

**Lecture 16 – Floating Point II**

2004-10-06



Lecturer PSOE Dan Garcia

www.cs.berkeley.edu/~ddgarcia

VP debate... →

Cheney & Edwards

took off their gloves and both came out strong. Will you just answer the question, please? Both are seasoned orators...



cs61c.com/2004/ALLPOLITICS/10/05/debate.mai.n/  
 CS 61C L16 - Floating Point II (1)

Garcia, Fall 2004 © UCB

**Example: Representing 1/3 in MIPS**

- 1/3
  - = 0.33333...<sub>10</sub>
  - = 0.25 + 0.0625 + 0.015625 + 0.00390625 + ...
  - = 1/4 + 1/16 + 1/64 + 1/256 + ...
  - = 2<sup>-2</sup> + 2<sup>-4</sup> + 2<sup>-6</sup> + 2<sup>-8</sup> + ...
  - = 0.0101010101... 2<sup>-2</sup> 2<sup>0</sup>
  - = 1.0101010101... 2<sup>-2</sup> 2<sup>-2</sup>
- Sign: 0
- Exponent = -2 + 127 = 125 = 01111101
- Significand = 0101010101...



0 0111 1101 0101 0101 0101 0101 0101 0101 010

CS 61C L16 - Floating Point II (3)

Garcia, Fall 2004 © UCB

**Representation for ± ∞**

- In FP, divide by 0 should produce ± ∞, not overflow.
- Why?
  - OK to do further computations with ∞
  - E.g., X/0 > Y may be a valid comparison
  - Ask math majors
- IEEE 754 represents ± ∞
  - Most positive exponent reserved for ∞
  - Significands all zeroes



CS 61C L16 - Floating Point II (4)

Garcia, Fall 2004 © UCB

**Representation for 0**

- Represent 0?
  - exponent all zeroes
  - significand all zeroes too
  - What about sign?
    - +0: 0 00000000 000000000000000000000000
    - -0: 1 00000000 000000000000000000000000
- Why two zeroes?
  - Helps in some limit comparisons
  - Ask math majors



CS 61C L16 - Floating Point II (5)

Garcia, Fall 2004 © UCB

**Special Numbers**

- What have we defined so far? (Single Precision)
- | Exponent | Significand | Object        |
|----------|-------------|---------------|
| 0        | 0           | 0             |
| 0        | nonzero     | ???           |
| 1-254    | anything    | +/- fl. pt. # |
| 255      | 0           | +/- ∞         |
| 255      | nonzero     | ???           |

• Professor Kahan had clever ideas; "Waste not, want not"

• Exp=0,255 & Sig!=0 ...



CS 61C L16 - Floating Point II (6)

Garcia, Fall 2004 © UCB

**Representation for Not a Number**

- What is sqrt(-4.0) or 0/0?
  - If ∞ not an error, these shouldn't be either.
  - Called **Not a Number (NaN)**
  - Exponent = 255, Significand nonzero
- Why is this useful?
  - Hope NaNs help with debugging?
  - They contaminate: op(NaN, X) = NaN



CS 61C L16 - Floating Point II (7)

Garcia, Fall 2004 © UCB

### Representation for Denorms (1/2)

• **Problem:** There's a gap among representable FP numbers around 0

• Smallest representable pos num:

$$a = 1.0 \dots_2 * 2^{-126} = 2^{-126}$$

• Second smallest representable pos num:

$$b = 1.000 \dots_2 * 2^{-126} = 2^{-126} + 2^{-149}$$

$$a - 0 = 2^{-126}$$

$$b - a = 2^{-149}$$

Normalization and implicit 1 is to blame!



RQ answer!



CS 61C L16 - Floating Point II (8)

Garcia, Fall 2004 © UCB

### Representation for Denorms (2/2)

• **Solution:**

• We still haven't used Exponent = 0, Significant nonzero

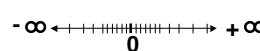
• Denormalized number: no leading 1, **implicit exponent = -126.**

• Smallest representable pos num:

$$a = 2^{-149}$$

• Second smallest representable pos num:

$$b = 2^{-148}$$



CS 61C L16 - Floating Point II (9)

Garcia, Fall 2004 © UCB

### Rounding

• Math on real numbers  $\Rightarrow$  we worry about rounding to fit result in the significant field.

RQ answer!

• FP hardware carries 2 extra bits of precision, and rounds for proper value

• Rounding occurs when converting...

- double to single precision
- floating point # to an integer



CS 61C L16 - Floating Point II (10)

Garcia, Fall 2004 © UCB

### IEEE Four Rounding Modes

• Round towards  $+\infty$

• ALWAYS round "up":  $2.1 \Rightarrow 3, -2.1 \Rightarrow -2$

• Round towards  $-\infty$

• ALWAYS round "down":  $1.9 \Rightarrow 1, -1.9 \Rightarrow -2$

• Truncate

• Just drop the last bits (round towards 0)

• Round to (nearest) even (default)

• Normal rounding, almost:  $2.5 \Rightarrow 2, 3.5 \Rightarrow 4$

• Like you learned in grade school

• Insures fairness on calculation

• Half the time we round up, other half down



CS 61C L16 - Floating Point II (11)

Garcia, Fall 2004 © UCB

### Integer Multiplication (1/3)

• Paper and pencil example (unsigned):

$$\begin{array}{r} \text{Multiplicand} \quad 1000 \quad 8 \\ \text{Multiplier} \quad \times 1001 \quad 9 \\ \hline \phantom{0000} 1000 \\ \phantom{0000} 0000 \\ \phantom{0000} 0000 \\ \phantom{0000} +1000 \\ \hline 01001000 \end{array}$$

• m bits x n bits = m + n bit product



CS 61C L16 - Floating Point II (13)

Garcia, Fall 2004 © UCB

### Integer Multiplication (2/3)

• In MIPS, we multiply registers, so:

• 32-bit value x 32-bit value = 64-bit value

• Syntax of Multiplication (signed):

• `mult register1, register2`

• Multiplies 32-bit values in those registers & puts 64-bit product in special result regs:

- puts product **upper half in hi**, **lower half in lo**

• **hi** and **lo** are 2 registers separate from the 32 general purpose registers

• Use **mfhi** register & **mflo** register to move from **hi**, **lo** to another register



CS 61C L16 - Floating Point II (14)

Garcia, Fall 2004 © UCB

### Integer Multiplication (3/3)

- Example:

- in C:  $a = b * c$ ;

- in MIPS:

- let b be \$s2; let c be \$s3; and let a be \$s0 and \$s1 (since it may be up to 64 bits)

```
mult $s2, $s3 # b*c
mfhi $s0      # upper half of
              # product into $s0
mflo $s1      # lower half of
              # product into $s1
```

- Note: Often, we only care about the lower half of the product.



### Integer Division (1/2)

- Paper and pencil example (unsigned):

```

          1001  Quotient
Divisor 1000 | 1001010  Dividend
          -1000
            10
            101
            1010
            -1000
              10  Remainder
              (or Modulo result)
```

- Dividend = Quotient x Divisor + Remainder



### Integer Division (2/2)

- Syntax of Division (signed):

- `div register1, register2`
- Divides 32-bit register 1 by 32-bit register 2:
- puts remainder of division in `hi`, quotient in `lo`

- Implements C division (/) and modulo (%)

- Example in C:  $a = c / d$ ;  
 $b = c \% d$ ;

- in MIPS:  $a \leftrightarrow \$s0$ ;  $b \leftrightarrow \$s1$ ;  $c \leftrightarrow \$s2$ ;  $d \leftrightarrow \$s3$

```
div $s2, $s3 # lo=c/d, hi=c%d
mflo $s0     # get quotient
mfhi $s1     # get remainder
```



### Unsigned Instructions & Overflow

- MIPS also has versions of `mult`, `div` for unsigned operands:

```
multu
```

```
divu
```

- Determines whether or not the product and quotient are changed if the operands are signed or unsigned.

- MIPS does not check overflow on ANY signed/unsigned multiply, divide instr

- Up to the software to check `hi`



### FP Addition & Subtraction

- Much more difficult than with integers (can't just add significands)

- How do we do it?

- De-normalize to match larger exponent
- Add significands to get resulting one
- Normalize (& check for under/overflow)
- Round if needed (may need to renormalize)

- If signs  $\neq$ , do a subtract. (Subtract similar)

- If signs  $\neq$  for add (or = for sub), what's ans sign?

- Question: How do we integrate this into the integer arithmetic unit? [Answer: We don't!]



### MIPS Floating Point Architecture (1/4)

- Separate floating point instructions:

- Single Precision:

```
add.s, sub.s, mul.s, div.s
```

- Double Precision:

```
add.d, sub.d, mul.d, div.d
```

- These are far more complicated than their integer counterparts

- Can take much longer to execute



### MIPS Floating Point Architecture (2/4)

- **Problems:**
  - Inefficient to have different instructions take vastly differing amounts of time.
  - Generally, a particular piece of data will not change FP  $\leftrightarrow$  int within a program.
    - Only 1 type of instruction will be used on it.
  - Some programs do no FP calculations
  - It takes lots of hardware relative to integers to do FP fast



CS 61C L16 - Floating Point II (21)

Garcia, Fall 2004 © UCB

### MIPS Floating Point Architecture (3/4)

- **1990 Solution: Make a completely separate chip that handles only FP.**
- **Coprocessor 1: FP chip**
  - contains 32 32-bit registers: \$f0, \$f1, ...
  - most of the registers specified in .s and .d instruction refer to this set
  - separate load and store: lwc1 and swc1 (“load word coprocessor 1”, “store ...”)
  - Double Precision: by convention, **even/odd pair** contain one DP FP number: \$f0/\$f1, \$f2/\$f3, ..., \$f30/\$f31
    - **Even register** is the name



CS 61C L16 - Floating Point II (22)

Garcia, Fall 2004 © UCB

### MIPS Floating Point Architecture (4/4)

- **1990 Computer actually contains multiple separate chips:**
  - Processor: handles all the normal stuff
  - Coprocessor 1: handles FP and only FP;
  - more coprocessors?... Yes, later
  - Today, FP coprocessor integrated with CPU, or cheap chips may leave out FP HW
- **Instructions to move data between main processor and coprocessors:**
  - mfc0, mtc0, mfc1, mtc1, etc.
- Appendix contains many more FP ops



CS 61C L16 - Floating Point II (23)

Garcia, Fall 2004 © UCB

### Peer Instruction

1. Converting float  $\rightarrow$  int  $\rightarrow$  float produces **same float number**
2. Converting int  $\rightarrow$  float  $\rightarrow$  int produces **same int number**
3. FP **add** is associative:  
(x+y)+z = x+(y+z)

	ABC
1:	FFF
2:	FFT
3:	FTF
4:	FTT
5:	TFF
6:	TFT
7:	TFB
8:	TTF



CS 61C L16 - Floating Point II (24)

Garcia, Fall 2004 © UCB

### As Promised, the way to remember #s

- What is  $2^{34}$ ? How many bits addresses (i.e., what's  $\text{ceil } \log_2 = \lg$  of) 2.5 TB?

- Answer!  $2^{XY}$  means...

X=0 $\Rightarrow$ 0	Y=0 $\Rightarrow$ 1
X=1 $\Rightarrow$ Kilo $\sim 10^3$	Y=1 $\Rightarrow$ 2
X=2 $\Rightarrow$ Mega $\sim 10^6$	Y=2 $\Rightarrow$ 4
X=3 $\Rightarrow$ Giga $\sim 10^9$	Y=3 $\Rightarrow$ 8
X=4 $\Rightarrow$ Tera $\sim 10^{12}$	Y=4 $\Rightarrow$ 16
X=5 $\Rightarrow$ Peta $\sim 10^{15}$	Y=5 $\Rightarrow$ 32
X=6 $\Rightarrow$ Exa $\sim 10^{18}$	Y=6 $\Rightarrow$ 64
X=7 $\Rightarrow$ Zetta $\sim 10^{21}$	Y=7 $\Rightarrow$ 128
X=8 $\Rightarrow$ Yotta $\sim 10^{24}$	Y=8 $\Rightarrow$ 256
	Y=9 $\Rightarrow$ 512



**MEMORIZE!**



CS 61C L16 - Floating Point II (26)

Garcia, Fall 2004 © UCB

### “And in conclusion...”

- Reserve exponents, significands:

Exponent	Significand	Object
0	0	0
0	<b>nonzero</b>	<b>Denorm</b>
1-254	anything	+/- fl. pt. #
255	0	+/- $\infty$
255	<b>nonzero</b>	<b>NaN</b>

- Integer mult, div uses hi, lo regs
  - mfhi and mflo copies out.
- Four rounding modes (to even default)
- MIPS FL ops complicated, expensive



CS 61C L16 - Floating Point II (27)

Garcia, Fall 2004 © UCB