

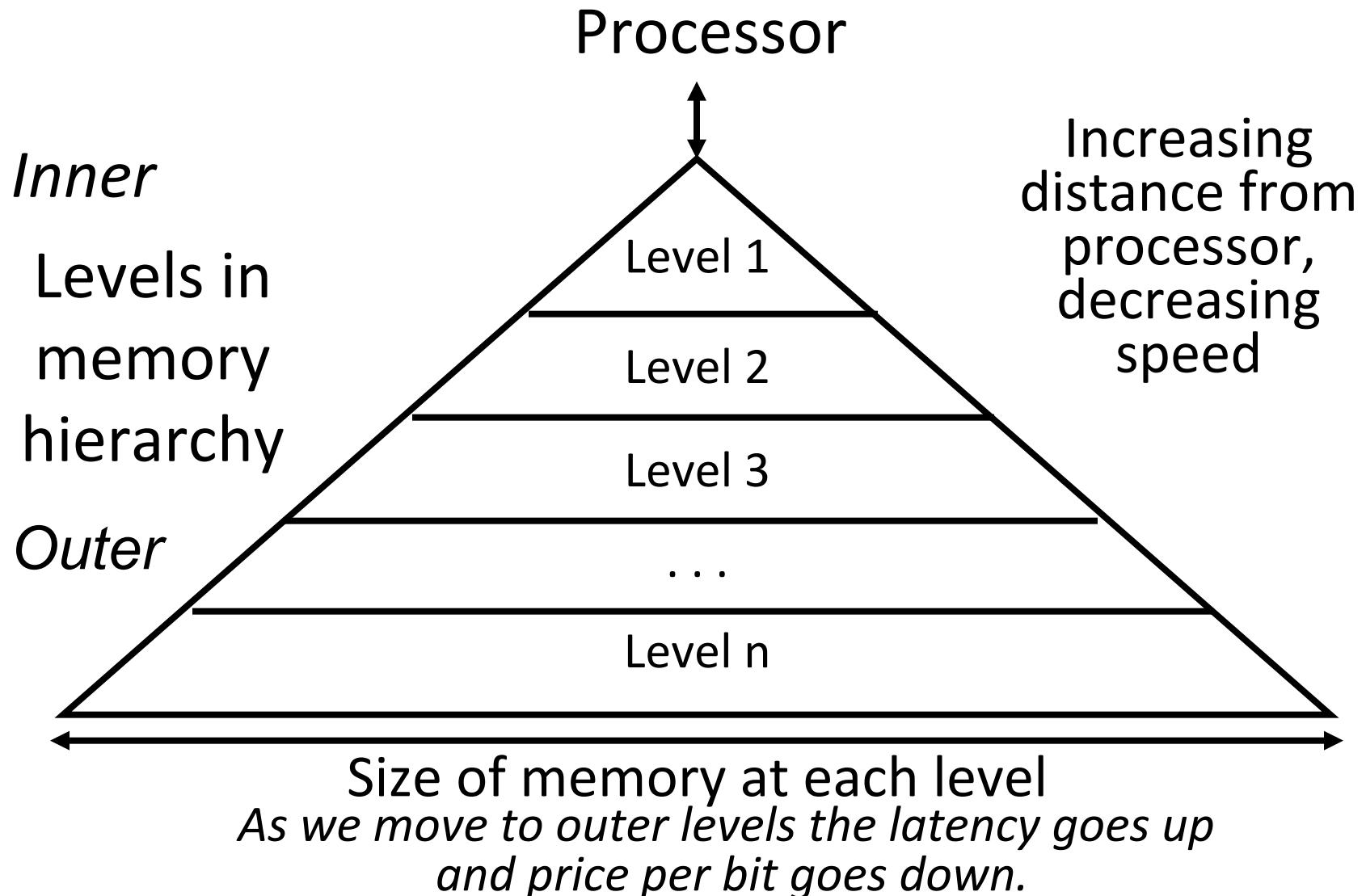
CS 61C: Great Ideas in Computer Architecture

Lecture 17: *Performance and Floating Point Arithmetic*

Instructor: Sagar Karandikar
sagark@eecs.berkeley.edu

<http://inst.eecs.berkeley.edu/~cs61c>

Review: Adding to the Memory Hierarchy



Review: Local vs. Global Miss Rates

- *Local miss rate* – the fraction of references to one level of a cache that miss

$$\text{Local Miss rate L2\$} = \frac{\$L2 \text{ Misses}}{L1\$ \text{ Misses}}$$

- *Global miss rate* – the fraction of references that miss in all levels of a multilevel cache

$$\text{L2\$ Global Miss rate} = \frac{L2\$ \text{ Misses}}{\text{Total Accesses}}$$

$$= \frac{L2\$ \text{ Misses}}{(L1\$ \text{ Misses} \times L1\$ \text{ Misses})} / \text{Total Accesses}$$

$$= \text{Local Miss rate L2\$} \times \text{Local Miss rate L1\$}$$

- AMAT = Time for a hit + Miss rate x Miss penalty
- AMAT = Time for a L1\\$ hit + (local) Miss rate L1\\$ x (Time for a L2\\$ hit + (local) Miss rate L2\\$ x L2\\$ Miss penalty)

Review: Analyzing Code Performance on a Cache

```
#define SZ 2048 // 2^11
int foo () {
    int a [SZ];
    int sum = 0;
    for (int i = 0; i < SZ; i++) {
        sum += a[i];
    }
    return sum;
}
```

1) Assume we only care about D\$

2) Assume array is “block aligned”
i.e. First byte of a[0] is the first byte in a block

3) Assume local vars live in registers

So what *do* we care about?
The sequence of accesses to array a

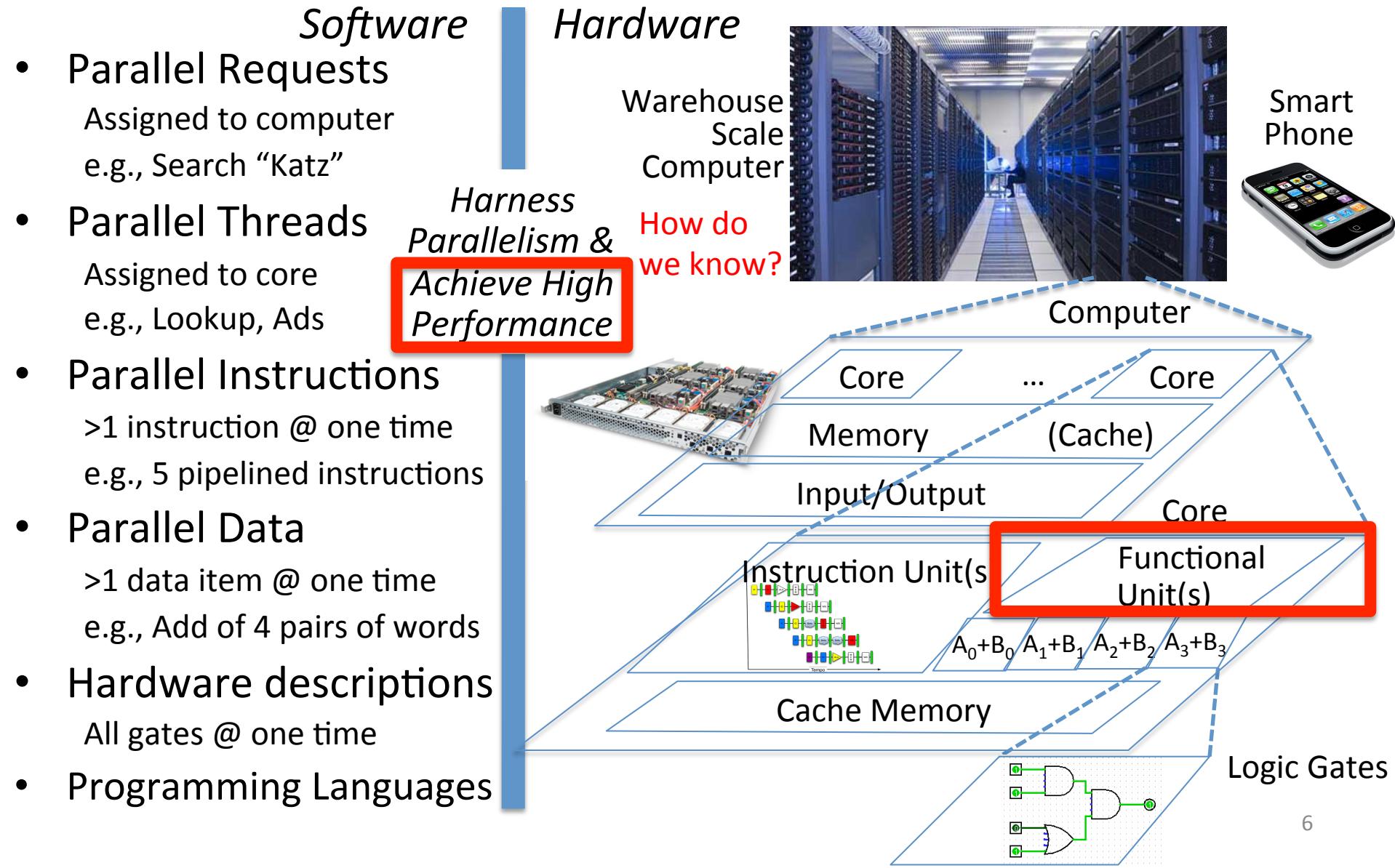
Review: Clicker Question

```
#define SZ 2048 // 2^11  
int foo() {  
    int a[SZ];  
    int sum = 0;  
    for (int i = 0; i < SZ; i+=256) sum += a[i];  
for (int i = 0; i < SZ; i+=256) sum += a[i];  
    return sum;  
}
```

Cache specs:
1 KiB Direct-Mapped
16 B Blocks

- What types of misses do we get in the first/second loops?
- A. Capacity/Capacity
 - B. Compulsory/Capacity
 - C. Compulsory/Conflict
 - D. Conflict/Capacity
 - E. Compulsory/Compulsory

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



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 satisfac- tion | 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?
 - Bugatti vs. School Bus?A photograph showing a dark blue and orange Bugatti Veyron 16.4 Super Sport on the left and a yellow Type D school bus on the right, parked side-by-side on a paved surface.
 - 2010 Bugatti Veyron 16.4 Super Sport
 - 2 passengers, 9.6 secs in quarter mile
 - 2013 Type D school bus
 - 54 passengers, quarter mile time?
- <https://youtu.be/KwyCoQuhUNA?t=4s>
- *Response Time/Latency*: e.g., time to travel $\frac{1}{4}$ mile
 - *Throughput/Bandwidth*: e.g., passenger-mi in 1 hour

Defining Relative CPU Performance

- $\text{Performance}_X = 1/\text{Program Execution Time}_X$
- $\text{Performance}_X > \text{Performance}_Y \Rightarrow$
 $1/\text{Execution Time}_X > 1/\text{Execution Time}_Y \Rightarrow$
 $\text{Execution Time}_Y > \text{Execution Time}_X$

Defining Relative CPU Performance

- Computer X is N times faster than Computer Y
 $\text{Performance}_X / \text{Performance}_Y = N$ or
 $\text{Execution Time}_Y / \text{Execution Time}_X = N$
- Bus is to Bugatti as 12 is to 9.6:
Bugatti is 1.25 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
 - CPU Time/Program
 - = Clock Cycles/Program
 - × Clock Cycle Time
 - Or
- CPU Time/Program
- $$= \text{Clock Cycles/Program} \div \text{Clock Rate}$$

CPU Performance Factors

- But a program executes instructions
- CPU Time/Program
 - = Clock Cycles/Program x Clock Cycle Time
 - = Instructions/Program
 - x Average Clock Cycles/Instruction
 - x Clock Cycle Time
- 1st term called *Instruction Count*
- 2nd term abbreviated *CPI* for average *Clock Cycles Per Instruction*
- 3rd term is 1 / Clock rate

Restating Performance Equation

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

The “Iron Law” of Processor Performance

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 ≈ 1.2 times faster than B
- B: Computer A is ≈ 4.0 times faster than B
- C: Computer B is ≈ 1.7 times faster than A
- D: Computer B is ≈ 3.4 times faster than A

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

SPEC

(System Performance Evaluation Cooperative)

- Computer Vendor cooperative for benchmarks, started in 1989
- SPECCPU2006
 - 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

SPECINT2006 on AMD Barcelona

| Description | Instruc- tion Count (B) | CPI | Clock cycle time (ps) | Execu- tion Time (s) | Refer- ence Time (s) | SPEC- ratio |
|-----------------------------------|-------------------------------|------|-----------------------------|----------------------------|----------------------------|----------------|
| 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 | 21 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!**

... 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

Summarizing SPEC Performance

- Varies from 6x to 22x faster than reference computer
- *Geometric mean* of ratios:
N-th root of product
of N ratios
 - Geometric Mean gives same relative answer no matter what computer is used as reference
- Geometric Mean for Barcelona is 11.7

$$\sqrt[n]{\prod_{i=1}^n \text{Execution time ratio}_i}$$

Administrivia

- Project 3-1 Out
 - Last week, we built a CPU together, this week, you start building your own!
- HW4 Out - Caches
- Guerrilla Section on Pipelining, Caches on Thursday, 5-7pm, Woz

Administrivia

- Midterm 2 is on 7/28 (one week from today)
 - In this room, at this time
 - Two double-sided 8.5"x11" handwritten cheatsheets
 - We'll provide a MIPS green sheet
 - No electronics
 - Covers up to and including **today's lecture (07/21)***
 - Review session is Friday, 7/24 from 1-4pm in HP Aud.

* This is one less lecture than initially indicated

Break

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..

Quote of the day

How Java's Floating-Point Hurts Everyone Everywhere

How Java's Floating-Point Hurts Everyone Everywhere

by

Prof. W. Kahan and Joseph D. Darcy
Elect. Eng. & Computer Science
Univ. of Calif. @ Berkeley

Originally presented 1 March 1998
at the invitation of the
ACM 1998 Workshop on Java for
High-Performance Network Computing
held at Stanford University

<http://www.cs.ucsb.edu/conferences/java98>

This document: <http://www.cs.berkeley.edu/~wkahan/JAVAhurt.pdf> or
<http://www.cs.berkeley.edu/~darcy/JAVAhurt.pdf>

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

Review of Numbers

- Computers are made to deal with numbers
- What can we represent in N bits?
 - 2^N things, and no more! They could be...
 - Unsigned integers:

0 to $2^N - 1$

(for N=32, $2^N - 1 = 4,294,967,295$)

- Signed Integers (Two's Complement)

$-2^{(N-1)}$ to $2^{(N-1)} - 1$

(for N=32, $2^{(N-1)} = 2,147,483,648$)

What about other numbers?

1. Very large numbers? (seconds/millennium)
⇒ $31,556,926,000_{10}$ ($3.1556926_{10} \times 10^{10}$)
2. Very small numbers? (Bohr radius)
⇒ $0.000000000529177_{10}\text{m}$ ($5.29177_{10} \times 10^{-11}$)
3. Numbers with both integer & fractional parts?
⇒ 1.5

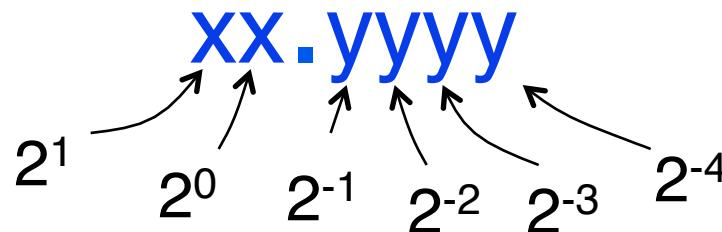
First consider #3.

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

Representation of Fractions

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

Example 6-bit representation:



$$10.1010_2 = 1 \times 2^1 + 1 \times 2^{-1} + 1 \times 2^{-3} = 2.625_{10}$$

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

0 to 3.9375 (almost 4)

Fractional Powers of 2

| i | 2^{-i} | |
|----|-------------------|--------|
| 0 | 1.0 | 1 |
| 1 | 0.5 | $1/2$ |
| 2 | 0.25 | $1/4$ |
| 3 | 0.125 | $1/8$ |
| 4 | 0.0625 | $1/16$ |
| 5 | 0.03125 | $1/32$ |
| 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 | |

Representation of Fractions with Fixed Pt.

What about addition and multiplication?

Addition is
straightforward:

$$\begin{array}{r} 01.100 \\ + 00.100 \\ \hline 10.000 \end{array} \quad \begin{array}{r} 1.5_{10} \\ 0.5_{10} \\ \hline 2.0_{10} \end{array}$$
$$\begin{array}{r} 01.100 \\ 00.100 \\ \hline 00.000 \end{array} \quad \begin{array}{r} 1.5_{10} \\ 0.5_{10} \\ \hline 0.5_{10} \end{array}$$

Multiplication a bit more complex:

$$\begin{array}{r} 00000 \\ 0110 \ 0 \\ 00000 \\ \hline 00000 \\ \hline 0000110000 \end{array}$$

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

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.

... 00000.001010100000...



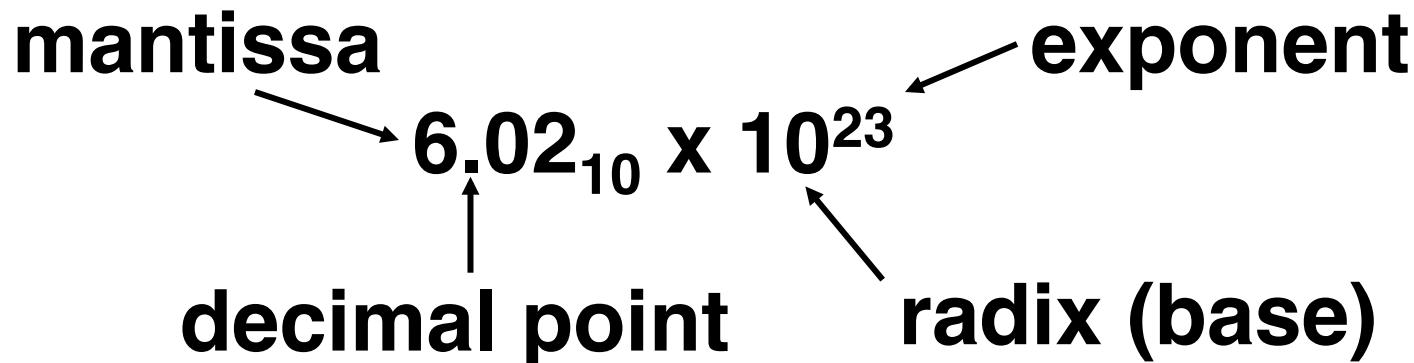
Store these bits and keep track of the binary point 2 places to the left of the MSB

Any other solution would lose accuracy!

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

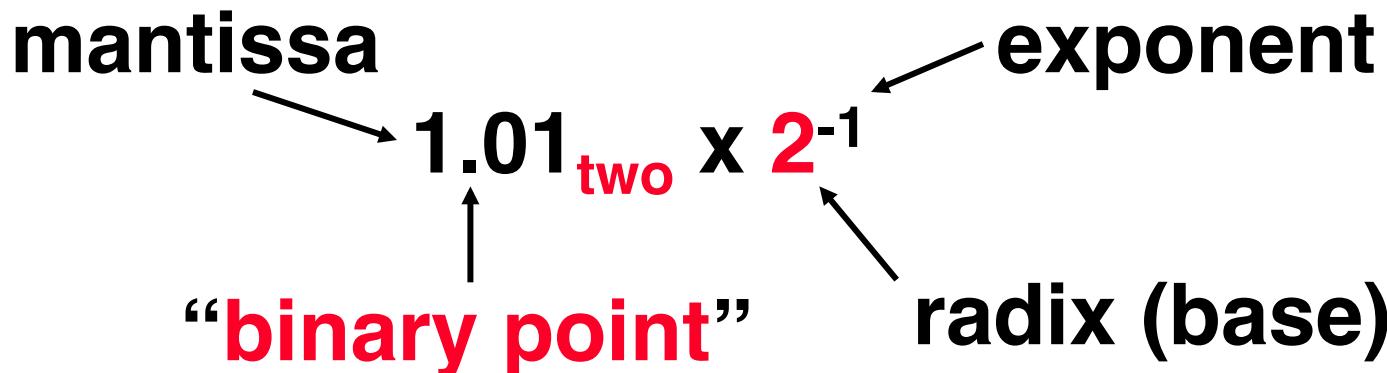
The binary point can be outside the stored bits, so very large and small numbers can be represented.

Scientific Notation (in Decimal)



- Normalized form: no leading 0s
(exactly one digit to left of decimal point)
- Alternatives to representing 1/1,000,000,000
 - Normalized: 1.0×10^{-9}
 - Not normalized: $0.1 \times 10^{-8}, 10.0 \times 10^{-10}$

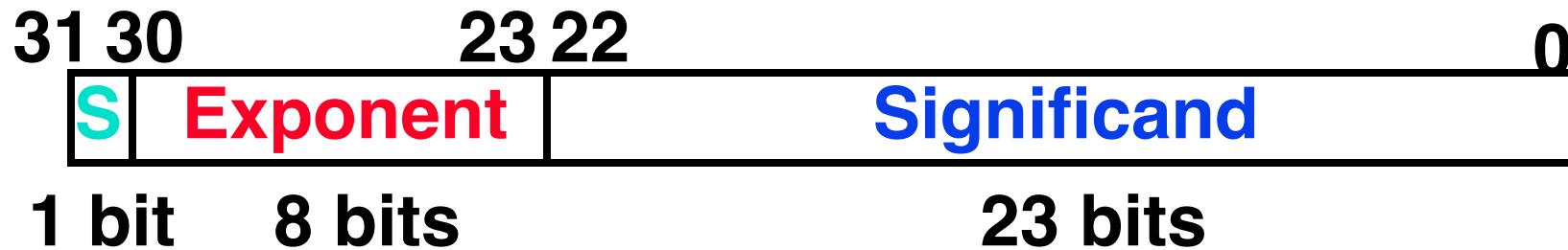
Scientific Notation (in Binary)



- 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.

Floating-Point Representation (1/2)

- Normal format: $+1.\text{xxx...x}_\text{two} * 2^{\text{yyy...y}_\text{two}}$
- Multiple of Word Size (32 bits)



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

- Represent numbers as small as 2.0×10^{-38} to as large as 2.0×10^{38}

Floating-Point Representation (2/2)

- What if result too large?

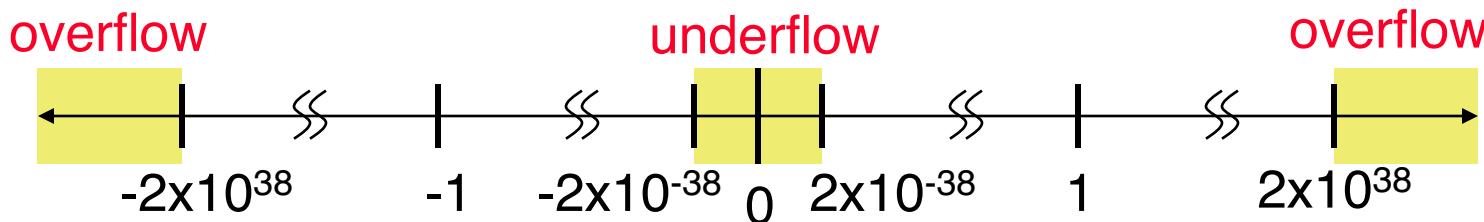
($> 2.0 \times 10^{38}$, $< -2.0 \times 10^{38}$)

- **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?

Clicker Question

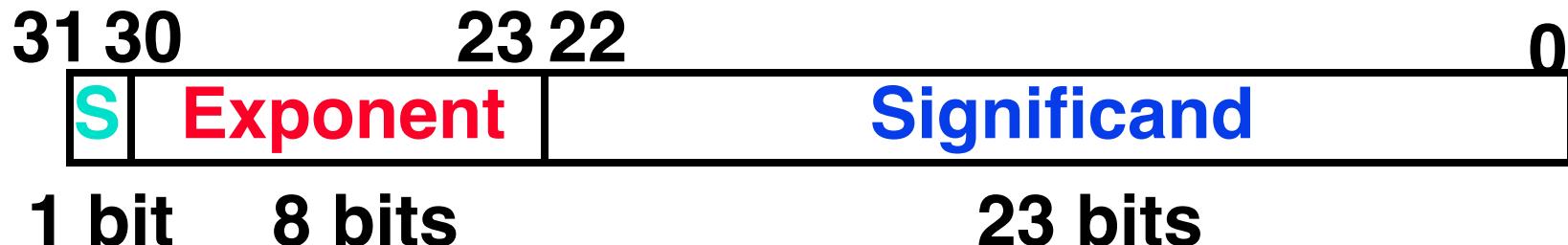
Convert 3.625 to binary point, in scientific notation:

- A: 1.0001×2^1
- B: 1.0100×2^{-1}
- C: 0.110
- D: 11.111
- E: 1.1101×2^1

Break

IEEE 754 Floating-Point Standard (1/3)

Single Precision (Double Precision similar):



- **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

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

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)

“Father” of the Floating point standard

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



**Prof. Kahan
Professor
(Emeritus)
UC Berkeley**

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 00000000000000000000000000000000

-0: 1 00000000 00000000000000000000000000000000

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

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

Representation for Denorms (1/2)

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

- Smallest representable pos num:

$$a = 1.0\ldots_2 \cdot 2^{-126} = 2^{-126}$$

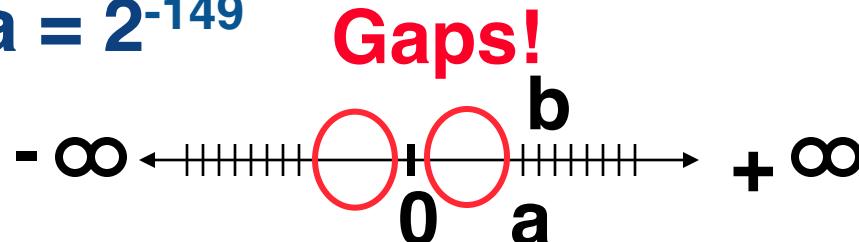
- Second smallest representable pos num:

$$\begin{aligned} b &= 1.000\ldots_2 \cdot 2^{-126} \\ &= (1 + 0.00\ldots_2) \cdot 2^{-126} \\ &= (1 + 2^{-23}) \cdot 2^{-126} \\ &= 2^{-126} + 2^{-149} \end{aligned}$$

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

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

Normalization
and implicit 1
is to blame!



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:
 $b = 2^{-148}$



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> |

Conclusion

- 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)

Exponent tells Significand how much (2^i) to count by (... , $1/4$, $1/2$, 1 , 2 , ...)

Can store NaN,
 $\pm \infty$

- Summary (single precision):

| | | | | |
|----|----------|-------------|----|---|
| 31 | 30 | 23 | 22 | 0 |
| S | Exponent | Significand | | |

1 bit 8 bits

23 bits

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

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

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