You are free to

▷ share, adapt or re-mix
▷ photograph, video or broadcast
▷ blog, live-blog or post-video

this presentation

Provided that

you attribute the work to its author and respect the rights and licenses associated with its components

# Research data are first-class citizens in science



advertising + text + **data** + code + version → science

*reproducibility spectrum*

0%                         100%

but who does really care about these data, to begin with?

# The short answer: everybody interested in (modern) science



Researchers AND machines need to find/discover data having features of interest, for which they will be using links, metadata, as well as actual data



Once found, machines need to access/retrieve data of interest (obtain a copy of the contents in some format)



Finally, researchers go to their computers and start to re-use/analyze the data; they need to have workflow tools to:
-aggregate data and run analyses,
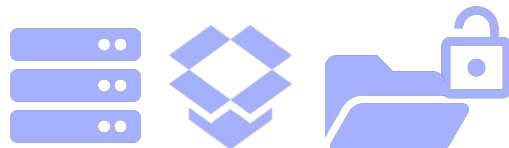-reproduce existing pipelines,
-validate sanity of data,
-….

# In the eScience ecosystem, enabling optimal use of research data and methods is a big challenge

**Researchers** ready to share their data and interpretations

**Professional data publishers** (data repositories, data journals)

**Data science community** analysing data to advance discovery

**Funding agencies** increasingly concerned about data stewardship

Providing **machine-readable data as the main substrate for Knowledge Discovery** is a big challenge in modern Science

The **best practice recommendation for research data** is to be as open and FAIR as possible
(while accounting for ethical, commercial and privacy constraints with sensitive data or proprietary data)

Force11 FAIR principles

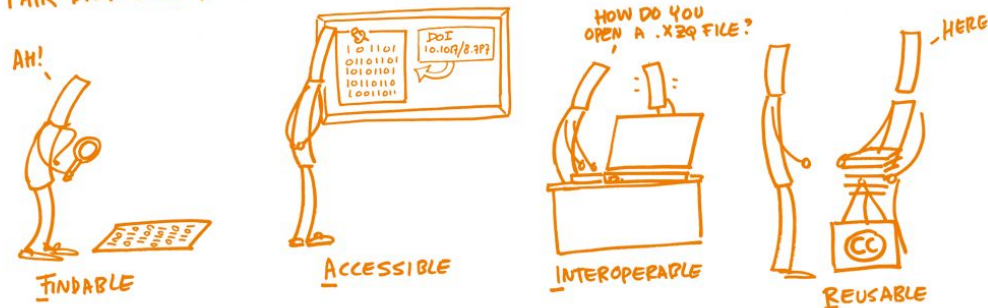# Research objects: scholarly information on the Web



The basic idea is that objects that belong together (*e.g.*, an article with its associated code, data and workflows) should have some means of being aggregated, so that all associated research objects can be discovered together.

Although this might seem to be obvious, as research objects are scattered across different repositories on the web, the connections between them are often lost.

# The FAIR principles provide guidance for scientific data management and stewardship



1) FAIR *does not* imply Open (but FAIR data can definitely be Open Data; the two concepts can overlap)
2) the FAIR facets are related to each other but technically somehow independent of each other
   (we can speak of degrees of FAIR-ness)
3) the FAIR principles are agnostic of technical implementations
   (FAIR-ness can be achieved with a wide range of technologies and implementations)

the principles describe characteristics that data should exhibit to **assist discovery and reuse through the web**

Findable: data and supplementary materials have sufficiently rich metadata and a unique and persistent identifier

Accessible: metadata and data are understandable to humans and machines and deposited in a trusted repository

Interoperable: metadata use a formal, accessible, shared, and broadly applicable language for knowledge representation

Reusable: data and collections have a clear usage licenses and provide accurate information on provenance

Open Science Training Handbook, Foster; Force11 FAIR principles

# Findable

| F | 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 |

FAIR DATA PRINCIPLES

AH!

FINDABLE

# Findable

| F | 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 |

DOI [ Digital Object Identifier ]
the most widely used PID for research data



A URL is not a PID

the dataset needs to be identified by a Persistent Identifier [ PID ]
so that it can be located by a machine

PID = a **globally unique** and **long-lasting** reference to something, *e.g.*, documents, files, books, people

A PID is **separated from location**: if a web document is moved, the PID points to the same object in the new location
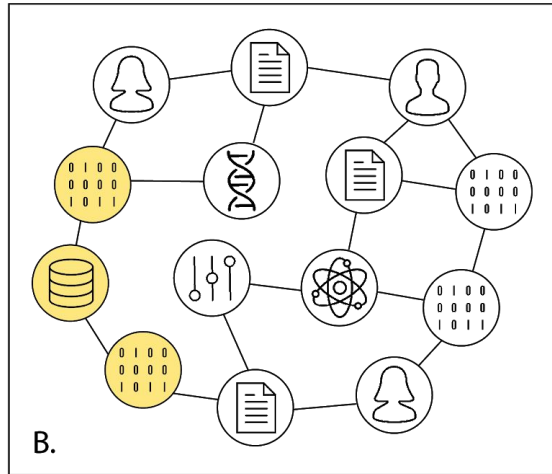
**Best practice**
1. Deposit your data to a domain-specific repository
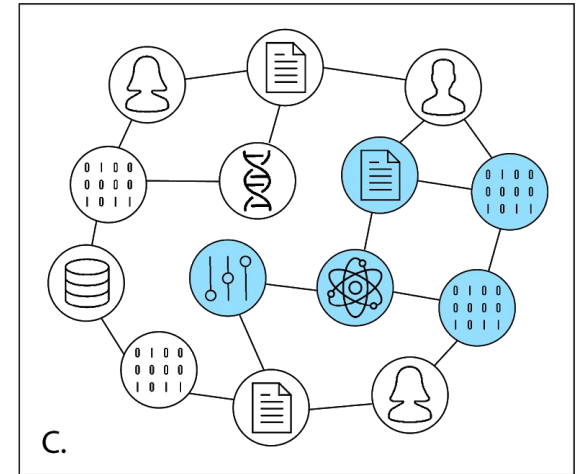2. If that does not exist (yet), use general-purpose repositories

# Towards a PID graph to enable discovery and impact assessment



A.

**Different versions of software code**

B.

**Datasets hosted by a particular repository**

C.

**All digital objects connected to a research object**

Data Sharing or Data Visitation?

Findable where?

https://fairsharing.org/databases/, https://www.re3data.org/

# Findable

| F | 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 |

the data are described by rich metadata, so they can be discovered by a human

*Your first collaborator is yourself, and your past self does not answer emails.*

Intrinsic metadata: immutable
the author of a book
the date a photo was taken

Extrinsic metadata: depend on the context
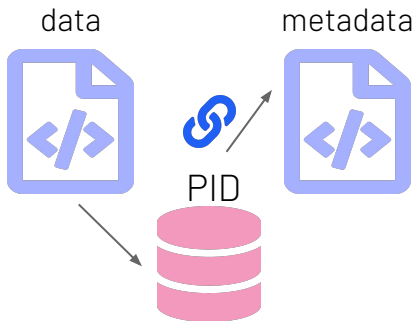the date a book was purchased
the publications a photo has appeared in

**Best practice**
1. Never presume that you know who will want to use your data, or for what purpose: be generous!
2. Whenever possible, use community standards
3. Have someone "naive" check your annotations

RDM promotion resources, list of metadata standards, FAIRsharing

# Findable

| F | 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 |

data          metadata

PID

, list of metadata standards, FAIRsharing

the data are described by rich metadata, so they can be discovered by a human

*Your first collaborator is yourself, and your past self does not answer emails.*

Intrinsic metadata: immutable
the author of a book
the date a photo was taken

Extrinsic metadata: depend on the context
the date a book was purchased
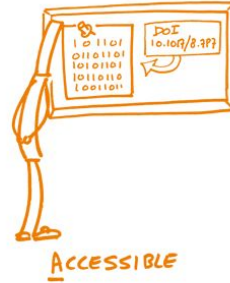the publications a photo has appeared in

**Best practice**
1. Never presume that you know who will want to use your data, or for what purpose: be generous!
2. Whenever possible, use community standards
3. Have someone "naive" check your annotations
4. IF you do F1 correctly, F3 comes for free!

# Accessible

| A | 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 |
| A4 | metadata are accessible, even when the data are no longer available |



ACCESSIBLE

# Accessible

| A | 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** |
| A4 | metadata are accessible, even when the data are no longer available |

limitations on and protocols for the use of data are made explicit

Data should be retrievable by anyone with a computer and an internet connection, if they are authorized, with a well-defined protocol.

Accessible data does not automatically imply open or free access: data published with restricted access can also be FAIR.

**Best practice**

1. Include **clear** licenses and conditions of use: who can access the data and in what way?
2. Remember: no license does not mean OK to access (when in doubt, the answer is NO)

# Interoperable

| I | 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 |



HOW DO YOU OPEN A .XZQ FILE?

INTEROPERABLE

# Interoperable

**Clearly, this is a very challenging requirement to meet.**

**Best practice**

1. Document as much as possible - be generous!
2. Use *preferred file formats*
3. If existing in your discipline, use standard vocabularies, ontologies and thesauri in your (meta)data, or provide mapping of your data to these vocabularies, ontologies and thesauri

| I | 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 |

| Type | Preferred format(s) | Non-preferred format(s) |
|---|---|---|
| Text documents | • PDF/A (.pdf) | • ODT (.odt)<br>• MS Word (.doc, .docx)<br>• RTF (.rtf)<br>• PDF (.pdf) |
| Plain text | • Unicode text (.txt) | • Non-Unicode text (.txt) |
| Markup language | • XML (.xml)<br>• HTML (.html)<br>• Related files: .css, .xslt, .js, .es | • SGML (.sgml) |
| Spreadsheets | • ODS (.ods)<br>• CSV (.csv) | • MS Excel (.xls, .xlsx)<br>• PDF/A (.pdf)<br>• OOXML (.docx, .docm) |

https://en.wikibooks.org/wiki/Choosing_The_Right_File_Format/Quick_Guide

lymphocyte of B lineage
CL_0000945

is a

B cell
CL_0000236

is a

precursor B cell
CL_0000817

# Metadata, Ontologies, and Vocabularies Standards help Achieve Interoperability Principles
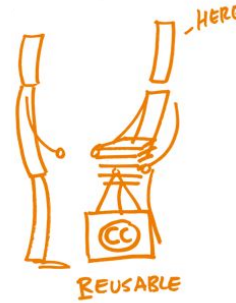
# Reusable

| R | 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 |

# Reusable

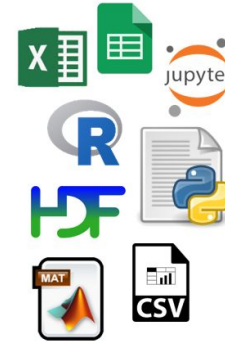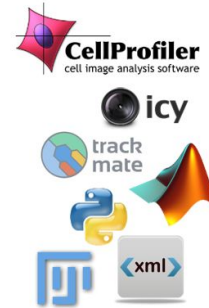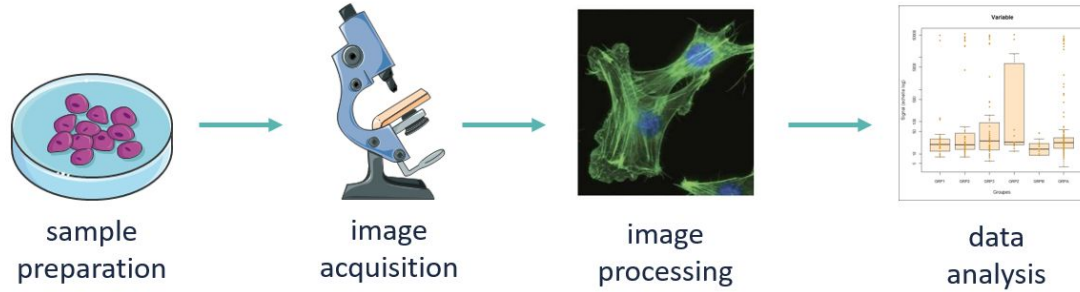| R | Reusable |
|---|---|
| R1 | meta(data) are richly described with a plurality of accurate and relevant attributes |
| **R1.1** | **(meta)data are released with a clear and accessible data usage license** |
| R1.2 | (meta)data are associated with detailed provenance |
| R1.3 | (meta)data meet domain-relevant community standards |

Reusers need clear signals from researchers on what they can and cannot do with their research data.

**Best practice**

1. Check for funder/data repository/local policy/institution obligations
2. If your obligations are non-exclusive, consider multiple-licensing (different versions)
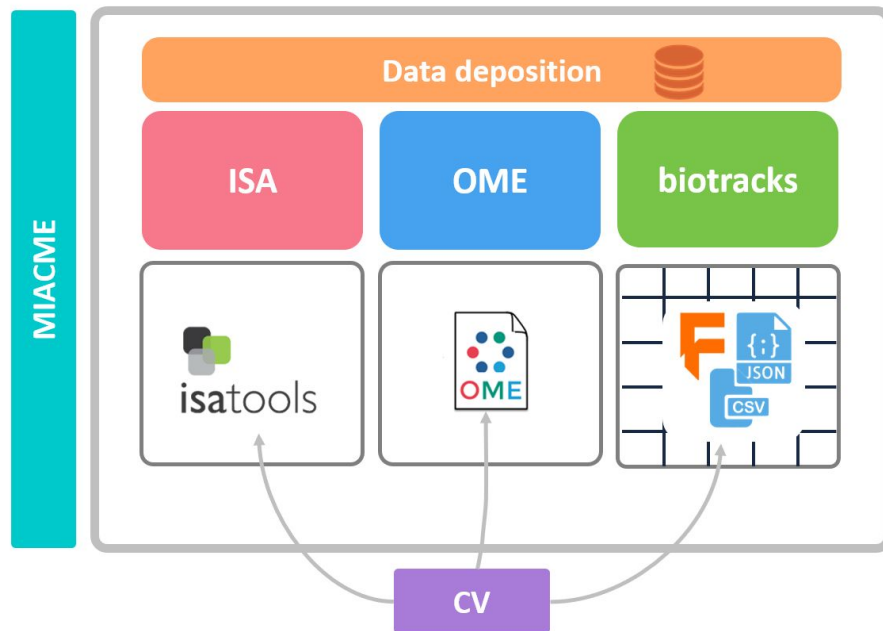2. Avoid bespoke licenses

# My story: towards FAIR cell migration data

# My story: towards FAIR cell migration data

MIACME: Minimum Information About a Cell Migration Experiment

ISA: Investigation Study Assay

OME: Open Microscopy Environment

biotracks: specialises the Tabular Data Package container format



Gonzales, Masuzzo et al., 2019; Frictionless Data

# A FAIR cell migration dataset

Create new file | Upload files | Find file | History

gsergeant added underscore to differentiate between column and row if both woul... ...     Latest commit ced8b3b on Jun 21, 2018

..

| | | |
|---|---|---|
| 📁 isa | Fixed code to load ISA-Tab datasets | a year ago |
| 📁 miacme | Renamed folder to include miacme files in cmsodataset0001-masuzzo folder | 2 years ago |
| 📁 trackmate/2_B/dp | added underscore to differentiate between column and row if both woul... | 9 months ago |
| 📄 9I5TT808.companion.ome | Rename CMSO d ataset 1 | 2 years ago |
| 📄 README.md | Fixes to the investigation files and changes in the README. | 2 years ago |

📖 README.md

Data from Masuzzo et al. 2017.

The above publication describes two experiments. Metadata included here refers to the first experiment on Ba/F3 cells. Image data (not included here) consists of 12 144-frame TIFF movies, each relative to a well in a 48-well plate.

Currently, we have:

- An OME-TIFF companion XML file that groups the 12 image data files, including information on the plate layout
- An ISA-tab dataset containing the experimental metadata, compliant with the MIACME minimum information guideline
- A biotracks datapackage containing TrackMate data

CMSO dataset

# A FAIR cell migration dataset

CMSO-datasets / cmsodataset0001-masuzzo /

| | Create new file | Upload files | Find file | History |

gsergeant added underscore to differentiate between column and row if both woul... ...    Latest commit ced8b3b on Jun 21, 2018

..

| 📁 isa | Fixed code to load ISA-Tab datasets | a year ago |
| 📁 miacme | Renamed folder to include miacme files in cmsodataset0001-masuzzo folder | 2 years ago |
| 📁 trackmate/2_B/dp | added underscore to differentiate between column and row if both woul... | 9 months ago |
| 📄 9I5TT808.companion.ome | Rename CMSO d ataset 1 | 2 years ago |
| 📄 README.md | Fixes to the investigation files and changes in the README. | 2 years ago |

📖 README.md ✏️

Data from Masuzzo et al. 2017.

The above publication describes two experiments. Metadata included here refers to the first experiment on Ba/F3 cells.
Image data (not included here) consists of 12 144-frame TIFF movies, each relative to a well in a 48-well plate.

Currently, we have:

- An OME-TIFF companion XML file that groups the 12 image data files, including information on the plate layout
- An ISA-tab dataset containing the experimental metadata, compliant with the MIACME minimum information guideline
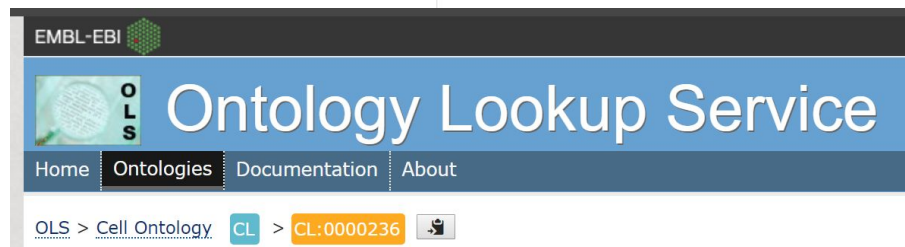- A biotracks datapackage containing TrackMate data

Annotate,
annotate,
annotate!

# A FAIR cell migration dataset

| MIACME | v0.3 |
|---|---|
| Experimental setup | |
| Cell type | Cell line [CLO_0000001] |
| Cell details | B-cell [CL:0000236] |
| | Ba/F3 cell [CLO:0001842] |
| Treatment (Experimental variable) | Bcr-Abl [PR:000044437] oncogene variants |
| Treatment (Experimental variable) | y-27632 [CHEBI:75393] (Rock inhibitor at 10 μM) |

Use standard ontologies and controlled vocabularies!

EMBL-EBI

## Ontology Lookup Service

Home    Ontologies    Documentation    About

OLS > Cell Ontology    CL  >  CL:0000236

Y-27632 Source: ChEBI (ID: 75393)

A monocarboxylic acid amide that is trans-[(1R)-1-aminoethyl]cyclohexanecarboxamide in which one of the nitrogens of the aminocarbony group is substituted by a pyridine nucleus. It has been shown to e ...

# A FAIR cell migration dataset

```xml
<Image ID="Image:2" Name="9I5TT808_F00000012.tif">
    <AcquisitionDate>2013-03-25T22:05:53</AcquisitionDate>
    <Description>Camera 1376x1038, Controler present, Lightsource off</Description>
    <Pixels BigEndian="false" DimensionOrder="XYCZT" ID="Pixels:0" Interleaved="false" PhysicalSizeX="0.322000000022352" PhysicalSizeXUnit="µm" PhysicalSizeY="0.322000000022352" PhysicalSizeYUnit="µm" SignificantBits="16" SizeC="1" SizeT="144" SizeX="1376" SizeY="1038" SizeZ="1" Type="uint16">
        <Channel ID="Channel:0:0" SamplesPerPixel="1">
            <LightPath/>
        </Channel>
        <TiffData FirstC="0" FirstT="0" FirstZ="0" IFD="0" PlaneCount="1">
            <UUID FileName="9I5TT808_F00000012.tif">urn:uuid:90e9dc02-94db-4e32-9759-2780958e6c6d</UUID>
        </TiffData>
```

Expose as rich metadata as possible!

## Note

Some data elements, for instance images and 'raw data' can not always be made machine-processable.

Being published with FAIR metadata is of very high value in its own right.

# Different levels of FAIR-ness can exist

**Findable** ⓘ

Does the dataset have any identifiers assigned?
> Web address (URL) ▾

Is the dataset identifier included in all metadata records/files describing the data?
> Yes ▾

How is the data described with metadata?
> Comprehensively, but in a text-based, non-standard format. ▾

What type of repository or registry is the metadata record in?
> Data is in one place but discoverable through several registries ▾

A little step can make a big difference

**Findable** ⓘ

Does the dataset have any identifiers assigned?
> Globally Unique, citable and persistent (e.g. DOI, PURL, ARK c ▾

Is the dataset identifier included in all metadata records/files describing the data?
> Yes ▾

How is the data described with metadata?
> Comprehensively (see suggestion) using a recognised formal m ▾

What type of repository or registry is the metadata record in?
> Data is in one place but discoverable through several registries ▾

https://www.ands-nectar-rds.org.au/fair-tool

# Ethical considerations and legal obligations guide the way

With appropriate **data management planning** much sensitive and proprietary data can be shared, reused, and FAIR.

The metadata can almost always be shared.

Guidance and best practices for sharing sensitive data are necessarily region-specific because of local regulations.
As general guidance, the EU RESPECT code for Ethics and data protection highlights 3 key aspects:

1) **Upholding scientific standards**
   when formulating your research questions, you do not pre-determine or prejudice the outcome through your choice of questions or actions

2) **Compliance with the law**
   be aware of all the relevant national and international laws that may affect your research project. With collaborative projects which cross legal borders, this may involve various laws. Ones of particular relevance will be in regards to data protection and intellectual property.

3) **Avoidance of social and personal harm**
   your research project should be designed responsibly and should consider participants throughout. For example, participation in the research project should be voluntary and on the basis of fully informed consent

Ethics and Data Protection

*A goal without a plan is just a wish*

[Antoine de Saint-Exupéry (1900 –1944)]

www.digitalbevaring.dk

# What is a Data Management Plan (DMP)?

A DMP is a term describing how you organize, structure, store and care for the information used in a research project

**How do you look at the data on a day-by-day base?**
Organize your data, store them and back them up.
Choose the right file formats.
Document your data.

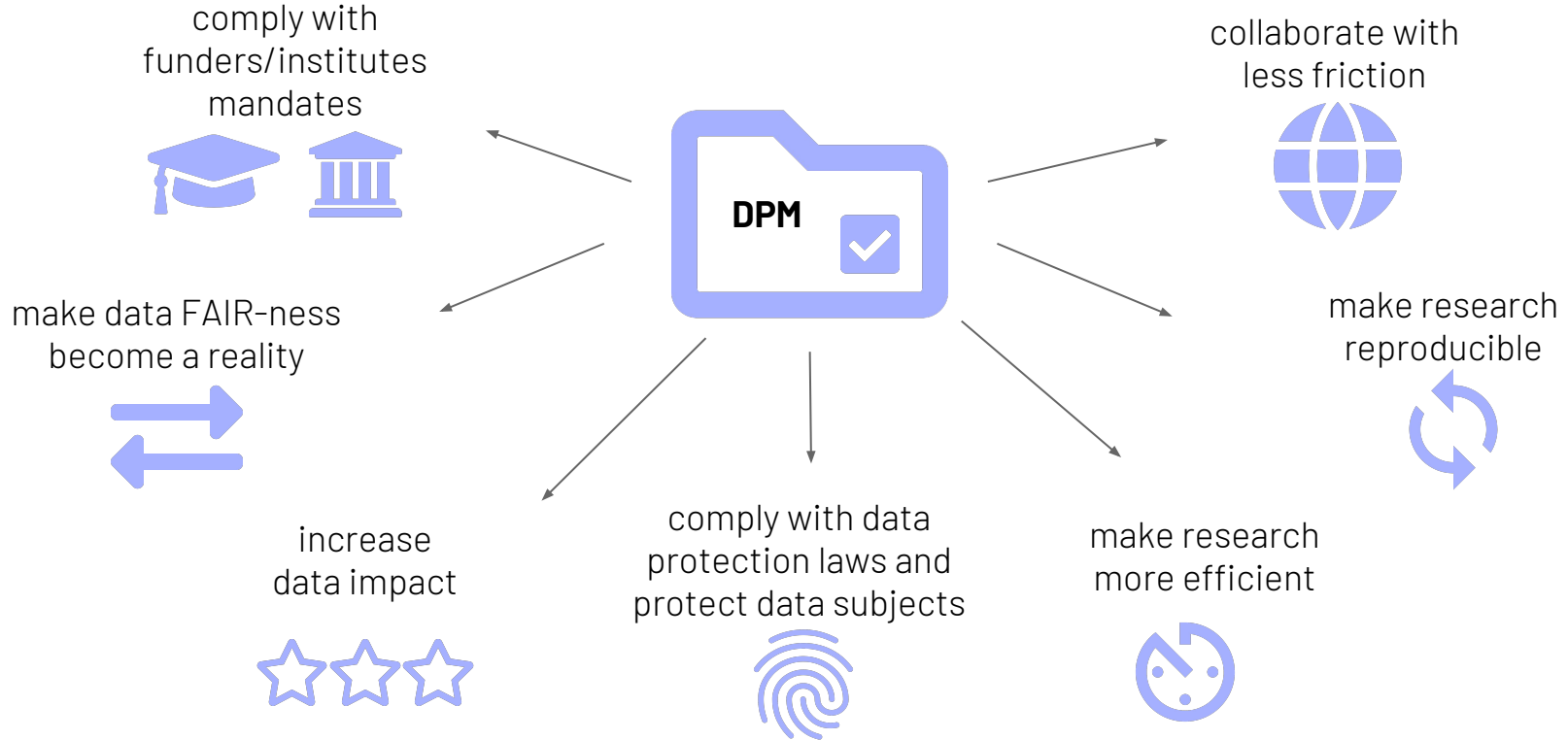**What happens to the data after the research project finishes?**
What data do you need to keep (and share)?
What data must not be kept (and shared)?
How are you going to achieve long-term data storage/access?

# Why should you develop a DMP?



comply with funders/institutes mandates

collaborate with less friction

make data FAIR-ness become a reality

DPM

make research reproducible

increase data impact

comply with data protection laws and protect data subjects

make research more efficient

# How do you create a DMP?

My personal message is: don't be afraid, go to your desk and try!

[DPMOnline](#) developed by the DCC has free-to-download [public templates](#) + your university might have a local version of the tool, which means you can use your institutional account to sign in and the templates available will make sure you don't miss out on relevant info/fields.

| |
|---|
| Data Collection (0 / 2)     ✚ |
| Documentation and Metadata (0 / 1)     ✚ |
| Ethics and Legal Compliance (0 / 2)     ✚ |
| Storage and Backup (0 / 2)     ✚ |
| Selection and Preservation (0 / 2)     ✚ |
| Data Sharing (0 / 2)     ✚ |
| Responsibilities and Resources (0 / 2)     ✚ |

[Data Management Checklist](#), [https://dmponline.dcc.ac.uk](https://dmponline.dcc.ac.uk)

# Where are you going to store your data?







central file shares managed by the
ICT Department of your University

live snapshots, backups!
**prevent disasters!**

# How are you going to organize your data?



use directories and folders hierarchy

name files with a bit of common sense: no special characters, no capitals, no spaces

**know what your file is before you double click it!**
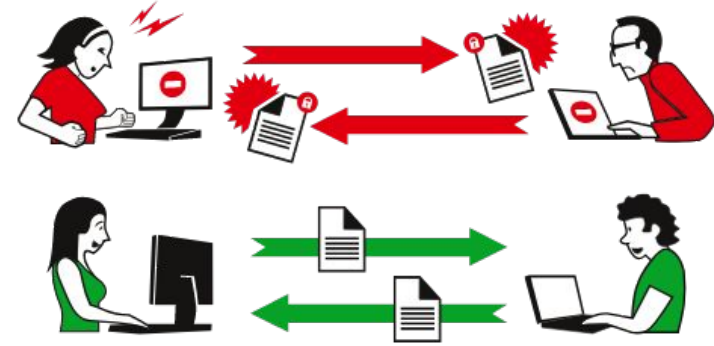
# What formats and what software will you use?

During your research you will most likely use formats that fit your scope, workflow and methodology

After research is completed, you should look for formats that are easy to share and re-use and that may last longer into the future
- open specifications
- widely used formats
- uncompressed
- ASCII formats
- exchange formats

tip: check if software you are using has an option to export into a more suitable format for sharing or long-term reuse

(sometimes 'crazy' formats are simply ASCII-based, and can be saved and then opened as regular .txt files!)

# THANKS!

You can find me at:

@pcmasuzzo

paola.masuzzo@gmail.com