

Reproducible Research

EES 4891/5891

Probability & Statistics for Geosciences

Jonathan Gilligan

Class #5: Tuesday, January 21 2025

Learning Goals

Learning Goals for Today

- Understand what reproducible research is, and why it's important
- Learn about version control
 - Know how to manage revisions with `git` and `GitHub`
- Understand how to use `Git` in RStudio projects
- Learn about using `RMarkdown` and `Quarto` for literate programming and reproducible research methods.

Getting Started

Go to GitHub Classroom and accept the project at
<https://classroom.github.com/a/Hg3m9DoE>



Reproducible Research

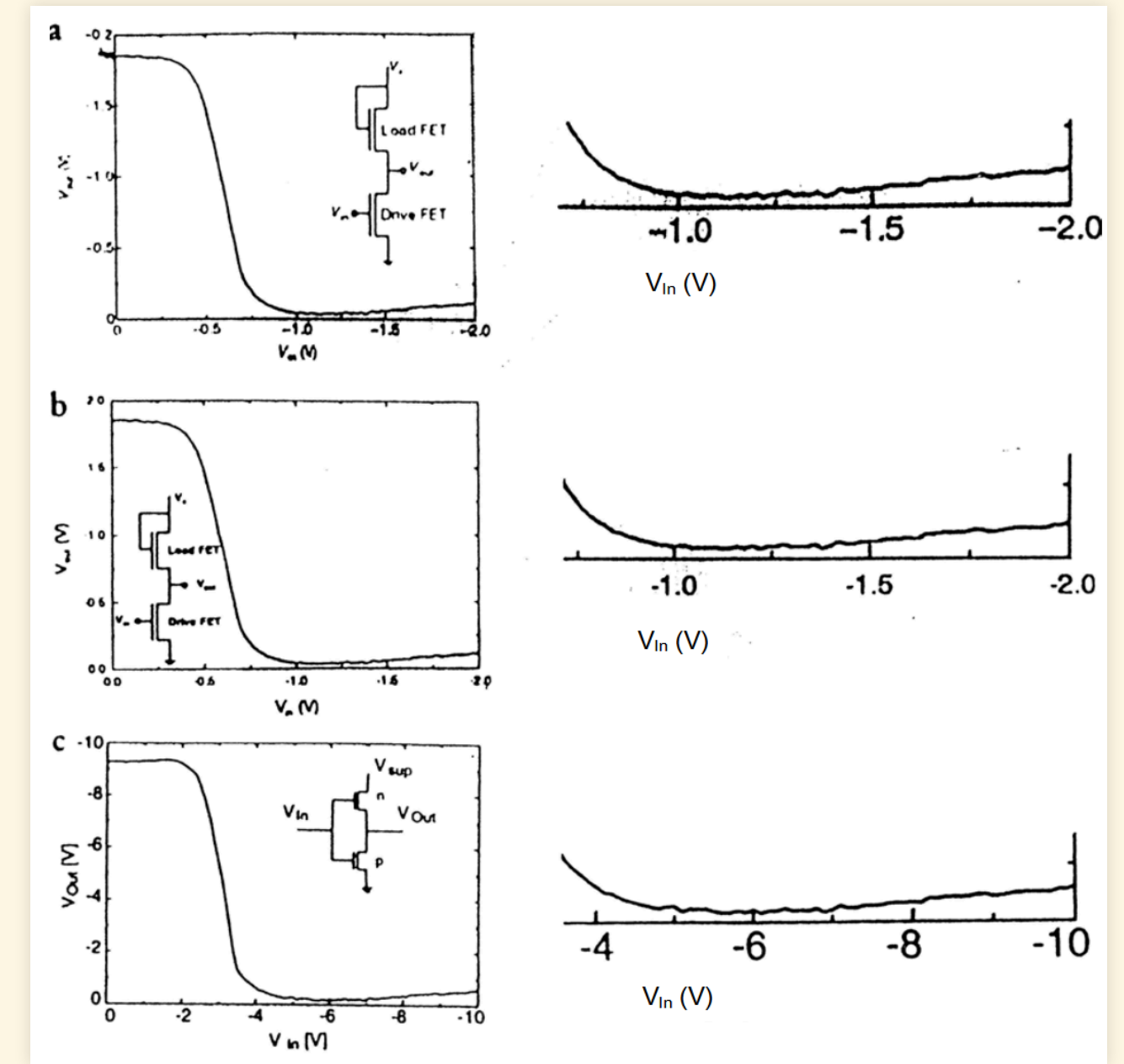
Reproducible Research

- Science is supposed to be reproducible
 - Scientific laws assume that anyone testing them will observe consistent results
 - Radioactive decay of ^{14}C occurs at the same rate in every laboratory.
 - Melting point of pure quartz is the same in every laboratory.
- Engineering, medicine, economics, business, ... are less precise.
 - But reproducible analyses are important:
 - Research should reproduce statistical patterns.
 - Applied work (clinical medicine, engineering designs, business records) should be accurate.
- Clarity & Transparency
 - Build trust by sharing all the steps of the work
 - Avoid errors by making it easy to check your work
 - Simplify correcting errors, by automating analysis and reporting
 - Speed up routine analysis and reports

Urgency of Reproducible Methods

- Fraud thrives where methods and details are secret

- Jan Hendrik Schön
- Physicist at AT&T Bell Labs 1997–2002
- Published dozens of papers per year
 - Claimed many major breakthroughs
- No one else could reproduce his results
- 2002: Researchers discovered he was re-using graphs
 - An investigation revealed massive fraud
- None of the other scientists working with him asked to witness his experiments or review his analyses.



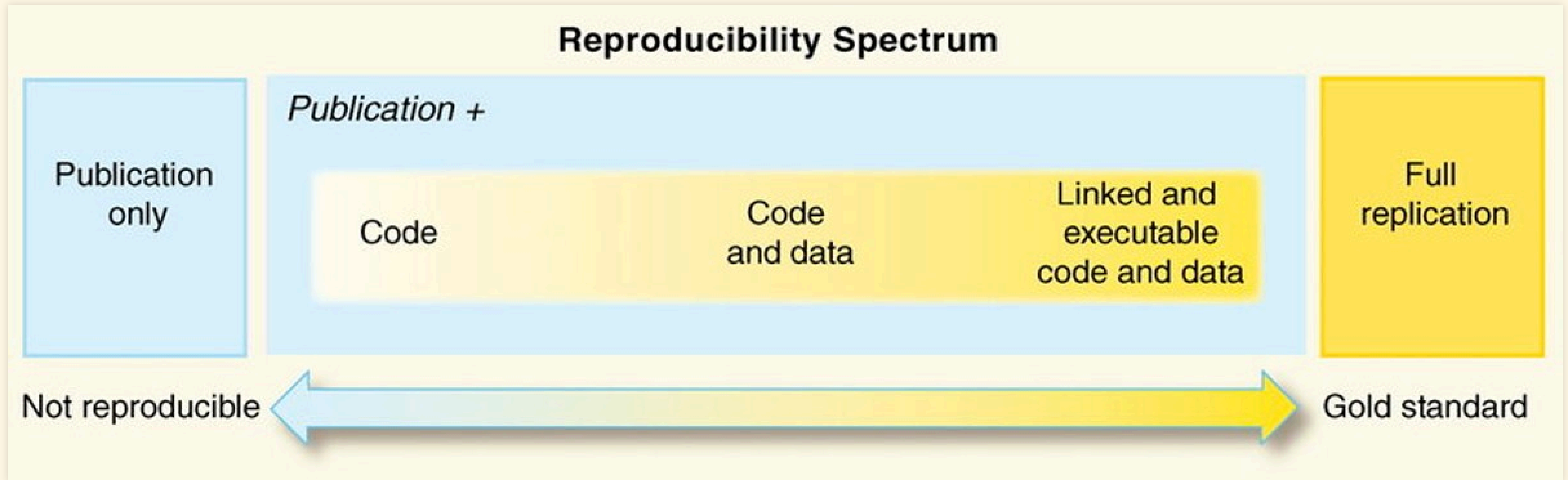
Honest Error

- London Whale: JP Morgan-Chase lost \$6.2 billion because of a spreadsheet error.
- UK Public Health Service undercounted 16,000 COVID cases.
- Many European countries made policy based on a flawed economics analysis.
- 16% of papers in 8 top medical journals have strong indications of fraud or incompetent analysis.
- Errors in computer code led scientists to think satellites showed global cooling.
- Psychology has been experiencing a “crisis of replication” where many major results can’t be reproduced.

Response: Reproducible Research

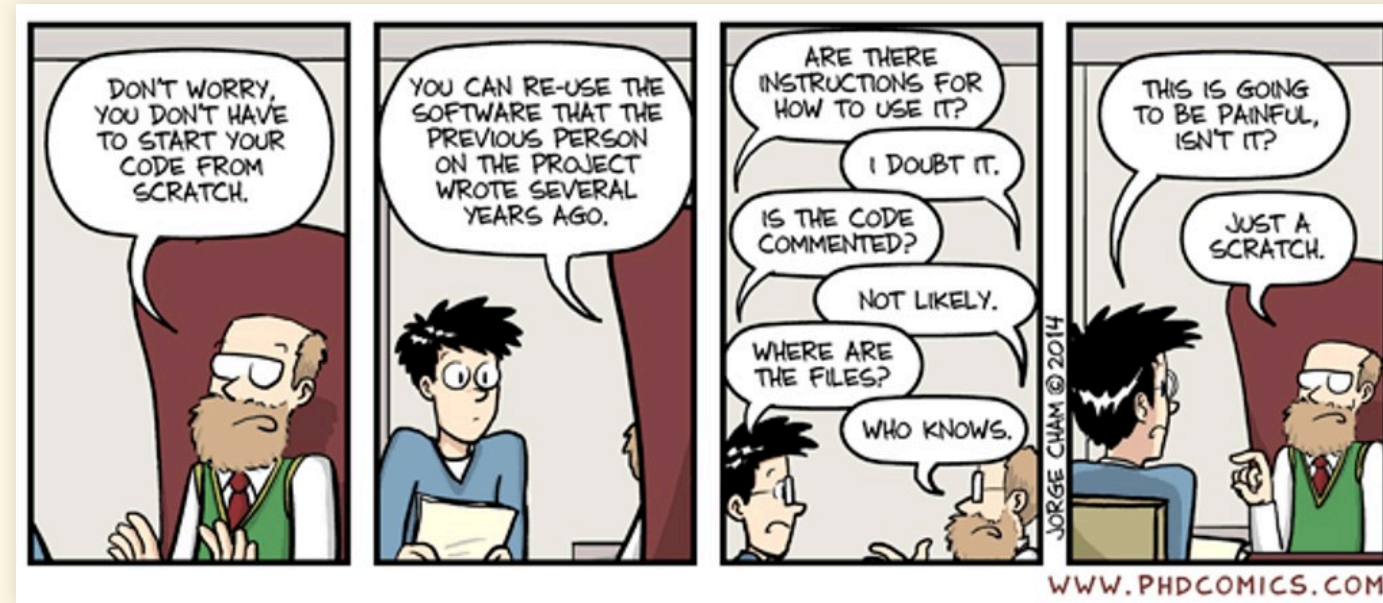
- To restore trust and confidence in research:
 1. Document everything clearly:
 - Methods for collecting data
 - Methods for analyzing data
 - Make sure everything is clear and detailed
 2. Make code and data publicly available
 3. Archive code and data to keep it available over time
- Reproducibility doesn't just mean reproducing the exact result:
 - It's also being able to make changes and exploring their consequences

Reproducibility Spectrum



Literate Programming

- Reading & understanding someone else's computer code is difficult



- Literate programming (D. Knuth, 1984)
 - Integrate code with text that clearly explains it
 - Explanations:
 - How it works
 - What it does
 - How to use it

Testing Code

- Unit testing
 - Assume your code has errors



- Write tests to catch errors

- Fake Data (Gelman, 2009)
 - Before analyzing research data, test your analysis with “fake data”
 - Generate data where you know what the result should be
 - Test whether your analysis code gets the right answer

Using `git` and GitHub

Tracking Changes

- Your data-analysis code works, and you edit it to add some new capabilities, and it breaks down.
- `git` can remember all the changes you make.
 - If you edit your code and it breaks, you can see what you changed.
 - You publish a paper using your code
 - You keep working on the code, adding new features.
 - Two years later, someone asks a question about the paper.
 - Can you recreate the code from that paper, to answer their question?



What `git` does

- Organize each project in a folder & subfolders
- **Workspace:** The files and folders in your project
- **Repository:** The records `git` keeps of your changes
 - You decide when to record changes
 - **Stage:** select changed files that you want to track
 - **Commit:** `git` makes a snapshot of the changes in all the *staged* files and records it in the repository
 - A **commit** records changes in multiple files
 - The repository contains a record of the current state of your project (as of the last commit), and all the previous commits.
 - You can easily recover all the files the way they were at any previous commit
- You can **push** or **pull** commits to synchronize two different repositories
 - A local repository on your computer and a remote repository somewhere else (e.g., on GitHub)

Git structure



Local Factory
(where you're making the product)



Forklift
(preparing to put it in the warehouse)



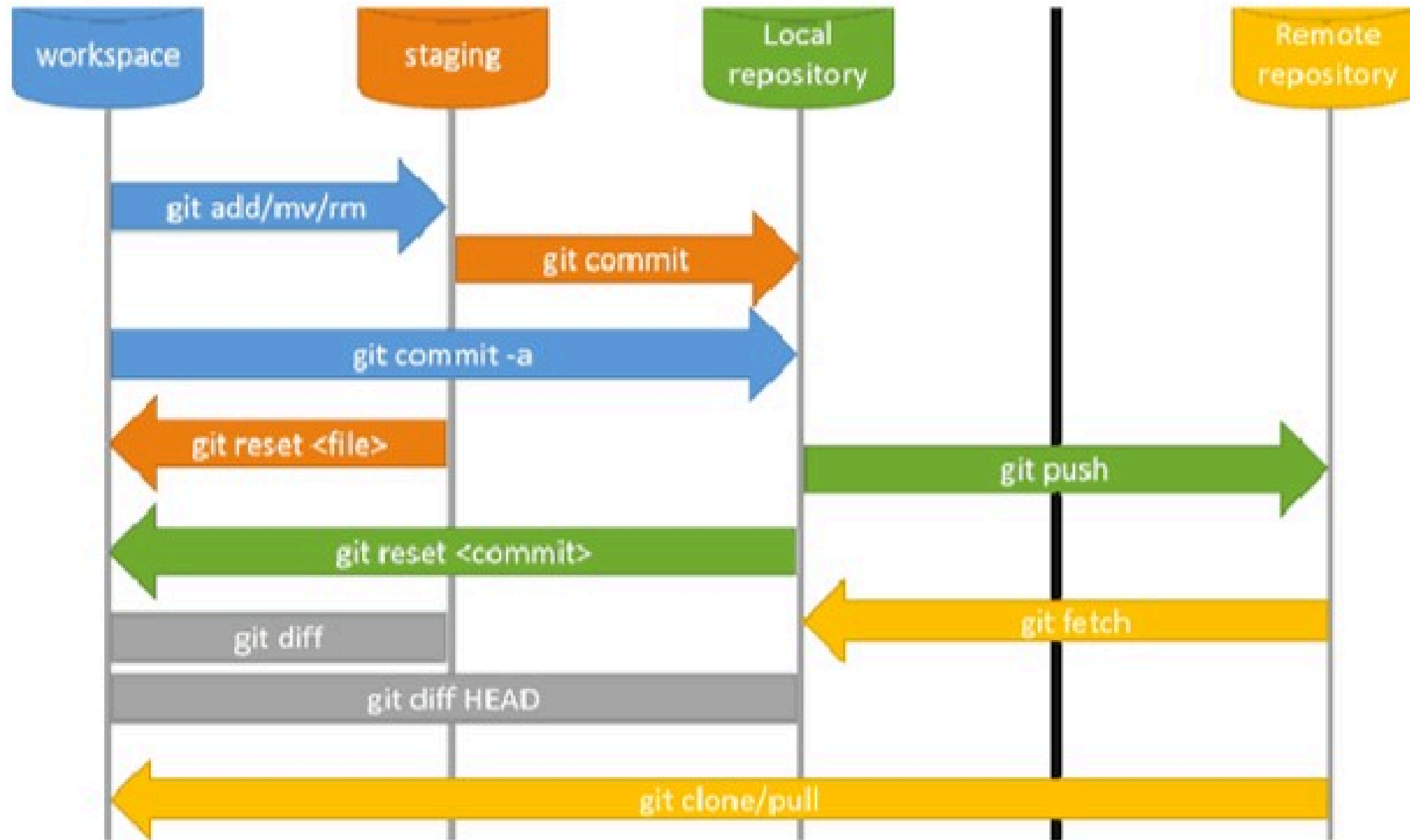
Private warehouse
(holds the product)



Distribution Center
(Can push the product to the distribution center)



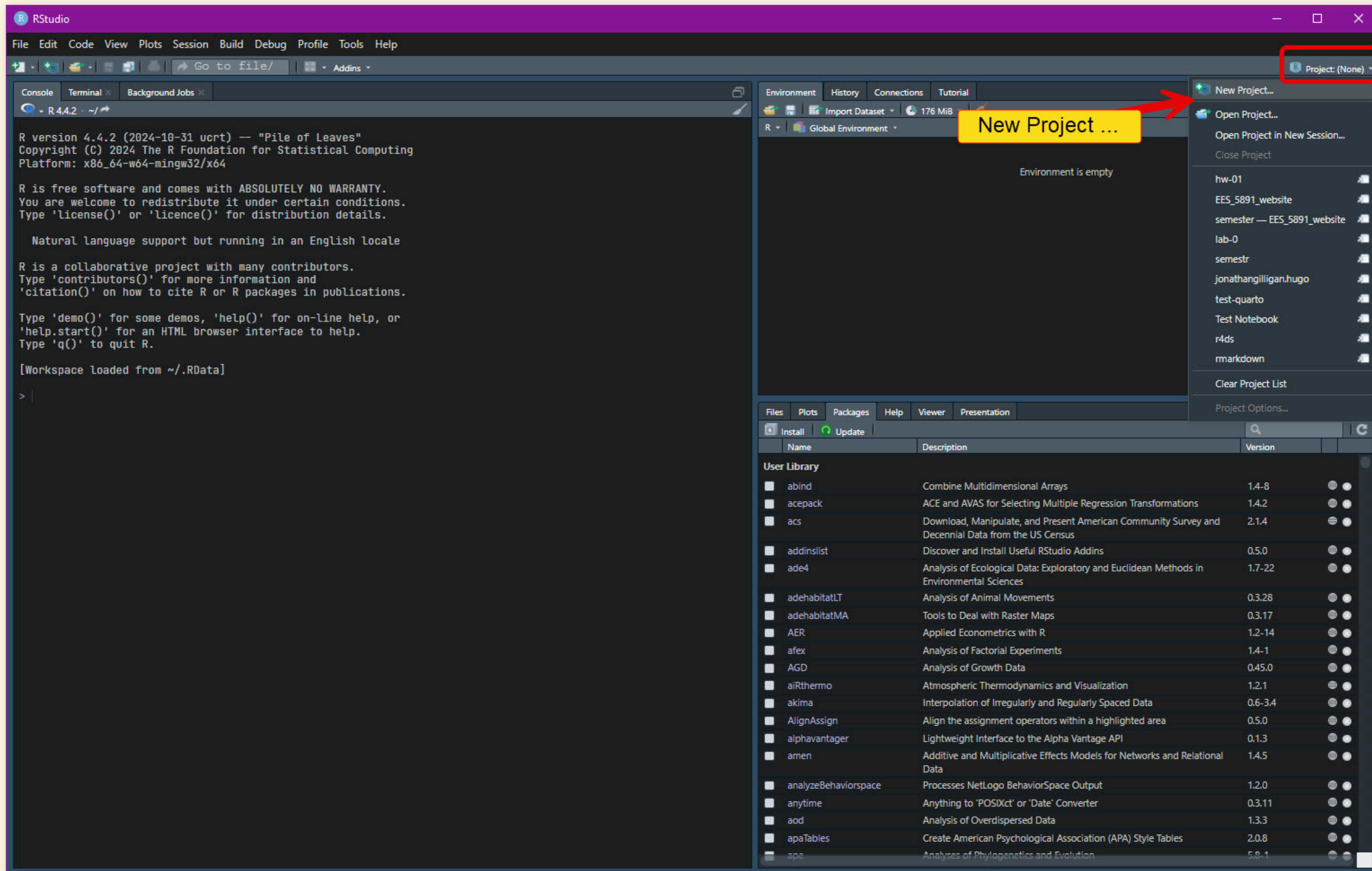
Git Process



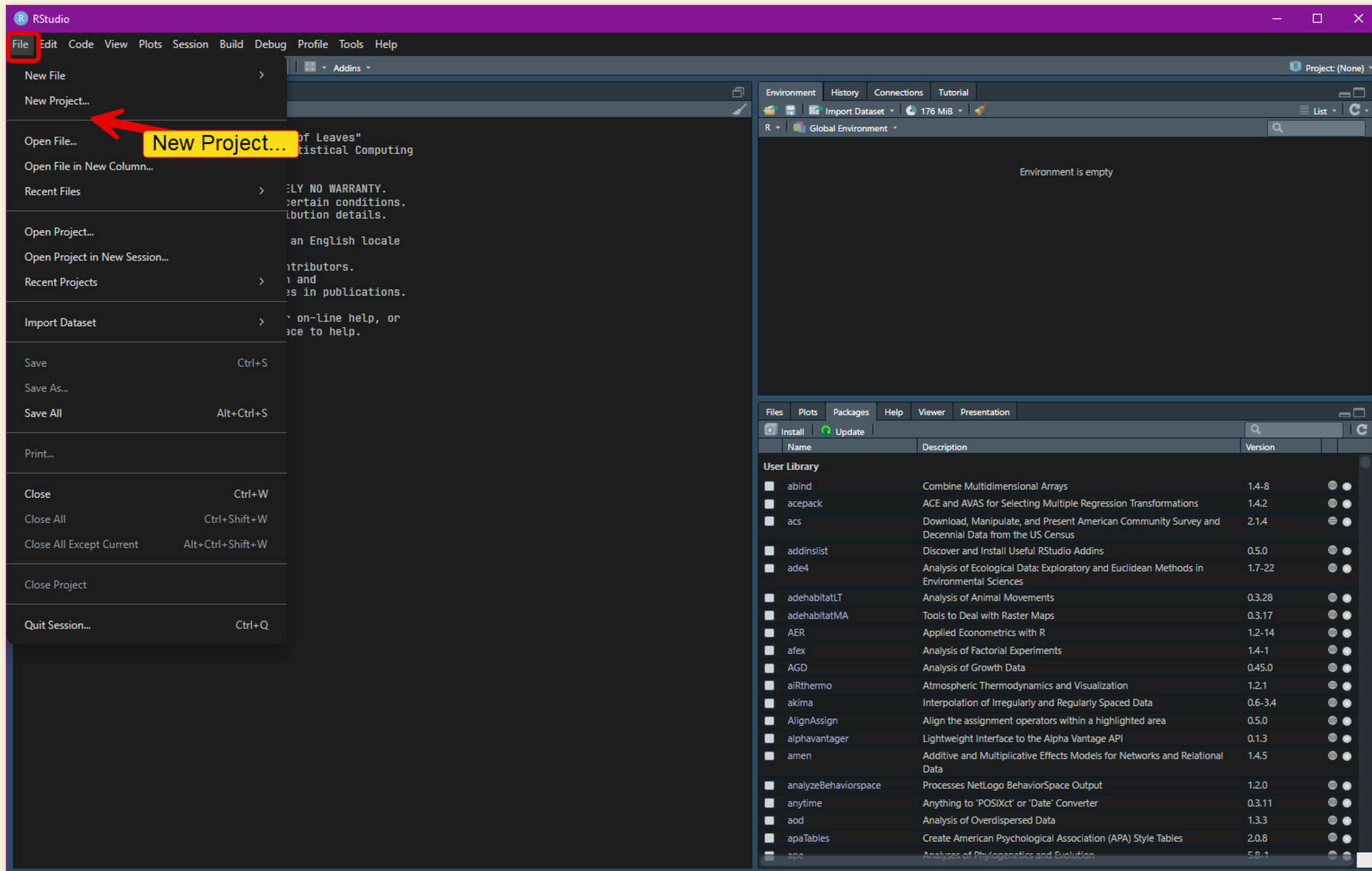
Working with `git` and RStudio

- RStudio works with **projects**: Collections of files in a folder and subfolders
- Each project you work on should be in an RStudio project
 - Use RStudio's "New Project" command to make a new project
- If you're creating a project from a GitHub repository, use the RStudio "New Project" and choose "Version Control", and then put the URL for the GitHub project into the box

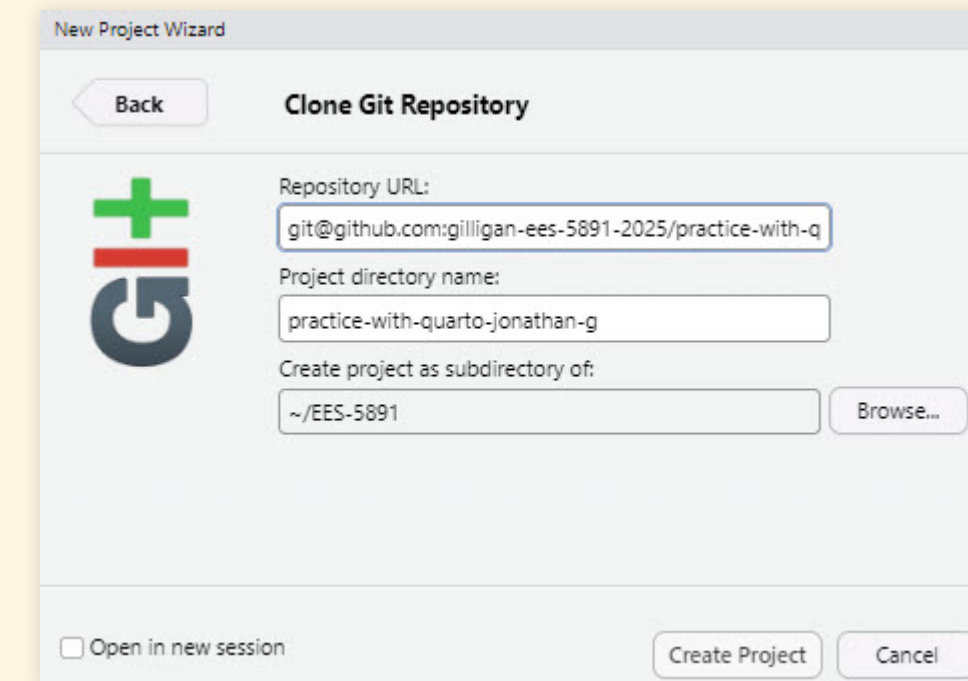
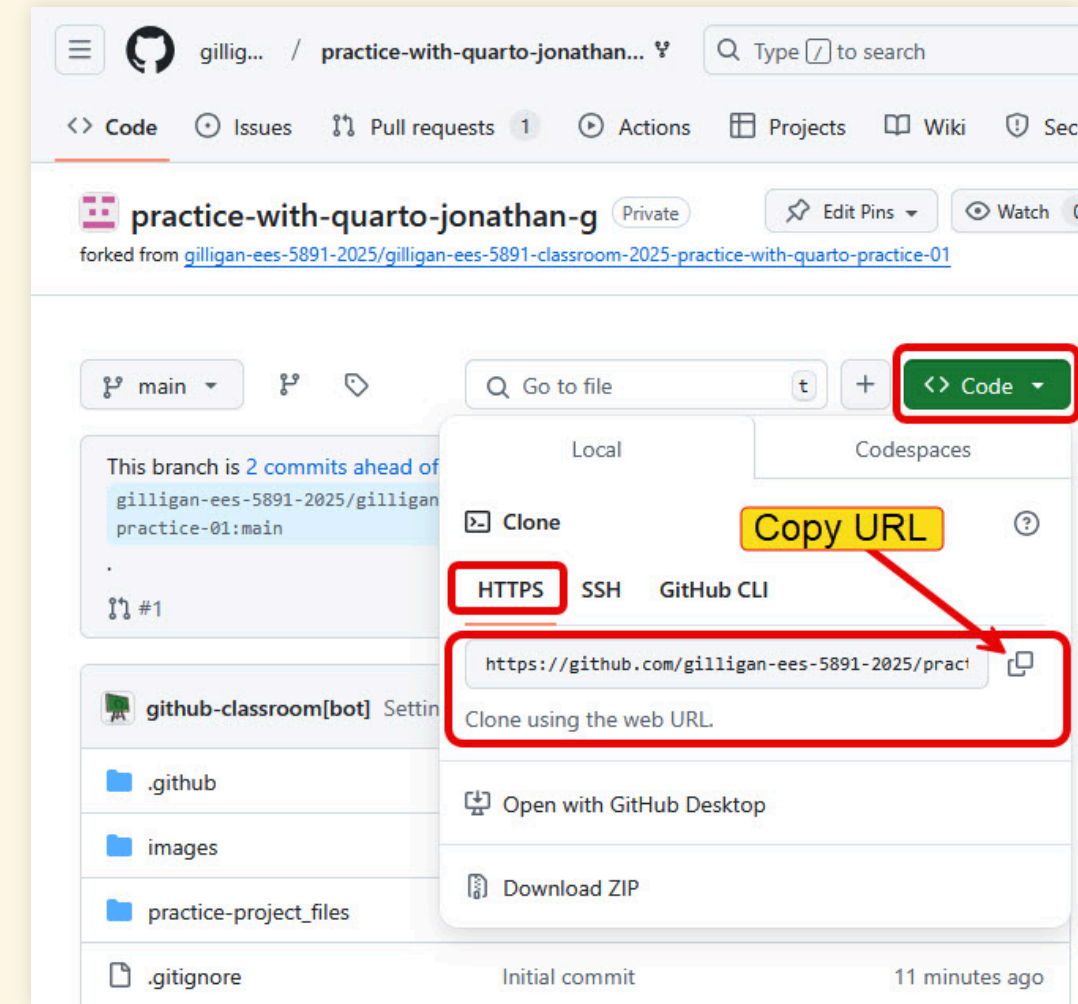
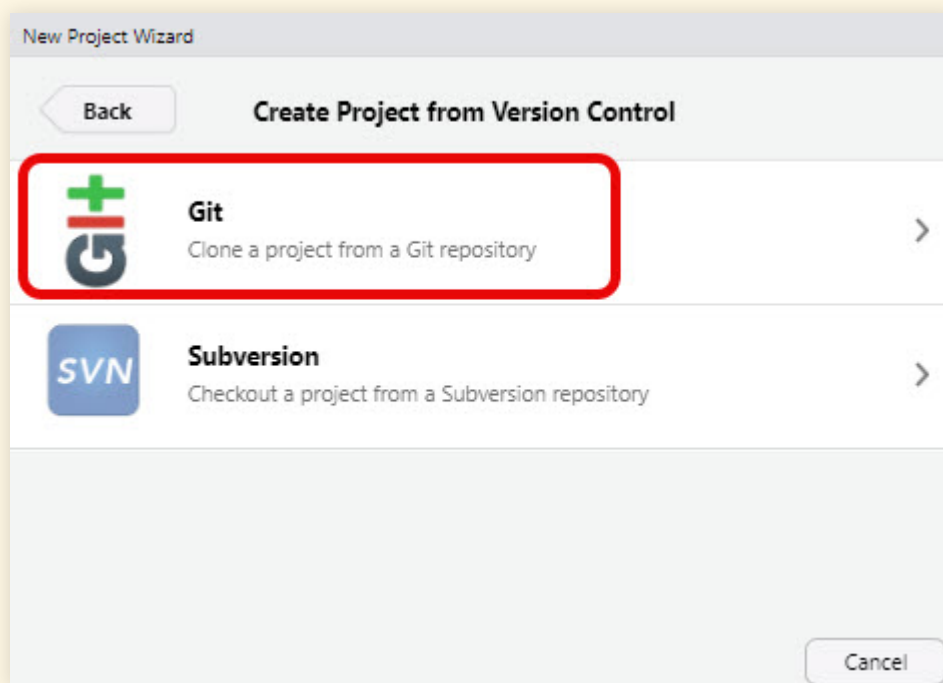
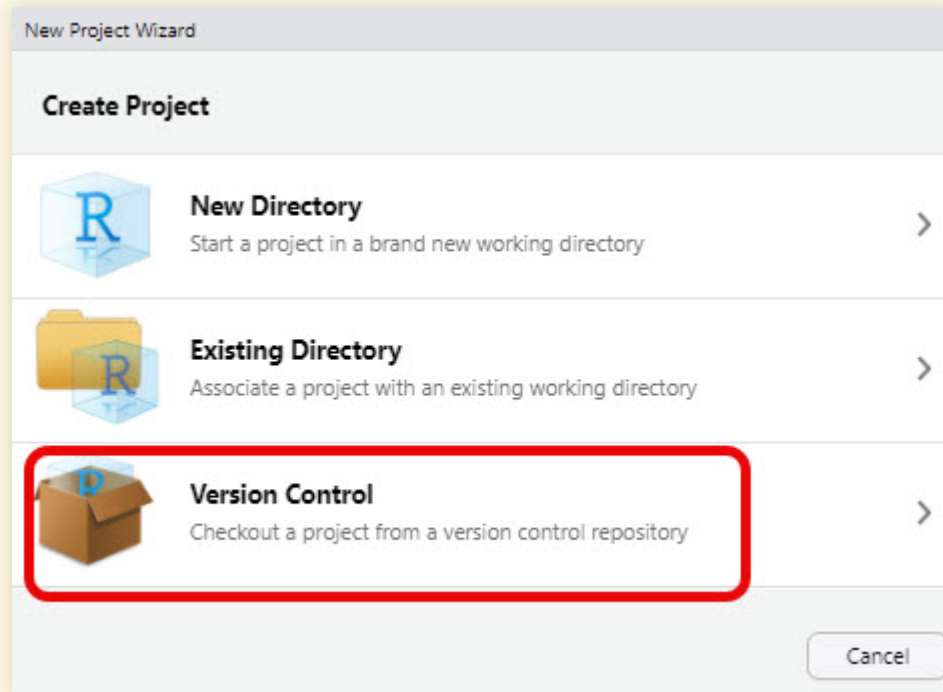
Creating a New Project







Creating a New Project

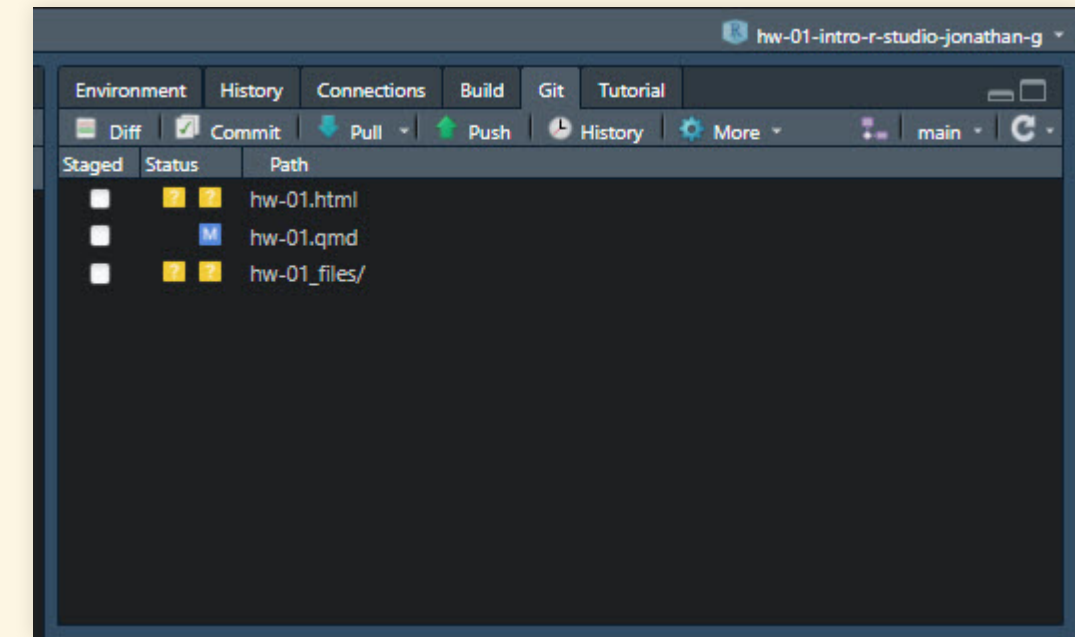


Creating an RStudio Project from GitHub



Working with Git in RStudio Projects

- The `git` panel shows which files have been changed:
 -  New (added) file
 -  Modified file
 -  Deleted file
 -  Renamed or moved file
- The “Diff” button examines changes in text files
- The “Commit” button lets you stage and commit changes



Viewing File Differences

The image shows the RStudio 'Review Changes' window. The top bar is purple with the R logo and the title 'RStudio: Review Changes'. Below the bar, there are tabs for 'Changes', 'History', and 'main'. To the right of these tabs are icons for 'Stage', 'Revert', and 'Ignore', and buttons for 'Pull' and 'Push'. The main area is divided into two panes. The left pane shows a list of files: 'hw-01.html', 'hw-01.qmd' (selected), and 'hw-01_files/'. The right pane is for the 'Commit message' and has a 'Commit' button. Below the file list, there are options to 'Show' 'Staged' or 'Unstaged' changes, a 'Context' dropdown set to '5 lines', and checkboxes for 'Ignore Whitespace', 'Stage All', and 'Discard All'. The main pane displays the differences for the selected file 'hw-01.qmd'. It shows a diff view with line numbers on the left and the diff content on the right. The diff is color-coded: red for deletions and green for additions. The changes are grouped into sections separated by '@@' markers. The first section shows changes to the title and author information. The second section shows changes to the content of the file, including a new line for 'best = read_csv("best.csv")'. The third section shows changes to the content of the file, including a new line for 'head(best, 10)'. The fourth section shows changes to the content of the file, including a new line for 'glimpse(best)'. The fifth section shows changes to the content of the file, including a new line for '# 2 Visualize the data'.

RStudio: Review Changes

Changes History main Stage Revert Ignore Pull Push

Staged Status Path

- hw-01.html
- hw-01.qmd
- hw-01_files/

Commit message

☐ Amend previous commit Commit

Show Staged Unstaged Context 5 lines Ignore Whitespace Stage All Discard All

@@ -1,9 +1,9 @@

1 1 ---

2 2 title: "HW-01: Data Wrangling and Visualization with R"

3 3 date: "2025-01-14"

4 4 author: "Your Name"

3 3 date: "2025-01-21"

4 4 author: "Jonathan Gilligan"

5 5 ---

6 6

7 7 # Homework #1: Data wrangling and visualization with R

8 8

9 9 In this assignment you will be getting started with R to do something very

@@ -166,11 +166,11 @@ gistemp = read_csv("gistemp.csv")

166 166

167 167

168 168 **Question 1.1** Load `best.csv` and assign it to a variable `best`

169 169

170 170 ```{r}

171 171 best = read_csv("best.csv")

172 172

173 173

174 174 ## 1.2 Inspect the data

175 175

176 176 When you read data from a spreadsheet or CSV file into R, it takes the form of

@@ -213,11 +213,13 @@ You can inspect a data frame or tibble in several ways:

213 213

214 214

215 215 **Question 1.2** Use `head()` and `glimpse()` to show `best`.

216 216

217 217 ```{r}

218 218 head(best, 10)

219 219

220 220 glimpse(best)

221 221

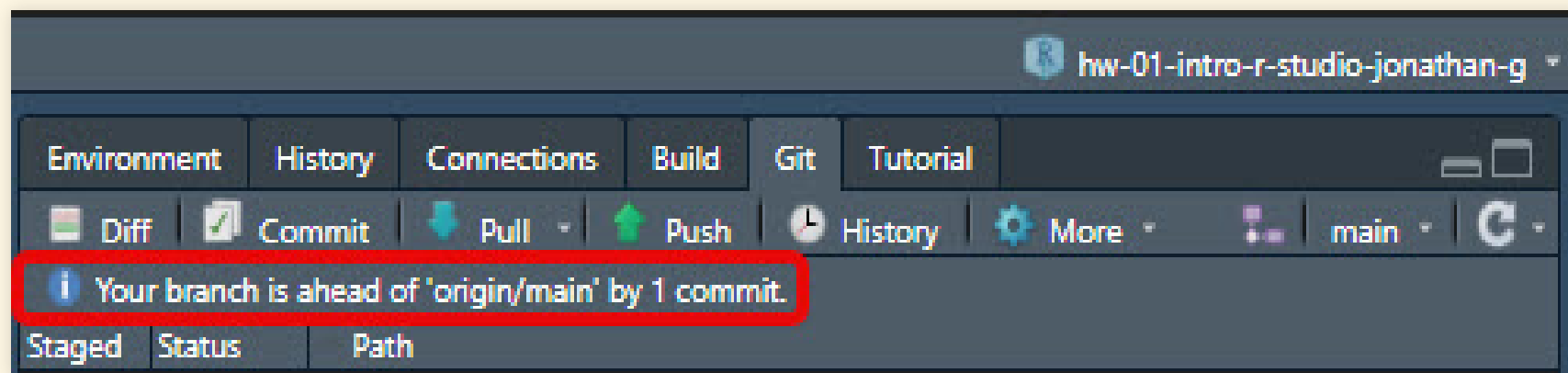
222 222

223 223 # 2 Visualize the data

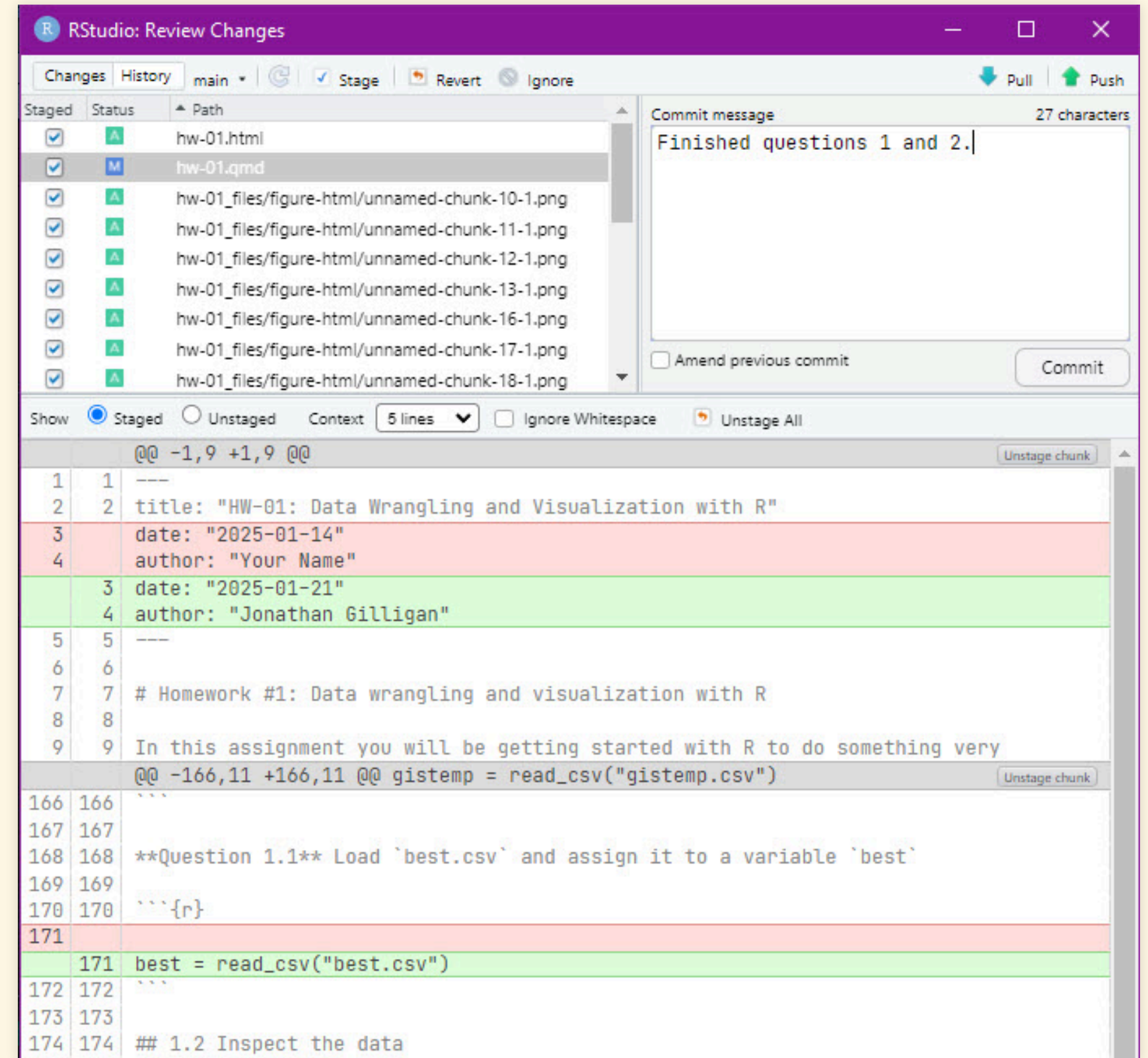
224 224

Staging and Committing Changes

- Stage individual files by checking boxes
- When you're ready to commit:
 - Write a comment in the comment box, describing the commit.
 - Click the "Commit" button
- After you commit, RStudio will remind you that your local repository is out of sync with GitHub



- Click on the "Push" button to send your recent commits to GitHub
- The "Pull" button will check whether GitHub has commits that aren't on your computer, and pull them down.



Using Quarto and RMarkdown

- Based on the Markdown standard for formatting with plain text.
- RStudio has added a visual WYSIWYG editor that lets you format without knowing Markdown codes.
 - You can toggle back and forth between “Code” and “Visual” modes.
- Click on “Render” to turn your Quarto document into a formatted HTML or PDF document.
 - It can also export to Word, Powerpoint, and other formats.

History:

- Yihui Xie, a graduate student in statistics. got interested in reproducible research and developed `knitr` to integrate different kinds of textual markup with R to make reproducible documents.
- Later, he was hired by RStudio and developed the `RMarkdown`, which allowed sophisticated integration of R into producing many kinds of documents.
- Many people wanted to use RMarkdown with other languages, such as Python and Julia,
 - and there was also interest in integrating RMarkdown with the Jupyter notebook system
- Quarto combined features of Jupyter, RMarkdown, and other systems.
 - It is very powerful and customizable

Some basic RMarkdown:

- Headings: #, ##, ###, ...
- Lists:

- Bulleted:

```
* blah blah blah
* blah blah blah
  * foo foo foo
* blah blah blah
```

- Numbered:

```
1. first item
2. second item
  a. Multiple levels
  b. And so forth
  b. Labels increment
    automatically
#. third item
#. fourth item
```

- Text formatting:
 - ****bold**** = **bold**
 - *_italic_* = *italic*

- Block quotations:

```
> This is a block quotation
> that goes on for several
> lines
>
> With multiple paragraphs
```

- inline R code ``r sqrt(2)``
- Hyperlinks:
 - `<https://vanderbilt.edu>` makes <https://www.vanderbilt.edu>
 - `[Vanderbilt]`
`(https://www.vanderbilt.edu)` makes Vanderbilt
- Images:
 - `![alt text](/path/to/image.jpg)`

Trying It Out

