



the **Italian** Common Language Resources and Technology Infrastructure



CLARIN ERIC and CLARIN-IT

Francesca Frontini - CLARIN Board of Directors

Istituto di Linguistica Computazionale - ILC CNR

24/03/2023



Common Language Resources and
Technology Infrastructure



Open Data & FAIR principles

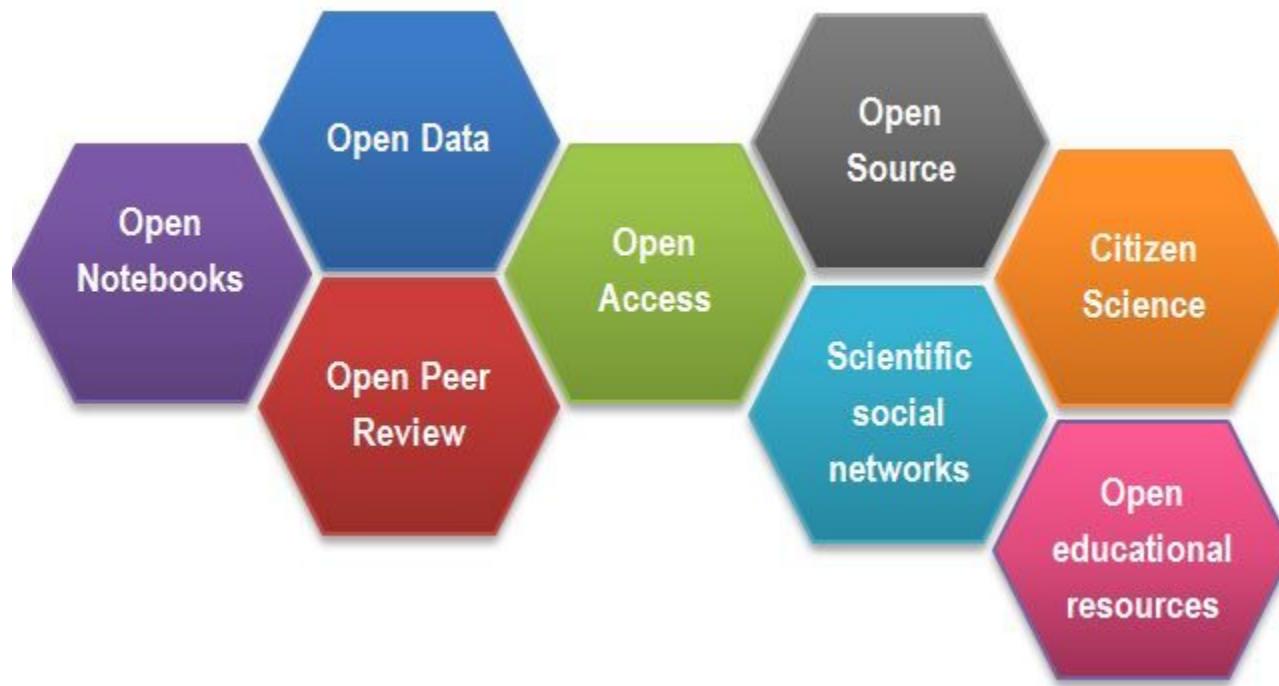
A New Paradigm



"Open Science represents a **new approach to the scientific process** based on cooperative work and new ways of diffusing knowledge by using digital technologies and new collaborative tools ([European Commission, 2016](#))"

Open science encompasses unhindered access to scientific articles, access to data from public research, and collaborative research enabled by ICT tools and incentives.
([OCSE, 2015](#))

Open Science Practices



Source of the image and further information:

<https://www.fosteropenscience.eu/learning/open-science-at-the-core-of-libraries/#/id/5a01e2d1c2af651d1e3b1b3c>

FAIR Data



FAIR Principles

Compliance



Findability

Resource and its metadata are easy to find by both, humans and computer systems. Basic machine readable descriptive metadata allows the discovery of interesting data sets and services.

- ✓ F1. Resource is uploaded to a public repository.

- ✓ F2. Metadata are assigned a globally unique and persistent identifier.



Accessibility

Resource and metadata are stored for the long term such that they can be easily accessed and downloaded or locally used by humans and ideally also machines using standard communication protocols.

- ✓ A1. Resource is accessible for download or manipulation by humans and is ideally also machine readable.

- ✓ A2. Publications and data repositories have contingency plans to assure that metadata remain accessible, even when the resource or the repository are no longer available.



Interoperability

Metadata should be ready to be exchanged, interpreted and combined in a (semi)automated way with other data sets by humans as well as computer systems.

- ✓ I1. Resource is uploaded to a repository that is interoperable with other platforms.

- ✓ I2. Repository meta- data schema maps to or implements the CG Core metadata schema.

- ✓ I3. Metadata use standard vocabularies and/or ontologies.



Reusability

Data and metadata are sufficiently well-described to allow data to be reused in future research, allowing for integration with other compatible data sources. Proper citation must be facilitated, and the conditions under which the data can be used should be clear to machines and humans.

- ✓ R1. Metadata are released with a clear and accessible usage license.

- ✓ R2. Metadata about data and datasets are richly described with a plurality of accurate and relevant attributes.

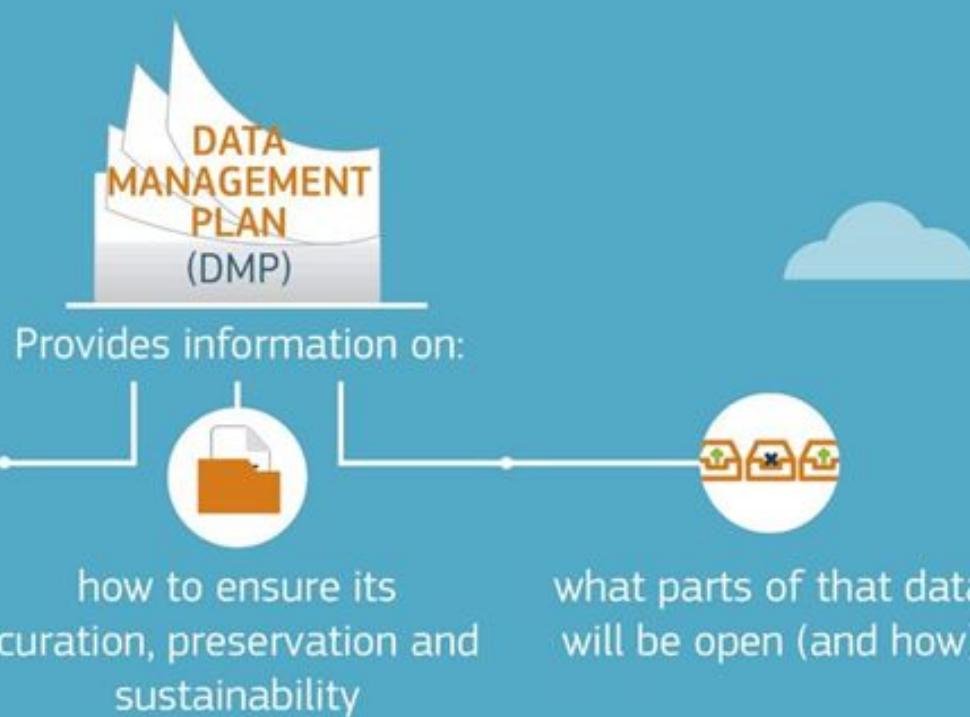
'FAIR Guiding Principles for scientific data management and stewardship', 2016
<https://www.go-fair.org/fair-principles/>

A European Strategy



RESEARCH DATA - OPEN BY DEFAULT

Projects must have



How to ensure ...



- long term deposit
- metadata quality
- findability
- citability
- clear licenses
- interoperability?



<https://www.youtube.com/watch?v=lfDWBaaAclw>

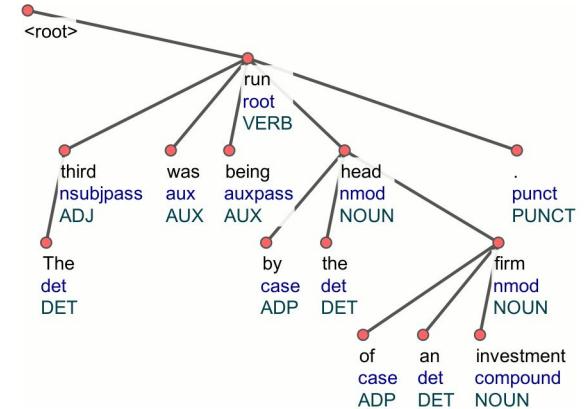
The CLARIN Research Infrastructure for data sharing

CLARIN & language resources



Common Language Resources and
Technology Infrastructure

<https://www.clarin.eu/>



CLARIN in a nutshell



- 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, video or multimodal form)
 - and **advanced tools** to discover, explore, exploit, annotate, analyse or combine them, wherever they are located
 - through a **single sign-on** environment
- serves as an ecosystem for **knowledge sharing and training**
- is an integral part of the **European Open Science Cloud**
 - See clarin.eu/eosc

CLARIN data and communities



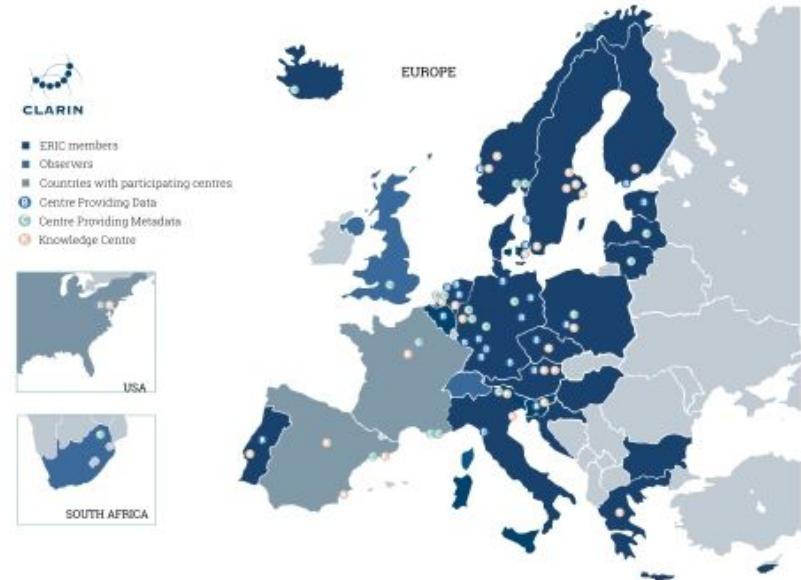
- Newspaper archives
 - Literary texts
 - Parliamentary records
 - Literary texts
 - Historical letters
 - Broadcast archives
 - Oral History data
 - Social Media data
 - L2 Learner Resources
 - Survey data
 - ...
- Digital humanities
 - Linguistics and Philology
 - Translation and Lexicography
 - Literary Studies
 - History
 - Political and Social Sciences
 - Media Studies
 - Culture, Folklore, Anthropology
 - Speech therapy
 - General Public
 - ...

For the CLARIN Resource Families initiative, see:
<https://www.clarin.eu/resource-families>

CLARIN today



- **22 members:** (AT, BE, BG, CY, CZ, DE, DK, EE, FI, GR, HR, HU, IS, IT, LT, LV, NL, NO, PL, PT, SE, SI)
- 3 observers: CH, UK, ZA
- > **70 centres**



Common Language Resources and
Technology Infrastructure

CLARIN types of centres



- **B-Centres**
 - These are the distributed data centres that universities or academic institutions can access to search for or deposit language resources, access services and expert knowledge. There are several [certified B-centres](#) that have passed the [CLARIN centre assessment procedure](#). The CLARIN-B centres are required to apply for the [CoreTrustSeal](#) certification.
- **C-Centres**
 - These are centres that provide metadata to CLARIN but they do not offer any other services.
- **K-Centres**
 - Knowledge centres that share knowledge and expertise on one or more aspects of the domains covered by the CLARIN infrastructure.

....Consortia and Centres offer Data, Tools, Services



Search Catalogue Education Projects Tools Services About ▾

Digital Research Infrastructure for the Language Technologies, Arts and Humanities

Catalogue Corpora Treebanks ČDK Bibliography

search by type of data

Search

e.g. [corpus](#) or [lexicon](#) or [editor](#)

<https://lindat.cz/>

....hosting repositories

The screenshot shows the homepage of the LINDAT CLARIAH-CZ repository. At the top, there is a navigation bar with links for Search, Catalogue, Education, Projects, Tools, Services, and About. Logos for Dariah-EU and CLARIN are also present. Below the navigation bar, there is a search interface featuring a magnifying glass icon and a search input field with a "Search" button. To the right of the search bar are the LINDAT and CLARIN logos. The main content area has sections for "Author", "Subject", and "Language (ISO)" with lists of items and "View More" links.

Find

Linguistic Data and NLP Tools
Citation Support (with Persistent IDs)

[Advanced Search](#)

Author	Subject	Language (ISO)
Veselý, Bohumil (787)	People (803)	Nolinguistic content (712)
Hajič, Jan (89)	Galerie osobností (787)	Czech (486)
Aktualita (78)	Places (555)	English (313)
Straka, Milan (67)	machine translation (63)	German (216)
Krátký film (63)	Český zvukový týdení ... (51)	French (113)
... View More	... View More	... View More

<https://lindat.mff.cuni.cz/repository/xmlui/?locale-attribute=en>

.... offering tools for exploration and processing

Tools for Natural Language Processing

Enrich your texts
Make them more searchable

?

Terms of use
Service Status

Phonetics, Phonology

- [UWebASR](#) (audio transcription, Czech)

Machine Translation

- [LINDAT translation](#)

Morphology and tagging

- [ElixirFM](#) (Arabic)
- [MorphoDiTa](#) (Czech)
- [UDPipe](#)

Natural language processing

- [Treex](#) (Czech, English, Latin)
- [UDPipe](#)

Syntactic parsing

- [Korektor](#) (spellchecker, Czech)
- [Parsito](#)
- [UDPipe](#)

Search

- [CzEngVallex](#) (lexicon, Czech, English)
- [EngVallex](#) (lexicon, English)
- [PDT-Vallex](#) (lexicon, Czech)
- [Dialogy.org](#) (audio-visual corpora search)
- [Internet Language Reference Book](#) (Czech)
- [KonText](#) (concordances, collocations, word frequencies)
- [TEITOK](#) (search, visualize, edit corpora)
- [PML Tree Query search](#) (treebank search)

Discourse, Pragmatics

- [Evald](#) (coherence evaluation, native speakers of Czech)
- [Evald](#) (coherence evaluation, non-native speakers of Czech)
- [KER](#) (keyword extraction, Czech, English)
- [NameTag](#) (named-entity recognition, Czech, English)

<https://lindat.cz/#tools>



Read more about CLARIN-UK

Data & Tools

[More](#)

CorCenCC

Corpus
Cenedlaethol
Cymraeg
Cyfoes – the
National Corpus
of
Contemporary
Welsh

Latest News

[More](#)

New
members of
the
CLARIN_UK
consortium

18 November
2020
Welcome
aboard!

Events

[More](#)

Lancaster
Symposium
on Innovation
in Corpus
Linguistics
2021

Wednesday 23
June, Online



Read more about CLARIN-UK

Data & Tools

[More](#)

CorCenCC

Corpus
Cenedlaethol
Cymraeg
Cyfoes – the
National Corpus
of
Contemporary
Welsh

#LancsBox

BNC

CKLD

CLAWS

CLiC

CorCenCC

CQPweb

ELAR

GATE

GATE Cloud

Hansard at
Huddersfield

[More](#)

Events

[More](#)

Lancaster
Symposium
on Innovation
in Corpus
Linguistics
2021

Wednesday 23
June, Online



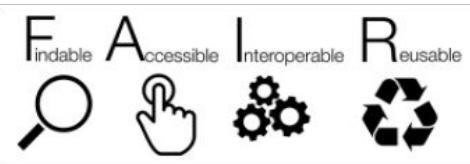
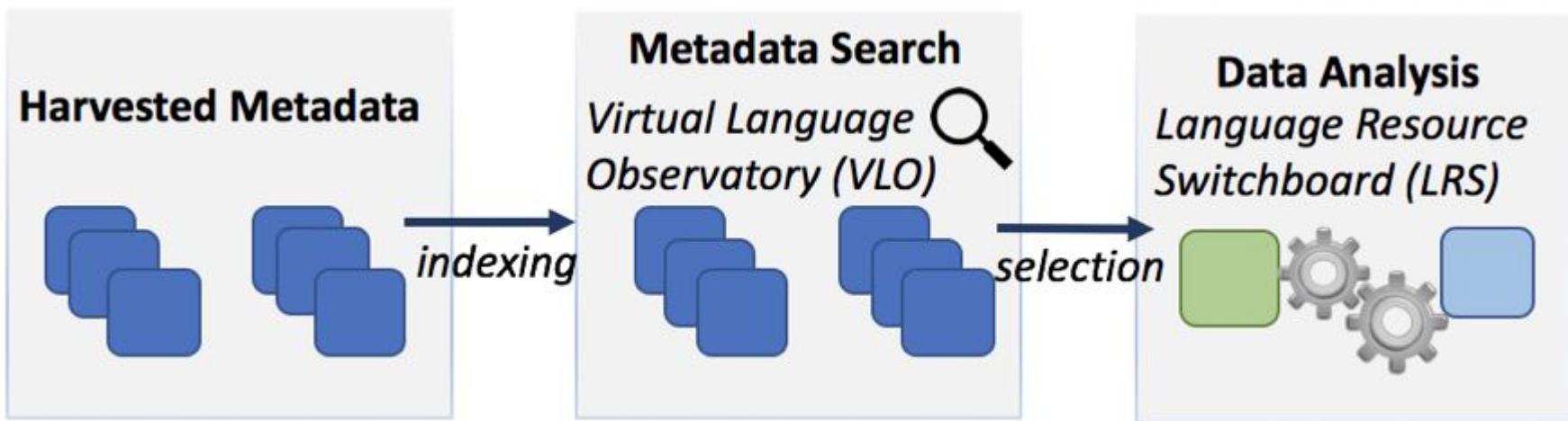
Corpus
Linguistics



The CLARIN portal: A single access point to the distributed infrastructure

- Make it easier for researchers, lecturers, students to **find** data, tools, training materials from centres
- Make it easier to **share and publish** data following open science practices
- Allow you to be **cited** for your work and acknowledge that of others
- Promote **knowledge sharing** and networking

The technical infrastructure



clarin.eu/fair



vlo.clarin.eu



switchboard.clarin.eu

Barack Obama's identity-building in the health care debate: A corpus-assisted discourse study

AUTHOR

Katherina Riesner

Summary, in English

In this study, I demonstrate that identity-building is an important discursive strategy for President Barack Obama in the seven-year long debate surrounding the Affordable Care Act (ACA). The data for the study comes from a 6-million word corpus of speeches that were held by Obama between January 2009 and January 2016, all published by the White House. The speeches are classified according to genre, audience, topic and date of delivery. Throughout the paper, I adopt the notion that identity is intentionally constructed by the speaker and strategically exploited for his communicative goals. With the help of two methodological approaches, I investigate what kind of identities Obama builds. The purely qualitative part of the study deals with three central corpus speeches from a discourse-analytic perspective. In the second, more quantitative part, I use a group of seven verbs with epistemic meaning to trace the usage of two predominant discursive identities in the ACA debate. The results suggest that President Obama repeatedly constructs the identities of father and teacher to persuade his audience. I argue that his use of these identities constitutes an attempt to reach the argumentative goals of effectiveness and reasonableness.

Department/s

Master's Programme: Language and Linguistics

Publishing year

2016

Language

English

Full text

[Available as PDF - 2 MB](#)

[Download statistics](#)

Document type

Student publication for Master's degree (two years)

Topic

Languages and Literatures

-  IMDI Corpora
-  Lund Corpora
-  Eline Visser
-  ESST
-  Eye-Tracked Frog Stories
-  LACOLA
-  LANG-KEY
-  LUNDIC
-  REaCHeS
-  SpaceH
-  Strömqvist-Richthoff
-  Swedia2000
-  Tactile Reading
-  Test
-  ThaiSweVideo
-  **The Barack Obama Corpus**
 -  the_barack_obama_corpus_information.txt
 -  2009
 -  2009
 -  2010
 -  2010
 -  2011
 -  2011
 -  2012
 -  2012
 -  2013
 -  2013
 -  2014
 -  2014
 -  2015
 -  2015
 -  2016
 -  2016
 -  USE
 -  VOKART

 METADATA SEARCH
 CONTENT SEARCH
 MANAGE ACCESS

 REQUEST ACCESS
 CITATION

Corpus

Name The Barack Obama Corpus
Title The Barack Obama Corpus

Description

the_barack_obama_corpus_information.txt

Description

The Barack Obama Corpus (BOC) consists of 6,215,948 words (tokens), which are sourced from nearly 3,500 different texts, dating from January 2009 to January 2016. The texts, all taken from the White House Archives, comprise all speeches held by Barack Obama in his official capacity as 44th President of the United States of America. The earliest speech in the BOC is President Obama's inauguration speech and the last is his final State of the Union speech (January 2016). In total, the corpus includes 34,967 word types, which leads to a type/token-ratio of 0.56.

The files, which display the original titles given to them by the White House, have been tagged for genre, audience type, date and location of delivery, and principal topics. The genres include remarks, addresses, statements, press conferences, debates and question-

Description

How to cite this resource:
 Riesner, Katherina (2017). The Barack Obama Corpus [Data set]. <http://hdl.handle.net/10050/00-0000-0000-0003-C53B-4@view>

Appropriate data citation and PID

Riesner, Katherina (2017). The Barack Obama Corpus [Data set].
<http://hdl.handle.net/10050/00-0000-0000-0003-C53B-4@view>

obama

Showing 5 results for obama  

Results per page:

10



Use the categories below to limit the search results to those matching the selected value(s).

Language



Collection



Modality



The Barack Obama Corpus

(Part of Lund University Humanities Lab)



the_barack_obama_corpus_information.txt; The Barack **Obama** Corpus (BOC) consists of 6,215,948 words (tokens), which are sourced from nearly 3,500 different texts, dating from January 2009 to January 2016. The texts, all taken from the White House Archives, comprise all speeches held by Barack **Obama** in his official capac...

 [Landing page for this record at corpora.humlab.lu.se](#)



metadata

[vlo.clarin.eu](#)

SWE-CLARIN

CLARIN Resource families



Corpora

- Computer-mediated communication corpora
- Corpora of academic texts
- Historical corpora
- L2 learner corpora
- Literary corpora
- Manually annotated corpora
- Multimodal corpora
- Newspaper corpora
- Parallel corpora
- Parliamentary corpora
- Reference corpora
- Spoken corpora

Lexical Resources

- Lexica
- Dictionaries
- Conceptual Resources
- Glossaries
- Wordlists

Tools

- Normalization
- Named entity recognition
- Part-of-speech tagging and lemmatization
- Tools for sentiment analysis

Spoken corpora in the CLARIN infrastructure

Corpora with transcriptions and audio recordings

Corpus	Language	Description	Availability
Arabic Speech Corpus	Arabic	The corpus is available for download from a dedicated webpage. For a relevant publication, see Halabi (2016) .	Download
DIALEKT v1: dialectal corpus with multi-tier transcription	Czech	This corpus contains traditional dialectological material, mostly unprepared monologue-type speech. The corpus is available download (upon request) and through the concordancer KonText. For a related publication, see Komrsková et al. (2018) .	Concordancer Download

CLARIN-IT (2015-onwards)



About Governance Consortium Centres Join Access Events Initiatives News **Home**

the Italian Common Language Resources and Technology Infrastructure

English ▾



ONLINE EVENTS

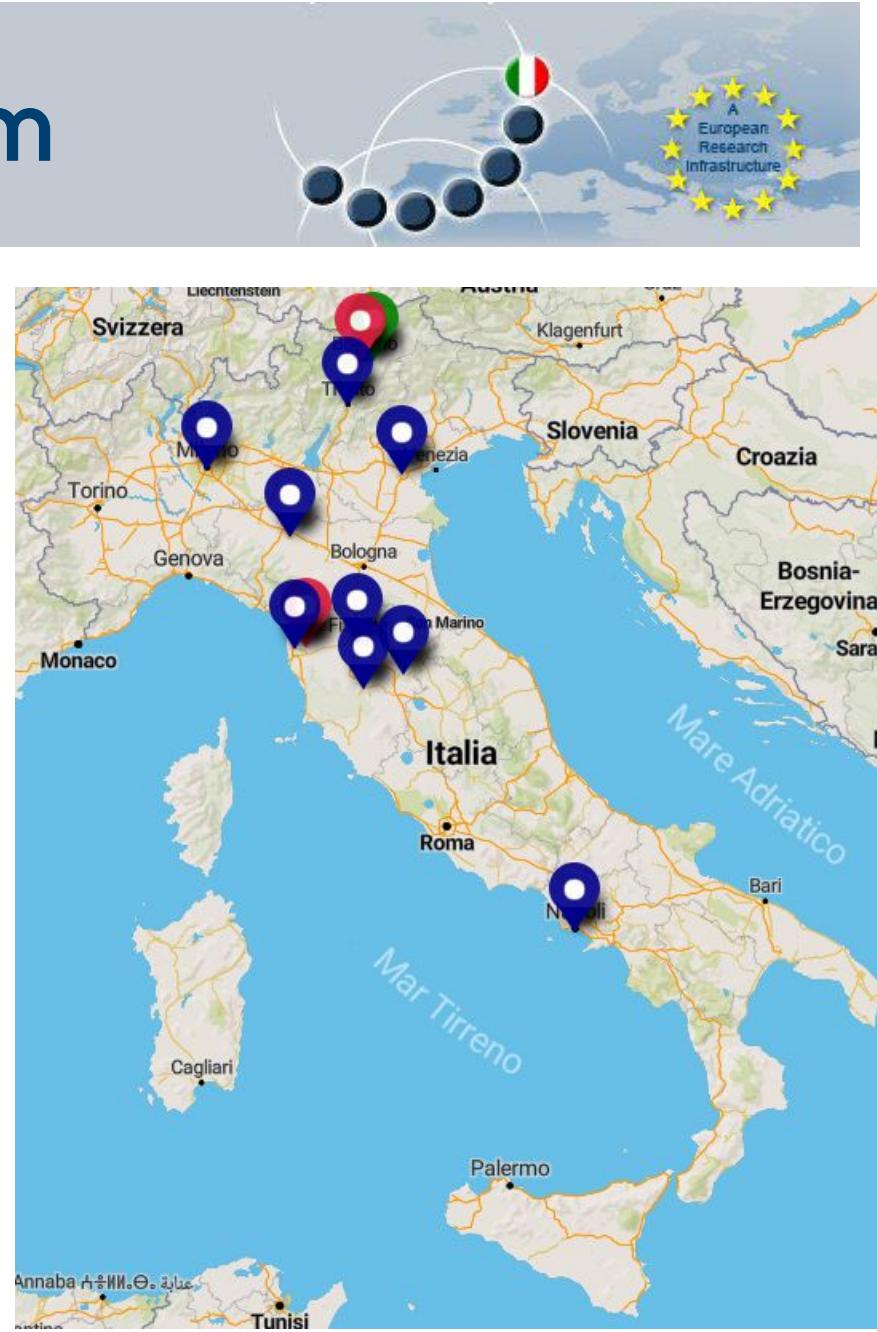
EUPORIA 2021 Webinar - Encoding a Critical Apparatus



07/12/2020 - [Registration](#) to follow the Webinar via Zoom

CLARIN-IT Consortium

- The Department of Education, Human Sciences and Intercultural Communication of the University of Siena (Arezzo)
- The Department of Philology and Literary Criticism of the University of Siena
- The Eurac Research Association (Bolzano)
- The Bruno Kessler Foundation (Trento)
- The Archival and Bibliographical Superintendence of Tuscany (Firenze)
- The Department of Electrical Engineering and Information Technology and the Interdepartmental Research Center "URBAN/ECO" of the University of Naples Federico II
- The Catholic University of the Sacred Heart (Milano)
- The University of Parma (Parma)
- The University of Padova (Padua)



ILC4CLARIN



[Chi siamo](#) [Organizzazione](#) [Repository](#) [Servizi](#) [Eventi CLARIN](#)



The background of this section is a dark grey image of a complex network graph, consisting of numerous teal-colored circular nodes connected by thin teal lines, representing a large dataset or research community.

ILC4CLARIN
REPOSITORY

Easy to be found | Easy to be cited

ILC4CLARIN: data types



CLARIN-IT Repository About CLARIN

ItalWordNet v.2

Please use the following text to cite this item or export to a predefined format:

Roventini, Adriana; Marinelli, hosted at Institute for Computer Science
<http://hdl.handle.net/11336/1000>

BIBTEX CMDI

ALIM Share Home Authors Item identifier Project URL Demo URL Date issued Type Size Language Description

ITA EN

ARCHIVIOVi.Vo. Conservazione e diffusione degli archivi orali e audiovisivi

REGIONE TOSCANA Archivioorale Università di Siena CASENTINO

3 youaignoz eccon

ILC4CLARIN: deposit



Deposit Free and Safe

License of your Choice (Open licenses encouraged)

Easy to Find

Easy to Cite



🔍

Search

Advanced Search

Author	Subject	Language (L)
Anonymous (79)	Latin (413)	Latin (41)
Nahli, Ouafae (31)	Middle Ages (405)	English (1)
Khalfi, Mustapha (30)	Prosa (318)	Arabic (3)
Zarghili, Arsalane (30)	Fonti Letterarie (310)	Italian (24)
Multiple Authors (11)	Storiografia (82)	Ancient Greek (to 1453) (10)
... View More	... View More	... View More

Exercise 1 - Explore ILC4CLARIN

- Find the ILC4CLARIN repository
 - Which is the most represented language in terms of records?
 - What kind of data can you find in Arabic?
 - How do you cite CophiWordNet?
- Try to log in with your institutional identifier, using the Login function (top right)

UDPipe

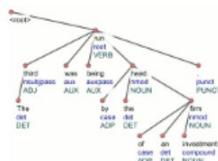
Please use the following text to cite this item or export to a predefined format:

[BIBTEX](#)[CMDI](#)

Straka, Milan and Straková, Jana, 2016, *UDPipe*, LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University, <http://hdl.handle.net/11234/1-1702>.



Share:  



Authors:

Milan Straka, Jana Straková

Description:

UDPipe is an trainable pipeline for tokenization, tagging, lemmatization and dependency parsing of CoNLL-U files. UDPipe is language-agnostic and can be trained given only annotated data in ConLL-U format. Trained models are provided for nearly all UD treebanks. UDPipe is available as a binary, as a library for C++, Python, Perl, Java, C#, and as a web service. UDPipe is a free software under Mozilla Public License 2.0 and the linguistic models are free for non-commercial use and distributed under CC BY-NC-SA license, although for some models the original data used to create the model may impose additional licensing conditions.

[Project home](#)[Run](#)

Tool Inventory

▼ Constituency Parsing



> WebLicht Const Parsing DE



> WebLicht Const Parsing EN

▼ Coreference Resolution



> Concraft -> Bartek

▼ Dependency Parsing



> Concraft -> DependencyParser



> MaltParser



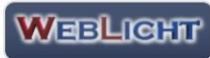
> Spacy (hosted by D4Science) - DE



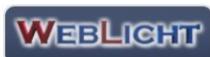
> Spacy (hosted by D4Science) - EN



> UDPipe



> WebLicht Dep Parsing DE



> WebLicht Dep Parsing EN



Exercise 2

- Go to clarin.eu
- Find the VLO (tip: it is a service that allows you to search for language resources...)
 - Take a look at the examples of queries
- Search for texts by Robert Louis Stevenson
 - What can you find?
 - Which format is the corpus in?
 - Where are they hosted?
- On the Links tab, use the three dots (...) to activate the Switchboard (see image)
 - Explore with Voyant or
 - Process with UDpipe
 - Process with Weblicht (after logging in with Single Sign On)

CLARIN for knowledge sharing

Knowledge Centres



About ▾ Language Resources ▾ Learn & Exchange ▾ Events News Contact

 CLARIN and Ukraine

TA

Home / Learn & Exchange / Knowledge Centres

Knowledge Centres

CLARIN Knowledge Infrastructure

CLARIN Knowledge Centres (K-centres) are a cornerstone of the CLARIN Knowledge Infrastructure (KI), one of the main components ensuring a continuous transfer of knowledge between all players involved in the construction, operation and use of the infrastructure. The mission of the CLARIN KI is to ensure that the available knowledge and expertise does not exist as a fragmented collection of unconnected bits and pieces, but is made accessible in an organised way to both the CLARIN community and the social sciences and humanities research community more widely.

The Role of K-Centres

The focus of CLARIN is on language resources (in all modalities, from all regions and with any topical orientation) and K-centres serve researchers and educators from any discipline where language plays one of its many roles, ranging from object of study, a means of communication or expression, a means to store and extract information, object of learning or teaching activities, to training source for data-driven analytics, and many others. K-centres share their knowledge and expertise on one or more aspects of the domain covered by the CLARIN infrastructure and can be mostly found in CLARIN countries, but also exist elsewhere, and they all have a virtual presence.

Areas of K-Centre Expertise

K-centres all have their own specific areas of expertise, which can belong to many different categories, such as

<https://www.clarin.eu/content/knowledge-centres>

The Knowledge
Centres of CLARIN
can be contacted
through their
Help Desks

Knowledge Centres

List of all 22 CLARIN K-centres with expertise in specific linguistic topics

Click on the full name of the K-centre to go to its landing page, and click on the acronym to see its full organisation details

ACE

CLARIN Knowledge Centre for Atypical Communication Expertise

Areas of competence	Atypical communication encompasses language and speech as encountered during (second) language acquisition and development, and in language disorders, but also more broadly in bilingual language development and in sign language. ACE is specialised in this type of research and concomitant infrastructural issues related to data acquisition, processing and sharing, which is typically highly characterised by sensitivity issues. For data storage and access the centre collaborates with MPI's TLA (The Language Archive) which is a CLARIN B Centre and also based in Nijmegen.
Audiences served	- linguists; - psychologists; - neuroscientists; - computer scientists; - speech and language therapists; - education specialists
Types of services	- how-to documents; - access to document templates; - Access to data; - Depositing; - FAQ; - Helpdesk; - Technical support
Is portal for language(s)	-
Other languages covered	-
Modalities covered	- Audio: speech; - Text; - Video: sign language
Linguistic topics	- Language acquisition (L1 and L2); - language disorders; - Language learning
Language processing	-
Data types	-
Resource families	- Spoken corpora; - Manually annotated corpora; - Multimodal corpora
Generic topics	- Critical Data Management; - Legal and ethical issues
Other keywords	- Language acquisition; - sign language; - language pathologies
Tour de CLARIN	Introduction Interview

K-Centres & DH: DiPText-KC



DiPText-KC

CLARIN Knowledge Centre for Digital and Public Textual Scholarship

DiPText-KC offers expertise on methods, data, instruments and technologies relevant in the field of Philological and Literary Studies, History, Art History and Cultural Heritage.

Its actions aim at:

- sharing information with scholars and students about the state of the art in digital scholarly editing and text annotation through domain-specific languages;
- supporting scholars and students in the creation and publication of digital scholarly editions and resources;
- organizing training activities (for instance webinars, workshops and summer schools).

DiPText-KC is one of the Centres of [CLARIN-IT](#), the Italian node of [CLARIN](#) (Common Language Resources and Technology Infrastructure), a digital infrastructure of pan-European interest identified by [ESFRI](#) (European Strategy Forum on Research Infrastructures) and classified as a Landmark Research Infrastructure for the Social Sciences and Humanities (ESFRI Landmarks SSH RI).



HIGHLIGHTS



The Digital and Public Textual Scholarship Knowledge Centre is focused on **digital philology**

<https://diptext-kc.clarin-it.it>

Activities of the DiPText-KC



Consortia, Associations, Centers

- Consortia
 - [Unicode Consortium](#)
 - [TEI Consortium](#)
- National and International Associations
 - [ADHO](#)
 - [EADH](#)
 - [AIUCD](#)
- DH Centres and Labs
 - Italy
 - [CIRCSE](#)
 - [LabCD](#)

Training

- Summer Schools
 - Venice Centre of Digital and Public Humanities, Department of Humanities, Ca' Foscari University of Venice
 - [Summer Camp 2020](#)
 - University of Pisa
 - [Digital Tools for Humanists 2022](#)
 - (past editions: [2021](#) | [2020](#) | [2019](#) | [2018](#) | [2017](#))
 - [Digitising, Cataloguing, Searching and Sharing the Medieval and Early-Modern Image](#)
- Seminars / Webinars
 - [VeDPH Seminars in Digital and Public Humanities – January-May 2022](#)
 - (past editions: [October 2019 – May 2020](#) | [September 2020 – December 2021](#))
 - [Humanities Horizons – History, Hacktivism and Genetic Criticism, Solstice Seminar in DPH](#)
 - [The Public Staging of Gender in Shakespearean Theatre Discussion with Pamela Allen Brown](#)

Digital Libraries

- Zotero Collections
 - [DiPText-KC Library](#)
 - [CLARIN Library](#)

CNR-ILC CoPhiLab @ GARR 2022
The Collaborative and Cooperative Philology Lab (CoPhiLab, CNR-ILC): data, applications, services and infrastructures (5:50:11-5:59:00)

BREAKING NEWS

- Fourth Appointment of the Workshop Cycle "Digital Philology meets Computational Linguistics"
- Third Appointment of the Workshop Cycle "Digital Philology meets Computational Linguistics"
- Concluded the First Cycle of the Permanent Seminar Series "A bridge between two worlds"
- CNR-ILC CoPhiLab @ GARR 2022
- Second Appointment of the Workshop Cycle "Digital Philology meets Computational Linguistics"

The Digital and Public Textual Scholarship Knowledge Centre keeps you informed on Consortia, Associations, Centres, Training Schools, and Digital Libraries relevant for digital philologists

<https://diptext-kc.clarin-it.it>

Learning Hub

Here you can access training materials, share best practices in teaching, and find out more about training events. If you would like to stay informed about new training materials or share your teaching experience, please subscribe to our dedicated mailing list for trainers.



Training Materials

Browse CLARIN's wide range of open-access teaching and learning resources designed by teachers, lecturers and trainers in our network.

[Explore →](#)

Training Events

Join our free workshops and training events to learn how to use the CLARIN services and natural language processing tools for language and linguistic research.

[Explore →](#)

Teaching Award

The Teaching with CLARIN Award is for teachers and lecturers who successfully integrate CLARIN resources into the classroom and university curricula.

[Explore →](#)

2022 Winners



What's on the Agenda? Topic Modelling Parliamentary Debates before and during the COVID-19 Pandemic

Ajda Pretnar Žagar, Kristina Pahor de Maiti & Darja Fišer
University of Ljubljana, Slovenia



Lithuanian Collocations: Usage, Teaching, Learning, and Translation

Jurgita Vaičenonienė
Vytautas Magnus University, Lithuania



Natural Language Processing Methods

Rachele Sprugnoli
Dipartimento di Discipline Umanistiche, Sociali e delle Imprese Culturali – Università degli Studi di Parma, Italy

Natural Language Processing Methods

Goals and Objectives

Natural Language Processing (*NLP*) is an interdisciplinary field whose goal is to create machines that understand natural languages.

This training material features: (i) an introduction to the key concepts and techniques of NLP; (ii) hands-on activities on some NLP tasks, such as lemmatization, part-of-speech tagging and named entity recognition using CLARIN-ERIC tools. The use case covered by the hands-on activities is particularly suited for trainees in the field of Digital Humanities given that the text is taken from a corpus of historical travel writings.

Learning Outcomes

After the first part of the lesson, trainees will acquire a basic knowledge of the key concepts and methods of NLP (e.g., pipelines, linguistic tasks, text annotation, machine learning, and evaluation metrics).

After the second part, trainees will be able to 1) automatically analyze texts with NLP tools; 2) choose the NLP tool that best fits their specific purposes among those provided by CLARIN-ERIC; 3) create a geographical visualization starting from a text processed with a Named Entity Recognition tool.

Author(s)

Rachele Sprugnoli

Researcher

Dipartimento di Discipline Umanistiche, Sociali e delle Imprese Culturali – Università degli Studi di Parma, Via D'Azeglio 85, 43125 Parma, Italy



Co-funded by the
Erasmus+ Programme
of the European Union



Welcome to the UPSKILLS block dedicated to
Introduction to Language Data: Standards and Repositories

Workload: 5 ECTS

Designers: Iulianna van der Lek and Darja Fišer with contributions from Francesca Frontini, Alexander König and Willem Elbers

Designers' affiliations: CLARIN ERIC

Prerequisites:

- Text Processing
- First steps into scientific research

Learning outcomes:

By the end of this unit block, students will be able to:

- Explain the main concepts related to research data repositories for language resources and technologies and the role they play in the linguistic research data lifecycle in the context of Open Science and FAIR
- Find and use certified research data repositories to discover, share, publish, and archive language and linguistic resources and datasets
- Find and use integrated repository services and tools to process, annotate, and analyse different types of corpora according to standards and formats used by the community
- Identify potential legal and ethical issues when collecting, sharing and reusing language resources

The materials in this unit block are primarily intended as learning content for lecturers in language-related courses, who are invited to view the materials, export them for use via the Moodle installation of their own institution, and adapt them as they see fit for their purposes.

The learning content in this unit block is general and can be reused in any programme or course related to *Computational Linguistics, Language Technologies, and/or Digital Humanities*.

Before proceeding, please consult the guide attached below, which explains how the materials are organised and how they can be used.

* Illustrations on tiles by Storyset

This work is licensed under a [Creative Commons Attribution-ShareAlike 4.0 International License](#).



Questions? Interested in
piloting the content?

Contact:
training@clarin.eu



1. Introduction to the Language Resource Management Lifecycle



2. How Research Data Repositories Make Language Data FAIR



3. Finding and (Re)using Language Resources in CLARIN Repositories



4. Citing Language and Linguistic Data



5. Archiving and Publishing Language Resources



6. Legal Issues in Language Resources Management (draft)



7. Example of a Student Project



8. Unit Glossary

CLARIN Internships – Meet the Interns

Submitted by Julia Misersky on 22 January 2023

CLARIN has launched a remote internship programme for students interested in learning how to use the research infrastructure to access and engage with digital language data with the help of advanced technologies. The internship aims to help the students enhance their language data handling and processing skills, and show them how to manage language resources according to the FAIR Data Principles.

Our first two interns are Elton Pistolia and Lesley Messori, who are second-year Translation Technology MA students at the University of Bologna.



Elton Pistolia

Elton has a BA in Foreign Languages for Tourism and International Mediation and is now in his second year of an MA in Translation Technology at the University of Bologna. He has a strong background in linguistics and translation and is keen to develop his programming and machine-learning skills. He speaks Italian, English, Greek, Albanian, German, Russian and Portuguese. Elton was selected as the best participant for the language combination EN-EL in the eMT Challenge

2022. After his MA, Elton plans to do a PhD in machine translation.



Lesley Messori

Lesley has a BA in Linguistic and Intercultural Mediation and is now in her second year of the MA in Translation Technology programme at the University of Bologna. During her MA, she acquired knowledge of corpus linguistics, terminology management, translation technology tools and learnt how to adapt machine translation systems. Lesley speaks five languages, Italian (mother tongue), English, Russian, Danish and Slovak and has experience in audiovisual translation, translation and post-editing.

In the future, she sees herself working as a freelance translator because she simply loves anything related to languages.

Project: Improving the discoverability of the language resources in the Virtual Language Observatory

- GitHub Notebook to parse the XML files from the VLO, detect languages and NER with spaCy
- Create a gold standard to evaluate the VLO metadata set
- Annotation of named entities
- Translation

Project: CLARIN Vocabulary to tag web content and academic outputs

- Corpus & extraction of terms candidates using SketchEngine
- Term validation guidelines & 3 external validators
- Selection of final terms, concept mapping, definitions and translations
- SKOS and TBX formats
- Sharing, depositing and archiving in CLARIN-IT



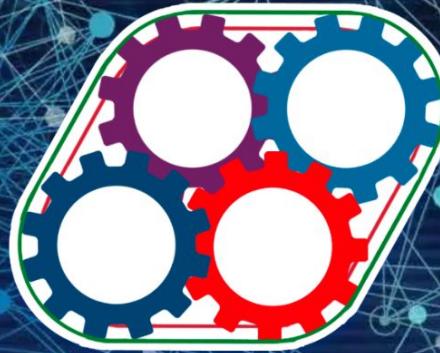
H2IOSC

Humanities and Heritage Italian Open Science Cloud

[Home](#)[News](#)[The Project](#)[Partner Institutions](#)[Targeted RIs](#)[Colophon](#)

H2IOSC

H2IOSC 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.



<https://www.h2iosc.cnr.it/>

CLARIN for researchers



- Contact CLARIN for help with your research
- Visit this page
 - <https://www.clarin.eu/content/clarin-researchers>
- Participate in CLARIN (virtual) events
 - <https://www.clarin.eu/events>
- Tour de CLARIN
 - <https://www.clarin.eu/Tour-de-CLARIN>
- Use CLARIN Training and videolectures

Contacts



- CLARIN ERIC
 - Newsletter <https://www.clarin.eu/news>
 - @CLARINERIC

- CLARIN IT
 - www.clarin-it.it
 - @CLARIN_IT
 - coordination@clarin-it.it