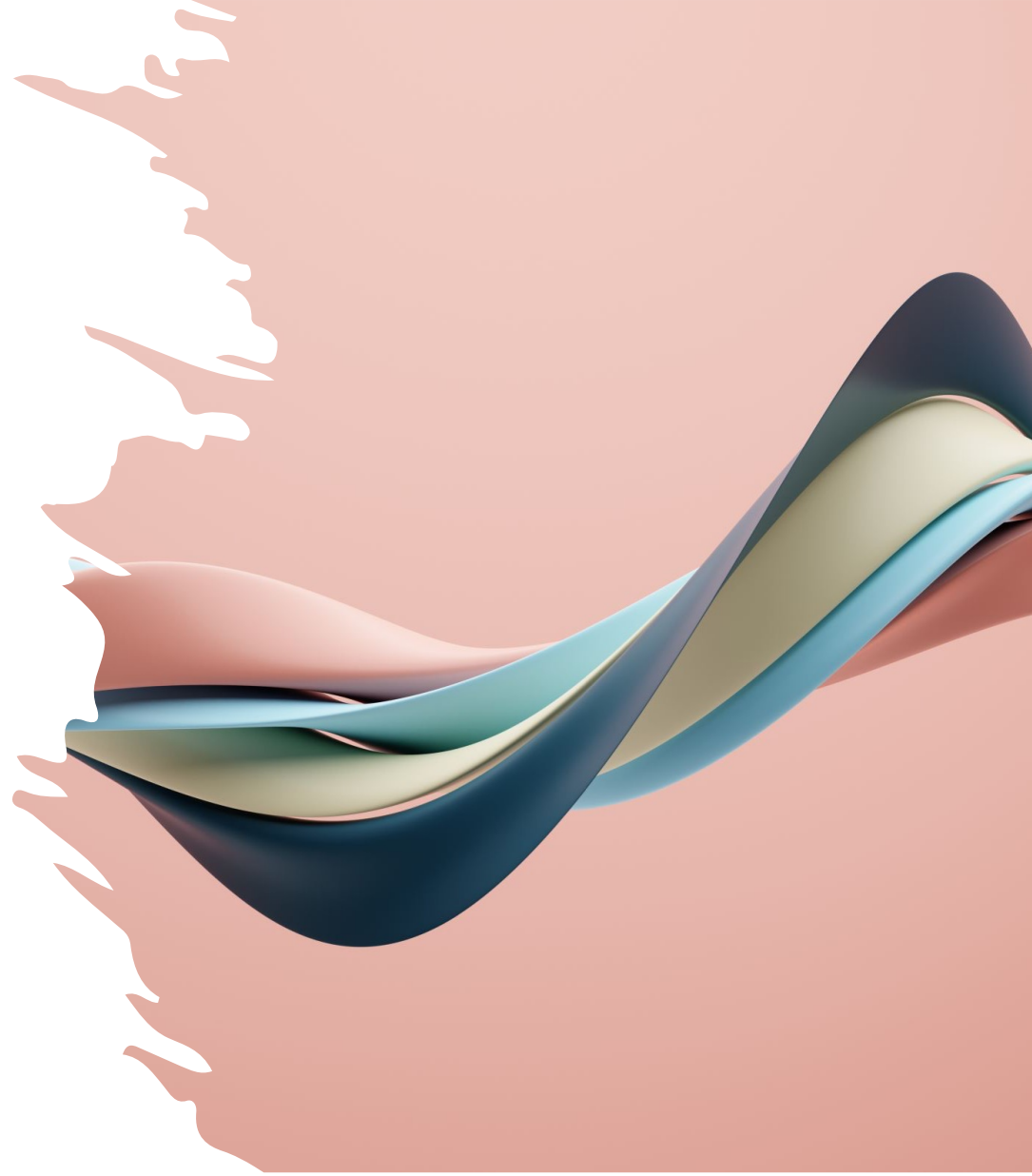


Journey of TreeBanking

STARLANG YAZILIM
DANIŞMANLIK



Outline

- Starlang TreeBanks
- Annotation Scheme
- Annotation Interface
- Future Work

TreeBanks

- Atis (Tr, En)
- KeNet
- Penn
- Tourism
- FrameNet



Atis TreeBank (2.9)

- Airline Travel Information System
- Parallel TreeBank: Atis English
- Mainly question and imperative sentences
- Small vocabulary, English and Turkish mixed
- 5.432 Sentences (4274 train, 586 test, 572 dev)
- 45.875 Tokens (61.879 Tokens in English)
- Contributors: Mehmet Köse⁴, Olcay Taner Yıldız²

KeNet TreeBank (2.8)

- Turkish WordNet KeNet
- Example sentences for synsets. Sentences are mainly taken from Turkish Language Association.
- Complex vocabulary
- 18.687 Sentences (9.345 train, 9.342 test)
- 178.658 Tokens
- Contributors: Aslı Kuzgun¹, Neslihan Cesur³, Olcay Taner Yıldız², Oğuzhan Kuyrukçu⁵, Arife Betül Yenice⁶, Bilge Nas Arıcan⁷, Ezgi Sanıyar⁸

Penn TreeBank (2.8)

- Largest Turkish TreeBank
- Translated sentences from subsample of Penn Constituency TreeBank
- Parallel TreeBank: Penn English
- Sentences of length < 20 words in English
- Complex vocabulary, English and Turkish mixed
- 16.396 Sentences
- 183.555 Tokens
- Contributors: Neslihan Cesur³, Aslı Kuzgun¹, Olcay Taner Yıldız², Büşra Marşan⁹, Neslihan Kara¹⁰, Bilge Nas Arıcan⁷, Merve Özçelik¹¹, Deniz Baran Aslan¹²

Tourism TreeBank (2.8)

- Customer reviews of a tourism company
- Small, sometimes irregular sentences, i.e. multiple sentences in one sentence, nominal sentences, etc.
- Small vocabulary, domain dataset.
- 19.833 Sentences
- 91.469 Tokens
- Contributors: Aslı Kuzgun¹, Neslihan Cesur³, Olcay Taner Yıldız², Oğuzhan Kuyrukçu⁵, Büşra Marşan⁹, Bilge Nas Arıcan⁷, Neslihan Kara¹⁰, Deniz Baran Aslan¹², Ezgi Sanıyar⁸, Cengiz Asmazoğlu¹³

FrameNet (2.8)

- Turkish FrameNet
- Example constructed sentences taken from Turkish FrameNet
- Complex vocabulary, basically one sentence per verb
- 2.698 Sentences
- 19.221 Tokens
- Contributors: Neslihan Cesur³, Aslı Kuzgun¹, Olcay Taner Yıldız², Büşra Marşan⁹, Oğuzhan Kuyrukçu⁵, Bilge Nas Arıcan⁷, Ezgi Sanıyar⁸, Neslihan Kara¹⁰, Merve Özçelik¹¹

Annotation Scheme

- Morphological Disambiguation Process
- Dependency Annotation Process
- Backed by NlpToolkit Library (8 Programming Languages, over 1 million lines of code)
 - MorphologicalAnalysis, MorphologicalDisambiguation
 - WordNet
 - DependencyParser, UniversalDependencyParser
 - AnnotatedSentence
 - DataCollector

Features

- Layered architecture: surface form, morphological analysis, semantics, ...
- Store intermediate representation, namely morphological analysis
- After morphological disambiguation, root, universal pos tag determined.
- Morphological analyzer is connected to Turkish WordNet, each root form is taken from Turkish WordNet (has a sense in WordNet)
- Universal features are extracted from morphological analysis.

Morphological Disambiguation

The screenshot shows a software interface for morphological disambiguation. At the top, there is a toolbar with six navigation buttons (back, forward, etc.) and a checkbox labeled "Auto Morphological Disambiguation". Below the toolbar, the file name "0006.train" is displayed. The main area shows a list of words and their morphological tags:

| Word | Morphological Tag |
|---------|---|
| Çocuk | çocuk+NOUN+A3SG+PNON+NOM |
| kalem | kale+NOUN+A3SG+P1SG+ACC kalem+NOUN+A3SG+P3SG+NOM kalem+NOUN+A3SG+PNON+ACC |
| masaya | masa+NOUN+A3SG+PNON+DAT |
| biraktı | birak+VERB+POS+PAST+A3SG |
| . | .+PUNC |

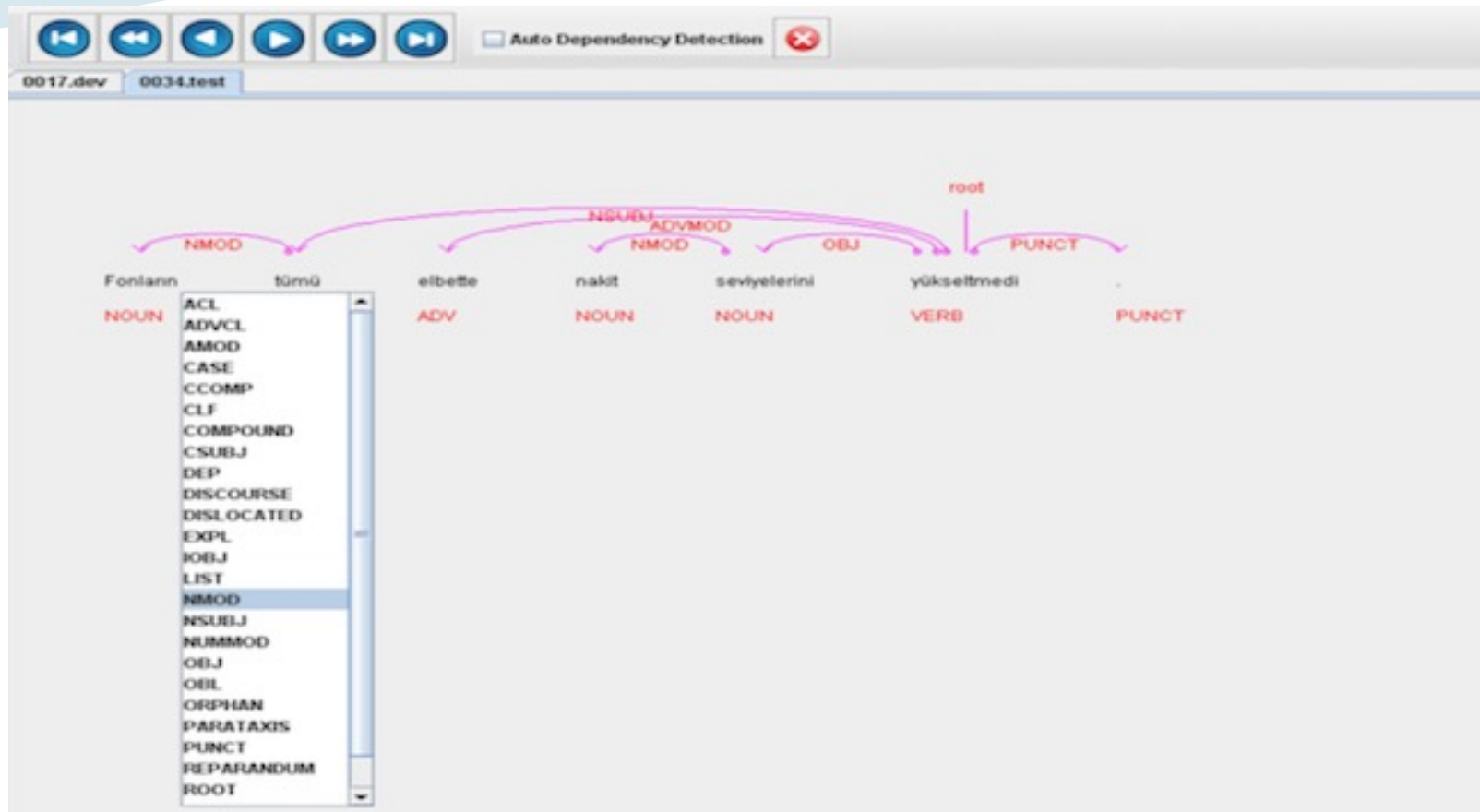
Editing Words

The screenshot shows a software interface for editing words. At the top, there is a "Project" window with a toolbar containing six navigation buttons (back, forward, etc.) and a checkbox for "Auto Dependency Detection" which is currently unchecked. Below the toolbar, there are two tabs: "0017.dev" and "0030.test". The main area displays a dependency diagram for the sentence "İnsanlar paniklemiyor .". The word "İnsanlar" is labeled as "NOUN". The word "paniklemiyor" is labeled as "PUNCT" and is highlighted with a blue box. The word "." is labeled as "PUNCT". A vertical line labeled "root" connects the "paniklemiyor" node to a central point. Two curved lines labeled "NSUBJ" and "PUNCT" connect the "İnsanlar" node to the "paniklemiyor" node and the "paniklemiyor" node to the "." node, respectively.

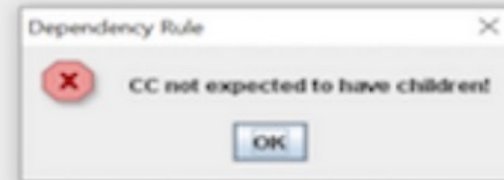
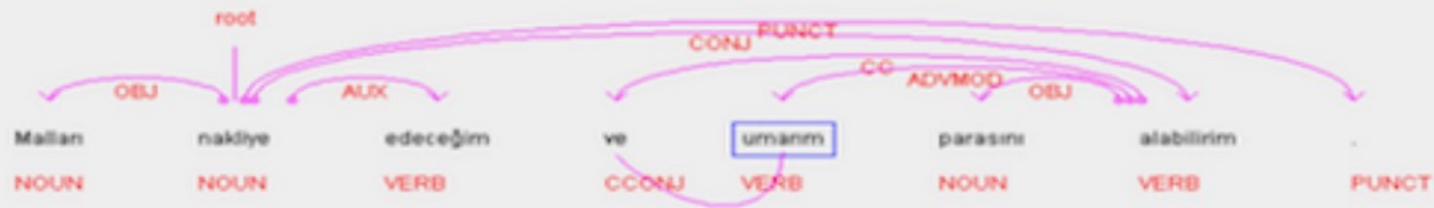
Annotation Comparison

| FileName | Index | Word | Morphological Analysis | Sentence |
|------------|-------|------------|--|---|
| 0000.dev | 1 | Devasa | devasa + ADJ | Devasa ölçekli yeni kanunda kullanılan karmaşık ve çetrefilli di kav... |
| 0000.dev | 2 | ölçekli | ölçek + NOUN + A3SG + PNON + NOM*DB + ADJ + WITH | Devasa ölçekli yeni kanunda kullanılan karmaşık ve çetrefilli di kav... |
| 0000.dev | 3 | yeni | yeni + ADJ | Devasa ölçekli yeni kanunda kullanılan karmaşık ve çetrefilli di kav... |
| 0000.dev | 4 | kanunda | kanun + NOUN + A3SG + PNON + LOC | Devasa ölçekli yeni kanunda kullanılan karmaşık ve çetrefilli di kav... |
| 0000.dev | 5 | kullanılan | kullan + VERB*DB + VERB + PASS + POS*DB + ADJ + PRESPT | Devasa ölçekli yeni kanunda kullanılan karmaşık ve çetrefilli di kav... |
| 0000.dev | 6 | karmaşık | karmaşık + ADJ | Devasa ölçekli yeni kanunda kullanılan karmaşık ve çetrefilli di kav... |
| 0000.dev | 7 | ve | ve + CONJ | Devasa ölçekli yeni kanunda kullanılan karmaşık ve çetrefilli di kav... |
| 0000.dev | 8 | çetrefilli | çetrefil + NOUN + A3SG + PNON + NOM*DB + ADJ + WITH | Devasa ölçekli yeni kanunda kullanılan karmaşık ve çetrefilli di kav... |
| 0000.dev | 9 | di | di + NOUN + A3SG + PNON + NOM | Devasa ölçekli yeni kanunda kullanılan karmaşık ve çetrefilli di kav... |
| 0000.dev | 10 | kavga | kavga + NOUN + A3SG + PNON + ACC | Devasa ölçekli yeni kanunda kullanılan karmaşık ve çetrefilli di kav... |
| 0000.dev | 11 | bulandı | bulan + VERB*DB + VERB + CAUS + POS + PAST + A3SG | Devasa ölçekli yeni kanunda kullanılan karmaşık ve çetrefilli di kav... |
| 0000.dev | 12 | - | + PUNC | Devasa ölçekli yeni kanunda kullanılan karmaşık ve çetrefilli di kav... |
| 0000.test | 1 | Hayır | hayır + ADV | Hayır , kara pazartesi değildi . |
| 0000.test | 2 | - | + PUNC | Hayır , kara pazartesi değildi . |
| 0000.test | 3 | kara | kara + ADJ | Hayır , kara pazartesi değildi . |
| 0000.test | 4 | pazartesi | pazartesi + NOUN + A3SG + PNON + NOM | Hayır , kara pazartesi değildi . |
| 0000.test | 5 | değildi | değil + ADJ*DB + VERB + ZERO + PAST + A3SG | Hayır , kara pazartesi değildi . |
| 0000.test | 6 | - | + PUNC | Hayır , kara pazartesi değildi . |
| 0000.train | 1 | Bayan | bayan + NOUN + A3SG + PNON + NOM | Bayan Haag Elanti oynamıyor . |
| 0000.train | 2 | Haag | haag + NOUN + PROP + A3SG + PNON + NOM | Bayan Haag Elanti oynamıyor . |
| 0000.train | 3 | Elanti | elanti + NOUN + PROP + A3SG + PNON + NOM | Bayan Haag Elanti oynamıyor . |
| 0000.train | 4 | oynuyor | oyna + VERB + POS + PROC1 + A3SG | Bayan Haag Elanti oynamıyor . |
| 0000.train | 5 | - | + PUNC | Bayan Haag Elanti oynamıyor . |

Dependency Annotation



Error Check



I'm going to ship and hope I get paid .

code 3:SHOULDNT_BE_OF_POS <--> AUX should not be VERB

code 5:SHOULDNT_BE_OF_POS <--> ADVMOD should not be VERB

Annotation Examples



Annotation Comparison

| FileName | Index | Word | Depend. | Dependency Type | Sentence |
|-----------|-------|------------|---------|-----------------|--|
| 0000 dev | 1 | Devasa | 2 | AMOO | Devasa ölçekđ yeni kanunda kulanılan kamađık ve çetrefilli di kavğay bulandırđ . |
| 0000 dev | 2 | ölçekđ | 4 | AMOO | Devasa ölçekđ yeni kanunda kulanılan kamađık ve çetrefilli di kavğay bulandırđ . |
| 0000 dev | 3 | yeni | 4 | AMOO | Devasa ölçekđ yeni kanunda kulanılan kamađık ve çetrefilli di kavğay bulandırđ . |
| 0000 dev | 4 | kanunda | 5 | OBL | Devasa ölçekđ yeni kanunda kulanılan kamađık ve çetrefilli di kavğay bulandırđ . |
| 0000 dev | 5 | kulanılan | 9 | ACL | Devasa ölçekđ yeni kanunda kulanılan kamađık ve çetrefilli di kavğay bulandırđ . |
| 0000 dev | 6 | kamađık | 9 | AMOO | Devasa ölçekđ yeni kanunda kulanılan kamađık ve çetrefilli di kavğay bulandırđ . |
| 0000 dev | 7 | ve | 8 | CC | Devasa ölçekđ yeni kanunda kulanılan kamađık ve çetrefilli di kavğay bulandırđ . |
| 0000 dev | 8 | çetrefilli | 6 | CONJ | Devasa ölçekđ yeni kanunda kulanılan kamađık ve çetrefilli di kavğay bulandırđ . |
| 0000 dev | 9 | di | 11 | NSUBJ | Devasa ölçekđ yeni kanunda kulanılan kamađık ve çetrefilli di kavğay bulandırđ . |
| 0000 dev | 10 | kavğay | 11 | OBJ | Devasa ölçekđ yeni kanunda kulanılan kamađık ve çetrefilli di kavğay bulandırđ . |
| 0000 dev | 11 | bulandırđ | 0 | ROOT | Devasa ölçekđ yeni kanunda kulanılan kamađık ve çetrefilli di kavğay bulandırđ . |
| 0000 dev | 12 | . | 11 | PUNCT | Devasa ölçekđ yeni kanunda kulanılan kamađık ve çetrefilli di kavğay bulandırđ . |
| 0000 test | 1 | Hayr | 4 | DISCOURSE | Hayr , kara pazarfesi deđildi . |
| 0000 test | 2 | . | 4 | PUNCT | Hayr , kara pazarfesi deđildi . |
| 0000 test | 3 | kara | 4 | AMOO | Hayr , kara pazarfesi deđildi . |
| 0000 test | 4 | pazarfesi | 0 | ROOT | Hayr , kara pazarfesi deđildi . |
| 0000 test | 5 | deđildi | 4 | ALX | Hayr , kara pazarfesi deđildi . |
| 0000 test | 6 | . | 4 | PUNCT | Hayr , kara pazarfesi deđildi . |

Planned TreeBanks (2.13 or 2.14)

- QuestionBank
 - 10.000 Questions taken from SQUAD data set.
 - Parallel TreeBank with English also annotated.
 - Will be largest Question Treebank in UD.
 - 73.579 tokens in Turkish, 99.834 tokens in English
- Blogs
 - 11.128 Blog Sentences in 28 Topics (About 400 sentences per topic)
 - 141.567 Tokens

Future Work

- Extending dependency with extra layers:
 - Semantic layer
 - PropBank layer
 - FrameNet layer
 - NER layer
 - Shallow Parse layer
- Conversion to other representations
 - AMR (Abstract Meaning Representation)
 - Constituency treebanks

A light blue brushstroke graphic that tapers from left to right, serving as a background for the text.

Questions?