

Floating Point Addition/ Subtraction

- $35.23142 + 0.00053$

$$X = 35.23142$$

X (Binary)
= 100011.0011101111

X (Binary Normalized) =
 $1.000110011101111 \times 2^5$

$$Y = 0.00053$$

Y (Binary)
= 0.000000000010001011

Y (Binary Normalized) =
 1.0001011×2^{11}

Y (Binary) =
 $0.00000000000000010001011 \times 2^5$

Rule: Match the Lower
Exponent with the
Higher Exponent

$$\begin{aligned} X + Y &= (1.000110011101111 \times 2^5) + (0.00000000000000010001011 \times 2^5) \\ &= (1.000110011101111 + 0.00000000000000010001011) \times 2^5 \\ &= 1.0001100111011110001011 \times 2^5 \\ &= 100011.00111011110001011 \\ &= 35.2342224121 \text{ (Decimal)} \end{aligned}$$

Floating Point Multiplication

- $5.234 \times (-0.003)$

$$X = 5.234$$

X (Binary)
= 101.0011101111

X (Binary Normalized) =
 $1.010011101111 \times 2^2$

$$Y = 0.003$$

Y (Binary)
= 0.0000000011000100101

Y (Binary Normalized) =
 $1.1000100101 \times 2^{-9}$

$$\begin{aligned} X \times Y &= - (1.010011101111 \times 2^2) \times (1.1000100101 \times 2^{-9}) \\ &= - (1.010011101111 \times 1.1000100101) \times 2^{(2+(-9))} \\ &= - (1.010011101111 \times 1.1000100101) \times 2^{(-7)} \\ &= - 10.000000101 \times 2^{(-7)} \\ &= - 0.0000010000000101 \times 2^{(0)} \\ &= - 0.0157012939 \text{ (Decimal)} \end{aligned}$$

Floating Point Arithmetic

- 51500000 – BA10A000

X = 51500000

X (Binary)

= 0101 0001 0101 0000 0000 0000 0000 [32 bit]

= 0 10100010 10100000000000000000000000000000

| (Biased) | | |
|----------|----------|--------------------|
| Sign bit | Exponent | Fraction/ Mantissa |
| 1 bit | 8 bit | 23 bit |

Find Out Exponent and Fraction

Biased Exponent = 10100010

Biased Exponent (Decimal) = 162 For Exponent being 8 bit, Bias = $2^{(8-1)} - 1 = 127$

Exponent (Decimal) = 162 - 127

= 35

Fraction/ Mantissa = 0. 10100000000000000000000000000000

X (Binary Normalized) = 1.Fraction x $2^{(\text{Exponent})}$

X (Binary Normalized) = 1. 10100000000000000000000000000000 x 2^{35} *Sign bit = 0 (Positive)

Floating Point Arithmetic

- 51500000 – BA10A000

Y = BA10A000

Y (Binary)

= 1011 1010 0001 0000 1010 0000 0000 0000 [32 bit]

= 1 01110100 001000101000000000000000

| (Biased) | | |
|----------|----------|--------------------|
| Sign bit | Exponent | Fraction/ Mantissa |
| 1 bit | 8 bit | 23 bit |

Find Out Exponent and Fraction

Biased Exponent = 01110100

Biased Exponent (Decimal) = 116 For Exponent being 8 bit, Bias = $2^{(8-1)} - 1 = 127$

Exponent (Decimal) = 116 - 127

= -11

Fraction/ Mantissa = 0. 00100010100000000000000000000000

Y (Binary Normalized) = 1.Fraction x $2^{(\text{Exponent})}$

Y (Binary Normalized) = - 1.0010001010000000000000000000000 x 2^{-11} *Sign bit = 1 (Negative)

Floating Point Addition/ Subtraction

- 51500000 – BA10A000

Rule: Match the Lower Exponent with the Higher Exponent

X = 51500000

Y = BA10A000

X (Binary Normalized) =

$$1.1010000000000000000000000000000 \times 2^{35}$$

Y (Binary Normalized) =

$$-1.0010000101000000000000000 \times 2^{-11}$$

Y (Binary Normalized) =

$$-0.[45\text{ Os..}] 1001000010100000000000000 \times 2^{35}$$

X - (-Y) =

X + Y =

$$= (1.1010000000000000000000000000000 \times 2^{35}) + (0.[45\text{ Os..}] 1001000010100000000000000 \times 2^{35})$$

$$= (1.1010000000000000000000000000000 + 0.[45\text{ Os..}] 1001000010100000000000000) \times 2^{35}$$

$$= 1.101[42\text{ Os..}] 1001000010100000000000000 \times 2^{35}$$

Floating Point Arithmetic

- 7ACD0000 + 5BCA0000

| (Biased) | | |
|----------|----------|--------------------|
| Sign bit | Exponent | Fraction/ Mantissa |
| 1 bit | 8 bit | 23 bit |

X = 7ACD0000

X (Binary)

$$= 0111\ 1010\ 1100\ 1101\ 0000\ 0000\ 0000\ 0000 \quad [32\text{ bit}]$$

$$= 0\ \underline{1110101}\ \underline{1001101000000000000000000}$$

Find Out Exponent and Fraction

Biased Exponent = 11110101

Biased Exponent (Decimal) = 245 For Exponent being 8 bit, Bias = $2^{(8-1)} - 1 = 127$

Exponent (Decimal) = 245 - 127

= 118

Fraction/ Mantissa = 0.1001101000000000000000000

X (Binary Normalized) = 1.Fraction x 2^(Exponent)

X (Binary Normalized) = 1. 1001101000000000000000000 x 2¹¹⁸ *Sign bit = 0 (Positive)

Floating Point Arithmetic

- 7ACD0000 + 5BCA0000

Y = 5BCA0000

Y (Binary)

= 0101 1011 1100 1010 0000 0000 0000 0000 [32 bit]

= 0 10110111 10010100000000000000000000000000

| (Biased) | | |
|----------|----------|--------------------|
| Sign bit | Exponent | Fraction/ Mantissa |
| 1 bit | 8 bit | 23 bit |

Find Out Exponent and Fraction

Biased Exponent = 10110111

Biased Exponent (Decimal) = 183 For Exponent being 8 bit, Bias = $2^{(8-1)} - 1 = 127$

Exponent (Decimal) = 183 - 127

= 56

Fraction/ Mantissa = 0. 1001010000000000000000000

Y (Binary Normalized) = 1.Fraction x 2^{Exponent}

Y (Binary Normalized) = 1. 1001010000000000000000000 x 2^{56}

Floating Point Addition/ Subtraction

- 7ACD0000 + 5BCA0000

Rule: Match the Lower Exponent with the Higher Exponent

X = 51500000

Y = 3A10A000

X (Binary Normalized) =

1. 1001101000000000000000000 x 2^{118}

Y (Binary Normalized) =

1. 1001010000000000000000000 x 2^{56}

Y (Binary Normalized) =

0.[61 0s..] 1001010000000000000000000 x 2^{118}

X + Y =

= 1. 1001101000000000000000000 x 2^{118} + 0.[61 0s..] 1001000010100000000000000 x 2^{118}

= (1. 1001101000000000000000000 + 0.[61 0s..] 1001000010100000000000000) x 2^{118}

= 1.1001101[54 0s..] 1001000010100000000000000 x 2^{118}

8. Suppose X=19.454 and Y=3.0124, perform X*Y using IEEE floating-point representation.

Answer8:

X= 19.454

X (Binary) = 10011.01110100

X (Normalized) = 1.001101110100 $\times 2^4$

Y= 3.012

X (Binary) = 11.0000001100

X (Normalized) = 1.10000001100 $\times 2^1$

X * Y =

(1.001101110100 $\times 2^4$) \times (1.10000001100 $\times 2^1$)

= (1.001101110100 \times 1.10000001100) $\times 2^{4+1}$

= 1.11010100101100101110000 $\times 2^5$

= 111010.100101100101110000

= 58.58734... (Decimal)

$$\begin{array}{r} 1.001101110100 \\ \times 1.10000001100 \\ \hline 00000000000000 \\ 1001101110100xx \\ 1001101110100xxx \\ 1001101110100 xxxxxxxxxx \\ 1001101110100 xxxxxxxxxxxx \\ \hline 1.11010100101100101110000 \end{array}$$

Practice

Consider the value 63.7813

a) Let's assume you have a 21-bit register having 6-bit for exponent. Now convert this value using IEEE floating-point representation. Also, convert this into hexadecimal form.

Ans: 49FC8

b) Let's assume you have a 12-bit register having 4-bit for exponent. Now convert this value using IEEE floating-point representation. Also, convert this into hexadecimal form

Ans: 67F

c) Suppose X= -9.435 and Y= 15.129, perform X*Y using IEEE floating-point representation.

Ans: -142.719955... (Decimal)

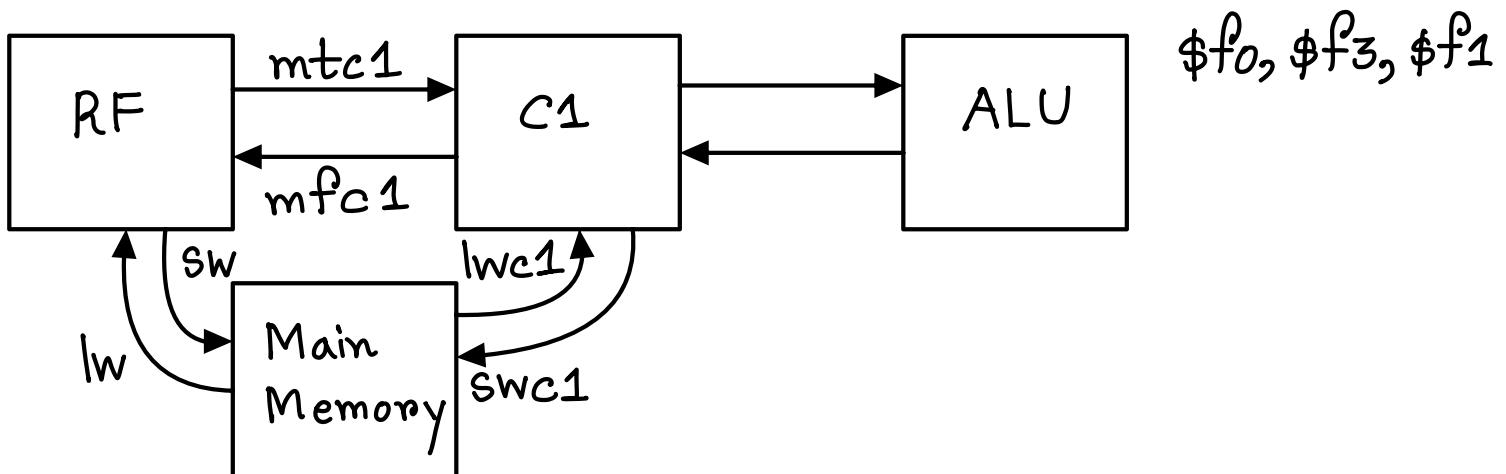
MIPS Division

- Use HI/LO registers for result
 - HI: 32-bit remainder
 - LO: 32-bit quotient
- Instructions
 - div rs, rt / divu rs, rt
 - No overflow or divide-by-0 checking
 - Software must perform checks if required
 - Use mfhi, mflo to access result

Now Written Part!

Next Quiz: Chapter → 3

Floating point Registers and Instruction:



$A[3] = \overbrace{\text{float}}^{\text{type}} x$

Where x is in $\$f_0$ and base address of A is in $\$s_1$.

$\text{swc1 } \$f_0, 12(\$s_1)$

$\text{float}(a) = B[2]$

Where a is in $\$f_0$ and base address of B is in $\$s_1$.

$\text{lwc1 } \$f_0, 8(\$s_1)$

Single Precision

$\$f_0 = \$f_1 + \$f_2$

$\text{add.s } \$f_0, \$f_1, \$f_2$

$\$f_0 = \$f_1 - \$f_2$

$\text{sub.s } \$f_0, \$f_1, \$f_2$

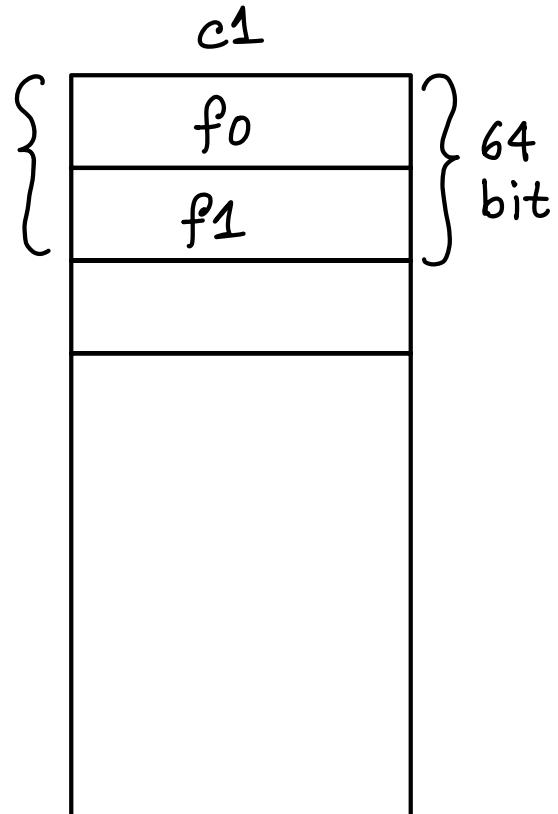
Double Precision

$\$f_0 = \$f_2 + \$f_4$

$\text{add.d } \$f_0, \$f_2, \$f_4$

$\$f_0 = \$f_2 - \$f_4$

$\text{sub.d } \$f_0, \$f_2, \$f_4$



Moving values from register

$\text{mtc1 } \$s_0, \f_1

$\text{mfcc1 } \$s_0, \f_1

Convert a register value from float to integer for single precision: cvt.w.s \$f₀, \$f₁

Convert a register value from integer to float for single precision:

cvt.s.w \$f₀, \$f₁

int

Double precision to single precision

cvt.s.d \$f₀, \$f₁

Single precision to double precision

cvt.d.s \$f₀, \$f₁

$$\text{float}(x) = (\text{int})y + \text{float}(z)$$

$$(\text{int})y = (\text{float})x - (\text{float})z$$

Suppose, x, y, z are in \$f₁, \$s₁, and \$f₂.

Mips Code:

mtc1 \$s₁, \$f₃

cvt.s.w \$f₃, \$f₃ (Possible)

add.s \$f₁, \$f₃, \$f₂

sub.s \$f₄, \$f₁, \$f₂

cvt.w.s \$f₄, \$f₄

mtc1 \$s₁, \$f₄