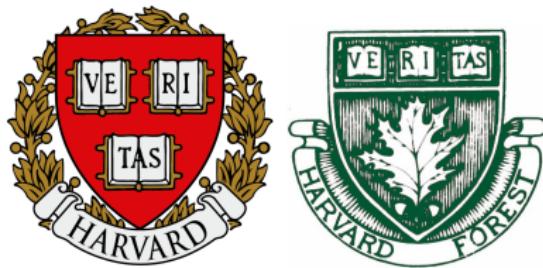


Opportunity in the Reproducibility Crisis

Computational tools to improve scientific benefaction

Matthew K. Lau, PhD



Overview:

1. Crisis (Baker, big picture, computational focus) (10 min)
 - ▶ Baker 2015
 - ▶ Sources of irreproducibility
 - ▶ Focus on computation
 - ▶ Just providing data is not enough (Stodden 2018)
 - ▶ Journal policies are shifting, motivated to make the review process easier (Data-PASS)
 - ▶ Crosas, White et al.
 - ▶ Statistical motivation: Being wrong for the right reasons $p = 0.05$
 - ▶ Crisis and opportunity (Kennedy)

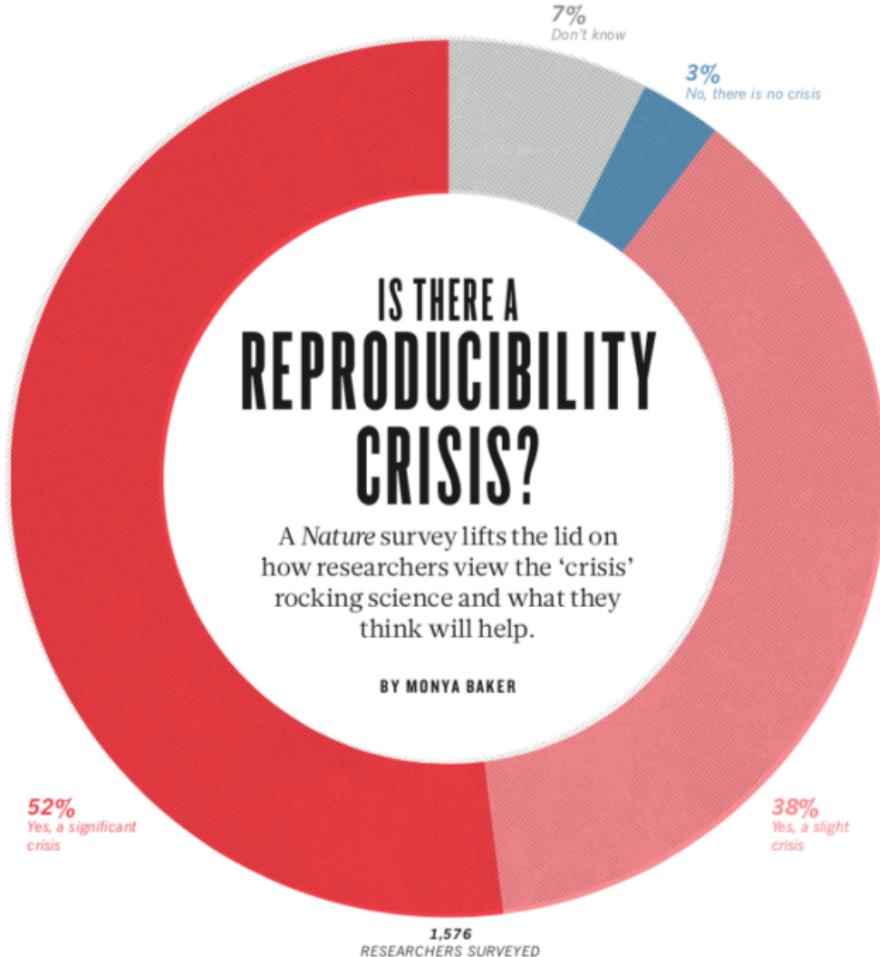
Crisis and Opportunity

The Chinese use two brush strokes to write the word 'crisis.' One brush stroke stands for danger; the other for opportunity. In a crisis, be aware of the danger—but recognize the opportunity.

— John F. Kennedy

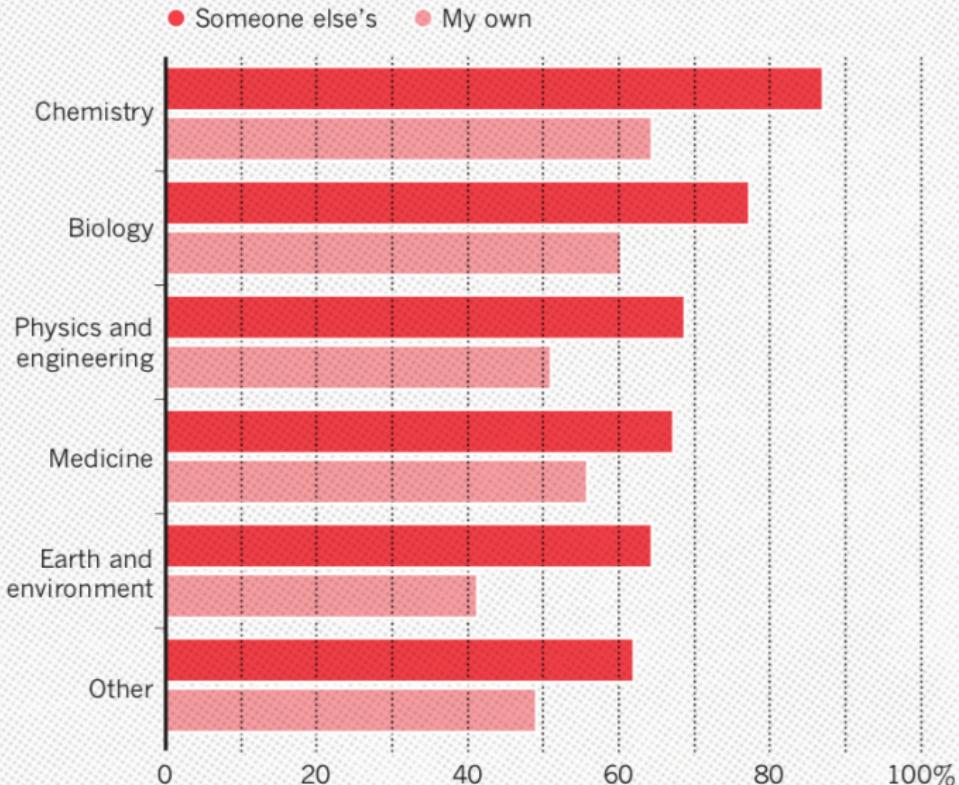
2. Opportunities (10 min)

- ▶ Benefaction, not just reproducibility
- ▶ Shifting practices will make it easier to share and access science (Benefaction and Synthesis)
- ▶ BREAK What can we do? (Ostrich)
- ▶ Computation is a significant limitation
- ▶ Sources of computational irreproducibility
- ▶ Poisot, Marwick
- ▶ *Tales of Analytical Terror* (Sharing code via email, versioning via filenames, Absolute paths)
- ▶ Pasquier, Lau, Trisovic et al. 2017
- ▶ R Rube Goldberg and The Wild Wild West of Code
- ▶ Analyses should not be a black box
- ▶ Every script is a perfectly unique snowflake that should be reusable



HAVE YOU FAILED TO REPRODUCE AN EXPERIMENT?

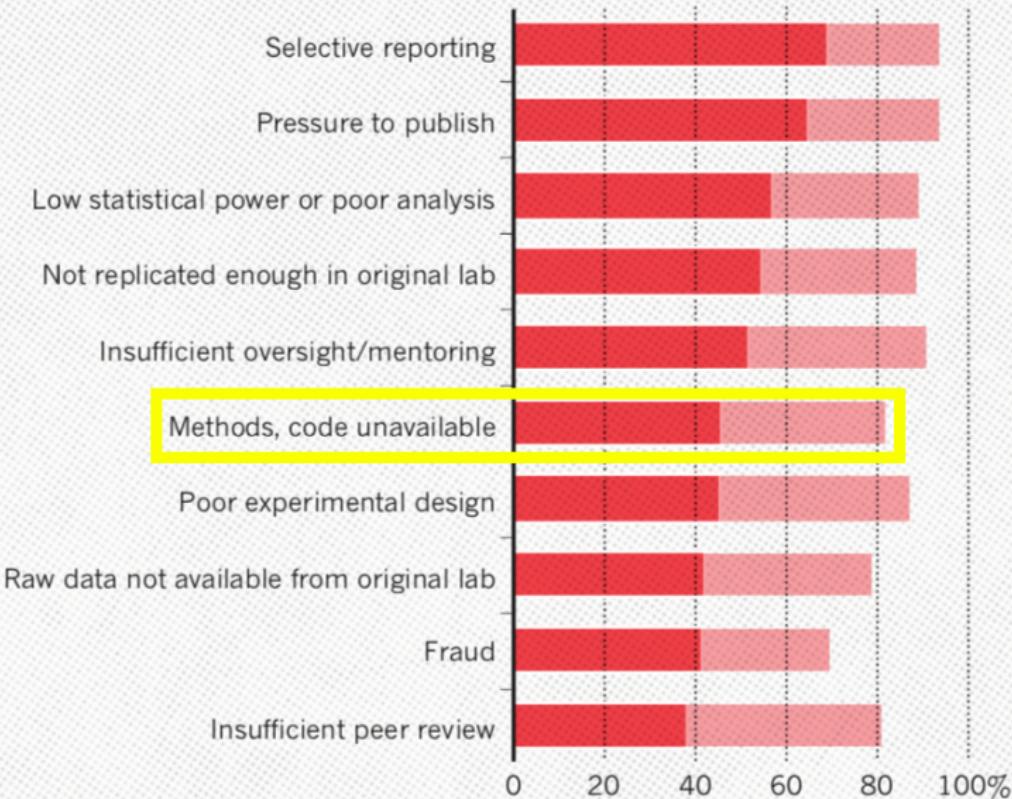
Most scientists have experienced failure to reproduce results.



WHAT FACTORS CONTRIBUTE TO IRREPRODUCIBLE RESEARCH?

Many top-rated factors relate to intense competition and time pressure.

- Always/often contribute
- Sometimes contribute



So, how bad is it, really?

We got one big thing wrong (dietetics) – In the 1950s there was an alarming increase in the rates of coronary heart disease in America. A prominent scientist argued that this was due to the increase of fats in our diet. This view won out. By the 1980s nobody would have come up with any other explanation. The US government, in their dietary guidelines painted fat as unhealthy and produced guidelines that would limit its intake. The American Heart Association agreed. Thus probably billions of research dollars were channeled into researching this. In the 1960s a number of researchers also suggested that carbohydrates/sugars were the cause, but they lost. And only in the last decade has this become an acceptable hypothesis again and it seems increasingly clear that the increase in carbohydrates in our diet (a trend present in the 1950s but actually accelerated by the avoidance of fats) is also highly problematic. At least as problematic as fat. Evidence on the relative balance is not yet clear, but it is at least credible that sugar is much worse for us than fat. The definition of an acceptable research agenda, research dollars, and public policy were all off target for half a century. Oops.

Computation is playing a significant role

Take-home

- ▶ There are major issues in reproducibility and even repeatability across science
- ▶ Analytics is playing a significant role

What can we do?

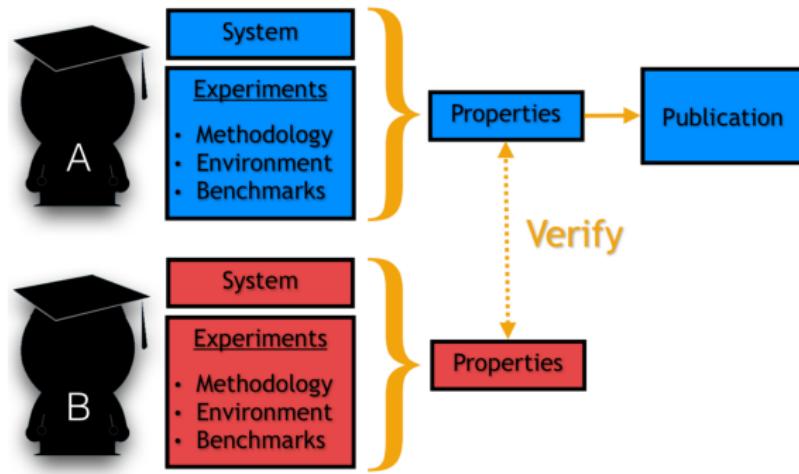


What can we do?

- ▶ Open-process
- ▶ Open-data
- ▶ Open-software

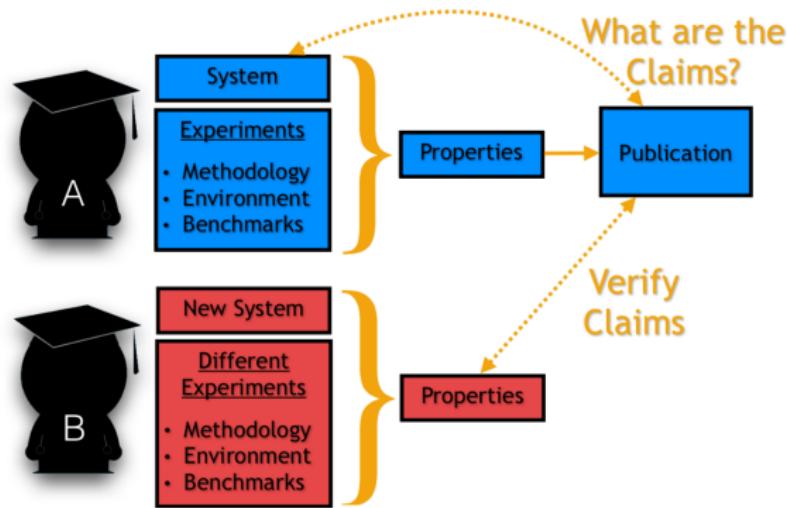
Benefaction: the real benefit

Opportunity: Benefaction not just reproducibility



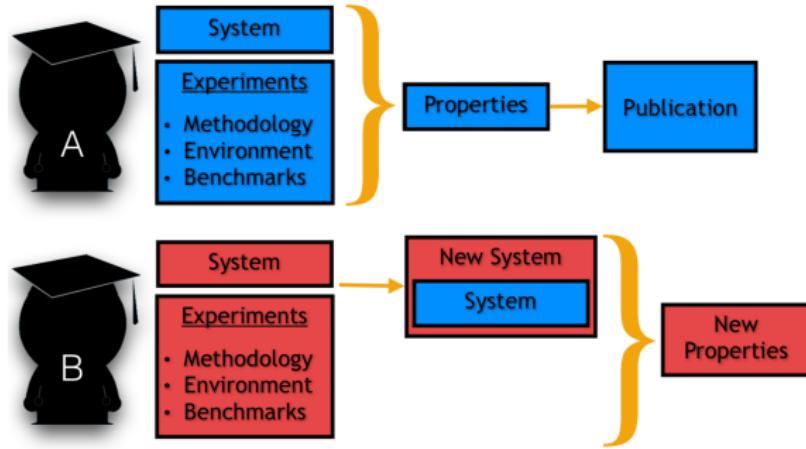
Colberg et al. 2015 Comm ACM

Opportunity: Benefaction not just reproducibility

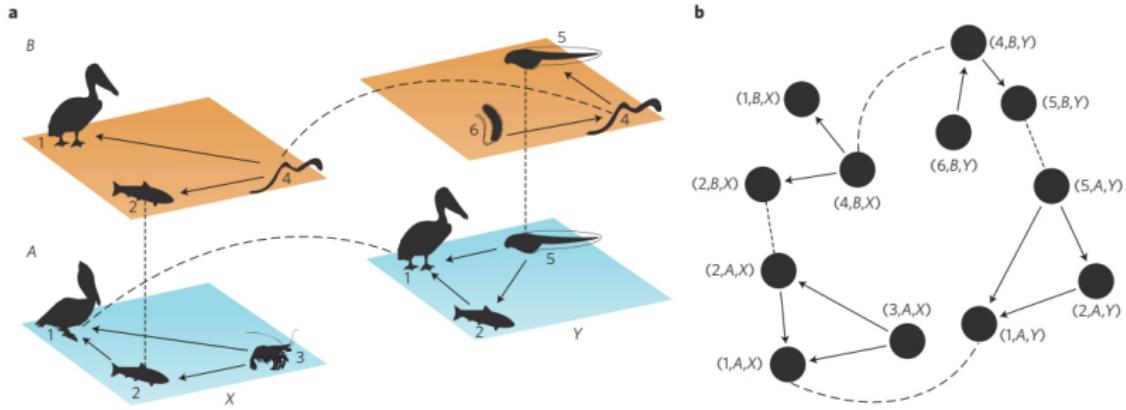


Colberg et al. 2015 Comm ACM

Opportunity: Benefaction not just reproducibility



Colberg et al. 2015 Comm ACM



Motivation: Code in Ecology



IDEAS IN ECOLOGY AND EVOLUTION 8: 55–57, 2015

doi:10.4033/iee.2015.8.9.c

© 2015 The Author. © Ideas in Ecology and Evolution 2015

Received 13 May 2015; Accepted 8 June 2015

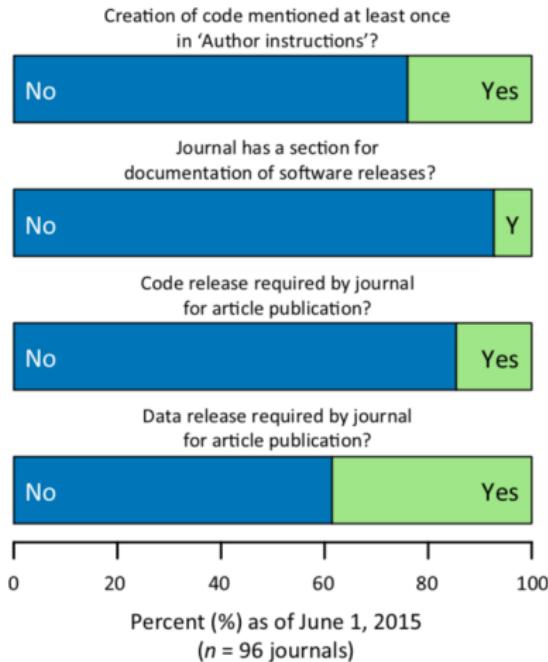
Commentary

Some thoughts on best publishing practices for scientific software

Ethan P. White

Ethan P. White (ethan@weecology.org), Department of Wildlife Ecology and Conservation, University of Florida, Gainesville, FL 32611-0430 and Department of Biology, Utah State University, Logan, UT 84322

Motivation: Ecology Journal Policies



Meeslan, Heer and White 2016 Trends in Eco Evo

Motivation: Social Science Journal Policies



Crosas et al. 2018 SocArXiv

Motivation: Journal Policy Impacts

PNAS



An empirical analysis of journal policy effectiveness for computational reproducibility

Victoria Stodden^{a,1}, Jennifer Seiler^b, and Zhaokun Ma^b

^aSchool of Information Sciences, University of Illinois at Urbana–Champaign, Champaign, IL 61820; and ^bDepartment of Statistics, Columbia University, New York, NY 10027

Edited by David B. Allison, Indiana University Bloomington, Bloomington, IN, and accepted by Editorial Board Member Susan T. Fiske January 9, 2018

(Received for review July 11, 2017)

A key component of scientific communication is sufficient information for other researchers in the field to reproduce published findings. For computational and data-enabled research, this has often been interpreted to mean making available the raw data from which results were generated, the computer code that generated the findings, and any additional information needed such as workflows and input parameters. Many journals are revising author guidelines to include data and code availability. This work evaluates the effectiveness of journal policy that requires the data and code necessary for reproducibility be made available postpublication by the authors upon request. We assess the effectiveness of such a policy by (i) requesting data and code from authors and (ii) attempting reproduction of the published findings. We chose a random sample of 204 scientific papers published in the journal *Science* after the implementation of their policy in February 2011. We found that we were able to obtain artifacts from 44% of our sample and were able to reproduce the findings for 26%. We find this policy—author remission of data and code postpublication upon request—an improvement over no policy, but currently insufficient for reproducibility.

computational reproducibility of published results. We use a survey instrument to test the availability of data and code for articles published in *Science* in 2011–2012. We then use the scientific communication standards from the 2012 Institute for Computational and Experimental Research in Mathematics (ICERM) workshop report to evaluate the reproducibility of articles for which artifacts were made available (11). We then assess the impact of the policy change directly, by examining articles published in *Science* in 2009–2010 and comparing artifact ability to our postpolicy sample from 2011–2012. Finally, we discuss possible improvements to journal policies for enabling reproducible computational research in light of our results.

Results

We emailed corresponding authors in our sample to request the data and code associated with their articles and attempted to replicate the findings from a randomly chosen subset of the articles for which we received artifacts. We estimate the artifact recovery rate to be 44% with a 95% bootstrap confidence interval of the proportion [0.36, 0.50], and we estimate the replication rate to be 26% with a 95% bootstrap confidence interval [0.20, 0.32].

reproducible research | data access | code access | reproducibility policy | open science

Motivation: Journal Policy Impacts

Table 1. Responses to emailed requests ($n = 180$)

Type of response	Count	Percent, %
Did not share data or code:		
Contact another person	20	11
Asked for reasons	20	11
Refusal to share	12	7
Directed back to supplement	6	3
Unfulfilled promise to follow up	5	3
Impossible to share	3	2
Shared data and code	65	36
Email bounced	3	2
No response	46	26

Motivation: Journal Policy Impacts

Table 2. ICERM implementation criteria for articles deemed likely to reproduce ($n = 56$)

ICERM criteria	Percent compliant, %
A precise statement of assertions to be made in the paper.	100
Full statement (or valid summary) of experimental results.	100
Salient details of data reduction & statistical analysis methods.	91
Necessary run parameters were given.	86
A statement of the computational approach, and why it constitutes a rigorous test of the hypothesized assertions.	8
Complete statements of, or references to, every algorithm used, and salient details of auxiliary software (both research and commercial software) used in the computation.	80
Discussion of the adequacy of parameters such as precision level and grid resolution.	79
Proper citation of all code and data used, including that generated by the authors.	79
Availability of computer code, input and output data, with some reasonable level of documentation.	77
Avenues of exploration examined throughout development, including information about negative findings.	68
Instructions for repeating computational experiments described in the article.	63
Precise functions were given, with settings.	41
Salient details of the test environment, including hardware, system software, and number of	13

Goal: Repeatability/Reproducibility

metadata + data + code + results + contact

Goal: Repeatability/Reproducibility

BestPractices(metadata + data + code + results + contact)

Goal: Repeatability/Reproducibility

BestPractices(metadata * data * code * results * contact)

Opportunity: Benefaction not just reproducibility

Synthesis = f(benefaction)

Research Pipeline

Data Thought	Data Collection	Data Processing	Analysis	Reporting
<i>Meta- Data + Provenance</i>	<i>Meta-Data</i>	<i>Provenance + Versioning</i>	<i>Versioning + Provenance</i>	<i>Lit Prog + Versioning</i>

Reproducibility Top Three List

1. Don't process your data manually
2. Make your pipelines transparent (data, code and notes)
3. Take credit!

Analytical Code Top Three List

1. Follow software best practices
2. Follow a consistent project architecture
3. Use version control for code and data

Resources:

- ▶ Databases List
- ▶ Tools List
- ▶ Reproducibility blogs
- ▶ RStudio
- ▶ ROpenSci
- ▶ Open Science Foundation
- ▶ Transparency in Ecology and Evolution Website and Blog
- ▶ BioRxiv

Reality: Common Ground

www.nature.com/scientificdata

SCIENTIFIC DATA

110110
0111101
11011110
011101101

OPEN

Comment: If these data could talk

Thomas Pasquier¹, Matthew K. Lau², Ana Trisovic^{3,4}, Emery Boose², Ben Couturier², Mercè Crosas⁵, Aaron M. Ellison², Valerie Gibson⁴, Chris Jones⁴ & Margo Seltzer¹

In the last few decades, data-driven methods have come to dominate many fields of scientific inquiry. Open data and open-source software have enabled the rapid implementation of novel methods to manage and analyze the growing flood of data. However, it has become apparent that many scientific fields exhibit distressingly low rates of repeatability and reproducibility. Although there are many dimensions to this issue, we believe that there is a lack of formalism used when describing end-to-end published results, from the data source to the analysis to the final published results. Even when authors do their best to make their research and data accessible, this lack of formalism reduces the clarity and efficiency of reporting, which contributes to issues of reproducibility. Data provenance aids both repeatability and reproducibility through systematic and formal records of the relationships among data sources, processes, datasets, publications and researchers.

Received: 12 April 2017

Accepted: 24 July 2017

Published: xx xxx 2017

Reality: Common Ground



Reality: Common Ground

- ▶ *Most scientists don't want to produce software, they want to do science.*

Reality: Common Ground

- ▶ *Most scientists don't want to produce software, they want to do science.*
- ▶ *Let's automate as much of the process as we can to lower activation energy, decrease error rates and increase sharing.*

Tools: Encapsulator

Sharing and Preserving Computational Analyses for Posterity with *encapsulator*

Thomas Pasquier

University of Cambridge

Matthew K. Lau and

Xueyuan Han

Harvard University

Elizabeth Fong and

Barbara S. Lerner

Mount Holyoke College

Emery R. Boose, Mercè

Crosas, Aaron M. Ellison,

and Margo Seltzer

Harvard University

Editors: Lorena A. Barba,

labarba@gwu.edu;

George K. Thiruvathukal,

gkt@cs.luc.edu

Reproducibility has become a recurring topic of discussion in many scientific disciplines.¹ Although it might be expected that some studies will be difficult to reproduce, recent conversations highlight important aspects of the scientific endeavor that could be improved to facilitate reproducibility. Open data and open source software are two important parts of a concerted effort to achieve reproducibility.² However, multiple publications point out these approaches' shortcomings,^{3,4} such as the identification of dependencies, poor documentation of the installation processes, "code rot," failure to capture dynamic inputs, and technical barriers.

In prior work,⁵ we pointed out that open data and open source software alone are insufficient to ensure reproducibility, as they do not capture information about the computational execution, that is, the "process" and context that produced the results using the data and code. In keep-

Tools: Encapsulator

Goal: Simplify computational reproducibility

1. Create a data “capsule” with code, data and environment

Tools: Encapsulator

Goal: Simplify computational reproducibility

1. Create a data “capsule” with code, data and environment
2. Increase transparency with “cleaned” code and workspace

Tools: Encapsulator

Goal: Simplify computational reproducibility

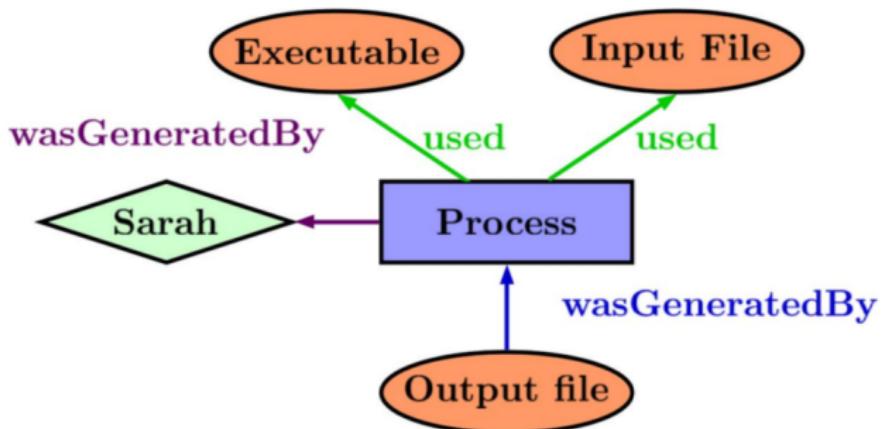
1. Capsule = all necessary software and data
2. Cleaned = organize files, remove non-essential code and re-format

Tools: Encapsulator

Basic Usage (current paradigm):

1. Code as usual in your normal environment while recording provenance
2. Run encapsulator from the console
3. List desired results
4. Product = Capsule containing essential code and data with a virtual machine

What is data provenance?



Tools: Encapsulator

Example: Messycode



Figure 4. Provenance graph corresponding to a small R script (approximately 60 lines of code).

Tools: Encapsulator

Example: Messycode

- ▶ near stream-of-consciousness coding that follows a train of thought in script development,
- ▶ output to console that is not written to disk,
- ▶ intermediate objects that are abandoned,
- ▶ library and new data calls throughout the script,
- ▶ output written to disk but not used in final documents,
- ▶ code that is not modularized,
- ▶ code that is syntactically correct but not particularly comprehensible.

Tools: Encapsulator

Example: Messycode

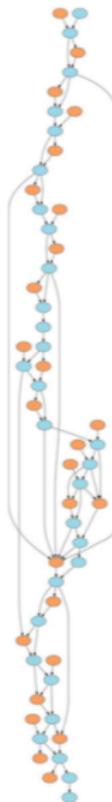


Figure 4. Provenance graph corresponding to a small R script (approximately 60 lines of code).

Tools: Encapsulator

Example: Messycode

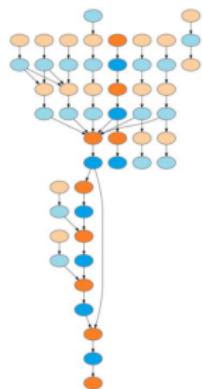


Figure 5. Data dependency transformation of the provenance graph shown in Figure 4.

Tools: Encapsulator

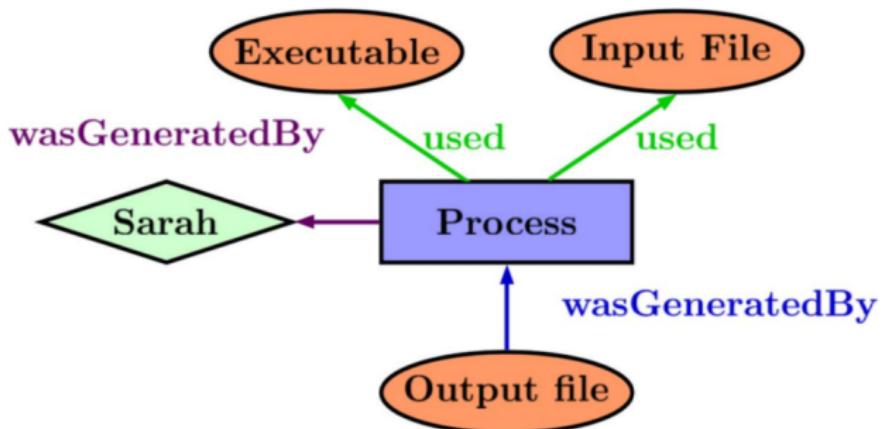
Requirements: Simplify computational reproducibility

1. The environment should present a user interface familiar to scientists.
2. Encapsulation and use (de-encapsulation) of time capsules must require minimal technical expertise.
3. The installation process itself must also require minimum intervention and technical knowledge.
4. Time capsules, their installation, and re-execution must be platform-independent.

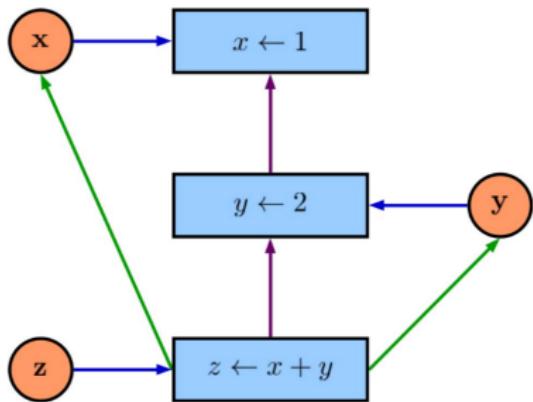
encapsulator(A Kit of Parts): Capsule creation

- ▶ Virtual Machine (encapsulator)
- ▶ Docker (containR)
- ▶ Literate computing notebook (Jupyter)
- ▶ Compressed (Reprozip)
- ▶ Capsule database (Code Ocean)

What is data provenance?



What is data provenance?



Data Provenance and R



Data Provenance and R



Data Provenance and R

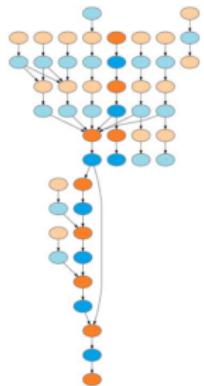


Figure 5. Data dependency transformation of the provenance graph shown in Figure 4.

encapsulator(A Kit of Parts): Provenance Details

- ▶ inputs
- ▶ outputs
- ▶ transient data objects and their values
- ▶ operations
- ▶ library dependencies

Encapsulator and benefaction

1. Eases and (potentially) improves project archiving
2. Increases clarity for re-use (others and self)

Conclusion: The next great challenge is synthesis

**** Software should not limit science ****

Conclusion: The next great challenge is synthesis



Questions and Discussion:

Possible discussion topics:



1. What checks are in place to verify and link dataverses?
2. Can provenance production become a part of the checking system?
3. What are the pros and cons of automated checking/verification and/or cleaning/encapsulation of dataverses?
4. I'm focused on R's wild-wild-west, but how does this translate to other languages?

Contact Info:

Email: *matthewklau@fas.harvard.edu*

Github: MKLau

Slack: MKLau

Much of this work was supported by NSF SSI-1450277 (End-to-End

Tools: Overview

	Code				GitHub & Gitbucket		Supplemental Material	
	Data	ocean	Zenodo	Bigshare	Dryad	PANGAEA	S3 bucket	
Meta Data	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Data Hosting	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Code Hosting	Yes	Yes	Yes	No	No	No	Yes	Yes
Versioning	No?	No?	Yes	No	No	No	Yes	No
Capsules	No	Yes	No	No	No	No	No	No
Assigns DOI	Yes	Yes	Yes	Yes	Yes	Yes	No	No
License	Flexible	Flexible	Flexible	MIT	CC0	CC-BY	Flexible	None
Cost	None	Possible	None	None	Possible	None	None	None

Adapted from Mislan, Heer & White 2016 Trends in Ecol Evol