Enough with trickle down reproducibility: scientists, open this gate! scientists, tear down this wall!

# How does reproducible research actually work in practice?

## Karthik Ram

# What exactly is reproducibility anyway?

# 4 kinds of reproducibility

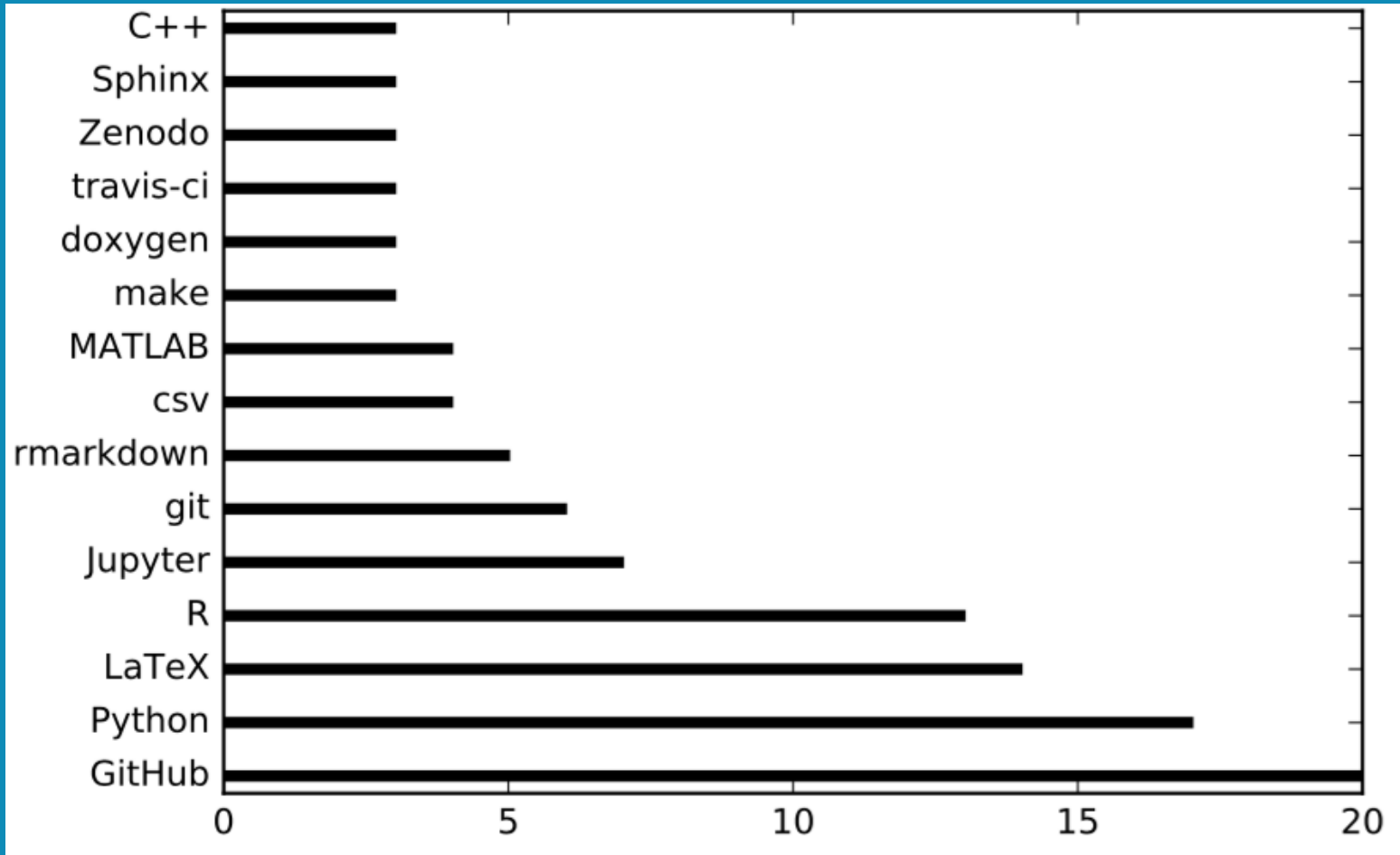*Computational reproducibility and transparency*

*Computational correctness*

*Scientific reproducibility & transparency*

*Statistical reproducibility*

*Millman et al , 2016*

# *What tools are scientists using?*

*Huff 2016*

What are some obstacles around making research reproducible?

## 1  Leveling up skills

Biggest bottleneck to adoption of reproducible research practices was related to diversity of skills

More homogeneity in tool familiarity = better reproducibility

# 2 Dependencies, build systems, and packaging

Scientific software often built on numerous dependencies

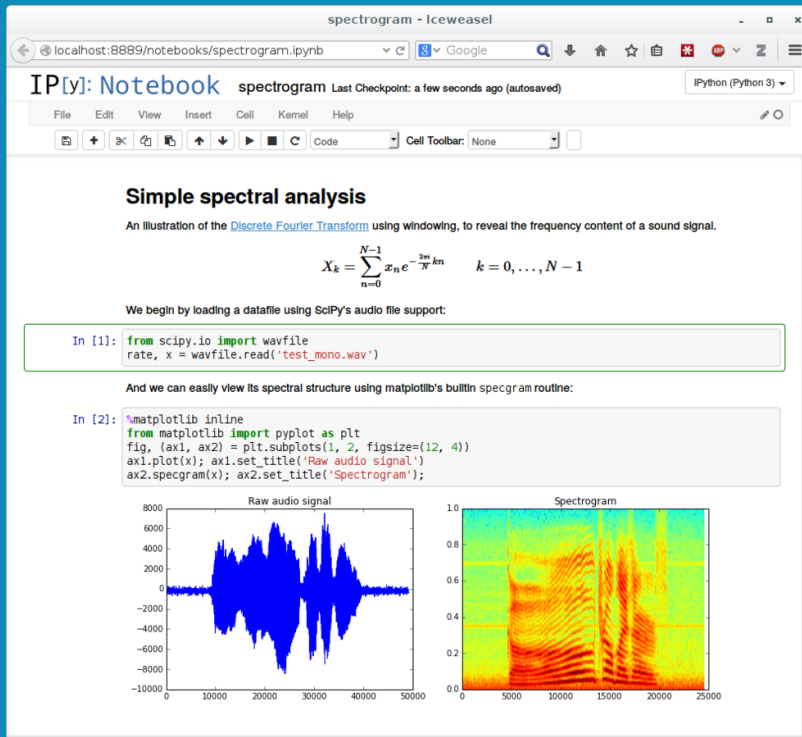Improved build systems for software, data & workflows

# 3 Testing

Code that went beyond simple script reported testing systematically

However, many scientists were discouraged by the perceived effort of unit testing

# 4 Publishing

There is still a need for publication formats that allow for effortless collaboration.

**5** Data sharing & versioning

Versioning data is hard, as is finding reliable places to archive them

# ROpenSci

| 65+ TOOLS | R, C++, NODE | LARGE CONTRIBUTOR COMMUNITY |
|---|---|---|

1. **Data retrieval** (APIs, data storage services, journals)

2. **Data visualization** (e.g. plot.ly)

3. **Data sharing** (figshare, Zenodo, dat)

4. **Reproducibility**

## 6 Time and incentives

*"time and efforts spent on creating reproducible research are not very well rewarded"*

*Ram & Marwick, 2016*

POINT OF VIEW

# How open science helps researchers succeed

**Abstract** Open access, open data, open source and other open scholarship practices are growing in popularity and necessity. However, widespread adoption of these practices has not yet been achieved. One reason is that researchers are uncertain about how sharing their work will affect their careers. We review literature demonstrating that open research is associated with increases in citations, media attention, potential collaborators, job opportunities and funding opportunities. These findings are evidence that open research practices bring significant benefits to researchers relative to more traditional closed practices.

10.7554/eLife.16800.001

# Journal of Open Source Software

JOSS   10.21105/joss.00037

The Journal of Open Source Software

Submit    Papers    About    **Karthik Ram** · Sign out

# tidytext: Text Mining and Analysis Using Tidy Data Principles in R

**Authors**

Julia Silge / David Robinson

Repository:                     Paper:                          Review:
Repository link »               PDF link »                      View review issue »

DOI:                            Status badge:                   Cite this paper:
http://dx.doi.org/10.21105/joss.00037    JOSS  10.21105/joss.00037   doi2bib

## Summary

The tidytext package (Silge, Robinson, and Hester 2016) is an R package (R Core Team 2016) for text mining using tidy data principles. As described by Hadley Wickham (Wickham 2014), tidy data has a specific structure:

joss.theoj.org

Will reproducibility always be this hard?

# Practices you can adopt now

**Version** your code
  **Automate** everywhere
  **Open** your **data**
**Document** your processes
**Test** everything

Avoid excessive dependencies

**DOIs** everywhere

Avoid spreadsheets *

Workflow and provenance
frameworks are hard to adopt

*

Hadley Wickham ✔
@hadleywickham

⚙ Following

To paraphrase @JennyBryan: teaching
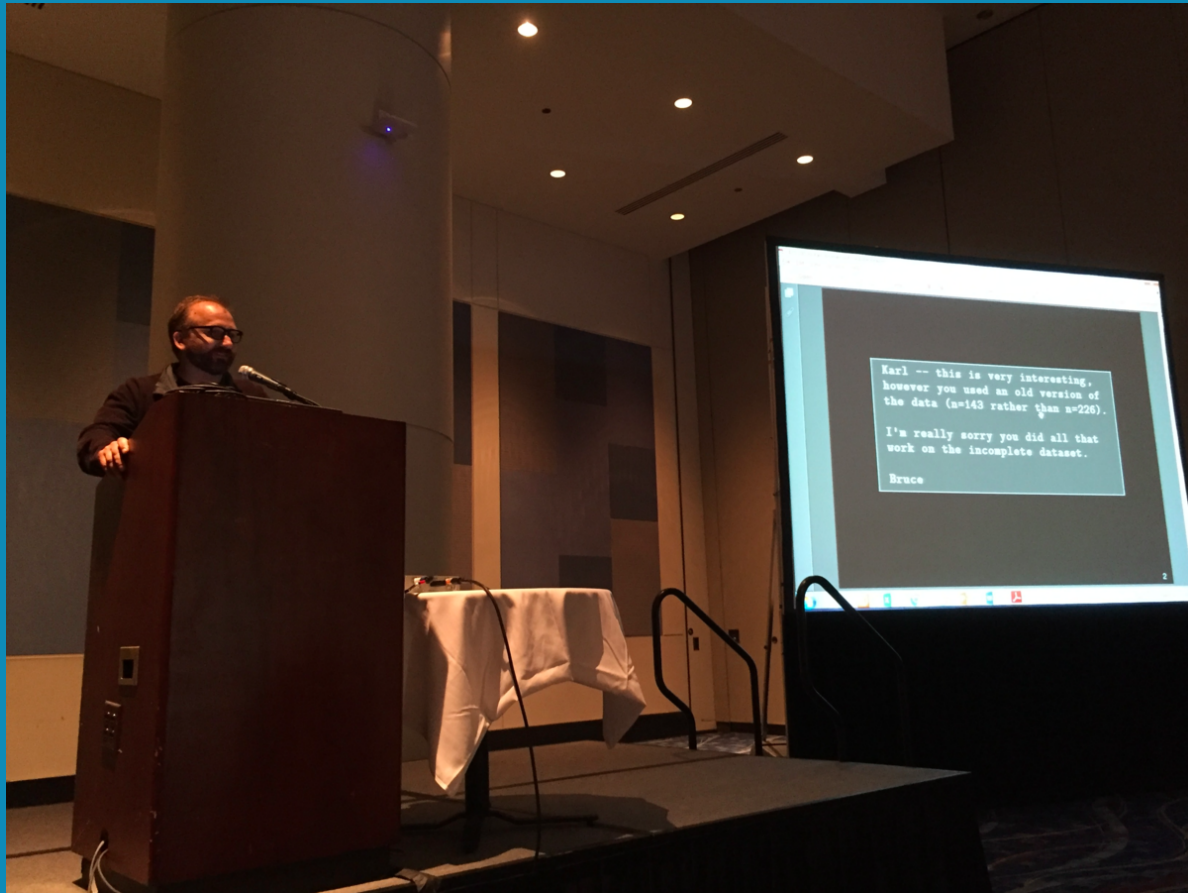spreadsheet abstinence doesn't work, need to
focus on harm reduction

# Partial reproducibility is better than nothing

Start small -- provide raw data, post any scripts,  and versions of programs you used

<div align="right">

**Karl Broman**

</div>

**What we need right now is** *scientists actually using stuff that already exists*, **not** *engineers building new stuff that no one will ever use*

<div align="right">

**C. Titus Brown**

</div>

See previous talk by Karl Broman
kbroman.org/steps2rr/

# BIDS

## BERKELEY INSTITUTE FOR DATA SCIENCE

## The Practice of Reproducible Research

*A collection of case studies to be published in spring 2017*

**inundata.org/talks/jsm2016**