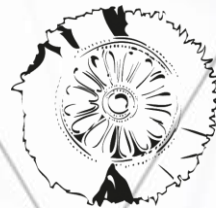# FAIR for Machine Learning;
# Building on the Lessons from FAIR Software

Fotis E. Psomopoulos
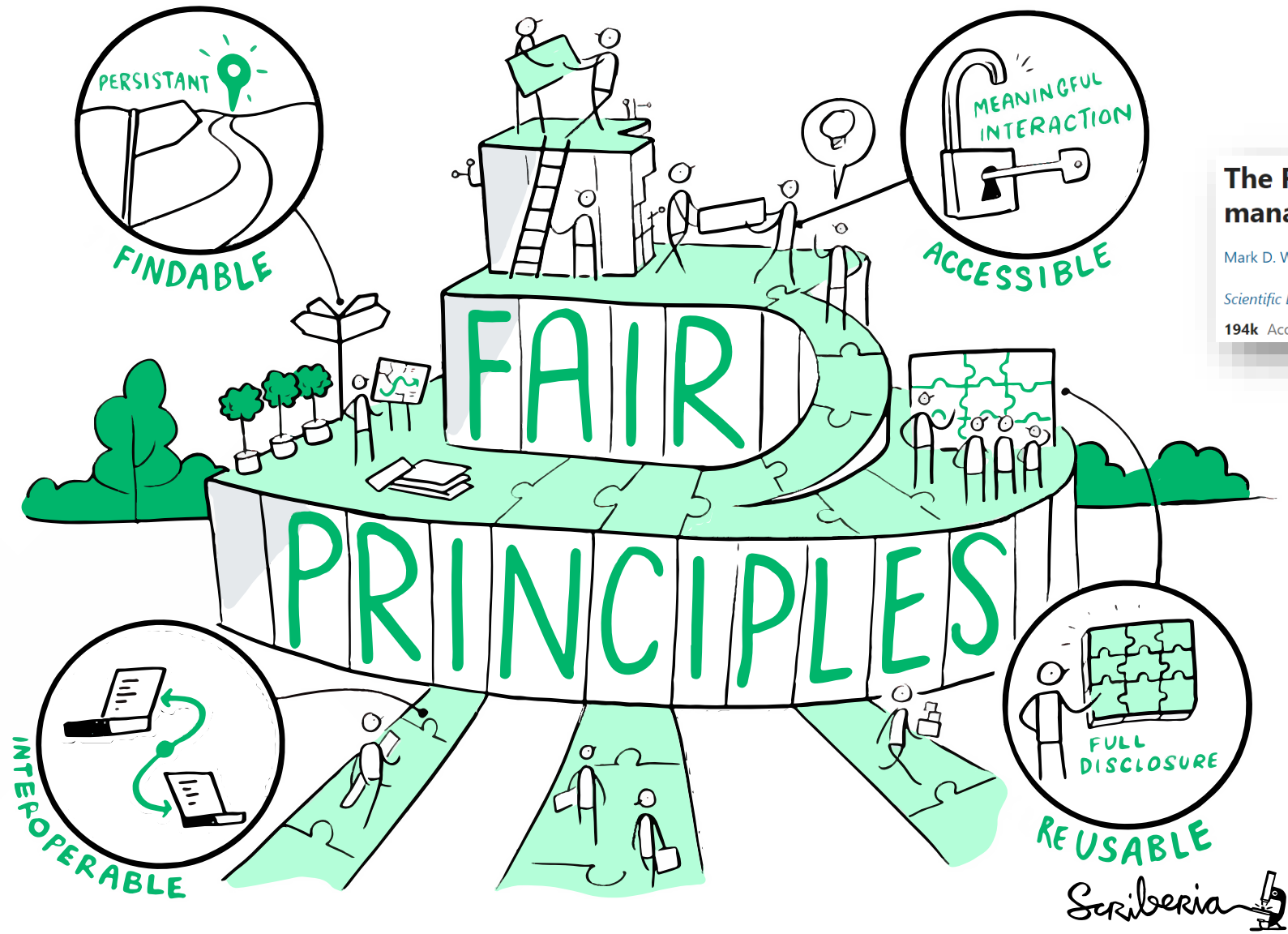
*Institute of Applied Biosciences, CERTH, Greece*

RDA
RESEARCH DATA ALLIANCE

elixir

CERTH
CENTRE FOR
RESEARCH & TECHNOLOGY
HELLAS

INAB
INSTITUTE OF APPLIED BIOSCIENCES
ΙΝΣΤΙΤΟΥΤΟ ΕΦΑΡΜΟΣΜΕΝΩΝ ΒΙΟΕΠΙΣΤΗΜΩΝ
CENTRE for RESEARCH and TECHNOLOGY-HELLAS

# FAIR for non-data objects: some context

- FAIR Principles, at a high level, are intended to **apply to all research objects**; both those used in research and those that are research outputs

- Text in principles often includes "(Meta)data …"
  - Shorthand for "metadata and data …"

- Principles applied via dataset creators and repositories, collectively responsible for creating, annotating, indexing, preserving, sharing the datasets and their metadata

- What about non-data objects?
  - While they can often be stored as data, they are not **just** data

- While high level goals (F, A, I, R) are mostly the same, the details and how they are implemented depend on
  - How objects are created and used
  - How/where the objects are stored and shared
  - How/where metadata is stored and indexed

  *Slide adapted from the presentation of the RDA FAIR4RS steering group at the International Funders Workshop (Nov 2022), https://zenodo.org/doi/10.5281/zenodo.7350198*

- Work needed to define, then implement, then adopt principles

# FAIR for non-data objects: an ongoing effort

## Introducing the FAIR Principles for research software

Michelle Barker, Neil P. Chue Hong, Daniel S. Katz, Anna-Lena Lamprecht, Carlos Martinez-Ortiz, Fotis Psomopoulos, Jennifer Harrow, Leyla Jael Castro, Morane Gruenpeter, Paula Andrea Martinez & Tom Honeyman

*Scientific Data* **9**, Article number: 622 (2022) | Cite this article

---

**DOI:** 10.15497/RDA00065

**Citation and download:** Chue Hong, N. P., Katz, D. S., Barker, M., Lamprecht, A.-L., Martinez, C., Psomopoulos, F. E., Harrow, J., Castro, L. J., Gruenpeter, M., Martinez, P. A., Honeyman, T., et al. (2021). FAIR Principles for Research Software (FAIR4RS Principles). *Research Data Alliance*. DOI: 10.15497/RDA00065

---

Breakout 7 | Data Infrastructures - Organisa... | The FAIR Agenda | WGs Getting started
WG FAIR for **Virtual** Research Environments: FAIR for VREs - The Path Forward
7:30 AM - 9:00 AM
Room E

---

January 01 2020

## FAIR Computational Workflows

Carole Goble , Sarah Cohen-Boulakia, Stian Soiland-Reyes, Daniel Garijo, Yolanda Gil, Michael R. Crusoe, Kristian Peters, Daniel Schober

> Author and Article Information

*Data Intelligence* (2020) 2 (1-2): 108–121.

https://doi.org/10.1162/dint_a_00033

---

## FAIR for AI: An interdisciplinary and international community building perspective

E. A. Huerta, Ben Blaiszik, L. Catherine Brinson, Kristofer E. Bouchard, Daniel Diaz, Caterina Doglioni, Javier M. Duarte, Murali Emani, Ian Foster, Geoffrey Fox, Philip Harris, Lukas Heinrich, Shantenu Jha, Daniel S. Katz, Volodymyr Kindratenko, Christine R. Kirkpatrick, Kati Lassila-Perini, Ravi K. Madduri, Mark S. Neubauer, Fotis E. Psomopoulos, Avik Roy, Oliver Rübel, Zhizhen Zhao & Ruike Zhu

*Scientific Data* **10**, Article number: 487 (2023) | Cite this article

---

## Ten simple rules for making training materials FAIR

Leyla Garcia, Bérénice Batut, Melissa L. Burke, Mateusz Kuzak, Fotis Psomopoulos, Ricardo Arcila, Teresa K. Attwood, Niall Beard, Denise Carvalho-Silva, Alexandros C. Dimopoulos, Victoria Dominguez del Angel, Michel Dumontier, Kim T. Gurwitz, [ ... ], Patricia M. Palagi [ view all ]

Published: May 21, 2020 • https://doi.org/10.1371/journal.pcbi.1007854

# On the road to Define FAIR for Research Software

- Efforts to adapt and adopt the FAIR principles to research software (RDA FAIR4RS)

Recommendation n°2 :

Make sure **the specific nature of software** is recognized and not considered as "just data" particularly in the context of discussion about the notion of FAIR data.

**2019:** the **Opportunity Note** by the French national Committee for Open Science's Free Software and Open Source Project Group

(Clément-Fontaine, 2019)

Recommendation n°5 :

Recognise that FAIR guidelines will require **translation for other digital objects** and support such efforts.

**2020: 'Six Recommendations for Implementation of FAIR Practice'**

(FAIR Practice Task Force EOSC, 2020)

# FAIR4RS Principles

- **Findable**: Software, and its associated metadata, is easy for both humans and machines to find.

- **Accessible**: Software, and its metadata, is retrievable via standardized protocols.

- *Interoperable*: *Software interoperates with other software by exchanging data and/or metadata, and/or through interaction via application programming interfaces (APIs), described through standards.*

- *Reusable*: *Software is both usable (can be executed) and reusable (can be understood, modified, built upon, or incorporated into other software).*

(key differences from FAIR data principles in *italics*)

### Introducing the FAIR Principles for research software

Michelle Barker ✉, Neil P. Chue Hong, Daniel S. Katz, Anna-Lena Lamprecht, Carlos Martinez-Ortiz, Fotis Psomopoulos, Jennifer Harrow, Leyla Jael Castro, Morane Gruenpeter, Paula Andrea Martinez & Tom Honeyman

Output of the FAIR principles for research software (FAIR4S) - joint Research Software Alliance (**ReSA**), Research Data Alliance (**RDA**), **FORCE11** Working Group/Task force
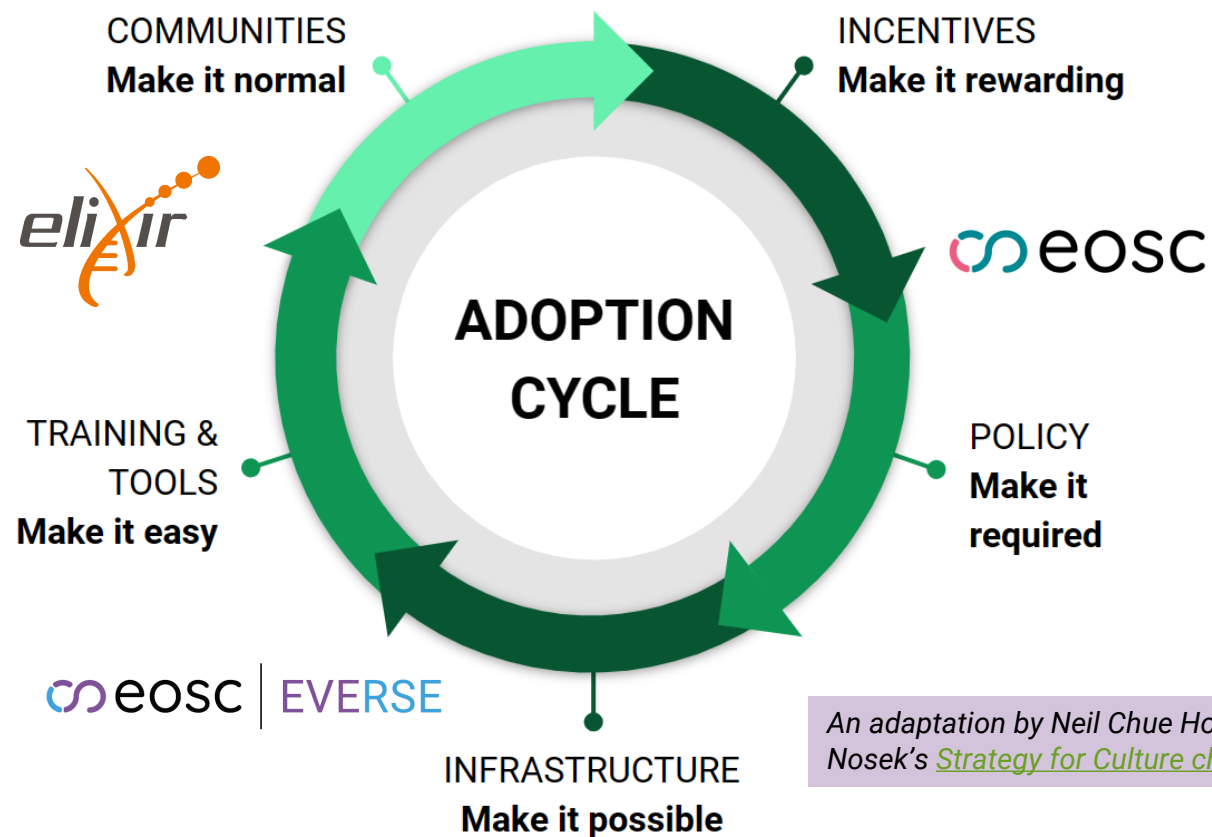
*Slide adapted from the presentation of the RDA FAIR4RS steering group at the International Funders Workshop (Nov 2022), https://zenodo.org/doi/10.5281/zenodo.7350198*

# Who is responsible for FAIR software?

Who is expected to apply FAIR?

➤ *The application of the FAIR4RS Principles is the responsibility of the **owners** (who are often the creators) of the software, not the users.*

➤ *The FAIR4RS Principles are also relevant to, and require support from, the **larger ecosystem** and various **stakeholders** that <u>support research software</u> (e.g., repositories and registries).*

*Slide adapted from the presentation of the RDA FAIR4RS steering group at the International Funders Workshop (Nov 2022), https://zenodo.org/doi/10.5281/zenodo.7350198*



COMMUNITIES
**Make it normal**

INCENTIVES
**Make it rewarding**

**ADOPTION CYCLE**

eosc

TRAINING & TOOLS
**Make it easy**

POLICY
**Make it required**

eosc | EVERSE

INFRASTRUCTURE
**Make it possible**

*An adaptation by Neil Chue Hong of Nosek's Strategy for Culture change*

# Managing (FAIR) Software



➢ helps to **implement best practices** during software development

➢ ensures that software is **accessible** and **reusable** in the short and longer term

➢ contributes to the **reproducibility** of results

➢ stimulates **collaborative** work on open-source software for research.

*Martinez-Ortiz, C. et al. (2022). Practical guide to Software Management Plans (1.0).* https://doi.org/10.5281/zenodo.7248877

Slides adapted from the "OrgMycology - eResearch NZ 2024" by Jonah Duckles (orgmycology)
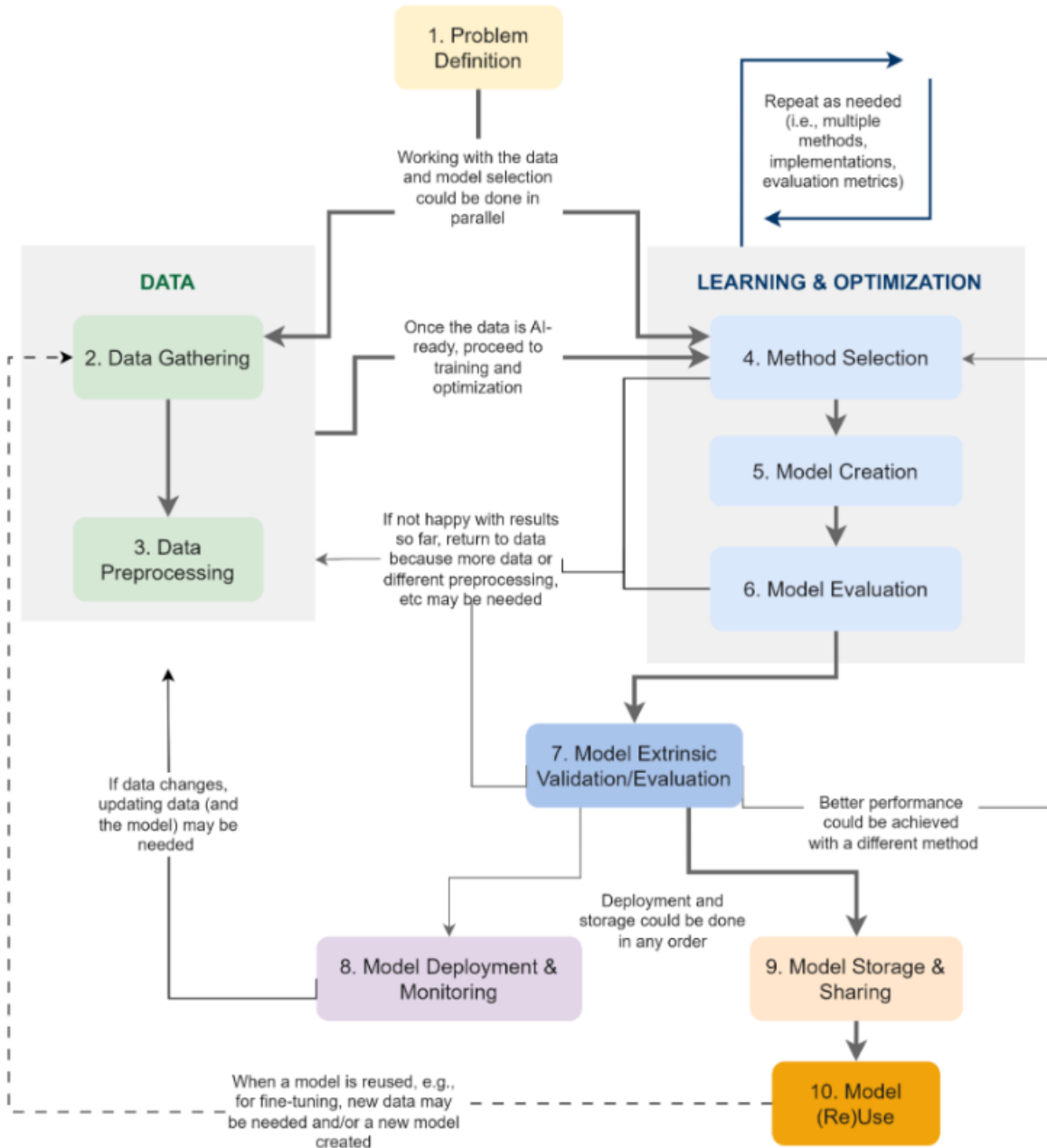
# How does FAIR fare (for Data / Software / ML)



Significant effort and push towards **FAIR data** (AI-ready datasets being a key demand)

**Software**: is only just beginning to get the support it needs as a first-class citizen in science

**ML/AI**: community just started realizing the challenges
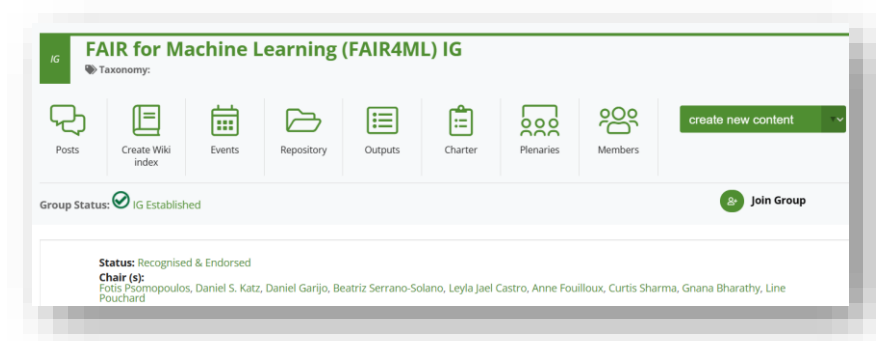
# FAIR in the ML Lifecycle

## Different FAIR principles apply to different aspects of the Life Cycle

| | | FAIR Principles | Best practices on reporting | Metadata schemas | Resources | What do you need to do here |
|---|---|---|---|---|---|---|
| 1 | Problem Definition | FAIR Data | | | | - Documentation |
| | | FAIR Software | | | | - Documentation |
| | | FAIR AI Models | | | | - Documentation |
| 2 | Data Gathering | FAIR Data (for training dataset) | DOME (D part) | | - Data Management, e.g., Data Stewardship Wizard (DSW)[7] [14]and Research Data Management Organizer (RDMO)[8] [15]<br>- Report data provenance and availability DOME registry[9] and BioImage Archive[10]<br>- SPDX licenses[11] | - Create a DMP<br>- Fill in information on the data in the DOME registry through the DOME Wizard |
| | | FAIR Software | | | | |
| | | FAIR AI Models | | | | - Fill in information on the data in the DOME registry through the DOME Wizard |
| 3 | Data Preprocessing | FAIR Data (for training dataset and splits) | DOME (D part) | ML Commons Croissant[12] | - Data Management, e.g., DSW and RDMO<br>- Report data splits DOME | - Create a DMP of the AI-ready data<br>- Report data features |

*Castro, L. J, et.al., & Zhang, Y. (2023). Lifecycle for FAIR Machine Learning..* https://doi.org/10.5281/zenodo.10407265

# FAIR in Machine Learning models

- ## What does FAIR apply to?
  - ### Are they data?
    - E.g., a set of parameters and options for a particular framework
  - ### Are they software?
    - E.g., an executable object that takes input and provides output
  - ### Are they something else?

- ## How does FAIR apply?
  - Searched and shared via repositories?
  - Searched and shared via executable platforms?
  - Searched and shared via something else? (e.g., DLHub, OpenML, HuggingFace…)
  - ## Models and training data are linked - should they be shared together?

*Slide adapted from various presentations of the RDA FAIR4ML interest group during Plenary events*

# New set of challenges

# Need for Community-led Standards and Best Practices (1/2)



ML Commons

mlcommons/
**croissant**

Croissant is a high-level format for machine learning datasets that brings together four rich layers.

ONNX 🤗 **Hugging Face**

Pistoia Alliance

Good Machine Learning Practices in the Modern Pharmaceutical Discovery Enterprise

_https://doi.org/10.31219/osf.io/kuz8p_

© 2023 World Scientific Publishing Company
https://doi.org/10.1142/9789811265679_0022

Chapter 22

A Roadmap for Defining Machine Learning Standards in Life Sciences

Fotis Psomopoulos*, Carole Goble†,**,
Leyla Jael Castro‡,††, Jennifer Harrow§,‡‡,
and Silvio C. E. Tosatto¶,§§

Artificial Intelligence for Science
A Deep Learning Revolution

editors
Alok Choudhary, Geoffrey Fox & Tony Hey

**DOME: recommendations for supervised machine learning validation in biology**

Ian Walsh, Dmytro Fishman, Dario Garcia-Gasulla, Tiina Titma, Gianluca Pollastri, ELIXIR Machine Learning Focus Group, Jennifer Harrow ✉, Fotis E. Psomopoulos ✉ & Silvio C. E. Tosatto ✉

_Nature Methods_ (2021) | Cite this article

4927 Accesses | 73 Altmetric | Metrics

elixir

**DOME Registry**
A database of annotations for published papers describing machine learning methods in biology.

(Giga)ⁿ Science

**DOME** adopted as part of the **submission system** for **GigaScience**
(_see example here:_ _http://gigadb.org/dataset/102404_)
**Online registry** of annotated papers: _https://registry.dome-ml.org_

# Need for Community-led Standards and Best Practices (2/2)

**FARR:** FAIR in ML, AI Readiness, & Reproducibility Research Coordination Network

## Ways to Get Involved

- **Input** on community needs, gaps & roadmap
- **Suggest use cases** and let us promote your project's use of AI and FARR-related practices
- Let us feature you in a **science story**

**Using FAIR to foster AI-readiness in Data Facilities: A resource list**

This work is supported through NSF award # 2226453.

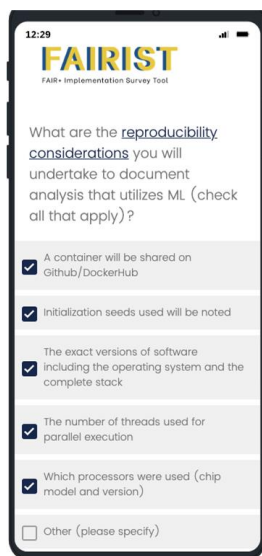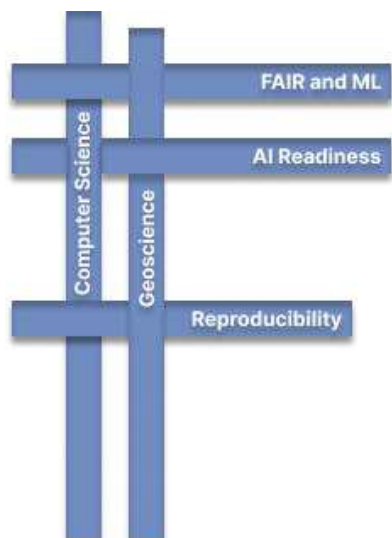**What is FAIR?**

- **A refresher on FAIR:** More than an acronym, it stands for 15 principles for making research objects more Findable, Accessible, Interoperable, Reusable https://www.go-fair.org/fair-principles/
- **Suggestions on how to implement FAIR:** https://bit.ly/implementFAIR

**Data repositories supporting AI with FAIR practices**

- **The geosciences:** https://www.hydroshare.org/
- **High energy physics:** https://bit.ly/AI-readyHEP
- **Materials science:** https://bit.ly/MLinMS

**Contact:**
https://www.farr-rcn.org/
community@farr-rcn.org

*This work is supported through the NSF award #2226453.*

**FAIRIST**
FAIR+ Implementation Survey Tool

What are the reproducibility considerations you will undertake to document analysis that utilizes ML (check all that apply)?

- ☑ A container will be shared on Github/DockerHub
- ☑ Initialization seeds used will be noted
- ☑ The exact versions of software including the operating system and the complete stack
- ☑ The number of threads used for parallel execution
- ☑ Which processors were used (chip model and version)
- ☐ Other (please specify)

## DEPARTMENT: LAST WORD

This article originally appeared in SECURITY & PRIVACY vol. 21, no. 6, 2023

## Trustworthy AI Means Public AI

Bruce Schneier, *Harvard University*

Today's generative AI systems are not trustworthy. We don't know how they are trained. We don't know their secret instructions. We don't know their biases, either accidental or deliberate. All we know is that they are created, at great expense, by corporations that will use every trick they can think of to make them as profitable as possible.

*If you want to go fast, go alone*
*If you want to go far, go together*

THANK YOU!
MERCI!
GRAZIE!
GRACIAS!
DANK JE WEL!

**CERTH**
CENTRE FOR
RESEARCH & TECHNOLOGY
HELLAS

**INAB**
INSTITUTE OF APPLIED BIOSCIENCES
ΙΝΣΤΙΤΟΥΤΟ ΕΦΑΡΜΟΣΜΕΝΩΝ ΒΙΟΕΠΙΣΤΗΜΩΝ
CENTRE for RESEARCH and TECHNOLOGY-HELLAS

elixir

*Slides available at: https://zenodo.org/records/10953108*

@fopsom@genomic.social

@fopsom