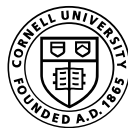Panel on Reproducibility and Transparency

# Transparency of Software and Code

Lars Vilhuber

Cornell University

# Benefits of Transparency and Reproducibility

- efficiency,
- innovation and progress,

- trust and confidence,
- and the value from the use of the data products

National Academies of Sciences, Engineering, and Medicine. 2022. Transparency in Statistical Information for the National Center for Science and Engineering Statistics and All Federal Statistical Agencies. Washington, DC: The National Academies Press. https://doi.org/10.17226/26360.

# Transparency and:

Transparency requires "**the provision of sufficiently detailed documentation**"

# Transparency and:

Transparency requires "**the provision of sufficiently detailed documentation**"

Blaise

LimeSurvey

Qualtrics

R

SAS

PostgreSQL

Python

Oracle

# Transparency and: *Proprietary Software*

Transparency requires "**the provision of sufficiently detailed documentation**"

Q: What to do when software is **proprietary**?

# Transparency and: ~~Proprietary Software~~

Transparency requires
"**the provision of sufficiently detailed documentation**"

Q: What to do when software is **proprietary**?

A: **Clearly describe use of software** *(accessibility, price, version)*

# Transparency and: *Code*

Code is
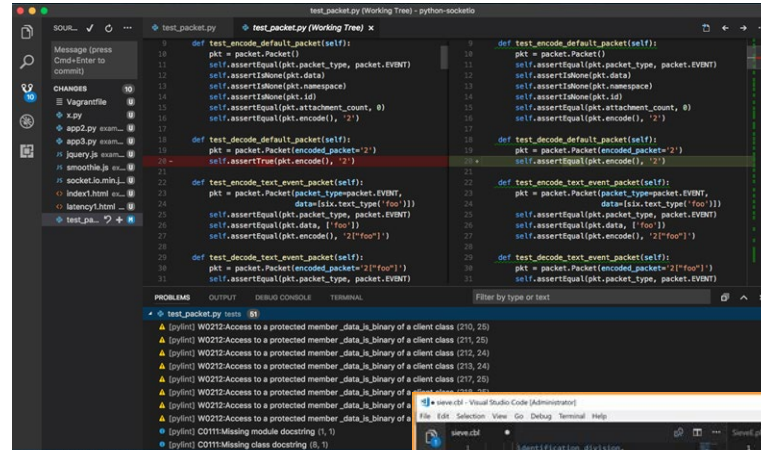**instructions to make software function**
(could also be source code for software)

- functionality of code
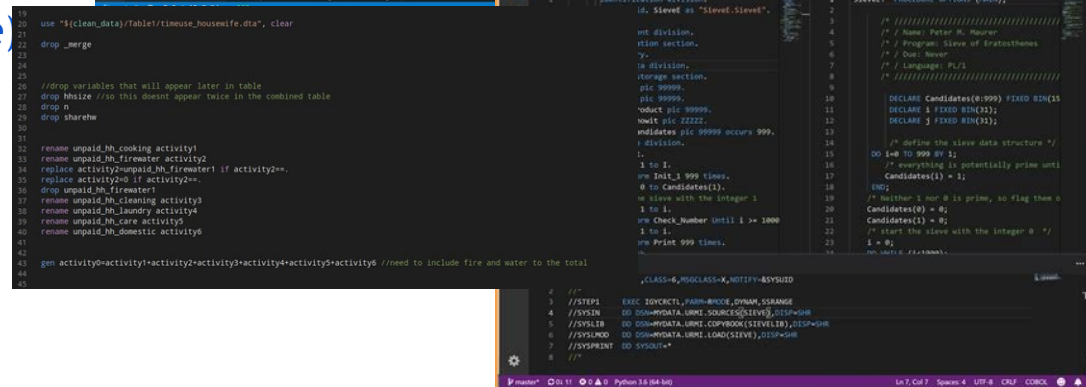- development of code (when, who)
- dependencies

# Transparency and: Code

Code is

**instructions to make software function**

(could also be source code for software)



The National Academies of SCIENCES ENGINEERING MEDICINE

# Transparency and: Code



- functionality of code

# Transparency and: *Code*

[master] 6c6faa5 My first commit - John Doe

[develop] 3e89ec8 Develop a feature - part 1 - John Doe

[develop] e188fa9 Develop a feature - part 2 - John Doe

[master] 665003d Fast bugfix - John Fixer

[myfeature] eaf618c New cool feature - John Feature

[master] 8f1e0e7 Merge branch `develop` into `master` - John Doe

[master] 6a3dacc Merge branch `myfeature` into `master` - John Doe

0.1  [master] abcdef0 Release of version 0.1 - John Releaser

- functionality of code
- development of code (when, who)
  - Includes release policy

The National Academies of | SCIENCES ENGINEERING MEDICINE

# Transparency and: *Code*



- functionality of code
- development of code (when, who)
- dependencies

# Transparency and Software

**Recommendation 4.1**

Agencies that produce federal statistics […]

- should review and make a priority of
- adopting modern information technology tools that assist in
- **collaborative software development and documentation of workflow and methodology.**

# Transparency and Software



**Recommendation 4.1**
- Use versioning systems
  - Use them broadly (not just selectively)
- Create and use style guides
  - Not just for "developers"

# Some particular notes

# Transparency and: *Surveys*

Agencies have exemplary tradition of **publishing questionnaires**

- Most often not in re-usable formats (such as DDI, Blaise)
- Transparency of "code" would allow greater reuse/consistency

The National Academies of | SCIENCES ENGINEERING MEDICINE

# Example: Survey of Earned Doctorates

*"Data collection. In 2020, for the first time, the SED data collection did not use the self-administered paper questionnaire. The SED was completed primarily by self-administered Web survey with a small number of nonrespondents contacted to complete computer-assisted telephone interviewing (CATI)."*

https://ncses.nsf.gov/pubs/nsf22300/technical-notes#survey-design



NCSES | Survey of Earned Doctorates | SED

| Data Tables | Technical Notes | Survey Description | Additional Resources | Downloads | Contact Us | How Do I? |

## Downloads

| DESCRIPTION | PDF | EXCEL | PNG | ALL |
| --- | --- | --- | --- | --- |
| Report | PDF (674 KB) | | | |
| Report (figures and tables) | PDF (.zip 3.1 MB) | XLSX (.zip 205 KB) | PNG (.zip 3.7 MB) | All (.zip 7.0 MB) |
| Data Tables and Resources | PDF (4.2 MB) | | | |
| Data Tables | PDF (.zip 5.8 MB) | XLSX (.zip 871 KB) | | All (.zip 6.7 MB) |
| Technical Notes | PDF (350 KB) | | | |
| Technical Tables | PDF (.zip 440 KB) | XLSX (.zip 58 KB) | | All (.zip 498 KB) |
| Additional Resources | PDF (70 KB) | | | |

Blaise

LimeSurvey

Qualtrics

# Transparency and: *Processing*

**All data cleaning and preparation is *(or should be)* done by code**

- Can such code be made available?
- The processed data is probably still confidential (PII)

# Analogy: Code for academic articles

**All data cleaning and preparation code *(and instructions)* MUST be provided.**

## AEA Data Editor

The AEA Data Editor's mission is to design and oversee the AEA journals' strategy for archiving and curating research data and promoting reproducible research.

Twitter

GitHub

### Some remarks on coding when data are confidential

6 minute read

📅 **Published:** April 13, 2022

Back in the fall, I made a few notes regarding how to prepare replication packages when data are confidential (here). What I did not address, and what comes up regularly, is how to **write code** when some code and/or data are confidential.

### What is confidential code, you say?

- In the United States, some variables on IRS databases are considered super-top-secret. So you can't name that-variable-that-you-filled-out-on-your-Form-1040 in your analysis code of same data. (They are often referred to in jargon as "Title 26 variables"). Not sure why that continues to be perceived as a problem, but until the law changes, that's one possible constraint.

- Your code contains the random seed you used to anonymize the sensitive identifiers. This might allow to reverse-engineer the anonymization, and is not a good idea to publish.

- You used a look-up table hard-coded in your Stata code to anonymize

The National Academies of SCIENCES ENGINEERING MEDICINE

# Do's and **don't's**

```stata
set seed 12345
use q2f q3e county using "/data/economic/cmf2012/extract.dta", clear
gen logprofit = log(q2f)
by county: collapse (count)  n=q3e (mean) logprofit
drop if n<10
graph twoway n logprofit
```

# Do's and **don't's**

```stata
set seed NNNNN
use <removed vars> county using "<removed path>", clear
gen logprofit = log(XXXX)
by county: collapse (count)  n=XXXX (mean) logprofit
drop if n<XXXX
graph twoway n logprofit
```

# Do's and don't's

Auxiliary file `include/confparms.do` (not released)

```
//=========== confidential parameters =============
global confseed    12345
global confpath    "/data/economic/cmf2012"
global confprofit  q2f
global confemploy  q3e
global confmincell 10
//=========== end confidential parameters =========
```

# Do's and don't's

Main file `main.do`:

```
//============ confidential parameters =============
capture confirm file "include/confparms.do"
if _rc == 0 {
    // file exists
    include "include/confparms.do"
} else {
    di in red "No confidential parameters found"
}
//============ end confidential parameters =========

//============ non-confidential parameters =========
global safepath "releasable"
cap mkdir "$safepath"

//============ end parameters =======================
```

# Do's and don't's



Creating reproducible packages
when data are confidential

Lars Vilhuber
2022-10-17

lars.vilhuber.com/p/fsrdc2022/

The National Academies of | SCIENCES ENGINEERING MEDICINE

# Transparency and: *Processing*

**Tracing of code execution** is hard (logging)

- **Already occurs for security purposes**
- May not need to be at the finest level for transparency
- May not need to be publicly available (but auditable)

# Transparency and: *Processing*

**Tracing of code execution** is hard (logging)

- **Already occurs for security purposes**
- May not need to be at the finest level for

be

(but



National Science Foundation
WHERE DISCOVERIES BEGIN

SEARCH

HOME   RESEARCH AREAS   FUNDING   AWARDS   DOCUMENT LIBRARY   NEWS   ABOUT NSF

**Awards**

Award Abstract # 2209629
**Collaborative Research: Elements: TRAnsparency CErtified (TRACE): Trusting Computational Research Without Repeating It**

Search Awards

NSF Org:   OAC
Office of Advanced Cyberinfrastructure (OAC)

The National
Academies of | SCIENCES
ENGINEERING
MEDICINE

# Transparency and: *Manual steps*

**Some manual steps** in processing may be **unavoidable**

Transparency implies
- that it be identified
- documented (instructions, training manuals)
- ideally publicly

# Transparency and: *Manual steps*

**Some manual steps** in processing may be **unavoidable**

Can include:

- Rules for manual edits
- Human edits

The National Academies of | SCIENCES ENGINEERING MEDICINE

# Transparency and: *Consequences*

**Transparency can be hard**

- 1,000 of people looking "over your shoulder"
- Errors **will** be found

# Transparency and: *Consequences*

**Transparency can be hard**
**... but valuable**

- Self-disciplining device
- Possible crowd-sourcing of solutions

# Transparency and: *Policies*

**Transparency needs frameworks**

- Internal policies on how to respond to (legitimate) criticism
- Support for the process

# Transparency and: *Policies*

**Transparency needs frameworks**

- Coding guides
- Continuous review process
  - For quality
  - For security
- Training!

NATIONAL ACADEMIES
Sciences
Engineering
Medicine

Consensus Study Report

**Transparency in Statistical Information for the National Center for Science and Engineering Statistics and All Federal Statistical Agencies**
(2022)

Download Free PDF    Read Free Online    Buy Paperback:$35.00

Widely available, trustworthy government statistics are essential for policy makers and program administrators at all levels of government, for private sector decision makers, for researchers, and for the media and the public. In the United States, principal statistical agencies as well as units and programs in many other agencies produce various key statistics in areas ranging from the science and engineering enterprise to education and economic welfare. Official statistics are often the result of complex data collection, processing, and estimation methods. These methods can be challenging for agencies to document and for users to understand.

[read full description]

**Contributor(s):** National Academies of Sciences, Engineering, and Medicine; Division of Behavioral and Social Sciences and Education; Committee on National Statistics; Panel on Transparency and Reproducibility of Federal Statistics for the National Center for Science and Engineering Statistics

VIEW LARGER COVER

doi.org/10.17226/26360

The National Academies of
SCIENCES
ENGINEERING
MEDICINE