# CLARIN per le Lingue Antiche

Workshop online

20 marzo 2024

Giulia Pedonese

# Sources

This presentation is the result of the adaptation and modification of the following sources:

- van der Lek, Iulianna; Fišer, Darja. (2023). *Introduction to Language Data: Standards and Repositories.* In UPSKILLS Learning Content. https://upskillsproject.eu/project/standards_repositories/. CC BY 4.0.

- van der Lek, I., Fišer, D., Samardzic, T., Simonovic, M., Assimakopoulos, S., Bernardini, S., Milicevic Petrovic, M., & Puskas, G. (2023). Integrating research infrastructures into teaching: Recommendations and best practices (Versione 2). Zenodo. https://doi.org/10.5281/zenodo.8114407. CC BY 4.0.

- CLARIN ERIC Official Website: https://www.clarin.eu/ . CC BY 2.0

# What Is CLARIN and How Can You Access It?

**CLARIN** stands for **Common Language Resources and Technology Infrastructure**

CLARIN is a distributed digital infrastructure, with participating centres all over Europe and further afield, which include universities, research centres, libraries and public archives. Tools and data from different centres are interoperable so that data collections can be combined and tools from different sources can be chained to perform operations at different levels of complexity, regardless of their location.

# CLARIN...

 Belongs to the Social Sciences and Humanities cluster and an integral part of the [European Open Science Cloud](#)

 Has the ESFRI ERIC status since 2012, Landmark since 2016

 Provides easy and sustainable access for scholars in the humanities and social sciences and beyond:

- to digital language data (in written, spoken or multimodal form)
- to advanced tools to discover, explore, exploit, annotate, analyse or combine them
- through a single sign-on environment

 Serves as an ecosystem for knowledge sharing and training

# CLARIN-IT: the Italian node of CLARIN

Italy became the 16th Full Member of CLARIN ERIC in 2015. The Founding member of the National Consortium is the Institute for Computational Linguistics "Antonio Zampolli" of the National Research Council of Italy.

CLARIN is involved in the H2IOSC project through its Italian national consortium CLARIN-IT.

# The H2IOSC project

The H2IOSC project aims at creating a federated and inclusive cluster of RIs in the ESFRI domain of Social and Cultural Innovation to allow researchers from various disciplines in the Humanities, Language technologies and the Cultural Heritage sectors collaborate in data and compute intensive research.

It encompasses the Italian nodes of four Research Infrastructures:

☐ CLARIN

☐ DARIAH, Digital Research Infrastructure for the Arts and Humanities

☐ E-RHIS, European Research Infrastructure for Heritage Science

☐ OPERAS, Open Scholarly Communication in the European Research Area for SSH

# The H2IOSC project



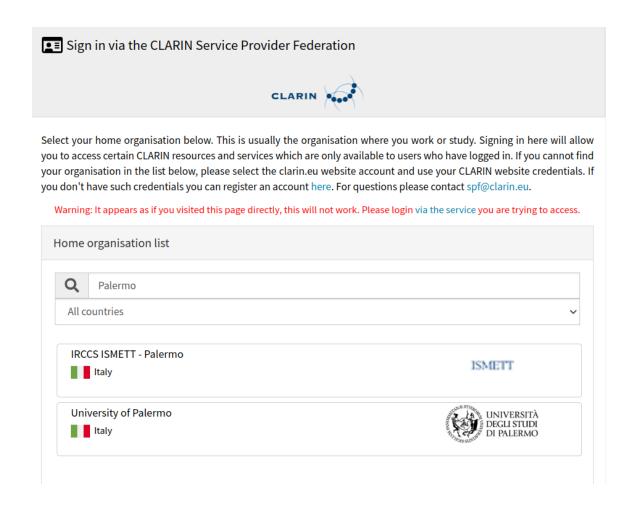## Research Infrastructures Targeted by the Project

**CLARIN**

The research infrastructure for language as social and cultural data

Learn more

DARIAH-EU
Digital Research Infrastructure for the Arts and Humanities
DARIAH

E-RIHS
EUROPEAN RESEARCH INFRASTRUCTURE FOR HERITAGE SCIENCE
E-RIHS

OPERAS
open scholarly communication in the european research area for social sciences and humanities
OPERAS

# CLARIN core services

- Depositing services to make sure that language resources can be archived and made available to the community in a reliable manner and to help researchers to store their resources in a sustainable way

- The Virtual Language Observatory provides an easy-to-use interface, allowing for a uniform search and discovery process for a large number of resources from a wide variety of domains.

- The Federated Content Search is a search engine that connects to the local data collections that are available in the centres

- The Language Resource Switchboard helps users to find a matching language processing web application for your data

- The Virtual Collection Registry provides a registry where scholars can create and publish their virtual collections

# How to access CLARIN services

- All users can freely explore the CLARIN core services to search for language resources and expertise

- Due to license restrictions, some resources are only available for academic users and login is required using your institutional credentials or CLARIN credentials

- Academic users in all participating countries can access and use the language resources available in CLARIN data centers with a single sign-on access through the CLARIN Service Provider Federation using their institutional credentials

# Sign in via the CLARIN Service Provider Federation

# CLARIN account registration

If your university or institute is not in the list of federated organisations, you can request a CLARIN account

# Depositing services

Many of the CLARIN centres offer a depositing service. They are willing to store the resources in their repository and assist with the technical and organisational details. This service ensures many advantages:

- Long-term archiving: a storage guarantee can be given for a long period (up to 50 years)

- Resources can be cited easily with a persistent identifier

- The resources and their metadata will be integrated into the infrastructure, making it possible to search them efficiently

- Password-protected resources can be made available via an institutional login

- Once resources are integrated in the CLARIN infrastructure, they can be analysed and enriched more easily with various linguistic tools

# The CLARIN Repositories

- Ensure long-term preservation and curation of language resources, datasets and tools

- Provide specific and specialised metadata to describe language resources

- Assign Persistent Identifiers (PID), e.g. Handle, to the deposited resources, which enable easy citation

- Offer advanced services to explore the language resources and their metadata, e.g. Virtual Language Observatory, Content Search, Switchboard tools, integrated corpus query engines (e.g. PML-TQ and Kontext)

- Academic users can access restricted resources with their institutional credentials

# How to use a repository: the example of ILC4CLARIN

ILC4CLARIN is CLARIN-IT B-centre hosted at Institute for Computational Linguistics, National Research Council, in Pisa. It offers depositing services for language datasets and tools for research, especially for Italian and classical languages via its repository

The ILC4CLARIN repository is a disciplinary repository certified by CoreTrust Seal and it offers advanced services to explore the language resources and their metadata (e.g. VLO, Switchboard ecc.)

# ILC4CLARIN repository: services

# ILC4CLARIN language search: Ancient Greek

# Example of a corpus in ILC4CLARIN Repository



- Citation information via handle

- Metadata fields describing the corpus

- The corpus is referenced in a journal

- The corpus is described, e.g.
texts available in UTF-8 format and TEI-XML format

- Publisher information

- Download instructions

- All the files can be