

Universal Dependencies: Principles and Tools

Daniel Zeman

📅 July 1, 2022



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

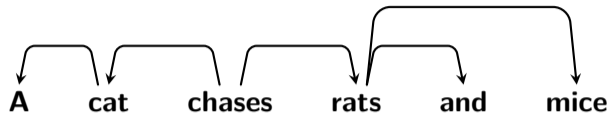


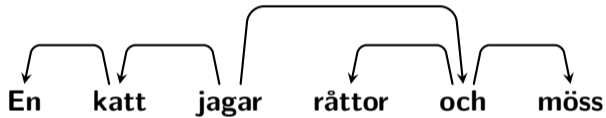
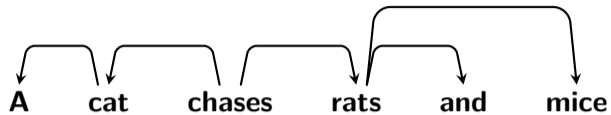
unless otherwise stated

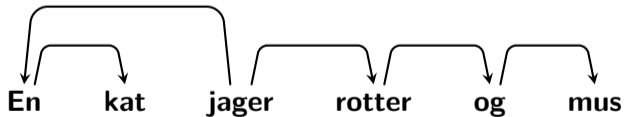
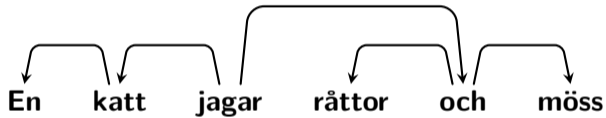
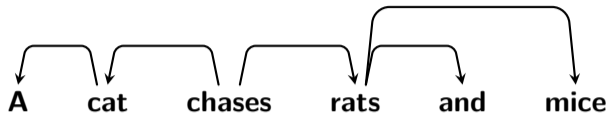
Introduction

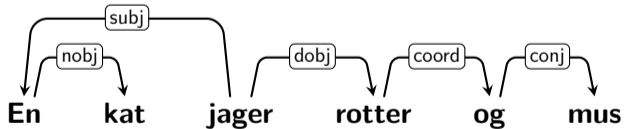
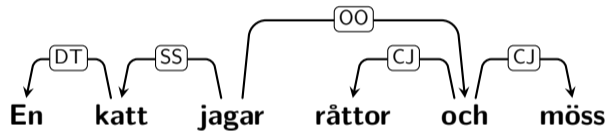
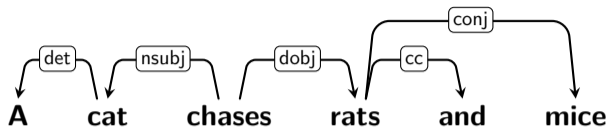
- 1 Introduction
- 2 Morphological Annotation in UD
- 3 Syntactic Annotation in UD
- 4 Core vs. Oblique
- 5 Enhanced Universal Dependencies
- 6 UD Tools

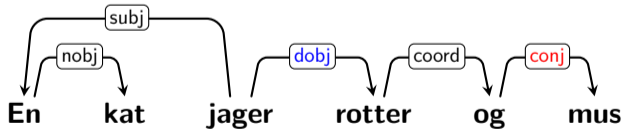
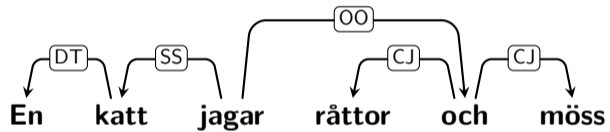
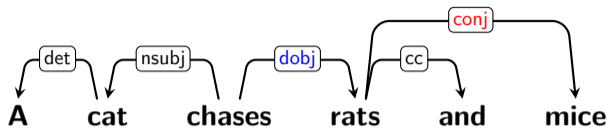
- Around 2010:
- Increasing interest in multilingual NLP
 - Multilingual evaluation campaigns to test generality
 - Cross-lingual learning to support low-resource languages
- Increasing awareness of methodological problems
 - Current NLP relies heavily on annotation
 - Annotation schemes vary across languages











Why was 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

- <https://universaldependencies.org/>
- Same things annotated same way across languages...
- ... while highlighting different **coding strategies**

Manning's Law

The secret to understanding UD is to realize that the design is a very subtle compromise between approximately 6 things:

- 1 UD must be satisfactory on linguistic analysis grounds for **individual languages**.



It's easy to come up with a proposal that improves UD on one of these dimensions. The interesting and difficult part is to improve UD while remaining sensitive to all these dimensions.

Manning's Law

The secret to understanding UD is to realize that the design is a very subtle compromise between approximately 6 things:

- 1 UD must be satisfactory on linguistic analysis grounds for **individual languages**.
- 2 UD must be good for linguistic **typology**, i.e., providing a suitable basis for bringing out cross-linguistic parallelism across languages and language families.



It's easy to come up with a proposal that improves UD on one of these dimensions. The interesting and difficult part is to improve UD while remaining sensitive to all these dimensions.

Manning's Law

The secret to understanding UD is to realize that the design is a very subtle compromise between approximately 6 things:

- 1 UD must be satisfactory on linguistic analysis grounds for **individual languages**.
- 2 UD must be good for linguistic **typology**, i.e., providing a suitable basis for bringing out cross-linguistic parallelism across languages and language families.
- 3 UD must be suitable for **rapid, consistent annotation** by a human annotator.



It's easy to come up with a proposal that improves UD on one of these dimensions. The interesting and difficult part is to improve UD while remaining sensitive to all these dimensions.

Manning's Law



The secret to understanding UD is to realize that the design is a very subtle compromise between approximately 6 things:

- 1 UD must be satisfactory on linguistic analysis grounds for **individual languages**.
- 2 UD must be good for linguistic **typology**, i.e., providing a suitable basis for bringing out cross-linguistic parallelism across languages and language families.
- 3 UD must be suitable for **rapid, consistent annotation** by a human annotator.
- 4 UD must be easily comprehended and used by a **non-linguist**, whether a language learner or an engineer with prosaic needs for language processing. ... it leads us to favor **traditional grammar** notions and terminology.

It's easy to come up with a proposal that improves UD on one of these dimensions. The interesting and difficult part is to improve UD while remaining sensitive to all these dimensions.

Manning's Law



The secret to understanding UD is to realize that the design is a very subtle compromise between approximately 6 things:

- 1 UD must be satisfactory on linguistic analysis grounds for **individual languages**.
- 2 UD must be good for linguistic **typology**, i.e., providing a suitable basis for bringing out cross-linguistic parallelism across languages and language families.
- 3 UD must be suitable for **rapid, consistent annotation** by a human annotator.
- 4 UD must be easily comprehended and used by a **non-linguist**, whether a language learner or an engineer with prosaic needs for language processing. ... it leads us to favor **traditional grammar** notions and terminology.
- 5 UD must be suitable for **computer parsing** with high accuracy.

It's easy to come up with a proposal that improves UD on one of these dimensions. The interesting and difficult part is to improve UD while remaining sensitive to all these dimensions.

Manning's Law

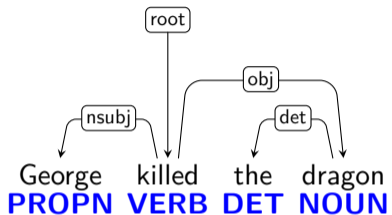


The secret to understanding UD is to realize that the design is a very subtle compromise between approximately 6 things:

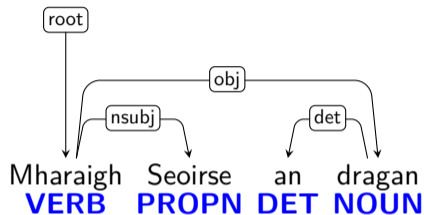
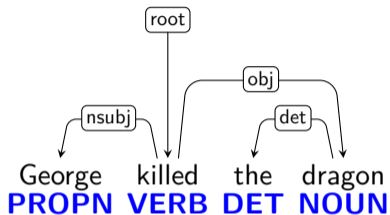
- 1 UD must be satisfactory on linguistic analysis grounds for **individual languages**.
- 2 UD must be good for linguistic **typology**, i.e., providing a suitable basis for bringing out cross-linguistic parallelism across languages and language families.
- 3 UD must be suitable for **rapid, consistent annotation** by a human annotator.
- 4 UD must be easily comprehended and used by a **non-linguist**, whether a language learner or an engineer with prosaic needs for language processing. ... it leads us to favor **traditional grammar** notions and terminology.
- 5 UD must be suitable for **computer parsing** with high accuracy.
- 6 UD must support well **downstream language understanding tasks** (relation extraction, reading comprehension, machine translation, ...)

It's easy to come up with a proposal that improves UD on one of these dimensions. The interesting and difficult part is to improve UD while remaining sensitive to all these dimensions.

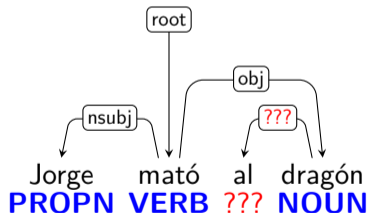
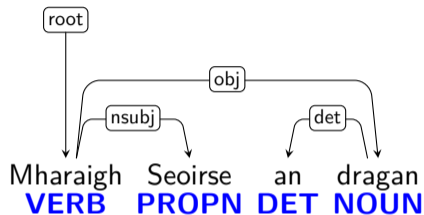
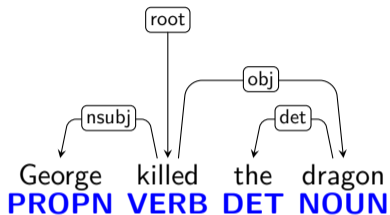
Same Thing Same Way



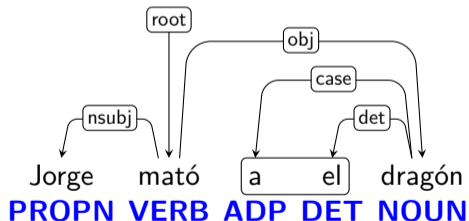
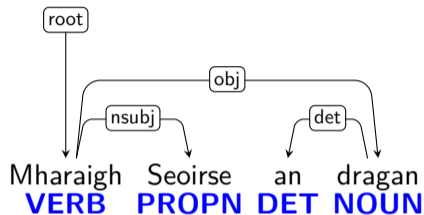
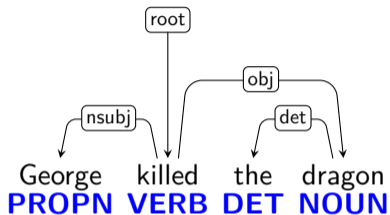
Same Thing Same Way



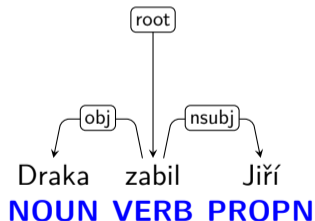
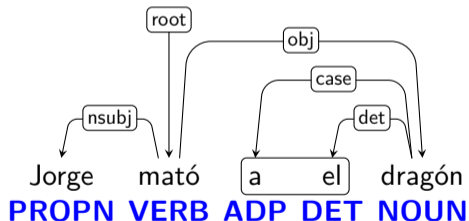
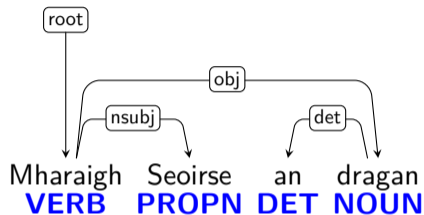
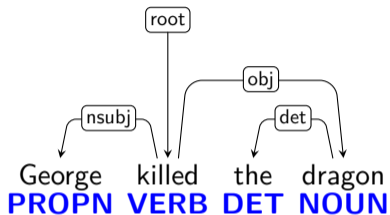
Same Thing Same Way



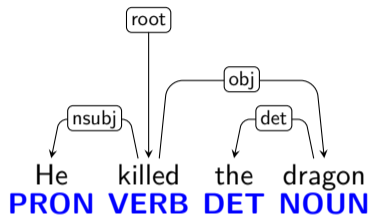
Same Thing Same Way



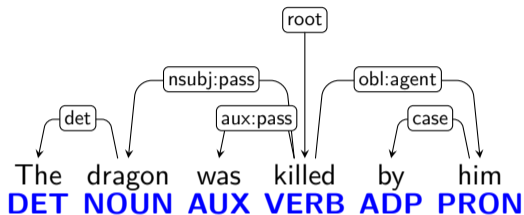
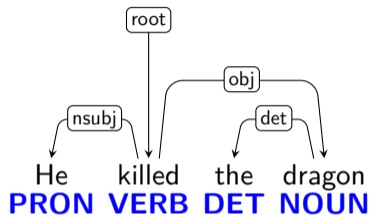
Same Thing Same Way



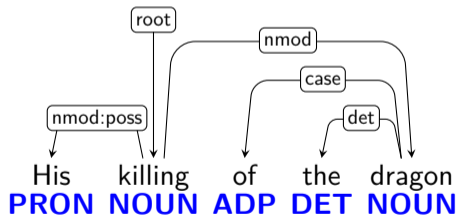
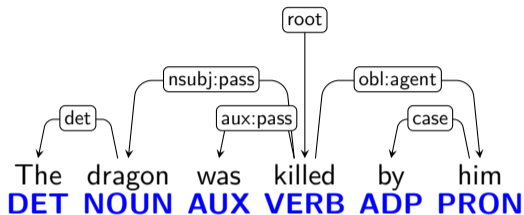
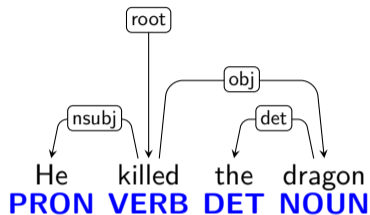
Same Meaning \neq Same Construction!



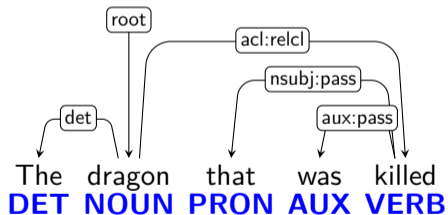
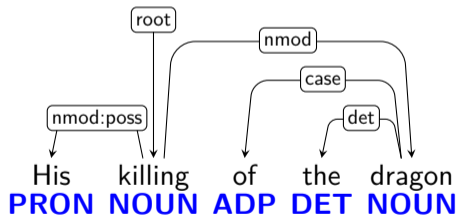
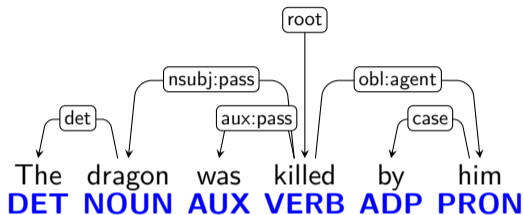
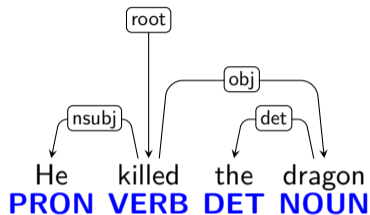
Same Meaning \neq Same Construction!



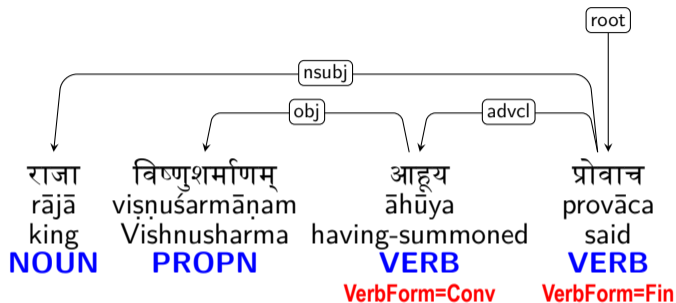
Same Meaning \neq Same Construction!



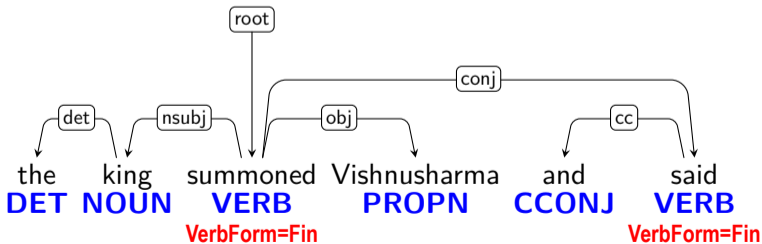
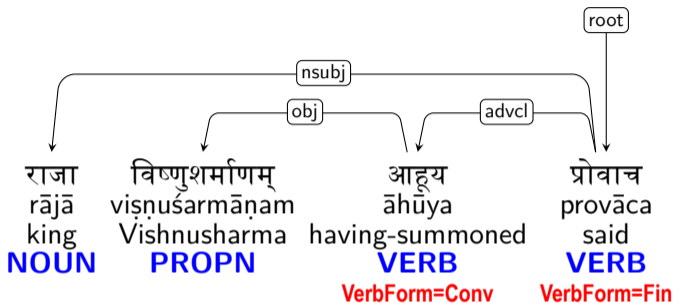
Same Meaning \neq Same Construction!



Language-specific Preferences



Language-specific Preferences



Morphological Annotation

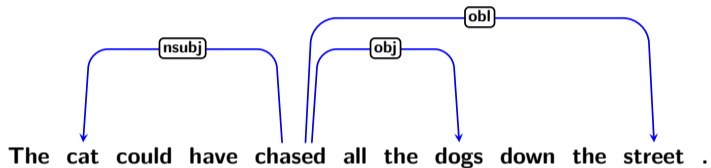
Le	chat	chasse	les	chiens	.
le	chat	chasser	le	chien	.
DET	NOUN	VERB	DET	NOUN	PUNCT
Definite=Def Gender=Masc Number=Sing	Gender=Masc Number=Sing	Mood=Ind Number=Sing Person=3 Tense=Pres VerbForm=Fin	Definite=Def Gender=Masc Number=Plur	Gender=Masc Number=Plur	

- Lemma representing the semantic content of a word
- Part-of-speech tag representing its grammatical class
- Features representing lexical and grammatical properties of the lemma or the particular word form

The cat could have chased all the dogs down the street .

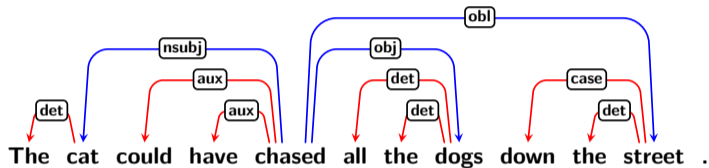
- Content words are related by dependency relations
- Function words attach to the content word they modify
- Punctuation attach to head of phrase or clause

Syntactic Annotation



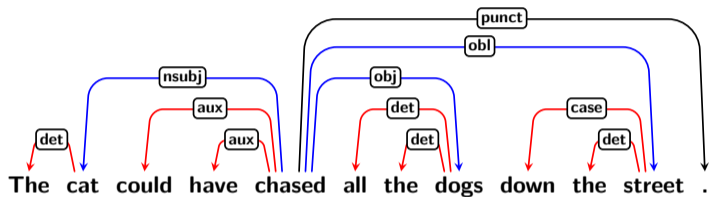
- Content words are related by dependency relations
- Function words attach to the content word they modify
- Punctuation attach to head of phrase or clause

Syntactic Annotation



- Content words are related by dependency relations
- Function words attach to the content word they modify
- Punctuation attach to head of phrase or clause

Syntactic Annotation



- Content words are related by dependency relations
- Function words attach to the content word they modify
- Punctuation attach to head of phrase or clause

CoNLL-U Format

ID	FORM	LEMMA	UPOS	XPOS	FEATS	HEAD	DEPREL	DEPS	MISC
1	Le	le	DET	-	-	2	det	-	-
2	chat	chat	NOUN	-	-	3	nsubj	-	-
3	boit	boire	VERB	-	-	0	root	-	-
4-5	du	-	-	-	-	-	-	-	-
4	de	de	ADP	-	-	6	case	-	-
5	le	le	DET	-	-	6	det	-	-
6	lait	lait	NOUN	-	-	3	obj	-	SpaceAfter=No
7	.	.	PUNCT	-	-	3	punct	-	-

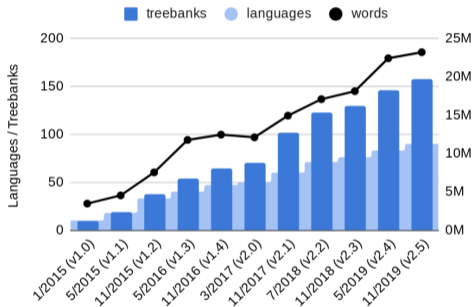
- Revised and extended version of CoNLL-X format
- Two-level segmentation and enhanced dependencies

Basic Universal Dependencies: 130 (128) Languages and Growing

▪ **I.-E.:**  Armenian (+West),  Greek (+Ancient),  Albanian,  Hittite,  Breton,  Irish,  Manx,  Scottish,  Welsh,  Afrikaans,  Danish,  Dutch,  English,  Faroese,  Frisian,  German,  Gothic,  Icelandic,  Low Saxon,  Norwegian,  Swedish,  Swiss German,  Catalan,  French,  Galician,  Italian,  Latin,  Ligurian,  Neapolitan,  Old French,  Portuguese,  Romanian,  Spanish,  Umbrian,  Belarusian,  Bulgarian,  Church Slavonic,  Croatian,  Czech,  Old Russian,  Polish,  Pomak,  Russian,  Serbian,  Slovak,  Slovenian,  Ukrainian,  Upper Sorbian,  Latvian,  Lithuanian,  Kurmanji,  Persian,  Khunsari,  Nayini,  Soi,  Urdu,  Hindi,  Kangri,  Bhojpuri,  Bengali,  Marathi,  Sanskrit ▪ **Dravidian:**  Tamil,  Telugu ▪ **Uralic:**  Erzya,  Estonian,  Finnish,  Hungarian,  Karelian,  Livvi,  Komi Permyak+Zyrian,  Moksha,  Sámi North+Skolt ▪ **Turkic:**  Kazakh,  Old Turkish,  Tatar,  Turkish,  Uyghur,  Yakut ▪  Buryat ▪  Xibe ▪  Korean ▪  Japanese ▪ **Sino-T.:**  Cantonese,  Classical Chinese,  Chinese ▪ **Tai-Kadai:**  Thai ▪ **Aus.-As.:**  Vietnamese ▪ **Austron.:**  Indonesian,  Javanese,  Tagalog,  Cebuano ▪ **Pama-Nyu.:**  Warlpiri ▪ **Chu.-Kam.:**  Chukchi ▪ **Esk.-Al.:**  Yupik ▪ **Mayan:**  Kiche ▪ **Arawakan:**  Apurinã ▪ **Arawan:**  Madi ▪ **Tupian:**  Akuntsu,  Guajajara,  Kaapor,  Makurap,  Mundurukú,  Tupinambá,  Mbyá,  Guaraní,  Teko ▪ **Af.-As.:**  Akkadian,  Amharic,  Arabic Standard+Levantine,  Assyrian,  Beja,  Coptic,  Hebrew (+Ancient),  Maltese ▪ **Niger-Congo:**

Where are we today?

- Brief history of UD:
 - First guidelines launched in October 2014
 - Treebank releases (roughly) **every six months**
 - Version 2 guidelines/treebanks in 2016–2017
 - New: guideline amendments since May 2022
 - Extensions: MWEs, PropBanks, Coreference
- UD in numbers:
 - 130 languages
 - 228 treebanks
 - 502 contributors
 - 150,000+ downloads
- Past and current UD events:
 - 4 CoNLL and IWPT shared tasks on UD parsing
 - UD workshops: next in Washington 2023
 - COST action: UniDive (since 2022)
 - Next release in November 2022 (v2.11)



Use in Digital Humanities



Linguists Can Search Treebanks

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

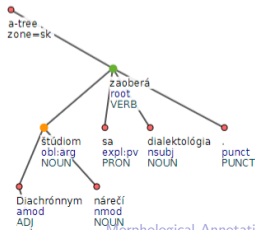
Relations Node Types Attributes Operators Functions

```
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.



Linguists Can Parse and Search New Data

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

The screenshot shows the Czech Wikipedia page for COVID-19. The main text discusses the identification of the virus in China in January 2020 and its subsequent spread. A sidebar on the left contains navigation links, and a right sidebar provides a classification table for the disease.

Classification Table:

Klasifikace	
MKN-10	U07.1 a U07.2
Statistické údaje – obě pohlaví	
Incidence	230 567 044 ^[1] (z toho 0 ^[1] uzdravených) ke dni 22. září 2021
Mortalita	2,23 % ^[20497] (celosvětový průměr: 2,44%) ^[20497]

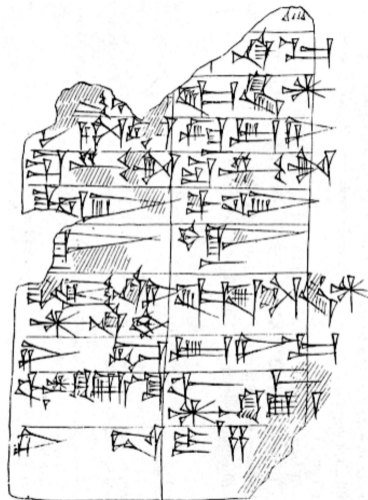
The screenshot shows the UDPipe web interface. It displays a syntactic tree for the sentence "První případ byl identifikován v čínském Wu-chanu v prosinci 2019." The tree is rooted at "root" and branches into "identifikován" (root ADJ), "případ" (případ nsubj pass NOUN), "byl" (byl aux:pass AUX), "v" (v case ADP), "čínském" (čínském amod ADJ), "Wu-chanu" (Wu amod ADJ), "v" (v case ADP), "prosinci" (prosinci obl NOUN), and "2019" (2019 nummod NUM). The interface also includes buttons for "Process Input", "Output Text", "Show Table", "Show Trees", and "Save Tree as SVG".

- Check grammar usage in the corpus
- Learner corpora

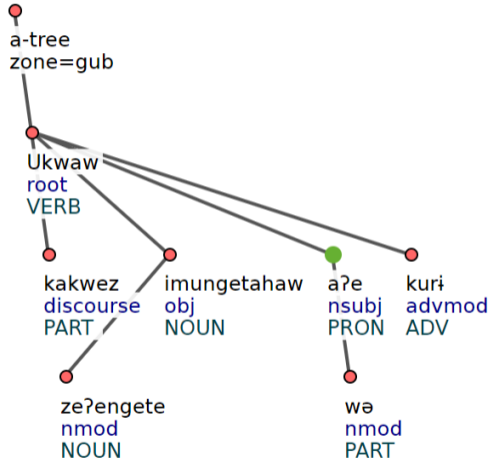


Historical Linguistics, Classical Languages

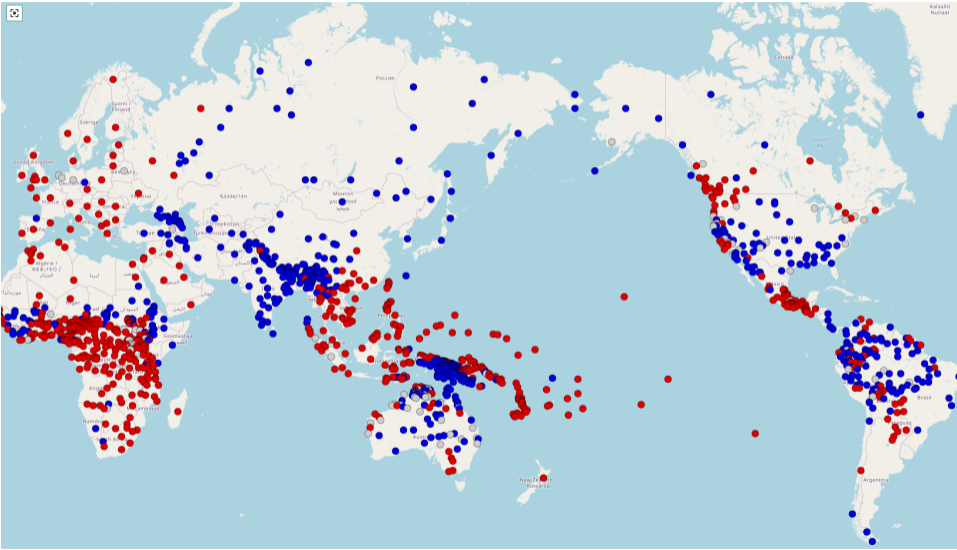
- Old Turkish
- Classical Chinese
- Sanskrit
- Hittite
- Akkadian
- Coptic
- Ancient Hebrew
- Ancient Greek
- Latin
- Old French
- Gothic
- Old Church Slavonic
- Old East Slavic



Documentation of Endangered Languages



Linguistic Typology



Linguistic Typology

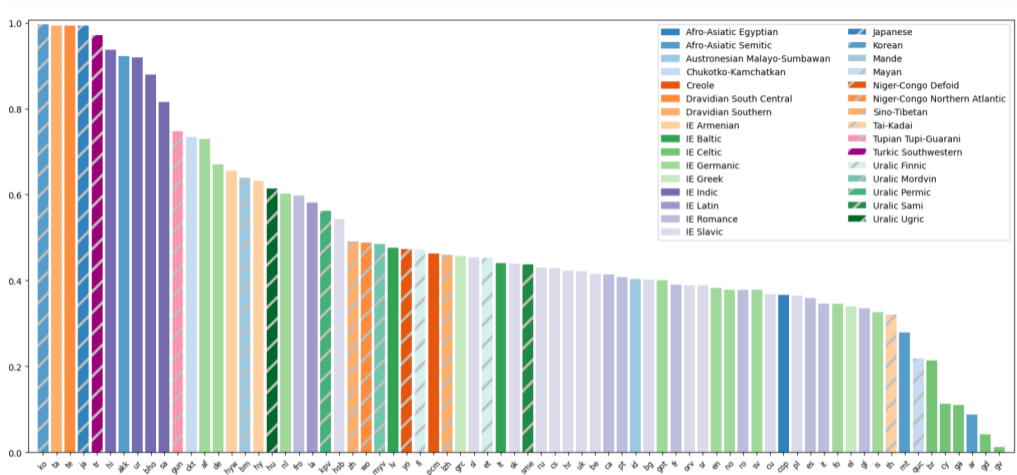


Figure 7 Percentage of head-final dependencies. Each bar is one language.

Morphological Annotation in UD

- 1 Introduction
- 2 Morphological Annotation in UD**
- 3 Syntactic Annotation in UD
- 4 Core vs. Oblique
- 5 Enhanced Universal Dependencies
- 6 UD Tools

- Tokenization / word segmentation
- Lemmatization (**LEMMA**)
- Universal part-of-speech tags (**UPOS**)
- Universal features (**FEATS**)
- Language-specific features

“María, I love you!” Juan exclaimed.

«¡María, te amo!», exclamó Juan.
X PRON X VERB X

« ¡ María , te amo ! » ,
PUNCT PUNCT PROPN PUNCT PRON VERB PUNCT PUNCT PUNCT

- Classic tokenization:
 - Separate punctuation from words
 - Recognize certain clusters of symbols like “...”
 - Perhaps keep together things like `user@mail.x.edu`

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)

- Lexicalist hypothesis:
 - Words (not morphemes) are the basic units in syntax
 - Words enter in dependency relations
 - Words are forms of lemmas and have morphological features

- Orthographic vs. syntactic word
 - Syntactically autonomous part of orthographic word
 - Contractions (*al = a + el*)
 - Clitics (*vámonos = vamos + nos*)
 - ¿A qué hora *nos vamos* mañana?
 - *Nos* despertamos a las cinco.
“We wake up at five.”
 - *Nuestro guía nos* despierta a las cinco.
“Our guide wakes us up at five.”

Contractions in Arabic (EXTRA)

He abdicated in favour of his son Baudouin.

يتنازل	عن	العرش	لابنه	بودوان
yatanāzalu	ʿan	al-ʿarši	li+ibni+hi	būdūān
surrendered	on	the throne	to son his	Baudouin
VERB	ADP	NOUN	ADP+NOUN+PRON	PROPN

We are now in Valencia.

現在我們在瓦倫西亞。

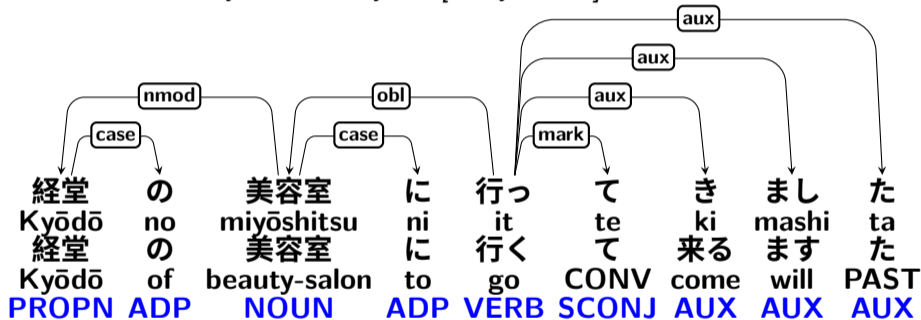
Xiànzài wǒ men zài wǎ lún xī yǎ.

We are now in Valencia.

現在	我們	在	瓦倫西亞	。
Xiànzài	wǒmen	zài	Wǎlúnxīyǎ	.
Now	we	in	Valencia	.
ADV	PRON	ADP	PROPN	PUNCT

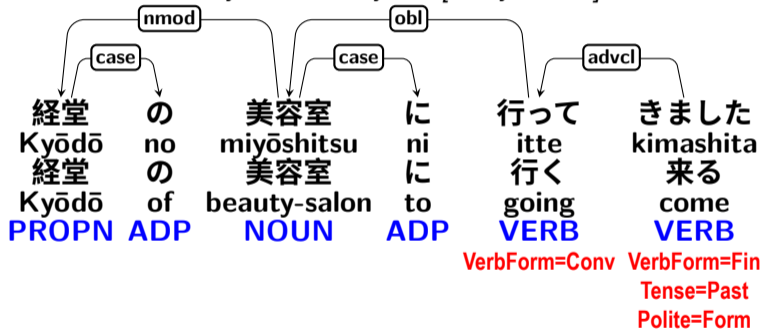
Words in Japanese (EXTRA)

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



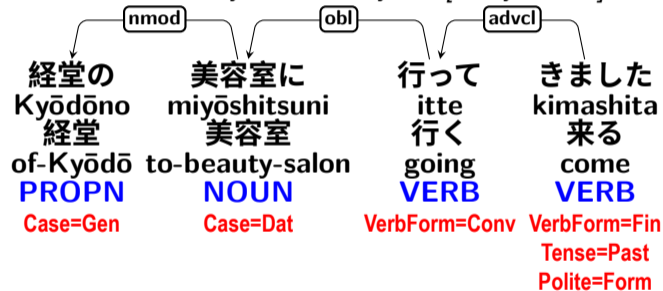
Words in Japanese (EXTRA)

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



Words in Japanese (EXTRA)

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



Vietnamese: Words with Spaces

All the concrete country roads are the result of...

Tất cả	đường	bê tông	nội đồng	là	thành quả	...
All	road	concrete	country	is	achievement	...
PRON	NOUN	NOUN	NOUN	AUX	NOUN	PUNCT

- Spaces delimit monosyllabic morphemes, not words.
- Multiple syllables without space occur in loanwords (*bê tông*).
- Spaces are allowed to occur word-internally in Vietnamese UD.

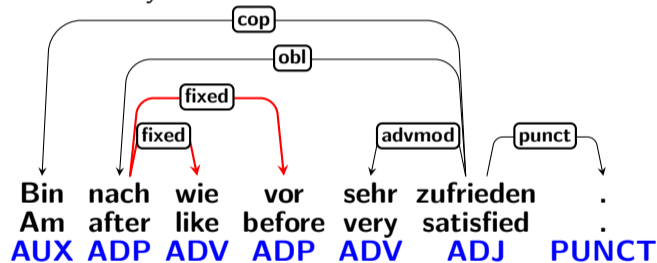
Numbers with Spaces

#	text = Il touche environ 100 000 sesterces par an.						
1	Il	il	PRON	...	2	nsubj	--
2	touche	toucher	VERB	...	0	root	--
3	environ	environ	ADV	...	4	advmod	--
4	100 000	100 000	NUM	...	5	nummod	--
5	sesterces	sesterce	NOUN	...	2	obj	--
6	par	par	ADP	...	7	case	--
7	an	an	NOUN	...	2	obl	_ SpaceAfter=No
8	.	.	PUNCT	...	2	punct	--

Fixed Expressions

One syntactic word spans several orthographic words?

I am still very satisfied.



- When to split?
 - Only part of the token involved in a relation to something outside the token? Split!

- When to split?
 - Only part of the token involved in a relation to something outside the token? Split!
 - Hard time finding POS tag? Split!

- When to split?
 - Only part of the token involved in a relation to something outside the token? Split!
 - Hard time finding POS tag? Split!
 - Hard time finding dependency relation? Don't split!
 - Or not hard time but the relation would be compound, flat, fixed or goeswith.

- When to split?
 - Only part of the token involved in a relation to something outside the token? Split!
 - Hard time finding POS tag? Split!
 - Hard time finding dependency relation? Don't split!
 - Or not hard time but the relation would be compound, flat, fixed or goeswith.
 - Border case? Keep orthographic words (if they exist).

Word Segmentation Summary (EXTRA)

- When to split?
 - Only part of the token involved in a relation to something outside the token? Split!
 - Hard time finding POS tag? Split!
 - Hard time finding dependency relation? Don't split!
 - Or not hard time but the relation would be compound, flat, fixed or goeswith.
 - Border case? Keep orthographic words (if they exist).
- Words with spaces
 - Vietnamese writing system
 - Very restricted set of exceptions (numbers)
 - Special relations elsewhere (fixed, compound)

Recoverability: CoNLL-U Format

text = Vámonos al mar.

text_en = Let's go to the sea.

ID	FORM	LEMMA	UPOS	...	HEAD	_	MISC
1-2	Vámonos	—	—	...	—	—	—
1	Vamos	ir	VERB	...	0	root	—
2	nos	nosotros	PRON	...	1	obj	—
3-4	al	—	—	...	—	—	—
3	a	a	ADP	...	5	case	—
4	el	el	DET	...	5	det	—
5	mar	mar	NOUN	...	1	obl	SpaceAfter=No
6	.	.	PUNCT	...	1	punct	—

Recoverability: CoNLL-U Format

text = Vámonos al mar.

text_en = Let's go to the sea.

ID	FORM	LEMMA	UPOS	...	HEAD	_	MISC
1-2	Vámonos	—	—	...	—	—	— —
1	Vamos	ir	VERB	...	0	root	— —
2	nos	nosotros	PRON	...	1	obj	— —
3-4	al	—	—	...	—	—	— —
3	a	a	ADP	...	5	case	— —
4	el	el	DET	...	5	det	— —
5-6	mar.	—	—	...	—	—	— —
5	mar	mar	NOUN	...	1	obl	— —
6	.	.	PUNCT	...	1	punct	— —

Tokenization vs. Multi-word Tokens (EXTRA)

- Parallelism among closely related languages
 - ca: **informar-se** sobre el patrimoni cultural
 - es: **informarse** sobre el patrimonio cultural
 - en: *learn about cultural heritage*

- ca: L'únic que veig és => **L' únic** que veig és
- en: don't => **do n't**

- No strict guidelines for tokenization (yet)
 - UD English: **non-stop**, **post-war**: single-word tokens
 - UD Czech: **non-stop** would be split to three tokens

Tokenization vs. Multi-word Tokens Summary (EXTRA)

- Punctuation involved? Low level!
 - Exceptions: Spanish-Catalan parallelism.

Tokenization vs. Multi-word Tokens Summary (EXTRA)

- Punctuation involved? Low level!
 - Exceptions: Spanish-Catalan parallelism.
- Boundary between two letters? Typically high level.
 - Exceptions: Chinese, Japanese.

Tokenization vs. Multi-word Tokens Summary (EXTRA)

- Punctuation involved? Low level!
 - Exceptions: Spanish-Catalan parallelism.
- Boundary between two letters? Typically high level.
 - Exceptions: Chinese, Japanese.
- Non-concatenative? High level!

- Basic or citation form (\Rightarrow it is an existing word in most cases)
- Disambiguating ids, if available, go to MISC
- Derivational vs. inflectional morphology (if participles are ADJ, their lemma should not be infinitive)

within a year Algeria will become an islamic state

13	do	do	ADP	...	Lld=do-1
14	roka	rok	NOUN	...	—
15	se	se	PRON	...	LGloss=(zvr._zájmeno/částice)
16	Alžírsko	Alžírsko	PROPN	...	—
17	stane	stát	VERB	...	Lld=stát-2
18	islámským	islámský	ADJ	...	—
19	státem	stát	NOUN	...	Lld=stát-1 LGloss=(státní_útvár) SpaceAfter=No

- Basic or citation form
- Disambiguating ids, if available, go to MISC

Part-of-Speech Tags

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

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

Part-of-Speech Tags (EXTRA)

- Traditionally a mixture of morphological, syntactic/distributional and semantic/notional criteria
- Prefer grammatical > semantic criteria
 - Language-particular definition of a category
- But the **name** of the category is universal
 - Translated words: overlapping categories, but not perfect match
 - UPOS of English *dog* is **NOUN**; so is French *chien* or Russian *собака*
- Preferably POS is encoded in lexicon, not heavily usage-dependent
 - But not for incompatible syntactic functions (e.g. **PRON** vs. **SCONJ**)

Universal Features

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

- **PronType** (*pronoun type*)
- **NumType** (*numeral type*)
- **Poss** (*possessive*)
- **Reflex** (*reflexive*)
- **Foreign** (*foreign*)
- **Abbr** (*abbreviation*)
- **Typo** (*typo*)
- **Gender** (*gender*)
- **Animacy** (*animacy*)
- **NounClass** (*noun class*)
- **Number** (*number*)
- **Case** (*case*)
- **Definite** (*definiteness*)
- **Degree** (*degree of comparison*)
- **VerbForm** (*((de)verbal form*)
- **Mood** (*mood*)
- **Tense** (*tense*)
- **Aspect** (*aspect*)
- **Voice** (*voice*)
- **Evident** (*evidentiality*)
- **Polarity** (*polarity*)
- **Person** (*person*)
- **Polite** (*politeness*)
- **Clusivity** (*clusivity*)

Features (EXTRA)

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

Three types of infinitives in Finnish:

Example: *olla* “to be”

1st	2nd	3rd
olla	ollessa ollen	olemassa olemaan olemasta olemalla olematta

Language-Specific Features (EXTRA)

Joku Someone PRON	yrittää tries VERB VerbForm=Fin Mood=Ind Tense=Pres	piristää to-uplift VERB VerbForm=Inf	itseään oneself PRON	värjäämällä by-staining VERB VerbForm=Inf3 Case=Ade	hiuksensa their-hair NOUN
---------------------------------------	--	---	--	---	---

Language-Specific Features (EXTRA)

Joku	yrittää	piristää	itseään	värjäämällä	hiuksensa
Someone	tries	to-uplift	oneself	by-staining	their-hair
PRON	VERB	VERB	PRON	VERB	NOUN
	VerbForm=Fin	VerbForm=Inf		VerbForm=Inf3	
	Mood=Ind			Case=Ade	
	Tense=Pres				

Joku	yrittää	piristää	itseään	värjäämällä	hiuksensa
Someone	tries	to-uplift	oneself	by-staining	their-hair
PRON	VERB	VERB	PRON	VERB	NOUN
	VerbForm=Fin	VerbForm=Inf		VerbForm=Inf	
	Mood=Ind	<u>InfForm=1</u>		<u>InfForm=3</u>	
	Tense=Pres			Case=Ade	

Czech adjectives agree with nouns in gender.

velký
big
ADJ

bratr
brother
NOUN

Gender=Masc **Gender=Masc**

velká
big
ADJ

sestra
sister
NOUN

Gender=Fem **Gender=Fem**

Layered Features (EXTRA)

Possessive adjectives: agreement gender vs. lexical gender

otcův father's ADJ Gender=Masc Gender[psor]=Masc	bratr brother NOUN Gender=Masc	matčin mother's ADJ Gender=Masc Gender[psor]=Fem	bratr brother NOUN Gender=Masc
otcova father's ADJ Gender=Fem Gender[psor]=Masc	sestra sister NOUN Gender=Fem	matčina mother's ADJ Gender=Fem Gender[psor]=Fem	sestra sister NOUN Gender=Fem

Multi-valued Features (Disjunction / Parallel Application) (EXTRA)

- Feature can have two or more values
- Interpreted as disjunction
- Example: in some languages, many pronouns function both as interrogative and relative, but some pronouns are only relative. The former will have **PronType=Int,Rel**
- In other cases, it is desirable to disambiguate by context. Polish *którym* (form of *który* “which”) can be **Case=Ins, Loc** in singular or **Dat** in plural but we do not want to annotate **Case=Dat,Ins,Loc!**
- All values of the feature/language? Omit the feature completely! Polish: **Gender=Fem,Masc,Neut**. Spanish: **Gender=Fem,Masc**

Multi-valued Features (Serial Application) (EXTRA)

- Currently used in Turkish (language-specific values)
- Two or more morphemes in chain, affecting the same feature
- Example: **Voice=CauPass** (causative + passive => someone is caused to do something)
 - *yanıl* “be wrong”
 - *yanılmışım* **Voice=Act** “I was wrong”
 - *okuru yanılttığını* **Voice=Cau** “mislead the reader”
 - *okurlar yanıltılmıştır* **Voice=CauPass** “readers were misled”

Multi-valued Features (Serial Application) (EXTRA)

- Currently used in Turkish (language-specific values)
- Two or more morphemes in chain, affecting the same feature
- Example: **Voice=CauPass** (causative + passive => someone is caused to do something)
 - *yanıl* “be wrong”
 - *yanılmışım* **Voice=Act** “I was wrong”
 - *okuru yanılttığını* **Voice=Cau** “mislead the reader”
 - *okurlar yanıltılmıştır* **Voice=CauPass** “readers were misled”
 - Hypothetical: **Voice=PassCau** (not used in Turkish) could mean “to cause something to be done by someone”

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

Participle Types (EXTRA)

некурящий человек
nekurjaščij čelovek
non-smoking person
ADJ NOUN

VerbForm=Part

Tense=Pres

Gender=Masc

Number=Sing

Case=Nom

Gender=Masc

Number=Sing

Case=Nom

начавшийся разговор
načavšijsja razgovor
that-has-started conversation
ADJ NOUN

VerbForm=Part

Tense=Past

Gender=Masc

Number=Sing

Case=Nom

Gender=Masc

Number=Sing

Case=Nom

- Sometimes features like **Tense** help distinguish participle types
- Not the same tense as with finite verbs (reference point)
- But useful because:
 - We use known UD primitives rather than language-specific labels such as **VerbForm=PastPart**, or even **ParticType=Past**
 - Reasonably close to the grammatical meaning

Conflicting Traditional Terminologies (EXTRA)

- If possible, stay compatible with traditional grammar
- Often it is not possible: terminology conflicts
- **VerbForm=Conv** – *converb, transgressive, adverbial participle, gerund*

Conflicting Traditional Terminologies (EXTRA)

- If possible, stay compatible with traditional grammar
- Often it is not possible: terminology conflicts
- **VerbForm=Conv** – *converb*, *transgressive*, *adverbial participle*, *gerund*
- *Gerund* (**VerbForm=Ger**)
 - English: close to verbal nouns (**VerbForm=Vnoun**)
 - Spanish: more like present participle (**VerbForm=Part | Tense=Pres**)
 - Slavic: *converb* (**VerbForm=Conv**)

Conflicting Traditional Terminologies (EXTRA)

- If possible, stay compatible with traditional grammar
- Often it is not possible: terminology conflicts
- **VerbForm=Conv** – *converb*, *transgressive*, *adverbial participle*, *gerund*
- *Gerund* (**VerbForm=Ger**)
 - English: close to verbal nouns (**VerbForm=Vnoun**)
 - Spanish: more like present participle (**VerbForm=Part | Tense=Pres**)
 - Slavic: converb (**VerbForm=Conv**)
- *Aorist*
 - Ancient Greek, Turkish: neutral non-past tense (they use a language-specific value **Tense=Aor**)
 - Slavic languages: simple past tense (**Tense=Past**)

Conflicting Traditional Terminologies (EXTRA)

A	ko	so	se	leta	1942	vračali	...
And	as	they-were	REFL	in-year	1942	returning	...
CCONJ	SCONJ	AUX	PRON	NOUN	NUM	VERB	
		VerbForm=Fin Tense=Pres				VerbForm=Part <u>Tense=Past?</u>	

Conflicting Traditional Terminologies (EXTRA)

A	ko	so	se	leta	1942	vračali	...
And	as	they-were	REFL	in-year	1942	returning	...
CCONJ	SCONJ	AUX	PRON	NOUN	NUM	VERB	
		VerbForm=Fin				VerbForm=Part	
		Tense=Pres				<u>Tense=Past?</u>	

da	ne	bi	v	Atene	prišli	...
that	not	would	in	Athens	they-come	...
SCONJ	PART	AUX	ADP	PROPN	VERB	
		VerbForm=Fin			VerbForm=Part	
		Mood=Cnd			<u>Tense=Past??</u>	

Conflicting Traditional Terminologies (EXTRA)

A	ko	so	se	leta	1942	vračali	...
And	as	they-were	REFL	in-year	1942	returning	...
CCONJ	SCONJ	AUX	PRON	NOUN	NUM	VERB	
		VerbForm=Fin				VerbForm=Part	
		Tense=Pres				<u>Tense=Past?</u>	

da	ne	bi	v	Atene	prišli	...
that	not	would	in	Athens	they-come	...
SCONJ	PART	AUX	ADP	PROPN	VERB	
		VerbForm=Fin			VerbForm=Part	
		Mood=Cnd			<u>Tense=Past??</u>	

v	prihodnje	ne	bodo	vozili	zgolj	les
in	future	not	they-will	drive	just	wood
ADP	NOUN	PART	AUX	VERB	PART	NOUN
			VerbForm=Fin	VerbForm=Part		
			Tense=Fut	<u>Tense=Past???</u>		

Conflicting Traditional Terminologies (EXTRA)

- West/South Slavic: **VerbForm=Part**
- Russian: **VerbForm=Fin** (past tense)
 - **Tense=Past** useful to distinguish from other participles (especially in Bulgarian)
 - But it is also used for the conditional (any tense)
 - In Slovenian even for the future tense!

Conflicting Traditional Terminologies (EXTRA)

- West/South Slavic: **VerbForm=Part**
- Russian: **VerbForm=Fin** (past tense)
 - **Tense=Past** useful to distinguish from other participles (especially in Bulgarian)
 - But it is also used for the conditional (any tense)
 - In Slovenian even for the future tense!
- Terminology – options:
 - cs “active participle” / “past tense”
 - ru “past tense” / “finite!”
 - Active participle is something else: *нарушивший* / *наруšivšij*
 - bg “participle + past (aorist) / imperfect” (two subtypes)
 - cu “participle + resultative aspect” (lang-spec)

Conflicting Traditional Terminologies (EXTRA)

- West/South Slavic: **VerbForm=Part**
- Russian: **VerbForm=Fin** (past tense)
 - **Tense=Past** useful to distinguish from other participles (especially in Bulgarian)
 - But it is also used for the conditional (any tense)
 - In Slovenian even for the future tense!
- Terminology – options:
 - cs “active participle” / “past tense”
 - ru “past tense” / “finite!”
 - Active participle is something else: *нарушивший* / *наруšivšij*
 - bg “participle + past (aorist) / imperfect” (two subtypes)
 - cu “participle + resultative aspect” (lang-spec)
- “I-participle”
 - But that would be a language-specific verb form.

Summary of Morphology

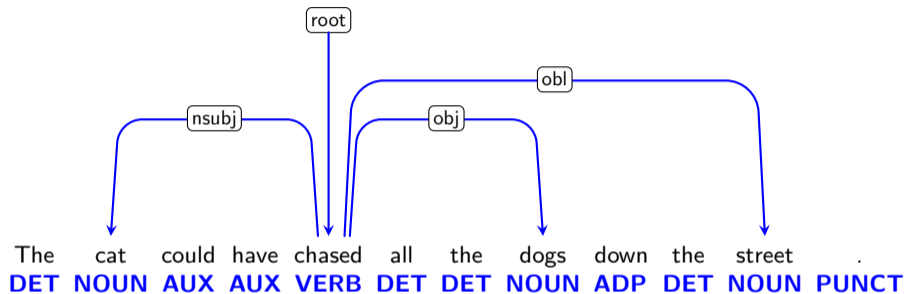
- Multi-word tokens: 1 orthographic token = N syntactic words
- Lemma = citation form of the word
- UPOS = universal part-of-speech tag (17 coarse-grained tags)
- Morphological features (feature-value pairs)
 - Universal feature-value pairs
 - Language-specific values or even features
- Lemmas, tags, and features apply to words (tree nodes), not to multi-word expressions and not to sub-word units (morphemes)
- (EXTRA)
 - Categories are **comparable** (but not identical) across languages
 - Layered features
 - Multi-valued features

<https://universaldependencies.org/>

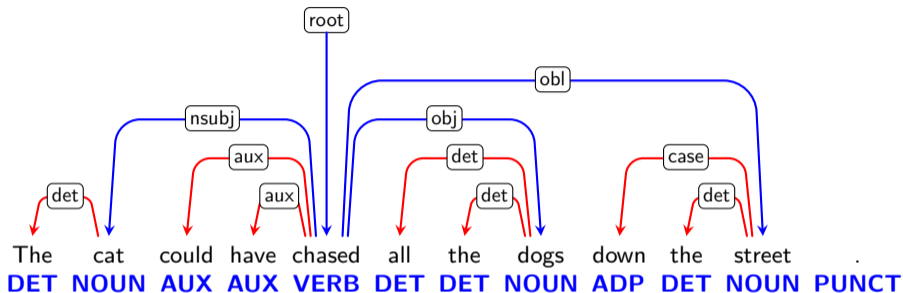
Syntactic Annotation in UD

- 1 Introduction
- 2 Morphological Annotation in UD
- 3 Syntactic Annotation in UD**
- 4 Core vs. Oblique
- 5 Enhanced Universal Dependencies
- 6 UD Tools

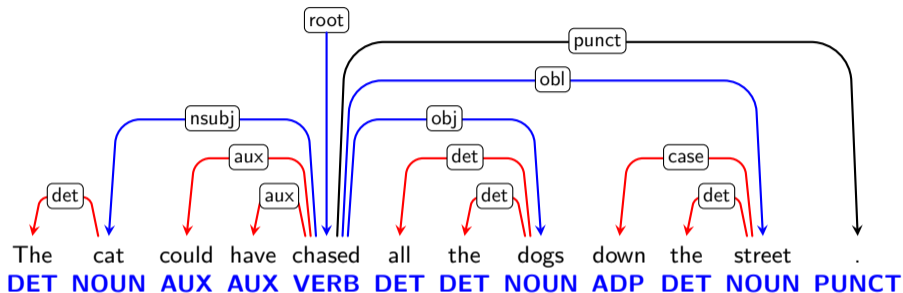
The cat could have chased all the dogs down the street .
DET NOUN AUX AUX VERB DET DET NOUN ADP DET NOUN PUNCT



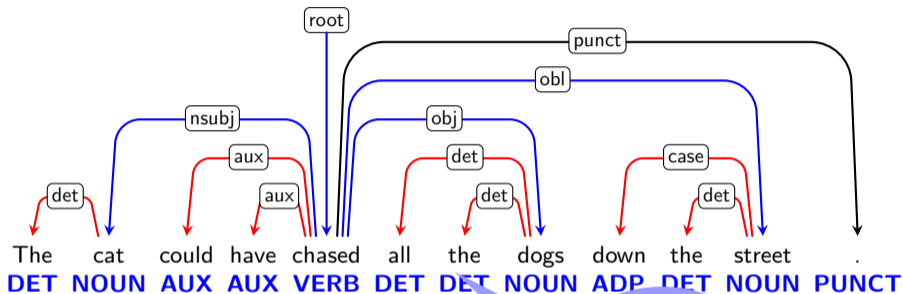
- Content words are related by dependency relations



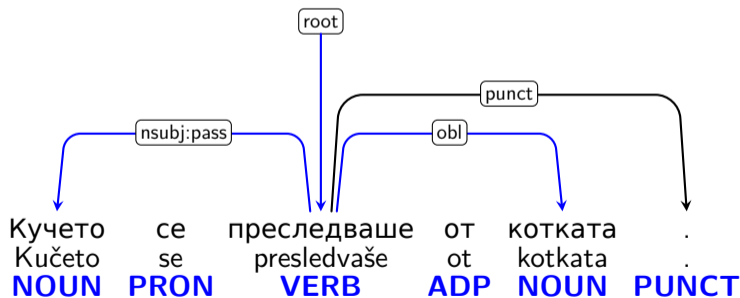
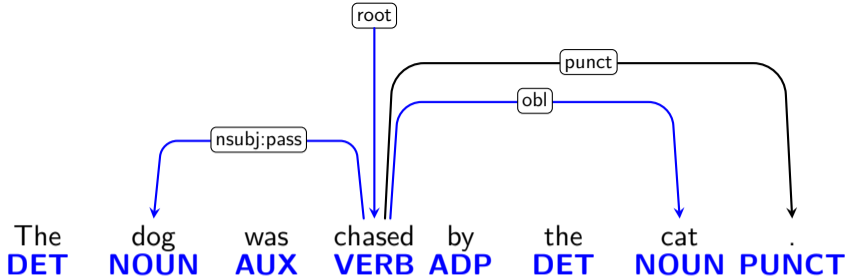
- Content words are related by dependency relations
- Function words attach to closest content words

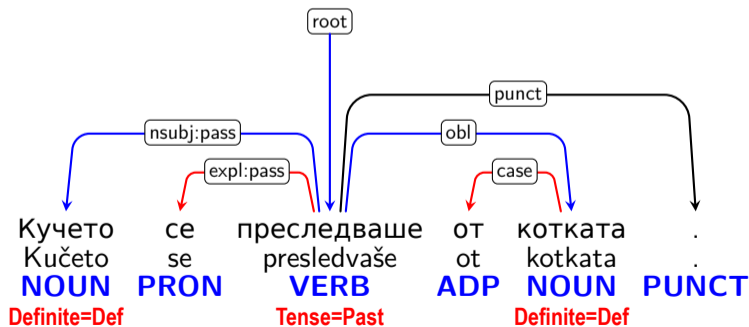
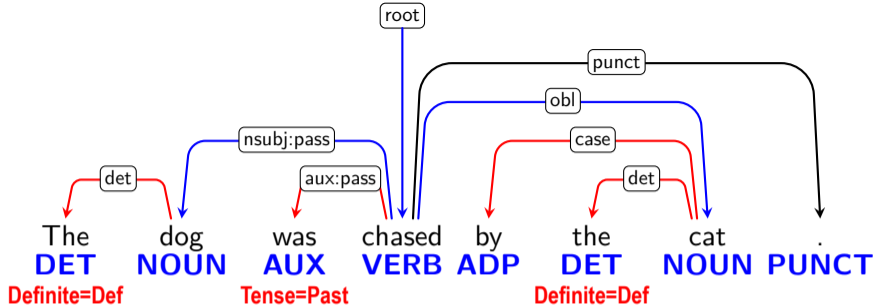


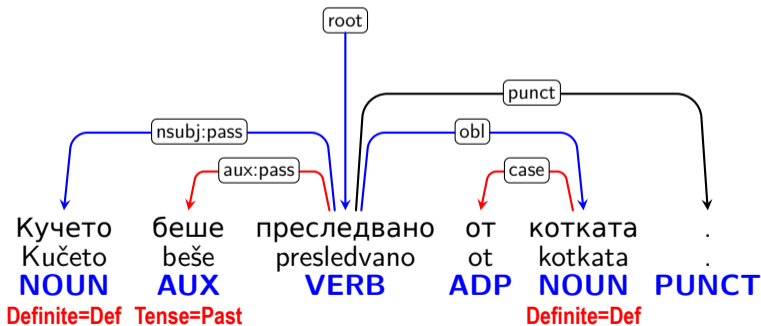
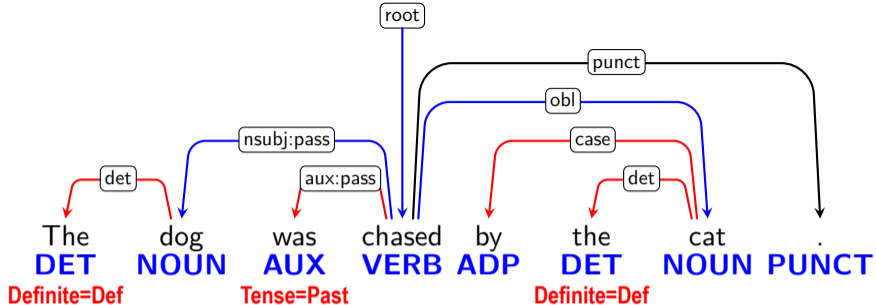
- Content words are related by dependency relations
- Function words attach to closest content words
- Punctuation attach to head of phrase or clause

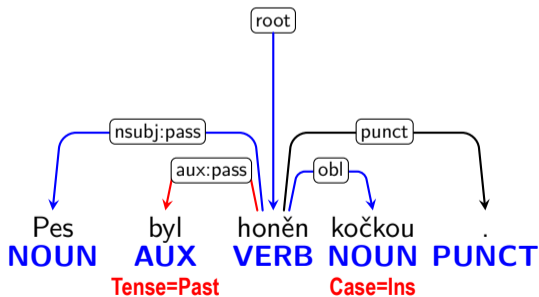
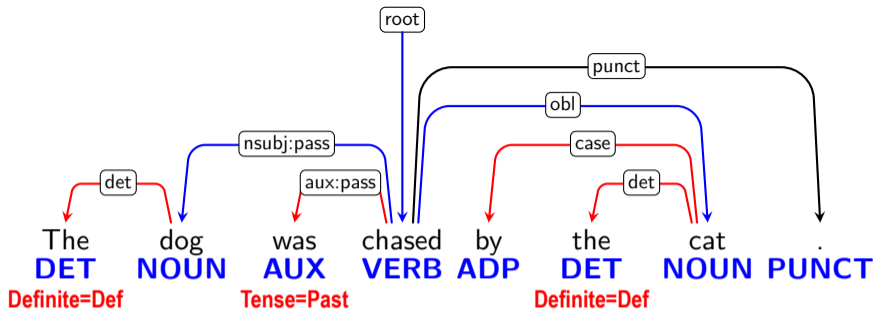


Not
"dependency"
in the strictly
syntactic
sense!









Dependents of Clauses (Verbal or Not)

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

Dependents of Verbs, Adjectives and Adverbs

	Nominal	Clausal	Modifier
Core	obj iobj	ccomp xcomp	
Non-Core	obl expl	advcl	advmod

Dependents of Nominals

Nominal	Clausal	Modifier	Function
nmod	acl	amod	det
case		nummod	case

Dependents of Nominals

Nominal

nmod

appos

compound

flat

Clausal

acl

Modifier

amod

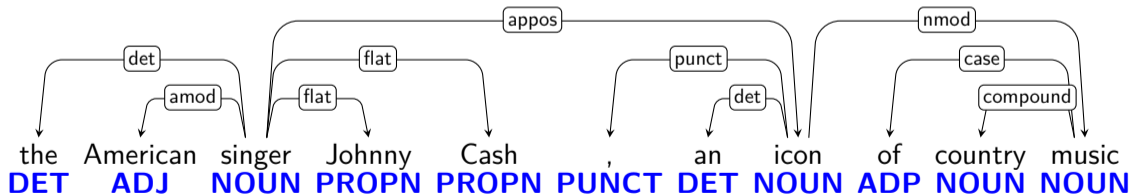
nummod

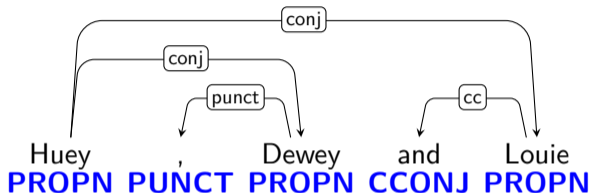
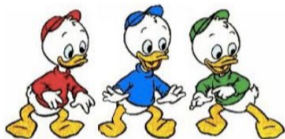
Function

det

case

clf





- Coordinate structures are headed by the first conjunct
 - Subsequent conjuncts depend on it via the **conj** relation
 - Conjunctions depend on the next conjunct via the **cc** relation
 - Punctuation marks depend on the next conjunct via the **punct** relation

Multiword Expressions

Relation	Examples
fixed	<i>as well, by and large, according to, more than</i>
flat	<i>president Havel, New York, four thousand</i>
compound	<i>phone book, dress up</i>
goeswith	<i>notwith standing, with out</i>

- UD annotation **almost** does not permit “words with spaces”
 - Multiword expressions are analyzed using special relations
 - The **fixed**, **flat** and **goeswith** relations are always head-initial
 - The **compound** relation reflects the internal structure
- Words with spaces allowed in exceptional cases:
 - Vietnamese (spaces delimit syllables, not words)
 - Numbers (“1 000 000”)
 - Possibly other approved cases, e.g. multi-word abbreviations

Other Relations

Relation	Explanation
parataxis	Loosely linked clauses of same rank
list	Lists without syntactic structure
orphan	Orphans in ellipsis linked together
reparandum	Disfluency linked to (speech) repair
dep	Unspecified dependency
root	The single syntactically independent element of the sentence

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

Relation	Explanation
acl:relcl	Relative clause (the boy who lived)
compound:prt	Verb particle (dress up)
nmod:poss	Possessive nominal (Mary 's book)
obl:agent	Agent in passive (saved by the bell)

Core vs. Oblique

- 1 Introduction
- 2 Morphological Annotation in UD
- 3 Syntactic Annotation in UD
- 4 Core vs. Oblique**
- 5 Enhanced Universal Dependencies
- 6 UD Tools

Dependents of Clauses (Verbal or Not)

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

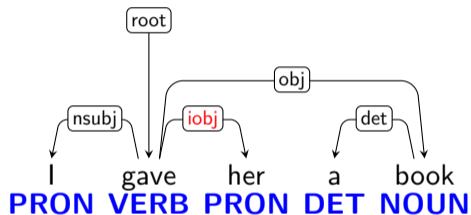
Dependents of Verbs, Adjectives and Adverbs

	Nominal	Clausal	Modifier
Core	obj iobj	ccomp xcomp	
Non-Core	obl expl	advcl	advmod

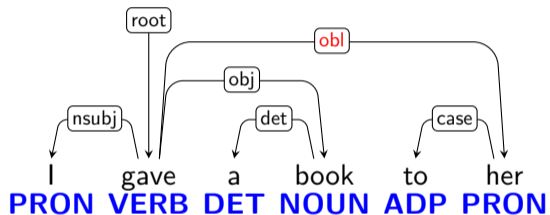
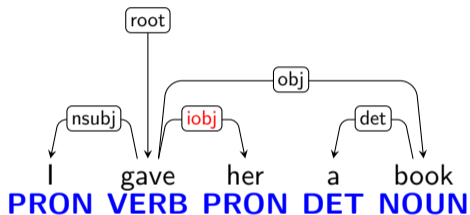
Dependents of Nominals

Nominal	Clausal	Modifier	Function
nmod	acl	amod	det
case		nummod	case

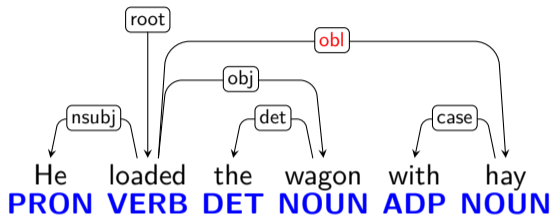
Information Packaging



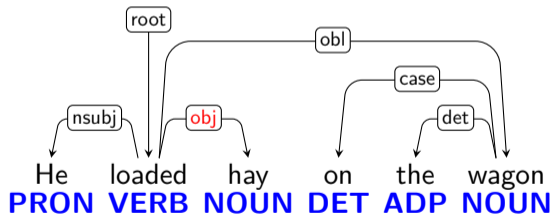
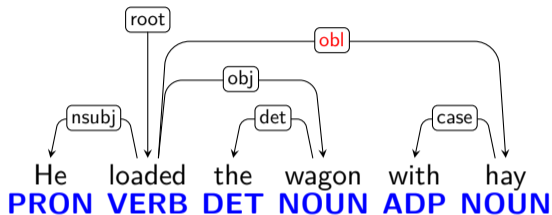
Information Packaging



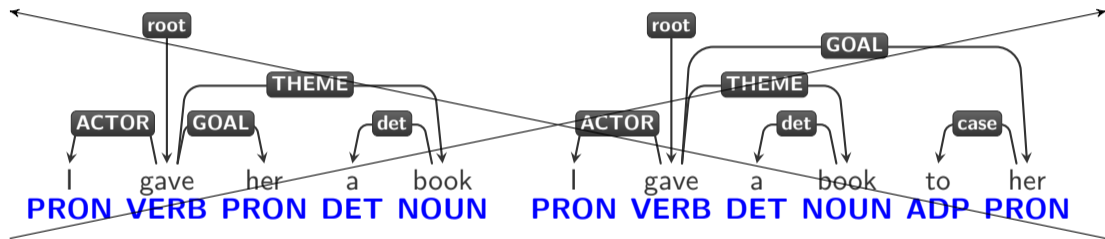
Information Packaging



Information Packaging



UD is NOT about Semantic Roles!



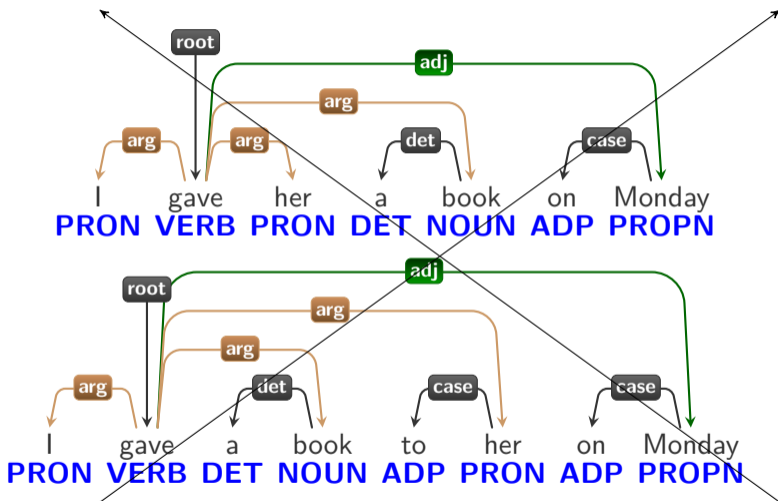
Manning's Law – What If We Do Semantic Roles?

The secret to understanding the design and current success of UD is to realize that the design is a very subtle compromise between approximately 6 things:

- 1 UD must be satisfactory on linguistic analysis grounds for **individual languages**.
- 2 UD must be good for linguistic **typology**, i.e., providing a suitable basis for bringing out cross-linguistic parallelism across languages and language families.
- 3 UD must be suitable for **rapid, consistent annotation** by a human annotator.
- 4 UD must be easily comprehended and used by a **non-linguist**, whether a language learner or an engineer with prosaic needs for language processing. ... it leads us to favor traditional grammar notions and terminology.
- 5 UD must be suitable for **computer parsing** with high accuracy.
- 6 UD must support well **downstream language understanding tasks** (relation extraction, reading comprehension, machine translation, ...)

It's easy to come up with a proposal that improves UD on one of these dimensions. The interesting and difficult part is to improve UD while remaining sensitive to all these dimensions.

UD Avoids Argument-Adjunct Distinction!



Avoiding an Argument-Adjunct Distinction

- From the guidelines:
 - Subtle, unclear, and frequently argued over
 - Questionable as a categorical distinction
 - Best practical solution is to eliminate it

Avoiding an Argument-Adjunct Distinction

- From the guidelines:
 - Subtle, unclear, and frequently argued over
 - Questionable as a categorical distinction
 - Best practical solution is to eliminate it
- **BUT:**
 - Cannot be eliminated completely
 - Some people/data have it and want to keep it
 - It aligns well with traditional grammars
 - \Rightarrow there is now a relation subtype `obl:arg`

Avoiding an Argument-Adjunct Distinction

- From the guidelines:
 - Subtle, unclear, and frequently argued over
 - Questionable as a categorical distinction
 - Best practical solution is to eliminate it
- BUT:
 - Cannot be eliminated completely
 - Some people/data have it and want to keep it
 - It aligns well with traditional grammars
 - \Rightarrow there is now a relation subtype `obl:arg`
- **AND** I will argue that
 - Core-oblique distinction is unclear and argued over too
 - (Though I will **not** propose to discard it.)

So What Is Core and Why?



THE CORE

- UD v1 guidelines took core-oblique for granted
- English (simplified):
 - Bare noun phrase \Rightarrow core argument (nsubj, obj, iobj)
 - Prepositional phrase \Rightarrow oblique argument or adjunct (obl)

- UD v1 guidelines took core-oblique for granted
- English (simplified):
 - Bare noun phrase \Rightarrow core argument (nsubj, obj, iobj)
 - Prepositional phrase \Rightarrow oblique argument or adjunct (obl)
- Other languages: not necessarily! (Spanish, Japanese)
 - But some people simply took the English rule...
 - Manning's law: non-linguists!

- UD v1 guidelines took core-oblique for granted
- English (simplified):
 - Bare noun phrase \Rightarrow core argument (nsubj, obj, iobj)
 - Prepositional phrase \Rightarrow oblique argument or adjunct (obl)
- Other languages: not necessarily! (Spanish, Japanese)
 - But some people simply took the English rule...
 - Manning's law: non-linguists!
- Clash with traditional terminology
 - Grammars of German, Czech etc. define **prepositional objects**
 - But these are not necessarily core...
 - Yet some people took their national definition of object...

Language-specific Coding Strategy

- Idea:
 - **Oblique** arguments are marked **similarly to adjuncts** (prepositions, certain morphological cases...)
 - Core arguments are marked differently
 - \Rightarrow easy for annotators and non-linguists!
- Why are core arguments special?
 - They tend to be **targeted by grammatical rules**
 - Passivization
 - Control verbs
 - Reflexives
 - ...

Language-specific Coding Strategy

- Core vs. oblique is not defined in traditional grammar
- How shall we define it?

Language-specific Coding Strategy

- Core vs. oblique is not defined in traditional grammar
- How shall we define it?
- Andrews, 2007 (In Shopen: Language Typology)
 - Identify **primary transitive predicates**
 - We need semantic roles for this! (One-time only.)

Language-specific Coding Strategy

- Core vs. oblique is not defined in traditional grammar
- How shall we define it?
- Andrews, 2007 (In Shopen: Language Typology)
 - Identify **primary transitive predicates**
 - We need semantic roles for this! (One-time only.)
 - Actor/agent = function **A**
 - Undergoer/patient = function **P**

Language-specific Coding Strategy

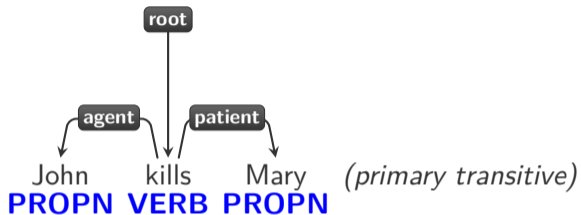
- Core vs. oblique is not defined in traditional grammar
- How shall we define it?
- Andrews, 2007 (In Shopen: Language Typology)
 - Identify **primary transitive predicates**
 - We need semantic roles for this! (One-time only.)
 - Actor/agent = function **A**
 - Undergoer/patient = function **P**
 - Note the way they are coded
 - Note other grammatical rules that target them
 - Generalize to other predicates with same coding and rules

Language-specific Coding Strategy

- Core vs. oblique is not defined in traditional grammar
- How shall we define it?
- Andrews, 2007 (In Shopen: Language Typology)
 - Identify **primary transitive predicates**
 - We need semantic roles for this! (One-time only.)
 - Actor/agent = function **A**
 - Undergoer/patient = function **P**
 - Note the way they are coded
 - Note other grammatical rules that target them
 - Generalize to other predicates with same coding and rules
- Then define:
 - function A \Rightarrow **nsubj**
 - function P \Rightarrow **obj**

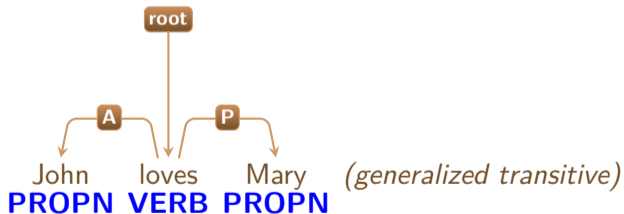
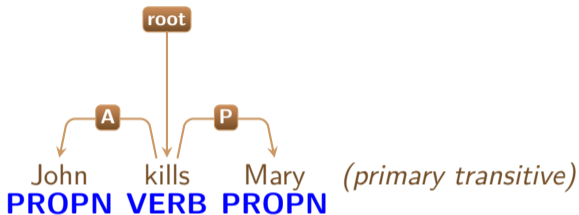


Transitive Predicates in English



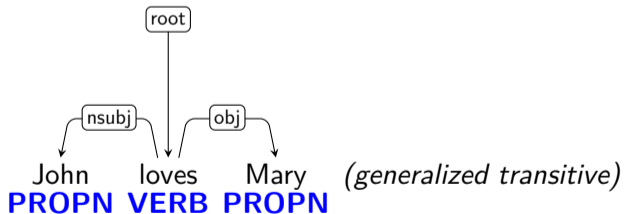
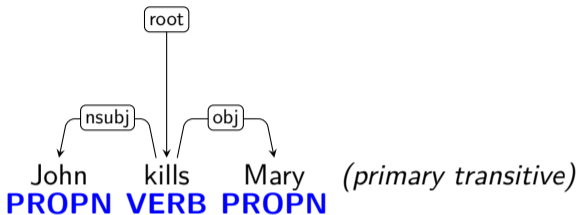


Transitive Predicates in English



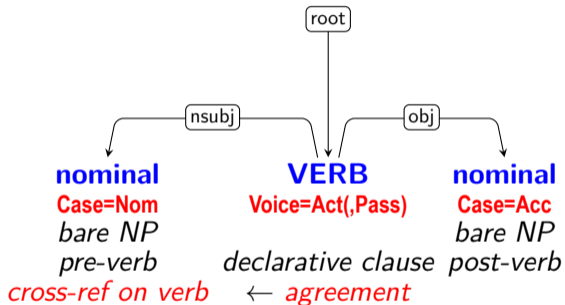


Transitive Predicates in English



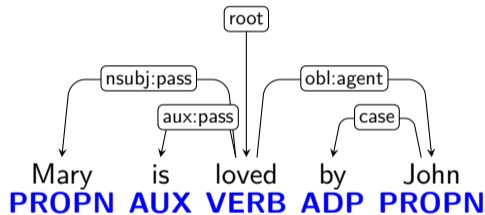
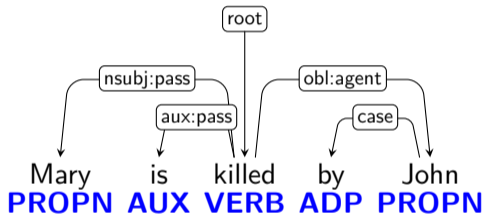


Transitive Predicates in English



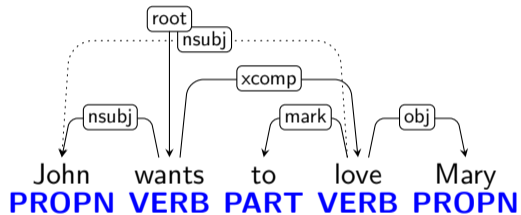
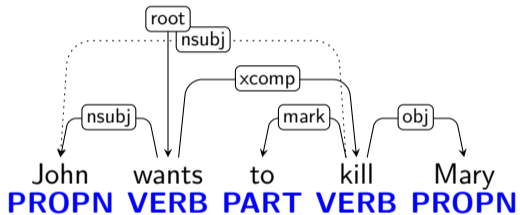


Passivization in English



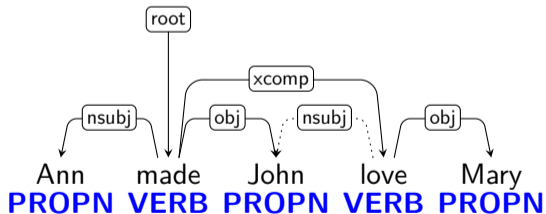
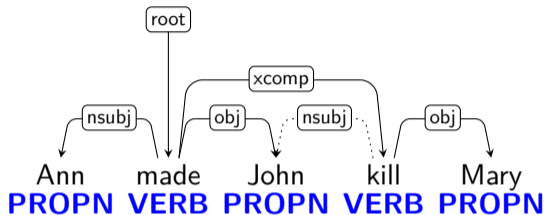


Subject Control in English





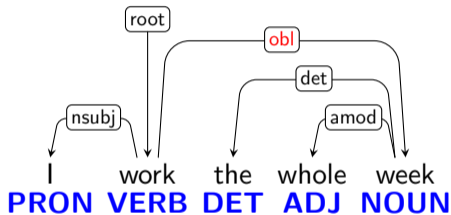
Object Control in English





Some Problems

- Some temporal adjuncts are bare noun phrases
 - *I work the whole week.*
 - *I work every Friday.*



- At least it cannot passivize:
 - **The whole week is worked by me.*
 - **Every Friday is worked by me.*
- But...



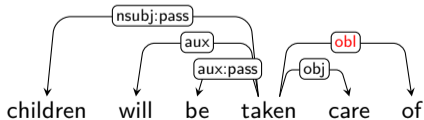
Some Problems

- Some transitive verbs cannot passivize
 - *John has a new car.*
 - **A new car is had by John.*
 - *Friday does not suit me.*
 - **I am not suited by Friday.*



Some Problems

- Some transitive verbs cannot passivize
 - John *has* a new car.
 - *A new car is had by John.
 - Friday does not *suit* me.
 - *I am not suited by Friday.
- Some prepositional verbs can passivize
 - You can *rely* on Ben.
 - Ben can be relied on.
 - They will *take* care of your children.
 - Your children will be taken care of.





Bare Temporal Adjuncts: Any Other Criteria?

- *I work **the whole week**.*
- *I work **every Friday**.*
- English has a fixed word order; adjuncts are less fixed than objects:
 - *I work every Friday in Paris.*
 - *I work in Paris every Friday.*
 - *I spend every Friday in Paris.*
 - **I spend in Paris every Friday.*
- Unlike objects, adjuncts cannot be replaced by pronouns:
 - *Where do you spend this Friday? I spend it in Paris.*
 - *Where do you work this Friday? *I work it in Paris.*

Tentative Summary?

- The borderline is inherently fuzzy
- No universally applicable and exact algorithm
- Better described in terms of probability



Tentative Summary?



- The borderline is inherently fuzzy
- No universally applicable and exact algorithm
- Better described in terms of probability

- Core coding **not favored by adjuncts**
- Oblique coding **similar to most adjuncts**
- Passivization etc. may help...
- ... but does **not** work as **strict criterion**

Tentative Summary?



- The borderline is inherently fuzzy
- No universally applicable and exact algorithm
- Better described in terms of probability

- Core coding **not favored by adjuncts**
- Oblique coding **similar to most adjuncts**
- Passivization etc. may help...
- ... but does **not** work as **strict criterion**

- Semantic roles needed when starting a new language
- Argument-adjunct might help with exceptions
 - Although we managed to explain *the whole week* without it



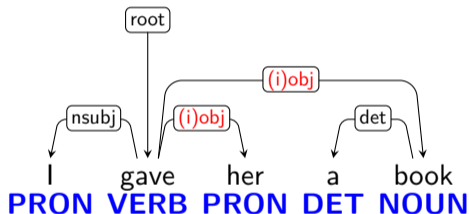
Intransitive Predicates

- Just one core argument
 - We already “know” how to find out if there are two
- \Rightarrow function **S**
 - Regardless of semantic role:
 - *John runs.*
 - *John sleeps.*
 - *John falls.*
- Then define:
 - function **S** \Rightarrow **nsubj**



Ditransitive Predicates

- Three core arguments
- Is one of them “least core”? \Rightarrow **iobj**
- (Alternatively, we could look at the semantic roles once again.)

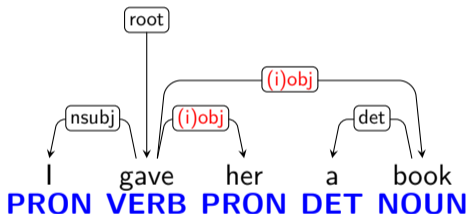


- Passivization:
 - *She* was given a book by me.
 - ?*A book* was given her by me.



Ditransitive Predicates

- Three core arguments
- Is one of them “least core”? \Rightarrow **iobj**
- (Alternatively, we could look at the semantic roles once again.)

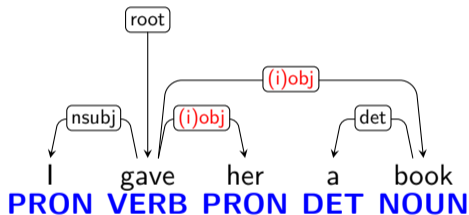


- Fronting in questions:
 - *What* did I give her?
 - **Who* did I give a book?



Ditransitive Predicates

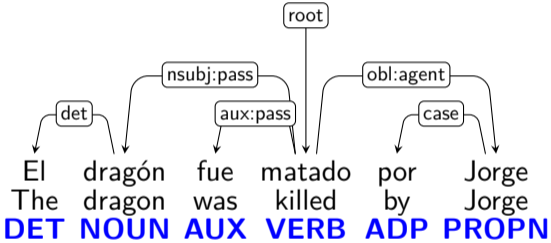
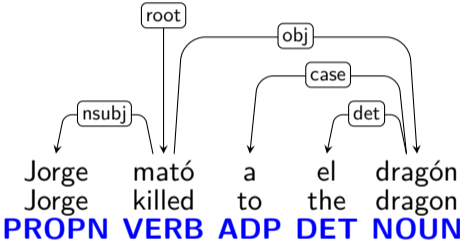
- Three core arguments
- Is one of them “least core”? \Rightarrow **iobj**
- (Alternatively, we could look at the semantic roles once again.)



- Andrews (2007): *the status of the notion of ‘indirect object’ is problematic and difficult to sort out. The top priority is to work out what properties recipients and themes do and do not share with P arguments of primary transitive verbs.*

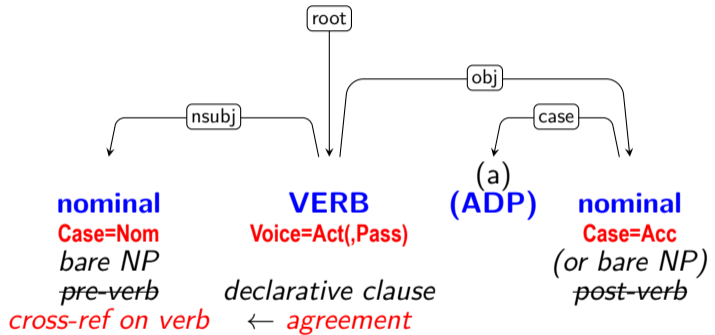


Spanish (EXTRA)



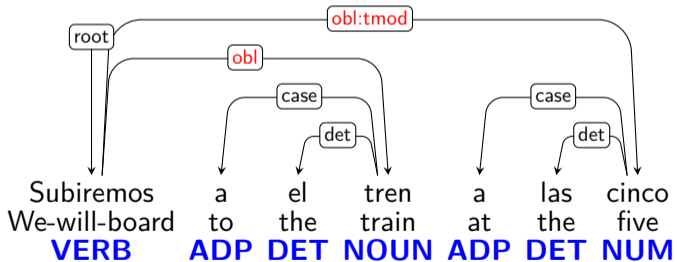
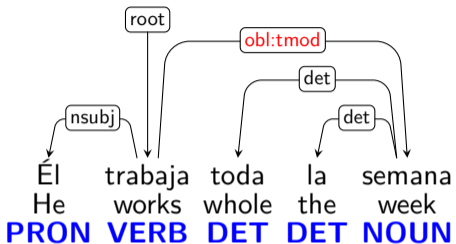


Spanish Transitive Clauses (EXTRA)



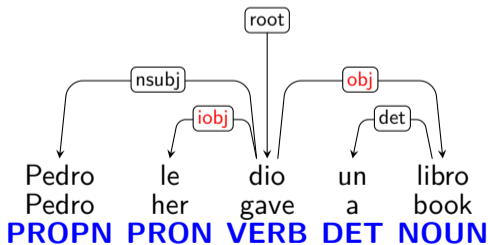
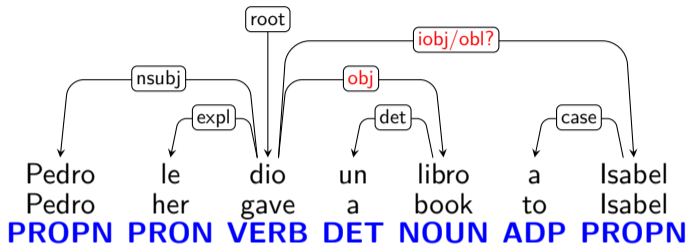


Spanish Adjunct Exceptions (EXTRA)



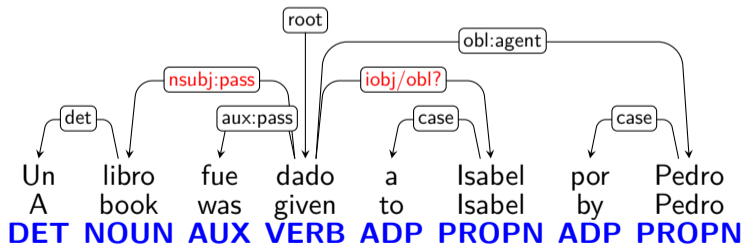
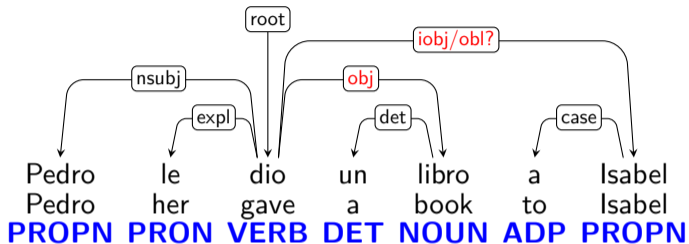


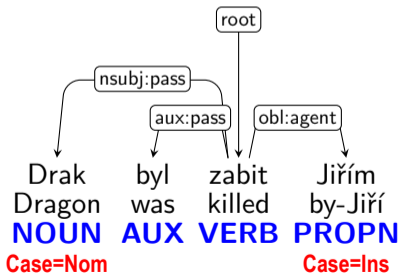
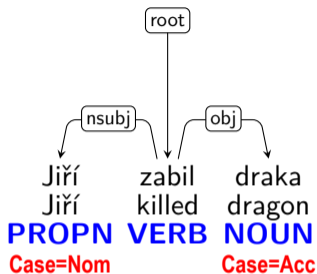
Spanish Ditransitive Clauses (EXTRA)

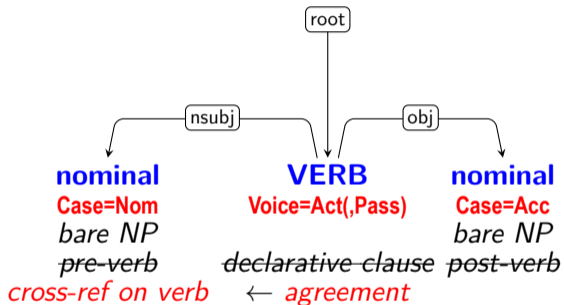




Spanish Ditransitive Clauses (EXTRA)

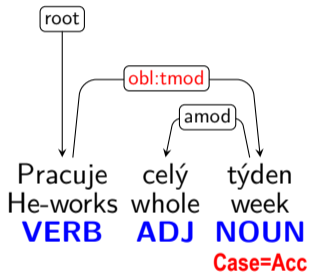


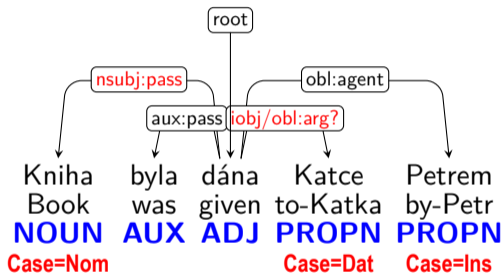
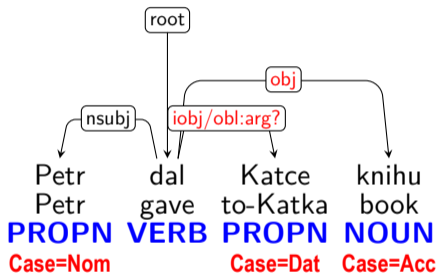






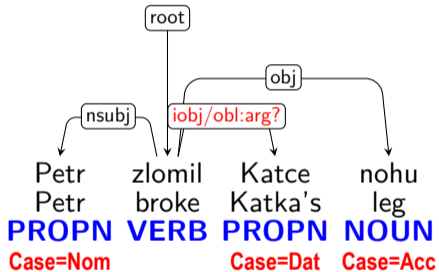
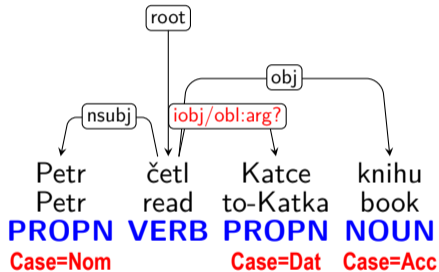
Czech Adjunct Exceptions (EXTRA)





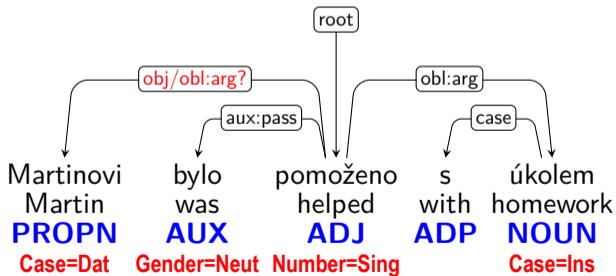
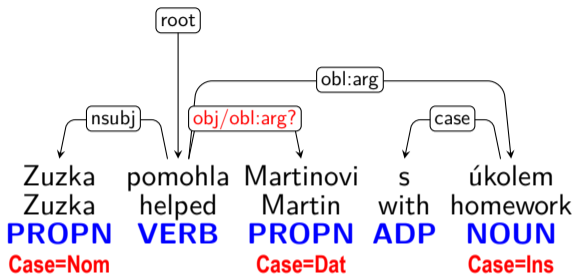


Dative: Recipient vs. Beneficiary (EXTRA)



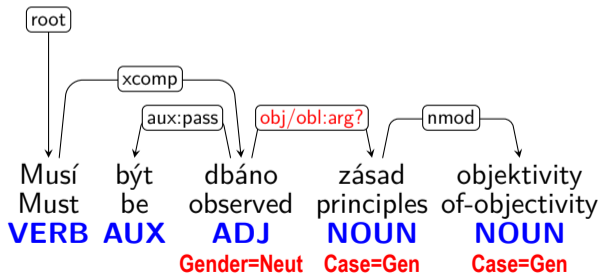
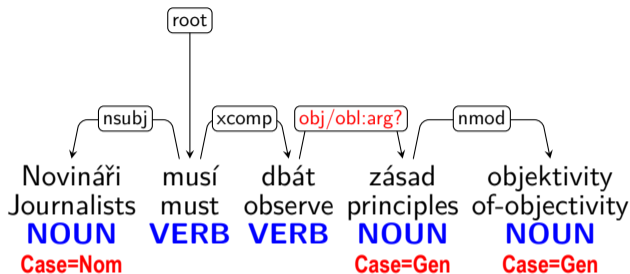


Monotransitive with Dative? (EXTRA)



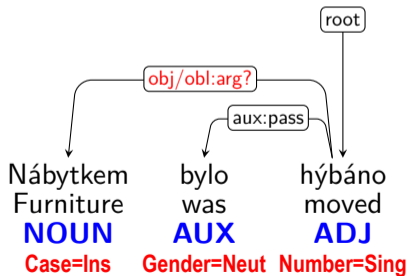
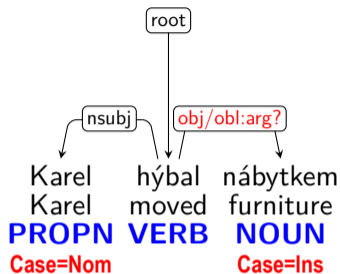


Monotransitive with Genitive? (EXTRA)



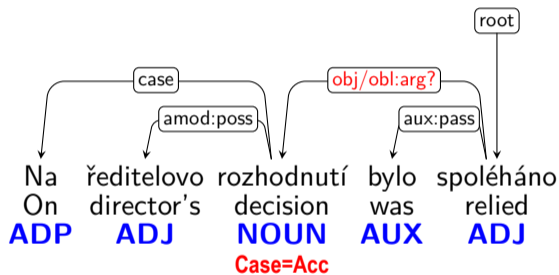
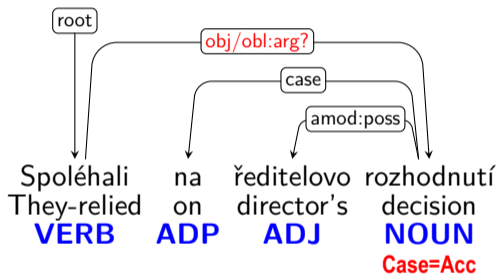


Monotransitive with Instrumental? (EXTRA)





Monotransitive with Preposition? (EXTRA)





- There is a core-oblique scale:
- **Nom** > **Acc** > **Gen,Dat** > **Ins** > **preposition**
- Where is the borderline?





- There is a core-oblique scale:
- **Nom** > **Acc** > **Gen,Dat** > **Ins** > **preposition**
- Where is the borderline?
- UD Czech 1.0: object = argument
 - **Nom, Acc, Gen, Dat, Ins, ADP** > “adverbial”



- There is a core-oblique scale:
- **Nom** > **Acc** > **Gen,Dat** > **Ins** > **preposition**
- Where is the borderline?
- UD Czech 1.0: object = argument
 - **Nom, Acc, Gen, Dat, Ins, ADP** > “adverbial”
- UD Czech 2.1–2.5: bare NP > PP
 - **Nom, Acc, Gen, Dat, Ins** > **ADP** + adjuncts



- There is a core-oblique scale:
- **Nom** > **Acc** > **Gen,Dat** > **Ins** > **preposition**
- Where is the borderline?
- UD Czech 1.0: object = argument
 - **Nom, Acc, Gen, Dat, Ins, ADP** > “adverbial”
- UD Czech 2.1–2.5: bare NP > PP
 - **Nom, Acc, Gen, Dat, Ins** > **ADP** + adjuncts
- UD Czech 2.6 (May 2020):
 - **Nom, Acc** > **Gen, Dat, Ins, ADP** + adjuncts



- There is a core-oblique scale:
- **Nom** > **Acc** > **Gen,Dat** > **Ins** > **preposition**
- Where is the borderline?
- UD Czech 1.0: object = argument
 - **Nom, Acc, Gen, Dat, Ins, ADP** > “adverbial”
- UD Czech 2.1–2.5: bare NP > PP
 - **Nom, Acc, Gen, Dat, Ins** > **ADP + adjuncts**
- UD Czech 2.6 (May 2020):
 - **Nom, Acc** > **Gen, Dat, Ins, ADP + adjuncts**
 - \Rightarrow No ditransitives in Czech!
 - (Exception: *učit* “to teach” takes two Acc.)

Summary of Basic Syntax

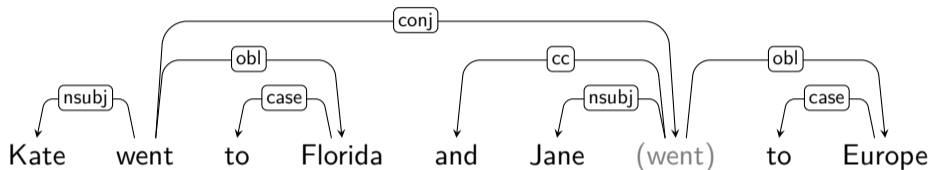
- Universal Dependencies
 - Content words higher than function words ... better parallelism
 - Clauses – nominals – modifier words
 - Distinguished both as heads and as dependents
 - Core arguments vs. oblique dependents
 - Language-specific subtypes of relations

<https://universaldependencies.org/>

Enhanced Universal Dependencies

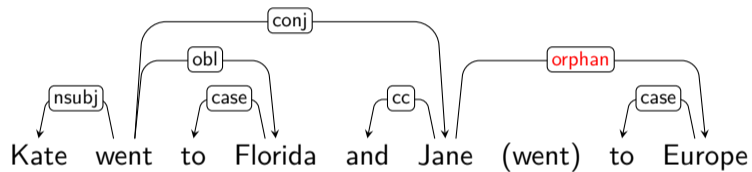
- 1 Introduction
- 2 Morphological Annotation in UD
- 3 Syntactic Annotation in UD
- 4 Core vs. Oblique
- 5 Enhanced Universal Dependencies**
- 6 UD Tools

Deleted Predicates in Coordination

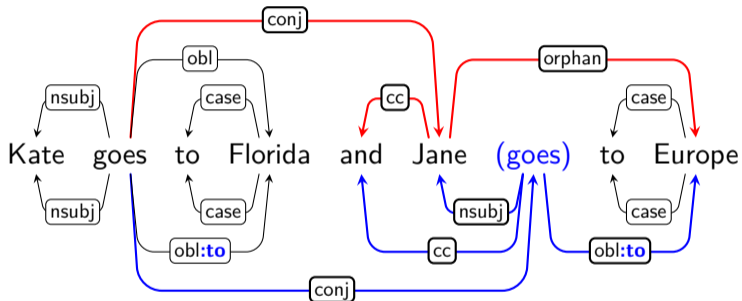


- Some treebanks would use an **empty node** to represent the second *went*.
- UD **enhanced representation** allows empty nodes!
- But the basic representation sticks with the overt words.

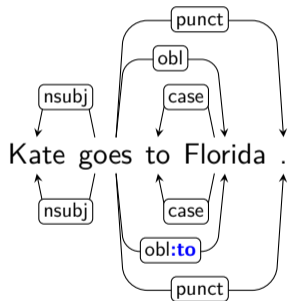
UD V2 Basic Dependencies: The orphan Relation



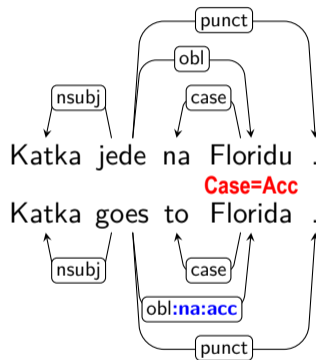
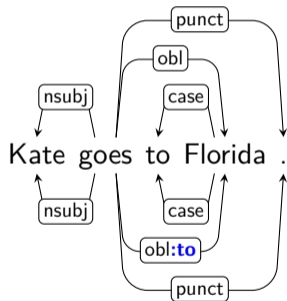
UD Enhanced Dependencies: Gapping and Stripping



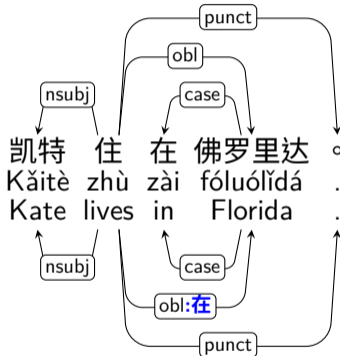
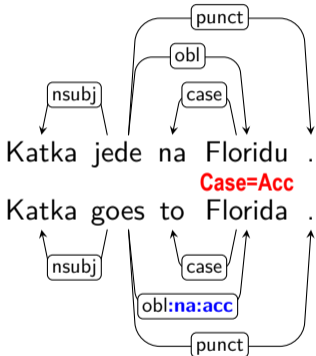
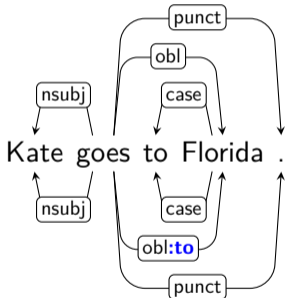
Enhanced UD: Case Information in Dependency Label



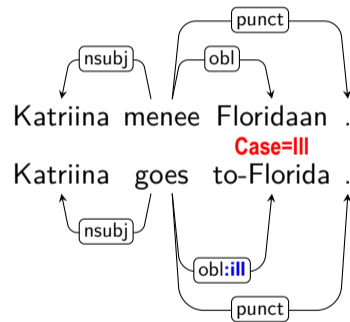
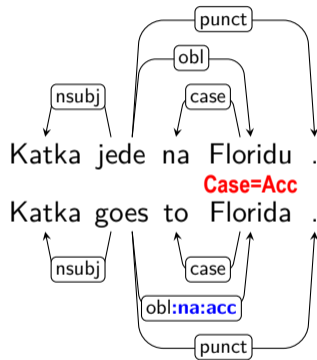
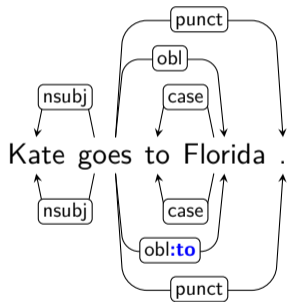
Enhanced UD: Case Information in Dependency Label



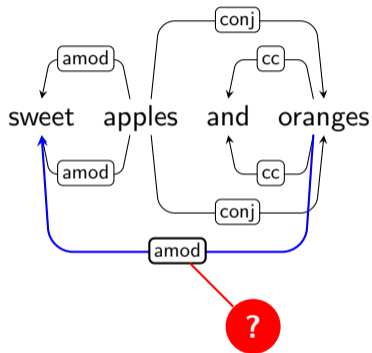
Enhanced UD: Case Information in Dependency Label



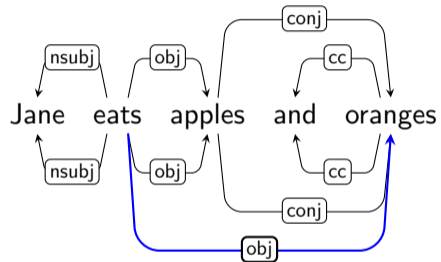
Enhanced UD: Case Information in Dependency Label



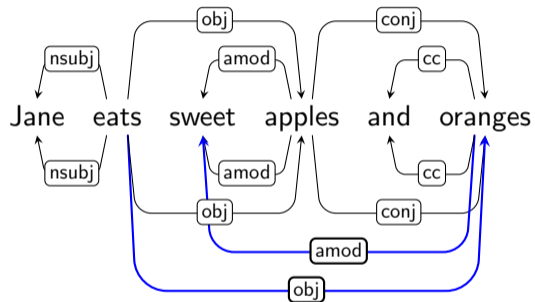
Enhanced UD: Shared Dependent of Coordination



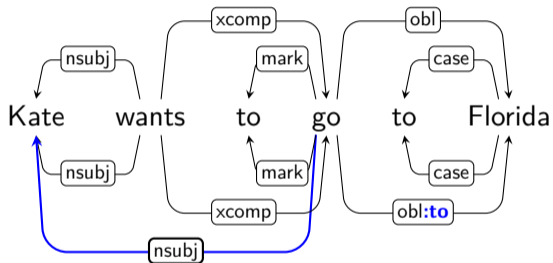
Enhanced UD: Parent of Coordination



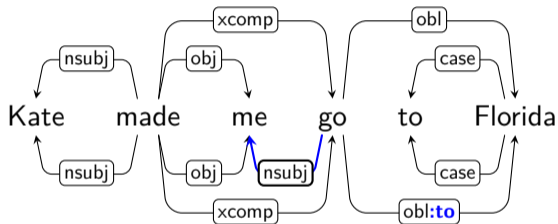
Enhanced UD: Coordination



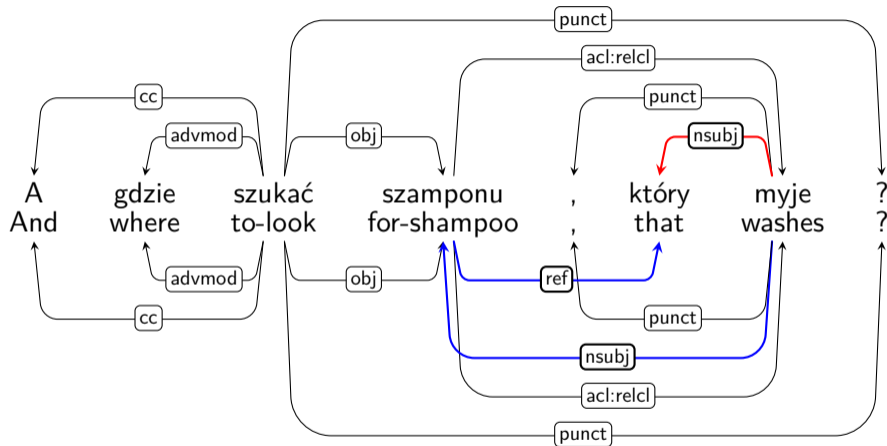
Enhanced UD: External Subject of Controlled Predicate



Enhanced UD: External Subject in Object-Control Construction

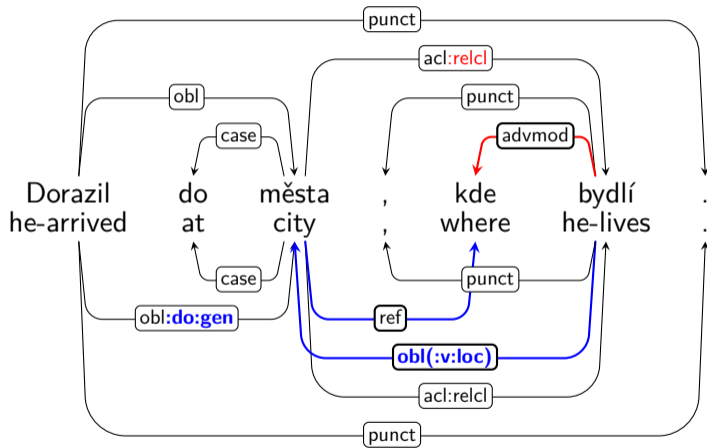


Enhanced UD: Relative Clauses

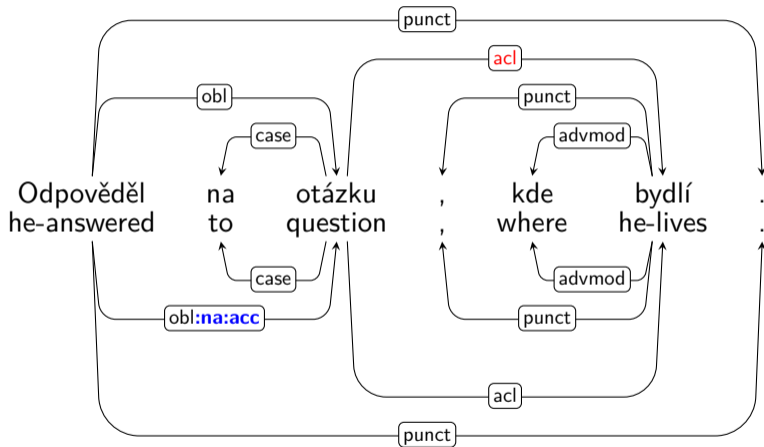


“And where to look for shampoo that works?”

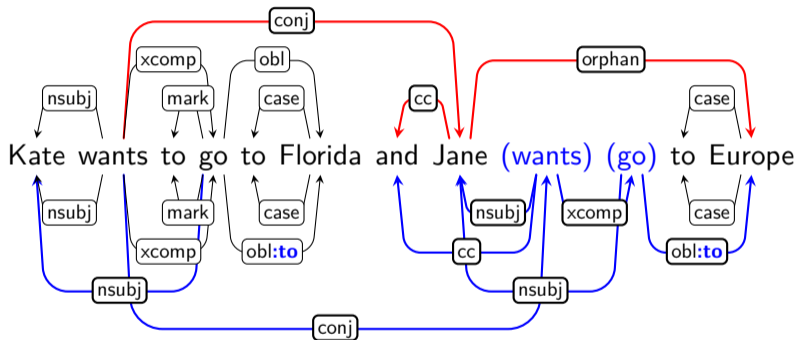
How to Recognize Relative Clauses (EXTRA)



How to Recognize Relative Clauses (EXTRA)



Enhanced UD: Gapping + Control Verb Construction



Enhanced UD: Six Enhancements (EXTRA)

- Null nodes for **gapping** (24 treebanks in UD 2.8)
- Dependency propagation in **coordination**
 - Common parent of coordination (28 treebanks)
 - Shared dependents of coordination (25 treebanks)
- External subjects of **controlled predicates** (21 treebanks)
- Cyclic dependencies to/from **relative clauses** (22 treebanks)
- **Case**-enhanced dependency labels (21 treebanks)

- All 6 types: 15 treebanks, 8 languages
- At least 1 type: 30 treebanks, 18 languages
- Only basic UD: 172 treebanks

- Part of Stanford CoreNLP (Java)
- Rules from basic to enhanced UD

- Part of Stanford CoreNLP (Java)
- Rules from basic to enhanced UD
- **Gapping**: embeddings for similarity of arguments

- Part of Stanford CoreNLP (Java)
- Rules from basic to enhanced UD
- **Gapping**: embeddings for similarity of arguments
- Coordination:
 - **Parent propagation**: deterministic
 - **Shared dependents**: heuristics (human desirable!)

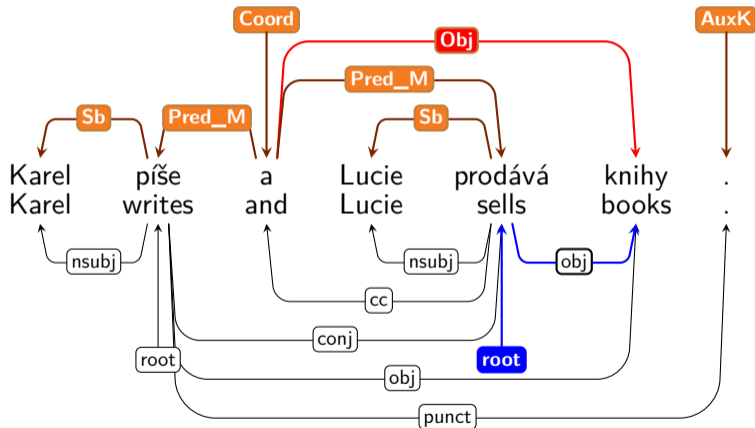
- Part of Stanford CoreNLP (Java)
- Rules from basic to enhanced UD
- **Gapping**: embeddings for similarity of arguments
- **Coordination**:
 - **Parent propagation**: deterministic
 - **Shared dependents**: heuristics (human desirable!)
- **External subjects**: heuristics (subject vs. object control)

- Part of Stanford CoreNLP (Java)
- Rules from basic to enhanced UD
- **Gapping**: embeddings for similarity of arguments
- **Coordination**:
 - **Parent propagation**: deterministic
 - **Shared dependents**: heuristics (human desirable!)
- **External subjects**: heuristics (subject vs. object control)
- **Relative clauses**: need `ac1:relcl` and list of relative pronouns

- Part of Stanford CoreNLP (Java)
- Rules from basic to enhanced UD
- **Gapping**: embeddings for similarity of arguments
- **Coordination**:
 - **Parent propagation**: deterministic
 - **Shared dependents**: heuristics (human desirable!)
- **External subjects**: heuristics (subject vs. object control)
- **Relative clauses**: need `ac1:relcl` and list of relative pronouns
- **Case-enhanced labels**: deterministic

Conversion from non-UD Data: Extra Information? (EXTRA)

- Analytical layer of Prague-style treebanks: **shared dependents of coordination** are known!



Summary of Enhanced Syntax

- Some deep-syntactic relations
- Directed graph (**not always tree**)
- Six enhancement types
 - Empty nodes for gapping
 - Shared parents in coordination
 - Shared dependents in coordination
 - External subjects in control/raising
 - Cyclic dependency in relative clauses
 - Case-enhanced modifier relations
- Some (not all) can be guessed from the basic tree

<https://universaldependencies.org/>

UD Tools

- 1 Introduction
- 2 Morphological Annotation in UD
- 3 Syntactic Annotation in UD
- 4 Core vs. Oblique
- 5 Enhanced Universal Dependencies
- 6 UD Tools**

- <https://universaldependencies.org/tools.html>

Linguists Can Search Treebanks

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

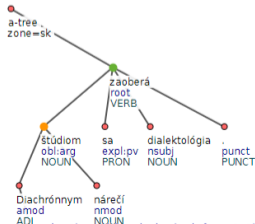
Relations Node Types Attributes Operators Functions

```
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.



Linguists Can Parse and Search New Data

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

cs.wikipedia.org/wiki/Covid-19

WIKIPEDIE
Otevřená encyklopedie

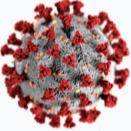
Covid-19

Tento článek reaguje na aktuální nebo nedávné události. Informace zde uvedené se vzhledem k neustálému vývoji mohou průběžně měnit. Je třeba je se zvýšenou péčí aktualizovat a doplňovat.

Tento článek pojednává o nemoci, kterou způsobuje koronavirus SARS-CoV-2. Možná hledáte: *Pandemie covidu-19, nebo SARS-CoV-2, nebo Pandemie covidu-19 v Česku.*

Covid-19 (těž COVID-19,^[zdroj 1] z anglického spojení *coronavirus disease 2019*, což česky znamená koronavirové onemocnění 2019; výslovnost: [kovid devatenáct]; podle ICD-11 označeno **XN109**) je vysoce infekční onemocnění, které je způsobeno koronavirem SARS-CoV-2. První případ byl identifikován v čínském Wu-chanu v prosinci 2019. Od té doby se virus rozšířil po celém světě, což způsobilo pletnávající pandemii. Příznaky nemoci covid-19 jsou různé, od bezpříznakového stavu až po závažné onemocnění, ale často zahrnují horečku, kašel, únavu, dýchací potíže a ztrátu čichu a chuti. Příznaky začínají jeden až čtrnáct dní po vystavení viru. U přibližně jednoho z pěti infikovaných jedinců se neobjeví žádné příznaky.^[2] Zatímco většina lidí má mírné příznaky, u některých lidí se vyvine syndrom akutní

Coronavirus disease 2019 covid-19



Koronavirus SARS-CoV-2 způsobující onemocnění

Klasifikace	
MKN-10	U07.1 a U07.2
Statistické údaje – obě pohlaví	
Incidence	230 567 044 ^[1] (z toho 0 ^[1] uzdravených) ke dni 22. září 2021
Mortalita	2,23 % ^[20497] (celosvětový průměr: 1,84%) ^[20498]

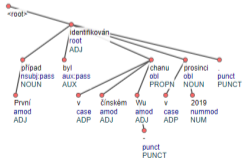
lindat.mff.cuni.cz/services/udpipe/

Process Input

Output Text Show Table Show Trees Save Tree as SVG

Previous 1 2 3 4 Next

První případ byl identifikován v čínském Wu - chanu v prosinci 2019 .



```

graph TD
    root["<root>"] --> identifikovan["identifikován<br/>root<br/>ADJ"]
    root --> punct1["punct<br/>PUNCT"]
    identifikovan --> pripad["případ<br/>noun<br/>pass<br/>NOUN"]
    identifikovan --> byl["byl<br/>aux<br/>pass<br/>AUX"]
    identifikovan --> v1["v<br/>case<br/>ADP"]
    identifikovan --> cinskem["čínském<br/>amod<br/>ADJ"]
    identifikovan --> wu["Wu<br/>amod<br/>ADJ"]
    identifikovan --> v2["v<br/>case<br/>ADP"]
    identifikovan --> prosinci["prosinci<br/>obl<br/>NOUN"]
    identifikovan --> 2019["2019<br/>nummod<br/>NUM"]
    pripad --> prvni["První<br/>amod<br/>ADJ"]
    v1 --> v1_text["v"]
    cinskem --> cinskem_text["čínském"]
    wu --> wu_text["Wu"]
    v2 --> v2_text["v"]
    prosinci --> prosinci_text["prosinci"]
    2019 --> 2019_text["2019"]
  
```

- UDPipe (<https://lindat.mff.cuni.cz/services/udpipe/>)
 - <https://ufal.mff.cuni.cz/udpipe>
- Stanza (<https://stanfordnlp.github.io/stanza/>)

- (Tred) <https://ufal.mff.cuni.cz/tred/>
- UD Annotatrix <https://github.com/jonorthwash/ud-annotatrix>

The screenshot shows the UD Annotatrix web interface. The top part displays the text and its morphological annotations in a table format. Below this, a dependency graph is shown with nodes for each word and edges representing grammatical relations.

Index	Word	POS	Case	Number	Gender	Case	Number	Gender	Case	Number	Gender	Gloss
1	جنگ	NOUN	-	-	-	9	nsbj	-	-	-	-	Gloss=girl Translit=jnk
2	ء	ADP	Case=Erg	-	-	1	case	-	-	-	-	Translit='a
3	ونی	PRON	-	-	-	4	nmod	-	-	-	-	Gloss=her Translit=üti
4	دزگوپار	NOUN	-	-	-	9	iobj	-	-	-	-	Gloss=friend Translit=dzgoöhär
5	ء	ADP	Case=Erg	-	-	4	case	-	-	-	-	Translit='a
6	را	ADP	Case=Dat	-	-	4	case	-	-	-	-	Translit=rä
7	ندی	NOUN	-	-	-	9	obj	-	-	-	-	Gloss=letter Translit=nndi
8	ے	X	-	-	-	7	case	-	-	-	-	Gloss=a Translit=ie

The dependency graph below shows the following relations:

- 1 (جنگ) is the subject of 9 (نیشته) via the relation <nsubj>.
- 2 (ء) is the case marker for 1 (جنگ) via the relation <case>.
- 3 (ونی) is the modifier of 4 (دزگوپار) via the relation <nmod>.
- 5 (ء) is the case marker for 4 (دزگوپار) via the relation <case>.
- 6 (را) is the case marker for 7 (ندی) via the relation <case>.
- 8 (ے) is the case marker for 7 (ندی) via the relation <case>.
- 9 (نیشته) is the object of 4 (دزگوپار) via the relation <iobj>.
- 9 (نیشته) is the object of 7 (ندی) via the relation <obj>.
- 10 (کت) is the modifier of 9 (نیشته) via the relation <compound>.
- 11 (.) is the punctuation at the end of the sentence via the relation <punct>.

- Morphology Spreadsheet (Excel / LibreOffice Calc)

Validation and Releasing

- (Tred) <https://ufal.mff.cuni.cz/tred/>
- UD Annotatrix <https://github.com/jonorthwash/ud-annotatrix>

→ ↻ ↗ ⚠ Nezábezpečeno | quest.ms.mff.cuni.cz/udvalidator/cgi-bin/unidep/validation-report.pl
punct 3; L3 Syntax rel-upos-punct 16; L3 Syntax right-to-left-apos 10; L3 Syntax too-many-subjects 4
UD_Hittite-HitTB: VALID
UD_Hungarian-Szeged: NEGLECTED; 2018-11-16 (TOTAL 318; L0 Repo lang-spec-doc 1; L3 Morpho
5; L3 Morpho goeswith-upos 5; L3 Syntax leaf-aux-cop 6; L3 Syntax leaf-cc 4; L3 Syntax leaf-mark-c
admod 75; L3 Syntax rel-upos-cc 1; L3 Syntax rel-upos-cop 1; L3 Syntax right-to-left-apos 43; L3 S
UD_Icelandic-IcePaHC: LEGACY; 2022-05-31 (TOTAL 81; L3 Syntax too-many-subjects 81) ([report](#))
UD_Icelandic-Modern: LEGACY; 2022-05-31 (TOTAL 16; L3 Syntax too-many-subjects 16) ([report](#))
UD_Icelandic-PUD: LEGACY; 2022-05-31 (TOTAL 6; L3 Syntax too-many-subjects 6) ([report](#))
UD_Indonesian-CSUI: VALID
UD_Indonesian-GSD: LEGACY; 2022-05-31 (TOTAL 58; L3 Syntax too-many-subjects 58) ([report](#))
UD_Indonesian-PUD: VALID
UD_Irish-Cadhan: VALID
UD_Irish-IDT: LEGACY; 2022-02-19 (TOTAL 82; L3 Morpho goeswith-feats 1; L3 Morpho goeswith-l
too-many-subjects 72) ([report](#))
UD_Irish-TwitIrish: LEGACY; 2022-02-19 (TOTAL 12; L3 Morpho goeswith-lemma 1; L3 Morpho gr
UD_Italian-ISDT: LEGACY; 2022-05-31 (TOTAL 6; L3 Syntax too-many-subjects 6) ([report](#)) The foll
unknown-edeprel
UD_Italian-MarkIT: LEGACY; 2022-05-31 (TOTAL 10; L3 Syntax too-many-subjects 10) ([report](#))
UD_Italian-PUD: VALID
UD_Italian-ParTUT: VALID
UD_Italian-PoSTWITA: LEGACY; 2022-02-19 (TOTAL 48; L3 Morpho goeswith-feats 3; L3 Morpho
Syntax too-many-subjects 20) ([report](#))
UD_Italian-TWITTIRO: LEGACY; 2022-02-19 (TOTAL 7; L3 Morpho goeswith-lemma 1; L3 Morpho
UD_Italian-VIT: ERROR; BACKUP 2.10 (TOTAL 71; L3 Syntax punct-causes-nonproj 1; L3 Syntax p
UD_Italian-Valico: VALID
UD_Japanese-BCCWJ: LEGACY; 2022-05-31 (TOTAL 1063; L3 Syntax too-many-subjects 1063) ([re](#)
UD_Japanese-BCCWJLUW: LEGACY; 2022-05-31 (TOTAL 1244; L3 Syntax too-many-subjects 12
UD_Japanese-GSD: LEGACY; 2022-05-31 (TOTAL 191; L3 Syntax too-many-subjects 191) ([report](#))
UD_Japanese-GSDLUW: LEGACY; 2022-05-31 (TOTAL 205; L3 Syntax too-many-subjects 205) ([re](#)
UD_Japanese-KTC: ERROR; DISCARD (TOTAL 40878; L0 Repo readme 1; L2 Metadata missing e

← → ↻ ↗ ⚠ quest.ms.mff.cuni.cz/udvalidator/cgi-bin/unidep/langspec/specify_auxiliary.pl?icode=it

Specify auxiliaries for Italian

Remember: Not everything that a traditional grammar labels as auxiliary is necessarily an [auxiliary in UD](#). In
auxiliary; the usual alternative in UD is treating one of the verbs as an [xcore](#) of the other, or in some languages
are grammatical rather than semantic: just because something has a modal or near-modal meaning does not m
Some verbs function as auxiliaries in some constructions and as full verbs in others (e.g., *to have* in English).

Remember: A language typically has at most one lemma for [copula](#). Exceptions include deficient paradigms (
Romance verbs *ser* and *estar* (both equivalents of “to be”). In contrast, equivalents of “to become, to stay, to k
such. In UD they should head an [xcore](#) relation instead. Existential “to be” can be copula only if it is the same
existential one is not a copula. A copula is normally tagged [AUX](#). Exception: in some languages a personal o
or [DET](#).

Edit or add auxiliaries

[andare](#) [avere](#) [dovere](#) [essere](#) [fare](#) [potere](#) [sapere](#) [stare](#) [venire](#) [volere](#)

Known auxiliaries for this and other languages

	Language	Total	Copula	Perfect	Past	Future	Passive	Conditional	N
Italian	it	10	essere	avere			andare venire		dovere
Neapolitan	nap	2	essè	avé					
Romanian	ro	9	fi hiu	am avea		voi vrea			trebui
Ligurian	lij	8	stá èse	avei			stá vegni		dovei
French	fr	3	être	avoir			être		
Middle French	frm	0							
Old French	fro	7	estre	avoir			estre		devoir

Documentation and Statistics

universaldependencies.org/cs/index.html

UD for Czech



Tokenization and Word Segmentation

- In general, words are delimited by whitespace characters. Description of exceptions follows.
- According to typographical rules, many punctuation marks are attached to a neighboring word. We always tokenize them as separate tokens (words): that holds even for hyphenated compounds such as česko-slovensky "Czech-Slovak" (three tokens) and for abbreviations such as atd "etc." (two tokens).
- A whitespace separating digits in a large number is not treated as a word separator. For example, 1 000 000 ("1,000,000" by English rules) is one token.
- There are several closed classes of contractions that are treated as multi-word tokens and segmented to individual syntactic words. The most prominent type is a subordinating conjunction fused with a conditional auxiliary: *kydyžch* = *kydyž* + *bych* "if". For more details, see [tokenization](#).

Morphology

Tags

This is an overview only. For more detailed discussion and examples, see the list of [Czech POS tags](#) and [Czech features](#).

- Czech uses all 17 universal POS categories, including particles ([tagset](#)). At present, more than 70 word types are tagged [tagset](#). This is a legacy of an existing Czech morphological analyzer and many of these words should probably belong to another category in UD; however, the exact list has yet to be worked out.
- The pronoun ([tagset](#)) vs. determiner ([tagset](#)) distinction is based on word lists because the traditional grammar does not define determiners. In general, words that inflect for gender, to be able to agree with a modified noun, are tagged [tagset](#), even if they act independently in a given sentence; that includes possessives. Pronominal quantifiers (which the traditional grammar includes in numerals) are [tagset](#) as well.
- Czech has just one auxiliary verb ([tagset](#)), *byť* ("to be"), but lemmas *byvat* and *byvatel* are also possible. They are in fact just variants of *byť*, but they are separate lemmas because the morphological process that relates them to *byť* is considered derivational. The auxiliary verb is used in

universaldependencies.org/treebanks/cs_poses/cs_poses-learn-Cases.html

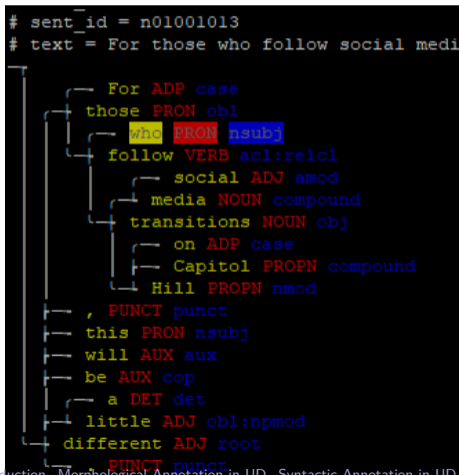
Paradigm <i>n</i>	Nom	Acc	Dat	Gen	Loc	Ins
Gender=Masc, Neut Number=Sing Person=2		memou				
Gender=Masc, Neut Number=Sing Person=3		emou			memě	memě
Gender=Masc Number=Sing Person=3	ji, jí, mě, tě, ně	emou; memou; mou; emou;		jeho, jejího	memě; emě;	memě; memě;
Gender=Masc Number=Dual Person=3	ji, j	memi		jejich, jejich	memě; emě;	memě; memě
Gender=Masc Number=Plur Person=3	ji, je, je	memi; memi		jejich, jejich; jich, jich;	memě; jichě	memi; memi; mi
Gender=Fem Number=Sing Person=3	ji, je, je, je	emě; memi; emě		jejich, jejich; jich, jich;	memě	memě; emě; memě; emě
Gender=Fem Number=Dual Person=3	ji	memi		jejich, jejich		memi
Gender=Fem Number=Plur Person=3	ji, je	memi; memě		jejich, jejich	memě	memi
Gender=Neut Number=Sing Person=3	o, je, o'			jeho, jejího	memě	memě
Gender=Neut Number=Dual Person=3				eho		
Gender=Neut Number=Plur Person=3	o, o'			jichě		

universaldependencies.org/treebanks/fa-comparison.html

Pronouns, Determiners, Quantifiers Pronouns, Determiners, Quantifiers Pronouns, Determiners, Quantifiers Pronouns, Determiners, Quantifiers

- PronType**
- Art**
 - DET: ly
- Con**
 - ADV: solum, alias
 - DET: alia, aliud, alius, alio, alterius, aliam, alii, sola, aliorum, utrumque
- Dem**
 - ADV: ecce, ita, tam, idem
 - DET: hoc, ipsum, haec, his, illud, ipsa, ipse, ipsius, ipso, idem
 - PRON: quod
- Dem**
 - ADV: taliter, tantum, hic, ita, tam, ecce, tantu, hece, hinc, ic
 - DET: ipsa, superscripta, ipsius, hanc, ipse, illa, ipso, superscripte, superscripto, hec
- Dem**
 - ADV: sic, tamen, t tantum, tunc, ita, hic, tanto, preterit
 - DET: hoc, illa, illud, ipsium, huius, his ipse, ille
 - PART: itaque, Eco

- <http://udapi.github.io/>
- <https://ufal.mff.cuni.cz/~zeman/vyuka/deptreebanks/NPFL075-working-with-UD.pdf>



Summary of Tools

- PML-TQ (tree query)
- UDPipe, Stanza (parsers)
- Annotatrix, Tred (annotation + tree view)
- UD web infrastructure (documentation, validation, statistics)
- Udapi (searching, batch processing, visualization)
- Many other tools listed on the UD website

<https://universaldependencies.org/>

Thanks!
Grazie!

<https://universaldependencies.org/>