# FAIR_bioinfo : Open Science and FAIR principles in a bioinformatics project

## How to make a bioinformatics project more reproducible

C. Hernandez[1]    T. Denecker[2]    J. Sellier[2]    G. Le Corguillé[2]
C. Toffano-Nioche[1]

[1]Institute for Integrative Biology of the Cell (I2BC)
UMR 9198, Université Paris-Sud, CNRS, CEA
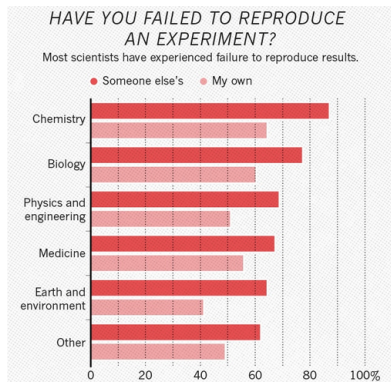91190 - Gif-sur-Yvette, France

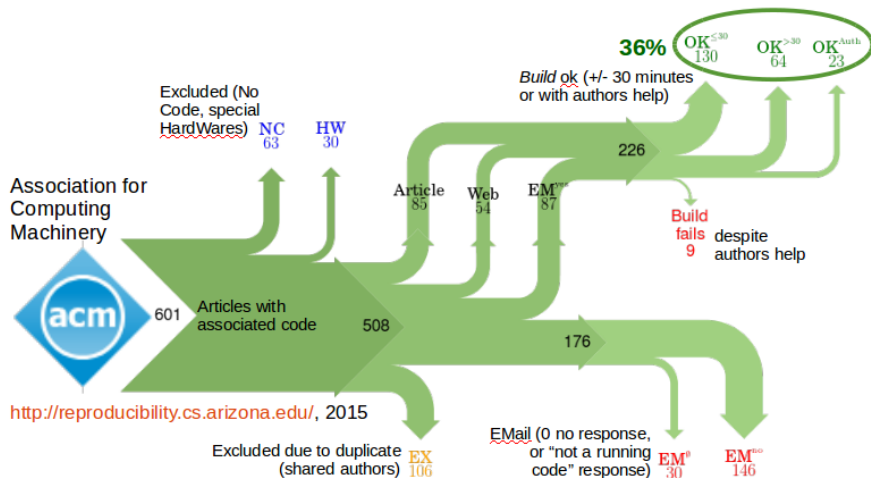[2]IFB Core Cluster taskforce

June 2021

# Reproducibility

# A reproducibility problem, Biology

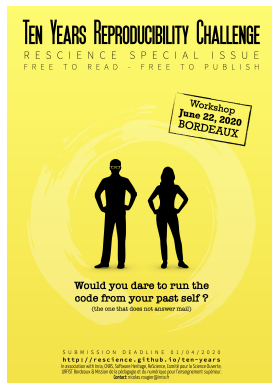70% of the analyses in Experimental Biology are not reproducible



Monya Baker, 1,500 scientists lift the lid on reproducibility, *Nature*, 2016

# A reproducibility problem, Computer Sciences
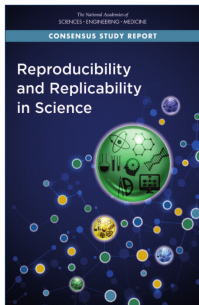
# A reproducibility problem, Bioinformatics



Ten-Year Reproducibility
Challenge, Konrad Hinsen
Can your 2009 code still run?
special issue of ReScience and
result comments in *Nature*

Who's never wanted to take over a protocol, a pipeline, or a tool without running into it?

- unable to install tools: not compatible OS, not availability of dependencies
- tool update $\Rightarrow$ codes unusable: python 2 vs. 3, change of function arguments (R)
- inability to reproduce the results of computational analysis: package versions, IDE: stable version of the language different according to the OS (Rstudio)

# Reproducibility in science

*Reproducible research, Repeatability, Replicability, Reproducibility, Replication:* overlapping semantics ⇒ a plethora of definitions![a]



National Academies of Sciences, Engineering, and Medicine (2019).[b]

ACM definition (2016):

Repeatability   Same team, same exp. setup

Replicability   Different team, same exp. setup

Reproducibility   Different team, different exp. setup

Whitaker's matrix of reproducibility (2017):[c]

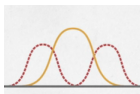| | | Data | |
|---|---|---|---|
| | | Same | Different |
| **Analysis** | Same | Reproducible | Replicable |
| | Different | Robust | Generalisable |

# FAIR_bionfo's finding

Depends on the object of study x
        what needs to be "memorized" to replay the experience:

 $\rightarrow$  $\rightarrow$ 

| Raw Data | Statistical or bioinformatic analysis | Validation |
|:---:|:---:|:---:|
| FAIR data principles & Data Management Plans | Codes - parameters - workflows | Publication: thesis, article, report, etc |

**How to gain in reproductibility?**

Focus on codes, parameters, and workflows used throughout the analysis process

Monya Baker, 1,500 scientists lift the lid on reproducibility, *Nature*, 2016

# A solution

# Divert FAIR data principles towards processes

| **F**indable | **A**ccessible | **I**nteroperable | **R**eusable |
|:---:|:---:|:---:|:---:|
| Third party tools used = ref. in their field | Available codes (Github, dockerhub) | Cooperation of tools (snakemake, docker) as well as locally than on servers (cloud or cluster) | Protocol replayable (snakemake) identically (Rshiny) in a virtual environment (docker) |
| . | . | | |
| Easy to find analysis protocol (Github pages) | Third party open source tools | | |

# Promote learning


FAIR bioinfo

## Our objective

FAIR raw data
+
FAIR scripts
=
FAIR processed data

## Course

Take your first steps with several companion tools to gain in reproducibility

## Example based

Just the beginning of an NGS analysis A full analysis is given as bonus (The NGS analysis is simply used as an example and not explained)

# Ressources

- [awesome](#) a curated list of reproducible research case studies, projects, tutorials, and media
- The Role of [Metadata](#) in Reproducible Computational Research
- [Towards reproducible computational biology](#)
- A very similar sweden [courses](#) with git, conda, snakemake, jupyter, r-markdown, docker, singularity