### **Advanced PCFG models**

**Computational Linguistics** 

Alexander Koller

5 December 2017

# The story so far

- Train PCFG with MLE on the Penn Treebank 02-21.
- Compute parse trees for PTB 23 using Viterbi-CKY.
- Trick against data sparseness in lexicon: delete words, train and test on sequences of POS tags.
- This yields labeled f-score in the low 70's.
  - Why so low?
  - How can we fix it?

## Fundamental problem of PCFGs

- Context-free grammar: One rule can only "see" parent and its children, not anything above or below.
- PCFG: Assumes statistical independence of all rewrite events.



 $NP \rightarrow PRP?$  $NP \rightarrow Det N?$ 

#### Independence assumptions



PTB statistics, from slides by Dan Klein

# Independence assumptions

- Accurate disambiguation of PP attachment requires lexical information.
  - I shot the elephant with a long trunk.
  - I shot the elephant with a long rifle.
- PP attachment influenced by choice of P.
  - Collins note: "workers dumped sacks into a bin"
  - *into*-PPs in PTB 9x more likely to attach to VP than to N
- PCFGs rely on nonterminals alone, cannot "see" lexical information.

### Directions

- Need to make nonterminals more informative to make PCFG rules sensitive to more context.
- Several approaches discussed today:
  - Johnson 98: Parent annotations
  - Collins 97: Lexicalized PCFGs
  - Klein & Manning 03: Unlexicalized PCFGs with nonterminals split by hand
  - Petrov & Klein 06: Unlexicalized PCFGs with automatically learned nonterminal splits

## Johnson 1998



- Discusses PTB preprocessing and impact of PTB representation changes.
- One key idea: *parent annotations*.
  - If parent of NP makes such a difference in how it should be expanded, why don't we encode the parent of the NP?
  - Replace nonterminal NP by NP^S (NP as child of S), NP^VP (NP as child of VP), and so on in PTB trees.
  - Train grammar on modified treebank.
     After parsing, remove annotations and compare to gold standard tree.

#### Example



Result: Labeled f-score on Section 22 jumps from 71.5 to 79.6. Number of production rules grows from 15,000 to 22,000.

# Lexicalized parsing



- Fundamental idea: If words are so important to distribution of rules, let's put them in the rules.
- Step 1: Mark each node in PTB with its *lexical head*.
  - identify head automatically using hand-written rules



### Lexicalized PCFGs

- Step 2: Read off lexicalized PCFG from treebank.
  - ▶ rules of the form S(examined)  $\rightarrow_2$  NP(lawyer) VP(examined)
  - "2" on arrow indicates that second child is head
- MLE and Viterbi-CKY adapt easily to new setting. So we're basically done!
- But! Number of rules multiplied by V<sup>r</sup> (V = vocabulary size, r = rank of rules).
  - ordinary rule × head word × heads of other children
  - increases number of parameters accordingly
  - astronomical sparse data problem

# Dealing with sparse data

- Horizontal Markovization:
  - break rules up into parts by generating children one by one
  - independence assumptions: child depends on limited context

![](_page_10_Figure_4.jpeg)

# **Dealing with sparse data**

- This helps a lot, but is still not enough for rare events.
- Need aggressive smoothing. Collins uses interpolation:
  - ▶  $p_1 = C(S \Rightarrow NP VP, H = examined) / C(S, H = examined)$
  - ▶  $p_0 = C(S \Rightarrow NP VP) / C(S)$
  - $P(S \rightarrow NP VP \mid S, examined) = \lambda p_1 + (1-\lambda) p_0$
  - $\bullet \ estimate \ \lambda \ from \ data$

Collins 1997 (with more complex lexicalization model): f-score 87.7 on PTB word strings of length  $\leq 40$ 

# Parsing speed

- Parsing slower than usual, because
  - grammar is much bigger
  - must be careful in managing head words
- Key insight: head word of (A,i,k) must be one of  $w_i, ..., w_{k-1}$ ; use pointers into input string.
  - ▶ this gives O(n<sup>5</sup>) parsing time with acceptable memory use
  - Eisner & Satta 99: can do it in O(n<sup>4</sup>) with clever algorithm
     still too slow in practice
  - use beam search to maintain only best hypotheses for each chart cell

# **Unlexicalized parsing**

![](_page_13_Picture_1.jpeg)

- Is lexicalization really as helpful as it seems?
  - Gildea 01: what counts is effect of head word on choice of subcategorization frame, not bilexical dependencies
  - Dubey & Keller 03: bilexical dependencies not useful when parsing German
  - Even lexicalized parsers (e.g. Collins 99, Charniak 00) make use of non-lexical splits of nonterminals.
- Klein & Manning 03: Perhaps usefulness of lexicalization is primarily in giving us more nonterminals? Can we get the same effect more cheaply?

#### Markovization

![](_page_14_Figure_1.jpeg)

Vertical Markovization: v = 2 is parent annotations v = 3 grandparent, etc.

			Horizon	ntal Mark	ov Order	
Vertical Order		h = 0	h = 1	$h \leq 2$	h = 2	$h = \infty$
v = 1	No annotation	71.27	72.5	73.46	72.96	72.62
		(854)	(3119)	(3863)	(6207)	(9657)
$v \leq 2$	Sel. Parents	74.75	77.42	77.77	77.50	76.91
		(2285)	(6564)	(7619)	(11398)	(14247)
v = 2	All Parents	74.68	77.42	77.81	77.50	76.81
		(2984)	(7312)	(8367)	(12132)	(14666)
$v \leq 3$	Sel. GParents	76.50	78.59	79.07	78.97	78.54
		(4943)	(12374)	(13627)	(19545)	(20123)
v = 3	All GParents	76.74	79.18	79.74	79.07	78.72
		(7797)	(15740)	(16994)	(22886)	(22002)

### **Rule-based state splitting**

![](_page_15_Figure_1.jpeg)

#### Results

	Cı	Indiv.		
Annotation	Size	F <sub>1</sub>	$\Delta F_1$	$\Delta F_1$
Baseline ( $v \le 2, h \le 2$ )	7619	77.77	—	—
UNARY-INTERNAL	8065	78.32	0.55	0.55
UNARY-DT	8066	78.48	0.71	0.17
UNARY-RB	8069	78.86	1.09	0.43
TAG-PA	8520	80.62	2.85	2.52
SPLIT-IN	8541	81.19	3.42	2.12
SPLIT-AUX	9034	81.66	3.89	0.57
SPLIT-CC	9190	81.69	3.92	0.12
SPLIT-%	9255	81.81	4.04	0.15
TMP-NP	9594	82.25	4.48	1.07
GAPPED-S	9741	82.28	4.51	0.17
POSS-NP	9820	83.06	5.29	0.28
SPLIT-VP	10499	85.72	7.95	1.36
BASE-NP	11660	86.04	8.27	0.73
DOMINATES-V	14097	86.91	9.14	1.42
RIGHT-REC-NP	15276	87.04	9.27	1.94

Compare against f-score 87-89 of lexicalized parsers. But much smaller grammars, simpler and faster parsing!

# State splitting

![](_page_17_Picture_1.jpeg)

• Can see all of these approaches as methods for refining the nonterminals of the PTB.

![](_page_17_Figure_3.jpeg)

• Petrov et al. 06: Can we automatically learn how to refine ("split") the nonterminals?

### **Split-Merge**

![](_page_18_Figure_1.jpeg)

#### Results

$\leq 40$ words	LP	LR	CB	0CB
Klein and Manning (2003)	86.9	85.7	1.10	60.3
Matsuzaki et al. (2005)	86.6	86.7	1.19	61.1
Collins (1999)	88.7	88.5	0.92	66.7
Charniak and Johnson (2005)	90.1	90.1	0.74	<b>70.1</b>
This Paper	90.3	90.0	0.78	68.5
all sentences	LP	LR	CB	0CB
all sentencesKlein and Manning (2003)	LP 86.3	LR 85.1	CB 1.31	0CB 57.2
all sentences Klein and Manning (2003) Matsuzaki et al. (2005)	LP 86.3 86.1	LR 85.1 86.0	CB 1.31 1.39	0CB 57.2 58.3
all sentences Klein and Manning (2003) Matsuzaki et al. (2005) Collins (1999)	LP 86.3 86.1 88.3	LR 85.1 86.0 88.1	CB 1.31 1.39 1.06	0CB 57.2 58.3 64.0
all sentences Klein and Manning (2003) Matsuzaki et al. (2005) Collins (1999) Charniak and Johnson (2005)	LP 86.3 86.1 88.3 89.5	LR 85.1 86.0 88.1 <b>89.6</b>	CB 1.31 1.39 1.06 <b>0.88</b>	0CB 57.2 58.3 64.0 <b>67.6</b>

("this paper" = Petrov et al. 06)

#### Some state-split POS tags

	VBZ				Γ	ЭТ					IN	
VBZ-0	gives	sells	takes	DT-0	the	The	а	] [	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	Some	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	а		IN-8	If	While	As
VBZ-9	's	is	remains	DT-9	no	any	some		IN-9	because	if	while
VBZ-10	has	's	is	DT-10	an	а	the		IN-10	whether	if	That
VBZ-11	does	Is	Does	DT-11	а	this	the		IN-11	that	like	whether
	N	INP			C	CD			IN-12	about	over	between
NNP-0	Jr.	Goldman	INC.	CD-0	1	50	100		IN-13	as	de	Up
NNP-1	Bush	Noriega	Peters	CD-1	8.50	15	1.2		IN-14	than	ago	until
NNP-2	J.	E.	L.	CD-2	8	10	20		IN-15	out	up	down
						20				RB		
NNP-3	York	Francisco	Street	CD-3	1	30	31				RB	
NNP-3 NNP-4	York Inc	Francisco Exchange	Street Co	CD-3 CD-4	1 1989	30 1990	31 1988	[	RB-0	recently	RB previously	still
NNP-3 NNP-4 NNP-5	York Inc Inc.	Francisco Exchange Corp.	Street Co Co.	CD-3 CD-4 CD-5	1 1989 1988	30 1990 1987	31 1988 1990		RB-0 RB-1	recently here	RB previously back	still now
NNP-3 NNP-4 NNP-5 NNP-6	York Inc Inc. Stock	Francisco Exchange Corp. Exchange	Street Co Co. York	CD-3 CD-4 CD-5 CD-6	1 1989 1988 two	30 1990 1987 three	31 1988 1990 five		RB-0 RB-1 RB-2	recently here very	RB previously back highly	still now relatively
NNP-3 NNP-4 NNP-5 NNP-6 NNP-7	York Inc Inc. Stock Corp.	Francisco Exchange Corp. Exchange Inc.	Street Co Co. York Group	CD-3 CD-4 CD-5 CD-6 CD-7	1 1989 1988 two one	30 1990 1987 three One	31 1988 1990 five Three		RB-0 RB-1 RB-2 RB-3	recently here very so	RB previously back highly too	still now relatively as
NNP-3 NNP-4 NNP-5 NNP-6 NNP-7 NNP-8	York Inc Inc. Stock Corp. Congress	Francisco Exchange Corp. Exchange Inc. Japan	Street Co Co. York Group IBM	CD-3 CD-4 CD-5 CD-6 CD-7 CD-8	1 1989 1988 two one 12	30 1990 1987 three One 34	31 1988 1990 five Three 14		RB-0 RB-1 RB-2 RB-3 RB-4	recently here very so also	RB previously back highly too now	still now relatively as still
NNP-3 NNP-4 NNP-5 NNP-6 NNP-7 NNP-8 NNP-9	York Inc Inc. Stock Corp. Congress Friday	Francisco Exchange Corp. Exchange Inc. Japan September	Street Co Co. York Group IBM August	CD-3 CD-4 CD-5 CD-6 CD-7 CD-8 CD-9	1 1989 1988 two one 12 78	30 1990 1987 three One 34 58	31 1988 1990 five Three 14 34		RB-0 RB-1 RB-2 RB-3 RB-4 RB-5	recently here very so also however	RB previously back highly too now Now	still now relatively as still However
NNP-3 NNP-4 NNP-5 NNP-6 NNP-7 NNP-7 NNP-8 NNP-9 NNP-10	York Inc Inc. Stock Corp. Congress Friday Shearson	Francisco Exchange Corp. Exchange Inc. Japan September D.	Street Co Co. York Group IBM August Ford	CD-3 CD-4 CD-5 CD-6 CD-7 CD-8 CD-9 CD-10	1 1989 1988 two one 12 78 one	30 1990 1987 three One 34 58 two	31 1988 1990 five Three 14 34 three		RB-0 RB-1 RB-2 RB-3 RB-4 RB-5 RB-6	recently here very so also however much	RB previously back highly too now Now far	still now relatively as still However enough
NNP-3 NNP-4 NNP-5 NNP-6 NNP-7 NNP-7 NNP-8 NNP-9 NNP-10 NNP-11	York Inc Inc. Stock Corp. Congress Friday Shearson U.S.	Francisco Exchange Corp. Exchange Inc. Japan September D. Treasury	Street Co Co. York Group IBM August Ford Senate	CD-3 CD-4 CD-5 CD-6 CD-7 CD-8 CD-9 CD-10 CD-11	1 1989 1988 two one 12 78 one million	30 1990 1987 three One 34 58 two billion	31 1988 1990 five Three 14 34 three trillion		RB-0 RB-1 RB-2 RB-3 RB-4 RB-5 RB-6 RB-7	recently here very so also however much even	RB previously back highly too now Now far well	still now relatively as still However enough then
NNP-3 NNP-4 NNP-5 NNP-6 NNP-7 NNP-7 NNP-8 NNP-9 NNP-10 NNP-11 NNP-12	York Inc Inc. Stock Corp. Congress Friday Shearson U.S. John	Francisco Exchange Corp. Exchange Inc. Japan September D. Treasury Robert	Street Co Co. York Group IBM August Ford Senate James	CD-3 CD-4 CD-5 CD-6 CD-7 CD-8 CD-9 CD-10 CD-11	1 1989 1988 two one 12 78 one million	30 1990 1987 three One 34 58 two billion RP	31 1988 1990 five Three 14 34 three trillion		RB-0 RB-1 RB-2 RB-3 RB-4 RB-5 RB-5 RB-6 RB-7 RB-8	recently here very so also however much even as	RB previously back highly too now Now far well about	still now relatively as still However enough then nearly
NNP-3 NNP-4 NNP-5 NNP-6 NNP-7 NNP-7 NNP-8 NNP-9 NNP-10 NNP-10 NNP-11 NNP-12 NNP-13	York Inc Inc. Stock Corp. Congress Friday Shearson U.S. John Mr.	Francisco Exchange Corp. Exchange Inc. Japan September D. Treasury Robert Ms.	Street Co Co. York Group IBM August Ford Senate James President	CD-3 CD-4 CD-5 CD-6 CD-7 CD-8 CD-9 CD-10 CD-11 PRP-0	1 1989 1988 two one 12 78 one million	30 1990 1987 three One 34 58 two billion RP He	31 1988 1990 five Three 14 34 three trillion		RB-0 RB-1 RB-2 RB-3 RB-4 RB-5 RB-6 RB-7 RB-8 RB-9	recently here very so also however much even as only	RB previously back highly too now Now far well about just	still now relatively as still However enough then nearly almost
NNP-3 NNP-4 NNP-5 NNP-6 NNP-7 NNP-7 NNP-7 NNP-8 NNP-9 NNP-10 NNP-10 NNP-11 NNP-12 NNP-13 NNP-14	York Inc Inc. Stock Corp. Congress Friday Shearson U.S. John Mr. Oct.	Francisco Exchange Corp. Exchange Inc. Japan September D. Treasury Robert Ms. Nov.	Street Co Co. York Group IBM August Ford Senate James President Sept.	CD-3 CD-4 CD-5 CD-6 CD-7 CD-8 CD-9 CD-10 CD-11 PRP-0 PRP-1	1 1989 1988 two one 12 78 one million Pl It it	30 $1990$ $1987$ three One $34$ $58$ two billion RP He he	31 1988 1990 five Three 14 34 three trillion		RB-0 RB-1 RB-2 RB-3 RB-4 RB-5 RB-6 RB-7 RB-8 RB-9 RB-10	recently here very so also however much even as only ago	RB previously back highly too now Now far well about just earlier	still now relatively as still However enough then nearly almost later
NNP-3 NNP-4 NNP-5 NNP-6 NNP-7 NNP-7 NNP-7 NNP-8 NNP-9 NNP-10 NNP-10 NNP-11 NNP-12 NNP-13 NNP-14 NNP-15	York Inc Inc. Stock Corp. Congress Friday Shearson U.S. John Mr. Oct. New	Francisco Exchange Corp. Exchange Inc. Japan September D. Treasury Robert Ms. Nov. San	Street Co Co. York Group IBM August Ford Senate James President Sept. Wall	CD-3 CD-4 CD-5 CD-6 CD-7 CD-8 CD-9 CD-10 CD-11 PRP-0 PRP-1 PRP-2	1 1989 1988 two one 12 78 one million P It it it	30 1990 1987 three One 34 58 two billion RP He he them	31 1988 1990 five Three 14 34 three trillion I they him		RB-0 RB-1 RB-2 RB-3 RB-4 RB-5 RB-6 RB-7 RB-8 RB-9 RB-10 RB-11	recently here very so also however much even as only ago rather	RB previously back highly too now Now far well about just earlier instead	still now relatively as still However enough then nearly almost later because
NNP-3 NNP-4 NNP-5 NNP-6 NNP-7 NNP-7 NNP-7 NNP-8 NNP-9 NNP-10 NNP-10 NNP-11 NNP-12 NNP-13 NNP-14 NNP-15	York Inc Inc. Stock Corp. Congress Friday Shearson U.S. John Mr. Oct. New	Francisco Exchange Corp. Exchange Inc. Japan September D. Treasury Robert Ms. Nov. San	Street Co Co. York Group IBM August Ford Senate James President Sept. Wall	CD-3 CD-4 CD-5 CD-6 CD-7 CD-8 CD-9 CD-10 CD-11 PRP-0 PRP-1 PRP-2	1 1989 1988 two one 12 78 one million Pl It it it R	30 1990 1987 three One 34 58 two billion RP He he them BR	31 1988 1990 five Three 14 34 three trillion I they him		RB-0 RB-1 RB-2 RB-3 RB-4 RB-5 RB-5 RB-6 RB-7 RB-8 RB-9 RB-10 RB-11 RB-12	recently here very so also however much even as only ago rather back	RB previously back highly too now Now far well about just earlier instead close	still now relatively as still However enough then nearly almost later because ahead
NNP-3 NNP-4 NNP-5 NNP-6 NNP-7 NNP-8 NNP-9 NNP-10 NNP-10 NNP-10 NNP-11 NNP-12 NNP-13 NNP-13 NNP-14 NNP-15	York Inc Inc. Stock Corp. Congress Friday Shearson U.S. John Mr. Oct. New	Francisco Exchange Corp. Exchange Inc. Japan September D. Treasury Robert Ms. Nov. San JJS latest	Street Co Co. York Group IBM August Ford Senate James President Sept. Wall biggest	CD-3 CD-4 CD-5 CD-6 CD-7 CD-8 CD-9 CD-10 CD-11 PRP-0 PRP-1 PRP-2 RBR-0	l 1989 1988 two one 12 78 one million Pl It it it it R further	30 1990 1987 three One 34 58 two billion RP He he them BR lower	31 1988 1990 five Three 14 34 three trillion I they him		RB-0 RB-1 RB-2 RB-3 RB-4 RB-5 RB-6 RB-7 RB-8 RB-9 RB-10 RB-11 RB-12 RB-13	recently here very so also however much even as only ago rather back up	RB previously back highly too now Now far well about just earlier instead close down	still now relatively as still However enough then nearly almost later because ahead off
NNP-3 NNP-4 NNP-5 NNP-6 NNP-7 NNP-7 NNP-7 NNP-8 NNP-9 NNP-10 NNP-10 NNP-10 NNP-11 NNP-12 NNP-13 NNP-13 NNP-14 NNP-15	York Inc Inc. Stock Corp. Congress Friday Shearson U.S. John Mr. Oct. New	Francisco Exchange Corp. Exchange Inc. Japan September D. Treasury Robert Ms. Nov. San JJS latest best	Street Co Co. York Group IBM August Ford Senate James President Sept. Wall biggest worst	CD-3 CD-4 CD-5 CD-6 CD-7 CD-8 CD-9 CD-10 CD-11 PRP-0 PRP-1 PRP-2 RBR-0 RBR-1	l 1989 1988 two one 12 78 one million P It it it it further more	30 1990 1987 three One 34 58 two billion RP He he them BR lower less	31 1988 1990 five Three 14 34 three trillion I they him higher More		RB-0 RB-1 RB-2 RB-3 RB-4 RB-5 RB-6 RB-7 RB-8 RB-9 RB-10 RB-11 RB-12 RB-13 RB-14	recently here very so also however much even as only ago rather back up not	RB previously back highly too now Now far well about just earlier instead close down Not	still now relatively as still However enough then nearly almost later because ahead off maybe

### Summary

- PCFGs that we read off of treebank suffer from overly strong independence assumptions.
- Improve parser accuracy by encoding context in nonterminal vocabulary.
  - parent annotations
  - lexicalization
  - rule-based and automatically computed state splitting
- Berkeley parser: f-score around 90.