# Reproducible science

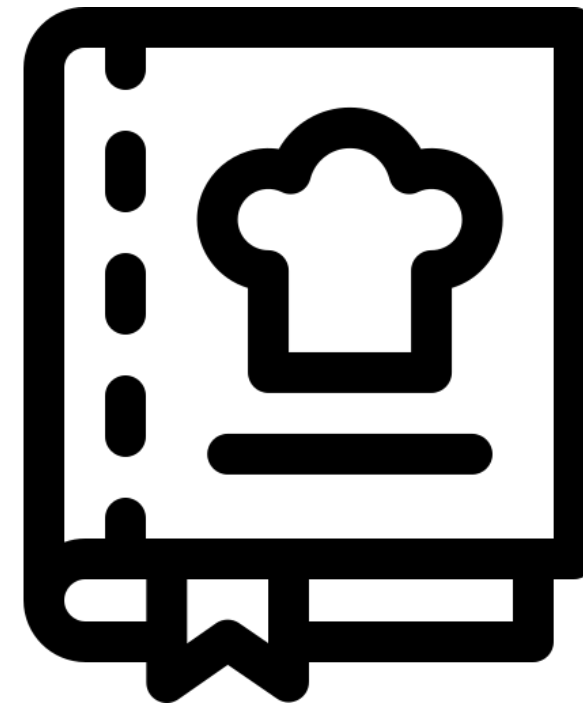## Tools and good practices

Viola Fanfani

January 2022

# Disclaimers

1. If you already know this stuff, I am sorry you have to hear about it again and please correct me. I am open to incorporate more/other things

2. These are things I have used and they have been helpful. We can mix and match, adopt things collectively or singularly, reshape structure according to the needs of the lab

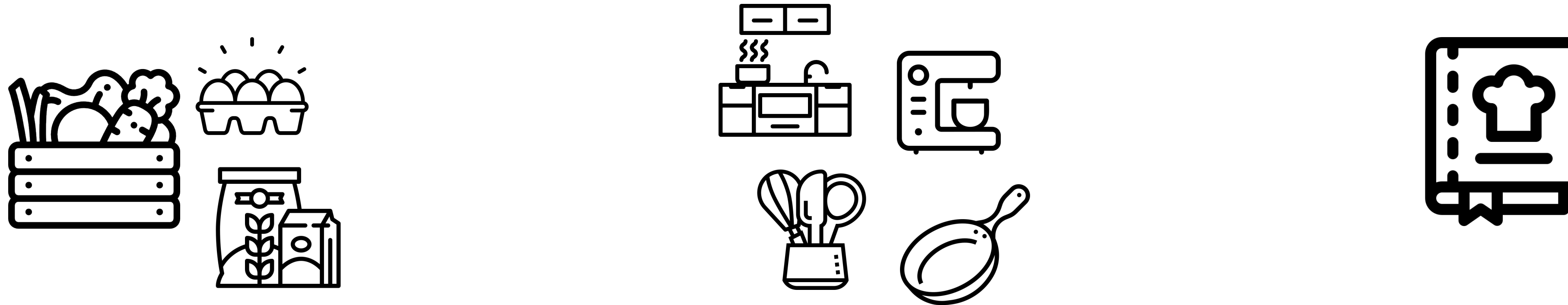3. Takes a moment to have everything up and running, but it saves time and headaches

# Digital project management

**What reproducible science looks like 90% of the time:**

# Digital project management

## What reproducible science should look like:

# Digital project management

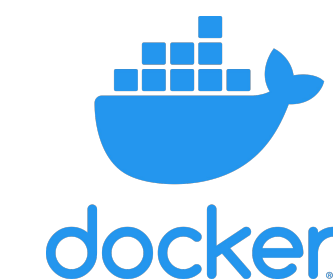It needs to be easy to prepare, to use, to distribute, and to reuse!

Think about a colleague (or yourself) who will work on your stuff 5 years from now
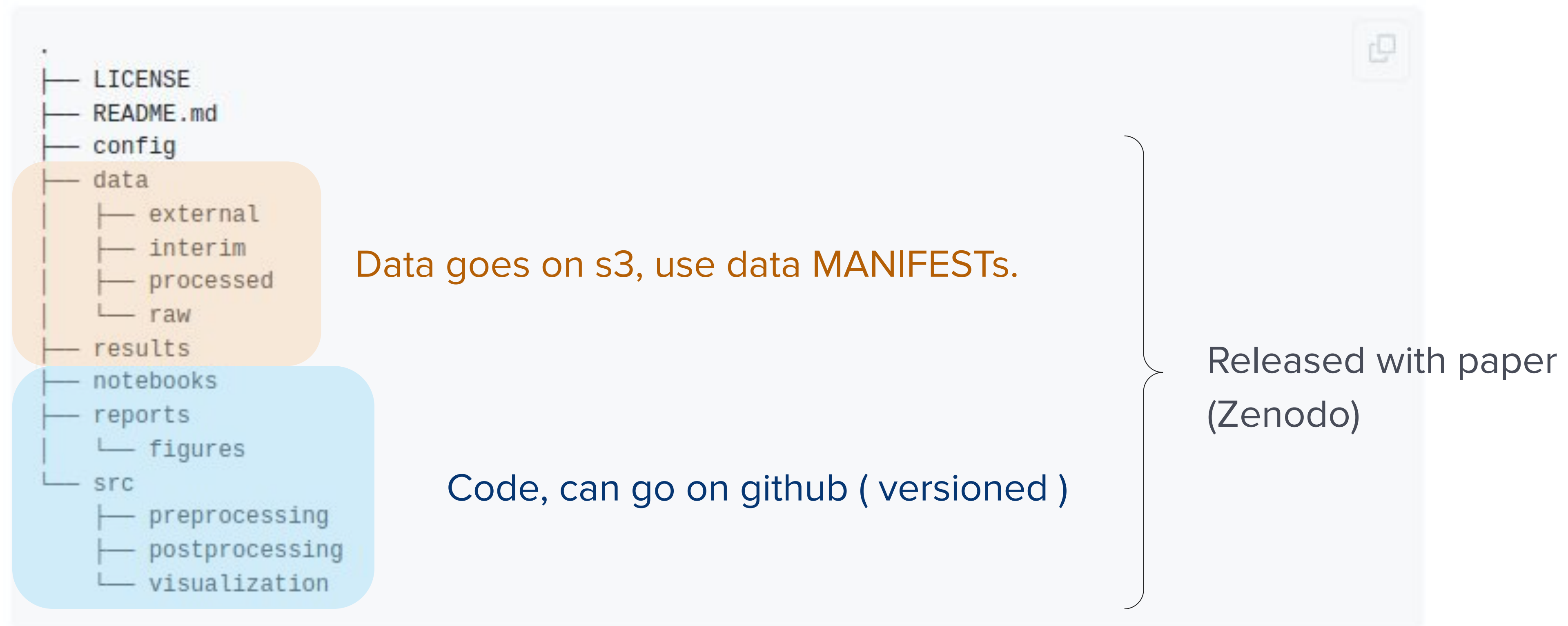
**Data:**

- safe
- portable
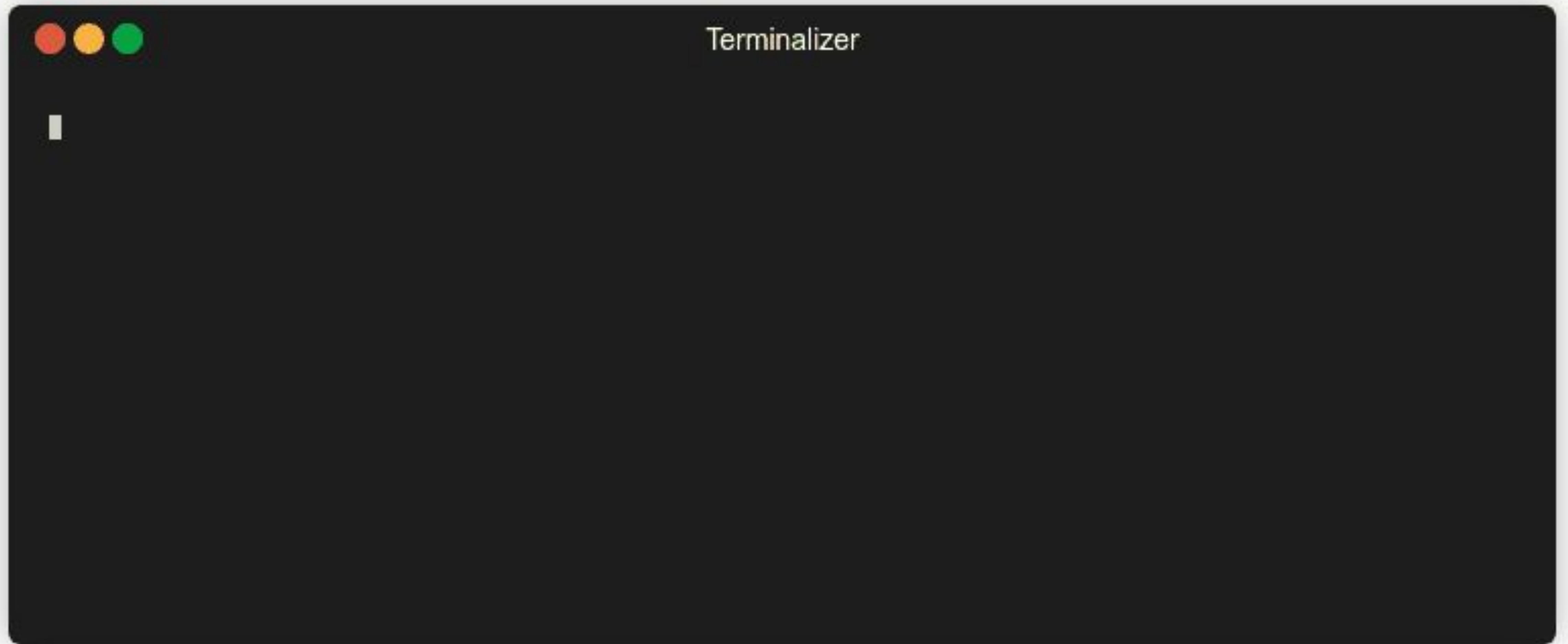- organized

**Code:**

- reproducible
- portable
- (scalable)

# Organise the project folders

Consistent structure, which is easy to navigate, maintain, and release



```
.
├── LICENSE
├── README.md
├── config
├── data
│   ├── external
│   ├── interim
│   ├── processed
│   └── raw
├── results
├── notebooks
├── reports
│   └── figures
└── src
    ├── preprocessing
    ├── postprocessing
    └── visualization
```

Data goes on s3, use data MANIFESTs.

Code, can go on github ( versioned )

Released with paper (Zenodo)

https://drivendata.github.io/cookiecutter-data-science/

# Cookiecutter

# Reproducible analysis

**Conda environments** allow to specify and to keep track of the packages and their versions. Whoever tries to run the same code should have the same output and you should be able to reinstall everything on any machine. They also work with jupyter notebooks.
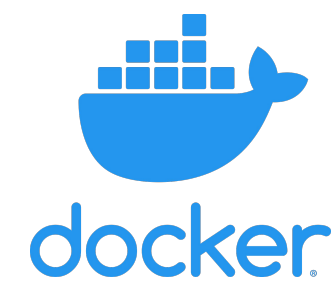
**environment.yml**

```
name: pygna
channels:
 - defaults
 - bioconda
 - conda-forge
dependencies:
 - python>=3.7
 - pandas
- numpy
 - scipy
 - matplotlib
 - pyyaml
 - pytables>=3.4.4
 - seaborn>=0.9
 - networkx=2.3
```

# Reproducible analysis



**Conda environments** allow to specify and to keep track of the packages and their versions. Whoever tries to run the same code should have the same output  and you should be able to reinstall everything on any machine. They also work with jupyter notebooks.



**Docker containers** allow to create 'virtual machines' with all the software you need. Not only you have all the conda packages, but you also set the system where the tools are installed. They make the analysis reproducible top to bottom.

**Running jupyter projects with docker:**

https://towardsdatascience.com/dockerizing-jupyter-projects-39aad547484a

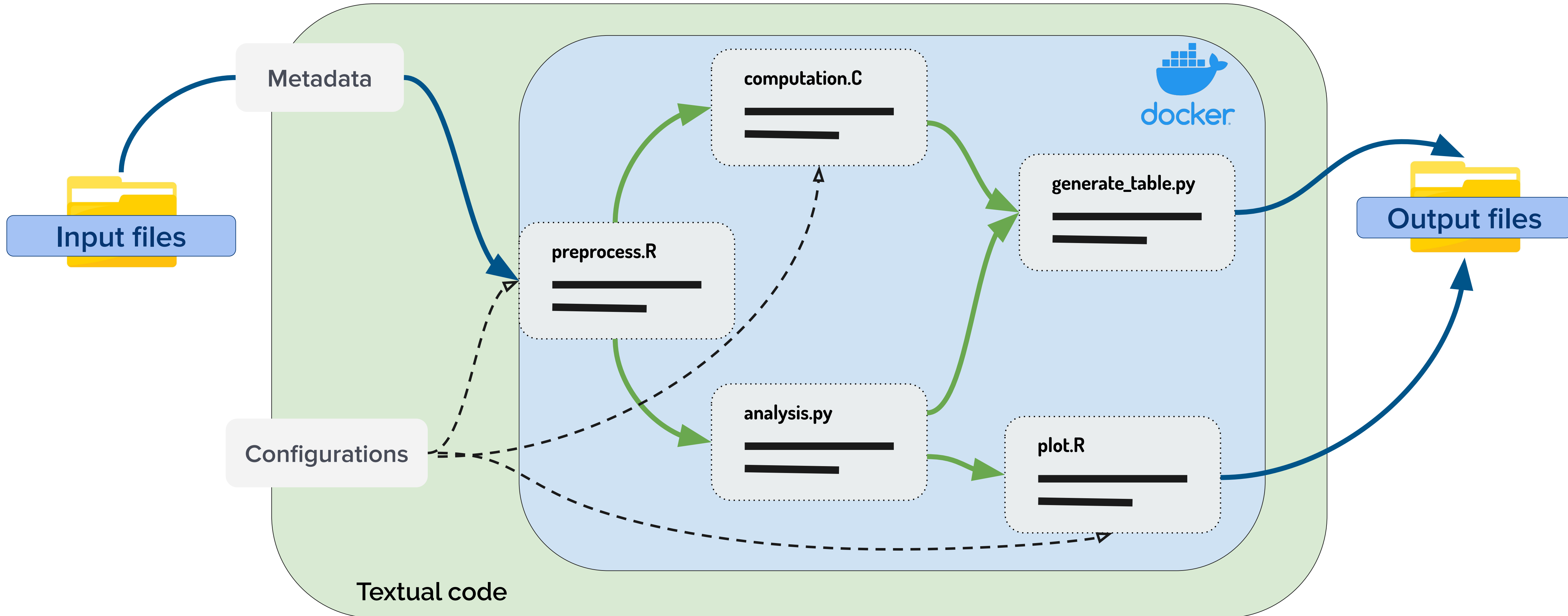# Reproducible workflows

**nextflow**

**Snakemake**

**Workflow management tools** allow to specify a pipeline for the analysis: entire analysis reproducible with 1 command. Portable: it uses docker containers and works by using configuration files, the same analysis can be done on any dataset and any machine.

*cookiecutter

Nextflow AWS ready: launches jobs by itself, reads and saves files in S3, installs the whole image every time ( no need to keep 100 AMIs )
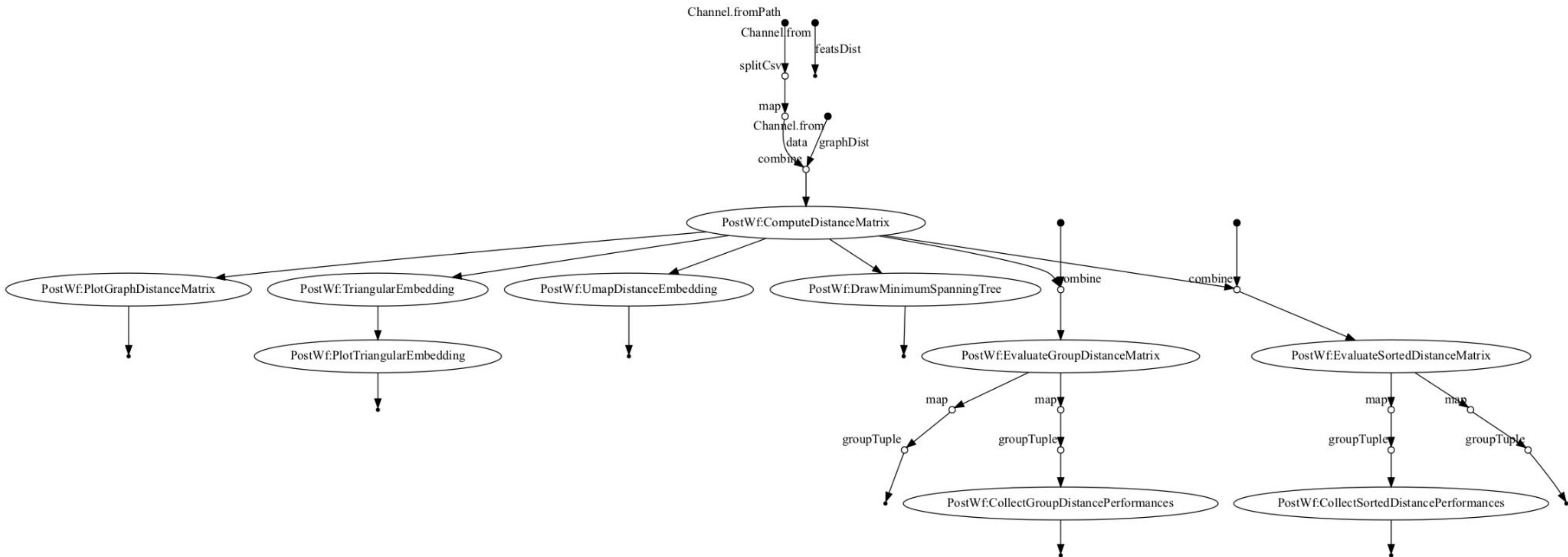
# Workflow Management Tools

# My workflow

> nextflow run **yeast-cell-cycle-nf/main.nf** –param <param–value> –resultsDir batch–yeast–20220406–expression ....

# My workflow

> nextflow run yeast-cell-cycle-nf/main.nf –param <param-value> –resultsDir batch-yeast-20220406-expression ....

```
Session ID : 8dcd6c85-4706-4f4f-accf-5107eca6afae
Workflow: post
Name : batch-yeastcc-20220406-expression
Data table: ./conf/tables/yeast-cc-metadata-expression.csv
Results dir: remove
Graph Distances : [euclidean, euclidean_normal, pearson, pearson_normal, spearman]
Feats Distances : []

./conf/tables/yeast-cc-metadata-expression.csv
executor >   local (126)
[8a/11b04e] process > PostWf:ComputeDistanceMatrix (14)            [100%] 15 of 15 ✔
[61/090829] process > PostWf:PlotGraphDistanceMatrix (15)          [100%] 15 of 15 ✔
[69/ea8b45] process > PostWf:EvaluateSortedDistanceMatrix (15)     [100%] 15 of 15 ✔
[fb/ae7b61] process > PostWf:CollectSortedDistancePerformances (3) [100%] 3 of 3 ✔
[13/308def] process > PostWf:EvaluateGroupDistanceMatrix (15)      [100%] 15 of 15 ✔
[ec/2aaf5a] process > PostWf:CollectGroupDistancePerformances (2)  [100%] 3 of 3 ✔
[3e/1ed42a] process > PostWf:TriangularEmbedding (15)              [100%] 15 of 15 ✔
[48/6ef36e] process > PostWf:PlotTriangularEmbedding (15)          [100%] 15 of 15 ✔
[62/0fd8a0] process > PostWf:UmapDistanceEmbedding (15)            [100%] 15 of 15 ✔
[ae/811ac9] process > PostWf:DrawMinimumSpanningTree (15)          [100%] 15 of 15 ✔
WARN: Task runtime metrics are not reported when using macOS without a container engine
Completed at: 08-Apr-2022 11:35:07
Duration    : 1m 14s
CPU hours   : 0.2
Succeeded   : 126
```

```
.
└── batch-yeastcc-20220406-expression
    ├── expr_cc
    │   ├── distances
    │   ├── figures
    │   └── tables
    ├── expr_cc1
    │   ├── distances
    │   ├── figures
    │   └── tables
    └── expr_cc2
        ├── distances
        ├── figures
        └── tables
```

13

# Nextflow
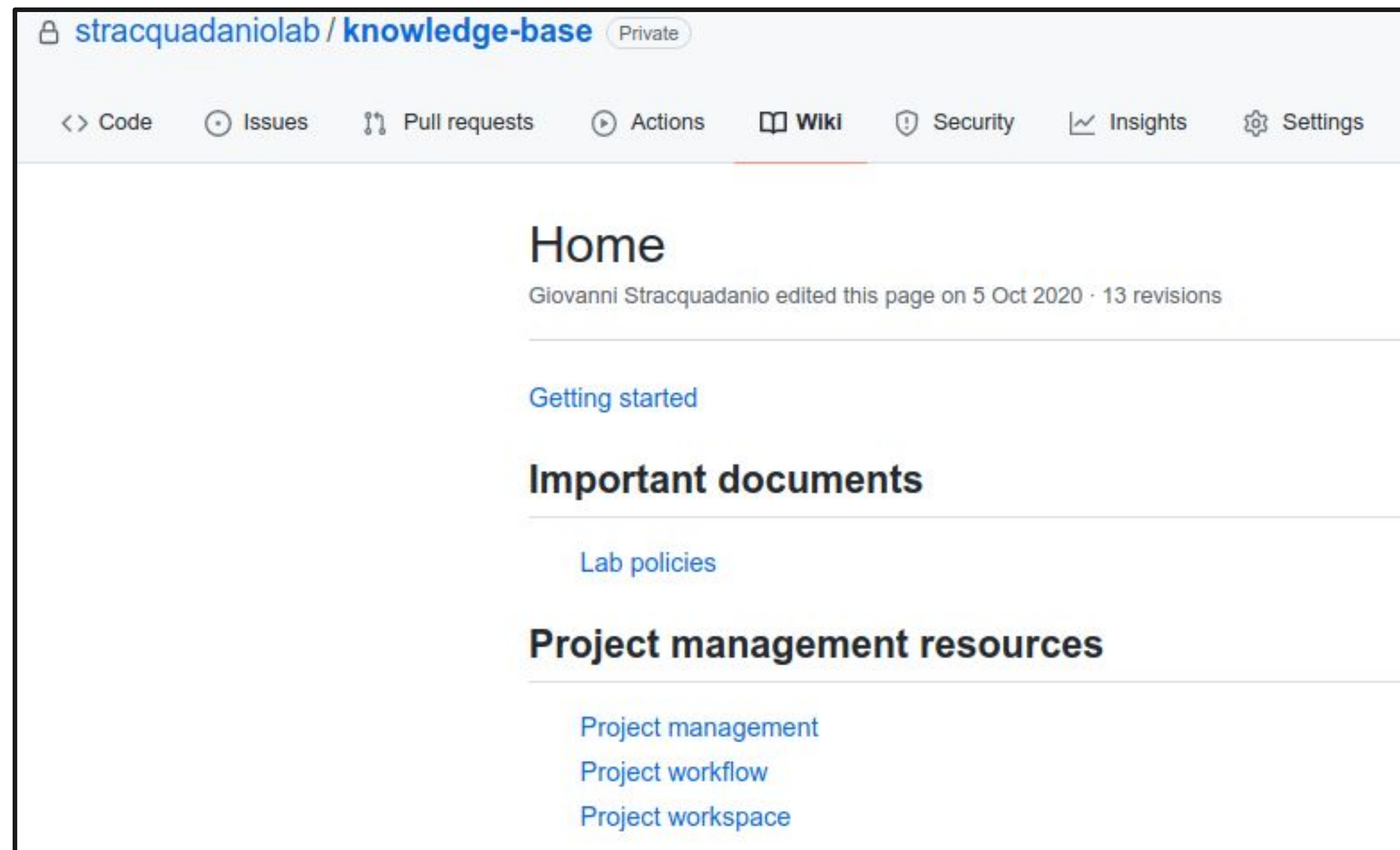


14

# Possible workflow for the lab

- Create data folder: data sync on S3 and code sync on github. Preprocessing, exploratory analyses, final postprocessing, examples all in here.

- Create workflow for the analysis: testable on your own computer and then usable on AWS

- Run analysis on AWS (possibly with AWS batches )

# Lab wiki

**Keep lab knowledge and expertise in a 'centralised' wiki.**

# Reproducible science

## Tools and good practices

Viola Fanfani

January 2022