# Tutorial on Universal Dependencies

**Adding a new language to UD**

Joakim Nivre[1]    Daniel Zeman[2]    Filip Ginter[3]    **Francis M. Tyers**[45]

[1]Department of Linguistics and Philology, Uppsala University, Sweden

[2]Institute of Formal and Applied Linguistics, Charles University, Prague, Czechia

[3]Department of Information Technology, University of Turku, Finland

[4]Giela ja kultuvrra instituhtta, UiT Norgga árktalaš universitehta, Tromsø, Norway

[5]Arvutiteaduse instituut, Tartu Ülikool, Estonia

You want your language in UD

↙                           ↘

Existing treebank        No existing treebank
You have permission      No permission/licence

↓                           ↓

**Treebank conversion**   **Building from scratch**
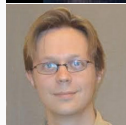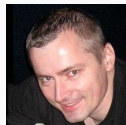
**First steps**

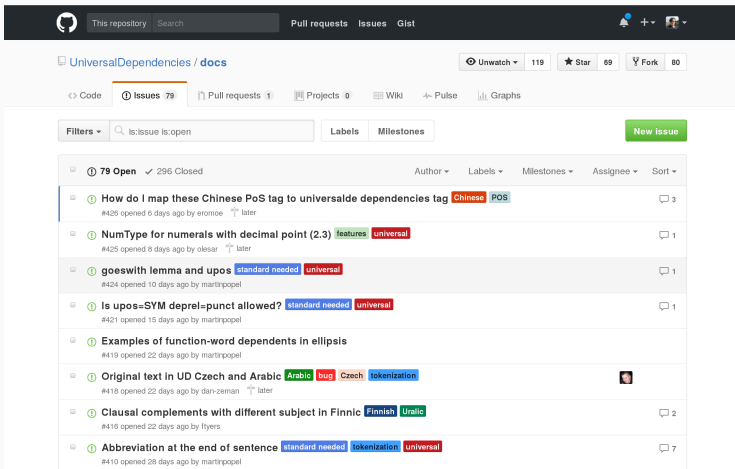- Get an account in Github
  - All development goes on here

**Get in contact**

- Ask someone from the release team to set up a module

- Get in contact with any other teams working on your language, or a related one

- Register for the mailing list *

* `http://stp.lingfil.uu.se/mailman/listinfo/ud`

**Release team**

# Linguistic Discussion

Linguistic discussion goes on under the *docs* module

# Linguistic Discussion

Annotation guidelines are discussed with examples

# Linguistic Discussion

## Annotation guidelines are discussed with examples

**MemduhG** commented on 30 Aug 2016    `Member`  + 😊  ✏

We came across a sentence the annotation of which seems to be tricky, with regards to subject, apposition and their relation.

*Hingê ez û xweha xwe Cûlya; em her du jî du-salî bûn*
So I and my sister Julia; we both were two-year-olds.

Either *ez û xweha xwe Cûlya* "I and my sister Julia" or *em her du* "we both" should be the subject, with the other one attached as `appos`. I am not sure which one should be the subject, as the part that would usually be considered an apposition is given before rather than after what would be the subject.

dan-zeman commented on 30 Aug 2016    `Member`  + 😊  ✏  ✕

This looks like a perfect example of the `dislocated` relation, http://universaldependencies.org/u/dep/dislocated.html

I would make *em her du* the subject and *ez û xweha xwe Cûlya* a dislocated dependent, both attached to the main predicate.

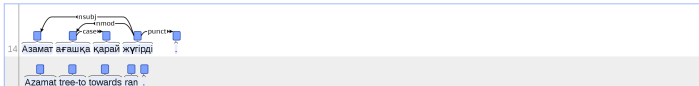🏷  🧑 **dan-zeman** added `dependencies` `question` labels on 30 Aug 2016

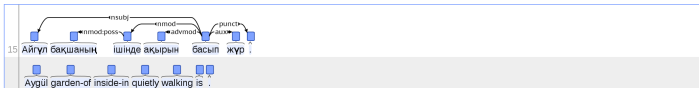✝  🧑 **dan-zeman** added this to the **lg-specific v2** milestone on 30 Aug 2016

3

# Write Documentation



- Documentation is written in Markdown and converted to HTML
- Not mandatory, but highly recommended
- Document as you write conversion rules/the annotation scheme

4

# Language-Family Documentation

## Slavic

- Introduction
- Tokenization
- Morphology
  - General principles
  - POS tags (single document)
  - Features (single document)
- Syntax
  - General principles
  - Specific constructions
  - Relations (single document)

## Quick links

- Pronominal words

## UD Treebanks

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| ▸ | Belarusian | 6K | | - | | | | |
| ▸ | Bulgarian | 140K | | | | ✔ | | |
| ▸ | Croatian | 183K | | - | | ✔ | | |
| ▸ | Czech | 1,330K | | | | ✔ | | |
| ▸ | Czech-CAC | 482K | | | | ✔ | | |
| ▸ | Czech-CLTT | 26K | | | | ✔ | | |
| ▸ | Old Church Slavonic | 47K | | - | | ✔ | | |
| ▸ | Polish | 72K | | - | | ✔ | | |
| ▸ | Russian | 87K | | | | ✔ | | |
| ▸ | Russian-SynTagRus | 988K | | | | ✔ | | |
| ▸ | Serbian | - | | - | | | | |
| ▸ | Slovak | 93K | | | | ✔ | | |
| ▸ | Slovenian | 126K | | | | ✔ | | |
| ▸ | Slovenian-SST | 19K | | | | ✔ | | |
| ▸ | Ukrainian | 12K | | | | ✔ | | |

**NEW!** Documentation by language family

# Treebank conversion

**Case study: Turkish**

**METU-Sabancı Treebank**

- Started in 2003
- Converted to CoNLL format for the 2006 shared task

**İTÜ-METU-Sabancı Treebank**

- 2016 reannotation of the METU-Sabancı treebank
- Morphology editted, dependencies from scratch
- 60k tokens in 2 months with 5 annotators (Sulubacak et al., 2016)

A loosely co-ordinated effort between:

- Çağrı Çöltekin (U. Tübingen)
- A team from İTÜ
    - Umut Sulubacak
    - Memduh Gökırmak
    - Gülşen Eryiğit
- Hüner Kaşıkara (U. Boğaziçi)
- Joakim Nivre
- Francis Tyers

Kickoff meeting in Uppsala (November, 2015)

**Method:**

- Go through reference grammar (Göksel & Kerslake, 2011)
- Document phenomena
- Convert treebank according to documentation

**Tools:**

- 6,000 lines of Java
  - Morphological synthesis
  - Collapse derivations
  - Remove multiwords
  - Distinguish clause from non-clause

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 1 | _ | yanlış | Adj | _ | _ | 2 | DERIV |
| 2 | _ | _ | Verb | _ | Acquire\|Pos | 3 | DERIV |
| 3 | Yanlışlanan | _ | _ | APresPart | _ | 4 | MODIFIER |
| 4 | kuramın | kuram | Noun | _ | A3sg\|Pnon\|Gen | 6 | POSSESSOR |
| 5 | _ | doğurgan | Adj | _ | _ | 6 | DERIV |
| 6 | doğurganlığı | _ | Noun | _ | Ness\|A3sg\|P3sg\|Nom | 8 | SUBJECT |
| 7 | burada | bura | Noun | _ | A3sg\|Pnon\|Loc | 8 | LOCATIVE |
| 8 | yatar | yat | Verb | _ | Pos\|Aor\|A3sg | 9 | SENTENCE |
| 9 | . | . | Punc | _ | _ | 0 | ROOT |

9

| 1 | Yanlışlanan | Yanlışlan | VERB | _ | Tense=Pres\|VerbForm=Part\|... | 2 | acl |
| 2 | kuramın | kuram | NOUN | _ | Case=Gen\|... | 3 | nmod:poss |
| 3 | doğurganlığı | doğurganlık | NOUN | _ | Case=Nom\|... | 5 | nsubj |
| 4 | burada | bura | NOUN | _ | Case=Loc\|... | 5 | obl |
| 5 | yatar | yat | VERB | _ | Tense=Aor\|VerbForm=Fin\|... | 0 | root |
| 6 | . | . | PUNCT | _ | _ | 5 | punct |

# From scratch

**Case study: Kazakh**

- Which annotation scheme?
- Where to get the data?
- How much data?
- How long will it take?

- **Non-UD:**
  - Perhaps there are existing treebanks for your language and you want to retain compatibility
- **UD:**
  - No need for any special conversion
  - …at least until v3.0 ;)
- **Mixed:**
  - Follow UD guidelines
  - Add information where you think it is useful
  - …providing it is easily convertible

**At the end of the day**: Do what is best for your language and your application

**Free text:**

- Plenty of options:
    - WikiMedia projects: Wikipedia, Wikinews, …
    - Public domain texts (varies by country)
        - Out of copyright (e.g. old literature, folktales)
        - Laws/state administrative texts

**Non-free text:**

- Contact copyright holders early on

- No minimum size
  - Smallest treebank: 1K tokens
  - Biggest treebank: 1.3M tokens

- CoNLL-2006, smallest treebank: 29K tokens

« *You can actually train a parser and get over 50% accuracy for many languages with just about 100 sentences.* » — Dan Zeman

# How Long Will It Take?

- How long is a piece of string ?
- Some approximate numbers:

| Language | Annotators | Tokens | Months |
|---------|-----------|--------|--------|
| Kazakh  | 2 | 4,500 | 1 |
| Buryat  | 1 | 10,000 | 3 |
| Irish   | 1 | 23,600 | 12 |

In all the above cases, annotation guidelines were developed from scratch by people with no prior exposure to UD.

# How We Made A Kazakh Treebank

**Two people:**

- Francis Tyers: Computational linguist, interests in Turkic languages and morphosyntax
- Jonathan North Washington: Phonologist, interests in Turkic languages, fluent speaker of Kazakh and Kyrgyz

**One month:**

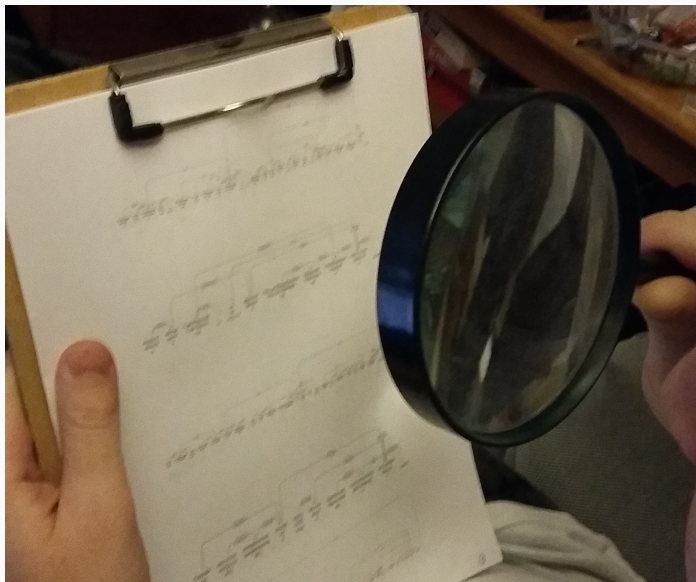- Summer holiday in Bloomington, Indiana

**Resources:**

- Morphological analyser and constraint grammar

The whole thing would have been impossible without the UD project.

- Guidelines were straightforward to apply
- Community was exceptionally helpful and welcoming

# Summary

**What you need to do**

- Join the project
- Start annotating or converting
- Ask if you get stuck!

**Can't wait to get started?**

- Come and talk to us!

# Questions?