



Scientific Data Management Training

5.12.2024





Jana Martínková

Jana.martinkova@cvut.cz 0000-0001-8575-6533 About me



Ing. Jana Martínková

- FIT ČVUT
- ELIXIR CZ
- OSTrails
- DSW Team, DMP Methodology Specialist

jana.martinkova@cvut.cz







- Motivation
- Research Data Management
- FAIR
- Data Management Plan
- Data Stewardship Wizard
 - Demonstration

Motivation



- Data management plans are necessary
- *DMP is a formal requirement*

 \rightarrow but they also gather valuable information about your research

• Writing a DMP is a tedious work

 \rightarrow answer questions instead

• *Researchers are not always aware of the best practices*

 \rightarrow provide guidance



- Required by organizations or funding agencies
- Planning resources and equipment
- Defines roles and responsibilities
- Identifies risks and its solutions
- Facilitates data sharing, reusability, and preservation
- Prevention vs. firefighting





Research Data Management

DATA STEWARDSHIP WIZARD



Data Life Cycle





https://rdmkit.elixir-europe.org

Data Life Cycle: Plan





- Planning on how we will work with data during (and after) the project
- Output should be a **Data Management Plan**

Data Life Cycle: Collect





- Collecting new data (methods differ according to the research area)
- Using existing data (e.g., from previous projects)
- Recording the origin of the data (samples, researchers, tools, etc.)
- Emphasis on the quality of recorded data

Data Life Cycle: Process





- Conversion of data from the recorded format to a format suitable for analysis
- Exclusion of bad data or data of low quality
- Pseudonymization/anonymization of sensitive data

Data Life Cycle: Analyse





- Exploration of collected data
- Main part of research gaining new knowledge
- The workflow used in the analysis should be reproducible
- Big data analysis may require significant computational power

Data Life Cycle: Preserve





- Ensuring long-term preservation of data after the project ends
- The possibility to verify the project's results even after several years
- Using the data in the future for other purposes (teaching, other research)

Data Life Cycle: Share





- Sharing data with others (e.g., in another research project)
- Sharing does not mean that the data must be publicly available; they can be shared with limited access
- Consideration of all ethical, legal, licensing, and other restrictions

Data Life Cycle: Reuse





- Using data for a purpose other than for which they were collected, for example:
 - As reference data for other research
 - Verification of the original research results
 - Combining results from multiple studies into meta-studies





- Researcher = work with data during the research project, create DMP for their project
 - PhD, grant applicants, project managers...

- **Data Steward** = takes care of data on several levels
 - Policy = cooperate with management and grant agencies; oversee processes, ethic and legal issues
 - **Research** = cooperation with the scientists, help with DMP creation
 - **Infrastructure** = cooperation with IT; IT solutions for RDM, infrastructure...











- What is FAIR?
- Principles applied on data to make them:



- The first step towards reusability is finding the data
 - Data are described with rich metadata
 - (Meta)data have a globally unique persistent identifier
 - (Meta)data are registered in a searchable resource





Findability

- After finding the data, it must be clear how they can be accessed
 - Accessible ≠ Open
 - Access to (meta)data via a standard communication protocol
 - Use of authentication and authorization if needed





Accessibility

- Ability to integrate data with other data
- Ability to process in various applications and workflows
- Use of standard formats, vocabularies, and ontologies (RDF, JSON-LD, OWL)





Interoperability

Reusability



- The main goal of FAIR is reusability
- Quality description of (meta)data
 - Licenses
 - Origin
 - Community standards in the given domain



Why should we apply FAIR Principles?



- Who will reuse the data?
 - You
 - Colleagues in the team
 - Colleagues within the institution

• Someone else





Data Management Plan

DATA STEWARDSHIP WIZARD



Data Management Plan (DMP)



- General project information
- Description of data
- Metadata and ontologies, documentation
- Data storage, security, data preservation
- Data sharing
- Costs and human resources
- Ethical and legal issues, licenses

DMP as pre-flight checklist









Data Stewardship Wizard

DATA STEWARDSHIP WIZARD



What is Data Stewardship Wizard?



- Expert system for data management planning and creating DMPs
 - "From burden to benefit"
 - "Plans are worthless, but planning is everything"
- Open-source tool
- Tool suitable for everyone (from beginners to data stewards)
- Serves as a check-list before starting the project
- Supports data management planning with respect to (current) best practices
- A recommended tool in the Horizon Europe Program Guide,
 ELIXIR Recommended Interoperability Resources



DSW Main Ideas



- Minimum of writing = DMP is not an essay, as little writing as possible
- Guidance = DSW guides users through the smart Questionnaire
- Flexibility = easy to edit the content and integrate with other services
- Openness = anybody can use it and create own content
- User-oriented = DSW development is strongly user feedback driven

Demonstration





Data Stewardship Wizard



Knowledge Model



Knowledge Model (KM)



- Contains knowledge about what to ask for and how
- **Template** for structured questionnaire
- Tree structure build of chapters, questions, answers, follow-up questions and other resources



Data Stewardship Wizard



Knowledge Model



Questionnaire



- Can be filled in any order
- Only relevant questions are asked based on the previous answers
- Links to additional resources
- Metrics indications

VIII.2 Will you be collecting physical samples?	+	2	S
Vill you be collecting artefacts like specimens, minerals, biological samples?			
 Desirable: Before Submitting the Proposal Data Stewardship for Open Science: <u>kuz</u> 			
O a. No			
● b. Yes :Ξ			
D Clear answer			
Answered less than a minute ago by Isaac Newton.			
III.2.b.1 Do your samples need to be submitted to a public repository?	+	2	9
In some fields, it is considered good practice to deposit samples to a public repository before deriving any digital data	from the	em.	
Desirable: Before Submitting the Proposal			
○ a. No 🗄			
O b. Yes			
X III.2.b.2 How will the samples be identified?	+	•	S
○ a. Samples will only get an internal code Reusability 0%			
O b. Samples will receive a universally unique persistent identifier (PID) Reusability 100%			

Data Stewardship Wizard





Document Template



- Template development requires some technical knowledge
- Templates are composed of JSON metadata, Jinja2 templates and other files
- **DSW Template Development Kit** (TDK)

EXPLORER ····	{} template.	json ×
> OPEN EDITORS	{} template	e.json >
\vee MADMP-TEMPLATE-MA	1 {	
\sim content	2	"organizationId": "dsw",
≣ _mapping.j2	3	"templateId": "rda-madmp",
≣ _uuids.j2	4	"version": "1.4.0",
≡ madmp.json.j2	5	"name": "maDMP (RDA DMP Comm
≡ madmp.ttl.j2	6	"description": "Machine-acti
• .gitignore	7	"recommendedPackageId": "dsw
	8	"license": "Apache-2.0",
README.ma	9	"metamodelVersion": 3.
	10	"allowedPackages": [
	11	
	12	"orgId": "dsw".
	13	"kmId": "root".
	14	"minVersion": "2.3.0".
	15	"maxVersion": null
	16	},

Data Stewardship Wizard





Data Stewardship Wizard





Generated Documents



		Projects		Section A: Data Collection
		We will be w described in	orking on the following projects and for those this DMP.	1. What data will you collect or create?
Data Mana Science	agement Plan e Europe Example	Arsenic an Start date: End date: Funding: The main ge produced fro Pfibram, Kut herbivorous will be condu	d Selenium Speciation Using Hyphena 1.1.2021 31.12.2021 Grantová Agentura České Republiky: grant (planned) val of this study is to determine whether i om the meadows in the vicinity of old meti ná Hora, and Nalžovské Hory (Czech Republi herds. Total and speciation analysis of As and ucted using a hyphenated technique of HPLC z	Instrument datasets The following instrument datasets will be acquired in the p • HPLC This dataset will be collected by experts in the project equipment. The equipment is very well described and known. • ICP-MS This dataset will be collected by experts in the project equipment. The equipment is very well described and known.
Contact person: Based on: Created by:	Jana Freeman (jana,freeman@ds-wizc 0000-0000-0000-0001) Czech Technical University in Prague. Common DSW Knowledge Model, 2.3.0 Jana Freeman (jana freeman@ds.wizz			Re-used datasets We will use the following reference datasets: • Chemical Component Dictionary (http://dx.doi.org/d We will use the following already existing non-reference d • Previous in-house Arsenic and Selenium analysis We already have a copy of this dataset. Data formats and types We will be using the following data formats and types:
cicultu by.	DSW			Chemistry vocabulary
Generated on:	24 Jun 2021			It is a standardized format. This is a suitable format f We will have only a small amount of data stored in th

"dmp": {

```
"title": "DMP in a planning phase",
"description": "Example of a DMP describing a project in which source code will be crea
"created": "2019-02-22T13:20:15.5",
"modified": "2019-02-22T13:20:15.5",
"project": [],
"contact": {
    "mbox": "TMiksa@sba-research.org",
   "name": "Tomasz Miksa",
    "contact_id": {
        "identifier": "https://orcid.org/0000-0000-0000",
        "type": "orcid"
"language": "eng",
"ethical_issues_exist": "no",
"dmp_id": {
   "identifier": "https://doi.org/10.0000/00.0.1234",
    "type": "doi"
"dataset": [
        "title": "Source Code",
        "description": "Proof of (
```

Data Stewardship Wizard





Data Stewardship Wizard



Data Steward



Researcher

Researchers





Researchers Workflow









- Questionnaire
- Metrics
- Preview
- Documents
- Settings

View Current Phase Before Submitting the Proposal Chapters I. Administrative information	I. Admin	histrative i tributors DMP Horizon Euro	information	TODOs	Comments	Version history
Current Phase Before Submitting the Proposal ~ Chapters I. Administrative information	I. Admin	tributors	information			
Before Submitting the Proposal Chapters I. Administrative information	Horizon 2020	tributors				
Chapters I. Administrative information	Horizon 2020	tributors				+ •
Chapters I. Administrative information	Horizon 2020	OMP Horizon Euro				
I. Administrative information	Each person		rope DMP Science Euro	ope DMP maDMP		
	A project proj	contributing to cre	eating or executing the	e data management	plan should be	added as a contrib
Q Contributors	A project proi	ably should have a	a Contact Person, and	a Data Curator.		
♀ Research Project(s)	☑ Desirable:	Before Submitting	g the DMP			
\mathcal{O} To execute the DMP, is additiona	l sp + Add					
Do you require hardware or softw Describe national / funder / softw	var					
	Jila					
II. Re-using data	1 📝 I.2 Res	earch Project((s)			+ 🗨
III. Creating and collecting data	6 Horizon 2020	OMP Horizon Euro	rope DMP Science Euro	ope DMP maDMP		
IV. Processing data	3 Add each of	he research proje	ect(s) that you are (or	r will be) working o	n and for which	n the data and work
	described in	lis DiviP. Give eaci	on project a small ident	urying name for you	sen.	
V. Interpreting data	2 Desirable:	Before Submitting	g the Proposal			
VI. Preserving data	4 + Add					
	2					

Researchers Workflow











Question types



- Value
- Options
- Multi-Choice
- List of items
- Integration
- Item Select
- File



- Connection to **external database** or service on the Knowledge Model level
- Answer is not just a text but also a **link to selected item**





✓ II.1.b.1.a.1 Reference database or dataset	+ *	2	0
Horizon Europe DMP			
Give the name of the database or dataset. You will be shown suggestions of data bases from FAIR: you can also type the name of a dataset that is not in FAIRsharing	Sharii	ng,	but
☑ Desirable: Before Submitting the DMP			
protein			
The Protein Database (Protein) type knowledgebase The Protein database is a collection of sequences from several sources, including translations from annotated coding regions in GenBank, RefSeq and TPA, as well as records from UniProt/SwissProt PRF, and PDB. Protein sequences are the fundamental determinants of biological structure and function.	n , PIR,		*
LocustMine (LocustMine) type knowledgebase An integrated Omics data warehouse for Locust, Locusta migratoria.			
DistiLD Database: Diseases and Traits in Linkage Disequilibrium Blocks (DistiLD) type knowledgebase			
The Dictil D database aims to increase the usage of evicting geneme wide association studies (G)	(A C)		













FAIR metrics



- Each answer can affect **metrics**
- The resulting value is calculated as a weighted average of all responses that affect the given metric



Online cooperation







TODOs & Comments



- TODOs when you're unsure how to answer a question, you can add TODO
- **Comments** for discussion about the questions
 - Comments vs Editor notes

Comments 1 Albert Einstein	Editor notes
	Comments 1 Albert Einstein 16. 8. 2022, 10:22 (edited)

project documentation.

Create a new comment...

Reply...

Researchers Workflow





Creating Data Management Plan



- Created from a filled questionnaire using a document template
- Document formats based on template
- Saved in the DSW, can be downloaded

H Horizon Europe DMP Data Management Plar	1.2.1 n according to the Horizon Europe template
mat	
O 🛃 HTML	PDF Document
O MS Word Document	

Researchers Workflow





- List of all changes made while filling in the questionnaire
- Named versions
- Revert to version
- Create document from an older version





Version History





Possibilities for Institutions



Customizations and use of DSW

- ELIXIR infrastructure (provided by CESNET, managed by us):
 - ELIXIR is an intergovernmental organisation that brings together life science resources from across Europe.
 - <u>researchers.dsw.elixir-europe.org</u> can be simply used by researchers, not customized
 - Institutional instance can be customized by institution (admins), they manage it
- Self-hosted (on-premise / own cloud solution)
 - DSW is open-source, documentation and Docker images available
 - Institution can host their own instance, manage it, update it, keep the data in-house
- Commercial cloud (<u>https://fair-wizard.com</u>)
 - Enterprise-ready with additional services for institutions and commercial subjects





Acknowledgements



OStrails



This presentation was supported by the EU-funded OSTrails project (Grant Agreement No 101130187) More information on Data Stewardship Wizard support and funding is available at <u>https://ds-wizard.org/#acknowledgements</u>





Discussion and question