

# UD Syntax for the ELEXIS-WSD Parallel Sense-Annotated Corpus: A Pilot Study



Carole Tiberius<sup>1</sup>, Jaka Čibej<sup>2</sup>, Jelena Kallas<sup>3</sup>, Kertu Saul<sup>3</sup>, Kadri Muischnek<sup>4</sup>, Simon Krek<sup>5</sup>

<sup>1</sup>Instituut voor de Nederlandse Taal, The Netherlands, <sup>2</sup>Faculty of Arts, University of Ljubljana, Slovenia, <sup>3</sup>Institute of the Estonian Language, Estonia, <sup>4</sup>University of Tartu, Estonia, <sup>5</sup>Jožef Stefan Institute, Slovenia

## What is the ELEXIS corpus?

An entirely manually-curated lexical-semantic resource available in ten European languages combining corpora and sense inventories (dictionaries and WordNets)

- Origin: WikiMatrix
- Number of sentences per language: 2,024
- Identical sentences in 10 languages
- Manually checked translations
- Manually checked annotations

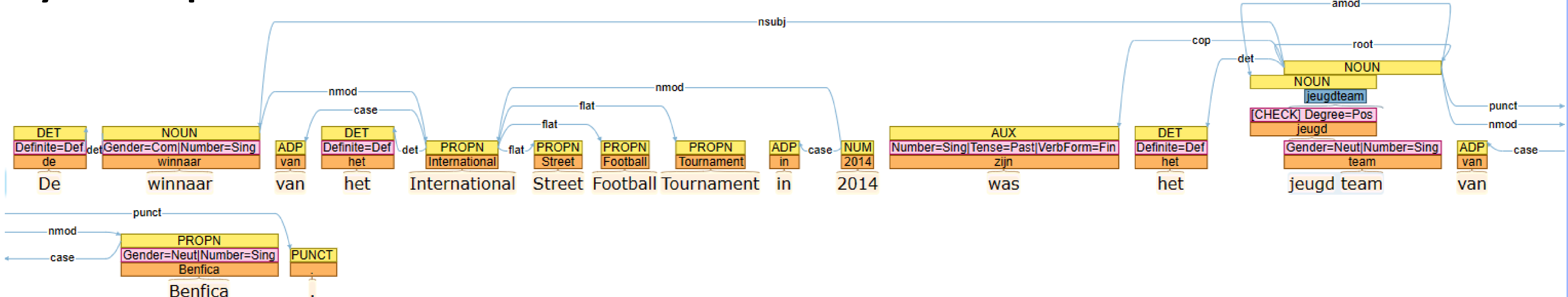
## ELEXIS corpus and Unidive

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

## Annotation in INCEpTION

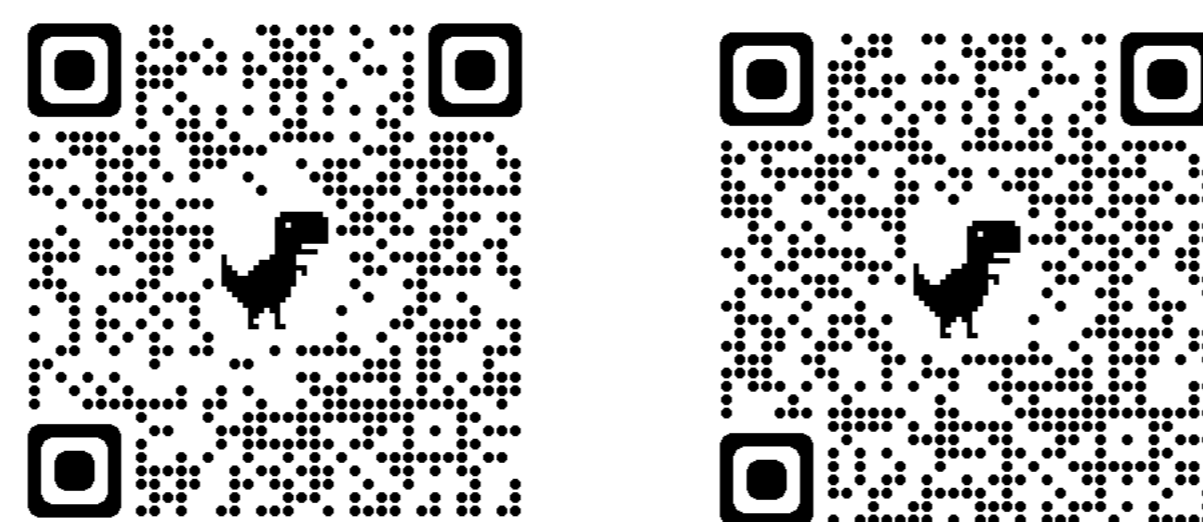
### Syntactic parse structure



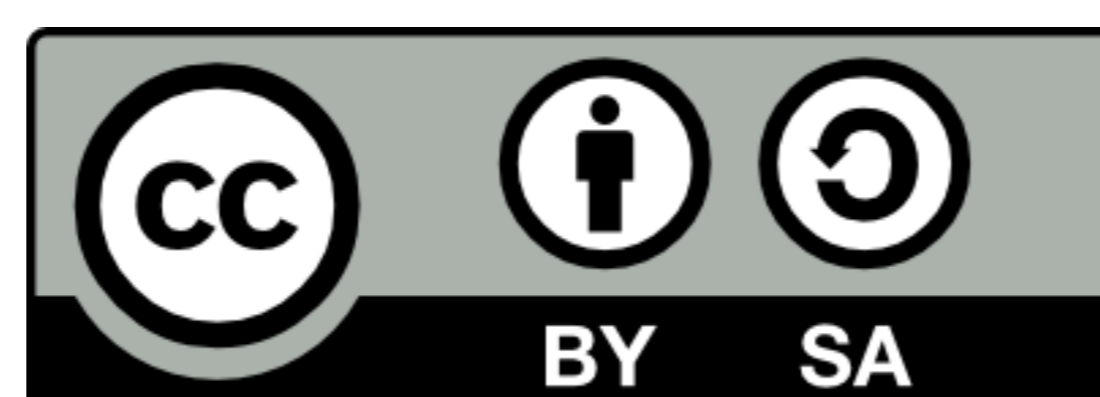
The winner of the International Street Football Tournament in 2014 was the Benfica junior team.

## Availability (repository & online concordancer)

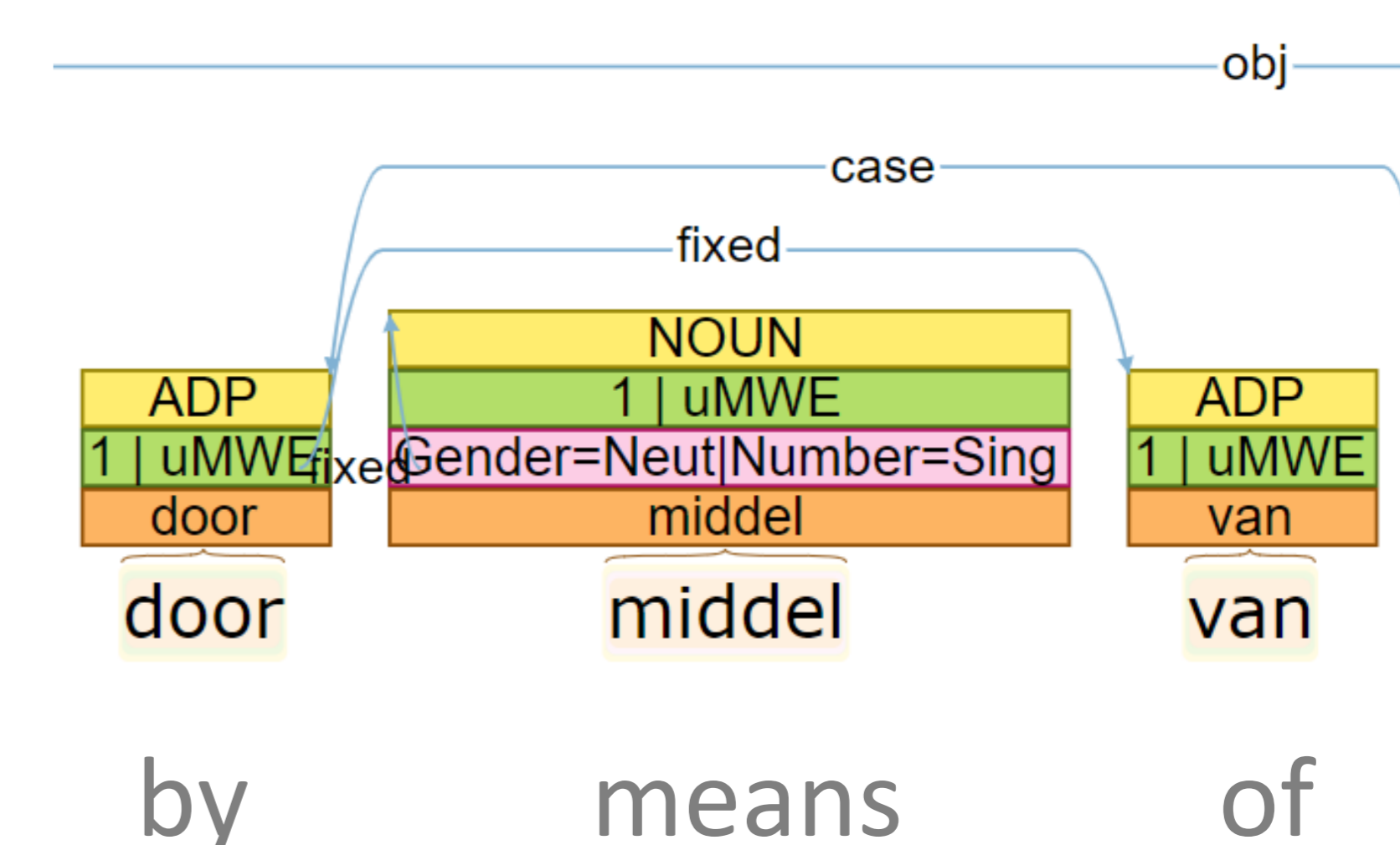
Martelli, Federico, et al. 2022. **Parallel sense-annotated corpus ELEXIS-WSD 1.0**, Slovenian language resource repository CLARIN.SI, ISSN 2820-4042.



- CC BY-SA includes the following elements:
- BY – Credit must be given to the creator
  - SA – Adaptations must be shared under the same terms



## Multiword expressions



Annotation layers: lemma (orange), morphosyntactic features (pink), POS-tags (yellow), MWEs (green), and UD-syntax (blue links).

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

## Design of the dataset

Martelli, Federico, et al. 2021. **Designing the ELEXIS Parallel Sense-Annotated Dataset in 10 European Languages**. In Electronic lexicography in the 21st century. Proceedings of the eLex 2021 conference. 377-395.

