# Data Management:
# The First Step in Reproducible Research

---

*Harvard Chan Bioinformatics Core | Tools for Reproducible Research*
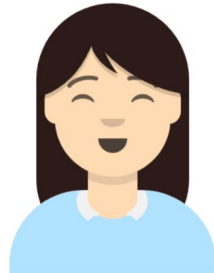
August 6, 2024

Julie Goldman | julie_goldman@harvard.edu

# Learning Objectives

- Understand the impact of creating reproducible research

- Examine challenges of creating reproducible research data

- Discuss foundational data management practices

- Review available tools that facilitate reproducible research data

# Defining Reproducibility

**Source**: ITCR Training Network (ITN). 2024. "Intro to Reproducibility in Cancer Informatics."
https://jhudatascience.org/Reproducibility_in_Cancer_Informatics
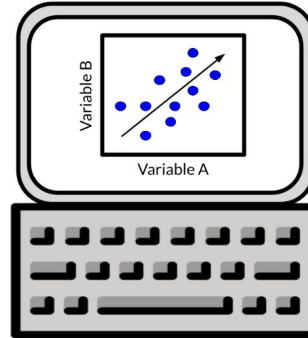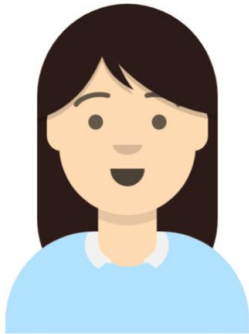
# Repeatability



Image created by Candace Savonen using Avataars.

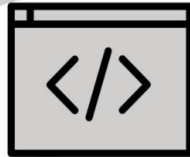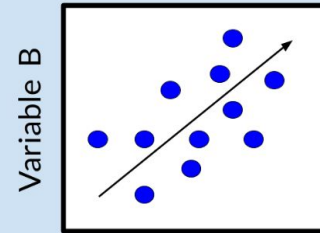**Source**: ITCR Training Network (ITN). 2024. "Intro to Reproducibility in Cancer Informatics."
https://jhudatascience.org/Reproducibility_in_Cancer_Informatics
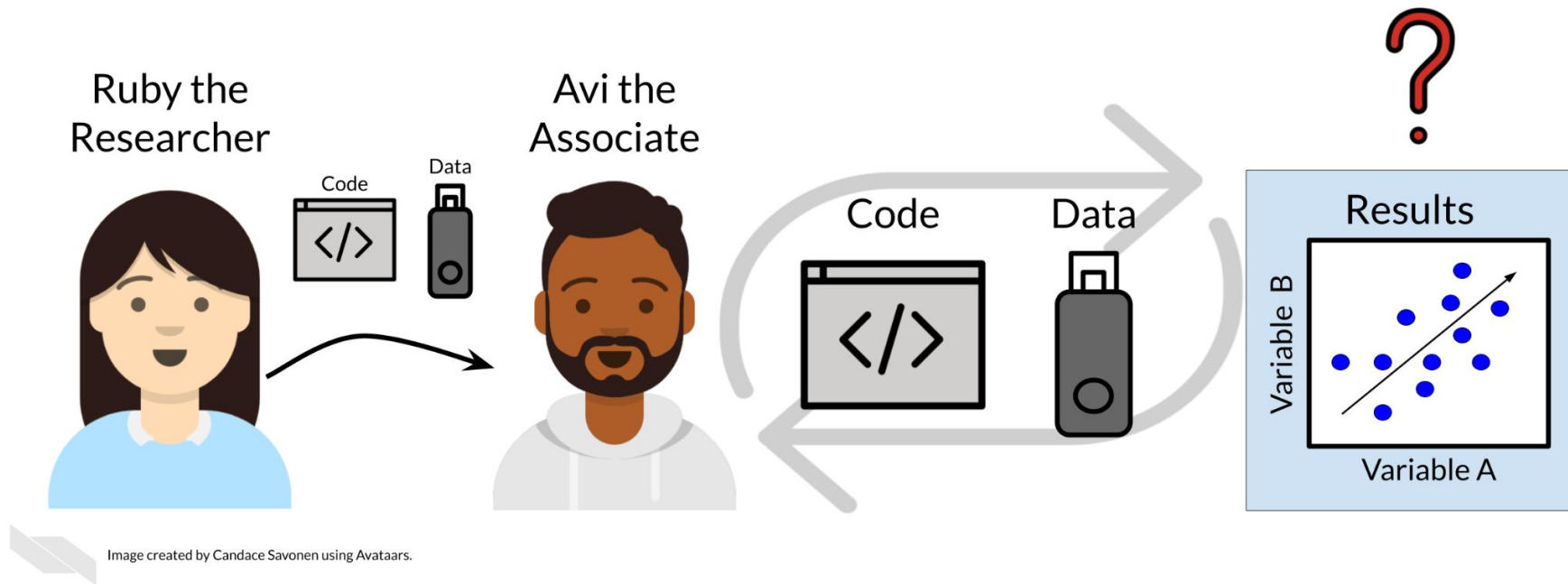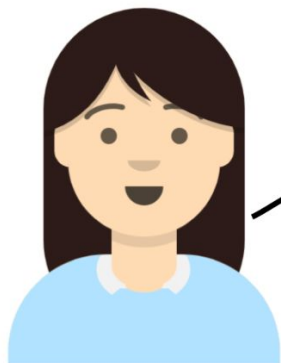
# Reproducibility



Image created by Candace Savonen using Avataars.

**Source**: ITCR Training Network (ITN). 2024. "Intro to Reproducibility in Cancer Informatics."
https://jhudatascience.org/Reproducibility_in_Cancer_Informatics

# Replicability



Image created by Candace Savonen using Avataars.

# Research process as a hierarchy

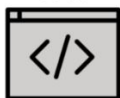# So, what's the issue?

# Reproducibility in daily life

**Source**: ITCR Training Network (ITN). 2024. "Intro to Reproducibility in Cancer Informatics."
https://jhudatascience.org/Reproducibility_in_Cancer_Informatics

# Reproducibility in daily life



Image created by Candace Savonen using Avataars.

# Reproducibility in daily life



**Now Ruby**

Ruby's code

**Future Ruby**

Ruby's code

ERROR

# Reproducibility is worth the effort!



Ruby's code - made reproducibly

# So, why not put in the effort?

I don't have enough time

Technical obsolescence

I don't have the skills

There's not enough incentive

I don't know where to start

Another reason

You can't have any sort of reproducibility without good data and project management.

# Research Data Management

Is the active and ongoing management of data through its lifecycle of interest and usefulness.

Ensures and facilitates the timely collection of complete and accurate protocol-required information.

Includes decisions that are agreed to at the beginning of a study and carried out to completion.



Research Data Lifecycle by LMA RDMWG

# Data Through the Research Lifecycle

**Final Data**



Are the published final data available for validation, reproduction or reuse?



Plan & Design

Publish & Reuse

Share & Disseminate

Store & Manage

Collect & Create

Evaluate & Archive

Analyze & Collaborate

What about experimental methods and measurement parameters?



**Raw Data**



**Intermediate Data**

Research Data Lifecycle by LMA RDMWG

# Data Management Practices for Reproducibility

## Organization

- Directory structure
- File naming
- Version control

## Documentation

- README File
- Data Dictionary
- Metadata

## Automation

- Scripts for workflows
- Computing environment
- Dependencies

## Dissemination
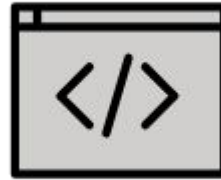
- Share in repository
- Get DOI for citation
- License and terms of use language

# Organization: What to avoid

# Organization: Better practice



I read my README to get me back up to speed with this project. Now I know that I can run a single command to call *run_analysis.sh* to re-run my analysis.

Ruby the Researcher

- raw-data
- README.md
- cleaned-data
- figures
- source-code
- run_analysis.sh
- 01-clean-data.R
- 02-create-plot.R

Image created by Candace Savonen using Avataars.

# Organization: Tips and tricks

- Make file names informative – avoid using spaces, quotes, or unusual characters

- Keep like-files together in their own directory – keep raw data separate from processed data or other results!

- Number scripts in the order that they are run

- Put source scripts and functions in their own directory

- Put output in its own directories like results and plots

- Have a central document (like a README) that describes the basic information about the project and analysis (see: documentation)

- Make a central script that re-runs everything (see: automation)

# Documentation: Good practice



Avi the Associate

I had no idea where to start with this analysis that Ruby sent me to review, but then I saw she included a **README** and that saved me so much time and effort in getting started!
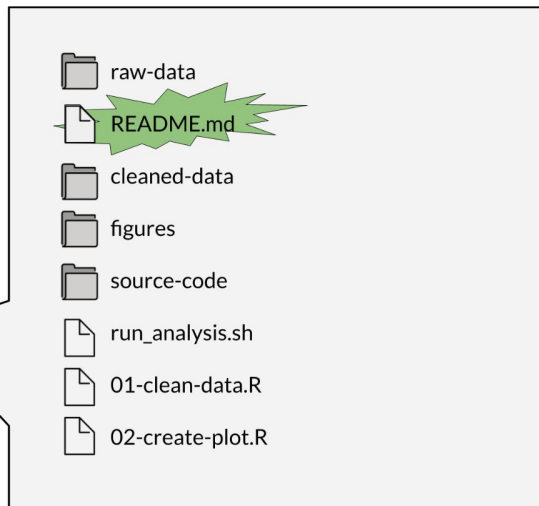
- raw-data
- README.md
- cleaned-data
- figures
- source-code
- run_analysis.sh
- 01-clean-data.R
- 02-create-plot.R

Image created by Candace Savonen using Avataars.

# Documentation: Good practice
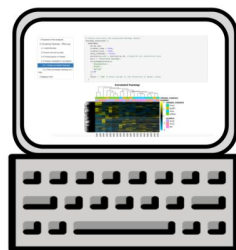


Source: ITCR Training Network (ITN). 2024. "Intro to Reproducibility in Cancer Informatics."
https://jhudatascience.org/Reproducibility_in_Cancer_Informatics

# Documentation: Useful tools



README File Example Template:
http://data.research.cornell.edu/content/readme

R Markdown: The Definitive Guide:
https://bookdown.org/yihui/rmarkdown/

# Automation: What to avoid



Ruby's local computing environment    Avi's local computing environment

Created by Candace Savonen

# Automation: Good practice



**Ruby's computing environment**

Ruby the Researcher

Avi the Associate

Slide from the CCDL adapted by Candace Savonen

**Source**: ITCR Training Network (ITN). 2024. "Advanced Reproducibility in Cancer Informatics."
https://jhudatascience.org/Adv_Reproducibility_in_Cancer_Informatics

# Automation: Tips and tricks

Create a script that can execute all of the various subcomponents of the entire workflow.

This simple example has three steps that can be performed automatically:

1. **`clean_data.R`** to generate the cleaned data table

2. **`analysis.R`** to perform the statistical test

3. **`runall.sh`** saved in the src directory to run the entire workflow process

```
|-- tomato_project
|    |-- data_raw
|    |    |-- raw_yield_data.csv
|    |    |-- README.txt
|    |-- src
|    |    |-- analysis.R
|    |    |-- clean_data.R
|    |    |-- runall.sh
```

# Dissemination: Good practice

# Dissemination: Better practice

"Just email me and
I'll send it to you"

1. See "supplemental materials"

GitHub

www.mywebsite.com/my-data/projectHelloWorld

Dropbox
Box.com
drive.google.com

Data repository

# Dissemination: Useful tools

| Disciplinary | General | Software | Methods |
|:---:|:---:|:---:|:---:|

# Putting it all together!



**Source**: Peng, Roger D. 2011. "Reproducible Research in Computational Science." *Science* 334 (6060): 1226-1227. https://doi.org/10.1126/science.1213847

It takes some effort to organize your research to be reproducible...the principal beneficiary is generally the author themself.
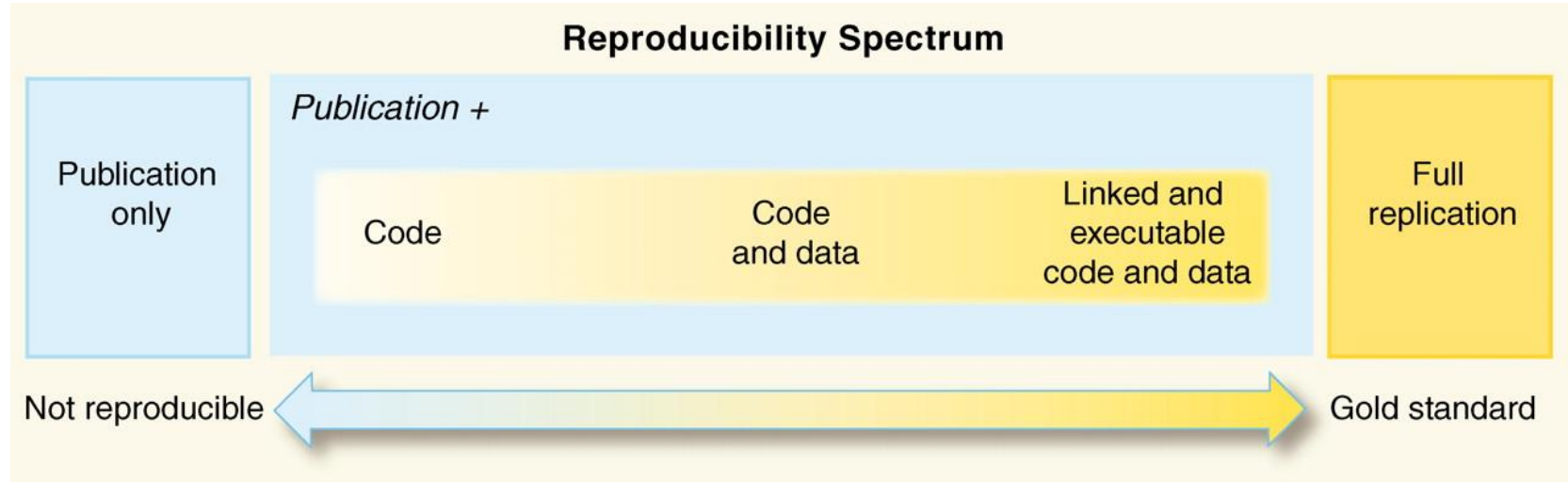
– Jon Claerbout
*Making Scientific Contributions Reproducible*

# Why Reproducibility? Think Selfishly!

# Open Researcher and Contributor ID

- ORCID: Provides a persistent digital identifier that distinguishes you from every other researcher and supports automated linkages between you and your professional activities ensuring that your work is recognized

- URI with a 16-digit number that is compatible with the ISO Standard (ISO 27729) or International Standard Name Identifier (ISNI), e.g. https://orcid.org/0000-0001-2345-6789



https://orcid.org

# Closing Remarks



**Image Source**: Taron Egerton & Richard Madden on "Carpool Karaoke" Season 2, Episode 18, March 21, 2019

# References & Resources

- Borghi, John, et al. 2018. "Support your data: A research data management guide for researchers." Research Ideas and Outcomes 4: e26439. https://doi.org/10.3897/rio.4.e26439

- Briney, Kristin A., Heather L. Coates, and Abigail Goben. 2020. "Foundational practices of research data management." Research Ideas and Outcomes 6: e56508. https://doi.org/10.3897/rio.6.e56508

- Kathawalla, Ummul-Kiram, Priya Silverstein, and Moin Syed. 2021. "Easing into open science: A guide for graduate students and their advisors." Collabra: Psychology 7 (1): 18684. https://doi.org/10.1525/collabra.18684

- McKiernan, Erin C., et al. 2016. "How open science helps researchers succeed." elife 5: e16800. https://doi.org/10.7554/eLife.16800

- Wilkinson, Mark D., et al. 2016. "The FAIR Guiding Principles for scientific data management and stewardship." Scientific Data 3: 160018. https://doi.org/10.1038/sdata.2016.18

- Wilson, Greg, et al. 2017. "Good enough practices in scientific computing." PLoS computational biology 13 (6): e1005510. https://doi.org/10.1371/journal.pcbi.1005510