



UPPSALA
UNIVERSITET



UNIVERSITÉ
DE GENÈVE

Universal Dependencies

A Framework for Cross-Linguistically
Consistent Grammatical Annotation

Joakim Nivre

Introduction

Introduction

Linguistic annotation is tremendously useful

- Computational linguistics use it for machine learning and evaluation
- Corpus linguistics use it for studying complex linguistic phenomena

Introduction

Linguistic annotation is tremendously useful

- Computational linguistics use it for machine learning and evaluation
- Corpus linguistics use it for studying complex linguistic phenomena

Linguistic annotation is available for many languages

- Multilingual evaluation to test generality of computational models
- Cross-lingual learning to support under-resourced languages
- Empirically grounded linguistic typology and comparative linguistics

Introduction

Linguistic annotation is tremendously useful

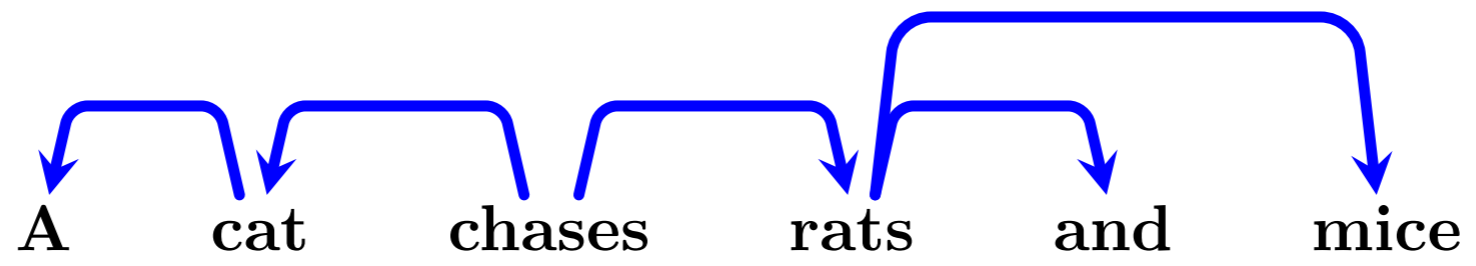
- Computational linguistics use it for machine learning and evaluation
- Corpus linguistics use it for studying complex linguistic phenomena

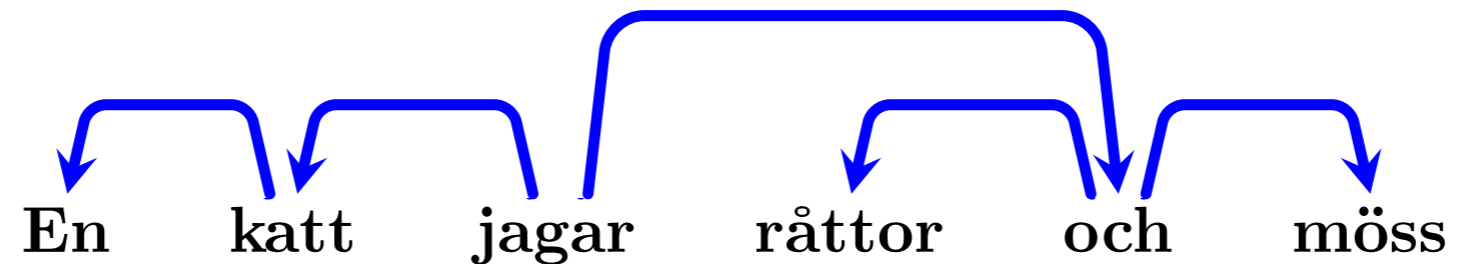
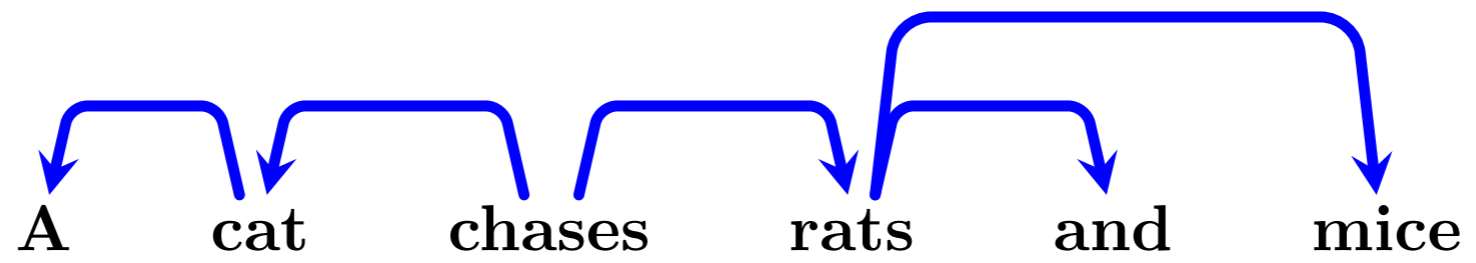
Linguistic annotation is available for many languages

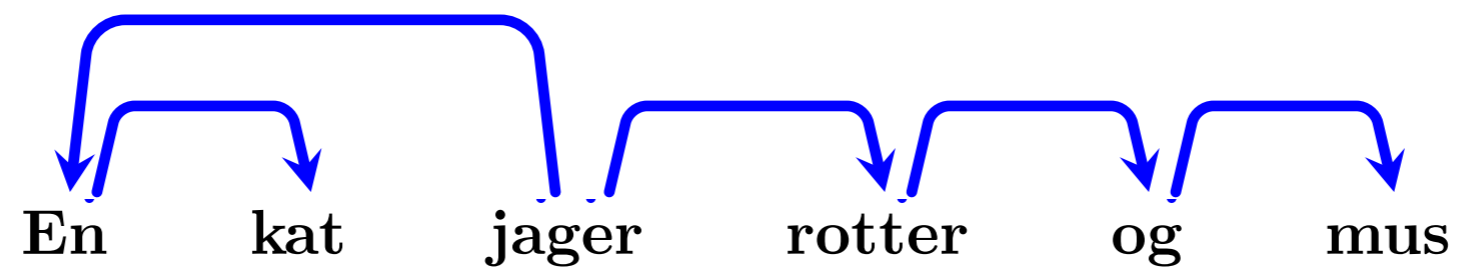
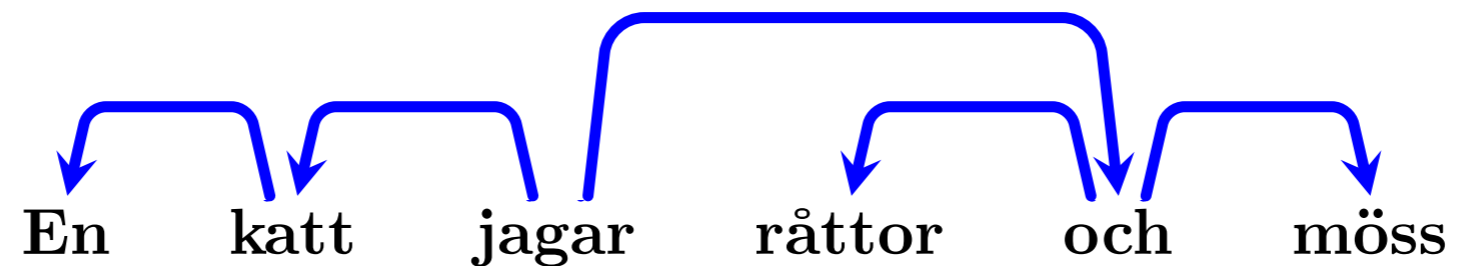
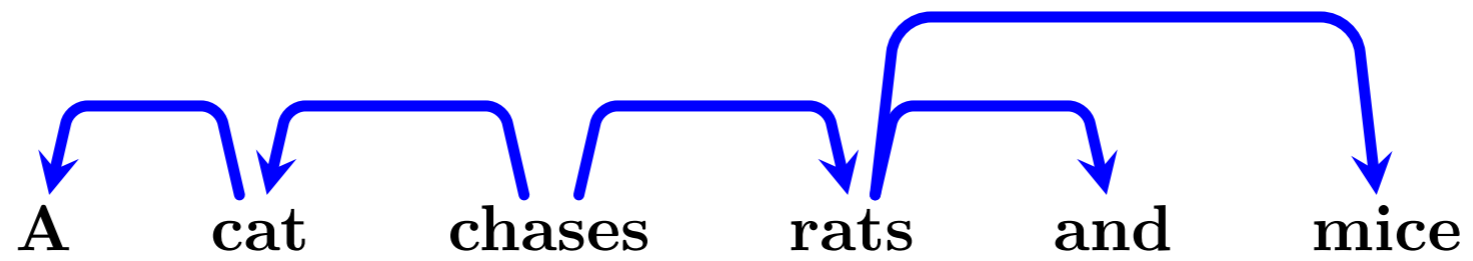
- Multilingual evaluation to test generality of computational models
- Cross-lingual learning to support under-resourced languages
- Empirically grounded linguistic typology and comparative linguistics

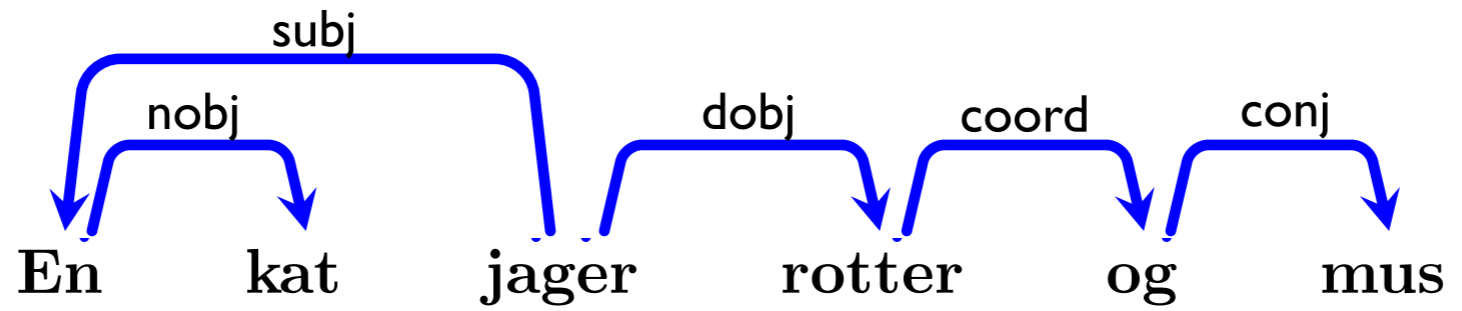
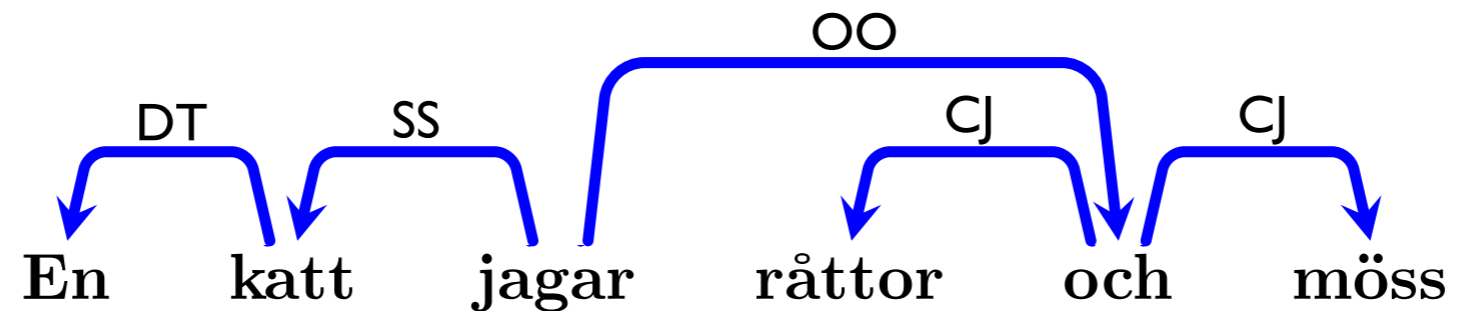
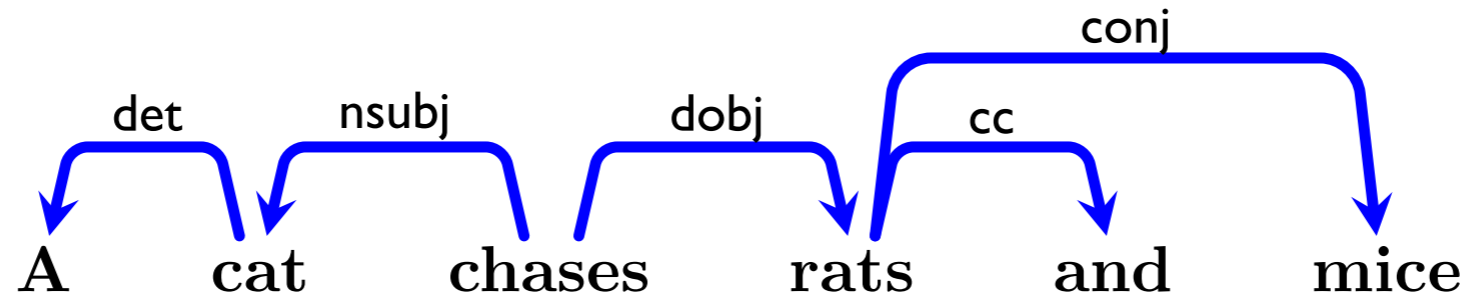
Linguistic annotation guidelines vary across languages

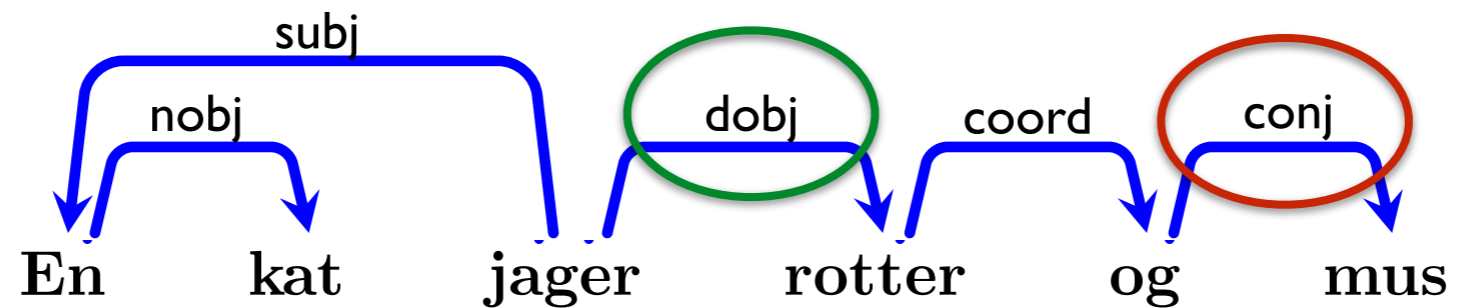
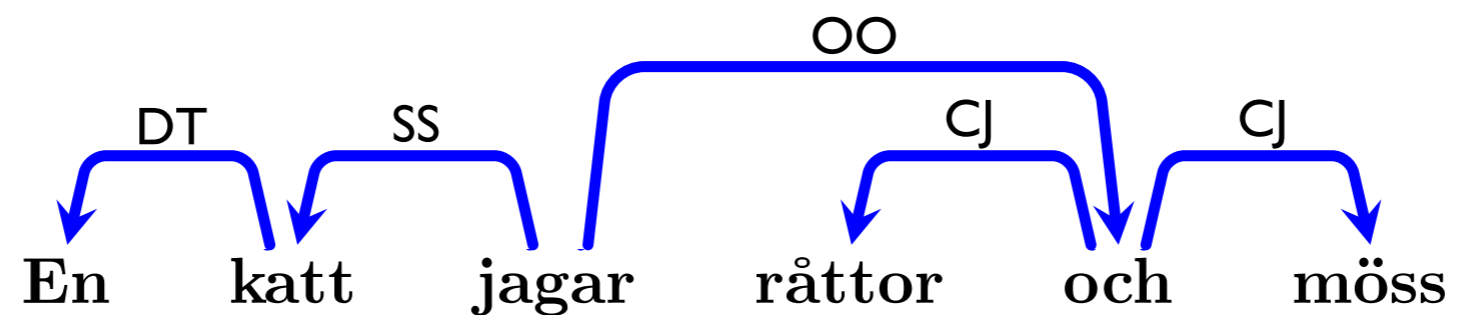
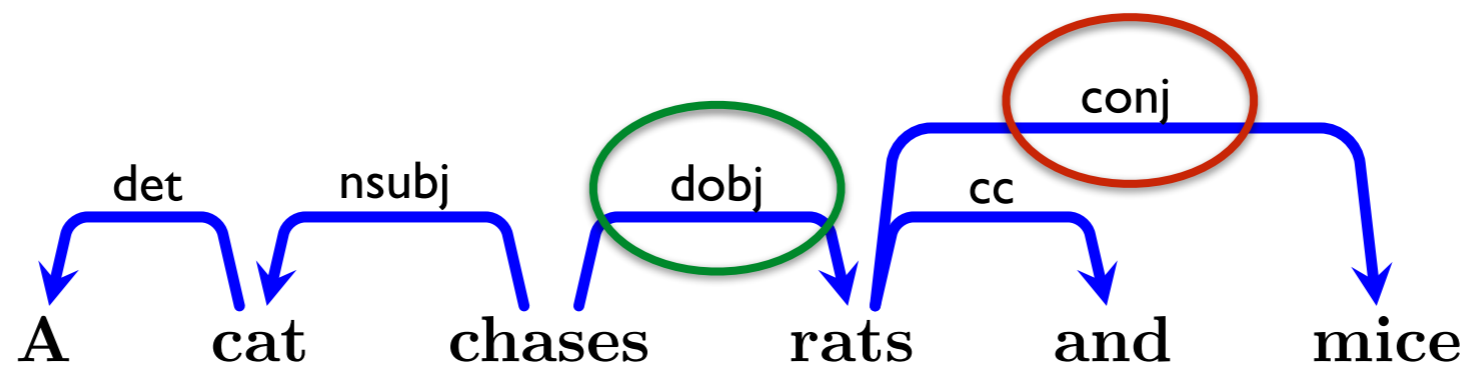












Why is this a problem?

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 learning

Why is this a problem?

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

Why is this a problem?

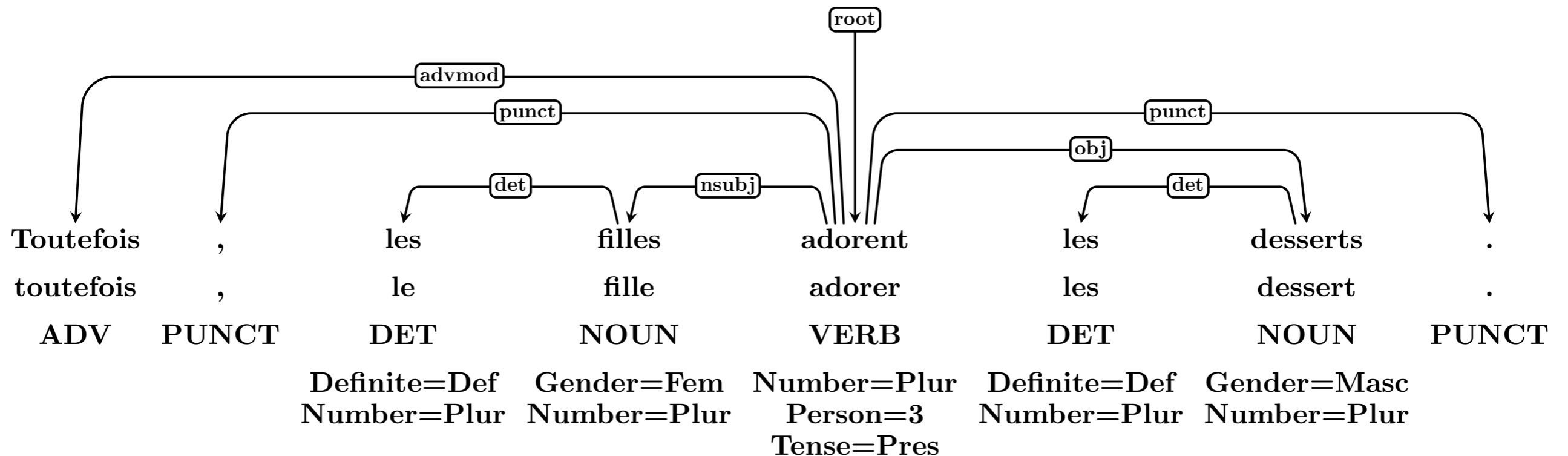
- Hard to compare empirical results across languages
- Hard to usefully do cross-lingual structure learning
- Hard to evaluate cross-lingual learning
- Hard to build and maintain multilingual systems
- Hard to make comparative linguistic studies
- Hard to validate linguistic typology

Why is this a problem?

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

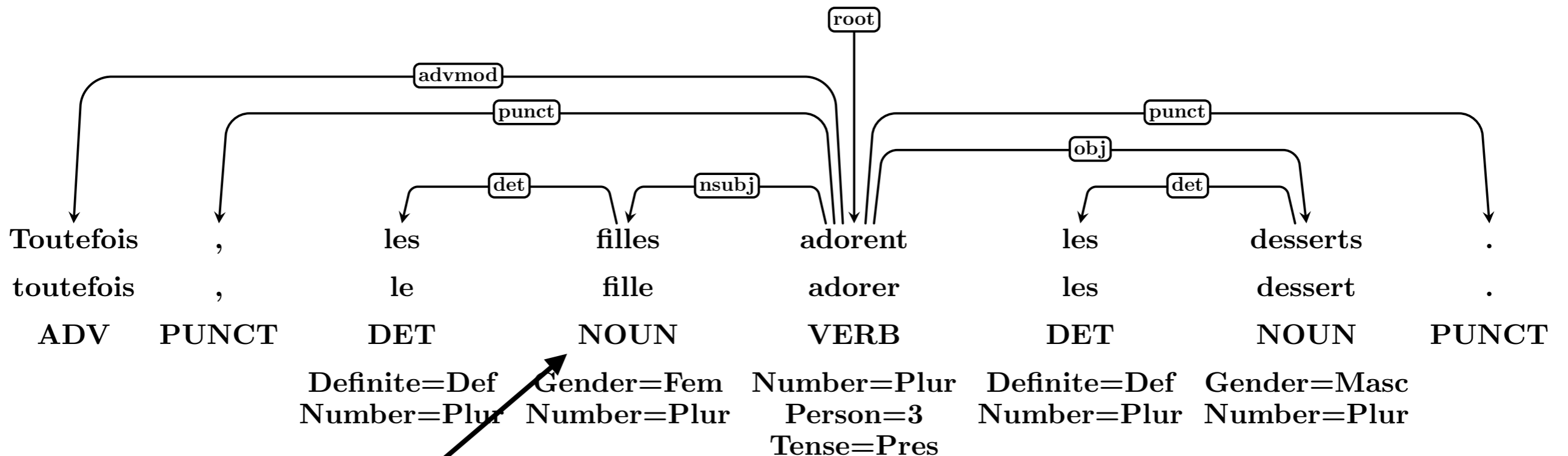
Universal Dependencies

<http://universaldependencies.org>



Universal Dependencies

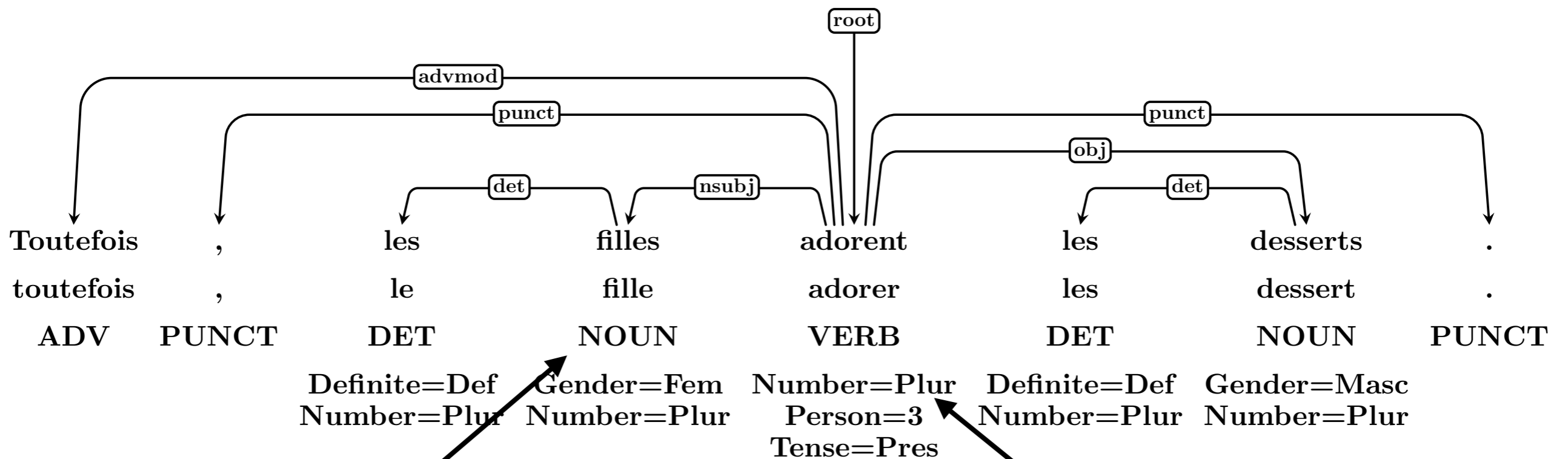
<http://universaldependencies.org>



Part-of-speech tags 

Universal Dependencies

<http://universaldependencies.org>



Part-of-speech tags 

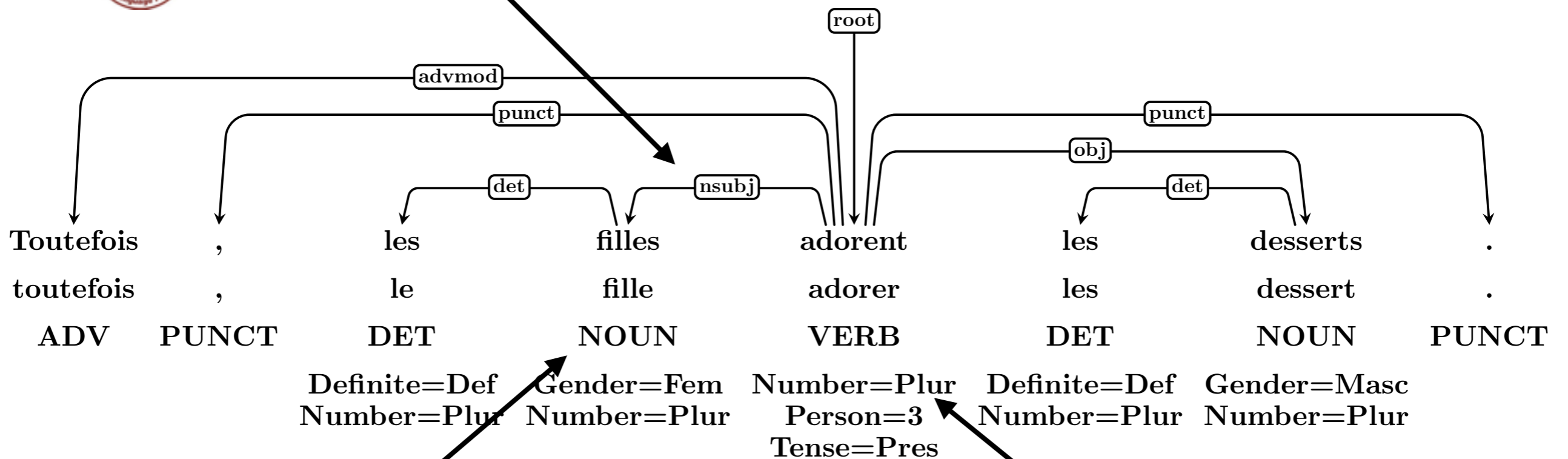
Morphological features 

Universal Dependencies

<http://universaldependencies.org>



Dependency relations



Part-of-speech tags

Morphological features

Universal Dependencies

<http://universaldependencies.org>

Universal Dependencies

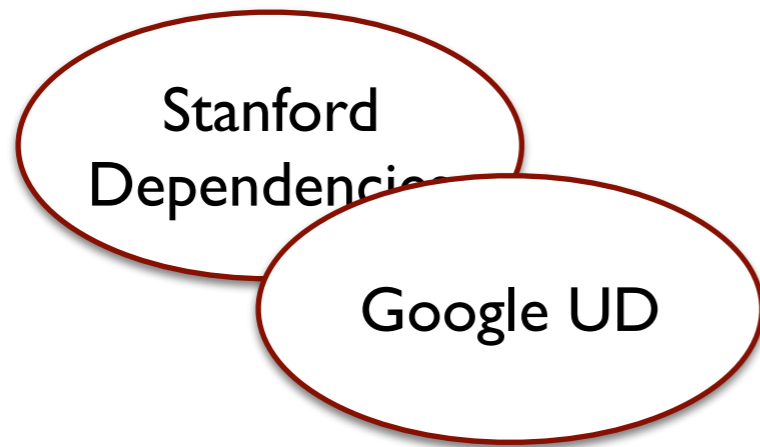
<http://universaldependencies.org>

The text "Stanford Dependencies" is enclosed in a red oval with a drop shadow.

Stanford
Dependencies

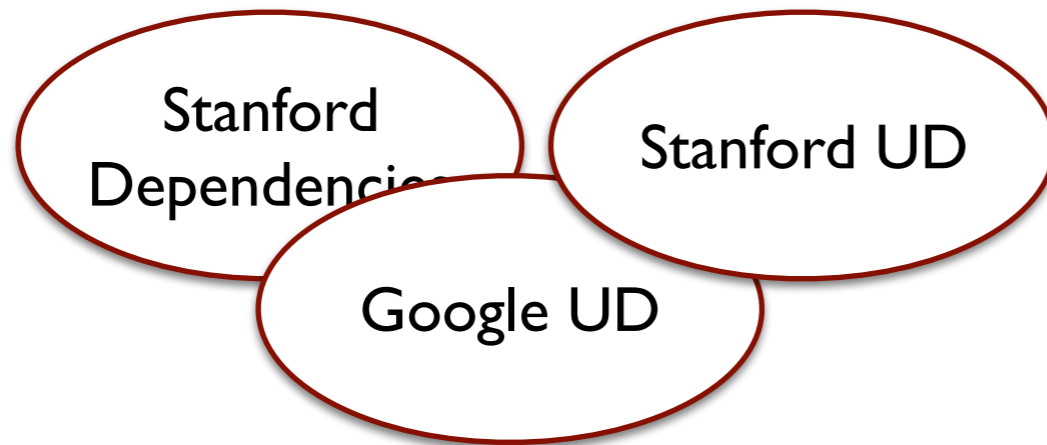
Universal Dependencies

<http://universaldependencies.org>



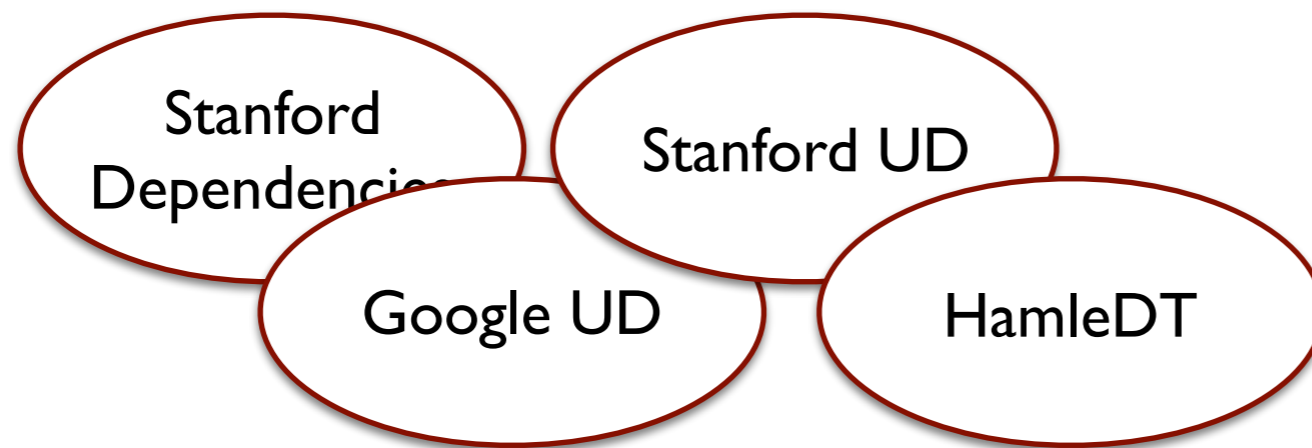
Universal Dependencies

<http://universaldependencies.org>



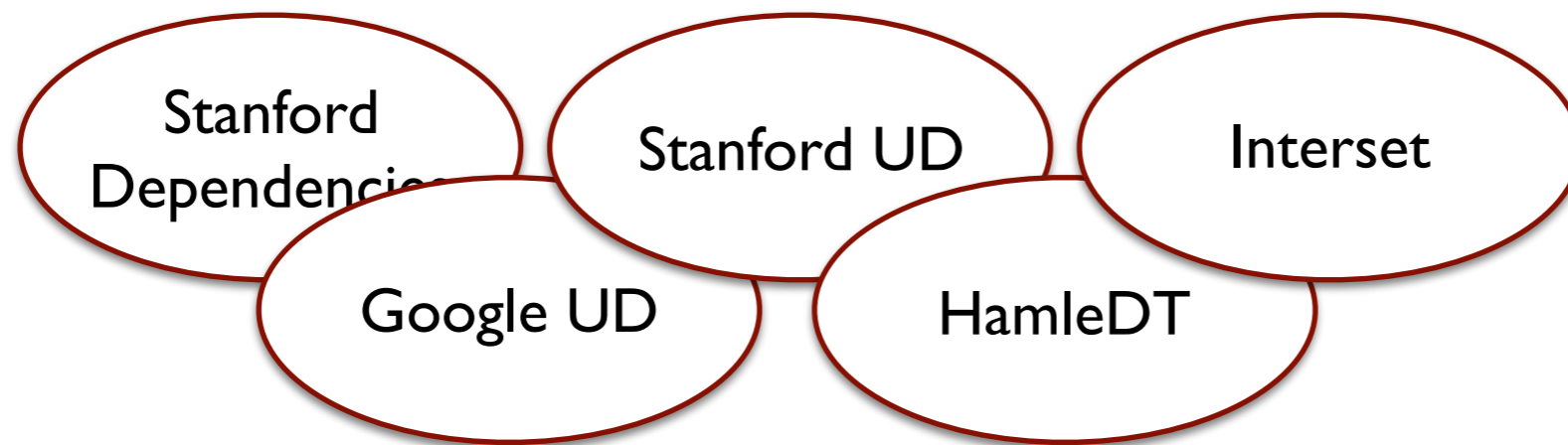
Universal Dependencies

<http://universaldependencies.org>



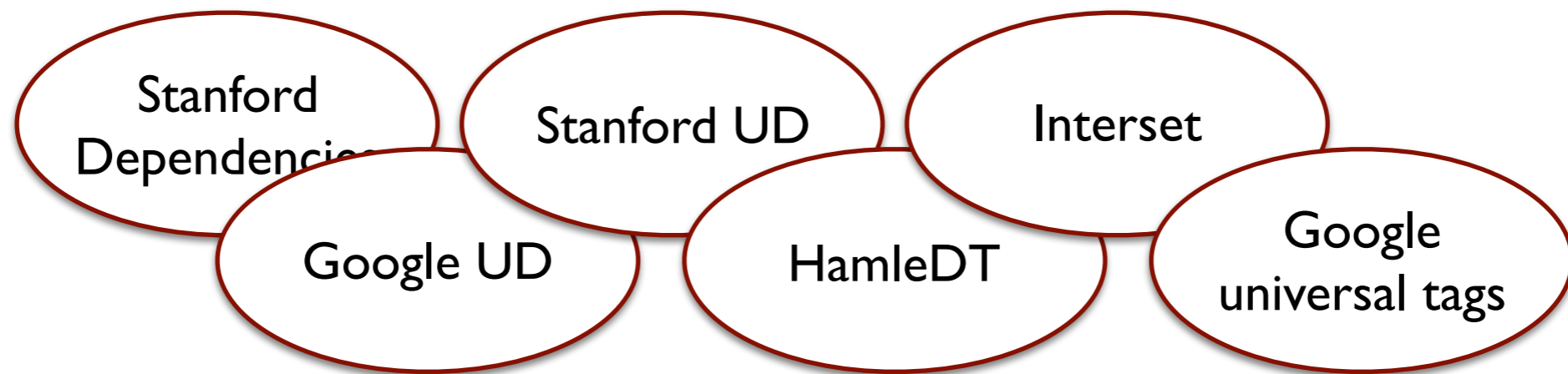
Universal Dependencies

<http://universaldependencies.org>



Universal Dependencies

<http://universaldependencies.org>



Universal Dependencies

<http://universaldependencies.org>




Universal Dependencies

Universal Dependencies

<http://universaldependencies.org>

Universal Dependencies

Milestones:

- 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	182K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Ancient Greek-PROIEL	198K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Arabic	217K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Arabic-NYUAD	629K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Basque	97K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Belarusian	6K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Bulgarian	140K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Catalan	472K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Chinese	111K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Coptic	3K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Croatian	183K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Czech	1,330K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Czech-CAC	482K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Czech-CLTT	26K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Danish	94K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Dutch	197K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Dutch-LassySmall	93K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		English	229K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		English-ESL	88K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		English-LinES	67K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		English-ParTUT	38K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Estonian	34K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Finnish	181K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Finnish-FTB	143K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		French	381K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		French-ParTUT	17K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		French-Sequola	58K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Galician	109K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Galician-TreeGal	14K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		German	277K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Gothic	45K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Greek	51K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Hebrew	106K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Hindi	316K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Hungarian	37K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Indonesian	110K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Irish	13K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Italian	195K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Italian-ParTUT	39K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Japanese	173K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Japanese-KTC	189K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Kazakh	<1K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Korean	63K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Korean-Sejong	89K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Latin	18K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Latin-ITTB	280K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Latin-PROIEL	159K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Latvian	44K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Lithuanian	40K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Norwegian-Bokmaal	280K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Norwegian-Nynorsk	276K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Old Church Slavonic	47K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Persian	135K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Polish	72K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Portuguese	201K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Portuguese-BR	268K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Romanian	202K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Russian	87K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Russian-SynTagRus	988K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Sanskrit	1K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Slovak	93K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Slovenian	126K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Slovenian-SST	19K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Spanish	411K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Spanish-AnCor	495K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Swedish	76K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Swedish-LinES	64K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Swedish Sign Language	<1K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Tamil	8K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Turkish	46K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Ukrainian	12K	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Urdu	123K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Uyghur	1K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️
▶		Vietnamese	31K	🗄️	-	🗄️	🗄️	🗄️	🗄️	🗄️	🗄️

March 1, 2017:

- 50 languages
- 70 treebanks
- 162 contributors
- 8000+ downloads

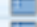
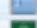


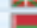










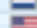

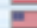































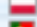






















UD Treebanks

▶		Ancient Greek	182K								
▶		Ancient Greek-PROIEL	198K								
▶		Arabic	217K								
▶		Arabic-NYUAD	629K								
▶		Basque	97K								
▶		Belarusian	6K								
▶		Bulgarian	140K								
▶		Catalan	472K								
▶		Chinese	111K								
▶		Coptic	3K								
▶		Croatian	183K								
▶		Czech	1,330K								
▶		Czech-CAC	482K								
▶		Czech-CLTT	26K								
▶		Danish	94K								
▶		Dutch	197K								
▶		Dutch-LassySmall	93K								
▶		English	229K								
▶		English-ESL	88K								
▶		English-LinES	67K								
▶		English-ParTUT	38K								
▶		Estonian	34K								
▶		Finnish	181K								
▶		Finnish-FTB	143K								
▶		French	381K								
▶		French-ParTUT	17K								
▶		French-Sequola	58K								
▶		Galician	199K								
▶		Galician-TreeGal	14K								
▶		German	277K								
▶		Gothic	45K								
▶		Greek	51K								
▶		Hebrew	106K								
▶		Hindi	316K								
▶		Hungarian	37K								
▶		Indonesian	110K								
▶		Irish	13K								
▶		Italian	195K								
▶		Italian-ParTUT	39K								
▶		Japanese	173K								
▶		Japanese-KTC	189K								
▶		Kazakh	<1K								
▶		Korean	63K								
▶		Korean-Sejong	89K								
▶		Latin	18K								
▶		Latin-ITTB	280K								
▶		Latin-PROIEL	159K								
▶		Latvian	44K								
▶		Lithuanian	40K								
▶		Norwegian-Bokmaal	280K								
▶		Norwegian-Nynorsk	276K								
▶		Old Church Slavonic	47K								
▶		Persian	135K								
▶		Polish	72K								
▶		Portuguese	201K								
▶		Portuguese-BR	268K								
▶		Romanian	202K								
▶		Russian	87K								
▶		Russian-SynTagRus	988K								
▶		Sanskrit	1K								
▶		Slovak	93K								
▶		Slovenian	126K								
▶		Slovenian-SST	19K								
▶		Spanish	411K								
▶		Spanish-AnCora	495K								
▶		Swedish	76K								
▶		Swedish-LinES	64K								
▶		Swedish Sign Language	<1K								
▶		Tamil	8K								
▶		Turkish	46K								
▶		Ukrainian	12K								
▶		Urdu	123K								
▶		Uyghur	1K								
▶		Vietnamese	31K								

March 1, 2017:

- 50 languages
- 70 treebanks
- 162 contributors
- 8000+ downloads

































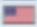
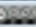





















































































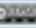














UD Treebanks

▶		Ancient Greek	182K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Ancient Greek-PROIEL	198K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Arabic	217K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Arabic-NYUAD	629K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Basque	97K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Belarusian	6K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Bulgarian	140K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Catalan	472K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Chinese	111K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Coptic	3K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Croatian	183K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Czech	1,330K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Czech-CAC	482K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Czech-CLTT	26K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Danish	94K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Dutch	197K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Dutch-LassySmall	93K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		English	229K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		English-ESL	88K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		English-LinES	67K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		English-ParTUT	38K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Estonian	34K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Finnish	181K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Finnish-FTB	143K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		French	381K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		French-ParTUT	17K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		French-Sequola	58K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Galician	109K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Galician-TreeCat	14K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		German	277K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Greek	51K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Hebrew	106K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Hindi	316K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Hungarian	37K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Indonesian	110K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Irish	13K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Italian	195K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Italian-ParTUT	39K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Japanese	173K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Japanese-KTC	189K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Kazakh	<1K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Korean	63K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Korean-Sejong	89K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Latin	18K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Latin-ITTB	280K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Latin-PROIEL	159K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Latvian	44K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Lithuanian	40K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Norwegian-Bokmaal	280K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Norwegian-Nynorsk	276K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Old Church Slavonic	47K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Persian	135K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Polish	72K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Portuguese	201K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Portuguese-BR	268K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Romanian	202K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Russian	87K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Russian-SynTagRus	988K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Sanskrit	1K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Slovak	93K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Slovenian	126K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Slovenian-SST	19K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Spanish	411K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Spanish-AnCor	495K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Swedish	76K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Swedish-LinES	64K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Swedish Sign Language	<1K	-	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Tamil	8K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Turkish	46K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Ukrainian	12K	🔒	🔒	🔒	🔒	🔒	🔒	🔒	🔒
▶		Urdu	123K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Uyghur	1K	-	-	🔒	🔒	🔒	🔒	🔒	🔒
▶		Vietnamese	31K	🔒	-	🔒	🔒	🔒	🔒	🔒	🔒

March 1, 2017:

- 50 languages
- 70 treebanks
- 162 contributors
- 8000+ downloads

UD Treebanks

▶		Ancient Greek	182K	00	0	0	✓		🌐
▶		Ancient Greek - PROIEL	198K	00	-	0	✓		🌐
▶		Arabic	217K	00	-	0	✓		🌐
▶		Arabic - NYUAD	629K	00	-	0	✓		🌐
▶		Basque	97K	00	0	0	✓		🌐
▶		Belarusian	6K	00	-	0	✓		🌐
▶		Bulgarian	140K	00	0	0	✓		🌐
▶		Catalan	472K	00	0	0	✓		🌐
▶		Chinese	111K	00	0	0	✓		W
▶		Coptic	3K	00	0	0	✓		🌐
▶		Croatian	183K	00	-	0	✓		🌐
▶		Czech	1,330K	00	0	0	✓		🌐
▶		Czech - CAC	482K	00	0	0	✓		🌐
▶		Czech - CLTT	26K	00	0	0	✓		🌐
▶		Danish	94K	00	0	0	✓		🌐
▶		Dutch	197K	00	-	0	✓		🌐
▶		Dutch - LassySmall	93K	00	-	0	✓		W
▶		English	229K	00	0	0	✓		🌐
▶		English - ESL	88K	0	0	0	✓		🌐
▶		English - LinES	67K	00	0	0	✓		🌐
▶		English - ParTUT	38K	00	0	0	✓		🌐
▶		Estonian	34K	00	-	0	✓		🌐
▶		Finnish	181K	00	0	0	✓		🌐
▶		Finnish - FTB	143K	00	-	0	✓		🌐
▶		French	381K	00	0	0	✓		🌐
▶		French - ParTUT	17K	00	0	0	✓		🌐
▶		French - Sequoia	58K	00	-	0	✓		🌐
▶		Galician	109K	0	0	0	✓		🌐
▶		Galician - TreeGal	14K	00	0	0	✓		🌐
▶		German	277K	00	-	0	✓		🌐
▶		Gothic	45K	00	-	0	✓		🌐
▶		Greek	51K	00	0	0	✓		🌐
▶		Hebrew	106K	0	-	0	✓		🌐
▶		Hindi	316K	00	-	0	✓		🌐
▶		Hungarian	37K	00	0	0	✓		🌐
▶		Indonesian	110K	0	-	0	✓		🌐
▶		Irish	13K	00	0	0	✓		🌐
▶		Italian	195K	00	0	0	✓		🌐
▶		Italian - ParTUT	39K	00	0	0	✓		🌐
▶		Japanese	175K	00	0	0	✓		🌐
▶		Japanese - KTC	189K	0	0	0	✓		🌐
▶		Kazakh	<1K	00	0	0	✓		W
▶		Korean	63K	0	0	0	✓		🌐
▶		Korean - Sejong	89K	0	-	0	?		🌐
▶		Latin	18K	00	0	0	✓		🌐
▶		Latin - ITTB	280K	00	-	0	✓		🌐
▶		Latin - PROIEL	159K	00	-	0	✓		🌐
▶		Latvian	44K	00	-	0	✓		🌐
▶		Lithuanian	40K	00	-	0	?		🌐
▶		Norwegian - Bokmaal	280K	00	0	0	✓		🌐
▶		Norwegian - Nynorsk	276K	00	0	0	✓		🌐
▶		Old Church Slavonic	47K	00	-	0	✓		🌐
▶		Persian	135K	00	0	0	✓		🌐
▶		Polish	72K	00	-	0	✓		🌐
▶		Portuguese	201K	00	0	0	✓		🌐
▶		Portuguese - BR	268K	0	-	0	✓		🌐
▶		Romanian	202K	00	0	0	✓		W
▶		Russian	87K	00	0	0	✓		W
▶		Russian - SynTagRus	988K	00	0	0	✓		🌐
▶		Sanskrit	1K	00	-	0	✓		🌐
▶		Slovak	93K	00	-	0	✓		🌐
▶		Slovenian	126K	00	0	0	✓		🌐
▶		Slovenian - SST	19K	00	0	0	✓		🌐
▶		Spanish	411K	00	0	0	✓		🌐
▶		Spanish - AnCor	495K	00	0	0	✓		🌐
▶		Swedish	76K	00	0	0	✓		🌐
▶		Swedish - LinES	64K	00	0	0	✓		🌐
▶		Swedish Sign Language	<1K	-	-	0	✓		🌐
▶		Tamil	8K	00	-	0	✓		🌐
▶		Turkish	46K	00	0	0	✓		🌐
▶		Ukrainian	12K	00	0	0	✓		🌐
▶		Urdu	123K	00	-	0	✓		🌐
▶		Uyghur	1K	-	-	0	✓		🌐
▶		Vietnamese	31K	00	-	0	✓		🌐

March 1, 2017:

- 50 languages
- 70 treebanks
- 162 contributors
- 8000+ downloads

UD Treebanks

Language	Count	Treebanks	Contributors	Downloads	Other
Ancient Greek	182K	00	0	0	0
Ancient Greek-PROIEL	198K	00	-	0	0
Arabic	217K	00	-	0	0
Arabic-NYUAD	629K	00	-	0	0
Basque	97K	00	0	0	0
Belarusian	6K	00	-	0	0
Bulgarian	140K	00	0	0	0
Catalan	472K	00	0	0	0
Chinese	111K	00	0	0	0
Coptic	3K	00	0	0	0
Croatian	183K	00	-	0	0
Czech	1,330K	00	0	0	0
Czech-CAC	482K	00	0	0	0
Czech-CLTT	26K	00	0	0	0
Danish	94K	00	0	0	0
Dutch	197K	00	-	0	0
Dutch-LassySmall	93K	00	-	0	0
English	229K	000	0	0	0
English-ESL	88K	0	0	0	0
English-LinES	67K	00	0	0	0
English-ParTUT	38K	00	0	0	0
Estonian	34K	00	-	0	0
Finnish	181K	000	0	0	0
Finnish-FTB	143K	00	-	0	0
French	381K	00	0	0	0
French-ParTUT	17K	00	0	0	0
French-Sequola	58K	00	-	0	0
Galician	109K	0	0	0	0
Galician-TreeGal	14K	00	0	0	0
German	277K	00	-	0	0
Gothic	45K	00	-	0	0
Greek	51K	00	0	0	0
Hebrew	106K	0	-	0	0
Hindi	316K	00	-	0	0
Hungarian	37K	00	-	0	0
Indonesian	110K	0	-	0	0
Irish	13K	00	0	0	0
Italian	195K	00	0	0	0
Italian-ParTUT	39K	00	0	0	0
Japanese	173K	00	0	0	0
Japanese-KTC	189K	0	0	0	0
Kazakh	<1K	00	0	0	0
Korean	63K	0	0	0	0
Korean-Sejong	89K	0	-	0	0
Latin	18K	00	0	0	0
Latin-ITTB	280K	00	-	0	0
Latin-PROIEL	159K	00	-	0	0
Latvian	44K	00	-	0	0
Lithuanian	40K	00	-	0	0
Norwegian-Bokmaal	280K	00	0	0	0
Norwegian-Nynorsk	276K	00	0	0	0
Old Church Slavonic	47K	00	-	0	0
Persian	135K	00	0	0	0
Polish	72K	00	-	0	0
Portuguese	201K	00	0	0	0
Portuguese-BR	268K	0	-	0	0
Romanian	202K	00	0	0	0
Russian	87K	00	0	0	0
Russian-SynTagRus	988K	00	0	0	0
Sanskrit	1K	00	-	0	0
Slovak	93K	00	-	0	0
Slovenian	126K	00	0	0	0
Slovenian-SST	19K	00	0	0	0
Spanish	411K	00	0	0	0
Spanish-AnCor	495K	00	0	0	0
Swedish	76K	00	0	0	0
Swedish-LinES	64K	00	0	0	0
Swedish Sign Language	<1K	-	-	0	0
Tamil	8K	00	-	0	0
Turkish	46K	00	0	0	0
Ukrainian	12K	00	0	0	0
Urdu	123K	00	-	0	0
Uyghur	1K	-	-	0	0
Vietnamese	31K	00	-	0	0

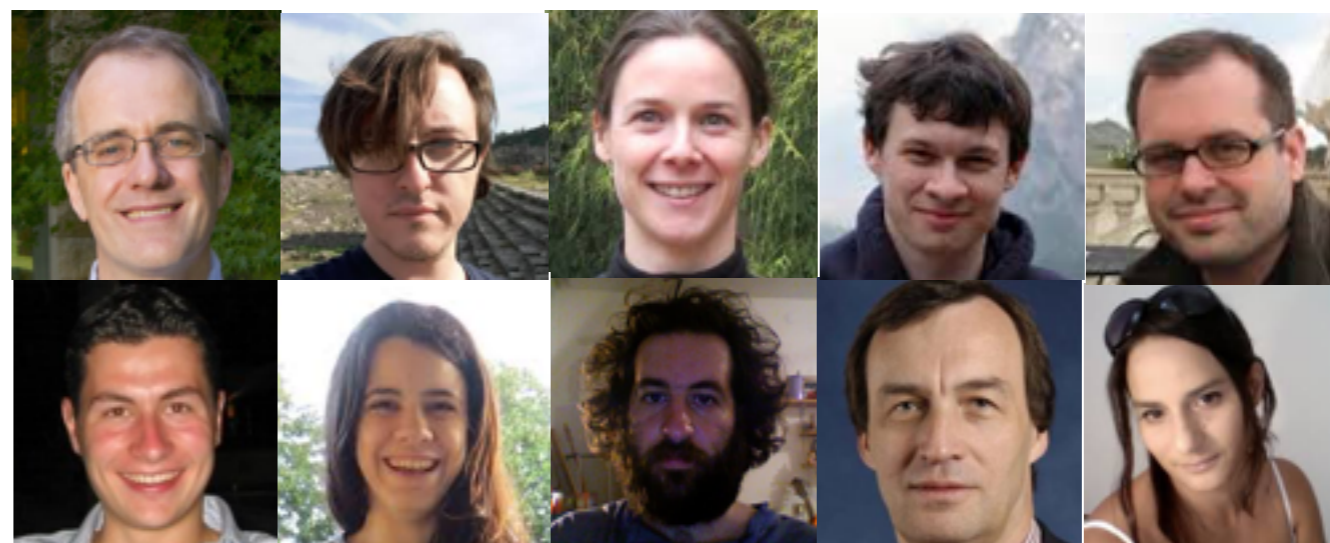
March 1, 2017:

- 50 languages
- 70 treebanks
- 162 contributors
- 8000+ downloads

Chief Cat Herder



Release and Documentation Task Force



Universal Guidelines Group

Goals and Requirements

Goals and Requirements

Cross-linguistically consistent grammatical annotation

Goals and Requirements

Cross-linguistically consistent grammatical annotation

Support multilingual research in NLP and linguistics

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

Goals and Requirements

Cross-linguistically consistent grammatical annotation

Support multilingual research in NLP and linguistics

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

Build on common usage and existing de facto standards

Goals and Requirements

Cross-linguistically consistent grammatical annotation

Support multilingual research in NLP and linguistics

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

Build on common usage and existing de facto standards

Complement – not replace – language-specific schemes

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

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

Recoverability

- Transparent mapping from input text to word segmentation

Word Segmentation

What is a word?

- Single part-of-speech tag
- Real syntactic relation

Two-level segmentation

- Represent orthographic tokens in addition to syntactic words

Word Segmentation

What is a word?

- Single part-of-speech tag
- Real syntactic relation

Two-level segmentation

- Represent orthographic tokens in addition to syntactic words

Text

Words

Word Segmentation

What is a word?

- Single part-of-speech tag
- Real syntactic relation

Two-level segmentation

- Represent orthographic tokens in addition to syntactic words

Text	Words
del	di il

Word Segmentation

What is a word?

- Single part-of-speech tag
- Real syntactic relation

Two-level segmentation

- Represent orthographic tokens in addition to syntactic words

Text	Words
del	di il
dámelo	da me lo

Word Segmentation

What is a word?

- Single part-of-speech tag
- Real syntactic relation

Two-level segmentation

- Represent orthographic tokens in addition to syntactic words

Text	Words
del	di il
dámelo	da me lo
ושמהשמש	ו ש מ ה ש מ ש

Word Segmentation

What is a word?

- Single part-of-speech tag
- Real syntactic relation

Two-level segmentation

- Represent orthographic tokens in addition to syntactic words

Text	Words
del	di il
dámelo	da me lo
ושמהשמש	ו ש מ ה ש מ ש
大阪国際会議場	大阪 国際 会議場

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

N

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

tiens

chien

NOUN

.

.

PUNCT

- Lemma representation 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 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 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

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

Le
le
DET
Definite=Def
Gender=Masc
Number=Sing

N
Gender
Number

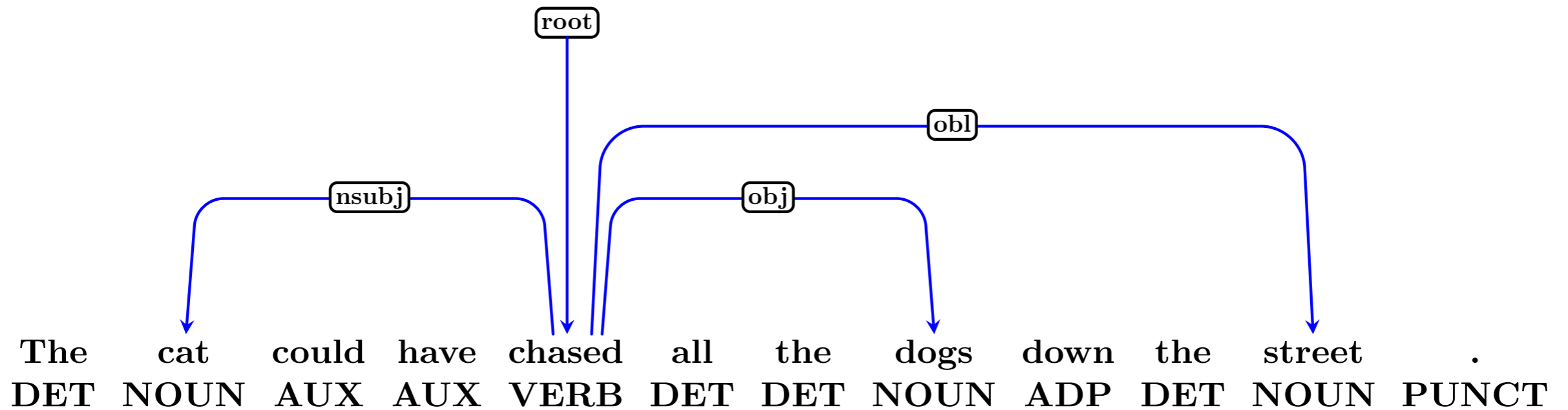
niens .
hien .
OUN **PUNCT**
er=Masc
er=Plur

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

Syntax

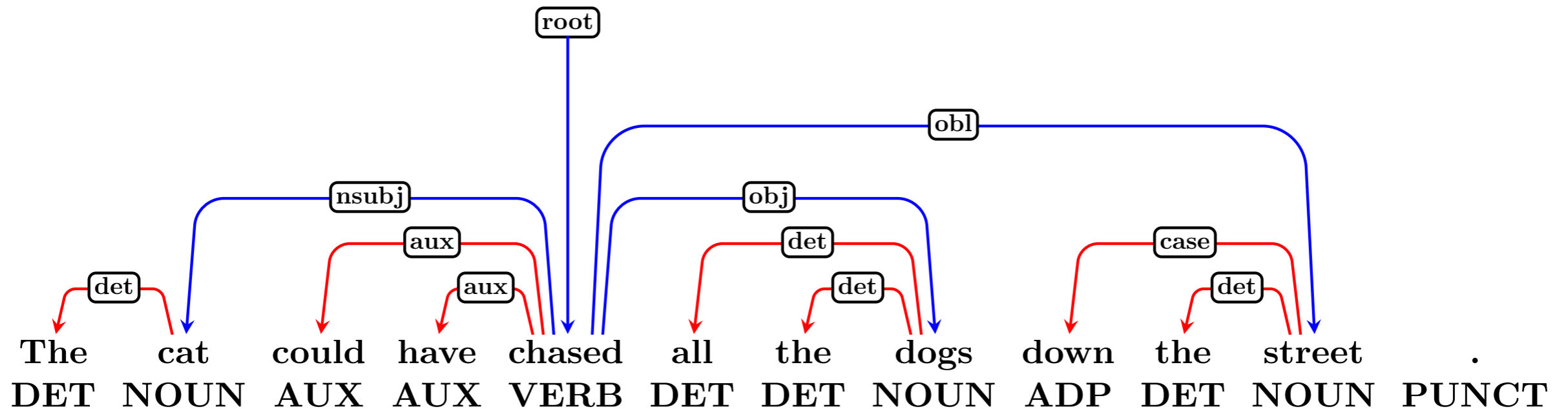
The cat could have chased all the dogs down the street .
DET NOUN AUX AUX VERB DET DET NOUN ADP DET NOUN PUNCT

Syntax



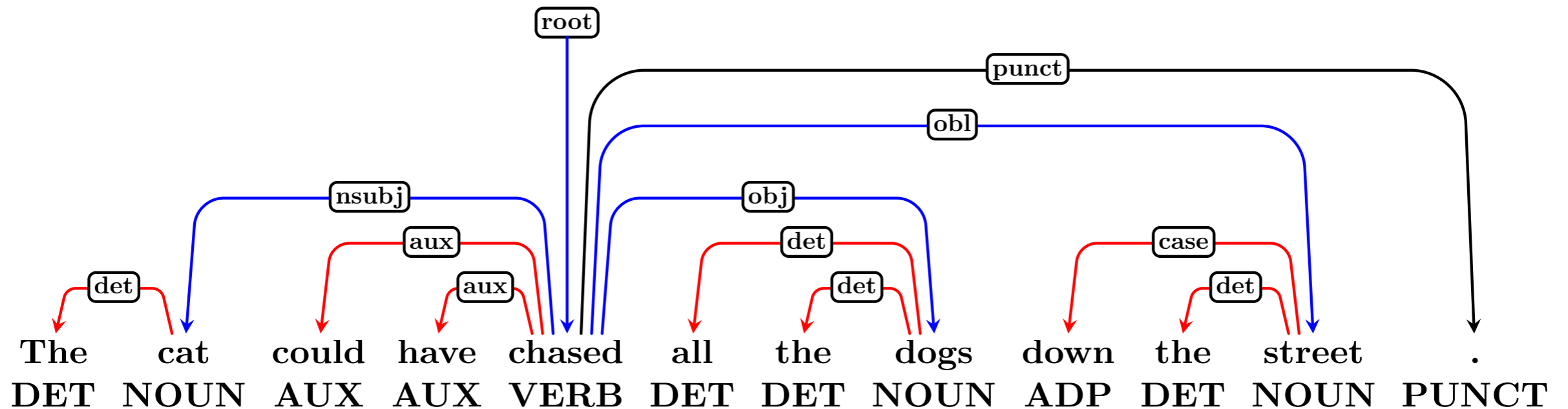
- Content words are related by dependency relations

Syntax

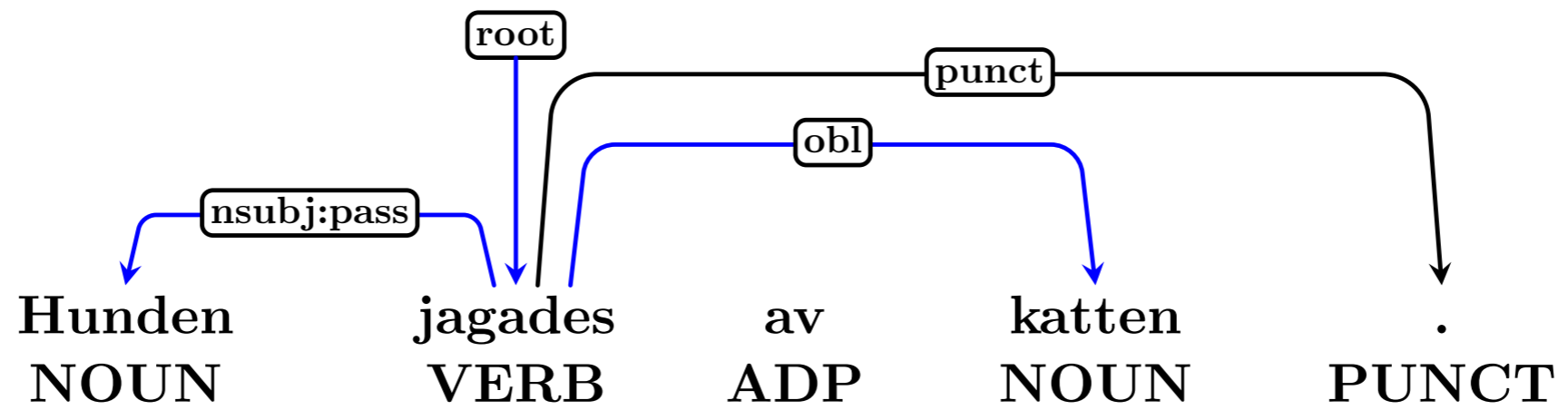
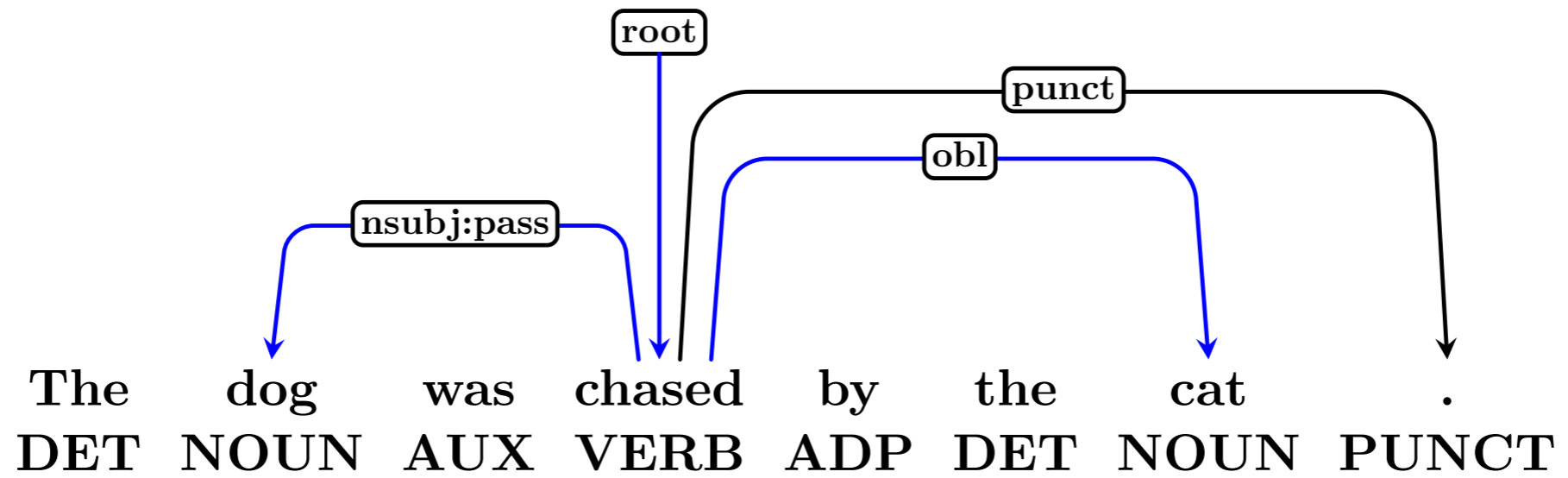


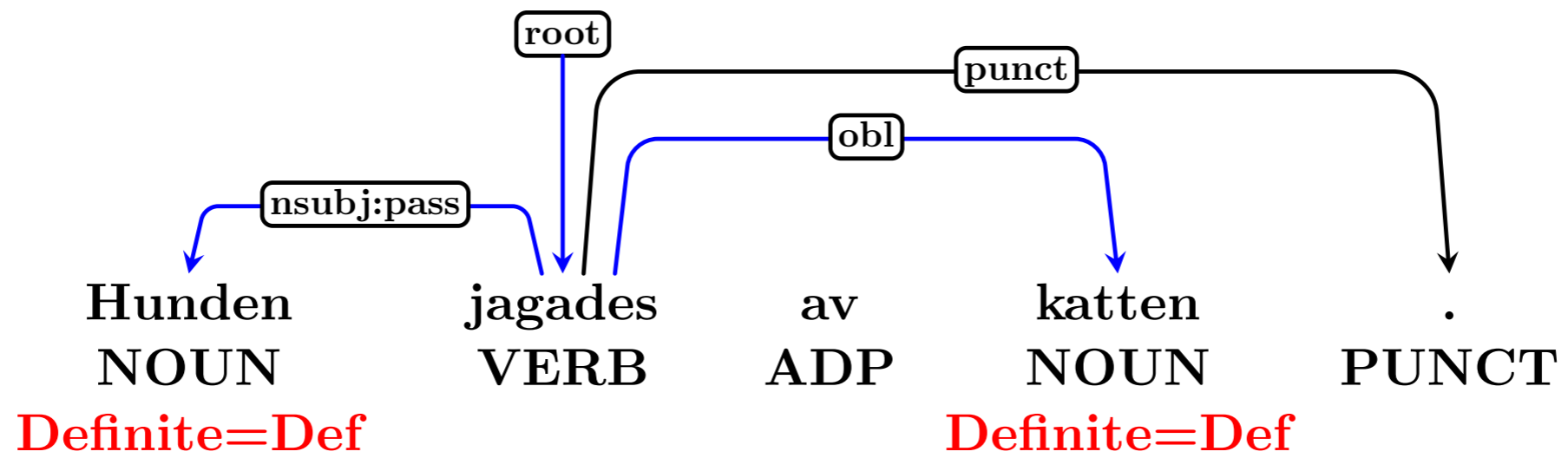
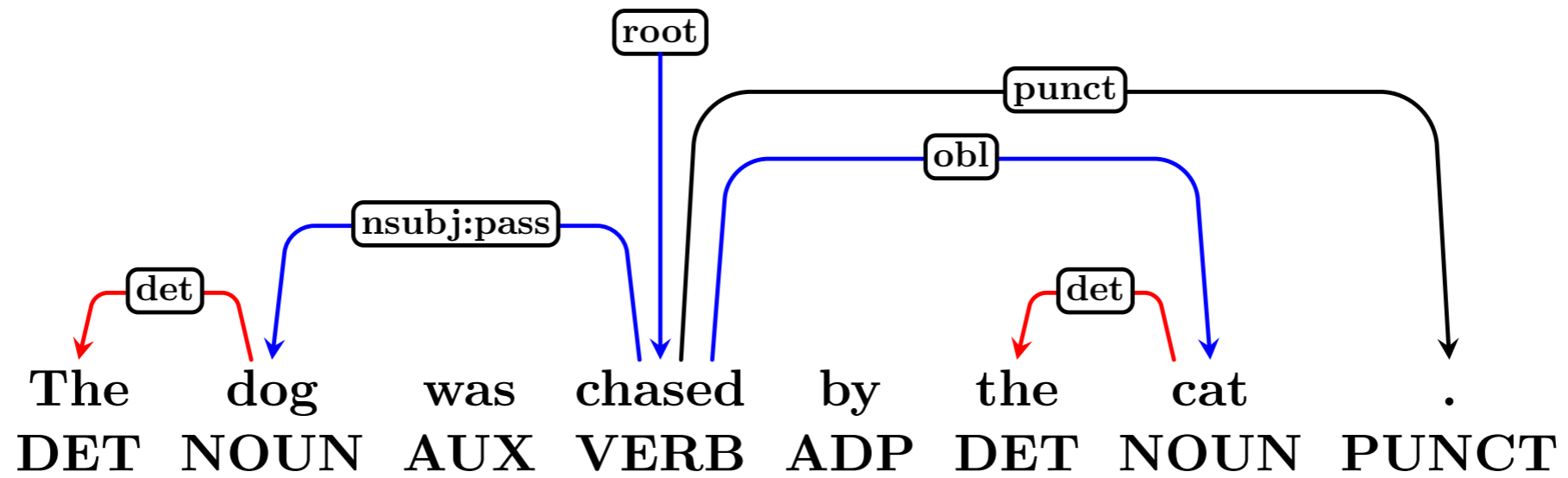
- Content words are related by dependency relations
- Function words attach to the content word they modify

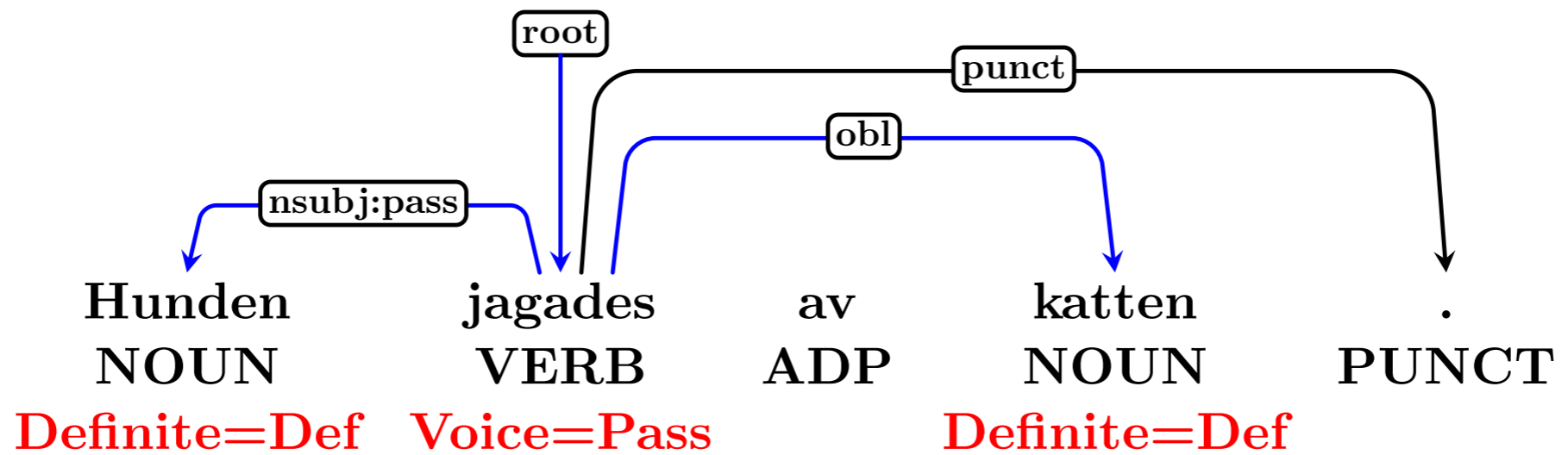
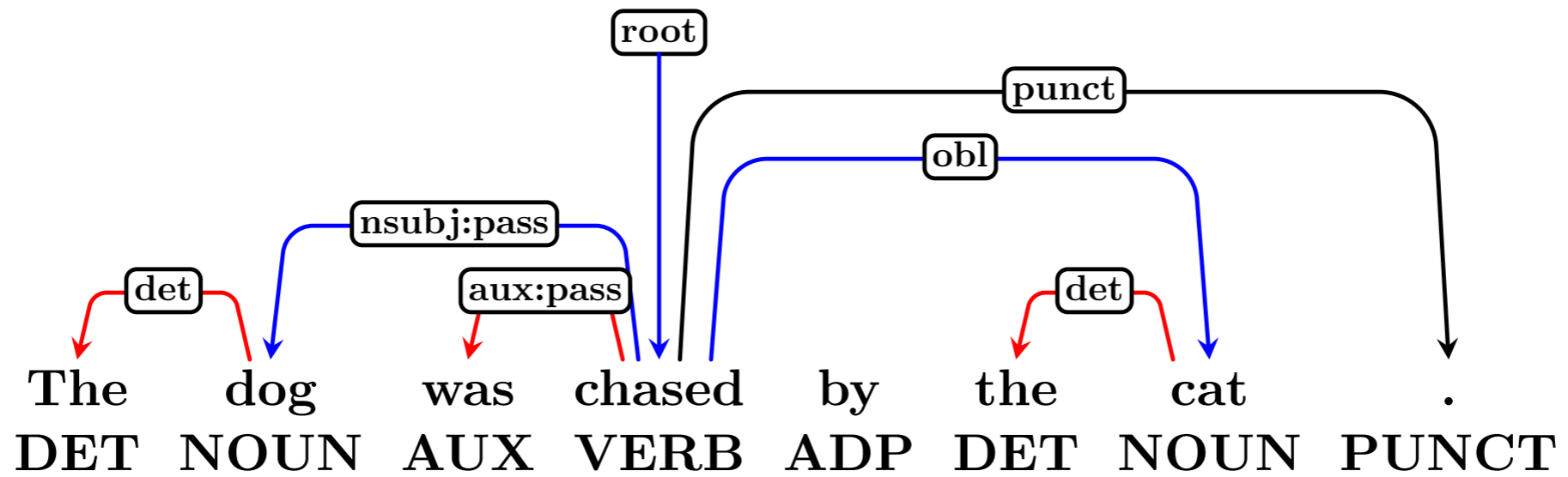
Syntax

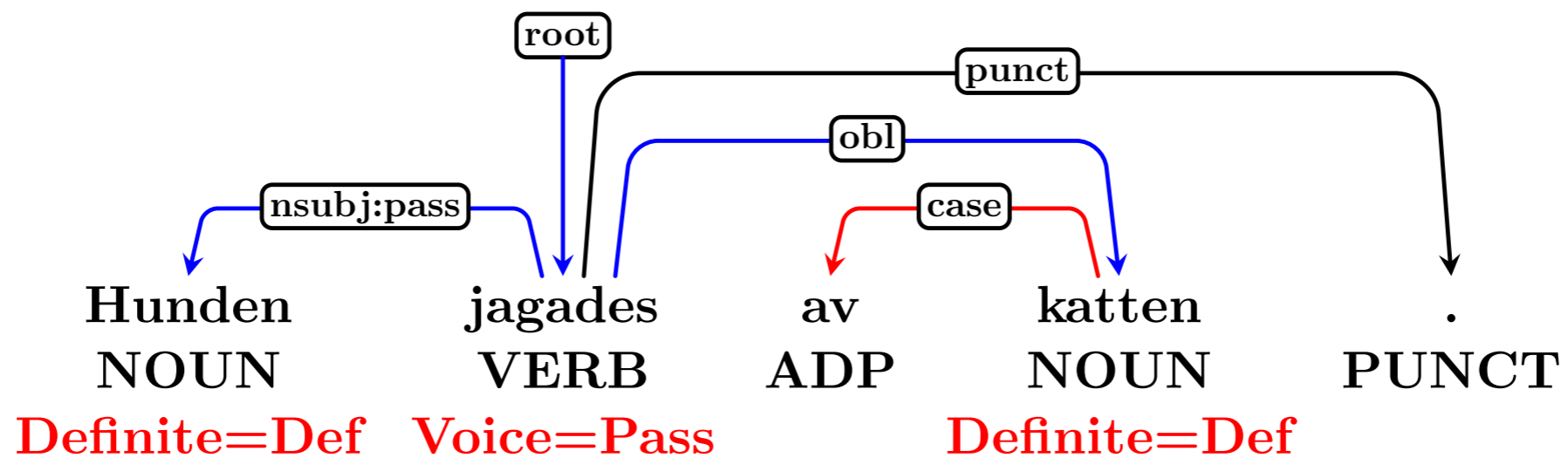
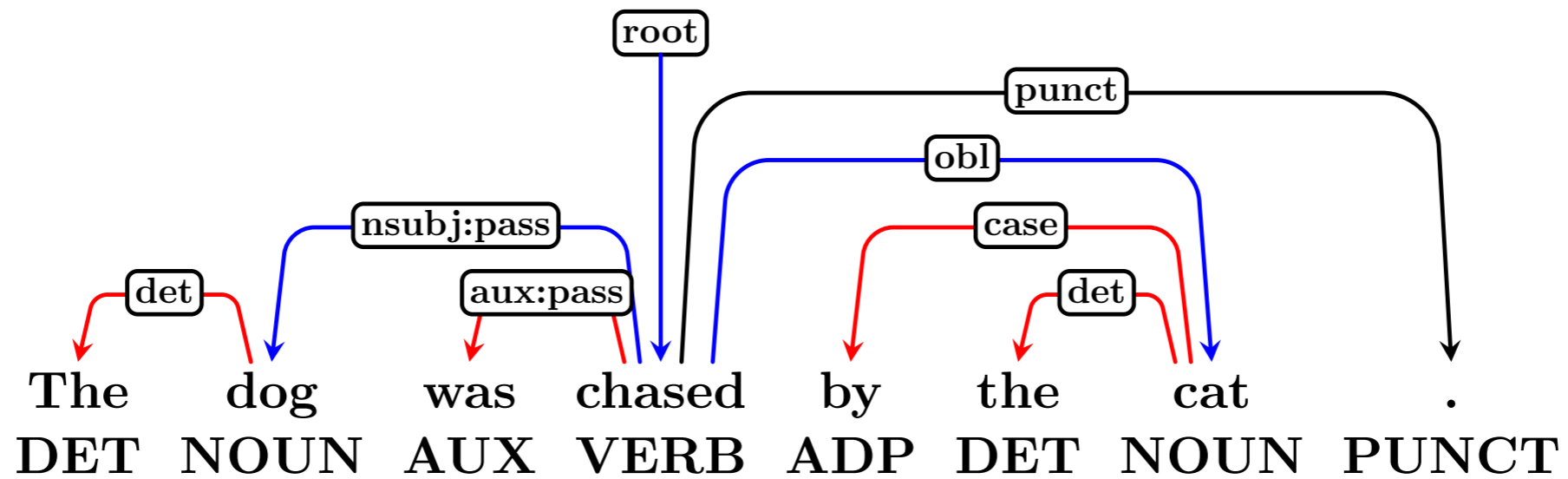


- Content words are related by dependency relations
- Function words attach to the content word they modify
- Punctuation attach to head of phrase or clause









Syntactic Relations

Syntactic Relations

Taxonomy of 37 universal syntactic relations

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

Syntactic Relations

Taxonomy of 37 universal syntactic relations

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

Basic and enhanced representations

- Basic dependencies form a (possibly non-projective) tree
- Additional dependencies in the enhanced representation

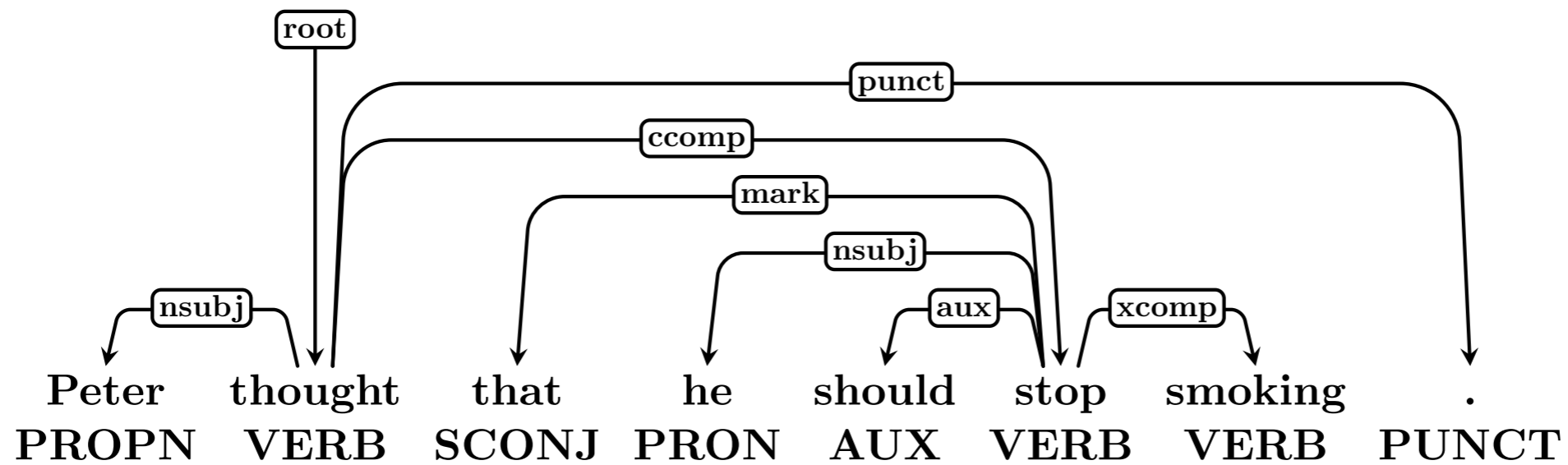
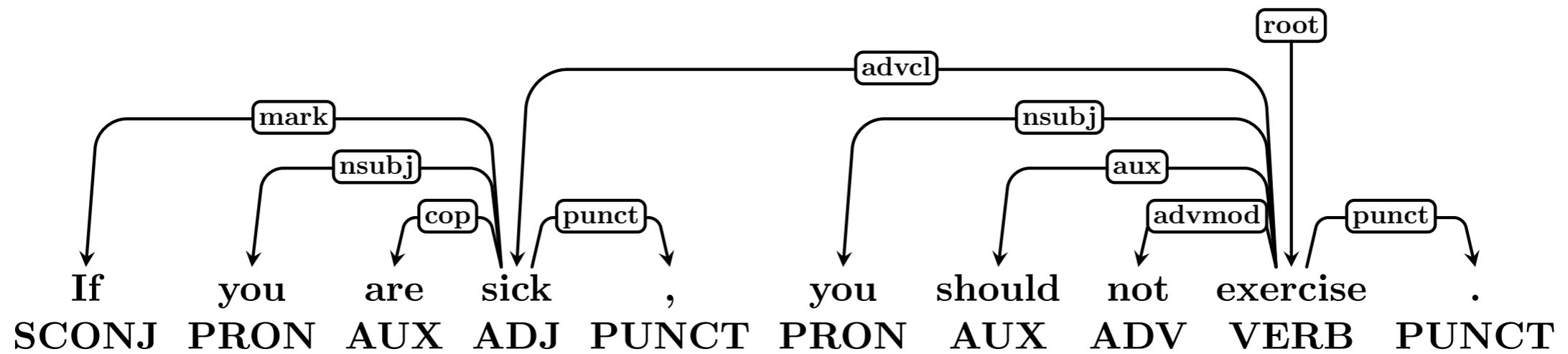
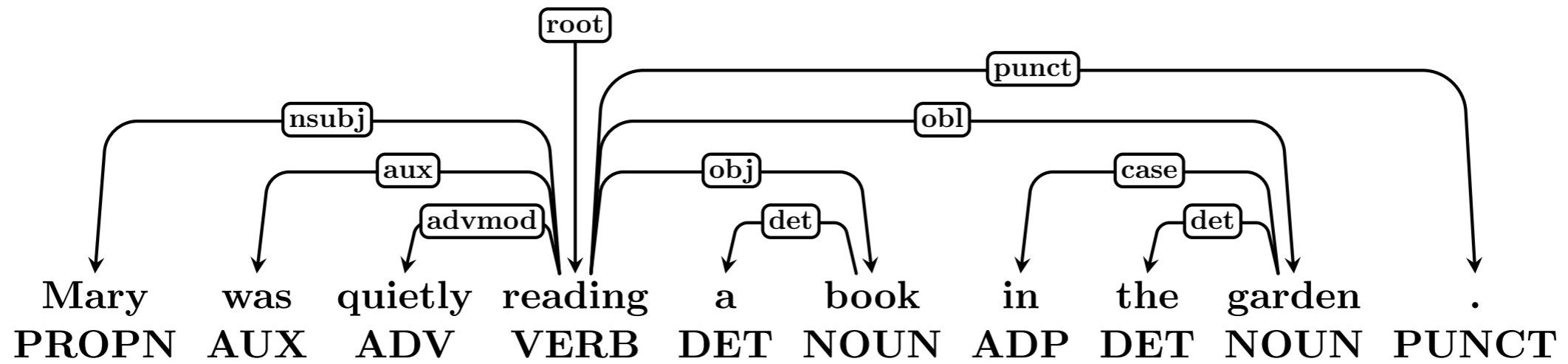
Syntactic Relations

	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

* Generalized modifier of predicates and (non-nominal) modifiers

Dependents of Clausal Predicates

	Nominal	Clause	Modifier Word	Function Word
Core	nsubj obj iobj	csubj ccomp xcomp		
Non-Core	obl vocative expl dislocated	advcl	advmod discourse	aux cop mark



Core Arguments

Core Arguments

Arguments of basic intransitive and transitive verbs

- Verbs usually only agree with core arguments
- Core arguments normally appear as bare nominals without adpositions
- Certain cases, traditionally called nominative, accusative, and absolutive are typically reserved core arguments
- Core arguments often occupy special positions in the clause
- Syntactic phenomena like control, relativization and passivization can be restricted to core arguments

Core Arguments

Arguments of basic intransitive and transitive verbs

- Verbs usually only agree with core arguments
- Core arguments normally appear as bare nominals without adpositions
- Certain cases, traditionally called nominative, accusative, and absolutive are typically reserved core arguments
- Core arguments often occupy special positions in the clause
- Syntactic phenomena like control, relativization and passivization can be restricted to core arguments

Do not confuse

- Core arguments vs. oblique dependents – encoding of grammatical function
- Arguments vs. adjuncts – valency or subcategorization

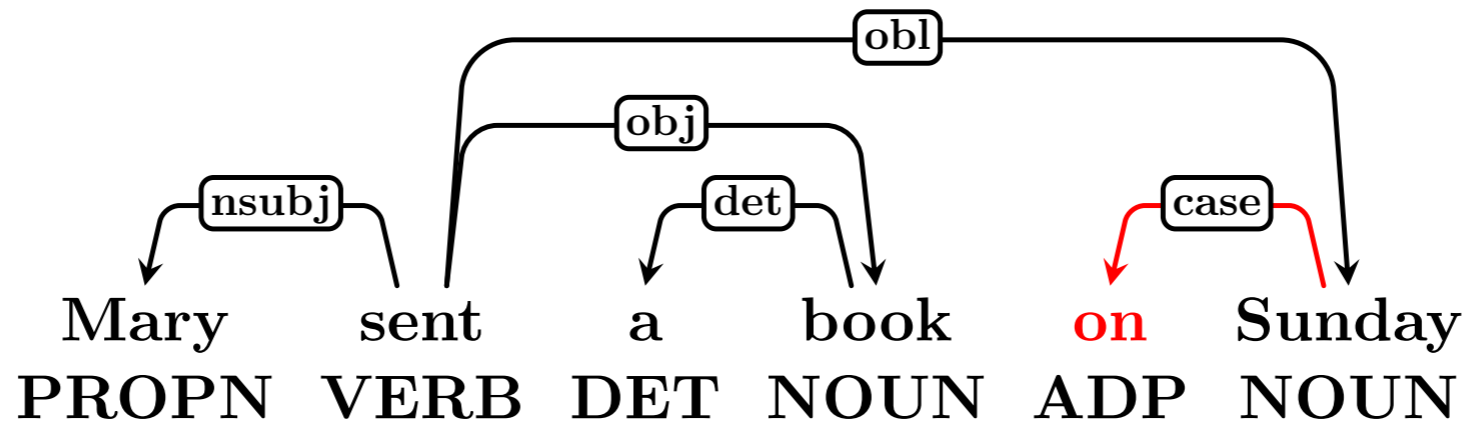
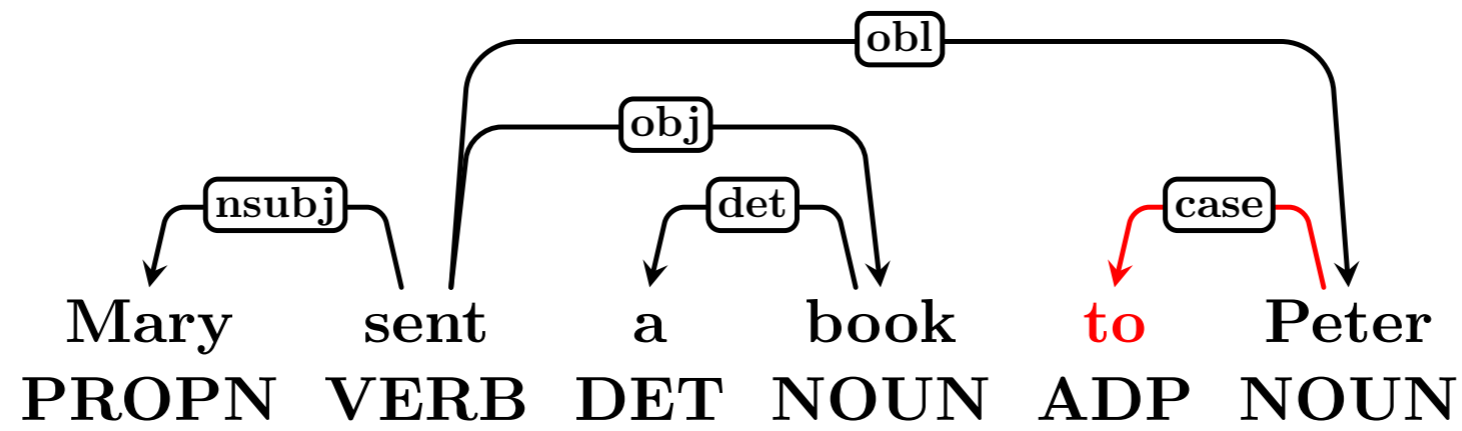
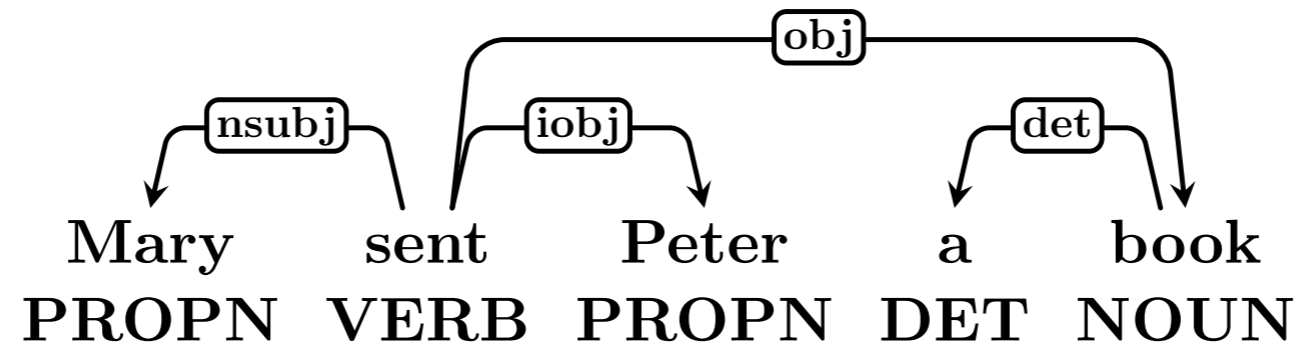
Core Arguments

Arguments of basic intransitive and transitive verbs

- Verbs usually only agree with core arguments
- Core arguments normally appear as bare nominals without adpositions
- Certain cases, traditionally called nominative, accusative, and absolutive are typically reserved core arguments
- Core arguments often occupy special positions in the clause
- Syntactic phenomena like control, relativization and passivization can be restricted to core arguments

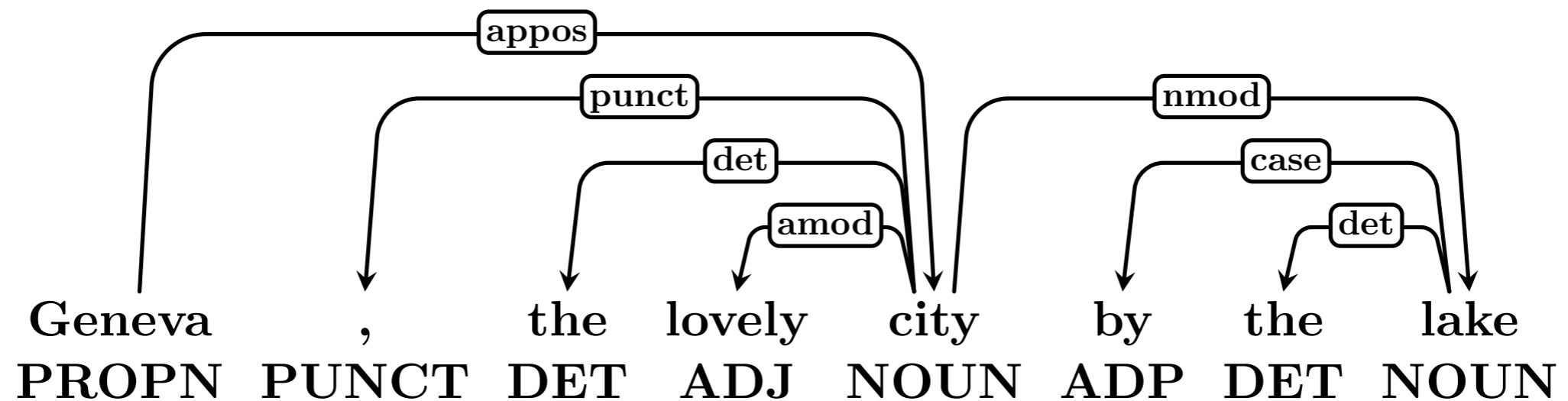
Do not confuse

- Core arguments vs. oblique dependents – encoding of grammatical function
- Arguments vs. adjuncts – valency or subcategorization

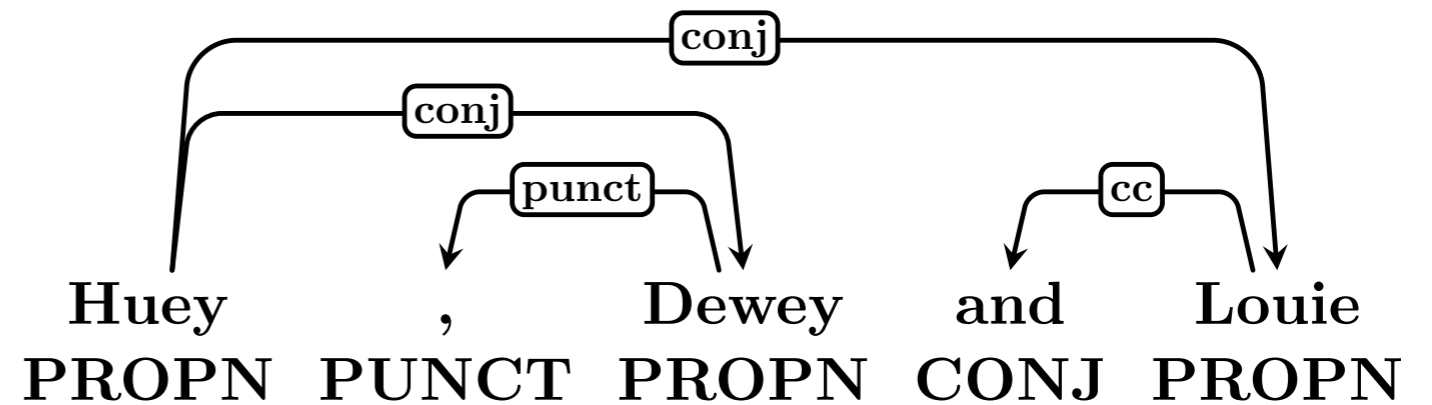
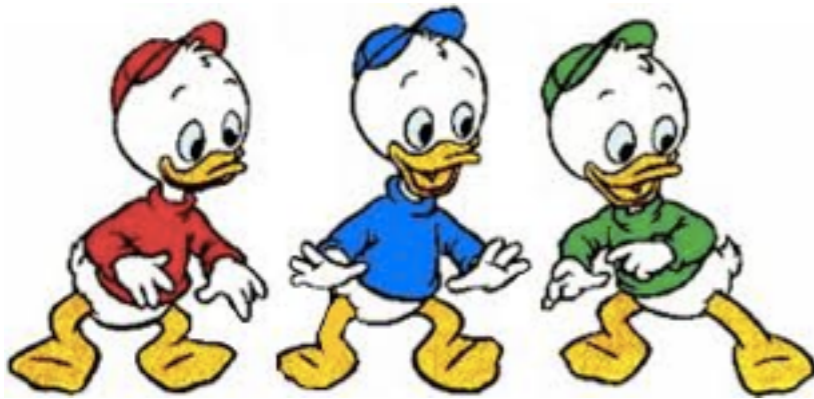


Dependents of Nominals

Nominal	Clause	Modifier Word	Function Word
nmod appos nummod	acl	amod	det clf case



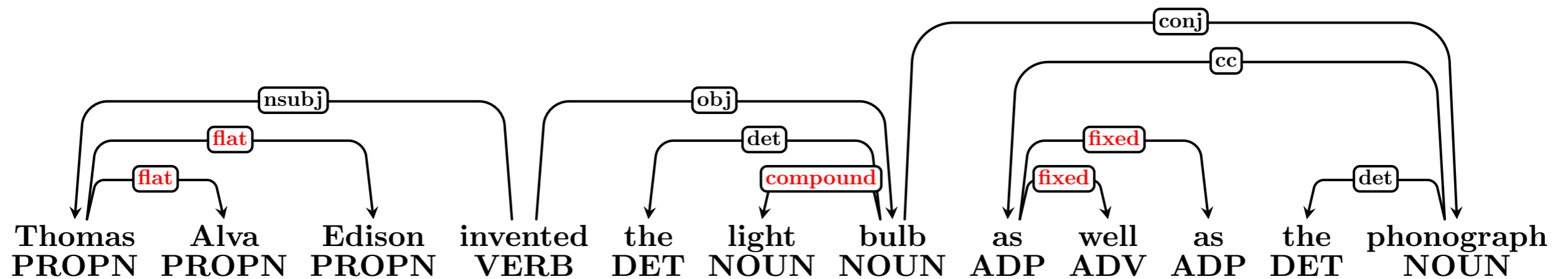
Coordination



Coordinate structures are headed by the first conjunct

- Subsequent conjuncts depend on it via the **conj** relation
- Conjunction depends on following conjunct via the **cc** relation
- Punctuation depends on following conjunct via the **punct** relation

Multiword Expressions



Only restricted classes of MWEs get special treatment:

- Fixed grammaticized expressions (**fixed**)
- Semi-fixed expressions with no clear head (**flat**)
- Lexical compounds – normally headed (**compound**)

Loose Joining Relations

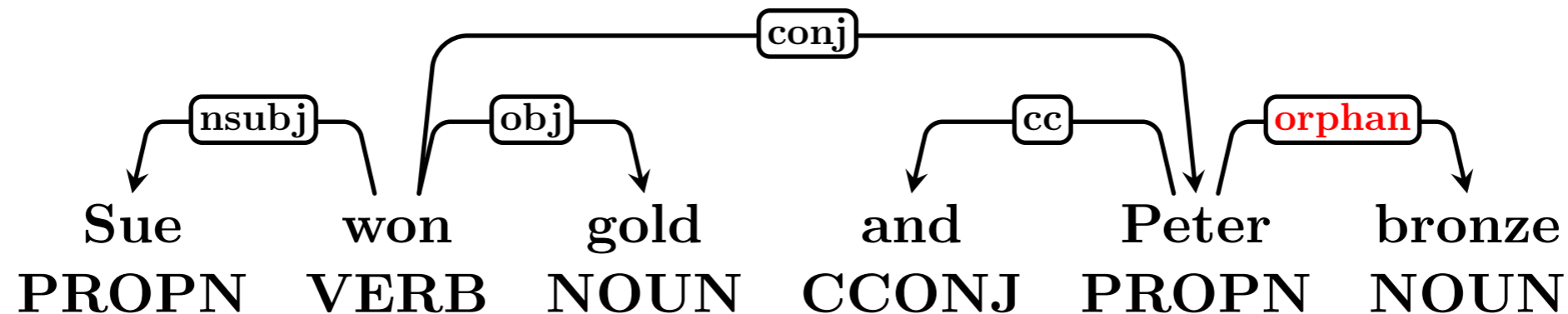
The **parataxis** relation:

- Side-by-side sentences (“run-on sentences”)
Bearded dragons are sight hunters, they need to see the food to move.
- Injective clauses (parentheticals)
Calafia has great fries (they are to die for!) and decent burgers.
- Certain types of reported speech
That guy, he said, left early this morning.
- Tag questions
It's not me, is it?

The **list** relation:

- Chains of comparable items
Steve Jones Phone: 555-9814 Email: jones@abc.edf

Ellipsis

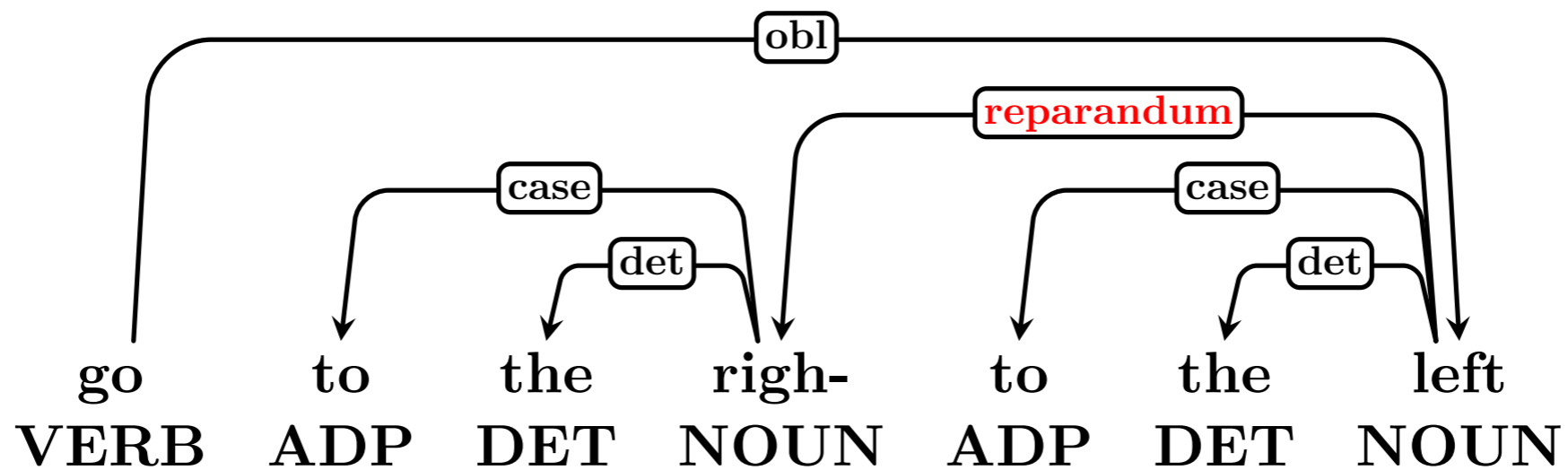


The UD approach to ellipsis (from v2):

1. If the elided word has no children, do nothing.
2. If the elided word has children, promote one of them to be the head.
3. If the elided word is a predicate and the new head a core argument, attach other non-functional elements with the **orphan** relation.

Implicit relations are recovered in enhanced dependencies

Disfluencies



The **reparandum** relation:

- Disfluencies that are overridden in a speech repair

The **goeswith** relation:

- Parts of words resulting from orthographic or editing mistakes

Punctuation

- A punctuation mark separating coordinated units is attached to the following conjunct.
- A punctuation mark preceding or following a subordinated unit is attached to this unit.
- Within the relevant unit, a punctuation mark is attached at the highest possible node that preserves projectivity.
- Paired punctuation marks should be attached to the same word unless that would create non-projectivity.

Special Relations

The **root** relation:

- The word at the root of the dependency tree
- Normally the predicate of the main clause
- Exactly one word in each tree

The **dep** relation:

- Unspecified syntactic relation (when all else fails)

A Two-Level Architecture

- Universal relations to allow cross-linguistic comparison
- Subtypes to capture language-specific phenomena

A Two-Level Architecture

- Universal relations to allow cross-linguistic comparison
- Subtypes to capture language-specific phenomena

Universal

Subtype

A Two-Level Architecture

- Universal relations to allow cross-linguistic comparison
- Subtypes to capture language-specific phenomena

Universal	Subtype
acl	acl:relcl

A Two-Level Architecture

- Universal relations to allow cross-linguistic comparison
- Subtypes to capture language-specific phenomena

Universal	Subtype
acl	acl:relcl
compound	compound:prt

A Two-Level Architecture

- Universal relations to allow cross-linguistic comparison
- Subtypes to capture language-specific phenomena

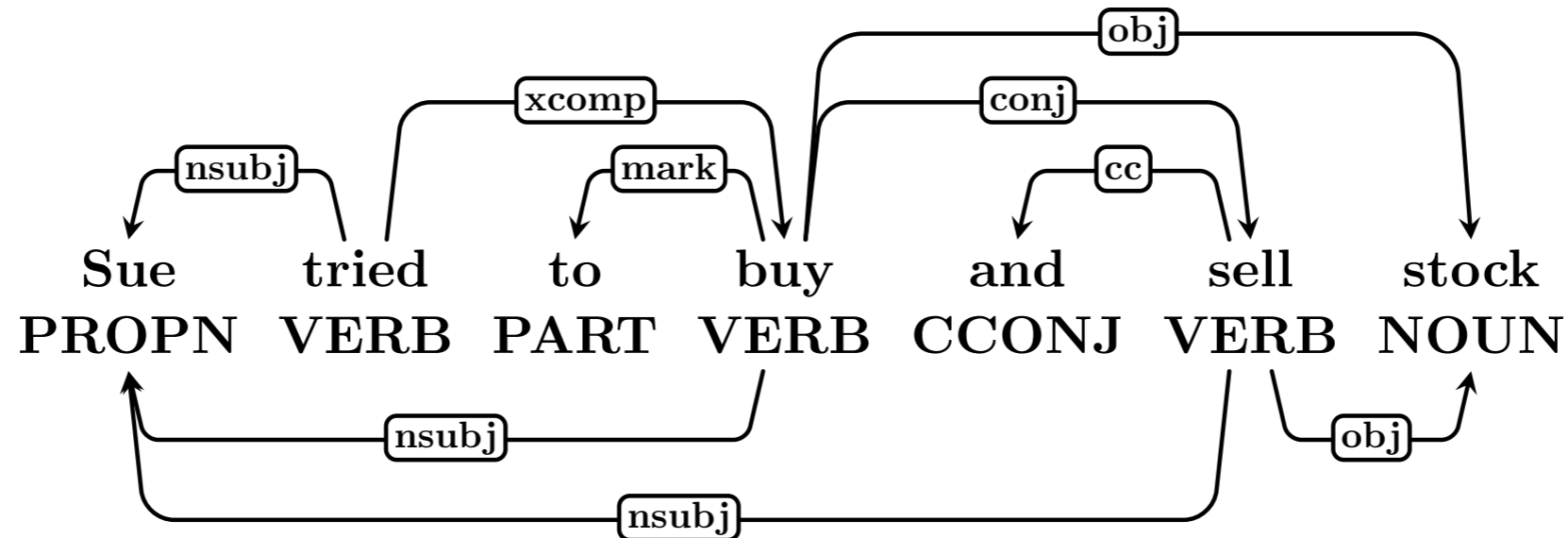
Universal	Subtype
acl	acl:relcl
compound	compound:prt
nmod	nmod:poss

A Two-Level Architecture

- Universal relations to allow cross-linguistic comparison
- Subtypes to capture language-specific phenomena

Universal	Subtype
acl	acl:relcl
compound	compound:prt
nmod	nmod:poss
flat	flat:name

Enhanced Dependencies



An extended dependency graph containing

- Null nodes for elided predicates
- Additional subject relations for control and raising constructions
- Propagation of dependents over coordination
- Coreference in relative clause constructions
- Labels augmented with function word information

CoNLL-U Format

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

CoNLL-U Format

ID
1-2
1
2
3-4
3
4
5
6

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

CoNLL-U Format

ID	FORM
1-2	Vámonos
1	Vamos
2	nos
3-4	al
3	a
4	el
5	mar
6	.

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

CoNLL-U Format

ID	FORM	LEMMA
1-2	Vámonos	—
1	Vamos	ir
2	nos	nosotros
3-4	al	—
3	a	a
4	el	el
5	mar	mar
6	.	.

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

CoNLL-U Format

ID	FORM	LEMMA	UPOSTAG
1-2	Vámonos	—	—
1	Vamos	ir	VERB
2	nos	nosotros	PRON
3-4	al	—	—
3	a	a	ADP
4	el	el	DET
5	mar	mar	NOUN
6	.	.	.

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

CoNLL-U Format

ID	FORM	LEMMA	UPOSTAG	XPOSTAG
1-2	Vámonos	—	—	—
1	Vamos	ir	VERB	—
2	nos	nosotros	PRON	—
3-4	al	—	—	—
3	a	a	ADP	—
4	el	el	DET	—
5	mar	mar	NOUN	—
6	.	.	.	—

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

CoNLL-U Format

ID	FORM	LEMMA	UPOSTAG	XPOSTAG	FEATS
1-2	Vámonos	—	—	—	—
1	Vamos	ir	VERB	—	Mood=Imp Number=Plur Person=1
2	nos	nosotros	PRON	—	PronType=Per Number=Plur Person=1
3-4	al	—	—	—	—
3	a	a	ADP	—	—
4	el	el	DET	—	Definite=Def Number=Sing
5	mar	mar	NOUN	—	Number=Sing Gender=Masc
6	.	.	.	—	—

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

CoNLL-U Format

ID	FORM	LEMMA	UPOSTAG	XPOSTAG	FEATS	HEAD
1-2	Vámonos	—	—	—	—	—
1	Vamos	ir	VERB	—	Mood=Imp Number=Plur Person=1	0
2	nos	nosotros	PRON	—	PronType=Per Number=Plur Person=1	1
3-4	al	—	—	—	—	—
3	a	a	ADP	—	—	5
4	el	el	DET	—	Definite=Def Number=Sing ~ ' ' ~	5
5	mar	mar	NOUN	—	Number=Sing Gender=Masc	1
6	.	.	.	—	—	1

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

CoNLL-U Format

ID	FORM	LEMMA	UPOSTAG	XPOSTAG	FEATS	HEAD	DEPREL
1-2	Vámonos	—	—	—	—	—	—
1	Vamos	ir	VERB	—	Mood=Imp Number=Plur Person=1	0	root
2	nos	nosotros	PRON	—	PronType=Per Number=Plur Person=1	1	expl
3-4	al	—	—	—	—	—	—
3	a	a	ADP	—	—	5	case
4	el	el	DET	—	Definite=Def Number=Sing ~ ' ' ~	5	det
5	mar	mar	NOUN	—	Number=Sing Gender=Masc	1	obl
6	.	.	.	—	—	1	punct

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

CoNLL-U Format

ID	FORM	LEMMA	UPOSTAG	XPOSTAG	FEATS	HEAD	DEPREL	DEPS
1-2	Vámonos	—	—	—	—	—	—	—
1	Vamos	ir	VERB	—	Mood=Imp Number=Plur Person=1	0	root	—
2	nos	nosotros	PRON	—	PronType=Per Number=Plur Person=1	1	expl	—
3-4	al	—	—	—	—	—	—	—
3	a	a	ADP	—	—	5	case	—
4	el	el	DET	—	Definite=Def Number=Sing ~ ' ' ~	5	det	—
5	mar	mar	NOUN	—	Number=Sing Gender=Masc	1	obl	—
6	.	.	.	—	—	1	punct	—

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

CoNLL-U Format

ID	FORM	LEMMA	UPOSTAG	XPOSTAG	FEATS	HEAD	DEPREL	DEPS	MISC
1-2	Vámonos	—	—	—	—	—	—	—	—
1	Vamos	ir	VERB	—	Mood=Imp Number=Plur Person=1	0	root	—	—
2	nos	nosotros	PRON	—	PronType=Per Number=Plur Person=1	1	expl	—	—
3-4	al	—	—	—	—	—	—	—	—
3	a	a	ADP	—	—	5	case	—	—
4	el	el	DET	—	Definite=Def Number=Sing ~ ' ' ~	5	det	—	—
5	mar	mar	NOUN	—	Number=Sing Gender=Masc	1	obl	—	—
6	.	.	.	—	—	1	punct	—	—

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

<http://universaldependencies.org>

So what exactly is UD?

So what exactly is UD?

A new linguistic theory?

- Not at all, but we like to think it is informed by linguistic theory and definitely useful also for linguistic studies

So what exactly is UD?

A new linguistic theory?

- Not at all, but we like to think it is informed by linguistic theory and definitely useful also for linguistic studies

A better parsing framework?

- Maybe not, since parsers seem to prefer function words as heads so we may have to tweak the representations for parsing – but stay tuned

So what exactly is UD?

A new linguistic theory?

- Not at all, but we like to think it is informed by linguistic theory and definitely useful also for linguistic studies

A better parsing framework?

- Maybe not, since parsers seem to prefer function words as heads so we may have to tweak the representations for parsing – but stay tuned

The ultimate annotation scheme?

- Not quite, more like a lingua franca for treebank developers and definitely useful for some annotation projects

So what exactly is UD?

A new linguistic theory?

- Not at all, but we like to think it is informed by linguistic theory and definitely useful also for linguistic studies

A better parsing framework?

- Maybe not, since parsers seem to prefer function words as heads so we may have to tweak the representations for parsing – but stay tuned

The ultimate annotation scheme?

- Not quite, more like a lingua franca for treebank developers and definitely useful for some annotation projects

A universal grammar?

- Not in any strong sense, but hopefully in the sense of providing comparable concepts for meaningful cross-linguistic comparison

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.

UD Events in 2017

Tutorial on Universal Dependencies

- Tutorial at EACL, April 4, 2017, Valencia, Spain



CoNLL-2017 Shared Task

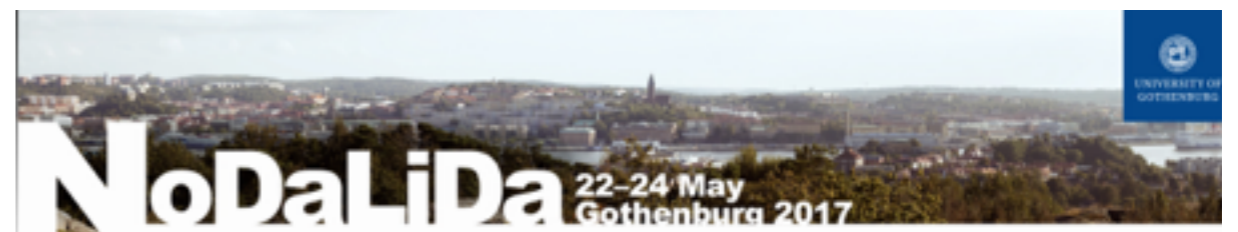
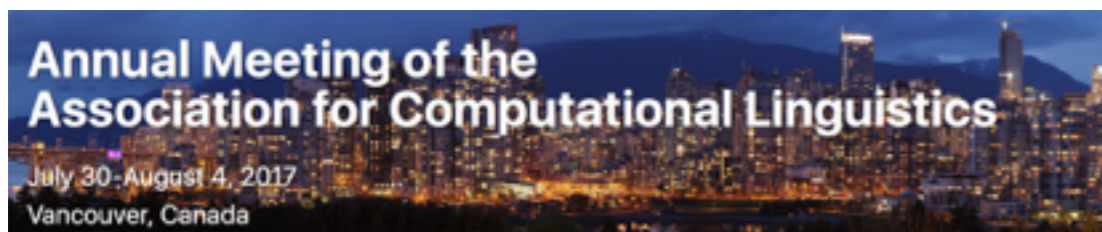
<http://universaldependencies.org/conll17/>

- Multilingual parsing from raw text to universal dependencies
- Collocated with ACL, August 3–4, 2017, Vancouver, Canada

First Workshop on Universal Dependencies

<http://universaldependencies.org/udw17/>

- Collocated with NoDaLiDa, May 20, 2017, Gothenburg, Sweden



Thanks to all UD contributors!

Željko Agić, Lars Ahrenberg, Maria Jesus Aranzabe, Masayuki Asahara, Aitziber Atutxa, Miguel Ballesteros, John Bauer, Kepa Bengoetxea, Yevgeni Berzak, Riyaz Ahmad Bhat, Eckhard Bick, Carl Börstell, Cristina Bosco, Gosse Bouma, Sam Bowman, Marie Candito, Gülşen Cebiroğlu Eryiğit, Giuseppe G.A. Celano, Fabricio Chalub, Minho Choi, Çağrı Çöltekin, Miriam Connor, Elizabeth Davidson, Marie-Catherine de Marneffe, Arantza Diaz de Ilarraza, Kaja Dobrovoljc, Timothy Dozat, Kira Droganova, Puneet Dwivedi, Marhaba Eli, Tomaž Erjavec, Richárd Farkas, Jennifer Foster, Claudia Freitas, Katarína Gajdošová, Daniel Galbraith, Marcos Garcia, Moa Gärdenfors, Sebastian Garza, Filip Ginter, Iakes Goenaga, Koldo Gojenola, Memduh Gökırmak, Yoav Goldberg, Xavier Gómez Guinovart, Berta González Saavedra, Matias Grioni, Normunds Grūzītis, Bruno Guillaume, Nizar Habash, Jan Hajič, Linh Hà Mỹ, Dag Haug, Barbora Hladká, Petter Hohle, Radu Ion, Elena Irimia, Anders Johannsen, Fredrik Jørgensen, Hüner Kaşıkara, Hiroshi Kanayama, Jenna Kanerva, Boris Katz, Jessica Kenney, Natalia Kotsyba, Simon Krek, Veronika Laippala, Lucia Lam, Phương Lê Hồng, Alessandro Lenci, Nikola Ljubešić, Olga Lyashevskaya, Teresa Lynn, Aibek Makazhanov, Christopher Manning, Cătălina Mărănduc, David Mareček, Héctor Martínez Alonso, André Martins, Jan Mašek, Yuji Matsumoto, Ryan McDonald, Anna Missilä, Verginica Mititelu, Yusuke Miyao, Simonetta Montemagni, Amir More, Keiko Sophie Mori, Shunsuke Mori, Bohdan Moskalevskyi, Kadri Muischnek, Nina Mustafina, Kaili Müürisep, Lương Nguyễn Thị, Huyền Nguyễn Thị Minh, Vitaly Nikolaev, Joakim Nivre, Hanna Nurmi, Petya Osenova, Robert Östling, Lilja Øvrelid, Valeria de Paiva, Elena Pascual, Marco Passarotti, Cenel-Augusto Perez, Guy Perrier, Slav Petrov, Jussi Piitulainen, Barbara Plank, Martin Popel, Lauma Pretkalniņa, Prokopis Prokopidis, Tiina Puolakainen, Sampo Pyysalo, Alexandre Rademaker, Loganathan Ramasamy, Livy Real, Laura Rituma, Rudolf Rosa, Manuela Sanguinetti, Shadi Saleh, Baiba Saulīte, Sebastian Schuster, Djamé Seddah, Wolfgang Seeker, Mojgan Seraji, Lena Shakurova, Mo Shen, Natalia Silveira, Maria Simi, Dmitry Sichinava, Radu Simionescu, Katalin Simkó, Mária Šimková, Kiril Simov, Aaron Smith, Carolyn Spadine, Alane Suhr, Umut Sulubacak, Zolt Szántó, Dima Taji, Takaaki Tanaka, Reut Tsarfaty, Francis Tyers, Sumire Uematsu, Larraitz Uria, Gertjan van Noord, Viktor Varga, Veronika Vincze, Lars Wallin, Jing Xian Wang, Jonathan North Washington, Mats Wirén, Zdeněk Žabokrtský, Amir Zeldes, Daniel Zeman, Hanzhi Zhu