**TROLLing**
The Tromsø Repository
of Language and Linguistics

# Archiving research data: Whys and hows

Meeting with AcqVA Aurora, 15 November 2021

Helene N. Andreassen & Leif Longva, UiT University Library

# Outline
## for this session
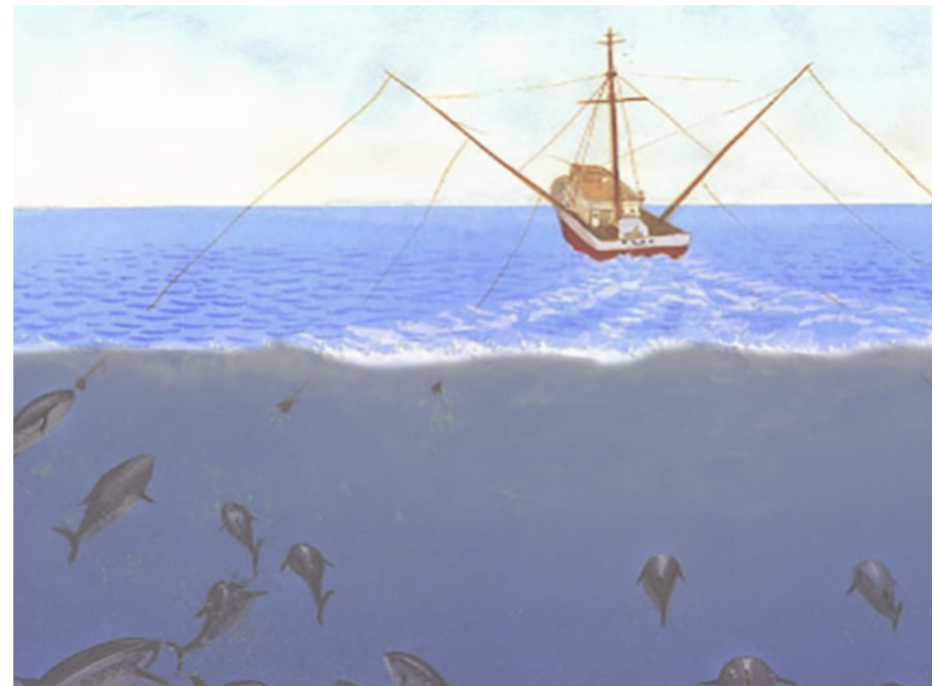
About TROLLing
> Key facts to repeat

Work process
> Preparation
> Creation
> Submission
> Curation
> Revision
> Publishing

# TROLLing: Scope

All subdisciplines of linguistics

International repository: Available to all for upload
and download

# TROLLing: Scope

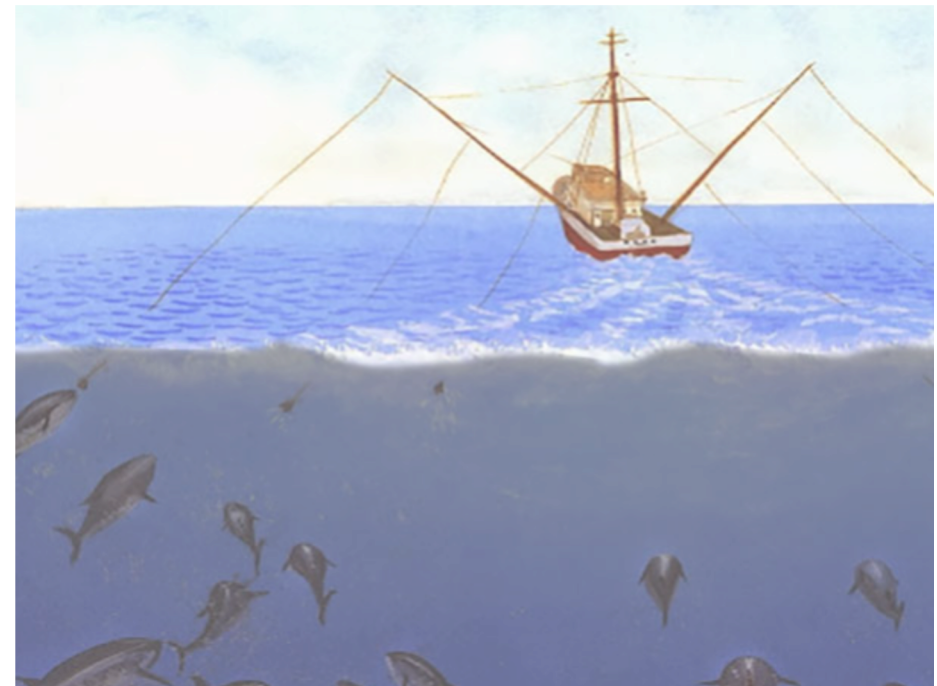All types of data (but open)
    Raw data
    Processed data

All types of supplementary material
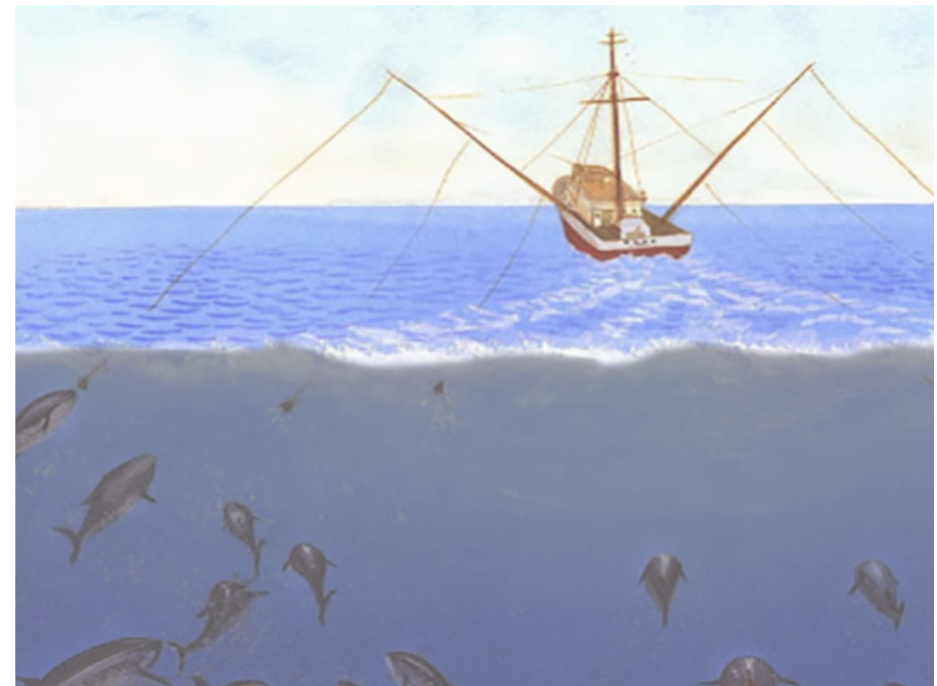    Scripts
    Experimental protocol
    ...

# TROLLing: Infrastructure

Based on the international data platform
Dataverse

Part of DataverseNO, a multi-institutional generic
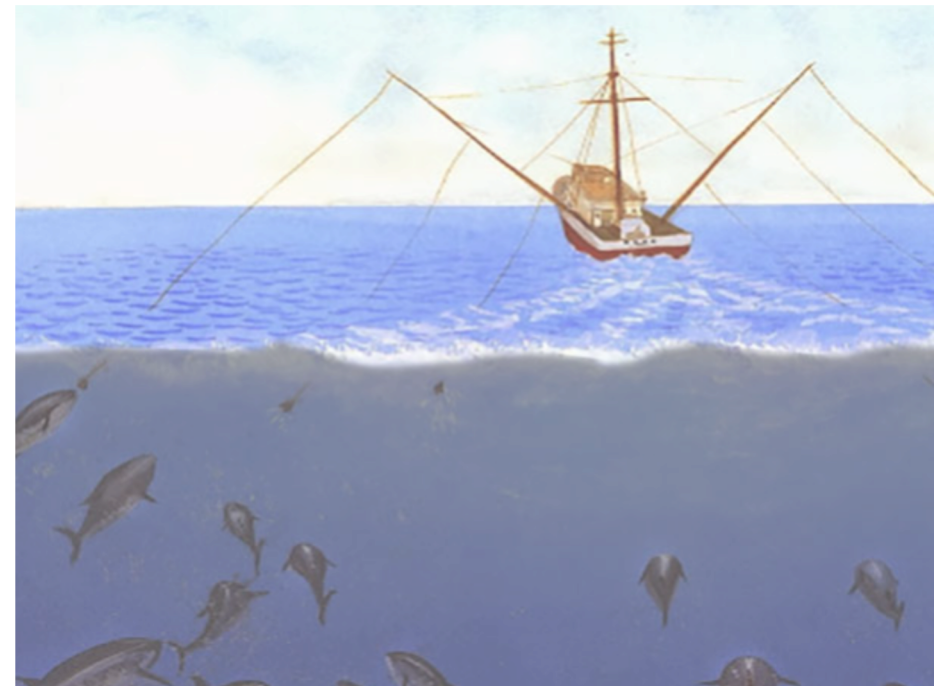repository for open research data

Developed and operated at UiT by the University
Library and the IT Department

# TROLLing: Infrastructure

Operated in line with the FAIR principles (Findable – Accessible – Interoperable – Reusable)

Since 2020, CoreTrustSeal certified as a sustainable and trusted research data repository

# Archiving datasets in TROLLing
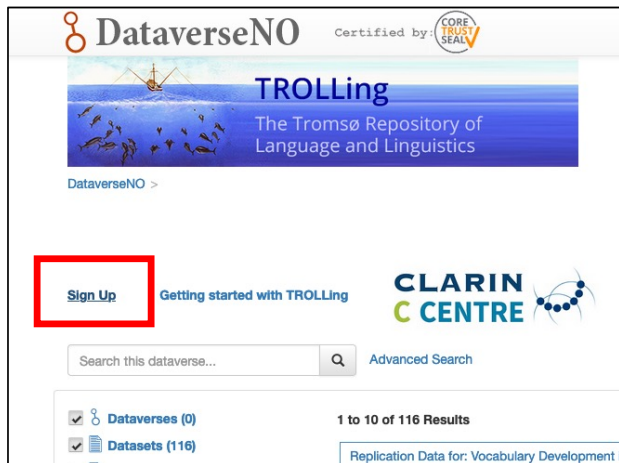
## The work process

# Preparation

Create TROLLing account
Read Deposit Guide
Write readme file
Check reuse restrictions
Format files

# Creating a TROLLing account

All who wish to contribute with data need to be granted status as contributor.

# TROLLing Deposit guide

Deposit guide common for the entire DataverseNO.

Useful to scroll through to see what needs to be prepared before your data can be deposited.
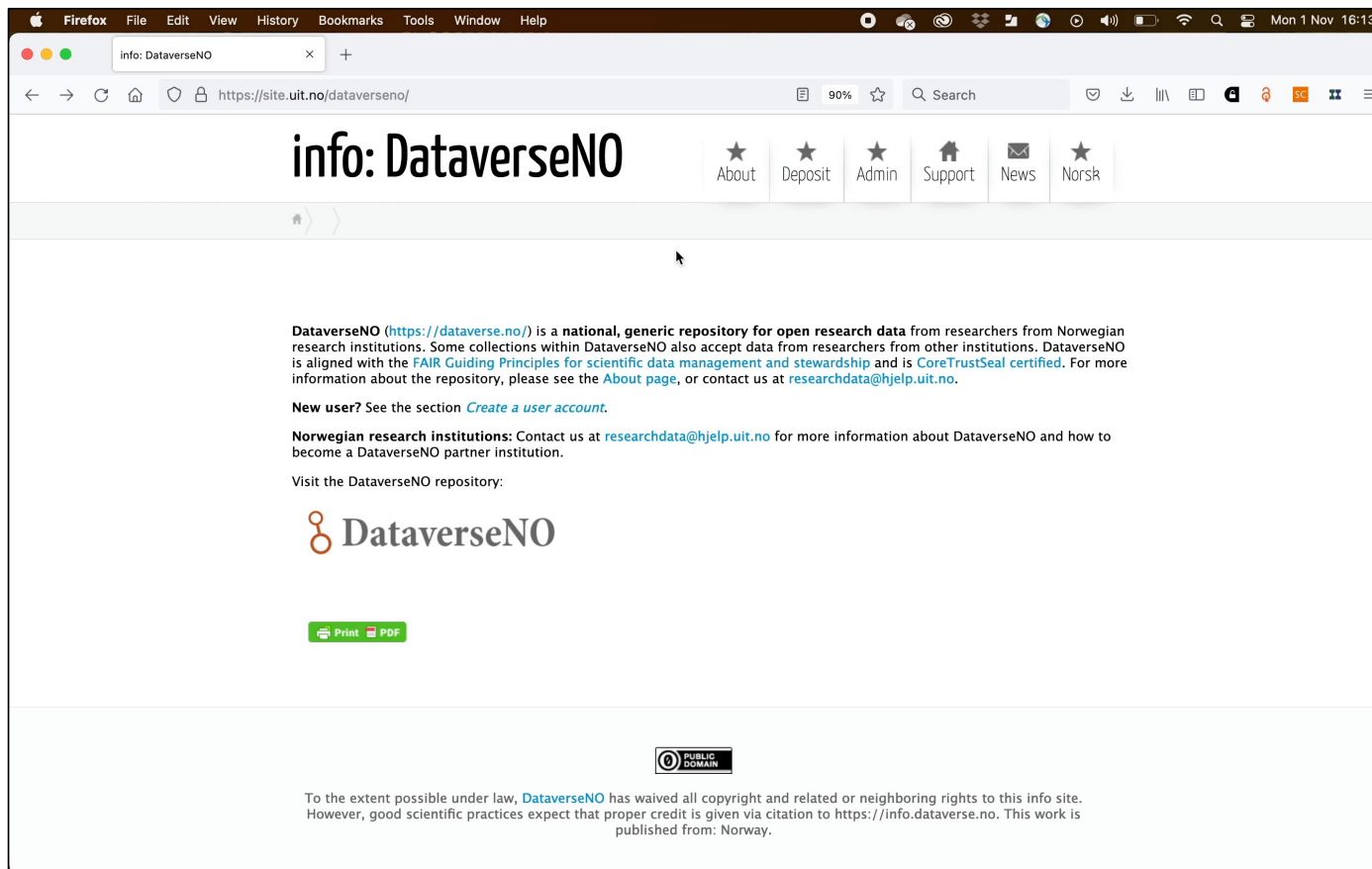
https://info.dataverse.no

# Readme file

- Enough information to understand correctly the data set
- A human readable introduction and description

- General information about the project
  - Title, project period, description, funding sources, project participants, contact information.
- About the dataset
  - License for reuse, related datasets, data sources.

- Methodology
  - Data collection, data processing, data quality etc.
- Content overview
  - File overview, relation between files, version of the dataset.
- Information about the data
  - Column headings, abbreviations, units of measure, contextual information for interpretation of the data set.

# Readme file template

# Readme file examples

Alaee, M., Rasekh-Mahand, M., & Tehrani-Doost, M. (2021). Replication Data for: Eye Behavior during Syntactic Movement Evidence for Processing Approach to Persian Syntax. https://doi.org/10.18710/TZBAOR, *DataverseNO*, V1.

Lybaert, C., De Clerck, B., Saelens, J., & De Cuypere, L. (2021). Replication Data for: A Corpus Based Analysis of V2 Variation in West Flemish and French Flemish Dialects. https://doi.org/10.18710/NSFN2B, *DataverseNO*, V1.

Wang, X. (2021). Replication Data for: Sound symbolism in Chinese children's literature. https://doi.org/10.18710/DCAFEP, *DataverseNO*, V1.

# Reuse restrictions

You are responsible for ensuring that the data uploaded to TROLLing can be shared

GDPR

Third-party data



**Download the transcriptions**
The transcriptions are downloadable, some of them in html format, some in text format.

Read or download the transcriptions:

- User license for the transcriptions
- Transcriptions from CANS

**Please refer to the corpus with this reference:**
Johannessen, Janne Bondi. 2015. The Corpus of American Norwegian Speech (CANS). In Béata Megyesi (ed.): *Proceedings of the 20th Nordic Conference of Computational Linguistics, NODALIDA 2015, May 11-13, 2015, Vilnius, Lithuania. NEALT Proceedings Series 23.*
Download.

**Please also add the corpus URL:**
CANS - Corpus of American Nordic Speech v.3.1: http://tekstlab.uio.no/norskiamerika/english /corpus.html

# File formatting

Characteristics of preferred file formats

- non-proprietary
- open, with documented international standards
- in common usage by the research community
- using standard character encodings (e.g. ASCII, UTF-8)
- uncompressed (space permitting)

# File formatting

| File type | Preferred file formats (examples) | Non-preferred file formats (examples) |
|---|---|---|
| Audio | → Uncompressed and lossless Wav or AIFF (.wav/.aiff)<br>→ Compressed and lossless FLAC (.flac)<br>→ Compressed and lossy Mp3 (.mp3) | → AAC (.m4a)<br>→ Monkey's Audio (.ape)<br>→ Ogg Vorbis (.ogg)<br>→ Windows Media Audio (.wma) |
| Container file | Container files are automatically unpacked when uploaded and should only be used to keep the folder structure in your dataset; see more in section Upload data files. | In case container files need to be archived as container files, use .zip. Note! In this case, files must be packed twice. That way, the inner container will be preserved when uploaded. |
| Image | → Uncompressed TIFF (.tif or .tiff)<br>→ Compressed and lossless PNG (.png)<br>→ Compressed and lossy JPEG (.jpg) | → Adobe Photoshop (.psd)<br>→ Apple Picture File (.pct)<br>→ Graphics Interchange Format (.gif)<br>→ Raw Image Data File (.raw)<br>→ Windows Bitmap (.bmp) |
| Text (slides, illustrations) | → PDF/A (.pdf) combined with original file | → PowerPoint (.pptx) |
| Text (tables) | → Tab separated Unicode plain text (.txt) | → Excel (.xlsx) |
| Text (text) | → Plain text (.txt, .md)<br><br>**If formatting needed:**<br><br>→ XML, PDF/A (.pdf) combined with original file | → Word (.docx)<br>→ HTML |
| Markup language | → XML (.xml)<br>→ HTML (.html)<br>→ Related files: .css, .xslt, .js, .es | → SGML (.sgml)<br>→ Markdown (.md) |
| | **File format:**<br><br>→ PDF/A (.pdf) combined with original file | **File format:** |

https://site.uit.no/dataverseno/deposit/prepare/#what-are-preferred-file-formats

# Creation

Add metadata
Upload data files
Inform about anonymisation (if applicable)

Add more
metadata

Upload files

# File naming

DATE
2016_interview_Inf1
2016_questionnaire_Inf1
2018_interview_Inf1
2018_questionnaire_Inf1

TYPE
Interview_Inf1_2016
Interview_Inf2_2016
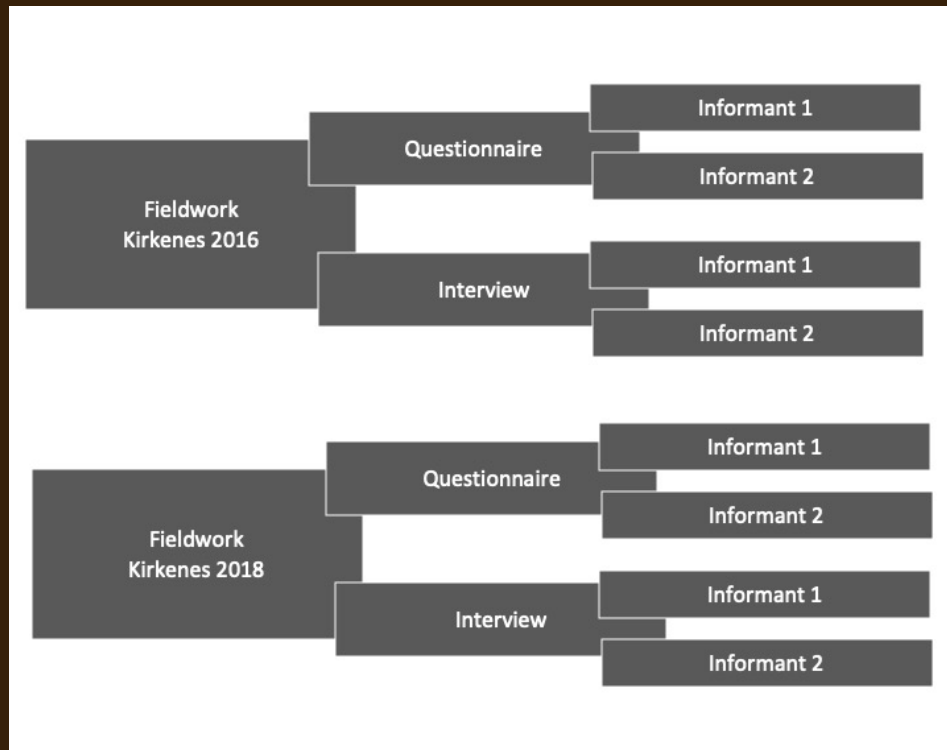Questionnaire_Inf1_2016
Questionnaire_Inf2_2016

INFORMANT (OR TOPIC)
Inf1_interview_2016
Inf1_questionnaire_2016
Inf2_interview_2016
Inf2_questionnaire_2016

FORCED ORDER
01_Overview_fieldwork_2016
02_Inf1_questionnaire_2016
02_Inf1_interview_2016
03_Inf1_questionnaire_2018

# Folder structuring



Creating a folder structure in TROLLing:
https://site.uit.no/dataverseno/deposit/deposit-your-data/#upload-data-files

# Folder structuring

Example:
Strand, Bror-Magnus S., 2021, "Replication data for: Playing with fire compounds", https://doi.org/10.18710/09GQFO, DataverseNO, V1, UNF:6:myhi9R8BuhLnHuc+kHflSw== [fileUNF]

Creating a folder structure in TROLLing:
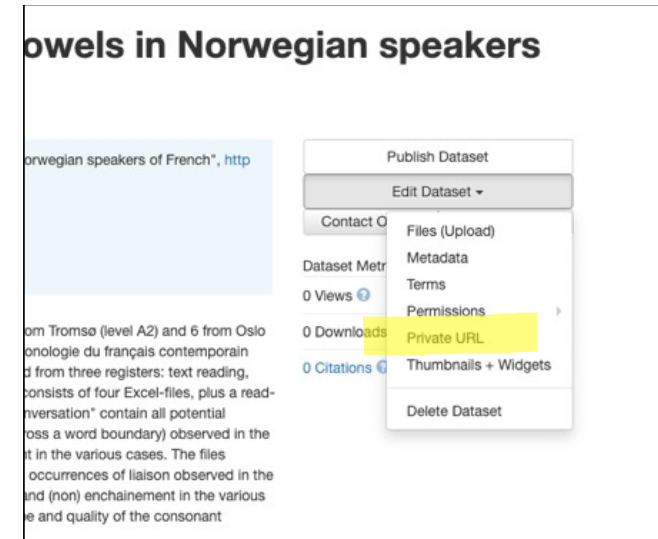https://site.uit.no/dataverseno/deposit/deposit-your-data/#upload-data-files

# Embargo?

Possibility to restrict access to data files for a certain period of time.

This is something you can administer when you upload data files.

# Anonymisation



- Submission of manuscript for double-blind peer review?

- In the *Related Publication* field:
  - Inform that you wish an anonymised version of your dataset.

- Once the dataset is curated:
  - We create an **anonymised version** with a **private URL** that you may send to the editors alongside your manuscript.
  - NB! Your dataset is NOT published yet, and may still be modified (or deleted).

# Submission

When your dataset is ready, you click *Submit for review*

Allow a few days for the data curation process.

# Curation

TROLLing curators: Helene N. Andreassen & Philipp Conzett

(Part of a larger group in the library working on research data management.)

Checking of

Metadata
Readme file
File format
License control

Creation of curator report with suggested modifications

# Revision

Make recommended changes
Re-submit dataset


Rarely more than 3 rounds before all
details are in place

# Publishing

Only curators can publish datasets.

All modifications after initial
publication need to go through/be
approved by us.
    - this includes removal of embargo
      on files.

For anonymised datasets, inform us
when the (non-anonymised) dataset
can be published.

# Publishing

To improve transparency, you should

- link to your published dataset via a
data availability statement

- cite it in your own paper the way you
cite other sources

# Citation of datasets

Citing the dataset in the paper allows the reader to easily access it.

Each TROLLing dataset comes with a fully fledged bibliographic reference.



Want to learn more about data citation in linguistics? Check out the *Tromsø Recommendations for citation of research data in linguistics* (https://doi.org/10.15497/rda00040)

# Plans

- AcqVA Aurora data as a sub-collection within TROLLing

- Own page with TROLLing info on the AcqVA Aurora webpages
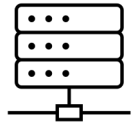  - Links
  - Templates
  - Tutorials

- Move TROLLing to its own installation (outside of DataverseNO)
- Market it to scientific publishers as a trustworthy, discipline-specific repository

## Repository

[trolling.uit.no](trolling.uit.no)

## User guide

[info.dataverse.no](info.dataverse.no)

## Info

[info.trolling.uit.no](info.trolling.uit.no)

## Support

[researchdata@hjelp.uit.no](researchdata@hjelp.uit.no)

*You may also contact [Helene](Helene) or [Philipp](Philipp) directly.*

# Archiving research data: Whys and hows

Meeting with AcqVA Aurora, 15 November 2021

Helene N. Andreassen & Leif Longva, UiT University Library