## Intellectual Contributions to Digitized Science: Implementing the Scientific Method

Victoria Stodden Department of Statistics Columbia University

SIAM Conference on Computational Science and Engineering "Verifiable, Reproducible Research and Computational Science" Mini Symposium Mar 4, 2011



1. Centrality of reproducibility to the scientific method, 2. The practice and tools of scientific investigation are changing, 3. Incentives: how to facilitate code and data sharing? 4. Legal barriers to reproducibility 5. Emergent incentive problems to be addressed: citation, peer review, ex-ivory tower scientific contributions. 6. Community response

### Agenda

## Reproducibility is Central to the Scientific Method

- - concept of the proof,
- ➡ Data and Code Sharing, with publication.

• Other branches of science incorporate reproducibility of results: - deductive branch (mathematics, formal logic): the well-defined

- inductive branch (experimental sciences): machinery of hypothesis testing, structured communication of methods and protocols.

• Computational Science must develop standards for reproducibility before it can be considered a third branch of the scientific method,





















Oxford Journals > Life Sciences > Nucleic Acids Research > Volume 33, Issue suppl 1 > Pp. D338-D343.

### Pseudomonas aeruginosa Genome Database and PseudoCAP: facilitating community-based, continually updated, genome annotation

Geoffrey L. Winsor, Raymond Lo, Shannan J. Ho Sui, Korine S.E. Ung, Shaoshan Huang, Dean Cheng, Wai-Kay Ho Ching, Robert E. W. Hancock<sup>1</sup> and Fiona S. L. Brinkman<sup>\*</sup> + Author Affiliations

"To whom correspondence should be addressed. Tel: +1 604 291 5646; Fax: +1 604 291 5583; Email: brinkman@sfu.ca

### Abstract

Using the Pseudomonas aeruginosa Genome Project as a test case, we have developed a database and submission system to facilitate a community-based approach to continually updated genome annotation (http://www.pseudomonas.com). Researchers submit proposed annotation updates through one of three web-based form options which are then subjected to review, and if accepted, entered into both the database and log file of updates with author acknowledgement. In addition, a coordinator continually reviews literature for suitable updates, as we have found such reviews to be the most efficient. Both the annotations database and updates-log database have Boolean search capability with the ability to sort results and download all data or search results as tab-delimited files. To complement this peer-reviewed genome annotation, we also provide a linked GBrowse view which displays alternate annotations. Additional tools and analyses are also integrated, including PseudoCyc, and knockout mutant information. We propose that this database system, with its focus on facilitating flexible queries of the data and providing access to both peer-reviewed annotations as well as alternate annotation information, may be a suitable model for other genome projects wishing to use a continually updated, community-based annotation approach. The source code is freely available under GNU General Public Licence.

Received August 12, 2004; Accepted September 28, 2004





OPEN ACCESS

```
This Article
```

```
Nucl. Acids Res. (2005) 33 (suppl

    D338-D343.

    doi: 10.1093/nar/gki047
```

```
Abstract
```

```
» Full Text (HTML)
Full Text (PDF)
```

```
    Classifications
```

```
Article
```

```
    Services
```

```
Alert me when cited
Alert me if corrected
Find similar articles
Similar articles in PubMed
Add to my archive
Download citation
Request Permissions
```

```
+ Citing Articles
```

```
+ Google Scholar
```

```
+ PubMed
```

```
+ Share
```

```
Navigate This Article
```

# Computation Central to the Scientific Endeavor

For example, in statistics,

JASA June	Computational A
1996	9 of 20
2006	33 of 35
2009	32 of 32

### Articles Code Publicly Available

0% 9% 16%

### A Crisis in Computational Science

- Computational methods becoming central to the scientific enterprise:
  - enormous, and increasing, amounts of data collection,
  - intellectual contributions now encoded in software,
  - typical scientific results rely on both data and code.
- Data and code typically not made available, rendering published results unverifiable, not reproducible.
- A Credibility Crisis

## Barriers to Data and Code Sharing in Computational Science

Survey of Machine Learning Community (Stodden, 2010):

Code 77% 52% 44% 40% 34% 30% 30% 20%

Time to docu Dealing with qu Not receiv Possibilit Legal Barrie Time to verify Potential loss of Competitors m Web/disk s

	Data
ment and clean up	54%
lestions from users	34%
ving attribution	42%
ty of patents	-
rs (ie. copyright)	41%
release with admin	38%
f future publications	35%
ay get an advantage	33%
pace limitations	29%

## Legal Barriers: Copyright

"To promote the Progress of Science and useful Arts, by securing for limited Times to Authors and Inventors the exclusive Right to their respective Writings and Discoveries." (U.S. Const. art. I, §8, cl. 8)

- (papers, code, figures, tables..)
- - reproduce the work

Exceptions and Limitations: Fair Use.

• Original expression of ideas falls under copyright by default

Copyright secures exclusive rights vested in the author to:

- prepare derivative works based upon the original - limited time: generally life of the author +70 years

## Responses Outside the Sciences I: Open Source Software

- Hundreds of open source software licenses:
  - GNU Public License (GPL) - (Modified) BSD License - MIT License - Apache 2.0 License

 Software with licenses that communicate alternative terms of use to code developers, rather than the copyright default.

- ... see <u>http://www.opensource.org/licenses/alphabetical</u>

## Responses Outside the Sciences 2: Creative Commons

- Adapts the Open Source Software approach to artistic and creative digital works
- Provides a suite of licensing options:

  - BY: if you use the work attribution must be provided, - NC: the work cannot be used for commercial purposes, - ND: no derivative works permitted, - SA: derivative works must carry the same license as the original

## Response from Within the Sciences

The Reproducible Research Standard (RRS) (Stodden, 2009)

- Release media components (text, figures) under CC BY, • Release code components under Modified BSD or similar, • Release data to public domain or attach attribution license.

Winner of the Access to Knowledge Kaltura Award 2008

- A suite of license recommendations for computational science:

- Remove copyright's barrier to reproducible research and,
- Realign the IP framework with longstanding scientific norms.



 Collaborative efforts in database building? • differential citation? (web vs article citation, microcitation) • citizen contributions? (Galaxy Zoo, Open Dinosaur Project) • Code development? review? Code maintainance for reproducibility, scientific reuse? platform building (DANSE, Wavelab, Sparselab) • open source software as a model?

## Incentives and Open Questions: Citation and Contributions

- database versioning (e.g. King and Altman 2007, Donoho and Gavish 2011)

## Challenges to Open Science

• nefarious uses? black boxes and opacity in software (why the traditional methods section is inadequate, massive codebases), Iock-in: calcification of ideas in software? Independent replication discouraged?

• "Taleb Effect" - scientific discoveries as (misused) black boxes,

- policy maker engagement: finding support for our norms,
- Commercial incentives for the scientist/university (Bayh-Dole).

## Yale Data and Code Sharing Roundtable 2009

- Nov 21, 2009:

  - Engineering, Sep/Oct 2010)
  - agencies, universities.

Roundtable on Data and Code Sharing in computational science

• gathered 30 computational scientists from a variety of fields, funding agency folks, publishers, librarians, university policy makers, lawyers...

Draft Position Statement (published in IEEE Computing in Science and

• recommendations for stakeholders: scientists, journal editors, funding

<u>http://www.stanford.edu/~vcs/Conferences/RoundtableNov212009/</u>



### References

- "Enabling Reproducible Research: Open Licensing for Scientific Innovation"
- "The Scientific Method in Practice: Reproducibility in the Computational Sciences"
- "Open Science: Policy Implications for the Evolving Phenomenon of Userled Scientific Innovation"
- <u>Reproducible Research: Tools and Strategies for Scientific Computing</u>, July 2011
- <u>Reproducible Research in Computational Science: What, Why and How,</u> Community Forum, July 2011

available at http://www.stanford.edu/~vcs

