

Implementing Semantic Annotation for a Ukrainian Corpus

Full-fledged semantic annotation is currently not available for any sizable Ukrainian corpus. We have designed a semantic tagset for Ukrainian which is being implemented to add semantic annotation to the General Regionally Annotated Corpus of Ukrainian (GRAC).¹ We have adopted a taxonomic approach to semantic annotation and are leveraging the proven tools that exist for Ukrainian corpora. Semantic tags are assigned to lemmas in the Large Electronic Dictionary of Ukrainian (VESUM)² from the r2u team, yielding a semantic lexicon. This enriched dictionary will then be used by the TagText tagger³ to add both POS and semantic tags to the GRAC corpus.

We discuss the particulars of implementation and analyze the challenges that arise in the process of developing and implementing semantic annotation for Ukrainian.

A combination of grammatical and semantic annotation adds valuable flexibility and functionality to the corpus, enabling a variety of finely-formulated search queries and expanding the possibilities for linguistic corpus research.

Keywords: semantic annotation, semantic lexicon, Ukrainian, corpus, GRAC, VESUM, TagText.

References

1. Shvedova, Maria, Ruprecht von Waldenfels, Sergiy Yarygin, Andriy Rysin, Vasyl Starko, Michał Woźniak, Mikhail Kruk et al. (2017-2020): *GRAC: General Regionally Annotated Corpus of Ukrainian*. Electronic resource: Kyiv, Lviv, Jena. Available at uacorporus.org.
2. Rysin, Andriy, Vasyl Starko, and the BrUK team, *Large Electronic Dictionary of Ukrainian (VESUM)*. Electronic resource. Available at: github.com/brown-uk/dict_uk
3. Rysin, Andriy. *TagText tagger for Ukrainian*, available at: github.com/brown-uk/nlp_uk