

FAIR_bioinfo for bioinformaticians

Introduction to the tools of reproducibility in bioinformatics

C. Hernandez¹ T. Denecker¹ J.Sellier² C. Toffano-Nioche¹

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

²Institut Français de Bioinformatique
à compléter

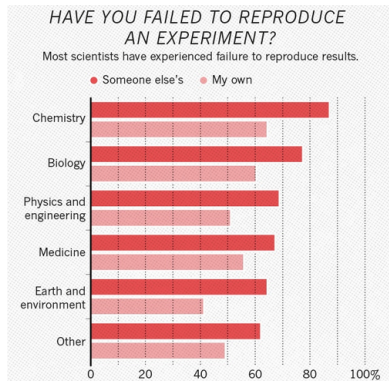
Sept. 2020



Introduction to reproducibility

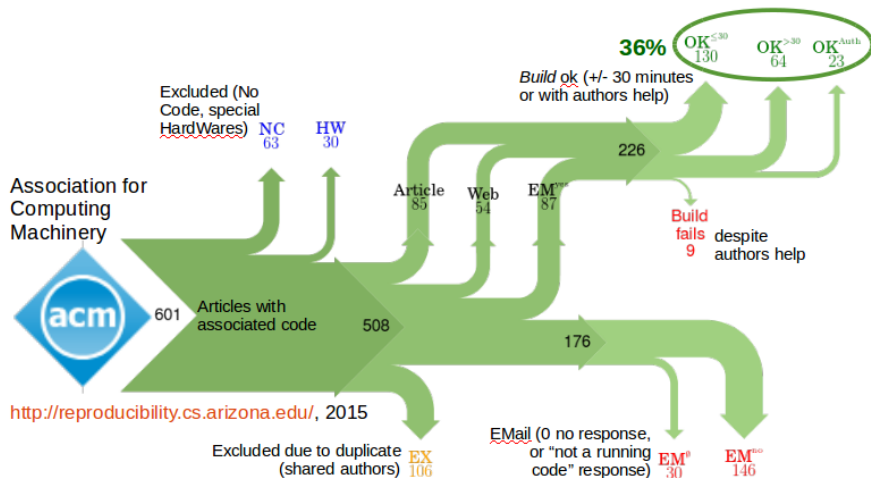
A reproducibility problem, Biology

70% of the analyses in Experimental Biology are **not** reproducible

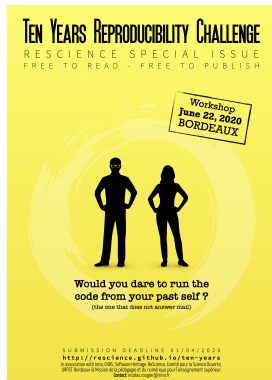


Monya Baker, 1,500 scientists lift the lid on reproducibility, *Nature*, 2016

A reproducibility problem, Computer Sciences



A reproducibility problem, Bioinformatics



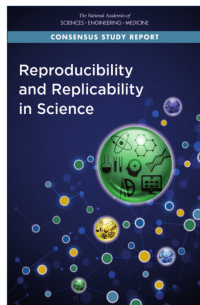
Ten-Year Reproducibility Challenge, Konrad Hinsén
Can your 2009 code still run?
special issue of [ReScience](https://www.re-science.org/)

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 \Rightarrow a plethora of definitions!^a



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

a: https://www.researchgate.net/publication/323118701_Terminologies_for_Reproducible_Research

b: National Academies of Sciences, Engineering, and Medicine. 2019. Washington DC. The National Academies Press, <https://www.nap.edu/read/25303/chapter/1>

c: <https://doi.org/10.6084/m9.figshare.5443201.v1>, Slide number 7

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

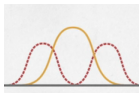
FAIR_bionfo's finding

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



Raw Data

FAIR data principles
& Data Management
Plans



Statistical or
bioinformatic analysis
Codes - algorithms -
workflows



Validation

Publication: thesis,
article, report, etc

How to gain in reproductibility?

Focus on codes, algorithms, workflows used throughout the process

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

A solution

Divert FAIR data principles towards processes

Findable



Third party tools
used = ref. in
their field

Easy to find
analysis protocol
(Github pages)

Accessible



Available codes
(Github,
dockerhub)

Third party open
source tools

Interoperable



Cooperation of
tools (snakemake,
docker) as well as
locally than on
servers (cloud or
cluster)

Reusable



Protocol
replayable
(snakemake)
identically
(Rshiny) in a
virtual
environment
(docker)

Promote learning



Our objective

FAIR raw data

+

FAIR scripts

=

FAIR processed data

Course

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

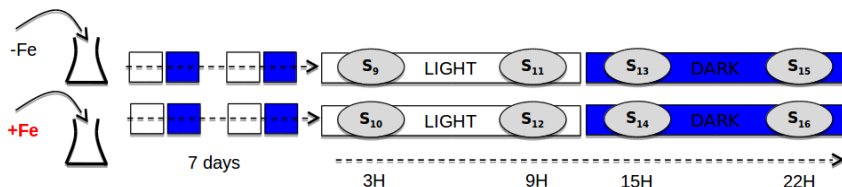
Example based

Classical RNA-seq analysis
(finding genes with differential expression
between 2 conditions)
used as an example (not explained)

Working example

The biological study example

- Study of the green alga *Ostreococcus tauri* response to iron deprivation.
- 16 RNAseq samples in triplicate, single-end of 100bp.
- Choice of the 9h point of the long-term adaptative response (s11 and s12 samples):



Lelandais G, Scheiber I, Paz-Yepes J, Lozano JC, Botebol H, Pilátová J, Žárský V, Léger T, Blaiseau PL, Bowler C, Bouget FY, Camadro JM, Sutak R, Lesuisse E.

Ostreococcus tauri is a new model green alga for studying iron metabolism in eukaryotic phytoplankton.

BMC Genomics. 2016 May 3;17:319. doi: 10.1186/s12864-016-2666-6.

Reduced RNAseq Data

Genome

- sequence: GCF_000214015.3_version_140606_genomic.fna
(https://www.ncbi.nlm.nih.gov/assembly/GCF_000214015.3/)
- annotation: GCF_000214015.3_version_140606_genomic.gff
- \Rightarrow 13.0328 Mb, 20 chromosomes, mitochondria, & chloroplast

RNAseq samples

- Project: PRJNA304086
- Selection samples 11 and 12: SRR3099585-87, SRR3105697-99
(fastq.gz \sim 360M each x 6 files)
- Reads selection to reduce data volume for the course (mapped on the smallest chromosome, chr18, NC_014443.2 + 100000 first) \Rightarrow
*_chr18.fastq.gz \sim 19M each (<https://zenodo.org/record/3997237>)
- Counts table, complete RNAseq: <https://zenodo.org/record/3997137>

Data

Access in the IFB ressources

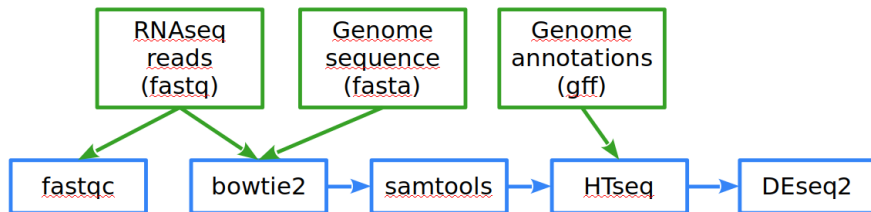
```
1 /shared/projects/fair_training2020/Data/
```

Or raw download in a local "Data" directory

```
1 mkdir Data ; cd Data
2 wget https://ftp.ncbi.nlm.nih.gov/genomes/all/GCF
   /000/214/015/GCF_000214015.3_version_140606/
   GCF_000214015.3_version_140606_genomic.fna.gz
3 wget https://ftp.ncbi.nlm.nih.gov/genomes/all/GCF
   /000/214/015/GCF_000214015.3_version_140606/
   GCF_000214015.3_version_140606_genomic.gff.gz
4 wget https://zenodo.org/record/3997237/files/
   FAIR_Bioinfo_data.tar.gz
5 wget https://zenodo.org/record/3997137/files/counts.txt
6 cd ..
```

RNAseq analysis

Analysis workflow



green=input, blue=tool

fastqc control quality of the input reads

bowtie2 reads mapping on the genome sequence

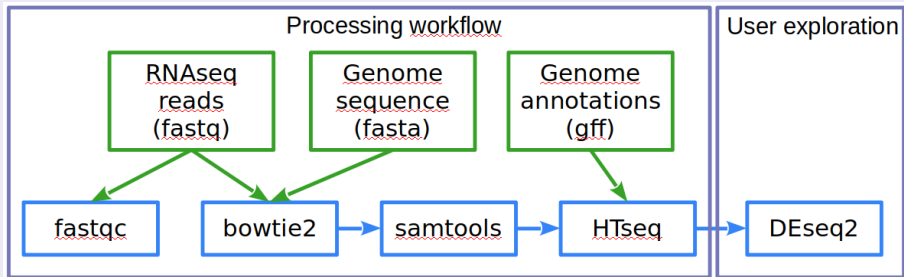
samtools mapped reads selection & formatting

HTseq count table of mapped reads on genes (annotations)

DEseq2 statistical analysis: genes list having differential expression

2 bioinformatician skills

Analysis workflow



green=input, blue=tool, violet=skill

Reproducibility

Processing workflow automatization, scripting

User exploration report choices (or import choices for further analysis)

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