# StatTag: A Practical Approach to Reproducibility in Clinical and Translational Science

ACTStat 2017

Baltimore, MD

Leah J. Welty, PhD

# Reproducible Research

Some confusion

## PERSPECTIVE

### SCIENTIFIC INTEGRITY

# What does research reproducibility mean?

**Steven N. Goodman,\* Daniele Fanelli, John P. A. Ioannidis**

The language and conceptual framework of "research reproducibility" are nonstandard and unsettled across the sciences. In this Perspective, we review an array of explicit and implicit definitions of reproducibility and related terminology, and discuss how to avoid potential misunderstandings when these terms are used as a surrogate for "truth."

**N Northwestern Medicine®**

# What is reproducible research?

## An Evolution in Biostatistics

Requirement "that data sets and computer code be made available to others for verifying published results and conducting alternative analyses."

- Peng, 2009, *Biostatistics*

Dynamic documents that combine manuscript with code and data
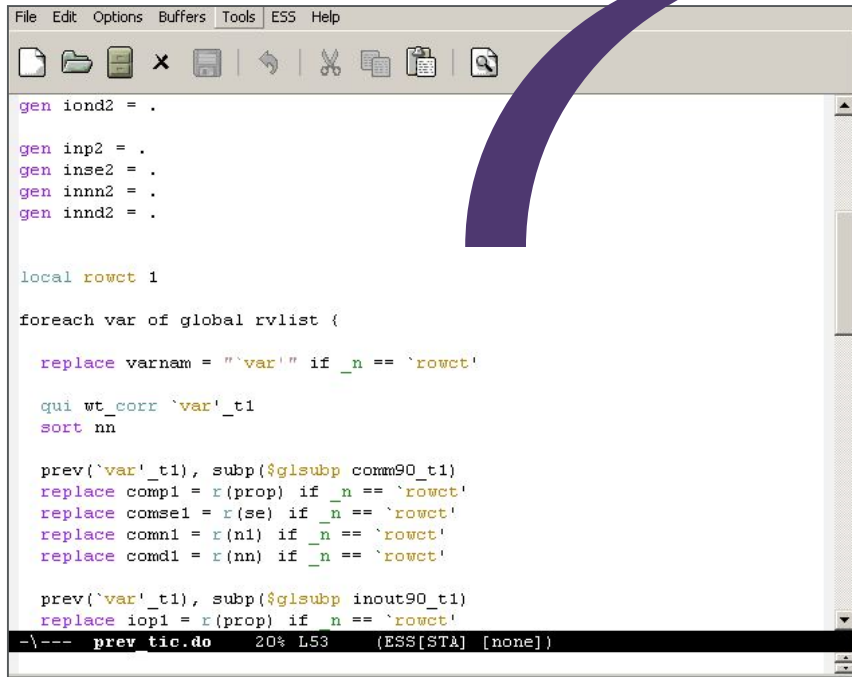e.g. Sweave (2002), R Markdown for literate programming

"We define reproducibility as the ability to re-compute data analytic results given an observed dataset and knowledge of the data analysis pipeline."

- Leek and Peng, 2015 *PNAS*

# Tools for Reproducible Research: Manuscript Prep
## Dynamic Documents



- Rather than results being hard coded in a manuscript, they can be updated automatically when data or models change.

# Tools for Reproducible Research: Manuscript Prep

## Existing tools for Dynamic Documents



Methods and results are documented and reproducible

Markdown
Ketchup
ODS
knitR
MarkDoc
Others ...

evolving rapidly ...

TEX
DOC
HTML
PDF
RTF

XLS
TXT
DAT

# Tools for Reproducible Research: Manuscript Prep

## Dynamic Documents with R Markdown

# Tools for Reproducible Research: Manuscript Prep

## The Problem with Dynamic Documents: Text Files and Collaborators

Current tools require writing within a text editor. For example, a Markdown document looks something like this:

```
3  author: "Leah Welty"
4  date: "July 27, 2006"
5  output: word_document
6  ---
7
8  ```{r setup, include=FALSE}
9  knitr::opts_chunk$set(echo = TRUE)
10 ```
11 You can use R Markdown from within RStudio.  You write in a simple text editor, using the (fairly simple) Markdown language to indicate *italics*
   output in the document.
12
13 For example, if I want to see a summary of the *cars* dataset that comes standard with R, I can insert R code that produces this:
14
15 ```{r cars}
16 summary(cars)
17 ```
18
19 I can also embed results directly in the text.  For example, the median speed is `r mean(cars$speed)`.
20
21 That's pretty nice, because if I change something about the data, then that number can be automatically updated.
22
23 Another recent thing is that I can actually call and run Stata code from this interface.  Neat, but I still have a problem ...
24
25
```

# Tools for Reproducible Research: Manuscript Prep
## The Problem with Dynamic Documents: Text Files and Collaborators

Current tools require writing within a text editor. For example, a Markdown document looks something like this:

```
3   author: "Leah Welty"
4   date: "July 27, 2006"
5   output: word_document
6   ---
7
8   ```{r setup, include=FALSE}
9   knitr::opts_chunk$set(echo = TRUE)
10  ```
11  You can use R Markdown from within RStudio.  You write in a simple text editor, using the (fairly simple) Markdown language to indicate *italics*
    output in the document.
12
13  For example, if I want to see a summary of the *cars* dataset that comes standard with R, I can insert R code that produces this:
14
15  ```{r cars}
16  summary(cars)
17  ```
18
19  I can also embed results directly in the text.  For example, the median speed is `r mean(cars$speed)`.
20
21  That's pretty nice, because if I change something about the data, then that number can be automatically updated.
22
23  Another recent thing is that I can actually call and run Stata code from this interface.  Neat, but I still have a problem ...
24
25
```

Do you have non-technical collaborators who are willing to work this way? My collaborators (primarily doctors and social scientists) are not.

I create a dynamic document, generate the Word file and send it to collaborators.

They send back:

# Tools for Reproducible Research: Manuscript Prep
## A Problem for Dynamic Documents: Track Changes

I create a dynamic document, generate the Word file and send it to collaborators.

They send back:

**Importance**: Substance abuse—among the most costly health problems in the United States—is prevalent among incarcerated juveniles. Most stays are brief; youth then become the responsibility of the community mental health system. This is the first large-scale study to examine the prevalence of substance use disorders (SUDs) in delinquent youth during adulthood and sex- and racial/ethnic differences in the types of drugs abused. However, no large-scale study has examined substance use disorders (SUDs) in delinquent youth during adulthood.

**Objective**: To examine sex and racial/ethnic differenceschanges in the prevalence of 9 SUDs (alcohol, marijuana, cocaine, hallucinogen/PCP, opiate, amphetamine, inhalant, sedative, and unspecified drug) during the 12 years after detention (up to median age 28), focusing on sex and racial/ethnic differences.

# Tools for Reproducible Research: Manuscript Prep
## A Problem for Dynamic Documents: Track Changes

I create a dynamic document, generate the Word file and send it to collaborators.

They send back:

**Importance**: Substance abuse—among the most costly health problems in the United States—is prevalent among incarcerated juveniles. Most stays are brief; youth then become the responsibility of the community mental health system. ~~This~~ is the first large-scale study to examine the prevalence of substance use disorders (SUDs) in delinquent youth during adulthood and sex- and racial/ethnic differences in the ~~types~~ of drugs abused. ~~However, no large-scale study has examined substance use disorders (SUDs) in delinquent youth during adulthood.~~

**Objective**: To examine sex and racial/ethnic differences~~changes~~ in the prevalence of 9 SUDs (alcohol, marijuana, cocaine, hallucinogen/PCP, opiate, amphetamine, inhalant, sedative, and unspecified drug) during the 12 years after detention (up to ~~median~~ age 28)~~, focusing on sex and racial/ethnic differences.~~

> I have two (bad) choices:
> 1. Continue in Word, and loose the dynamic nature of the document.
> 2. Re-enter all of their changes in my source file.

# Tools for Reproducible Research: Manuscript Prep

## Limitations of Existing Tools for Dynamic Documents



Methods and results are documented and reproducible

XLS
TXT
DAT

14
sas
R

DO
SAS
R

MarkDoc
Ketchup
ODS
knitR
Markdown
Others ...

TEX
DOC
HTML
PDF
RTF

Northwestern Medicine®

# Tools for Reproducible Research: Manuscript Prep

A Problem for Dynamic Documents: MS Word is Ubiquitous

The **NEW ENGLAND JOURNAL** *of* **MEDICINE**

*"All text...should be in one double-spaced electronic document (preferably a **Word Doc**)"*

**JAMA** The Journal of the American Medical Association

*"For submission and review, please submit the manuscript as a **Word document**. Do not submit your manuscript in PDF format."*

**Science** AAAS

*"Science **prefers to receive files in Word's .docx** format."*

# Reproducible research using Microsoft Word?

# Overview

- StatTag is a free plug-in for Microsoft Word (Windows) and app (Mac)
  - Connects Stata, SAS (Windows only) or R code and Word document
  - You and your collaborators can work from the same Word document without breaking links between the code and data
  - Can work separately on code and the Word document



- User-friendly, easy learning curve, still evolving
  - StatTag menu consistent with Word layout
  - EndNote:Citations as StatTag:Results

Northwestern Medicine®

[See stattag.org for videos]

# Software Agnostic

StatTag works with multiple code files of different types

- StatTag can connect to multiple .do, .sas and .r files.



- This is a departure from other tools that work with one specific program (e.g., R Markdown, SAS ODS)

# Document Sharing

How does StatTag work when sharing a Word document with collaborators?

| If I have… | I can… | | |
|---|---|---|---|
| | Review/edit manuscript text | View code associated with a tag | Insert or update a tag |
| Microsoft Word | ✔ | ✖ | ✖ |
| + StatTag and Stata/SAS/R code | ✔ | ✔ | ✖ |
| + Stata/SAS/R code and Data | ✔ | ✔ | ✔ |

# Windows and Mac Versions

Not just for R/SAS/Stata on Windows

**StatTag**

- **The first releases were for Windows and Stata/SAS/R.**
- A Mac version of StatTag for Stata and R is available in beta.

| | Stata | SAS | R |
|---|---|---|---|
| **Windows** | ✓ (July 2016) | ✓ (Sept 2016) | ✓ (May 2017) |
| **Mac** | ✓ (July 2017) *Beta* | X | ✓ (July 2017) *Beta* |

# Tag highlighting
## StatTag identifies tags in a document

**StatTag**

- Inserted tags are highlighted when they are clicked on.
- Future versions will include a "highlight all tags" function to quickly find any inserted tags in a document.

**CONCLUSIONS:**
Intervention X was not statistically significantly associated with a reduction in S placebo control. Longer term follow up may be needed to assess if interventio time.

Table 1. Participant Characteristics (N=120).

| Characteristic, N (%) | Control (N=56) | | Intervention (N=64) | | P-value* |
|---|---|---|---|---|---|
| Male | 29.00 | 0.24 | 31.00 | 0.26 | 0.71 |
| Female | 27.00 | 0.22 | 33.00 | 0.28 | . |
| 30-45 Years | 23.00 | 0.19 | 17.00 | 0.14 | 0.13 |
| 45-59 Years | 19.00 | 0.16 | 21.00 | 0.17 | . |
| 60+ Years | 14.00 | 0.12 | 26.00 | 0.22 | . |
| SBP Before** | 155.09 | 10.65 | 157.64 | 11.96 | 0.22 |
| SBP After** | 149.80 | 13.78 | 152.72 | 14.48 | 0.26 |
| SBP Change** | -5.29 | 15.51 | -4.92 | 17.82 | 0.91 |

* Chi-squared or t-test
** Presented as mean (sd)
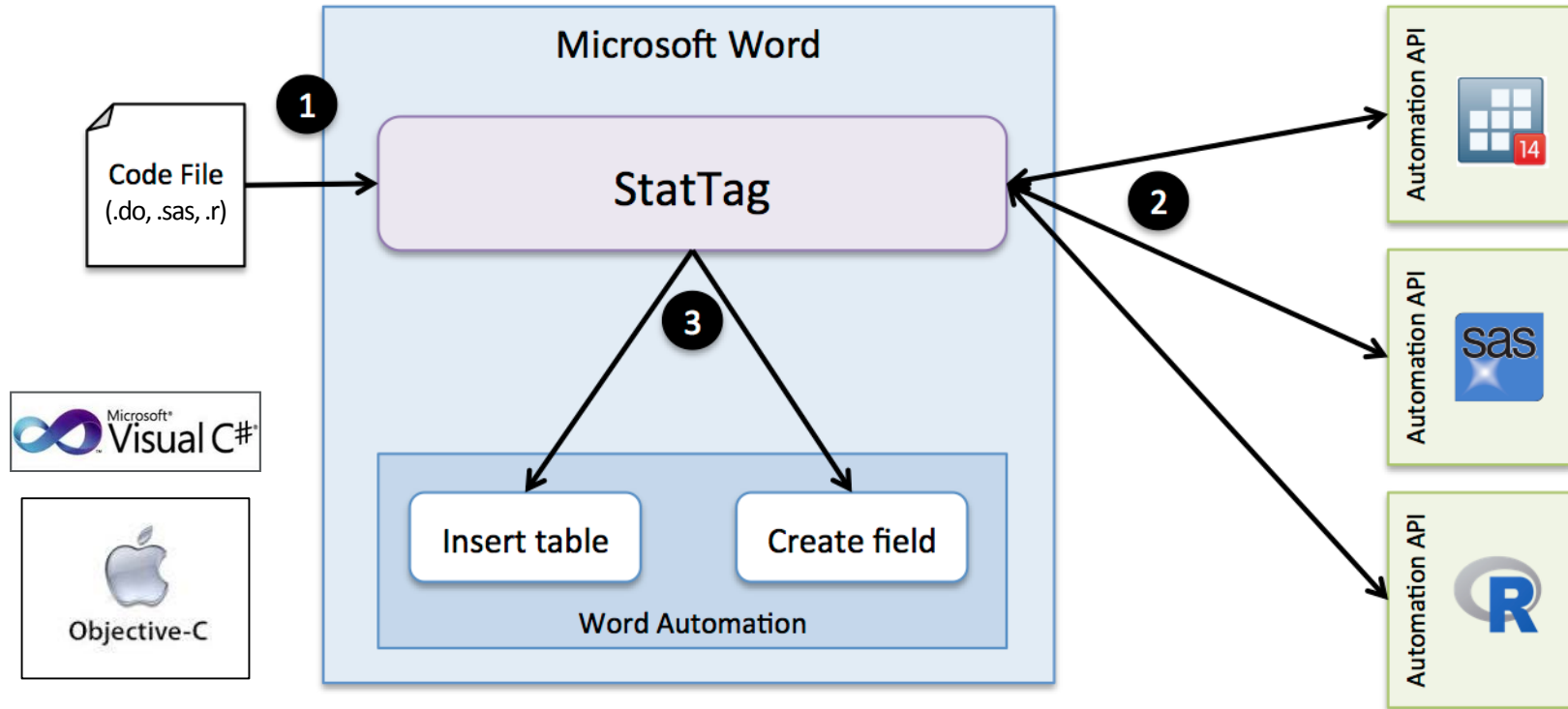
# Data Security

## Concerns over PHI, PII

- StatTag *doesn't* store a copy of the data.
- StatTag will eventually store a *read-only* copy of code files

# Architecture



1. Read the code file & parse out the tags
2. Send commands to the stat program and get individual results
3. Use Word automation to add results to the document (using native Word formatting for tables and fields).

# Getting StatTag

Freely available at stattag.org



## stattag.org



download stattag  /  user guide and tutorial  /  cite stattag  /  announcements  /  faq  /  contact

**STATTAG**

StatTag is a free software plug-in for conducting reproducible research. It facilitates the creation of dynamic documents using Microsoft Word documents and statistical software, such as Stata. Users can use StatTag to embed statistical output (estimates, tables and figures) into a Word document and then with one click individually or collectively update output with a call to the statistical program. What makes StatTag different from other tools for creating dynamic documents is that it allows for statistical code to be edited directly from Microsoft Word. Using StatTag means that modifications to a dataset or analysis no longer require transcribing or re-copying results into a manuscript or table.

Northwestern Medicine®

# (Really) Getting StatTag

GitHub repo at github.com/StatTag

StatTag

---

Features   Business   Explore   Marketplace   Pricing        This organization   Search        Sign in or Sign up

## StatTag

📍 Chicago, IL      🔗 http://stattag.org

📖 Repositories      👥 People 0

Search repositories...                              Type: All ▾    Language: All ▾

### Simple-Code-Examples
Easy code examples for StatTag

🔴 SAS    Updated on May 18

### StatTag
Windows version of StatTag

r    reproducible-research    stata    sas

#### Top languages
🟢 C#    🔴 SAS    🔴 C++

#### People                              0 ›
This organization has no public members. You must be a member to see who's a part of this organization.

**Northwestern Medicine**®

# Planned Enhancements and Future Work

- Compatibility with native code editors

- Improved table construction and formatting

- Tag highlighting and inventory

- Possible extensions to Matlab, SQL, other programs

- **Feedback – suggestions, complaints, comments – are very welcome!**
  **stattag@northwestern.edu**

# Citation and Acknowledgements

- We ask that anyone who uses StatTag to please cite:
  - Welty, L.J., Rasmussen, L.V., Baldridge, A.S, and Whitley E. (2016). *StatTag*. Chicago, Illinois, United States: Galter Health Sciences Library. doi:10.18131/G3K76

- StatTag is distributed under the MIT License

- StatTag was developed with funding through a Clinical Translational Sciences Award (CTSA) to Northwestern University. Tracking the impact of the award is a key metric in demonstrating effectiveness.

# Acknowledgements (continued)

- StatTag was inspired in part by the Stata Automation Report project:  Lo Magno, G.L. (2013). Sar: Automatic generation of statistical reports using Stata and Microsoft Word for Windows. *The Stata Journal*, 13(1); 39-64.

- StatTag makes use of the following open source projects:

  Scintilla - http://www.scintilla.org/

  ScintillaNET - https://github.com/jacobslusser/ScintillaNET

  Json.NET - http://www.newtonsoft.com/json

  *Use of these projects does not imply endorsement of StatTag by the respective project owners, or endorsement of the use of these projects by Northwestern University.*

# References

- Peng, R "Reproducible Research in Computational Science" (Science) 2 December 2011, vol. 334
- Leek and Peng "Opinion: Reproducible Research can still be wrong: Adopting a prevention approach" PNAS, February 10, 2015, vol. 112, no. 6, 1645–1646
- Matthias Schwab, Martin Karrenback, and Jon Claerbout "Making scientific computations reproducible" (2000) *Computing in Science and Engineering*, 2, pp. 61 – 67.
- Roger D. Peng. "Reproducible research and *Biostatistics."* (2009) *Biostatistics,* pp. 405-408.
- Paul Thompson and Andrew Burnett. "Reproducible Research" CORE Issues in Professional and Research Ethics, Volume 1, Paper 6, 2012. Accessed from http://nationalethicscenter.org/content/article/175
- Jonathan Buckheit and David Donoho. "WaveLab and Reproducible Research." (1995) Technical Report No. 474, Department of Statistics, Stanford Univeristy. Accessed from http://statistics.stanford.edu/~ckirby/techreports/NSF/EFS%20NSF%20474.pdf, February 2013.
- Babel: Introduction. OrgMode; Available from: http://orgmode.org/worg/org-contrib/babel/intro.html.
- Buckheit JB, Donoho DL. WaveLab and Reproducible Research. In: Antoniadis A, Oppenheim G, editors. Wavelets and Statistics. Springer New York; 1995. p. 55-81.
- Laine C, Goodman SN, Griswold ME, Sox HC. Reproducible research: moving toward research the public can really trust. Ann Intern Med. 2007;146(6):450-3.
- Thompson PA, Burnett A. Reproducible Research. CORE Issues [serial on the Internet]. 2012; Vol. 1, Paper 6: Available from: https://nationalethicscenter.org/content/article/175.
- Ware J. Reproducible Research Standards and Definitions. CTSPedia; 2010; Available from: http://www.ctspedia.org/do/view/CTSpedia/ReproducibleResearchStandards.
- Baggerly KA, Coombes KR. Deriving chemosensitivity from cell lines: Forensic bioinformatics and reproducible research in high-throughput biology. Annals of Applied Statistics. 2009;3(4):1309-34.
- Baggerly KA, Coombes KR. What information should be required to support clinical "omics" publications? Clin Chem. 2011;57(5):688-90.
- Coombes KR, Wang J, Baggerly KA. Microarrays: Retracing steps [4]. Nat Med. 2007;13(11):1276-7.
- Donoho DL. An invitation to reproducible computational research. Biostatistics. 2010;11(3):385-8.
- Gentleman R. Reproducible research: a bioinformatics case study. Stat Appl Genet Mol Biol. 2005;4:Article2.
- Announcement: Reducing our Irreproducibility. Nature. 2013;496:398 (Editorial).
- Code share: Papers in Nature journals should make computer code accessible where possible. Nature. 2014;514:536 (Editorial).
- Diggle PJ, Zeger SL. Embracing the concept of reproducible research. Biostatistics (Oxford, England). 2010;11(3):375.
- Godlee F, Groves T. The new BMJ policy on sharing data from drug and device trials. BMJ (Clinical research ed). 2012;345:e7888.
- Groves T, Godlee F. Open science and reproducible research. BMJ (Clinical research ed). 2012;344:e4383.
- McNutt M. Journals unite for reproducibility. Science (New York, NY). 2014;346(6210):679.

# Thank You!