



# Towards Better Internet Empiricism

Mark Allman

*International Computer Science Institute*

Passive and Active Measurement Conference  
March 2015

*“My job ain’t a job, it’s a damn good time;  
City to city I’m runnin’ my rhymes”*



# PAM 2002





# PAM 2002

## A Scalable System for Sharing Internet Measurements\*

Mark Allman  
BBN Technologies/NASA GRC  
mallman@grc.nasa.gov

Ethan Blanton  
Ohio University  
eblanton@irg.cs.ohiou.edu

Wesley M. Eddy  
Ohio University  
weddy@irg.cs.ohiou.edu

### Abstract

This paper proposes a system for storing and sharing Internet measurement data amongst researchers. The Scalable Internet Measurement Repository (SIMR) is centered around a database of measurements, tools, experiments, users and datasets. From this set of databases users can search for particular measurements, download the tools used to make and analyze those measurements, and quickly ascer-

age more scientists to share their data with their colleagues.

- Much of our understanding about the network is currently limited by our individual abilities to collect data. For instance, [All00] studies TCP connections to a single WWW server. While the data presented in such papers may be useful, the results would be stronger and more compelling if the conclusions were based on mea-

# Data Sharing Goals

- Broader datasets in everyone's hands ...  
... leads to more conclusive results
- Broadens participation beyond those who are data affluent
- Offers transparency to reproduce results
- Fosters longitudinal studies by building an aggregate global dataset

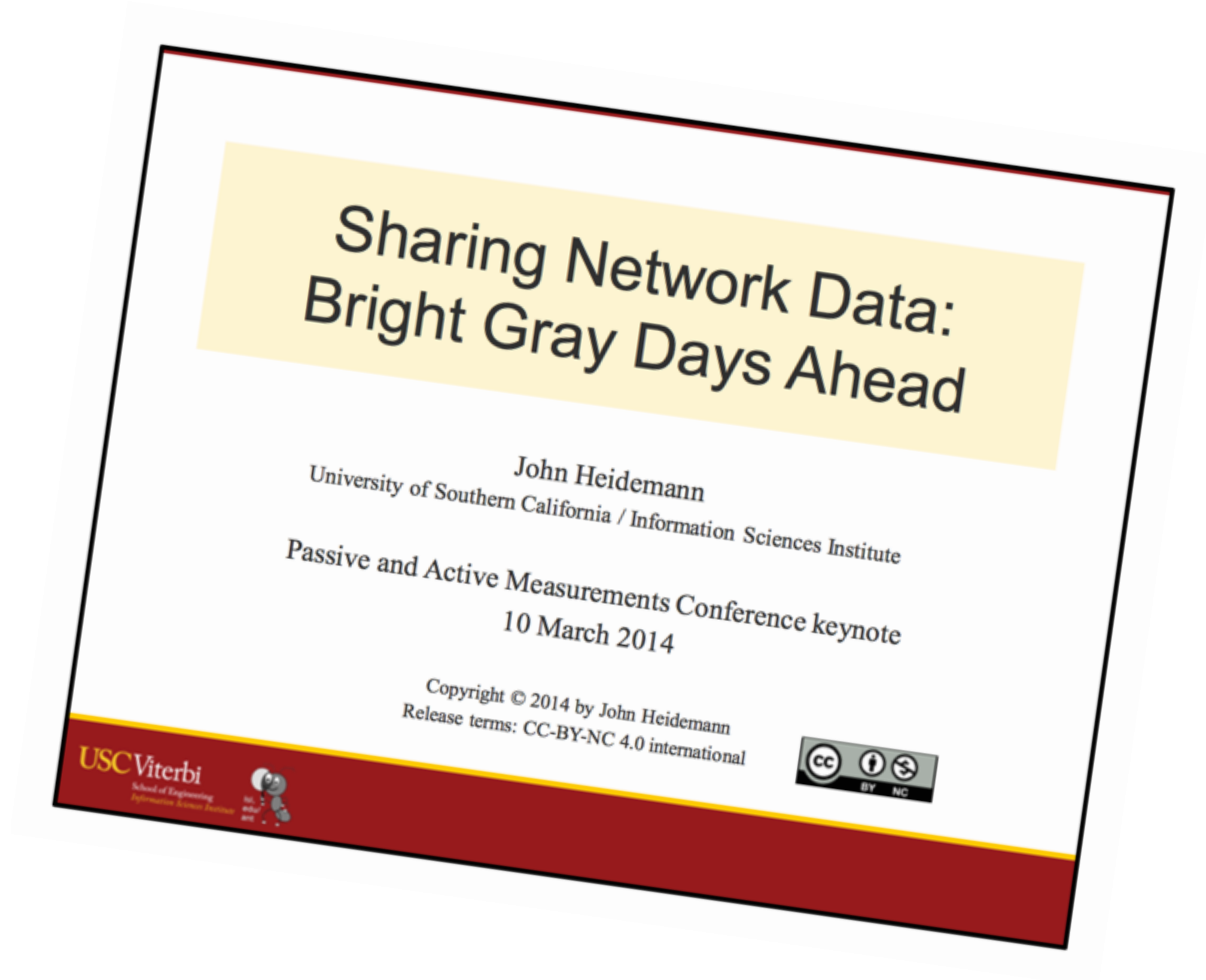
# PAM 2014



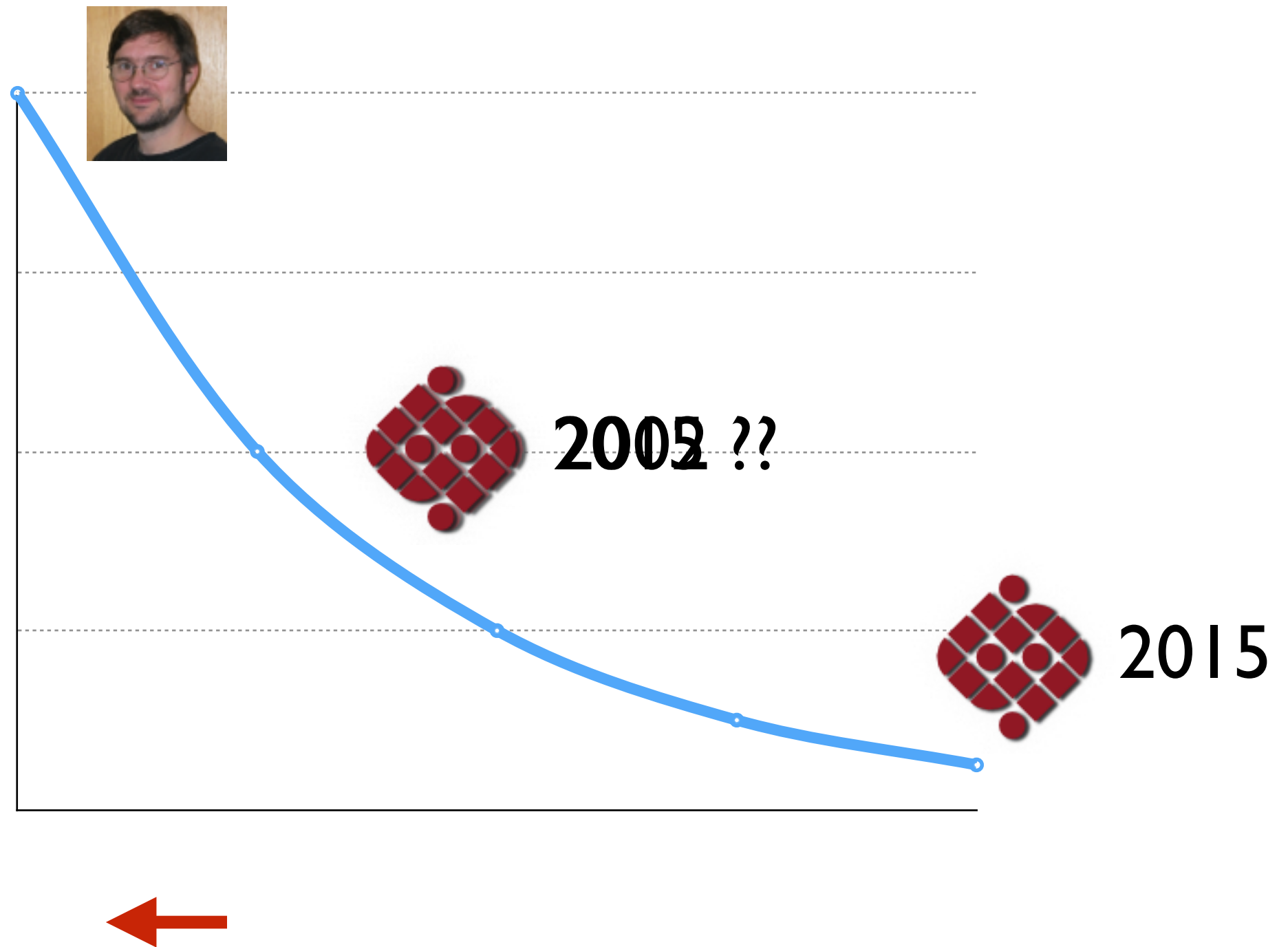
Allman



# PAM 2014



# State of Data Sharing



# Now What?



*“Some of the old songs I sing often, because they help me to reflect on where I've been and that's important for me to do - so I don't lose track of where I am going.”*  
—Johnny Cash



# PAM 2015



Allman

# PAM 2015

- Revisit the *broad notion* of bringing *better science* to bear to understand the Internet

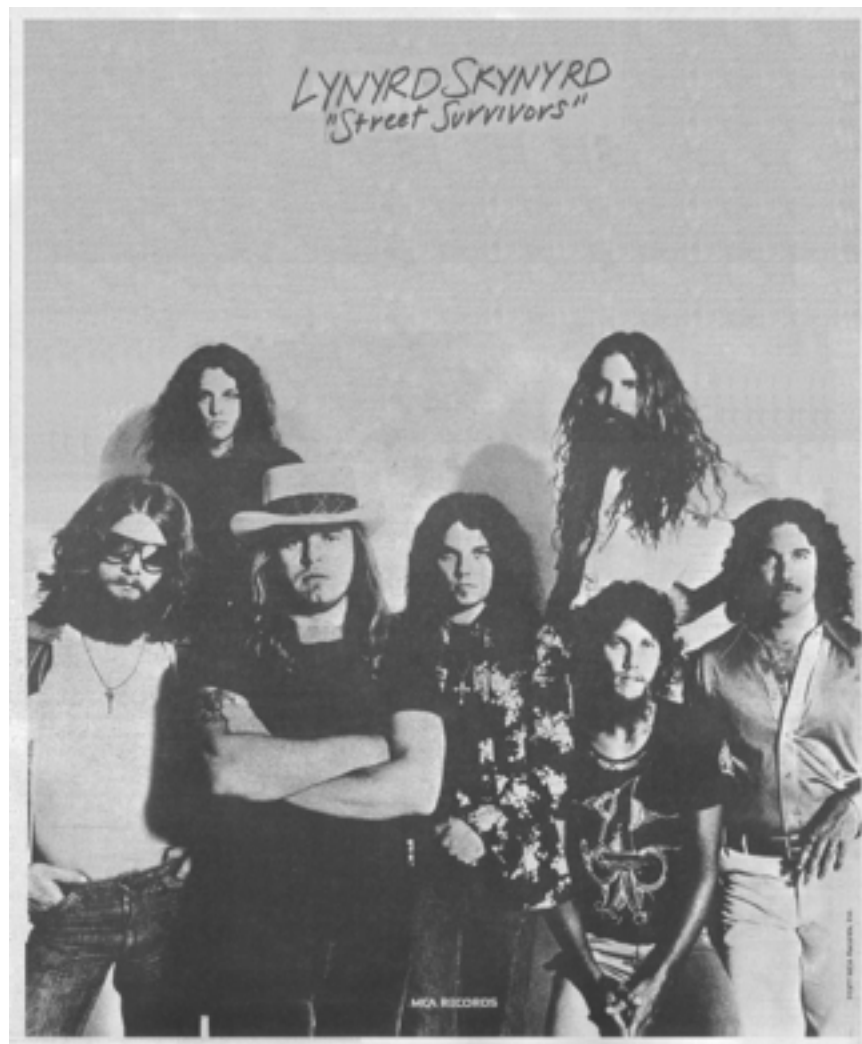
# Embrace Insight (not numbers)



# “I Know A Little”

*“The bigger the city, well the brighter the lights;  
The bigger the dog, well the harder the bite”*

—Skynyrd



# A Few Things I Know ...

- Most TCP connections are short ...  
... but most bytes belong to long connections
- Most paths have little packet reordering ...  
... but some paths reorder many packets
- Most routes are fairly stable ...  
... at least at the timescales of transactions
- Scanning is incessant
- Some countries censor their users

# A Few More Things I Know ...

- Open DNS resolvers are quite prevalent
- ECN is not used much
- DNS pre-fetching is prevalent
- CDNs carry much of the Internet's popular content
- Spoofing is possible from a non-trivial fraction of edge networks
- ...

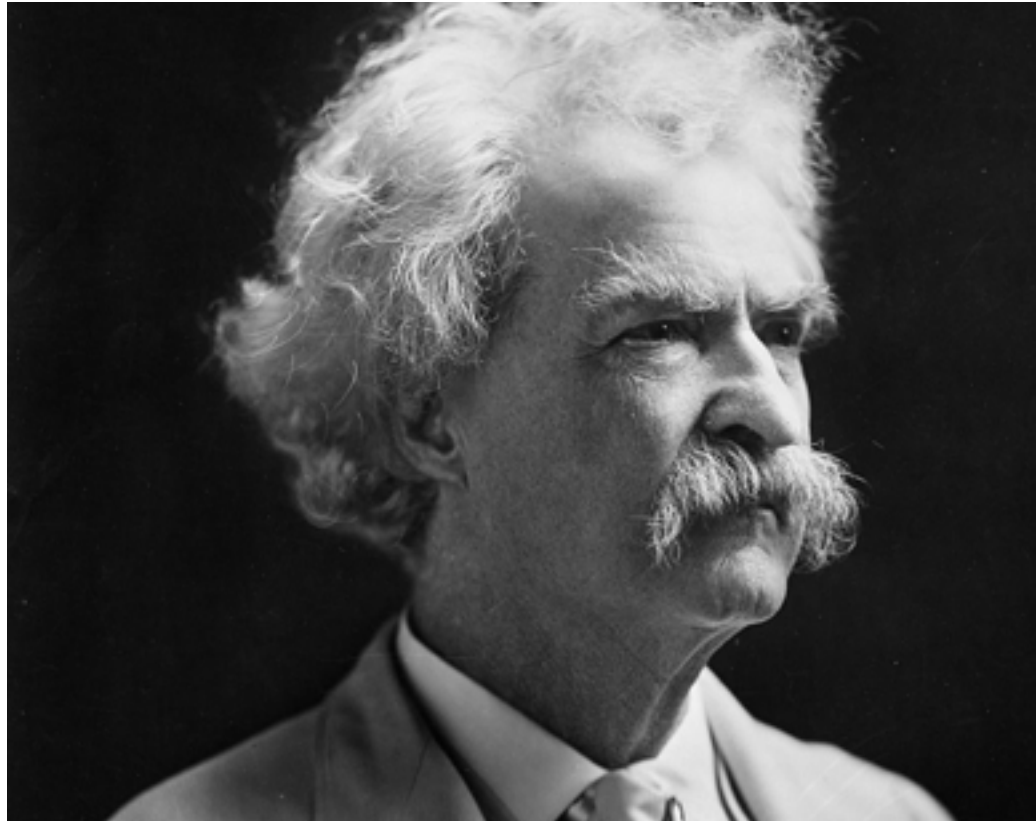


# Embracing Insight

- These knowns are all established via empirical observation
- And, none of these knowns contain *numbers*!

# Embrace Calibration

# “I Know A Little [More]”



*“There are three kinds of lies:  
lies,  
damned lies,  
and statistics.”*

—???

Networking Addendum:  
*... and traceroute.*



# Measurements vs. Reality

- Measurements are an *approximation of reality*

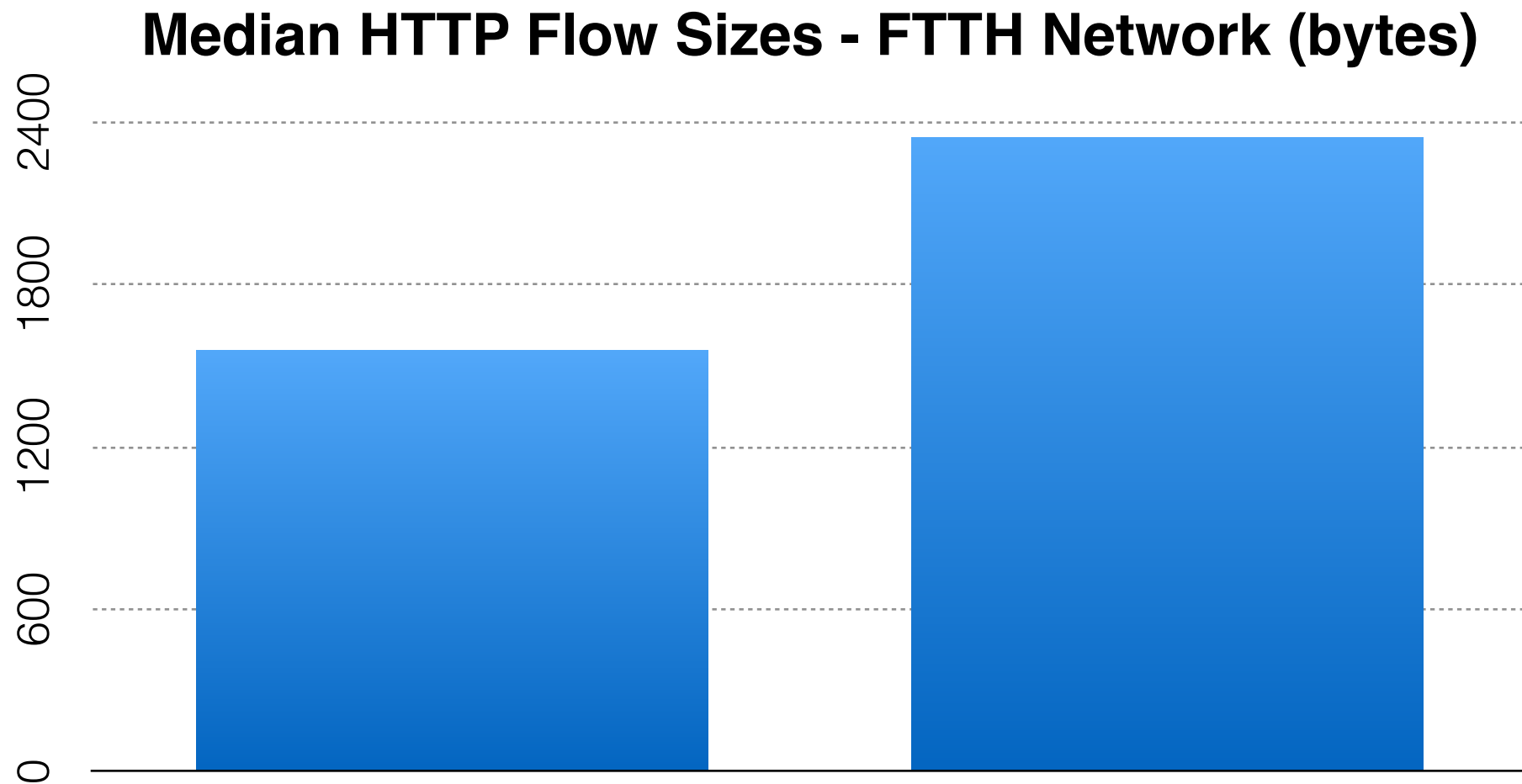
# Measurements vs. Reality

- We well understand that network behavior naturally varies
  - e.g., RTT varies over time
- But, we also need error bars for our collection and processing!

# Measurements vs. Reality

- Before we can get to the insight we need to understand the contours of the approximation
- I.e., we need to *calibrate* our data collection process

# Measurements vs. Reality



# Calibration

- Sound calibration of data lends credibility to results
- I.e., by developing the difference between network behavior and measurement artifacts
- I.e., by demonstrating carefulness



# Calibration Process

- Always validate basic properties in the data
  - “realness” of events, timing accuracy, etc.
- Calibration is a process, not an initial phase
  - Plan to re-visit calibration as necessary

# Calibration Process

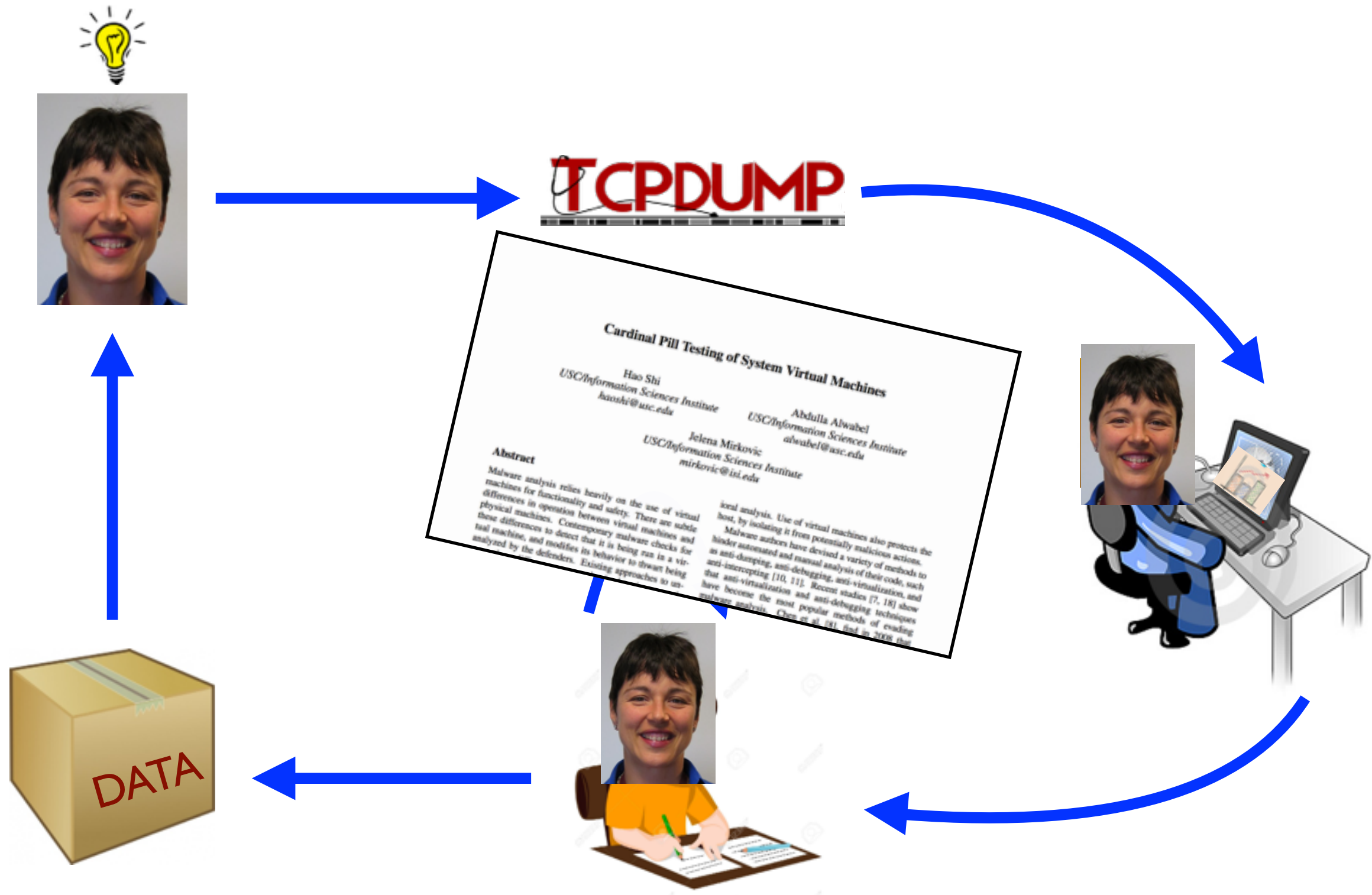
- Not always appreciated by reviewers ...
  - Mysterious to me ...
- *Stop beating each other up for being careful!*

# Embrace Hoarding

# How We Work ...

- We compartmentalize ...
  - ... a paper
  - ... a project
  - ... an internship
  - ... a semester
  - ... a student
  - ... a talk

# Circle of Research

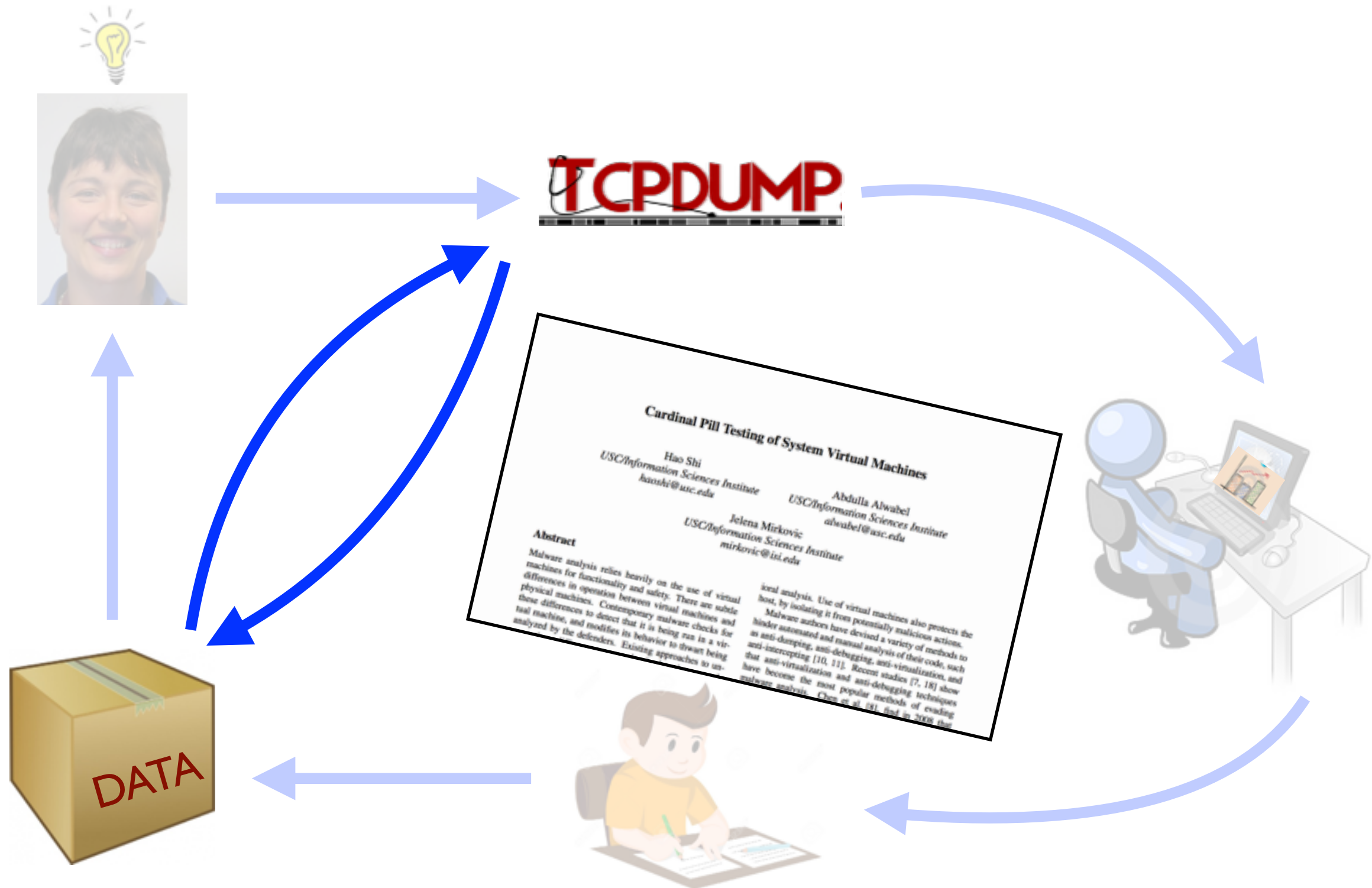




# How We Work ...

- Compartmentalization is good for organizing ourselves ...
- ... but suboptimal for data collection

# Hoarder's Circle of Research



# Why Hoard?

- When data is on hand ...
  - ... you can *quickly execute* on ideas
  - ... you can develop *longitudinal* understanding
  - ... you can act on *unexpected events*
- Anecdotally: colleagues who hoard get more done than colleagues who use purpose-driven data collection

# Hoarding Examples

- E.g., CAIDA's skitter (ark) traceroute data
- E.g., ISI's Internet census data
- E.g., Paxson's myriad use of LBL connection summaries
- E.g., Bailey's dark net data
- E.g., UCSD's spam feeds

# Hoarding

- Do not blindly expend effort collecting every possible thing you can think to collect
- If data is *cheap and easy* to collect, save it.
- If you have to expend energy to setup a collection, keep it going
  - ... even if there is no clear reason
  - ... even though it will take a bit of (ongoing) effort



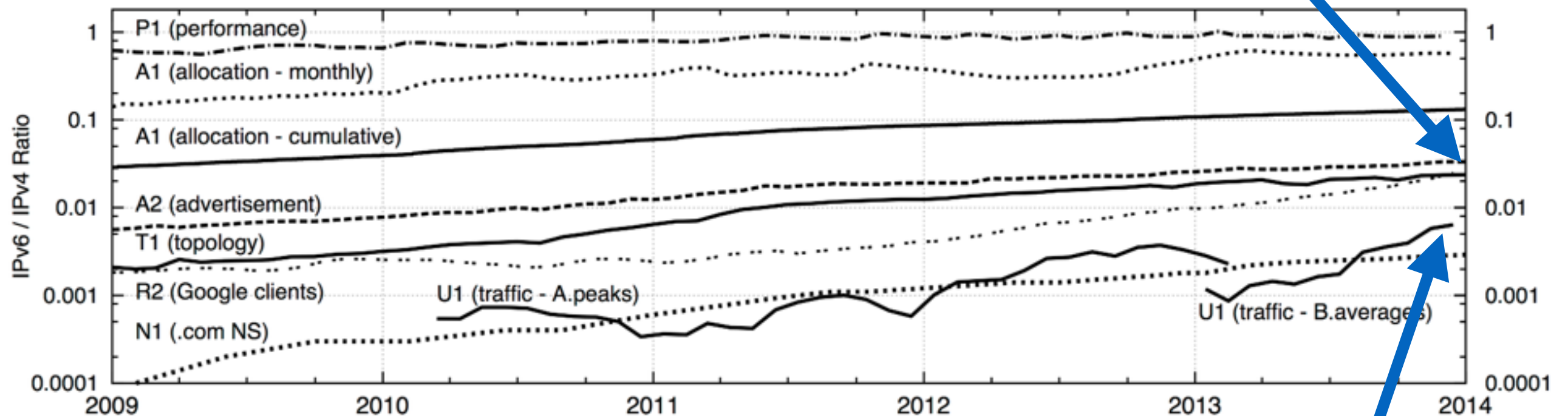
# Hoarding

- Do not collect data without ever looking at it
  - ... you'll some day find you have *NO* data!
  - ... simple automated analyses suffice
- Don't delegate hoarding to students
- Realize you'll still need to augment with purpose-driven data collection
- Realize hoarding is sometimes too expensive

# Embrace Heterogeneity

# Heterogeneous By Metric

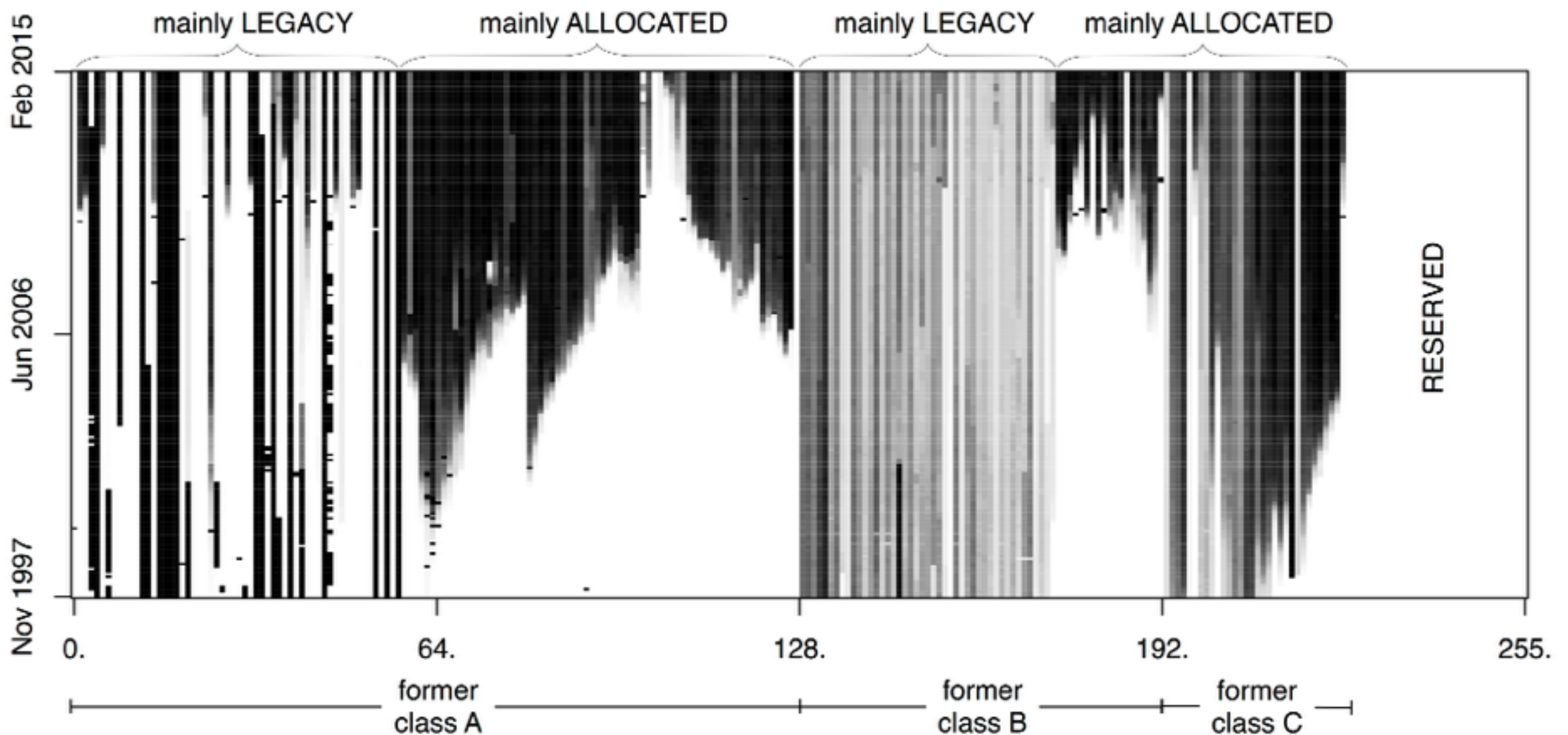
IPv6 Advertisement



IPv6 Traffic

Czyz, et. al.  
SIGCOMM 2014

# Heterogeneous By Time



# Heterogeneous By Path

Change	Both	Fwd	Rev	Flows	Affected
HICCUPS not capable	68	0	2	10360	0.68%
NAT	7704	0	0	10281	74.93%
ISN translation	924	178	0	10290	10.71%
IPID change	0	0	0	10290	0.00%
RCVWIN change	0	0	0	10290	0.00%
ECN IP add	26	0	0	10270	0.25%
ECN IP change	16	1342	48	10283	13.67%
ECN TCP add	16	0	0	10261	0.16%
ECN TCP change	19	46	0	10285	0.63%
MSS add	119	47	1036	10258	11.72%
MSS480 change	21	0	1132	10281	11.21%
MSS1460 change	1113	0	0	10275	10.83%
MSS1600 change	1105	157	0	10294	12.26%
SACK Permit changed	1	24	0	10123	0.25%
Timestamps add	12	0	0	10267	0.12%
Timestamps change	26	2	0	10279	0.27%
Window Scaling add	45	0	0	10265	0.44%
Window Scaling change	24	0	0	10279	0.23%
MPCAPABLE change	24	837	0	10267	8.39%
Exp. option change	20	884	0	10266	8.81%

# We All Know A Little

- We all accept heterogeneity ...



# How Much Is Enough?

## On the Marginal Utility of Network Topology Measurements

Paul Barford, Azer Bestavros, John Byers, Mark Crovella

*Abstract—*

The cost and complexity of deploying measurement infrastructure in the Internet for the purpose of analyzing its structure and behavior is considerable. Basic questions about the *utility* of increasing the number of measurements and measurement sites have not yet been addressed which has led to a “more is better” approach to wide-area measurement studies. In this paper, we step toward a more quantifiable understanding of the marginal utility of performing wide-area measurements in the context of Internet topology discovery. We characterize the observable topology in terms

simplicity in its internal components; for this reason, measurements made at network endpoints are especially attractive. An example of this approach is the use of `traceroute` [17] for the discovery of network connectivity and routing.

While `traceroute` is remarkably flexible and informative, it is an open question how useful `traceroute` is for uncovering topological information about the Internet. In this paper we study the use of `traceroute` as a tool for Internet topology discovery. We consider the common case.

Internet Measurement  
Workshop, 2001

# We All Know A Little

- We all accept heterogeneity ...
- ... right up until the moment we walk into the PC meeting!



# But, We Forget A Little

- Among my least favorite review comments ...
  - “data comes only from one university ...”
  - “authors watch only one mobile carrier ...”
  - “the data only encompasses one day ...”
  - “the user population is ‘small’ ...”
- “... so, *the study is not representative.*”

# On Representativeness

- These review comments are simultaneously ...
  - ... correct
  - ... vacuous
  - ... hypocritical
- *We need to stop beating each other up for conducting sound empirical work!*

# Embrace Reappraisal

# Validating Results

A large percentage of the Internet measurement studies currently published are not verified by the community due to the inability of researchers to access others' data and measurement/analysis tools.



Allman, et.al.  
PAM 2002

# So You Wanna Reappraise?

## A First Look at Modern Enterprise Traffic\*

Ruoming Pang<sup>†</sup>, Mark Allman<sup>‡</sup>, Mike Bennett<sup>¶</sup>, Jason Lee<sup>¶</sup>, Vern Paxson<sup>‡,¶</sup>, Brian Tierney<sup>¶</sup>  
<sup>†</sup>Princeton University, <sup>‡</sup>International Computer Science Institute,  
<sup>¶</sup>Lawrence Berkeley National Laboratory (LBNL)

### Abstract

While wide-area Internet traffic has been heavily studied for many years, the characteristics of traffic *inside* Internet enterprises remain almost wholly unexplored. Nearly all of the studies of enterprise traffic available in the literature are well over a decade old and focus on individual LANs rather than whole sites. In this paper we give a broad overview of internal enterprise traffic

within a site [2]. The only broadly flavored look at traffic within modern enterprises of which we are aware is the study of OSPF routing behavior in [21]. Our aim is to complement that study with a look at the make-up of traffic as seen at the packet level within a contemporary enterprise network.

One likely reason why enterprise traffic has gone unexamined for so long is that it is technically difficult to mea-

examine briefly (or not at all) in this paper. Towards this end, we are releasing anonymized versions of our traces to the community [1].

[1] LBNL Enterprise Trace Repository, 2005.  
<http://www.icir.org/enterprise-tracing/>.

IMC 2005

# So You Wanna Reappraise?

Share  
your data  
so the  
community  
can validate  
your  
results.

Give me  
your data so I  
can use it for  
my stuff!



# Why I Know What I Know

- I believe in heavy tails ...
  - ... not because Leland told me
  - ... or because Willinger told me
  - ... or because Feldmann told me
  - ... or because Crovella told me
  - ... or because Paxson told me
- ... but because *they all told me!*

# Why I Know What I Know

- I believe in heavy tails ...
  - ... but not because a bunch of researchers *re-produced* numbers from some dataset
  - ... but, rather, because a bunch of researchers arrived at the same insights across datasets!



# Reappraisal

- Insights are stronger when they come from multiple ...
  - researchers
  - datasets
  - vantage points
  - methodologies
- Yet, in general we do little reappraisal

# We Hate Reappraisal

- We are novelty junkies ...
  - ... reappraisal is boring
  - ... so, reappraisal is unwelcome
  - ... even when it is!
- Reappraisal is not viewed as on par with new contributions
- But, it does strengthen our understanding
- *We should stop beating up sound empirical work!*

# Embrace Risk

# Ground Truth

- Ground truth in Internet data is elusive
- This makes understanding the efficacy of inferences difficult (or impossible)

# Risky Business

## Capturing Ghosts: Predicting the Used IPv4 Space by Inferring Unobserved Addresses

Sebastian Zander  
CAIA, Swinburne University of  
Technology  
Melbourne, Australia  
szander@swin.edu.au

Lachlan L. H. Andrew  
Faculty of IT,  
Monash University  
Melbourne, Australia  
lachlan.andrew@monash.edu

Grenville Armitage  
CAIA, Swinburne University of  
Technology  
Melbourne, Australia  
garmitage@swin.edu.au

### ABSTRACT

The pool of unused routable IPv4 prefixes is dwindling, with less than 4% remaining for allocation at the end of June 2014. Yet the adoption of IPv6 remains slow. We demonstrate a new capture-recapture technique for improved estimation of the size of “IPv4 reserves” (allocated yet unused IPv4 addresses or routable prefixes) from multiple incomplete data sources. A key contribution of our approach is the plausible estimation of both observed and unobserved-yet-active (ghost) IPv4 address space. This signifi-

IPv4 address markets, requires plausible estimates of actual IPv4 address use – particularly the efficiency with which allocated prefixes are filled with actively-used addresses. Ideally, our estimation techniques should also help the community track progressive exhaustion once all routable IPv4 prefixes are allocated. Prior studies that, among other things, analysed IPv4 space growth [2–4] and a port scan census from 2012 [5] used mainly active probing (“pinging”). Yet pinging alone will under-count, as many hosts do not respond or their responses are filtered (e.g., by

IMC 2014

# Accepting Risk

- We should be willing to accept some risk
- ... or we will never make progress

# Risk Mitigation

- We should expect *some* validation ...
  - lab-based experiments
  - using alternate methodologies to cross-check results
  - small amounts of ground truth

# Do Not Embrace Sloppiness

- Risk from sloppy work is always unacceptable



# Risk Should Not Be ...



... a crutch  
for the lazy



...a beating  
for the diligent

# Embrace The Whole Bag Of Tricks

# Bag Of Tricks

Insight	Calibration	Hoarding
Heterogeneity	Reappraisal	Risk

- Individual notions have merit
- Also, feed on each other

# Comment #1

Insight

Calibration

Hoarding

- These provide a foundation for solid individual pieces of work

# Comment #2

Heterogeneity

Reappraisal

Risk

- These let us start thinking beyond individual pieces of solid work
- Rather, we lean on a *community-based body of work*

# Composition #1

Calibration

Risk

- Building confidence that we know the difference between network properties and measurement artifacts reduces risk

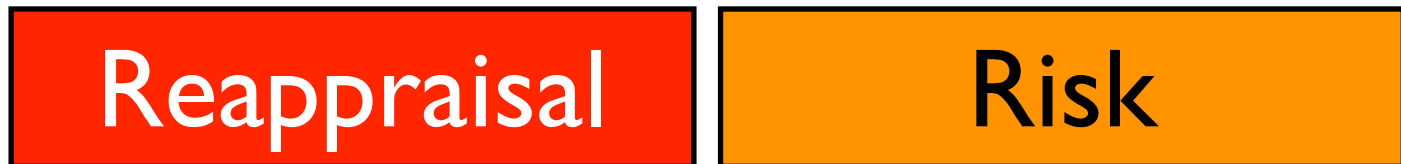
# Composition #2

Calibration

Hoarding

- Hoarding with a given measurement apparatus leads to easier calibration

# Composition #3



- Fostering reappraisal reduces the overall risk of any particular study being wrong



# Composition #4

Insight

Reappraisal

- Focusing on insight makes reappraisal easier

# Composition #5

Hoarding

Reappraisal

- Pack ratting data leads to more opportunities to reappraise
  - data is on hand
- Reduces the cost to, for instance, make reappraisal a first student project

# Composition #6

Heterogeneity

Reappraisal

- Reappraisal allows for a variety of perspectives from a variety of researchers
- Ultimately we have stronger and “more representative” insights

# Keeping Our Eye on the Ball



*“Some of the old songs I sing often, because they help me to reflect on where I've been and that's important for me to do - so I don't lose track of where I am going.”*  
—Johnny Cash

# Thanks!

- Jelena Mirkovic & the PAM organizers
- Mike Bailey, Paul Barford, Rob Beverly, Ethan Blanton, kc claffy, Wes Eddy, Sally Floyd, John Heidemann, Christian Kreibich, Boris Nechaev, Shawn Ostermann, Craig Partridge, Vern Paxson, Matt Roughan, Walter Willinger, ...



# Questions? Comments?



Mark Allman, [mallman@icir.org](mailto:mallman@icir.org)  
<http://www.icir.org/mallman/>