

Data Management Plans

Sarah Jones Digital Curation Centre, Glasgow <u>sarah.jones@glasgow.ac.uk</u>

Twitter: @sjDCC

Data Management Plan (DMP) workshop, e-infrastructures Austria, Vienna, 17 November 2016

What is a DMP?

A DMP is a brief plan to define:

- how the data will be created
- how it will be documented
- who will be able to access it
- where it will be stored
- who will back it up
- whether (and how) it will be shared & preserved

DMPs are often submitted as part of grant applications, but are useful whenever researchers are creating data.

Why manage data?

NON PECUNIAE INVESTIGATIONIS CURATORE SED VITAE FACIMUS PROGRAMMAS DATORUM PROCURATIONIS

(Not for the research funder, but for life we make data management plans)

- Make your research easier
- Stop yourself drowning in irrelevant stuff
- Save data for later
- Avoid accusations of fraud or bad science
- Write a data paper
- Share your data for re-use
- Get credit for it

Don't undervalue research data



PUBLICATIONS AND DATA

Benefits of DMPs for institutions

- Opportunity to engage with researchers and improve RDM practice
- Raise awareness of support available
- Collate information to inform service delivery
- Ensure the University is not exposed to risk
- Ability to recover costs via grants

Research data lifecycle



Planning trick 1: think backwards

What data organisation would a re-user like?



Data organisation

Meaningful file names white_data_20140708.csv Below are tips on meaningful and consistent file names. Read more in 'Choosing a file name'. (2) Make sure to use consistent file names. When you use a date in the blue_data_20140708.docx file name, choose a notation (for instance, YYYYMMDD of yymmdd). Do not use strange characters like ?\!@*%{[<> in the file name. Use traceable file names, such red data 20140708.R as Project_Instrument_locatie_YYYYMMDD.ext. Make sure to only use each file once in the folder structure. If you store a file in more than one place, several versions of the same file red data 20140708 v02.R can unwillingly be created. See also version management. It is good practice to note the file naming and its meaning in a File naming and version management readme.txt.

Even if a researcher is well underway with his project consistent file naming is still an option by using a <u>bulk file</u> <u>rename utility</u>.⁽³⁾ It is important, however, to check if this bulk renamer delivers on its promises.

http://datasupport.researchdata.nl/en/start-de-cursus/iii-onderzoeksfase/organising-data

Planning trick 2: include RDM stakeholders



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











DCC support on DMPs

- Webinars and training materials
- How-to guides and other advisory documents
- Checklist on what to cover in DMPs
- Example DMPs
- DMPonline

www.dcc.ac.uk/resources/data-management-plans





What is DMPonline?

A web-based tool to help researchers write DMPs

Includes a template for Horizon 2020

My plan (Horizon 2020 DMP)	0/9 questions answered approx- 16% of <mark>a</mark> vailable space used
Plan details Initial DMP Detailed DMP Final review DMP Share Export 1. Data summary (1 question, 0 answered) 2. FAIR data (4 questions, 0 answered)	+ +
3. Allocation of resources (1 question, 0 answered)	
 Explain the allocation of resources, addressing the following issues: Estimate the costs for making your data FAIR. Describe how you intend to cover these costs Clearly identify responsibilities for data management in your project Describe costs and potential value of long term preservation 	Guidance Share note EC Guidance - Note that costs related to open access to research data are eligible as part of the Horizon 2020 grant (if compliant with the Grant Agreement conditions). - Costs are eligible for reimbursement during the duration of the project under the conditions defined in the H2020 Grant Agreement, in particular Article 6 and Article 6.2 D.3, but also other articles relevant for the cost category chosen. Glasgow Uni guidance on Resourcing * DCC guidance on Responsibilities +

https://dmponline.dcc.ac.uk

Main features in DMPonline

- Templates for different requirements (funder or institution)
- Tailored guidance (funder, institutional, discipline-specific etc)
- Ability to provide examples and suggested answers
- Supports multiple phases (e.g. pre- / during / post-project)
- Granular read / write / share permissions
- Customised exports to a variety of formats
- Shibboleth authentication \rightarrow eduGAIN

Guidance in DMPonline



Options for unis to customise DMPonline

W ^{NIVE} R _F	Signed in as Sarah Jones. •
The University has a <u>Research Data Management policy</u>	and provides many services to support data management and sharing. See the <u>RDM webpages</u> for more information.
ONLINE View plans Create plan About Roadmap Help	
My plan (UoE Data Management Plan)	0/10 questions answered approx. 25% of available space used
Plan details Default UoE plan Share Export	
Data Capture (2 questions, 0 answered)	+
Data Management (2 questions, 0 answered)	
How will the data be documented to ensure it can be understood? Example of answer Metadata will be tagged in XML using the Data Documentation Initiative (DDI) format. The codebook will information on study design, sampling methodology, fieldwork, variable-level detail, and all information ne for a secondary analyst to use the data accurately and effectively. From the ICPSR <u>Framework for Creating a Data Management Plan</u> .	Guidance Share note contain Producing good documentation and metadata provides context for your data, and makes it easier to find and use in the long term. The amount of effort put into documenting your data will depend on the intended lifespan and how broadly you intend to share it.
B I IIII IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII	anisations can: Add their own template(s) Customise existing funder templates Provide example and suggested answers Local guidance with links to support and services Include their own logo and text in a banner Review basic statistics

A single platform for all things DMP

Agreed to converge on a single codebase, based on DMPonline with additional features from DMPTool

Bring together features and strengths of each tool

Co-manage, co-develop and issue joint roadmap

DMPRoadmap: https://github.com/DMPRoadmap



A FAIR approach to DMPs

Findable

- Assign persistent IDs, provide metadata, register in a searchable resource...

Accessible

 Retrievable by their ID using a standard protocol, metadata remain accessible even if data aren't...

Interoperable

 Use formal, broadly applicable languages, use standard vocabularies, qualified references...

Reusable

- Rich metadata, clear licences, provenance, use of community standards...

www.force11.org/group/fairgroup/fairprinciples

H2020 template

- 1. Data summary
- 2. FAIR data

2.1 Making data findable, including provisions for metadata

2.2 Making data openly accessible

2.3 Making data interoperable

2.4 Increase data re-use (through clarifying licences)

- 3. Allocation of resources
- 4. Data security
- 5. Ethical aspects
- 6. Other issues

http://ec.europa.eu/research/participants/data/ref/h2020/grants_manual/hi /oa_pilot/h2020-hi-oa-data-mgt_en.pdf

Key differences in H2020

- The Commission does NOT require applicants to submit a DMP at the proposal stage. It's a deliverable (due by month 6).
- A DMP is therefore NOT part of the evaluation
- Optional section on data management in proposal is worth doing, especially to help justify costs
- A DMP is a living or "active" document that should be updated

Example H2020 DMPs in Zenodo

Helix Nebula – High Energy Physics example <u>https://zenodo.org/record/48171#.WATexnriF40</u>

Tweether – engineering (micro-electronics) example <u>https://zenodo.org/record/55791#.WATei3riF40</u>

AutoPost – ICT example <u>https://zenodo.org/record/56107#.WATefXriF40</u>

Example: OpenMinTed

1. Table of Contents

OpenMinTed aims to create an infrastructure for Text and Data Mining (TDM) of scientific and scholarly content

Have adopted their own structure to create a 'Data and Software Management Plan'

2.	INTRODUCTION	5
2.1	PROJECT BACKGROUND	5
2.2	GUIDELINES AND REVISIONS	5
3.	DATA IN OPENMINTED	6
3.1	RESEARCH DATASETS	6
3.2	OTHER DATA	
3.2.1	POLL DATA AND ANALYTICS FROM COMMUNITY SURVEYS	
3.2.2	ONLINE KNOWLEDGE BASE OF TDM ISSUES	10
3.2.3	ANALYTICS FOR VISITORS AND TRENDS	
4.	SOFTWARE IN OPENMINTED	12
4.1	INTRODUCTION	
4.1	INTRODUCTION	12 12
4.2	INTRODUCTION	
4.1 4.2 4.2.1 4.2.2	INTRODUCTION	
4.2 4.2.1 4.2.2 4.2.3	INTRODUCTION	
4.2 4.2.1 4.2.2 4.2.3 4.3	INTRODUCTION	12 12 12 13 13 13 14
4.2 4.2.1 4.2.2 4.2.3 4.3 4.4	INTRODUCTION SOFTWARE ASSETS TO BE PRODUCED	12 12 12 13 13 13 14 14
4.2 4.2.1 4.2.2 4.2.3 4.3 4.4 4.4.1	INTRODUCTION SOFTWARE ASSETS TO BE PRODUCED	12 12 12 13 13 13 14 14 15 15
4.2 4.2.1 4.2.2 4.2.3 4.3 4.4 4.4.1 4.5	INTRODUCTION SOFTWARE ASSETS TO BE PRODUCED SOFTWARE TO BE PRODUCED WHAT EXISTING THIRD-PARTY SOFTWARE WILL YOU USE? WHAT IS THE PROCESS FOR DOCUMENTING AND TRACKING SOFTWARE ASSETS AND DEPENDENCIES? INTELLECTUAL PROPERTY GOVERNANCE WHAT IS THE GOVERNANCE MODEL FOR YOUR SOFTWARE PROJECT? Access, SHARING AND REUSE	12 12 13 13 13 14 15 15 15 15
4.2 4.2.1 4.2.2 4.2.3 4.3 4.4 4.4.1 4.5 4.5.1	INTRODUCTION SOFTWARE ASSETS TO BE PRODUCED SOFTWARE TO BE PRODUCED WHAT EXISTING THIRD-PARTY SOFTWARE WILL YOU USE? WHAT IS THE PROCESS FOR DOCUMENTING AND TRACKING SOFTWARE ASSETS AND DEPENDENCIES? INTELLECTUAL PROPERTY GOVERNANCE WHAT IS THE GOVERNANCE MODEL FOR YOUR SOFTWARE PROJECT?	12 12 12 13 13 13 14 15 15 15 15 15 16
4.2 4.2.1 4.2.2 4.2.3 4.3 4.4 4.4.1 4.5 4.5.1 4.6	INTRODUCTION SOFTWARE ASSETS TO BE PRODUCED	12 12 13 13 13 14 14 15 15 15 15 16 17
4.1 4.2 4.2.1 4.2.2 4.2.3 4.3 4.4 4.4.1 4.5 4.5.1 4.6 4.7	INTRODUCTION SOFTWARE ASSETS TO BE PRODUCED	12 12 12 13 13 13 14 14 15 15 15 15 16 16 17 18



http://openminted.eu

Example: OpenMinTed – Data chapter

Six high-level datasets identified:

- 1. Scholarly publications
- 2. Language and knowledge resources
- 3. Services and workflows
- 4. Automatically and manually generated annotations
- 5. Consortium publications
- 6. Metadata

Described in a table per dataset (see illustration)

Information Needed	Answer
Data set reference and name	ID or PID (Persistent Identifier) or DOI as a long-lasting reference to a publication or set of publications, depending on usage (eg. temporary data used for processing should be assigned a simple OpenMinTeD ID)
Data set description	Collections of publications stored for searching and sharing as well as input to be processed by workflows and services
Standards and metadata	To be decided, in accordance with our Interoperability guidelines (See also Table 6.)
Data sharing	Available under permissive licenses, (CC-BY 4.0, CC-0 or comparable) but certain conditions (e.g. Noncommercial use=NC) and/or exceptions may also apply
Archiving and preservation	To be decided. Depending on usage, data used for processing will be stored in PITHOS, persistent storage provided by GRNET and follow its respective archiving and preservation policy

OpenMinTed – Software examples

4.3 Intellectual Property		
Information Needed	Answer	
What type of software license have you chosen?	We plan to release our code under the Apache Software License 2.0 (<u>http://www.apache.org/licenses/LICENSE-2.0</u>) if possible. Alternative OSI-approved open source licenses if necessary due to external constraints, e.g. due to copyleft conditions on used third-party libraries. In some cases, external constraints may also mandate the release under proprietary licenses.	
Have you chosen an OSI- approved open source license	yes	
Is the license valid under your national laws	yes	
Is your license acceptable to all partners	yes	
Where will you publish your license	A copy of the license must be included with every released artifact, e.g. in e Java JAR file, in a ZIP archive containing the software. Furthermore, each source file shall bear a suitable license header if this is technically feasible	

How you will track who does	GitHub's issue tracker also provides full tracking and history of
and has done what	changes is also available at any time. Due to the distributed nature
	of Git, copies of the full history are also typically maintained in
	clones of the canonical repository on Github, e.g. in the copies that
	individual developers keep on their machines or on the build server.
	Strong cryptographic measures ensure that changes to history (if
	made at all) do not go unnoticed.

Helix Nebula: access policy

"The 4 LHC experiments have policies for making data available, including reasonable embargo periods, together with the provision of the necessary software, documentation and other tools for re-use."

"The meta-data catalogues are typically experiment-specific although globally similar. The "open data release" policies foresee the available of the necessary metadata and other "knowledge" to make the data usable"

"Re-use of the data is made by theorists, by the collaborations themselves, by scientists in the wider context as well as for education and outreach."

Helix Nebula: open data

"Data releases through the CERN Open Data Portal (<u>http://opendata.cern.ch</u>) are published with accompanying software and documentation. A dedicated education section provides access to tailored datasets for self-supported study or use in classrooms."

"All materials are shared with Open Science licenses (e.g. CCO or CC-BY) to enable others to build on the results of these experiments. All materials are also assigned a persistent identifier and come with citation recommendations."

"The data behind plots in publications has been made available since many decades via an online database: <u>http://hepdata.cedar.ac.uk</u>"

Plan to share data from the outset

Decisions made early on affect what you can do later

- Negotiation on licenses and consent agreement may preclude later sharing if not careful
- Costings can't be included retrospectively
- Useful to consider data issues at the consortium negotiation stage to make sure potential issues are identified and sorted asap

Key messages

• Data management is part of good practice whether you plan to make the data open or not

– it benefits you!

- The process of planning is the most important aspect. Think about the desired end result and plan for this.
- Approach DMPs in whatever way best fits your project. Don't just let funder requirements drive things.

Thanks for listening

DCC resources on DMPs

www.dcc.ac.uk/resources/data-management-plans

Follow us on twitter: @DMPonline and #DMPonline