

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



Lecturer SOE
Dan Garcia

Lecture 36 – Performance
2010-04-23

How fast is your computer?

CRAY XT5-HE IS FASTEST SUPERCOMPUTER!

Every 6 months (Nov/June), the fastest supercomputers in the world face off. The fastest computer is now the Jaguar, a Linux Cray XT5-HE Opteron Six Core 2.6 GHz, with 224,256 cores, achieves 2.3 PFlops. They use LINPACK floating point benchmark ($A \times B$).



www.top500.org/lists/2009/11
www.nccs.gov/computing-resources/jaguar/

Why Performance? Faster is better!

- **Purchasing Perspective:** given a collection of machines (or upgrade options), which has the
 - best performance ?
 - least cost ?
 - best performance / cost ?
- **Computer Designer Perspective:** faced with design options, which has the
 - best performance improvement ?
 - least cost ?
 - best performance / cost ?
- **All require basis for comparison and metric for evaluation!**
 - Solid metrics lead to solid progress!



Two Notions of "Performance"

Plane	DC to Paris	Top Speed	Passengers	Throughput (pmp)
Boeing 747	6.5 hours	610 mph	470	286,700
BAD/Sud Concorde	3 hours	1350 mph	132	178,200

- Which has higher performance?
 - Interested in time to deliver 100 passengers?
 - Interested in delivering as many passengers per day as possible?
- In a computer, time for one task called **Response Time** or **Execution Time**
- In a computer, tasks per unit time called **Throughput** or **Bandwidth**



Definitions

- Performance is in units of things per sec
 - bigger is better
- If mostly concerned with response time
 - $\text{performance}(x) = \frac{1}{\text{execution_time}(x)}$
- “ **F**(ast) is *n* times faster than **S**(low) ” means:

$$n = \frac{\text{performance}(F)}{\text{performance}(S)} = \frac{\text{execution_time}(S)}{\text{execution_time}(F)}$$



Example of Response Time v. Throughput

- **Time of Concorde vs. Boeing 747?**
 - Concorde is 6.5 hours / 3 hours
= 2.2 times faster
 - Concorde is 2.2 times ("120%") faster in terms of flying time (response time)
- **Throughput of Boeing vs. Concorde?**
 - Boeing 747: 286,700 pmph / 178,200 pmph
= 1.6 times faster
 - Boeing is 1.6 times ("60%") faster in terms of throughput
- We will focus primarily on **response time**.



Words, Words, Words...

- Will (try to) stick to “**n times faster**”; its less confusing than “m % faster”
- As faster means both decreased execution time and increased performance, to reduce confusion we will (and you should) use “improve execution time” or “improve performance”



What is Time?

- **Straightforward definition of time:**
 - Total time to complete a task, including disk accesses, memory accesses, I/O activities, operating system overhead, ...
 - "real time", "response time" or "elapsed time"
- **Alternative: just time processor (CPU) is working only on your program (since multiple processes running at same time)**
 - "CPU execution time" or "CPU time"
 - Often divided into system CPU time (in OS) and user CPU time (in user program)



How to Measure Time?

- Real Time \Rightarrow Actual time elapsed
- CPU Time: Computers constructed using a clock that runs at a constant rate and determines when events take place in the hardware
 - These discrete time intervals called clock cycles (or informally clocks or cycles)
 - Length of clock period: clock cycle time (e.g., $\frac{1}{2}$ nanoseconds or $\frac{1}{2}$ ns) and clock rate (e.g., 2 gigahertz, or 2 GHz), which is the inverse of the clock period; use these!



Measuring Time using Clock Cycles (1/2)

- CPU execution time for a program

- Units of [seconds / program] or [s/p]

= Clock Cycles for a program x Clock Period

- Units of [s/p] = [cycles / p] x [s / cycle] = [c/p] x [s/c]

- Or

= Clock Cycles for a program [c / p]
Clock Rate [c / s]



Measuring Time using Clock Cycles (2/2)

- One way to define clock cycles:

Clock Cycles for program [c/p]

= Instructions for a program [i/p]
(called "Instruction Count")

× Average Clock cycles Per Instruction [c/i]
(abbreviated "CPI")

- CPI one way to compare two machines with **same** instruction set, since Instruction Count would be the same



Performance Calculation (1/2)

- CPU execution time for program [s/p]
= **Clock Cycles for program** [c/p]
x **Clock Cycle Time** [s/c]

- Substituting for clock cycles:

$$\begin{aligned} &\text{CPU execution time for program [s/p]} \\ &= (\text{Instruction Count [i/p]} \times \text{CPI [c/i]}) \\ &\quad \times \text{Clock Cycle Time [s/c]} \end{aligned}$$

$$= \underline{\text{Instruction Count}} \times \underline{\text{CPI}} \times \underline{\text{Clock Cycle Time}}$$



Performance Calculation (2/2)

$$\text{CPU time} = \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Cycles}}{\text{Instruction}} \times \frac{\text{Seconds}}{\text{Cycle}}$$

$$\text{CPU time} = \frac{\cancel{\text{Instructions}}}{\text{Program}} \times \frac{\text{Cycles}}{\cancel{\text{Instruction}}} \times \frac{\text{Seconds}}{\text{Cycle}}$$

$$\text{CPU time} = \frac{\cancel{\text{Instructions}}}{\text{Program}} \times \frac{\cancel{\text{Cycles}}}{\cancel{\text{Instruction}}} \times \frac{\text{Seconds}}{\cancel{\text{Cycle}}}$$

$$\text{CPU time} = \frac{\text{Seconds}}{\text{Program}}$$

Product of all 3 terms: if missing a term, can't predict time, the real measure of performance



How Calculate the 3 Components?

- **Clock Cycle Time:** in specification of computer (Clock Rate in advertisements)
- **Instruction Count:**
 - Count instructions in loop of small program
 - Use simulator to count instructions
 - Hardware counter in spec. register
 - (Pentium II,III,4)
- **CPI:**
 - Calculate:
$$\frac{\text{Execution Time}}{\text{Instruction Count}} \div \text{Clock cycle time}$$
 - Hardware counter in special register (PII,III,4)



Calculating CPI Another Way

- First calculate CPI for each individual instruction (add, sub, and, etc.)
- Next calculate frequency of each individual instruction
- Finally multiply these two for each instruction and add them up to get final CPI (the weighted sum)



Example (RISC processor)

Op	Freq _i	CPI _i	Prod	(% Time)
ALU	50%	1	.5	(23%)
Load	20%	5	1.0	(45%)
Store	10%	3	.3	(14%)
Branch	20%	2	.4	(18%)

2.2

Instruction Mix

(Where time spent)

- What if Branch instructions twice as fast?



What Programs Measure for Comparison?

- Ideally run typical programs with typical input before purchase, or before even build machine
 - Called a “workload”; For example:
 - Engineer uses compiler, spreadsheet
 - Author uses word processor, drawing program, compression software
- In some situations its hard to do
 - Don't have access to machine to “benchmark” before purchase
 - Don't know workload in future



Benchmarks

- Obviously, apparent speed of processor depends on code used to test it
- Need industry standards so that different processors can be fairly compared
- Companies exist that create these **benchmarks**: “typical” code used to evaluate systems
- Need to be changed every ~5 years since designers could (and do!) target for these standard benchmarks



Example Standardized Benchmarks (1/2)

- **Standard Performance Evaluation Corporation (SPEC) SPEC CPU2006**
 - CINT2006 **12** integer (perl, bzip, gcc, go, ...)
 - CFP2006 **17** floating-point (povray, bwaves, ...)
 - All relative to base machine (which gets **100**)
Sun Ultra Enterprise 2 w/296 MHz UltraSPARC II
 - They measure
 - System speed (SPECint2006)
 - System throughput (SPECint_rate2006)
 - www.spec.org/osg/cpu2006/



Example Standardized Benchmarks (2/2)

■ SPEC

- Benchmarks distributed in source code
- Members of consortium select workload
 - 30+ companies, 40+ universities, research labs
- Compiler, machine designers target benchmarks, so try to change every 5 years
- SPEC CPU2006:

CINT2006

perlbench	C	Perl Programming language
bzip2	C	Compression
gcc	C	C Programming Language Compiler
mcf	C	Combinatorial Optimization
gobmk	C	Artificial Intelligence : Go
hmmcr	C	Search Gene Sequence
sjeng	C	Artificial Intelligence : Chess
libquantum	C	Simulates quantum computer
h264ref	C	H.264 Video compression
omnetpp	C++	Discrete Event Simulation
astar	C++	Path-finding Algorithms
xalancbmk	C++	XML Processing

CFP2006

bwaves	Fortran	Fluid Dynamics
gamess	Fortran	Quantum Chemistry
milc	C	Physics / Quantum Chromodynamics
zeusmp	Fortran	Physics / CFD
gromacs	C, Fortran	Biochemistry / Molecular Dynamics
cactusADM	C, Fortran	Physics / General Relativity
leslie3d	Fortran	Fluid Dynamics
namd	C++	Biology / Molecular Dynamics
deall1	C++	Finite Element Analysis
soplex	C++	Linear Programming, Optimization
povray	C++	Image Ray-tracing
calculix	C, Fortran	Structural Mechanics
GemsFDTD	Fortran	Computational Electromagnetics
tonto	Fortran	Quantum Chemistry
lbm	C	Fluid Dynamics
wrf	C, Fortran	Weather
sphinx3	C	Speech recognition



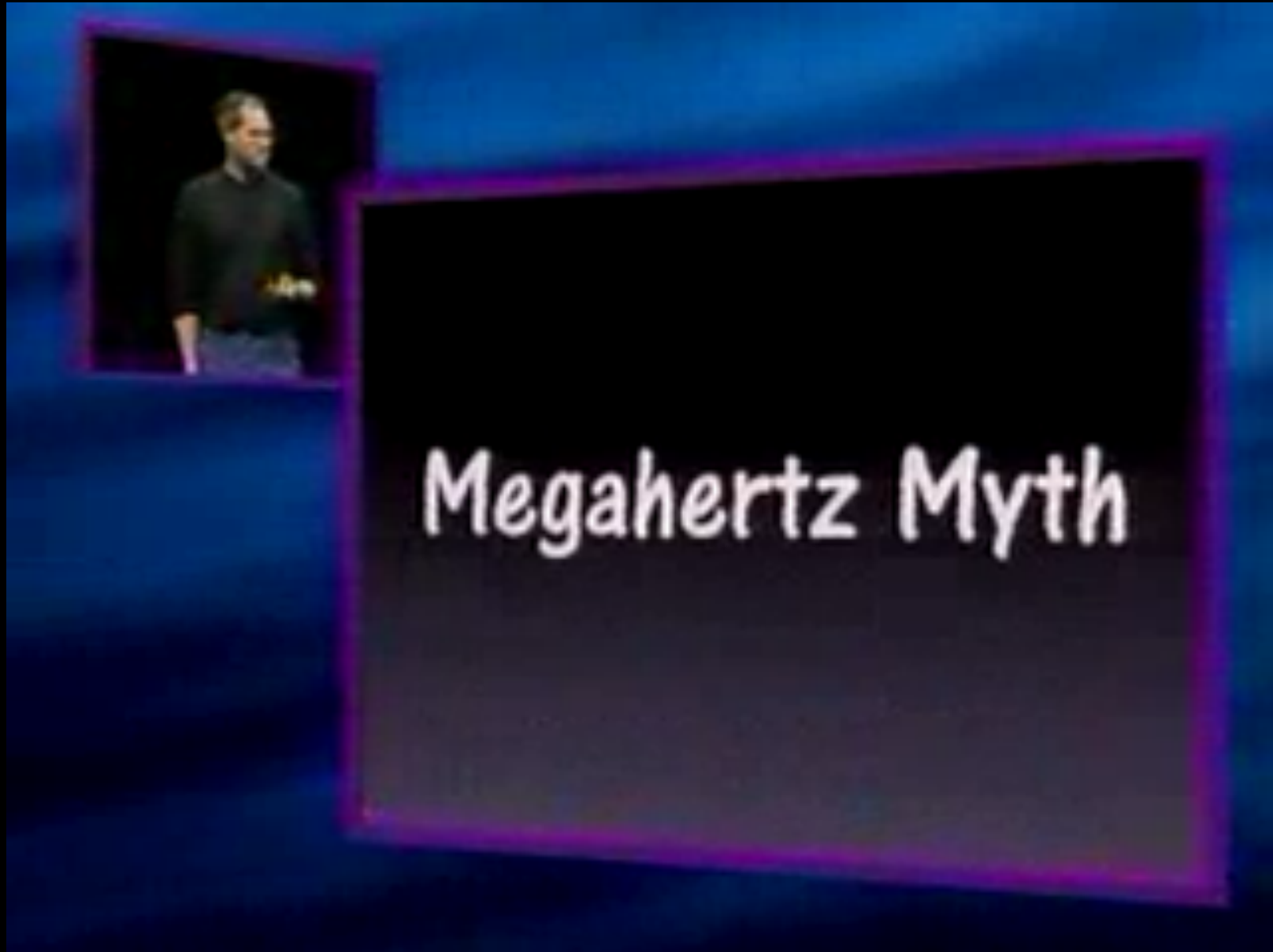
Performance Evaluation: The Demo

If we're talking about performance, let's discuss the ways shady salespeople have fooled consumers (so you don't get taken!)

1. Play a movie
2. Preprocess all available data
3. Run it on a stock machine in which "no expense was spared"
4. Only run the demo through a script
5. Never let the user touch it



Megahertz Myth Marketing Movie



Peer Instruction

- 1) The Sieve of Eratosthenes and Quicksort were early effective benchmarks.
- 2) A program runs in 100 sec. on a machine, `mult` accounts for 80 sec. of that. If we want to make the program run 6 times faster, we need to up the speed of `mults` by AT LEAST 6.

	12
a)	FF
b)	FT
c)	TF
d)	TT



Peer Instruction Answers

- 1) The Sieve of Eratosthenes and Quicksort were early effective benchmarks. **FALSE**
- 2) A program runs in 100 sec. on a machine, mult accounts for 80 sec. of that. If we want to make the program run 6 times faster, we need to up the speed of mults by AT LEAST 6. **FALSE**

1. Early benchmarks? Yes.
Effective? No. Too simple!
2. 6 times faster = 16 sec.
mults must take -4 sec!
I.e., impossible!

	12
a)	FF
b)	FT
c)	TF
d)	TT



“And in conclusion...”

$$\text{CPU time} = \frac{\text{Instructions}}{\text{Instruction}} \times \frac{\text{Cycles}}{\text{Cycle}} \times \frac{\text{Seconds}}{\text{Seconds}}$$

Program Instruction Cycle

- Latency v. Throughput
- **Performance doesn't depend on any single factor:** need Instruction Count, Clocks Per Instruction (CPI) and Clock Rate to get valid estimations
- **User Time:** time user waits for program to execute: depends heavily on how OS switches between tasks
- **CPU Time:** time spent executing a single program: depends solely on design of processor (datapath, pipelining effectiveness, caches, etc.)
- **Benchmarks**
 - Attempt to predict perf, Updated every few years
 - Measure everything from simulation of desktop graphics programs to battery life
- **Megahertz Myth**
 - **MHz ≠ performance, it's just one factor**

