

# Introduction to Reproducible Research

Roger D. Peng  
*@rdpeng, @simplystats, simplystatistics.org*

Biostatistics 140.776

How do you know if a data analysis is successful?

When has a data analysis failed?

# Parable

ARTICLES

• Retracted •

nature  
medicine

## Genomic signatures to guide the use of chemotherapeutics

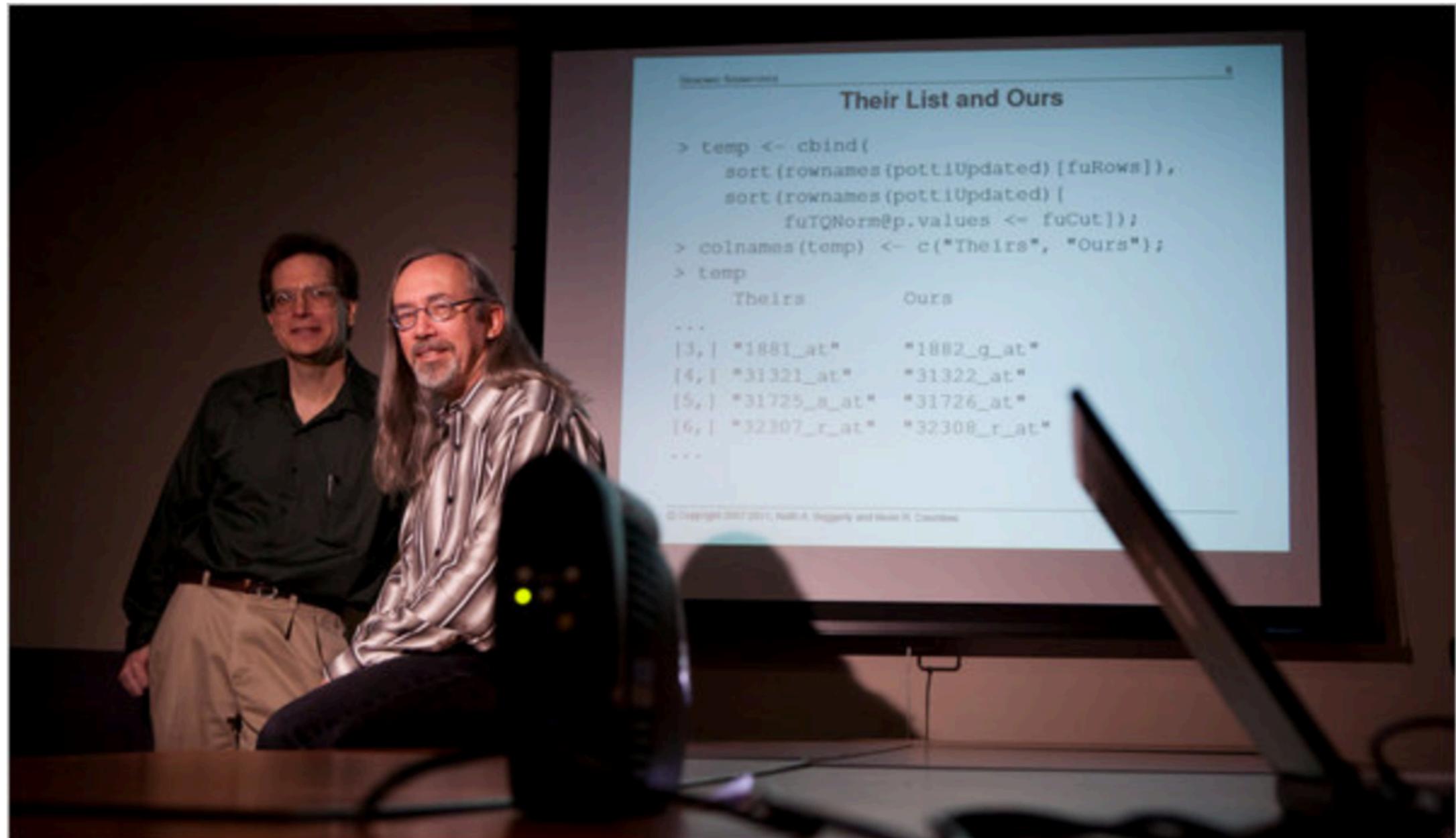
Anil Potti<sup>1,2</sup>, Holly K Dressman<sup>1,3</sup>, Andrea Bild<sup>1,3</sup>, Richard F Riedel<sup>1,2</sup>, Gina Chan<sup>4</sup>, Robyn Sayer<sup>4</sup>,  
Janiel Cragun<sup>4</sup>, Hope Cottrill<sup>4</sup>, Michael J Kelley<sup>2</sup>, Rebecca Petersen<sup>5</sup>, David Harpole<sup>5</sup>, Jeffrey Marks<sup>5</sup>,  
Andrew Berchuck<sup>1,6</sup>, Geoffrey S Ginsburg<sup>1,2</sup>, Phillip Febbo<sup>1-3</sup>, Johnathan Lancaster<sup>4</sup> &  
Joseph R Nevins<sup>1-3</sup>

# Deception at Duke

The image is a screenshot of a YouTube video player. At the top, a red banner features a stopwatch and the text "60 MINUTES". Below this is a navigation menu with links: HOME, UP NEXT, 60 OVERTIME, NEWSMAKERS, POLITICS, SCIENCE, BUSINESS, and ENTERTA. The main video frame shows a man in a suit standing in front of a backdrop. The backdrop has the "Deception At Duke" logo on the left and text on the right that reads "Produced By Kyra Darnton". The video player interface includes a progress bar at 0:52 / 13:46, a "SHARE" button, and social media sharing options: "23 Comments", "Share this Video:", "Recommend" (473), "Tweet" (49), and a red share icon (363). Below the player, the video title "Deception at Duke" is displayed, along with the date and time "February 12, 2012 4:00 PM" and a description: "Were some cancer patients at Duke University given experimental treatments based on fabricated data? Scott Pelley reports."

# “Rock Star” Statisticians

How Bright Promise in Cancer Testing Fell Apart



Michael Stravato for The New York Times

# Brief Summary of Problems

- Off-by-one table row labels
- Inadvertent switching of outcome labels
- Duplicated observations
- Genes identified not on microarray used
- Completely arbitrary statistical formulas used

# Lessons?



# Institute of Medicine Committee

REPORT BRIEF  MARCH 2012

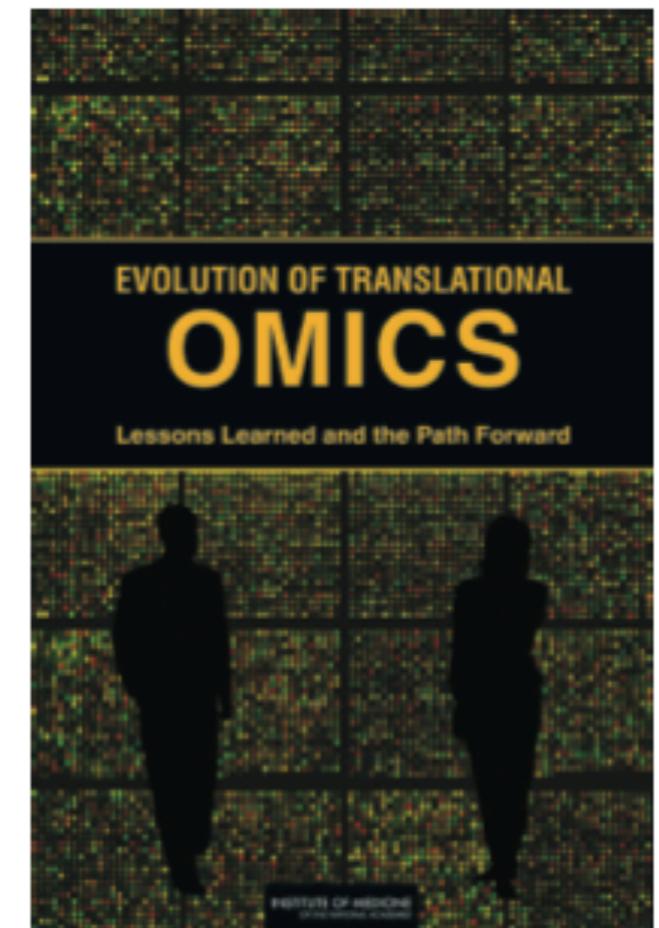
INSTITUTE OF MEDICINE  
OF THE NATIONAL ACADEMIES

Advising the nation • Improving health

For more information visit [www.iom.edu/translationalomics](http://www.iom.edu/translationalomics)

## Evolution of Translational Omics

Lessons Learned and the  
Path Forward



# The IOM Report

- **Data/metadata** used to develop test should be made publicly available
- The **computer code** and fully specified computational procedures used or development of the omics-based test should be made available
- Ideally, the computer code that is released will **encompass all of the steps** of computational analysis, including all data preprocessing steps

# Replication and Reproducibility

- **Replication**

- Focuses on the validity of the *scientific claim*
- “Is this claim true?”
- Ultimate standard for scientific evidence
- New investigators, data, analytic methods, labs, instruments, etc.
- Important in studies that can impact policy or regulation

- **Reproducibility**

- Focuses on the quality of the *data analysis*
- “Can we trust this analysis?”
- A minimum standard
- New investigators, same data, same methods
- Important when replication is impossible

# What's Wrong with Replication?

- Nothing, but...
- Some studies cannot be replicated
  - No time, opportunistic
  - No money
  - Unique
- **Reproducible Research:** Make analytic data and code available so that others may reproduce findings

# Upon Seeing Your Work...

Information Required

Minimum

Maximum

Do Nothing

Full Replication

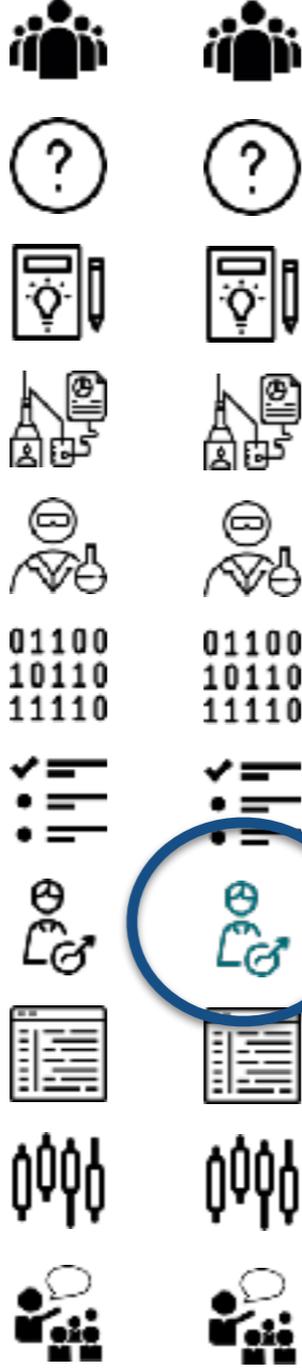
Reproducibility

Most  
common  
activity

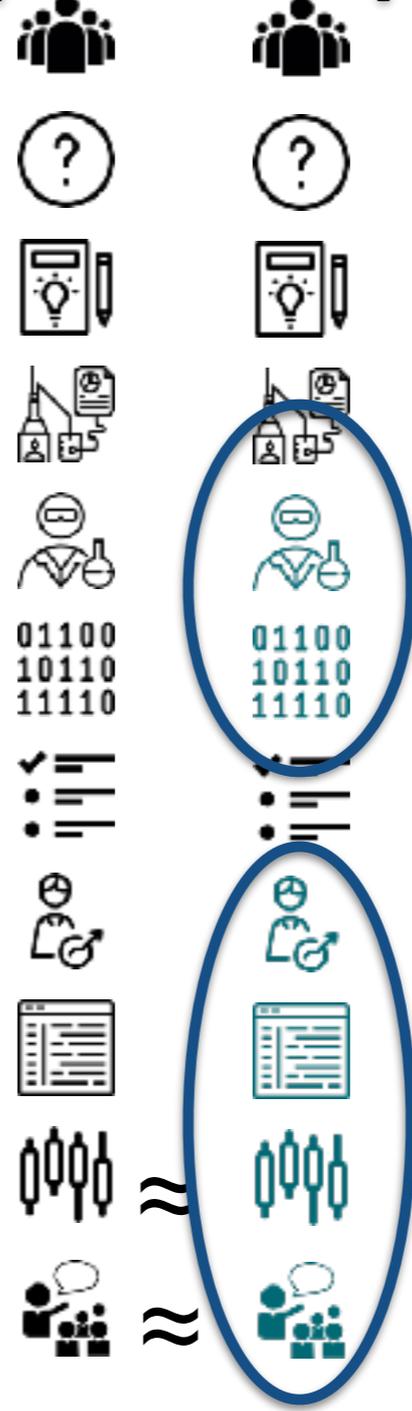
# Reproducible

# Replicable

- Population
- Question
- Hypothesis
- Experimental Design
- Experimenter
- Data  
01100  
10110  
11110
- Analysis Plan
- Analyst
- Code
- Estimate
- Claim



Original study    New study



Original study    New study

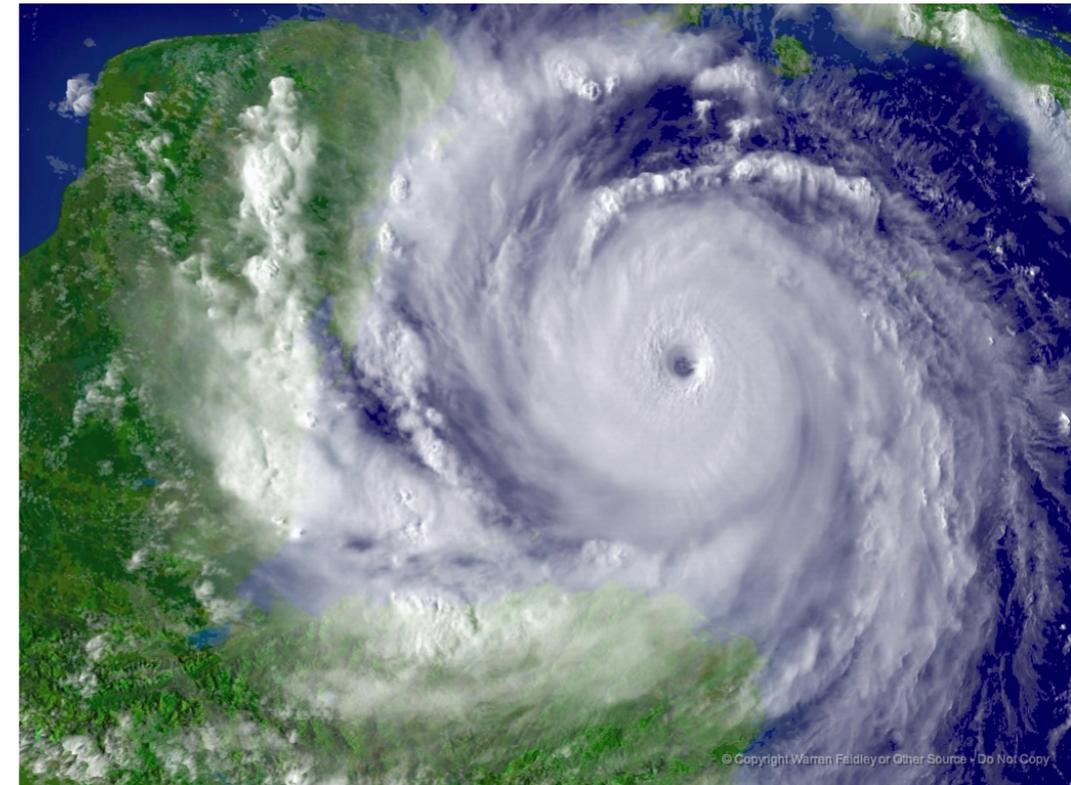
- Observed
- Missing
- Different Value
- Incorrectly reported

# Why Do We Need Reproducible Research?

- New technologies increasing data collection throughput
- Data are more complex and high dimensional
- Existing databases can be merged into new and bigger databases
- Computing power is greatly increased, allowing more sophisticated/complicated analyses
- For every field “X” there is a field “Computational X”

# Air Pollution and Health: A Perfect Storm?

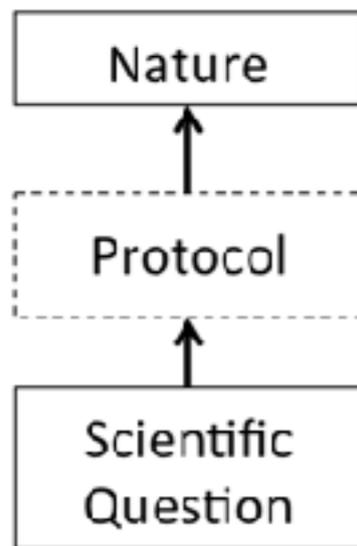
- Estimating small health effects in the presence of much stronger signals
- Results inform substantial policy decisions and affect many stakeholders
- EPA regulations can cost billions of dollars
- Complex statistical methods are needed and subjected to intense scrutiny



# The End Result

- Basic analyses can be difficult to describe
- Heavy computational requirements are thrust upon people without adequate training in statistics and computing
- Errors are more easily introduced into long and complex analysis pipelines
- Knowledge transfer is limited
- Complicated analyses cannot be trusted

# What is Reproducible Research?



Published  
Article

# What is Reproducible Research?

Author



Published  
Article

Nature



Protocol



Scientific  
Question



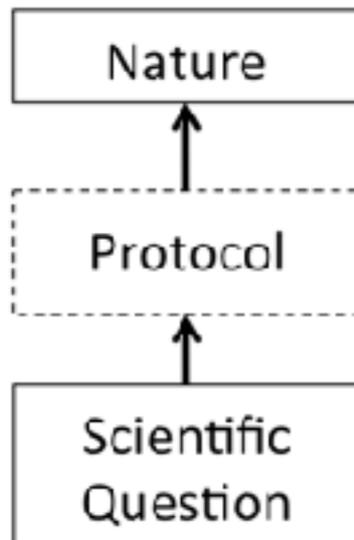
Reader

# What is Reproducible Research?

Author



Published  
Article

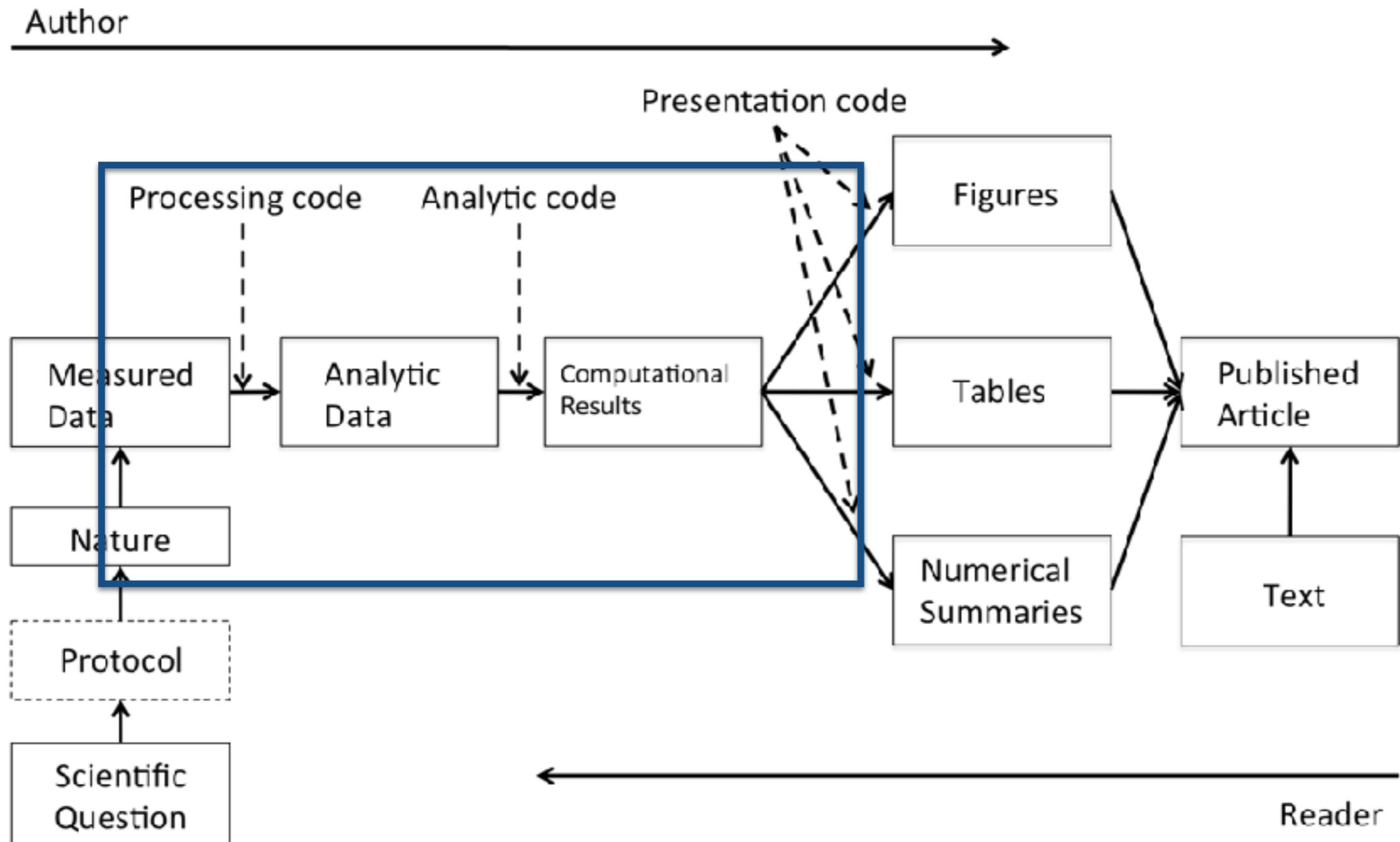


Express train to nature



Reader

# What is Reproducible Research?



# What is Reproducible Research?

- Analytic data are available
- Analytic (and preprocessing) code are available
- Documentation of code and data
- Standard means of distribution

# What is Reproducible Research?

- Authors
  - Want to make their research reproducible
  - Want tools for RR to make their lives easier (or at least not much harder)
- Readers
  - Want to reproduce (and perhaps expand upon) interesting findings
  - Want tools for RR to make their lives easier

# Challenges

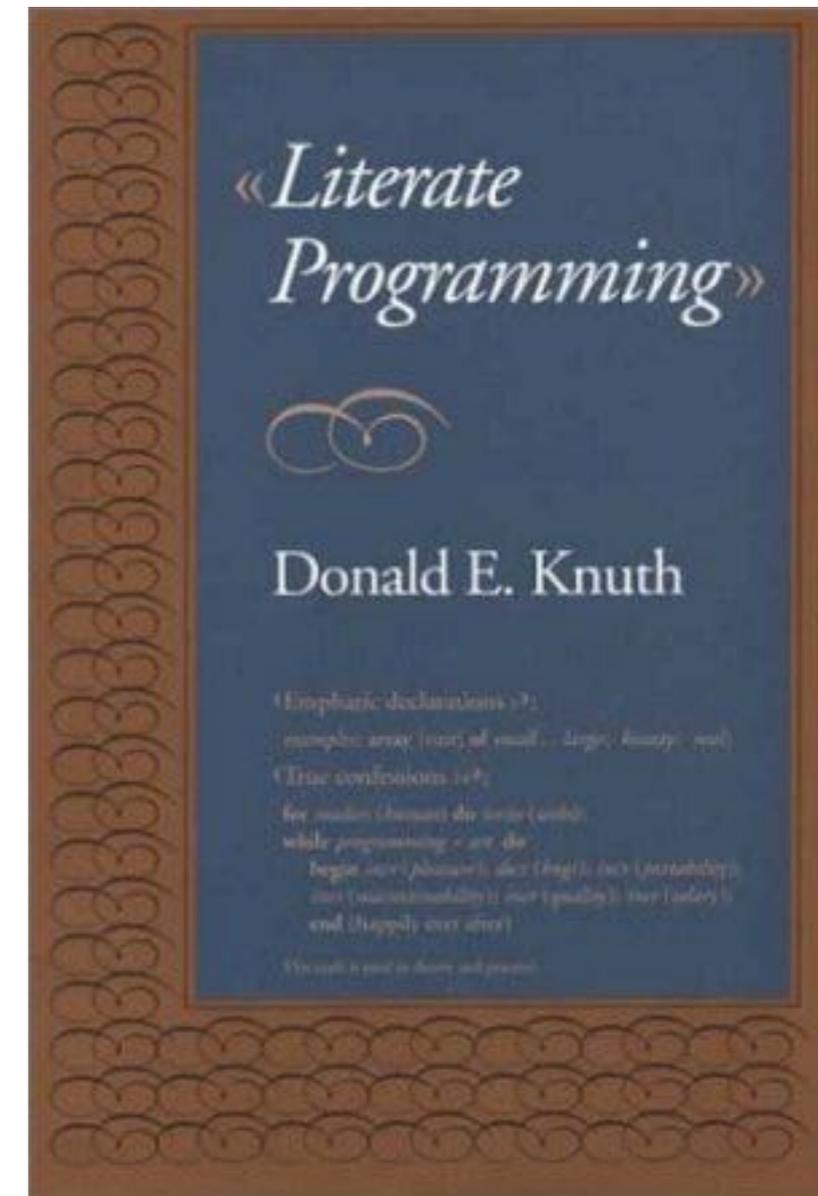
- Authors must undertake considerable effort to put data and results on the web (may not have resources like a web server)
- Readers must download data/results individually and piece together which data go with which code sections, etc.
- Readers may not have the same resources as authors
- Few tools to help authors/readers (although toolbox is growing!)

# Recent Developments

- **Software:** Jupyter Notebooks, knitr, markdown, LONI, Galaxy
- **Repositories:** GitHub, NCBI, ICPSR, Dataverse, Open Science Framework, *Google Dataset Search*
- **Policy:** *Science*, *Nature*, *PLOS ONE*, OSTP, NIH

# Literate Statistical Programming

- An article/report is a stream of text and code
- Analysis code is divided into text and code “chunks”
- Each code chunk loads data and computes results
- Presentation code formats results (tables, figures, etc.)
- Article text explains what is going on
- Literate programs can be **weaved** to produce human-readable documents and **tangled** to produce machine-readable documents
- See *Literate Programming* by Donald Knuth



# Literate Statistical Programming

- Literate programming is a general concept that requires
  - A documentation language (human readable)
  - A programming language (machine readable)
- Sweave uses LaTeX and R as the documentation and programming languages
- Sweave was developed by Friedrich Leisch (member of the R Core) and is maintained by R core
- Main web site: <http://www.statistik.lmu.de/~leisch/Sweave>

# Literate Statistical Programming

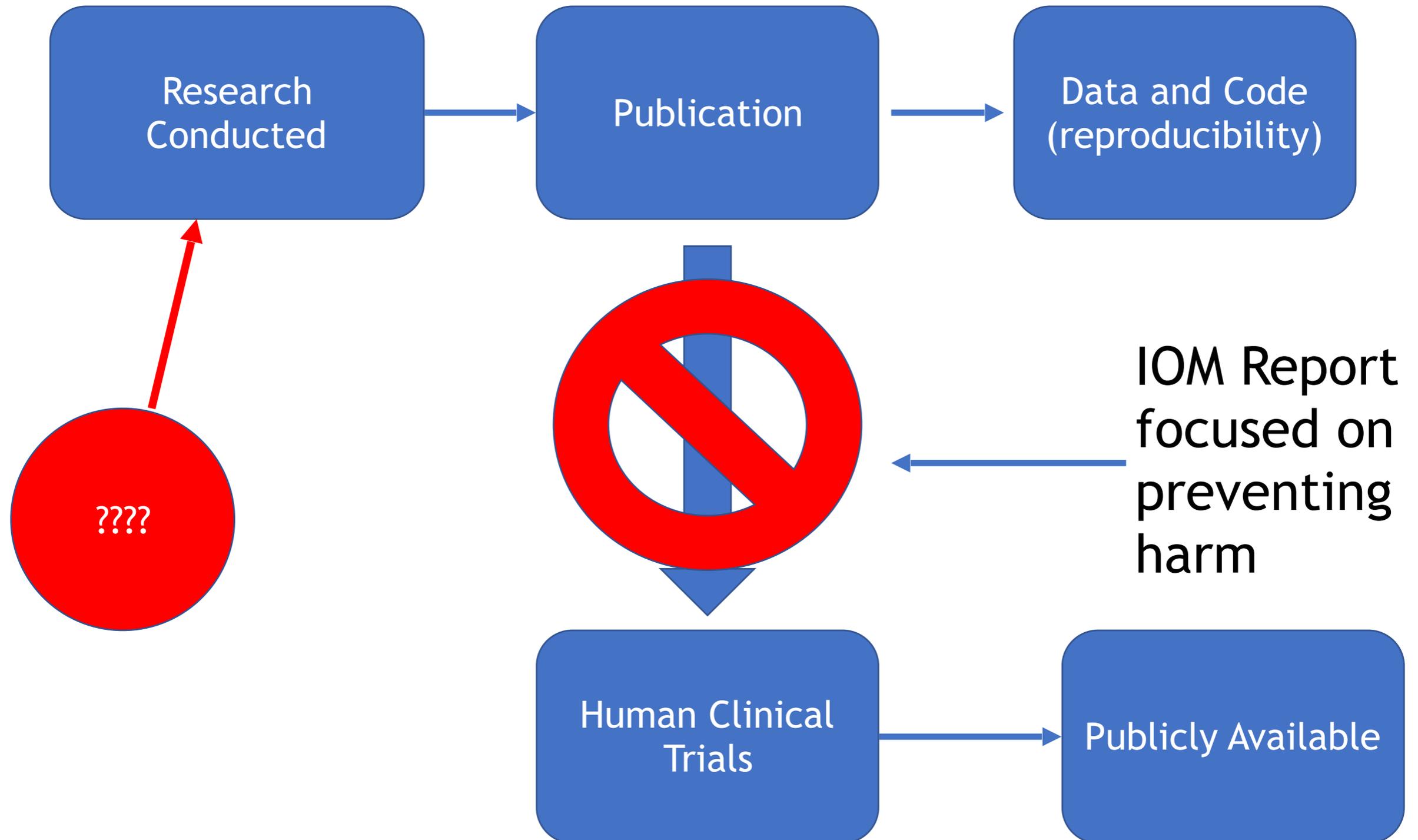
- knitr is package that brings together many features added on to Sweave to address limitations
- knitr uses the R programming language (although others are allowed) and variety of documentation languages
  - LaTeX, Markdown, HTML
- Built into RStudio pipeline
- See <http://yihui.name/knitr/>



# What Problem Does Reproducibility Solve?

- What we get
  - Transparency / Improved knowledge transfer
  - Data availability
  - Software / Methods
- What we do NOT get
  - Validity / Correctness of the analysis

# Where to Intervene?



# Lessons?

- Reproducibility
- Expertise and training
- Publication pressure; glamour journals
- Funding, conflicts of interest

# New Details Emerge (Jan 2015)

*“In raising  
these concerns,  
I have nothing  
to gain and  
much to lose.”*

*— Bradford  
Perez*



# The Perez Memo (cont'd)

“At this point, I believe that the situation is serious enough that **all further analysis should be stopped** to evaluate what is known about each predictor and it should be reconsidered which are appropriate to continue using and under what circumstances.... I would argue that at this point nothing...should be taken for granted. **All claims of predictor validations should be independently and blindly performed.**” [emphasis added]

*-Memo from Bradford Perez, April 2008*

# Lessons Learned

- Not all problems have a technical solution
- Analyses were not “too complicated” in that there was insufficient expertise; problems were readily recognized
- Lab/institute cultural problems lead to unwillingness to communicate obvious problems
- From the analyst perspective, a breakdown in communication is an early warning sign of potential data analytic problems
- Making analyses more reproducible likely would not have made much difference

Communicating a  
Data Analysis?



**FOUNTAINS OF WAYNE**

**JOHN COLTRANE** **GIANT STEPS**



Verse 1:

F Bb F  
It's you and me on a beach/ In nineteen-ninety eight  
Bb Dm Bb F  
Leaning into the breeze/ From the willows and rivermen grace  
Am Bb F C  
Are reborn in this place/ I'm assured the procedure is painless  
F Bb F  
The taxicab with no brakes/ around the mountain pass  
Bb Dm Bb F  
Keep your head in your hands/ if anybody asks what you mean when  
Am Bb F C  
You were picking a fight/ you were only complimenting the waitress

Chorus 1:

Am Bb  
Give us a room, with a mountain view:  
F C F C Bb  
A tiny cabana by the water, yeah by the water and I,  
C F Bb  
Got a rental for an hour or two, for a ride up the coast  
C F  
And a dip in the ocean.

(Repeat Intro)

Verse 2:

F Bb F  
The waterfront is alight/ with Citronella flame  
Bb Dm  
Tourists flash in the night/ from the grottoes and  
Bb F Am Bb F  
Gathering now on the heel-worn planks for a drunken  
C  
Form another, a mumble (?)  
F Bb F  
And lovers paddle a boat/ on the molten bay  
Bb Dm  
Peering into the the reeds/ on a ripple  
Bb F Am Bb  
And playing it cool/ in a bar by the pool  
F C  
With a Caribbean Kiss Amaretto

Verse 1:  
 F Bb  
 It's you and me on a beach/ In nineteen-n  
 Bb Dm  
 Leaning into the breeze/ From the willows  
 Am Bb F  
 Are reborn in this place/ I'm assured the  
 F Bb  
 The taxicab with no brakes/ around the mo  
 Bb Dm  
 Keep your head in your hands/ if anybody  
 Am Bb F  
 You were picking a fight/ you were only c

Chorus 1:  
 Am Bb  
 Give us a room, with a mountain view:  
 F C F C  
 A tiny cabana by the water, yeah by the w  
 C F Bb  
 Got a rental for an hour or two, for a ri  
 C F  
 And a dip in the ocean.

(Repeat Intro)

Verse 2:  
 F Bb  
 The waterfront is alight/ with Citronella  
 Bb Dm  
 Tourists flash in the night/ from the gro  
 Bb F Am Bb  
 Gathering now on the heel-worn planks for  
 C  
 Form another, a mumble (?)  
 F Bb F  
 And lovers paddle a boat/ on the molten b  
 Bb Dm  
 Peering into the the reeds/ on a ripple  
 Bb F Am Bb  
 And playing it cool/ in a bar by the pool  
 F C  
 With a Caribbean Kiss Amaretto

Allegro Impetuoso. 1

B♭-Klarinette in B  
 1, 2, 3, 4. Fagott.  
 Kontra-Fagott.  
 1, 2, 3, 4. Trompete in F.  
 1, 2, 3, 4. Posaune.  
 Pauken.  
 Maracas.  
 Orgel.  
 Pedal.  
 1, 2. Sopran.  
 1, 2. Alt.  
 Tenor.  
 Bariton.  
 Bass.  
 Klavierchor.  
 Sopran.  
 Alt.  
 Tenor.  
 Bass.  
 1. CHOR.  
 Sopran.  
 Alt.  
 Tenor.  
 Bass.  
 2. CHOR.  
 Sopran.  
 Alt.  
 Tenor.  
 Bass.  
 1. Violine.  
 2. Violine.  
 Bratsche.  
 Violoncell.  
 Kontrabaß.

Allegro Impetuoso.

Copyright 1935 by Universal Edition.

## GIANT STEPS

BY JOHN COLTRAN

C VERSION

FAST SWING

Chords: B $\natural$ 7, D7, G $\natural$ 7, B $\flat$ 7, E $\flat$ 7, A $\natural$ 7, D7

Chords: G $\natural$ 7, B $\flat$ 7, E $\flat$ 7, F $\sharp$ 7, B $\natural$ 7, F $\natural$ 7, B $\flat$ 7

Chords: E $\flat$ 7, A $\natural$ 7, D7, G $\natural$ 7, C $\sharp$ 7, F $\sharp$ 7

Chords: B $\natural$ 7, F $\natural$ 7, B $\flat$ 7, E $\flat$ 7, C $\sharp$ 7, F $\sharp$ 7

TO CODA

SOLOS / (11 CHORDS)

Chords: B $\natural$ 7, D7, G $\natural$ 7, B $\flat$ 7, E $\flat$ 7, A $\natural$ 7, D7, G $\natural$ 7, B $\flat$ 7

Chords: E $\flat$ 7, F $\sharp$ 7, B $\natural$ 7, F $\natural$ 7, B $\flat$ 7, E $\flat$ 7, A $\natural$ 7, D7, G $\natural$ 7

Chords: C $\sharp$ 7, F $\sharp$ 7, B $\natural$ 7, F $\natural$ 7, B $\flat$ 7, E $\flat$ 7, C $\sharp$ 7, F $\sharp$ 7

D.C. AL COI  
WITH REPE

CODA

Chords: F $\natural$ 7, B $\flat$ 7, E $\flat$ 7

The screenshot displays a DAW mixer interface with several tracks and a detailed channel strip on the right. The tracks include:

- 27 Stereo:** Utility, Gain, Master.
- 28 Group:** Chorus, Dry/Wet.
- Funky Mel (3 tracks):** Utility, Gain, Mixer, Speaker On.
- A Reverb:** Post.
- B Delay:** Post.
- Master:** Utility, Device On.

The channel strip on the right shows settings for each track, including gain levels (e.g., -3.3, 0, 0) and various processing options like 'Auto' and 'Off'. The mixer faders are visible at the bottom of each track.

The **Chorus** plugin interface features the following controls:

- Delay 1:** Highpass filter (FF1, FF2, FB), 1.00 kHz, 3.00 ms.
- Delay 2:** Modulation (Off, Fix, Mod), 7.00 ms.
- Modulation:** Amount (2.01 ms), Rate (20, 2.64 Hz).
- Polarity:** Inverted (-).
- Feedback:** 0.0%.
- Dry/Wet:** 17%.

The **EQ Eight** plugin interface shows a frequency response curve with the following settings:

- Mode:** Stereo.
- Adapt. Q:** On.
- Scale:** 100%.
- Gain:** 0.00 dB.
- Frequency:** 10.8 kHz.
- Gain:** 4.09 dB.
- Q:** 0.71.

The frequency response curve shows a boost in the high-frequency range.

The **Reverb** plugin interface includes the following parameters:

- Input Processing:** Lo Cut, Hi Cut, Spin.
- Global Quality:** High, 4.50 Hz, 0.70.
- Diffusion Network:** Size (100.00), Low (90.0 Hz), 0.75.
- Predelay:** 2.50 ms.
- Shape:** 0.50.
- Stereo:** 100.00.
- Decay Time:** 6.24 s.

# The Central Problem

**Data Analysis = ???**

# What's Next?

- Reproducibility is critical for *communicating* a data analysis
- One cannot sufficiently describe an analysis in words
- General consensus about its importance
- Infrastructure for making all research reproducible is not there yet, but things are ever improving

# How Do You Know if a Data Analysis is Successful?

- Reproducible
- Uses the best available statistical methods of analysis