

Transparency and Reproducibility: Case Studies, Formalisms, and Structured Guidance in Computational Social Science Applications

Victoria Stodden
University of Illinois at Urbana-Champaign

SIAM Conference on Parallel Processing for Scientific Computing (PP20)
*Session: “Transparency, Reproducibility, Sustainability, and Security: The Four Pillars of the
Next Generation Scientific Software Stack”*

Seattle, WA
February 13, 2020

Agenda

1. Setting the Stage: Research Reproducibility

- National Academies of Science, Engineering, and Medicine report

2. A Tour of Three Examples

- Container-based Reproducible Data Science with the Whole Tale project
- The “Time/Value Tradeoff” for Reproducibility: Execution in the Long Run
- Reproducibility Journal Policy: Who Re-executes the Research? Where?

3. A “Lifecycle of Data Science” Approach Includes Security

1. Setting the Stage: Research Reproducibility

Reproducibility Definitions: National Academies

In 2019 the “Reproducibility and Replication in Science” committee published consensus report (I was a committee member).

Produced key definitions and several recommendations.

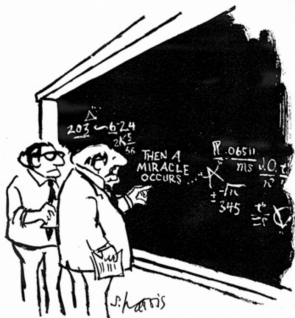
- *Reproducibility* is obtaining consistent results using the same input data, computational steps, methods, and code, and conditions of analysis. This definition is synonymous with “computational reproducibility.”
- *Replicability* is obtaining consistent results across studies aimed at answering the same scientific question, each of which has obtained its own data. Two studies may be considered to have replicated if they obtain consistent results given the level of uncertainty inherent in the system under study.

Some Reproducibility Efforts

"Setting the Default to Reproducible" in Computational Science Research

By Victoria Stodden, Jonathan M. Borwein and David H. Bailey

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.



"I THINK YOU SHOULD BE MORE EXPLICIT HERE IN STEP TWO."

Courtesy of S. Harris, ScienceCartoonsPlus.com.

recommendations emerged from the workshop discussions:

Computation is now central to the scientific enterprise, and the emergence of powerful computational hardware, combined with a vast array of computational software, presents novel opportunities for researchers. Unfortunately, the scientific culture surrounding computational work has evolved in ways that make it difficult to verify findings, efficiently build on past research, or even apply the basic tenets of the scientific method to computational procedures.

As a result, computational science is facing a credibility crisis [1,2,4,5]. The enormous scale of state-of-the-art scientific computations, using tens or hundreds of thousands of processors, presents unprecedented challenges. Numerical reproducibility is a major issue, as is hardware reliability. For some applications, even rare interactions of circuitry with stray subatomic particles matter.

In December 2012, more than 70 computational scientists and stakeholders, such as journal editors and funding agency officials, gathered at Brown University for the ICERM Workshop on Reproducibility in Computational and Experimental Mathematics. This workshop gave a broad cross section of computational scientists their first opportunity to discuss these issues and brainstorm ways to improve on current practices; the result was a series of recommendations for establishing really reproducible computational science as a standard [13]. Three main

New NISO Project: Badging Scheme for Reproducibility in the Computational and Computing Sciences

January 2019

Call for Participation

NISO voting members have approved a new project, Recommended Practice: Toward a Compatible Taxonomy, Definitions, and Recognition Badging Scheme for Reproducibility in the Computational and Computing Sciences. As publishers and researchers are placing greater emphasis on the practice of reproducibility as an essential ingredient of the scientific research process, it is critical to make compatible the taxonomies used to define the various levels of reproducibility

FROM THE EDITOR'S DESK

On Reproducible Research

By Brian Vanderborght

During the IEEE Panel of Editors meeting held this past April in Los Angeles, California, I was invited, as editor-in-chief of IEEE Robotics and Automation Magazine, to participate in a panel discussing reproducible research under the lead of

the leading role of our associate editor, Fabio Bonsignorio and colleagues, who began work on this topic ten years ago and has since organized several related workshops. Remember, we had a special issue in



Again, I offer a warm call to submit articles with reproducible research. Together with the authors, we need to improve the procedure over time, but we can only learn from experience.

With full awareness that it requires time and effort that may be in opposition to current publish-or-perish pressure, we need to increase efforts to reward the positive behavior of authors contributing to reproducible research. One initiative suggests highlighting paper with reproducible research content, perhaps even granting awards dedicated to such papers. Discussions are in progress with the Association for

EDITOR IRON VETTER
University of North Carolina Wilmington
vetter@uncw.edu



SPOTLIGHT ON TRANSACTIONS



The Reproducibility Initiative

Manish Parashar, Rutgers University

This installment of Computer's series highlighting the work published in IEEE Computer Society journals comes from IEEE Transactions on Parallel and Distributed Systems.

Reproducibility is a foundation of solid scientific and technical research. The ability to repeat research is key to confirming the validity of a... the EIC a review copy of the compute capsule, which is passed on to the assigned reproducibility associate editor for the article.

SIAM News 2013

- Editorial Policies and Badging
- Pilot Partnerships: Code Ocean

2. A Tour of Three Examples

1. Data Science in the Whole Tale Project

- Building an **open platform for computational reproducibility**
 - Create and publish **executable research objects** ("*Tales*")
- Simplify process of creating & verifying reproducible computational artifacts for scientific discovery

Easy-to-access
cloud-based computing
environments



Transparent access to
research **data**

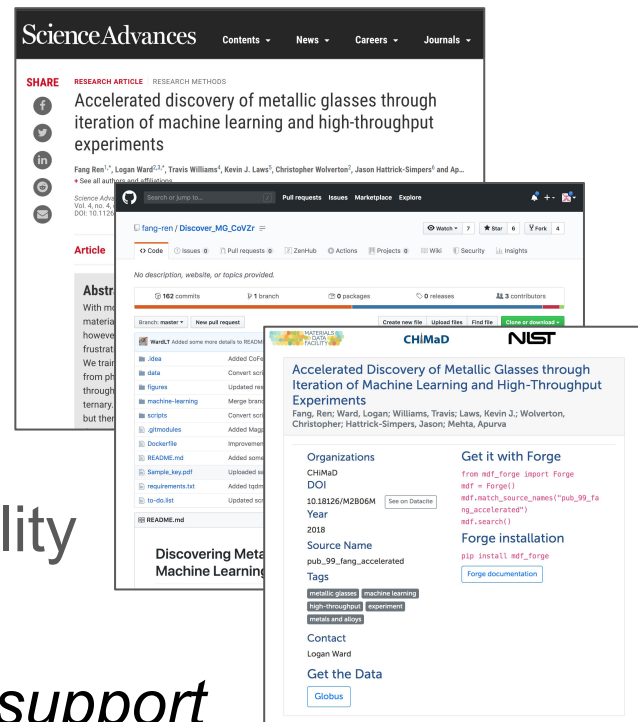


Export and **publish**
executable research
objects



Use case: Ren et al. (2018)

- ML experiments in materials science
- Published in *Science Advances*
- Code in Github
- Data published to Materials Data Facility



How can we publish the code and data to support computational reproducibility and reuse/exploration?

- Reproducibility implemented in Whole Tale

A Proposed Formalism: The “Tale”

What information do we need to reproduce and verify computational findings?

- **Manuscript**
 - source or reference
- **Documentation**
 - README, codebook, install instructions, user guide, etc.
 - License, copyright, permissions
- **Code**
 - Preprocessing, analysis, workflow
- **Data**
 - By copy, by reference, data access protocol
- **Results**
 - Output, figures, tables
- **Environment**
 - Hardware, OS, compilers, dependent software
 - Runtime, image, container
- **Provenance**
 - Computational, archival
- **Metadata**
 - Identifiers, related artifacts, Domain metadata
 - Badges
- **Version**

Tale Packaging for Sharing, Dissemination, Archiving

- Research Object
 - Beyond PDFs and datasets -- include code, workflows
 - Distributed elements
- Interoperability between systems
 - Archives/repositories
 - Active compute platforms
- BagIt serialized "Research Object" bundle
 - Zip archive + metadata + JSON-LD
 - <https://github.com/ResearchObject/bagit-ro> (=> ro-crate)



researchobject.org

2. Reproducibility Journal Policy: Who Re-executes the Research? Where?

2. Reproducibility Standards Development

Reproducibility requires community adoption and standards development.

Example: AAAS 2016 Workshop on Code and Modeling Reproducibility recommended:

- **Share** data, software, workflows, and details of the computational environment that generate published findings in open trusted repositories.
- **Persistent links** should appear in the published article and include a permanent identifier for data, code, and digital artifacts upon which the results depend.
- To enable credit for shared digital scholarly objects, **citation** should be standard practice.
- To facilitate reuse, adequately **document** digital scholarly artifacts.
- Use **Open Licensing** when publishing digital scholarly objects.
- Funding agencies should instigate new research programs and pilot studies.
- Journals should conduct a **reproducibility check** as part of the publication process.

→ NASEM 2019 “Reproducibility and Replication in Science” report recommendations.

REPRODUCIBLE RESEARCH

ADDRESSING THE NEED FOR DATA AND CODE SHARING IN COMPUTATIONAL SCIENCE

By the Yale Law School Roundtable on Data and Code Sharing

Roundtable participants identified ways of making computational research details readily available, which is a crucial step in addressing the current credibility crisis.

Progress in computing is often hampered by the inability to reproduce or verify results. Attendees at the Yale Law School RoundtableNov21/22 a set of steps that agencies, and journals improve the situation those steps here, also for best practices available options term goals for the de tools and standards.

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 Experimentation. Experimental high-performance computing questions in pure and automatic theorem proving of computational reproducibility bounds in very high performance.



Science cartoon by Ph.D. student.

INSIGHTS | POLICY FORUM

REPRODUCIBILITY

Enhancing reproducibility for computational methods

Data, code, and workflows should be available and cited

By Victoria Stodden, Marcia McNutt, David H. Bailey, Eva Deelman, Yolanda Gil, Brooke Hanson, Michael A. Heroux, John P.A. Ioannidis, Michela Taufer

Over the past two decades, computational methods have radically changed the ability of researchers from all areas of scholarship to process and analyze data and to simulate complex systems. But with these advances come challenges that are contributing to broader concerns over irreproducibility in the scholarly literature, among them the lack of transparency in disclosure of computational methods. Current reporting methods are often uneven, incomplete, and still evolving. We present a novel set of Reproducibility Enhancement Principles (REP) targeting disclosure challenges involving computation. These recommendations, which build upon more general

to understanding how computational results were derived and to reconciling any differences that might arise between independent replications (4). We thus focus on the ability to rerun the same computational steps on the same data the original authors used as a minimum dissemination standard (5, 6), which includes workflow information that explains what raw data and intermediate results are input to which computations (7). Access to the data and code that underlies discoveries can also enable downstream scientific contributions, such as meta-analyses, reuse, and other efforts that include results from multiple studies.

RECOMMENDATIONS

Share data, software, workflow, and details of the computational environment that generate published findings in open trusted repositories. The minimal components that enable

Sufficient metadata should be provided for someone in the field to use the shared digital scholarly objects without resorting to contacting the original authors (i.e., <http://bit.ly/2VqP1H1>). Software metadata should include, at a minimum, the title, authors, version, language, license, Uniform Resource Identifier (DOI), software description (including purpose, inputs, outputs, dependencies), and execution requirements.

To enable credit for shared digital scholarly

Stodden, McNutt, Bailey, Deelman, Gil, Hanson, Heroux, Ioannidis, Taufer (2016). Enhancing Reproducibility for Computational Methods. Science.

Transparency and Openness Promotion (TOP) and Open Problems

- Responsibility for verification; 3rd party re-execution of codes?
- JASA-ACS Reproducibility Editors? Cloud infrastructure (Whole Tale?)? Automation?
- Documentation and meta-data for data and code: transparency and liability

Summary of the eight standards and three levels of the TOP guidelines

Levels 1 to 3 are increasingly stringent for each standard. Level 0 offers a comparison that does not meet the standard.

	LEVEL 0	LEVEL 1	LEVEL 2	LEVEL 3
Citation standards	Journal encourages citation of data, code, and materials—or says nothing.	Journal describes citation of data in guidelines to authors with clear rules and examples.	Article provides appropriate citation for data and materials used, consistent with journal's author guidelines.	Article is not published until appropriate citation for data and materials is provided that follows journal's author guidelines.
Data transparency	Journal encourages data sharing—or says nothing.	Article states whether data are available and, if so, where to access them.	Data must be posted to a trusted repository. Exceptions must be identified at article submission.	Data must be posted to a trusted repository, and reported analyses will be reproduced independently before publication.
Analytic methods (code) transparency	Journal encourages code sharing—or says nothing.	Article states whether code is available and, if so, where to access them.	Code must be posted to a trusted repository. Exceptions must be identified at article submission.	Code must be posted to a trusted repository, and reported analyses will be reproduced independently before publication.
Research materials transparency	Journal encourages materials sharing—or says nothing.	Article states whether materials are available and, if so, where to access them.	Materials must be posted to a trusted repository. Exceptions must be identified at article submission.	Materials must be posted to a trusted repository, and reported analyses will be reproduced independently before publication.
Design and analysis transparency	Journal encourages design and analysis transparency or says nothing.	Journal articulates design transparency standards.	Journal requires adherence to design transparency standards for review and publication.	Journal requires and enforces adherence to design transparency standards for review and publication.
Preregistration of studies	Journal says nothing.	Journal encourages preregistration of studies and provides link in article to preregistration if it exists.	Journal encourages preregistration of studies and provides link in article and certification of meeting preregistration badge requirements.	Journal requires preregistration of studies and provides link and badge in article to meeting requirements.
Preregistration of analysis plans	Journal says nothing.	Journal encourages preanalysis plans and provides link in article to registered analysis plan if it exists.	Journal encourages preanalysis plans and provides link in article and certification of meeting registered analysis plan badge requirements.	Journal requires preregistration of studies with analysis plans and provides link and badge in article to meeting requirements.
Replication	Journal discourages submission of replication studies—or says nothing.	Journal encourages submission of replication studies.	Journal encourages submission of replication studies and conducts blind review of results.	Journal uses Registered Reports as a submission option for replication studies with peer review before observing the study outcomes.

3. A “Lifecycle of Data Science” Includes Security

The Lifecycle of Data Science

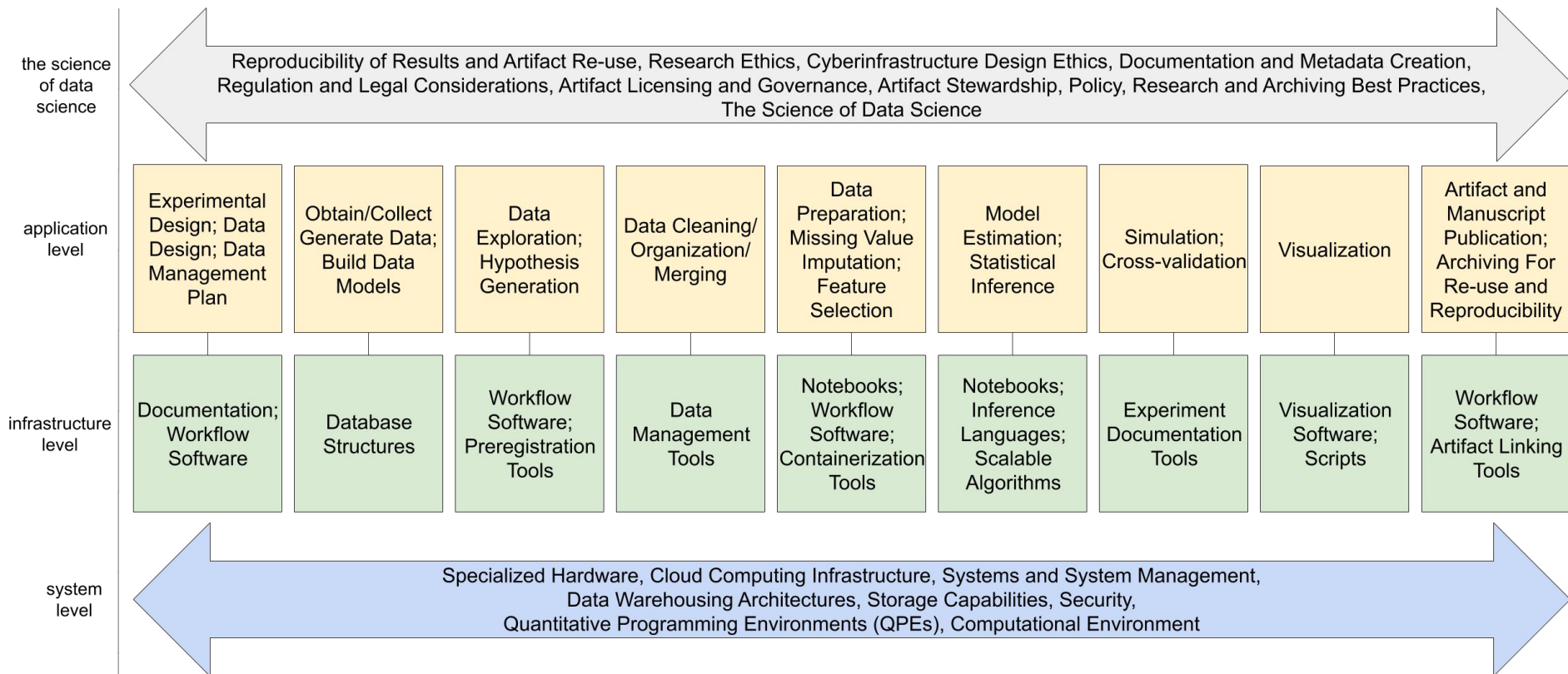
“Lifecycle of Data” is an abstraction from the Information Sciences

- Describes and relates actors in the ecosystem of data use and re-use.

What if we applied this idea to Data Science?

- **Clarify steps** in data science projects: people/skills involved, tools and infrastructure, and reproducibility through the cycle.
- **Guide implementations:** infrastructure, ethics, reproducibility and sources of uncertainty, curricula, training, and other programmatic initiatives.
- **Develop and reward contributing areas.**

A Proposal: Lifecycle of Data Science



Leveraging the Lifecycle of Data Science

An abstraction that organizes the computational pipeline.. and so recognizes different contributions including from e.g.:

- Ethicists
- Knowledge and data managers
- Compute resources and cyberinfrastructure

Goals:

- Improve understanding of Data Science advancement.
- Permit the comparison of results.
- Improve research output and social impact.

Conclusion

Two (ordinarily antagonistic) trends are converging:

Research will become **massively more compute and data intensive**,
and

Research computing will become **dramatically more transparent**.

These are reinforcing trends, which can admit exciting new opportunities:

- greater understanding of norms and social structures for discovery,
- enabling **efficiency**, **productivity**, and **discovery**.

Security issues pervasive and of ongoing importance with cyberinfrastructure development.