

## CS 61C: Great Ideas in Computer Architecture (Machine Structures) Performance and Floating-Point Arithmetic

Instructors:  
 Krste Asanovic & Vladimir Stojanovic  
<http://inst.eecs.berkeley.edu/~cs61c/>

### New-School Machine Structures (It's a bit more complicated!)

- Parallel Requests**  
Assigned to computer  
e.g., Search "Katz"
- Parallel Threads**  
Assigned to core  
e.g., Lookup, Ads
- Parallel Instructions**  
>1 instruction @ one time  
e.g., 5 pipelined instructions
- Parallel Data**  
>1 data item @ one time  
e.g., Add of 4 pairs of words
- Hardware descriptions**  
All gates @ one time
- Programming Languages**

**Software** | **Hardware**

Warehouse Scale Computer | Smart Phone

**Harness Parallelism & Achieve High Performance** | **How do we know?**

### What is Performance?

- Latency (or response time or execution time)**
  - Time to complete one task
- Bandwidth (or throughput)**
  - Tasks completed per unit time

### Cloud Performance: Why Application Latency Matters

Server Delay (ms)	Increased time to next click (ms)	Queries/user	Any clicks/user	User satisfaction	Revenue/User
50	--	--	--	--	--
200	500	--	-0.3%	-0.4%	--
500	1200	--	-1.0%	-0.9%	-1.2%
1000	1900	-0.7%	-1.9%	-1.6%	-2.8%
2000	3100	-1.8%	-4.4%	-3.8%	-4.3%

Figure 6.10 Negative impact of delays at Bing search server on user behavior [Brutlag and Schurman 2009].

- Key figure of merit: application responsiveness
  - Longer the delay, the fewer the user clicks, the less the user happiness, and the lower the revenue per user

### Defining CPU Performance

- What does it mean to say X is faster than Y?
- Ferrari vs. School Bus?
- 2013 Ferrari 599 GTB
  - 2 passengers, 11.1 secs in quarter mile
- 2013 Type D school bus
  - 54 passengers, quarter mile time?
- <http://www.youtube.com/watch?v=KwyCoQuhUNA>
- Response Time/Latency:** e.g., time to travel ¼ mile
- Throughput/Bandwidth:** e.g., passenger-mi in 1 hour

### Defining Relative CPU Performance

- $Performance_x = 1/Program\ Execution\ Time_x$
- $Performance_x > Performance_y \Rightarrow 1/Execution\ Time_x > 1/Execution\ Time_y \Rightarrow Execution\ Time_y > Execution\ Time_x$
- Computer X is N times faster than Computer Y  
 $Performance_x / Performance_y = N$  or  
 $Execution\ Time_y / Execution\ Time_x = N$
- Bus is to Ferrari as 12 is to 11.1:  
 Ferrari is 1.08 times faster than the bus!

### Measuring CPU Performance

- Computers use a clock to determine when events takes place within hardware
- **Clock cycles:** discrete time intervals
  - aka clocks, cycles, clock periods, clock ticks
- **Clock rate or clock frequency:** clock cycles per second (inverse of clock cycle time)
- 3 GigaHertz clock rate
  - => clock cycle time =  $1/(3 \times 10^9)$  seconds
  - clock cycle time = 333 picoseconds (ps)

### CPU Performance Factors

- To distinguish between processor time and I/O, **CPU time** is time spent in processor
- $\text{CPU Time/Program} = \text{Clock Cycles/Program} \times \text{Clock Cycle Time}$
- Or
  - $\text{CPU Time/Program} = \text{Clock Cycles/Program} \div \text{Clock Rate}$

### CPU Performance Factors

- But a program executes instructions
- $\text{CPU Time/Program} = \text{Clock Cycles/Program} \times \text{Clock Cycle Time}$ 
  - =  $\text{Instructions/Program} \times \text{Average Clock Cycles/Instruction} \times \text{Clock Cycle Time}$
- 1<sup>st</sup> term called **Instruction Count**
- 2<sup>nd</sup> term abbreviated **CPI** for average **Clock Cycles Per Instruction**
- 3<sup>rd</sup> term is 1 / Clock rate

### Restating Performance Equation

- $\text{Time} = \frac{\text{Seconds}}{\text{Program}}$ 
  - =  $\frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Clock cycles}}{\text{Instruction}} \times \frac{\text{Seconds}}{\text{Clock Cycle}}$

### What Affects Each Component? Instruction Count, CPI, Clock Rate


Hardware or software component?	Affects What?
Algorithm	
Programming Language	
Compiler	
Instruction Set Architecture	

### What Affects Each Component? Instruction Count, CPI, Clock Rate

Hardware or software component?	Affects What?
Algorithm	Instruction Count, CPI
Programming Language	Instruction Count, CPI
Compiler	Instruction Count, CPI
Instruction Set Architecture	Instruction Count, Clock Rate, CPI

**Clickers:**  
 Computer A clock cycle time 250 ps,  $CPI_A = 2$   
 Computer B clock cycle time 500 ps,  $CPI_B = 1.2$   
 Assume A and B have same instruction set  
 Which statement is true?

- A: Computer A is  $\approx 1.2$  times faster than B
- B: Computer A is  $\approx 4.0$  times faster than B
- C: Computer B is  $\approx 1.7$  times faster than A
- D: Computer B is  $\approx 3.4$  times faster than A



13

### Workload and Benchmark

- **Workload:** Set of programs run on a computer
  - Actual collection of applications run or made from real programs to approximate such a mix
  - Specifies programs, inputs, and relative frequencies
- **Benchmark:** Program selected for use in comparing computer performance
  - Benchmarks form a workload
  - Usually standardized so that many use them

14

### SPEC (System Performance Evaluation Cooperative)

- Computer Vendor cooperative for benchmarks, started in 1989
- SPEC CPU2006
  - 12 Integer Programs
  - 17 Floating-Point Programs
- Often turn into number where bigger is faster
- **SPECratio:** reference execution time on old reference computer divide by execution time on new computer to get an effective speed-up

15

### SPECINT2006 on AMD Barcelona

Description	Instruction Count (B)	CPI	Clock cycle time (ps)	Execution Time (s)	Reference Time (s)	SPECratio
Interpreted string processing	2,118	0.75	400	637	9,770	15.3
Block-sorting compression	2,389	0.85	400	817	9,650	11.8
GNU C compiler	1,050	1.72	400	724	8,050	11.1
Combinatorial optimization	336	10.0	400	1,345	9,120	6.8
Go game	1,658	1.09	400	721	10,490	14.6
Search gene sequence	2,783	0.80	400	890	9,330	10.5
Chess game	2,176	0.96	400	837	12,100	14.5
Quantum computer simulation	1,623	1.61	400	1,047	20,720	19.8
Video compression	3,102	0.80	400	993	22,130	22.3
Discrete event simulation library	587	2.94	400	690	6,250	9.1
Games/path finding	1,082	1.79	400	773	7,020	9.1
XML parsing	1,058	2.70	400	1,143	6,900	6.0

### Summarizing Performance ...

System	Rate (Task 1)	Rate (Task 2)
A	10	20
B	20	10

Clickers: Which system is faster?

**A: System A**  
**B: System B**  
**C: Same performance**  
**D: Unanswerable question!**

17

### ... Depends Who's Selling

System	Rate (Task 1)	Rate (Task 2)	Average
A	10	20	15
B	20	10	15

Average throughput

System	Rate (Task 1)	Rate (Task 2)	Average
A	0.50	2.00	1.25
B	1.00	1.00	1.00

Throughput relative to B

System	Rate (Task 1)	Rate (Task 2)	Average
A	1.00	1.00	1.00
B	2.00	0.50	1.25

Throughput relative to A

18

## Summarizing SPEC Performance

- Varies from 6x to 22x faster than reference computer
- Geometric mean of ratios:  
N-th root of product of N ratios
 
$$\sqrt[n]{\prod_{i=1}^n \text{Execution time ratio}_i}$$
  - Geometric Mean gives same relative answer no matter what computer is used as reference
- Geometric Mean for Barcelona is 11.7

19

## Administrivia

- Project 2, Part 2 due 3/22
- No assigned work over spring break
- Next assignment, HW5, due 04/05
- Midterm II is 04/09
  - Conflict? Email Sagar
  - DSP to receive email about accommodations soon

20

## Quote of the day

“95% of the folks out there are completely clueless about floating-point.”

James Gosling  
Sun Fellow  
Java Inventor  
1998-02-28



The same has been said about Java's FP design..

CS61C L15 Floating Point I (21)

Garcia, Fall 2014 © UC Berkeley

## Review of Numbers

- Computers are made to deal with numbers
- What can we represent in N bits?
  - 2<sup>N</sup> things, and no more! They could be...
    - Unsigned integers:  
0 to 2<sup>N</sup> - 1  
(for N=32, 2<sup>N</sup>-1 = 4,294,967,295)
    - Signed Integers (Two's Complement)  
-2<sup>(N-1)</sup> to 2<sup>(N-1)</sup> - 1  
(for N=32, 2<sup>(N-1)</sup> = 2,147,483,648)



CS61C L15 Floating Point I (22)

Garcia, Fall 2014 © UC Berkeley

## What about other numbers?

- Very large numbers? (seconds/millennium)  
⇒ 31,556,926,000<sub>10</sub> (3.1556926<sub>10</sub> × 10<sup>10</sup>)
- Very small numbers? (Bohr radius)  
⇒ 0.000000000529177<sub>10</sub>m (5.29177<sub>10</sub> × 10<sup>-11</sup>)
- Numbers with both integer & fractional parts?  
⇒ 1.5

First consider #3.

...our solution will also help with 1 and 2.



CS61C L15 Floating Point I (23)

Garcia, Fall 2014 © UC Berkeley

## Representation of Fractions

“Binary Point” like decimal point signifies boundary between integer and fractional parts:

Example 6-bit representation:  $xx.yyyy$

$$10.1010_2 = 1x2^1 + 1x2^0 + 1x2^{-1} + 1x2^{-3} = 2.625_{10}$$

If we assume “fixed binary point”, range of 6-bit representations with this format:

0 to 3.9375 (almost 4)



CS61C L15 Floating Point I (24)

Garcia, Fall 2014 © UC Berkeley

### Fractional Powers of 2

i	$2^{-i}$
0	1.0
1	0.5
2	0.25
3	0.125
4	0.0625
5	0.03125
6	0.015625
7	0.0078125
8	0.00390625
9	0.001953125
10	0.0009765625
11	0.00048828125
12	0.000244140625
13	0.0001220703125
14	0.00006103515625
15	0.000030517578125

*Cal* CS61C L15 Floating Point I (25) Garcia, Fall 2014 © UC

### Representation of Fractions with Fixed Pt.

What about addition and multiplication?

Addition is straightforward:

$$\begin{array}{r} 01.100 \\ + 00.100 \\ \hline 10.000 \end{array} \quad \begin{array}{r} 1.5_{10} \\ 0.5_{10} \\ \hline 2.0_{10} \end{array}$$

Multiplication a bit more complex:

$$\begin{array}{r} 01.100 \\ \times 00.100 \\ \hline 00000 \\ 01100 \\ \hline 00000 \\ 00000 \\ \hline 0000110000 \end{array}$$

Where's the answer, 0.11? (need to remember where point is)

*Cal* CS61C L15 Floating Point I (26) Garcia, Fall 2014 © UC

### Representation of Fractions

So far, in our examples we used a "fixed" binary point what we really want is to "float" the binary point. Why?

Floating binary point most effective use of our limited bits (and thus more accuracy in our number representation):

example: put 0.1640625 into binary. Represent as in 5-bits choosing where to put the binary point.  
 ... 000000.001010100000...

Store these bits and keep track of the binary point 2 places to the left of the MSB

Any other solution would lose accuracy!

With floating-point rep., each numeral carries an exponent field recording the whereabouts of its binary point.

The binary point can be outside the stored bits, so very large and small numbers can be represented.

*Cal* CS61C L15 Floating Point I (27) Garcia, Fall 2014 © UC

### Scientific Notation (in Decimal)

mantissa: 6.02<sub>10</sub> x 10<sup>23</sup> exponent

decimal point: 6.02<sub>10</sub> radix (base): 10<sup>23</sup>

- Normalized form: no leading 0s (exactly one digit to left of decimal point)
- Alternatives to representing 1/1,000,000,000
  - Normalized: 1.0 x 10<sup>-9</sup>
  - Not normalized: 0.1 x 10<sup>-8</sup>, 10.0 x 10<sup>-10</sup>

*Cal* CS61C L15 Floating Point I (28) Garcia, Fall 2014 © UC

### Scientific Notation (in Binary)

mantissa: 1.01<sub>two</sub> x 2<sup>-1</sup> exponent

"binary point": 1.01<sub>two</sub> radix (base): 2<sup>-1</sup>

- Computer arithmetic that supports it called **floating point**, because it represents numbers where the binary point is not fixed, as it is for integers
  - Declare such variable in C as `float`
    - Or `double` for double precision.

*Cal* CS61C L15 Floating Point I (29) Garcia, Fall 2014 © UC

### Floating-Point Representation (1/2)

- Normal format: +1.xxx...x<sub>two</sub> \* 2<sup>yy...y</sup><sub>two</sub>
- Multiple of Word Size (32 bits)

31 30      23 22      0

S Exponent      Significand

1 bit   8 bits      23 bits

- S represents Sign
- Exponent represents y's
- Significand represents x's

- Represent numbers as small as 2.0 x 10<sup>-38</sup> to as large as 2.0 x 10<sup>38</sup>

*Cal* CS61C L15 Floating Point I (30) Garcia, Fall 2014 © UC

### Floating-Point Representation (2/2)

- What if result too large? ( $> 2.0 \times 10^{38}$ ,  $< -2.0 \times 10^{38}$ )
  - **Overflow!**  $\Rightarrow$  Exponent larger than represented in 8-bit Exponent field
- What if result too small? ( $> 0$  &  $< 2.0 \times 10^{-38}$ ,  $< 0$  &  $> -2.0 \times 10^{-38}$ )
  - **Underflow!**  $\Rightarrow$  Negative exponent larger than represented in 8-bit Exponent field

- What would help reduce chances of overflow and/or underflow?

CS61C L15 Floating Point | (31) Garcia, Fall 2014 © UC Berkeley

### IEEE 754 Floating-Point Standard (1/3)

Single Precision (Double Precision similar):

31 30                      23 22                      0

S	Exponent	Significand
---	----------	-------------

1 bit    8 bits                      23 bits

- Sign bit: 1 means negative, 0 means positive
- Significand in *sign-magnitude* format (not 2's complement)
  - To pack more bits, leading 1 implicit for normalized numbers
  - 1 + 23 bits single, 1 + 52 bits double
  - always true:  $0 < \text{Significand} < 1$  (for normalized numbers)

**Note:** 0 has no leading 1, so reserve exponent value 0 just for number 0

CS61C L15 Floating Point | (32) Garcia, Fall 2014 © UC Berkeley

### IEEE 754 Floating Point Standard (2/3)

- IEEE 754 uses “biased exponent” representation.
  - Designers wanted FP numbers to be used even if no FP hardware; e.g., sort records with FP numbers using integer compares
  - Wanted bigger (integer) exponent field to represent bigger numbers.
  - 2's complement poses a problem (because negative numbers look bigger)
  - We're going to see that the numbers are ordered EXACTLY as in sign-magnitude
    - I.e., counting from binary odometer 00...00 up to 11...11 goes from 0 to +MAX to -0 to -MAX to 0

CS61C L15 Floating Point | (33) Garcia, Fall 2014 © UC Berkeley

### IEEE 754 Floating Point Standard (3/3)

- Called **Biased Notation**, where bias is number subtracted to get real number
  - IEEE 754 uses bias of 127 for single prec.
  - Subtract 127 from Exponent field to get actual value for exponent
  - 1023 is bias for double precision

• **Summary (single precision):**

31 30                      23 22                      0

S	Exponent	Significand
---	----------	-------------

1 bit    8 bits                      23 bits


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

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

CS61C L15 Floating Point | (34) Garcia, Fall 2014 © UC Berkeley

### “Father” of the Floating point standard

## IEEE Standard 754 for Binary Floating-Point Arithmetic.



Prof. Kahan

1989  
ACM Turing  
Award Winner!

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

CS61C L15 Floating Point | (35) Garcia, Fall 2014 © UC Berkeley

### Representation for $\pm \infty$


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

CS61C L15 Floating Point | (36) Garcia, Fall 2014 © UC Berkeley

### Representation for 0

- Represent 0?
  - exponent all zeroes
  - significand all zeroes
  - What about sign? Both cases valid.

+0: 0 00000000 00000000000000000000000000000000  
 -0: 1 00000000 00000000000000000000000000000000




CS61C L15 Floating Point I (37) Garcia, Fall 2014 © UCB

### 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"
- Wanted to use Exp=0,255 & Sig!=0



CS61C L15 Floating Point I (38) Garcia, Fall 2014 © UCB

### 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:  $\text{op}(\text{NaN}, X) = \text{NaN}$
  - Can use the significand to identify which!



CS61C L15 Floating Point I (39) Garcia, Fall 2014 © UCB

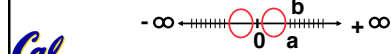

### Representation for Denorms (1/2)

- 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:  $b = 1.000 \dots 1_2 * 2^{-126} = (1 + 0.00 \dots 1_2) * 2^{-126} = (1 + 2^{-23}) * 2^{-126} = 2^{-126} + 2^{-149}$

**Normalization and implicit 1 is to blame!**

$a - 0 = 2^{-126}$   
 $b - a = 2^{-149}$

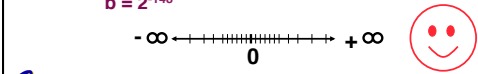

**Gaps!**

CS61C L15 Floating Point I (40) Garcia, Fall 2014 © 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 L15 Floating Point I (41) Garcia, Fall 2014 © UCB

### 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>	+/- ∞
255	<u>nonzero</u>	<u>NaN</u>



CS61C L15 Floating Point I (42) Garcia, Fall 2014 © UCB

www.n-comput.net/floatapplet/ieee754.html

### Conclusion

- **Floating Point lets us:** Exponent tells Significant how much (2) to count by (... , 1/4, 1/2, 1, 2, ...)
- 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)


Can store NaN, ±∞

• **Summary (single precision):**

31	30	23	22	0
S	Exponent		Significand	
1 bit	8 bits	23 bits		

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

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

 CSRIC L15 Floating Point | 43 Garcia, Fall 2014 © UCB

### And In Conclusion, ...

- Time (seconds/program) is measure of performance

$$= \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Clock cycles}}{\text{Instruction}} \times \frac{\text{Seconds}}{\text{Clock Cycle}}$$

- Floating-point representations hold approximations of real numbers in a finite number of bits

44