



FROM WORDS TO LINGUISTIC STRUCTURE

Sandeep Soni

04/06/2024

WORDS

They can fish

WORDS

They can fish

Is this sentence about the
ability to catch fish?

WORDS

They can fish

... Or is the sentence about
storing fish?

PARTS OF SPEECH

Categories of words that have similar morphological and syntactic properties

Verb Verb

Noun Noun Noun

They can fish

CONTEXTUAL DEPENDENCE

To find the syntactic role that a word has in a sentence,
more context is helpful

Verb Verb

Noun Noun Noun

They can fish

Noun Verb Noun

They can fish in a cold room

PART OF SPEECH TAGGING

The task of assigning a tag (category) to every word in a sequence.

PRP VB NN IN DT JJ NN

They can fish in a cold room

TAG SET

The number of categories depends on language, genre, etc although there are some core categories

PRP VB NN IN DT JJ NN

They can fish in a cold room

Nouns	People, places, things, dates, etc. Depend on quantity (singular Vs Plural)
Adjectives	Properties or qualities of noun phrases
Verbs	Actions, processes. Depends on tense, number, person, etc
Adverbs	Modifiers of verbs; qualify the action (e.g., quickly stopped)
Determiner	Beginning of a noun phrase; emphasizes specificity (e.g., a film)
Prepositions	Indicates spatial/direction/temporal relationship with a noun phrase (e.g., in a room)
Conjunctions	Connectives between phrases, clauses, sentences (e.g., Jack and Jill)
Pronouns	References to noun phrases (e.g., he, she, it)

Tag	Description	Examples	%
Nominal, Nominal + Verbal			
N	common noun (NN, NNS)	books someone	13.7
O	pronoun (personal/WH; not possessive; PRP, WP)	it you u meeee	6.8
S	nominal + possessive	books' someone's	0.1
[^]	proper noun (NNP, NNPS)	lebron usa iPad	6.4
Z	proper noun + possessive	America's	0.2
L	nominal + verbal (= <i>I don't know</i>)	he's book'll iono	1.6
M	proper noun + verbal	Mark'll	0.0
Other open-class words			
V	verb incl. copula, auxiliaries (V*, MD)	might gonna ought couldn't is eats	15.1
A	adjective (J*)	good fav lil	5.1
R	adverb (R*, WRB)	2 (i.e., <i>too</i>)	4.6
!	interjection (UH)	lol haha FTW yea right	2.6
Other closed-class words			
D	determiner (WDT, DT, WP\$, PRP\$)	the teh its it's	6.5
P	pre- or postposition, or subordinating conjunction (IN, TO)	while to for 2 (i.e., to) 4 (i.e., <i>for</i>)	8.7
&	coordinating conjunction (CC)	and n & + BUT	1.7
T	verb particle (RP)	out off Up UP	0.6
X	existential <i>there</i> , predeterminers (EX, PDT)	both	0.1
Y	X + verbal	there's all's	0.0
Twitter/online-specific			
#	hashtag (indicates topic/category for tweet)	#acl	1.0
@	at-mention (indicates another user as a recipient of a tweet)	@BarackObama	4.9
~	discourse marker, indications of continuation of a message across multiple tweets	RT and : in retweet construction RT @user : hello	3.4
U	URL or email address	http://bit.ly/xyz	1.6
E	emoticon	:-) :b (: <3 o_O	1.0

Miscellaneous			
\$	numeral (CD)	2010 four 9:30	1.5
,	punctuation (#, \$, ' ', (,)) , , . , :, ` `)	!!! ?!?	11.6
G	other abbreviations, foreign words, possessive endings, symbols, garbage (FW, POS, SYM, LS)	ily (<i>I love you</i>) wby (<i>what about you</i>) 's ♪ --> awesome...I'm	1.1

SEQUENCE LABELING

- Classification: Predict the tag for each word independently

The old man the boat

SEQUENCE LABELING

- Structured prediction: Predict the entire tag sequence

$$f([x_1, x_2, \dots, x_n]) \rightarrow [y_1, y_2, \dots, y_n]$$

Where can part of speech tagging be used?

Sarcasm as Contrast between a Positive Sentiment and Negative Situation

**Ellen Riloff, Ashequl Qadir, Prafulla Surve, Lalindra De Silva,
Nathan Gilbert, Ruihong Huang**

School Of Computing

University of Utah

Salt Lake City, UT 84112

{riloff,asheq,alnds,ngilbert,huangrh}@cs.utah.edu, prafulla.surve@gmail.com

- (a) Oh how I love *being ignored*. #sarcasm
- (b) Thoroughly enjoyed *shoveling the driveway* today! :) #sarcasm
- (c) Absolutely adore it when *my bus is late* #sarcasm
- (d) I'm so pleased mom *woke me up* with *vacuuming my room* this morning. :) #sarcasm

Identified positive sentiment phrases and negative situation phrases using part of speech tags to build a sarcasm recognizer

PARSING

- Beyond annotating every word with a category or mapping spans to categories, more linguistic structure is enforced by how spans are combined in a sequence

GRAMMAR

$S \rightarrow NP\ VP$

$VP \rightarrow V\ NP \mid VP\ PP$

$PP \rightarrow IN\ NP$

$V \rightarrow eat$

$IN \rightarrow with$

$NP \rightarrow NP\ PP \mid we \mid sushi \mid chopsticks$

These are called the production rules.
Infinite sequences can be generated
by following these rules

SYNTAX TREE

$S \rightarrow NP\ VP$

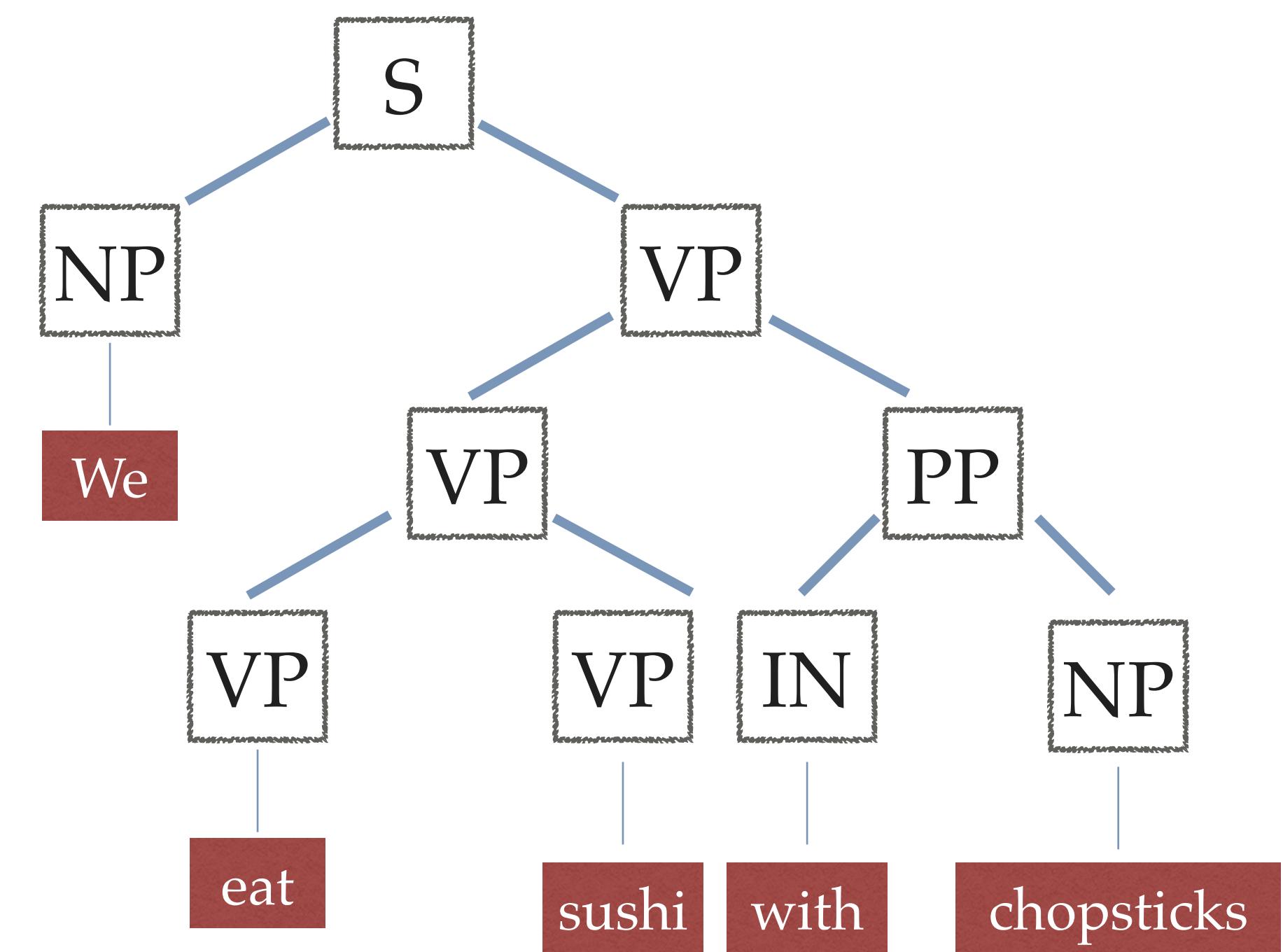
$VP \rightarrow V\ NP \mid VP\ PP$

$PP \rightarrow IN\ NP$

$V \rightarrow eat$

$IN \rightarrow with$

$NP \rightarrow NP\ PP \mid we \mid sushi \mid chopsticks$



PARSING

$S \rightarrow NP\ VP$

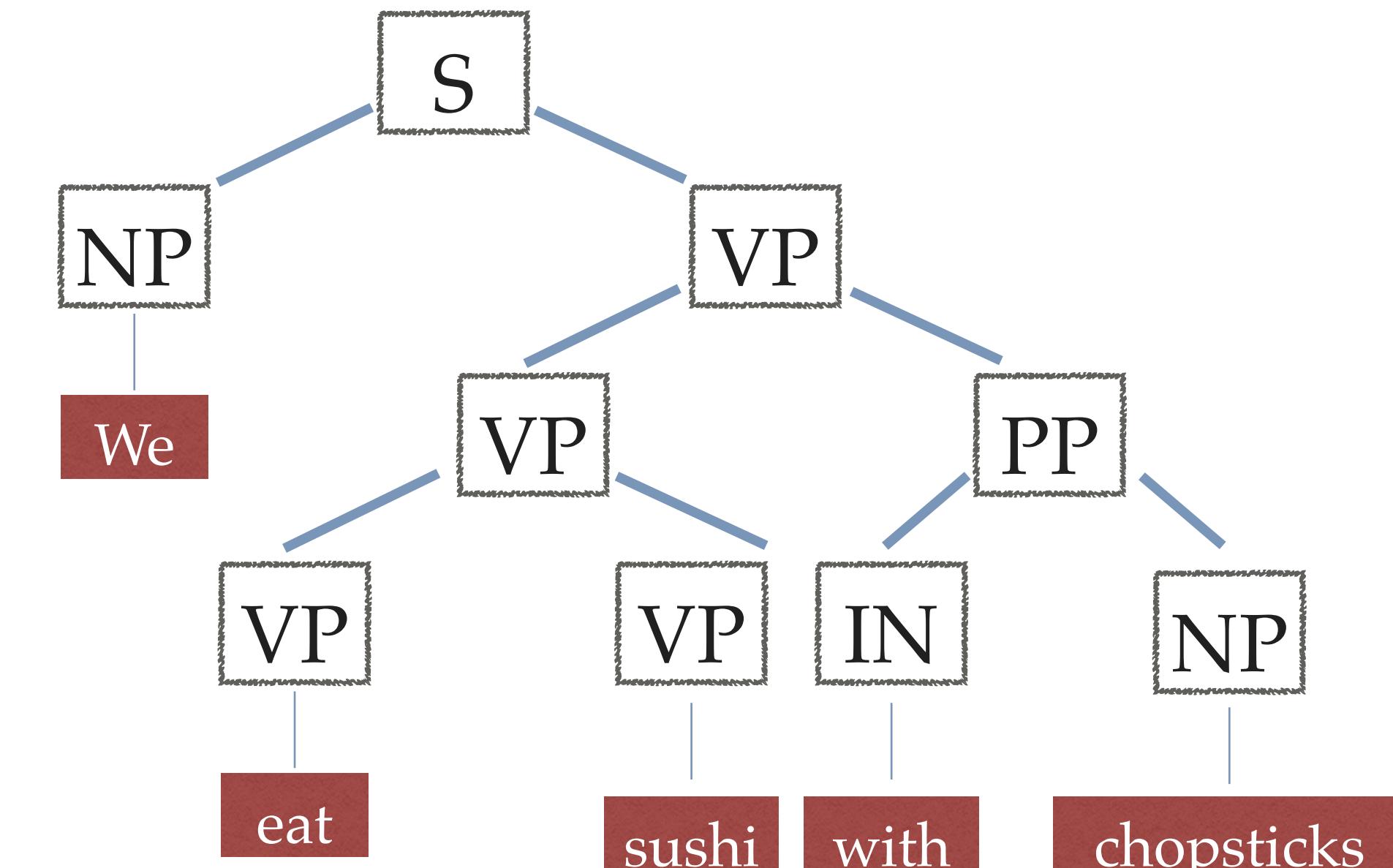
$VP \rightarrow V\ NP \mid VP\ PP$

$PP \rightarrow IN\ NP$

$V \rightarrow eat$

$IN \rightarrow with$

$NP \rightarrow NP\ PP \mid we \mid sushi \mid chopsticks$

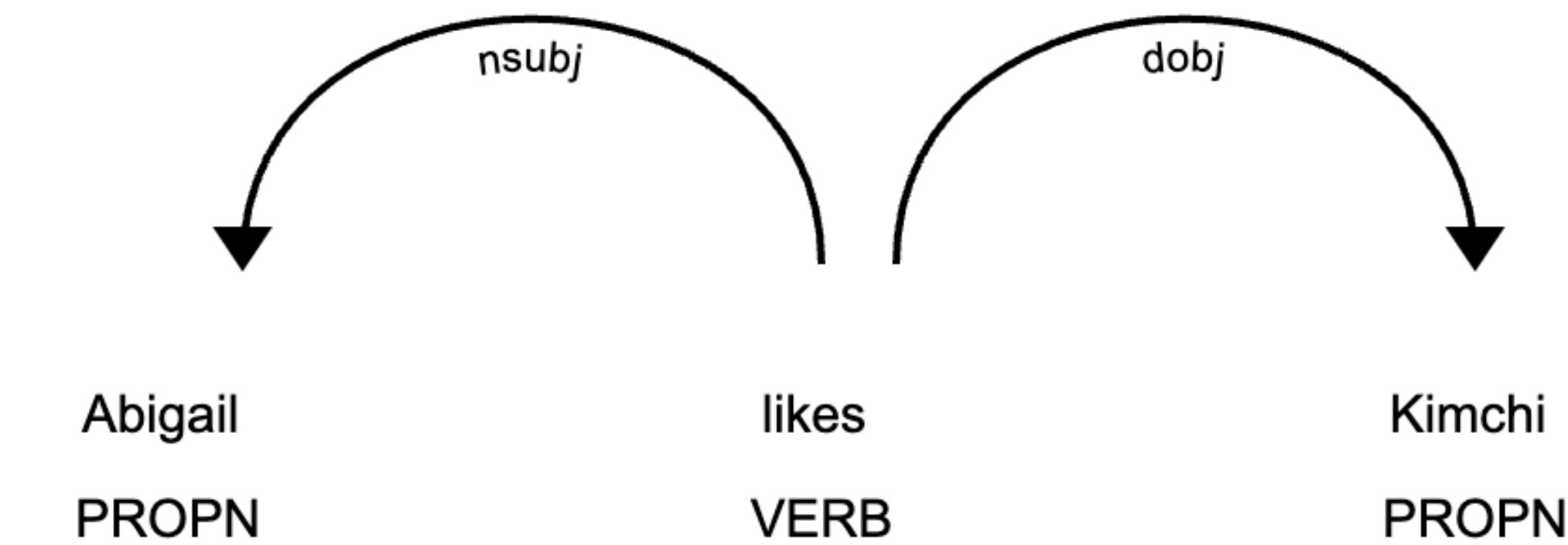


The goal of parsing is to learn the production rules and infer the most likely parse tree for a given sequence

DEPENDENCY PARSE

- Instead of decomposing a sequence into constituent subtrees, one can markup the sequence by labeling:
 - Important constituents (heads)
 - Their modifiers (dependents)
 - The type of relationship between the heads and dependents

EXAMPLE



- What is happening? Something is being liked
- Who is liking? Abigail (the subject)
- What is being liked? Kimchi (the object)

Core dependents of clausal predicates	Non-core dependents of clausal predicates	Special clausal dependents
<p><i>Nominal dep</i> <i>Predicate dep</i></p> <p><u>nsubj</u> <u>csubj</u></p> <p>↳ <u>nsubj:pass</u> ↳ <u>csubj:pass</u></p> <p>↳ <u>nsubj:outer</u> ↳ <u>csubj:outer</u></p> <p><u>obj</u> <u>ccomp</u></p> <p><u>iobj</u></p>	<p><i>Nominal dep</i> <i>Predicate dep</i> <i>Modifier word</i></p> <p><u>obl</u> <u>advcl</u> <u>advmod</u></p> <p>↳ <u>obl:nmod</u> ↳ <u>advcl:relcl</u></p> <p>↳ <u>obl:tmod</u></p>	<p><i>Nominal dep</i> <i>Auxiliary</i> <i>Other</i></p> <p><u>vocative</u> <u>aux</u> <u>mark</u></p> <p><u>discourse</u> ↳ <u>aux:pass</u></p> <p><u>expl</u> <u>cop</u></p>
<p>Noun dependents</p> <p><i>Nominal dep</i> <i>Predicate dep</i> <i>Modifier word</i></p> <p><u>nummod</u> <u>acl</u> <u>amod</u></p> <p>↳ <u>acl:relcl</u></p> <p><u>appos</u> <u>det</u></p> <p>↳ <u>det:predet</u></p> <p><u>nmod</u></p> <p>↳ <u>nmod:nmod</u></p> <p>↳ <u>nmod:tmod</u></p> <p>↳ <u>nmod:poss</u></p>	<p>Compounding and unanalyzed</p> <p><u>compound</u> <u>flat</u></p> <p>↳ <u>compound:prt</u> ↳ <u>flat:foreign</u></p> <p><u>fixed</u> <u>goeswith</u></p>	<p>Coordination</p> <p><u>conj</u> <u>cc</u></p> <p>↳ <u>cc:preconj</u>.</p>
<p>Case-marking, prepositions, possessive</p> <p><u>case</u></p>	<p>Loose joining relations</p> <p><u>list</u> <u>parataxis</u> <u>orphan</u></p> <p><u>dislocated</u> <u>reparandum</u></p>	<p>Other</p> <p><i>Sentence head</i> <i>Punctuation</i> <i>Unspecified dependency</i></p> <p><u>root</u> <u>punct</u> <u>dep</u></p>

There are many syntactic dependencies that can be identified (example taken from universal dependencies treebank)

What can we do with dependency parse?



“Mistakes
were
made.”

by whom?
by you?
by us?!
by

Source: NYTimes

Drivers of English Syntactic Change in the Canadian Parliament

Liwen Hou

Khoury College of Computer Sciences
Northeastern University
Boston, MA
hou.l@northeastern.edu

David A. Smith

Khoury College of Computer Sciences
Northeastern University
Boston, MA
dasmith@ccs.neu.edu

Abstract

Corpus linguists have long noted the “colloquialization” of many genres of English. While the average decline in many features of formal speech is obvious in aggregate, we are better able to disentangle drivers of change by examining Canadian parliamentary speeches coded for characteristics of individual speakers across more than 100 years—much longer than previous studies of individuals’ language change in a common environment. While many language changes proceed by cohort replacement and often originate with female speakers, the Canadian Hansard shows that

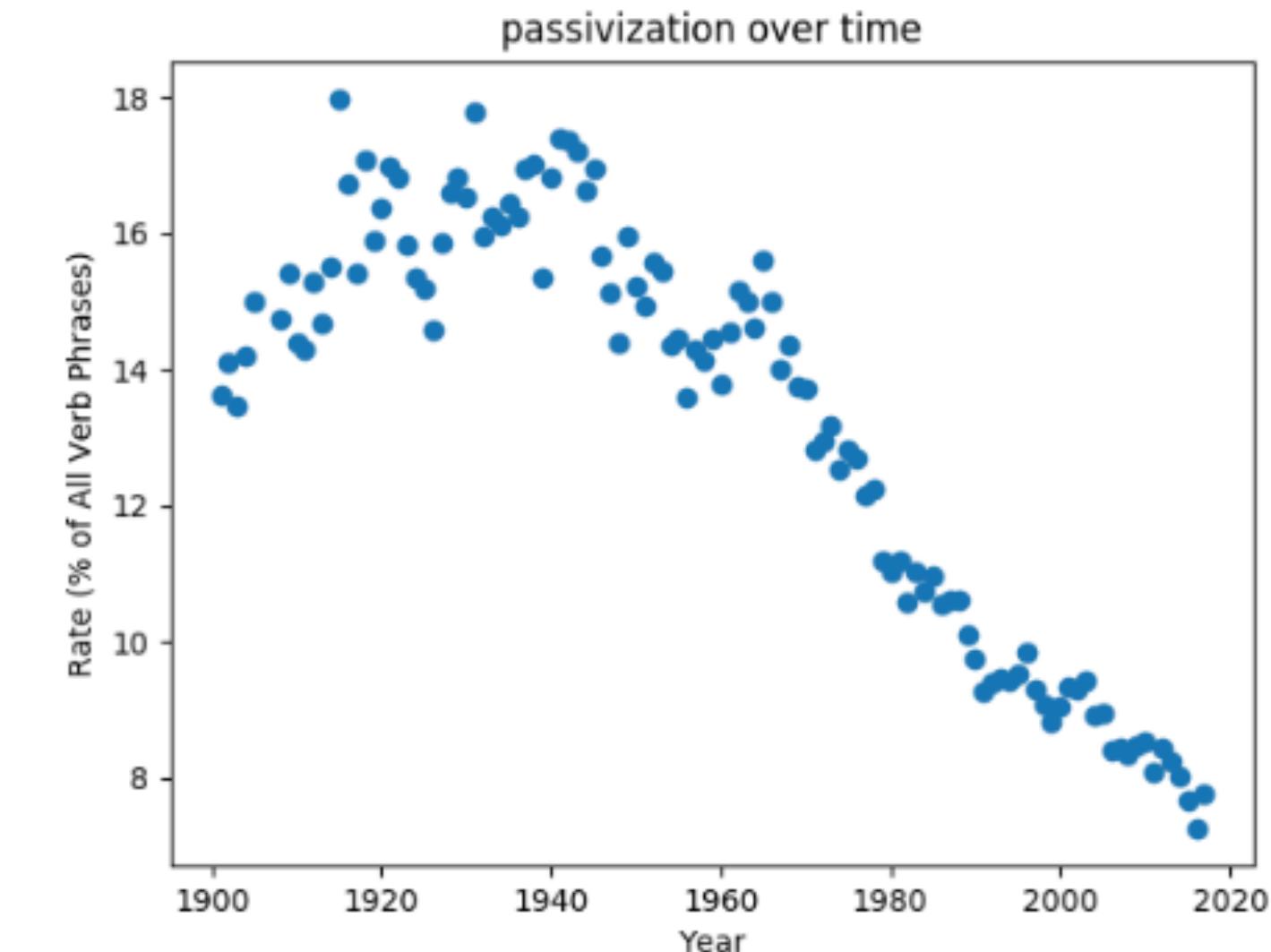
type of language change, we focus on syntactic aspects of colloquialization in this work.

Specifically, we find that the *be*-passive became rarer, that the progressive became more common, that pied-piping became less frequent, that modals such as “shall” also decreased in frequency, and that semi-modals such as “have to” rose in frequency. See [Section 3](#) (and in particular [Table 1](#)) for examples of the five aforementioned features.

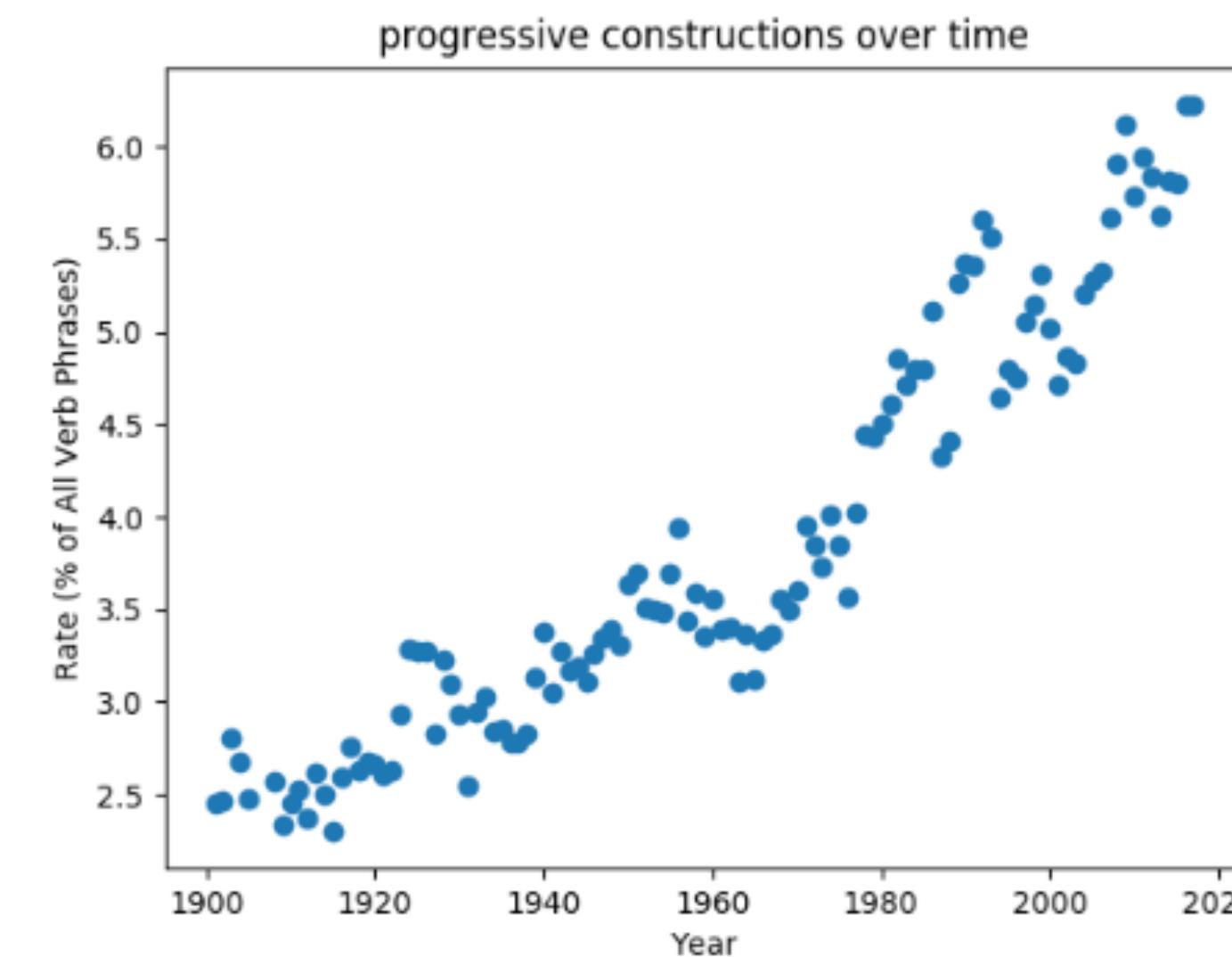
We are particularly interested in the forces driving language change. One natural hypothesis is that newer parliamentary members repeatedly introduce more colloquial language in their speeches, so that

SYNTACTIC CHANGES

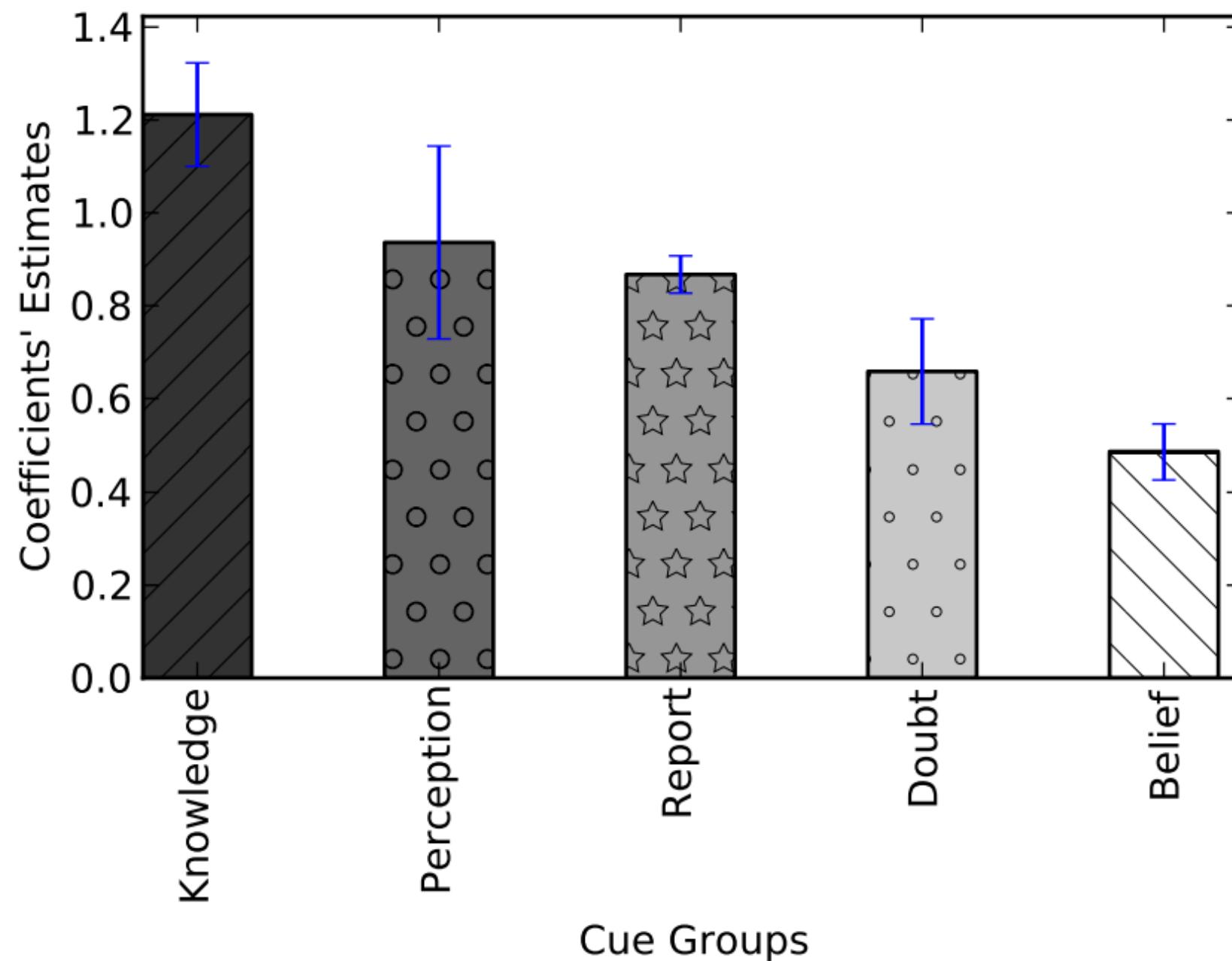
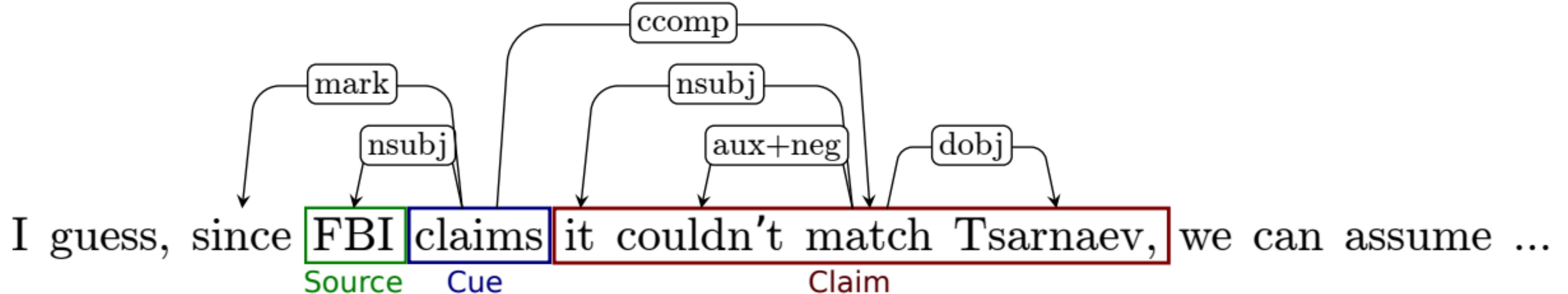
- Colloquialization increasing in political language
- Governing party less colloquial
- Women leaders less colloquial



(a) Passivization

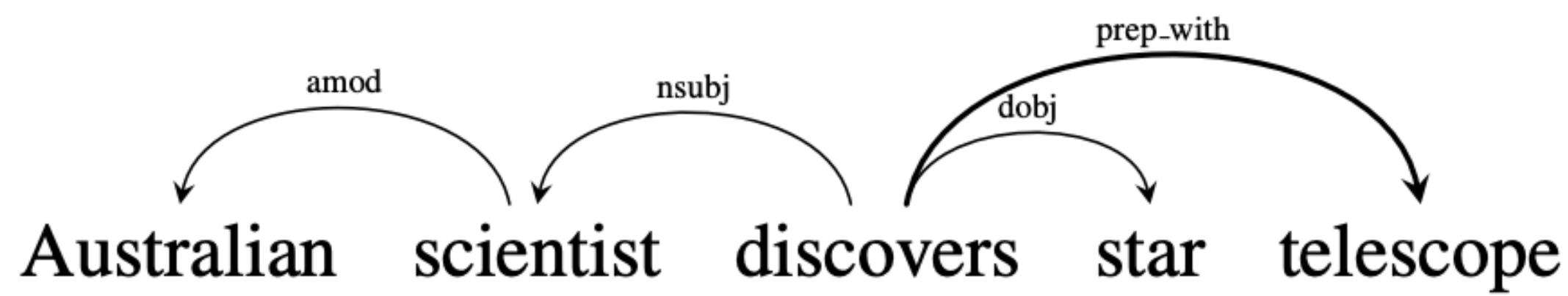
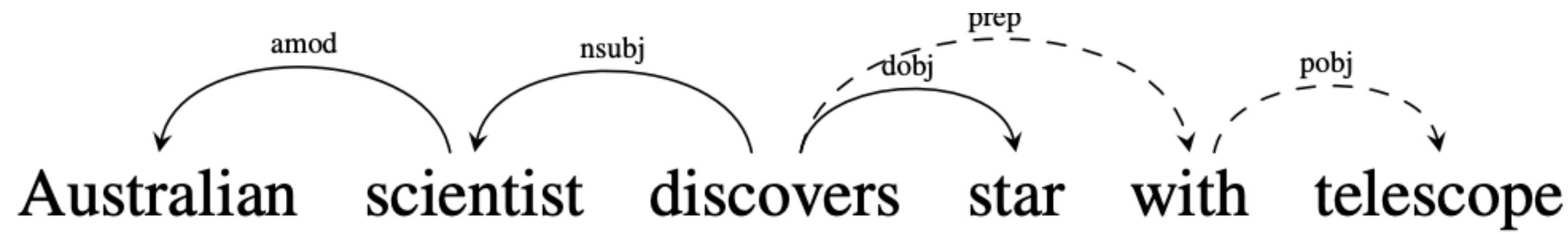


(b) Progressives



Using the dependency parse, we can identify
the different parts of a sentence

- Sandeep Soni, Tanushree Mitra, Eric Gilbert, and Jacob Eisenstein. 2014. [Modeling Factuality Judgments in Social Media Text](#). In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 415–420, Baltimore, Maryland. Association for Computational Linguistics.



WORD	CONTEXTS
australian	scientist/amod ⁻¹
scientist	australian/amod, discovers/nsubj ⁻¹
discovers	scientist/nsubj, star/dobj, telescope/prep_with
star	discovers/dobj ⁻¹
telescope	discovers/prep_with ⁻¹

florida	gainesville fla jacksonville tampa lauderdale	fla alabama gainesville tallahassee texas	texas louisiana georgia california carolina
---------	---	---	---

Similar words (last column is based on dependent parsing contexts)

IN CLASS

- Sequence labeling demo