

Raziskovalna infrastruktura CLARIN.SI

Tomaž Erjavec

Odsek za tehnologije znanja, Institut "Jožef Stefan"

Delavnica Odprta znanost
Pravna fakulteta
23. januar 2020



<http://www.clarin.si/>

- Začetek dela v 2014
- Institut "Jožef Stefan":
 - Odsek za tehnologije znanja (E8)
 - Laboratorij za umetno inteligenco (E3)
 - Center za mrežno infrastrukturo (CMI)
- Organiziran kot konzorcij 11 partnerjev:
 - 4 univerze: Ljubljana, Maribor, Nova Gorica, Primorska
 - 3 raziskovalni inštituti: ZRC SAZU, IJS, INZ
 - 2 podjetji: Amebis, Alpineon
 - 1 društvo + 1 zavod: SDJT, Trojina

Trije stebri:

- Certificiran repozitorij jezikovnih virov in orodij
 - dolgotrajno hranjenje
 - avtentikacija in avtorizacija
 - stalni identifikatorji
 - eksplisitni pogoji uporabe in licence
 - principi FAIR
- Spletne storitve
 - dva konkordančnika (spletne analize korpusov)
 - orodja za označevanje besedil
 - itd.
- Podpora:
 - financiranje priprave virov za vključitev v repozitorij
 - večji projekti: 30.000 EUR letno: 7 projektov v 2018, 6 v 2019
 - dogodki: JOTA, JT-DH 2018, EURALEX 2018, TSD 2019, ...
 - CLASLA K-centre: CLARIN.SI center znanja za računalniško obdelavo južnoslovanskih jezikov

Konkordančnika

- noSketch Engine + KonText
- Obdelava velikih korpusov (več milijard besed)
- Korpsi so lahko bogato označeni:
 - strukture (besedilo, odstavek, termin, ...)
 - metapodatki (datum, vrsta besedila, spol avtorja, standardnost, ...)
 - atributi pojavnic (oblikoskladenjska oznaka, lema, normalizirana oblika, ...)
- Bogat poizvedovalni jezik: CQL
- Raznovrstni izpisi in analize
- RESTful, tj. poizvedbe prek URLjev
- Dostopnih 70 korpusov v 27 jezikih in s prek 15 milijard besed

noSketch Engine

The screenshot shows a web browser window for the CLARIN.si CLASSLA knowledge centre. The URL is https://www.clarin.si/noske/run.cgi/first_form?corpusname=gfida20_dedup;align=. The page title is "Gigafida v2.0 DeDup (referenčni, dedupliciran)". On the left, there's a sidebar with links: Home, Search, Word list, Corpus Info, My jobs, and User guide. The main area contains a search form for the "Gigafida v2.0 DeDup (referenčni, dedupliciran)" corpus. The search term "krava" is entered in the "Simple query" field. Below it, the "Query type" is set to "simple". Other fields include "Lemma:", "Phrase:", "Word form:", "Character:", and "CQL:". Buttons for "Make Concordance" and "Tagset summary CQL builder" are present. At the bottom of the search form are "Make Concordance" and "Clear All" buttons. In the bottom right corner, there are logos for CLARIN.si, Lexical Computing, and the version number 2.36.7-open-2.158.8-open-3.105.1.

- Odprtokodna različica Sketch Engine
- Prijava ni niti potrebna niti mogoča

KonText

The screenshot shows the KonText web application running in a browser window. The title bar reads "Concordance - First query". The URL in the address bar is "Vamo | https://www.clarin.si/kontext/first_form?corpname=janes". The top navigation bar includes links for "Repository", "About", "Contact", and the CLARIN logo. A user profile is shown on the right with the name "Tomaz;Tomaž Erjavec" and a "logout" button.

The main interface features a logo with "kontext" and a navigation menu with links for "Query", "Corpora", "Save", "Concordance", "Filter", "Frequency", "Collocations", "View", and "Help". Below the menu, it says "Corpus: Janes (družbena omrežja)".

A large green search form is centered. It has fields for "Corpus" (set to "Janes (družbena omrežja)"), "Query Type" (set to "Basic"), and "Query" (containing the word "krava"). There are two expandable sections: "Specify context" and "Specify query according to the meta-information", both currently collapsed. At the bottom of the form is a blue "Search" button.

- Razvil ga je češki CLARIN
- Prijava: shranjene poizvedbe, nastavitev zaslona, podkorpsi
- Nima nekaterih funkcionalnosti noSketch Engine

- Zaenkrat najbolj pomembna storitev CLARINA
- Stalna in varna hramba jezikovnih virov ([https](https://), Nagios)
- Eksplisitni pogoji uporabe (ToS, licenca), etični kodeks (CoCo)
- Standarden zapis metapodatkov: Component Metadata Infrastructure (CMDI) Dublin Core (DC)
- Žetev metapodatkov (metadata harvesting)
- Večinoma standarden zapis podatkov (XML, TEI)
- Trenutno vsebuje 159 virov, od tega 121 slovenskih

Anatomija vnosa, 1

The screenshot shows a web browser window for the CLARIN.SI repository. The URL is <https://www.clarin.si/repository/xmlui/handle/11356/1236>. The page title is "Slovenian parliamentary corpus siParl 1.0 (1990-2018)". The main content area displays the following information:

- Za citiranje vnosa uporabite naslednjo referenco ali jo izvitezte v prednastavljeno obliko:**
Pančur, Andrej; Erjavec, Tomaz; Ojsteršek, Mihael; Šorn, Mojca and Blaj Hribar, Neja, 2019, *Slovenian parliamentary corpus siParl 1.0 (1990-2018)*, Slovenian language resource repository CLARIN.SI, <http://hdl.handle.net/11356/1236>.
- Ta vir je integriran tudi v naslednje storitve:**
[KonText](#) [noSketch](#)
- Avtorji:** Pančur, Andrej ; Erjavec, Tomaž ; Ojsteršek, Mihael ; Šorn, Mojca ; Blaj Hribar, Neja
- Item identifier:** <http://hdl.handle.net/11356/1236>
- URL projekta:** <https://github.com/DARIAH-SI/siParl/commit/c6e7942b9fb2199a85e60de6dd30679ce735cf1a>
- Demo URL:** <http://exist.sistory.si/exist/apps/parla/index.html>
- Datum objave:** 2019-05-03

The right sidebar contains links for **CLARIN.SI Data & Tools**, **Brskaj**, **Moj račun**, **Prijava**, **Statistike**, and **Statistika Piwik**. There is also a **BETA** button.

- osnovni metapodatki; citiranje; integracija s storitvami; povezave na projekt, demo & objave
- lokalizacija, orodna vrstica, prijava; iskanje, osnovne informacije o repozitoriju, brskanje; Piwik

Anatomija vnosa, 2

| | | |
|------------|--|-------------------------|
| Vrsta | corpus | Splošne informacije |
| Velikost | 11351 texts, 1083233 utterances, 227896145 tokens | Prijava |
| Jezik(i) | Slovenian | O vnosu v repozitorij |
| Opis | <p>The siParl corpus contains minutes of the Assembly of the Republic of Slovenia for 11th legislative period 1990-1992, minutes of the National Assembly of the Republic of Slovenia from the 1st to the 7th legislative period 1992-2018, minutes of the working bodies of the National Assembly of the Republic of Slovenia from the 2nd to the 7th legislative period 1996-2018, and minutes of the Council of the President of the National Assembly from the 2nd to the 7th legislative period 1996-2018. The corpus comprises over a million speeches or 195 million words. The corpus contains basic meta-data about the speakers, a typology of sessions etc. and structural and editorial annotations.</p> <p>This item comprises three datasets:</p> <ul style="list-style-type: none">- the corpus in TEI (module Transcriptions of speech);- the corpus in TEI with added automatic linguistic annotation: tokenisation, MSD tagging and lemmatisation;- the linguistically annotated corpus in vertical format used by various concordancers, e.g. CWB and Sketch Engine; this format is simpler and smaller but does not contain all the information from the source TEI. <p>A preliminary version of this resource is presented in the paper: Pančur, Andrej, Mojca Šom and Tomaz Erjavec (2018). "SlovParl 2.0: The Collection of Slovene Parliamentary Debates from the Period of Secession." Darja Fišer and María Eskovich and Francisca de Jong (eds.), Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018), Miyazaki, Japan, May 2018. http://lrec-conf.org/workshops/lrec2018/W2/summaries/4_W2.html</p> | Citiranje |
| Izdajatelj | Institute of Contemporary History | Življenski ciklus vnosa |
| | | Pogosta vprašanja |
| | | O repozitoriju |
| | | Pomoč uporabnikom |

- vrsta, velikost, jezik, opis, založnik
- več informacij o repozitoriju

Anatomija vnosa, 3

Subject(s)

parliamentary debates

Slovenian Parliament

TEI

Collection(s)

CLARIN.SI data & tools

Show full item record

Files in this item



Download instructions for command line

This item is **Publicly Available** and licensed under:
Creative Commons - Attribution 4.0 International (CC BY 4.0)



| | |
|-------------|----------------------------------|
| Name | siParl.zip |
| Size | 477.97 MB |
| Format | application/zip |
| Description | Corpus in TEI format |
| MD5 | 08c83cecc1ac42e2cf5a28652aac996d |



[Download file](#) [Preview](#)

| | |
|-------------|---|
| Name | siParl-ana.zip |
| Size | 1.18 GB |
| Format | application/zip |
| Description | Corpus in TEI format, with linguistic annotations |
| MD5 | 91929e140d2d7fa1426748700d21ebcb |



[Download file](#) [Preview](#)

- ključne besede, vsi metapodatki
- licenca, prevzem podatkov

Nekaj vnosov

- **ccGigafida:** Referečni korpus standardne slovenščine
vzorčeni odstavki korpusa Gigafida (100 mil. besed)
- **Janes:** Korpusi družbenih omrežij
tviti, blogi, komentarji (250 mil. besed)
- **IMP:** Jezikovni viri starejše slovenščine
600 knjig, 1700 - 1918, (17 mil. besed)
- **siParl:** Slovenski parlament
zapisniki slovenskega parlamenta 1990-2018 (195 mil. besed).
- **Gos VideoLectures:** Govorni korpus
55 predavanj, 22 ur govora (180,000 besed)

Zaključki

- Namen CLARIN(.SI) je spodbujati raziskave, ki potrebujejo dostop do jezikovnih podatkov
 - digitalna humanistika in družboslovje
 - jezikovne tehnologije
 - vse ostale vede, kjer je jezik pomemben
- Odprt dostop do virov, orodij in storitev
- Avtentikacija AAI, citiranje (handle)
- Prek 70 korpusov na konkordančnikih
- Prek 150 virov v repozitoriju
- Slovenski raziskovalci imajo dostop tudi do storitev CLARIN ERIC in drugih nacionalnih konzorcijev CLARIN

Raziskovalna infrastruktura CLARIN.SI

Tomaž Erjavec

Odsek za tehnologije znanja, Institut "Jožef Stefan"

Delavnica Odprta znanost
Pravna fakulteta
23. januar 2020