# You will not hear any sound until the webinar starts.
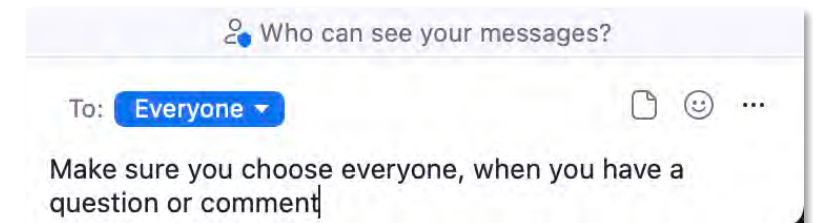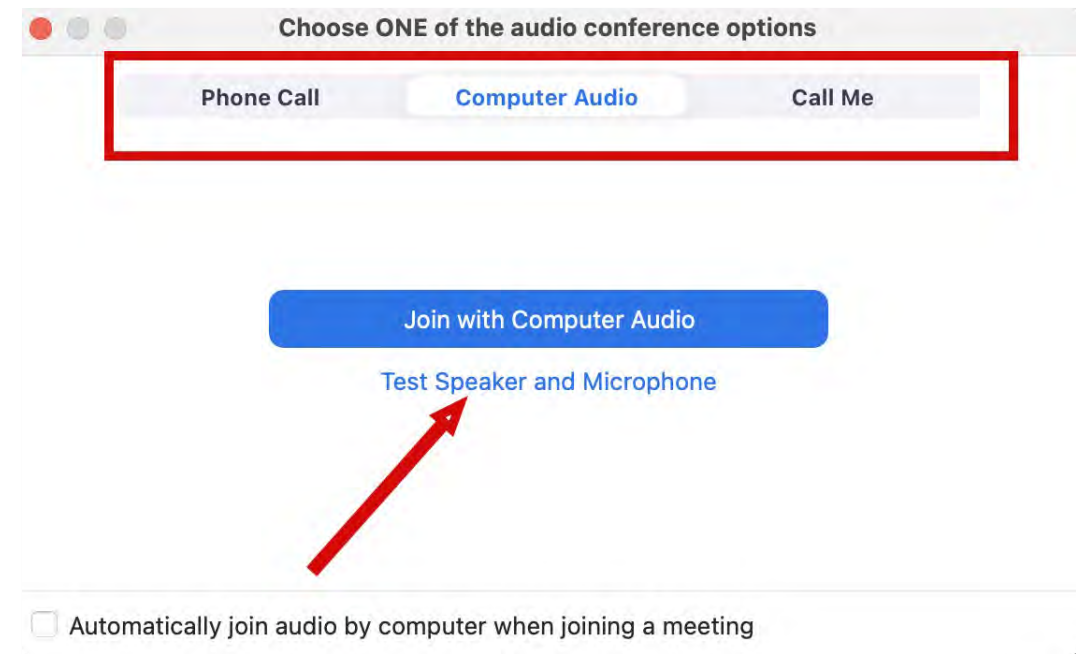
## Connect Audio
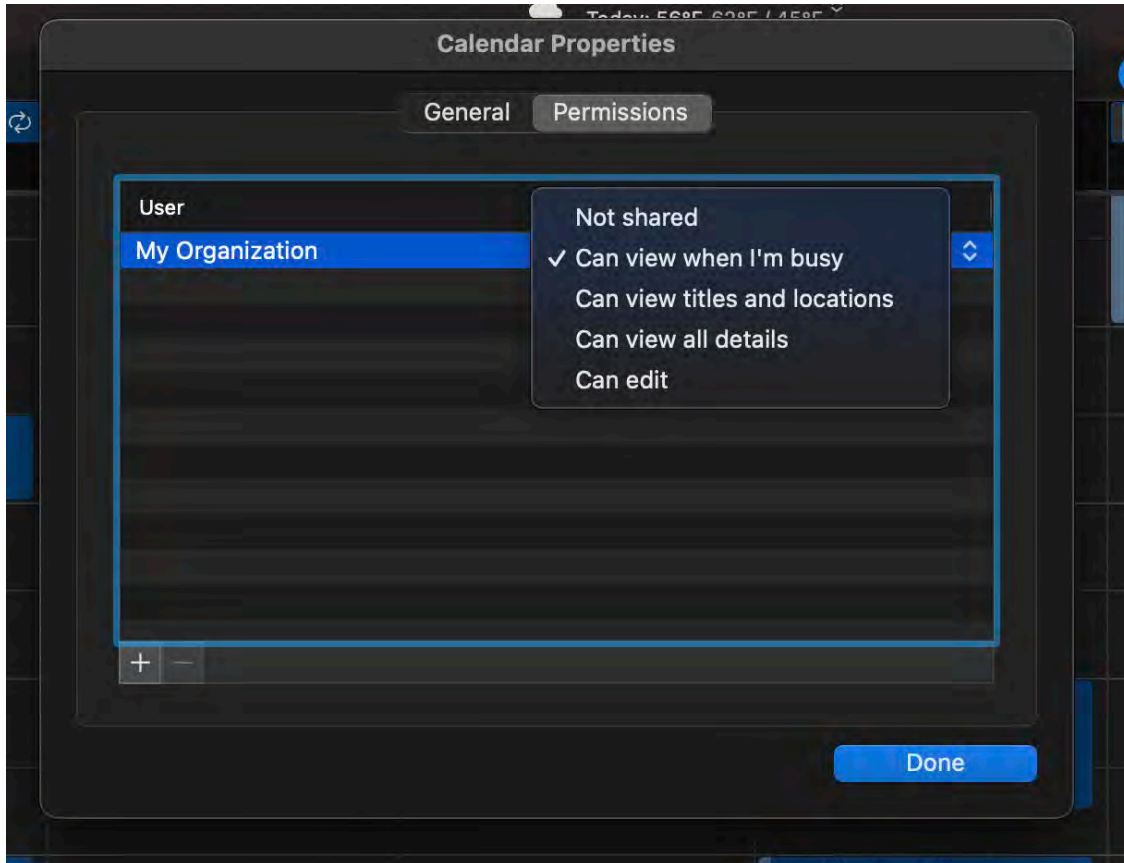
1. When you join Zoom, the *Join Audio* preferences box pops-up (Phone Call, Computer Audio, or Call Me)
2. Choose an option that works best for you
3. Join using that option
4. Use Test Speakers and Microphone option to optimize your webinar experience

## Chat

Please send your chat to *Everyone* to make sure the monitor sees your question

Choose ONE of the audio conference options

Phone Call · Computer Audio · Call Me

Join with Computer Audio

Test Speaker and Microphone

☐ Automatically join audio by computer when joining a meeting

Who can see your messages?

To: Everyone ▾

Make sure you choose everyone, when you have a question or comment
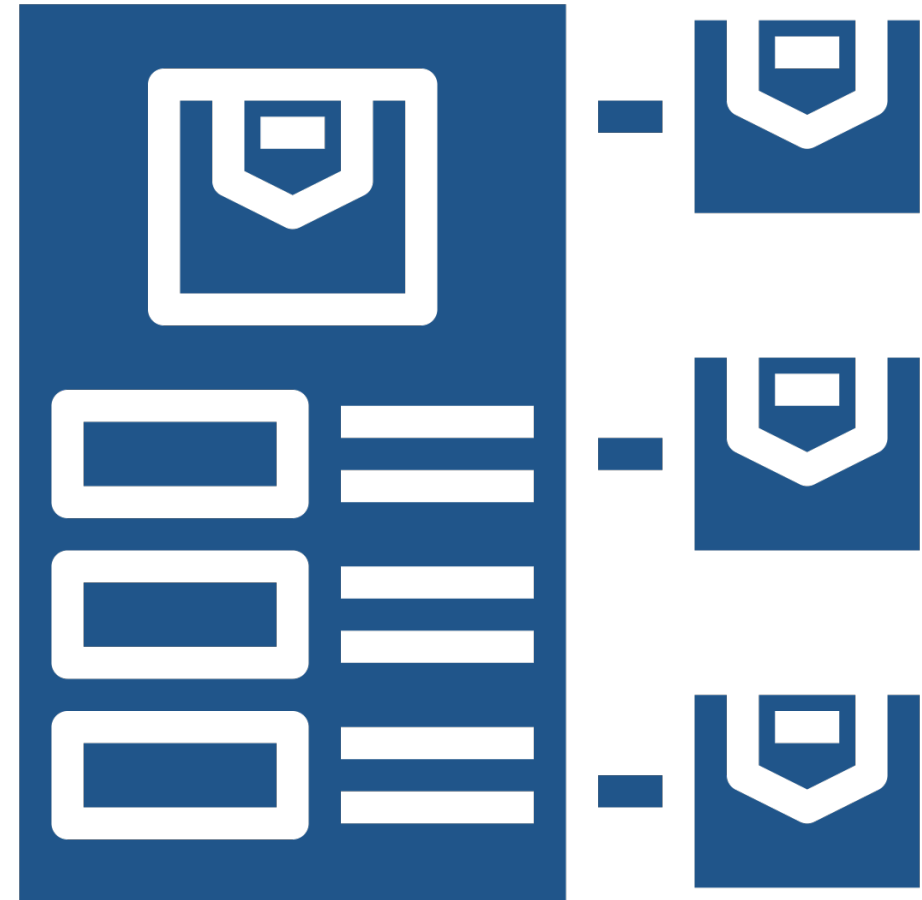
Please rename yourself, so we can:

- Send you the student version of the PowerPoint
- Send your training certificate
- Add you to our list-serve

# Project Management and Reproducibility In RStudio

**Doug Joubert**

**2023-02-28**

- Focuses on data and project management using R and Rstudio

- Some familiarity or experience in R and RStudio is recommended but not required

- Define scientific reproducibility
- Discuss best practices for organizing data in an RStudio project
- Discuss the importance of using a data dictionary and read me files
- Ensure that their data is machine readable

# Configuration for Exercises

- R is a programming language that is especially powerful for data exploration, visualization

- RStudio is an integrated development environment (IDE) that makes using R easier

- R and RStudio are two separate pieces of software

- **Must install R before you install RStudio**

1. Download R from the [CRAN website](#)

2. Run the .exe file that was just downloaded

1. Go to the RStudio [download page](#)

2. Under Installers select RStudio x.yy.zzz - Windows Vista/7/8/10 (where x, y, and z represent version numbers)

3. Double click the file to install it

1. Download R from the CRAN website
2. Select the .pkg file for the latest R version
3. Double click on the downloaded file to install R
4. It is also a good idea to install XQuartz (needed by some packages)

1. Go to the RStudio [download page](#)

2. Under Installers select RStudio x.yy.zzz - Mac OS X 10.6+ (64-bit) (where x, y, and z represent version numbers)

3. Double click the file to install RStudio

- Open Power Shell on Windows or Terminal in RStudio
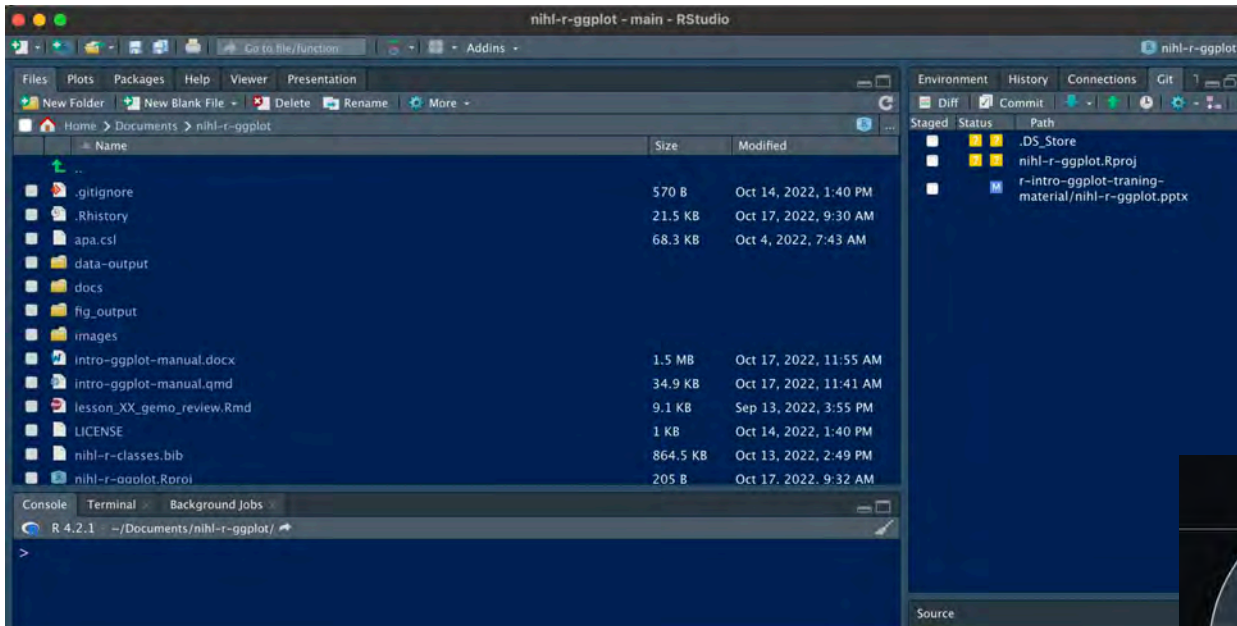- type: `where git`

- Open Power Shell on Windows or Terminal in RStudio
- type: `which git`

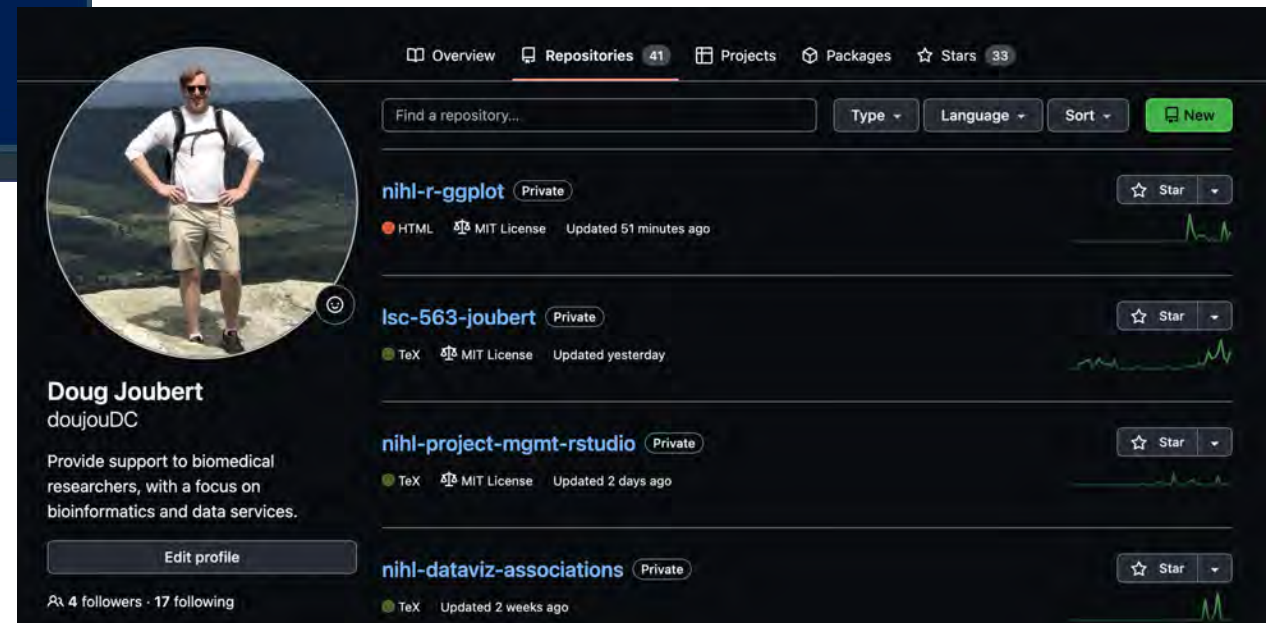- [Git installation instructions](#) to install Git, if not on your computer

# Project Configuration

- Follows RStudio project framework
- Code in .qmd file
- Wilson, G., Bryan, J., Cranston, K., Kitzes, J., Nederbragt, L., & Teal, T. K. (2017). Good enough practices in scientific computing. PLOS Computational Biology, 13(6), e1005510-e1005510.

- All class repos sync'd to GitHub
- Includes:
  - Class data
  - Folder structure
  - Training material
  - Exercises

- Introduction to R and RStudio
- Introduction to R Data Types
- Data Wrangling in R
- Introduction to Project Management in RStudio (A)
- Reproducibility in RStudio: Basic Markdown
- Introduction to Data Visualization in R: ggplot (A)

- Reproducibility in RStudio: Advanced Markdown
- Working with Git in RStudio
- Introduction to Data Visualization in R: Customization in ggplot

According to the U.S. National Science Foundation (NSF) subcommittee on replicability in science (2015):

Science should routinely evaluate the reproducibility of findings that enjoy a prominent role in the published literature. To make reproduction possible, efficient, and informative, **researchers should sufficiently document the details of the procedures used to collect data**, to **convert observations into analyzable data**, and to **perform data analysis**.

Bollen, K., et al (2015)

7%
Don't know

3%
No, there is no crisis

IS THERE A
REPRODUCIBILITY
CRISIS?

A *Nature* survey lifts the lid on
how researchers view the 'crisis'
rocking science and what they
think will help.

BY MONYA BAKER

52%
Yes, a significant
crisis

38%
Yes, a slight
crisis

1,576
RESEARCHERS SURVEYED

Baker, M. (2016).

- 1,576 researchers took online questionnaire
- More than 70% of researchers have tried and failed to reproduce another scientist's experiments
- More than half have failed to reproduce their own experiments
- Specific factors in handout

# Reproducibility Crisis?

7%
Don't know

3%
No, there is no crisis

IS THERE A REPRODUCIBILITY CRISIS?

A *Nature* survey lifts the lid on how researchers view the 'crisis' rocking science and what they think will help.

BY MONYA BAKER

52%
Yes, a significant crisis

38%
Yes, a slight crisis

1,576
RESEARCHERS SURVEYED

Baker, M. (2016).

- Sometimes-contradictory attitudes towards reproducibility:
  - > 52% agree that there is a significant crisis
  - < 31% think that failure to reproduce published results means results are wrong
  - Most still trust the published literature

# Problems in Reproducibility

WHAT FACTORS CONTRIBUTE TO
IRREPRODUCIBLE RESEARCH?

Many top-rated factors relate to intense competition and time pressure.

● Always/often contribute   ● Sometimes contribute

Selective reporting

Pressure to publish
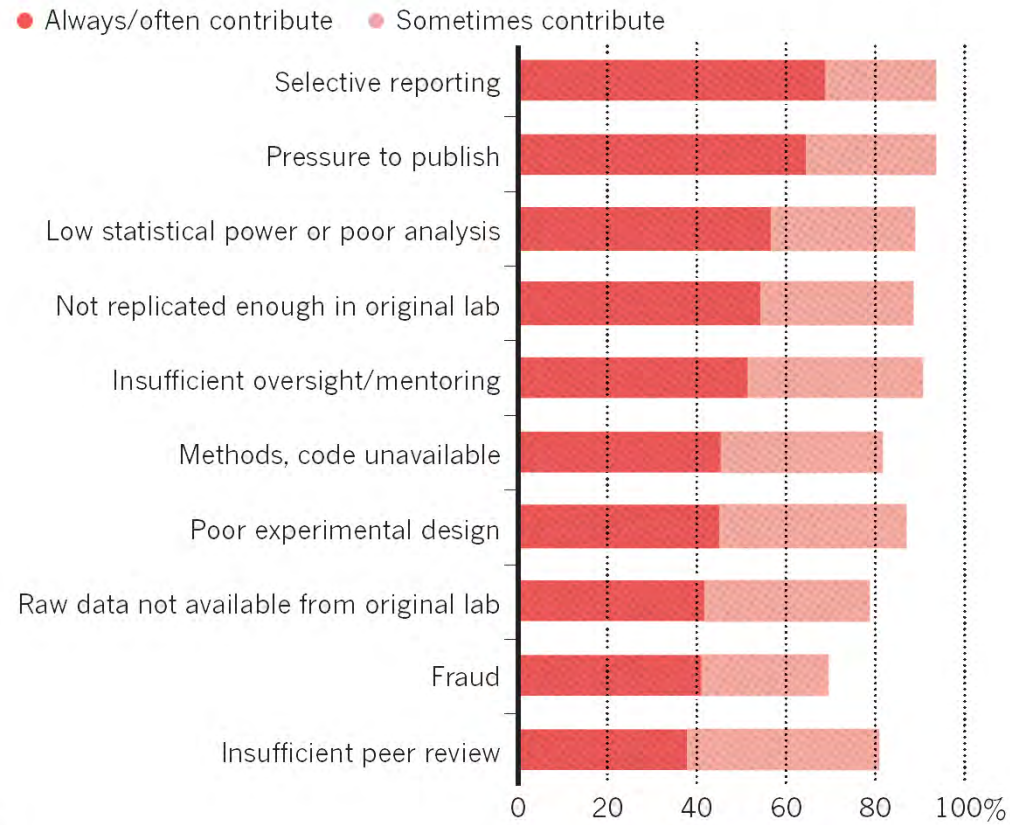
Low statistical power or poor analysis

Not replicated enough in original lab

Insufficient oversight/mentoring

Methods, code unavailable

Poor experimental design

Raw data not available from original lab

Fraud

Insufficient peer review

0   20   40   60   80   100%

Baker, M. (2016).

- > 60% said that two factors were problems:
  - Pressure to publish
  - Selective reporting
- > 50% pointed to:
  - Insufficient replication in the lab
  - Poor oversight
  - Low statistical power
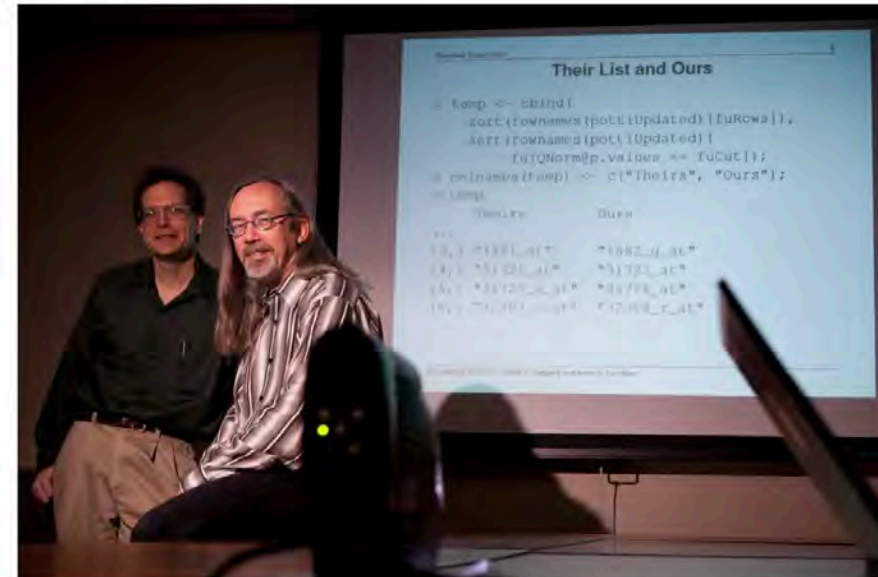
# Embracing Reproducibility Practices

- 5 selfish reasons to use reproducibility practices:
  - Helps to avoid data loss and disaster
  - Makes it easier to write papers
  - Helps reviewers see it your way
  - Enables continuity of your work
  - Helps to build your reputation

Markowetz, F. (2015).



### How Bright Promise in Cancer Testing Fell Apart

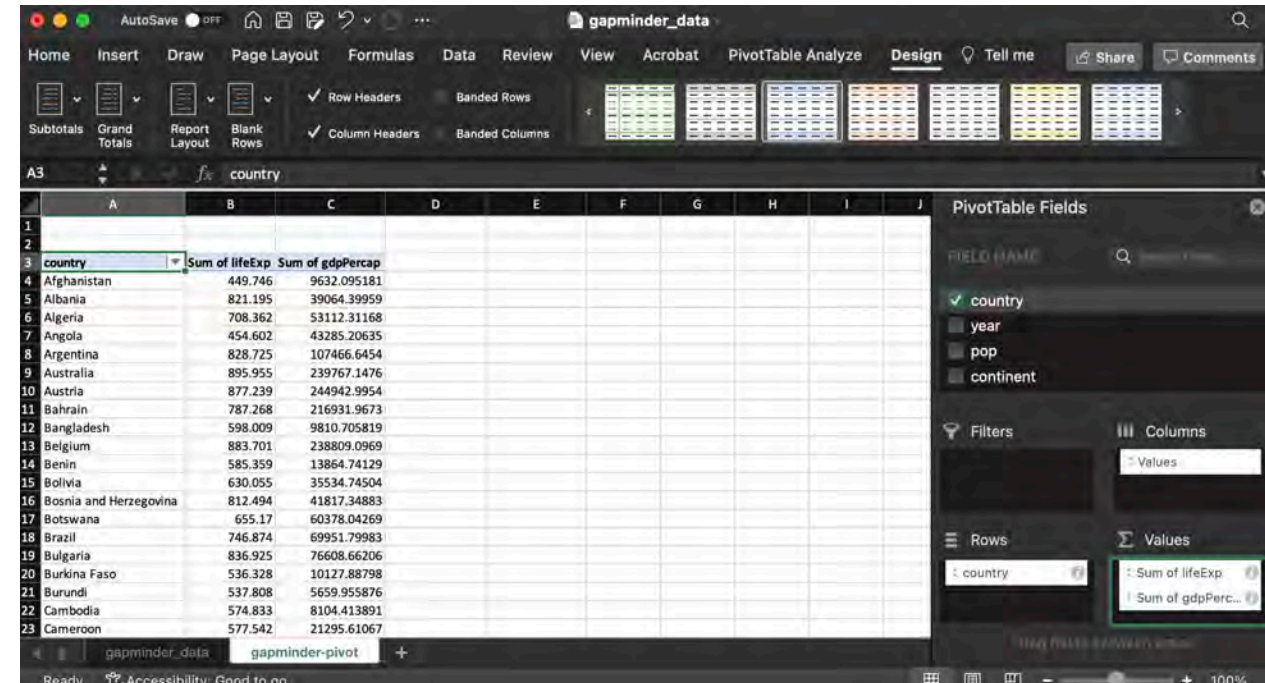Give this article   75

Their List and Ours

Keith Baggerly, left, and Kevin Coombes, statisticians at M. D. Anderson Cancer Center, found flaws in research on tumors.   Michael Stravato for The New York Times
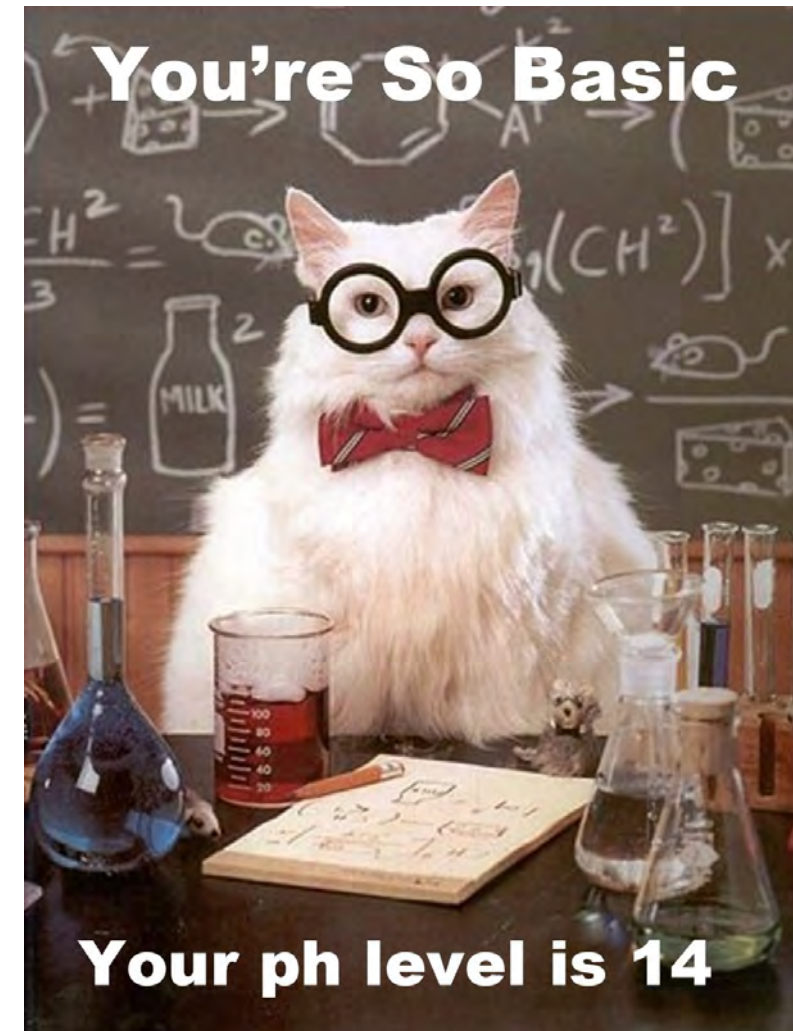
Image: NY Times

- Why should I learn R when Pivot Tables do the same thing?
- For example, Figure 4 is showing a Pivot Table for [gapminder](#) data
  - Mean
  - Median
  - Min value
  - Max value

- With R, using 2 lines of code
  - Group data by country and run one regression model
  - Display statistics like the coefficients or the R2 value
  - I agree with Chemistry Cat, Excel is so basic

- This gets at a more philosophical reason for using R

- Not everyone has access to Office 365

- Tableau or Stata skills are great, but products are very expensive

- If sharing analysis results in Excel, not everyone will be able to open that file

(Heiss, 2022)

# Why Learn R? – Reproducibility

- Only way this would be reproducible is if you write down all the steps for:
  - Every menu you clicked
  - Every cell you clicked on to add a formula or changed the formatting
  - Have you ever seen an Excel spreadsheet with an accompanying set of instructions?

(Heiss, 2022)

- A 2016 study found gene name errors in 20% of the papers that they reviewed
- Web of Science search using the terms "Reproducible statistical analyses" OR "Reproducibility" resulted in over 7,000 published papers, in the last 10 years

(Ziemann, Eren, and EL-Osta)

Septin 2

Membrane-Associated Ring Finger (C3HC4) 1

2310009E13

| | A | B |
|---|---|---|
| 1 | Actual value | What Excel turns it into |
| 2 | SEPT2 | 2-Sep |
| 3 | MARCH1 | 1-Mar |
| 4 | 2310009E13 | 2.31E+19 |

- **Don't Touch the Raw Data**: no analysis on the original data, or, if you do, then explain what you did to the data.

- **Self-documenting and Reproducible Code**: consider writing your reports or papers in markdown. Markdown combines text with code.

- **Use Open Formats**: open formats as much as possible. That means sharing your data in csv or tab-delimited format.

(Heiss, 2022)

# RStudio Project Demo

- Creating a project in RStudio
- Brief discussion about "best practices"

# Best Practices for Managing Projects

- RStudio is an integrated development environment (IDE) for R and Python:
  - Free and open-source
  - Designed to make it easy to write and reuse code
  - Convenient to view and interact with the objects stored in your environment
  - Collaboration and Publishing Tools
  - Documents using R Markdown
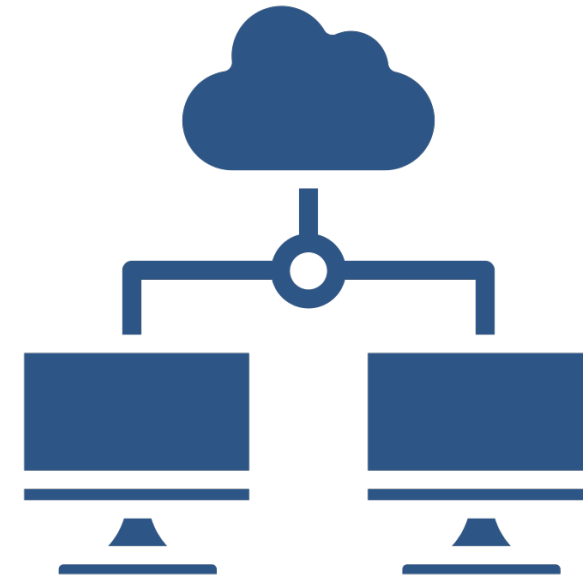- We have a separate [class](#) that focuses on Intro to R

- Using the poll, list some problems that you have encountered when dealing with files and folders

# Files and folders

- Cannot find your files on your computer (or your cloud storage)
- Multiple versions of files with names such as "finaldraft_4.txt"
- Path issues when trying to run code
- Reviewers or colleagues cannot re-run your code/analyses

- Using the poll, list some problems that you have encountered when dealing with storage or sharing

- Files are only saved to your computer

- Collaborators don't share the files needed

- Files are shared via email attachments

- Difficult to know if you have the latest version of documents

- [Good Enough Practices for Scientific Computing](#) gives the following recommendations for project organization:
  - Each project in its own directory, which is named after the project
  - Text documents associated with the project in the doc directory
  - Raw data and metadata in the data directory (raw-data)
  - Files generated during cleanup and analysis in a results directory

Wilson., et al, 2017

- [Good Enough Practices for Scientific Computing](#) gives the following recommendations for project organization:
  - Project's scripts and programs in the src directory
  - Programs brought in from elsewhere or compiled locally in the bin directory
  - Name all files to reflect their content or function

Wilson., et al, 2017

# Practice Good File Organization

- Additional filles to include:
  - [README file](), to communicate important information about your project
  - [LICENSE file](), so that others are free to use, change, and distribute the software
  - [CITATION.cff file](), to let others know how you would like them to cite your work
- Student version of the PowerPoint has more resources to explore

Wilson., et al, 2017

1. Machine readable
2. Human readable
3. Plays well with default ordering

Wilson., et al, 2017

# Naming Files: Why This is Important

- Globbing: using wildcard characters to request or evaluate sets of files with the same partial names or sets of character
- Regular expression friendly
- "Findability"
- Global name changes

- **Machine-readable**
  - No spaces, unsupported punctuation, accented characters, or case-sensitive file names
  - Deliberate use of delimiters (i.e. for splitting file names)
  - Consistently use the same delimiters: data-analyses-fig1.R as an example
- **Human-readable**
  - Name contains brief description of content: i.e. anova-analyses-control.R

- With **chronological ordering**, file name starts with date:
  - 2022-02-26-BRAFWTNEGASSAY-FFPEDNA-CRC-1-41-AO2.csv
  - 2022-02-26-BRAFWTNEGASSAY-FFPEDNA-CRC-1-41-AO3.csv
  - 2022-02-26-BRAFWTNEGASSAY-FFPEDNA-CRC-1-41-AO4.csv

Consider using [ISO 8601 date standard](#)

- With **logical ordering**, the filename starts with a number or keyword/number combo.
  - 01-marshall-data.r *see code directory*
  - 02-pre-dea-filtering.r *see code directory*
  - 03-explore-dea-limma-voom.r
  - 04-exploe-dea-results.r helper
  - 01-load-counts.rmd
  - helper02-load-exp-des.r
  - helper03-extract-and-tidy.r

■ As illustrated on the previous slide, left-pad your numbers to facilitate sorting. If you do not do this, your data sorts like this…which is really sad

 – 01-marshall-data.r *see code directory*

 – 04-exploe-dea-results.r

 – 2-explore-dea-limma-voom.r

 – 3-helper-extract-and-tidy.r

 – helper01-load-counts.rmd

- Configuration for exercises
- Introduction to [Reproducible Publications with RStudio](#)
- R for [Reproducible Scientific Analysis](#)

- Reproducible [Research Data and Project Management in R](#)
- Using [Projects in Rstudio](#)
- email me for a copy: [douglas.joubert@nih.gov](mailto:douglas.joubert@nih.gov)
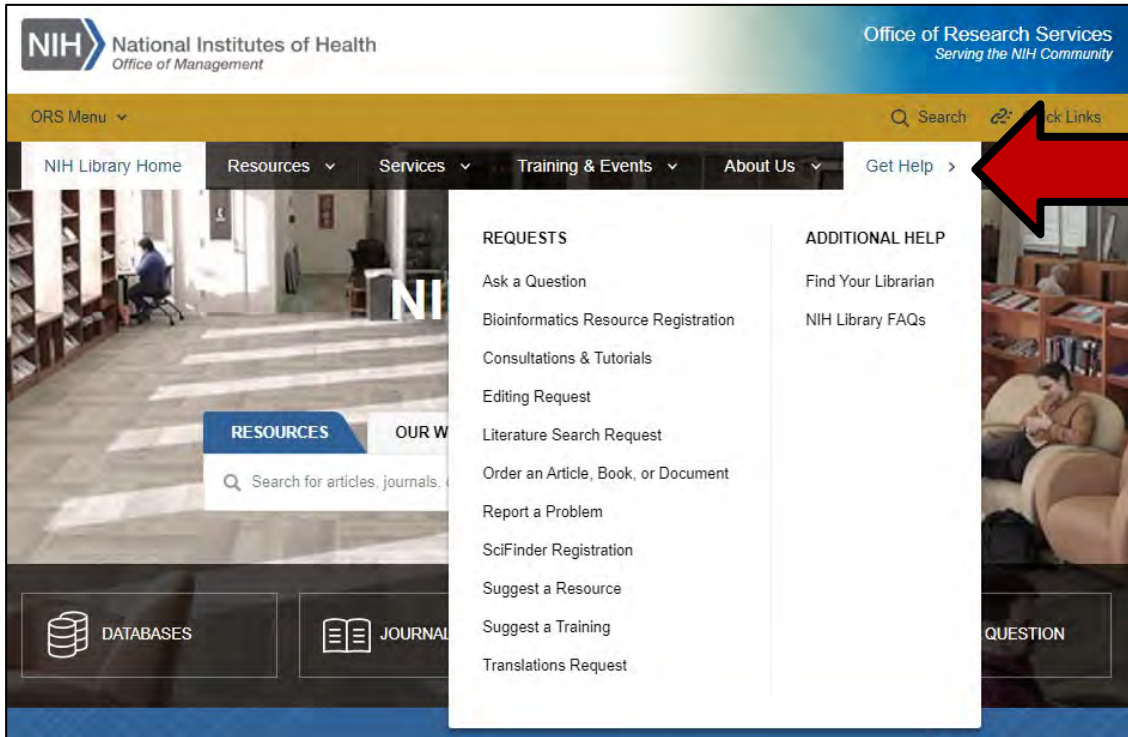
- [Classes](#) on a variety of data-related topics, including:
  - Data management
  - Data visualization
  - Data analysis
  - R and RStudio
- [Computers](#) which offers a suite of tools for data analysis, processing, and visualization

# Contact Us for Ongoing Support

**Doug Joubert**

Bioinformatics Support Program

301-827-3829

douglas.joubert@nih.gov

**NIH Library Help Desk**
(301) 496-1080

- **Ask a Question**: https://www.nihlibrary.nih.gov/get-help/ask-question
- **Request a Tutorial**: https://www.nihlibrary.nih.gov/get-help/consultations-tutorials
- **Sign up for Additional Classes**: https://www.nihlibrary.nih.gov/training/calendar