

# To tree or not to tree

Zdeněk Žabokrtský

📅 November 6, 2020



EUROPEAN UNION  
European Structural and Investment Fund  
Operational Programme Research,  
Development and Education

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



unless otherwise stated

**Troubles with choosing an annotation  
scheme: a case study on problematic  
corpus/treebank design decisions**

# Corpus annotation criticism

- some critics: an annotated corpus is worse than a raw corpus because of forced interpretations
  - one has to struggle with different linguistic traditions of different national schools
  - example: part of speech categories
- relying on annotation might be misleading if the quality is low (errors or inconsistencies)

# Variability of PoS tag sets

## Penn Treebank POS tagset (for English)

<b>CC</b> coordinating conjunction ( <i>and</i> )	<b>PRP\$</b> possessive pronoun ( <i>my, his</i> )
<b>CD</b> cardinal number ( <i>1, third</i> )	<b>RB</b> adverb ( <i>however, usually, naturally, here, good</i> )
<b>DT</b> determiner ( <i>the</i> )	<b>RBR</b> adverb, comparative ( <i>better</i> )
<b>EX</b> existential there ( <i>there is</i> )	<b>RBS</b> adverb, superlative ( <i>best</i> )
<b>FW</b> foreign word ( <i>d'hœvre</i> )	<b>RP</b> particle ( <i>give up</i> )
<b>IN</b> preposition/subordinating conjunction ( <i>in, of, like</i> )	<b>TO</b> to ( <i>to go, to him</i> )
<b>JJ</b> adjective ( <i>green</i> )	<b>UH</b> interjection ( <i>uhhuhhuhh</i> )
<b>JJR</b> adjective, comparative ( <i>greener</i> )	<b>VB</b> verb, base form ( <i>take</i> )
<b>JJS</b> adjective, superlative ( <i>greenest</i> )	<b>VBD</b> verb, past tense ( <i>took</i> )
<b>LS</b> list marker ( <i>1</i> )	<b>VBG</b> verb, gerund/present participle ( <i>taking</i> )
<b>MD</b> modal ( <i>could, will</i> )	<b>VBN</b> verb, past participle ( <i>taken</i> )
<b>NN</b> noun, singular or mass ( <i>table</i> )	<b>VBP</b> verb, sing. present, non-3d ( <i>take</i> )
<b>NNS</b> noun plural ( <i>tables</i> )	<b>VBZ</b> verb, 3rd person sing. present ( <i>takes</i> )
<b>NNP</b> proper noun, singular ( <i>John</i> )	<b>WDT</b> wh-determiner ( <i>which</i> )
<b>NNPS</b> proper noun, plural ( <i>Vikings</i> )	<b>WP</b> wh-pronoun ( <i>who, what</i> )
<b>PDT</b> predeterminer ( <i>iġboth/iġ the boys</i> )	<b>WP\$</b> possessive wh-pronoun ( <i>whose</i> )
<b>POS</b> possessive ending ( <i>friend's</i> )	<b>WRB</b> wh-adverb ( <i>where, when</i> )
<b>PRP</b> personal pronoun ( <i>I, he, it</i> )	

# Variability of PoS tag sets, cont.

## Negra Corpus POS tagset (for German)

<b>ADJA</b> Attributives Adjektiv	<b>KOKOM</b> Vergleichspartikel, ohne Satz	<b>PRF</b> Reflexives Personalpronomen	<b>VVIZU</b> Infinitiv mit zu, voll
<b>ADJD</b> Adverbiales oder prädikatives Adjektiv	<b>NN</b> Normales Nomen	<b>PWS</b> Substituierendes Interrogativpronomen	<b>VVPP</b> Partizip Perfekt, voll
<b>ADV</b> Adverb	<b>NE</b> Eigennamen	<b>PWAT</b> Attribuierendes Interrogativpronomen	<b>VAFIN</b> Finites Verb, aux
<b>APPR</b> Präposition; Zirkumposition links	<b>PDS</b> Substituierendes Demonstrativpronomen	<b>PWAV</b> Adverbiales Interrogativ- oder Relativpronomen	<b>VAIMP</b> Imperativ, aux
<b>APPRART</b> Präposition mit Artikel	<b>PDAT</b> Attribuierendes Demonstrativpronomen	<b>PROAV</b> Pronominaladverb	<b>VAINF</b> Infinitiv, aux
<b>APPO</b> Postposition	<b>PIS</b> Substituierendes Indefinitpronomen	<b>PTKZU</b> zu vor Infinitiv	<b>VAPP</b> Partizip Perfekt, aux
<b>APZR</b> Zirkumposition rechts	<b>PIAT</b> Attribuierendes Indefinitpronomen	<b>PTKNEG</b> Negationspartikel	<b>VMFIN</b> Finites Verb, modal
<b>ART</b> Bestimmer oder unbestimmter Artikel	<b>PIDAT</b> Attribuierendes Indefinitpronomen mit Determiner	<b>PTKQZ</b> Abgetrennter Verbsatz	<b>VMINF</b> Infinitiv, modal
<b>CARD</b> Kardinalzahl	<b>PPER</b> Irreflexives Personalpronomen	<b>PTKANT</b> Antwortpartikel	<b>VMPP</b> Partizip Perfekt, modal
<b>FM</b> Fremdsprachliches Material	<b>POSS</b> Substituierendes Possessivpronomen	<b>PTKA</b> Partikel bei Adjektiv oder Adverb	<b>XY</b> Nichtwort, Sonderzeichen
<b>ITJ</b> Interjektion	<b>PPOSAT</b> Attribuierendes Possessivpronomen	<b>TRUNC</b> Kompositions-Erstglied	<b>§</b> , Komma
<b>KOUI</b> Unterordnende Konjunktion mit zu und Infinitiv	<b>PRELS</b> Substituierendes Relativpronomen	<b>VVFIN</b> Finites Verb, voll	<b>§.</b> Satzbeendende Interpunktion
<b>KOUS</b> Unterordnende Konjunktion mit Satz	<b>RELAT</b> Attribuierendes Relativpronomen	<b>VVIMP</b> Imperativ, voll	<b>\$(</b> Sonstige Satzzeichen; Satzintern
<b>KON</b> Nebenordnende Konjunktion		<b>VVINF</b> Infinitiv, voll	<b>NNE</b> Verbindung aus Eigennamen und normalen Nomen

## Variability of PoS tag sets, cont.

Prague Dependency Treebank morphologitagset (for Czech), several thousand combinations using 15-character long positional tags

Form	Lemma	Morphological tag
<i>Některé</i>	<i>některý</i>	PZFP1-----
<i>kontury</i>	<i>kontura</i>	NNFP1-----A----
<i>problému</i>	<i>problém</i>	NNIS2-----A----
<i>se</i>	<i>se_^(zvr._zájmeno/částice)</i>	P7-X4-----
<i>však</i>	<i>však</i>	J^-----
<i>po</i>	<i>po-1</i>	RR--6-----
<i>oživení</i>	<i>oživení_^(*3it)</i>	NNNS6-----A----
<i>Havlovým</i>	<i>Havlův;S_^(*3el)</i>	AUIS7M-----
<i>projevem</i>	<i>projev</i>	NNIS7-----A----
<i>zdají</i>	<i>zdat</i>	VB-P---3P-AA---
<i>být</i>	<i>být</i>	Vf-----A----
<i>jasnější</i>	<i>jasný</i>	AAFP1-----2A----
.	.	Z:-----



## Why trees: Initial thoughts

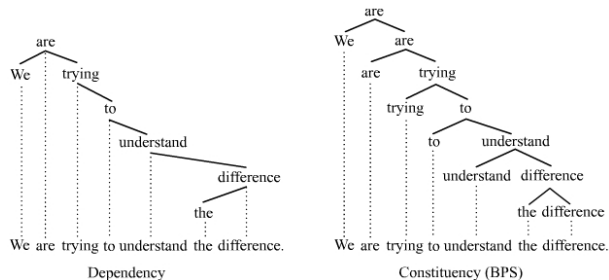
1. Honestly: trees are irresistibly attractive data structures.
2. We believe sentences can be reasonably represented by discrete units and relations among them.
3. Some relations among sentence components (such as some word groupings) make more sense than others.
4. In other words, we believe there is an latent but identifiable discrete structure hidden in each sentence.
5. The structure must allow for various kinds of nestedness (*...a já mu řek, že nejsem Řek, abych mu řek, kolik je v Řecku řeckých řek ...*).
6. This resembles recursivity. Recursivity reminds us of trees.
7. Let's try to find such trees that make sense linguistically and can be supported by empirical evidence.
8. Let's hope they'll be useful in developing NLP applications such as Machine Translation.



## So what kind of trees?

There are two types of trees broadly used:

- constituency (phrase-structure) trees
- dependency trees



Credit: Wikipedia

Constituency trees simply don't fit to languages with freer word order, such as Czech. Let's use dependency trees.

## BTW how do we know there is a dependency between two words?

- There are various clues manifested, such as
  - word order (juxtaposition): “...*přijdu* *zítra* ...”
  - agreement: “...*novými*<sub>.pl.instr</sub> *knihami*<sub>.pl.instr</sub>...”
  - government: “...*slíbil* *Petrovi*<sub>.dative</sub>...”
- Different languages use different mixtures of morphological strategies to express relations among sentence units.

## Basic assumptions about building units

If a sentence is to be represented by a dependency tree, then we need to be able to:

- identify **sentence boundaries**.
- identify **word boundaries** within a sentence.

# Basic assumptions about dependencies

If a sentence is to be represented by a dependency tree, then:

- there must be a **unique parent word** for each word in each sentence, except for the root word
- there are **no loops** allowed.

## Even the most basic assumptions are violated

- Sometimes **sentence boundaries are unclear** – generally in speech, but e.g. in written Arabic too, and in some situations even in written Czech (e.g. direct speech)
- Sometimes **word boundaries are unclear**, (Chinese, “ins” in German, “abych” in Czech).
- Sometimes its **unclear which words should become parents** (A preposition or a noun? An auxiliary verb or a meaningful verb? ...).
- Sometimes there are too many relations (“Zahlédla ho bosého.”), which implies **loops**.

Life's hard. Let's ignore it and insist on trees.

## Counter-examples revisited

If we cannot find linguistically justified decisions, then make them at least consistent.

- Sometimes sentence boundaries are unclear (generally in speech, but e.g. in written Arabic too...)
  - **OK, so let's introduce annotation rules for sentence segmentation.**
- Sometimes word boundaries are unclear, (Chinese, “ins” in German, “abych” in Czech).
  - **OK, so let's introduce annotation rules for tokenization.**
- Sometimes it's not clear which word should become parent (e.g. a preposition or a noun?).
  - **OK, so let's introduce annotation rules for choosing parent.**
- Sometimes there are too many relations (“Zahlédla ho bosého.”), which implies loops.
  - **OK, so let's introduce annotation rules for choosing tree-shaped skeleton.**

# Trebanking

- Is our dependency approach viable? Can we check it?
- Let's start by building the trees manually.
- a treebank - a collection of sentences and associated (typically manually annotated) dependency trees
- for English: Penn Treebank [Marcus et al., 1993]
- for Czech: Prague Dependency Treebank [Hajič et al., 2001]
  - layered annotation scheme: morphology, surface syntax, deep syntax
  - dependency trees for about 100,000 sentences
- high degree of design freedom and local linguistic tradition bias
- different treebanks  $\implies$  different annotation styles

# An example of a treebank variability cause: the case of coordination

- coordination structures such as “*lazy dogs, cats and rats*” consists of
  - conjuncts
  - conjunctions
  - shared modifiers
  - punctuation tokens
- 16 different annotation styles identified in 26 treebanks (and many more possible)
- different expressivity, limited convertibility, limited comparability of experiments...
- **harmonization of annotation styles badly needed!**

Main family	Prague family (code IP) [14 treebanks]	Moscow family (code IM) [5 treebanks]	Stanford family (code IS) [6 treebanks]
<b>Choice of head</b>			
Head on left (code hL) [10 treebanks]			
Head on right (code hR) [14 treebanks]			
Mixed head (code hM) [1 treebank]	A mixture of hL and hR		
<b>Attachment of shared modifiers</b>			
Shared modifier below the nearest conjunct (code sN) [15 treebanks]			
Shared modifier below head (code sH) [11 treebanks]			
<b>Attachment of coordinating conjunction</b>			
Coordinating conjunction below previous conjunct (code cP) [2 treebanks]	—		
Coordinating conjunction below following conjunct (code cF) [1 treebank]	—		
Coordinating conjunction between two conjuncts (code cB) [8 treebanks]	—		
Coordinating conjunction as the head (code cH) is the only applicable style for the Prague family [14 treebanks]	—	—	—
<b>Placement of punctuation</b>			
values pP [7 treebanks], pF [1 treebank] and pB [15 treebanks] are analogous to cP, cF and cB (but applicable also to the Prague family)			



## Btw how many treebanks are there out there?

- growing interest in dependency treebanks in the last decade or two
- existing treebanks for about 100 languages now (but roughly 7,000 languages in the world)
- UFAL participated in several treebank unification efforts:
  - 13 languages in CoNLL in 2006
  - 29 languages in HamleDT in 2011
  - 37 languages in Universal Dependencies in 2015:
  - 70+ languages in UD in 2019

# Conclusion

- one should keep in mind that there's no straightforward “God's truth” when it comes to language data resources
- all resources are heavily influenced by numerous design choices, for which no perfect answers exists
- examples of trade-offs:
  - the bigger data the better, but you can't remove all noise from really big data
  - parallel annotation reduces the amount of annotation errors, but increases costs
  - linguistically-based annotation brings interpretability, but at the same time we risk being trapped in some suboptimal traditions that are possibly not useful beyond a given language family
  - a better quality/coverage is sometimes achievable by integrating more resources focused on a same task, but their licenses might be incompatible