

### **Research Data Management - AVRE Training**

# UNIVERSITÉ DE MONS



#### Sébastien Hoyas

Data Project Manager & Business Analyst

Place du Parc, 22 Office 1.12 7000 Mons +32 65 37 34 54

sebastien.hoyas@umons.ac.be



### Outline

- Scope of this training session
- What are (research) data?
- What is data management?
- Data management advices
- Data management plan (DMP)
- **Closing remarks**
- Q&A

### Typical journeys in research... 🕲

### Start of a project

### Typical journeys in research... ©



### Typical journeys in research... ©



### Typical journeys in research... 😊



All these problems are related to *data* and their *management* 

#### Scope of this training:

To try help you manage your data Make you aware about resources to help you

#### What I will <u>not</u> be able to do:

Provide a unique answer to manage your data

### What are (research) data?

**Research data**: information that is <u>collected</u>, <u>observed</u>, or <u>created</u> for purposes of analysis to produce original research results:

- Survey answers
- Code
- Measurements
- Samples
- ...

**Quantitative**: numerical (number of heartbeats per minute)

VS.

**Qualitative**: textual, visual (observations of a patient)

**Different formats**: text, numerical, image, audio, video, physical, etc.

### What is data management?

Data management is the **handling** of research data **during** <u>AND</u> after a research activity:

- Collection
- Documentation
- Organization
- Storage
- Sharing

Good data management helps to ensure that researchers share data in a FAIR way (see later)

Research organizations/funders increasingly require their researchers to *plan* how they will manage their data to ensure that all aspects are considered from the start and they can get a return on investment (public goods)

**1.Efficiency**: Proper data management can make your research process more efficient, saving you time and effort in the long run.

- **2.Reproducibility**: Crucial for reproducibility of research. It ensures that you and others can understand and replicate your work in the future.
- **3.Data Integrity**: It helps maintain the integrity of your data and reduces the risk of data loss.
- **4.Collaboration**: It makes it easier to share and collaborate with others, both within and outside your research team by setting proper licences.
- 5.Compliance: Many funding agencies require DMP and data sharing (see later).
- **6.Preservation**: It ensures your data is preserved for future use and reuse.
- **7.Return on investment**: Making data discoverable, accessible, and reusable maximizes the research potential of the data and provides greater returns on public investments in research.

### Why is it so important?

Most scientific results are difficult, even **impossible**, to reproduce and/or replicate

Research integrity is not favoured (publish or perish)

Avoid reiventing the wheel





Many **different actors** in a research project with **different needs**:

- Primary Researcher or Principal Investigator: Creates and uses data
- Institution: Sets internal data management policy (including Data Ambassadors, Promoters)
- Data Repository: Curates and provides access to data
- User: Uses 3rd party data
- Funder: Provides the resources to support a research project
- **Publisher**: Disseminates discoveries and maintains the scientific records

Data management is not easy (who said research is easy?) and heavily **depends on** *you*, *your project* and *your peers* 

To help you, general advices in:

- Data organization
- Data description & documentation
- Data storage
- Data sharing

All these « categories » are linked together

It is better to think in advance about *how* you will organize your data, but you can also change your mind during your research as long as you keep track of those changes.

Take into consideration that you may not be the only one working on those data, so **make it clear for** *anyone*!

In general, most research data are digitalized, but it also applies to physical data (samples, etc.).

### <u>Tips & tricks</u>

**Existing procedures**: check if there are already established ways to organize your data in your team (and check if it suits your needs!).

**File organization**: there is no universal answer  $\circledast$  You need to apply an organization that is compatible with you, your project and your team.

Above all, you must be **CONSISTENT** accross your project.

#### File organization: examples



#### File organization: examples





**Folder names:** keep folder names short (max 15-20 characters) and make them descriptive of what is inside, without being redundant with the folder structure.

• Bad example





#### NOTES

- Avoid using spaces, dots and special characters (&, ?, etc.)
- Use hyphens (-) or underscores (\_) to separate elements
- Use a minimum of two leading zeros for padding (001, 002 ,etc.) to properly sort folders by names

**File names:** keep file names brief and explicit, without being redundant with the folder structure.

- Avoid using spaces, dots and special characters (&, ?, ;, etc.)
- Use hyphens (-) or underscores (\_) to separate elements (easier to recover in the OS)
- Use a minimum of two leading zeros for padding (001, 002, etc.) to properly sort files by names
- Use an extension that matches the file format
- If your files cannot be integrated in a versioning tool like <u>Git</u>, include a version number at the end. Keep a logfile where you briefly state changes in each of the new versions
- Include elements such as the date (YYYYMMDD format, best to sort) at the beginning of the file name
- Avoid starting to name files with « draft », « final » or the version number

#### Example: 20230130\_RDM\_Training\_V001.pptx

**Versioning:** keep older versions in a separate folder, and do not delete them unless you are absolutely sure you can. Keep a logfile that briefly explains the changes in each version.

*Key considerations*: there are bad practices, but there are *no unique answer* or best method to organize data. You must be *consistent* throughout your project so that you and your team can work on.



### **Data description & documentation**

### Data description and documentation

Once you know how you will organize your data, you can start <u>collecting</u> them. You should also start describing them using **metadata** and **documentation**.

- **Metadata** are data that describe data... ③ (make data findable)
- Documentation of your data should include the method(s) you used to obtain them, how they were analyzed, processed, where you can find them, etc.

Those information should be stored in distinct files in the relevant folders.

## **Data description**

### Data description and documentation

### Standard metadata $\rightarrow$ depends on your discipline/data type/purpose

Metadata are important to **find** the data, ensure **reproducibility** and **reuse** 

#### **Common elements**

- **Title**: The name given to the dataset.
- Author: The main researchers involved in producing the data.
- Date of creation
- Identifier: A unique code assigned to the dataset. May be added later, when *sharing* the dataset.

#### Standards commonly used:

- <u>Dublin Core</u>: 15 properties for describing a wide range of resources (general purpose) Check <u>this website</u> to generate your metadata file !
- Digital Curation Center: social sciences
- Biology, earth sciences, physical sciences

### **Data documentation**

### Data description and documentation

**Data documentation:** document that explains your data. It will help others that would use your data, but also yourself to remember how you obtained and processed those data (*try to do it on the fly*). It also helps you and others to reproduce your results.

- Folder organization: you should explain how data are organized so that anyone starting to collaborate with you can understand what is going on.
- Data collection: explain how you obtained the data (from known datasets or from an experiment, survey, simulations, etc.).
- **Data cleaning**: explain your investigation and why you removed part of them (errors, inaccuracies, etc.).
- **Data analysis:** how you analyzed the cleaned data (which software, parameters, results of the analysis, etc.).
- **Plan for Change**: keep in mind that your data may change over the course of your project. Plan for how you will document and manage these changes.
- **Consider Your Audience**: remember that the description of your data may be read by people who are not experts in your field. Try to write in a way that is accessible to non-experts.

### Data description and documentation

*Key considerations*: Metadata and data documentation files should *always accompany your data* (same folder or dedicated subfolder and the dataset). The more descriptive the better for future reuse and reproducibility.

Data Storage

- Reliability: Use <u>reliable storage solutions</u> that ensure data integrity and availability NOT YOUR LOCAL LAPTOP/DESKTOP!
- Accessibility: Ensure that data is easily accessible to authorized users.
- Scalability: Choose storage solutions that can grow with the size of your data.

#### **NOTE: Backup is different than preservation**

- **Backup** = periodic snapshots in case current version is lost or destroyed (*cloud, NAS, etc.*)
- Preservation = archival, usually the final version of a dataset, stored for long-term and further use (*data repositories*)

Backup:

- Frequency: Regularly back up data according to the importance and frequency of change.
- <u>3-2-1 Rule</u>: Keep at least three copies of your data, on two different media, with one backup offsite<sup>1</sup>.
- **Verification**: Regularly verify the integrity of backup copies.
- **Disaster Recovery**: Have a disaster recovery plan in place to restore data if needed.

Nice **free open source** software for Linux, Windows and MacOS: <u>FreeFileSync: Open Source</u> <u>File Synchronization & Backup Software</u>

#### Backup: where?

- **Cloud**: UMONS provides 1To/user or team (see <u>sharepoint</u>) for free
- Cloud: if you need extra space or you want to backup somewhere else, you must pay for another service (not offered by UMONS)
- Local server: you can setup a Network Attached Server (<u>NAS</u>) for you and your team with backup services (not offered by UMONS)

#### Long-term preservation: *data repositories*

A trusted digital repository provides reliable long-term access to managed digital resources to its designated community, now and in the future!

Only completed datasets with the purpose to *publish*, *share* and/or *preserve* them should be uploaded (not all research data).

Typical files included in a dataset repository: inputs, outputs, method and metadata, <u>not all</u> <u>intermediary data</u>! With only those information, anyone should be able to obtain the same outputs as you.

#### Examples:

- <u>Zenodo</u>: general purpose repository
- <u>SODHA</u>: the federal Belgian data archive for social sciences and the digital humanities

#### **Dataverse:** our institutional repository under development to host your clean datasets



Data repositories also allow you to *share* your data.

Your data should meet the FAIR data principles:

- Findable → Metadata
  - Data should be easy to find for both humans and computers. Machine-readable metadata are essential for automatic discovery of datasets and services.

#### Accessible $\rightarrow$ Data repository and unique identifier (DOI

• Once the user finds the required data, they need to know how they can obtain them. The data might be publicly accessible, or access may be restricted but metadata should remain accessible.

#### Interoperable → Metadata, documentation

 The data usually need to be integrated with other data regardless of the systems or tools being used → non-proprietary data formats (.txt, .csv, .md, .pdf, etc.)

#### **Reusable** → **Documentation**

• The ultimate goal of FAIR is to optimize the reuse of data. To achieve this, metadata and data should be well-described so that they can be replicated and/or combined in different settings.



*Key take-away:* To effectively share data, resolve any data ownership or intellectual property rights issues early. Consult your institution (<u>AVRE</u>) to determine what policies might affect data ownership and sharing.

- **Understand Your Goals**: What do you want others to be able to do with your data? This can help guide your choice of license.
- **Public Domain**: If you want to give people the most freedom, consider a public domain license like CC0.
- Attribution: If you want to allow free use but also want to be credited, consider a license that requires attribution, like CC-BY.
- Share-Alike: If you want any derivatives of your work to be licensed under the same terms, consider a share-alike license, like CC-BY-SA.
- Non-Commercial: If you want to restrict the commercial use of your data, consider a non-commercial license, like CC-BY-NC.
- **Understand the Implications**: Make sure you understand the implications of the license you choose. Some licenses may have implications for how your data can be used or shared.
- Seek Legal Advice: If you're unsure, consider seeking legal advice. Licensing can be complex, and it's important to get it right.

### Data sharing and licenses

#### **Common licenses for datasets**

- <u>CC0 Public Domain Dedication</u>
- Open Database License
- Open Licence Etalab
- Open Data Commons Public Domain Dedication and License
- Open Data Commons Attribution License

In more details: <u>SPDX License List | Software Package Data Exchange (SPDX)</u>

### **Budget for storage**

Data production can be costly (equipment, products, time).

In addition to the costs of data collection, data management, curation, documentation, storage can be expensive.

 $\rightarrow$  Need to think about all costs (IT, server, etc.).

At the moment, 1 TB/person for free with Microsoft OneDrive, but may not be enough or practical for your project.

These costs may be eligible in some calls for projects.

Examples:

- Cloud storage on Microsoft Azure: 250 GB/month → min. 50€ → 4 years = 2400€ only for your research
- Network attached storage: fixed cost <u>NAS</u> 500€ + UPS 200€ + 4TB disks 120€ \* 3 = 1060€
  → multiple users, research, etc. (*no support from the university yet...*)

### **Data Management**

Common element in previous advices about data management: *planning* 

To help you managing your data, it is encourage (and even more often than before required) to create a *Data Management Plan* (DMP).

**DMP**: Document that describes the data produced/used in all stages of a research project and outlines the strategy to manage your data, before and during the project.

This document is *not static* and can *evolve* during your research (research can be unpredictive 😳).

### Data lifecycle in research

To create a DMP, you need to understand the various stages data goes through during a research project == "*data lifecycle*".



Common elements:

- Planning
- Collecting
- Analyzing
- Organizing
- Preserving data for future use and share

More and more funding agencies require to submit a **DMP** to ensure that data will be accessible and usable on the long run  $\rightarrow$  return on investment

Better to write a DMP **BEFORE** starting a project, but also beneficial to do it **NOW**.

#### Two helpful and free DMP tools:

- <u>DMPTool</u>, created by the University of California Curation Center of the California Digital Library
- <u>DMPOnline</u>, developed by the UK Digital Curation Centre  $\rightarrow$  <u>Mandatory for European</u> <u>Projects!</u>
- Both of these tools provide guidances and templates for creating DMPs in compliance with institutional and funder requirements.

### Data management plan

#### Typical questions in a DMP:

- What data will you collect or create?
- How will the data be collected or created?
- What documentation and metadata will accompany the data?
- How will you manage any ethical issues?
- How will you manage copyright and intellectual property rights issues?
- How will the data be stored and backed up during research?
- How will you manage access and security?
- Which data should be retained, shared, and/or preserved?
- What is the long-term preservation plan for the dataset?
- How will you share the data?
- Are any restrictions on data sharing required?
- Who will be responsible for data management?
- What resources will you require to implement your plan?

Live demo: <u>Welcome to DMPonline.be</u>

#### Log in using your institutional credentials



#### Welcome to DMPonline.be

We can help you write and maintain data management plans for your research.

This instance of DMPonline is provided by the DMPbelgium Consortium, which was founded in 2017 by:

- Instituut voor Natuur- en Bosonderzoek
- Université Libre de Bruxelles
- Universiteit Antwerpen
- Universiteit Gent
- Universiteit Hasselt
- Vrije Universiteit Brussel
- Wetenschappelijk Instituut Volksgezondheid Institut Scientifique de Santé Publique (Sciensano)

In 2018 they were joined by:

- Université Catholique de Louvain
- Université de Liège
- Université de Mons
- Université de Namur
- Vlaamse Instelling voor Technologisch Onderzoek

Since then, the Consortium has been joined by:

- Arteveldehogeschool
- Instituut voor Landbouw-, Visserij- en Voedingsonderzoek
- Universitair Ziekenhuis Gent
- Vlaams Instituut voor de Zee
- Vlerick Business School
- Hogeschool Gent

#### Interested in joining the Consortium?

#### Sign in with your institutional account

Flanders Make (Belgium) (flandersmake.be)

Flanders Marine Institute

Flemish Institute for Technological Research

Ghent University

Ghent University (UZGent)

Hasselt University

Hogeschool Gent (HOGENT)

Hogeschool VIVES (VIVES)

IMEC

KU Leuven (KUL)

Research Institute for Agriculture, Fisheries and Food (ILVO)

Royal Institute for Cultural Heritage (kikirpa.be)

Royal Library of Belgium (KBR)

Royal Observatory of Belgium

Sciensano

Thomas More Hogeschool

Université catholique de Louvain (UCLouvain)

Université de Liège

Université de Mons

Université de Namur

#### Create a new plan

Before you get started, we need some information about your research project to set you up with the best DMP template for your needs.



					Request feedback	Download
* Project title						
test					Select Gu	uidance
mock project for	testing, practic	ce, or educational purpo	oses		To help you write yo show you guidance organisations.	our plan, DMPonline.be o from a variety of
	· 8 = ·				Select up to 6 orga guidance.	nisations to see their
			ju ka		and guida, from below See the full list Save	additional organisations
Project Start		Project End				
<b>Project Start</b> jj / mm / aaaa	Ö	<b>Project End</b> jj/mm/aaaa	Ö			
Project Start jj/mm/aaaa		Project End jj / mm / aaaa				
Project Start jj/mm/aaaa ID 133671	Ë	Project End jj / mm / aaaa				
Project Start jj/mm/aaaa ID 133671 Funder	Ö	Project End jj/mm/aaaa	Ë			
Project Start jj/mm/aaaa ID 133671 Funder European Commis	Sion (Horizon)	Project End jj/mm/aaaa				
Project Start jj / mm / aaaa ID 133671 Funder European Commis	Sion (Horizon)	Project End jj/mm/aaaa				
Project Start jj/mm/aaaa ID 133671 Funder European Commis Funding status	ssion (Horizon)	Project End jj/mm/aaaa				
Project Start jj/mm/aaaa ID 133671 Funder European Commis Funding status - Please select one	ssion (Horizon)	Project End jj/mm/aaaa	÷			
Project Start jj/mm/aaaa ID 133671 Funder European Commis Funding status - Please select one Grant number/url	ssion (Horizon) e -	Project End jj/mm/aaaa	~			
Project Start jj/mm/aaaa ID 133671 Funder European Commis Funding status - Please select one Grant number/url	e -	Project End jj/mm/aaaa	~			

52

#### Funder Templates

Templates are provided by a funder.

Templates for data management plans are based on the specific requirements listed in funder policy documents. DMPonline.be maintains these templates, however, researchers should always consult the funder guidelines directly for authoritative information.

Template Name 🔶	Download	Organisation Name	Last Updated \$	Funder Links	Sample Plans (if available)
BELSPO DMP +	w k	Belgian Federal Science Policy Office (BELSPO)	27-09-2021		
ERC DMP +	w k	European Research Council (ERC)	27-09-2021		
DCC Template	w k	Digital Curation Centre	27-09-2021		
Horizon 2020 FAIR DMP +	w k	European Commission (Horizon)	27-09-2021		
BRAIN 2.0	w L	Belgian Federal Science Policy Office (BELSPO)	13-05-2022	www.belspo.be	
FNRS DMP	w k	onds National de la Recherche Scientifique (FNRS)	19-05-2022		
VLAIO cSBO DMP (Flemson Standard DMP)		Vlaams Agentschap Innoveren & Ondernemen (VLAIO)	02-09-2022		
Horizon Europe DMP +	w k	European Commission (Horizon)	12-10-2022		
FWO DMP (Flemish Standard DMP)	w k	Fonds voor Wetenschappelijk Onderzoek - Research Foundation Flanders (FWO)	24-10-2022		

### **Resources from the Data Ambassadors Network**

You may need help processing those information or have questions about RDM:

The **Data Ambassadors Network** is there for you!

#### **UMONS Data Ambassadors:**

VISEUR	Robert	Business and Economics	Robert.VISEUR@umons.ac.be
GALLAS	Mohamed-Anis	Architecture	Mohamed-Anis.GALLAS@umons.ac.be
COPPEE	Frédérique	Medicine	Frederique.COPPEE@umons.ac.be
GROSJEAN	Philippe	Sciences	Philippe.GROSJEAN@umons.ac.be
DUPONT	Nicolas	Applied Sciences	Nicolas.DUPONT@umons.ac.be
MEYERS	Charlène	Languages	Charlene.MEYERS@umons.ac.be
RIVIERE LORPHEVRE	Edouard	Applied Sciences	Edouard.RIVIERELORPHEVRE@umons.ac.be
SIMOES LOUREIRO	Isabelle	Psychology	Isabelle.SIMOESLOUREIRO@umons.ac.be

+ webinars (data anonymization, how to archive data, etc.)

### **Closing remarks**

- Research data management is not an easy task but is not impossible either.
- There is no unique way to properly manage your data, as long as your are *consistent*, *descriptive* and organize in a way that *anyone can understand*.
- Different tools are available (DMPTool, DMP online) to help you to plan how to manage your data.
- Data Ambassadors and AVRE people are also here to help you.

### Q&A If I was not able to answer your question, please fill in this <u>Microsoft Forms for Q&A</u> and I will contact you later



Credits to <u>Céline Thillou</u>, <u>Judith Biernaux</u>, <u>David Lhoir</u> & <u>Edouard Rivière-Lorphevre</u> <u>UMONS</u> <u>UMONS</u> <u>UMONS</u> <u>UMONS</u> <u>UMONS</u> <u>UMONS</u>

### **Resources for data management**

Share personal data through a repository

Facilitating FAIR practices in Research Methods, Data, And Software in Natural and Engineering Sciences

Complete training about RDM by Macalester College Library

DocFetcher - Fast Document Search (sourceforge.io) to index your files and quickly find their content

<u>Understanding Research Data Management – University of Pittsburgh</u>

https://libereurope.eu/event/data-management-plans-use-and-reuse-webinar/; coming webinar

Challenges in RDM

Completed DMP for « PURE » European Union's Horizon 2020 project <u>Pure Project Data Management Plan</u> (<u>zenodo.org</u>)

Other DMPs for different disciplines: Example DMPs and guidance | DCC

Do's and Don'ts of DMP

### **Resources to find data for your research**

EOSC (Europen Open Science Cloud): <u>https://ec.europa.eu/research/openscience/index.cfm?pg=open-science-cloud</u>

Mendeley Data website: <u>https://data.mendeley.com/datasets</u>

OpenAire : <u>https://explore.openaire.eu/search/find</u>

re3data.org: <u>https://www.re3data.org/:</u> Harvesting several data repositories

Google: <u>https://toolbox.google.com/datasetsearch</u>

FigShare: figshare - credit for all your research

Zenodo: Zenodo - Research. Shared.