# Sharing and archiving research data

Live Håndlykken Kvale and Agata Bochynska

May 20, 2021

9:00 AM–10:30 AM, Zoom

# What types of data do you work with?

# Report

# The Availability of Research Data Declines Rapidly with Article Age

Timothy H. Vines,[1,2,*] Arianne Y.K. Albert,[3] Rose L. Andrew,[1]
Florence Débarre,[1,4] Dan G. Bock,[1] Michelle T. Franklin,[1,5]
Kimberly J. Gilbert,[1] Jean-Sébastien Moore,[1,6]
Sébastien Renaut,[1] and Diana J. Rennison[1]
[1]Biodiversity Research Centre, University of British Columbia,
6270 University Boulevard, Vancouver, BC V6T 1Z4, Canada
[2]Molecular Ecology Editorial Office, 6270 University
Boulevard, Vancouver, BC V6T 1Z4, Canada
[3]Women's Health Research Institute, 4500 Oak Street,
Vancouver, BC V6H 3N1, Canada
[4]Centre for Ecology & Conservation Biosciences, University of
Exeter, Cornwall Campus, Tremough, Penryn TR10 9EZ, UK
[5]Institute for Sustainable Horticulture, Kwantlen Polytechnic
University, 12666 72nd Avenue, Surrey, BC V3W 2M8, Canada
[6]Department of Biology, Université Laval, 1030 Avenue de la
Médecine, Laval, QC G1V 0A6, Canada

## Summary

Policies ensuring that research data are available on public

sets (23%) were confirmed as extant. Table 1 provides a breakdown of the data by year.

We used logistic regression to formally investigate the relationships between the age of the paper and (1) the probability that at least one e-mail appeared to work (i.e., did not generate an error message), (2) the conditional probability of a response given that at least one e-mail appeared to work, (3) the conditional probability of getting a response that indicated the status of the data (data lost, data exist but unwilling to share, or data shared) given that a response was received, and, finally, (4) the conditional probability that the data were extant (either "shared" or "exists but unwilling to share") given that an informative response was received.

There was a negative relationship between the age of the paper and the probability of finding at least one apparently working e-mail either in the paper or by searching online (odds ratio [OR] = 0.93 [0.90–0.96, 95% confidence interval (CI)], p < 0.00001). The odds ratio suggests that for every year since publication, the odds of finding at least one apparently working e-mail decreased by 7% (Figure 1A). Since we searched for e-mails in both the paper and online, four factors

# Why would you like to share your data?



support the results

mandatory

follow up study
future research
optimise knowledge growth

credibility

open research

training material

reproducibility

it's messy
open science
transparancy

transparency

sensitive data

funding org requirement

fair

requirement

privacy

inform other research

alternative analyses

collaboration

expectation of funders

progress in research

# Why would you not share your data?

not good enough

takes time to prepare

technical problems

personal data protection

supervisor does not allow

loosing context

publications

exclusivity

sensitive data

problems with anonymity

unpublished work

gdpr

working on it actively

identifiability

data of economical value

it's messy
workload

patents

privacy

# Reasons for sharing your data

**Career Benefits**

- Increased visibility
- More reuse
- Increased citations

**Norms**

- «This is how we do it here»

**External Factors**

- Funder requirements
- Publisher requirements

**Scientific Progress**

- More robust research
- Enables new collaborations
- Opens up for new uses of data
- Avoids duplication
- Builds links to younger researchers
- Easier to use data in teaching

Source: https://www.cessda.eu/Training/Training-Resources/Library/Data-Management-Expert-Guide/6.-Archive-Publish/Towards-archiving-publication

# Open science

**Open science** means transparency and knowledge-sharing in research processes to make knowledge accessible across academic groups, sectors and national boundaries. The concept of open science encompasses the entire research process [...].

# What are the advantages of open science?

| reproducibility | Visibility | Replicability and reproducibility |
| Reproducibility | Interdisciplinary research | It's fair, transparent, fuels new research analyses. |
| Transparency | Saves cost | Accessibility to more knowledge |

19

# What are the advantages of open science?

better for the world that we share and help each other

Collaboration

Transparency reusability

transparency

scientific progress

Share, avoid duplication, more citations, more collaborations...

Not re-doing the same stuff over and over again

Science is based on sharing knowledge

Truth-seeking

19

# What are the advantages of open science?

closed science is no science

19

Source: FORCE11.org

Data

Metadata

Filename:
Tadzik.jpg
Author:
Piotr Kononow
Date:
August 15, 2016

https://dataedo.com/kb/data-glossary/what-is-metadata

# Descriptive metadata in a research data context

→ Tittle (on project and files)

→ Author (creator, copyright holder)

→ Publication year (and year(s) of data collection)

→ Persistent identifier (DOI)

→ Location (preferably coordinates)

→ Which publication(s) the datasets are used in

→ And so on...

# Selecting data for archiving

→ Does your dataset have a potential for reuse?

→ (Inter-)national or historical importance

→ Quality

→ Uniqueness or originality

→ Size, scale, cost

→ Innovativeness

# Preparing for archiving

→ Determine scientific relevance and need for archiving long-term

→ Consistent, meaningful, and compatible file naming

→ Choose accessible, patent-free, and open file formats

→ Make sure you have the necessary documentation (and metadata)

→ Reduce complexity by grouping large groups of similar files in zip bundles to make upload and download easier

→ Presence of personal or confidential data can affect choice of archive

→ Consider size limitations when choosing an archive (e.g. some archives have a limit of 10 to 50 GB per dataset)

| Type | Preferred format(s) | Non-preferred format(s) |
|---|---|---|
| Text documents | PDF/A (.pdf)<br>ODT (.odt) | Microsoft Word (.doc)<br>Office Open XML (.docx)<br>Rich Text File (.rtf)<br>PDF other than PDF/A (.pdf) |
| Plain text | Unicode text (.txt) | Non-Unicode text (.txt) |
| Markup language | XML (.xml)<br>HTML (.html)<br>Related files: .css, .xslt, .js, .es | SGML (.sgml)<br>Markdown (.md) |
| Programming languages | MATLAB<br>NetCDF<br>TextFabric | |
| Spreadsheets | ODS (.ods)<br>CSV (.csv) | Microsoft Excel (.xls)<br>Office Open XML Workbook (.xlsx)<br>PDF/A (.pdf) |
| Databases | SQL (.sql)<br>SIARD (.siard)<br>CSV (.csv) | Microsoft Access (.mdb, .accdb)<br>dBase (.dbf)<br>HDF5 (.hdf5, .he5, .h5) |
| Statistical data | SPSS (.dat/.sps)<br>STATA (.dat/.DO)<br>R | SPSS Portable (.por)<br>SPSS (.sav)<br>STATA (.dta)<br>SAS (.7dat; .sd2; .tpt) |

https://dans.knaw.nl/en/about/services/easy/information-about-depositing-data/before-depositing/file-formats

# Licensing your data

A license agreement is a legal arrangement between the creator/depositor of the data set and the data repository, signifying what a user is allowed to do with the data.

Creative Commons licenses are often used

**CC BY NC**
- Credit must be given to the creator
– Only noncommercial uses of the work are permitted

**CC BY SA**
– Credit must be given to the creator
– Adaptations must be shared under the same terms

**CC BY**
- Credit must be given to the creator

**PUBLIC DOMAIN**
- is a public dedication tool, which allows creators to give up their copyright and put their works into the worldwide public domain.

Persistent identifiers and data citation explained

Reusable code

# Selecting an archive

→ Should the data be openly available?

→ Should the metadata be openly available?

→ What can the relevant archives offer for long term perspectives?

→ Does the archive offer curation – control of metadata and updating of formats?

# Types of data archives

→ Domain-specific

→ General purpose

→ Institutional

→ Supplementary material to an article

→ Data paper

# 7

minutes break

# Case 1

A shares movement

# A investigates humans interacting with music

→ Professional musicians participating with self-compound music

→ Data collection from multiple sources

→ Audience informed that the concerts are also research projects

→ A wishes to do science as open as possible

# Data from A's project consists of:

→ video recordings form the concert

→ sensor-data from someone in the audience

→ sensor-data from three musicians

→ audio recordings of the music

→ survey responses from the audience

→ survey questions

→ analysis from "live" data jockeying during the event

→ notes

→ photos

→ music scores

→ code

# What legal challenges might A encounter?

identifiable personal inf

consent

music copyright

privacy

test

anonymity difficult

videos

music scores

copyright

anonymisation

photos

privacy of the videos

# Which data types are likely to be sensitive/special categories?

**17**
video recordings form the concert

**15**
sensor-data from someone in the audience

**10**
sensor-data from the three musicians

**6**
audio recordings of the music

**9**
survey responses from the audience

**0**
survey questions

**5**
analysis from "live" data jockeying during the event

**2**
notes

**19**
photos

**5**
music scores

**1**
code

# The approach chosen by A

https://www.uio.no/ritmo/english/news-and-events/events/musiclab/2019/utopia/index.html

# Ø conducted interviews with adolescents

→ Ø is a specialist in psychiatry and a researcher at U

→ The adolescents are in a vulnerable situation

→ The interviews are unique and of huge value for both research and as historical documents

→ Because it is difficult to talk about the trauma, other researchers want to reuse the data and not interview the adolescents again

→ Ø claims he has a unique right to the material, and does not want to share it with anyone

# What ethical challenges occur in this situation?

It is not Ø's data, it's the participants'

Sensitive data!

Ø might have right to use the data, but he does not own the data

Ø does not own the data, his institution does

do you mean sharing data while he is still working on the data?

Confidentiality agreement

He keeps data private thought its collection has been funded by a public institution.

risk of further trauma in participants

have the people beining inteviewed agreed to the data being used by other researches?

11

# What ethical challenges occur in this situation?

How to provide sufficient context for transparency to be true?

Making the data non-identifiable

11

# 22. juli-forskere kan bli overkjørt av kunnskapsdeling

Terrorofrene fra 2011 skal ikke belastes unødig av forskere. Men å måtte dele dataene som samles inn, mener noen forskere er overkjørende.

Siw Ellen Jakobsen, frilansjournalist | De nasjonale forskningsetiske komiteene

16.6 2013 05:00

ANNONSE

A thorough ethical evaluation conclude that the data should be shared with other researchers

# Regarding reuse and sharing of data

1) The coordination group recommends that in future research common data platforms should be created and used across subjects and institutions.

2) It also recommends that the metadata for July 22nd research from different studies is gathered at NSD, and that a portal for 22nd of July research is created trough CRIStin. This can be done within the research environment at NSD, with some support to cover the cost of metadata creation.
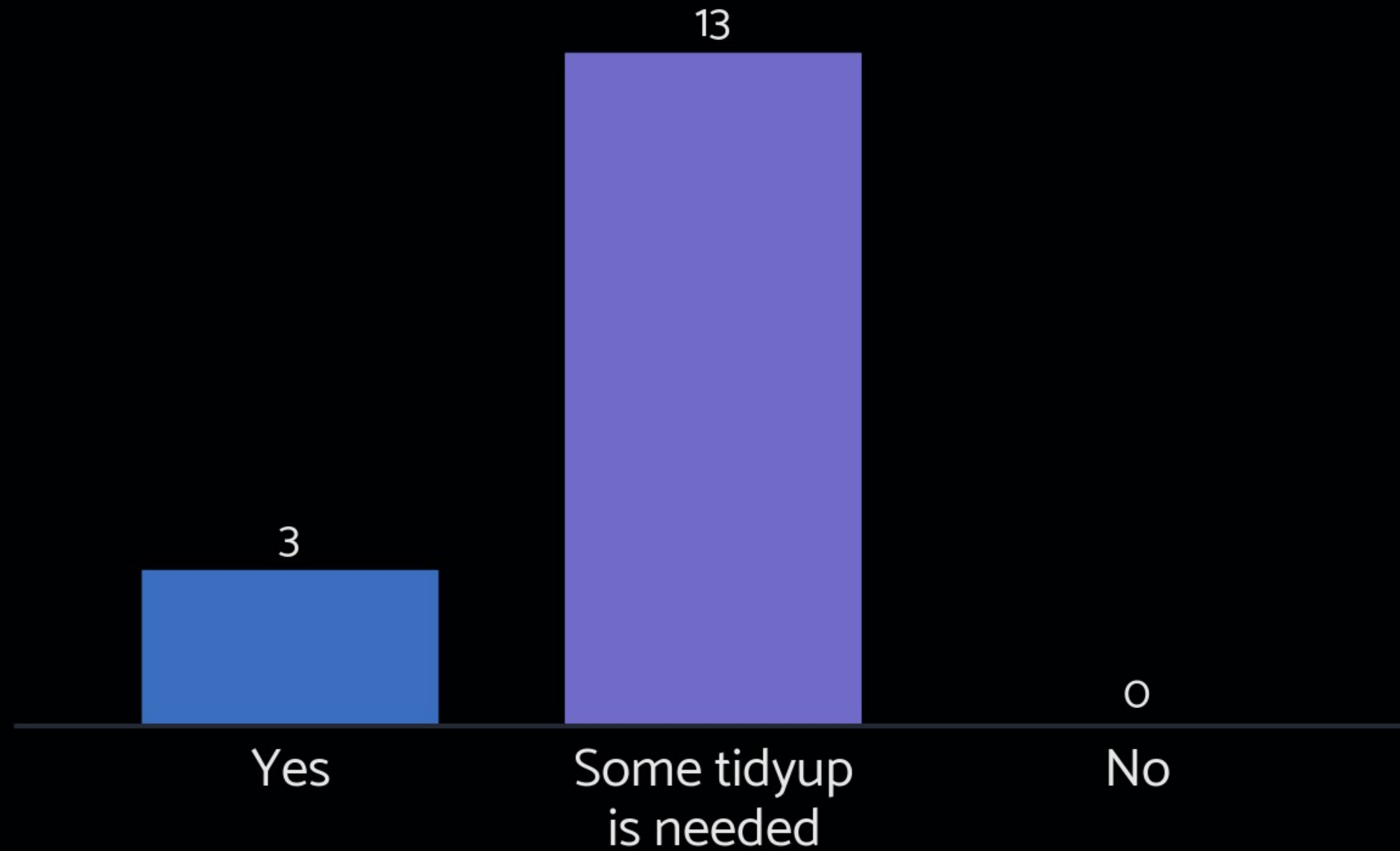
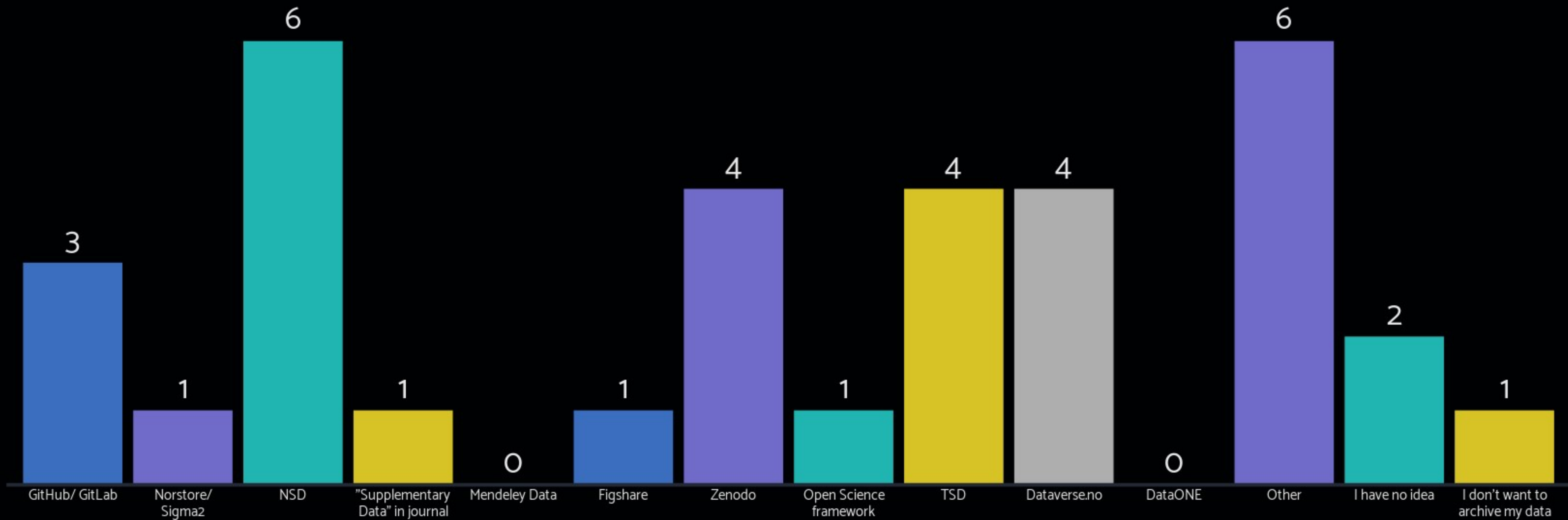# Research data on July 22nd events are today archived at NSD

https://nsd.no/nsddata/22juli/datasett.html?a=/nsddata/22juli/datasett/datasett0009.html

Where do you plan to archive your research data?

# Ask us

0 questions
0 upvotes

# Thank you for your attention

Contact research-data@uio.no for questions

Agata Bochynska, Ivana Malovic, Solveig
Sørbø, Live H. Kvale

# Sources

→ Bjerknes Centre. Bjerknes Climate Data Centre. https://www.bcdc.no/

→ Center for Open Science. Open Science Framework. https://osf.io/

→ Creative Commons. https://creativecommons.org/

→ DataverseNO. https://dataverse.no/

→ DataverseNO. Prepare your data. https://site.uit.no/dataverseno/deposit/prepare/

→ Mons, Barend et al. "Cloudy, increasingly FAIR; revisiting the FAIR Data guiding principles for the European Open Science Cloud." Inf. Services and Use 37 (2017): 49-56. https://doi.org/10.3233/ISU-170824

→ Re3data. https://www.re3data.org/

→ Sigma2. NIRD research data archive. https://archive.sigma2.no/

→ Vines, Timothy H. et al. «The Availability of Research Data Declines Rapidly with Article Age." Current biology 24 (2014): 94-97. http://dx.doi.org/10.1016/j.cub.2013.11.014

→ Zenodo. https://zenodo.org/

→ UB digital https://bibsys-almaprimo.hosted.exlibrisgroup.com/primo-explore/collectionDiscovery?vid=UIO

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

→ Licenses for code: https://github.com/coderefinery/social-coding/blob/main/talk.md

→ Photos from Pixabay if not otherwise indicated. Scientist: RAEng_Publications.

→ 22nd of July research: www.forskningsetikk.no/globalassets/dokumenter/x22.-juli-forskning/sluttrapport-koordineringsgruppa-for-22.-juli-forsking.pdf

→ The Practice of Reproducible Research: http://www.practicereproducibleresearch.org