# A Collaborative Web Tool for Linguistic Annotation

## Salih Furkan Akkurt

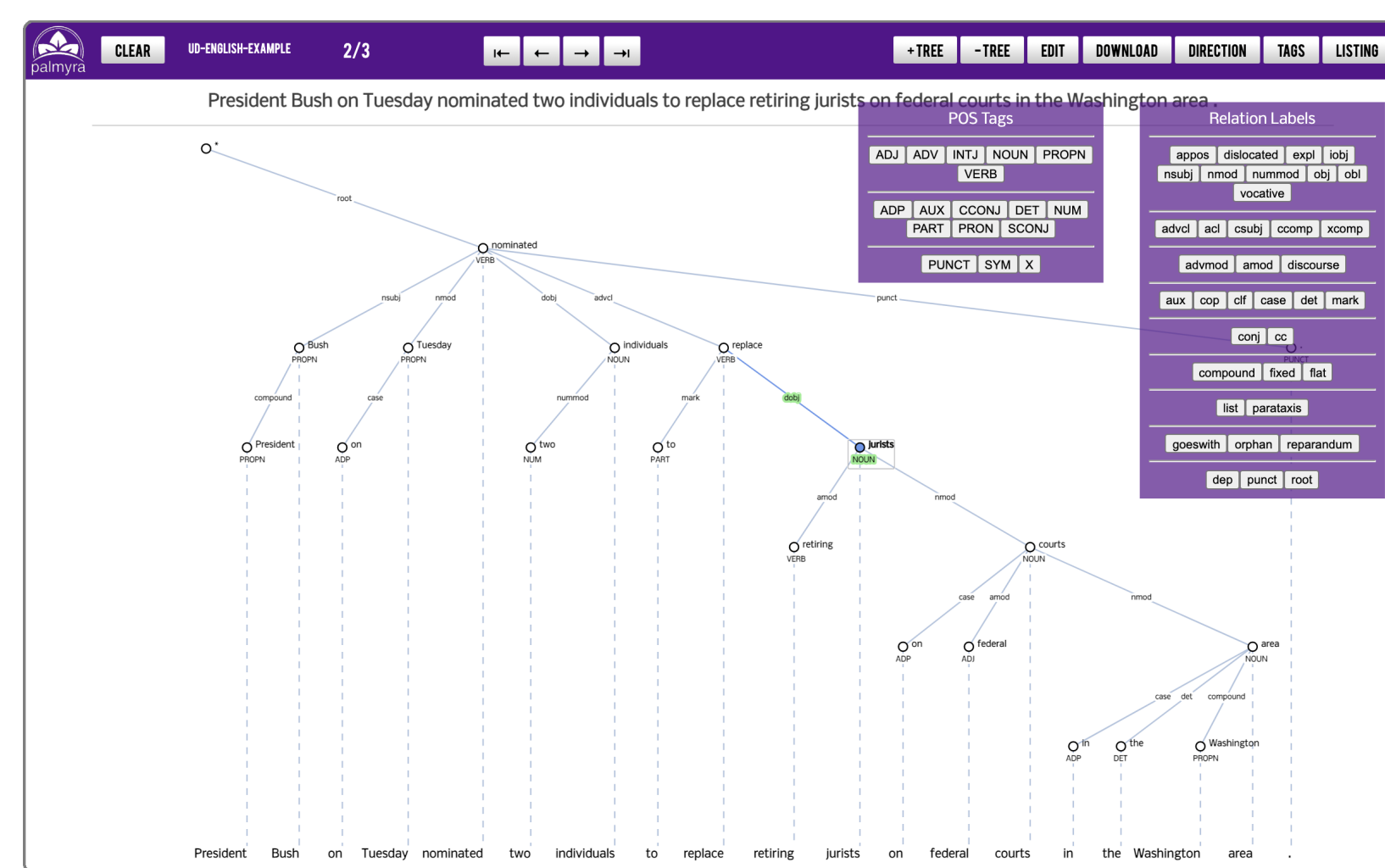Boğaziçi University - CMPE 492 Project - Advisor: Suzan Uskudarli

## Problem Statement

To design and develop a web-based tool that supports grammatical annotations of morphologically rich languages (MRLs), centers around annotator's user experience and is collaborative, open-source and accessible.

## Motivation

Large annotated datasets are crucial to NLP solutions. Creation of quality treebanks takes time and effort. Existing tools are not sufficient for annotating MRLs. Also, treebanks are annotated with multiple annotators, making collaborative tools necessary.

## Related Work



Drag-drop interface for annotation [1]

**Boğaziçi University Annotation Tool (BoAT):**
An annotation tool released in 2019 by the CMPE Department of Boğaziçi University. It was designed to annotate MRLs and has been used to annotate the BOUN Treebank, which has 9,761 sentences and is published on UD[2]. It's a standalone application.

## Requirements Elicitation & Management

Based on the experience with the BOUN Treebank and numerous meetings with an annotator with linguistic background, we decided to redesign BoAT with these 3 main themes: (1) Addressing annotator requests and improving existing functionality; (2) Extending features with a collaborative experience; (3) New architecture that supports persistence, accessibility and maintainability.

## Sentence Example for Agglutinative vs Analytical Languages

MRLs require extensive and complex annotations, illustrated with the following two examples:

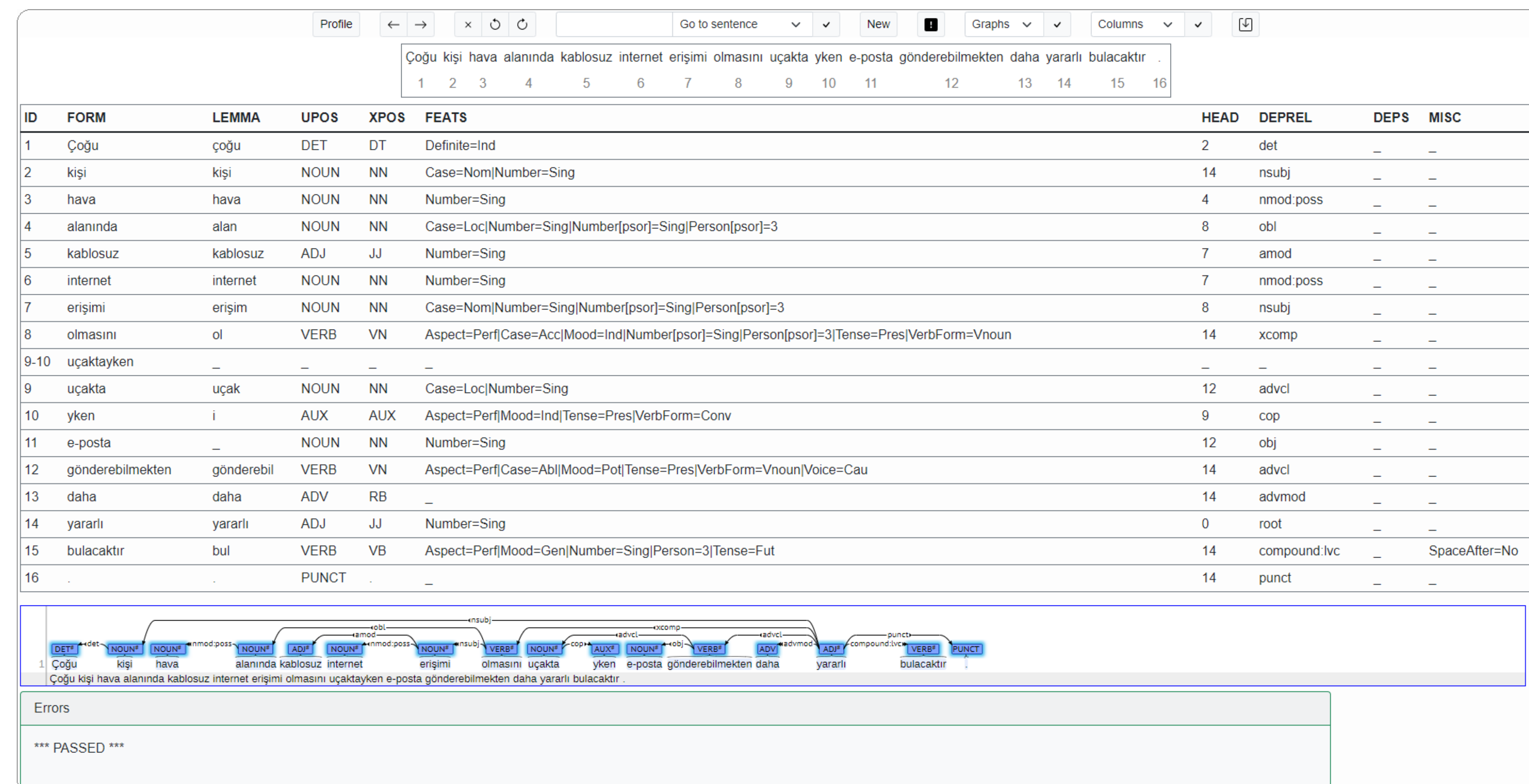English sentence: Most people would find airport wi-fi more useful than being able to send emails *on a plane*.[3]

| ID | FORM | LEMMA | UPOS | XPOS | FEATS | HEAD | DEPREL |
|----|------|-------|------|------|-------|------|--------|
| 15 | on | on | ADP | IN | _ | 18 | case |
| 16 | a | a | DET | DT | Definite=Ind\|PronType=Art | 18 | det |
| 17 | plane | plane | NOUN | NN | Number=Sing | 14 | obl |

Turkish sentence: Çoğu kişi hava alanında kablosuz internet erişimi olmasını *uçaktayken* e-posta gönderebilmekten daha yararlı bulacaktır.[4]

| ID | FORM | LEMMA | UPOS | XPOS | FEATS | HEAD | DEPREL |
|-----|-----------|-------|------|------|-------|------|--------|
| 9-10 | uçaktayken | _ | _ | _ | _ | _ | _ |
| 9 | uçakta | uçak | NOUN | NN | Case=Loc\|Number=Sing | 12 | advcl |
| 10 | yken | i | AUX | AUX | Aspect=Perf\|Mood=Ind\|Tense=Pres\|VerbForm=Conv | 9 | cop |

## Solution

We implemented a tool for MRLs like Turkish.



Annotation of "Çoğu kişi hava alanında kablosuz internet erişimi olmasını uçaktayken e-posta gönderebilmekten daha yararlı bulacaktır.", using BoAT v2



Results from searching the "form" field as "erişim" for the treebank *UD-Turkish-PUD*

## Result

Presented at The International Conference on Agglutinative Language Technologies as a Challenge of Natural Language Processing (ALTNLP) in Jun 2022 [5], which was well-received, resulting in several downloads.
Currently, acceptance testing is being performed with annotators with linguistic background. Preliminary results are promising.
A dockerized version is available on GitHub [6].

## Future Directions

- Endangered languages in the Linguistics Department of Boğaziçi University.
- Uzbek: Offer came to collaborate during ALTNLP 2022.
- Automated partial annotations.

## Acknowledgements

I am grateful to the NLP team of folks from both CMPE and LING departments. A special thanks is due to Büşra Marşan who gave constant feedback and taught me about the annotation process. I am grateful, most of all, to my advisor Ms Uskudarli for spending countless hours with and for me. 🤗

## References

[1] https://universaldependencies.org/tools.html#palmyra

[2] https://universaldependencies.org/treebanks/tr_boun/index.html

[3] https://universaldependencies.org/treebanks/en_pud/index.html

[4] https://universaldependencies.org/treebanks/tr_pud/index.html

[5] Akkurt, S. F., Marşan B., Uskudarli S. (2022) 'BoAT v2 - A Web-Based Dependency Annotation Tool with Focus on Agglutinative Languages', *The International Conference on Agglutinative Language Technologies as a Challenge of Natural Language Processing.* Slovenia, 7-8 June.

[6] https://github.com/furkanakkurt1335/boat

Contact: furkan.akkurt@boun.edu.tr