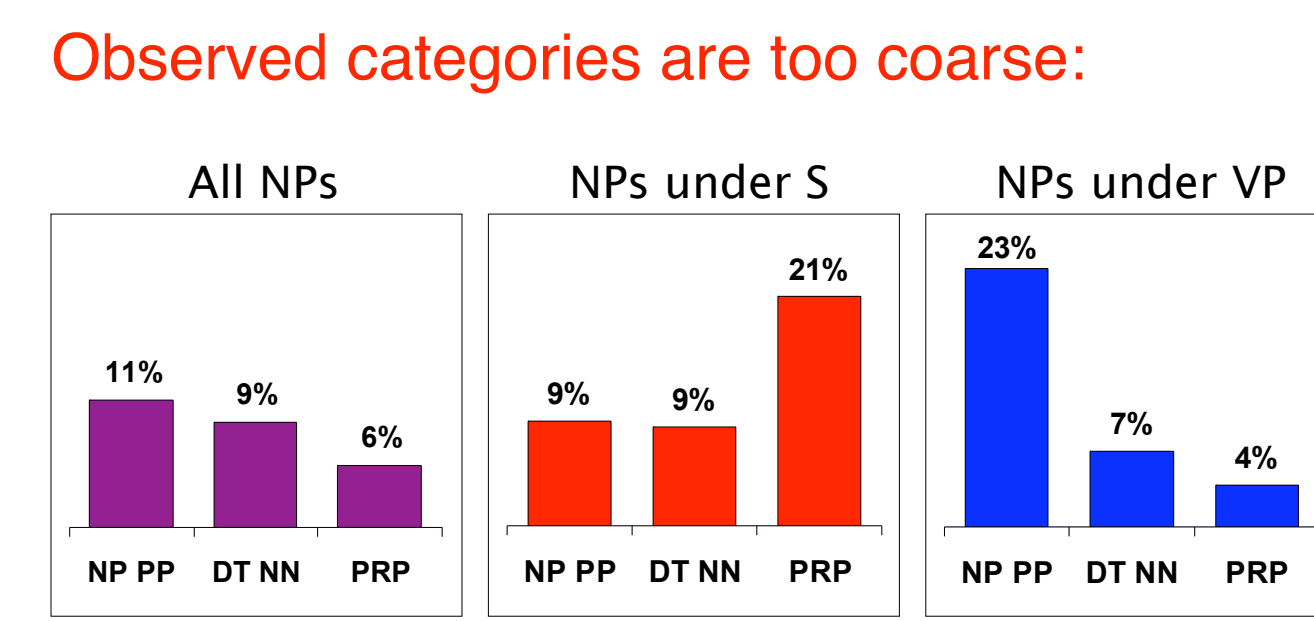
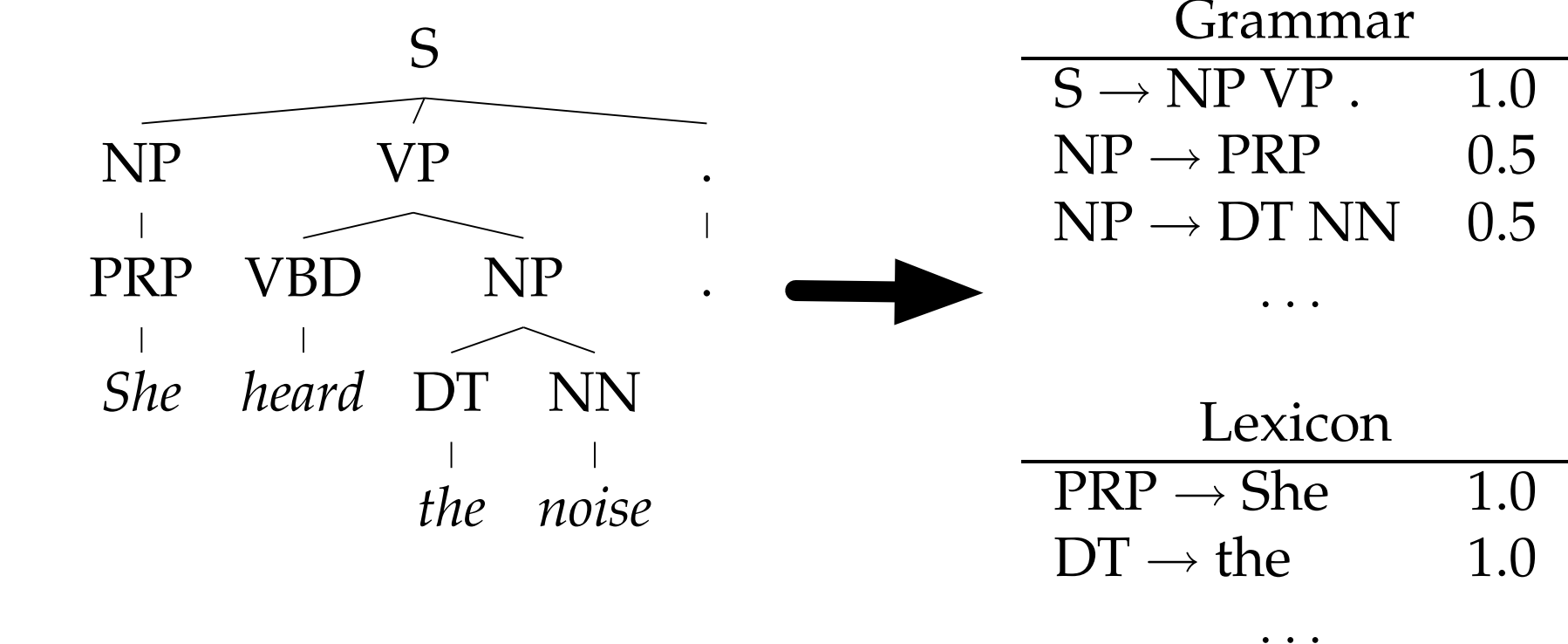


Learning and Inference for Hierarchically Split Probabilistic Context-Free Grammars

Slav Petrov and Dan Klein
University of California, Berkeley

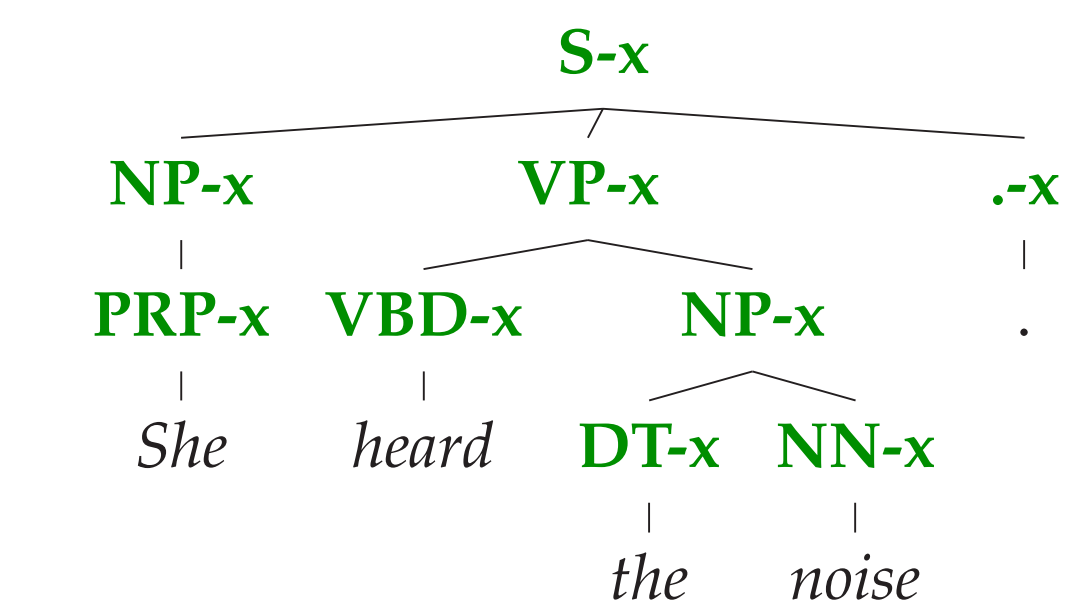
Learning

Grammar Extraction:



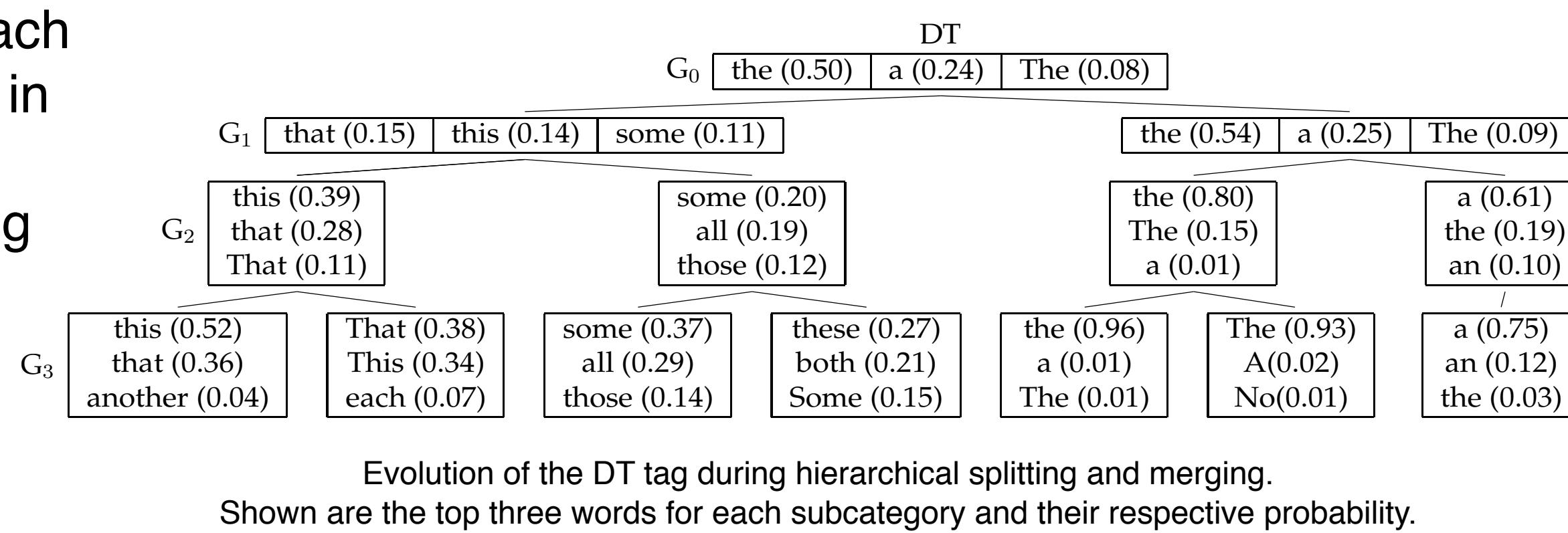
Adaptive Grammar Refinement:

Split each category in k subcategories and fit grammar with the EM algorithm.



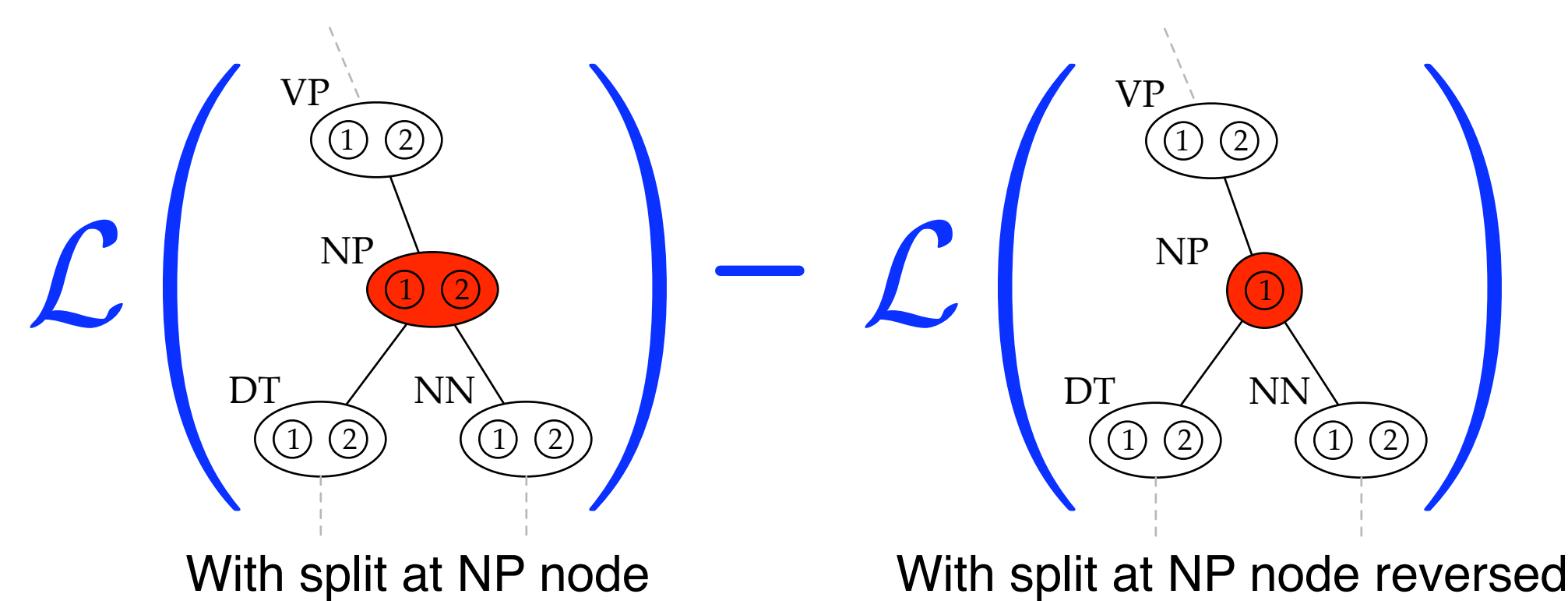
Hierarchical Splitting:

Repeatedly split each annotation symbol in two and retrain the grammar, initializing with the previous grammar.



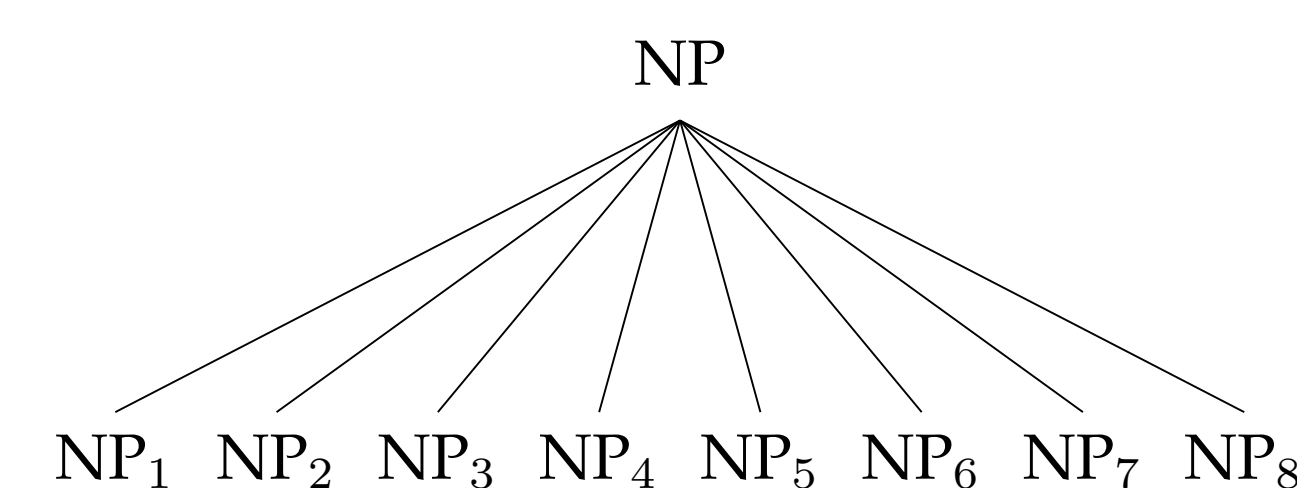
Merging:

Roll back the least useful splits in order to allocate complexity only where needed.



Smoothing:

Reduce overfitting by shrinking the productions of each subcategory towards their common base category.

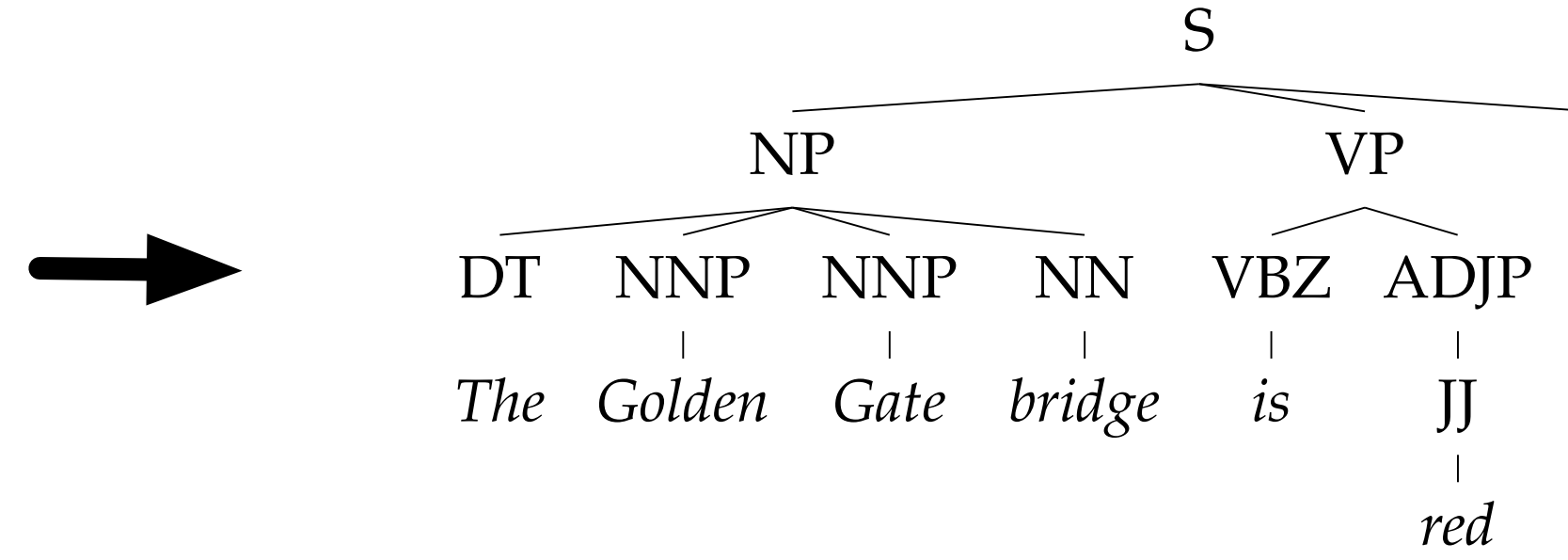


Reference:
Slav Petrov, Leon Barrett, Romain Thibaux and Dan Klein,
"Learning accurate, compact and interpretable tree annotation", in **ACL-COLING '06**

Inference

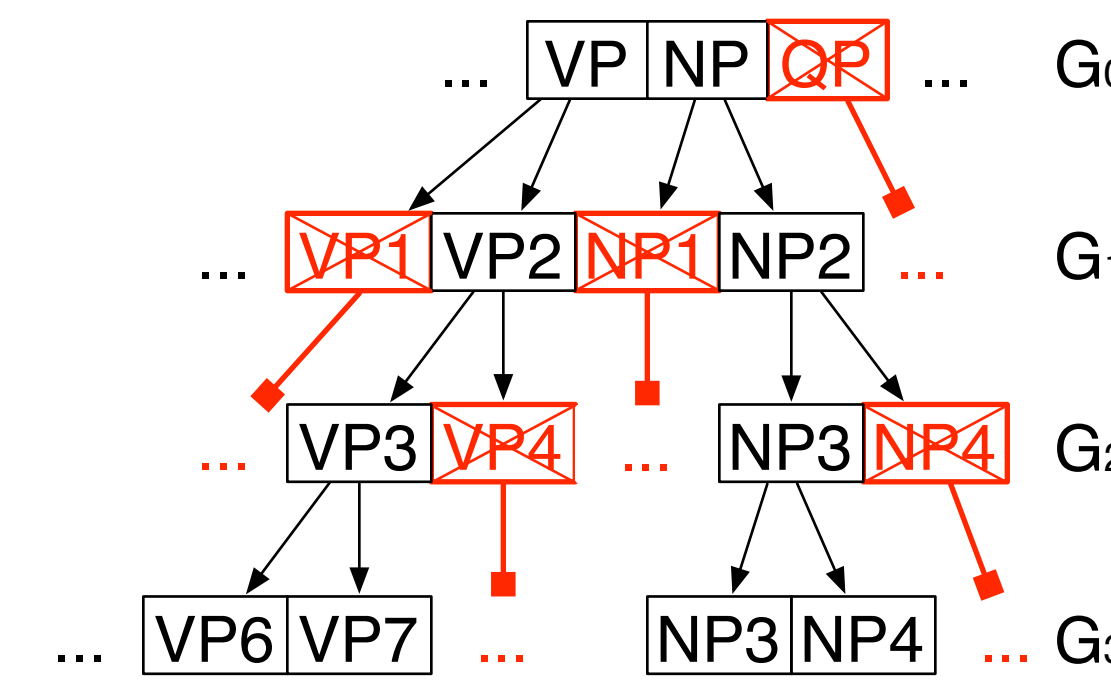
Parsing:

The Golden Gate bridge is red.

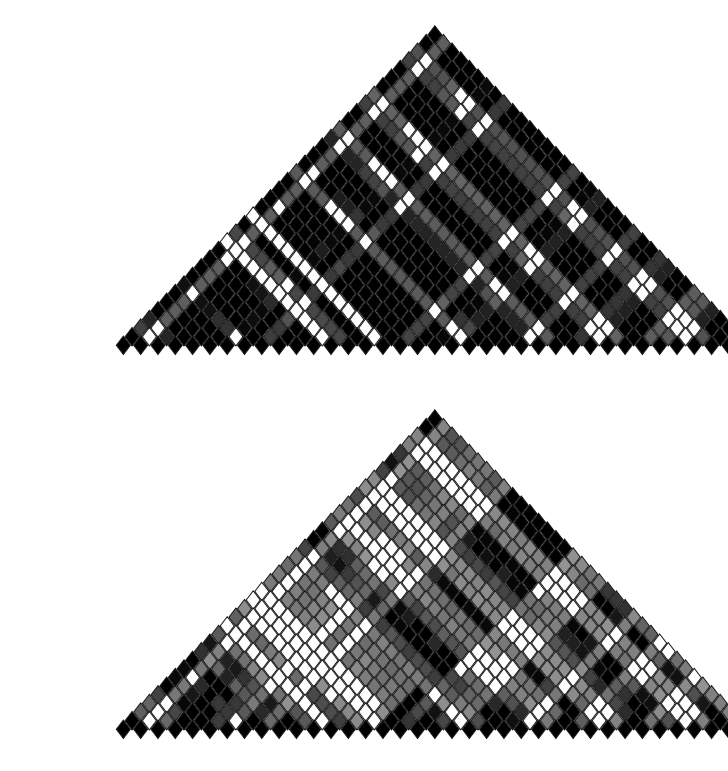


Parse Efficiency:

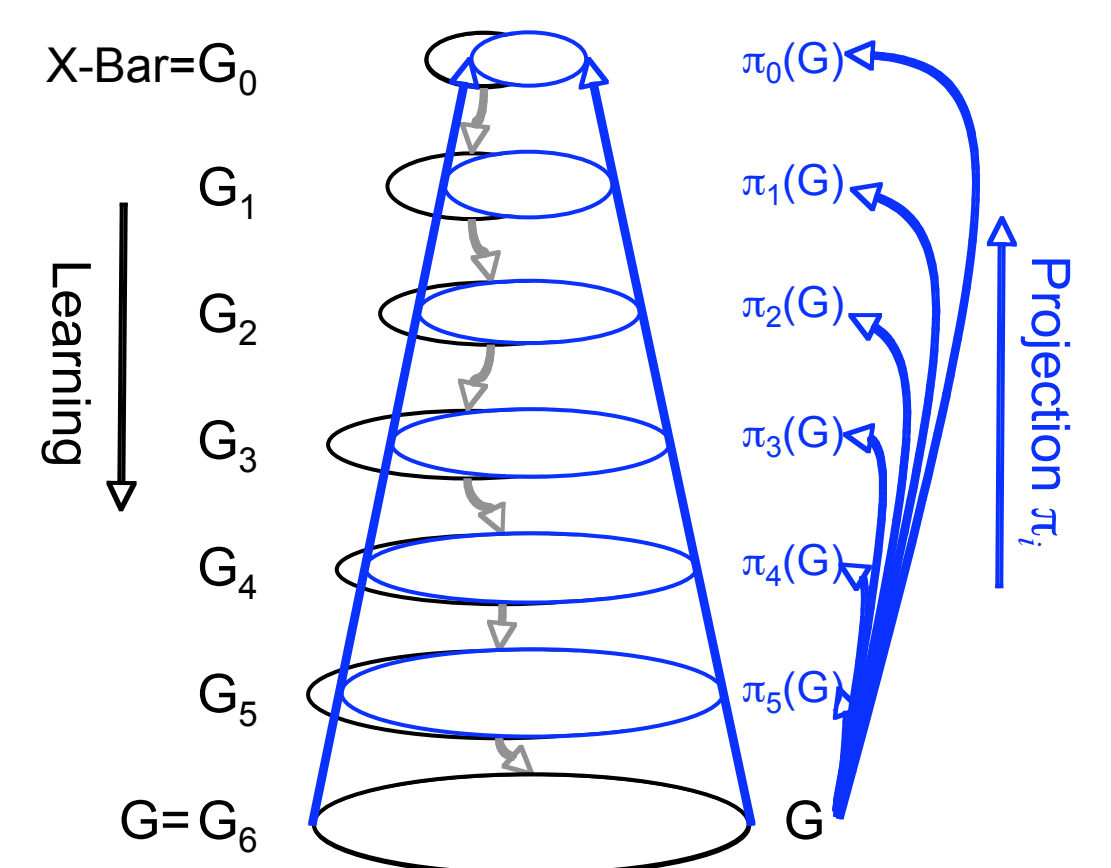
Rapidly pre-parse the sentence in a hierarchical coarse-to-fine fashion pruning away unlikely chart items.



Bracket posterior probabilities during coarse-to-fine decoding

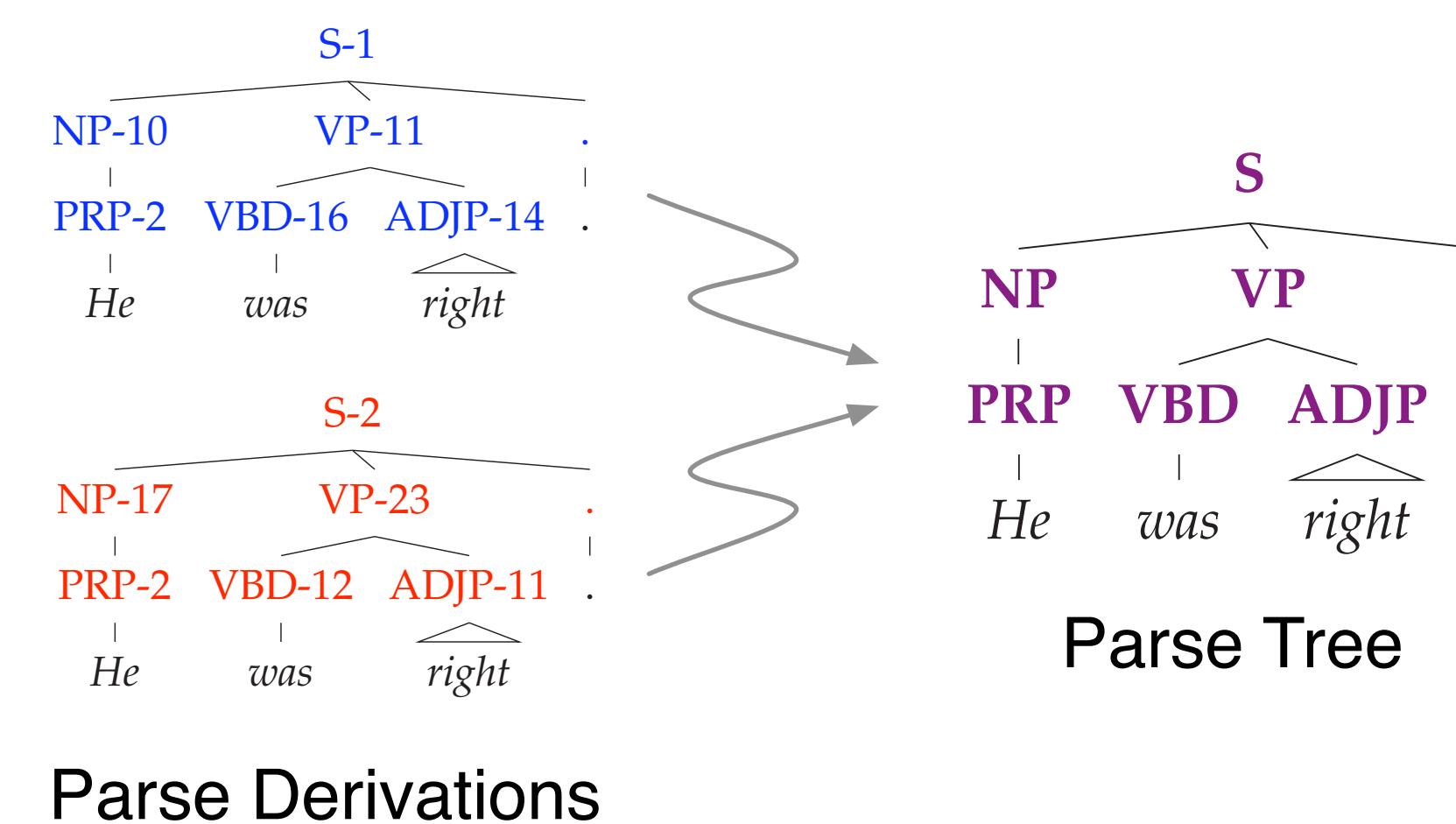


Compute grammars of intermediate complexity by projecting the most refined grammar.



Parse Selection:

Use a variational approximation to select the tree with the maximum number of expected correct rules (since computing the best parse tree is intractable and selecting the best derivation is a poor approximation).



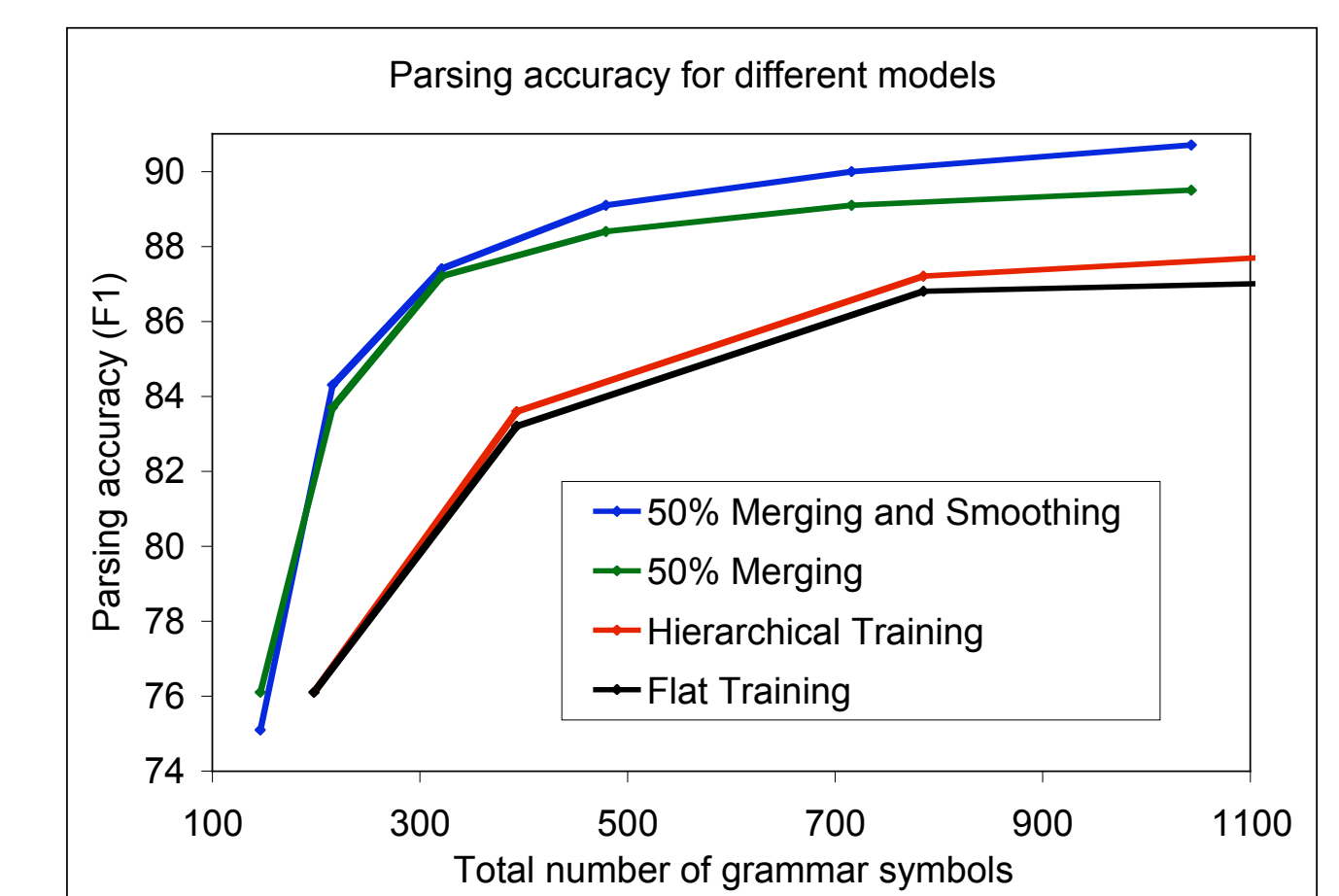
Reference:
Slav Petrov and Dan Klein,
"Improved Inference for Unlexicalized Parsing", in **NAACL-HLT '07**

Results

General technique for learning refined, structured models when only the trace of a complex underlying process is observed.
Learns compact and accurate grammars from a treebank without additional human input.
Gives best known parsing accuracy on a variety of languages, while being extremely efficient.

Interactive demo and download at <http://nlp.cs.berkeley.edu>

Parser	≤ 40 words F ₁	all F ₁
ENGLISH		
Charniak & Johnson '05	90.1	89.6
This Work	90.6	90.1
GERMAN		
Dubey '05	76.3	-
This Work	80.8	80.1
CHINESE		
Chiang & Bikel '02	80.0	76.6
This Work	86.3	83.4



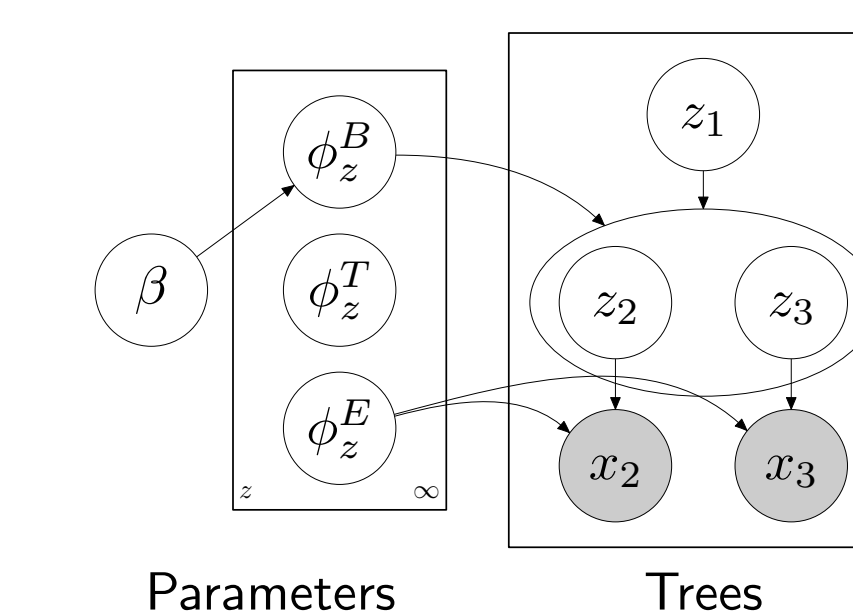
Learned grammars are compact and interpretable:

VBZ				DT				IN			
VBZ-0	gives	sells	takes	DT-0	the	The	a	IN-0	In	With	After
VBZ-1	comes	goes	works	DT-1	A	An	Another	IN-1	In	For	At
VBZ-2	includes	owns	is	DT-2	The	No	This	IN-2	in	for	on
VBZ-3	puts	provides	takes	DT-3	The	Same	These	IN-3	of	for	on
VBZ-4	says	adds	Says	DT-4	all	those	some	IN-4	from	on	with
VBZ-5	believes	means	thinks	DT-5	some	these	both	IN-5	at	for	by
VBZ-6	expects	makes	calls	DT-6	That	This	each	IN-6	by	in	with
VBZ-7	plans	expects	wants	DT-7	this	that	each	IN-7	for	with	on
VBZ-8	is	's	gets	DT-8	the	The	a	IN-8	If	White	As
VBZ-9	's	is	remains	DT-9	no	any	some	IN-9	because	if	That
VBZ-10	has	's	is	DT-10	an	a	the	IN-10	whether	if	That
VBZ-11	does	is	Does	DT-11	a	this	the	IN-11	that	like	whether
								IN-12	about	over	between
								IN-13	as	de	Up
								IN-14	than	ago	until
								IN-15	out	up	down
NNP				CD				RB			
NNP-0	Jr.	Goldman	INC.	CD-0	1	50	100	RB-0	recently	previously	still
NNP-1	Bush	Noriega	Peters	CD-1	8.50	15	1.2	RB-1	here	back	now
NNP-2	J.	F.	L.	CD-2	8	10	20	RB-2	very	highly	relatively
NNP-3	York	Francisco	Street	CD-3	1	30	31	RB-3	so	too	as
NNP-4	Inc	Exchange	Co	CD-4	1989	1990	1988	RB-4	also	now	still
NNP-5	Inc.	Corp.	Co.	CD-5	1988	1987	1990	RB-5	however	Now	However
NNP-6	Stock	Exchange	York	CD-6	two	three	five	RB-6	much	far	enough
NNP-7	Corp.	Inc.	Group	CD-7	one	One	Three	RB-7	even	well	then
NNP-8	Congress	Japan	IBM	CD-8	12	34	14	RB-8	as	about	nearly
NNP-9	Friday	September	August	CD-9	78	58	34	RB-9	only	just	almost
NNP-10	Shearson	D.	Ford	CD-10	one	two	three	RB-10	ago	earlier	later
NNP-11	U.S.	Treasury	Senate	CD-11	million	billion	trillion	RB-11	rather	instead	because
NNP-12	John	Robert	James					RB-12	back	close	ahead
NNP-13	Mr.	Ms.	President					RB-13	up	down	off
NNP-14	Oct.	Nov.	Sept.					RB-14	not	Not	maybe
NNP-15	New	San	Wall					RB-15	n't	not	also
JJS				RBR							
JJS-0	largest	latest	biggest	RBR-0	further	lower	higher				
JJS-1	least	best	worst	RBR-1	more	less	More				
JJS-2	most	Most	least	RBR-2	earlier	Earlier	later				

The most frequent three words in the subcategories of several part-of-speech tags.

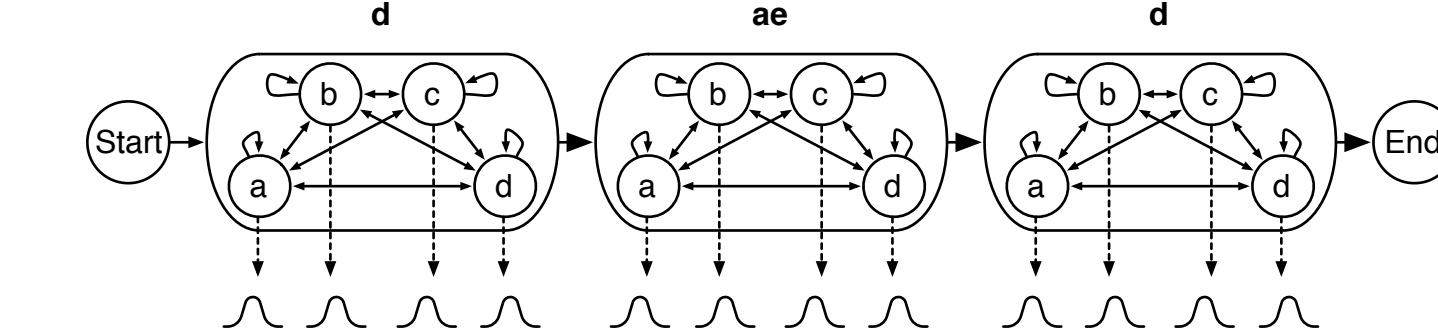
Extensions

Hierarchical Dirichlet Processes as a nonparametric Bayesian alternative to split and merge:



Reference:
Percy Liang, Slav Petrov, Michael Jordan and Dan Klein,
"The Infinite PCFG using Hierarchical Dirichlet Processes",
in **EMNLP-CoNLL '07**

Automatic refinement of acoustic models learns phone-internal structure as well as phone-external context:



Reference:
Slav Petrov, Adam Pauls and Dan Klein,
"Learning Structured Models for Phone Recognition",
in **EMNLP-CoNLL '07**