# CS388: Natural Language Processing Lecture 10: Syntax I
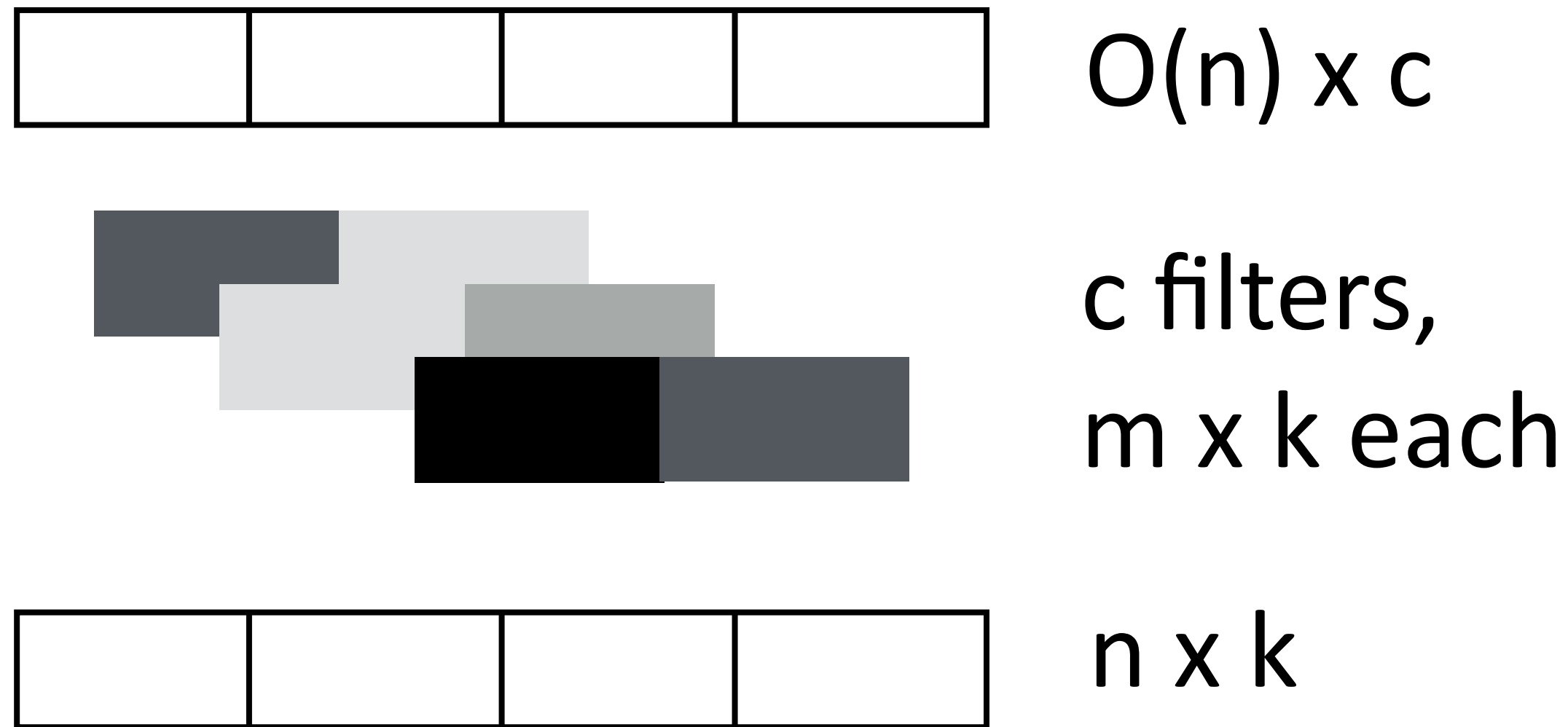


## Greg Durrett
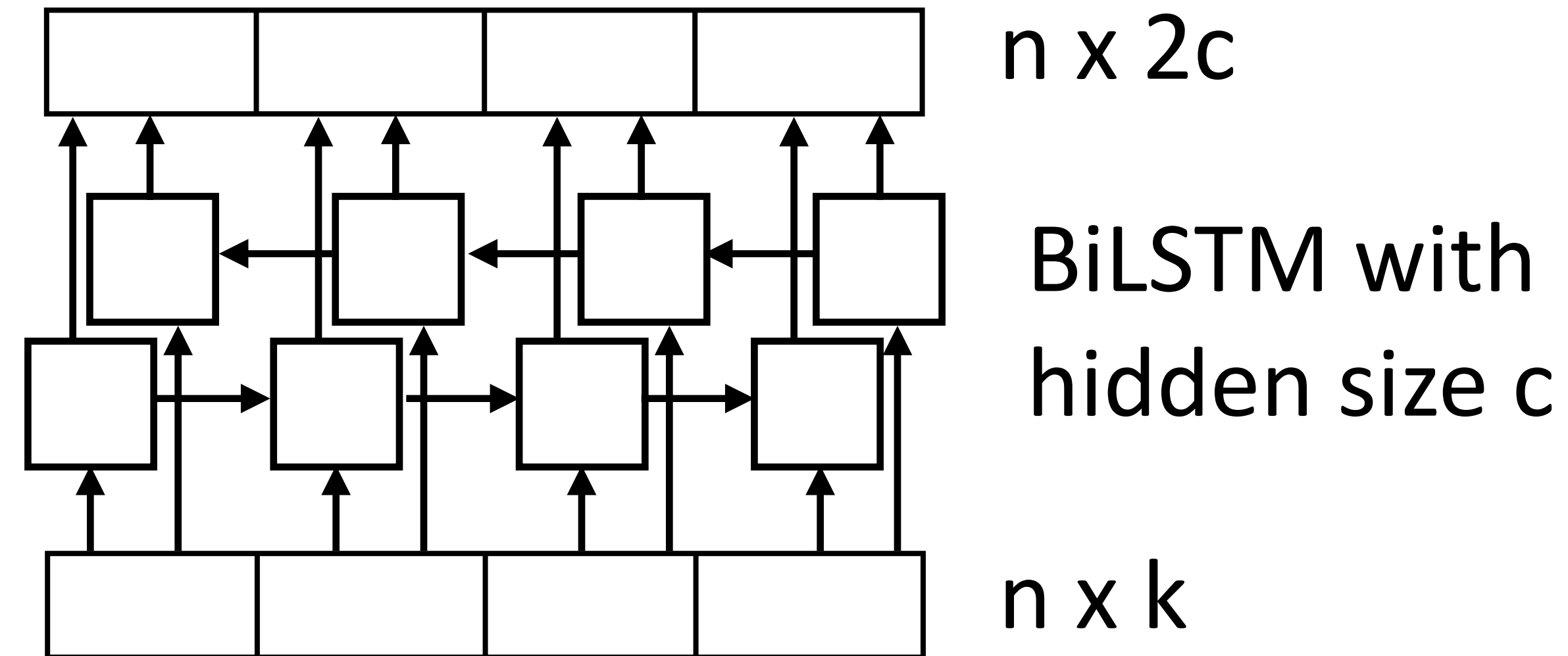
Slides adapted from Dan Klein, UC Berkeley

# Recall: CNNs vs. LSTMs

O(n) x c

c filters,
m x k each

n x k

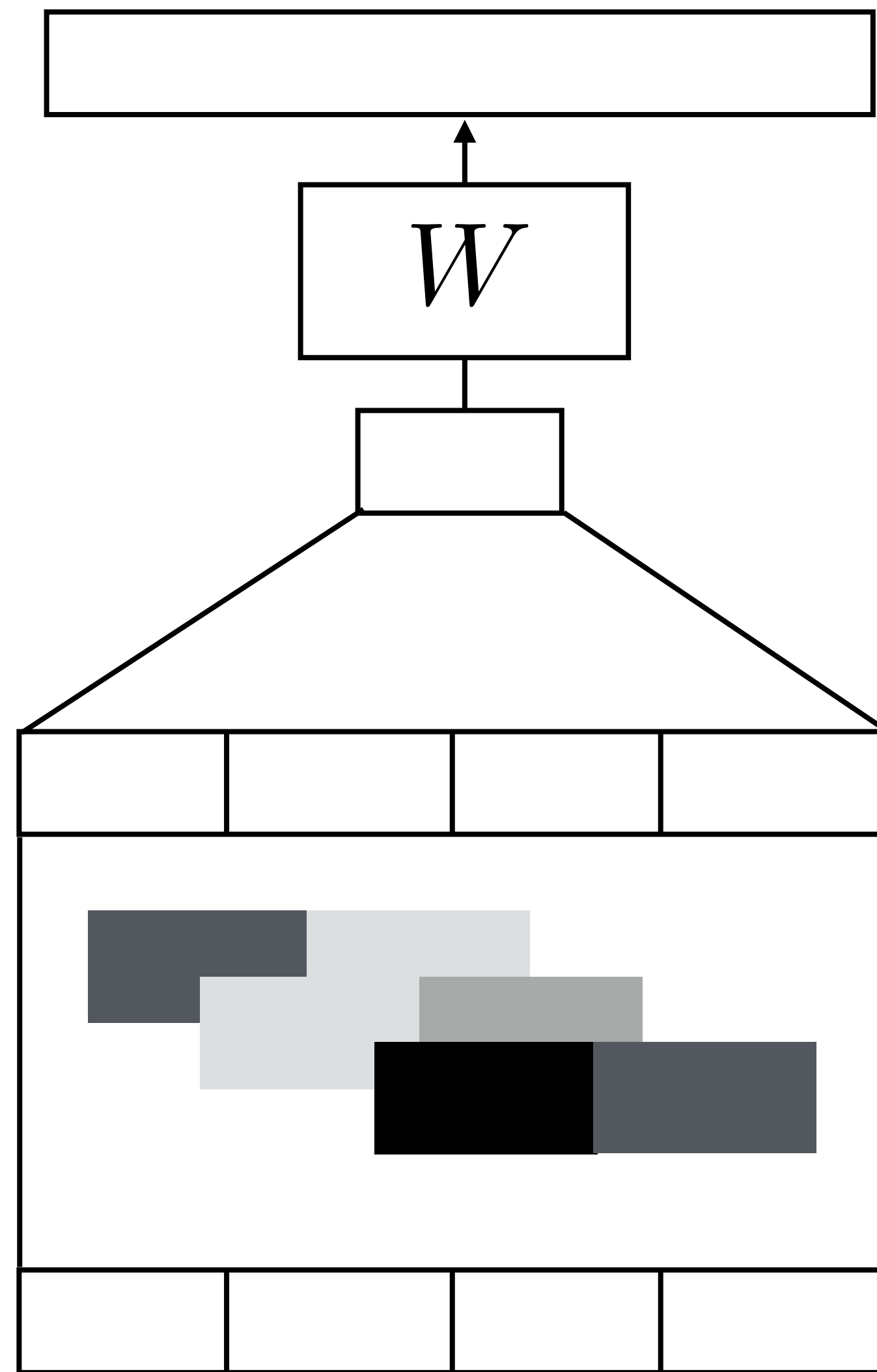the movie was good

n x 2c

BiLSTM with
hidden size c

n x k

the movie was good

▸ Both LSTMs and convolutional layers transform the input using context

▸ LSTM: "globally" looks at the entire sentence (but local for many problems)

▸ CNN: local depending on filter width + number of layers

# Recall: CNNs

$P(y|\mathbf{x})$

$W$

projection + softmax

c-dimensional vector

max pooling over the sentence

n x c

c filters,
m x k each

n x k

the movie was good

▸ Max pooling: return the max activation of a given filter over the entire sentence; like a logical OR (sum pooling is like logical AND)

# Recall: Neural CRFs

B-PER   I-PER        O      O     O     B-LOC        O      O    O B-ORG     O      O
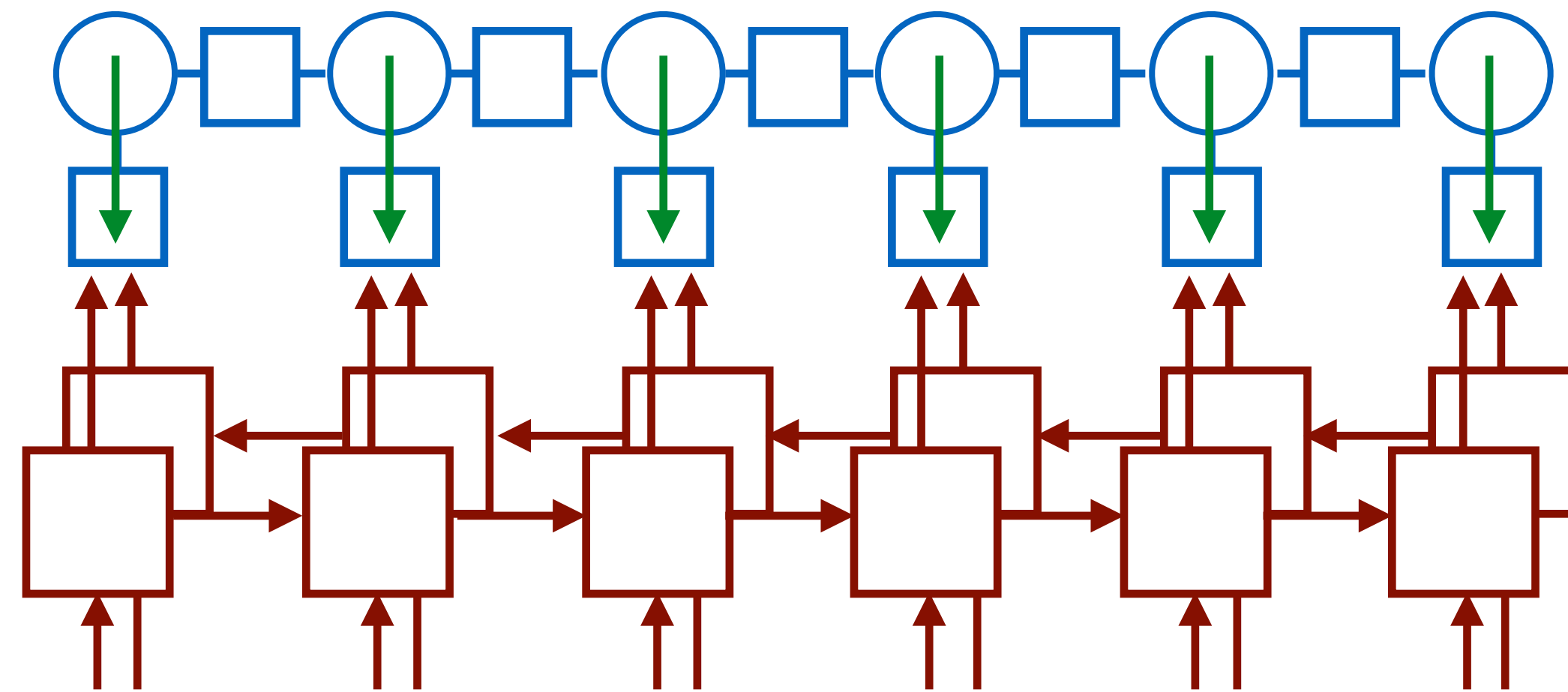
*Barack Obama will travel to Hangzhou today for the G20 meeting .*

PERSON                                          LOC                            ORG

Barack Obama will travel   to Hangzhou

2) Run forward-backward

3) Compute error signal

1) Compute f(**x**)

4) Backprop (no knowledge of sequential structure required)

# This Lecture

▸ Constituency formalism

▸ Context-free grammars and the CKY algorithm

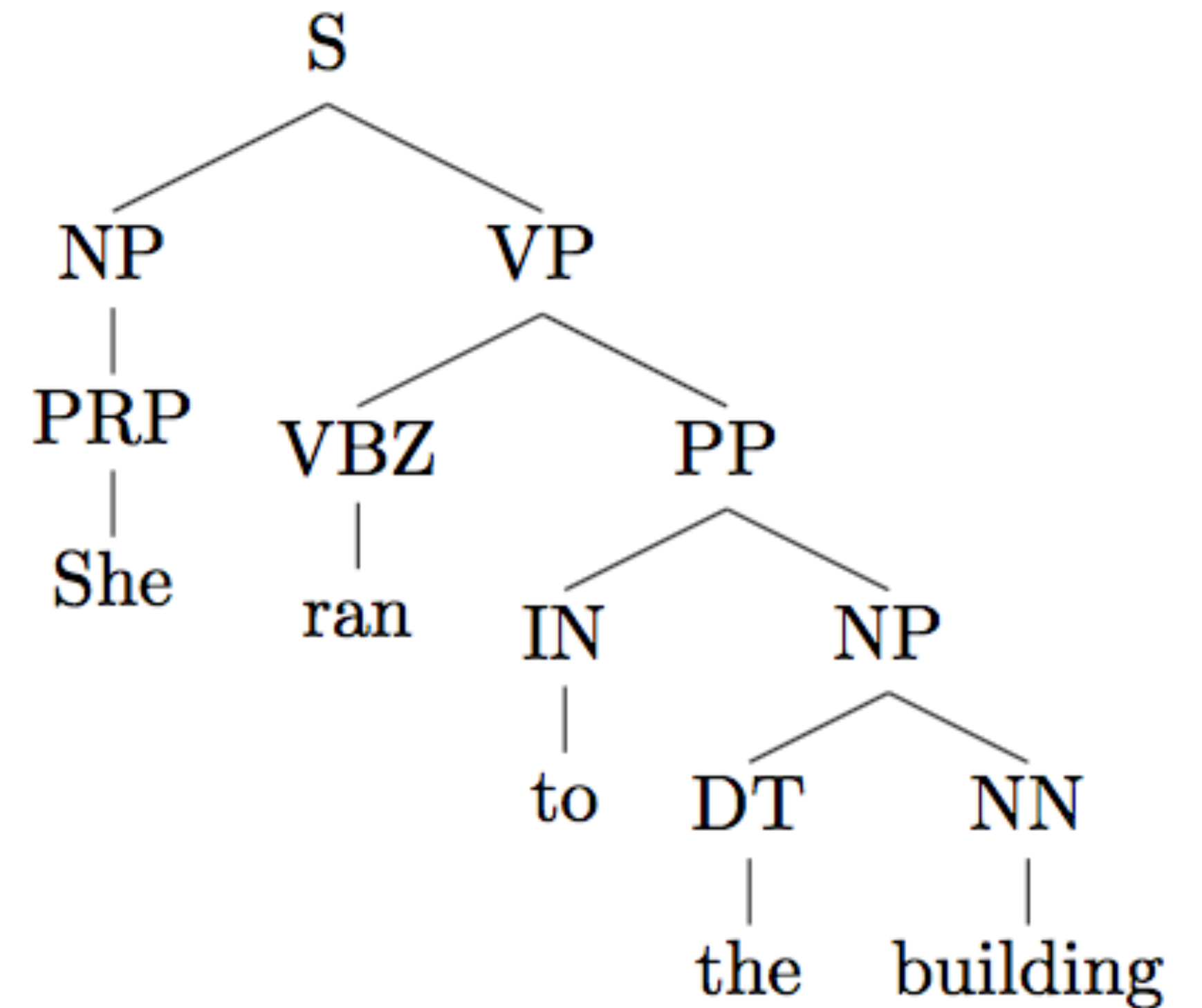▸ Refining grammars

▸ Discriminative parsers

# Constituency

# Syntax
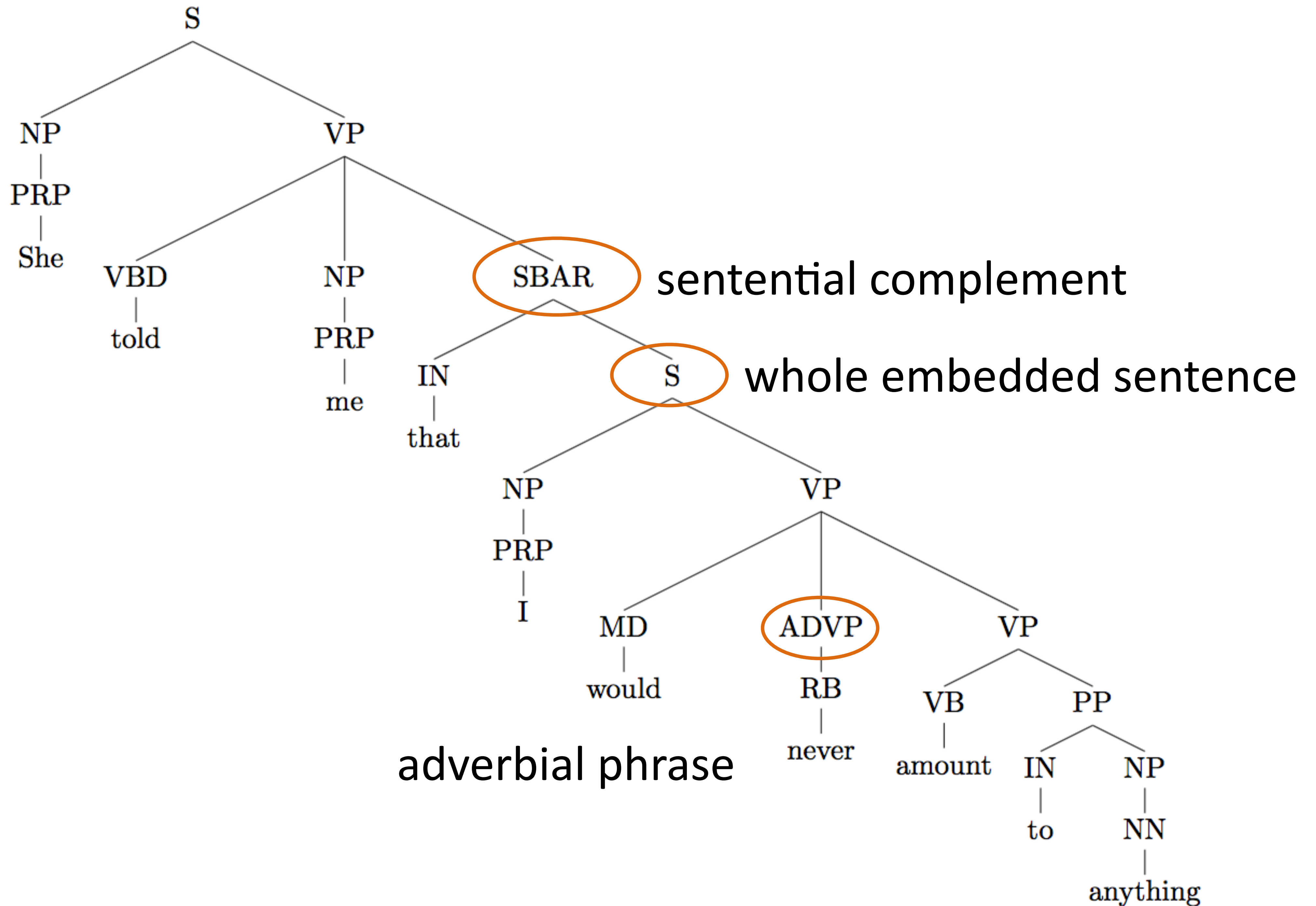
▸ Study of word order and how words form sentences

▸ Why do we care about syntax?

    ▸ Multiple interpretations of words (noun or verb? *Fed raises*... example)

    ▸ Recognize verb-argument structures (who is doing what to whom?)

    ▸ Higher level of abstraction beyond words: some languages are SVO, some are VSO, some are SOV, parsing can canonicalize

# Constituency Parsing

▸ Tree-structured syntactic analyses of sentences

▸ Common things: noun phrases,
verb phrases, prepositional phrases

▸ Bottom layer is POS tags

▸ Examples will be in English. Constituency
makes sense for a lot of languages but
not all

sentential complement

whole embedded sentence

adverbial phrase
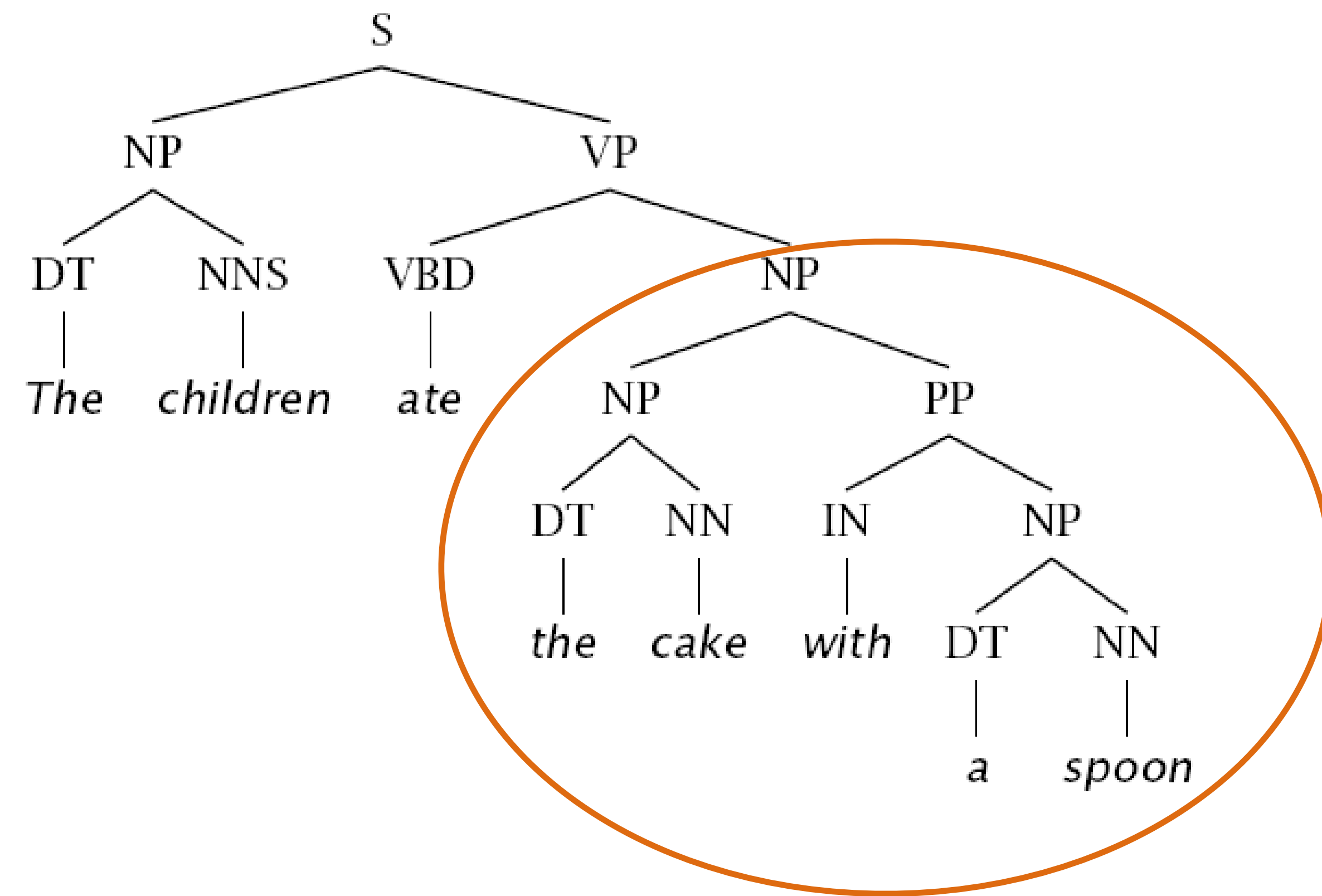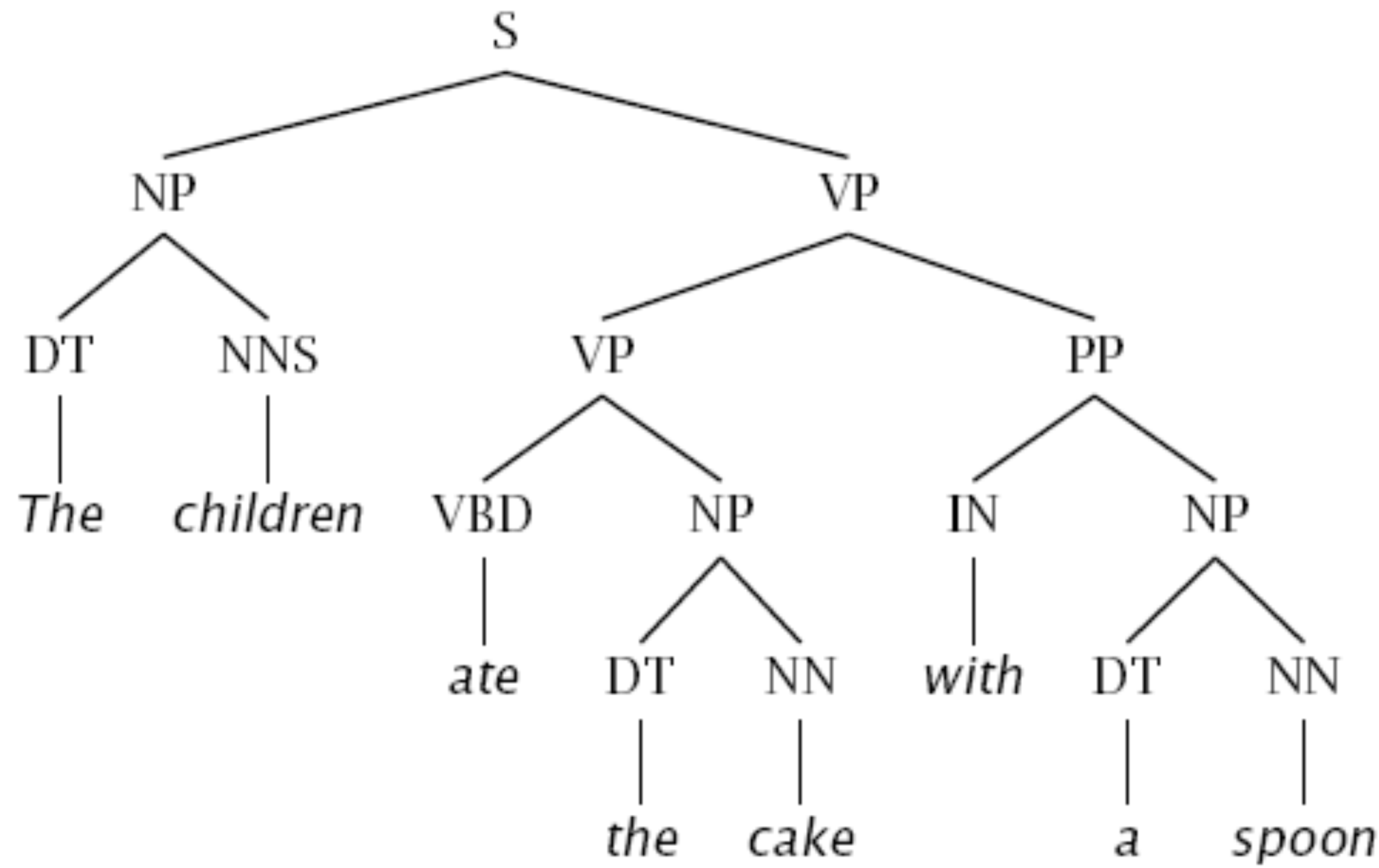
# Constituency Parsing

The rat the cat chased squeaked

I raced to Indianapolis , unimpeded by traffic
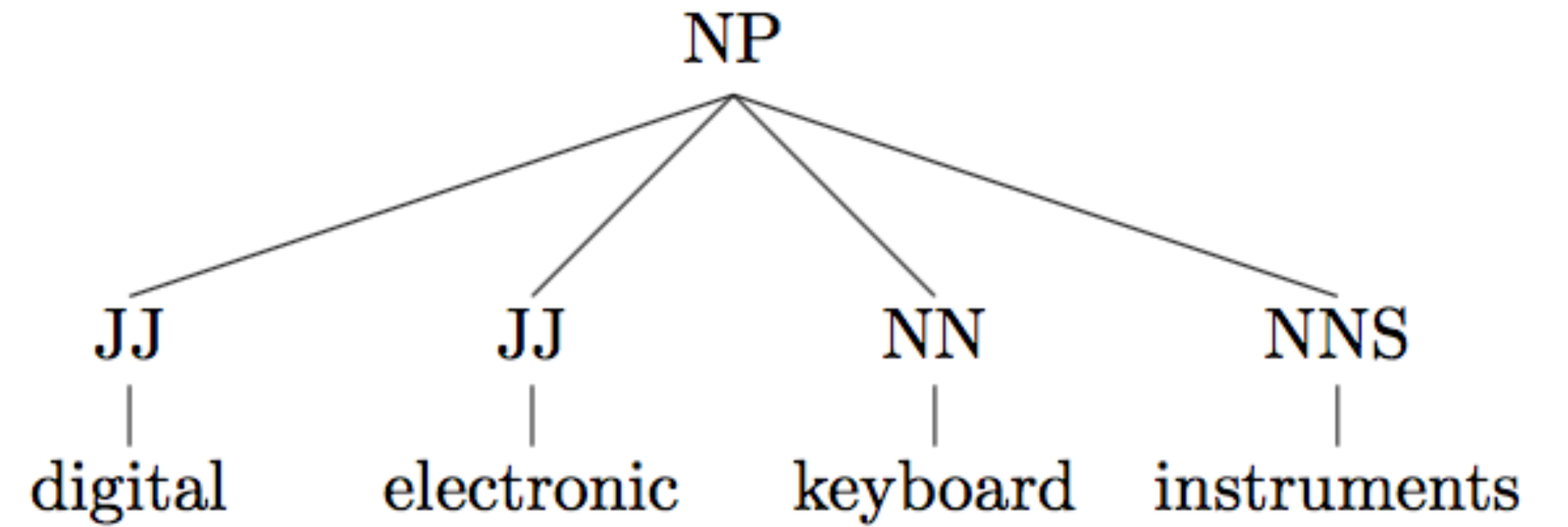
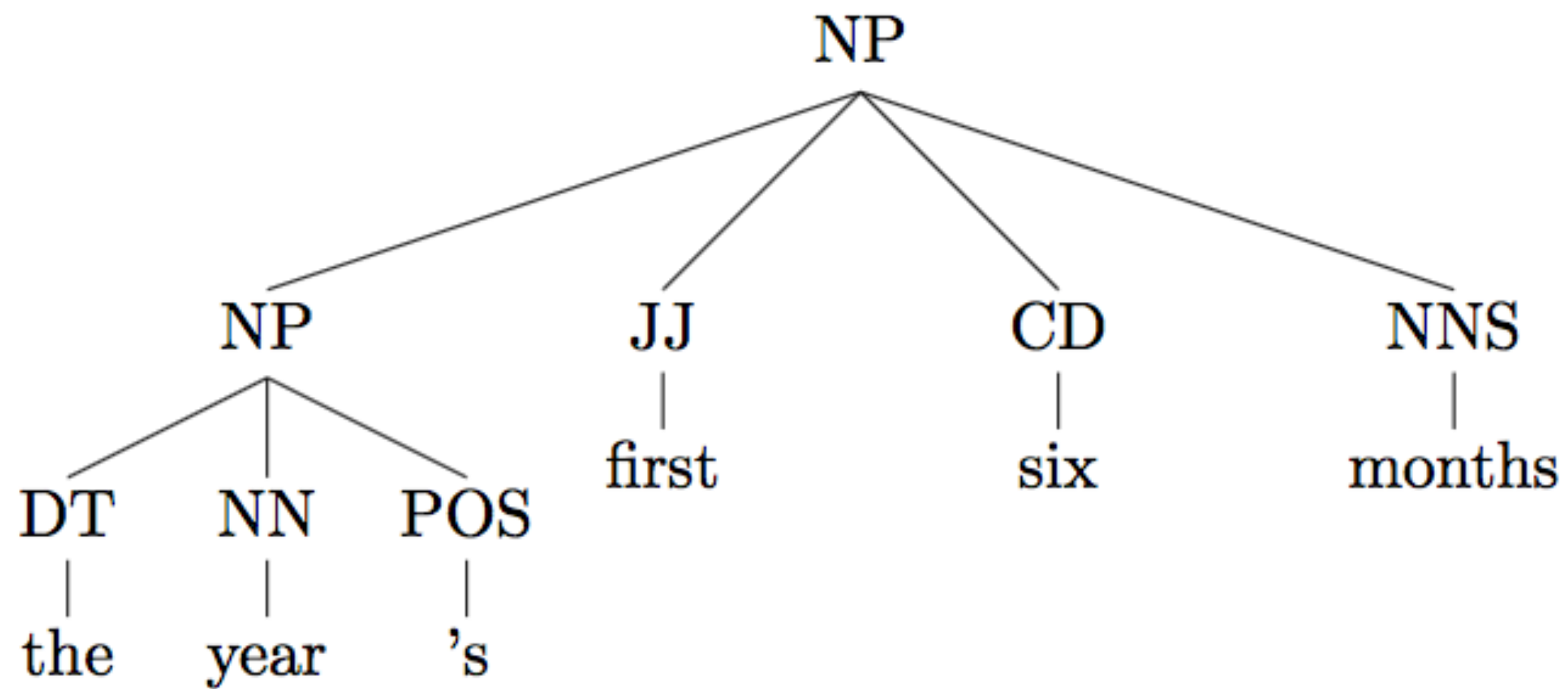# Challenges

▸ PP attachment



same parse as "the cake with some icing"

# Challenges

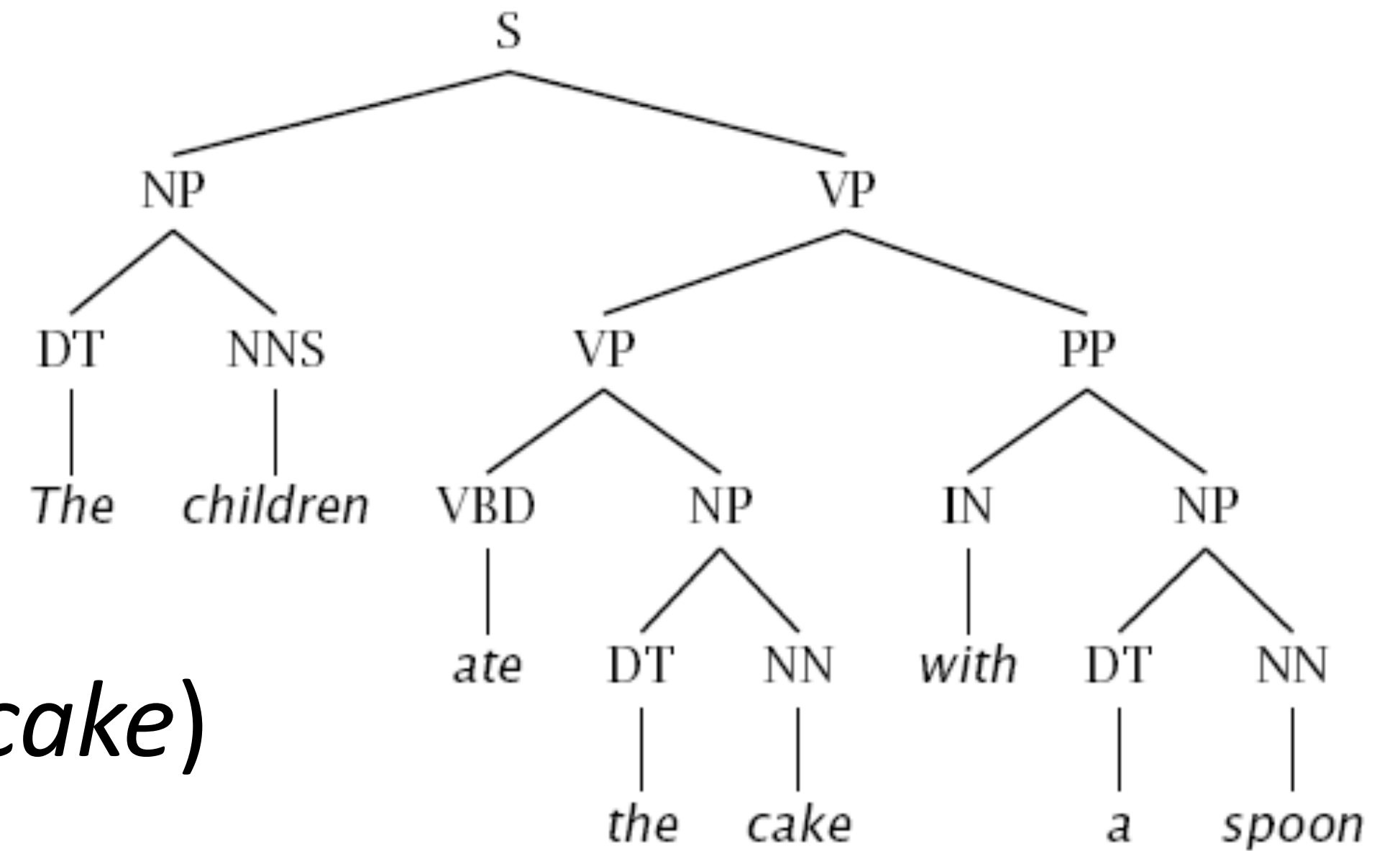▸ NP internal structure: tags + depth of analysis

# Constituency

▸ How do we know what the constituents are?

▸ Constituency tests:

   ▸ Substitution by *proform* (e.g., pronoun)

   ▸ Clefting (*It was with a spoon that...*)

   ▸ Answer ellipsis (What did they eat? *the cake*)
   (How? *with a spoon*)



▸ Sometimes constituency is not clear, e.g., coordination: *she went to and bought food at the store*

# Context-Free Grammars, CKY

# CFGs and PCFGs

## Grammar (CFG)

| | | | |
|---|---|---|---|
| ROOT → S | 1.0 | NP → NP PP | 0.3 |
| S → NP VP | 1.0 | VP → VBP NP | 0.7 |
| NP → DT NN | 0.2 | VP → VBP NP PP | 0.3 |
| NP → NN NNS | 0.5 | PP → IN NP | 1.0 |

## Lexicon

| | |
|---|---|
| NN → interest | 1.0 |
| NNS → raises | 1.0 |
| VBP → interest | 1.0 |
| VBZ → raises | 1.0 |

▸ Context-free grammar: symbols which rewrite as one or more symbols

▸ Lexicon consists of "preterminals" (POS tags) rewriting as terminals (words)

▸ CFG is a tuple (N, T, S, R): N = nonterminals, T = terminals, S = start symbol (generally a special ROOT symbol), R = rules

▸ PCFG: probabilities associated with rewrites, normalize by source symbol

# Estimating PCFGs

▸ Tree *T* is a series of rule applications *r*.
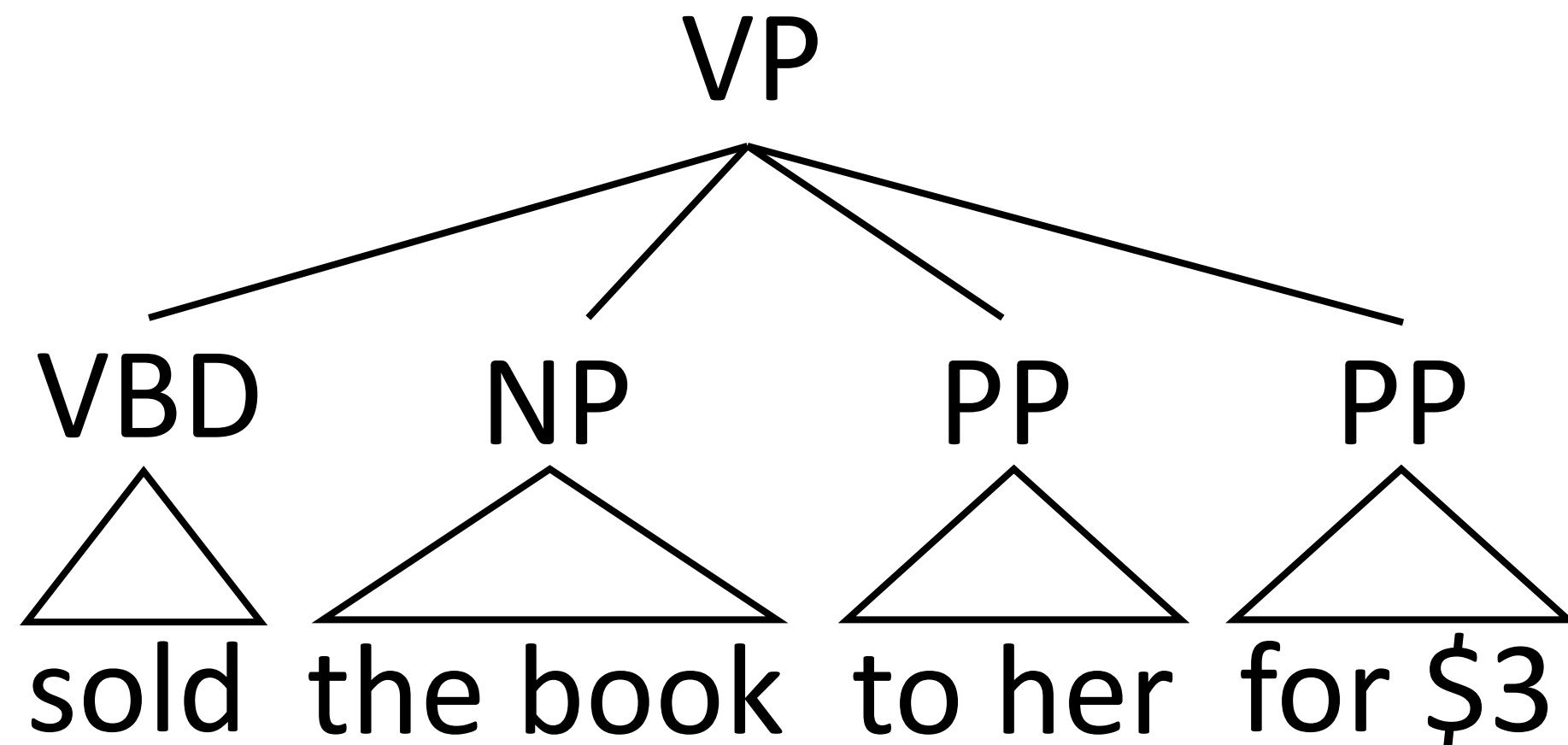
$$P(T) = \prod_{r \in T} P(r | \operatorname{parent}(r))$$



$\longrightarrow$

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

▸ Maximum likelihood PCFG: count and normalize! Same as HMMs / Naive Bayes

# Binarization

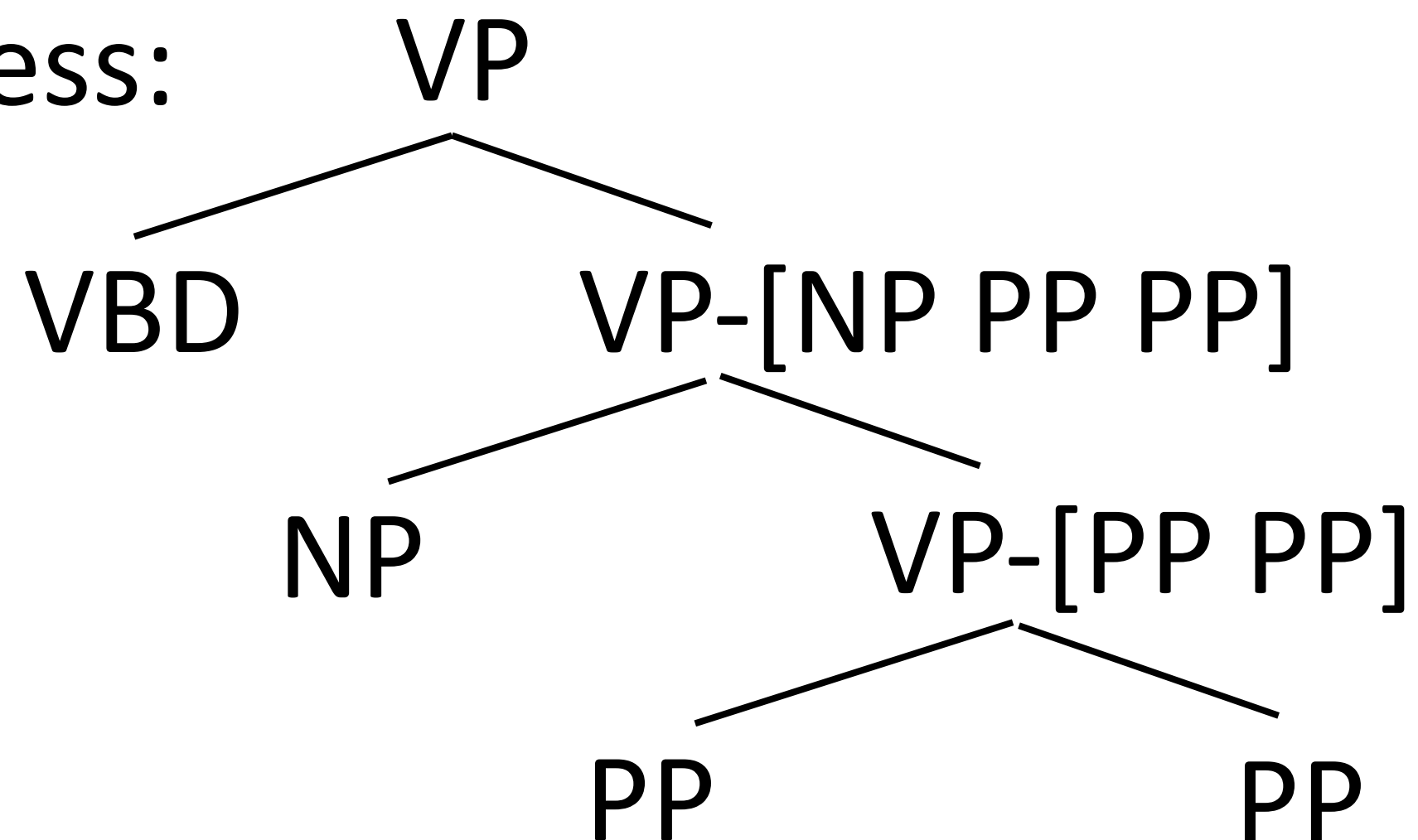▶ To parse efficiently, we need our PCFGs to be at most binary (not CNF)

VP
├── VBD — sold
├── NP — the book
├── PP — to her
└── PP — for $3

$P(VP \rightarrow VBD\ NP\ PP\ PP) = 0.2$

$P(VP \rightarrow VBZ\ PP) = 0.1$

...

▶ Lossless:

VP
├── VBD
└── VP-[NP PP PP]
    ├── NP
    └── VP-[PP PP]
        ├── PP
        └── PP

▶ Lossy:

VP
├── VBD
└── VP
    ├── NP
    └── VP
        ├── PP
        └── PP

# Chomsky Normal Form

▸ Lossless:



P(VP → VBD VP-[NP PP PP]) = 0.2

P(VP-[NP PP PP] → NP VP-[PP PP]) = 1.0

P(VP-[PP PP] → PP PP) = 1.0

▸ Deterministic symbols make this the same as before

▸ Lossy:



P(VP → VBD VP) = 0.2

P(VP → NP VP) = 0.03

P(VP → PP PP) = 0.001

▸ Makes different independent assumptions, not the same PCFG

# CKY

- Find argmax P(T|x) = argmax P(T, **x**)

- Dynamic programming: chart maintains the best way of building symbol X over span (i, j)

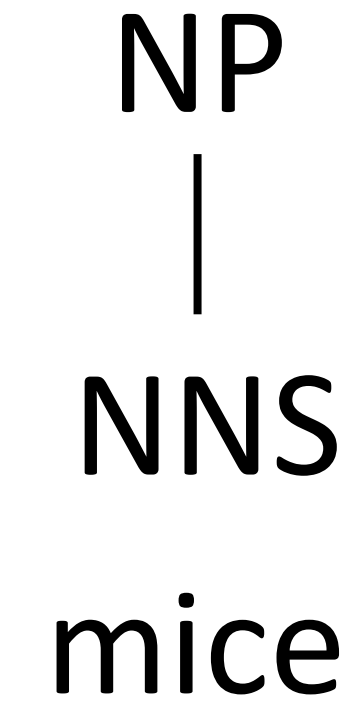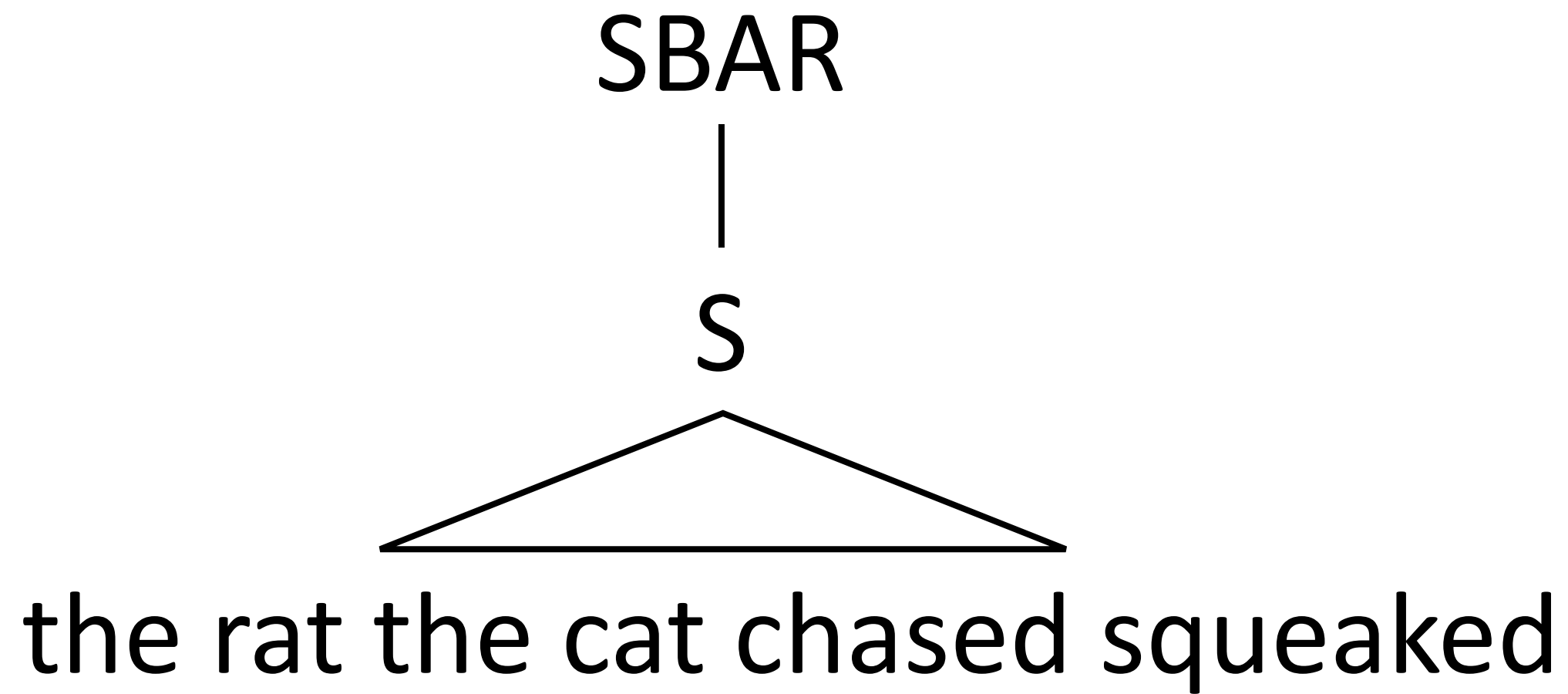- Loop over all split points k, apply rules X -> Y Z to build X in every possible way

- CKY = Viterbi, also an algorithm called inside-outside = forward-backward



He wrote a long report on Mars

Cocke-Kasami-Younger

# Unary Rules

SBAR
|
S
◿◺
the rat the cat chased squeaked

NP
|
NNS

mice

▸ Unary productions in treebank need to be dealt with by parsers

▸ Binary trees over n words have at most n-1 nodes, but you can have unlimited numbers of nodes with unaries (S → SBAR → NP → S → …)

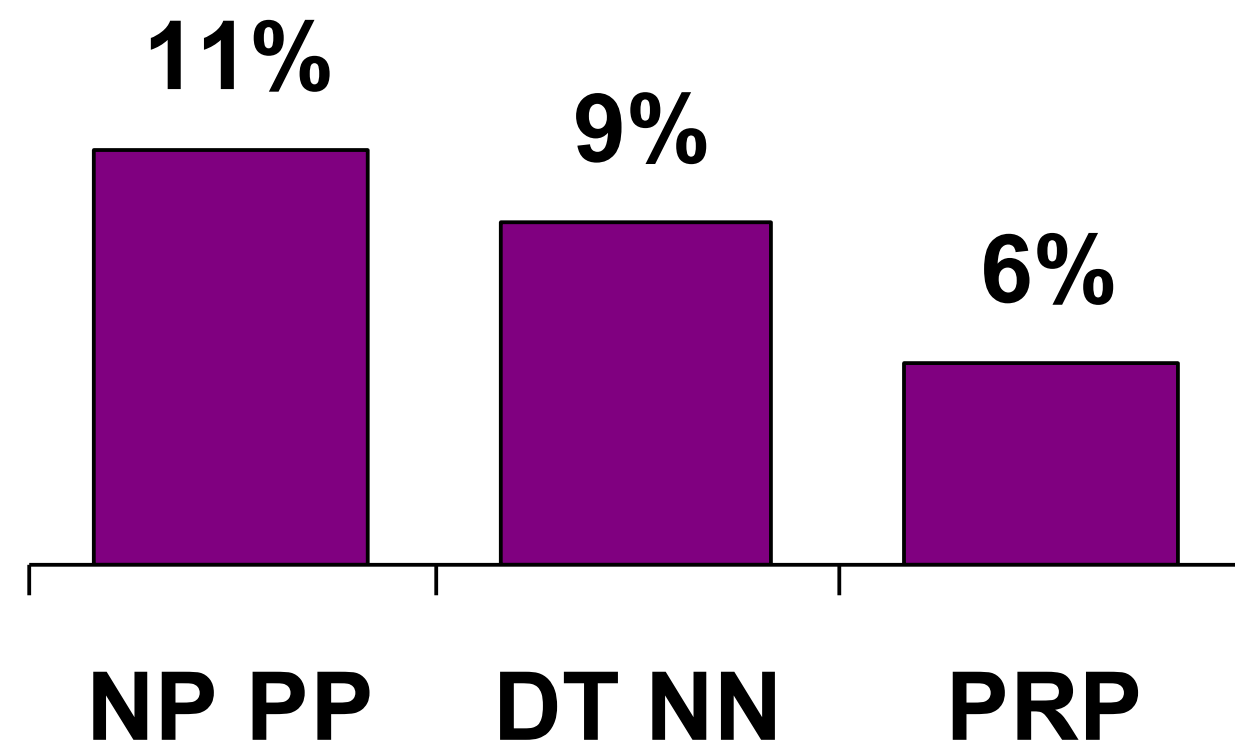▸ In practice: enforce at most one unary over each span, modify CKY accordingly

# Results

- Standard dataset for English: Penn Treebank (Marcus et al., 1993)

  - Evaluation: F1 over labeled constituents of the sentence

- Vanilla PCFG: ~75 F1

- Best PCFGs for English: ~90 F1

- SOTA: 95 F1

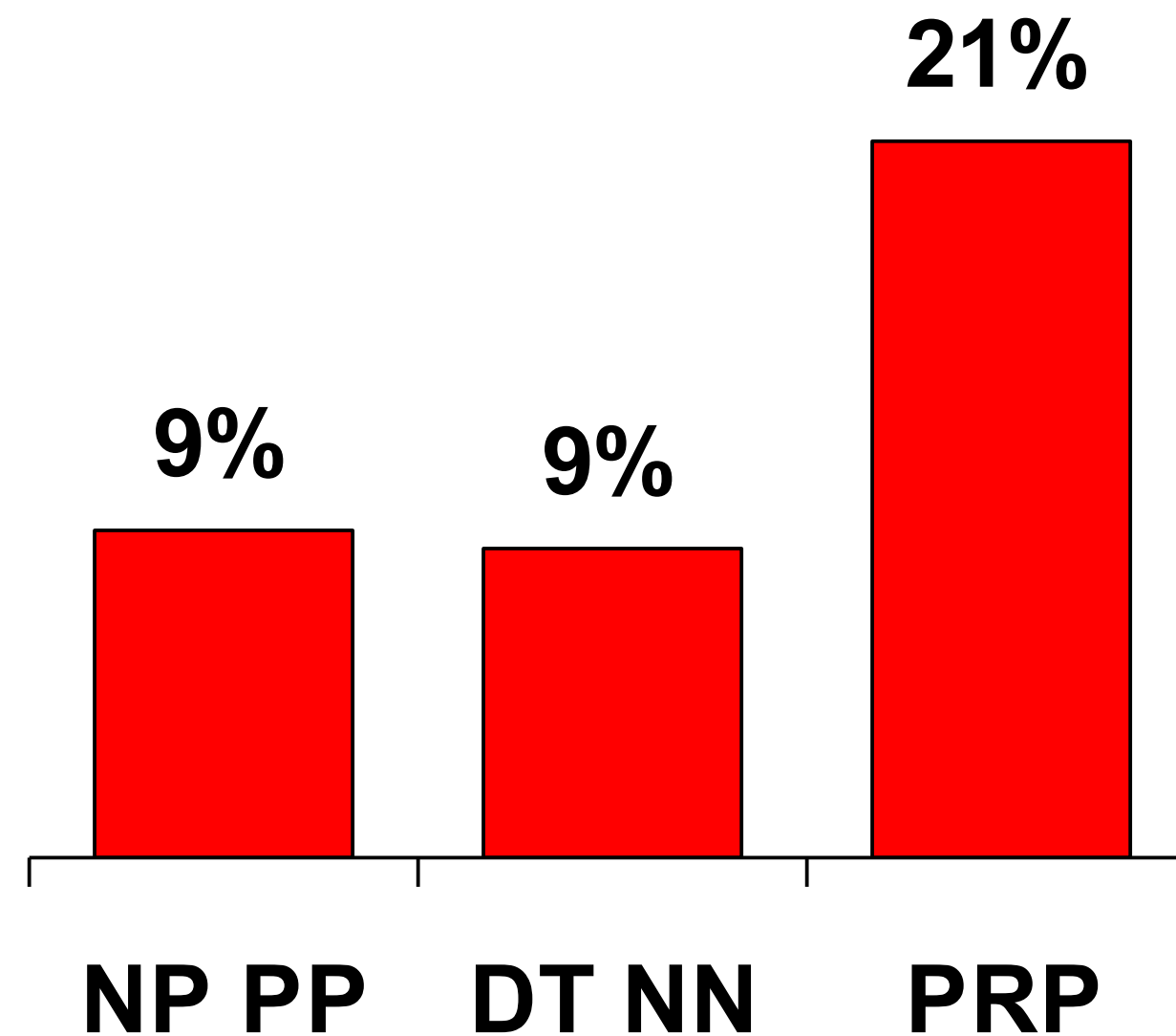- Other languages: results vary widely depending on annotation + complexity of the grammar

Klein and Manning (2003)

# Refining Generative Grammars

# PCFG Independence Assumptions

### All NPs

**11%** NP PP  
**9%** DT NN  
**6%** PRP

### NPs under S

**9%** NP PP  
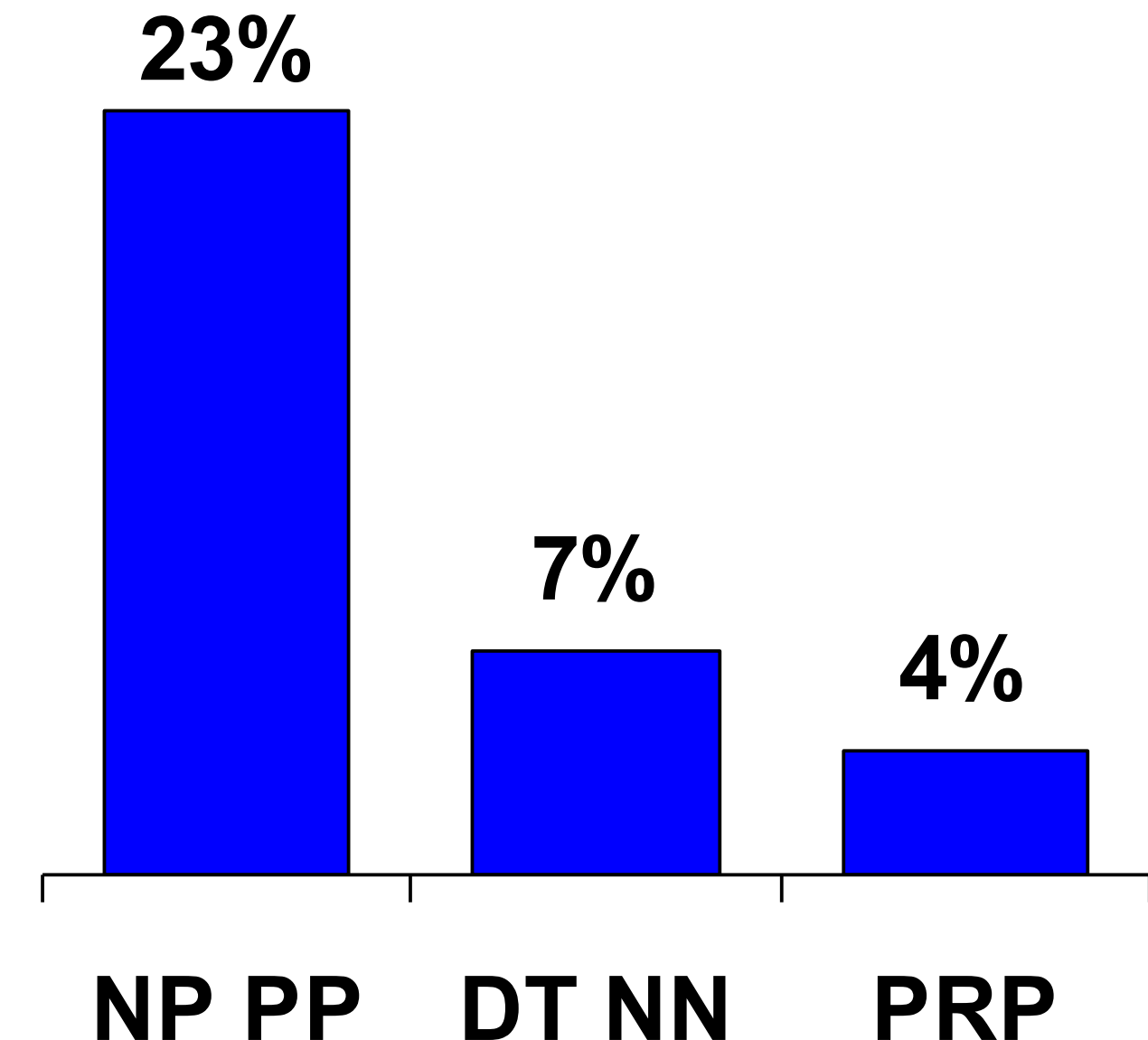**9%** DT NN  
**21%** PRP

### NPs under VP

**23%** NP PP  
**7%** DT NN  
**4%** PRP
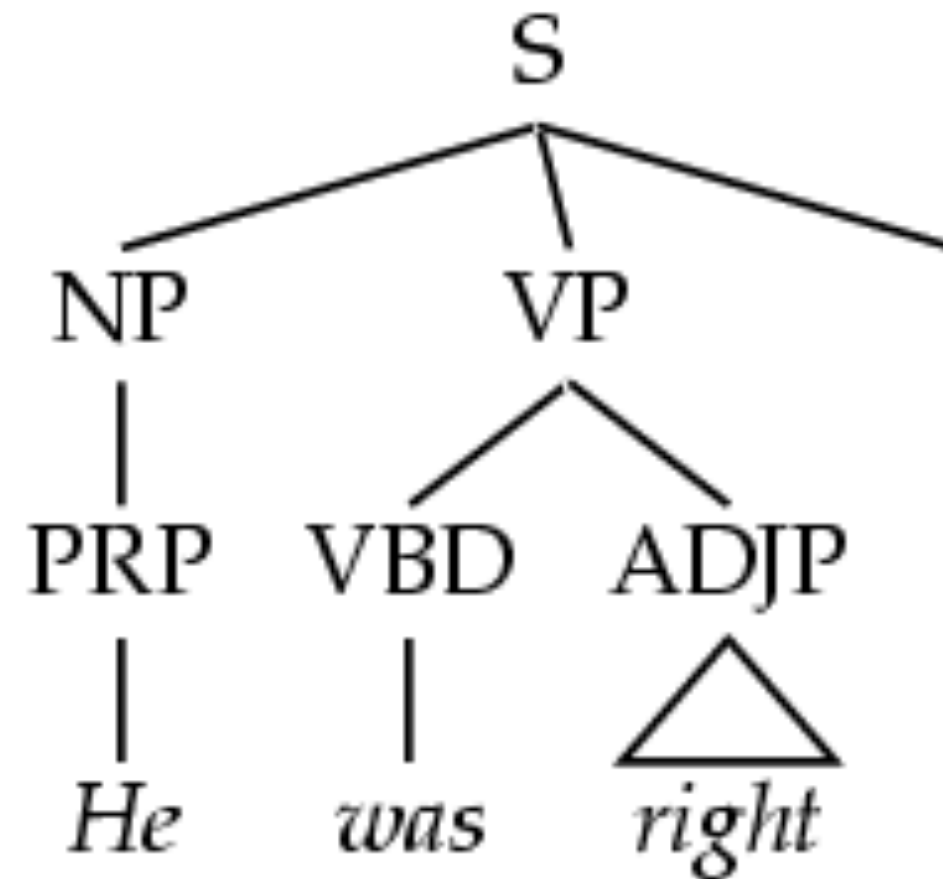
▸ Language is not context-free: NPs in different contexts rewrite differently
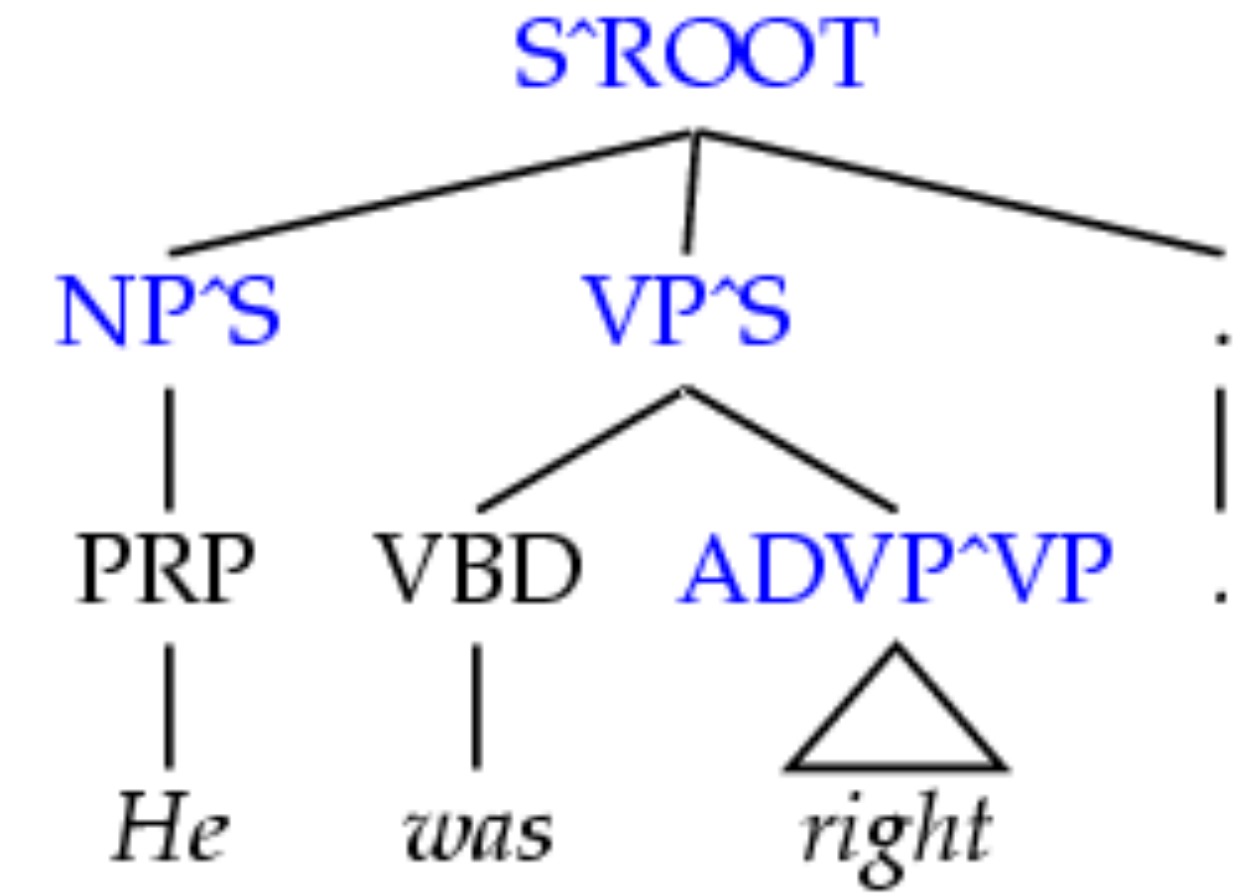
▸ Can we make the grammar "less context-free"?

# Rule Annotation

- Like a trigram HMM tagger, incorporates more context

- Vertical (parent) annotation: add the parent symbol to each node, can do grandparents too

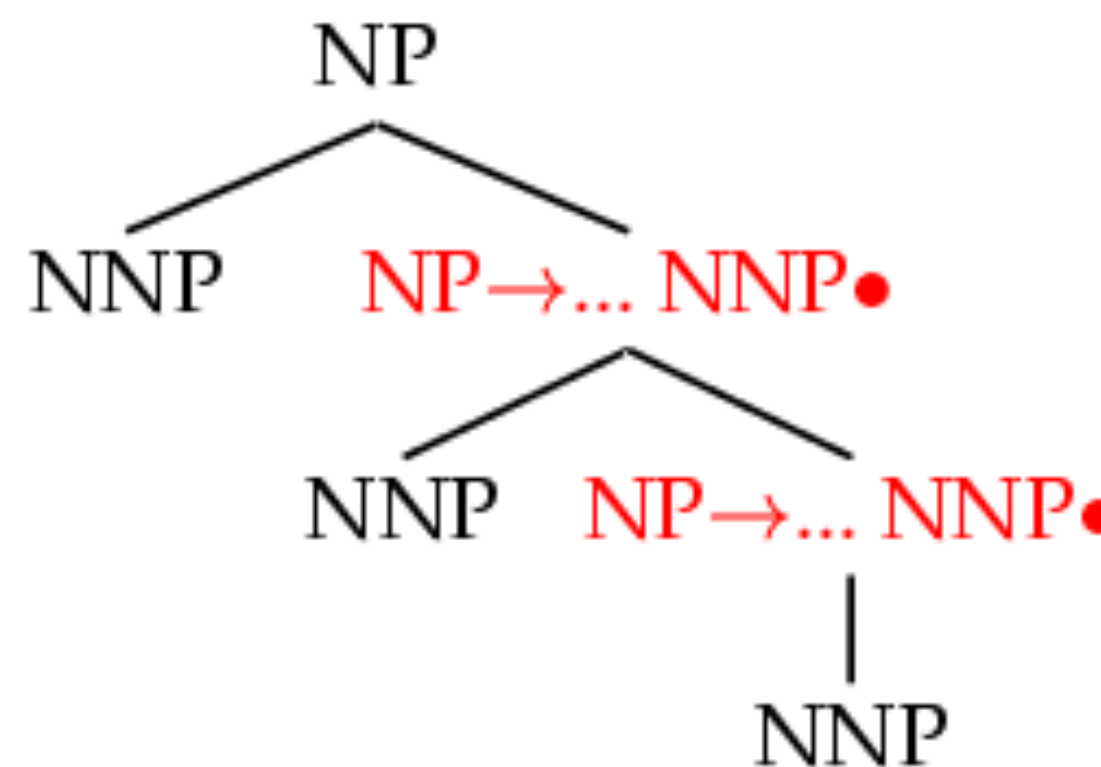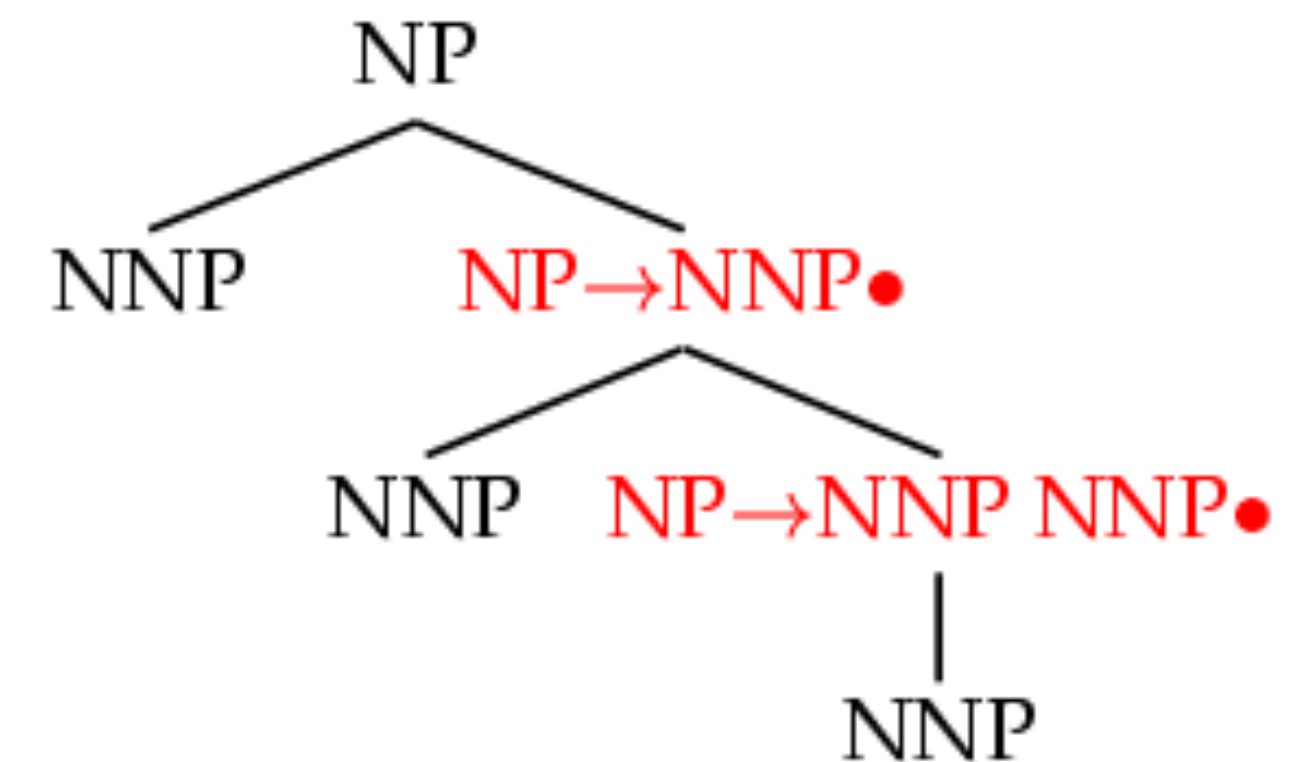- Horizontal annotation: remember the states of multi-arity rules during binarization

Order 1

Order 2



Order 1

Order ∞

# Annotated Tree



▸ 75 F1 with basic PCFG => 86.3 F1 with this highly customized PCFG (SOTA was 90 F1 at the time, but with more complex methods)

Klein and Manning (2003)

# Lexicalized Parsers



▸ Even with parent annotation, these trees have the same rules. Need to use the words

# Lexicalized Parsers

▸ Annotate each grammar symbol with its "head word": most important word of that constituent

▸ Rules for identifying headwords (e.g., the last word of an NP before a preposition is typically the head)

▸ Collins and Charniak (late 90s): ~89 F1 with these

# Discriminative Parsers

# CRF Parsing

# CRF Parsing

$$\text{score}\left( \begin{array}{c} \text{NP} \\ \diagup \quad \diagdown \\ \text{NP} \qquad \text{PP} \\ \diagup\diagdown \quad \diagup\diagdown \\ \text{He wrote a long report on Mars .} \\ {}_{2} \qquad\qquad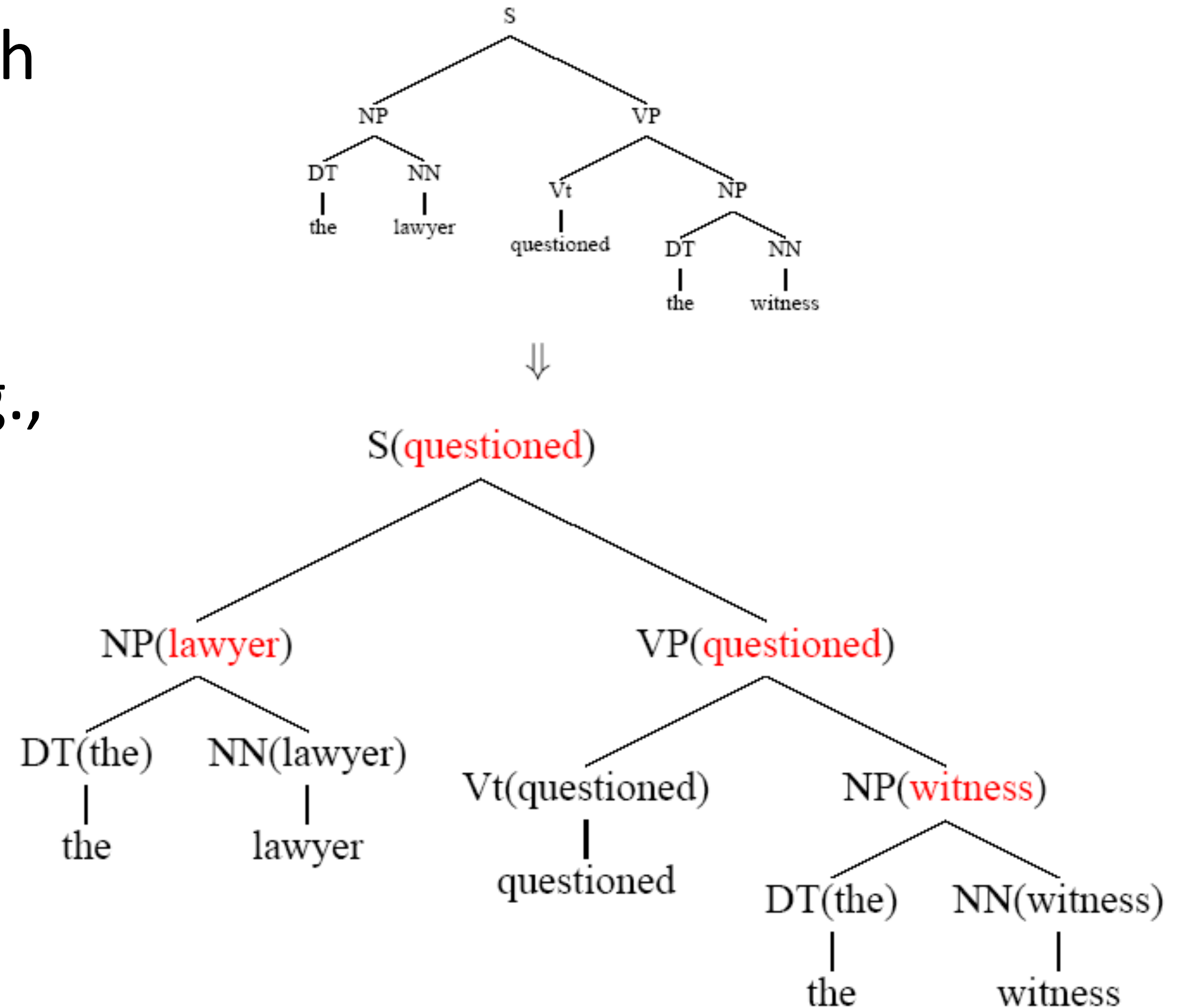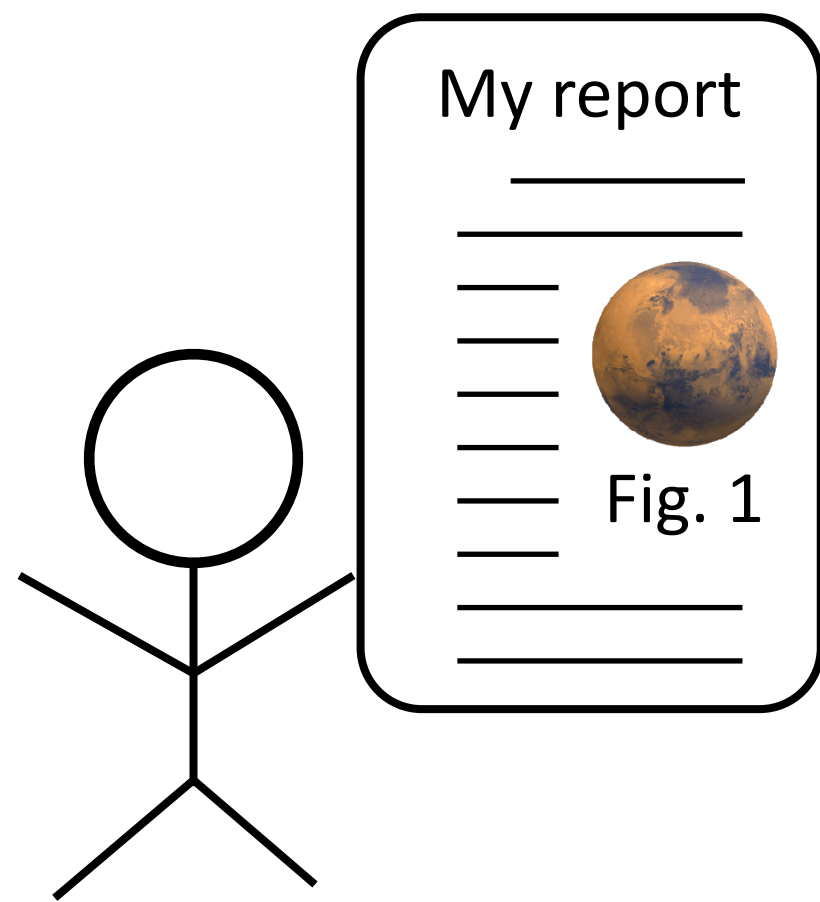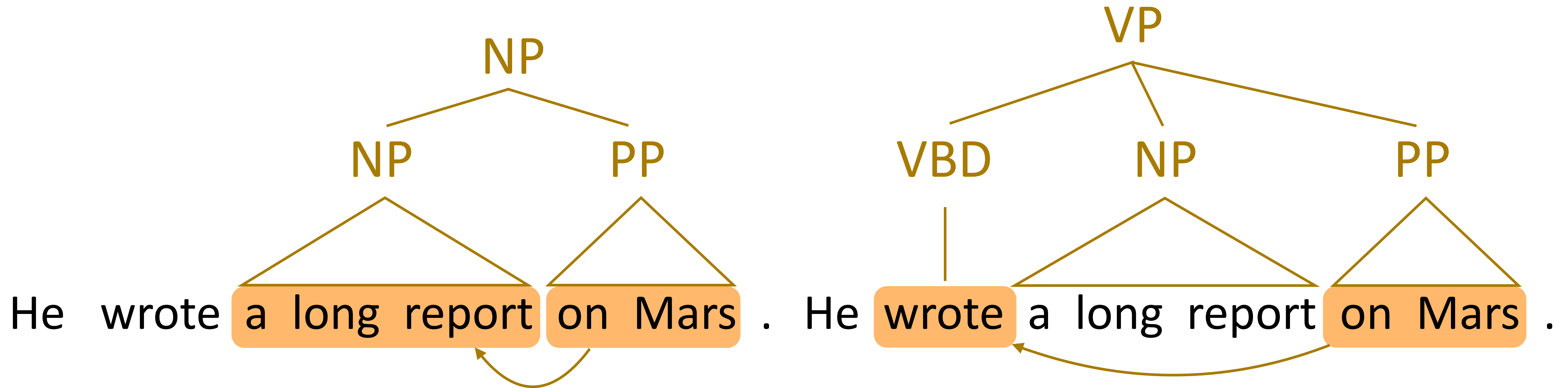 {}_{5} \qquad {}_{7} \end{array} \right) = w^{\top} f \left( \begin{array}{c} \text{NP} \\ \frown \\ {}_{2}\text{NP}_{5}\text{PP}_{7} \\ {}_{\text{wrote a long report on Mars .}} \end{array} \right)$$

$$f \left( \begin{array}{c} \text{NP} \\ \frown \\ {}_{2}\text{NP}_{5}\text{PP}_{7} \\ {}_{\text{wrote a long report on Mars .}} \end{array} \right) = \boxed{\bullet\circ\circ\bullet\circ\bullet\circ\circ\circ\bullet}$$

Left child last word = *report* $\wedge$ $\begin{array}{c} \text{NP} \\ \diagup\diagdown \\ \text{NP} \quad \text{PP} \end{array}$

▸ Can learn that we *report* [PP], which is common due to *reporting on* things

▸ Can "neuralize" this as well like neural CRFs for NER

Taskar et al. (2004)

Hall, Durrett, and Klein (2014)

Durrett and Klein (2015)

# Joint Discrete and Continuous Parsing

▸ Chart remains discrete!



Discrete + Continuous    Discrete + Continuous    …

NP

NP            PP

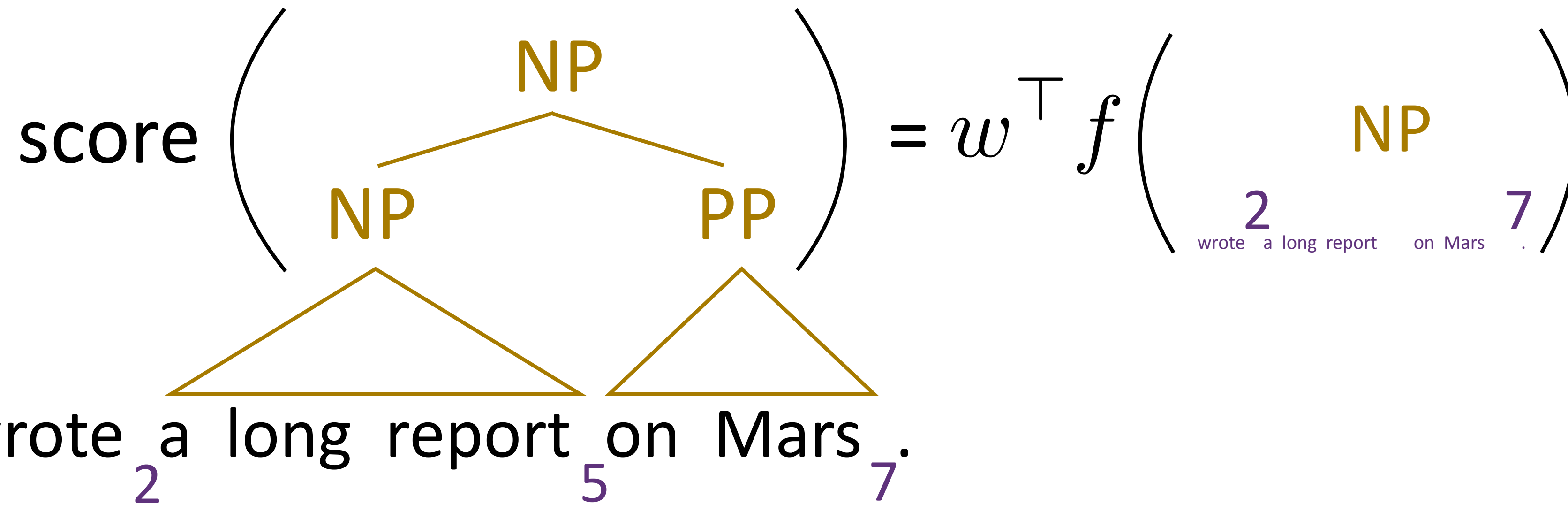He   wrote   a   long   report   on   Mars

Parsing a sentence:

▸ Feedforward pass on nets

▸ Discrete feature computation

▸ Run CKY dynamic program

Durrett and Klein (ACL 2015)

# Neural CRF Parsing

▸ Simpler version: score *constituents* rather than rule applications

$$\text{score} \left( \begin{array}{c} \text{NP} \\ \text{NP} \quad \text{PP} \end{array} \right) = w^{\top} f \left( \begin{array}{c} \text{NP} \\ 2 \qquad 7 \\ \text{wrote} \quad \text{a long report} \quad \text{on Mars} \quad . \end{array} \right)$$

He wrote a long report on Mars .
$2$ $5$ $7$

▸ Use BiLSTMs (Stern) or self-attention (Kitaev) to compute span embeddings

▸ 91-93 F1, 95 F1 with ELMo (SOTA). Great on other langs too!

Stern et al. (2017),
Kitaev et al. (2018)

# Takeaways

- PCFGs estimated generatively can perform well if sufficiently engineered

- Neural CRFs work well for constituency parsing

- Next time: revisit lexicalized parsing as *dependency parsing*

# Survey

▸ Write one thing you like about the class

▸ Write one thing you don't like about the class