# Polynomial time parsing of PCFGs



Gerald Penn

(some slides from Pi-Chuan Chang and Christopher Manning)
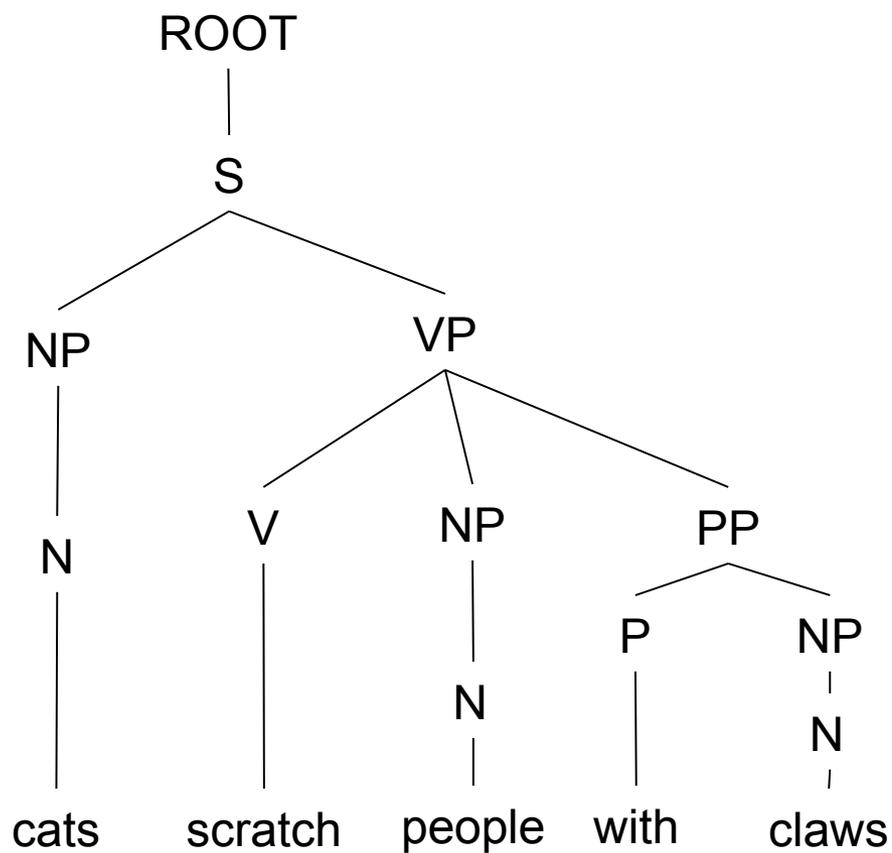
# 0. Chomsky Normal Form

- All rules are of the form $X \rightarrow Y\ Z$ or $X \rightarrow w$.
- A transformation to this form doesn't change the weak generative capacity of CFGs.
  - With some extra book-keeping in symbol names, you can even reconstruct the same trees with a detransform
  - Unaries/empties are removed recursively
  - $n$-ary rules introduce new nonterminals ($n > 2$)
    - $VP \rightarrow V\ NP\ PP$  becomes  $VP \rightarrow V\ @VP\text{-}V$  and  $@VP\text{-}V \rightarrow NP\ PP$
- In practice it's a pain
  - Reconstructing $n$-aries is easy
  - Reconstructing unaries can be trickier
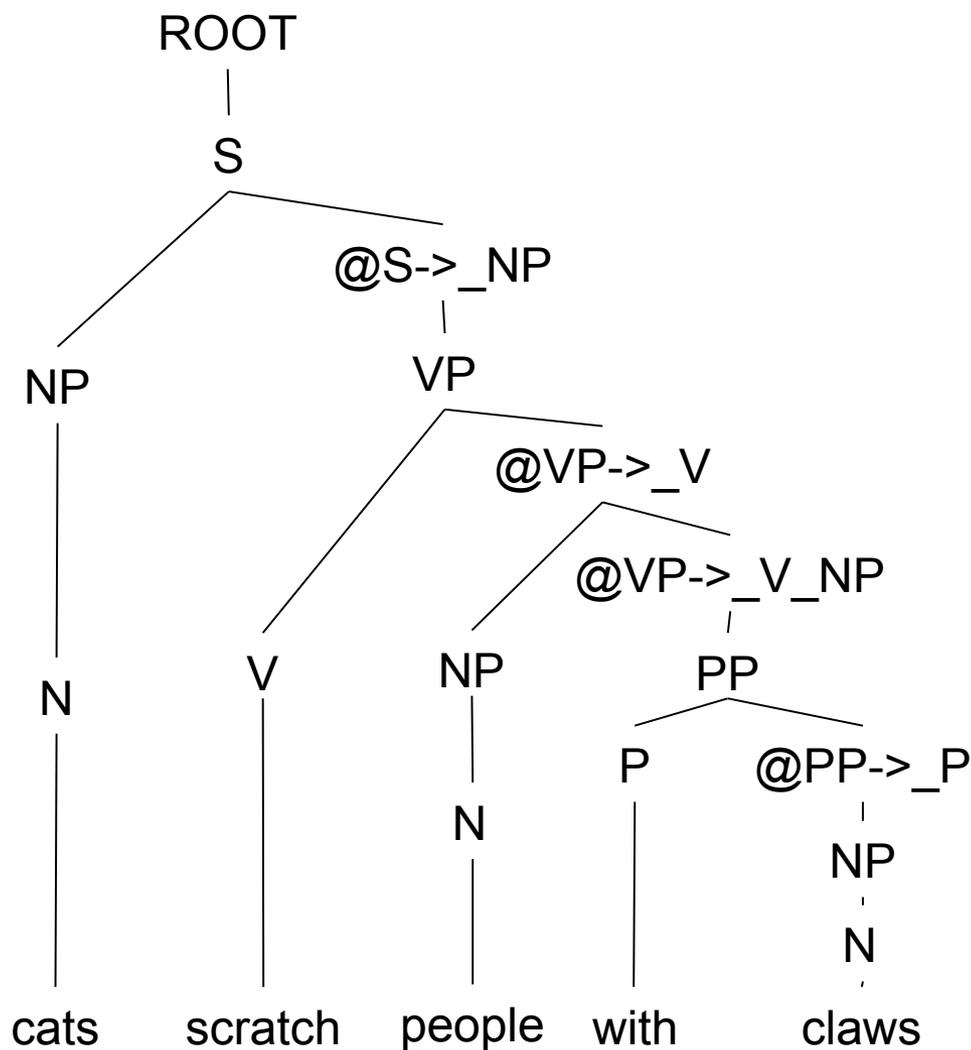- But it makes parsing easier/more efficient
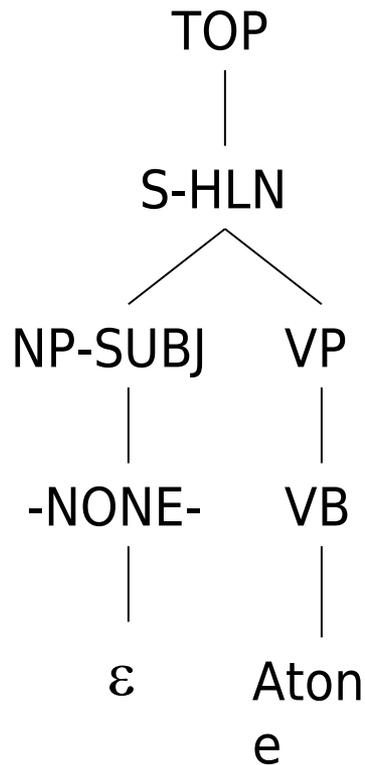
# An example: before binarization…

# After binarization…

ROOT
- S
  - NP
    - N
      - cats
  - @S->_NP
    - VP
      - V
        - scratch
      - @VP->_V
        - NP
          - N
            - people
        - @VP->_V_NP
          - PP
            - P
              - with
            - @PP->_P
              - NP
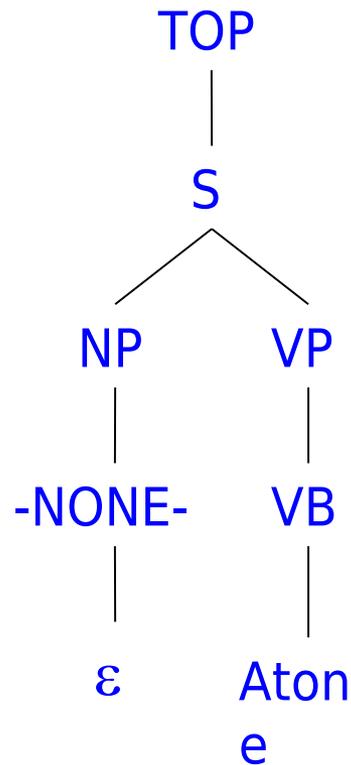                - N
                  - claws

# Treebank: empties and unaries



PTB Tree     NoFuncTags     NoEmpties     High     Low

NoUnaries

# Constituency Parsing



**PCFG**

**Rule Probs $\theta_i$**

$\theta_0$: S → NP VP

$\theta_1$: NP → NN NNS

...

$\theta_{42}$: NN→Factory

$\theta_{43}$: NNS→payrolls

...

# 1. Cocke-Kasami-Younger (CKY) Constituency Parsing



Factory payrolls  fell     in   September

# Viterbi (Max) Scores

NP $1.87 \times 10^{-7}$

NN    0.0023
NNP    0.001

NNS    0.0014

Factory          payrolls

NP→NN NNS  0.13
$i_{NP} =$
$(0.13)(0.0023)(0.0014)$
$= 1.87 \times 10^{-7}$

NP→NNP NNS 0.056
$i_{NP} =$
$(0.056)(0.001)(0.0014)$
$= 7.84 \times 10^{-8}$

# Extended CKY parsing

- Unaries can be incorporated into the algorithm
  - Messy, but doesn't increase algorithmic complexity
- Empties can be incorporated
  - Use fenceposts
  - Doesn't increase complexity; essentially like unaries

- Binarization is *vital*
  - All sorts of optimizations depend on this
    - Binarization may be an explicit transformation or implicit in how the parser works (Early-style dotted rules), but it's almost always there.

# The CKY algorithm (1960/1965) ... generalized

```
function CKY(words, grammar) returns most probable parse/prob
  score = new double[#(words)+1][#(words)+][#(nonterms)]
  back = new Pair[#(words)+1][#(words)+1][#nonterms]]
  for i=0; i<#(words); i++
    for A in nonterms
      if A -> words[i] in grammar
        score[i][i+1][A] = P(A -> words[i])
    //handle unaries
    boolean added = true
    while added
      added = false
      for A, B in nonterms
        if score[i][i+1][B] > 0 && A->B in grammar
          prob = P(A->B)*score[i][i+1][B]
          if(prob > score[i][i+1][A])
            score[i][i+1][A] = prob
            back[i][i+1] [A] = B
            added = true
```

# The CKY algorithm (1960/1965) ... generalized

```
for span = 2 to #(words)
  for begin = 0 to #(words)- span
    end = begin + span
    for split = begin+1 to end-1
      for A,B,C in nonterms
        prob=score[begin][split][B]*score[split][end][C]*P(A->BC)
        if(prob > score[begin][end][A])
          score[begin]end][A] = prob
          back[begin][end][A] = new Triple(split,B,C)
      //handle unaries
      boolean added = true
      while added
        added = false
        for A, B in nonterms
          prob = P(A->B)*score[begin][end][B];
          if(prob > score[begin][end] [A])
            score[begin][end] [A] = prob
            back[begin][end] [A] = B
            added = true
  return buildTree(score, back)
```

|  | cats 1 | scratch 2 | walls 3 | with 4 | claws 5 |
|---|---|---|---|---|---|
| 0 | score[0][1] | score[0][2] | score[0][3] | score[0][4] | score[0][5] |
| 1 | | score[1][2] | score[1][3] | score[1][4] | score[1][5] |
| 2 | | | score[2][3] | score[2][4] | score[2][5] |
| 3 | | | | score[3][4] | score[3][5] |
| 4 | | | | | score[4][5] |
| 5 | | | | | |

|  | cats ₁ | scratch ₂ | walls ₃ | with ₄ | claws ₅ |

| | cats | 1 | scratch | 2 | walls | 3 | with | 4 | claws | 5 |
|---|---|---|---|---|---|---|---|---|---|---|

0

N→cats
P→cats
V→cats

1

N→scratch
P→scratch
V→scratch

2

N→walls
P→walls
V→walls

3

N→with
P→with
V→with

4

```
for i=0; i<#(words); i++
    for A in nonterms
    if A -> words[i] in grammar
        score[i][i+1][A] = P(A -> words[i]);
```

N→claws
P→claws
V→claws

5

|   | cats | 1 | scratch | 2 | walls | 3 | with | 4 | claws | 5 |
|---|------|---|---------|---|-------|---|------|---|-------|---|
| 0 | N→cats<br>P→cats<br>V→cats<br>NP→N<br>@VP->V→NP<br>@PP->P→NP | | | | | | | | | |
| 1 | | | N→scratch<br>P→scratch<br>V→scratch<br>NP→N<br>@VP->V→NP<br>@PP->P→NP | | | | | | | |
| 2 | | | | | N→walls<br>P→walls<br>V→walls<br>NP→N<br>@VP->V→NP<br>@PP->P→NP | | | | | |
| 3 | | | | | | | N→with<br>P→with<br>V→with<br>NP→N<br>@VP->V→NP<br>@PP->P→NP | | | |
| 4 | | | | | | | | | N→claws<br>P→claws<br>V→claws<br>NP→N<br>@VP->V→NP<br>@PP->P→NP | |
| 5 | | | | | | | | | | |

// handle unaries

|  | cats | 1 | scratch | 2 | walls | 3 | with | 4 | claws | 5 |

| | cats 1 | scratch 2 | walls 3 | with 4 | claws 5 |
|---|---|---|---|---|---|
| 0 | N→cats<br>P→cats<br>V→cats<br>NP→N<br>@VP->V→NP<br>@PP->P→NP | PP→P @PP->_P<br>VP→V @VP->_V | | | |
| 1 | | N→scratch<br>P→scratch<br>V→scratch<br>NP→N<br>@VP->V→NP<br>@PP->P→NP | PP→P @PP->_P<br>VP→V @VP->_V | | |
| 2 | | | N→walls<br>P→walls<br>V→walls<br>NP→N<br>@VP->V→NP<br>@PP->P→NP | PP→P @PP->_P<br>VP→V @VP->_V | |
| 3 | | | | N→with<br>P→with<br>V→with<br>NP→N<br>@VP->V→NP<br>@PP->P→NP | PP→P @PP->_P<br>VP→V @VP->_V |
| 4 | | | | | N→claws<br>P→claws<br>V→claws<br>NP→N<br>@VP->V→NP<br>@PP->P→NP |
| 5 | | | | | |

prob=score[begin][split][B]*score[split][end][C]*P(A->BC)
prob=score[0][1][P]*score[1][2][@PP->_P]*P(PP→P @PP->_P)

For each A, only keep the "A->BC" with highest prob.

|  | cats | 1 | scratch | 2 | walls | 3 | with | 4 | claws | 5 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | N→cats<br>P→cats<br>V→cats<br>NP→N<br>@VP->V→NP<br>@PP->P→NP | | PP→P @PP->_P<br>VP→V @VP->_V<br>@S->_NP→VP<br>@NP->_NP→PP<br>@VP->_V_NP→PP | | | | | | | |
| 1 | | | N→scratch<br>P→scratch<br>V→scratch<br>NP→N<br>@VP->V→NP<br>@PP->P→NP | | PP→P @PP->_P<br>VP→V @VP->_V<br>@S->_NP→VP<br>@NP->_NP→PP<br>@VP->_V_NP→PP | | | | | |
| 2 | | | 0.0859 | | N→walls<br>P→walls<br>V→walls<br>NP→N<br>@VP->V→NP<br>@PP->P→NP | | PP→P @PP->_P<br>VP→V @VP->_V<br>@S->_NP→VP<br>@NP->_NP→PP<br>@VP->_V_NP→PP | | | |
| 3 | | | | | 0.2514 | | N→with<br>P→with<br>V→with<br>NP→N<br>@VP->V→NP<br>@PP->P→NP | | PP→P @PP->_P<br>VP→V @VP->_V<br>@S->_NP→VP<br>@NP->_NP→PP<br>@VP->_V_NP→PP | |
| 4 | // handle unaries | | | | | | 0.0859 | | N→claws<br>P→claws<br>V→claws<br>NP→N<br>@VP->V→NP<br>@PP->P→NP | |
| 5 | | | | | | | | | 0.3611 | |

. . . . . . . . .

cats 1    scratch 2    walls 3    with 4    claws 5

| | cats (1) | scratch (2) | walls (3) | with (4) | claws (5) |
|---|---|---|---|---|---|
| **0** | N→cats 0.5259<br>P→cats 0.0725<br>V→cats 0.0967<br>NP→N 0.4675<br>@VP->V→NP 0.3116<br>@PP->P→NP 0.4675 | PP→P @PP->_P 0.0062<br>VP→V @VP->_V 0.0055<br>@S->_NP→VP 0.0055<br>@NP->_NP→PP 0.0062<br>@VP->_V_NP→PP 0.0062 | @VP->_V→NP @VP->_V_NP 0.0030<br>NP→NP @NP->_NP 0.0010<br>S→NP @S->_NP 0.0727<br>ROOT→S 0.0727<br>@PP->_P→NP 0.0010 | PP→P @PP->_P 5.187E-6<br>VP→V @VP->_V 2.074E-5<br>@S->_NP→VP 2.074E-5<br>@NP->_NP→PP 5.187E-6<br>@VP->_V_NP→PP 5.187E-6 | @VP->_V→NP @VP->_V_NP 1.600E-4<br>NP→NP @NP->_NP 5.335E-5<br>S→NP @S->_NP 0.0172<br>ROOT→S 0.0172<br>@PP->_P→NP 5.335E-5 |
| **1** | | N→scratch 0.0967<br>P→scratch 0.0773<br>V→scratch 0.9285<br>NP→N 0.0859<br>@VP->V→NP 0.0573<br>@PP->P→NP 0.0859 | PP→P @PP->_P 0.0194<br>VP→V @VP->_V 0.1556<br>@S->_NP→VP 0.1556<br>@NP->_NP→PP 0.0194<br>@VP->_V_NP→PP 0.0194 | @VP->_V→NP @VP->_V_NP 2.145E-4<br>NP→NP @NP->_NP 7.150E-5<br>S→NP @S->_NP 5.720E-4<br>ROOT→S 5.720E-4<br>@PP->_P→NP 7.150E-5 | PP→P @PP->_P 0.0010<br>VP→V @VP->_V 0.0369<br>@S->_NP→VP 0.0369<br>@NP->_NP→PP 0.0010<br>@VP->_V_NP→PP 0.0010 |
| **2** | | | N→walls 0.2829<br>P→walls 0.0870<br>V→walls 0.1160<br>NP→N 0.2514<br>@VP->V→NP 0.1676<br>@PP->P→NP 0.2514 | PP→P @PP->_P 0.0074<br>VP→V @VP->_V 0.0066<br>@S->_NP→VP 0.0066<br>@NP->_NP→PP 0.0074<br>@VP->_V_NP→PP 0.0074 | @VP->_V→NP @VP->_V_NP 0.0398<br>NP→NP @NP->_NP 0.0132<br>S→NP @S->_NP 0.0062<br>ROOT→S 0.0062<br>@PP->_P→NP 0.0132 |
| **3** | | | | N→with 0.0967<br>P→with 1.3154<br>V→with 0.1031<br>NP→N 0.0859<br>@VP->V→NP 0.0573<br>@PP->P→NP 0.0859 | PP→P @PP->_P 0.4750<br>VP→V @VP->_V 0.0248<br>@S->_NP→VP 0.0248<br>@NP->_NP→PP 0.4750<br>@VP->_V_NP→PP 0.4750 |
| **4** | | | | | N→claws 0.4062<br>P→claws 0.0773<br>V→claws 0.1031<br>NP→N 0.3611<br>@VP->V→NP 0.2407<br>@PP->P→NP 0.3611 |
| **5** | | | | | |

Call buildTree(score, back) to get the best parse

# Unary rules:
# alchemy in the land of treebanks

# Same-Span Reachability



NoEmpties

TOP

SQ    X    RRC

NX

ADJP ADVP
FRAG INTJ NP
PP PRN QP S
SBAR UCP VP
WHNP

LST

CONJP

NAC

SINV    PRT

WHADJP    SBARQ    WHPP

WHADVP

# Efficient CKY parsing

- CKY parsing can be made very fast (!), partly due to the simplicity of the structures used.
  - But that means a lot of the speed comes from engineering details
  - And a little from cleverer filtering

  - Store chart as (ragged) 3 dimensional array of float (log probabilities)
    - score[start][end][category]
      - For treebank grammars the load is high enough that you don't really gain from lists of things that were possible
      - 50 wds: (50x50)/2x(1000 to 20000)x4 bytes = 5–100MB for parse triangle. Large. (Can move to beam for span[i][j].)
  - Use int to represent categories/words (Index)

# Efficient CKY parsing

- Provide efficient grammar/lexicon accessors:
  - E.g., return list of rules with this left child category
  - Iterate over left child, check for zero (Neg. inf.) prob of X:[i,j] (abort loop), otherwise get rules with X on left

  - Some X:[i,j] can be filtered based on the input string
    - Not enough space to complete a long flat rule?
    - No word in the string can be a CC?
      - Using a lexicon of possible POS for words gives a lot of constraint rather than allowing all POS for words
    - Cf. later discussion of figures-of-merit/A* heuristics

# Quiz Question!

PP → IN              0.002
NP → NNS NNS         0.01
NP → NNS NP 0.005
NP → NNS PP 0.01
VP → VB PP              0.045
VP → VB NP              0.015

a) NP (2.3e-9)
b) NP (4.6e-6)
c) PP (.0004)
d) VP (.0002)

PP        0.2
IN        0.0014
NNS        0.0001

NNS      0.0023
VB          0.001

Which constituent (with probability) can you make?

runs                    down

# 3. Evaluating Parsing Accuracy

- Most sentences are not given a completely correct parse by any currently existing parser.
- For Penn Treebank parsing, the standard evaluation is over the number of correct constituents (labeled spans).

- [ label, start, finish ]

- A constituent is a triple, which must be exact in the true parse for the constituent to be marked correct.
- The LP/LR $F_1$ is the micro-averaged harmonic mean of labeled constituent precision and recall
- This isn't necessarily a great measure … many people think dependency accuracy or raw data likelihood would be better.

(a)

```
                                          ROOT
                                           |
                                           S
          ┌────────────────────────────────┼──────────────────────────────┐
          NP                               VP                              NP      .
      ┌────┴────┐              ┌────────────┼──────────────┐                |      |
     NNS        NNS           VBD                  VP                       NN   . 11
      |          |             |         ┌──────────┼──────────┐            |
   0 Sales 1  executives 2 were        VBG         NP          PP       yesterday 10
                                        |      ┌────┴────┐  ┌────┴────┐
                                   3 examining  DT       NNS  IN        NP
                                                |         |    |    ┌────┴────┐
                                             4 the  5 figures 6 with JJ        NN
                                                                     |         |
                                                                 7 great   8 care 9
```

(b)   Brackets in gold standard tree (a.):
      **S-(0:11)**, **NP-(0:2)**, VP-(2:9), VP-(3:9), **NP-(4:6)**, PP-(6-9), NP-(7,9), *NP-(9:10)

(c)   Brackets in candidate parse:
      **S-(0:11)**, **NP-(0:2)**, VP-(2:10), VP-(3:10), NP-(4:10), **NP-(4:6)**, PP-(6-10), NP-(7,10)

(d)

| | | | |
|---|---|---|---|
| Precision: | 3/8 = 37.5% | Crossing Brackets: | 0 |
| Recall: | 3/8 = 37.5% | Crossing Accuracy: | 100% |
| Labeled Precision: | 3/8 = 37.5% | Tagging Accuracy: | 10/11 = 90.9% |
| Labeled Recall: | 3/8 = 37.5% | | |

# How good are PCFGs?

- Robust (usually admit everything, but with low probability)
- Partial solution for grammar ambiguity: a PCFG gives some idea of the plausibility of a sentence
- But not so good because the independence assumptions are too strong
- Give a probabilistic language model
  - But in a simple case it performs worse than a trigram model
- WSJ parsing accuracy: about 73% LP/LR F1
- The problem seems to be that PCFGs lack the lexicalization of a trigram model

# Putting words into PCFGs

- A PCFG uses the actual words only to determine the probability of parts-of-speech (the preterminals)
- In many cases we need to know about words to choose a parse
- The head word of a phrase gives a good representation of the phrase's structure and meaning
  - Attachment ambiguities

    *The astronomer saw the moon with the telescope*
  - Coordination

    *the dogs in the house and the cats*
  - Subcategorization frames

    *put* versus *like*

# (Head) Lexicalization

- *put* takes both an NP and a VP
  - *Sue put [ the book ]$_{NP}$ [ on the table ]$_{PP}$*
  - *\* Sue put [ the book ]$_{NP}$*
  - *\* Sue put [ on the table ]$_{PP}$*

- *like* usually takes an NP and not a PP
  - *Sue likes [ the book ]$_{NP}$*
  - *\* Sue likes [ on the table ]$_{PP}$*

- We can't tell this if we just have a VP with a verb, but we can if we know *what* verb it is

# 4. Accurate Unlexicalized Parsing: PCFGs and Independence

- The symbols in a PCFG define independence assumptions:

  S → NP VP

  NP → DT NN

  

  - At any node, the material inside that node is independent of the material outside that node, given the label of that node.
  - Any information that statistically connects behavior inside and outside a node must flow through that node.

# Non-Independence I

- Independence assumptions are often too strong.



| All NPs | NPs under S | NPs under VP |
|---------|-------------|--------------|
| NP PP: 11%, DT NN: 9%, PRP: 6% | NP PP: 9%, DT NN: 9%, PRP: 21% | NP PP: 23%, DT NN: 7%, PRP: 4% |

- Example: the expansion of an NP is highly dependent on the parent of the NP (i.e., subjects vs. objects).

# Michael Collins (2003, COLT)

## Independence Assumptions

- PCFGs



73% accuracy

- Lexicalized PCFGs



88% accuracy

# Non-Independence II

- Who cares?
  - NB, HMMs, all make false assumptions!
  - For generation/LMs, consequences would be obvious.
  - For parsing, does it impact accuracy?
- Symptoms of overly strong assumptions:
  - Rewrites get used where they don't belong.
  - Rewrites get used too often or too rarely.



*In the PTB, this construction is for possessives*

# Breaking Up the Symbols

- We can relax independence assumptions by encoding dependencies into the PCFG symbols:

Parent annotation
[Johnson 98]

Marking
possesive NPs



- What are the most useful features to encode?

# Annotations

- Annotations split the grammar categories into sub-categories.

- Conditioning on history vs. annotating
  - P(NP^S → PRP) is a lot like P(NP → PRP | S)
  - P(NP-POS → NNP POS) isn't history conditioning.

- Feature grammars vs. annotation
  - Can think of a symbol like NP^NP-POS as
    NP [parent:NP, +POS]

- After parsing with an annotated grammar, the annotations are then stripped for evaluation.

# Experimental Setup

- Corpus: Penn Treebank, WSJ

Training: sections 02-21
Development: section 22 (first 20 files)
Test: section 23

- Accuracy – F1: harmonic mean of per-node labeled precision and recall.

- Size – number of symbols in grammar.
  - Passive / complete symbols: NP, NP^S
  - Active / incomplete symbols: NP $\rightarrow$ NP CC $\bullet$

# Experimental Process

- We'll take a highly conservative approach:
  - Annotate as sparingly as possible
  - Highest accuracy with fewest symbols
  - Error-driven, manual hill-climb, adding one annotation type at a time

# Lexicalization

- Lexical heads are important for certain classes of ambiguities (e.g., PP attachment):



- Lexicalizing grammar creates a much larger grammar.
    - Sophisticated smoothing needed
    - Smarter parsing algorithms needed
    - More data needed

- How necessary is lexicalization?
    - Bilexical vs. monolexical selection
    - Closed vs. open class lexicalization

# Unlexicalized PCFGs

- 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^S-CC is fine
  - Closed vs. open class words (NP^S-the)
    - Long tradition in linguistics of using function words as features or markers for selection
    - Contrary to the bilexical idea of semantic heads
    - Open-class selection really a proxy for semantics

- Honesty checks:
  - Number of symbols: keep the grammar very small
  - No smoothing: over-annotating is a real danger
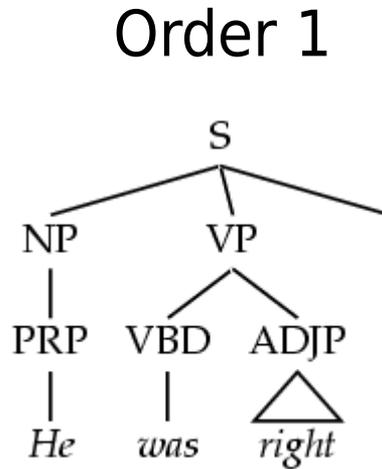
# Horizontal Markovization

- Horizontal Markovization: Merges States

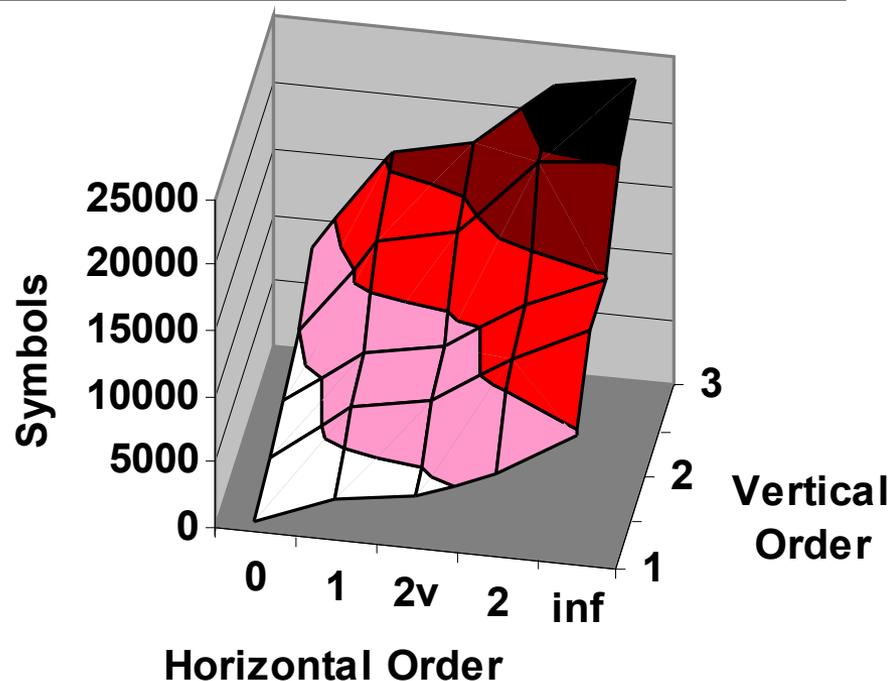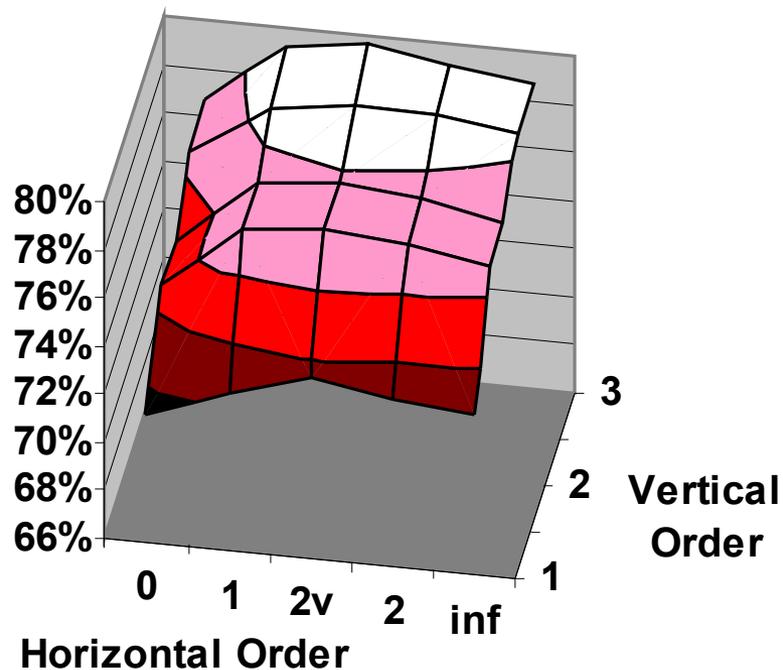# Vertical Markovization

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

Order 1

Order 2
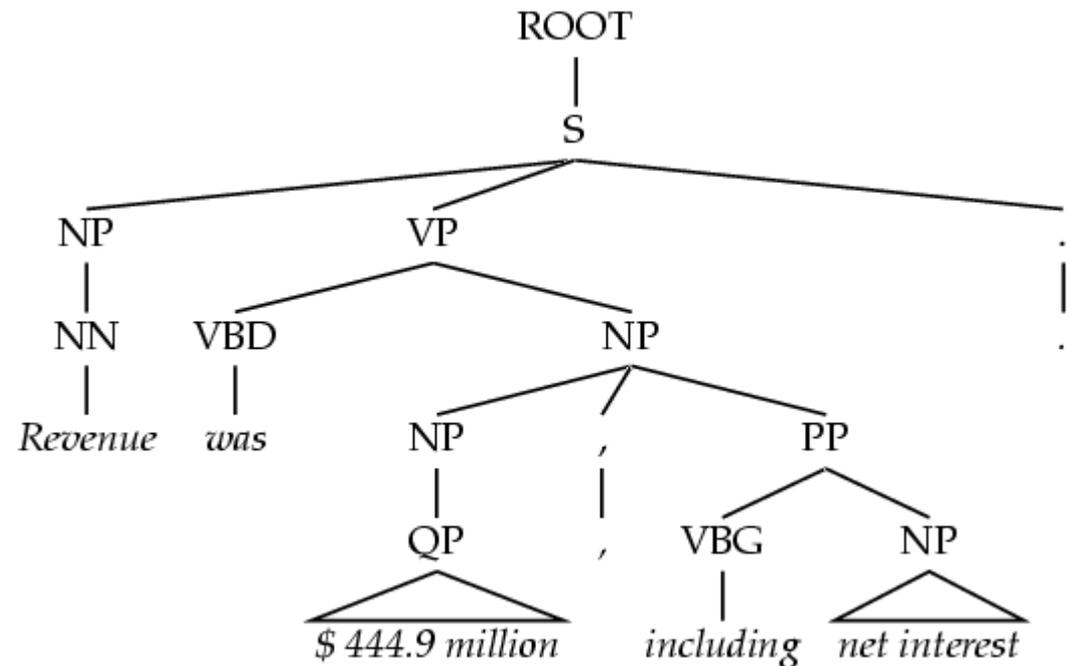




Vertical Markov Order

Vertical Markov Order

# Vertical and Horizontal





- Examples:
  - Raw treebank:        v=1, h=∞
  - Johnson 98: v=2, h=∞
  - Collins 99:   v=2, h=2
  - Best F1:                v=3, h=2v

| Model | F1 | Size |
|---|---|---|
| Base: v=h=2v | 77.8 | 7.5K |

# Unary Splits

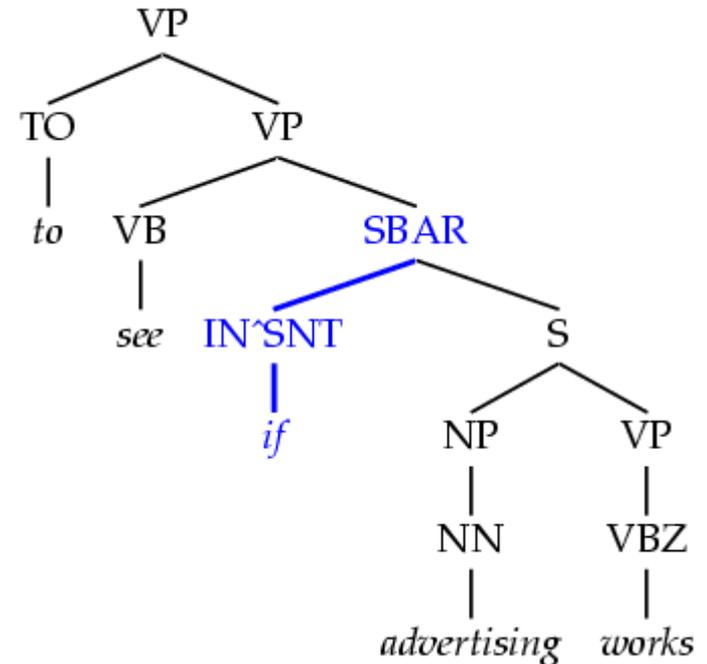- Problem: unary rewrites used to transmute categories so a high-probability rule can be used.



- 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: Sentential, PP, and other prepositions are all marked IN.

- Partial Solution:
  - Subdivide the IN tag.

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

# Other Tag Splits

| F1 | Size |
|---|---|
| 80.4 | 8.1K |
| 80.5 | 8.1K |
| 81.2 | 8.5K |
| 81.6 | 9.0K |
| 81.7 | 9.1K |
| 81.8 | 9.3K |

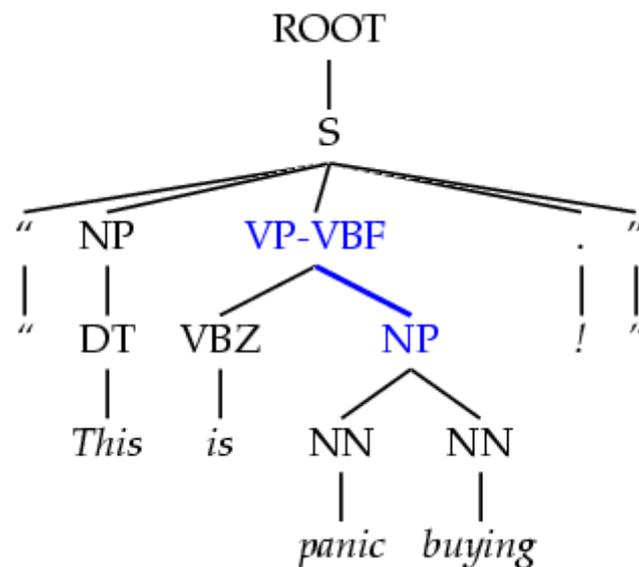- UNARY-DT: mark demonstratives as DT^U ("the X" vs. "those")
- UNARY-RB: mark phrasal adverbs as RB^U ("quickly" vs. "very")
- TAG-PA: mark tags with non-canonical parents ("not" is an RB^VP)
- SPLIT-AUX: mark auxiliary verbs with –AUX [cf. Charniak 97]
- SPLIT-CC: separate "but" and "&" from other conjunctions
- SPLIT-%: "%" gets its own tag.

# 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!

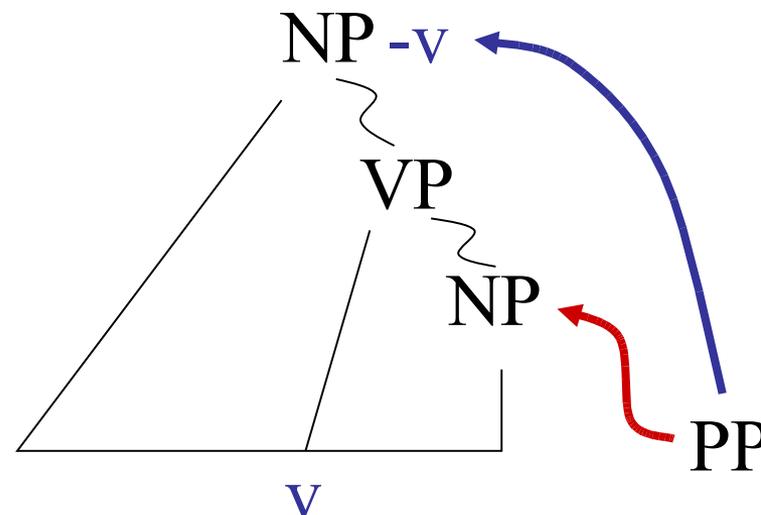- Solution: annotate future elements into nodes.



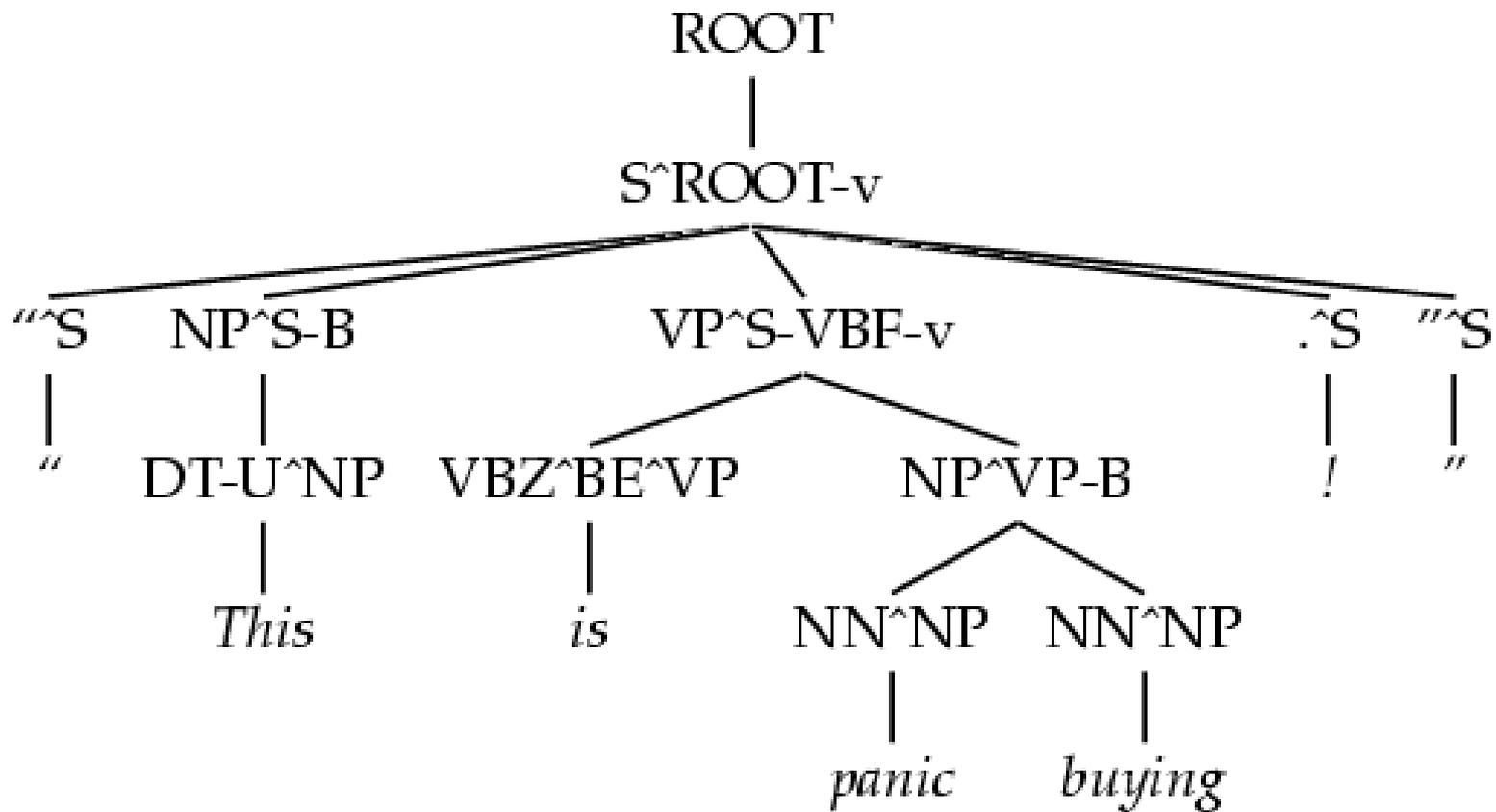| Annotation | F1 | Size |
|---|---|---|
| Previous | 82.3 | 9.7K |
| POSS-NP | 83.1 | 9.8K |
| SPLIT-VP | 85.7 | 10.5K |

# 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

NP -v

VP

NP

PP

v

| 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** | CB | 0 CB |
|--------|------|------|----------|------|------|
| Magerman 95 | 84.9 | 84.6 | **84.7** | 1.26 | 56.6 |
| Collins 96 | 86.3 | 85.8 | **86.0** | 1.14 | 59.9 |
| Klein & M 03 | 86.9 | 85.7 | **86.3** | 1.10 | 60.3 |
| Charniak 97 | 87.4 | 87.5 | **87.4** | 1.00 | 62.1 |
| Collins 99 | 88.7 | 88.6 | **88.6** | 0.90 | 67.1 |

- Beats "first generation" lexicalized parsers.