



UPPSALA
UNIVERSITET

The Logic of Universal Dependencies

Joakim Nivre

Uppsala University

Department of Linguistics and Philology

Based on collaborative work with Marie-Catherine de Marneffe, Filip Ginter, Yoav Goldberg, Jan Hajic, Christopher Manning, Ryan McDonald, Natalia Silveira, Slav Petrov, Sampo Pyysalo, Sebastian Schuster, Reut Tsarfaty, Francis Tyers, Daniel Zeman and many others

Introduction

Introduction

Growing interest in multilingual and cross-lingual NLP

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

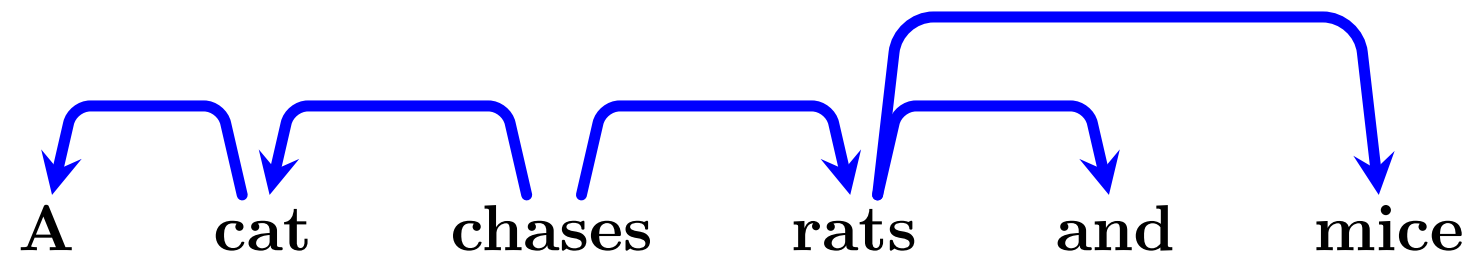
Introduction

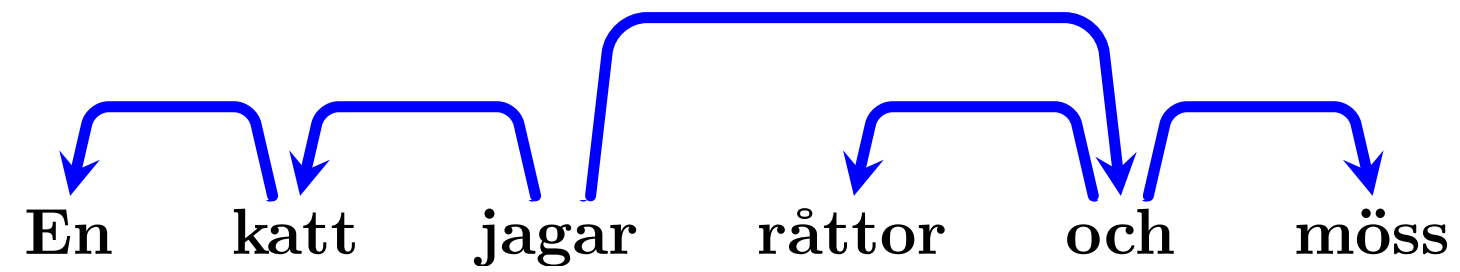
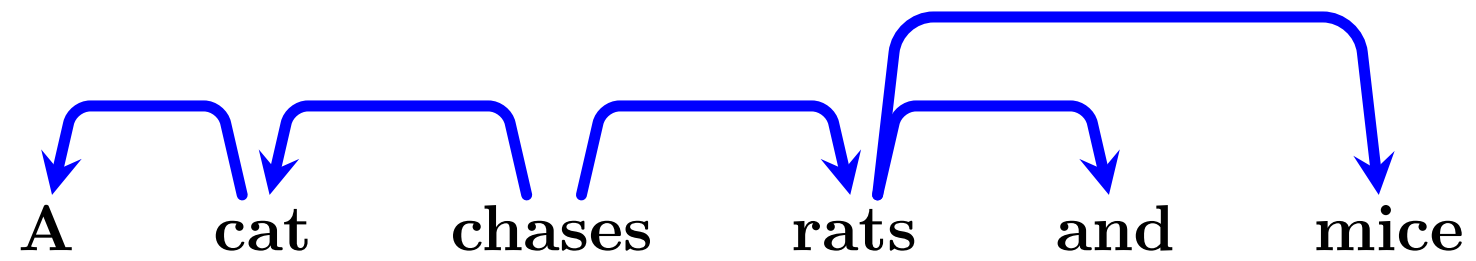
Growing interest in multilingual and cross-lingual NLP

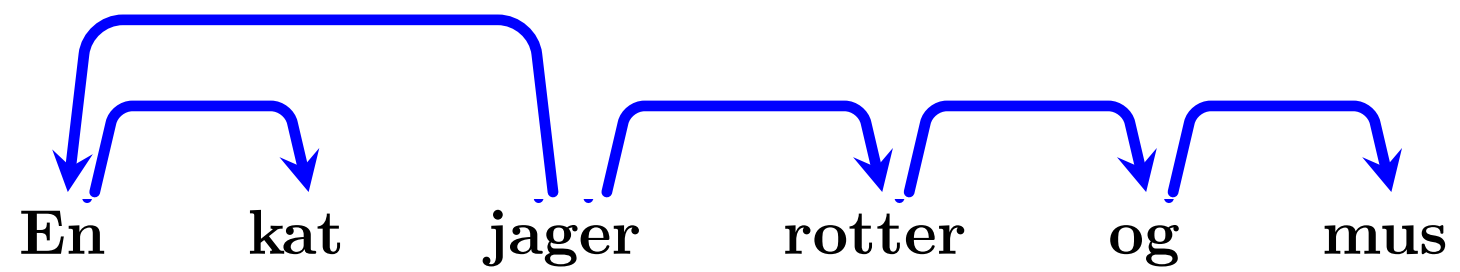
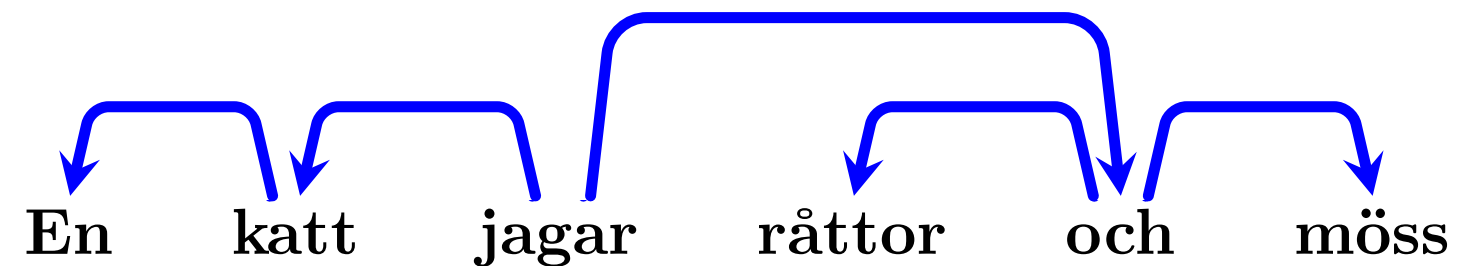
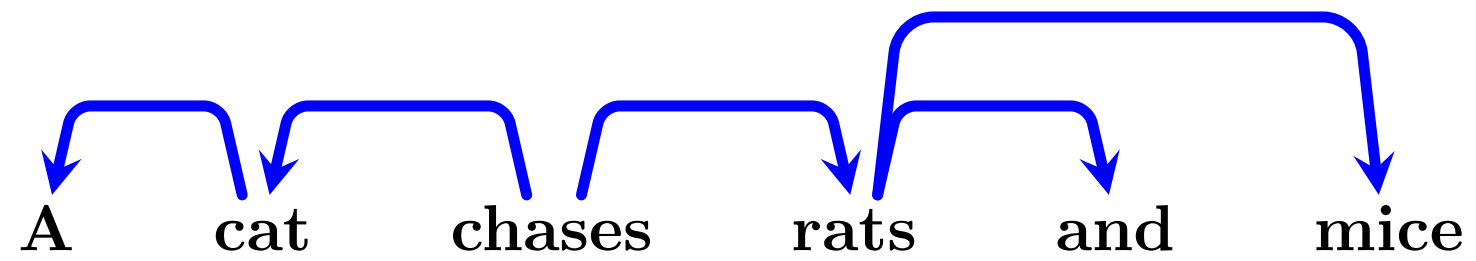
- Multilingual evaluation campaigns to test generality of approaches
- Cross-lingual learning to support low-resource languages

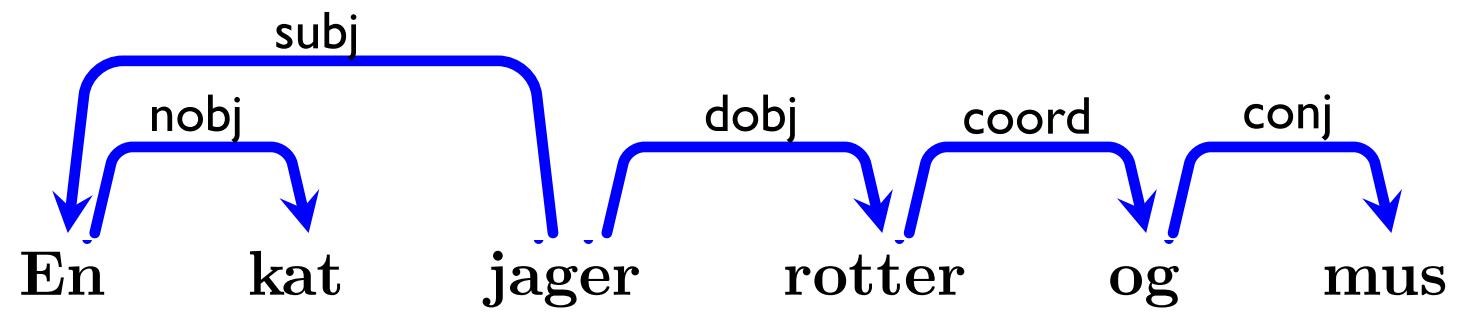
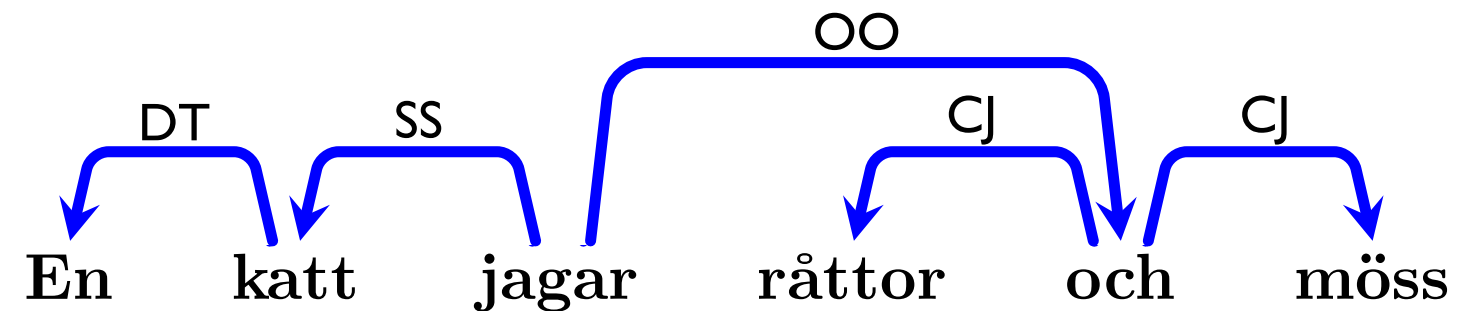
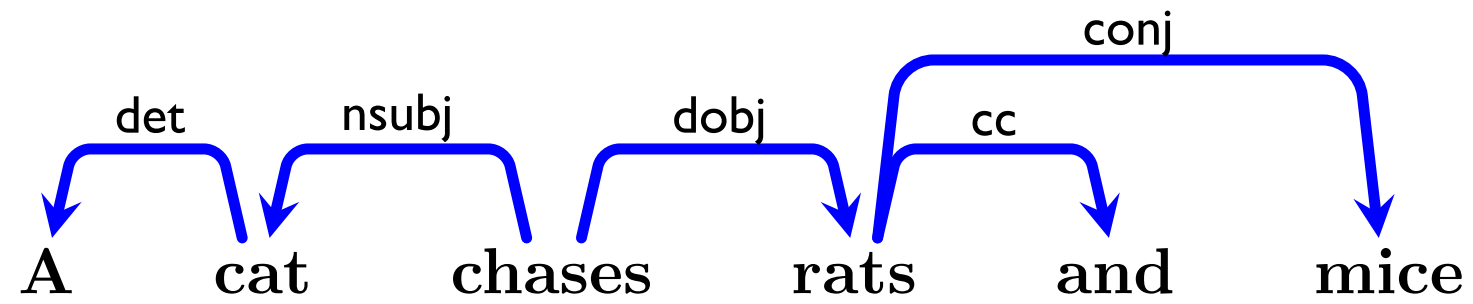
Growing awareness of methodological problems

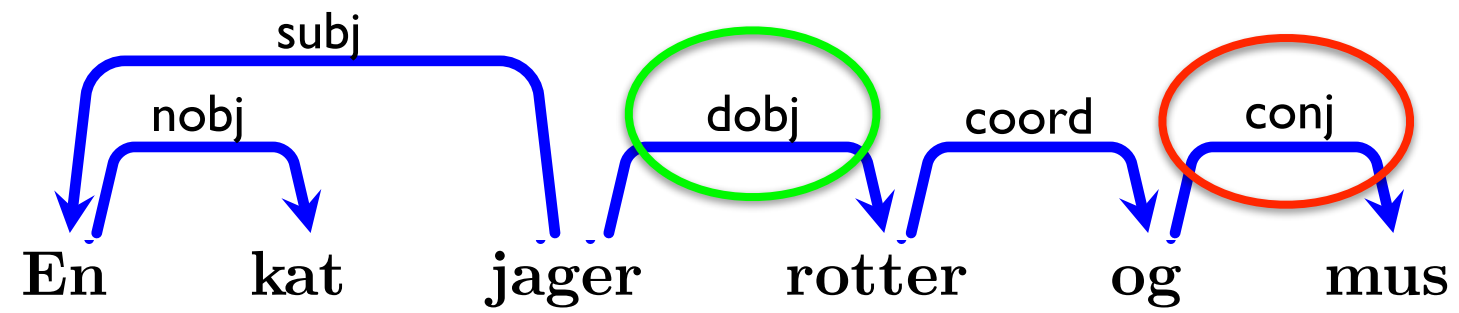
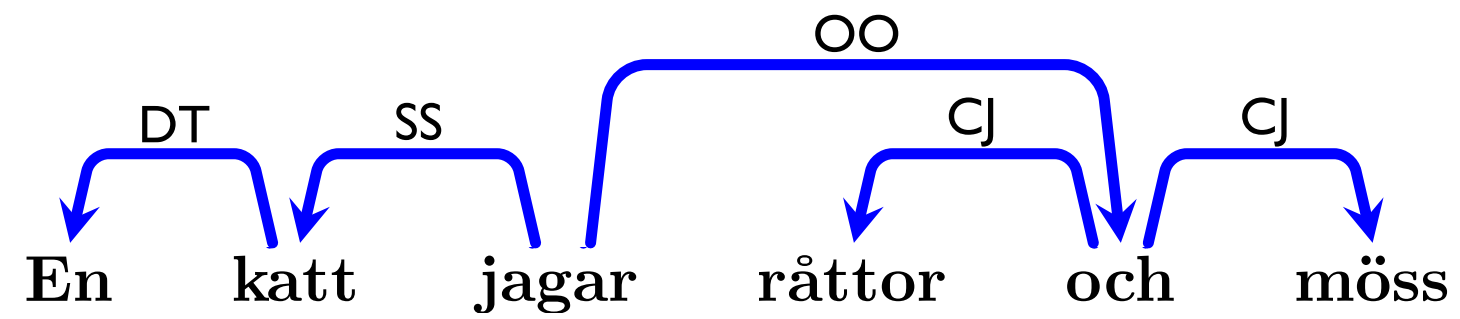
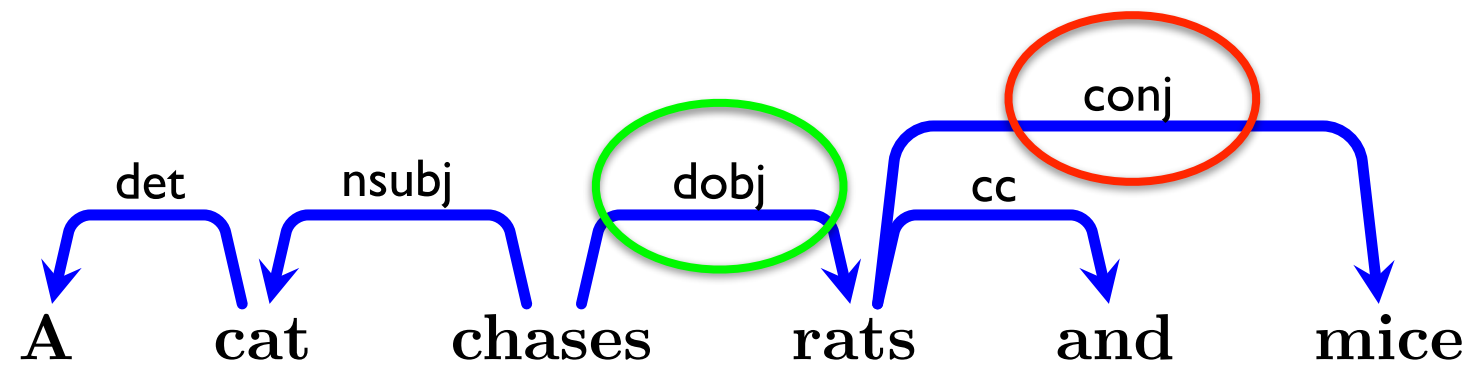
- Current NLP relies heavily on linguistic annotation
- Annotation guidelines vary across languages











Why is this a problem?

Why is this a problem?

- Hard to compare empirical results across languages

Why is this a problem?

- Hard to compare empirical results across languages
- Hard to usefully do cross-lingual structure transfer

Why is this a problem?

- Hard to compare empirical results across languages
- Hard to usefully do cross-lingual structure transfer
- Hard to evaluate cross-lingual learning

Why is this a problem?

- Hard to compare empirical results across languages
- Hard to usefully do cross-lingual structure transfer
- Hard to evaluate cross-lingual learning
- Hard to build and maintain multilingual systems

Why is this a problem?

- Hard to compare empirical results across languages
- Hard to usefully do cross-lingual structure transfer
- Hard to evaluate cross-lingual learning
- Hard to build and maintain multilingual systems
- Hard to make comparative linguistic studies

Why is this a problem?

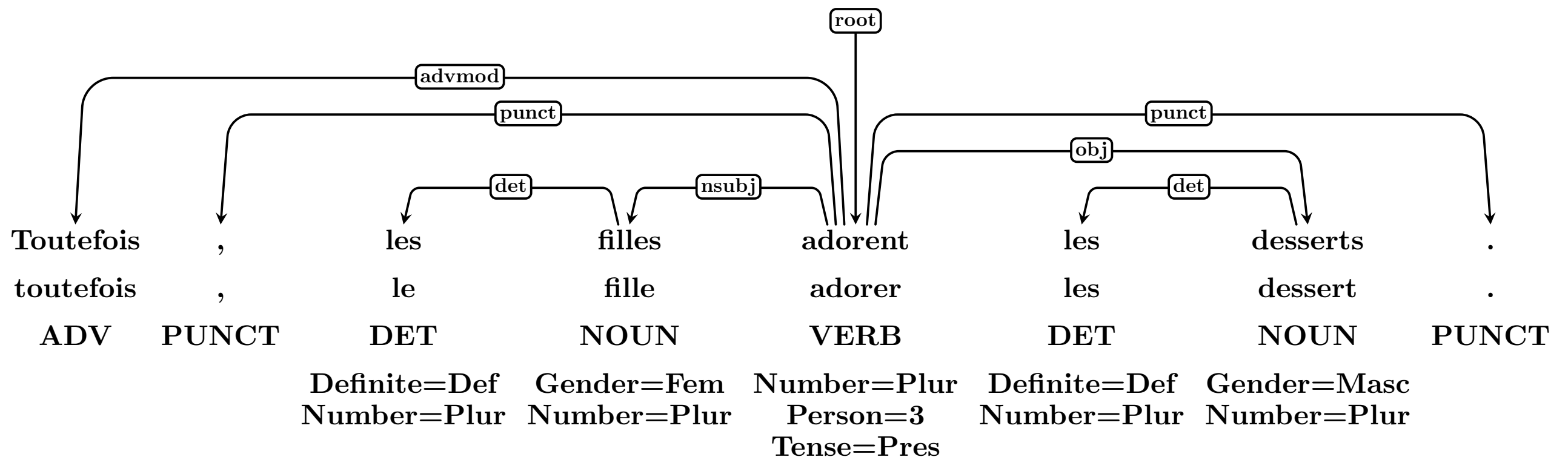
- Hard to compare empirical results across languages
- Hard to usefully do cross-lingual structure transfer
- Hard to evaluate cross-lingual learning
- Hard to build and maintain multilingual systems
- Hard to make comparative linguistic studies
- Hard to validate linguistic typology

Why is this a problem?

- Hard to compare empirical results across languages
- Hard to usefully do cross-lingual structure transfer
- Hard to evaluate cross-lingual learning
- Hard to build and maintain multilingual systems
- Hard to make comparative linguistic studies
- Hard to validate linguistic typology
- Hard to make progress towards a universal parser

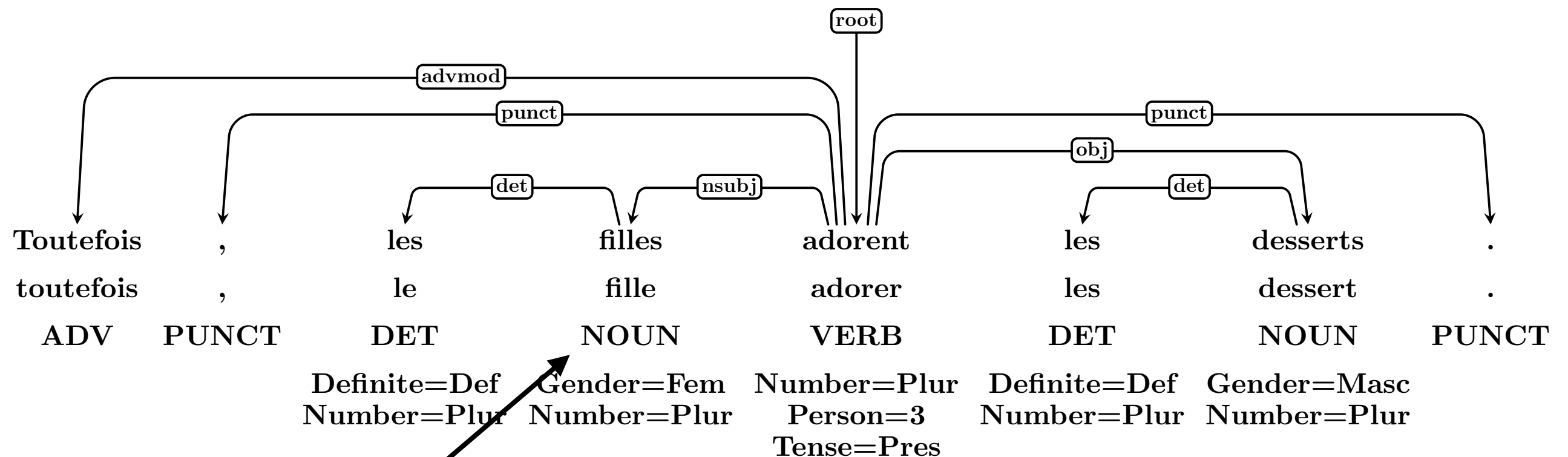
Universal Dependencies

<http://universaldependencies.org>



Universal Dependencies

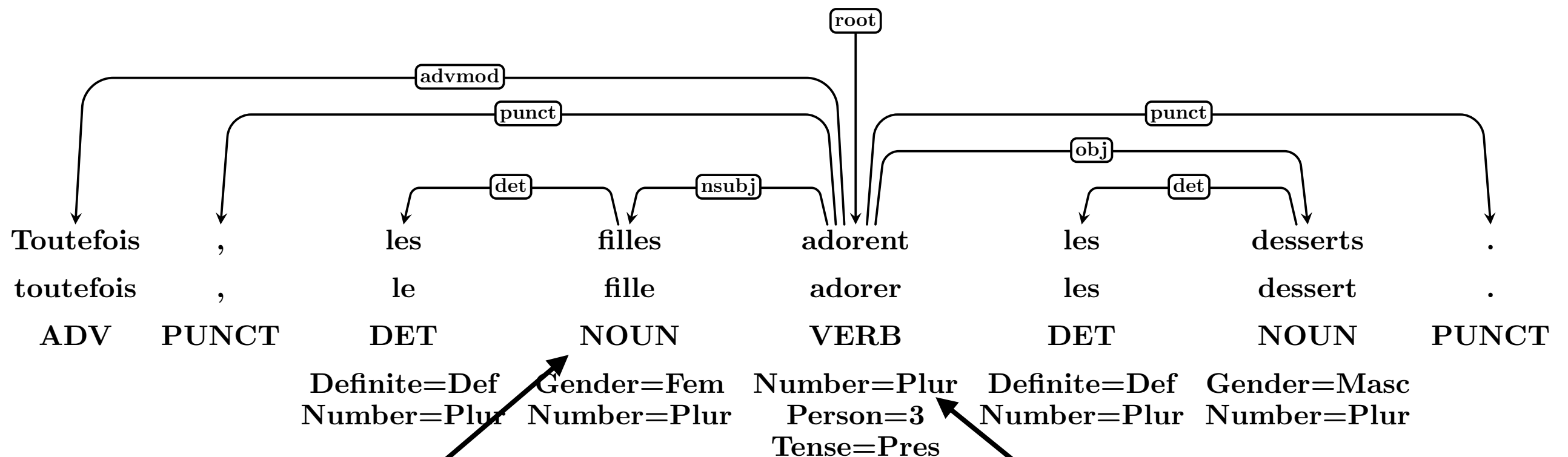
<http://universaldependencies.org>



Part-of-speech tags 

Universal Dependencies

<http://universaldependencies.org>



Part-of-speech tags

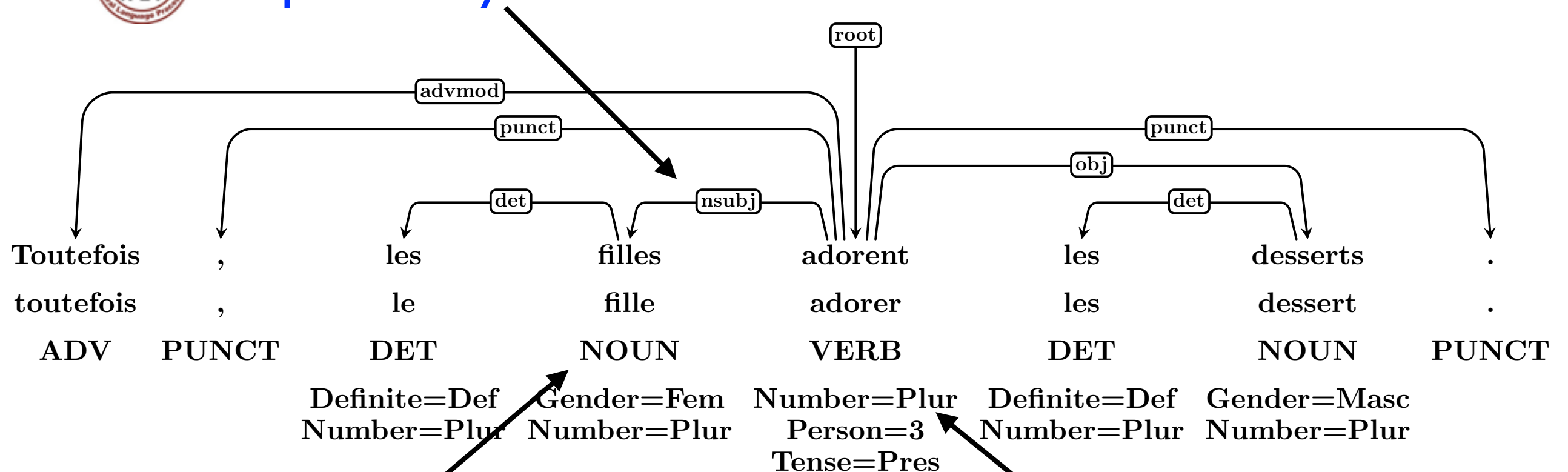
Morphological features

Universal Dependencies

<http://universaldependencies.org>



Dependency relations



Part-of-speech tags

Morphological features

Universal Dependencies

<http://universaldependencies.org>

Brief history of UD:


- Kick-off meeting at EACL in Gothenburg, April 2014
- Guidelines v1, October 2014
- Treebank releases every 6 months (v1.0–v1.4)
- Guidelines v2, December 2016
- Treebank release v2.0, March 2017

Open community effort – anyone can contribute!

Universal Dependencies


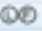
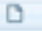
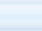
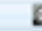
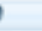
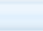








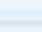
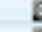

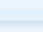


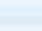


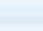








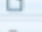

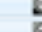




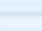
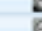

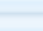


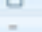
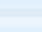


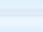


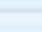


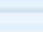


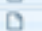
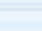

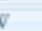
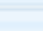


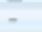
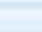


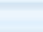

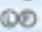
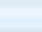


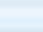


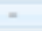
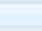
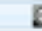
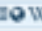
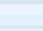


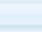
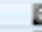
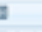
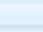










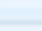

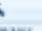
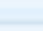



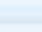


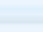


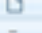
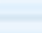


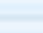


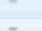
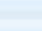
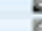
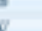
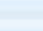


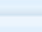
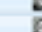

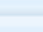

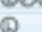
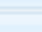

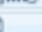
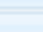

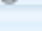

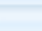
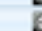

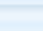



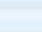

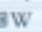
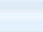




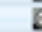



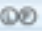
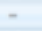
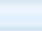


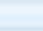



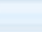

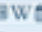
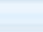

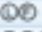
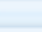
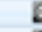
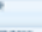
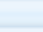


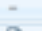












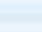


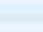


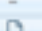
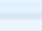
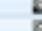

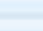


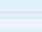


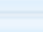


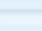


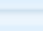








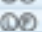
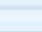


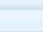


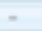
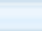


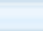


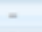
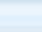

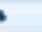
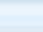


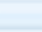
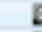
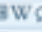
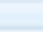

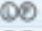
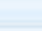

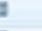
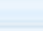















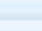


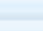


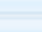
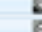

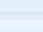


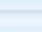


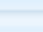




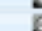




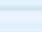


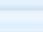

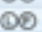
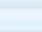


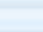

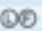

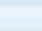


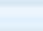



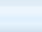


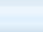


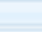

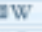
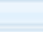

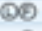







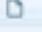

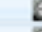




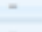
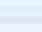


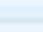




<http://universaldependencies.org>

Brief history of UD:

- Kick-off meeting at EACL in Gothenburg, April 2014
- Guidelines v1, October 2014
- Treebank releases every 6 months (v1.0–v1.4)
- Guidelines v2, December 2016
- Treebank release v2.0, March 2017 

Open community effort – anyone can contribute!

UD Treebanks

▶		Ancient Greek	202K						
▶		Ancient Greek-PROIEL	211K		-				
▶		Arabic	242K		-				
▶		Arabic-NYUAD	629K		-				
▶		Arabic-PUD	20K		-				
▶		Basque	121K						
▶		Belarusian	6K		-				
▶		Bulgarian	156K						
▶		Buryat	10K		-				
▶		Catalan	530K						
▶		Chinese	123K						
▶		Chinese-PUD	21K		-				
▶		Coptic	3K						
▶		Croatian	197K		-				
▶		Czech	1,330K						
▶		Czech-CAC	493K						
▶		Czech-CLTT	37K						
▶		Czech-PUD	18K						
▶		Danish	100K						
▶		Dutch	209K		-				
▶		Dutch-LassySmall	101K		-				
▶		English	254K						
▶		English-ESL	88K						
▶		English-LinES	82K						
▶		English-PUD	21K						
▶		English-ParTUT	49K						
▶		Estonian	47K		-				
▶		Finnish	202K						
▶		Finnish-FTB	159K		-				
▶		Finnish-PUD	15K		-				
▶		French	391K						
▶		French-FTB	556K		-				
▶		French-PUD	24K		-				
▶		French-ParTUT	27K						
▶		French-Sequoia	68K		-				
▶		Galician	138K						
▶		Galician-TreeGal	23K						
▶		German	287K		-				
▶		German-PUD	20K		-				
▶		Gothic	55K		-				
▶		Greek	61K						
▶		Hebrew	115K		-				
▶		Hindi	351K		-				
▶		Hindi-PUD	23K		-				
▶		Hungarian	42K						
▶		Indonesian	121K		-				
▶		Indonesian-PUD	25K		-				
▶		Irish	23K						
▶		Italian	273K						
▶		Italian-PUD	22K		-				
▶		Italian-ParTUT	39K						
▶		Japanese	186K						
▶		Japanese-KTC	189K						
▶		Japanese-PUD	26K		-				

UD Treebanks

▶	🇬🇷	Ancient Greek	202K	👤	📄	🔍	📊	📁
▶	🇬🇷	Ancient Greek-PROIEL	211K	👤	-	🔍	📊	📁
▶	🇸🇦	Arabic	242K	👤	-	🔍	📊	📁
▶	🇸🇦	Arabic-NYUAD	629K	👤	-	🔍	📊	📁
▶	🇸🇦	Arabic-PUD	20K	👤	-	🔍	📊	📁
▶	🇪🇸	Basque	121K	👤	📄	🔍	📊	📁
▶	🇧🇪	Belarusian	6K	👤	-	🔍	📊	📁
▶	🇧🇬	Bulgarian	156K	👤	📄	🔍	📊	📁
▶	🇲🇳	Buryat	10K	👤	-	🔍	📊	📁
▶	🇪🇸	Catalan	530K	👤	📄	🔍	📊	📁
▶	🇨🇳	Chinese	123K	👤	📄	🔍	📊	📁
▶	🇨🇳	Chinese-PUD	21K	👤	-	🔍	📊	📁
▶	🇪🇬	Coptic	3K	👤	📄	🔍	📊	📁
▶	🇭🇷	Croatian	197K	👤	-	🔍	📊	📁
▶	🇨🇪	Czech	1,330K	👤	📄	🔍	📊	📁
▶	🇨🇪	Czech-CAC	493K	👤	📄	🔍	📊	📁
▶	🇨🇪	Czech-CLTT	37K	👤	📄	🔍	📊	📁
▶	🇨🇪	Czech-PUD	18K	👤	📄	🔍	📊	📁
▶	🇩🇰	Danish	100K	👤	📄	🔍	📊	📁
▶	🇳🇱	Dutch	209K	👤	-	🔍	📊	📁
▶	🇳🇱	Dutch-LassySmall	101K	👤	-	🔍	📊	📁
▶	🇺🇸	English	254K	👤	📄	🔍	📊	📁
▶	🇺🇸	English-ESL	88K	👤	📄	🔍	📊	📁
▶	🇺🇸	English-LinES	82K	👤	📄	🔍	📊	📁
▶	🇺🇸	English-PUD	21K	👤	📄	🔍	📊	📁
▶	🇺🇸	English-ParTUT	49K	👤	📄	🔍	📊	📁
▶	🇪🇪	Estonian	47K	👤	-	🔍	📊	📁
▶	🇫🇮	Finnish	202K	👤	📄	🔍	📊	📁
▶	🇫🇮	Finnish-FTB	159K	👤	-	🔍	📊	📁
▶	🇫🇮	Finnish-PUD	15K	👤	-	🔍	📊	📁
▶	🇫🇷	French	391K	👤	📄	🔍	📊	📁
▶	🇫🇷	French-FTB	556K	👤	-	🔍	📊	📁
▶	🇫🇷	French-PUD	24K	👤	-	🔍	📊	📁
▶	🇫🇷	French-ParTUT	27K	👤	📄	🔍	📊	📁
▶	🇫🇷	French-Sequoia	68K	👤	-	🔍	📊	📁
▶	🇪🇸	Galician	138K	👤	📄	🔍	📊	📁
▶	🇪🇸	Galician-TreeGal	23K	👤	📄	🔍	📊	📁
▶	🇩🇪	German	287K	👤	-	🔍	📊	📁
▶	🇩🇪	German-PUD	20K	👤	-	🔍	📊	📁
▶	🇩🇪	Gothic	55K	👤	-	🔍	📊	📁
▶	🇬🇷	Greek	61K	👤	📄	🔍	📊	📁
▶	🇮🇱	Hebrew	115K	👤	-	🔍	📊	📁
▶	🇮🇳	Hindi	351K	👤	-	🔍	📊	📁
▶	🇮🇳	Hindi-PUD	23K	👤	-	🔍	📊	📁
▶	🇮🇳	Hungarian	42K	👤	📄	🔍	📊	📁
▶	🇮🇩	Indonesian	121K	👤	-	🔍	📊	📁
▶	🇮🇩	Indonesian-PUD	25K	👤	-	🔍	📊	📁
▶	🇮🇪	Irish	23K	👤	📄	🔍	📊	📁
▶	🇮🇹	Italian	273K	👤	📄	🔍	📊	📁
▶	🇮🇹	Italian-PUD	22K	👤	-	🔍	📊	📁
▶	🇮🇹	Italian-ParTUT	39K	👤	📄	🔍	📊	📁
▶	🇯🇵	Japanese	186K	👤	📄	🔍	📊	📁
▶	🇯🇵	Japanese-KTC	189K	👤	📄	🔍	📊	📁
▶	🇯🇵	Japanese-PUD	26K	👤	-	🔍	📊	📁
▶	🇰🇿	Kazakh	10K	👤	📄	🔍	📊	📁
▶	🇰🇷	Korean	74K	👤	📄	🔍	📊	📁
▶	🇰🇷	Korean-PUD	22K	👤	-	🔍	📊	📁
▶	🇰🇷	Korean-Sejong	89K	👤	-	🔍	📊	📁
▶	🇹🇲	Kurmanji	10K	👤	-	🔍	📊	📁
▶	🇵🇪	Latin	29K	👤	📄	🔍	📊	📁
▶	🇵🇪	Latin-ITTB	291K	👤	-	🔍	📊	📁
▶	🇵🇪	Latin-PROIEL	171K	👤	-	🔍	📊	📁
▶	🇱🇻	Latvian	54K	👤	-	🔍	📊	📁
▶	🇱🇮	Lithuanian	5K	👤	-	🔍	📊	📁
▶	🇲🇹	Maltese	2K	👤	-	🔍	📊	📁
▶	🇳🇴	North Sami	55K	👤	-	🔍	📊	📁
▶	🇳🇴	Norwegian-Bokmaal	310K	👤	📄	🔍	📊	📁
▶	🇳🇴	Norwegian-Nynorsk	301K	👤	📄	🔍	📊	📁
▶	🇷🇺	Old Church Slavonic	57K	👤	-	🔍	📊	📁

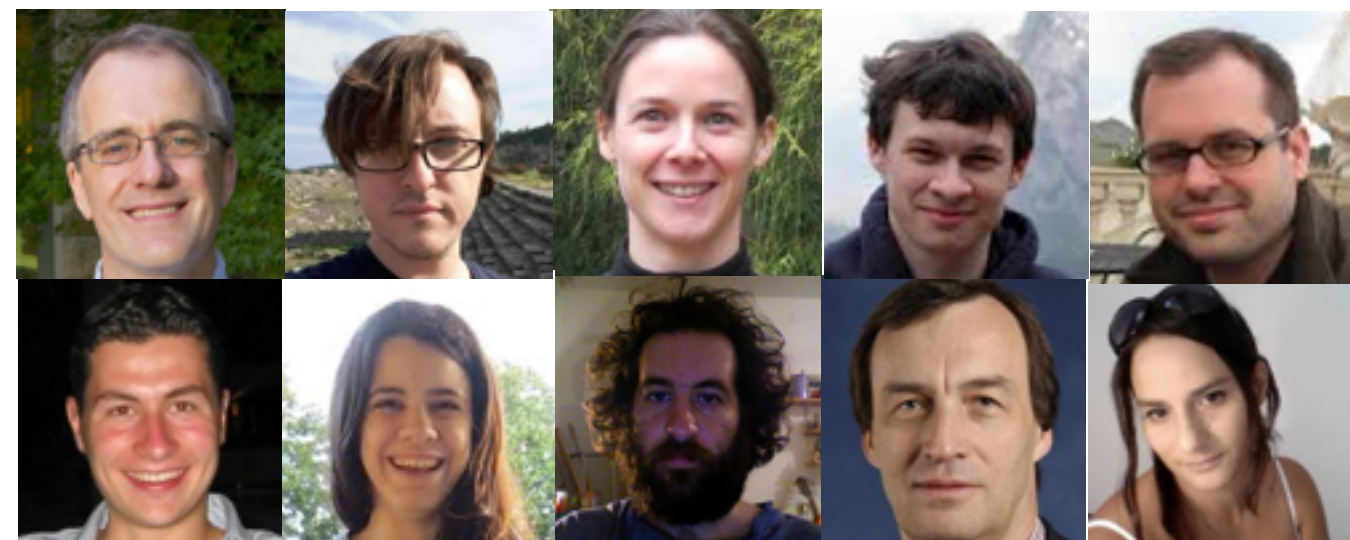
June, 2017:

- 56 languages
- 95 treebanks
- 13.9 million words
- 202 contributors
- 10791 downloads

Chief Cat Herder



Release and Documentation Task Force



Universal Guidelines Group

The UD Philosophy

The UD Philosophy

Maximize parallelism – but don't overdo it

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

The UD Philosophy

Maximize parallelism – but don't overdo it

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

Universal taxonomy with language-specific elaboration

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

Design Principles

Design Principles

Dependency

- Widely used in practical NLP systems
- Available in treebanks for many languages

Design Principles

Dependency

- Widely used in practical NLP systems
- Available in treebanks for many languages

Lexicalism

- Basic annotation units are words – syntactic words
- Words have morphological properties
- Words enter into syntactic relations

Morphology

Le chat chasse les chiens .

Morphology

Le	chat	chasse	les	chiens	.
le	chat	chasser	le	chien	.

- Lemma representing the semantic content of the word

Morphology

Le	chat	chasse	les	chiens	.
le	chat	chasser	le	chien	.
DET	NOUN	VERB	DET	NOUN	PUNCT

- Lemma representing the semantic content of the word
- Part-of-speech tag representing its grammatical class

Morphology

Le
le
DET

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

tiens
hien
PUNCT

- Lemma rep of the word
- Part-of-speech tag representing its grammatical class

Morphology

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

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

Morphology

Le
le
DET

Definite=Def
Gender=Masc
Number=Sing

Num

Lexical	Inflectional Nominal	Inflectional Verbal
PronType	Gender	VerbForm
NumType	Animacy	Mood
Poss	Number	Tense
Reflex	Case	Aspect
Foreign	Definite	Voice
Abbr	Degree	Evident
		Polarity
		Person
		Polite

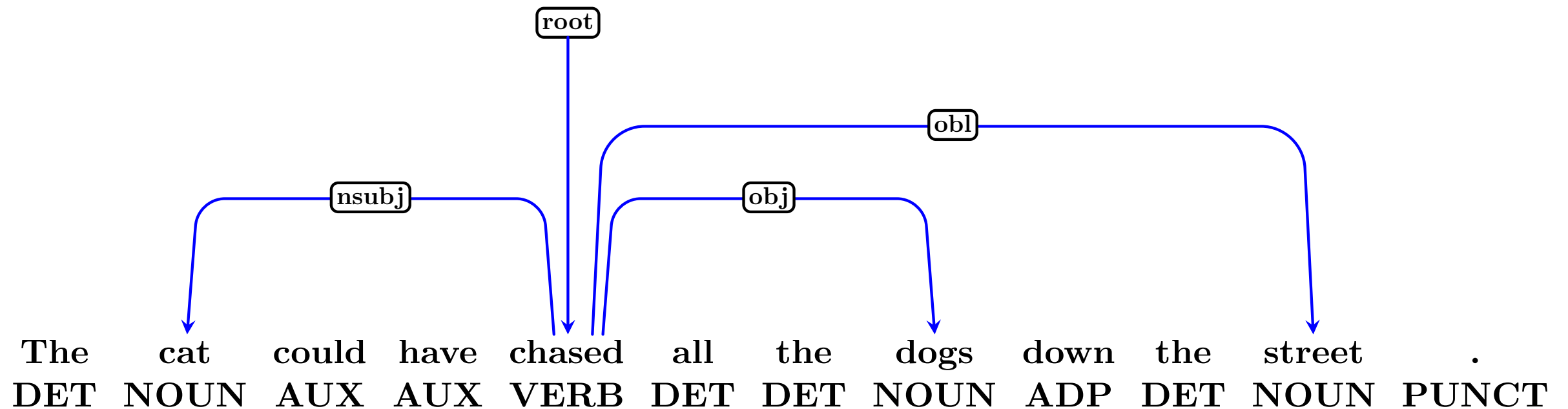
niens .
hien .
NOUN **PUNCT**
er=Masc
er=Plur

- Lemma representation
- Part-of-speech of the word
- Grammatical class
- Features representing lexical and grammatical properties of the lemma or the particular word form

Syntax

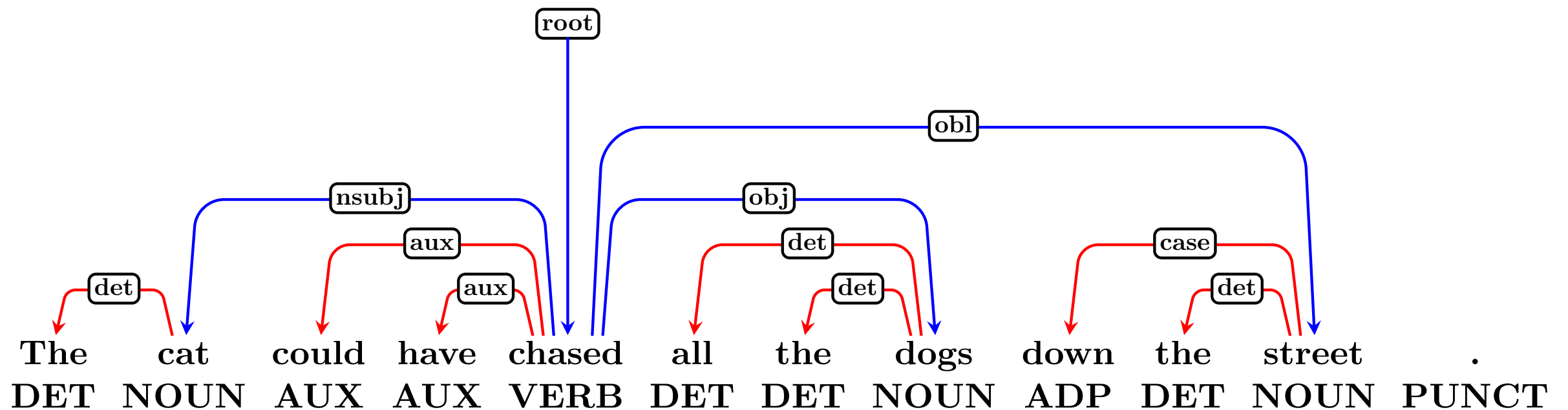
The	cat	could	have	chased	all	the	dogs	down	the	street	.
DET	NOUN	AUX	AUX	VERB	DET	DET	NOUN	ADP	DET	NOUN	PUNCT

Syntax



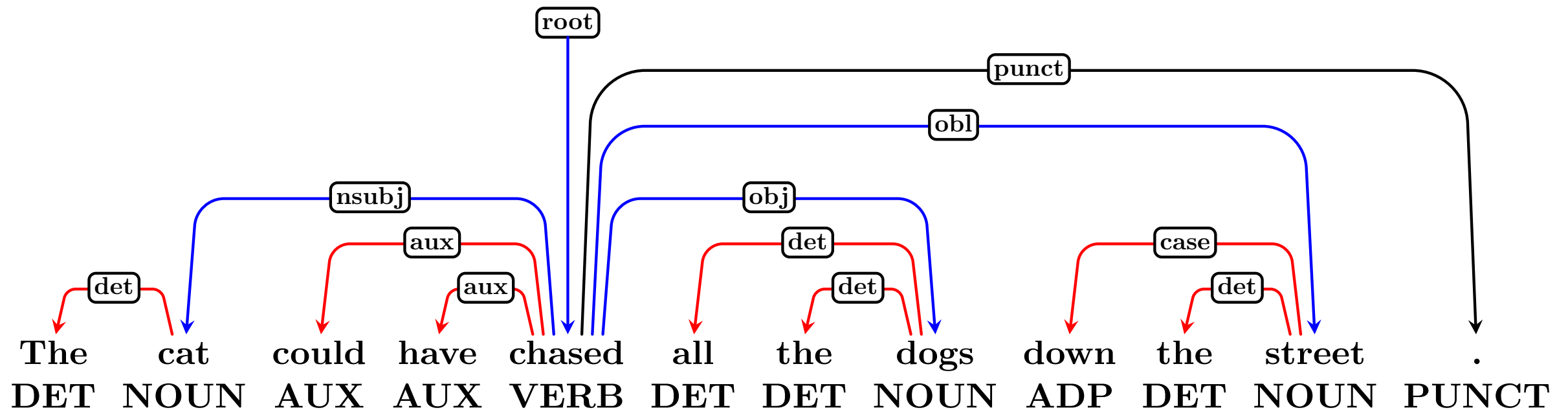
- Content words are related by dependency relations

Syntax



- Content words are related by dependency relations
- Function words attach to the content word they modify

Syntax



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

Syntactic Relations

Syntactic Relations

Taxonomy of 37 universal syntactic relations

- Three types of structures: nominals, clauses, modifiers
- Core arguments vs. other dependents (**not** complements vs. adjuncts)

Syntactic Relations

Taxonomy of 37 universal syntactic relations

- Three types of structures: nominals, clauses, modifiers
- Core arguments vs. other dependents (**not** complements vs. adjuncts)

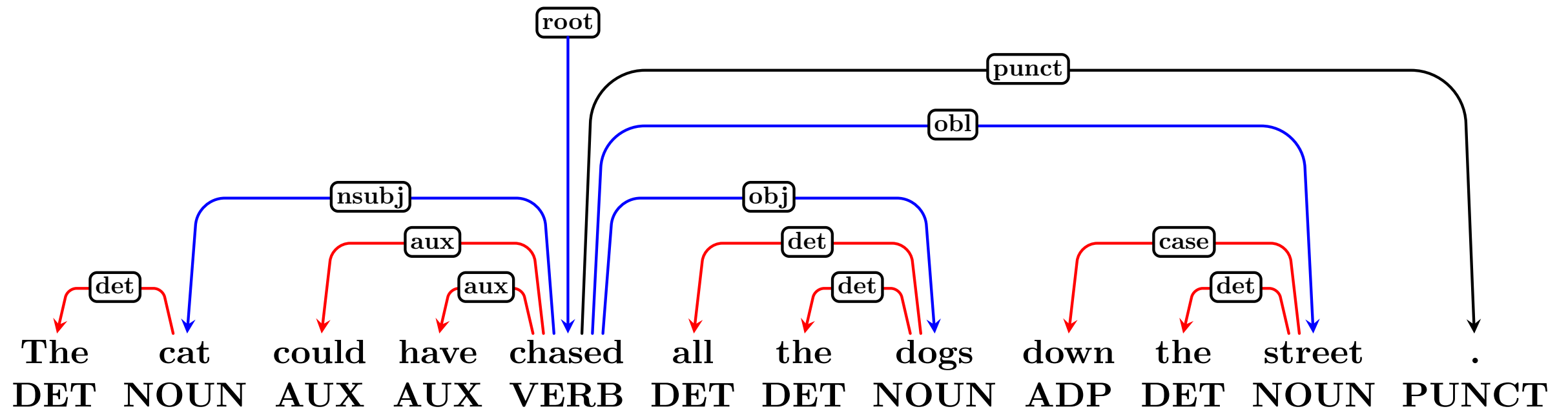
A two-level architecture

- Broad universal categories to allow cross-linguistic comparison
- Subtypes to capture language-specific phenomena

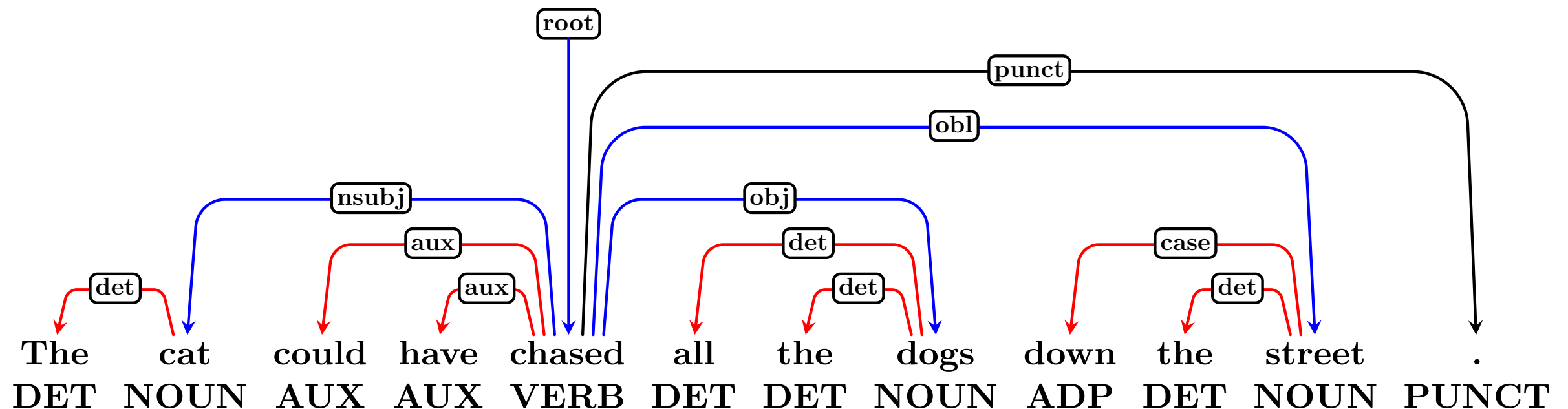
Syntactic Relations

	Nominal	Clause	Modifier Word	Function Word
Core Predicate Dep	nsubj obj iobj	csubj ccomp xcomp		
Non-Core Predicate Dep	obl vocative expl dislocated	advcl	advmod* discourse	aux cop mark
Nominal Dep	nmod appos nummod	acl	amod	det clf case
Coordination	MWE	Loose	Special	Other
conj cc	fixed flat compound	parataxis list	orphan goeswith reparandum	punct root dep

The Primacy of Content Words



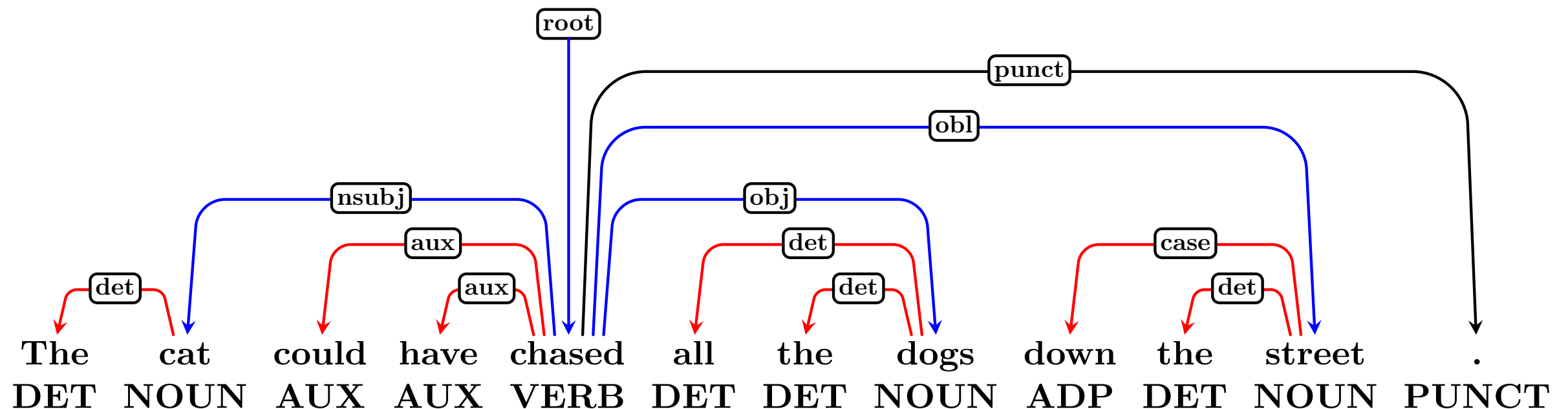
The Primacy of Content Words



Dubious Linguistics?

“Such an approach to the syntax of natural languages is contrary to most work in theoretical syntax in the past 35 years, regardless of whether this work is constituency- or dependency-based.” (Groß and Osborne, 2015)

The Primacy of Content Words



Bad for parsing?

“It is now fairly well known that, while dependency representations in which content words are made heads tend to help semantically oriented downstream applications, dependency parsing numbers are higher if you make auxiliary verbs heads [...] and if you make prepositions the head of prepositional phrases.” (De Marneffe et al., 2014)

Manning's Law



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

- 1 UD needs to be satisfactory on **linguistic analysis** grounds for individual languages.
- 2 UD needs to 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 suitable for **computer parsing** with high accuracy.
- 5 UD must be **easily comprehended** and used by a non-linguist, whether a language learner or an engineer with prosaic needs for language processing.
- 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.

Manning's Law



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

- 1 UD needs to be satisfactory on **linguistic analysis** grounds for individual languages.
- 2 UD needs to 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 suitable for **computer parsing** with high accuracy.
- 5 UD must be **easily comprehended** and used by a non-linguist, whether a language learner or an engineer with prosaic needs for language processing.
- 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.

Manning's Law



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

- 1 UD needs to be satisfactory on **linguistic analysis** grounds for individual languages.
- 2 UD needs to 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 suitable for **computer parsing** with high accuracy.
- 5 UD must be **easily comprehended** and used by a non-linguist, whether a language learner or an engineer with prosaic needs for language processing.
- 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.

Manning's Law

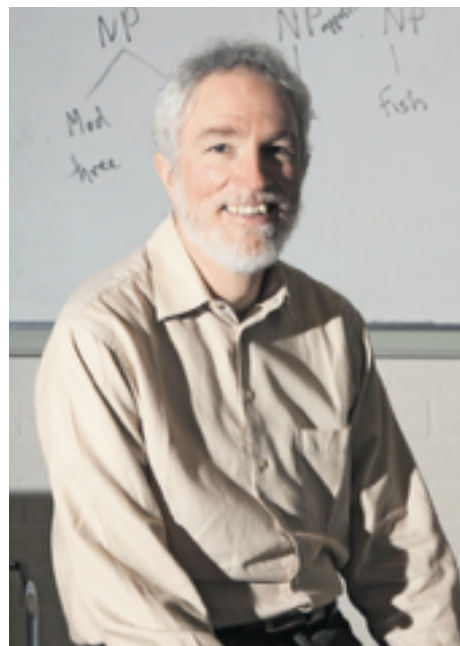


The secret to understanding the design of UD is to realize that it is a very subtle compromise between approximately 6 things:

- 1 UD needs to be satisfactory on **linguistic analysis** grounds for individual languages.
- 2 UD needs to 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 suitable for **computer parsing** with high accuracy.
- 5 UD must be **easily comprehended** and used by a non-linguist, whether a language learner or an engineer with prosaic needs for language processing.
- 6 UD must support well **downstream language understanding tasks** (relation extraction, reading comprehension, machine translation, ...).

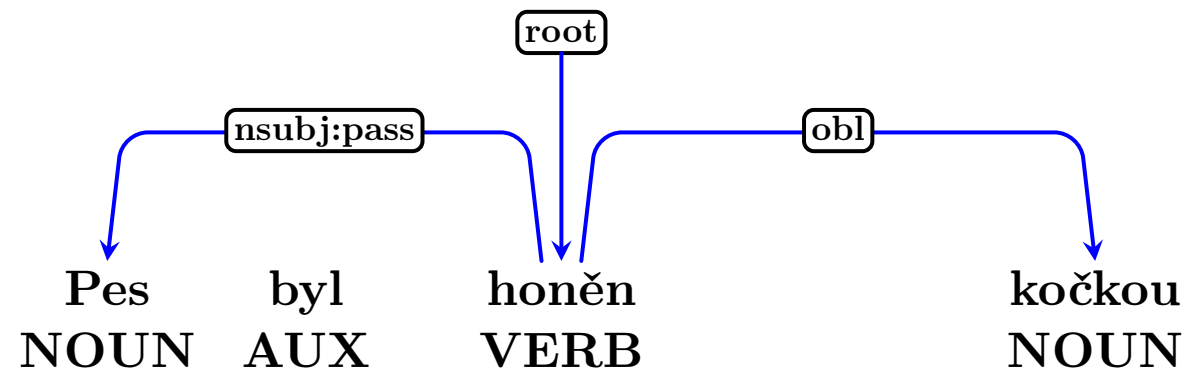
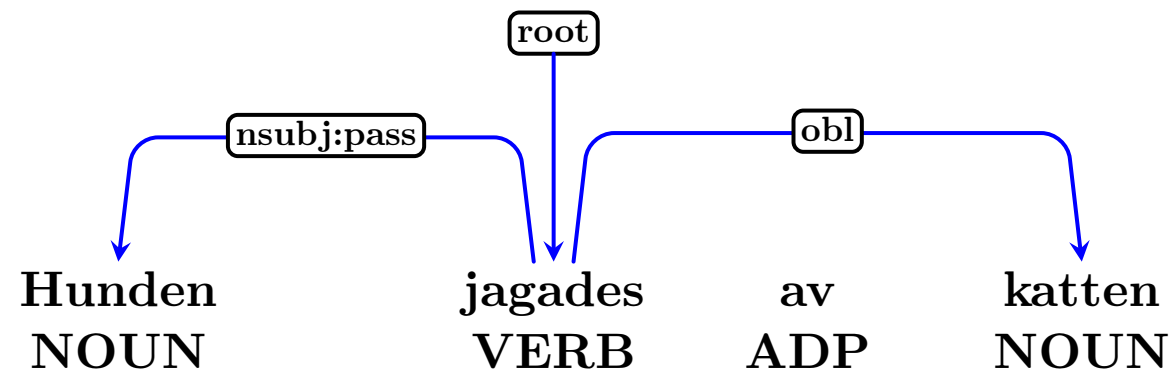
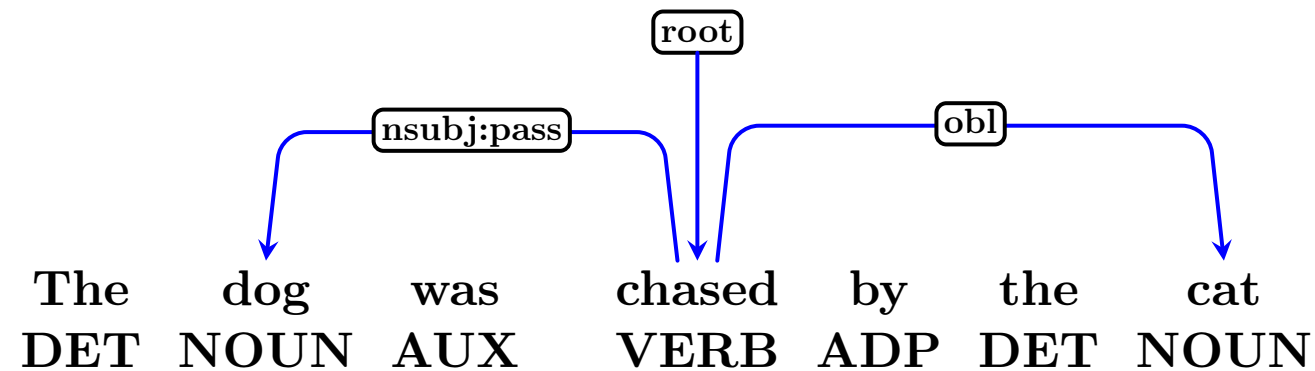
It's easy to come up with a proposal that improves UD on one of these dimensions. The interesting and difficult part is to improve UD while remaining sensitive to all these dimensions.

Linguistic Typology

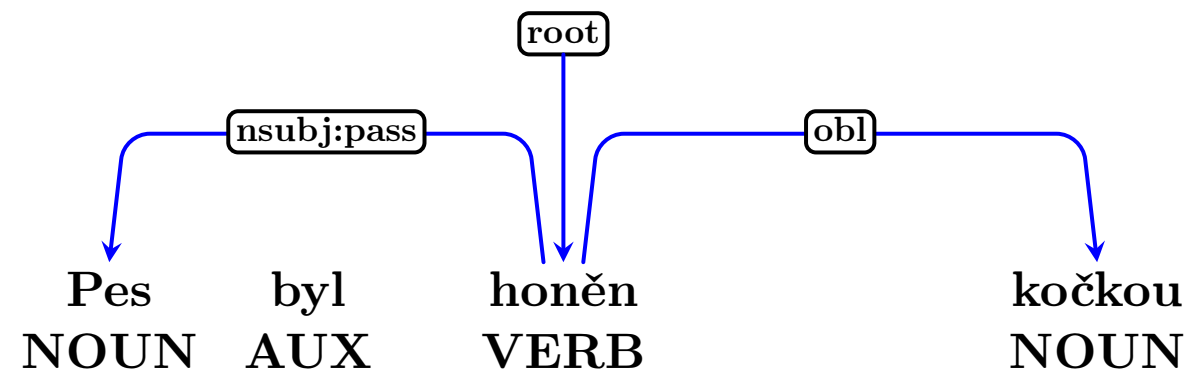
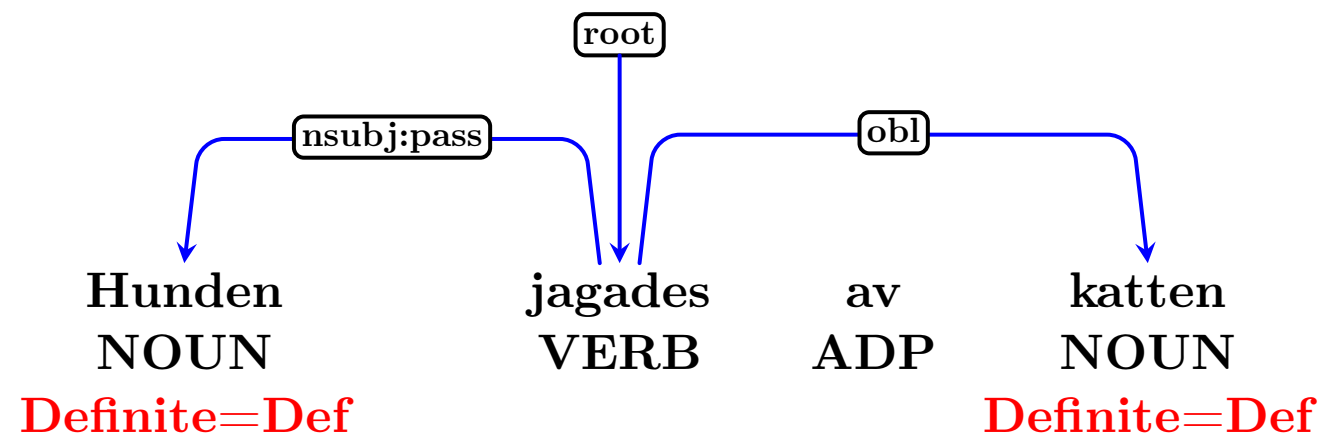
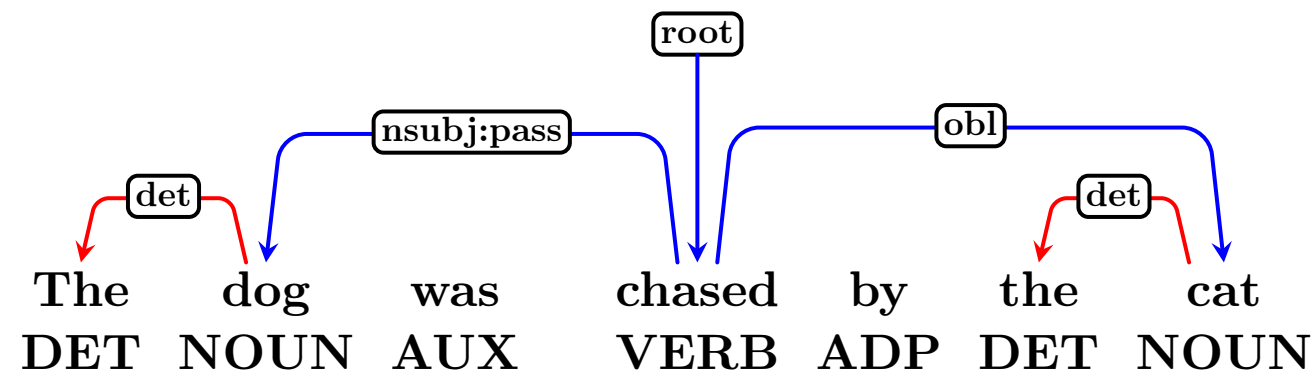


With contributions by William Croft

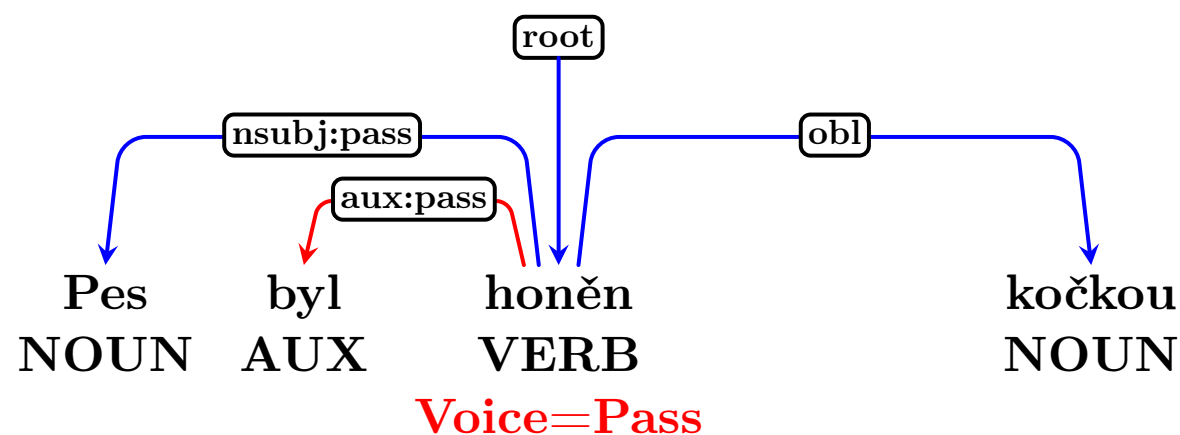
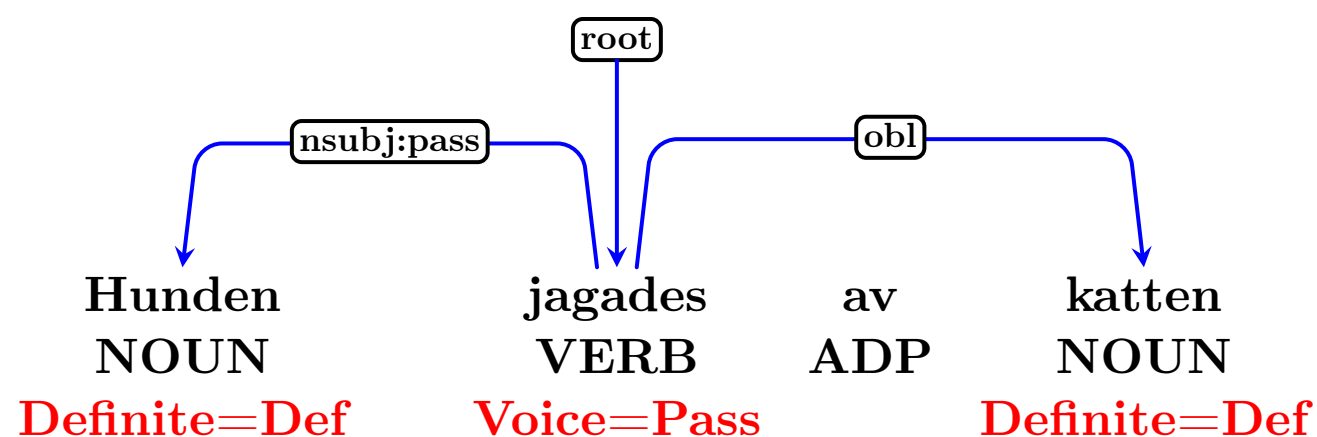
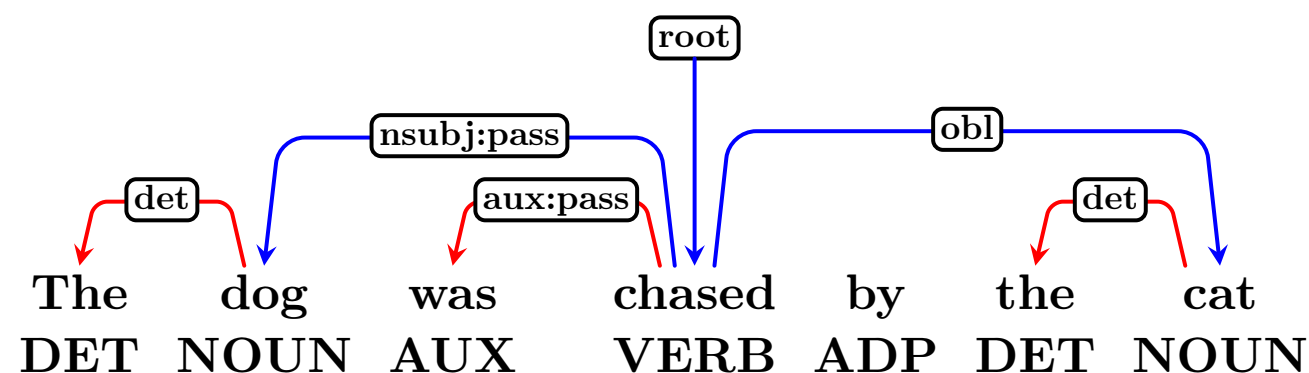
Similarities and Differences



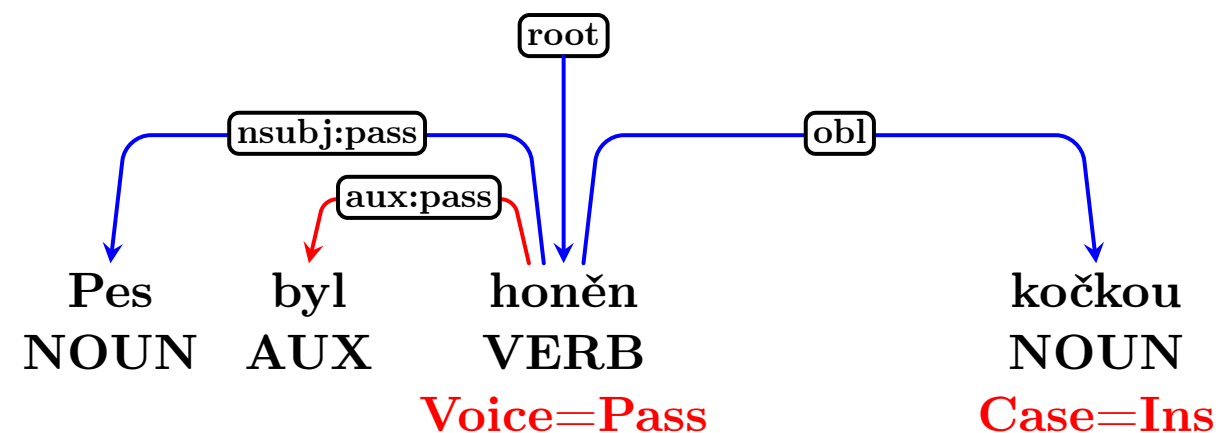
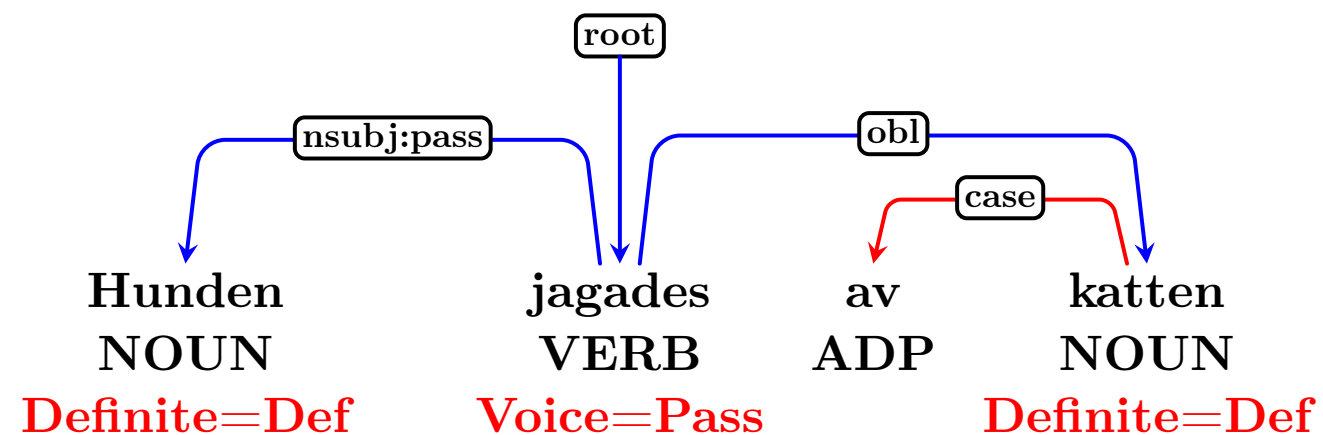
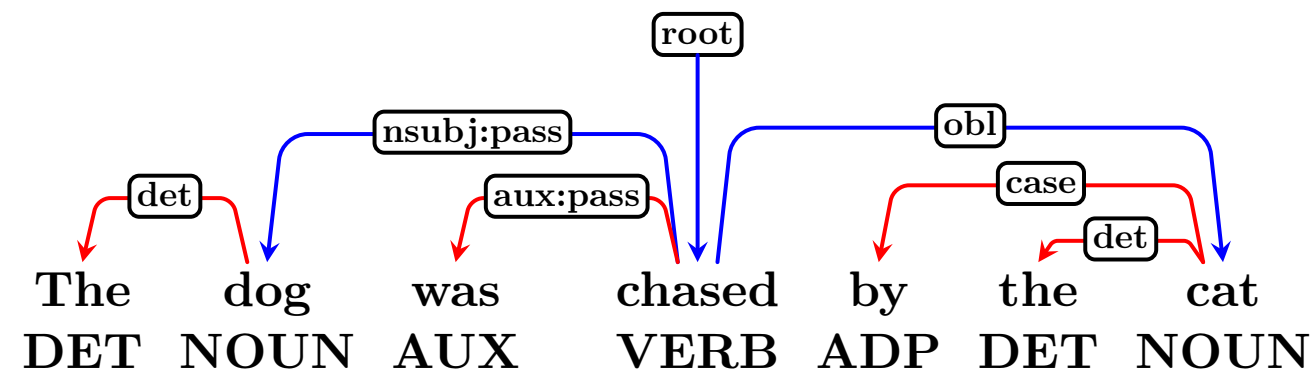
Similarities and Differences



Similarities and Differences



Similarities and Differences



What is universal?

What is universal?



- Descriptive categories – language-specific, cannot be equated across languages
- Comparative concepts – created by typologists for cross-linguistic comparison

Haspelmath (2010) Comparative Concepts and Descriptive Categories in Crosslinguistic Studies

What is universal?



- Descriptive categories – language-specific, cannot be equated across languages
- Comparative concepts – created by typologists for cross-linguistic comparison

Haspelmath (2010) Comparative Concepts and Descriptive Categories in Crosslinguistic Studies



- Construction – whatever structure is used to express a function, universal by definition
- Strategy – a specific, cross-linguistically definable structure used to express a function

Croft et al. (2017) Linguistic Typology Meets Universal Dependencies

Construction

Strategies

predication
of object
concept

inflected
copula

English:
Ivan is the best dancer.

Russian:
Ivan lučšij tancor

zero
copula/zero
inflection





Croft et al. (2017) Linguistic Typology Meets Universal Dependencies

- A universal annotation scheme should have a classification of constructions, not strategies, as its universal foundational layer



Croft et al. (2017) Linguistic Typology Meets Universal Dependencies

- A universal annotation scheme should have a classification of constructions, not strategies, as its universal foundational layer
- However, common recurrent strategies can and should be included as well



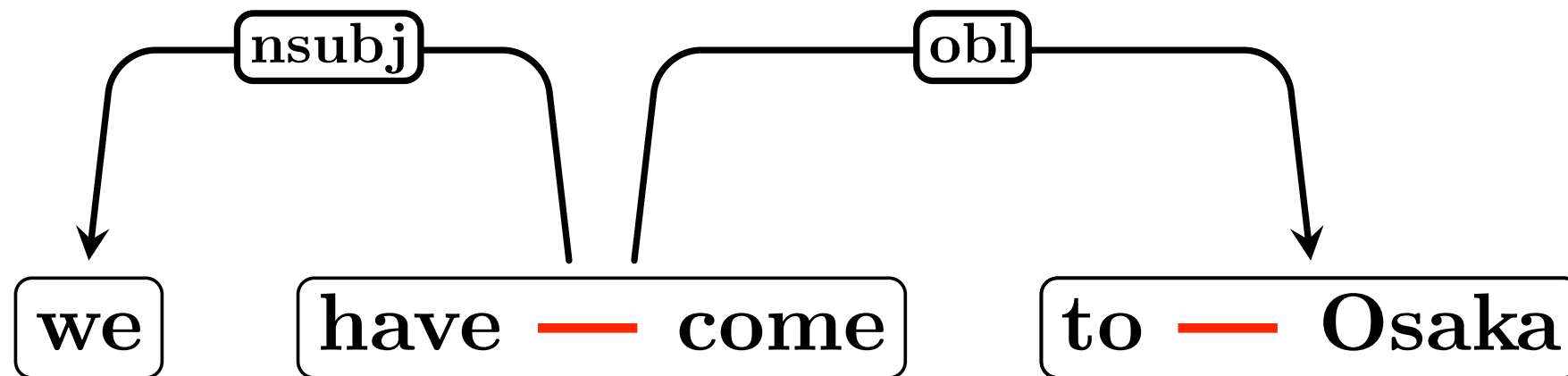
Croft et al. (2017) Linguistic Typology Meets Universal Dependencies

- A universal annotation scheme should have a classification of constructions, not strategies, as its universal foundational layer
- However, common recurrent strategies can and should be included as well
- The primacy of content words in UD is good, because function words such as copulas are strategies and not found in every language



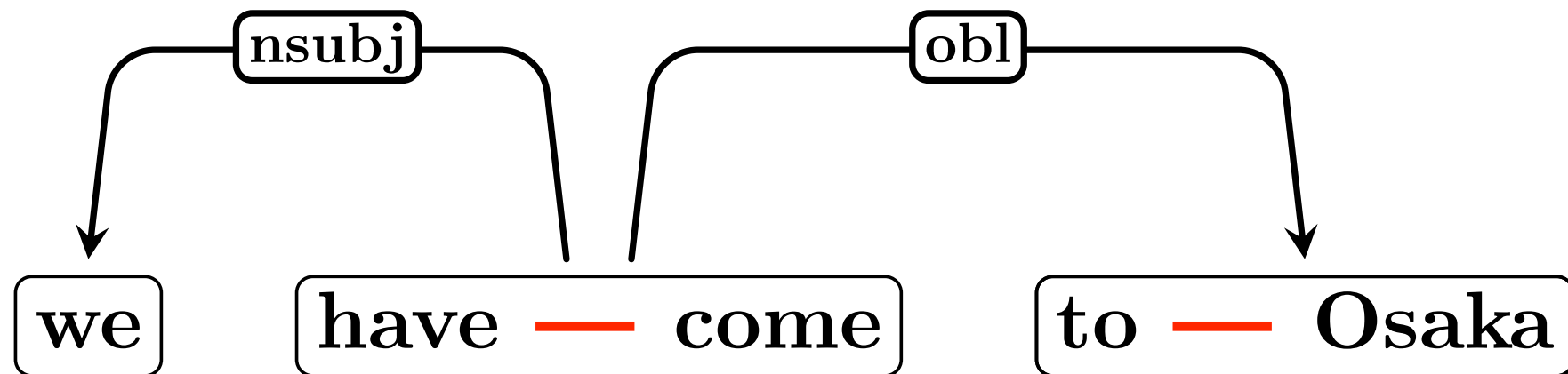
- A universal annotation scheme should have a classification of constructions, not strategies, as its universal foundational layer
- However, common recurrent strategies can and should be included as well
- The primacy of content words in UD is good, because function words such as copulas are strategies and not found in every language
- This means that the topology of dependency trees in UD are basically fine from a typological-universal point of view



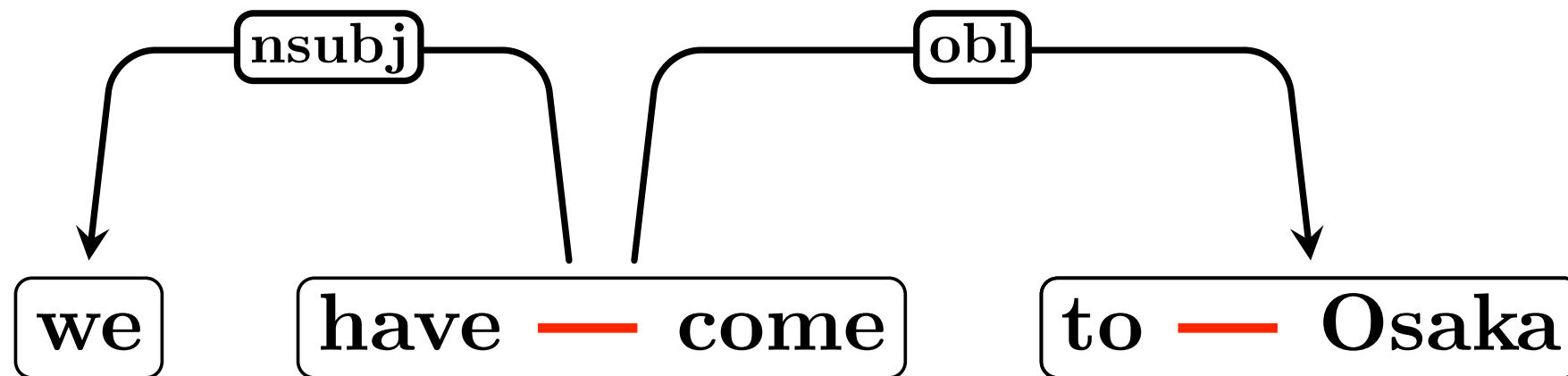


dependency

nucleus



dependency	nucleus
karaka	vibhakti



dependency	nucleus
karaka	vibhakti
kakariuke	bunsetsu

Conclusion – Part I

The primacy of content words

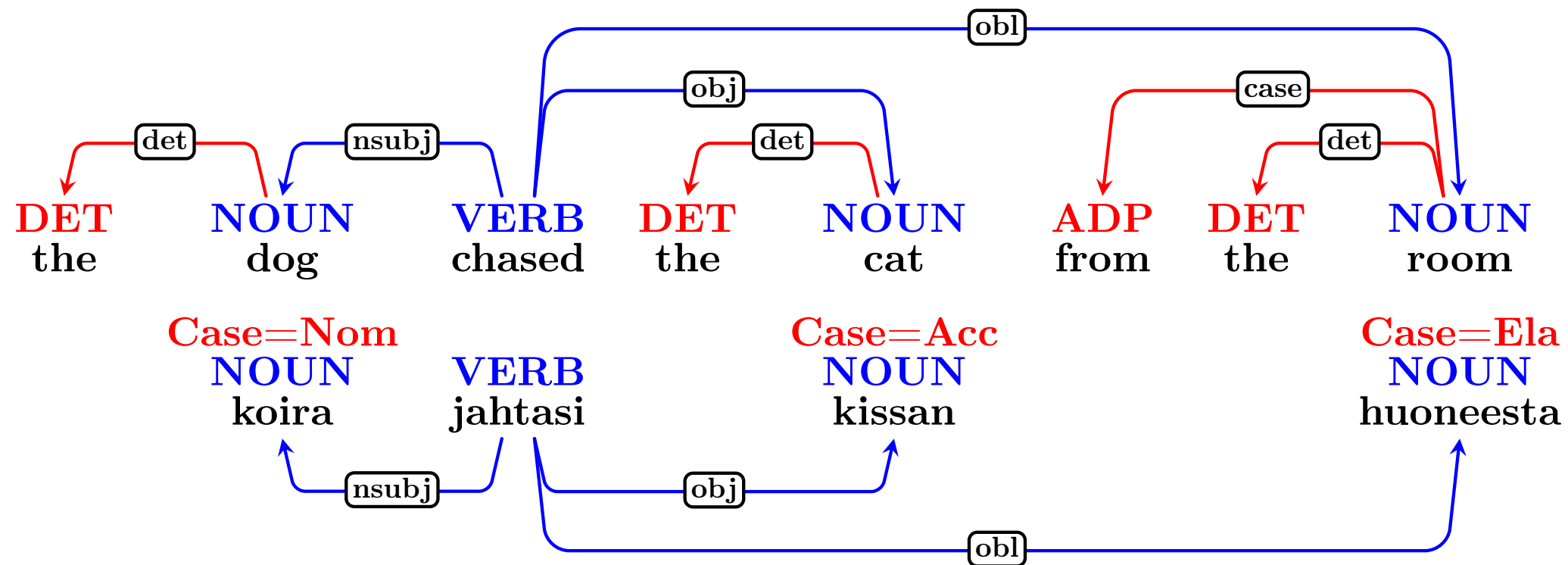
- Perhaps incompatible with some approaches to theoretical syntax
- Compatible with current approaches to linguistic typology
- Consistent with older traditions of dependency grammar

Natural Language Understanding



With contributions (and slides) by Siva Reddy et al.

From Syntax to Semantics



- Direct relations from predicates to arguments/modifiers
- Cross-linguistically consistent patterns
- Universal syntax-semantics interface?

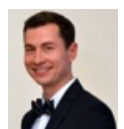
UDepLambda



UDepLambda

Goals

- From dependencies to logical forms
- Compositional, language-agnostic conversion
- Dependency tree dictates the semantics



UDepLambda

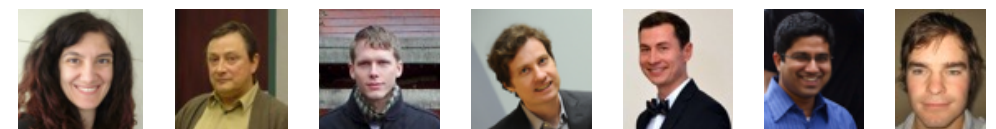
Goals

- From dependencies to logical forms
- Compositional, language-agnostic conversion
- Dependency tree dictates the semantics

Compositionality

The semantics of a **complex expression** is determined by the semantics of its **constituent expressions** and the **rules** used to combine them

- Complex expressions are dependency trees
- Constituent expressions are subtrees
- Rules are dependency labels



Why Dependencies?

Dependency Tree to Semantics



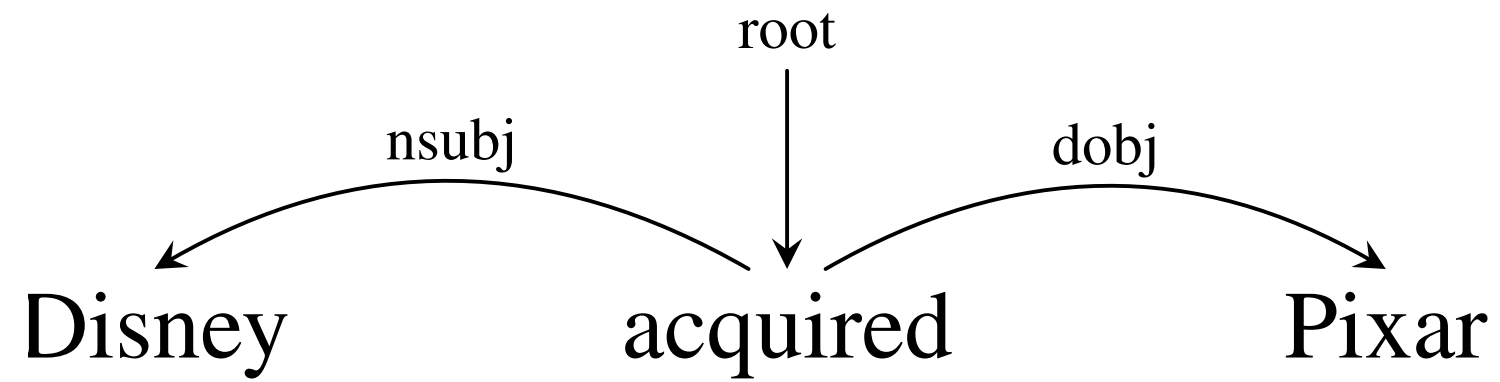
Dependencies **lack** a formal theory of semantics

CCG

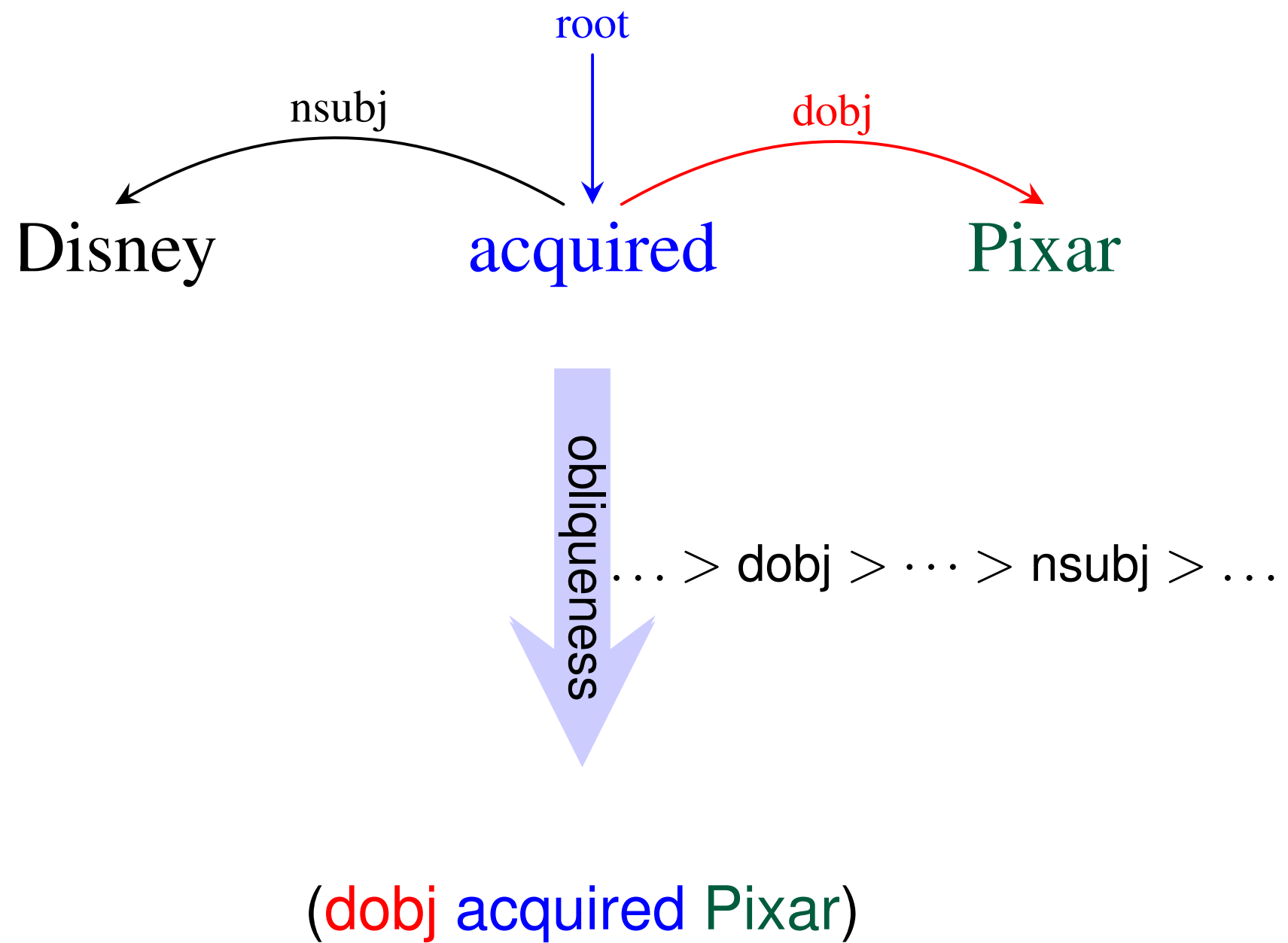
$$\begin{array}{c}
 \text{Disney} \quad \text{acquired} \quad \text{Pixar} \\
 \hline
 NP \quad S \backslash NP / NP \quad NP \\
 \text{Disney} \quad \lambda y \lambda x \lambda e. \text{acquired}(e) \quad \text{Pixar} \\
 \quad \quad \quad \wedge \text{arg}_1(e, x) \\
 \quad \quad \quad \wedge \text{arg}_2(e, y) \\
 \hline
 S \backslash NP \quad \rightarrow \\
 \quad \quad \quad \lambda x \lambda e. \text{acquired}(e) \\
 \quad \quad \quad \wedge \text{arg}_1(e, x) \wedge \text{arg}_2(e, \text{Pixar}) \\
 \hline
 S \quad \leftarrow \\
 \lambda e. \text{acquired}(e) \wedge \text{arg}_1(e, \text{Disney}) \wedge \text{arg}_2(e, \text{Pixar})
 \end{array}$$

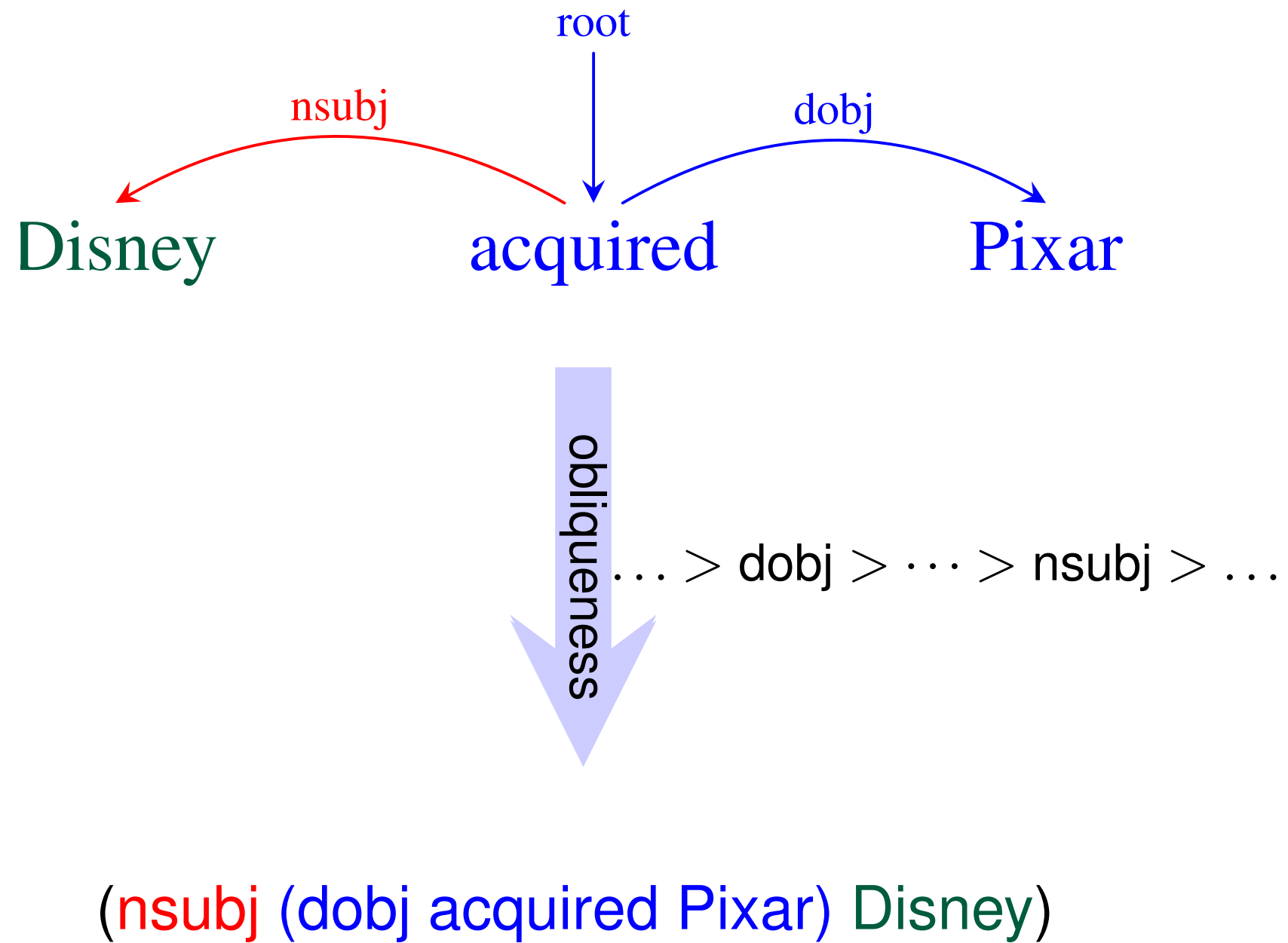
Typing and Combinator Rules allow
Synchronous Syntax-Semantics interface

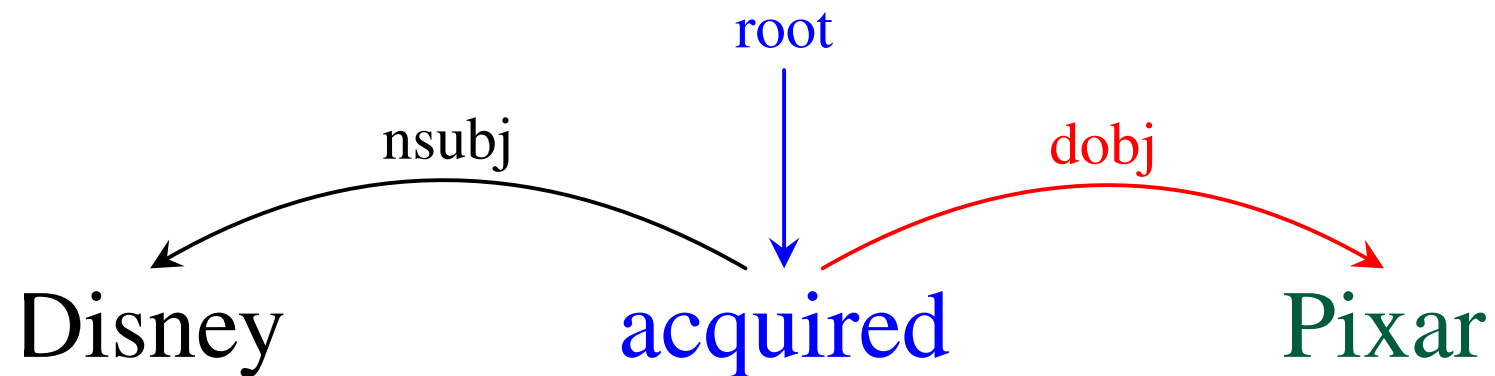
- Easy to annotate (with UD annotation scheme)
- Treebanks exist in many languages
- Robust, efficient, accurate parsers



Dependency labels drive the composition







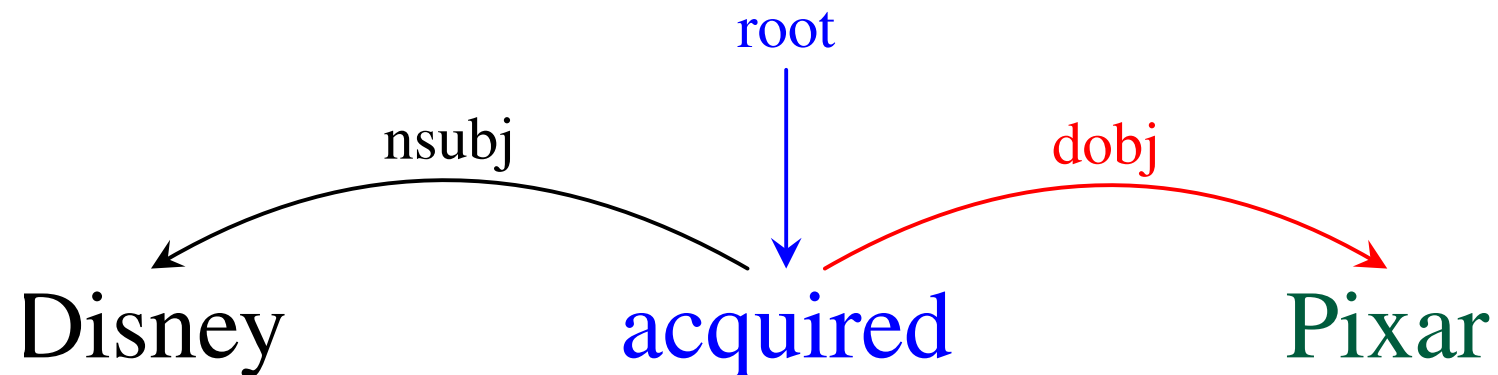
Lambda Expression for words

acquired $\Rightarrow \lambda \mathbf{x_a} x_e. \text{acquired}(x_e) \quad \Rightarrow \text{TYPE} = \mathbf{Ind} \times \mathbf{Event} \rightarrow \mathbf{Bool}$

Pixar $\Rightarrow \lambda x_a \mathbf{x_e}. \text{Pixar}(x_a) \quad \Rightarrow \text{TYPE} = \mathbf{Ind} \times \mathbf{Event} \rightarrow \mathbf{Bool}$

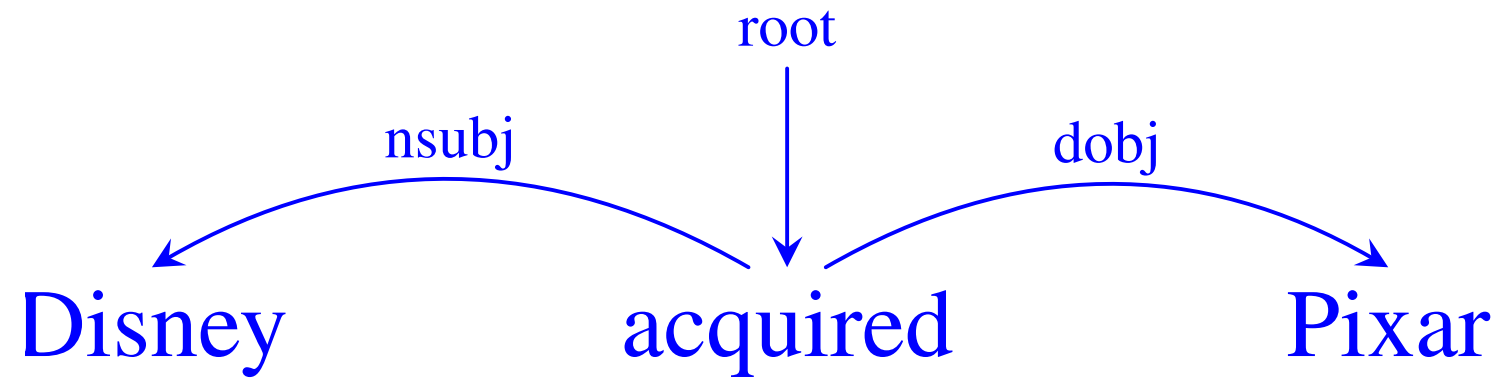
All constituents are of the same lambda expression type

$\text{TYPE}[\text{acquired}] = \text{TYPE}[\text{Pixar}] = \text{TYPE}[(\text{dobj acquired Pixar})]$



Lambda Expression for dependency labels

$$\text{dobj} \Rightarrow \lambda \mathbf{f} \lambda \mathbf{g} \lambda \mathbf{z} . \exists \mathbf{x} . \mathbf{f}(\mathbf{z}) \wedge \mathbf{g}(\mathbf{x}) \wedge \mathbf{arg}_2(\mathbf{z}_e, \mathbf{x}_a)$$



(nsubj (dobj acquired Pixar) Disney)

$$\lambda z. \exists xy. \text{acquired}(z_e) \wedge \text{Pixar}(y_a) \wedge \text{Disney}(x_a) \wedge \arg_1(z_e, x_a) \wedge \arg_2(z_e, y_a)$$

Comparison with CCG

CCG

Lexicalized semantics

Words drive composition

Argument and adjunct
distinction

Complex types are powerful

UDepLambda

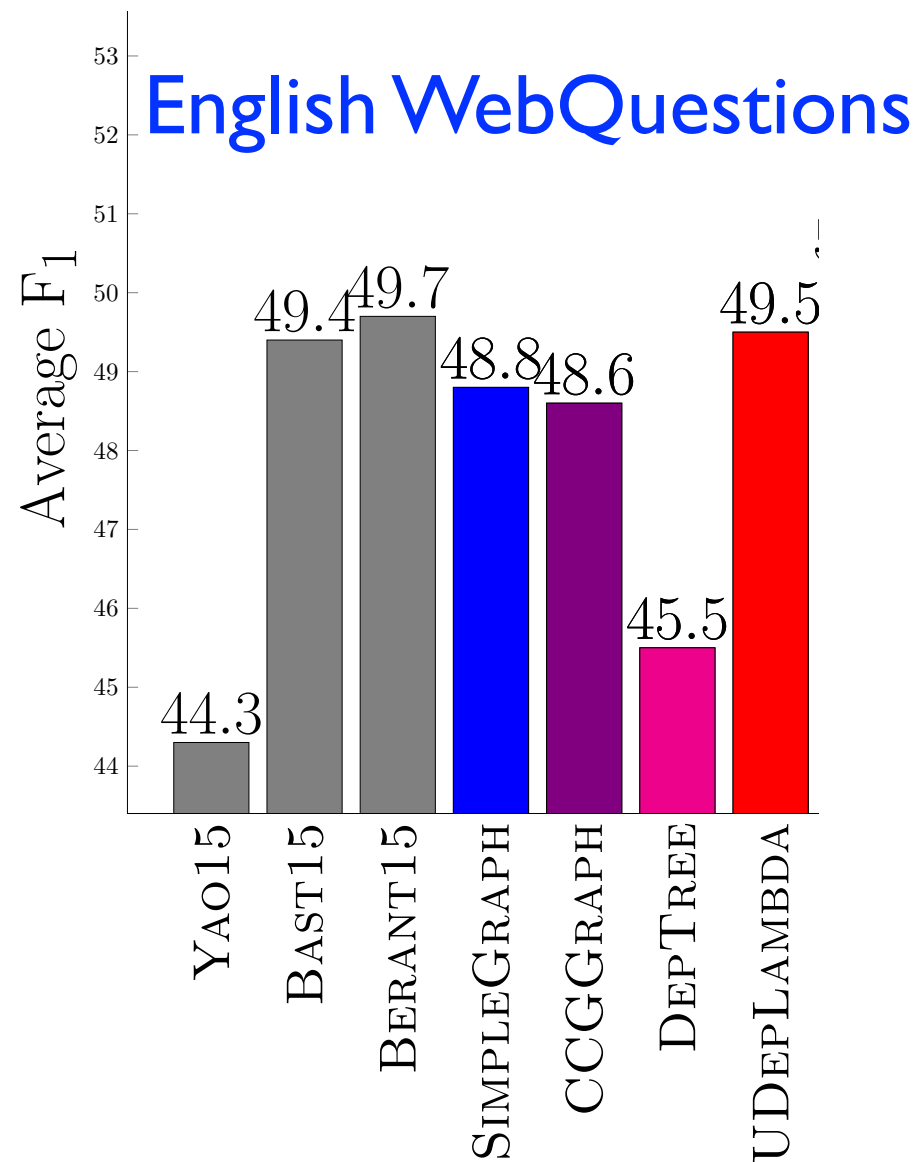
Simple lexical semantics

Dependencies drive
composition

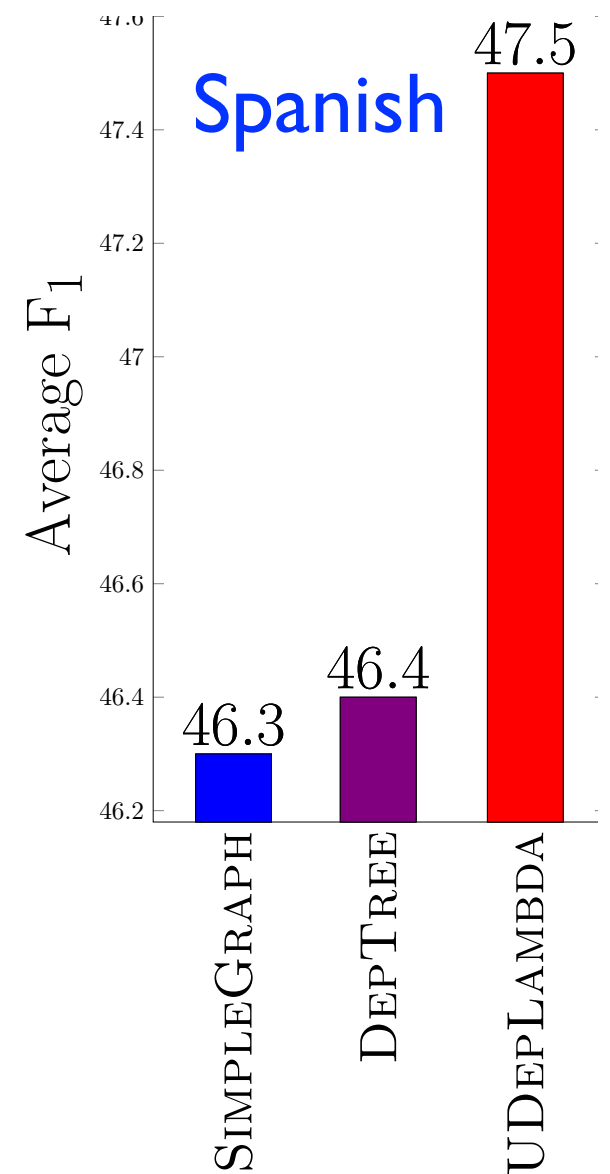
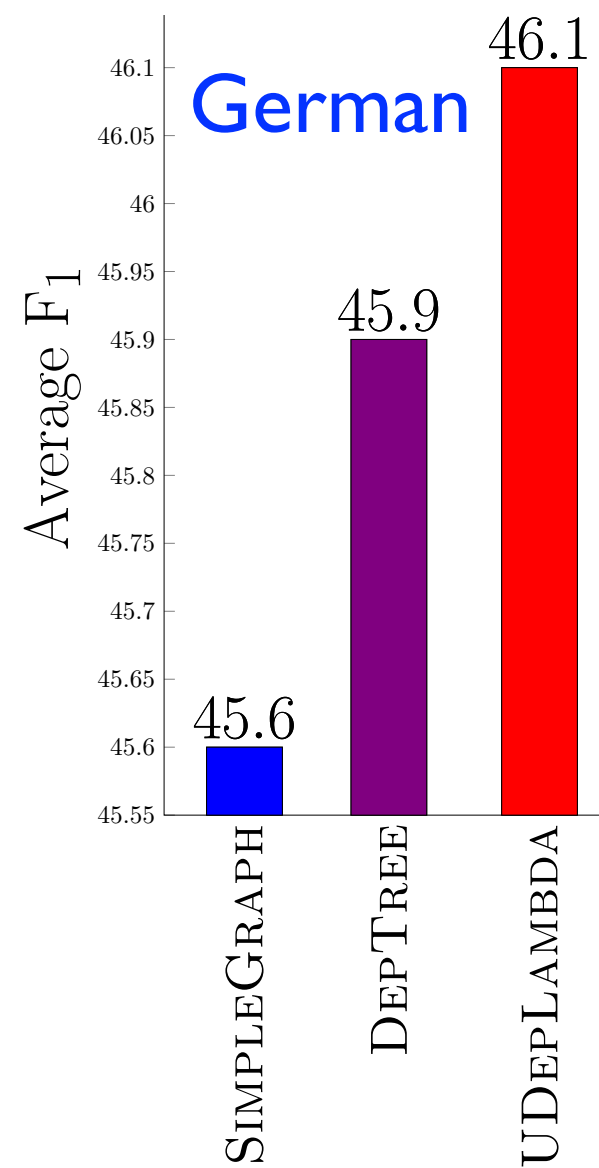
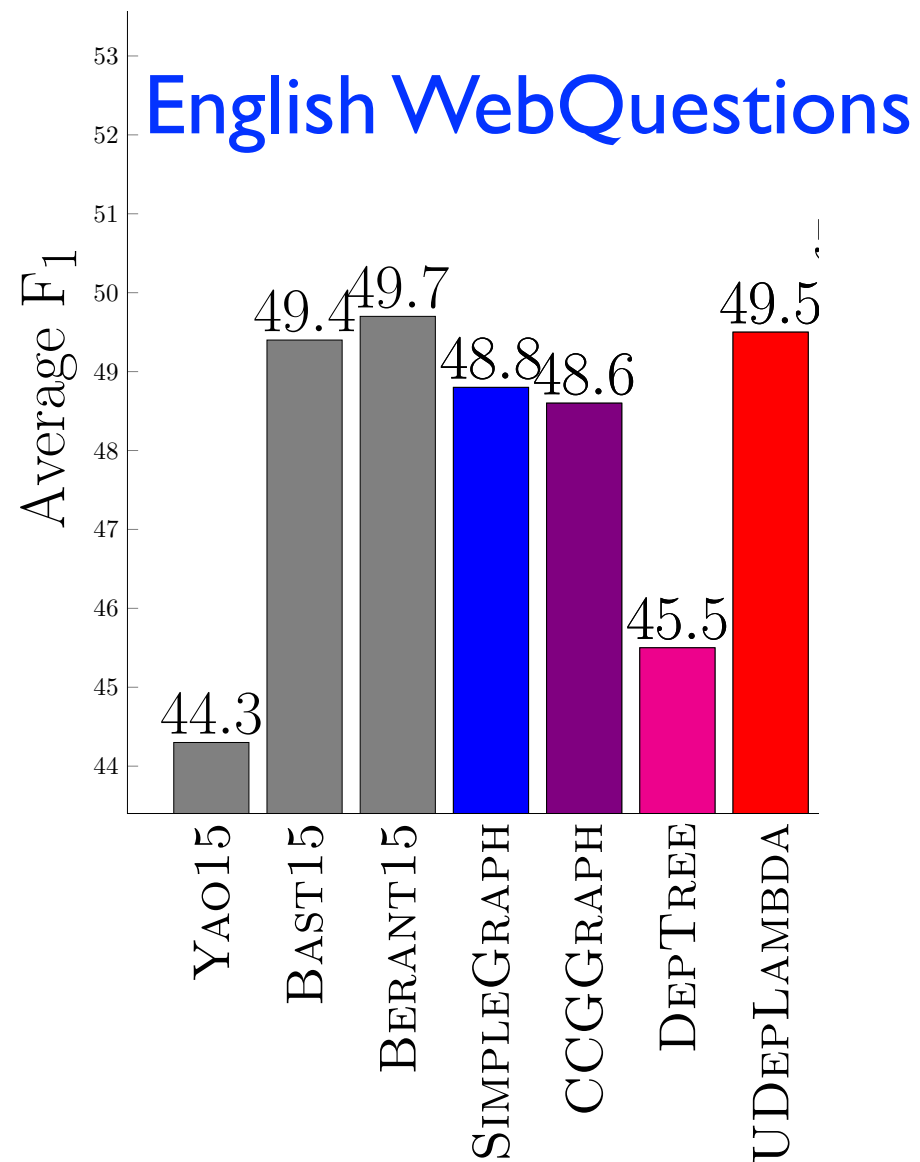
Every dependent is
an adjunct

Simplicity gives robustness

Freebase Semantic Parsing



Freebase Semantic Parsing

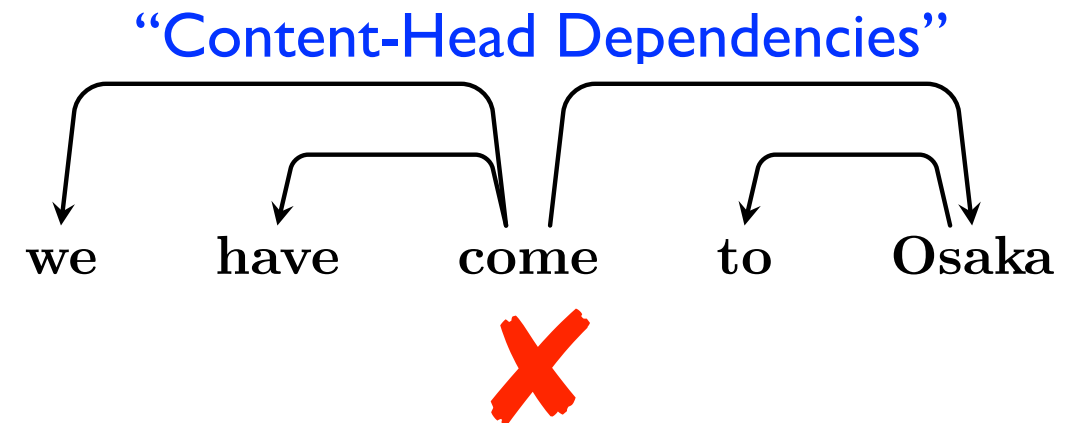
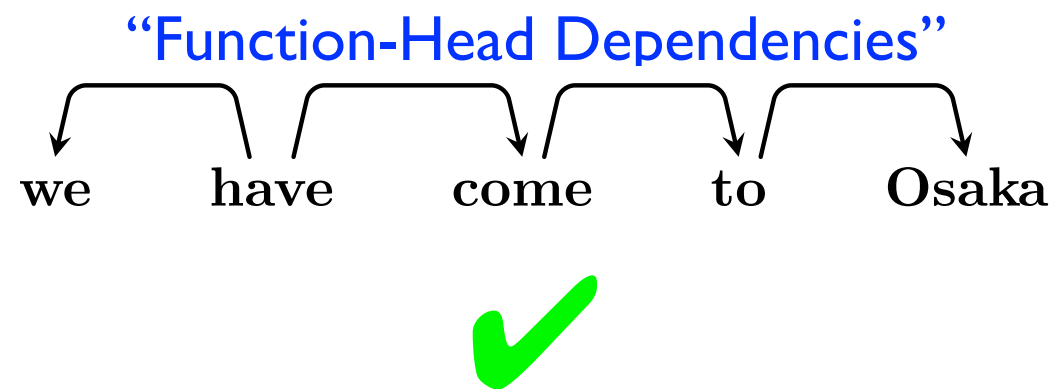


Conclusion – Part II

The primacy of content words

- Supports dependency-driven compositional semantics
- Generalizes across languages – universal semantic parsing?
- Useful for downstream natural language understanding tasks

Syntactic Parsing



Conventional wisdom

- Dependency trees with “function heads” are easier to parse
- Possibly related to dependency length minimisation
- Based on scattered empirical (and anecdotal) evidence

Paper

Language

the	cat	for	her	has	put	is	nice	to	go	and	that
det		case		aux		cop		mark		cc	

Paper	Language	<div><div>the</div><div>cat</div><div>for</div><div>her</div><div>has</div><div>put</div><div>is</div><div>nice</div><div>to</div><div>go</div><div>and</div><div>that</div></div>					
		det	case	aux	cop	mark	cc
Schwartz et al. (2012)	English						

Paper	Language	the cat det	for her case	has put aux	is nice cop	to go mark	and that cc
Schwartz et al. (2012)	English						
Silveira and Manning (2015)	English						

Paper	Language	<div>the cat</div> <div>det</div>	<div>for her</div> <div>case</div>	<div>has put</div> <div>aux</div>	<div>is nice</div> <div>cop</div>	<div>to go</div> <div>mark</div>	<div>and that</div> <div>cc</div>
		det	case	aux	cop	mark	cc
Schwartz et al. (2012)	English						
Silveira and Manning (2015)	English						
Attardi et al. (2015)	Italian						

Paper	Language	<div>the cat</div> <div>det</div>	<div>for her</div> <div>case</div>	<div>has put</div> <div>aux</div>	<div>is nice</div> <div>cop</div>	<div>to go</div> <div>mark</div>	<div>and that</div> <div>cc</div>
		det	case	aux	cop	mark	cc
Schwartz et al. (2012)	English						
Silveira and Manning (2015)	English						
Attardi et al. (2015)	Italian						
Rosa (2015)	Multi						
De Lhoneux and Nivre (2016)	Multi						

Paper	Language	<div> <div>det</div> <div>the</div> <div>cat</div> </div>	<div> <div>case</div> <div>for</div> <div>her</div> </div>	<div> <div>aux</div> <div>has</div> <div>put</div> </div>	<div> <div>cop</div> <div>is</div> <div>nice</div> </div>	<div> <div>mark</div> <div>to</div> <div>go</div> </div>	<div> <div>cc</div> <div>and</div> <div>that</div> </div>
		det	case	aux	cop	mark	cc
Schwartz et al. (2012)	English						
Silveira and Manning (2015)	English						
Attardi et al. (2015)	Italian						
Rosa (2015)	Multi						
De Lhoneux and Nivre (2016)	Multi						
Wisniewski and Lacroix (2017)	Multi						

Conclusion – Part III

The primacy of content words

- Not universally harmful for monolingual syntactic parsing
- Arguably superior for cross-lingual parsing
- Neural dependency parsers less sensitive to representations?

Universal Dependencies

Designed to fulfill multiple purposes/requirements

- Meaningful linguistic analysis within and across languages
- Syntactic parsing in monolingual and cross-lingual settings
- Useful information for downstream language understanding tasks

Gives priority to dependencies between content words

Thanks to all UD contributors!

Željko Agić, Lars Ahrenberg, Maria Jesus Aranzabe, Masayuki Asahara, Luma Ateyah, Mohammed Attia, Aitziber Atutxa, Elena Badmaeva, Miguel Ballesteros, Esha Banerjee, Sebastian Bank, John Bauer, Kepa Bengoetxea, Riyaz Ahmad Bhat, Eckhard Bick, Cristina Bosco, Gosse Bouma, Sam Bowman, Aljoscha Burchardt, Marie Candito, Gauthier Caron, Gülşen Cebiroğlu Eryiğit, Giuseppe G.A. Celano, Savas Cetin, Fabricio Chalub, Jinho Choi, Yongseok Cho, Silvie Cinková, Çağrı Çöltekin, Miriam Connor, Elizabeth Davidson, Eric de la Clergerie, Marie-Catherine de Marneffe, Valeria de Paiva, Arantza Diaz de Ilarraza, Kaja Dobrovoljc, Timothy Dozat, Kira Droganova, Puneet Dwivedi, Marhaba Eli, Ali Elkahky, Tomaž Erjavec, Richárd Farkas, Hector Fernandez Alcalde, Jennifer Foster, Cláudia Freitas, Katarína Gajdošová, Daniel Galbraith, Marcos Garcia, Filip Ginter, Iakes Goenaga, Koldo Gojenola, Memduh Gökırmak, Yoav Goldberg, Xavier Gómez Guinovart, Berta Gonzáles Saavedra, Matias Grioni, Normunds Grūzītis, Bruno Guillaume, Nizar Habash, Jan Hajič, Jan Hajič jr., Linh Hà Mỹ, Kim Harris, Dag Haug, Barbora Hladká, Jaroslava Hlaváčová, Petter Hohle, Radu Ion, Elena Irimia, Anders Johannsen, Fredrik Jørgensen, Hüner Kaşıkara, Hiroshi Kanayama, Jenna Kanerva, Jenna Kannerva, Tolga Kayadelen, Václava Kettnerová, Jesse Kirchner, Natalia Kotsyba, Simon Krek, Sookyoung Kwak, Veronika Laippala, Lorenzo Lambertino, Tatiana Lando, John Lee, Phươg Lê Hồng, Alessandro Lenci, Saran Lertpradit, Herman Leung, Cheuk Ying Li, Josie Li, Keying Li, Nikola Ljubešić, Olga Loginova, Olga Lyashevskaya, Teresa Lynn, Vivien Macketanz, Aibek Makazhanov, Michael Mandl, Christopher Manning, Ruli Manurung, Cătălina Măranduc, David Mareček, Katrin Marheinecke, Héctor Martínez Alonso, André Martins, Jan Mašek, Yuji Matsumoto, Ryan McDonald, Gustavo Mendonça, Anna Missilä, Verginica Mititelu, Yusuke Miyao, Simonetta Montemagni, Amir More, Laura Moreno Romero, Shunsuke Mori, Bohdan Moskalevskyi, Kadri Muischnek, Nina Mustafina, Kaili Müürisep, Pinkey Nainwani, Anna Nedoluzhko, Lươg Nguyễn Thị, Huyền Nguyễn Thị Minh, Vitaly Nikolaev, Rattima Nitisaroj, Joakim Nivre, Hanna Nurmi, Stina Ojala, Petya Osenova, Lilja Øvrelid, Elena Pascual, Marco Passarotti, Cenel-Augusto Perez, Guy Perrier, Slav Petrov, Jussi Piitulainen, Emily Pitler, Barbara Plank, Martin Popel, Lauma Pretkalniņa, Prokopis Prokopidis, Tiina Puolakainen, Sampo Pyysalo, Alexandre Rademaker, Loganathan Ramasamy, Vinit Ravishankar, Livy Real, Siva Reddy, Georg Rehm, Larissa Rinaldi, Laura Rituma, Rudolf Rosa, Mike Rosner, Davide Rovati, Benoit Sagot, Shadi Saleh, Manuela Sanguinetti, Baiba Saulīte, Yanin Sawanakunanon, Sebastian Schuster, Djamé Seddah, Wolfgang Seeker, Mojgan Seraji, Lena Shakurova, Mo Shen, Atsuko Shimada, Muh Shohibussirri, Dmitry Sichinava, Natalia Silveira, Maria Simi, Radu Simionescu, Katalin Simkó, Mária Šimková, Kiril Simov, Aaron Smith, Antonio Stella, Milan Straka, Jana Strnadová, Alane Suhr, Umut Sulubacak, Zolt Szántó, Dima Taji, Takaaki Tanaka, Anna Trukhina, Reut Tsarfaty, Francis Tyers, Sumire Uematsu, Zdeňka Urešová, Larraitz Uria, Hans Uszkoreit, Gertjan van Noord, Viktor Varga, Veronika Vincze, Jonathan North Washington, Zhuoran Yu, Zdeněk Žabokrtský, Amir Zeldes, Daniel Zeman, Hanzhi Zhu