# Algorithms for NLP



# Parsing III

Maria Ryskina – CMU

Slides adapted from: Dan Klein – UC Berkeley
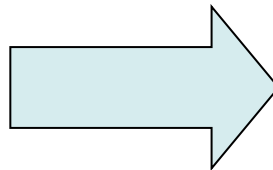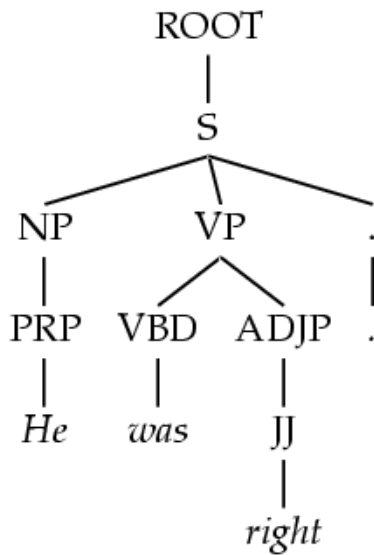
Taylor Berg-Kirkpatrick, Yulia Tsvetkov – CMU

# Learning PCFGs

# Treebank PCFGs

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

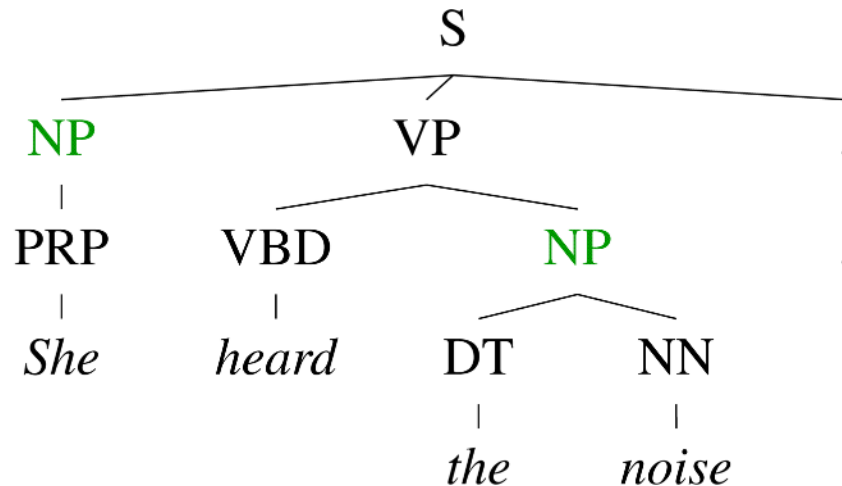ROOT → S                1

S → NP VP .             1

NP → PRP                1

VP → VBD ADJP           1

…..

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

# Conditional Independence?



S
├── NP
│   └── PRP
│       └── *She*
├── VP
│   ├── VBD
│   │   └── *heard*
│   └── NP
│       ├── DT
│       │   └── *the*
│       └── NN
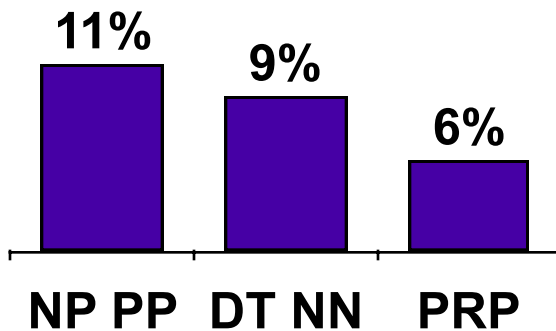│           └── *noise*
└── .
    └── .

- Not every NP expansion can fill every NP slot
  - A grammar with symbols like "NP" won't be context-free
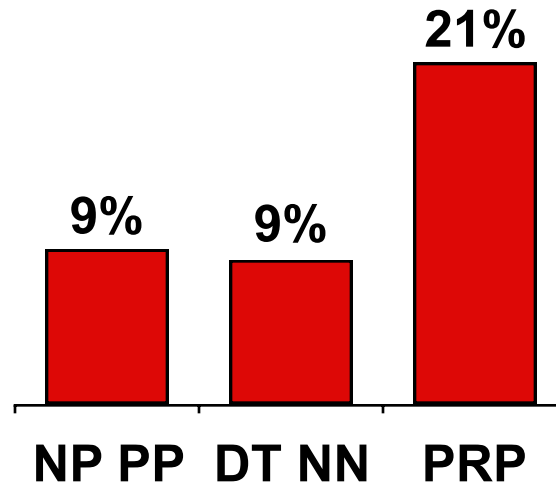  - Statistically, conditional independence too strong

# Non-Independence

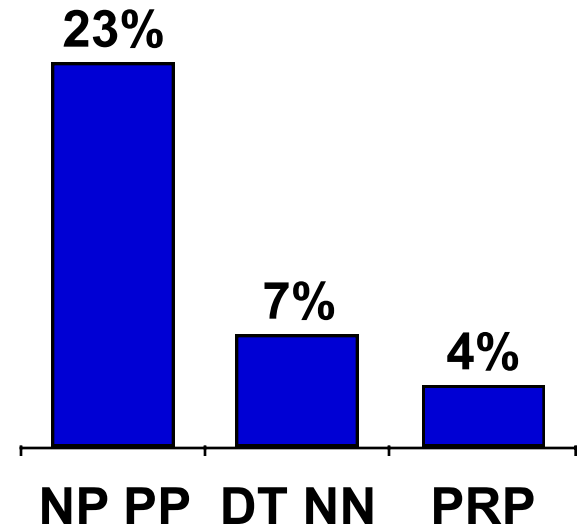- Independence assumptions are often too strong.

### All NPs

### NPs under S

### NPs under VP



| All NPs | | |
|---|---|---|
| 11% | 9% | 6% |
| NP PP | DT NN | PRP |

| NPs under S | | |
|---|---|---|
| 9% | 9% | 21% |
| NP PP | DT NN | PRP |

| NPs under VP | | |
|---|---|---|
| 23% | 7% | 4% |
| NP PP | DT NN | PRP |

- 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!

- Example: PP attachment

*They*

*raised*

*a    point    of    order*

- Example: PP attachment

They

V
|
raised

a     point     of     order

# Grammar Refinement

- Example: PP attachment

They

V                          NP

*raised*       DT    NN

              *a*    *point*      *of*   *order*

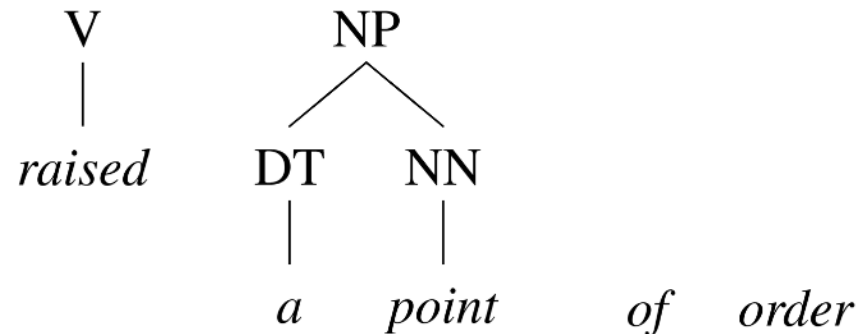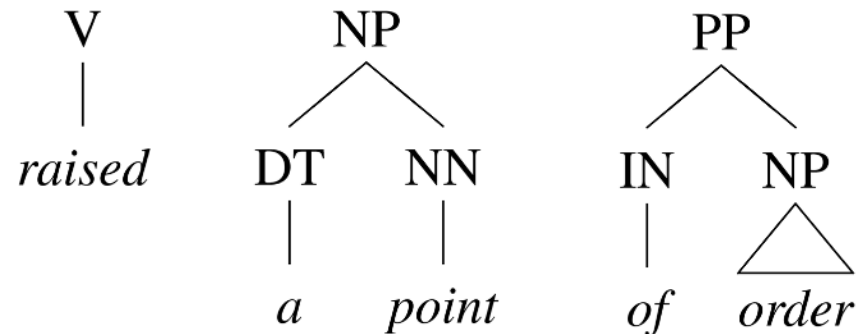# Grammar Refinement

- Example: PP attachment



They

V — raised

NP
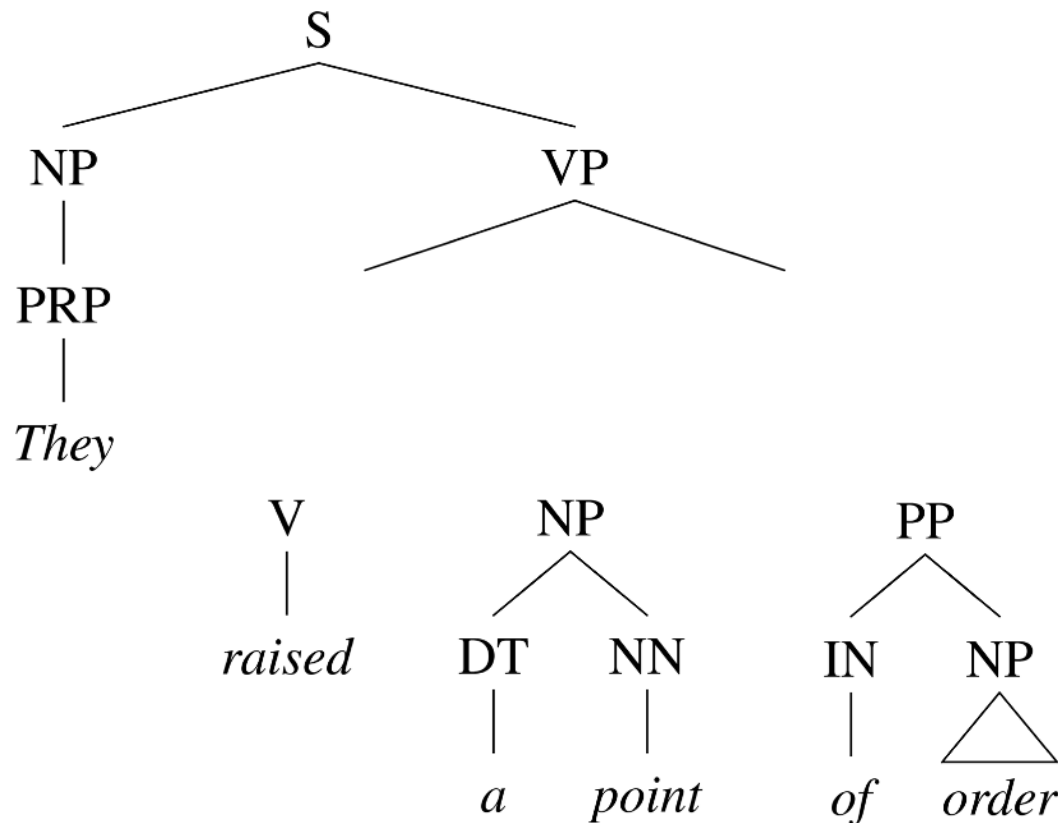├── DT — a
└── NN — point

PP
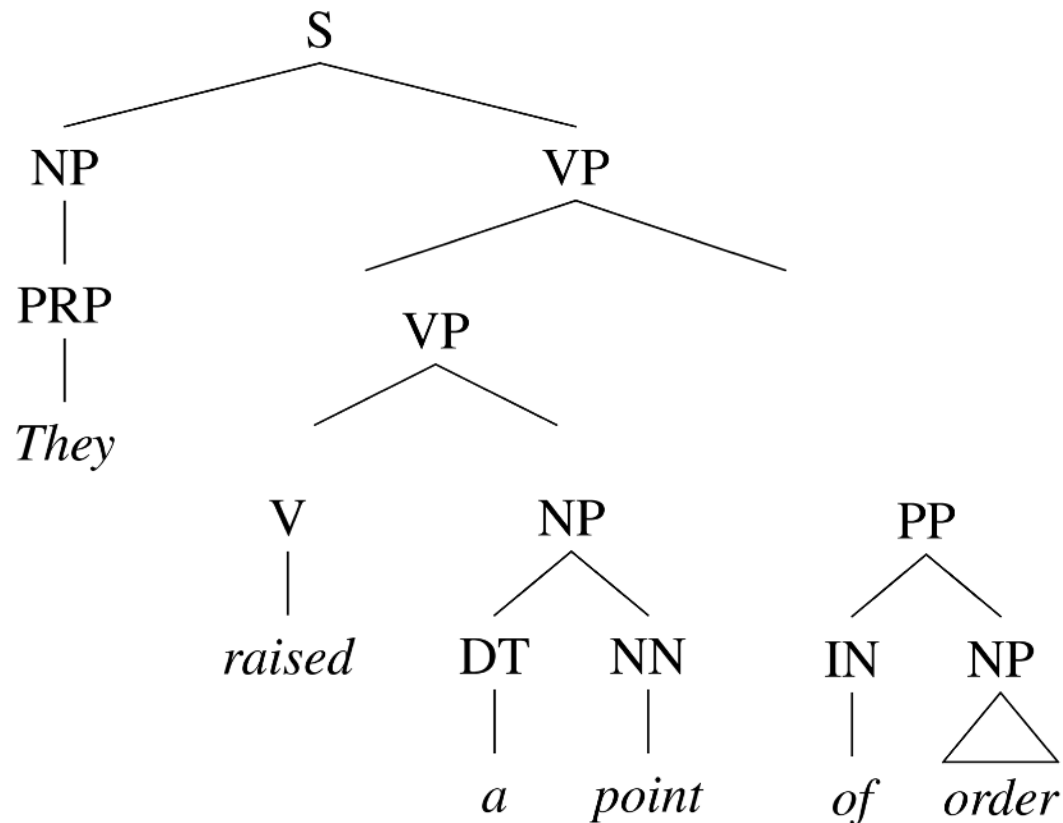├── IN — of
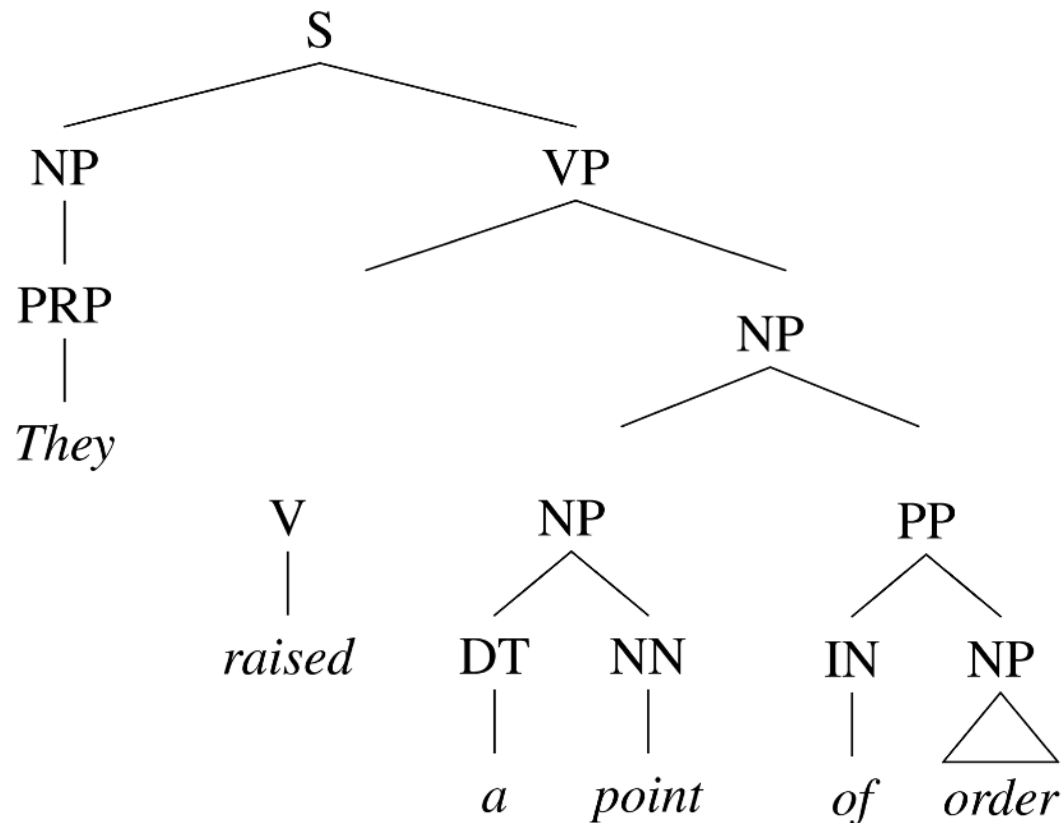└── NP — order

# Grammar Refinement

- Example: PP attachment

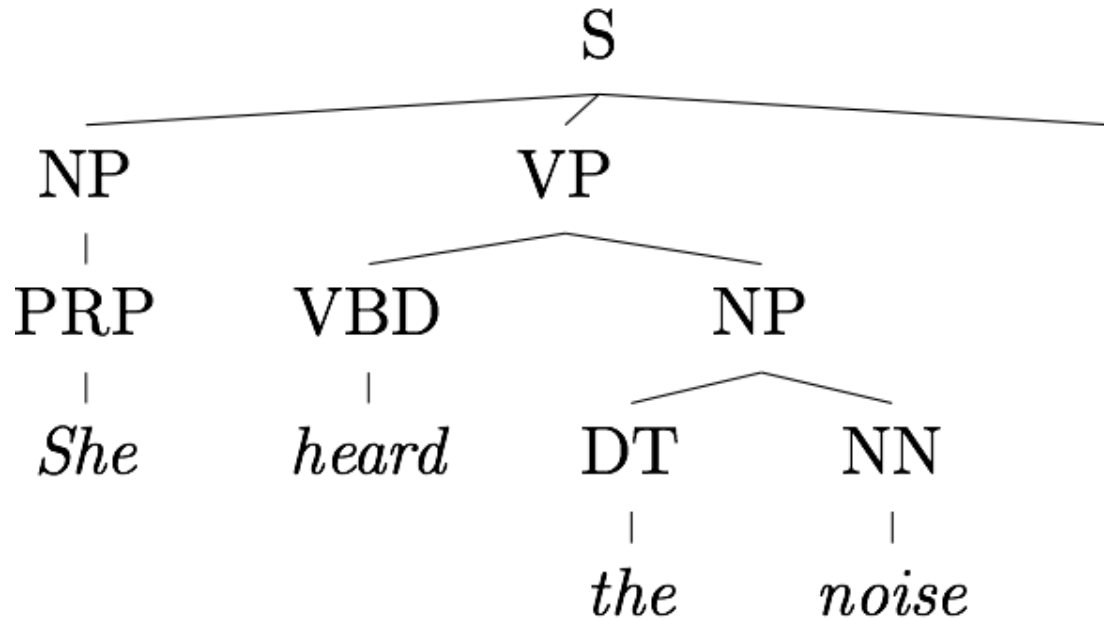# Grammar Refinement

- Example: PP attachment

# Grammar Refinement
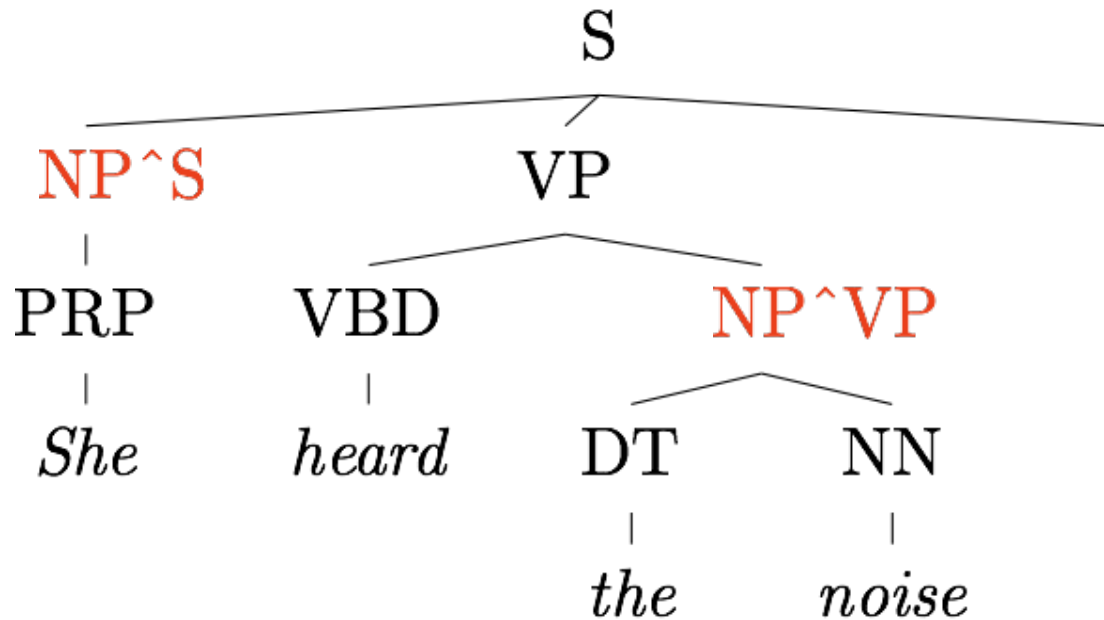
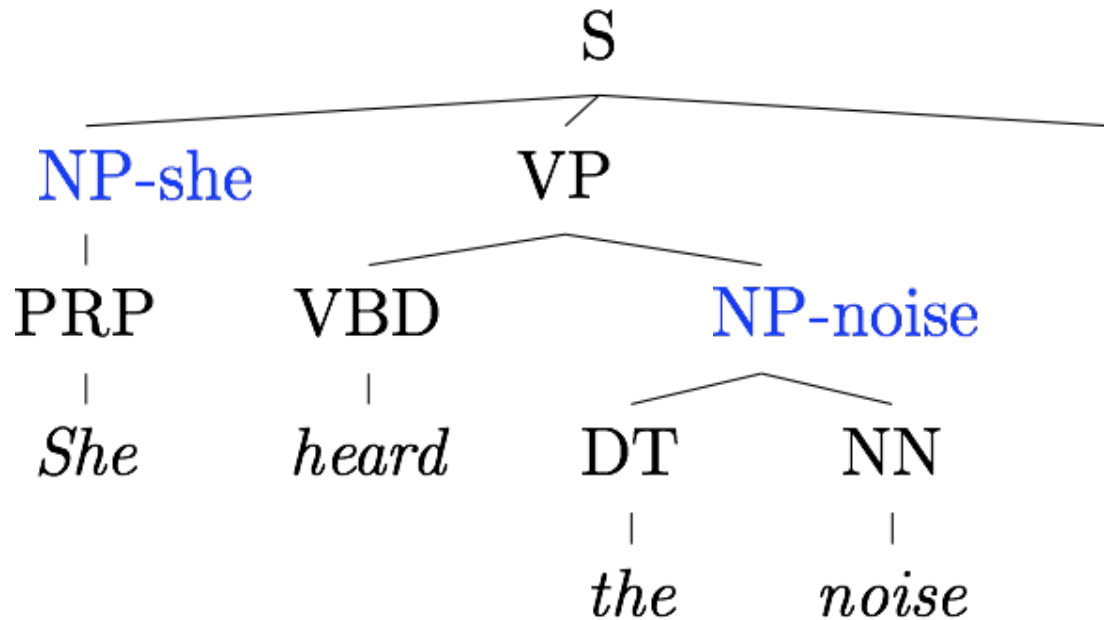- Example: PP attachment

# Grammar Refinement

# Grammar Refinement



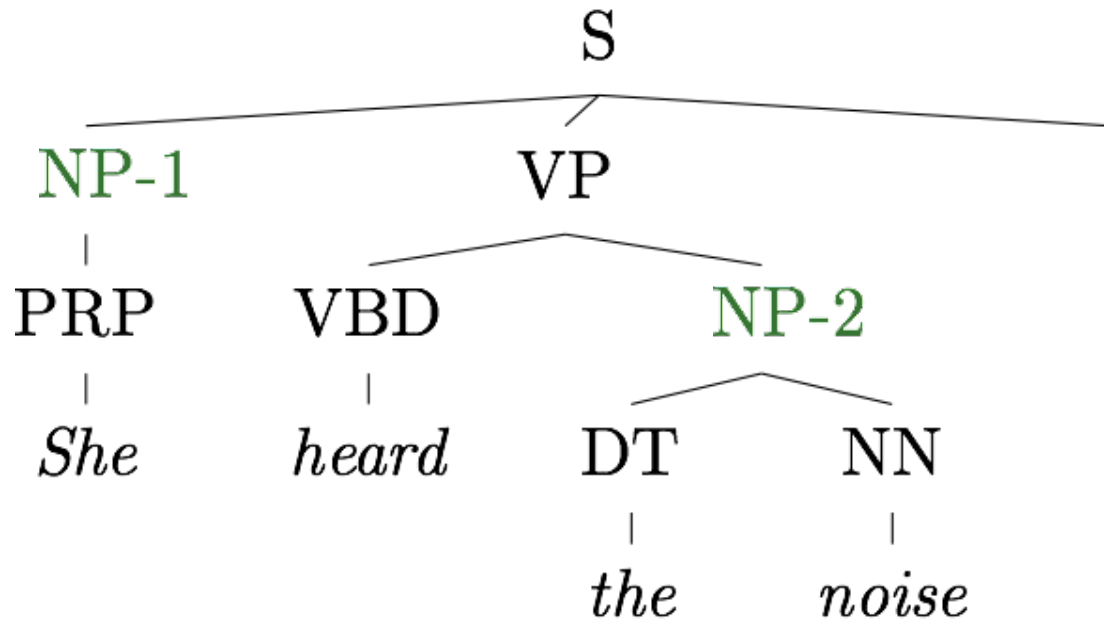- Structural Annotation [Johnson '98, Klein&Manning '03]

# Grammar Refinement



- Structural Annotation [Johnson '98, Klein&Manning '03]
- Lexicalization [Collins '99, Charniak '00]
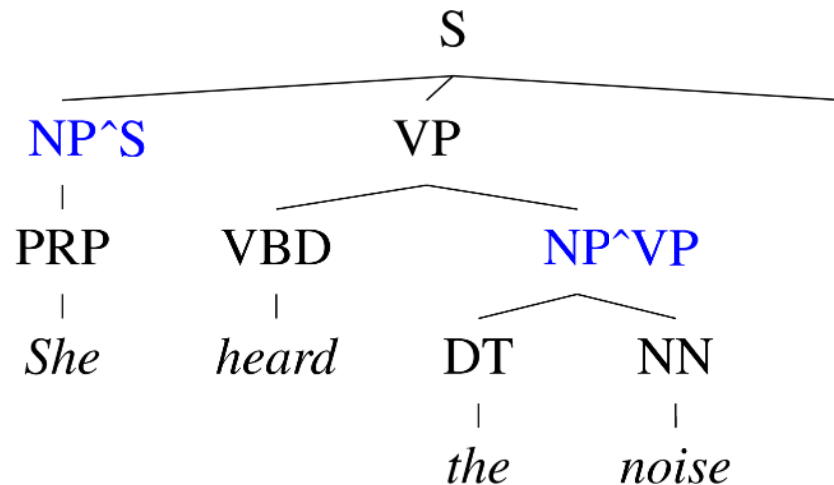
# Grammar Refinement



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

# Structural Annotation

# The Game of Designing a Grammar



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

# Typical Experimental Setup

- Corpus: Penn Treebank, WSJ

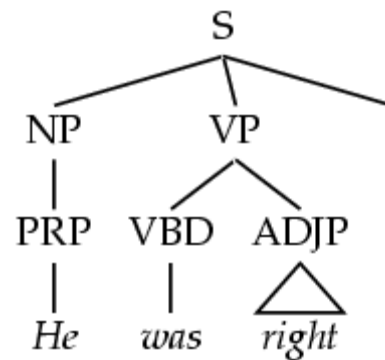| Training: | sections | 02-21 |
| Development: | section | 22 (here, first 20 files) |
| Test: | section | 23 |

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

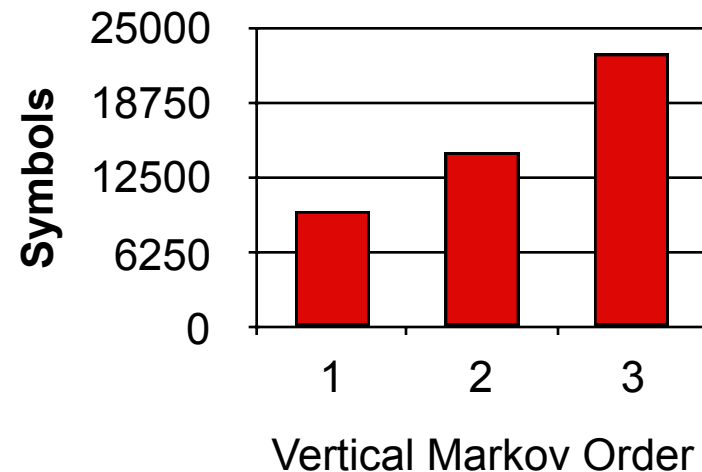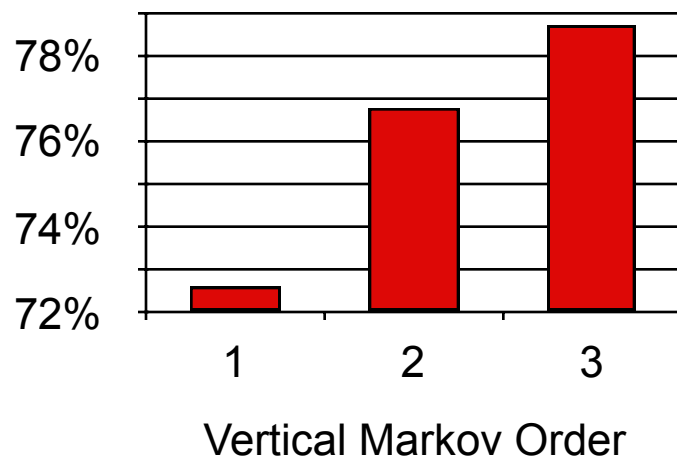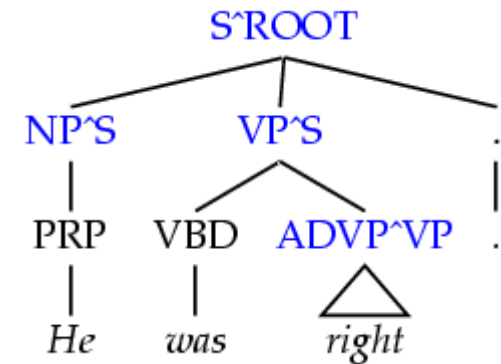- Here: also size – number of symbols in grammar.

# Vertical Markovization

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

Order 1

Order 2





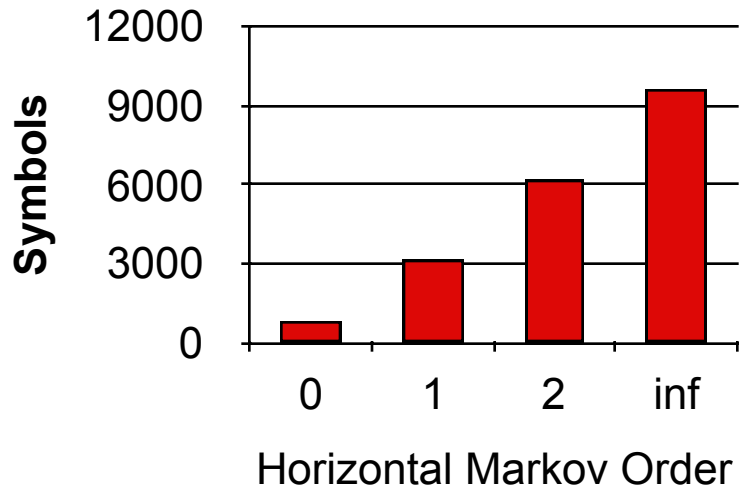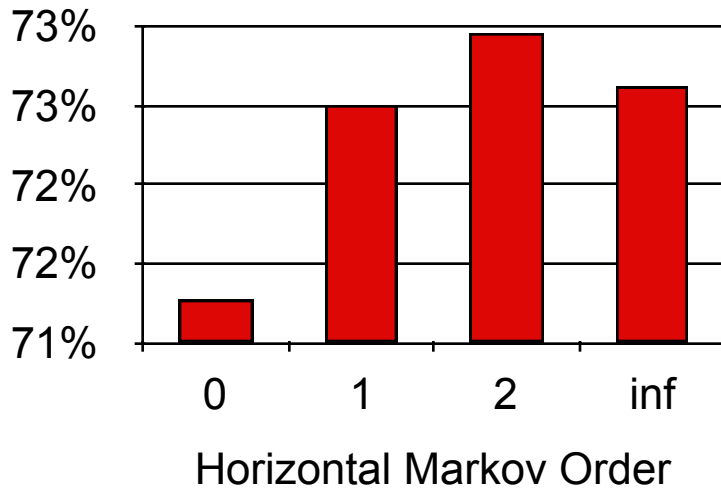Vertical Markov Order

Vertical Markov Order

# Horizontal Markovization

Order 1  Order ∞

NP

DT JJNN NN

NP

DT  JJNN  NN

v=1,h=∞

NP

DT  JJNN  NN

v=1,h=∞

NP

NP

DT JJNN NN

v=1,h=∞

NP

DT

NP

DT JJNN NN

$v=1, h=\infty$

NP

DT    @NP[DT]

NP

DT JJ NN NN

v=1, h=∞

NP

DT @NP[DT]

JJ @NP[DT,JJ]

# Binarization / Markovization

NP
DT  JJNN  NN

v=1,h=∞

NP
DT  @NP[DT]
JJ  @NP[DT,JJ]
NN  @NP[DT,JJ,NN]

# Binarization / Markovization

NP
├── DT
├── JJ
├── NN
└── NN

v=1,h=∞

NP
├── DT
└── @NP[DT]
    ├── JJ
    └── @NP[DT,JJ]
        ├── NN
        └── @NP[DT,JJ,NN]
            └── NN

NP

DT JJNN NN

v=1,h=∞

NP

DT @NP[DT]

JJ @NP[DT,JJ]

NN @NP[DT,JJ,NN]

NN

v=1,h=1

NP

DT @NP[DT]

JJ @NP[...,JJ]

NN @NP[...,NN]

NN

# Binarization / Markovization

NP
├─ DT
├─ JJ
├─ NN
└─ NN

**v=1,h=∞**

NP
├─ DT
└─ @NP[DT]
   ├─ JJ
   └─ @NP[DT,JJ]
      ├─ NN
      └─ @NP[DT,JJ,NN]
         └─ NN

**v=1,h=1**

NP
├─ DT
└─ @NP[DT]
   ├─ JJ
   └─ @NP[...,JJ]
      ├─ NN
      └─ @NP[...,NN]
         └─ NN

**v=1,h=0**

NP
├─ DT
└─ @NP
   ├─ JJ
   └─ @NP
      ├─ NN
      └─ @NP
         └─ NN

# Binarization / Markovization

NP

DT  JJNN  NN

v=2,h=∞

NP^VP

DT^NP  @NP^VP[DT]

JJ^NP  @NP^VP[DT,JJ]

NN^NP  @NP^VP[DT,JJ,NN]

NN^NP

v=2,h=1

NP^VP

DT^NP  @NP^VP[DT]

JJ^NP  @NP^VP[…,JJ]

NN^NP  @NP^VP[…,NN]

NN^NP

v=2,h=0

NP^VP

DT^NP  @NP

JJ^NP  @NP

NN^NP  @NP

NN^NP

# Unary Splits

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



| Annotation | F1 | Size |
|---|---|---|
| Base | 77.8 | 7.5K |
| UNARY | 78.3 | 8.0K |

# 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 |

# 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 |

# 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 |

# 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 |

# A Fully Annotated (Unlex) Tree

# 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.

# Efficient Parsing for Structural Annotation

# Grammar Projections

Coarse Grammar

```
        NP
       /  \
     DT    @NP
          /   \
        JJ     @NP
              /   \
            NN     @NP
                    |
                    NN
```

Fine Grammar

```
           NP^VP
          /     \
    DT^NP        @NP^VP[DT]
                /          \
          JJ^NP             @NP^VP[...,JJ]
                          /              \
                    NN^NP                 @NP^VP[...,NN]
                                            |
                                          NN^NP
```

# Grammar Projections

## Coarse Grammar

NP
├─ DT
└─ @NP
   ├─ JJ
   └─ @NP
      ├─ NN
      └─ @NP
         └─ NN

$$NP \rightarrow DT\ @NP$$

## Fine Grammar

NP^VP
├─ DT^NP
└─ @NP^VP[DT]
   ├─ JJ^NP
   └─ @NP^VP[...,JJ]
      ├─ NN^NP
      └─ @NP^VP[...,NN]
         └─ NN^NP

# Grammar Projections

Coarse Grammar

Fine Grammar

NP
├── DT
└── @NP
    ├── JJ
    └── @NP
        ├── NN
        └── @NP
            └── NN

NP^VP
├── DT^NP
└── @NP^VP[DT]
    ├── JJ^NP
    └── @NP^VP[...,JJ]
        ├── NN^NP
        └── @NP^VP[...,NN]
            └── NN^NP

NP → DT @NP

NP^VP → DT^NP @NP^VP[DT]

# Grammar Projections

Coarse Grammar

Fine Grammar



$$NP \rightarrow DT\ @NP$$

$$NP^\wedge VP \rightarrow DT^\wedge NP$$
$$@NP^\wedge VP[DT]$$

*Note: X-Bar Grammars are projections with rules like XP → Y @X or XP → @X Y or @X → X*

# Grammar Projections

Coarse Symbols

NP

@NP

DT

Fine Symbols

NP^VP
NP^S
@NP^VP[DT]
@NP^S[DT]
@NP^VP[...,JJ]
@NP^S[...,JJ]
DT^NP

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

$$\frac{\mathrm{P_{IN}}(X, i, j) \cdot \mathrm{P_{OUT}}(X, i, j)}{\mathrm{P_{IN}}(root, 0, n)}$$

E.g. consider the span 5 to 12:

coarse:　　　　… | QP | NP | VP | …

# Coarse-to-Fine Pruning

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

$$\frac{P_{IN}(X,i,j) \cdot P_{OUT}(X,i,j)}{P_{IN}(root,0,n)} \quad < \quad \textit{threshold}$$

E.g. consider the span 5 to 12:

coarse:    … | QP | NP | VP | …

# Coarse-to-Fine Pruning

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

$$\frac{\text{P}_{\text{IN}}(X, i, j) \cdot \text{P}_{\text{OUT}}(X, i, j)}{\text{P}_{\text{IN}}(root, 0, n)} \quad < \quad \textit{threshold}$$

E.g. consider the span 5 to 12:

coarse:        ...   | QP | NP | VP |   ...

fine:   | | | | | | | | | | | | | | | | |

# Coarse-to-Fine Pruning

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

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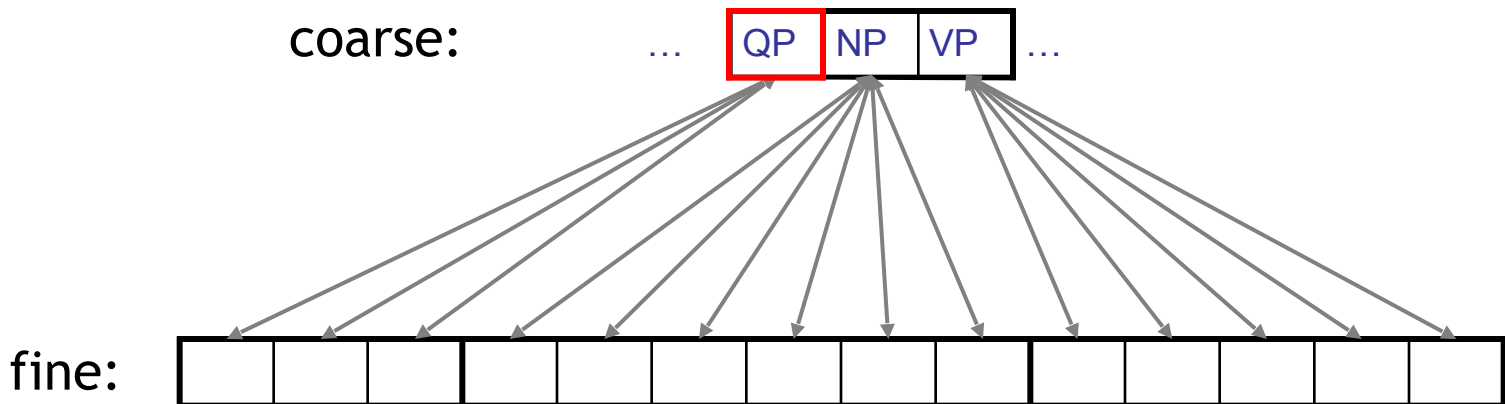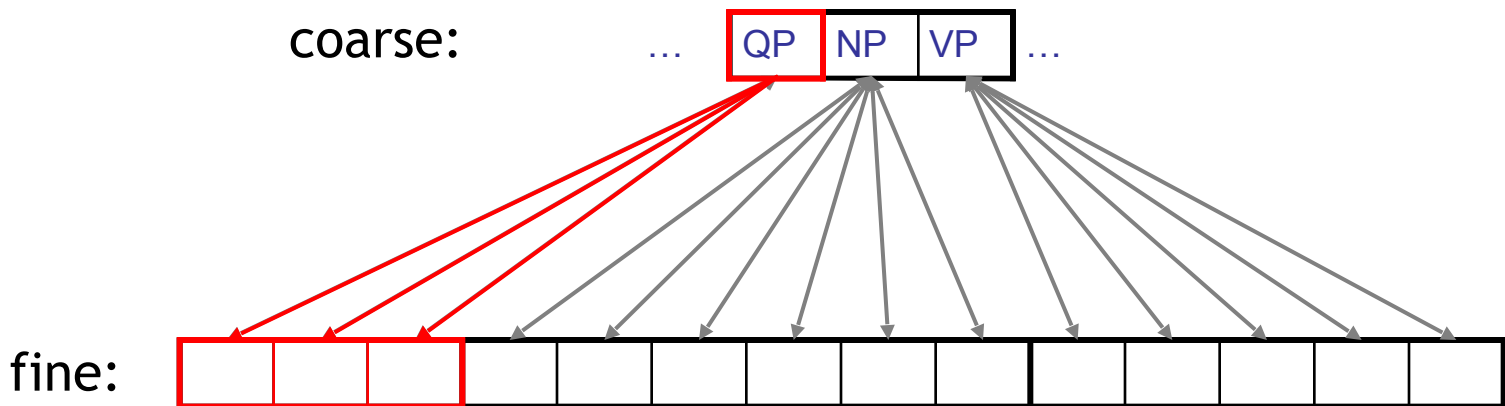
E.g. consider the span 5 to 12:

# Coarse-to-Fine Pruning

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

$$\frac{\mathrm{P_{IN}}(X, i, j) \cdot \mathrm{P_{OUT}}(X, i, j)}{\mathrm{P_{IN}}(root, 0, n)} \quad < \quad \textit{threshold}$$

E.g. consider the span 5 to 12:

| NP→DET N | 0.8 | NP→N | 0.2 |
|----------|-----|------|-----|
| DET→a | 0.6 | DET→the | 0.4 |
| N→apple | 0.8 | N→orange | 0.2 |

$$\beta_{DET}(1,1) = P(the \mid DET_{11}, G) = P(DET \to the \mid G) = 0.4$$

$$\beta_N(2,2) = P(N \to orange \mid G) = 0.2$$

$$\beta_{NP}(1,2) = P(NP \to DET \cdot N)\beta_{DET}(1,1)\beta_N(2,2)$$
$$= 0.8 \qquad \times 0.4 \qquad \times 0.2$$
$$\beta_{NP}(1,2) = 0.064$$

NP, DET, N, the, orange

$$\beta_S(1,m) = P(S \to w_1, \ldots, w_m \mid G)$$

# Calculating outside probability

▶ The joint probability corresponding to the yellow, red and blue areas, <span style="color:red">assuming</span> $N^j$ was the <span style="color:red">L</span> child of some non-terminal:

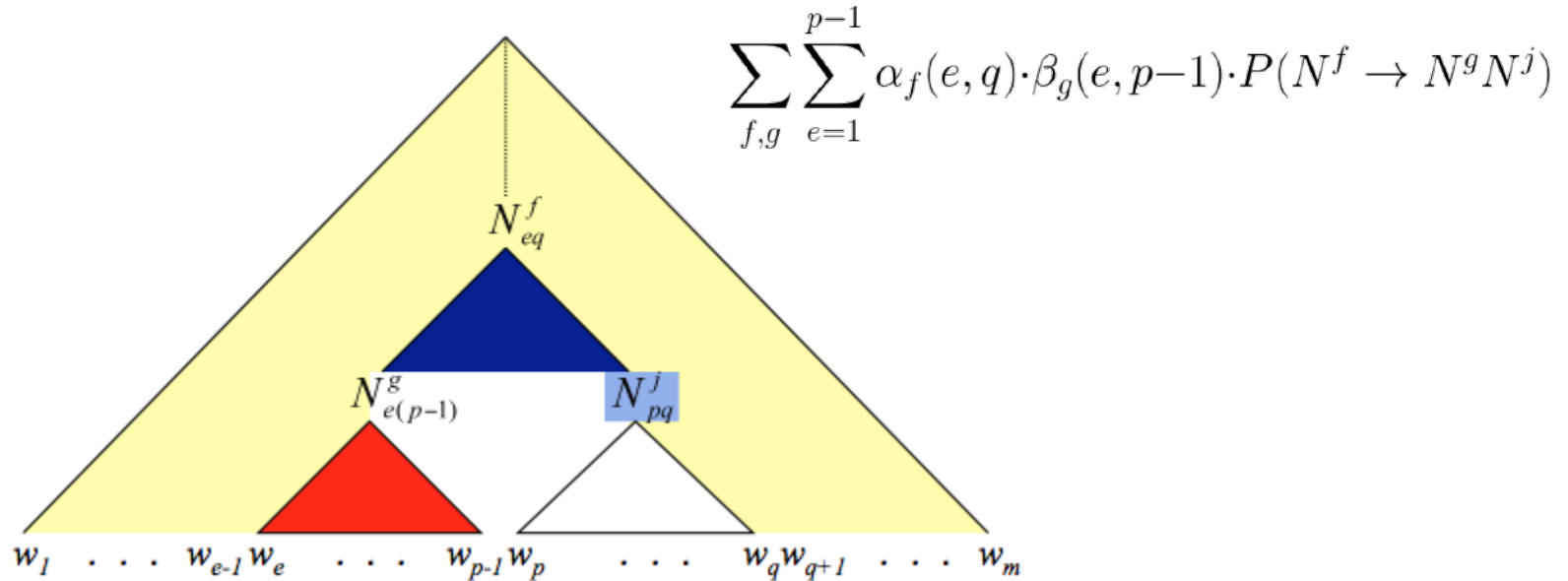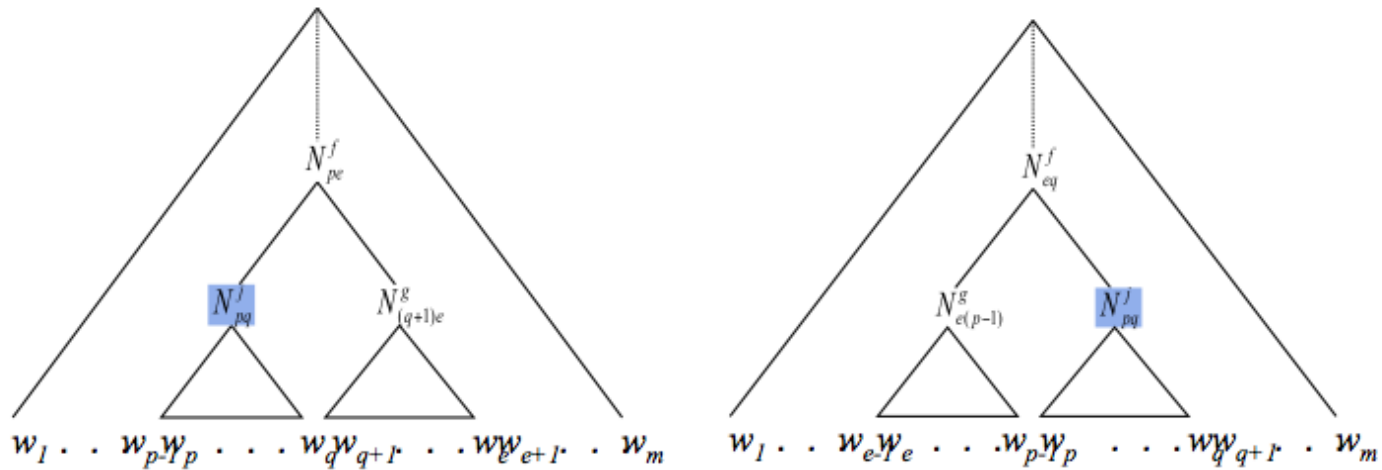$$\sum_{f,g} \sum_{e=q+1}^{m} \alpha_f(p,e) \cdot \beta_g(q+1,e) \cdot P(N^f \to N^j N^g)$$

▶ The joint probability corresponding to the yellow, red and blue areas, assuming $N^j$ was the R child of some non-terminal:

$$\sum_{f,g}\sum_{e=1}^{p-1}\alpha_f(e,q)\cdot\beta_g(e,p-1)\cdot P(N^f \rightarrow N^g N^j)$$

# Calculating outside probability

▶ The joint final joint probability (the sum over the L and R cases):

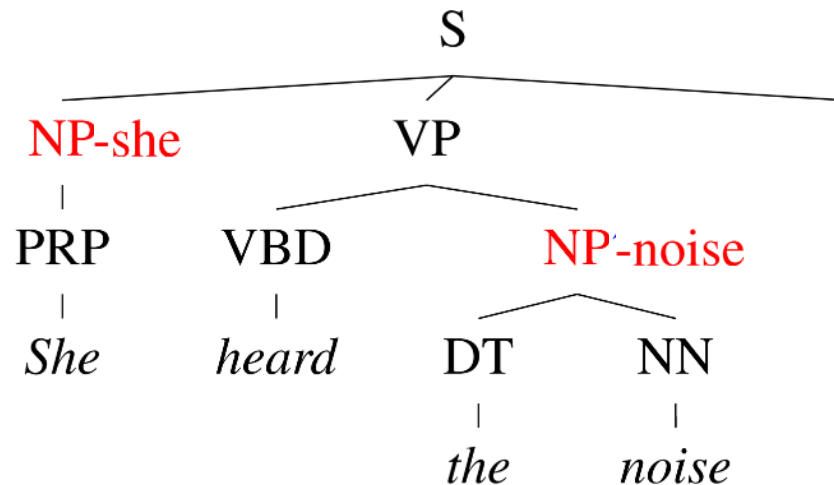

$$\alpha_j(p,q) = \sum_{f,g} \sum_{e=q+1}^{m} \alpha_f(p,e) \cdot \beta_g(q+1,e) \cdot P(N^f \to N^j N^g) + \sum_{f,g} \sum_{e=1}^{p-1} \alpha_f(e,q) \cdot \beta_g(e,p-1) \cdot P(N^f \to N^g N^j)$$
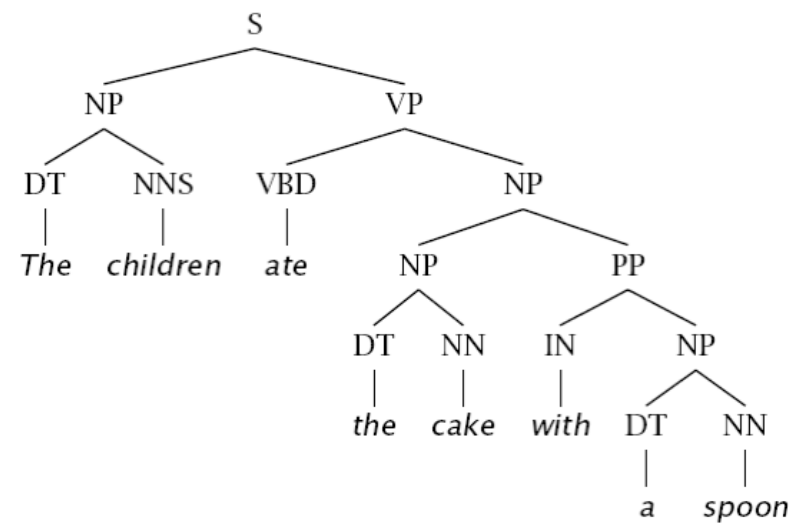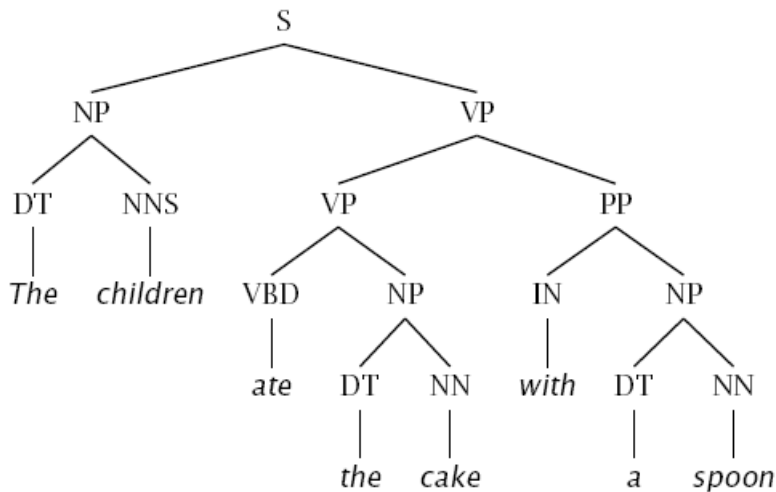
# Lexicalization

# 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]

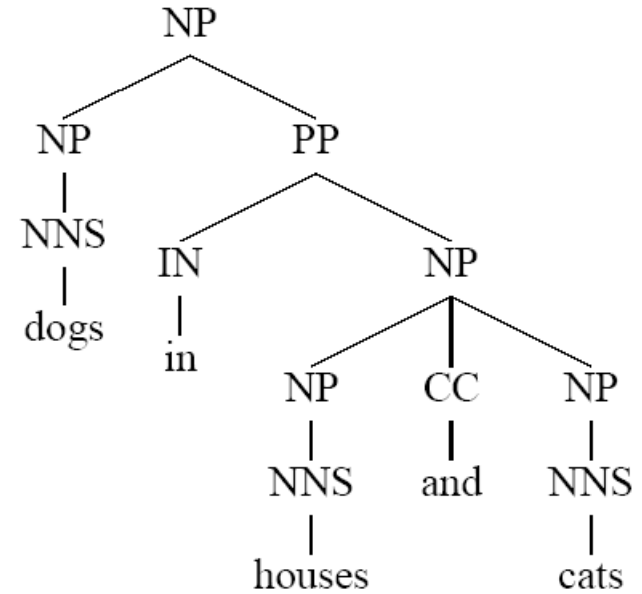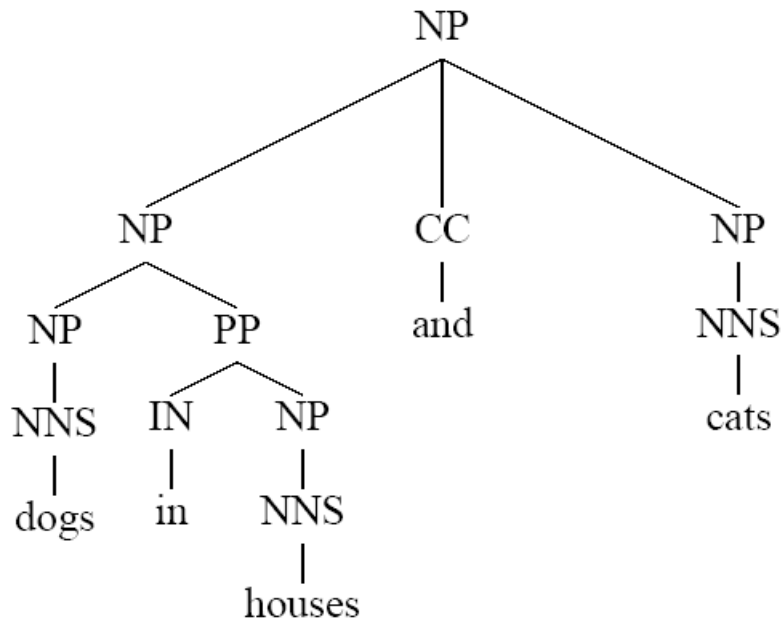# Problems with PCFGs



- If we do no annotation, these trees differ only in one rule:
  - VP → VP PP
  - NP → NP PP
- Parse will go one way or the other, regardless of words
- We addressed this in one way with unlexicalized grammars (how?)
- Lexicalization allows us to be sensitive to specific words

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

# Lexicalized Trees

- Add "head words" 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

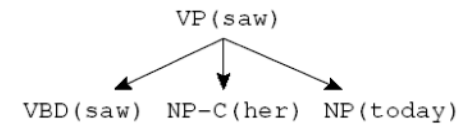$$\text{VP(saw)} \rightarrow \text{VBD(saw)} \ \text{NP-C(her)} \ \text{NP(today)}$$

- 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]

VP(saw)

VBD(saw)

Choose a head tag and word

VP(saw)

VBD(saw)        {NP-C(    )}

Choose a complement bag

VP(saw)

VBD(saw)   NP-C(    )   NP(        )

Generate children (incl. adjuncts)

VP(saw)

VBD(saw)   NP-C(her)   NP(today)

Recursively derive children

# Lexicalized CKY

`(VP->VBD...NP •)[saw]`

`(VP->VBD •)[saw]`    `NP[her]`

X[h]

Y[h]  Z[h']

i    h    k    h'    j

```
bestScore(X,i,j,h)

  if (j = i+1)

    return tagScore(X,s[i])

  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)
```

# Efficient Parsing for Lexical Grammars

# Quartic Parsing

- Turns out, you can do (a little) better [Eisner 99]



X[h]

Y[h]    Z[h']

i       h       k       h'      j

X[h]

Y[h]    Z

i       h       k               j

- Gives an $O(n^4)$ algorithm
- Still prohibitive in practice if not pruned

# 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>.
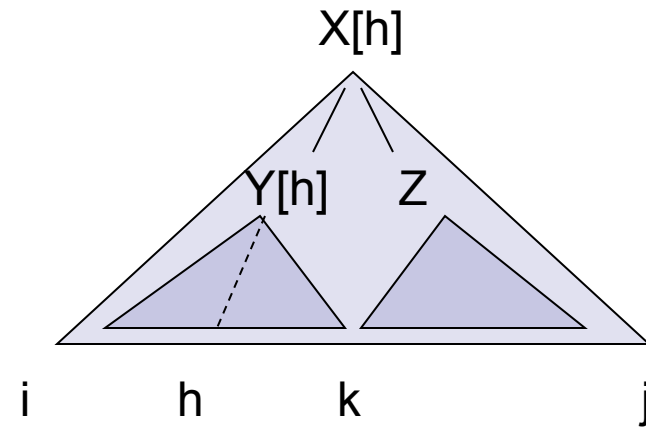  - If we keep K hypotheses at each span, then we do at most $O(nK^2)$ work per span (why?)
  - Keeps things more or less cubic (and in practice is more like linear!)

- 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, coarse-to-fine 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

# 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)
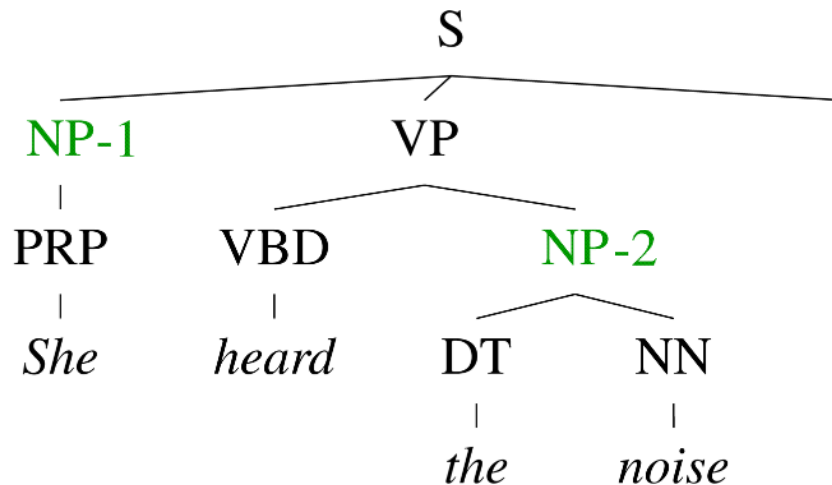  - McClosky et al 06 – 92.1 F1 (gen + rerank + self-train)

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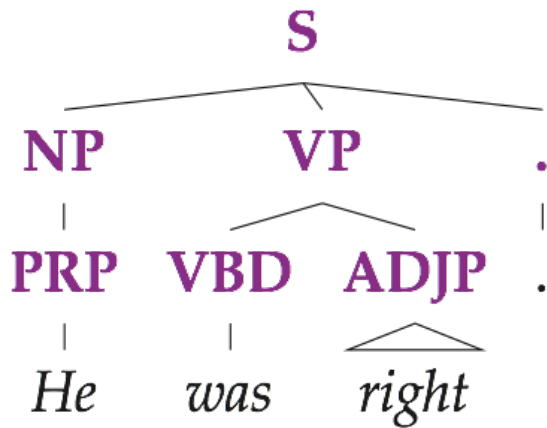
# Latent Variable PCFGs

- Annotation refines base treebank symbols to improve statistical fit of the grammar
  - Parent annotation [Johnson '98]
  - Head lexicalization [Collins '99, Charniak '00]
  - Automatic clustering?

```
                    S
        ┌───────────┼───────────┐
       NP          VP           .
        │       ┌───┴───┐       │
       PRP    VBD     ADJP      .
        │       │     ┌─△─┐
       He      was    right
```

Parse Tree $T$
Sentence $w$

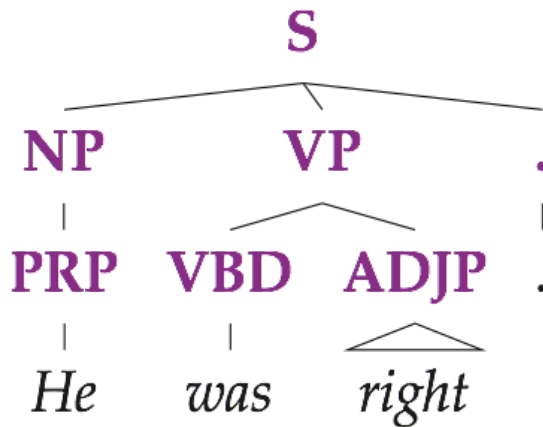# Latent Variable Grammars


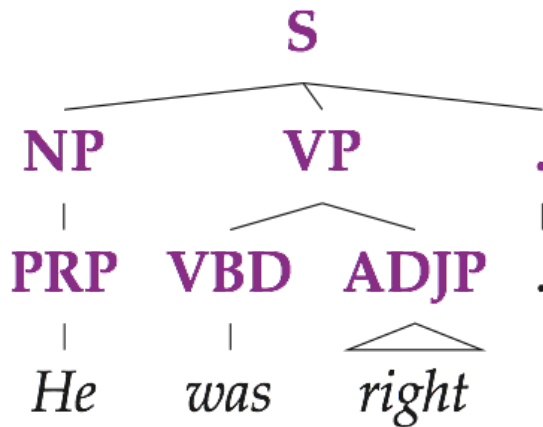
Parse Tree $T$
Sentence $w$

Derivations $t : T$

# Latent Variable Grammars



Parse Tree $T$
Sentence $w$

Derivations $t : T$

Parameters $\theta$

EM algorithm:

$$\text{S}[X_1]$$

$$\text{NP}[X_2] \quad\quad \text{VP}[X_4] \quad\quad .[X_7]$$

$$\text{PRP}[X_3] \quad \text{VBD}[X_5] \quad \text{ADJP}[X_6] \quad\quad .$$

$$He \quad\quad was \quad\quad right$$

# Learning Latent Annotations

EM algorithm:

- Brackets are known
- Base categories are known
- Only induce subcategories

$$S[X_1]$$

$$NP[X_2] \quad VP[X_4] \quad .[X_7]$$

$$PRP[X_3] \quad VBD[X_5] \quad ADJP[X_6] \quad .$$

*He*     *was*     *right*

# Learning Latent Annotations

**EM algorithm:**

- Brackets are known
- Base categories are known
- Only induce subcategories

# Learning Latent Annotations

EM algorithm:

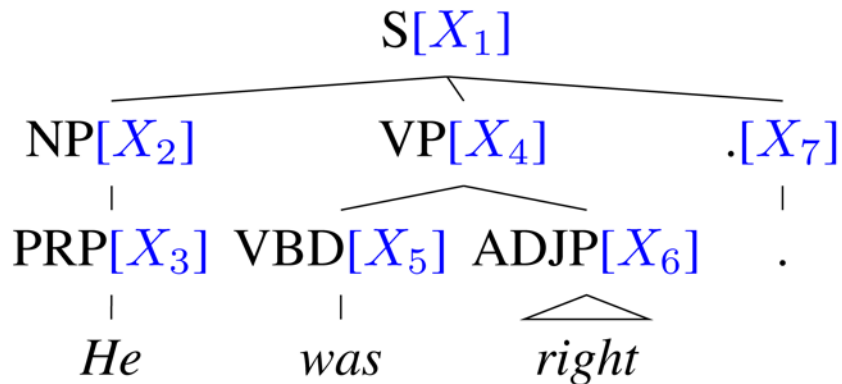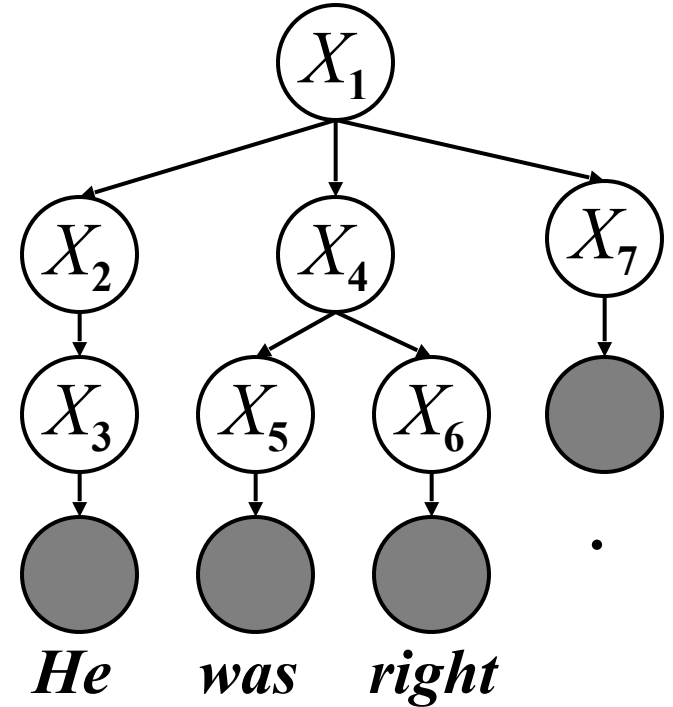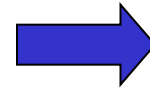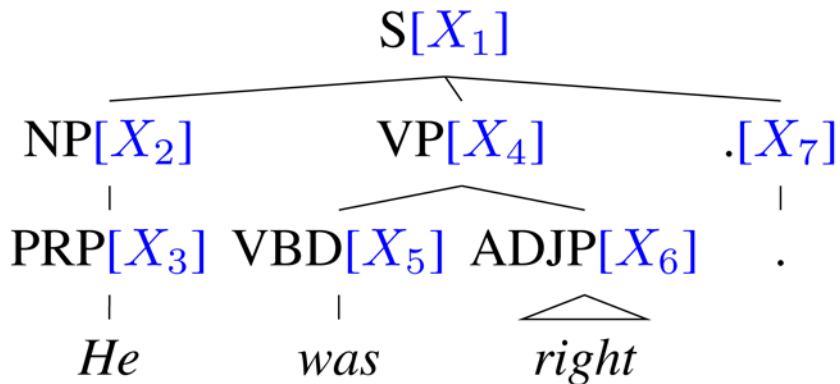- Brackets are known
- Base categories are known
- Only induce subcategories

$$\text{S}[X_1]$$

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

DT

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

# Refinement of the DT tag

DT

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

| 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) |

DT-1        DT-2        DT-3        DT-4

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

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)

# Hierarchical Estimation Results



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

- Splitting all categories equally is wasteful:

# 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

# Adaptive Splitting Results

# Adaptive Splitting Results

# Adaptive Splitting Results



Parsing accuracy (F1) vs. Total Number of grammar s[...]

Legend:
- 50% Merging (blue)
- Hierarchical Training (red)
- Flat Training (black)

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

# Number of Phrasal Subcategories

# Number of Lexical 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 |

# Final Results (Accuracy)

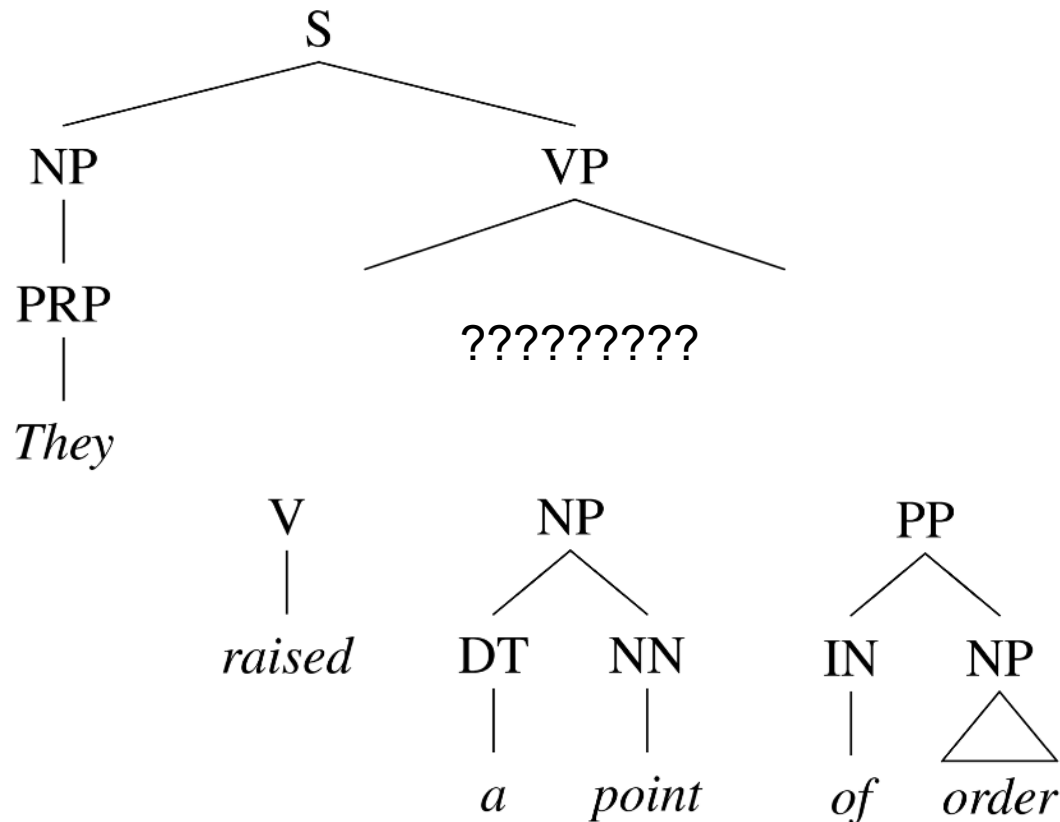|   |   | ≤ 40 words F1 | all F1 |
|---|---|---|---|
| **ENG** | Charniak&Johnson '05 (generative) | 90.1 | 89.6 |
|   | **Split / Merge** | **90.6** | **90.1** |
| **GER** | Dubey '05 | 76.3 | - |
|   | **Split / Merge** | **80.8** | **80.1** |
| **CHN** | Chiang et al. '02 | 80.0 | 76.6 |
|   | **Split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

Still higher numbers from reranking / self-training methods

# Efficient Parsing for Hierarchical Grammars

# Coarse-to-Fine Inference

- Example: PP attachment

coarse:          ...  | QP | NP | VP |  ...

coarse:          …  QP  NP  VP  …

# Hierarchical Pruning

coarse:

... | QP | NP | VP | ...

split in two:

... | QP1 | QP2 | NP1 | NP2 | VP1 | VP2 | ...

coarse:  … | QP | NP | VP | …

split in two:  … | QP 1 | QP 2 | NP 1 | NP2 | VP1 | VP2 | …

# Hierarchical Pruning



coarse: ... | QP | NP | VP | ...

split in two: ... | QP 1 | QP 2 | NP 1 | NP2 | VP1 | VP2 | ...

split in four: ... | QP 1 | QP 1 | QP 3 | QP 4 | NP 1 | NP2 | NP3 | NP4 | VP1 | VP2 | VP3 | VP4 | ...

# Hierarchical Pruning



coarse:  ... | QP | NP | VP | ...

split in two:  ... | QP1 | QP2 | NP1 | NP2 | VP1 | VP2 | ...

split in four:  ... | QP1 | QP1 | QP3 | QP4 | NP1 | NP2 | NP3 | NP4 | VP1 | VP2 | VP3 | VP4 | ...

# Hierarchical Pruning



coarse: ... | QP | NP | VP | ...

split in two: ... | QP1 | QP2 | NP1 | NP2 | VP1 | VP2 | ...

split in four: ... | QP1 | QP1 | QP3 | QP4 | NP1 | NP2 | NP3 | NP4 | VP1 | VP2 | VP3 | VP4 | ...

split in eight: ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ...

Influential members of the House Ways and Means Committee introduced legislation that would restrict how the new s&l bailout agency can raise capital ; creating another potential obstacle to the government 's sale of sick thrifts .
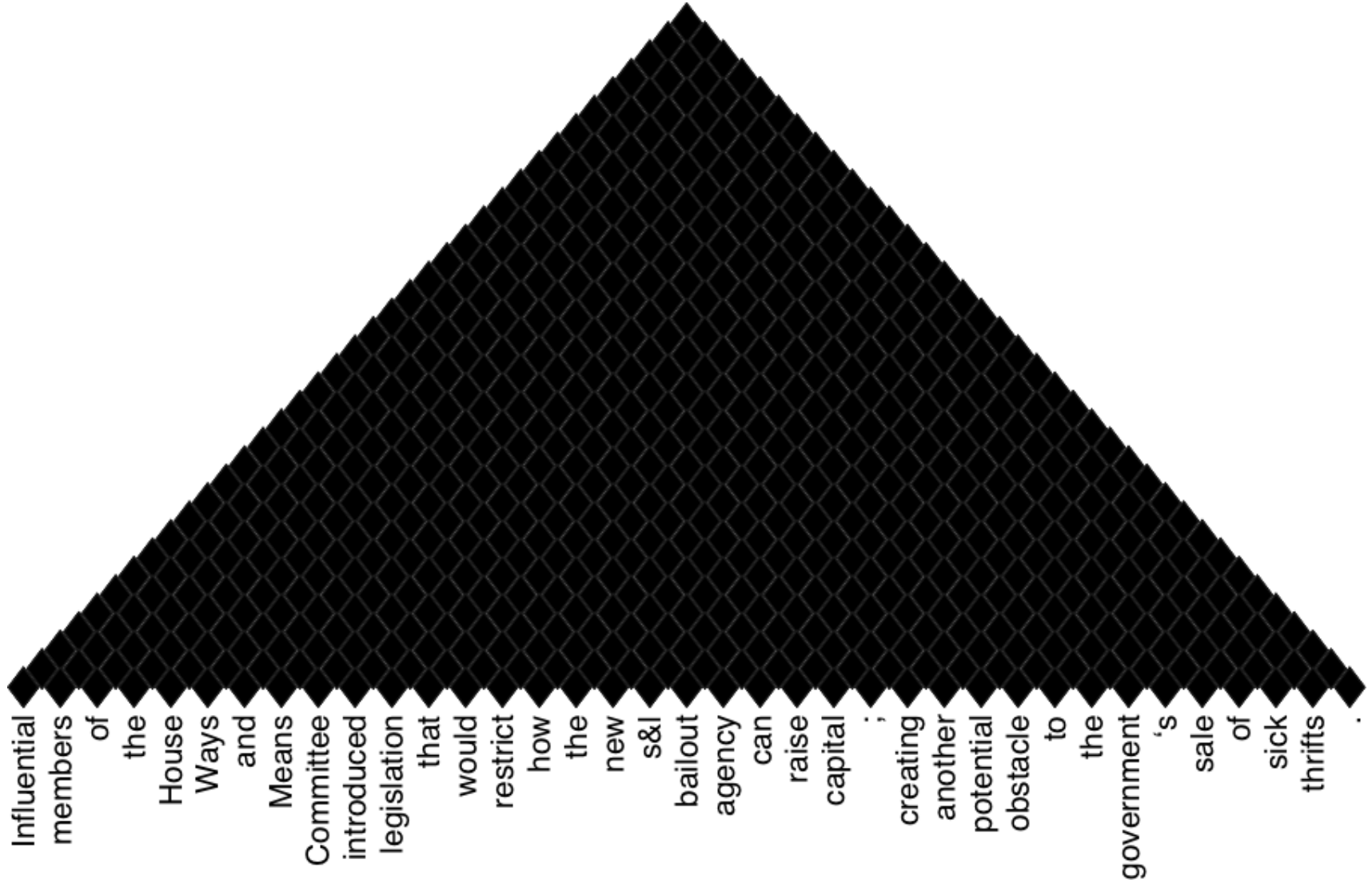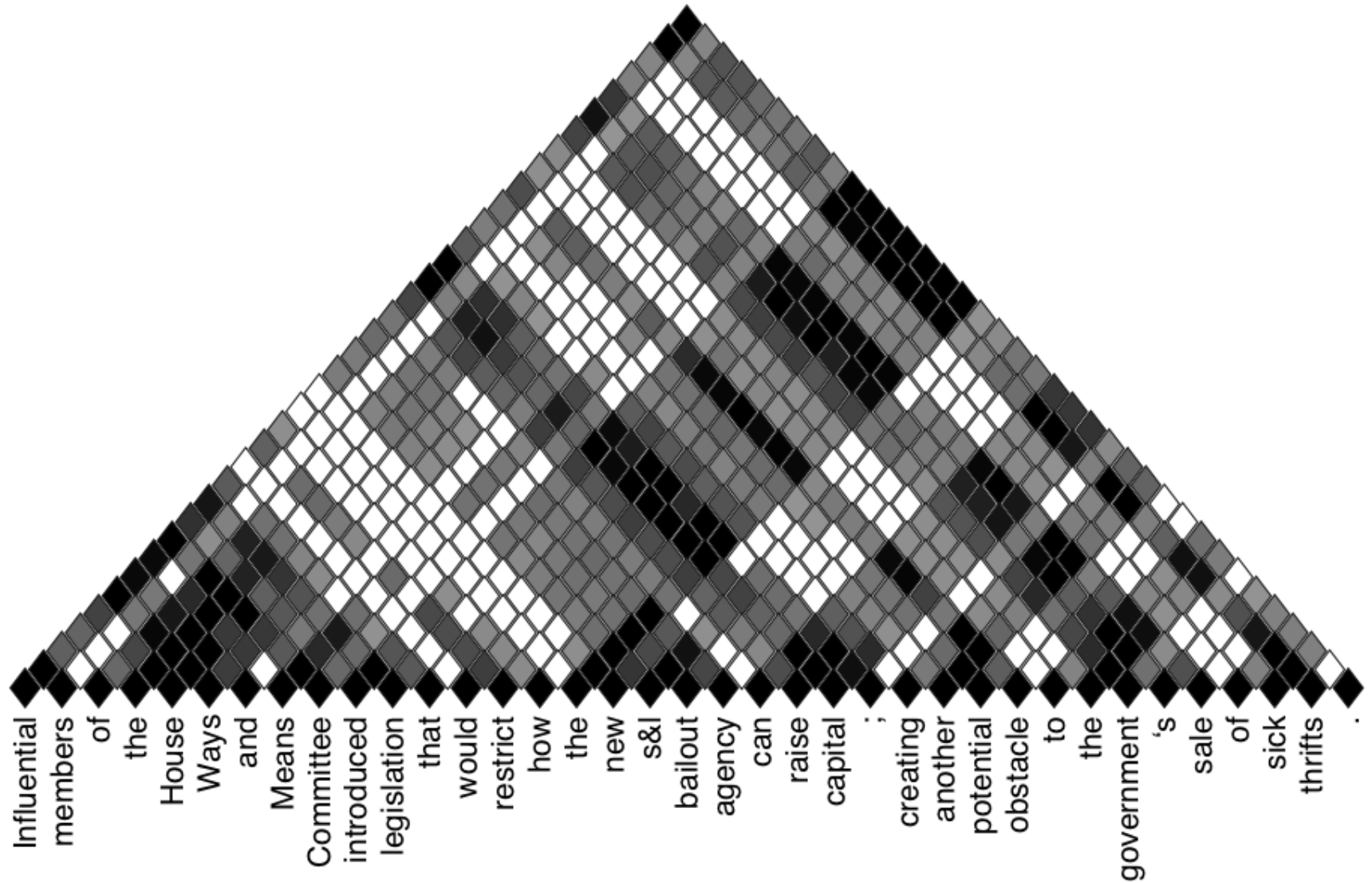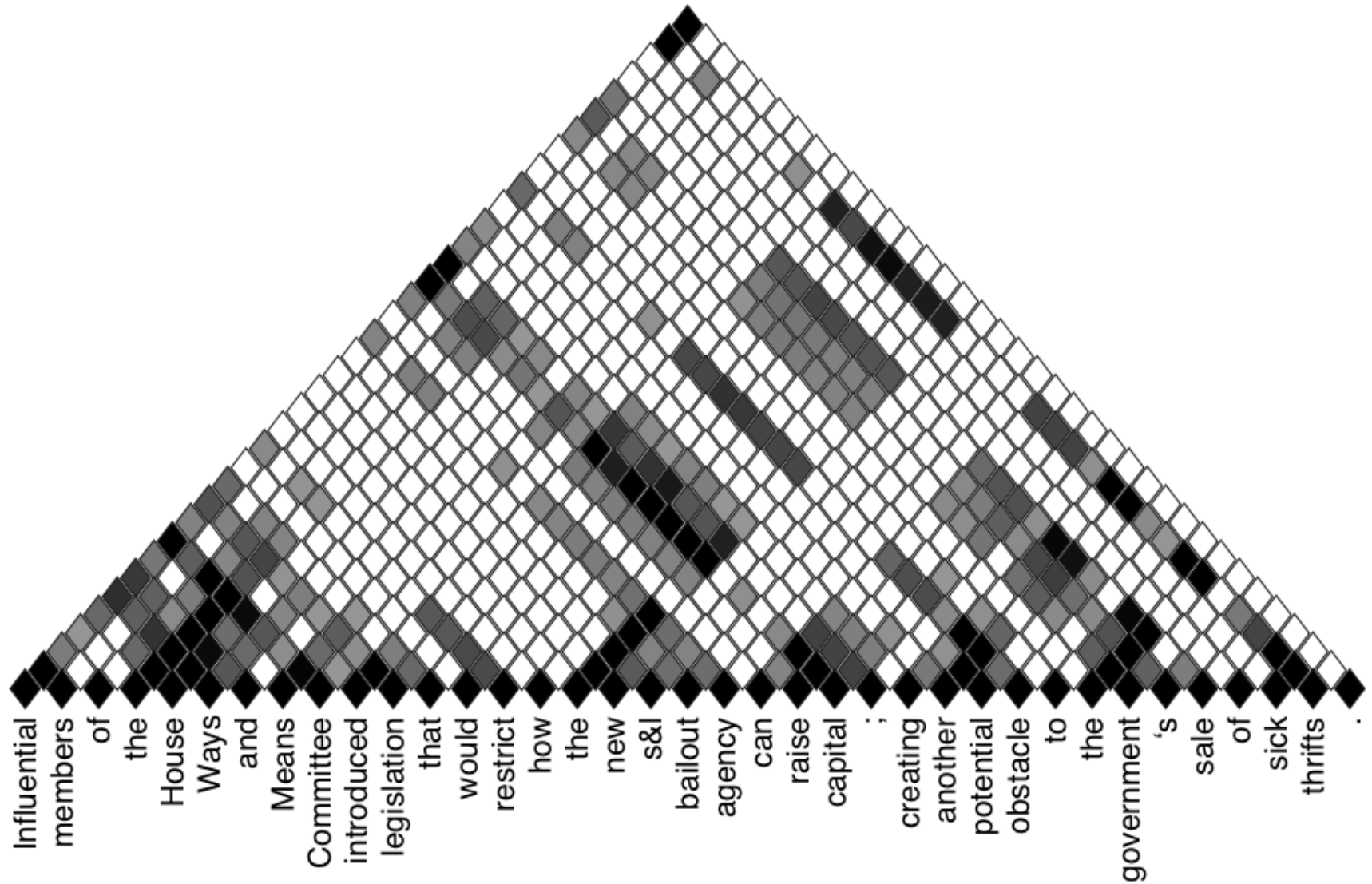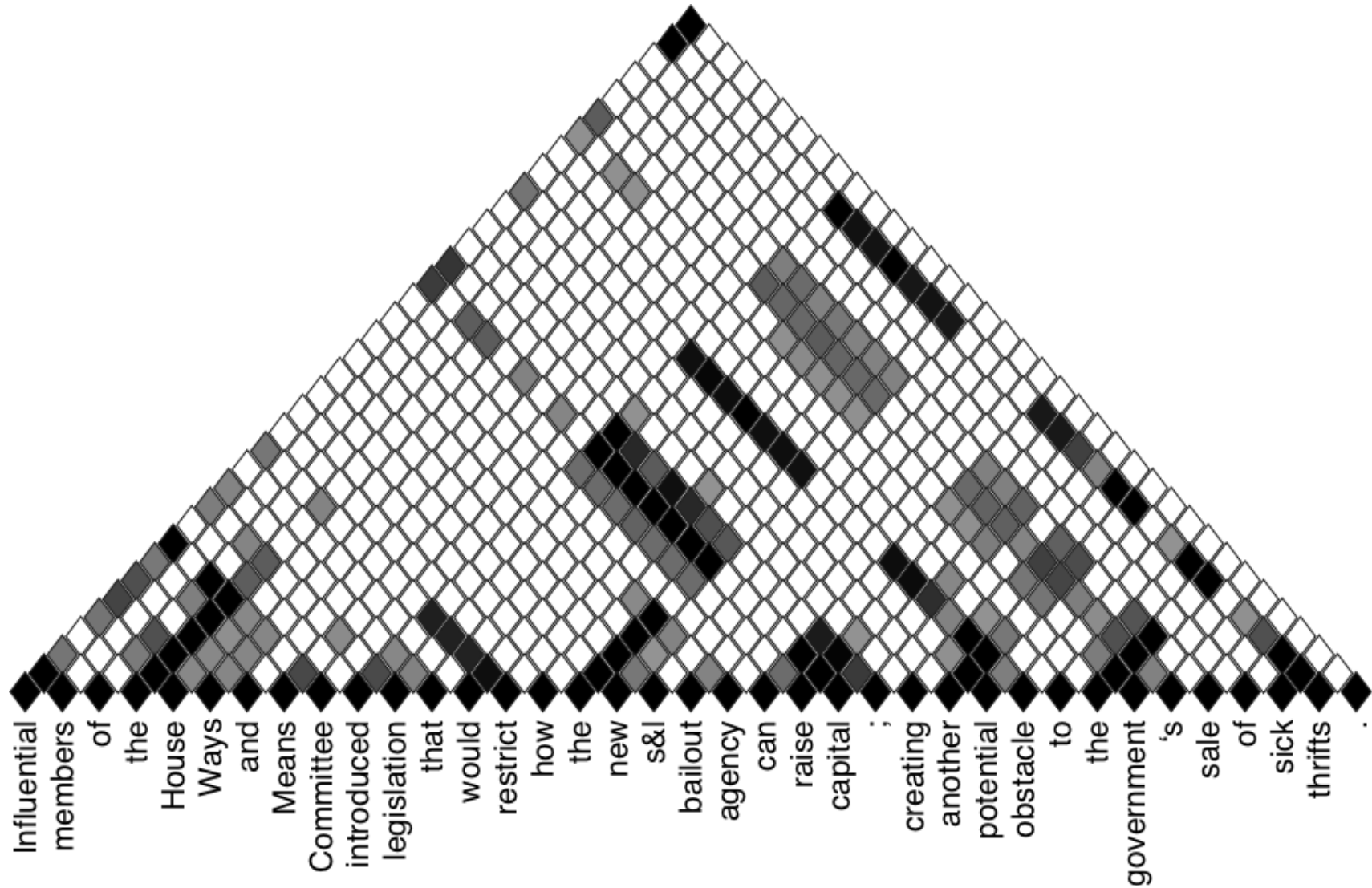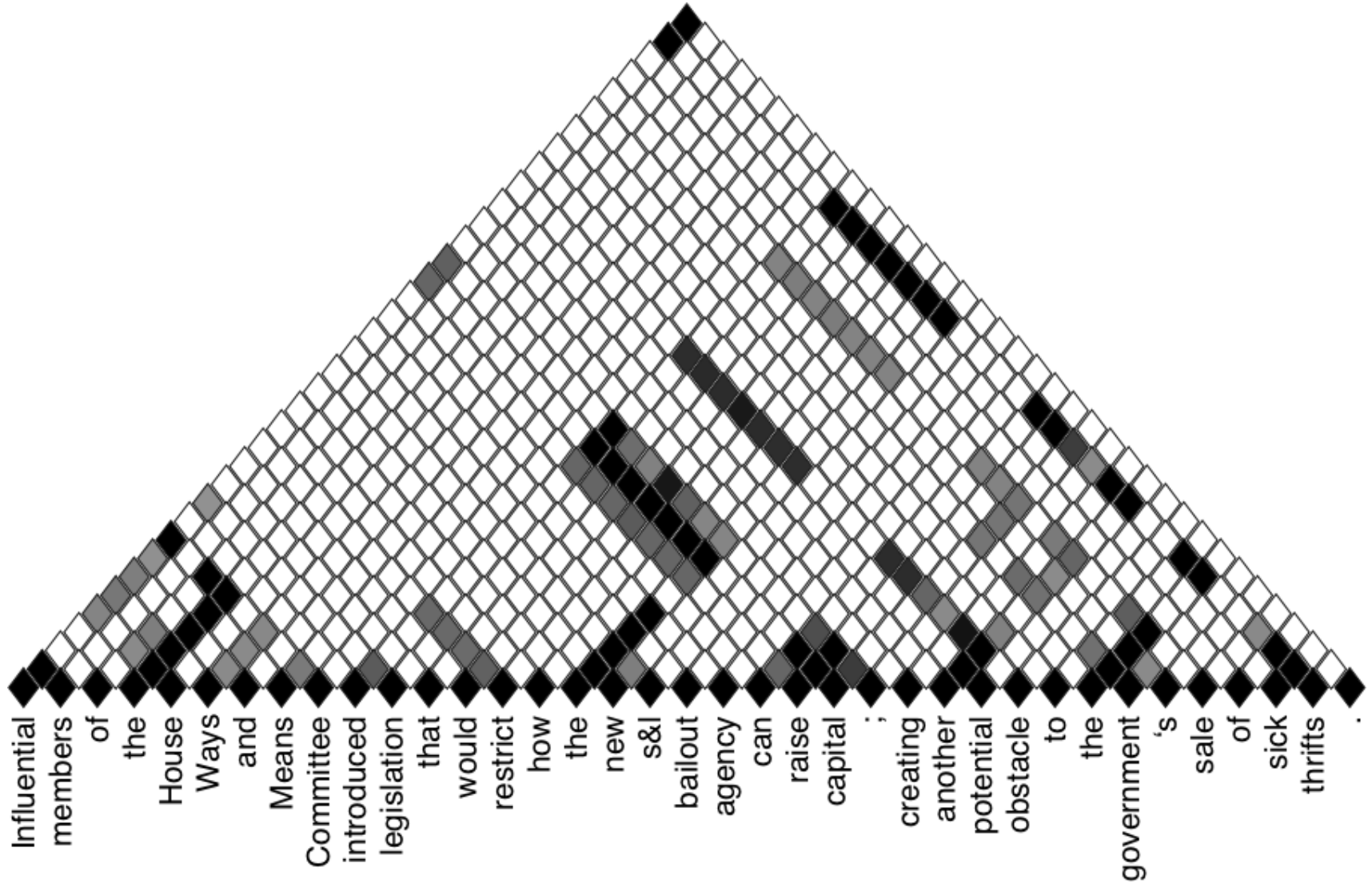
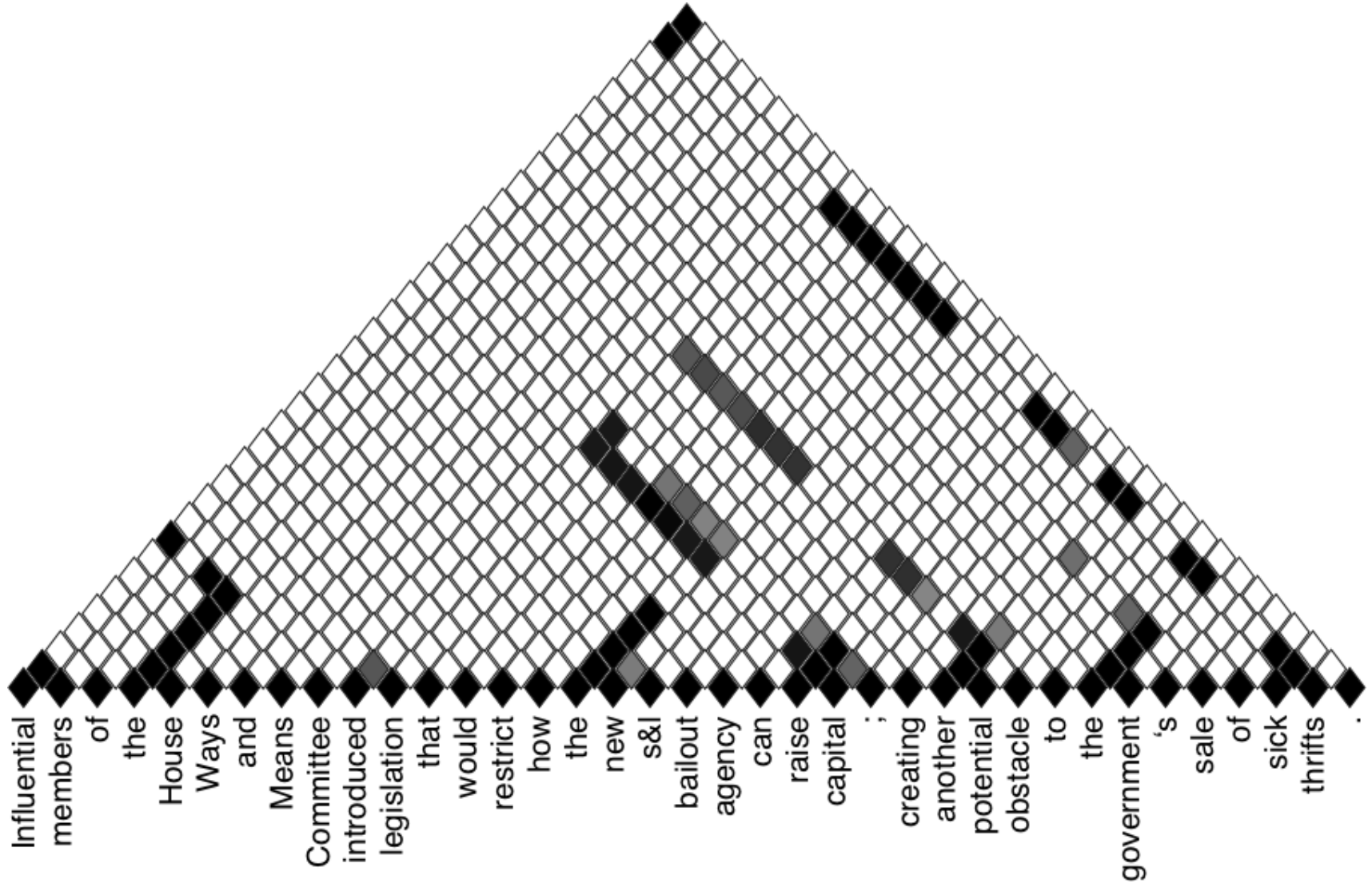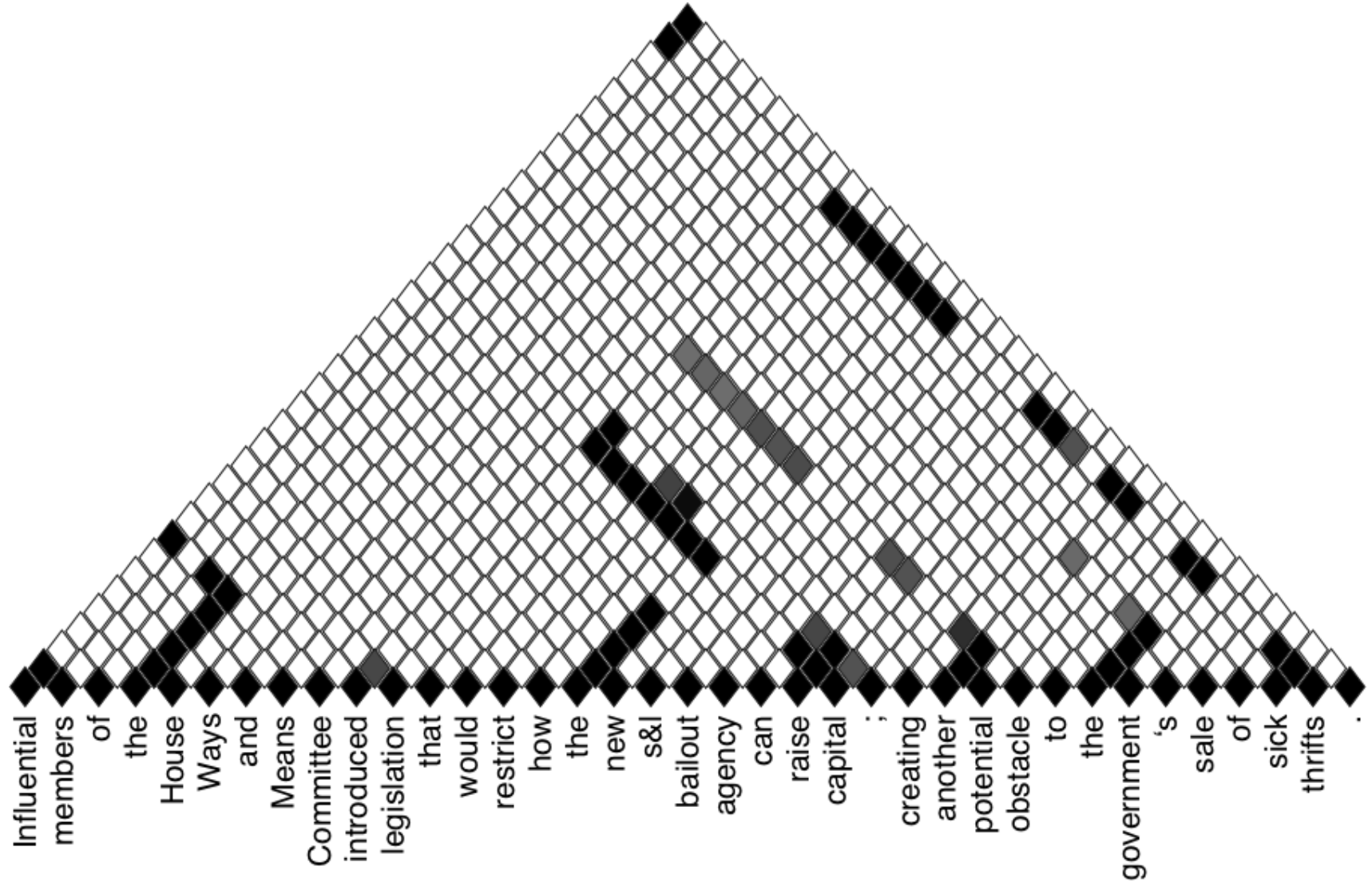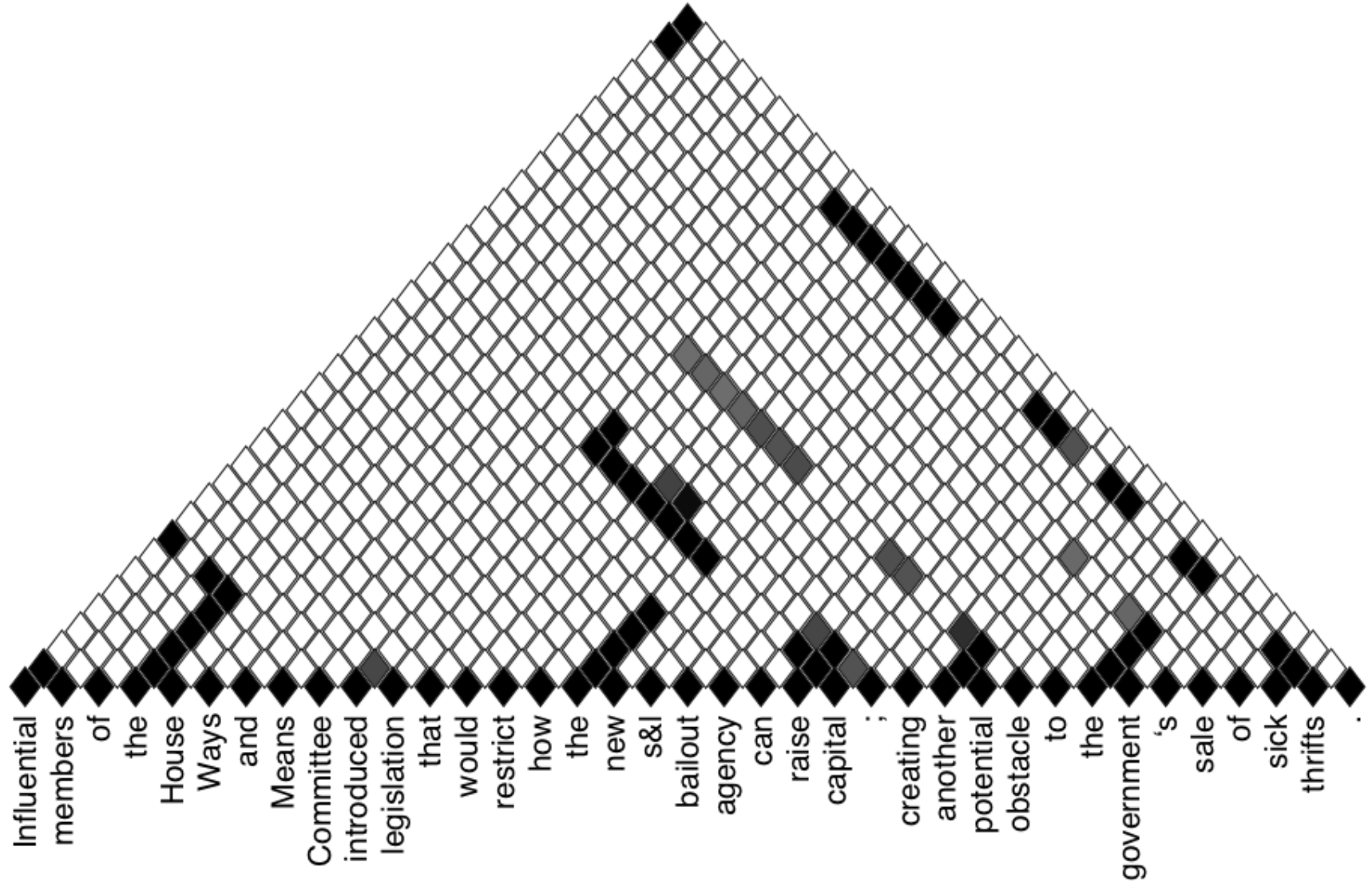# Bracket Posteriors

# Bracket Posteriors

# Bracket Posteriors



Influential members of the House Ways and Means Committee introduced legislation that would restrict how the new s&l bailout agency can raise capital , creating another potential obstacle to the government 's sale of sick thrifts .

# Bracket Posteriors



Influential members of the House Ways and Means Committee introduced legislation that would restrict how the new s&l bailout agency can raise capital ; creating another potential obstacle to the government 's sale of sick thrifts .