# SimRelUz: Similarity and Relatedness scores as a Semantic Evaluation Dataset for Uzbek Language

**Ulugbek Salaev[1], Elmurod Kuriyozov[2], Carlos Gómez-Rodríguez[2]**

[1]Urgench State University, Uzbekistan (*ulugbek0302@gmail.com*)

[2]Universidade da Coruña, Spain (*{e.kuriyozov, carlos.gomez}@udc.es*)

## Semantic Analysis

The process of drawing language-independent meaning from text;
Answers following questions:

- What is the computational meaning of individual words/phrases in context? (Lexical Semantics);
- How can we learn semantic representations from data? (Distributional Semantics).

## Terminology

- **Semantic Similarity:**
  - Sense of relatedness that is dependent on the amount of shared properties (degree of synonymy)
  - **Example :** Bus - Train

- **Semantic Relatedness:**
  - General sense of semantic proximity or semantic association, regardless of the causes of the connection humans can perceive
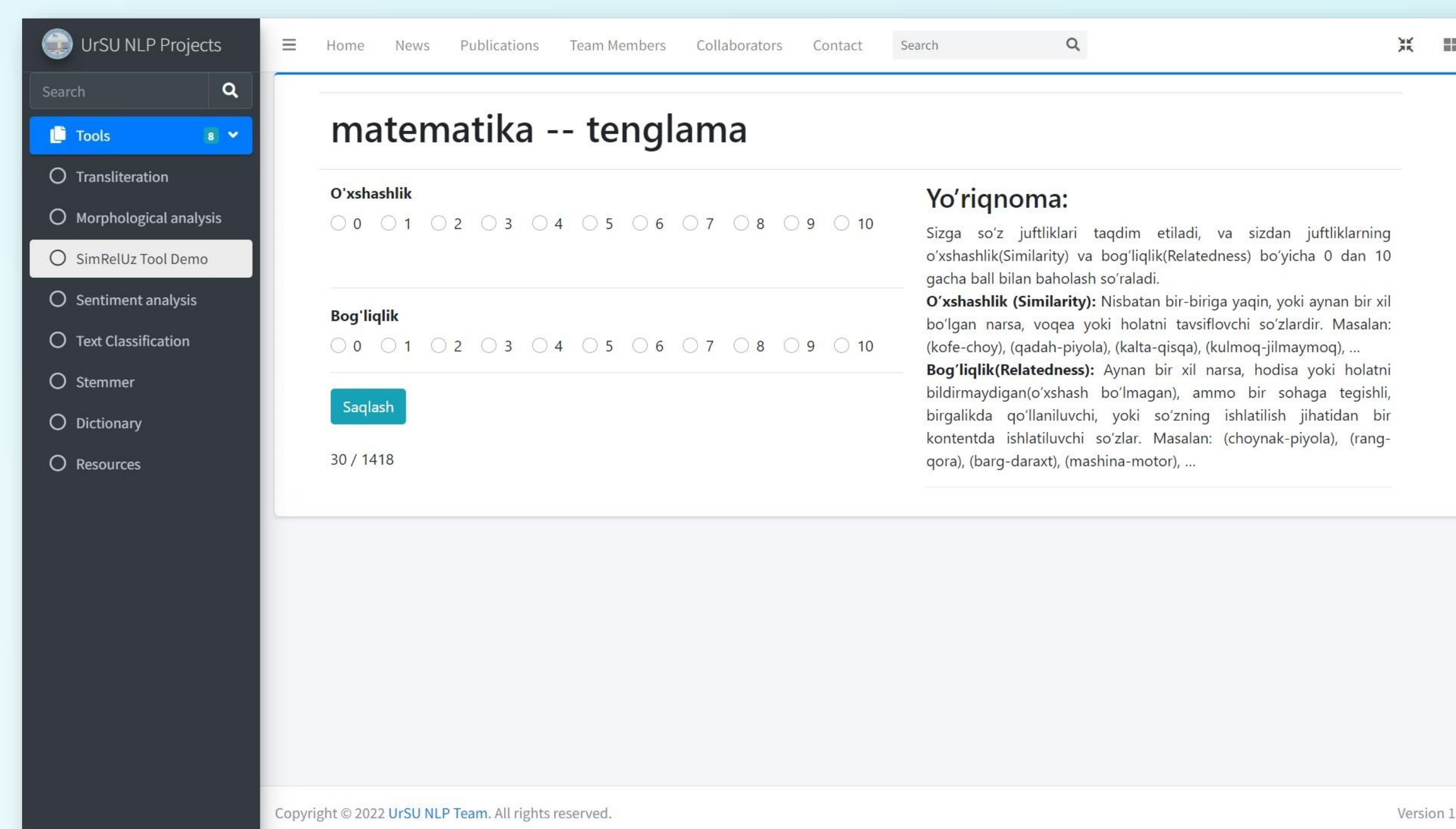  - **Example :** Coffee – Cup

## Uzbek Language

- **Official language of Uzbekistan;**
- **Native to All CA Countries, Russia, China;**
- **Spoken by more than 30 million people.**

- **Turkic Family**
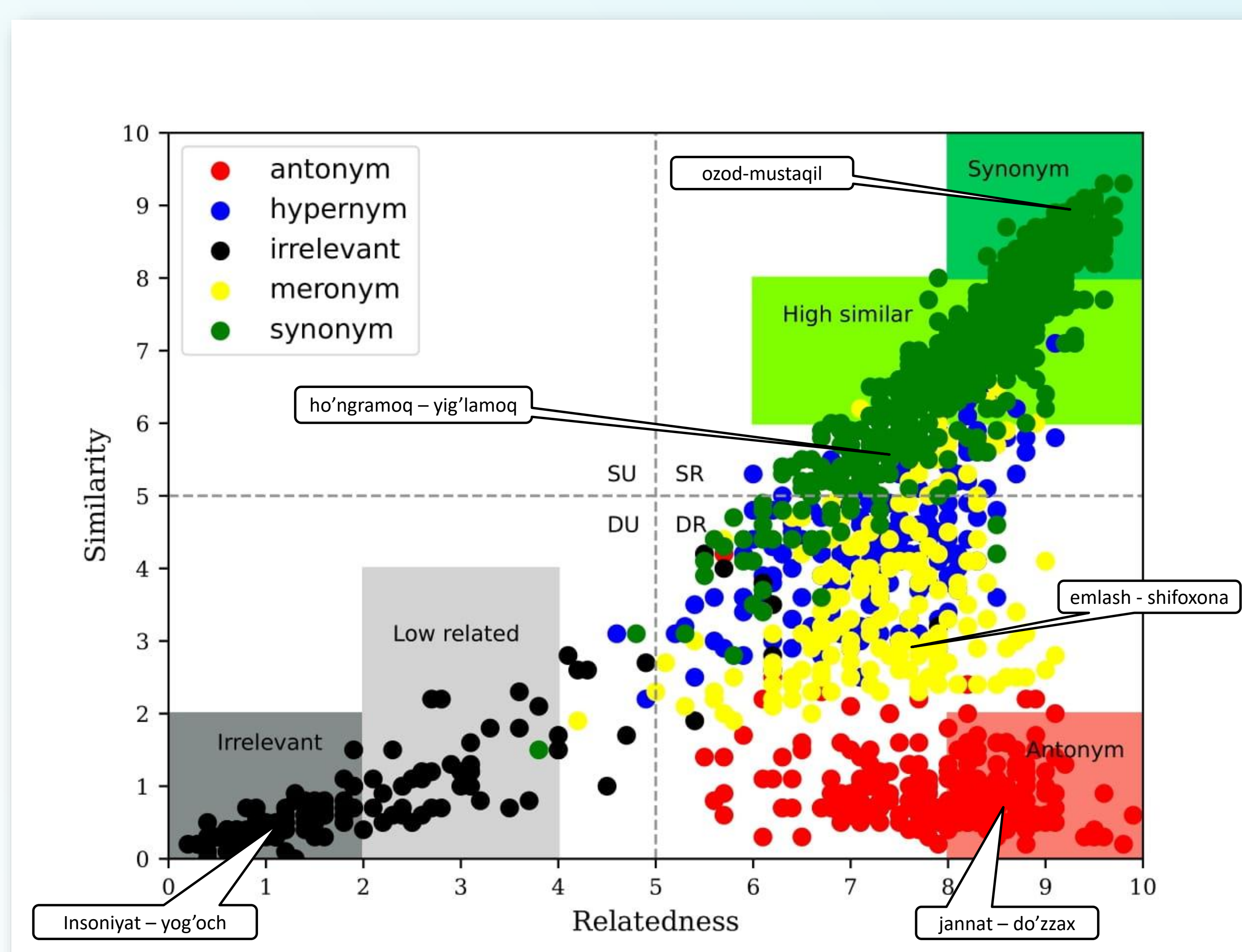- **Agglutinative**
- **Null Subject**
- **No Articles**
- **No Gender**
- **SOV**

## Semantic scores Annotation Tool



- **Open-source, web-based, multi-user tool;**
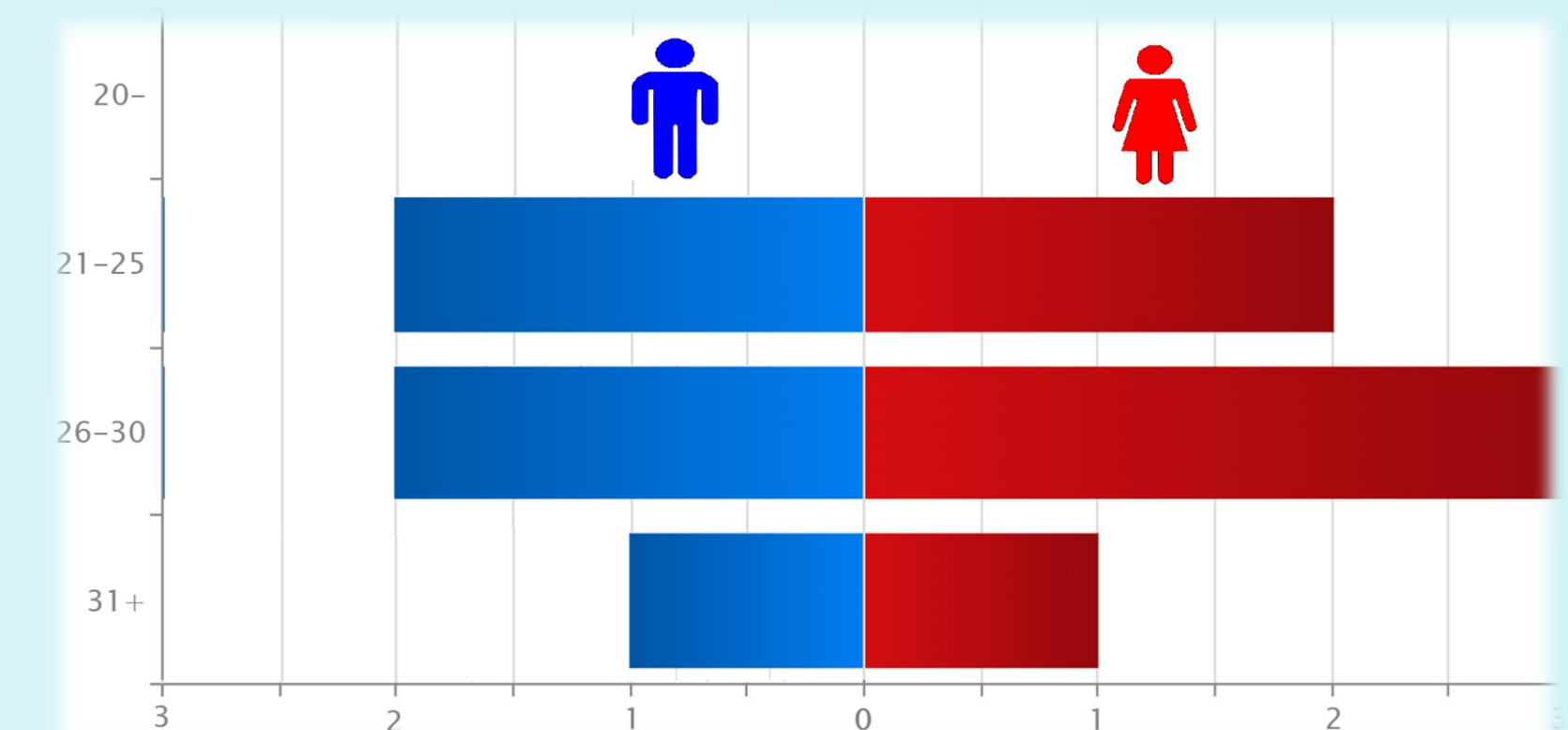- **Demo available:** *https://simrel.urdu.uz*
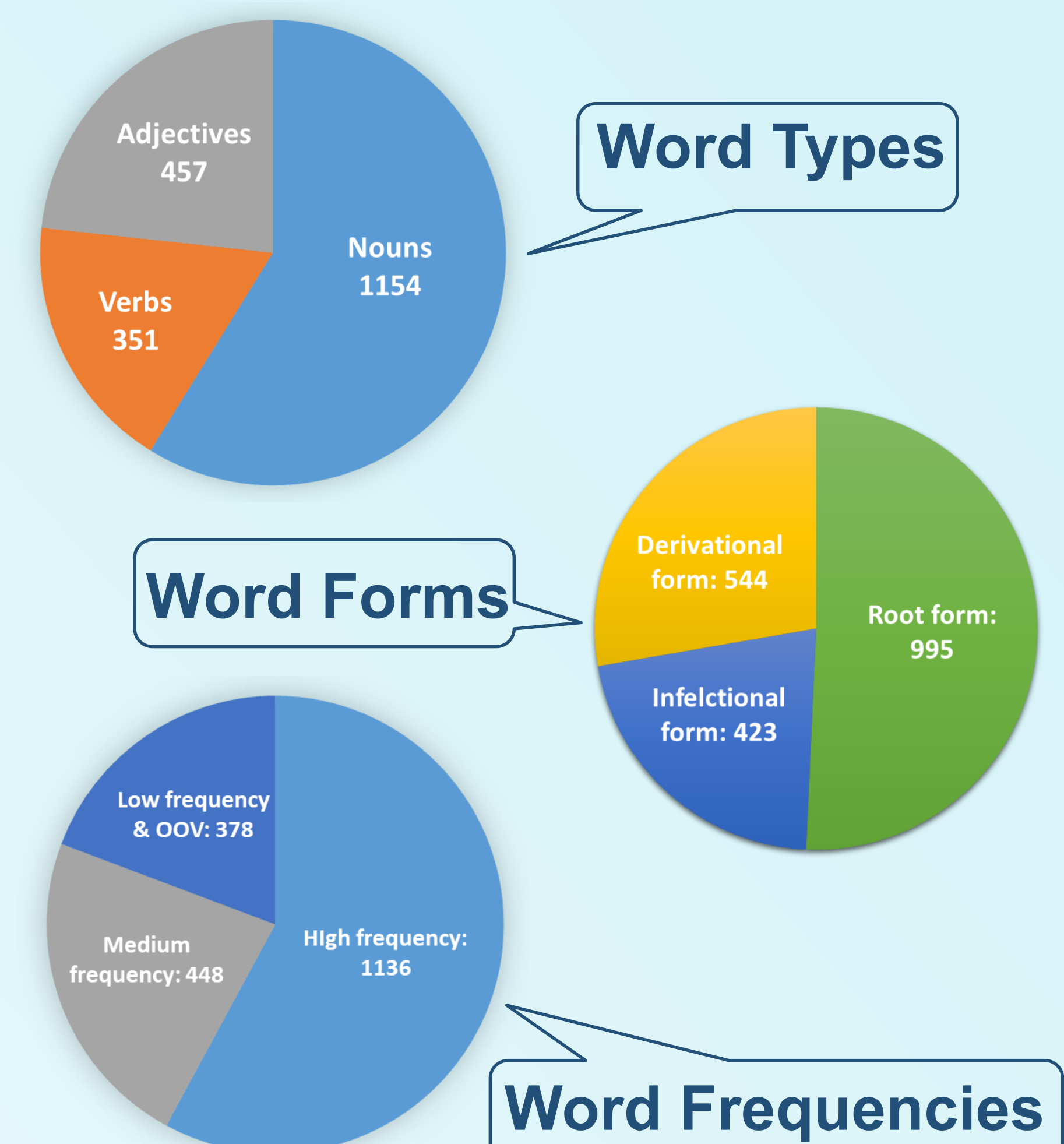
## Dataset Results in a 2D Visualization



- *Each dot represents a word-pair in the dataset;*
  - *X-axis: Relatedness scores, Y-axis: Similarity scores;*
  - *SU – Similar-Unrelated, SR – Similar-Related,*
  - *DU – Dissimilar-Unrelated, DR – Dissimilar-Related.*

## Annotation

- **Total number of annotators: 11**
  - High inter-annotator agreement score



## Results



**Word Types**

Adjectives 457
Verbs 351
Nouns 1154

**Word Forms**

Derivational form: 544
Inflectional form: 423
Root form: 995

**Word Frequencies**

Low frequency & OOV: 378
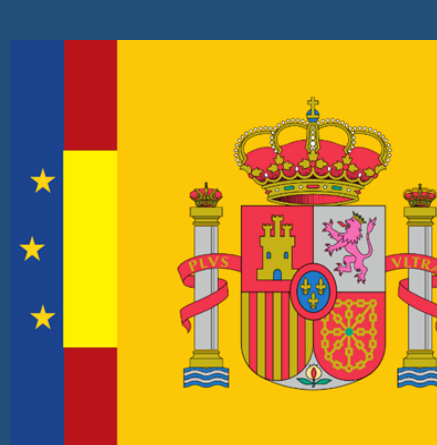Medium frequency: 448
High frequency: 1136

## Conclusion

In this project, we presented the SimRelUz, a novel semantic evaluation dataset for the Uzbek language that consists of similarity and relatedness scores for word-pairs.

We have also presented an open-source web-based annotation tool designed for multiple-user semantic scores annotation.

## Data

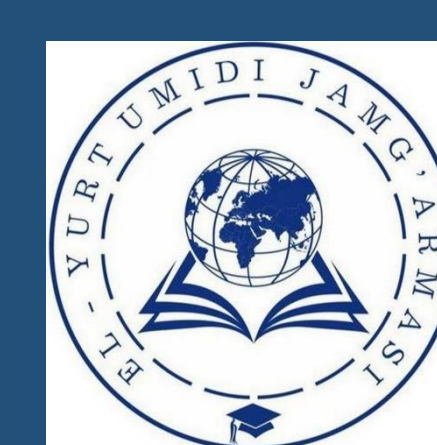The Dataset and the code for the Annotation tool are available: