# FAIR Bioinfo 2022
# Les principes FAIR dans un projet de bioinformatique

# Conclusion

- Introduction to FAIR & Open Science
- Project management (data, repository)
- Traceability with notebooks (jupyterlab)
- History management (Git)
- Introduction to encapsulation
- Environment management (Conda)
- Containerisation (docker)
- Share & disseminate, code & project (GitHub)
- HPC : cluster (Slurm)
- Analysis workflow (snakemake)
- HPC : containerisation (singularity)
- Conclusion

- Introduction to FAIR & Open Science
- Project management (data, repository)
- Traceability with notebooks (jupyterlab)
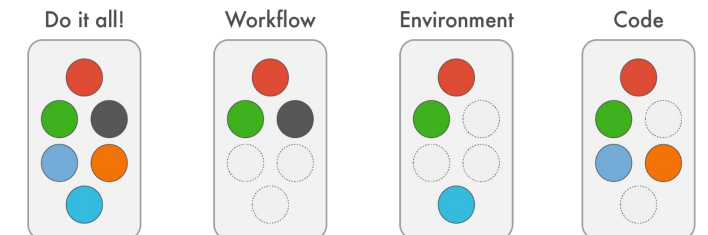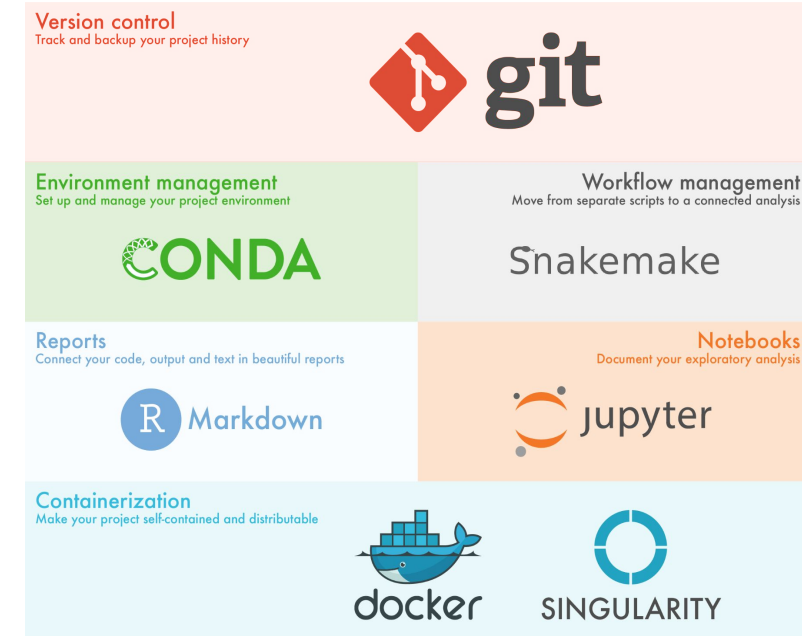- History management (Git)

e-labbook

- Introduction to encapsulation
- Environment management (Conda)
- Containerisation (docker)
- Share & disseminate, code & project (GitHub)

code development

- HPC : cluster (Slurm)
- Analysis workflow (snakemake)
- HPC : containerisation (singularity)

HPC

- Conclusion

- Introduction to FAIR & Open Science
- Project management (data, repository)
- Traceability with notebooks (jupyterlab)
- History management (Git)
- Introduction to encapsulation
- Environment management (Conda)
- Containerisation (docker)
- Share & disseminate, code & project (GitHub)
- HPC : cluster (Slurm)
- Analysis workflow (snakemake)
- HPC : containerisation (singularity)
- Conclusion

- Project management (data, repository)
- Traceability with notebooks (jupyterlab)
- History management (Git)
- Environment management (Conda)
- Containerisation (docker)
- Share & disseminate, code & project (Git***)
- HPC : cluster (Slurm)
- Analysis workflow (snakemake)
- HPC : containerisation (singularity)

| Céline (NGS) | Claire (NGS analysis) | Emilie (NGS) | Gildas (admin) |
| Julien (Dev) | Pauline (WGA) | Hugo (Dev & NGS) | Thomas (Dev) |

*** Lab !

- With a tool:
  - download, install, config, run
  - create new (env_config, snakefile, dockerfile)
  - participate to the development of the tool (bioconda)
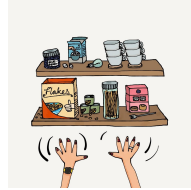
**F**indable  **A**ccessible  **I**nteroperable  **R**eusable

A virtuous cycle

**Objective:**

FAIR raw data
+
FAIR scripts
=
FAIR processed data

**Code** avoid workflows based on **point-and-click interfaces** (eg. Excel), enshrine computations and data manipulation in code

**Document** how code works, define parameters and computational environment required: comments, **notebooks** and **README**

**Record** key parameters (eg. the 'seed' values of a random-number generator)

**Test** functions using positive and negative control **data sets**, run those tests throughout development

**Guide** with master script (eg. 'run.sh') that downloads data sets and executes workflow

**Archive** with long-term stability services such as Zenodo, Figshare and Software Heritage (GitHub is impermanent online repository). Track the project's history with a **version-control** tools (eg. Git).

**Note** (tag) which version you used to create each result

**Package** with ready-to-use computational environments using **containerization** tools (eg. Docker, Singularity), web services (Code Ocean, Gigantum, Binder) or **virtual-environment** managers (Conda)

**Simplify** and avoid niche or hard-to-install third-party code libraries

**Verify** your code's portability by running it in a range of computing environments

**Automate** the test of your code with **continuous-integration** services(eg. Travis CI)

https://www.nature.com/articles/d41586-020-02462-7

## Unit test: test a part of the code

```
## module 1
sum <- function(x, y){
    return (x+y)
}


# Unit test
sum(2,2) == 4
```

```
1 ## module 2
2 power <- function(x, y){
3     return (x**y)
4 }
5
6 # Unit test
7 power(2,2) == 4
```

## Functional test: test all the code

```
# Functional test
power(sum(2,2),2) == 16
```

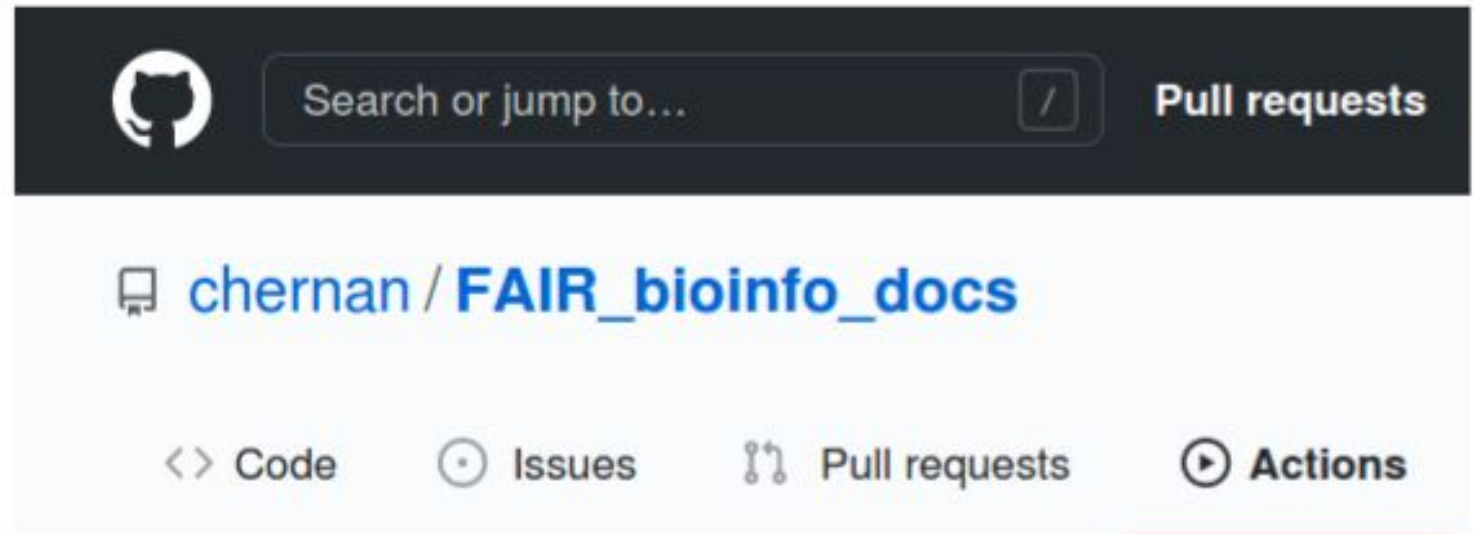Automated verification each time the source code is modified that the modifications do not produce:

- any regression in the developed application
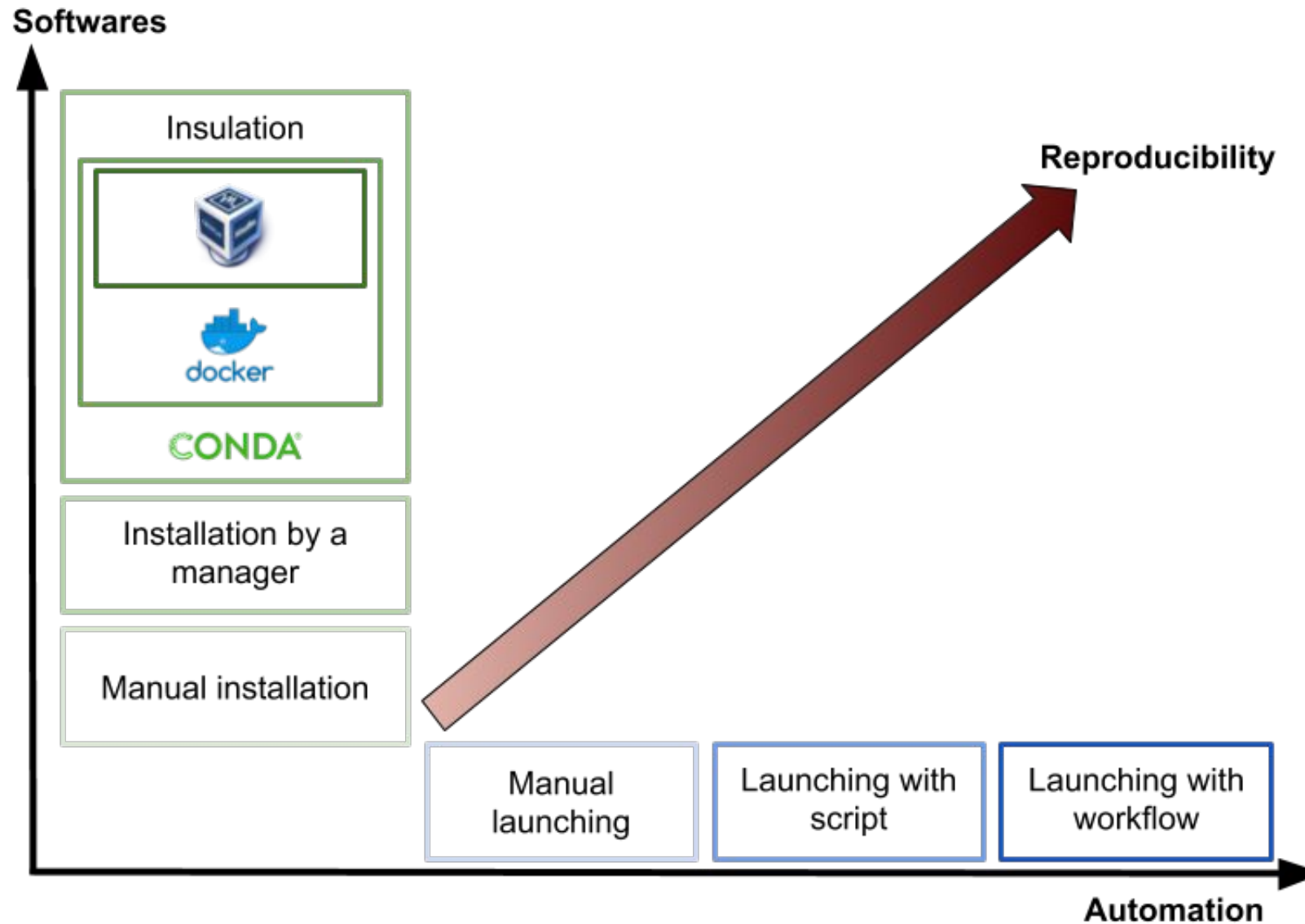- any change in the results obtained

## Automation

Manual command lines

↓

Write a shell script

↓

Use a workflow manager

↓

Tests and continuous integration (*)

## User analysis (trial-and-error)

Offer a GUI (eg. with R-Shiny) (*)

↓

Save and re-import choices (*)

## Softwares

Local installation
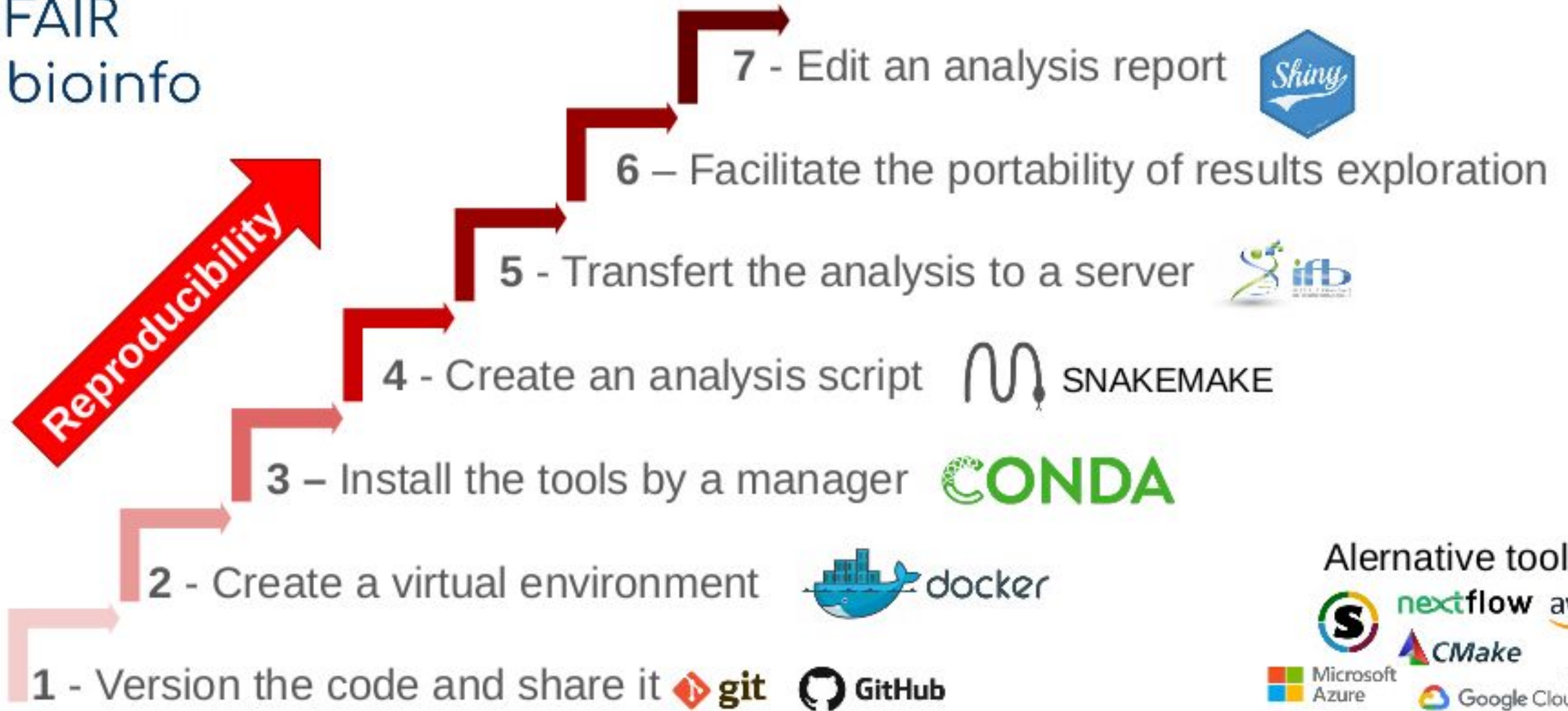
↓

Package manager

↓

Conda environment

↓

Image / container

↓

Virtual machine (*)

(*) not in the course

Reproducibility to the exact bit?
container uses some resources of the
support machine
⇒ version control of the env.
(Nix, Guix)

HPC and parallelization?
loss of computational order,
multithreading, identical
hardware?
⇒ …

Pedagogical team (our guardian angels): Yousra, Hélène

IFB Core Cluster taskforce: Julien, Gildas, and all those who provide in the shadows

Helpers: Emilie, Pauline, Hugo

Organisations: CNRS, INRAE, IFB, I2BC, Paris Saclay University