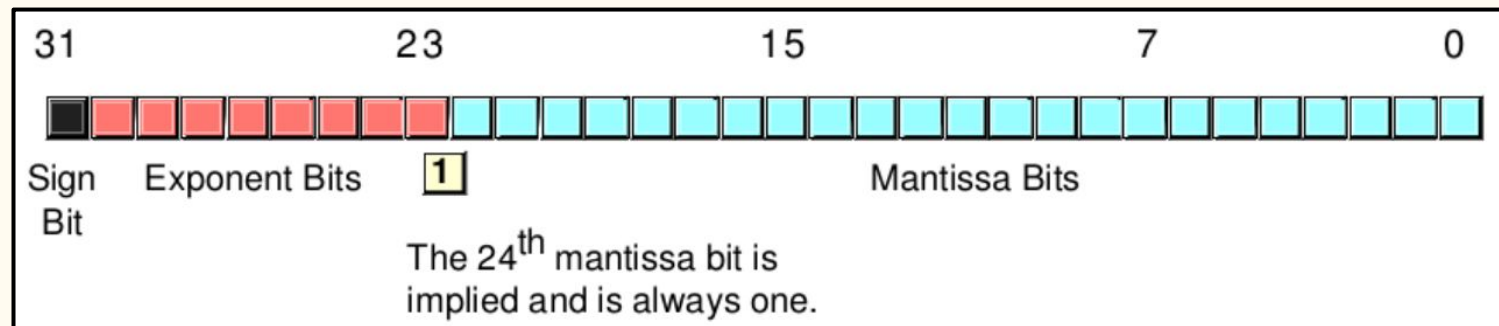
$$m = 00000000_2 - 127_{10}$$

$$M = 11111111_2 - 127_{10}$$

$$m = -127$$

$$M = 128$$

Formato IEEE 754



Se todos os bits do expoente são 1

Se todos os bits da mantissa são 0

Representa $\pm\text{infinity}$

Ex: 1/0

Senão

Representa $\pm\text{NaN}$

Ex: $\log(-1)$

Se todos os bits do expoente são 0

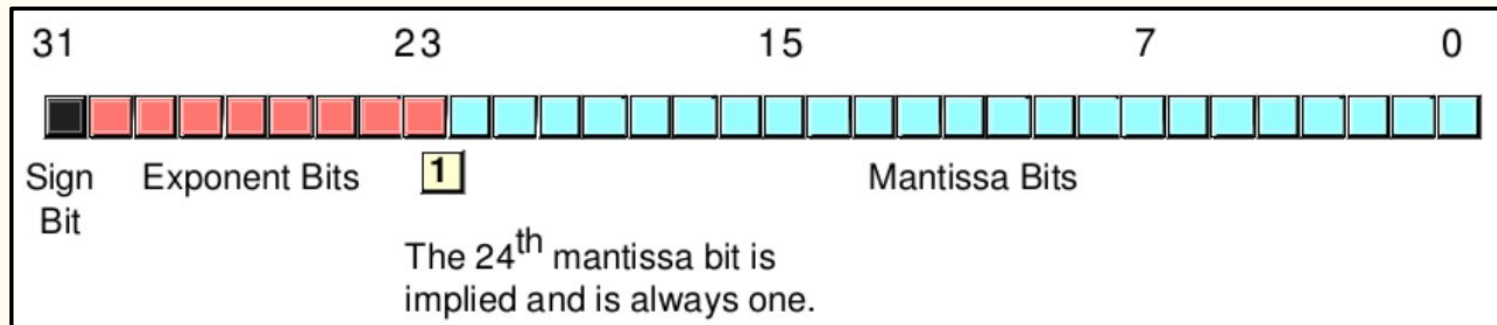
Se todos os bits da mantissa são 0

Representa ± 0

Senão

Representa um número não normalizado

Formato IEEE 754



Se todos os bits do expoente são 1

Se todos os bits da mantissa são 0

Representa $\pm\text{infinity}$

Ex: 1/0

Senão

Representa $\pm\text{NaN}$

Ex: $\log(-1)$

Se todos os bits do expoente são 0

~~$$(-1)^{\text{Sign}} \times (1 + \text{Mantissa}) \times 2^{(\text{Exponent} - \Delta)}$$~~

$$(-1)^{\text{Sign}} \times (0 + \text{Mantissa}) \times 2^{(-126)}$$

Senão

Representa um número não normalizado

Atenção!

A ordem das operações pode afetar a acurácia do resultado (operações sucessivas que acumulam erros)

$$SPF(\beta, t, m, M) = SPF(10, 3, -9, 9)$$

○ Ex: $x = 1,23 \times 10^3$, $y = 1,00 \times 10^0$

```
z = 0;  
for (int i=0; i<10; ++i)  
    z += y  
z += x
```

```
z = x;  
for (int i=0; i<10; ++i)  
    z += y
```

Atenção!

A ordem das operações pode afetar a acurácia do resultado (operações sucessivas que acumulam erros)

$SPF(\beta, t, m, M) = SPF(10, 3, -9, 9)$

○ Ex: $x = 1,23 \times 10^3$, $y = 1,00 \times 10^0$

```
z = 0;  
for (int i=0; i<10; ++i)  
    z += y  
z += x
```

```
z = x;  
for (int i=0; i<10; ++i)  
    z += y
```

$i = 0$

$$z = 0,00 \times 10^0 + 1,00 \times 10^0$$

$$z = 1,00 \times 10^0$$

$$z = 1,23 \times 10^3 + 1,00 \times 10^0$$

$$z = 1,23 \times 10^3$$

Atenção!

A ordem das operações pode afetar a acurácia do resultado (operações sucessivas que acumulam erros)

$SPF(\beta, t, m, M) = SPF(10, 3, -9, 9)$

Ex: $x = 1,23 \times 10^3$, $y = 1,00 \times 10^0$

```
z = 0;  
for (int i=0; i<10; ++i)  
    z += y  
z += x
```

```
z = x;  
for (int i=0; i<10; ++i)  
    z += y
```

$i = 0$

$$z = 0,00 \times 10^0 + 1,00 \times 10^0$$

$$z = 1,00 \times 10^0$$

$$z = 1,23 \times 10^3 + 1,00 \times 10^0$$

$$z = 1,23 \times 10^3$$

$$(1,23 + 0,001) \times 10^3$$

$$1,231 \times 10^3$$

Atenção!

A ordem das operações pode afetar a acurácia do resultado (operações sucessivas que acumulam erros)

$SPF(\beta, t, m, M) = SPF(10, 3, -9, 9)$

○ Ex: $x = 1,23 \times 10^3$, $y = 1,00 \times 10^0$

```
z = 0;  
for (int i=0; i<10; ++i)  
    z += y  
z += x
```

```
z = x;  
for (int i=0; i<10; ++i)  
    z += y
```

$i = 1$

$$z = 1,00 \times 10^0 + 1,00 \times 10^0$$

$$z = 2,00 \times 10^0$$

$$z = 1,23 \times 10^3 + 1,00 \times 10^0$$

$$z = 1,23 \times 10^3$$

Atenção!

A ordem das operações pode afetar a acurácia do resultado (operações sucessivas que acumulam erros)

$SPF(\beta, t, m, M) = SPF(10, 3, -9, 9)$

○ Ex: $x = 1,23 \times 10^3$, $y = 1,00 \times 10^0$

```
z = 0;  
for (int i=0; i<10; ++i)  
    z += y  
z += x
```

```
z = x;  
for (int i=0; i<10; ++i)  
    z += y
```

$i = 2$

$$z = 2,00 \times 10^0 + 1,00 \times 10^0$$

$$z = 3,00 \times 10^0$$

$$z = 1,23 \times 10^3 + 1,00 \times 10^0$$

$$z = 1,23 \times 10^3$$

Atenção!

A ordem das operações pode afetar a acurácia do resultado (operações sucessivas que acumulam erros)

$SPF(\beta, t, m, M) = SPF(10, 3, -9, 9)$

○ Ex: $x = 1,23 \times 10^3$, $y = 1,00 \times 10^0$

```
z = 0;  
for (int i=0; i<10; ++i)  
    z += y  
z += x
```

```
z = x;  
for (int i=0; i<10; ++i)  
    z += y
```

$i = 3$

$$z = 3,00 \times 10^0 + 1,00 \times 10^0$$

$$z = 4,00 \times 10^0$$

$$z = 1,23 \times 10^3 + 1,00 \times 10^0$$

$$z = 1,23 \times 10^3$$

Atenção!

A ordem das operações pode afetar a acurácia do resultado (operações sucessivas que acumulam erros)

$SPF(\beta, t, m, M) = SPF(10, 3, -9, 9)$

○ Ex: $x = 1,23 \times 10^3$, $y = 1,00 \times 10^0$

```
z = 0;  
for (int i=0; i<10; ++i)  
    z += y  
z += x
```

```
z = x;  
for (int i=0; i<10; ++i)  
    z += y
```

$i = 4$

$$z = 4,00 \times 10^0 + 1,00 \times 10^0$$

$$z = 5,00 \times 10^0$$

$$z = 1,23 \times 10^3 + 1,00 \times 10^0$$

$$z = 1,23 \times 10^3$$

Atenção!

A ordem das operações pode afetar a acurácia do resultado (operações sucessivas que acumulam erros)

$SPF(\beta, t, m, M) = SPF(10, 3, -9, 9)$

○ Ex: $x = 1,23 \times 10^3$, $y = 1,00 \times 10^0$

```
z = 0;  
for (int i=0; i<10; ++i)  
    z += y  
z += x
```

```
z = x;  
for (int i=0; i<10; ++i)  
    z += y
```

$i = 5$

$$z = 5,00 \times 10^0 + 1,00 \times 10^0$$

$$z = 6,00 \times 10^0$$

$$z = 1,23 \times 10^3 + 1,00 \times 10^0$$

$$z = 1,23 \times 10^3$$

Atenção!

A ordem das operações pode afetar a acurácia do resultado (operações sucessivas que acumulam erros)

$SPF(\beta, t, m, M) = SPF(10, 3, -9, 9)$

○ Ex: $x = 1,23 \times 10^3$, $y = 1,00 \times 10^0$

```
z = 0;  
for (int i=0; i<10; ++i)  
    z += y  
z += x
```

```
z = x;  
for (int i=0; i<10; ++i)  
    z += y
```

$i = 6$

$$z = 6,00 \times 10^0 + 1,00 \times 10^0$$

$$z = 7,00 \times 10^0$$

$$z = 1,23 \times 10^3 + 1,00 \times 10^0$$

$$z = 1,23 \times 10^3$$

Atenção!

A ordem das operações pode afetar a acurácia do resultado (operações sucessivas que acumulam erros)

$SPF(\beta, t, m, M) = SPF(10, 3, -9, 9)$

○ Ex: $x = 1,23 \times 10^3$, $y = 1,00 \times 10^0$

```
z = 0;  
for (int i=0; i<10; ++i)  
    z += y  
z += x
```

```
z = x;  
for (int i=0; i<10; ++i)  
    z += y
```

$i = 7$

$$z = 7,00 \times 10^0 + 1,00 \times 10^0$$

$$z = 8,00 \times 10^0$$

$$z = 1,23 \times 10^3 + 1,00 \times 10^0$$

$$z = 1,23 \times 10^3$$

Atenção!

A ordem das operações pode afetar a acurácia do resultado (operações sucessivas que acumulam erros)

$SPF(\beta, t, m, M) = SPF(10, 3, -9, 9)$

○ Ex: $x = 1,23 \times 10^3$, $y = 1,00 \times 10^0$

```
z = 0;  
for (int i=0; i<10; ++i)  
    z += y  
z += x
```

```
z = x;  
for (int i=0; i<10; ++i)  
    z += y
```

$i = 8$

$$z = 8,00 \times 10^0 + 1,00 \times 10^0$$

$$z = 9,00 \times 10^0$$

$$z = 1,23 \times 10^3 + 1,00 \times 10^0$$

$$z = 1,23 \times 10^3$$

Atenção!

A ordem das operações pode afetar a acurácia do resultado (operações sucessivas que acumulam erros)

$SPF(\beta, t, m, M) = SPF(10, 3, -9, 9)$

○ Ex: $x = 1,23 \times 10^3$, $y = 1,00 \times 10^0$

```
z = 0;  
for (int i=0; i<10; ++i)  
    z += y  
z += x
```

```
z = x;  
for (int i=0; i<10; ++i)  
    z += y
```

$i = 9$

$$z = 9,00 \times 10^0 + 1,00 \times 10^0$$

$$z = 1,00 \times 10^1$$

$$z = 1,23 \times 10^3 + 1,00 \times 10^0$$

$$z = 1,23 \times 10^3$$

Atenção!

A ordem das operações pode afetar a acurácia do resultado (operações sucessivas que acumulam erros)

$SPF(\beta, t, m, M) = SPF(10, 3, -9, 9)$

○ Ex: $x = 1,23 \times 10^3$, $y = 1,00 \times 10^0$

```
z = 0;  
for (int i=0; i<10; ++i)  
    z += y  
z += x
```

```
z = x;  
for (int i=0; i<10; ++i)  
    z += y
```

$$z = 1,00 \times 10^1 + 1,23 \times 10^3$$

$$z = 1,24 \times 10^3$$

$$z = 1,23 \times 10^3$$

Atenção!

Cadeia de cálculos envolvendo $+$, $-$, \times e \div  efetuar \times e \div primeiro.

Atenção!

Cadeia de cálculos envolvendo $+$, $-$, \times e \div  efetuar \times e \div primeiro.


$$x = 1,14 \times 10^1$$

$$y = 3,18 \times 10^0$$

$$z = 5,05 \times 10^0$$

SPF($\beta=10$, $t=3$, $m=-5$, $M=5$)

Atenção!

Cadeia de cálculos envolvendo $+$, $-$, \times e \div  efetuar \times e \div primeiro.

$$x = 1,14 \times 10^1$$

$$y = 3,18 \times 10^0$$

$$z = 5,05 \times 10^0$$

SPF($\beta=10$, $t=3$, $m=-5$, $M=5$)

$$y \times (z + x)$$

$$y \times (0,505 + 1,14) \times 10^1$$

$$(3,18 \times 10^0) \times (1,64 \times 10^1)$$

$$(3,18 \times 1,64) \times 10^{(0+1)}$$

$$5,21 \times 10^1$$

Atenção!

Cadeia de cálculos envolvendo $+$, $-$, \times e \div  efetuar \times e \div primeiro.

$$\begin{aligned}x &= 1,14 \times 10^1 \\y &= 3,18 \times 10^0 \\z &= 5,05 \times 10^0\end{aligned}$$

SPF($\beta=10$, $t=3$, $m=-5$, $M=5$)

$$y \times (z + x)$$

$$y \times (0,505 + 1,14) \times 10^1$$

$$(3,18 \times 10^0) \times (1,64 \times 10^1)$$

$$(3,18 \times 1,64) \times 10^{(0+1)}$$

$$5,21 \times 10^1$$

$$(y \times z) + (y \times x)$$

$$(3,18 \times 10^0) \times (5,05 \times 10^0) + (3,18 \times 10^0) \times (1,14 \times 10^1)$$

$$(3,18 \times 5,05) \times 10^{(0+0)} + (3,18 \times 1,14) \times 10^{(0+1)}$$

$$1,60 \times 10^1 + 3,62 \times 10^1$$

$$(1,60 + 3,62) \times 10^1$$

$$5,22 \times 10^1$$

Atenção!

Cadeia de cálculos envolvendo $+$, $-$, \times e \div  efetuar \times e \div primeiro.

$$x = 1,14 \times 10^1$$

$$y = 3,18 \times 10^0$$

$$z = 5,05 \times 10^0$$

52,311

SPF($\beta=10$, $t=3$, $m=-5$, $M=5$)

$$y \times (z + x)$$

$$y \times (0,505 + 1,14) \times 10^1$$

$$(3,18 \times 10^0) \times (1,64 \times 10^1)$$

$$(3,18 \times 1,64) \times 10^{(0+1)}$$

$$5,21 \times 10^1$$

$$(y \times z) + (y \times x)$$

$$(3,18 \times 10^0) \times (5,05 \times 10^0) + (3,18 \times 10^0) \times (1,14 \times 10^1)$$

$$(3,18 \times 5,05) \times 10^{(0+0)} + (3,18 \times 1,14) \times 10^{(0+1)}$$

$$1,60 \times 10^1 + 3,62 \times 10^1$$

$$(1,60 + 3,62) \times 10^1$$

$$5,22 \times 10^1$$

Atenção!

Multiplicação e Divisão de conjuntos de números ➡ multiplicar números grandes com número pequenos
e dividir números com magnitude semelhantes

Comparando números em Ponto Flutuante

SPF($\beta=10$, $t=3$, $m=-5$, $M=5$)

$$x = 1,14 \times 10^1$$

$$y = 3,18 \times 10^0$$

$$z = 5,05 \times 10^0$$

$$a = y \times (z + x)$$

$$a = 5,21 \times 10^1$$

$$b = (y \times z) + (y \times x)$$

$$b = 5,22 \times 10^1$$

Comparando números em Ponto Flutuante

SPF($\beta=10$, $t=3$, $m=-5$, $M=5$)

$$x = 1,14 \times 10^1$$

$$y = 3,18 \times 10^0$$

$$z = 5,05 \times 10^0$$

$$a = y \times (z + x)$$

$$a = 5,21 \times 10^1$$

$$b = (y \times z) + (y \times x)$$

$$b = 5,22 \times 10^1$$

```
float a = y*(z+x);  
float b = (y*z)+(y*x);  
  
if (a == b)  
    return 1;  
return 0;
```

Comparando números em Ponto Flutuante

SPF($\beta=10$, $t=3$, $m=-5$, $M=5$)

$x = 1,14 \times 10^1$
 $y = 3,18 \times 10^0$
 $z = 5,05 \times 10^0$

$a = y \times (z + x)$
 $a = 5,21 \times 10^1$

$b = (y \times z) + (y \times x)$
 $b = 5,22 \times 10^1$

```
float a = y*(z+x);  
float b = (y*z)+(y*x);
```

```
if (a == b)  
    return 1;  
return 0;
```

```
float a = y*(z+x);  
float b = (y*z)+(y*x);  
  
if (abs(a-b) <= erro)  
    return 1;  
return 0;
```

Comparando números em Ponto Flutuante

- `= if abs(x-y) <= error`
- `≠ if abs(x-y) > error`
- `< if (x-y) < error`
- `≤ if (x-y) <= error`
- `> if (x-y) > error`
- `≥ if (x-y) >= error`

Comparando números em Ponto Flutuante

- `= if abs(x-y) <= error`
- `≠ if abs(x-y) > error`
- `< if (x-y) < error`
- `≤ if (x-y) <= error`
- `> if (x-y) > error`
- `≥ if (x-y) >= error`

Usando erro proporcional aos números

- `= if abs(x-y) <= abs(x+y)*erro`

Epsilon da Máquina

É a diferença entre 1,0 e o menor valor representável maior que 1,0.

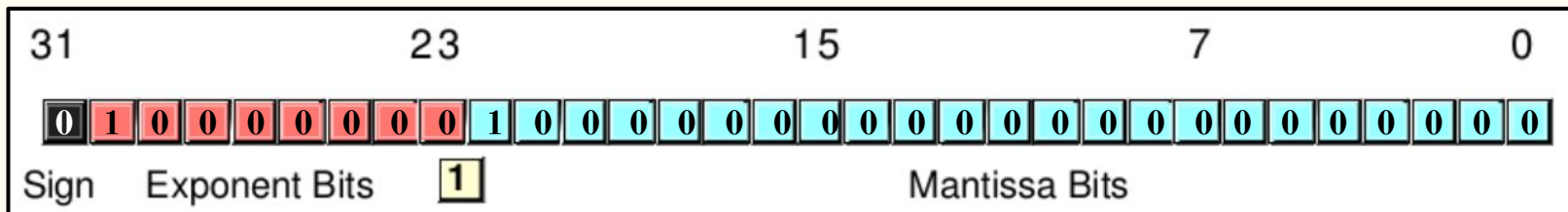
- float.h: FLT_EPSILON, DBL_EPSILON

```
int AlmostEqualRelative(float A, float B)
{
    // Calculate the difference.
    float diff = fabs(A - B);
    A = fabs(A);
    B = fabs(B);
    // Find the largest
    float largest = (B > A) ? B : A;

    if (diff <= largest * FLT_EPSILON)
        return 1;
    return 0;
}
```

ULP

(Units in the Last Place) - Quantos números podem ser representados entre dois números quaisquer.

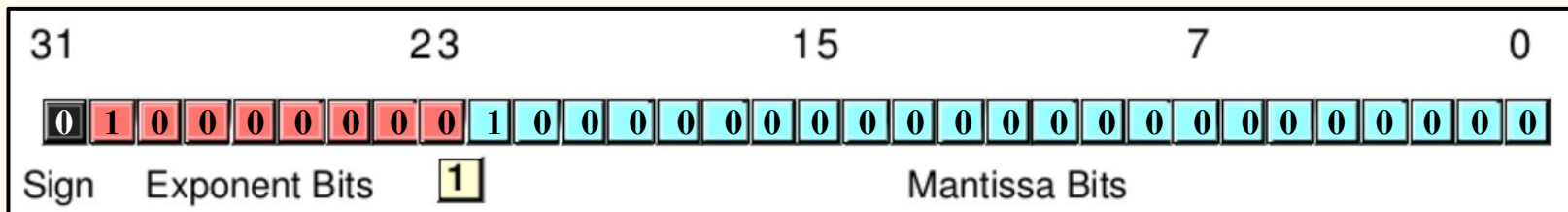


Ponto Flutuante: $1,1 \times 2^{128-127} = 3_{10}$
1077936128₁₀

Inteiro: $2^{30} + 2^{22} =$

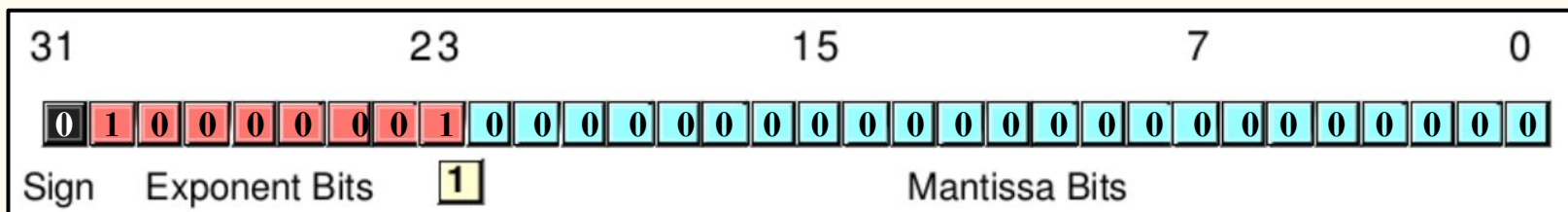
ULP

(Units in the Last Place) - Quantos números podem ser representados entre dois números quaisquer.



Ponto Flutuante: $1,1 \times 2^{128-127} = 3_{10}$
1077936128₁₀

Inteiro: $2^{30} + 2^{22} =$

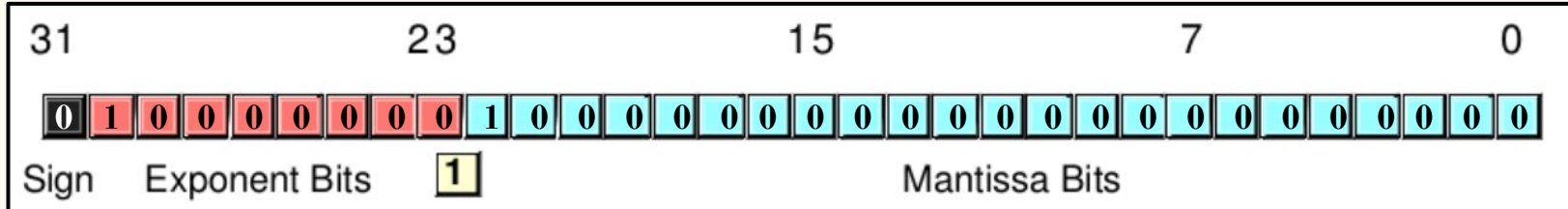


Ponto Flutuante: $1,0 \times 2^{129-127} = 4_{10}$
1082130432₁₀

Inteiro: $2^{30} + 2^{23} =$

ULP

(Units in the Last Place) - Quantos números podem ser representados entre dois números quaisquer.



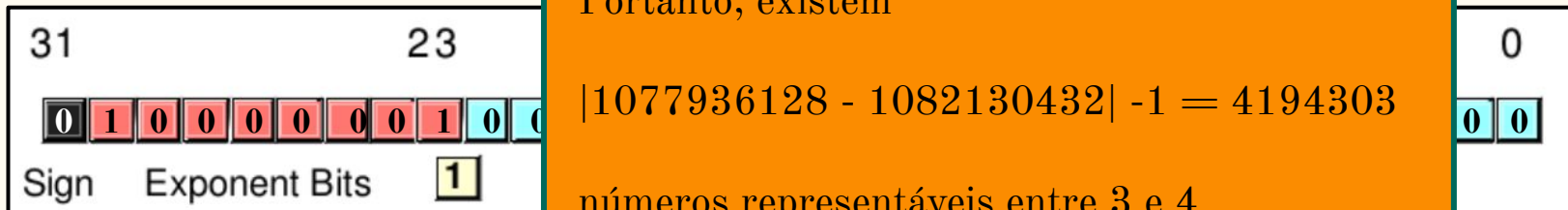
Ponto Flutuante: $1,1 \times 2^{128-127} = 3_{10}$
1077936128₁₀

Inteiro: $2^{30} + 2^{22} =$

Portanto, existem

$$|1077936128 - 1082130432| - 1 = 4194303$$

números representáveis entre 3 e 4

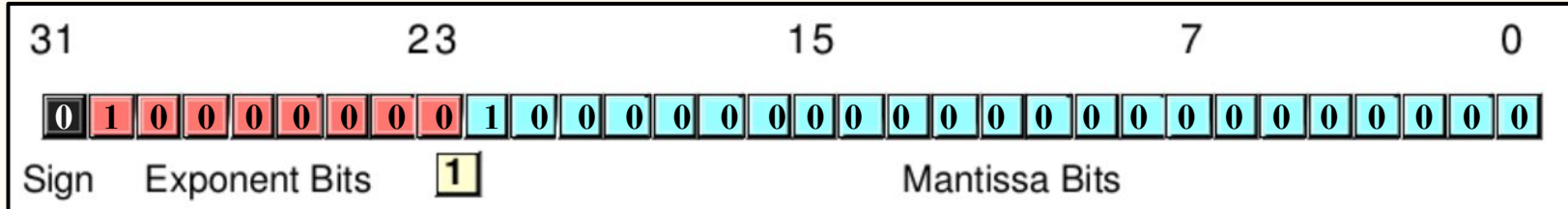


Ponto Flutuante: $1,0 \times 2^{129-127} = 4_{10}$
1082130432₁₀

Inteiro: $2^{30} + 2^{22} =$

ULP

(Units in the Last Place) - Quantos números podem ser representados entre dois números quaisquer.



Ponto Flutuante: $1,1 \times 2^{128-127} = 3_{10}$
1077936128₁₀

Inteiro: $2^{30} + 2^{22} =$

Portanto, existem

$|1077936128 - 1082130432| - 1 = 4194303$

números representáveis entre 3 e 4

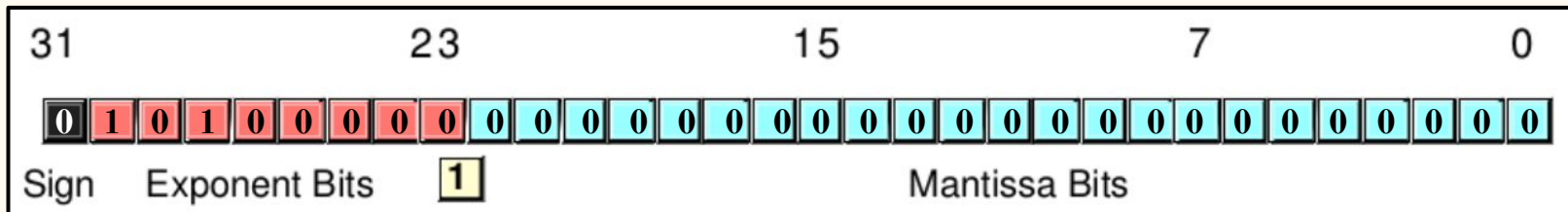
EA = 1
ER = 0,33333

Ponto Flutuante: $1,0 \times 2^{129-127} = 4_{10}$
1082130432₁₀

Inteiro: $2^{30} + 2^{22} =$

ULP

(Units in the Last Place) - Quantos números podem ser representados entre dois números quaisquer.

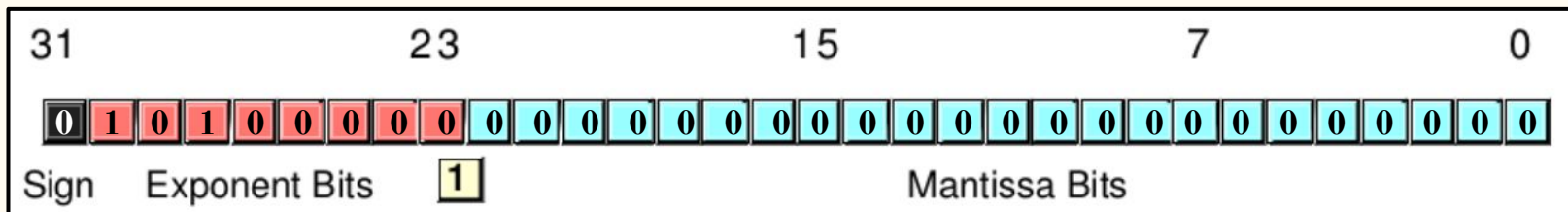


Ponto Flutuante: $1,0 \times 2^{160-127} \cong 8589934592$
1342177280

Inteiro: $2^{30} + 2^{28} =$

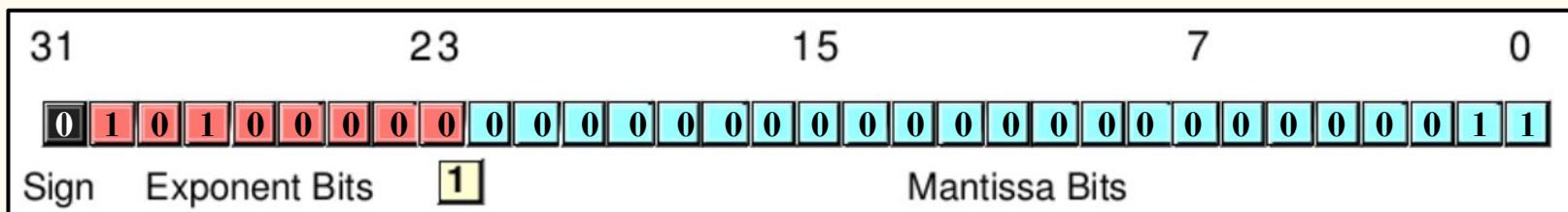
ULP

(Units in the Last Place) - Quantos números podem ser representados entre dois números quaisquer.



Ponto Flutuante: $1,0 \times 2^{160-127} \cong 8589934592$
1342177280

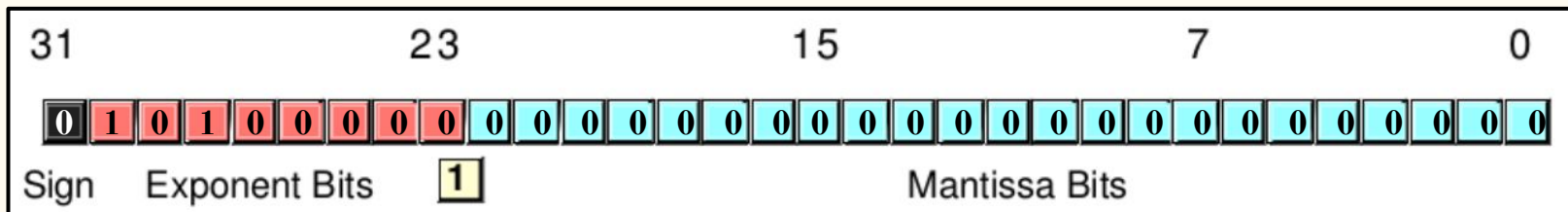
Inteiro: $2^{30} + 2^{28} =$



Ponto Flutuante: $1,0...11 \times 2^{160-127} \cong 8589937664$ Inteiro: $2^{30} + 2^{28} + 2^1 + 2^0 = 1342177283$

ULP

(Units in the Last Place) - Quantos números podem ser representados entre dois números quaisquer.



Ponto Flutuante: $1,0 \times 2^{160-127} \cong 8589934592$
1342177280

Inteiro: $2^{30} + 2^{28} =$

Portanto, existem

$$|1342177280 - 1342177283| - 1 = 2$$

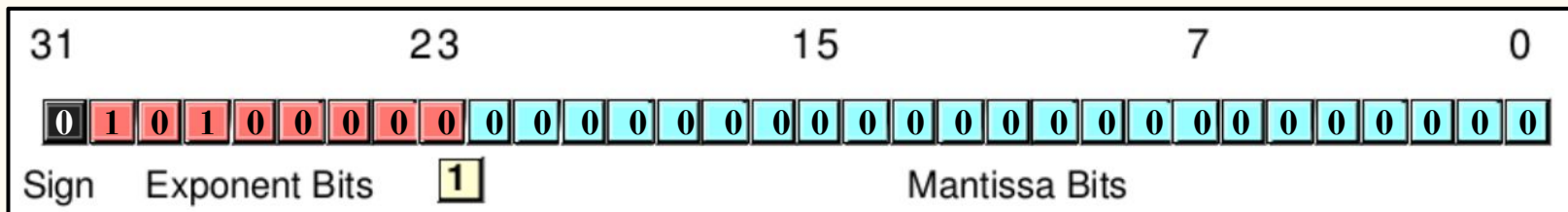
números representáveis entre

8589934592 e 8589937664

Ponto Flutuante: $1,0...11 \times 2^{160-127} \cong 8589937664$ Inteiro: $2^{30} + 2^{28} + 2^1 + 2^0 = 1342177283$

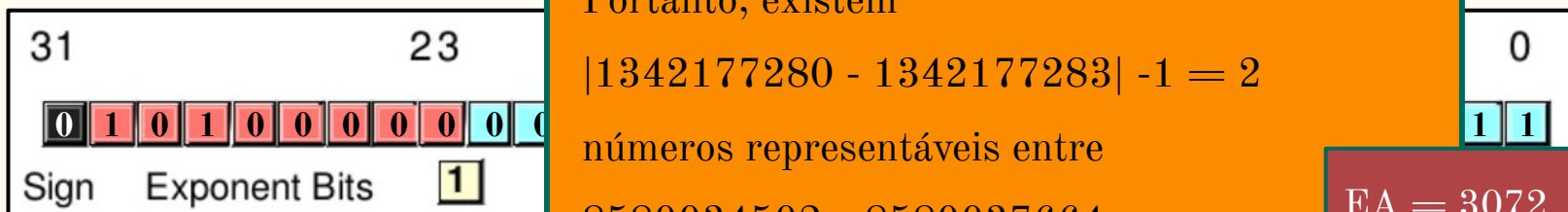
ULP

(Units in the Last Place) - Quantos números podem ser representados entre dois números quaisquer.



Ponto Flutuante: $1,0 \times 2^{160-127} \cong 8589934592$
1342177280

Inteiro: $2^{30} + 2^{28} =$



Portanto, existem

$$|1342177280 - 1342177283| - 1 = 2$$

números representáveis entre

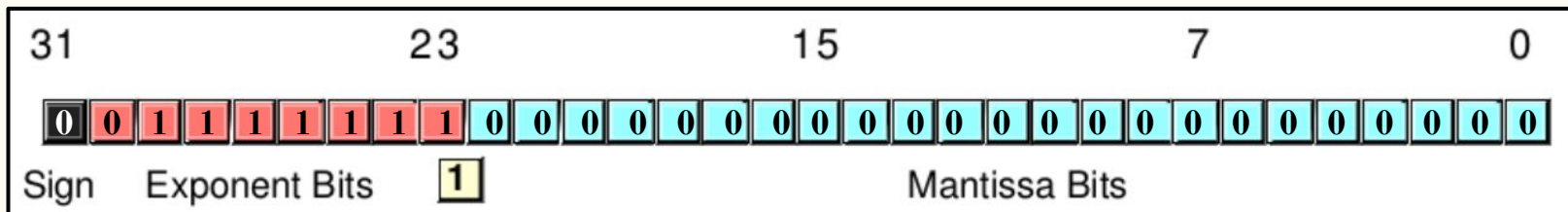
8589934592 e 8589937664

EA = 3072
ER < 0,000001

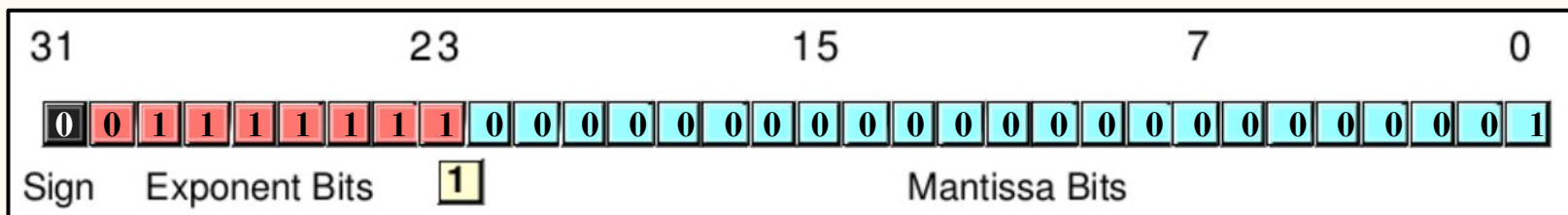
Ponto Flutuante: $1,0...11 \times 2^{160-127}$
1342177283

Epsilon da Máquina

É a diferença entre 1,0 e o menor valor representável maior que 1,0.



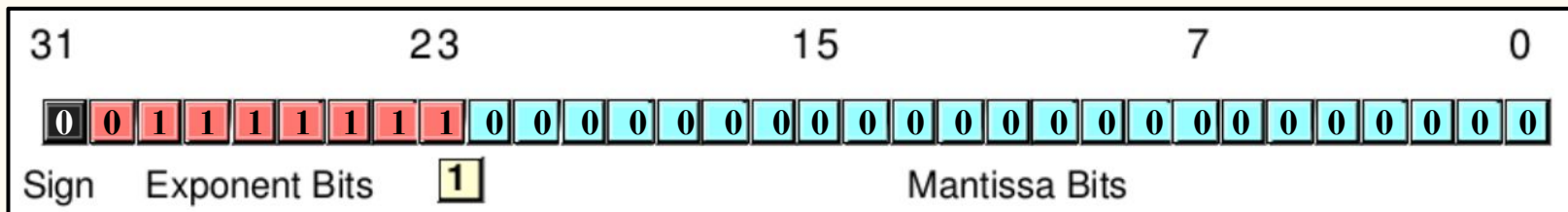
Ponto Flutuante: $1,0 \times 2^{127-127} = 1,0$



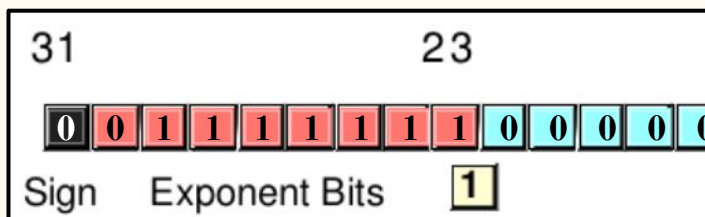
Ponto Flutuante: $1,0...01 \times 2^{127-127} = 1,00000011920928955078125$

Epsilon da Máquina

É a diferença entre 1,0 e o menor valor representável maior que 1,0.



Ponto Flutuante: $1,0 \times 2^{127-127} = 1,0$

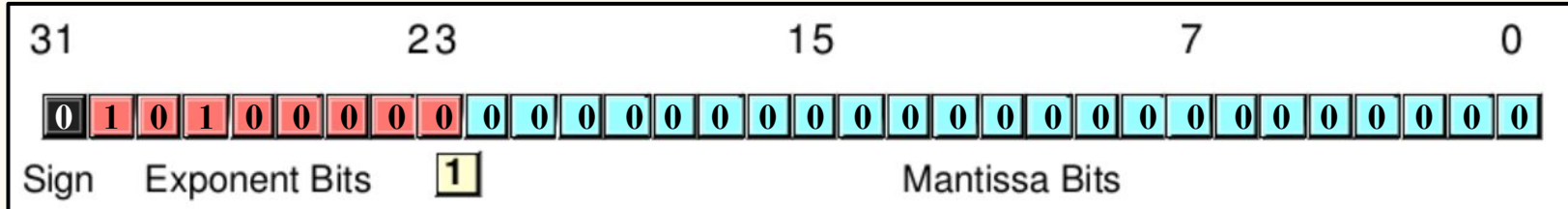


Portanto, $\varepsilon = 0,00000011920928955078125$

Ponto Flutuante: $1,0...01 \times 2^{127-127} = 1,00000011920928955078125$

ULP

(Units in the Last Place) - Quantos números podem ser representados entre dois números quaisquer.



Ponto Flutuante: $1,0 \times 2^{160-127} \cong 8589934592$
1342177280

Inteiro: $2^{30} + 2^{28} =$

$ER < 0,00000012$

Portanto, existem

$|1342177280 - 1342177283| - 1 = 2$

números representáveis entre

8589934592 e 8589937664



$EA = 3072$
 $ER < 0,000001$

Ponto Flutuante: $1,0...11 \times 2^{160-127} \cong 8589937664$
1342177283

Inteiro: $2^{30} + 2^{28} + 2^{27} =$

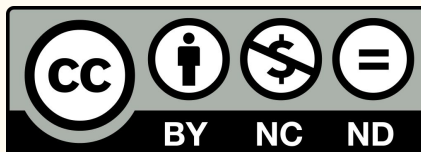
Referências

- Daniel Weingaertner; notas de aula da disciplina **Introdução à Computação Científica** (UFPR/DINF)
- A. Kaw, E. Kalu; **Numerical Methods with Applications**. Disponível em <https://nm.mathforcollege.com/textbook-numerical-methods-with-applications/>
- Bruce Dawson; **Comparing Floating Point Numbers, 2012 Edition**. Disponível em <https://randomascii.wordpress.com/2012/02/25/comparing-floating-point-numbers-2012-edition/>
- Sérgio Peters e Julio Felipe Szeremeta; **Cálculo Numérico Computacional**. Editora UFSC. Disponível em <http://sergiopeters.prof.ufsc.br/livro-calculo-numerico-computacional/>

Créditos

Este documento é de autoria do Prof. Guilherme Alex Derenievicz (UFPR/DINF), para uso na disciplina Introdução à Computação Científica (CI1164).

Compartilhe este documento de acordo com a licença abaixo



Este documento está licenciado com uma Licença Creative Commons **Atribuição-NãoComercial-SemDerivações** 4.0 Internacional.

<https://creativecommons.org/licenses/by-nc-sa/4.0/>