

Data science lifecycle & Exploratory data analysis using visualization

ds4owd - data science for openwashdata

Lars Schöbitz

2023-11-07

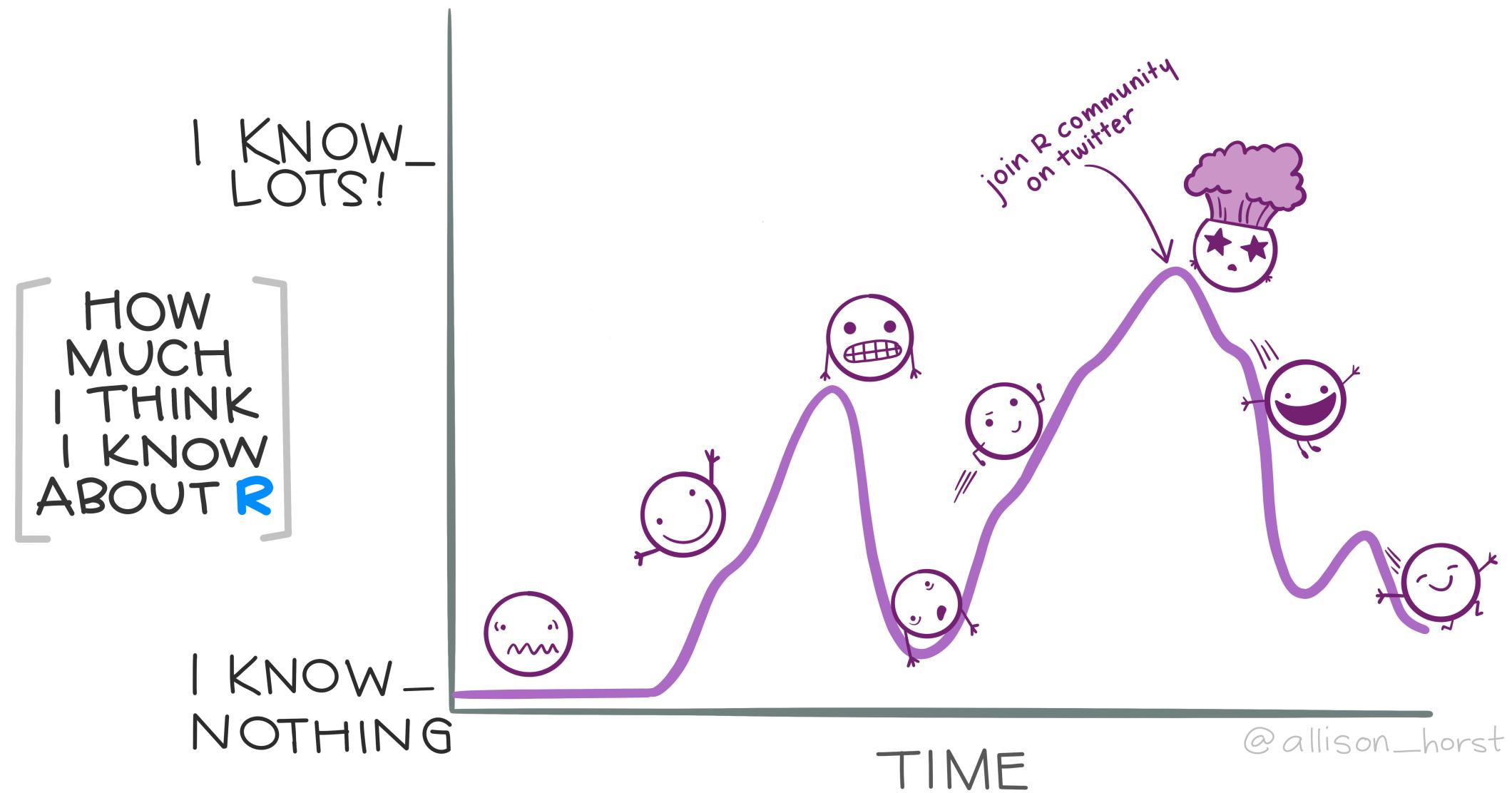


Q: How do I successfully complete the course?

You successfully complete the course and you will receive a certificate of completion if you:

hand in a complete capstone project report that uses a dataset of your choice by 30 January 2024 (instructions will follow)

This is the only requirement to successfully complete the course, independent of how many classes you attended or how many homework assignments you completed.



Solving coding problems

Tipps for search engines

- Use actionable verbs that describe what you want to do
- Be specific
- Add R to the search query
- Add the name of the R package name to the search query
- Scroll through the top 5 results (don't just pick the first)

Example: “How to remove a legend from a plot in R ggplot2”

Stack Overflow

What is it?

- The biggest support network for (coding) problems
- Can be intimidating at first
- Up-vote system

Workflow

- First, briefly read the question that was posted
- Then, read the answer marked as “correct”
- Then, read one or two more answers with high votes
- Then, check out the “Linked” posts
- Always give credit for the solution

Tipps for AI tools

- Use actionable verbs that describe what you want to do
- Be specific
- Add R to the search query
- Add the name of the R package name to the search query

Example: “How to remove a legend from a plot in R ggplot2”

Other sources for help

- Posit Community Forum:
<https://community.rstudio.com/>
- Documentation websites:
<https://ggplot2.tidyverse.org/>
- Mastodon tag: [#rstats](#)
- Quarto GitHub Discussion:
<https://github.com/quarto-dev/quarto-cli/discussions>





Homework assignment module 1

on GitHub Organisation

Bookmark this link in your browser!

github.com/ds4owd-001

ds4owd-001

github.com/ds4owd-001?q=rainbow-train&type=all&language=&sort=

website Public
Forked from cven5873-ss23/website
Website for data science for openwashdata course 001
JavaScript

Top languages
HTML JavaScript

Repositories

rainbow-train Type Language Sort New

1 result for all repositories matching rainbow-train sorted by last updated Clear filter

md-01-assignments-rainbow-train Private
0 stars 0 forks 2 issues 0 pull requests Updated 4 days ago

© 2023 GitHub, Inc. Terms Privacy Security Status Docs Contact GitHub Pricing API Training Blog About

on your repository

A screenshot of a GitHub repository page for 'ds4owd-001/md-01-assignments-rainbow-train'. The repository is private and was generated from 'ds4owd/md-01-assignments'. The page shows one branch ('main') and zero tags. The repository contains four files: '.gitignore', 'git-configuration.qmd', 'md-01-assignments.Rproj', and a README file. The README file content is: 'Help people interested in this repository understand what you are working on.' A large green box highlights the repository name 'md-01-assignments-rainbow-train' in the header. A pink arrow labeled 'Step 1' points to the 'Code' dropdown menu, which is open to show cloning options via 'HTTPS', 'SSH', or 'GitHub CLI'. Another pink arrow labeled 'Step 2' points to the 'Watch' button in the top right corner of the repository header.

on Posit Cloud

Bookmark this link in your browser!

posit.cloud/spaces/426916/content/

Posit Cloud Posit Cloud

posit.cloud/spaces/426916/content/all?sort=name_asc

ds4owd-001 Global Health Engineering

All Content (3)

Step 1 → New Project

TYPE * ACCESS * SORT + New Project from Template

md-01-assignments-rainbow-train

R RStudio Project RT Rainbow Train Private Created Nov 2, 2023 2:5 jupyter New Jupyter Project

Step 2 → New Project from Git Repository

md-01-exercises CONTINUE

R RStudio Project LS Lars Schöbitz Space members Created Oct 31, 2023 11:18 AM 1 derived project

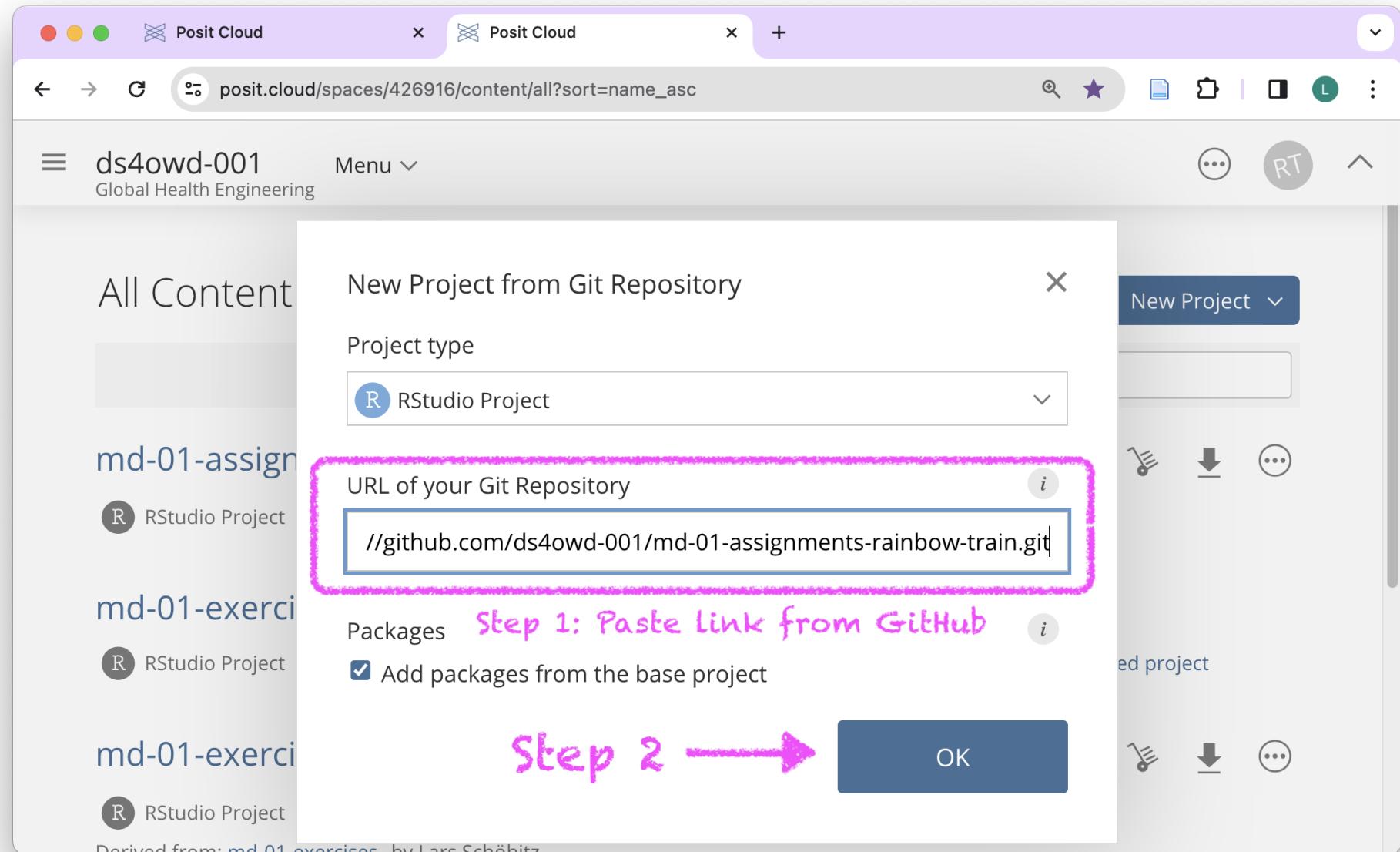
md-01-exercises

R RStudio Project RT Rainbow Train Private Created Oct 31, 2023 2:33 PM

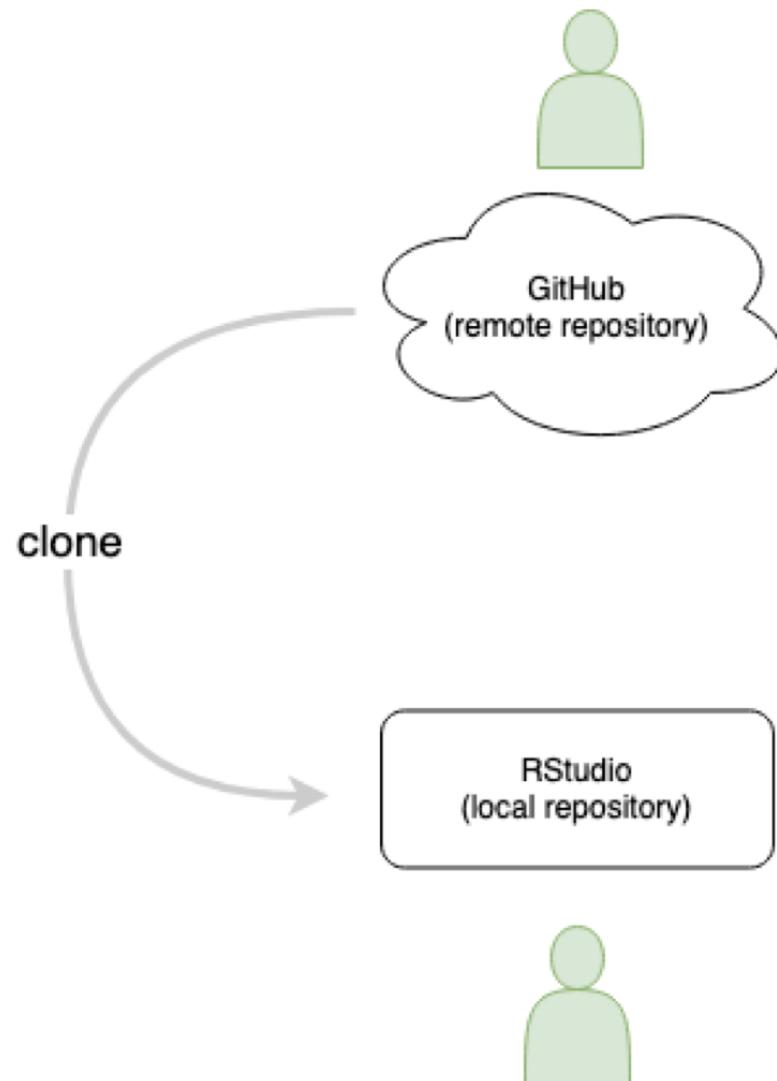
Derived from: md-01-exercises by Lars Schöbitz

on Posit Cloud

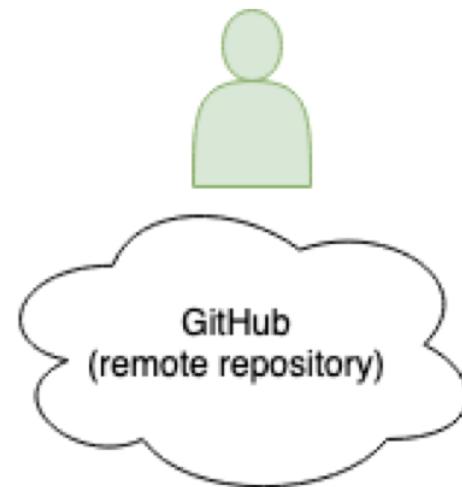
posit.cloud/spaces/426916/content/

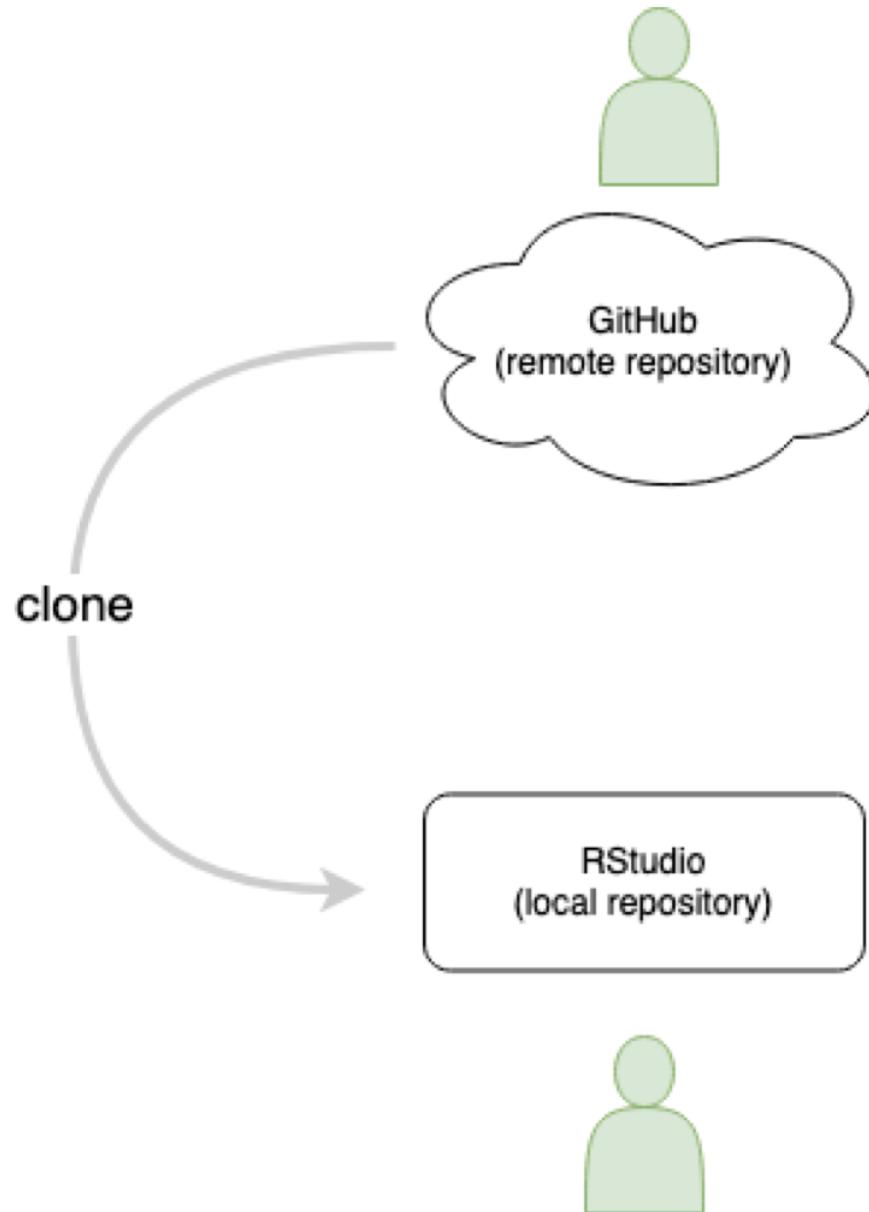


command: git clone



Version Control - Terminology

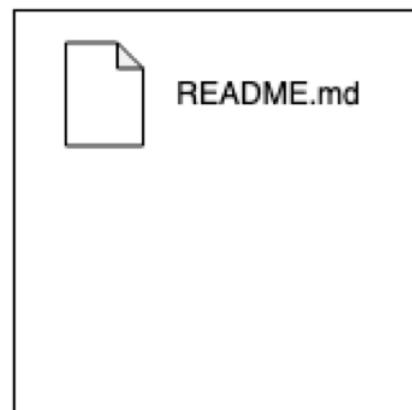






Create README.md

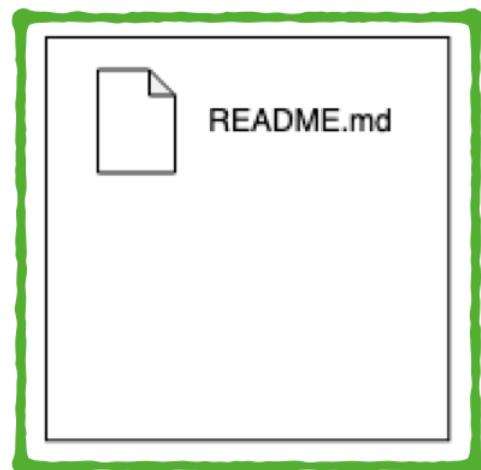
75aa637





Create README.md

75aa637



Repo(sitory)

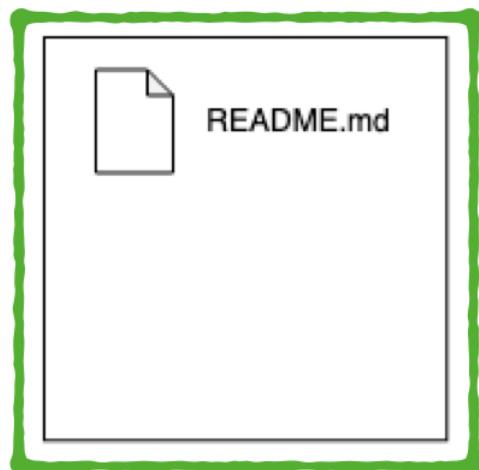


Commit
message

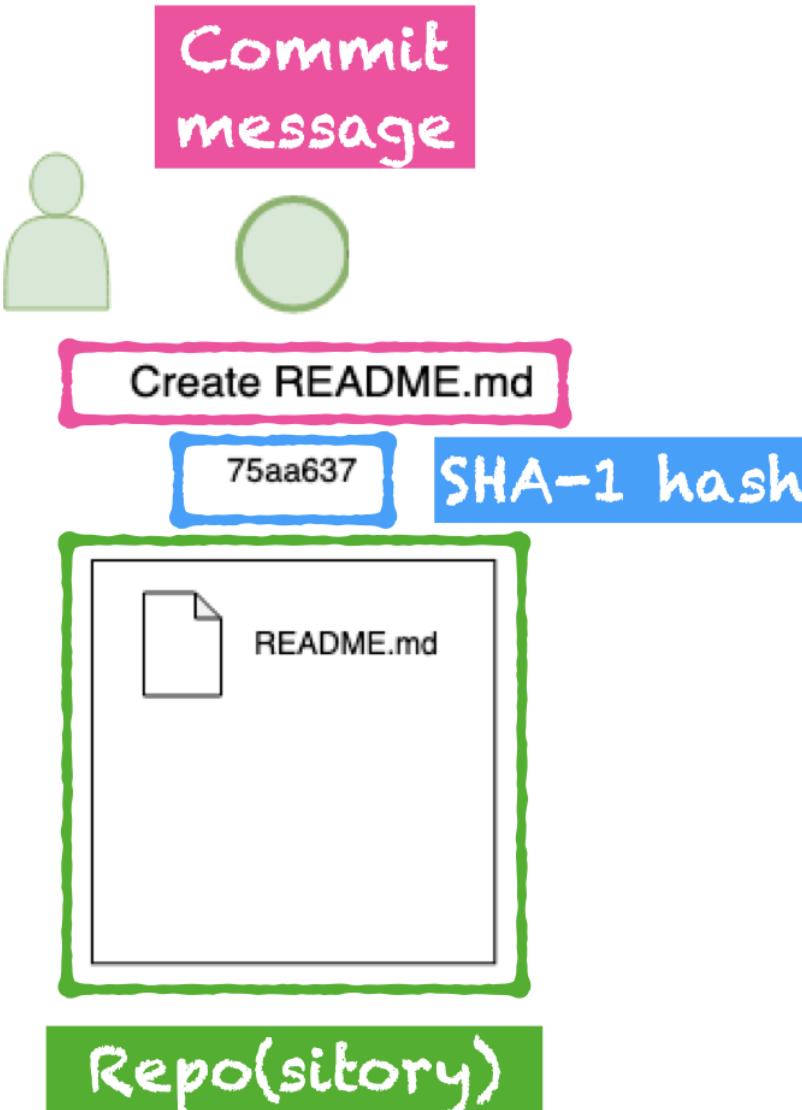


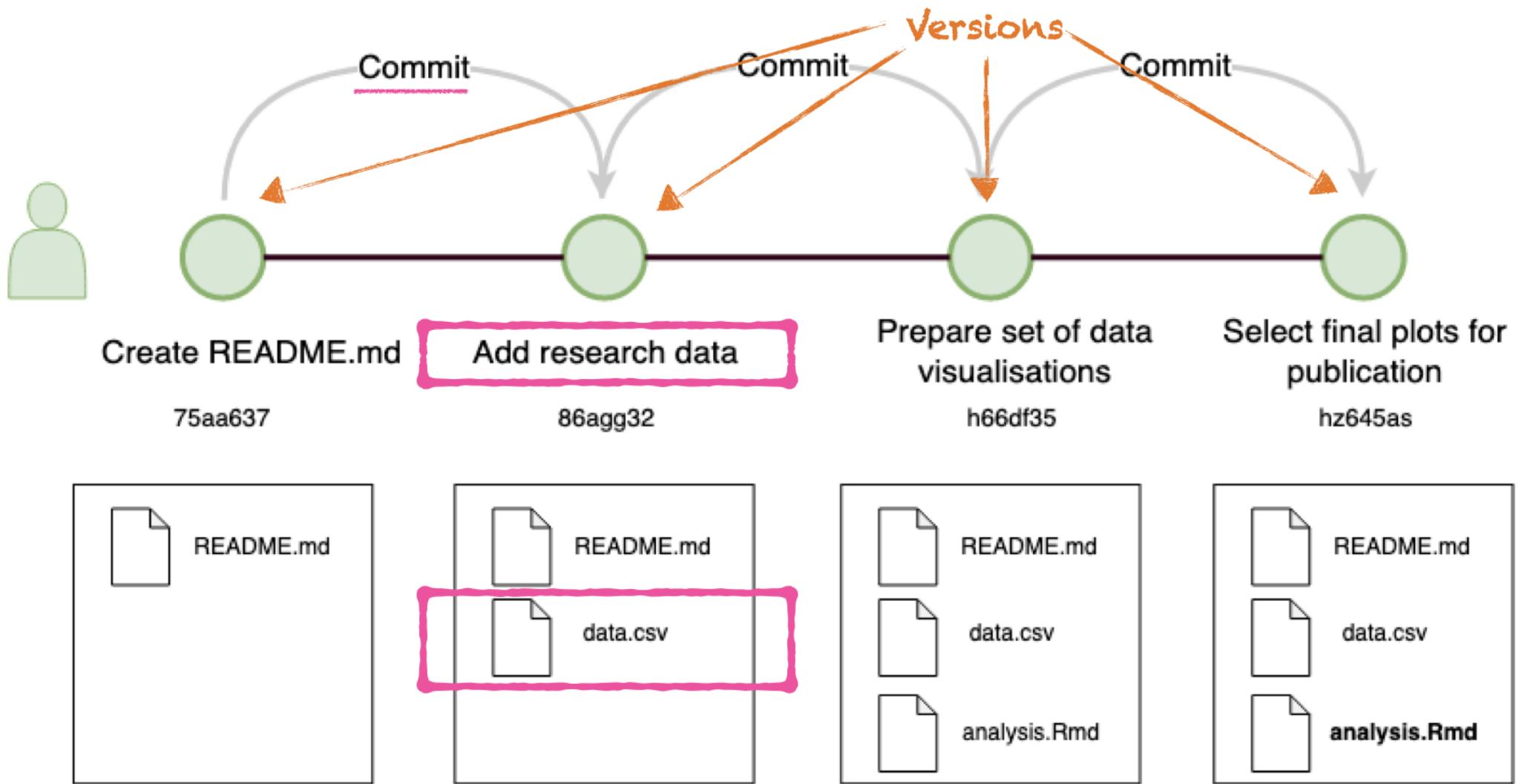
Create README.md

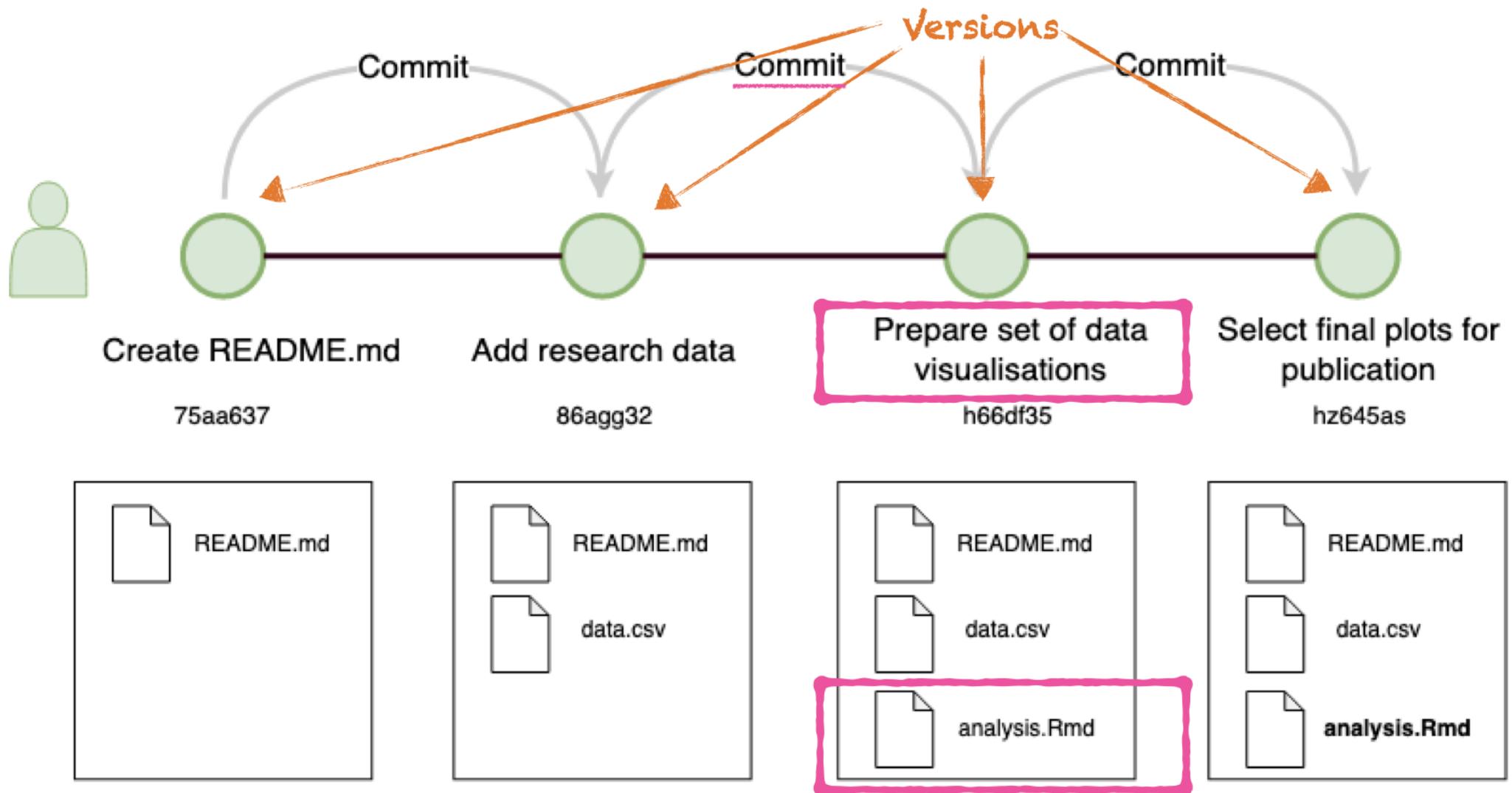
75aa637

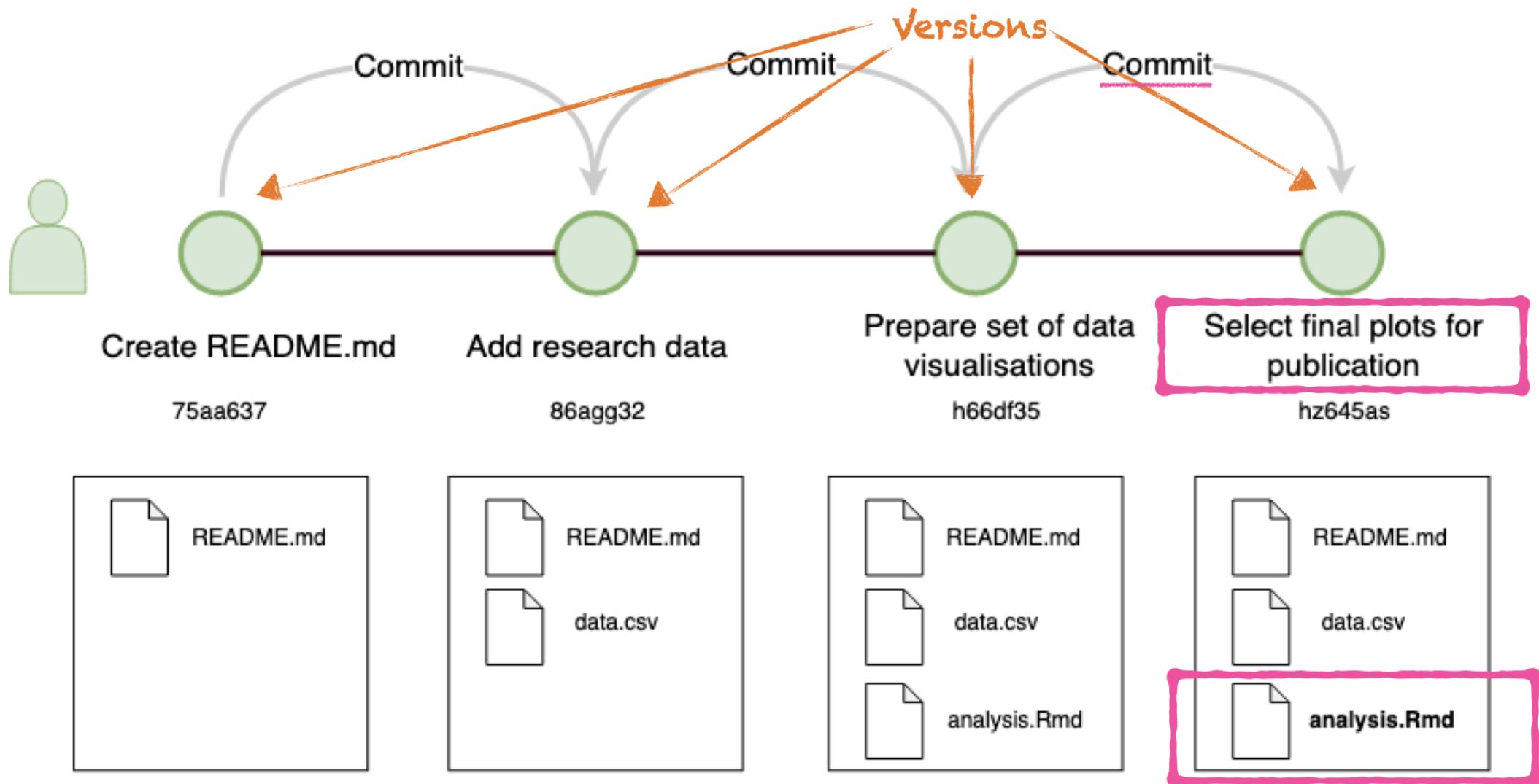


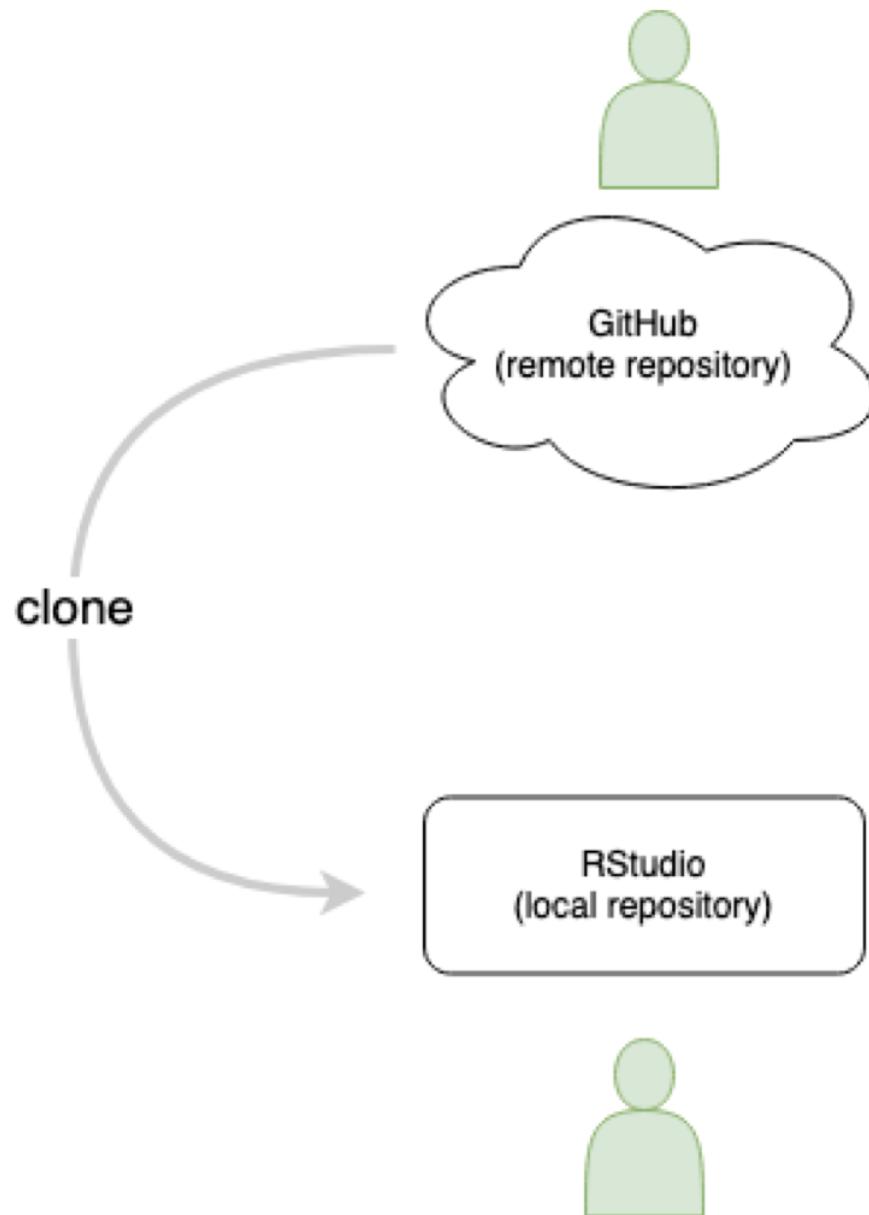
Repo(sitory)

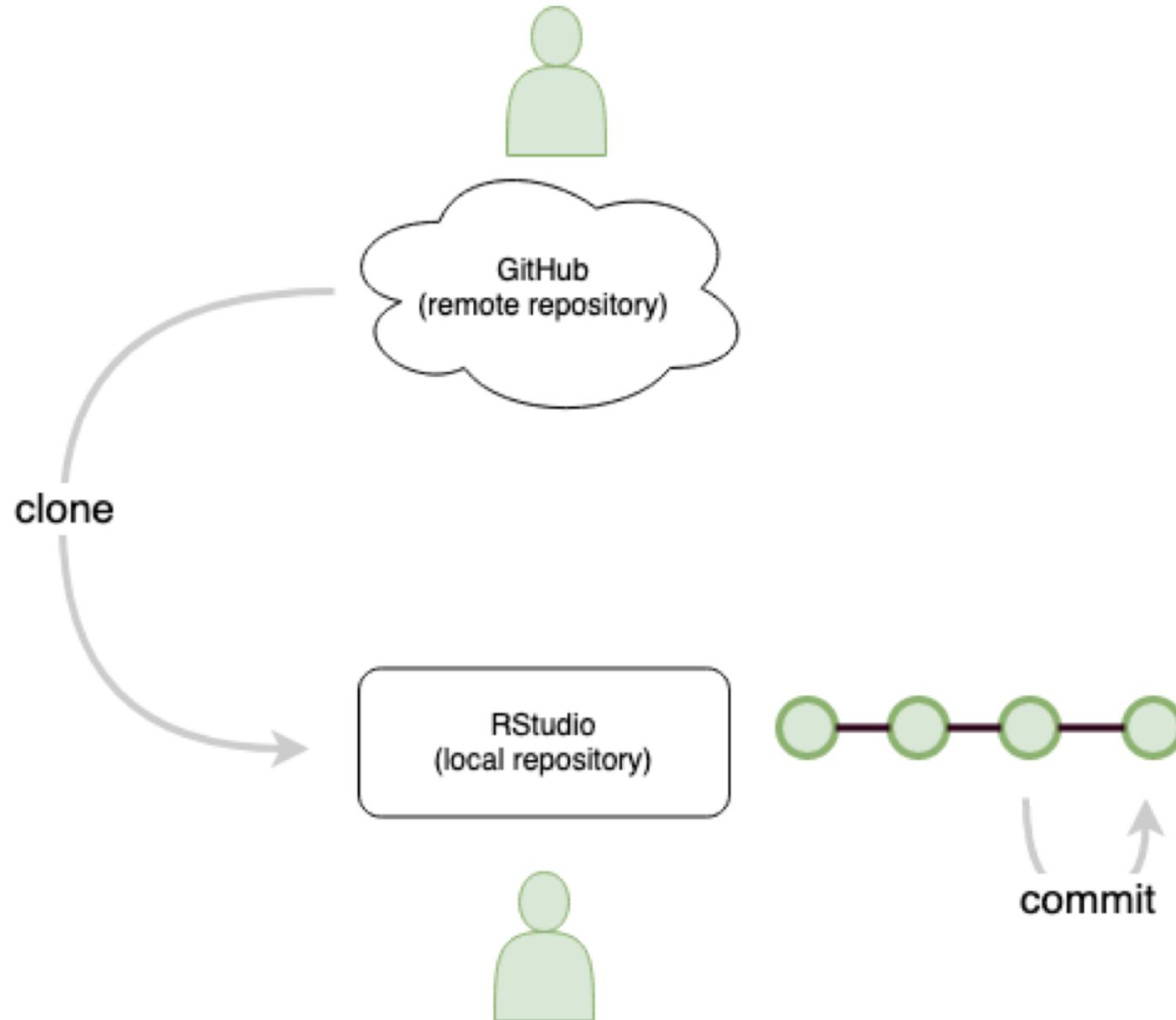


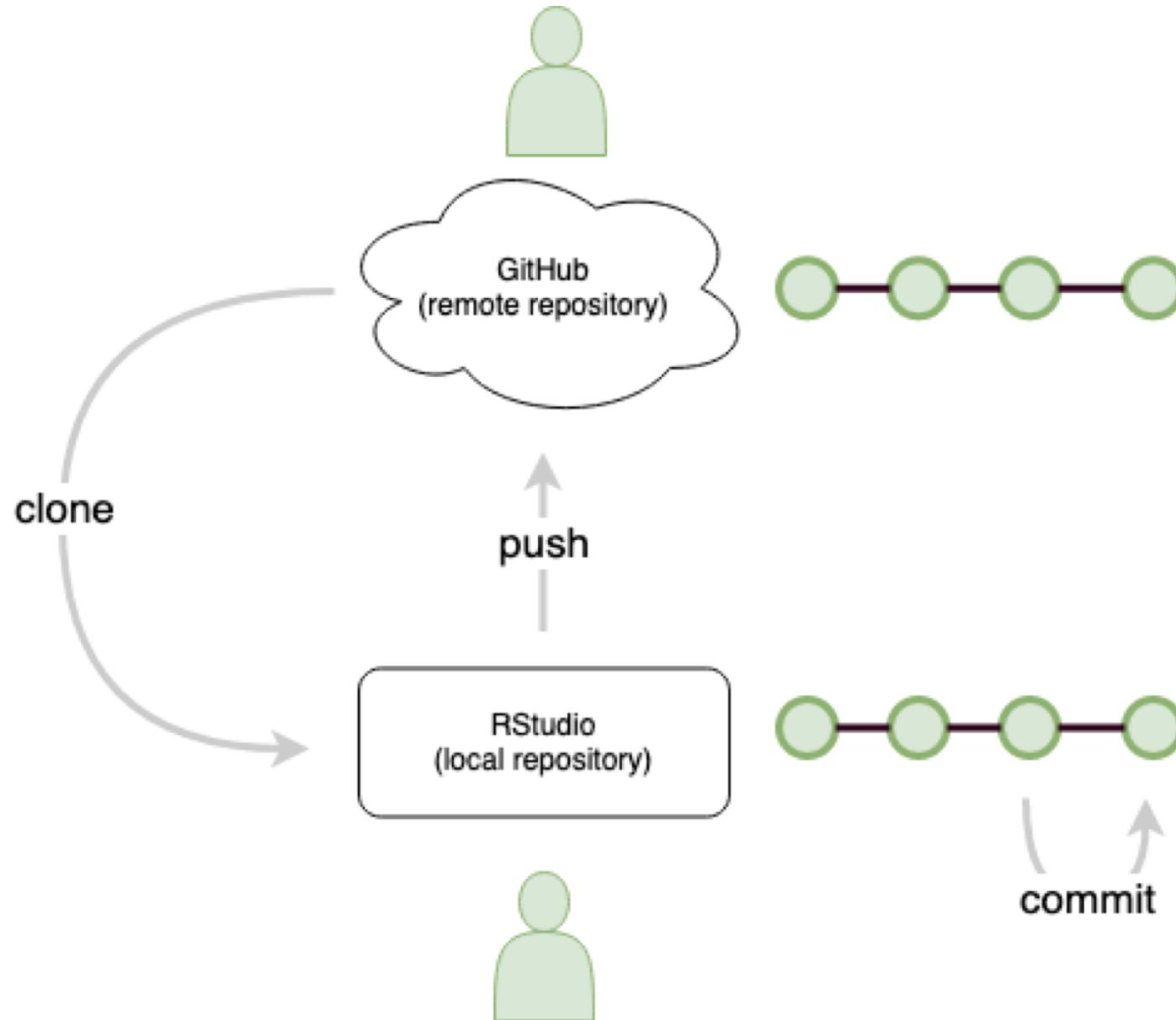




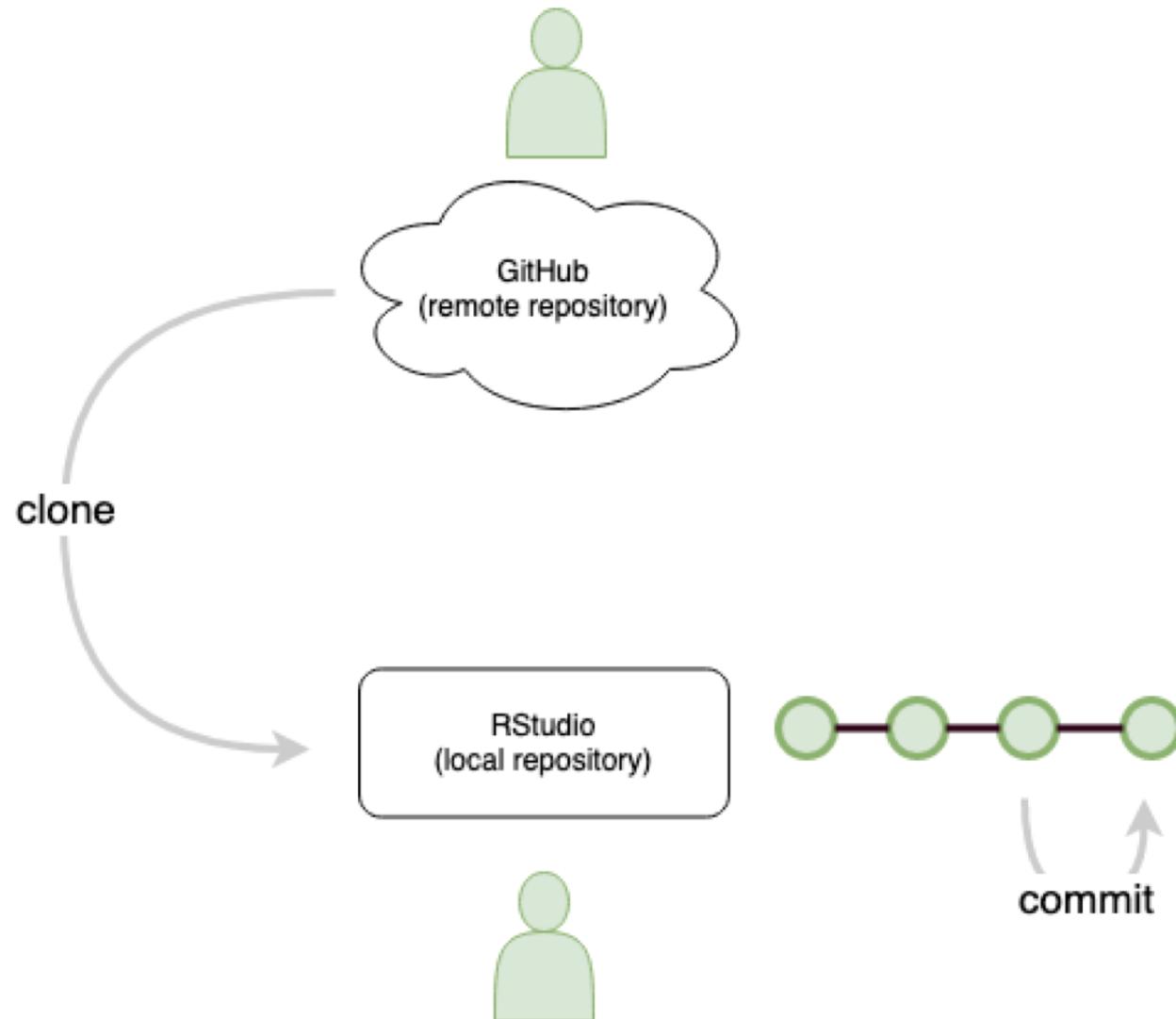




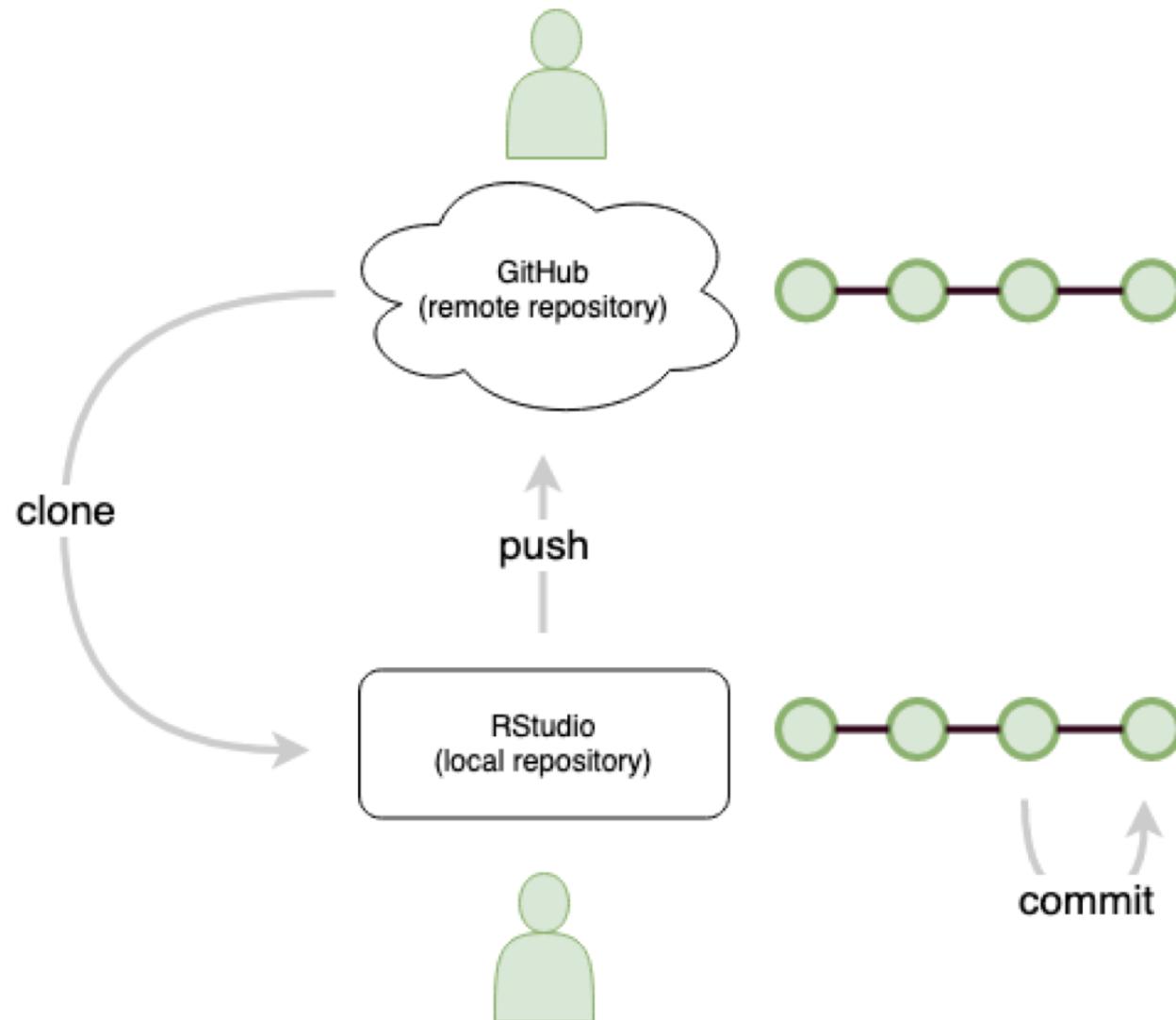




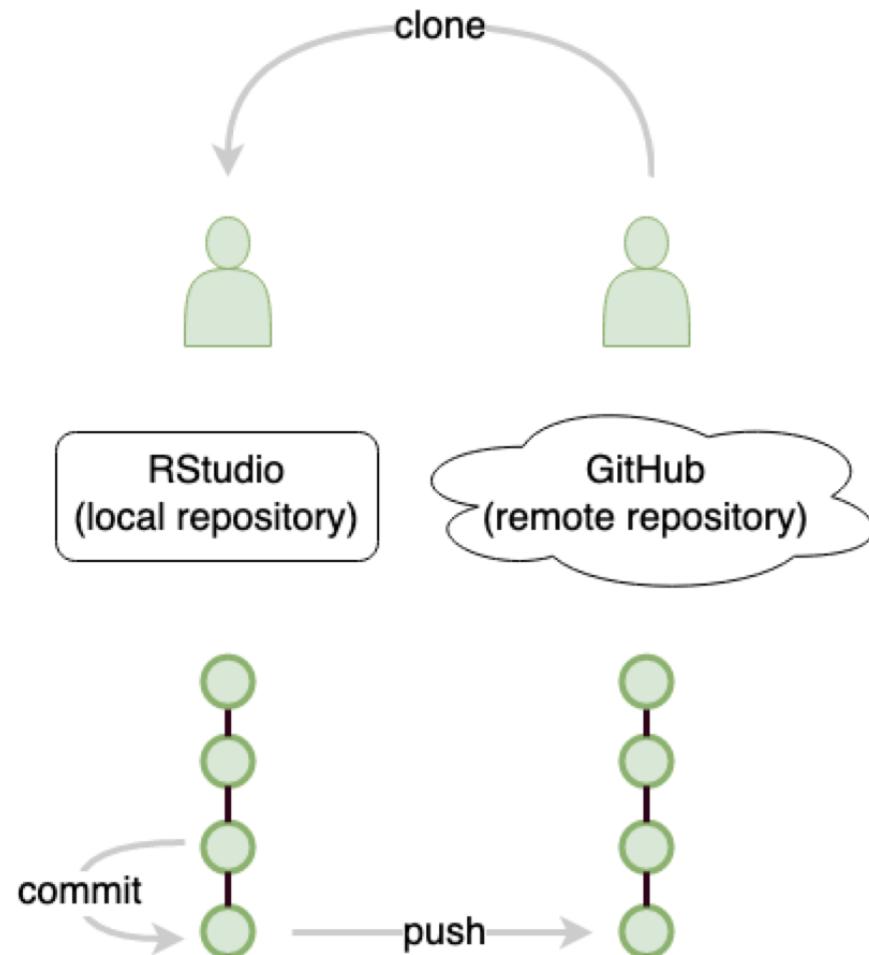
remember: git commit



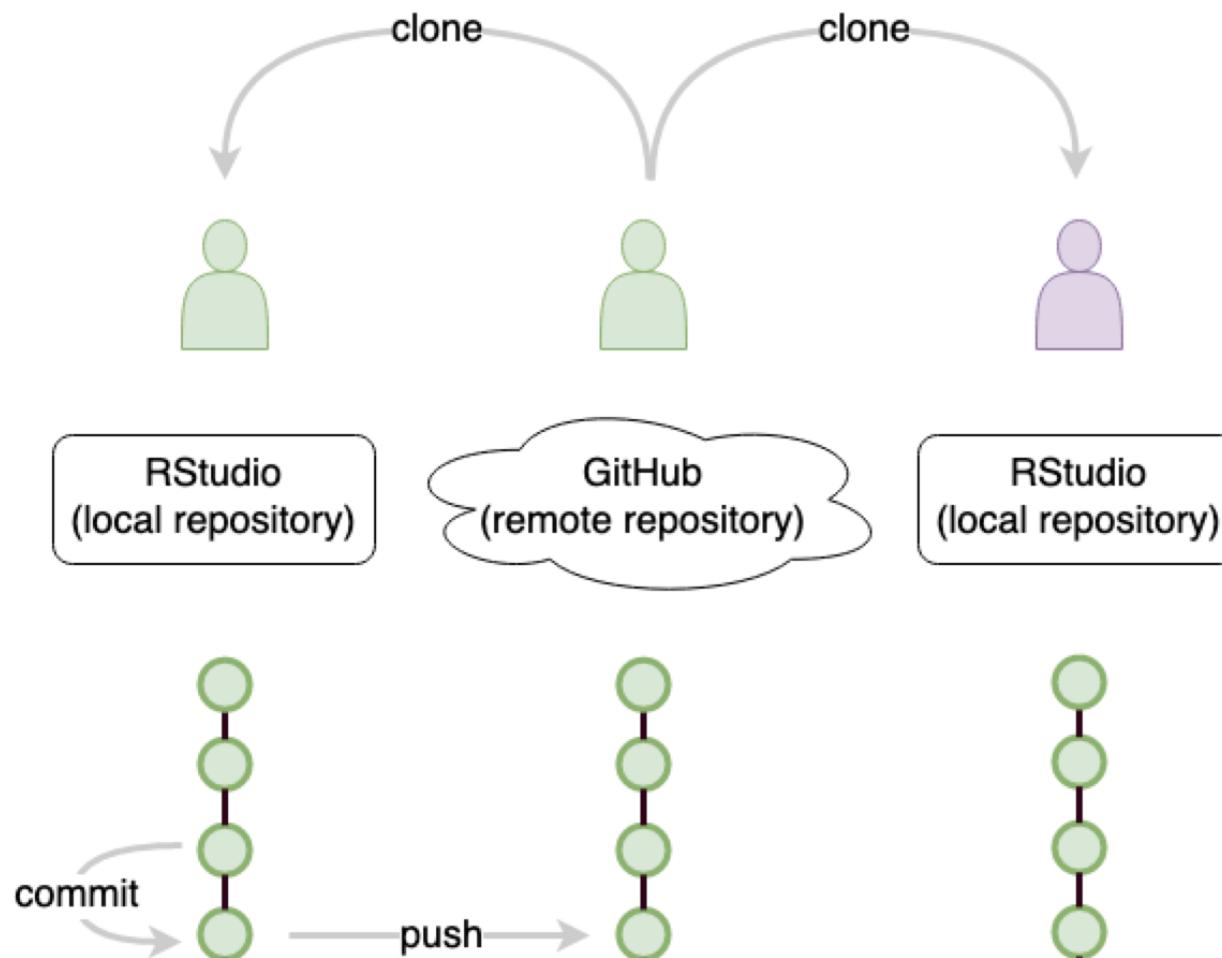
remember: git push



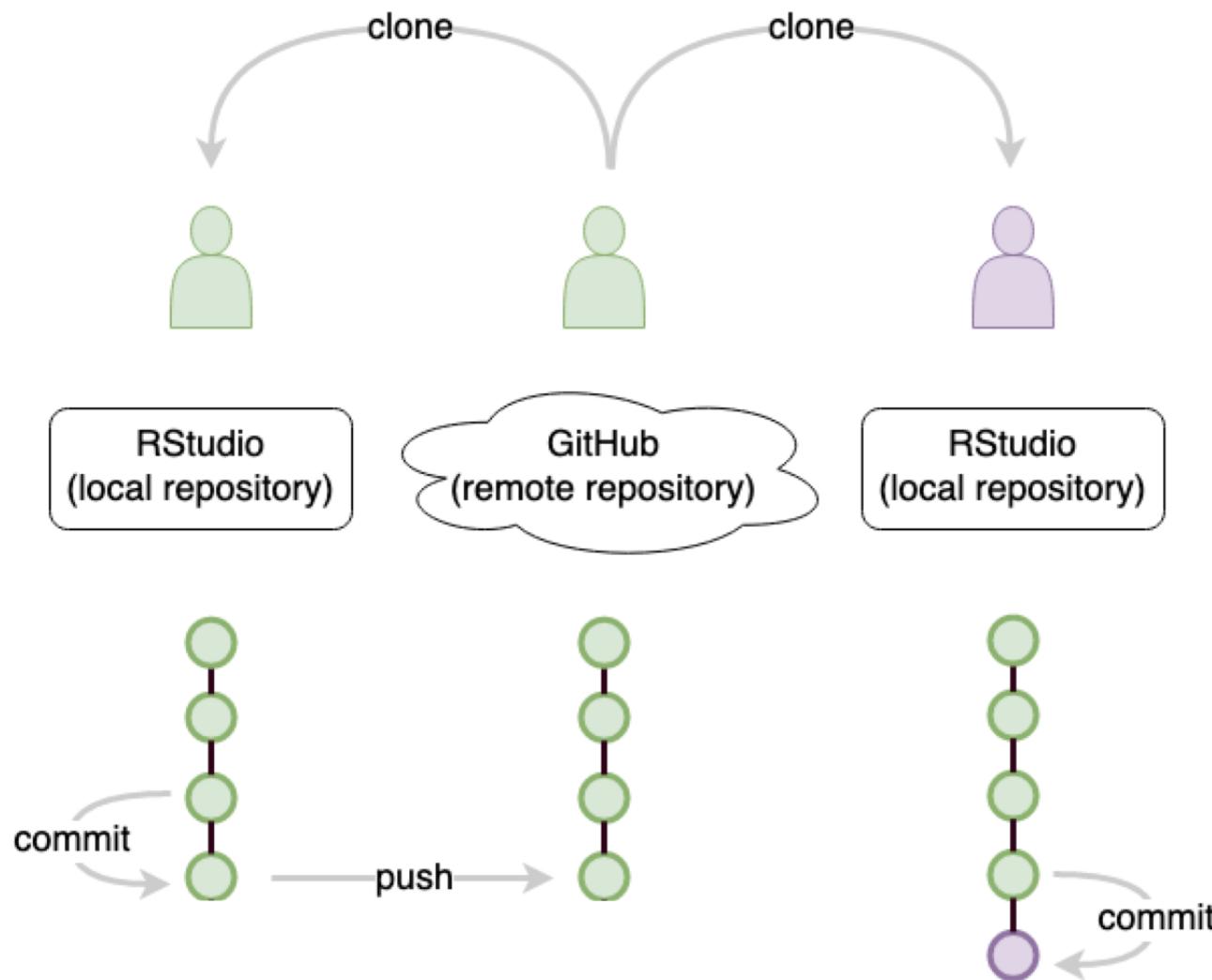
remember: git push



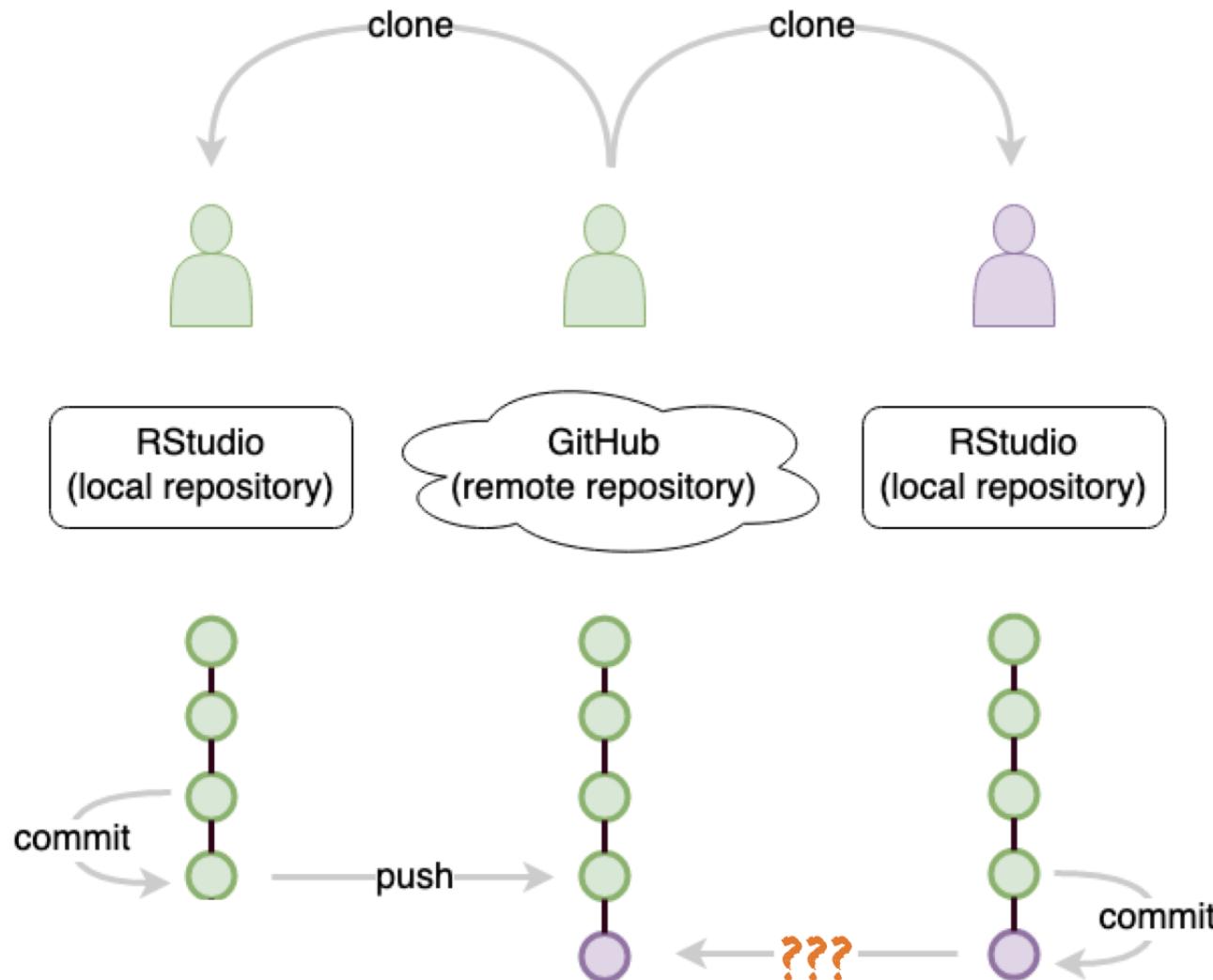
collaborate: git clone



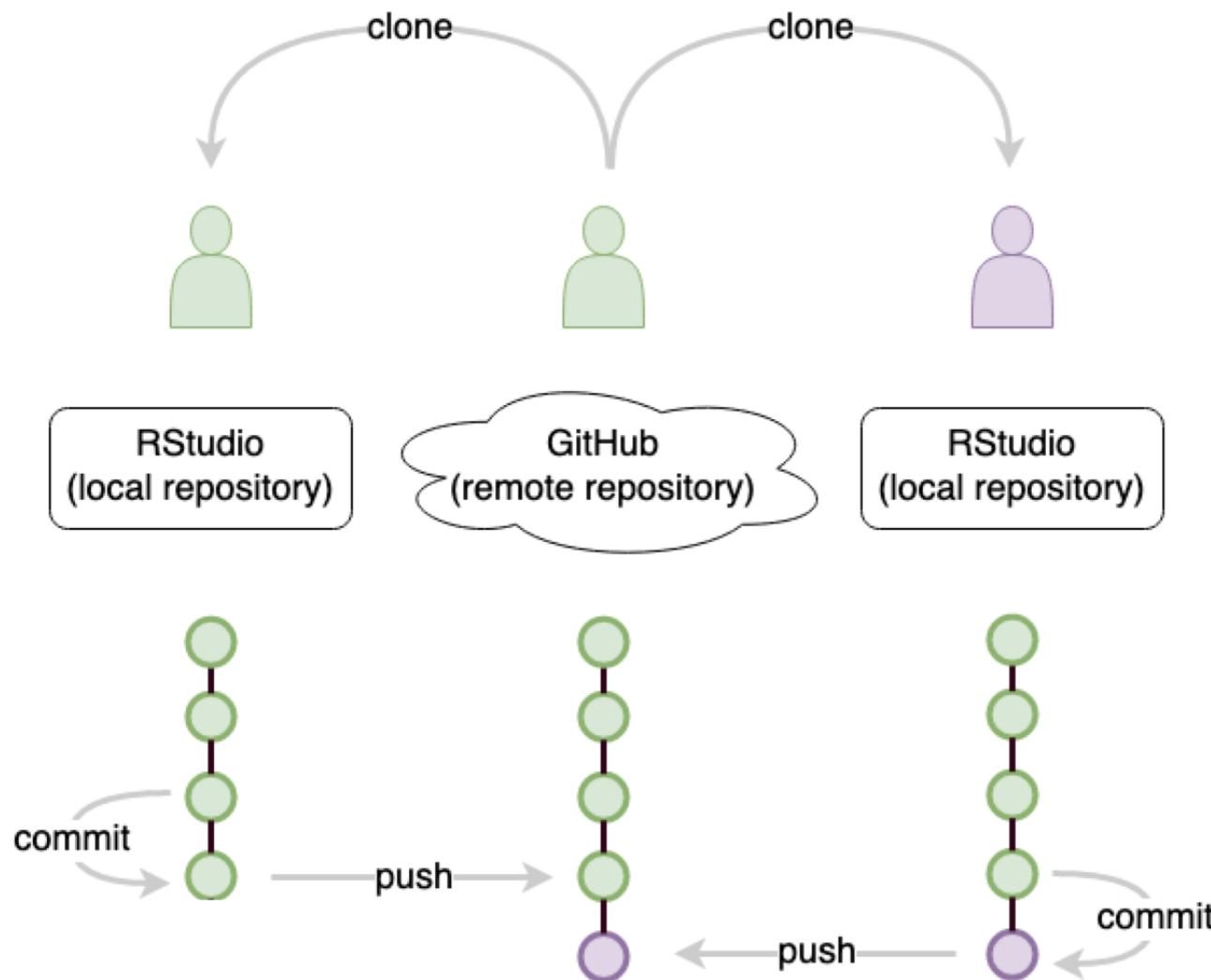
track work: git commit



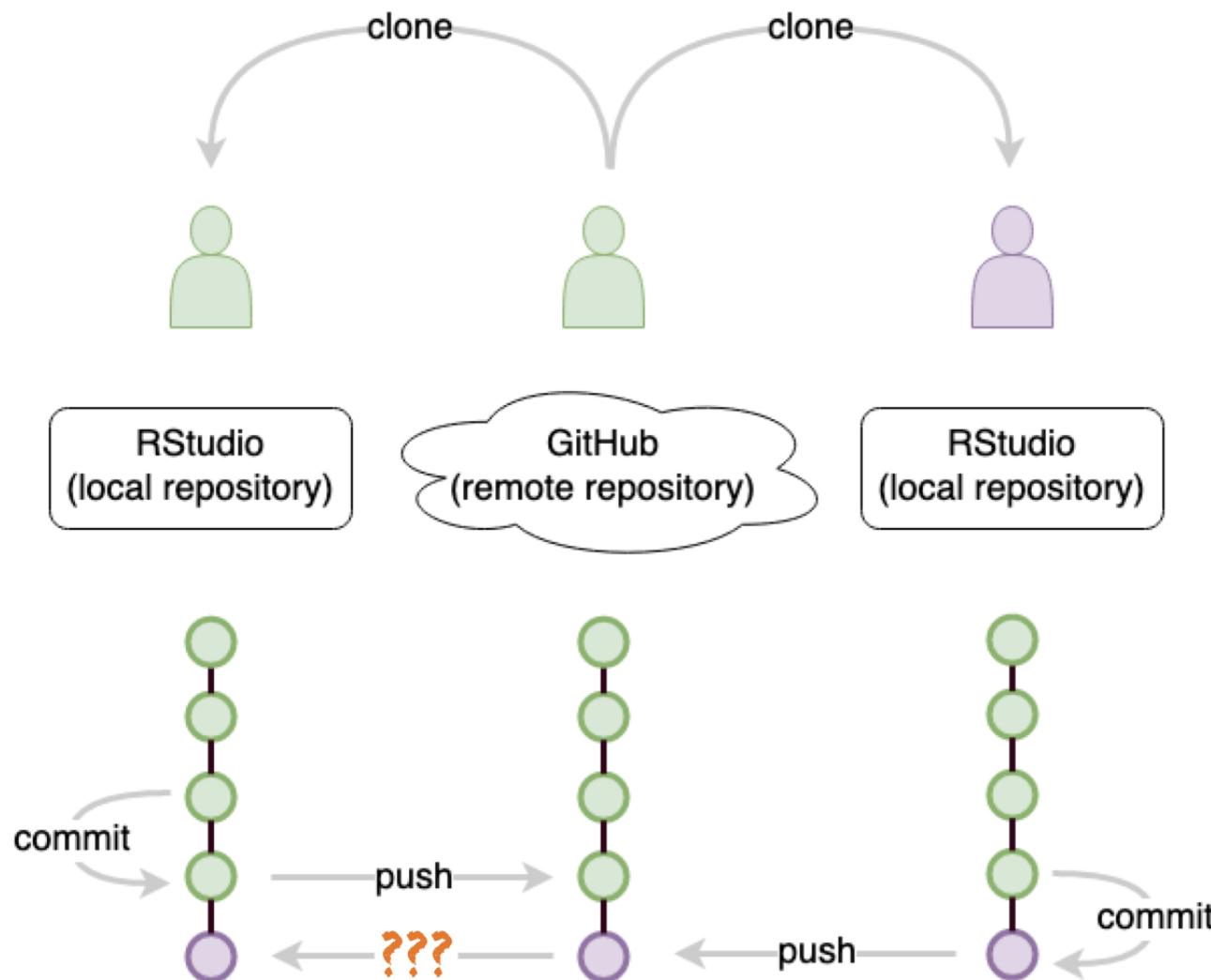
update: git ???



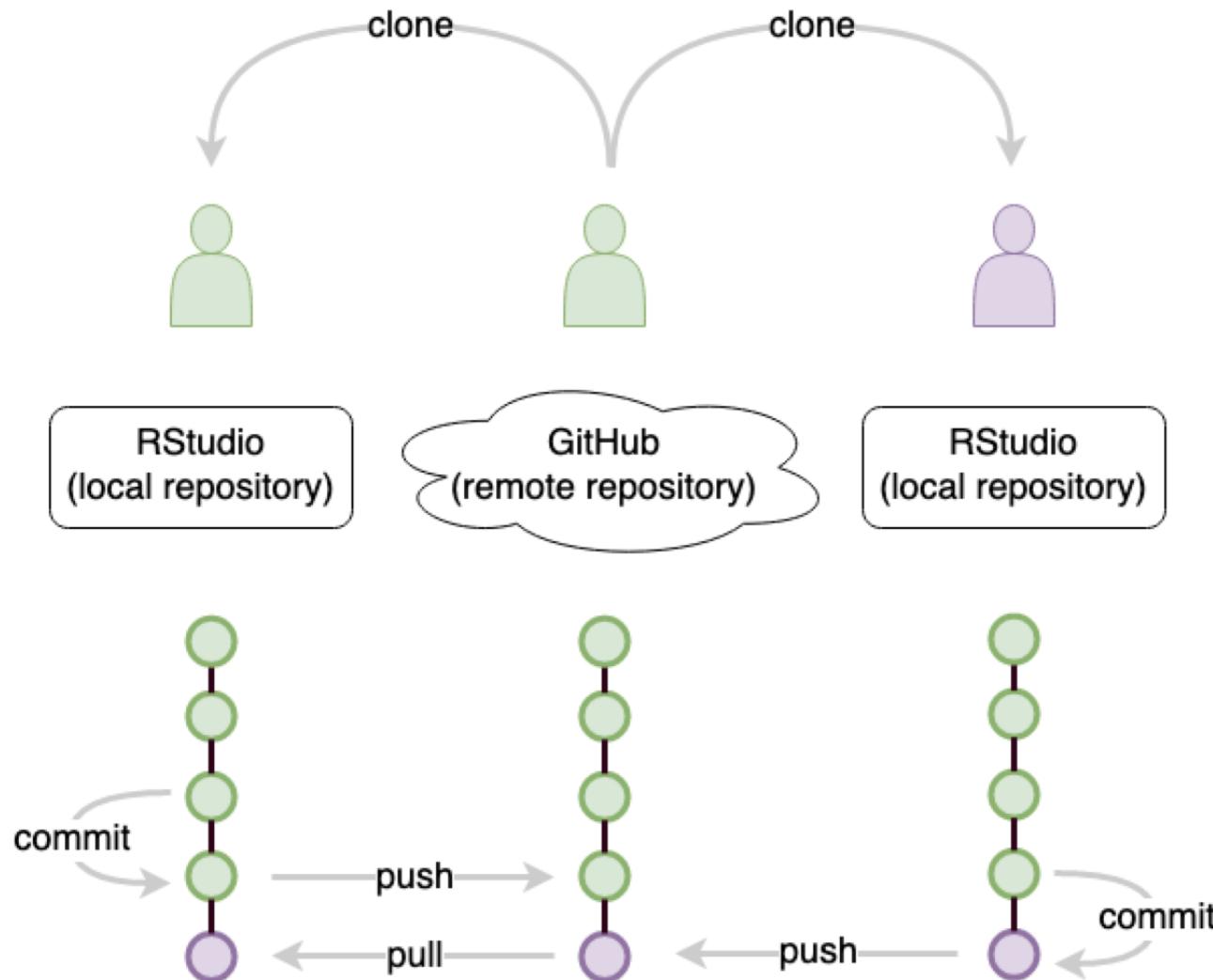
update: git push



git ???



new: git pull



Learning Objectives (for this week)

1. Learners can list the six elements of the data science lifecycle.
2. Learners can describe the four main aesthetic mappings that can be used to visualise data using the ggplot2 R Package.
3. Learners can control the colour scaling applied to a plot using colour as an aesthetic mapping.
4. Learners can compare three different geoms (bar/col, histogram, point) and their use case.

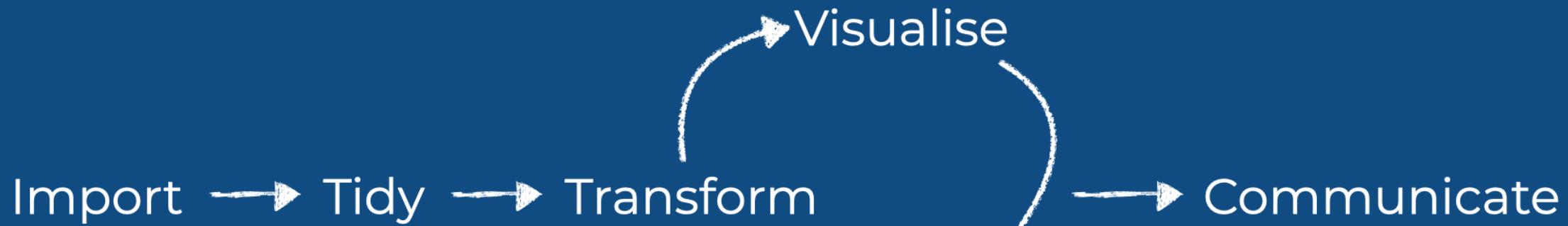
Data Science Lifecycle

Deep End

[via GIPHY](#)

 ds4owd-001.github.io/website/

Data Science Lifecycle

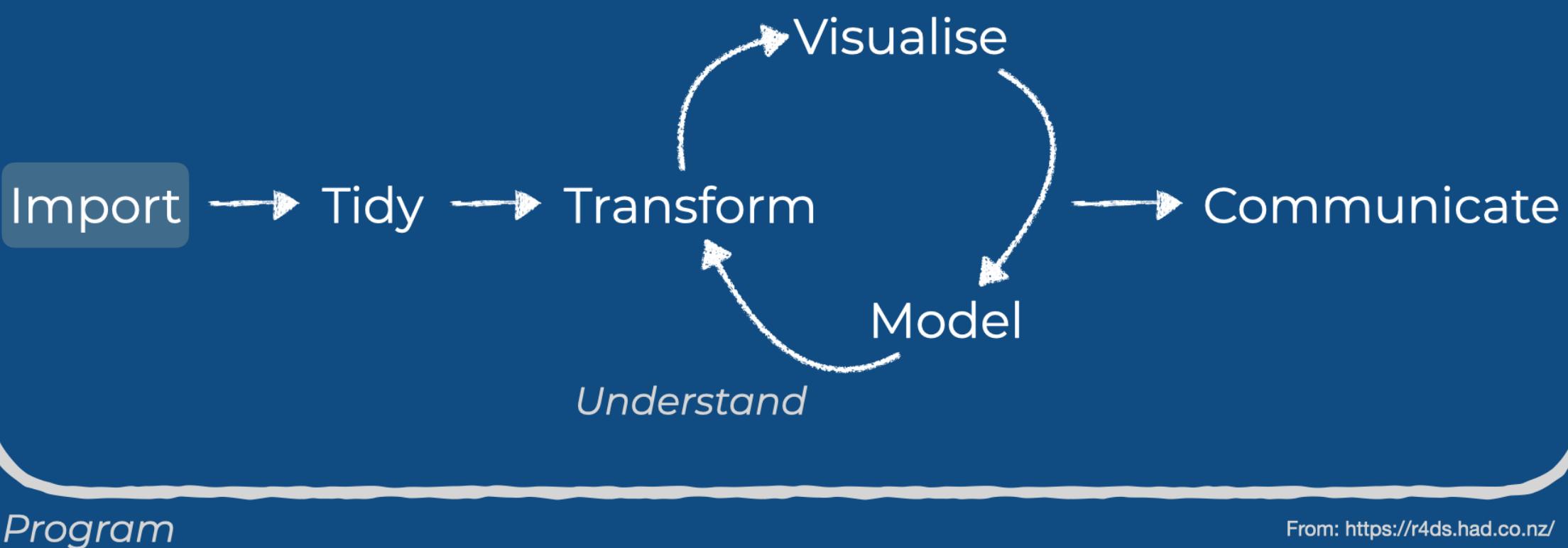


Program

From: <https://r4ds.had.co.nz/>

Data Science Lifecycle

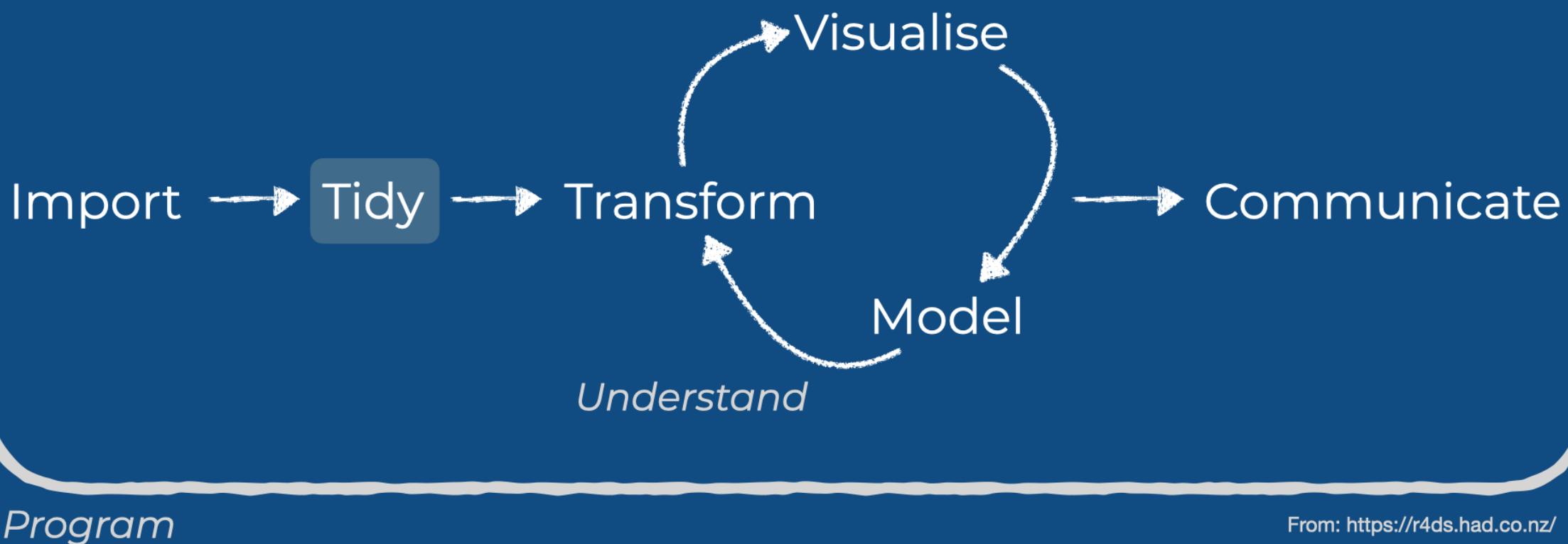
Get your data into R



From: <https://r4ds.had.co.nz/>

Data Science Lifecycle

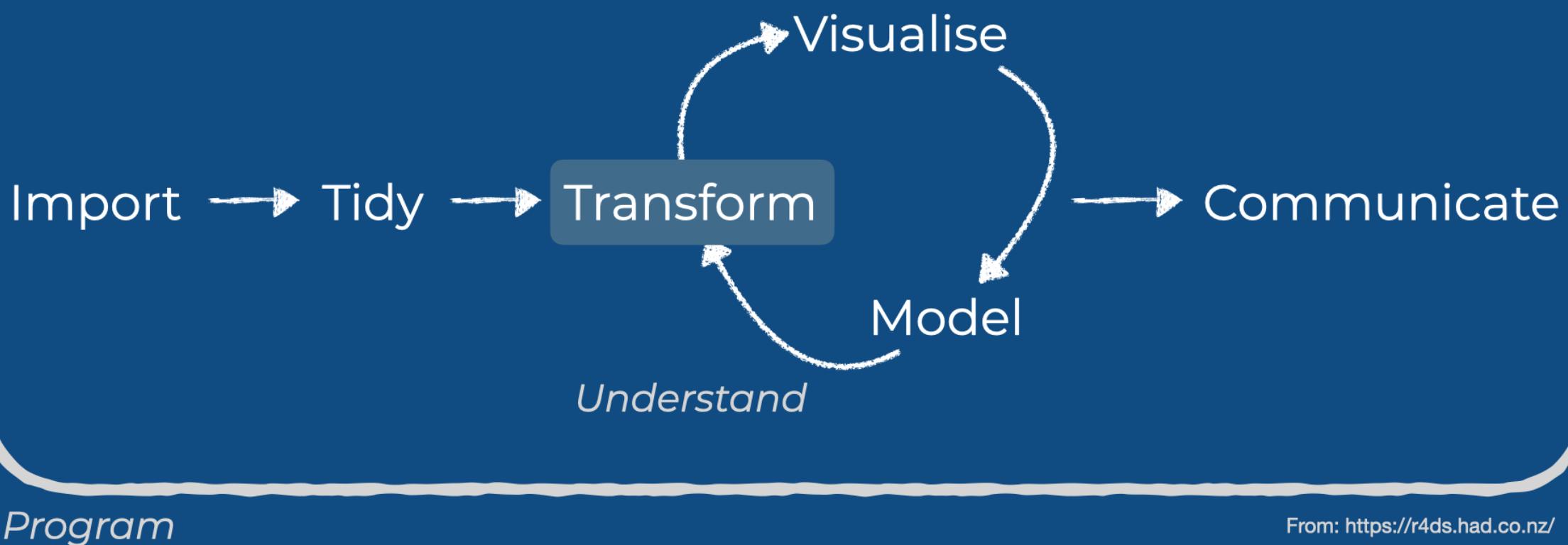
Store your data in a consistent form



From: <https://r4ds.had.co.nz/>

Data Science Lifecycle

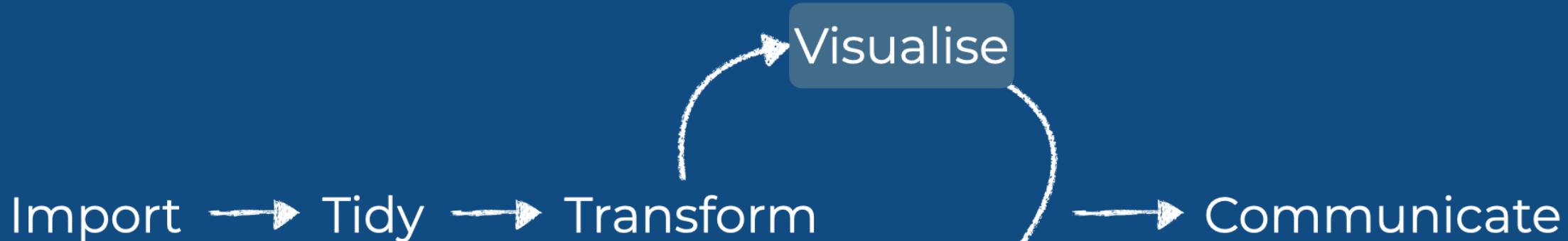
Narrow down + Create new variables + Summary stats



From: <https://r4ds.had.co.nz/>

Data Science Lifecycle

Explore your with visual representations

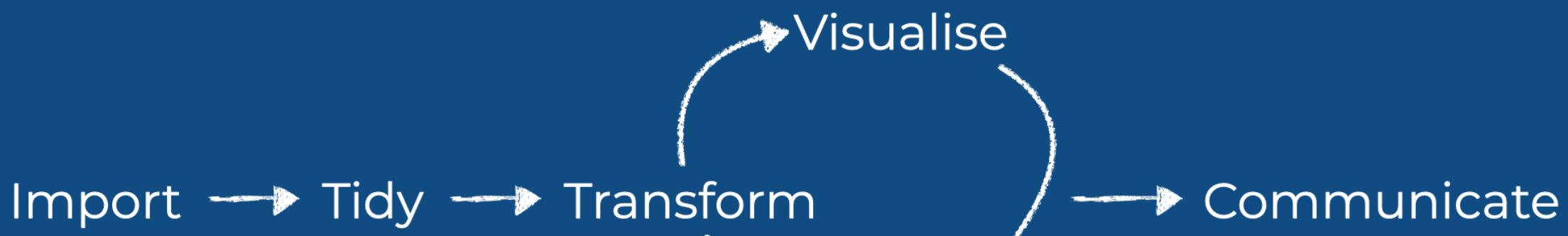


Program

From: <https://r4ds.had.co.nz/>

Data Science Lifecycle

Explore your with visual representations

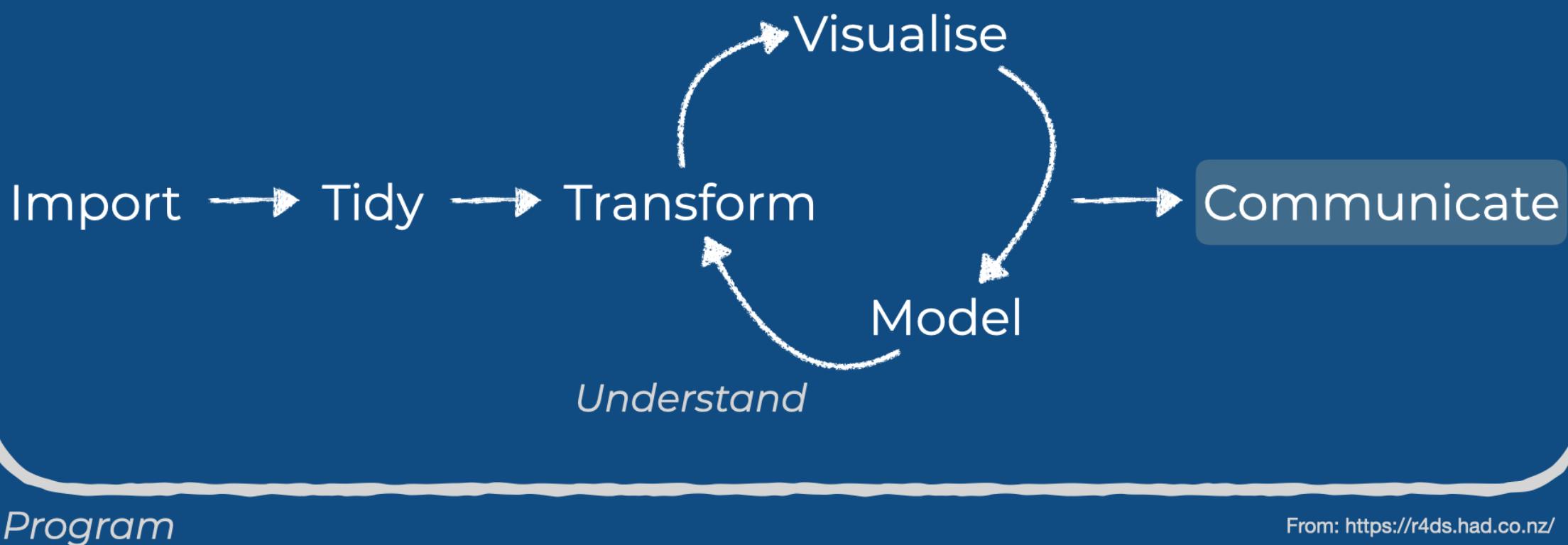


Program

From: <https://r4ds.had.co.nz/>

Data Science Lifecycle

Share your findings with others

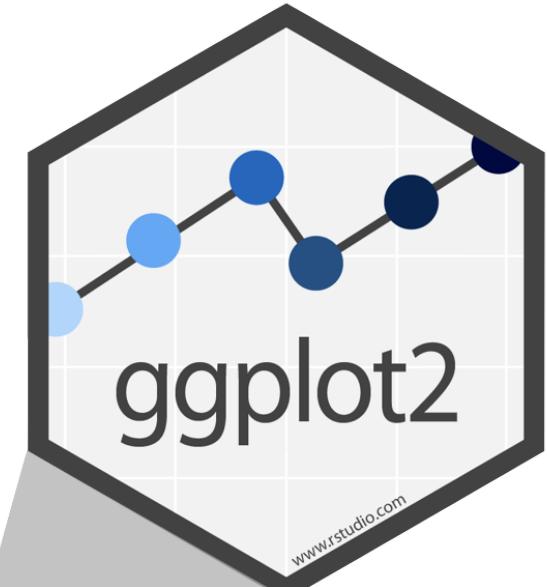


From: <https://r4ds.had.co.nz/>

Exploratory Data Analysis with `ggplot2`

R Package ggplot2

- **ggplot2** is tidyverse's data visualization package
- **gg** in ggplot2 stands for Grammar of Graphics
- Inspired by the book **Grammar of Graphics** by Leland Wilkinson
- **Documentation:**
<https://ggplot2.tidyverse.org/>
- **Book:** <https://ggplot2-book.org>



My turn: Working with R

Sit back and enjoy!

Code structure

- `ggplot()` is the main function in `ggplot2`
- Plots are constructed in layers
- Structure of the code for plots can be summarized as

```
1 ggplot(data = [dataset],  
2         mapping = aes(x = [x-variable],  
3                             y = [y-variable])) +  
4     geom_xxx() +  
5     other options
```

Code structure

```
1 ggplot()
```

Code structure

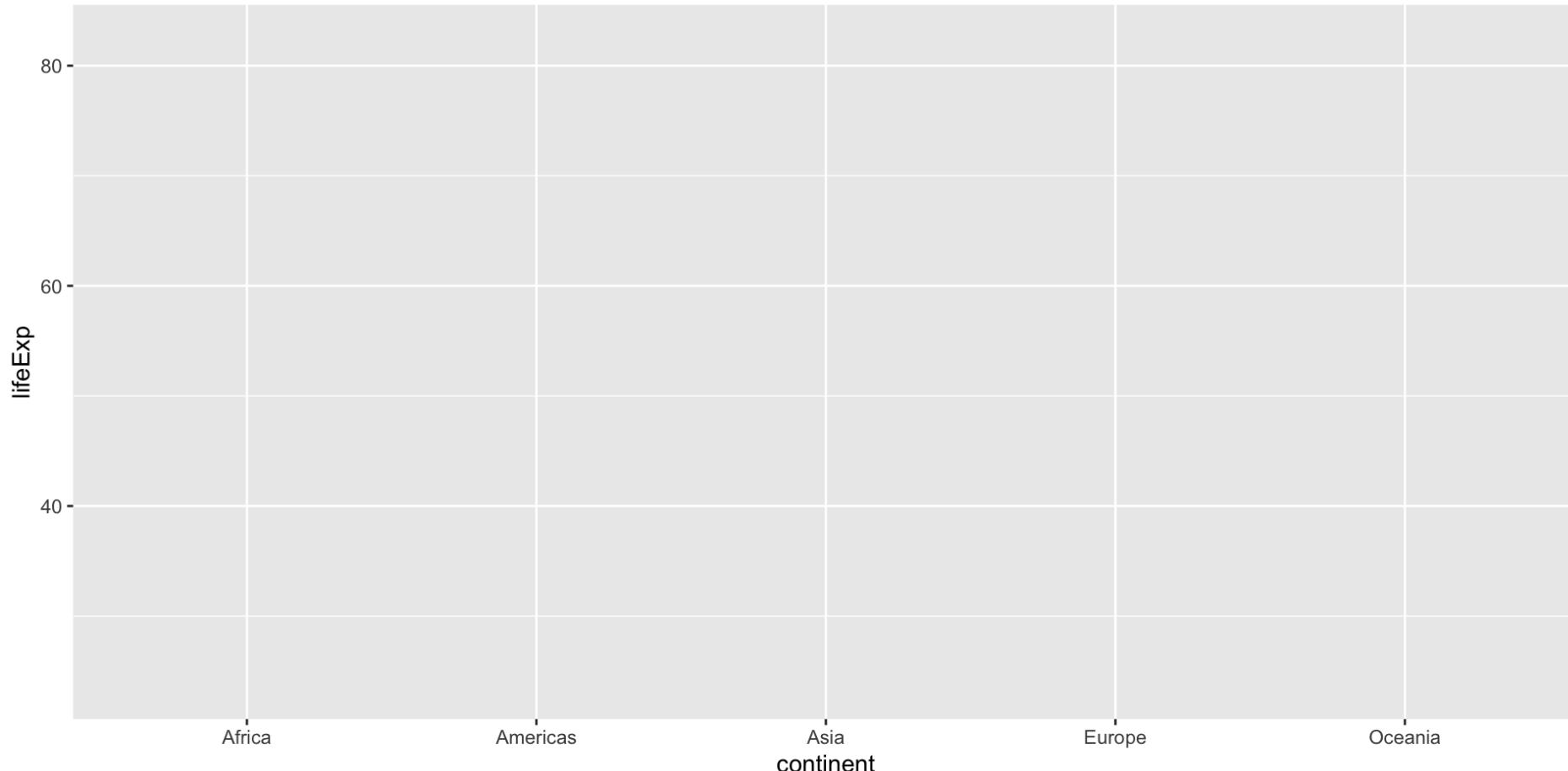
```
1 ggplot(data = gapminder)
```

Code structure

```
1 ggplot(data = gapminder,  
2         mapping = aes())
```

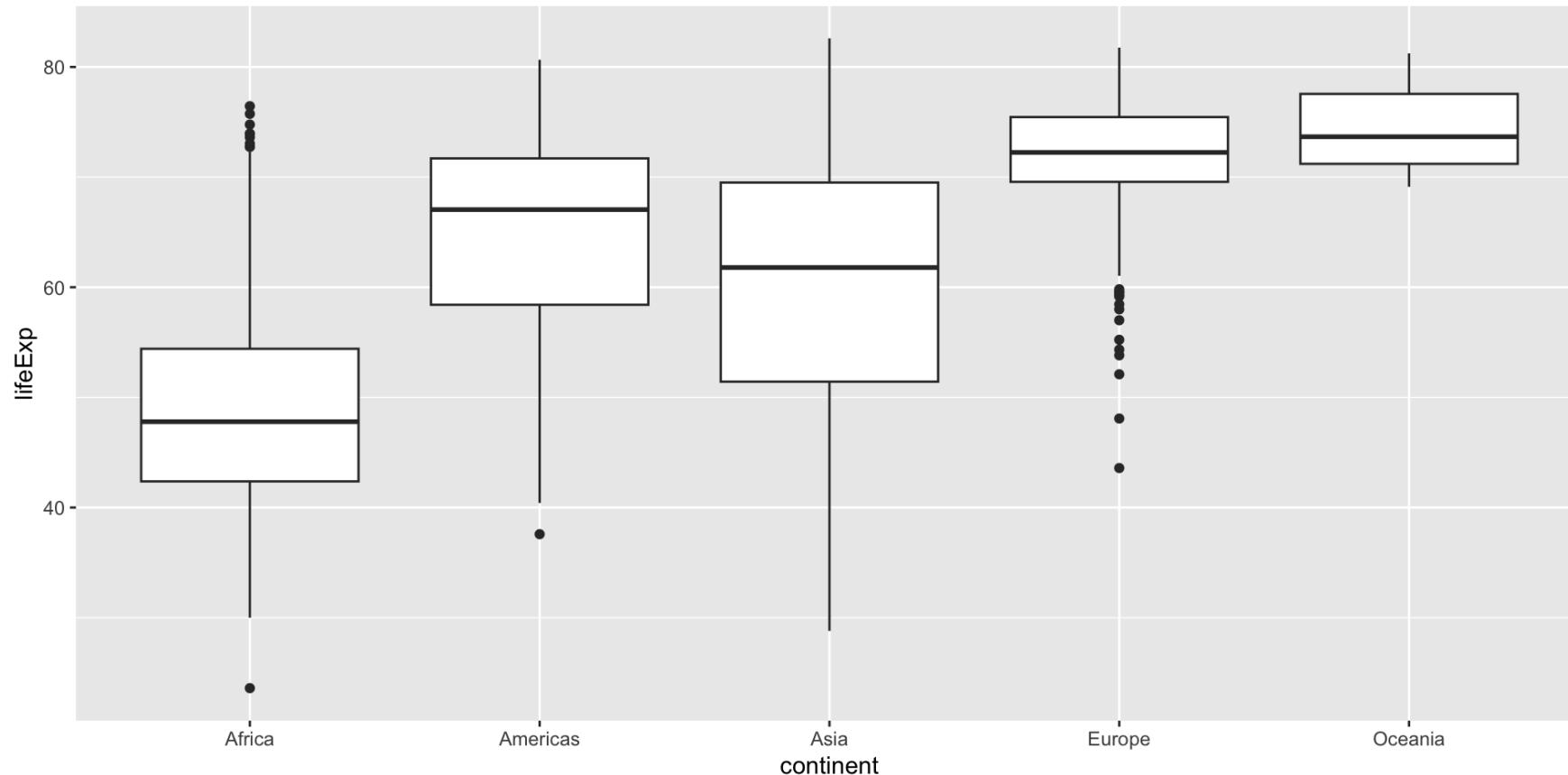
Code structure

```
1 ggplot(data = gapminder,  
2         mapping = aes(x = continent,  
3                           y = lifeExp))
```



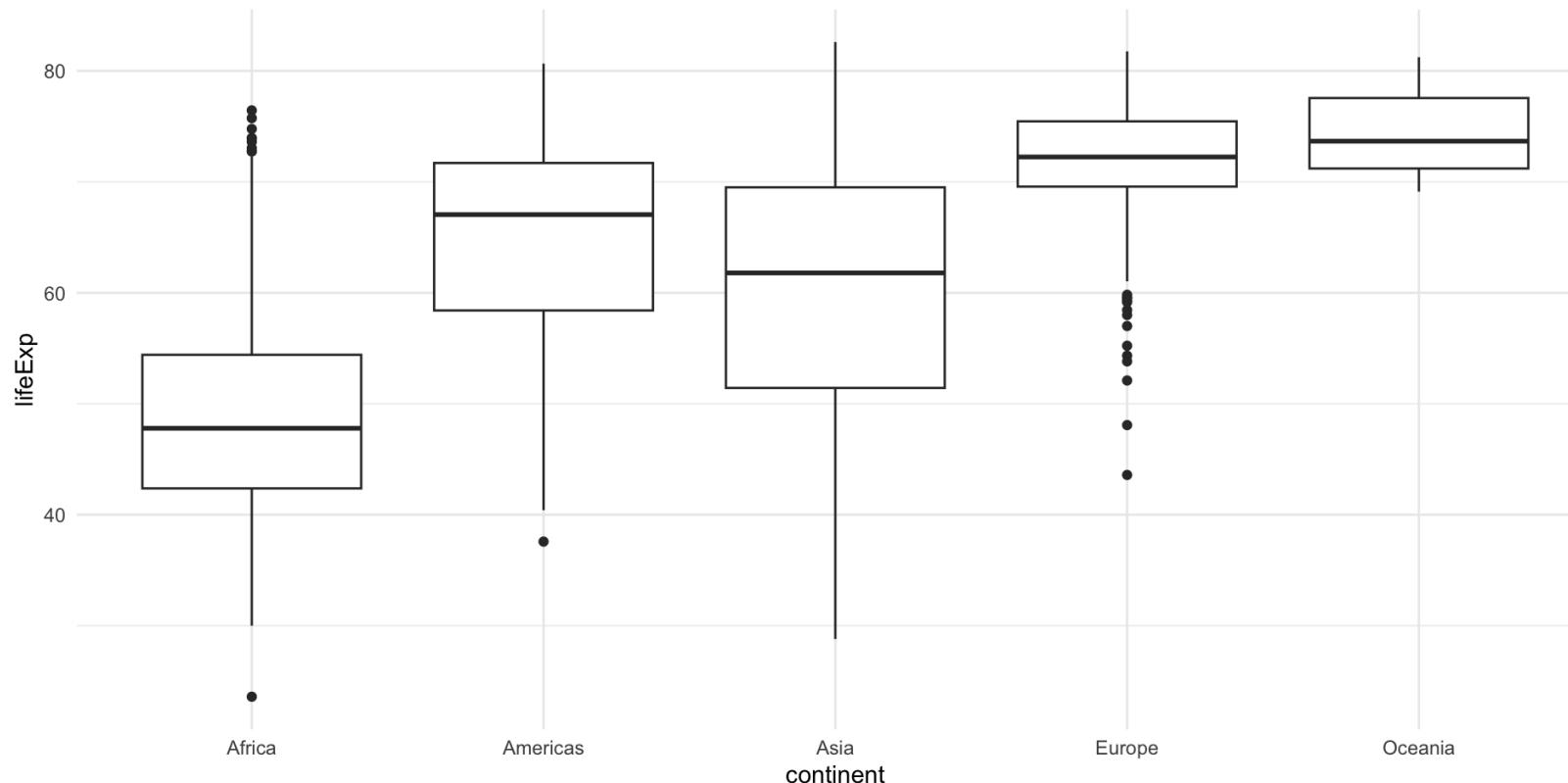
Code structure

```
1 ggplot(data = gapminder,  
2         mapping = aes(x = continent,  
3                           y = lifeExp)) +  
4     geom_boxplot()
```



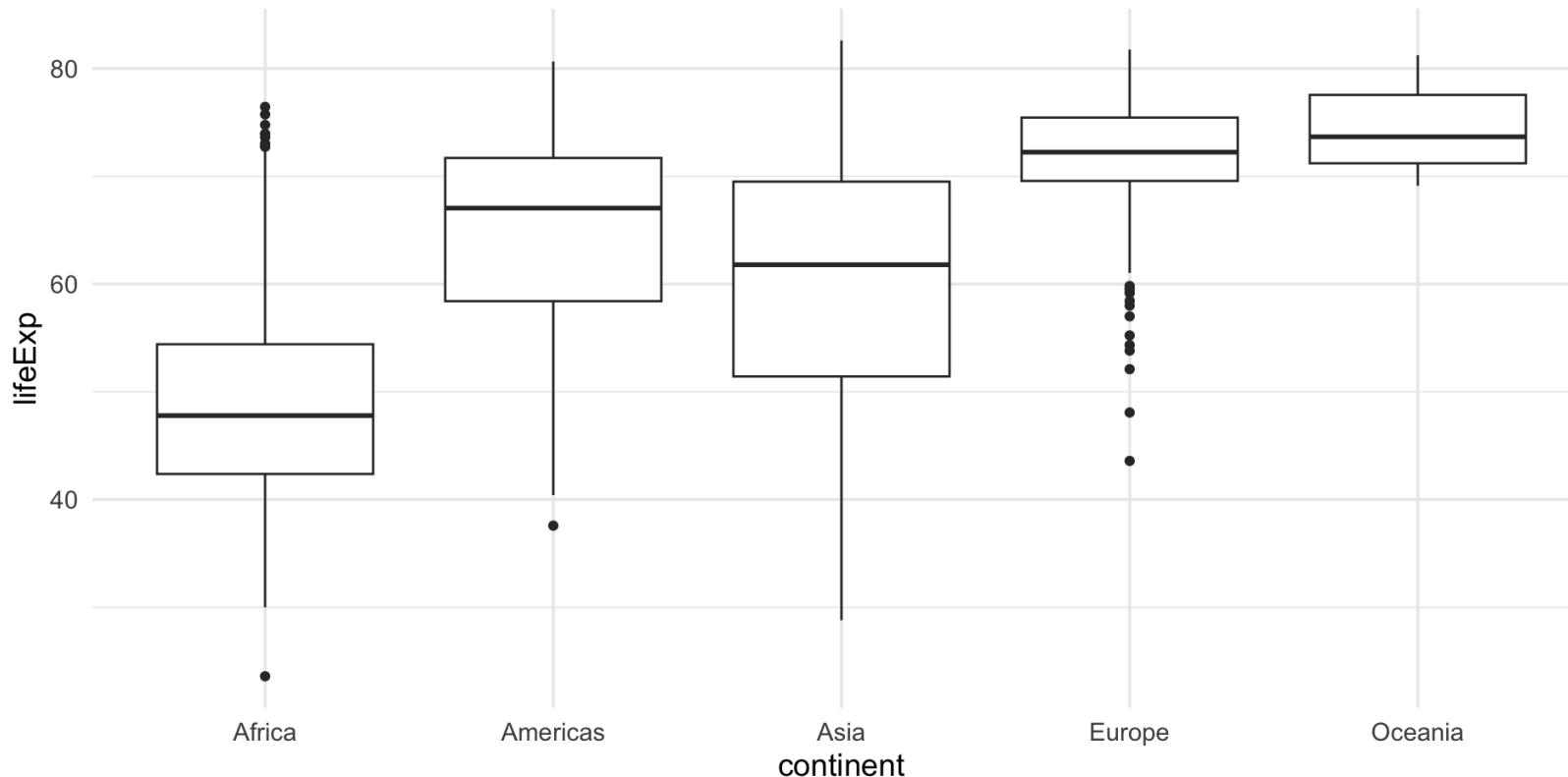
Code structure

```
1 ggplot(data = gapminder,  
2         mapping = aes(x = continent,  
3                           y = lifeExp)) +  
4   geom_boxplot() +  
5   theme_minimal()
```



Code structure

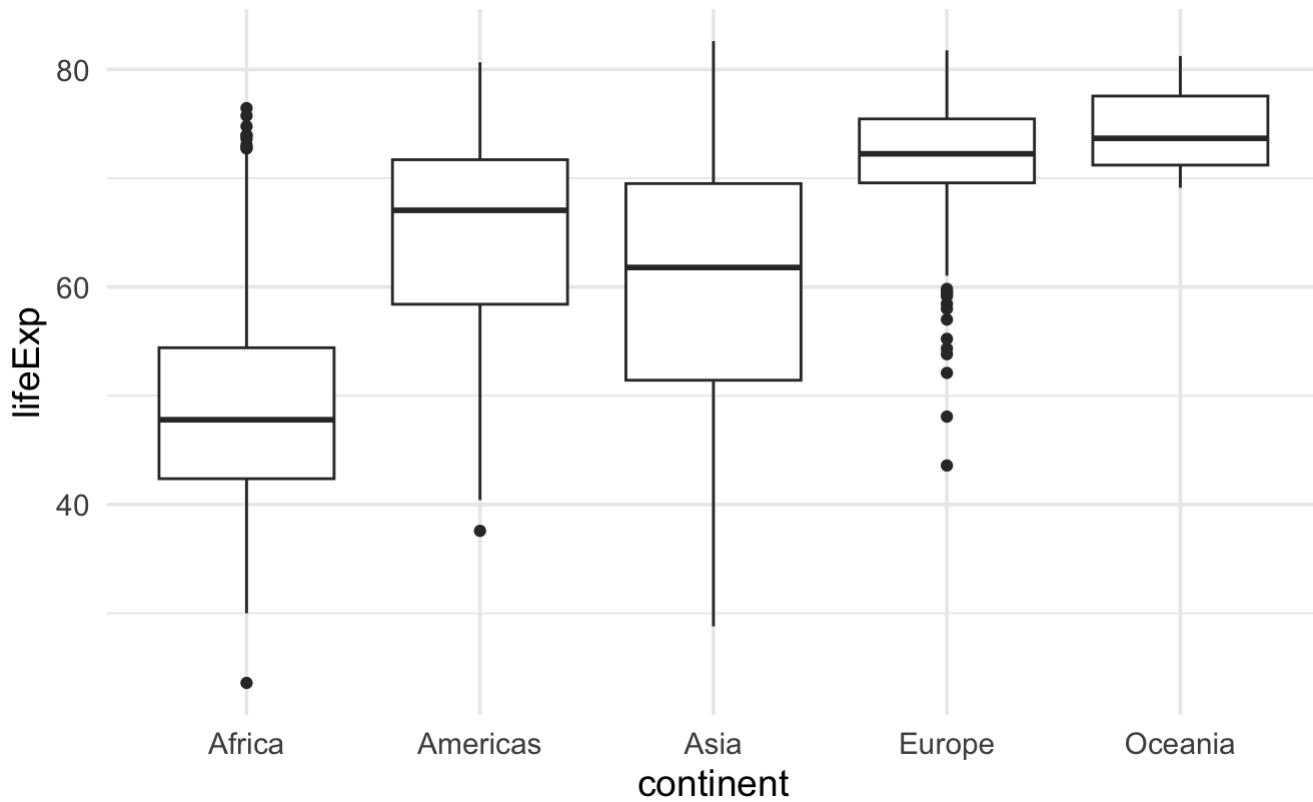
```
1 ggplot(data = gapminder,  
2         mapping = aes(x = continent,  
3                           y = lifeExp)) +  
4   geom_boxplot() +  
5   theme_minimal(base_size = 14)
```



Polls

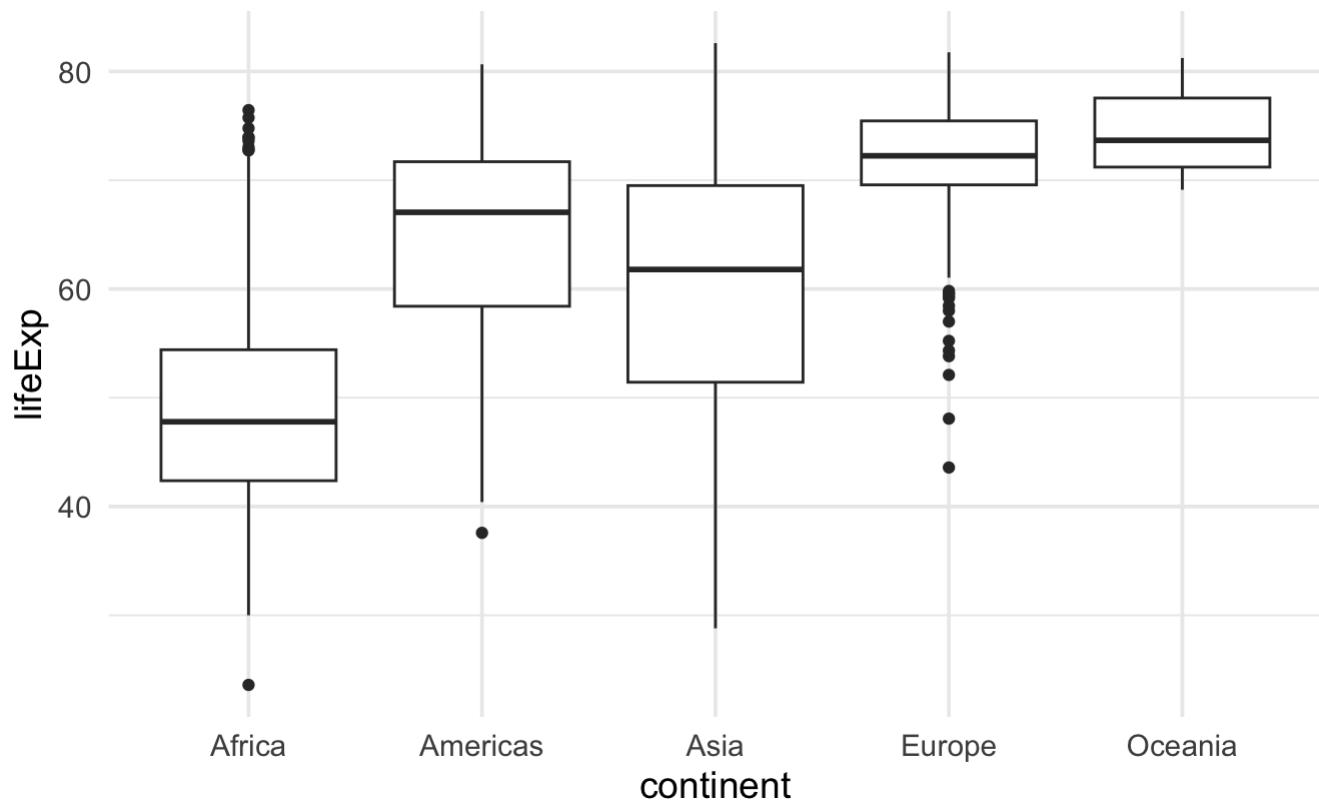
Poll 1: What does the thick line inside the box of a boxplot represent?

1. the mean of the observations
2. the middle of the box
3. the median of the observations
4. none of the above



Poll 2: What percentage of observations are contained inside the box of a boxplot (interquartile range)?

- 1. 25%
- 2. depends on the median
- 3. 50%
- 4. none of the above



Poll 3: What is the median of a set of observations?

1. The median is the most frequently occurring value in a dataset.
2. The median is the sum of all values in a dataset divided by the number of observations.
3. The median is the point above and below which half (50%) of the observations falls.
4. The median is the square root of the sum of the squares of each value in a dataset.

Poll 4: If you have the values: 1, 2, 3, and 10: which statistical measure best represents the “true” value?

1. The mean
2. The standard deviation
3. The median
4. The interquartile range

Boxplot, explained

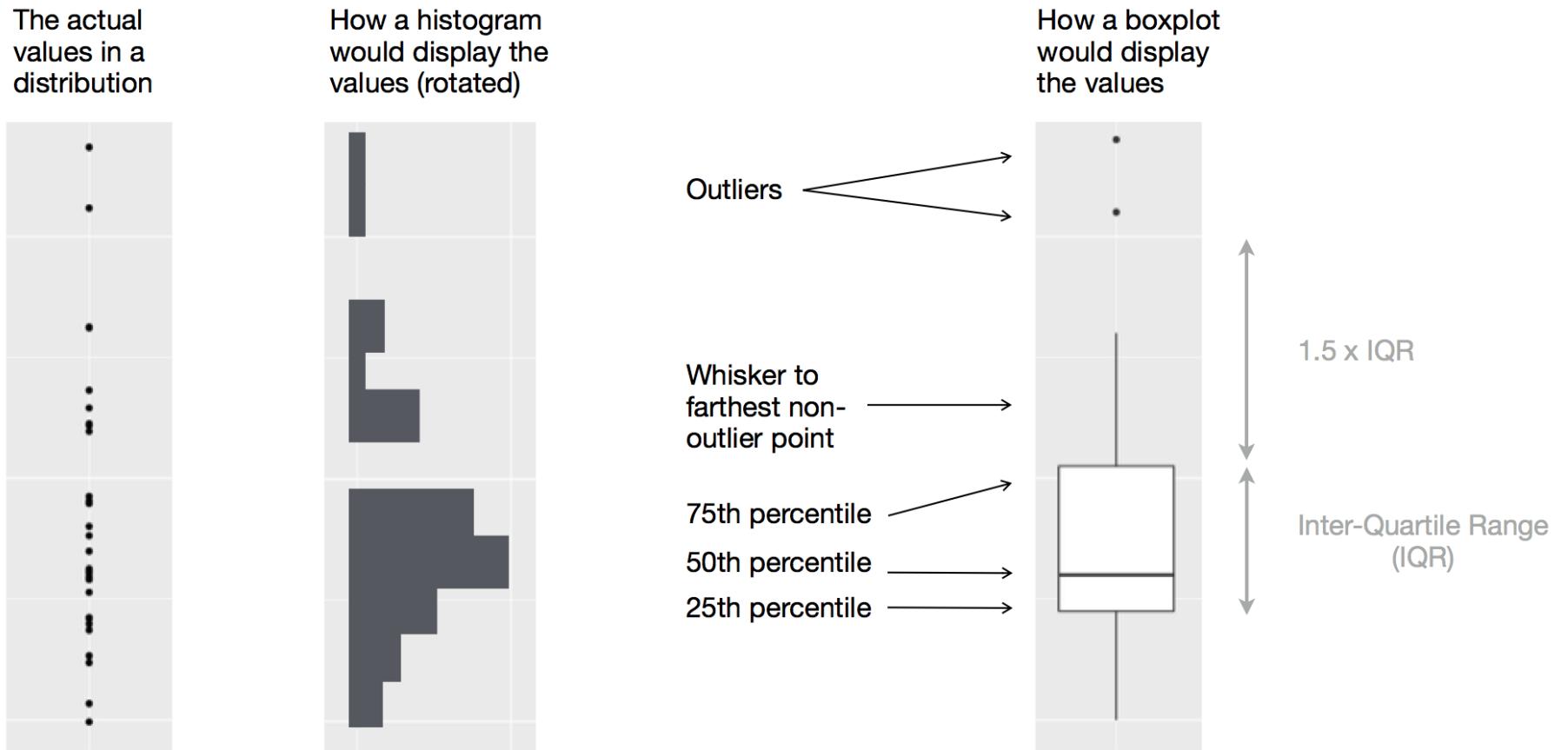


Figure 1: Diagram depicting how a boxplot is created.

Our turn: md-02-exercises

1. Open [posit.cloud](#) in your browser (use your bookmark).
2. Open the [ds4owd workspace](#) for the course.
3. Click [Start](#) next to [md-02-exercises](#).
4. In the File Manager in the bottom right window, locate the [md-02b-data-visualization.qmd](#) file and click on it to open it in the top left window.

Visualising data

Types of variables

numerical

discrete variables

- non-negative
- whole numbers
- e.g. number of students, roll of a dice

continuous variables

- infinite number of values
- also dates and times
- e.g. length, weight, size

non-numerical

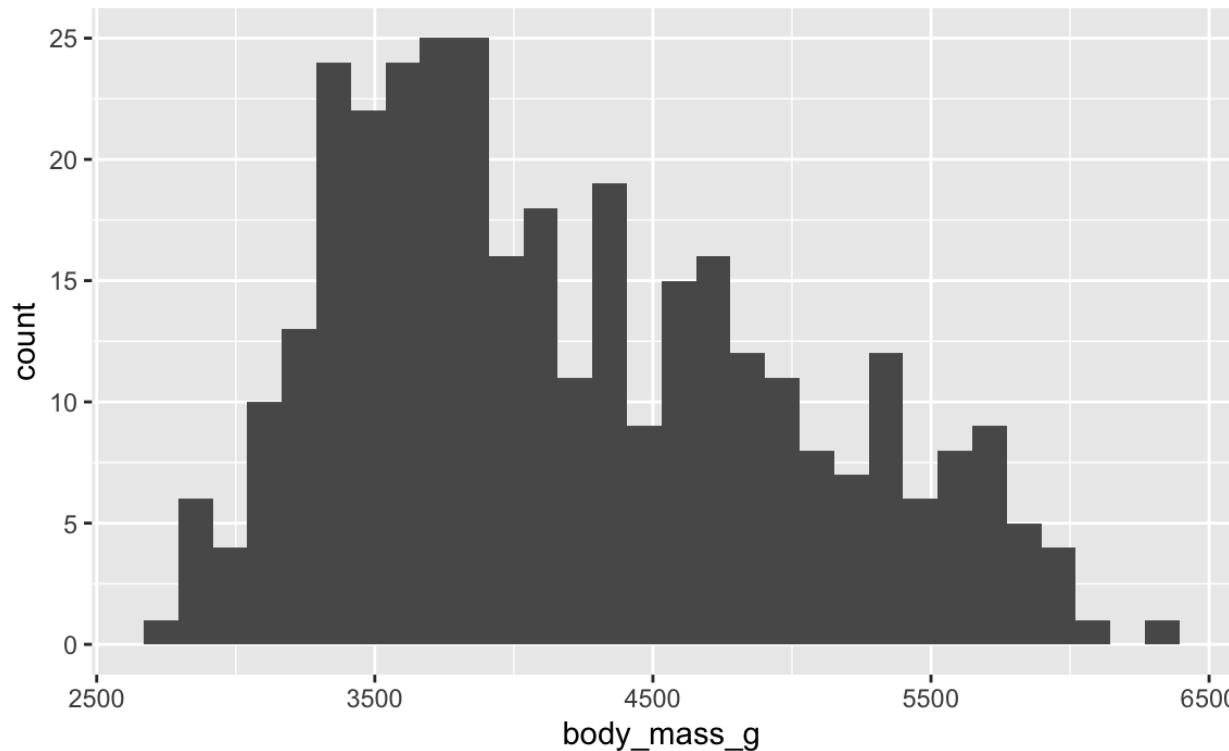
categorical variables

- finite number of values
- distinct groups (e.g. EU countries, continents)
- ordinal if levels have natural ordering (e.g. week days, school grades)

Histogram

- for visualizing distribution of continuous (numerical) variables

```
1 ggplot(data = penguins,  
2         mapping = aes(x = body_mass_g)) +  
3     geom_histogram()
```



Barplot

- for visualizing distribution of categorical (non-numerical) variables

```
1 ggplot(data = penguins,  
2         mapping = aes(x = species)) +  
3     geom_bar()
```

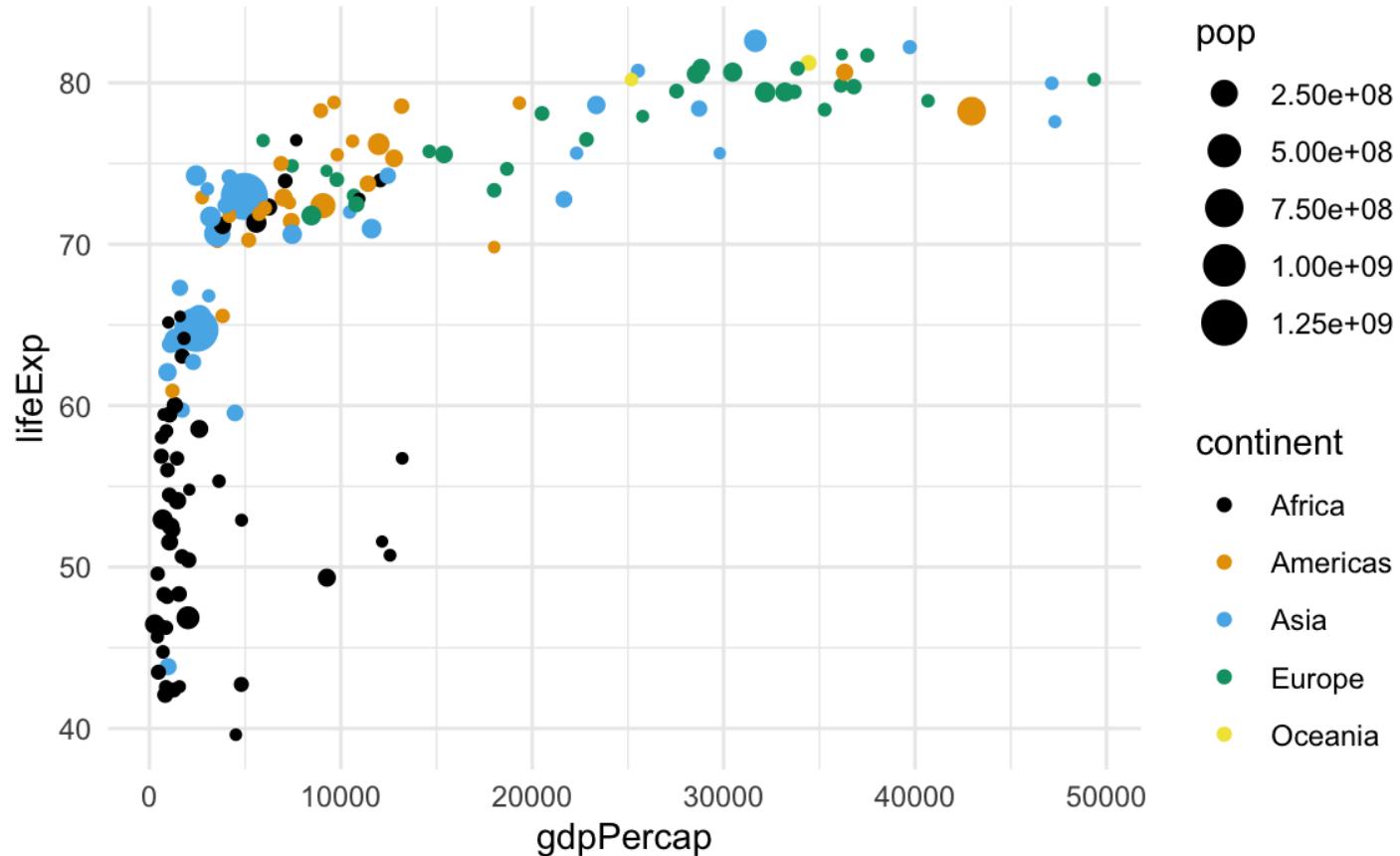
Scatterplot

- for visualising relationships between two continuous (numerical) variables

```

1 ggplot(data = gapminder_2007,
2         mapping = aes(x = gdpPercap,
3                          y = lifeExp,
4                          size = pop,
5                          color = continent)) +
6   geom_point() +
7   scale_color_colorblind() +
8   theme_minimal()

```



Your turn: md-02-exercises

1. Open [posit.cloud](#) in your browser (use your bookmark).
2. Open the [ds4owd workspace](#) for the course.
3. In the File Manager in the bottom right window, locate the [md-02c-make-a-plot.qmd](#) file and click on it to open it in the top left window.
4. Follow instructions in the file

Homework assignments

module 2

Module 2 documentation

ds4owd-001.github.io/website/modules/md-02.html

Module 2

Data science lifecycle & Quarto

◎ Learning Objectives

1. Learners can render a Quarto file to an output file in HTML, PDF and DOCX format.
2. Learners can list the six elements of the data science lifecycle.
3. Learners can identify four components of a Quarto file (YAML, code chunk, R code, markdown).

Slides

- In preparation

Readings

Homework due date

- Homework assignment due: Monday, November 13th
- Correction & feedback phase up to: Thursday, November 16th

Wrap-up

Thanks! 🌻

Slides created via revealjs and Quarto:

<https://quarto.org/docs/presentations/revealjs/> Access slides as
[PDF on GitHub](#)

All material is licensed under [Creative Commons Attribution Share Alike 4.0 International.](#)

