https://moodle.france-bioinformatique.f r/course/view.php?id=38









FAIR principles applied to bioinformatics

MESSAK Imane









Hello !









Part 1 : Context

- Open Science
- The crisis of reproducibility

Part 2 : Basic concepts

- The definitions of reproducibility
- FAIR Principles

Part 3 : A solution proposal

- A solution proposal
- FAIR bioinfo training





Open Science, the crisis of reproducibility







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





... to create more knowledge



PHPKB, Defining Knowledge, Information, Data phpkb.com/kb/article/defining-knowledge-information-data-239

Knowledge

Is subjective Has meaning for a specific purpose Is processed and understood

Is not quantifiable, there is no knowledge overload

Information

Should be objective

Has a meaning

Is processed

Is quantifiable, there can be information overload

Data

Is objective

Has no meaning

Is unprocessed

Is quantifiable, there can be data overload





But how to produce good data?



Repeatable	Replicable		Reproducible			
* <u>(j. (j.</u>						
same lab	dif	ferent labs	ch	differe	nt labs	chos



doi.org/10.1038/s41597-020-0486-7

From collecting to analyzing data

→ From preserving to re-using data





INISTELE FRANCAIS DE BIOINFORMATION

UNESCO Recommendations



Objective: Make research accessible to everyone

- Not just access to the knowledge itself
- The entire process of its creation and dissemination
- The possibility of **reuse**
- Open dialogue with all stakeholders, interdisciplinarity
- Commitment to and from society



Ambitions

- Democratize access to knowledge
- Make science more cumulative, strongly supported by data, and more transparent
- Increase research efficiency by avoiding duplicated efforts and reusing data or scientific material
- Promote scientific advancements and innovation
- Foster public trust in science







In biology,

76 %

of the researchers surveyed

failed to reproduce results.



Monya Baker, 2016 https://doi.org/10.1038/533452a







Collberg et al. 2015







nature genetics

DRITHARY Wylie Vale AVIAN INFLUENZA Shift expertise FARTH SYSTEMS Past climate STORY OF SCIENCE Descartes

lost letter tracked using

Google p.540

give valuable clues to future

warming p.537

	1		
		1.2	
		2.81	
- Party and the			

Many landmark findings in preclinical oncology research are not reproducible, in part because of inadequate cell lines and animal models.

Raise standards for preclinical cancer research

C. Glenn Begley and Lee M. Ellis propose how methods, publications and incentives must change if patients are to benefit.

fforts over the past decade to characterize the genetic alterations in human cancers have led to a better understanding of molecular drivers of this complex set of diseases. Although we in the cancer field hoped that this would lead to more effective drugs, historically, our ability to translate cancer research to clinical success has been remarkably low1. Sadly, clinical

to track mutations where

they emerge p.534

trials in oncology have the highest failure rate compared with other therapeutic areas. Given the high unmet need in oncology, it is understandable that barriers to clinical development may be lower than for other disease areas, and a larger number of drugs with suboptimal preclinical validation will enter oncology trials. However, this low success rate is not sustainable or acceptable, and

investigators must reassess their approach to translating discovery research into greater clinical success and impact. Many factors are responsible for the high failure rate, notwithstanding the inherently difficult nature of this disease. Certainly, the limitations of preclinical tools such as inadequate cancer-cell-line and mouse models2 make it difficult for even

and an elusive stress

hormone p.542

29 MARCH 2012 | VOL 483 | NATURE | 531 © 2012 Macmillan Publishers Limited. All rights reserved



Repeatability of published microarray gene expression analyses

John P A Ioannidis¹⁻³, David B Allison⁴, Catherine A Ball⁵, Issa Coulibaly⁴, Xiangqin Cui⁴, Aedín C Culhane^{6,7}, Mario Falchi^{8,9}, Cesare Furlanello¹⁰, Laurence Game¹¹, Giuseppe Jurman¹⁰, Jon Mangion¹¹, Tapan Mehta⁴, Michael Nitzberg⁵, Grier P Page^{4,12}, Enrico Petretto^{11,13} & Vera van Noort¹⁴



Alsheikh-Ali et al. PLoS ONE (2011) Nekrutenko & Taylor, Nature Genetics (2012) Begley & Ellis Nature (2012)







An interesting article ...



Christophe Monnet 🖾, Valentin Loux, Jean-François Gibrat, Eric Spinnler, Valérie Barbe, Benoit Vacherie, Frederick Gavory Edith Gourbeyre, Patricia Siguier, Michaël Chandler, Rayda Elleuch, Françoise Iriinger 🔯, Tatiana Vallaeys 🔯

Published: November 24, 2010 • https://doi.org/10.1371/journal.pone.0015489

... but a deceptive M&M

The genome sequences of *Arthrobacter aurescens* TC1 [24] (accession number CP000474), *Arthrobacter* sp. strain FB24 (accession number CP000454), a strain isolated from a xylene and chromate enriched soil microcosm [25], and *A. chlorophenolicus* A6 (accession number CP001341), a strain capable of degrading high concentrations of 4-chlorophenol [26], were used for comparative genomic analyses. The sequence data from the strains FB24 and A6 were produced by the US Department of Energy Joint Genome Institute (<u>http://www.jgi.doe.gov/</u>) in collaboration with the user community. Genome comparisons were performed using Origami, an <u>In-house</u> tool developed for microbial genome comparison. Orthologs were defined as reciprocal best hits with an e-value lower than 10^{-3} . Transposases were excluded from the analysis. Core genes were defined as orthologs shared between the four *Arthrobacter* strains. Synteny was studied using an <u>In-house</u> developed tool, Align, using dynamic programming to search conserved gene trains allowing gaps and "mismatches" (homology relation instead of orthology). Circular representation of the genome was produced using the Circos software [27].







EVOLVING TOWARDS AN ERA OF OPEN RESEARCH











Michener (1997)

MISSING DATA

As research articles age, the odds of their raw data being extant drop dramatically.















52.39%

3.52%

Indicator #4 0.04%

0.24%

Likely cost of not having FAIR research data

The basic concepts of reproducibility





Reproducible research, Repeatability, Replicability, Reproducibility, Replication...

Overlapping concepts

 \Rightarrow many definitions!

Definitions from the computing Machinery (2016) :

Repeatability: same team, same experimental design Reproducibility: different teams, same experimental design Replicability: different teams, different experimental designs

https://www.researchgate.net/publication/323118701 _Terminologies_for_Reproducible_Research https://www.acm.org/publications/policies /artifact-review-and-badging-current





Matrix for the Reproducibility of Whitaker (2017):

		Data		
		Same	Different	
de	Same	Reproducible	Replicable	
°	Different	Robust	Generalisable	

https://doi.org/10.6084/m9.figshare.5443201.v1, Slide 7



The Turing Way project illustration by Scriberia. Used under a CC-BY 4.0 licence. DOI: <u>The Turing Way Community & Scriberia (2024)</u>.





Data & Code : Well known for data









Reminder : FAIR data principles

In brief, FAIR data should be:



Findable: The first step in (re)using data is to find it! Descriptive metadata (information about the data such as keywords) is essential.

Accessible: Once the user finds the data and software they need to know how to access it. Data could be openly available but it is also possible that authentication and authorisation procedures are necessary.



Interoperable: Data needs to be integrated with other data and interoperate with applications or workflows.



Reusable: Data should be well-described so that they can be used, combined, and extended in different settings.



The Turing Way project illustration by Scriberia. Used under a CC-BY 4.0 licence. DOI: <u>The Turing Way Community & Scriberia (2024)</u>.







- \bigstar Open License
- \star machine REadable
- \star Open Format
- \star Uniform Resource Identifier
- \star Linked Data

Bernes-Lee







Data & Code : but for the code ?











Like data, software is backed up, but it's not "just" data. They are **living** and **complex** entities.

	Data	Software
are	facts, observations	creations
produce	evidence	an executable tool
change?	no (unless underlying hardware changes)	yes: continuous maintenance & updates
lifespan	long-term necessary when the experiment cannot be reproduced (cost of collection or validation)	short: often rebuilt to use other software (complex dependencies), replaceable (another performs the task better)

Katz DS, Niemeyer KE, Smith AM, Anderson WL, Boettiger C, Hinsen K, Hooft R, Hucka M, Lee A, Löffler F, Pollard T, Rios F. 2016. Software vs. data in the context of citation. PeerJ Preprints 4:e2630v1 <u>https://doi.org/10.7287/peerj.preprints.2630v1</u>





FAIR principles



Findable (for humans and machines)

- Unique identifier (e.g., a DOI with Zenodo)
- Metadata describing the analysis and tools; they are FAIR, searchable, and indexable (e.g., README)
- Available on a forge (GitHub, GitLab) and in an archive repository (SWH)

Accessible

- Retrievable using a standardized protocol
- Open, free and universally implementable
- Metadata accessible, even when the software is no longer available (SWH)

Interoperable

- Use a formal, accessible, shared and broadly applicable language
- Tool cooperation (snakemake, conda, and docker) locally or on a server (cloud or cluster)

B Reusable

- Clear licensing and access rights (according to employer, DPO)
- Have a detailed provenance
- Follows community standards







OpenAire





.

A suggested solution to make your project reproducible









The Turing Way project illustration by Scriberia. Used under a CC-BY 4.0 licence. DOI: <u>The Turing Way Community & Scriberia (2024)</u>.







The pillars of reproducibility

	Reproducibility						
Versioning Sharing Archiving	Create a virtual environment	Installing tools using a package manager	Create analysis scripts	Offloading analysis to a server	Portability of results exploration	Editing the analyses reports	

based on

The five pillars of computational reproducibility: Bioinformatics and beyond Mark Ziemann, Pierre Poulain, Anusuiya Bora, OSF preprint, 2023





- Having the right version of the code _
- **Temporal insight** _
- Openness to the community _

"FINAL".doc







FINAL.doc!





FINAL_rev.6.COMMENTS.doc





FINAL_rev.8.comments5. CORRECTIONS.doc







FINAL_rev.18.comments7. FINAL_rev.22.comments49. corrections9.MORE.30.doc corrections.10.#@\$%WHYDID ICOMETOGRADSCHOOL????.doc

WWW. PHDCOMICS. COM





Pillar 1 - Versioning, sharing, and archiving code



DOI 10.



INSTITUT FRANCAIS DE BIOINFORMATIQUE



FRANCE

. . .



GitHub

Bitbucket





ady 38, 2022		Autors Realision		
Rtools4Emergen			101.03	n
🕽 linisinal, Driana 🌍 lifensala, imame, 🍘 can hitilare, Jacopar	O Interior, Transa			
wikasikapi ta ampilyitata analara uf two tribiside ce a	rabee		Extinuities Extintion Permit	11111
Proven		~	national consenture for genomic auroritance and	
Commences of the			research about COVID-18 and other amonging	
D. Roldgrov D. pilgrov D. pilgrov D. pilgrov D. pilgrov D. pilgrov D.Alwaved saturation Occommon saturation Occommon Decommon		28 Oyen Hitchyns 1,8 48 10,3 68 1,3 48 1,8 48	History Constant Indite Parquis de Sonternatique OFBLI E(ORD/Transe	
COMMEDMOE Defection, by camping, does A Defection, by camping, does A Defection, by camping, before A Defection, by camping, before A Defection, part, month? Defection, part, participant Defection, part, participant Defection, part, participant Defection, part, participant Defection, part, participant		903 Nyan 13 40 12 45 43 46 42 46 62 46	31 * ress ter man de	0 Arcontext
 Class, series per report? Classifies cons.8 		2348		
 Plan second and the 		414	Indexed #	
File course		÷	Open	AIRE
kana	Sec. 1			
staalakeenergevol 3.55ar	243.88	Atomic		-
with dear of the antide time to			45y 54, 3600	
	1.150	Access.	an a	•
tablement of the st	1310	Around	Report(s)	/
-due-District Statistics (1)-9			Intrative Constrained	ana ke de
Automorphed 23.29	1210	alors disting	and an analysis of	
			1944	
			CR Constant Community, Million	dan 12 Meruhani
Obstant to the version	IS COMMONS	~ ~	Versions	
No ch	riore.		Second 23 10109, Second Method	A# 19,3003
	/		Charl content" for second allo Solid and all shape would be for local and all shape would be for local	men traing to im ngaata di anima, an Inalman
/			Chare Chare	•
			CEP 44 Arternal, Olarea, Hironak, Irina Januper, & Orrender, Thomas Hindor/Energies (2008), Den Mijoc/Min.org/10.0013/Denerg	ne, van Helslen, - (2022) India 6-6907 (2017
	_			

1 Dat Plancks (1) B kins		
# Rando Mass - Holdon 1		Grany Americal Brane
 Ty write are distributed as a second standard and a second second	ri i; Tama Denke oʻnilapatana, es	22400
Phe .	1544	2m
B1		
in inc		
• no		
R size		
B date		
in sports		
8 Multiple	242.0	of lyins
R ellever	1911	te bea
R gibbiged		1248
6 accuracy		678
8 mouth	797.0	0.03
B states	494.6	1048
E SINGHT		648
& Malletand		20dates
8 Kolschept fort	797.0	anigen
READ/ME/red		
Rtools4Emergen		
Warning		
The participle is a sourcedow, it sims at providing accurate althoughploates his programming propriorities of announline of feature although pro-	Ingenet Indocelle (NGE)=08	index.bdtjazzanticke prz. ot bai
Intro		
The Random and Antoin Strations and a Different Academic	Manager di Sana Schernarigan	0-61072147
DESCRIPTION in the location of the advanced providence of providence of the second providence of	encentantly in beaching expression	confusion project (MIRCO).
And a state of the		



How?

- Freezing the environment
- Sharing the environment



Advantages / Disadvantages

- + Code backup
- + Easy for sharing
- + Automatic version management

- Not easy for novices





How?

- Freezing the environment
- Sharing the environment



Advantages / Disadvantages

- + Fast and lightweight
- + Portable
- + Easy to share and deploy

CONDA



- With an up-to-date system
- [docker] Accepted in your structure?





How?

- Having the right version of the tools used
- Installing them simply



ΡΙΧΙ

Advantages / Disadvantages

- + Simple manager to install
- + Simple package installation
- + Version management
- Can be heavy (miniconda solution)
- Missing packages (R)





- Having a reproducible analysis script
- Not redoing what has already been done
- Parallelizing



How?

nextflow

Advantages / Disadvantages

- + Job management
- + Powerful and fast
- + Capable of using Conda environments
- + Parallelizable
 - A logic to learn (syntax less simple than shell script)





- Controlled environment
- Offloading analysis

INSTITUT FRANÇAIS DE BIOINFORMATIQUE

Advantages / Disadvantages

- + Easy to set up
- + Increased power (cloud or cluster)
- + For everyone



- Not easy for novices
- Attention to sensitive data





How?

Pillar 6 - Portability of results exploration

Why?

- Making exploration simple
- Easy to share



How?



quarto[°]

Advantages / Disadvantages

- + Portable (HTML)
- + Accessible everywhere
- + Interactive (configurable, dynamic graphs, etc.)
 - Mix of languages





- Having a record of the analysis (date, time, parameters, etc.)
- Storing tool versions



How?



quarto[°]

Advantages / Disadvantages

- + Simple syntax (Markdown)
- + Sharing (PDF, HTML, etc.)
 - Rare visualization problems using LATEX



_





Conclusion



1. Data

- a. Do I do even it if the data are not 100% FAIR?
- b. How to manage large volumes of data?
- c. How to manage data and metadata updates?

2. Code

- a. How to ensure that the code will always be accessible?
- b. Is it acceptable to make adaptations? To what extent?
- c. Provide all the code? (valorization, creation of a start-up, etc.)

3. Computation time (days, months, years)

4. Skill and sensitivity

- a. Willingness but technical inability to do so
- b. "Why bother, it's useless"
- c. too long

5. The coverage

a. Should everything be made reproducible?

6. When to do it?

- a. At the beginning? But what if it doesn't work?
- b. At the end?







Proposal of a solution that helps make any analysis protocol reproducible.

Reproducibility is an added value for your work !













And you? How many pillars support your reproducibility?

Reproducibility







Trainings by the IFB





IFB trainings on the theme of FAIR

- FAIR principles for research data management in life sciences (*FAIRdata*)
- The FAIR principles in a Bioinformatics project (*FAIRbioinfo*)



https://moodle.france-bioinformatique.fr/









Thank you for your attention !









