

Syntactic Parsing

Human Language from a Computational Perspective
April 25, 2018

Language models reminder

Use an n -gram model to predict the next token:

* ~~MY ONLY WISH I~~

MY ONLY **WISH IS**

Bigram counts

(starting with `WISH`):

WISH I	8
WISH IS	6
WISH THEY	4
WISH WAS	4
WISH THAT	2
WISH YOU	1

Lexical ambiguity

The word *WISH* is ambiguous

WISH (verb): **לאחל, לבקש**

WISH (noun): **משאלה**

Some context helps

Verb:

How I **WISH** YOU WERE HERE
CAREFUL WHAT YOU **WISH** FOR
WISH YOU A HAPPY BIRTHDAY

Noun:

YOUR **WISH** IS MY COMMAND
IF YOU COULD HAVE ONE **WISH**
MAKE A **WISH**

But sometimes it doesn't

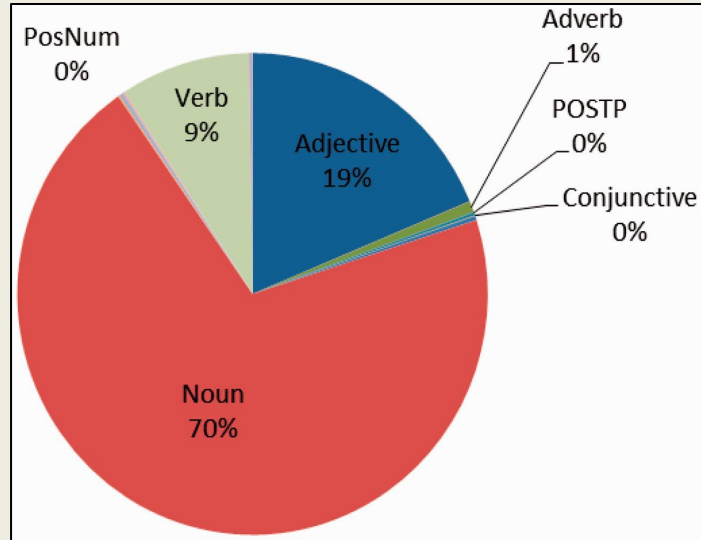
SQUAD HELPS DOG **BITE** VICTIM

EYE **DROPS** OFF SHELF

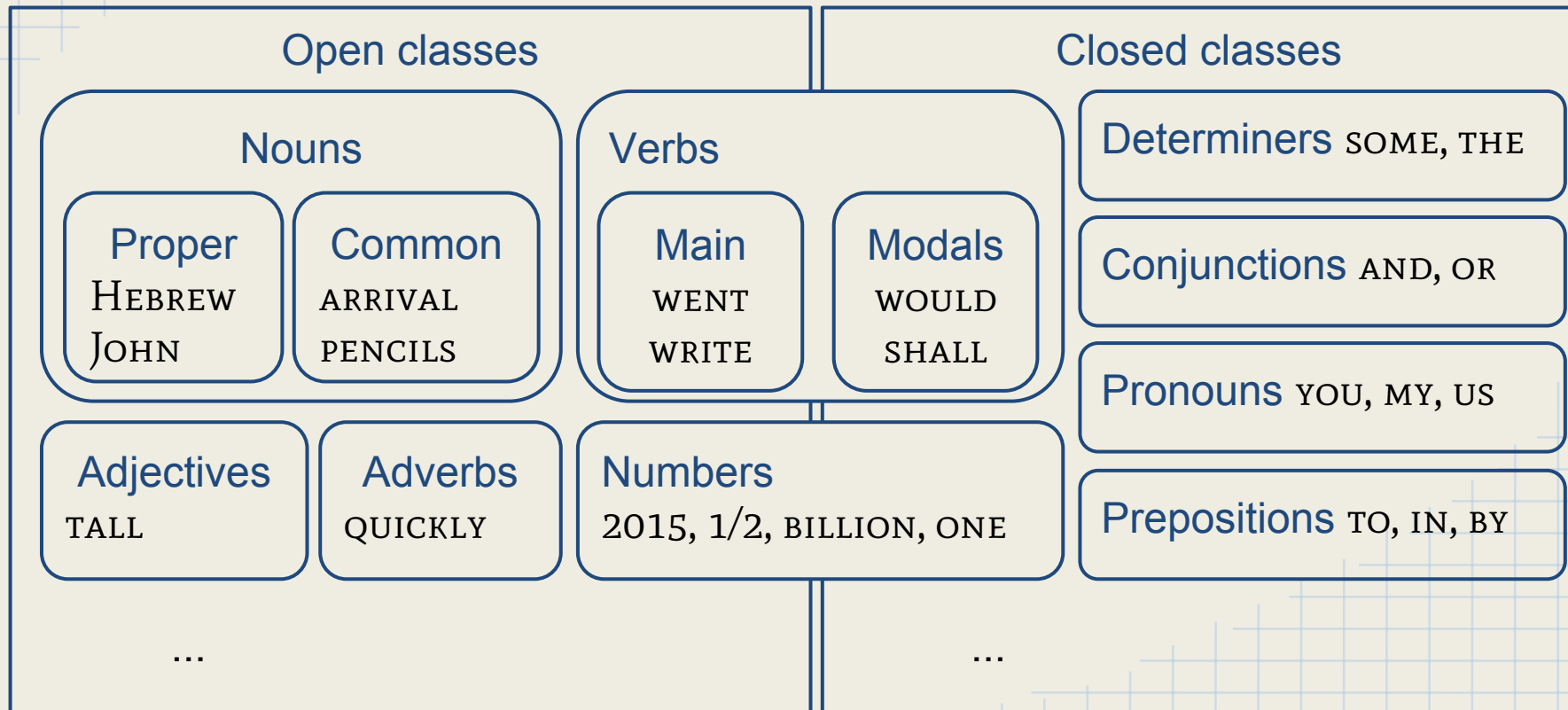
REAGAN WINS ON BUDGET, BUT MORE **LIES** AHEAD

Parts of speech (POS)

Words can roughly be divided into **distributional categories** based on their **syntactic roles**.



Part-of-speech hierarchy



Part-of-speech tags

Tag	Description	Example	Tag	Description	Example
CC	Coordin. Conjunction	<i>and, but, or</i>	SYM	Symbol	<i>+, %, &</i>
CD	Cardinal number	<i>one, two, three</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>	VBN	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, singular	<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>
PP	Personal pronoun	<i>I, you, he</i>	(Left parenthesis	<i>([, (, { , <</i>
PP\$	Possessive pronoun	<i>your, one’s</i>)	Right parenthesis	<i>(] ,) , } , ></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>			

Tag guide:

<https://catalog ldc.upenn.edu/docs/LDC99T42/tagguid1.pdf>

Penn Treebank Part-of-Speech Tags
for English

Language variations

AD	adverb	还
AS	aspect marker	着
BA	把 in ba-construction	把, 将
CC	coordinating conjunction	和
CD	cardinal number	一
CS	subordinating conjunction	虽然
DEC	的 in a relative-clause	的
DEG	associative 的	的
DER	得 in V-de const. and V-de-R	得
DEV	地 before VP	地
DT	determiner	这
ETC	for words 等, 等等	等, 等等
FW	foreign words	I S O
IJ	interjection	啊
JJ	other noun-modifier	男, 共同
LB	被 in long bei-const	被, 给
LC	localizer	里
M	measure word	个
MSP	other particle	所

NN	common noun	书
NR	proper noun	美国
NT	temporal noun	今天
OD	ordinal number	第一
ON	onomatopoeia	哈哈, 哗哗
P	preposition excl. 被 and 把	从
PN	pronoun	他
PU	punctuation	、?。
SB	被 in short bei-const	被, 给
SP	sentence-final particle	吗
VA	predicative adjective	红
VC	是	是
VE	有 as the main verb	有
VV	other verb	走

Penn Treebank Part-of-Speech Tags
for Mandarin Chinese

Part-of-speech tagging

Tag the following text for POS:

ALICE WAS BEGINNING TO GET VERY TIRED

NNP

VBD

VBG

TO

VB

RB

JJ

Statistical POS tagging

We can use counts from the corpus to tag text for POS,
but it requires **annotation**:
just the text is not enough.

Annotated corpus example

Alice/NNP was/VBD beginning/VBG to/TO get/VB very/RB tired/JJ of/IN sitting/VBG by/IN her/PRP\$ sister/NN on/IN the/DT bank/NN ,/, and/CC of/IN having/VBG nothing/NN to/TO do/VB ./.. Once/RB or/CC twice/RB she/PRP had/VBD peeped/VBN into/IN the/DT book/NN her/PRP\$ sister/NN was/VBD reading/VBG ,/, but/CC it/PRP had/VBD no/DT pictures/NNS or/CC conversations/NNS in/IN it/PRP ,/, "/" and/CC what/WDT is/VBZ the/DT use/NN of/IN a/DT book/NN ,/, "/" thought/VBD Alice/NNP ,/, "/" without/IN pictures/NNS or/CC conversations/NNS ?/. "/" So/CC she/PRP

Word/tag counts

Simple method: count
the times each word
occurred with each
POS in the corpus

THE	DT	1527
WELL	RB	37
WELL	NN	3
SLEEP	NN	4
SLEEP	VBP	2
THAT	IN	197
THAT	DT	50

Algorithm to count word/tag

count(L, T):

$Cwt \leftarrow [0]$

$i \leftarrow 1$

while $i \leq \text{len}(L)$:

$Cwt[L[i], T[i]] \leftarrow Cwt[L[i], T[i]] + 1$ ▷ increment count

$i \leftarrow i + 1$

return Cwt

- ▷ L: list of tokens, T: list of correct tags
- ▷ create a table of zeros
- ▷ assign 1 to i
- ▷ repeat while i is at most $\text{len}(L)=\text{len}(T)$
- ▷ increment i
- ▷ output is the counts table

POS tagging algorithm

Find POS sequence of token sequence

{BOB WENT OUT FOR A SWIM .} → {NNP VBD IN IN DT NN .}

This is the wanted result, but no algorithm is perfect

POS tagging algorithm 1

tag1(L, Cwt):

$T \leftarrow []$

$i \leftarrow 1$

while $i \leq \text{len}(L)$:

$T[i] \leftarrow \text{argmax}(\text{Cwt}[L[i], \cdot])$

$i \leftarrow i + 1$

return T

- ▷ L: tokens, Cwt: word-tag counts
- ▷ create empty list of tags
- ▷ initialize i to 1
- ▷ repeat for all tokens
- ▷ most common tag for L[i]
- ▷ increment i
- ▷ output is list of tags

POS tagging algorithm 1

tag1(L, Cwt):

$T \leftarrow []$

$i \leftarrow 1$

while $i \leq \text{len}(L)$:

$T[i] \leftarrow \text{argmax}(\text{Cwt}[L[i], \cdot])$

$i \leftarrow i + 1$

return T

- ▷ L: tokens, Cwt: word-tag counts
- ▷ create empty list of tags
- ▷ initialize i to 1
- ▷ repeat for all tokens
- ▷ most common tag for L[i]
- ▷ increment i
- ▷ output is list of tags

The same algorithm as prediction with bigram

Surprising accuracy

This simple approach actually gets about 90% of the POS tags correctly!

Most words almost always appear with the same POS.

Problem 1: variability



Use the most common POS for each word

THE

FISH

SLEEP

IN

THAT

WELL

DT

~~NN~~

~~NN~~

IN

~~IN~~

~~RB~~

But the correct tags are:

DT

NNS

VBP

IN

DT

NN

State of the art

The best methods today get slightly more than 97% accuracy, so 90% is not so bad.

Problem 2: unknown words

'T WAS BRILLIG , AND THE SLITHY TOVES

? VBD ? , CC DT ? ?

DID GYRE AND GIMBLE IN THE WABE ;

VBD ? CC ? IN DT ? :

ALL MIMSY WERE THE BOROGOVES ,

DT ? VBD DT ? ,

AND THE MOME RATHS OUTGRABE .

CC DT ? ? ? .

First stanza of
Jabberwocky
from *Through
the
Looking-Glass,
and What Alice
Found There*
(1871) by
Lewis Carroll

Solutions

- Context (above the word level)
- Morphology (below the word level)

Transition counts

Count the times each tag follows another tag.

These are **tag bigram** counts (transition counts).

NN	NN	312
NN	IN	690
NN	DT	113
IN	NN	262
DT	NN	1256
PRP	VBD	847
VBD	DT	464

Algorithm to count tag pairs

count(T):

$Ct2 \leftarrow [0]$

$i \leftarrow 2$

while $i \leq \text{len}(L)$:

$Ct2[T[i - 1], T[i]] \leftarrow Ct2[T[i - 1], T[i]] + 1$

$i \leftarrow i + 1$

return $Ct2$

▷ T: list of correct tags from corpus

▷ create a table of zeros

▷ assign 2 to i

▷ repeat while i is at most len(L)

▷ increment

▷ increment i

▷ output is the counts table

The same algorithm as for counting word bigrams

POS tagging algorithm 2

tag2(L, Cwt, Ct2):

$T \leftarrow []$

$T[1] \leftarrow \operatorname{argmax}(\text{Cwt}[L[1], \cdot])$ \triangleright most common tag for first token

$i \leftarrow 2$

\triangleright initialize i to 2

while $i \leq \text{len}(L)$:

\triangleright repeat for all tokens

$T[i] \leftarrow \operatorname{argmax}(\text{Cwt}[L[i], \cdot] \times \text{Ct2}[T[i-1], \cdot])$ \triangleright multiply counts

$i \leftarrow i + 1$

\triangleright increment i

return T

\triangleright output is list of tags

POS tagging algorithm 2

```
tag2(L, Cwt, Ct2):  
    T ← []  
    T[1] ← argmax(Cwt[L[1], ·])  
    i ← 2  
    while i ≤ len(L):  
        T[i] ← argmax(Cwt[L[i], ·] × Ct2[T[i - 1], ·])  
        i ← i + 1  
    return T
```

- ▷ L: tokens, Cwt: word-tag counts, Ct2: tag bigram counts
- ▷ most common tag for first token
- ▷ initialize i to 2
- ▷ repeat for all tokens
- ▷ multiply counts
- ▷ increment i
- ▷ output is list of tags

Multiply corresponding elements in the two tables

Combining counts

Cwt

THE	DT	1527
WELL	RB	37
WELL	NN	3
SLEEP	NN	4
SLEEP	VBP	2

Ct2

NN	NN	312
NN	IN	690
NN	RB	113
IN	NN	262
DT	NN	1256

Tag the short sentence: SLEEP WELL

Combining counts

Cwt

THE	DT	1527
WELL	RB	37
WELL	NN	3
SLEEP	NN	4
SLEEP	VBP	2

Ct2

NN	NN	312
NN	IN	690
NN	RB	113
IN	NN	262
DT	NN	1256

Tag the short sentence: SLEEP WELL
NN

Combining counts

Cwt

THE DT 1527

WELL RB 37

WELL NN 3

SLEEP NN 4

SLEEP VBP 2

Ct2

NN NN 312

NN IN 690

NN RB 113

IN NN 262

DT NN 1256

NN: $3 \times 312 = 936$

RB: $37 \times 113 = 4181$

Tag the short sentence: SLEEP WELL
NN RB

Phrases

Parts of speech are for single words, but multi-word phrases may have similar syntactic roles.

Noun phrases (NP)

I SAW A [DOG]

[SMALL DOG]

[SMALL DOG WITH A BLACK TAIL]

Verb phrases (VP)

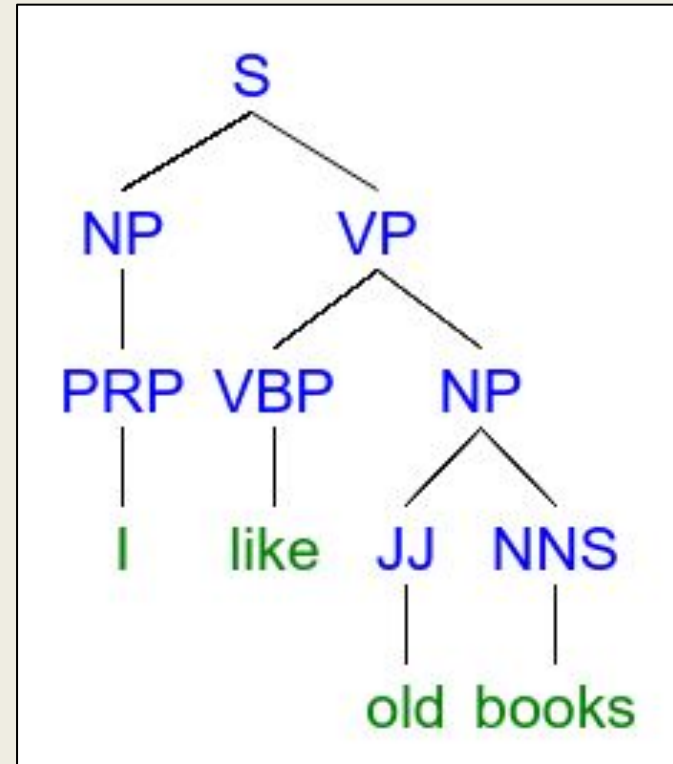
I [WALK]

I [WALK HOME]

I [WALK HOME QUICKLY BUT SURELY]

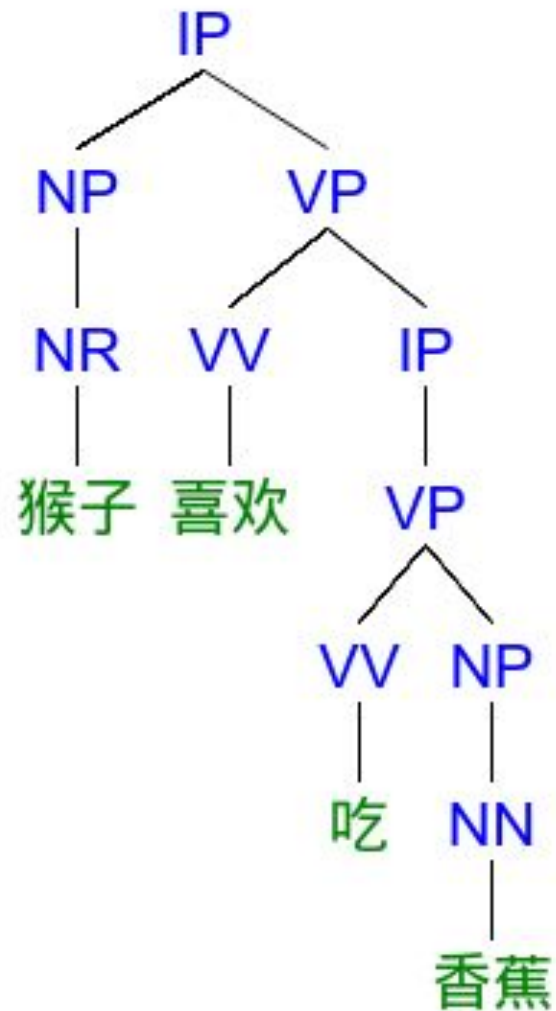
Phrase structure

Represents text structure as
a **tree**: tokens are leaves

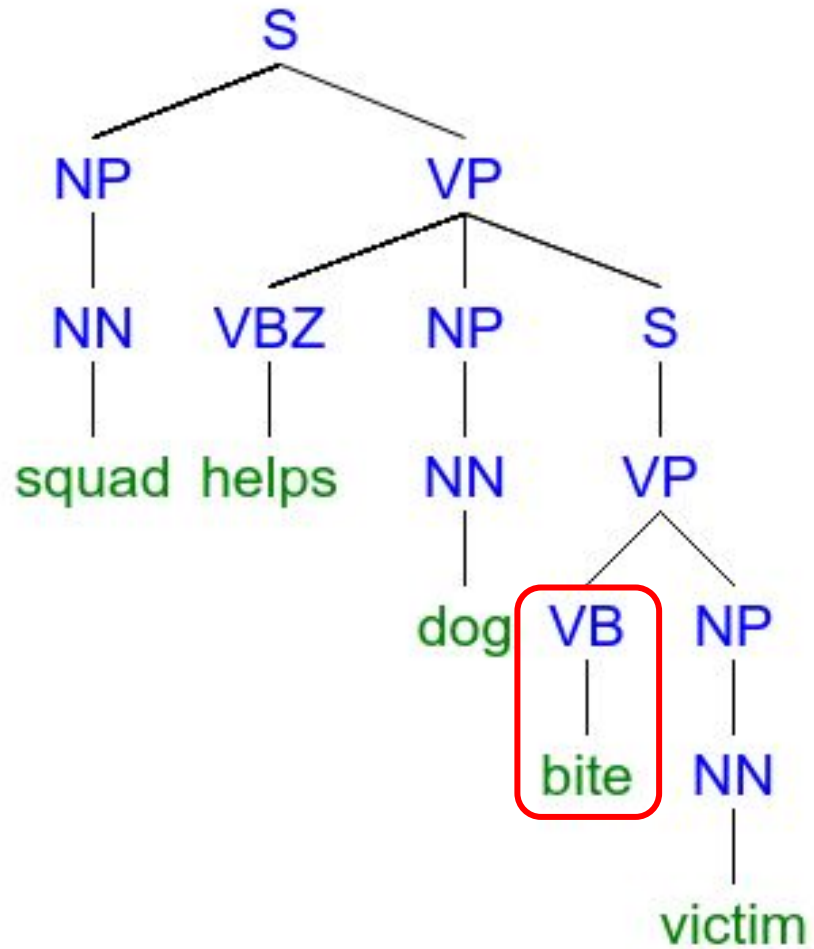
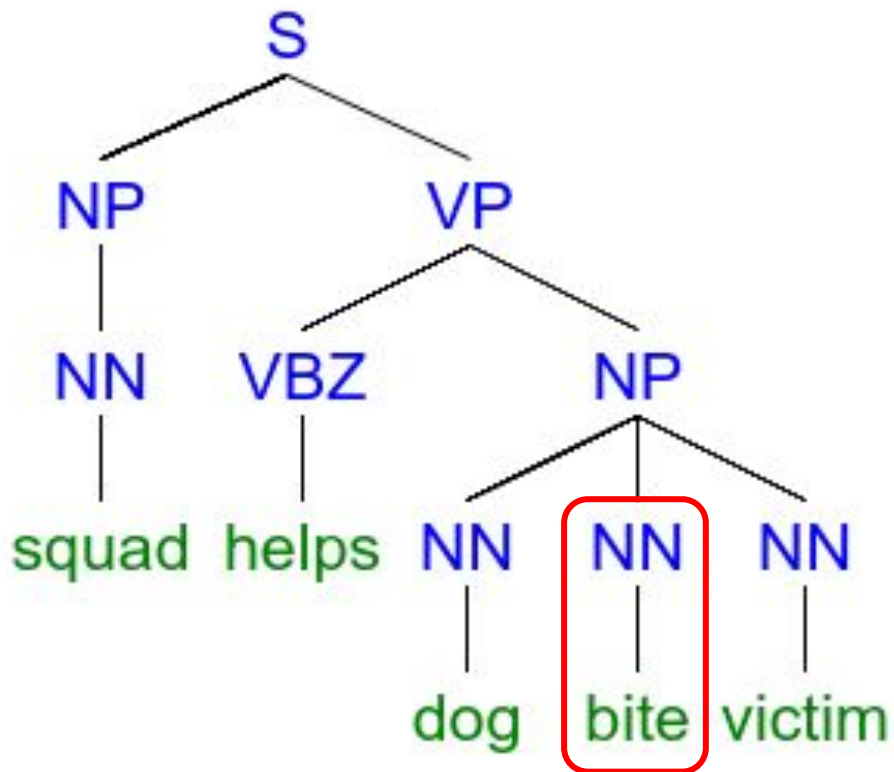


Chinese example

Different rules/labels are
used for different languages



Lexical ambiguity



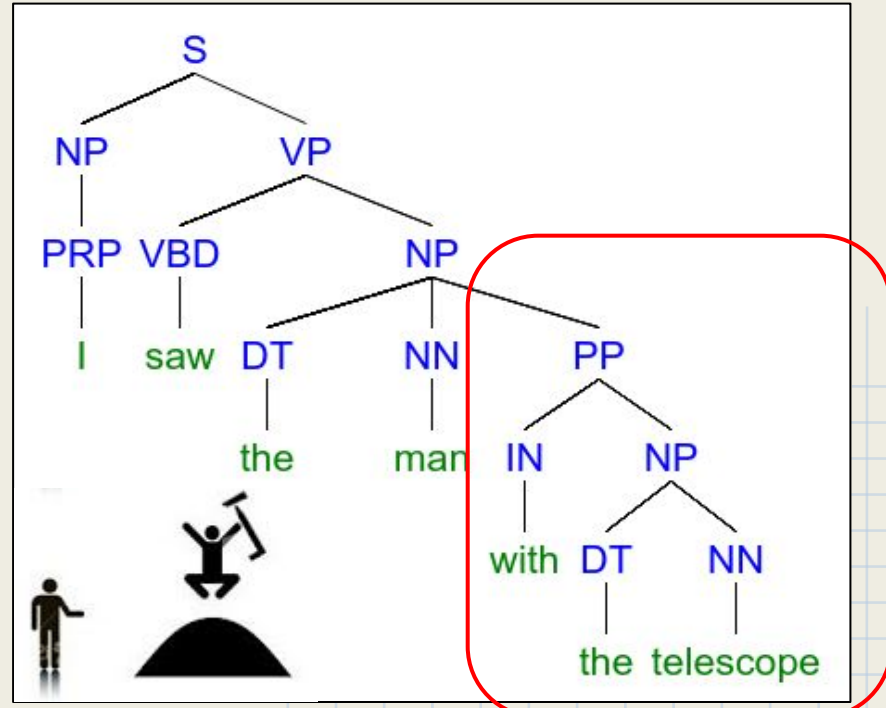
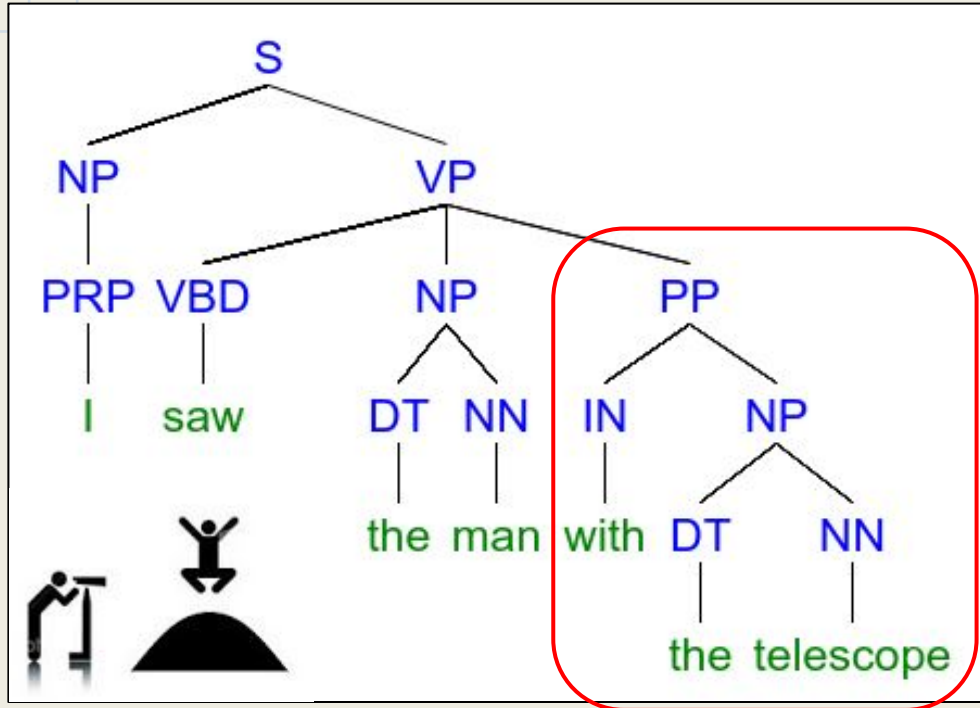
Prepositional phrases (PP)

CATS FALL [ON THEIR FEET]

I'M WEARING THE SHIRT [FROM ITALY]

I'M TAKING THE BUS [FROM TEL AVIV]

Syntactic ambiguity



Adjective phrases (AP)

THIS CAR IS [FAST]

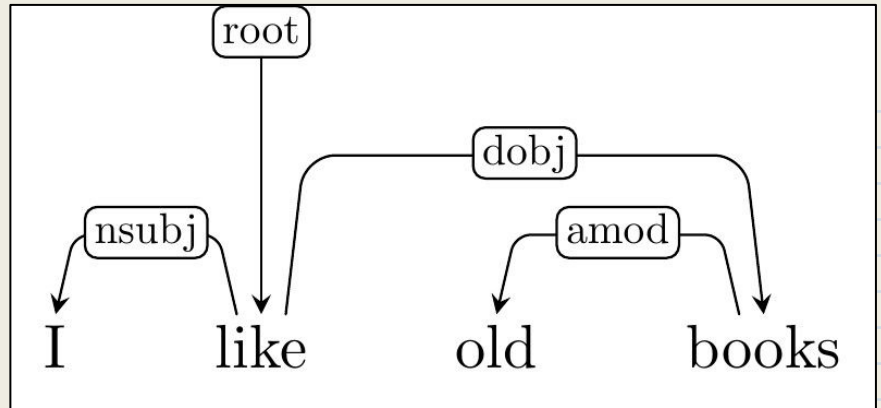
THIS CAR IS [REALLY VERY FAST]

THIS CAR IS [FASTER THAN MY OLD ONE]

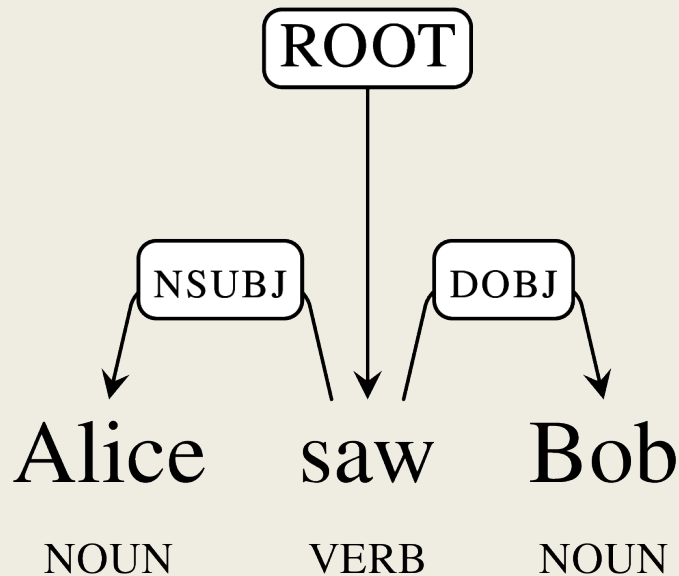
Dependency parsing

Represents text structure as a **tree**:

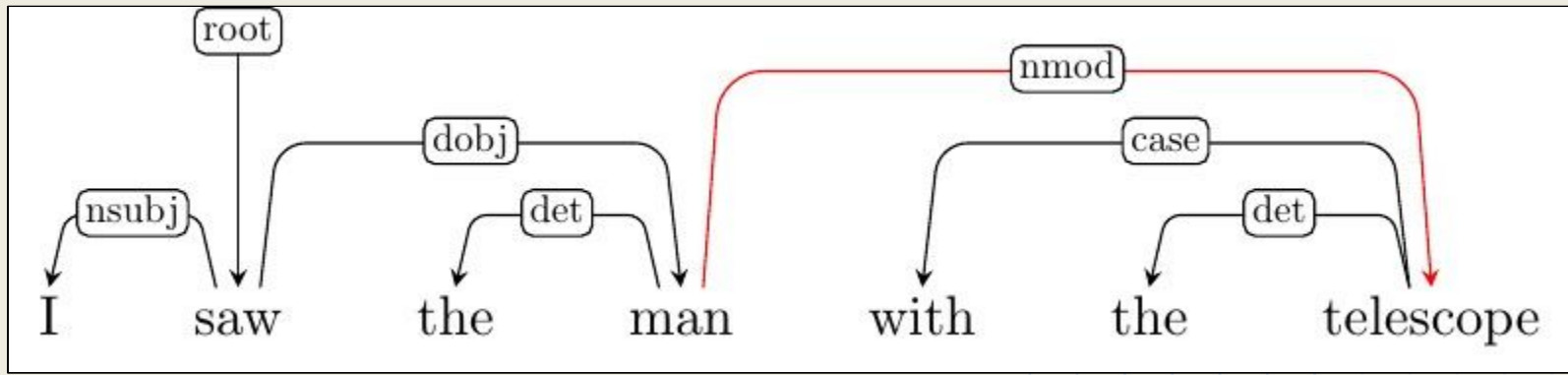
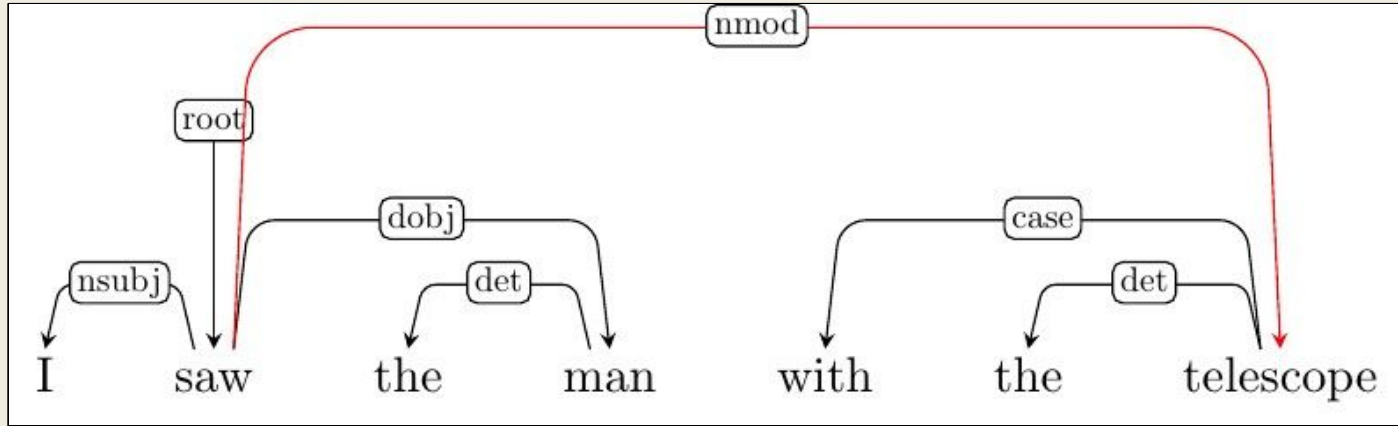
tokens are all the nodes (not just leaves)



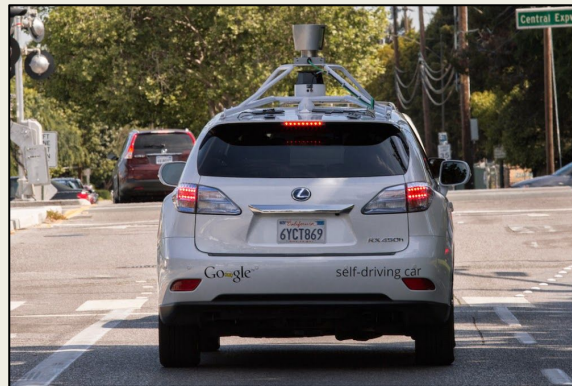
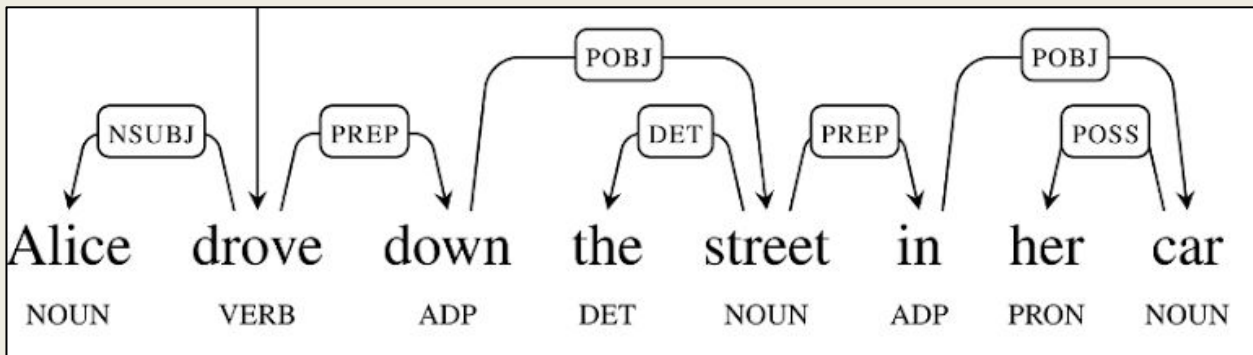
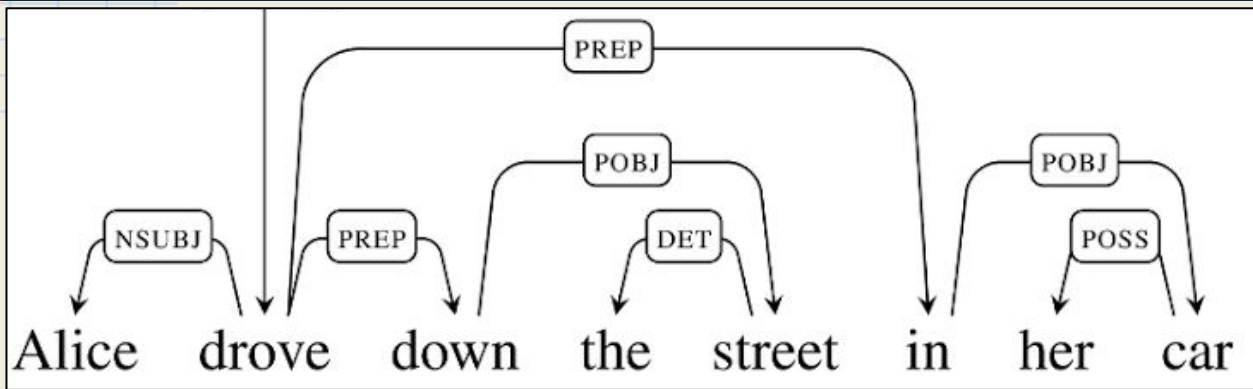
Dependency Parsing



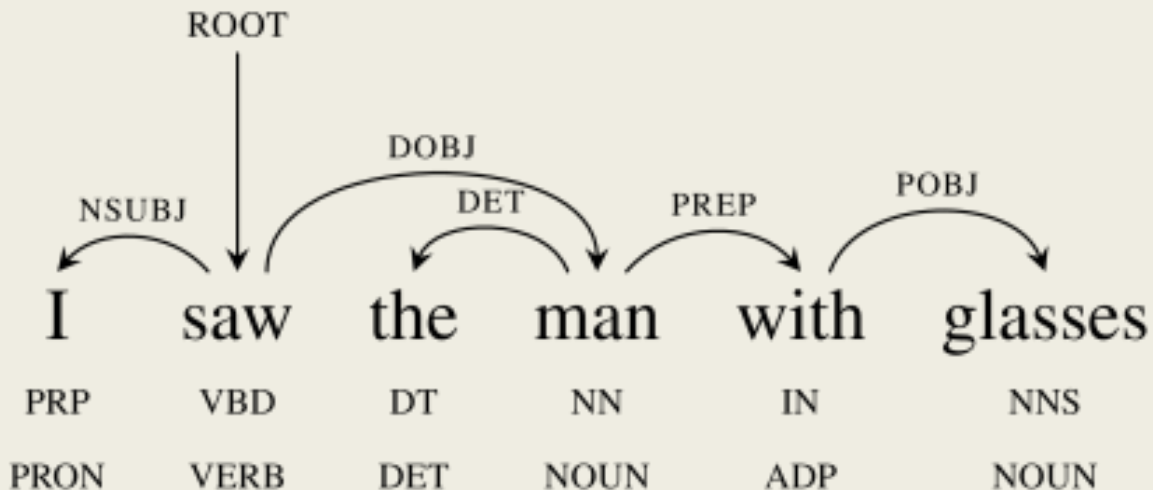
Syntactic ambiguity



Syntactic Ambiguity



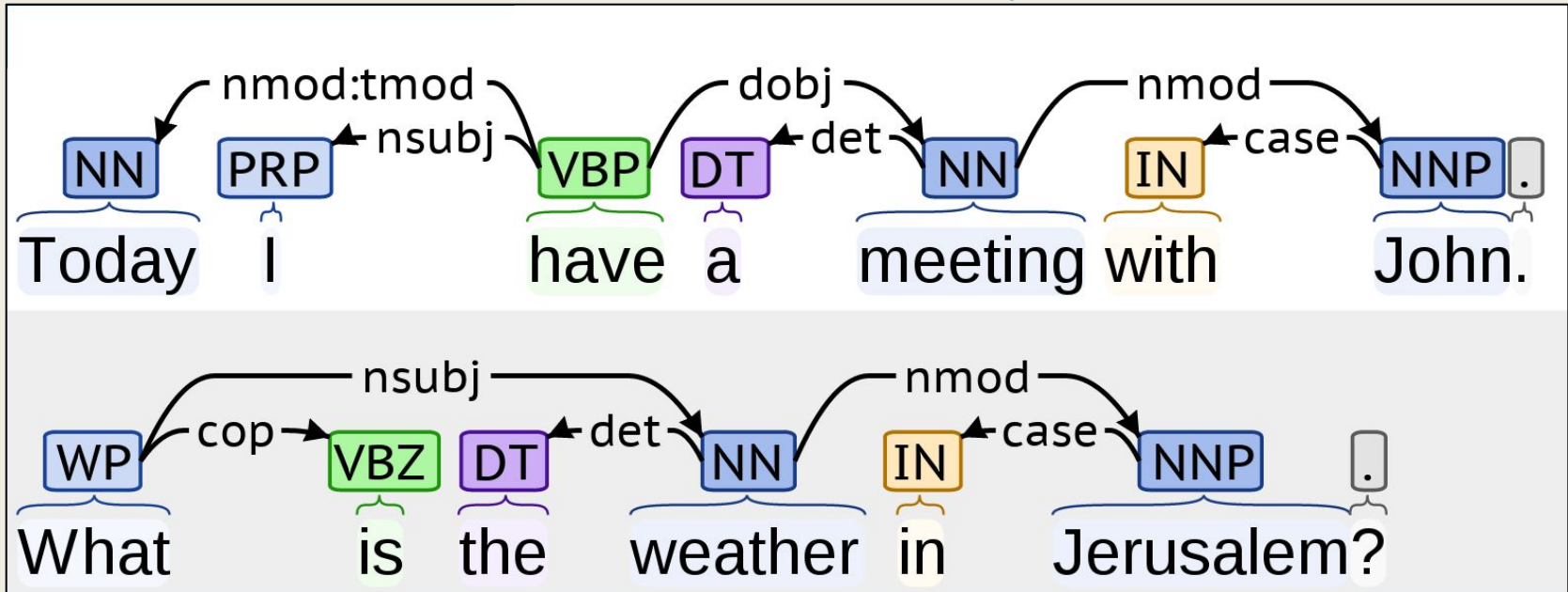
Syntactic Ambiguity



Who had the glasses?

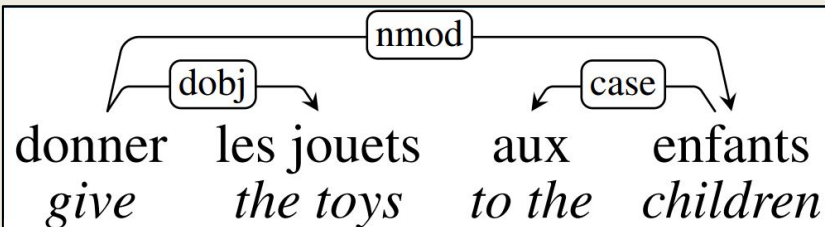
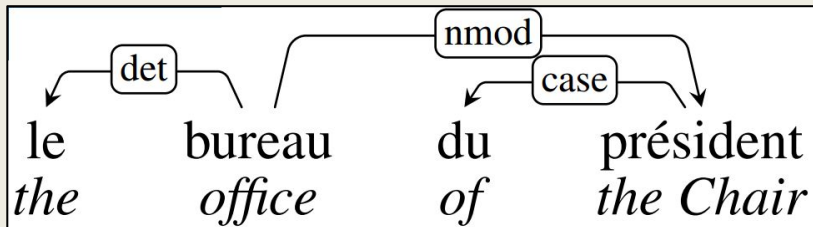
Representation

Natural Language Understanding: who did what to whom and where and when and how and why?

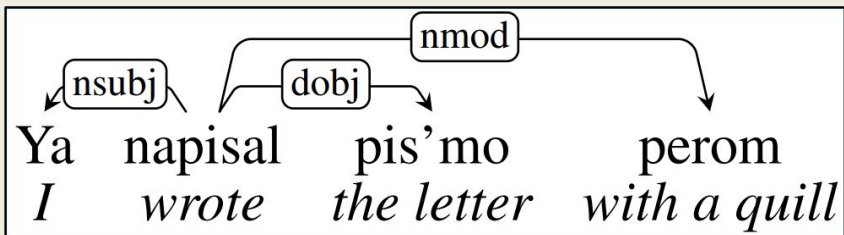
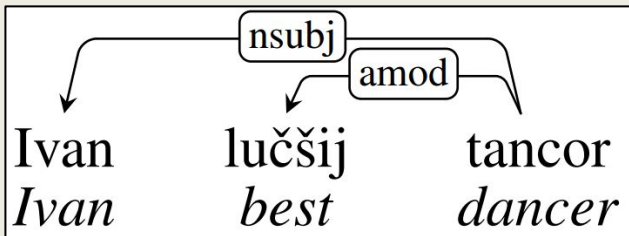


Cross-Linguistic Examples

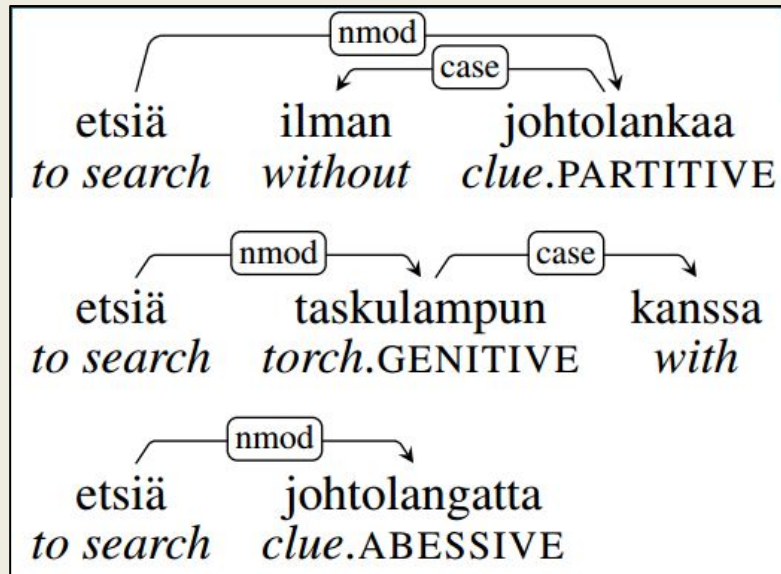
French



Russian



Finnish

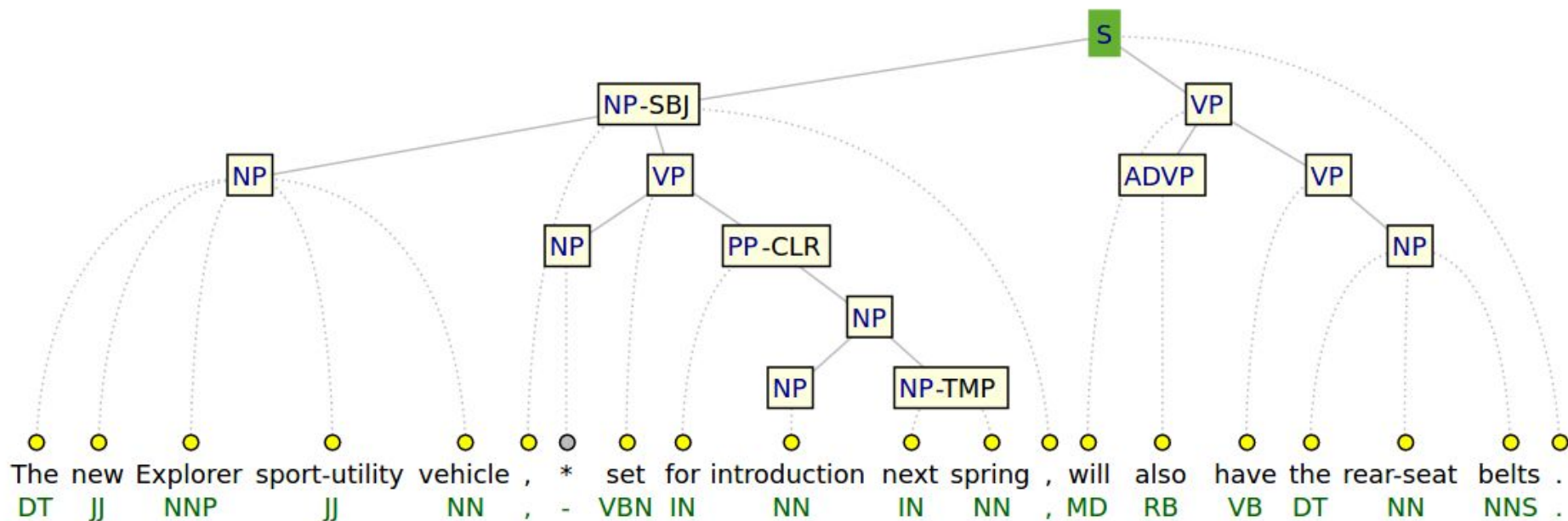


Resources: Treebanks

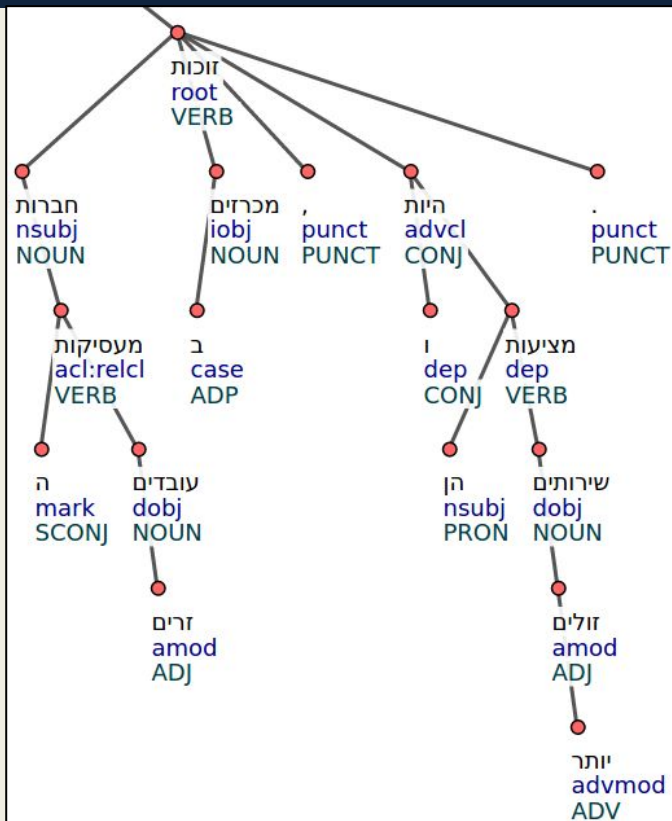
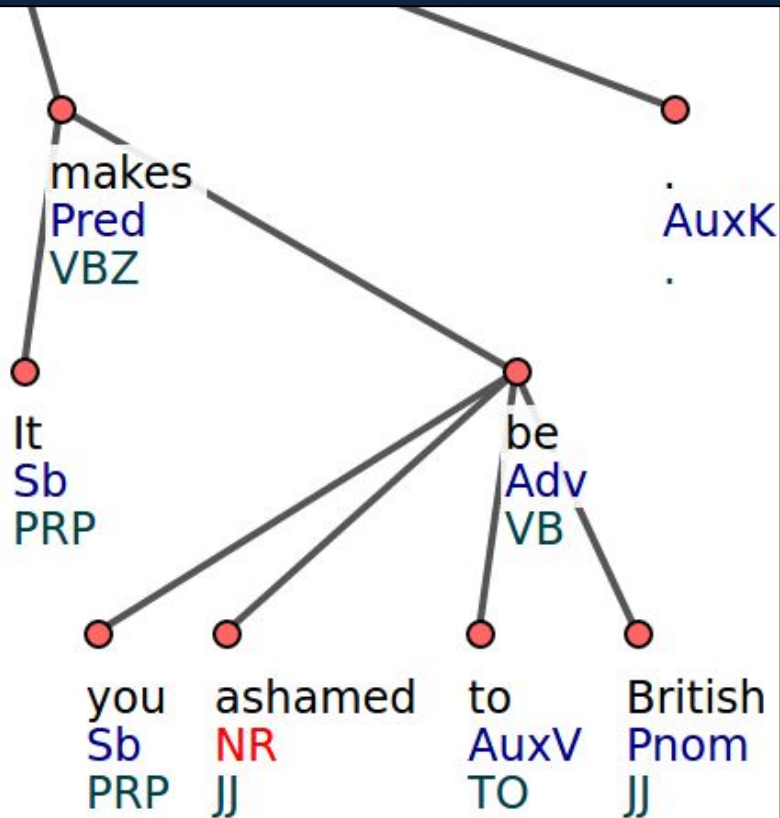
- Many text corpora parsed by humans
- Used for training automatic parsers



Penn Treebank (Constituency)

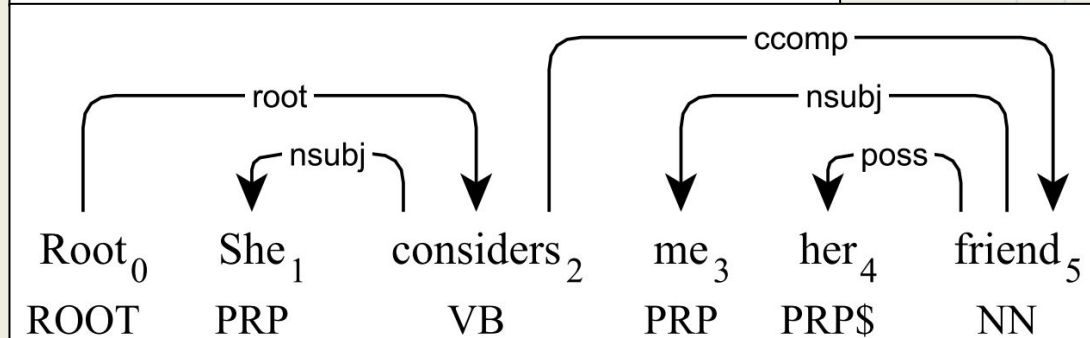
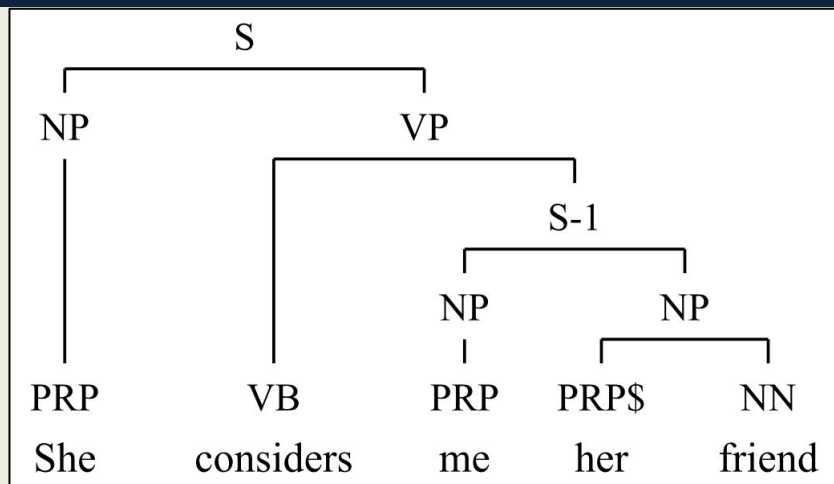


Dependency Treebanks



Converted Treebanks

Trees can be automatically converted to save manual work



Evaluation

Labeled Attachment Score (LAS):

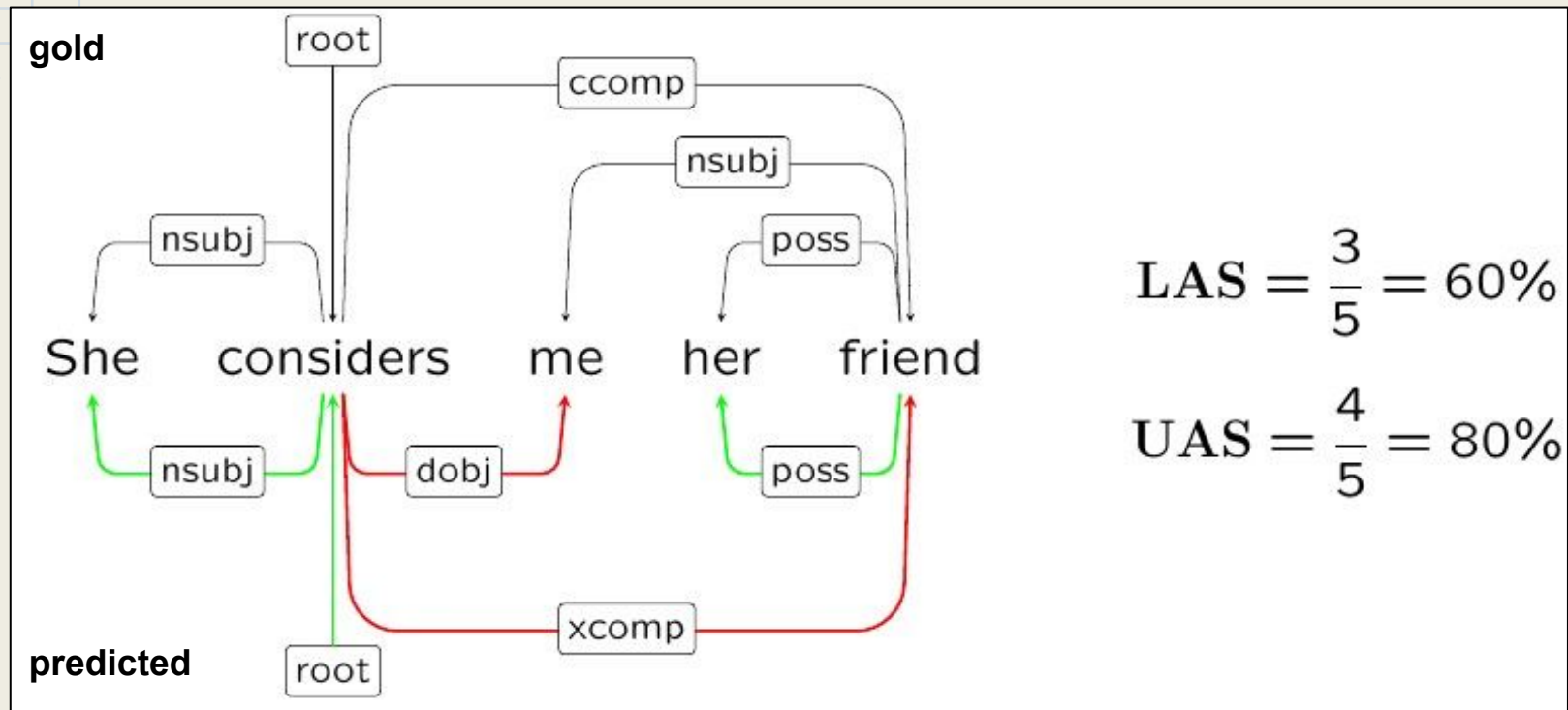
% of edges both in predicted tree and in gold tree

Unlabeled Attachment Score (UAS):

same as LAS, but ignoring edge label

$$0 \leq \text{LAS} \leq \text{UAS} \leq 100\%$$

Evaluation Example



Parser Scores

Parser	UAS	LAS
MaltParser	90.93	88.95
MSTParser	92.17	89.86
ZPar	92.93	91.28
TurboParser	93.80	92.00
Parsey McParseface	94.41	92.55

Summary

	Dependency	Constituency
Structure	tree	tree
Tokens are	all nodes	only leaves
Labels on	edges	nodes

References

- NLP class on Coursera: class.coursera.org/nlp
- Parts of speech: en.wikipedia.org/wiki/Part_of_speech
- Jurafsky, Daniel, and James H. Martin. 2009. Speech and Language Processing: An Introduction to Natural Language Processing, Speech Recognition, and Computational Linguistics. 2nd edition. Prentice-Hall. Pg. 295.

I DON'T MEAN TO GO ALL LANGUAGE
NERD ON YOU, BUT I JUST LEGIT
ADVERBED "LEGIT," VERBED "ADVERB,"
AND ADJECTIVED "LANGUAGE NERD."

