

# FAIR Bioinfo 2022

Best practice in your bioinformatic projects



P. Marin, M. Hiriart, P. Ruiz & N. Goué  
aubi@uca.fr

Université Clermont Auvergne, AuBi, Mésocentre

25 novembre 2022



. This work is based on the IFB and I2BC formation offer

## Essay

# Why Most Published Research Findings Are False

John P.A. Ioannidis

## Summary

There is increasing concern that most current published research findings are false. The probability that a research claim is true may depend on study power and bias, the number of other studies on the same question, and, importantly, the ratio of true to no relationships among the relationships probed in each scientific field. In this framework, a research finding is less likely to be true when the studies conducted in a field are smaller; when effect sizes are smaller; when there is a greater number and lesser preselection of tested relationships; where there is greater flexibility in designs, definitions, outcomes, and analytical modes; when there is greater financial and other interest and prejudice; and when more teams are involved in a scientific field in chase of statistical significance. Simulations show that for most study designs and settings, it is more likely for a research claim to be false than true. Moreover, for many current scientific fields, claimed research findings may often be simply accurate measures of the prevailing bias. In this essay, I discuss the implications of these problems for the conduct and interpretation of research.

factors that influence this problem and some corollaries thereof.

## Modeling the Framework for False Positive Findings

Several methodologists have pointed out [9–11] that the high rate of nonreplication (lack of confirmation) of research discoveries is a consequence of the convenient, yet ill-founded strategy of claiming conclusive research findings solely on the basis of a single study assessed by formal statistical significance, typically for a *p*-value less than 0.05. Research is not most appropriately represented and summarized by *p*-values, but, unfortunately, there is a widespread notion that medical research articles

## It can be proven that most claimed research findings are false.

should be interpreted based only on *p*-values. Research findings are defined here as any relationship reaching formal statistical significance, e.g., effective interventions, informative predictors, risk factors, or associations. “Negative” research is also very useful.

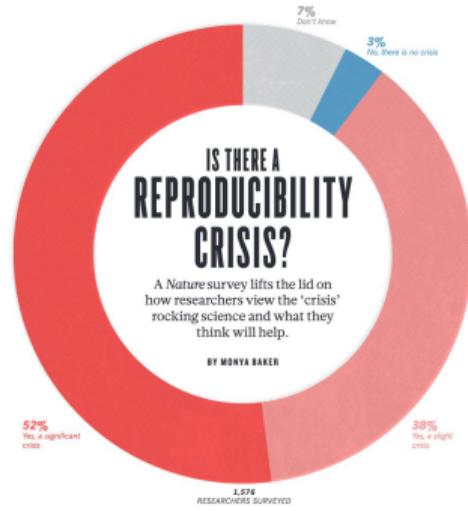
is characteristic of the field and can vary a lot depending on whether the field targets highly likely relationships or searches for only one or a few true relationships among thousands and millions of hypotheses that may be postulated. Let us also consider, for computational simplicity, circumscribed fields where either there is only one true relationship (among many that can be hypothesized) or the power is similar to find any of the several existing true relationships. The pre-study probability of a relationship being true is  $R/(R+1)$ . The probability of a study finding a true relationship reflects the power  $1-\beta$  (one minus the Type II error rate). The probability of claiming a relationship when none truly exists reflects the Type I error rate,  $\alpha$ . Assuming that *c* relationships are being probed in the field, the expected values of the  $2 \times 2$  table are given in Table 1. After a research finding has been claimed based on achieving formal statistical significance, the post-study probability that it is true is the positive predictive value, PPV. The PPV is also the complementary probability of what Wacholder et al. have called the false positive report probability [10]. According to the  $2 \times 2$  table, one gets  $PPV = (1-\beta)R/(1-\beta)R + \alpha(1-R)$ .

## Crisis elements

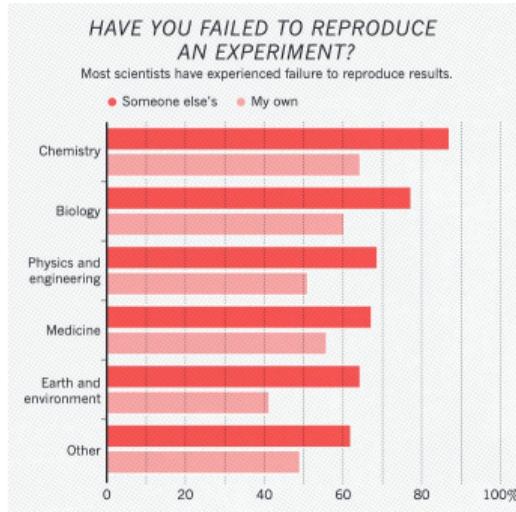
- Highlighted around 2005
- Since 2010 more articles related to the non reproducibility
- Medicine is one of the most impacted discipline

# Reproducibility crisis

2016



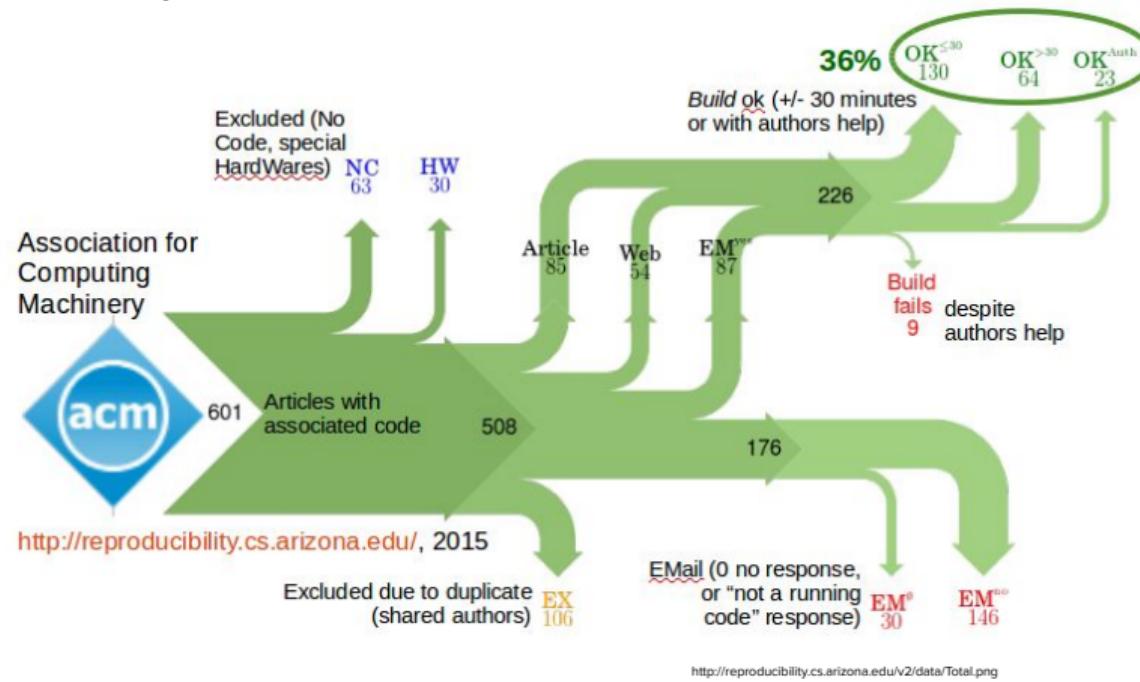
Baker, M. 1,500 scientists lift the lid on reproducibility. *Nature* 533, 452–454 (2016). <https://doi.org/10.1038/533452a>



4

<https://doi.org/10.1038/533452a>

## Also in computer sciences



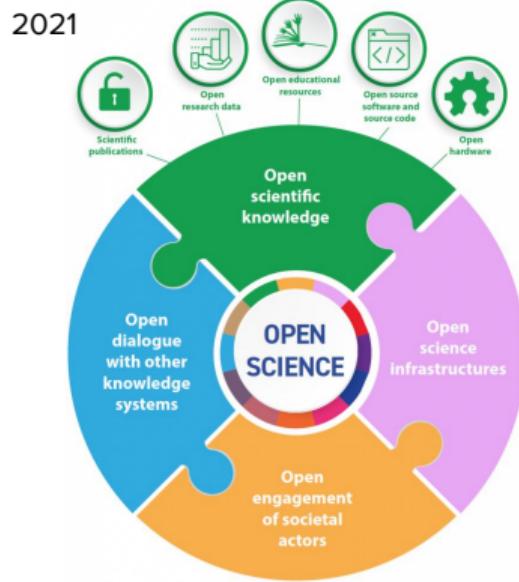
# Long term negative impact of retracted papers

Article	Year of retraction	Citing Articles before retraction	Citing Articles after retraction	Total cites (journals indexed by Web of Science)
1. Primary Prevention of Cardiovascular Disease with a Mediterranean Diet. N ENGL J MED; APR <b>2013</b> . Estruch R, et al.	2018	1919	816	2735
2. Ileal-lymphoid-nodular hyperplasia, non-specific colitis, and pervasive developmental disorder in children. LANCET; FEB 28 <b>1998</b> . Wakefield AJ, et al.	2010	642	867	1509
3. Visfatin: A protein secreted by visceral fat that mimics the effects of insulin. SCIENCE; JAN <b>2005</b> . Fukuhara A, et al.	2007	232	1192	1424
4. An enhanced transient expression system in plants based on suppression of gene silencing by the p19 protein of tomato bushy stunt virus. PLANT J; MAR <b>2003</b> . Voinnet O, et al.	2015	896	375	1271
5. Lysyl oxidase is essential for hypoxia-induced metastasis. NATURE; APR <b>2006</b> . Erler JT, et al.	2020	977	81	1058

Retraction Watch : Top 10 most highly cited retracted papers  
<https://retractionwatch.com/the-retraction-watch-leaderboard/top-10-most-highly-cited-retracted-papers/>

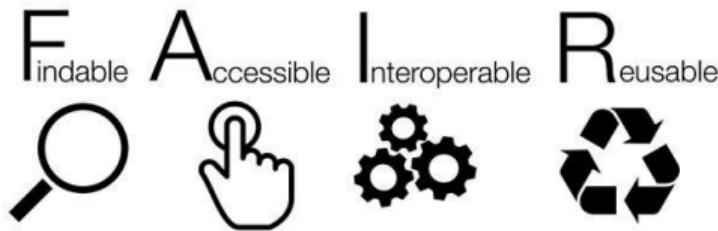
6

# A way out: Open science and FAIR principles



Graphic on page 11. [UNESCO Recommendation on Open Science](#). CC BY IGO 3.0 C. Green

2016



Wilkinson, M., Dumontier, M., Aalbersberg, I. et al. The FAIR Guiding Principles for scientific data management and stewardship. *Sci Data* **3**, 160018 (2016).  
<https://doi.org/10.1038/sdata.2016.18>

7

## FAIR history

- Born in 2016 with *The FAIR Guiding Principles for scientific data management and stewardship*
- How to build, stock, share, use and publish data
- Make criteria to better use our data

. <https://doi.org/10.1038/sdata.2016.18>

# SCIENTIFIC DATA

Amended: Addendum

OPEN

## SUBJECT CATEGORIES

- » Research data
- » Publication characteristics

Received: 10 December 2015

Accepted: 12 February 2016

Published: 15 March 2016

## Comment: The FAIR Guiding Principles for scientific data management and stewardship

Mark D. Wilkinson et al.<sup>#</sup>

There is an urgent need to improve the infrastructure supporting the reuse of scholarly data. A diverse set of stakeholders—representing academia, industry, funding agencies, and scholarly publishers—have come together to design and jointly endorse a concise and measurable set of principles that we refer to as the FAIR Data Principles. The intent is that these may act as a guideline for those wishing to enhance the reusability of their data holdings. Distinct from peer initiatives that focus on the human scholar, the FAIR Principles put specific emphasis on enhancing the ability of machines to automatically find and use the data, in addition to supporting its reuse by individuals. This Comment is the first formal publication of the FAIR Principles, and includes the rationale behind them, and some exemplar implementations in the community.

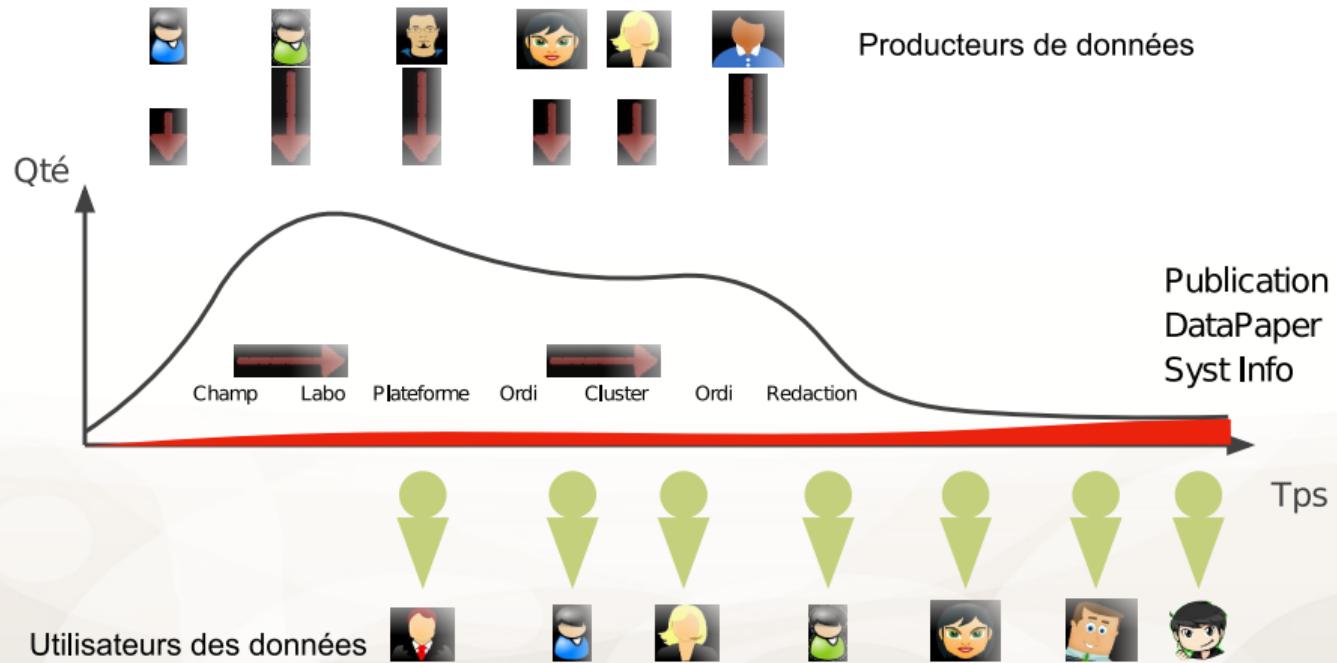
<https://doi.org/10.1038/sdata.2016.18>

## Apply FAIR TO

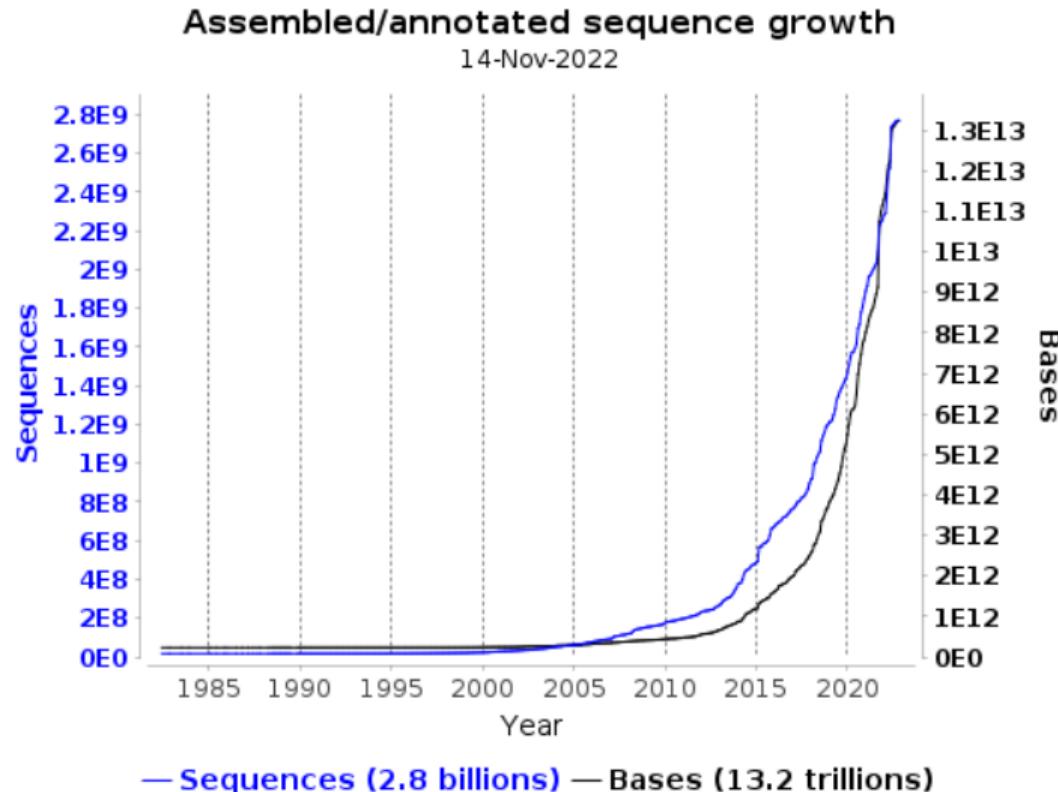
### ■ Your DATA

- Data lifecycle
- Data Management Plan (DMP)
- Metadata
- Data storage

# Focus on the data



# Focus on the data



<https://www.ebi.ac.uk/ena/browser/about/statistics>

## Apply FAIR TO

- Your DATA

- Data lifecycle
- Data Management Plan (DMP)
- Metadata
- Data storage

- Your scripts, environment...

- Objective of this training

# FAIR principles

F  
indable



By 維基小霸王 - Own work, CC BY-SA 4.0,  
<https://commons.wikimedia.org/w/index.php?curid=88894774>

PID  
Repository

8

<https://doi.org/10.1038/sdata.2016.18>

## To be Findable

- (meta)data are assigned a globally unique and persistent identifier
- data are described with rich metadata
- metadata clearly and explicitly include the identifier of the data it describes
- (meta)data are registered or indexed in a searchable resource

# FAIR principles

F  
indable



By 維基小霸王 - Own work, CC BY-SA 4.0,  
<https://commons.wikimedia.org/w/index.php?curid=88894774>

A  
ccessible



<https://nitsfirstworldproblems.tumblr.com/post/147555650875/i-can-t-reach-the-top-shelves-of-the-kitchen>

PID  
Repository

Protocols  
(free, open, auth.)

9  
<https://doi.org/10.1038/sdata.2016.18>

## To be Accessible

- (meta)data are retrievable by their identifier using a standardized communication protocol
- the protocol is open, free, and universally implementable
- the protocol allows for an authentication and authorization procedure, where necessary
- metadata are accessible, even when the data are no longer available

# FAIR principles

**F**indable



By 維基小霸王 - Own work, CC BY-SA 4.0,  
<https://commons.wikimedia.org/w/index.php?curid=88894774>

**A**ccessible



<https://nitsfirstworldproblems.tumblr.com/post/147555650875/i-can-t-reach-the-top-shelves-of-the-kitchen>

**I**nteroperable



By Unknown author - Popular Science Monthly Volume 88, Public Domain  
<https://commons.wikimedia.org/w/index.php?curid=22614407>

PID  
Repository

Protocols  
(free, open, auth.)

Standards  
(format, vocabulary)

10  
<https://doi.org/10.1038/sdata.2016.18>

## To be Interoperable

- (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.
- (meta)data use vocabularies that follow FAIR principles
- (meta)data include qualified references to other (meta)data

# FAIR principles

**F**indable



By 維基小霸王 - Own work, CC BY-SA 4.0,  
<https://commons.wikimedia.org/w/index.php?curid=88894774>

**A**ccessible



<https://nillsfirstworldproblems.tumblr.com/post/147555650875/i-can-t-reach-the-top-shelves-of-the-kitchen>

**I**nteroperable



By Unknown author - Popular Science Monthly Volume 88, Public Domain  
<https://commons.wikimedia.org/w/index.php?curid=22614407>

**R**eusable



By Sun Ladder - Own work, CC BY-SA 3.0,  
<https://commons.wikimedia.org/w/index.php?curid=5746428>

PID  
Repository

Protocols  
(free, open, auth.)

Standards  
(format, vocabulary)

Metadata  
License  
Origin

11  
<https://doi.org/10.1038/sdata.2016.18>

## To be Reusable

- meta(data) are richly described with a plurality of accurate and relevant attributes
- (meta)data are released with a clear and accessible data usage license
- (meta)data are associated with detailed provenance
- (meta)data meet domain-relevant community standard

# FAIR tools

**F**indable



**A**ccessible



**I**nteroperable



**R**eusable



Data

Software  
and  
analyses

# How to integrate FAIR concepts in my work ?

Some tools for reproducible research



# FAIR session with AuBi

## Objectives

- Discover FAIR practices
- Discover tools for best practices
- Learn tools and best practices

# FAIR session with AuBi

## Objectives

- Discover FAIR practices
- Discover tools for best practices
- Learn tools and best practices
- 5 sessions for courses and practices
  - Day 1 : Introduction to FAIR training and Git
  - Day 2 : Git practice
  - Day 3 : Encapsulation course (conda and Docker)and training
  - Day 4 : Encapsulation course (Singularity) and workflow course and practice
  - Day 5 : Documentation course and training

## Contents

- Introduction to FAIR practices

## Contents

- Introduction to FAIR practices
- Code control using Git 

  - Git environment
  - Gitlab and Github  

## Contents

- Introduction to FAIR practices
- Code control using Git 

  - Git environment
  - Gitlab and Github 

- Encapsulation process
  - Conda environment and packages use 
  - Containers as docker & singularity 
  - Reproducible workflow using snakemake 

## Contents

- Introduction to FAIR practices
- Code control using Git 

  - Git environment
  - Gitlab and Github 

- Encapsulation process
  - Conda environment and packages use 
  - Containers as docker & singularity 
  - Reproducible workflow using snakemake 
- Literate programming and documentation
  - Markdown syntax 
  - Rmarkdown for R 
  - Jupyterlab for Python 

## Good practice on...

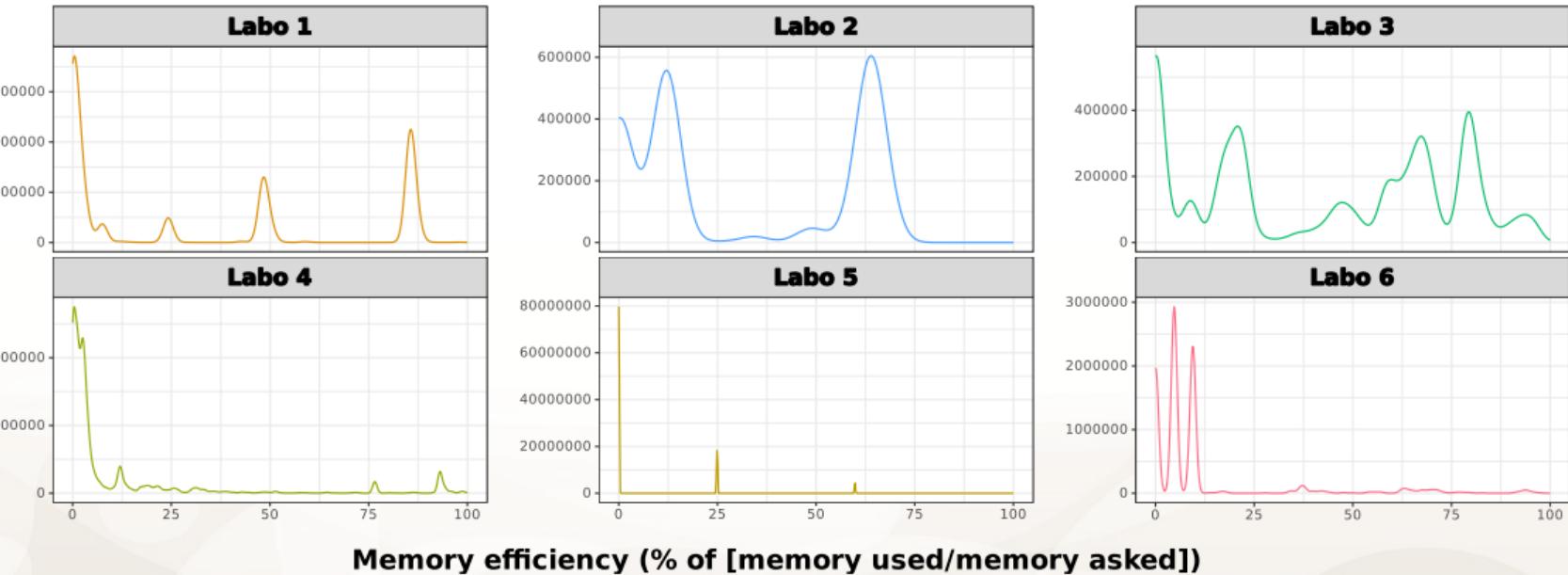
- Must have good practices on data usage
- Must have good practices on tool & code usage

Then...

Must have good practices on computing ressources usage ?

## Efficiency of jobs on used RAM per account

Count of jobs



## Efficiency of jobs on used CPU per account

Count of jobs

