# Adapting flexible metadata support in Dataverse to the needs of domain-specific repositories
# –
# the case of The Tromsø Repository of Language and Linguistics (TROLLing)

24 November 2021

Philipp Conzett
Helene N. Andreassen

University Library
UiT The Arctic University of Norway

ISKO UK
Knowledge Organization Research Observatory

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

@TROLLingRepo
@PhilippConzett
@n_andreassen

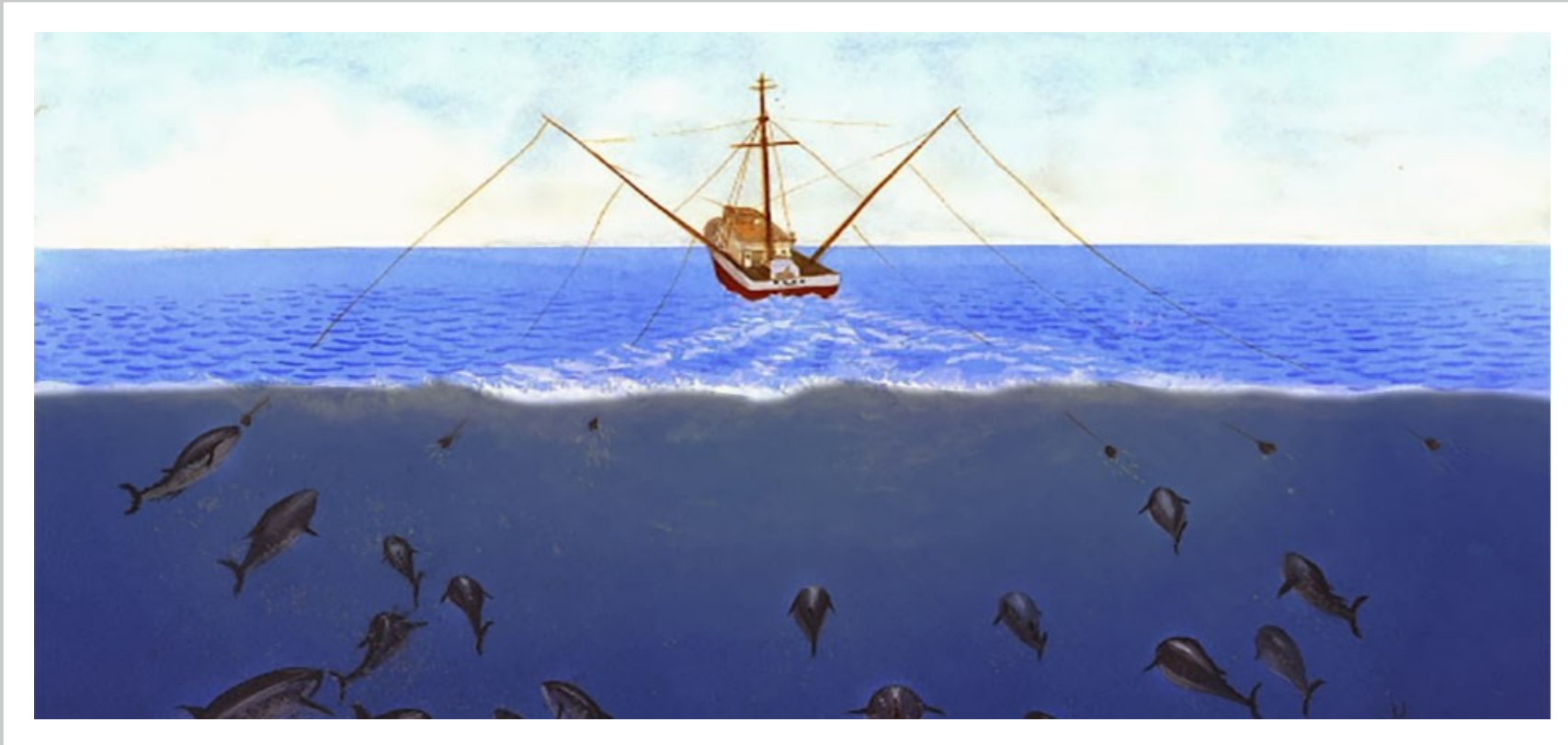CLARIN

SSHOC
social sciences & humanities open cloud

# Outline of the presentation

1. What is TROLLing?
   History, scope, infrastructure, support, numbers
2. Current metadata support in TROLLing
3. Future metadata support in TROLLing

# Part 1: What is TROLLing?

The Tromsø Repository of Language and Linguistics

trolling.uit.no

# TROLLing: history

Pre 2013: UiT University Library providing Open Access publication support.

Fall 2013: The UiT Library was contacted by Laura Janda and Tore Nesset, professors of Russian language at UiT asking for help to create a community-driven repository where linguists worldwide could archive and share their data and code to support the transparency and reproducibility of linguistic studies.

Establishment of working group and development of TROLLing; user guidelines, curation routines, outreach.

June 2014: TROLLing was launched, as (one of) the first open repository for linguistic research data.

**SPARC✳ Europe**

**European Open Data Champions**

Inspiration from influential European academics and information professionals on Open Data

Home » Champions » By sharing our data, and doing this in an open, public, community fashion, we can determine the best practices for our field

**Name:** Prof Laura A. Janda
**Position:** Professor of Russian Linguistics
**Institution:** UiT The Arctic University of Norway
**Country:** Norway
**More info:** Home Page; Other
**ORCID ID:** http://orcid.org/0000-0001-5047-1909

*"By sharing our data, and doing this in an open, public, community fashion, we can determine the best practices for our field"*

# TROLLing: scope

All subdisciplines of linguistics

The international community

All types of data (but open)
    Raw data and processed data
    Text, image, audio, video, …

All types of supplementary material
    Code/scripts
    Experimental protocol
    …



TROLLing
The Tromsø Repository
of Language and Linguistics

# TROLLing: the infrastructure

Based on the community-driven Dataverse software

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

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

For historical reasons still part of DataverseNO, an institution-based national generic repository for open research data. Will be moved to its own Dataverse installation in 2022.

# TROLLing: the infrastructure

Being part of DataverseNO, TROLLing has since 2020 been CoreTrustSeal certified as a sustainable and trusted research data repository.

# Some main technical features:

Baten, Kristof; Van Hiel, Silke; De Cuypere, Ludovic, 2021, "Replication Data for: Vocabulary Development in a CLIL Context: A Comparison between French and English L2.", https://doi.org/10.18710/PXJX1F, DataverseNO, V1

Cite Dataset ▾        Learn about Data Citation Standards.

- ✓ automatically generated reference, including a
- ✓ Permanent identifier (DOI)

| Files | Metadata | Terms | Versions |
|-------|----------|-------|----------|

| Dataset | Summary | Contributors | Published |
|---------|---------|--------------|-----------|
| 2.0 | **Citation Metadata:** Related Publication (1 Changed); **Files (Replaced: 1);** View Details | Philipp Conzett, Tobias Ungerer | Sep 22, 2021 |
| 1.0 | This is the first published version. | Tobias Ungerer, Philipp Conzett | Dec 1, 2020 |

- ✓ Version control

ℹ **Unpublished Dataset Private URL** – Privately share this dataset before it is published:

- ✓ Private URL

## 2_Values_plosives.csv
Comma Separated Values - 47.9 KB
Published Dec 18, 2020
3 Downloads
MD5: 7e4...2b2
Measures for plosives /p t k b d g/ produced in word-initial and word-final position, in a reading task and a repetition task, by informants from the Tromsø and Oslo corpora. The file also contains measures of plosives produced by a native francophone speaker serving as model in the repetition task.

Not available until 2021-04-01

- ✓ Embargo file access

# TROLLing: the infrastructure

Since 2018, TROLLing has been a CLARIN C Centre, and basic citation metadata from TROLLing is harvested by the CLARIN Virtual Language Observatory (VLO).

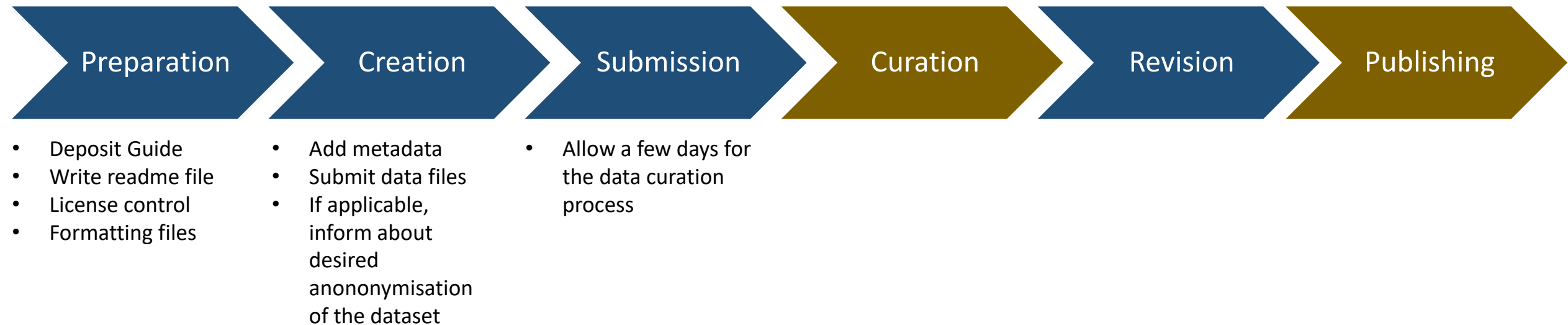Descriptive metadata harvested by more generic search engines such as Google Dataset Search and BASE Bielefeld.

Other search engines:

https://search.datacite.org/
http://b2find.eudat.eu/

# TROLLing: publishing process

- 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.

- Address recommended changes
- Alternatively, explain why you don't agree
- Re-submit

- Metadata quality control
- File format and readme check
- License control

| Preparation | Creation | Submission | Curation | Revision | Publishing |
|---|---|---|---|---|---|

- Deposit Guide
- Write readme file
- License control
- Formatting files

- Add metadata
- Submit data files
- If applicable, inform about desired anononymisation of the dataset

- Allow a few days for the data curation process

# Deposit support

info.dataverse.no

## Deposit guide:

## README file template:

# TROLLing: repository managers and curators

**Helene N. Andreassen**

PhD in French Phonology

Responsible for the UiT training program in research data management

Co-chair of the Linguistics Data Interest Group (Research Data Alliance)

**Philipp Conzett**

MA in Nordic Linguistics

Part of the repository management of DataverseNO

Member of the Steering Committee of the Global Dataverse Community Consortium

Photo: private

**Linguistics Data IG**
☐ **Taxonomy:** Humanities

# TROLLing collaboration

**CLARIN** – Common Language Resources and Technology Infrastructure, a European Research Infrastructure Consortium (ERIC)

**COST** – European Cooperation in Science and Technology: European network for Web-centred linguistic data science

**SSHOC** – Social Sciences and Humanities Open Cloud – a Horizon 2020 project

**RDA** – Research Data Alliance Linguistics Data Interest Group

# TROLLing: numbers

## Contributors

(as of 30 January 2021, when TROLLing reached 100 published datasets)

82 contributing authors

Representing a total of 42 research organizations

From 17 countries in 4 continents

# TROLLing: numbers
(as of 24 November 2021)

## Data

116 datasets containing 3 026 files

39 languages represented

Mostly supporting / replication data (articles and books)

Data from PhD and MA dissertations

Several datasets anonymised and shared with editors/peer reviewers together with a submitted journal or book manuscript

Published datasets

| Year | Value |
|------|-------|
| 2014 | 21 |
| 2015 | 12 |
| 2016 | 17 |
| 2017 | 12 |
| 2018 | 10 |
| 2019 | 7 |
| 2020 | 21 |
| 2021 | 19 |

# TROLLing: numbers
(as of 17 November 2021)

## Usage

In total, 2302 dataset downloads

At average 4.25 downloads per dataset

**Dataset downloads per year**

| Year | Downloads |
|------|-----------|
| 2014 | 225 |
| 2015 | 125 |
| 2016 | 112 |
| 2017 | 121 |
| 2018 | 82 |
| 2019 | 124 |
| 2020 | 286 |
| 2021 | 1227 |

# Part 2: Current metadata support in TROLLing

# Metadata registration in Dataverse

Metadata are registered in two rounds:
  Round 1: all **mandatory (M)** and a few **recommended (R)** fields
  Round 2: other recommended fields and optional fields (e.g. Social Science and Humanities Metadata)
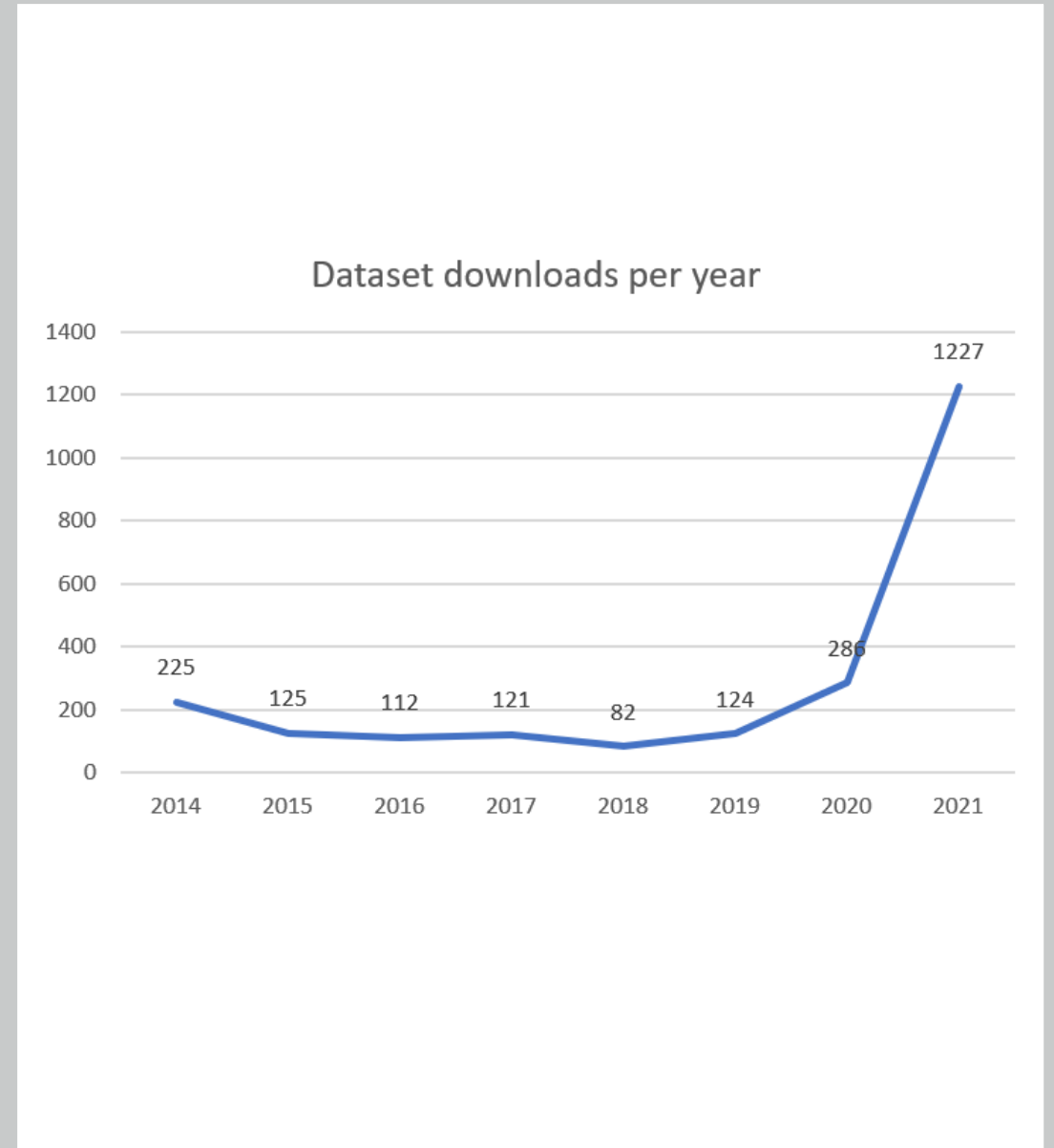
Deposit Guidelines contain more information about the mandatory and recommended fields.

Round 1:

*Citation Metadata:*
- ❏ Title (M)
- ❏ Author (M), ORCID (R)
- ❏ Contact (M)
- ❏ Description (M)
- ❏ Keyword (M)
- ❏ Related Publication (R)

Round 2:

*Citation Metadata:*
- ❏ Language (R)
- ❏ Contributor (R)
- ❏ Grant Information (R)
- ❏ Time Period Covered (R)
- ❏ Date of Collection (R)
- ❏ Kind of Data (R)
- ❏ Related Material (R)
- ❏ Related Dataset (R)
- ❏ Data Sources (R)

*Geospatial Metadata:*
- ❏ Geographic Coverage (R)
- ❏ Geographic Bounding Box (R)

# Need for more domain-specific metadata support

Example 1:

**Language**

Currently: only language of description

Need: also language that is investigated (currently added as keyword)

# Need for more domain-specific metadata support

Example 2:

**Contributor**

<u>Currently:</u> only general/academic contributor roles

<u>Need:</u> also language research-specific roles, e.g., the OLAC Role Vocabulary, as recommended, e.g., in Tromsø Recommendations for Citation of Research Data in Linguistics (https://doi.org/10.15497/rda00040)

**Dataverse Contributor Roles:**
- Data Collector
- Data Curator
- Data Manager
- Editor
- Funder
- Hosting Institution
- Project Leader
- Project Manager
- Project Member
- Related Person
- Researcher
- Research Group
- Rights Holder
- Sponsor
- Supervisor
- Work Package Leader
- Other

**OLAC Role Vocabulary:**
- annotator
- author
- compiler
- consultant
- data_inputter
- depositor
- developer
- editor
- illustrator
- interpreter
- interviewer
- participant
- performer
- photographer
- recorder
- researcher
- research_participant
- responder
- signer
- singer
- speaker
- sponsor
- transcriber
- translator

(Source: http://www.language-archives.org/REC/role.html)

# Need for more domain-specific metadata support

Example 3:

**CMDI compatibility**

<u>Currently:</u> only some basic citation metadata is harvested by CLARIN Virtual Language Observatory (VLO)
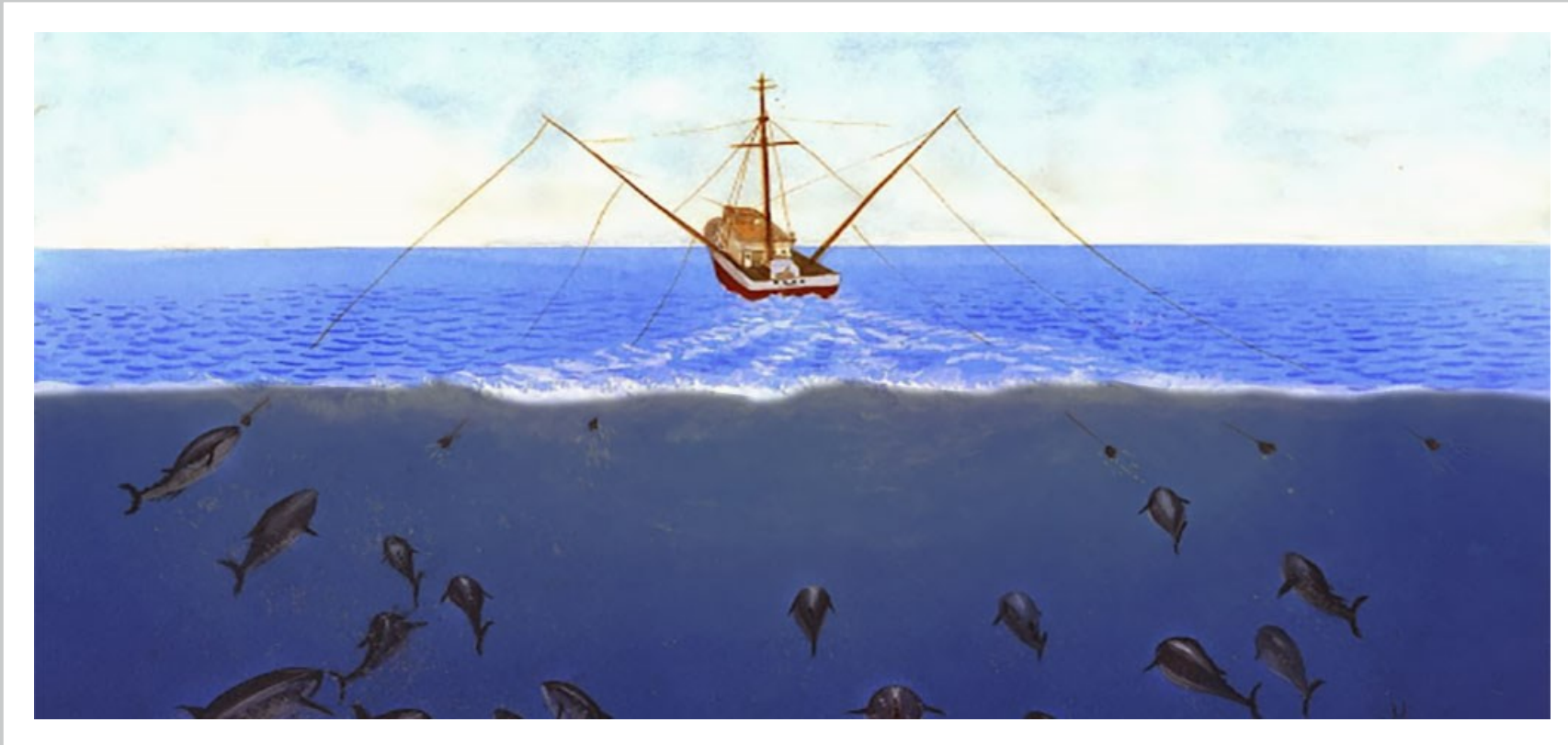
<u>Need:</u> full CMDI compatibility

TROLLing dataset in VLO:

Fully CMDI-compatible dataset in VLO:

# Part 3: Future metadata support in TROLLing

# Domain-specific metadata schema(s)

**Metadata Fields**

Choose the metadata fields to use in dataset templates and when adding a dataset to this dataverse.

☑ Citation Metadata (Required)  [+] View fields + set as hidden, required, or optional

☑ Geospatial Metadata  [+] View fields + set as hidden, required, or optional

☑ Social Science and Humanities Metadata  [+] View fields + set as hidden, required, or optional

☐ Astronomy and Astrophysics Metadata  [+] View fields

☐ Life Sciences Metadata  [+] View fields

☐ Journal Metadata  [+] View fields

☑ Language and Linguistic Metadata
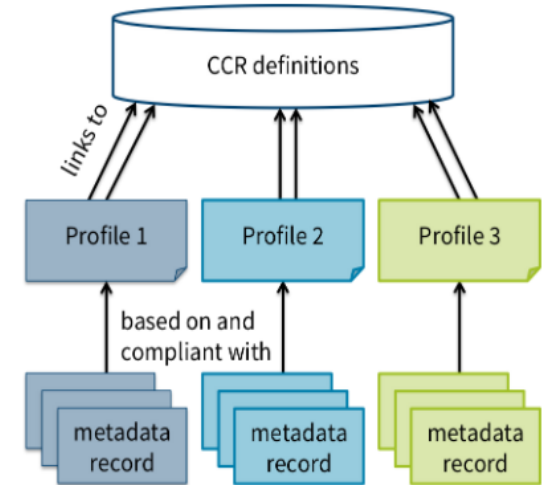
# Language and Linguistic metadata

Examples:

**CLARIN Core Metadata**

CMDI compatible
Recommended by CLARIN metadata WG
(work in progress)

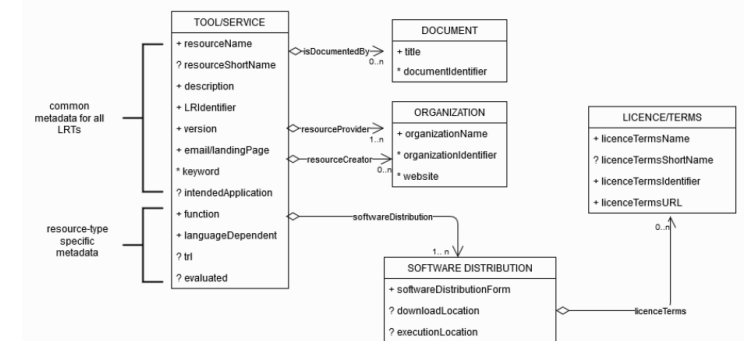**European Language Grid (ELG) Metadata Schema**

ELG = "primary platform for Language Technology in Europe"

CLARIN CMDI



Source:
https://www.clarin.eu/content/component-metadata

ELG Metadata Schema



Source: https://european-language-grid.readthedocs.io/en/release1.1.1/all/A1_Metadata/Metadata.html

# Language and Linguistic metadata

Examples:

**External Controlled Vocabularies**

- OLAC Role Vocabulary
- META-SHARE Ontology, e.g., modalityType
- ...

META-SHARE Ontology:
modalityType



Source: http://w3id.org/meta-share/meta-share

# Challenges

1. How to implement complex metadata schemas (e.g. ELG)?

2. How to ensure maintenance of (complex) metadata schemas?

3. How to ensure sustainability of external controlled vocabulary services?

4. How to support interoperability on file-level?



1036 metadata fields.
Only one out of 10(?) ELG sub-schemas!

# Possible approaches

1. Use CLARIN Core Metadata for (small) supporting/replication datasets; use ELG Metadata for larger resources such as corpora.

2. Formalize and strengthen the role of the Global Dataverse Community Consortium (GDCC) to maintain Dataverse-related resources.

3. Use recognized vocabulary services, or if not available, have them run them by CLARIN, GDCC or another suitable organization.

4. For tabular data, consider adopting the Cross-Linguistic Data Formats initiative (CLDF).



### Cross-Linguistic Data Formats, advancing data sharing and re-use in comparative linguistics

Robert Forkel ✉, Johann-Mattis List ✉, Simon J. Greenhill, Christoph Rzymski, Sebastian Bank, Michael Cysouw, Harald Hammarström, Martin Haspelmath, Gereon A. Kaiping & Russell D. Gray

*Scientific Data* **5**, Article number: 180205 (2018)    | Cite this article

# Thank you for your attention!

Philipp Conzett
Helene N. Andreassen

University Library
UiT The Arctic University
of Norway

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

CLARIN

SSHOC
social sciences & humanities open cloud