



Easily Accessible Language Technologies for Slovene, Croatian and Serbian

Nikola Ljubešić,^{1,2} Tomaž Erjavec,¹ Darja Fišer,^{1,3} Tanja Samardžić⁴,
Maja Miličević⁵, Filip Klubička², Filip Petkovski⁶

¹ Department of Knowledge Technologies, Jožef Stefan Institute, Ljubljana

² Faculty of Humanities and Social Sciences, University of Zagreb

³ Faculty of Arts, University of Ljubljana

⁴ CorpusLab, University of Zürich

⁵ Faculty of Philology, University of Belgrade

⁶ Freelance developer, Macedonia

Projects

- ReLDI (Regional Linguistic Data Initiative) – SNSF
- JANES (Jezikovna Analiza NEstandardne Slovenščine) – ARRS

Overview

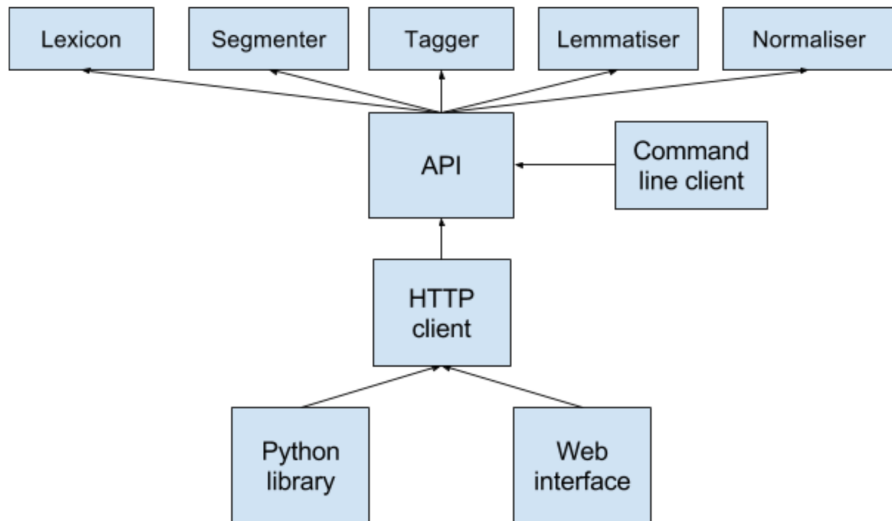
Projects

- ReLDI (Regional Linguistic Data Initiative) – SNSF
- JANES (Jezikovna Analiza NEstandardne Slovenščine) – ARRS

Included technologies

- Morphological lexicons (Sloleks 101k lexemes, hrLex 187k lexemes, srLex 193k lexemes)
- Segmentation
- Diacritic restoration (99% standard, 98% non-standard language)
- Morphosyntactic tagging (92-93% accuracy)
- Lemmatisation (98% accuracy)

Architecture



The near future

Technologies to follow

- Dependency parsing (<http://universaldependencies.org>)
- Named entity recognition (PERS, LOC, ORG, MISC)
- Text normalisation via character-level SMT – non-standard and historical language (error reduction of 30-80%, downstream processing error reduction of (up to) 50%)
- Semantic Role Labeling (Slovene-Croatian bilateral project)

The near future

Technologies to follow

- Dependency parsing (<http://universaldependencies.org>)
- Named entity recognition (PERS, LOC, ORG, MISC)
- Text normalisation via character-level SMT – non-standard and historical language (error reduction of 30-80%, downstream processing error reduction of (up to) 50%)
- Semantic Role Labeling (Slovene-Croatian bilateral project)

Migration to CLARIN

- All web services are developed to be easily included in WebLicht
- Received funding from CLARIN to add existing and emerging technologies, workshop in Ljubljana in November



Easily Accessible Language Technologies for Slovene, Croatian and Serbian

Nikola Ljubešić,^{1,2} Tomaž Erjavec,¹ Darja Fišer,^{1,3} Tanja Samardžić⁴,
Maja Miličević⁵, Filip Klubička², Filip Petkovski⁶

¹ Department of Knowledge Technologies, Jožef Stefan Institute, Ljubljana

² Faculty of Humanities and Social Sciences, University of Zagreb

³ Faculty of Arts, University of Ljubljana

⁴ CorpusLab, University of Zürich

⁵ Faculty of Philology, University of Belgrade

⁶ Freelance developer, Macedonia