# 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
  - Assumes separate and sequential creator/publisher (repository) roles
- 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
- Work needed to define, then implement, then adopt principles

# Need for FAIR for non-data objects

- FAIR Principles, are intended to apply to all digital objects ([Wilkinson et al. 2016](#))

**FAIR Practice Task Force EOSC, "[Six Recommendations for Implementation of FAIR Practice](#)," 2020:**

*Recommendation 5*:

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

https://doi.org/10.5281/zenodo.TBD

# FAIR and ML Models

- As previously stated, original FAIR principles
  - Claim to apply to "scholarly digital research objects"
  - But actually focus on metadata and data
- FAIR for Research Software work and FAIR Workflows focusing on how to translate/interpret the principles for research software & workflows
- What about machine learning (ML) models?
  - 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 a combination of data+ software + workflows?
  - Are they something else?

https://doi.org/10.5281/zenodo.TBD

# How does FAIR apply?

- Large elements of FAIR for data are dependent on archival repositories (e.g. Zenodo, re3data.org)
  - Hold data and/or metadata, provide search and access capabilities
- Software is different, since it typically isn't shared via archival repositories but instead via social coding platform (e.g., GitHub) and package management systems (e.g. PyPI, CRAN)
- What about ML models?
  - Searched and shared via repositories?
  - Searched and shared via executable platforms?
  - Searched and shared via something else? (e.g., DLHub, OpenML, …)
- Models and training data are linked - should they be shared together?

# Work to-date and going forward

- [Poster](#) at RDA VP16 (Nov 2020):
- [BoF](#) at RDA VP17 (Apr 2021):
- [FAIR for Machine Learning Models](#) (Jun 2021), FAIR Festival
- 1st Community call (Jul 2021)
- [DaMaLOS talk](#) (24 Oct 2021)
- [BoF](#) at RDA VP18 (4 & 9 Nov 2021)
- [BoF](#) at SC21 (18 Nov 2021)
- [BoF](#) at RDA P19 (23 Jun 2022)

- Discussing a possible new interest group
- Discussing a potential white paper on FAIR 4 ML



**Defining FAIR for Machine Learning (ML)**

*Home*

https://doi.org/10.5281/zenodo.TBD

# Relevant projects & stakeholders

- Platforms
  - DLHub - Find, share, publish, and run machine learning models and discover training data for science
  - Kipoi - API & repository of ready-to-use trained models for genomics
  - OpenML - Build open source tools to discover (and share) open data, draw them into machine learning environments, build models, analyse results, get advice on better models
- Communities
  - Pistoia Alliance - a global, not-for-profit members' organization working to lower barriers to innovation in life science and healthcare R&D through pre-competitive collaboration
  - ELIXIR - An intergovernmental organisation that brings together life science resources (including databases, software tools, training materials, cloud storage and supercomputers) from across Europe
  - CLAIRE - Confederation of Laboratories for Artificial Intelligence Research in Europe
- Projects
  - FAIR4HEP - Using high-energy physics (HEP) as the science driver, developing a FAIR framework to advance understanding of AI, applying AI techniques, and exploring approaches to AI
  - HPC-FAIR - Providing a generic HPC data management framework to make both training data and AI models of scientific applications FAIR, focusing on the domain of program analyses/optimizations using AI/ML
- Others?  Please contact Dan

I ILLINOIS NCSA

https://doi.org/10.5281/zenodo.TBD