

# UNIVERSITY OF OSLO

# Sharing and archiving

research data

Live Håndlykken Kvale & Agata Bochynska  
Open Research, University of Oslo Library  
CC-BY-SA-4.0 2023



# Agenda

---

- **Why archive research data?**

- Requirements
- Reasons for sharing
- Which data should be shared?

- **Data repositories**

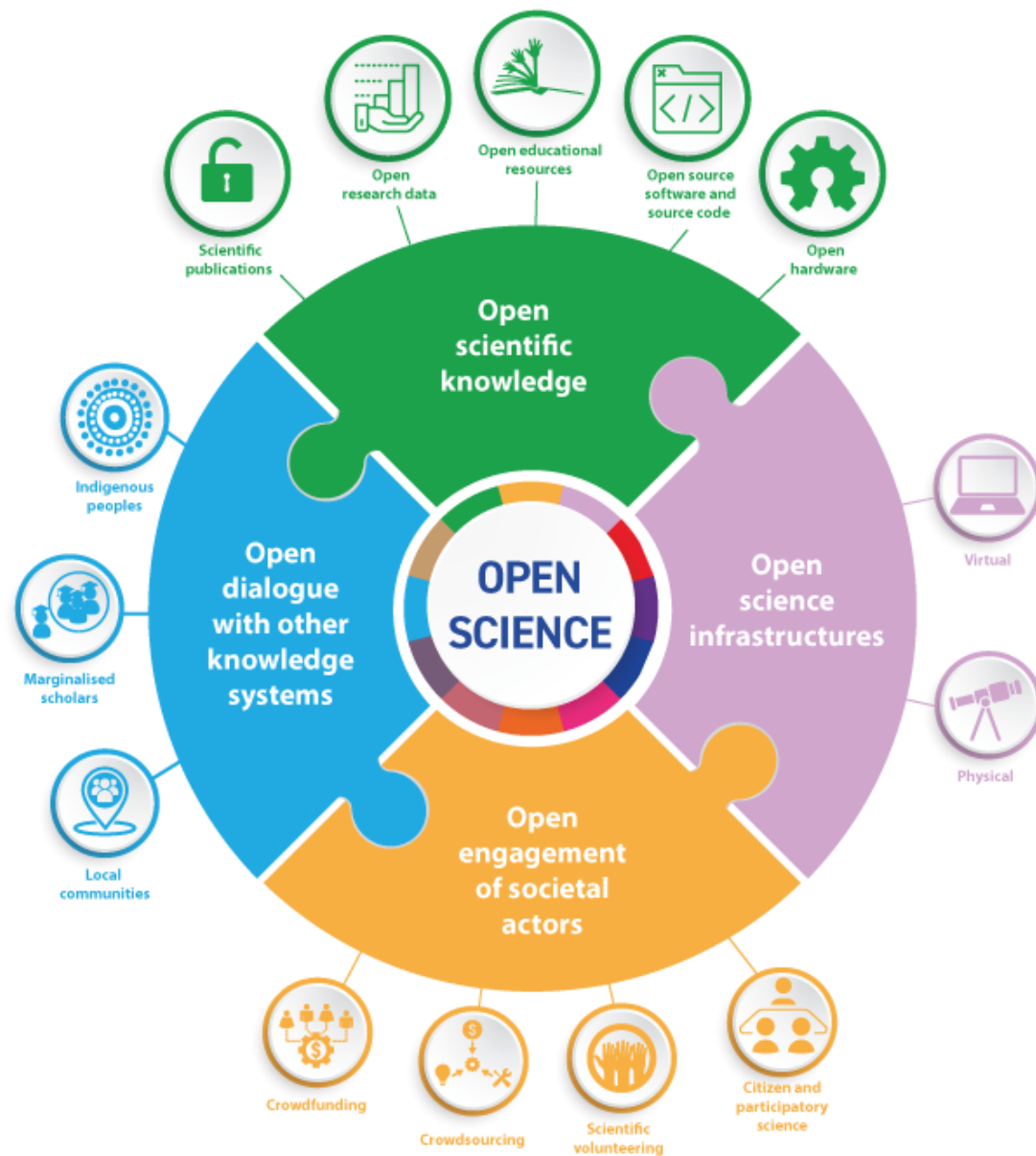
- Selecting repositories
- Levels of FAIR
- Types of repositories
- DataverseNO at UiO
- Data journals
- Archiving of code
- Finding a repository

- **What to consider**

- Preparing data
- What cannot be open
- PIDs
- Certification
- Licenses

- **Part 2:**

- **Examples**
- **Menti**
- **Key takeaways**
- **Q&A**



# Research Data in Horizon Europe

---

- You must provide open access to research data under the principle 'as open as possible, as closed as necessary'.
- In general, you should deposit data generated or collected by the project as soon as possible after data production/generation or after adequate processing and quality control have taken place (for dynamic data, a snapshot of the data is enough).
- This should happen at the latest by the end of the project, and does not entail that data are immediately open, but rather that they have been deposited so that metadata information is available and hence information about the data is findable.
- Provide information via the repository about any research output or any other tools and instruments needed to re-use or validate the data.
- It is important that you check before depositing your research data that your chosen repository is technically capable of accepting the required metadata.

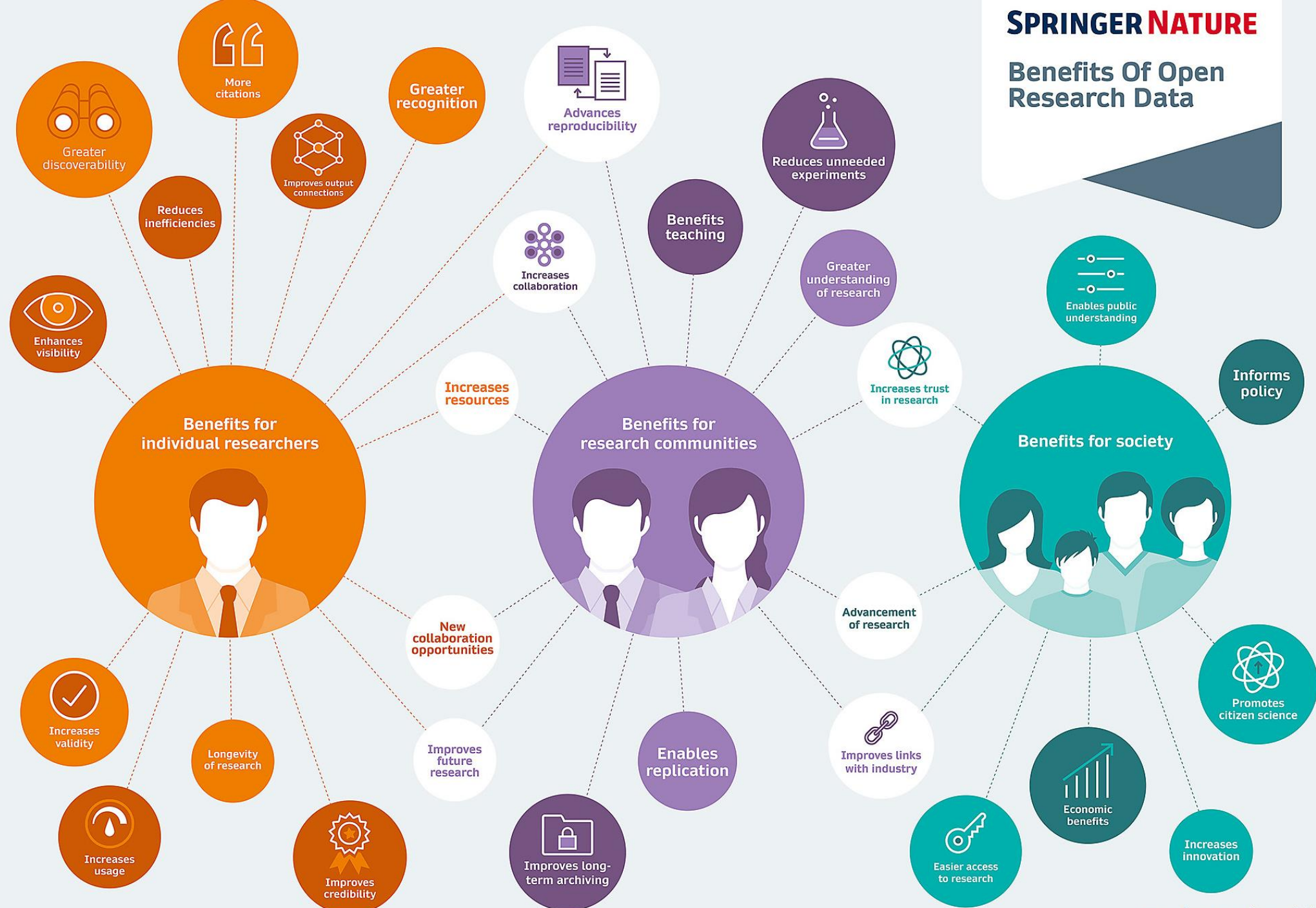
# Local Requirements

Research data at the University of Oslo shall:

- be made openly available for further usage
- be made available at an early stage
- have a data management plan
- have metadata and be documented
- must be securely archived
- have licenses for access, reuse and redistribution
- made freely available  
(but the actual distribution cost should be covered)



# Benefits Of Open Research Data



# Reasons for sharing research data

---

## **External Factors**

- Funder and publisher requirements
- Institutional requirements
- New assessment systems

## **Career Benefits**

- Increased visibility
- More data reuse
- New collaborations
- Increased citations

## **Scientific Progress**

- More robust research
- Enables verification of results
- Enables new collaborations across disciplines and borders
- Opens up for new uses of data
- Avoids duplication
- Easier to use data in teaching

# Scientific Misconduct and the Myth of Self-Correction in Science

Wolfgang Stroebe<sup>1,2</sup>, Tom Postmes<sup>2</sup>, and Russell Spears<sup>2</sup>

<sup>1</sup>Utrecht University, The Netherlands, and <sup>2</sup>University of Groningen, The Netherlands

## Abstract

The recent Stapel fraud case came as a shattering blow to the scientific community of psychologists and their image in the media and their collective self-esteem. The field responded with suggestions of how to prevent it. However, the Stapel fraud is only one among many cases. Before basing recommendations on this case would be informative to study other cases to assess how these frauds were discovered. The authors analyzed a sample of fraud cases to see whether (social) psychology is more susceptible to fraud than other disciplines. We evaluate whether the peer review process and replications work well in practice to detect fraud. There is evidence that psychology is more vulnerable to fraud than the biomedical sciences, and most frauds are detected through whistleblowers with inside information. On the basis of this analysis, the authors suggest a number of strategies to reduce the risk of scientific fraud.

## Keywords

fraud, scientific misconduct, research integrity, replication, peer review

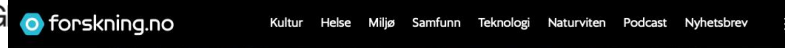
The news that the highly respected Dutch social psychologist Diederik Stapel had been accused of scientific misconduct and had admitted large-scale research fraud came as a terrible shock to the scientific community of social psychologists. Stapel was internationally renowned, and his work had received prestigious awards from the European Association of Social Psychology and the U.S. Society of Experimental Social Psychology. This scandal provided a field day for the international press, and psychology was portrayed as being highly vulnerable to scientific misconduct. The field responded with suggestions on how the risk of fraud could be reduced in the future (e.g., Crocker & Cooper, 2011; Mummendey, 2012; Roediger, 2012). However, the Stapel case, although very high profile, is only one of many fraud cases that were discovered in recent years. Instead of proposing changes on the basis of one case it

the way that most of these frauds have been discovered demonstrate that the idea of the self-correction in science is a myth.

## Notorious Cases of Research Fraud: A Review

The National Science Foundation (2001) defines scientific misconduct as fabrication, falsification, or plagiarism, performing, or reviewing research or in reporting results. Such misconduct must be committed knowingly, or in disregard of accepted practices. Fabrication of data involves totally inventing information that refers to manipulation of equipment or procedures that the research is not accurately representing.

Perspectives on Psychological Science  
7(6) 670–688  
© The Author(s) 2012  
Reprints and permission:  
sagepub.com/journalsPermissions.nav  
DOI: 10.1177/1745691612460687  
http://pps.sagepub.com



NMBU brukte åresvis på å granske mulig juks i en doktorgradsavhandling, men gjorde aldri jobben godt nok. Ikke før nå. (Foto: Lise Åserud / NTB)

# Bortkommen banankasse funnet. Etter ni år er mulig juksesak avsluttet

KOMMENTAR: En bortkommen banankasse ble funnet i et kott. Tapte data var allikevel ikke slettet. NMBU har tatt en ny runde i sin ni år lange gransking av mulig juks i en doktorgradsavhandling.



Nina Kristiansen  
JOURNALIST

Torsdag 22. september 2022 - 04:30



Det startet for ni år siden. En maraton av en sak som aldri fikk en avslutning. Verken granskerne eller NMBU gjorde en god nok jobb. Saken har versert i fagmiljøene som en verkebyll.

Nå kan saken ha fått en endelig slutt – i en ny gransking der data alle trodde var forsvunnet, har kommet til rette.

forskning.no rullet opp saken for to år siden. Da hadde rektor avsluttet saken, selv om mange i fagmiljøet mente den ikke var undersøkt grundig nok.

- Les hele saken fra 2019: [De så tegn til juks i en doktorgrad, men i løpet av tre granskninger har ingen undersøkt hele avhandlingen](#)

## Onsiktsekkende resultater

Ulyo\_fyYBBDUzLXce42n9UsAv\_9X5mPxcqRlRiInTuRQmrtkoy17ZmrojVlKXw2QxHaY69qYvheJ...aTBOIU5kz5FKKDrE092KuSxMqwb6dYLeDt4WQaJqovE57JdwmvGDCJfUHF78tXDRl4w\_Tk0

## Utgave 1 – 2007

### ARTIKLER

- > **Research misconduct: lessons to be learned?**
- > Investigation of scientific misconduct – some personal reflections
- > From Darsee to Sudbe: NLM's role in the retraction process
- > Can editors police scientific misconduct?
- > Playing by the rules – Scientific misconduct in a legal perspective

Magne Nylenna

## Research misconduct: lessons to be learned?

Michael 2007;4:7–9

«It can never happen here» has been the traditional saying in Norway when incidents of scientific dishonesty have been disclosed around the world. In a small country with a limited number of medical researchers, traditions for transparency and a strong belief in honesty, there has been a more or less naïve attitude to research misconduct.

In January 2006, on Friday 13th (1), the news was broken that a Norwegian scientist at Rikshospitalet-ospitalet, Jon Sudbø, had admitted to research misconduct in a recently published paper in *The Lancet*. It became a national sensation. The case made headline news in all major newspapers and television news, more than 330 media reports were registered over the first two weeks and the case received national attention.

At an early stage it became evident that the actual case, widely known as the Sudbø case, included fabrication, and a special Commission was appointed on 18 January to conduct an independent investigation. The Commission chaired by the Swedish epidemiologist, Professor Anders Ekbo, then presented an interim report on 30 June 2006 (2).

Most of Jon Sudbø's scientific publications are invalid due to the fabrication and manipulation of the original data material», read the main conclusion of the Commission. Based on investigations into the details of Sudbø's scientific work, 38 published papers, the Commission found several breaches of scientific practice. Jon Sudbø, a dentist and physician, had been doing research on the early stages of oral cancer. One of his main questions was whether and to what extent different types of leukoplakia could be a risk for developing oral cancer. Sudbø's results had been published in high-profile international journals (1,3,4) and formed the basis for his PhD thesis. A series of flaws were, however, found in his data and the summing up by the Commission was harsh: «The Commission is of the opinion that the identified defects that have been exposed are too numerous, too great and too obvious to be attributed to errors, incompetence or the like; and that the raw data therefore appear to have been fabricated, falsified and adapted to the desired findings»(2).

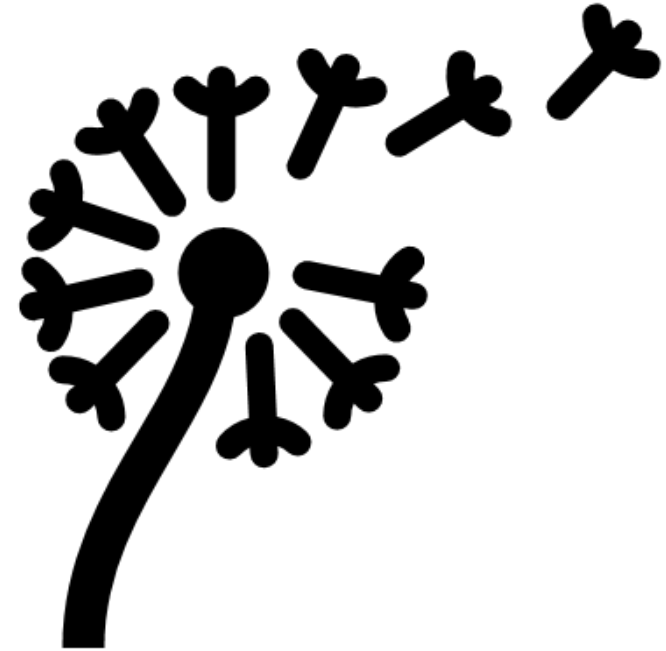
The Sudbø case has been intensively discussed within the health care sector in Norway over the last year, and undoubtedly led to an increase in the awareness of research misconduct. Many institutions have reviewed their research programmes and routines. Supervisory and regulatory systems have been strengthened.

The Sudbø case is also of interest from an international perspective. Learning from adverse events is a way to improve the quality in all parts of medicine – research as well as patient treatment. What lessons can be learned from this and other revealed cases of scientific fraud for researchers, research institutions, scientific journals and other parties? Is a more detailed bureaucratic regulation of research the inevitable consequence of research misconduct being prevented through information campaigns? And who is really responsible for the prevention of published research?

- <https://journals.sagepub.com/doi/pdf/10.1177/1745691612460687>
- <https://www.michaeljournal.no/article/2007/05/Research-misconduct-lessons-to-be-learned->
- <https://forskning.no/forskningsetikk-veterinaermedisin/bortkommen-banankasse-funnet-etter-ni-ar-er-mulig-juksesak-avsluttet/2081272>



# Strategy for archiving

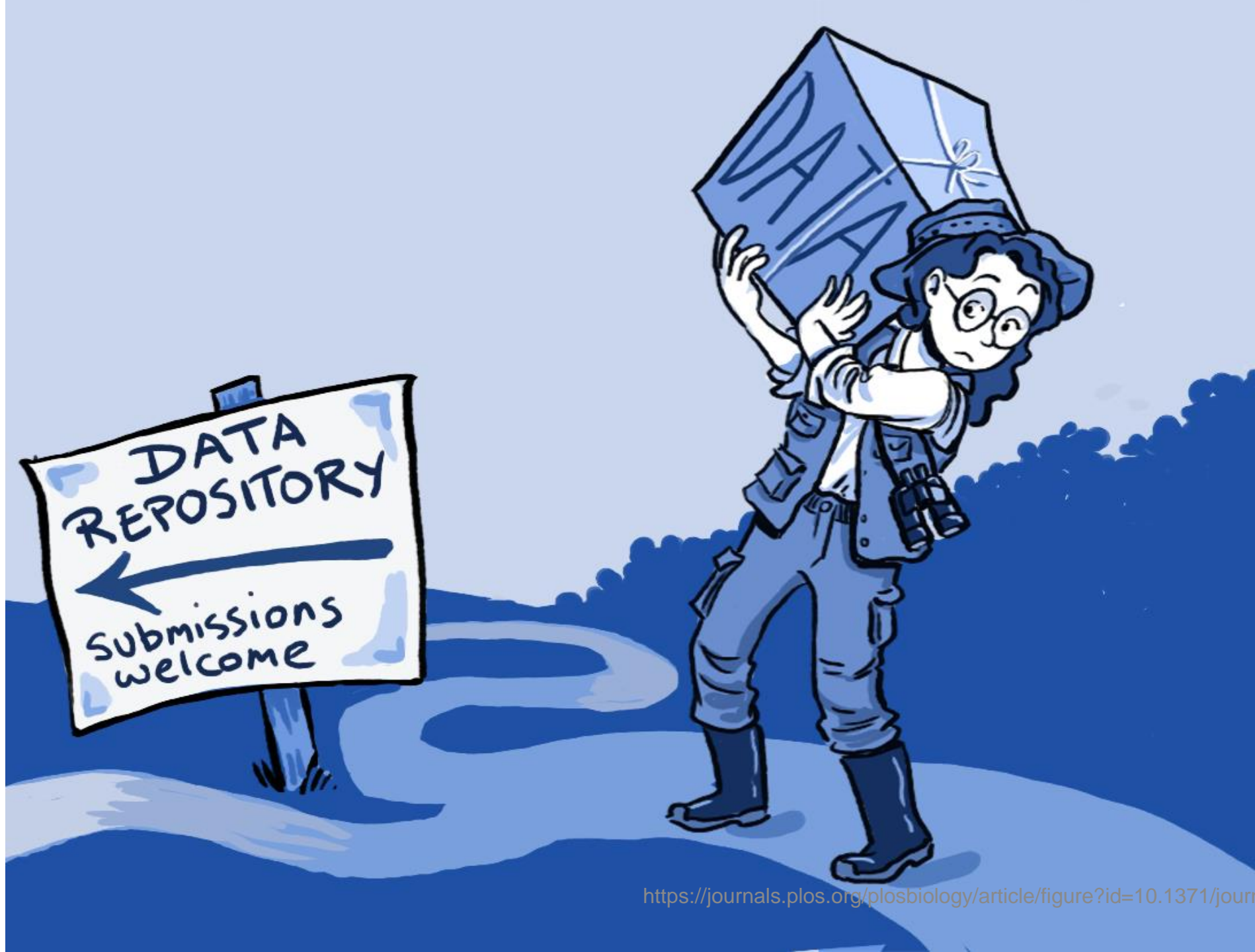


# Selecting data for archiving

---

- Does your dataset have a potential for reuse?
- (Inter-)national or historical importance
- Data quality
- Uniqueness or originality
- Size, scale, cost
- Innovativeness

How do I share data?



**YES**, a repository is the  
best place to archive data

# select a repository which...

---

- is domain specific if this exist in your field
- is a certified as trusted repository.
- supports persistent identifiers.
- offers curation.
- offers an informative landing page with metadata.
- attaches a licence.
- provide usage statistics.
- matches your particular data needs (formats, size, openness)
- provides guidance on how to cite deposited data.

# Levels of FAIR

---

- Open data, but not FAIR
- FAIR metadata
- FAIR data Restricted access
- FAIR data Open access
- FAIR linked data Restricted access
- FAIR linked open data

METADATA

DATA



# Levels of FAIR

- Open data, but not FAIR
- **FAIR metadata**
- FAIR data Restricted access
- FAIR data Open access
- FAIR linked data Restricted access
- FAIR linked open data

## METADATA



## DATA



# Levels of FAIR

- Open data, but not FAIR
- FAIR metadata
- **FAIR data Restricted access**
- FAIR data Open access
- FAIR linked data Restricted access
- FAIR linked open data

## METADATA



## DATA



# Levels of FAIR

- Open data, but not FAIR
- FAIR metadata
- FAIR data Restricted access
- **FAIR data Open access**
- FAIR linked data Restricted access
- FAIR linked open data

## METADATA



## DATA



# Levels of FAIR

- Open data, but not FAIR
- FAIR metadata
- FAIR data Restricted access
- FAIR data Open access
- **FAIR linked data Restricted access**
- FAIR linked open data

## METADATA



## DATA



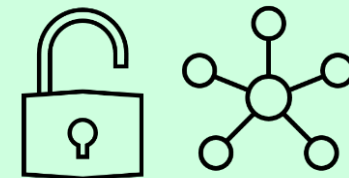
# Levels of FAIR

- Open data, but not FAIR
- FAIR metadata
- FAIR data Restricted access
- FAIR data Open access
- FAIR linked data Restricted access
- **FAIR linked open data**

## METADATA



## DATA



# Domain-specific data repositories



HEPData



# General-purpose data repositories



Open Science Framework



# National or institutional data archives



**NIRD** RESEARCH DATA ARCHIVE

# Institutional data archive





# DataverseNO at UiO

---

- **Follow the guidance for the mandatory metadata fields**
  - The curator team will go through the metadata, and normally make suggestions for improvements.
- **Use ORCID as identifier for all authors**
- **The data gets a DOI**
  - A DOI is reserved at the draft stage
  - The DOI will not be activated and work until after the dataset is published.
  - If it turns out that the dataset cannot be published the DOI will never be activated.

ORCID



# DataverseNO at UiO

---

- **DataverseNO is for Open data only**
  - Once a dataset is published, it is NOT possible to delete.
- Make sure the data does not contain **personal or confidential information.**

Notify the curator team if you:

- Need to share access with fellow researchers before publishing.
- Need to share the data for double-blind peer review.
- Have other needs regarding access.

# DataverseNO at UiO

---

- **Use keywords from controlled vocabularies to describe your data and link to these.**
  - The curator team might make smaller suggestions to improve the interoperability of the metadata.
- **Follow the curation guide**
  - Use the recommended file formats.
  - Use existing standard for describing data whenever possible.

The curator team can advise on file formats and data provenance.

# DataverseNO at UiO

---

- **Be detailed when writing the readme file.**
  - The curator team will read your readme file, and normally make suggestions for improvement.
- **Choose an appropriate license.**
  - CC-0 is currently the default license in DataverseNO.

Contact the curator team if you:

- Are uncertain whether you hold the rights to the data.
- Need advice on which license to assign.
- Need to assign different licenses to different parts of the dataset.



# The curation process

- After you submit your dataset, you will receive a curation report, a standardized document used by all DataverseNO institutions.
- Here we will inform you about required and/or recommended changes to be done before publication.
- We will return the dataset to you at the same time as you receive the curation report.
- You must resubmit the dataset after you have completed the revisions.
- We will publish the dataset when revisions have ensured the dataset is ready for publication.

DataverseNO Curation Report	
<b>Author(s):</b>	Kvale, Live
<b>Dataset:</b>	«Curation_Report_Replication_data_for-xxx_2020_xxxxx»
<b>Collection:</b>	University of Oslo
<b>Curator:</b>	Elin Frøshaug
<b>Date:</b>	09.01.2023

DataverseNO aims to make published datasets as FAIR (Findable, Accessible, Interoperable, Reusable) as possible. In order for other researchers to be able to find, understand and reuse your data, it is important that you describe them in a good way before they are published. There are particularly two places in DataverseNO where such documentation is important:

1. In the metadata schema, you should enter as much relevant information as possible so that your dataset can be found via search engines such as Google Dataset Search.
2. The ReadMe file should provide an overview of your dataset and explain how you have collected and processed your data. This documentation serves as a guide to your dataset and enables others to reuse your data.

Below you will find suggestions for changes that will make your dataset more in line with the DataverseNO guidelines (see the [Deposit Guidelines](#)) and thus increase its value and the chance that it will be found and reused.

To carry out the changes, first navigate to the draft of your dataset (DRAFT) and then click **Edit**, and select **Metadata, Files or Terms**. After making the changes, click **Save Changes**. Before uploading an edited file, you need to delete the file which needs to be replaced. After uploading the new file, click **Save Changes**. When you have done the necessary changes in your dataset, click **Submit for Review** so I can review the dataset again.

METADATA (see the section "Metadata" in the <a href="#">Deposit Guidelines</a> )	
<b>Citation Metadata</b>	
<b>Title:</b>	Replication_data_for-xxx_2020_xxxxx»
<b>Author – Name:</b>	You should write your name with your family name first: "Kvale, Live"
<b>Author – Identifier:</b>	We recommend adding your ORCID in the <i>Identifier</i> field (e.g. 0000-0001-1234-5678). Using an ORCID ensures that your research results are unambiguously linked to you as a researcher. Learn more about ORCID in <a href="#">this video</a> , and get your own ORCID at <a href="http://orcid.org/">http://orcid.org/</a> .



# UNIVERSITY OF OSLO

DataverseNO >

[✉ Contact](#) [🔄 Share](#)



[Advanced Search](#)

 **Dataverses (0)**

 **Datasets (21)**

 **Files (437)**

### Publication Year

2022 (7)

2021 (7)

2020 (5)

2018 (2)

### Subject

[Earth and Environmental Sciences \(13\)](#)

[Social Sciences \(6\)](#)

[Physics \(3\)](#)

[Agricultural Sciences \(1\)](#)

[Mathematical Sciences \(1\)](#)

[More...](#)

### Keyword Term

[Faults \(7\)](#)

[Svalbard \(7\)](#)

[Billefjorden Group \(3\)](#)

[Devonian \(3\)](#)

[Ebbadalen \(3\)](#)

[More...](#)

1 to 10 of 21 Results

↑↓ Sort ▾


[Structural field measurements in Proterozoic basement, Devonian, and Carboniferous rocks in Kongsfjorden, September 2022](#) 

Nov 2, 2022



Koehl, Jean-Baptiste P.; Stokmo, Eirik M. B., 2022, "Structural field measurements in Proterozoic basement, Devonian, and Carboniferous rocks in Kongsfjorden, September 2022", <https://doi.org/10.18710/APGAWL>, DataverseNO, V1

Structural field measurements (bedding surfaces, brittle to ductile fault surfaces and associated brittle fault lineations, foliation surfaces, fold axes, fold axial planes) from Proterozoic basement, Devonian, and Carboniferous rocks from September 2022 expedition to Kongsfjorde...

[Replication data for the MULTICLIM project "Pesticide effects on the abundance of springtails and mites in field mesocosms at an agricultural site"](#) 

Nov 2, 2022



Konestabo, Heidi Sjursen, 2022, "Replication data for the MULTICLIM project "Pesticide effects on the abundance of springtails and mites in field mesocosms at an agricultural site"", <https://doi.org/10.18710/QWIDIT>, DataverseNO, V1

This dataset was collected as part of the MULTICLIM project at the University of Oslo: <https://www.mn.uio.no/ibv/english/research/sections/aqua/research-projects/144612/> The data has been published in Ecotoxicology, doi:10.1007/s10646-022-02599-3. The main aim of the study for wh...

[Field photographs Kongsfjorden September 2022](#) 

Oct 26, 2022



Koehl, Jean-Baptiste P.; Stokmo, Eirik M. B., 2022, "Field photographs Kongsfjorden September 2022", <https://doi.org/10.18710/KEB2MM>, DataverseNO, V1

Geological fieldwork photographs of September 2022 excursion to Kongsfjorden. Day 0: flight journey onboard Lufttransport plane from Longyearbyen to Ny-Ålesund. Day 1: trip by zodiac boat from Ny-Ålesund to southern Blomstrandhalvøya; landing site: London. Then walking southeastw...

scientific **data**

# Data journals



**Research Data Journal  
for the Humanities  
and Social Sciences**

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


<https://www.sciencedirect.com/journal/data-in-brief>

<https://brill.com/view/journals/rdj/rdj-overview.xml>

[nature](#) > scientific data

**Dynamic World, Near real-time global 10m land use land cover mapping**

Christopher F. Brown, Steven P. Brumby ... Alexander M. Tait  
Data Descriptor | 09 June 2022

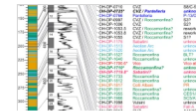


**Featured**

Data Descriptor  
Open Access  
02 Sept 2021

**Lake Ohrid's tephrochronological dataset reveals 1.36 Ma of Mediterranean explosive volcanic activity**

Niklas Leicher, Biagio Giaccio ... Bernd Wagner



Data Descriptor  
Open Access  
12 Aug 2021

**OPERA tau neutrino charged current interactions**

N. Agafonova, A. Alexandrov ... C. S. Yoon



**Announcements**

**Open data in the COVID-19 pandemic**

A collection presenting a series of rapidly evolving resources that aggregate and bring cohesion to the massive volume of data being generated in the COVID-19 crisis



Search Scientific Data

All Subjects Q

[Find out more about Scientific Data](#)



[Find the right repository for your data](#)

26.04.2023





# Archiving of Code



# Finding archives/repositories

---

Search

Browse ▾

Suggest

Resources ▾

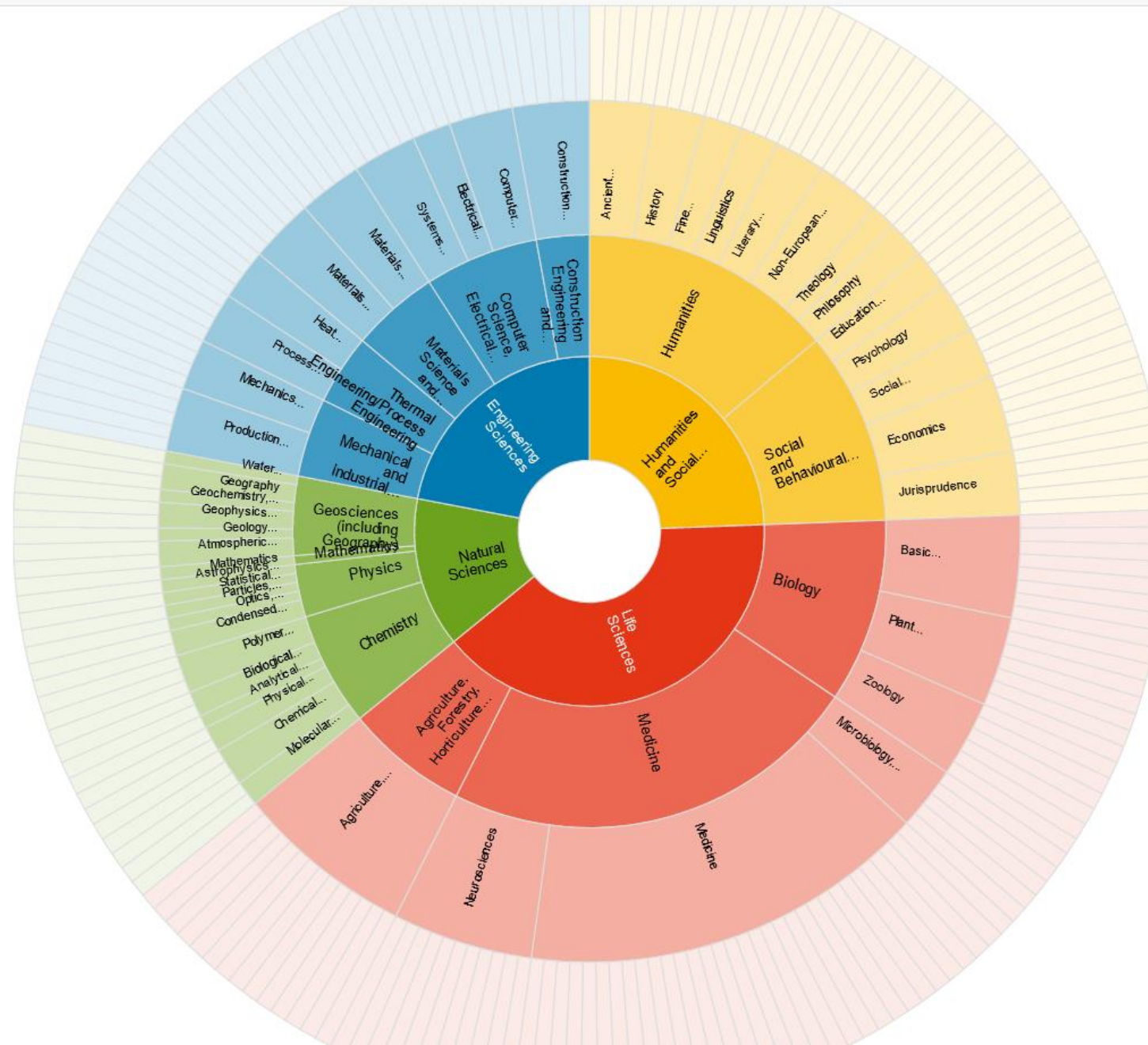
Contact



**re3data.org**  
REGISTRY OF RESEARCH DATA REPOSITORIES

Search...

🔍 Search



# Finding archives

---

STANDARDS

**DATABASES**

POLICIES

COLLECTIONS

ORGANISATIONS

ADD CONTENT

## Databases

A registry of knowledgebases and repositories of data and other digital assets.



**Dr. Irene Salinas Remiro** @DrSalinasLab · 27 Jan



Is there a public **repository** to submit 3D microscopy images for **data sharing** purposes? Zebrafish folks, what do you usually do? Thanks!



8



6



20



13.6K



**Victoria E. Abraira** @VAbraira · 23 Jan



Hey Neuro colleagues: what is your favorite **repository** for **data sharing**?



2



6



2,061



# Selecting an archive

---

- Should the data be openly available?
- Should the metadata be openly available?
- Presence of personal or confidential data can affect choice of archive / preservation solution
- What can the relevant archives offer for long term perspectives?
- Does the archive offer curation – control of metadata and updating of formats?

# Preparing for archiving

---

- Do you have permission to share the data?
- Consistent, meaningful, and compatible file naming
- Choose accessible, patent-free, and open file formats
- Write a readme file where you describe the data.

## Prepare your data

Before depositing your data in DataverseNO (including the different collections, e.g. UiT Open Research Data, TROLLing, etc.) you have to make sure your dataset(s) comply with our guidelines below. DataverseNO accepts only research data in digital formats. In brief, good practice for preparing research data for archiving may be summarized as follows:

- Use consistent and comprehensible file names (see section 1 below).
- Save your data in a preferred file format(s) (see section 2 below).
- Describe your data in a ReadMe file (see section 3 below).

For more detailed guidelines, see below:

### ▼ 1 File naming and organization

### ▼ 2 Preferred file formats

### ^ 3 How to describe your data



In order for other researchers to be able to understand and reuse your data, it is essential that you describe them in a comprehensible and consistent manner before they are published. In DataverseNO, this kind of documentation must be provided in two ways, in the **metadata fields**, and in a separate **ReadMe file** which **must** be uploaded together with your data files:

### ▼ Metadata

### ^ ReadMe file



A **ReadMe file** is a more detailed user guide to your dataset so that other researchers are able to interpret, understand, and reuse your data, including information about how the dataset was created, how complete it is, and what kind of restrictions it has. The ReadMe file must minimally contain the following:

- Title of the dataset, DOI, contact information
- Methods
- Data and file overview
- Data-specific information
- Terms of Reuse





The Turing Way

Search this book...

Welcome

**Guide for Reproducible Research**

Overview

Open Research

Version Control

Licensing

**Research Data Management**

Research Data

Data Management Plan

The FAIR Principles and Practices

Personal data management

Data Storage and Organisation

**Data Organisation in Spreadsheets**

Documentation and Metadata

Data Curation

Sharing and Archiving Data

Data Article

Research Data Management Toolkits

Personal Impact Stories

Checklist

Resources

Reproducible Environments

BinderHub

Code quality

Code Testing

Code Reviewing Process

Reusable Code

Continuous Integration (CI)

# Data Organisation in Spreadsheets

Spreadsheets, such as Microsoft Excel files, google sheets, and their Open Source alternative (for instance) LibreOffice, are commonly used to collect, store, manipulate, analyse, and share research data. Spreadsheets are convenient and easy-to-use tools for organising information into an easy to write and easy to read forms for humans. However, one should use them with caution, as the use of an inappropriate spreadsheet is a major cause of mistakes in the data analysis workflow. There is a collection of [horror-stories](#) that tells how the use of spreadsheets can ruin analysis-based studies due to unexpected behaviour of the spreadsheet or error-prone editing processes. Some of these mishaps are not unique to spreadsheets, but many, such as [this](#) and [this](#), are.

Fortunately, most problems can be avoided with the following recommendations:

- Use spreadsheet in a text-only format (.csv or .tsv),
- Create tidy spreadsheets,
- Make spreadsheets consistent (with each other) and implement rules for data entries, and
- Avoid manipulating and analysing data in spreadsheet software (this includes copy-paste).

Spreadsheets are a powerful tool only if the dataset is collected and organised in specific formats that are usable for both the computers and researchers.

## 1. Avoid Non-Data Content

Spreadsheets are used for organising data in a tabular form. The subject, the object and the relationship between them are transformed into rows, cells and columns, respectively. For example, the subject: `experiment`, relationship: `was performed on the date`, and the object: `2020-06-06` gives one row for each experiment, one column for `date of experiment`, and the value `2020-06-06` in the cell. Unfortunately, spreadsheet programs allow you to add other kinds of contents to this, like color to specific cells. While it may help the researchers at some point, one needs to remember that this kind of **cell modification should not be considered as data**, primarily because they cannot be exported to other software.

As a simple rule, what can be exported in a text-only format, comma-separated values (CSV), or tab-separated values (TSV), can be considered as the data. Other functions should be avoided when using these programs for research data. This includes:

- changing font, color or borders,
- using functions,
- merging cells (this one is particularly problematic),
- using specific cell formats (especially dates, see below).

As a test for your spreadsheet compatibility with reproducible research, export your data from the spreadsheet to the CSV format and reopen it. If you can still get all the information that you stored in your sheet, then your data is fine.

### Tip

If you want to use color to help with a rapid highlight in your document, create a new column to indicate which cells are highlighted (it becomes a part of your data). In addition to the visual feedback, you can now also use this information to filter or sort your data and get the highlighted cells quickly.



Contents

1. Avoid Non-Data Content
  2. Tidy Format For Spreadsheets
  3. Consistent Values
  4. Data Manipulation and Analysis
- Other Tips
- Summary

<https://the-turing-way.netlify.app/reproducible-research/rdm/rdm-spreadsheets.html>



### Text documents

### Plain text

### Markup language

### Programming languages

### Spreadsheets

### Databases

### Statistical data

#### • Preferred format(s)

- PDF/A (.pdf)
- ODT (.odt)
  
- Unicode text (.txt)
  
- XML (.xml)
- HTML (.html)
- Related files: .css, .xslt, .js, .es
  
- MATLAB
- NetCDF
- TextFabric
  
- ODS (.ods)
- CSV (.csv)
  
- SQL (.sql)
- SIARD (.siard)
- CSV (.csv)
  
- SPSS (.dat/.sps)
- STATA (.dat/.DO)
- R

#### • Non-preferred format(s)

- Microsoft Word (.doc)
- Office Open XML (.docx)
- Rich Text File (.rtf)
- PDF other than PDF/A (.pdf)
  
- Non-Unicode text (.txt)
  
- SGML (.sgml)
- Markdown (.md)
  
  
  
  
  
- Microsoft Excel (.xls)
- Office Open XML Workbook (.xlsx)
- PDF/A (.pdf)
  
- Microsoft Access (.mdb, .accdb)
- dBase (.dbf)
- HDF5 (.hdf5, .he5, .h5)
  
- SPSS Portable (.por)
- SPSS (.sav)
- STATA (.dta)
- SAS (.7dat; .sd2; .tpt)

As open as possible,  
as closed as necessary

EDITORIAL

## Open science and sharing personal data widely – legally impossible for Europeans?

Giske Ursin<sup>a,b,c</sup> and Heidi Beate Bentzen<sup>a,d</sup> 

<sup>a</sup>Cancer Registry of Norway, Oslo, Norway; <sup>b</sup>Institute of Basic Medical Sciences, University of Oslo, Oslo, Norway; <sup>c</sup>Department of Preventive Medicine, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA; <sup>d</sup>Norwegian Research Center for Computers and Law, Faculty of Law, University of Oslo, Oslo, Norway

**ARTICLE HISTORY** Received 15 October 2021; accepted 15 October 2021

A requirement for having a research paper published in many medical journals is that the authors include a data sharing statement. Although the requirement from the International Committee of Medical Journal Editors is not very strict, simply requiring a statement [1], interpretation varies. Some journals essentially require that data must be *readily available* for other researchers for the paper to be accepted.

While most of us eagerly welcome open science and reuse of data to ensure reproducible science, the EU General Data Protection Regulation (GDPR) provides strong protection of privacy and rather restricts and counteracts open sharing of personal data [2]. Some editors will accept that data are not readily sharable with others than peer reviewers for legal reasons. However, editors of non-European journals will often object to a GDPR-compatible data sharing statement and, consequently and often at the last minute, reject the research paper.

Why is this an issue? How difficult is it for European researchers to share data with researchers in other parts of the world?

supplementary measures in place to protect the data. The European research institution will in collaboration with the data importer need to conduct a thorough assessment of the importer's country's laws to ensure that an EU level of data protection is obtained. Such assessments require sound knowledge of the EU Charter of Fundamental Rights, the GDPR, the Court of Justice of the European Union *Schrems II* judgment, and subsequent guidance from the European Data Protection Board, which comprises all Data Protection Authorities in the European Economic Area (EEA) [4]. Finally, the data exporter must be willing to take the risk that the national Data Protection Authority agrees that all requirements have indeed been met, as fines can be high if the institution makes a mistake. Whenever such transfer is possible to achieve, you are lucky! The only cost is that the legal and administrative work on your end has quadrupled compared to the pre-GDPR era.

### The federal challenge



Info Tags Related

Item Type Journal Article  
Title Open science and sharing personal data widely – legally impossible for Europeans?  
Author Ursin, Giske  
Author Bentzen, Heidi Beate  
Abstract  
Publication Acta Oncologica  
Volume 60  
Issue 12  
Pages 1555–1556  
Date 2021–12–02 y m d  
Series  
Series Title  
Series Text  
Journal Abbr  
Language  
DOI 10.1080/0284186X.2021.1995894  
ISSN 0284–186X  
Short Title  
URL <https://doi.org/10.1080/0284186X.2021....>  
Accessed 02/08/2022, 14:30:18  
Archive  
Loc. in Archive  
Library Catalogue Taylor and Francis+NEJM  
Call Number  
Rights  
Extra Publisher: Taylor & Francis  
\_eprint: <https://doi.org/10.1080/0284186X.2021.1995894>  
PMID: 34797207  
Date Added 02/08/2022, 14:30:18  
Modified 02/08/2022, 14:30:18

# Persistent identifiers (PIDs)



<http://urn.nb.no/URN:NBN:no-5678>

ORCID

ROR

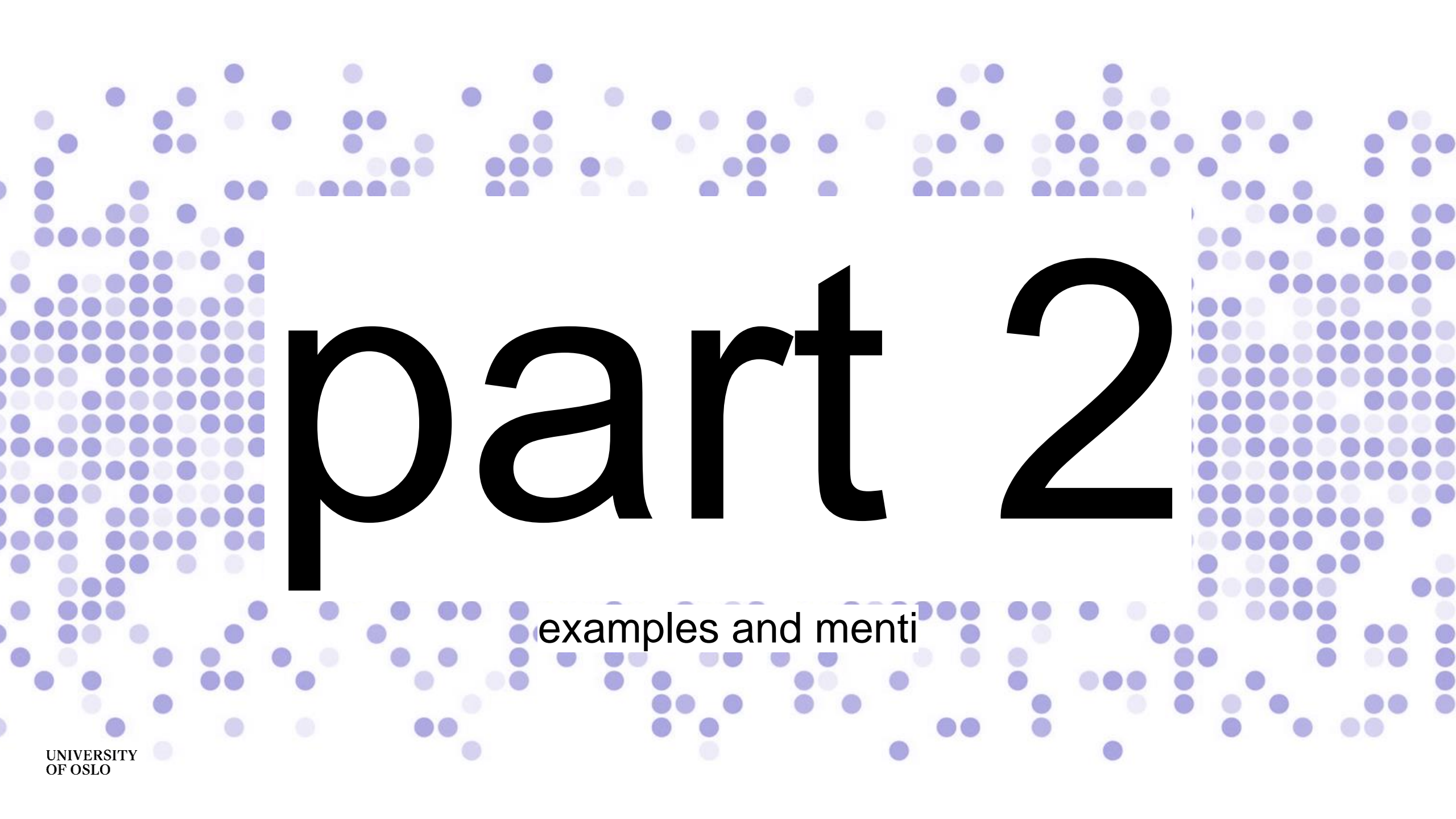
DOI – (digital object identifier) Commonly used for datasets and publications  
URN – Persistent identifiers used in DUO for theses and self-archived publications.  
ORCID – identifies the researcher  
ROR - Identifies the University

# Certification



international, community based, non-governmental, and non-profit organization promoting sustainable and trustworthy data infrastructures

**License** makes data reusable



# part 2

examples and menti



Contact us at [research-data@uio.no](mailto:research-data@uio.no)



# Key takeaways:

---

1. If you put your data in an archive, you get a **DOI**
2. Not all archives offer **curation** – if the data have a long-term value, choose an archive with curation
3. Add a **license** to your dataset and code so that others know what they can and cannot do with the data
4. If you share **code**, archive a frozen, citeable version (e.g. via GitHub - Zenodo connection)
5. Without metadata, documentation and structure, your data are not reusable

# Links:

FAIR: <https://www.force11.org/fairprinciples>

Selecting an archive: <https://www.ub.uio.no/english/writing-publishing/data-archiving/selecting-archive.html>

DataverseNO: <https://dataverse.no/dataverse/uio>

NSD/Sikt: <https://sikt.no/en/archiving-research-data>

NIRD: <https://www.sigma2.no/research-data-archive>

Zenodo: <https://zenodo.org>

Figshare: <https://figshare.com>

OSF <https://osf.io>

Prepare data: <https://site.uit.no/dataverseno/deposit/prepare/>

Github citable code: <https://guides.github.com/activities/citable-code/>

Re3data: <https://www.re3data.org/>

Creative Commons Licenses: <https://creativecommons.org>

CoreTrustSeal: <https://www.coretrustseal.org/>

UiO Research data management <https://www.uio.no/english/for-employees/support/research/research-data-management/>