

The ELEXIS parallel sense-annotated corpus



Simon Krek¹, Carole Tiberius², Kaja Dobrovoljc¹, Jaka Čibej¹, Polona Gantar³, Jelena Kallas⁴, Kristina Koppel⁴, Svetla Koeva⁵, Veronika Lipp⁶, László Simon⁶

¹Jožef Stefan Institute, Slovenia, ²Instituut voor de Nederlandse Taal, The Netherlands, ³Faculty of Arts, University of Ljubljana, Slovenia, ⁴Institute of the Estonian Language, Estonia, ⁵Institute for Bulgarian Language, Bulgarian Academy of Sciences, Bulgaria, ⁶Hungarian Research Centre for Linguistics, Hungary



Manually-curated lexical-semantic resource combining corpora and sense inventories

Corpora and sense inventories

- Corpora
 - Origin: WikiMatrix
 - Sentences per language: 2,024
 - Identical sentences in 10 languages
 - Manually checked translations
 - Manually checked annotations
- Sense inventories
 - Dictionaries: Bulgarian, Estonian, Hungarian, Portuguese, Slovenian, Spanish
 - WordNets: Danish, Dutch, English, Italian

ELEXIS WSD dataset & UniDive

Relevance: **WG1** and **WG2**

Possible extension of the current dataset with additional languages and additional annotation layers:

- annotation of multiword expressions following the PARSEME annotation guidelines
- annotation of named entities
- syntactic parse structure following Universal Dependencies