# The CLARIN.SI repository

Tomaž Erjavec

Department of Knowledge Technologies, Jožef Stefan Institute

FRI, 29th April, 2020

# Repository

- Currently the most important service of CLARIN.SI
- Long-term and safe archiving of language resources (https, Nagios)
- Explicit terms of use (terms of service, licence)
- Ethical codex (Code of conduct)
- Most resources in standard encoding (Unicode, XML)
- Certified as CLARIN Centre B
  - Digital Seal of Approval, DSA
  - Currently being (re)certified for Core Trust Seal, CTS

# Repository platform

- Based on the DSpace platform, developed for open digital repositories
- DSpace modified for use by CLARIN repositories: CLARIN/DSpace developed by Czech CLARIN
- Used by Czech, Italian, Norwegian, Polish, and Slovenian CLARIN
- Maintenance on GitHub, Slovenian fork on GitLab

# Metadata

- Standard encoding of metadata:
  - Component Metadata Infrastructure (CMDI)
  - Dublin Core (DC)
- Metadata always CC0
- Meta-data harvesting:
  - CLARIN Virtual Language Observatory (VLO)
  - also:

# Permanent identifiers

- Each repository entry is assigned a PID
- So, even if the repository platform is changed, the PIDs stay the same (but need to be reconfigured)
- Best known solution: DOI
- CLARIN used the Handle system
- **http://hdl.handle.net/11356/1044** → https://www.clarin.si/repository/xmlui/handle/11356/1044
- Important for citation of repository items:

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

Krsnik, Luka; Dobrovoljc, Kaja and Robnik-Šikonja, Marko, 2019, *Dependency tree extraction tool STARK 1.0*, Slovenian language resource repository CLARIN.SI, http://hdl.handle.net/11356/1284.

# Top level page

# Top level page

**Author(s):**
Pančur, Andrej ; Erjavec, Tomaž ; Ojsteršek, Mihael ; Šorn, Mojca ; Blaj Hribar, Neja

**Description:**
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 ...

🔗 **This item contains 6 files (11.43 GB).**

**Publicly Available** ⓒ ①

**LexicalConceptualResource**　　　　　　　　　　　　**CLARIN.SI Data & Tools**

## Consonant-vowel structures in the GOS 1.0 corpus

**Author(s):**
Čibej, Jaka ; Arhar Holdt, Špela ; Dobrovoljc, Kaja ; Krek, Simon

**Description:**
The lists contain consonant-vowel structures of all lemmas, word forms, and normalized word forms in the GOS 1.0 Corpus of Spoken Slovene (http://hdl.handle.net/11356/1040). In each unit, its characters were converted as ...

🔗 **This item contains 7 files (3.6 MB).**

**Publicly Available** ⓒ ① ⓪

**LexicalConceptualResource**　　　　　　　　　　　　**CLARIN.SI Data & Tools**

## Consonant-vowel structures in the Gigafida 2.0 corpus

**Author(s):**
Čibej, Jaka ; Arhar Holdt, Špela ; Dobrovoljc, Kaja ; Krek, Simon

**Description:**
The lists contain consonant-vowel structures of all lemmas and word forms in the Gigafida 2.0 corpus. In each unit, its characters were converted as follows: C - consonant (in lists with finegrained character categorizations, ...

🔗 **This item contains 5 files (141.75 MB).**

**Publicly Available** ⓒ ① ⓪

◎ **Browse**
> All of the Repository

👤 **My Account**
➡ Login

ℹ **General Information**
⬆ Deposit
❝ Cite
🔄 Submission Lifecycle
❓ FAQ
ⓘ About
✉ Help Desk

🔊 **RSS Feed**
🔊 RSS 1.0
🔊 RSS 2.0
🔊 Atom

# Log-in

- Necessary for making a new repository item, for accessing non-CC items and for using some advanced functions
- For those without EduGain, CLARIN ERIC also gives accounts

# How to find interesting resources

- Browsing by language, type of resource, keywords, author, etc.
- Search (fuzzy matching)
- Advanced facet search
- Currently 163 items (without prior versions)

| Type | |
| --- | --- |
| corpus | 78 |
| languageDescription | 2 |
| lexicalConceptualResource | 66 |
| toolService | 16 |

# Landing page of a repository item 1.

# Landing page of a repository item 2

**📄 Description**

The ssj500k training corpus contains about 500,000 tokens manually annotated on the levels of tokenisation, sentence segmentation, morphosyntactic tagging, and lemmatisation. About half of the corpus is also manually annotated with syntactic dependencies, named entities, and verbal multiword expressions. About a quarter of the corpus is annotated with semantic role labels. The morphosyntactic tags and syntactic dependencies are included both in the JOS/MULTEXT-East framework, as well as in the framework of Universal Dependencies.

The annotations of the ssj500k corpus follow (1) the MULTEXT-East V6 morphosyntactic specifications for Slovene, http://nl.ijs.si/ME/V6/msd/, (2) the JOS dependency schema, http://nl.ijs.si/jos/bib/jos-skladnja-navodila.pdf, the Universal Dependencies morphosyntactic specifications and syntactic dependencies for Slovene-SSJ, https://universaldependencies.org/, (4) the Janes annotation guidelines for Slovenian named entities, http://nl.ijs.si/janes/wp-content/uploads/2017/09/SlovenianNER-eng-v1.1.pdf, and (5) the Guidelines of the PARSEME shared task on verbal multiword expressions, http://parsemefr.lif.univ-mrs.fr/parseme-st-guidelines/1.1/
The vocabulary of (1) and (2) is provided in the back element and (3), (4), and (5) in the teiHeader of the TEI encoded corpus. The semantic role labels are also documented in the teiHeader.

In contrast to the previous version 2.1, this version corrects various errata in spacing and text metadata and adds UD morphological and (where it was possible to do so automatically) dependency annotations to the corpus. Note that the UD annotations are not included in the vertical file.

**📄 Publisher**

Centre for Language Resources and Technologies, University of Ljubljana

**📄 Acknowledgement**

Ministry of Education, Science and Sport 3311-08-986003 "Communication in Slovene"
ARRS (Slovenian Research Agency) P2-103 "Knowledge Technologies"
ARRS (Slovenian Research Agency) J6-8256 "New grammar of contemporary standard Slovene: sources and methods"
ARRS (Slovenian Research Agency) MR-37487 "Young Researcher Programme"
ARRS (Slovenian Research Agency) P6-0411 "Language Resources and Technologies for Slovene"

**🏷 Subject(s)**

part-of-speech tagging | dependency treebank | parsing | named entities | tokenisation | manual annotation | TEI | verbal multiword expressions | semantic role labelling | CONLL-U

# Landing page of a repository item 3

⑁ **Other versions**

[ List all versions ▾ ]

Show full item record

🔖 Files in this item

| | |
|---|---|
| ⬇ **Download instructions for command line** | ⬇ **Download all files in item (40.95 MB)** |

This item is **Publicly Available** and licensed under:
Creative Commons - Attribution-NonCommercial-ShareAlike 4.0 International (CC BY-NC-SA 4.0)
ⓒ ① ⓢ ⟲

| | |
|---|---|
| **Name** | ssj500k.conllu.zip |
| **Size** | 10 MB |
| **Format** | application/zip |
| **Description** | Corpus in CONLL-U format, complete corpus with UD morphology and separately the UD syntactically annotated part, also split into train/dev/test. |
| **MD5** | f65ae2995a2a7acfe43b1a5aa3140dca |

⊙ Download file    👁 Preview

| | |
|---|---|
| **Name** | ssj500k-en.TEI.zip |
| **Size** | 11.92 MB |

# DC metadata

| | |
|---|---|
| dc.contributor.author | Holz, Nanika |
| dc.contributor.author | Zupan, Katja |
| dc.contributor.author | Gantar, Polona |
| dc.contributor.author | Kuzman, Taja |
| dc.contributor.author | Čibej, Jaka |
| dc.contributor.author | Arhar Holdt, Špela |
| dc.contributor.author | Kavčič, Teja |
| dc.contributor.author | Škrjanec, Iza |
| dc.contributor.author | Marko, Dafne |
| dc.contributor.author | Jezeršek, Lucija |
| dc.contributor.author | Zajc, Anja |
| dc.date.accessioned | 2019-01-26T20:37:28Z |
| dc.date.available | 2019-01-26T20:37:28Z |
| dc.date.issued | 2019-01-26 |
| dc.identifier.uri | http://hdl.handle.net/11356/1210 |
| dc.description | The ssj500k training corpus contains about 500,000 tokens manually annotated on the levels of tokenisation, sentence segmentation, morphosyntactic tagging, and lemmatisation. About half of the corpus |

# Piwik



**Statistics**

Piwik Statistics  BETA

> View Usage Statistics

✉ **Get a monthly report for this item via email**  **Subscribe**

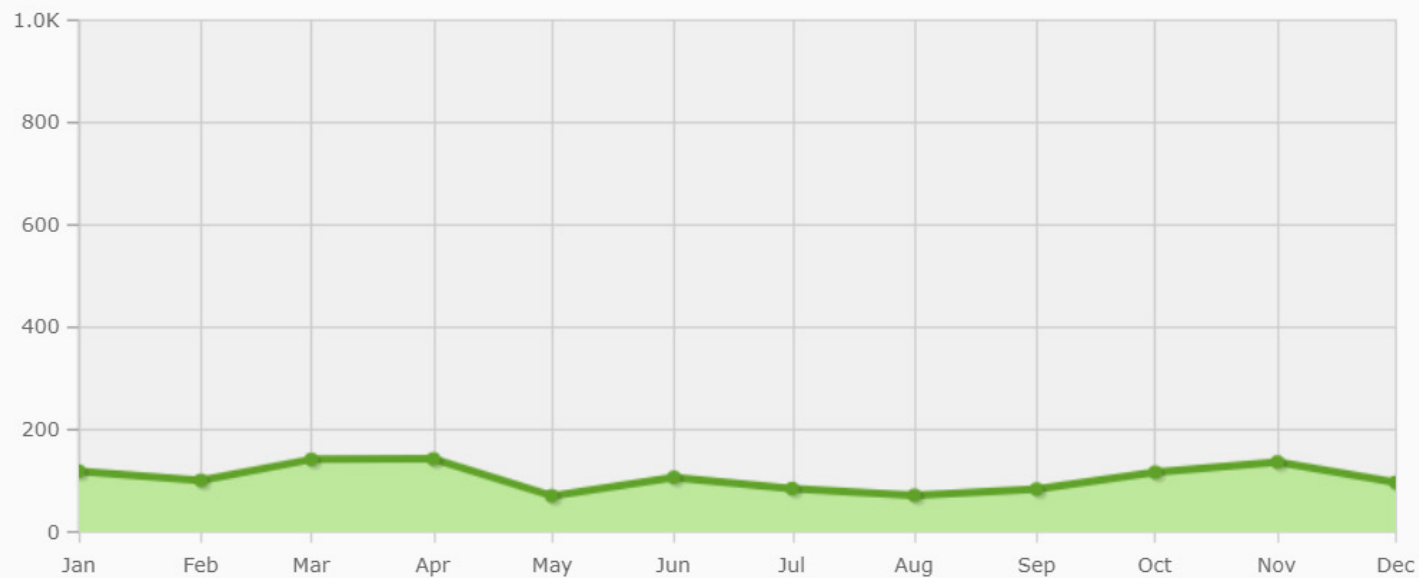ⓘ Click on a data point to summarize by year / month.

👁 Views **1275**

⬇ Downloads **18297**

**Statistics for the year 2018**  **Back**

# Repository: download

- Most entries available under Creative Commons licences
- Others: AAI login, agree to licence conditions: CLARIN.SI Licence ACA ID-BY-NC-INF-NORED
- Cite the resource!
- Always use the handle, not the URL!

Emoji Sentiment Ranking 1.0

66 **Please use the following text to cite this item or export to a predefined format:**   BIBTEX   CMDI

Kralj Novak, Petra; Smailović, Jasmina; Sluban, Borut and Mozetič, Igor, 2015, *Emoji Sentiment Ranking 1.0*, Slovenian language resource repository CLARIN.SI, **http://hdl.handle.net/11356/1048**.

Share: 

CLARIN.SI Data & Tools

| ✏ Authors | Kralj Novak, Petra ; Smailović, Jasmina ; Sluban, Borut ; Mozetič, Igor |
| --- | --- |
| ➦ Item identifier | http://hdl.handle.net/11356/1048 |
| ☑ Demo URL | http://kt.ijs.si/data/Emoji_sentiment_ranking/ |
| ⚲ Referenced by | https://doi.org/10.1371/journal.pone.0144296 |

# Repository: deposit

- AAI log-in
- Fairly simple workflow:

Item submission



| 1. | 2. | 3. | 4. | 5. | 6. | 7. | 8. |
| Basic Info | Who's involved | Describe | Upload | License | Note | Review | Complete |

- However: follow best practices, look at other entries!
- Choosing the licence
- Possibility to embargo the data
- Can enter a new version of a resource
- Option to share item editing
- The entry is checked by one of the CLARIN.SI editors

# Data Formats

- Media files: wav, avi, jpg, png…
- Data with tabular structure: TSV, CSV, XML
- Hierarchical data: XML
  - using one of the „standard" schemas, e.g. TEI, …, TRS, ELAN
  - using your own schema
  - in all cases: schema + documentation must be part of the entry
- + Derived formats + Documentation (can be PDF)
- Deposit as ZIP

📎 Datoteke v tem vnosu

| ⬇ Prenesi navodila za ukazno vrstico | ⬇ Prevzem vseh datotek v viru (52.35 MB) |
|---|---|

Ta vnos je **Publicly Available** z licenco:
Creative Commons - Attribution-NonCommercial-ShareAlike 4.0 International (CC BY-NC-SA 4.0)
🅭🅯🅏🄍

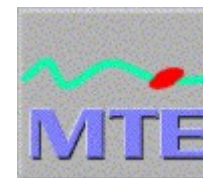| | |
|---|---|
| Ime | Sloleks_v1.2.zip |
| Velikost | 26.96 MB |
| Format | application/zip |
| Opis | Sloleks in LMF XML format, PoS tags in Slovenian. |
| MD5 | 3b15dc1a094e3ff2f4f1ef69702e6f41 |

⬇ Prevzem datoteke   👁 Predogled

| | |
|---|---|
| Ime | sloleks-sl.tbl_v1.2.zip |
| Velikost | 12.65 MB |
| Format | application/zip |
| Opis | Sloleks in tabular format, PoS tags in Slovenian. |
| MD5 | 0c2886e88558df5a6b9cdde4c7e20fb3 |

⬇ Prevzem datoteke   👁 Predogled

# What is in the repository?

- Language models, training data, lexicons for tagging, lemmatisation, syntax of Slovene, Croatian, Serbian, for standard and non-standard language

- Various types of word embeddings

- Large annotated corpora of various text types

- Speech corpora

- Multilingual corpora

- Machine readable dictionaries

- Some software

# Tools

- "Language models, training data, lexicons
  for tagging, lemmatisation, syntax of
  Slovene, Croatian, Serbian,
  for standard and non-standard language"
- tokenisers, morphosyntactic taggers for sl, hbs
- parsers (also UD-PIPE)
- Viewers and editors for linguistic annotation

# CLARINSI@GitHub

**clarin-dspace**

Forked from ufal/clarin-dspace

LINDAT/CLARIN digital repository based on DSpace

● Java  942  ★ 0  ① 8  0  Updated 6 days ago

**mte-msd**

MULTEXT-East morphosyntactic specifications

● HTML  0  ★ 1  ① 0  0  Updated 21 days ago

**babushka-bench**

Benchmarking NLP tools on Slovene, Croatian and Serbian

● Python  1  ★ 2  ① 1  0  Updated on Mar 19

**classla-stanfordnlp**

Forked from stanfordnlp/stanza

CLASSLA Fork of the Official Stanford NLP Python Library for Many Human Languages

● Python  505  ★ 2  ① 0  0  Updated on Mar 19

**reldi-tokeniser**
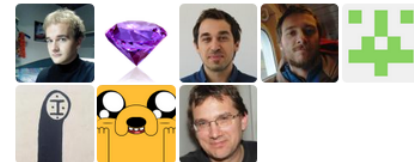
## Top languages

● Python  ● Java  ● HTML  ● Shell
● C

## People  8 >

Invite your teammates...

Invite