## Introduction to Networking and Systems Measurements

Reproducible Experiments



**Dr Andrew W. Moore** 

andrew.moore@cl.cam.ac.uk

**Dr Noa Zilberman** 

noa.zilberman@cl.cam.ac.uk

## Reproducibility vs Repeatability

- Repeatability measures the variation in measurements taken by a single instrument or person under the same conditions
- Reproducibility measures whether an entire study or experiment can be reproduced in its entirety.

## Why?

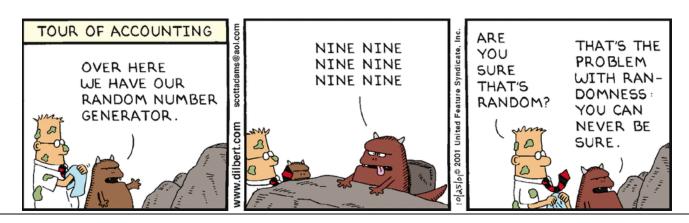
- Establish variance Repeatability
- Establish reliability Repeatability
- Evaluate a new method Reproducability
- Eval. a new environment Reproducability
- Evaluate a new approach Reproducability

#### Variables and Constants

- Why will tell us what we want to vary and
- Why what we need to hold constant

#### Random-ish

- Rarely do we want true-random
- Typically we want pseudo-random
- Often we want to specify the seed(s)
- Knuth gets this right take care everywhere



#### Method and Environment

- Simualtion?
- Emulation?
- Implementation driven evaluation?
- Deployment?
- Partial emulation?
- Partial implementation driven evaluation?

## Software tools: Scripts, Make, etc

- We have some quite useful repeatability tools:
- Make (links dependencies)
- Scripts sort of documents what you actually need to do to get from (A) to (B)

So please use them.

## Machines (/Hardware)

- Memory? CPU? Disk type and config?
- Hyperthreading and temperature controls?
- Which slots where stuff in?
- Switch config? Switch hardware? Which actual Switch?
- Which transceivers? NICs? cables?

- Tell me again which disk did you dump data to?
- (Oh you forgot to mention the periodic process that moved the data from your machine to another machine so the local disk didn't overrun....)

#### Workloads

- Why is this workload the right one?
  - Stress testing?
- Did you use the workload-generator correctly?
- Record everything from command line options to software and library versions.

#### **Benchmarks**

Often well equipped to run with good reproducibility

- Often not representative of what you want
- Benchmarks might exercise, but just like in fitness training: exercise is not competition.

## Logged data

- Lets talk about time....
  - ➤ No god clock
  - Many representations
  - TimeZone is fun
  - UNIX time is fine, sometimes...
- Text records are nice (for humans)
- Binary records are nice (for programmes)

#### So what is meta-data?

The other stuff needed to repeat precisely the same experiment

- Make and Model (and firmware and config)
- DNS (yep the whole damn thing ok just the entries for your systems)
- Bootp/dhcp/activedirectory all state

### Examples where I screwed up

#### Things I forgot:

- DNS (yep the whole damn thing ok just the DNS entries for all hosts in the trace)
- Snapshot of the code base of the executable we used
- Photo of the setup

# Try stuff! (don't be hipster flanders)



## Other peoples work

## To reproduce other peoples work You must get inside other peoples heads

Very few *high-bars* in reproducibility

http://www.cl.cam.ac.uk/research/srg/netos/qjump/repro.html