



Conference of Bioinformatics  
and Computational Biology

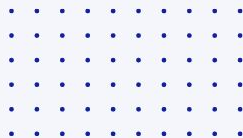


# Leveraging **Open Science in Machine Learning and Bioinformatics**

**Batool Almarzouq**



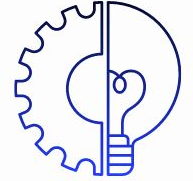
Open Science Community Saudi Arabia  
مجتمع العلوم المفتوحة في المملكة العربية السعودية





# Thank you!

---



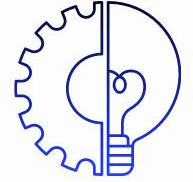
- Yo Yehudi
- Piv Gopalasingam
- John Ogunsola
- Emily Angiolini
- Turing Way



THANK YOU



# Batool Almarzouq



Computational Biologist



Honorary research Fellow



UNIVERSITY OF  
LIVERPOOL

Community Lead



Education Committee



Core Contributor



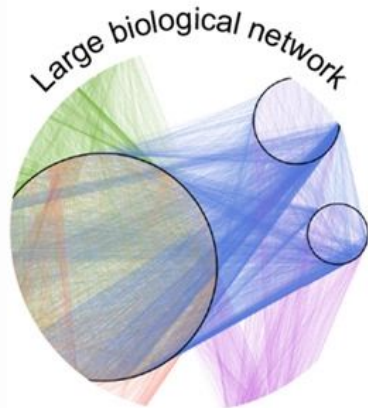
Global Team



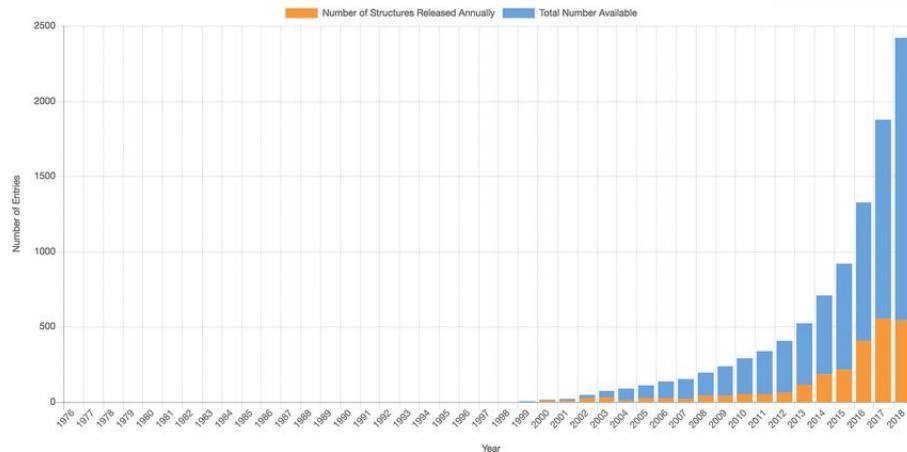
Subject Matter Experts (SME)



# Biology is **data-intensive** science



Descriptive biology  
(data integration)



Source: Miquel Duran-Frigola *et al.* Formatting biological big data for modern machine learning in drug discovery



@batool664, @OpenSciSaudi



BatoolMM



<https://batool-almazouq.netlify.app>

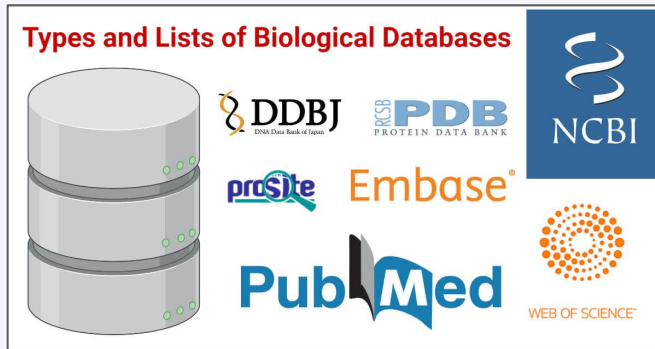




# Can we **go beyond** that?



With the improvements in computing power, storage and exponential growth in biological data, can we **answer what we failed to solve** using conventional methods?



Biological Databases. Image Source: [The Biology Notes](#)



Credit: [The supercomputer Hochleistungsrechner Karlsruhe\\*](#)  
(Author: Amadeus Bramsiepe, KIT)



@batool664, @OpenSciSaudi



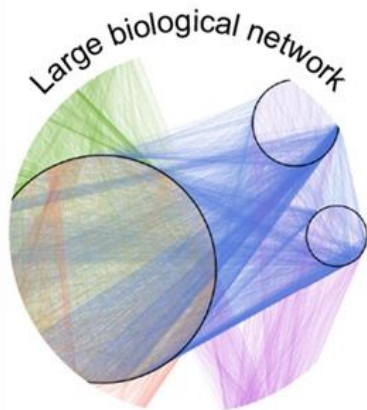
BatoolMM



<https://batool-almazouq.netlify.app>



# Superpower!



Descriptive biology  
(data integration)

Source: Miquel Duran-Frigola *et al.*, [Formatting biological big data for modern machine learning in drug discovery](#)



@batool664, @OpenSciSaudi



BatoolMM

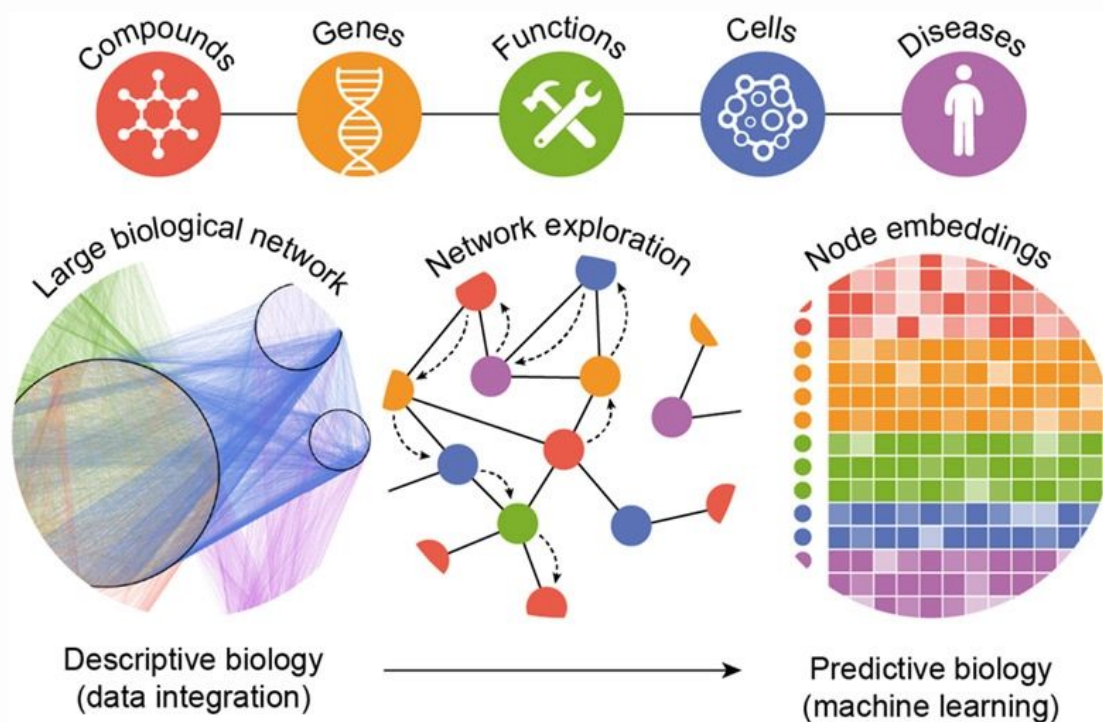


<https://batool-almazouq.netlify.app>



creative commons

# Superpower!



Source: Miquel Duran-Frigola *et al.*, Formatting biological big data for modern machine learning in drug discovery



@batool664, @OpenSciSaudi



BatoolMM

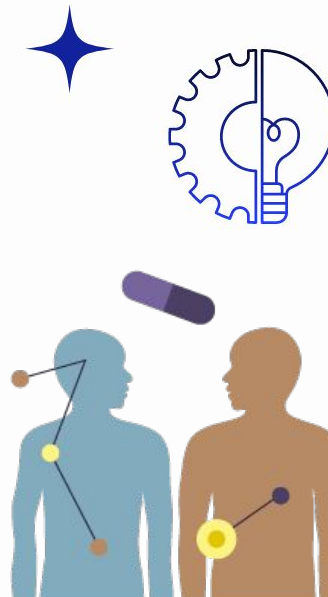


<https://batool-almazouq.netlify.app>



# Success!

	Protein structure prediction	Protein function prediction	Genome engineering	Systems biology and data integration	Phylogenetic inference
Paradigm shifting	✓				
Major success		✓	✓		
Moderate success				✓	
Minor success					✓



**DeepGOPlus: Improved protein function prediction from sequence**



Source: Sapoval, N., Aghazadeh, A., Nute, M.G. et al. [Current progress and open challenges for applying deep learning across the biosciences.](#)



@batool664, @OpenSciSaudi



BatoolMM



<https://batool-almazouq.netlify.app>



creative commons



# If that's the case!



[Source: Techlady](#)



@batool664, @OpenSciSaudi



BatoolMM



<https://batool-almazouq.netlify.app>



# Challenges we have seen so far! ✨



Challenge
Biased results
High infrastructure costs
Lack of explainability
Limited training data
Overfitting
Poor performance on novel data



Source: Sapoval, N., Aghazadeh, A., Nute, M.G. et al. **Current progress and open challenges for applying deep learning across the biosciences**

Cartoon Credit: [Maclo](#).



@batool664, @OpenSciSaudi



BatoolMM



<https://batool-almazouq.netlify.app>



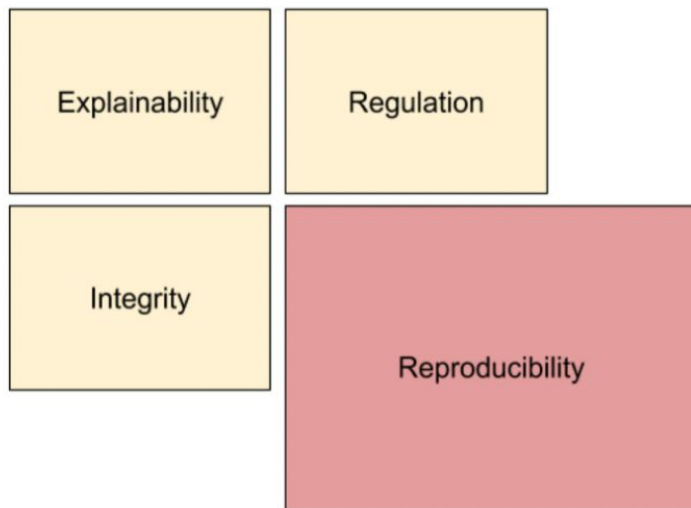
creative commons



# What is trustworthy ML?



## Trustworthy ML



If deep learning is to solve real scientific and technological problems, the community needs to build on each **others' tools and data** by implementing **Open Science practices**.

Resource: [Reproducible Machine Learning by Preeti Hemant](#)



# Reproducibility Crisis!



## The \$28 Billion a Year Research Reproducibility Crisis

📅 October 19, 2015

The biggest crisis in the future of healthcare has talk about, such as payment systems and access conditions like heart disease, dementia, diabetes: treatments that rely upon genomics, nanotechnology that at least half of the research money we invest implications? What are the solutions? We'll show

### Session 2: Bioimaging and Artificial Intelligence

**Session Lead:** Jean-Marie Burel, Senior Software Architect, University of Dundee

**Talk 1:** Yang Zhang, Project Scientist & Project Manager, Carnegie Mellon University (US)

*Nucleome Browser: An integrative and multimodal data navigation platform for 4D Nucleome*

**Talk 2:** Matthew Hartley, Biologie Archive Team Leader, European Bioinformatics Institute (EMBL-EBI)

*Open bioimaging data at scale: publication, analysis and reuse*

**Talk 3:** Josh Moore, Senior Software Architect, University of Dundee

*OME-NGFF (next-generation file format): Zarr as a cloud-native solution for FAIRer bioimaging data*

Alina Goldenberg, Anshul Kundaje, Casey S. Greene, Tamara Broderick, Michael W. Norman, Jeremy T.

Leek, Keegan Korthauer, Wolfgang Huber, Alvis Brazma, Joelle Pineau, Robert Tibshirani, Trevor Hastie, John P. A. Ioannidis, John Quackenbush & Hugo J. W. L. Aerts

[Nature](#) 586, E14–E16 (2020) | [Cite this article](#)

16k Accesses | 85 Citations | 530 Altmetric | [Metrics](#)

Comment | Published: 30 August 2021

## Reproducibility standards for machine learning in life sciences

[Benjamin J. Heil](#), [Michael M. Hoffman](#), [Florian Markowitz](#), [Su-In Lee](#), [Casey S. Greene](#) & [Steph Hicks](#)

[Nature Methods](#) 18, 1132–1135 (2021) | [Cite this article](#)

15k Accesses | 16 Citations | 267 Altmetric | [Metrics](#)



@batool664, @OpenSciSaudi



BatoolMM



<https://batool-almazouq.netlify.app>



creative commons



# Reproducibility Crisis in ML?

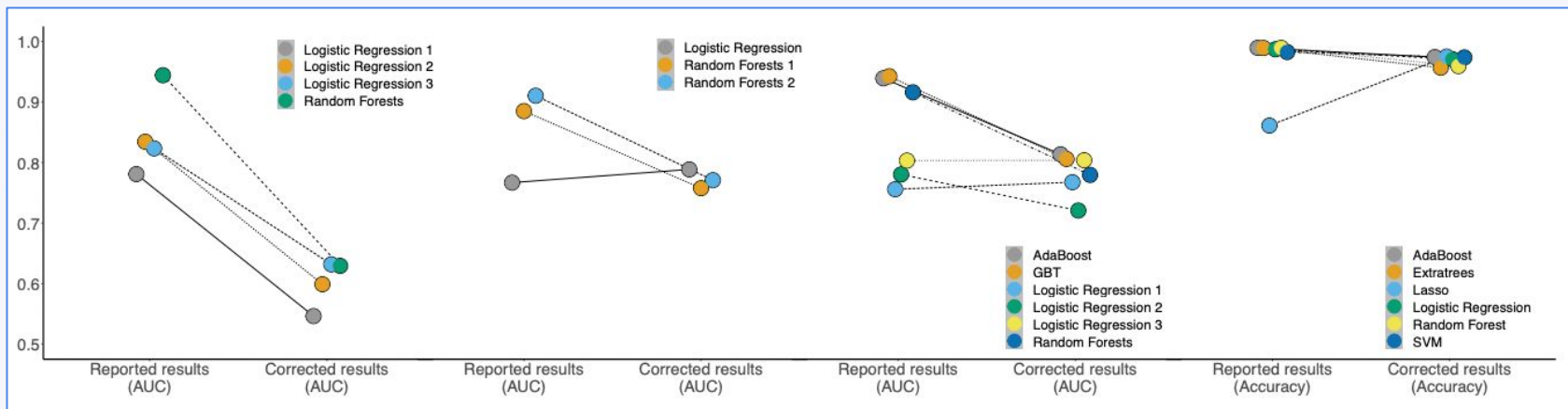


Field	Paper	Number of papers reviewed	Number of papers with pitfalls	[L1.1] No test set	[L1.2] Pre-proc. on train-test	[L1.3] Feature sel. on train-test	[L1.4] Duplicates	[L2] Illegitimate features	[L3.1] Temporal leakage	[L3.2] Non-ind. b/w train-test	Comput. reproducibility issues	Data quality issues	Metric choice issues	Standard dataset used?
Medicine	Bouwmeester et al. (2012)	71	27	o							o			
Neuroimaging	Whelan & Garavan (2014)	-	14	o	o									
Autism Diagnostics	Bone et al. (2015)	-	3						o		o	o	o	o
Bioinformatics	Blagus & Lusa (2015)	-	6	o										
Nutrition Research	Ivanescu et al. (2016)	-	4	o							o	o		
Software Eng.	Tu et al. (2018)	58	11			o			o	o	o	o	o	o
Toxicology	Alves et al. (2019)	-	1		o						o	o		
Satellite Imaging	Nalepa et al. (2019)	17	17					o			o	o	o	o
Tractography	Poulin et al. (2019)	4	2	o							o	o	o	o
Clinical Epidem.	Christodoulou et al. (2019)	71	48		o						o			
Brain-computer Int.	Nakanishi et al. (2020)	-	1	o										o
Histopathology	Oner et al. (2020)	-	1					o						
Neuropsychiatry	Poldrack et al. (2020)	100	53	o	o						o	o		
Medicine	Vandewiele et al. (2021)	24	21		o			o	o	o	o	o	o	o
Radiology	Roberts et al. (2021)	62	62	o	o				o	o				o
IT Operations	Lyu et al. (2021)	9	3				o							o
Medicine	Filho et al. (2021)	-	1			o								
Neuropsychiatry	Shim et al. (2021)	-	1		o					o				
Genomics	Barnett et al. (2022)	41	23		o						o			
Computer Security	Arp et al. (2022)	30	30	o	o	o	o	o	o	o	o	o	o	o

Source: [Leakage and the Reproducibility Crisis in ML-based Science](#) by Sayash Kapoor, Arvind Narayanan



# We should talk about **data leak!**



Source: [Leakage and the Reproducibility Crisis in ML-based Science by Sayash Kapoor, Arvind Narayanan](#)



@batool664, @OpenSciSaudi



BatoolMM



<https://batool-almazouq.netlify.app>



creative commons

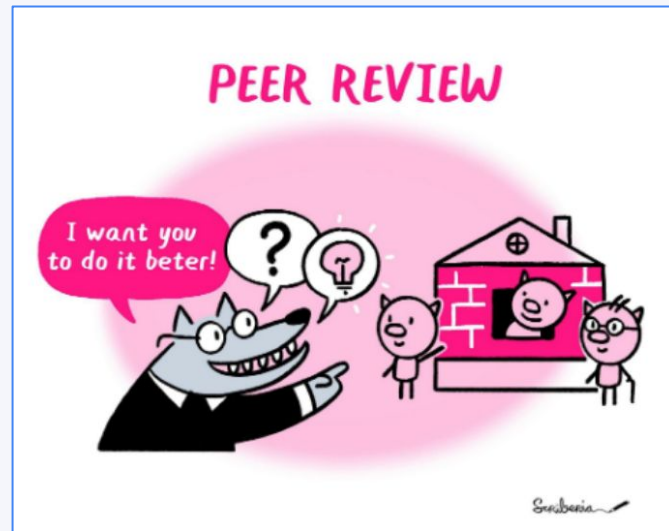


# Peer Review is **not** enough!



“To consult the [experts] **after** an experiment is finished is often merely to ask to **conduct a postmortem examination**. [...] can perhaps say what the experiment died of”

- Ronald Fisher



This Slide is inspired from [Talk by Malvika Sharan](#)

The Turing Way project illustration by Scriberia. Used under a CC-BY 4.0 licence. DOI: [10.5281/zenodo.3332807](https://doi.org/10.5281/zenodo.3332807)



@batool664, @OpenSciSaudi



BatoolMM

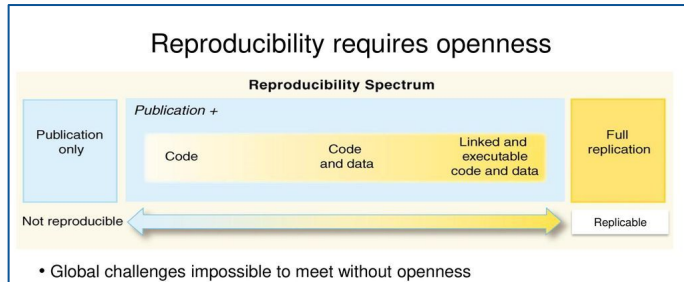
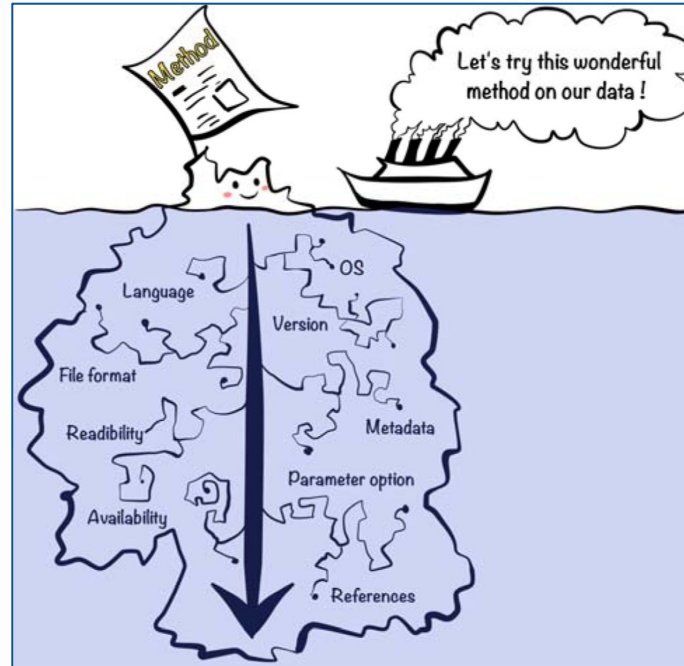


<https://batool-almazouq.netlify.app>



creative commons

# Reproducibility require **Openness!**



[Peng \(2011\)](#)

Credit: Yang-Min Kim, Jean-Baptiste Poline, Guillaume Dumas, Experimenting with reproducibility: a case study of robustness in bioinformatics, *GigaScience*, Volume 7, Issue 7, July 2018, giy077, <https://doi.org/10.1093/gigascience/giy077>



@batool664, @OpenSciSaudi



BatoolMM



<https://batool-almazouq.netlify.app>





# What is Open Science?



'**Open Science**' is an **umbrella** term that encompasses various practices to increase the **visibility, the transparency, useability** of the scientific work.

NEWS | 03 February 2021

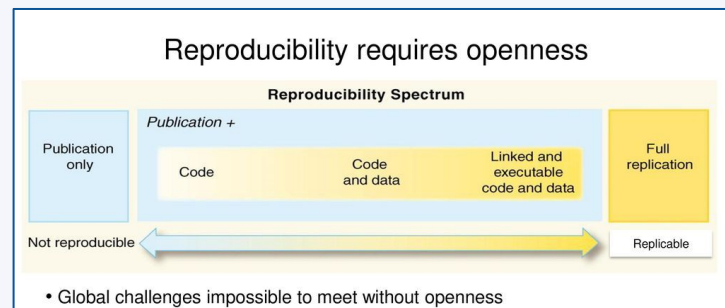
## Scientists call for fully open sharing of coronavirus genome data

Other researchers say that restrictions at the largest SARS-CoV-2 genome platform encourage fast sharing while protecting data providers' rights.

Richard Van Noorden



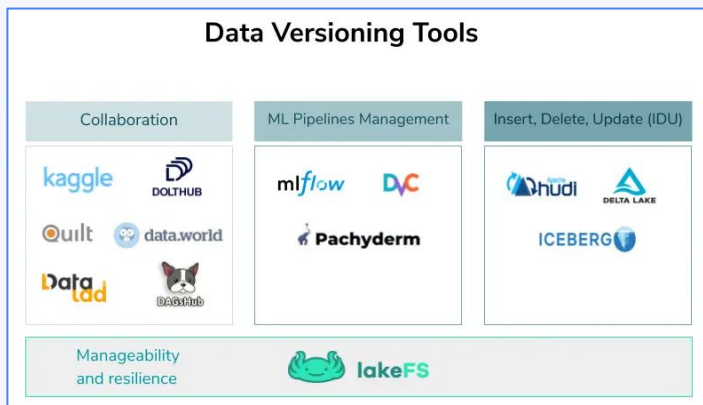
# What can we do to implement Open Science?



[Peng \(2011\)](#)



# Versioning Data and Model



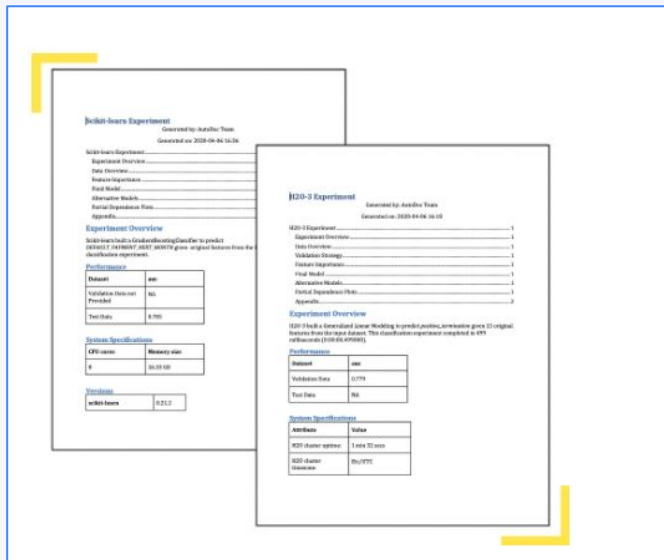
Source: [KDnuggets](#)



Source: [MLOps Guide](#)



# Documentation for Model Training



## Documenting details of:

- How the model was trained will ensure repeatable results.
- Feature transformations, order of features,
- Hyperparameters and the method to select them





# Community-based Projects!



**OpenML**  
A worldwide machine learning lab

Machine learning research should be easily accessible and reusable. OpenML is an open platform for sharing datasets, algorithms, and experiments - to learn how to learn better, together.

I shared a new data set  
I found a better model!

Search OpenML Datasets ▾

[Sign Up](#) to start tracking and sharing your own work. OpenML is open and free to use.

- AI-ready data**  
All datasets are uniformly formatted, have rich, consistent metadata, and can be loaded directly into your favourite environments.
- ML library integrations**  
Pipelines and models can be shared directly from your favourite machine learning libraries. No manual steps required.
- A treasure trove of ML results**  
Learn from millions of reproducible machine learning experiments on thousands of datasets to make informed decisions.

Website: [OpenML](https://openml.org)



@batool664, @OpenSciSaudi



BatoolMM



<https://batool-almazouq.netlify.app>



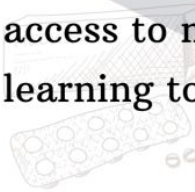
creative commons

# Community-based Projects!



## MISSION

Strengthen the research capacity against infectious and neglected diseases by democratising the access to machine learning tools.



## VISION

A world with egalitarian research capacity and access to healthcare.



# Ersilia

Reference: Ersilia Open Source Initiative – Strategic Plan 2021 – 2023 Main Points (Outreachy Contribution)



@batool664, @OpenSciSaudi



BatoolMM

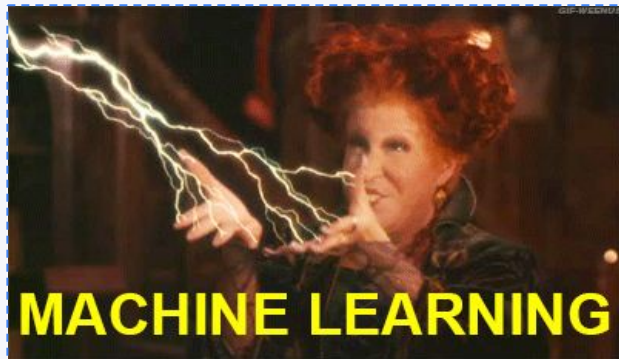


<https://batool-almazouq.netlify.app>



creative commons





“So far, each research community has independently rediscovered these pitfalls. **Without fundamental changes** to research and reporting practices, **we risk losing public trust owing to the severity and prevalence of the reproducibility crisis across disciplines**”

– Sayash Kapoor and Arvind Narayanan



# Discussion!



- What **concerns** do you have about sharing your bioinformatics research and are there any research objects you should not share?
- How to **support researchers and the community** creating a comprehensive open source machine learning environment?
- What is the **potential for innovation** when adopting open science in bioinformatics?



# References



- [Ersilia Open Source Initiative](#)
- [Reproducible Machine Learning by Preeti Hemant](#)
- [Leakage and the Reproducibility Crisis in ML-based Science by Savash Kapoor, Arvind Narayanan](#)
- [AI slipping on tiles: data leakage in digital pathology by Nicole Bussola, Alessia Marcolini, Valerio Maggio, Giuseppe Jurman, Cesare Furlanello](#)
- [Experimenting with reproducibility: a case study of robustness in bioinformatics \(GigaScience\) by Yang-Min Kim, Jean-Baptiste Poline, Guillaume Dumas](#)
- [The Turing Way](#)







**Thank You**

