# Data Management Planning

*Pinar Alper*

*Training on Research Data Management*
*25-26 June 2019*
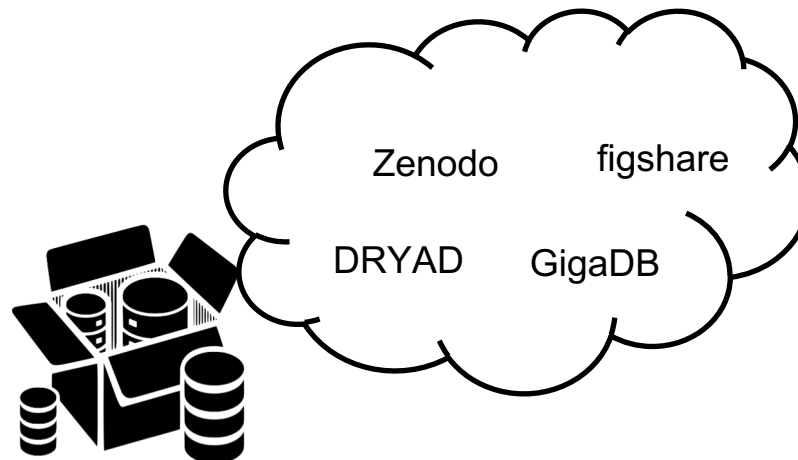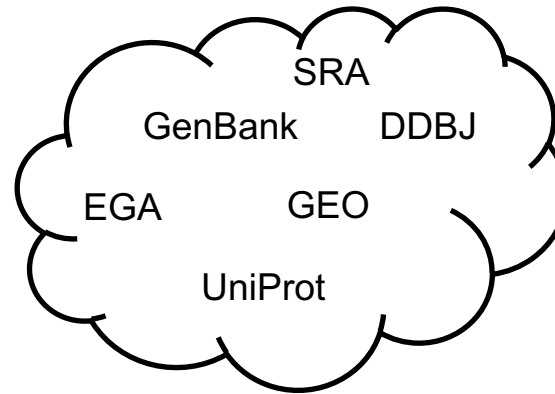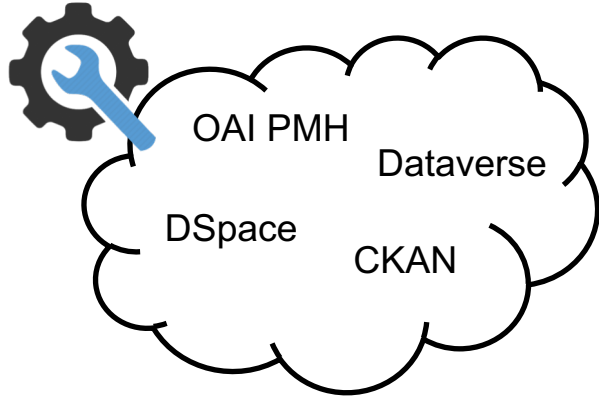*Luxembourg Learning Centre*

# Why DMPs?

## - Funder requirement

- "Ensure data management via DMPs."

- "Researchers are accountable for how data is treated during and after the project."

- timely release of data - once patents are filed or on (acceptance for) publication

- (open) data sharing - minimal or no restrictions if possible

- preservation of data - typically 5-10+ years

*Introduction to Data Management.* Joy Davidson. UK Digital Curation centre. 2015

*Turning FAIR into reality.* Final Report and Action Plan from the EC Expert Group on FAIR Data. November 2018

*Funding research data management and related infrastructures.* Knowledge Exchange and Science Europe briefing paper. May 2016

# Making data a first class-citizen in research
## — decades of efforts, varying levels of "FAIR" ness

OAI PMH

Dataverse

DSpace

CKAN

SRA

GenBank        DDBJ

EGA        GEO

UniProt

Zenodo        figshare

DRYAD        GigaDB

# Making data a first class-citizen in research

- A change in research culture and funding

- DMP is one such intervention

- Future interventions; reporting standards, code, computational reproducibility?

# What is a DMP?

- formal document

- awareness tool



Project        Eternity

# A DMP is shaped by

(Funder) Requirements

Host Infrastructure

Consortium/ Project

Your Research

DMP

tip

An example goes a long way.

# The DMP worldview



*Adapted from "Ten Simple Rules for Creating a Good Data Management Plan". W Michener*

# Identification

- What constitutes " data "

  Primary/Derived data

  Research Record

  Accompanying documents

**tip**    Be thorough when identifying data.

- Type, structure, format, size

Identify data as early as possible, ideally during consortium setup.

**tip**

Use of proprietary formats must be justified. Normally, a big no-no for preservation.

# Collection or generation

- Reference data

- Newly generated data

  - Capture instrument, software, method

  - Period of capture and update

- Sources, Subjects

- Cohort/Study/Site description

# Storage and Preservation

- Platforms

  - Institutional, project-specific, other

- Transfer channels

- Backup policy

- Which data will be retained during and after the project? For how long?

Storing data on hard-drive/server/labbook is not preservation.

Estimate time and effort to prepare data for preservation e.g. export from DB

*tip*

Not all data may need to be retained, e.g. re-computable results.

Check whether institutional resources will suffice for your backup and recovery requirements.

# Organization

- Data Management

- Data Organization

  - File organization/naming

  - Version management

  - Identifiers; accession, doi

  - Searchability

*tip* Paid data management software also goes in to the DMP budget.

# Quality Assurance and Control

- Automated or manual QA/QC measures for validity and integrity

  - Tool/pipeline

  - Training

  - Standards

  - Instrument Calibration

  - Visualisation/Analytics platform



Data Science Analytics
Machine learning
Discovery, New algorithms

Data stewardship
*Standardisation, Harmonisation, Annotation and enrichment, Maintaining access, preserving*

Software stewardship
*Updates, versions, porting*

Prep & Processing
*Data wrangling & curation
Instrument pipelines
Simulation sweeps*



The New York Times

**For Big-Data Scientists, 'Janitor Work' Is Key Hurdle to Insights**

Monica Rogati, Jawbone's vice president for data science, with Brian Wilt, a senior data scientist.
Peter DaSilva for The New York Times

**tip** Be very conservative in your assumptions of source data quality.

Building the FAIR Research Commons: A Data Driven Society of Scientists. C A Goble.

# Ethics , Legal Complaince

- What are the applicable regulations to your data?

- Is ethics/irb approval required?

- For Sensitive Human data, what is your Data Protection Concept?

  - Have you gained consent also for data <u>sharing</u> and <u>preservation?</u>

  - What type of agreements are needed during and after the project. – consortium agreement should cover sharing!

  - What are the protective measures for data and subject privacy?

*tip*

Seek support from Legal Office and Ethics Board. ...Standard clauses..

# Security

- What are the security measures?

  - Encryption, Access Control, Password Management



**tip**

Refer to standards, institutional policy/procedures.

e.g. ISO 27001, SAML SSO
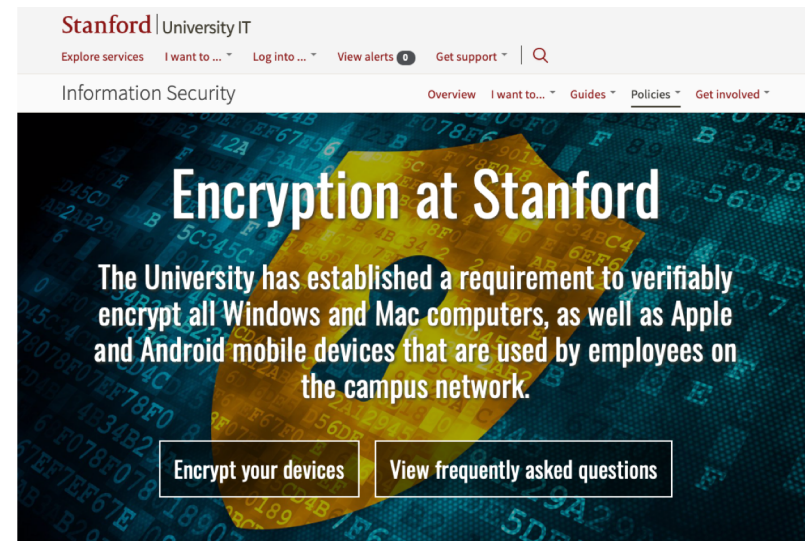
# Documentation

- **Metadata**: Information enabling the read and interpretation of data

- What documentation and metadata will accompany the data

- Required for publicly shared data

- How will you capture/create it

    - Basics e.g. title, source, license

    - Standards e.g. Minimum information guidelines

tip

Estimate effort for documentation before commiting to it in the DMP.

# Ownership, Intellectual rights, Access

- Will data stay in the public domain?

- Or are there IP restrictions for ownership (consortium owns the data), copyright commercialisation.

- What will be the license for the data?

- May be added during DMP update

*tip*   https://eudat.eu/services/userdoc/b2share-usage
http://www.dcc.ac.uk/resources/how-guides/license-research-data

# Dissemination

- If, when and how will you share data?

   Data Paper, Repository, Database

- For sensitive human data what is the access orchestration process?

- May be added during DMP update

**tip** Prefer repositories that handle preservation and access-orchestration.

# DMP as a living document

- You will be expected to update your DMP

**First version**

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.
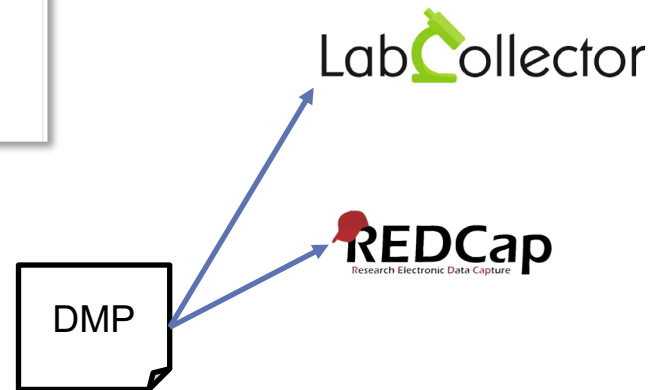
**Updates**

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.
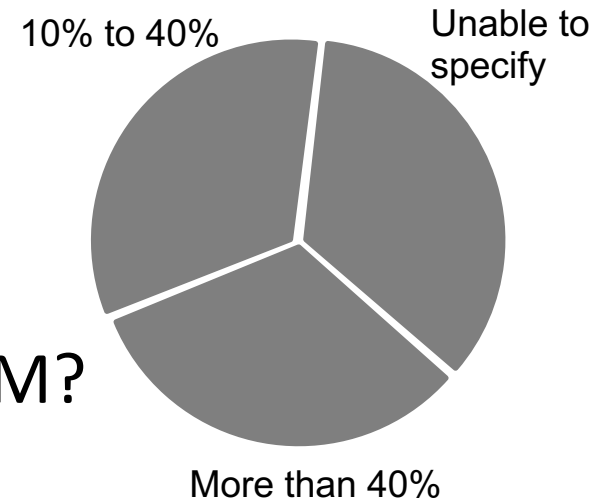- Furthermore, the consortium can define a timetable for review in the DMP itself.

LabCollector

REDCap
Research Electronic Data Capture

DMP

*tip*   Annex DMP with executable specifications rather than static listings.

# Budgeting for Data Management
— 5% of total project costs

**An 2016 survey:**

- Infra providers, libraries, universities

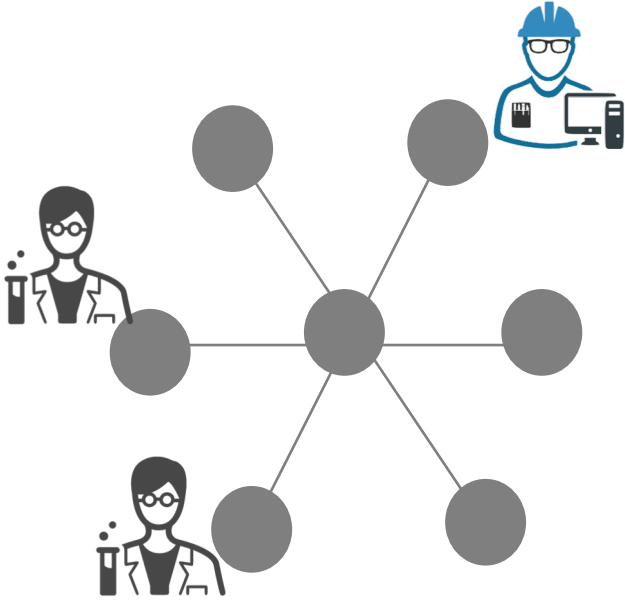- What percentage of total budget of your organization is allocated for RDM?



10% to 40%

Unable to specify

More than 40%

**A commonly cited recommendation:**

- An overall average of **5% of the total project costs** ......to sustain and share data"

# Roles and responsibilities

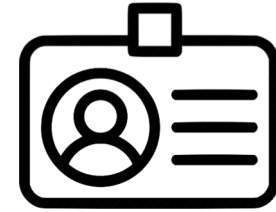— DMP via "data specialist" partner

- Main author of the DMP
- Task planning and follow up.

**tip**

DMP tools may be of help to the "data special" partner
Or to researchers when writing their DMP for the very first time.

# Administrative information

- DMP ID

- Funder, Grant ID

- Project name description

- PI , Data Contact name and ID

- Dates

- Related policies

# DMP tools

- UK DCC **DMP**ONLINE

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

  https://dmponline.be

  roadmap

  https://github.com/DMPRoadmap/

- University of California DCC **DMP**Tool
  Build your Data Management Plan

  https://dmptool.org

- Dutch Techcenter for Life Sciences, ELIXIR-NL  DS Wizard

  https://app.ds-wizard.org/welcome

# DMP Online

## Pinar's Plan

| Project Details | Plan overview | Write Plan | Share | Download |
|---|---|---|---|---|

**\* Project title**

Pinar's Plan

☑ mock project for testing, practice, or educational purposes

**Funder**

Biotechnology and Biological Sciences Research Council (BBSRC)

**Grant number**

345345345

**Project abstract**

TEst abstract

> Briefly summarise your research project to help others understand the purposes for which the data are being collected or created.

**ID**

PROJECT-ID

**Principal Investigator**

**Name**

Pinar Alper

**ORCID iD**

0000-0002-2224-0780

**Email**

pinarpink@googlemail.com

**Phone**

342342324234

**Data Contact Person**

☑ Same as Principal Investigator

Save

**Select**

To help you w
guidance fror

Select up to 6

☑ Digital
☑ Univers

Find guidance

See the full lis

Save

- Document generator app

- 17000 + Users

- 23000 + Plans

- 50 DMPs publicly available

# DMP Online

**2. If you will be applying for funding from multiple sources who else will you be applying to?**

- [ ] Not applicable
- [ ] Arts and Humanities Research Council
- [ ] Cancer Research UK
- [ ] Economic and Social Research Council
- [ ] Engineering and Physical Sciences Research Council
- [ ] European Research Council
- [x] Horizon2020
- [ ] Medical Research Council
- [ ] Natural Environment Research Council
- [ ] Science and Technology Facilities Council
- [ ] Wellcome Trust

**Additional Information**

| B | I | ☰ ▾ | ☰ ▾ | 🔗 | ▦ ▾ |

dfsdfsd

**Save**

`Answered 4 days ago by pinarpink@googlemail.com`

---

Guidance     **Comments**

**Manchester**

If your funder isn't listed, please enter in the free text box provided.

---

**Add comments to share with collaborators**

| B | I | ☰ ▾ | ☰ ▾ | 🔗 | ▦ ▾ |

**Save**

---

**3. Is The University of Manchester the lead institution for this project?**

- ( ) Yes – only institution involved
- ( ) Yes – leading a collaboration
- ( ) No (please provide details of the lead institution below and your role in the project)

**Additional Information**

| B | I | ☰ ▾ | ☰ ▾ | 🔗 | ▦ ▾ |

**Save**

Guidance     **Comments**

**Manchester**

If The University of Manchester is not the lead institution for this project please only consider data that Manchester researchers will be responsible for whilst filling in this form. (Please also ensure that you have discussed with your collaborators how data stored elsewhere will be managed.)

---

**4. What data will you use in this project (please select all that apply)?**

- [x] Acquire new data
- [ ] Re-use existing data (please list below)
- [ ] Generate textual supporting information only
- [ ] Not acquire or re-use data (please provide details)

**Additional Information**

| B | I | ☰ ▾ | ☰ ▾ | 🔗 | ▦ ▾ |

Guidance     **Comments**

**Manchester**

If you are not sure if you have data please refer to the definition at the top of this section. If your

# DMP Online

# ELIXIR DS Wizard



Data Stewardship Mind Map (Rob Hooft)

- "Insights collected over 4 years"

- "Guidance over 300+items "

- Awareness tool

- Under active development

  https://github.com/ds-wizard/ds-wizard/issues

  info@ds-wizard.org

# Common DSW Knowledge Model

**Common DSW Knowledge Model** 1.4.0

`dsw:root:1.4.0` · 🔴 Data Stewardship Wizard · DSW Knowledge Model originated from mindmap made by Rob Hooft

## Knowledge Model

**Name**

Core DS Knowledge Model

**Chapters**

Design of experiment

Data design and planning

Data Capture/Measurement

Data processing and curation

Data integration

Data interpretation

Information and insight



DATA STEWARDSHIP
FOR
OPEN SCIENCE

Implementing FAIR Principles

BAREND MONS

WITH VITALSOURCE EBOOK

CRC Press
Taylor & Francis Group
A CHAPMAN & HALL BOOK

"All models are wrong but some are useful."

*George E. P. Box*

# ELIXIR DS Wizard

# ELIXIR DS Wizard

# ELIXIR DS Wizard

# Should you use a DMP tool?

- DMP tools are not a match to state of collaborative document authoring tools.

- Embodies data management expertise.

- Not integrated to grant submission systems.

DMP ONLINE
https://DMPOnline.be

DMP Tool
Build your Data Management Plan

Pooling DMP templates and instances

DS Wizard

Questionnaires to help the data specialist

*tip*

DMP tools may be of help when there is no "data expert" partner.

Pooling DMP content has proved useful for some research communities.

Be mindful when re-using DMPs; document agreed-upon decisions not possibilities.

"Keep it short and specific"

# Data Stewardship Wizard
# -Practical-

# Scenario

You are the data specialist in project X, which works with human biosamples and data. In the DMP you had specified that the project will deposit its output data to relevant repositories. You had listed GigaDB and Nature Scientific Data as possible venues.

Project completion date is fast approaching and a few manuscripts are already under preparation. You check the submission requirements of above-mentioned venues, they require extra information in the case of human data. You need to collect this information from various consortium members, as data is acquired from various cohorts that collaborators have access.

You decide to encode Nature Scientific Data's submission checklist as a questionnaire for the consortium, and collect information this way.

# Practical

Application URL:

https://learning.ds-wizard.org

Usernames:

datasteward_XX@example.com

researcher_XX@example.com

*(see handout for your number)*

Password:

**<to-be-provided>**

Materials:

**<to-be-provided>**

# Acknowledgements

Jan Slifka, Vojtech Knaisl - ELIXIR Czech Republic

Celia van Gelder, Mateusz Kuzak - ELIXIR Netherlands

Wei Gu, Roland Krause - ELIXIR Luxembourg