# Universal Dependencies

Sara Stymne
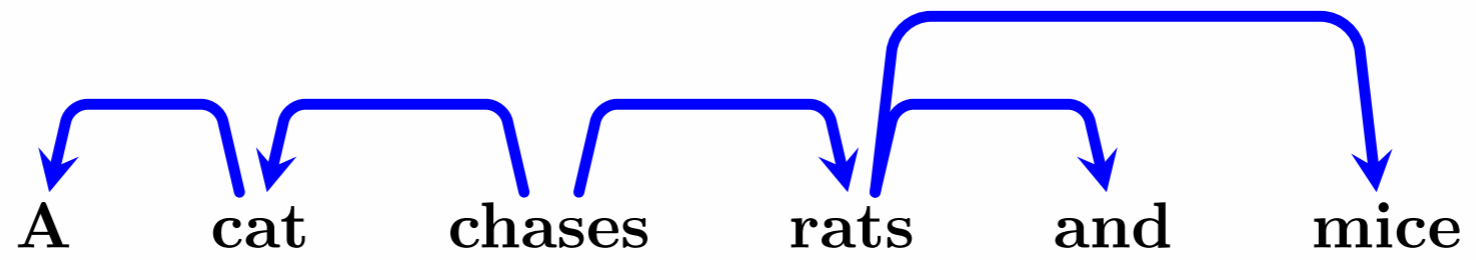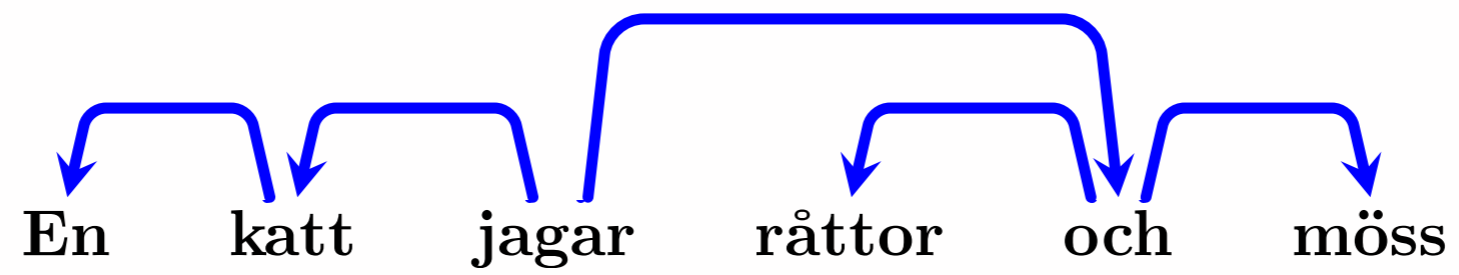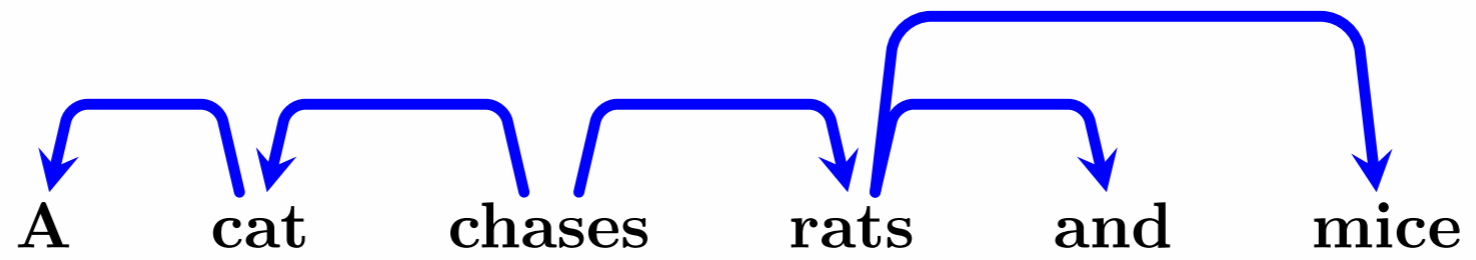
# Introduction

Growing interest in multilingual and cross-lingual NLP

- Multilingual evaluation campaigns to test generality of approaches
- Cross-lingual learning to support low-resource languages

Growing awareness of methodological problems

- Current NLP relies heavily on annotation
- Annotation guidelines vary across languages

A cat chases rats and mice

A    cat    chases    rats    and    mice

En    katt    jagar    råttor    och    möss

A cat chases rats and mice

En katt jagar råttor och möss

En kat jager rotter og mus

## Sentence 1

det — A → cat
nsubj — cat → chases
dobj — chases → rats
cc — rats → and
conj — rats → mice

**A cat chases rats and mice**

## Sentence 2

DT — En → katt
SS — katt → jagar
OO — jagar → och
CJ — och → råttor
CJ — och → möss

**En katt jagar råttor och möss**

## Sentence 3

subj — En → jager
nobj — En → kat
dobj — jager → rotter
coord — rotter → og
conj — og → mus

**En kat jager rotter og mus**

Sentence 1: A cat chases rats and mice
- det: A → cat
- nsubj: cat → chases
- dobj: chases → rats
- cc: rats → and
- conj: rats → mice

Sentence 2: En katt jagar råttor och möss
- DT: En → katt
- SS: katt → jagar
- OO: jagar → och
- CJ: och → råttor
- CJ: och → möss

Sentence 3: En kat jager rotter og mus
- subj: En → jager
- nobj: En → kat
- dobj: jager → rotter
- coord: rotter → og
- conj: og → mus

# Why is this a problem?

- Hard to compare empirical results across languages

- Hard to usefully do cross-lingual structure transfer

- Hard to evaluate cross-lingual learning

- Hard to build and maintain multilingual systems

- Hard to make comparative linguistic studies

- Hard to validate linguistic typology

- Hard to make progress towards a universal parser

# Universal Dependencies v2:
# A Multilingual Treebank Collection

Joakim Nivre, Marie-Catherine de Marneffe, Filip Ginter

Yoav Goldberg, Jan Hajic, Christopher D. Manning

Ryan McDonald, Slav Petrov, Sampo Pyysalo

Natalia Silveira, Reut Tsarfaty, Daniel Zeman and many others

# Universal Dependencies

http://universaldependencies.org

# Universal Dependencies

http://universaldependencies.org



Part-of-speech tags

# Universal Dependencies

http://universaldependencies.org



Part-of-speech tags

Morphological features

# Universal Dependencies

http://universaldependencies.org

Dependency relations

Part-of-speech tags

Morphological features

# The UD Philosophy

Maximize parallelism – but don't overdo it

- Don't annotate the same thing in different ways
- Don't make different things look the same
- Don't annotate things that are not there

Universal taxonomy with language-specific elaboration

- Languages select from a universal pool of categories
- Allow language-specific extensions

# Morphology

Toutefois , les filles adorent les desserts .

# Morphology

| Toutefois | , | les | filles | adorent | les | desserts | . |
|-----------|---|-----|--------|---------|-----|----------|---|
| toutefois | , | le | fille | adorer | le | dessert | . |

- Lemma representing the semantic content of the word

# Morphology

| Toutefois | , | les | filles | adorent | les | desserts | . |
|-----------|---|-----|--------|---------|-----|----------|---|
| toutefois | , | le | fille | adorer | le | dessert | . |
| ADV | PUNCT | DET | NOUN | VERB | DET | NOUN | PUNCT |

- Lemma representing the semantic content of the word

- Part-of-speech tag representing the abstract lexical category associated with the word

# Morphology

| Open | Closed | Other |
|------|--------|-------|
| ADJ | ADP | PUNCT |
| ADV | AUX | SYM |
| INTJ | CONJ | X |
| NOUN | DET | |
| PROPN | NUM | |
| VERB | PART | |
| | PRON | |
| | SCONJ | |

Toutefois ,
toutefois ,
ADV PUNCT

desserts .
dessert .
NOUN PUNCT

- Lemma rep                                    of the word

- Part-of-speech tag representing the abstract lexical category associated with the word

# Morphology

| Toutefois | , | les | filles | adorent | les | desserts | . |
|-----------|---|-----|--------|---------|-----|----------|---|
| toutefois | , | le | fille | adorer | le | dessert | . |
| ADV | PUNCT | DET | NOUN | VERB | DET | NOUN | PUNCT |
| | | Definite=Def Number=Plur | Gender=Fem Number=Plur | Number=Plur Person=3 Tense=Pres | Definite=Def Number=Plur | Gender=Masc Number=Plur | |

- Lemma representing the semantic content of the word

- Part-of-speech tag representing the abstract lexical category associated with the word

- Features representing lexical and grammatical properties associated with the lemma or the particular word form

# Morphology

| Toutefois | , | les | filles | adorent | les | desserts | . |
|-----------|---|-----|--------|---------|-----|----------|---|
| toutefois | , | le | fille | adorer | le | dessert | . |
| ADV | PUNCT | DET | NOUN | VERB | DET | NOUN | PUNCT |
| | | Definite=Def Number=Plur | Gender=Fem Number=Plur | Number=Plur Person=3 Tense=Pres | Definite=Def Number=Plur | Gender=Masc Number=Plur | |

## Lexicalism

- Basic annotation units are words

- Morphological analysis through features, not segmentation

- Clitics and contractions are segmented if required by syntax

# Syntax

| The | cat | could | have | chased | all | the | dogs | down | the | street | . |
|-----|-----|-------|------|--------|-----|-----|------|------|-----|--------|---|
| DET | NOUN | AUX | AUX | VERB | DET | DET | NOUN | ADP | DET | NOUN | PUNCT |

# Syntax



- Content words are related by dependency relations

# Syntax



- Content words are related by dependency relations

- Function words attach to the content word they modify

# Syntax



- Content words are related by dependency relations

- Function words attach to the content word they modify

- Punctuation attach to head of phrase or clause

The dog was chased by the cat .
DET NOUN AUX VERB ADP DET NOUN PUNCT

root — chased
nsubjpass — dog
nmod — cat
punct — .

Hunden jagades av katten .
NOUN VERB ADP NOUN PUNCT

root — jagades
nsubjpass — Hunden
nmod — katten
punct — .

The dog was chased by the cat .
DET NOUN AUX VERB ADP DET NOUN PUNCT

- det: The
- nsubjpass: dog
- root: chased
- punct: .
- nmod: cat
- det: the

Hunden jagades av katten .
NOUN VERB ADP NOUN PUNCT
Definite=Def          Definite=Def

- nsubjpass: Hunden
- root: jagades
- punct: .
- nmod: katten

The dog was chased by the cat .
DET NOUN AUX VERB ADP DET NOUN PUNCT

root
nsubjpass
det
auxpass
nmod
punct
det

Hunden jagades av katten .
NOUN VERB ADP NOUN PUNCT
**Definite=Def** **Voice=Pass** **Definite=Def**

root
nsubjpass
nmod
punct

**Sentence 1:**

| The | dog | was | chased | by | the | cat | . |
|-----|-----|-----|--------|-----|-----|-----|---|
| DET | NOUN | AUX | VERB | ADP | DET | NOUN | PUNCT |

Dependencies: det, nsubjpass, auxpass, root, case, det, nmod, punct

**Sentence 2:**

| Hunden | jagades | av | katten | . |
|--------|---------|-----|--------|---|
| NOUN | VERB | ADP | NOUN | PUNCT |
| Definite=Def | Voice=Pass | | Definite=Def | |

Dependencies: nsubjpass, root, case, nmod, punct

# Dependency Relations

## Taxonomy of 40 universal grammatical relations

- Three types of structures: nominals, clauses, modifiers

- Core arguments vs. other dependents (not complements vs. adjuncts)

## A two-level architecture

- Universal: broad categories to allow cross-linguistic comparison

- Language-specific: subtypes to capture language-specific phenomena

| Universal | Subtype |
|-----------|---------|
| acl | acl:relcl |
| compound | compound:prt |
| nmod | nmod:poss |

# Release v2

102 treebanks

61 languages

145 contributors

http://universaldependencies.org

## UD Treebanks

| Language | Size | | | | | | |
|---|---|---|---|---|---|---|---|
| Amharic | – | | – | ? | – | | |
| Ancient Greek | 244K | | | | ✔ | | |
| Ancient Greek-PROIEL | 206K | | – | | ✔ | | |
| Arabic | 242K | | – | | ✔ | | |
| Basque | 121K | | | | ✔ | | |
| Bulgarian | 156K | | | | ✔ | | |
| Catalan | 530K | | | | | | |
| Chinese | 123K | | | | | | |
| Croatian | 87K | | – | | ✔ | | W |
| Czech | 1,503K | | | | ✔ | | |
| Czech-CAC | 493K | | | | | | |
| Czech-CLTT | 35K | | | | | | |
| Danish | 100K | | | | ✔ | | |
| Dutch | 209K | | – | | ✔ | | |
| Dutch-LassySmall | 98K | | – | | | | W |
| English | 254K | | | | ✔ | | |
| English-ESL | 97K | | | | | | |
| English-LinES | 82K | | | | | | |
| Estonian | 234K | | – | | ✔ | | |
| Finnish | 181K | | | | ✔ | | |
| Finnish-FTB | 159K | | – | | ✔ | | |
| French | 390K | | | | ✔ | | W |
| Galician | 138K | | | | | | |
| German | 293K | | – | | ✔ | | W |
| Gothic | 56K | | – | | ✔ | | |
| Greek | 59K | | | | ✔ | | W |
| Hebrew | 115K | | – | | ✔ | | |
| Hindi | 351K | | – | | ✔ | | |
| Hungarian | 42K | | | | | | |
| Indonesian | 121K | | – | | ✔ | | |
| Irish | 23K | | | | ✔ | | |
| Italian | 252K | | | | ✔ | | W |
| Japanese-KTC | 267K | | | | ✔ | | |
| Kazakh | 4K | | | | | | W |
| Korean | – | | – | – | – | | |
| Latin | 47K | | – | | ✔ | | |
| Latin-ITTB | 291K | | – | | ✔ | | |
| Latin-PROIEL | 165K | | – | | ✔ | | |
| Latvian | 20K | | – | | | | |
| Norwegian | 311K | | | | ✔ | | |
| Old Church Slavonic | 57K | | – | | ✔ | | |
| Persian | 151K | | | | ✔ | | |
| Polish | 83K | | – | | ✔ | | |
| Portuguese | 226K | | – | | ✔ | | |
| Portuguese-BR | 298K | | – | | | | |
| Romanian | 145K | | | | ✔ | | W |
| Russian | 99K | | | | | | W |
| Russian-SynTagRus | 1,032K | | | | ✔ | | |
| Slovenian | 140K | | | | ✔ | | |
| Slovenian-SST | 29K | | | | | | |
| Spanish | 423K | | | | ✔ | | W |
| Spanish-AnCora | 547K | | | | | | |
| Swedish | 96K | | | | ✔ | | |
| Swedish-LinES | 79K | | | | | | |
| Tamil | 8K | | – | | ✔ | | |
| Turkish | 56K | | | | | | |
| Ukrainian | – | | – | | – | | |

**Release v2**

102 treebanks

61 languages

145 contributors

http://universaldependencies.org

**Size:**

<1K to 1,5M tokens

**UD Treebanks**

| Language | Size | | | | | | |
|---|---|---|---|---|---|---|---|
| Amharic | – | | – | ? | – | | |
| Ancient Greek | 244K | | | | ✔ | | |
| Ancient Greek-PROIEL | 206K | | – | | ✔ | | |
| Arabic | 242K | | – | | ✔ | | |
| Basque | 121K | | | | ✔ | | |
| Bulgarian | 156K | | | | ✔ | | |
| Catalan | 530K | | | | | | |
| Chinese | 123K | | | | | | |
| Croatian | 87K | | – | | ✔ | | |
| Czech | 1,503K | | | | ✔ | | |
| Czech-CAC | 493K | | | | | | |
| Czech-CLTT | 35K | | | | | | |
| Danish | 100K | | | | ✔ | | |
| Dutch | 209K | | – | | ✔ | | |
| Dutch-LassySmall | 98K | | – | | | | |
| English | 254K | | | | ✔ | | |
| English-ESL | 97K | | | | | | |
| English-LinES | 82K | | | | | | |
| Estonian | 234K | | – | | ✔ | | |
| Finnish | 181K | | | | ✔ | | |
| Finnish-FTB | 159K | | – | | ✔ | | |
| French | 390K | | | | ✔ | | |
| Galician | 138K | | | | | | |
| German | 293K | | – | | ✔ | | |
| Gothic | 56K | | – | | ✔ | | |
| Greek | 59K | | | | ✔ | | |
| Hebrew | 115K | | – | | ✔ | | |
| Hindi | 351K | | – | | ✔ | | |
| Hungarian | 42K | | | | | | |
| Indonesian | 121K | | – | | ✔ | | |
| Irish | 23K | | | | ✔ | | |
| Italian | 252K | | | | ✔ | | |
| Japanese-KTC | 267K | | | | ✔ | | |
| Kazakh | 4K | | | | | | |
| Korean | – | | – | – | – | | |
| Latin | 47K | | – | | ✔ | | |
| Latin-ITTB | 291K | | – | | ✔ | | |
| Latin-PROIEL | 165K | | – | | ✔ | | |
| Latvian | 20K | | – | | | | |
| Norwegian | 311K | | | | ✔ | | |
| Old Church Slavonic | 57K | | – | | ✔ | | |
| Persian | 151K | | | | ✔ | | |
| Polish | 83K | | – | | ✔ | | |
| Portuguese | 226K | | – | | ✔ | | |
| Portuguese-BR | 298K | | – | | | | |
| Romanian | 145K | | | | ✔ | | |
| Russian | 99K | | | | ✔ | | |
| Russian-SynTagRus | 1,032K | | | | | | |
| Slovenian | 140K | | | | ✔ | | |
| Slovenian-SST | 29K | | | | | | |
| Spanish | 423K | | | | ✔ | | |
| Spanish-AnCora | 547K | | | | | | |
| Swedish | 96K | | | | ✔ | | |
| Swedish-LinES | 79K | | | | ✔ | | |
| Tamil | 8K | | – | | ✔ | | |
| Turkish | 56K | | | | | | |
| Ukrainian | – | | – | | – | | |

# Release v2

102 treebanks

61 languages

145 contributors

http://universaldependencies.org

## Size:

<1K to 1,5M tokens

## Annotation:

lemmas (83), features (94)

**UD Treebanks**

| Language | Size | | | | | | |
|---|---|---|---|---|---|---|---|
| Amharic | – | | – | ? | – | | |
| Ancient Greek | 244K | | | | ✔ | | |
| Ancient Greek-PROIEL | 206K | | – | | ✔ | | |
| Arabic | 242K | | – | | ✔ | | |
| Basque | 121K | | | | ✔ | | |
| Bulgarian | 156K | | | | ✔ | | |
| Catalan | 530K | | | | | | |
| Chinese | 123K | | | | | | |
| Croatian | 87K | | – | | ✔ | | |
| Czech | 1,503K | | | | ✔ | | |
| Czech-CAC | 493K | | | | | | |
| Czech-CLTT | 35K | | | | | | |
| Danish | 100K | | | | ✔ | | |
| Dutch | 209K | | – | | ✔ | | |
| Dutch-LassySmall | 98K | | – | | | | |
| English | 254K | | | | ✔ | | |
| English-ESL | 97K | | | | | | |
| English-LinES | 82K | | | | ✔ | | |
| Estonian | 234K | | – | | ✔ | | |
| Finnish | 181K | | | | ✔ | | |
| Finnish-FTB | 159K | | – | | ✔ | | |
| French | 390K | | | | ✔ | | |
| Galician | 138K | | | | | | |
| German | 293K | | – | | ✔ | | |
| Gothic | 56K | | – | | ✔ | | |
| Greek | 59K | | | | ✔ | | |
| Hebrew | 115K | | – | | ✔ | | |
| Hindi | 351K | | – | | ✔ | | |
| Hungarian | 42K | | | | ✔ | | |
| Indonesian | 121K | | – | | ✔ | | |
| Irish | 23K | | | | ✔ | | |
| Italian | 252K | | | | ✔ | | |
| Japanese-KTC | 267K | | | | ✔ | | |
| Kazakh | 4K | | | | | | |
| Korean | – | | – | – | – | | |
| Latin | 47K | | – | | ✔ | | |
| Latin-ITTB | 291K | | – | | ✔ | | |
| Latin-PROIEL | 165K | | – | | ✔ | | |
| Latvian | 20K | | – | | | | |
| Norwegian | 311K | | | | ✔ | | |
| Old Church Slavonic | 57K | | – | | ✔ | | |
| Persian | 151K | | | | ✔ | | |
| Polish | 83K | | – | | ✔ | | |
| Portuguese | 226K | | – | | ✔ | | |
| Portuguese-BR | 298K | | – | | | | |
| Romanian | 145K | | | | ✔ | | |
| Russian | 99K | | | | ✔ | | |
| Russian-SynTagRus | 1,032K | | | | ✔ | | |
| Slovenian | 140K | | | | ✔ | | |
| Slovenian-SST | 29K | | | | | | |
| Spanish | 423K | | | | ✔ | | |
| Spanish-AnCora | 547K | | | | ✔ | | |
| Swedish | 96K | | | | ✔ | | |
| Swedish-LinES | 79K | | | | ✔ | | |
| Tamil | 8K | | – | | ✔ | | |
| Turkish | 56K | | | | | | |
| Ukrainian | – | | – | | ✔ | – | |

**Release v2**

102 treebanks

61 languages

145 contributors

http://universaldependencies.org

**Size:**

<1K to 1,5M tokens

**Annotation:**

lemmas (83), features (94)

**Language-specific guidelines:**

complete (14), partial (34)



**UD Treebanks**

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Amharic | – | | – | ? | – | | |
| Ancient Greek | 244K | | | | ✔ | | |
| Ancient Greek-PROIEL | 206K | | – | | ✔ | | |
| Arabic | 242K | | – | | ✔ | | |
| Basque | 121K | | | | ✔ | | |
| Bulgarian | 156K | | | | ✔ | | |
| Catalan | 530K | | | | | | |
| Chinese | 123K | | | | | | |
| Croatian | 87K | | – | | ✔ | | |
| Czech | 1,503K | | | | ✔ | | |
| Czech-CAC | 493K | | | | | | |
| Czech-CLTT | 35K | | | | | | |
| Danish | 100K | | | | ✔ | | |
| Dutch | 209K | | – | | | | |
| Dutch-LassySmall | 98K | | – | | | | |
| English | 254K | | | | ✔ | | |
| English-ESL | 97K | | | | | | |
| English-LinES | 82K | | | | | | |
| Estonian | 234K | | – | | ✔ | | |
| Finnish | 181K | | | | ✔ | | |
| Finnish-FTB | 159K | | – | | ✔ | | |
| French | 390K | | | | ✔ | | |
| Galician | 138K | | | | | | |
| German | 293K | | – | | ✔ | | |
| Gothic | 56K | | – | | ✔ | | |
| Greek | 59K | | | | ✔ | | |
| Hebrew | 115K | | – | | ✔ | | |
| Hindi | 351K | | – | | ✔ | | |
| Hungarian | 42K | | | | ✔ | | |
| Indonesian | 121K | | – | | ✔ | | |
| Irish | 23K | | | | ✔ | | |
| Italian | 252K | | | | ✔ | | |
| Japanese-KTC | 267K | | | | ✔ | | |
| Kazakh | 4K | | | | | | |
| Korean | | | | – | – | | |
| Latin | 47K | | – | | ✔ | | |
| Latin-ITTB | 291K | | – | | ✔ | | |
| Latin-PROIEL | 165K | | – | | ✔ | | |
| Latvian | 20K | | – | | | | |
| Norwegian | 311K | | | | ✔ | | |
| Old Church Slavonic | 57K | | – | | ✔ | | |
| Persian | 151K | | | | ✔ | | |
| Polish | 83K | | – | | ✔ | | |
| Portuguese | 226K | | – | | ✔ | | |
| Portuguese-BR | 298K | | – | | | | |
| Romanian | 145K | | | | ✔ | | |
| Russian | 99K | | | | ✔ | | |
| Russian-SynTagRus | 1,032K | | | | ✔ | | |
| Slovenian | 140K | | | | ✔ | | |
| Slovenian-SST | 29K | | | | | | |
| Spanish | 423K | | | | ✔ | | |
| Spanish-AnCora | 547K | | | | | | |
| Swedish | 96K | | | | ✔ | | |
| Swedish-LinES | 79K | | | | ✔ | | |
| Tamil | 8K | | – | | ✔ | | |
| Turkish | 56K | | | | | | |
| Ukrainian | – | | – | | ✔ | – | | |

**Release v2**
102 treebanks
61 languages
145 contributors

http://universaldependencies.org

**Size:**
<1K to 1,5M tokens

**Annotation:**
lemmas (83), features (94)

**Language-specific guidelines:**
complete (14), partial (34)

**Manual annotation/ validation:**
complete (34), partial (31)

**UD Treebanks**

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Amharic | – | | – | ? | – | | |
| Ancient Greek | 244K | | | | ✔ | | |
| Ancient Greek-PROIEL | 206K | | – | | ✔ | | |
| Arabic | 242K | | – | | ✔ | | |
| Basque | 121K | | | | ✔ | | |
| Bulgarian | 156K | | | | ✔ | | |
| Catalan | 530K | | | | | | |
| Chinese | 123K | | | | | | |
| Croatian | 87K | | – | | ✔ | | |
| Czech | 1,503K | | | | ✔ | | |
| Czech-CAC | 493K | | | | | | |
| Czech-CLTT | 35K | | | | | | |
| Danish | 100K | | | | ✔ | | |
| Dutch | 209K | | – | | | | |
| Dutch-LassySmall | 98K | | – | | | | |
| English | 254K | | | | ✔ | | |
| English-ESL | 97K | | | | | | |
| English-LinES | 82K | | | | ✔ | | |
| Estonian | 234K | | – | | ✔ | | |
| Finnish | 181K | | | | ✔ | | |
| Finnish-FTB | 159K | | – | | ✔ | | |
| French | 390K | | | | ✔ | | |
| Galician | 138K | | | | | | |
| German | 293K | | – | | ✔ | | |
| Gothic | 56K | | – | | ✔ | | |
| Greek | 59K | | | | ✔ | | |
| Hebrew | 115K | | – | | ✔ | | |
| Hindi | 351K | | – | | ✔ | | |
| Hungarian | 42K | | | | ✔ | | |
| Indonesian | 121K | | – | | ✔ | | |
| Irish | 23K | | | | ✔ | | |
| Italian | 252K | | | | ✔ | | |
| Japanese-KTC | 267K | | | | ✔ | | |
| Kazakh | 4K | | | | | | |
| Korean | – | | – | – | – | | |
| Latin | 47K | | – | | ✔ | | |
| Latin-ITTB | 291K | | – | | ✔ | | |
| Latin-PROIEL | 165K | | – | | ✔ | | |
| Latvian | 20K | | – | | | | |
| Norwegian | 311K | | | | ✔ | | |
| Old Church Slavonic | 57K | | – | | ✔ | | |
| Persian | 151K | | | | ✔ | | |
| Polish | 83K | | – | | ✔ | | |
| Portuguese | 226K | | – | | ✔ | | |
| Portuguese-BR | 298K | | – | | | | |
| Romanian | 145K | | | | ✔ | | |
| Russian | 99K | | | | ✔ | | |
| Russian-SynTagRus | 1,032K | | | | ✔ | | |
| Slovenian | 140K | | | | ✔ | | |
| Slovenian-SST | 29K | | | | | | |
| Spanish | 423K | | | | ✔ | | |
| Spanish-AnCora | 547K | | | | ✔ | | |
| Swedish | 96K | | | | ✔ | | |
| Swedish-LinES | 79K | | | | ✔ | | |
| Tamil | 8K | | – | | ✔ | | |
| Turkish | 56K | | | | | | |
| Ukrainian | – | | – | | ✔ | – | |

# Release v2

102 treebanks

61 languages

145 contributors

http://universaldependencies.org

## Size:

<1K to 1,5M tokens

## Annotation:

lemmas (83), features (94)

## Language-specific guidelines:

complete (14), partial (34)
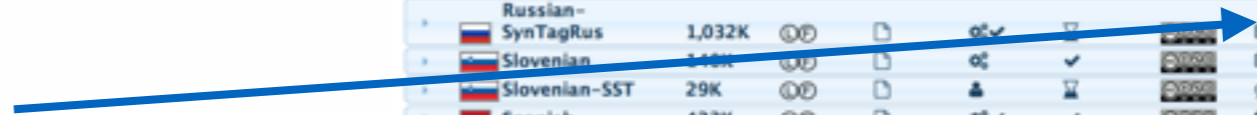
## Manual annotation/ validation:

complete (34), partial (31)

## Genres:

bible, blog, fiction, grammar examples, legal text, medical text, news, non-fiction, reviews, spoken, social, web, wikipedia

**UD Treebanks**

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | Amharic | – | | – | ? | – | | |
| | Ancient Greek | 244K | | | | ✓ | | |
| | Ancient Greek-PROIEL | 206K | | – | | ✓ | | |
| | Arabic | 242K | | – | | ✓ | | |
| | Basque | 121K | | | | ✓ | | |
| | Bulgarian | 156K | | | | ✓ | | |
| | Catalan | 530K | | | | | | |
| | Chinese | 123K | | | | | | |
| | Croatian | 87K | | – | | ✓ | | |
| | Czech | 1,503K | | | | ✓ | | |
| | Czech-CAC | 493K | | | | | | |
| | Czech-CLTT | 35K | | | | | | |
| | Danish | 100K | | | | ✓ | | |
| | Dutch | 209K | | – | | | | |
| | Dutch-LassySmall | 98K | | – | | | | |
| | English | 254K | | | | ✓ | | |
| | English-ESL | 97K | | | | | | |
| | English-LinES | 82K | | | | | | |
| | Estonian | 234K | | – | | ✓ | | |
| | Finnish | 181K | | | | ✓ | | |
| | Finnish-FTB | 159K | | – | | ✓ | | |
| | French | 390K | | | | ✓ | | |
| | Galician | 138K | | | | | | |
| | German | 293K | | – | | ✓ | | |
| | Gothic | 56K | | – | | ✓ | | |
| | Greek | 59K | | | | ✓ | | |
| | Hebrew | 115K | | – | | ✓ | | |
| | Hindi | 351K | | – | | ✓ | | |
| | Hungarian | 42K | | | | ✓ | | |
| | Indonesian | 121K | | – | | ✓ | | |
| | Irish | 23K | | | | ✓ | | |
| | Italian | 252K | | | | ✓ | | |
| | Japanese-KTC | 267K | | | | ✓ | | |
| | Kazakh | 4K | | | | | | |
| | Korean | – | | – | – | – | | |
| | Latin | 47K | | – | | ✓ | | |
| | Latin-ITTB | 291K | | – | | ✓ | | |
| | Latin-PROIEL | 165K | | – | | ✓ | | |
| | Latvian | 20K | | – | | | | |
| | Norwegian | 311K | | | | ✓ | | |
| | Old Church Slavonic | 57K | | – | | ✓ | | |
| | Persian | 151K | | | | ✓ | | |
| | Polish | 83K | | – | | ✓ | | |
| | Portuguese | 226K | | – | | ✓ | | |
| | Portuguese-BR | 298K | | – | | | | |
| | Romanian | 145K | | | | ✓ | | |
| | Russian | 99K | | | | ✓ | | |
| | Russian-SynTagRus | 1,032K | | | | | | |
| | Slovenian | 140K | | | | ✓ | | |
| | Slovenian-SST | 29K | | | | | | |
| | Spanish | 423K | | | | ✓ | | |
| | Spanish-AnCora | 547K | | | | ✓ | | |
| | Swedish | 96K | | | | ✓ | | |
| | Swedish-LinES | 79K | | | | ✓ | | |
| | Tamil | 8K | | – | | ✓ | | |
| | Turkish | 56K | | | | | | |
| | Ukrainian | – | | – | | ✓ | – | |

**Release v2**

102 treebanks

61 languages

145 contributors

http://universaldependencies.org

**Size:**

<1K to 1,5M tokens

**Annotation:**

lemmas (83), fe____ (94)

**Language-sp_____uidelines:**

co____

**Manual**
**validation**

complete__

**Genres:**

bible, bl__ fiction, grammar exa___es,
legal text, medical text, news, non-fiction,
reviews, spoken, social, web, wikipedia

**Over 50% of the world's native speakers**

**UD Treebanks**

| | Language | Size | | | | | | |
|---|---|---|---|---|---|---|---|---|
| ► | Amharic | – | | – | ? | – | | |
| ► | Ancient Greek | 244K | | | | ✔ | | |
| ► | Ancient Greek-PROIEL | 206K | | – | | ✔ | | |
| ► | Arabic | 242K | | – | | ✔ | | |
| ► | Basque | 121K | | | | ✔ | | |
| ► | Bulgarian | 156K | | | | ✔ | | |
| ► | Catalan | 530K | | | | | | |
| ► | Chinese | 123K | | | | | | W |
| ► | Croatian | 87K | | – | | ✔ | | W |
| ► | Czech | 1,503K | | | | ✔ | | |
| ► | Czech-CAC | 493K | | | | | | |
| ► | Czech-CLTT | 35K | | | | | | |
| ► | Danish | 100K | | | | ✔ | | |
| ► | Dutch | 209K | | – | | ✔ | | |
| ► | Dutch-LassySmall | 98K | | – | | | | W |
| ► | English | 254K | | | | ✔ | | |
| ► | English-ESL | 97K | | | | | | |
| ► | English-LinES | 82K | | | | | | |
| ► | Estonian | 234K | | – | | ✔ | | |
| ► | Finnish | 181K | | | | ✔ | | |
| ► | Finnish-FTB | 159K | | – | | ✔ | | |
| ► | French | 390K | | | | ✔ | | |
| ► | Galician | 138K | | | | | | |
| ► | German | 293K | | – | | ✔ | | |
| ► | Gothic | 56K | | – | | ✔ | | |
| ► | Greek | 59K | | | | ✔ | | |
| ► | Hebrew | 115K | | – | | ✔ | | |
| ► | Hindi | 351K | | – | | ✔ | | |
| ► | Hungarian | 42K | | | | | | |
| ► | Indonesian | 121K | | – | | ✔ | | |
| ► | Irish | 23K | | | | ✔ | | |
| ► | Italian | 252K | | | | ✔ | | |
| ► | Japanese-KTC | 267K | | | | | | |
| ► | Kazakh | 4K | | | | | | W |
| ► | Korean | – | | – | – | – | | |
| ► | Latin | 47K | | – | | ✔ | | |
| ► | Latin-ITTB | 291K | | – | | ✔ | | |
| ► | Latin-PROIEL | 165K | | – | | ✔ | | |
| ► | Latvian | 20K | | – | | | | |
| ► | Norwegian | 311K | | | | ✔ | | |
| ► | Old Church Slavonic | 57K | | – | | ✔ | | |
| ► | Persian | 151K | | | | ✔ | | |
| ► | Polish | 83K | | – | | ✔ | | |
| ► | Portuguese | 226K | | – | | ✔ | | |
| ► | Portuguese-BR | 298K | | – | | ✔ | | |
| ► | Romanian | 145K | | | | ✔ | | W |
| ► | Russian | 99K | | | | ✔ | | W |
| ► | Russian-SynTagRus | 1,032K | | | | | | |
| ► | Slovenian | 140K | | | | | | |
| ► | Slovenian-SST | 29K | | | | | | |
| ► | Spanish | 423K | | | | ✔ | | |
| ► | Spanish-AnCora | 547K | | | | ✔ | | |
| ► | Swedish | 96K | | | | ✔ | | |
| ► | Swedish-LinES | 79K | | | | ✔ | | |
| ► | Tamil | 8K | | – | | ✔ | | |
| ► | Turkish | 56K | | | | | | |
| ► | Ukrainian | – | | – | | ✔ | – | |

**Release v2**

102 treebanks

61 languages

145 contributors

http://universaldependencies.org

**Size:**

<1K to 1,5M tokens

**Annotation:**

lemmas (83), features (94)

**Language-specific guidelines:**

complete (14), partial (34)
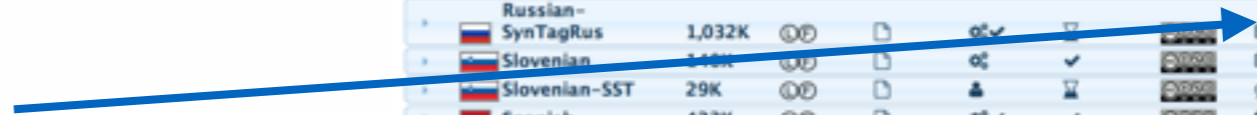
**Manual annotation/ validation:**

complete (34), partial (31)

**Genres:**

bible, blog, fiction, grammar examples, legal text, medical text, news, non-fiction, reviews, spoken, social, web, wikipedia

**UD Treebanks**

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | Amharic | – | | – | ? | – | | |
| | Ancient Greek | 244K | | | | ✓ | | |
| | Ancient Greek-PROIEL | 206K | | – | | ✓ | | |
| | Arabic | 242K | | – | | ✓ | | |
| | Basque | 121K | | | | ✓ | | |
| | Bulgarian | 156K | | | ✓ | ✓ | | |
| | Catalan | 530K | | | ✓ | | | |
| | Chinese | 123K | | | | | | |
| | Croatian | 87K | | – | | ✓ | | |
| | Czech | 1,503K | | | | ✓ | | |
| | Czech-CAC | 493K | | | | | | |
| | Czech-CLTT | 35K | | | | | | |
| | Danish | 100K | | | | ✓ | | |
| | Dutch | 209K | | – | | | | |
| | Dutch-LassySmall | 98K | | – | | | | |
| | English | 254K | | | | ✓ | | |
| | English-ESL | 97K | | | | | | |
| | English-LinES | 82K | | | ✓ | | | |
| | Estonian | 234K | | – | ✓ | ✓ | | |
| | Finnish | 181K | | | ✓ | ✓ | | |
| | Finnish-FTB | 159K | | – | ✓ | ✓ | | |
| | French | 390K | | | ✓ | ✓ | | |
| | Galician | 138K | | | ✓ | | | |
| | German | 293K | | – | | ✓ | | |
| | Gothic | 56K | | – | | ✓ | | |
| | Greek | 59K | | | | ✓ | | |
| | Hebrew | 115K | | – | | ✓ | | |
| | Hindi | 351K | | – | | ✓ | | |
| | Hungarian | 42K | | | | ✓ | | |
| | Indonesian | 121K | | – | | ✓ | | |
| | Irish | 23K | | | ✓ | ✓ | | |
| | Italian | 252K | | | ✓ | ✓ | | |
| | Japanese-KTC | 267K | | | | ✓ | | |
| | Kazakh | 4K | | | | | | |
| | Korean | – | | – | – | – | | |
| | Latin | 47K | | – | | ✓ | | |
| | Latin-ITTB | 291K | | – | | ✓ | | |
| | Latin-PROIEL | 165K | | – | | ✓ | | |
| | Latvian | 20K | | – | | | | |
| | Norwegian | 311K | | | | ✓ | | |
| | Old Church Slavonic | 57K | | – | | ✓ | | |
| | Persian | 151K | | | ✓ | ✓ | | |
| | Polish | 83K | | – | | ✓ | | |
| | Portuguese | 226K | | – | | ✓ | | |
| | Portuguese-BR | 298K | | – | | | | |
| | Romanian | 145K | | | ✓ | ✓ | | |
| | Russian | 99K | | | ✓ | ✓ | | |
| | Russian-SynTagRus | 1,032K | | | ✓ | | | |
| | Slovenian | 140K | | – | | ✓ | | |
| | Slovenian-SST | 29K | | | | | | |
| | Spanish | 423K | | | ✓ | ✓ | | |
| | Spanish-AnCora | 547K | | | ✓ | ✓ | | |
| | Swedish | 96K | | | ✓ | ✓ | | |
| | Swedish-LinES | 79K | | | ✓ | | | |
| | Tamil | 8K | | – | | ✓ | | |
| | Turkish | 56K | | | | | | |
| | Ukrainian | – | | – | ✓ | – | | |

# Research Problems (examples)

- Annotation
  - Produce guidelines and/or annotation for language X
  - Study the annotation of construction Y across languages

- Parsing
  - Develop and/or evaluate a parser for language X
  - Study cross-lingual transfer learning and/or annotation projection

- Typology
  - Study word order patterns in language X
  - Compare the realisation of construction Y across languages

- Applications
  - Investigate UD-based features for your favorite NLP problem

# Questions?