

# 05 Lematizacija

Lema (osnovna besedna oblika) je pri označevanju pripisana vsem besednim oblikam v besedilu, kar omogoča njihovo nadaljnje enovito procesiranje. Sistem lematizacije je bil razvit v projektu Sporazumevanje v slovenskem jeziku (Holozan et al. 2008) in sledi sistemu oblikoskladenjskega označevanja MULTEXT-East v4 oz. JOS pri določanju besedne vrste, velike začetnice in nekaterih drugih značilnosti.

- [Označevalne smernice](#)
- [Reference in povezave](#)

# Označevalne smernice

V tem poglavju so zbrane označevalne smernice za lematizacijo. Smernice so razvrščene od nastarejše različice do zadnje, ažurne različice.

## **Različica 2.0 (25-02-2023)**

### **projekt [Razvoj slovenščine v digitalnem okolju](#)**

HOLOZAN, Peter, KREK, Simon, PIVEC, Matej, RIGAČ, Simon, ROZMAN, Simon, VELUŠČEK, Aleš, PORI, Eva, ARHAR HOLDT, Špela, 2023: *Specifikacije za učni korpus: lematizacija in MSD*. Različica 2.0. Rezultat projekta Razvoj slovenščine v digitalnem okolju.

[\[DOCX\]](#) [\[PDF\]](#)

## **Različica 1.0 (2008)**

### **projekt [Sporazumevanje v slovenskem jeziku](#)**

HOLOZAN, Peter, KREK, Simon, PIVEC, Matej, RIGAČ, Simon, ROZMAN, Simon, VELUŠČEK, Aleš, 2008: *Specifikacije za učni korpus*. Različica 1.0. Kazalnik K2 projekta Sporazumevanje v slovenskem jeziku. [\[PDF\]](#)

## **Različica 1.0 za nestandardno slovenščino (21-12-2016)**

### **projekt [Janes - Viri, orodja in metode za raziskovanje nestandardne spletne slovenščine](#)**

ČIBEJ, Jaka, ARHAR HOLDT, Špela, ERJAVEC, Tomaž, FIŠER, Darja, ZUPAN, Katja, 2016: *Smernice za označevanje računalniško posredovane komunikacije: tokenizacija, stavčna segmentacija, normalizacija, lematizacija in oblikoskladenjsko označevanje*. Različica 1.0. Rezultat projekta Viri, orodja in metode za raziskovanje nestandardne spletne slovenščine. [\[PDF\]](#)

# Reference in povezave

V tem poglavju so zbrane relevantne reference in povezave na projekte, v katerih se je označevalni sistem razvijal ter uporabljal.

## **Projekti, na katerih se je sistem razvijal oz. uporabljal:**

[JOS - Jezikoslovno označevanje slovenskega jezika: metode in viri](#)

[Sporazumevanje v slovenskem jeziku](#)

[Janes - Viri, orodja in metode za raziskovanje nestandardne spletne slovenščine](#)

[Razvoj slovenščine v digitalnem okolju](#)

## **Učni korpusi z ročno pregledanimi lemmi**

Arhar Holdt, Špela; Krek, Simon; Dobrovoljc, Kaja; Erjavec, Tomaž; Gantar, Polona; Čibej, Jaka; Pori, Eva; Terčon, Luka; Munda, Tina; Žitnik, Slavko; Robida, Nejc; Blagus, Neli; Može, Sara; Ledinek, Nina; Holz, Nanika; Zupan, Katja; Kuzman, Taja; Kavčič, Teja; Škrjanec, Iza; Marko, Dafne; Jezeršek, Lucija; Zajc, Anja, 2022, Training corpus SUK 1.0, Slovenian language resource repository

CLARIN.SI, ISSN 2820-4042, <http://hdl.handle.net/11356/1747>.

Lenardič, Jakob; Čibej, Jaka; Arhar Holdt, Špela; Erjavec, Tomaž; Fišer, Darja; Ljubešič, Nikola; Zupan, Katja; Dobrovoljc, Kaja, 2022, CMC training corpus Janes-Tag 3.0, Slovenian language resource repository CLARIN.SI, ISSN 2820-4042, <http://hdl.handle.net/11356/1732>.

## **Reference**

Erjavec, Tomaž; Fišer, Darja; Krek, Simon in Ledinek, Nina. 2010. The JOS Linguistically Tagged Corpus of Slovene. V: Proceedings of the Seventh conference on International Language Resources and Evaluation (LREC'10), Valeta, Malta, Maj. European Language Resources Association (ELRA).

[http://www.lrec-conf.org/proceedings/lrec2010/pdf/139\\_Paper.pdf](http://www.lrec-conf.org/proceedings/lrec2010/pdf/139_Paper.pdf) [\[PDF\]](#)

Erjavec, Tomaž. 2012. MULTEXT-East: morphosyntactic resources for Central and Eastern European languages. Language Resources and Evaluation, 46(1): 131–142. DOI: [10.1007/s10579-011-9174-8](https://doi.org/10.1007/s10579-011-9174-8)

[\[PDF\]](#)

Krek, Simon; Erjavec, Tomaž; Dobrovoljc, Kaja; Gantar, Polona; Arhar Holdt, Špela; Čibej, Jaka in Brank, Janez. The ssj500k training corpus for Slovene language processing. V: Fišer, D. in Erjavec, T. Jezikovne tehnologije in digitalna humanistika: zbornik konference: 24.-25. september 2020, Ljubljana, Slovenija. Ljubljana: Inštitut za novejšo zgodovino, 2020. Str. 24–33.

[http://nl.ijs.si/jtdh20/pdf/JT-DH\\_2020\\_Krek-et-al\\_The-ssj500k-Training-Corpus-for-Slovene-Language-Processing.pdf](http://nl.ijs.si/jtdh20/pdf/JT-DH_2020_Krek-et-al_The-ssj500k-Training-Corpus-for-Slovene-Language-Processing.pdf) [\[PDF\]](#)