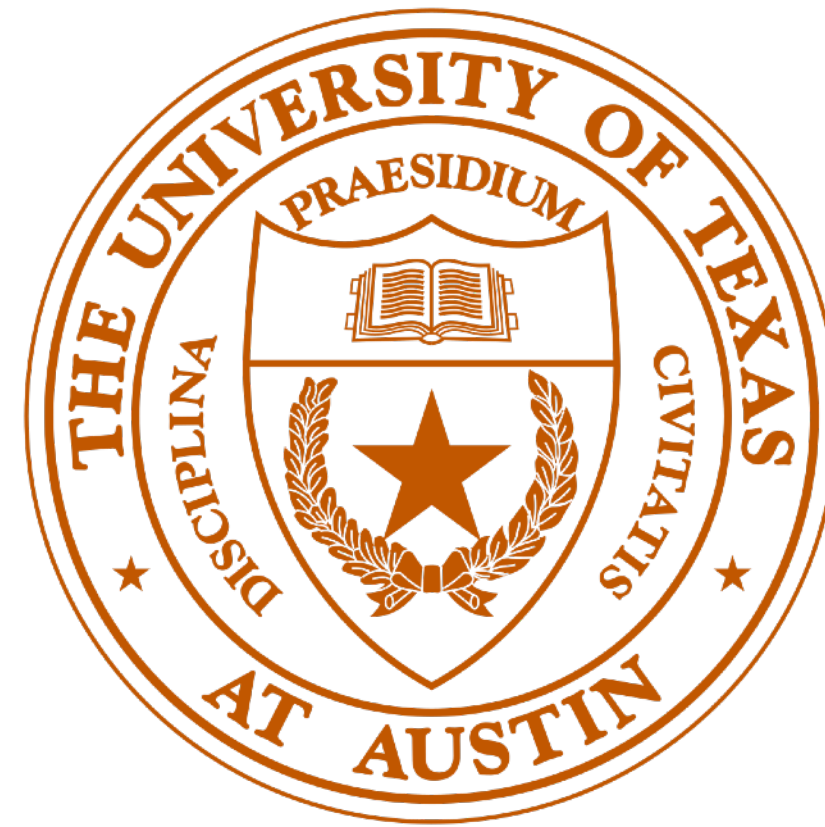
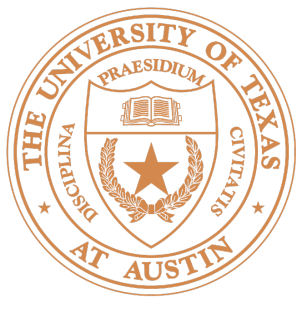


# CS378: Natural Language Processing

## Lecture 10: Seq 3 / Syntax I



Greg Durrett



# Announcements

---

- ▶ A2 due today
- ▶ A3 out tomorrow
- ▶ Midterm: list of topics next week. Covers content up to March 7
  - ▶ CRFs will NOT be on the midterm, a couple other topics too



# Today

---

- ▶ Conditional random fields
- ▶ Named entity recognition
- ▶ Syntax and constituency parsing

# CRFs and NER



# Named Entity Recognition

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 .*

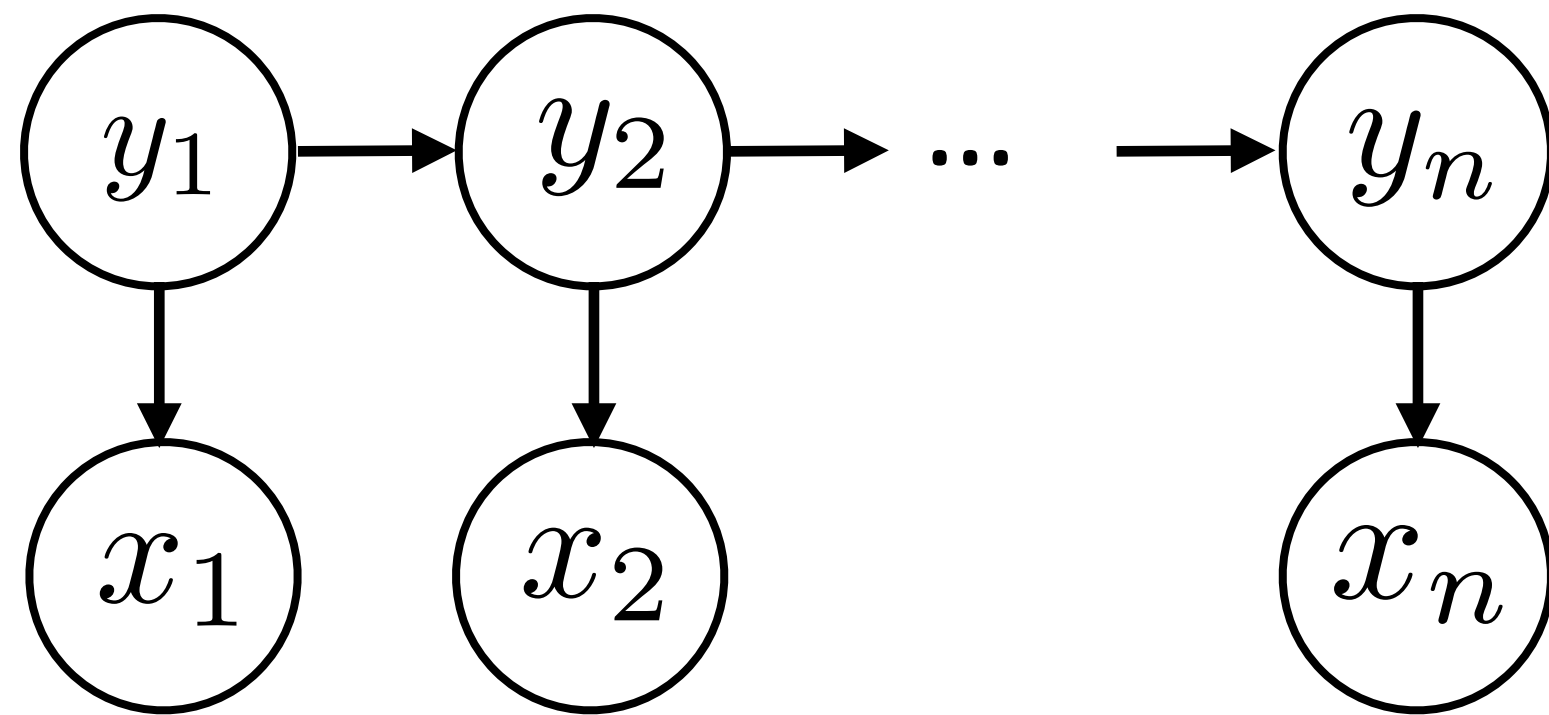
PERSON LOC ORG

- ▶ Frame as a sequence problem with a BIO tagset: begin, inside, outside
- ▶ Why might an HMM not do so well here?
  - ▶ Lots of O's, so tags aren't as informative about context
  - ▶ Need sub-word features on unknown words
- ▶ CRFs are discriminative models that will solve these problems



# Conditional Random Fields

- ▶ HMMs are expressible as Bayes nets (factor graphs)



- ▶ This reflects the following decomposition:

$$P(\mathbf{y}, \mathbf{x}) = P(y_1)P(x_1|y_1)P(y_2|y_1)P(x_2|y_2) \dots$$

- ▶ Locally normalized model: each factor is a probability distribution that normalizes



# Conditional Random Fields

- ▶ HMMs:  $P(\mathbf{y}, \mathbf{x}) = P(y_1)P(x_1|y_1)P(y_2|y_1)P(x_2|y_2) \dots$
- ▶ CRFs: discriminative models with the following globally-normalized form:

$$P(\mathbf{y}|\mathbf{x}) = \frac{\prod_k \exp(\phi_k(\mathbf{x}, \mathbf{y}))}{\sum_{\mathbf{y}'} \prod_k \exp(\phi_k(\mathbf{x}, \mathbf{y}'))}$$

any real-valued scoring function of its arguments

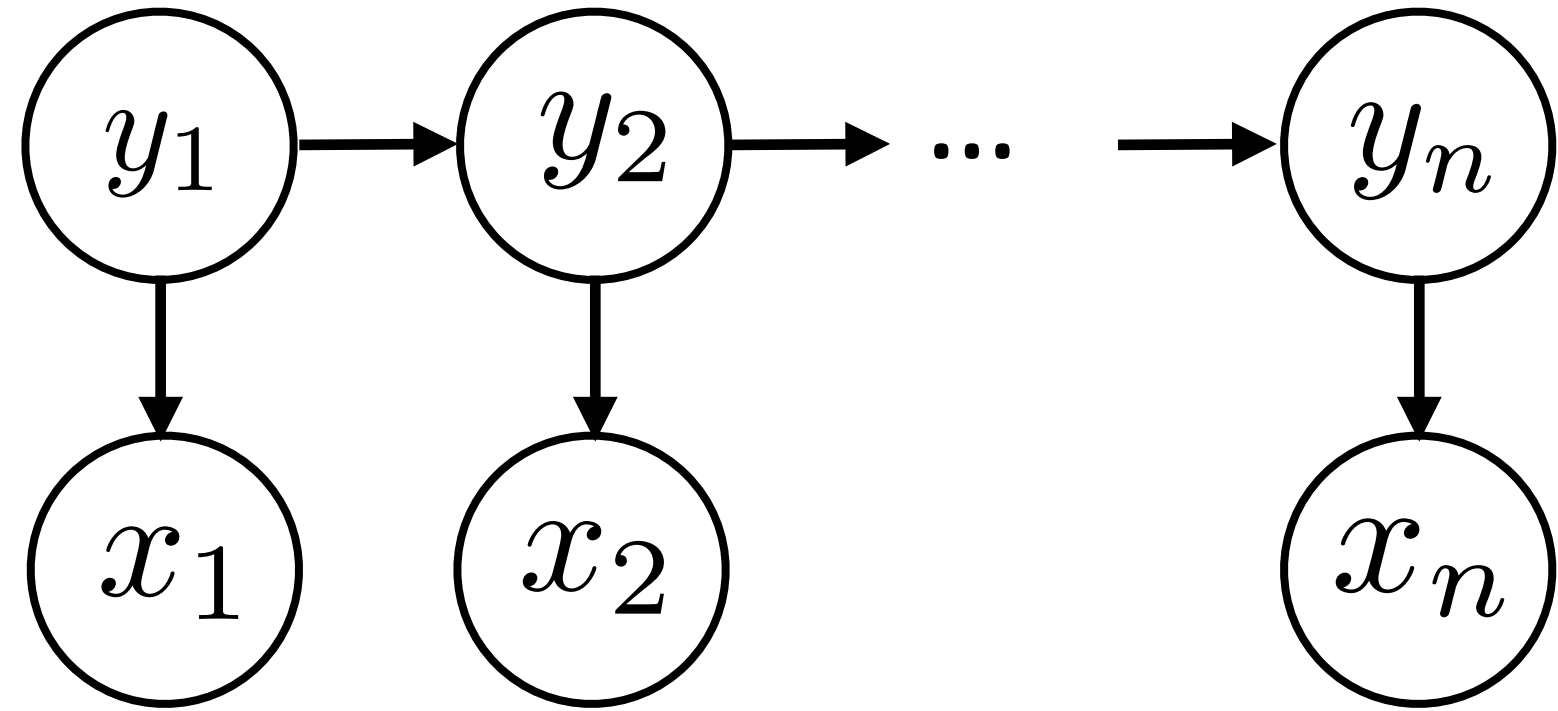
normalizer  $Z$

- ▶ Naive Bayes : logistic regression :: HMMs : CRFs  
local vs. global normalization <-> generative vs. discriminative
- ▶ How do we max over  $\mathbf{y}$ ? Requires considering an exponential number of sequences in general



# Sequential CRFs

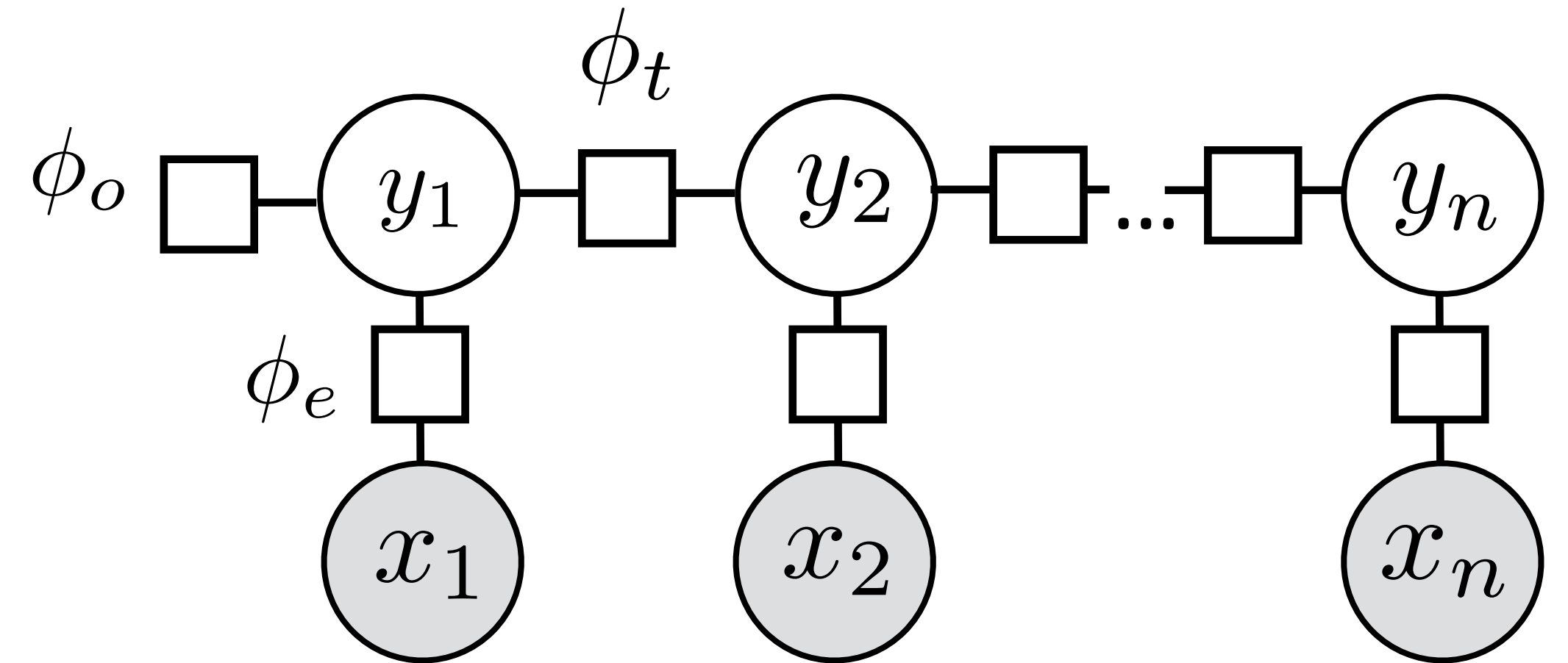
- ▶ HMMs:  $P(\mathbf{y}, \mathbf{x}) = P(y_1)P(x_1|y_1)P(y_2|y_1)P(x_2|y_2) \dots$



- ▶ CRFs:

$$P(\mathbf{y}|\mathbf{x}) \propto \prod_k \exp(\phi_k(\mathbf{x}, \mathbf{y}))$$

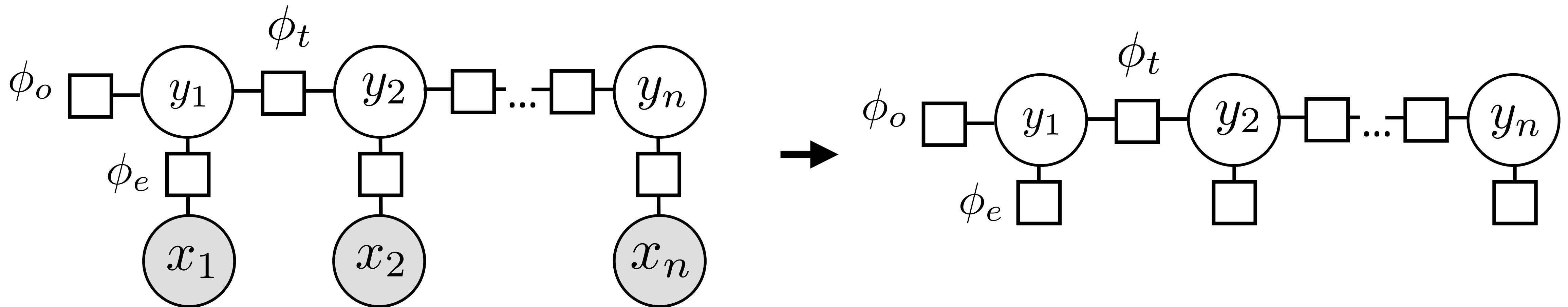
$$P(\mathbf{y}|\mathbf{x}) \propto \exp(\phi_o(y_1)) \prod_{i=2}^n \exp(\phi_t(y_{i-1}, y_i)) \prod_{i=1}^n \exp(\phi_e(x_i, y_i))$$







# Sequential CRFs



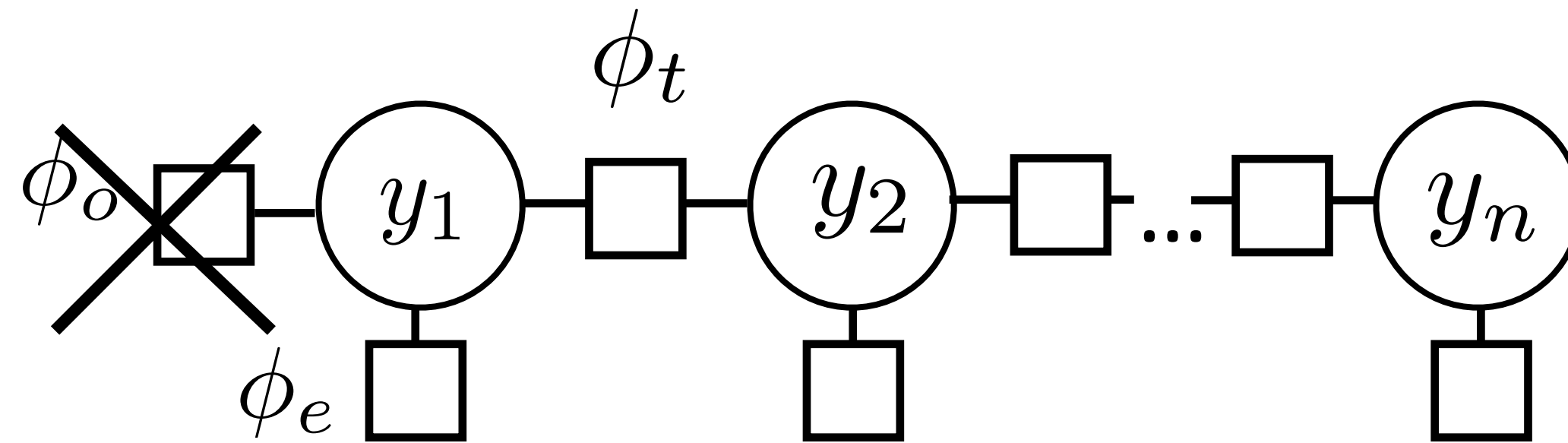
$$P(\mathbf{y}|\mathbf{x}) \propto \exp(\phi_o(y_1)) \prod_{i=2}^n \exp(\phi_t(y_{i-1}, y_i)) \prod_{i=1}^n \exp(\phi_e(x_i, y_i))$$

$\prod_{i=1}^n \exp(\phi_e(y_i, i, \mathbf{x}))$

- ▶ We condition on  $\mathbf{x}$ , so every factor can depend on all of  $\mathbf{x}$
  - ▶  $\mathbf{y}$  can't depend arbitrarily on  $\mathbf{x}$  in a generative model
- token index — lets us look at current word



# Sequential CRFs



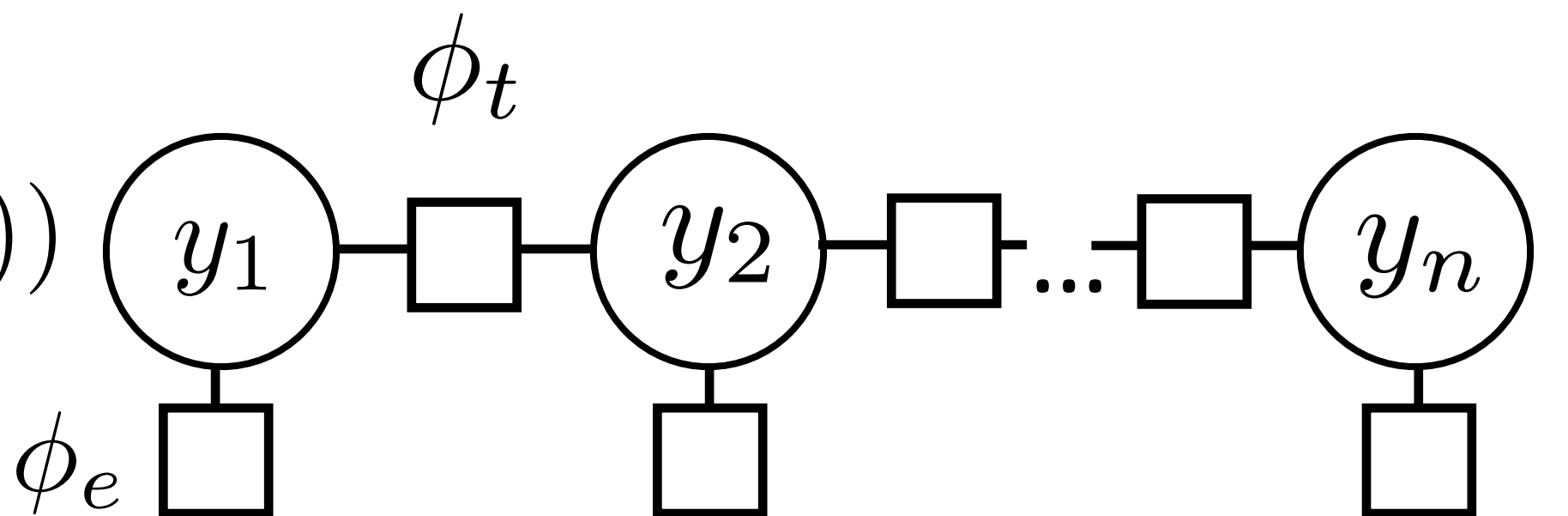
- Don't include initial distribution, can bake into other factors

Sequential CRFs:

$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z} \prod_{i=2}^n \exp(\phi_t(y_{i-1}, y_i)) \prod_{i=1}^n \exp(\phi_e(y_i, i, \mathbf{x}))$$



# Feature Functions

$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z} \prod_{i=2}^n \exp(\phi_t(y_{i-1}, y_i)) \prod_{i=1}^n \exp(\phi_e(y_i, i, \mathbf{x}))$$


- This can be almost anything! Here we use linear functions of sparse features

$$\phi_e(y_i, i, \mathbf{x}) = w^\top f_e(y_i, i, \mathbf{x}) \quad \phi_t(y_{i-1}, y_i) = w^\top f_t(y_{i-1}, y_i)$$

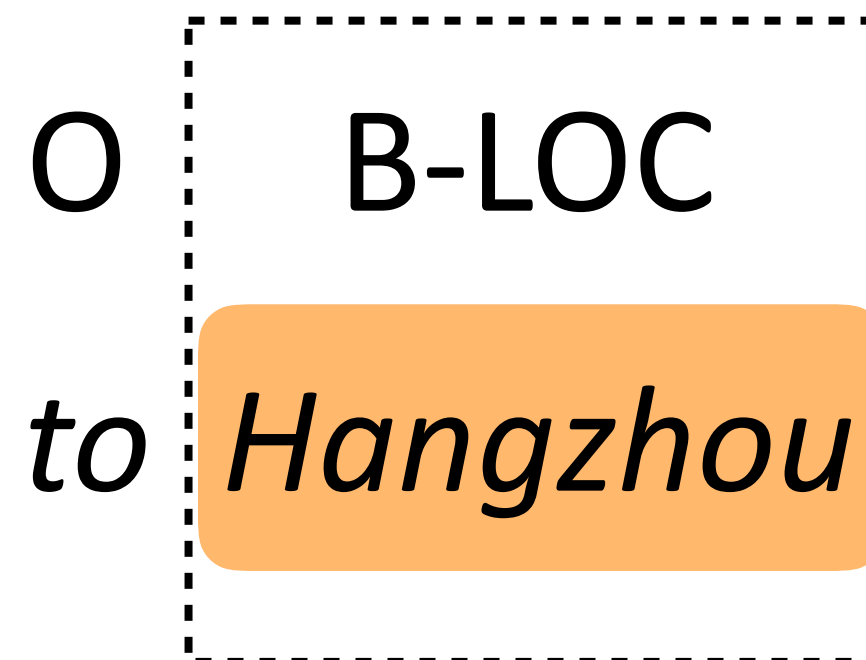
$$P(\mathbf{y}|\mathbf{x}) \propto \exp w^\top \left[ \sum_{i=2}^n f_t(y_{i-1}, y_i) + \sum_{i=1}^n f_e(y_i, i, \mathbf{x}) \right]$$

- Looks like our single weight vector multiclass logistic regression model



# Basic Features for NER

$$P(\mathbf{y}|\mathbf{x}) \propto \exp w^\top \left[ \sum_{i=2}^n f_t(y_{i-1}, y_i) + \sum_{i=1}^n f_e(y_i, i, \mathbf{x}) \right]$$



*Barack Obama will travel to Hangzhou today for the G20 meeting .*

Transitions:  $f_t(y_{i-1}, y_i) = \text{Ind}[y_{i-1} \ \& \ y_i] = \text{Ind}[O \text{ — } B\text{-LOC}]$

Emissions:  $f_e(y_6, 6, \mathbf{x}) = \text{Ind}[B\text{-LOC} \ \& \ \text{Current word} = \textit{Hangzhou}]$   
 $\text{Ind}[B\text{-LOC} \ \& \ \text{Prev word} = \textit{to}]$



# CRFs Outline

► **Model:** 
$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z} \prod_{i=2}^n \exp(\phi_t(y_{i-1}, y_i)) \prod_{i=1}^n \exp(\phi_e(y_i, i, \mathbf{x}))$$

$$P(\mathbf{y}|\mathbf{x}) \propto \exp w^\top \left[ \sum_{i=2}^n f_t(y_{i-1}, y_i) + \sum_{i=1}^n f_e(y_i, i, \mathbf{x}) \right]$$

- Inference:  $\operatorname{argmax} P(\mathbf{y}|\mathbf{x})$  from Viterbi
- Learning: requires running sum-product Viterbi to compute posterior probabilities  $P(y | \mathbf{x})$  at each step  $i$



# Features for NER

- ▶ Word features (can use in HMM)
  - ▶ Capitalization
  - ▶ Word shape
  - ▶ Prefixes/suffixes
  - ▶ Lexical indicators
- ▶ Context features (can't use in HMM!)
  - ▶ Words before/after
  - ▶ Tags before/after
- ▶ Word clusters
- ▶ Gazetteers

*Leicestershire*

*Boston*

*Apple* released a new version...

According to the *New York Times*...





# Evaluating NER

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

- ▶ Prediction of all Os still gets 66% accuracy on this example!
- ▶ What we really want to know: how many named entity *chunk* predictions did we get right?
  - ▶ Precision: of the ones we predicted, how many are right?
  - ▶ Recall: of the gold named entities, how many did we find?
  - ▶ F-measure: harmonic mean of these two



# NER

- ▶ CRF with lexical features can get around 85 F1 on this problem
- ▶ Other pieces of information that many systems capture
- ▶ World knowledge:

The delegation met the president at the airport, **Tanjug** said.

**Tanjug**

From Wikipedia, the free encyclopedia

**Tanjug** (/ˈtʌnjʊɡ/) ([Serbian Cyrillic](#): Танјуг) is a Serbian state news agency based in [Belgrade](#).<sup>[2]</sup>





# Nonlocal Features

The news agency **Tanjug** reported on the outcome of the meeting.

ORG?

PER?

The delegation met the president at the airport, **Tanjug** said.

- ▶ More complex factor graph structures can let you capture this, or just decode sentences in order and use features on previous sentences



# How well do NER systems do?

	System	Resources Used	$F_1$
+	LBJ-NER	Wikipedia, Nonlocal Features, Word-class Model	90.80
-	(Suzuki and Isozaki, 2008)	Semi-supervised on 1G-word unlabeled data	89.92
-	(Ando and Zhang, 2005)	Semi-supervised on 27M-word unlabeled data	89.31
-	(Kazama and Torisawa, 2007a)	Wikipedia	88.02
-	(Krishnan and Manning, 2006)	Non-local Features	87.24
-	(Kazama and Torisawa, 2007b)	Non-local Features	87.17
+	(Finkel et al., 2005)	Non-local Features	86.86

Ratinov and Roth (2009)

Lample et al. (2016)

LSTM-CRF (no char)	90.20
LSTM-CRF	<b>90.94</b>
S-LSTM (no char)	87.96
S-LSTM	90.33

BiLSTM-CRF + ELMo  
Peters et al. (2018) **92.2**



# Takeaways

---

- ▶ CRFs are structured feature-based models
- ▶ Efficient to do inference and learning using dynamic programs
- ▶ Looks like logistic regression, but requires more effort to implement

# Constituency Parsing



# 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



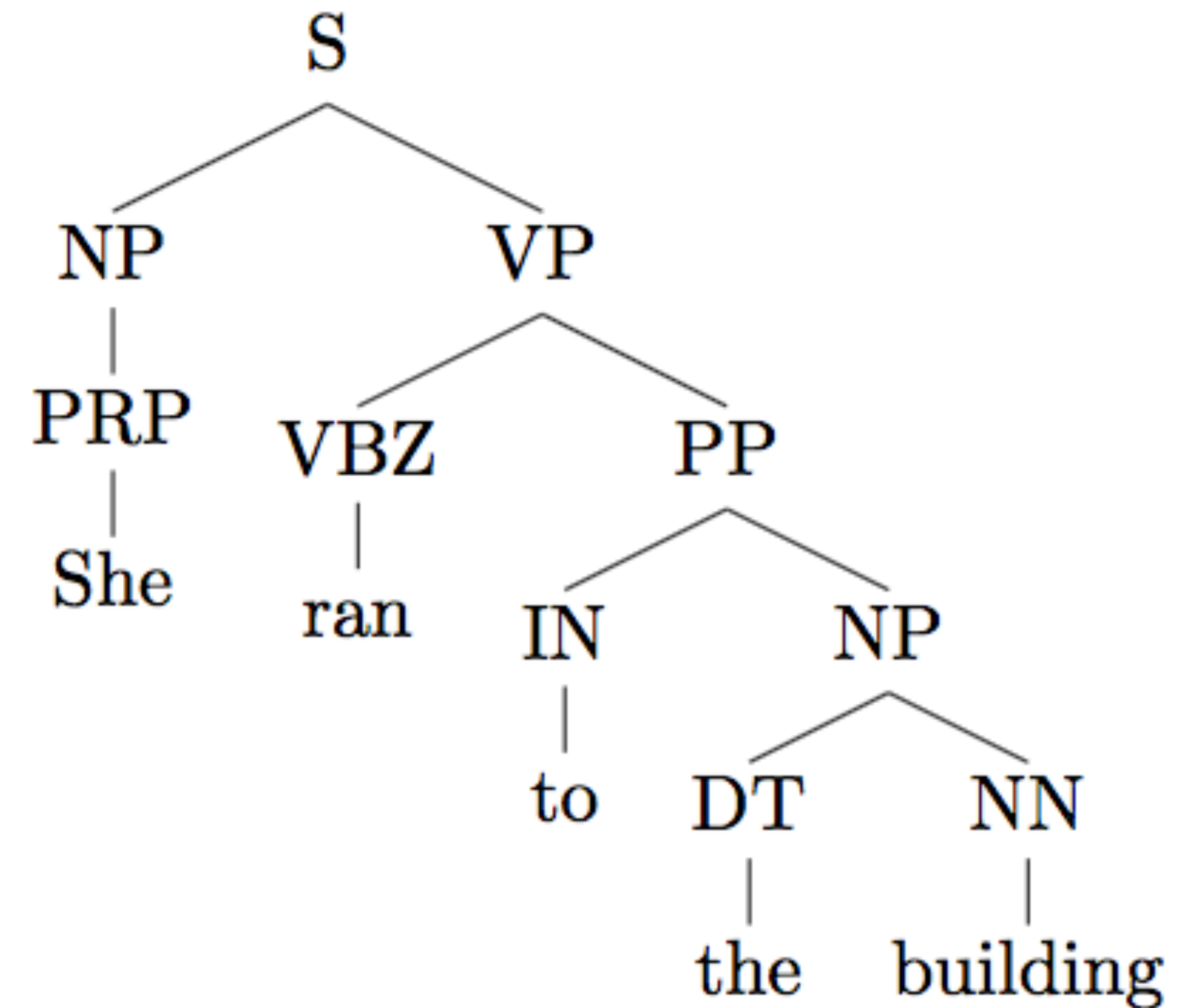
# Constituency Parsing

- ▶ Tree-structured syntactic analyses of sentences

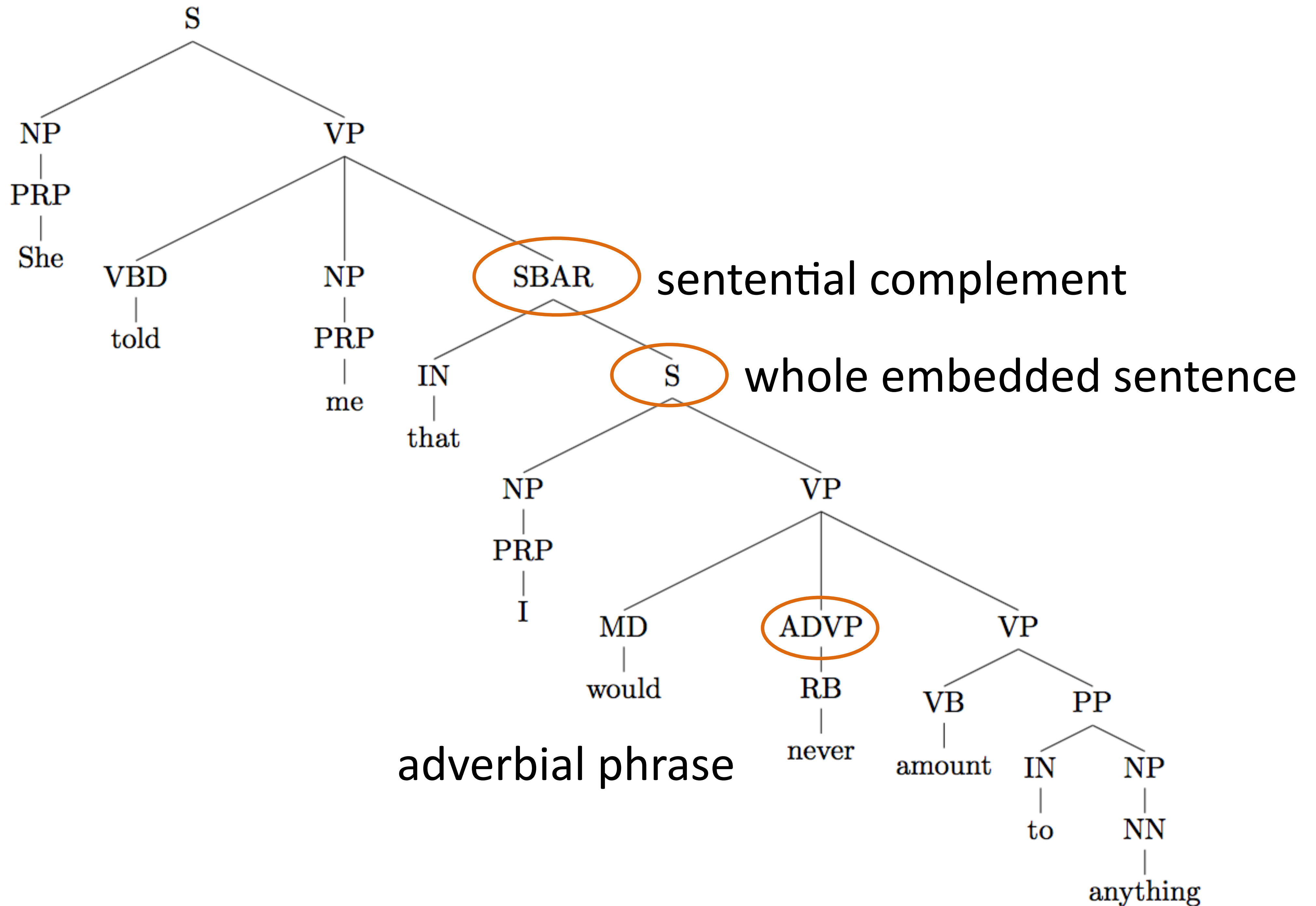
- ▶ Common things: noun phrases, verb phrases, prepositional phrases

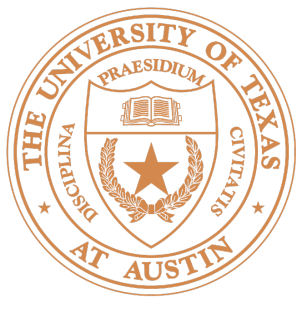
- ▶ Bottom layer is POS tags

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







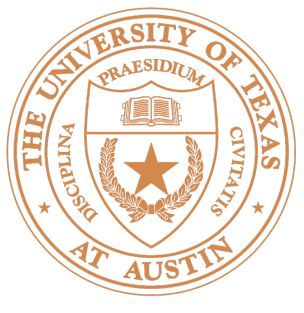


# Constituency Parsing

---

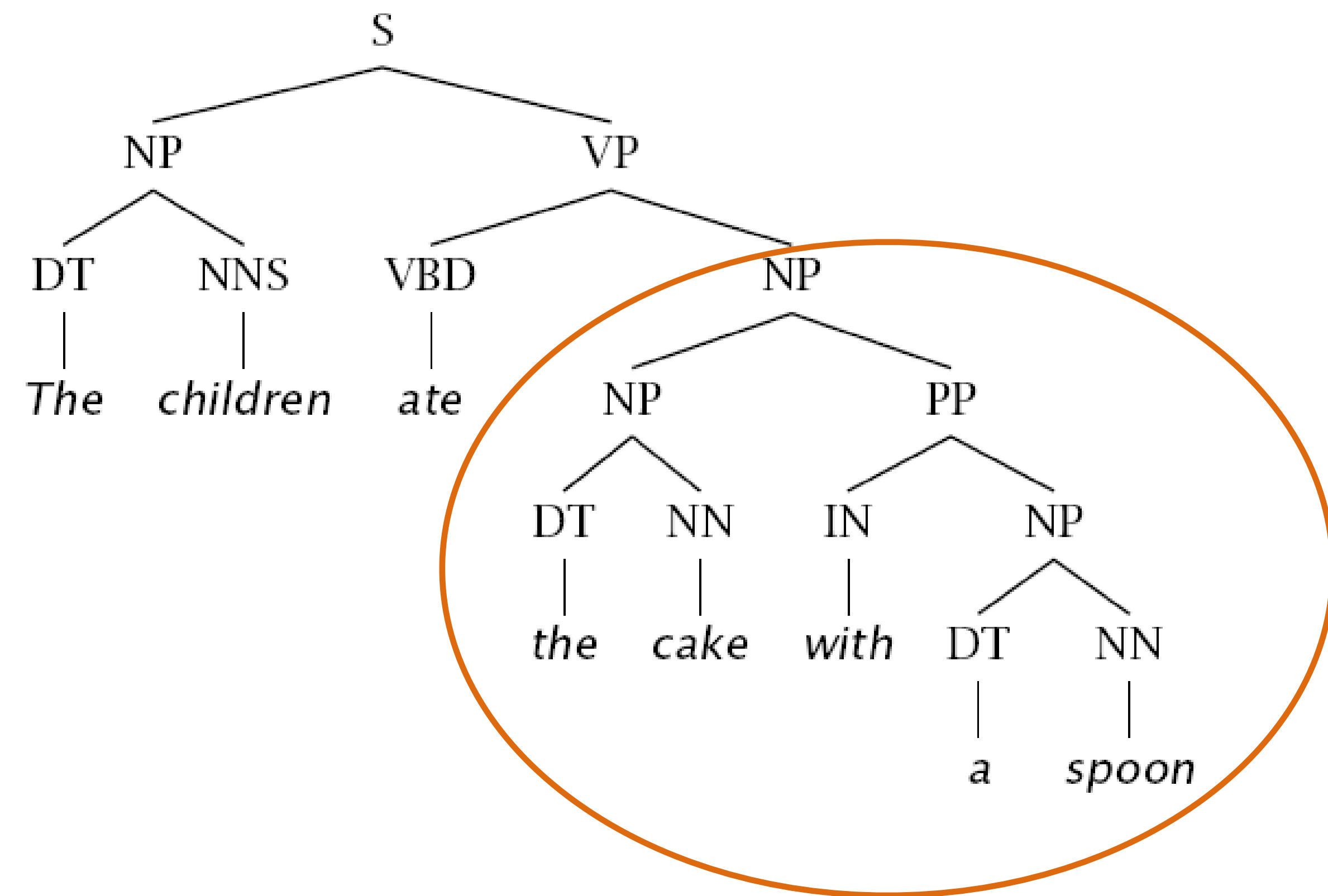
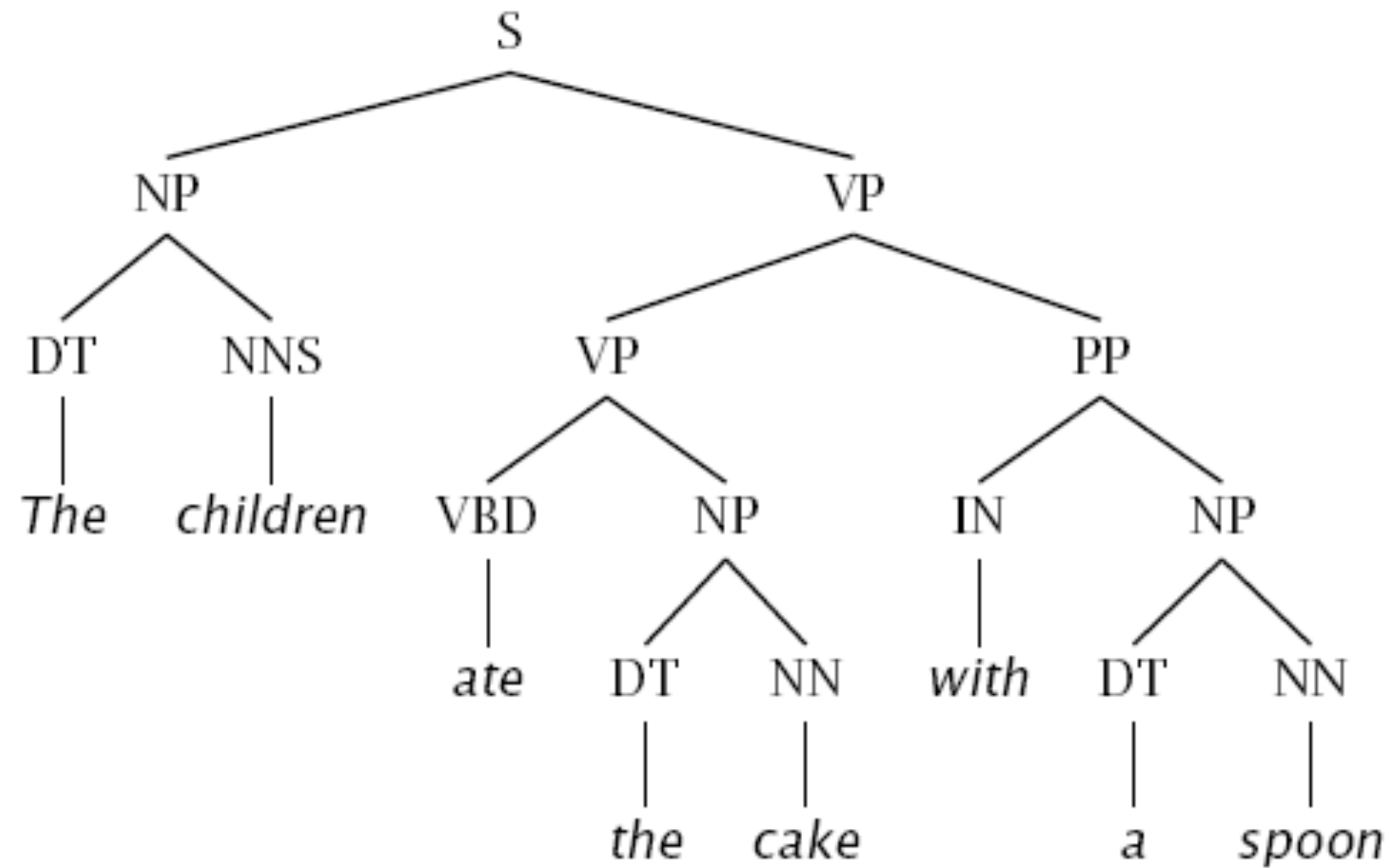
Examples



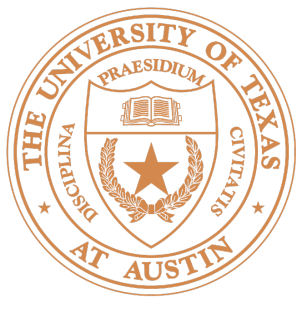


# Challenges

## ► PP attachment

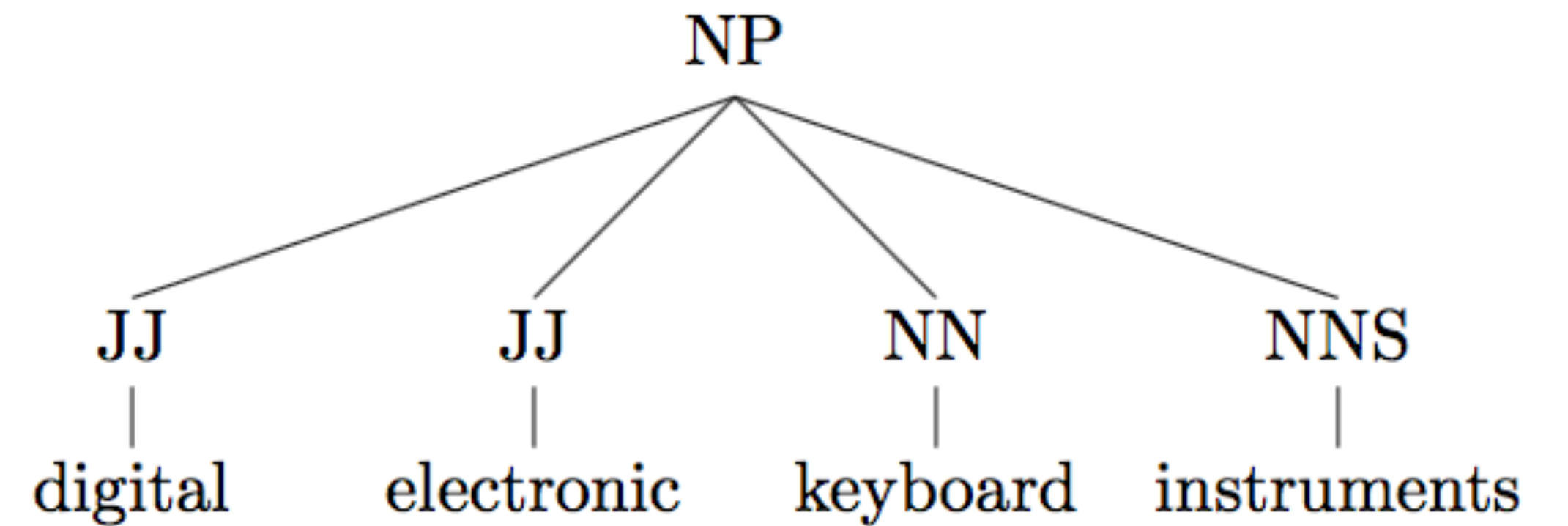
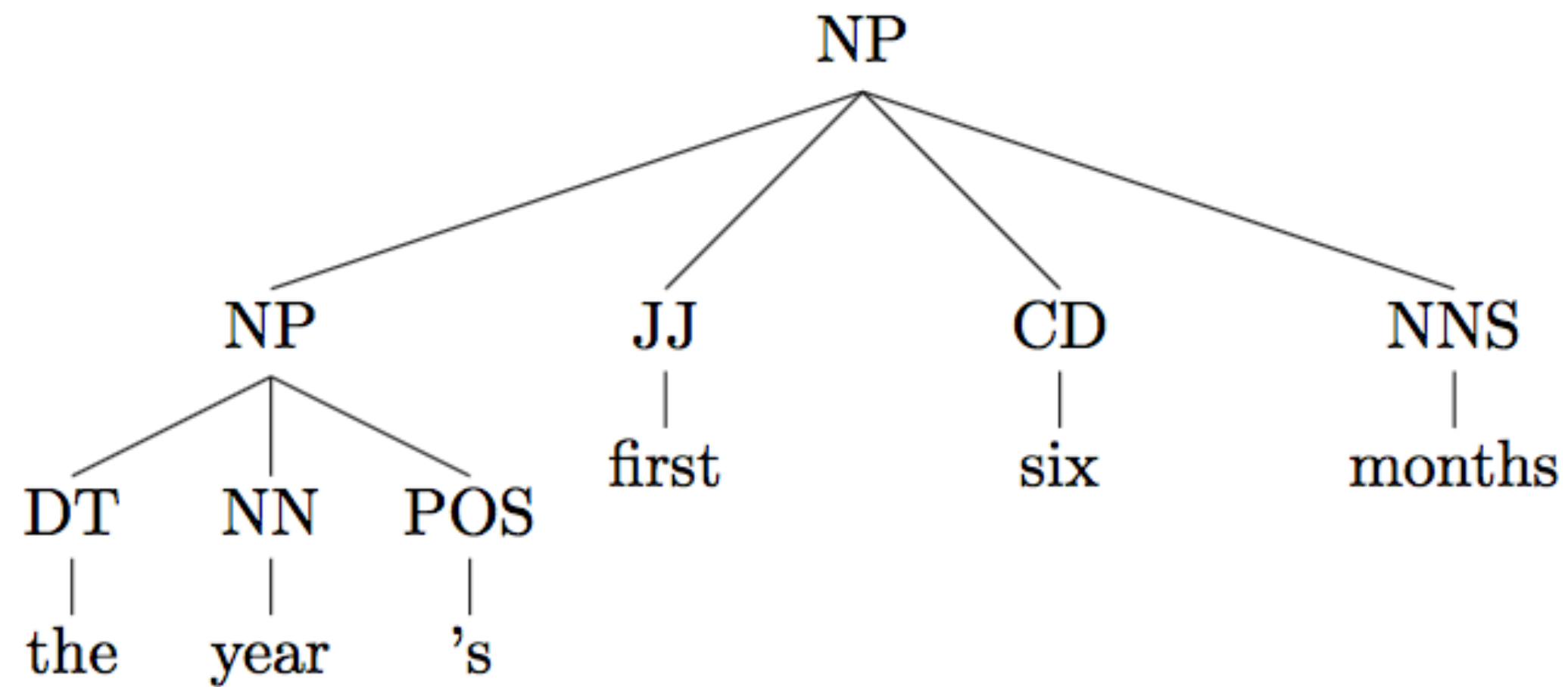


same parse as "the cake with some icing"



# Challenges

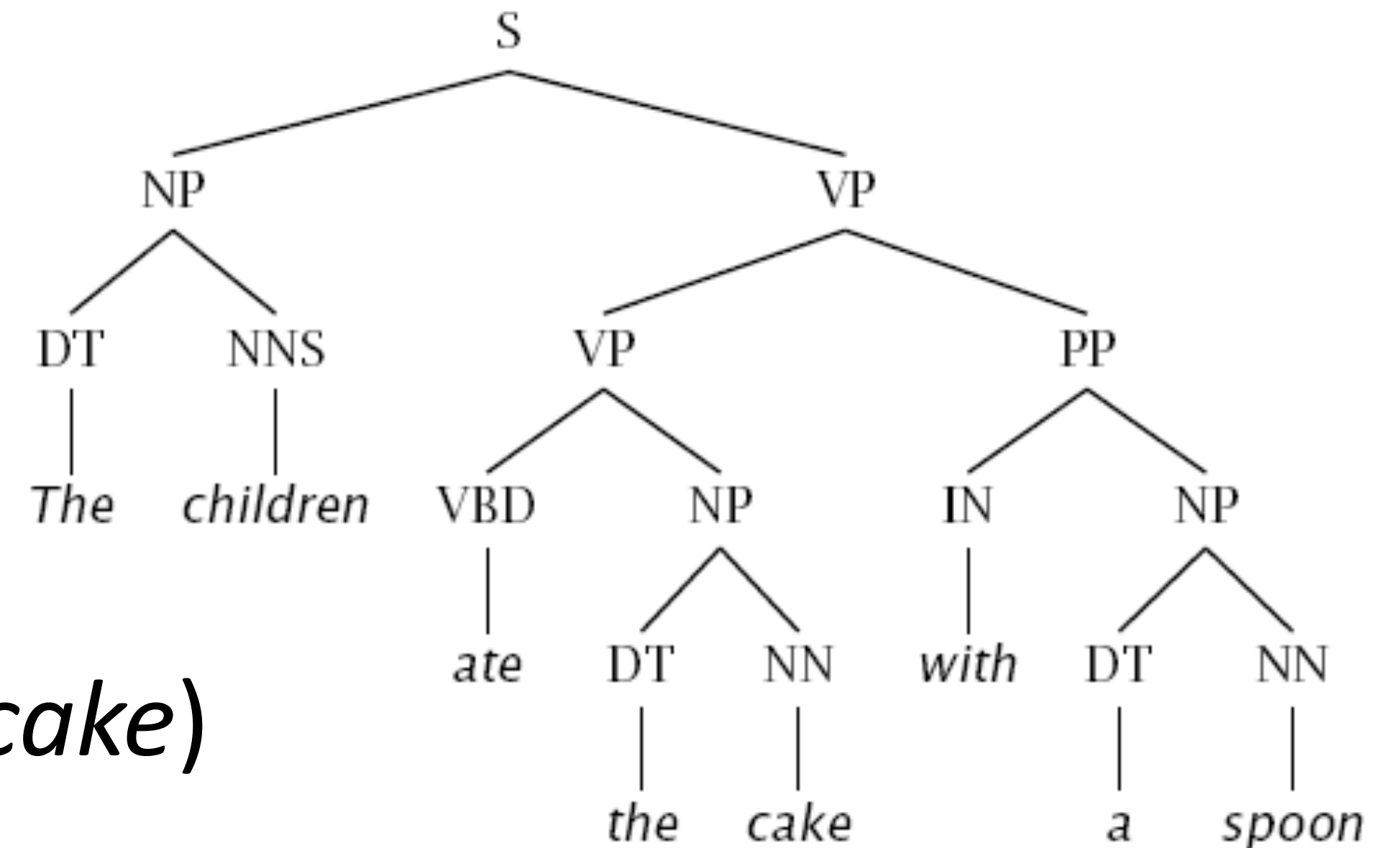
- NP internal structure: tags + depth of analysis



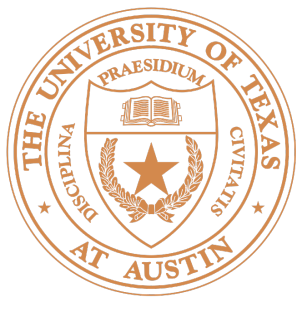


# Constituency

- ▶ 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



# Survey

---

- ▶ 1. The pace of the first few lectures (naive Bayes, logistic regression, perceptron, etc.) was [too fast/too slow/just right]
- ▶ 2. The pace of the last few lectures (tagging, Viterbi, parsing) was [too fast/too slow/just right]
- ▶ 3. The homeworks overall are [too hard/too easy/just right]
- ▶ 4. I would prefer A3 be due on [Friday March 8 / Monday March 11]  
(midterm is on Thursday, March 14)
- ▶ 5. Other comments (likes/dislikes)