

The Slovene CLARIN Infrastructure

Tomaž Erjavec

Department of Knowledge Technologies
Jožef Stefan Institute
Ljubljana, Slovenia

CLARIN.SI



The objective of the Slovene research infrastructure CLARIN.SI is to facilitate research in the humanities and social sciences by offering researchers authorised access to its platform, which will integrate Slovene language resources and advanced tools for processing of Slovene. The Slovene CLARIN is based at the Jožef Stefan Institute and has received initial funding in 2013-2014. In 2014 CLARIN.SI has been established as a consortium of all the main institutions involved in language technologies and linguistic research. Currently, CLARIN.SI offers access to the platform of the project "Communication in Slovene" and to a number of Slovene corpora via a web-based concordancer and has set up its web pages and a test version of the LINDAT repository.

History of Slovene CLARIN

- **2011:** Slovene Government adopts the Research infrastructure Roadmap 2011-2020, which includes CLARIN, as well as DARIAH and CESSDA;
- **2012:** DARIAH and CESSDA establish Slovene infrastructures SI-DIH and ADP;
- **October 2013:** surprise (minimal) funding for CLARIN @ Jožef Stefan Institute;
- **February 2014:** Slovenian Language Technologies Society publishes a call for the creation of the consortium of Slovenian CLARIN.SI;
- **March 2014:** transfer of the www.slovenscina.eu portal to JSI infrastructure
- **June 2014:** establishment of the CLARIN.SI Consortium;
- **July 2014:** extraordinary elections in Slovenia, meaning that the signing of the agreement for Slovenia to join CLARIN ERIC is, once again, postponed.

CLARIN.SI Consortium

Universities

University of Ljubljana

Univerza v Ljubljani



At the University of Ljubljana, corpus linguistics and language technology research is carried out at the Centre for Language Technologies, at the Faculty of Arts (primarily the Department of Translation), Faculty of Social Sciences (Research Centre for the Terminology of Social Sciences and Journalism), Faculty of Electrical Engineering (Laboratory of Artificial Perception, Systems and Cybernetics) and the Faculty of Computer and Information Science (Laboratory for Cognitive Modeling).

University of Maribor

Univerza na Mariboru



Research on language and speech technologies at the University of Maribor is undertaken mainly at the Faculty of Electrical Engineering and Computer Sciences, in the scope of the Institute for Electronics and Telecommunications and the Institute for Computer Science (Laboratory for Heterogeneous Computer systems).

University of Primorska

Univerza na Primorskem



Language technology research and corpus linguistics are undertaken at the Faculty for Mathematics, Natural Sciences and Information Technologies, mainly at the Department of Information Sciences and Technologies.

Companies

Alpineon, d.o.o.

alpineon

The development of hardware and software in the field of language and speech technologies: speech recognition and synthesis, machine translation, voice portals, SMS and e-mail readers.

Amebis, d.o.o.

Kamnik



Language technology software: Slovenian spell and syntax checkers, machine translations, speech synthesis, corpora, online dictionaries, virtual agents, etc.

Institutes

Jožef Stefan Institute

Institut "Jožef Stefan" Ljubljana, Slovenija



At JSI several departments (E3, E8, E9) are involved in language technology research, working in the fields of linguistic annotation of texts, compiling language corpora and other Slovene language resources, text mining and text analytics, machine translation, speech synthesis, etc.

Institute of Contemporary History

Institut za novejšo zgodovino



The institute is, together with the SRC SASA, in charge of the SI-DIH portal, the Slovenian branch of the European research infrastructure for the Humanities DARIAH.

Scientific and Research Centre of the Slovenian Academy of Sciences and Arts



At SRC SASA, the Fran Ramovš Institute for Slovenian Language is the national centre for systematic monitoring and description of Slovene language materials. The results of this research are dictionaries, collections and reference works for Slovene language and scientific papers.

Societies

Slovenian Language Technologies Society

SDJT



The Society was founded in 1998 and joins people working on language technologies from the scientific, educational or user perspectives. Its activities are aimed at promoting the development of language technologies for the Slovenian language.

Domestic Research Society



Domestic Research Society manages the first free online dictionary of the Slovene spoken language, "The Unleashed Tongue", developed in 2004.

Trojina, Institute for Applied Slovene Studies

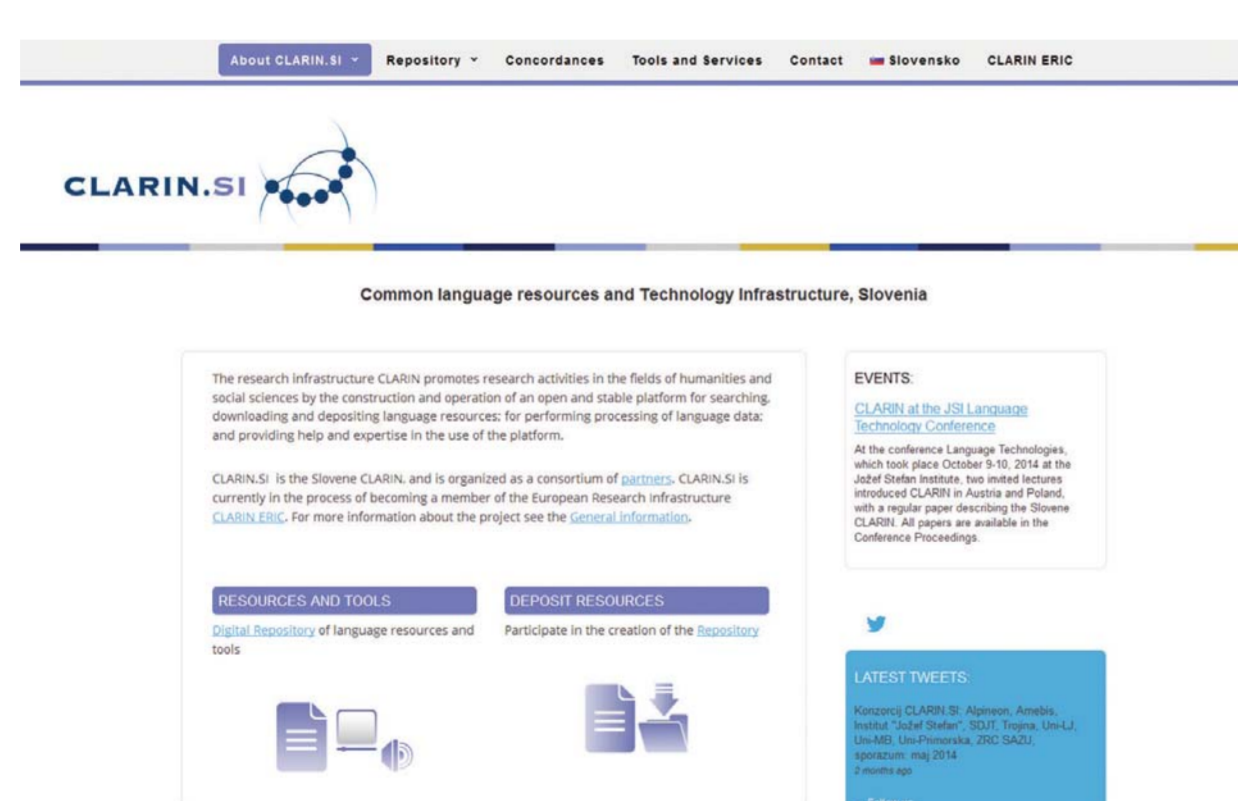
trojina



Trojina is undertaking projects aimed at modern, targeted linguistic research and at increasing the confidence of speakers in public and private use of the Slovene language.

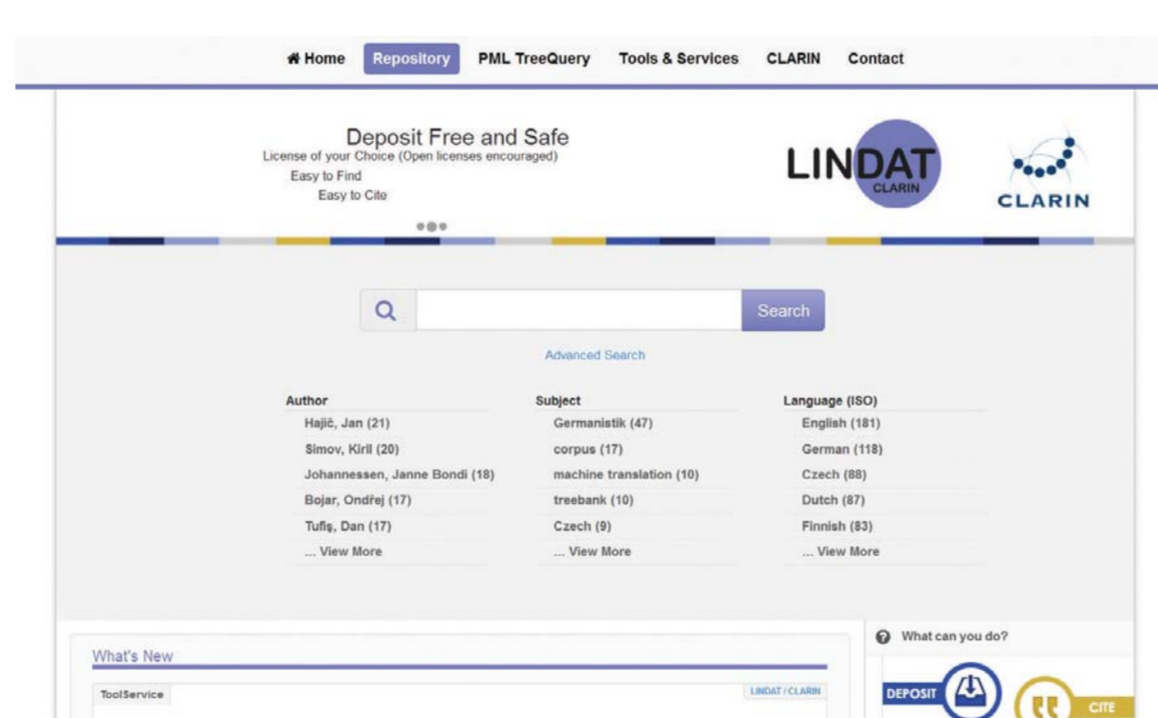
Web page

- Web pages introduce the infrastructure and consortium partners
- Information in Slovene and English
- Links to the repository and services



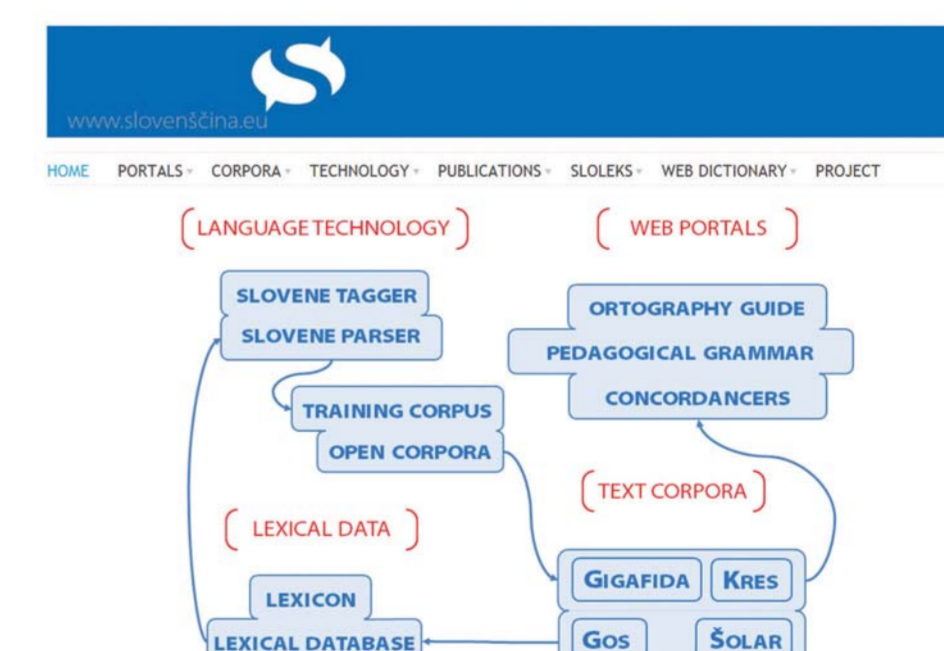
LINDAT Repository

- CLARIN.SI uses the Czech UFAL CLARIN LINDAT repository platform
- Based on DSpace, available on GitHub
- Supports handles and AAI
- Supports handles and AAI
- Currently in testing



www.slovenscina.eu

- CLARIN.SI hosts the platform of the project »Communication in Slovene«
- Web concordancers to reference corpora of Slovene
- Web services for tagging, lemmatising and parsing of Slovene
- Downloadable training corpora, morphological lexicon and lexical database



Searching Slovene Corpora

- Freely available concordancing over many corpora
- Reference and specialised monolingual and parallel corpora
- Linguistic annotations: normalised word-forms, PoS tags, lemmas
- Corpora mounted on noSketchEngine and CUWI concordancers
- Powerful CWB queries and fast processing
- Localised interfaces

Overview of corpora						
Corpus	Accession	Language	Year	Words	Texts	URL
CLARIN.SI	1	Slovene	2014	100,000,000	10,000	http://www.clarin.si
CLARIN.SI	2	Slovene	2014	100,000,000	10,000	http://www.clarin.si
CLARIN.SI	3	Slovene	2014	100,000,000	10,000	http://www.clarin.si
CLARIN.SI	4	Slovene	2014	100,000,000	10,000	http://www.clarin.si
CLARIN.SI	5	Slovene	2014	100,000,000	10,000	http://www.clarin.si
CLARIN.SI	6	Slovene	2014	100,000,000	10,000	http://www.clarin.si
CLARIN.SI	7	Slovene	2014	100,000,000	10,000	http://www.clarin.si
CLARIN.SI	8	Slovene	2014	100,000,000	10,000	http://www.clarin.si
CLARIN.SI	9	Slovene	2014	100,000,000	10,000	http://www.clarin.si
CLARIN.SI	10	Slovene	2014	100,000,000	10,000	http://www.clarin.si

