#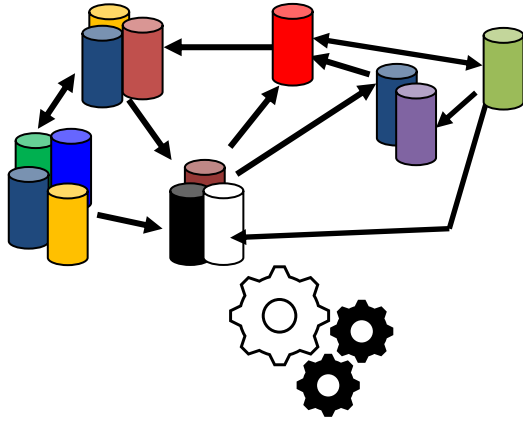 10 ans après la crise de la reproductibilité en bioinformatique : de la reproduction à la réutilisationde workflows scientifiques

Sarah Cohen-Boulakia November, 2025
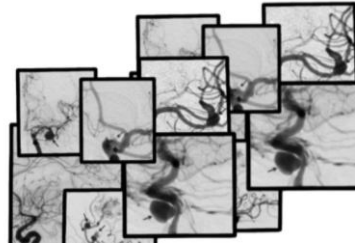
Laboratoire Interdisciplinaire des Sciences du Numérique – LISN
Université Paris-Saclay

# Biological and Biomedical data analysis



**Public sources**
Distributed
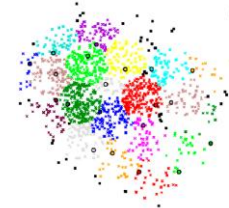Heterogeneous Network
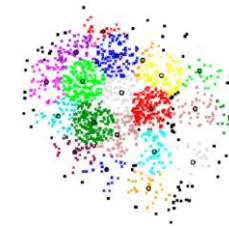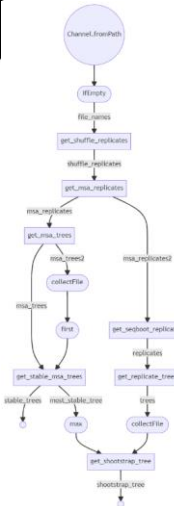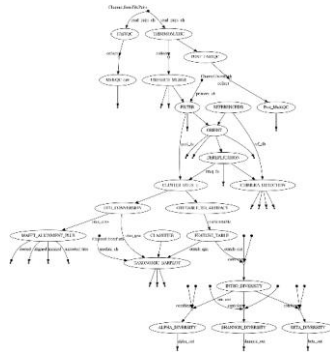> 1,500 (NAR)

**Tools – Scripts**
Distributed
Heterogeneous
> 23,000 (bio.tools)

*Which exact dataset did I use? Which tool version? Which parameters...?*

**Analysis pipelines**
Combining multiple tools
Heterogeneous
Various environnements
& platforms
> 1,000 workflows
(GitHub)

## Nekrutenko & Taylor - Nature Genetics (2012)

50 papers 2011 using the Burrows-Wheeler Aligner

31/50 (62%) provide no information

no version of the tool + no parameters + no genomic ref sequence

7/50 (14%) provide all the necessary details

## Alsheikh-Ali et al, PLoS one (2011)

10 papers in the top-50 IF journals → 500 papers

149 (30%) were not subject to any data availability policy

(0% data available)

Of the remaining 351 papers

208 papers (59%) did not adhere to the data availability instructions
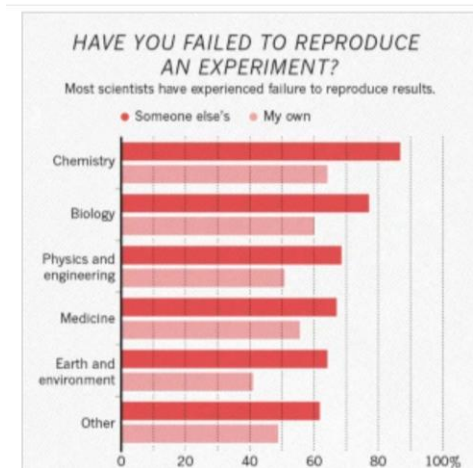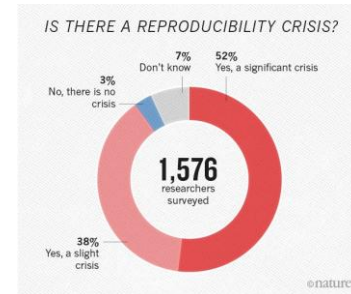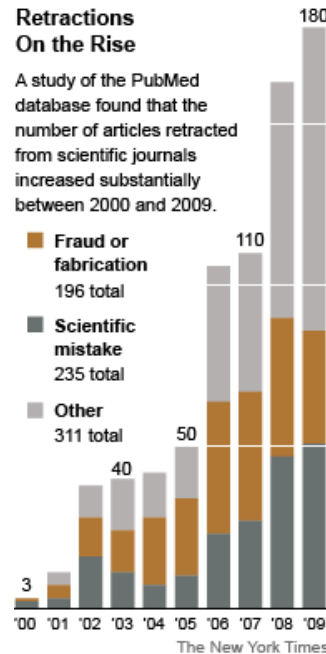
143 make a statement of *willingness* to share

47 papers (9%) deposited full primary raw data online

→ Computational reproducibility

## Nekrutenko & Taylor - Nature Genetics (2012)

50 papers 2011 using the Burrows-Wheeler Aligner

31/50 (62%) provide no information

  no version of the tool + no parameters + no genomic ref sequence

7/50 (14%) provide all the necessary details

## Alsheikh-Ali et al, PLoS one (2011)

10 papers in the top-50 IF journals → 500 papers

149 (30%) were not subject to any data availability policy

(0% data available)

Of the remaining 351 papers

208 papers (59%) did not adhere to the data availability instructions

143 make a statement of *willingness* to share

47 papers (9%) deposited full primary raw data online

### → Computational reproducibility

Scientific Workflow Management Systems
"Data analysis pipeline "
Data flow driven

**WF specification**
connected *processors*
**steps** of the analysis
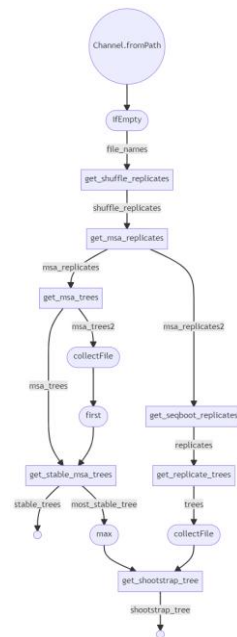
**WF execution**
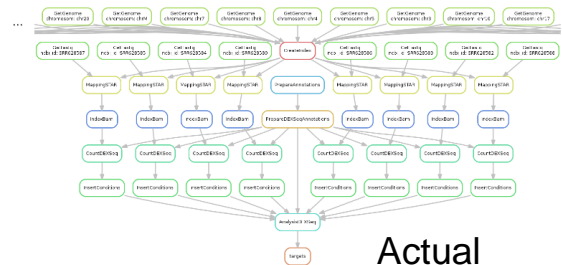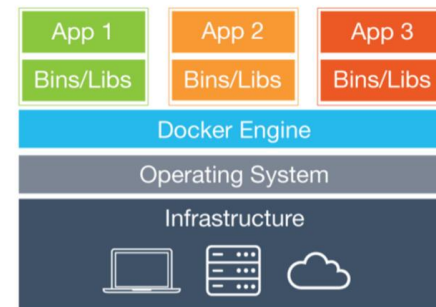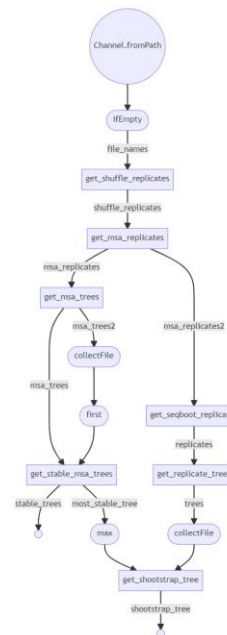**data consumed/produced**
Provenance modules
**Scheduling** …



**Specification Steps**



Possible dataflow

**Executed Steps**



Actual dataflow

Scientific Workflow Management Systems
"Data analysis pipeline "
Data flow driven

**WF specification**
connected *processors*
steps of the analysis

**WF execution**
data consumed/produced
Provenance modules
Scheduling …

**WF environment**
everything needed to run Libraries, dependencies…
Coupling WF with Docker/AppTainer
BioContainer

# Levels of computational reproducibility

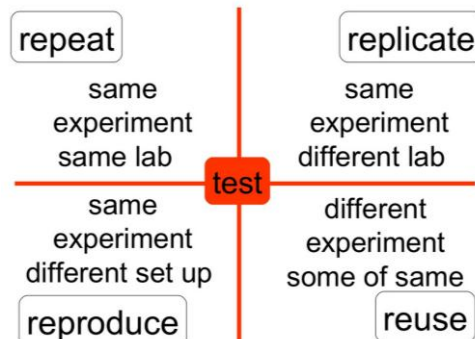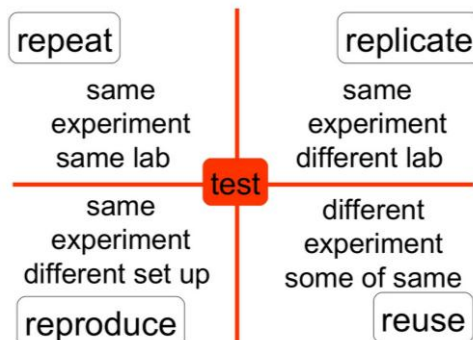*From identical output to the same scientific results*

**Repeat**

*Redo* - exact same context

**Same workflow, execution setting, environment**

Identical *output*

Aim = proof for reviewers ☺

**Replicate**

Variation allowed in the workflows execution setting environment

Similar *output*

Aim = robustness

*From identical output to the same scientific results*

**Repeat**

*Redo* - exact same context

**Same workflow, execution setting, environment**

Identical *output*

Aim = proof for reviewers ☺



repeat — same experiment same lab

replicate — same experiment different lab

reproduce — same experiment different set up

reuse — different experiment some of same

Drummond C Replicability is not Reproducibility: Nor is it Good Science, online
Peng RD, Reproducible Research in Computational Science *Science 2 Dec 2011: 1226-1227.*

**Replicate**

Variation allowed in the workflows execution setting environment

Similar *output*

Aim = robustness

**Reproduce**

Same *scientific result*

But the means used may be changed

Different workflows, execution setting, environment

Different output but in accordance with the result

**Reuse**

Adapt to new needs

Possibly different scientific result

Reuse in part existing workflows

→ **Cumulative science**

→ **No Reuse without Repeat!**

French Reproducibility
Network born in 2023
270+ colleagues

Next Scientific Days in Lyon
April 3-4th

https://www.recherche-reproductible.fr/index-en

**MINISTÈRE DE L'ENSEIGNEMENT SUPÉRIEUR ET DE LA RECHERCHE**
*Liberté*
*Égalité*
*Fraternité*

**Global Networks**
Outside the UK? Find a Reproducibility Network in your area

See full Global Networks Statement

**Global Reproducibility Networks**

A Reproducibility Network (RN) is a national, peer-led consortium of researchers that aims to promote and ensure rigorous research practices by establishing appropriate training activities, designing and evaluating research improvement efforts, disseminating best practice and working with stakeholders to coordinate efforts across the sector. RNs aim for broad disciplinary representation and an intensive interdisciplinary dialogue (e.g., with funding agencies, publishers, learned societies and other sectoral organisations, as well as researchers from all disciplines and across all career stages).

To reach as many researchers as possible, and to operate as efficiently as possible, we are keen to support other countries interested in creating similar networks. If you are interested in setting up a national RN, or finding out who in your country is working towards this, please email: contact@ukrn.org

Africa — africamn.org
Australia — aus-rn.org
Belgium — reproducibilitynetwork.be
Brazil — reprodutibilidade.org
Canada — carn-recar.ca
Croatia — crorin.hr
Denmark — danish-repro.github.io
Finland — finnish-rn.org
France — recherche-reproductible.fr
Germany — reproducibilitynetwork.de
Italy — itrn.org
Luxembourg — uni.lu
Netherlands — reproducibilitynetwork.nl
Norway — norrn.no
Portugal — ptrn.pt
Slovakia — slovakrn.wixsite.com/skrn
Spain — sprn.es
Sweden — swern.org
Switzerland — swissrn.org
United Kingdom — ukrn.org

The reproducibility landscape

Status of workflow reuse

How to improve workflow reuse?

ZOOM on 2 current contributions

Conclusion

Distinct processors

1,700 Taverna workflows
10 242 processors (analysis steps)
Centralized in myExperiment

Re-use rates have a Zipf-like distribution

The top ten authors published 62% of all workflows

2,443 workflows Nextflow & Snakemake
15 540 processors (analysis steps)
Distributed in github
nf-core: repro of carefully checked Nextflow processors

Marine **Djaffardjy**

Still low reuse

The top ten authors published only 15% of all workflows

Higher reuse (black box) of processors in nf-core

The reproducibility landscape

Status of workflow reuse

How to improve workflow reuse?

ZOOM on 2 current contributions

Conclusion

# ShareFAIR

**ShareFAIR** : Sharing reliable workflows to transform datasets into gold standards: Application to Neuro-Vascular Pathologies

**PEPR Santé Numérique**

**Coordinator**
- Sarah Cohen-Boulakia

**Managment structure**
- Université Paris-Saclay

**Partners**
- Université Paris Dauphine PSL
- Institut Pasteur
- Université Lyon
- Université Rennes
- Inria
- CEA
- CNRS, INSERM (ITX)

**Duration**: 48 mois

*Define **an interoperable framework** for the design, annotation, and sharing of **reproducible and reuseable workflows***

Design a **language** to query workflows, and their executions

**Capture provenance** (executions) in an optimized manner

Develop a **FAIRification method** for datasets

Guide workflow developers in **discovering and comparing** workflows

Execution traces: provenance of produced data

*Complete standards* **for uniformly annotating workflows in terms of analysis tools and input/output datasets to produce FAIR datasets**

Identification of a set of existing standards for annotating protocols, workflows, and datasets: EDAM, EDAM-BioImaging, MONDO, DUO, etc.

Identification of overlaps, varying levels of precision, and missing concepts

Construction of annotation guides

Development of a Knowledge base



Annotated data set 1

Annotated tool 1

Annotated tool 2

16

***Augment the set of workflows*** **by extracting them from text (litterature) and large datasets from communauties**

**Develop NLP models** to extract workflow description from the scientific literature

**Learn protocols form clinical data** collected all along the care activity (patient data). Comparision between learnt and declared protocols.

**Extract protocols and workflows** from large shared datasets in neuroimagery

nature

Explore content   About the journal   Publish with us   Subscribe

nature > articles > article

Article | Published: 20 May 2020

**Variability in the analysis of a single neuroimaging dataset by many teams**

17

# Outline

The reproducibility landscape

Status of workflow reuse

How to improve workflow reuse

ZOOM on 2 current contributions

Conclusion

# Exploiting the 3 forms of the workflow



**Code**: GitHub - data science

# Exploiting the 3 forms of the workflow



**Code**: GitHub - data science

**Text**: scientific papers

# Exploiting the 3 forms of the workflow



**Code**: GitHub - data science          **Graph**: workflow structure          **Text**: scientific papers

# Exploiting the 3 forms of the workflow



**Code**: GitHub - data science        **Graph**: workflow structure        **Text**: scientific papers
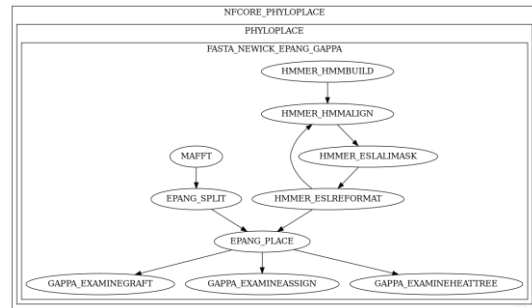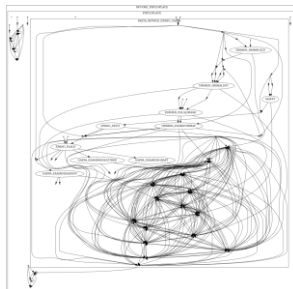
# BioFlow-Insight

**Aim**

Provide an overview of the main step of a workflow

Coherent with the code...



*Workflow Code*

George
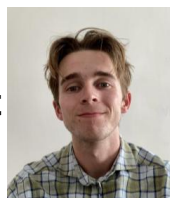**Marchment**

# BioFlow-Insight

## Aim

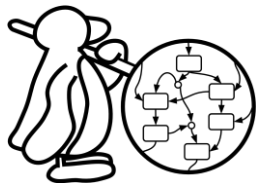Provide an overview of the main step of a workflow

Coherent with the code...

*Workflow Code*
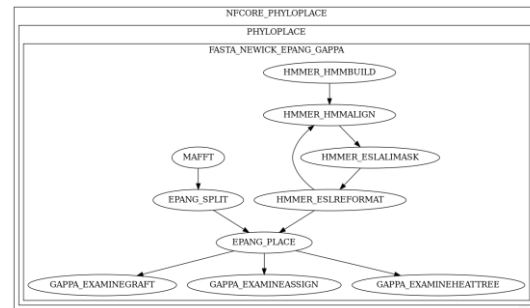
George
**Marchment**

## Functionalities

**Analyses** the code of Nextflow workflows

**Generates** visual graphs at different **simplification** levels depicting the workflow's structure

**Detects errors** in the workflow code

Open source: command line or web service

https://bioflow-insight.pasteur.cloud/
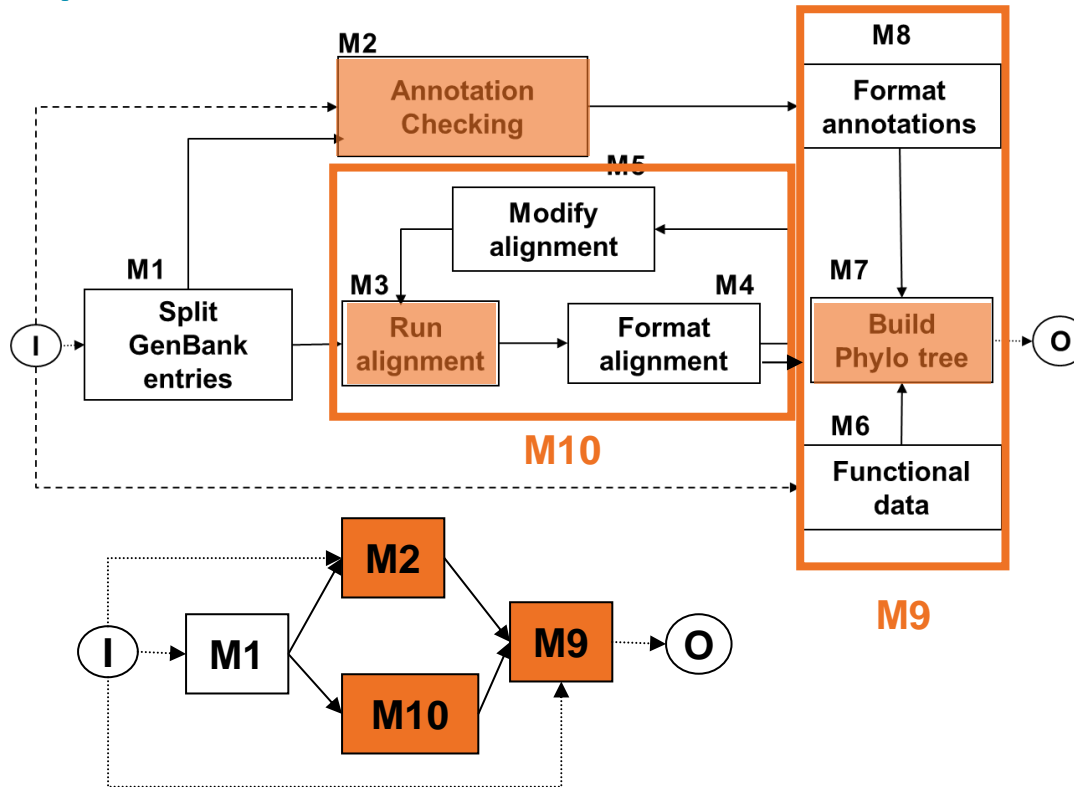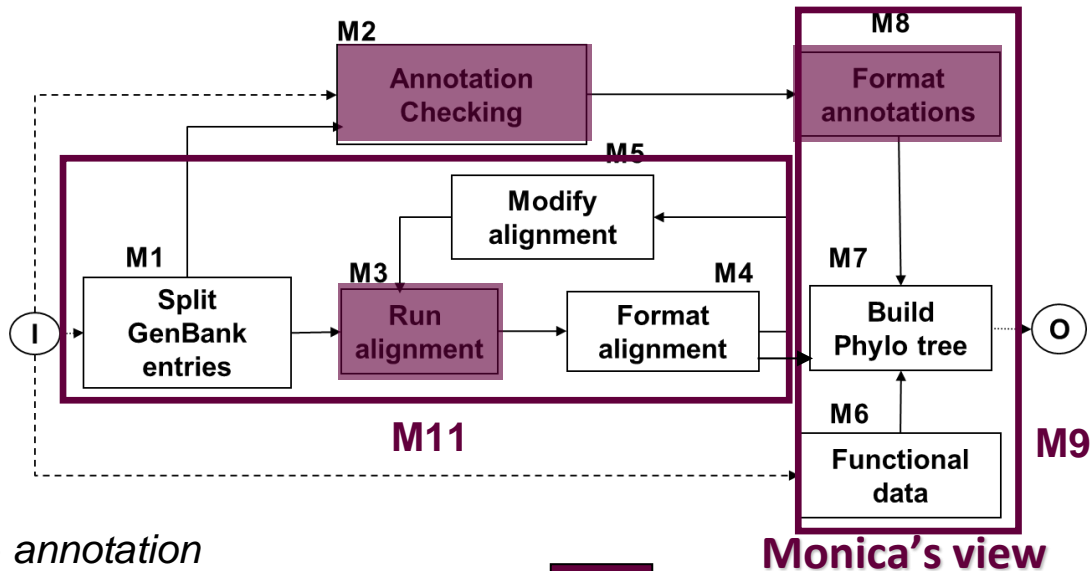
Choose **relevant modules**

Each **composite** module is constructed around one **relevant module**
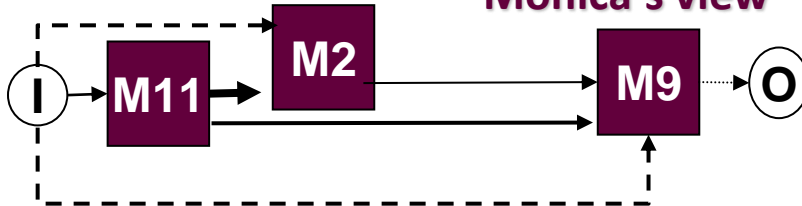


The composite module takes the meaning of the relevant module it contains

Grouping should preserve the relationships between relevant modules



The *annotation checking* module does not need input from the *run alignment* module!

Formalization of the set of properties to be preserved

Property 1: Given Gw and R ⊆ N relevant modules, U is well-formed iff every composite module in U contains at most one element of R.

Property 2: A user view U preserves dataflow iff every edge in Gw that induces an edge on an nr-path from C(r) to C(r') in U(Gw) lies on an nr-path r to r' in Gw.

Property 3: A user view U is complete w.r.t dataflow iff for every edge e on an nr-path from r to r' in Gw that induces an edge e' in U(Gw), e' lies on an nr-path from C(r) to C(r').

**_Theorem_** _ZOOM_ is a polynomial-time which _preserves Properties 1- 3 and produces a minimal user view_

Implementation of ZOOM

**New!** Rewritting the workflow code to implement user views

# Exploiting the 3 forms of the workflow



**Code**: GitHub - data science          **Graph**: workflow structure          **Text**: scientific papers

# Entity extraction from scientific papers

Low-resource extraction task
  Bioinformatics workflows underrepresented in NLP
  No Corpus available
**Several strategies tested**

**Schema of entities** representing
key entities related to workflow

**Annotated corpus of 52 full papers** describing
Snakemake and Nextflow workflows
  7 annotators - Inter annot. Agreement (0.70)



IDA 2025

Neuronal approach biLSTM-CRF
  Nlstruct python lib
  SciBert & BioBert

**F1-score ~0.70 all entities (0.77 on Tools)**
Accepted at IDA 2025

Clémence **Sebe**

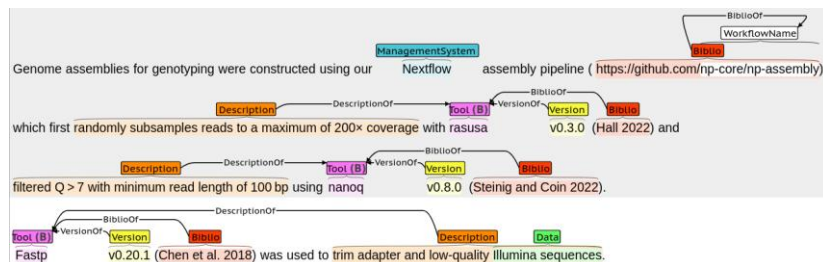Lack of reproducibility hurts **cumulative science**
We are still living the **reproducibility crisis**
Several technical solutions exist to help redo/reexecute

Challenges lie at **the reuse level**: repurposing workflow analyses, adaptating to own needs

ShareFAIR ambitions to provide a **proof-of-concept of workflow sharing** by providing a **reuse platform**

Current work on
      Coupling **workflow code & workflow papers** (NLP-code)
      **Abstracting** workflow **graph structures** (code-graph)

April 3-4 2025
Lyon

# Thanks!