

Code != Data

**Why Reproducible Research Needs Both
but Should Not Treat Them the Same**

Ian M. Mitchell

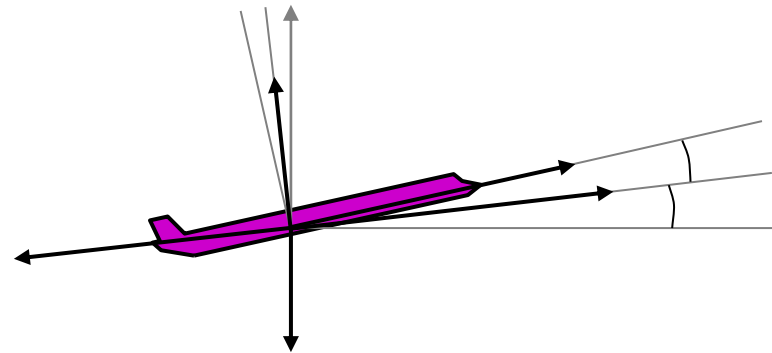
Department of Computer Science
University of British Columbia

research supported by
the Natural Science and Engineering Research Council of Canada



Outline

- Motivation: Some lessons in Irreproducible Research
- Computation and Data Science
- Reproducible Research: Changing the Culture
- Code \neq Data



June 2014



Ian Mitchell, University of British Columbia



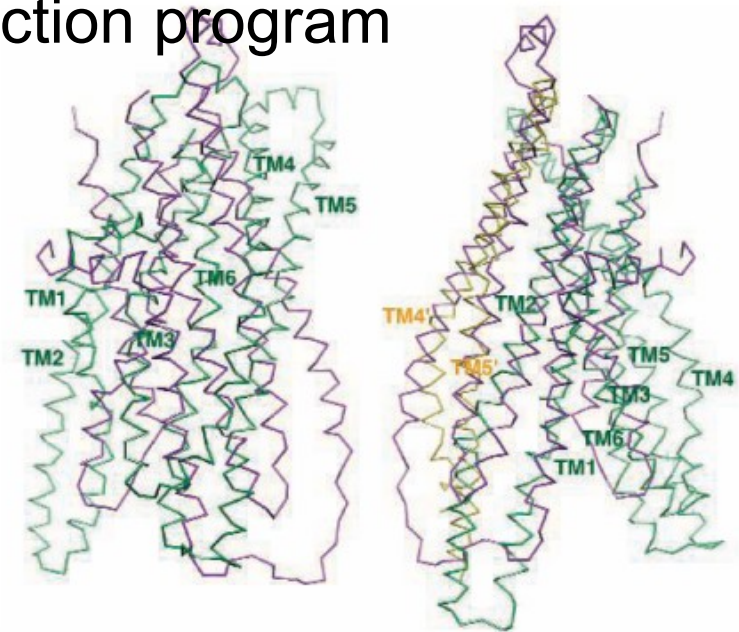
"Your recent Amazon purchases, Tweet score and location history makes you 23.5% welcome here."

[Thierry Gregorius](#)

Accurate to within \pm One (Hundred Percent)

- 2001–2005: Geoffrey Chang and colleagues published a number of high profile protein structures
 - 2001 paper on MsbA cited 360+ times by 2006
- September 2006: A dramatically different structure for a related protein is published
- December 2006: Chang et al retract five papers because “An in-house data reduction program introduced a change in sign...”

Image from:
Miller, "[A Scientist's Nightmare: Software Problem leads to Five Retractions](#)" in *Science* 314(5807): 1856-1857 (22 December 2006)



Flipping fiasco. The structures of MsbA (purple) and Sav1866 (green) overlap little (left) until MsbA is inverted (right).

A Simple Labelling Mistake?

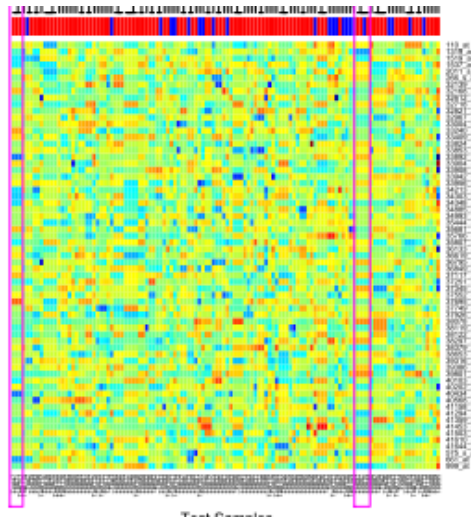
- 2006: Anil Potti and colleagues announce method for predicting patient response to chemotherapy drugs based on gene microarray data
 - 200+ citations by 2009
- 2007: Clinical trials begin
- 2007–2009: Baggerly, Coombes and colleagues try to reproduce results, but find frequent inconsistencies
- 2010–2011: Trials stopped, Potti resigns, 7+ retractions

Images from:
Baggerly & Coombes, "[Deriving chemosensitivity from cell lines: Forensic bioinformatics and reproducible research in high-throughput biology](#)" in *Annals of Applied Statistics* 3(4): 1309-1334 (2009)

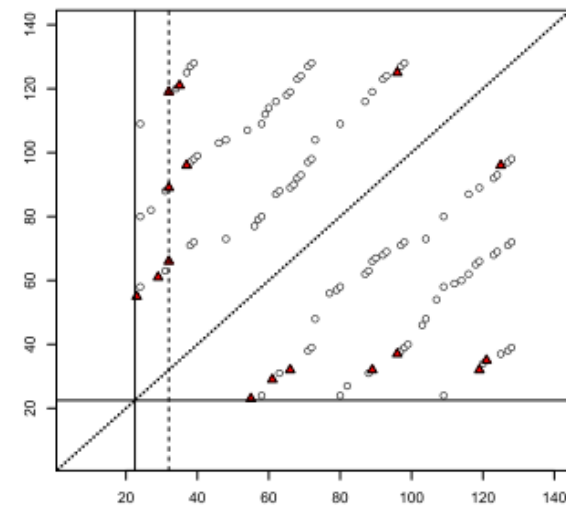
June 2014

Ian Mit

Response Labelling +
Gene Expression Heatmap



Repeated Columns
(Δ : inconsistent labels)



It's Only the Global Economy

- 2010: Reinhart & Rogoff:
 - “...whereas the link between growth and debt seems relatively weak at ‘normal’ debt levels, median growth rates for countries with public debt over roughly 90% of GDP are about one percent lower than otherwise; average (mean) growth rates are several percent lower.”
 - Common justification for austerity measures
- 2013: Herndon, Ash & Pollin, unable to recreate results from raw data receive original spreadsheet from RR
 - Discover several discrepancies including that the first five “advanced economies” (alphabetically) were omitted from first calculation

Images from:

Reinhart & Rogoff, “[Growth in a Time of Debt](#),” in *American Economic Review* 100:573-578 (2010) and Herndon, Ash & Pollin, “[Does High Public Debt Consistently Stifle Economic Growth? A Critique of Reinhart and Rogoff](#)” Political Economy Research Institute Working Paper (April 2013)

Ian Mitchell, University of British Columbia

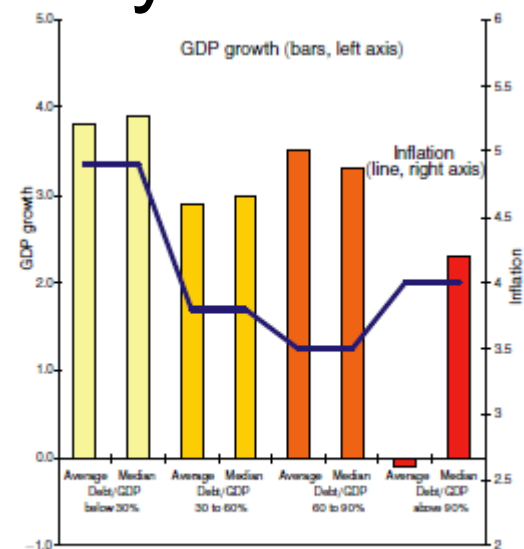
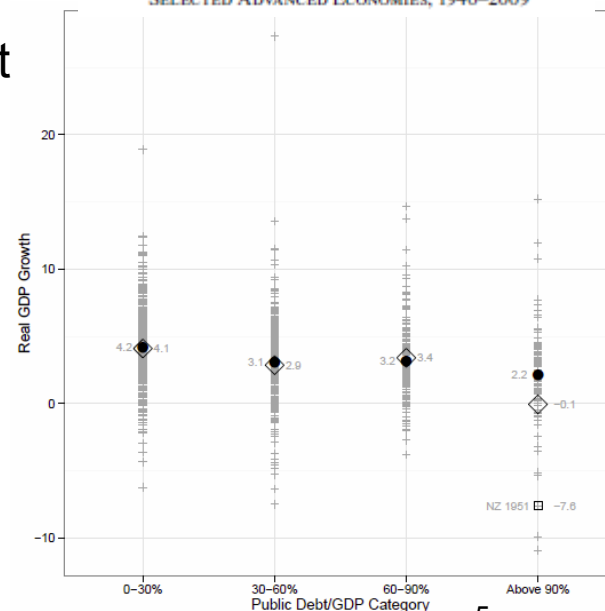


FIGURE 2. GOVERNMENT DEBT, GROWTH, AND INFLATION: SELECTED ADVANCED ECONOMIES, 1946-2009

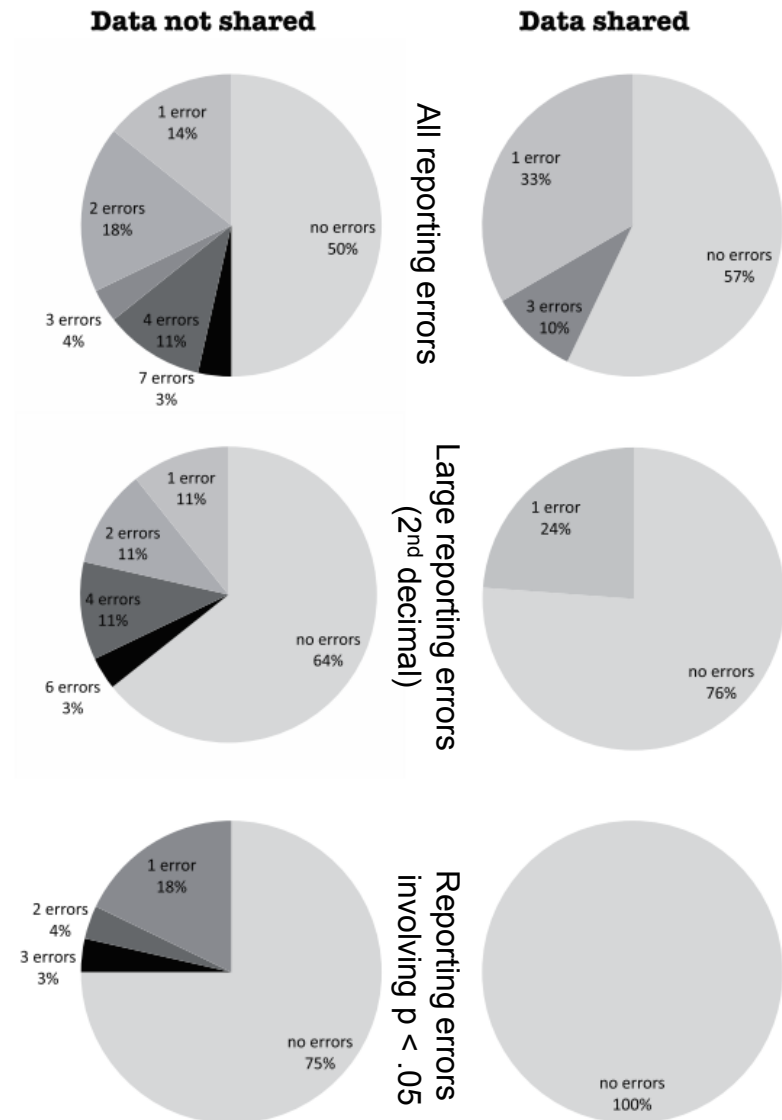


Why so Secretive?

- 2005: Wicherts and colleagues requested data from 49 papers recently published in two highly ranked American Psychological Association journals (part of a larger study)
 - Corresponding authors had signed publication form agreeing to share data
 - 21 shared some data, 3 refused (lost or inaccessible data), 12 promised to later but did not, and 13 never responded
- 2011: Wicherts and colleagues analyze internal consistency of p-values reported from null hypothesis tests
 - Willingness to share is correlated with fewer reporting errors and relatively stronger evidence against NH

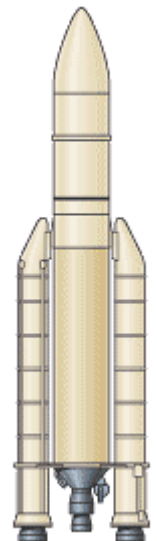
Image from:

Wicherts, Bakker & Molenaar, "[Willingness to share research data is related to the strength of the evidence and the quality of reporting of statistical results](#)" in *PLoS ONE* 6(11), Nov. 2011.



Disasters in Numerical Computing

- Feb 25, 1991: Patriot missile battery fails to track an incoming Scud missile
 - Error caused by rounding error in 24 bit timer
- August 23, 1991, Sleipner A oil platform collapses and sinks when first submerged
 - Error in finite element analysis of the strength of key concrete support structures
- June 4, 1996: maiden Ariane 5 rocket's guidance fails leading to self-destruct
 - Error caused by overflow stemming from sloppy software reuse and parameter modification

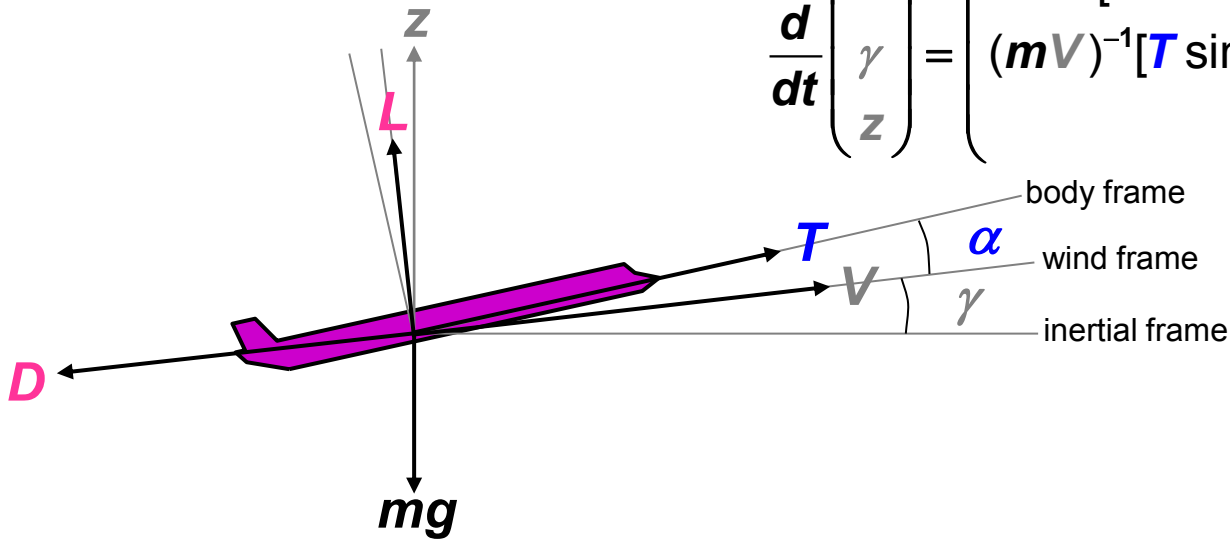


Examples from Douglas N. Arnold
<http://www.ima.umn.edu/~arnold/disasters>

A Personal Example

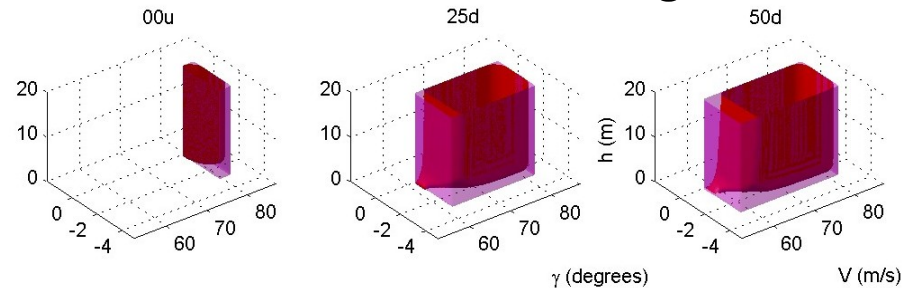
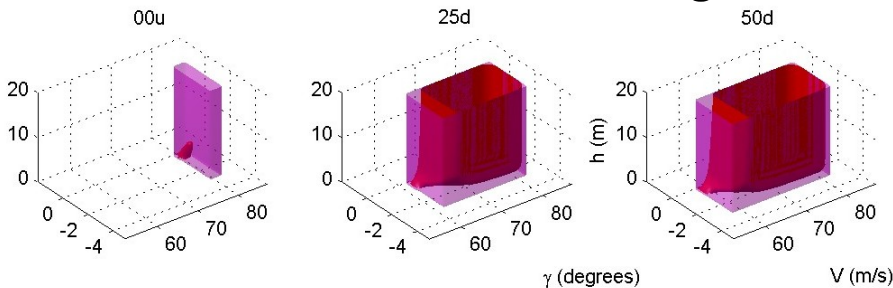
- Study of safe flap settings during aircraft final approach to runway
 - Publication: Bayen, Mitchell, Oishi & Tomlin, “Aircraft Autolander Safety Analysis Through Optimal Control-Based Reach Set Computation” in *AIAA Journal of Guidance, Control & Dynamics*, 30(1): 68–77 (2007).

$$\frac{d}{dt} \begin{pmatrix} V \\ \gamma \\ z \end{pmatrix} = \begin{pmatrix} m^{-1}[T \cos \alpha - D(\alpha, V) - mg \sin \gamma] \\ (mV)^{-1}[T \sin \alpha + L(\alpha, V) - mg \cos \gamma] \\ V \sin \gamma \end{pmatrix}$$



without mode switching

with mode switching



Could You Send Me the Code?

- Which directory was that in?

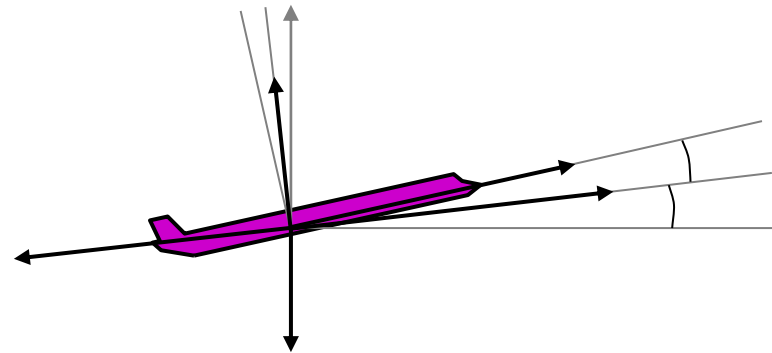
```
~/OldStanfordGagarin/Cyghome/Source/HS01/Landing/  
~/OldStanfordGagarin/Cyghome/Source/Projection/Working/  
~/OldStanfordGagarin/Cyghome/Source/JCP/  
~/OldStanfordGagarin/Cyghome/Papers/AIAA02/Source/  
~/OldStanfordGagarin/Winhome/VisualStudioProjects/LandingHighD/  
~/OldVonBraun/CygHome/Papers/AIAA03/Landing/Source  
~/OldVonBraun/CygHome/Papers/AIAA03/Landing/Shriram
```

- Which parameters did I use?

```
// 70% of 160e3 is 112e3  
// assumes fixed thrust at minT (see Flow::hamiltonian() function)  
//const GradValue ModeMinT = 0e3;  
//const GradValue ModeMaxT = 160e3;  
//const GradValue ModeMinT = 32e3;  
//const GradValue ModeMaxT = 32e3;
```

Outline

- Motivation: Some lessons in Irreproducible Research
- **Computation and Data Science**
- Reproducible Research: Changing the Culture
- Code != Data



June 2014



Ian Mitchell, University of British Columbia

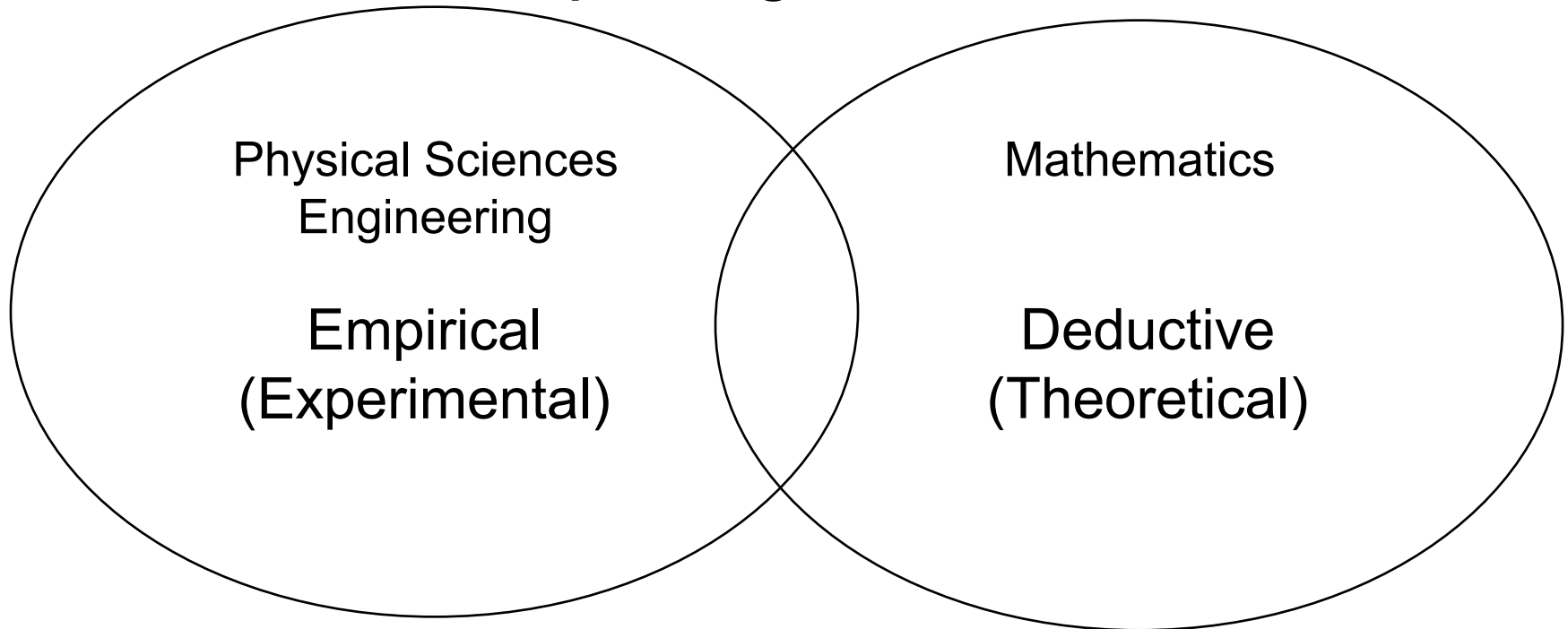


"Your recent Amazon purchases, Tweet score and location history makes you 23.5% welcome here."

Thierry Gregorius

10

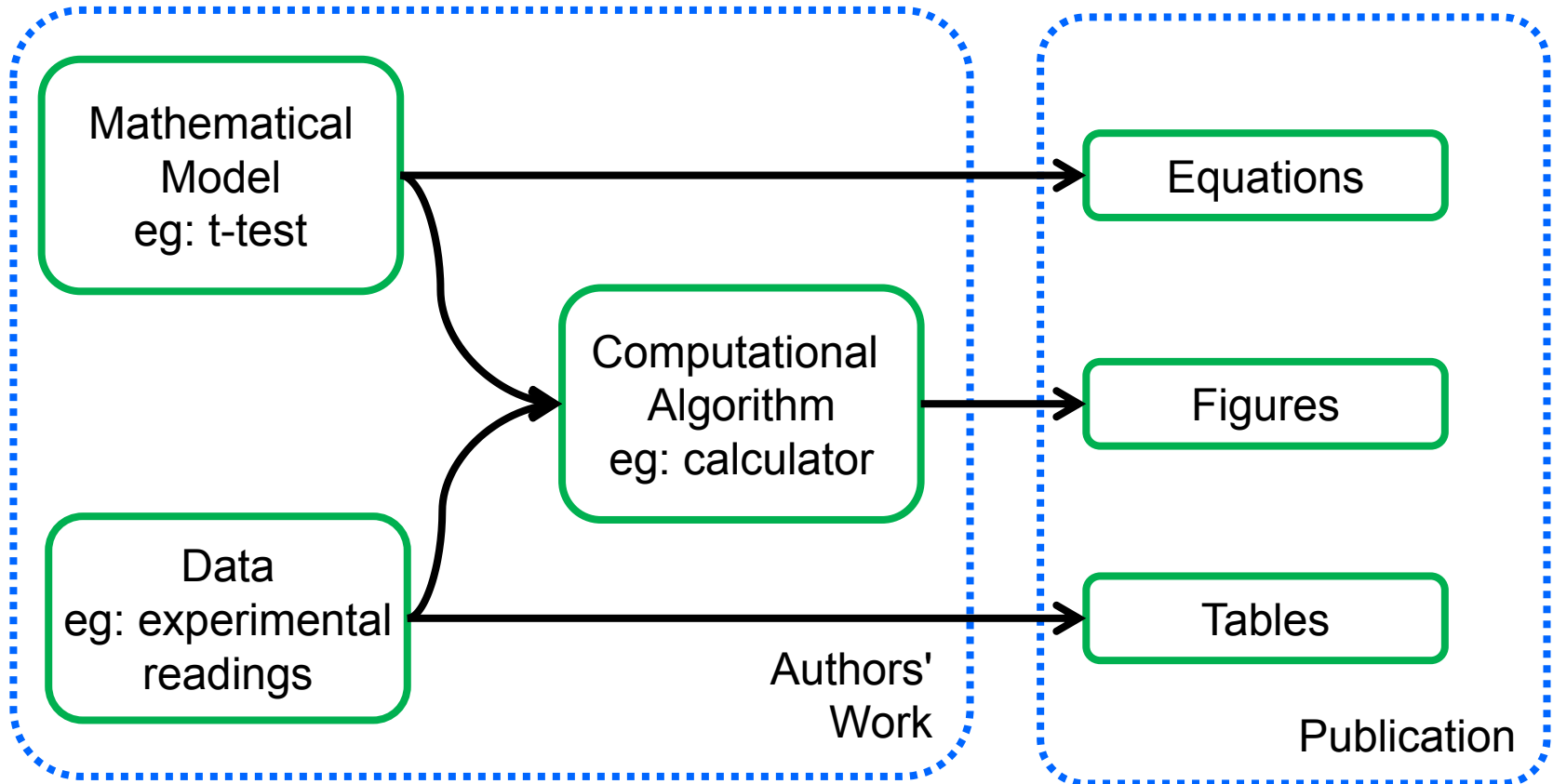
Exploring the World



- Traditionally, scientists used two approaches to build knowledge about the world
 - Data was gathered and processed by hand through simple procedures (eg: statistical summaries)

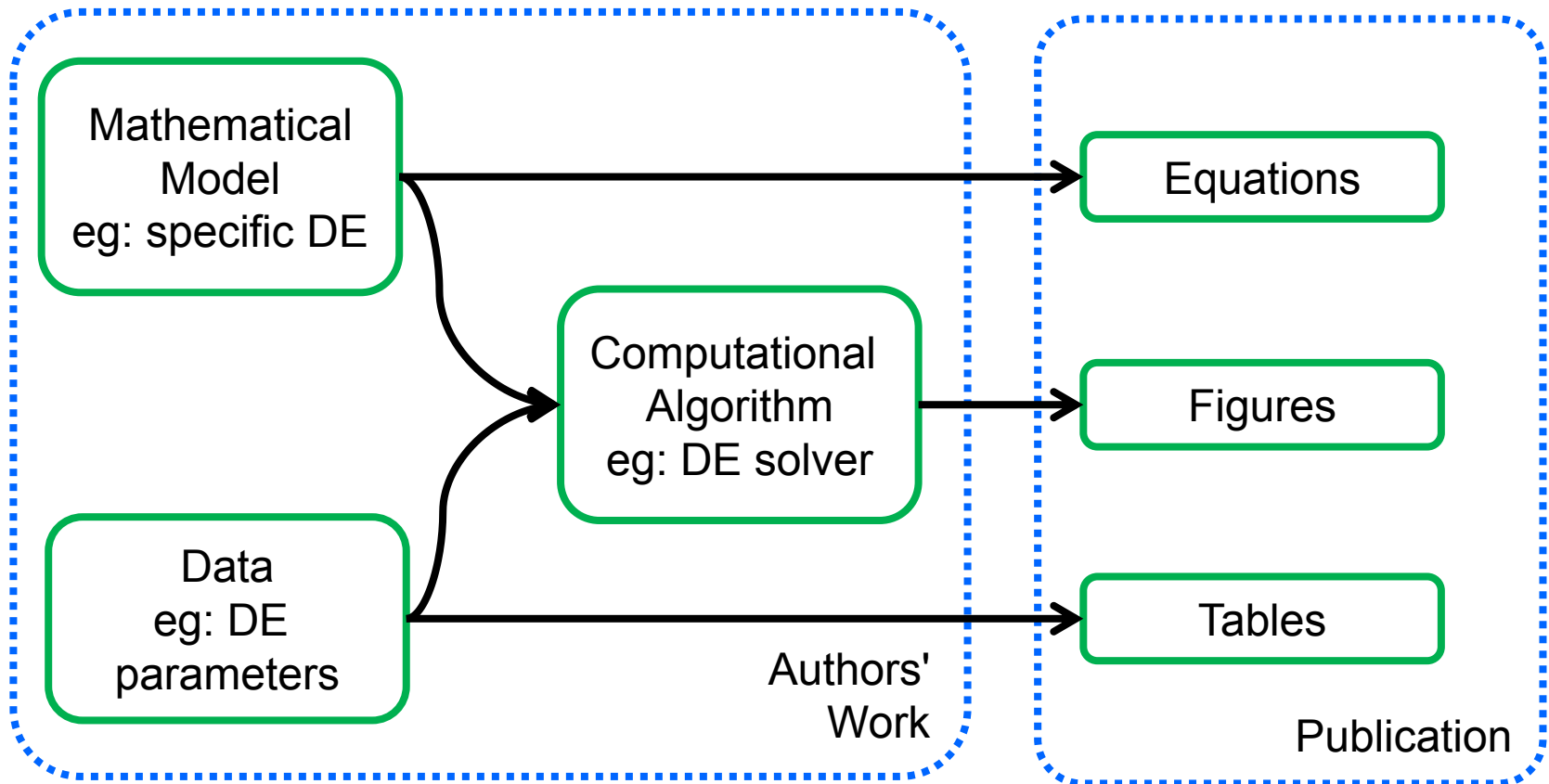
What Came Before

- Computational support for experimental analysis
 - Example: Is my hypothesis valid?



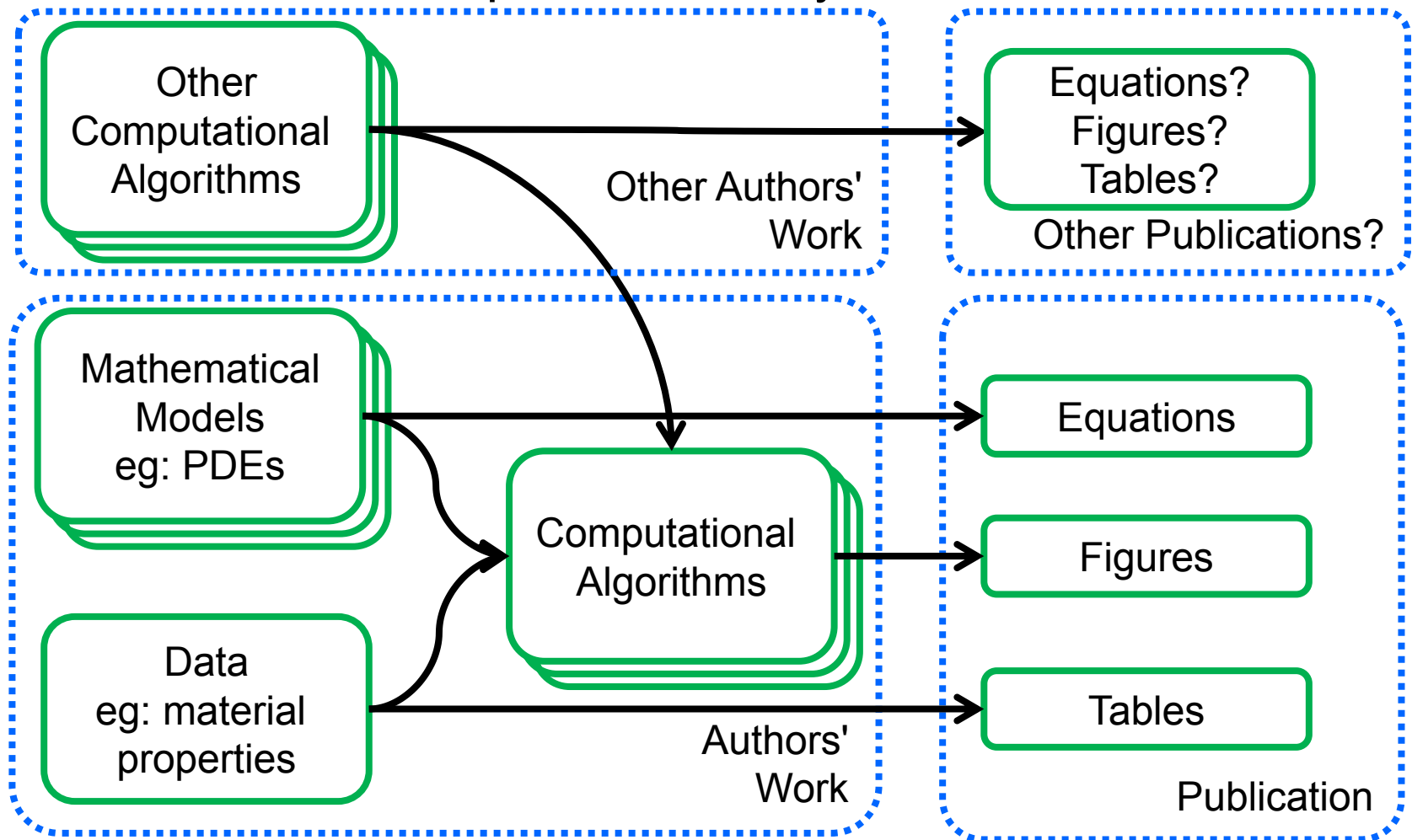
What Came Before

- Computational support for theoretical analysis
 - Example: Is my differential equation (DE) solver stable?



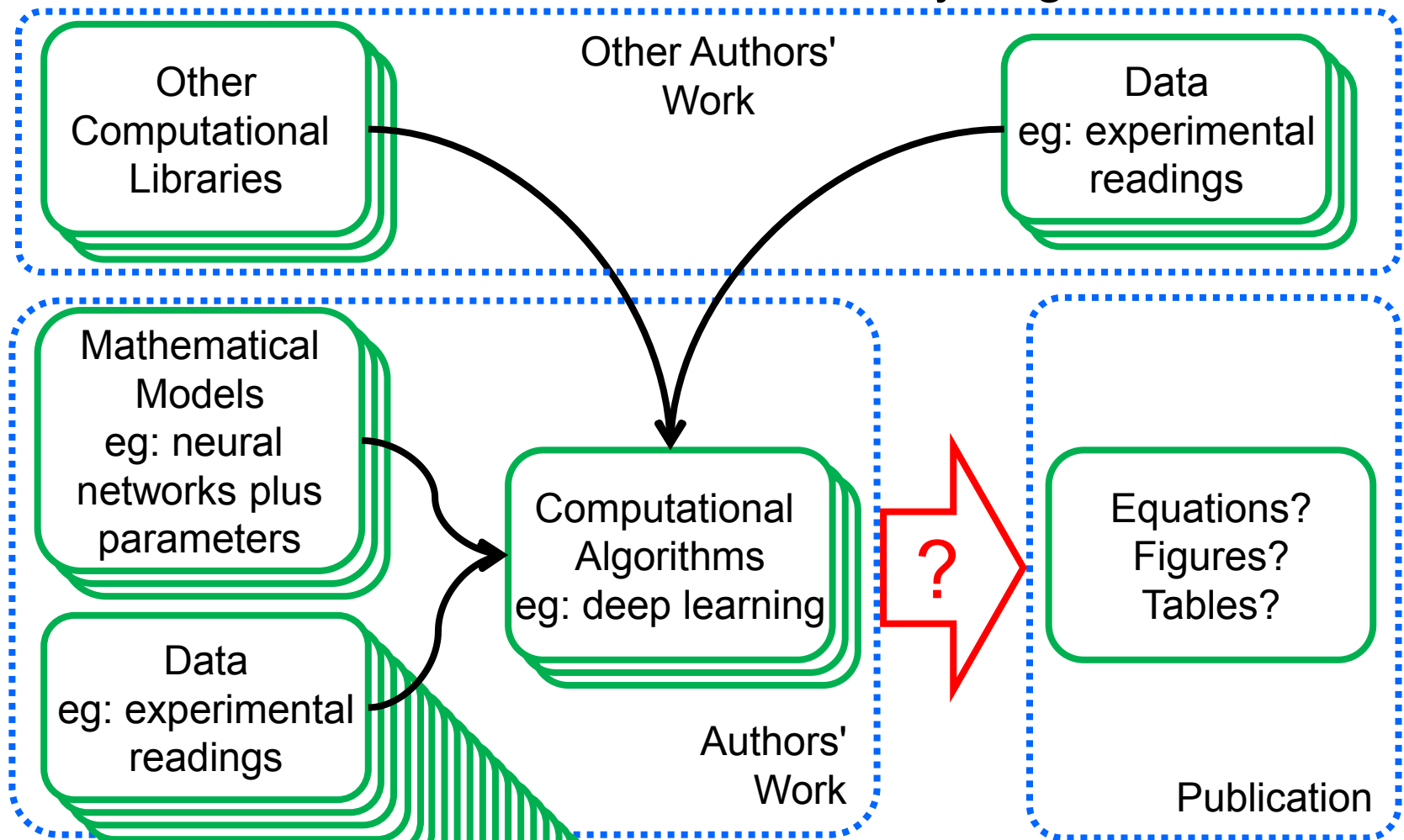
Computational Science & Engineering

- Simulation beyond the bounds of traditional theoretical or experiment analysis

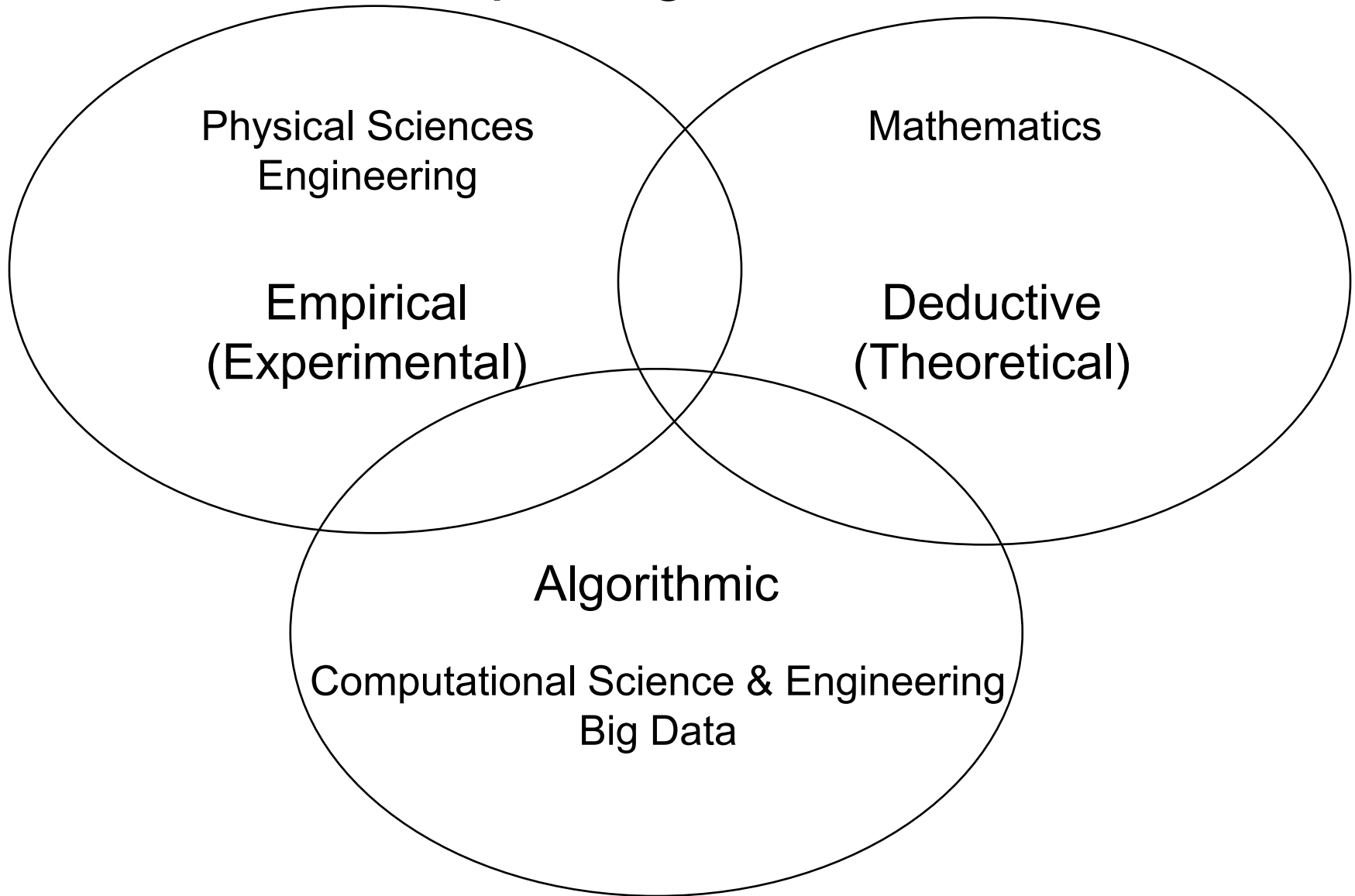


Big Data

- Algorithmically identifying and characterizing features, correlations, etc. from very large data sets



Exploring the World



What's the Big Deal?

- We must face the ubiquity of error
 - Logic (eg: in proofs)
 - Resolution (eg: accuracy, precision, sensitivity)
 - Observation (eg: calibration, misalignment, noise)
 - Transcription (eg: recording / copying the data)
 - Modeling (eg: one vs two sided t-tests)
 - Tuning (eg: choosing parameters)
 - Implementation (eg: coding the algorithm)
 - Provenance (eg: getting the right data / software)
 - Execution (eg: different hardware / software platforms)
 - Analysis (eg: drawing conclusions)
- These sources of error have always existed
- The scientific method seeks to root out such error
 - Open publication of peer reviewed manuscripts
 - Expectation of reproducibility / repeatability

The Goal: Reproducible Research

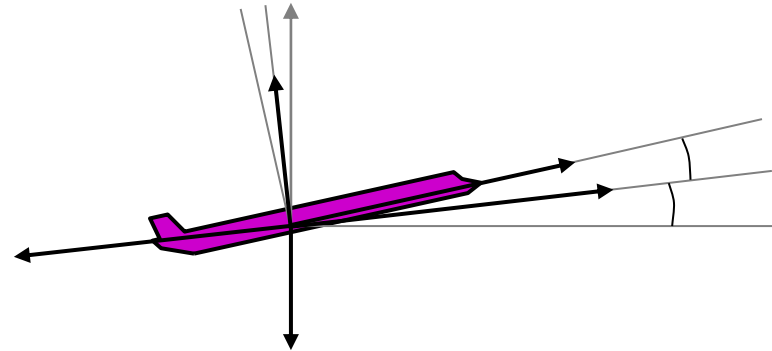
- Our current approach evolved in an age when
 - All critical details could be recorded in a manuscript
 - A single person could reasonably vet them for correctness
- As automation grows, this is no longer true
 - We can work with data at scales, speeds and efficiencies far beyond manual human oversight
 - Even the details which drive the automation (eg: code and parameters) are often more than a peer reviewer can handle
- The *reproducible research* community seeks to overcome these challenges:

“[a]n 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 software development environment and the complete set of instructions which generated the figures.”

[Jon Claerbout, as quoted by Buckheit & Donoho, 1995]

Outline

- Motivation: Some lessons in Irreproducible Research
- Computation and Data Science
- **Reproducible Research: Changing the Culture**
- Code != Data



June 2014



Ian Mitchell, University of British Columbia



"Your recent Amazon purchases, Tweet score and location history makes you 23.5% welcome here."

Thierry Gregorius

19

Changing the Culture

- Special Issue on Reproducible Research
 - Computing in Science & Engineering (July/August 2012)
 - Articles drawn from workshop and community forum held at UBC in July 2011
 - Co-organized with Victoria Stodden & Randall J. LeVeque

Cover Image of [Computing in Science & Engineering](#), 14:4



Three Themes from the Workshop

- Reproducibility of computational and data-driven science must be improved
- Challenges of encouraging reproducibility
 - How can we define, interpret, review, reduce barriers to, improve incentives for and provide examples of reproducible research?
- Development of tools & strategies to enhance and simplify reproducibility
 - Need to capture the computational environment, the provenance and the scientific narrative

Two Discussions at a Community Forum

- Journals & Publishers
 - Unclear whether computational and data science artifacts need traditional journal services (eg: managing peer review, formatting, dissemination, archiving)
 - Not clear to what extent code peer review is feasible
 - Policies can be used to encourage reproducibility, both directly (requiring code and data submission) and indirectly (eg: enforcing consistent citation)
- Funding Agencies
 - NSF data management plan requirements depend on research community
 - Short-term grant funding at odds with archival requirements
 - Include code and data sharing in CVs to provide credit
 - Computational scientists must become involved with discussions around open science

A Call to Arms

- "Next Steps" from the special issue:
 - All computational scientists should practice reproducibility, even if only privately and for the benefit of current and future research efforts
 - All interested computational scientists should tackle institutional and community challenges: train students, publish examples, request code during reviews, audit data management plans, etc.
 - All stakeholders must "consider code a vital part of the digitization of science"



World War I recruiting poster
[US Library of Congress Collection](#)

Setting the Default to Reproducible

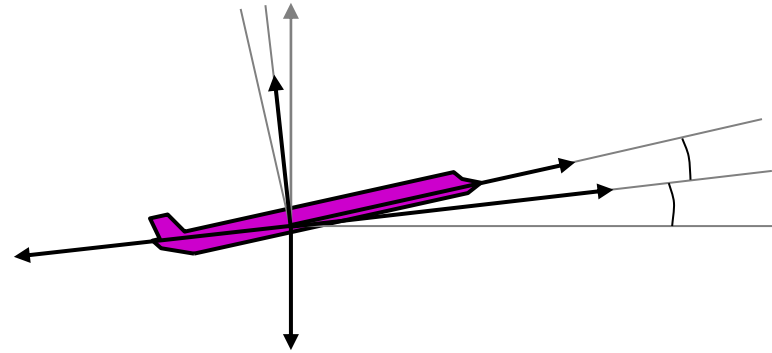
- Workshop at ICERM in December 2012 produced three recommendations:
 1. It is important to promote a culture change that will integrate computational reproducibility into the research process.
 2. Journals, funding agencies, and employers should support this culture change.
 3. Reproducible research practices and the use of appropriate tools should be taught as standard operating procedure in relation to computational aspects of research.

Making Progress

- Some top journals and conferences are allowing / encouraging / requiring elements of reproducibility
 - Nature (April 2013): Key features of data collection and statistical analysis must be specified plus data deposition mandatory for some data types, strongly recommended for many others, availability of code must be specified
 - Science (Jan 2014): Key features of data collection must be specified
 - Computer science conferences (SIGMOD, OOPSLA, ESEC/FSE, SAS, ECOOP, CAV, HSCC) have begun to optionally accept and evaluate supplemental "artifacts"
 - ACM Digital Library supports linking of both reviewed and unreviewed supplemental material to papers
 - Software Carpentry project is teaching dozens of "bootcamps" on code and data management around the world

Outline

- Motivation: Some lessons in Irreproducible Research
- Computation and Data Science
- Reproducible Research: Changing the Culture
- Code != Data



June 2014



Ian Mitchell, University of British Columbia



"Your recent Amazon purchases, Tweet score and location history makes you 23.5% welcome here."

[Thierry Gregorius](#)

26

Code vs Data in Reproducible Research

Treating code as a form of supplementary data ignores important features of code as an information storage artifact

- Relating to the practice of science
- Relating to the management of code
- Relating to the interaction of code with society

Code and the Practice of Science

- Code is a mechanism for generating data and hence a source of error
 - Digital data formats introduce no error (except when they do)
 - Errors are not smooth: The size of the mistake has little relationship to the size of the resultant error
 - Errors are not well characterized
- Scientists at all levels are not trained to manage code (and its errors)
 - At UBC: Physical science undergraduates take two courses in programming, life science undergraduates take none
 - Little incentive for giving or receiving instruction

Management of Code

- Almost always evolving
 - Bug fixes, refactoring, new features
 - Application programming interface (API) attempts to hide internal details from users
- Inverted data to metadata ratio
 - The code written to support a particular analysis may be short, but it draws upon libraries, compilers, operating systems, drivers, etc.
- Readable by both machine and people
- Many practices and tools have been developed to manage code
 - Version control systems and ecosystems (eg: github^h)
 - Virtual machines
 - Lints, automated testing, debuggers, profilers, ...
 - Extensive opportunities for training



Software Carpentry

- www.software-carpentry.org
- Dedicated to teaching basic software and data management skills to scientists
- Bootcamps: Intensive two-day, hands-on session covers:
 - Programming basics (Python or R)
 - Version control (git or subversion)
 - Unit testing
 - Using shell to automate tasks
 - Optional topics: databases & SQL, regular expressions, debugging, numerical packages, ...
- Screencasts covering many more topics are available from the website



Code and Society

- Privacy is not an issue
- Intellectual property rights are a huge issue
 - Afforded strong copyright protection, possibly also patents
 - Most companies and some universities restrict researchers' ability to release code
 - Proprietary platforms / libraries restrict ability to capture metadata and reproduce results
 - Broad legal consensus that open code should be treated differently than open data or open creative works
 - Huge open source community provides examples of and demonstrates benefits of open code, although size is critical to success

Conclusions

- Increasing dependence on poorly shared code and data threatens the credibility of research throughout the sciences
- Reproducible research is a broad and diffuse effort to counteract this threat
 - Overlaps with but is distinct from open access, open science, open source, etc.
 - Many exploratory efforts underway to change the culture
- The "big data revolution" cannot ignore the code
 - Automation is critical to managing the data glut
 - Code can and must be managed differently than other types of data

Reproducible Research Citations

- Reproducible research
 - Stodden, Leisch & Peng (eds.), *Implementing Reproducible Research*, CRC Press (2014)
 - Stodden, Borwein & Bailey, “[Setting the Default to Reproducible in Computational Science Research](#)” in *SIAM News*, June 2013
 - Leveque, “[Top Ten Reasons to Not Share Your Code \(and why you should anyway\)](#)” in *SIAM News*, April 2013
 - LeVeque, Mitchell & Stodden, “[Reproducible Research for Scientific Computing: Tools and Strategies for Changing the Culture](#)” in *Computing in Science and Engineering* 14(4): 13–17 (2012)
 - Stodden, “[Enabling Reproducible Research: Licensing for Scientific Innovation](#)” in *Int. J. Communications Law & Policy* 13 (winter 2009)

Code != Data

Why Reproducible Research Needs Both but Should Not Treat Them the Same

For more information contact

Ian M. Mitchell

Department of Computer Science
The University of British Columbia
`mitchell@cs.ubc.ca`

`http://www.cs.ubc.ca/~mitchell`

