# Big Data and the Credibility Crisis

Victoria Stodden
School of Graduate and Library Science
University of Illinois at Urbana-Champaign

Advisory Committee for the Education and Human Resources Directorate of the NSF
Washington, D. C. (remotely)
Nov 6, 2014

# Framing: The Scientific Method

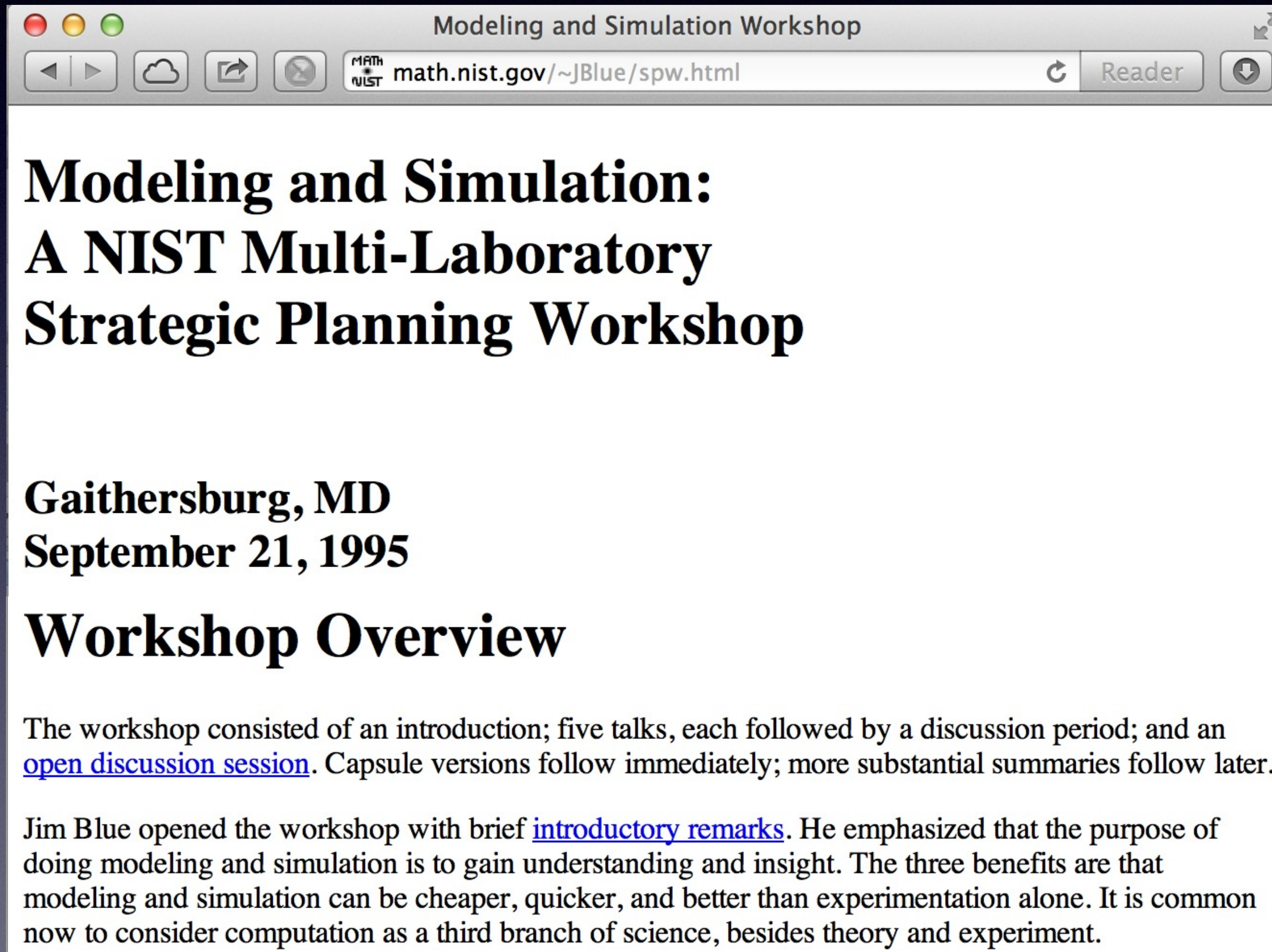Traditionally two branches to the scientific method:

- Branch 1 (deductive): mathematics, formal logic,

- Branch 2 (empirical): statistical analysis of controlled experiments.

Now, new branches due to technological changes?

- Branch 3,4? (computational): large scale simulations / data driven computational science.

Argument: computation presents only a *potential* third/fourth branch of the scientific method (Donoho et al 2009).

# New Paradigms for Discovery?



"It is common now to consider computation as a third branch of science, besides theory and experiment."

"This book is about a new, fourth paradigm for science based on data-intensive computing."

# The Impact of Technology

1. *Big Data / Data Driven Discovery*: high dimensional data, $p \gg n$,

2. *Computational Power*: simulation of the complete evolution of a physical system, systematically varying parameters,

3. Deep intellectual contributions now encoded only in *software*.



CSHL Keynote; Dr. Lior Pachter, UC Berkeley

The software contains "ideas that enable biology..."
*Stories from the Supplement*, 2013.
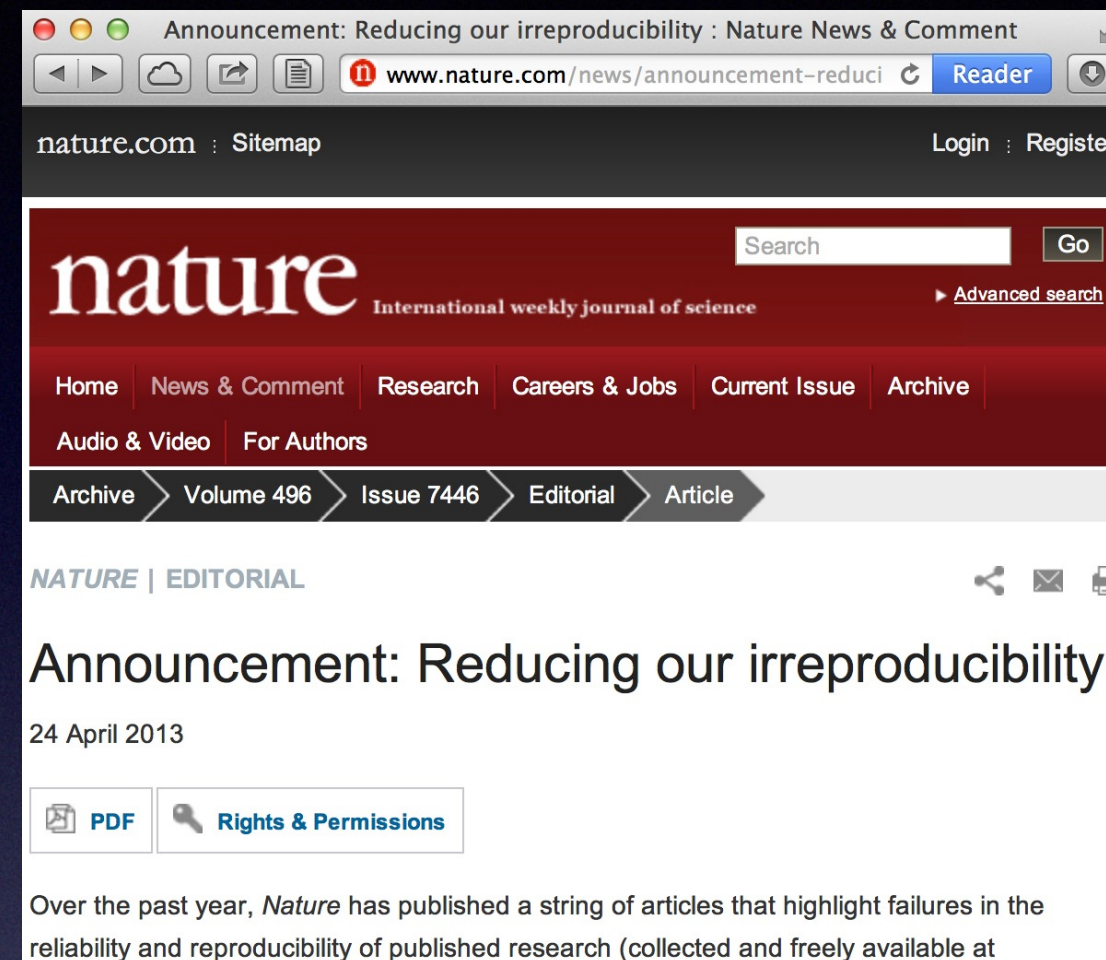
# The Ubiquity of Error

The central motivation for the scientific method is to root out error:

- Deductive branch: the well-defined concept of the proof,

- Empirical branch: the machinery of hypothesis testing, appropriate statistical methods, structured communication of methods and protocols.
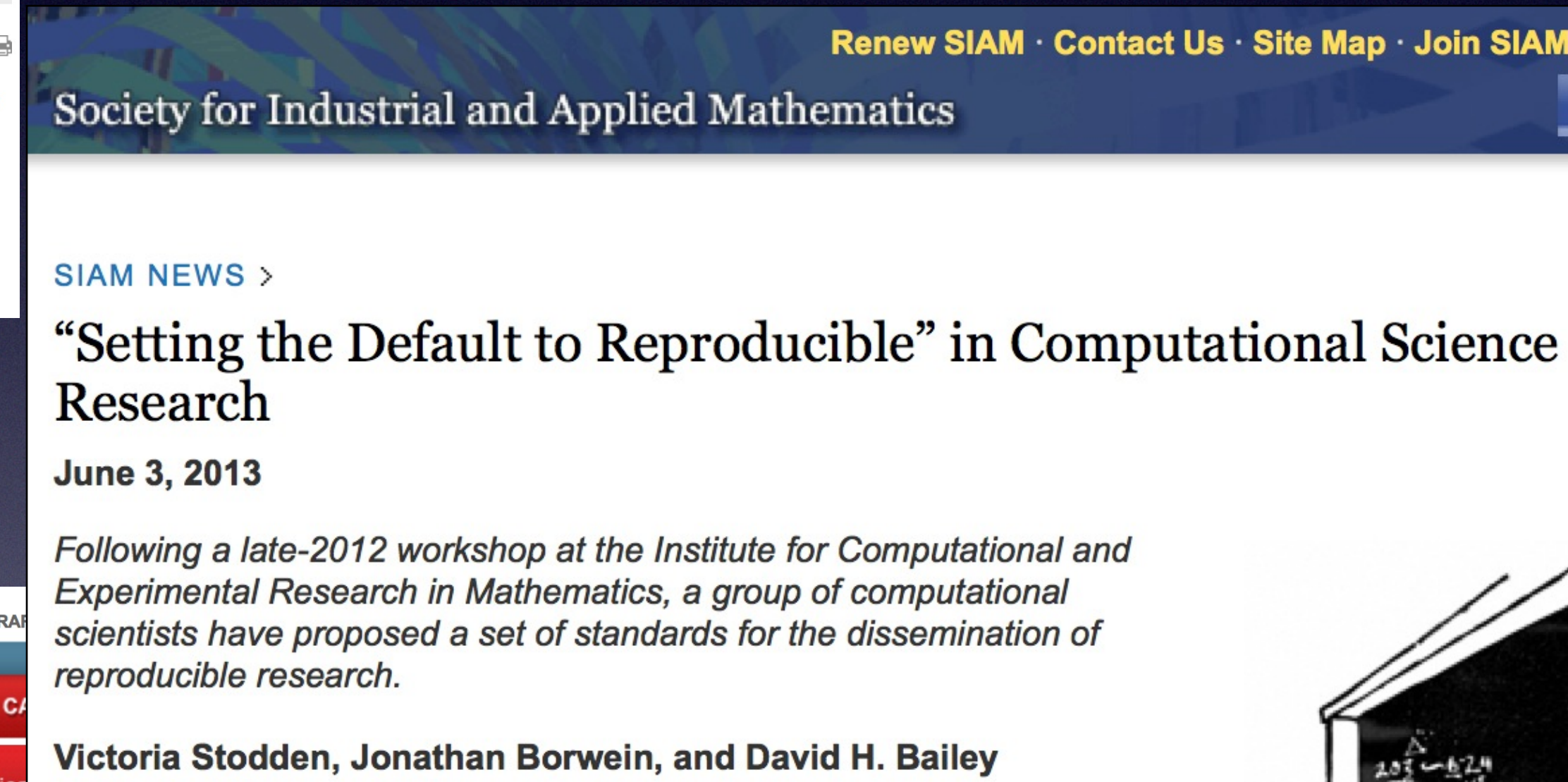
**Claim**: Computation presents only a potential third/fourth branch of the scientific method (Donoho, Stodden, et al. 2009), until the development of comparable standards.
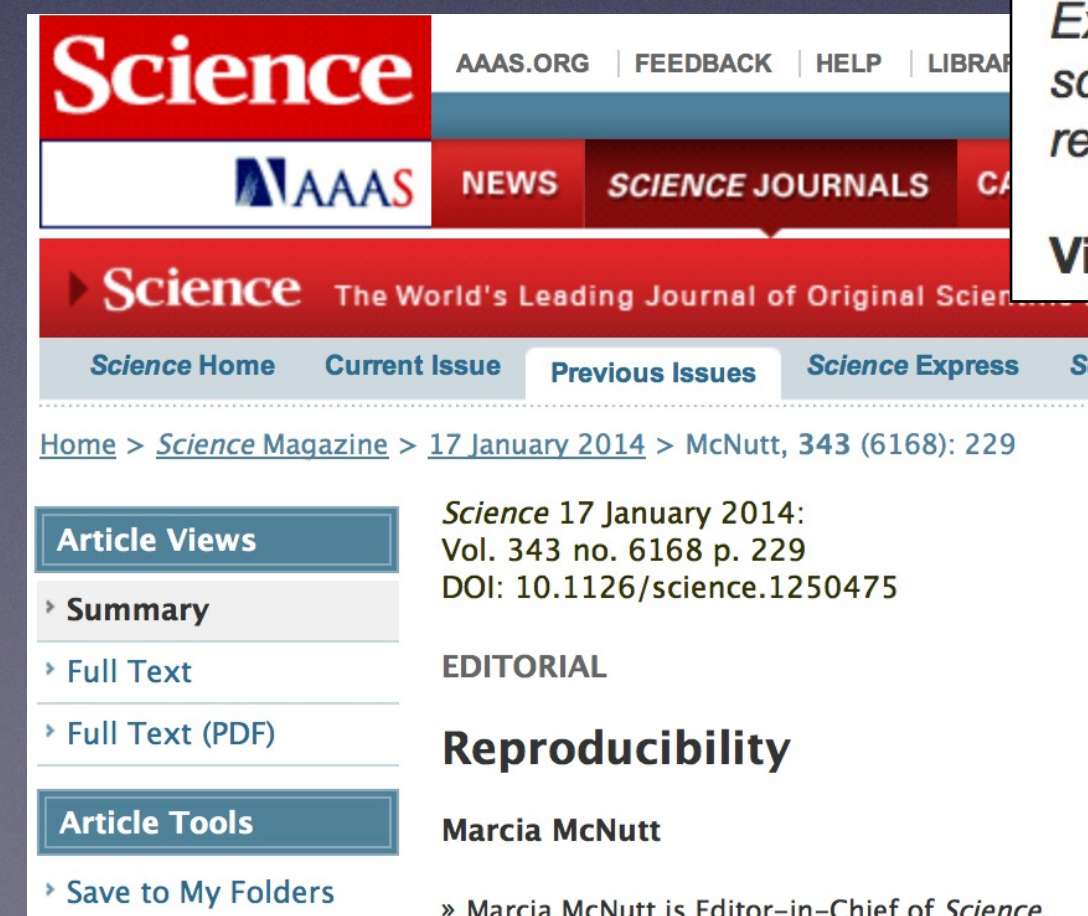
# Parsing Reproducibility

"Empirical Reproducibility"

"Computational Reproducibility"

"Statistical Reproducibility"

V. Stodden, IMS Bulletin (2013)

# Computational Reproducibility

"Really Reproducible Research" pioneered by Stanford Professor Jon Claerbout:

"The idea is: An article about computational science in a scientific publication is *not* the scholarship itself, it is merely *advertising* of the scholarship. The actual scholarship is the complete ... set of instructions [and data] which generated the figures."

paraphrased by David Donoho, 1998.

# Reproducibility is a Statistical Issue

- False discovery, chasing significance, p-hacking (Simonsohn 2012), overuse and mis-use of p-values,

- Multiple testing, file drawer problem, sensitivity analysis, poor reporting/tracking practices,

- Data preparation, treatment of outliers,

- Poor statistical methods (nonrandom sampling, inappropriate methods,..)

- Model robustness to parameter changes and data perturbations,

- Investigator bias toward previous findings; conflicts of interest.

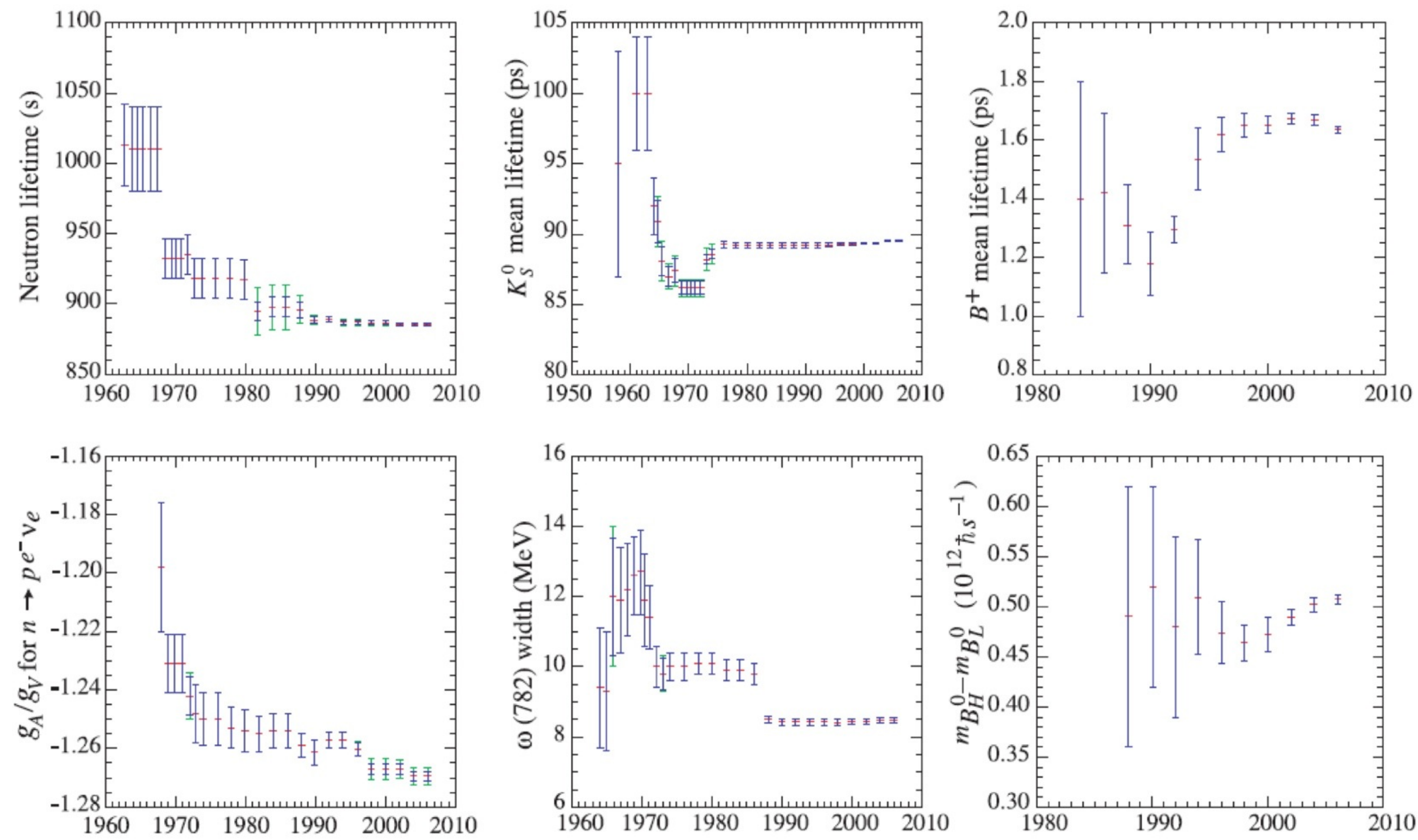# Experimental Bias



**Experimental biases:**

Figure 2: Historical record of values of some particle properties published over time, with quoted error bars (Particle Data Group).
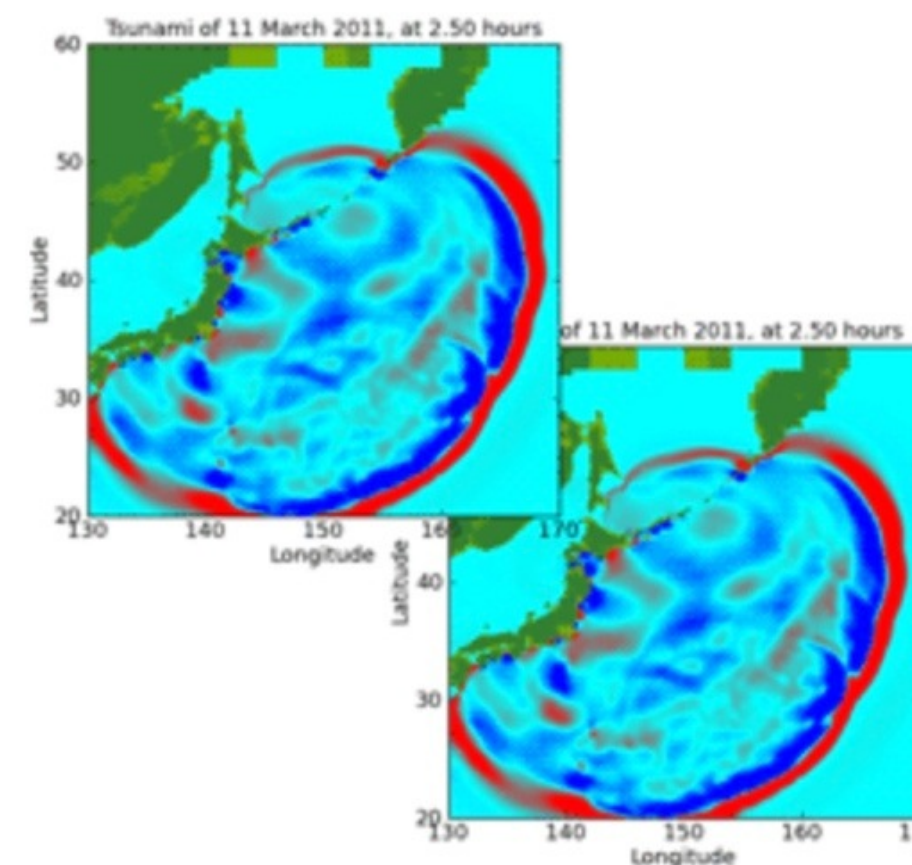
Figure courtesy of James Berger

# ICERM Workshop

# ICERM Workshop Report

## Setting the Default to Reproducible

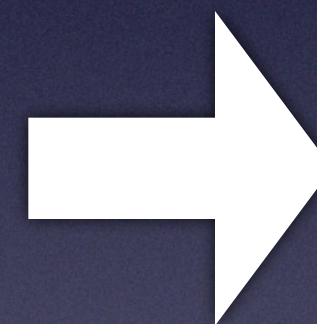### Reproducibility in Computational and Experimental Mathematics

Developed collaboratively by the ICERM workshop participants[1]

Compiled and edited by the Organizers

V. Stodden, D. H. Bailey, J. Borwein, R. J. LeVeque, W. Rider, and W. Stein

**Abstract**

Science is built upon foundations of theory and experiment validated and improved through open, transparent communication. With the increasingly central role of computation in scientific discovery this means communicating all details of the computations needed for others to replicate the experiment, i.e. making available to others the associated data and code. The "reproducible research" movement recognizes that traditional scientific research and publication practices now fall short of this ideal, and encourages all those involved in the production of computational science – scientists who use computational methods and the institutions that employ them, journals and dissemination mechanisms, and funding agencies – to facilitate and practice really reproducible research.

## Set the Default to "Open"

**Reproducible Science in the Computer Age.** Conventional wisdom sees computing as the "third leg" of science, complementing theory and experiment. That metaphor is outdated. Computing now pervades all of science. Massive computation is often required to reduce and analyze data; simulations are employed in fields as diverse as climate modeling and astrophysics. Unfortunately, scientific computing culture has not kept pace. Experimental researchers are taught early to keep notebooks or computer logs of every work detail: design, procedures, equipment, raw results, processing techniques, statistical methods of analysis, etc. In contrast, few computational experiments are performed with such care. Typically, there is no record of workflow, computer hardware and software configuration, or parameter settings. Often source code is lost. While crippling reproducibility of results, these practices ultimately impede the researcher's own productivity.

**The State of Experimental and Computational Mathematics.** Experimental mathematics[1]—application of high-performance computing technology to research questions in pure and applied mathematics, including

"It says it's sick of doing things like inventories and payrolls, and it wants to make some breakthroughs in astrophysics."

ScienceCartoonsPlus.com

physicists, legal scholars, journal editors, and funding agency officials representing academia, government labs, industry research, and all points in between. While

Renew SIAM · Contact Us · Site Map · Join SIAM

Society for Industrial and Applied Mathematics

SIAM NEWS >

## "Setting the Default to Reproducible" in Computational Science Research

June 3, 2013

*Following a late-2012 workshop at the Institute for Computational and Experimental Research in Mathematics, a group of computational scientists have proposed a set of standards for the dissemination of reproducible research.*

**Victoria Stodden, Jonathan Borwein, and David H. Bailey**

# Statistical Issues in Software

The challenge of reproducible computational science:

- shareable encoding of good statistical practices,

- permitting independent verification and comparison,

- extending statistical notions of integrity to statistical software practices,

Foundational research needed..

# Supporting Computational Science

- Dissemination Platforms:

  ResearchCompendia.org     IPOL     Madagascar

  MLOSS.org     thedatahub.org     nanoHUB.org

  Open Science Framework     RunMyCode.org

- Workflow Tracking and Research Environments:

  VisTrails     Kepler     CDE     IPython Notebook

  Galaxy     GenePattern     Paper Mâché

  Sumatra     Taverna     Pegasus

- Embedded Publishing:

  Verifiable Computational Research     SOLE     knitR

  Collage Authoring Environment     SHARE     Sweave

# Open Science from the Whitehouse

- Feb 22, 2013: <u>Executive Memorandum</u> directing federal funding agencies to develop plans for public access to data and publications.

- May 9, 2013: <u>Executive Order</u> directing federal agencies to make their data publicly available.

- July 29, 2014: <u>Notice of Request for Information</u> "Strategy for American Innovation"

# Request for Input: "Strategy for American Innovation"

- "to guide the Administration's efforts to promote lasting economic growth and competitiveness through policies that support transformative American innovation in products, processes, and services and spur new fundamental discoveries that in the long run lead to growing economic prosperity and rising living standards."

- "(11) Given recent evidence of the irreproducibility of a surprising number of published scientific findings, how can the Federal Government leverage its role as a significant funder of scientific research to most effectively address the problem?"

# Science Policy in Congress

- America COMPETES due to be reauthorized, drafting underway.

- Sensenbrenner introduced "Public Access to Science," Sept 19, 2013.

- Hearing on Research Integrity and Transparency by the House Science, Space, and Technology Committee (March 5, 2013).

- *Reproducibility cannot be an unfunded mandate.*

# NAS Data Sharing Report



- <u>Sharing Publication-Related Data and Materials: Responsibilities of Authorship in the Life Sciences</u>, (2003)

- "Principle 1. Authors should include in their publications the data, algorithms, or other information that is central or integral to the publication—that is, whatever is necessary to support the major claims of the paper and would enable one skilled in the art to verify or replicate the claims."

# National Science Board Report



"Digital Research Data Sharing and Management," December 2011.

http://www.nsf.gov/nsb/publications/2011/nsb1124.pdf

# We need:

Standards for reproducibility of computational findings:

1. data access, software access, persistent linking to publications.

2. innovation around data and code access for privacy protection and scale.

3. robust methods, producing stable results, emphasis on reliability and reproducibility.

Example: Google Flu Trends results: worked at first, but what happened? (Lazer et al. "The Parable of Google Flu: Traps in Big Data Analysis" *Science*, 2014)

# References

"Toward Reproducible Computational Research: An Empirical Analysis of Data and Code Policy Adoption by Journals," PLoS ONE, June 2013

"Reproducible Research," guest editor for Computing in Science and Engineering, July/August 2012.

"Reproducible Research: Tools and Strategies for Scientific Computing," July 2011.

"Enabling Reproducible Research: Open Licensing for Scientific Innovation," 2009.

available at http://www.stodden.net

# Credibility Crisis