

Floating Point

Tarcísio Fischer

tarcisio.fischer.cco@gmail.com

github.com/tarcisiofischer

Tópicos

- Introdução
- Resultados curiosos
- Representação
- Comparação aproximada
- Exemplos

Introdução

Problema: Representação de números Reais

Exemplo: 123.456

Introdução

Problema: Representação de números Reais

Exemplo: 123.456

Solução - Escolher um número fixo de valores antes e depois da vírgula
valor = 123456

(Ou seja, $123.456 = 123456 / 1000$)

Introdução

Problema: Representação de números Reais

Exemplo: 123.456

Outra solução - Utilizar notação científica ($m \cdot B^e$), com B fixo
(Exemplo para B=10):

$m = 1.23456$

$e = 2$

(Ou seja, $123.456 = 1.23456 * 10^2$)

floats & doubles utilizam essa notação, com B=2

Introdução

Conversão decimal para binário

$$1.625 = 1 + 0.625$$

$$0.625 \quad *2= 1.25 \quad \begin{matrix} . \\ = 1 + 0.25 \end{matrix}$$

$$0.25 \quad *2= 0.5 \quad = 0 + 0.5$$

$$0.5 \quad *2= 1.0 \quad = 1 + 0.0$$

Resultado

$$1.625 = 1.101$$

Introdução

Conversão binário para decimal

$$\begin{aligned}1.101b &= 1 * 2^{(0)} + 1 * 2^{(-1)} + 0 * 2^{(-2)} + 1 * 2^{(-3)} \\&= 1.0 + 0.5 + 0.0 + 0.125 \\&= \mathbf{1.625}\end{aligned}$$

Introdução

```
#include <stdio.h>
int main() {
    float x = 123.456;
    printf("%3.3f\n", x); // 123.456
}
```


Resultados Curiosos

Resultados Curiosos

```
tarcisio@Fermat: /partition1/workspace
tarcisio@Fermat:/partition1/workspaces$ cat example1.cc
#include <stdio.h>
int main() {
    printf("%1.40f\n", 0.1 + 0.2);
}
tarcisio@Fermat:/partition1/workspaces$ g++ example1.cc && ./a.out
0.300000000000000000444089209850062616169453
tarcisio@Fermat:/partition1/workspaces$
tarcisio@Fermat:/partition1/workspaces$
tarcisio@Fermat:/partition1/workspaces$ python -c "print(0.1 + 0.2)"
0.3000000000000000004
tarcisio@Fermat:/partition1/workspaces$
```

Resultados Curiosos



Resultados Curiosos

```
int main() {  
    float sum = 0.0;  
    for (int i = 0; i < 100; ++i)  
        sum += 0.01;  
    printf("%.60f\n", sum);  
}
```


Resultados Curiosos

```
int main() {  
    float a = 10e+5;  
    float b = -10e+5;  
    float c = 0.01;  
    float result1 = (a + b) + c;  
    float result2 = a + (b + c);  
    printf("result1 = %.60f\n", result1);  
    printf("result2 = %.60f\n", result2);  
}
```


Resultados Curiosos

```
int main() {  
    float a = 1.5f;  
    for (int i = 0; i < 10000000; ++i) {  
        a = a * 1.5f;  
        a = a / 2.5f;  
    }  
}
```

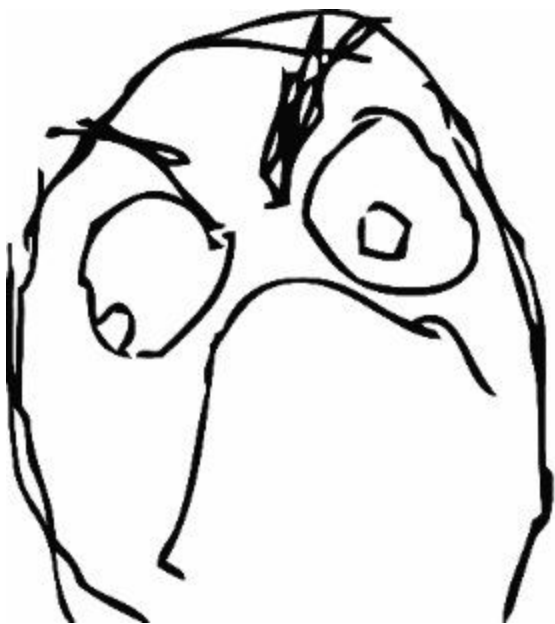
Performance (média de 10 rodadas): **1.036s**

Resultados Curiosos

```
int main() {  
    float a = 1.5f;  
    for (int i = 0; i < 10000000; ++i) {  
        a = a * 1.5f;  
        a = a / 2.5f;  
        a = a + 0.1f;  
        a = a - 0.1f;  
    }  
}
```



Performance (média de 10 rodadas): **0.151s**



Representação

Representação (float & double)

Notação científica

$$m \cdot B^e$$

m = mantissa

e = expoente ou ordem de grandeza

B = base

Exemplos (B = 10)

$$0.1 = 1.0 \cdot 10^{-1}$$

$$0.052 = 5.2 \cdot 10^{-2}$$

$$891.7 = 8.917 \cdot 10^2$$

Em binário (B = 2):

$$[\text{bin}] \cdot 2^e$$

Representação (float & double)

“Squeezing infinitely many real numbers into a finite number of bits requires an approximate representation.”

Referência: What Every Computer Scientist Should Know About Floating-Point Arithmetic

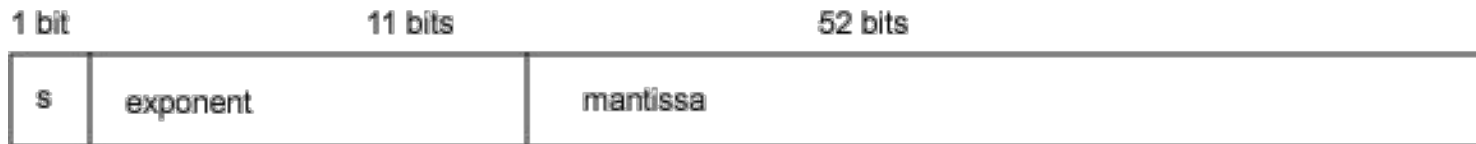
Representação (float & double)

IEEE 754 standard

IEEE Floating Point Representation



IEEE Double Precision Floating Point Representation



Referência: <http://www.ibm.com/developerworks/library/j-jtp0114/float.gif>

Representação (float & double)

“Floating-point representations are not necessarily unique.
For example, both 0.01×10^1 and 1.00×10^{-1} represent 0.1.”

Referência: What Every Computer Scientist Should Know About Floating-Point Arithmetic

Representação (float & double)

Normal form: **1.[m]** * B**e

Representação (float & double)

float $-1^{\textcolor{red}{s}} * 2^{(\textcolor{blue}{e} - 127)} * (\textcolor{brown}{1.0} + \textcolor{brown}{m}), \quad 0.0 \leq m < 1.0$

double $-1^{\textcolor{red}{s}} * 2^{(\textcolor{blue}{e} - 1023)} * (\textcolor{brown}{1.0} + \textcolor{brown}{m}), \quad 0.0 \leq m < 1.0$

IEEE Floating Point Representation



IEEE Double Precision Floating Point Representation



Representação (float & double)

Tem exemplo?



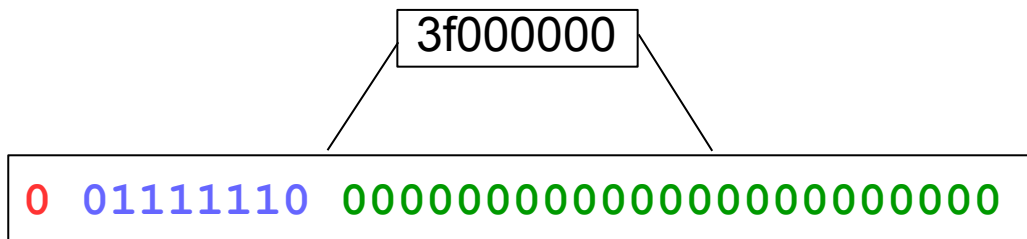
Representação (float & double)

Exemplo: Representação de **0.5** em float & double

```
int main() {  
    float f = 0.5f;  
    printf("%.40f\n", f); —————→ 0.50000000000000000000000000000000...  
    int f_bits = *(int*)&f;  
    printf("%8x\n", f_bits); —————→ 3f000000  
  
    double d = 0.5;  
    printf("%.60f\n", d); —————→ 0.500000000000000000000000000000000000...  
    long d_bits = *(long*)&d;  
    printf("%16lx\n", d_bits); —————→ 3fe0000000000000  
}
```

Representação (float & double)

Exemplo: Representação de **0.5** em **float**

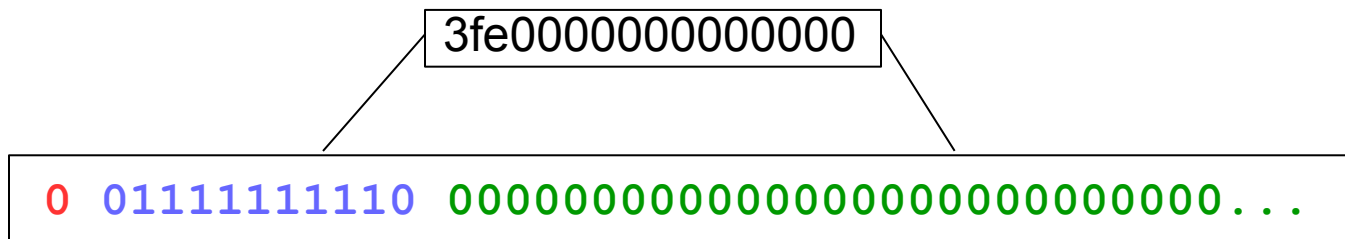


$$\begin{aligned} s &= 0 & = & \Rightarrow s = + \\ e &= 01111110 & = 126 & \Rightarrow e - 127 = -1 \\ m &= 000\dots000 & = 0 & \Rightarrow m = 0 \end{aligned}$$

$$(-1)^{**0} 1.0 * 2^{**(-1)} = 0.5$$

Representação (float & double)

Exemplo: Representação de **0.5** em **double**



$$\begin{array}{lll} s = 0 & = & \Rightarrow s = + \\ e = 011111110 & = 1022 & \Rightarrow e - 1023 = -1 \\ m = 000\dots000 & = 0 & \Rightarrow m = 0 \end{array}$$

$$(-1)^{**0} 1.0 * 2^{**(-1)} = 0.5$$

Representação (float & double)

Exemplo: Representação de **0.01** em **float**

```
#include <stdio.h>
int main()
{
    float f = 0.01f; —————→ 0.00999999977648258209228515625...
    printf("%.40f\n", f);
    int f_bits = *(int*)&f; —————→ 3c23d70a
    printf("%8x\n", f_bits);
}
```

Representação (float & double)

Exemplo: Representação de **0.01** em **float**

(Os exemplos 1 e 2 tem relação com essa característica!)

3c23d70a

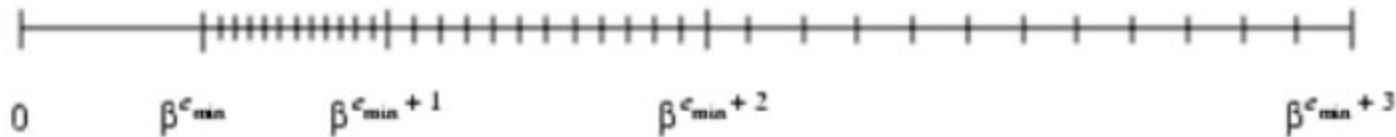
0 01111000 01000111101011100001010

$s = 0 \quad \Rightarrow \quad s = +$
 $e = 01111110 \quad = 120 \quad \Rightarrow \quad e - 127 = -7$
 $m = 010...010 \quad = \dots \quad \Rightarrow \quad m = 1.2799999\dots$

$(-1)^{**0} 1.m * 2^{**(-7)} =$
0.00999999977648258209228515625

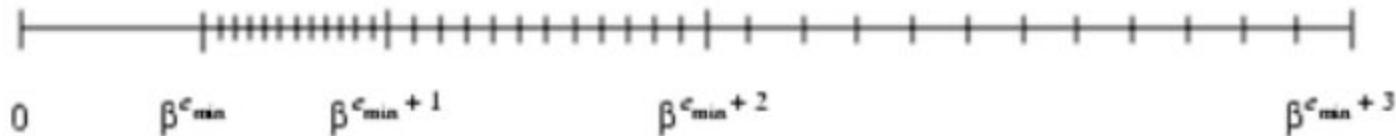


Representação (float & double)



Nota: Quanto maior o expoente, maior o “passo” entre dois valores representáveis utilizando **float**

Representação (float & double)



Exemplo:

(Números pequenos possuem “passo” pequeno)

0.00003051759267691522836685180664062500

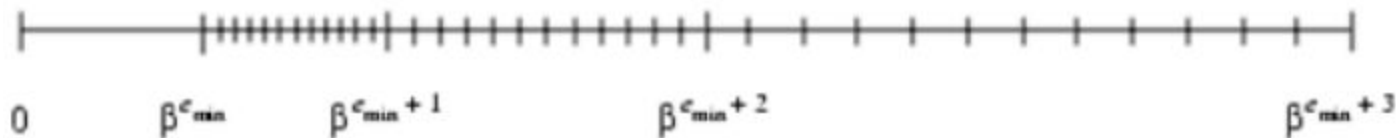
0.00003051759631489403545856475830078125 >> **Passo de $\sim 10^{-12}$**

(Números grandes possuem “passo” grande)

32768.01562500

32768.01953125 >> **Passo de $\sim 10^{-3}$**

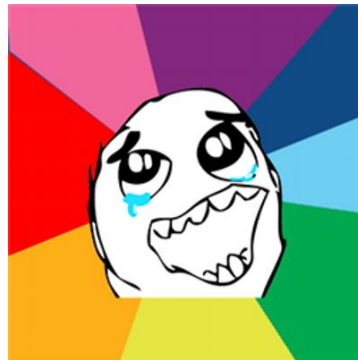
Representação (float & double)



Isso significa que operações envolvendo números com grandezas diferentes podem resultar em valores diferentes, ou seja,

não é verdade que
 $(a + b) + c == a + (b + c)$
para todo a, b, c

O exemplo 3 tem relação com essa característica



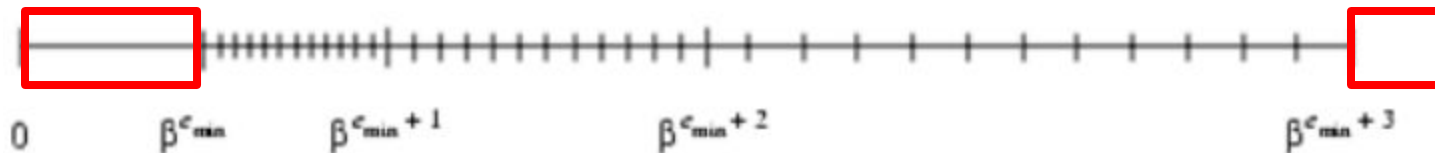
Overflow & Underflow

Overflow

Uma operação pode produzir um valor tão alto que não pode ser representado em um float/double.

Underflow

Uma operação pode produzir um valor tão próximo de 0.0 que não pode ser representado por um float/double.



Representação - Casos especiais

+0	0 00000000 000000000000000000000000
-0	1 00000000 000000000000000000000000
+inf	0 11111111 000000000000000000000000
-inf	1 11111111 000000000000000000000000
NaN	<any signal> 11111111 <At least one 1 in m>

Representação - Casos especiais

“Just as **NaNs** provide a way to continue a computation when expressions like **0/0** or **sqrt(-1)** are encountered, **infinities** provide a way to continue when an **overflow** occurs.”

Referência: What Every Computer Scientist Should Know About Floating-Point Arithmetic

```
int main() {  
    printf("%.40f\n", sqrt(-1)); —————> NaN  
  
    printf("%.40f\n", 0.0f/0.0f); —————> NaN  
    printf("%.40f\n", 1.0f/0.0f); —————> Inf  
}
```

Representação - Números Denormais/Subnormais

```
int main() {  
    // 1.175494e-38  
    printf("%e\n", std::numeric_limits<float>::min());  
    // 1.175494e-38  
    printf("%e\n", FLT_MIN);  
  
    float really_small = std::numeric_limits<float>::min() / 2f;  
    // 5.877472e-39 -> Como é possível?  
    printf("%e\n", really_small);  
    // 2.938736e-39 -> Como é possível?  
    printf("%e\n", really_small / 2.0f);  
}
```

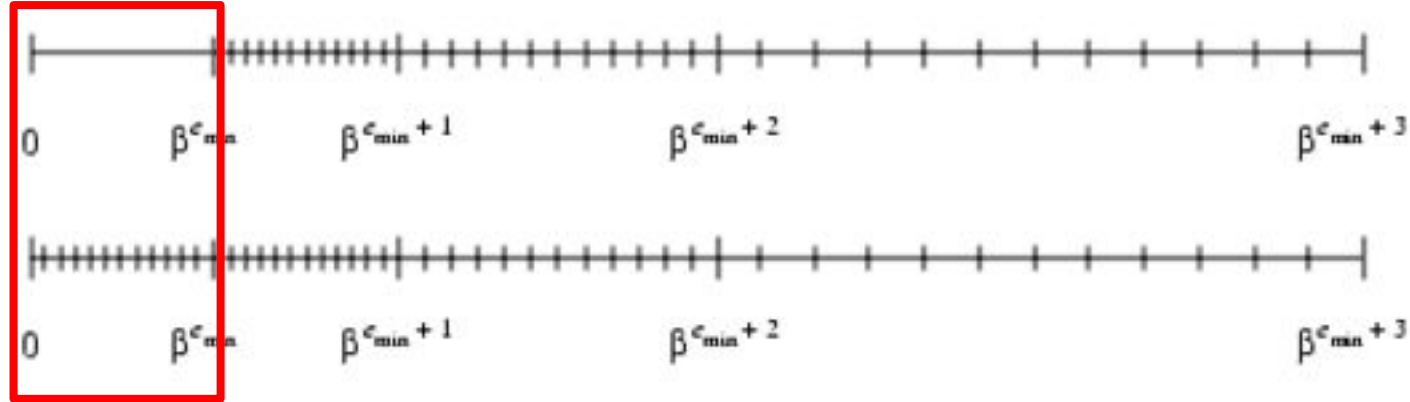
Representação - Números Denormais/Subnormais

+0	0 00000000 000000000000000000000000
-0	1 00000000 000000000000000000000000
+inf	0 11111111 000000000000000000000000
-inf	1 11111111 000000000000000000000000
NaN	<any signal> 11111111 <At least one 1 in m>
Denormal	0 00000000 <Any except all zeroes>

0 ou 1

0 . [m] * 2 ** -126

Representação - Números Denormais/Subnormais



Representação - Números Denormais/Subnormais

Mas.....

“(...) the speed of computation is significantly reduced on many modern processors”

Referência: https://en.wikipedia.org/wiki/Denormal_number

(O exemplo 4 tem relação com essa característica!)

Resultados Curiosos | Números Denormais

```
int main() {  
    float a = 1.5f;  
    for (int i = 0; i < 10000000; ++i) {  
        a = a * 1.5f;  
        a = a / 2.5f;  
  
        a = a + 0.1f; // Evita que números denormais sejam gerados  
        a = a - 0.1f; // (Mas note que o resultado será DIFERENTE)  
    }  
}
```

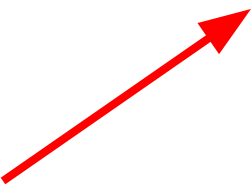


Performance (média de 10 rodadas): **0.151s**

Comparações com floats/doubles

Comparando floating points - Diferença absoluta

```
int main() {  
    float sum = 0.0f;  
    for (int i = 0; i < 100; ++i){  
        sum += 0.01f;  
    }  
  
    float eps = 1e-8;  
    float b = 0.01f * 100.0f;  
    if (fabs(sum - b) < eps) { printf("Equal!\n"); }  
    else { printf("Different!\n"); }  
}
```

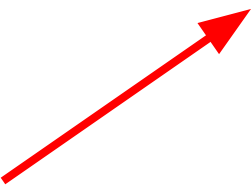

$$|a - b| < \epsilon$$

Se a e b forem muito grandes, sua diferença pode ser sempre um valor maior que epsilon

Comparando floating points - Diferença relativa

```
int main() {  
    float sum = 0.0f;  
    for (int i = 0; i < 100; ++i){  
        sum += 0.01f;  
    }
```

```
    float eps = 1e-8;  
    float b = 0.01f * 100.0f;  
    if (fabs((sum - b) / b) < eps) { printf("Equal!\n"); }  
    else { printf("Different!\n"); }  
}
```


$$\frac{|a - b|}{b} < \epsilon$$

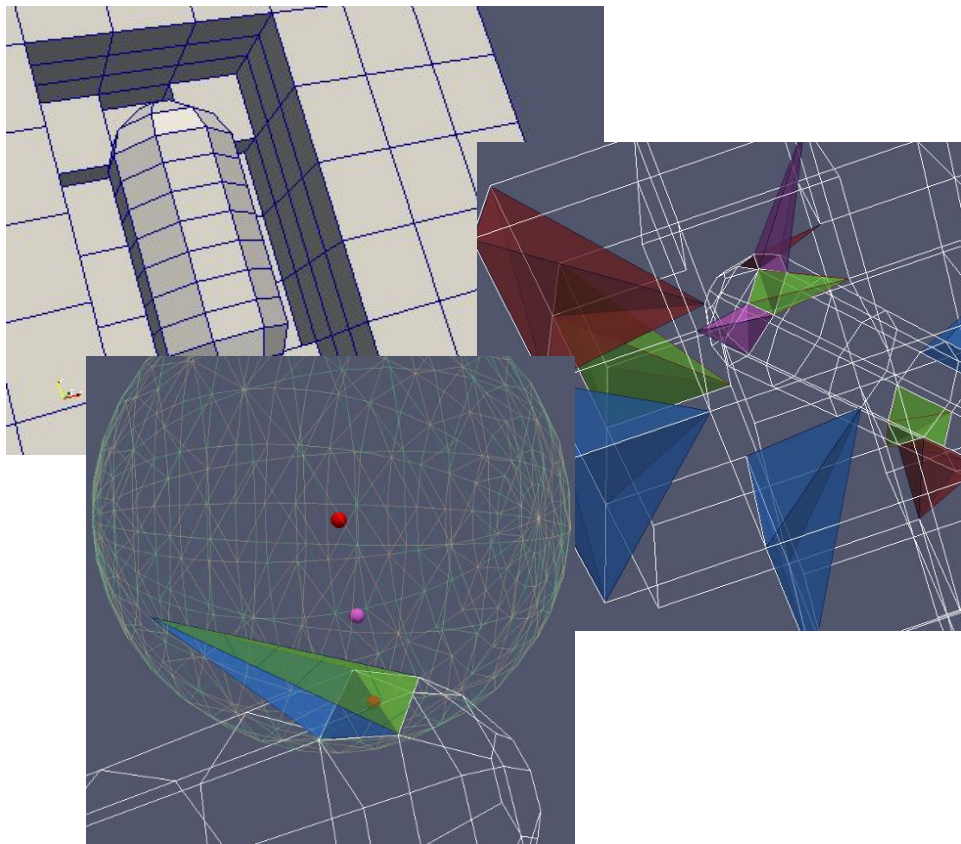
E se $b == 0$?

E se $a == b == 0$?

Comparar esses casos
separadamente?

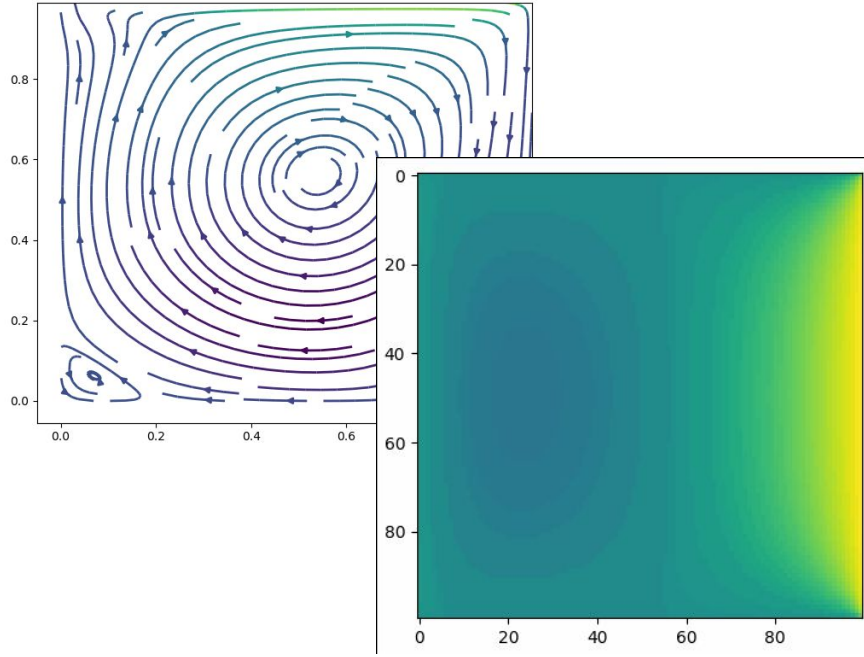
Exemplos reais

Geração automatizada de malhas tridimensionais



- Utiliza ponto flutuante para a representação dos nós no espaço 3D (x, y, z)
- Algoritmos de limpeza de nós muito próximos
- Algoritmos de detecção de elementos com volume “muito pequeno”
- Algoritmos de detecção de elementos muito “finos”

Simulação (Difusão de calor / Escoamento fluido)



- Utilização de ponto flutuante para a representação das grandezas físicas (Temperatura, Pressão, Velocidade...)
- Algoritmos numéricos iterativos para **aproximar** soluções de equações diferenciais
- Cuidados ao testar resultados (Utilização de métodos de comparação que levam em conta a grandeza esperada das soluções)

Perguntas?

Referências

“What Every Computer Scientist Should Know About Floating-Point Arithmetic” - DAVID GOLDBERG

<https://www.floating-point-gui.de/>

<https://randomascii.wordpress.com/>



RECOMENDADO!

Representação - Casos especiais

“The rule for determining the result of an operation that has **infinity** as an operand is simple: Replace **infinity** with a finite number x and take the **limit as $x \rightarrow \infty$** ”

“When the limit doesn't exist, the result is a NaN”

Referência: What Every Computer Scientist Should Know About Floating-Point Arithmetic

```
#include <stdio.h>
#include <math.h>
int main()
{
    float inf = 1.0f / 0.0f;
    printf("%.40f\n", 1.0f / inf);
    printf("%.40f\n", 4.0f - inf);
    printf("%.40f\n", sqrt(inf));
}
```

$$\lim_{x \rightarrow \infty} \{1/x\} = 0.00000000$$

$$\lim_{x \rightarrow \infty} \{4 - x\} = -\text{inf}$$

$$\lim_{x \rightarrow \infty} \{\text{sqrt}(x)\} = +\text{inf}$$

Representação - Casos especiais

“Since the sign bit can take on two different values, there are **two zeros, +0 and -0.**”

“If zero did not have a sign, then the relation $1/(1/x) = x$ would fail to hold when $x = +/-\infty$ ”

Referência: What Every Computer Scientist Should Know About Floating-Point Arithmetic

```
#include <stdio.h>
int main()
{
    float inf = 1.0f / 0.0f;
    printf("%.40f\n", 1.0f/+inf);           // +0.0
    printf("%.40f\n", 1.0f/(1.0f/+inf));    // +inf
    printf("%.40f\n", 1.0f/-inf);           // -0.0
    printf("%.40f\n", 1.0f/(1.0f/-inf));    // -inf
}
```