# Introduction to Data Management Plans

**S. Venkataraman, Training Officer, OpenAIRE**

s.venkataraman@openaire.eu

*1st August 2023, ICTP, Trieste*

or Go to
[www.menti.com](www.menti.com)
and enter code
83359905

# Bringing together what you've learnt

- Make informed decisions to anticipate and avoid problems.

- Avoid duplication, data loss and security breaches.

- Develop procedures early on for consistency.

- Ensure data are accurate, complete, reliable and secure.

- Save time and effort to make your life easier!

- Useful both to researchers and institutions

Schiermeier, Q. "Data management made simple" *Nature* **555**, 403-405 (2018).
https://www.nature.com/articles/d41586-018-03071-1
doi: 10.1038/d41586-018-03071-1

# Data Management Plans

A formal document explaining how research data will be handled throughout the data lifecycle.

- Mandatory in e.g. European Commission funded projects
- National and institutional mandates
- Increasingly often required as part of PhD requirements in some places

**Deliverable and "living" document**
Documents processes undertaken throughout data management lifecycle, including costs

**What a DMP is <u>not</u>?**
Research assessment method

| Data description | Documentation and metadata | Storage and backup | Legal and ethical issues | Data sharing and long-term preservation | Responsibilities and resources |

# Data Management Plans

- Depends on the funder/institution requirements
- Differences in research communities
  - Formats, standards, documentation etc.

**Minimum requirements:** Science Europe – **DDPs** (Domain Data Protocols)

# Common themes in DMPs

1. Description of data to be collected / created (i.e. content, type, format, volume...).
2. Standards/methodologies for data collection & management.
3. Ethics and Intellectual Property (highlight any restrictions on data sharing e.g. embargoes, confidentiality).
4. Plans for data sharing and access (i.e. how, when, to whom).
5. Strategy for long-term preservation.

# Planning trick 1: think backwards

What data organisation would a re-user like?



Design how you will organise data in the project (folder structure, file naming convention, ...)

# Planning trick 2: include RDM stakeholders



Commercial partners

Publishers Data Availability policy

**Researchers**

**Front office**

**Back office** data centers

Institution RDM policy Facilities

Research funders

- Information and awareness
- Training
- Storage

www.openaire.eu/briefpaper-rdm-infonoads

# Planning trick 3: ground your plan in reality

Base plans on available skills, support and good practice for the field – show it's feasible to implement.

# Horizon Europe DMP Template

- Guidance: issues to be covered
- Structured approach
- Living document

Initial DMP → Updates → Periodic reporting

https://ec.europa.eu/research/participants/docs/h2020-funding-guide/cross-cutting-issues/open-access-data-management/data-management_en.htm

https://enspire.science/wp-content/uploads/2021/09/Horizon-Europe-Data-Management-Plan-Template.pdf

The Horizon Europe Model Grant Agreement requires that a data management plan ('DMP') is established and regularly updated. The use of this template is recommended for Horizon Europe beneficiaries. In completing the sections of the template the requirements for research data management of Horizon Europe as described in article 17 and analysed in the Annotated Grant Agreement, article 17, must be addressed.

## 1. Data Summary

Will you re-use any existing data and what will you re-use it for? State the reasons if re-use of any existing data has been considered but discarded.

What types and formats of data will the project generate or re-use?

What is the purpose of the data generation or re-use and its relation to the objectives of the project?

What is the expected size of the data that you intend to generate or re-use?

What is the origin/provenance of the data, either generated or re-used?

To whom might your data be useful ('data utility'), outside your project?

## 2. FAIR data

### 2.1. Making data findable, including provisions for metadata

Will data be identified by a persistent identifier?

Will rich metadata be provided to allow discovery? What metadata will be created? What disciplinary or general standards will be followed? In case metadata standards do not exist in your discipline, please outline what type of metadata will be created and how.

Will search keywords be provided in the metadata to optimize the possibility for discovery and then potential re-use?

Will metadata be offered in such a way that it can be harvested and indexed?

### 2.2. Making data accessible

Repository:

Will the data be deposited in a trusted repository?

Have you explored appropriate arrangements with the identified repository where your data will be deposited?

Does the repository ensure that the data is assigned an identifier? Will the repository resolve the identifier to a digital object?

Data:

Will all data be made openly available? If certain datasets cannot be shared (or need to be shared under restricted access conditions), explain why, clearly separating legal and contractual reasons from intentional restrictions. Note that in multi-beneficiary projects it is also possible for specific beneficiaries to keep their data closed if opening their data goes against their legitimate interests or other constraints as per the Grant Agreement.

If an embargo is applied to give time to publish or seek protection of the intellectual property (e.g. patents), specify why and how long this will apply, bearing in mind that research data should be made available as soon as possible.

Will the data be accessible through a free and standardized access protocol?

If there are restrictions on use, how will access be provided to the data, both during and after the end of the project?

How will the identity of the person accessing the data be ascertained?

Is there a need for a data access committee (e.g. to evaluate/approve access requests to personal/sensitive data)?

Metadata:

Will metadata be made openly available and licenced under a public domain dedication CC0, as per the Grant Agreement? If not, please clarify why. Will metadata contain information to enable the user to access the data?

# Assistance: tools

- Wizards
  - Web-based tools
  - Usually free for individual researchers
  - Inbuilt templates
  - Customizable (for institutions)
  - Machine readable DMPs
  - Integration with repositories
- Checklists
  - Swedish National Data Service DMP checklist
  - Customised versions are often found on the websites of academic libraries

# DMP Online

- Established by DCC
- Free for individual researchers
- Institutional instances
- Templates
- Public DMPs
- Export: xml, csv, txt, PDF, html

https://dmponline.dcc.ac.uk/

# Argos

- Established by OpenAIRE
- Based on open-source software
- Free for individual researchers
- User interface in many languages
- Templates
- Public DMPs and datasets
- Machine actionable (based on the RDA Common Standard for machine-actionable Data Management Plans)
- Export: JSON, xml
- Publish on Zenodo

- Strong community support – monthly calls: https://www.openaire.eu/argos-community-calls



https://argos.openaire.eu

# Best practice: DMP examples and use cases

- DMP Use Case Project (OpenAIRE Austria): a collection of public DMPs of EC-funded projects: https://hdl.handle.net/11353/10.1140797

- Public DMPs in Argos, DMP Online etc.

- Search in Zenodo

# DCC Checklist for a DMP

- The DCC assessed existing funder requirements, DMP templates and other best practice to see what should be included in plans. This was synthesised down into common themes and questions.

- 13 questions on what's asked across the board.

- Prompts/pointers to help researchers get started.

- Guidance on how to answer.

# Swedish National Data Service DMP checklist

https://snd.gu.se/en/manage-data/guides/dmp-checklist

- "The SND DMP checklist is designed so that it can be used as a complete and comprehensive plan for an entire research project. It can be used for different research areas, data types, phases of the research process, requirements from funding bodies, as well as various legal requirements."

- Very useful to those working with sensitive data

# Discipline-specific DMPs



- Guidance Document Presenting a Framework for Discipline-specific Research Data Management (Science Europe, 2018), https://www.scienceeurope.org/our-resources/guidance-document-presenting-a-framework-for-discipline-specific-research-data-management

- Training materials in the SSH Open Marketplace: https://marketplace.sshopencloud.eu/search?order=score&q=data+management+plan&categories=training-material

# Assessment frameworks for DMPs

- DMP Evaluation Rubric (part of Practical Guide to the International Alignment of Research Data Management)
- "Reverse engineering"
- Structured approach
- Customised versions available on the websites of academic libraries

| DMP question | DMP guidance | Performance level | |
| --- | --- | --- | --- |
| | | Sufficiently addressed | Insufficiently addressed |
| | | | |

Science Europe. (2021). Practical Guide to the International Alignment of Research Data Management - Extended Edition. https://doi.org/10.5281/zenodo.4915862

Science Europe. (2021). Practical Guide to the International Alignment of Research Data Management - Extended Edition. https://doi.org/10.5281/zenodo.4915862

| 2b | | Sufficiently Addressed The DMP... | Insufficiently Addressed The DMP... |
|---|---|---|---|
| **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. | • Clearly describes the approach taken to ensure and document quality control in the collection of data during the lifetime of the project. | • Provides no information or only a vague mention on how data quality is controlled and documented during the lifetime of the project. |

**3  STORAGE AND BACKUP DURING THE RESEARCH PROCESS**

| **Guidance for Researchers** | | **Sufficiently Addressed The DMP...** | **Insufficiently Addressed The DMP...** |
|---|---|---|---|
| **3a**<br><br>**How will data and metadata be stored and backed up during the research?** | • Describe where the data will be stored and backed up during research activities and how often the backup will be performed. It is recommended to store data in least at two separate locations.<br>• Give preference to the use of robust, managed storage with automatic backup, such as provided by IT support services of the home institution. Storing data on laptops, stand-alone hard drives, or external storage devices such as USB sticks is not recommended. | • Clearly (even if briefly) describes:<br>› The location where the data and backups will be stored during the research activities.<br>› How often backups will be performed.<br>› The use of robust, managed storage with automatic backup (for example storage provided by the home institution).<br><br>or<br><br>• Explains why institutional storage will not be used (and for what part of the data) and describes the (additional) locations, storage media, and procedures that will be used for storing and backing up data during the project. | • Provides no information or very vague reference to how data will be stored and backed up during the project. |
| **3b**<br><br>**How will data security and protection of sensitive data be taken care of during the research?** | • Explain how the data will be recovered in the event of an incident.<br>• Explain who will have access to the data during the research and how access to data is controlled, especially in collaborative partnerships. | • Clearly explains:<br>› How the data will be recovered in the event of an incident.<br>› Which institutional and/or national data protection policies are in place and provides a link to where they can be accessed.<br>› Who will have access to the data during the research. | • Provides little or no details on how the data will be recovered in the event of an incident, which institutional data protection policies are in place, and who will have access to the data during the research. |

# What makes a good DMP?

o Clear, detailed information that is relevant to the science:
  o adopting recognised standards.
  o practices in line with norms for that field.
  o use of support services e.g. university storage, subject repositories…
o Realistic approach that is feasible to implement.

o Evidence of consultation and seeking advice.

o Proper justification of restrictions and costs.

o **Have you taken time to reflect on what to do?**

# Is the information specific enough?

*"we will use suitable formats to ensure that our data can be preserved and sustained over the long term"*

o Which standards? Name them!

o Show that you know which are suitable.

o Does your chosen repository have preferences?

# Are decisions justified?

*"data will be made available upon request to bona fide medieval historians"*

o Why is it restricted?

o Could other communities not reuse the data?

o Will the research team be around to handle access requests in the future?

# A better response...

*"We will provide MP3 audio files for online dissemination. While this is not an open format, it is well-established and the most widely supported. High-resolution WAV files will be used for the archival master recordings."*

o Be clear, specific and detailed.

o Justify decisions.

# Example plans

o Plans from several funders and disciplines via DCC www.dcc.ac.uk/resources/data-management-plans/guidance-examples

o Scientific DMPs submitted to the NSF (USA) provided by DataOne https://www.dataone.org/data-management-planning

o DMPs published in RIO journal http://riojournal.com/browse_user_collection_documents.php?collection_id=3&journal_id=17

o Share yours! - www.dcc.ac.uk/share-DMPs

# Data description examples

The final dataset will include self-reported demographic and behavioural data from interviews with the subjects and laboratory data from urine specimens provided.

From NIH data sharing statements

Every two days, we will subsample E. affinis populations growing under our treatment conditions. We will use a microscope to identify the life stage and sex of the subsampled individuals. We will document the information first in a laboratory notebook and then copy the data into an Excel spreadsheet.  The Excel spreadsheet will be saved as a comma separated value (.csv) file.

From DataOne – E. affinis DMP example

# Metadata examples

Metadata will be tagged in XML using the Data Documentation Initiative (DDI) format. The codebook will contain information on study design, sampling methodology, fieldwork, variable-level detail, and all information necessary for a secondary analyst to use the data accurately and effectively.

From ICPSR Framework for Creating a DMP

We will first document our metadata by taking careful notes in the laboratory notebook that refer to specific data files and describe all columns, units, abbreviations, and missing value identifiers.  These notes will be transcribed into a .txt document that will be stored with the data file.  After all of the data are collected, we will then use EML (Ecological Metadata Language) to digitize our metadata. EML is one of the accepted formats used in ecology, and works well for the types of data we will be producing. We will create these metadata using Morpho software, available through KNB. The metadata will fully describe the data files and the context of the measurements.

From DataOne – E. affinis DMP example

# Data sharing examples

The videos will be made available via the bristol.ac.uk website (both as streaming media and downloads) HD and SD versions will be provided to accommodate those with lower bandwidth. Videos will also be made available via Vimeo, a platform that is already well used by research students at Bristol. Appropriate metadata will also be provided to the existing Vimeo standard.

All video will also be available for download and re-editing by third parties. To facilitate this Creative Commons licenses will be assigned to each item. In order to ensure this usage is possible, the required permissions will be gathered from participants (using a suitable release form) before recording commences.

From University of Bristol Kitchen Cosmology DMP

We will make the data and associated documentation available to users under a data-sharing agreement that provides for: (1) a commitment to using the data only for research purposes and not to identify any individual participant; (2) a commitment to securing the data using appropriate computer technology; and (3) a commitment to destroying or returning the data after analyses are completed.

From NIH data sharing statements

# Restrictions examples

Because the STDs being studied are reportable diseases, we will be collecting identifying information. Even though the final dataset will be stripped of identifiers prior to release for sharing, we believe that there remains the possibility of deductive disclosure of subjects with unusual characteristics. Thus, we will make the data and associated documentation available to users only under a data-sharing agreement.

From NIH data sharing statements

1. Share data privately within 1 year.
   *Data will be held in Private Repository, but metadata will be public*

2. Release data to public within 2 years.
   *Encouraged after one year to release data for public access.*

3. Request, in writing, data privacy up to 4 years.
   *Extensions beyond 3 years will only be granted for compelling cases.*

4. Consult with creators of private CZO datasets prior to use.
   *Pis required to seek consent before using private data they can access*

From Boulder Creek Critical Zone Observatory DMP

# Archiving examples

The investigators will work with staff at the UKDA to determine what to archive and how long the deposited data should be retained. Future long-term use of the data will be ensured by placing a copy of the data into the repository.

From ICPSR Framework for Creating a DMP

Data will be provided in file formats considered appropriate for long-term access, as recommended by the UK Data Service. For example, SPSS Portal format and tab-delimited text for qualitative tabular data and RTF and PDF/A for interview transcripts. Appropriate documentation necessary to understand the data will also be provided. Anonymised data will be held for a minimum of 10 years following project completion, in compliance with LSHTM's Records Retention and Disposal Schedule. Biological samples (output 3) will be deposited with the UK BioBank for future use.

From Writing a Wellcome Trust Data Management and Sharing Plan

# Thank you!

Questions?

(Please get in touch!)