

[illegible]

# The Return of Unlexicalized PCFGs



# Accurate Unlexicalized Parsing

[Klein and Manning 1993]

- What do we mean by an “unlexicalized” PCFG?
  - Grammar rules are not systematically specified down to the level of lexical items
    - NP-stocks is not allowed
    - NP<sup>S</sup>-CC is fine
  - Closed vs. open class words
    - Long tradition in linguistics of using function words as features or markers for selection (VB-have, SBAR-if/whether)
    - Different to the bilexical idea of semantic heads
    - Open-class selection is really a proxy for semantics
- Thesis
  - Most of what you need for accurate parsing, and much of what lexicalized PCFGs actually capture *isn't* lexical selection between content words but just basic grammatical features, like verb form, finiteness, presence of a verbal auxiliary, etc.



# Experimental Approach

- Corpus: Penn Treebank, WSJ; iterate on small dev set



Training: sections 02-21

Development: section 22 (first 20 files) ←

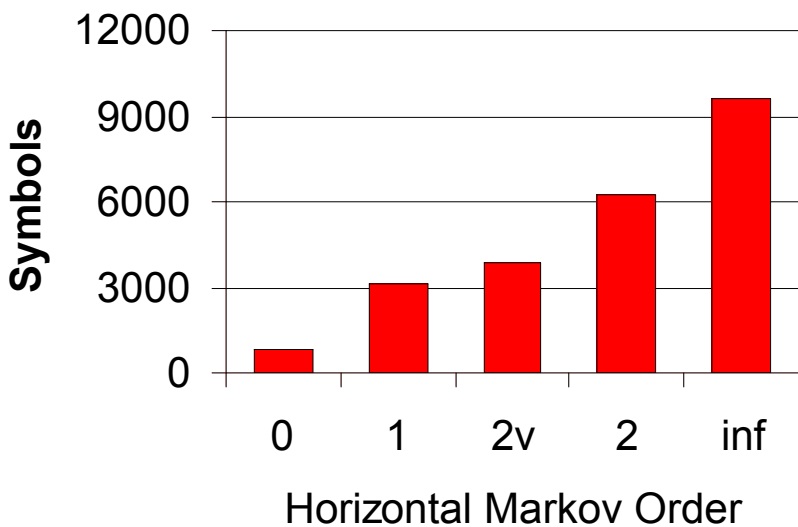
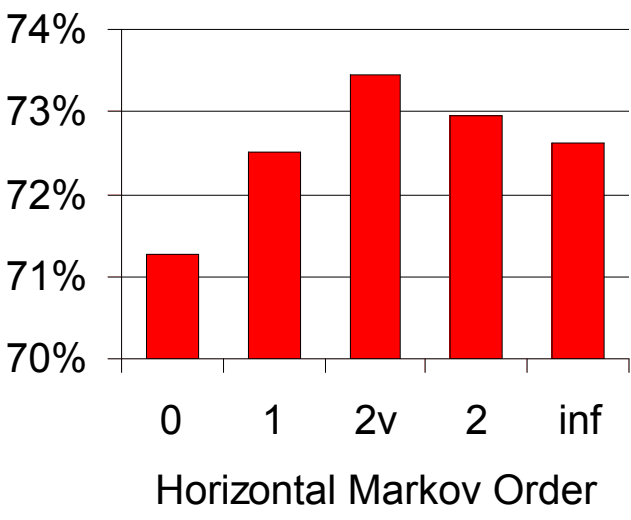
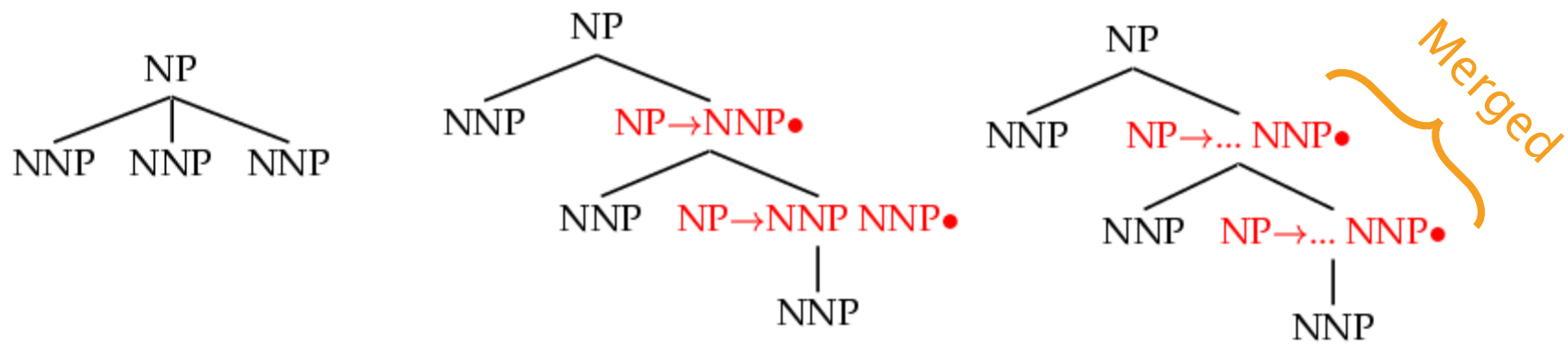
Test: section 23

- Size – number of symbols in grammar.
  - Passive / complete symbols: NP, NP<sup>S</sup>
  - Active / incomplete symbols: @NP\_NP\_CC [from binarization]
- We state-split as sparingly as possible
  - Highest accuracy with fewest symbols
  - Error-driven, manual hill-climb, one annotation at a time



# Horizontal Markovization

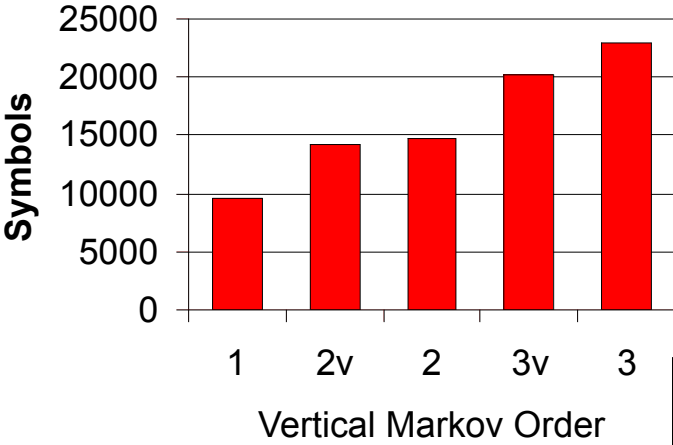
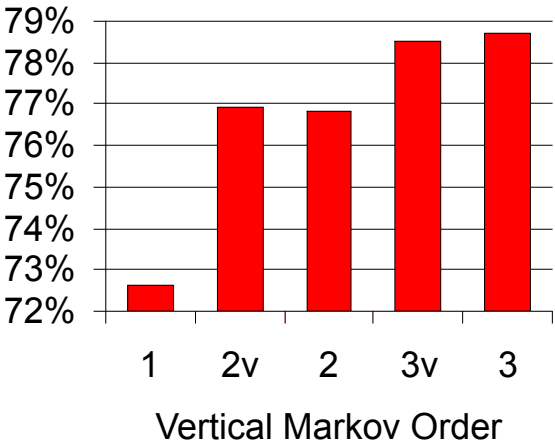
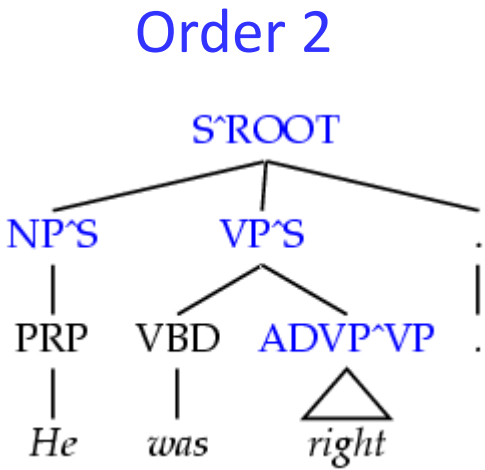
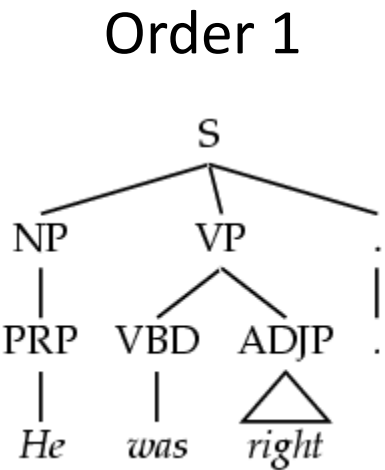
- Horizontal Markovization: Merges States





# Vertical Markovization

- Vertical Markov order: rewrites depend on past  $k$  ancestor nodes. (i.e., parent annotation)

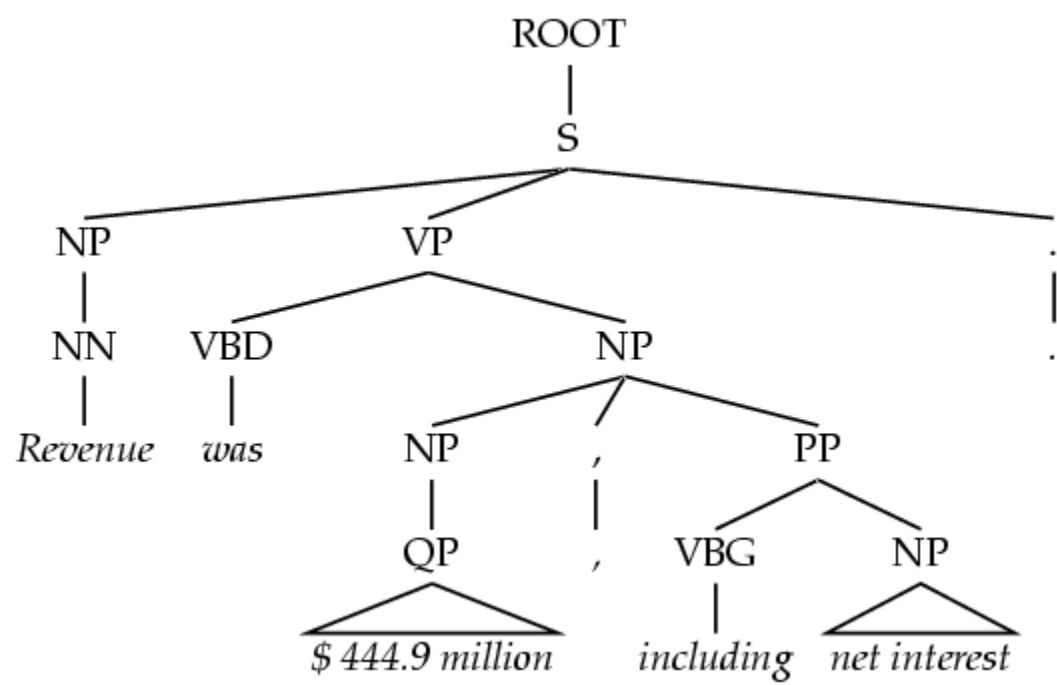


Model	F1	Size
v=h=2v	77.8	7.5K



# Unary Splits

- Problem: unary rewrites are used to transmute categories so a high-probability rule can be used.
- Solution: Mark unary rewrite sites with -U

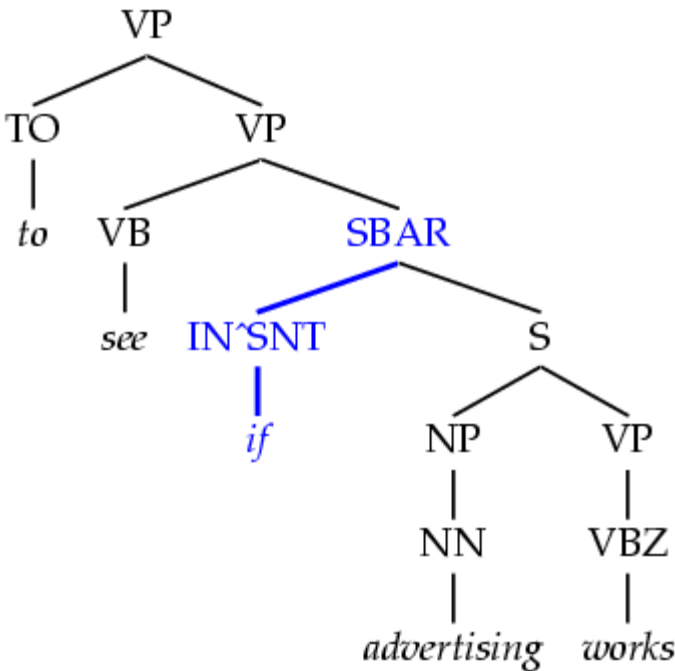


Annotation	F1	Size
Base	77.8	7.5K
UNARY	78.3	8.0K



# Tag Splits

- Problem: Treebank tags are too coarse.
- Example: SBAR sentential complementizers (*that*, *whether*, *if*), subordinating conjunctions (*while*, *after*), and true prepositions (*in*, *of*, *to*) are all tagged IN.
- Partial Solution:
  - Subdivide the IN tag.

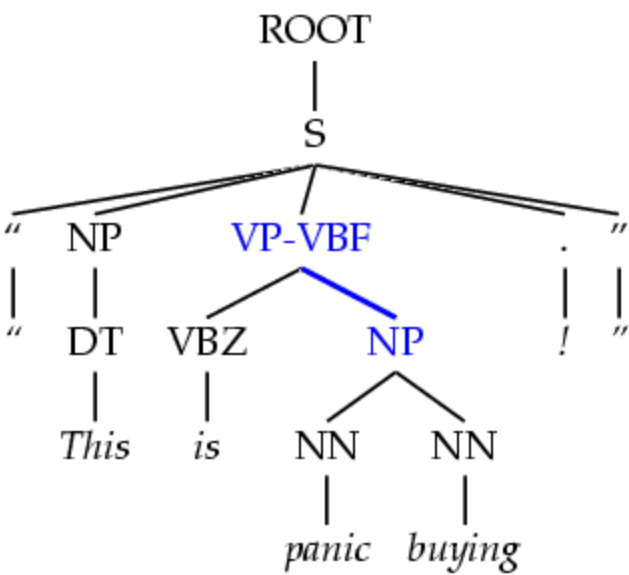


Annotation	F1	Size
Previous	78.3	8.0K
SPLIT-IN	80.3	8.1K



# Yield Splits

- Problem: sometimes the behavior of a category depends on something inside its future yield.
- Examples:
  - Possessive NPs
  - Finite vs. infinite VPs
  - Lexical heads!
- Solution: annotate future elements into nodes.



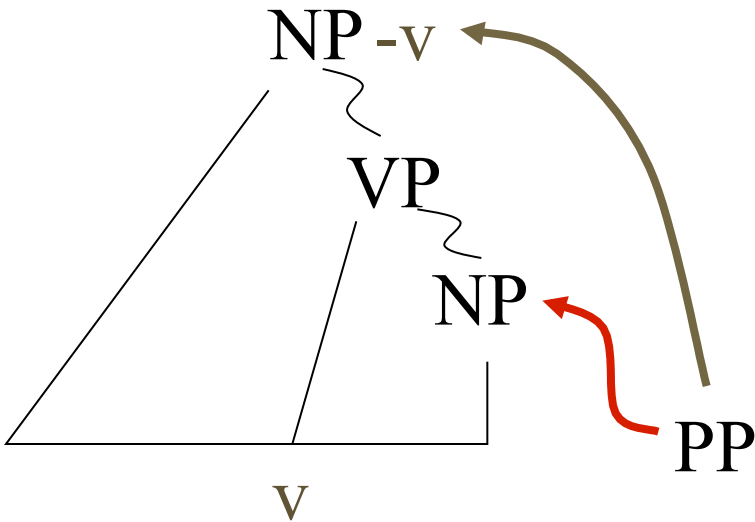
Annotation	F1	Size
tag splits	82.3	9.7K
POSS-NP	83.1	9.8K
SPLIT-VP	85.7	10.5K





# Distance / Recursion Splits

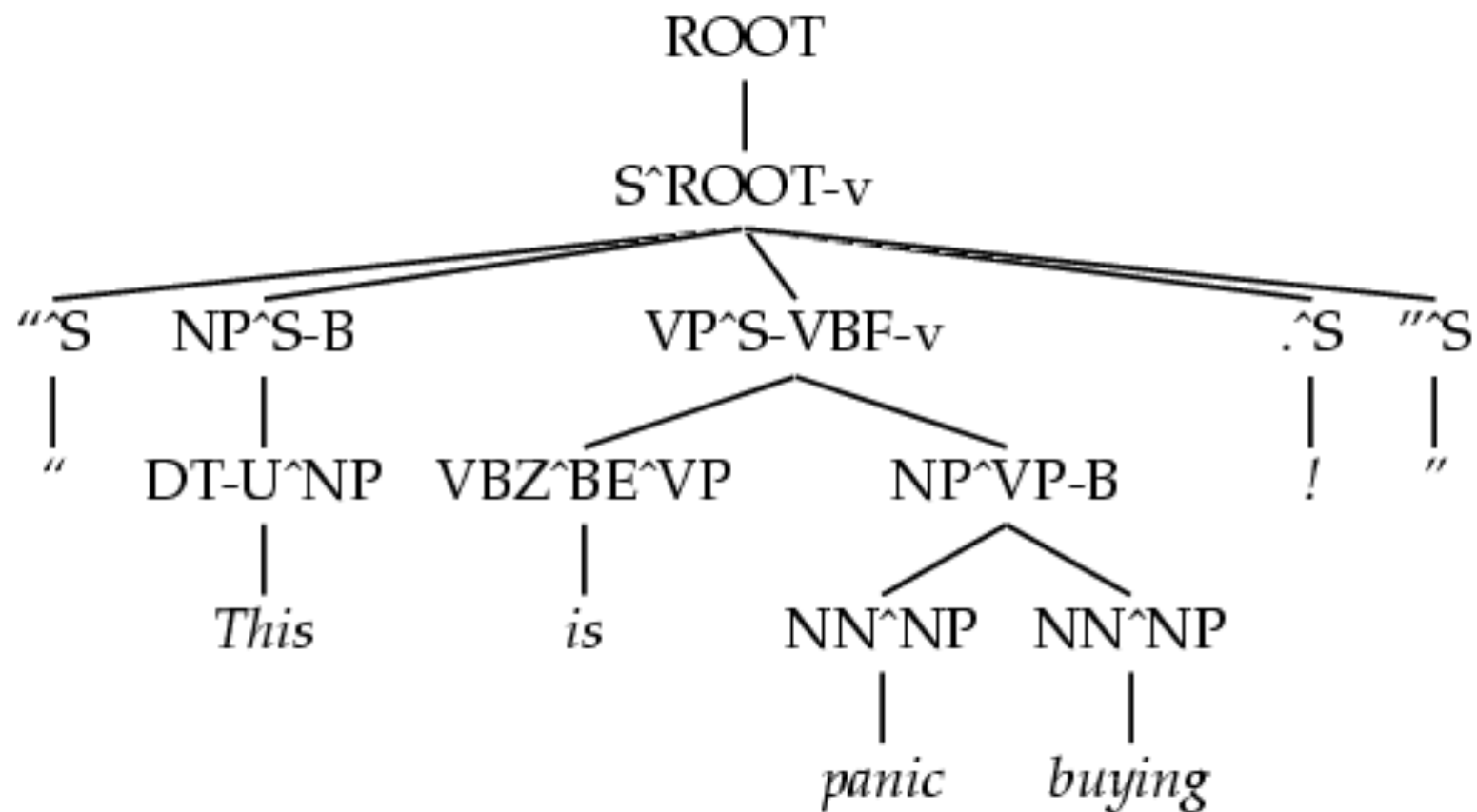
- Problem: vanilla PCFGs cannot distinguish attachment heights.
- Solution: mark a property of higher or lower sites:
  - Contains a verb.
  - Is (non)-recursive.
  - Base NPs [cf. Collins 99]
  - Right-recursive NPs



Annotation	F1	Size
Previous	85.7	10.5K
BASE-NP	86.0	11.7K
DOMINATES-V	86.9	14.1K
RIGHT-REC-NP	87.0	15.2K



# A Fully Annotated Tree





# Final Test Set Results

Parser	LP	LR	F1
Magerman 95	84.9	84.6	<b>84.7</b>
Collins 96	86.3	85.8	<b>86.0</b>
Klein & Manning 03	86.9	85.7	<b>86.3</b>
Charniak 97	87.4	87.5	<b>87.4</b>
Collins 99	88.7	88.6	<b>88.6</b>

- Beats “first generation” lexicalized parsers

# The Return of Unlexicalized PCFGs