

# Arquitectura de Computadores I

Vírgula flutuante - standard IEEE754

Miguel Barão

**Vírgula fixa**

**Vírgula flutuante**

**Formato binário IEEE754**

**Precisões simples e dupla**

**Overflow e underflow**

**Problemas numéricos no cálculo em vírgula flutuante**

**Vírgula flutuante em RISC-V**

## Revisão: Conversão de números fraccionários para binário

$$14.375 = 14 + 0.375 =$$

Parte inteira converte-se com divisões por 2:

14	
<hr/>	
7	0
3	1
1	1
0	1

$$14_{\text{dec}} = 1110_{\text{bin}}$$

## Revisão: Conversão de números fraccionários para binário

$$14.375 = 14 + 0.375 = 1110.011\bar{0}_{\text{bin}}$$

Parte inteira converte-se com divisões por 2:

14	
7	0
3	1
1	1
0	1

$$14_{\text{dec}} = 1110_{\text{bin}}$$

Parte fraccionária convete-se com multiplicações por 2:

0.	375
0.	75
1.	5
1.	0
0.	0 (rep)

$$0.625_{\text{dec}} = 0.011\bar{0}_{\text{bin}}$$

# Vírgula fixa

## Vírgula fixa

**Vírgula fixa** representa números com número fixo de algarismos à esquerda e direita da vírgula

$a_2$	$a_1$	$a_0.$	$a_{-1}$	$a_{-2}$	$a_{-3}$	$a_{-4}$
-------	-------	--------	----------	----------	----------	----------

Simples, pode ser implementada recorrendo a números inteiros.

**Vírgula fixa** representa números com número fixo de algarismos à esquerda e direita da vírgula

$a_2$	$a_1$	$a_0.$	$a_{-1}$	$a_{-2}$	$a_{-3}$	$a_{-4}$
-------	-------	--------	----------	----------	----------	----------

Simple, pode ser implementada recorrendo a números inteiros.

Por exemplo, com 2 casas decimais

$$\begin{aligned}12.30 \times 6.02 &= (1230 \times 10^{-2}) \times (602 \times 10^{-2}) \\&= 740460 \times 10^{-4} \\&= 74.0460 \\&\approx 74.05\end{aligned}$$

Tipicamente, usa-se vírgula fixa em base decimal para representar valores monetários.

# Vírgula flutuante



Vírgula flutuante consiste na representação de números no formato

$$\text{significando} \times \text{base}^{\text{expoente}}$$

onde

- base é normalmente 10 ou 2.
- expoente é um número inteiro.
- significando é um número  $1 \leq x < \text{base}$ .

Vírgula flutuante consiste na representação de números no formato

$$\text{significando} \times \text{base}^{\text{expoente}}$$

onde

- base é normalmente 10 ou 2.
- expoente é um número inteiro.
- significando é um número  $1 \leq x < \text{base}$ .

Exemplos de vírgula flutuante em base decimal:

- $3.14159265359 \longrightarrow 3.14159265359 \times 10^0$
- $-0.0012345 \longrightarrow -1.2345 \times 10^{-3}$
- $1024 \longrightarrow 1.024 \times 10^3$

■  $0.375 =$

■  $0.375 = 0.011_{\text{bin}} =$

- $0.375 = 0.011_{\text{bin}} = 1.1_{\text{bin}} \times 2^{-2}$

- $128 =$

## Vírgula flutuante em binário

- $0.375 = 0.011_{\text{bin}} = 1.1_{\text{bin}} \times 2^{-2}$
- $128 = 10000000_{\text{bin}} = 1.0_{\text{bin}} \times 2^7$
- $3.25 =$

## Vírgula flutuante em binário

- $0.375 = 0.011_{\text{bin}} = 1.1_{\text{bin}} \times 2^{-2}$
- $128 = 10000000_{\text{bin}} = 1.0_{\text{bin}} \times 2^7$
- $3.25 = 11.01_{\text{bin}} = 1.101_{\text{bin}} \times 2^1$
- $0.1 =$

## Vírgula flutuante em binário

- $0.375 = 0.011_{\text{bin}} = 1.1_{\text{bin}} \times 2^{-2}$
- $128 = 10000000_{\text{bin}} = 1.0_{\text{bin}} \times 2^7$
- $3.25 = 11.01_{\text{bin}} = 1.101_{\text{bin}} \times 2^1$
- $0.1 = \underbrace{0.0001100110011 \dots}_{\text{dízima infinita}}_{\text{bin}} =$



## Vírgula flutuante em binário

- $0.375 = 0.011_{\text{bin}} = 1.1_{\text{bin}} \times 2^{-2}$
- $128 = 10000000_{\text{bin}} = 1.0_{\text{bin}} \times 2^7$
- $3.25 = 11.01_{\text{bin}} = 1.101_{\text{bin}} \times 2^1$
- $0.1 = \underbrace{0.0001100110011 \dots}_{\text{dízima infinita}}_{\text{bin}} = 1.100110011 \dots_{\text{bin}} \times 2^{-4}$

## Vírgula flutuante em binário

$$\blacksquare 0.375 = 0.011_{\text{bin}} = 1.1_{\text{bin}} \times 2^{-2}$$

$$\blacksquare 128 = 10000000_{\text{bin}} = 1.0_{\text{bin}} \times 2^7$$

$$\blacksquare 3.25 = 11.01_{\text{bin}} = 1.101_{\text{bin}} \times 2^1$$

$$\blacksquare 0.1 = \underbrace{0.0001100110011 \dots}_{\text{dízima infinita}}_{\text{bin}} = 1.100110011 \dots_{\text{bin}} \times 2^{-4}$$

$$\blacksquare 0.2 = 0.001100110011 \dots_{\text{bin}} = 1.1001100110011 \dots_{\text{bin}} \times 2^{-3}$$

## Vírgula flutuante em binário

- $0.375 = 0.011_{\text{bin}} = 1.1_{\text{bin}} \times 2^{-2}$
- $128 = 10000000_{\text{bin}} = 1.0_{\text{bin}} \times 2^7$
- $3.25 = 11.01_{\text{bin}} = 1.101_{\text{bin}} \times 2^1$
- $0.1 = \underbrace{0.0001100110011 \dots}_{\text{dízima infinita}}_{\text{bin}} = 1.100110011 \dots_{\text{bin}} \times 2^{-4}$
- $0.2 = 0.001100110011 \dots_{\text{bin}} = 1.1001100110011 \dots_{\text{bin}} \times 2^{-3}$
- $0.3 = 0.010011001100 \dots_{\text{bin}} = 1.0011001100110 \dots_{\text{bin}} \times 2^{-2}$

### A maioria dos números não têm representação exacta!

- Há muitos números com dízima finita em decimal, mas com dízima infinita em binário!

## Vírgula flutuante em binário

- $0.375 = 0.011_{\text{bin}} = 1.1_{\text{bin}} \times 2^{-2}$
- $128 = 10000000_{\text{bin}} = 1.0_{\text{bin}} \times 2^7$
- $3.25 = 11.01_{\text{bin}} = 1.101_{\text{bin}} \times 2^1$
- $0.1 = \underbrace{0.0001100110011 \dots}_{\text{dízima infinita}}_{\text{bin}} = 1.100110011 \dots_{\text{bin}} \times 2^{-4}$
- $0.2 = 0.001100110011 \dots_{\text{bin}} = 1.1001100110011 \dots_{\text{bin}} \times 2^{-3}$
- $0.3 = 0.010011001100 \dots_{\text{bin}} = 1.0011001100110 \dots_{\text{bin}} \times 2^{-2}$

### A maioria dos números não têm representação exacta!

- Há muitos números com dízima finita em decimal, mas com dízima infinita em binário!
- Apenas os números que são somas de potências de 2 têm dízima finita em binário.

# Formato binário IEEE754

O standard IEEE754 especifica como os números em vírgula flutuante são codificados:

`float` precisão simples: 32 bits

`double` precisão dupla: 64 bits

O standard IEEE754 especifica como os números em vírgula flutuante são codificados:

**float** precisão simples: 32 bits

**double** precisão dupla: 64 bits

Adicionalmente são definidos os símbolos:

**$\pm\text{Inf}$**  representa infinito, e.g. devido a overflow.

**NaN** not-a-number representa resultados indefinidos.

**$\pm 0.0$**  zeros são números especiais.

## Precisões simples e dupla



## Precisão simples: float (32 bits)

$\pm$ (1 bit)	expoente (8 bits)	significando (23 bits + 1)
---------------	-------------------	----------------------------

Expoente:

- entre  $00000001_{\text{bin}}$  e  $11111110_{\text{bin}}$ , i.e. de 1 a 254 decimal.
- expoente é **biased**:  $2^0$  é codificado com  $127 = 01111111_{\text{bin}}$ .

Significando:

- Na realidade são 24 bits (um dos bits está escondido).
- O primeiro bit é sempre 1 e não é representado (é o bit escondido).

## Precisão simples: float (32 bits)

$\pm$ (1 bit)	expoente (8 bits)	significando (23 bits + 1)
---------------	-------------------	----------------------------

Expoente:

- entre  $00000001_{\text{bin}}$  e  $11111110_{\text{bin}}$ , i.e. de 1 a 254 decimal.
- expoente é **biased**:  $2^0$  é codificado com  $127 = 01111111_{\text{bin}}$ .

Significando:

- Na realidade são 24 bits (um dos bits está escondido).
- O primeiro bit é sempre 1 e não é representado (é o bit escondido).

---

$$0 \ 01111111 \ 000000000000000000000000 = +1.0 \times 2^0$$

## Precisão simples: float (32 bits)

$\pm$ (1 bit)	expoente (8 bits)	significando (23 bits + 1)
---------------	-------------------	----------------------------

Expoente:

- entre  $00000001_{\text{bin}}$  e  $11111110_{\text{bin}}$ , i.e. de 1 a 254 decimal.
- expoente é **biased**:  $2^0$  é codificado com  $127 = 01111111_{\text{bin}}$ .

Significando:

- Na realidade são 24 bits (um dos bits está escondido).
- O primeiro bit é sempre 1 e não é representado (é o bit escondido).

---

0	01111111	000000000000000000000000	$= +1.0 \times 2^0$
1	01111111	010000000000000000000000	$= -1.01 \times 2^0$

## Precisão simples: float (32 bits)

$\pm$ (1 bit)	expoente (8 bits)	significando (23 bits + 1)
---------------	-------------------	----------------------------

Expoente:

- entre  $00000001_{\text{bin}}$  e  $11111110_{\text{bin}}$ , i.e. de 1 a 254 decimal.
- expoente é **biased**:  $2^0$  é codificado com  $127 = 01111111_{\text{bin}}$ .

Significando:

- Na realidade são 24 bits (um dos bits está escondido).
- O primeiro bit é sempre 1 e não é representado (é o bit escondido).

---

0 01111111 000000000000000000000000	$= +1.0 \times 2^0$
1 01111111 010000000000000000000000	$= -1.01 \times 2^0$
0 01111101 00110011001100110011001	$= 1.00110011001100110 \dots \times 2^{-2} \approx 0.3$

## Precisão simples: float (32 bits)

$\pm$ (1 bit)	expoente (8 bits)	significando (23 bits + 1)
---------------	-------------------	----------------------------

Expoente:

- entre  $00000001_{\text{bin}}$  e  $11111110_{\text{bin}}$ , i.e. de 1 a 254 decimal.
- expoente é **biased**:  $2^0$  é codificado com  $127 = 01111111_{\text{bin}}$ .

Significando:

- Na realidade são 24 bits (um dos bits está escondido).
- O primeiro bit é sempre 1 e não é representado (é o bit escondido).

---

0	01111111	000000000000000000000000	$= +1.0 \times 2^0$
1	01111111	010000000000000000000000	$= -1.01 \times 2^0$
0	01111101	00110011001100110011001	$= 1.00110011001100110 \dots \times 2^{-2} \approx 0.3$
1	10000000	110000000000000000000000	$= -1.11 \times 2^1 = -3.5$

## Precisão simples: float (32 bits)

$\pm$ (1 bit)	expoente (8 bits)	significando (23 bits + 1)
---------------	-------------------	----------------------------

Expoente:

- entre  $00000001_{\text{bin}}$  e  $11111110_{\text{bin}}$ , i.e. de 1 a 254 decimal.
- expoente é **biased**:  $2^0$  é codificado com  $127 = 01111111_{\text{bin}}$ .

Significando:

- Na realidade são 24 bits (um dos bits está escondido).
- O primeiro bit é sempre 1 e não é representado (é o bit escondido).

---

0	01111111	000000000000000000000000	$= +1.0 \times 2^0$
1	01111111	010000000000000000000000	$= -1.01 \times 2^0$
0	01111101	00110011001100110011001	$= 1.00110011001100110 \dots \times 2^{-2} \approx 0.3$
1	10000000	110000000000000000000000	$= -1.11 \times 2^1 = -3.5$
0	10000000	10010010000111111011011	$\approx 3.14159265358979$

---

## Infinito: Inf

O infinito é gerado em várias situações:

- Quando uma operação resulta em overflow.
- Divisão por zero:  $1.0/0.0 = \text{Inf}$  ou  $-1.0/0.0 = -\text{Inf}$ .
- Funções podem definir esse resultado:  $\log(0.0) = -\text{Inf}$ .

Formato:

0	11111111	00000000000000000000000000000000	+Inf
1	11111111	00000000000000000000000000000000	-Inf

Algumas propriedades:

- $x + \text{Inf} = \text{Inf}$ , onde  $x$  é um número finito.
- $x/\text{Inf} = 0.0$ , se  $x \neq \pm\text{Inf}, \text{NaN}$ .
- $\exp(-\text{Inf}) = 0.0$
- $1.0/ - 0.0 = -\text{Inf}$

## Not-a-Number: NaN

O Not-a-Number é gerado em várias situações:

- $0.0/0.0 = \text{NaN}$
- $0.0 * \text{Inf} = \text{NaN}$
- $\text{Inf} - \text{Inf} = \text{NaN}$ .

Formato:

0	11111111	10000000000000000000000000000000	+NaN
1	11111111	10000000000000000000000000000000	-NaN

Algumas propriedades:

- Os Not-a-Number tem tendência a propagar-se e contaminar os resultados seguintes.
- $x + \text{NaN} = \text{NaN}$ .
- A comparação  $\text{NaN} = \text{NaN}$  é *FALSA*.
- A comparação  $\text{NaN} \neq \text{NaN}$  é *VERDADEIRA*.
- $1.0^{\text{NaN}} == 1.0$ .



## Zero: 0.0

Existem **dois** zeros:

- $+0.0$
- $-0.0$

Formato:

0	00000000	000000000000000000000000	$+0.0$
1	00000000	000000000000000000000000	$-0.0$

Algumas propriedades:

- O zero não é um número em vírgula flutuante regular.
- Não se considera o bit escondido no significando.
- A comparação  $0.0 = -0.0$  é *VERDADEIRA*.
- $1.0/0.0 = +\text{Inf}$  enquanto que  $1.0/(-0.0) = -\text{Inf}$ .
- $0.0^{0.0} = 1.0$ .

## Precisão dupla: double (64 bits)

$\pm$ (1 bit)	expoente (11 bits)	significando (52 bits + 1)
---------------	--------------------	----------------------------

Expoente:

- entre  $00000000001_{\text{bin}}$  e  $11111111110_{\text{bin}}$ , i.e. de 1 a 2046.
- expoente é *biased*:  $2^0$  é codificado com  $1023 = 01111111111_{\text{bin}}$ .

Significando:

- Na realidade são 53 bits (há 1 bit escondido).
- O primeiro bit é sempre 1 e não é representado (é o bit escondido).

## Precisão dupla: double (64 bits)

$\pm$ (1 bit)	expoente (11 bits)	significando (52 bits + 1)
---------------	--------------------	----------------------------

Expoente:

- entre  $00000000001_{\text{bin}}$  e  $11111111110_{\text{bin}}$ , i.e. de 1 a 2046.
- expoente é *biased*:  $2^0$  é codificado com  $1023 = 01111111111_{\text{bin}}$ .

Significando:

- Na realidade são 53 bits (há 1 bit escondido).
- O primeiro bit é sempre 1 e não é representado (é o bit escondido).

[illegible]

Que números são estes?

# Overflow e underflow

**Overflow** ocorre quando o expoente atinge 128 em precisão simples, ou 1024 em precisão dupla. Nesta situação deixa de ser possível representar o número em vírgula flutuante, passando o resultado a ser representado por  $\pm\text{Inf}$ .

**Underflow** ocorre quando o expoente atinge -127 em precisão simples, ou -1023 em precisão dupla, e o significando é nulo. Nesta situação deixa de ser possível representar o número em vírgula flutuante, passando o resultado a ser representado por  $\pm 0.0$ .

**Underflow progressivo** ocorre quando o expoente atinge -127 ou -1023, em precisão simples ou dupla, e o significando ainda não é nulo. Nesta situação, o número deixa de ser normal, passando a ser um número **não normalizado** ou **subnormal**. Estes números perdem progressivamente precisão à medida que se aproximam de um underflow.

# Problemas numéricos no cálculo em vírgula flutuante

- 1 Calcule  $(0.1 + 0.2) + 0.3$
- 2 Calcule  $0.1 + (0.2 + 0.3)$

- 1 Calcule  $(0.1 + 0.2) + 0.3$
- 2 Calcule  $0.1 + (0.2 + 0.3)$

Resultado em precisão dupla:

[illegible]

A adição é associativa?



- 1 Calcule  $(0.1 + 0.2) + 0.3$
- 2 Calcule  $0.1 + (0.2 + 0.3)$

Resultado em precisão dupla:

[illegible]

## A adição é associativa?

- Em matemática **Sim**
- Em vírgula flutuante **Não...**

Caso notável:

$$x^2 - y^2 = (x + y)(x - y)$$

Será válido em vírgula flutuante?

Suponhamos que  $x$  e  $y$  são números próximos. Por exemplo,

$$x = 0.9999, \quad y = 0.9998$$

obtemos

$$x^2 - y^2 \approx 0.00019997000000004928$$

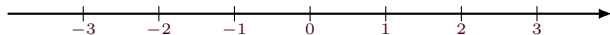
$$(x - y)(x + y) \approx 0.000199969999999978$$

valor correcto: 0.00019997

- Operações em decimal e binário dão resultados diferentes.
- Conversão para binário introduz erros.
- Propriedades das operações (e.g. associatividade) não são satisfeitas.
- Erros acumulam-se à medida que se realizam mais operações.
- Não usar vírgula flutuante para valores monetários (se necessário, implementar cálculos decimais em software em vez de usar a vírgula flutuante do cpu.)

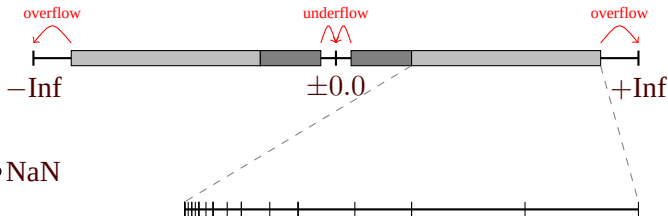
# Distância entre números

## ■ Inteiros:



A distância entre os números inteiros é constante em toda a escala.

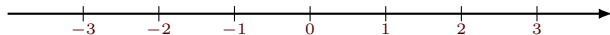
## ■ Vírgula flutuante:



● NaN

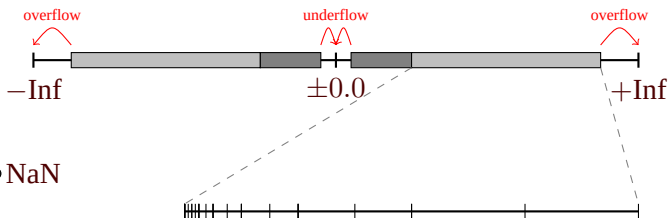
# Distância entre números

## ■ Inteiros:



A distância entre os números inteiros é constante em toda a escala.

## ■ Vírgula flutuante:



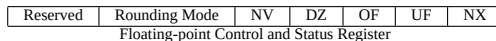
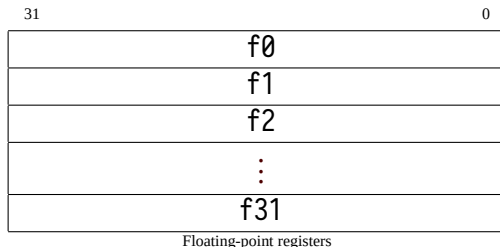
### ● NaN

A resolução de um número em vírgula flutuante varia com o expoente.

# Vírgula flutuante em RISC-V

# Vírgula flutuante na arquitectura RISC-V

- O suporte de vírgula flutuante é opcional na arquitectura RISC-V.
- A extensão F adiciona registos específicos para vírgula flutuante:
  - 32 registos de 32 bits f0–f31.
  - 1 registo fcsr (*Floating-point Control and Status Register*).



## Instruções de vírgula flutuante RISC-V

São adicionadas várias instruções novas:

<code>flw f0, 0(t0)</code>	<i># float load word</i>
<code>fsw f0, 4(t0)</code>	<i># float store word</i>
<code>fadd.s f2, f0, f1</code>	<i># soma em vírgula flutuante</i>
<code>fsub.s f2, f0, f1</code>	<i># subtração</i>
<code>fmul.s f2, f0, f1</code>	<i># multiplicação</i>
<code>fdiv.s f2, f0, f1</code>	<i># divisão</i>
<code>fmin.s f2, f0, f1</code>	<i># menor de dois números</i>
<code>fmax.s f2, f0, f1</code>	<i># maior de dois números</i>
<code>fsqrt.s f2, f0</code>	<i># raiz quadrada</i>
<code>fmadd.s f3, f0, f1, f2</code>	<i># multiply-add <math>f3 = f0*f1+f2</math></i>
<code>fmsub.s f3, f0, f1, f2</code>	<i># multiply-sub <math>f3 = f0*f1-f2</math></i>

e mais algumas...