# Syntax and Parsing I

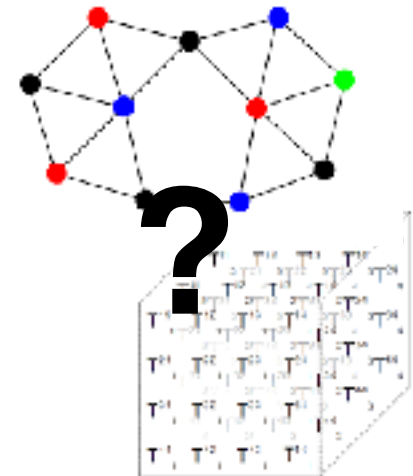## Constituency Parsing

## Slav Petrov – Google

Thanks to:

Dan Klein, Ryan McDonald, Alexander Rush, Joakim Nivre, Greg Durrett, David Weiss, Luheng He, Timothy Dozat

Lisbon Machine Learning School 2018

# Why Parsing?

On January 13, 2018, a false ballistic missile alert was issued via the Emergency Alert System and Commercial Mobile Alert System over television, radio, and cellphones in the U.S. state of Hawaii. The alert stated that there was an incoming ballistic missile threat to Hawaii, advised residents to seek shelter, and concluded "This is not a drill". The message was sent at 8:07 a.m. local time.
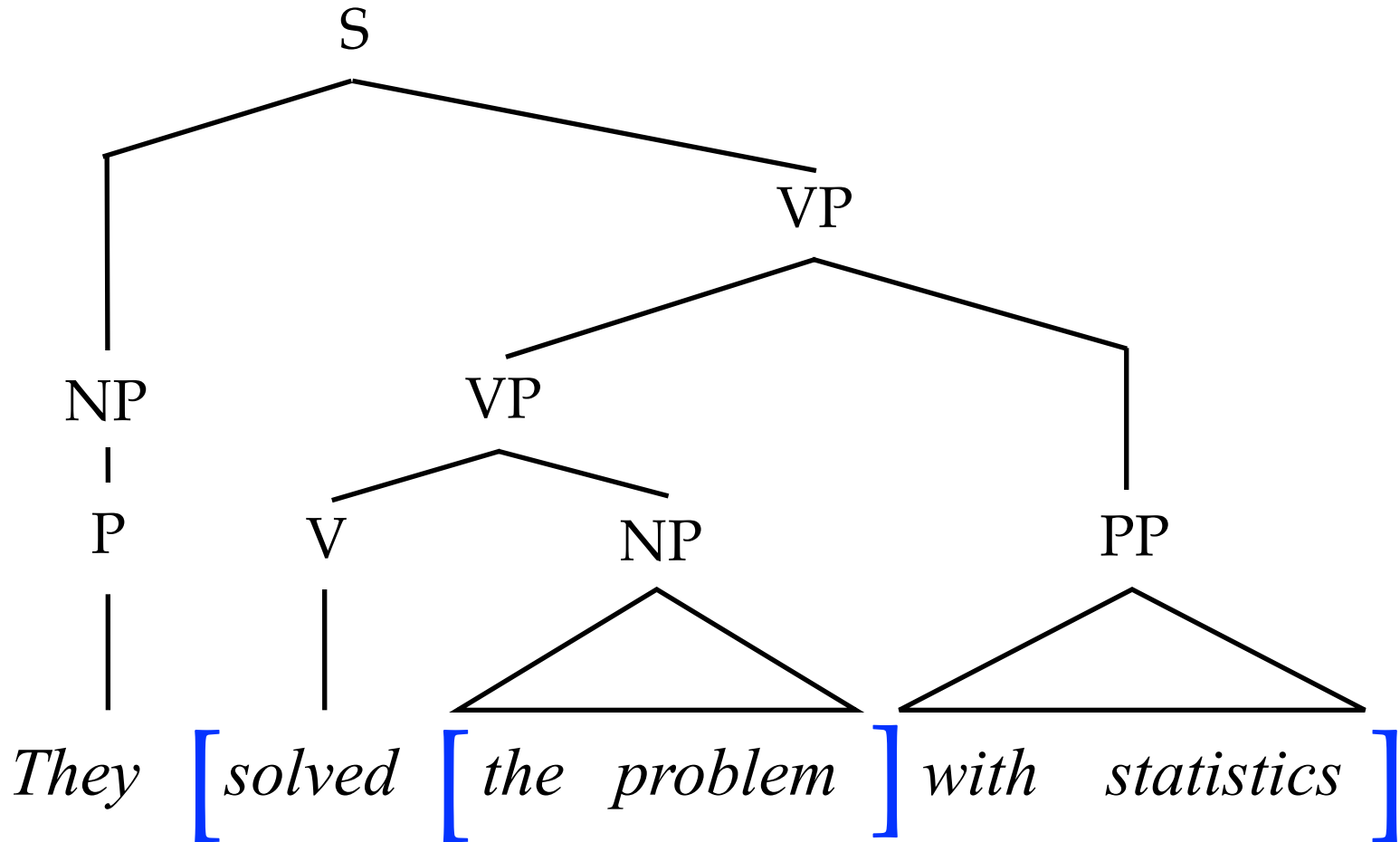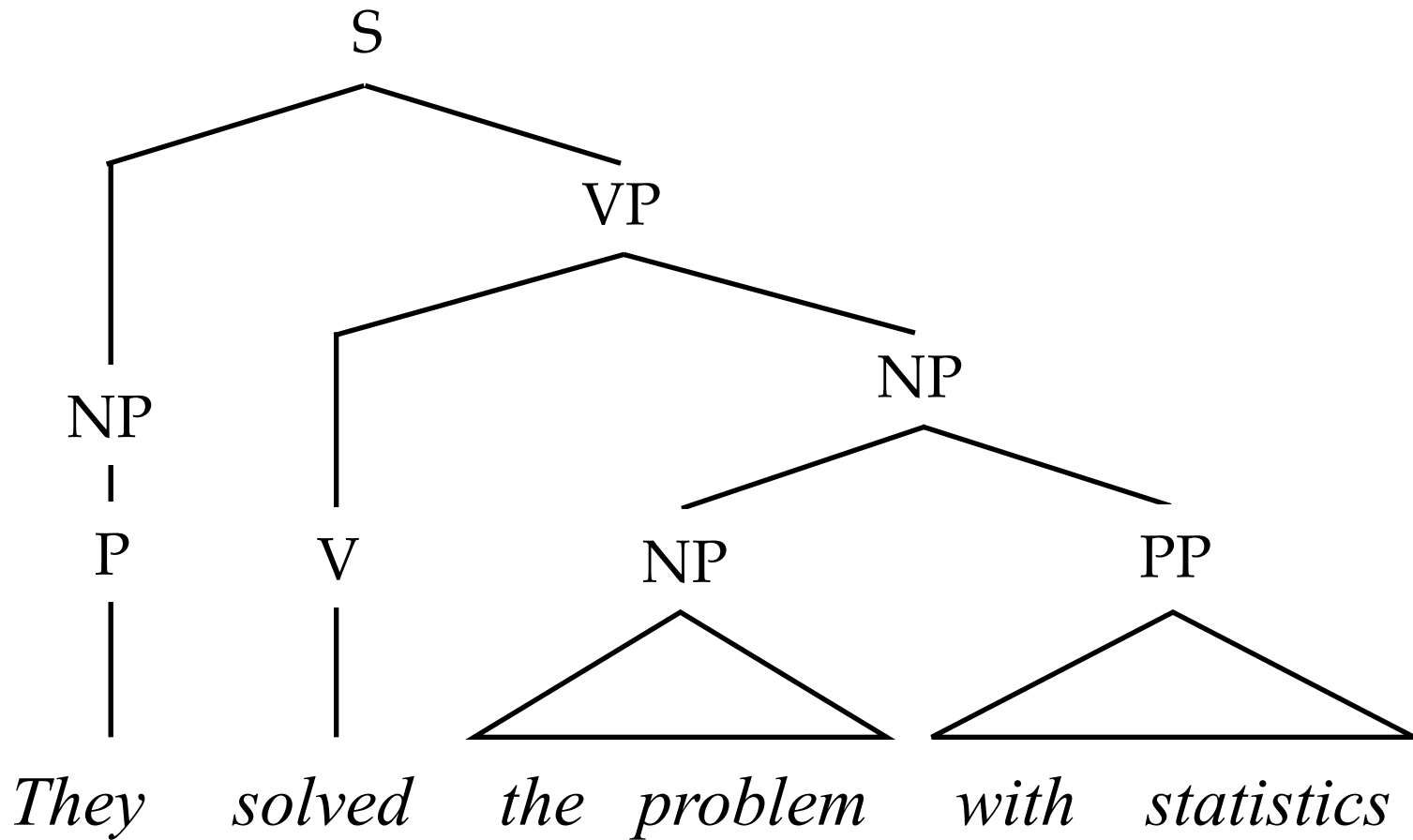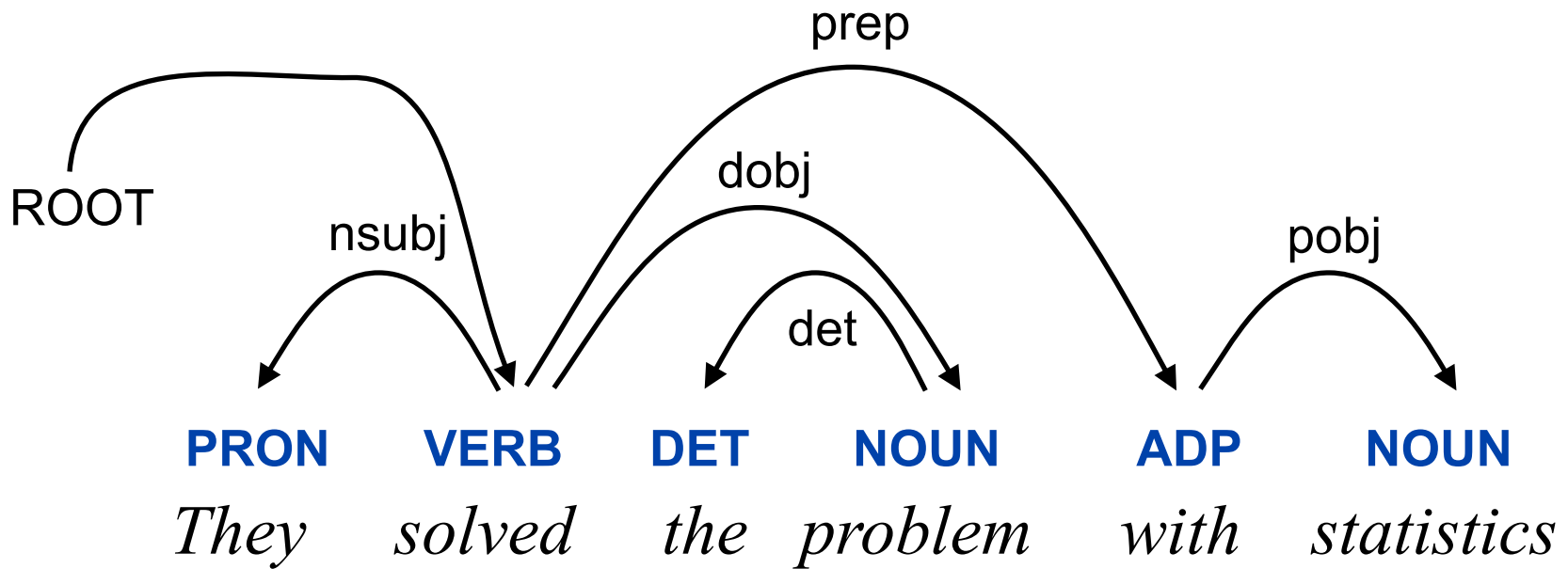
**Input Document**

**Structure:**
Syntactic,
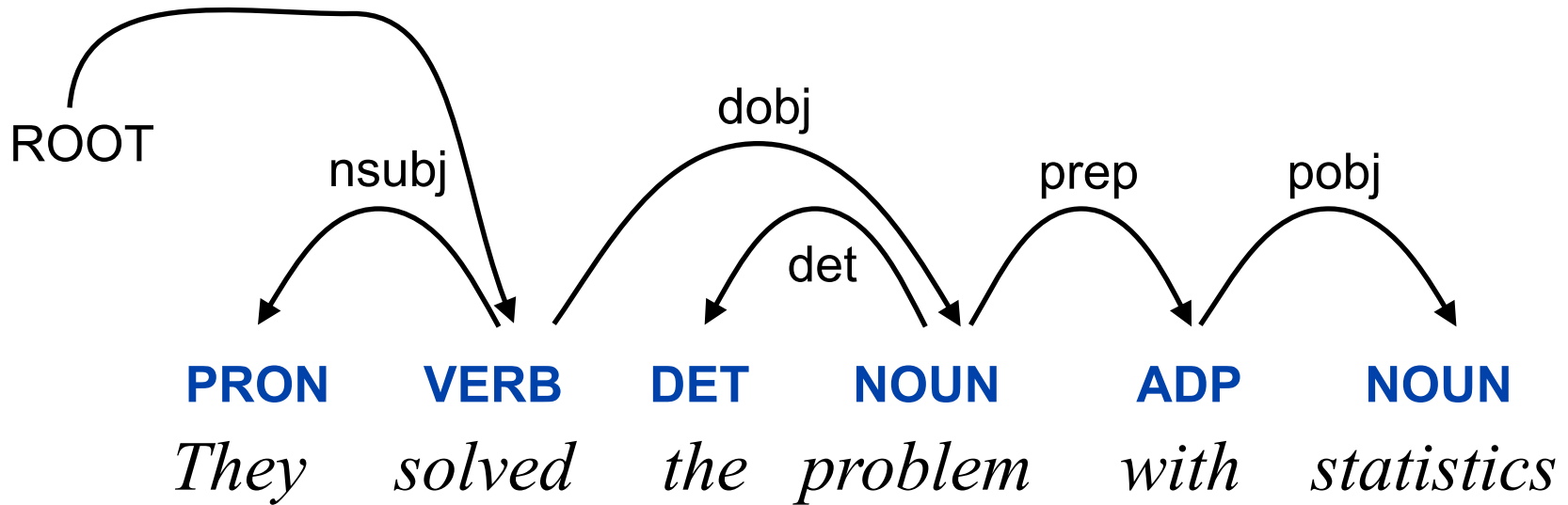Semantic,
etc.

# Analyzing Natural Language

# Syntax and Semantics



```
                    S
          ┌─────────┴─────────┐
          │                   VP
          │          ┌────────┴────────┐
          NP         │                 NP
          │          │          ┌──────┴──────┐
          P          V          NP            PP
          │          │         ╱│╲           ╱│╲
        They       solved    the problem   with statistics
```
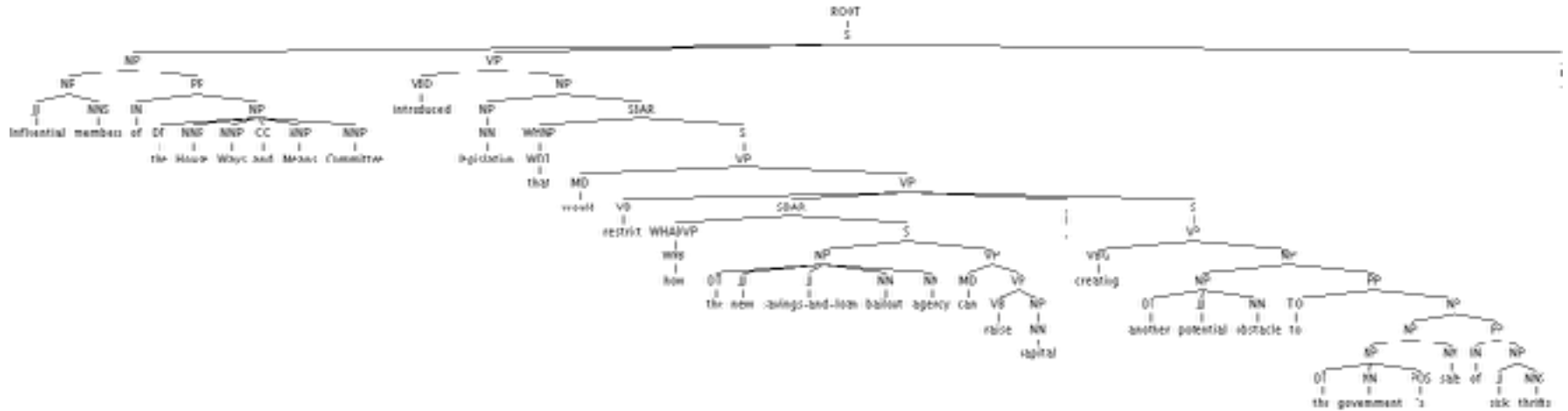
# Constituency and Dependency
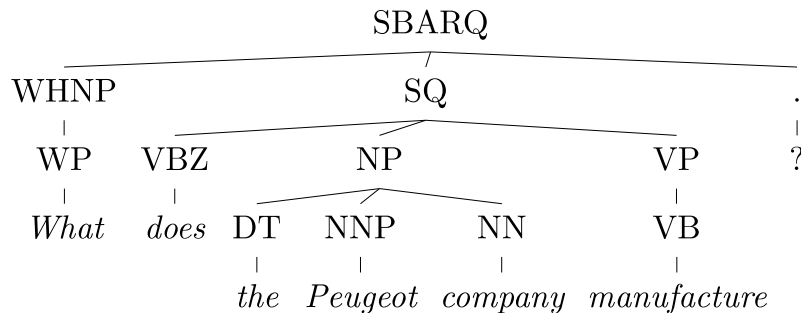
# Constituency and Dependency
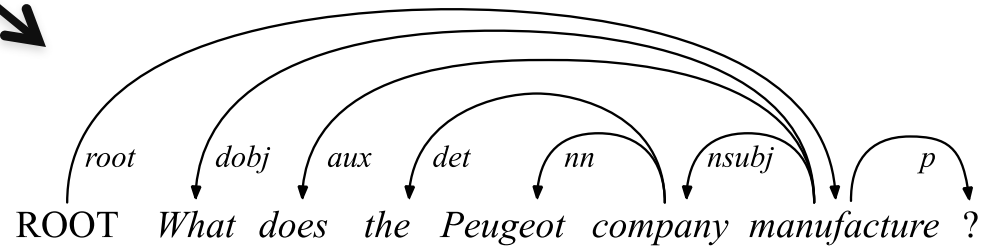
# A "real" Sentence



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

# Dependency Parsing

SBARQ

WHNP      SQ      .

WP   VBZ    NP    VP   ?

*What*   *does*   DT   NNP   NN   VB

*the*   *Peugeot*   *company*   *manufacture*

- Directed edges between pairs of word (head, dependent)

- Can handle free word-order languages

- Very efficient decoding algorithms exist

- Second part of today's lecture

*root*   *dobj*   *aux*   *det*   *nn*   *nsubj*   *p*

ROOT   *What*   *does*   *the*   *Peugeot*   *company*   *manufacture*   ?

# Attachments

- I cleaned the dishes from dinner

- I cleaned the dishes with detergent

- I cleaned the dishes in my pajamas

- I cleaned the dishes in the sink

# Classical NLP: Parsing

- Write symbolic or logical rules:

| | VBD | | VB | | | |
| | VBN | VBZ | VBP | VBZ | | |
| | NNP | NNS | NN | NNS | CD | NN |

Fed raises interest rates 0.5 percent

---

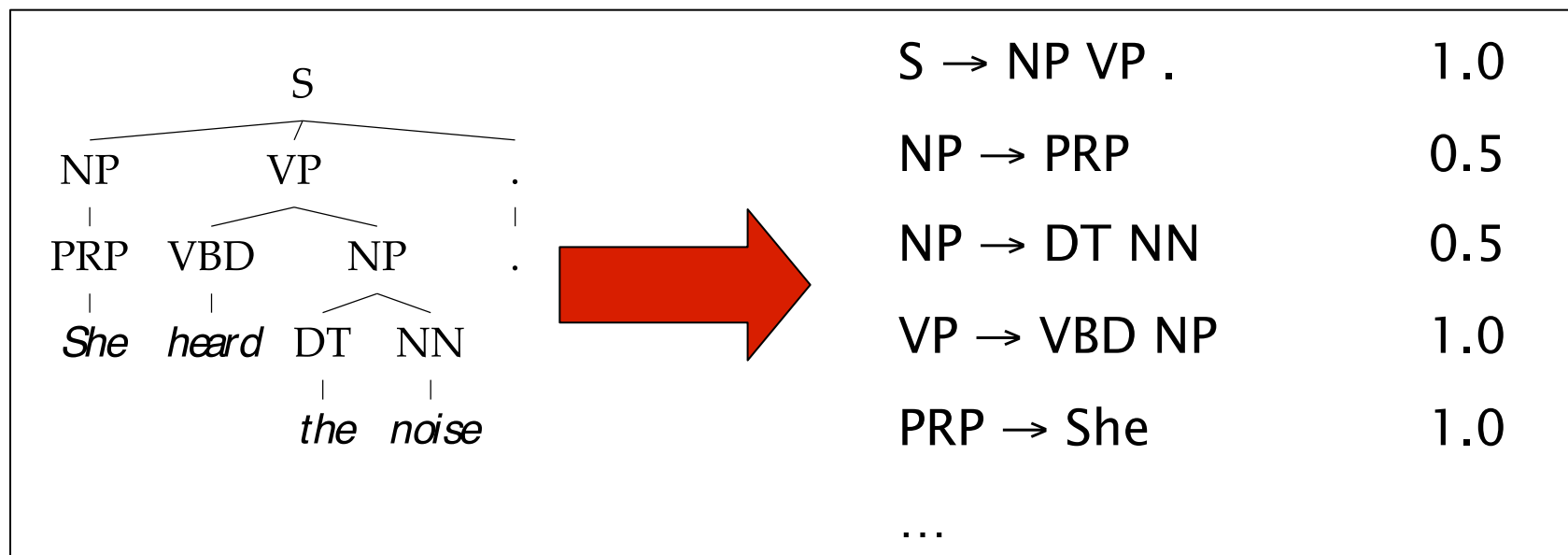| **Grammar (CFG)** | | **Lexicon** |
| ROOT → S | NP → NP PP | NN → interest |
| S → NP VP | VP → VBP NP | NNS → raises |
| NP → DT NN | VP → VBP NP PP | VBP → interest |
| NP → NN NNS | PP → IN NP | VBZ → raises |

- Use deduction systems to prove parses from words
  - Minimal grammar on "Fed raises" sentence: 36 parses
  - Real-size grammar: many millions of parses
- This scaled very badly, didn't yield broad-coverage tools

# Probabilistic Context-Free Grammars

- A context-free grammar is a tuple <N, T, S, R>
  - N : the set of non-terminals
    - Phrasal categories: S, NP, VP, ADJP, etc.
    - Parts-of-speech (pre-terminals): NN, JJ, DT, VB
  - T : the set of terminals (the words)
  - S : the start symbol
    - Often written as ROOT or TOP
    - Not usually the sentence non-terminal S
  - R : the set of rules
    - Of the form $X \rightarrow Y_1\ Y_2 \ldots Y_k$, with $X, Y_i \in N$
    - Examples: $S \rightarrow NP\ VP$, $VP \rightarrow VP\ CC\ VP$
    - Also called rewrites, productions, or local trees
- A PCFG adds:
  - A top-down production probability per rule $P(Y_1\ Y_2 \ldots Y_k \mid X)$

# Treebank Grammars

- Need a PCFG for broad coverage parsing.

- Can take a grammar right off the trees (doesn't work well):



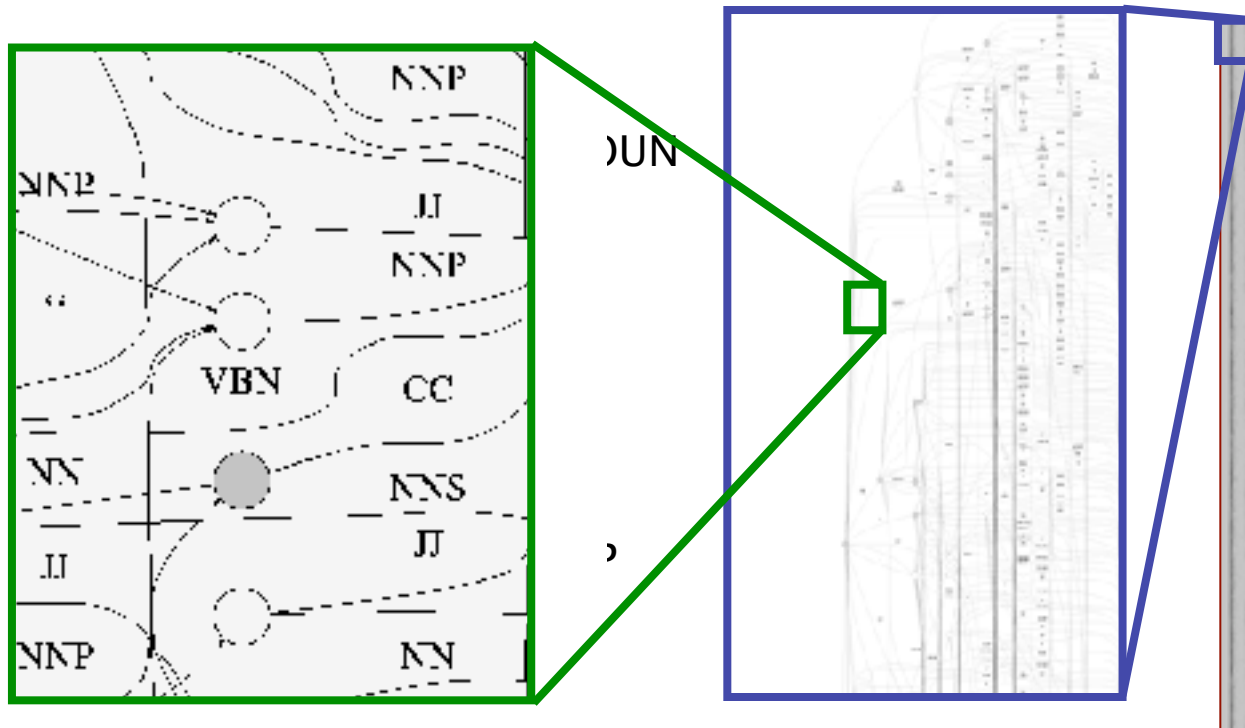| | |
|---|---|
| S → NP VP . | 1.0 |
| NP → PRP | 0.5 |
| NP → DT NN | 0.5 |
| VP → VBD NP | 1.0 |
| PRP → She | 1.0 |
| … | |

- Better results by enriching the grammar (e.g., lexicalization).

- Can also get reasonable parsers without lexicalization.

# Treebank Grammar Scale

- Treebank grammars can be enormous
  - As FSAs, the raw grammar has ~10K states, excluding the lexicon
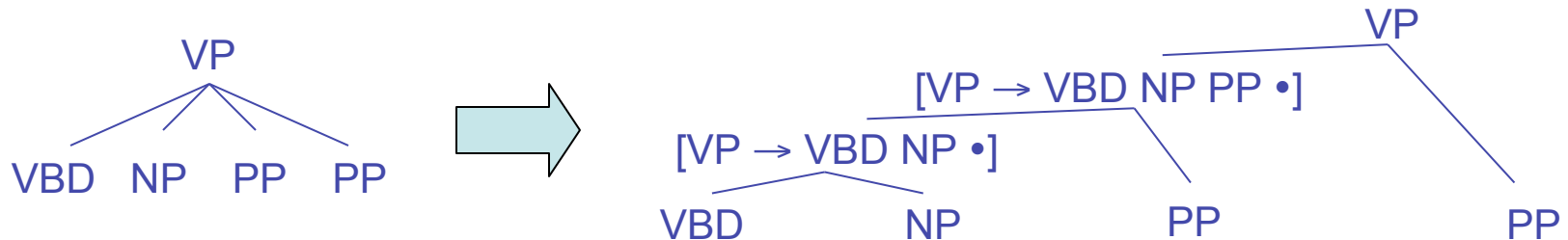  - Better parsers usually make the grammars larger, not smaller

NP

# Chomsky Normal Form

- ## Chomsky normal form:
  - All rules of the form X → Y Z or X → w
  - In principle, this is no limitation on the space of (P)CFGs
    - N-ary rules introduce new non-terminals



    - Unaries / empties are "promoted"
  - In practice it's kind of a pain:
    - Reconstructing n-aries is easy
    - Reconstructing unaries is trickier
    - The straightforward transformations don't preserve tree scores
  - Makes parsing algorithms simpler!

# A Recursive Parser

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

- Will this parser work?
- Why or why not?
- Memory requirements?

# A Memoized Parser

- One small change:

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

# A Bottom-Up Parser (CKY)

• Can also organize things bottom-up

```
bestScore(s)
  for (i : [0,n-1])
      for (X : tags[s[i]])
      score[X][i][i+1] =
          tagScore(X,s[i])
  for (diff : [2,n])
    for (i : [0,n-diff])
      j = i + diff
      for (X->YZ : rule)
        for (k : [i+1, j-1])
      score[X][i][j] = max score[X][i][j],
                          score(X->YZ) *
                          score[Y][i][k] *
                          score[Z][k][j]
```
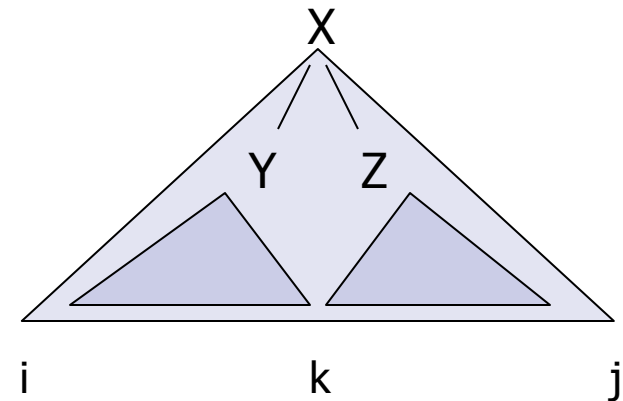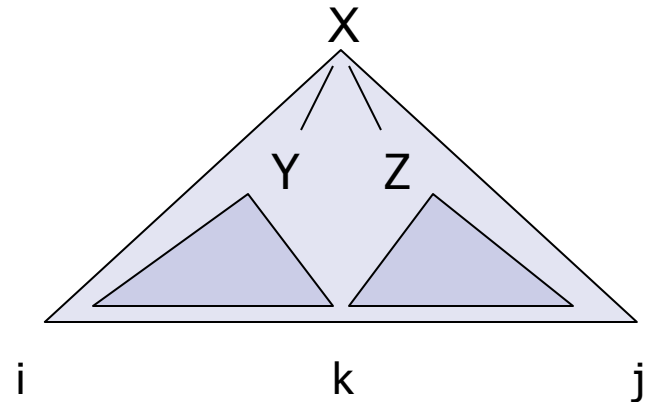
# Time: Theory

- ## How much time will it take to parse?

    - ### For each diff (<= n)
        - #### For each i (<= n)
            - ##### For each rule X → Y Z
                - ###### For each split point k
                    Do constant work

X

Y    Z

i          k          j

    - ## Total time: |rules|*n3

    - ## Something like 5 sec for an unoptimized parse of a 20-word sentences, or 0.2sec for an optimized parser
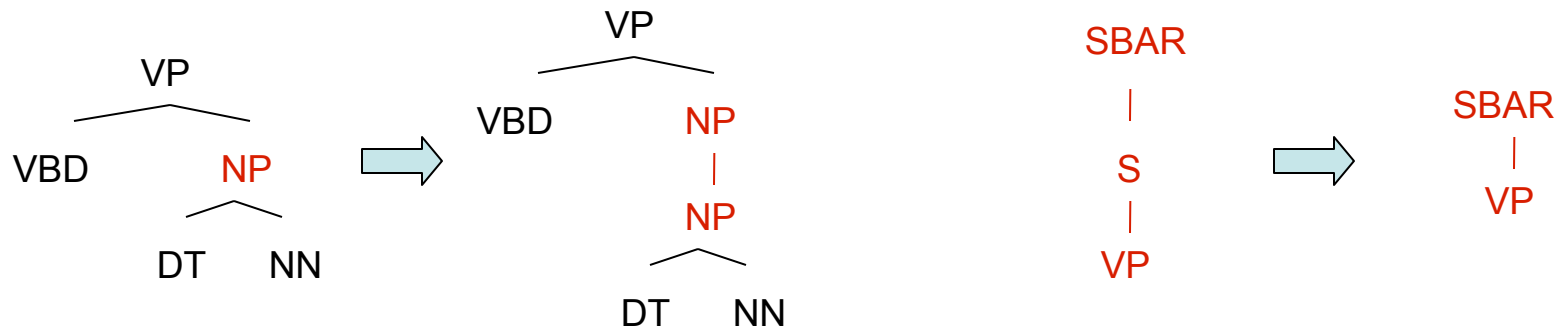
# Unary Rules

- Unary rules?

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

# CNF + Unary Closure

- ## We need unaries to be non-cyclic

  - Can address by pre-calculating the unary closure
  - Rather than having zero or more unaries, always have exactly one

VP
VBD    NP
    DT    NN

⟹

VP
VBD        NP
         NP
      DT    NN

SBAR
 |
 S
 |
 VP

⟹

SBAR
 |
 VP

  - Alternate unary and binary layers
  - Reconstruct unary chains afterwards

# Alternating Layers

```
bestScoreB(X,i,j,s)

     return max max score(X->YZ) *
     bestScoreU(Y,i,k) *
     bestScoreU(Z,k,j)



bestScoreU(X,i,j,s)

    if (j = i+1)
     return tagScore(X,s[i])
    else
     return max max score(X->Y) *
                    bestScoreB(Y,i,j)
```

# Treebank Grammars

- Need a PCFG for broad coverage parsing.

- Can take a grammar right off the trees (doesn't work well):

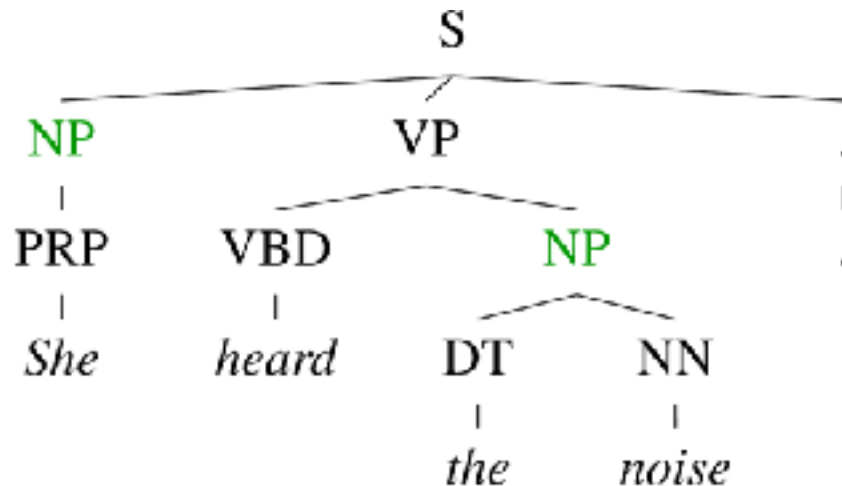| Tree | Rules | Prob |
|------|-------|------|
| S<br>NP VP .<br>PRP VBD NP .<br>She heard DT NN<br>the noise | S → NP VP . | 1.0 |
| | NP → PRP | 0.5 |
| | NP → DT NN | 0.5 |
| | VP → VBD NP | 1.0 |
| | PRP → She | 1.0 |
| | … | |

- Better results by enriching the grammar (e.g., lexicalization).
- Can also get reasonable parsers

| Model | F1 |
|-------|-----|
| Charniak '96 | 72.0 |

# Conditional Independence?

- Not every NP expansion can fill every NP slot



- A grammar with symbols like "NP" won't be context-free
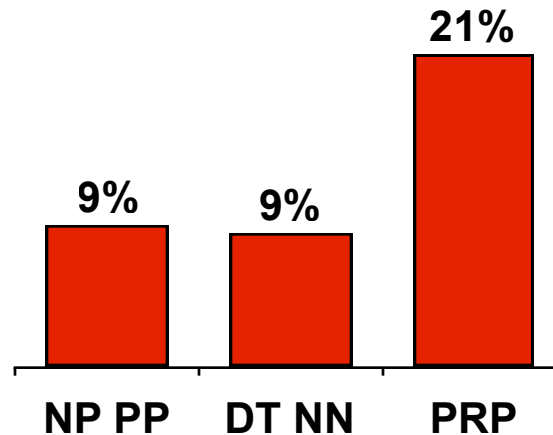
- Statistically, conditional independence too strong

# Non-Independence

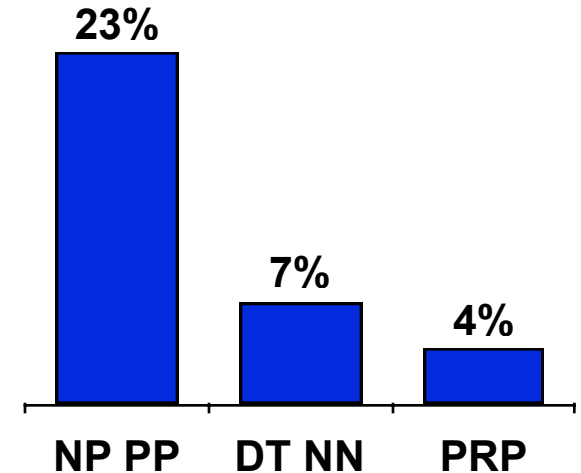- Independence assumptions are often too strong.

## All NPs

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

## NPs under S

| | | |
|---|---|---|
| 9% | 9% | 21% |
| NP PP | DT NN | PRP |

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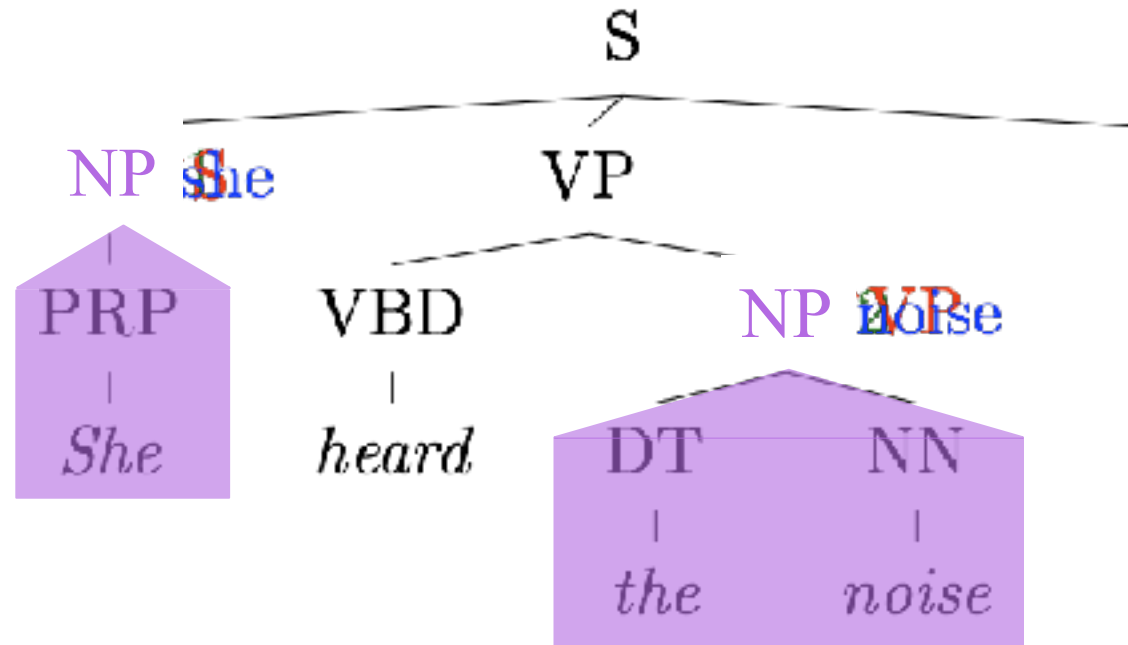
| | | |
|---|---|---|
| 23% | 7% | 4% |
| NP PP | DT NN | PRP |

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

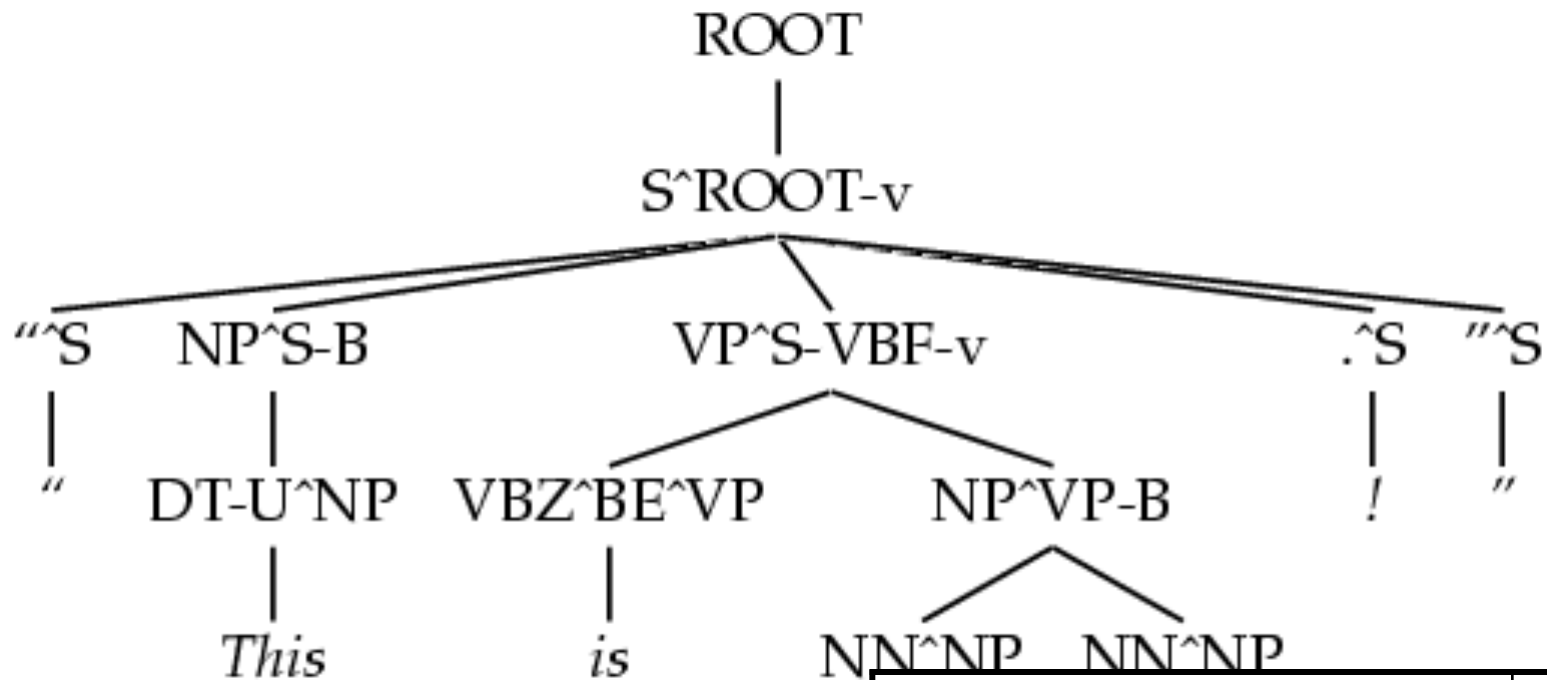- Also: the subject and object expansions are correlated!

# The Game of Designing a Grammar



- Structure Annotation [Johnson '98, Klein & Manning '03]
- Lexicalization [Collins '99, Charniak '00]
- Latent Variables [Matsuzaki et al. 05, Petrov et al. '06]
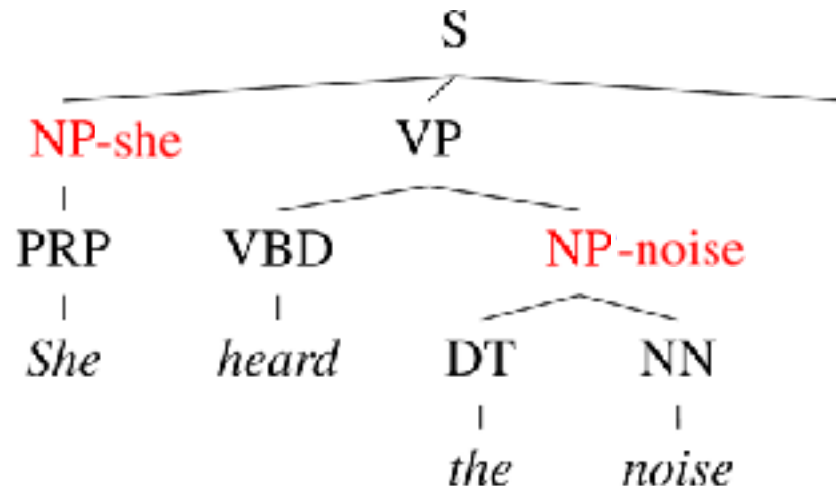- (Neural) CRF Parsing [Hall et al. '14, Durrett & Klein '15]

# A Fully Annotated (Unlexicalized) Tree
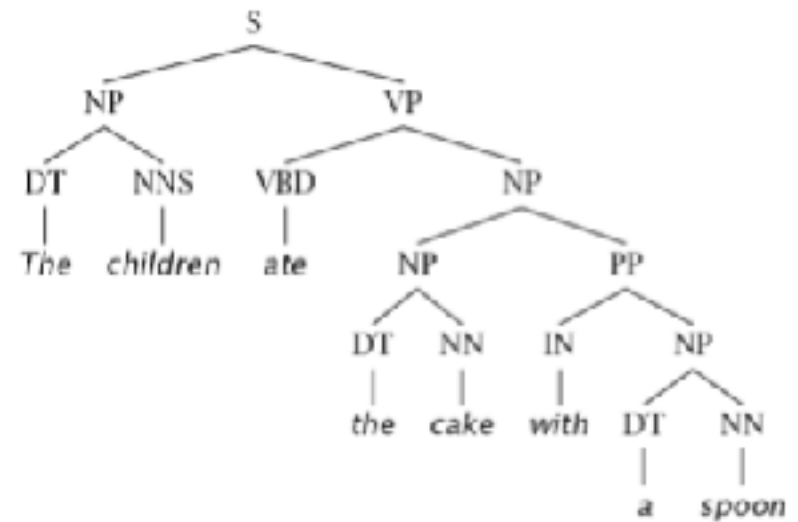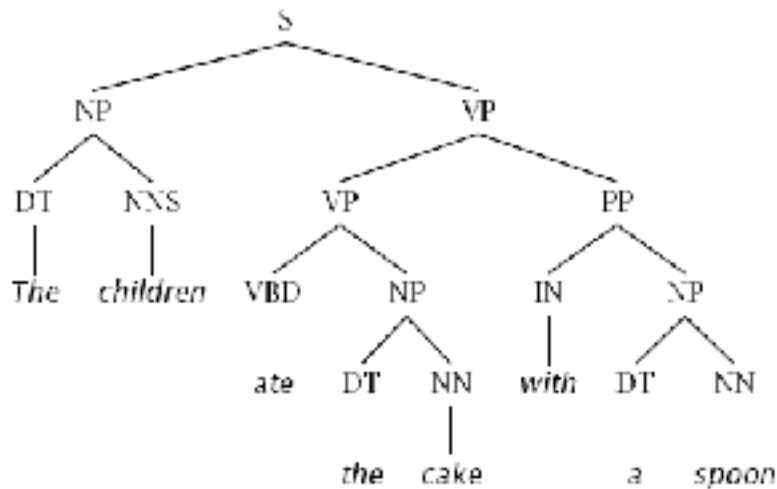
[Klein & Manning '03]



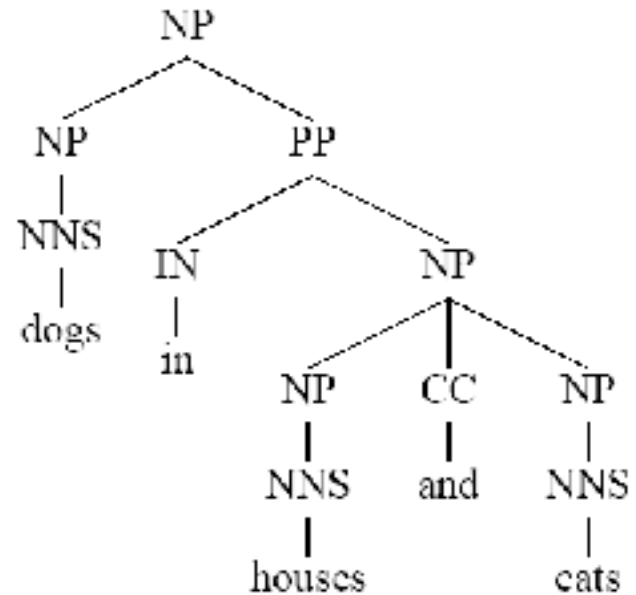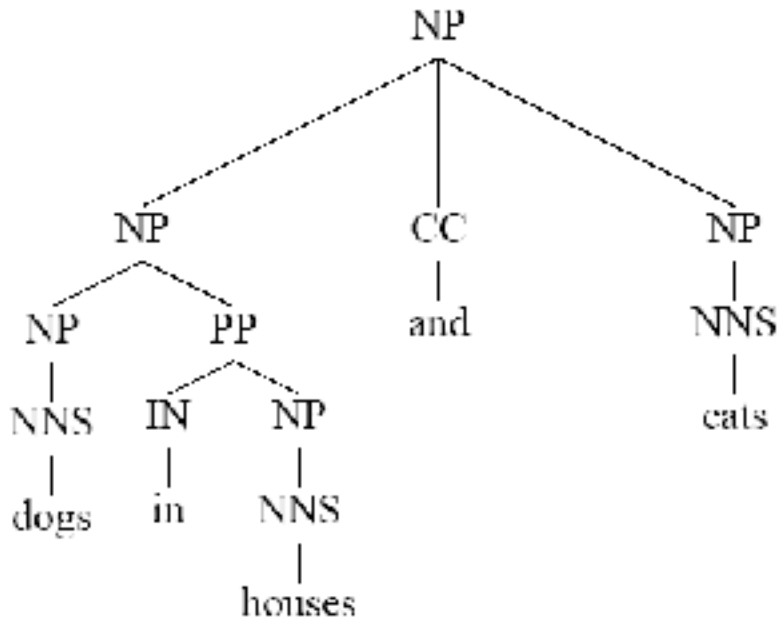| Model | F1 |
|---|---|
| Charniak '96 | 72.0 |
| Klein&Manning '03 | 86,3 |

# The Game of Designing a Grammar



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

# Problems with PCFGs



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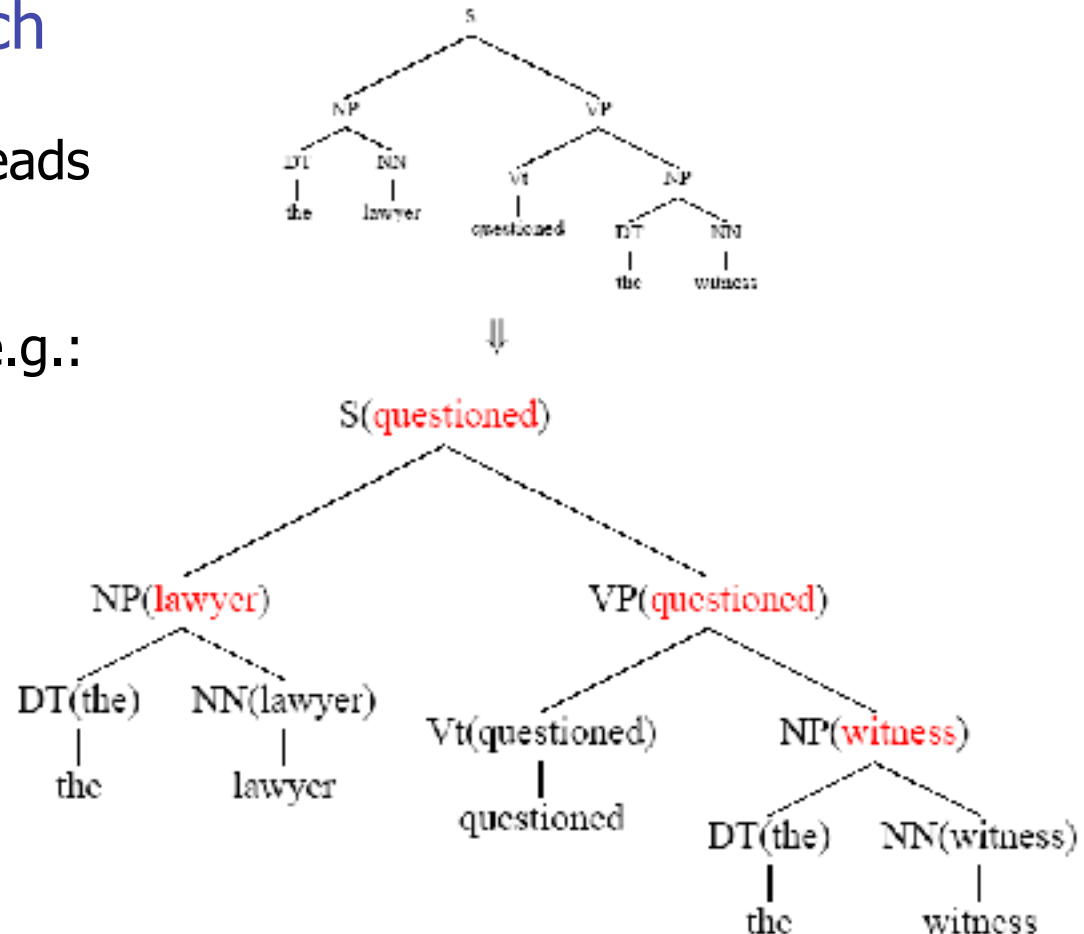- What (lexical) correlations need to be scored?

# Lexicalized Trees

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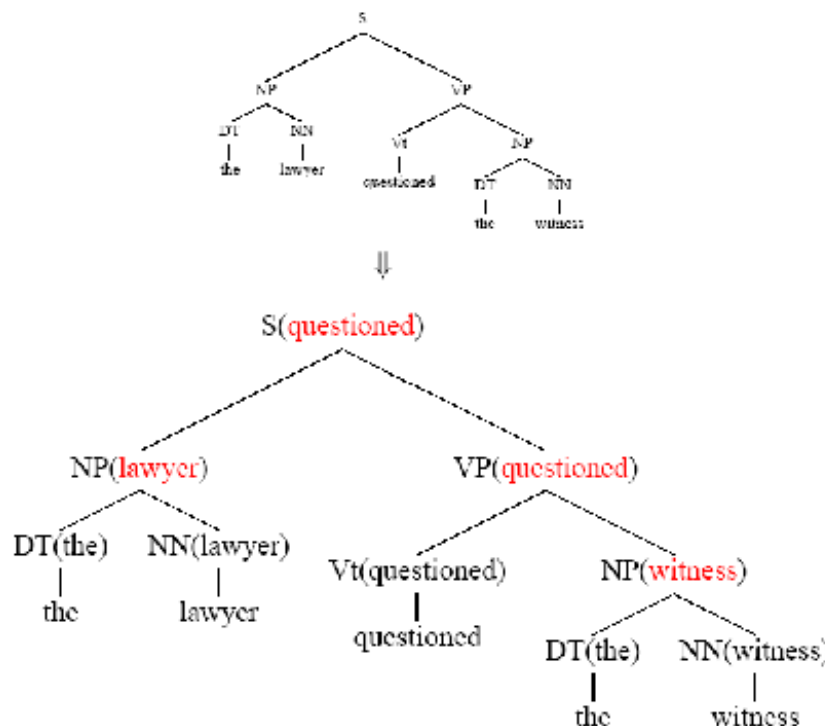    - NP:
      - Take leftmost NP
      - Take rightmost N*
      - Take rightmost JJ
      - Take right child
    - VP:
      - Take leftmost VB*
      - Take leftmost VP
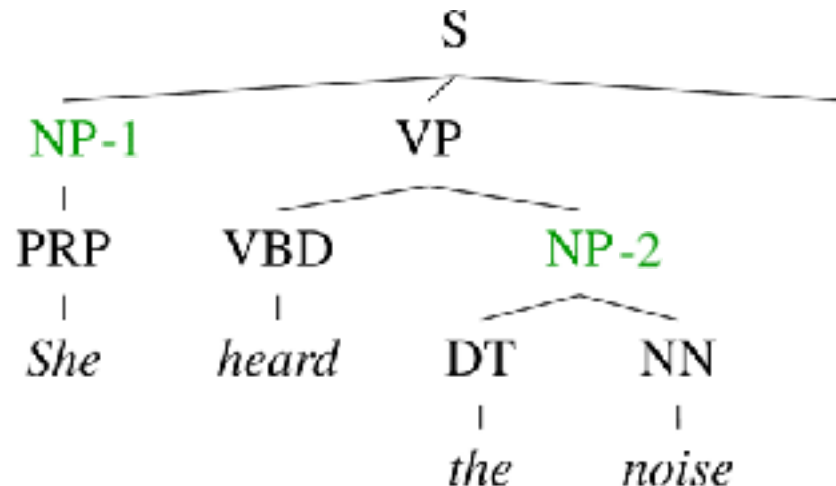      - Take left child

# Lexicalized Grammars

- Challenges:
  - Many parameters to estimate: requires sophisticated smoothing techniques
  - Exact inference is too slow: requires pruning heuristics
  - Difficult to adapt to new languages: At least head rules need to be specified, typically more changes needed

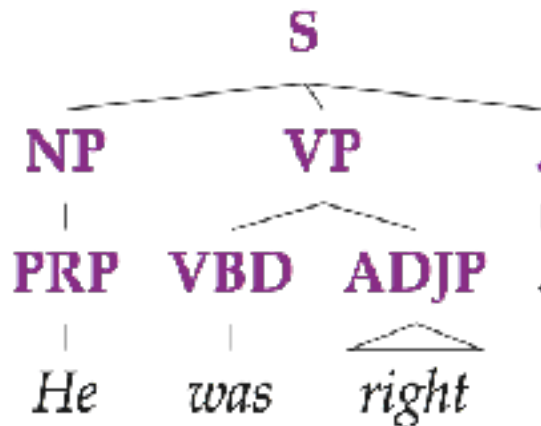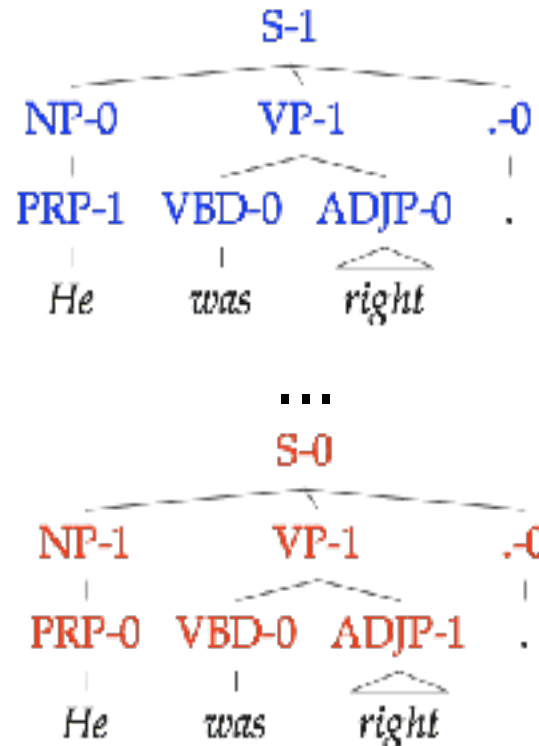| Model | F1 |
|---|---|
| Klein&Manning '03 | 86,3 |
| Charniak '00 | 89,7 |

# The Game of Designing a Grammar



- Annotation refines base treebank symbols to improve statistical fit of the grammar
  - Automatic clustering

# Latent Variable Grammars

[Matsuzaki et al. '05, Petrov et al. '06]



Parse Tree $T$
Sentence $w$

Derivations $t : T$

Parameters $\theta$

# 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-Backwa for HMMs.**

**Forward**



| Model | F1 |
|---|---|
| Charniak '00 | 89,7 |
| Petrov et al. '06 | 90,1 |

# The Game of Designing a Grammar



- Annotation refines base treebank symbols to improve statistical fit of the grammar
  - CRF Parsing (+Neural Network Representations)

# Generative vs. Discriminative

## Generative

Maximize joint likelihood
of gold tree **and** sentence

$w_1 w_2 \ldots w_n$ → EM-algorithm →

$W_1 W_2 \ldots W_n$

**EASY:** expectations over **observed** trees

[Matsuzaki et al. '05, Petrov et al. '06]

## Discriminative

Maximize conditional likelihood
of gold tree **given** sentence

$w_1 w_2 \ldots w_n$ → Gradient-based algorithm →

**HARD:** expectations over **all** trees

[Petrov & Klein '07, '08]

# Objective Functions

Generative Objective Function:

$$\max_\theta \ \mathcal{L}_\theta(\text{🌳}, w_1 \ldots w_n)$$

[Petrov, Barrett, Thibaux & Klein '06]

Discriminative Objective Function:

$$\max_\theta \ \mathcal{L}_\theta(\text{🌳} \,|\, w_1 \ldots w_n)$$

[Petrov & Klein '08, Finkel et. al '08]

Bayesian Objective Function:

$$\max_\theta \ \mathcal{P}(\theta \,|\, \text{🌳}) \mathcal{L}_\theta(\text{🌳}, w_1 \ldots w_n)$$

[Liang, Petrov, Jordan & Klein '07]

# (Neural) CRF Parsing

[Taskar et al. '04. Petrov & Klein '07, Hall et al. '14, Durrett et al. '15]



Score of VP over this span

Be a tree

$w \cdot f_s$

dense neural network

$w \cdot f_s$

sparse log-linear model

He    gave    a    speech

# CRF Parsing Sparse Features

$$P(T|x) \propto \prod_{r \in T} \exp\left(\text{score}(r)\right)$$

$$\text{score}(_2\text{NP}_7 \to {_2}\text{NP}_4\ {_4}\text{PP}_7) = w^\top f(_2\text{NP}_7 \to {_2}\text{NP}_4\ {_4}\text{PP}_7)$$

FirstWord = a  & NP → NP PP

PrevWord = gave  & NP → NP PP

AfterSplit = on  & NP → NP PP

FirstWord = a  & NP

…

# Neural CRF Model

$$\text{score}(_2\text{NP}_7 \rightarrow {}_2\text{NP}_4\ {}_4\text{PP}_7) =$$

$$W \odot \left( f_s \left( {}_2\text{X}_7 \rightarrow {}_2\text{X}_4\ {}_4\text{X}_7 \right) f_o^\top (\text{NP} \rightarrow \text{NP PP}) \right)$$

$f_s$ •••••••••••••••••••••••••••

$f_s$ ○●○●○○○○●   Sparse

$f_s = g(Hv)$    (arbitrary neural network)

$v$ •••• •••• •••• •••• •••• ••••

He  gave a  speech  on  fo

| Model | F1 |
|---|---|
| Petrov et al. '06 | 90,1 |
| Durrett et. al. '15 | 91,1 |

0  1      2  3        4

# LSTM Parsing [Vinyals et al. '15]

- Treat parsing as a sequence-to-sequence prediction problem
- Completely ignores tree structure, uses LSTMs as black boxes
- No global normalization, only local normalization



$$P(y_1, \ldots, y_{T'} | x_1, \ldots, x_T) = \prod$$

| Model | F1 |
|---|---|
| Durrett et. al. '15 | 91,1 |
| Vinyals et al. '15 | 88.6* |

# Parsing with Self-Attention

[Kitaev & Kleint '18]

Output   …(VP(VBD fled)(NP(DT the)(NN market))…

Decoder

Encoder

Input   …   and   fled   the   market   in   …
             CC    VBD    DT    NN       IN

| Model | F1 |
|---|---|
| Durrett et. al. '15 | 91,1 |
| Kitaev & Klein '18 | 93,6 |

8 layers

LayerNorm

Feed Forward

LayerNorm

Multi-Head Attention

wor

# Detailed English Results (Old)



| | Single Parser | Self-Trained | Reranker | Product | Combination |
|---|---|---|---|---|---|

**Single Parser**
- 89,7 [Charniak '00]
- 90,1 [Petrov et al. '06]
- 91,0 [Carreras et al. '08]
- 89,2 [Hall '12]
- 91,1 [Durrett et al. '15]
- 91,3 [Zhu et al. '13]
- **93,6** [Dyer et al. '16]

**Self-Trained**
- 91,4 [Huang & Harper '08]
- 91,6 [Huang & Harper, Petrov '10]

**Reranker**
- 91,5 [Charniak & Johnson '05]
- 91,7 [Huang '08]
- 92,3 [McClosky et al. '06]

**Product**
- 91,8 [Petrov '10]
- 92,4 [Huang & Harper, Petrov '10]

**Combination**
- 92 [Sagae & Lavie '06]
- 92,4 [Fossum & Knight '09]
- 92,6 [Zhang et al. '09]
- 92,8 [Vinyals et al. '16]

# Detailed English Results (New)



| | |
|---|---|
| Single Parser | Multi-Modal / Additional Data |

- 89,7 [Charniak '00]
- 90,1 [Petrov et al. '06]
- 91,0 [Carreras et al. '08]
- 88,3 [Vinyals et al. '16]
- 91,3 [Cross & Huang '16]
- 92,1 [Gaddy et al. '18]
- 92,6 [Stern et al. '17]
- 93,6 [Kitaev & Klein '18]
- 91,8 [Petrov '10]
- 92,4 [Huang et al. '10]
- 91,1 [Durrett & Klein '15]
- 92,8 [Vinyals et al. '15]
- 93,3 [Dyer et al. '16]
- 93,8 [Choe & Charniak '16]
- 94,2 [Liu & Chang '17]
- 94,7 [Fried et al. '17]
- 95,1 [Kitaev & Klein '18]

# Multi-Lingual Results