## CS/ECE 3810: Computer Organization

# Lecture 6: Floating Point

Anton Burtsev September, 2022 Motivation...

# Floating Point Standard

- Defined by IEEE Std 754-1985
  - Institute of Electrical and Electronics Engineers
- Developed in response to divergence of representations
  - Portability issues for scientific code
- Now almost universally adopted
- Two representations
  - Single precision (32-bit)
  - Double precision (64-bit)

 Normalized scientific notation: single non-zero digit to the left of the decimal (binary) point – example: 3.5 x 10<sup>9</sup>

- 1.010001 x  $2^{-5}_{two} = (1 + 0 \times 2^{-1} + 1 \times 2^{-2} + ... + 1 \times 2^{-6}) \times 2^{-5}_{ten}$
- A standard notation enables easy exchange of data between machines and simplifies hardware algorithms – the IEEE 754 standard defines how floating point numbers are represented

#### Sign and Magnitude Representation

Sign	Exponent	Fraction
<u>1 bi</u> t	8 bits	23 bits
S	E	F

- More exponent bits wider range of numbers (not necessarily more numbers – recall there are infinite real numbers)
- More fraction bits 📥 higher precision
- Register value =  $(-1)^{s} \times F \times 2^{E}$
- Since we are only representing normalized numbers, we are guaranteed that the number is of the form 1.xxxx..
   Hence, in IEEE 754 standard, the 1 is implicit
   Register value = (-1)<sup>s</sup> x (1 + F) x 2<sup>E</sup>

#### Sign and Magnitude Representation

Sign Exponent		Fraction
<u>1 bi</u> t	8 bits	23 bits
S	E	F

- Largest number that can be represented: 2.0 x 2<sup>128</sup> = 2.0 x 10<sup>38</sup> (not really – see upcoming details)
- Smallest number that can be represented: 1.0 x 2<sup>-127</sup> = 2.0 x 10<sup>-38</sup> (not really – see upcoming details)
- Overflow: when representing a number larger than the max; Underflow: when representing a number smaller than the min
- Double precision format: occupies two 32-bit registers: Largest: Smallest:



6

#### Details

- The number "0" has a special code so that the implicit 1 does not get added: the code is all 0s
   (it may seem that this takes up the representation for 1.0, but given how the exponent is represented, that's not the case)
   (see discussion of denorms in the textbook)
- The largest exponent value (with zero fraction) represents +/- infinity
- The largest exponent value (with non-zero fraction) represents
   NaN (not a number) for the result of 0/0 or (infinity minus infinity)
- Note that these choices impact the smallest and largest numbers that can be represented

#### **Exponent Representation**

- To simplify sort, sign was placed as the first bit
- For a similar reason, the representation of the exponent is also modified: in order to use integer compares, it would be preferable to have the smallest exponent as 00...0 and the largest exponent as 11...1
- This is the biased notation, where a bias is subtracted from the exponent field to yield the true exponent
- IEEE 754 single-precision uses a bias of 127 (since the exponent must have values between -127 and 128)...double precision uses a bias of 1023

Final representation: (-1)<sup>s</sup> x (1 + Fraction) x 2<sup>(Exponent - Bias)</sup>



Same rules as above, but the sign bit is 1 Same magnitudes as above, but negative numbers

## **Denormal Numbers**

• Exponent =  $000...0 \Rightarrow$  hidden bit is 0

$$x = (-1)^{S} \times (0 + Fraction) \times 2^{-Bias}$$

- Smaller than normal numbers
  - allow for gradual underflow, with diminishing precision
- Denormal with fraction = 000...0

$$x = (-1)^{S} \times (0+0) \times 2^{-Bias} = \pm 0.0$$
Two representations
of 0.0!

# Infinities and NaNs

- Exponent = 111...1, Fraction = 000...0
  - ±Infinity
  - Can be used in subsequent calculations, avoiding need for overflow check
- Exponent = 111...1, Fraction ≠ 000...0
  - Not-a-Number (NaN)
  - Indicates illegal or undefined result
     e.g., 0.0 / 0.0
  - Can be used in subsequent calculations

# **Floating-Point Example**

- Represent –0.75
  - $-0.75 = (-1)^{1} \times 1.1_{2} \times 2^{-1}$
  - S = 1
  - Fraction =  $1000...00_2$
  - Exponent = -1 + Bias
    - Single:  $-1 + 127 = 126 = 01111110_2$
    - Double:  $-1 + 1023 = 1022 = 0111111110_2$
- Single: 10111110100...00
- Double: 1011111110100...00

# Floating-Point Example

- What number is represented by the singleprecision float
  - 1100000101000...00
    - S = 1
    - Fraction =  $01000...00_2$
    - Fxponent =  $1000001_2 = 129$

• 
$$x = (-1)^{1} \times (1 + 01_{2}) \times 2^{(129 - 127)}$$
  
=  $(-1) \times 1.25 \times 2^{2}$   
=  $-5.0$ 

# **Floating-Point Addition**

- Consider a 4-digit decimal example
  - $9.999 \times 10^{1} + 1.610 \times 10^{-1}$
- 1. Align decimal points
  - Shift number with smaller exponent
  - $9.999 \times 10^{1} + 0.016 \times 10^{1}$
- 2. Add significands
  - $9.999 \times 10^{1} + 0.016 \times 10^{1} = 10.015 \times 10^{1}$
- 3. Normalize result & check for over/underflow
  - $1.0015 \times 10^{2}$
- 4. Round and renormalize if necessary
  - $1.002 \times 10^{2}$

# **Floating-Point Addition**

- Now consider a 4-digit binary example
  - $1.000_2 \times 2^{-1} + -1.110_2 \times 2^{-2} (0.5 + -0.4375)$
- 1. Align binary points
  - Shift number with smaller exponent
  - $1.000_2 \times 2^{-1} + -0.111_2 \times 2^{-1}$
- 2. Add significands
  - $1.000_2 \times 2^{-1} + -0.111_2 \times 2^{-1} = 0.001_2 \times 2^{-1}$
- 3. Normalize result & check for over/underflow
  - $1.000_2 \times 2^{-4}$ , with no over/underflow
- 4. Round and renormalize if necessary
  - $1.000_2 \times 2^{-4}$  (no change) = 0.0625

# FP Adder Hardware

- Much more complex than integer adder
- Doing it in one clock cycle would take too long
  - Much longer than integer operations
  - Slower clock would penalize all instructions
- FP adder usually takes several cycles
  - Can be pipelined

## **FP Adder Hardware**



#### Chapter 3 — Arithmetic for Computers — 17

# **Floating-Point Multiplication**

- Consider a 4-digit decimal example
  - $1.110 \times 10^{10} \times 9.200 \times 10^{-5}$
- 1. Add exponents
  - For biased exponents, subtract bias from sum
  - New exponent = 10 + -5 = 5
- 2. Multiply significands
  - $1.110 \times 9.200 = 10.212 \implies 10.212 \times 10^{5}$
- 3. Normalize result & check for over/underflow
  - $1.0212 \times 10^{6}$
- 4. Round and renormalize if necessary
  - 1.021 × 10<sup>6</sup>
- 5. Determine sign of result from signs of operands
  - +1.021 × 10<sup>6</sup>

# **Floating-Point Multiplication**

- Now consider a 4-digit binary example
  - $1.000_2 \times 2^{-1} \times -1.110_2 \times 2^{-2} (0.5 \times -0.4375)$
- 1. Add exponents
  - Unbiased: -1 + -2 = -3
  - Biased: (-1 + 127) + (-2 + 127) = -3 + 254 127 = -3 + 127
- 2. Multiply significands
  - $1.000_2 \times 1.110_2 = 1.1102 \implies 1.110_2 \times 2^{-3}$
- 3. Normalize result & check for over/underflow
  - $1.110_2 \times 2^{-3}$  (no change) with no over/underflow
- 4. Round and renormalize if necessary
  - 1.110<sub>2</sub> × 2<sup>-3</sup> (no change)
- 5. Determine sign: +ve × –ve  $\Rightarrow$  –ve
  - $-1.110_2 \times 2^{-3} = -0.21875$

# **FP** Arithmetic Hardware

- FP multiplier is of similar complexity to FP adder
  - But uses a multiplier for significands instead of an adder
- FP arithmetic hardware usually does
  - Addition, subtraction, multiplication, division, reciprocal, square-root
  - $FP \leftrightarrow$  integer conversion
- Operations usually takes several cycles
  - Can be pipelined

# **FP** Instructions in MIPS

- FP hardware is coprocessor 1
  - Adjunct processor that extends the ISA
- Separate FP registers
  - 32 single-precision: \$f0, \$f1, ... \$f31
  - Paired for double-precision: \$f0/\$f1, \$f2/\$f3, ...
    Release 2 of MIPs ISA supports 32 × 64-bit FP reg's
- FP instructions operate only on FP registers
  - Programs generally don't do integer ops on FP data, or vice versa
  - More registers with minimal code-size impact
- FP load and store instructions
  - lwc1, ldc1, swc1, sdc1
    - e.g., ldc1 \$f8, 32(\$sp)

# **FP** Instructions in MIPS

- Single-precision arithmetic
  - add.s, sub.s, mul.s, div.s
     e.g., add.s \$f0, \$f1, \$f6
- Double-precision arithmetic
  - add.d, sub.d, mul.d, div.d
     e.g., mul.d \$f4, \$f4, \$f6
- Single- and double-precision comparison
  - c.*xx*.s, c.*xx*.d (*xx* is eq, lt, le, ...)
  - Sets or clears FP condition-code bit
     e.g. c.lt.s \$f3, \$f4
- Branch on FP condition code true or false
  - bclt, bclf
    - e.g., bclt TargetLabel

# x86 FP Architecture

- Originally based on 8087 FP coprocessor
  - 8 × 80-bit extended-precision registers
  - Used as a push-down stack
  - Registers indexed from TOS: ST(0), ST(1), ...
- FP values are 32-bit or 64 in memory
  - Converted on load/store of memory operand
  - Integer operands can also be converted on load/store
- Very difficult to generate and optimize code
  - Result: poor FP performance

# FP Example: °F to °C

• C code:

```
float f2c (float fahr) {
    return ((5.0/9.0)*(fahr - 32.0));
}
```

- fahr in \$f12, result in \$f0, literals in global memory space
- Compiled MIPS code:

```
f2c: lwc1 $f16, const5($gp)
    lwc2 $f18, const9($gp)
    div.s $f16, $f16, $f18
    lwc1 $f18, const32($gp)
    sub.s $f18, $f12, $f18
    mul.s $f0, $f16, $f18
    jr $ra
```



- FP operations are much slower than integer ops
- Fixed point arithmetic uses integers, but assumes that every number is multiplied by the same factor
- Example: with a factor of 1/1000, the fixed-point representations for 1.46, 1.7198, and 5624 are respectively 1460, 1720, and 5624000
- More programming effort and possibly lower precision for higher performance

- ALUs are typically designed to perform 64-bit or 128-bit arithmetic
- Some data types are much smaller, e.g., bytes for pixel RGB values, half-words for audio samples
- Partitioning the carry-chains within the ALU can convert the 64-bit adder into 4 16-bit adders or 8 8-bit adders
- A single load can fetch multiple values, and a single add instruction can perform multiple parallel additions, referred to as subword parallelism

#### Thank you!

## FP Example: Array Multiplication

- $X = X + Y \times Z$ 
  - All 32 × 32 matrices, 64-bit double-precision elements
- C code:

  - }
    - Addresses of x, y, z in \$a0, \$a1, \$a2, and i, j, k in \$s0, \$s1, \$s2

## FP Example: Array Multiplication

MIPS code:

	li	\$t1,	32		#	\$t1	= 32 (row size/loop end)
	li	\$s0,	0		#	i =	0; initialize 1st for loop
L1:	li	\$s1,	0		#	j =	0; restart 2nd for loop
L2:	li	\$s2,	0		#	k =	0; restart 3rd for loop
	sll	\$t2,	\$s0,	5	#	\$t2	= i * 32 (size of row of x)
	addu	\$t2,	\$t2,	\$s1	#	\$t2	= i * size(row) + j
	sll	\$t2,	\$t2,	3	#	\$t2	<pre>= byte offset of [i][j]</pre>
	addu	\$t2,	\$a0,	\$t2	#	\$t2	<pre>= byte address of x[i][j]</pre>
	l.d	\$f4,	0(\$t2	2)	#	\$f4	= 8 bytes of x[i][j]
L3:	sll	\$t0,	\$s2,	5	#	\$t0	= k * 32 (size of row of z)
	addu	\$t0,	\$t0,	\$s1	#	\$t0	= k * size(row) + j
	sll	\$t0,	\$t0,	3	#	\$t0	<pre>= byte offset of [k][j]</pre>
	addu	\$t0,	\$a2,	\$t0	#	\$t0	<pre>= byte address of z[k][j]</pre>
	l.d	\$f16	, 0(\$t	:0)	#	\$f16	5 = 8 bytes of $z[k][j]$

...

### FP Example: Array Multiplication

# $t0 = i*32$ (size of row of y)
# \$t0 = i*size(row) + k
# \$t0 = byte offset of [i][k]
<pre># \$t0 = byte address of y[i][k]</pre>
# \$f18 = 8 bytes of y[i][k]
6 # \$f16 = y[i][k] * z[k][j]
# f4=x[i][j] + y[i][k]*z[k][j]
# \$k k + 1
# if (k != 32) go to L3
# x[i][j] = \$f4
# \$j = j + 1
# if (j != 32) go to L2
# \$i = i + 1
# if (i != 32) go to L1
5