

Communication / Dissemination / Exploitation

Communication: Promote your action and results

Inform, promote and communicate your activities and results



Citizens, the media, stakeholders



- Having a well-designed strategy
- · Conveying clear messages
- · Using the right media channels



From the start of the action until the end



- Engage with stakeholders
- · Attract the best experts to your team
- · Generate market demand
- · Raise awareness of how public money is spent
- · Show the success of European collaboration
- . Legal obligation: Article 38.1 of the Grant Agreement

Dissemination: Make your results public

Open Science: knowledge and results (free of charge) for others to use



Only to scientists?

Not only but also to others that can learn from the results: authorities, industry, policymakers, sectors of interest, civil society



A How?

Publishing your results on:

- · Scientific magazines
- · Scientific and/or targeted conferences
- Databases



At any time, and as soon as the action has results



- · Maximise results' impact
- · Allow other researchers to go a step forward
- · Contribute to the advancement of the state of the art
- · Make scientific results a common good
- Legal obligation: Article 29 of the Grant Agreement

Exploitation: Make concrete use of results

Commercial, Societal, Political Purposes



Only by researchers?

Not only, but also:

- · Industry including SMEs
- · Those that can make good use of them: authorities, industrial authorities, policymakers, sectors of interest, civil society



- · Creating roadmaps, prototypes, softwares
- · Sharing knowledge, skills, data



Towards the end and beyond, as soon as the action has exploitable results



- · Lead to new legislation or recommendations
- · For the benefit of innovation, the economy and the society
- Help to tackle a problem and respond to an existing demand Legal obligation: Article 28 of the Grant Agreement

EPFL Importance of sharing (!)

1976 – Experiments on supercooled water (cooled far below its freezing point) showed a **critical point** at −20°C: its structure fluctuates widely between high- and low-density forms

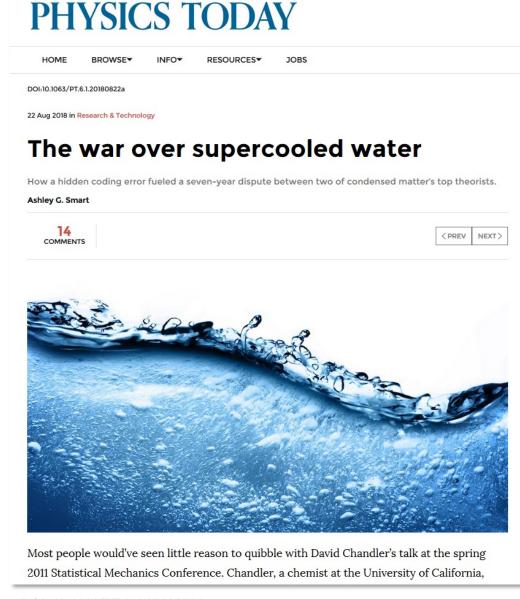
2011 – Seeking a unified theory of water, simulations on supercooled water by two world-leading groups revealed:

- Chandler et al.: no critical point (resembles ordinary water)
- Debenedetti et al.: critical point (morphs between two forms)

2014 – Debenedetti et al. **published their code** openly

2016 – At first, Chandler et al. only shared data, then revealed where to find its code and, after lot of reverse engineering ...

2018 – ... the trouble stemmed from an algorithmic trick the Chandler's team used to speed up their code!



DOI: 10.1063/PT.6.1.20180822a

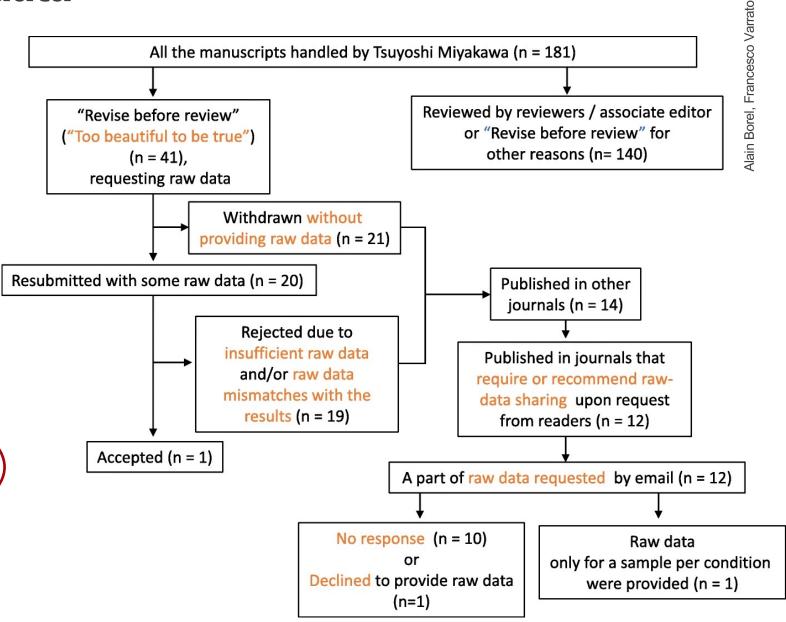
EPFL Importance of raw data

EDITORIAL

No raw data, no science: another possible source of the reproducibility crisis

Tsuyoshi Miyakawa

- Lack of raw data: another possible cause of irreproducibility
- Many researchers did not provide the raw data
- Data fabrication: raw data may not even exist in some cases
- Good faith: the insufficiency or mismatch between raw data and results can be honest mistakes
- Systematic review and metaanalysis: estimated that 1.97% of authors admitted to have **fabricated**, **falsified**, **or modified data** or results at least once [...] the admission rate was 14.12% for falsification when asked about the colleagues



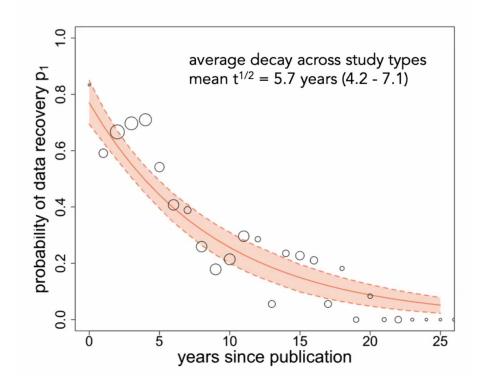
EDCH 2024 / Hands-on with Research Data Management in Chemistry

Importance of raw data

EDITORIAL

No raw data, no science: another possible source of the reproducibility crisis

Tsuyoshi Miyakawa



Data & analysis script availability (prevalence estimates)

	Data	Analysis scripts
Psychology (2014-2017) ¹	2% [1-4%]*	1% [0-1%]
Social Sciences (2014-2017) ²	7% [2-13%]	1% [0-3%]

¹Hardwicke et al. (2021)

*[95% confidence intervals]

²Hardwicke et al. (2020)

Data availability on request (selected studies)

141 articles published in four major APA journals (2004)³ 27%

516 ecology articles published (1991-2011)⁴ 20%

111 most highly-cited psychology & psychiatry articles (2006-2016)⁵

³Wicherts et al. (2006)

⁴Vines et al. (2014)

5Hardwicke & Ioannidis (2018)

Bibliothèque

Data shared

14%

Open Data Decision Tree

EDCH 2024 / Hands-on with Research Data Management in Chemistry



No Has a description of the data been Is it sensitive data that cannot be published or the data refereed? openly shared? Yes No No Yes **IDEAL** Were permissions obtained to share the data within regulatory requirements? Yes No Outcome explanation Sources: QR Code generator library: Project Nayuki doi.org/10.5281/zenodo.581415

Open Data: Can I make my data open?

Version: 0.4

Authors: Scholarly Commons Subworking

group 3

Yes

Can the data be digitized? <

No

Does the data you wish to

share already exist?

Yes

Do you have the rights to

No

See: Making data open by design

> Is the data in a proprietary (closed)

> > format?

Yes

Yes

Are the data digital?

No

FORCE11

Yes

Can the data be

converted into an

open format?



Legal constraints for open data publication

- Tests on animals / humans
- Handle personal data
 - Federal Act on Data Protection (FADP), Human Research Act (HRA), GDPR
 - name, identification number, location data, online identifier, ...
 - factors specific to physical, physiological, genetic, mental, economic, cultural or social identity
- → check the EPFL <u>Human Research Ethics Committee</u> (AREC + HREC)
 - **Data from 3rd party** sources? (e.g. commercial datasets, research cooperations, etc.)
- → check out the **contract** for data usage / sharing ... Or make one!
- Want to potentially submit a **patent**?
- → check the TTO (Technology Transfer Office) ... Choose the data license + tell in the DMP!



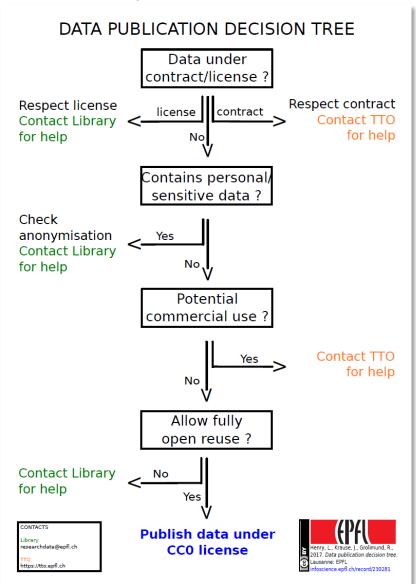
Live poll: Who owns my data/code?

- ☐ Me
- Thesis supervisor / PI
- Publisher / Platform (once published)
- □ 3rd party (obtained from provider)
- EPFL (ETH Domain)
- I don't know

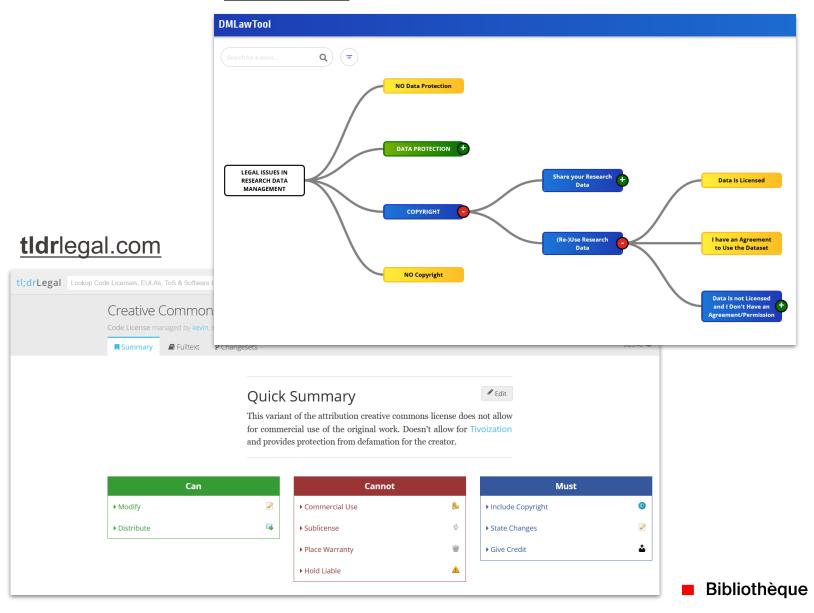
NOTE: ownership ≠ authorship ≠ licensee

Potential commercial use?

infoscience.epfl.ch/record/230281



DMLawTool



Let's find out [5']

Check out these licenses:

CC0

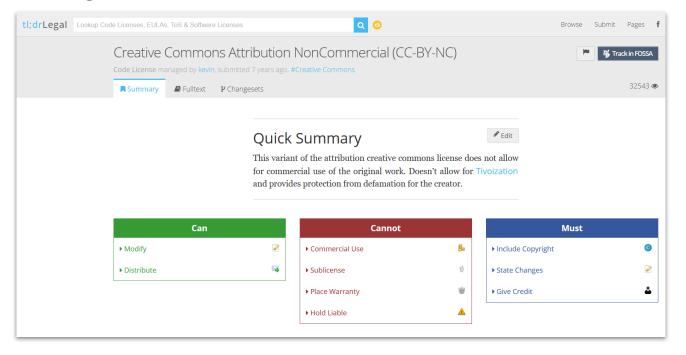
CC BY-4.0

MIT

CC BY-NC-ND



tldrlegal.com





Data / Code licences

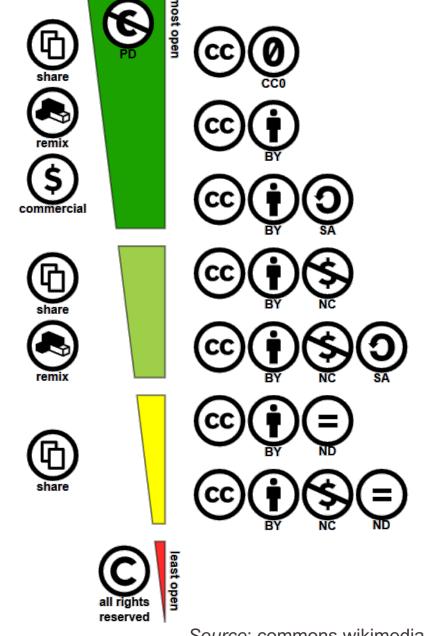
COMMONS LICENSES FOR DATA, TEXT & MULTIMEDIA

Creative Commons:

- Enforced by the author
- Check platform's policy
- On datasets (no data points)

The **Open Data Commons** can be a viable option

The 96|9|EC Directive protects only vs. "substantial" copies of datasets



Source: commons.wikimedia.org

EPFL

Data / Code licences

MORE SPECIFIC FOR CODE





- Apache2.0 (smaller codes, libraries)
 - Permissive
 - No share-alike clause
 - Preservation of copyright notice
- BSD-3clause Similar
- MIT License GPL compatible







Importance of licensing (Yes, again!)

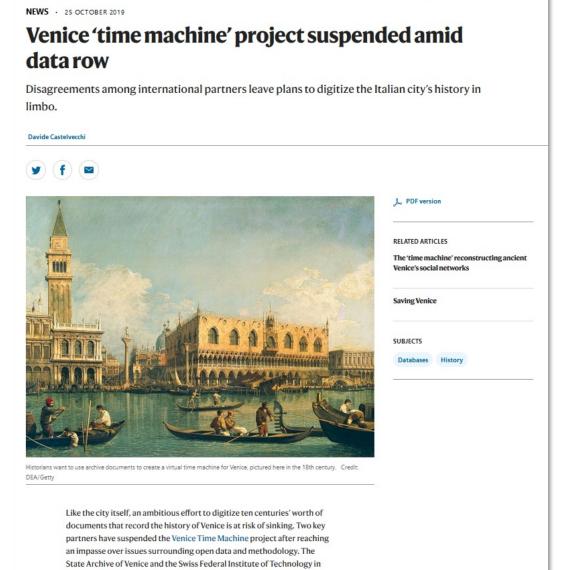
2012 – Project of officially **launched**: Venice's State Archive + Ca' Foscari Univ. + EPFL (DHLAB)

2014 – Non-binding agreement signed. But ... **didn't specify the licensing** that would regulate researchers' use of the digitized data

2017 – At stake: 1,000 years of records in dynamic digital form: special high-speed scanners, thousands HD images per hour

2019 – Allegedly, the digitization of ~190,000 documents (**8 TB**) didn't follow a common metadata policy: <u>archival-science guidelines</u> (require records of provenance for each document)

2019 – ... data collection has been paused, amid doubts on the usability of the data already collected!



Lausanne (EPFL) say they have had to pause data collection, and the archive's director has raised questions about the usability of the 8

nature

DOI: 10.1038/d41586-019-03240-w

Discussion: sensitive data in research?

- Do you collect, process or store data which is... sensitive?
- Do you collect, process or store information on... identifiable persons?
- How do you inform persons/subjects on what you will be doing?

Discussion [5']

Discussion: sensitive data in research?

Personal data

Information that relates to an identified or identifiable individual.

Not only data that can directly identify a person (e.g., name) is considered personal data, but also data that can make a person identifiable through the combination of data (e.g., combining age, e-mail address and information related to usage of social networks may allow identifying a person).

Consider that information together with the means reasonably likely to be used by either you or any other person to identify that individual.

Sensitive data

Information that includes but is not limited to religious, ideological, political or trade union-related views or activities; **health, genetic, biometric**, or concerning the intimate sphere or the racial origin; ethnic data or social security measures; administrative or criminal proceedings and sanctions.

See also Art. 3.c of the FADP Swiss law.

EPFL

Personal / Sensitive data processing





Any operation with personal data [...] in particular

- the collection
- storage
- use
- revision
- disclosure
- archiving
- or destruction of data

Swiss Federal Act on Data Protection (FADP) (Loi sur la Protection des Données LPD), Art. 5

Protection: Personal data must be protected against unauthorised processing through adequate technical and organisational measures

FADP Art. 7-9

Disclosure: Making personal data accessible, for example:

- by permitting access
- transmission
- or publication

FADP Art. 5e

EPFL Co

Collecting consent (online & offline)

Art. 7 Consent

- ¹ Research involving human beings may only be carried out if, in accordance with the provisions of this Act, the persons concerned have given their informed consent or, after being duly informed, have not exercised their right to dissent.
- ² The persons concerned may withhold or revoke their consent at any time, without stating their reasons.

Human Research Act, Art. 7

The consent must be:

- Simple
- Understandable
- Adapted to the subject (child, teenager...)

HRA, Art. 21-22

EPFL

HREC Review Procedures

Standard Procedure

- Projects with higher risks
- For example, involvement of vulnerable participants, sensitive data, physical risks, data protection risks

Simplified Procedure

- Projects with lower risks For example, collection of non-sensitive personal data, testing of prototype
- Modification of authorized research projects, if they raise minor/ specific ethical, scientific or legal issues
- Sub-projects covered by an already approved general protocol

Presidential Decision

- Further use of data obtained with informed consent
- Research projects that do not raise specific ethical, scientific or legal issues

Source: three different review procedures on ReO webpage, https://drive.google.com/file/d/104o9imHvJ3tLp6c8gtQPZV7qvZJBPYuN/view

EPFL

Data masking techniques

Pseudonymization

(working data, reversible)



PSEUDONYMIZATION

Replace data by identifiers. The key is kept separately & securely

ENCRYPTION

Encrypt the data & keep the key secure. Also for long-term preservation, not data publishing

Some tools:

R package: <u>sdcMicro</u>

Java application: ARX Data Anonymization Tool

Java application: <u>ARGUS</u>

Platform: Amnesia

Anonymization

(published data, irreversible)



GENERALIZATION

Diminish granularity by generalizing the variables. Appropriate for data too specific or unique records

SUPPRESSION

Suppress data or part of the outlier records. Appropriate for processing identifiers

ADD FAKE DATA

To prevent the identification of specific records, add fake data while preserving correlations

SHUFFLE

Shuffle data over one / several columns without compromising the utility of the data

(Other: Differential Privacy, T-closeness, ...)

Images:

- https://www.flaticon.com/packs/general
- https://www.flaticon.com/packs/hawcons-documents-filled

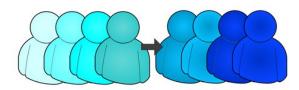
Try it out later!

Amnesia

Amnesia is a data anonymization tool, that allows to remove identifying information from data. Amnesia not only removes direct identifiers like names, SSNs etc but also transforms secondary identifiers like birth date and zip code so that individuals cannot be identified in the data. Amnesia supports k-anonymity and k^m -anonymity.

version: 1.0.7 (release date: 21/02/2019)









- Implements data anonymization techniques from the field of Privacy Preserving Data Publishing (PPDP)
- Transforms original data to anonymized data by using generalization and suppression
- Anonymization not limited to the removal of direct identifiers; it also includes removing secondary information (e.g. like age, zipcode, etc.) that might indirectly lead to identify an individual
- Focuses on *k*-anonymity: guarantees that every record will be indistinguishable from other k-1 records
- Supports 2 algorithms for k-anonymity, <u>Incognito</u> and a parallel version of the <u>Flash</u> <u>algorithm</u>.



Deletion of identifying data

name	gender	city	age	disease
KELLER Anna	f	Basel	32	no diabetes
BRUNNER Emilia	f	Basel	37	diabetes 2
DURANT Pierre	f	Basel	44	no diabetes
GRAF Julia	f	Basel	45	diabetes 2
GERBER Fritz	m	Basel	20	diabetes 1
FISCHER Urs	m	Basel	23	diabetes 1
WYSS Emilien	m	Geneva	24	no diabetes
STEINER Leo	m	Geneva	28	no diabetes
ROTH Christian	m	Geneva	42	no diabetes
WYSS Rudolf	m	Geneva	48	diabetes 2

	name	gender	city	age	disease
0	*	f	Basel	30 - 39	no diabetes
1	*	f	Basel	30 - 39	diabetes 2
2	*	f	Basel	40 - 49	no diabetes
3	*	f	Basel	40 - 49	diabetes 2
4	*	m	Basel	20 - 29	diabetes 1
5	*	m	Basel	20 - 29	diabetes 1
6	*	m	Geneva	20 - 29	no diabetes
7	*	m	Geneva	20 - 29	no diabetes
8	*	m	Geneva	40 - 49	no diabetes
9	*	m	Geneva	40 - 49	diabetes 2

K-anonymity 2



Deletion of identifying data

name	gender	city	age	disease
KELLER Anna	f	Basel	32	no diabetes
BRUNNER Emilia	f	Basel	37	diabetes 2
DURANT Pierre	f	Basel	44	no diabetes
GRAF Julia	f	Basel	45	diabetes 2
GERBER Fritz	m	Basel	20	diabetes 1
FISCHER Urs	m	Basel	23	diabetes 1
WYSS Emilien	m	Geneva	24	no diabetes
STEINER Leo	m	Geneva	28	no diabetes
ROTH Christian	m	Geneva	42	no diabetes
WYSS Rudolf	m	Geneva	48	diabetes 2

	name	gender	city	age	disease
0	*	f	*	30 - 39	no diabetes
1	*	f	*	30 - 39	diabetes 2
2	*	f	*	40 - 49	no diabetes
3	*	f	*	40 - 49	diabetes 2
4	*	m	*	20 - 29	diabetes 1
5	*	m	*	20 - 29	diabetes 1
6	*	m	*	20 - 29	no diabetes
7	*	m	*	20 - 29	no diabetes
8	*	m	*	40 - 49	no diabetes
9	*	m	*	40 - 49	diabetes 2

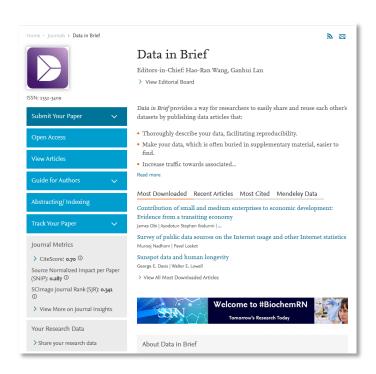
L-diversity 2

EPFL

Data publication: journal pathway

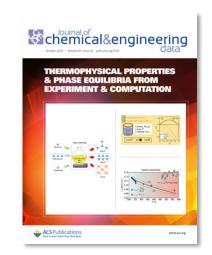
Data Papers

A data paper is a peer reviewed document describing a dataset, published in a peer reviewed journal. It takes effort to prepare, curate and describe data (GBIF, 2019)



Data Journals

Data papers are supported by many journals, some of which are "pure", i.e. they are dedicated to publish data papers only, while others – the majority – are "mixed", i.e. they publish a number of articles types including data papers. (Wikipedia, 01.04.2019)







EPFL

Data / Code journals

go.epfl.ch/datajournals













etc.

JOURNAL	OA TYPE	FOCUS	DISCIPLINE(S)	APC (* = Waiver policy)	EPFL LIBRARY AGREEMENT
Computational Engineering and Physical Modeling (Pouyan Press, Iran)	Diamond	Code	Electrical engineering, Electronics, Nuclear engineering, Computer engineering	<u>o</u>	
<u>Digital Humanities Quarterly</u> (Alliance of Digital Humanities Organizations, Netherlands)	Diamond	Code	Humanities, Philology, Linguistics, Mass media	<u>o</u>	
Earth system science data (Copernicus Publications, Germany)	Diamond	Data	Environmental science, Earth science, Geology	<u>o</u>	
Image Processing On Line (Image Processing On Line, France)	Diamond	Code	Mathematics, Computer science	<u>o</u>	
International Journal of Data and Network Science (Growing Science, Canada)	Diamond	Both	Social Sciences, Management, Industrial management	<u>o</u>	
Journal of Data and Information Science (Sciendo, Poland)	Diamond	Code	Technology, Industrial engineering, Management engineering, Information technology, Mathematics, Computer science	<u>o</u>	
Journal of data science (School of Statistics, Renmin University of China, China)	Diamond	Code	Mathematics, Computer science, Mathematical statistics	<u>0</u>	
<u>Journal of Statistics and Data</u> <u>Science Education</u> (Taylor & Francis Group, USA)	Diamond	Code	Mathematics, Mathematical statistics, Education	<u>o</u>	
Journal of Open Source Software, (Journal of Open Source Software,	Diamond	Code	Science, Computer Science	<u>0</u>	

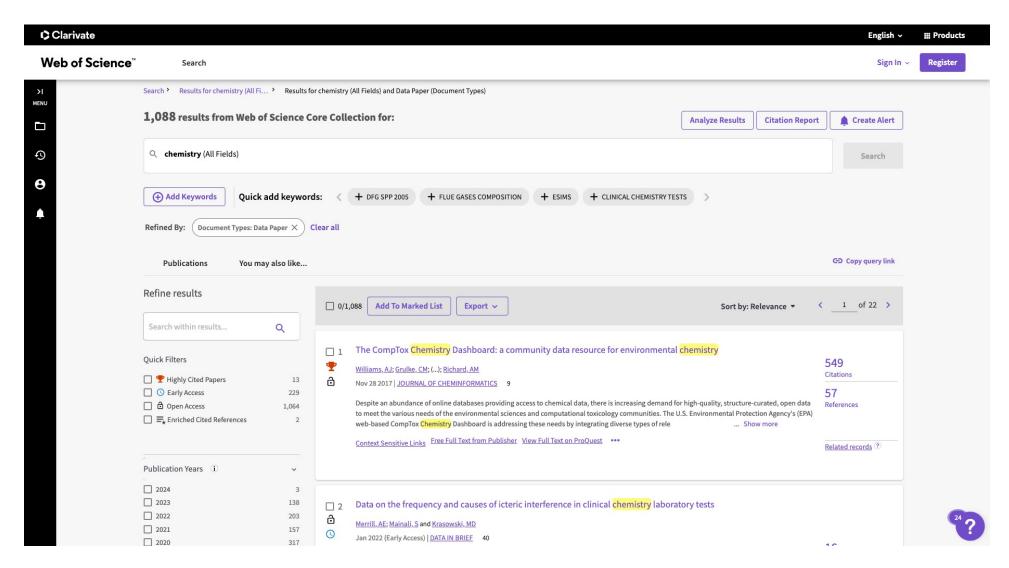
Data publication

Data papers can be highly cited

3,667 documents found ✓ Analyze results ¬							
All	✓ Export ✓ Download Citation overview ••• More	Show all abstracts	Sort by Cited by (high	hest) 🗸	⊞ ⊨		
	Document title	Authors	Source	Year	Citations		
	Data Paper • Open access ERA5-Land: A state-of-the-art global reanalysis dataset for land applications	Muñoz-Sabater, J., Dutra, E., Agustí-Panareda, A., Buontempo, C., Thépaut, JN.	Earth System Science Data, 13(9), pp. 4349– 4383	2021	987		
	Show abstract ✓ Full text at EPFL Library View at Publisher	¬ Related documents					
	Data Paper • Open access						
2	Global Carbon Budget 2021	Friedlingstein, P., Jones, M.W., O'Sullivan, M., Zaehle, S., Zeng, J.	Earth System Science Data, 14(4), pp. 1917– 2005	2022	590		
	Show abstract ✓ Full text at EPFL Library View at Publisher	¬ Related documents					
3	Data Paper • Open access China CO ₂ emission accounts 2016–2017	Shan, Y., Huang, Q., Guan, D., Hubacek, K.	Scientific Data, 7(1), 54	2020	513		
	Show abstract ✓ Full text at EPFL Library View at Publisher	¬ Related documents					

Data publication

Data papers can be highly cited



Data repositories: Publication and/or Preservation

Search here						
PLATFORM (*= Institutional)	TYPE(S) (* = For-Profit)	RECOMMENDED BY SNSF/EU	DISCIPLINE(S)	* HOSTING	Ĉ DOI	
ACOUA* EPFL	Archive		All	CH	1	10 TB
<u>ArrayExpress</u>	Data Repository	SNSF, EU	Genetics, Biology, Life Sciences	USA/EU/UK	0	N/A
BORIS Portal*	Archive, Data Repository	SNSF	All	СН	1	No limit
c4Science* EPFL	Code Repository		All	СН	0	N/A
CERN Open Data*	Data Repository		High Energy Physics, Condensed Matter Physics, Physics	EU	1	N/A
Channelpedia EPFL	Databank		Electrophysiology, Physiology	СН	0	N/A
Copernicus	Data Repository		Geosciences, Ecology, Atmospheric Science	EU	0	N/A
<u>DaSCH</u>	Data Repository	SNSF	Humanities, Social Sciences, Linguistics	СН	0	N/A
<u>dbGap</u>	Data Repository		Genetics	USA	0	N/A

go.epfl.ch/datarepo

- + SNSF Open Data criteria
- + Europe-approved repositories
- + List of Nature's Recommended Data Repositories per discipline
- + (Non-exhaustive) list of repositories approved by some publishers for hosting data alongside the articles



Which data repository? Try re3data [5']

Apply filters: subj. (Chemistry) & country (EU + CH)

Found 3 result(s)	
CARIBIC Civil Aircraft for the Regular Investigation of the	e atmosphere Based on an Instrument Container
Subject(s)	Atmospheric Science Geodesy, Photogrammetry, Remote Sensing, Geoinformatics, Cartogaphy Analytical Chemistry, Method Development Geosciences (Including Geography) Geophysics and Geodesy Chemistry
Content type(s)	Plain text Raw data Scientific and statistical data formats other
Country	Germany United Kingdom Switzerland France European Union Netherlands Sweden
	tudy and monitor important chemical and physical processes in the Earth's atmosphere. Detailed and extensive measurements are made air and particle (aerosol) inlet underneath the aircraft. We use an Airbus A340-600 from Lufthansa since December 2004.
Network for the Detection of Atmo	ospheric Composition Change
Subject(s)	Geosciences (including Geography) Particles, Nuclei and Fields Geodesy, Photogrammetry, Remote Sensing, Geoinformatics, Cartogaphy Atmospheric Science and Oceanography Natural Sciences Physics Geophysics and Geodesy Chemistry
Content type(s)	Plain text other
Country	United States Germany Belgium Switzerland European Union International
	Composition Change (NDACC), a major contributor to the worldwide atmospheric research effort, consists of a set of globally distributed re UV radiation reaching the Earth's surface, and physical parameters, centered around the following priorities.
Rhea	
Subject(s)	Basic Biological and Medical Research Bioinformatics and Theoretical Biology Metabolism, Biochemistry and Genetics of Microorganisms Microbiology, Virology and Immunology Medicine Chemistry Natural Sciences
Content type(s)	Standard office documents Plain text Scientific and statistical data formats Structured graphics Structured text Software application
Country	European Union Switzerland
and metabolic network reconstruction. There a Biological Interest) which provides detailed info	resource of expert-curated biochemical reactions. It has been designed to provide a non-redundant set of chemical transformations for appare three types of reaction participants (reactants and products): Small molecules, Rhea polymers, Generic compounds. All three types of remation about structure, formula and charge. Rhea provides built-in validations that ensure both mass and charge balance of the reaction latabases), extending it with additional known reactions of biological interest. While the main focus of Rhea is enzyme-catalysed reactions





SNSF recommends using re3data.org vertical search engine

How to choose a data repository

- 1. Listed on re3data: for peace of mind
- 2. DOI or other Persistent IDentifier
- 3. Non-profit: SNSF doesn't reimburse
- 4. Good licenses choice: reuse & compliancy
- **5.** Cross-linking: dataset / code ↔ article
- 6. Target public: field-specific and/or generic
- 7. Max upload: size matters

GitHub is **not** a data repository Your own website is **not** a data repository You can choose both a GENERIC & field-SPECIFIC data repository

EPFL

Data repositories (Example)

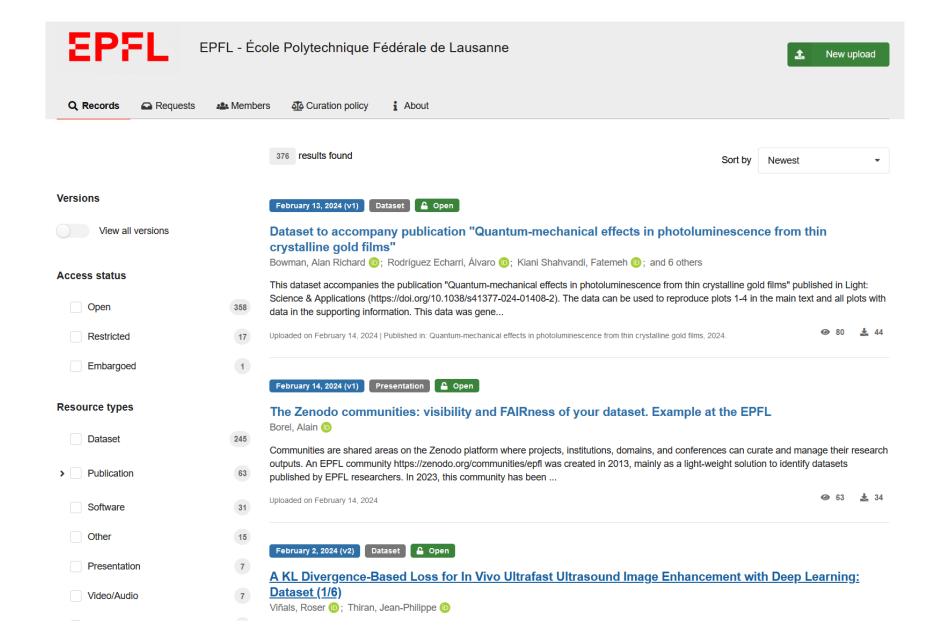
- Hosted by the CERN
- Free of charges
- Max 50GB/dataset
- Unlimited datasets
- Automated DOI assignment

- OpenAIRE integration (EC reporting)
- GitHub integration
- ORCID integration
- All file formats accepted
- Usage statistics interface
- OAI-PMH protocol (content harvesting)
- 18 petabytes disk cluster
- Each file has 2 replicas on different servers
- 2 independent MD5 checksums per file
- Metadata 12-hourly backup cycle
- ...

Zenodo.org has an EPFL Community!

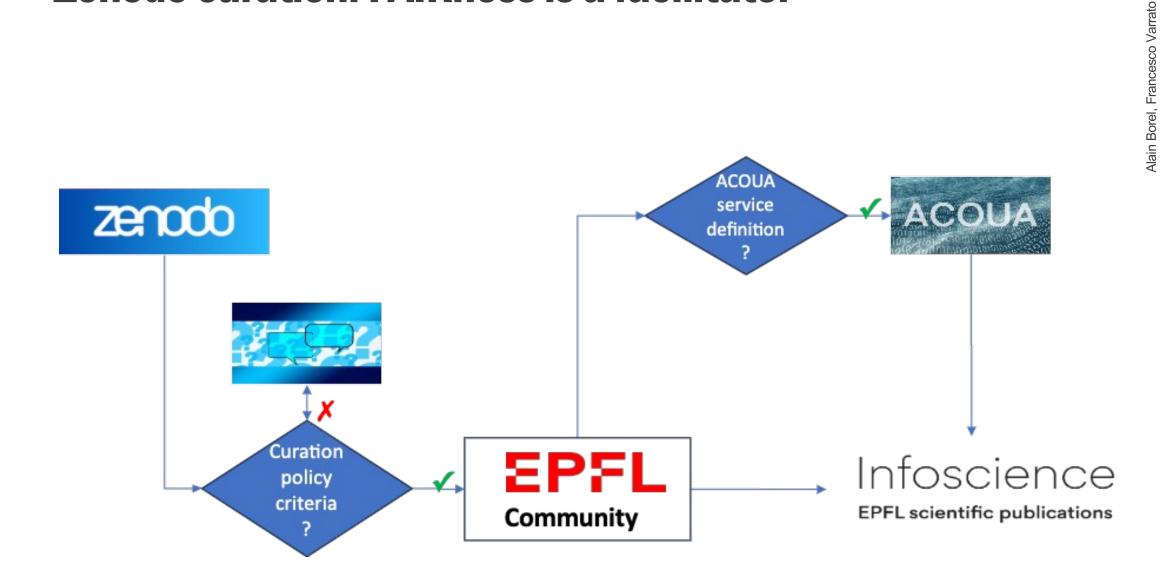
EPFL

Data repositories (Example)





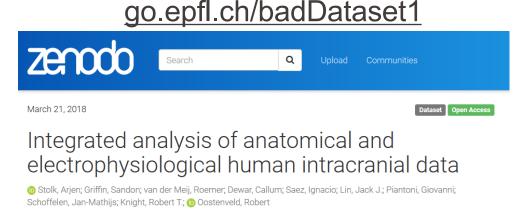
Zenodo curation: FAIRness is a facilitator



EPFL

Discussion: Publication (bad dataset)

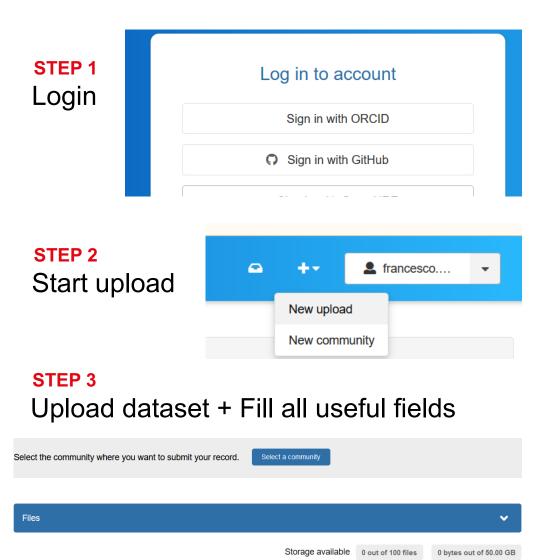
Look at this dataset ...



- 1. Do you think that it's reusable under a ...
 - ... legal standpoint? (license, sensitive data, ...)
 - ... technical standpoint? (formats, interoperability, code, ...)
 - scientific standpoint? (data quality, documentation, ...)



Hands-on: **Zenodo** sandbox



- or -

Upload files

Drag and drop files



You can use a dataset from Day 2 on Moodle

Preservation ≠ **Publication**

ACOUA: Long-term preservation

Why using it

- archive the entire dataset underlying a publication
- archive datasets of a finished research project
- archive datasets of a collaborator leaving EPFL
- get space for large datasets that need preservation
- preserve raw data useful during your research
- get expert support for data curation

What you get

- trustworthy, safe and EPFL-backed environment
- free for EPFL researchers
- up to **10TB** per archived dataset
- help in data curation prior to archival
- periodic integrity audits of your datasets
- periodic **reports** on your preserved datasets
- referral of your datasets on Infoscience
- can publish your datasets on <u>Zenodo</u> (size <u>limits</u>)

Data access sustainability

More than 60% of links to astronomy datasets are **broken after 10 years**

The bibliography of 1 out of every 5 is impacted by this phenomenon



A "good example of a large-scale research endeavour in which an openly accessible data repository is being used successfully" [OECD]

REPORT | VOLUME 24, ISSUE 1, P94-97, JANUARY 06, 2014

The Availability of Research Data Declines Rapidly with Article Age

Open Archive • Published: December 19, 2013 • DOI: https://doi.org/10.1016/j.cub.2013.11.014 •



Highlights

- We examined the availability of data from 516 studies between 2 and 22 years old
- The odds of a data set being reported as extant fell by 17% per year
- Broken e-mails and obsolete storage devices were the main obstacles to data sharing
- Policies mandating data archiving at publication are clearly needed

EPFL

Data degradation problem

CERN

A <u>2007 study</u> showed that a bitrot error ratio of 10^{-7} (over 2 months) Ex.: $\sim 10^9 \cdot 10^{-7} = 10^2 = 100$ bytes of bitrot every 1GB (1024MB)

US FDA

In 2017 the agengy <u>added data integrity requirements</u> for the drugs industry (FDA 21 CFR, 11 & 211)

Data integrity failure (Possible causes)

Processing

CPU heat, encryption errors, ...

Transfer

Network failures, backup errors, ...

Read / Write

Single bits errors at RAM or ROM levels

Storage

Aging, background radiation, ...

Countermeasures (Data repositories)

Redundant hardware

Uninterruptible power supply

Certain types of RAID arrays

Radiation hardened chips

Error-correcting memory

Clustered file system

File systems with block level checksums



Access sustainability without degradation



Researchers must share data "according to the FAIR Data Principles on publicly accessible, digital repositories."



Researchers must "deposit research data [...], including associated metadata, in the repository as soon as possible."

EPFL

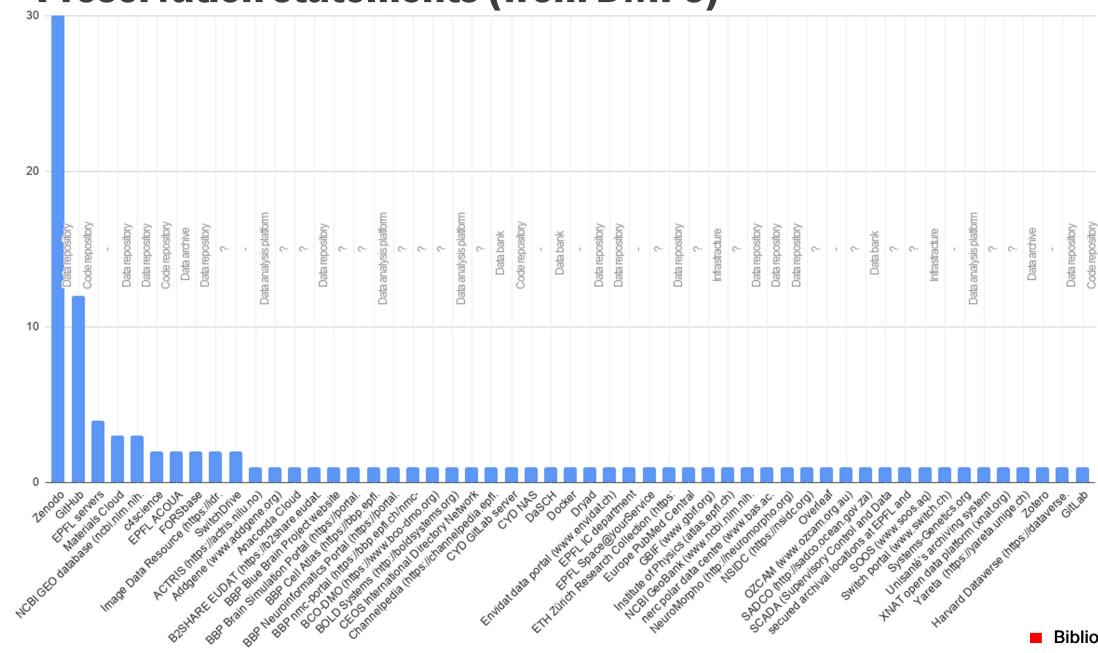
Back-up ≠ Publication ≠ Preservation

	BACKUP	PUBLICATION	LONG-TERM PRESERVATION
ACTIVE DATA	✓	*	*
DATA RECOVERY	√	✓	✓
INTEGRITY (monitoring, repair, authenticity)	?	?	✓
APPRAISAL (what & for how long)	×	√ ×	✓
PERMANENT IDENTIFIERS	*	√	√ x
DESCRIPTION (metadata)	*	✓	✓
RENDERABILITY (format migration, virtualization)	*	?	✓

Bibliothèque

EPFL

Preservation statements (from DMPs)



EPFL

How to cite data(sets)?

Same as any other citation:

- Author(s) of the dataset
- Title of the dataset / study
- Year of online publication
- Publisher responsible for distributing the dataset
- Edition / Version number associated with the dataset
- Persistent identifier(s) as URI, DOI, ORCID, ...
- Link to related objects (paper, poster, other datasets, code)



EPFL

Exploring data citations (examples)

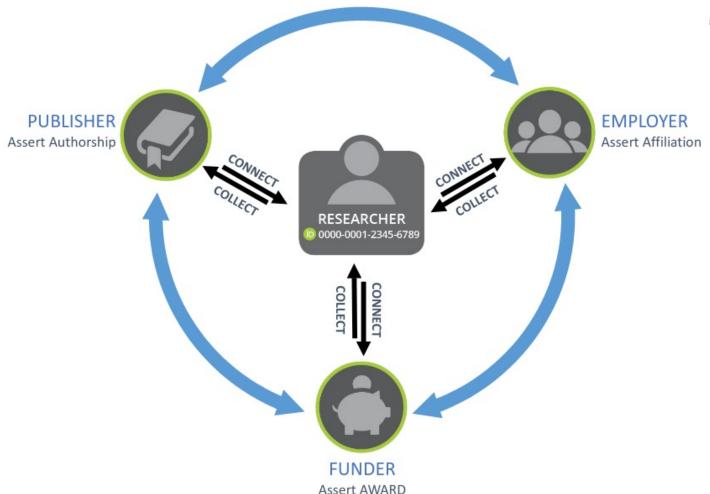
Google Dataset Search

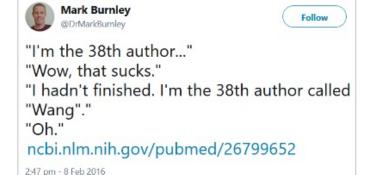
ScholeXplorer



EPFL ORCID

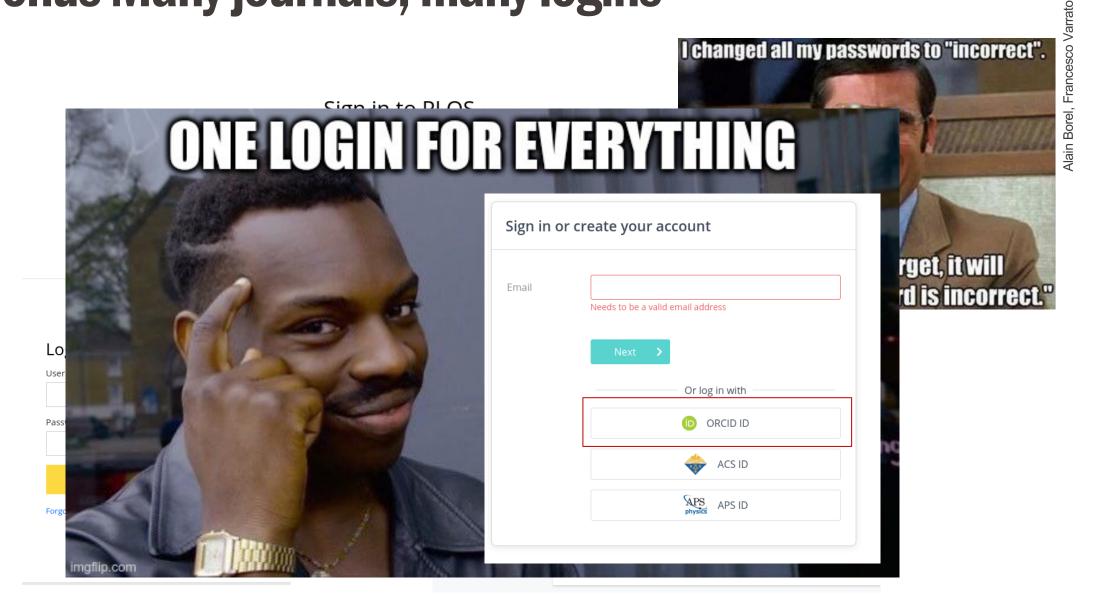
A unified personal id for scientists



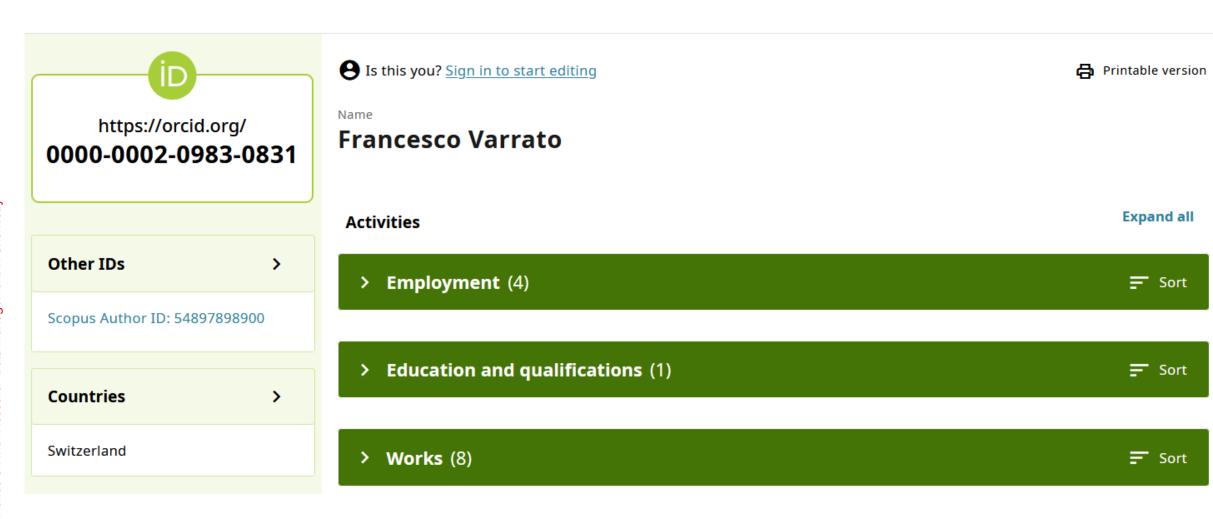


Vyas JM, Waeber C, Walker CL, Walker MJ, Walter J, Wan L, Wan X, Wang B, Wang C, Wang C, Wang C, Wang C, Wang D, Wang F, Wang F, Wang G, Wang HJ, Wang HG, Wang H, Wang HD, Wang J, Wang J, Wang M, Wang MQ, Wang PY, Wang P, Wang RC, Wang S, Wang TF, Wang X, Wang XJ, Wang XW, Wang X, Wang X, Wang Y, Wang

Bonus Many journals, many logins



ORCID



If you change name, gender, have weird umlauts, submit an SNSF project proposal,..., ORCID is the solution.

FAQ End of thesis

Where should I archive the data and code that support my thesis? ...

Recommendation

You are recommended (1, 2 and 3) to:

- archive the Research Data necessary to make your thesis reproducible
- whenever legally possible, provide the jury president and members with access to the archived datasets.

What

Research Data includes code and <u>is defined</u> as evidence that underpins the answer to the research question, and can be used to validate findings regardless of its form (e.g. print, digital, or physical). Datasets can be composed by research data necessary to validate the findings of your thesis, such as:

- raw data
- pre-processed data
- processed data
- plots
- source code
- executables documentation (e.g., README files, protocols, parameter files, log files, etc.).

www.epfl.ch/education/phd/regulations/int ernal-regulations/edoc-fag-end-of-thesis

How

Research Data should be saved in a digital archive that allows for its long-term preservation and retrieval. To ensure that the datasets will remain usable in the future, data curation prior to archiving is recommended:

- cleaning the datasets
- documenting
- enriching metadata
- converting proprietary formats into open formats
- restructuring the dataset and naming etc.

Where

EPFL offers a free archiving service for Research Data: ACOUA, the ACademic OUtput Archive. For support and information, contact the Research Data team of EPFL Library.

How long

Archived datasets can be safely stored and retrieved for many years. A basic recommendation is to make datasets preserved for at least 10 years. Depending on the research funder (SNSF, ERC, etc.) there might be specific duration requirements.

If your datasets include personal data or health data (e.g. clinical trials), both the archiving and the associated retention duration can be legal requirements. For more information in such cases, contact the EPFL Human Research Ethics Committee (HREC).