

# Universal Dependencies

Where do we stand and where do we want to go from here?

Daniel Zeman

📅 July 3, 2023



Charles University  
Faculty of Mathematics and Physics  
Institute of Formal and Applied Linguistics



unless otherwise stated

# What is Universal Dependencies?

“Universal Dependencies (UD) is a project developing cross-linguistically consistent treebank annotation for many languages, with the goal of facilitating multilingual parser development, cross-lingual learning, and parsing research from a language typology perspective.”

(<https://universaldependencies.org/introduction.html>)

# What is Universal Dependencies?

“Universal Dependencies (UD) is a project developing cross-linguistically consistent treebank annotation for many languages, with the goal of facilitating multilingual parser development, cross-lingual learning, and parsing research from a language typology perspective.”

(<https://universaldependencies.org/introduction.html>)

A project – but also an annotation scheme, a data repository and a community.

# What is Universal Dependencies?

“Universal Dependencies (UD) is a project developing cross-linguistically consistent treebank annotation for many languages, with the goal of facilitating multilingual parser development, cross-lingual learning, and parsing research from a language typology perspective.”

(<https://universaldependencies.org/introduction.html>)

A project – but also an **annotation scheme**, a data repository and a community.

# What is Universal Dependencies?

“Universal Dependencies (UD) is a project developing cross-linguistically consistent treebank annotation for many languages, with the goal of facilitating multilingual parser development, cross-lingual learning, and parsing research from a language typology perspective.”

(<https://universaldependencies.org/introduction.html>)

A project – but also an **annotation scheme**, a **data repository** and a community.

# What is Universal Dependencies?

“Universal Dependencies (UD) is a project developing cross-linguistically consistent treebank annotation for many languages, with the goal of facilitating multilingual parser development, cross-lingual learning, and parsing research from a language typology perspective.”

(<https://universaldependencies.org/introduction.html>)

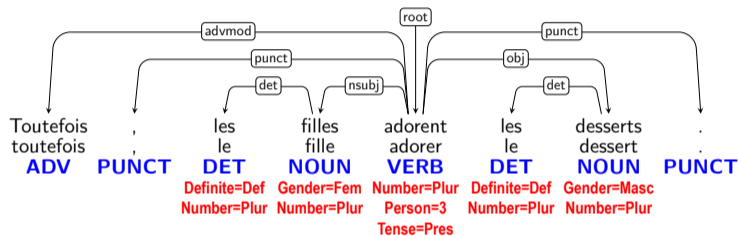
A project – but also an **annotation scheme**, a **data repository** and a **community**.

# The UD Annotation Scheme

Basic design principles:

- Words as basic units
- Binary syntactic relations

Built on de facto standards



J. Nivre, M.-C. de Marneffe, F. Ginter, Y. Goldberg, J. Hajič, C. D. Manning, R. McDonald, S. Petrov, S. Pyysalo, N. Silveira, R. Tsarfaty, and D. Zeman (2016) Universal Dependencies v1: A Multilingual Treebank Collection. In Proceedings of LREC, 1659–1666.

J. Nivre, M.-C. de Marneffe, F. Ginter, J. Hajič, C. D. Manning, S. Pyysalo, S. Schuster, F. Tyers, and D. Zeman (2020) Universal Dependencies v2: An Evergrowing Multilingual Treebank Collection. In Proceedings of LREC, 4034–4043.

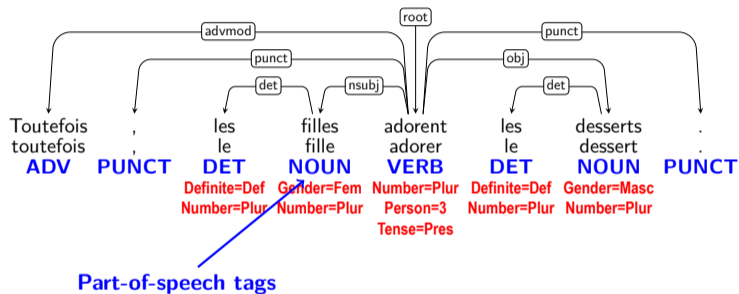
M.-C. de Marneffe, C. D. Manning, J. Nivre, and D. Zeman (2021) Universal Dependencies. Computational Linguistics, 47(2):255–308, 2021.

# The UD Annotation Scheme

Basic design principles:

- Words as basic units
- Binary syntactic relations

Built on de facto standards



J. Nivre, M.-C. de Marneffe, F. Ginter, Y. Goldberg, J. Hajič, C. D. Manning, R. McDonald, S. Petrov, S. Pyysalo, N. Silveira, R. Tsarfaty, and D. Zeman (2016) Universal Dependencies v1: A Multilingual Treebank Collection. In Proceedings of LREC, 1659–1666.

J. Nivre, M.-C. de Marneffe, F. Ginter, J. Hajič, C. D. Manning, S. Pyysalo, S. Schuster, F. Tyers, and D. Zeman (2020) Universal Dependencies v2: An Evergrowing Multilingual Treebank Collection. In Proceedings of LREC, 4034–4043.

M.-C. de Marneffe, C. D. Manning, J. Nivre, and D. Zeman (2021) Universal Dependencies. Computational Linguistics, 47(2):255–308, 2021.

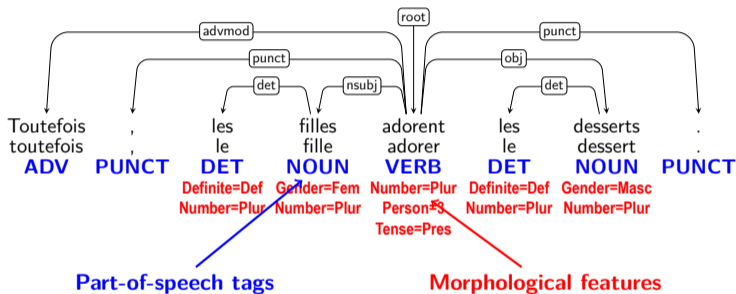


# The UD Annotation Scheme

Basic design principles:

- Words as basic units
- Binary syntactic relations

Built on de facto standards



J. Nivre, M.-C. de Marneffe, F. Ginter, Y. Goldberg, J. Hajič, C. D. Manning, R. McDonald, S. Petrov, S. Pyysalo, N. Silveira, R. Tsarfaty, and D. Zeman (2016) Universal Dependencies v1: A Multilingual Treebank Collection. In Proceedings of LREC, 1659–1666.

J. Nivre, M.-C. de Marneffe, F. Ginter, J. Hajič, C. D. Manning, S. Pyysalo, S. Schuster, F. Tyers, and D. Zeman (2020) Universal Dependencies v2: An Evergrowing Multilingual Treebank Collection. In Proceedings of LREC, 4034–4043.

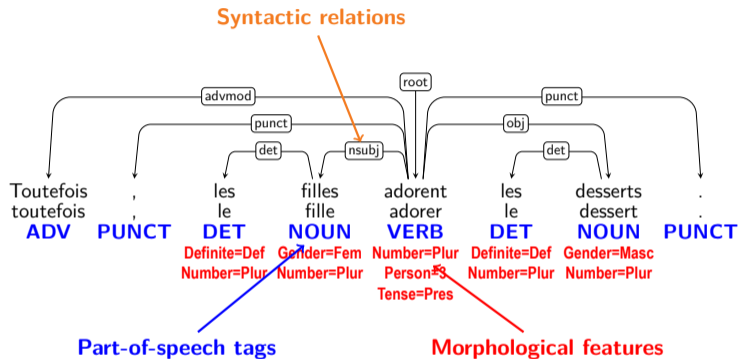
M.-C. de Marneffe, C. D. Manning, J. Nivre, and D. Zeman (2021) Universal Dependencies. Computational Linguistics, 47(2):255–308, 2021.

# The UD Annotation Scheme

Basic design principles:

- Words as basic units
- Binary syntactic relations

Built on de facto standards



J. Nivre, M.-C. de Marneffe, F. Ginter, Y. Goldberg, J. Hajič, C. D. Manning, R. McDonald, S. Petrov, S. Pyysalo, N. Silveira, R. Tsarfaty, and D. Zeman (2016) Universal Dependencies v1: A Multilingual Treebank Collection. In Proceedings of LREC, 1659–1666.

J. Nivre, M.-C. de Marneffe, F. Ginter, J. Hajič, C. D. Manning, S. Pyysalo, S. Schuster, F. Tyers, and D. Zeman (2020) Universal Dependencies v2: An Evergrowing Multilingual Treebank Collection. In Proceedings of LREC, 4034–4043.

M.-C. de Marneffe, C. D. Manning, J. Nivre, and D. Zeman (2021) Universal Dependencies. Computational Linguistics, 47(2):255–308, 2021.

## UD v2.12:

- 30 language families
- 141 languages
- 245 treebanks
- 553 contributors
- 1.8 million sentences
- 31 million words

## Current UD Languages

Information about language families (and genera for families with multiple branches) is mostly taken from [WALS Online](#) (IE = Indo-European).

▶		Abaza	1	<1K		Northwest Caucasian
▶		Afrikaans	1	49K		IE, Germanic
▶		Akkadian	2	25K		Afro-Asiatic, Semitic
▶		Akuntsu	1	1K		Tupian, Tupari
▶		Albanian	1	<1K		IE, Albanian
▶		Amharic	1	10K		Afro-Asiatic, Semitic
▶		Ancient Greek	2	416K		IE, Greek
▶		Ancient Hebrew	1	39K		Afro-Asiatic, Semitic
▶		Apurina	1	<1K		Arawakan
▶		Arabic	3	1,042K		Afro-Asiatic, Semitic
▶		Armenian	2	94K		IE, Armenian
▶		Assyrian	1	<1K		Afro-Asiatic, Semitic
▶		Bambara	1	13K		Mande
▶		Basque	1	121K		Basque
▶		Beja	1	<1K		Afro-Asiatic, Cushitic
▶		Belarusian	1	305K		IE, Slavic
▶		Bengali	1	<1K		IE, Indic
▶		Bhojpuri	1	6K		IE, Indic
▶		Bororo	1	<1K		Bororoan
▶		Breton	1	10K		IE, Celtic
▶		Bulgarian	1	156K		IE, Slavic
▶		Buryat	1	10K		Mongolic
▶		Cantonese	1	13K		Sino-Tibetan
▶		Catalan	1	553K		IE, Romance
▶		Cebuano	1	1K		Austronesian, Central Philippine
▶		Chinese	6	287K		Sino-Tibetan
▶		Chukchi	1	6K		Chukotko-Kamchatkan
▶		Classical Chinese	1	433K		Sino-Tibetan
▶		Coptic	1	55K		Afro-Asiatic, Egyptian
▶		Croatian	1	199K		IE, Slavic
▶		Czech	5	2,247K		IE, Slavic
▶		Danish	1	100K		IE, Germanic
▶		Dutch	2	306K		IE, Germanic
▶		English	10	726K		IE, Germanic
▶		Erzya	1	20K		Uralic, Mordvin

## Treebank Users

### Treebank Developers

#### Universal Guidelines Group



# Word Segmentation

*Let's go to the sea.*

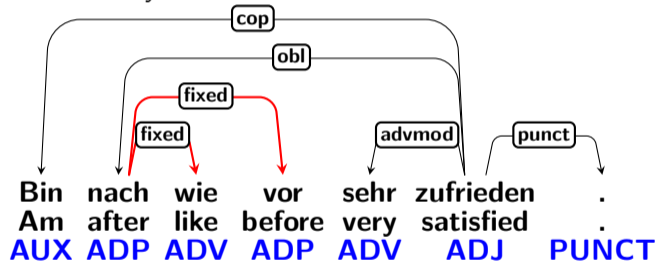
Vámonos al mar .      Vamos nos a el mar .  
VERB? X NOUN PUNCT    VERB PRON ADP DET NOUN PUNCT

- **Syntactic word** vs. orthographic word
- **Multi-word tokens**
- Two-level scheme:
  - Tokenization (low level, punctuation, concatenative)
  - Word segmentation (higher level, not necessarily concatenative)

# Fixed Expressions

One syntactic word spans several orthographic words?

*I am still very satisfied.*



# Part-of-Speech Tags

<http://universaldependencies.org/u/pos/index.html>

Open		Closed		Other	
<b>NOUN</b>	common noun	<b>PRON</b>	pronoun	<b>PUNCT</b>	punctuation
<b>PROPN</b>	proper noun	<b>DET</b>	determiner	<b>SYM</b>	symbol
<b>VERB</b>	verb	<b>AUX</b>	auxiliary	<b>X</b>	unknown
<b>ADJ</b>	adjective	<b>NUM</b>	numeral		
<b>ADV</b>	adverb	<b>ADP</b>	adposition		
<b>INTJ</b>	interjection	<b>SCONJ</b>	subordinator		
		<b>CCONJ</b>	coordinator		
		<b>PART</b>	particle		

- Taxonomy of 17 universal POS tags
- All languages use the same inventory
  - Not all tags have to be used by all languages
  - Need extensions? Use features!

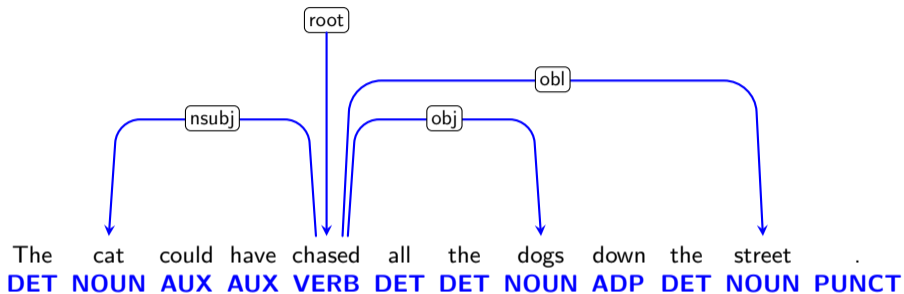
# Features

Lexical	Inflectional (“Nominal”)	Inflectional (“Verbal, Pronominal”)
PronType	Gender	VerbForm
NumType	Animacy	Mood
Poss	NounClass	Tense
Reflect	Number	Aspect
Foreign	Case	Voice
	Definite	Evident
	Degree	Polarity
Abbr		Person
Typo		Polite
		Clusivity

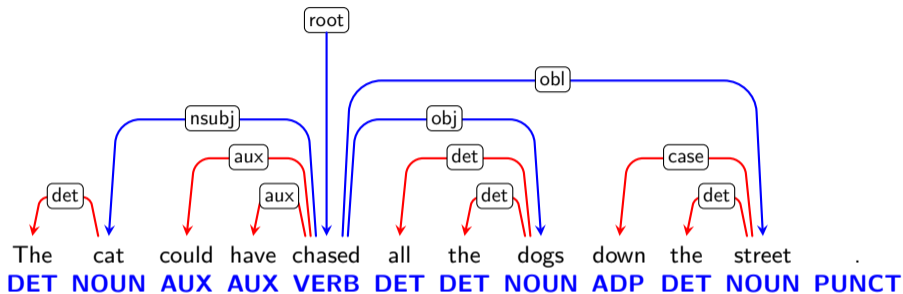
- 24 features, each with a number of possible *values*
- Languages select relevant features
- May add language-specific features or values



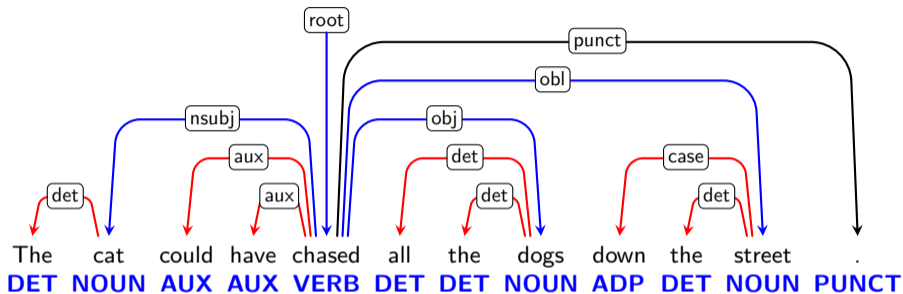
The cat could have chased all the dogs down the street .  
DET NOUN AUX AUX VERB DET DET NOUN ADP DET NOUN PUNCT



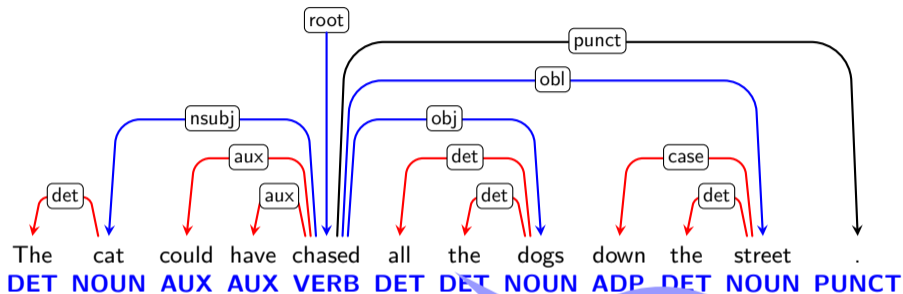
- Content words are related by dependency relations



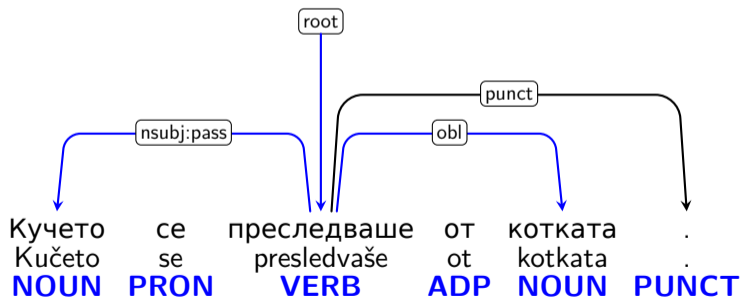
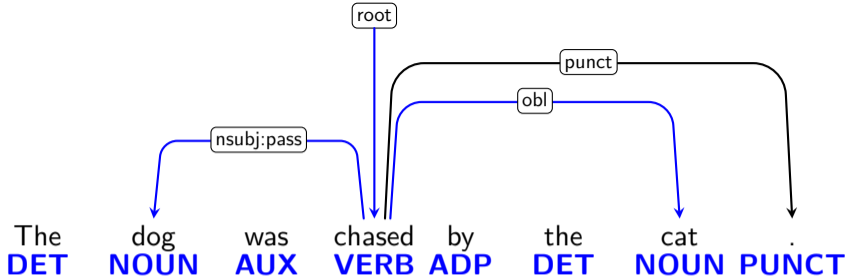
- Content words are related by dependency relations
- Function words attach to closest content words

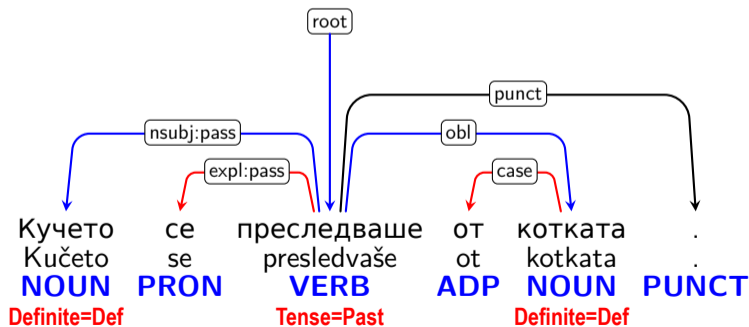
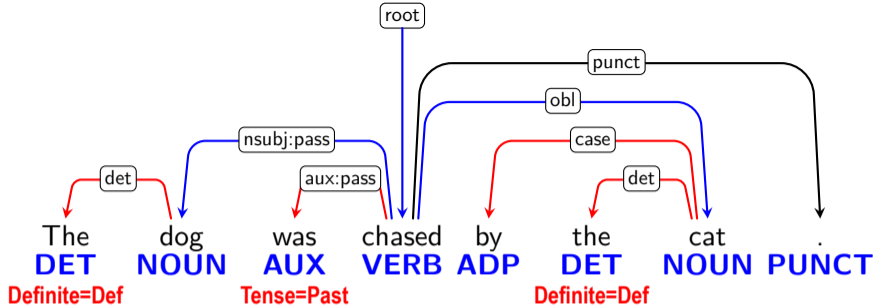


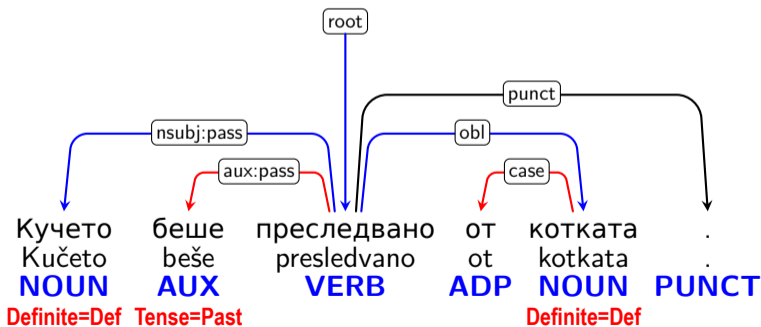
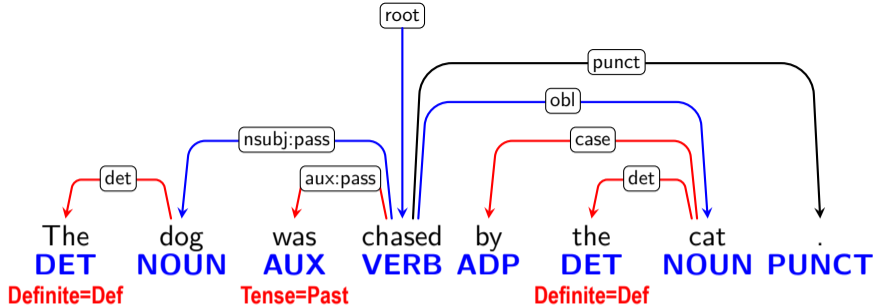
- Content words are related by dependency relations
- Function words attach to closest content words
- Punctuation attach to head of phrase or clause



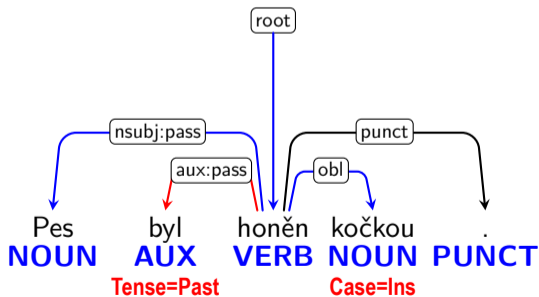
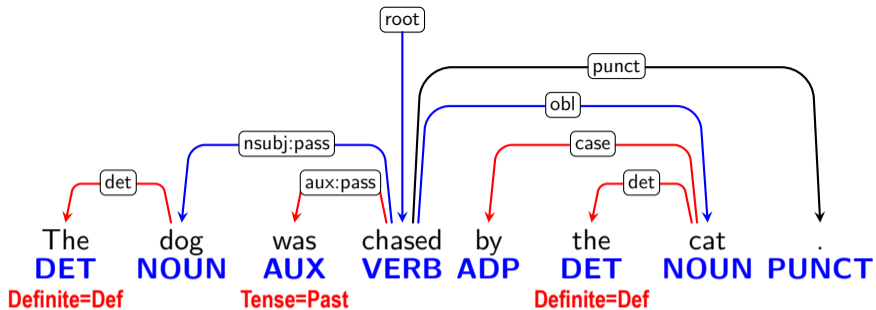
Not  
"dependency"  
in the strictly  
syntactic  
sense!











## Dependents of Clauses (Verbal or Not)

	<b>Nominal</b>	<b>Clausal</b>	<b>Modifier</b>	<b>Function</b>
<b>Core</b>	<b>nsubj</b>	csubj		
<b>Non-Core</b>	<b>obl</b> vocative dislocated expl	advcl	advmod discourse	aux cop mark

## Dependents of Verbs, Adjectives and Adverbs

	<b>Nominal</b>	<b>Clausal</b>	<b>Modifier</b>
<b>Core</b>	<b>obj</b> <b>iobj</b>	ccomp xcomp	
<b>Non-Core</b>	<b>obl</b> expl	advcl	advmod

## Dependents of Nominals

	<b>Nominal</b>	<b>Clausal</b>	<b>Modifier</b>	<b>Function</b>
	nmod appos	acl	amod nummod	det case

## Dependents of Nominals

### Nominal

nmod

appos

compound

flat

### Clausal

acl

### Modifier

amod

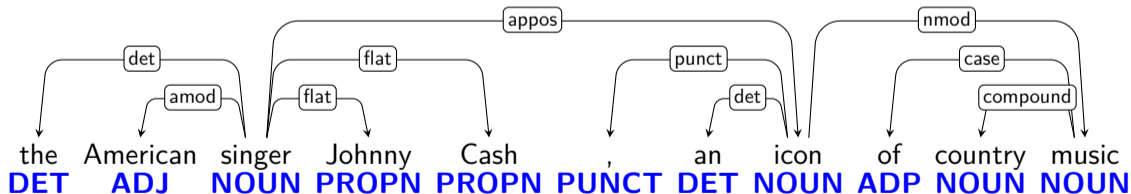
nummod

### Function

det

case

clf



# Language-specific Relation Subtypes

- Language-specific relations are **subtypes** of universal relations added to capture important phenomena
- Subtyping permits us to “back off” to universal relations

## Language-specific Relation Subtypes

<b>Relation</b>	<b>Explanation</b>
acl:relcl	Relative clause (the boy <b>who lived</b> )
compound:prt	Verb particle (dress <b>up</b> )
nmod:poss	Possessive nominal ( <b>Mary 's</b> book)
obl:agent	Agent in passive (saved <b>by the bell</b> )
cc:preconj	Preconjunction ( <b>both</b> ... and)
det:predet	Predeterminer ( <b>all</b> those ...)

# Use in Digital Humanities



# Linguists Can Search Treebanks

<https://lindat.mff.cuni.cz/services/pmltq/>

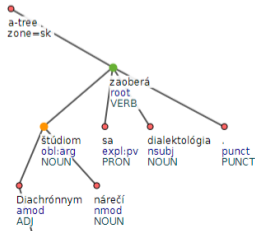
Relations Node Types Attributes Operators Functions

```
a-node $v := [  
  tag="VERB",  
  child a-node $o := [deprel="obl:arg", iset/case="ins", &empty; child a-node [deprel="case"]]  
];
```

Execute query w/o Filters Suggest (0)

Result: 3 / 100

[sk] Diachrónnym a synchrónnym štúdiom nárečí sa zaoberá dialektológia.





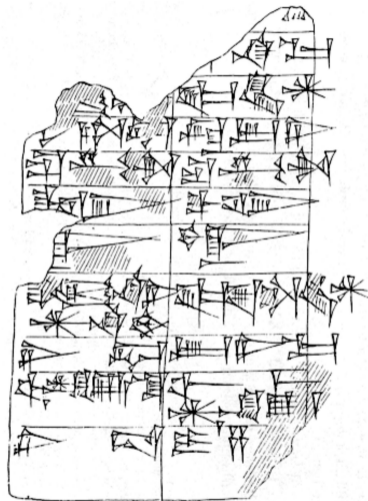
- Check grammar usage in the corpus
- Learner corpora



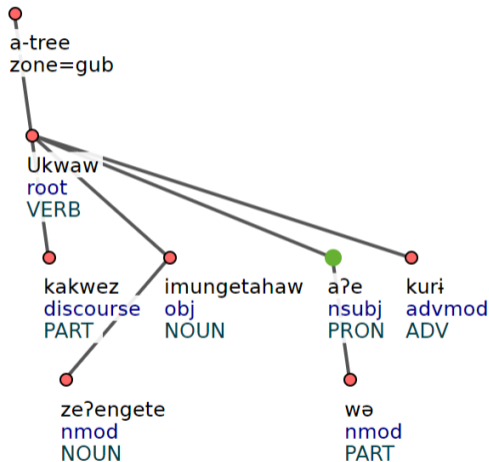


# Historical Linguistics, Classical Languages

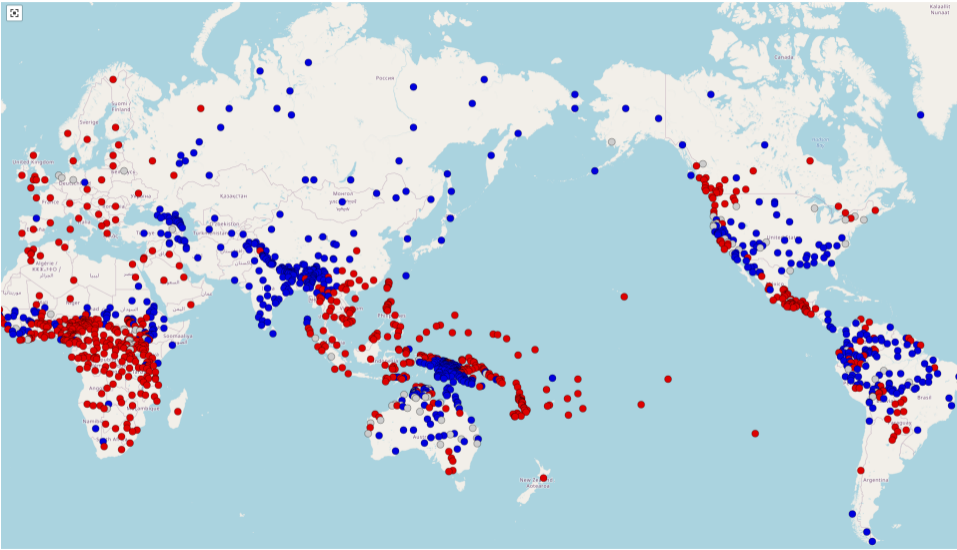
- Old Turkish
- Classical Chinese
- Sanskrit
- Akkadian
- Ancient Hebrew
- Coptic
- Ancient Greek
- Latin
- Old French
- Old Irish
- Gothic
- Old Church Slavonic



# Documentation of Endangered Languages



# Linguistic Typology



# Linguistic Typology

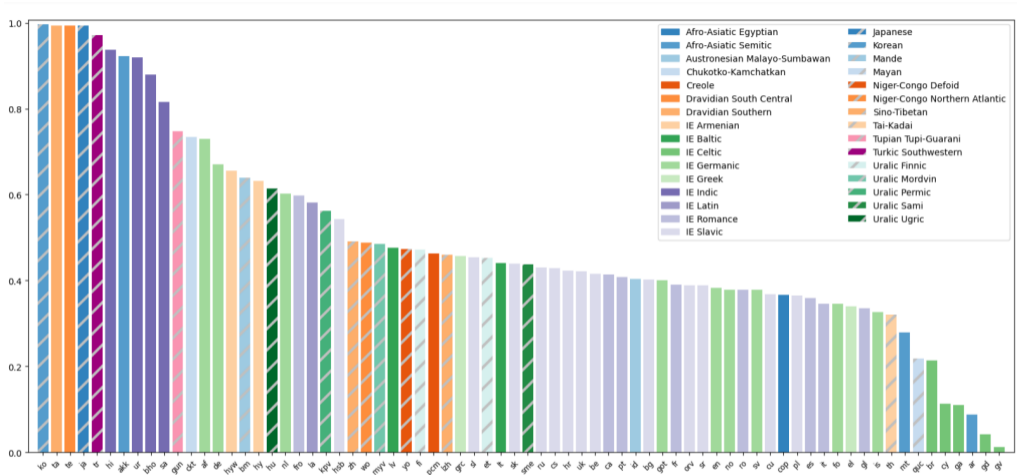


Figure 7 Percentage of head-final dependencies. Each bar is one language.

# So Where Are We Now?

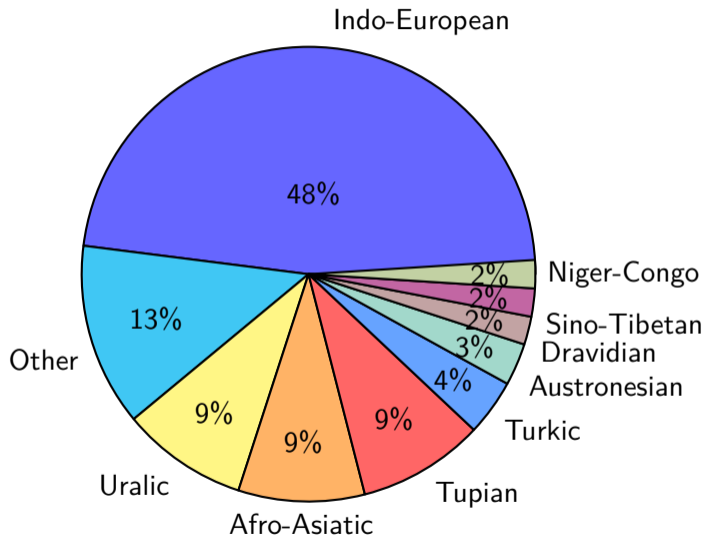
## Achievements:

- Amazing growth in terms of contributors and treebanks
- Data sets widely used for research in NLP and linguistics

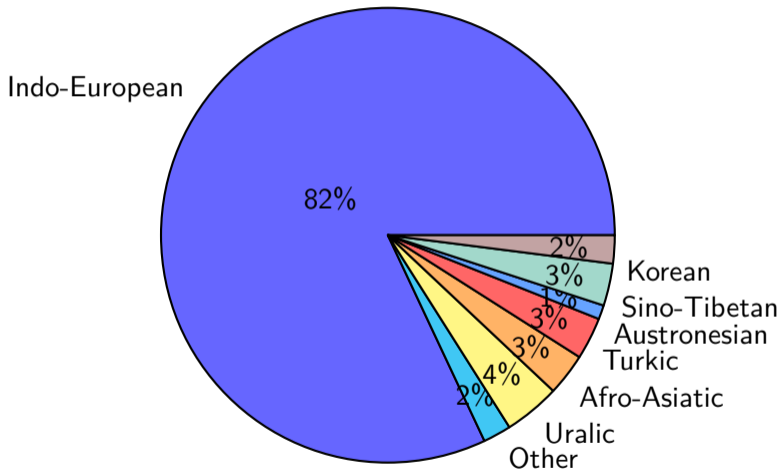
## Challenges:

- Annotation scheme and documentation can be improved
- Data sets are uneven in terms of quality and quantity
- Extensions: what to do with demand to annotate extra stuff?

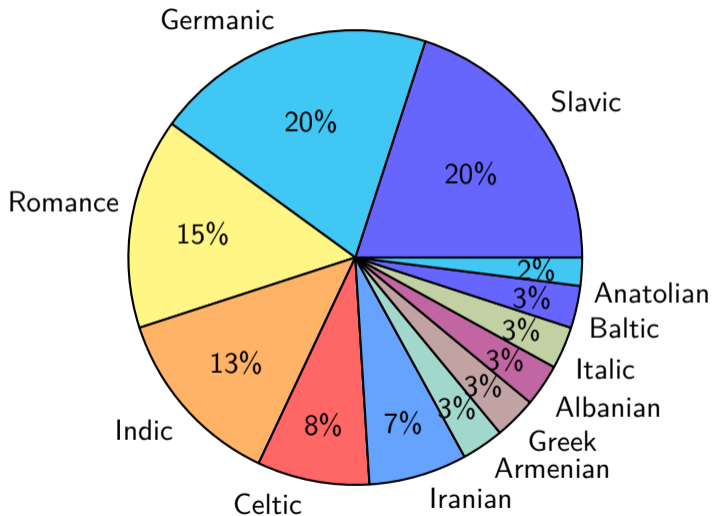
# Languages per Families



# Words per Families

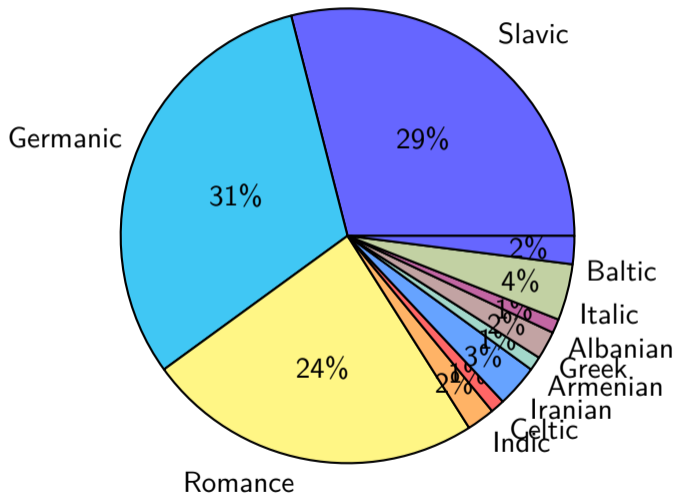


# Languages per Indo-European Genuses





# Words per Indo-European Genuses



- Enhanced Universal Dependencies



- Enhanced Universal Dependencies
- Deep Universal Dependencies



- Enhanced Universal Dependencies
- Deep Universal Dependencies
- Universal Proposition Banks



- Enhanced Universal Dependencies
- Deep Universal Dependencies
- Universal Proposition Banks
- CorefUD
- Named Entities



Universal  
**PropBank**



- Enhanced Universal Dependencies
- Deep Universal Dependencies
- Universal Proposition Banks
- CorefUD
- Named Entities
- Surface Syntactic UD (SUD)



Universal  
**PropBank**



- Enhanced Universal Dependencies
- Deep Universal Dependencies
- Universal Proposition Banks
- CorefUD
- Named Entities
- Surface Syntactic UD (SUD)
- PARSEME



Universal  
**PropBank**



P A R S E M E



- Enhanced Universal Dependencies
- Deep Universal Dependencies
- Universal Proposition Banks
- CorefUD
- Named Entities
- Surface Syntactic UD (SUD)
- PARSEME
- UniMorph+UD, sub- and superwords



Universal  
**PropBank**



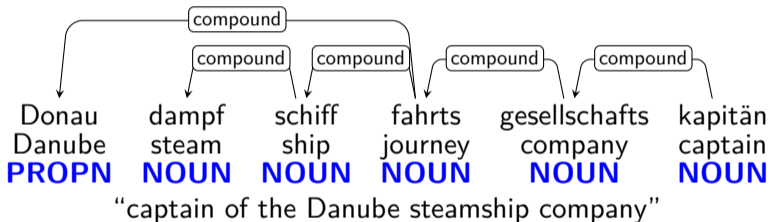
P A R S E M E





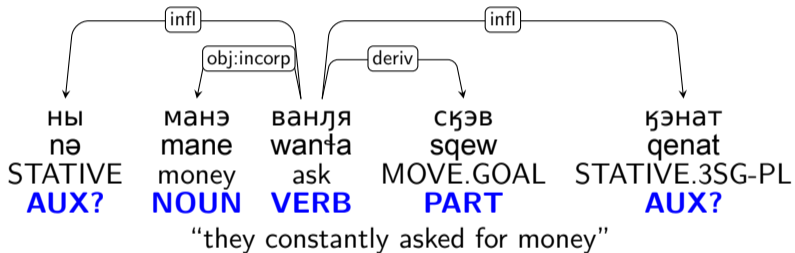


## Subword Relations: Compounds

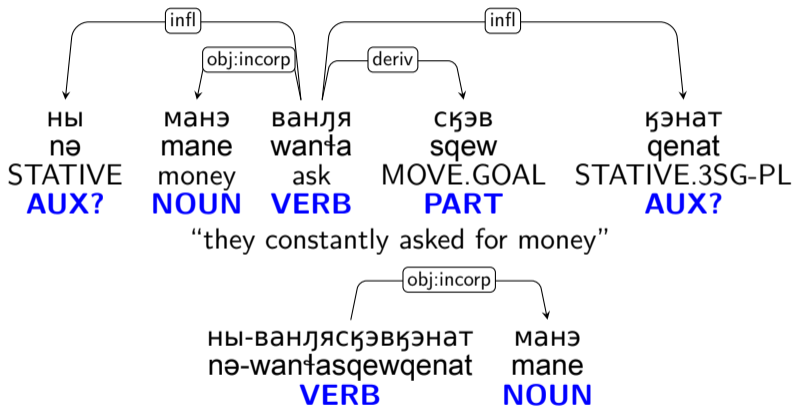


PUD sentence id w05003050

- [en] Under his administration, **the female vote** was approved.
- [sv] Under hans administration godkändes **den kvinnliga rösträtten**.
- [de] Unter seiner Regierung wurde **das Frauenwahlrecht** eingeführt.
- [es] En su gestión se decretó **el sufragio femenino**.
  - sufragio = suffrage
- [tr] Yönetimi sırasında **kadınların oy kullanması** kabul edilmiş.
  - = Administration during **women-of vote to-use** acceptance been-taken
- [cs] Pod jeho správou směly **volit ženy**.
  - (Pod jeho správou bylo zavedeno **volební právo pro ženy**.)

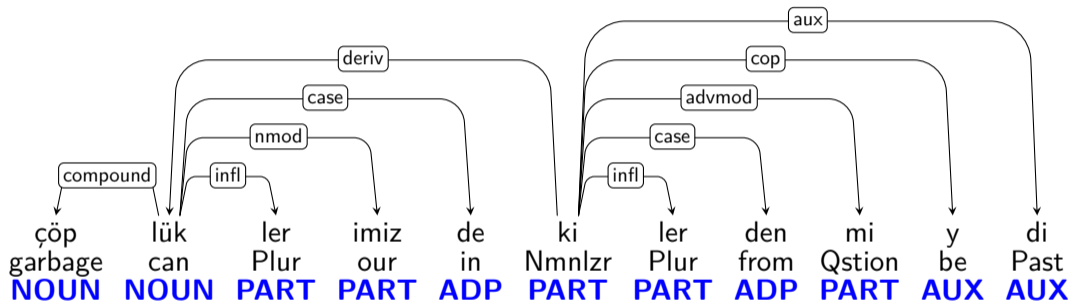


# Subword Relations: Incorporation





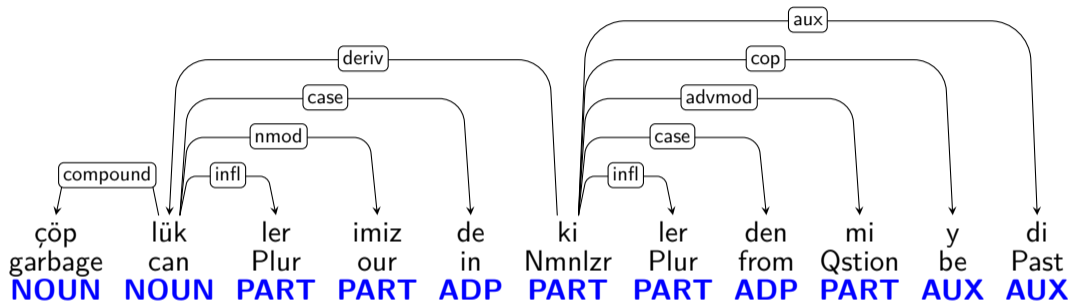
# Subword Relations: Agglutinative Languages



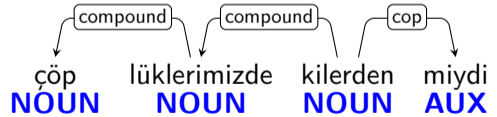
“was it from those that were in our garbage cans?”



# Subword Relations: Agglutinative Languages

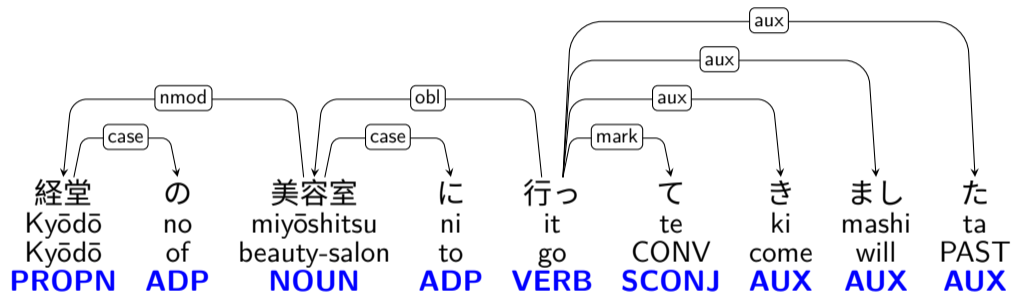


"was it from those that were in our garbage cans?"





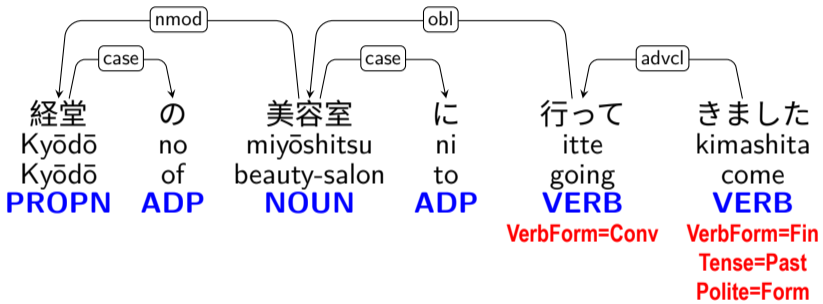
# Subword Relations: Long and Short Words in Japanese



"I went to the beauty salon of Kyōdō [, Beyond-R.]"



# Subword Relations: Long and Short Words in Japanese

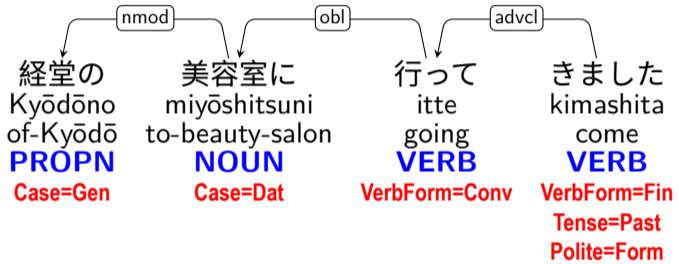


"I went to the beauty salon of Kyōdō [ , Beyond-R.]"





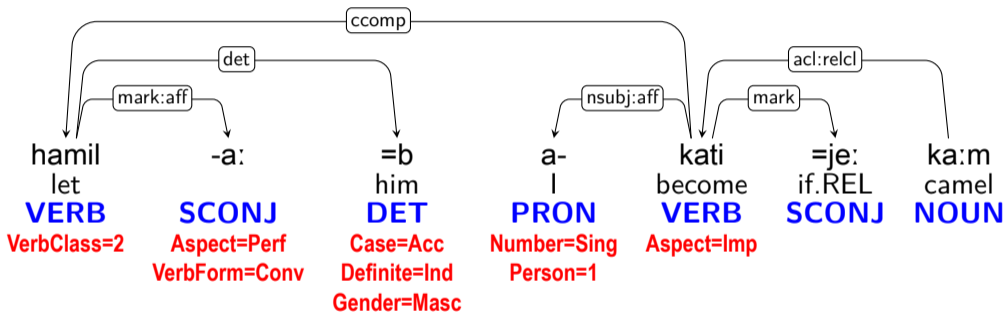
# Subword Relations: Long and Short Words in Japanese



"I went to the beauty salon of Kyōdō [, Beyond-R.]"



# Subword Relations: Fieldwork Glosses



“a camel that I had let loose”



# Superword Features

- In UD, morphological features apply to individual words.
  - Future tense in Spanish and German: no **Tense=Fut** in German!

Dormirá  
He-will-sleep  
**VERB**

VerbForm=Fin  
Mood=Ind  
Tense=Fut  
Number=Sing  
Person=3

Er  
He  
**PRON**

PronType=Prs  
Number=Sing  
Person=3  
Gender=Masc  
Case=Nom

wird  
will  
**AUX**

VerbForm=Fin  
Mood=Ind  
Tense=Pres  
Number=Sing  
Person=3

schlafen  
sleep  
**VERB**

VerbForm=Inf



# Superword Features

- In UD, morphological features apply to individual words.
  - Future tense in Spanish and German: no **Tense=Fut** in German!
- Can we also add features to “phrases”?

Dormirá  
He-will-sleep

**VERB**

**VerbForm=Fin**

**Mood=Ind**

**Tense=Fut**

**Number=Sing**

**Person=3**

Er

He

**PRON**

**PronType=Prs**

**Number=Sing**

**Person=3**

**Gender=Masc**

**Case=Nom**

wird schlafen

will sleep

**VERB**

**VerbForm=Fin**

**Mood=Ind**

**Tense=Fut**

**Number=Sing**

**Person=3**

## Summary

- Universal Dependencies
  - Unified annotation for all languages
    - Language-specific extensions
  - Content words higher than function words ... better parallelism
  - Words as the main unit (tree nodes)
- Extension of UD: above and below word level
  - UD-like relations between subword units
  - UD-like features for groups of words (“phrases”)

<https://universaldependencies.org/>

Thanks!  
Hvala!

<https://universaldependencies.org/>