# Tutorial on Universal Dependencies

**Infrastructure, resources and tools for UD**

Joakim Nivre[1]    Daniel Zeman[2]    **Filip Ginter**[3]    Francis M. Tyers[4][5]

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

[2]Institute of Formal and Applied Linguistics, Charles University in Prague, Czech Republic

[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

**How many?**

- Languages: **50**
- Treebanks: **72**
- Trees: **642,000**
- Words: **12,400,000**

**Can I use them?**

- Creative Commons and GPL-like: **30**
- Creative Commons Non-Commercial: **42**

**Where from?**

- http://universaldependencies.org
- Official release preferred over GitHub
- Currently officially released: **70 treebanks**
- Twist: test sets currently withheld

# UD Treebanks Come in Many Flavors and Sizes
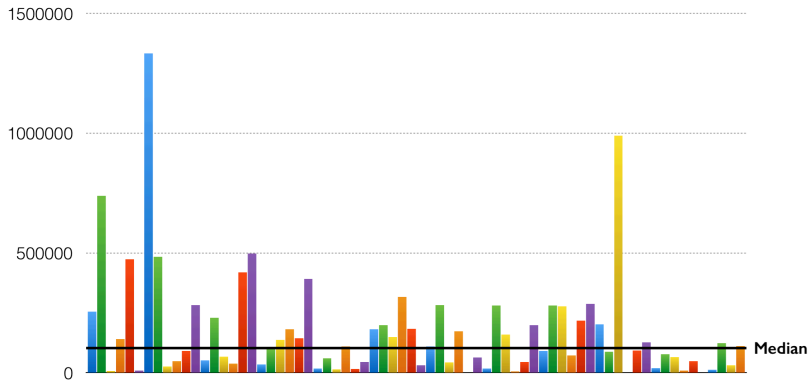
**Annotation:**

- POS and base dependency relations compulsory: 72 treebanks
- ...and additionally:
  - Forms + Features + Lemmas: 58
  - Forms - Features + Lemmas: 4
  - Forms - Features - Lemmas: 7
  - No Forms: 3 (Arabic-NYUAD, English-ESL, Japanese-KTC)
    — licensing
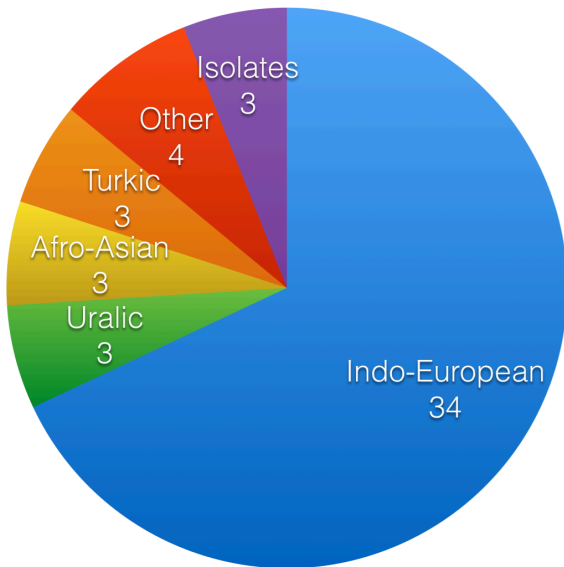
**Size:**

- Smallest: approx. 1000 words — Swedish Sign Language, Kazakh, Sanskrit
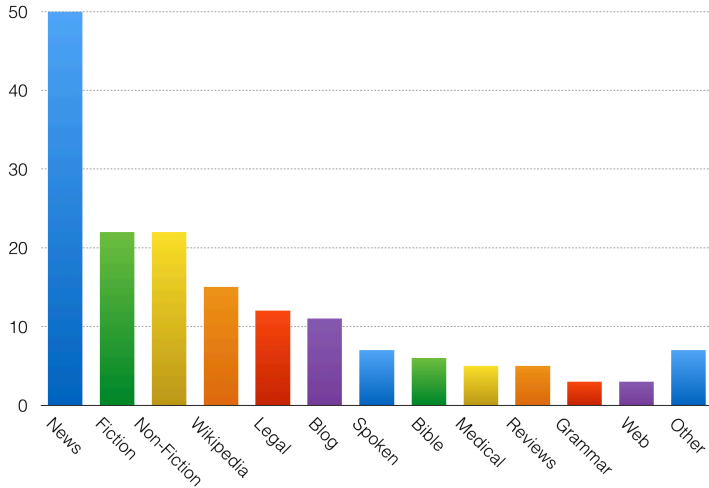- Largest: Czech with 1.3M words, Russian with 980K words

**Treebank Size**

3

Language Family

Isolates 3
Other 4
Turkic 3
Afro-Asian 3
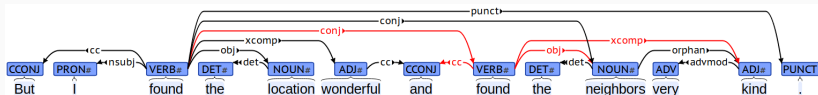Uralic 3
Indo-European 34

4

**Genre**

# CoNLL-U Format

- Derived from CoNLL-X, overall logic same, details differ
- `ID FORM LEMMA UPOS XPOS FEATS HEAD DEPREL DEPS MISC`
- Only `ID UPOS HEAD DEPREL` compulsory

**Distinguishing features:**

- Sentence-level metadata part of the format
- Explicit (and compulsory!) representation of the original text
- `DEPS` field encodes the enhanced dependencies (non-tree structure)
- `MISC` field allows arbitrary data stored for every word
- Empty nodes — only referred to from the enhanced representation
- Words as opposed to tokens

```
# sent_id = reviews-044427-0003
# text = But I found the location wonderful and the neighbors very kind.
1      But         but        CCONJ    CC                                                 3    cc          _          _
2      I           I          PRON     PRP    Case=Nom|Number=Sing|Person=1|PronType=Prs  3    nsubj       _          _
3      found       find       VERB     VBD    Mood=Ind|Tense=Past|VerbForm=Fin            0    root        _          _
4      the         the        DET      DT     Definite=Def|PronType=Art                   5    det         _          _
5      location    location   NOUN     NN     Number=Sing                                 3    obj         _          _
6      wonderful   wonderful  ADJ      JJ     Degree=Pos                                  3    xcomp       _          _
7      and         and        CCONJ    CC                                                 6    cc          _          _
7.1    found       find       VERB     VBD    Mood=Ind|Tense=Past|VerbForm=Fin            _    _           3:conj     _
8      the         the        DET      DT     Definite=Def|PronType=Art                   9    det         _          _
9      neighbors   neighbor   NOUN     NNS    Number=Plur                                 3    conj        7.1:obj    _
10     very        very       ADV      RB                                                 11   advmod      _          _
11     kind        kind       ADJ      JJ     Degree=Pos                                  9    orphan      7.1:xcomp  SpaceAfter=No
12     .           .          PUNCT    .                                                  3    punct       _          _
```

```
16     it       it     PRON    PRP   _  17  nsubj    _  _
17-18  hadn't   _      _       _     _  5   _        _  SpaceAfter=No
17     had      have   VERB    VBD   _  5   ccomp    _  _
18     n't      not    PART    RB    _  17  advmod   _  _
19     .        .      PUNCT   .     _  5   punct    _  _
```

- 83 treebank repositories
- 100+ contributors
- Online documentation consisting of roughly 14,000 web-pages
- Guidelines, universal and language-specific
- Discussions, decision making, validation
- Regular, carefully checked official releases
- A comparatively small group of core "staff" running the show
- Budget: $0

- GitHub in use from Day 1
- Documentation and data first
- Followed exclusive use of the issue tracker for discussions and proposals
  - Before: many email chains — chaos
- Practically *everything* happens openly

# UD is Open



11

- A GitHub repository for every treebank
  - UD_{Language}-{Treebank}
  - **master** branch holds the most recent official release
  - **dev** branch holds development data, not guaranteed to be valid
  - Some teams use GitHub for development, others only to "submit" their data prior to the release
  - No strict requirements on the workflow

- **Official release:** LINDAT, May & November, all treebanks which contain valid data

# Docs

- One set of documentation for every language (not treebank)
- A GitHub repository holding mostly markdown pages
- Special care taken to make it easy to add tree visualizations and examples
- Stubs pre-generated when adding a new language
- 11,000+ commits from 80+ contributors
- Automatically regenerated on every push and published on GitHub pages
- The issue tracker for the *docs* repository is where all the UD activity is happening
  - Hundreds of issues, thousands of replies
- Documentation system: `http://spyysalo.github.io/annodoc/`

## Workflow and Organization

- Highly ~~chaotic~~ distributed
- All contributors given broad edit rights to all data, docs, and tools repositories
- Fully trust-based setup, `git` giving a safety net
- Joakim holds the honorary title of *Chief Cat Herder* and looks after the project as a whole — is obeyed unconditionally

# Validation

- Script to validate treebank data
- Passing is compulsory
- Format validation
- Runs automatically every time a treebank is updated
- Indispensable especially close to an official release date
- Contributors: do we validate?
- Release team: whom to help next?

`http://universaldependencies.org/validation.html`

# Content Validation

- Runs automatically every time a treebank is updated
- Reports "suspicious" syntactic constructions
- Passing not compulsory at the moment
- Contributors: Is there anything odd-looking in my data?
- Release team: Overview of guideline adoption

### Aux chain

Auxiliary dependencies should not form a chain.

Search expression: _ <aux (_ <aux _)

Correct example:



1. Do you think that he will have left when we come ?

Incorrect example:



2. Do you think that he will have left when we come ?

Link to documentation

| | | |
|---|---|---|
| ▸ | Hit overview | |
| ▸ | UD_Basque | 3 hits |
| ▸ | UD_Galician | 10 hits |
| ▸ | UD_Italian | 2 hits |
| ▸ | UD_Japanese | 1 hits |
| ▸ | UD_Persian | 2 hits |
| ▸ | UD_Urdu | 2341 hits |

### Flat is right-headed

Flat relations should be left-headed, not right.

Search expression: _ <flat@R _

Correct example:



3. Carl XVI Gustaf

Incorrect example:



4. Carl XVI Gustaf

`http://universaldependencies.org/svalidation.html`

UD is **not just the treebanks**

- Parsers trained on UD data
- Large multilingual parsebanks
- Query tools for treebanks and parsebanks
- Libraries for handling CoNLL-U
- Tree visualization tools
- Annotation tools

- UDPipe and SyntaxNet
- State-of-the-art parsers, free
- Full-stack parsers: raw text in - parses out
- Models trained on all of UD
- UDPipe — demo & Web API
- UDPipe Web API — get parsed text with a simple HTTP request

# UDPipe

```
ginter@dg:~/eacl17tutorial$ curl -s -F 'data=@test_input.txt' -F 'model=english' -F 'tokenizer=' -F 'tagger=' -F 'parser='
http://lindat.mff.cuni.cz/services/udpipe/api/process | python -c "import sys,json; print json.load(sys.stdin)['result']"
# newdoc
# newpar
# sent_id = 1
# text = This is for the tutorial, so please do try to get it right!
1    This      _    PRON    DEM-SG     _    2     nsubj       _    _
2    is        _    VERB    PRES       _    0     root        _    _
3    for       _    ADP     _          _    5     case        _    _
4    the       _    DET     DEF        _    5     det         _    _
5    tutorial  _    NOUN    SG-NOM     _    2     obl         _    SpaceAfter=No
6    ,         _    PUNCT   Comma      _    2     punct       _    _
7    so        _    ADV     _          _    8     advmod      _    _
8    please    _    INTJ    _          _    10    discourse   _    _
9    do        _    VERB    IMP        _    10    aux         _    _
10   try       _    VERB    INF        _    2     advcl       _    _
11   to        _    PART    _          _    12    mark        _    _
12   get       _    VERB    INF        _    10    xcomp       _    _
13   it        _    PRON    PERS-SG    _    12    obj         _    _
14   right     _    ADJ     POS        _    12    xcomp       _    SpaceAfter=No
15   !         _    PUNCT   ExclMark   _    2     punct       _    _

# sent_id = 2
# text = And I really mean this.
1    And       _    CCONJ   _             _    4    cc       _    _
2    I         _    PRON    PERS-P1SG-NOM _    4    nsubj    _    _
3    really    _    ADV     _             _    4    advmod   _    _
4    mean      _    VERB    PRES          _    0    root     _    _
5    this      _    PRON    DEM-SG        _    4    obj      _    SpaceAfter=No
6    .         _    PUNCT   Period        _    4    punct    _    SpacesAfter=\n
```
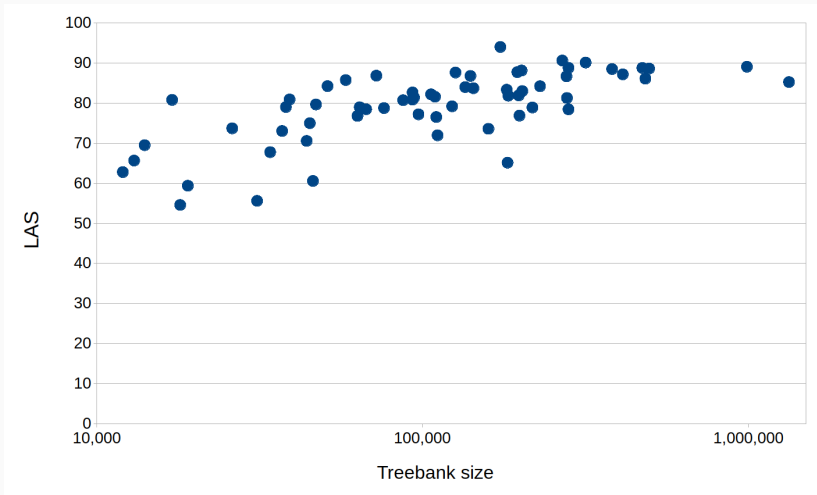
20

- Major improvement upon *SyntaxNet's Parsey's cousins*
- Considerably improved models released mid-March 2017
- `http://tiny.cc/psaurus` — description
- `http://tiny.cc/psaurus-base` — numbers

Average=78% Median=81%

## Parsebanks

- UD-parsed corpora for 45 languages
- Data: CommonCrawl + Wiki + Perseus
- Parses: UDPipe
- Over 90B words total, 630GB zipped CoNLL-U files

Ancient Greek, Arabic, Basque, Bulgarian, Catalan, ChineseT, Croatian, Czech, Danish, Dutch, English, Estonian, Finnish, French, Galician, German, Greek, Hebrew, Hindi, Hungarian, Indonesian, Irish, Italian, Japanese, Kazakh, Korean, Latin, Latvian, Norwegian-Bokmaal, Norwegian-Nynorsk, Old Church Slavonic, Persian, Polish, Portuguese, Romanian, Russian, Slovak, Slovenian, Spanish, Swedish, Turkish, Ukrainian, Urdu, Uyghur, and Vietnamese

## Syntactic Query

- dep_search
- `http://bionlp-www.utu.fi/dep_search`
- Relatively expressive query language, especially geared towards dependencies and rich morphology
- Indexed:
    - Latest UD official release
    - 'dev' branches - reindexed on every push
    - Up to 2 million trees for every language from the UD Parsebanks
- Web and API access
- Used by some during annotation
- Also serves as content validation back-end

# Syntactic Query

# Syntactic Query

```
ginter@dg:~/eacl17tutorial$ curl -sL 'http://epsilon-it.utu.fi/dep_search_webapi?search=_%20%3Ecop%20_%20%21%3Ensubj%20_&db=English-
pbank&case=True&retmax=5000&dl' | tee data.conllu | head -n 27
# db-name: /home/ginter/conll17_idx/English/trees_00000.db
# graph id: 12
# db-name: /home/ginter/conll17_idx/English/trees_00000.db
# graph id: 28
# graph id: 12
# visual-style   14          bgColor:lightgreen
# hittoken:       14          carcinogen      carcinogen      NOUN    NN        Number=Sing    3         conj    _          SpaceAfter=No
# text = Ochratoxin is damaging to the kidneys and liver and is also a suspected carcinogen.
1       Ochratoxin  Ochratoxin  PROPN   NNP     Number=Sing      3         nsubj   _
2       is          be          AUX     VBZ     Mood=Ind|Number=Sing|Person=3|Tense=Pres|VerbForm=Fin    3         aux     _          _
3       damaging    damaging    VERB    VBG     Tense=Pres|VerbForm=Part         0         root    _          _
4       to          to          ADP     IN      6         case
5       the         the         DET     DT      Definite=Def|PronType=Art        6         det     _          _
6       kidneys     kidney      NOUN    NNS     Number=Plur      3         obl     _          _
7       and         and         CCONJ   CC      8         cc
8       liver       liver       NOUN    NN      Number=Sing      6         conj    _          _
9       and         and         CCONJ   CC      14        cc
10      is          be          AUX     VBZ     Mood=Ind|Number=Sing|Person=3|Tense=Pres|VerbForm=Fin    14        cop     _          _
11      also        also        ADV     RB      14        advmod
12      a           a           DET     DT      Definite=Ind|PronType=Art        14        det     _          _
13      suspected   suspect     VERB    VBN     Tense=Past|VerbForm=Part         14        amod    _          _
14      carcinogen  carcinogen  NOUN    NN      Number=Sing      3         conj    _          SpaceAfter=No
15      .           .           PUNCT   .       3         punct   _          _

# db-name: /home/ginter/conll17_idx/English/trees_00000.db
# graph id: 42
```

## Syntactic Query

- PML Tree Query
- `http://lindat.mff.cuni.cz/services/pmltq/`
- A very expressive query language
- Indexed: official UD releases

# Syntactic Query



```
a-root [descendant a-node $z := [descendant a-node $x := [order-precedes $z, parent a-node $y := [order-follows $z]]]]
```

Relations ▾    Node Types ▾    Attributes ▾    Operators ▾    Functions ▾

🔍 Execute query    ▼ w/o Filters

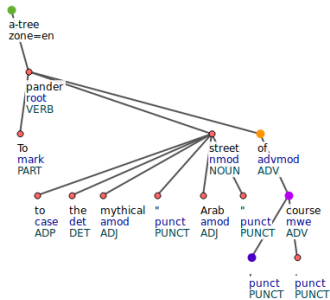← Previous    | 1 |  of 30    Next →    ● 1 a-root    ● 2 a-node

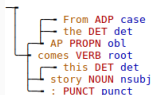[en] To pander to the mythical "Arab street", of course.

## Udapi

- A library and command line tool for processing UD data
  - **Python**, Java, Perl
- Format conversions
- Initial v1-v2 conversion
- Validation tests
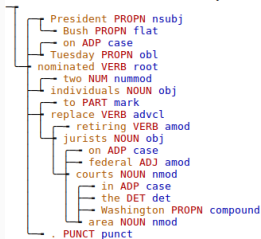- Evaluation, filtering, statistics
- Tree visualization
- `https://udapi.github.io`

```
cat en-ud-dev.conllu | udapy -T | less -R
docname = weblog-blogspot.com_nominations_20041117172713_ENG_20041117_172713
# sent_id = weblog-blogspot.com_nominations_20041117172713_ENG_20041117_172713-0001
# text = From the AP comes this story :

┌──────┌── From ADP case
│      └── the DET det
│  ┌── AP PROPN obl
└── comes VERB root
   ├── this DET det
   ├── story NOUN nsubj
   └── . PUNCT punct

# sent_id = weblog-blogspot.com_nominations_20041117172713_ENG_20041117_172713-0002
# text = President Bush on Tuesday nominated two individuals to replace retiring jurists on federal courts in the Washington area

┌──────── President PROPN nsubj
│      ├── Bush PROPN flat
│   ┌── on ADP case
│  ┌── Tuesday PROPN obl
└── nominated VERB root
   ├── two NUM nummod
   ├── individuals NOUN obj
   │  ┌── to PART mark
   ├── replace VERB advcl
   │     ┌── retiring VERB amod
   │  ┌── jurists NOUN obj
   │  │     ┌── on ADP case
   │  │    ┌── federal ADJ amod
   │  └── courts NOUN nmod
   │        ┌── in ADP case
   │       ├── the DET det
   │       ├── Washington PROPN compound
   │      └── area NOUN nmod
   └── . PUNCT punct

# sent_id = weblog-blogspot.com_nominations_20041117172713_ENG_20041117_172713-0003
# text = Bush nominated Jennifer M. Anderson for a 15-year term as associate judge of the Superior Court of the District of Colum
```
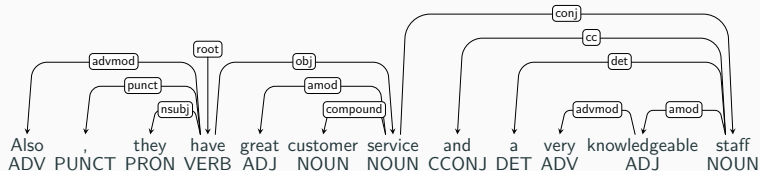
30

# Tree Visualization Tools

```
cat en-ud-dev.conllu | udapy write.Tikz
```
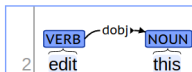
# Tree Visualization Tools

http://spyysalo.github.io/conllu.js/

http://spyysalo.github.io/annodoc/sdparse.html



Input (editable):

```
1       edit    edit    VERB    VERB    _       0       root    _       _
2       this    this    NOUN    NOUN    _       1       dobj    _       _
```

## Annotation Tools

- No official annotation tool (yet)
- A list of tools:
  http://universaldependencies.org/tools.html
- At present, none downright outstanding

**Questions?**