

## CS 700: Quantitative Methods & Experimental Design in Computer Science

Sanjeev Setia  
Dept of Computer Science  
George Mason University

### Logistics

- ❑ Grade: 35% project, 25% Homework assignments  
20% midterm, 20% take home final
- ❑ Slides, assignments, reading material on class web page <http://www.cs.gmu.edu/~setia/cs700/>
- ❑ Several small assignments related to material discussed in class
  - Not all will be graded, but we will go over solutions in class
- ❑ Term project
  - should involve experimentation (measurement, simulation)
  - select a topic in your research area if possible
  - apply techniques discussed in this class

## Readings

- ❑ Textbook
  - David Lilja, "Measuring Computer Performance: A Practitioner's Guide" OR
  - Raj Jain, "Art of Computer Systems Performance Analysis"
- ❑ Class notes, slides
- ❑ Relevant research articles (links on class web site)
- ❑ Books on probability and statistics for engineers (see syllabus)
- ❑ More specialized topics
  - Cohen "Empirical techniques in AI"
  - Crovella et al "Internet Measurement"

3

## Course Topics

- ❑ Basic techniques in "experimental" computer science
  - Basic measurement tools and techniques
  - Simulation
  - Design of experiments
- ❑ Quantitative Methods
  - Use of statistical techniques in design of experiments
  - Use of statistical techniques in comparing alternatives
  - Characterizing and interpreting measured data
- ❑ Simple analytical modeling

Most examples will be from performance measurement of computer systems and networks, but techniques are applicable in all fields of CS

4

## Experimental Science

### Scientific Method

1. Identify a problem and form hypothesis
  - ◆ Hypothesis must be testable and refutable
2. Design an experiment
3. Conduct the experiment
4. Perform hypothesis testing
  - ◆ Use statistical techniques

What about Computer Science?

5

## Experimental Computer Science

Three definitions (Feitelson, 2006)

- Building of systems, hardware or software
  - Counterpart to "theoretical CS"
- Experimentation as a feedback step in engineering loop
- Evaluation of computer systems using the methodologies of the natural sciences, i.e. rigorous methodologies
  - Focus of this class

Feitelson makes the case that there is a place for observation of the real world, as in the natural sciences, e.g., analyzing measured network traffic

6

## The Role of Experimentation in CS

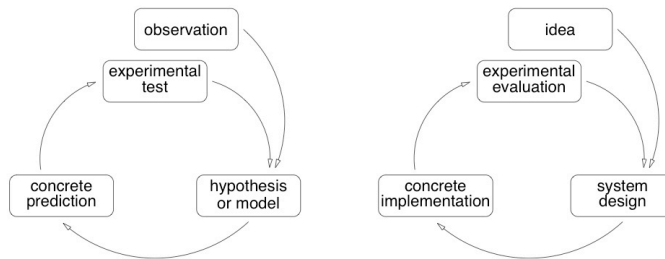


Figure 1: A comparison of the scientific method (on the left) with the role of experimentation in system design (right).

7

## Schedule

- Introduction
- Metrics
- Summarizing Measured Data
- Measurement Techniques
- Simulation
- Comparing Alternatives
- Linear Regression Models
- Design of experiments
- Interpreting & characterizing measured data
- Analytical Modeling

8

## Course Goals

- ❑ Understand the inherent trade-offs involved in using simulation, measurement, and analytical modeling.
- ❑ Rigorously compare computer systems/networks/software/artifacts/... often in the presence of measurement noise
  - Usually compare performance but in many fields of CS, "quality" of the output is more important than raw performance, e.g. face recognition software
- ❑ Determine whether a change made to a system has a statistically significant impact

9

## Course Goals

- ❑ Use statistical tools to reduce the number of simulations that need to be performed of a computer system.
- ❑ Design a set of experiments to obtain the most information for a given level of effort.

10

## Course Goals

- Provide intuitive conceptual background for some standard statistical tools.
  - Draw meaningful conclusions in presence of noisy measurements.
  - Allow you to correctly and intelligently apply techniques in new situations.

11

## Course Goals

- Present techniques for aggregating and interpreting large quantities of data.
    - Obtain a big-picture view of your results.
    - Obtain new insights from complex measurement and simulation results.
- E.g. How does a new feature impact the overall system?

12

## Agenda

### □ Today

- Overview of course
- Performance metrics
  - Characteristics of good metrics
  - Standard processor and system metrics
  - Speedup and relative change

13

## Agenda (cont.)

### □ Measurement tools and techniques

- Fundamental strategies
- Interval timers
- Cycle counters
- Program profiling
- Tracing
- Indirect measurement
- Measuring network delays and bandwidth

14

## Agenda (cont.)

### □ Simulation

- Types of simulations
- Random number generation
- Verification and validation

15

## Agenda (cont.)

### □ Statistical interpretation of measured data

- Arithmetic, harmonic, geometric means
- Sources of measurement error
- Confidence intervals
- Statistically comparing alternatives

16



## Agenda (cont.)

- Design of experiments
  - Terminology
  - One-factor analysis of variance (ANOVA)
  - Two-factor ANOVA
  - Generalized  $m$ -factor experiments
  - Fractional factorial designs
  - Multi-factorial designs (Plackett and Berman)

17

## Agenda (cont'd)

- Characterizing Measured Data
  - Workload Characterization
  - Fitting Measured Data to Known Distributions
    - Q-Q plots
    - Chi-square, K-S tests

18

## Agenda (cont'd)

- Simple analytical modeling
  - Simple queuing models
    - Single queue models
    - Networks of queues
  - Operational analysis
    - Little's Law

19

## Readings

- Dror Feitelson, "Experimental Computer Science: The need for a cultural change"

20

## Metrics

### Performance metrics

- ❑ What is a performance metric?
- ❑ Characteristics of good metrics
- ❑ Standard processor and system metrics
- ❑ Speedup and relative change



## What is a performance metric?

- ❑ Count
  - Of how many times an event occurs
- ❑ Duration
  - Of a time interval
- ❑ Size
  - Of some parameter
- ❑ A value derived from these fundamental measurements

23

## Time-normalized metrics

- ❑ "Rate" metrics
  - Normalize metric to common time basis
    - Transactions per second
    - Bytes per second
  - $(\text{Number of events}) \div (\text{time interval over which events occurred})$
- ❑ "Throughput"
- ❑ Useful for comparing measurements over different time intervals

24

## What makes a "good" metric?

- ❑ Allows accurate and detailed comparisons
- ❑ Leads to correct conclusions
- ❑ Is well understood by everyone
- ❑ Has a quantitative basis
- ❑ A good metric helps avoid erroneous conclusions

25

## Good metrics are ...

- ❑ Linear
  - Fits with our intuition
  - If metric increases 2x, performance should increase 2x
  - Not an absolute requirement, but very appealing
    - dB scale to measure sound is nonlinear

26

## Good metrics are ...

### □ Reliable

- If metric A > metric B
  - Then, Performance A > Performance B
- Seems obvious!
- However,
  - MIPS(A) > MIPS(B), but
  - Execution time (A) > Execution time (B)

27

## Good metrics are ...

### □ Repeatable

- Same value is measured each time an experiment is performed
- Must be *deterministic*
- Seems obvious, but not always true...
  - E.g. Time-sharing changes measured execution time

28

## Good metrics are ...

### □ Easy to use

- No one will use it if it is hard to measure
- Hard to measure/derive
  - → less likely to be measured correctly
- A *wrong value* is worse than a bad metric!

29

## Good metrics are ...

### □ Consistent

- Units and definition are constant across systems
- Seems obvious, but often not true...
  - E.g. MIPS, MFLOPS
- Inconsistent → impossible to make comparisons

30

## Good metrics are ...

### □ Independent

- A lot of \$\$\$ riding on performance results
- Pressure on manufacturers to *optimize* for a particular metric
- Pressure to influence *definition* of a metric
- But a good metric is independent of this pressure

31

## Good metrics are ...

- Linear -- *nice, but not necessary*
- Reliable -- *required*
- Repeatable -- *required*
- Easy to use -- *nice, but not necessary*
- Consistent -- *required*
- Independent -- *required*

32



## Clock rate

- ❑ Faster clock == higher performance
  - 1 GHz processor always better than 2 GHz
- ❑ But is it a proportional increase?
- ❑ What about architectural differences?
  - Actual operations performed per cycle
  - Clocks per instruction (CPI)
  - Penalty on branches due to pipeline depth
- ❑ What if the processor is not the bottleneck?
  - Memory and I/O delays

33

## Clock rate

- ❑ (Faster clock)  
≠ (better performance)
- ❑ A good first-order metric

Linear	
Reliable	
Repeatable	☺
Easy to measure	☺
Consistent	☺
Independent	☺

34

## MIPS

- ❑ Measure of computation "speed"
- ❑ *Millions of instructions executed per second*
- ❑  $MIPS = n / (T_e * 1000000)$ 
  - $n$  = number of instructions
  - $T_e$  = execution time
- ❑ Physical analog
  - Distance traveled per unit time

35

## MIPS

- ❑ But how much actual computation per instruction?
  - E.g. CISC vs. RISC
  - Clocks per instruction (CPI)
- ❑ *MIPS = Meaningless Indicator of Performance*

Linear	
Reliable	
Repeatable	☺
Easy to measure	☺
Consistent	
Independent	☺

36

## MFLOPS

- ❑ Better definition of "distance traveled"
- ❑ 1 unit of computation (~distance)  $\equiv$  1 floating-point operation
- ❑ *Millions of floating-point ops per second*
- ❑  $\text{MFLOPS} = f / (T_e * 1000000)$ 
  - $f$  = number of floating-point instructions
  - $T_e$  = execution time
- ❑ GFLOPS, TFLOPS,...

37

## MFLOPS

- ❑ Integer program = 0 MFLOPS
  - But still doing useful work
  - Sorting, searching, etc.
- ❑ How to count a FLOP?
  - E.g. transcendental ops, roots
- ❑ Too much flexibility in definition of a FLOP
- ❑ Not consistent across machines

Linear	
Reliable	
Repeatable	☺
Easy to measure	☺
Consistent	
Independent	

38

## SPEC

- ❑ System Performance Evaluation Coop
- ❑ Computer manufacturers select "representative" programs for benchmark suite
- ❑ Standardized methodology
  - Measure execution times
  - Normalize to standard basis machine
  - SPECmark = geometric mean of normalized values

39

## SPEC

- ❑ Geometric mean is inappropriate (more later)
- ❑ SPEC rating does not correspond to execution times of non-SPEC programs
- ❑ Subject to tinkering
  - Committee determines which programs should be part of the suite
  - Targeted compiler optimizations

Linear	
Reliable	
Repeatable	☺
Easy to measure	$\frac{1}{2}$ ☺
Consistent	☺
Independent	

40

## Execution time

- ❑ Ultimately interested in **time** required to execute *your* program
- ❑ Smallest total execution time == highest performance
- ❑ Compare times directly
- ❑ Derive appropriate rates
- ❑ Time == *fundamental metric* of performance
  - If you can't measure time, you don't know anything

41

## Execution time

- ❑ "Stopwatch" measured execution time

```
Start_count = read_timer();
    Portion of program to be measured
Stop_count = read_timer();
Elapsed_time = (stop_count - start_count)
    * clock_period;
```
- ❑ Measures "wall clock" time
  - Includes I/O waits, time-sharing, OS overhead,
  - ...
- ❑ "CPU time" -- include only processor time

42

## Execution time

- ❑ Best to report both wall clock and CPU times
- ❑ Includes system noise effects
  - Background OS tasks
  - Virtual to physical page mapping
  - Random cache mapping and replacement
  - Variable system load
- ❑ Report both mean and variance (more later)

Linear	😊
Reliable	😊
Repeatable	≈ 😊
Easy to measure	😊
Consistent	😊
Independent	😊

43

## Performance metrics summary

	Clock	MIPS	MFLOPS	SPEC	TIME
Linear					😊
Reliable					≈ 😊
Repeatable	😊	😊	😊	😊	😊
Easy to measure	😊	😊	😊	$\frac{1}{2}$ 😊	😊
Consistent	😊			😊	😊
Independent	😊	😊			😊

44

## Other metrics

- ❑ Response time
  - Elapsed time from request to response
- ❑ Throughput
  - Jobs, operations completed per unit time
  - E.g. video frames per second
- ❑ Bandwidth
  - Bits per second
- ❑ *Ad hoc* metrics
  - Defined for a specific need

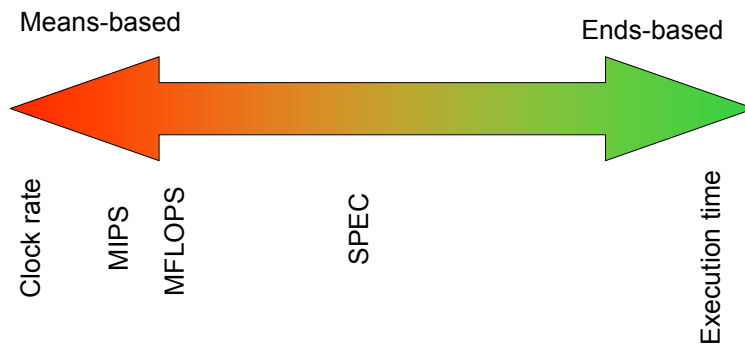
45

## Means vs. ends metrics

- ❑ Means-based metrics
  - Measure what was done
  - Whether or not it was useful!
    - o Nop instructions, multiply by zero, ...
  - Produces unreliable metrics
- ❑ Ends-based metrics
  - Measures progress towards a goal
  - Only counts what is actually accomplished

46

## Means vs. ends metrics



47

## Speedup

- Speedup of System 1 w.r.t System 2
  - $S_{2,1}$  such that:  $R_2 = S_{2,1} R_1$
  - $R_1$  = "speed" of System 1
  - $R_2$  = "speed" of System 2
- *System 2 is  $S_{2,1}$  times faster than System 1*

48



## Speedup

- "Speed" is a rate metric
  - $R_i = D_i / T_i$
  - $D_i \sim$  "distance traveled" by System i
- Run the same benchmark program on both systems
  - $\rightarrow D_1 = D_2 = D$
- Total work done by each system is defined to be "execution of the same program"
  - Independent of architectural differences

49

## Speedup

$$\begin{aligned} S_{2,1} &= \frac{R_2}{R_1} = \frac{D_2 / T_2}{D_1 / T_1} = \frac{D / T_2}{D / T_1} \\ &= \frac{T_1}{T_2} \end{aligned}$$

50

## Speedup

$S_{2,1} > 1 \Rightarrow$  System 2 is faster than System 1

$S_{2,1} < 1 \Rightarrow$  System 2 is slower than System 1

51

## Relative change

- Performance of System 2 relative to System 1 as *percent change*

$$\Delta_{2,1} = \frac{R_2 - R_1}{R_1}$$

52

### Relative change

$$\begin{aligned}\Delta_{2,1} &= \frac{R_2 - R_1}{R_1} = \frac{D/T_2 - D/T_1}{D/T_1} \\ &= \frac{T_1 - T_2}{T_2} \\ &= S_{2,1} - 1\end{aligned}$$

53

### Relative change

$\Delta_{2,1} > 0 \Rightarrow$  System 2 is faster than System 1

$\Delta_{2,1} < 0 \Rightarrow$  System 2 is slower than System 1

54

## Important Points

- Metrics can be
  - Counts
  - Durations
  - Sizes
  - Some combination of the above

55

## Important Points

- Good metrics are
  - Linear
    - More "natural"
  - Reliable
    - Useful for comparison and prediction
  - Repeatable
  - Easy to use
  - Consistent
    - Definition is the same across different systems
  - Independent of outside influences

56