



Tutorial on Universal Dependencies

Introduction

Joakim Nivre¹ Daniel Zeman² Filip Ginter³ Francis M. Tyers^{4,5}

¹Department of Linguistics and Philology, Uppsala University, Sweden

²Institute of Formal and Applied Linguistics, Charles University in Prague, Czech Republic

³Department of Information Technology, University of Turku, Finland

⁴Giela ja kultuvrra instituhtta, UiT Norgga árktalaš universitehta, Tromsø, Norway

⁵Arvutiteaduse instituut, Tartu Ülikool, Estonia

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





A cat chases rats and mice

```
graph TD; Root[ ] --- A[A]; Root --- Node1[ ]; Node1 --- cat[cat]; Node1 --- Node2[ ]; Node2 --- chases[chases]; Node2 --- Node3[ ]; Node3 --- rats[rats]; Node3 --- Node4[ ]; Node4 --- and[and]; Node4 --- mice[mice];
```

En katt jagar råttor och möss

```
graph TD; Root[ ] --- En[En]; Root --- Node1[ ]; Node1 --- katt[katt]; Node1 --- Node2[ ]; Node2 --- jagar[jagar]; Node2 --- Node3[ ]; Node3 --- rattor[råttor]; Node3 --- Node4[ ]; Node4 --- och[och]; Node4 --- moss[möss];
```

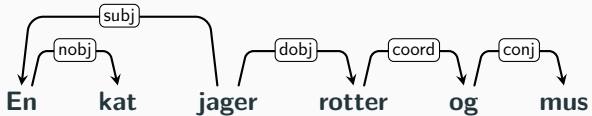
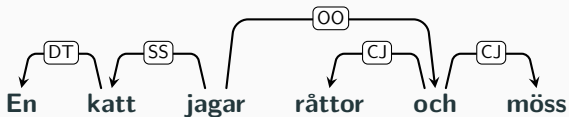
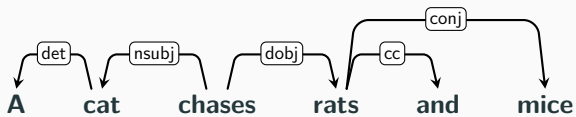


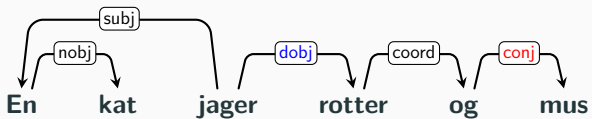
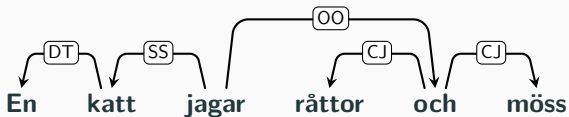
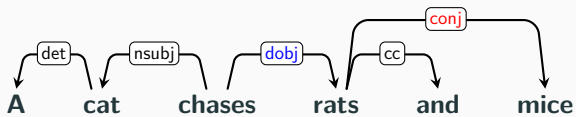
A cat chases rats and mice

En katt jagar råttor och möss

En kat jager rotter og mus







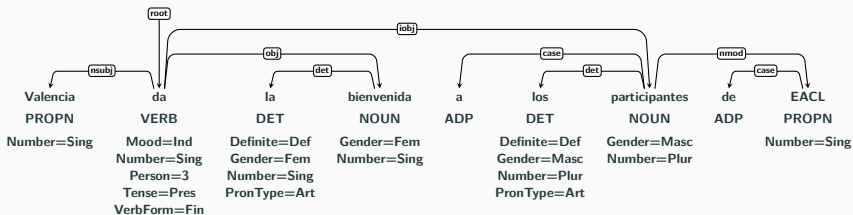
Why is this a problem?

- 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>



- Part-of-speech tags
- Morphological features
- Syntactic dependencies



Goals and Requirements

- Cross-linguistically consistent grammatical annotation
- Support multilingual NLP and linguistic research
- Build on common usage and existing de facto standards
- Complement – not replace – language-specific schemes
- Open community effort – anyone can contribute!



- 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



- 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



Morphological Annotation

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

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

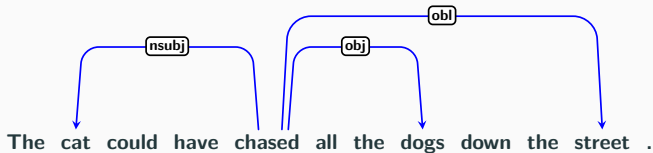


The cat could have chased all the dogs down the street .

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



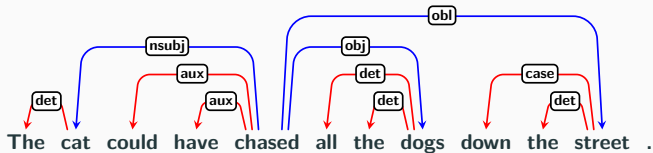
Syntactic Annotation



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



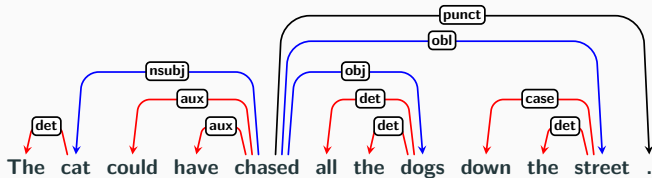
Syntactic Annotation



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



Syntactic Annotation



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



CoNLL-U Format

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

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



Where are we today?

- Brief history of UD:
 - First guidelines launched in October 2014
 - Treebank releases (roughly) every six months
 - Version 2 in December 2016 (guidelines) and March 2017 (treebanks)
- UD in numbers:
 - 50 languages
 - 70 treebanks
 - 164 contributors
 - 8000+ downloads
- Current UD events:
 - CoNLL shared task on UD parsing
 - First UD workshop (Gothenburg, May 22)
 - Next release in November (v2.1)



1. Introduction [Joakim]
2. Cross-linguistically consistent syntactic annotation [Fran]
3. Word segmentation and morphological annotation [Dan]

BREAK

4. Infrastructure, resources and tools for UD [Filip]
5. Making use of UD in NLP and linguistics [Joakim]
6. Adding a new language to UD [Fran]
7. CoNLL shared task on UD parsing [Dan]



Questions?

