

inst.eecs.berkeley.edu/~cs61c
UC Berkeley CS61C : Machine Structures



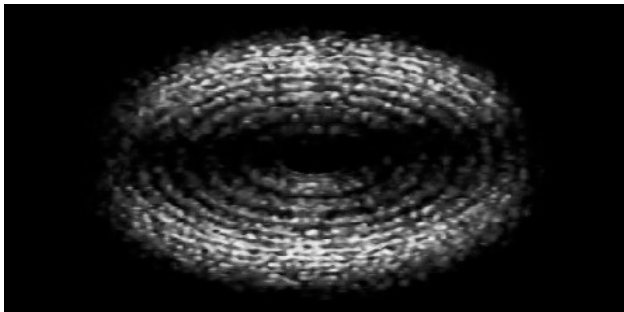
Lecture 16
Floating Point II

2008-02-29

TA Oddinaire Keaton Mowery

www.cs.berkeley.edu/~cs61c-tk

Electron Filmed for the First Time



“Scientists have filmed an electron in motion for the first time, using a new technique that will allow researchers to study the tiny particle’s movements directly.”

“Extremely short flashes of light are necessary to capture an electron in motion. A technology developed within the last few years can generate short pulses of intense laser light, called attosecond pulses, to get the job done.”

<http://www.livescience.com/strangenews/080225-electron-movie.html>



“Father” of the Floating point standard

IEEE Standard 754 for Binary Floating-Point Arithmetic.

1989 ACM Turing
Award Winner!



Prof. Kahan

[www.cs.berkeley.edu/~wkahan/
.../ieee754status/754story.html](http://www.cs.berkeley.edu/~wkahan/.../ieee754status/754story.html)



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 $\pm \infty$

- In FP, divide by 0 should produce $\pm \infty$, 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 $\pm \infty$
 - 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 000000000000000000000000000000

-0: 1 00000000 000000000000000000000000000000



Special Numbers

- **What have we defined so far? (Single Precision)**

Exponent	Significand	Object
0	0	0
0	<u>nonzero</u>	<u>???</u>
1-254	anything	+/- fl. pt. #
255	0	+/- ∞
255	<u>nonzero</u>	<u>???</u>

- **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 `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`



Representation for Denorms (1/2)

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

- **Smallest representable pos num:**

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

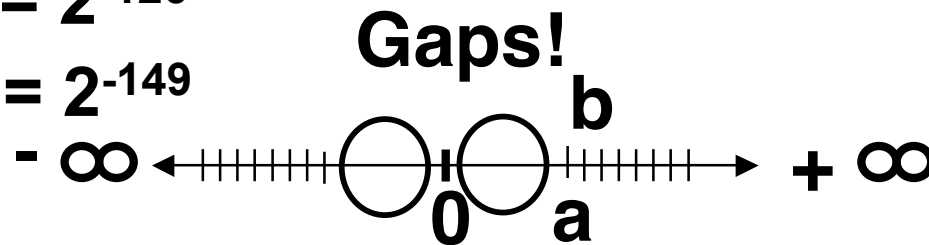
- **Second smallest representable pos num:**

$$\begin{aligned} b &= 1.000\dots1_2 * 2^{-126} \\ &= (1 + 0.00\dots1_2) * 2^{-126} \\ &= (1 + 2^{-23}) * 2^{-126} \\ &= 2^{-126} + 2^{-149} \end{aligned}$$

Normalization and implicit 1 is to blame!

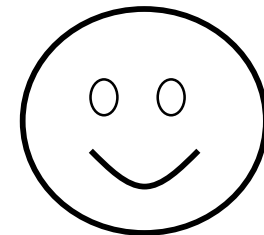
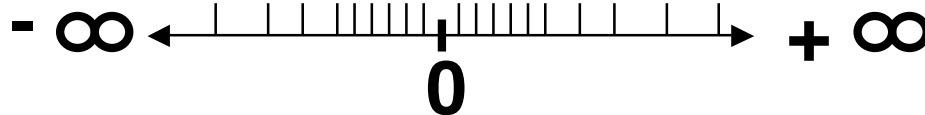
$$a - 0 = 2^{-126}$$

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



Representation for Denorm (2/2)

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



Special Numbers Summary

- Reserve exponents, significands:

Exponent	Significand	Object
0	0	0
0	<u>nonzero</u>	<u>Denorm</u>
1-254	Anything	+/- fl. Pt #
255	<u>0</u>	<u>+/- ∞</u>
255	<u>nonzero</u>	<u>NaN</u>



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**



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**



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: TTF
8: TTT



Peer Instruction Answer

1. **F A L S E** Converting float -> int float produces same float number

2. **F A L S E** Converting int -> float int produces same int number

3. **F A L S E** FP add is associative.
 $(x+y)+z = x+(y+z)$

1. 3.14 -> 3 -> 3

2. 32 bits for signed int,
but 24 for FP mantissa?

3. $x =$ biggest pos #,
 $y = -x$, $z = 1$ ($x \neq \text{inf}$)

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



Peer Instruction

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

$$1: f(1, 2) < f(2, 3)$$

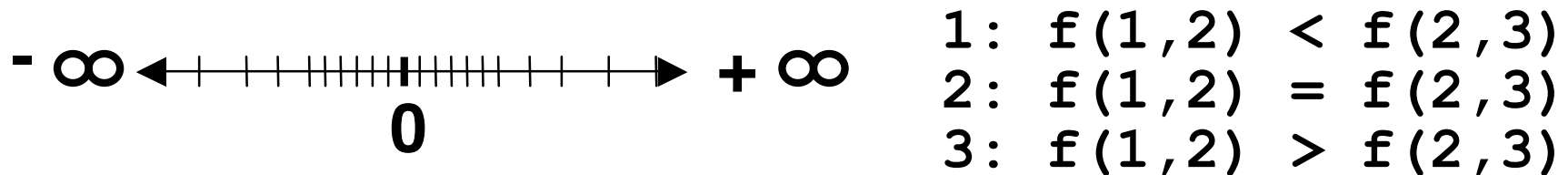
$$2: f(1, 2) = f(2, 3)$$

$$3: f(1, 2) > f(2, 3)$$



Peer Instruction Answer

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



“And in conclusion...”

- Reserve exponents, significands:

Exponent	Significand	Object
0	0	0
0	<u>nonzero</u>	<u>Denorm</u>
1-254	Anything	+/- fl. Pt #
255	<u>0</u>	<u>+/- ∞</u>
255	<u>nonzero</u>	<u>NaN</u>

- 4 Rounding modes (default: unbiased)
- MIPS FI ops complicated, expensive



Bonus slides

- **These are extra slides that used to be included in lecture notes, but have been moved to this, the “bonus” area to serve as a supplement.**
- **The slides will appear in the order they would have in the normal presentation**

Bonus



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.**



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.**



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**



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: $\$f0, \$f1, \dots$
 - 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: $\$f0/\$f1, \$f2/\$f3, \dots, \$f30/\$f31$



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.**



Example: Representing 1/3 in MIPS

- $1/3$
 - = $0.33333\dots_{10}$
 - = $0.25 + 0.0625 + 0.015625 + 0.00390625 + \dots$
 - = $1/4 + 1/16 + 1/64 + 1/256 + \dots$
 - = $2^{-2} + 2^{-4} + 2^{-6} + 2^{-8} + \dots$
 - = $0.01010101\dots_2 * 2^0$
 - = $1.01010101\dots_2 * 2^{-2}$
 - Sign: 0
 - Exponent = $-2 + 127 = 125 = 01111101$
 - Significand = $01010101\dots$

0	0111 1101	0101 0101 0101 0101 0101 010
---	-----------	------------------------------



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;
```



int → float → int

```
if (i == (int)((float) i)) {  
    printf("true");  
}
```

- **Will not always print “true”**
- **Most large values of integers don’t have exact floating point representations!**
- **What about double?**



float \rightarrow int \rightarrow float

```
if (f == (float)((int) f)) {  
    printf("true");  
}
```

- **Will not always print “true”**
- **Small floating point numbers (<1) don't have integer representations**
- **For other numbers, rounding errors**



Floating Point Fallacy

- **FP add associative: FALSE!**

- $x = -1.5 \times 10^{38}$, $y = 1.5 \times 10^{38}$, and $z = 1.0$

- $x + (y + z) = -1.5 \times 10^{38} + (1.5 \times 10^{38} + 1.0)$
 $= -1.5 \times 10^{38} + (1.5 \times 10^{38}) = \underline{0.0}$

- $(x + y) + z = (-1.5 \times 10^{38} + 1.5 \times 10^{38}) + 1.0$
 $= (0.0) + 1.0 = \underline{1.0}$

- **Therefore, Floating Point add is not associative!**

- **Why? FP result approximates real result!**

- **This example: 1.5×10^{38} is so much larger than 1.0 that $1.5 \times 10^{38} + 1.0$ in floating point representation is still 1.5×10^{38}**

