

Part-of-Speech Tagging, Syntactic Parsing

Human Language from a Computational Perspective
April 5, 2016

Part-of-Speech Tagging

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): לבקש, לאחל

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

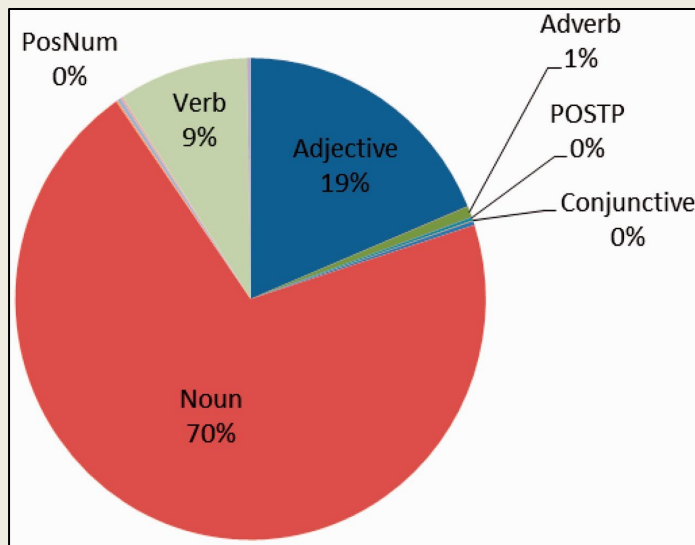
SQUAD HELPS DOG BITE VICTIM



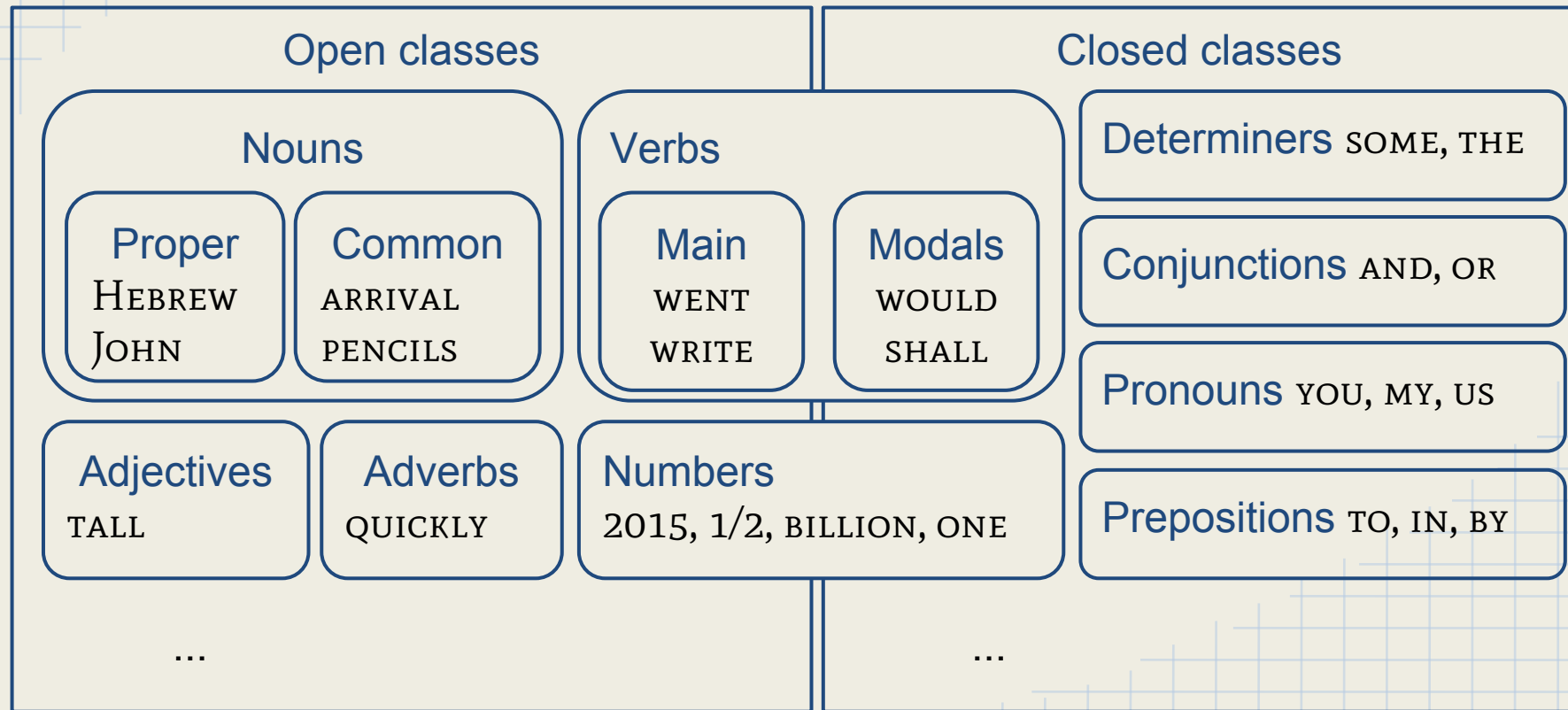
בסוף המרוץ הוא נפל מהסוס. הוא לא פרש.

Parts of Speech

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, Jurafsky & Martin 2009

Language Variations

AD	adverb	还
AS	aspect marker	着
BA	把 in ba-construction	把, 将
CC	coordinating conjunction	和
CD	cardinal number	一, 百
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 SO
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, Xia 2000

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

The/AT grand/JJ jury/NN commented/VBD on/IN a/AT
number/NN of/IN other/AP topics/NNS ,/, AMONG/IN them/PPO
the/AT Atlanta/NP and/CC Fulton/NP-tl County/NN-tl
purchasing/VBG departments/NNS which/WDT it/PPS said/VBD
``/`` ARE/BER well/QL operated/VBN and/CC follow/VB
generally/RB accepted/VBN practices/NNS which/WDT
inure/VB to/IN the/AT best/JJT interest/NN of/IN both/ABX
governments/NNS ''/'' ./.

Lexical POS 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

POS Tagging Algorithm

Find POS sequence of token sequence

Given: ["WHAT", "IS", "THE", "ANSWER", "?"]

Return: ["WP", "VBZ", "DT", "NN", "."]

POS Tagging Algorithm 1

TAG_POS_SIMPLE(TOKENS, COUNTSL):

FOR INDEX \leftarrow 1 TO LENGTH(TOKENS):

 TAGS_{INDEX} \leftarrow **Max2**(TOKENS_{INDEX}, COUNTSL)

RETURN TAGS

TOKENS is a sequence of strings

COUNTSL is a table of [string, string] \rightarrow number

Returns a sequence of strings

Auxiliary Algorithm

MAX2(TOKEN, COUNTSL):

MAX \leftarrow 0

FOR EACH $[T_1, T_2]$ IN COUNTSL:

IF $(T_1 = \text{TOKEN})$ AND $(\text{COUNTSL}[T_1, T_2] > \text{MAX})$:

MAX \leftarrow COUNTSL $[T_1, T_2]$

BEST $\leftarrow T_2$

RETURN BEST

TOKEN is a string
COUNTSL is a table of [string, string] \rightarrow number
Returns a string

Auxiliary Algorithm

MAX2 is very similar to the **BIGRAM** algorithm for LM text generation.

It just returns the T_2 with the highest count among entries with $T_1 = \text{TOKEN}$.

Surprising Accuracy

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

Most words almost always appear with the same POS.

Problem: 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: 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

POS Tagging Algorithm 2

Tag_POS(TOKENS, COUNTSL, COUNTST):

$\text{TAGS}_1 \leftarrow \mathbf{Max2}(\text{TOKENS}_1, \text{COUNTSL})$

FOR INDEX $\leftarrow 2$ TO LENGTH(TOKENS):

$\text{TAGS}_{\text{INDEX}} \leftarrow \mathbf{Max3}(\text{TOKENS}_{\text{INDEX}}, \text{TAGS}_{\text{INDEX} - 1}, \text{COUNTSL}, \text{COUNTST})$

RETURN TAGS

TOKENS is a sequence of strings
COUNTSL and COUNTST are tables of [string, string]→number
Returns a sequence of strings

Combining the Counts

How to implement **MAX3**?

Multiply lexical count (from COUNTSL)
with transition count (from COUNTST)

Auxiliary Algorithm

Max3(TOKEN, TAG, COUNTSL, COUNTST):

MAX \leftarrow 0

FOR EACH $[T_1, T_2]$ IN COUNTSL:

IF ($T_1 = \text{TOKEN}$):

SCORE \leftarrow COUNTSL $[T_1, T_2]$ \times COUNTST[TAG, T_2]

IF (SCORE > MAX):

MAX \leftarrow SCORE

BEST \leftarrow T_2

RETURN BEST

TOKEN and TAG are strings
COUNTSL and COUNTST are tables of [string, string] \rightarrow number
Returns a string

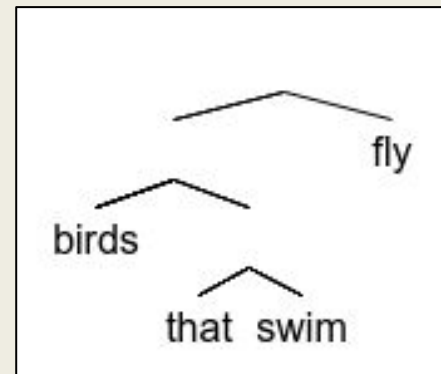


Syntactic Parsing

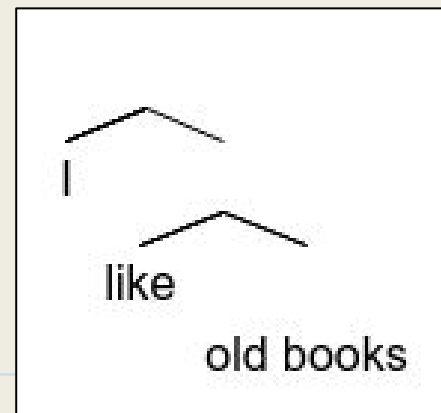
Unlabeled Bracketing

A hierarchy of **constituents**

[[BIRDS [THAT SWIM]] FLY]



[I [LIKE [OLD BOOKS]]]



Labeled Bracketing

Annotates each constituent with a **label**,
and each token with a **part-of-speech tag**

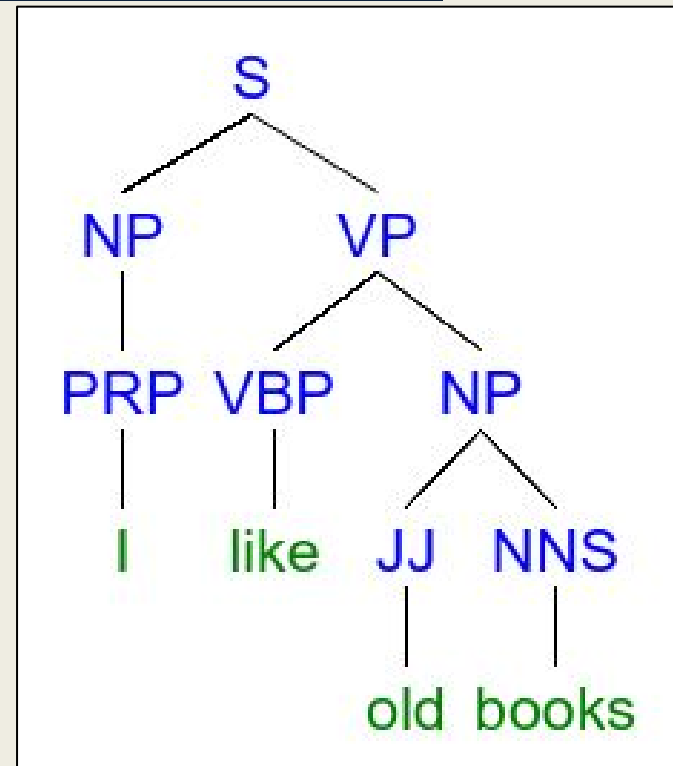
```
[S [NP [PRP I]] [VP [VBP LIKE] [NP [JJ  
OLD] [NNS BOOKS]]]]]
```

Phrase Structure (Constituency Parsing)

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

(Equivalent to labeled bracketing)

```
[S [NP [PRP I]] [VP [VBP LIKE]  
[NP [JJ OLD] [NNS BOOKS]]]]]
```

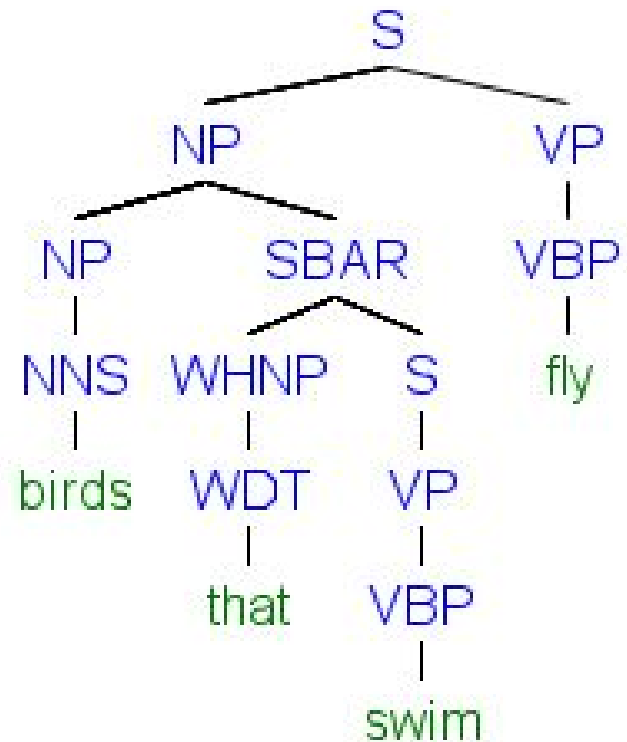


More Complicated Example

Relative clauses introduce
SBAR nodes and are parsed too

<http://cs.jhu.edu/~jason/465/hw-parse/treebank-manual.pdf>

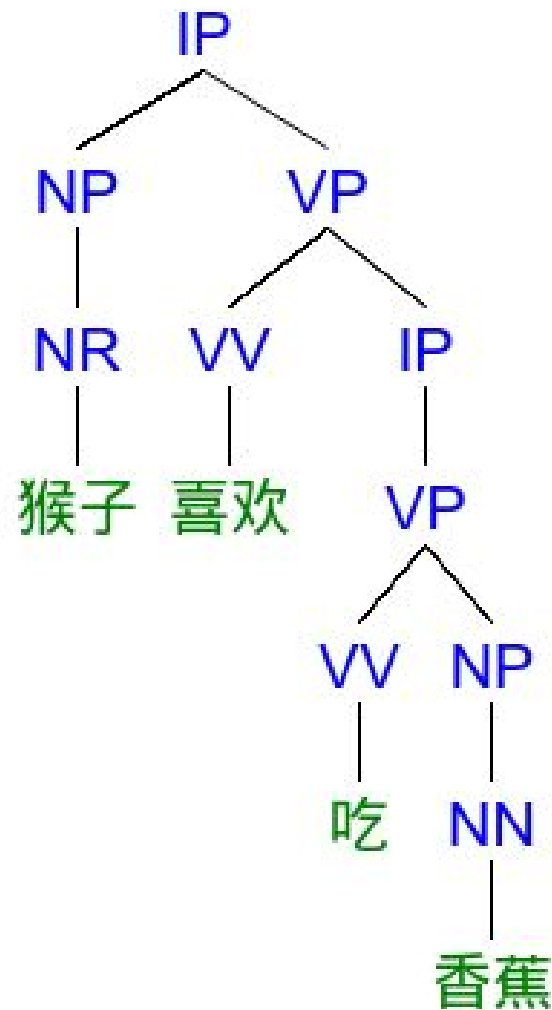
**[S [NP [NP [NNS BIRDS]] [SBAR
[WHNP [WDT THAT]] [S [VP
[VBP SWIM]]]]] [VP [VBP FLY]]]]**



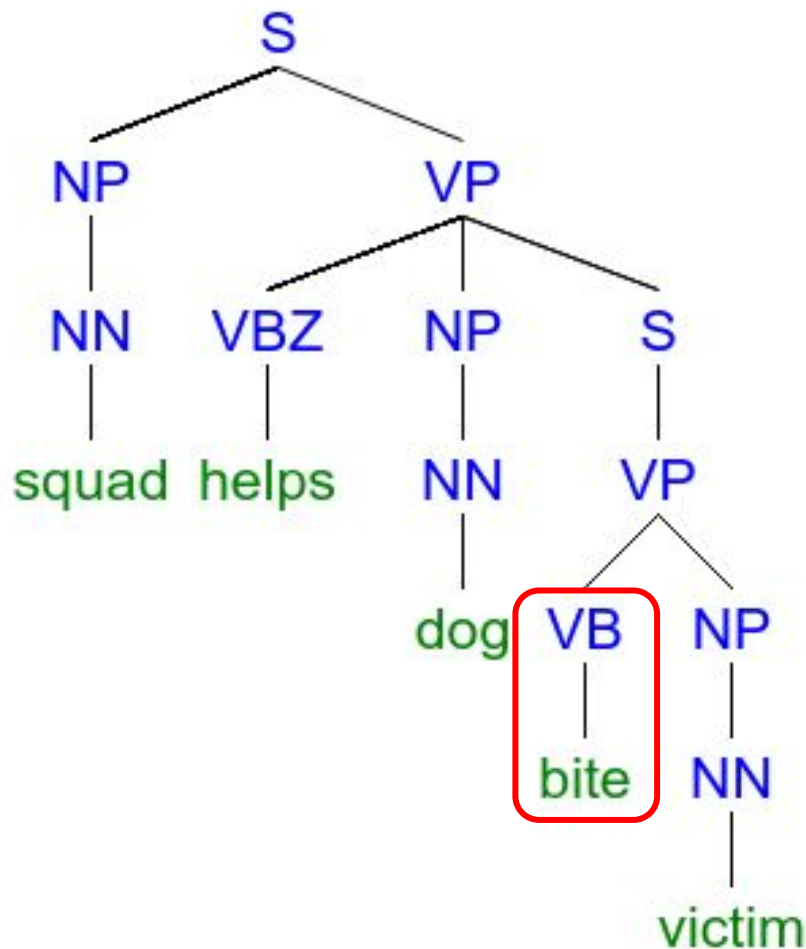
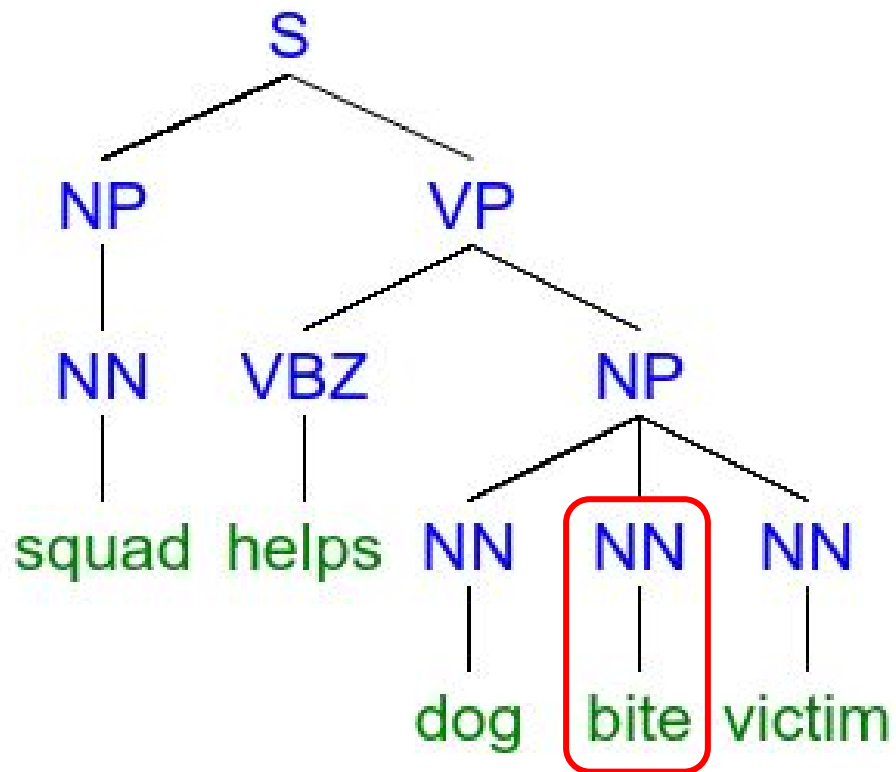
Chinese Example

Different rules/labels are
used for different languages

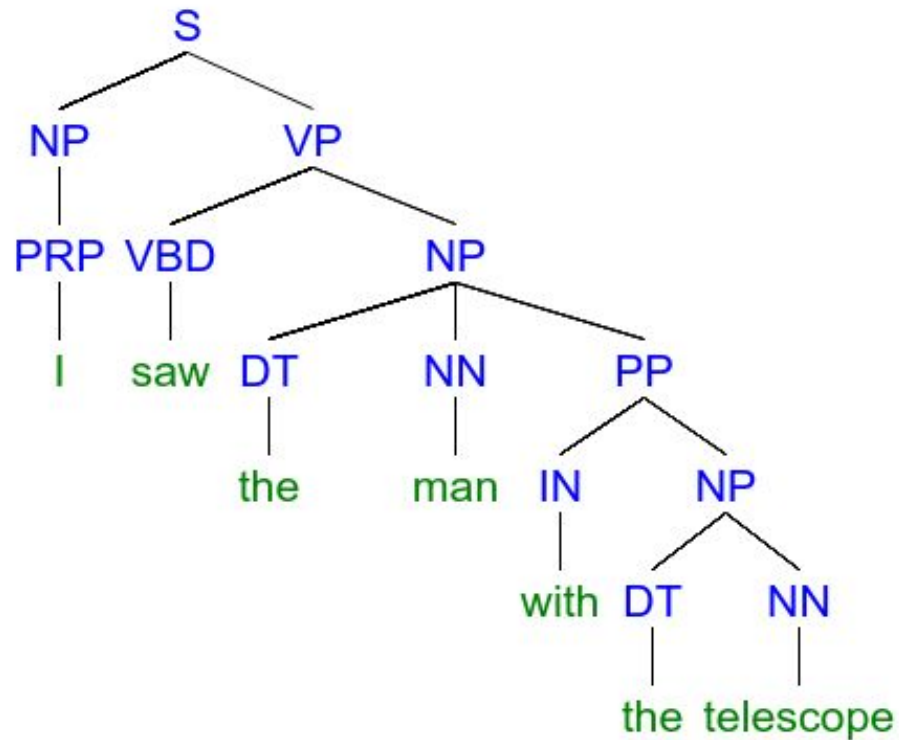
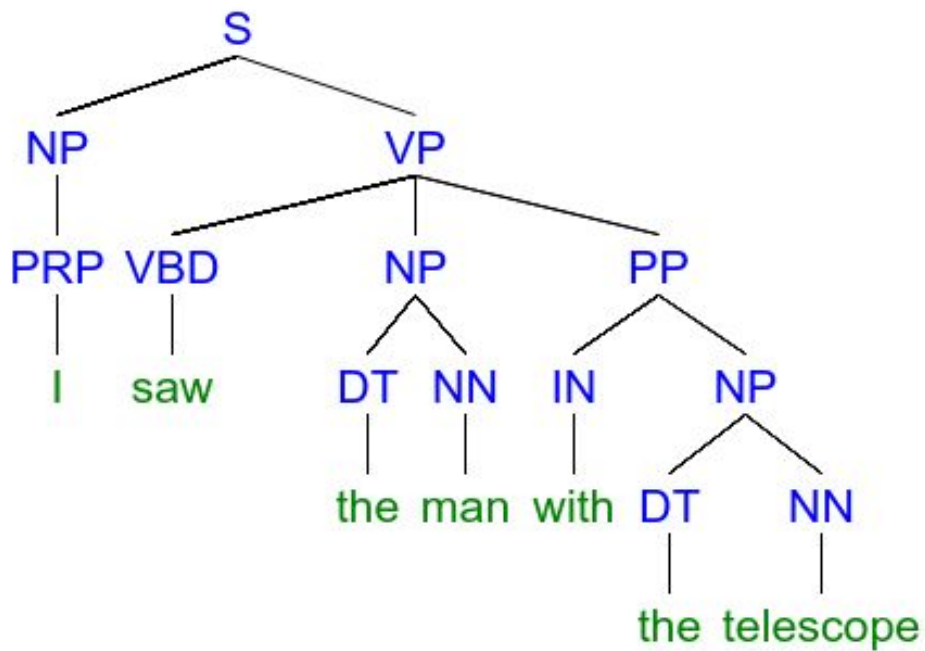
[IP [NP [NR 猴子]] [VP [VV 喜欢]
[IP [VP [VV 吃] [NP [NN 香蕉]]]]]]]



Lexical Ambiguity



Syntactic Ambiguity



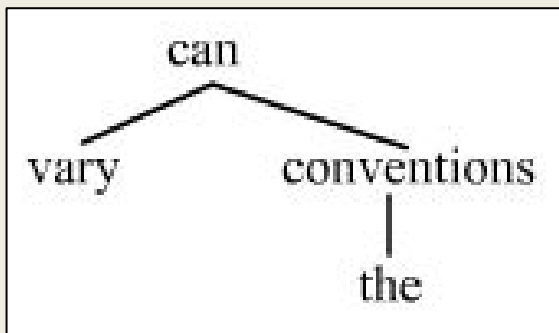
Penn Treebank

- Started in 1989
- More than 4.5 million tokens
- Mostly *Wall Street Journal*
- Constituency-parsed by humans
- Used to train/test parsers

Dependency Parsing

Represents text structure as a **tree**:

tokens are all the nodes (not just leaves)



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."

