# FAIR and Open Data

Silvio Peroni

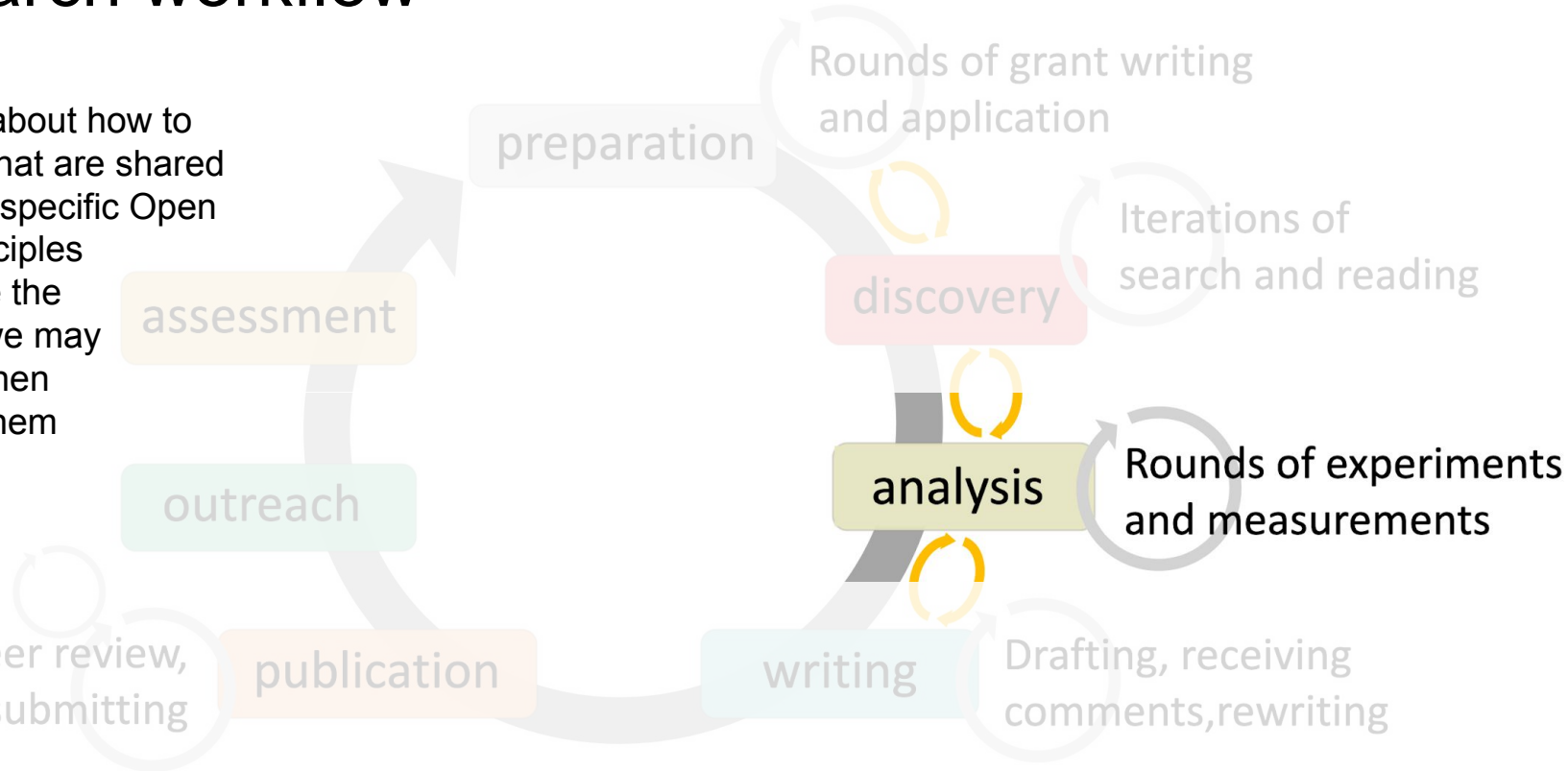silvio.peroni@unibo.it – https://orcid.org/0000-0003-0530-4305 – @essepuntato

Open Science (A.Y. 2022/2023)
Second Cycle Degree in Digital Humanities and Digital Knowledge
Alma Mater Studiorum - Università di Bologna

DIPARTIMENTO DI FILOLOGIA CLASSICA E ITALIANISTICA

# Research workflow

We discuss about how to create data that are shared according to specific Open Science principles and what are the Issues that we may encounter when processing them

Rounds of grant writing and application

preparation

Iterations of search and reading

discovery

assessment

analysis

Rounds of experiments and measurements

outreach

Submit, peer review, rejection, resubmitting

publication

writing

Drafting, receiving comments, rewriting

Kramer, B., & Bosman, J. (2015, June 18). The good, the efficient and the open—Changing research workflows and the need to move from Open Access to Open Science. CERN Workshop on Innovations in Scholarly Communication (OAI9), University of Geneva, Geneva, Switzerland. https://www.slideshare.net/bmkramer/the-good-the-efficient-and-the-open-oai9

# What are the data?

Metadata about the publication

Metadata about the data

Used in

A publication (https://doi.org/10.7717/peerj.4375) accepted in a journal

The related data on which the entire publication is based on

Piwowar, H., Priem, J., Larivière, V., Alperin, J. P., Matthias, L., Norlander, B., Farley, A., West, J., & Haustein, S. (2017). Data From: The State Of Oa: A Large-Scale Analysis Of The Prevalence And Impact Of Open Access Articles [Data set]. Zenodo. https://doi.org/10.5281/zenodo.837901

# FAIR data principles

Findability, Accessibility, Interoperability, and Reusability: these four principles, proposed by the FORCE11 community, serve to guide data producers and publishers for helping them to maximize the added-value gained by contemporary, formal scholarly digital publishing

They represent goals and desiderata of good data management and stewardship

Even if they have been devised originally for data, they have been proposed considering all scholarly digital research objects in mind, since all components of the research process must be available to ensure transparency, reproducibility, and reusability

Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G., Axton, M., Baak, A., Blomberg, N., Boiten, J.-W., da Silva Santos, L. B., Bourne, P. E., Bouwman, J., Brookes, A. J., Clark, T., Crosas, M., Dillo, I., Dumon, O., ... Mons, B. (2016). The FAIR Guiding Principles for scientific data management and stewardship. Scientific Data, 3, 160018. https://doi.org/10.1038/sdata.2016.18

# Findable

F1. (Meta)data are assigned a globally unique and persistent identifier

F2. Data are described with rich metadata (see also R1)

F3. Metadata clearly and explicitly include the identifier of the data they describe

F4. (Meta)data are registered or indexed in a searchable resource

# Accessible

A1. (Meta)data are retrievable by their identifier using a standardised communications protocol

    A1.1 The protocol is open, free, and universally implementable

    A1.2 The protocol allows for an authentication and authorisation procedure, where necessary

A2. Metadata are accessible, even when the data are no longer available

# Interoperable

I1. (Meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation

I2. (Meta)data use vocabularies that follow FAIR principles

I3. (Meta)data include qualified references to other (meta)data

# Reusable

R1. (Meta)data are richly described with a plurality of accurate and relevant attributes

R1.1. (Meta)data are released with a clear and accessible data usage license

R1.2. (Meta)data are associated with detailed provenance

R1.3. (Meta)data meet domain-relevant community standards

GO FAIR. (2018). FAIR Principles. https://www.go-fair.org/fair-principles/

# Different types of data

The agnosticism of generic principles like FAIR requires the various communities' specificities (i.e. proper to Informatics, Humanities, Medicine, Statistics, etc.) to be taken into account

General rule: within the FAIR framework, <span style="color:red">everything is data</span>

Data should be "intended to be as universal as possible, including datasets, publications, software, etc."

Avanço, K., Balula, A., Błaszczyńska, M., Buchner, A., … Wieneke, L. (2021). Future of Scholarly Communication—Forging an inclusive and innovative research infrastructure for scholarly communication in Social Sciences and Humanities (p. 46). Digital Humanities Centre at the Institute of Literary Research of the Polish Academy of Sciences. https://doi.org/10.5281/zenodo.5017705

# Some examples: what is data in Informatics?

A number of research communities and groups have been considering how to apply aspects of FAIR to research software since 2017

The adoption of FAIR principles enables the transparency, reproducibility, and reusability of research – and part of this story applies also to software that needs to be well-described (i.e. metadata), inspectable, documented and appropriately structured so that it can be executed, replicated, built-upon, combined, reinterpreted, reimplemented, and/or used in different settings

| **F: Software, and its associated metadata, is easy for both humans and machines to find.** |
| --- |
| F1. Software is assigned a globally unique and persistent identifier.<br>    ● F1.1. Components of the software representing levels of granularity are assigned distinct identifiers.<br>    ● F1.2. Different versions of the software are assigned distinct identifiers.<br>F2. Software is described with rich metadata.<br>F3. Metadata clearly and explicitly include the identifier of the software they describe.<br>F4. Metadata are FAIR, searchable and indexable. |
| **A: Software, and its metadata, is retrievable via standardized protocols.** |
| A1. Software is retrievable by its identifier using a standardized communications protocol.<br>    ● A1.1. The protocol is open, free, and universally implementable.<br>    ● A1.2. The protocol allows for an authentication and authorization procedure, where necessary.<br>A2. Metadata are accessible, even when the software is no longer available. |
| **I: Software interoperates with other software by exchanging data and/or metadata, and/or through interaction via application programming interfaces (APIs), described through standards.** |
| I1. Software reads, writes and exchanges data in a way that meets domain-relevant community standards.<br>I2. Software includes qualified references to other objects. |
| **R: Software is both usable (can be executed) and reusable (can be understood, modified, built upon, or incorporated into other software).** |
| R1. Software is described with a plurality of accurate and relevant attributes.<br>    ● R1.1. Software is given a clear and accessible license.<br>    ● R1.2. Software is associated with detailed provenance.<br>R2. Software includes qualified references to other software.<br>R3. Software meets domain-relevant community standards. |

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. (2022). FAIR Principles for Research Software (FAIR4RS Principles) [Recommendations with RDA Endorsement in Process]. Research Data Alliance. https://doi.org/10.15497/RDA00068

# Some examples: what is data in the Humanities?

Humanities (at least, those related to studies in philology, literary criticism, language, linguistics, history of art and archival and library studies) uses and/or produces several kinds of data

Probably, one of the most surprisingly data types is <span style="color:red">events</span>, meaning "any one-off gathering of people organised as a result of a research project, to share ideas, offer training, or present something to the public", including:

- Conferences
- Exhibitions
- Webinars
- Guided tours
- Teacher training

| No. | Data types | Produced by (tot. 19) | Used by (tot. 19) |
|-----|-----------|-----------------------|-------------------|
| i | Publications | 18 | 15 |
| ii | Other primary sources (e.g. manuscripts and artworks) | 0 | 18 |
| iii | Digital representation of cultural objects (e.g. facsimiles and photos) | 4 | 8 |
| iv | Catalogues, databases and other search tools | 2 | 9 |
| v | Events (e.g. conferences and exhibitions) | 6 | 0 |
| vi | Websites | 4 | 0 |
| vii | Software | 2 | 2 |
| viii | Documentation | 3 | 0 |
| ix | Digital infrastructures (e.g. mobile apps and web platforms) | 3 | 0 |
| x | Personal data | 2 | 0 |
| xi | Corpora | 2 | 0 |
| xii | Standards | 0 | 2 |
| xiii | Born-digital artefacts (e.g. tags, associations and texts) | 1 | 1 |

# Sharing: some guidelines

Many researchers do not share their data simply because they do not know how

Data should be published alongside detailed metadata

Data sharing will increase opportunities for collaboration



**sharing rights**
form an **agreement**
check your **library** for resources
follow authors' **guidelines**

**scooping**
you **know** your data
ideas are **plentiful**
open data = **more citations**

**lack of time**
sharing data **saves time**
create a **data management plan**

**transient storage**
avoid **proprietary** formats
share **as soon as possible**
use **stable repositories**

**lack of incentives**
open data = **more citations**
scientific **community recognition**

**sensitive content**
**aggregate** and **anonymize**
provide **sample data**
generate **synthetic datasets**

**insecurity**
share with **trusted colleagues**
recognize **no 'perfect code'**
emphasize **growth** and **learning**

**inappropriate use**
write detailed **metadata**
be willing to **help**
set data **governance plans**

**data too large**
**split data** into smaller chunks
share **properties of data**
**advocate** for storage funding

**unclear value**
value is **subjective**
perspectives are **limitless**
opportunities for **synthesis**

**unclear process**
check with your **library**
many **resources** exist
check **data templates**

**complex workflow**
write a detailed **readme**
use **graphics** to explain
**automate** where possible

reuse concerns · disincentives · Knowledge barriers

**data and code sharing**
*perceived barriers and solutions*

Gomes, D. G. E., Pottier, P., Crystal-Ornelas, R., Hudgins, E. J., Foroughirad, V., …, & Gaynor, K. M. (2022). Why don't we share data and code? Perceived barriers and benefits to public archiving practices. Proceedings of the Royal Society B: Biological Sciences, 289(1987), 20221113. https://doi.org/10.1098/rspb.2022.1113

# What about Open?

Warning: FAIR data does not imply Open Data

Open as defined in the Open Definition:

> "Open means anyone can freely access, use, modify, and share for any purpose (subject, at most, to requirements that preserve provenance and openness)"

| License | Attribution (i.e. preserving provenance) | Forcing openness |
|---|---|---|
| Creative Commons Attribution (CC-BY) | yes | no |
| Creative Commons Attribution - Share-Alike (CC-BY-SA) | yes | yes |
| Creative Commons CCZero (CC0) | no | no |
| Open Data Commons Public Domain Dedication and Licence (PDDL) | no | no |
| Open Data Commons Attribution License (ODC-BY) | yes | no |
| Open Data Commons Open Database License (ODbL) | yes | yes |

# Issues to the path towards the Open: personal data

GDPR is a regulation in EU law which concerns data protection and privacy, with particular regard to the gathering and processing of personal data of individuals

As stated in the Article 4 of the GDPR, personal data is any information that enable the identification of a natural person (i.e. the data subject), in particular by reference to an identifier such as a name, an identification number, etc.

Personal data can be gathered under specific conditions and cannot be published as Open Data – derogation to the this rule may exist, e.g. personal data in bibliographic information (authors' names, ORCIDs, etc.)

Note: the principles of data protection do not apply to anonymous information

# Anonymisation and pseudonymisation

It is possible to publish personal data if they are somehow changed in a way that do not allow the identification of a natural person anymore

Two ways:

1. Anonymisation – rendering personal data anonymous in such a manner that the data subject is not or no longer identifiable (there exists tools that try to semi-automate this activity, such as OpenAIRE Amnesia)
2. Pseudonymisation – rendering personal data in such a manner that they can no longer be attributed to a specific data subject without the use of additional information, provided that such additional information is kept separately and is subject to technical and organisational measures to ensure that the personal data are not attributed to an identified or identifiable natural person

Applying one of these approaches enable one to publish Open Data also when the original gathering of such data included personal data

# How to publish data

Generally speaking, data should be published <span style="color:red">as open as possible, as closed as necessary</span>

In choosing the license, in particular if considering or not the attribution clause, one has to think about the final goal to reach – that, within the Open Science ecosystem, should be fostering maximum reuse

Recently, the European Union has provided a [clear guideline](clear guideline) to follow for what it concerns the publication of data:

> "raw data, metadata or other documents of comparable nature may alternatively be distributed under the provisions of the Creative Commons Universal Public Domain Dedication deed (CC0 1.0)"

Landi, A., Thompson, M., Giannuzzi, V., Bonifazi, F., Labastida, I., da Silva Santos, L. O. B., & Roos, M. (2020). The "A" of FAIR – As Open as Possible, as Closed as Necessary. Data Intelligence, 2(1–2), 47–55. https://doi.org/10.1162/dint_a_00027

# The (senseless) scholars' fear

Misconception: *I want to share my data, but it is important to me that I'm cited when people use it. I prefer CC BY to CC0 because this kind of attribution is what I care about most*

If the issue is to be cited and acknowledged for a work, problem solved: citation practices <span style="color:red">are built around ethical norms</span>, not around legal requirements

Using and not citing and acknowledging an existing work relates to <span style="color:red">plagiarism</span>, and this act has a legal consequence independently from the license used to protect the data

Thus, the general rule for sharing data compliantly with Open Science practices is to use CC0 license instead of CC BY to foster maximum reuse

Wolfe, M. (2017, August 9). CC0 and Data Citation. https://www.library.ucdavis.edu/news/cc0-and-data-citation/

# Trustworthy data

Provenance is a record that describes the people, institutions, entities, and activities involved in producing, influencing, or delivering data, and it is crucial in deciding whether information is to be trusted, how it should be integrated with other diverse information sources, and how to give credit to its originators when reusing it

Provenance can help to make trust judgements about a piece of data, since it , which can be used to form assessments about its quality, reliability or trustworthiness

Minimal provenance information:

- When a piece of data has been created and/or invalidated
- Who created a piece of data
- What is the source that has been used to create a piece of data

Belhajjame, K., B'Far, R., Cheney, J., Coppens, S., Cresswell, S., Gil, Y., Groth, P., Klyne, G., Lebo, T., McCusker, J., Miles, S., Myers, J., Sahoo, S., & Tilmes, C. (2013). PROV-DM: The PROV Data Model (L. Moreau & P. Missier, Eds.). World Wide Web Consortium. https://www.w3.org/TR/prov-dm/

# Trustworthy long-term preservation of data

Guaranteeing that the data one publish are reusable by humans and machines, i.e. FAIR, is not enough, since a good part of the story concerns also to think ways to trustworthily preserve such data in the long term

You should take into consideration such digital repositories which are compliant with the TRUST principles:

- Transparency – repository services and data holdings are publicly verifiable
- Responsibility – authenticity+integrity of data and reliability+persistence of services
- User Focus – meet norms and expectations of target user communities
- Sustainability – preserve data and services for the long-term
- Technology – support secure, persistent, and reliable services

Lin, D., Crabtree, J., Dillo, I., Downs, R. R., Edmunds, R., Giaretta, D., De Giusti, M., L'Hours, H., Hugo, W., Jenkyns, R., Khodiyar, V., Martone, M. E., Mokrane, M., Navale, V., Petters, J., Sierman, B., Sokolova, D. V., Stockhause, M., & Westbrook, J. (2020). The TRUST Principles for digital repositories. Scientific Data, 7(1), 144. https://doi.org/10.1038/s41597-020-0486-7

# Data Management Plan

A data management plan (DMP) is the tool one has to use to ensure that all the aspects introduced in the previous slides will be addressed appropriately – i.e. before one starts to gather the data

A DMP is a document that describes how you will treat your data during a project and what happens with the data after the project ends

What it is covered in a DMP:

- data discovery, collection and organization
- quality assurance/quality control
- documentation
- data preservation and sharing

Michener, W. K. (2015). Ten Simple Rules for Creating a Good Data Management Plan. PLOS Computational Biology, 11(10), e1004525. https://doi.org/10.1371/journal.pcbi.1004525

# End

## FAIR and Open Data

Silvio Peroni

silvio.peroni@unibo.it – https://orcid.org/0000-0003-0530-4305 – @essepuntato

Open Science (A.Y. 2022/2023)
Second Cycle Degree in Digital Humanities and Digital Knowledge
Alma Mater Studiorum - Università di Bologna

DIPARTIMENTO DI FILOLOGIA CLASSICA E ITALIANISTICA