

Aritmética Computacional

Parte III

Histórico de revisões

2

Revisão	Data	Responsável	Descrição
0.1	03/2016	Prof. Cesar Zeferino	Primeira versão
0.2	04/2017	Prof. Cesar Zeferino	Atualização do modelo

Observação: Este material foi produzido por pesquisadores do Laboratório de Sistemas Embarcados e Distribuídos (LEDS – Laboratory of Embedded and Distributed Systems) da Universidade do Vale do Itajaí e é destinado para uso em aulas ministradas por seus pesquisadores.

Introdução

3

❑ Objetivo

- ❑ Conhecer a representação em formato de ponto flutuante

❑ Conteúdo

- ❑ Representação em ponto flutuante
- ❑ Operações em ponto flutuante

Introdução

❑ Bibliografia

- ❑ PATTERSON, David A.; HENNESSY, John L. Abstrações e tecnologias computacionais. *In*: _____. **Organização e projeto de computadores: a interface hardware/software**. 4. ed. Rio de Janeiro: Campus, 2014. cap. 3. Disponível em: <<http://www.sciencedirect.com/science/article/pii/B9788535235852000032>>. Acesso em: 25 abr. 2017.

- ❑ Edições anteriores
 - ❑ Patterson & Hennessy (2000, p. 160-172)
 - ❑ Patterson & Hennessy (2005, p. 142-158)

2 Ponto flutuante

☐ Números reais

3,141159265_{dez} (π)

0,000000001_{dez} = 1,0_{dez} $\times 10^{-9}$ (segundos em um nanossegundo)

3.155.760.000_{dez} = 3,15576_{dez} $\times 10^9$ (segundos em um século)

Notação Científica

☐ Notação científica normalizada

☐ Apenas um dígito (diferente de 0) à esquerda do ponto decimal

☐ Exemplos

1,0_{dez} $\times 10^{-9}$

Normalizado

0,1_{dez} $\times 10^{-8}$

Não Normalizado

10,0_{dez} $\times 10^{-10}$

Não Normalizado

2 Ponto flutuante

6

❑ Notação científica normalizada em binário

❑ Apenas um dígito à esquerda do ponto “binário”

❑ Exemplo

$$1, \textcolor{red}{xxxxxxxx}_{\text{dois}} \times 2^{\textcolor{red}{yyyyyy}}$$

2 Ponto flutuante

❑ Representação em ponto flutuante

$$(-1)^s \times F \times 2^E$$

❑ Onde

- ❑ s : sinal (0: positivo, 1: negativo)
- ❑ F : mantissa
- ❑ E : expoente

❑ Padrão IEEE 754

- ❑ Utilizado em quase todos os computadores
- ❑ A mantissa armazena apenas a parte fracionária
- ❑ Assume que o dígito à esquerda da vírgula é igual a 1

$$(-1)^s \times (1 + F) \times 2^E$$

- ❑ O número 0 é indicado fazendo $E = 0$

❏ Precisão simples (*float*)

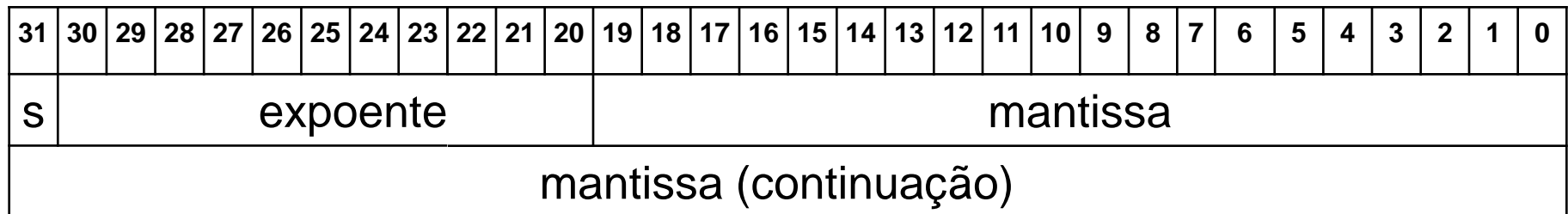
- ❑ 8 bits para o expoente
- ❑ 23 bits para a mantissa

[illegible]

2 Ponto flutuante

☐ Precisão dupla (*double*)

- ☐ 11 bits para o expoente
- ☐ 52 bits para a mantissa



- ☐ Intervalo de representação: $[2,0_{\text{dez}} \times 10^{-308}, 2,0_{\text{dez}} \times 10^{+308}]$

2 Ponto flutuante

❑ Porque o expoente é colocado antes da mantissa?

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
s	expoente								mantissa																						

❑ Para agilizar as comparações em ordenações

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0	0	0	0	0	0	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1

é maior que

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1

porque o seu expoente é maior

2 Ponto flutuante

- ❑ **Solução:** uso de notação com peso
- ❑ **IEEE 754 usa o peso 127 para precisão simples**
 - ❑ Expoente com peso = $127 + \text{Expoente original}$, ou seja
 - ❑ Expoente -1 é representado como 126
 - ❑ Expoente 0 é representado como 127
 - ❑ Expoente $+1$ é representado como 128
- ❑ **Representação**

$$(-1)^s \times (1 + \textit{Mantissa}) \times 2^{(\textit{Expoente} - \textit{Peso})}$$

- ❑ **IEEE 754 usa o peso 1.023 para precisão dupla**

2 Ponto flutuante

□ Potências de 2

$$\square 2^3_{(\text{dez})} = 8_{(\text{dez})} = 1000_{(\text{dois})}$$

$$\square 2^2_{(\text{dez})} = 4_{(\text{dez})} = 0100_{(\text{dois})}$$

$$\square 2^1_{(\text{dez})} = 2_{(\text{dez})} = 0010_{(\text{dois})}$$

$$\square 2^0_{(\text{dez})} = 1_{(\text{dez})} = 0001_{(\text{dois})}$$

$$\square 2^{-1}_{(\text{dez})} = 0,5000_{(\text{dez})} = 0,1000_{(\text{dois})}$$

$$\square 2^{-2}_{(\text{dez})} = 0,2500_{(\text{dez})} = 0,0100_{(\text{dois})}$$

$$\square 2^{-3}_{(\text{dez})} = 0,1250_{(\text{dez})} = 0,0010_{(\text{dois})}$$

$$\square 2^{-4}_{(\text{dez})} = 0,0625_{(\text{dez})} = 0,0001_{(\text{dois})}$$

□ Exemplos

$$8,5_{(\text{dez})} = 8_{(\text{dez})} + 0,5_{(\text{dez})} = 1000_{(\text{dois})} + 0,1000_{(\text{dois})} = 1000,1000_{(\text{dois})}$$

$$\begin{aligned} 10,375_{(\text{dez})} &= 10_{(\text{dez})} + 0,25_{(\text{dez})} + 0,125_{(\text{dez})} \\ &= 1010_{(\text{dois})} + 0,0100_{(\text{dois})} + 0,0010_{(\text{dois})} = 1010,0110_{(\text{dois})} \end{aligned}$$

2 Ponto flutuante

❑ **Exercícios:** Represente os números abaixo como ponto flutuante de precisão simples no padrão IEEE 754

(a) $+ 1,0_{(\text{dez})}$

(b) $+ 1,5_{(\text{dez})}$

(c) $+ 2,75_{(\text{dez})}$

(d) $- 9,4375_{(\text{dez})}$

(e) $+ 14,1875_{(\text{dez})}$

(f) $- 255,5_{(\text{dez})}$

2 Ponto flutuante

□ **Exercícios:** Represente os números abaixo como ponto flutuante de precisão simples no padrão IEEE 754

$$\begin{aligned} \text{(a)} \quad +1,0_{(\text{dez})} &= (-1)^0 \times (1 + 0,0) \times 2^{(127-127)} \\ &= 0_01111111_000000000000000000000000_{(\text{dois})} \end{aligned}$$

$$\begin{aligned} \text{(b)} \quad +1,5_{(\text{dez})} &= (-1)^0 \times (1 + 0,5) \times 2^{(127-127)} \\ &= 0_01111111_100000000000000000000000_{(\text{dois})} \end{aligned}$$

$$\begin{aligned} \text{(c)} \quad +2,75_{(\text{dez})} &= +1,375_{(\text{dez})} \times 2^1 \\ &= (-1)^0 \times (1 + 0,375) \times 2^{(128-127)} \\ &= 0_10000000_011000000000000000000000_{(\text{dois})} \end{aligned}$$

ou

$$\begin{aligned} &= +2,75_{(\text{dez})} = +10,11_{(\text{dois})} = +1,011_{(\text{dois})} \times 2^1_{(\text{dez})} \\ &= (-1)^0 \times (1 + 0,011_{(\text{dois})}) \times 2^{(128-127)} \\ &= 0_10000000_011000000000000000000000_{(\text{dois})} \end{aligned}$$

2 Ponto flutuante

❑ **Exercícios:** Represente os números abaixo como ponto flutuante de precisão simples no padrão IEEE 754

$$\begin{aligned}
 \text{(d)} - 9,4375_{(\text{dez})} &= -1001,0111_{(\text{dois})} = -1, 0010111_{(\text{dois})} \times 2^3_{(\text{dez})} \\
 &= (-1)^1 \times (1 + 0,0010111) \times 2^{(130-127)} \\
 &= 1_10000010_001011100000000000000000_{(\text{dois})}
 \end{aligned}$$

$$\begin{aligned}
 \text{(e)} + 14,1875_{(\text{dez})} &= +1110,0011_{(\text{dois})} = +1, 1100011_{(\text{dois})} \times 2^3_{(\text{dez})} \\
 &= (-1)^0 \times (1 + 0,1100011) \times 2^{(130-127)} \\
 &= 0_10000010_110001100000000000000000_{(\text{dois})}
 \end{aligned}$$

$$\begin{aligned}
 \text{(f)} - 255,5_{(\text{dez})} &= -11111111,1_{(\text{dois})} = -1, 11111111_{(\text{dois})} \times 2^7_{(\text{dez})} \\
 &= (-1)^1 \times (1 + 0,11111111) \times 2^{(134-127)} \\
 &= 1_10000110_111111110000000000000000_{(\text{dois})}
 \end{aligned}$$

2 Ponto flutuante

17

❑ **Exercícios:** Obtenha o equivalente decimal dos seguintes números representados em ponto flutuante com precisão simples no padrão IEEE 754

(a) 0 10000000 11110000...0000

(b) 1 10000001 10101000...0000

(c) 0 01111110 01100000...0000

2 Ponto flutuante

❑ **Exercícios:** Obtenha o equivalente decimal dos seguintes números representados em ponto flutuante com precisão simples no padrão IEEE 754

(a) 0 10000000 11110000...0000

$$(-1)^0 \times (1 + 0,1111) \times 2^{(128-127)} = 1,1111 \times 2^1 = 11,111 = 3,875_{(\text{dez})}$$

(b) 1 10000001 10101000...0000

$$(-1)^1 \times (1 + 0,10101) \times 2^{(129-127)} = -1,10101 \times 2^2 = -110,101 = -6,625_{(\text{dez})}$$

(c) 0 01111110 01100000...0000

$$(-1)^0 \times (1 + 0,011) \times 2^{(126-127)} = 1,011 \times 2^{-1} = 0,1011 = 0,6875_{(\text{dez})}$$

2 Adição e subtração em ponto flutuante

19

- ❑ Como funciona no sistema decimal?
- ❑ Exemplo: $9,999_{(dez)} \times 10^1 + 1,610_{(dez)} \times 10^{-1}$
- ❑ Restrições
 - ❑ Só pode armazenar 04 dígitos da mantissa
 - ❑ Só pode armazenar 02 dígitos do expoente

2 Adição e subtração em ponto flutuante

❑ Passo 1 – Alinhar o ponto decimal

$$1,610_{(\text{dez})} \times 10^{-1} = 0,1610_{(\text{dez})} \times 10^0 = 0,0161_{(\text{dez})} \times 10^1$$

Como só pode armazenar 4 dígitos: $0,016_{(\text{dez})} \times 10^1$

❑ Passo 2 – Somar as mantissas

$$\begin{array}{r} 9,999_{(\text{dez})} \\ + \underline{0,016_{(\text{dez})}} \\ \hline 10,015_{(\text{dez})} \end{array}$$

2 Adição e subtração em ponto flutuante

21

- ❑ **Passo 3 – Normalizar a soma e verificar a ocorrência de overflow**

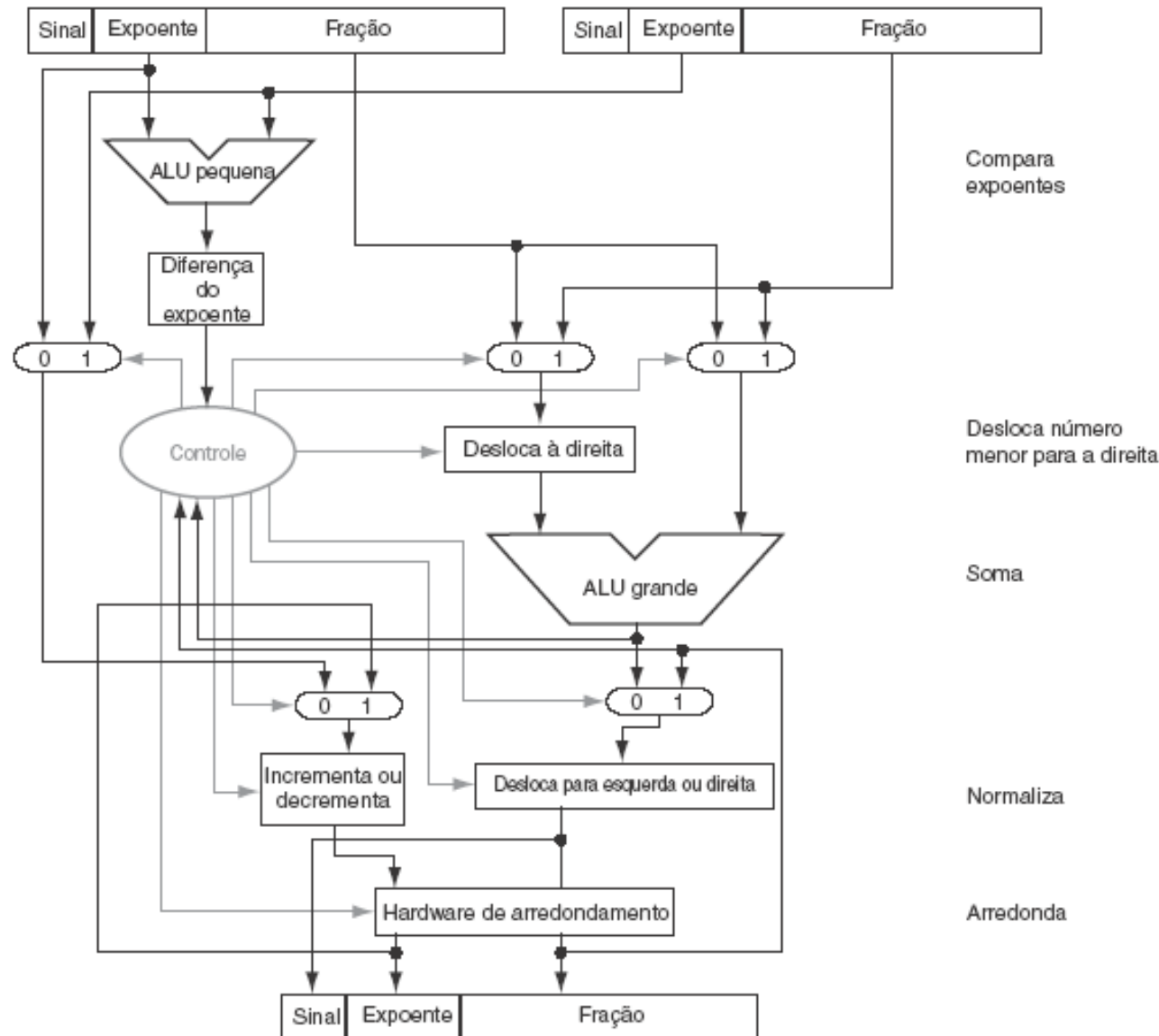
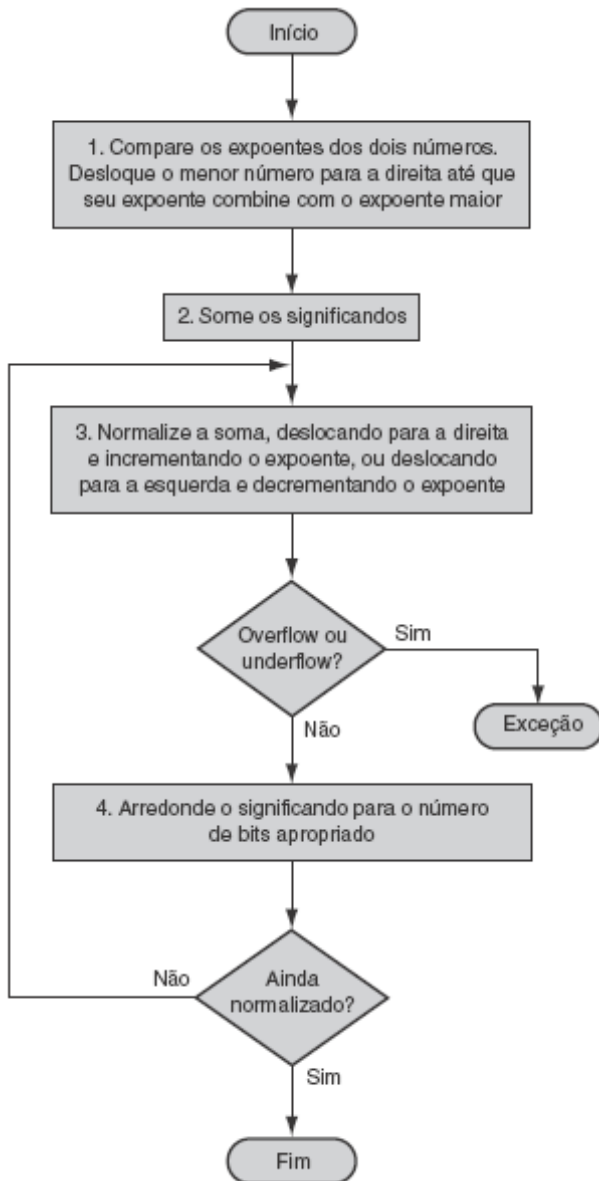
$$10,015_{(dez)} \times 10^1 = 1,0015_{(dez)} \times 10^2$$

- ❑ **Passo 4 – Arredondar o número**

$$1,0015_{(dez)} \times 10^2 = 1,002_{(dez)} \times 10^2$$

Se o resultado não estiver normalizado, retornar ao Passo 3

2 Adição e subtração em ponto flutuante



2 Adição e subtração em ponto flutuante

- ❑ Para suportar instruções em ponto flutuante, o MIPS possui um coprocessador (CP1) com 32 registradores de ponto flutuante: \$f0 a \$f31
- ❑ Instruções de soma e subtração

Categoria	Instrução	Exemplo	Significado
Aritmética	Adição de precisão simples	<code>add.s \$f2, \$f4, \$f6</code>	$\$f2 = \$f4 + \$f6$
	Adição de precisão dupla	<code>add.d \$f2, \$f4, \$f6</code>	$\$f2 = \$f4 + \$f6$
	Subtração de precisão simples	<code>sub.s \$f2, \$f4, \$f6</code>	$\$f2 = \$f4 - \$f6$
	Subtração de precisão dupla	<code>sub.d \$f2, \$f4, \$f6</code>	$\$f2 = \$f4 - \$f6$

2 Adição e subtração em ponto flutuante

24

❑ Instruções especiais para transferência de dados com os registradores do Coprocessador 1

Categoria	Instrução	Exemplo	Significado
Transferência de dados	Carga	<code>lwc1 \$f1, 100 (\$s2)</code>	<code>\$f1 = Mem[\$s2+100]</code>
	Armazenamento	<code>swc1 \$f1, 100 (\$s2)</code>	<code>Mem[\$s2+100] = \$f1</code>

2 Adição e subtração em ponto flutuante

25

❑ Instruções especiais para desvio e comparação

Categoria	Instrução	Exemplo	Significado
Desvio condicional	Desvio em FP verdadeiro	<code>bclt 25</code>	Se (cond == 1) vá para PC = PC + 4 + 100
	Desvio em FP falso	<code>bclf 25</code>	Se (cond == 0) vá para PC = PC + 4 + 100
	Comparação FP precisão simples (eq, ne, lt, le, gt, ge)	<code>c.lt.s \$f2, \$f4</code>	Se ($\$f2 < \$f4$) cond = 1 Senão cond = 0
	Comparação FP dupla precisão (eq, ne, lt, le, gt, ge)	<code>c.lt.d \$f2, \$f4</code>	Se ($\$f2 < \$f4$) cond = 1 Senão cond = 0