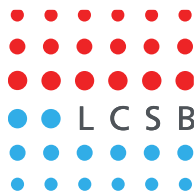




CC BY 4.0

# FAIR Principles for Research Data

*Pinar Alper*



*Training on "Best practices in research data management and stewardship"*

*14 June 2021*

# Definitions

" Research data management (RDM) concerns the organisation of data, from its entry to the research cycle through to the dissemination and archiving of valuable results. It aims to ensure reliable verification of results, and permits new and innovative research built on existing information. "

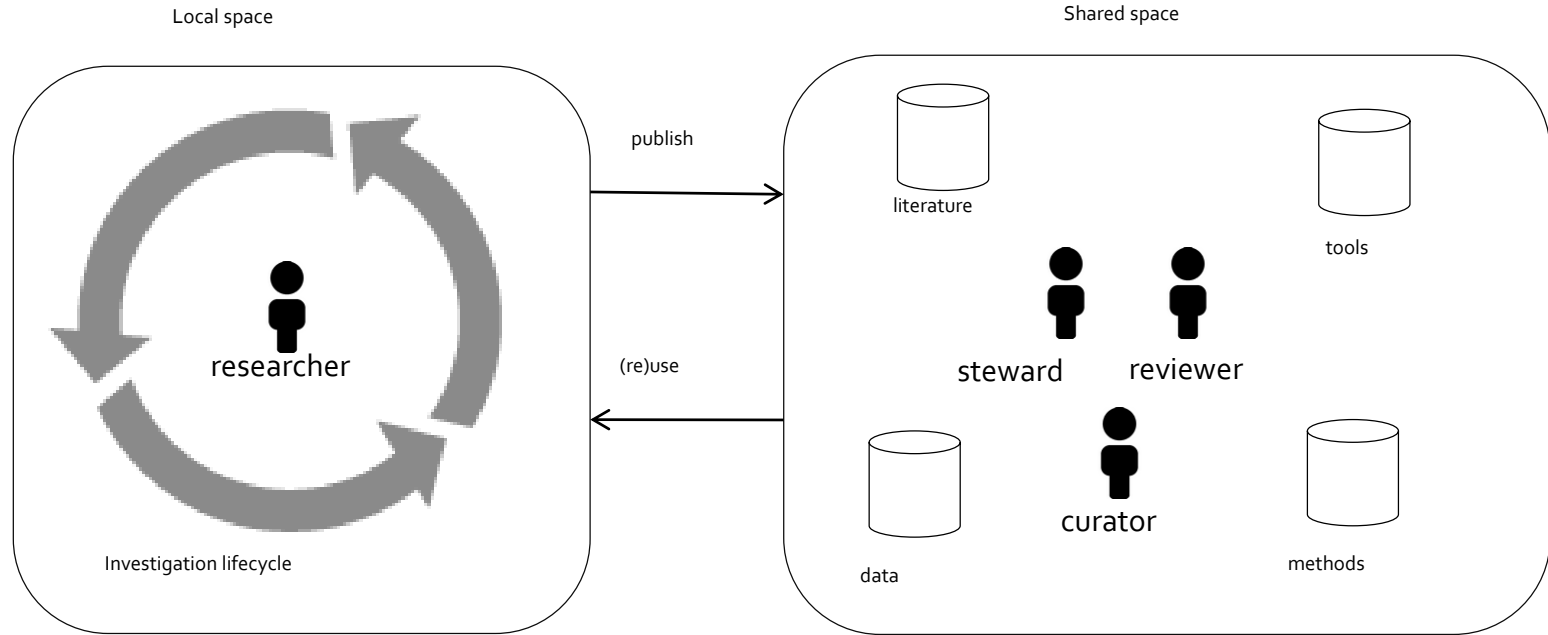
- Examples:
  - Day to day data handling during project , e.g. using consistent file naming conventions.
  - Publishing and sharing after the project completion e.g. depositing the data in a community repository.

# Data Stewardship = RDM ++

" Beyond proper collection, annotation, and archival, data stewardship includes the notion of 'long-term care' of valuable digital assets, with the goal that they should be discovered and re-used for downstream investigations, either alone, or in combination with newly generated data."

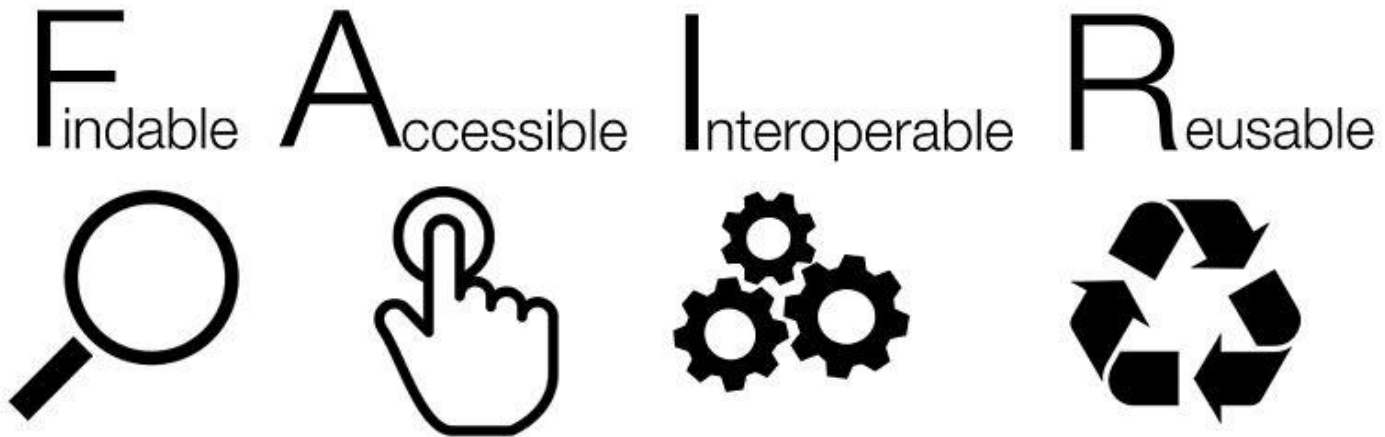
# Data Stewardship

— aims to enable long term care and re-use of data



# End goal

— FAIR principles for research data

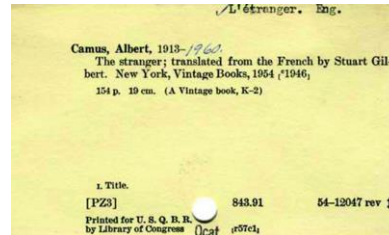


Wilkinson M, Dumontier M et al. Nature Scientific Data 2016. "The FAIR Guiding Principles for scientific data management and stewardship"

# FAIR principles for research data

## Findable

- (Meta)data
- Unique and eternal identifiers for (meta)data
- Indexed in a searchable resource
- Metadata contains data identifier



# FAIR principles for research data

## Findable

- Metadata
- Unique and external identifiers for (meta)data
- Indexed in a searchable resource

Hibsh, D., Schori, H., Efroni, S. & Shefi, O. *Figshare*  
<http://dx.doi.org/10.6084/m9.figshare.1289242> (2015).

NCBI Sequence Read Archive [SRP059260](#) (2015).



### Epi4K: Gene Discovery in 4,000 Epilepsy Genomes

dbGaP Study Accession: phs000653.v3.p1

Study Variables Documents Analyses **Datasets** Molecular Data

#### Dataset Name and Accession

**Dataset Name:** Epi4K\_Subject\_Phenotypes

**Dataset Accession:** pht008354.v1.p1

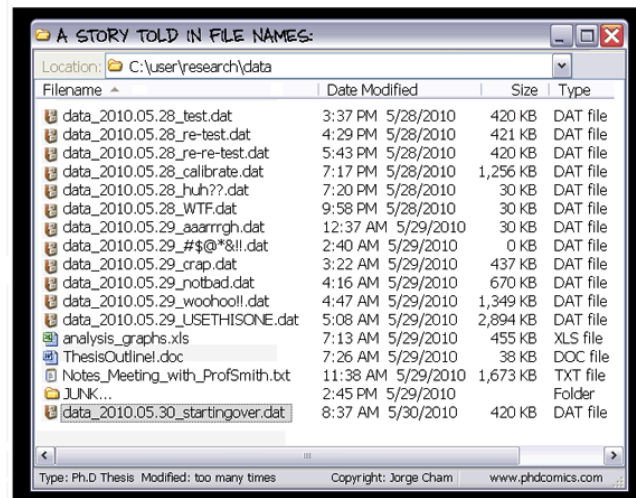
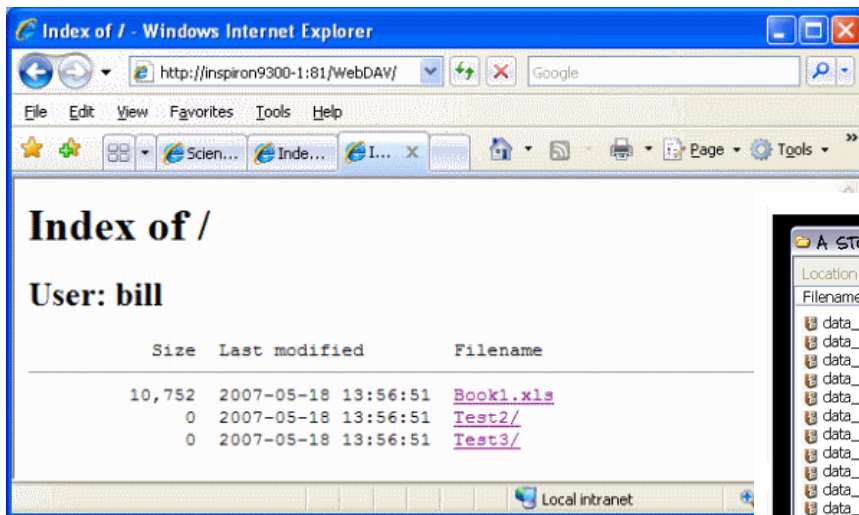
#### Dataset Description

Phenotype Data (All Sub-Studies of phs000653): Includes description of epilepsy category, e.g. infantile spasms, Lennox-Gastaut syndrome, generalized idiopathic epilepsy, focal epilepsies, etc., but also healthy parents, etc, plus gender and ethnicity data.

# FAIR principles for research data

Not (so)Findable

- Metadata
- Identifiers for (meta)data
- Indexed in a searchable resource





# FAIR principles for research data

## Accessible



- (Meta)data are retrievable by a protocol
- Open, free universally implementable
- Authentication/Authorization
- Metadata available even when data is not

Hibsh, D., Schori, H., Efroni, S. & Shefi, O. *Figshare*  
<http://dx.doi.org/10.6084/m9.figshare.128924> (2015).

The image shows the top part of the DOI (Digital Object Identifier) resolution interface. It features the DOI logo (a yellow circle with 'doi' in black) and a navigation bar with links: HOME | HANDBOOK | FACTSHEETS | FAQs | RESOURCES | USERS | NEWS | MEMBERS AREA. Below the navigation bar is a section titled 'Resolve a DOI Name' with a text input field labeled 'doi:' and a 'Go' button.

A DOI is a unique persistent identifier for a published digital object

# FAIR principles for research data

## Accessible

- (Meta)data are retrievable by a protocol
- Open, free universally implementable
- Authentication/Authorization
- Metadata available even when data is not



Moved data

10.1004/123456	URL	<a href="http://www.pub.com/">http://www.pub.com/</a>
	URL	<a href="http://www.pub2.com/">http://www.pub2.com/</a>

Data versions



### The Y4 Seismic Network, 2014–2015

Network code: Y4  
 Restricted: No  
 Network KML file: [K](#)  
 Seismic metadata: [fdsnws-station](#)  
 Institution(s): GFZ  
 Creator(s): Roessler, Dirk; Passarelli, Luigi; Govoni, Aladino; Bautz, Ralf; Dahm, Torsten; Maccaferri, Francesco; Rivalta, Eleonora; Schierjott, Jana; Woith, Heiko  
 Description\*: Extended Pollino Seismic Experiment, GFZ Potsdam (FEFI, CCMP-Pompei, NERA projects)  
 Abstract: The temporary Extended Pollino; The temporary Extended Pollino Seismic Experiment (FDSN network code Y4) monitored the earthquake swarm in the Pollino Range region, Italy, between September 2014 and April 2015.

Type: Temporary  
 Archived at: GFZ  
 Identifier: [doi:10.14470/L9180569](#)  
 (Citation information)  
 DataCite metadata: [HTML](#) | [JSON](#) | [XML](#) | [INSPIRE](#)

Time Range: 2014–2015

Embargoed data

the swarm sequence. Waveform data will be fully open after April 2017. [1] Pollino Seismic Experiment, 2012-2014, doi:10.14470/L9180569

\* Description is taken from seismic metadata, and may not match the preferred title for citations.

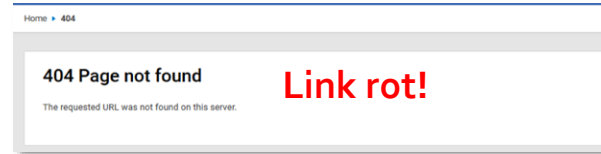
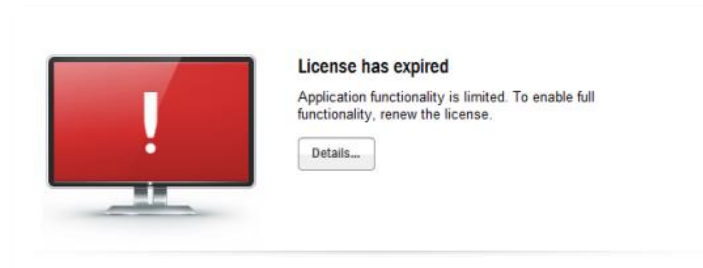
Sponsor(s): CCMP-POMPEI

# FAIR principles for research data

Not so accessible

- (Meta)data are retrievable by a protocol
- Open, free universally implementable
- Authentication/Authorization
- Metadata available even when data is not

Data are available on request due to privacy or other restrictions

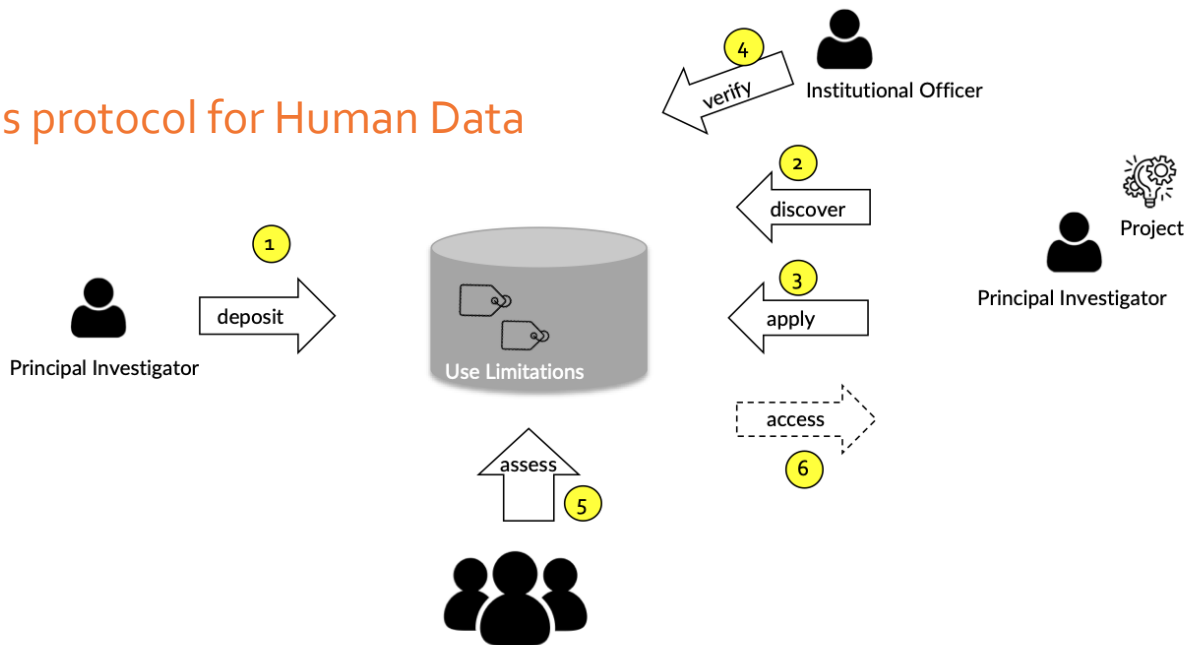
A screenshot of a login form. The title is "Log in to intranet.uni.lux:443". Below the title, it says "Your login information will be sent securely." There are two input fields: "User Name" and "Password". Below the "Password" field, there is a checkbox labeled "Remember this password". At the bottom right, there are two buttons: "Cancel" and "Log In".

# FAIR principles for research data

Accessible  $\neq$  Unrestricted for all

Accessible  $\neq$  Free(  )

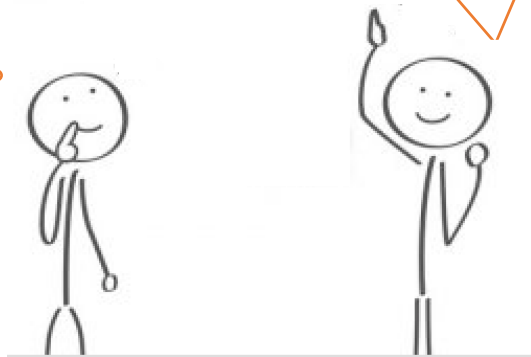
## Access protocol for Human Data



# FAIR principles for research data

Let's apply our technique for the deep learning of **gene mutation**  $\leftarrow \rightarrow$  **drug** relations from literature to **epilepsy**. Here's data from a German cohort, which we can curate.

Good bye  
Christmas break.



# FAIR principles for research data

Let's apply our technique of deep learning of **gene mutation**  $\leftrightarrow$  **drug** relations from literature to **epilepsy**. Here's a German cohort's data we can curate.

ID	Sample	Gender	Race	Vater	Mutter	Birthdate	Disease	Mutation1	Mutation2	Mutation3	Family ID	Comment
101	L1001	male	German	104	103	26/2/1992	Idiopathic Generalized Epilepsy	A37V	F73T	Y301V	K1	AID1029 Available
102	L102	m	German	104	103	29/2/1992	Unaffected	D40V	F73T	Y301V	K1	AID1381 Available
102	L1003	m	German	"	"	"		D124V		D124T	K1	
102	L1003	m	German	"	"	"		E345T	F78I	I98N	K1	
103	L1004	m	Caucasian			12/11/70		E12T	F73T	I98T	K1	AID2738 available: no
104	L1005	male	Caucasian			11-12-70	No	D45H			K1	AID2731 Available
104	L1103										K1	AID2735
106	T1007	f	Asian	107	108	1975-12-10	IGE	C98K	F73T		F2	AID1291 Sample lost
108	T1008	f	Chinese			1972-10-01	Genetic Generalized, Father had GGE	K76V			F2	AID2389
109	L2987	male	Chinese			5.4.1970	Unaffected	V98G	F73T	Y301V	L4	AID3849 Vater von Kevin
110	L4002	male	Russian	Horst	Inga	1/1/1990	Focal epilepsy		F73T		L4	yes
111	L1872	Female	East-German			31.4.1970	Temporal lobe epilepsy	C27G	F73T	Y301V	L4	AID8782 yes



Manual interoperation.



C

`mutation` table

id	gene name	acc no ncbi	nt seq	amplicon start	amplicon end	m1 name	m2 name	allele
[...]	[...]	[...]	[...]	[...]	[...]	[...]	[...]	[...]
274	HvSNAC1	JF796130	ATGGGG...	1	1132	7009	1	.n
275	HvSNAC1	JF796130	ATGGGG...	1	1132	1825	1	.o
[...]	[...]	[...]	[...]	[...]	[...]	[...]	[...]	[...]



homo hetero	nt before	nt position	nt after	localisation	mutation type	aa seq	aa before	aa position	aa after
[...]	[...]	[...]	[...]	[...]	[...]	[...]	[...]	[...]	[...]
homo	G	883	A	coding	silent	MGMPAA...	G	295	G
homo	G	965	A	coding	missense	MGMPAA...	G	276	S
[...]	[...]	[...]	[...]	[...]	[...]	[...]	[...]	[...]	[...]

# FAIR principles for research data

Let's apply our technique of deep learning of **gene mutation**  $\leftrightarrow$  **drug** relations from literature to **epilepsy**. Here's a German cohort's data we can curate.

ID	Sample	Gender	Race	Vater	Mutter	Birthdate	Disease	Mutation1	Mutation2	Mutation3	Family ID	Comment
101	L1001	male	German	104	103	26/2/1992	Idiopathic Generalized Epilepsy	A37V		Y301V	K1	AID1029 Available
102	L102	m	German	104	103	29/2/1992	Unaffected	D45H		Y301V	K1	AID1381 Available
102	L1003	m	German	"	"	"		E345T		D124T	K1	
102	L1003	m	German	"	"	"		E345T		I98N	K1	
103	L1004	m	Caucasian			12/11/70		E12T		I98T	K1	AID2738 available: no
104	L1005	male	Caucasian			11-12-70	No	D45H			K1	AID2731 Available
104	L1103										K1	AID2735 Available
106	T1007	f	Asian	107	108	1975-12-10	IGE	C98K	F73T		F2	AID1291 Sample lost
108	T1008	f	Chinese			1972-10-01	Genetic Generalized, Father had GGE	K76V			F2	AID2389
109	L2987	male	Chinese			5.4.1970	Unaffected	V98G	F73T	Y301V	L4	AID3849 Vater von Kevin
110	L4002	male	Russian	Horst	Inga	1/1/1990	Focal epilepsy		F73T		L4	yes
111	L1872	Female	East-German			31.4.1970	Temporal lobe epilepsy	C27G	F73T	Y301V	L4	AID8782 yes



Automated interoperation.



C

`mutation` table

id	gene	acc no ncbi	nt seq	amplicon start	amplicon end	m1 name	allele
271	HVSNA1	JF796130	ATGGGG...	1	1132		.n
275	HVSNA1	JF796130	ATGGGG...	1	1132	1625	.o
[...]	[...]	[...]	[...]	[...]	[...]	[...]	[...]

homo hetero	nt before	nt position	nt after	localisation	mutation type	aa seq	aa before	aa position	aa after
[...]	[...]	[...]	[...]	[...]	[...]	[...]	[...]	[...]	[...]
homo	G	883	A	coding	silent	MGMPAA...	G	295	G
homo	G	965	A	coding	missense	MGMPAA...	G	276	S
[...]	[...]	[...]	[...]	[...]	[...]	[...]	[...]	[...]	[...]

# FAIR principles for research data

## Interoperable



What machine sees

What we expect to see in Data Integration/Analysis tool

- (Meta)data represented in formal, shared language
- Machine-actionable
- Controlled vocabularies  
Tumour ≠ Tumor
- Community formats & standards  
e.g. CDISC, HL7, ISA Tab...



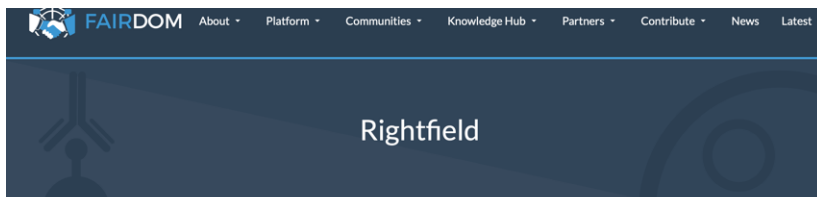
Source	Name	Organism	Age	Unit	Sample Name	Protocol REF	Labeled Extract Name	...	Protocol REF	Data File
H1		H. Sapiens	35	Years	H1.sample1	Labeling	H1.sample1.labeled		Scanning	h1-s1.cel
H1		H. Sapiens	35	Years	H1.sample2				Scanning	h1-s2.cel
H2		H. Sapiens	33	Years	H2.sample1	Labeling	H2.sample1.labeled		Scanning	h2-s1.cel



# FAIR principles for research data

Interoperable

- resources towards achieving interoperable data



Rightfield is an open-source tool for adding ontology term selection to Excel spreadsheets. Rightfield is used by a 'Template Creator' to create semantically aware Excel spreadsheet templates. The Excel templates are then reused by Scientists to collect and annotate their data; without any need to understand, or even be aware of, Rightfield or the ontologies used. Rightfield embedded templates are used within the [Samples](#) framework of the [SEEK](#).



FAIRsharing.org  
standards, databases, policies



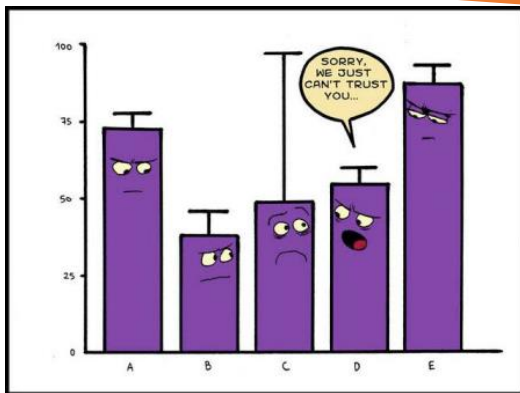
1412 Standards

Terminology Artifact	784
Model/Format	414
Reporting Guideline	169
Identifier Schema	15
FAIR metrics	30

[View all](#)

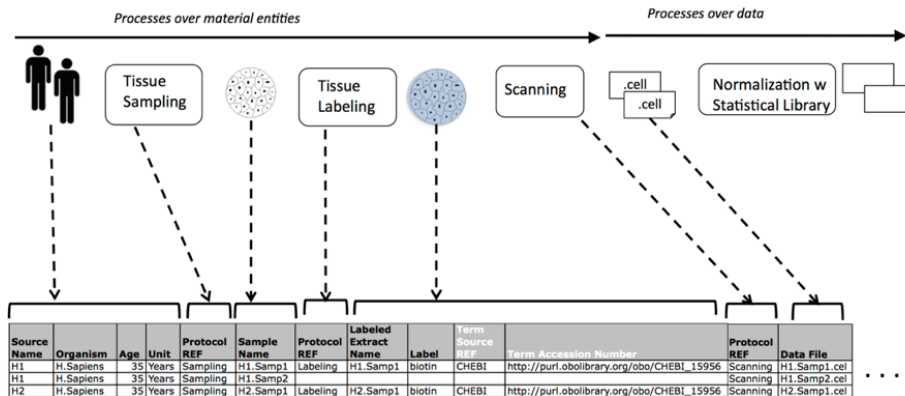
# FAIR principles for research data

## Reusable



- Multitude of metadata attributes
- Following community guideline
- Provenance

- cancer
- lung cancer
- lung cancer, 300 cases 200 controls
- lung cancer, 300 cases 200 controls, phenotyping, epigenetics, protocol A, platform B, ....



Ex-post-facto provenance collection = pain!

# FAIR principles for research data

Reusable

- Descriptive metadata, following community guideline
- Provenance of data
- Clear and accessible data use license

“generalist repositories”

(GIGA)<sup>n</sup>DB

figshare

“domain repositories”

european  
genome-phenome  
archive

dbGaP  
GENOTYPES and PHENOTYPES

More metadata, more transparency, more likelihood of re-use

# Why FAIR data?

## — Political pressure

- Increased push by public funders for maximum use of research results



## H2020

### 3. Open access to research data (Extended Open Research Data Pilot)

#### What?

Beneficiaries of actions that participate in the Open Research Data Pilot must give **open, free-of-charge access** to the end-user to **digital research data** generated during the action (⚠️ new in Horizon 2020).

## NIH

...

# Why FAIR data?

## — Scientific value

- Pooling results, for improving results or new questions
- Validation of models/methods over other data
- Accelerated inter-lab exchange of knowledge



### COVID-19 Disease Map

We aim to establish a knowledge repository of molecular mechanisms of COVID-19 as a broad community-driven effort. Here we share resources and best practices to develop a COVID-19 Disease Map of these mechanisms. The COVID-19 Disease Map is an assembly of molecular interaction diagrams, established based on literature evidence. We focus on host-pathogen interactions specific to the SARS-CoV-2 virus.

[FAIRDOMHub project space](#) lists our contributors, data and computational models, and literature.

### Publication

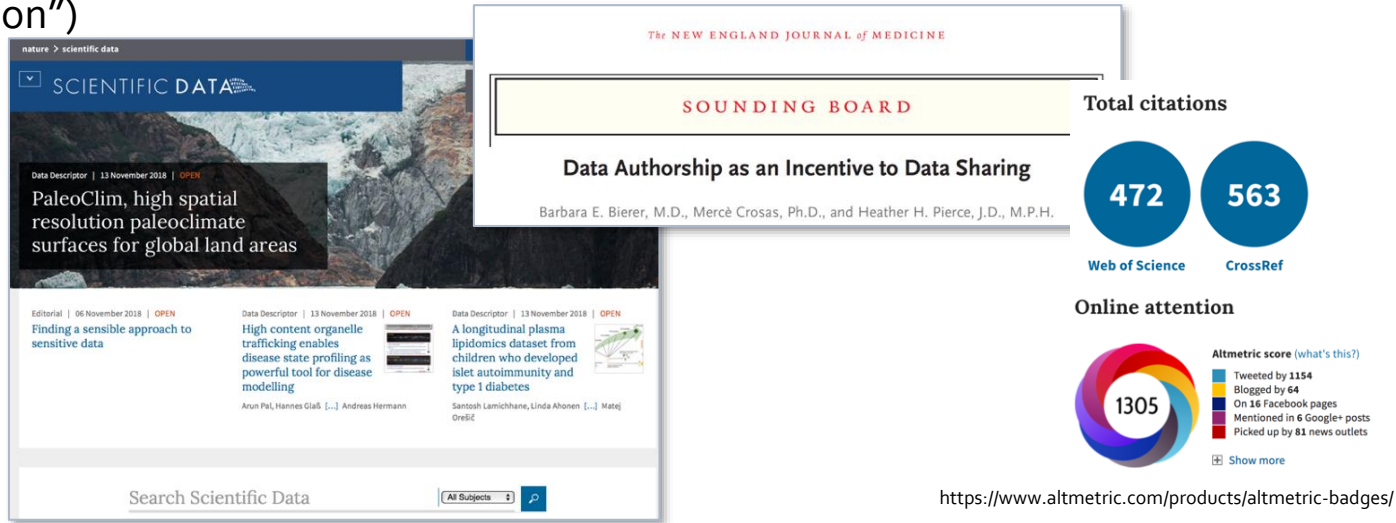
COVID-19 Disease Map, building a computational repository of SARS-CoV-2 virus-host interaction mechanisms. Ostaszewski M, Mazein A, Gillespie ME, Kuperstein I, Niarakis A, Hermjakob H, Pico AR, Willighagen EL, Evelo CT, Hasenauer J, Schreiber F, Dräger A, Demir E, Wolkenhauer O, Furlong LI, Barillot E, Dopazo J, Orta-Resendiz A, Messina F, Valencia A, Funahashi A, Kitano H, Auffray C, Balling R, Schneider R. Sci Data. 2020 May 5;7(1):136. doi: [10.1038/s41597-020-0477-8](https://doi.org/10.1038/s41597-020-0477-8). PubMed PMID: 32371892.

<https://fairdomhub.org/projects/190#models>

# Why FAIR data?

## — New incentives for scientists

- Increased visibility, attracts new collaborations
- Data sharing increases research citation 9%
- FAIR data is being incorporated in the scholarly communication system (“data paper”, “data citation”)



<https://www.nature.com/sdata/>

Piwowar, H. A. and Vision, T. J. et al. Data reuse and the open data citation advantage. PeerJ 2013 Volume 1

# Why FAIR data?

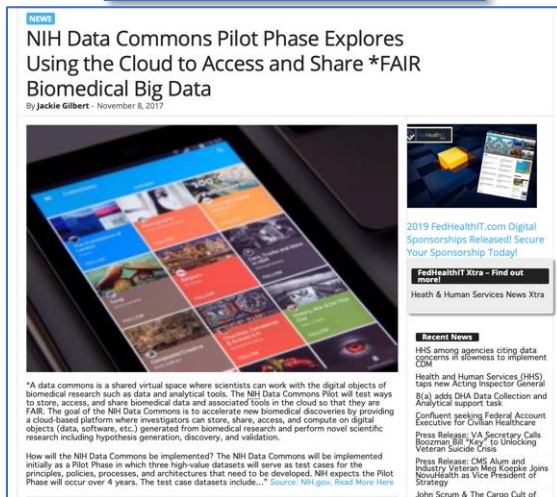
— Improved research quality, reproducibility

- Data + code + documentation

I can't send you the original data because I don't remember what my excel file names mean anymore [#overlyhonestmethods](#)

You can download our code from the URL supplied. Good luck downloading the only postdoc that can get it to run, though [#OverlyHonestMethods](#)

# FAIR is spread across the lands...





EUROPEAN COMMISSION  
DIRECTORATE-GENERAL FOR RESEARCH & INNOVATION

The Director-General

Brussels, 10 July 2017

## EOSC Declaration

RECOGNISING the challenges of data driven research in pursuing excellent science;

GRANTING that the vision of European Open Science is that of a research data commons, widely inclusive of all disciplines and Member States, sustainable in the long-term,

CONFIRMING that the implementation of the EOSC is a process, not a project, by its nature iterative and based on constant learning and mutual alignment;

UPHOLDING that the EOSC Summit marked the beginning and not the end of this process, one based on continuous engagement with scientific stakeholders, the European Commission,

PROPOSES that all EOSC stakeholders consider sharing the following intents and will actively support their implementation in the respective capacities:

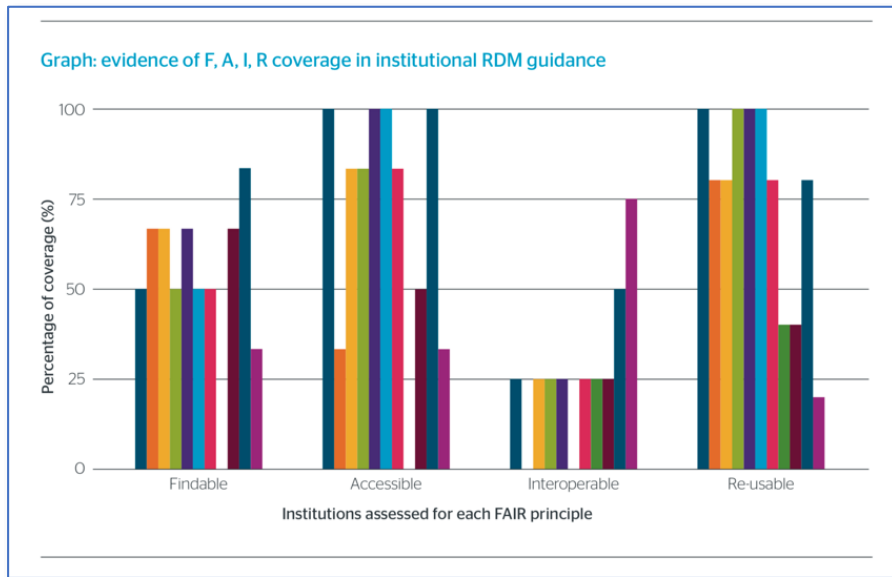
### Data culture and FAIR data

➤ [Data culture] European science must be grounded in a common culture of data stewardship, so that research data is recognised as a significant output of research and is appropriately curated throughout and after the period conducting the research. Only a considerable cultural change will enable long-term reuse for science and for innovation of data created by research activities: no disciplines, institutions or countries must be left behind.



# ...but not necessarily across all the peoples

## Stakeholder FAIR awareness



Government,  
Funder,  
Publisher,  
National &  
International  
Infrastructures...



Institutional



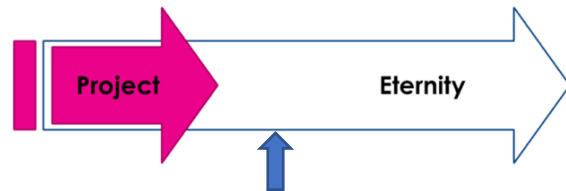
Researchers



.....has evolved into "FAIR fatigue" before "FAIR adoption"

# Achieving FAIR'ness

— Posthoc “FAIRification”



- Applied for datasets of higher value/re-usability potential
- Costly process
- Requires FAIRification experts
- **Assumption:** a percentage of research data will inevitably be non-FAIR at project end.

# FAIRification is an expertise

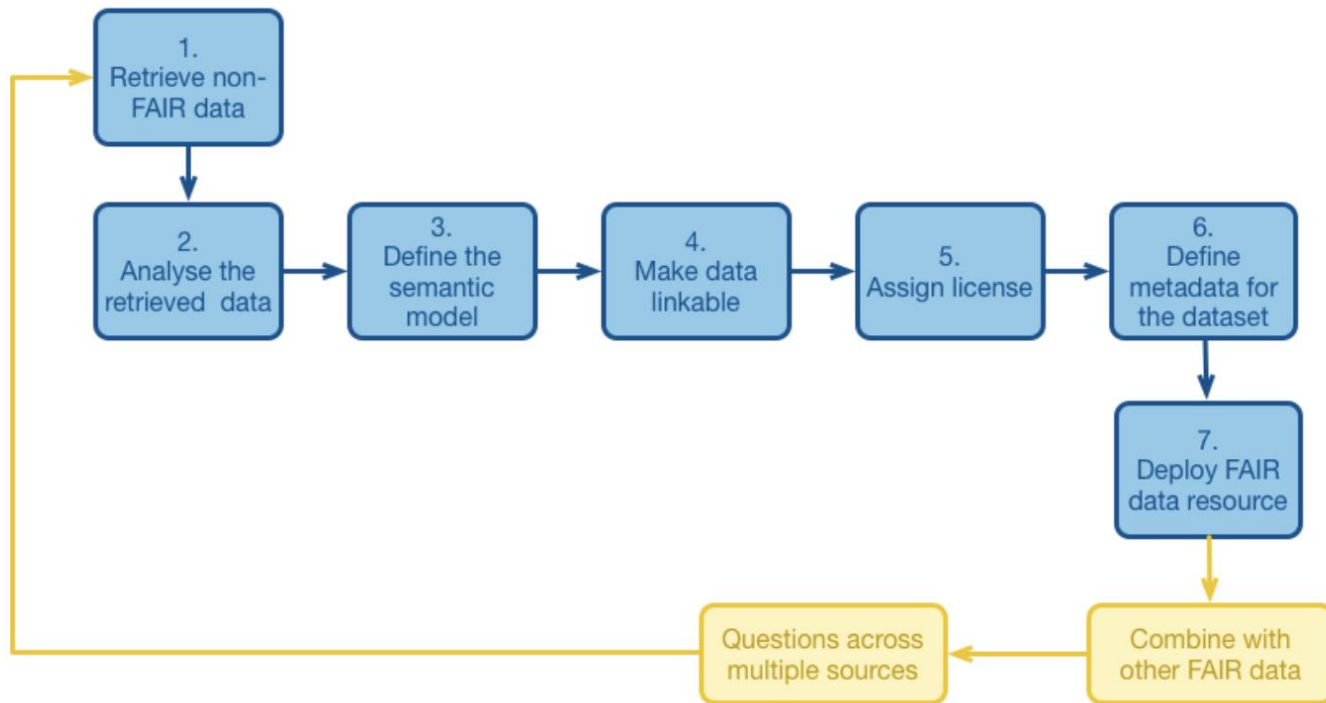
## — FAIRPlus

The screenshot shows the FAIR Cookbook website. The left sidebar contains the 'FAIR Cookbook' logo and a table of contents with sections: 1. FAIR Cookbook, 2. Infrastructure for FAIR, 3. Findability (with sub-items 3.1 to 3.6), 4. Accessibility, and 5. Interoperability. The main content area is titled 'Recipe 1: How to assess FAIRness'. It features several metadata cards: 'Recipe metadata' (identifier: RX.x, version: v0.1), 'Difficulty level' (3 out of 5 flames), 'Reading Time' (15 minutes), 'Recipe Type' (Hands-on), and 'Executable Code' (Yes). A right sidebar titled 'ON THIS PAGE' lists 'INGREDIENTS:', 'OBJECTIVES:', 'STEP BY STEP PROCESS:' (STEP1 to STEP5), 'REFERENCE:', 'AUTHORS:', and 'LICENSE:'. Below the ingredients list, it shows 'Data Managers' and 'Data Scientists'. At the bottom, an 'Ingredients:' section contains a table with 3 rows and 3 columns: Ingredient, Type, and Comment.

Ingredient	Type	Comment
HTTP1.1 protocol	data communication protocol	
guidance on persistent resolvable identifiers	policy	
Persistent Uniform Resource Locators - PURL	redirection service	

# FAIRification is an expertise

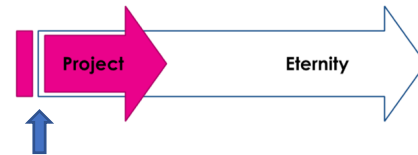
— GO FAIR



<https://www.go-fair.org/go-fair-initiative/>

# Achieving FAIR'ness

— Data that born FAIR



- Use research infrastructures, FAIR tools, standards and practices from Day 1
- **Assumption:** FAIR can only be achieved at scale by good data management practice.



# Good RDM practice during the entire project lifecycle

## ELIXIR RDMkit

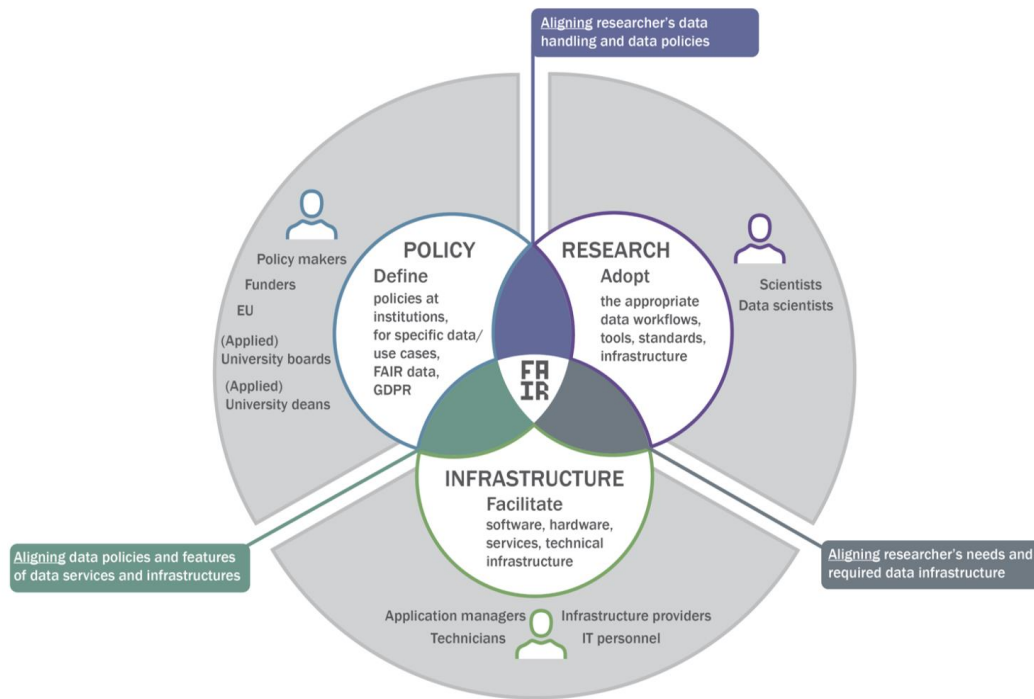


RDMkit

<https://rdmkit.elixir-europe.org>

# Professionalizing RDM Support

— “Professionalising data stewardship” in the Netherlands.



“Professionalising data stewardship in the Netherlands. Competences, training and education. Dutch roadmap towards national implementation of FAIR data stewardship” M Jetten et al [10.5281/zenodo.4320504](https://doi.org/10.5281/zenodo.4320504).

# Conclusions

97 %



- Data is a valuable research output
- You can optimize its value by performing RDM
- Managed data is FAIR!





# Thank you