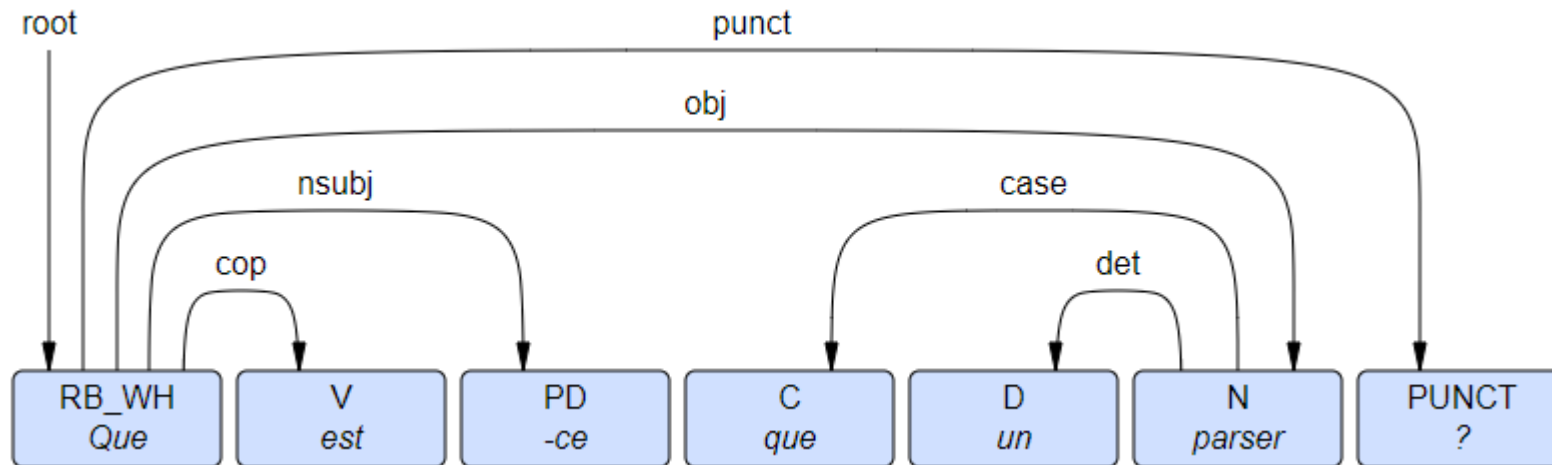


# Part-of-speech tagging, dependency parsing, and named entity recognition



Prof Dr Marko Robnik-Šikonja

Natural Language Processing, Version 2024

# Contents

- POS tagging
- Tag sets
- Dependency parsing
- Universal dependencies
- Named entity recognition

# Basic text processing

- document → paragraphs → sentences → words
- words and sentences ← **POS tagging**
- sentences ← **syntactical and grammatical analysis**

# An Example

WORD	LEMMA	TAG
the	the	+DET
girl	girl	+NOUN
kissed	kiss	+VPAST
the	the	+DET
boy	boy	+NOUN
on	on	+PREP
the	the	+DET
cheek	cheek	+NOUN

# First step: lemmatization

- Lemmatization is the process of grouping together the different inflected forms of a word so they can be analyzed as a single item.
- Lemmatization difficulty is language dependent, i.e. it depends on morphology
- *English*
  - *walk, walked, walking, walks, ne pa walker*
  - *go, goes, going, gone, went*
- *Slovene*
  - *priti, pridem, prideš, pride, prideva, prideta, pridejo, pridemo, pridete, pridejo,*  
but not *prihod, prihodnost, prihajanje, prišlec*
  - *vlak, vlaka, vlaku, vlakom, vlakov, vlakoma, vlakih, vlaki, vlake*
  - *jaz, mene, meni, mano*
  - *Gori na gori gori!*
  - *Gori, na gori gori!*

# Approaches to lemmatization

- Rules, dictionaries, lexicons, machine learning models
- Ambiguity resolution may be difficult

Meni je vzel z mize (zapestnico). Zaradi vrata ni mogel odpreti vrat.

- Quick solutions and heuristics, in English just remove suffixes: *-ing, -ation, -ed, ...*
- Essential approach for morphologically rich languages (Slavic, Arabic, Turkish, Spanish, etc)

# Part-of-Speech Tagging

- Assigning a part-of-speech to each word in a text.
- Words often have more than one POS.
- **book:**
  - VERB: (***Book** that flight*)
  - NOUN: (*Hand me that **book***).

# POS tagging

- Assigning the correct part of speech (noun, verb, etc.) to words
- Helps in recognizing phrases, names, terminology
- Helps in information retrieval, advanced search, named entity recognition, word sense disambiguation, coreference resolution, pronunciation, additional information for many classification tasks, useful heuristic for some tasks
- Helps in linguistic analyses such as verb valence, detection of multi-word expressions, semantic role labelling (SRL)
- Uses machine learning models



# POS tagging for speech

- Speech synthesis:
  - How to pronounce “lead”? /li:d/ or /led/
  - INsult                      insult                      noun: /'ɪnsʌlt/                      verb: /ɪn'sʌlt/
  - OBject                      obJECT
  - OVERflow                      overFLOW
  - DIScount                      disCOUNT
  - CONtent                      content
- In Slovene
  - peti (to sing)                      peti (the fifth)
- Machine translation
  - The meaning of a particular word depends on its POS tag
- Sentiment analysis
  - Adjectives are the major opinion holders (good vs. bad, excellent vs. terrible)

# Morphosyntactical tagging

- POS tagging
- Basic categories from old Greek
  - noun, verb, pronoun, preposition, adjective/adverb, conjunction, participle, and article
  - samostalnik, glagol, zaimек, predlog, pridevnik/prislov, veznik, deležnik, členek
- Many additional features with important information: gender, tense, conjugation, etc.
- Tags defined based on
  - word morphology, e.g., suffixes and prefixes
  - distributional properties, i.e. neighborhood words, role in sentence
- Important part of disambiguation

# POS examples

- N        noun        chair, bandwidth, pacing
- V        verb        study, debate, munch
- ADJ     adjective    purple, tall, ridiculous
- ADV     adverb        unfortunately, slowly,
- P        preposition   of, by, to
- PRO     pronoun     I, me, mine
- DET     determiner   the, a, that, those

# Open and closed class words

- Closed class: a relatively fixed membership
  - Prepositions: of, in, by, ...
  - Auxiliaries: may, can, will had, been, ...
  - Pronouns: I, you, she, mine, his, them, ...
  - Usually **function words** (short common words which play a role in grammar)
- Open class: new ones can be created all the time
  - English has 4: Nouns, Verbs, Adjectives, Adverbs
  - Many languages have all 4, but not all!
  - In Lakhota and possibly Chinese, what English treats as adjectives act more like verbs.
  - New nouns and verbs like *iPhone* or *to fax*

# Open class words

- Nouns
  - Proper nouns (Columbia University, New York City, Arthi Ramachandran, Metropolitan Transit Center). English capitalizes these.
  - Common nouns (the rest). German capitalizes these.
  - Count nouns and mass nouns
    - Count: have plurals, get counted: goat/goats, one goat, two goats
    - Mass: don't get counted (fish, salt, communism)  
(\*two fishes refers to two species of fish)
- Adverbs: tend to modify things
  - Unfortunately, John walked home extremely slowly yesterday
  - Directional/locative adverbs (here, home, downhill)
  - Degree adverbs (extremely, very, somewhat)
  - Manner adverbs (slowly, slinkily, delicately)

# Open class words

- Verbs:
  - In English, they have morphological affixes (eat/eats/eaten)
  - Actions (walk, ate) and states (be, exude)
  - Many subclasses, e.g.
    - eats/VBZ, eat/VB, eat/VBP, eats/VBZ, ate/VBD, eaten/VBN, eating/VBG, ...
    - Reflect morphological form & syntactic function

## Open class ("content") words

### Nouns

#### Proper

*Janet*  
*Italy*

#### Common

*cat, cats*  
*mango*

### Verbs

#### Main

*eat*  
*went*

#### Auxiliary

*can*  
*had*

### Adjectives

*old green tasty*

### Adverbs

*slowly yesterday*

### Numbers

*122,312*  
*one*

Interjections *Ow hello*

*... more*

## Closed class ("function")

Determiners *the some*

Conjunctions *and or*

Pronouns *they its*

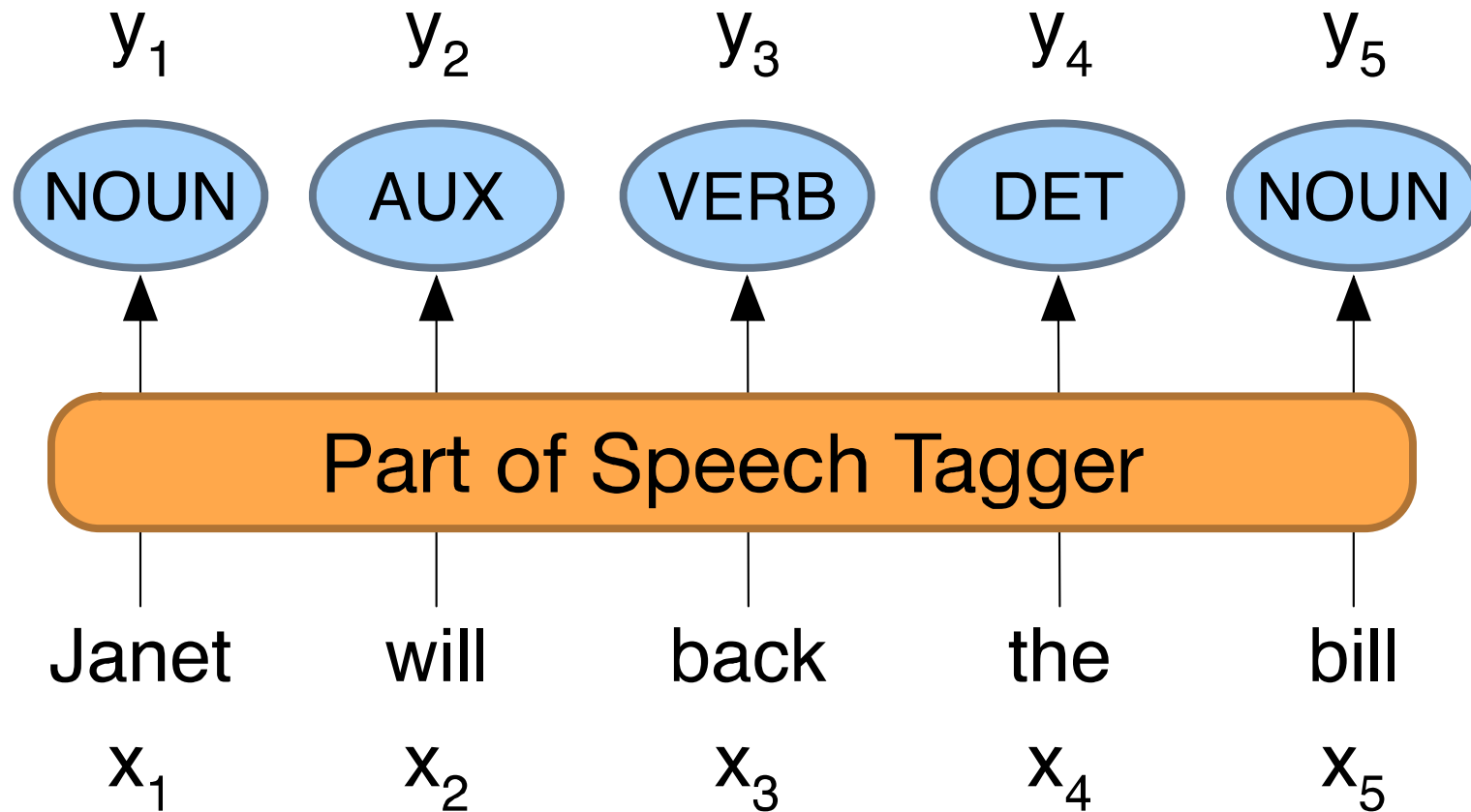
Prepositions *to with*

Particles *off up*

*... more*

# Part-of-Speech Tagging

Map from sequence  $x_1, \dots, x_n$  of words to  $y_1, \dots, y_n$  of POS tags





# Word classes: tag sets

- Vary in number of tags: for English from a dozen to over 200
- Size of tag sets depends on language, objectives and purpose
- We have to agree on a standard inventory of word classes
  - Taggers are trained on a labeled corpora
  - The tag set needs to capture semantically or syntactically important distinctions that can easily be made by trained human annotators

# Tag set example

- e.g., Penn-Treebank tag set
- between 45 and 70 tags

Tag	Description	Example	Tag	Description	Example
CC	coordin. conjunction	<i>and, but, or</i>	SYM	symbol	<i>+, %, &amp;</i>
CD	cardinal number	<i>one, two</i>	TO	“to”	<i>to</i>
DT	determiner	<i>a, the</i>	UH	interjection	<i>ah, oops</i>
EX	existential ‘there’	<i>there</i>	VB	verb base form	<i>eat</i>
FW	foreign word	<i>mea culpa</i>	VBD	verb past tense	<i>ate</i>
IN	preposition/sub-conj	<i>of, in, by</i>	VBG	verb gerund	<i>eating</i>
JJ	adjective	<i>yellow</i>	VCN	verb past participle	<i>eaten</i>
JJR	adj., comparative	<i>bigger</i>	VBP	verb non-3sg pres	<i>eat</i>
JJS	adj., superlative	<i>wildest</i>	VBZ	verb 3sg pres	<i>eats</i>
LS	list item marker	<i>1, 2, One</i>	WDT	wh-determiner	<i>which, that</i>
MD	modal	<i>can, should</i>	WP	wh-pronoun	<i>what, who</i>
NN	noun, sing. or mass	<i>llama</i>	WP\$	possessive wh-	<i>whose</i>
NNS	noun, plural	<i>llamas</i>	WRB	wh-adverb	<i>how, where</i>
NNP	proper noun, sing.	<i>IBM</i>	\$	dollar sign	<i>\$</i>
NNPS	proper noun, plural	<i>Carolinas</i>	#	pound sign	<i>#</i>
PDT	predeterminer	<i>all, both</i>	“	left quote	<i>‘ or “</i>
POS	possessive ending	<i>’s</i>	”	right quote	<i>’ or ”</i>
PRP	personal pronoun	<i>I, you, he</i>	(	left parenthesis	<i>[, (, {, &lt;</i>
PRP\$	possessive pronoun	<i>your, one’s</i>	)	right parenthesis	<i>], ), }, &gt;</i>
RB	adverb	<i>quickly, never</i>	,	comma	<i>,</i>
RBR	adverb, comparative	<i>faster</i>	.	sentence-final punc	<i>. ! ?</i>
RBS	adverb, superlative	<i>fastest</i>	:	mid-sentence punc	<i>: ; ... - -</i>
RP	particle	<i>up, off</i>			

# "Universal Dependencies" Tagset

Nivre et al. 2016

	Tag	Description	Example
Open Class	<b>ADJ</b>	Adjective: noun modifiers describing properties	<i>red, young, awesome</i>
	<b>ADV</b>	Adverb: verb modifiers of time, place, manner	<i>very, slowly, home, yesterday</i>
	<b>NOUN</b>	words for persons, places, things, etc.	<i>algorithm, cat, mango, beauty</i>
	<b>VERB</b>	words for actions and processes	<i>draw, provide, go</i>
	<b>PROPN</b>	Proper noun: name of a person, organization, place, etc..	<i>Regina, IBM, Colorado</i>
	<b>INTJ</b>	Interjection: exclamation, greeting, yes/no response, etc.	<i>oh, um, yes, hello</i>
Closed Class Words	<b>ADP</b>	Adposition (Preposition/Postposition): marks a noun's spacial, temporal, or other relation	<i>in, on, by under</i>
	<b>AUX</b>	Auxiliary: helping verb marking tense, aspect, mood, etc.,	<i>can, may, should, are</i>
	<b>CCONJ</b>	Coordinating Conjunction: joins two phrases/clauses	<i>and, or, but</i>
	<b>DET</b>	Determiner: marks noun phrase properties	<i>a, an, the, this</i>
	<b>NUM</b>	Numeral	<i>one, two, first, second</i>
	<b>PART</b>	Particle: a preposition-like form used together with a verb	<i>up, down, on, off, in, out, at, by</i>
	<b>PRON</b>	Pronoun: a shorthand for referring to an entity or event	<i>she, who, I, others</i>
	<b>SCONJ</b>	Subordinating Conjunction: joins a main clause with a subordinate clause such as a sentential complement	<i>that, which</i>
Other	<b>PUNCT</b>	Punctuation	<i>; , ()</i>
	<b>SYM</b>	Symbols like \$ or emoji	<i>\$, %</i>
	<b>X</b>	Other	<i>asdf, qwfg</i>

# Public tag sets in English

- Brown corpus - Francis and Kucera 1961
  - 500 samples, distributed across 15 genres in rough proportion to the amount published in 1961 in each of those genres
  - 87 tags
- [Penn Treebank](#) - Marcus et al. 1993
  - Hand-annotated corpus of Wall Street Journal, 1M words
  - 45 tags, a simplified version of Brown tag set
  - Standard for English now
    - Most statistical POS taggers are trained on this tagset
- Universal Dependencies (UD) – introduced later

# Example of Penn Treebank Tagging of Brown Corpus Sentence

- The/DT grand/JJ jury/NN commented/VBD on/IN a/DT number/NN of/IN other/JJ topics/NNS ./.

- VB DT NN .  
Book that flight .

- VBZ DT NN VB NN ?  
Does that flight serve dinner ?

# The Problem

- Words often have more than one word class: *this*
  - *This* is a nice day = PRP (personal pronoun)
  - *This* day is nice = DT (determiner)
  - You can go *this* far = RB (adverb)
- *Back*
  - The *back* door (adjective)
  - On my *back* (noun)
  - Promised to *back* the bill (verb)

# Buffalo example

- A grammatically correct (but lexically ambiguous) sentence in American English:  
**Buffalo buffalo Buffalo buffalo buffalo buffalo Buffalo buffalo.**
- [Dmitri Borgmann](#), 1967. [\*Beyond Language: Adventures in Word and Thought\*](#).
- The sentence employs three distinct meanings of the word *buffalo*:
  - as a proper noun to refer to a specific place named Buffalo, the city of [Buffalo, New York](#), being the most notable;
  - as a verb (uncommon in regular usage) *to buffalo*, meaning "to bully, harass, or intimidate" or "to baffle"; and
  - as a noun to refer to the animal, [bison](#) (often called *buffalo* in North America). The plural is also *buffalo*.
- An expanded form of the sentence which preserves the original word order is:  
"Buffalo bison, that other Buffalo bison bully, also bully Buffalo bison."

# How difficult is POS tagging in English?

- Roughly 15% of word types are ambiguous
- Hence 85% of word types are unambiguous
- *Janet* is always PROPN, *hesitantly* is always ADV
- But those 15% tend to be very common.
- So ~60% of word tokens are ambiguous
- E.g., *back*  
earnings growth took a *back*/ADJ seat  
a small building in the *back*/NOUN  
a clear majority of senators *back*/VERB the bill  
enable the country to buy *back*/PART debt  
I was twenty-one *back*/ADV then



# How much ambiguity is there?

- Statistics of word-tag pair in Brown Corpus and Penn Treebank

		87-tag Original Brown		45-tag Treebank Brown	
<b>Unambiguous (1 tag)</b>		<b>44,019</b>		<b>38,857</b>	
<b>Ambiguous (2–7 tags)</b>		<b>5,490</b>	<b>11%</b>	<b>8844</b>	<b>18%</b>
Details:	2 tags	4,967		6,731	
	3 tags	411		1621	
	4 tags	91		357	
	5 tags	17		90	
	6 tags	2	( <i>well, beat</i> )	32	
	7 tags	2	( <i>still, down</i> )	6	( <i>well, set, round, open, fit, down</i> )
	8 tags			4	( <i>'s, half, back, a</i> )
	9 tags			3	( <i>that, more, in</i> )

# POS tagging baselines

- Default classifier:
  - each word is assigned the most probable category,
  - probabilities are computed from manually tagged corpus,
  - in English around 92% classification accuracy
- Human expert accuracy is around 98%

# POS tagging performance in English

- How many tags are correct? (Tag accuracy)
  - About 97%
    - Slight improvement in the last 10+ years
    - HMMs, CRFs, BERT perform similarly .
    - Human accuracy about the same
- But baseline is 92%!
  - Baseline is performance of stupidest possible method
    - "Most frequent class baseline" is an important baseline for many tasks
      - Tag every word with its most frequent tag
      - (and tag unknown words as nouns)
  - Partly easy because
    - Many words are unambiguous

# Is POS tagging a solved problem?

- Baseline
  - Tag every word with its most frequent tag
  - Tag unknown words as nouns
- Accuracy
  - Word level: 90%
  - Sentence level
    - Average English sentence length 14.3 words
    - $0.9^{14.3} = 22\%$

## *Accuracy of better POS Tagger*

- *Word level: 97%*
- *Sentence level:  $0.97^{14.3} = 65\%$*

# Sources of information for POS tagging

Janet *will* back the *bill*

**AUX/NOUN/VERB?**

**NOUN/VERB?**

- Prior probabilities of word/tag
  - "*will*" is usually an AUX
- Identity of neighboring words
  - "*the*" means the next word is probably not a verb
- Morphology and wordshape:
  - Prefixes            *unable*:    *un-* → ADJ
  - Suffixes            *importantly*: *-ly* → ADJ
  - Capitalization    *Janet*:        *CAP* → PROP

# Standard algorithms for POS tagging

- Supervised Machine Learning Algorithms:
- Hidden Markov Models
- Conditional Random Fields (CRF)/ Maximum Entropy Markov Models (MEMM)
- Neural sequence models (RNNs or Transformers)
- Large Language Models (like BERT), finetuned
- All required a hand-labeled training set, all about equal performance (97% on English)
- All make use of information sources we discussed
- Via human created features: HMMs and CRFs
- Via representation learning: Neural LMs

# Classical ML models

- SVM
- Conditional Random Fields (CRF)
- Approach:
  - define a set of useful features
  - train a ML model
- Let us illustrate this approach on Slovene

# Morphosyntactical tagging for Slovene

- Slovene is morphologically rich language
- Large set of tags (1902 tags), why?
- Free word order means that certain taggers do not work well, e.g., HMM
- History of tagging
  - MULTEXT-East
    - Around 100.000 words
    - Very homogenous source, a single novel (George Orwell: 1984)
  - JOS 100k / 1M
    - Around 100.000 / 1.000.000 words
    - More heterogeneous
    - Manually labelled 100k corpus / corpus of 1M words partially manually labelled (estimate: 96% accurate tags)
    - Based on FidaPLUS corpus containing 620 million words



# Current Slovene POS datasets

- ssj500k
  - 600k words manually labelled corpus
  - Analysis of common errors (mostly due to underrepresentation of certain tags in the corpus), e.g., je
- SUK (2023)
  - superset of ssj500k
  - 1M words
  - seem to be sufficient for standard language
  - planned extensions for non-standard language domains
  - <https://www.clarin.si/repository/xmlui/handle/11356/1747#>

# An example in Slovene

- JOS ToTaLe text analyzer for Slovene: morphosyntactical tagging, (old variant available at <http://www.slovenscina.eu/>)

*Nekega dne sem se napotil v naravo. Že spočetka me je žulil čevelj, a sem na to povsem pozabil, ko sem jo zagledal. Bila je prelepa. Povsem nezakrita se je sončila na trati ob poti. Pritisk se mi je dvignil v višave. Popoln primerek kmečke lastovke!*

- Tags are standardized for East European languages in Multext-East specification, e.g.,

dne; tag Somer = Samostalnik, obče ime, moški spol, ednina, rodilnik; lema: dan

- *Nekega dne sem se napotil v naravo. Že spočetka me je žulil čevelj, a sem na to povsem pozabil, ko sem jo zagledal. Bila je prelepa. Povsem nezakrita se je sončila na trati ob poti. Pritisk se mi je dvignil v višave. Popoln primerek kmečke lastovke!*

1	<b>beseda</b> <b>lema</b> <b>oznaka</b>	Nekega dne sem se napotil v naravo . Že spočetka me je nek dan biti se napotiti v narava že spočetka jaz biti Zn-mer Somer Gp-spe-n Zp-----k Ggdd-em Dt Sozet . L Rsn Zop-et--k Gp-ste-n
2	<b>beseda</b> <b>lema</b> <b>oznaka</b>	žulil čevelj , a sem na to povsem pozabil , ko sem jo zagledal žuliti čevelj a biti na ta povsem pozabiti ko biti on zagledati Ggnd-em Somei , Vp Gp-spe-n Dt Zk-set Rsn Ggdd-em , Vd Gp-spe-n Zotzet--k Ggdd-em
3	<b>beseda</b> <b>lema</b> <b>oznaka</b>	. Bila je prelepa . Povsem nezakrita se je sončila na trati biti biti prelep povsem nezakrit se biti sončiti na trata . Gp-d-ez Gp-ste-n Ppnzei . Rsn Ppnzei Zp-----k Gp-ste-n Ggvd-ez Dm Sozem
4	<b>beseda</b> <b>lema</b> <b>oznaka</b>	ob poti . Pritisk se mi je dvignil v višave . Popoln ob pot pritisk se jaz biti dvigniti v višava popoln Dm Sozem . Somei Zp-----k Zop-ed--k Gp-ste-n Ggdd-em Dt Sozmt . Ppnmein
5	<b>beseda</b> <b>lema</b> <b>oznaka</b>	primerek kmečke lastovke ! primerek kmečki lastovka Somei Ppnzer Sozer !

# TEI-XML format

```
<TEI xmlns="http://www.tei-c.org/ns/1.0">
  <text>
    <body>
      <p>
        <s>
          <w msd="Zn-mer" lemma="nek">Nekega</w>
          <S/>
          <w msd="Somer" lemma="dan">dne</w>
          <S/>
          <w msd="Gp-spe-n" lemma="biti">sem</w>
          <S/>
          <w msd="Zp-----k" lemma="se">se</w>
          <S/>
          <w msd="Ggdd-em" lemma="napotiti">napotil</w>
          <S/>
          <w msd="Dt" lemma="v">v</w>
          <S/>
          <w msd="Sozet" lemma="narava">naravo</w>
          <c>.</c>
          <S/>
        </s>
        ...
      </p>
    </body>
  </text>
</TEI>
```

# MSD tags for Slovene

- Multext-East 4.0 specification
- example: dne;  
tag Somer = Samostalni<sup>k</sup>,  
obče ime, moški spol, ednina,  
rodilnik; lema: dan
- below top level tags there  
are many informative  
features
- example for verb

P	atribut	vrednost	koda	atribut	vrednost	koda
0	glagol		G	Verb		V
1	vrsta	glavni	g	Type	main	m
		pomožni	p		auxiliary	a
2	vid	dovršni	d	Aspect	perfective	e
		nedovršni	n		imperfective	p
		dvovidski	v		biaspectual	b
3	oblika	nedoločnik	n	VForm	infinitive	n
		namenilnik	m		supine	u
		deležnik	d		participle	p
		sedanjik	s		present	r
		prihodnjik	p		future	f
		pogojnik	g		conditional	c
		velelnik	v		imperative	m
4	oseba	prva	p	Person	first	1
		druga	d		second	2
		tretja	t		third	3
5	število	ednina	e	Number	singular	s
		množina	m		plural	p
		čvojina	d		dual	d
6	spol	moški	m	Gender	masculine	m
		ženski	z		feminine	f
		srednji	s		neuter	n
7	nikalnost	nezanikani	n	Negative	no	n
		zanikani	d		yes	y

# Parsing: finding linguistic structure

1. Constituency parsing
2. Dependency parsing

# Parsing reduces ambiguity

Scientists count whales from space



Scientists count whales from space



# Constituency parsing

- Dependency structure shows which words depend on (modify or are arguments of) which other words.
- *Look in the large crate in the kitchen by the door*
- We need to understand sentence structure in order to be able to interpret language correctly
- Humans communicate complex ideas by composing words together into bigger units to convey complex meanings
- We need to know what is connected to what



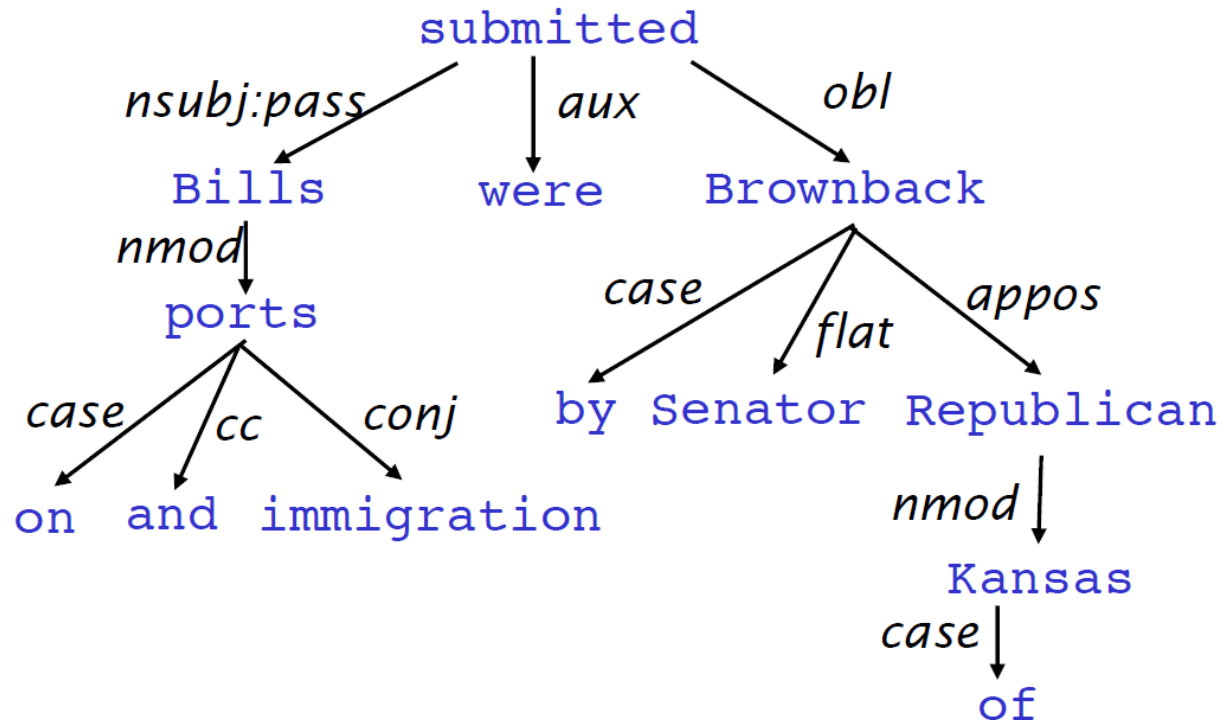
# Constituency parsing

- Phrase structure organizes words into nested constituents
- Starting unit: words are given a category (part of speech = pos)  
the, cat, cuddly, by, door
- Words combine into phrases with categories  
the cuddly cat, by the door
- Phrases can combine into bigger phrases recursively  
the cuddly cat by the door  
Det Adj N P Det N
- Words combine into phrases with categories  
the cuddly cat, by the door  
NP → Det Adj N NP → Det N PP → P NP
- Phrases can combine into bigger phrases recursively  
the cuddly cat by the door NP → NP PP

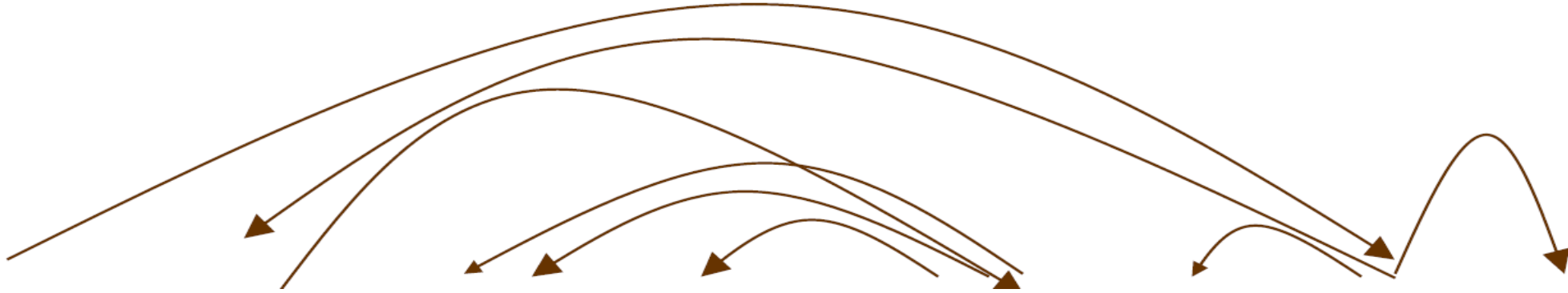
# Dependency parsing

- Dependency syntax postulates that syntactic structure consists of relations between lexical items, normally binary asymmetric relations (“arrows”) called dependencies

The arrows are commonly **typed** with the name of grammatical relations (subject, prepositional object, apposition, etc.)



# Dependency Grammar and Dependency Structure



ROOT Discussion of the outstanding issues was completed .

- Some people draw the arrows one way; some the other way!
- Usually add a fake ROOT so every word is a dependent of precisely 1 other node

# Advantages of dependency parsing

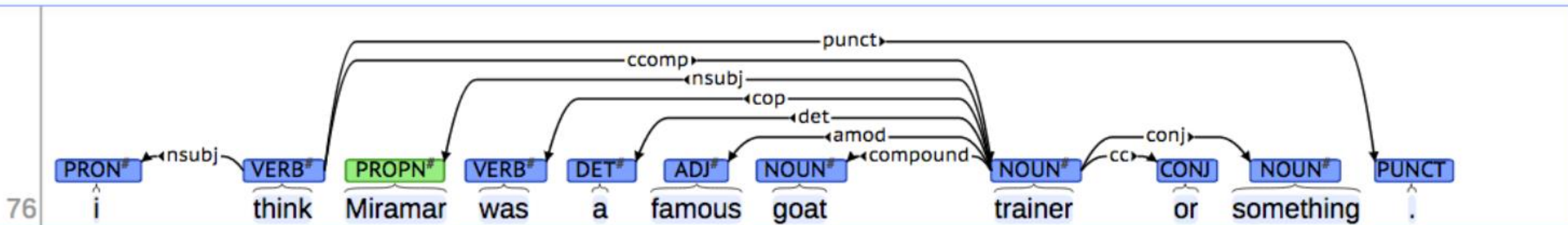
- Better handling of free word order (less-Anglo-centric)
- Node simplicity
- Clean mapping to semantic predicate-argument structure
- Easier to develop multilingual systems

# Role of dependency parsing in NLP

- Semantic role labeling
- Relation extraction,
- Machine translation,
- Helps in explanation
- Important role in the linguistic analysis

# Treebanks

- The rise of annotated data: Universal Dependencies treebanks
- <http://universaldependencies.org/>
- Earlier: Marcus et al. 1993, The Penn Treebank, *Computational Linguistics*

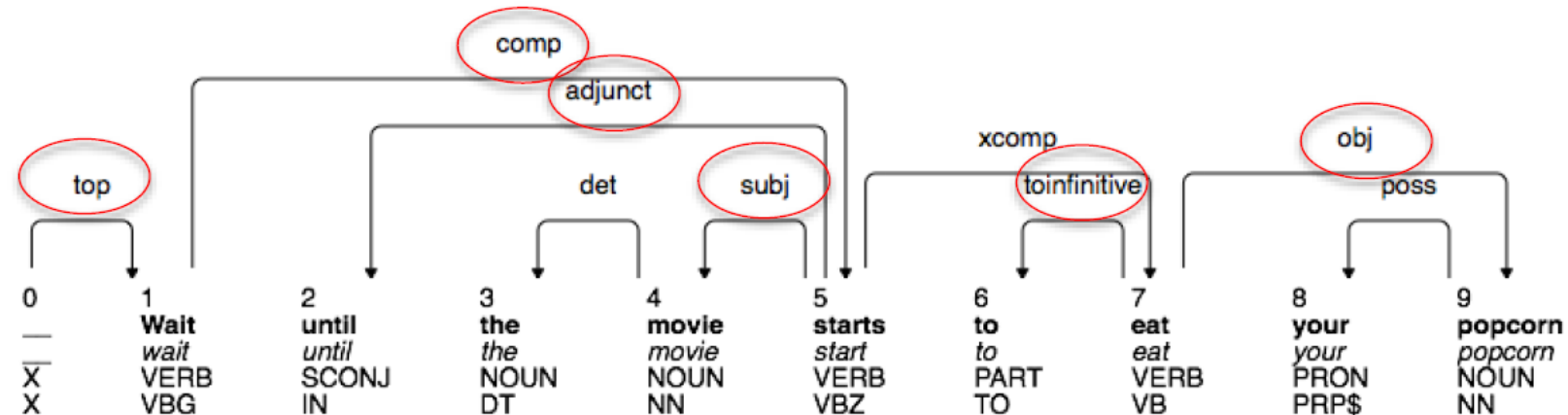
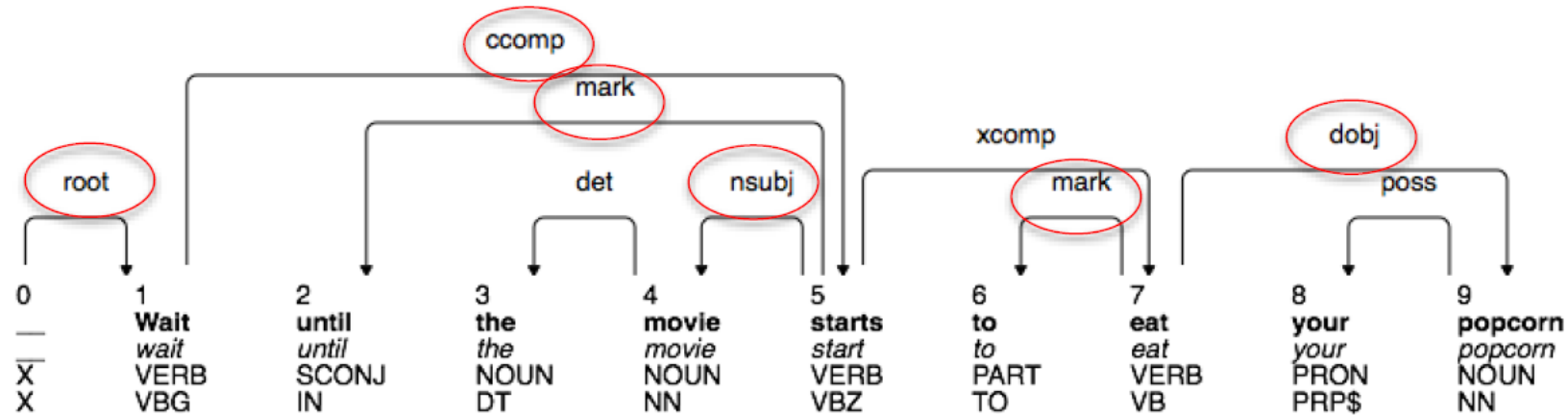


# Treebank

- Collection of parsed sentences (trees)
- Annotated with a pre-defined part-of-speech tagset (Noun, Verb, etc.)
- Pre-defined annotation scheme (list of prescribed labels)
- Pre-defined linguistic structure
- Used to develop statistical parsers (train, test, and bootstrap)

# Variation in labelling

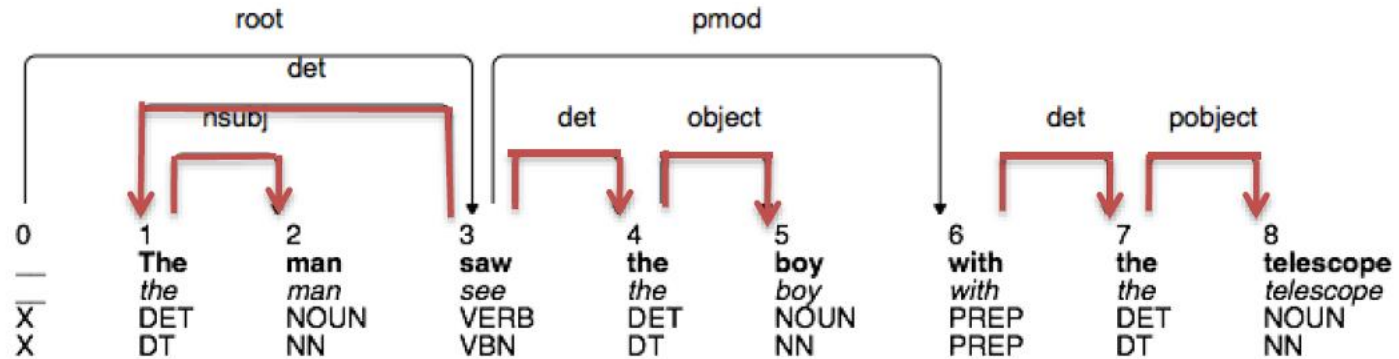
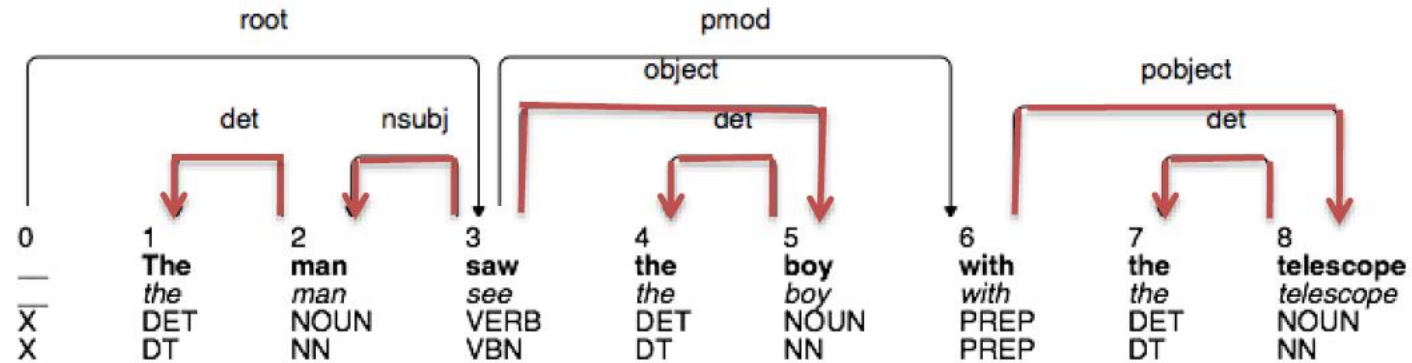
Varying labelling conventions:





# Variation in structure

Varying structural analyses:

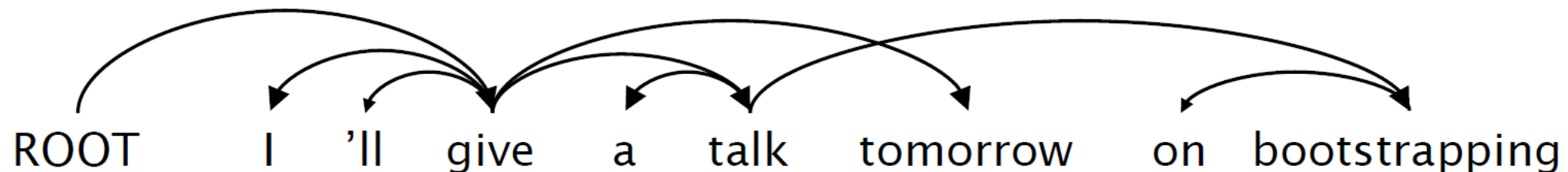


# Building treebank

- Building a treebank seems a lot slower and less useful than building a grammar
- But a treebank gives us many things
  - Reusability of the labor
    - Many parsers, part-of-speech taggers, etc. can be built on it
    - Valuable resource for linguistics
  - Broad coverage, not just a few intuitions
  - Frequencies and distributional information
  - A way to evaluate systems

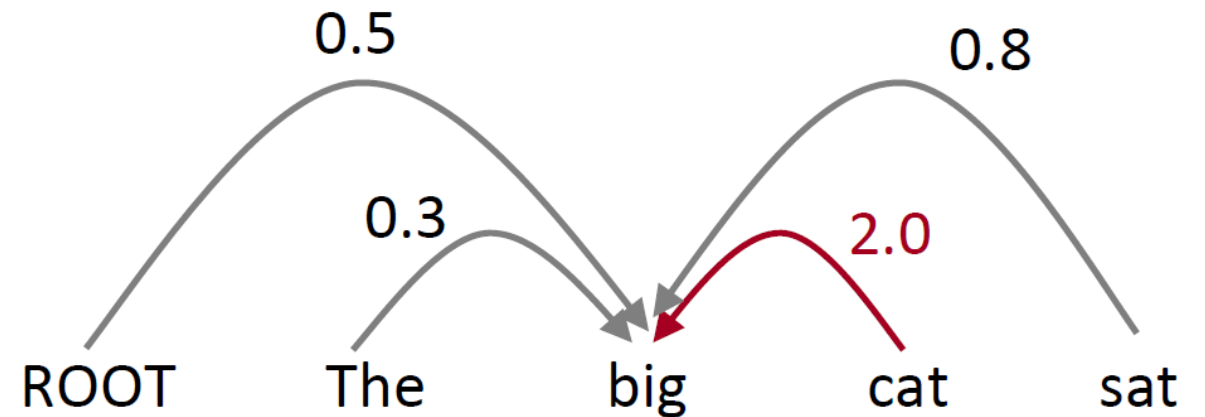
# Dependency parsing

- A sentence is parsed by choosing for each word what other word (including ROOT) is it a dependent of
- Usually some constraints:
  - Only one word is a dependent of ROOT
  - Don't want cycles  $A \rightarrow B, B \rightarrow A$
  - This makes the dependencies a tree
  - Final issue is whether arrows can cross (non-projective) or not



# Graph-based dependency parsers

- Compute a score for every possible dependency for each word
- Then add an edge from each word to its highest-scoring candidate head
- And repeat the same process for each other word
- E.g., picking the head for “big”



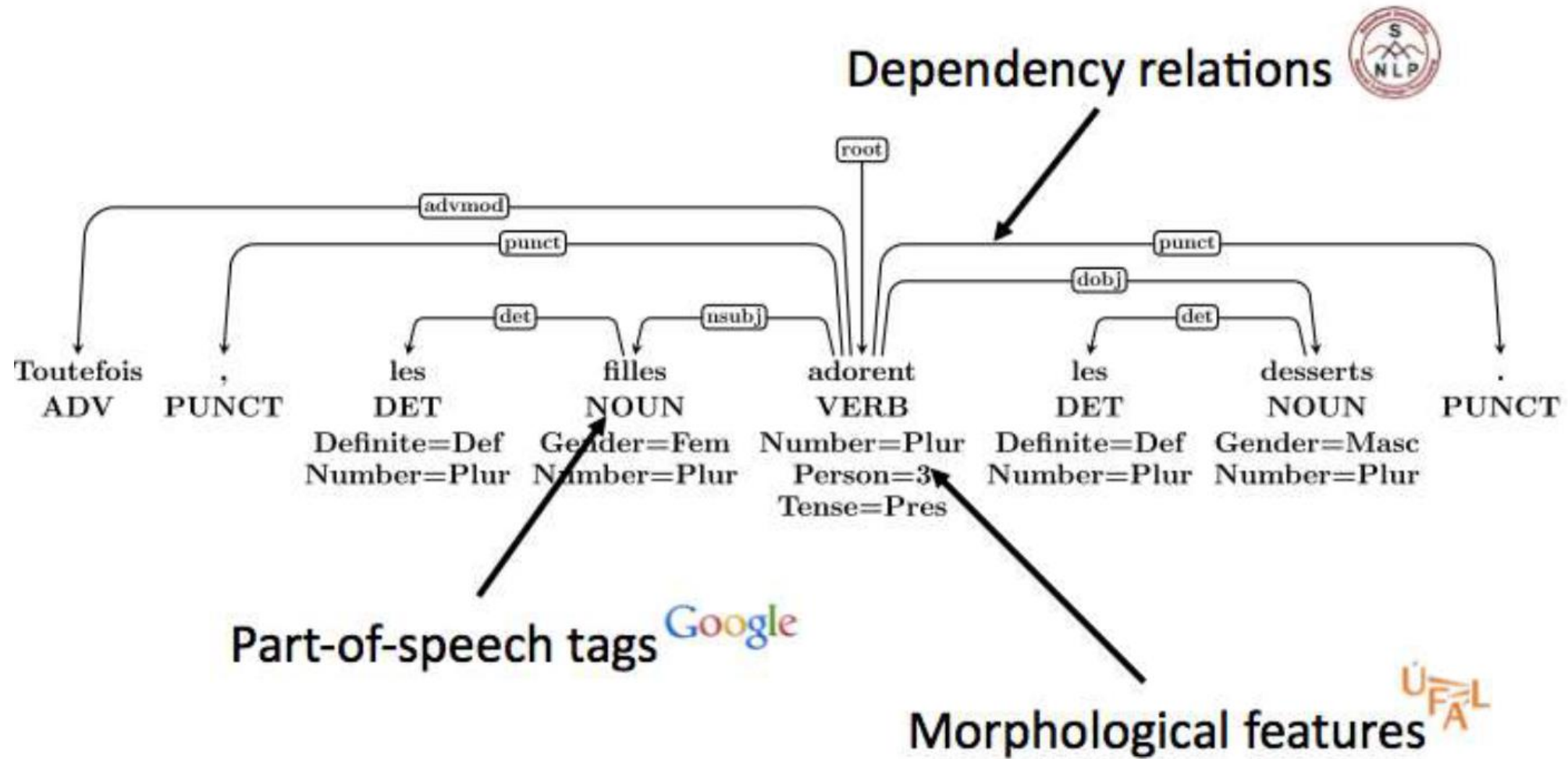
# Variation between languages

- **Problems** with variations
- Difficult to do cross-lingual analysis
- Difficult to compare parser performance
- Difficult to do cross-lingual transfer (using data from one language to help another)
- Difficult to build and evaluate multilingual systems

# Solution: Universal Dependencies

- **Premise:**
  - no Universal Grammar, but:
  - “all languages share fundamental similarities” (linguistic universals)
- **Goals:**
  - develop a set of harmonized dependency treebanks
  - design a universal annotation scheme
  - enable comparison of treebanks
  - enable comparison of parsing results
  - improve multilingual processing

# UD project



# UD POS tags

- Taxonomy of 17 universal part-of-speech tags, expanding on the Google Universal Tagset (Petrov et al., 2012)
- All languages use the same inventory, but not all tags have to be used by all languages

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

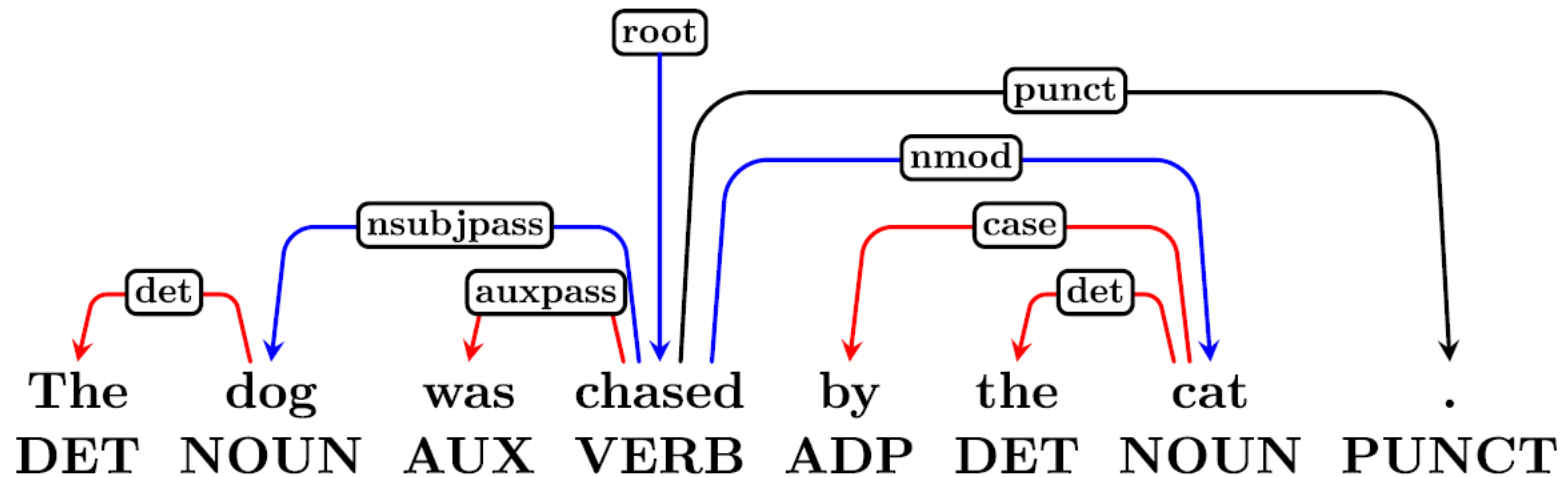


# Slovene UD POS tags

- ADJ: adjective
- ADP: adposition
- ADV: adverb
- AUX: auxiliary verb
- CONJ: coordinating conjunction
- DET: determiner
- INTJ: interjection
- NOUN: noun
- NUM: numeral
- PART: particle
- PRON: pronoun
- PROPN: proper noun
- PUNCT: punctuation
- SCONJ: subordinating conjunction
- SYM: symbol
- VERB: verb
- X: other

# UD syntax

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



# UD relations

- 40 universal grammatical relations (de Marneffe et al., 2014) (aim to address linguistic universals across languages)
- Language-specific subtypes may be added

# UD Features

- Standardized inventory of morphological features, based on the Intersect system (Zeman, 2008)
- Languages select relevant features and can add language-specific features or values with documentation

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

# Slovene UD features

- **POS Tags**

[ADJ](#) – [ADP](#) – [ADV](#) – [AUX](#) – [CCONJ](#) – [DET](#) – [INTJ](#) – [NOUN](#) – [NUM](#) – [PART](#) – [PRON](#) – [PROPN](#) – [PUNCT](#) – [SCONJ](#) – [VERB](#) – [X](#)

- **Features**

[Animacy](#) – [Aspect](#) – [Case](#) – [Definite](#) – [Degree](#) – [Foreign](#) – [Gender](#) – [Gender\[psor\]](#) – [Mood](#) – [Number](#) – [Number\[psor\]](#) – [NumForm](#) – [NumType](#) – [Person](#) – [Polarity](#) – [Poss](#) – [PronType](#) – [Tense](#) – [Variant](#) – [VerbForm](#)

- **Relations**

[acl](#) – [advcl](#) – [advmod](#) – [amod](#) – [appos](#) – [aux](#) – [case](#) – [cc](#) – [cc:preconj](#) – [ccomp](#) – [conj](#) – [conj:extend](#) – [cop](#) – [csubj](#) – [dep](#) – [det](#) – [discourse](#) – [discourse:filler](#) – [dislocated](#) – [expl](#) – [fixed](#) – [flat](#) – [flat:foreign](#) – [flat:name](#) – [goeswith](#) – [iobj](#) – [mark](#) – [nmod](#) – [nsubj](#) – [nummod](#) – [obj](#) – [obl](#) – [orphan](#) – [parataxis](#) – [parataxis:discourse](#) – [parataxis:restart](#) – [punct](#) – [reparandum](#) – [root](#) – [vocative](#) – [xcomp](#)

- [https://universaldependencies.org/treebanks/sl\\_sst/index.html](https://universaldependencies.org/treebanks/sl_sst/index.html)

# Modern POS and dependency parsing pipelines

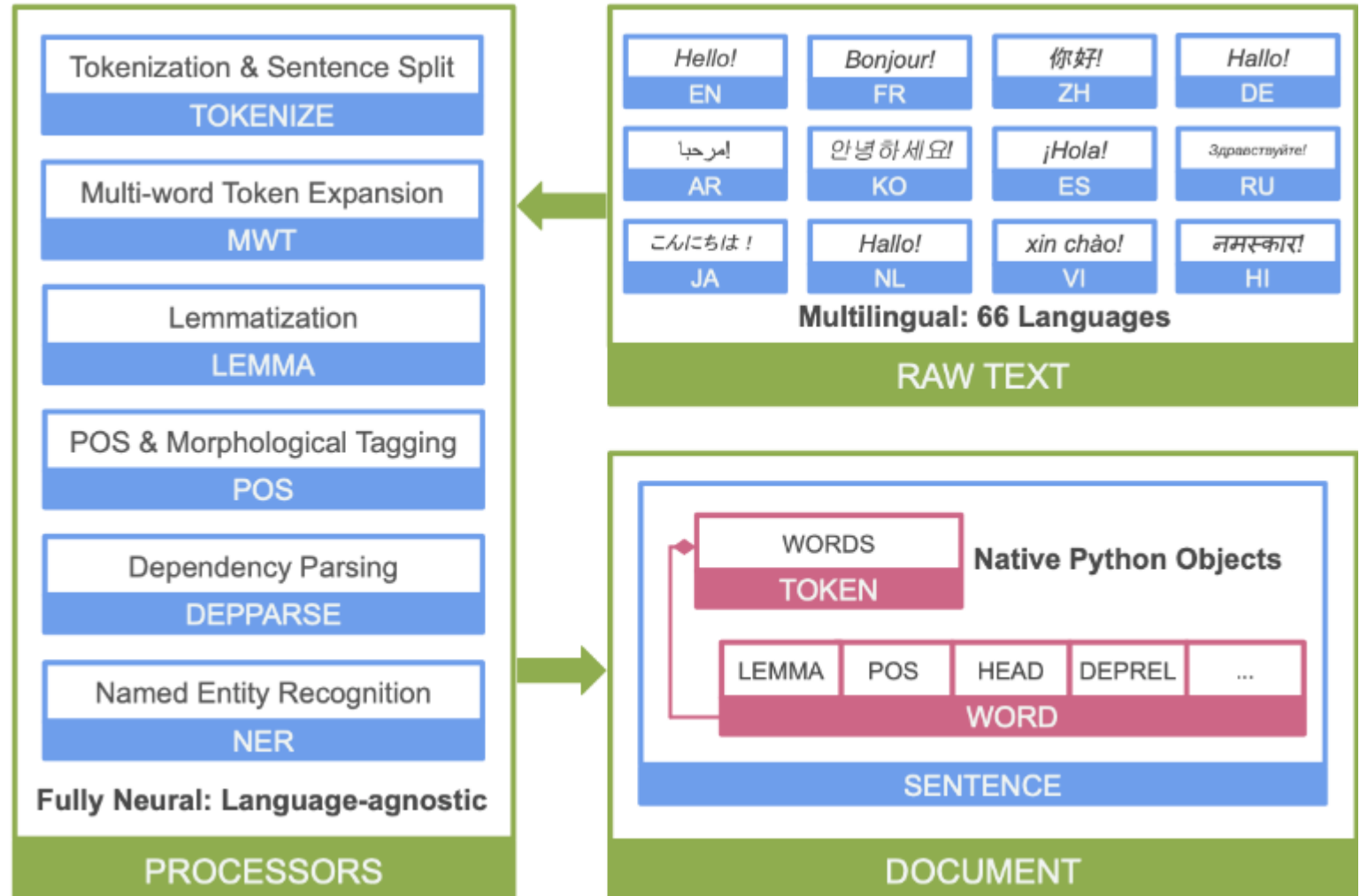
- A single neural pipeline for all bottom layer tasks
- Tokenization, sentence and word segmentation, part-of-speech (POS)/morphological features (UFeats) tagging, lemmatization, dependency parsing, and named entity recognition (NER)
- Predominant approach for many languages

Qi, P., Dozat, T., Zhang, Y. and Manning, C.D., 2018, October. Universal Dependency Parsing from Scratch. In *Proceedings of the CoNLL 2018 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies* (pp. 160-170).

Kaja Dobrovoljc, Luka Krsnik, Marko Robnik-Šikonja. 2023. STARK: A Tool for Dependency Tree Extraction and Analysis. UniDive 2023

# Stanford Stanza pipeline

- <https://stanfordnlp.github.io/stanza/>
- Given a document of raw text,
- The tokenizer/sentence segmenter/MWT expander splits it into sentences of syntactic words;
- The tagger assigns UPOS, XPOS and UFeat tags to each word;
- The lemmatizer takes the predicted word and UPOS tag and outputs a lemma;
- The parser takes all annotations as input and predicts the head and dependency label for each word
- NER is added into the pipeline



# Quality of tools in Slovene

tool	distributional information	Slovenian	Croatian	Serbian
reldi-tagger	Brown clusters	94.21	91.91	92.03
stanfordnlp	CoNLL w2v embeddings	96.45	93.85	94.78
stanfordnlp	CLARIN.SI w2v embeddings	<b>96.79</b>	<b>94.18</b>	94.91
stanfordnlp	CLARIN.SI fT embeddings	96.72	94.13	<b>95.23</b>

Table 1: F1 results in morphosyntactic annotation with the traditional and neural tool and different distributional information.

tool	morphosyntax	Slovenian	Croatian	Serbian
reldi-tagger	gold	99.46	98.17	97.89
reldi-tagger	reldi-tagger	98.35	96.82	96.44
reldi-tagger	stanfordnlp	98.77	97.22	97.26
stanfordnlp	gold	97.75	96.22	95.29
stanfordnlp	stanfordnlp	97.51	95.85	95.18
stanfordnlp+lex	gold	99.30	98.11	97.78
stanfordnlp+lex	stanfordnlp	98.74	97.22	97.13

Table 3: F1 results in lemmatisation with the traditional and neural tool and different upstream processing.

Ljubešić, N. and Dobrovoljc, K., 2019. What does Neural Bring? Analysing Improvements in Morphosyntactic Annotation and Lemmatisation of Slovenian, Croatian and Serbian. In *Proceedings of the 7th Workshop on Balto-Slavic Natural Language Processing* (pp. 29-34).



# Slovene Classla (Stanza) pipeline results

## 10 January 2022

METRIC	PRECISION	RECALL	F1 SCORE	ALIGNDACC
TOKENS	99.97	99.95	99.96	
SENTENCES	99.58	99.47	99.52	
WORDS	99.97	99.95	99.96	
UPOS	98.70	98.69	98.69	98.73
XPOS	97.39	97.37	97.38	97.42
UFEATS	97.01	96.99	97.00	97.04
ALLTAGS	96.33	96.31	96.32	96.36
LEMMAS	99.17	99.16	99.17	99.20
UAS	94.06	94.04	94.05	94.08
LAS	92.05	92.04	92.05	92.08
CLAS	89.34	90.04	89.69	90.09
MLAS	85.08	85.76	85.42	85.80
BLEX	88.75	89.45	89.10	89.50

# Broader POS-tagging comparison for Slovene: CoLLU Shared task 2018

Modeli	Tokeni	Stavki	UPOS	XPOS	Lema	UAS	LAS	Avg
SpaCy (brez vektorjev)	99,29	97,60	96,42	89,91	94,25	85,54	77,82	90,26
SpaCy (fastText.cc)	99,29	<b>98,80</b>	96,15	89,75	94,26	85,76	78,07	<b>90,47</b>
SpaCy (Clarin)	99,29	97,79	96,20	89,42	93,97	85,80	78,19	90,23
SpaCy (cbow – navadni)	99,29	98,37	96,18	89,30	93,94	85,55	77,87	90,20
SpaCy (cbow – floret)	99,29	97,76	96,25	89,55	94,16	85,44	77,87	90,17
SpaCy (skipgram – navadni)	99,29	98,26	96,22	89,47	93,80	85,38	77,55	90,11
SpaCy (skipgram – floret)	99,29	98,06	96,22	89,66	94,08	85,54	77,93	90,25
SpaCy (SloBERTa 2.0) [t]	99,29	97,79	<b>98,39</b>	<b>97,35</b>	<b>96,96</b>	<b>93,95</b>	<b>87,98</b>	<b>95,40</b>
CLASSLA (stand.) [3]	99,92	99,57	98,69	97,81	<b>99,20</b>	92,68	90,87	96,47
Stanza [4]	99,90	98,10	98,33	95,13	97,07	92,72	90,97	95,39
Trankit (large) [24] [t]	<b>99,97</b>	<b>100</b>	<b>99,24</b>	<b>97,83</b>	97,55	<b>96,91</b>	<b>96,06</b>	<b>97,93</b>
Trankit (base) [t]	99,93	99,81	99,03	96,70	97,49	95,94	94,99	97,33
UDPIPE 2.10 [39] [t]	98,95	99,94	98,97	96,97	98,58	93,99	92,60	96,84

# Broader DP comparison for Slovene: CoLLU Shared task 2018

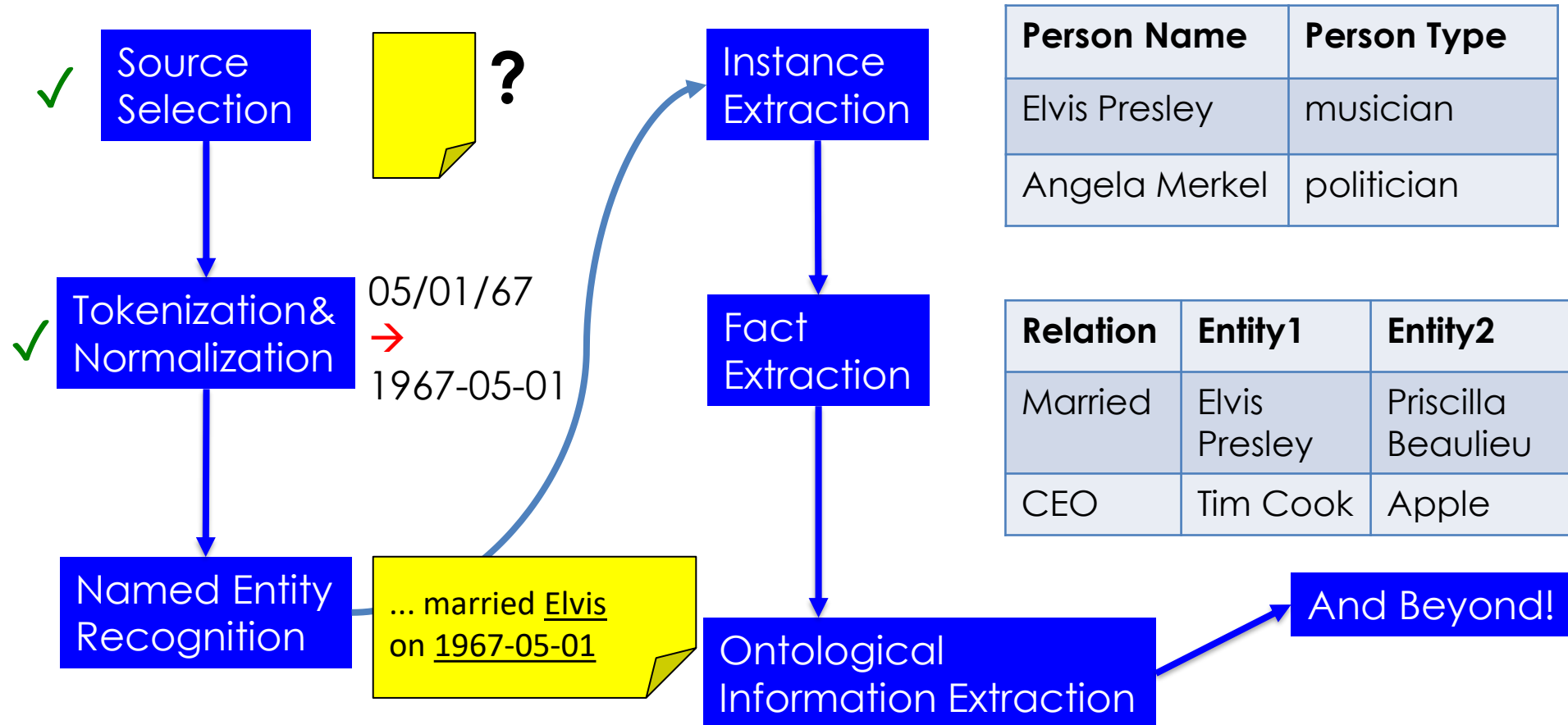
Modeli (ssj500k)	LAS	MLAS	BLEX
SpaCy (brez vektorjev)	73,99	64,22	68,96
SpaCy (fastText.cc)	74,15	64,20	<b>69,09</b>
SpaCy (Clarin)	<b>74,36</b>	64,13	68,96
SpaCy (cbow – navadni)	73,77	63,29	68,44
SpaCy (cbow – floret)	74,01	63,88	68,83
SpaCy (skipgram – navadni)	73,50	63,15	67,95
SpaCy (skipgram – floret)	74,12	<b>64,28</b>	68,95
SpaCy (SloBERTa 2.0)	<b>86,46</b>	<b>84,08</b>	<b>83,75</b>
CLASSLA (stand.)	88,60	85,38	87,92
Stanza	88,37	83,21	84,98
Trankit (large)	<b>94,88</b>	<b>91,78</b>	<b>91,37</b>
Trankit (base)	93,53	89,09	90,12
UDPIPE 2.10	x	86,83	88,91

# Named entity recognition (NER)

- Recently, NER was added to the the basic linguistic annotation pipeline
- Why?

# Information Extraction

**Information Extraction** (IE) is the process of extracting **structured information** from unstructured machine-readable documents



# Relation Extraction: Disease Outbreaks

May 19 1995, Atlanta -- The Centers for Disease Control and Prevention, which is in the front line of the world's response to the deadly Ebola epidemic in Zaire , is finding itself hard pressed to cope with the crisis...

**Information  
Extraction System**

<i>Date</i>	<i>Disease Name</i>	<i>Location</i>
Jan. 1995	Malaria	Ethiopia
July 1995	Mad Cow Disease	U.K.
Feb. 1995	Pneumonia	U.S.

# Named entity recognition

- A **named entity** is anything that can be referred to with a **proper name**:
  - a person, a location, an organization.
- **Named entity recognition** (NER) aims to find spans of text that constitute proper names and tag the type of NER entity.
- Four common entity tags:
  - **PER** (person), **LOC** (location), **ORG** (organization), or **GPE** (geo-political entity), **OTHER** (everything else)
- Commonly extended to dates, times, other temporal expressions, numerical expressions like prices.
- Also events, movie and book names, etc.

Type	Tag	Sample Categories	Example sentences
People	PER	people, characters	<b>Turing</b> is a giant of computer science.
Organization	ORG	companies, sports teams	The <b>IPCC</b> warned about the cyclone.
Location	LOC	regions, mountains, seas	<b>Mt. Sanitas</b> is in <b>Sunshine Canyon</b> .
Geo-Political Entity	GPE	countries, states	<b>Palo Alto</b> is raising the fees for parking.

# NER output

Citing high fuel prices, [ORG **United Airlines**] said [TIME **Friday**] it has increased fares by [MONEY **\$6**] per round trip on flights to some cities also served by lower-cost carriers. [ORG **American Airlines**], a unit of [ORG **AMR Corp.**], immediately matched the move, spokesman [PER **Tim Wagner**] said. [ORG **United**], a unit of [ORG **UAL Corp.**], said the increase took effect [TIME **Thursday**] and applies to most routes where it competes against discount carriers, such as [LOC **Chicago**] to [LOC **Dallas**] and [LOC **Denver**] to [LOC **San Francisco**].



# NER usefulness

- A useful first stage in question answering,
- Linking text to information in structured knowledge sources like Wikipedia.
- Natural language understanding
- Building semantic representations, like extracting events and the relationship between participants.

# NER problems

- Ambiguity

[PER Washington] was born into slavery on the farm of James Burroughs.  
[ORG Washington] went up 2 games to 1 in the four-game series.  
Blair arrived in [LOC Washington] for what may well be his last state visit.  
In June, [GPE Washington] passed a primary seatbelt law.

- Conceptual dilemmas:  
Republicans were angry because of the reform.  
PER (people of that conviction) or PER (members of that party) or ORG (Republican party) – shall all be labelled at all?
- More complications and ambiguities if PRODUCT is added as a category,
- e.g., Economist (as a physical newspaper or an organization)

# BIO Tagging

- How can we turn this structured problem into a sequence problem like POS tagging, with one label per word?
- [PER Jane Villanueva] of [ORG United] , a unit of [ORG United Airlines Holding] , said the fare applies to the [LOC Chicago ] route.

# BIO Tagging

- [PER Jane Villanueva] of [ORG United] , a unit of [ORG United Airlines Holding] , said the fare applies to the [LOC Chicago ] route.

Words	BIO Label
Jane	B-PER
Villanueva	I-PER
of	O
United	B-ORG
Airlines	I-ORG
Holding	I-ORG
discussed	O
the	O
Chicago	B-LOC
route	O
.	O

Now we have one tag per token!!!

# BIO Tagging

- B: token that *begins* a span
- I: tokens *inside* a span
- O: tokens outside of any span
- # of tags (where  $n$  is #entity types):
  - 1 O tag,
  - $n$  B tags,
  - $n$  I tags
  - total of  $2n+1$

Words	BIO Label
Jane	B-PER
Villanueva	I-PER
of	O
United	B-ORG
Airlines	I-ORG
Holding	I-ORG
discussed	O
the	O
Chicago	B-LOC
route	O
.	O

# NER is a sequence tagging task

- IO, BIO, and BIOES tagging
- for  $n$  different tags, the number of labels is:  $IO=n+1$   $BIO=2n+1$   $BIOES=4n+1$

[**PER Jane Villanueva**] of [**ORG United**], a unit of [**ORG United Airlines Holding**], said the fare applies to the [**LOC Chicago**] route.

Words	IO Label	BIO Label	BIOES Label
Jane	I-PER	B-PER	B-PER
Villanueva	I-PER	I-PER	E-PER
of	O	O	O
United	I-ORG	B-ORG	B-ORG
Airlines	I-ORG	I-ORG	I-ORG
Holding	I-ORG	I-ORG	E-ORG
discussed	O	O	O
the	O	O	O
Chicago	I-LOC	B-LOC	S-LOC
route	O	O	O
.	O	O	O

# Standard algorithms for NER

- Supervised Machine Learning given a human-labeled training set of text annotated with tags
  - Hidden Markov Models
  - Conditional Random Fields (CRF)/ Maximum Entropy Markov Models (MEMM)
  - Neural sequence models (RNNs or Transformers)
  - Large Language Models (like BERT), finetuned

# NER Evaluation

Comparison to the **gold standard** (i.e. manually labelled or checked output).

Algorithm output:

O = {Einstein, Bohr, Planck, Clinton, Obama}

✓        ✓        ✓        ✗        ✗

Gold standard:

G = {Einstein, Bohr, Planck, Heisenberg}

✓        ✓        ✓        ✗

Precision:

What proportion of the  
output is correct?

$$\frac{|O \cap G|}{|O|}$$

Recall:

What proportion of the  
gold standard did we get?

$$\frac{|O \cap G|}{|G|}$$



# Performance measures

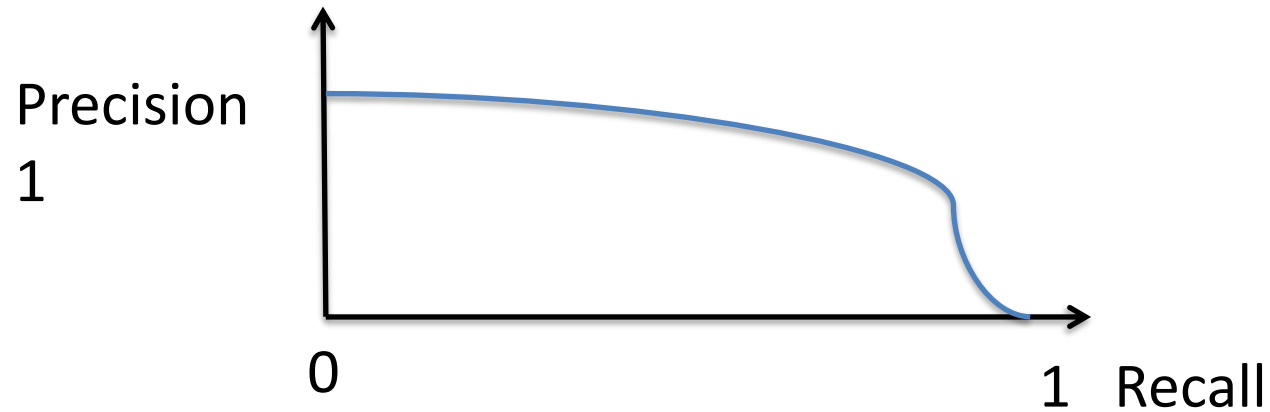
- A contingency table for the analysis of precision and recall

	Relevant	Non-relevant	
Retrieved	$a$	$b$	$a + b = m$
Not retrieved	$c$	$d$	$c + d = N - m$
	$a + c = n$	$b + d = N - n$	$a + b + c + d = N$

- $N$  = number of all tokens in the dataset
- $n$  = number of relevant tags
- $m$  = number of retrieved tags
- the system returns  $m$  tags including  $a$  relevant ones
- Precision  $P = a/m$   
proportion of relevant tags in the returned ones
- recall  $R = a/n$   
proportion of relevant tags in all relevant tags

# F1- Measure

You can't get it all...



The F1-measure combines precision and recall as the harmonic mean:

$$F1 = 2 * \text{precision} * \text{recall} / (\text{precision} + \text{recall})$$

# NER evaluation dilemmas

- How to treat partial matches?
  - entity may be composed of more than one labelled token
  - training loss (tag based) might not be the same as the test loss (entity based)
- Precision and recall assume two class problems, NER has several tags (at least four)
- The F1 score have to be adapted (micro and macro average variant)
- Micro-average F1: you sum up the individual true positives, false positives, and false negatives of the system for different sets and average them
  - compute several one-versus-all scores and average
  - assumes all instances are equally important
  - works well in balanced class case
- Macro-average F1: just take the average of the precision and recall of the system on different set
  - computes TP, FP, TN, FN for each class separately and then compute the measure
  - assumes all classes are equally important
  - works better in imbalanced class case
  - The *Other* tag is often ignored

# Micro and macro averaging example

- Let us compute precision  $P = TP / (TP + FP)$ .
- Let us assume multi-class classification system with four classes and the following numbers when tested:
- Class A: 1 TP and 1 FP
- Class B: 10 TP and 90 FP
- Class C: 1 TP and 1 FP
- Class D: 1 TP and 1 FP
- $P(A) = P(C) = P(D) = 0.5$ , whereas  $P(B) = 0.1$ .
- A macro-averaged precision:  $P_{macro} = (0.5 + 0.1 + 0.5 + 0.5) / 4 = 0.4$
- A micro-averaged precision:  $P_{micro} = (1 + 10 + 1 + 1) / (2 + 100 + 2 + 2) = 0.123$