

CENTRO TECNOLÓGICO  
DEPARTAMENTO DE INFORMÁTICA

Arquitetura de Computadores I – Turmas 01 e 02 (EARTE) – 2021/2  
Prof. Rodolfo da Silva Villaça – [rodolfo.villaca@ufes.br](mailto:rodolfo.villaca@ufes.br)

**Laboratório V – Aritmética de Ponto Flutuante**

**Integrantes do Grupo:**

- Dionatas Brito
- Maria Julia Damasceno
- Otávio Sales

**1. Objetivo**

- Conhecer as operações de ponto flutuante e sua representação numérica;
- Usar as instruções de ponto flutuante.

**2. Introdução**

O padrão IEEE-754 reserva vários padrões de bits para ter um significado especial. Em outras palavras, nem todos os padrões de bits representam algum número, conforme quadro a seguir.

**Special bit patterns in IEEE-754**

Sign bit	Exponent	Significand	Comment
x	0..0	0..0	Zero
x	0..0	not all zeros	Denormalized number
0	1..1	0..0	Plus infinity (+inf)
1	1..1	0..0	Minus infinity (-inf)
x	1..1	not all zeros	Not a Number (NaN)

Infinito significa algo grande demais para ser representado. Um estouro de representação pode retornar um + inf ou um -inf. Algumas operações no infinito retornam outro infinito como resultado. Um resumo dessa representação encontra-se na Tabela 1.

**Table 1: Some operations with infinity**

Operation	Result	Comment
$x + (\text{inf})$	inf	x finite
$x - (+\text{inf})$	-inf	x finite
$(+\text{inf}) + (+\text{inf})$	+inf	
$(-\text{inf}) + (-\text{inf})$	-inf	
$x * (+\text{inf})$	+inf if $x > 0$ , -inf otherwise	x nonzero

Pode haver um zero positivo se o bit de sinal for 0 e um zero negativo (bit de sinal é 1). Números desnormalizados estão incluídos no padrão IEEE-754 para lidar com casos de estouro negativo de expoente (números muito pequenos). Um NaN (às vezes denotado por nan) é usado para representar um resultado indeterminado. Existem dois tipos de NaNs: sinalização e silêncio. O padrão de bits no significando é usado para diferenciar entre eles, e é dependente da implementação. Um NaN de sinalização pode ser usado, por exemplo, para variáveis não inicializadas. Observe que qualquer operação em um NaN de sinalização terá

**CENTRO TECNOLÓGICO**  
**DEPARTAMENTO DE INFORMÁTICA**

como resultado, um NaN silencioso. Operar em um NaN silencioso simplesmente retorna outro NaN sem gerar nenhuma exceção, vide Tabela 2.

**Table 2: Some operations which produce a quiet NaN**

Operation	Comment
$x + (\text{NaN})$	Any operation on a quiet NaN (addition in this example)
$(+\text{inf}) + (-\text{inf})$	
$0 * (\text{inf})$	
$0/0$	
$\text{inf}/\text{inf}$	
$x\%0$	The remainder of division by 0
$\sqrt{x}, x < 0$	

**Atividade 1**

Usando o simulador MARS:

- Declare as variáveis Zero.s, PlusInf.s, MinusInf.s, PlusNaN.s, MinusNaN, inicializadas com os padrões de bits correspondendo a zero, mais infinito, menos infinito, NaN positivo, NaN negativo em representação de precisão simples.
- Declare as variáveis Zero.d, PlusInf.d, MinusInf.d, NaN.d, inicializadas com os padrões de bits correspondendo a zero, mais infinito, menos infinito, NaN positivo, NaN negativo em representação de precisão dupla.
- Carregue essas variáveis em registradores de ponto flutuante, começando com \$f0
- Imprima, começando com \$f0, o conteúdo dos registradores onde as variáveis foram carregadas; imprime um caractere de nova linha (\n) após cada valor.

Execute o programa e preencha a tabela a seguir com o nome de cada variável e o respectivo valor impresso após a execução:

Variável	Saída impressa
Zero.s e Zero.d	0.0
PlusInf.s PlusInf.d	Infinity
MinusInf.s MinusInf.d	-Infinity
PlusNaN.s PlusNaN.d	NaN
MinusNaN.s MinusNaN.d	NaN

Qual é o padrão de bits para o maior número possível de ponto flutuante de precisão única? Escreva em hexadecimal.

O padrão de bits é todos os bits setados como 1 menos o bit de sinal, que precisa ser 0. O número em hexadecimal é 7FFFFFFF.

**CENTRO TECNOLÓGICO**  
**DEPARTAMENTO DE INFORMÁTICA**

**Atividade 2**

Nesta atividade você deverá criar um programa que calcula o fatorial de um número inteiro, representado usando notação de números de ponto flutuante, conforme passos a seguir:

- Solicite ao usuário que insira um número inteiro;
- Verifique se o número inserido é negativo: se for negativo, imprima uma mensagem de erro e solicite o usuário para entrar novamente com um número inteiro positivo
- Realizar a operação 'FactorialSingle', cujo parâmetro será o número lido do usuário convertido em ponto flutuante de precisão simples. A operação deve retornar o fatorial (em PF precisão simples) desse número
- Imprimir o valor retornado por 'FactorialSingle'

Execute o 'FactorialSingle' para completar o plano de testes a seguir. Use notação científica normalizada para representar a saída impressa por 'FactorialSingle'.

Número	Fatorial (Inteiro)	Fatorial PF (single)
0	1	1.0
5	120	120.0
10	3628800	3628800.0
15	2004310016	1.30767428E12
20	-2102132736	2.4329023E18
40	0	Infinity

A seguir, responda:

a) Por que alguns dos resultados impressos são negativos?

A operação de multiplicação muda o bit de sinal, fazendo com que o número seja interpretado como negativo.

b) Quais são os valores máximos da entrada para os quais a saída correta ainda é impressa?

Para o fatorial inteiro, o máximo é 12 (479001600). 13 já gera um overflow e imprime a saída errada.

c) Usando algum outro método ou linguagem de programação, calcule o valor exato de 20 e compare com o valor impresso pelo seu programa no MARS. Quais são esses valores? Por quê os valores são diferentes?

2.432902E18 ou 2432902008176640000. Os valores são diferentes devido a aproximações e limitações na quantidade de números trabalhados.

**Atividade 3**

**CENTRO TECNOLÓGICO**  
**DEPARTAMENTO DE INFORMÁTICA**

Na Atividade 2 você usou a conversão de inteiro para ponto flutuante (cvt.s.w). Vamos agora tentar uma conversão de PF para inteiro, conforme instruções a seguir:

- Após de imprimir o valor retornado por 'FactorialSingle', converta esse valor em um número inteiro e imprima-o também.

Execute esse novo programa e conclua o próximo plano de teste. Anote o valor impresso para o fatorial como está, não use notação científica neste momento.

Número	Fatorial (Impressão como Inteiro)	Fatorial (Impressão como PF simples)
0	0	0.0
5	120	120.0
10	3628800	3628800.0
11	39916800	3.99168E7
12	479001600	4.790016E8
13	2147483647	6.2270208E9
14	2147483647	8.7178289E10
15	2147483647	1.30767428E12

Destaque os casos em que a saída de inteiro é um número diferente da saída de ponto flutuante. **A** operação de conversão produz um inteiro com ou sem sinal? Explique

A saída é diferente para os casos acima de 12 (13 em diante, como dito na questão anterior). A operação de conversão produz um inteiro com sinal, uma vez que o limite é 2147483647, o valor máximo de um inteiro com sinal.

#### 4. Execução

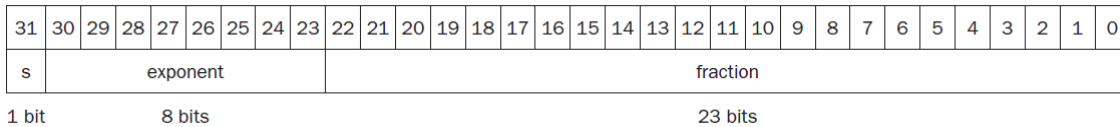
- Grupos de até 3 (três) alunos;
- Submissão até 17/12 (23:59h);

#### 5. Apêndice – PF no MARS

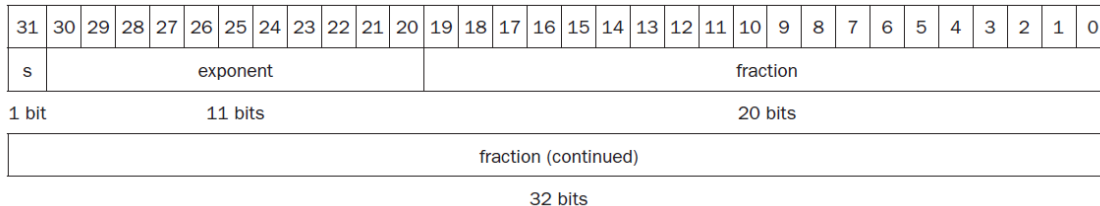
A aritmética de Ponto Flutuante é implementada pelo coprocessador 1 (*Coproc 1* no MARS) da arquitetura MIPS. O coprocessador possui 32 registradores de 32 bits, numerados de 0 a 31 (\$f0 a \$f31). Os valores armazenados nestes registradores seguem o padrão IEEE-754 (ver material do curso).

Para permitir o armazenamento de valores em precisão dupla (64 bits), a arquitetura usa o artifício de agrupar pares de registradores subsequentes, isto é, \$f0 com \$f1, \$f2 com \$f3, e assim, sucessivamente. Assim, a arquitetura oferece “virtualmente” 16 registradores de 64 bits, nos quais o valor armazenado terá sempre seus 32 bits de mais alta ordem armazenados num registrador par e os 32 de mais baixa ordem num registrador ímpar. A representação em IEEE-754 para 32 e 64 bits é mostrada abaixo:

**CENTRO TECNOLÓGICO  
DEPARTAMENTO DE INFORMÁTICA**



**Representação IEEE-754 em precisão simples (32 bits)**



**Representação IEEE-754 em precisão dupla (64 bits)**

Para simplificar as coisas, operações em ponto flutuante usam sempre registradores numerados e podem ser executadas em precisão simples (32 bits) ou precisão dupla (64 bits). A diferença das instruções MIPS é determinada por um sufixo após o mnemônico da instrução, por exemplo, as instruções `add.s` e `add.d` representam adições em precisão simples (.s) e precisão dupla (.d), respectivamente.

O conjunto de instruções MIPS oferece, além das instruções aritméticas, comparações, desvios, movimentação de dados da (*load*) e para (*store*) a memória, conversões de formatos de ponto-flutuante (32 para 64 bits e vice-versa) e conversão de inteiro para ponto flutuante e vice-versa. Enquanto nas comparações envolvendo inteiros um dos registradores de propósito geral pode ser definido como destino da execução, no caso de ponto flutuante, uma comparação irá implicitamente determinar ("setar") uma *flag*. Esta *flag* poderá então ser usada para testar se um desvio é realizado ou não numa instrução de desvio condicional.

Instruction	Comment
<code>mfc1 Rdest, FPsrc</code>	Move the content of floating-point register <code>FPsrc</code> to <code>Rdest</code>
<code>mtc1 Rsrc, FPdest</code>	Integer register <code>Rsrc</code> is moved to floating-point register <code>FPdest</code>
<code>mov.x FPdest, FPsrc</code>	Move floating-point register <code>FPsrc</code> to <code>FPdest</code>
<code>lwc1 FPdest, address</code>	Load word from address in register <code>FPdest</code> <sup>a</sup>
<code>swc1 FPsrc, address</code>	Store the content of register <code>FPsrc</code> at address <sup>b</sup>
<code>add.x FPdest, FPsrc1, FPsrc2</code>	Add single precision

**CENTRO TECNOLÓGICO  
DEPARTAMENTO DE INFORMÁTICA**

Instruction	Comment
<code>sub.x FPdest, FPsrc1, FPsrc2</code>	Subtract <code>FPsrc2</code> from <code>FPsrc1</code>
<code>mul.x FPdest, FPsrc1, FPsrc2</code>	Multiply
<code>div.x FPdest, FPsrc1, FPsrc2</code>	Divide <code>FPsrc1</code> by <code>FPsrc2</code>
<code>abs.s FPdest, FPsrc</code>	Store the absolute value of <code>FPsrc</code> in <code>FPdest</code>
<code>neg.x FPdest, FPsrc</code>	Negate number in <code>FPsrc</code> and store result in <code>FPdest</code>
<code>c.eq.x FPsrc1, FPsrc2</code>	Set the floating-point condition flag to true if the two registers are equal
<code>c.le.x FPsrc1, FPsrc2</code>	Set the floating-point condition flag to true if <code>FPsrc1</code> is less than or equal to <code>FPsrc2</code>
<code>c.lt.x FPsrc1, FPsrc2</code>	Set the floating-point condition flag to true if <code>FPsrc1</code> is less than <code>FPsrc2</code>
<code>bc1t label</code>	Branch if the floating-point condition flag is true
<code>bc1f label</code>	Branch if the floating-point condition flag is false
<code>cvt.x.w FPdest, FPsrc</code>	Convert the integer in <code>FPsrc</code> to floating-point
<code>cvt.w.x FPdest, FPsrc</code>	Convert the floating-point number in <code>FPsrc</code> to integer
<code>cvt.d.s FPdest, FPsrc</code>	Convert the single precision number in <code>FPsrc</code> to double precision and put the result in <code>FPdest</code>
<code>cvt.s.d FPdest, FPsrc</code>	Convert the double precision number in <code>FPsrc</code> to single precision and put the result in <code>FPdest</code>