# The role of CLARIN in supporting SSH scholars

Darja Fišer

21 November 2025

National EOSC Tripartite Event

Ljubljana, Slovenia

CLARIN

# The CLARIN Dream

*A humanities researcher in Maribor consults the CLARIN catalogue and:*

- finds relevant text collections in data repositories in Slovenia, Austria and Hungary
- a named-entity recognition tool in the Czech Republic
- performs research on this material, with support from a centre of expertise in the Czech Republic
- uses them to build an enriched version of the text collections, benefitting from the CLARIN interoperability standards
- deposits the annotated text collections in a data repository in Slovenia for later use by others

**… and can do all this from their own desk**

# CLARIN in a Nutshell

1. *'Common Language Resources and Technology Infrastructure'*

2. **ESFRI** roadmap (2006), **ESFRI** ERIC status (2012), Landmark (2016)

3. Easy and sustainable access for scholars in Social Sciences and Humanities (SSH):

   4. **digital language data** (written, spoken, video or multimodal)

   5. **tools** to discover, analyse, combine data wherever they are located

   6. **single sign-on** environment (**you all can get an account**)

7. Ecosystem for **knowledge exchange**

8. Member of the **SSH Open Cluster**

# CLARIN and Open Science

- Promotion of sharing and re-use of language data through **sustainable data registries**
- Adherence to **FAIR data principles**
- Enhancement and deployment of **interoperability** of **language data & services**
  - **common metadata framework**
  - distributed network of FAIR **certified data repositories** for language data

- Promotion of
  - studying language data from a **comparative** perspective
  - **multidisciplinary** collaboration
  - **transnational** research
  - **responsible** data science
- Support for **linguistic diversity**
  - **data** covering many languages
  - **tools** for many languages
  - language resources in **all modalities**
  - **discipline-** and **language-agnostic**

# CLARIN's Countries

A distributed **European Infrastructure Consortium(ERIC)** consisting of:

- **27** members

# CLARIN's Centres (August 2025)

**CLARIN B-centre**

## 22 CTS-certified data centres

Strong focus on **FAIRness** & **interoperability**:

- Federated login
- Central metadata harvesting for easy discovery
- Chained services

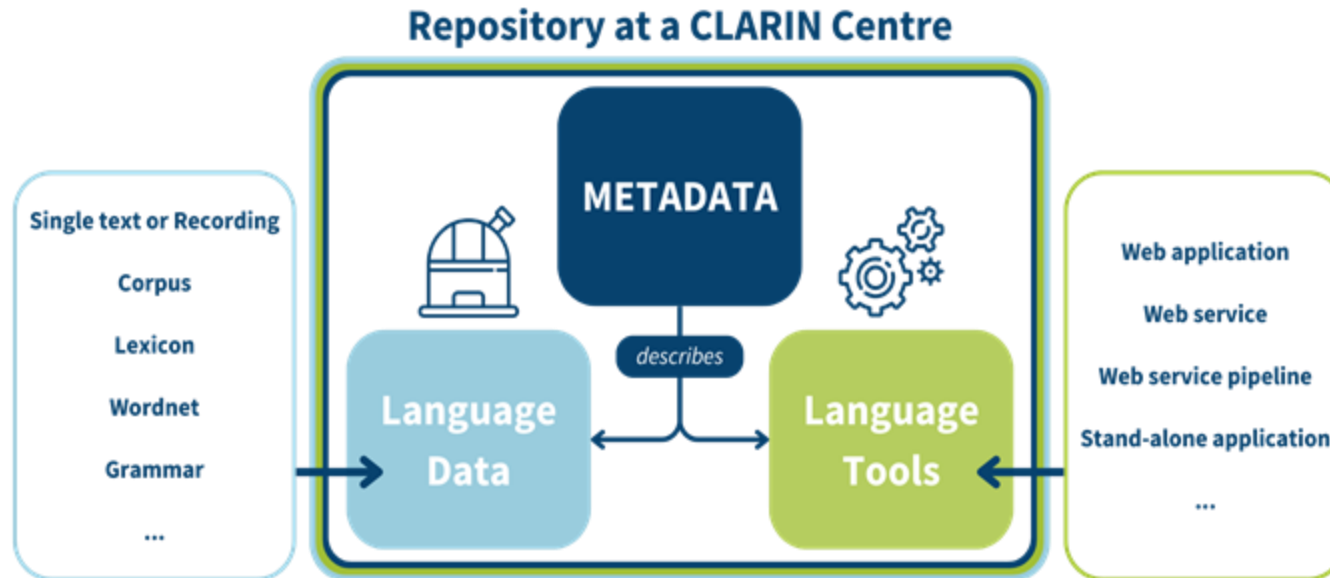👉 **Search for B-centres**

**CLARIN K-centre**

## 39 Knowledge centres

Operated by a single institute/group or as a distributed structure **covering a large number of research topics, languages and resource types.**

👉 **Search for K-centres**

# How CLARIN Works



Repository at a CLARIN Centre

METADATA

describes

Language Data

Language Tools

Single text or Recording

Corpus

Lexicon

Wordnet

Grammar

…

Web application

Web service

Web service pipeline

Stand-alone application

…

# How CLARIN Works

Home / Web Applications

# Web Applications

These shortcuts provide quick access to the central CLARIN web applications:

**Content Search**

Federated Data Search Engine

**Language Resource Switchboard**

Find a suitable tool to process language data

**Virtual Language Observatory**

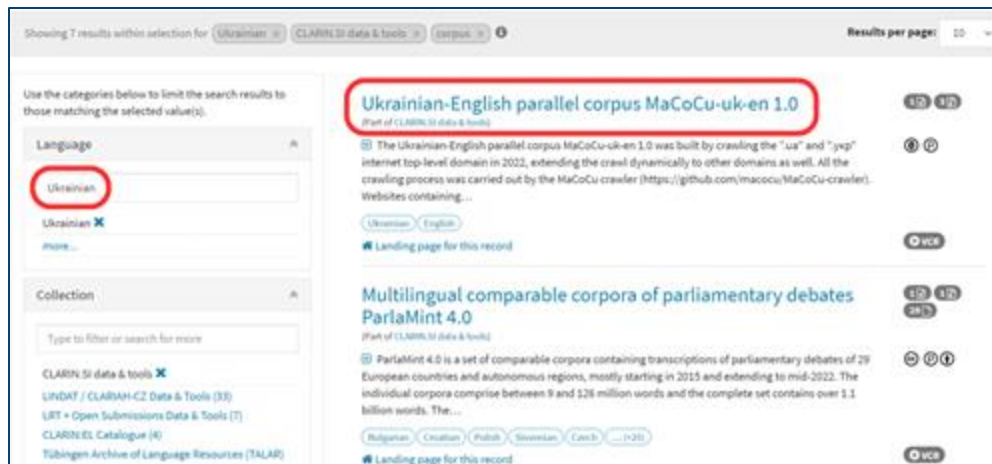Metadata search interface

**Virtual Collection Registry**

Publish and access digital bookmarks

# Virtual Language Observatory (VLO)

– Facet search to explore **>500,000 language resources, tools, services**

– Metadata to describe resources

– **Persistent Identifiers (HDL)** to access the landing pages of the resources and cite them

– Availability and/or **Licencing**

– Technical details

– Process with **Switchboard**

– Compile search results into a **Virtual Collection**

# Switchboard



https://voyant-tools.org

# Flagship Project - Parlamint

Parliamentary proceedings of 29 European parliaments

Corpora linguistically annotated, named entities, rich metadata

Great example of CLARIN collaboration and interoperability

The ParlaMint corpora are **interoperable**:

- annotated within a specifically developed TEI-based schema

The ParlaMint corpora are **comparable:**

- the same overlapping periods: 2015–2022
- the same metadata (related to speakers, parties, sessions, reactions, etc.)
- the same types of linguistic annotations

Interoperability and comparability are further increased also by **machine translations** into English

# Thank you!
# Questions & comments welcome!

www.clarin.eu
darja.fiser@clarin.eu

CLARIN