



# Latent Variable PCFGs

# Extending the idea to induced syntactico-semantic classes

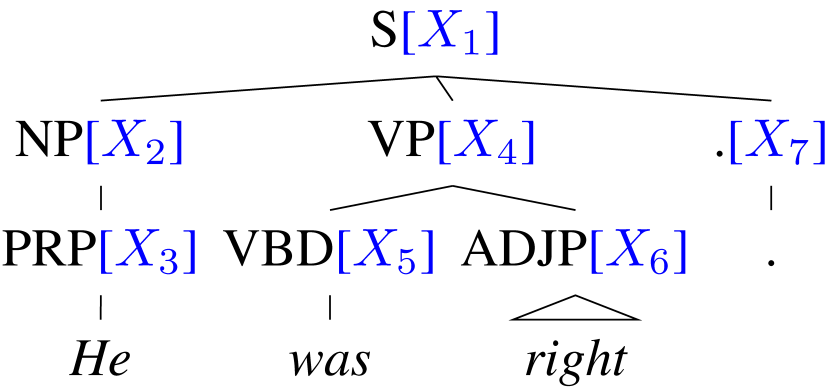


# Learning Latent Annotations

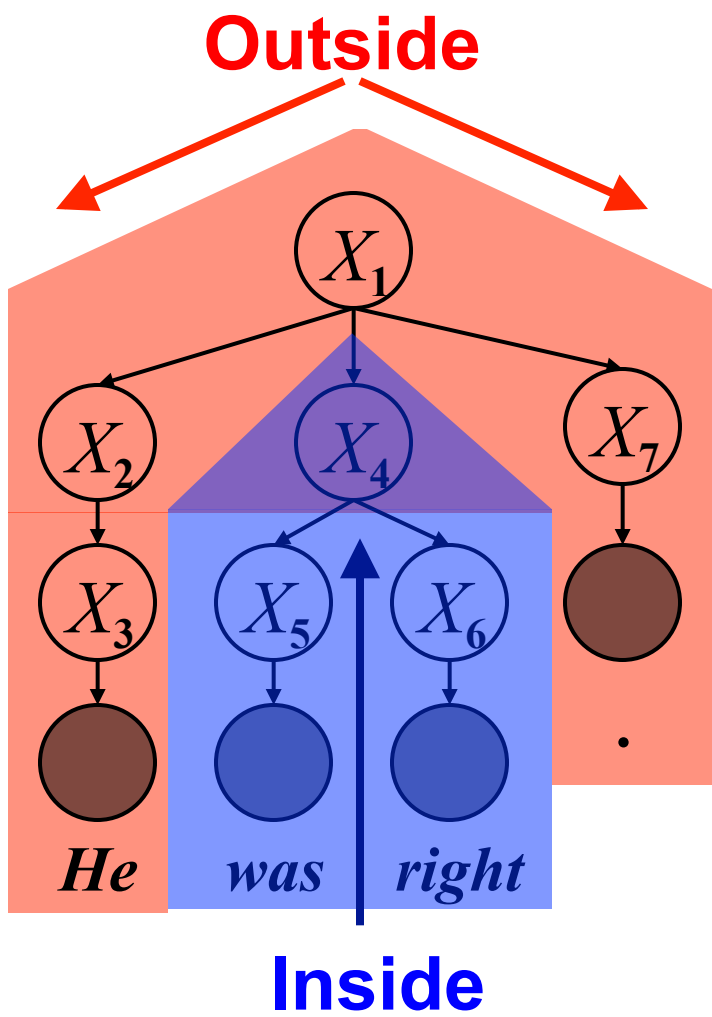
[Petrov and Klein 2006, 2007]

Can you automatically find good symbols?

- Brackets are known
- Base categories are known
- Induce subcategories
- Clever split/merge category refinement



EM algorithm, like Forward-Backward for HMMs, but constrained by tree





# POS tag splits' commonest words: effectively a semantic class-based model

- Proper Nouns (NNP):

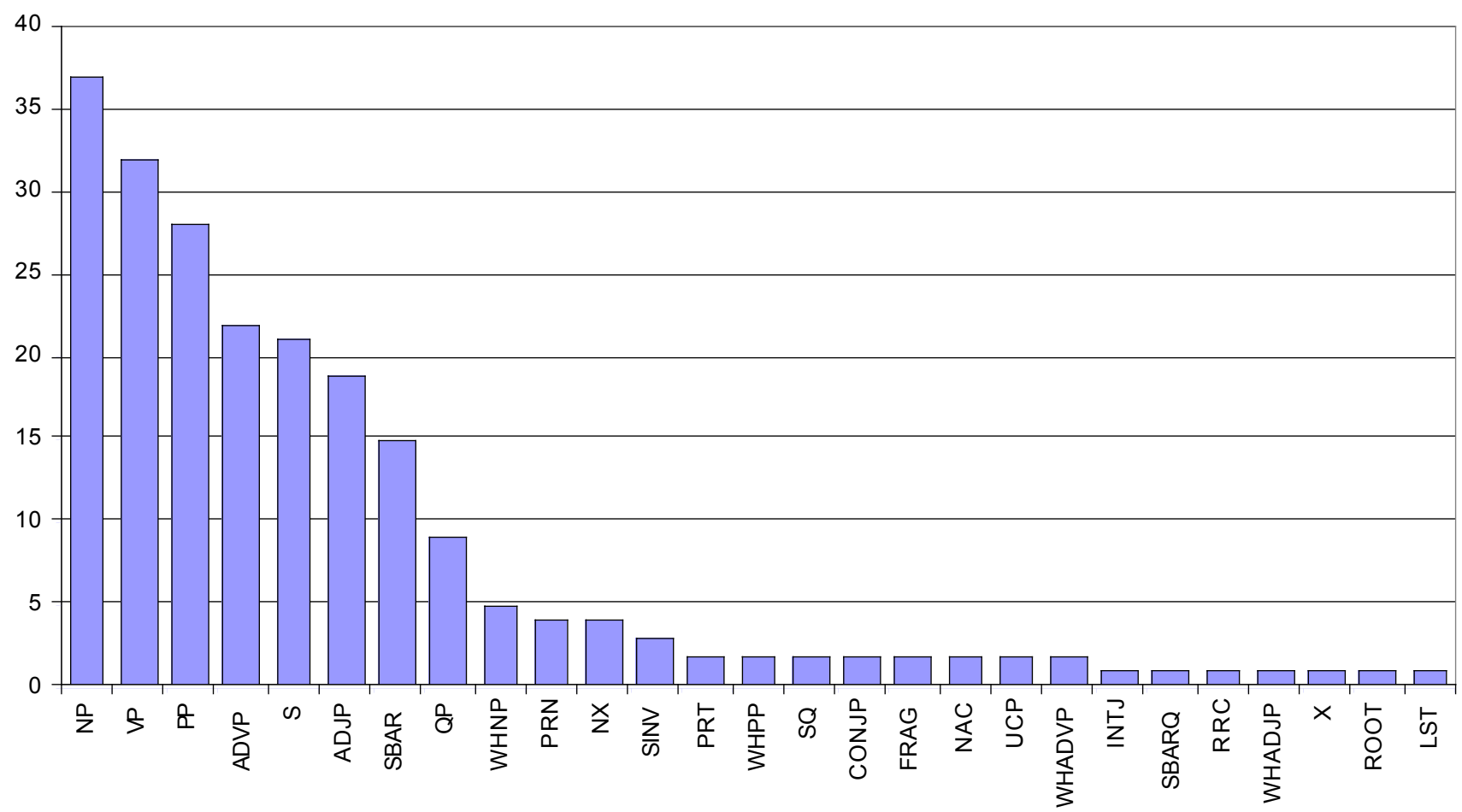
NNP-14	Oct.	Nov.	Sept.
NNP-12	John	Robert	James
NNP-2	J.	E.	L.
NNP-1	Bush	Noriega	Peters
NNP-15	New	San	Wall
NNP-3	York	Francisco	Street

- Personal pronouns (PRP):

PRP-0	It	He	I
PRP-1	it	he	they
PRP-2	it	them	him



# Number of phrasal subcategories





# The Latest Parsing Results... (English PTB3 WSJ train 2-21, test 23)

<i>Parser</i>	<i>F1 ≤ 40 words</i>	<i>F1 all words</i>
Klein & Manning unlexicalized 2003	86.3	85.7
Matsuzaki et al. simple EM latent states 2005	86.7	86.1
Charniak generative, lexicalized (“maxent inspired”) 2000	90.1	89.5
Petrov and Klein NAACL 2007	90.6	90.1
Charniak & Johnson discriminative reranker 2005	92.0	91.4
Fossum & Knight 2009 combining constituent parsers		<b>92.4</b>



# Latent Variable PCFGs

# Extending the idea to induced syntactico-semantic classes