# Statistical NLP
## Spring 2010

University of
California

C A L
N L P

Berkeley

## Lecture 14: PCFGs

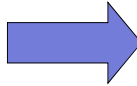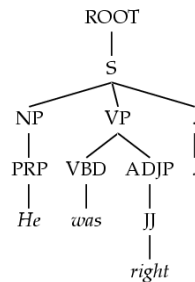Dan Klein – UC Berkeley

---

University of
California
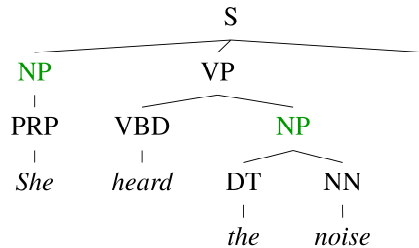C A L
N L P
Berkeley

# Treebank PCFGs          [Charniak 96]

- Use PCFGs for broad coverage parsing
- Can take a grammar right off the trees (doesn't work well):

ROOT
|
S

NP    VP    .
|
PRP   VBD   ADJP   .
|     |     |
He    was   JJ
|
right

$ROOT \rightarrow S$         1

$S \rightarrow NP\ VP\ .$       1

$NP \rightarrow PRP$           1

$VP \rightarrow VBD\ ADJP$      1

…..

| Model | F1 |
|---|---|
| Baseline | 72.0 |

# Conditional Independence?

```
                      S
         ┌────────────┼──────────────┐
        NP           VP               .
         │      ┌─────┼─────┐         │
        PRP    VBD         NP         .
         │      │      ┌────┴────┐
        She   heard   DT        NN
                       │         │
                      the      noise
```
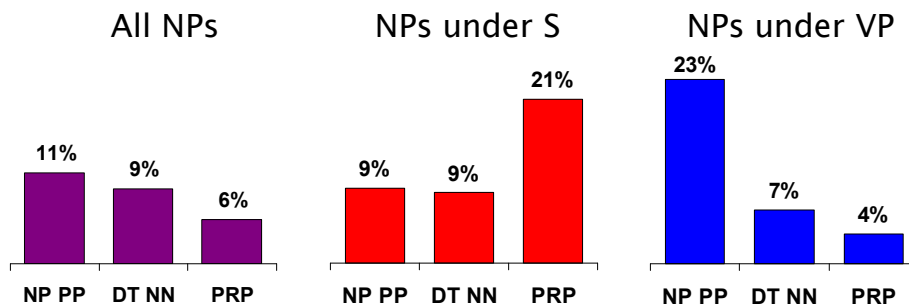
- Not every NP expansion can fill every NP slot
  - A grammar with symbols like "NP" won't be context-free
  - Statistically, conditional independence too strong
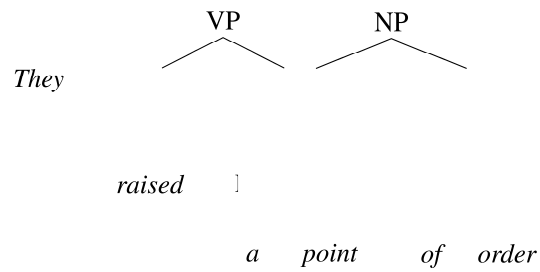
---

# Non-Independence

- Independence assumptions are often too strong.
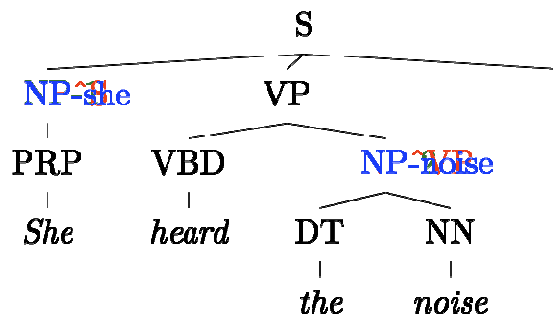


| All NPs | NPs under S | NPs under VP |

- Example: the expansion of an NP is highly dependent on the parent of the NP (i.e., subjects vs. objects).
- Also: the subject and object expansions are correlated!
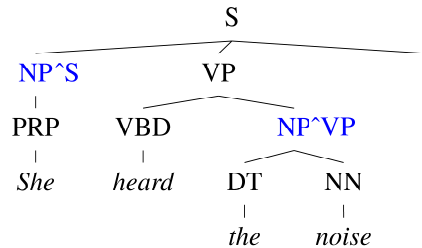
# Grammar Refinement

- Example: PP attachment

VP      NP

*They*

*raised*

*a point of order*

# Grammar Refinement

S

NP^~~the~~she    VP    .

PRP    VBD    NP-~~VP~~noise    .

*She*    *heard*    DT    NN

*the*    *noise*

- Structure Annotation [Johnson '98, Klein&Manning '03]
- Lexicalization [Collins '99, Charniak '00]
- Latent Variables [Matsuzaki et al. 05, Petrov et al. '06]

# The Game of Designing a Grammar

```
                    S
        ┌───────────┼───────────┐
      NP^S          VP           .
       │        ┌────┴────┐      │
      PRP      VBD      NP^VP     .
       │        │      ┌──┴──┐
      She     heard   DT    NN
                       │     │
                      the   noise
```

- Annotation refines base treebank symbols to improve statistical fit of the grammar
  - Structural annotation

# Typical Experimental Setup

- Corpus: Penn Treebank, WSJ
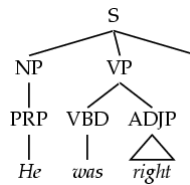
| Training: | sections | 02-21 |
| Development: | section | 22 (here, first 20 files) |
| Test: | section | 23 |

- Accuracy – F1: harmonic mean of per-node labeled precision and recall.
- Here: also size – number of symbols in grammar.
  - Passive / complete symbols: NP, NP^S
  - Active / incomplete symbols: NP $\rightarrow$ NP CC •
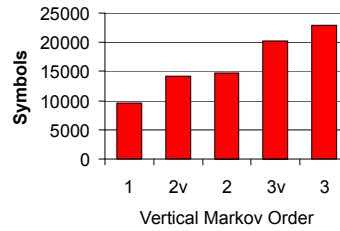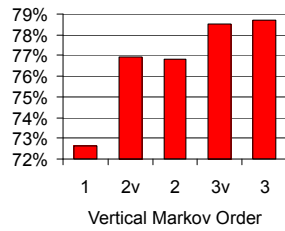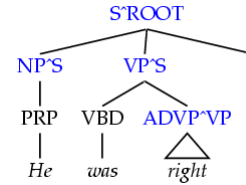
# Vertical Markovization

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

### Order 1

```
        S
      / | \
   NP   VP   .
   |   / \   |
  PRP VBD ADJP .
   |   |   /\
   He  was right
```
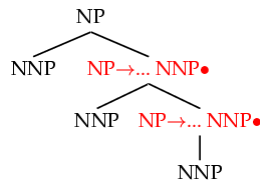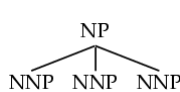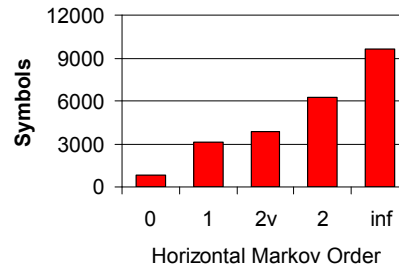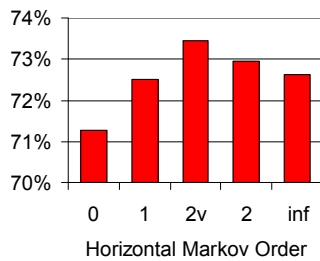
### Order 2

```
         S^ROOT
        /   |   \
    NP^S   VP^S   .
     |    /  \    |
    PRP  VBD  ADVP^VP .
     |    |    /\
     He  was  right
```

79% 78% 77% 76% 75% 74% 73% 72%

1  2v  2  3v  3
Vertical Markov Order

25000 20000 15000 10000 5000 0

Symbols

1  2v  2  3v  3
Vertical Markov Order

# Horizontal Markovization

```
       NP
     / | \
  NNP NNP NNP
```

### Order 1

```
        NP
       /    \
    NNP    NP→... NNP•
           /    \
         NNP   NP→... NNP•
                 |
                NNP
```

### Order ∞

```
        NP
       /    \
    NNP    NP→NNP•
           /    \
         NNP   NP→NNP NNP•
                 |
                NNP
```

74% 73% 72% 71% 70%

0  1  2v  2  inf
Horizontal Markov Order

12000 9000 6000 3000 0

Symbols

0  1  2v  2  inf
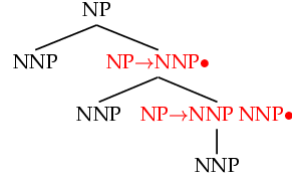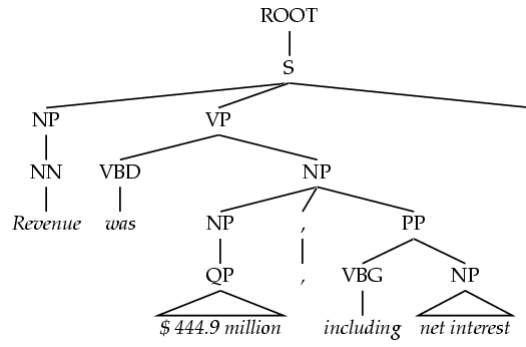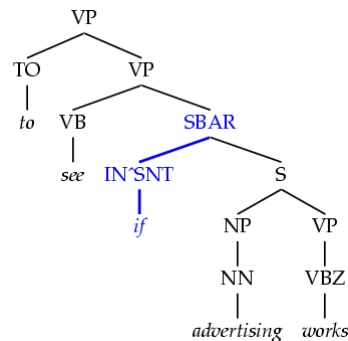Horizontal Markov Order

# 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.
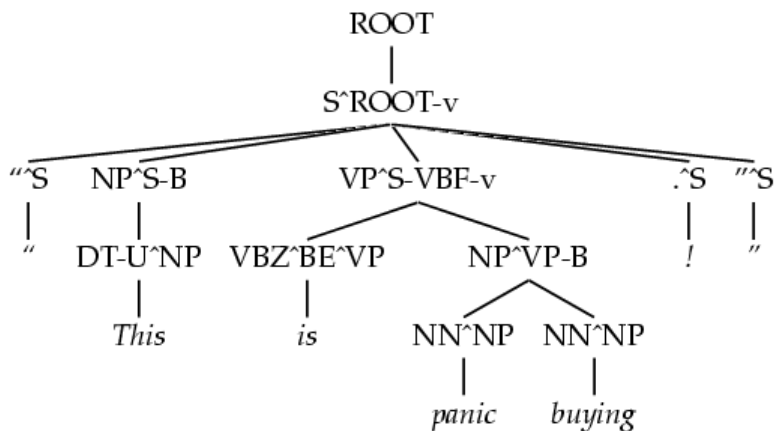
- Partial Solution:
  - Subdivide the IN tag.



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

# Other Tag Splits

| | F1 | Size |
|---|---|---|
| UNARY-DT: mark demonstratives as DT^U ("the X" vs. "those") | 80.4 | 8.1K |
| UNARY-RB: mark phrasal adverbs as RB^U ("quickly" vs. "very") | 80.5 | 8.1K |
| TAG-PA: mark tags with non-canonical parents ("not" is an RB^VP) | 81.2 | 8.5K |
| SPLIT-AUX: mark auxiliary verbs with –AUX [cf. Charniak 97] | 81.6 | 9.0K |
| SPLIT-CC: separate "but" and "&" from other conjunctions | 81.7 | 9.1K |
| SPLIT-%: "%" gets its own tag. | 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.

# A Fully Annotated (Unlex) Tree



7

# Some 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 |
| Unlexicalized | 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.
- Lots of room to improve – more complex models next.
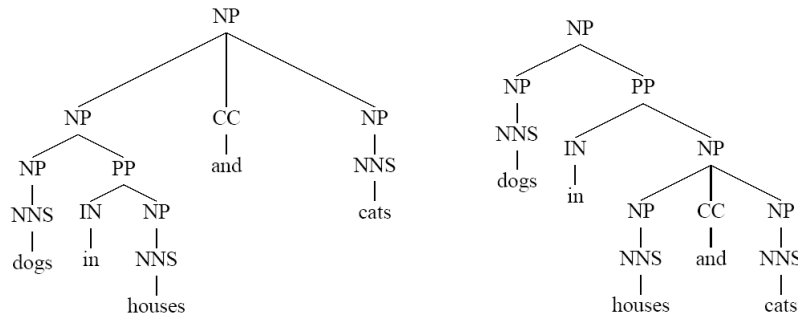
# The Game of Designing a Grammar

- Annotation refines base treebank symbols to improve statistical fit of the grammar
  - Structural annotation [Johnson '98, Klein and Manning 03]
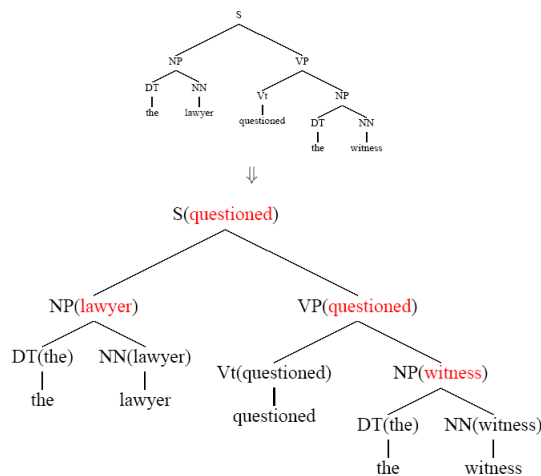  - Head lexicalization [Collins '99, Charniak '00]

# Problems with PCFGs



- What's different between basic PCFG scores here?
- What (lexical) correlations need to be scored?

# Lexicalized Trees

- Add "headwords" to each phrasal node
  - Syntactic vs. semantic heads
  - Headship not in (most) treebanks
  - Usually *use head rules*, e.g.:
    - NP:
      - Take leftmost NP
      - Take rightmost N*
      - Take rightmost JJ
      - Take right child
    - VP:
      - Take leftmost VB*
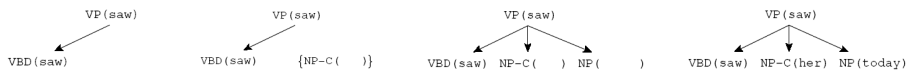      - Take leftmost VP
      - Take left child

# Lexicalized PCFGs?
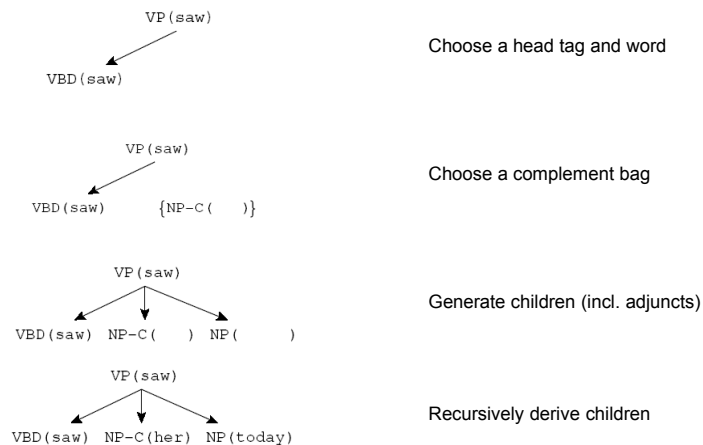
- Problem: we now have to estimate probabilities like

```
VP(saw) -> VBD(saw) NP-C(her) NP(today)
```

- Never going to get these atomically off of a treebank

- Solution: break up derivation into smaller steps



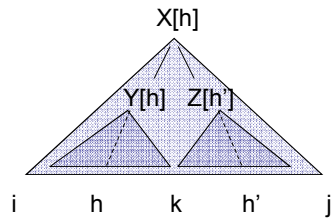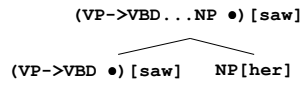# Lexical Derivation Steps
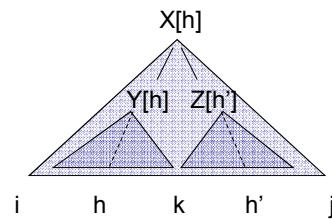
- A derivation of a local tree [Collins 99]



Choose a head tag and word

Choose a complement bag

Generate children (incl. adjuncts)

Recursively derive children

# Lexicalized CKY

```
                    (VP->VBD...NP •)[saw]                          X[h]
                   ____/\____
        (VP->VBD •)[saw]   NP[her]                              Y[h]  Z[h']

bestScore(X,i,j,h)

  if (j = i+1)
    return tagScore(X,s[i])                      i     h     k    h'    j
  else
    return
      max max score(X[h]->Y[h] Z[h']) *
      k,h',X->YZ
              bestScore(Y,i,k,h) *
              bestScore(Z,k,j,h')
        max score(X[h]->Y[h'] Z[h]) *
      k,h',X->YZ
              bestScore(Y,i,k,h') *
              bestScore(Z,k,j,h)
```

---

# Pruning with Beams

- The Collins parser prunes with per-cell beams [Collins 99]
  - Essentially, run the O($n^5$) CKY
  - Remember only a few hypotheses for each span <i,j>.
  - If we keep K hypotheses at each span, then we do at most O($nK^2$) work per span (why?)
  - Keeps things more or less cubic

- Also: certain spans are forbidden entirely on the basis of punctuation (crucial for speed)

```
                X[h]

           Y[h]  Z[h']

    i     h     k     h'    j
```

# Pruning with a PCFG
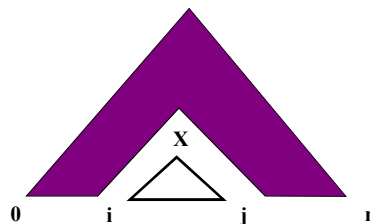
- The Charniak parser prunes using a two-pass approach [Charniak 97+]
    - First, parse with the base grammar
    - For each X:[i,j] calculate $P(X|i,j,s)$
        - This isn't trivial, and there are clever speed ups
    - Second, do the full $O(n^5)$ CKY
        - Skip any X :[i,j] which had low (say, < 0.0001) posterior
    - Avoids almost all work in the second phase!

- Charniak et al 06: can use more passes
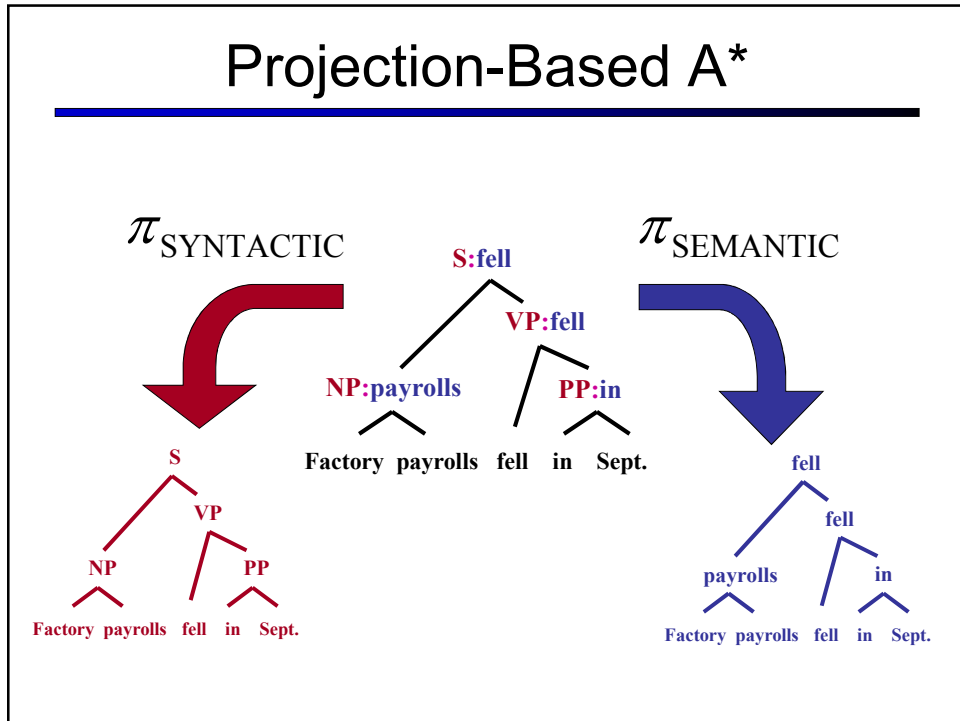- Petrov et al 07: can use many more passes
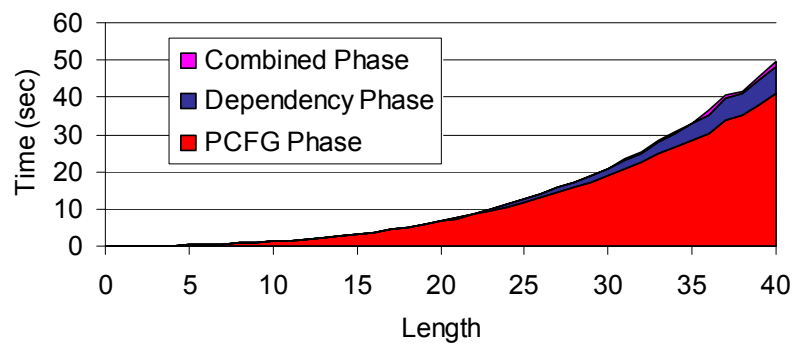

# Pruning with A*

- You can also speed up the search without sacrificing optimality
- For agenda-based parsers:
    - Can select which items to process first
    - Can do with any "figure of merit" [Charniak 98]
    - If your figure-of-merit is a valid A* heuristic, no loss of optimiality [Klein and Manning 03]

# Projection-Based A*

$\pi_{\text{SYNTACTIC}}$     **S:fell**     $\pi_{\text{SEMANTIC}}$

**VP:fell**

**NP:payrolls**     **PP:in**

Factory payrolls   fell   in   Sept.

**S**

**VP**

**NP**     **PP**

Factory payrolls   fell   in   Sept.

fell

fell

payrolls     in

Factory payrolls   fell   in   Sept.

---

# A* Speedup



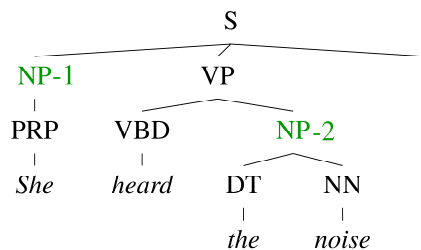- Combined Phase
- Dependency Phase
- PCFG Phase

Time (sec) — Length

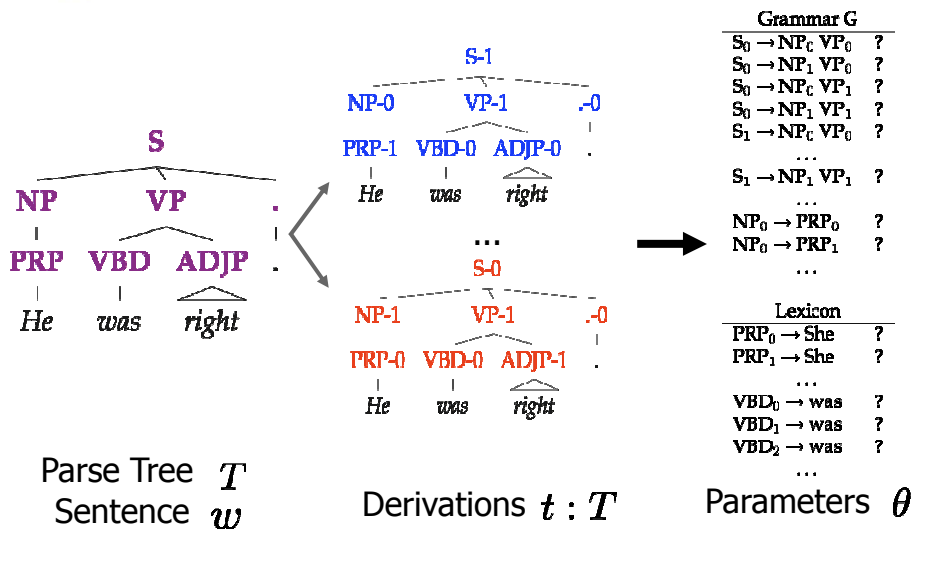- Total time dominated by calculation of A* tables in each projection… $O(n^3)$

# Results

- Some results
  - Collins 99 – 88.6 F1 (generative lexical)
  - Charniak and Johnson 05 – 89.7 / 91.3 F1 (generative lexical / reranked)
  - Petrov et al 06 – 90.7 F1 (generative unlexical)
  - McClosky et al 06 – 92.1 F1 (gen + rerank + self-train)

- However
  - Bilexical counts rarely make a difference (why?)
  - Gildea 01 – Removing bilexical counts costs < 0.5 F1

# The Game of Designing a Grammar

```
                    S
         ┌──────────┼──────────┐
      NP-1          VP          .
        │      ┌────┼─────┐     │
      PRP     VBD        NP-2   .
        │      │      ┌───┴───┐
      She    heard    DT      NN
                       │       │
                      the    noise
```

- Annotation refines base treebank symbols to improve statistical fit of the grammar
  - Structural annotation
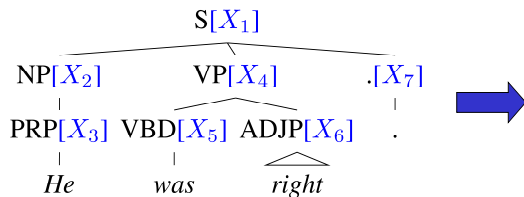  - Head lexicalization
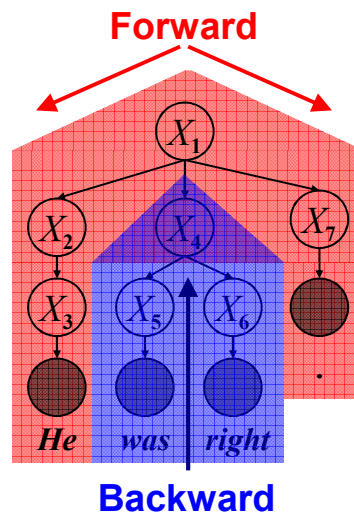  - Automatic clustering?

# Latent Variable Grammars

**Parse Tree** $T$
**Sentence** $w$

**Derivations** $t : T$

**Parameters** $\theta$

Grammar G

| | |
|---|---|
| $S_0 \to NP_0\ VP_0$ | ? |
| $S_0 \to NP_1\ VP_0$ | ? |
| $S_0 \to NP_0\ VP_1$ | ? |
| $S_0 \to NP_1\ VP_1$ | ? |
| $S_1 \to NP_0\ VP_0$ | ? |
| ... | |
| $S_1 \to NP_1\ VP_1$ | ? |
| ... | |
| $NP_0 \to PRP_0$ | ? |
| $NP_0 \to PRP_1$ | ? |
| ... | |

Lexicon

| | |
|---|---|
| $PRP_0 \to$ She | ? |
| $PRP_1 \to$ She | ? |
| ... | |
| $VBD_0 \to$ was | ? |
| $VBD_1 \to$ was | ? |
| $VBD_2 \to$ was | ? |
| ... | |

---

# Learning Latent Annotations

EM algorithm:
- Brackets are known
- Base categories are known
- Only induce subcategories

$S[X_1]$
$NP[X_2]$ $VP[X_4]$ $.[X_7]$
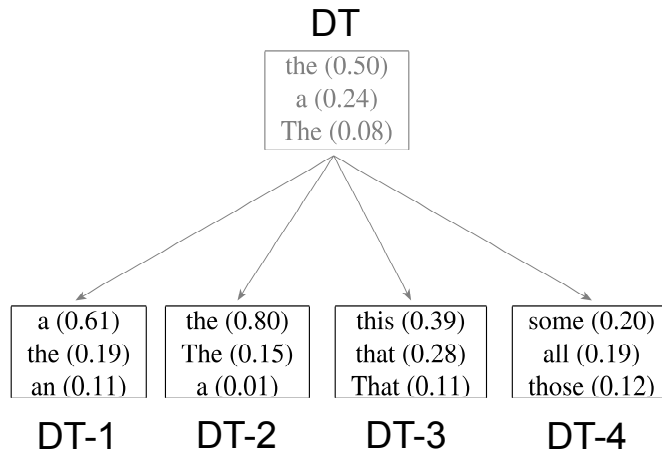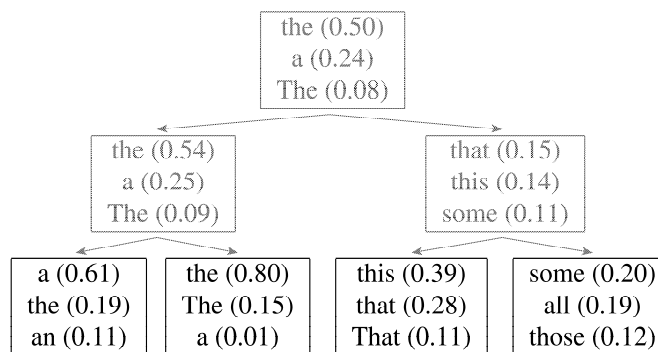$PRP[X_3]$ $VBD[X_5]$ $ADJP[X_6]$ .
*He* *was* *right*

**Just like Forward-Backward for HMMs.**

**Forward**

**Backward**

15 footer

# Refinement of the DT tag

DT

the (0.50)
a (0.24)
The (0.08)

| DT-1 | DT-2 | DT-3 | DT-4 |
|------|------|------|------|
| a (0.61) | the (0.80) | this (0.39) | some (0.20) |
| the (0.19) | The (0.15) | that (0.28) | all (0.19) |
| an (0.11) | a (0.01) | That (0.11) | those (0.12) |


# Hierarchical refinement

the (0.50)
a (0.24)
The (0.08)

the (0.54)
a (0.25)
The (0.09)

that (0.15)
this (0.14)
some (0.11)

| a (0.61) | the (0.80) | this (0.39) | some (0.20) |
|------|------|------|------|
| the (0.19) | The (0.15) | that (0.28) | all (0.19) |
| an (0.11) | a (0.01) | That (0.11) | those (0.12) |

# Hierarchical Estimation Results
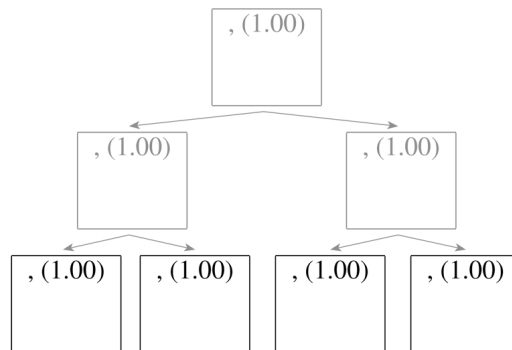


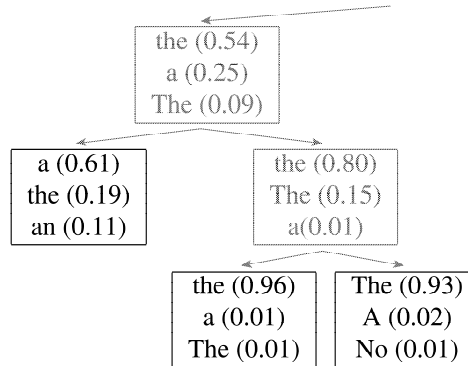| Model | F1 |
|---|---|
| Flat Training | 87.3 |
| Hierarchical Training | 88.4 |

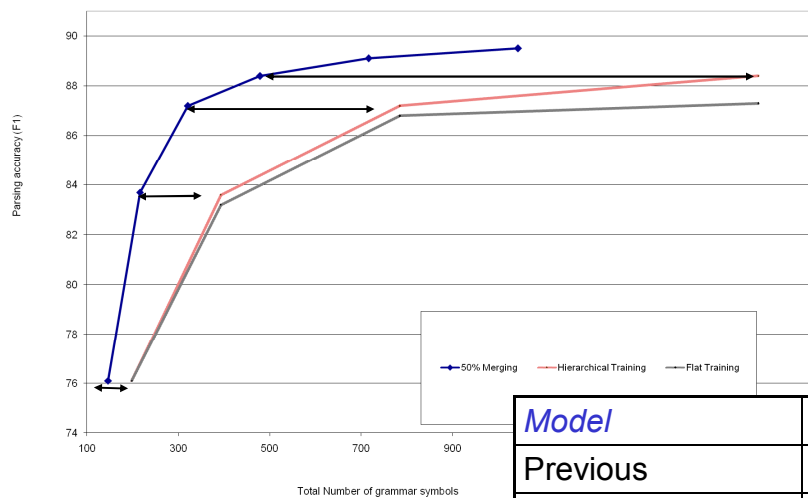# Refinement of the , tag

- Splitting all categories equally is wasteful:

# Adaptive Splitting

- Want to split complex categories more
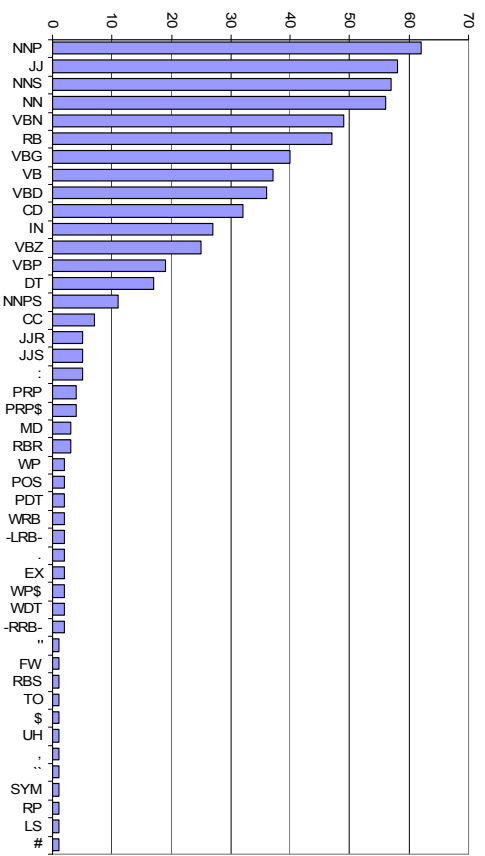- Idea: split everything, roll back splits which were least useful

the (0.54)
a (0.25)
The (0.09)

a (0.61)
the (0.19)
an (0.11)

the (0.80)
The (0.15)
a(0.01)

the (0.96)
a (0.01)
The (0.01)

The (0.93)
A (0.02)
No (0.01)

# Adaptive Splitting Results



| Model | F1 |
|---|---|
| Previous | 88.4 |
| With 50% Merging | 89.5 |

Number of Lexical Subcategories



Number of Phrasal Subcategories

# Learned Splits

- 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 |

# Learned Splits

- Relative adverbs (RBR):

| RBR-0 | further | lower | higher |
|---|---|---|---|
| RBR-1 | more | less | More |
| RBR-2 | earlier | Earlier | later |

- Cardinal Numbers (CD):

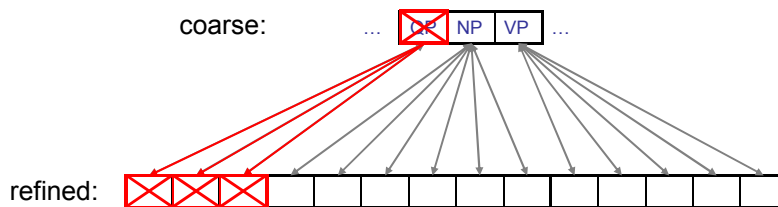| CD-7 | one | two | Three |
|---|---|---|---|
| CD-4 | 1989 | 1990 | 1988 |
| CD-11 | million | billion | trillion |
| CD-0 | 1 | 50 | 100 |
| CD-3 | 1 | 30 | 31 |
| CD-9 | 78 | 58 | 34 |

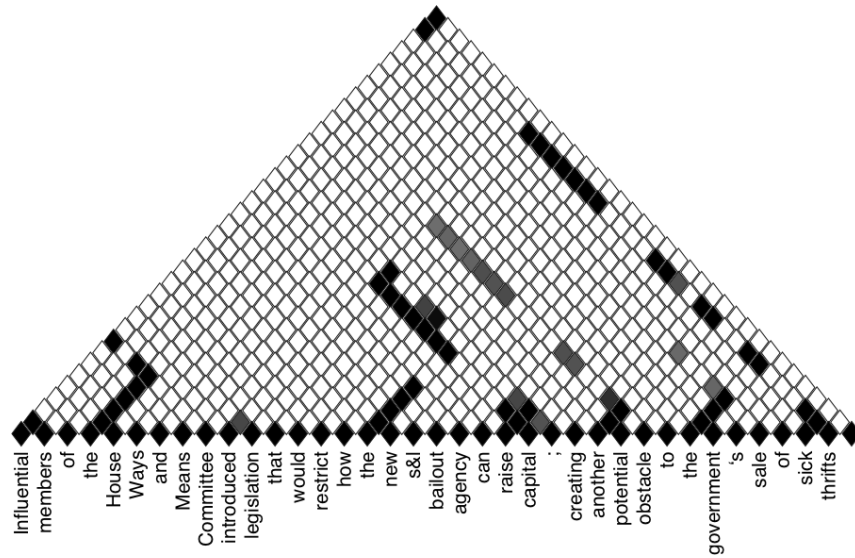# Coarse-to-Fine Inference

- Example: PP attachment



# Prune?

For each chart item $X[i,j]$, compute posterior probability:

$$\frac{P_{\text{IN}}(X, i, j) \cdot P_{\text{OUT}}(X, i, j)}{P_{\text{IN}}(root, 0, n)} \quad < \quad \textit{threshold}$$
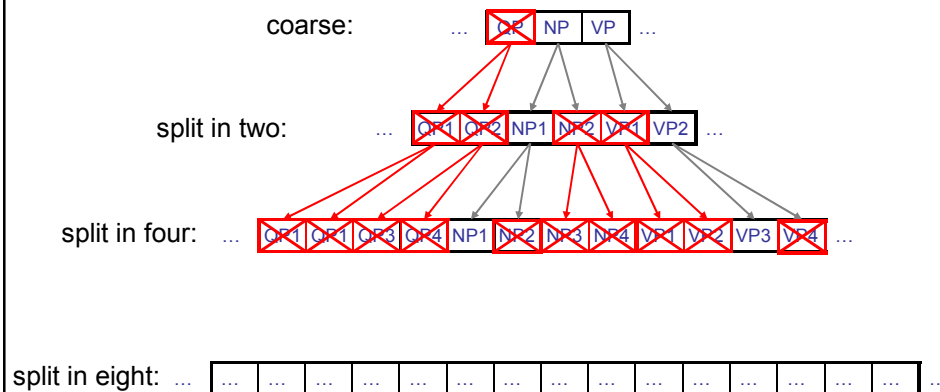
E.g. consider the span 5 to 12:

# Bracket Posteriors



# Hierarchical Pruning

# Final Results (Accuracy)

| | | ≤ 40 words F1 | all F1 |
|---|---|---|---|
| **ENG** | Charniak&Johnson '05 (generative) | 90.1 | 89.6 |
| | **Split / Merge** | **90.6** | **90.1** |
| **GER** | Dubey '05 | 76.3 | - |
| | **Split / Merge** | **80.8** | **80.1** |
| **CHN** | Chiang et al. '02 | 80.0 | 76.6 |
| | **Split / Merge** | **86.3** | **83.4** |

Still higher numbers from reranking / self-training methods