

# The State of Universal Dependencies

Marie-Catherine de Marneffe, Joakim Nivre, Daniel Zeman

📅 March 16, 2023



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



unless otherwise stated

# 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

# Word Segmentation

*Let's go to the sea.*

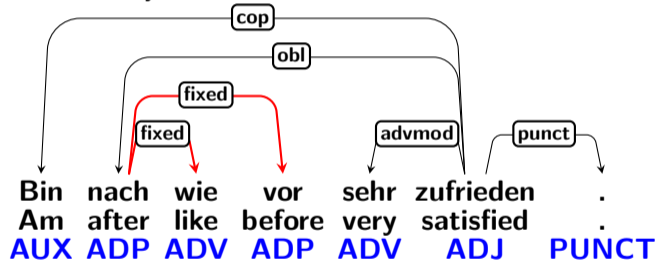
Vámonos al mar . Vamos nos a el mar .  
VERB? X NOUN PUNCT VERB PRON ADP DET NOUN PUNCT

- **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		

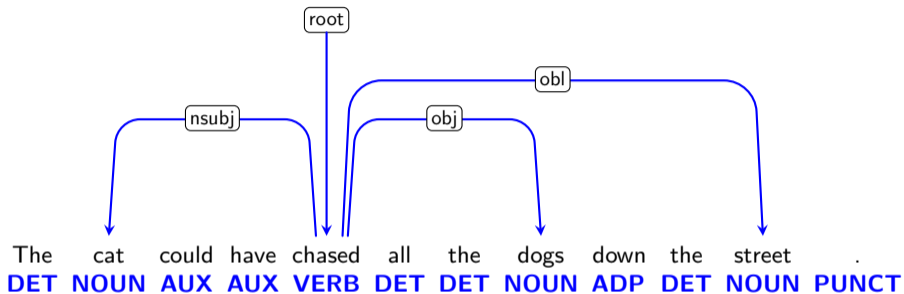
- Taxonomy of 17 universal POS tags
- All languages use the same inventory
  - Not all tags have to be used by all languages
  - Need extensions? Use features!

# 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

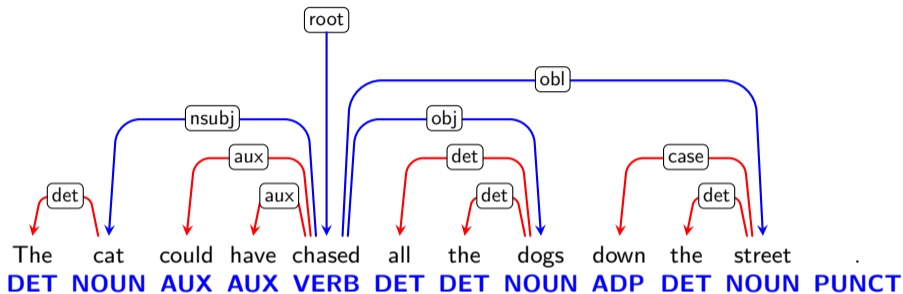
- 24 features, each with a number of possible *values*
- Languages select relevant features
- May add language-specific features or values

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

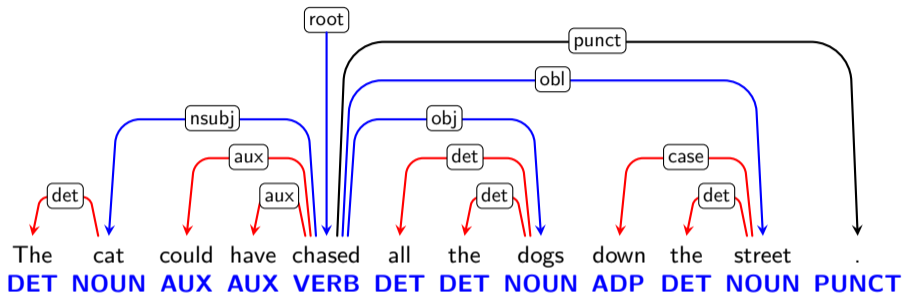


- Content words are related by dependency relations

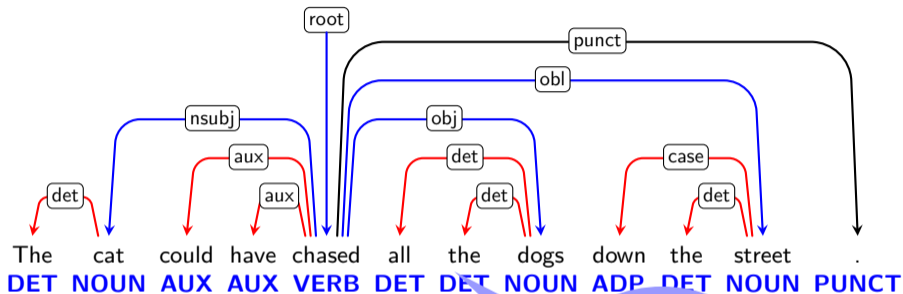




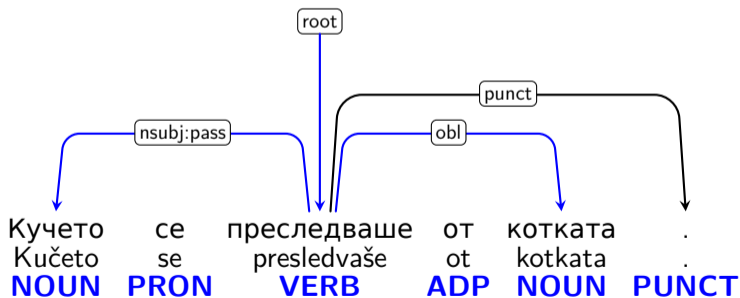
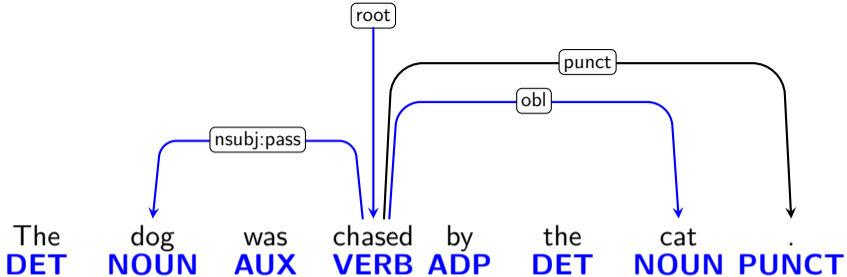
- Content words are related by dependency relations
- Function words attach to closest content words

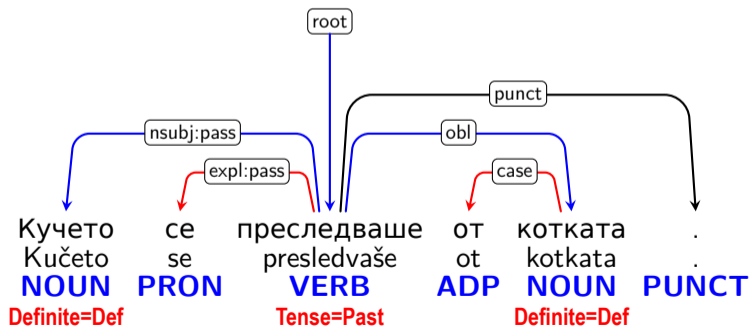
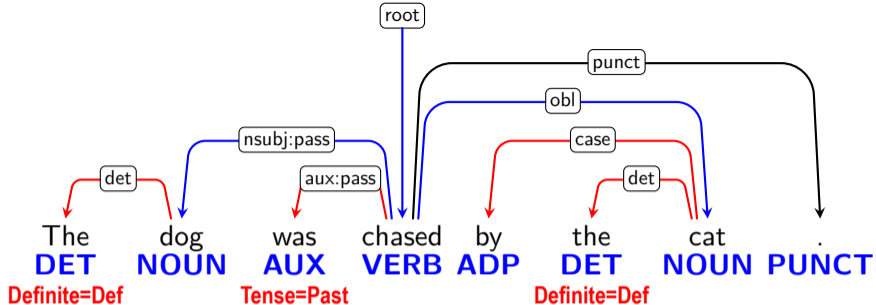


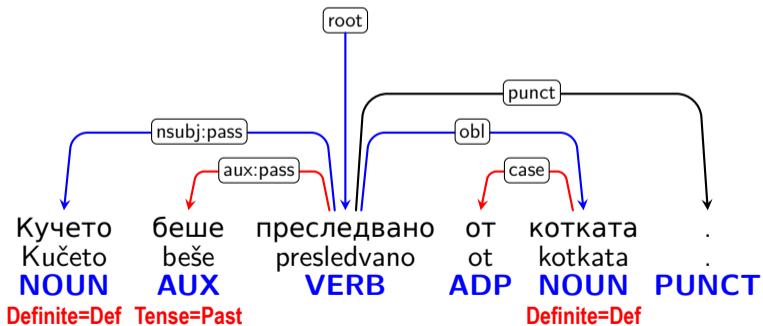
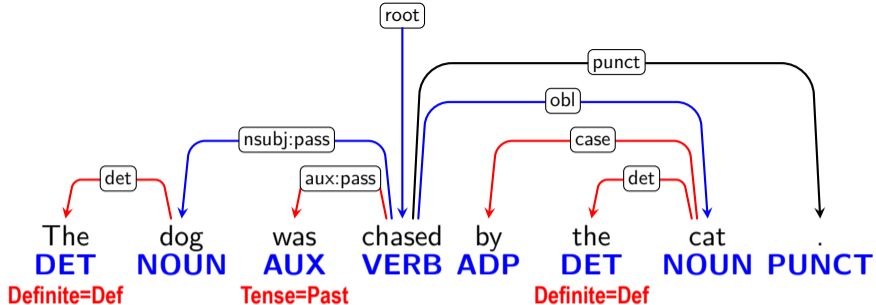
- Content words are related by dependency relations
- Function words attach to closest content words
- Punctuation attach to head of phrase or clause

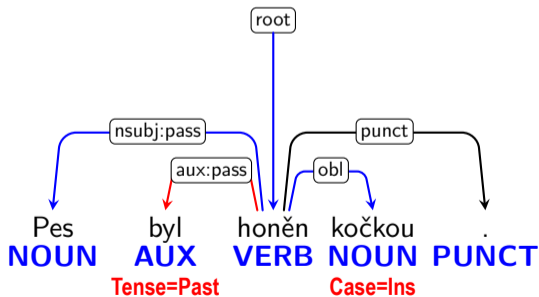
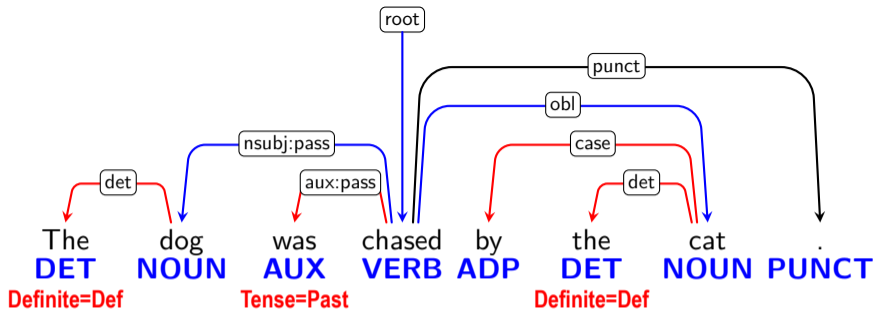


Not  
"dependency"  
in the strictly  
syntactic  
sense!









## 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

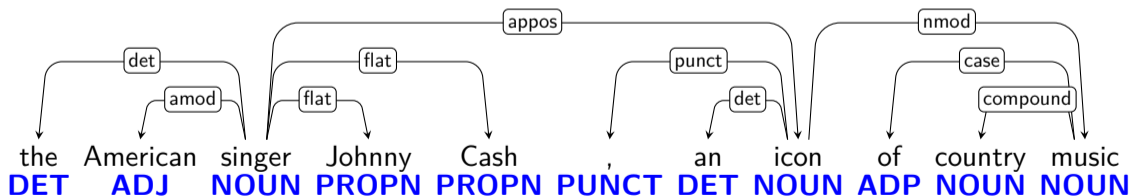
nummod

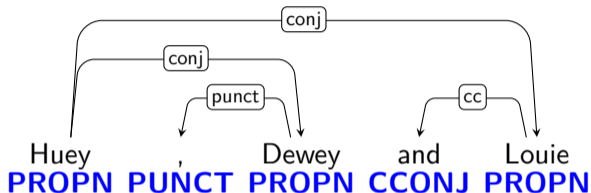
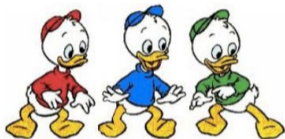
### Function

det

case

clf





- 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

<b>Relation</b>	<b>Explanation</b>
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

<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 ...)