

**Lecture 16
 Floating Point II**

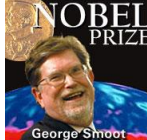
2006-10-04



Lecturer SOE Dan Garcia

www.cs.berkeley.edu/~ddgarcia

Prof Smoot given Nobel!! =>
 Prof George Smoot joins the ranks of the most distinguished faculty in the world (Cal has 7 active winners) when he was awarded the **Nobel Prize in Physics** for images that confirmed the Big Bang Theory.

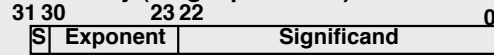


Review

Exponent tells Significant how much (2) to count by (... , 1/4, 1/2, 1, 2, ...)

- Floating Point lets us:
 - Represent numbers containing both integer and fractional parts; makes efficient use of available bits.
 - Store **approximate** values for very large and very small #s.
- **IEEE 754 Floating Point Standard** is most widely accepted attempt to standardize interpretation of such numbers (Every desktop or server computer sold since ~1997 follows these conventions)

• **Summary (single precision):**



$(-1)^S \times (1 + \text{Significand}) \times 2^{(\text{Exponent}-127)}$

- Double precision identical, except with exponent bias of 1023 (half, quad similar)

Precision and Accuracy

Don't confuse these two terms!

Precision is a count of the number bits in a computer word used to represent a value.

Accuracy is a measure of the difference between the actual value of a number and its computer representation.

High precision permits high accuracy but doesn't guarantee it. It is possible to have high precision but low accuracy.

Example: `float pi = 3.14;`
 pi will be represented using all 24 bits of the significant (highly precise), but is only an approximation (not accurate).

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

Representation for 0

- Represent 0?
 - exponent all zeroes
 - significand all zeroes
 - What about sign? Both cases valid.
- +0: 0 00000000 000000000000000000000000
 -0: 1 00000000 000000000000000000000000

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"
- We'll talk about Exp=0,255 & Sig!=0 later

Representation for Not a Number

- What do I get if I calculate $\text{sqrt}(-4.0)$ 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: $\text{op}(\text{NaN}, X) = \text{NaN}$

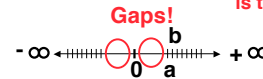


CS61C L16 Floating Point II (7)

Garcia, Fall 2006 © UCB

Representation for Denorms (1/2)

- Problem: There's a gap among representable FP numbers around 0
 - Smallest representable pos num: $a = 1.0..._2 * 2^{-126} = 2^{-126}$
 - Second smallest representable pos num: $b = 1.000...1_2 * 2^{-126} = 2^{-126} + 2^{-149}$
 - $a - 0 = 2^{-126}$
 - $b - a = 2^{-149}$



Normalization and implicit 1 is to blame!

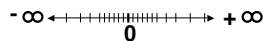


CS61C L16 Floating Point II (8)

Garcia, Fall 2006 © UCB

Representation for Denorms (2/2)

- Solution:
 - We still haven't used Exponent = 0, Significand nonzero
 - **Denormalized number**: no (implied) leading 1, **implicit exponent = -126**.
 - Smallest representable pos num: $a = 2^{-149}$
 - Second smallest representable pos num: $b = 2^{-148}$



CS61C L16 Floating Point II (9)

Garcia, Fall 2006 © UCB

Special Numbers Summary

- Reserve exponents, significands:
- | Exponent | Significand | Object |
|----------|----------------|---------------|
| 0 | 0 | 0 |
| 0 | nonzero | Denorm |
| 1-254 | anything | +/- fl. pt. # |
| 255 | 0 | +/- ∞ |
| 255 | nonzero | NaN |



CS61C L16 Floating Point II (10)

Garcia, Fall 2006 © UCB

Administrivia

- HW4 due tonight
- Project 2 was up on Monday, due next Fri
 - Pretend it is due on Wed so you have more time to study. We've made it easier than past project 2s...
- There are bugs on the Green sheet!
 - Check the course web page for details



CS61C L16 Floating Point II (11)

Garcia, Fall 2006 © UCB

Rounding

- When we perform math on real numbers, we have to worry about rounding to fit the result in the significant field.
- The FP hardware carries two extra bits of precision, and then round to get the proper value
- Rounding also occurs when converting:
 - double to a single precision value, or
 - floating point number to an integer



CS61C L16 Floating Point II (12)

Garcia, Fall 2006 © UCB

IEEE FP Rounding Modes

Examples in decimal (but, of course, IEEE754 in binary)

- Round towards $+\infty$
 - ALWAYS round "up": 2.001 \rightarrow 3, -2.001 \rightarrow -2
- Round towards $-\infty$
 - ALWAYS round "down": 1.999 \rightarrow 1, -1.999 \rightarrow -2
- Truncate
 - Just drop the last bits (round towards 0)
- Unbiased (default mode). Midway? Round to even
 - Normal rounding, almost: 2.4 \rightarrow 2, 2.6 \rightarrow 3, 2.5 \rightarrow 2, 3.5 \rightarrow 4
 - Round like you learned in grade school (nearest int)
 - Except if the value is right on the borderline, in which case we round to the nearest EVEN number
 - Insures fairness on calculation
 - This way, half the time we round up on tie, the other half time we round down. Tends to balance out inaccuracies



CS61C L16 Floating Point II (13)

Garcia, Fall 2006 © UCB

Casting floats to ints and vice versa

`(int) floating_point_expression`

Coerces and converts it to the nearest integer (C uses truncation)

```
i = (int) (3.14159 * f);
```

`(float) integer_expression`

converts integer to nearest floating point

```
f = f + (float) i;
```



CS61C L16 Floating Point II (14)

Garcia, Fall 2006 © UCB

FP Addition

- More difficult than with integers
- Can't just add significands
- How do we do it?
 - De-normalize to match exponents
 - Add significands to get resulting one
 - Keep the same exponent
 - Normalize (possibly changing exponent)
- Note: If signs differ, just perform a subtract instead.



CS61C L16 Floating Point II (15)

Garcia, Fall 2006 © UCB

MIPS Floating Point Architecture (1/4)

- MIPS has special instructions for floating point operations:
 - Single Precision:
`add.s, sub.s, mul.s, div.s`
 - Double Precision:
`add.d, sub.d, mul.d, div.d`
- These instructions are far more complicated than their integer counterparts. They require special hardware and usually they can take much longer to compute.



CS61C L16 Floating Point II (16)

Garcia, Fall 2006 © UCB

MIPS Floating Point Architecture (2/4)

- Problems:
 - It's inefficient to have different instructions take vastly differing amounts of time.
 - Generally, a particular piece of data will not change from FP to int, or vice versa, within a program. So only one type of instruction will be used on it.
 - Some programs do no floating point calculations
 - It takes lots of hardware relative to integers to do Floating Point fast



CS61C L16 Floating Point II (17)

Garcia, Fall 2006 © 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, $f1, ...`
 - most registers specified in `.s` and `.d` instruction refer to this set
 - separate load and store: `lwc1` and `swc1` ("load word coprocessor 1", "store ...")
 - Double Precision: by convention, even/odd pair contain one DP FP number: `$f0/$f1, $f2/$f3, ... , $f30/$f31`



CS61C L16 Floating Point II (18)

Garcia, Fall 2006 © 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, cheap chips may leave out FP HW
- Instructions to move data between main processor and coprocessors:
 - mfc0, mtc0, mfc1, mtc1, etc.
- Appendix pages A-70 to A-74 contain many, many more FP operations.



CS61C L16 Floating Point II (19)

Garcia, Fall 2006 © UCB

Peer Instruction

1 1000 0001 111 0000 0000 0000 0000 0000

What is the decimal equivalent of the floating pt # above?

- 1: -1.75
- 2: -3.5
- 3: -3.75
- 4: -7
- 5: -7.5
- 6: -15
- 7: $-7 * 2^{129}$
- 8: $-129 * 2^7$



CS61C L16 Floating Point II (20)

Garcia, Fall 2006 © 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: | FTF |
| 8: | TTT |



CS61C L16 Floating Point II (22)

Garcia, Fall 2006 © UCB

Peer Instruction

- Let $f(1, 2) = \#$ of floats between 1 and 2
- Let $f(2, 3) = \#$ of floats between 2 and 3

- 1: $f(1, 2) < f(2, 3)$
- 2: $f(1, 2) = f(2, 3)$
- 3: $f(1, 2) > f(2, 3)$



CS61C L16 Floating Point II (24)

Garcia, Fall 2006 © UCB

“And in conclusion...”

- Reserve exponents, significands:

| Exponent | Significand | Object |
|----------|-------------|---------------|
| 0 | 0 | 0 |
| 0 | nonzero | Denorm |
| 1-254 | anything | +/- fl. pt. # |
| 255 | 0 | +/- ∞ |
| 255 | nonzero | NaN |
- 4 rounding modes (default: unbiased)
- MIPS FL ops complicated, expensive



CS61C L16 Floating Point II (26)

Garcia, Fall 2006 © UCB

Example: Representing 1/3 in MIPS

- 1/3
 - = $0.33333..._{10}$
 - = $0.25 + 0.0625 + 0.015625 + 0.00390625 + \dots$
 - = $1/4 + 1/16 + 1/64 + 1/256 + \dots$
 - = $2^{-2} + 2^{-4} + 2^{-6} + 2^{-8} + \dots$
 - = $0.0101010101..._2 * 2^0$
 - = $1.0101010101..._2 * 2^{-2}$
- Sign: 0
- Exponent = $-2 + 127 = 125 = 01111101$
- Significand = $0101010101...$



0 0111 1101 0101 0101 0101 0101 0101 0101

CS61C L16 Floating Point II (27)

Garcia, Fall 2006 © UCB