

# OPEN SCIENCE DALLA A ALLA Z 6-DATI FAIR E DATA MANAGEMENT PLAN

Elena Giglia

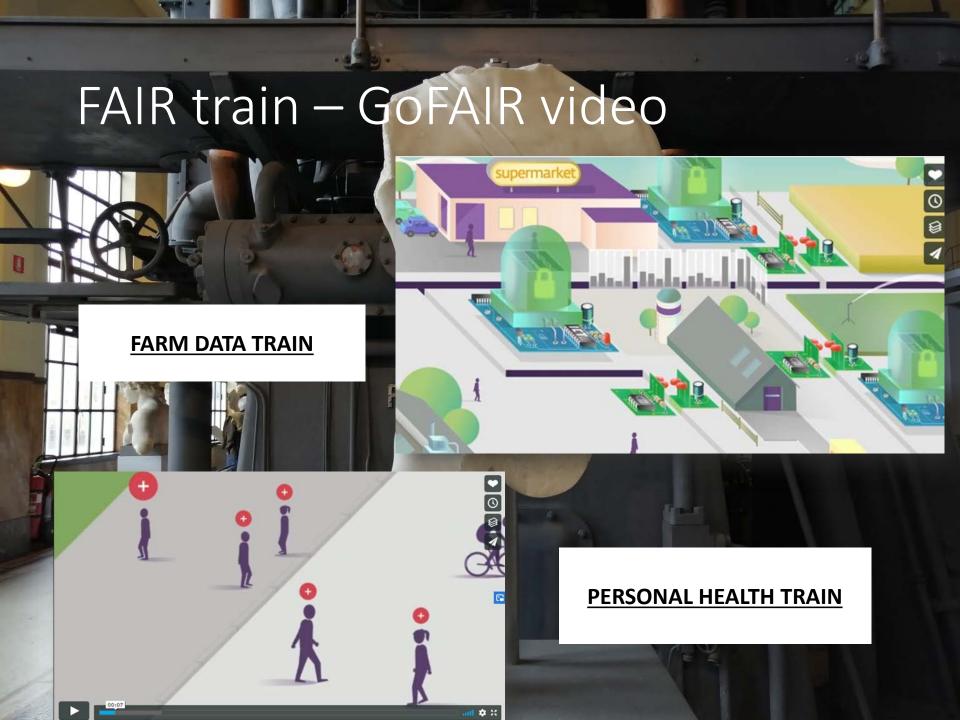


# Impareremo in questo modulo

- 1. Cosa significa FAIR nella ricerca quotidiana
- 2. come scrivere un DMP (Data Management Plan)

### MESSAGGI CHIAVE

- FAIR è il futuro (se no si resta tagliati fuori
- È più facile di quanto sembri, basta provarci



# ...FAIR SIGNIFICA [anche e soprattutto per le macchine]

### **FINDABLE**

- IDENTIFICATIVI
  - METADATI

### **INTEROPERABLE**

- STANDARDS
- ONTOLOGIE

### **A**CCESSIBLE

- DOVE SONO CONSERVATI E A QUALI CONDIZIONI DI ACCESSO
  - NON SIGNIFICA «OPEN»
    - FORMATI APERTI

### **R**EUSABLE

- LICENZE D'USO
- DOCUMENTAZIONE
- LEGGIBILI DALLE MACCHINE



Nov. 20, 2018





### **TURNING FAIR INTO**



### Define

### Implement

### Embed and sustain

Concepts for FAIR implementation

Rec. 1: Define FAIR for implementation

Rec. 2: Implement a Model for FAIR Digital Objects

Rec. 3: Develop components of a FAIR ecosystem

broadly

Rec. 17: Align and harmonise FAIR and Open data policy

FAIR culture

Rec. 4: Develop Interoperability frameworks

Rec. 5: Ensure data management via DMPs

Rec. 6: Recognise & reward FAIR data & stewardship

FAIR ecosystem

Rec. 7: Support semantic technologies

Rec. 8: Facilitate automated processing

> Rec. 9: Certify FAIR services

Skills for FAIR

Rec. 10: Professionalise data science & stewardship roles

Rec. 11: Implement curriculum frameworks and training

Above line = priority recommendations

Below line = supporting recommendations

Rec. 25: Implement and monitor metrics

Rec. 26: Support data citation and next generation metrics

Incentives and metrics for FAIR data and services

Rec. 12: Develop metrics for FAIR Digital Objects

Rec. 13: Develop metrics to certify FAIR services

Investment in FAIR

Rec. 14: Provide strategic and coordinated funding

Rec. 15: Provide sustainable funding

Rec. 27: Open EOSC to all providers but ensure services are FAIR

Rec. 16: Apply FAIR

Rec. 18: Cost data management

Rec. 19: Select and prioritise FAIR digital objects

Rec. 20: Deposit in Trusted Digital Repositories

Rec. 21: Incentivise reuse of FAIR outputs

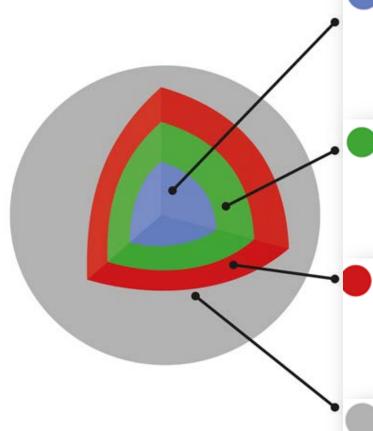
Rec. 22: Use information held in DMPs

Rec. 23: Develop components to meet research needs

Rec. 24: Incentivise research infrastructures to support FAIR data



# edeal FAIR object



### DIGITAL OBJECT

### Data, code and other research outputs

At its most basic level, data or code is a bitstream or binary sequence. For this to have meaning and to be FAIR, it needs to be represented in standard formats and be accompanied by Persistent Identifiers (PIDs), metadata and documentation. These layers of meaning enrich the object and enable reuse.

### **IDENTIFIERS**

### Persistent and unique (PIDs)

Digital Objects should be assigned a unique and persistent identifier such as a DOI or URN. This enables stable links to the object and support citation and reuse to be tracked. Identifiers should also be applied to other related concepts such as the data authors (ORCIDs), projects (RAIDs), funders and associated research resources (RRIDs).

### STANDARDS & CODE

### Open, documented formats

Digital Objects should be represented in common and ideally open file formats. This enables others to reuse them as the format is in widespread use and software is available to read the files. Open and well-documented formats are easier to preserve. Data also need to be accompanied by the code use to process and analyse the data.

### METADATA

### Contextual documentation

In order for Digital Objects to be assessable and reusable, they should be accompanied by sufficient metadata and documentation.

Basic metadata will enable data discovery, but much richer information and provenance is required to understand how, why, when and by whom the objects were created. To enable the broadest reuse, they should be accompanied by a plurality of relevant attributes and a clear and accessible usage license.

D, 2015



Technology

### **FAIR Principles**

Sci. Data 3:160018 doi: 10.1038/sdata.2016.18 (2016)



### Findable:

F1 (meta)data are assigned a globally unique and persistent identifier;

F2 data are described with rich metadata;

F3 metadata clearly and explicitly include the identifier of the data it describes;

F4 (meta)data are registered or indexed in a searchable resource;

### Interoperable:

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

12 (meta)data use vocabularies that follow FAIR principles;

13 (meta)data include qualified references to other (meta)data;

### Accessible:

A1 (meta)data are retrievable by their identifier using a standardized communications protocol;

A1.1 the protocol is open, free, and universally implementable;

A1.2 the protocol allows for an authentication and authorization procedure, where necessary;

A2 metadata are accessible, even when the data are no longer available;

### 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;

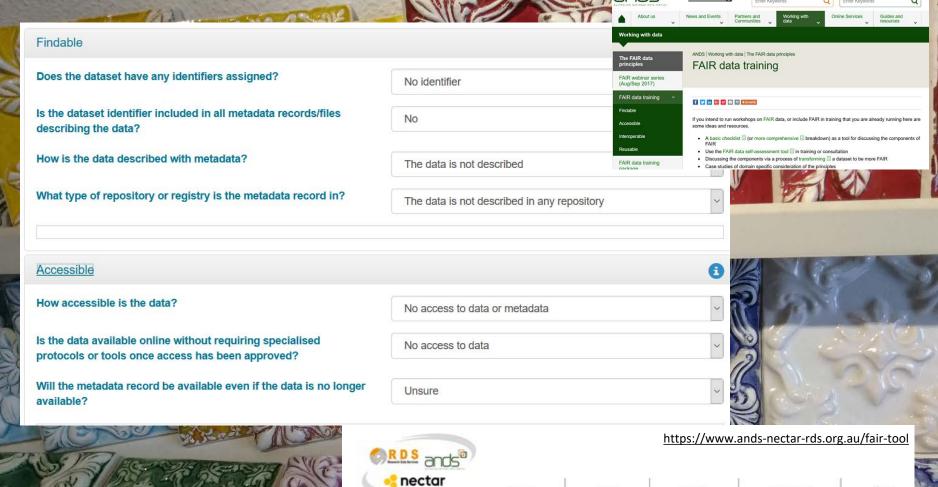
R1.2 (meta)data are associated with detailed provenance;

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

Normal CC BY Erik Schultes

https://www.go-fair.org/wp-content/uploads/2018/11/26102018 Country meeting GFISCO staff presentation.pdf

# ... sfumature di FAIR



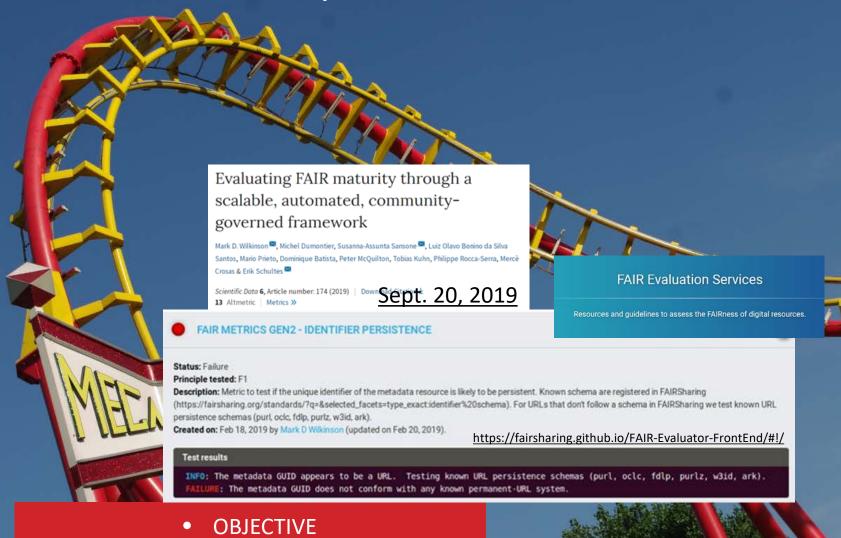
### FAIR self-assessment tool

Welcome to the ARDC FAIR Data self-assessment tool. Using this tool you will be able to assess the 'FAIRness' of a dataset and determine how to enhance its FAIRness (where applicable).

**Training** 

# FAIR maturity evaluator

MACHINE-READABLE ... AS FAIR DATA ARE



# Top 10 FA

Research Software, p. 15

Research Libraries, p. 20

Research Data Management Support, p. 25

International Relations, p. 30

Humanities: Historical Research, p. 34

Geoscience, p. 42

Biomedical Data Producers, Steward

Biodiversity, p. 59

Australian Government Data/Collec





### Top 10 FAIR Data & Software Things

February 1, 2019

Sprinters:

Reid Otsuii, Stephanie Labou, Rvan Johnson, Guilherme Castelao, Bia Villas Boas, Anna

on a man see an emmpre marte a specific rocabulary (the ro



### Sprinters:

Geoscience

ohn Brown, Janice Chan, Niamh Quigley (Curtin Univer

Audience:

Researchers

#### Things

#### Findable

Thing 1: Data sharing and discovery

Thing 6: Vocabularies for data description

Thing 7: Identifiers and linked data

Thing 10: Spatial data

#### Accessible

Thing 2: Long-lived data: curation & preservation

Thing 3: Data citation for access & attribution

markup a dataset can be found here. The dataset is part of a project to reconstruct the domestic market for colonial goods in the Dutch Republic.

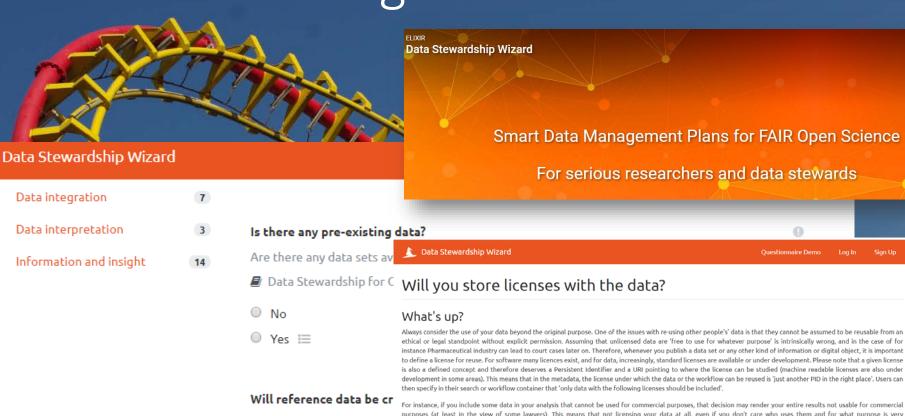
#### ACTIVITIES:

1. Try to find one or two terms that are relevant to your research using the resources that are mentioned above. You can aso use Swoogle to search for vocabularies related to your research. 2. Search for a term related to your research in the CIDOC Conceptual Reference Model (CRM) concept search. Were you able to find it? Tip 1: Search for "person" to get an idea of how the thesaurus works. Tip 2: All the terms used can be found in the last release of the model: http://www.cidoc-crm.org/get-last-official-release.

### Thing 7: FAIR data modelling

The fourth and the fifth star in Berner Lee's model can be awarded when the data are stored in a format in which the topics their properties and their characteristics are identified using URIs whenever possible. More concretely, it implies that you record your data using the Resource Description Framework (RDF) format. RDF, simply put, is a technology which enables you to publish the contents of a database via the web. It is based on a simple data model which assumes that all statements about resources can be reduced to a basic form.

# FAIR Data management wizard



Will any of the data that y others)?

- Data Stewardship for C
- No

ethical or legal standpoint without explicit permission. Assuming that unlicensed data are 'free to use for whatever purpose' is intrinsically wrong, and in the case of for instance Pharmaceutical industry can lead to court cases later on. Therefore, whenever you publish a data set or any other kind of information or digital object, it is important to define a license for reuse. For software many licences exist, and for data, increasingly, standard licenses are available or under development. Please note that a given license is also a defined concept and therefore deserves a Persistent Identifier and a URI pointing to where the license can be studied (machine readable licenses are also under development in some areas). This means that in the metadata, the license under which the data or the workflow can be reused is 'just another PID in the right place'. Users can

purposes (at least in the view of some lawyers). This means that not licensing your data at all, even if you don't care who uses them and for what purpose is very counterproductive and will severely undermine the actual reuse of your data by others and in particular by industry. It will also lower the attribution-rate (usually part of the

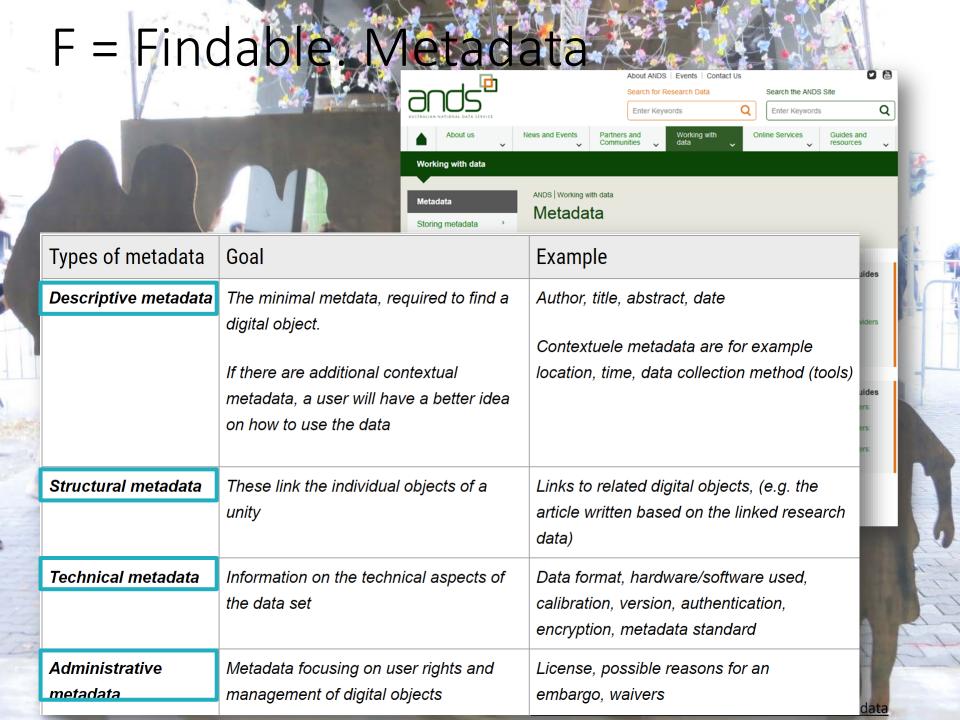
license conditions) and thus the citation and the impact score of your data.

#### Do

- · Always carefully choose a license to be attached to your data upon publication.
- · Include and clearly mark the licences PID as a concept + attributes in the metadata.
- · Store and 'expose' the license as part of the metadata in Open Access environments where search engines can easily find the license, even of the data they describe are not (yet) FAIR or even highly restricted in access. The 'fact' that a data set with a specific license is 'out there' is a first step toward effective reuse of your data or information source.
- · Make sure, especially when you restrict use of your data, that you are able to enforce the license you choose. Licenses that are not enforceable make no sense. (please note that the enforcement is usually not done by an individual research group but at institutional or repository level)

#### Don't

- · Ever publish data without a license attached or choose a license lightly, without considerations of anticipated reuse of your data.
- Choose a license that is not transitive (i.e can not be transferred with subsets of the data), but make sure its transitivity does not unduly restrict the reuse of your data.
- · Choose an unnecessary complicated license with many clauses and wherever possible one that is already widely adopted in the research community for either software Will you be storing samples:



# F = Findable. Metadata standards

### Metadata

RDA | Metadata Directory

Edit this page

View the standards

View the extensions

View the tools

View the use cases

Browse by subject areas

Contribute

Add standards

Add extensions

Add tools

Add use cases

### Arts and Humanities © Edit

- Creative art and design & Edit

- Law & Edit
- Music & Edit

### Engineering © Edit

- Architecture & Edit

### Life Sciences © Edit

### Physical Sciences & Mathematics © Ent

- Astronomy & Edit
- · Astrophysics & Edit
- Chemistry & Edit
- · Climatology & Edit
- Crystallography & Edit
- Environmental Science & Edit
- · Geology & Edit
- Geoscience © Edit
- Glaciology & Edit
- Hydrogeology Edit
- Hydrography & Edit
- · Hydrology & Edit
- Marine Science © Edit
- . Maritime Geography & Edit
- Meteorology & Edit
- . Minerology & Edit
- Nuclear and Particle Physics Edit
- Oceanography & Edit
- · Palaeontology & Edit
- . Physics (K Edit
- Planetary science Edit
- · Remote Sensing & Edit
- Soil Science & Edit

### Social and Behavioral Sciences © Edit

- Demography & Edit
- Economics © Edit
- Health Policy & Edit
- Planning (Urban, Rural and Regional) C Edit
- Politics & Edit
- Sociology & Edit

### General Research Data © Edit

# F = findable: Metadata tools

CEDAR CENTER FOR EXPANDED DATA ANNOTATION AND RETRIEVAL

PURPOSE

RESEARCH

Better data for better science

TOOLS | TRAINING

What CEDAR does

https://metadatacenter.org/

The CEDAR Workbench, as we refer to the suite of CEDAR tools, makes it easy to collect and use metadata. Eventually our tools will metadata record is created to its eventual processing, and even enhancement, by users and analysts. But for now, CEDAR tools held to users, and download the information that users have provided.

What can CEDAR do for me already?

As of its production release, in February 2017, CEDAR addresses these scenarios:

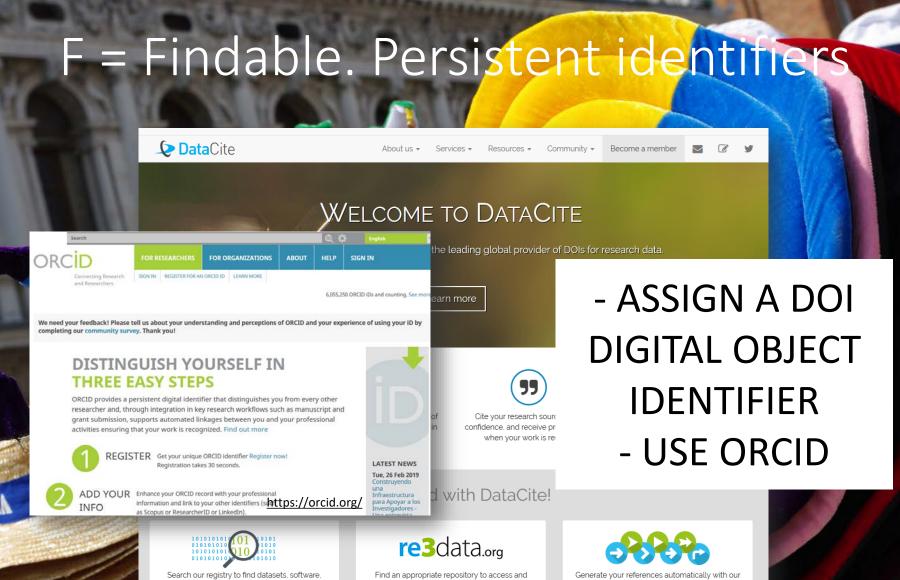
- · create user-friendly, shareable forms for collecting metadata, with features like
  - o nested and repeatable elements and fields
  - o reusable elements
  - o control over tool tips, field titles, and field descriptions
- · share your forms and metadata
  - o provide a link to your metadata editors, so they can enter metadata responses based on your forms
  - o share your forms and other content with individuals or a group
  - create and manage groups to make permissions simpler
- · associate your questions (fields) and possible answers (values) with controlled terms
  - select any term or collection of terms from the NCBO BioPortal semantic repository
  - combine different terms from different controlled vocabularies into a single set of options
  - o create your own terms, or term lists ('value sets') that can be re-used
- view responses meeting your (simple) search critieria, in several forms
  - o CEDAR Metadata Editor's metadata view
  - o an in-line JSON-LD format, used by CEDAR for all its metadata instances
  - download of JSON-LD files via the CEDAR REST API, for offline integration with your workflow
- · use the Workbench Desktop interface to manage your content
  - use My Workspace to see your items, or Shared with Me to see other items you can access
  - o select an item and control-click or use the 3-dot menu in the upper right to share it, copy it, delete it, or get info on it
- enable intelligent metadata suggestions in your template by using a field's Suggestions tab
  - CEDAR keeps track of metadata entered for that field
  - o users will see a drop down list of the most popular metadata entries, and can select from them
- remotely access CEDAR content and capabilities using the CEDAR REST API

B cell repertoire in myasthenia gravis but whether their generation is associated with broader defects in the B cell repertoire is unknown. To address this question, we performed deep sequencing of the B cell receptor repertoire of AChR-MG, MuSK-MG and healthy subjects. Project ID PRJNA338795 -Scope Multispecies Experiment Type 0 Select option... Genome sequencing and assembly Raw sequence reads Genome sequencing Sample Assembly Clone ends Epigenomics CANCEL

Let's pick a scope and an experiment type.



With these capabilities, you can capture simple or rich metadata for your project, build a repository of project metadata, or design particular needs. Advanced users can even submit metadata entries through CEDAR's REST API.



images, and other research material

deposit research data with re3data.org





### F = Findable. Persistent identifiers

#### Here are some identifier schemes:

- ARK (Archival Resource Key) a URL with extra features allowing you to ask for descriptive and archival metadata and to recognize certain kinds of relationships between identifiers. ARKs are used by memory organizations such as libraries, archives, and museums. They are resolved at "http://www.nt2.net". Resolution depends on HTTP redirection and can be managed through an API or a user interface.
- DOI (Digital Object Identifier) an identifier that becomes actionable when embedded in a URL. DOIs are very popular in academic journal publishing. They are resolved at "http://dx.doi.org". Resolution depends on HTTP redirection and the Handle identifier protocol, and can be managed through an API or a user interface.
- Handle an identifier that becomes actionable when embedded in a URL. Handles are resolved at "http://www.handle.net/". Resolution depends on HTTP redirection and the Handle protocol, and can be managed through an API or a user interface.
- InChI (IUPAC International Chemical Identifier) a non-actionable identifier for chemical substances that can be used in printed and electronic data sources, thus enabling easier linking of diverse data compilations.
- LSID (Life Sciences Identifier) a kind of URN that identifies a biologically significant resources, including species names, concepts, occurrences, and genes or proteins, or data objects that encode information about them. Like other URNs, it becomes actionable when embedded in a URL.
- NCBI (National Center for Biotechnology Information) ACCESSION
- a non-actionable number in use by NCBI.
- PURL (Persistent Uniform Resource Locator) a URL that is always redirected through a hostname (often purl.org). Resolution depends on HTTP redirection and can be managed through an API or a user interface.
- URL (Uniform Resource Locator) the typical "address" of web content. It is a kind of URI (Uniform Resource Identifier) that begins with "http://" and consists of a string of characters used to identify or name a resource on the Internet. Such identification enables interaction with representations of the resource over a network, typically the World Wide Web, using the HTTP protocol. Well-managed URL redirection can make URLs as persistent as any identifier. Resolution depends on HTTP redirection and can be managed through an API or a user interface.
- URN (Uniform Resource Name) an identifier that becomes actionable when embedded in a URL. Resolution depends on HTTP redirection and the DDDS protocol, and can be managed through an API or a user interface. A browser plug-in can save you from typing a hostname in front of it.

  https://dmptool.org/general\_guidance#persistent-identifiers

# ... [SE NON POTETE USARE UN DOI]

**PUBLICATIONS** -



0.00514

COMMUNITY -

FORCE11 » Community News » Introducing a new standard for the citation of research data

RESOURCES -

INTRODUCING A NEW STANDARD FOR THE CITATION OF RESEARCH DATA

Posted by Jennifer McLennan | May 8, 2018 | Sign In or Join Now! to post comments

#### Rules, registry and recommendations

Compact Identifiers. A "compact identifier" is a string constructed by concatenating a namespace prefix, a separating colon, and a locally unique identifier (LUI), e.g. pdb:2gc4.

NEWS + BLOGS →

**Provider Specification.** To specify a specific provider, where multiple providers exist, prepend the provider code and a "/" to the compact identifier, e.g. rcsb/pdb:2gc4.

**Provider Default.** Where multiple providers exist, and the provider is not specified in the compact identifier, the resolver will determine where to resolve the request based on its own rules, e.g., taking into account uptime availability, regional preference, or other criteria.

Redirect Rule. A URL template associated with the provider code is maintained in the namespace registry, defining how to forward compact identifiers to any specific provider (see 4.2.3 below).

SCIENTIFIC DATA

\_\_\_\_\_\_Altmetric: 20

May 8, 2018

More detail

Article | OPE

Uniform resolution of compact identifiers for biomedical data

Sarala M. Wimalaratne, Nick Juty, John Kunze, Greg Janée, Julie A. McMurry, Niall Beard, Rafael Jimenez, Jeffrey S. Grethe, Henning Hermjakob, Maryann E. Martone & Tim Clark

A reasonable solution to the identifiers problem is to assign Digital Object Identifiers (DOIs) to identify datasets. DOIs are already widely used in the publishing world as persistent identifiers for scholarly publications. They have been adopted by generalist data repositories such as Dryad, FigShare, Zenodo and Dataverse, as well as by domain data repositories outside of biomedicine. Handles<sup>14</sup>, which underlie the DOI system, may also be used directly. The DataCite consortium provides a robust central means for assigning DOIs to data.

However, DOIs are not commonly used for biomedical data, which is partitioned across over 600 autonomous repositories that are independently funded. Instead, in biomedicine there has been a tember 2017 uary 2018

### A = Accessible

### ACCESSIBLE≠OPEN «ACCESSO» PUÒ ANCHE ESSERE RISERVATO O SOTTO EMBARGO

### Open access

Data that can be accessed by any user whether they are registered or not. Data in this category shouldn't contain personal information (unless consent is given (see 'Informed consent').

### · Access for registered users (safeguarded)

Data that is accessible only to users who have registered with the archive. This data contains no direct identifiers but there may be a risk of disclosure through the linking of indirect identifiers.

### Restricted access

Access is limited and can only be granted upon request. This access category is for the most sensitive data that may contain disclosive information.

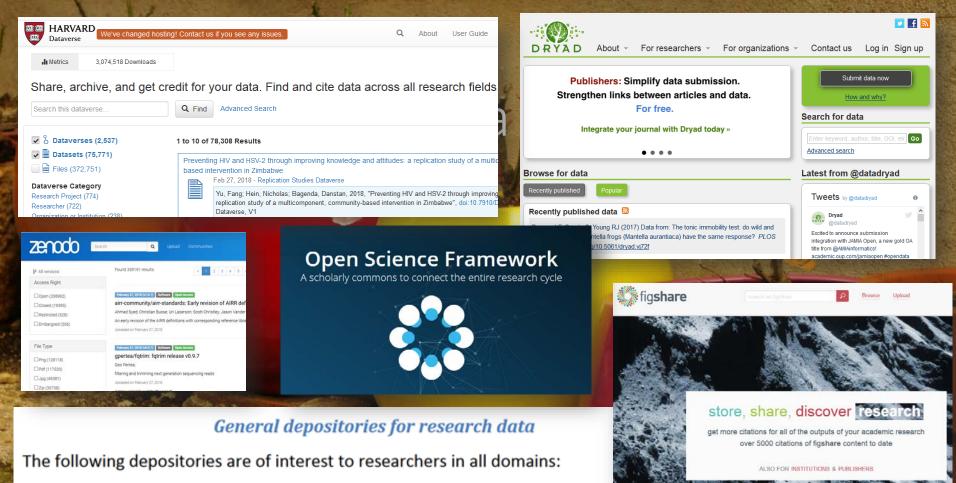
Restricted access requires long-term commitment of the researcher or person responsible for the data to handle the upcoming the permission requests.

### Embargo

Besides offering the opportunity for restricted access 'for eternity' most data repositories allow you to place a temporary embargo on your data. During the embargo period, only the description of the dataset is published. The data themselves will become available in open access after a certain period of time.

CESSDA Guide





- Zenodo (not-for-profit, hosted by CERN): <a href="https://zenodo.org">https://zenodo.org</a>:
- Dryad (not-for-profit membership organisation): <a href="http://www.datadryad.org">http://www.datadryad.org</a>
- Figshare (free service provided by private company): <a href="https://figshare.com">https://figshare.com</a>
- Open Science Framework (not-for-profit, developed and maintained by the Center for Open Science<sup>1</sup>): <a href="https://osf.io">https://osf.io</a>
- Harvard Dataverse (not-for-profit, hosted by the Institute for Quantitative Social Studies IQSS at Harvard University): <a href="https://dataverse.harvard.edu">https://dataverse.harvard.edu</a>

# A = Accessible. Data repositories



Value of data increases up the tiers:

from individual to community to

Each higher tier brings greater responsibility and demands for

sustainability and provenance.

social value.

access.

#### 3. Will the data be safe in legal terms?

For this criterion we first consider the basic legal terms and conditions to chec Here only one capability level is given, as a repository either will or will not comply. Then we consider licensing, disclosure risk and access control, when repository may offer different levels of capability you can match to your needs Note that Re3data provides relevant details for the repositories it lists, under

Legal terms and conditions

#### Level 1

Personal data or data which may identify individuals when linked to other data should not be stored outside the European Economic Area, unless in a legal jurisdiction that ensures personal data is adequately

By agreeing to the terms and conditions the depositor will not be breaching other Data Protection principles, or the terms of any confidentiality agreement with data subjects or owners (e.g. consent form, consortium agreement)

The checklist that follows addresses the five key questions posed in this quide:

- is the repository reputable?
- it take the data you want to deposit?
- will it be safe in legal terms?
- will the repository sustain the data value?
- will it support analysis and track data usage?



Tier 3 institutional repository

Individual collections

Tier 2 national data centre "Where to keep research datasearch Data > Where Keep Research Data Where to keep research data Version 1.1 of the DCC checklist for evaluating data repositories By Angus Whyte, Published: 28 December 2015. Updated 22 January 2016.

Findable, accessible and interoperable

Level 2

Repository publishes

disciplinary discovery

Supports assignment

of related persistent

Provides metadata

elements to enable

broader discovery

(e.g. geo-spatial) to

Home Digital curation About us News Events Resources Training Projects Cor

IDs per dataset/

collection

other pertinent

information as

enhance cross-

metadata fields to

Level 3

Metadata is

catalogued to

enhance reuse

to fulfil domain-

according to sector-

specific purposes

Supports assignment

of multiple persistent

IDs at different levels

of granularity within

dataset/ collection

Exposes discovery

metadata as Linked

optimise automatic

Open Data to

leading standards, or

Level 1

Metadata publishing:

Data collections are

repository according to

funder expectations so

discoverable by title.

creator, and date of deposition

Stable identifiers:

Enables a DOI or other

open standard identifier

to be assigned to a

ingested dataset/

metadata:Provides

Datacite mandatory

In this section

How-to Guides & Checklists Appraise & Select Research

Gite Datasets and Link to

metadata and exposes

collection

Discovery

landing page for each

catalogued in a

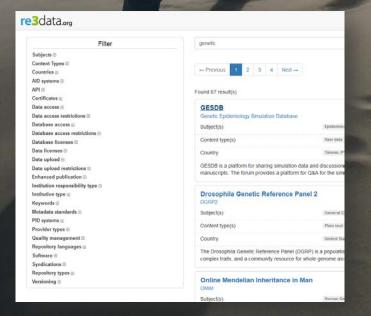
that they are

As infrastructure increases so must the attention given to standards,

Tier 4

# A = Accessible. Cercate un archivio?







### 2,000 Data Repositories and Science Europe's Framework for Discipline-specific Research Data Management

By offering detailed information on more than 2,000 research data repositories, re3data has become the most comprehensive source of reference for research data infrastructures globally. Through the development and advocacy of a framework for discipline...

Read more

### Three new DOI Fabrica features to simplify account management

Last month month we launched DOI Fabrica, the modernized version of the DataCite Metadata Store (MDS) web frontend. It is the one place for DataCite providers and their clients to create, find, connect and track every single DOI from their organization...

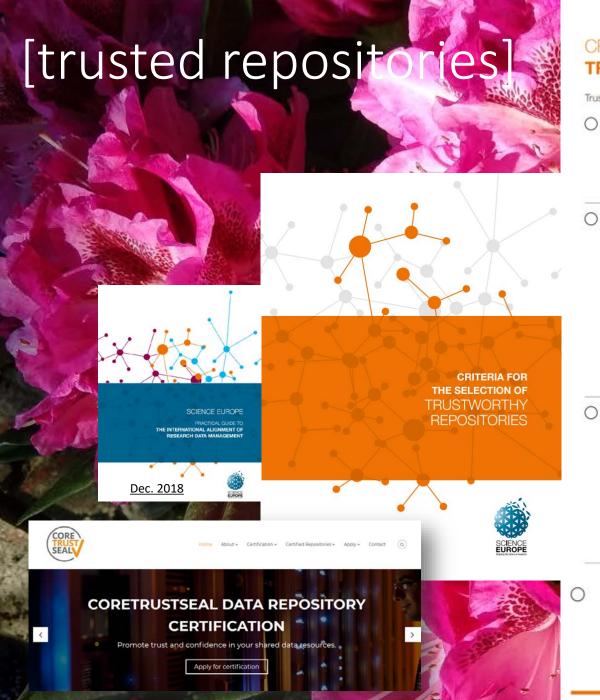
Read more

### One step closer towards instant DOI search results

Art Art? You might be wondering, what this pink and green picture illustrates? A few months ago we couldn't show you this picture; the data that we used to created it, did not exist. And the answer to what this illustrates – this is simply a distorted...

Read more

https://www.re3data.org/





### CRITERIA FOR THE SELECTION OF TRUSTWORTHY REPOSITORIES

Trustworthy repositories should meet the following minimum criteria:

- Provision of Persistent and Unique Identifiers (PIDs)
  - a. Allow data discovery and identification
  - b. Enable searching, citing, and retrieval of data
  - c. Provide support for data versioning
- O 2. Metadata
  - Enable finding of data
  - Enable referencing to related relevant information, such as other data and publications
  - Provide information that is publicly available and maintained, even for non-published, protected, retracted, or deleted data
  - d. Use metadata standards that are broadly accepted (by the scientific community)
  - e. Ensure that metadata are machine-retrievable
- 3. Data access and usage licences
  - a. Enable access to data under well-specified conditions
  - b. Ensure data authenticity and integrity
  - Enable retrieval of data
  - d. Provide information about licensing and permissions (in ideally machine-readable form)
  - Ensure confidentiality and respect rights of data subjects and creators
  - 4. Preservation
    - Ensure persistence of metadata and data
    - Be transparent about mission, scope, preservation policies, and plans (including governance, financial sustainability, retention period, and continuity plan)



### **Data Journals**

Hier entsteht eine Liste von Data Journals, die vorwiegend Data Papers

- Atomic Data and Nuclear Data Tables 
   (Elsevier)
- Biodiversity Data Journal 
   (Pensoft Publishers)
- Biomedical Data Journal 
   (Procon Ltd.)

- Data in Brief (Elsevier)

- Open Data Journal for Agricultural Research (diverse)
- Open Journal of Bioresources 

   (Ubiquity Press)
- Research Data Journal for the Humanities and Social Sciences ☑ (Brill)

#### Dataset Description

#### Object Name

- walkers three files providing the data, metadata and field type definitions (.csv, .txt, .csvt respectively) for records made by individual walkers during stage-one fieldwalking.
- counts three files providing the data, metadata and field type definitions (.csv, .txt, .csvt respectively) for potsherds countedduring stage-one fieldwalking.
- pottery three files providing the data, metadata and field type definitions (.csv, .txt, .csvt respectively) for the main pottery database, assembled various artefact specialists.
- petrography three files providing the data, metadata and field type definitions (.csv, .txt, .csvt respectively) for those sherds sampled for thin section petrography.
- lithics three files providing the data, metadata and field type definitions (.csv, .txt, .csvt respectively) for the main lithics database.
- other three files providing the data, metadata and field type definitions (.csv, .txt, .csvt respectively) for the main database of all non-ceramic and non-lithic finds.
- structs three files providing the data, metadata and field type definitions (.csv, .txt, .csvt respectively) for the main database of all standing remains, except for terraces.
- coast a vector polygon dataset (.shp and associated files) with the shape of Antkythera's coastline.
- geology –a vector polygon dataset (.shp and associated files) with the main bedrock units on Antkythera.
- tracts a vector polygon dataset (.shp and associated files) with the main stage-one survey units.
- $\cdot$  grids a vector polygon dataset (.shp and associated files) with the main stage-two survey units.
- terraces vector line dataset (.shp and associated files)
   with all observable agricultural terraces (i.e. the location)

IPER SKUME

- other primarily Andrew Bevan (UCL), with further assistance from James Conolly (Trent)
- geology a combination of fieldwork by Ruth Siddall (UCL) and remote sensing by Andrew Bevan (UCL)

#### Repository Location

I don't need

UK Archaeology Data Service Collection 1115 (doi: 10.5284/1012484)

Publication Date 05/02/2012

#### Language

English (a Greek language summary of the project methods and results can be found at www.ucl.ac.uk/asp/ or www.tuarc.trentu.ca/asp/).

#### License

Creative Commons CC-BY 3.0

#### Reuse Potential

Due to their unusual coverage of an entire landscape, these datasets would provided a good basis for developing a tutorial on survey, GIS and/or spatial analysis in archaeology. They also lend themselves to the comparative analysis of evidence from other intensive Mediterranean surveys that are in the public domain (ε ~ \_\_ttp://dx.doi.org/10.5384/1000371

public domain (¢ http://dx.doi.org/ org/10.5284/100 dx.doi.org/10.528 to the fact that th cal. The ASP data locations, dates a ally in the databas structures and ten

### **Data journals**

### Panayiota Polydoratou

Alexander Technological Educational Institute of Thessaloniki

European Commission Workshop

Alternative Open Access Publishing Models: Exploring New Territories in

Brussels, 12 October 2015

# A = Accessible. Formati

Data Archiving and Networked Service

HOME

DEPOSIT



Text documents

Plain text

Markup language

Spreadsheets

Databases

Statistical data

Raster images

Preferred format(s)

• PDF/A (.pdf)

Unicode text (.txt)

• XML (.xml)

HTML (.html)

• Related files: .css, .xslt, .js, .es

ODS (.ods)

CSV (.csv)

• SQL (.sql)

• SIARD (.siard)

• DB tables (.csv)

• SPSS Portable (.por)

• SPSS (.sav)

• STATA (.dta)

DDI (.xml)

• data (.csv) + setup (.txt)

JPEG (.jpg, .jpeg)

• TIFF (.tif, .tiff)

• PNG (.png)

• JPEG 2000 (.jp2)

Non-preferred format(s)

• ODT (.odt)

• MS Word (.doc, .docx)

• RTF (.rtf)

• PDF (.pdf)

• Non-Unicode text (.txt)

• SGML (.sgml)

MS Excel (.xls, .xlsx)

• PDF/A (.pdf)

• OOXML (.docx, .docm)

 MS Access (.mdb, .accdb) (v. 2000 or later)

dBase (.dbf)

HDF5 (.hdf5, .he5, .h5)

SAS (.7dat; .sd2; .tpt)

• R (\* under examination)

DICOM (.dcm) (by mutual agreement)

Type of data	Recommended formats	Acceptable formats				DECLINATE
Tabular data with extensive metada variable labels, co labels, and defined missing values	delimited text and command ('setup') file (SPSS, Stata,	proprietary formats of statistical packages: SPSS (.sav), Stata (.dta), MS Access (.mdb/.accdb)	fo	Textual data	Rich Text Format (.rtf)  plain text, ASCII (.txt)  eXtensible Mark-up  Language (.xml) text  according to an appropriate  Document Type Definition  (DTD) or schema	Hypertext Mark-up Language (.html) widely-used formats: MS Word (.doc/.docx) some software-specific formats: NUD*IST, NVivo and ATLAS.ti
Tabular data with minimal metadata column headings, variable names	a (.csv)	delimited text (.txt) with characte not present in data used as delimiters widely-used formats: MS Excel (.xls/.xlsx), MS Access (.mdb/.accdb), dBase (.dbf), OpenDocument Spreadsheet (.ods)	o)	Image data	TIFF 6.0 uncompressed (.tif)	JPEG (.jpeg, .jpg, .jp2) if original created in this format  GIF (.gif)  TIFF other versions (.tif, .tiff)  RAW image format (.raw)  Photoshop files (.psd)  BMP (.bmp)  PNG (.png)  Adobe Portable Document Format (PDF/A, PDF) (.pdf)
Geospatial data vector and raster data	ESRI Shapefile (.shp, .shx, .dbf, .prj, .sbx, .sbn optional) geo-referenced TIFF (.tif, .tfw) CAD data (.dwg)	Keyhole Mark-up Language (.km				
	Geography Markup Language (.gml) bir	Adobe Illustrator (.ai), CAD data (.dxf or .svg) binary formats of GIS and CAD packages.	ī	Audio data	Free Lossless Audio Codec (FLAC) (.flac)	MPEG-1 Audio Layer 3 (.mp3) if original created in this format  Audio Interchange File Format
		Video data		mp4) o (.ogv, .ogg) EG 2000 (.mj2)	AVCHD video (.avchd)	(.aif)  Waveform Audio Format (.wav)
UK Data Service	About us Get data Use data  Home > Manage data > Format your data > Recommende  Recommended formats  asservice.ac.uk/manage-data/fo	Documentation and scripts	PDF/UA, F XHTML or .htm)	Format (.rtf) PDF/A or PDF (.pdf) HTML (.xhtml,	plain text (.txt) widely-used formats: MS Word (.doc/.docx), MS Excel (.xls/.xlsx) XML marked-up text (.xml) according to an appropriate DTD or schema, e.g. XHMTL 1.0	

h

# - Accessible Conservazione

### LUNGO O BREVE TERMINE?



Software for Digital Preservation

Download version 3.0.1, released 25 March 2014 AEST

Checksum Checker is free and open source software developed by the National Archives of Australia. Checksum Checker is a piece of software that is used to monitor the contents of a digital archive for data loss or corruption.

Checksum Checker is a component of the Digital Preservation Software Platform (DPSP).

As part of the Digital Preservation Recorder (DPR) workflow, checksums are generated for each Archival Information Package (AIP). Checksum Checker generates a new checksum for each AIP and compares it against the stored checksum. If the checksums do not match, then the AIP is flagged as being corrupt.

Checksum Checker incorporates the following features:

- · Checksum Checker functions as a service.
- · Checksum Checker sends automated emails to a nominated administrator email address, coinciding with certain events (such as the start of a checking run or when an error is

Contact Us

Download

F.A.O

Licensing

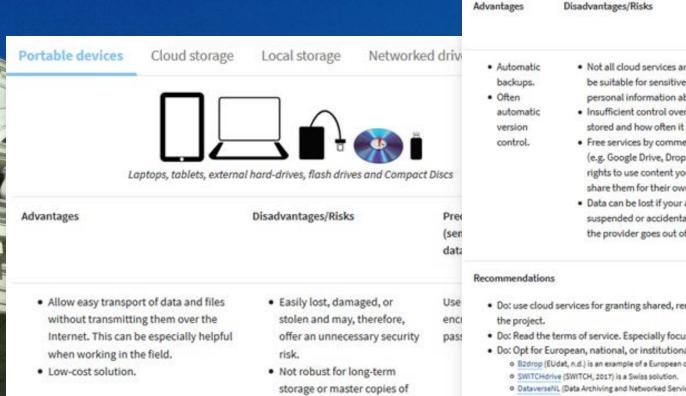
External

Storage Solutions Advantages Disadvantages Suitable for Personal Computer & Laptop Alwavs available Drive may fail Temporary storage Portable Laptop may be stolen Networked drives Regularly backed up Costs Master copy of vour data File servers managed by your Stored securely in a university, research group or facilities single place (if enough like a NAS-server storage space is provided ..) Easily damaged or lost External storage devices Temporary Low cost storage USB flash drive, DVD/CD, external Portability hard drive Cloud services Automatic It's not sure whether data security Data sharing synchronization is taken care of between folders and You don't have direct influence on how often backups take place and Easy to access and by whom

http://checksumchecker.sourceforge.net/

Organize and document research data. Make digital versions of paper data documentation in a PDF/A format (suitable for long-term storage).

# A = Accessible. Conservazione



your data and files.

· Possible quality control

issues due to version

confusion.

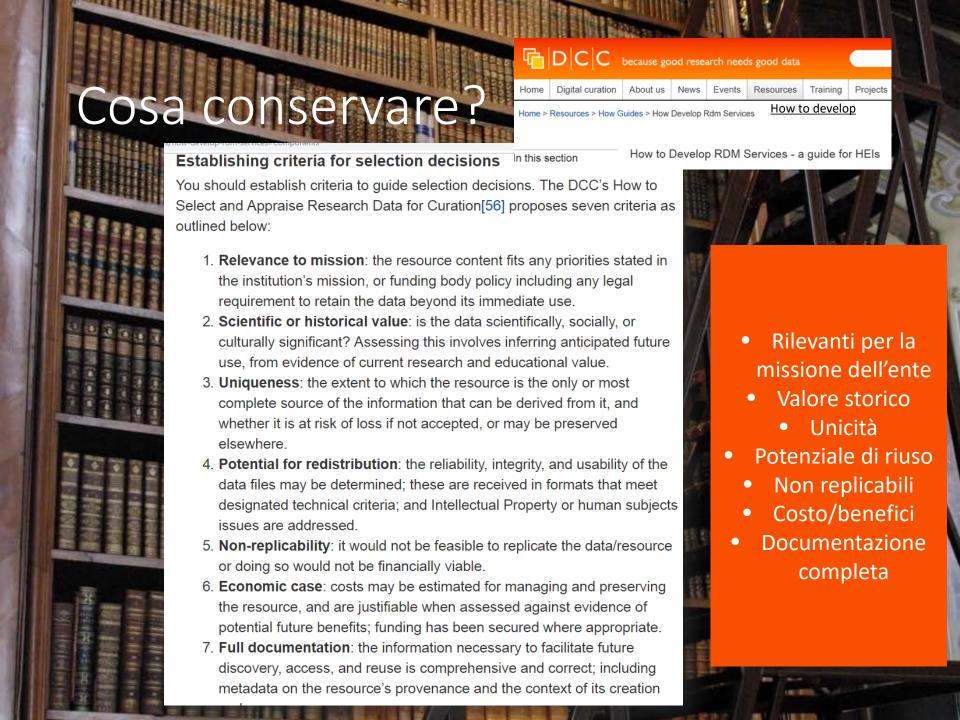
Precautions for (sensitive) personal

- · Not all cloud services are secure. May not be suitable for sensitive data containing personal information about EU citizens.
- · Insufficient control over where the data is stored and how often it is backed up.
- Free services by commercial providers (e.g. Google Drive, Dropbox) may claim rights to use content you manage and share them for their own purposes.
- . Data can be lost if your account is suspended or accidentally deleted, or if the provider goes out of business.
- · Encrypt all (sensitive) personal data before uploading it to the cloud. This is particularly important to avoid conflict with European data protection regulations if you do not know in which countries servers used for storage and backup are located (see 'Security' for more information on encryption; also see 'Protecting data').

- · Do: use cloud services for granting shared, remote and easy access to data and other files to all involved in
- Do: Read the terms of service. Especially focus on rights to use content given to the service provider.
- Do: Opt for European, national, or institutional cloud services which store data in Europe if possible.
  - B2drop (EUdat, n.d.) is an example of a European cloud storage solution.
  - DataverseNL (Data Archiving and Networked Services, 2017) is an example of a service for Dutch researchers that allows the storage and sharing of data both during and after the research period.
- Don't: make this your only storage and backup solution.
- Don't: use for unencrypted (sensitive) personal data.

**CESSDA Guide** 

Differenti bisogni, differenti strumenti. Durante l'esperimento, dovete poter condividere con il team



### I = Interoperable Standards



**PARTHENOS** 

TRAINING MODULES FOR TRAINER WHAT ARE KNOWLEDGE REPRESENTATION SYSTEMS AND **'ONTOLOGIES'?** 

### WHAT ARE STANDARDS?

Even perfect metadata may not allow data to become interoperable if a different standard commonly as an 'ontology'. Before the digital age, philosophers referred to an ontology as "the study of used. A "standard" refers to a system that structures what types of information are capture kinds of things that exist. Ontologies are similar to taxonomies, another knowledge organisation item in a collection. In our .mp3 library system, a standard is expressed in the header categories such as 'name,' 'time,' 'artist,' and 'album' are listed, with every entry having this filled in. Standards are used to ensure that metadata is as useful as possible for organising collection, ensuring that common questions (how many songs are there on the album "Big B can be easily and accurately answered.

addition to metadata and standardised metadata schemas, research infrastructures can also use ther forms of "knowledge representation system" to enhance the researcher's experience of the teroperable data they present. When we talk about 'Knowledge Representation Systems' in research frastructures, we usually mean a specific category of hierarchical systems of terms known more amework you probably remember from early lessons in biology.



### How Many Standards Are There and Who Decides Which One To Use?

Different standards have arisen in different kinds of cultural heritage institution: the most common standards in museums are different from those in archives, and those common in libraries are different again.

### What are Standards?

What Are Knowledge Representation Systems and 'Ontologies'?

Sustainability

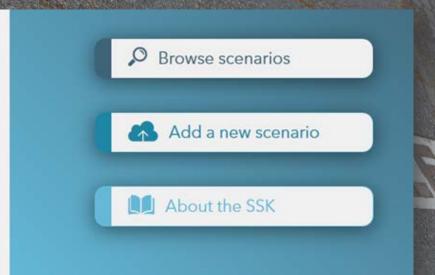
Methods and Tools

Networks

# I = Interoperable. Standards

## Standardization Survival Kit

A collection of research use case scenarios illustrating best practices in Digital Humanities and Heritage research



Increase efficiency, interoperability and sustainability by using standards

Incorporating standards in all the steps of your research process will make it last longer, easier to update, improve and share. Standards are non legally binding documents produced by an organisation ensuring:

# I = Interoperable

To speed up discovery and uncover new intights, research data should be easily combined with other datasets by humans as well as computer systems.

### INTEROPERABLE



Well documented and machine-actionable APIs - a set of subroutine definitions, protocols, and tools for building application software - allow for automatic indexing, retrieval and combining of (metaldata from different data repositories.



Document APIs well and make it possible to deliver the schema of the (meta)data model. Consider showing examples of how to successfully mine data from different endpoints and combine them into new data sets usable for new research.



The description of metadata elements should follow community guidelines that use open, well defined and well known vocabularies. Such vocabularies describe the exact meaning of the concepts and qualities that the data represent.



Use vocabularies relevant to your field, and enrich and structure your research output accordingly from the start of your research project.



Give examples of vocabularies the research community may use, based on research domain specifics.

### 12 Document metadata models

Clearly documenting metadata models helps developers to compare and make mappings between metadata.



Publish the metadata models in use in your research infrastructure. Document technical specifications and define classes (groups of things that have common properties) and properties (elements that express the attributes of a metadata section as well as the relationships between different parts of the metadata). For metadata mapping purposes, list the mandatory and recommended properties.



### Prescribe and use interoperable data standards

Using a data standard backed up by a strong community, increases the possibility to share, reuse and combine data collections.



Check with the repository where you want to deposit your data what data standardsthey use. Structure your data collection in this format from the start of your research project.



Clearly specify which data standard your institution uses, pool a community arround them and maintain them especially with a perspective on interoperability. Good examples are CMDI (language studies) and the SKB0102 Standard (archaeology).



#### Establish processes to enhance data quality

To boost (meta)data quality and, therefore, interoperability, establish (automatic) processes that clean up, derive and enrich (meta)data.



Establish procedures to minimise the risk of mistakes in collecting data.

E.g. choose a date from a calendar instead of filling it in by hand.



Invest in tools to help clean up (meta)data and to convert data into standardised and interoperable data formats. Combine efforts to develop workflows and software solutions for such automatic processes, e.g. by using machine learning tools.



#### Prescribe and use future-proof file formats

All data files held in a data repository should be in an open, international, standardised file format to ensure long-term interoperability in terms of usability, accessibility and sustainability.







**FAIRDOM Platform** 

# I= Interoperable – ac



### RightField 🧽

Rightfield is an open-source tool for adding ontology term selection to Excel spreadsheets. Rightfield is used by a 'Template Creator' to create semantically aware Excel spreadsheet templates. The Excel templates are then reused by Scientists to collect and annotate their data; without any need to understand, or even be aware of, Rightfield or the ontologies used. Rightfield embedded templates are used within the Samples framework of the SEEK.

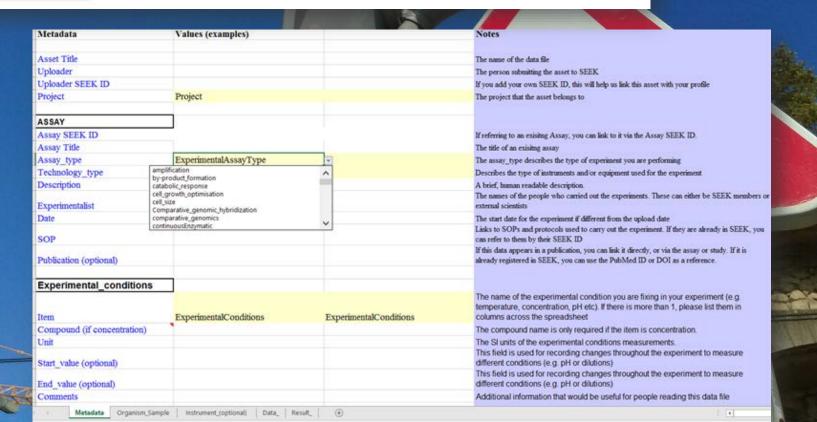
MORE INFORMATION

https://fair-dom.org/platform/rightfield/

fic research datasets, models or simulations, processes and nation about the people and organisations involved. The s) based on the ISA-Tools format. When paired with our ction through to publication. Norwegian users benefit from Pt simplifies upload and download of files.

ata sharing within groups and consortia. In addition,

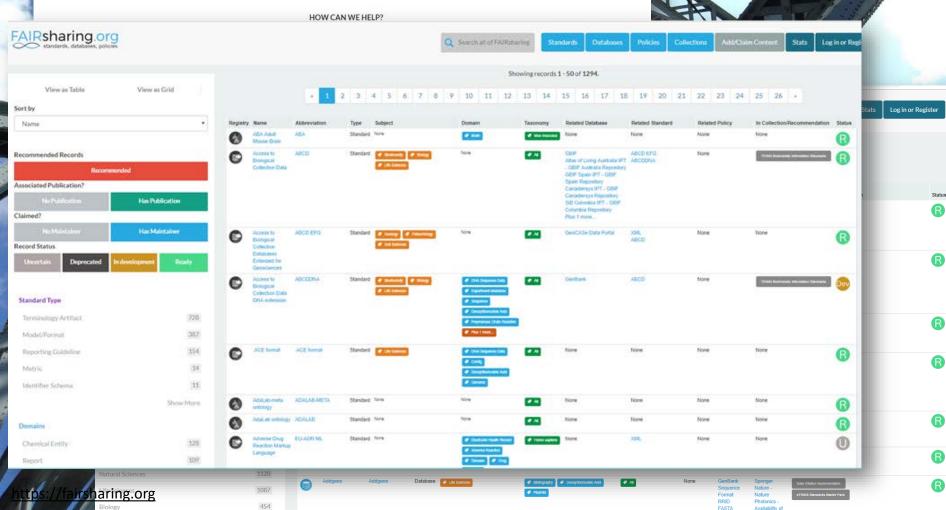
https://fair-dom.org/platform/





FAIR Sharing org
tradicely, detables, policies

A curated, informative and educational resource on data and metadata standards, inter-related to databases and data policies.



# R = Reusable. Documentazione

### Project-level documentation





letail:

Project-level documentation explains the aims of the study, what the research questions/hypotheses are, what methodologies were being used, what instruments and measures were being used, etc. In the accordion the questions which your project-level documentation should answer are stated in more

- ① 1. For what purpose was data created
- ① 2. What does the dataset contain
- 3. How was data collected
- # 4. Who collected the data and when
- ① 5. How was the data processed
- 🕀 6. What possible manipulations were done to the data
- + 7. What were the quality assurance procedures
- 3. How can data be accessed

### Data-level documentation

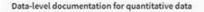
Data-level or object-level documentation provides information at the level of individual objects such as pictures or interview transcripts or variables in a database. You can embed data-level information in data files. For example, in interviews, it is best to write down the contextual and descriptive information about each



interview at the beginning of each file. And for quantitative data variable and value names can be embedded within the data file itself.

### O Quantitative data

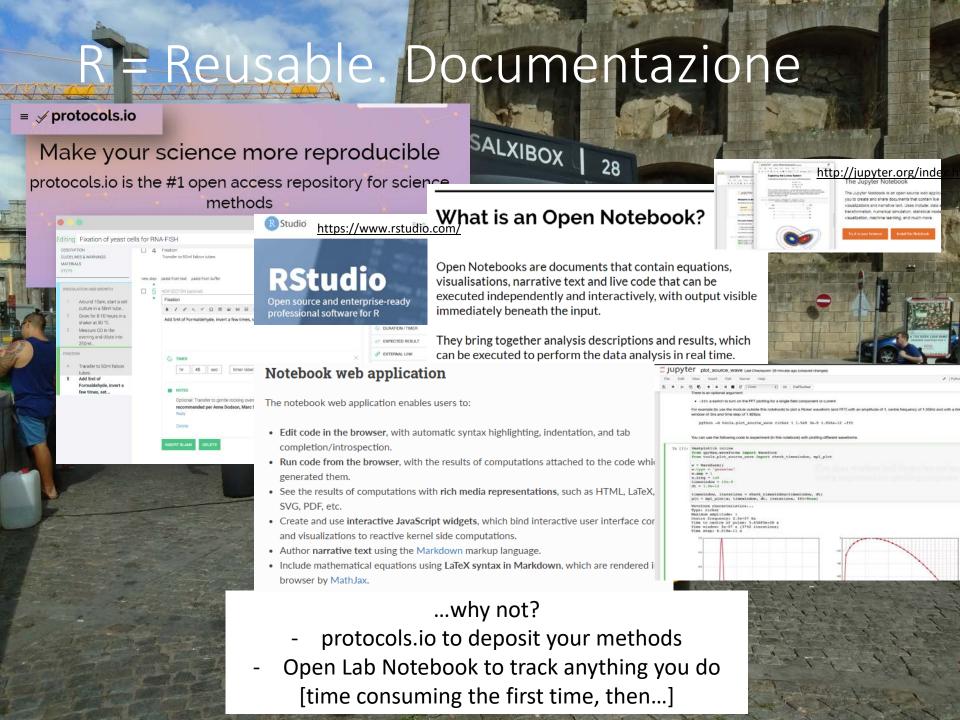
Variable-level annotation should be embedded within a data file itself. If you need to compile an extensive variable level documentation that can be created by using a structured metadata format.



For quantitative data document the following:

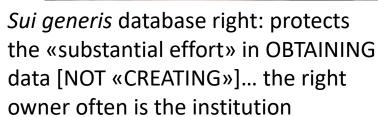
- Information about the data file
   Data type, file type and format, size, data processing scripts.
- Information about the variables in the file
   The names, labels and descriptions of variables, their values, a description of derived





# R=Reusable. Licenze

Copyright: protects the STRUCTURE, selection or arrangement of their contents" (Art. 3) NOT THE DATA





AND

works, data or other materials arranged in a systematic or methodical way (Art.1)

Database=a collection of independent

Simone Aliprandi

2014

RICORDA: NESSUN COPYRIGHT SUI DATI (NON CREATIVI) DIRECTIVE 96/9/EC OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL

on the legal protection of databases

NCIL OF THE EUROPEAN UNION

an Community, and in particular Article 57 (2), 66 and 100a thereof,

a QUALI DIRITTI SUI DATI?



semplici dati e informazioni

nessuna tutela



solo diritto sui generis



livello diritto d'autore

livello diritto sui generis



diritto sui generis + diritto d'autore



della ricerca

**OpenAIRE** 

**Thomas Margoni** University of Glasgow - CREATe OpenAIRE project

### **TRAINING**

Webinars

# [non suoniamo tutti la stessa musica]

OLA E MEGLIO

① Czech Republic

⊕ Finland

Obstacles to the trans-European archiving and sharing of research data

Making research data as openly available as possible is a widely recognised goal. For researchers working on an interdisciplinary project involving several countries, it can be difficult to fully comprehend in which ways open access to research data can be legally obtained. European national laws still diverge.

· Diversity in copyright owner

If protection applies, the right holder's consent is required for sharing the data. However, the designation of the copyright owner is also different in different jurisdictions. Although in many cases the maker of the work will be considered to be the author and therefore the right holder, only Dutch and UK law designate the employer as the right holder if the work was made in the course of employment.

A report from Knowledge Exchange (Knowledege Exchange, 2011) concludes that it will remain difficult to predict when particular files of research data are protected because of:

Diversity in copyright protection

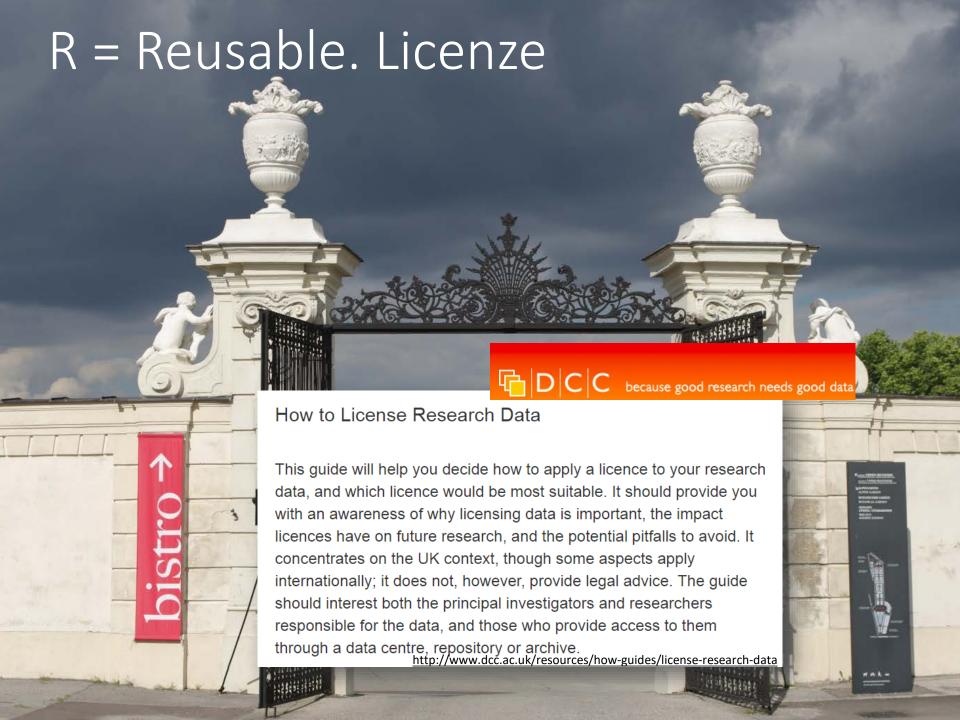
Even though most research data will fail to meet the criteria for copyright protection because they are not likely to be considered as "works" (they mainly concern facts), the lack of harmonisation of the criteria for copyright protection in Europe is tricky. E.g., whereas Germany, Denmark and the Netherlands have a relatively similar (higher) originality standard, the UK has a very low standard (skill, judgment and labour) making it

⊕ Switzerland

CESSDA guide

⊕ uk





This guide will help you decide how to apply a licence to your research data, and which licence would be most suitable. It should provide you with an awareness of why licensing data is important, the impact licences have on future research, and the potential pitfalls to avoid. It concentrates on the UK context, though some aspects apply internationally; it does not, however, provide legal advice. The guide should interest both the principal investigators and researchers responsible for the data, and those who provide access to them

through a data centre, repository or archive.

### R = Reusable. Licenze

### Creative Commons at a glance

### Good for

- very simple, factual datasets
- data to be used automatically

### Watch out for

- versions: use v. 4 or later
- attribution stacking
- the NC condition: only use with dual licensing
- the SA condition as it reduces interoperability
- the ND condition as it severely restricts reuse

### ODC-By at a glance

### Good for

- most databases and datasets
- data to be used automatically
- data to be used for generating non-data products

### Watch out for

attribution stacking

### Public domain at a glance

### Good for

- most databases and datasets
- data to be used by anyone or any tool
- data to be used for any purpose

### Watch out for

- lack of control over how database is reused
- lack of protection against unfair competition

### ODC-ODbL at a glance

### Good for

- most databases and datasets
- data to be used automatically
- data to be used for generating non-data products

### Watch out for

- attribution stacking
- the copyleft condition as it reduces interoperability

clause as it may put off some reusers

the DRM clause as it may put off some reusers

ack of protection against unfall charactition





### **FACT SHEET ON** CREATIVE COMMONS & OPEN SCIENCE

### https://doi.org/10.5281/zenodo.840651

### What is Open Science?

Open Science is the movement to make scientific research and data accessible to all for knowledge dissemination and

How should I licence my data for the purposes of Open Science?

We recommend you use the CCO Public Domain Dedication, which is first and foremost a waiver, but can act as a cence when a waiver is not possible

CC ZERO LICENCE, 'NO RIGHTS



one to freely reuse your data as they see fit by waiving (giving up) your copyright and related

ations in which data is not protected as a matter of law. Such data can include facts, names, num pers - things that are considered 'non-original' and part of the public domain thus not subject to copyright protections. Similarly, your database. which is a structured collection of data) might be considered 'non-original' and thus ineligible for copyright, and it might additionally be excluded

from other forms of protection (like the EUs

licence such as a CC BY could signal to users that you claim a copyright in the non-original data

Finally, if your data is in the public domain worldwide, you might state simply and obviously on the material that no restrictions attach to the reuse of your data and apply a Public Dor

**PUBLIC DOMAIN MARK LOGO** 



When in doubt, consider which use may be appro priate according to the chart below

CCO & PUBLIC DOMAIN LICENCES WHICH LICENSE TO USE AND WHEN





original; the authorized acknowledges this and communicat the data is in the

ah nins

### What is Open Science?

Open Science is the movement to make scientific research and data accessible to all for knowledge dissemination and public reuse.

### How should I licence my data for the purposes of Open Science?

We recommend you use the CCO Public Domain Dedication, which is first and foremost a waiver, but can act as a licence when a waiver is not possible.

**CC ZERO LICENCE, 'NO RIGHTS RESERVED' LOGO** 



By applying CCO to your data you enable everyone to freely reuse your data as they see fit by waiving (giving up) your copyright and related rights in that data.

You should keep in mind that there are many situations in which data is not protected as a matter of law. Such data can include facts, names, numbers - things that are considered 'non-original' and part of the public domain thus not subject to copyright protections. Similarly, your database (which is a structured collection of data) might be considered 'non-original' and thus ineligible for copyright, and it might additionally be excluded

from other forms of protection (like the EU sui generis database right, also known as the 'SGDR', for non-original databases).

In these cases, using a Creative Commons licence such as a CC BY could signal to users that you claim a copyright in the non-original data despite the law, and perhaps despite your real intention.

Finally, if your data is in the public domain worldwide, you might state simply and obviously on the material that no restrictions attach to the reuse of your data and apply a Public Domain Mark.

### **PUBLIC DOMAIN MARK LOGO**



When in doubt, consider which use may be appropriate according to the chart below:

### **CCO & PUBLIC DOMAIN LICENCES** WHICH LICENSE TO USE AND WHEN



'Creative arrangement' of data is original, but any copyright has been waived and content is made available copyright-free

TOX MINORY to million



'Creative arrangement' of data is not original; the author acknowledges this and communicates the data is in the public domain

But I would like attribution when others use my dataset. In that case, shouldn't I use a CC BY licence?

We recommend that you avoid using a CC BY licence. Here's why:

While attribution is a genuine, recognisable concern, not only might using a CC BY licence be legally unenforceable when no underlying copyright or SGDR protects the work, but it may also communicate the wrong message to the world. A better solution is to use CCO and simply ask for credit (rather than require attribution), and provide a citation for the dataset that others can copy and paste with ease. Such requests are consistent with scholarly norms for citing source materials.

Legally speaking, datasets that are *not* subject to copyright or related rights (and are thus in the public domain) cannot be the object of a copyright licence. Despite this, agreements based in contract law may be enforceable. Creative Commons licences, however, are copyright licences. Therefore, where the conditions for a copyright or related right are not triggered, copyright licences, such as the CC BY licence, are unenforceable.

In some cases, however, rights may exist (like the sui generis database right previously mentioned), and permission for others to use your dataset may be legally required. These rights are meant to protect the maker's investment, rather than originality. As such, database rights do not include the moral right of attribution. So by using a CC BY licence, you signal to users that you restrict access to your dataset beyond the protections provided by the law. We are not saying that this cannot be done, we are just saying that if you choose to do this, you should make sure you fully understand what it entails.

### mons e Op

### **USE A CCO**

- THEN ASK FOR CREDIT
- PROVIDE A
  CITATION TO C&P
- BEAR IN MIND IT'S
  BAD SCIENCE NOT
  TO CITE THE
  SOURCE
  - CC0 DOES NOT
    MEAN ACADEMIC
    UNPOLITENESS

It sounds like you're really pushing for the use of CCO for open science datasets.

Exactly. Data is only open if anyone is free to use, reuse, and distribute it. This means it must be made available for both commercial and non-commercial purposes under non-discriminatory conditions that allow for it to be modified.

When data is made available for all reuse, others can create new knowledge from combining it. This leads to the enrichment of open datasets and further dissemination of knowledge. Accordingly, CCO is ideal for open science as it both protects and promotes the unrestricted circulation of data.

And remember, it's bad science not to cite the source of data you use. To help others cite your data include a citation that users can copy and paste to give you credit for your hard work.

cannot be done, we are just saying that if you choose to do this, you should make sure you fully understand what it entails.

I'm uncomfortable with others using my research for commercial purposes. Should I use a non-commercial licence for my dataset?

We recommend you avoid using a non-commercial licence. Here's why:

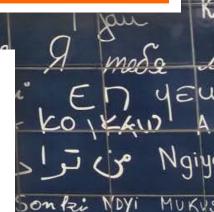
For legal purposes, drawing a line between what is and is not 'commercial' can be tricky; it's not as black and white as you might think. For example, if you release a dataset under a non-commercial licence, it would clearly prohibit an organisation

I'm uncomfortable permitting use of my research for any and all purposes. Should I use a 'No Derivatives' (ND) licence for my dataset?

We recommend you avoid using a 'No Derivatives' licence. Here's why:

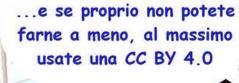
Similar to how a non-commercial licence might restrict meaningful reuse of your dataset, a ND licence can have the same effect: it may prevent someone from recombining and reusing your data for new research. For data to be truly Open Access, it must permit these important types of

Tak Ilmogu to million





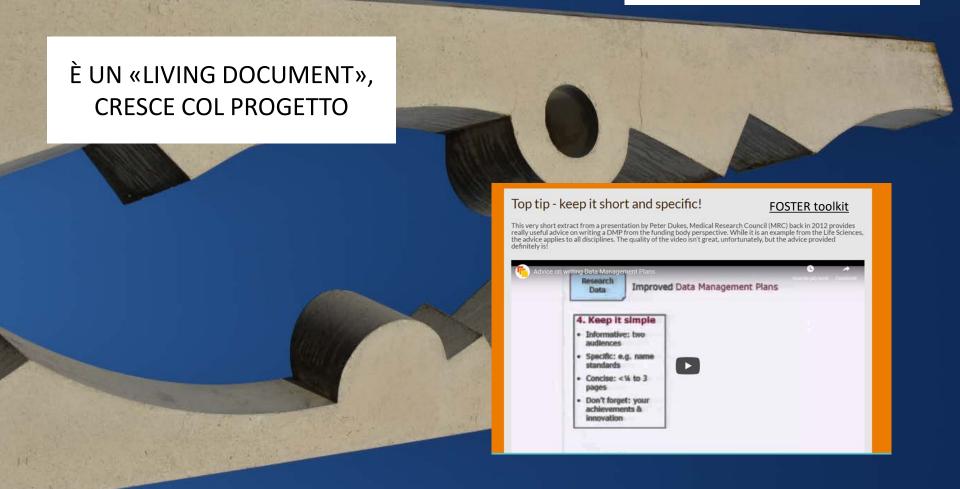
Ci sono solo TRE modi per rilasciare open data senza creare inutili complicazioni: CC Zero, CC Zero e CC Zero



S. Aliprandi

UN MODO STRUTTURATO DI PENSARE AI DATI

REGOLE CHIARE=MENO ERRORI DA SUBITO



DOVE METTERE TUTTE QUESTE INFORMAZIONI? NEL DATA MANAGEMENT PLAN

## Vantaggi di un DMP

A Data Management Plan

Useful tool to think ahead

Allows for easy project management

Clarifies needed budget

Makes data FAIRer

Shows accountability

Benefit 3. Clarifies needed budget

Data management is not free. You do not want to find yourself running out of funding before the end of the project because you have ignored or underestimated the cost of structured, detailed, and safe data management. Therefore, an important aspect of a DMP is its use in calculating how much money will be required for managing your research data during your research project.

A DMP can be useful in the process of applying for funding. Grant applications should not only include time and resources for collecting, analysing, and publishing on data in their budget, time and resources for careful documentation as well as server space, backup solutions, and documentation software need to be included as well. A DMP is also useful once funding is granted to plan and manage your expenses. Many research funders require a DMP as part of the application and decision-making process. The arguments for making data available are several, the most popular being that the data produced by public funds should be used to the greatest extent possible and available to the public. Unless there are legal, ethical or commercial barriers, data should also be openly available so that research results can be verified, replicated and reused.

Examples of Data Management cost assessments are given by the <u>University of Utrecht</u> (n.d.) and the Dutch Landelijk Coördinatiepunt Research Data Management (<u>LCRDM</u>, n.d.) inspired by the 'Data management costing tool' by UK Data Service, 2013.

**CESSDA Guide** 

È FONDAMENTALE PER STIMARE I COSTI DI GESTIONE

# DMP Core Requirements CORE REQUIREMENTS FOR DATA MANAGEMENT PLANS

When developing solid data management plans, researchers are required to deal with the following topics and answer the following questions:

- 1. Data description and collection or re-use of existing data
  - a. How will new data be collected or produced and/or how will existing data be re-used?
  - b. What data (for example the kinds, formats, and volumes) will be collected or produced?
- 2. Documentation and data quality
  - a. What metadata and documentation (for example the methodology of data collection and way of organising data) will accompany data?
  - b. What data quality control measures will be used?
- 3. Storage and backup during the research process
  - a. How will data and metadata be stored and backed up during the research process?
  - b. How will data security and protection of sensitive data be taken care of during the research?
- Legal and ethical requirements, codes of conduct
  - If personal data are processed, how will compliance with legislation on personal data and on data security be ensured?
  - b. How will other legal issues, such as intellectual property rights and ownership, be managed? What legislation is applicable?
  - c. How will possible ethical issues be taken into account, and codes of conduct followed?



SCIENCE EUROPE

- Data sharing and long-term preservation
  - How and when will data be shared? Are there possible restrictions to data sharing or embargo reasons?
  - How will data for preservation be selected, and where will data be preserved long-term (for example a data repository or archive)?
  - c. What methods or software tools will be needed to access and use the data?
  - d. How will the application of a unique and persistent identifier (such as a Digital Object Identifier (DOI)) to each data set be ensured?
- 6. Data management responsibilities and resources
  - a. Who (for example role, position, and institution) will be responsible for data management (i.e. the data steward)?
  - b. What resources (for example financial and time) will be dedicated to data management and ensuring that data will be FAIR (Findable, Accessible, Interoperable, Re-usable)?







DMP Core Requirements

### Translating the Core Requirements into a DMP template

The following example of a data management plan template is based on the core requirements for DMPs. These core requirements should be considered as a minimum standard, leaving the flexibility to formulate additional guidelines according to the needs of specific domains or to national or local legislation.

The template presented below refers to the 15 questions covering six core requirements for good data management. Additional guidance and explanations are provided to help researchers fill out such a template and to assure that all relevant aspects of research data management are covered. The below table is an example of how the core requirements can be transformed into a DMP template. It will be up to the individual organisations and disciplines to develop templates that fit their needs.

### **GENERAL INFORMATION**

Administrative information

 Provide information such as name of applicant, project number, funding programme, version of DMP.

### 1 DATA DESCRIPTION AND COLLECTION OR RE-USE OF EXISTING DATA

1a

How will new data be collected or produced and/or how will existing data be re-used?

- Explain which methodologies or software will be used if new data are collected or produced.
- State any constraints on re-use of existing data if there are any.
- Explain how data provenance will be documented.
- Briefly state the reasons if the re-use of any existing data sources has been considered but discarded.



SCIENCE EURO

PRACTICAL GUIDE TO THE INTERNATIONAL ALIGNMENT OF RESEARCH DATA MANAGEMENT





### 2 DOCUMENTATION AND DATA QUALITY

2

What metadata and documentation (for example the methodology of data collection and way of organising data) will accompany the data?

- Indicate which metadata will be provided to help others identify and discover the data.
- Indicate which metadata standards (for example DDI, TEI, EML, MARC, CMDI) will be used.
- Use community metadata standards where these are in place.
- Indicate how the data will be organised during the project, mentioning for example conventions, version control, and folder structures. Consistent, well-ordered research data will be easier to find, understand, and
- Consider what other documentation is needed to enable re-use. This may include information on the methodology used to collect the data, analytical and procedural information, definitions of variables, units of measurement, and so on.
- Consider how this information will be captured and where it will be recorded for example in a database with links to each item, a 'readme' text file, file headers, code books, or lab notebooks.

2b

What data quality control measures will be used?

Explain how the consistency and quality
of data collection will be controlled and
documented. This may include processes
such as calibration, repeated samples or
measurements, standardised data capture,
data entry validation, peer review of data, or
representation with controlled vocabularies.

# DMP questions

Adapt your Data Management Plan

A list of Data Management Questions based on the Expert Tour Guide on Data Management





### ORGANISE & DOCUMENT

### Overview

Title of the project

Date of this plan

Description of the project

- . What is the nature of the project?
- What is the research question?
- What is the project time line?

Origin of Data

- . What kind of data will be used during the project?
- If you are reusing existing data: What is the scope, volume and format? How are different data sources integrated?
- . If you are collecting new data can you clarify why this is necessary?

Principal researchers

- . Who are the main researchers involved?
- · What are their contact details?

Collaborating researchers (if applicable)

. What are their contact details and their roles in the project?

Funder (if applicable)

. If funding is granted, what is the reference number of the funding granted?

Data producer

. Which organisation has the administrative responsibility for the data?

Project data contact

. Who can be contacted about the project after it has finished?

Data owner(s)

. Which organisation(s) own(s) the data?

https://www.cessda.eu/content/download/4302/48656

. If several organisations are involved, which organisation owns what data?

Roles

- . Who is responsible for updating the DMP and making sure that it's followed?
- · Do project participants have any specific roles?
- · What is the project time line?

Costs

- Are there costs you need to consider to buy specific software or hardware?
- . Are there costs you need to consider for storage and backup?
- · Are potential expenses for (preparing the data for) archiving covered?

### Organising and documenting your data

### Data collection

- · How will the data be collected?
- · Is specific software or hardware or staff required?
- Who will be responsible for the data collection?
- During which period will the data be collected?
- · Where will the data be collected?

### Data organisation

- How will you organise your data?
- Will the data be organised in simple files or more complex databases?
- How will the data quality during the project be ensured?
- If data consists of many different file types (e.g. videos, text, photos), is it possible to structure the data in a logical way?

#### Data type and size

- · What type(s) of data will be collected?
- What is the scope, quantity and format of the material?
- After the project: What is the total amount of data collected (in MB/GB)?

### File format

- . In what format will your data be?
- Does the format change from the original to the processed/final data?
- Will your (final) data be available in an open format?

### Folder structure and names

• How will you structure and name your folders?

### File structure and names

. How will you structure and name your files?

#### Documentation

- What documentation will be created during the different phases of the project?
- How will the documentation be structured?

### Metadata

ile/TTT DO DMPExpertGuide v1.2.pc

- What metadata will be provided with the collected/ generated/ reused data?
- How will metadata for each object be created?
- Is there any program that can be used to document the data?
- Can metadata be added directly into the files or will the metadata be produced in another program or document?

#### Metadata standard (if applicable)

What metadata standard(s) will you use?

# DMP questions

cessda



Processing your data

### Versioning

- . What is your strategy concerning versioning your data files (and scripts) during the project?
- . Will you create and/or follow a convention for versioning your data?
- Who will be responsible for securing that a "Masterfile" will be maintained, documented and versioned according to the project guidelines?
- · How can different versions of a data file be distinguished?

### Interoperability

 Will you make use of established software and hardware? If not, how does the software and hardware you use relate to other research?

### If applicable:

- Will you make use of established terminologies/ontologies (i.e. structured controlled vocabularies) in the project? If not, how do your terminologies relate to established ones?
- Which coding is used (if any)? Will you build on established coding schemes? If not, how does your coding relate to other research?

STORE PR

o

### Storing your data

#### Storage

- . How and where will the data be stored during the project?
- . For how long will the data be stored?

#### Backup

- . How, where and at what intervals will the data be backed-up?
- How will data be recovered in the case of a data loss incident?

### Security

- . How will sensitive data be protected? (if applicable)
- . How will data access be managed?

### Protecting your data

Ethical review (if applicable)

. Does your project require approval by a local ethics committee?

Informed consent (if applicable)

- . Do you require informed consent for your project?
- . If so, how will permission be obtained?
- . How are consent files organised and stored?

(sensitive) Personal data /confidential information (if applicable)

- . How will access to (sensitive) personal data during the project be controlled?
- How will collaborators be granted access to the data in a secure way?
- If the research project is going to have data that includes confidential information or information that requires informed consent, is there a requirement to notify a privacy officer?
- Is there any confidential information within the material that requires special treatment and/or limits the access to it during/after the project?
- . How will the material be protected during/after the project?
- . How will permissions and restrictions be enforced?

### Intellectual property rights (IPR)/Copyrights

- · Are there IPR or copyright issues to consider?
- . Will permission be needed to collect/reuse the data?
- Will these rights be transferred to another organisation for data distribution and archiving?

Agreements (if applicable)

. What are the agreements with other stakeholders?

Restrictions (if applicable)

. Are there any other restrictions that need to be considered?



### Archiving and publishing your data

### Archiving

- How and where will the data be stored after the project's completion?
- Will you archive your data in a trusted data repository?
- Will your data receive a persistent identifier?

### Data formats

- What formats will you provide your data in for archiving (and sharing)?
- Will specific software be required to process your data? Can this software be deposited with the data?

### Access (if applicable)

- Will your data be available (Open Access)?
- Will all data or only parts of it be published?
- · What licenses do you need for your data?
- · How should your data be cited when reused?
- Will there be an embargo period for (all or some of) the data?
- Are there other agreements or restrictions (see above) that need to be considered?
- Are there any legal/ethical restrictions that prevents the publication of all the material?
- Will these restrictions mean that action must be taken before the material can be made available?
- Is there a risk of delayed publication/making data available (all or parts of)?
   If so what might be needed to do to avoid this?

## DISCOVER

### **Discovering data**

### Identification of needs

- Do you plan to use existing data for your research?
- · What is the purpose for which you need the data?
- What do you want to learn from the data?
- What type of data do you need?

### Search for data

- Do you know where the data may be located?
- How do you plan to search for the data?

### Evaluation of data quality

- What is the minimal required quality of the data (in terms of origin, contents, scope, size, methods, etc.)?
- How do you plan to evaluate data quality (evaluation of metadata, tests, analysis, comparisons)?

### Gaining access to data

- What are the (expected) terms and conditions for data access and use?
- What is the (expected) process for gaining access to the data?
- What is the (expected) time-span of the process for gaining access to the data?
- What are the (expected) costs for data access and use?

### Basic Information.

- · State the purpose of the data collection/generation.
- · Explain the relation to the objectives of the project
- · Consider what data will be collected or created as part of the study (RAW data).
- Consider what data will be produced by processing the RAW data (Secondary, processed data).
- · Specify if existing data is being re-used (if any)
- · Specify the origin of the data
- . Specify the types and formats you plan to use for the data generated/collected (raw, processed, published).
- Consider what data will be published as the result of your study (Published data).

### Volume and Life Cycle of the Data.

If you are using FAIRDOM, we will look after data that will be retained and potentially exchanged by your projects. It will help with local storage for temporarily-held local data prior to processing.

### For RAW data, please consider the following:

- . How much RAW data you think will be produced (Estimates, per month, year, full project duration)?
- Will all of the RAW data be kept for the duration of the study or will the RAW data be deleted once it is processed?
- · For large scale RAW data (images, sequence) have you planned the local storage capacity necessary for processing?
- . Do you require help to organise a suitable local management system for RAW data?
- Do you have policies that govern the management and usage of RAW data?
- How long will RAW data be kept? Will there be a long-term archive?

### For Secondary and Published data, please consider the following:

- · What data processing is foreseen in the project?
- · How much processed data will be produced, and stored (can you make estimates per month, year, full project)?
- . How much of this data will be published? (Estimates per month, year, full project)?
- · Does your institution, or the project funders, have policies governing the access and usage of processed data?

### Additional for personally sensitive data (e.g medical data)

- . When looking at the data flow through the project, define what data is:
  - · aggregated (typically safe to share, if names cannot be recovered)
  - · anonymized (name cannot be recovered from the data)
  - · pseudonymized (name can be recovered by some)
  - non-anonymized (name linked to data)
- · Determine which organisational boundaries have to be traversed by which data.
- Make sure with your "local" data protection officer and ethics commission that the data can be shared with your partners along the flow described with the anonymisation levels as described. Why local? Some
  laws change across surprising boundaries. E.g. in Germany Universities and other public organisations are subject to another data protection law than enterprises. Why seek advice? In some cases you may be
  required to be able to recover the name-data-relation, e.g. to enable study participants to "leave" a study.



### Data Management Checklist

https://fair-dom.org/knowledgehub/data-management-checklist/



### Data Management Checklist

### Making data findable (documentation and metadata management)

- What documentation and metadata will accompany the data (assist its discoverability)? (Details on methodology, definitions, procedures, SOPs, vocabularies, units, dependencies, etc)
- What information is needed for the data to be read and interpreted in the future?
- What naming conventions will be used?
- How will you approach versioning your data?
- How will you capture / create this documentation and metadata?
- How do you ensure the completeness of the captured data?

### Making Data Accessible

- Specify which data will be made openly available taking into consideration
  - What ethics and legal compliance issues do you have if any? Do you need consent for data preservation and sharing? Do you have to protect certain data? Is any data sensitive?
  - . Do you think you might have Intellectual Property Rights issues? Have you considered ownership of the data, licensing, restrictions on use?
  - · Do you think you will need to embargo any data?
- How will you make the data available? (consider the platforms you will use: databases, repositories, etc)
- What methods or software tools are needed to access the data? You should list where the software can be obtained. You should also document how to use the software to access the data. The documentation should be as complete as possible, including examples. If you distribute your system, include the access software and its documentation as part of any distribution.
- . If there are any restrictions on accessibility, how will you provide access?

### Making Data Interoperable



- How do you address data and model quality? What validation steps do you foresee?
- Will you use standardised vocabulary for all data types to allow inter-disciplinary interoperability?
- Where you can not used standardised vocabulary for all types of data, can you map to more commonly used ontologies?

### Making data Re-usable

- · How will you licence your data to permit the widest re-use possible?
- · When will the data be made available for re-use? Does this include an embargo period? (if yes, please detail why)
- Which data will be available for re-use during/after the project? For data that is not re-usable, please detail why
- What are your data quality assurance processes?
- How long do you expect your data to remain re-usable?



### **PERSONALIZZABILE**

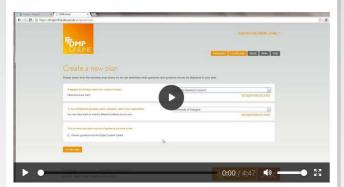




Future plans

Change language

### Screencast on how to use DMPonline





### Sign in

### Veteran tapes

expand all | collapse all

Write Plan

13/13 answered

### What data will you collect or create?

B I ≒ - ⊨ - ₽ ≡-

The "Veteran tape " project will collect and generate different types of datasets:

Type of data	Volume	Format	Storage format
Video recordings	600 x 1Gb	.mkv	.mkv
Transcriptions	600 x 1500Kb	MS Word	.txt
Structured interview text	1 x 500Kb	MS word	.txt

For the video recordings the selected format is .mkv; the same .mkv format will be used for the long-term preservation .

Transcriptions will be written in MS Word and then stored as .txt files.

We checked the format compatibility against EASY File format https://dans.knaw.nl/en/deposit/information-about-depositing-data/before-depositing/file-

As the total volume of data is greater than 50Gb, DANS requires a fee for the storage. We are currently in touch with EASY to determine the costs of archiving.

Guidance

### Guidance

Questions to consider:

- · What type, format and volume of data?
- Do your chosen formats and software enable sharing and long-term access to
- · Are there any existing data that you can

### Guidance:

Give a brief description of the data, including any existing data or third-party sources that will be used, in each case noting its content, type and coverage. Outline and justify your choice of format and consider the implications of data format and data volumes in terms of storage, backup and access.



# ... servono esempi [NON modelli]?...



https://dmptool.org/public plans

### Public DMPs

Public DMPs are plans created using the DMPTool service and shared publicly by their owners. They are not vetted for quality, completeness, or adherence to funder quidelines.

Q	Search

Project Title	Template \$	Organizatio
CAPERS - Cloud Applications and Platforms Enhancing Research Scalability	NSF-CISE: Computer and Information Science and Engineering	Georgia Tec
sample	NSF-SBE: Social, Behavioral, Economic Sciences	Binghamton University
Generic NSF DMP	NSF-GEN: Generic	Binghamton University
Sample NEH Plan	NEH-ODH: Office of Digital Humanities	Binghamton University
sample DMP Plan for workshop	NSF-CHE: Chemistry Division	Binghamton University
Sample NIH-GEN dmp	NIH-GEN: Generic	Binghamton University
As bibliotecas universitárias diante das mudanças contemporâneas	IMLS: Datasets	São Paulo S University (L
Gestão de acervos audiovisuais em repositórios	Digital Curation Centre	Non Partner Institution
Scientific Journals publishing	Digital Curation Centre	Non Partner Institution
The contribution of water retention, nutrient loading and microbial community to mosquito breeding and West Nile virus transmission in Spokane County	U.S. Geological Survey (USGS)	Non Partner Institution

# View all

### Horizon 2020

### AtlantOS Data Management Plan &

A 2016 DMP covering the concept of FAIR data for an AltaIntic Ocean observation systems project

Centre of Excellence in Simulation of Weather and Climate in Europe №(2016) DMP provided by the German Climate Centre (DKRZ)

### FREME Data Management Plan 🗗

Plan developed by the Open framework of e-services for multilingual and semantic enrichment of digital content project

### Helix Nebula open science cloud &

A High Energy Physics example covering the Large Hadron Collider experiments

### Open Sea Operating Experience to Reduce Wave Energy Costs @

First version of a DMP by the opera project

### Tweether &

An engineering based micro-electronics example

### AutoPost 🚱

An example DMP from an industry-driven innovation action that will deliver ICT-based solutions to enhance established post-production workflows.

### Evaluation rubric for assessing H2020 plans

ARESE FAMONT. GRECTHEWAY, respondence of the property of the p

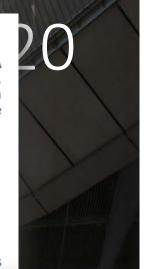
### Data Management Plan – general definition

Data Management Plans (DMPs) are a *key element* of good data management. A DMP describes the data management life cycle for the data to be collected, processed and/or generated by a Horizon 2020 project. As part of making research data findable, accessible, interoperable and re-usable (FAIR), a DMP should include information on:

- · the handling of research data during and after the end of the project
- what data will be collected, processed and/or generated
- which methodology and standards will be applied
- whether data will be shared/made open access and
- how data will be curated and preserved (including after the end of the project).

A DMP is required for all projects participating in the extended ORD pilot, unless they opt out of the ORD pilot. However, projects that opt

submit a DMP on a voluntary basis.



### H2020 Programme

Guidelines on
FAIR Data Management in Horizon 2020

**Guide 2016** 

Version 3.0

### 4. Research data management plans during the project life cycle

Once a project has had its funding approved and has started, you **must submit a first version of your DMP** (as a deliverable) within the first 6 months of the project. The Commission provides a DMP template in annex, the use of which is recommended but voluntary.

The DMP needs to be **updated** over the course of the project whenever significant changes arise, such as (but not limited to):

- new data
- changes in consortium policies (e.g. new innovation potential, decision to file for a patent)
- changes in consortium composition and external factors (e.g. new consortium members joining or old members leaving).

The DMP should be updated as a minimum in time with the periodic evaluation/assessment of the project. If there are no other periodic reviews foreseen within the grant agreement, then such an update needs to be made in time for the final review at the latest. Furthermore, the consortium can define a timetable for review in the DMP itself.

### Periodic reporting

For general information on periodic reporting please check the following sections of the online manual

- How to fill in reporting tables for publications, deliverables
- Process for continuous reporting in the grant management system.



# Open access to data

GRANT AGREEMENT ART. 29.3 pag. 248



### 3. Open access to research data (Extended Open Research Data Pilot)

### What?

Beneficiaries of actions that participate in the Open Research Data Pilot must give **open, free-of-charge access** to the end-user to **digital research data** generated during the action ( new in Horizon 2020).

As of the Work Programme 2017, the Open Research Data pilot has been extended to all thematic areas of Horizon 2020 (except ERC PoC actions, SME instrument Ph1 actions, ERA-NET Cofund actions that do not produce data, EJP Cofund actions, and prizes).

Participation is therefore now in principle **the default**. However, actions may **opt out** at any stage — both before signing the GA and afterwards (through an amendment; see Article 55) —, if:

- participation is incompatible with the obligation to protect results (see Article 27)
- participation is incompatible with the security obligations (see Article 37)
- participation is incompatible with rules on protection of personal data
- participation would mean that the project's main aim might not be achieved
- the project will not generate/collect any research data or
- there are other legitimate reasons not to take part.

- PILOT PROJECT EXTENDED TO ALL DISCIPLINES

OPT OUT CLAUSES

PRINCIPLE: «AS OPEN AS POSSIBLE, AS CLOSED AS NECESSARY»





### **ELIGIBLE COSTS:**

Rain au levain

- DATA CURATION
- DATA STORAGE
- DATA MANAGEMENT

Michel

### Open data — H2020

### How?

Open access to digital research data involves 3 steps:

### Procedure for open access (research data):

- Step 1 Deposit the digital research data, preferably in a research data repository.
- Step 2 Provide **open access** by taking measures to enable users to access, mine, exploit, reproduce and disseminate the data free of charge (e.g. for databases: by attaching an appropriate creative commons licence (CC-BY or CCO tool) to the data; if the access/use is not subject to any rights: by indicating that no licence is needed).

Open access must not be given immediately; for data needed to validate the results presented in scientific publications, as soon as possible; for other data, beneficiaries are free to specify embargo periods for their data in the data management plan (as appropriate in their scientific area).

Step 3 — Provide **information**, via the repository, about **tools and instruments** for validating the results.

Where possible, the beneficiaries should provide those tools and instruments (e.g. specialised software or software code, algorithms, analysis protocols, etc.).

- 1) DEPOSITARE IN UN DATA REPOSITORY
- RENDERE I DATI OPEN ACCESS CON LA LICENZA PIÙ APERTA POSSIBILE
   POSSIBILE EMBARGO
  - 4) FORNIRE ANCHE TUTTE LE INFORMAZIONI UTILI A VALIDARE I DATI



H2020 Programme

AGA - Annotated Model Grant Agreement

https://goo.gl/sryNTg

