# Organising your projects in Rstudio for reproducibility

Marina Vabistsevits

# What this session is about

1. Organising your projects with .Rproj aka **project-oriented workflow**
2. Rstudio efficiency tips
3. .Rproj with git
4. (.Rproj with .renv)

*We are going to try things out
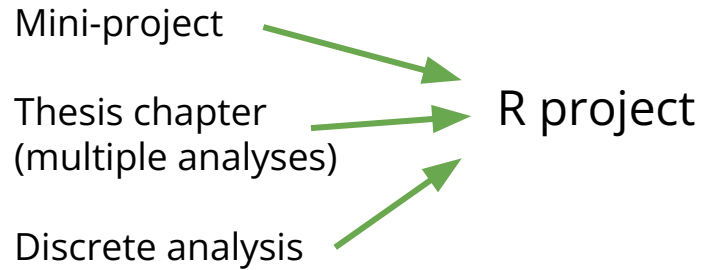in mini-breaks after each part*

Icon for R

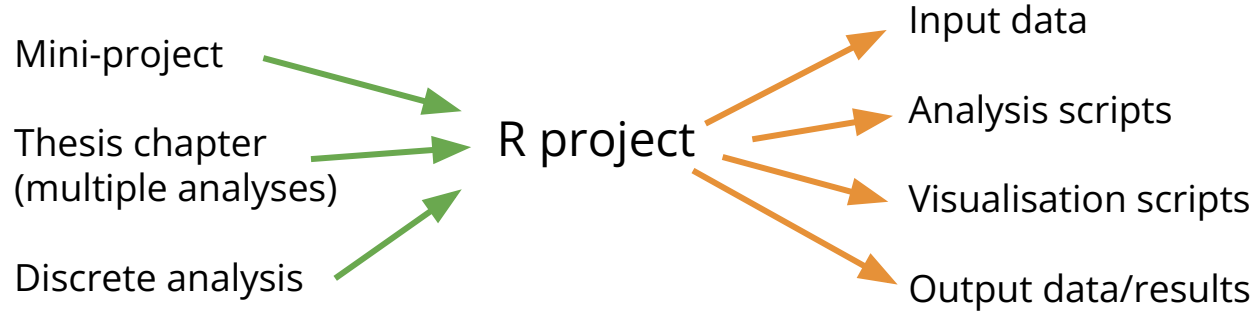Icon for RStudio

posit™

software
company who
develop
Rstudio
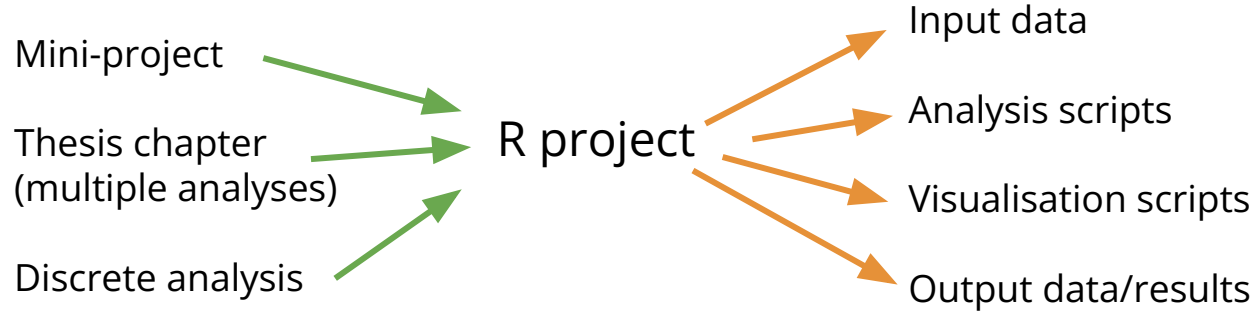
*The slides will be shared after!*

# 1. Project-oriented workflow

# R projects

Mini-project

Thesis chapter
(multiple analyses)

Discrete analysis

R project

# R projects

Mini-project

Thesis chapter
(multiple analyses)

Discrete analysis

R project

Input data

Analysis scripts

Visualisation scripts

Output data/results

# R projects

Mini-project

Thesis chapter
(multiple analyses)

Discrete analysis

R project
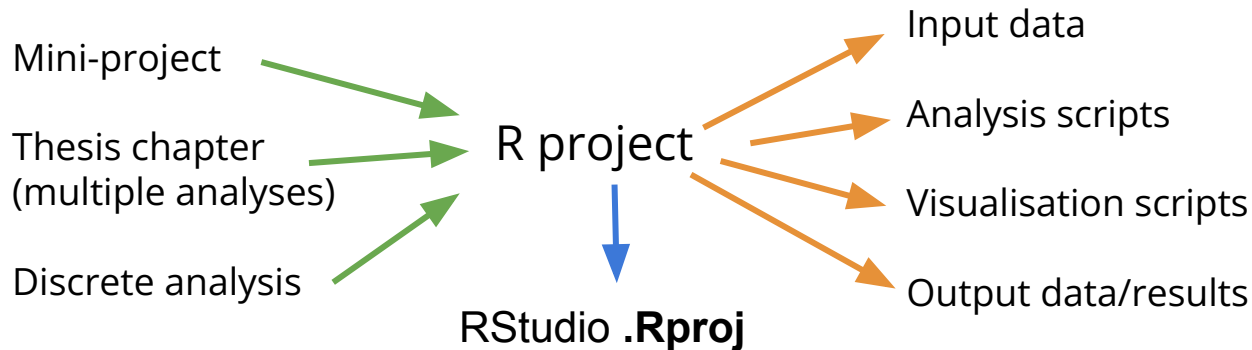
Input data

Analysis scripts

Visualisation scripts

Output data/results

How you do organise your projects? scripts? folders?

Where do you store and access your data? results? figures?

Does it matter? Are your analysis/projects reproducible?

Can you organise your projects better and make life easier for future self/colleagues?

# R projects

Mini-project →

Thesis chapter (multiple analyses) → **R project** → Input data

Discrete analysis →

↓

RStudio **.Rproj**

→ Analysis scripts

→ Visualisation scripts

→ Output data/results

How you do organise your projects? scripts? folders?

Where do you store and access your data? results? figures?

Does it matter? Are your analysis/projects reproducible?

Can you organise your projects better and make life easier for future self/colleagues?

Reproducibility can be enhanced through intentionally organising projects with

# Project-oriented workflow

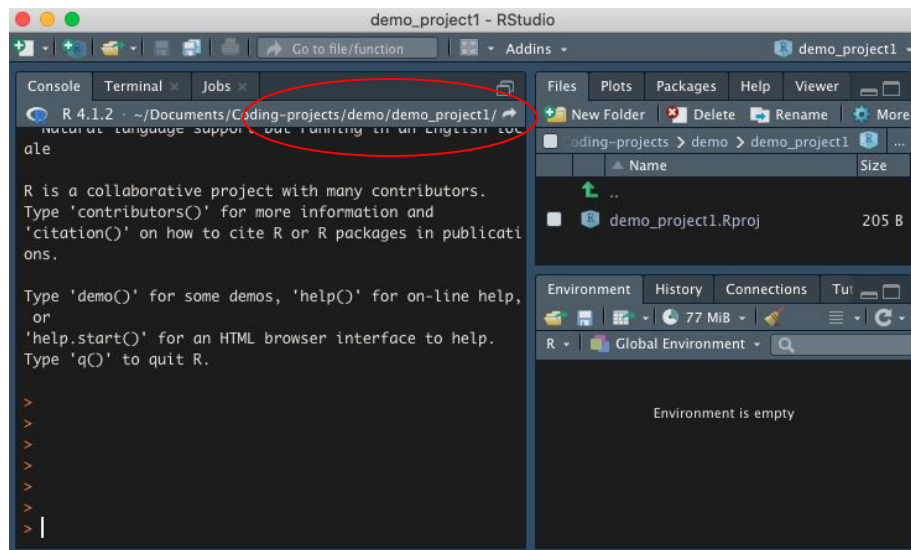Use Rstudio / .Rproj for your data analysis projects



*This means that you are essentially compartmentalizing your current project*
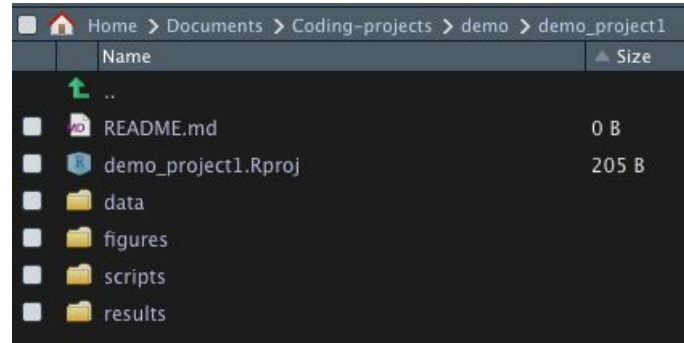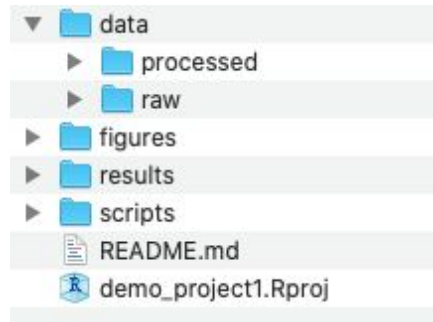
# Project-oriented workflow

## Use Rstudio / .Rproj for your data analysis projects

- Project directory stores all your data, scripts
- The working directory is set to the project directory (e.g. demo_project1), so don't need to specify full paths to data (only internal subfolders)


- The project creates everything it needs, within its workspace/folder, and touches nothing it did not create
- Any scripts are written assuming they will be run from a fresh R session within the project
- The project folder can be moved _anywhere_, and everything will still work (no paths will be broken)
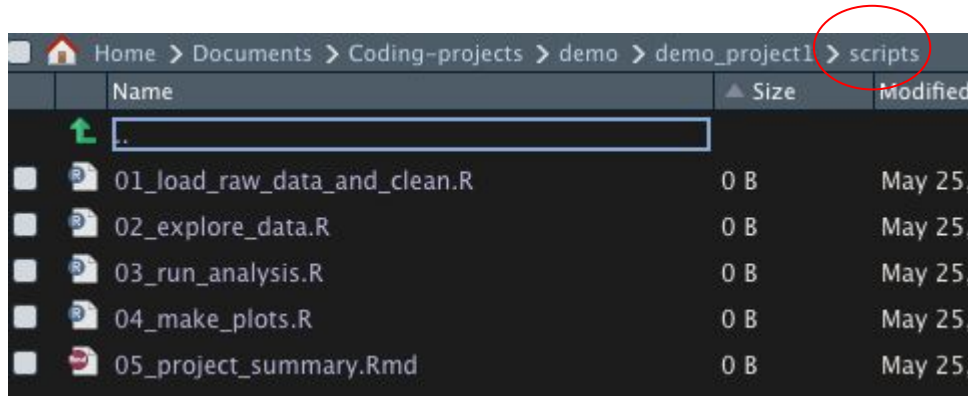
# Organise your projects intentionally

# Take advantage of default ordering

01_load_raw_data_and_clean.R
02_explore_data.R
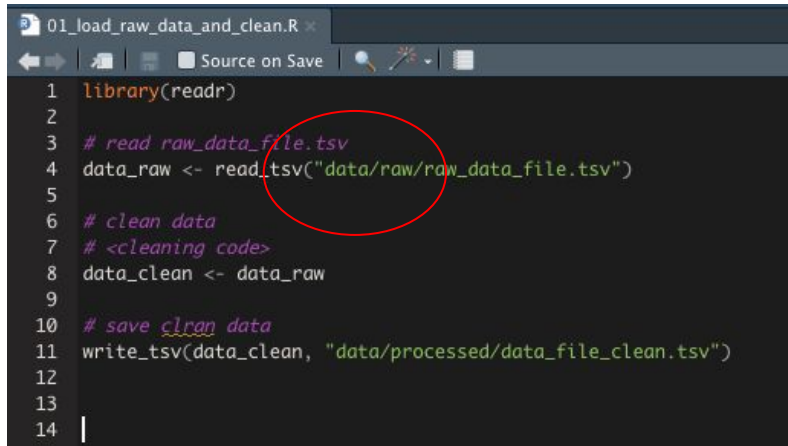03_run_analysis.R
04_make_plots.R
05_project_summary.Rmd

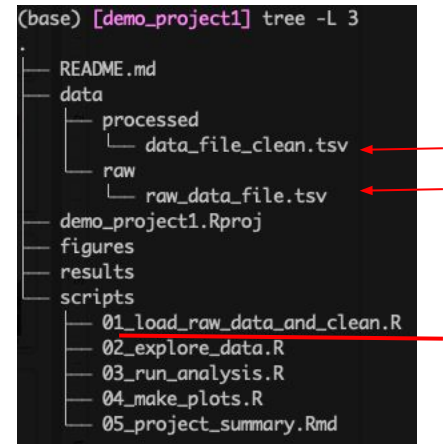Can have many parts of the analysis separately - save interim results as files and re-read then in the next script

Home > Documents > Coding-projects > demo > demo_project1 > scripts

| | Name | Size | Modified |
|---|---|---|---|
| | . | | |
| | 01_load_raw_data_and_clean.R | 0 B | May 25, |
| | 02_explore_data.R | 0 B | May 25, |
| | 03_run_analysis.R | 0 B | May 25, |
| | 04_make_plots.R | 0 B | May 25, |
| | 05_project_summary.Rmd | 0 B | May 25, |

# Don't use `setwd()`

Keeping your work as an .**Rproj** will help you manage your file paths

```
setwd("path/that/only/works/on/my/machine")
```





No need to hardcode paths when using Rproj

# Don't use `rm(list = ls())`

Restart R daily (or every time you start working after a break) to ensure a clean environment



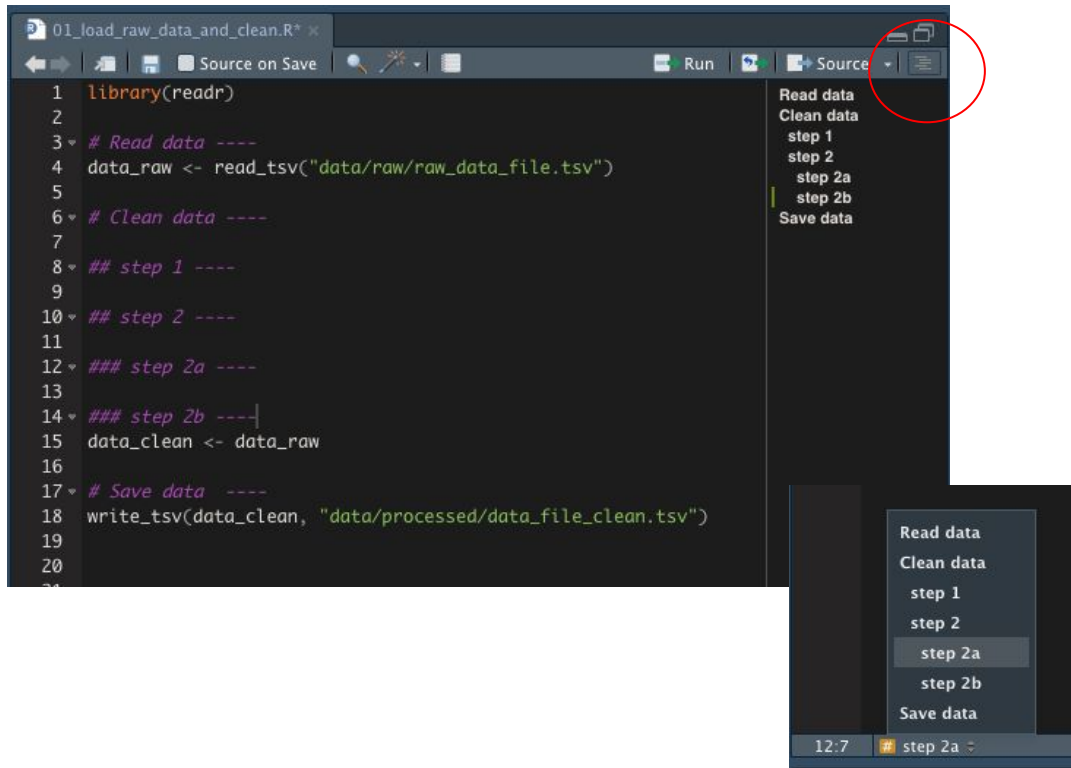And !! do not save your .Rdata workspace

(untick and select 'never')





Save your 'real' work, delete the rest

5 mins to try it

# 2. Various efficiency tips for R

# Name your code sections and use then for quick navigation



- Use section headings:

```
# section    ----

## subsection    ----

### subsubsection    ----
```
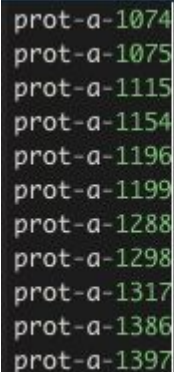
- Great for navigating in long scripts
- Can fold sections

# Vertical selection

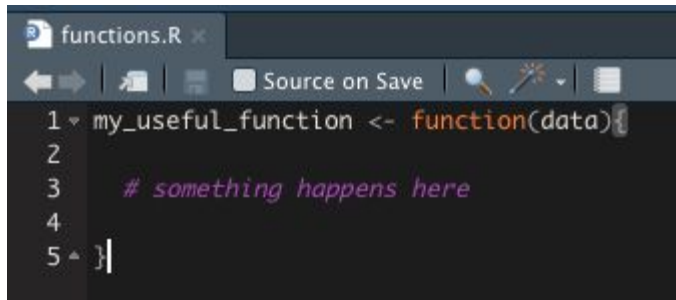(hold *option* or *alt* and drag cursor down to select vertically)



Great for e.g.
- commenting out a block of code with #
- adding " " around a column of ids

# Jump to function definition or open data frame



Cmd + mouse click on the name

(opens in a new window)



> also works on functions from external packages (if you want to check what they do internally)

# Keyboard shortcuts



-  (option + Enter)
- `<-` (option/alt + " - ")
- `%>%` (control + shift + M)
- ```` ```{r} ```` (control + shift + I)

  ```` ``` ````

Tools   Window   Help

Install Packages...
Check for Package Updates...

Version Control ▶

Shell...
Terminal ▶
Jobs ▶
Addins ▶
Memory ▶

Keyboard Shortcuts Help   ⌥⇧K
Modify Keyboard Shortcuts...
Edit Code Snippets...
Show Command Palette   ⇧⌘P

Project Options...   ⇧⌘,

Global Options...   ⌘,

Keyboard Shortcuts

Show: ● All  ○ Customized   🔍 pipe ⊗              ? Customizing Keyb

| Name | Shortcut | Scope |
|---|---|---|
| Insert Pipe Operator | Ctrl+Shift+M | Editor |

# Move all libraries to the top



Install *packup* add-in:

```
devtools::install_github("milesmcbain/packup")
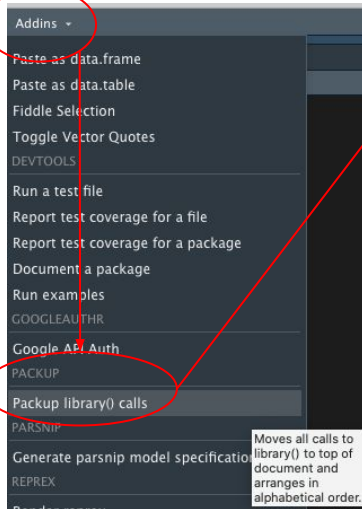```

Call it from Addins menu:



Moves all calls to library() to top of document and arranges in alphabetical order.

https://github.com/MilesMcBain/packup

# Any other R tricks to share?

# 3. Rproj for git users

*https://happygitwithr.com/*

# Git brief intro

**Git** - tool for code version control

- Tracking code changes
- Keeping older versions of the script
- Coding collaboration (with others or yourself in multiple locations)
- Tracking who and when made changes

If used consistently!

## Why use it?

- Can facilitate data/code integrity
- Improves reproducibility (e.g. keeps record of changes)
- Enables code sharing with colleagues / as a part of publication
- Important skill for anyone working with data (e.g. can showcase your work during job applications)

## Github / Gitlab etc

- Platforms where you can store your coding projects using git

# Why might you want to use git?



**Problem:**
- The same project is both locations (e.g due to data restrictions / tool availability)
- Don't want to copy over scripts all the time, as it might get messy!

Server: slade

Local Mac

# Why might you want to use git?

**Problem:**
- The same project is both locations (e.g due to data restrictions / tool availability)
- Don't want to copy over scripts all the time, as it might get messy!

Server: slade

Need a third place to sync two parts of the project

Local Mac

# Why might you want to use git?



Server: slade

Local Mac

# Why might you want to use git?



**Solution:**
- Use git to "collaborate with yourself"
- This enhances project reproducibility
- Keeps everything in one project folder
- Compatible with .Rproj

Server: slade

Local Mac

# Basic git commands / actions

git local repo will keep a record of the new changes that you want save

**Local**

**Remote**

Analysis folder on your local computer

| working directory | staging area | local repo | remote repo |
|---|---|---|---|

Your analysis folder repo on Github remote repo with all changes tracked, with a commit message that explains what has changes

git add

Tell git which files you changed/created and want to track with git

git commit

Send your updates to the remote repo on Github

Tell git what your new changes mean/do via a commit message

git push

git pull

git checkout

# Using git with your .Rproj

- Rstudio user interface makes it easy to get started with git
- (but also can use the terminal to run git commands)

# Create new repo on github:

**Repositories -> new**

## Create a new repository

A repository contains all project files, including the revision history. Already have a project repository elsewhere? **Import a repository.**

**Owner** *     **Repository name** *

[mvab ▾] / [demo-project2]  ✓

Great repository names are short and memorable. Need inspiration? How about **fuzzy-robot**?

**Description** (optional)

[                                                                    ]

◉  🖫  **Public**
         Anyone on the internet can see this repository. You choose who can commit.

○  🔒  **Private**
         You choose who can see and commit to this repository.

**Initialize this repository with:**

Skip this step if you're importing an existing repository.

☑  **Add a README file**
     This is where you can write a long description for your project. **Learn more.**

**Add .gitignore**

Choose which files not to track from a list of templates. **Learn more.**

[ .gitignore template: R ▾ ]

**1**

*.gitignore* and *README* come from github; *Rproj* was added by creating an R project

**2** Add folders etc

**3**

Add **data/** folder to *.gitignore* file so that your data files (if large or sensitive) are not committed to your project repo on Github

Git button - >
Commit

Tick files to commit:

Commit changes:                              Add message:



1

2

3

4

Your changes on Github:

# Adding a specific change:

5 mins to try it

# Using .Rproj for organising work

- "Work in a project" means:
  - **File system discipline:** all files related to a single project are stored in a designated folder;
  - **Working directory discipline**: intentionally work in project directory when opening Rproj
  - **File path discipline:** all paths are relative to the project directory (not hard-coded full paths!)
  - **Daily work habit:** Restarting R very often and re-running your script under development from the top will help you catch issues early on
  - **Using git for version control**


- Practising these habits together will give you the biggest pay-off
  - Reproducing your analyses will be easy
  - Organising your projects will help you make sense of them in 6/12/etc months
  - Can move your project anywhere or share it with anyone without changing paths

https://rstats.wtf/project-oriented-workflow.html#work-in-a-project

# Final thoughts / disclaimers

- Project-oriented workflow in R is not suitable/applicable to every scenario
  - Sometimes data is stored externally and can't be/too big to move (so can't use within-project paths)
- Not all work is done interactively in Rstudio
  - Some people use R from the terminal on the server - again, because of data access/size
  - Some analyses are computation-heavy and require to be submitted as scripts / run in parallel on server

- If your current workflow with `setwd()` works for you and your colleagues, consider future-proofing! ;)

# Recommended and used resources

https://www.tidyverse.org/blog/2017/12/workflow-vs-script/

https://richpauloo.github.io/2018-10-17-How-to-keep-your-R-projects-organized/

https://www.rforecology.com/post/organizing-your-r-studio-projects/

https://kkulma.github.io/2018-03-18-Prime-Hints-for-Running-a-data-project-in-R/

https://rstats.wtf/project-oriented-workflow.html

https://appsilon.com/rstudio-shortcuts-and-tips/

https://datacornering.com/my-favorite-rstudio-tips-and-tricks/

https://happygitwithr.com/

https://rstudio.github.io/renv/articles/renv.html

# .renv for managing package versions within projects



**renv** package helps you create reproducible environments for your R projects

- Keeps track of all your installed packages
- Installs them in a new location (i.e. easy to share)
- Useful for working on DNAnexus/AoU platforms

**See tutorial - https://rstudio.github.io/renv/articles/renv.html**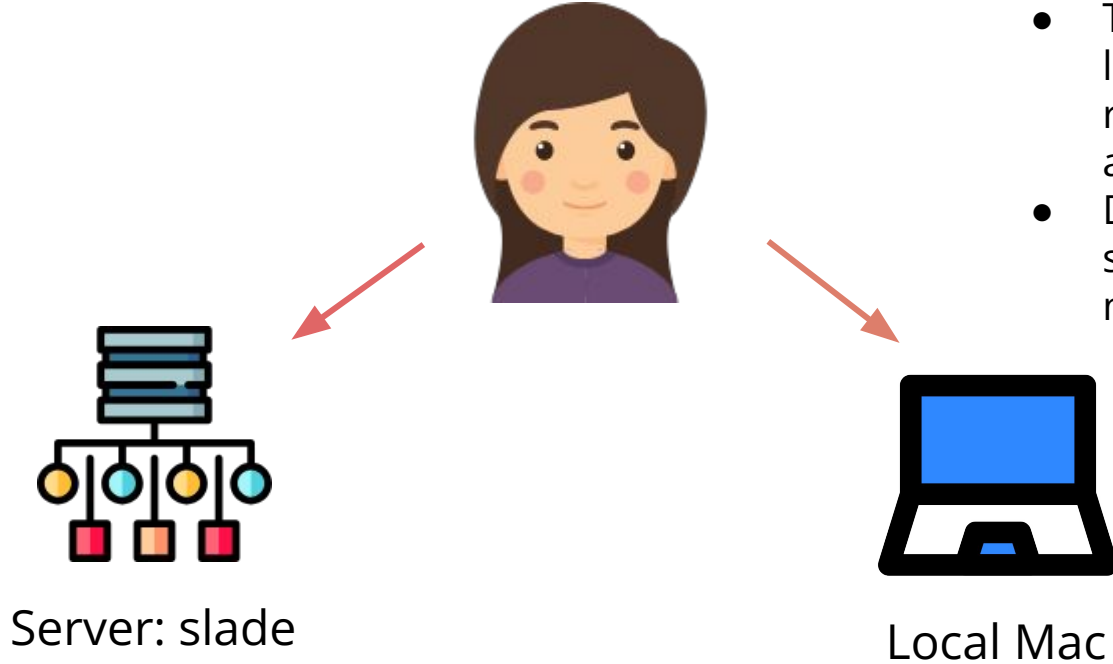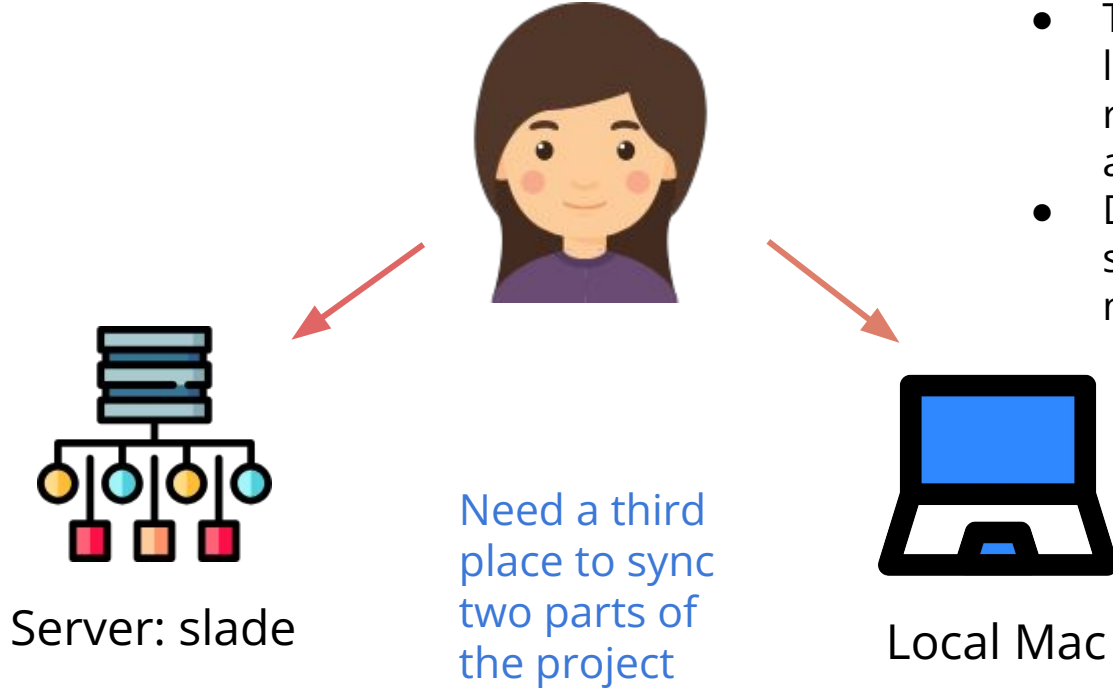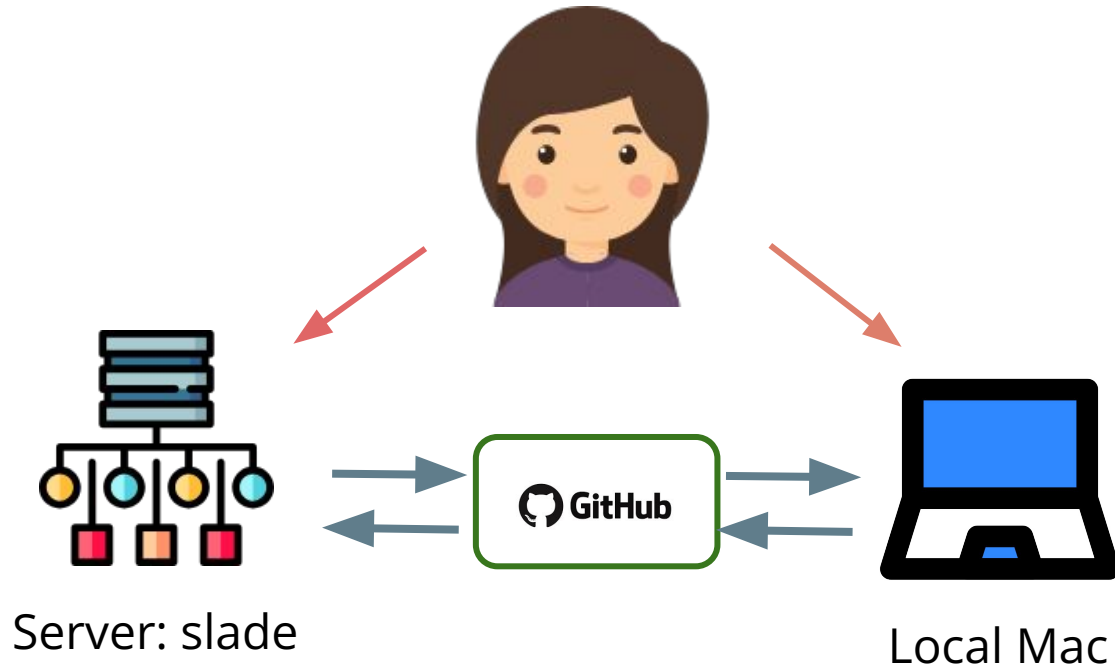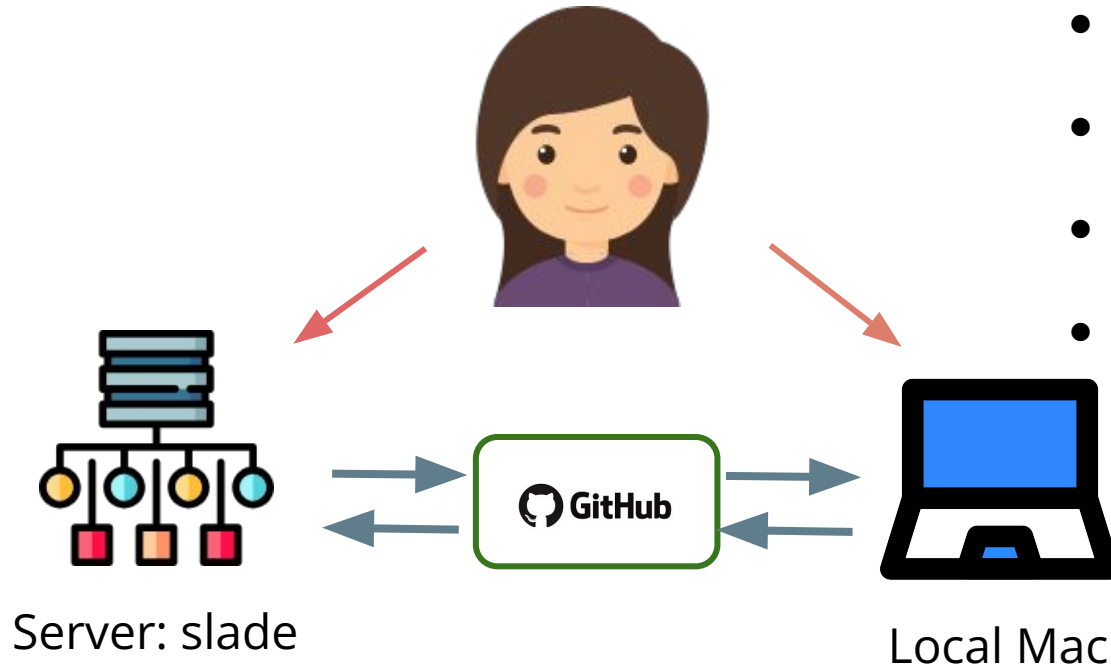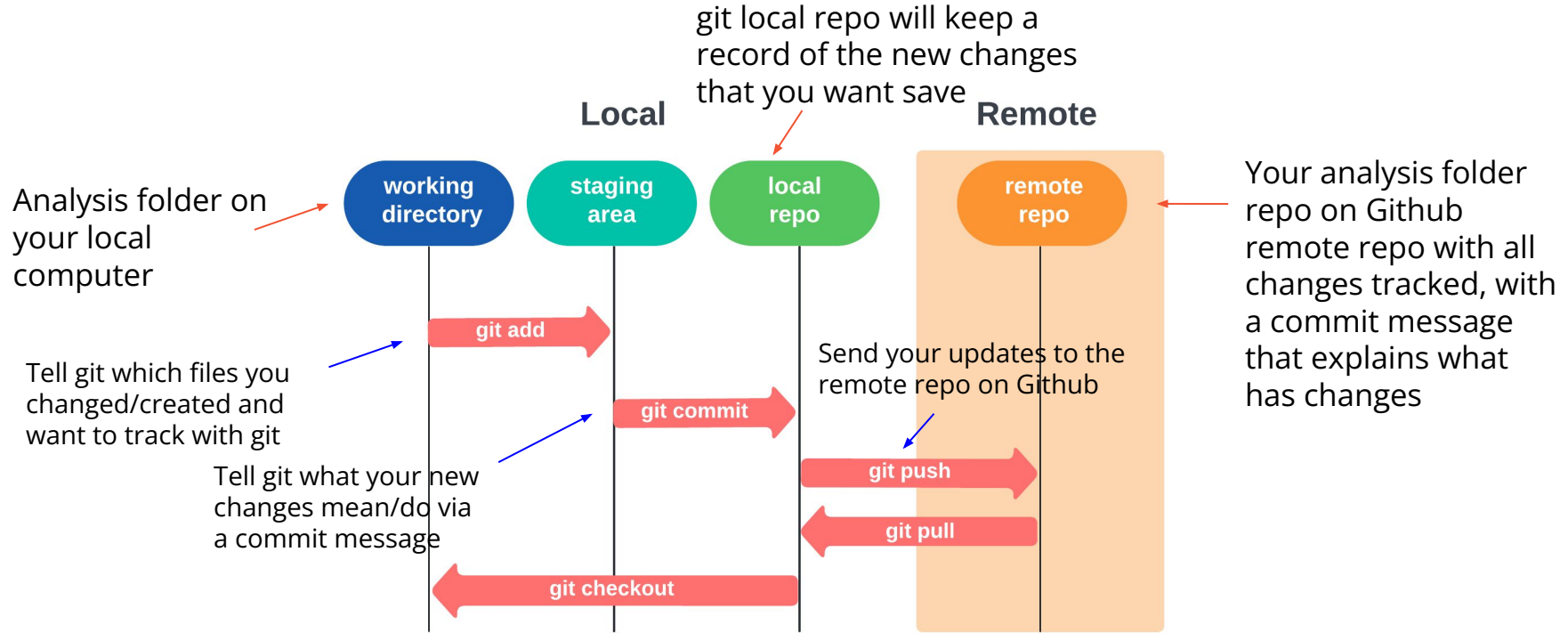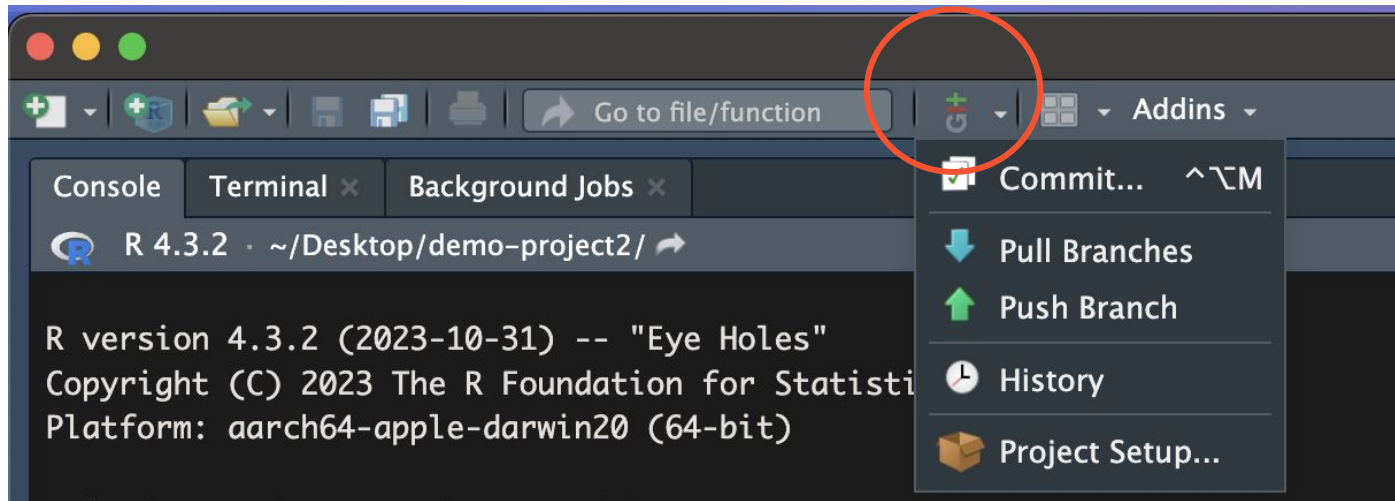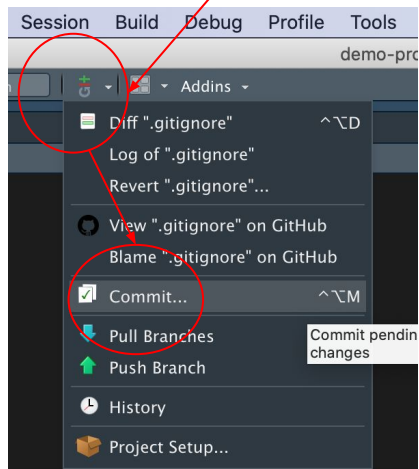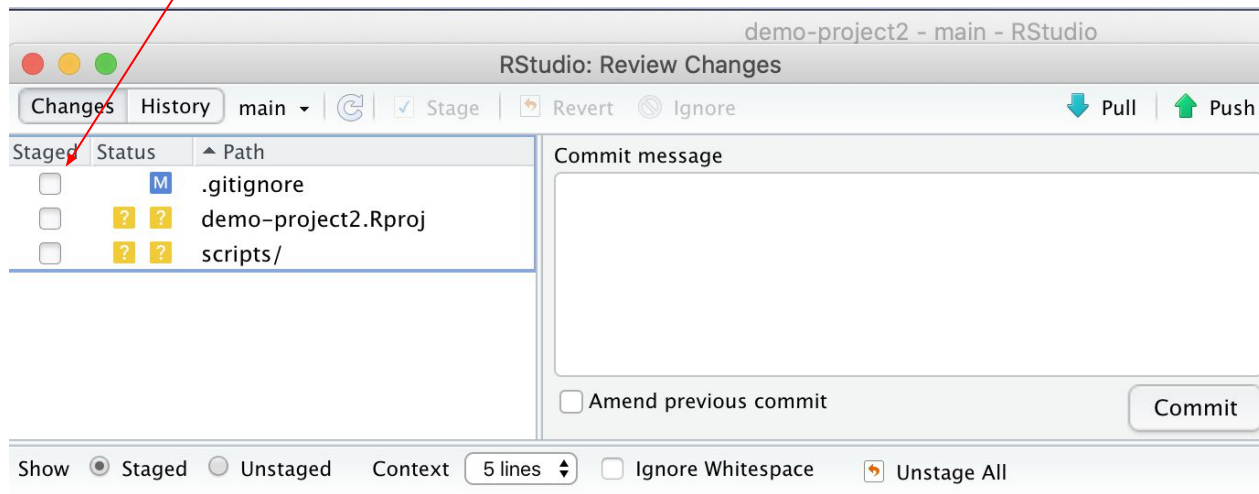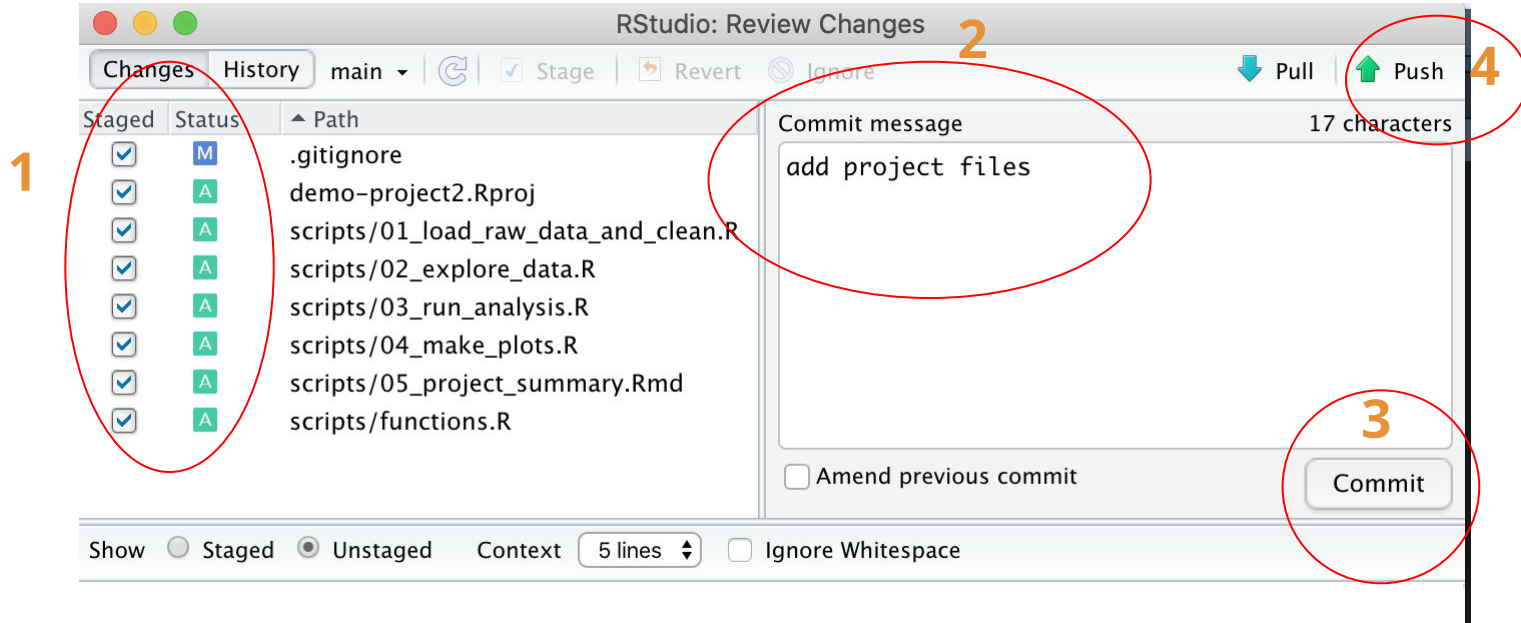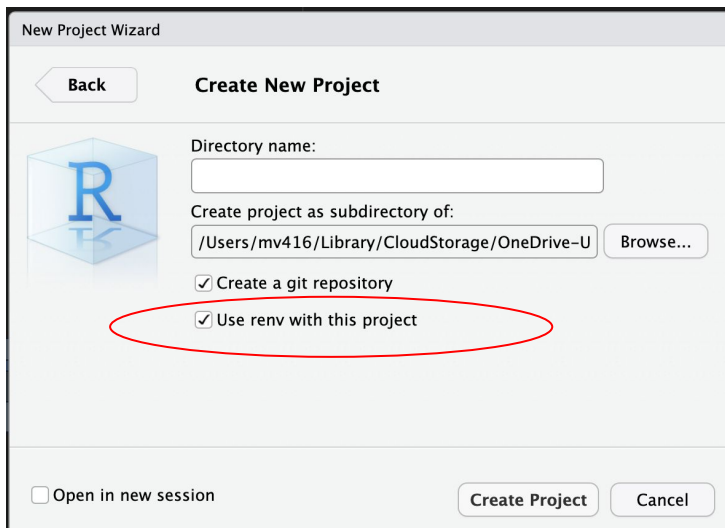