

A conversation on tools for more reproducible and interactive data analysis

with Jupyter, R Markdown, and Binder demos



Dana Miller

Aims

- Discuss how researchers and collaborators in our field can take advantage of increasingly common tools to supplement and share their research and teaching materials
- Hands-on demos
- Get feedback and iterate

→ What tools or problems are *you* are most interested in?

Overview

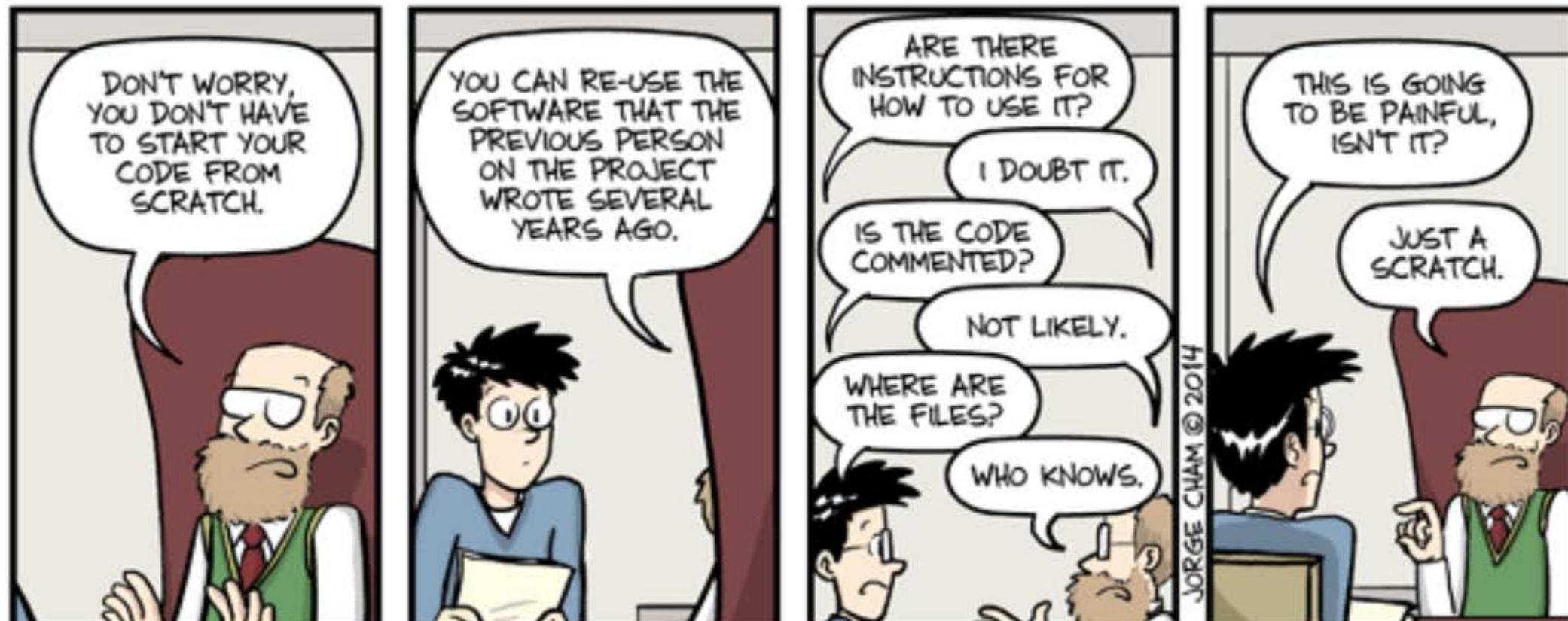
- What is reproducible research?
 - Motivating examples
 - Definitions
 - State of practice
- Literate programming tools
 - R Notebooks
 - History, key features
 - Demo
 - Jupyter Notebooks
 - History, key features
 - Demo
- Towards interactive + portable code
 - MyBinder + Docker
 - Demo
 - Cloud computer resources
- What next?

Not explicitly covered

- Version control+ collaboration (eg with git + GitHub)
 - [Separate workshop](#)
- Latest version of every tool (rapidly evolving ecosystem)
 - JupyterLab

Why reproducible research?

Motivating example 1 - productivity



Motivating example 2 - posterity

Email from grad student to a professor in 2017:

“...The citation on the slide is [your final project report on a state government website]. I downloaded that PDF and do not see [important statistic] anywhere in it. Table 3 seems to have [different value] , Table 15 has [related value] in California...

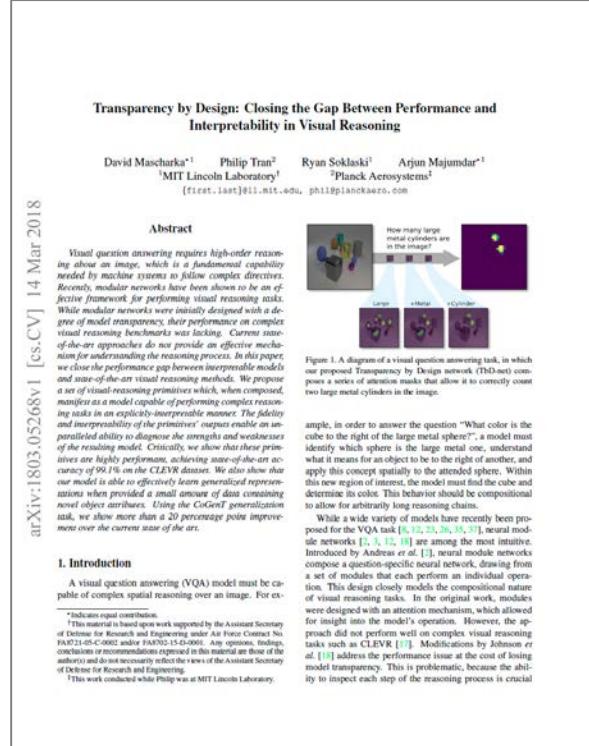
...The closest match that I've found is [other statistic] from Table 15... Is it possible there was an error in transcription? If not, can you help me figure out what's going on here? What does [important statistic] on [slide] represent?”

Response from professor:

“I don't remember the details. [Former student], do you? Thanks in advance.”

[Former student] never responded

Motivating example 3 – interactive publications!



New paper
(in this case, a preprint)

+

Website explaining paper
(has code to replicate experiments + plots and explicit software versions)

The image shows a screenshot of a GitHub repository page titled "Transparency-by-Design networks (TbD-nets)". It includes a "Launch Binder" button, Python 3.5.36, Python 3.6.5 (64-bit), and PyTorch 0.2.0 (64-bit) buttons. The repository description is "This repository contains code for replicating the experiments and visualizations from the paper Transparency by Design: Closing the Gap Between Performance and Interpretability in Visual Reasoning". It lists authors: David Mascharka, Philip Tran, Ryan Soklaski, Arjun Majumdar. The codebase provides a model architecture and code to replicate experiments. A diagram illustrates a visual reasoning task where the model counts large metal cylinders in an image. Below the diagram, there are three attention maps labeled "Large", "+Metal", and "+Cylinder". A note at the bottom encourages citation of the paper.

<https://github.com/davidmascharka/tbd-nets>

Interactive Jupyter notebook with code to apply their method to your own images

(with all required files and software ready to run in your browser with no installation required, thanks to Binder)

The image shows a Jupyter notebook interface with the title "jupyter full-vqa-example (unsaved changes)". The notebook code demonstrates the VQA pipeline using TbD-Net. It loads models, generates programs, and performs visual reasoning. A specific cell (In [1]) shows the code for loading the vocabulary and initializing the TbD-Net model. Another cell (In [2]) shows the code for generating a program and loading its checkpoint. The notebook concludes with a cell (In [3]) containing the function "run_question_and_image", which takes a question and image path as input and returns the image and intermediate attention masks.

<https://mybinder.org/v2/gh/davidmascharka/tbd-nets/binder?filepath=full-vqa-example.ipynb>

Additional reference: [Link](#) to tweet by first author publicizing this work

Demo – reproduce figures from paper on the previous slide

1. See the paper: <https://arxiv.org/abs/1803.05268>

- Note the PDF it includes link to a GitHub repo

2. Inspect the code

- Open the associated GitHub repository containing the Jupyter notebook used to produce the figures for the paper
 - Repo: <https://github.com/davidmascharka/tbd-nets>
 - Notebook: <https://github.com/davidmascharka/tbd-nets/blob/master/visualize-output.ipynb>

3. Reproduce the visualizations without installing anything!

- Click the  button on the main page of the repo

Interactive publications – in R too!

New paper
(open access)

The screenshot shows the journal article "From noise to knowledge: how randomness generates novel phenomena and reveals information" published in Ecology Letters. The article is categorized under "REVIEW AND SYNTHESIS". It includes an abstract, keywords, and a reference section. The abstract discusses the concept of noise in ecology and its role in generating new phenomena and revealing information. The keywords include Coloured noise, demographic noise, environmental noise, quasi-cycles, stochasticity, tipping points.

<https://onlinelibrary.wiley.com/doi/epdf/10.1111/ele.13085>

Corresponding “compendium”
(with text, code, data)

The screenshot shows the GitHub repository "noise-phenomena" by cboettig. The repository has 134 commits, 2 branches, and 4 releases. A preview of the "noise-phenomena" compendium is shown, which is an interactive R Markdown notebook. The notebook includes sections like "noise-phenomena compendium", "Overview", and "Appendix A: Stochastic simulations of noisy phenomena". It contains R code, mathematical equations, and explanatory text. The code includes functions for the Gillespie algorithm and a patch model.

<https://github.com/cboettig/noise-phenomena>

Interactive R Markdown notebook to explore code
(with all required files and software ready
to run in RStudio in your browser with no
installation required, thanks to Binder)

The screenshot shows an RStudio session running in a web browser (RStudio Server). The notebook "appendixA.Rmd" is open, displaying R code and its output. The code implements the Gillespie algorithm and a patch model. The RStudio interface shows the environment, history, and connections panes. The file browser sidebar lists various files related to the project, such as "appendix.pdf", "appendixA.Rmd", "Dai-Figure.R", and "gillespie.csv".

Seriously, why reproducible research?

As summarized by the ROpenSci Project:

- Show evidence of the correctness of our results
- Enable others to make use of our methods and results

<http://ropensci.github.io/reproducibility-guide/sections/introduction/>

Echoed by why Fernando Perez created Jupyter:

- ❖ **Ethical:** openness as fairness
- ❖ **Human/social:** openness fosters collaboration.
- ❖ **Epistemological:** proprietary science is an oxymoron.
- ❖ **Technical:** Python was cool :)

<https://speakerdeck.com/fperez/project-jupyter-architecture-and-evolution-of-an-open-platform-for-modern-data-science>

Attempts to establish a vocabulary for reproducibility

- ***Reproducibility*** (also called *repeatability* or *preproducibility*) [6, 7]
 - *Code, data (subject to privacy restrictions), software environment, and statistical methods are available so data analysis can be repeated and get a similar result* [1, 2, 3]
- ***Replicability*** (also called *reproducibility*)
 - *Another researcher is able to conduct the same experiment and get similar results*
- ***Correctness***
 - *Reproducible research can still be wrong!* [4]

→ There is an entire fascinating paper comparing how words used for types of reproducibility by decline [6]

1 - Peng Science 2011, DOI: 10.1126/science.1213847

2 - <http://ropensci.github.io/reproducibility-guide/sections/introduction/>

3 - <https://simplystatistics.org/2014/06/06/the-real-reason-reproducible-research-is-important/>

4 - Reproducible research can still be wrong, Leek, Peng PNAS 2015, DOI: 10.1073/pnas.1421412111

6 – Barba 2018 <https://arxiv.org/abs/1802.03311>

7 – Stark 2018 <https://www.nature.com/articles/d41586-018-05256-0>

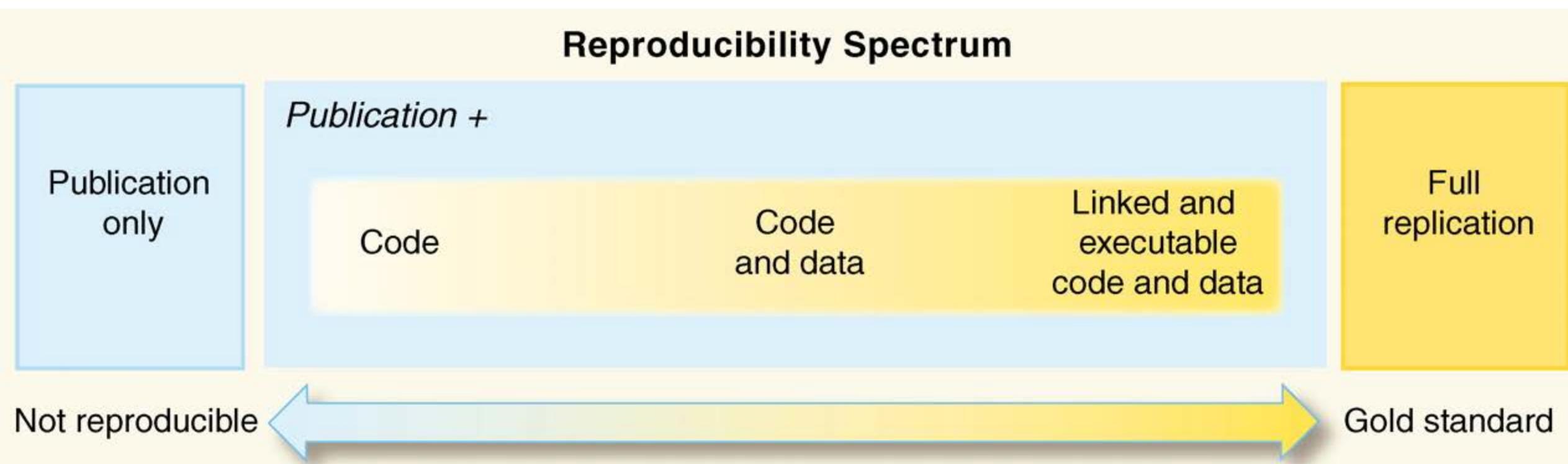
Attempts to establish a vocabulary for reproducibility within one field

Classification created by the American Society for Cell Biology [5]:

1. **analytic replication** - reproduce results by reanalysing original data
2. **direct replication** – repeat original experiment
3. **systematic replication** - repeat with different experimental conditions
(eg different cell line)
4. **conceptual replication**– repeat validity of a concept
(eg in different organisms)

5- <https://www.nature.com/news/muddled-meanings-hamper-efforts-to-fix-reproducibility-crisis-1.20076#/b1>, citing <https://www.ascb.org/wp-content/uploads/2015/11/How-can-scientist-enhance-rigor.pdf> citing Schmidt, Stefan. (2009). Shall We Really Do It Again? The Powerful Concept of Replication Is Neglected in the Social Sciences. *Review of General Psychology*. 13. 90-100. 10.1037/a0015108.

A reproducibility spectrum for computational research



Reference – even more vocabulary for types of reproducibility

“Computational reproducibility:

- when detailed information is provided about code, software, hardware and implementation details.

Empirical reproducibility:

- when detailed information is provided about non-computational empirical scientific experiments and observations. In practise this is enabled by making data freely available, as well as details of how the data was collected.

Statistical reproducibility:

- when detailed information is provided about the choice of statistical tests, model parameters, threshold values, etc. This mostly relates to pre-registration of study design to prevent p-value hacking and other manipulations.”

Reference – another reproducibility spectrum - any of these steps are helpful

Adopted from [Stodden et al 2013:](#)

Reviewable Research.

- The descriptions of the research methods can be independently assessed and the results judged credible.

Replicable Research.

- Tools are made available (might only be made available to referees or only upon request.) that would allow one to duplicate the results of the research, for example by running the authors' code to produce plots in the publication

Confirmable Research.

- The main conclusions of the research can be attained independently without the use of software provided by the author. (But using the complete description of algorithms and methodology provided in the publication and any supplementary materials.)

Auditable Research.

- Sufficient records (including data and software) have been archived so that the research can be defended later if necessary or differences between independent confirmations resolved. The archive might be private, as with traditional laboratory notebooks.

Open or Reproducible Research.

- Auditable research made openly available. This comprised well-documented and fully open code and allow one to (a) fully audit the computational procedure, (b) replicate and also independently reproduce the results of the research, and (c) extend the results or apply the method to new problems.

Reference— what is the state of practice?

PNAS

An empirical analysis of journal policy effectiveness for computational reproducibility

Victoria Stodden^{a,1}, Jennifer Seiler^b, and Zhaokun Ma^b

^aSchool of Information Sciences, University of Illinois at Urbana-Champaign, Champaign, IL 61820; and ^bDepartment of Statistics, Columbia University, New York, NY 10027

Edited by David B. Allison, Indiana University Bloomington, Bloomington, IN, and accepted by Editorial Board Member Susan T. Fiske January 9, 2018
(received for review July 11, 2017)



[Stodden, Seiler, Ma, PNAS 2018](#)

“This work evaluates the effectiveness of journal policy that requires the data and code necessary for reproducibility be made available postpublication by the authors upon request.

We assess the effectiveness of such a policy by

- (i) requesting data and code from authors and
- (ii) attempting replication of the published findings.

We chose a random sample of 204 scientific papers published in the journal Science after the implementation of their policy in February 2011.

We found that we were able to obtain artifacts from 44% of our sample and were able to reproduce the findings for 26%.”

Reference— what is the state of practice?

PNAS



An empirical analysis of journal policy effectiveness for computational reproducibility

Victoria Stodden^{a,1}, Jennifer Seiler^b, and Zhaokun Ma^b

^aSchool of Information Sciences, University of Illinois at Urbana-Champaign, Champaign, IL 61820; and ^bDepartment of Statistics, Columbia University, New York, NY 10027

Edited by David B. Allison, Indiana University Bloomington, Bloomington, IN, and accepted by Editorial Board Member Susan T. Fiske January 9, 2018
(received for review July 11, 2017)

[Stodden, Seiler, Ma, PNAS 2018](#)

Table 5. Changes in disclosure practices

Disclosure practice	2009–2010, %	2011–2012, %
Citations to data and/or code in references	25	29
Data location given in acknowledgements	29	48
Code location given in acknowledgements	4	5

Table 6. Materials availability via inspection

Materials availability	2009–2010, %	2011–2012, %
Most or all relevant data locations given	52	75
Most or all relevant software locations given	43	54
Some data and software locations given	25	45
All major software and data locations given	15	25
Code, scripts, parameters, documentation	10	12
No supporting materials available	4	1

Reference— what is the state of practice?

PNAS

An empirical analysis of journal policy effectiveness for computational reproducibility

Victoria Stodden^{a,1}, Jennifer Seiler^b, and Zhaokun Ma^b

^aSchool of Information Sciences, University of Illinois at Urbana-Champaign, Champaign, IL 61820; and ^bDepartment of Statistics, Columbia University, New York, NY 10027

Edited by David B. Allison, Indiana University Bloomington, Bloomington, IN, and accepted by Editorial Board Member Susan T. Fiske January 9, 2018
(received for review July 11, 2017)



[Stodden, Seiler, Ma, PNAS 2018](#)

Table 1. Responses to emailed requests ($n = 180$)

Type of response	Count	Percent, %
Did not share data or code:		
Contact another person	20	1 ⁺
Asked for reasons	20	1
Refusal to share	12	
Directed back to supplement	6	
Unfulfilled promise to follow up	5	
Impossible to share	3	2
Shared data and code	65	36
Email bounced	3	2
No response	46	26

~ 2% not possible to share
~ 35 % responded, did not share
~ 36 % shared
~ 28% unable to contact/ no response

Reference— what is the state of practice?

PNAS

An empirical analysis of journal policy effectiveness for computational reproducibility

Victoria Stodden^{a,1}, Jennifer Seiler^b, and Zhaokun Ma^b

^aSchool of Information Sciences, University of Illinois at Urbana-Champaign, Champaign, IL 61820; and ^bDepartment of Statistics, Columbia University, New York, NY 10027

Edited by David B. Allison, Indiana University Bloomington, Bloomington, IN, and accepted by Editorial Board Member Susan T. Fiske January 9, 2018
(received for review July 11, 2017)



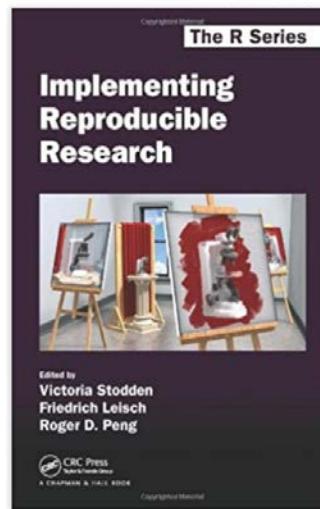
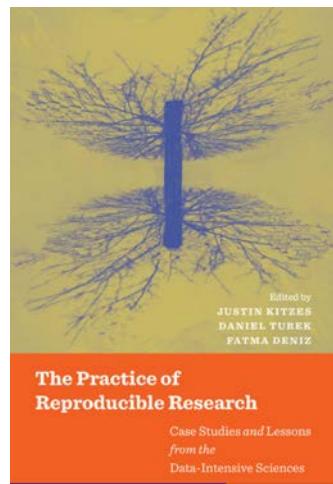
[Stodden, Seiler, Ma, PNAS 2018](#)

Table 4. Classification of reproducibility effort ($n = 22$)

Classification	Percent, %
Impossible to reproduce (missing essential code, data, or methodology)	5
Nearly impossible to reproduce (specialized hardware, intense computation requirements, sensitive data, human study, or other unavoidable reasons)	14
Difficult to reproduce because of unavoidable inherent complexity (e.g., requiring 300 million Markov chain Monte Carlo steps on each dataset, or needing months to do runs)	14
Reproducible with substantial tedious effort (e.g., individual download of a large number of datasets, hand coding of data into a new format, i.e., from an image, many archiving steps required)	5
Reproducible with substantial intellectual effort (e.g., methods well defined but required some knowledge of jargon or understanding of the field; or down the rabbit hole references to past articles required to reproduce; etc.)	5
Could reproduce with fairly substantial skill and knowledge (e.g., required GPU programming abilities to run code that wasn't given; translating complex models into MATLAB code; pseudo code with functions not detailed described in text into code; missing scripts)	23
Reproducible after tweaking (e.g., missing parameters required fiddling to find, missing modified code lines, missing arguments required for differing architecture; missing minor method step)	5
Minor difficulty in reproducing (e.g., installing a specialized library, converting to a different computational system)	18
Straightforward to reproduce with minimal effort	14

An abundance of resources + recommendations for reproducibility...

Books



Free resources

eg Software Carpentry

[R for Reproducible Scientific Analysis](#) (English & [Spanish](#))

Eg PyData videos (include R)

[CarpentryCon West @](#)

UC Davis June 29-July 1,
Includes Binder workshop

Coursera

This Course Video Transcript



Johns Hopkins University

Reproducible Research

★★★★★ 2478 ratings

Course 5 of 10 in the Specialization [Data Science](#)

<https://www.coursera.org/learn/reproducible-research/lecture/tVCoy/reproducible-research-concepts-and-ideas-part-1>

Video tutorials

(eg Pat Schloss, U Michigan)

[Riffomonas](#)

Training modules Updates About

Reproducible Research Tutorial Series

This is a series of tutorials on improving the reproducibility of data analysis for those doing microbial ecology research. Although the materials focus on issues in microbial ecology, the principles are broadly applicable. Also, this series of tutorials is not designed to teach you R or mothur. Again, although the tutorials use R and mothur, you could use other tools (e.g. Python, QIIME) to achieve the same goals. This workshop will focus on the importance of command line practices (e.g. bash), scripting languages (e.g. mothur, R), version control (e.g. git), automation (e.g. make), and literate programming (e.g. Rmarkdown). These are the tools that are used in the Schloss lab to help improve the reproducibility of our manuscripts. By completing the activities in the tutorials you will be listed on the [Reproducible Research Tutorial Honor Roll](#), which provides a certification of your training.

Tutorials

- 1. Introduction
- 2. Issues in reproducible research
- 3. First steps towards reproducibility
- 4. Using high performance computers

Reading

Much has been written on reproducibility over the past few years. These short papers provide a useful background for the overall scope of these materials and should be read before starting:

https://github.com/riffomonas/reproducible_research

Scholarly tools

The screenshot shows the Whole Tale homepage with the title 'Reproducibility Simplified' and the subtext 'Use Whole Tale to empower and share your research'. A large green button labeled 'Access' with the Whole Tale logo is prominent.

<https://dashboard.wholetale.org/>

<http://wholetale.org/runthrough.html>

statements from scholarly groups, and more...

The DataCamp logo is on the left, featuring a stylized head icon and the word 'DataCamp'. To its right is a circular badge with the text 'REPORTING WITH R MARKDOWN' and 'REPORTING RESULTS'.

Two very common (and free) interactive notebooks:

Jupyter

Sampling from the generative model

In this notebook, we will use the generative model of the HDHP (Hierarchical Dirichlet-Hawkes Process) in order to sample e-mail messages. We will start with a predefined number of users, say 10, and we will attempt to model their behavior as they are posting questions and answers on an online platform. For simplicity, our "vocabulary" will be dummy.

We start by importing all the libraries that will be required.

```
%matplotlib inline
import datetime
import string
import hdhp
import notebook_helpers
import seaborn as sns
```

Now, let us set some parameters for our model. These fall under two categories; the ones relevant to the content and then the ones relevant to the process.

```
: vocabulary = ['word' + str(i) for i in range(100)] # the `words` of our documents
doc_min_length = 5
doc_length = 10
words_per_pattern = 50

alpha_0 = (2.5, 0.75)
mu_0 = (2, 0.5)

: overlap = notebook_helpers.compute_pattern_overlap(process)
sns.distplot(overlap, kde=True, norm_hist=True, xlabel='Content overlap')

Average overlap: 0.338826769742
```



Other interactive notebooks:
Beaker (Two Sigma)
[Zeppelin](#) (Apache)
Databricks (Databricks)

R Markdown

Unit 3: Fisheries Collapse Module

This module will focus on understanding and replicating fisheries stock assessment data and fisheries collapse.

The Database

We will use data from the [RAM Legacy Stock Assessment Database](#)

First, load in the necessary libraries. Note that this time we need a package we haven't used before `readxl`. This package is useful for reading in .xls or .xlsx files. As always if you want more info on a package run `?readxl` after loading it.

```
library("tidyverse")
library("readxl")
library("scales") # For y-axis labels not in scientific notation - is there a better way to do this since 2012?
```

Reading in the tables

```
download.file("https://depts.washington.edu/ramlegac.wordpress/databaseVersions/RLSADB_v3.0_(assessment_data_only)_excel.zip",
              "ramlegacy.zip")
path <- unzip("ramlegacy.zip") # unzip the .xls files
sheets <- readxl::excel_sheets(path) # use the readxl package to identify sheet names

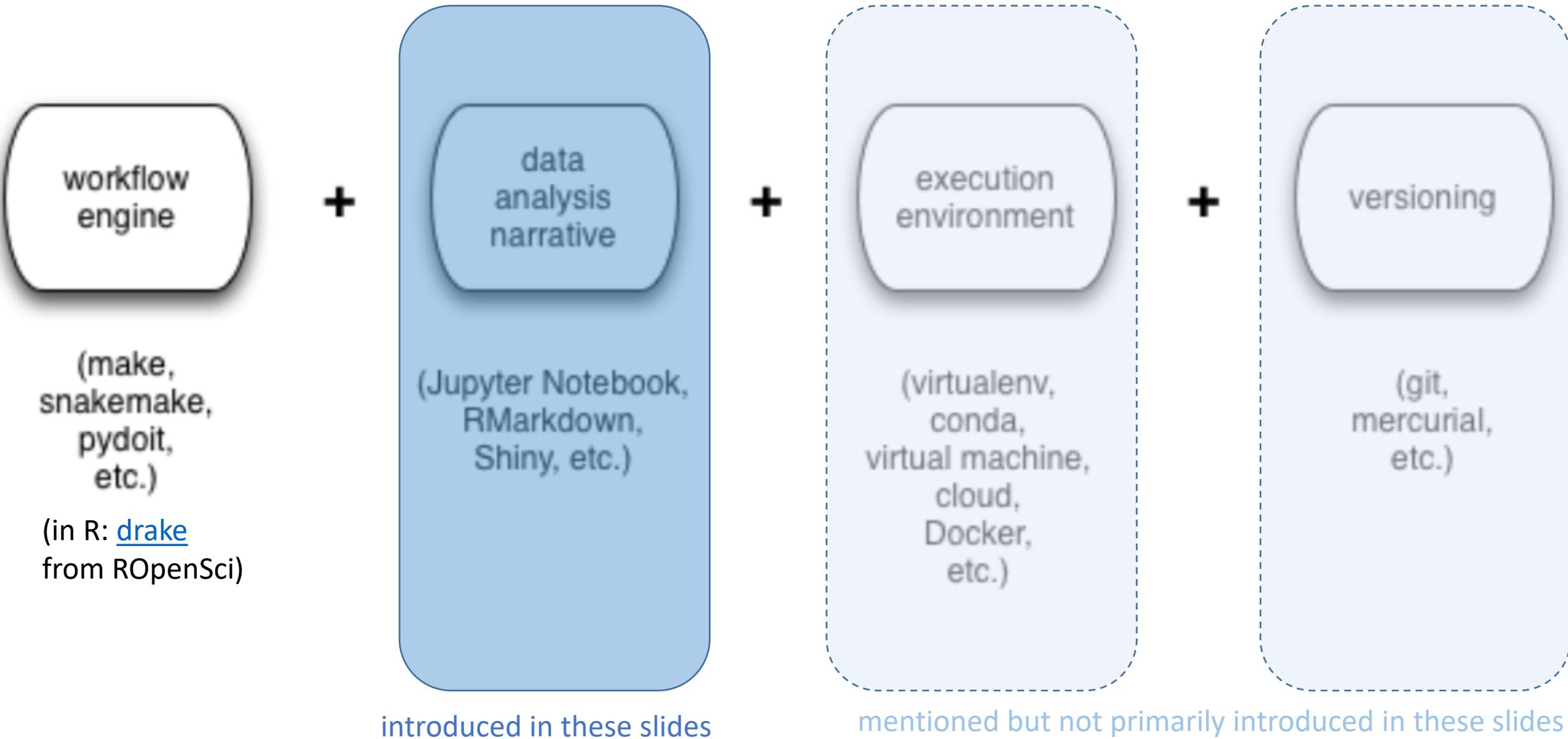
# purrr::map is the tidyverse version of lapply
ram <- lapply(sheets, readxl::read_excel, path = path) # read the data from all 3 sheets into a list
names(ram) <- sheets # give the list of datatables their assigned sheet names

## check the names
names(ram)
```

```
## [1] "area"                  "assessment"
## [3] "assessmethod"           "assessor"
## [5] "biometrics"             "bioparams"
## [7] "bioparams_ids_views"    "bioparams_units_views"
## [9] "bioparams_values_views" "management"
## [11] "stock"                  "taxonomy"
## [13] "timeseries"             "timeseries_ids_views"
## [15] "timeseries_units_views" "timeseries_values_views"
## [17] "tsmetrics"
```

```
## check our data
```


Interactive notebooks are **one part** of a toolset for reproducibility



Background– what is literate programming?



https://en.wikipedia.org/wiki/Donald_Knuth#/media/File:KnuthAtOpenContentAlliance.jpg

D. E. Knuth; Literate Programming, The Computer Journal, Volume 27, Issue 2, 1 January 1984, Pages 97–111,
<https://doi.org/10.1093/comjnl/27.2.97>

“Instead of imagining that our main task is to instruct a computer what to do, let us concentrate rather on explaining to *human beings* what we want a computer to do.”
- Donald Knuth (1984)

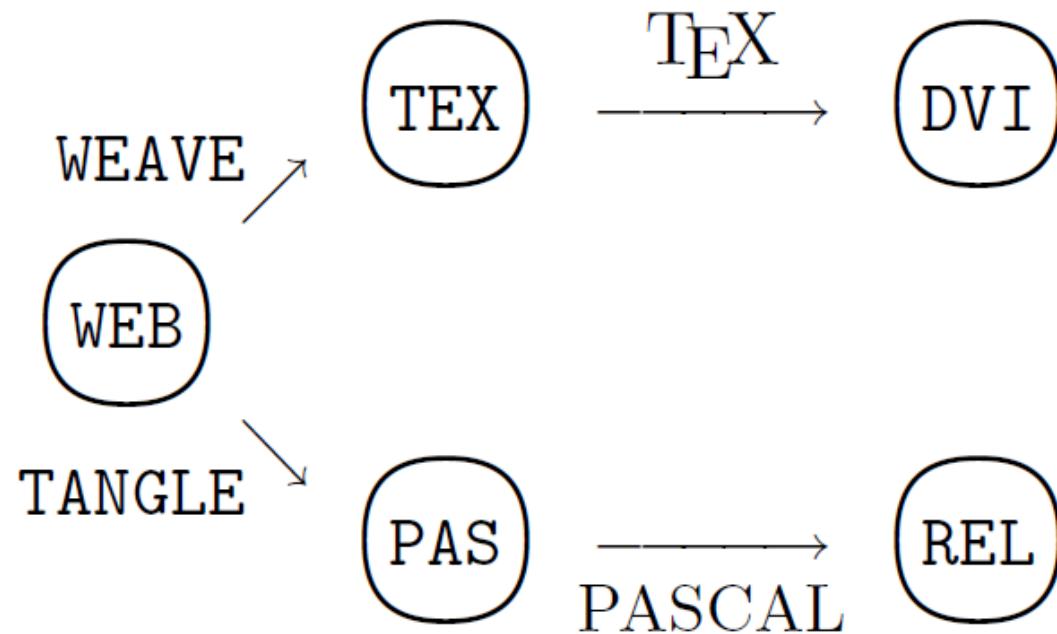


Figure 1. Dual usage of a WEB file.

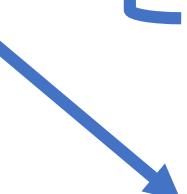
Literate programming example with R Markdown

Starting with one .Rmd file with both R code and text

This button or command:



executes all R code
(output shown in console)



The screenshot shows the RStudio interface with an R Markdown file named 'fish-assignment.Rmd'. The code includes a title, author information, and a GitHub document output. It describes the 'Overfishing Exercise' and the 'Unit 3: Fisheries Collapse Module'. The 'tidyverse', 'readxl', and 'scales' packages are loaded. The 'readxl' package is used to read an Excel file, and the 'scales' package is used for axis labeling. The 'tidyverse' package is attached, and its conflicts with 'dplyr' and 'grid' are noted. The 'readxl' package is also attached. The 'scales' package is mentioned again for axis labeling.

The screenshot shows the generated HTML document titled 'Overfishing Exercise' by 'Chelsea Andreozzi and Dana Miller'. It contains the module title 'Unit 3: Fisheries Collapse Module' and a brief description: 'This module will focus on understanding and replicating fisheries stock assessment data and fisheries collapse.' Below this, the 'The Database' section is shown, which mirrors the R code from the R Markdown file. The 'tidyverse', 'readxl', and 'scales' packages are listed, along with their respective descriptions and the 'readxl' package being used to read an Excel file. The 'scales' package is mentioned again for axis labeling.

This button or command:



renders formatted document with executed code and text (.md, .pdf, .doc etc)



R Markdown documents – key features

History

Rstudio ~ 2011

knitr ~ 2012

Markdown rendering ~ 2013

R Notebook ~ 2016

R Markdown (.Rmd)
document mixing
code and text

Code can be executed in-line
(above) or in console (below)

The screenshot shows the RStudio interface. On the left, the R Markdown file 'fish-assignment.Rmd' is open, displaying R code for data manipulation and analysis. On the right, the 'Console' tab shows the execution of the same code, with output results like 'TOTAL BIOMASS' and 'year country scientificname commonname SSB SSBUITS Total_Catch Total_Catch_Units'. A blue arrow points from the text 'Code can be executed in-line (above) or in console (below)' to the console window.

See file explorer, help documentation, plots, packages

Convert file to .md, .doc, .pdf, html, html notebook, slides,
book, blog, dashboards, Shiny app, etc

See [R Markdown output formats](#)

Version control
tracker
See [happygitwithr](#)

Exercise 1: Investigating the North-Atlantic Cod

First, we seek to replicate the following figure from the Millennium Ecosystem Assessment Project using the RAM data.

Fish landings in tons



Task 1: Joining the necessary data

To replicate this plot, we need a table with the following columns: "country", "ssb_unit", "catch_landings_unit", "scientificname", "commonname", "year", "ssb", and "TC".

Using the `select()` and `join()` functions you were introduced to in in Module 1, build a tidy table with the desired columns.

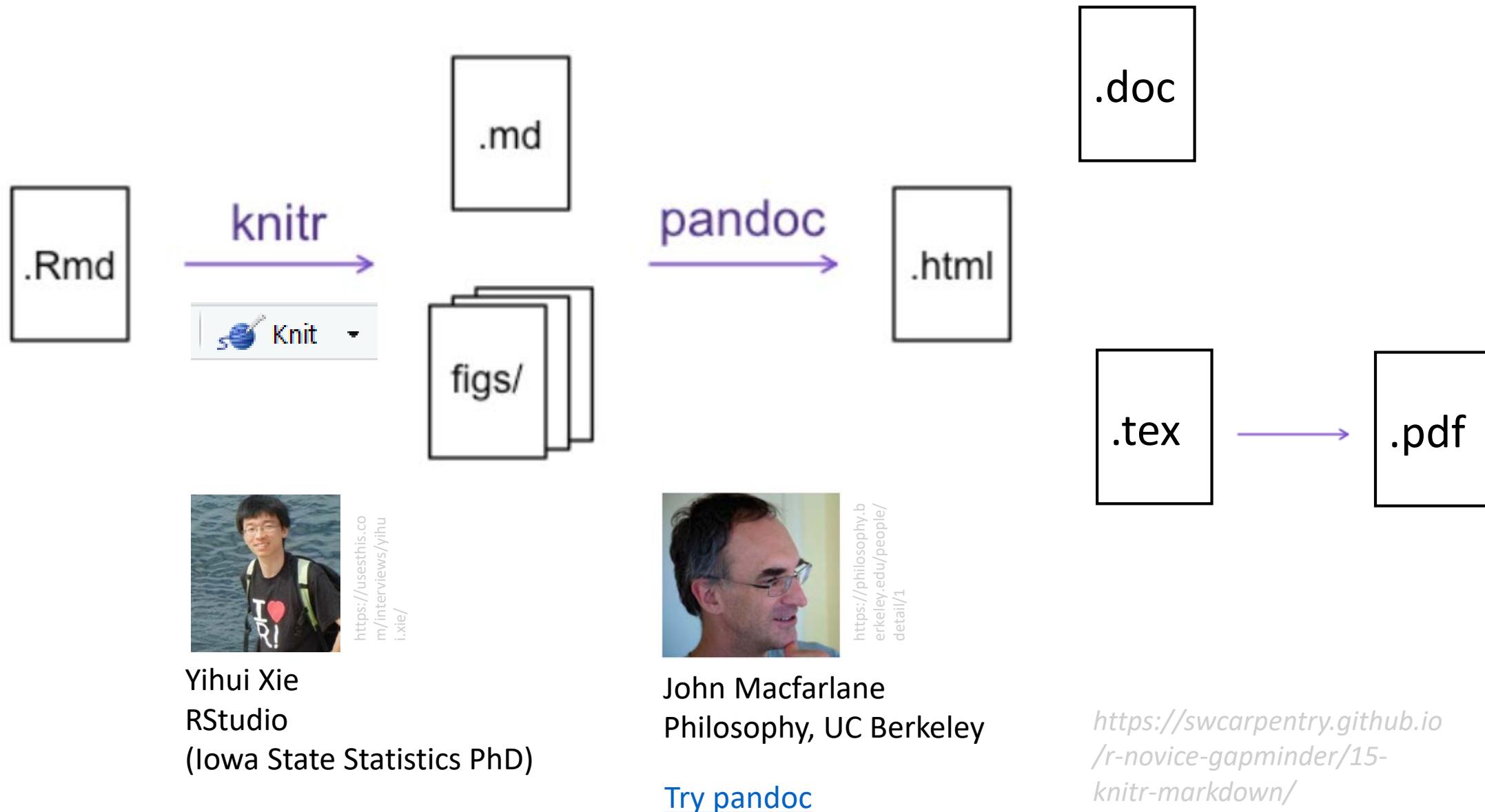
Suggestion: separate blocks within each category added first, e.g. timeseriesunits <- select(ts=TC....)

Code to look at categories of all the data tables

What is markdown?

Text using Markdown syntax	Corresponding HTML produced by a Markdown processor	Text viewed in a browser
<p>Heading =====</p> <p>## Sub-heading</p> <p>Paragraphs are separated by a blank line.</p> <p>Two spaces at the end of a line produces a line break.</p> <p>Text attributes <code>_italic_</code>, <code>**bold**</code>, <code>`monospace`</code>.</p> <p>Horizontal rule:</p> <p>---</p> <p>Bullet list:</p> <ul style="list-style-type: none">* apples* oranges* pears	<pre><h1>Heading</h1> <h2>Sub-heading</h2> <p>Paragraphs are separated by a blank line.</p> <p>Two spaces at the end of a line
 produces a line break.</p> <p>Text attributes italic, bold, <code>monospace</code>. </p></pre> <p><p>Horizontal rule:</p></p> <p><hr /></p> <p><p>Bullet list:</p></p> <pre> apples oranges pears</pre>	<p>Heading [edit]</p> <p>Sub-heading [edit]</p> <p>Paragraphs are separated by a blank line.</p> <p>Two spaces at the end of a line produces a line break.</p> <p>Text attributes <i>italic</i>, bold, <code>monospace</code>.</p> <p>Horizontal rule:</p> <p>Bullet list:</p> <ul style="list-style-type: none">• apples• oranges• pears <p>Numbered list:</p> <ol style="list-style-type: none">1. wash2. rinse3. repeat

Whoa, how does knitting work?



Demo: R Markdown

Learning objectives:

1.

- How to create a new R Markdown document in Rstudio
- Make edits
- Knit it
- Share it somewhere
- Have someone else download and run it

2.

- Example of assignments written in R Markdown

Reference to set up git w/ R: <http://happygitwithr.com>

Limitations of R + RMarkdown

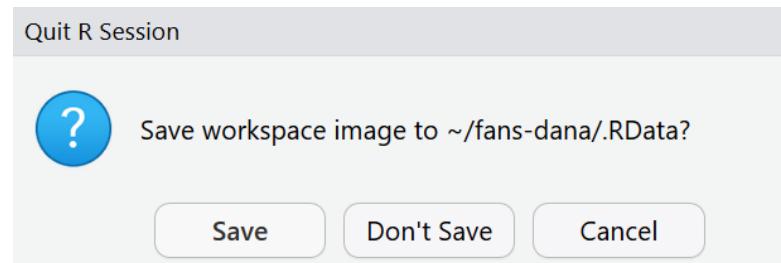
- Not a “general purpose programming language”
- By default, loads entire data into memory (not w/ sparklyr)
- Less widely used in industry than python
- Possible to run chunks out of order (but can “run all”)

Three R tips learners can start using today:

- here::here (use with R Project)
 - Freedom from hardcoded file paths!

```
library(here)
db <- read_csv(here("subfolder_name", "file_name.csv"))
```

- read_csv
 - tidyverse version of read.csv
 - won't coerce strings to factors, outputs a dataframe (tibble)
- “Save workspace image?”
 - Don’t save!



More R resources

R for Data Science book:
<https://r4ds.had.co.nz>

ROpenSci tools:
<https://ropensci.org/>

Reference for R with git:
<http://happygitwithr.com>

Stat 545:
<http://stat545.com/topics.html>

Software Carpentry
reproducible research lesson:
<https://swcarpentry.github.io/r-novice-gapminder/15-knitr-markdown/>

Packrat (package manager)

Thesis writing:
<https://eddjberry.netlify.com/post/writing-your-thesis-with-bookdown/>

<https://rosannavanhespenresearch.wordpress.com/2016/02/03/writing-your-thesis-with-r-markdown-1-getting-started/>

Website w/ blogdown:
<https://twitter.com/dsquintana/status/993410504570888192>

R notebooks online: <https://rpubs.com/>

Dplyr cheatsheet
http://stat545.com/bit001_dplyr-cheatsheet.html

Journal R templates:
<https://github.com/rstudio/rtitles>

Talk on teaching R:
<http://michaellevy.name/blog/useR-talk-on-teaching-R/>

Bonus:
One of many grad stats textbooks in R
<http://xcelab.net/rm/statistical-rethinking/>

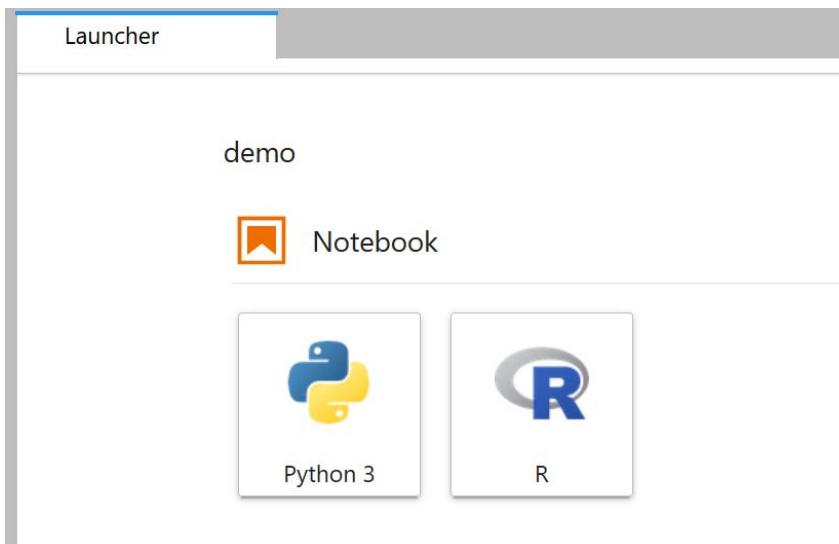
Workflow manager (like make) for R:
<https://github.com/ropensci/drake>

What about R and Python? They can be friends

Interoperability – though mastery of either takes plenty of time

Run R in Jupyter

- Run R inside python with rpy2 (formally rmagics)
- Or run R kernel (supported in Jupyter and JupyterLab)
 - [R in Jupyter demo](#)



Run Python in R

- Reticulate



```
13
14 * ````{python}
15 import pandas
16 flights = pandas.read_csv("flights.csv")
17 flights = flights[flights['dest'] == "ORD"]
18 flights = flights[['carrier', 'dep_delay', 'arr_delay']]
19 flights = flights.dropna()
20 ```
21
22 * ````{r, fig.width=7, fig.height=3}
23 library(ggplot2)
24 ggplot(py$flights, aes(carrier, arr_delay)) + geom_point() + geom_jitter()
25
26
```

<https://rstudio.github.io/reticulate/>

R and Python can be friends

Thankfully, key RStudio, python, and Jupyter developers have no patience for “language wars”

 **Fernando Perez** @fperez_org · Apr 21
Polyglot Data Science: Python, R, [@JuliaLanguage](#) and Fortran cooperating in-process, through [@IPythonDev](#); demo code: gist.github.com/fperez/5b49246..

Language wars among open tools are pointless, they complement each other.

Video/slides of the talk [@UCBIDS](#):


Project Jupyter: Architecture and Evolution of an O...
This lecture was presented at BIDS on Tuesday, April 17, 2018. Click here to view Fernando's slides for this lecture.
bids.berkeley.edu

3 replies 145 retweets 347 likes

 **Wes McKinney** ✅ @wesmckinn Follow
Very disappointing. The future for R and Python (IMHO) is collaborating on shared performance computing libraries using C/C++/LLVM

 **Hadley Wickham** ✅ @hadleywickham
Really disappointed by this @infoworld article comparing #rstats and #python: infoworld.com/article/318755...

7:23 AM - 7 Apr 2017
61 Retweets 118 Likes

9 replies 61 retweets 118 likes

 **Brian Ray** @brianray · Apr 13
As requested, I will be blogging on "Python vs R for Data Science". Meanwhile what are some points I should cover? I truly don't see this as religious war as they seem to be friendly and dually supported by most Data Scientists today as well as @ProjectJupyter

16 replies 12 retweets 46 likes

 **Hadley Wickham** ✅ @hadleywickham Follow
Replies to @brianray @ProjectJupyter
Replace "vs" with "and"

4:23 PM - 13 Apr 2018
17 Retweets 159 Likes

2 replies 17 retweets 159 likes

URSA LABS

Founded in 2018, Ursa Labs is an industry-funded development group specializing in open source data science tools. It is dedicated to advancing the state of the art in high-productivity, high-performance, cross-language software for data scientists.

Leadership



Wes McKinney
Director

Wes created the pandas project in 2008 and wrote the book *Python for Data Analysis*, helping popularize the use of Python for data science. He is a Member of The Apache Software Foundation and is a PMC member for Apache



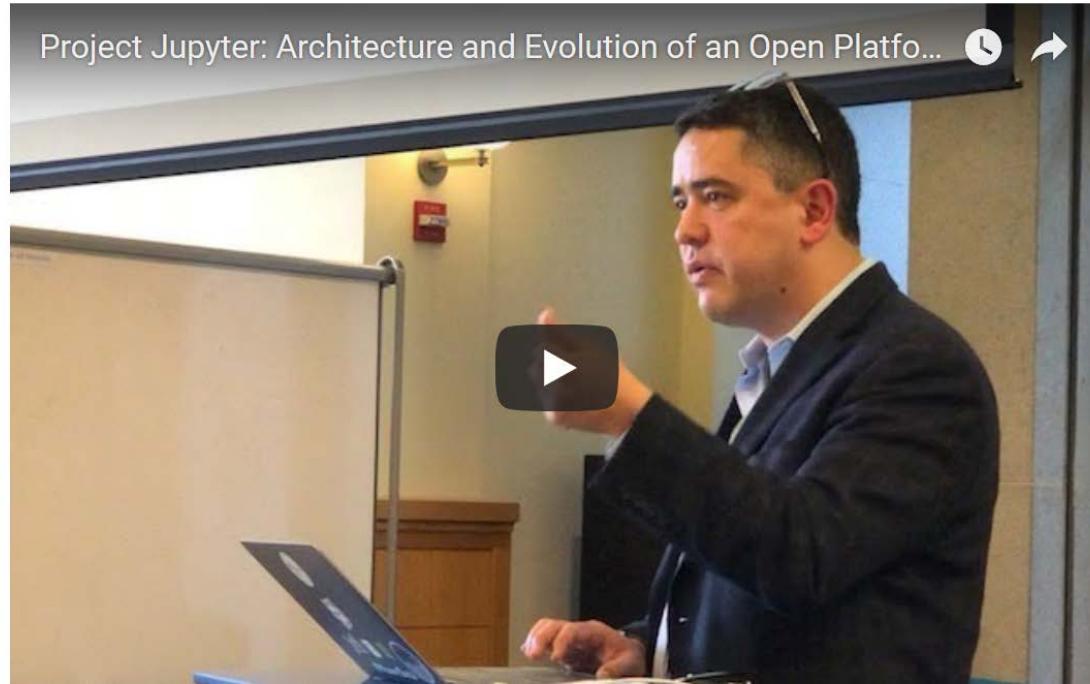
Hadley Wickham
Advisor

Hadley is the creator of many of the most widely-used R packages for data science, such as *ggplot2*, *dplyr*, and many others. He has written several books about R, such as *R for Data Science* and *Advanced R*. Hadley is the Chief Scientist at



<https://ursalabs.org/about/>

Jupyter notebooks – origins



Watch the talk:

<https://bids.berkeley.edu/events/project-jupyter-architecture-and-evolution-open-platform-modern-data-science>

Slides:

<https://speakerdeck.com/fperez/project-jupyter-architecture-and-evolution-of-an-open-platform-for-modern-data-science>

Jupyter notebooks – key features

Narrative Text

Notebook title and introduction

Description of model parameters

Description of need to profile data

Sampling from the generative model

In this notebook, we will use the generative model of the HDHP (Hierarchical Dirichlet-Hawkes Process) in order to sample events. We will start with a predefined number of users, say 10, and we will attempt to model their behavior as they are posting questions in an online platform. For simplicity, our "vocabulary" will be dummy.

We start by importing all the libraries that will be required.

```
In [1]: %matplotlib inline
import datetime
import string
import hdhp
import notebook_helpers
import seaborn as sns
```

Now, let us set some parameters for our model. These fall under two categories; the ones relevant to the content and then ones relevant to the time dynamics. Starting with the first set, we need to decide on:

- the vocabulary: a dummy set of 100 words, i.e. word0, word1, ..., word99.
- the minimum and maximum length of a question
- the number of words of each pattern

As far as the time dynamics is concerned, we need to set:

- α_0 : the parameters of the Gamma prior for the time kernel of each pattern
- μ_0 : the parameters of the Gamma prior for the user activity rate
- ω : the time decay parameter

Finally, in order to make the generative process more user-friendly, we can pre-set the number of patterns that our users can sample from.

```
In [2]: vocabulary = ['word' + str(i) for i in range(100)] # the 'words' of our documents
doc_min_length = 5
doc_length = 10
words_per_pattern = 50

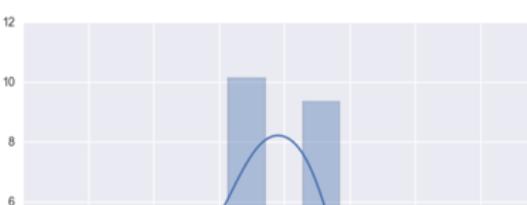
alpha_0 = (2.5, 0.75)
mu_0 = (2, 0.5)
omega = 3.5

num_patterns = 10

process = hdhp.HDHPProcess(num_patterns=num_patterns, alpha_0=alpha_0,
                           mu_0=mu_0, vocabulary=vocabulary,
                           omega=omega, words_per_pattern=words_per_pattern,
                           random_state=12)
```

Before generating any questions, we can take a look at the patterns that we initialized our process with, and look at the content distribution of each pattern. Although each pattern has a different word distribution, we can still plot the overlap (Jaccard similarity) between the words that have non-zero probability for each pattern. Since we used a limited number of patterns, the distribution of the overlap will not be smooth.

```
In [3]: overlap = notebook_helpers.compute_pattern_overlap(process)
sns.distplot(overlap, kde=True, norm_hist=True, xlabel='Content overlap')
Average overlap: 0.338826769742
Out[3]: <matplotlib.axes._subplots.AxesSubplot at 0x10de8ca50>
```



Code and Visualizations

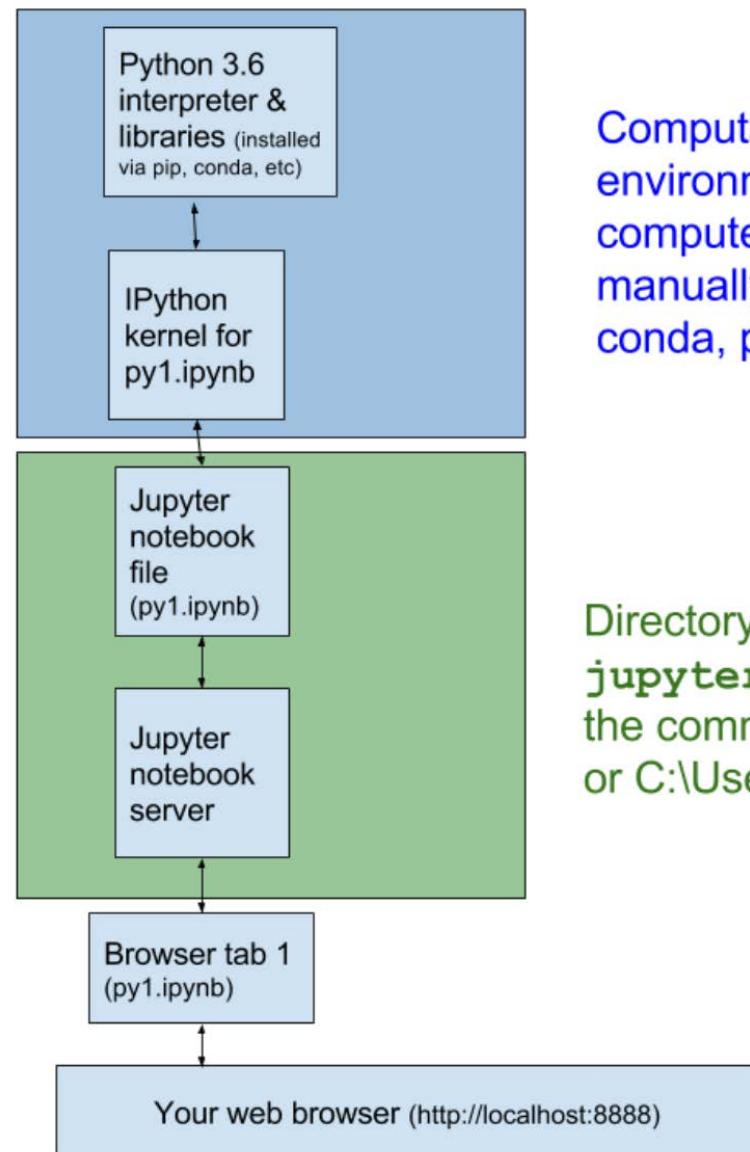
Importing external packages

Implementation of parameters

Profile plotting code

Inline plot

How do notebooks work?



Computational environment on your computer (installed manually or managed by conda, pip, apt-get, etc)

Directory where you run `jupyter notebook` from the command line (/home/stu/ or C:\Users\Stu)

Ways to share Jupyter notebooks:

- Static
 - nbviewer
 - GitHub
- Interactive
 - Download and run
 - Share with Binder

Course content in Jupyter

eg Homework in Jupyter notebooks

The screenshot shows a Jupyter Notebook interface with the title "HW2_Skeleton (unsaved changes)". The notebook contains two code cells:

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.integrate import odeint
from scipy import interp
from scipy import signal
%matplotlib inline
from __future__ import division
import pandas as pd

fs = 15; # Font Size for plots

In [ ]: ## Load Data
data=np.asarray(pd.read_csv("HW2_Data.csv",header=None))

t = data[:,0]      # t    : time vector [min]
thr = data[:,1]    # thr  : time vector [hours]
T = data[:,2]      # T    : Home interior temperature [deg C]
T_A = data[:,3]    # T_A  : Ambient outdoor temperature [deg C]
T_B = data[:,4]    # T_B  : Boiler temperature [deg C]

In [ ]: ## Problem 2(b) - Persistence of Excitation

#### OPTION with 2-D parameter vector
phi = ???? <---- enter signals for 2D parametric model
t_end = t[-1]
PE_mat = np.zeros(shape=(2,2))

phi_sq = np.zeros(shape=(len(t),2,2))
for k in range(0, len(t)):
    phi_sq[k,:,:] = phi[:,k]*(phi[:,k].T)

PE_mat[0,0] = 1/t_end * np.trapz(phi_sq[:,0,0],x=t)
PE_mat[0,1] = 1/t_end * np.trapz(phi_sq[:,0,1],x=t)
PE_mat[1,0] = ??? % <---- similarly compute other PE matrix elements
PE_mat[1,1] = ???

PE_lam_min = ??? % <---- NEED TO COMPUTE MINIMUM EIGENVALUE OF PE_mat
```

Numerous entire classes shared online:

[An introductory notebook on uncertainty quantification and sensitivity analysis \(for cardiovascular modelling\)](#)

[CFD Python: 12 steps to Navier-Stokes](#)

[Pytherm – applied thermodynamics](#)

[UCB Data 8 - Foundations of Data Science](#)

More tools for educators + researchers:

[nbgrader – assigning + grading Jupyter notebooks](#)

The screenshot shows the homepage of "The Journal of Open Source Education". The header includes the logo, the journal name, and navigation links for "Submit", "Papers", "About", and "Sign in". The main content features a large title "The Journal of Open Source Education" and a subtitle "An educator friendly journal for publishing computational learning modules and educational software."

Demo: Jupyter notebooks

Learning objectives:

1. Distinguish between static and interactive formats of sharing and viewing Jupyter notebooks

- **Static:** nbviewer or preview on GitHub
 - Example [publication](#) with [data + Jupyter notebook](#)
 - Example [notebook](#) accompanying publication
- **Interactive:** in Binder, or run locally

2. Run a Jupyter notebook locally

- *Note: requires having Python installed*

3. Write some code, save + share it

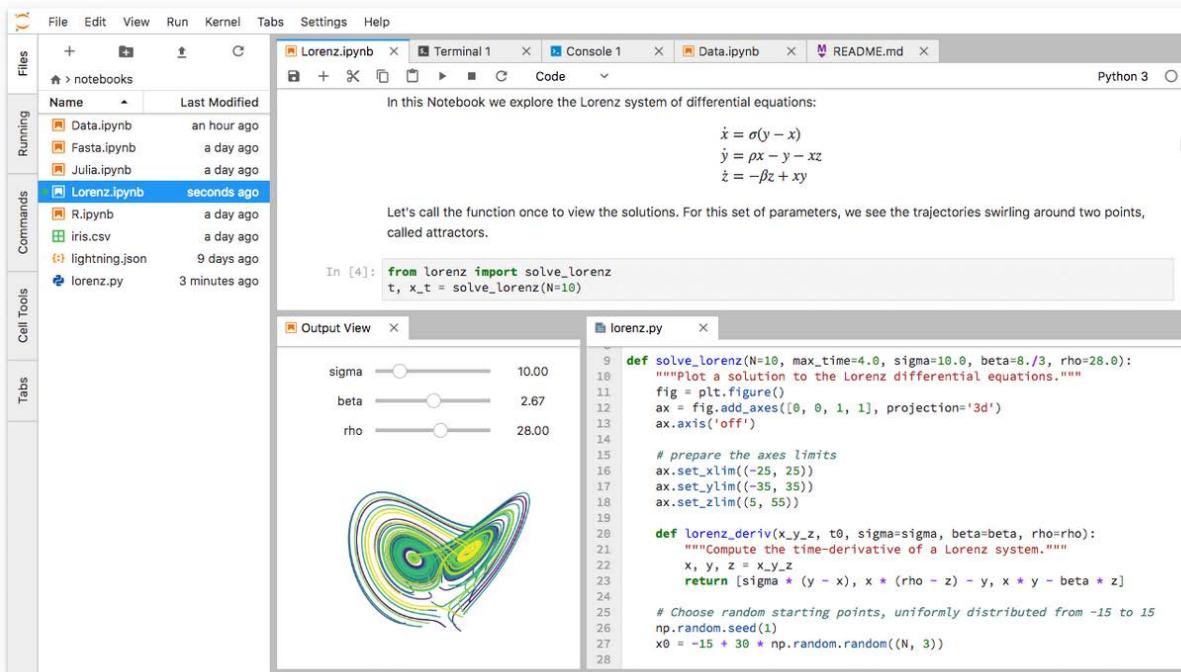
[Current citation for Jupyter](#) (2016)

Limitations of Jupyter Notebooks

- Challenging to version control
- No linting (autocorrect)
- Linear narrative structure
- Possible to run cells out of order
- Tricky to manage add-ons ([plugins](#))
- Not an environment for software engineering

The future of Jupyter

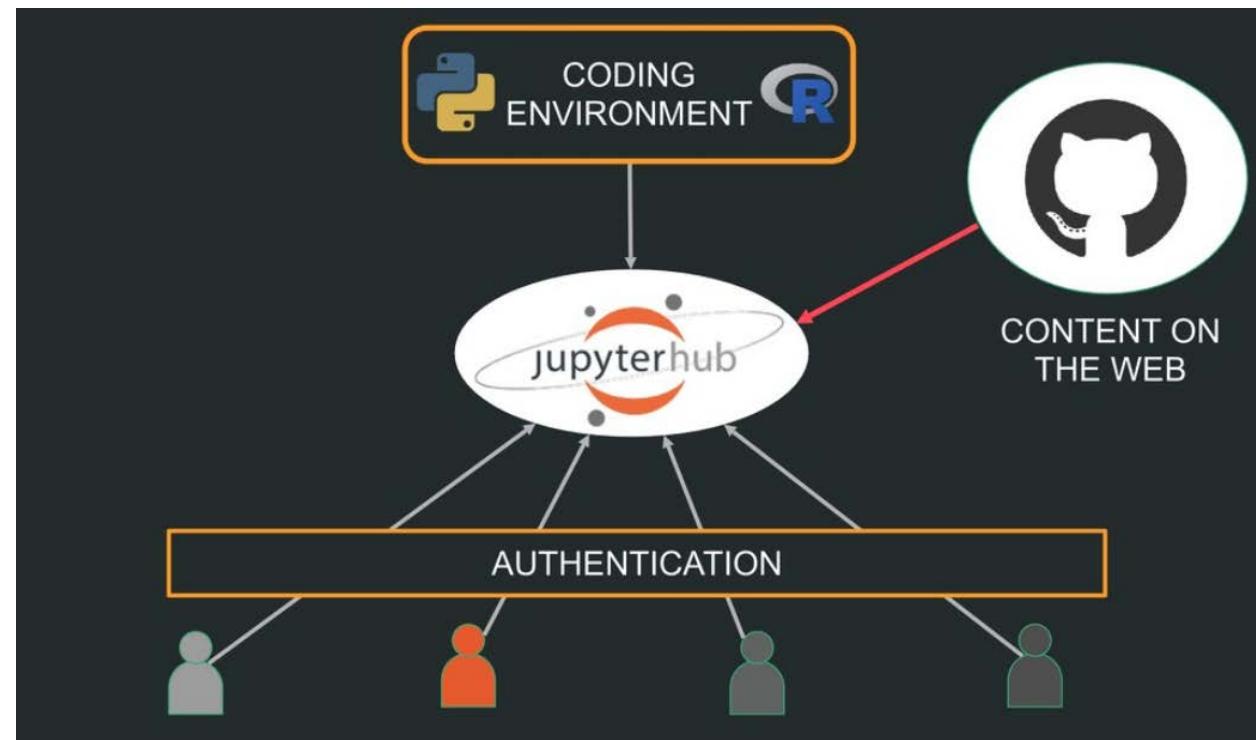
Jupyter Lab (notebook will eventually retire)



A screenshot of the Jupyter Lab interface. On the left, a sidebar shows a file tree with notebooks like 'Lorenz.ipynb' (seconds ago), 'Data.ipynb' (an hour ago), 'Fasta.ipynb' (a day ago), 'Julia.ipynb' (a day ago), 'R.ipynb' (a day ago), 'iris.csv' (a day ago), 'lightning.json' (9 days ago), and 'lorenz.py' (3 minutes ago). The main area has tabs for 'Lorenz.ipynb', 'Terminal 1', 'Console 1', 'Data.ipynb', and 'README.md'. A text cell discusses the Lorenz system of differential equations, showing the equations $\dot{x} = \sigma(y - x)$, $\dot{y} = \rho x - y - xz$, and $\dot{z} = -\beta z + xy$. Below it, a code cell runs 'lorenz.py' which imports 'solve_lorenz' and plots trajectories. A slider interface allows adjusting parameters sigma (10.00), beta (2.67), and rho (28.00). A 3D plot shows the resulting Lorenz attractor.

Launched Feb 2018 – includes demo in Docker

Jupyter Hub (no local install required, used in Data8)



Slide from Fernando Perez's talk at BIDS:
<https://speakerdeck.com/fperez/project-jupyter-architecture-and-evolution-of-an-open-platform-for-modern-data-science>

[Gallery of JupyterHub deployments](#)

Demo: Jupyter Lab

1. Go to <https://jupyter.org/try>

- Open a Jupyter Lab demo with the button shown below

- Learning objectives

- Open a notebook

- Try rearranging panes

- (can be similar to Rstudio layout)

- Appreciate scrolling smoothly through

- demo/big.csv, a file with too many rows

- (>1.1 Million) to open in Excel

Try JupyterLab



JupyterLab is the new interface for Jupyter notebooks and is ready for testing. Give it a try!

A rough comparison:

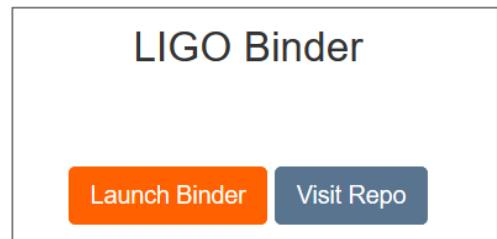
		Jupyter Notebook (.ipynb)	R Markdown document (.Rmd)
Interactive execution	In browser (after download, or hosted remotely)		In RStudio (after download, or hosted remotely)
Local install	Managing python install more complicated (use conda)		Usually straightforward
Format	<u>JSON</u>		Markdown
Version control	Improving (line-by-line diffs less nice, nbdime released)		Seamlessly integrated (line-by-line diffs, GUI in RStudio)
Environment	Python and packages on computer (local or remote), or per install with virtualenv		R and packages on computer (local or remote), or per project with packrat
Convert to output	Can download as pdf, markdown, HTML, etc or use nbconvert (requires installs)		Conversion to more types built in to RStudio Configure rendering with YAML header
Become a website	Yes, eg nbviewer, GitHub		Yes, “can publish to any static web host service”, eg RPubs, GitHub
Other languages	Over 100 kernels, including R, Or run R in python with r2py		R, python, bash, C++, SQL, can call functions from Julia, Fortran...
IDE	Notebook – no, Lab – yes, based on RStudio		Beautiful, Includes linting, integrated help
Render on GitHub	.ipynb – code and output		.Rmd – code only, not output .md + .png files – code, output, and images
Progress bar	Not by default (available as plugin)		Yes – shows progress when executing code

Demo: Notebooks in containers on cloud compute resources with Binder

1. Open one of these links to data from a famous gravity wave experiment:

<https://notebooks.gesis.org/>, alternate link

2. Click “Launch Binder”



3. Open index.ipynb

- *Ctrl + enter to run cells of code*

<https://mybinder.org/>



already ~50,000 users/week ([source](#))

Try it: <https://mybinder.org/>

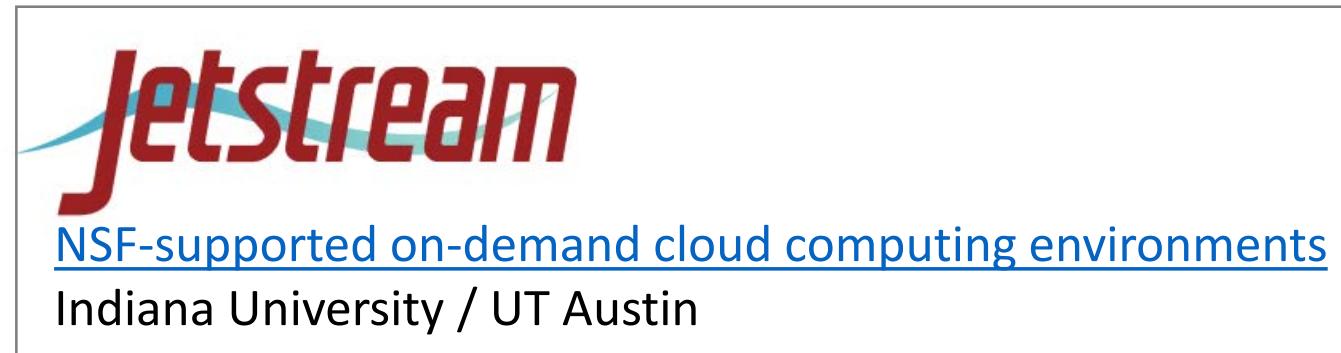
More examples:

<http://ivory.idyll.org/blog/tag/reproducibility.html>

Demo: RStudio running on cloud compute resources (eg AWS, Jetstream)

RStudio in Binder demo:

- [Link](#) from recent publication



Demo above uses [Rocker project](#):

- “provides a widely-used suite of Docker images with customized R environments for particular tasks”
- Over 3 Million installs!

image	description	size	metrics	build status
r-ver	Specify R version in docker tag. Builds on debian:stable	221.8MB 8 layers	docker pulls 23k	docker build automated
rstudio	Adds rstudio	0B 17 layers	docker pulls 1M	docker build automated

<https://www.rocker-project.org/images/>

An Introduction to Rocker: Docker Containers for R
[Carl Boettiger, Dirk Eddelbuettel, 2017, arXiv:1710.03675](#)

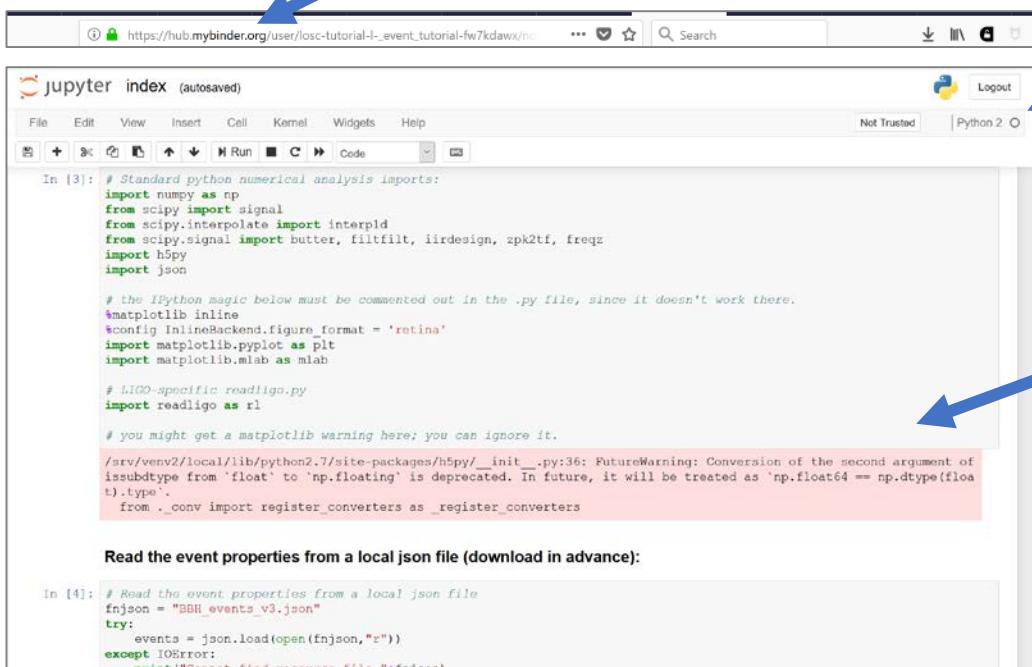
How does Binder work?

“Binder allows you to create custom computing environments that can be shared and used by many remote users”

<https://mybinder.readthedocs.io/en/latest/>

Visible in your browser,
but running on a remote machine

Your browser



A screenshot of a web browser window displaying a Jupyter notebook. The URL in the address bar is https://hub.mybinder.org/user/lsc-tutorial-l-event_tutorial-fw7kdawx/notebooks/index.ipynb. The browser title bar shows "jupyter index (autosaved)". The main content area displays a Jupyter notebook with two code cells. The first cell contains Python code for importing various scientific computing libraries like numpy, scipy, and h5py. The second cell contains code for reading event properties from a local JSON file. A red box highlights the output of the second cell, which shows a warning message about deprecated dtype conversion.

```
In [3]: # Standard python numerical analysis imports:  
import numpy as np  
from scipy import signal  
from scipy.interpolate import interp1d  
from scipy.signal import butter, filtfilt, iirdesign, zpk2tf, freqz  
import h5py  
import json  
  
# the IPython magic below must be commented out in the .py file, since it doesn't work there.  
%matplotlib inline  
config InlineBackend.figure_format = 'retina'  
import matplotlib.pyplot as plt  
import matplotlib.mlab as mlab  
  
# LIGO-specific readligo.py  
import readligo as rl  
  
# you might get a matplotlib warning here; you can ignore it.  
  
/srv/venv2/local/lib/python2.7/site-packages/h5py/_init_.py:36: FutureWarning: Conversion of the second argument of isubdtype from 'float' to 'np.floating' is deprecated. In future, it will be treated as 'np.float64 == np.dtype(float).type'.  
from .conv import register_converters as _register_converters  
  
Read the event properties from a local json file (download in advance):  
  
In [4]: # Read the event properties from a local json file  
fnjson = "GBH_events_v3.json"  
try:  
    events = json.load(open(fnjson,"r"))  
except IOError:  
    print("Cannot find resource file: %s" % fnjson)
```

Specified version
of Python

You see the code in
the .ipynb file

How does Binder work?

Code for notebook and list of dependencies and files available on GitHub (right now Binder only works with public repositories)



Your browser

The screenshot shows a Jupyter notebook interface running on a web browser. On the left, a code cell (In [3]) contains Python code for numerical analysis imports, including NumPy, SciPy, and Matplotlib. A warning message is present about the use of IPython magic in a .py file. On the right, a GitHub commit history for the 'ligo-binder' repository is displayed, showing several commits related to the LIGO GW150914 tutorial, such as adding data files and committing Dockerfiles. At the bottom, another code cell (In [4]) shows code for reading event properties from a local JSON file.

```
# Standard python numerical analysis imports:  
import numpy as np  
from scipy import signal  
from scipy.interpolate import interp1d  
from scipy.signal import butter, filtfilt, iirdesign, zpk2tf, freqz  
import h5py  
import json  
  
# the IPython magic below must be commented out in the .py file, since it doesn't work there.  
%matplotlib inline  
config InlineBackend.figure_format = 'retina'  
import matplotlib.pyplot as plt  
import matplotlib.mlab as mlab  
  
# LIGO-specific readligo.py  
import readligo as rl  
  
# you might get a matplotlib warning here; you can ignore it.  
  
/srv/venv2/local/lib/python2.7/site-packages/h5py/_init_.py:36: FutureWarning: Conversion of the second argument of isubdtype from 'float' to 'np.floating' is deprecated. In future, it will be treated as 'np.float64 == np.dtype(float).type'.  
from .conv import register_converters as _register_converters
```

Read the event properties from a local json file (download in advance):

```
# Read the event properties from a local json file  
fnjson = "GBH_events_v3.json"  
try:  
    events = json.load(open(fnjson,"r"))  
except IOError:  
    print("Cannot find resource file: %s" % fnjson)
```

How does Binder work?



builds Docker image based on repo and generates URL for public access

Code and dependencies
in public repository



Your browser

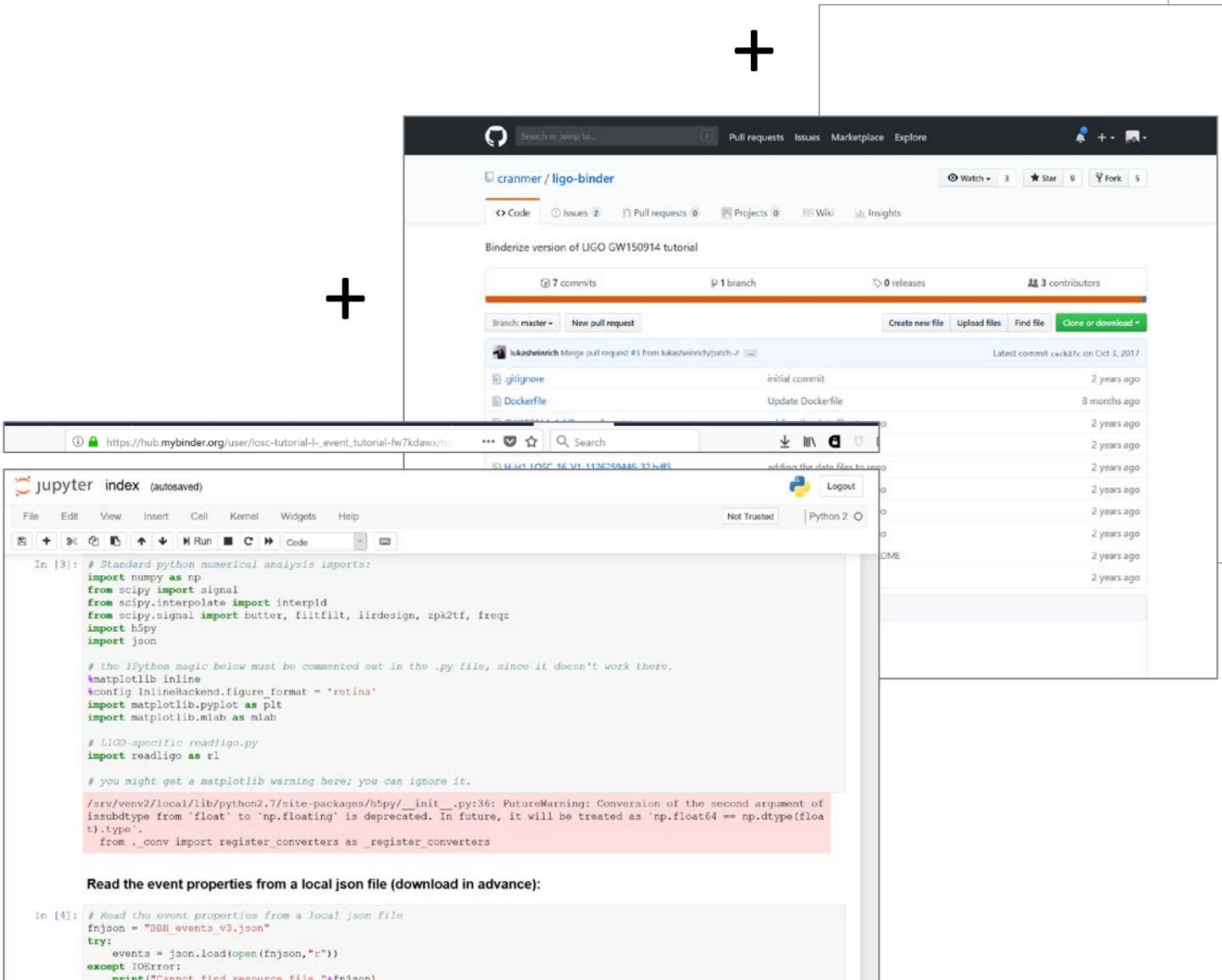
```
# Standard python numerical analysis imports:  
import numpy as np  
from scipy import signal  
from scipy.interpolate import interp1d  
from scipy.signal import butter, filtfilt, iirdesign, zpk2tf, freqz  
import h5py  
import json  
  
# the IPython magic below must be commented out in the .py file, since it doesn't work there.  
%matplotlib inline  
config InlineBackend.figure_format = 'retina'  
import matplotlib.pyplot as plt  
import matplotlib.mlab as mlab  
  
# LIGO-specific readligo.py  
import readligo as rl  
  
# you might get a matplotlib warning here; you can ignore it.  
  
/srv/venv2/local/lib/python2.7/site-packages/h5py/_init_.py:36: FutureWarning: Conversion of the second argument of isubdtype from 'float' to 'np.floating' is deprecated. In future, it will be treated as 'np.float64 == np.dtype(float).type'.  
from .conv import register_converters as _register_converters  
  
Read the event properties from a local json file (download in advance):  
  
In [4]: # Read the event properties from a local json file  
fnjson = "BBH_events_v3.json"  
try:  
    events = json.load(open(fnjson,"r"))  
except IOError:  
    print("Cannot find resource file: %s" % fnjson)
```



required
software
+
packages

included
data

How does Binder work?



Generates URL for temporary instance



required
software
+
packages

included
data

Generated binders are hosted by a hub that deploys, scales, and manages compute resources using [Kubernetes](#)

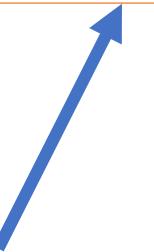
(can run on any cloud platform provider)
More details: [BinderHub](#)



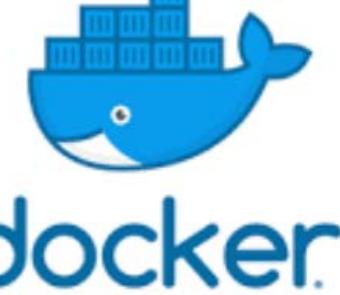
JupyterHub



Kubernetes



Background: What are (software) containers?



- Commonly used metaphor of global shipping industry + standardized shipping containers
- **A container contains “code, configurations, and dependencies” [3],**
- like ‘putting down a cloth over your table (Operating System) before installing anything’
- Enables multiple isolated sets of code, each of which can only access the files and resources allocated to it, all able to run on a shared operating system kernel

1- https://en.wikipedia.org/wiki/Operating-system-level_virtualization

2 - <https://cloud.google.com/containers/>

3 - <https://aws.amazon.com/what-are-containers/>

A few more resources...

- “Our path to better science in less time using open data science tools”, 2017,
<https://www.nature.com/articles/s41559-017-0160>
- “Developing a modern data workflow for living data”
<https://www.biorxiv.org/content/early/2018/06/12/344804>

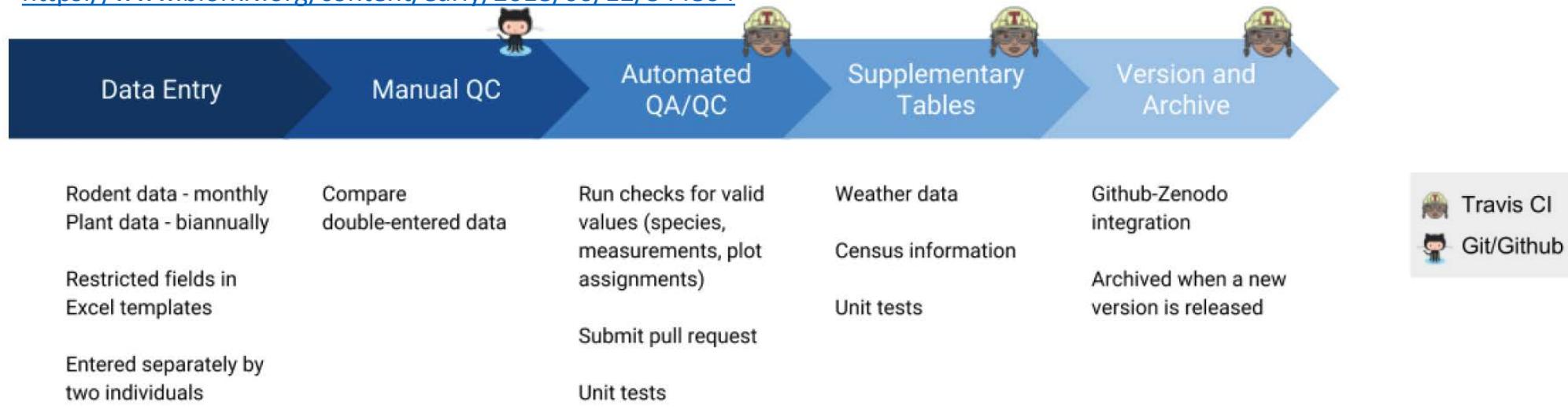


Figure 1: Our data workflow

- “Good enough practices in scientific computing” 2017, <http://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1005510>
- “Ten simple rules for making research software more robust”, 2017, <http://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1005412>
- “Ten quick tips for teaching programming” 2018 <http://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1006023>
- UO Psychology data science tutorials gallery https://robchavez.github.io/datascience_gallery/

Sidebar – who maintains or pays for all this software?

Key Python libraries for science + data analysis:



Screenshot from
Prof. Kelle Cruz's (Hunter
College/CUNY) Simons
Foundations Lecture
(link to [slides](#))

One model of support: an umbrella non-profit (mostly python)



OPEN CODE = BETTER SCIENCE

- Funding
- Stability
- Independence
- Logistical support
 - Fiscal + Legal
 - Community

[FISCALLY SPONSORED PROJECTS]



“From Netflix to NASA, researchers use our open source tools to solve the most challenging problems.”

Another model:  employs >25 software engineers/developers

Shiny, Rpubs



Joe Cheng • 3rd
CTO at RStudio, Inc.
Greater Seattle Area

R Graphics Cookbook



Winston Chang • 3rd
Software Developer at RStudio
Greater Minneapolis-St. Paul Area

R for Data Science co-author



Garrett Grolemund •
Data Scientist and Master
Greater Nashville Area, TN

Tidyverse, R for Data Science



Hadley Wickham • 3rd
Chief Scientist at RStudio, Inc.
Houston, Texas Area

knitr, bookdown, blogdown



Yihui Xie • 3rd
Software Engineer at RStudio, Inc.
Greater Omaha Area

RMarkdown, R + Tensorflow, reticulate (python in R)



JJ Allaire • 3rd
CEO at RStudio
Greater Boston Area



Jenny Bryan
STAT545 resources,
readxl

StationaRy



Richard Iannone • 2nd
Software Engineer at RStudio, Inc.

Sparklyr (Spark in R)



Javier Luraschi
javierluraschi

Source: LinkedIn as of May 2018



has been supported by grants, employers, and donations

Jupyter steering council, 2018



Damian Avila
Anaconda, Inc.
[@damianavila](#) on GitHub



Matthias Bussonnier
UC Berkeley
[@carreau](#) on GitHub



Sylvain Corlay
QuantStack
[@sylvaincorlay](#) on GitHub



Brian Granger
Cal Poly, San Luis Obispo
[@ellisonbg](#) on GitHub



Jason Grout
Bloomberg
[@jasongrout](#) on GitHub



Jessica Hamrick
DeepMind
[@jhamrick](#) on GitHub



Paul Ivanov
Bloomberg
[@ivanov](#) on GitHub



Kyle Kelley
Netflix
[@rgbkrk](#) on GitHub



Thomas Kluyver
University of Southampton
[@takluyver](#) on GitHub



Peter Parente
Valassis Digital
[@parente](#) on GitHub



Fernando Perez
UC Berkeley
[@fperez](#) on GitHub



Min Ragan-Kelley
Simula Research Lab
[@minrk](#) on GitHub



Ana Ruvalcaba
Cal Poly, San Luis Obispo
[@ruv7](#) on GitHub



Steven Silvester
JPMorgan Chase
[@blink1073](#) on GitHub



Carol Willing
Cal Poly
[@willingc](#) on GitHub

Sponsors

Project Jupyter receives direct funding from the following sources:

THE LEONA M. AND HARRY B.
HELMESLEY
CHARITABLE TRUST

ALFRED P. SLOAN
FOUNDATION

GORDON AND BETTY
MOORE
FOUNDATION

RACKSPACE

fastly

Horizon 2020
European Union Funding
for Research & Innovation

Microsoft

Google

Institutional Partners

Institutional Partners are organizations that support the project by employing Jupyter Steering Council members. Current Institutional Partners include:

ANACONDA

Bloomberg

NETFLIX

CAL POLY
SAN LUIS OBISPO

Berkeley
UNIVERSITY OF CALIFORNIA

QuantStack
Scientific Computing

JPMORGAN CHASE & CO.

Donations

Jupyter will always be 100% open source software, free for all to use and released under the liberal terms of the [modified BSD license](#). If you have found Project Jupyter to be useful in your work, research or company, please consider making a donation to the project commensurate with your resources.

All donations will be used strictly to fund the development of Project Jupyter's open source software, documentation and community. Our donations are managed by the [NumFOCUS Foundation](#), which is the legal and fiscal umbrella for the project. NumFOCUS is a 501(c)3 non-profit foundation; if you are subject to US Tax law, your contributions will be tax-deductible.

NUMFOCUS
OPEN CODE = BETTER SCIENCE

Support Project Jupyter

Where do we go from here?

ROpenSci's perspective:

- **Train students** by putting homework, assignments & dissertations on the reproducible research spectrum
- **Publish examples** of reproducible research in our field
- **Request code & data** when reviewing
- **Submit to & review for journals** that support reproducible research
- **Critically review & audit** data management plans in grant proposals
- Consider reproducibility wherever possible in **hiring, promotion & reference letters.**

<http://ropensci.github.io/reproducibility-guide/sections/introduction/> - Leveque et al

Where do we go from here?

Prof. Kelle Cruz's perspective

The image shows a presentation slide with a dark background. At the top, the word "Conclusions" is written in a large, white, sans-serif font. Below it is a bulleted list of six items, also in white text. To the right of the list are three logos: the GitHub logo (a black cat icon), the Python logo (a stylized "P" in blue and yellow), and the NumFOCUS logo (the word "NUMFOCUS" in red and blue with "OPEN CODE = BETTER SCIENCE" below it). The slide has a thin white border.

- Science is software
- GitHub facilitates open development and has changed the way we interact with software.
- Python is rapidly becoming the common language of science.
- NumFOCUS needs to be recognized as a vital part of science infrastructure.
- Software development needs to be taught and integrated into the formal STEM curriculum.

Screenshot from
Prof. Kelle Cruz's (Hunter
College/CUNY) Simons
Foundations Lecture
(link to [slides](#))

Where do we go from here?

Prof. Kelle Cruz's perspective

Implications for the future of academic science

- Software development skills will be considered fundamental & essential.
- Software will be considered a research product, similar to a journal article.
- Code review will part of peer review process.
- All code will be expected to be open and accessible.
- Contributing to open, collaborative projects will be the norm.



Screenshot from
Prof. Kelle Cruz's (Hunter
College/CUNY) Simons
Foundations Lecture
(link to [slides](#))

Discuss: Where do we go from here?

Bonus slides!

(for reference)

Preprint: “Terminologies for reproducible research”, Barba 2018

<https://arxiv.org/abs/1802.03311>

“...authors either,

A—make no distinction between the words reproduce and replicate, or

B—use them distinctly.

If B, then they are commonly divided in two camps.

In a spectrum of concerns that starts at a minimum standard of

“same data + same methods=same results,” to

“new data and/or new methods in an independent study=same findings,”

group 1 calls the minimum standard reproduce, while group 2 calls it
replicate...”

[spacing added for clarity]

Preprint: “Terminologies for reproducible research”, Barba 2018

<https://arxiv.org/abs/1802.03311>

Conclusions

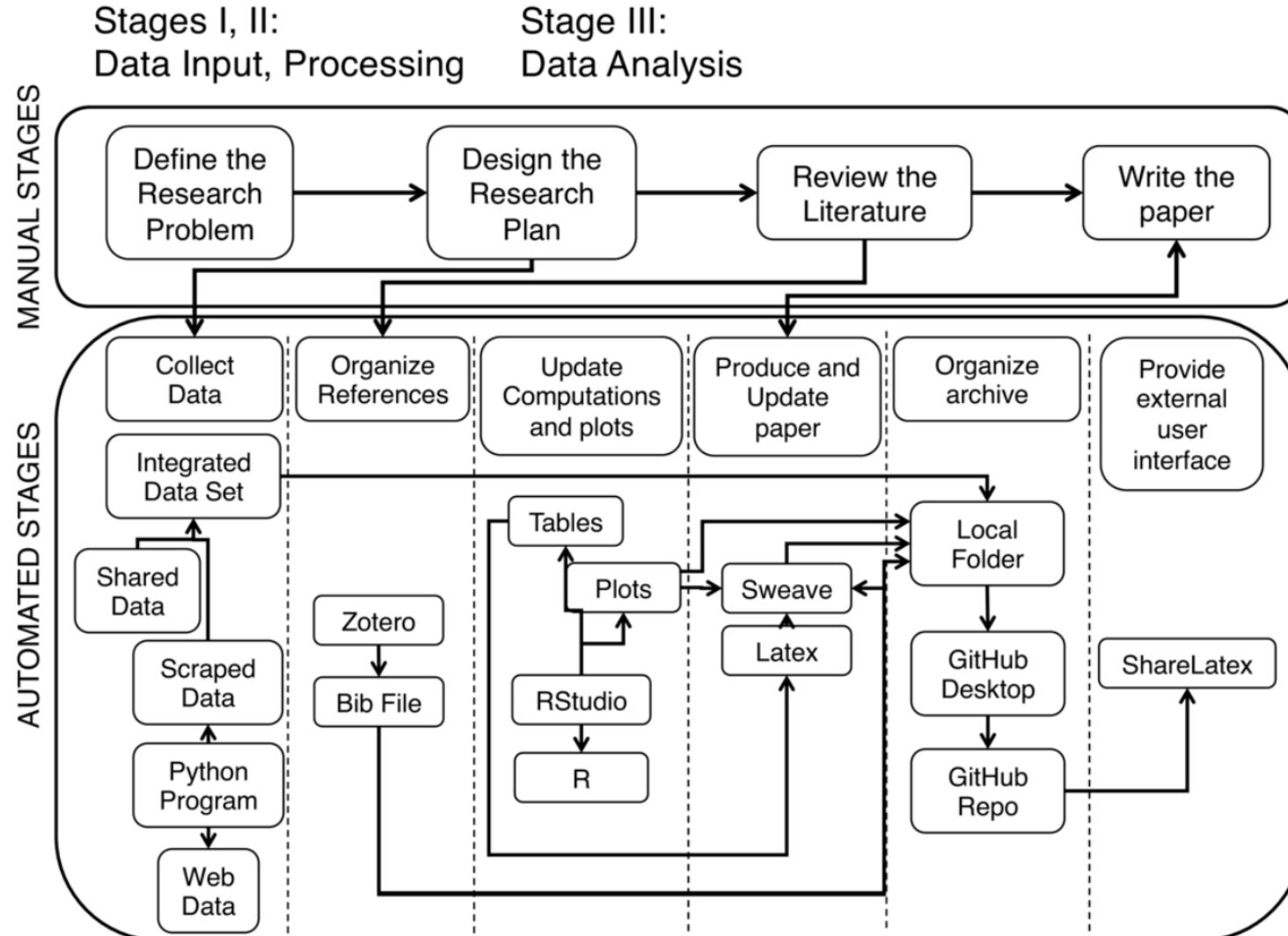
The Claerbout/Donoho/Peng terminology is broadly disseminated across disciplines (see Table 2). But the recent adoption of an opposing terminology by two large professional groups—ACM and FASEB—make standardization awkward. The ACM publicizes its rationale for adoption as based on the International Vocabulary of Metrology, but a close reading of the sources makes this justification tenuous. The source of the FASEB adoption is unclear, but there’s a chance that Casadevall and Fang (2010) had an influence there. They, in turn, based their definitions on the emphatic but essentially flawed work of Drummond (2009).

Table 2: Grouping of terminologies, as in Table 1, but by discipline.

A	B1	B2
political science	signal processing	microbiology, immunology (FASEB)
economics	scientific computing econometry epidemiology clinical studies internal medicine physiology (neuro)	computer science (ACM)
	computational biology biomedical research statistics	

Book: “The Practice of Reproducible Research”

<https://www.practicereproducibleresearch.org/core-chapters/1-intro.html>



From one case study by José Manuel Magallanes

<https://www.practicereproducibleresearch.org/case-studies/jmMagallanes.html>

“Identifying and Overcoming Threats to Reproducibility, Replicability, Robustness, and Generalizability in Microbiome Research”, Schloss 2018

Exercise 2. Imagine that a graduate student is really excited about an analysis that you performed in your most recent paper and would like to replicate the analysis with their own data. But first, they want to make sure that they reproduce your results. What steps are likely to cause the student problems?

If it is not clear to you what problems they might face, find your favorite figure from a paper by a different research group than your own. Can you reproduce the figure? What is standing in your way?

Exercise 5. Many of the threats to reproducibility and replicability are a product of scientific culture: methods sections are terse or vague, original data are not available, analyses rely on expensive and proprietary software, analysis scripts are available “upon request from the authors,” and papers are published behind paywalls. Some might give into despair thinking that one person or research group can have only a minor impact.

Have a discussion within your group about why things are this way, whether your group’s practices should change, and what would be the easiest and most impactful thing to change.

TABLE 2 An aspirational rubric for evaluating the practices that host-associated microbiome researchers might use to increase the reproducibility and replicability of their work^a

Practice	Good	Better	Best
Handling of confounding variables	Prior to generating data, did we identify a list of possible confounding variables—biological and technical—that may obscure the interpretation of our results?	Do we indicate the level of randomization and experimental blocking that we performed to minimize the effect of the confounding variables?	Does the interpretation of our results limit itself to only those variables that are not obviously confounded?
Sex/gender as confounding variables	Do we indicate the sex/gender of research animals/participants?	Do we provide a justification for the lack of even representation?	Is there equitable representation of sexes/genders? Do we account for them as a variable?
Experimental design considerations	Do we have an active collaboration with a statistician who helps with experimental design and analysis?	Do we indicate the number of hypothesis tests that we performed and have we corrected any <i>P</i> values for multiple comparisons?	For our primary research questions, have we run a power analysis to determine the necessary sample size?
Data analysis plan	Before starting an analysis, have we articulated a set of primary and secondary research questions?	Has someone else reviewed our data analysis plan prior to analyzing the data?	Have we registered our data analysis plan with a third party before starting the project?
Provenance of reagents	Is there a table of reagents such as cell lines, strains, and primer sequences that were used?	Where possible, have we obtained reagents from certified entities like the American Type Culture Collection (ATCC)?	Is there a statement indicating how we know the provenance and purity of each cell line and strain?
Controlling for initial microbiota	Are mice obtained from a breeding facility that allows us to track their pedigree?	Where possible, are mice from different treatment groups cohoused to control for differences in initial microbiota?	Are comparisons between mice with different genotypes made using mice that are the result of matings between animals that are heterozygous for that genotype?
Clarity of software descriptions	Are all methods, databases, and software tools cited? Do we follow the relevant licensing requirements of each tool?	Do we indicate dates and version numbers of websites that were used to obtain data, code, and other third-party resources?	Are detailed methods registered on a website like protocols.io or GitHub?
DNA contamination	Did we quantify the background DNA concentration in our reagents? Did we sequence an extraction control?	Are we taking steps to minimize reagent contamination?	What methods do we take to confirm a result that a sequencing result may be clouded by contaminating DNA?
Availability of data products	Are all of the raw data publicly available?	Are intermediate and final data files publicly available?	Are tools like Amazon Machine Images (AMIs) used to make a snapshot of our working directory?
Availability of metadata	Are all of the metadata necessary to repeat any analyses that we performed publicly available?	Have we adhered to standards in releasing the minimum amount of metadata about our samples?	Did we go beyond the minimum to incorporate other pieces of metadata that will inform future studies?
Data analysis organization	Are all data, code, results, and documentation housed within a monophyletic folder structure on our computer?	Is this project contained within a single directory on our computer, and does it separate our raw and processed data, code, documentation, and results?	Is this folder structure under version control? Is the project's repository publicly available? Are there assurances that this repository will remain accessible?

TABLE 2 An aspirational rubric for evaluating the practices that host-associated microbiome researchers might use to increase the reproducibility and replicability of their work^a

TABLE 2 (Continued)

Practice	Good	Better	Best
Use of random number generator	Do we know whether any of the steps in our data analysis workflow depend on the use of a random number generator?	For analyses that utilize a random number generator, have we noted the underlying random seed?	Have we repeated our analysis with multiple seeds to show that the results are insensitive to the choice of the seed?
Defensive data analysis	Is our data analysis pipeline flexible enough to add new data?	Does our code include tests to confirm that it does what we think it does?	Did we make use of automated tests and continuous integration tools to ensure internal reproducibility?
Ensuring short- and long-term reproducibility	Did we release the underlying code and new data at the time of submitting a paper with their DOIs and accession numbers?	Did we include a reproducibility statement or declaration at the end of the manuscript? Are ORCID identifiers provided for all authors?	What mechanisms are in place to ensure that our analysis remains accessible and reproducible in 5 years?
Open science to foster reproducibility	Have we released any embargoes on our code repository and raw data prior to submitting the manuscript?	Did we post a preprint version of our manuscript prior to submission?	Have we published under a Creative Commons license? Is a permissive reuse license posted with our code?
Transparency of data analysis	Is it clear where one would go to find the data and processing steps behind any of our figures?	Are electronic notebooks publicly accessible, and do they accompany the manuscript?	Were literate programming tools used to generate summary statistics, tables, and figures?

^aAlthough many of the questions can be thought of as having a yes-or-no answer, a better approach would be to see the questions as being open ended with the real question being "What can we do to improve the status of our project on this point?" With this in mind, a researcher is unlikely to have a project that satisfies the "Best" column for each line of the table. Researchers are encouraged to adapt and modify the categories to suit their own needs.

Background: R Notebook vs R Markdown documents

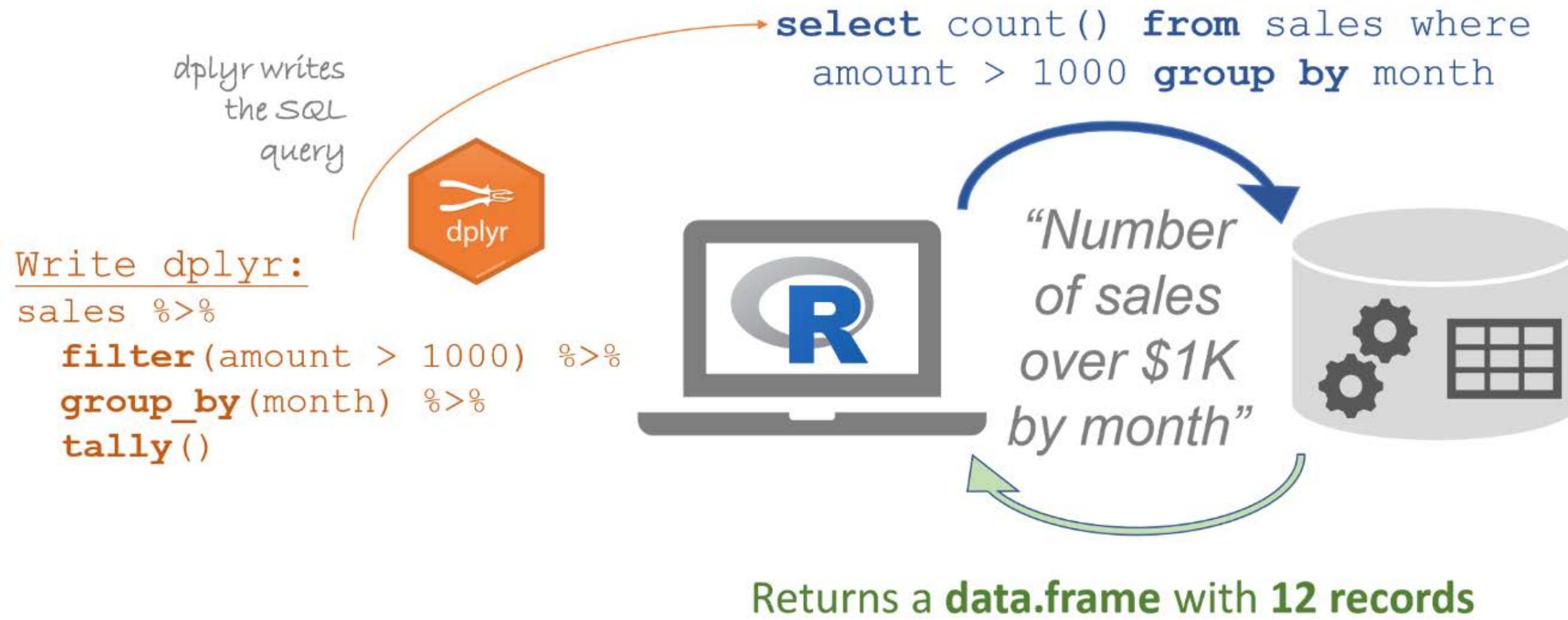
A notebook, `html_notebook`, is a variation on a `html_document`. The rendered outputs are very similar, but the purpose is different. A `html_document` is focussed on communicating with decision makers, while a notebook is focussed on collaborating with other data scientists. These different purposes lead to using the HTML output in different ways. Both HTML outputs will contain the fully rendered output, but the notebook also contains the full source code. That means you can use the `.nb.html` generated by the notebook in two ways:

1. You can view it in a web browser, and see the rendered output. Unlike `html_document`, this rendering always includes an embedded copy of the source code that generated it.
2. You can edit it in RStudio. When you open an `.nb.html` file, RStudio will automatically recreate `.Rmd` file that generated it. In the future, you will also be able include supporting files (e.g. `.csv` data files), which will be automatically extracted when needed.

Emailing `.nb.html` files is a simple way to share analyses with your colleagues. But things will get painful as soon as they want to make changes. If this starts to happen, it's a good time to learn Git and GitHub. Learning Git and GitHub is definitely painful at first, but the collaboration payoff is huge. As mentioned earlier, Git and GitHub are outside the scope of the book, but there's one tip that's useful if you're already using them: use both `html_notebook` and `github_document` outputs:

<http://r4ds.had.co.nz/r-markdown-formats.html>

Ideally, analyze in place, using dplyr



Where do we go from here?

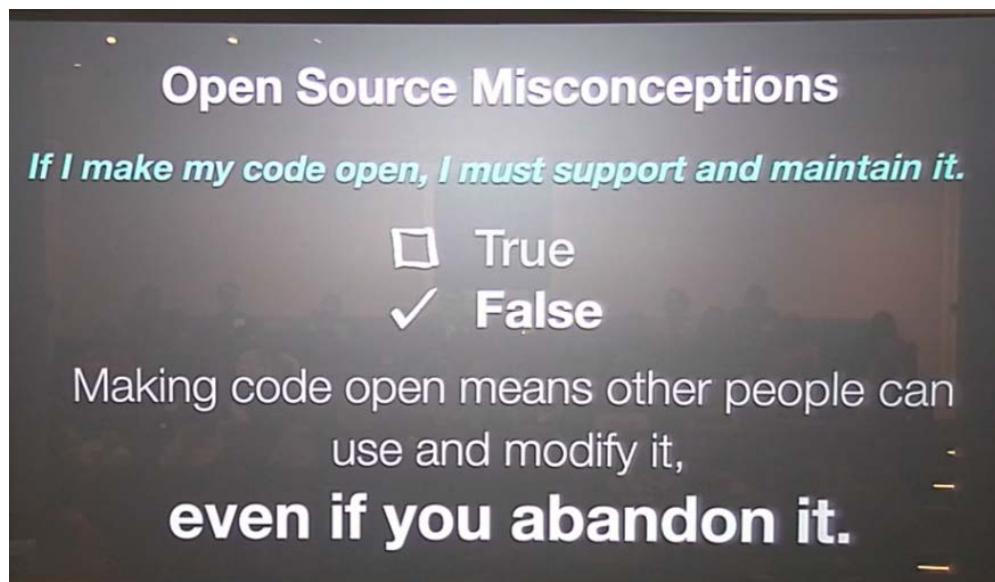
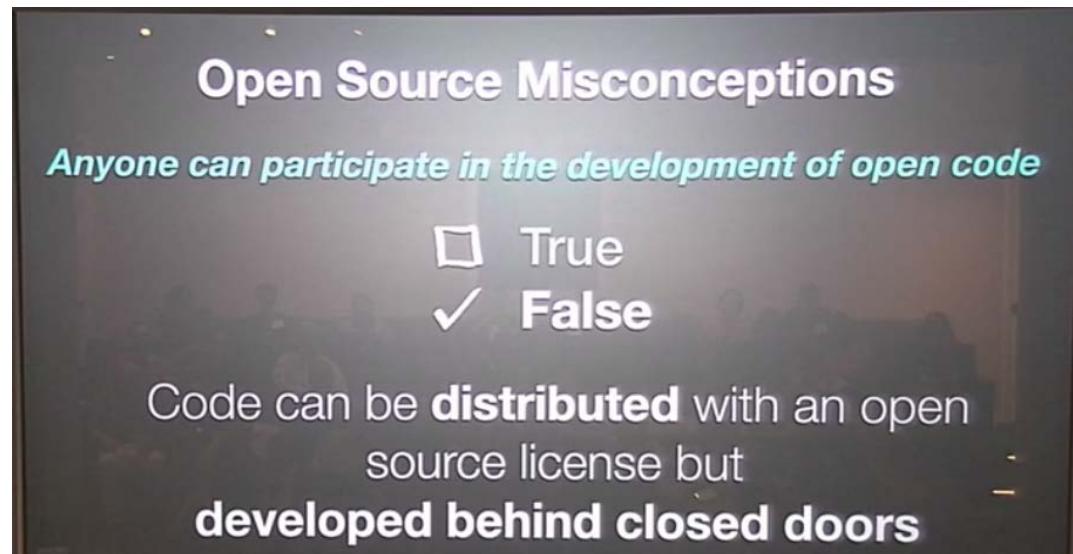
Prof. Kelle Cruz's perspective

The image shows a presentation slide with a dark background. At the top, the word "Conclusions" is written in a large, white, sans-serif font. Below it is a bulleted list of six items, also in white text. To the right of the list are three logos: the GitHub logo (a black cat icon), the Python logo (a stylized "P" in blue and yellow), and the NumFOCUS logo (the word "NUMFOCUS" in red and blue with "OPEN CODE = BETTER SCIENCE" below it). The slide has a thin white border.

- Science is software
- GitHub facilitates open development and has changed the way we interact with software.
- Python is rapidly becoming the common language of science.
- NumFOCUS needs to be recognized as a vital part of science infrastructure.
- Software development needs to be taught and integrated into the formal STEM curriculum.

Screenshot from
Prof. Kelle Cruz's (Hunter
College/CUNY) Simons
Foundations Lecture
(link to [slides](#))

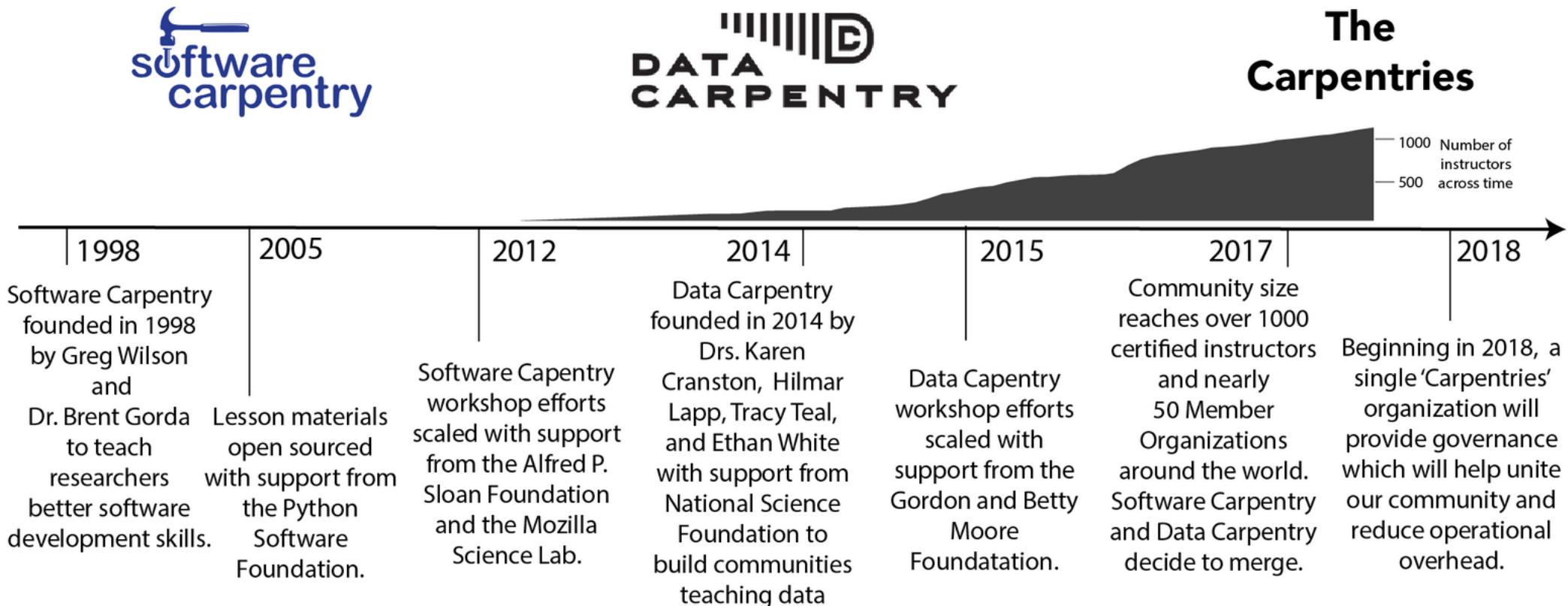
Open source misconceptions



Screenshots from
Prof. Kelle Cruz's (Hunter College/CUNY)
Simons Foundations Lecture
(link to [slides](#))

History of Software Carpentry (now: The Carpentries)

Check out upcoming workshops: <https://carpentries.org>



BinderHub Architecture (high-level details)

