# BITSS

BERKELEY INITIATIVE FOR TRANSPARENCY
IN THE SOCIAL SCIENCES

# Software and Workflow for Reproducible Research
## Discussion

Garret Christensen[1]

[1]UC Berkeley: Berkely Initiatiative for Transparency in the Social Sciences
Berkeley Institute for Data Science

APHRC, Summer 2015

- What are problems associated with reproducibility?
- What are solutions to these problems?
- What are practical tools to implement these solutions?

- Publication bias (see previous talk)
- Specification Searching (see previous talk)
- Data not available
- Code not available/unintelligible
- Code and data cannot reproduce original results

- Even with the original authors' help, you can't get the data to reproduce the published results. Or you just can't find the data to begin with.
- *Journal of Money, Credit, and Banking* Project. (Dewald et al., AER 1986)
- Martin Feldstein on Social Security and private savings, Reinhart and Rogoff on debt and GDP growth.

- Study Registry (see previous talk)
- Pre-Analysis Plan (see previous talk)
- Reproducible Workflow

- Study Registry (see previous talk)
- Pre-Analysis Plan (see previous talk)
- Reproducible Workflow

- Study Registry (see previous talk)
- Pre-Analysis Plan (see previous talk)
- Reproducible Workflow

- Literate Programing
- R Markdown and R Studio to write dynamic documents.
- Version control with Github or OSF.
- Data Sharing
    - Harvard's Dataverse

First, you should be *programming*. Working in Excel is not reproducible. Reinhart and Rogoff. If you are using SPSS, there is a 'syntax' command to record all the commands you run. Similarly in Stata, commandlog.

Better is to write scripts. R, Stata, SAS, Python, whatever. Open source has some advantages (being free, for one) but you're going to use what everyone in your field uses.

Second, literate programming. Write code to be read by a human being, with the code for the computer secondary.

Use version control.
DCVS distributed version control system
gentzkow and shapiro's RA guide

PLOS · ELSEVIER · Trello · SAGE · zotero · Twitter · MENDELEY · Impactstory · Google Scholar · GRANTS.GOV · SONA SYSTEMS · SurveyMonkey · qualtrics · amazon mechanical turk · Project Implicit · github · R · Dropbox · DRYAD · amazon web services · The Dataverse Project · figshare · DataONE · Google Docs · YouTube

OSF

GitHub and OSF Examples

A dynamic document includes your data, code, analysis, and output all in one place. Fully automated, you can guarantee no mistakes from copying and pasting.
Do this with R Markdown and R Studio or Ketchup in Stata.

R Studio Example
Stata Example

Make all your data and code publicly available.
Put it in a place where people can find it.
For APHRC that might be the APHRC repository.

# Conclusion

Simple tools exist to help you transparently and reproducibly take your research from beginning to end.

- Open Science Framework
- Trial Registries
- Version Control
- Dynamic Documents
- Trusted Public Data Archive

Read more in my *Manual of Best Practices in Transparent Social Science Research* on GitHub.