



CLARIN services for FAIR language resources

Franciska de Jong
executive director CLARIN ERIC

f.m.g.dejong@uu.nl

AIUCD2021 - DHs for society, 21 January 2021

From the call text of



AIUCD 2021

Digital humanists, thanks to the knowledge and expertise acquired in recent years and experimented in numerous research activities, have been among the often silent protagonists of the change of society, pervaded by information technologies in continuous development.



Part I

DHs for society and the role of FAIR data

Part II

CLARIN's service infrastructure – a bird's eye view

Part III

Examples of CLARIN service offer for DH tracks with potential for societal impact

Part IV

Concluding remarks

DH for Society

AIUCD

- Digital Humanities as a silent protagonist of societal change in the era of IT

Open Science

- Research as potential protagonist of innovation
- Investment in research with public funding should benefit society at large
- No monopoly for commercial parties

Digital turn -> Open Science +

- Adherence to principles of responsible data science

Aspects of Open Science

Open Access

Open Source

Open Data

Aspects of Open Science

Open Access

- Focus on reuse of academic publications
- Berlin Declaration on Open Access to Knowledge in the Sciences and Humanities (2003)
- <https://openaccess.mpg.de/Berlin-Declaration>



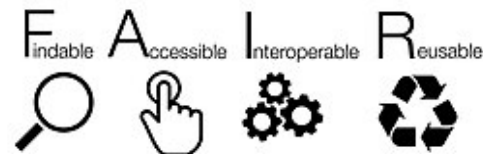
Open Source

- Focus on reuse of software
- OSS declaration (1998)



Open Data

- Focus on reuse of data resulting from research
- More recently (2016) promoted as FAIRness of data:
Findable, **A**ccessible, **I**nteroperable, **R**eusable
- <https://www.force11.org/group/fairgroup/fairprinciples>



Open Science is calling for an Open Mind

- Open Science is calling for open access to data. However, in many domains
(underlining that not all data can be open, and that Open Science is not to be imposed as a unified framework for all).
- Data and algorithms are surrounded by an aura of accuracy, objectivity, and truth. However,
(calling for attention for values such as democratic society, fundamental rights and the rule of law, through examination of epistemic and ethical impact and regulations that promote and enable good data practices)

Open Science is calling for an Open Mind

- Open Science is calling for open access to data. However, in many domains (underlining that not all data can be open, and that Open Science is not to be imposed as a unified framework for all).

Data open if it can, closed if required.

- Data and algorithms are surrounded by an aura of accuracy, objectivity, and truth. However, (calling for attention for values such as democratic society, fundamental rights and the rule of law, through examination of epistemic and ethical impact and regulations that promote and enable good data practices)

Data-driven research is to follow the principles of Responsible Data Science, taking into account context and diversity

Europe's multilinguality is key*



- to our understanding of how language affects identity, culture, society
- to our understanding of diversity across boundaries of time and regions
- and therefore for comparative studies

* image from https://www.coe.int/t/dg4/linguistic/jel_en.asp

Language data as pillar for data science for society

- Europe's **multilinguality** is a basis for **comparative research** of societal and cultural phenomena, that are reflected in language use
- Some examples:
 - Migration patterns
 - Intellectual history
 - Language variation across period and region
 - Dynamics in mental health conditions / speech pathology
 - Parliamentary discourse

Data reuse and sharing: disciplinary dynamics

History

- from fear for subjectivity to interest in multiperspectivity

Linguistics

- shift in attention from language proper to language in context
 - Identity shaping
 - Language variation
 - Language as vehicle for memories
 - Language as facilitating a dialogue
 - Language as part of multimodal expression

Social sciences

- Integration of surveys and interviews

CLARIN and Open Science

- Promoting the sharing and re-use of data through sustainable data registries
- All integrated datasets available in open access for research purposes
- Support for linguistic diversity
 - Data covering more than 1500 languages
 - Tools for many languages
 - Language resources in all modalities
- Adherence to FAIR data principles
 - Interoperability through a common metadata framework
- Promotion of responsible data science
- Strengthening the support for 500.000 professional SSH researchers

[CLARIN: Towards FAIR and Responsible Data Science Using Language Resources.](#) " In: *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, May 2018, 3259-3264.

Part II

A bird's eye view on the CLARIN services

CLARIN in eight bullets

- **CLARIN** is the Common Language Resources and Technology Infrastructure
- **ESFRI** ERIC status since 2012, Landmark since 2016
- that 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
- that serves as an ecosystem for **knowledge exchange**
- and: with some services already **integrated in EOSC** (European Open Science Cloud; [link](#))

CLARIN ERIC in members and centres

A consortium of:

- 21 members:

AT, BG, CY, CZ, DE,
DK, EE, FI, GR, HR, HU, IS, IT,
LT, LV, NL, NO, PL, PT, SE, SI

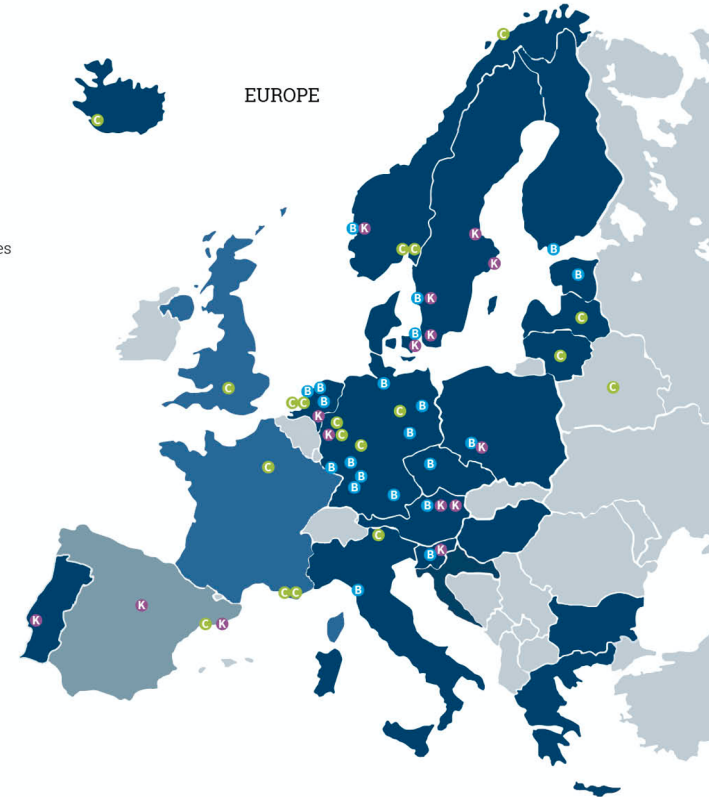
- 3 observers:

FR, UK, ZA

>64 centres, of which
>24 CTS* certified data centres



- ERIC members
- Observers
- Countries with participating centres
- ⓑ Centre Providing Data
- ⓒ Centre Providing Metadata
- Ⓚ Knowledge Centre



*CTS = Core Trust Seal: www.coretrustseal.org

CLARIN in data types

- Newspaper archives
- Literary texts
- Social Media data
- Parliamentary records
- Historical letters
- Oral History data
- Disciplinary libraries
- Institutional archival data
- Broadcast archives
- ...

See also the info on the CLARIN Resource Families initiative: <https://www.clarin.eu/resource-families>

CLARIN's communities of use

- Linguistics and Philology
- Digital Humanities
- Translation and Lexicography
- Literary Studies
- History
- Political and Social Sciences
- Media Studies
- Culture heritage experts
- Speech therapy
- Teachers
- General Public
- ...

CLARIN in central services



CLARIN portal

Get an example-based impression of what's currently available



Depositing services

Store language resources in a sustainable repository at a CLARIN centre



Virtual Language Observatory

Discover language resources using a faceted browser or a map



Easy access to protected resources

Get easy access to protected resources, with your institutional username and password.



Language Resource Switchboard

Explore and analyze language data with a wide variety of tools



Virtual Collections

Create your own digital bookmarks, ideal for citing data sets.



Language Resource Inventory

Submit and access information about language resources relevant to your research.



Content Search (prototype)

Search different corpora with a single search engine

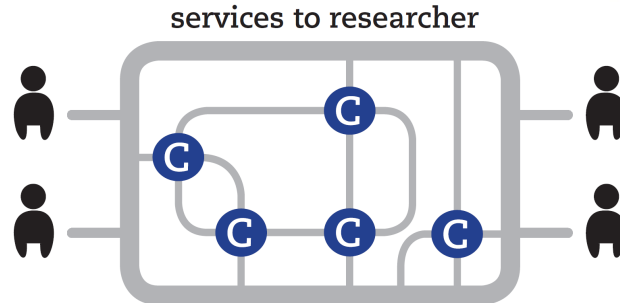


Questions & Answers

Searching for a specific data set or application? Wondering how CLARIN can assist your research? Feel free to contact us!

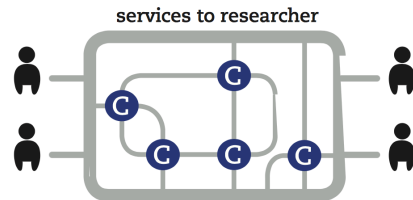
CLARIN: network of centres

- **Distributed architecture:**
(http-accessible) files, web applications and web services spread all over centres in Europe and beyond
- Metadata **harvesting** to enable search and access from one central platform
- Currently:
 - 24 certified **B Centres**
 - over 60 registered centres in total
- Tools and data from different CLARIN centres are **interoperable**, so that data collections can be combined and tools from different sources can be chained to perform complex operations to support researchers in their work.

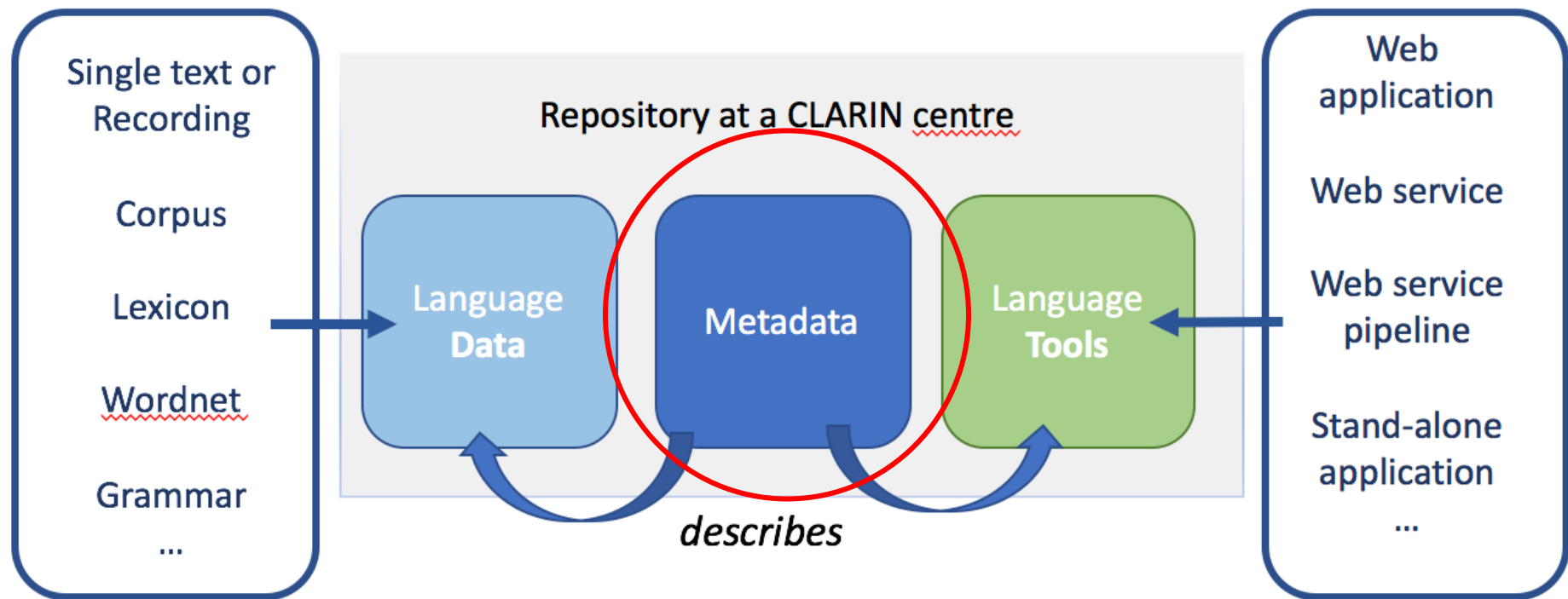


What CLARIN Centres offer

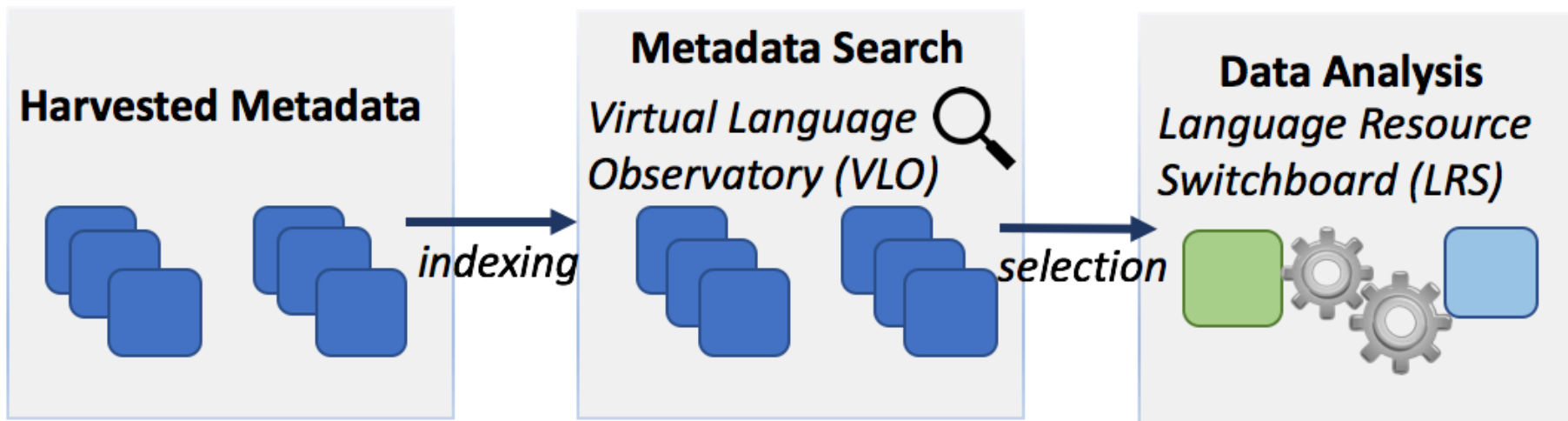
- Repository
 - library of linguistic data and tools
 - search for data and tools and easily use them online or download them
 - deposit your data and be sure it is safely stored, everyone can find it, and correctly cite it
- Federated single sign-on
 - log in once with your existing institutional credentials
 - get access to protected resources
- Metadata
 - describe content, provenance and formats of linguistic data and tools
 - facilitate preservation and dissemination of linguistic data and tools
- Persistent Identifier (PID or handle)
 - a special permanent URL that provides a permanent link to linguistic data and tools
 - will resolve correctly even if in some distant future the data is moved
 - to be used as URL in citations
- Licensing, aligned with diverse conditions
 - Public
 - Academic
 - Restricted
- Preservation
 - committed to long-term care of items in the repository



The CLARIN data architecture: *centre repositories*



The CLARIN data architecture: *central processing of metadata*



Examples of what CLARIN can do for you

- Datasets to be found via central discovery service
- Network of open access repositories
- Interoperability through common metadata framework
- Information on rights for reuse through licence details
- Virtual Language Observatory: search based on facets
- SwitchBoard: matching tool for selecting suitable tools
- Federated Content Search: content based search across the centre network
- CLARIN Resource Families: overviews of comparable datasets in multiple languages
- Virtual Collection Registry: create, cite, share coherent set of digital objects

Examples of what CLARIN can do for you

- Datasets to be **F**ound via central discovery service
- Network of open **A**ccess repositories
- **I**nteroperability through common metadata framework
- Information on rights for **R**euse through licence details
- Virtual Language Observatory: search based on facets
- SwitchBoard: matching tool for selecting suitable tools
- Federated Content Search: content based search across the centre network
- CLARIN Resource Families: overviews of comparable datasets in multiple languages
- Virtual Collection Registry: create, cite, share coherent set of digital objects

VLO facets selected:

Italian
TalkBank

Showing 1 to 10 of 892 results within selection for TalkBank Italian i Results per page: 10

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

Language ⌵

Italian ✕

Collection ⌵

TalkBank ✕

Resource type ⌵

Modality ⌵

Format ⌵

⌵

Availability ⌵

Search options ⌵

<< < 1 2 3 4 5 6 7 8 9 10 > >>

CABank MOVIN Corpus - MOVIN
(Part of TalkBank)

⊞ Political Talk Shows, CA transcripts from 9 languages - Corpus: MOVIN

German Danish Italian English Estonian ... (+4)

[Landing page for this record at *ca.talkbank.org*](#)

Bilingual ProjectS Corpus - ProjectS
(Part of TalkBank)

⊞ child learning English, Italian, and Scottish Gaelic - Corpus: ProjectS

English Italian Gaelic; Scot...

[Landing page for this record at *childes.talkbank.org*](#)

Roggero Italian-Japanese Bilingual Corpus - Roggero
(Part of TalkBank)

⊞ Italian and Japanese in Italy - Corpus: Roggero

Japanese Italian

[Landing page for this record at *childes.talkbank.org*](#)

Frogs Italian Roma Corpus - Italian-Roma
(Part of TalkBank)

⊞ Italian frog story narrations - Corpus: Italian-Roma

Italian

[Landing page for this record at *childes.talkbank.org*](#)

VLO facets selected:

ERCC

Showing 1 to 10 of 11 results within selection for ERCC: Learner Corpora or ERCC: Web Corpora or ERCC: Other ⓘ Results per page: 10 ▾

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

Language ▾

Collection ▾
ERCC: Learner Corpora ✕
OR ERCC: Web Corpora ✕
OR ERCC: Other ✕

Resource type ▾

Format ▾
▾

Availability ▾

Search options ▾

<< < 1 2 > >>

ACTER (Annotated Corpora for Term Extraction Research) v1.3

(Part of ERCC: Other)

▢ The ACTER (Annotated Corpora for Term Extraction Research) is an annotated dataset for term extraction. Terms and Named Entities have been manually annotated in specialised comparable corpora covering 3 languages (English, French, and Dutch), and 4 domains (corruption, dressage, heart failure, and wind energy).

English French Dutch

🏠 Landing page for this record

KrdWrđ CANOLA Corpus 1.0

(Part of ERCC: Web Corpora)

▢ The CANOLA Corpus is a visually annotated English web corpus for training classification engines to remove boiler plate on unseen Web pages. It was harvested, annotated and evaluated by the tools and infrastructure of the KrdWrđ Project.

English

🏠 Landing page for this record

KrdWrđ CANOLA Corpus 1.1

(Part of ERCC: Web Corpora)

▢ The CANOLA Corpus is a visually annotated English web corpus for training classification engines to remove boiler plate on unseen Web pages. It was harvested, annotated and evaluated by the tools and infrastructure of the KrdWrđ Project.

English

🏠 Landing page for this record

PAISÀ Corpus of Italian Web Text

(Part of ERCC: Web Corpora)

▢ The Paisà corpus is a large collection of Italian web texts, licensed under Creative Commons (Attribution-ShareAlike and Attribution-Noncommercial-ShareAlike). It has been created in the context of the project PAISÀ. All documents were selected in two different ways. A part of the corpus was constructed using a meth...

Italian

🏠 Landing page for this record

26

VLO facets selected:

ILC4CLARIN
ALIM

Showing 1 to 10 of 401 results within selection for ILC4CLARIN : ILC Data & Tools or ILC4CLARIN : OPEN Data & Tools or ALIM Literary Sources Results per page: 10

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

Language

Collection

alim

ILC4CLARIN : ILC Data & Tools ✕
OR ILC4CLARIN : OPEN Data & Tools ✕
OR ALIM Literary Sources ✕
[more...](#)

Resource type

Format

Availability

Search options

<< < 1 2 3 4 5 6 7 8 9 10 > >>

De controversia mensium

(Part of ALIM Literary Sources)

Edizione critica a stampa a cura di G. Orlandi esportazione in formato TEI XML livello ALIM2_0 effettuata da Chiara Casali Edizione in formato TEI XML, livello ALIM2_1, a cura di Jan Ctibor HomePage del progetto: <http://it.alim.unisi.it/il-progetto/> Documentazione: <http://alim.unisi.it/documentazione>

Latin

Landing page for this record

Carmina

(Part of ALIM Literary Sources)

Edizione critica a stampa a cura di E. Dümmler esportazione in formato TEI XML livello ALIM2_0 effettuata da Edoardo Ferrarini Edizione in formato TEI XML, livello ALIM2_1, a cura di Jan Ctibor HomePage del progetto: <http://it.alim.unisi.it/il-progetto/> Documentazione: <http://alim.unisi.it/documentazione>

Latin

Landing page for this record

The search results include 3 records with the same title.

Sermo de s. Anastasio confessore (BHL 407b)

(Part of ALIM Literary Sources)

Edizione critica a stampa a cura di E. D'Angelo esportazione in formato TEI XML livello ALIM2_0 effettuata da Paolo Monella Edizione in formato TEI XML, livello ALIM2_1, a cura di Jan Ctibor HomePage del progetto: <http://it.alim.unisi.it/il-progetto/> Documentazione: <http://alim.unisi.it/documentazione>

Latin

Landing page for this record



Image copied from the flash presentation at
AIUCD2021 by Marina Buzzoni on ALIM:
“A tea for two”

Cocoon – collection of stories in multiple language variants

The hare and the jackal

> Voir la notice complète

"The hare and the jackal" 2007. Nepali; Achhame. Tamata, Lal Shobha (speaker); Lecomte-Tilouine, Marie (researcher); Lecomte-Tilouine, Marie (interviewer); Michailovsky, Boyd (researcher); Michailovsky, Boyd (depositor); Agence Nationale de la Recherche ANR-06-CORP-030-01 (sponsor). Editeur(s): Laboratoire de langues et civilisations à tradition orale; Centre d'Etudes Himalayennes.

Humorous story recorded to illustrate a dialect of Nepali spoken in Achham. The narrator is a woman over 50 years old of an artisan caste.

CLARIN Resource Families (CRF)



Currently:

12 corpora families

5 families of lexical resources

4 tool families

Read more :

<https://www.clarin.eu/resource-families>

Computer-Mediated Communication corpora (16)

<p>SoNaR New Media</p> <p>Size: 35 million tokens</p> <p>Annotation: tokenised, PoS-tagged, lemmatised</p>	Dutch	<p>This corpus contains tweets, chats and SMS from 2005 to 2012.</p> <p>The corpus is available for download from the Dutch Language Institute and for searching online through the OpenSONAR environment.</p> <p>For the relevant publication, see Sanders (2012).</p>	<p>Concordancer</p> <p>Download</p>
<p>NTAP English</p> <p>Size: 660,798,199 tokens</p>	English	<p>This corpus contains blog posts that are related to climate change issues across science, politics, and the environment. The vast majority of the posts are from 2005 onwards.</p> <p>The corpus is available for searching online through the Corpuscle concordancer.</p> <p>For the relevant publication, see Salway et al. (2016).</p>	<p>Concordancer</p>
<p>DIDI - The DiDi Corpus of South Tyrolean CMC 1.0.0</p> <p>Size: 600,000 tokens</p> <p>Annotation: tokenised, PoS-tagged, lemmatised</p> <p>Licence: ACA-BY-NC-NORED 1.0</p>	English, German, Italian, Ladino	<p>The corpus consists of Facebook posts gathered from 136 Facebook users from South Tyrol. All texts are anonymised.</p> <p>The corpus is available for download from the EURAC Research CLARIN repository.</p> <p>For the relevant publication, see Frey et al. (2016).</p>	<p>Download</p>
<p>The Mixed Corpus: New Media</p> <p>Size: 25 million tokens</p> <p>Annotation: tokenised</p>	Estonian	<p>This corpus contains chat room messages, forum posts and news comments from 2000 to 2008.</p> <p>The corpus is available for download from a dedicated webpage associated with CLARIN Estonia and through a dedicated concordancer.</p>	<p>Concordancer</p> <p>Download</p>

Virtual Collection Registry

- <https://collections.clarin.eu/public?0>

Virtual Collection Registry

Browse

Create

Help

Login

CLARIN

Name	Type	Created
CLARIN services in the European Open Science Cloud	extensional	2020-11-13
SSHOC Webinar: CLARIN Hands-on Tutorial on Transcribing Interview Data	extensional	2020-10-14
Extremadura Buenas Noches	extensional	2020-10-02
PubMed references to articles on "face mask" and "influenza" as of May 2020	extensional	2020-07-17
Pirineos La Nuit	extensional	2020-07-16
DuFLOR - Dubbed Films	extensional	2018-06-13
VLO search results: BAS	extensional	2017-08-29
Treebanks for the RDA Collaboration Project	extensional	2016-06-24
Exploring genealogical blends: the Surinamese Creole Cluster and the Virgin Island Dutch Creole Cluster	extensional	2015-11-09
Henrik Ibsen: works and analyses	extensional	2014-10-20
Absolute spatial deixis and proto-toponyms in Kata Kolok	extensional	2014-09-26
The Trobriand Islanders' Ways of Speaking	extensional	2014-09-22

About
v1.5.1

Service provided by CLARIN
Hosted by Leibniz-Institut für Deutsche Sprache

Contact

CLARIN and service integration into EOSC portal

- The European Open Science Cloud is emerging (www.eosc.eu)
- An EOSC portal is up and running already.
- More and more discipline-specific services are being integrated (registered)
- At the launch of the EOSC portal the integration of the VLO and the Language Resource Switchboard in EOSC was demonstrated: <https://www.clarin.eu/eosc>
- 135,000 cultural heritage objects from Europeana (news papers) visible through VLO and therefore also in EOSC: <https://www.clarin.eu/tags/europeana>

Part III

- CLARIN service offer for DH work with potential for societal impact

CLARIN Resource Family: Parliamentary corpora

- Characteristics
 - essentially transcriptions of spoken language produced in highly controlled and regulated settings
 - rich in invaluable (sociodemographic) metadata
 - rich in links to other research data (e.g., legislation, data from news channels)
- Requirements
 - easily findable and accessible
 - encoded according to international standards / recommendations
 - equipped with rich and reliable annotations and metadata
- Scenarios of scholarly and public use and benefits
 - fuels multidisciplinary research, diachronic and transnational comparative analyses
 - enables more informed participation in public debate
 - supports effective functioning of democratic systems
 - facilitates focused comparative studies on public debates on global crises

Indication of diversity of corpora integrated in CLARIN

- UK's Hansard Corpus is the largest corpus (1.6 billion tokens, 7.6 million speeches made by around 40,000 different speakers) and spanning the longest time period (1803-2005)
- Corpora from other countries are significantly smaller (10 - 100 million tokens) and cover shorter periods (mostly from the 1970s onwards)
- 25 European national parliaments
- The corpus of the European parliament proceedings comes in 21 languages

Parliamentary corpora as shown in CLARIN

[Hansard corpus](#)

English

Size: 1.6 billion tokens

Annotation: tokenised, PoS-tagged, lemmatised, semantic tagging

The corpus contains British parliamentary debates from 1803 to 2005. It is semantically tagged with the [USAS semantic tagger](#) and the [Historical Thesaurus Semantic Tagger](#) (HTST).

[Concordancer](#)

The corpus is available through a dedicated concordancer.

For the relevant publication, see [Rayson et al. \(2015\)](#).

HANSARD CORPUS (BRITISH PARLIAMENT)

7.6 MILLION SPEECHES, 1.6 BILLION WORDS, 1803-2005

DISPLAY

☒ LIST ☐ CHART ☐ KWIC ☐ COMPARE

SEARCH STRING

WORD(S)

COLLOCATES

POS LIST

SEMANTIC CATEGORIES | WORDS

☒ SHOW ☐ DECADE ☐ SPEAKER

1

--IGNORE--
2000
1990
1980
1970
1960
1950

2

--IGNORE--
2000
1990
1980
1970
1960
1950

SORTING AND LIMITS

SORTING

MINIMUM

CLICK TO SEE OPTIONS

1800	1810	1820	1830	1840	1850	1860	1870	1880	1890	1900	1910	1920	1930	1940	1950	1960	1970	1980	1990	2000

INTRODUCTION [Help / information / contact](#)

OVERVIEW OF SEARCHES

This website allows you to quickly and easily search more than 1.6 billion words of text, in more than 7,600,000 speeches from nearly 40,000 individual speakers in the British Parliament from 1803-2005. You can search for **words**, **phrases**, **collocates** (nearby words), and even **grammatical constructions**. You can also search the corpus "semantically", when you want to find all words with a particular meaning.

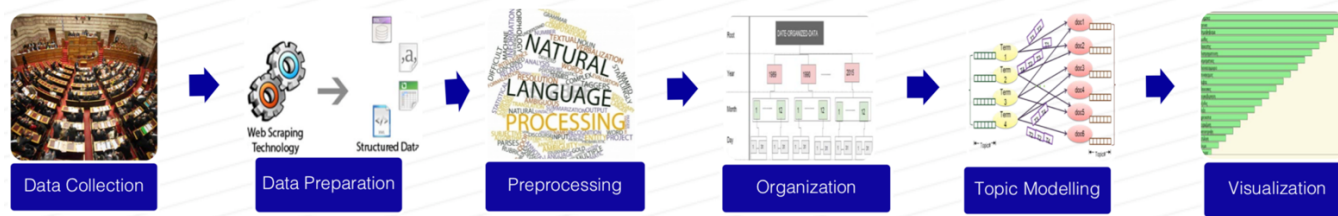
For any of these searches, you can see the frequency over time (1803-2005), either by decade or by year. You can also compare one set of decades to another. You can quickly and easily create "virtual corpora". These might include just the speeches of a certain member of parliament (or a group of MPs), perhaps during just a specific time period.

The corpus was created as part of the **SAMUELS** project (Semantic Annotation and Mark-Up for Enhancing Lexical Searches), 2014-2016. The corpus architecture and web interface was created by **Mark Davies**, and it is related to **other large corpora** of English.

The 5-10 minute **introduction** provides a very good overview of the main different types of searches that one can do. We invite you to start with that web page, if you haven't done so already.

Encoding and analysis of parliamentary data

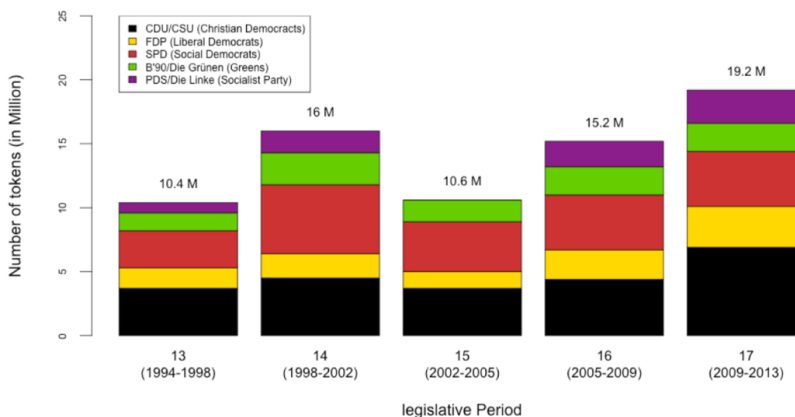
- format, quality and structure of the source files are essential
- traditionally, only transcriptions of parliamentary sessions have been made available; now increasingly released in audio and video as well
- key issues in the encoding of parliamentary data
 1. Source files
 2. Structural elements
 3. Encoding standards
 4. Metadata
 5. Linguistic annotations
 6. Text enrichment
- example analysis workflow



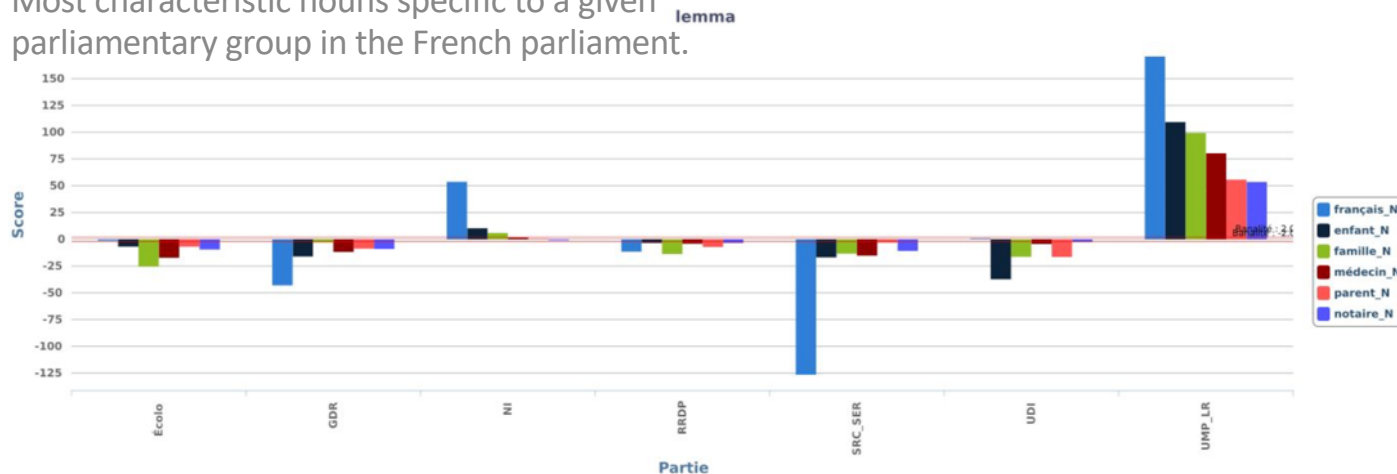
Metadata already tell a story

- Parliamentary records are rich in extralinguistic markup/metadata
 - on the speech level
 - parliamentary session
 - date
 - meeting item
 - speaking time
 - on the speaker level
 - speaker's name
 - date of birth
 - gender, education
 - party affiliation)
- that can be used as variables in analysis or for fine-grained filtering of the research data

Text production in the German parliament over time by political party



Most characteristic nouns specific to a given parliamentary group in the French parliament.



Representation and text production of female members of the Danish parliament

Age	% Female MPs	% Female speeches
20-29	56.1	52.2
30-39	47.9	47.0
40-49	42.7	35.4
50-59	29.8	22.0
60-69	24.7	15.0
70-79	33.3	30.8
Total	38.7	34.4

Case study: Gender in the Danish parliament

- Hansen et al. (2018) investigated gender differences in the revised transcripts of speeches from the Danish Parliament 2009-2017
 - active participation of women in politics is historically relatively new and women are still underrepresented in the Danish parliament
 - analysis:
 - number of the speakers, their age, party and role in the party
 - speech frequencies & speech lengths
 - topics addressed
 - results:
 - general
 - there were relatively more male spokespersons than female ones in the period covered by the corpus but the number of female MPs under 29 is larger than the number of male MPs from the same age group
 - in general, women speak less frequently and for a shorter time than male MPs in proportion to their seats in Parliament, the difference in speaking time between female and male MPs is statistically significant
 - role
 - women belonging to left-wing parties speak less frequently than women from right-wing parties compared to their seats in Parliament
 - female ministers and spokespersons speak more frequently than ordinary MPs
 - female ministers under a male prime minister give fewer speeches than female ministers under a female prime minister even though their percentages in the two periods are similar
 - topics
 - female MPs more often spoke about “softer” political areas, while in the speeches of male MPs “harder” subjects prevailed

Links CLARIN activities and publications for parliamentary data

- Workshop 2017: <https://www.clarin.eu/event/2017/clarin-plus-workshop-working-parliamentary-records>
- ParlaCLARIN workshop I @LREC 2018: <https://www.clarin.eu/ParlaCLARIN>
- ParlaFormat: <https://www.clarin.eu/event/2019/parlaformat-workshop>
- ParlaCLARIN workshop II @LREC 2020: <https://www.clarin.eu/ParlaCLARIN-II>
- ParlaMint project (2020-2021): <https://www.clarin.eu/content/parlamint-towards-comparable-parliamentary-corpora>
- CLARIN Café on ParlaMint: <https://www.clarin.eu/event/2020/clarin-cafe-join-our-parliamentary-flavoured-coffee-parlamint>

CLARIN Resource Family: CMC corpora from Italy

- [**DIDI - The DiDi Corpus of South Tyrolean CMC 1.0.0**](#) (Facebook comments in English, German, Italian, and the minority language “[Ladino](#)”)
- [**PAISÀ Corpus**](#) (Italian web text)

Aim DIDI: to provide scientific answers to current issues of language and education policy as well as to economic and social questions at both the local and international level” for South Tyrol, where “multilingualism involves geographic, institutional, social and personal aspects”.

CLARIN Resource Family: Learner corpora from Italy

- [Italian PAROLE corpus](#), in TalkBank

Aim: the PAROLE corpus presents storytelling-type tasks involving 23 learners of Italian (all French-L1; average age 19), and the transcriptions are marked up for prosodic phenomena relevant for L2-language learning (“unbroken sequences of hesitation phenomena, such as silent pauses, filled pauses, and certain paralinguistic noises”)

- [MERLIN Written Learner Corpus](#) for Czech, German, Italian

Aim: support language learning for speakers of Italian as an L1 language

CLARIN as Knowledge Infrastructure 1

Promotion of data harmonization and methodology development

- Activities related to CLARIN Resource Families
- TwinTalk workshops at DH events
- CLARIN Hackaton on the detection of misinformation related to Covid-19:

Coordination and support of collaborative work with potential for societal impact

- Workshops and webinars aimed at the development of a support chain for creating and processing oral history data.
Impact: providing instruments for minorities to voice their perspective on society
- Workshop series on how to share and deposit speech with communication disorders (aka DELAD).
Impact: speech pathology, early signalling of Alzheimer disease

CLARIN as Knowledge Infrastructure 2:



- Aims
 - increase the visibility of the national consortia and centres
 - reveal the richness of the CLARIN landscape
 - display the full range of activities throughout the network
- Highlights per chapter
 - consortium / centre
 - resource
 - tool
 - UI event
 - researcher
- Links:
 - <https://www.clarin.eu/Tour-de-CLARIN>
 - <https://www.clarin.eu/sites/default/files/tour-de-clarin-italy.pdf> (vol II)



Part IV

Concluding remarks

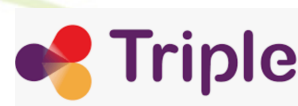
Open Science and multidisciplinary

Language as a relevant resource for all SSH disciplines

- *Social survey data*: interviews, questionnaires
- *Social media data*: content + metadata
- *Psychology data*: oral history, interviews, ego documents
- *Humanities data*: archives, libraries, web content
- *Digital cultural heritage*: free text metadata, content related to artefacts, oral reports documenting excavations, etc.

European research infrastructures rooted in SSH disciplines: CESSDA, CLARIN, DARIAH, ESS, SHARE, ERIHS;
all involved in collaborative projects (H2020 funding) towards better service alignment

- SSHOC (=SSH Open Cloud): www.sshopencloud.eu
SSH cluster activities aiming at integration of practices and services into the European Open Science Cloud
- TRIPLE: www.triple.eu
innovative platform and practices for multidisciplinary exploration of SSH artefacts



CLARIN TwinTalk workshops at DH events

Dynamics in Research Assessment: DORA

San Francisco Declaration on Research Assessment (DORA)

- Promotes recognition of the need to improve the ways in which researchers and the outputs of scholarly research are evaluated
- <https://sfdora.org/>
- >2000 organisations signed
- >600 times signed from Italy

Recognition aspects

- Publications in other channels than journals
- Data publications
- Multiliguality
- Impact beyond research

Links

- Previous events: <https://www.clarin.eu/events>
- Funding for mobility, development of training, CRF projects, event organisation: www.clarin.eu/funding
- Monthly Newsletter: [link](#)
- Teaching and training materials:
 - Recorded lectures: <http://videolectures.net/clarin/>
 - DH Course Registry: <https://dhcr.clarin-dariah.eu/>
 - more to come soon
- CLARIN2021 – Annual Conference, 27-29 September 2021
(PC chair: Monica Monachini, Institute of Computational Linguistics, Pisa)
<https://www.clarin.eu/event/2021/clarin-annual-conference-2021>

see you @

www.clarin.eu

or

f.m.g.dejong@uu.nl

Acknowledgement

Several colleagues generously provided ideas and slides for this presentation:
Darja Fiser, Francesca Frontini, Jakob Leonardic, Dieter Van Uytvanck