HOW TO FAIRIFY DATA IN YOUR OWN LAB

Marta Dembska, DLR Institute of Data Science MSE Day, 14.11.2023



Reality of laboratory data



Data are diverse...



Image Ernesto Enslava on Pixabay

Image by Tatiana on Pixabay



Image by Benis Arapovic on Vecteezy

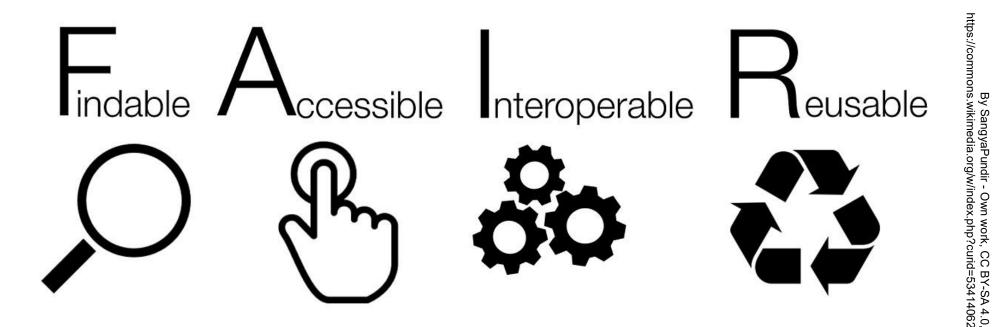




Photo by Campain Creators on Unslpash



Can someone else make use of your data?



FAIRification of your laboratory data starts with making them (internally or externally) available

Findability

F1

(Meta)data are assigned a globally unique and persistent identifier

F2

Data are described with rich metadata (defined by R1)

F3

Metadata clearly and explicitly include the identifier of the data they describe



Photo by Tobias Fischer on Unsplash

F4

(Meta)data are registered or indexed in a searchable resource

Findability



F1 Globally unique and pers	If y	Unique identifiers are the key for publishing FAIR data your data cannot be found, it's like it didn't exist! Keeping a web link active is not for free		
qu	 Generous and extensive – incl. information on context, quality and conditions, or characteristics of the data Make data location based on their metadata possible 			
F3 Connection between data and their metadata		 Metadata and datasets are often separate files! 		
		 Connection should be annotated in a formal manner – there are tools that can be used 		
F4 Metadata searchable	 Findability not gua 	aranteed by rich metadata and identifiers		
Marta Dembska, MSE Day, 14.11.2023	If nobody knows that your data resources exist, they may go unused! • Domain specific registries exist and are ready to use			

Accessibility



A1

(Meta)data are retrievable by their identifier using a standardised communications protocol

A1.1

The protocol is open, free, and universally implementable

A1.2

The protocol allows for an authentication and authorisation procedure, where necessary



Photo by Jonathon Young on Unsplash

A2

A2. Metadata are accessible, even when the data are no longer available

Accessibility



A1 Standard communication protocols retrieval

A1.1 Open and free protocol

A1.2 Authentication and authorisation where necessary

- Avoid barriers to access (e.g. unusual protocols, poor documentation)
- Provide contact person for sensitive data
- Access to at least metadata for (almost) every user
- o Clear instructions on authentication

Accessible does not automatically mean open or free!

Machine actionability of access requirements

A2 Metadata accessible after data are not

Metadata easier and cheaper to store
Metadata important for replication

Interoperability



11

(Meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation

12

(Meta)data use vocabularies that follow FAIR principles

13

(Meta)data include qualified references to other (meta)data



Photo by Kaleidico on Unsplash

Interoperability



I1 Common language for knowledge representation	 Precise syntax and grammar 	
(Meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation	 Language specifications accessible More than one scenario in mind in design process 	
unique and persistent ic	 Controlled vocabularies documented and using unique and persistent identifiers Documentation findable and accessible 	
I3 Meaningful links between (meta)data resouces	 Cross-reference between datasets Connection to enrich the contextual knowledge about them 	

Reusability

R1

(Meta)data are richly described with a plurality of accurate and relevant attributes

R1.1

(Meta)data are released with a clear and accessible data usage license

R1.2

(Meta)data are associated with detailed provenance

R1.3

(Meta)data meet domain-relevant community standards



Photo by <u>Ravin Rau</u> on <u>Unsplash</u>

Reusability



R1 Labels describe the context

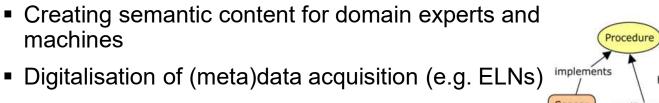
R1.1 Clear usage rights

R1.2 Detailed origin of data

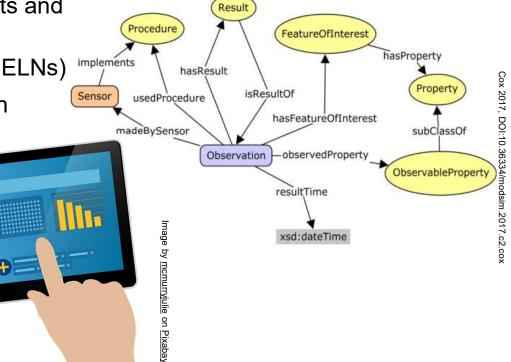
R1.3 Alignment with domain community standards

- Enable user to decide if the data is useful for them (incl. experimental protocols, brand of equipment, laboratory conditions etc.)
- o Metadata as general as possible
- Clear and commonly used licence attached to data
- Provenance of data (e.g. authorship and other roles in creation, workflow description, processing details)
- Similar datasets are easier to reuse
- Follow best practise standards of the community

Enabling interoperability with domain level ontologies



- Knowledge representation in a given domain
- Formalisation of procedures
- Reproducibility of experiments
- Data life cycle
- Quality measures

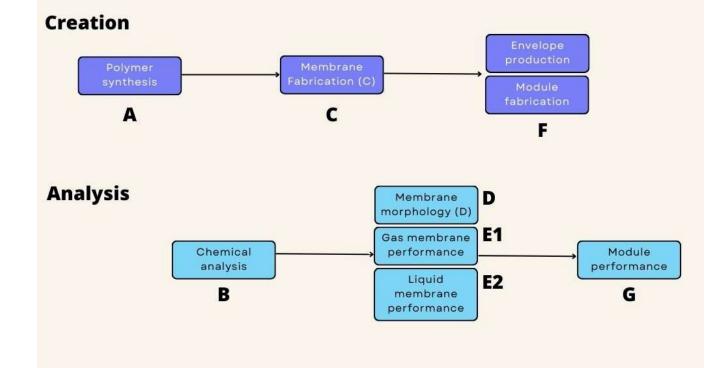


Goals

- Combine knowledge engineering and domain modelling through collaborative work of ontology engineers and domain experts
- Ontology development through use of formal languages and tooling

Domain knowledge

Laboratory work



ELNs

- Herbie (in house solution)
- Chemotion (external)

Vocabulary

- domain experts
- ontology reuse

Process modelling

- documentation
- recipes
- data cycle (metadata capture)

Ontology requirements specification



POLYMAT

Use cases:

- Knowledge representation in the fields of polymer and membrane research
- Documentation of scientific work and laboratory processes
- Ontology used in Herbie
- Knowledge representation in modelling laboratory processes

Research application

(X) creation

(Y) analysis

(Z) database, models and ELN

Example competency questions

- ❑ What polymers were used to fabricate a given membrane? (X)
- What equipment was used in a given experiment or measurement? (X, Y)
- □ Where are the results of a given calculation stored? (Z)
- □ What computational models were used to perform the experiment? (X, Y, Z)
- □ What (physical, chemical) calculated characteristics are

derived from a measured characteristic? (Y)

Ontology requirements specification



POLYLAB

Use cases:

- · Reproducibility of a given experiment
- Cross referencing of protocol sub elements in Herbie

Example competency questions

- □ Who is trained to perform a given experiment or measurement (S, M)
- What is the order number of a given measurement request?
- □ What substances are used in the experiment? (S)
- □ What properties of the equipment are important for a given measurement? (M)

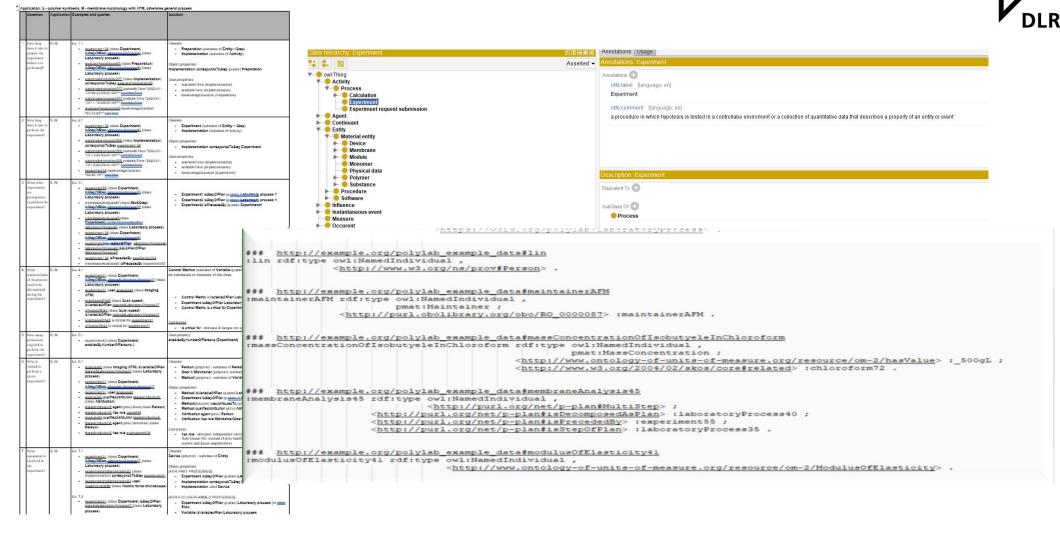
Process

(S) polymer synthesis

(M) membrane morphology with AFM (atomic force microscope)

Otherwise general process

Ontology implementation – conceptualization and encoding



Ontology implementation – evaluation



<u>polytetrafluoroethyleneMembrane5</u> (class: **Membrane for liquids**, **FlatSheetMembrane**) - instance of class **FlatSheetMembrane** <u>polytetrafluoroethylene1</u> (class: **Homopolymer**) - instance of class **Homopolymer**

olytetrafluoroethyleneMembrane5 (class: Membrane for liquids, FlatShee omopolymer, Poly(tetrafluoroethylene))	t Membrane) <i>contains</i> <u>polytetrafluoroethylene1</u> (class:
<u>olytetrafluoroethylene1</u>	CQ1: What polymers were used to fabricate a given membrane?
<u>olytetrafluoroethylene1</u> has IUPAC name Poly(1,1,2,2-tetrafluoroethylene) <u>olytetrafluoroethylene1</u> has product ID MHp201203\$	Q1.1 Return polymers that a specific membrane contains
	PREFIX rdf: <http: 02="" 1999="" 22-rdf-syntax-ns#="" www.w3.org=""> PREFIX rdfs: <http: 01="" 2000="" rdf-schema#="" www.w3.org=""></http:></http:>
	PREFIX pmat: <https: polymat="" w3id.org=""></https:> PREFIX bfo: <http: obo="" purl.obolibrary.org=""></http:>
	PREFIX ex: <http: example.org="" polymat_example_data#=""></http:>
	SELECT DISTINCT ?membrane ?polymer WHERE {
	?membrane_type rdfs:subClassOf* pmat:Membrane. ?membrane rdf:type ?membrane_type . ?membrane bfo:RO_0001019 ?polymer . ?polymer_type rdfs:subClassOf* pmat:Polymer. ?polymer rdf:type ?polymer_type
	FILTER (?membrane=ex:polytetrafluoroethyleneMembrane5>) }
	Simple view Ellipse Filter query results
membrane	polymer
1 <http: example.org="" polymat_example_data#polytetrafluoroethylenemembrane5=""></http:>	http://example.org/polymat_example_data#polytetrafluoroethylene1

Showing 1 to 1 of 1 entries



 $_{\odot}$ Crucial role of people in FAIRification

Domain specialist and knowledge engineers are one team

ELNs as a medium to facilitate digitalisation and FAIRification in the labs

F well described data are findable
 A clear documentation of access
 I shared vocabularies facilitate interoperability
 R reuse, reproduce and avoid repetitions

Contact: <u>Marta.Dembska@dlr.de</u> DLR Institute of Data Science, Jena