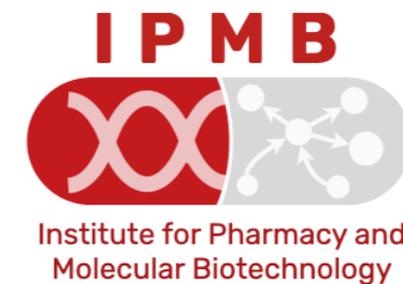


# Data Analysis Projects

## MoBi 4. FS - SoSe 2023

19.04.2023



UNIVERSITÄT  
HEIDELBERG  
ZUKUNFT  
SEIT 1386

- Allgemeine Vorstellung des Moduls (C. Herrmann ~15-20 Minuten)
- Vorstellung der 5 Themen (ca. 10 Minuten/Thema)
- Kurze Einführung in GitHub

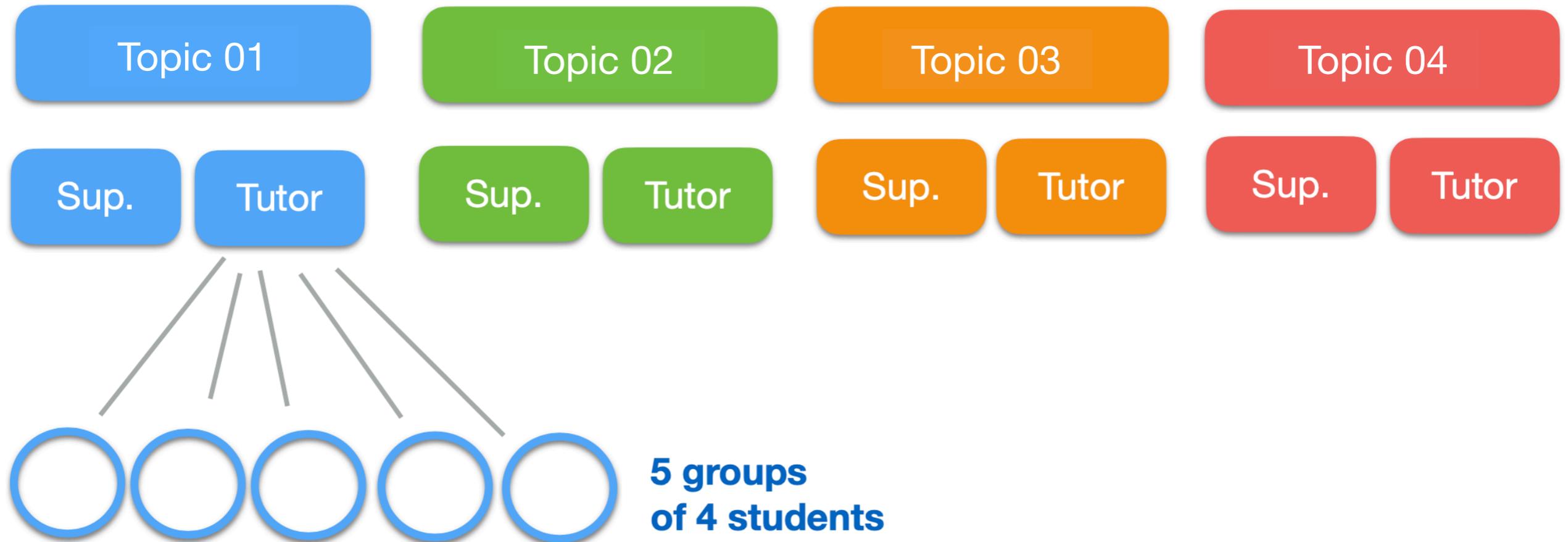
# Konzept

- Fortsetzung der **Vorlesung Datenanalyse** aus dem 3. FS
- **Projekt-orientiertes Lernen:** Erfahrung in der konkreten programmatischen Umsetzung der Methoden der Datenanalyse
- **Ziele**
  - Erfahrung in der **Anwendung der Methoden der Datenanalyse** anhand von reellen Datensätzen und wissenschaftlichen Fragen!
  - Erfahrung in der **Teamarbeit**
  - Erlernen des **Umgangs mit modernen Werkzeugen** der Datenanalyse (R / Python / Markdown / GitHub / ...)

***Learning by doing***

# Themen / Projekte

- **5 Forschungsthemen**
- Für jedes Thema gibt es bis zu 5 sub-Projekte
- Jedes Projekt wird durch ein Team von 4 Studierenden bearbeitet
- **Pro Thema: wissenschaftliche Betreuer und Tutor/Tutorin**
- **Aufgabe der TutorInnen:**
  - wöchentliche Treffen mit den Gruppen die an den Projekten eines Themas arbeiten  
(z.B. Mittwochs 10-13h)
- TutorInnen:  
Hannah Winter; Benedikt Wolf; Luise Nottmeyer; Ana Luisa Costa



R Markdown

from  R Studio

# Themen / Projekte

- **Topic 01: *Biomedical Image Analysis***  
(Karl Rohr / Leonid Kostrykin; Tutor: Hannah Winter)
  - Data types: MNIST images / cell nuclei images
- **Topic 02: *Deep Mutational Scans***  
(Dominik Niopek; Tutor: Benedikt Wolf)
  - Data types: mutation data
- **Topic 03: *Proteome screen***  
(Maiwen Caudron-Herger; Tutor: Fabio Rauscher)
  - Data types: mass spectrometry data
- **Topic 04: *Climate impact on dengue infection***  
(Marina Treskova; tutor Luise Nottmeyer)
  - Data types: Climate data / infection data
- **Topic 05: *Drug repurposing in cancer treatment***  
(Carl Herrmann; tutor Ana Luisa Costa)
  - Data types: Drug screen data / expression data / genetic data

# Zeitplan

*Was wir machen...*

**19/04**

Vorstellung der Projekte  
Intro zu GitHub

**26/04**

Plenum lineare Regression

*Was ihr macht!*

**21/04**

Auswahl der Projekte  
und Teams (Google)

**17/05**

Präsentation des  
project proposal  
(10 + 10 min)

**Schließung der GitHub repos**

**17/07 (8pm)**

**19/07**

Finale Präsentation und  
Bericht  
(15+10 minutes)



# Project proposal (17/05)

- In der Präsentation des **project proposals** solltet ihr...
  - einige der angegebenen Referenzen im Kontext des Projektes vorstellen
  - Die allgemeine Fragestellung / Herausforderungen erklären
  - Die Daten beschreiben
  - Die Ziele eures Projektes definieren
  - Ungefährer Zeitplan
    - ▶ milestones = important steps in the analysis
    - ▶ when these milestones should be achieved
- Mündliche Präsentation vor Betreuer / Tutoren
  - 10 Minuten Präsentation
  - 10 Minuten Diskussion / Fragen
- *Alle Mitglieder des Teams sollten aktiv beitragen!*

# Finale Präsentation und Bericht

I P M B



UNIVERSITÄT  
HEIDELBERG  
ZUKUNFT  
SEIT 1386

- **Finale Präsentation** am 19.07
  - 15 Minuten Präsentation
  - 10 Minuten Fragen
  - Vorstellung der wichtigsten Ergebnisse
- **Bericht (auf Englisch!)**
  - pdf Bericht (im GitHub repo als pdf)
  - **10 Seiten max.**
  - Aufbau: Introduction / material and methods / results/ discussion
  - Wichtig: sorgfältige Auswahl der Plots; Beschriftung der Plots wichtig!
- **Github repo**
  - Repo vor deadline aufräumen !!
  - Bitte ein klares README erstellen

***GitHub repo schließt am Montag 17.07 um 20h  
Bericht sollte bis dahin fertig sein!***

# Benotung

- Project proposal = 30%
- Finale Präsentation = 30%
- Bericht = 40%
- **Kriterien**
  - Qualität der Einführung in die wissenschaftliche Frage (insbesondere Vorstellung der Literatur)
  - Klare Definition der Forschungsfragen
  - Qualität der Plots und Beschriftungen
  - Qualität der Besprechung der Ergebnisse

# Themen / Projekte

- **Topic 01: *Biomedical Image Analysis***  
(Karl Rohr / Leonid Kostrykin; Tutor: Hannah Winter)
  - Data types: MNIST images / cell nuclei images
- **Topic 02: *Deep Mutational Scans***  
(Dominik Niopek; Tutor: Benedikt Wolf)
  - Data types: mutation data
- **Topic 03: *Proteome screen***  
(Maiwen Caudron-Herger; Tutor: Fabio Rauscher)
  - Data types: mass spectrometry data
- **Topic 04: *Climate impact on dengue infection***  
(Marina Treskova; tutor Luise Nottmeyer)
  - Data types: Climate data / infection data
- **Topic 05: *Drug repurposing in cancer treatment***  
(Carl Herrmann; tutor Ana Luisa Costa)
  - Data types: fitness data / ...

# Getting started...

# Auswahl der Projekte / Registrierung



- Vorstellung der Projekte/Themen; Beschreibung der Projekte lesen
- Webseite: <https://www.hdsu.org/teaching/data2023.html>
- Wenn ihr das Projekt und das Team zusammengestellt habt, bitte hier registrieren:

[https://docs.google.com/spreadsheets/d/](https://docs.google.com/spreadsheets/d/1AX28mauyWvg0_7wTtYKxHwl2MN2PVUcK6uyULB8g8uo/edit?usp=sharing)

[1AX28mauyWvg0\\_7wTtYKxHwl2MN2PVUcK6uyULB8g8uo/edit?](https://docs.google.com/spreadsheets/d/1AX28mauyWvg0_7wTtYKxHwl2MN2PVUcK6uyULB8g8uo/edit?usp=sharing)

[usp=sharing](https://docs.google.com/spreadsheets/d/1AX28mauyWvg0_7wTtYKxHwl2MN2PVUcK6uyULB8g8uo/edit?usp=sharing)

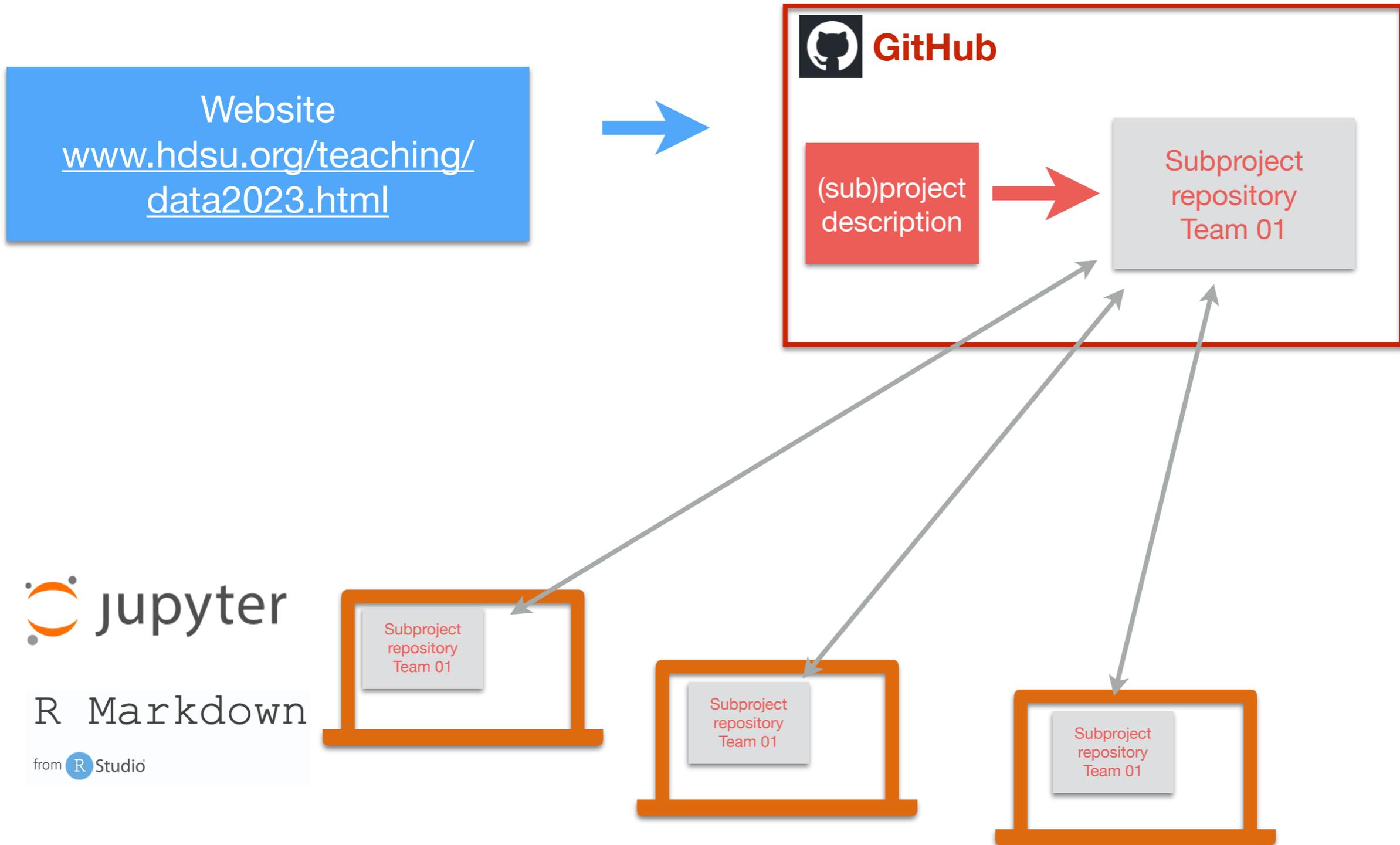
- ***Auswahl der Projekte***
  - ***startet Freitag 21.04 um 13 h***
  - ***Sollte bis Montag 24.04 abgeschlossen sein.***

# Wie sollte man den/die TutorIn (nicht) benutzen?

- Zeit für das wöchentliche Treffen ist **beschränkt (30 Minuten/Woche)**;  
→ nutzt die Zeit sinnvoll!!
- TutorInnen sind NICHT dazu da, euer Code zu debuggen!
- Bereitet das wöchentliche Treffen sorgfältig vor:
  - was haben wir seit letzter Woche erreicht?
  - welche Probleme/Fragen haben wir?
  - Ziele für die kommende Woche

***TutorInnen sind nicht dazu da, WhatsApp Nachrichten um Mitternacht zu beantworten...***

# Arbeit organisieren



# Für R-Projekte

- Es gibt viele tolle Software Libraries in R (Bioconductor / CRAN / ...)
- **Bitte keine Methoden/Tools benutzen, die ihr nicht absolut versteht!**
- Uns sind einfach Lösungen "per Hand" lieber als "black-boxes"!
- Wenn ihr externe libraries benutzt solltet ihr auch imstande sein, diese genau zu erklären!!

# Brief intro to Git(Hub)

# Git(Hub)

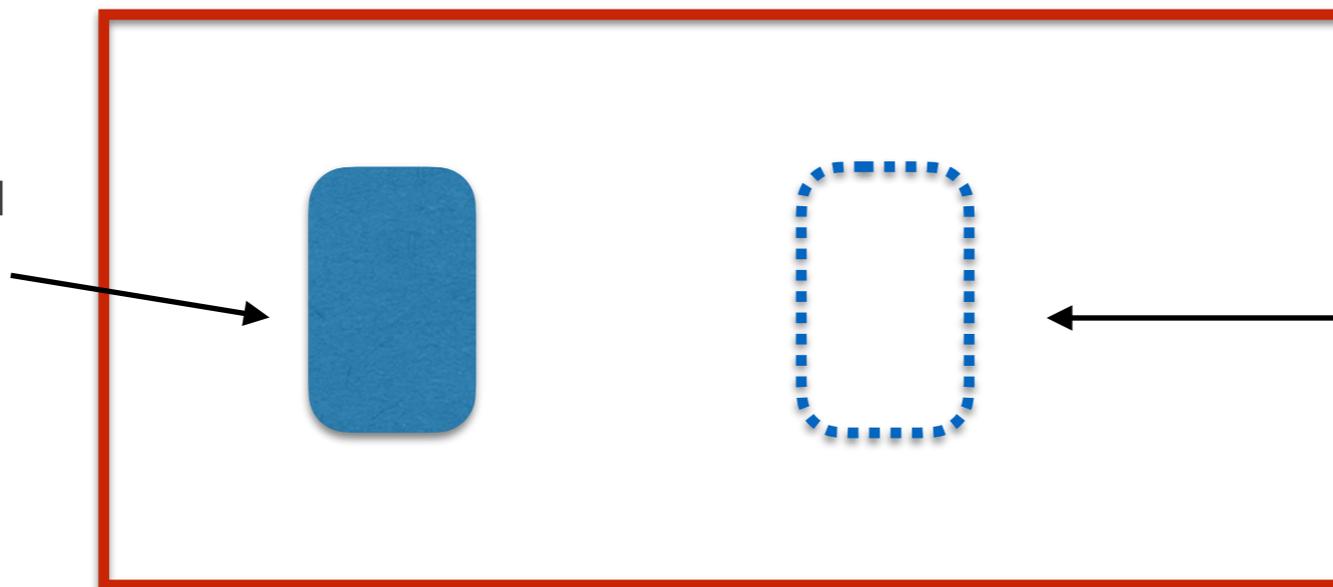
- Git is a **version control system**:
  - allows simultaneous work of different people on the same project
  - tracks the changes (**‘commits’**) made by each member
  - helps solve the **conflicts** between various versions
- GitHub is a platform which hosts Git projects (‘repositories’)
  - is free to use
  - required to create a (free) account
  - can be used in command line or using GUI tools (‘GitHub Desktop’)

# Git repository

## repository



this file is registered  
in the database  
(‘committed’)



this file exists, but  
has not been  
registered yet  
OR  
the file has been modified, but  
the changes have not been  
registered yet

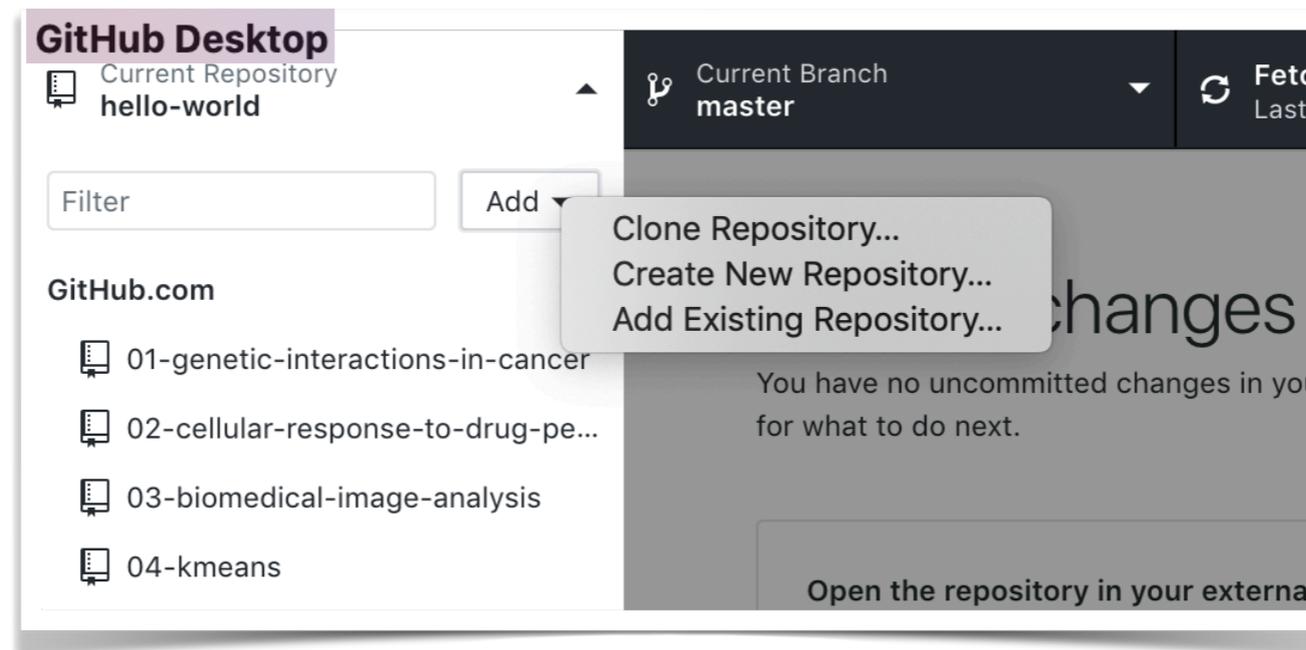
# 1. Cloning an existing repository

User's computer



 **GitHub**

Team's repository



Current Repository **hello-world** | Current Branch **master** | Fetch origin (Last fetched 19 minutes ago)

Changes | History

0 changed files

## No local changes

You have no uncommitted changes in your repository! Here are some friendly suggestions for what to do next.

- Open the repository in your external editor**  
Configure which editor you wish to use in [preferences](#) | [Open in Visual Studio Code](#)  
Repository menu or ⌘ ↑ A
- View the files in your repository in Finder**  
Repository menu or ⌘ ↑ F | [Show in Finder](#)
- Open the repository page on GitHub in your browser**  
Repository menu or ⌘ ↑ G | [View on GitHub](#)

sd Summary (required)

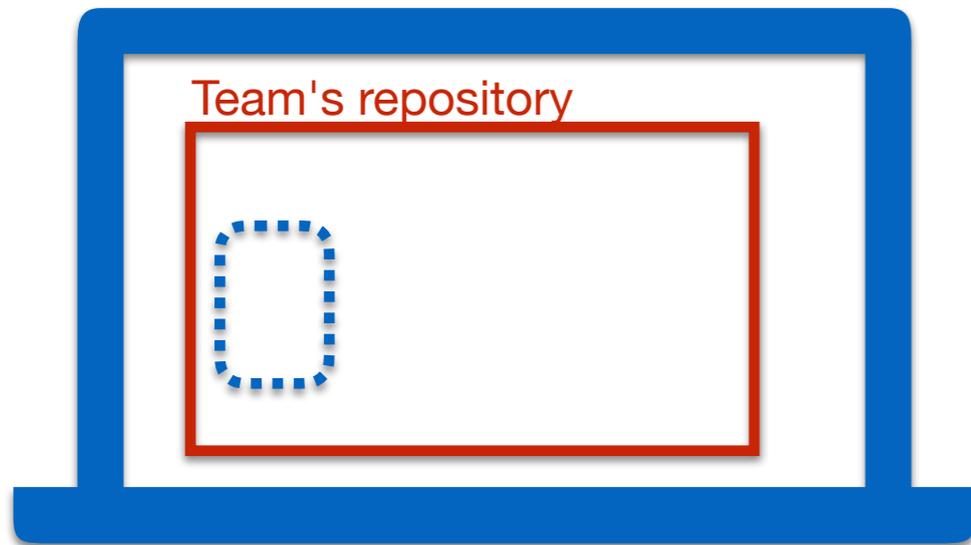
Description

+

[Commit to master](#)

## 2. creating a local file

### User's computer



### Team's repository



- When a new file is added / modified in the local folder, it is not yet registered in the git database!
- it first needs to be **committed**

Current Repository: **hello-world**
Current Branch: **master**
Fetch origin  
Last fetched 22 minutes ago

Changes **1**
History
my\_markdown.Rmd +

1 changed file

- my\_markdown.Rmd +

  
*new file created locally*

 Create my\_markdown.Rmd

Description

 +

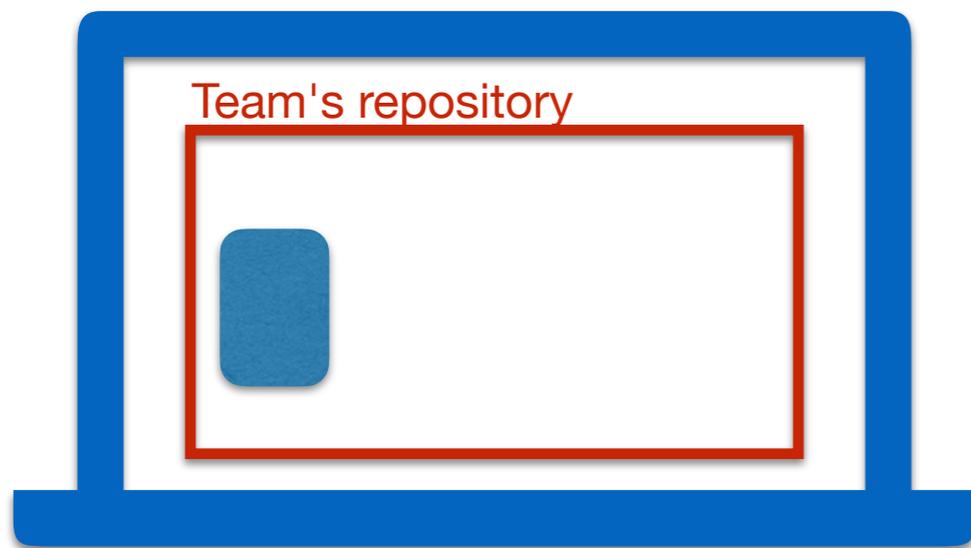
Commit to master

```

@@ -0,0 +1,30 @@
1 +---
2 +title: "My first markdown"
3 +author: "Carl Herrmann"
4 +date: "4/23/2019"
5 +output: html_document
6 +---
7 +
8 +```{r setup, include=FALSE}
9 +knitr::opts_chunk$set(echo = TRUE)
10 +```
11 +
12 +## R Markdown
13 +
14 +This is an R Markdown document. Markdown is a simple formatting syntax for authoring
15 +HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.
16 +
17 +When you click the Knit button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:
18 +
19 +```{r cars}
20 +summary(cars)
21 +```
22 +
23 +## Including Plots
24 +
25 +You can also embed plots, for example:
          
```

## 2. Adding a file

### User's computer



### Team's repository



- When a new file is added / modified in the local folder, it is not yet registered in the git database!
- it first needs to be **committed**

Current Repository: **hello-world**
Current Branch: **master**
Fetch origin  
Last fetched 22 minutes ago

Changes <span style="background-color: #e6e6e6;">1</span>	History	my_markdown.Rmd <span style="float: right;">+</span>
1 changed file		
<input checked="" type="checkbox"/> my_markdown.Rmd <span style="float: right;">+</span>		<pre> @@ -0,0 +1,30 @@ 1 +--- 2 +title: "My first markdown" 3 +author: "Carl Herrmann" 4 +date: "4/23/2019" 5 +output: html_document 6 +--- 7 + 8 +```{r setup, include=FALSE} 9 +knitr::opts_chunk\$set(echo = TRUE) 10 +``` 11 + 12 +## R Markdown 13 + 14 +This is an R Markdown document. Markdown is a simple formatting syntax for authoring 15 +HTML, PDF, and MS Word documents. For more details on using R Markdown see &lt;http://rmarkdown.rstudio.com&gt;. 16 + 17 +When you click the <b>Knit</b> button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this: 18 + 19 +```{r cars} 20 +summary(cars) 21 +``` 22 + 23 +## Including Plots 24 + 25 +You can also embed plots, for example:                 </pre>

*indicate the type of changes made and commit*

↓

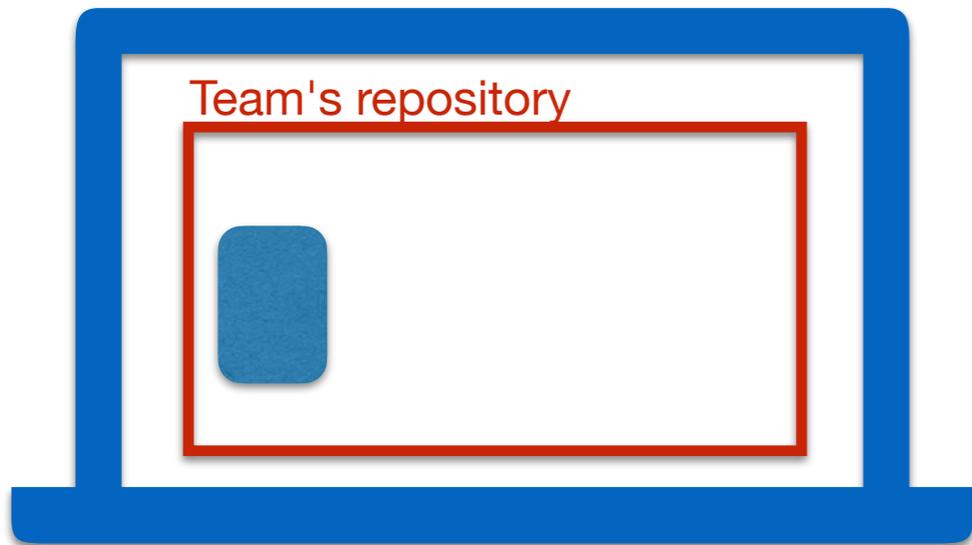
Description

+ Add people

Commit to master

# 2. Adding a file

## User's computer



## Team's repository



- the file is now committed to the local git repository
- it needs to be pushed to the remote repository on GitHub

Current Repository **hello-world** | Current Branch **master** | Push origin (1 ↑) Last fetched 30 minutes ago

Changes | History

0 changed files

# No local changes

You have no uncommitted changes in your repository! Here are some friendly suggestions for what to do next.

**Push 1 commit to the origin remote**  
You have one local commit waiting to be pushed to GitHub  
Always available in the toolbar when there are local commits waiting to be pushed or ⌘ P Push origin

**Open the repository in your external editor**  
Configure which editor you wish to use in [preferences](#) Open in Visual Studio Code  
Repository menu or ⌘ ↑ A

**View the files in your repository in Finder** Show in Finder  
Repository menu or ⌘ ↑ F

**Open the repository page on GitHub in your browser** View on GitHub  
Repository menu or ⌘ ↑ G

Summary (required)

Description

Commit to master

Committed just now  
Create my\_markdown.Rmd Undo

## test repository

Edit

[Manage topics](#)

2 commits

1 branch

0 releases

1 contributor

Branch: master ▾

New pull request

Create new file

Upload files

Find File

Clone or download ▾

 **carlherrmann** Create my\_markdown.Rmd

Latest commit b2b0bdf 2 minutes ago

 [README.md](#)

Initial commit

37 minutes ago

 [my\\_markdown.Rmd](#)

Create my\_markdown.Rmd

2 minutes ago

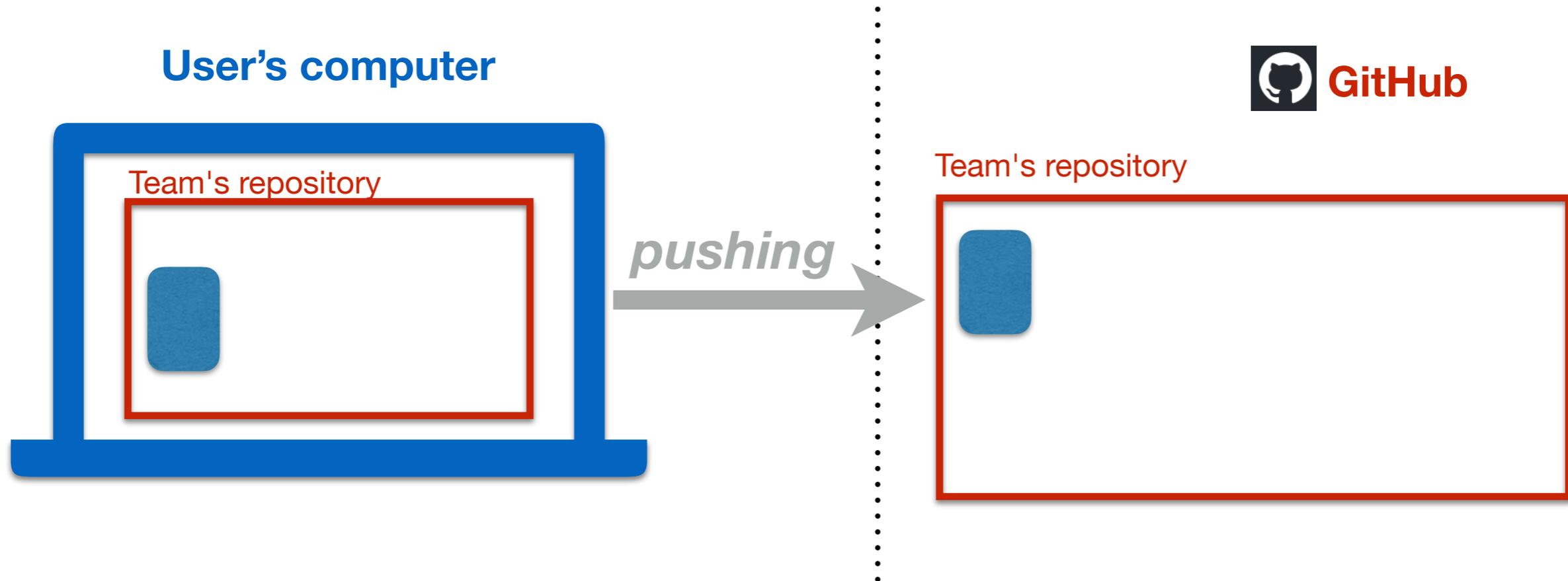
 [README.md](#)



# hello-world

test repository

## 2. Adding a file



- the file is not committed to the local git repository
- it needs to be pushed to the remote repository on GitHub

# 3. Pulling from the remote repository



- Someone (probably one of your team mates) has added a new file into the remote repository
- It is not yet in your local repository and need to be **pulled**

test repository Edit

[Manage topics](#)

3 commits    1 branch    0 releases    1 contributor

Branch: master ▾    [New pull request](#)    [Create new file](#)    [Upload files](#)    [Find File](#)    [Clone or download ▾](#)

Author	Message	Time
carlherrmann	Newly created file	Latest commit 4cff9b2 2 minutes from now
	Initial commit	39 minutes ago
	Create my_markdown.Rmd	4 minutes ago
	Newly created file	just now

[new\\_file\\_added](#)

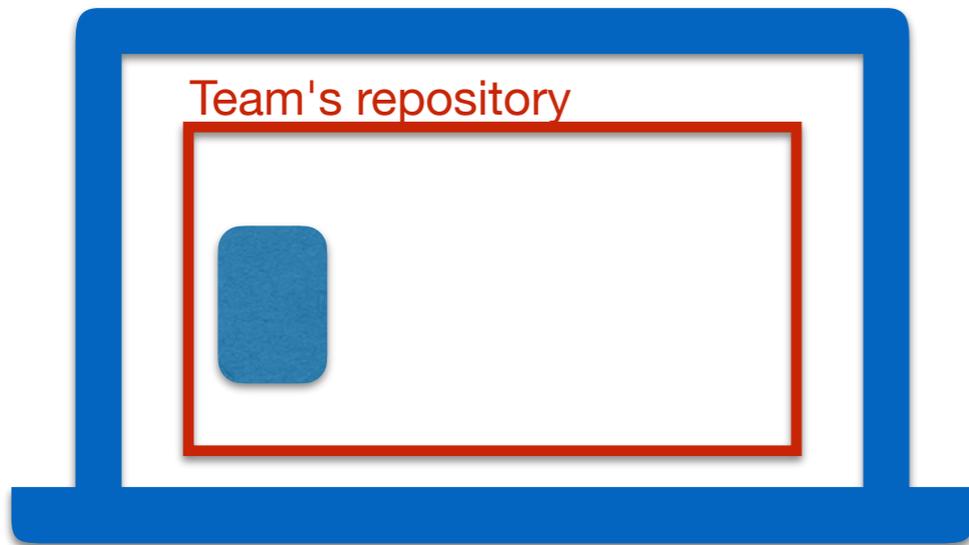
[README.md](#)

## hello-world

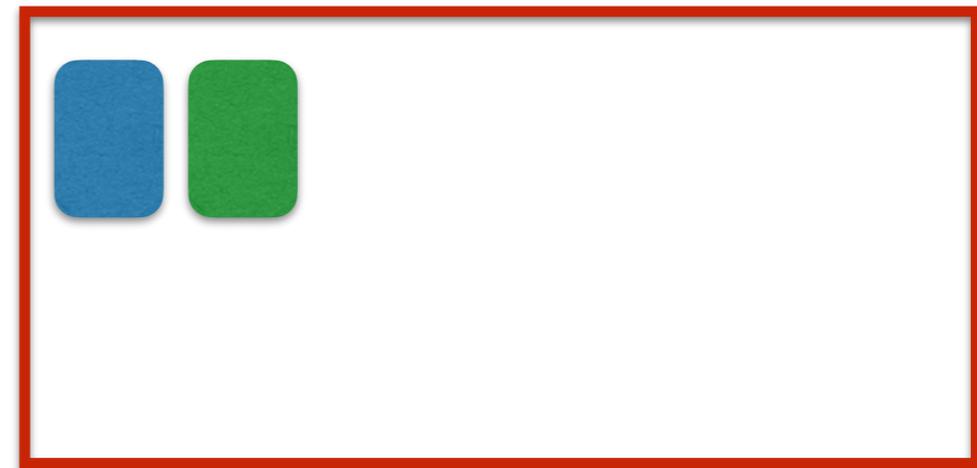
test repository

# 3. Pulling from the remote repository

## User's computer



## Team's repository



Current Repository **hello-world** | Current Branch **master** | **Fetch origin** (Last fetched 4 minutes ago)

Changes | History

0 changed files

# No local changes

You have no uncommitted changes in your repository! Here are some friendly suggestions for what to do next.

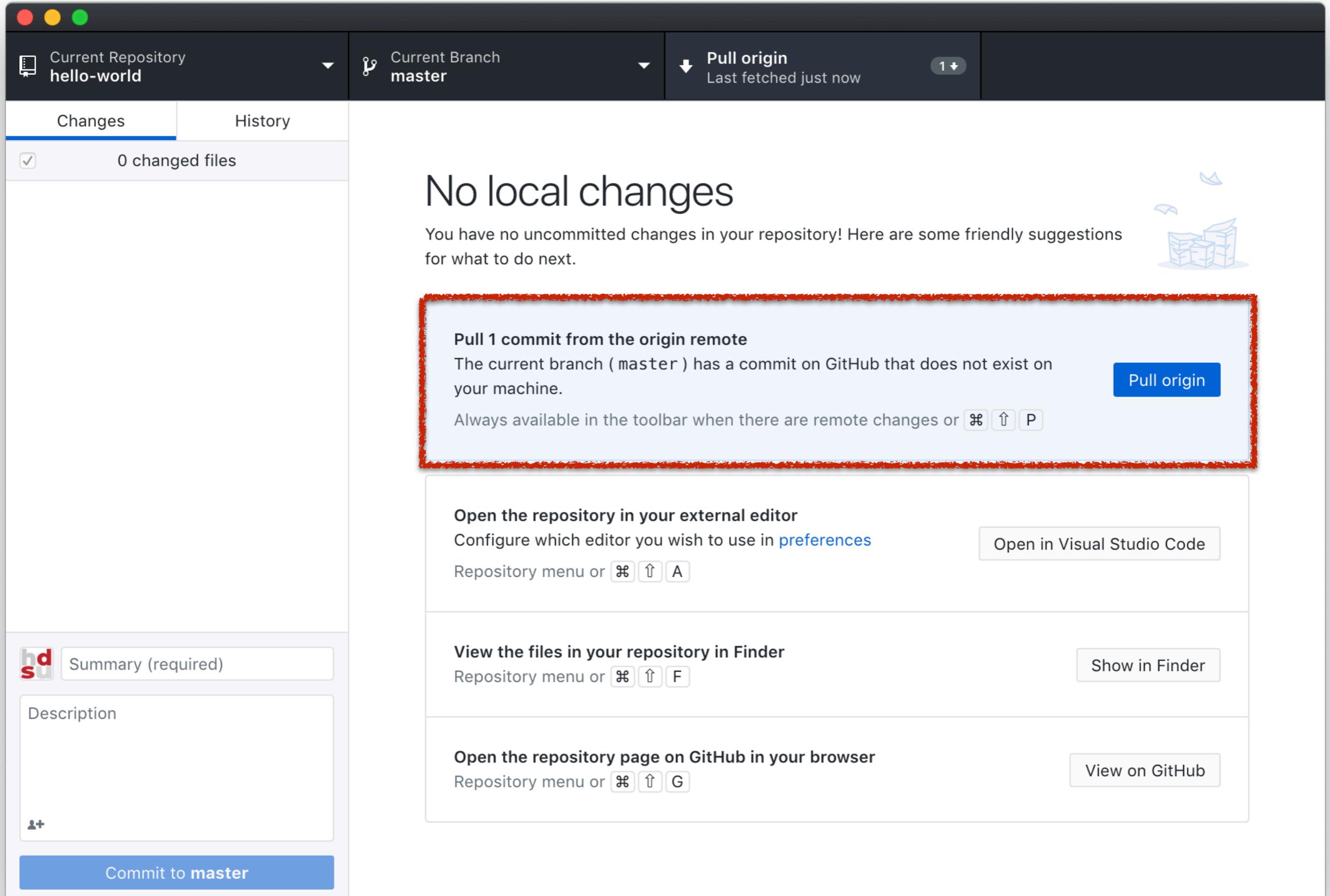
- Open the repository in your external editor**  
Configure which editor you wish to use in [preferences](#) | [Open in Visual Studio Code](#)  
Repository menu or `⌘ ↑ A`
- View the files in your repository in Finder**  
Repository menu or `⌘ ↑ F` | [Show in Finder](#)
- Open the repository page on GitHub in your browser**  
Repository menu or `⌘ ↑ G` | [View on GitHub](#)

sd Summary (required)

Description

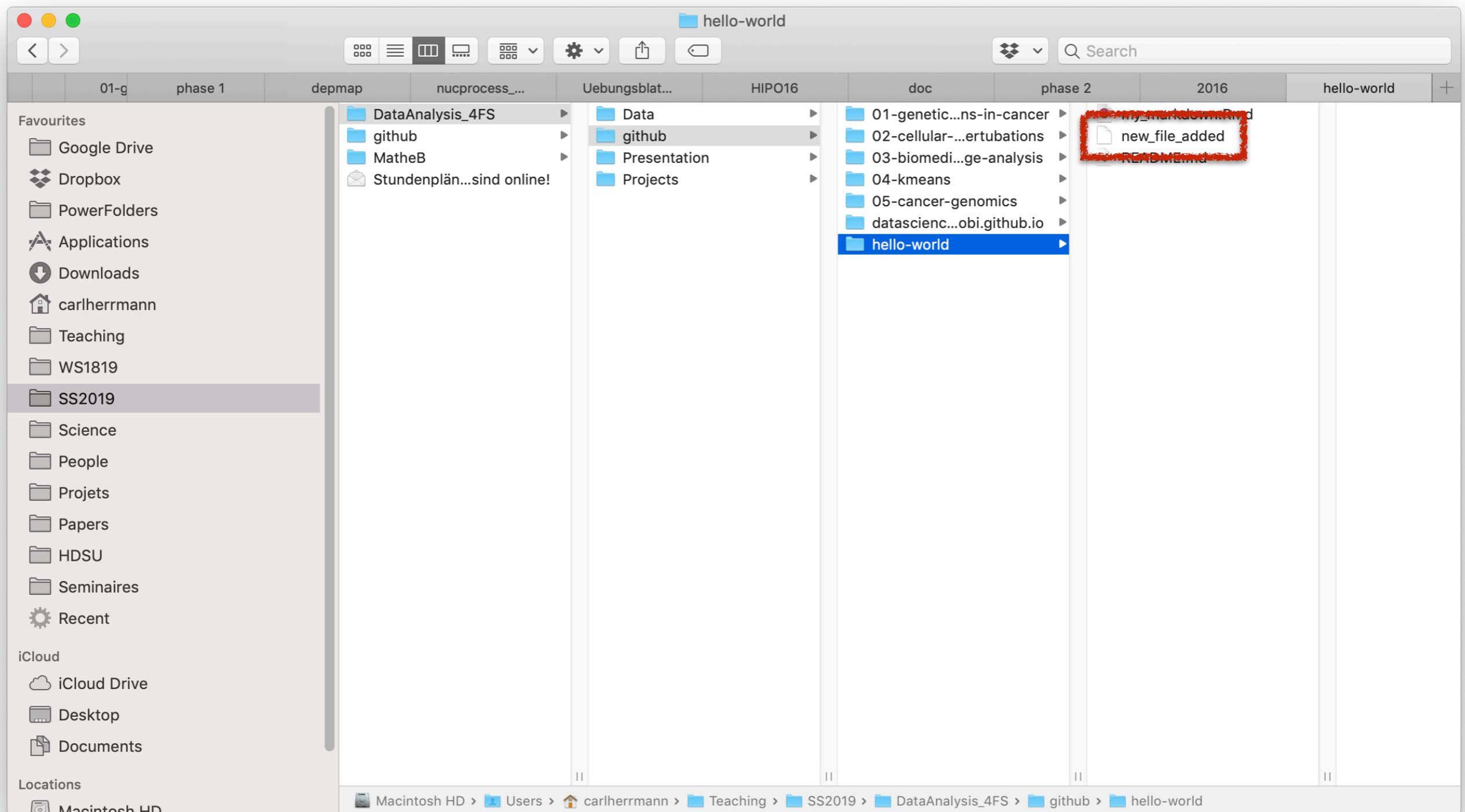
[+](#)

[Commit to master](#)

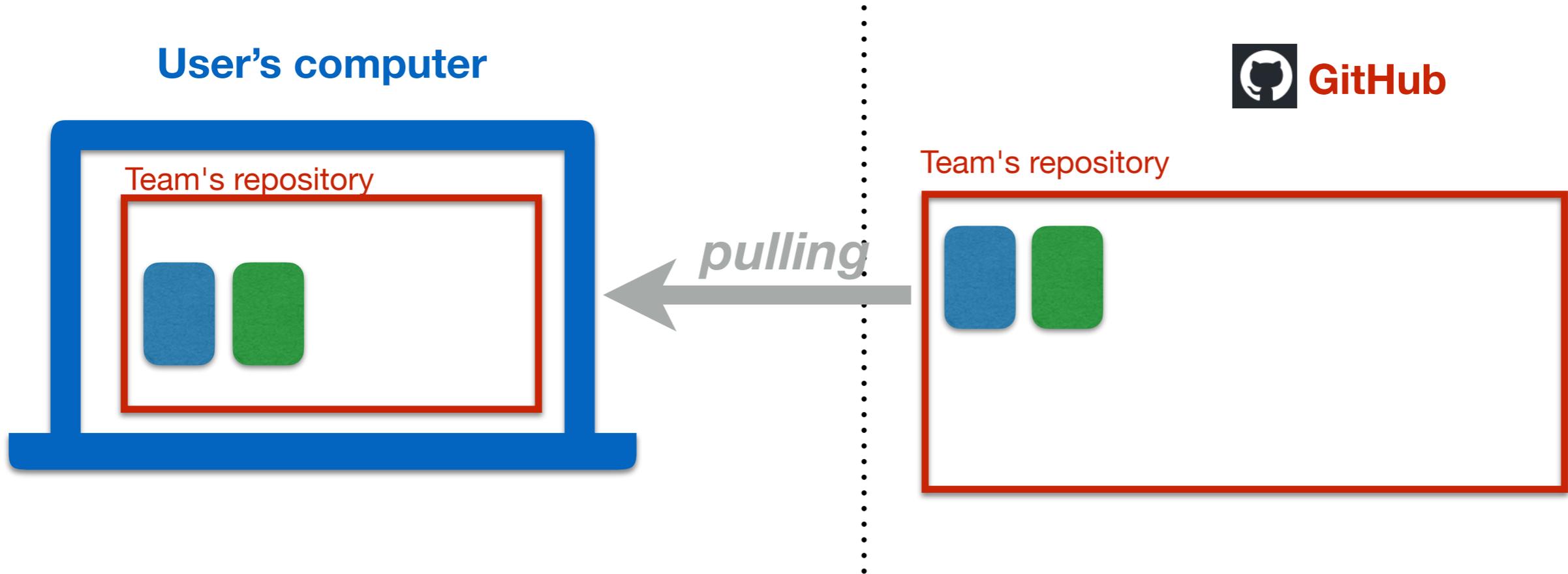


The screenshot shows the GitHub Desktop application interface. At the top, the 'Current Repository' is 'hello-world', the 'Current Branch' is 'master', and there is a 'Pull origin' button with a '1' icon and the text 'Last fetched just now'. Below the toolbar, there are two tabs: 'Changes' and 'History'. The 'Changes' tab is active and shows '0 changed files'. The main area displays 'No local changes' with a message: 'You have no uncommitted changes in your repository! Here are some friendly suggestions for what to do next.' A red dashed box highlights a suggestion: 'Pull 1 commit from the origin remote'. The text below this suggestion states: 'The current branch ( master ) has a commit on GitHub that does not exist on your machine.' and 'Always available in the toolbar when there are remote changes or ⌘ ↑ P'. A blue 'Pull origin' button is present. Other suggestions include 'Open the repository in your external editor' (with 'Open in Visual Studio Code' button), 'View the files in your repository in Finder' (with 'Show in Finder' button), and 'Open the repository page on GitHub in your browser' (with 'View on GitHub' button). On the left sidebar, there is a 'Summary (required)' field, a 'Description' field, and a 'Commit to master' button at the bottom.

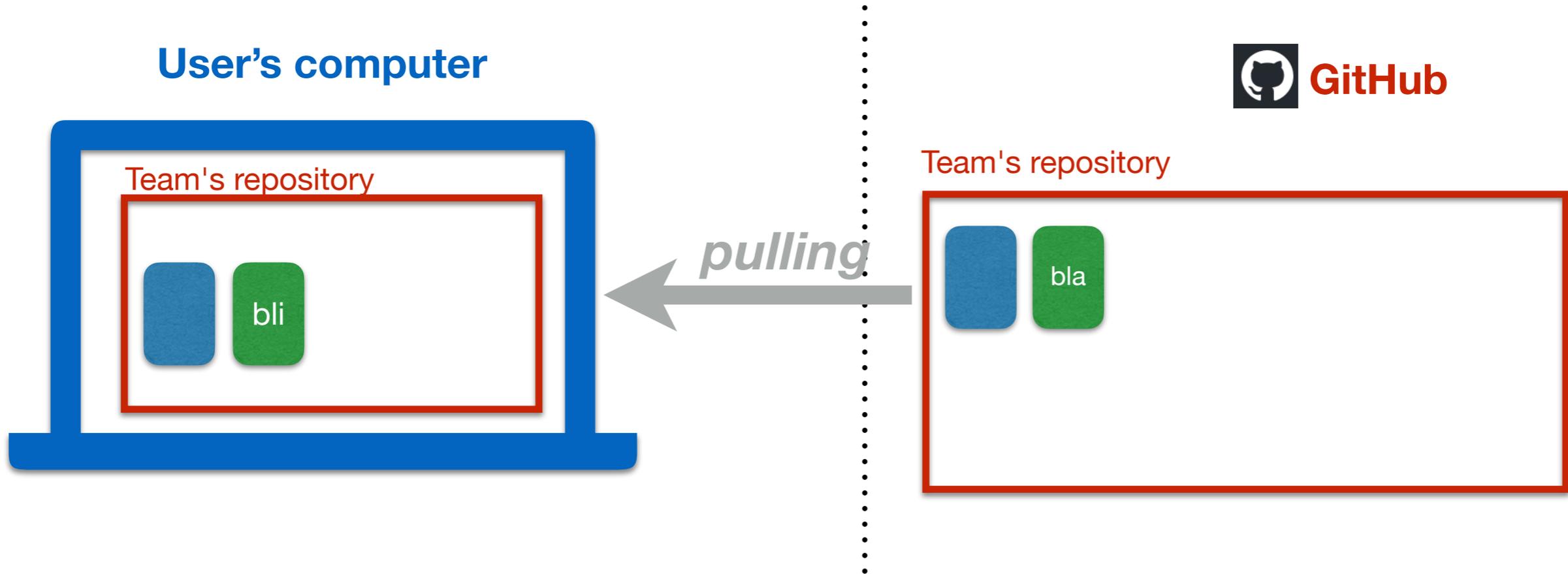
- Once the remote repository is pulled, the new file(s) are available locally

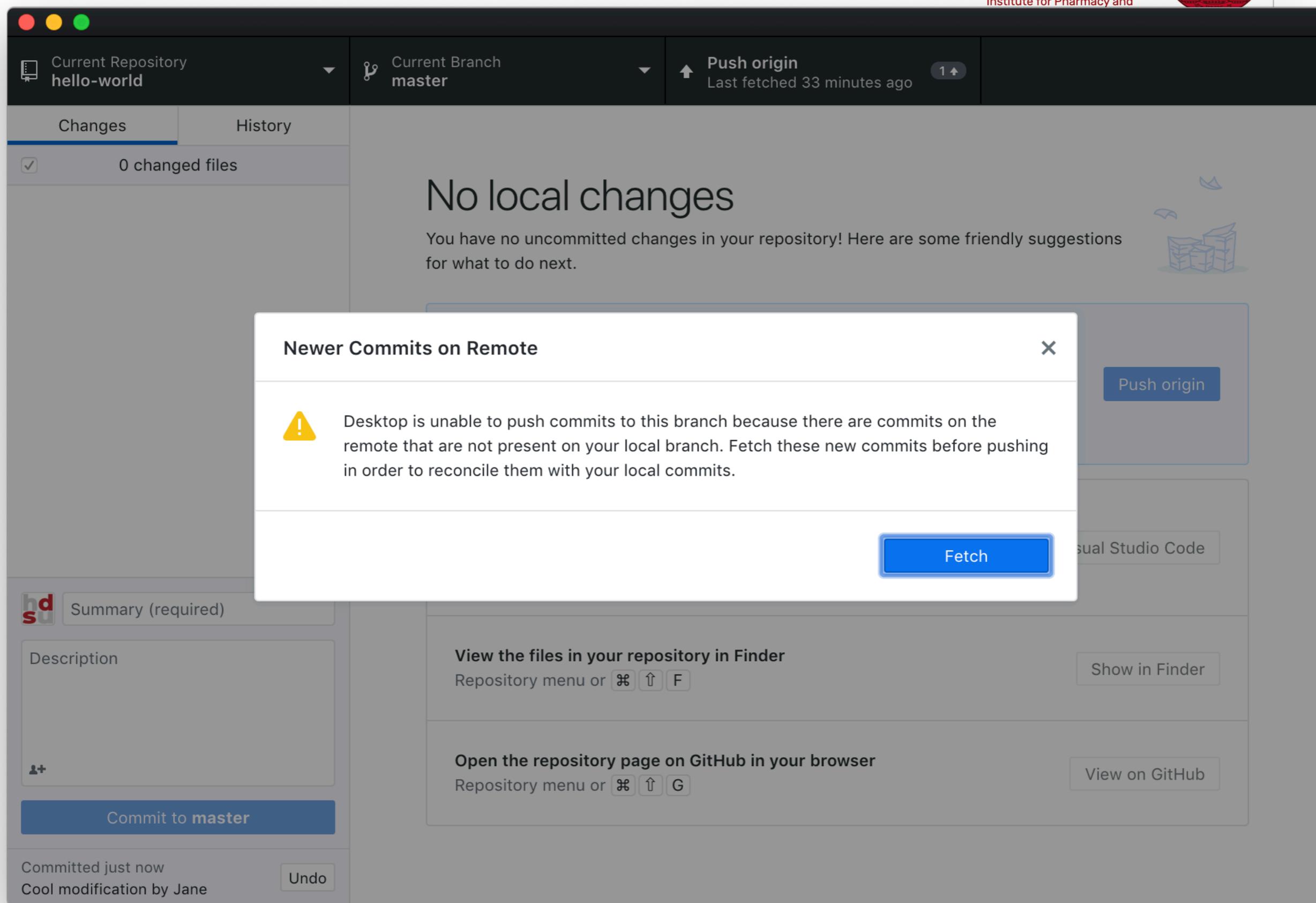


# 3. Pulling from the remote repository

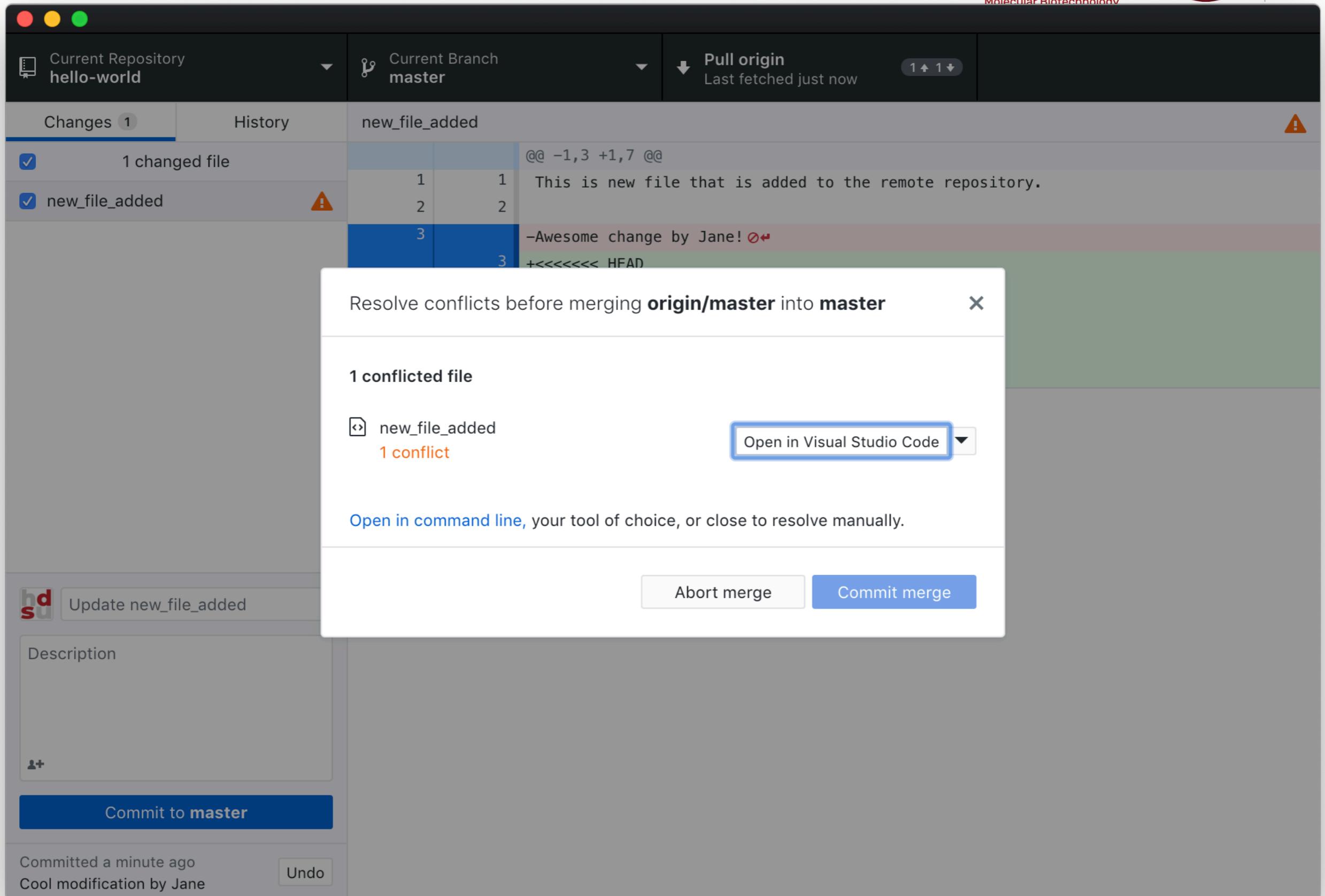


# 4. conflicting changes





The screenshot shows a Git GUI interface for a repository named 'hello-world' on the 'master' branch. The main area displays 'No local changes' and suggests actions like 'Push origin'. A modal dialog titled 'Newer Commits on Remote' is open, containing a yellow warning icon and the text: 'Desktop is unable to push commits to this branch because there are commits on the remote that are not present on your local branch. Fetch these new commits before pushing in order to reconcile them with your local commits.' A blue 'Fetch' button is visible at the bottom right of the dialog. The background interface includes a 'Commit to master' button, a 'Push origin' button, and a 'Show in Finder' button.



The screenshot shows a Git GUI interface with a merge conflict resolution dialog. The background interface includes:

- Current Repository: hello-world
- Current Branch: master
- Pull origin: Last fetched just now
- Changes: 1 changed file
- new\_file\_added: 1 changed file
- Diff view showing a conflict in 'new\_file\_added'.
- Buttons: Update new\_file\_added, Commit to master, Undo.

The foreground dialog is titled "Resolve conflicts before merging origin/master into master". It lists the conflicted file "new\_file\_added" with "1 conflict". A button "Open in Visual Studio Code" is highlighted with a blue box. Below the dialog, there are "Abort merge" and "Commit merge" buttons.

# 4. conflicting changes

- Conflicting changes can be resolved with a text editor
- options depend on which editor is used

```

You, a few seconds ago | 2 authors (You and others)
This is new file that is added to the remote repository.

Accept Current Change | Accept Incoming Change | Accept Both Changes | Compare Changes
<<<<<<< HEAD (Current Change)
Awesome change by Jane!
=====
This is a great new modification by Joe!
>>>>>>> af5e9c9981b21a39ed11d09f468bea576d669191 (Incoming Change)
  
```

local change

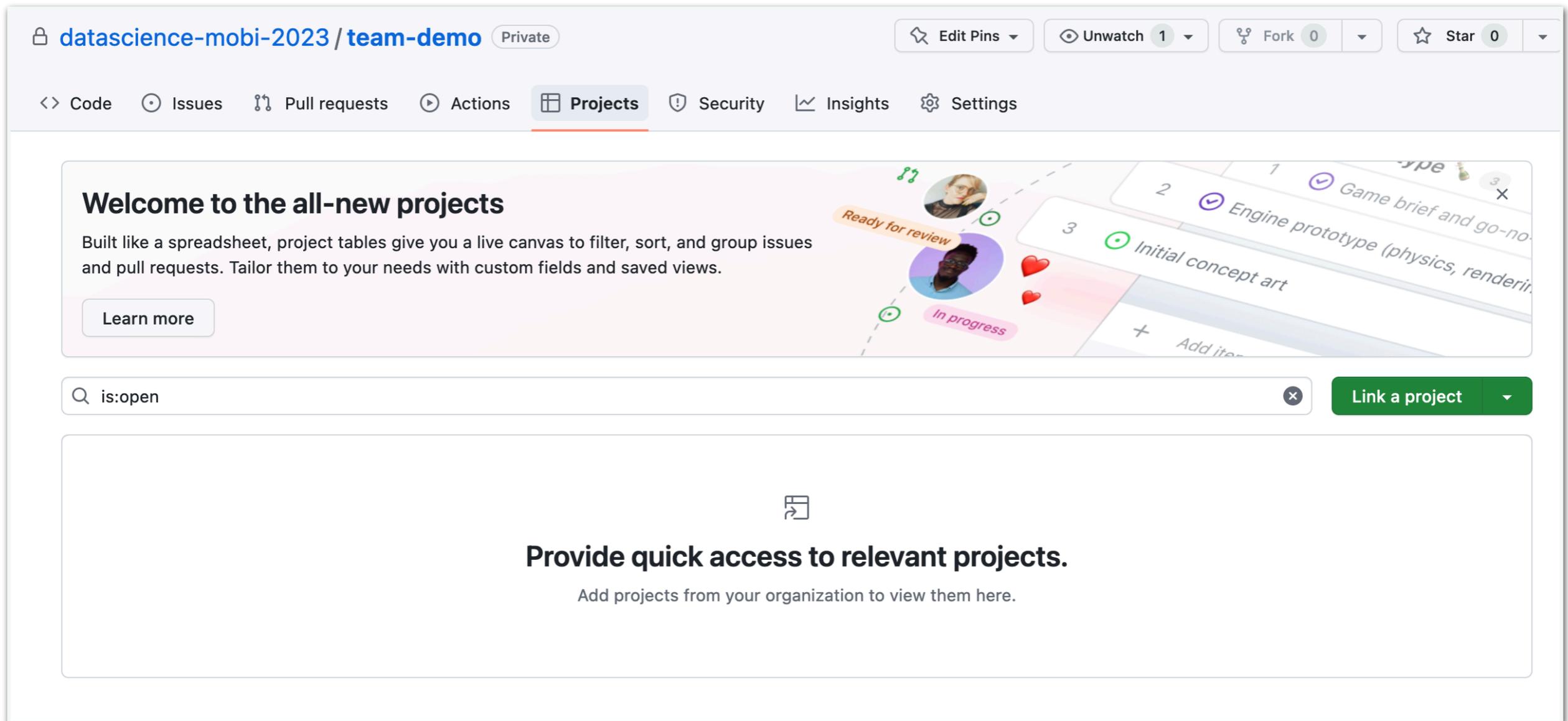
changes in the  
remote file

# To do

- Create your own personal GitHub account
- Register your Github user name into the Google Sheet
- all team members will be added to the corresponding GitHub repo
  - Project 03 - Team 02 → **project-03-group-02**

# Using projects

- You can define tasks using the GitHub project tool

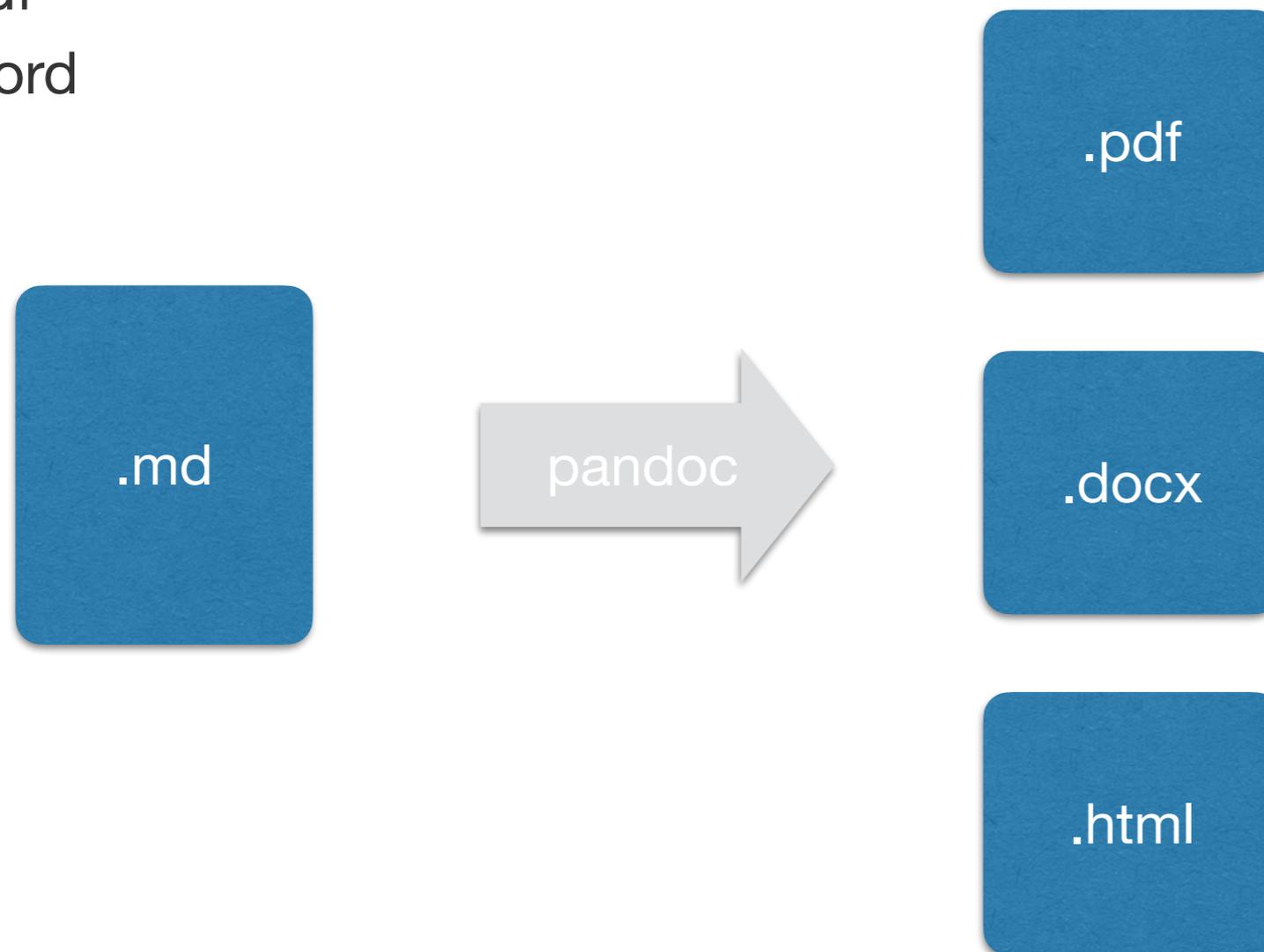


The screenshot shows the GitHub Projects interface for a repository named 'datascience-mobi-2023 / team-demo'. The page is set to 'Private'. At the top, there are navigation options: 'Code', 'Issues', 'Pull requests', 'Actions', 'Projects' (which is selected and highlighted), 'Security', 'Insights', and 'Settings'. Below the navigation, there are interaction buttons: 'Edit Pins', 'Unwatch 1', 'Fork 0', and 'Star 0'. The main content area features a 'Welcome to the all-new projects' section with a description: 'Built like a spreadsheet, project tables give you a live canvas to filter, sort, and group issues and pull requests. Tailor them to your needs with custom fields and saved views.' A 'Learn more' button is provided. Below this is a search bar with the filter 'is:open' and a 'Link a project' button. The bottom section is a large empty box with a folder icon and the text: 'Provide quick access to relevant projects. Add projects from your organization to view them here.'

# (R)markdown

# Markdown

- Markdown is a way to format plain text with a simple text editor
- Markdown documents can be converted with a **renderer** into
  - html
  - pdf
  - word



# Rendering markdown

markdown

```
# My document

## this is a header

In the text we can highlight or put in bold.

## making lists

We can make numbered lists:

1. first item
2. second item

or unordered lists

* first item
* second item
  + subitem
  + subitem
* third item

This is code which can be put inline

```bash
this is bash code
```

```python
this is python code
```
```

pdf

My document

this is a header

In the text we can *highlight* or put in **bold**.

making lists

We can make **numbered lists**:

1. first item
2. second item

or unordered lists

- first item
- second item
- subitem
- subitem
- third item

This is `code` which can be put inline

this is bash code

this **is** python code

html

**My document**

**this is a header**

In the text we can *highlight* or put in **bold**.

**making lists**

We can make **numbered lists**:

1. first item
2. second item

or unordered lists

- first item
- second item
- subitem
- subitem
- third item

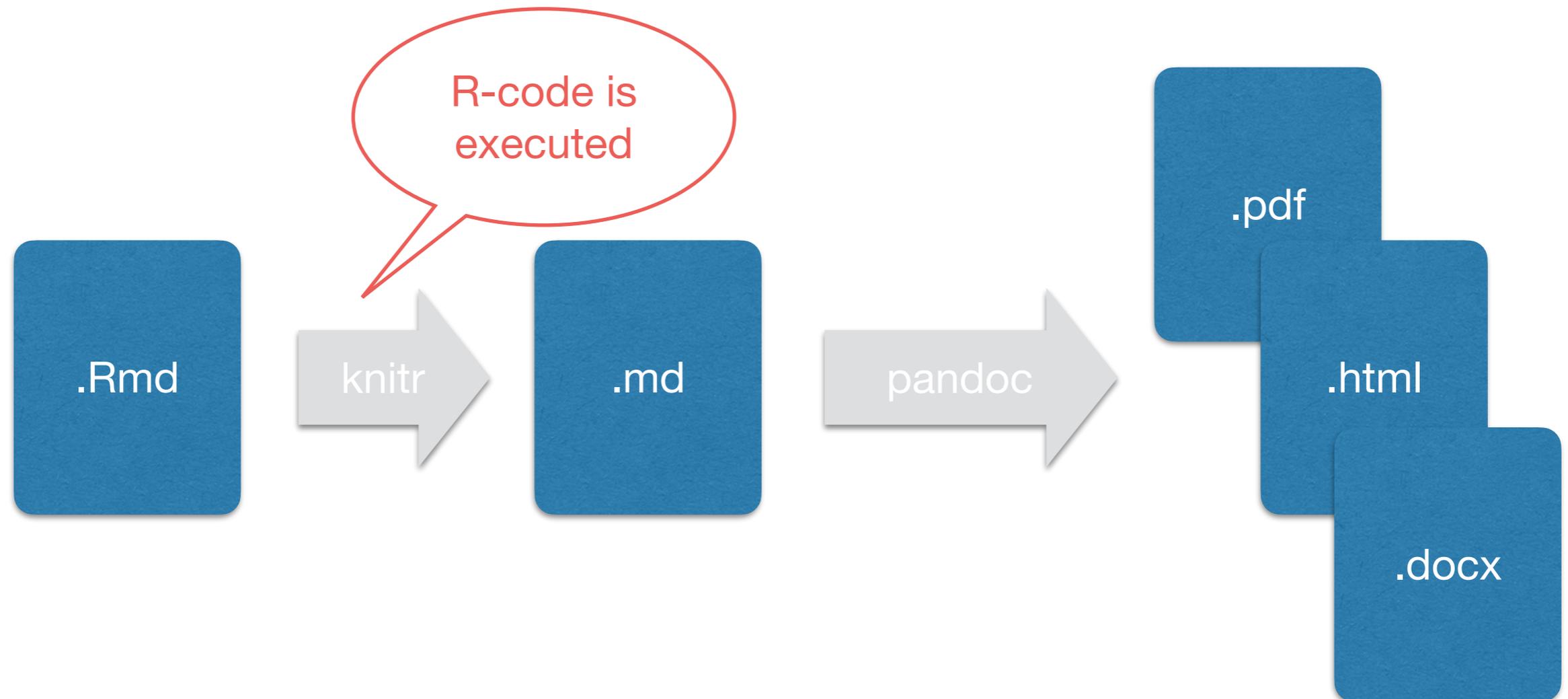
This is `code` which can be put inline

this is bash code

this is python code

# Rmarkdown

- With Rmarkdown, R-code parts can be included into the markdown document
- the R-code will be executed, the result integrated into markdown



# Rmarkdown format

```
---  
title: "Project 01"  
author: "Carl Herrmann"  
date: "4/17/2019"  
output:  
  html_document:  
  |   keep_md: yes  
  pdf_document: default  
---  
# A Rmarkdown tutorial  
  
This is a brief tutorial on how to use Rmarkdown to create dynamic documents  
  
```${r setup, include=FALSE}  
knitr::opts_chunk$set(echo = TRUE)  
knitr::opts_knit$set(root.dir='/Users/carlherrmann/Teaching/SS2019/DataAnalysis_4FS')  
````  
  
## Load the dataset  
  
```${r read_data}  
allDepMapData = readRDS('Data/depmap/DepMap19Q1_allData.RDS')  
````  
  
Now plot the distribution of the cell lines according to the tissue type  
  
```${r plot_data}  
freq = sort(table(allDepMapData$annotation$Primary.Disease))  
par(las=2,mar=c(3,8,3,3));barplot(freq,horiz=TRUE, col='lightgrey')  
````
```

header: set options

R code chunks

text in markdown

# Rmarkdown chunk options

- Display options can be set for each chunk individually, or for all chunks at the beginning of the document

```
```${r setup, include=FALSE}  
knitr::opts_chunk$set(echo = TRUE)  
knitr::opts_chunk$set(cache = TRUE)
```

valid for all chunks

- echo=TRUE : R-code is displayed in final document
- cache = TRUE : results of all chunks are cached

```
```${r plot_data, fig.height=12, fig.width=12}  
freq = sort(table(allDepMapData$annotation$Primary.Disease))  
par(las=2, mar=c(3, 8, 3, 3)); barplot(freq, horiz=FALSE, col='lightgrey')  
````
```

valid for **this** chunks

- set height and width of output figure

# Reference

- <https://rmarkdown.rstudio.com/>
- <https://www.rstudio.com/wp-content/uploads/2016/03/rmarkdown-cheatsheet-2.0.pdf>