



GitFun!

Ana Martinovici

20200206

Fun,
but first plan

- News
- **Why** we should use version control
- **How** to use version control

News

source: Antonio Schettino (OSCR)

- Talk by Daniel Lakens – February 11 at 15.00 in room T13-67
 - <https://www.openscience-rotterdam.com/2020/01/dpecs-lakens-feb2020/>
- Open science festival! August 27 2020, Wageningen
 - <https://opensciencefestival.nl/>
- How to use R and the tidyverse to clean, plot, and analyze data:
 - https://github.com/aschetti/MPI2020_intro_tidyverse

WHAT DO YOU MEAN

THERE'S NO VERSION CONTROL?

makeameme.org

WHY

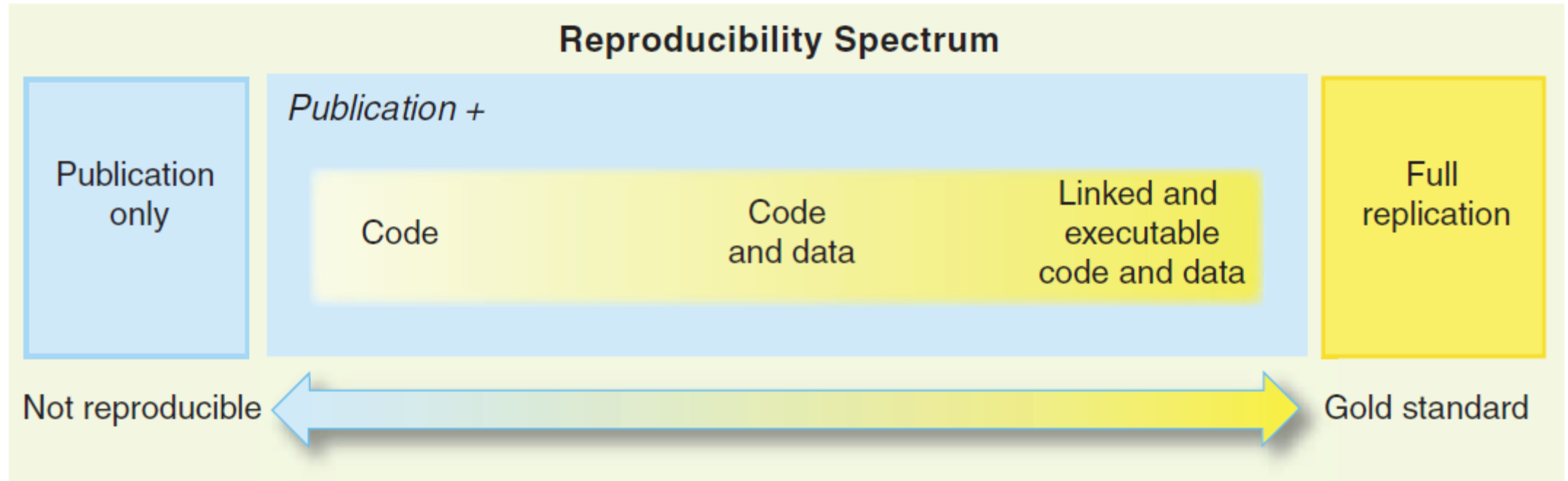
(some) Research goals

- Collaborate with co-authors
- Contribute to ongoing projects
- Share your work with the world
- Reproducible results

Reproducibility?

A study is reproducible if you can take the original data and the computer code used to analyze the data and reproduce all of the numerical findings from the study.

Reproducibility? Yes, No, Maybe



Source: Roger D. Peng (2011) "Reproducible Research in Computational Science", Science

Objective

Reproduce **ALL** results and **NOTHING BUT** results included in:

- Current version of the paper
- Previous analyses

as fast and easy as possible.

Best case scenario

Step 1

Raw
data

Cleaning /
processing

Step 2

Processed
data

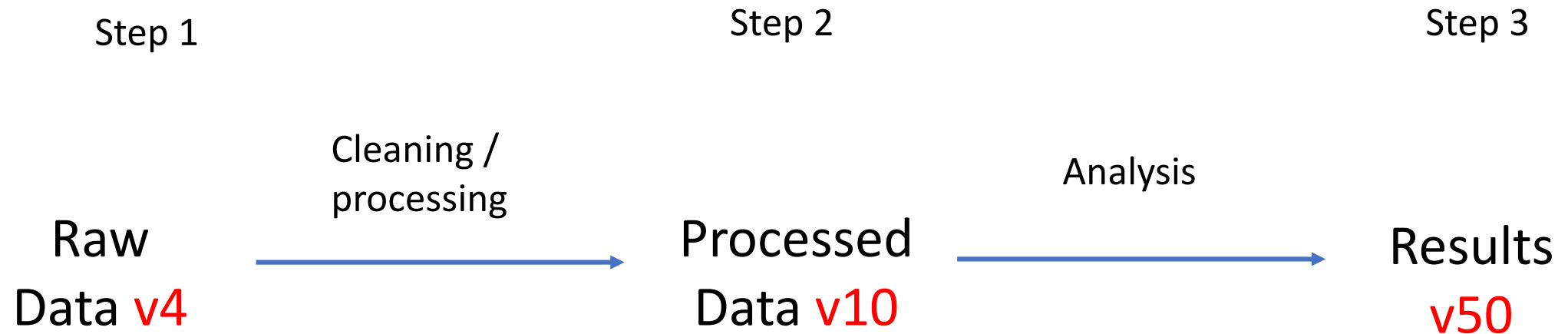
Analysis

Step 3

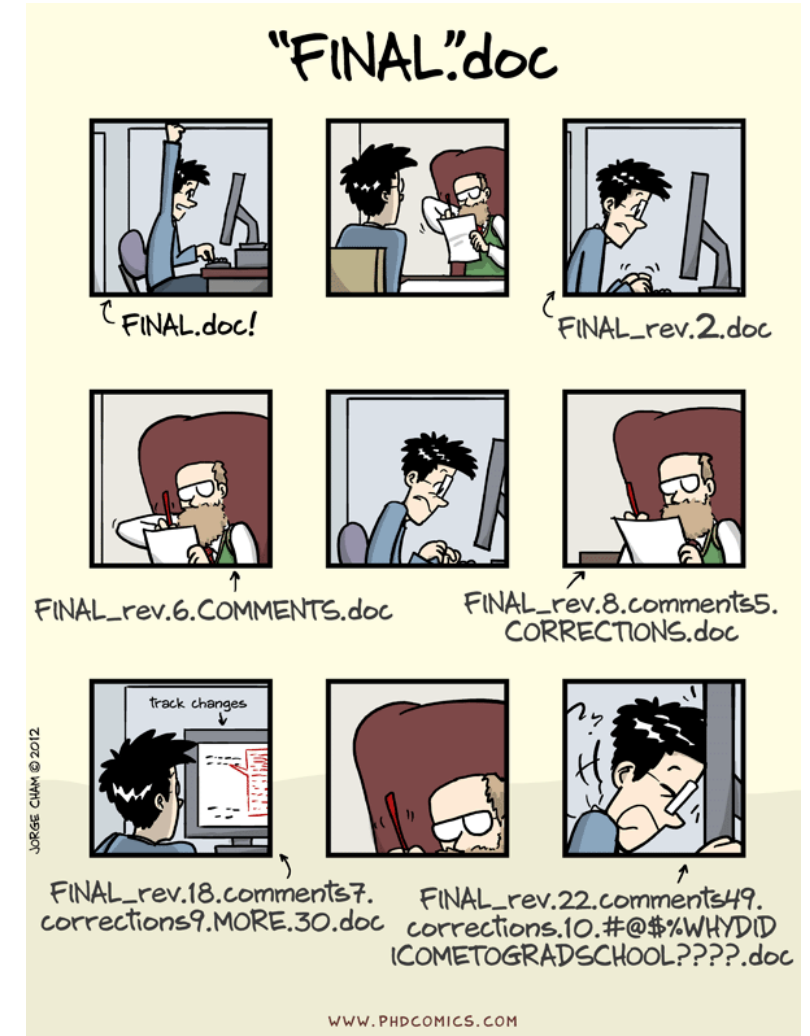
Results

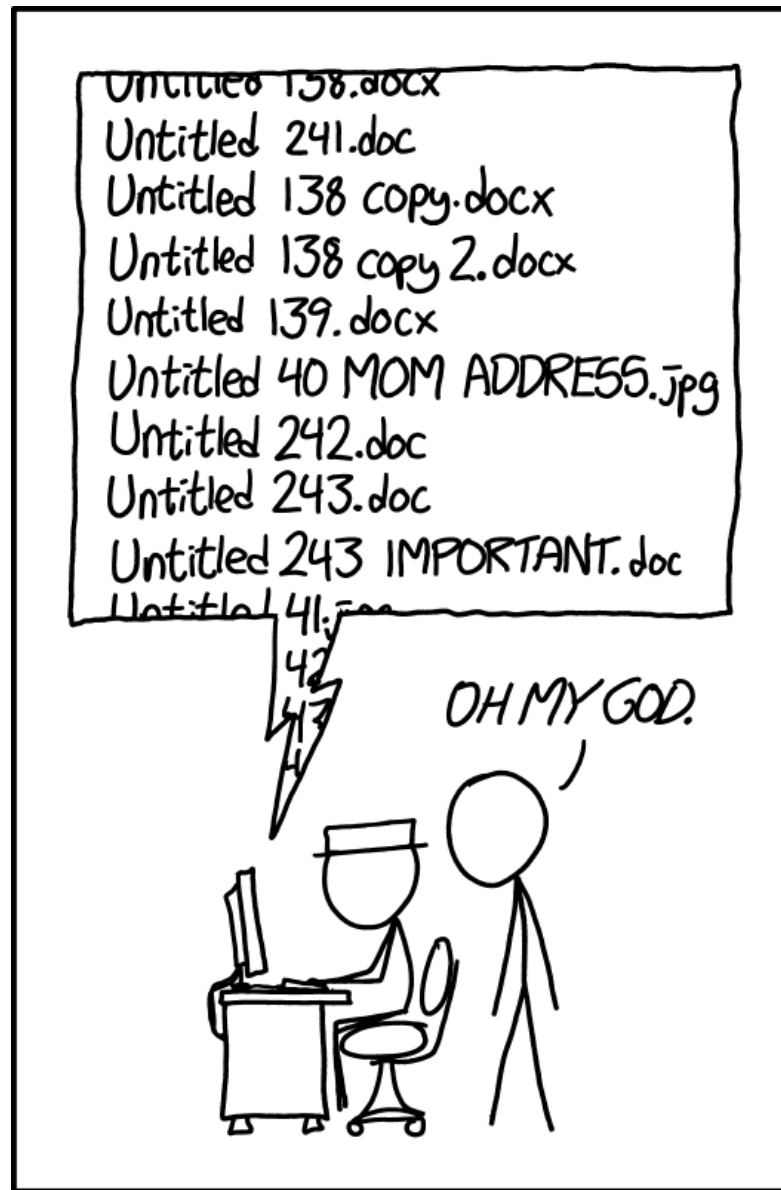


Most likely scenario



Final is never final





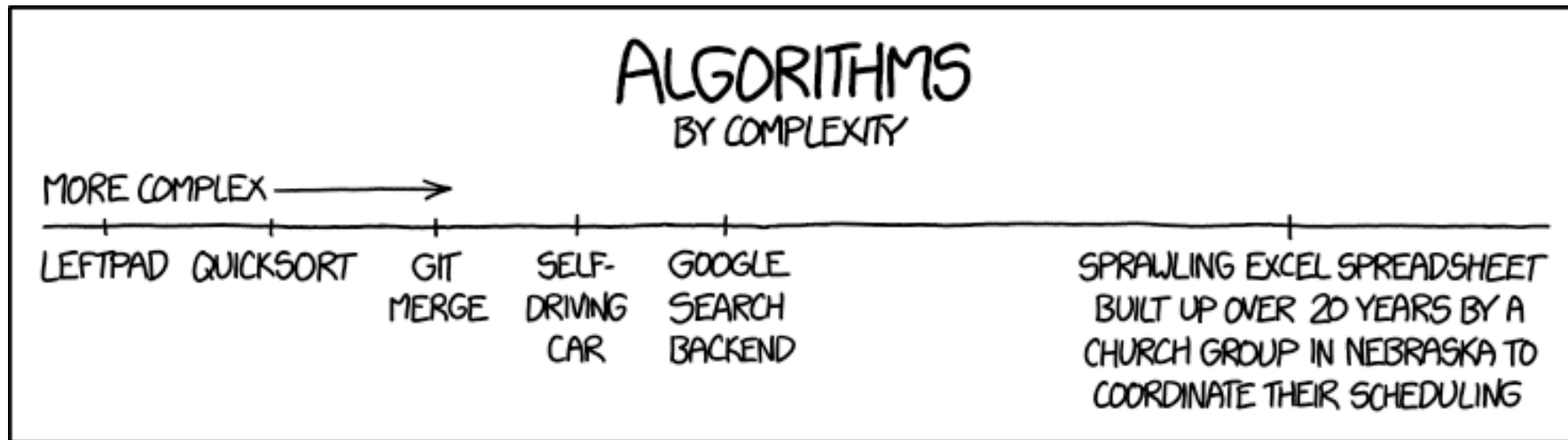
PROTIP: NEVER LOOK IN SOMEONE
ELSE'S DOCUMENTS FOLDER.

Things that can go wrong

- Your PC/laptop/external HD explode
- Overwriting files, but also not overwriting files
- Forgetting which of the “final” files is really final
- Change file X, but forget to update all the other files/results that depend on it
- Software changes

Things that can go wrong

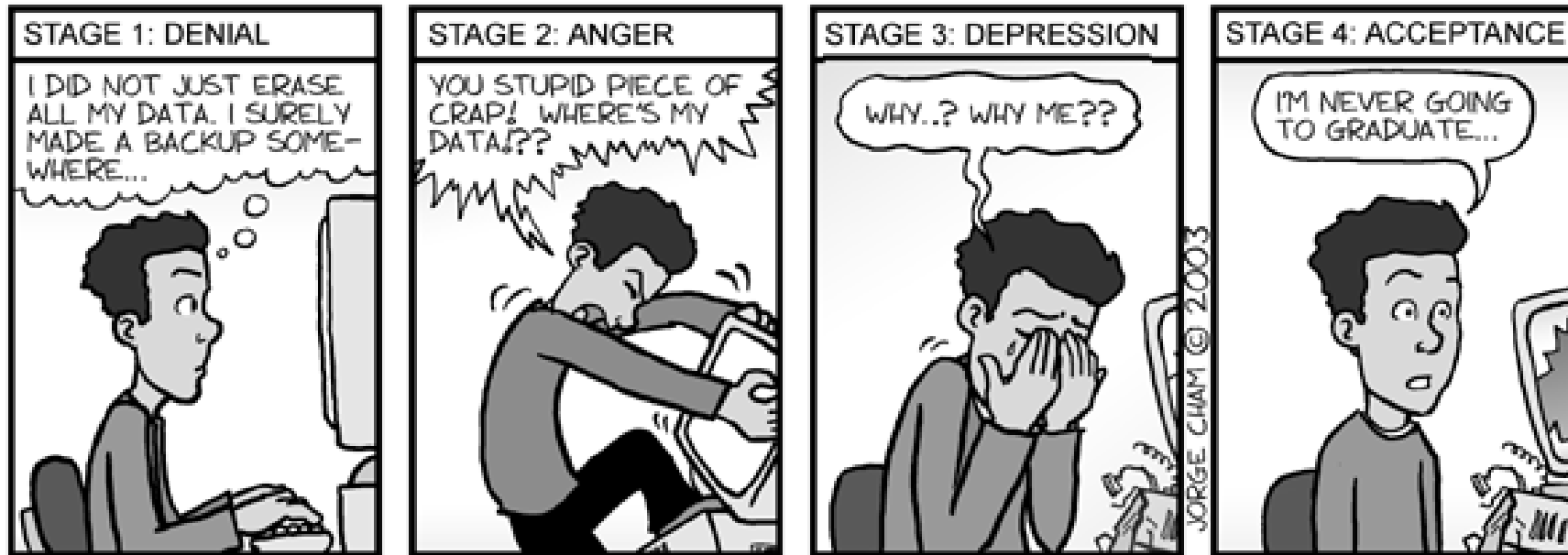
- Excel (in so many ways)



Things that can go wrong

THE FOUR STAGES OF DATA LOSS

DEALING WITH ACCIDENTAL DELETION OF MONTHS OF
HARD-EARNED DATA



www.phdcomics.com

Version control:



- Time travel
- Safe place for your data and code
- Notes to your future self

What is git?

- Formal version control system
- Developed by Linus Torvalds
 - Used to manage the course code for Linux
- Tracks content:
 - Source code
 - Data analysis projects
 - Websites
 - Presentations
 - Manuscripts

What is GitHub?

- A home for git repositories
- Interface for exploring public git repositories
- A 'safe place' to keep code and data
- Additional benefits: issues, projects, wiki, insights

Why use GitHub?

- Facilitates
 - Exploring code
 - Tracking issues
 - Learning from others
- Lowers the barriers to collaboration
 - Email “there’s a typo in your code in file X, line 30” vs
 - Pull request “here’s a correction to your code”
- Free for researchers and students

I SHOULD USE GITHUB



HOW

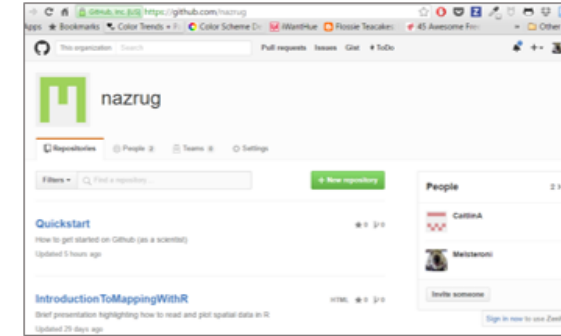
\$195 for 3h Intro in NYC

- <https://www.nobledesktop.com/classes/git-classes-nyc>

General idea

<https://jules32.github.io/2016-07-12-Oxford/git/>

REMOTE (aka Github website)



Clone (i.e., copy)
repository to your
computer (a one
time event)

Pull remote
changes

Push local
changes



LOCAL
(aka your computer)

Vocabulary

- **repository (repo):** (*noun*) folder containing all tracked files as well as the version control history
- **commit:** (*noun*) a snapshot of changes made to the staged file(s); (*verb*) to save a snapshot of changes made to the staged file(s)
- **stage:** (*noun*) the staging area holds the files to be included in the next commit; (*verb*) to mark a file to be included in the next commit
- **track:** (*noun*) a tracked file is one that is recognized by the Git repository
- **branch:** (*noun*) a parallel version of the files in a repository

Vocabulary



- **local:** (*noun*) the version of your repository that is stored on your personal computer
- **remote:** (*noun*) the version of your repository that is stored on a remote server; for instance, on GitHub
- **clone:** (*verb*) to create a local copy of a remote repository on your personal computer
- **fork:** (*noun*) a copy of another user's repository on GitHub; (*verb*) to copy a repository; for instance, from one user's GitHub account to your own

Vocabulary

- **merge:** (*verb*) to update files by incorporating the changes introduced in new commits
- **pull:** (*verb*) to retrieve commits from a remote repository and merge them into a local repository
- **push:** (*verb*) to send commits from a local repository to a remote repository
- **pull request:** (*noun*) a message sent by one GitHub user to merge the commits in their remote repository into another user's remote repository

First use of git

```
$ git config --global user.name "Vlad Dracula"
```

```
$ git config --global user.email "vlad@tran.sylvan.ia"
```

<https://swcarpentry.github.io/git-novice/02-setup/index.html>

Challenges

- Data: big, small, too easy to get, too hard to get
- Changes during the review process
- Dependencies (e.g. R packages)
- System specifications
- Lisa/clusters



Checklist

- Have I done anything by hand?
 - If so, are those parts precisely documented?
 - Is that documentation saved in a 'safe' place?
 - Is the documentation a complete, correct, and specific description of what was done by hand?
- Have I coded as many of the steps as I could?
- Am I using version control?
- Have I documented the software environment?
- Have I saved any output that I cannot reconstruct from the original data?

Not a good idea: Doing things by hand

- Editing spreadsheets of data to “clean it up”
 - Removing outliers
 - Rescaling (reverse coding)
 - Create new variables (dummy variables, intervals, categories)
- Edit tables or figures (e.g. rounding, formatting)
- Move/split/rename data files on your computer
- “I’m only doing this once...”

Things done by hand need to be precisely documented (harder than it sounds).

Not a good idea: point and click

- Many data processing / statistical analysis packages have graphical user interfaces (GUIs)
- GUIs are convenient but the actions you take can be difficult to reproduce
- Some GUIs produce a log file that can be saved for later examination

Not a good idea: saving output only

- Avoid saving data analysis (intermediary) output except perhaps temporarily for efficiency purposes
- Intermediate files can be ok as long as:
 - (1) there is clear documentation of how they were created and
 - (2) the links to and from these files are working
- Save the data + code that generated the output, rather than the output itself

Homework fun!

To 'push' you to practice the clone -> stage -> commit -> push steps 😊

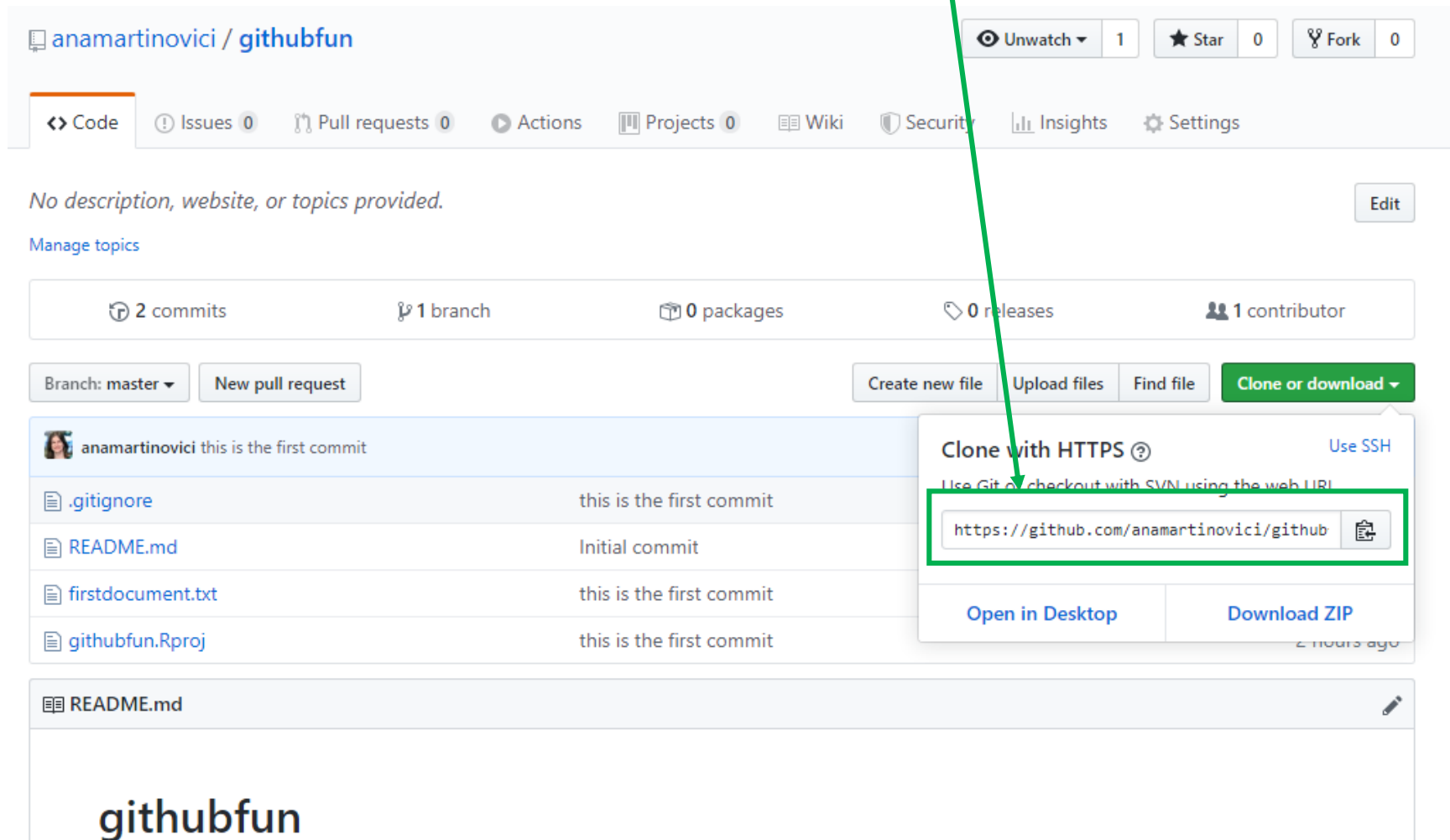
- Go to <https://github.com/anamartinovici/githubfun>
- Make sure you are a collaborator (check next slide)
- Clone the repository
- Add a txt file with your name ("name.txt") and some text that you're comfortable sharing with the world
- Push your changes

How to make sure you are a collaborator

I don't know all your GitHub user names, so I can't add all of you immediately. So this is what you can do:

- Go to the repository page: <https://github.com/anamartinovici/githubfun>
- Click on Issues and create a new one. After all, it is an issue if you can't contribute to this repository.
- Then, I will add you as a collaborator 😊

On GitHub.com: copy this link



The screenshot shows the GitHub interface for the repository 'anmartinovici / githubfun'. The repository has 1 watch, 0 stars, and 0 forks. The 'Code' tab is selected, showing a list of files: .gitignore, README.md, firstdocument.txt, and githubfun.Rproj. A modal window titled 'Clone with HTTPS' is open, displaying the URL 'https://github.com/anmartinovici/githubfun' which is highlighted with a green box. A green arrow points from the text 'On GitHub.com: copy this link' to this URL. The modal also includes options for 'Open in Desktop' and 'Download ZIP'.

anmartinovici / githubfun

Unwatch 1 Star 0 Fork 0

Code Issues 0 Pull requests 0 Actions Projects 0 Wiki Security Insights Settings

No description, website, or topics provided. Edit

Manage topics

2 commits 1 branch 0 packages 0 releases 1 contributor

Branch: master New pull request Create new file Upload files Find file Clone or download

anmartinovici this is the first commit

File	Commit Message
.gitignore	this is the first commit
README.md	Initial commit
firstdocument.txt	this is the first commit
githubfun.Rproj	this is the first commit

README.md

githubfun

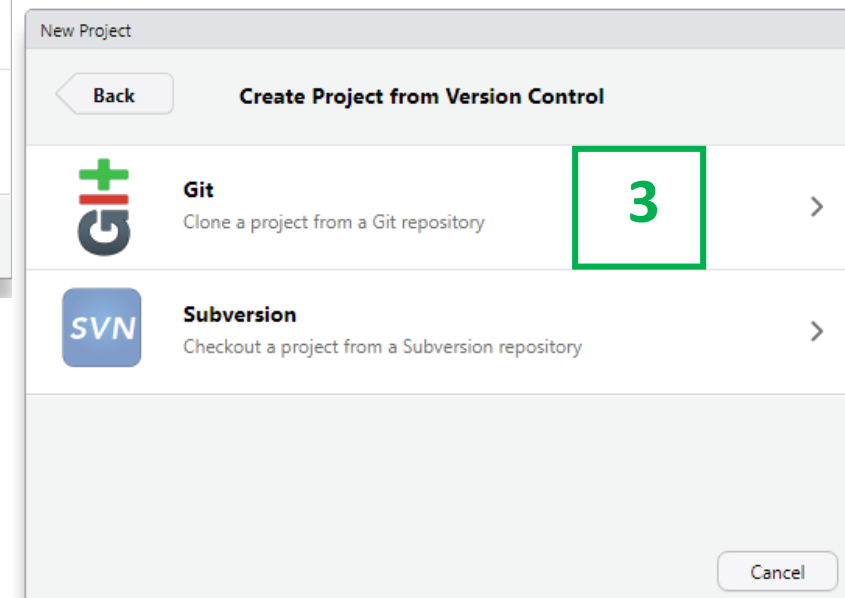
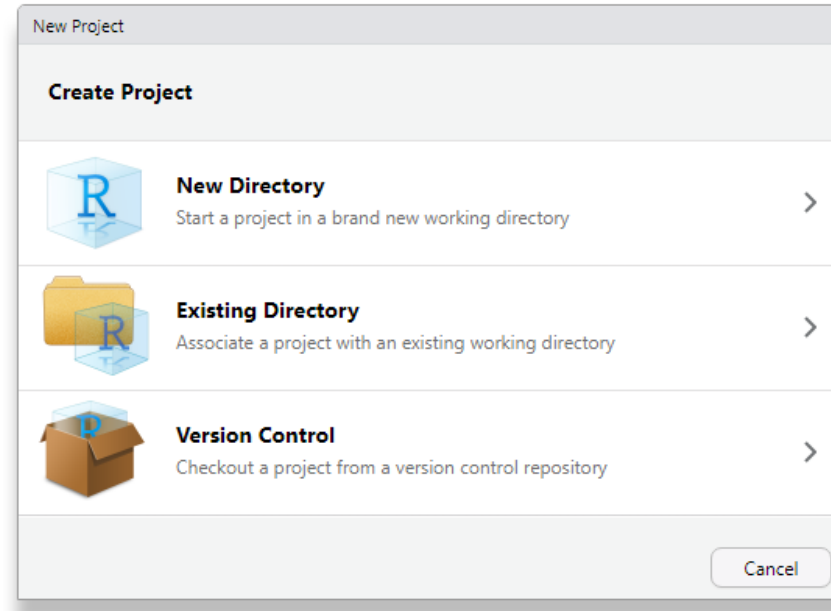
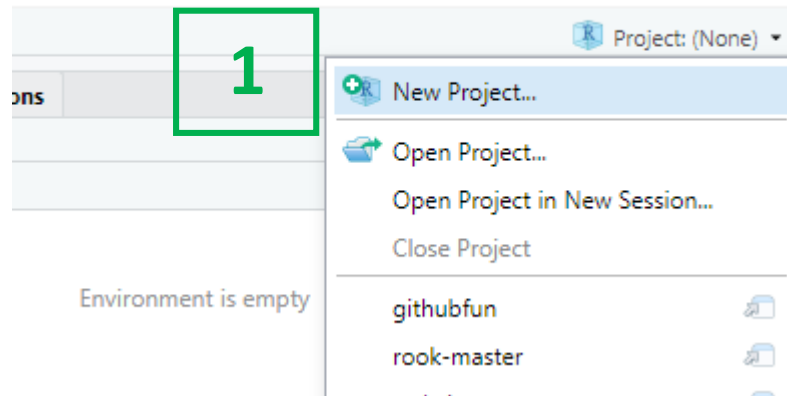
Clone with HTTPS Use SSH

Use Git to checkout with SVN using the web URL

https://github.com/anmartinovici/githubfun

Open in Desktop Download ZIP

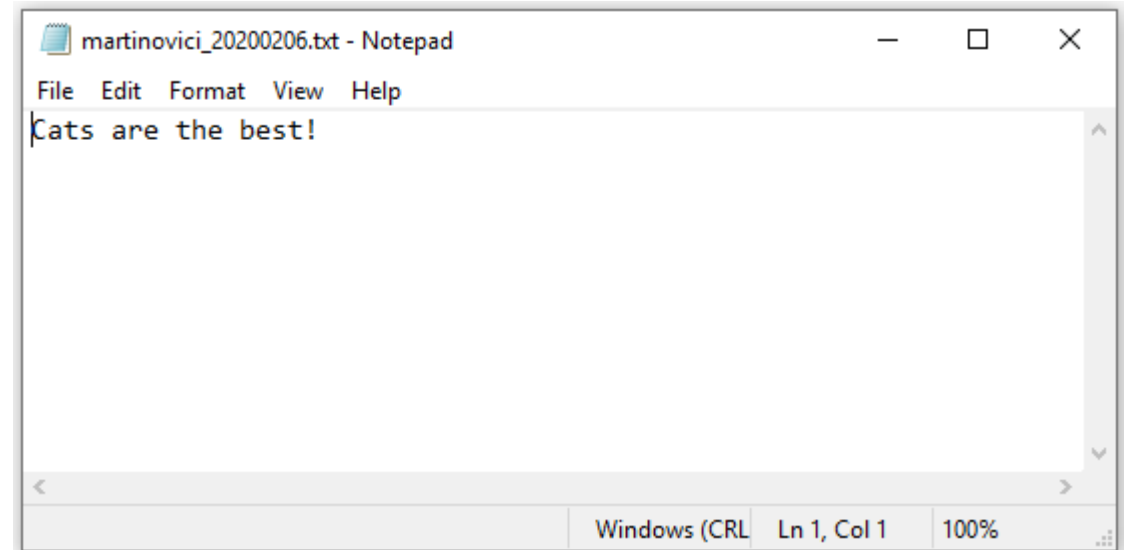
In Rstudio: new version control project (git)



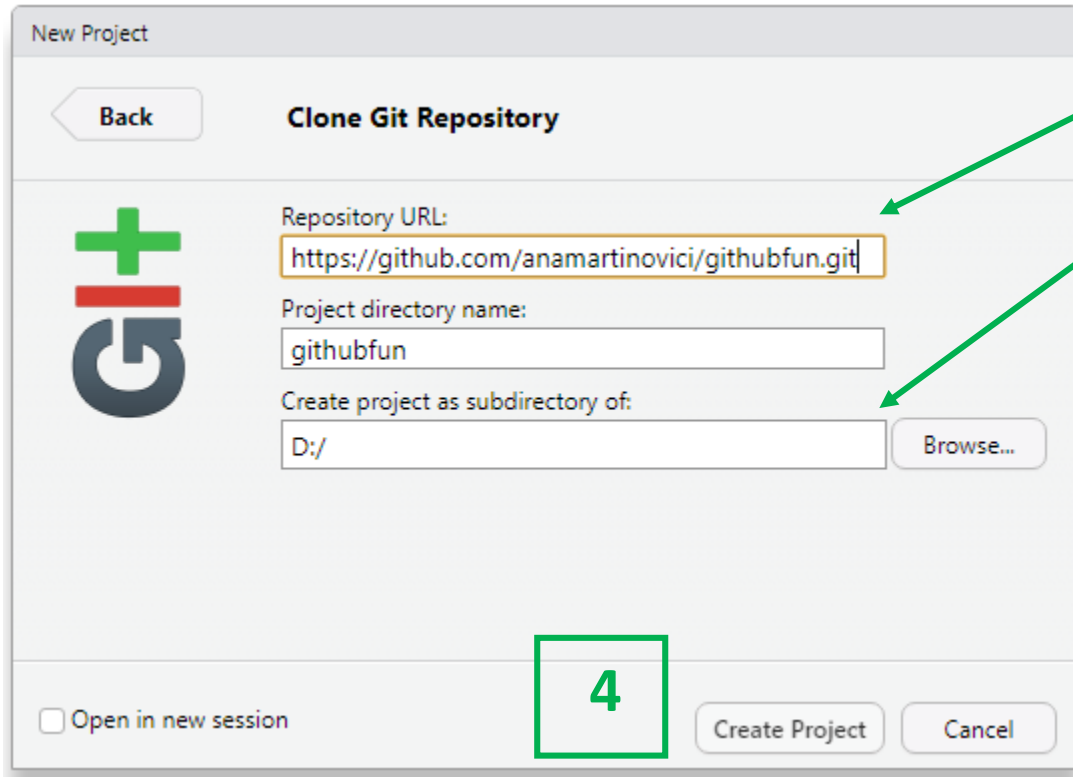
On your laptop/PC: create a file with your name

DATADRIVE1 (D:) > githubfun

Name	Date modified	Type	Size
.git	6-2-2020 14:52	File folder	
.Rproj.user	6-2-2020 14:11	File folder	
.gitignore	6-2-2020 14:11	Text Document	1 KB
githubfun.Rproj	6-2-2020 14:51	R Project	1 KB
martinovici_20200206.txt	6-2-2020 14:52	Text Document	1 KB
README.md	6-2-2020 14:11	MD File	1 KB



In Rstudio: new version control project (git)



New Project

Back Clone Git Repository

Repository URL:
https://github.com/anamartinovici/githubfun.git

Project directory name:
githubfun

Create project as subdirectory of:
D:/ Browse...

☐ Open in new session

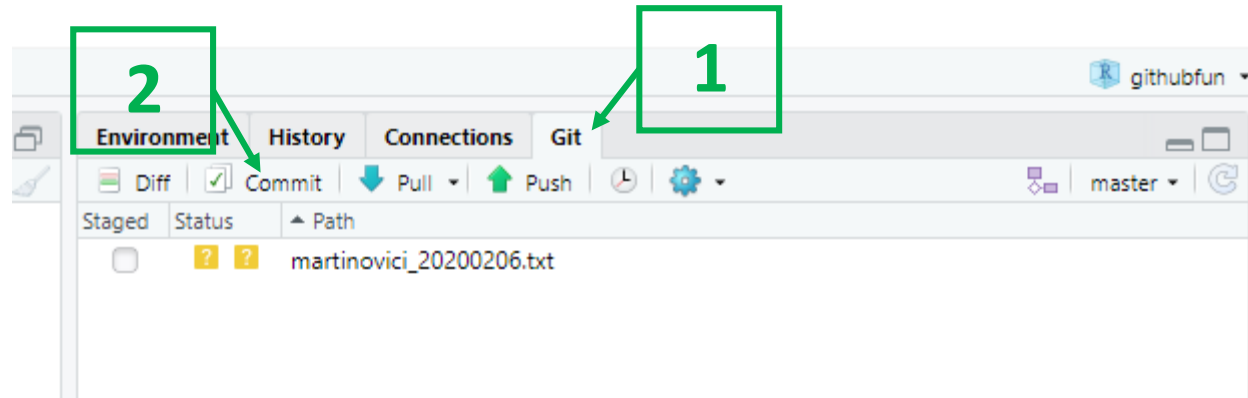
4 Create Project Cancel

- Paste the link you've previously copied

IMPORTANT:

- If you use a device from the university (laptop, PC), create the project on a local drive (C:/ or D:/). Don't create it on what appears to be the desktop ("C:\some numbers and letters\Desktop"), as that's in fact a network location.
- If you use your own device, then you can save it on Desktop.
- Regardless of which device you use, do NOT save a repository in Dropbox. This will create sync problems that are better avoided. Dropbox and GitHub are useful, but for different purposes.

In Rstudio: open the commit window

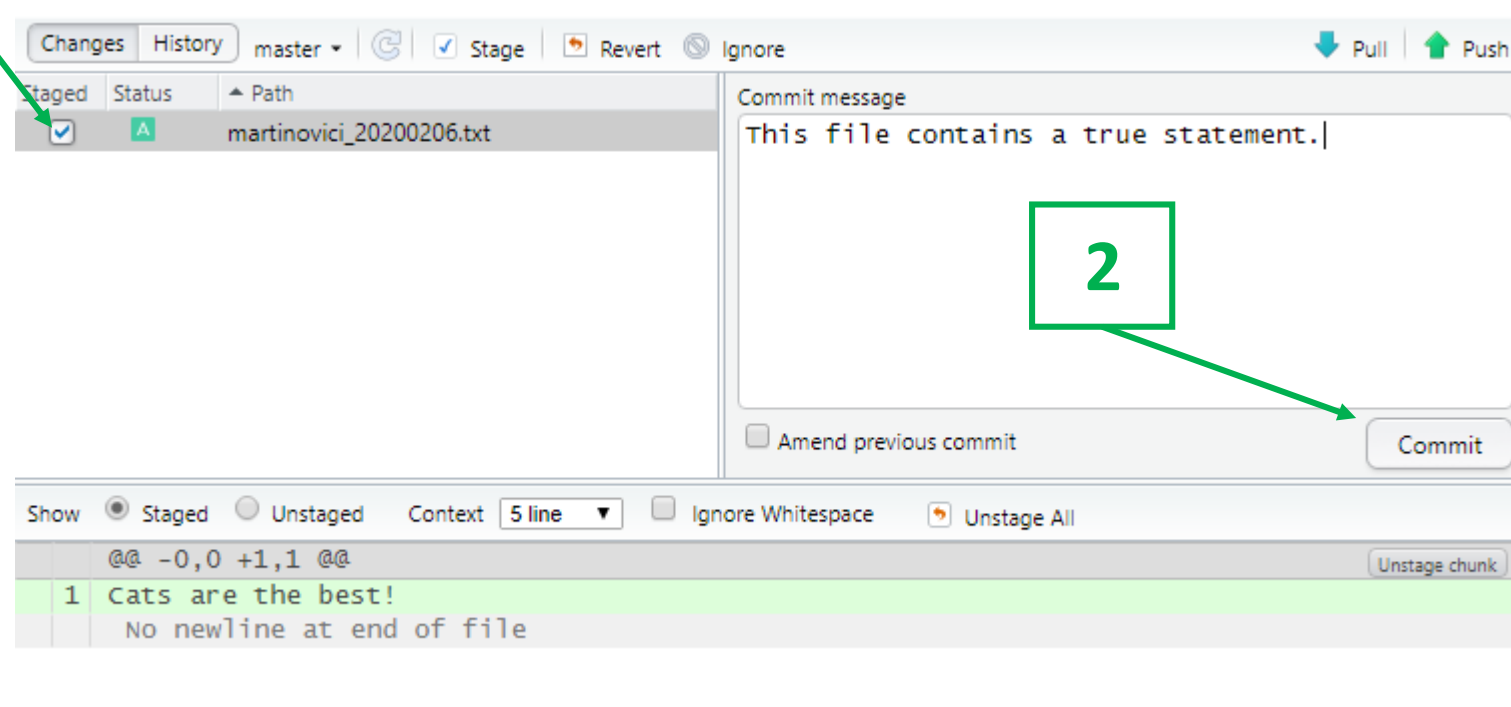


In Rstudio: stage -> commit -> push

1

3

2



On GitHub: refresh the page for “githubfun”

The screenshot shows the GitHub repository page for 'githubfun' by user 'anamartinovici'. The repository has 1 star, 0 forks, and 0 issues. The 'Code' tab is selected, showing a list of files and commits. A green box highlights the 'New file' button, which is located next to the 'Clone or download' button. The file list includes .gitignore, README.md, githubfun.Rproj, and martinovici_20200206.txt. The commit history shows four commits, with the latest commit being 'this is the first commit' 1 minute ago.

anamartinovici / githubfun

Unwatch 1 Star 0 Fork 0

Code Issues 0 Pull requests 0 Actions Projects 0 Wiki Security Insights Settings

No description, website, or topics provided. Edit

Manage topics

4 commits 1 branch 0 packages 0 releases 1 contributor

Branch: master New pull request Create new file Upload files Find file Clone or download

anamartinovici This file contains a true statement. Latest commit 0cd5f0b 1 minute ago

.gitignore	this is the first commit	2 hours ago
README.md	Initial commit	2 hours ago
githubfun.Rproj	this is the first commit	2 hours ago
New file martinovici_20200206.txt	This file contains a true statement.	1 minute ago

README.md

References and additional resources

- https://github.com/DataScienceSpecialization/courses/blob/master/05_ReproducibleResearch/
- <https://github.com/jasonmtroos/rook>
- <http://kbroman.org/Tools4RR/>
- <http://blogs.nature.com/naturejobs/2018/06/11/git-the-reproducibility-tool-scientists-love-to-hate/>
- <https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1004668>
- <https://swcarpentry.github.io/git-novice/>
- <https://jules32.github.io/2016-07-12-Oxford/git/>
- <https://desktop.github.com/>