

# Paper Presentation & Summary

- 15% of total grade
- Presentation (3%):
  - 12-minute talk (template [here](#)) + 3-minute Q&A
  - You will be punished (-3%) if you sign up but don't present
- Paper summary (12%):
  - For each topic, submitting summaries of at least two papers ( $2 * 6$  classes = 12)
  - Submission: students filling google forms for each paper (due the same week)

# Invited Presentations

- Oct 27: Pre-training & Prompting (Section 2)
  - Lecture by Dr. Pengfei Liu (<http://pfliu.com/>), postdoc at CMU LTI
  - <http://pretrain.nlpedia.ai/>
- TBD: “No Language Left Behind” by Meta/Facebook AI (Section 1)
- TBD: An industrial talk from Deltek

CS678 Advanced Natural Language Processing

# Structured Prediction 1: Sequence Tagging & Syntactic Parsing

Antonis Anastasopoulos & Ziyu Yao



<https://nlp.cs.gmu.edu/course/cs678-fall22>

*Acknowledgement: Many slides are taken from Greg Durrett CS388@UT Austin and  
Antonios Anastasopoulos CS695@GMU*

# Structured Prediction

- Sequence labeling
  - POS Tagging & Hidden Markov Models (HMMs)
  - Named Entity Recognition & Conditional Random Fields (CRFs)
- Syntactic Parsing
  - Dependency parsing
  - Constituency parsing (next lecture)
- Semantic Parsing (next lecture)

# Part-of-Speech Tagging

# Linguistic Structures

- Language has hierarchical structure, can represent with trees



- Understanding syntax fundamentally requires trees — the sentences have the same shallow analysis. But the first step we'll take towards understanding this is understanding parts of speech

NN NNS VBZ NNS  
Teacher strikes idle kids

# POS Tagging

## Open class (lexical) words

### Nouns

#### Proper

*IBM*

*Italy*

#### Common

*cat / cats*

*snow*

### Verbs

#### Main

*see*

*registered*

### Adjectives

*yellow*

### Adverbs

*slowly*

### Numbers

*122,312*

*one*

*... more*

## Closed class (functional)

### Determiners

*the some*

### Conjunctions

*and or*

### Pronouns

*he its*

### Auxiliary

*can*

*had*

### Prepositions

*to with*

### Particles

*off up*

*... more*

# POS Tagging

Tag	Description	Example	Tag	Description	Example	Tag	Description	Example
CC	coord. conj.	<i>and, but, or</i>	NNP	proper noun, sing.	<i>IBM</i>	TO	“to”	<i>to</i>
CD	cardinal number	<i>one, two</i>	NNPS	proper noun, plu.	<i>Carolinas</i>	UH	interjection	<i>ah, oops</i>
DT	determiner	<i>a, the</i>	NNS	noun, plural	<i>llamas</i>	VB	verb base	<i>eat</i>
EX	existential ‘there’	<i>there</i>	PDT	predeterminer	<i>all, both</i>	VBD	verb past tense	<i>ate</i>
FW	foreign word	<i>mea culpa</i>	POS	possessive ending	<i>'s</i>	VBG	verb gerund	<i>eating</i>
IN	preposition/ subordin-conj	<i>of, in, by</i>	PRP	personal pronoun	<i>I, you, he</i>	VBN	verb past participle	<i>eaten</i>
JJ	adjective	<i>yellow</i>	PRP\$	possess. pronoun	<i>your, one's</i>	VBP	verb non-3sg-pr	<i>eat</i>
JJR	comparative adj	<i>bigger</i>	RB	adverb	<i>quickly</i>	VBZ	verb 3sg pres	<i>eats</i>
JJS	superlative adj	<i>wildest</i>	RBR	comparative adv	<i>faster</i>	WDT	wh-determ.	<i>which, that</i>
LS	list item marker	<i>1, 2, One</i>	RBS	superlatv. adv	<i>fastest</i>	WP	wh-pronoun	<i>what, who</i>
MD	modal	<i>can, should</i>	RP	particle	<i>up, off</i>	WP\$	wh-possess.	<i>whose</i>
NN	sing or mass noun	<i>llama</i>	SYM	symbol	<i>+, %, &amp;</i>	WRB	wh-adverb	<i>how, where</i>

**Figure 8.2** Penn Treebank part-of-speech tags.

# POS Tagging

*Ghana's ambassador should have set up the big meeting in DC yesterday.*

NNP POS NN MD VB VBN RP DT JJ NN IN NNP NN .

Tag	Description	Example	Tag	Description	Example	Tag	Description	Example
CC	coord. conj.	<i>and, but, or</i>	NNP	proper noun, sing.	<i>IBM</i>	TO	"to"	<i>to</i>
CD	cardinal number	<i>one, two</i>	NNPS	proper noun, plu.	<i>Carolinas</i>	UH	interjection	<i>ah, oops</i>
DT	determiner	<i>a, the</i>	NNS	noun, plural	<i>llamas</i>	VB	verb base	<i>eat</i>
EX	existential 'there'	<i>there</i>	PDT	predeterminer	<i>all, both</i>	VBD	verb past tense	<i>ate</i>
FW	foreign word	<i>mea culpa</i>	POS	possessive ending	's	VBG	verb gerund	<i>eating</i>
IN	preposition/ subordin-conj	<i>of, in, by</i>	PRP	personal pronoun	<i>I, you, he</i>	VBN	verb past partici- ple	<i>eaten</i>
JJ	adjective	<i>yellow</i>	PRP\$	possess. pronoun	<i>your, one's</i>	VBP	verb non-3sg-pr	<i>eat</i>
JJR	comparative adj	<i>bigger</i>	RB	adverb	<i>quickly</i>	VBZ	verb 3sg pres	<i>eats</i>
JJS	superlative adj	<i>wildest</i>	RBR	comparative adv	<i>faster</i>	WDT	wh-determ.	<i>which, that</i>
LS	list item marker	<i>1, 2, One</i>	RBS	superlatv. adv	<i>fastest</i>	WP	wh-pronoun	<i>what, who</i>
MD	modal	<i>can, should</i>	RP	particle	<i>up, off</i>	WP\$	wh-possess.	<i>whose</i>
NN	sing or mass noun	<i>llama</i>	SYM	symbol	<i>+, %, &amp;</i>	WRB	wh-adverb	<i>how, where</i>

Figure 8.2 Penn Treebank part-of-speech tags.

# POS Tagging

VBD            VB  
VBN **VBZ**    VBP    VBZ  
**NNP** NNS    **NN**    **NNS** CD **NN**  
Fed raises interest rates 0.5 percent



VBD            VB  
VBN VBZ    **VBP**    **VBZ**  
**NNP** **NNS**    **NN**    **NNS** CD **NN**  
Fed raises interest rates 0.5 percent



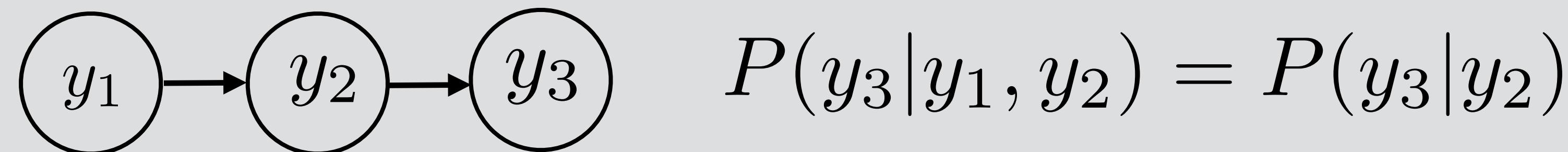
- Other paths are also plausible but even more semantically weird...
- What governs the correct choice? **Word + context**
  - Word identity: most words have  $<=2$  tags, many have one (percent, the)
  - Context: nouns start sentences, nouns follow verbs, etc.

# What is this good for?

- Text-to-speech: record, lead
- Preprocessing step for syntactic parsers or other tasks
- (Very) shallow information extraction

# Hidden Markov Models

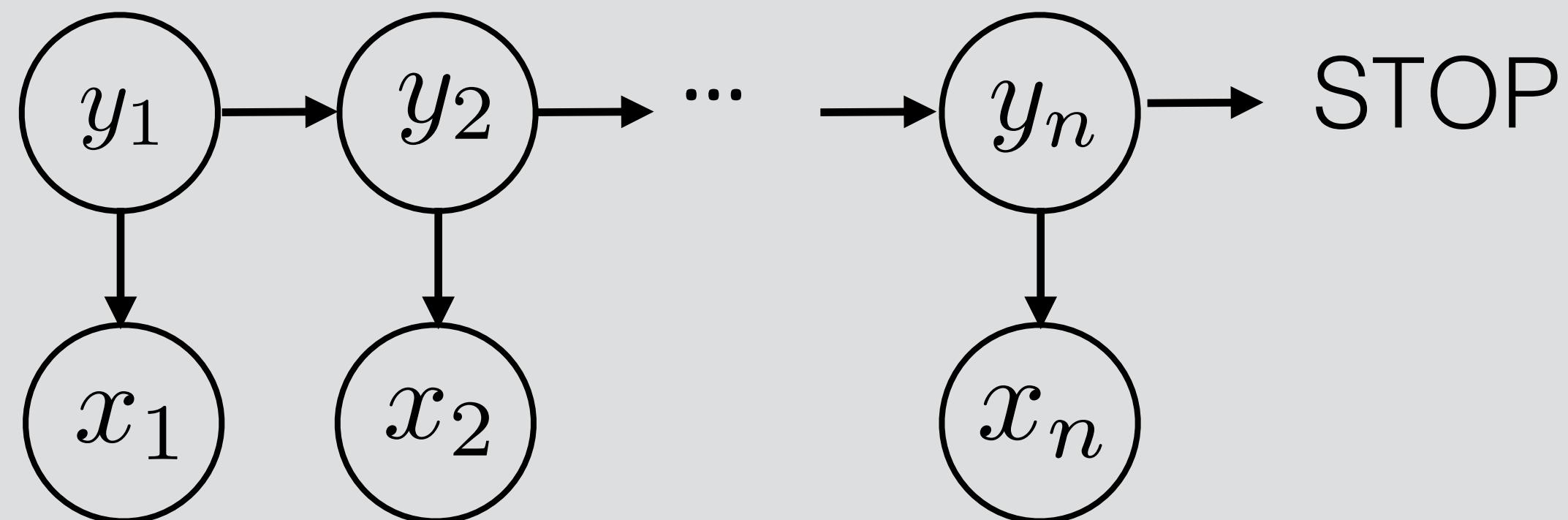
- Input  $\mathbf{x} = (x_1, \dots, x_n)$       Output  $\mathbf{y} = (y_1, \dots, y_n)$
- Model the sequence of tags  $\mathbf{y}$  over words  $\mathbf{x}$  as a Markov process
- Markov property: future is conditionally independent of the past given the present



- If  $\mathbf{y}$  are tags, this roughly corresponds to assuming that the next tag only depends on the current tag, not anything before

# Hidden Markov Models

- Input  $\mathbf{x} = (x_1, \dots, x_n)$       Output  $\mathbf{y} = (y_1, \dots, y_n)$   $y_i \in T$  = set of possible tags  
(including STOP);  
 $x_i \in V$  = vocab of words



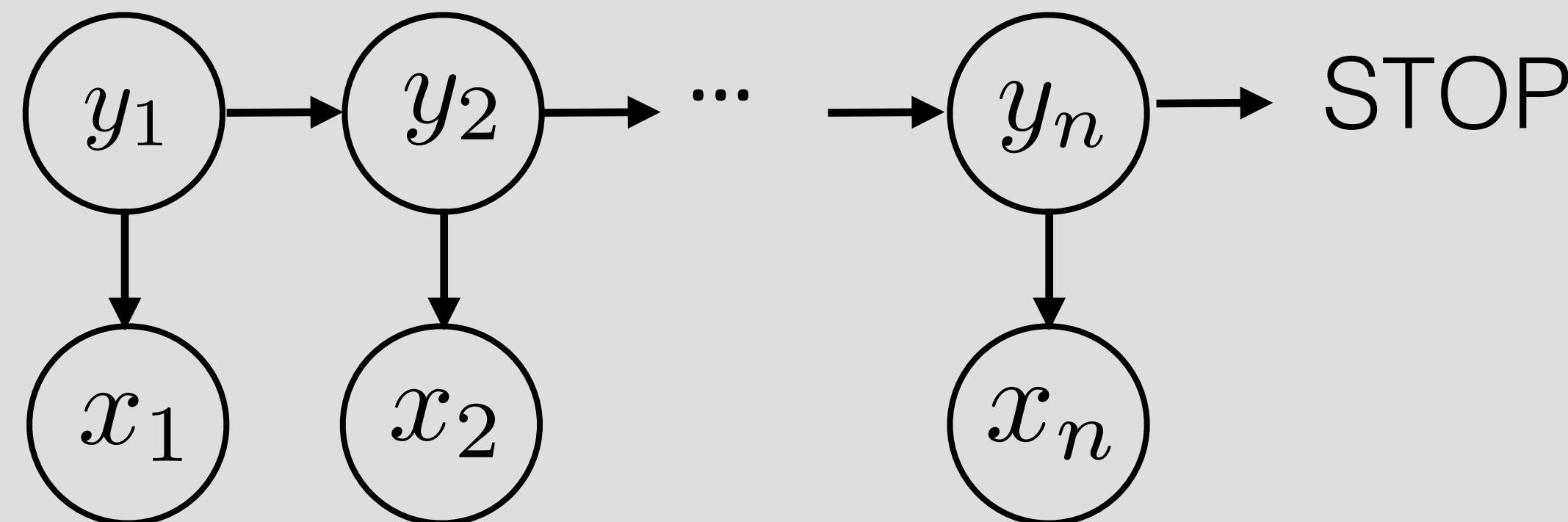
$$P(\mathbf{y}, \mathbf{x}) = P(y_1) \underbrace{\prod_{i=2}^n P(y_i | y_{i-1})}_{\text{Initial distribution}} \underbrace{\prod_{i=1}^n P(x_i | y_i)}_{\text{Transition probabilities}} \underbrace{\quad}_{\text{Emission probabilities}}$$

Observation ( $x_i$ ) depends  
only on current state ( $y_i$ )

# HMM: Parameters

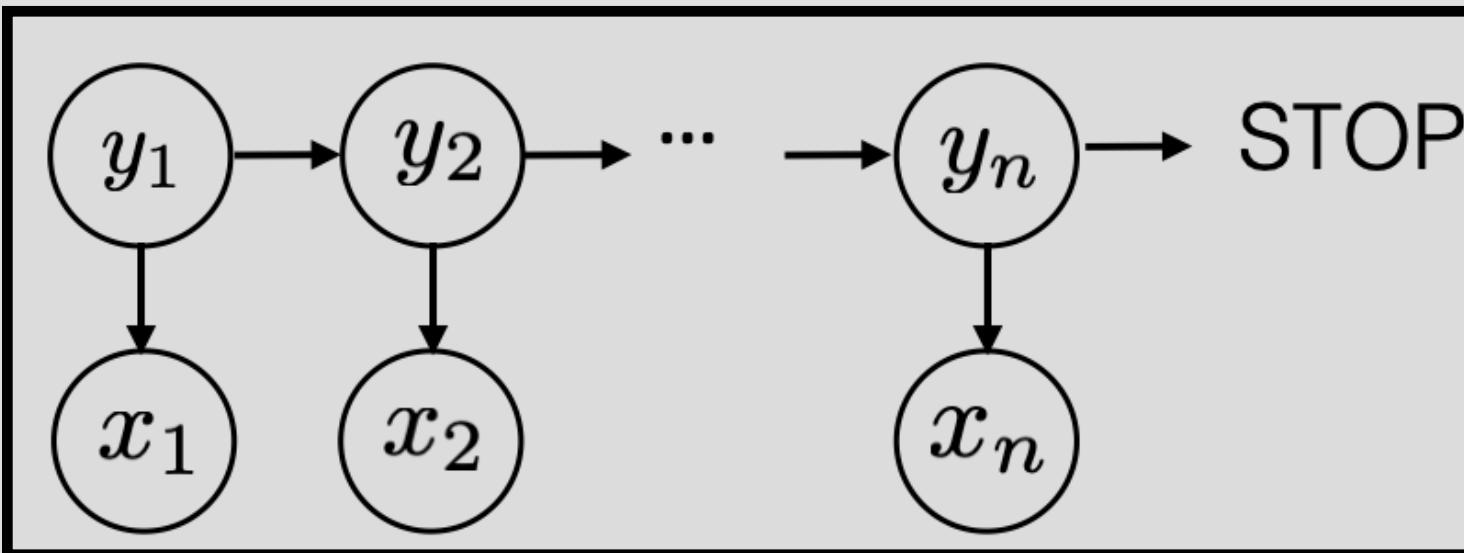
- Input  $\mathbf{x} = (x_1, \dots, x_n)$       Output  $\mathbf{y} = (y_1, \dots, y_n)$

$y_i \in T$  = set of possible tags (including STOP);  
 $x_i \in V$  = vocab of words



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

- Initial distribution:  $|T| \times 1$  vector (distribution over initial states)
- Emission distribution:  $|T| \times |V|$  matrix (distribution over words per tag)
- Transition distribution:  $|T| \times |T|$  matrix (distribution over next tags per tag)



# HMMs Example

Tags = {N, V, STOP}

Vocabulary = {they, can, fish}

Initial		Transition			Emission		
		$y_i$			$x_i$		
		N	V	STOP			
$y_1$	N	1.0					
	V	0					
	STOP	0					
			$y_{i-1}$	N	1/5	3/5	1/5
				V	1/5	1/5	3/5
			$y_i$	they	can	fish	
				N	1	0	0
				V	0	1/2	1/2

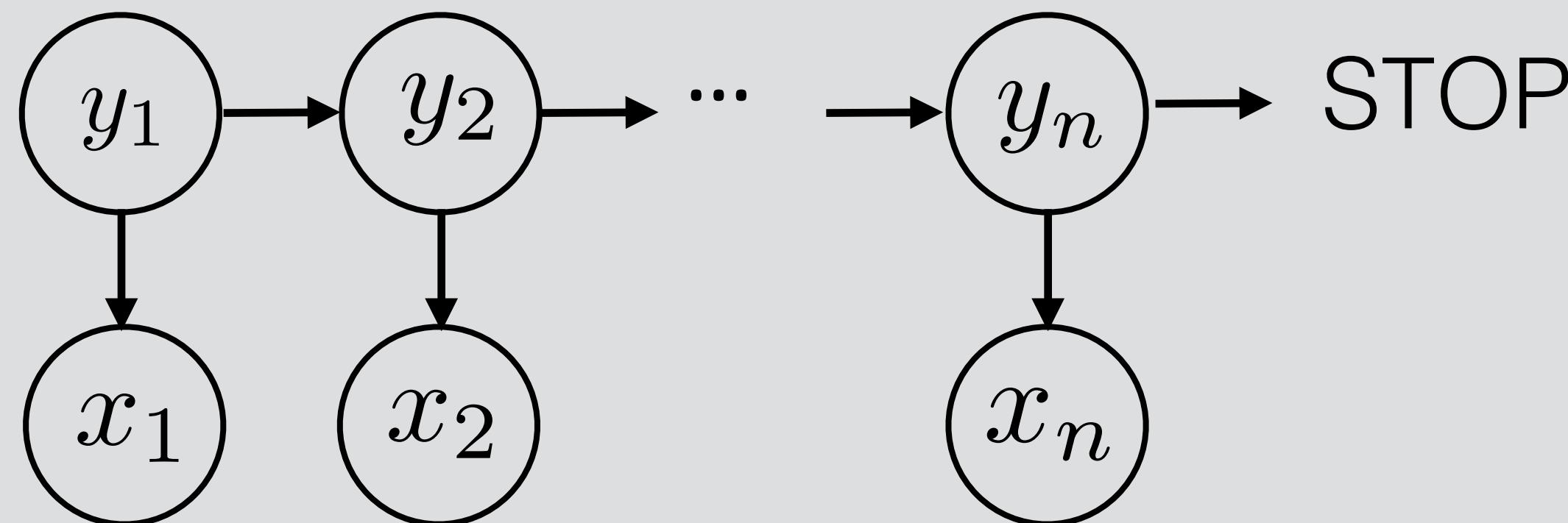
How are these probabilities learned?

# Training HMMs

- Maximum Likelihood Estimation, essentially read counts + normalize
- Transitions: Count up all pairs  $(y_i, y_{i+1})$  in the training data
  - Count up occurrences of what tag  $T$  can transition to
  - Normalize to get a distribution for  $P(\text{ next tag} \mid T)$
  - Need to *smooth* this distribution, won't discuss here
- Emissions: similar scheme, but trickier smoothing!
- Initials: even easier, just count up the currencies of possible first-position tags + doing smoothing

# Inference in HMMs

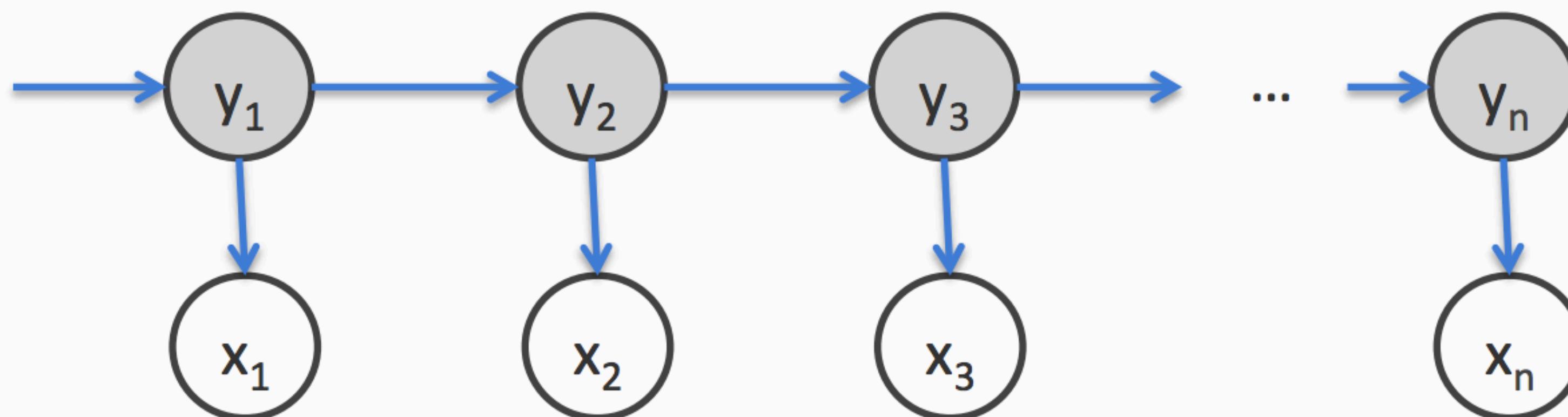
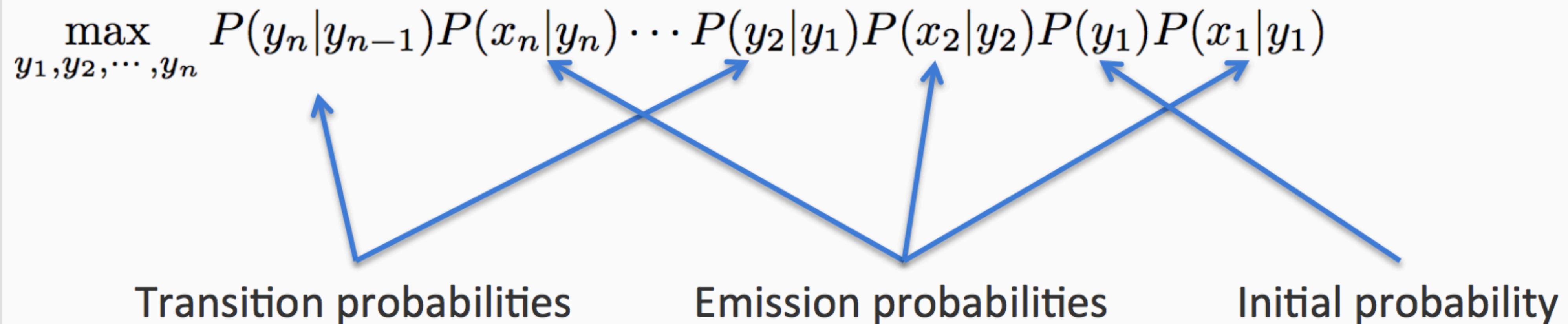
- Input  $\mathbf{x} = (x_1, \dots, x_n)$       Output  $\mathbf{y} = (y_1, \dots, y_n)$



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

- Inference problem:  $\text{argmax}_{\mathbf{y}} P(\mathbf{y}|\mathbf{x}) = \text{argmax}_{\mathbf{y}} \frac{P(\mathbf{y}, \mathbf{x})}{P(\mathbf{x})}$
- Exponentially many possible  $\mathbf{y}$  here!
- Solution: dynamic programming (possible because of **Markov structure!**)

$$P(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1}|y_i) \prod_{i=1}^n P(x_i|y_i)$$

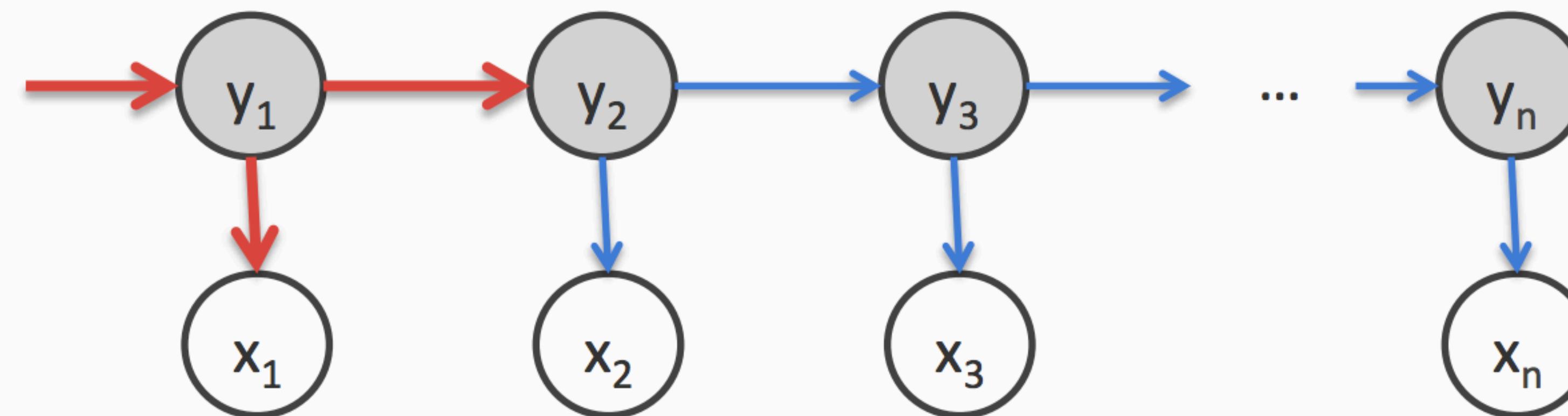


slide credit: Vivek Srikumar

$$P(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1}|y_i) \prod_{i=1}^n P(x_i|y_i)$$

$$\begin{aligned} & \max_{y_1, y_2, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots P(y_2|y_1)P(x_2|y_2)P(y_1)P(x_1|y_1) \\ &= \max_{y_2, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots \max_{y_1} P(y_2|y_1)P(x_2|y_2)P(y_1)P(x_1|y_1) \end{aligned}$$

The only terms that depend on  $y_1$



slide credit: Vivek Srikumar

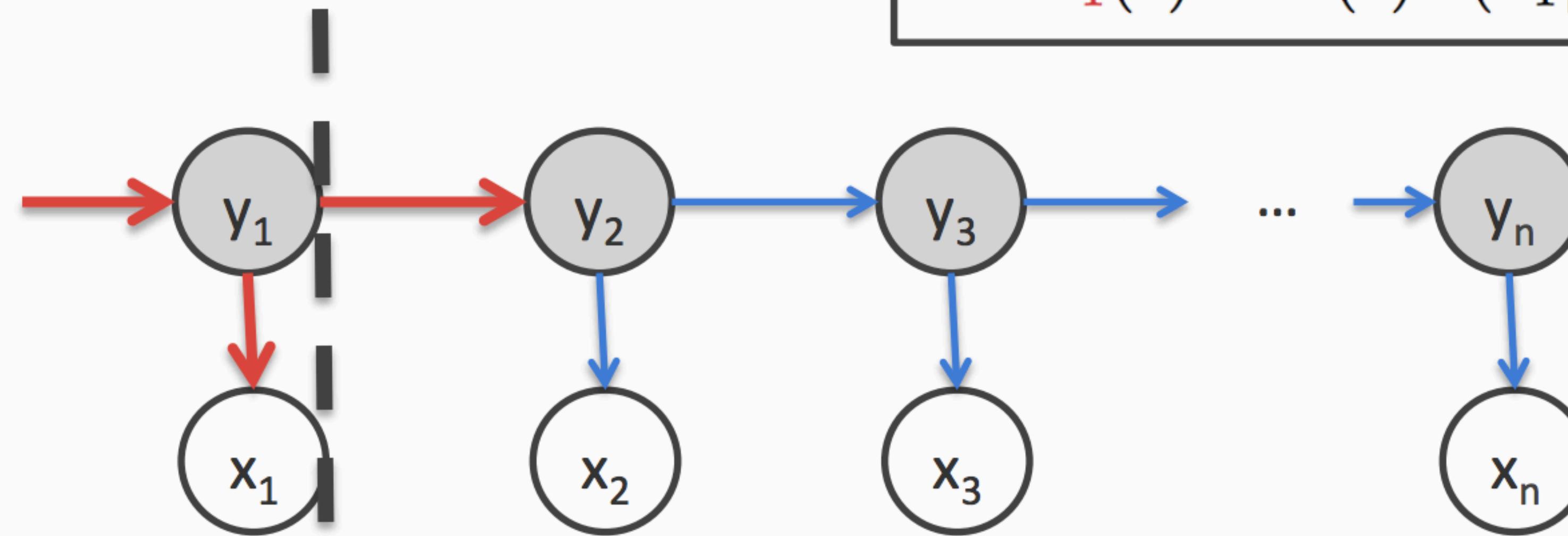
$$P(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1}|y_i) \prod_{i=1}^n P(x_i|y_i)$$

$$\begin{aligned} & \max_{y_1, y_2, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots P(y_2|y_1)P(x_2|y_2)P(y_1)P(x_1|y_1) \\ &= \max_{y_2, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots \max_{y_1} P(y_2|y_1)P(x_2|y_2)P(y_1)P(x_1|y_1) \\ &= \max_{y_2, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots \max_{y_1} P(y_2|y_1)P(x_2|y_2) \text{score}_1(y_1) \end{aligned}$$

Abstract away the score for all decisions till here into **score**

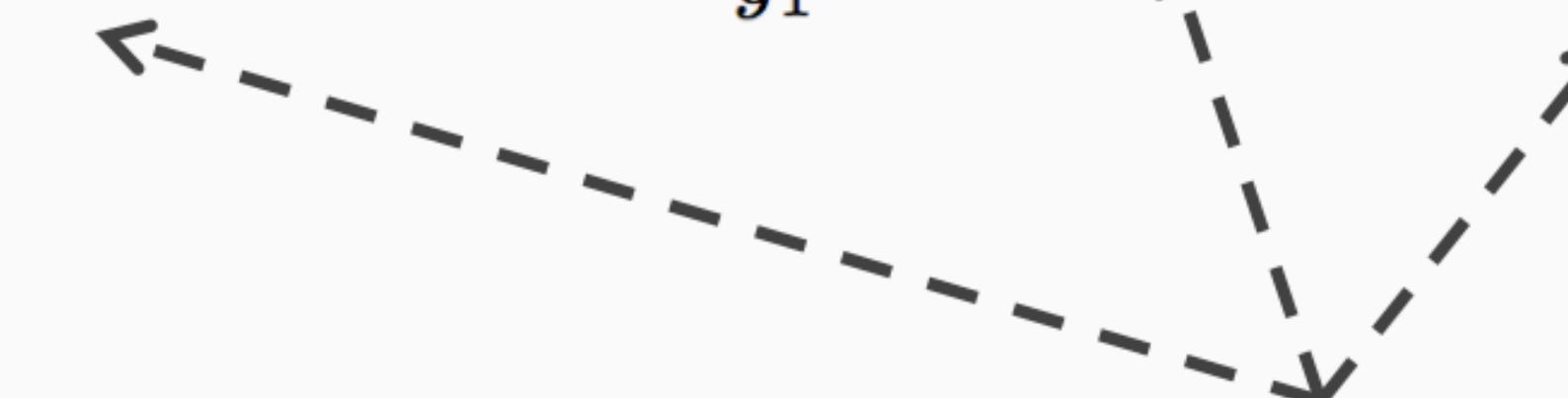
- ▶ Best (partial) score for a sequence ending in state  $s$

$$\text{score}_1(s) = P(s)P(x_1|s)$$

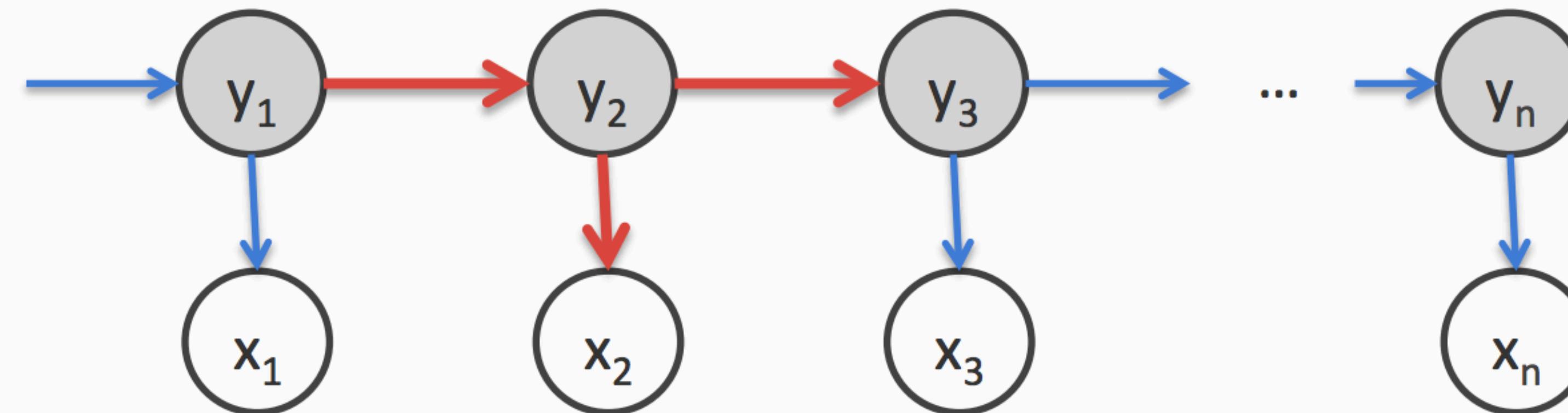


$$P(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1}|y_i) \prod_{i=1}^n P(x_i|y_i)$$

$$\begin{aligned} & \max_{y_1, y_2, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots P(y_2|y_1)P(x_2|y_2)P(y_1)P(x_1|y_1) \\ &= \max_{y_2, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots \max_{y_1} P(y_2|y_1)P(x_2|y_2)P(y_1)P(x_1|y_1) \\ &= \max_{y_2, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots \max_{y_1} P(y_2|y_1)P(x_2|y_2)\text{score}_1(y_1) \\ &= \max_{y_3, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots \max_{y_2} P(y_3|y_2)P(x_3|y_3) \max_{y_1} P(y_2|y_1)P(x_2|y_2)\text{score}_1(y_1) \end{aligned}$$



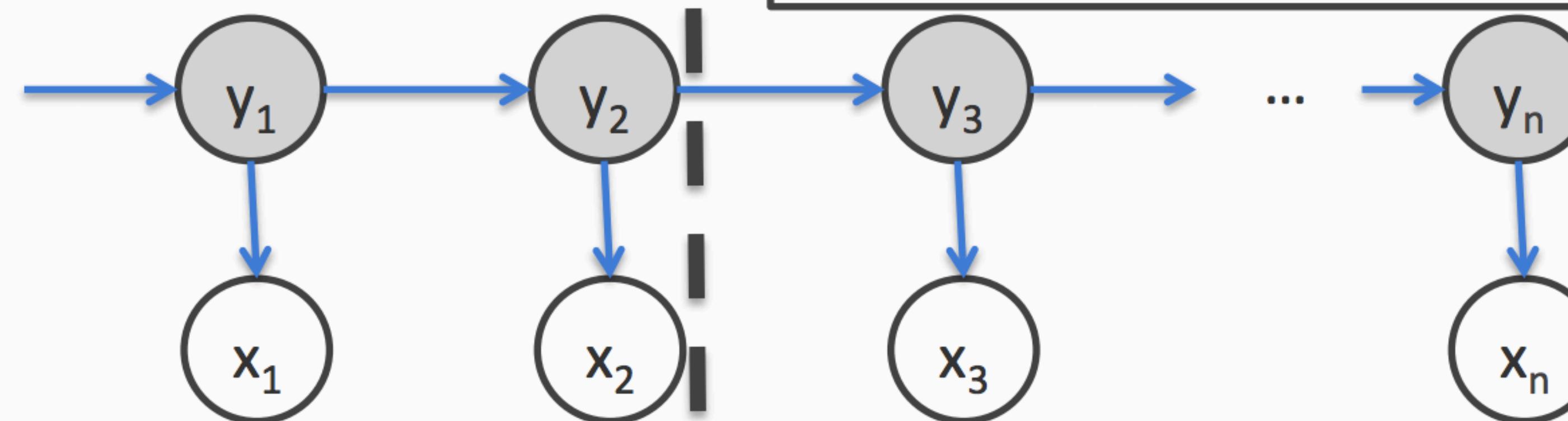
Only terms that depend on  $y_2$



$$P(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1}|y_i) \prod_{i=1}^n P(x_i|y_i)$$

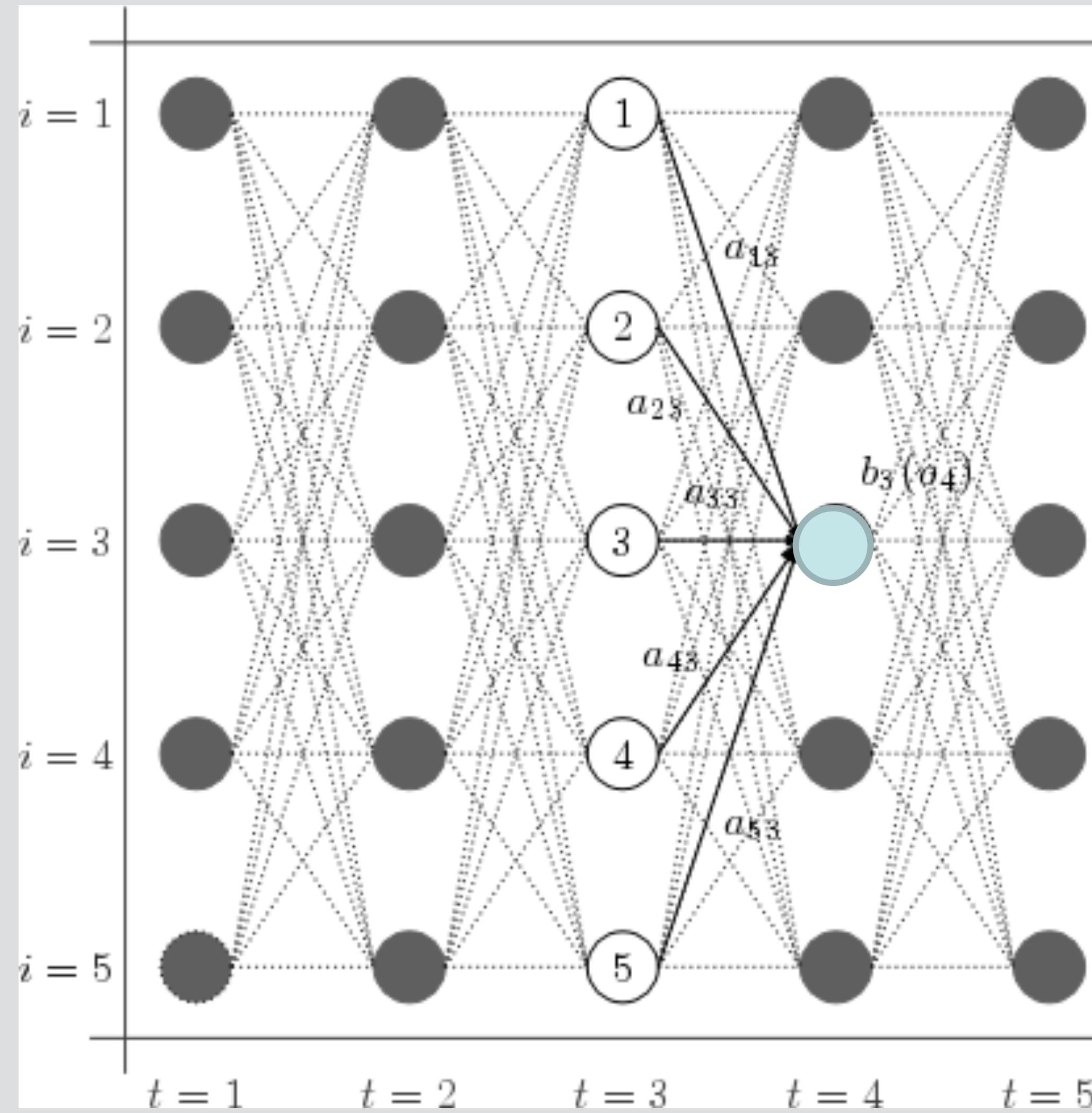
$$\begin{aligned} & \max_{y_1, y_2, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots P(y_2|y_1)P(x_2|y_2)P(y_1)P(x_1|y_1) \\ &= \max_{y_2, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots \max_{y_1} P(y_2|y_1)P(x_2|y_2)P(y_1)P(x_1|y_1) \\ &= \max_{y_2, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots \max_{y_1} P(y_2|y_1)P(x_2|y_2) \text{score}_1(y_1) \\ &= \max_{y_3, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots \max_{y_2} P(y_3|y_2)P(x_3|y_3) \max_{y_1} P(y_2|y_1)P(x_2|y_2) \text{score}_1(y_1) \\ &= \max_{y_3, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots \max_{y_2} P(y_3|y_2)P(x_3|y_3) \text{score}_2(y_2) \end{aligned}$$

$\boxed{\text{score}_i(s) = \max_{y_{i-1}} P(s|y_{i-1})P(x_i|s)\text{score}_{i-1}(y_{i-1})}$



Abstract away the score for all decisions till here into **score**

# Viterbi Algorithm



- “Think about” all possible immediate prior state values. Everything before that has already been accounted for by earlier stages.

$$P(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1}|y_i) \prod_{i=1}^n P(x_i|y_i)$$

$$\max_{y_1, y_2, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots P(y_2|y_1)P(x_2|y_2)P(y_1)P(x_1|y_1)$$

$$= \max_{y_2, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots \max_{y_1} P(y_2|y_1)P(x_2|y_2)P(y_1)P(x_1|y_1)$$

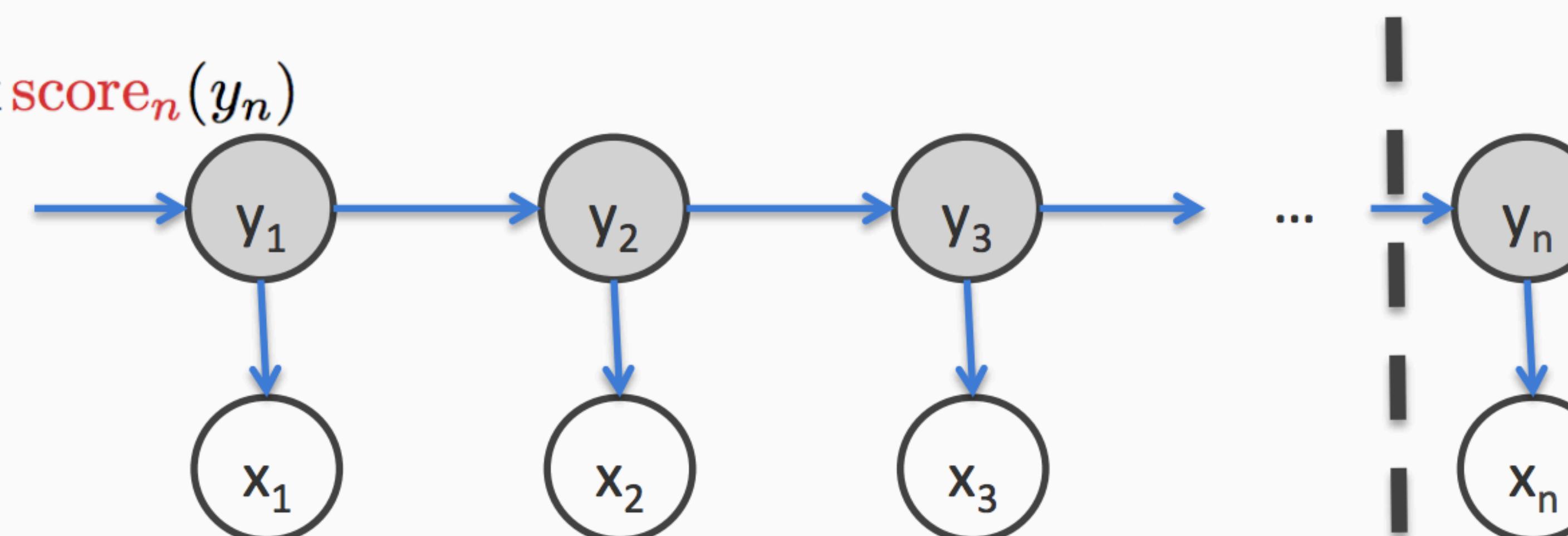
$$= \max_{y_2, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots \max_{y_1} P(y_2|y_1)P(x_2|y_2)\text{score}_1(y_1)$$

$$= \max_{y_3, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots \max_{y_2} P(y_3|y_2)P(x_3|y_3) \max_{y_1} P(y_2|y_1)P(x_2|y_2)\text{score}_1(y_1)$$

$$= \max_{y_3, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots \max_{y_2} P(y_3|y_2)P(x_3|y_3)\text{score}_2(y_2)$$

:

$$= \max_{y_n} \text{score}_n(y_n)$$



Abstract away the score for all decisions till here into **score**

slide credit: Vivek Srikumar

$$P(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1}|y_i) \prod_{i=1}^n P(x_i|y_i)$$

$$\begin{aligned} & \max_{y_1, y_2, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots P(y_2|y_1)P(x_2|y_2)P(y_1)P(x_1|y_1) \\ &= \max_{y_2, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots \max_{y_1} P(y_2|y_1)P(x_2|y_2)P(y_1)P(x_1|y_1) \\ &= \max_{y_2, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots \max_{y_1} P(y_2|y_1)P(x_2|y_2)\text{score}_1(y_1) \\ &= \max_{y_3, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots \max_{y_2} P(y_3|y_2)P(x_3|y_3) \max_{y_1} P(y_2|y_1)P(x_2|y_2)\text{score}_1(y_1) \\ &= \max_{y_3, \dots, y_n} P(y_n|y_{n-1})P(x_n|y_n) \cdots \max_{y_2} P(y_3|y_2)P(x_3|y_3)\text{score}_2(y_2) \\ &\vdots \\ &= \max_{y_n} \text{score}_n(y_n) \end{aligned}$$

$$\text{score}_1(s) = P(s)P(x_1|s)$$

$$\text{score}_i(s) = \max_{y_{i-1}} P(s|y_{i-1})P(x_i|s)\text{score}_{i-1}(y_{i-1})$$

1. **Initial:** For each state  $s$ , calculate

$$\text{score}_1(s) = P(s)P(x_1|s) = \pi_s B_{x_1,s}$$

2. **Recurrence:** For  $i = 2$  to  $n$ , for every state  $s$ , calculate

$$\begin{aligned}\text{score}_i(s) &= \max_{y_{i-1}} P(s|y_{i-1})P(x_i|s)\text{score}_{i-1}(y_{i-1}) \\ &= \max_{y_{i-1}} A_{y_{i-1},s}B_{s,x_i}\text{score}_{i-1}(y_{i-1})\end{aligned}$$

3. **Final state:** calculate

$$\max_{\mathbf{y}} P(\mathbf{y}, \mathbf{x}|\pi, A, B) = \max_s \text{score}_n(s)$$

$\pi$ : Initial probabilities

$A$ : Transitions

$B$ : Emissions

This only calculates the max. To get final answer (*argmax*),

- keep track of which state corresponds to the max at each step
- build the answer using these back pointers

# HMM POS Tagging

- Baseline: assign each word its most frequent tag: ~90% accuracy
- HMM-based POS Taggers:
  - Trigram HMM: ~95% accuracy / 55% on unknown words
  - TnT tagger (Brants 1998, tuned HMM): 96.2% accuracy / 86.0% on unks
- State-of-the-art (BiLSTM-CRFs): 97.5% / 89%+

# Named Entity Recognition

# 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

- BIO tag set: begin, inside, outside
- Sequence of tags — should we use an HMM?
- Why might an HMM not do so well here?
  - Lots of O's
  - Insufficient features/capacity with multinomials (especially for UNKs)

# Conditional Random Fields

- Flexible *discriminative model* for tagging tasks that can use arbitrary features of the input. Similar to logistic regression, but structured

B-PER I-PER

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



**Curr\_word=Barack & Label=B-PER**

**Next\_word=Obama & Label=B-PER**

**Curr\_word\_starts\_with\_capital=True & Label=B-PER**

**Posn\_in\_sentence=1st & Label=B-PER**

...

# Tagging with Logistic Regression

- Logistic regression over each tag individually:

$$P(y_i = y | \mathbf{x}, i) = \frac{\exp(\mathbf{w}^\top \mathbf{f}(y, i, \mathbf{x}))}{\sum_{y' \in \mathcal{Y}} \exp(\mathbf{w}^\top \mathbf{f}(y', i, \mathbf{x}))}$$

Probability of the  $i$ -th word getting assigned tag  $y$  (e.g., B-PER)

# Tagging with Logistic Regression

- Logistic regression over each tag individually:

$$P(y_i = y | \mathbf{x}, i) = \frac{\exp(\mathbf{w}^\top \mathbf{f}(y, i, \mathbf{x}))}{\sum_{y' \in \mathcal{Y}} \exp(\mathbf{w}^\top \mathbf{f}(y', i, \mathbf{x}))}$$

Look like HMM's  
“emission” modeling

- Over all tags:

$$P(\mathbf{y} = \tilde{\mathbf{y}} | \mathbf{x}) = \prod_{i=1}^n P(y_i = \tilde{y}_i | \mathbf{x}, i) = \frac{1}{Z} \exp \left( \sum_{i=1}^n \mathbf{w}^\top \mathbf{f}(\tilde{y}_i, i, \mathbf{x}) \right)$$

- Set Z equal to the product of denominators
- Score of a prediction: sum of weights dot features over each individual predicted tag (this is a simple CRF but not the general form)

# Example: “Emission Features” $f_e$

B-PER    I-PER    O    O

*Barack Obama will travel*

$$\text{feats} = f_e(\text{B-PER}, i=1, x) + f_e(\text{I-PER}, i=2, x) + f_e(\text{O}, i=3, x) + f_e(\text{O}, i=4, x)$$

[CurrWord=*Obama* & label=I-PER, PrevWord=*Barack* & label=I-PER,  
CurrWordIsCapitalized & label=I-PER, ...]

B-PER    B-PER    O    O

*Barack Obama will travel*

$$\text{feats} = f_e(\text{B-PER}, i=1, x) + \textcolor{blue}{f_e(\text{B-PER}, i=2, x)} + f_e(\text{O}, i=3, x) + f_e(\text{O}, i=4, x)$$

# Adding Structure

$$P(\mathbf{y} = \tilde{\mathbf{y}} | \mathbf{x}) = \frac{1}{Z} \exp \left( \sum_{i=1}^n \mathbf{w}^\top \mathbf{f}(\tilde{y}_i, i, \mathbf{x}) \right)$$

- We want to be able to learn that some tags don't follow other tags — want to have features on *tag pairs*

$$P(\mathbf{y} = \tilde{\mathbf{y}} | \mathbf{x}) = \frac{1}{Z} \exp \left( \sum_{i=1}^n \mathbf{w}^\top \mathbf{f}_e(\tilde{y}_i, i, \mathbf{x}) + \sum_{i=2}^n \mathbf{w}^\top \mathbf{f}_t(\tilde{y}_{i-1}, \tilde{y}_i, i, \mathbf{x}) \right)$$

- Score: sum of weights dot  $\mathbf{f\_e}$  features over each predicted tag (“emissions”) plus sum of weights dot  $\mathbf{f\_t}$  features over tag pairs (“transitions”)
- This is a *sequential CRF*

# Example

B-PER I-PER O O

*Barack Obama will travel*

$$\begin{aligned} \text{feats} = & f_e(\text{B-PER}, i=1, x) + f_e(\text{I-PER}, i=2, x) + f_e(\text{O}, i=3, x) + f_e(\text{O}, i=4, x) \\ & + f_t(\text{B-PER}, \text{I-PER}, i=1, x) + f_t(\text{I-PER}, \text{O}, i=2, x) + f_t(\text{O}, \text{O}, i=3, x) \end{aligned}$$

B-PER B-PER O O

*Barack Obama will travel*

$$\begin{aligned} \text{feats} = & f_e(\text{B-PER}, i=1, x) + \color{blue}{f_e(\text{B-PER}, i=2, x)} + f_e(\text{O}, i=3, x) + f_e(\text{O}, i=4, x) \\ & + \color{red}{f_t(\text{B-PER}, \text{B-PER}, i=1, x)} + f_t(\text{B-PER}, \text{O}, i=2, x) + f_t(\text{O}, \text{O}, i=3, x) \end{aligned}$$

- *Obama* can start a new named entity ([emission feats](#) look okay), but we're not likely to have two PER entities in a row ([transition feats](#))

# CRFs Outline

- Model:

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

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

**(Linear feature)**

- Inference
- Learning

# Inference

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

- $\text{argmax}_{\mathbf{y}} P(\mathbf{y}|\mathbf{x})$ : can use Viterbi exactly as in HMM case

$$\max_{y_1, \dots, y_n} e^{\phi_t(y_{n-1}, y_n)} e^{\phi_e(y_n, n, \mathbf{x})} \dots e^{\phi_e(y_2, 2, \mathbf{x})} e^{\phi_t(y_1, y_2)} e^{\phi_e(y_1, 1, \mathbf{x})}$$

- $\exp(\phi_t(y_{i-1}, y_i))$  and  $\exp(\phi_e(y_i, i, \mathbf{x}))$  play the role of the Ps now, use the exact same Viterbi dynamic program

# 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:  $\text{argmax } P(\mathbf{y}|\mathbf{x})$  from Viterbi
- Learning

# Learning

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

- Given a training set of  $(\mathbf{x}, \mathbf{y}^*)$  pairs, maximize  $\mathcal{L}(\mathbf{y}^*, \mathbf{x}) = \log P(\mathbf{y}^*|\mathbf{x})$
- Gradient is analogous to logistic regression: gold feats - expected feats

$$\frac{\partial}{\partial w} \mathcal{L}(\mathbf{y}^*, \mathbf{x}) = \sum_{i=2}^n f_t(y_{i-1}^*, y_i^*) + \sum_{i=1}^n f_e(y_i^*, i, \mathbf{x})$$

intractable!

$$-\mathbb{E}_{\mathbf{y}} \left[ \sum_{i=2}^n f_t(y_{i-1}, y_i) + \sum_{i=1}^n f_e(y_i, i, \mathbf{x}) \right]$$

# Training CRFs

$$\begin{aligned}\frac{\partial}{\partial w} \mathcal{L}(\mathbf{y}^*, \mathbf{x}) = & \sum_{i=2}^n f_t(y_{i-1}^*, y_i^*) + \sum_{i=1}^n f_e(y_i^*, i, \mathbf{x}) \\ & - \mathbb{E}_{\mathbf{y}} \left[ \sum_{i=2}^n f_t(y_{i-1}, y_i) + \sum_{i=1}^n f_e(y_i, i, \mathbf{x}) \right]\end{aligned}$$

- Let's consider the emission feature expectation first

$$\begin{aligned}\mathbb{E}_{\mathbf{y}} \left[ \sum_{i=1}^n f_e(y_i, i, \mathbf{x}) \right] &= \sum_{\mathbf{y} \in \mathcal{Y}} P(\mathbf{y} | \mathbf{x}) \left[ \sum_{i=1}^n f_e(y_i, i, \mathbf{x}) \right] = \sum_{i=1}^n \sum_{\mathbf{y} \in \mathcal{Y}} P(\mathbf{y} | \mathbf{x}) f_e(y_i, i, \mathbf{x}) \\ &= \sum_{i=1}^n \sum_s P(y_i = s | \mathbf{x}) f_e(s, i, \mathbf{x})\end{aligned}$$

How to calculate this marginal probability?

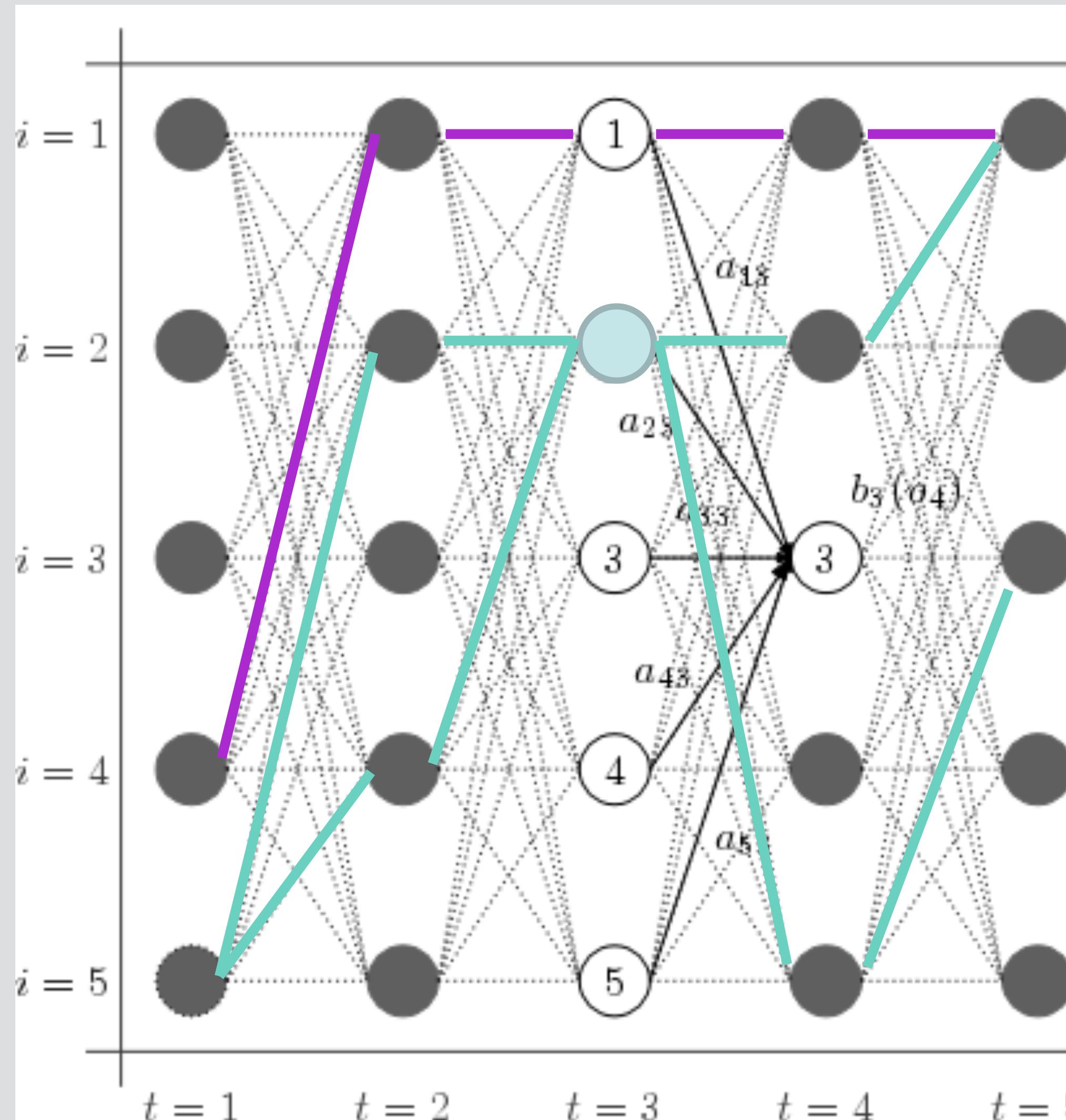
# Forward-Backward Algorithm

- How do we compute these marginals  $P(y_i = s | \mathbf{x})$  ?

$$P(y_i = s | \mathbf{x}) = \sum_{y_1, \dots, y_{i-1}, y_{i+1}, \dots, y_n} P(\mathbf{y} | \mathbf{x})$$

- What did Viterbi compute?  $P(\mathbf{y}_{\max} | \mathbf{x}) = \max_{y_1, \dots, y_n} P(\mathbf{y} | \mathbf{x})$
- Can compute marginals with dynamic programming as well using the forward-backward algorithm

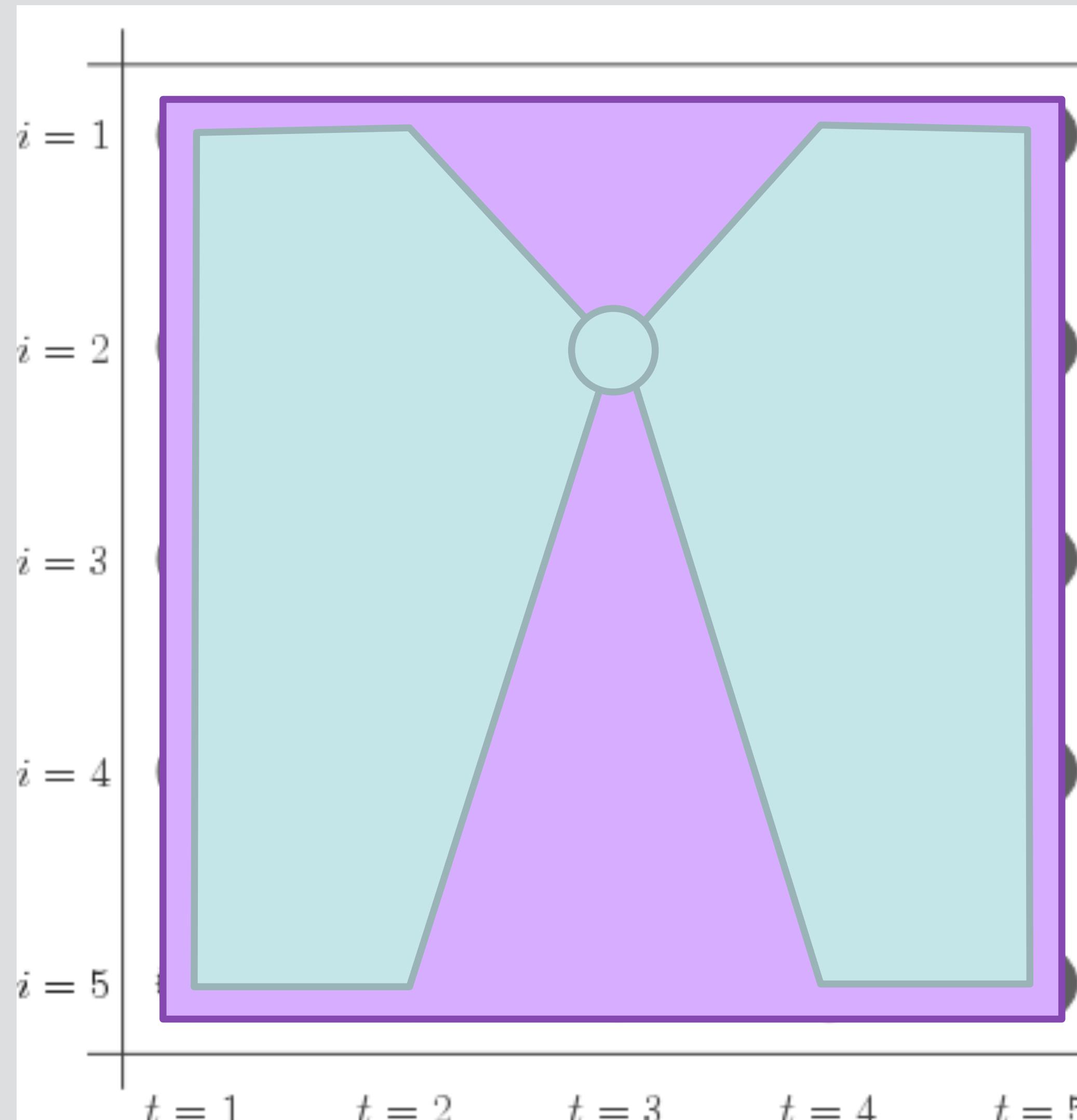
# Forward-Backward Algorithm



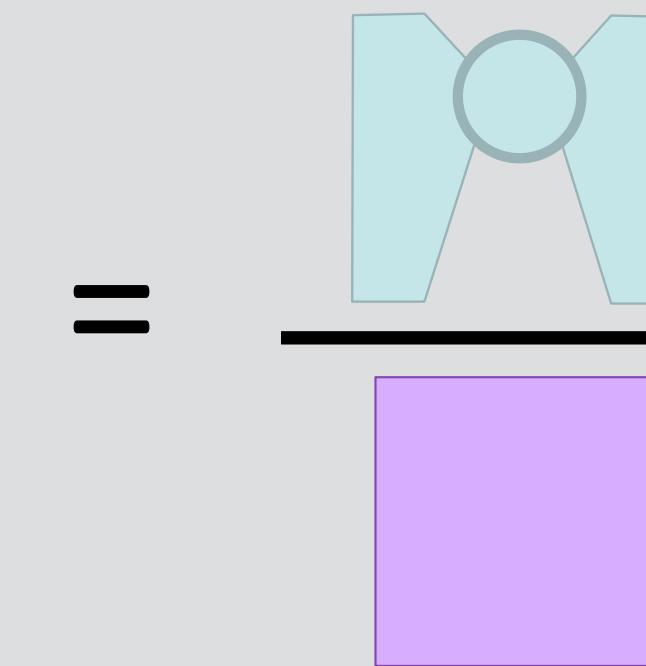
$$P(y_3 = 2 | \mathbf{x}) = \frac{\text{sum of all paths through state 2 at time 3}}{\text{sum of all paths}}$$

“2” is like a candidate NER tag (e.g., B-PER)

# Forward-Backward Algorithm



$$P(y_3 = 2 | \mathbf{x}) = \frac{\text{sum of all paths through state 2 at time 3}}{\text{sum of all paths}}$$



- Easiest and most flexible to do one pass to compute and one to compute

# Training CRFs

- For emission features:

$$\frac{\partial}{\partial w} \mathcal{L}(\mathbf{y}^*, \mathbf{x}) = \sum_{i=1}^n f_e(y_i^*, i, \mathbf{x}) - \sum_{i=1}^n \sum_s P(y_i = s | \mathbf{x}) f_e(s, i, \mathbf{x})$$

- Need compute  $P(y_i = s | \mathbf{x})$
- For transition features: similarly compute  $P(y_i = s_1, y_{i+1} = s_2 | \mathbf{x})$  using forward-backward algorithm
- ...but you can build a pretty good system without learned transition features (use heuristic weights, or just enforce constraints like B-PER -> I-ORG is illegal)

# 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: argmax  $P(\mathbf{y}|\mathbf{x})$  from Viterbi
- Learning: run forward-backward to compute the marginals

$$\frac{\partial}{\partial w} \mathcal{L}(\mathbf{y}^*, \mathbf{x}) = \sum_{i=1}^n f_e(y_i^*, i, \mathbf{x}) - \sum_{i=1}^n \sum_s P(y_i = s | \mathbf{x}) f_e(s, i, \mathbf{x})$$

# Learning: Pseudocode

- for each epoch
  - for each example
    - extract features on each emission and transition (look up in cache)
    - compute marginal probabilities with forward-backward
    - compute phis ( $\phi_t, \phi_e$ ) based on features and weights
    - accumulate gradient over all emissions and transitions
    - Perform gradient ascent

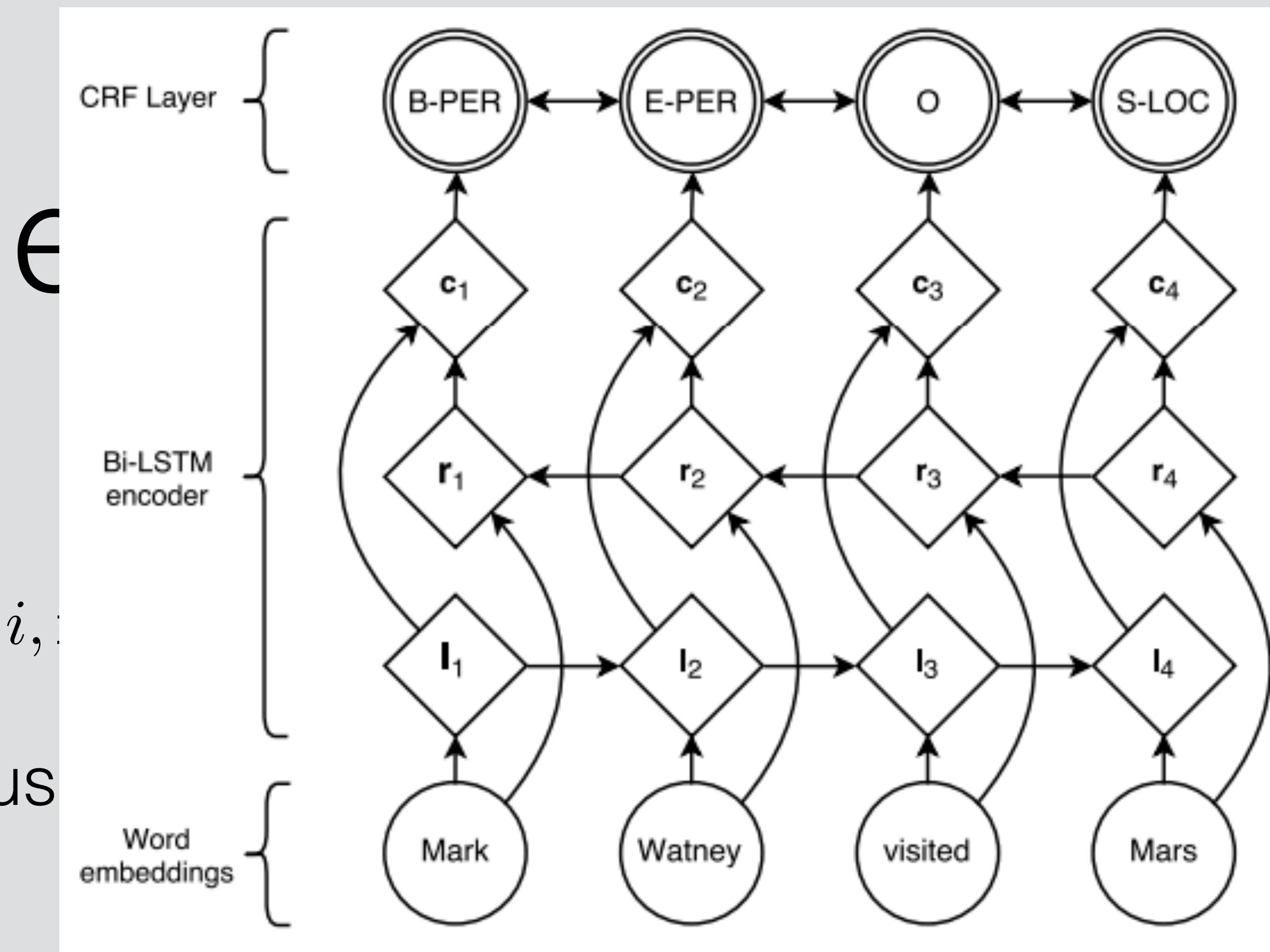
$$\frac{\partial}{\partial w} \mathcal{L}(\mathbf{y}^*, \mathbf{x}) = \sum_{i=1}^n f_e(y_i^*, i, \mathbf{x}) - \sum_{i=1}^n \sum_s P(y_i = s | \mathbf{x}) f_e(s, i, \mathbf{x})$$

# Last Note: Neural Nets

- So far, linear features

$$P(\mathbf{y} = \tilde{\mathbf{y}} | \mathbf{x}) = \frac{1}{Z} \exp \left( \sum_{i=1}^n \mathbf{w}^\top \mathbf{f}_e(\tilde{y}_i, i, \mathbf{x}) \right)$$

- How about nonlinear features, e.g., using neural nets?



- e.g., can we use RNN or LSTM to learn neural features for CRF?
- Lample et al., 2016 (F1 vs. SOTA using hand-crafted features): 90.94 (en, vs. 91.2), 78.76 (de, vs. 76.22), 81.74 (nl, vs. 82.84), 85.75 (es, vs. 82.95)

# Structured Prediction

- Sequence labeling
  - POS Tagging & Hidden Markov Models (HMMs)
  - Named Entity Recognition & Conditional Random Fields (CRFs)
- Syntactic Parsing
  - Dependency parsing
  - Constituency parsing (next lecture)
- Semantic Parsing (next lecture)

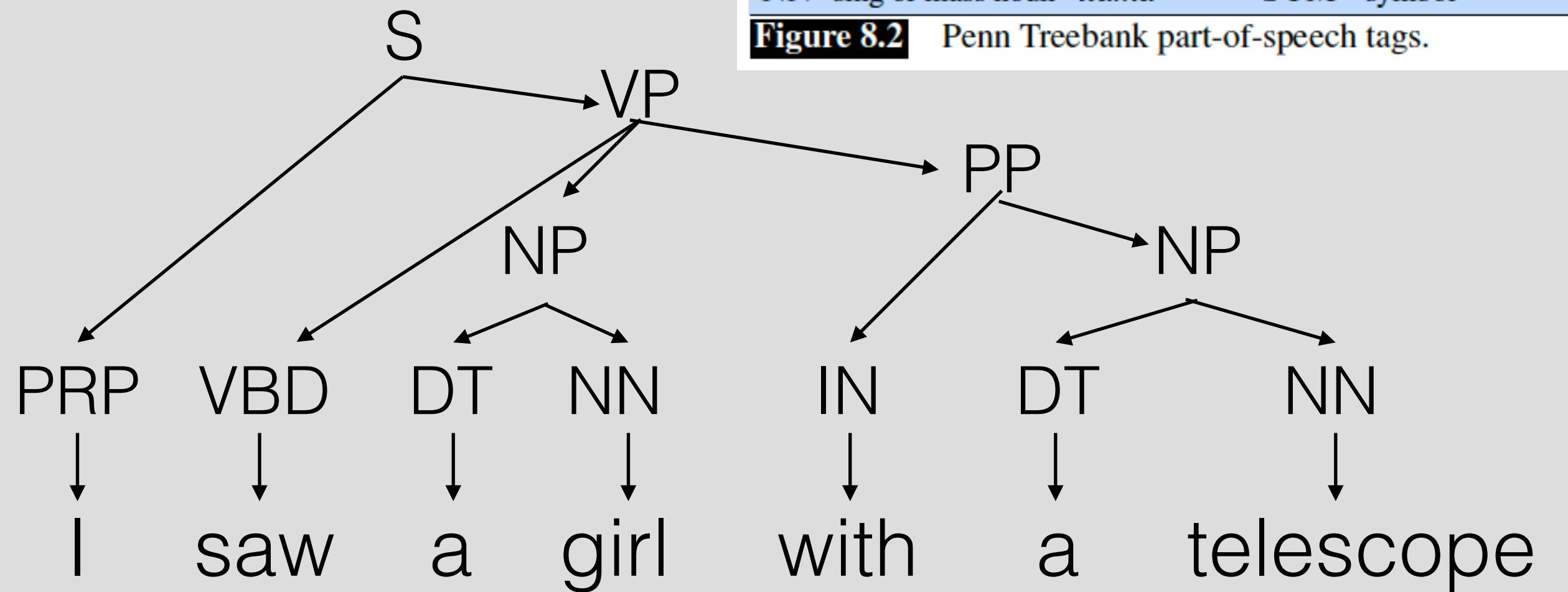
# Syntactic Parsing

# Syntax

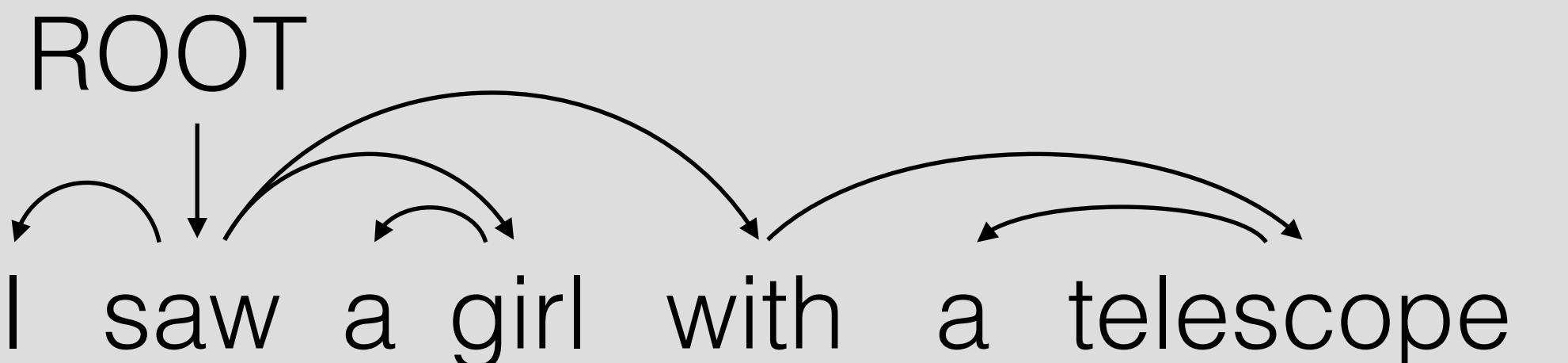
- Study of word order and how words form sentences
- Why do we care about syntax?
  - Recognize verb-argument structures (who is doing what to whom?)
  - Multiple interpretations of words (noun or verb?)
  - Higher level of abstraction beyond words: some languages are SVO, some are VSO, some are SOV, parsing can canonicalize

# Two Ty Linguistic

- **Phrase structure:** focus on the struc



- **Dependency:** focus on relations between words



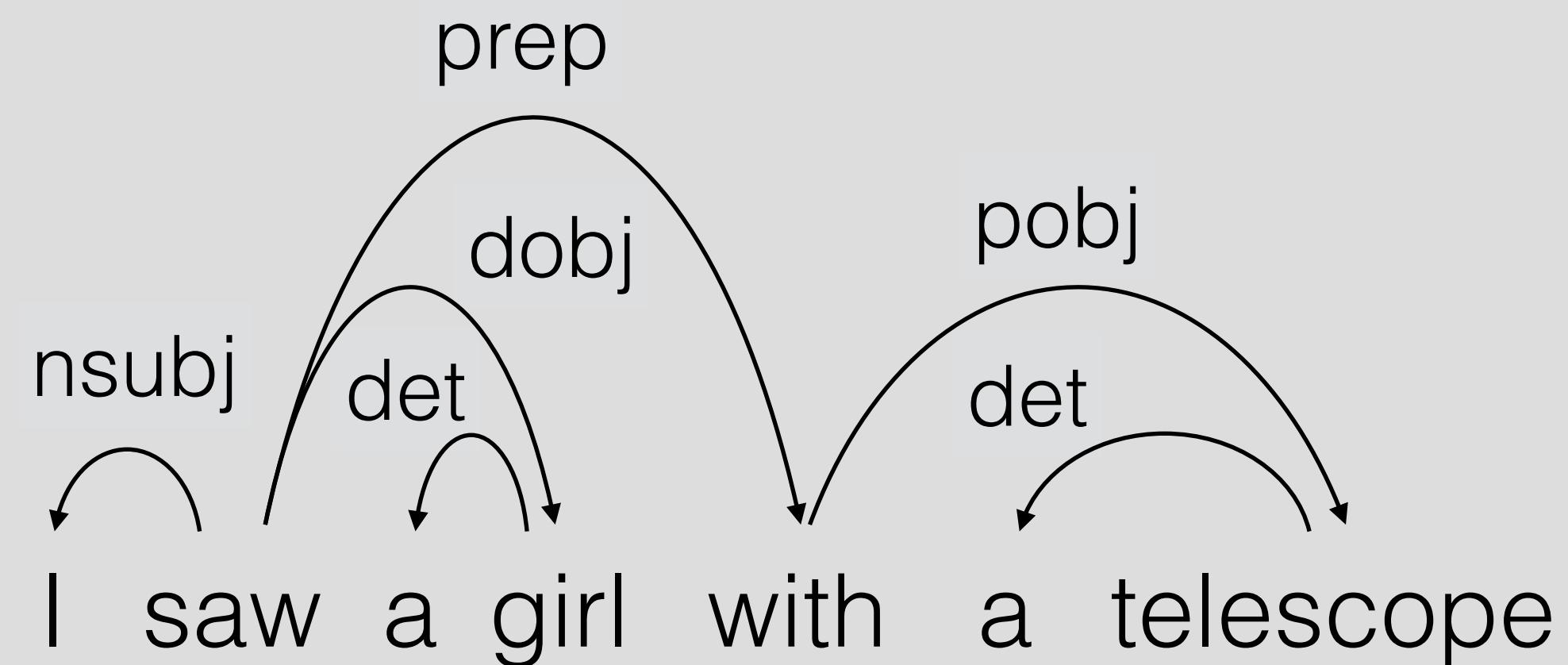
Tag	Description	Example	Tag	Description	Example	Tag	Description	Example
CC	coord. conj.	<i>and, but, or</i>	NNP	proper noun, sing.	<i>IBM</i>	TO	“to”	<i>to</i>
CD	cardinal number	<i>one, two</i>	NNPS	proper noun, plu.	<i>Carolinias</i>	UH	interjection	<i>ah, oops</i>
DT	determiner	<i>a, the</i>	NNS	noun, plural	<i>llamas</i>	VB	verb base	<i>eat</i>
EX	existential ‘there’	<i>there</i>	PDT	predeterminer	<i>all, both</i>	VBD	verb past tense	<i>ate</i>
FW	foreign word	<i>mea culpa</i>	POS	possessive ending	<i>'s</i>	VBG	verb gerund	<i>eating</i>
IN	preposition/ subordin-conj	<i>of, in, by</i>	PRP	personal pronoun	<i>I, you, he</i>	VBN	verb past participle	<i>eaten</i>
JJ	adjective	<i>yellow</i>	PRP\$	possess. pronoun	<i>your, one's</i>	VBP	verb non-3sg-pr	<i>eat</i>
JJR	comparative adj	<i>bigger</i>	RB	adverb	<i>quickly</i>	VBZ	verb 3sg pres	<i>eats</i>
JJS	superlative adj	<i>wildest</i>	RBR	comparative adv	<i>faster</i>	WDT	wh-determ.	<i>which, that</i>
LS	list item marker	<i>1, 2, One</i>	RBS	superlatv. adv	<i>fastest</i>	WP	wh-pronoun	<i>what, who</i>
MD	modal	<i>can, should</i>	RP	particle	<i>up, off</i>	WP\$	wh-possess.	<i>whose</i>
NN	sing or mass noun	<i>llama</i>	SYM	symbol	<i>+%, &amp;</i>	WRB	wh-adverb	<i>how, where</i>

Figure 8.2 Penn Treebank part-of-speech tags.

# Syntactic Parsing 1: Dependency Parsing

# Why Dependencies?

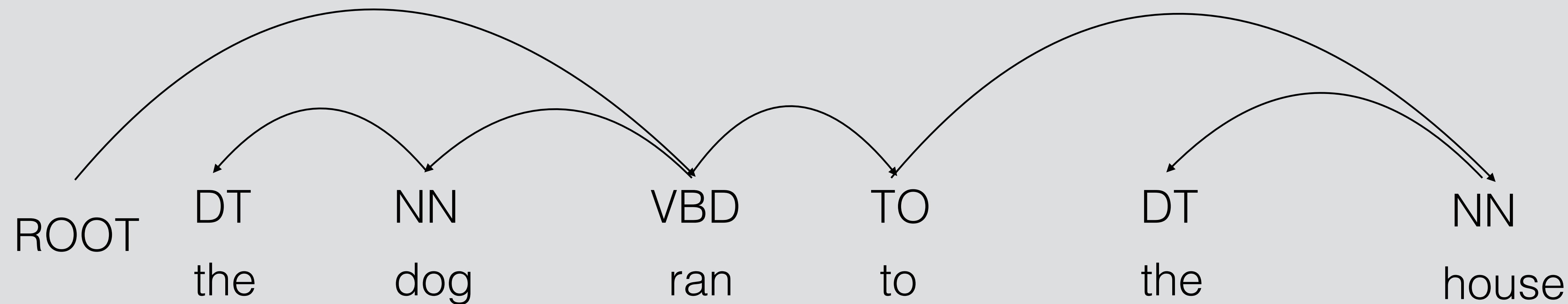
- Demonstrate the relationships between words in a straightforward way



- Particularly good for multilinguality, e.g. phrase-structure can be hard to define in languages with free word order

# Dependency Parsing

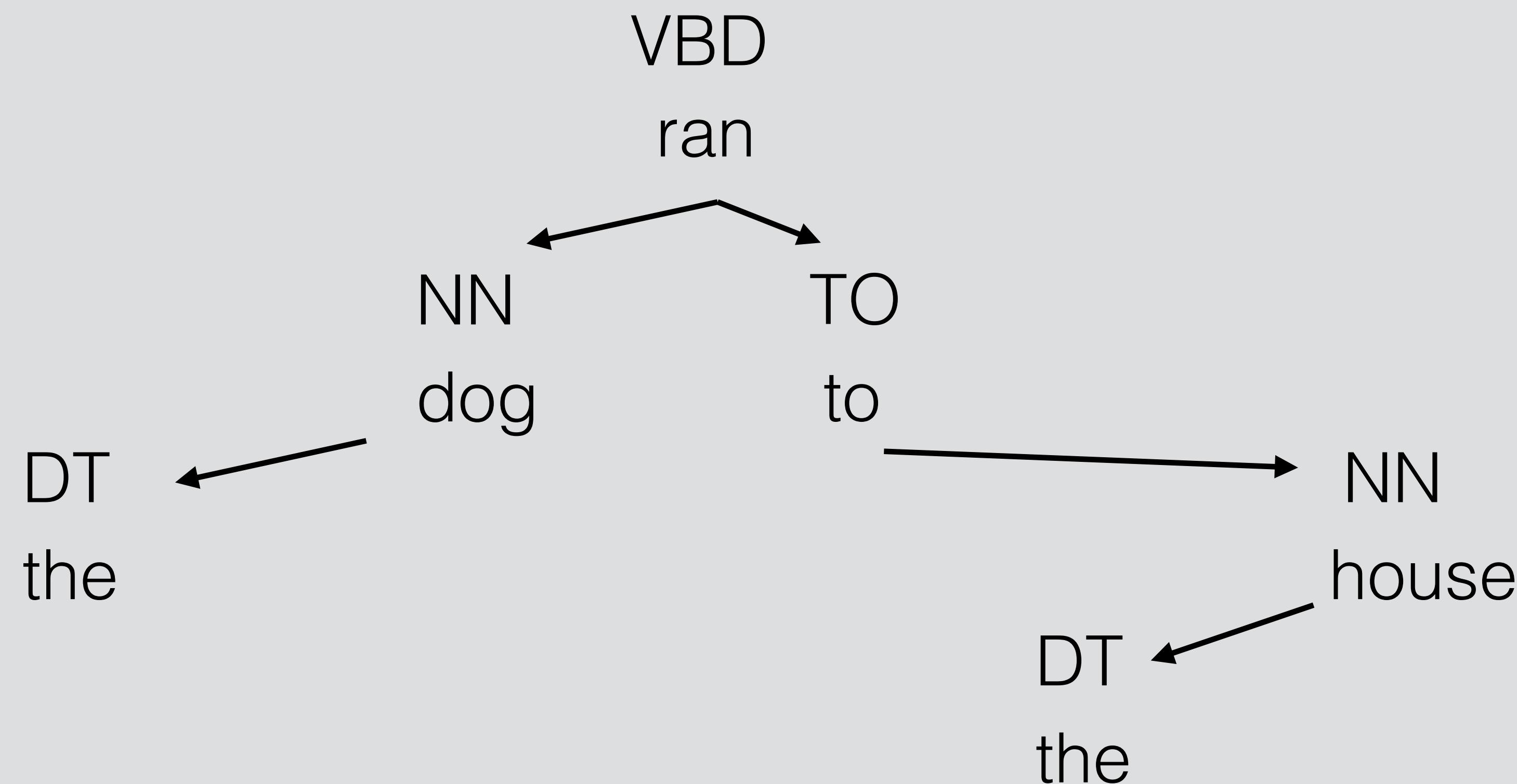
- Dependency syntax: syntactic structure is defined by these arcs
  - *Head (parent, governor)* connected to *dependent (child, modifier)*
  - Each word has exactly one parent except for the *ROOT* symbol, dependencies must form a directed acyclic graph



- POS tags same as before, usually run a tagger first as preprocessing

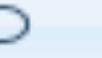
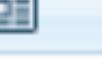
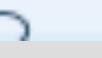
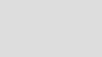
# Dependency Parsing

- Still a notion of hierarchy! Subtrees often align with constituents



# Universal Dependencies Treebank

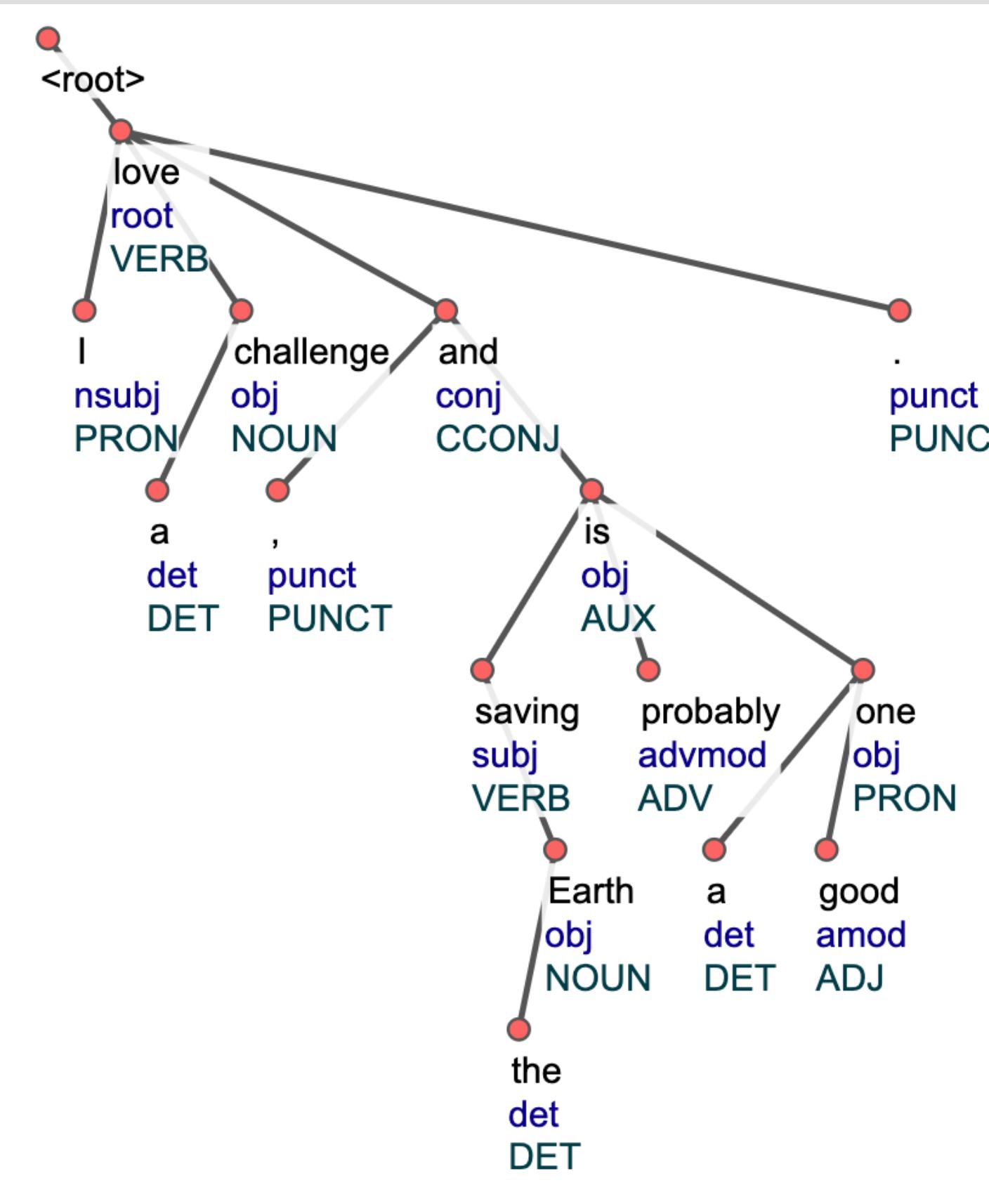
- Standard format for parse trees in many languages

▶  <b>Abaza</b>	1	3K		<b>Northwest Caucasian</b>
▶  <b>Afrikaans</b>	1	49K		<b>IE, Germanic</b>
▶  <b>Akkadian</b>	1	1K		<b>Afro-Asiatic, Semitic</b>
▶  <b>Albanian</b>	1	<1K		<b>IE, Albanian</b>
▶  <b>Amharic</b>	1	10K		<b>Afro-Asiatic, Semitic</b>
▶  <b>Ancient Greek</b>	2	416K		<b>IE, Greek</b>
▶  <b>Arabic</b>	3	1,042K		<b>Afro-Asiatic, Semitic</b>
▶  <b>Armenian</b>	1	52K		<b>IE, Armenian</b>
▶  <b>Assyrian</b>	1	<1K		<b>Afro-Asiatic, Semitic</b>
▶  <b>Bambara</b>	1	13K		<b>Mande</b>
▶  <b>Basque</b>	1	121K		<b>Basque</b>
▶  <b>Belarusian</b>	1	13K		<b>IE, Slavic</b>
▶  <b>Bhojpuri</b>	2	6K		<b>IE, Indic</b>
▶  <b>Breton</b>	1	10K		<b>IE, Celtic</b>
▶  <b>Bulgarian</b>	1	156K		<b>IE, Slavic</b>
▶  <b>Buryat</b>	1	10K		<b>Mongolic</b>
▶  <b>Cantonese</b>	1	13K		<b>Sino-Tibetan</b>

<https://universaldependencies.org/>

# Semantic and Syntactic Dependencies

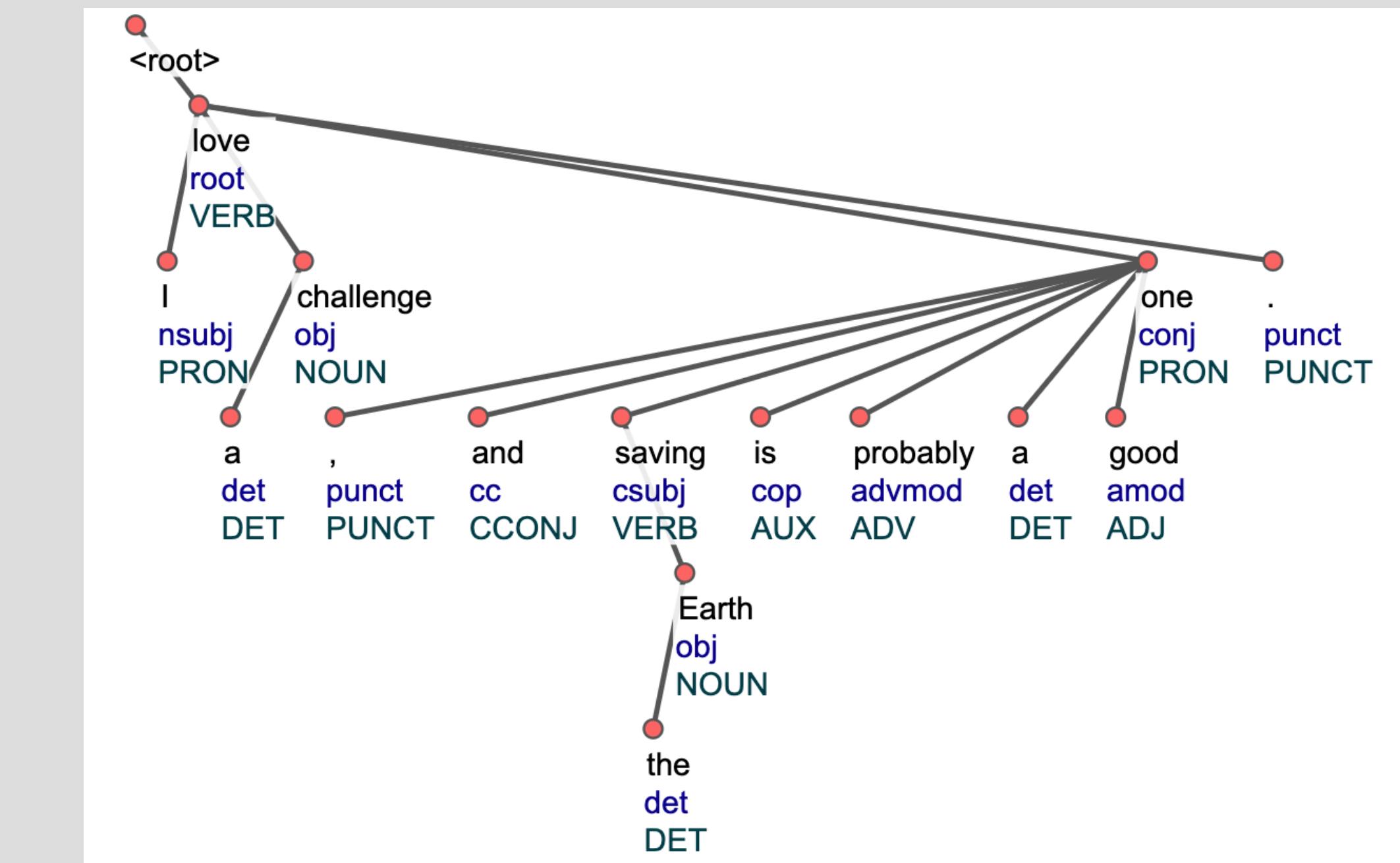
## Syntactic (SUD)



Deeper, reflect phrase structure,  
more function word heads

<https://surfacesyntacticud.github.io/>

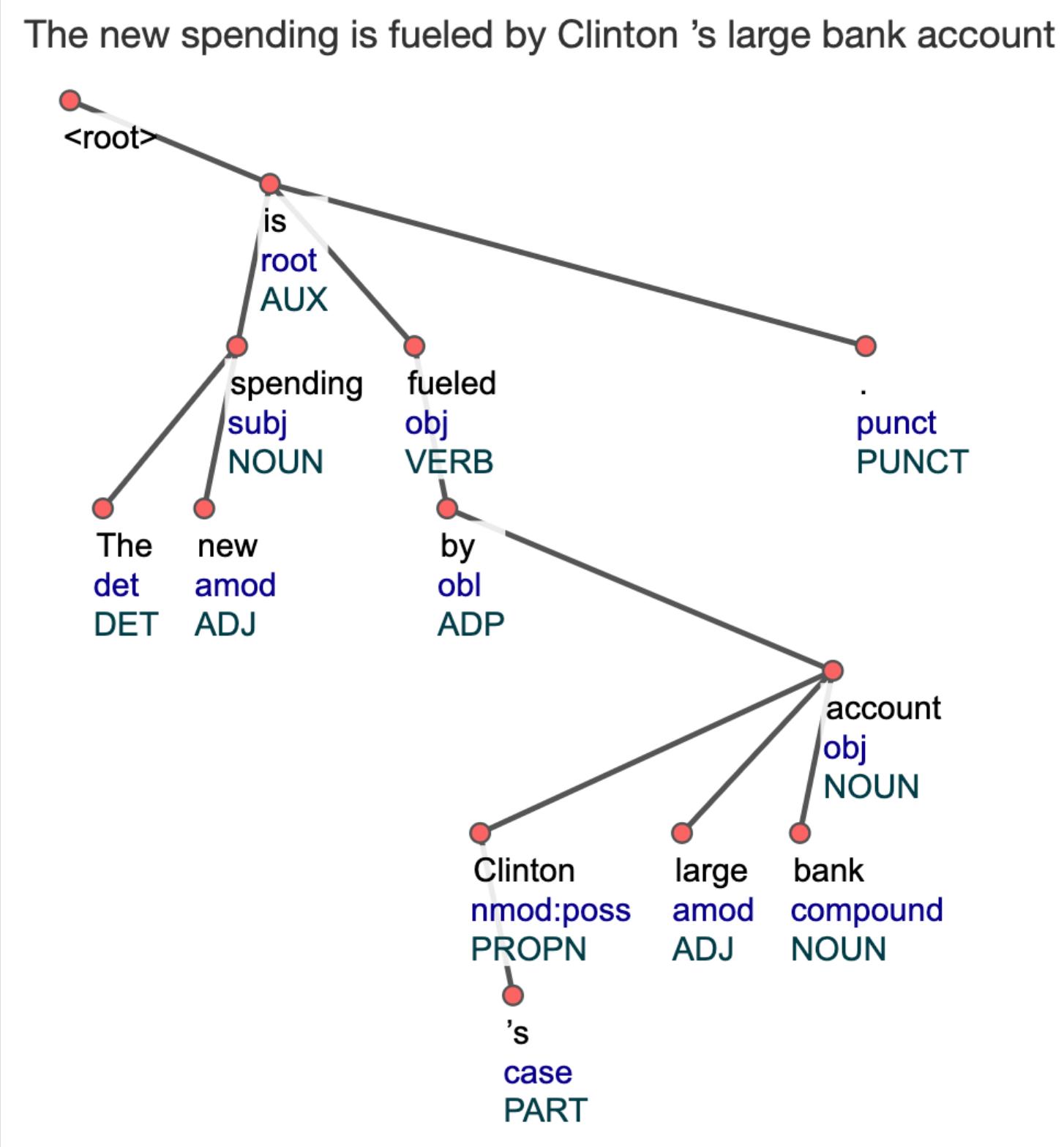
## Semantic (UD)



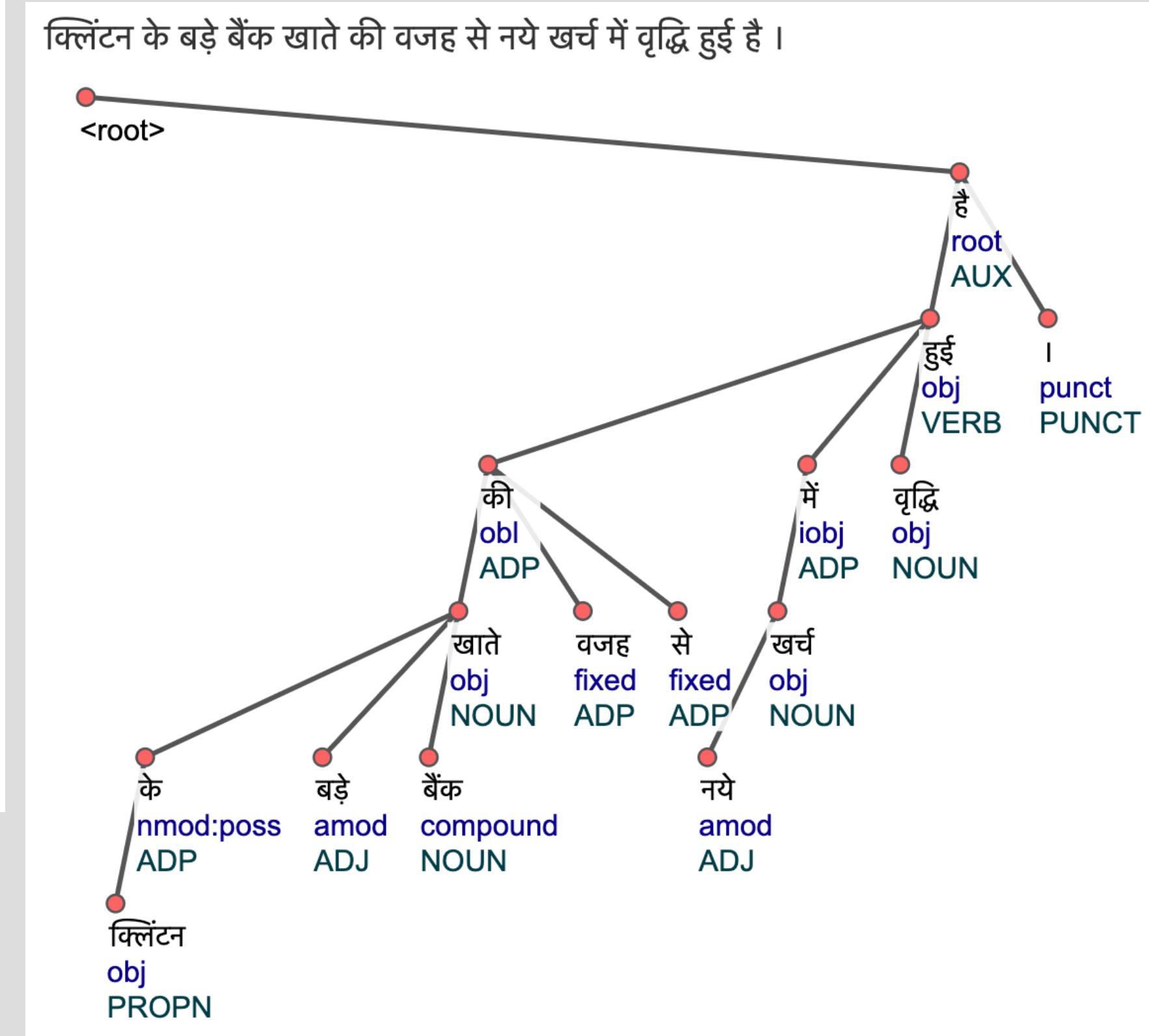
Flatter, semantically related words closer,  
more content word heads

# Cross-lingual Differences In Structure

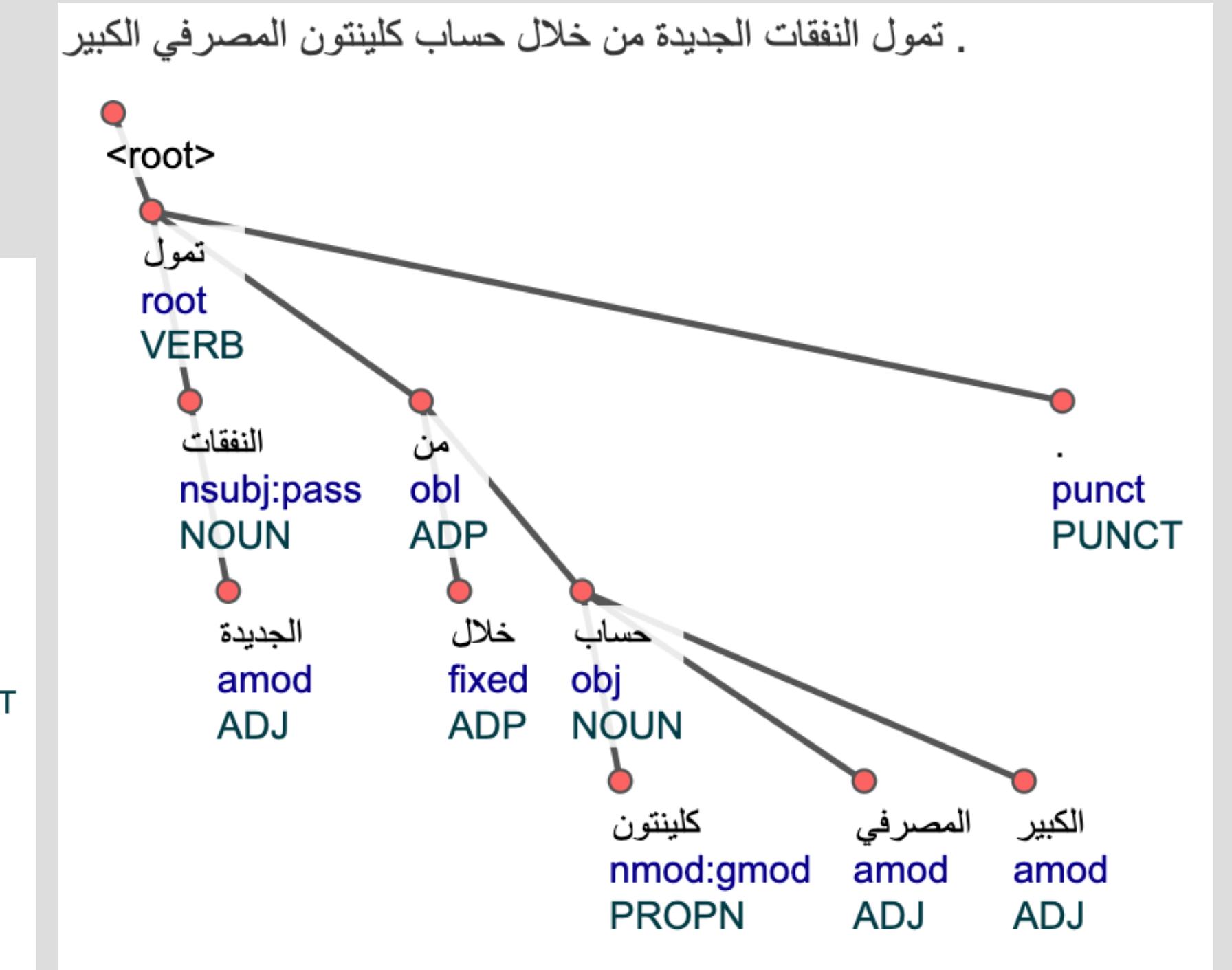
## English: SVO



## Hindi: Verb Final



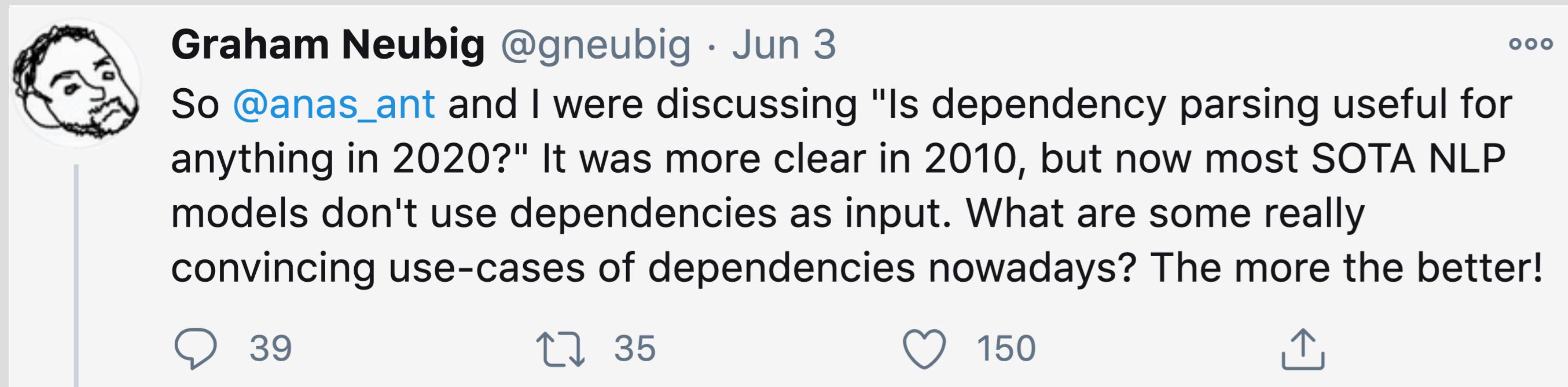
## Arabic: Verb Initial



What Can we Do w/  
Dependencies?

# Use Cases of Dependencies?

- Previously, used for feature engineering in systems (and still useful in some cases)
- Now: more useful for human-facing applications



**Graham Neubig** @gneubig · Jun 3

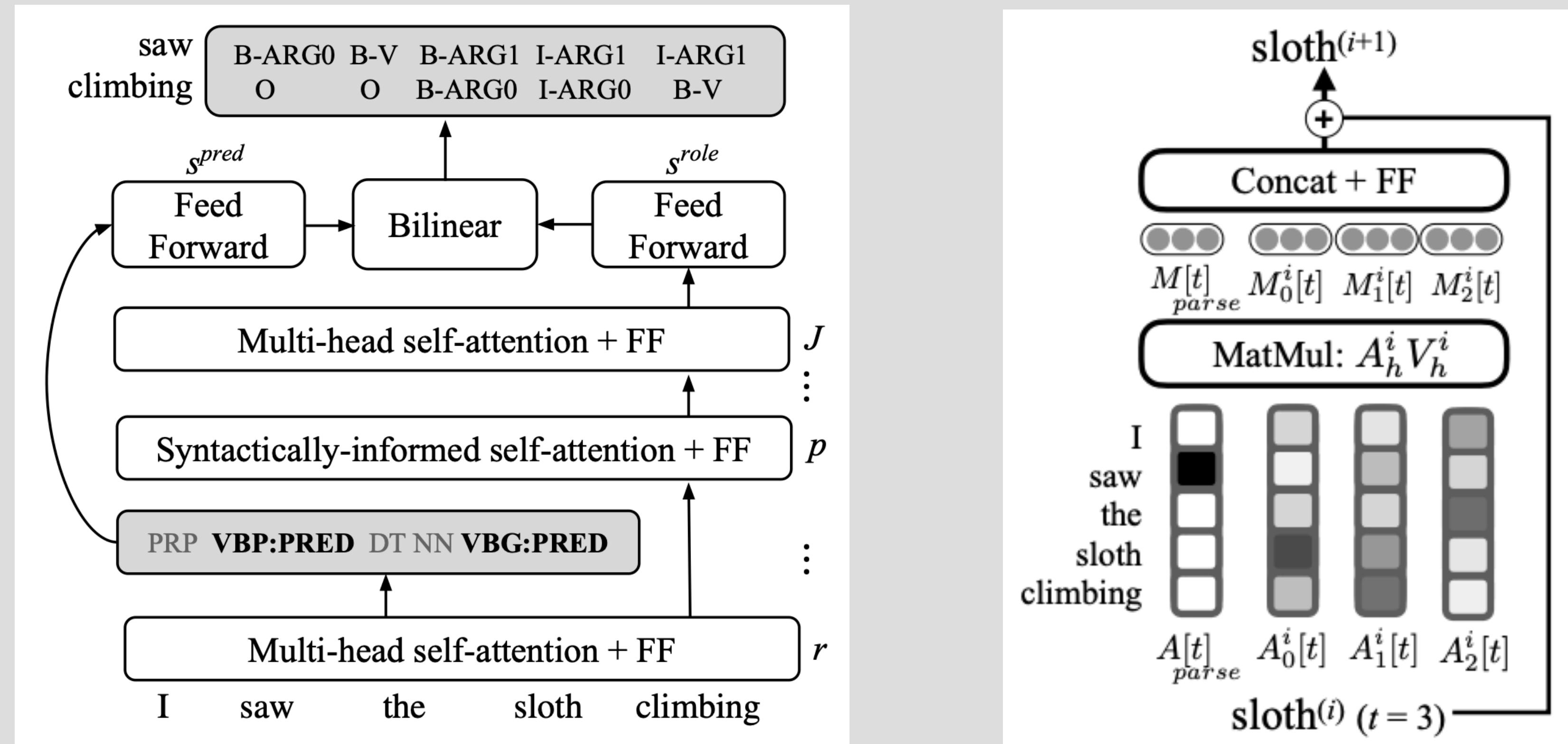
So [@anas\\_ant](#) and I were discussing "Is dependency parsing useful for anything in 2020?" It was more clear in 2010, but now most SOTA NLP models don't use dependencies as input. What are some really convincing use-cases of dependencies nowadays? The more the better!

39 35 150

<https://twitter.com/gneubig/status/1268238606101032962?lang=en>

# Example 1: Adding Inductive Bias to Neural Models

- Bias self attention to follow syntax



# Example 2: Searching over Parsed Corpora

- Search using "syntactic regex"

Syntactic Search [?](#) [?](#)

Query  
<>founder:[e]Paul was a t:[w]founder of <>entity:[e]Microsoft

[LOAD EXAMPLE ▾](#)

Result:

2321036 | Anderson who is the founder and director of the World Education Foundation , currently engages in research and implementation of sustainable developmental projects , globally .

2349619 | Anderson and co-organized with entrepreneur Robin Bates who is the founder and CEO of Caf e de la Soul and La Jolie Noire Media , and co-founder of Black Paris Divas .

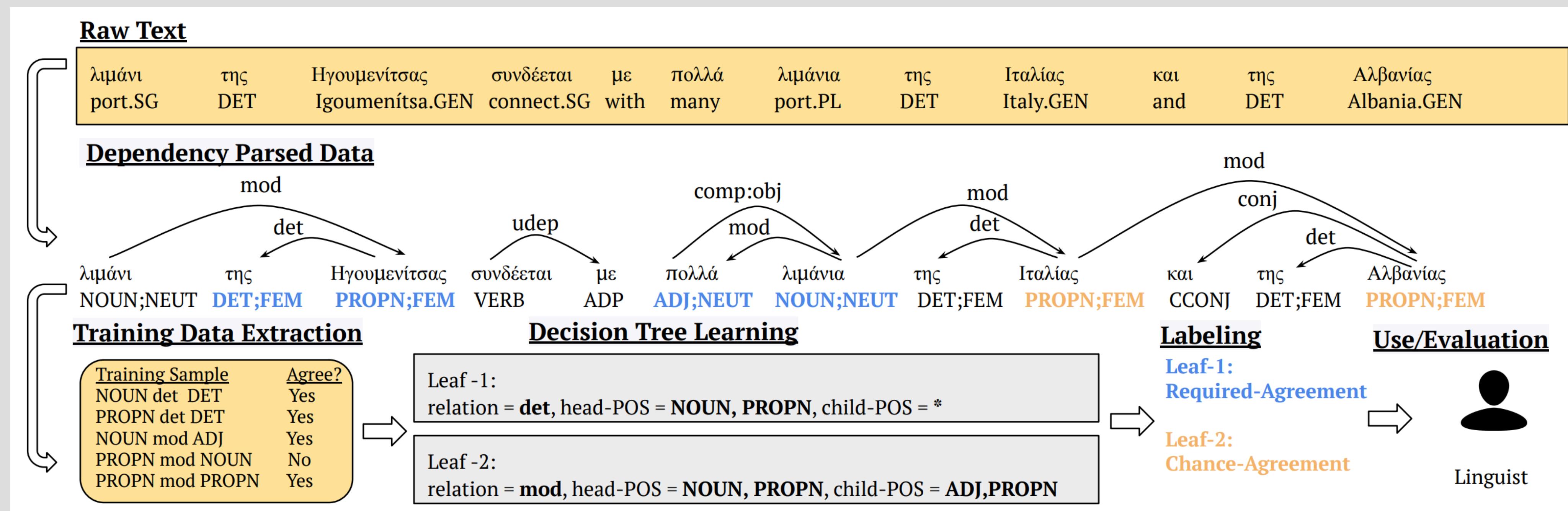
2352872 | Anderson is founder and chairman of Interface Inc .

2497762 | Ananda Kar is the founder of the Hemlock Society , that teaches aspirants how to successfully commit suicide .

The screenshot shows a syntactic search interface. At the top, there's a 'Syntactic Search' header with a question mark icon and a green square icon. Below it is a 'Query' input field containing the regex-like expression '<>founder:[e]Paul was a t:[w]founder of <>entity:[e]Microsoft'. To the right of the input are three buttons: a magnifying glass for search, a download icon, and a 'LOAD EXAMPLE ▾' button. The main area displays four search results, each with a unique ID (2321036, 2349619, 2352872, 2497762) and an 'i' icon. The results are presented in a list format, with each entry containing a snippet of text and its corresponding entity annotations. The annotations are color-coded: blue for entities (PERSON, ORGANIZATION), green for tokens (t), and red for other entities (entity). The interface also includes a dependency parse diagram at the top, showing the relationships between 'Paul', 'founder', and 'Microsoft'.

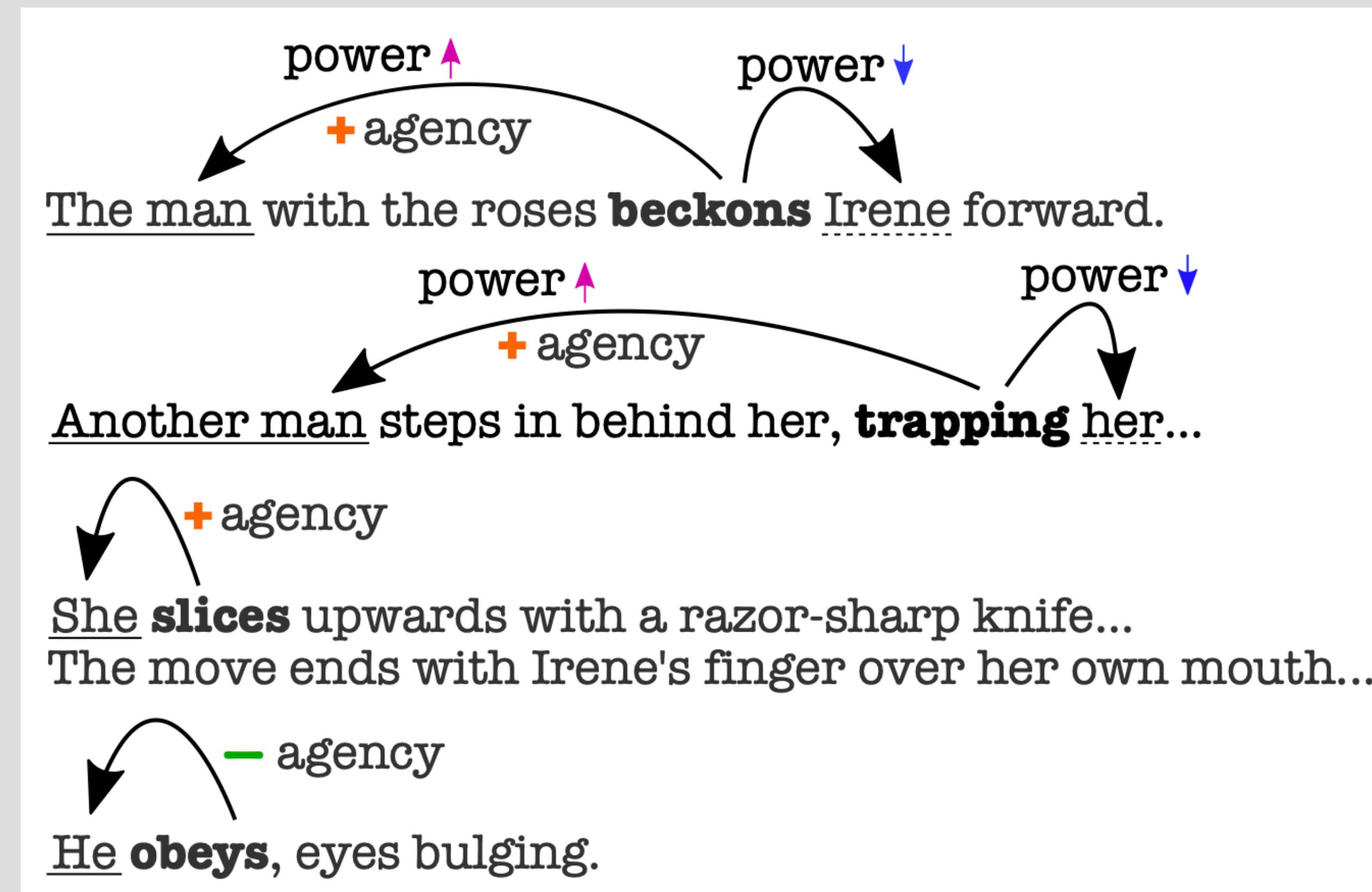
# Example 3: Understanding Language Structure

- Example of extracting morphological agreement rules using dependency relations



# Example 4: Analysis of Other Linguistic Phenomena

- Examining power and agency in film scripts



# Parsing

- Predicting linguistic structure from input sentence
- **Transition-based models**
  - step through actions one-by-one until we have output
  - modeling “what transition action to take in each step”
- **Graph-based models**
  - calculate probability of each edge/constituent, and perform some sort of dynamic programming
  - modeling “what edge to build for connecting the constituents”

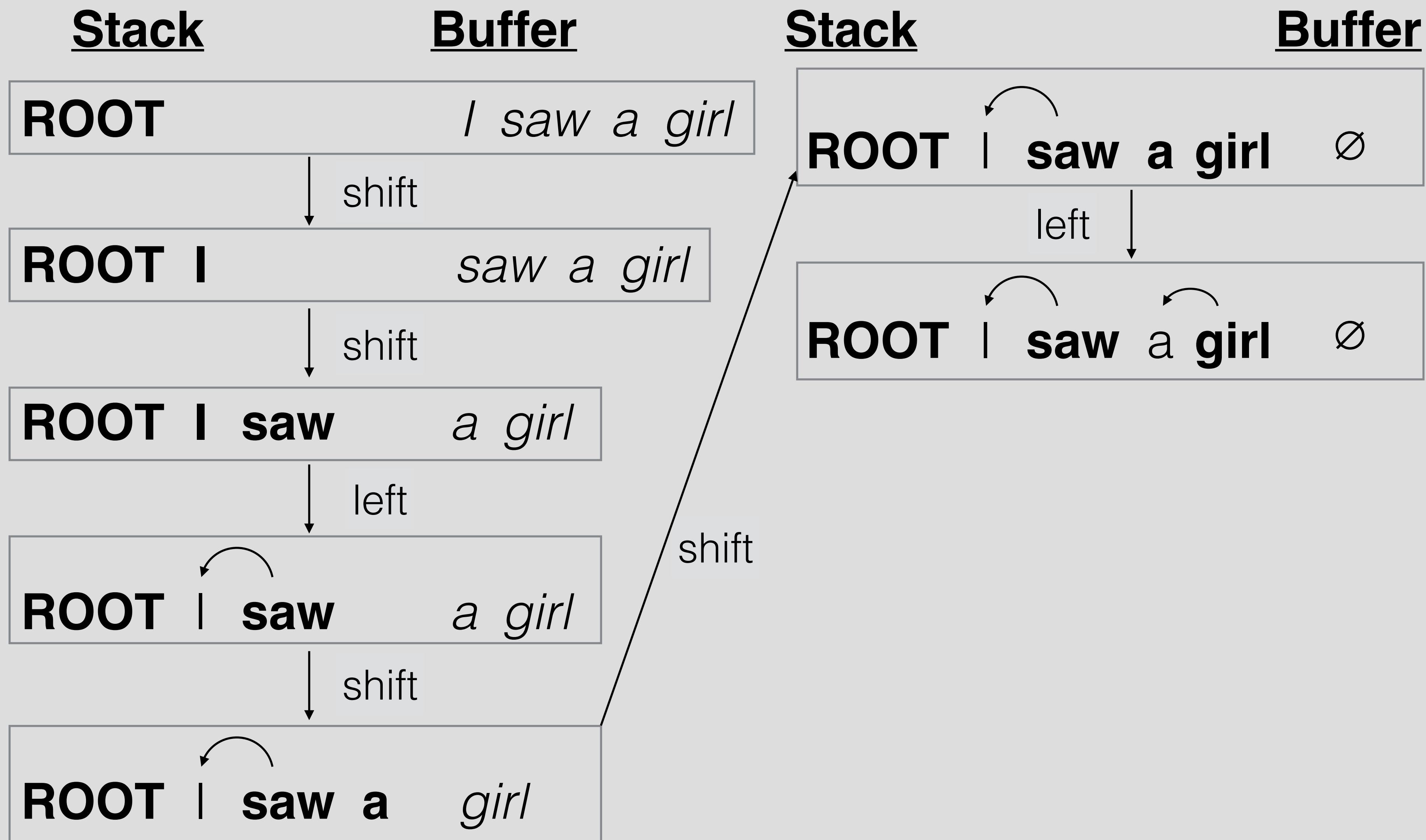
# Shift-reduce Parsing

# Arc Standard Shift-Reduce Parsing

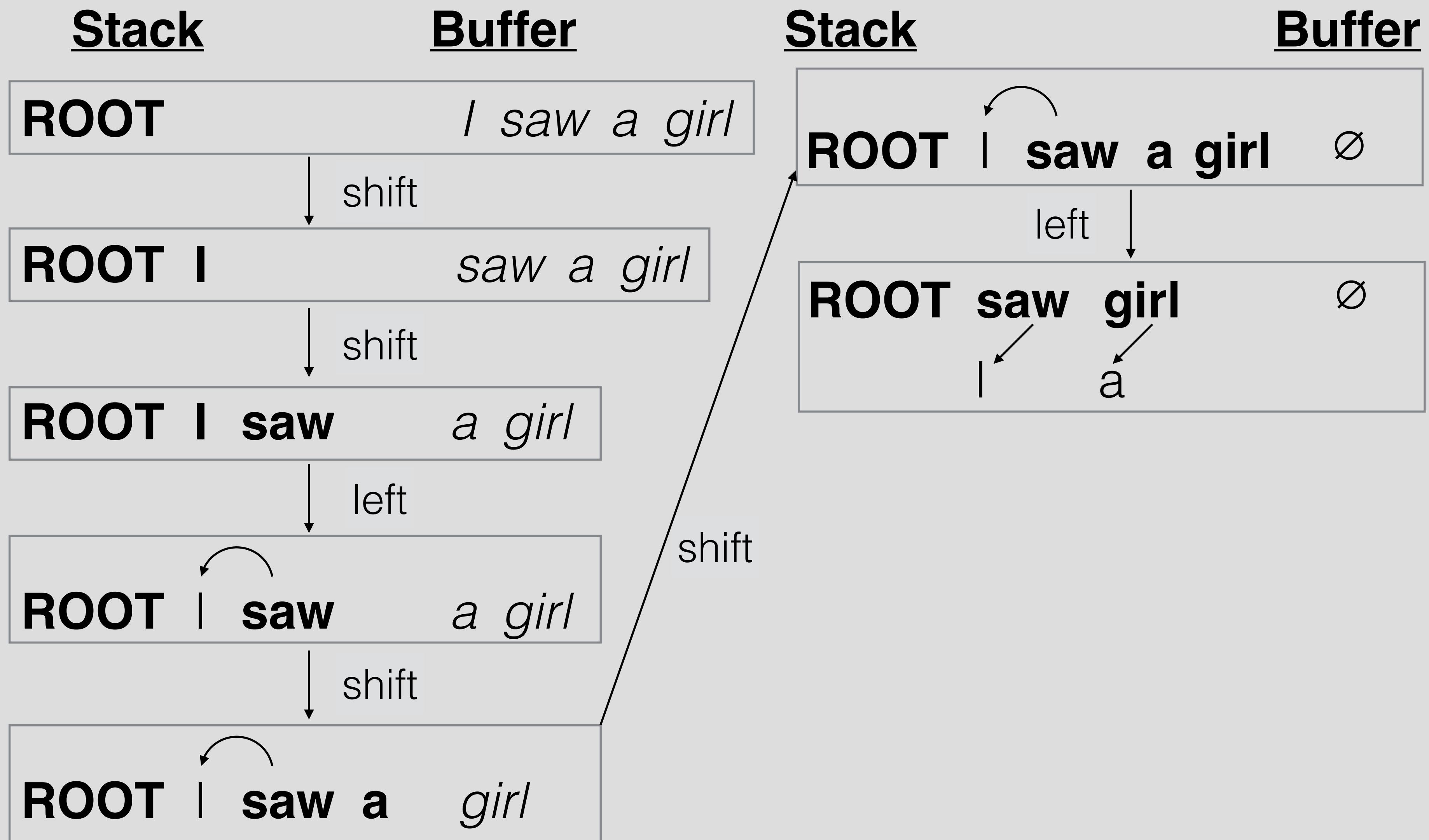
(Yamada & Matsumoto 2003, Nivre 2003)

- Process words one-by-one left-to-right
- Two data structures
  - **Queue/buffer:** of unprocessed words
  - **Stack:** of partially processed words
- At each point choose
  - **shift:** move one word from queue to stack
  - **reduce left:** top word on stack is head of second word
  - **reduce right:** second word on stack is head of top word
- Learn how to choose each action with a classifier

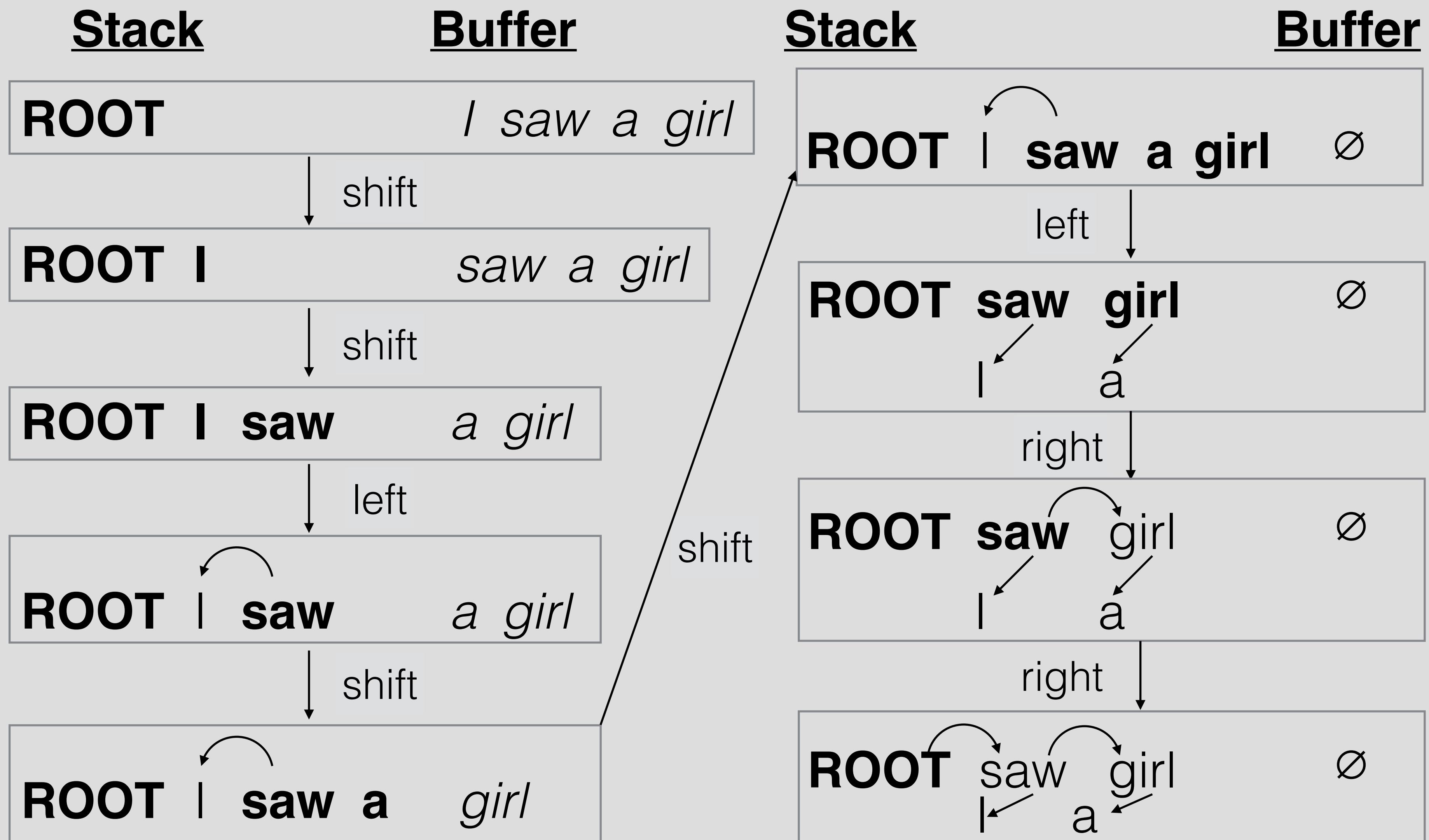
# Shift Reduce Example



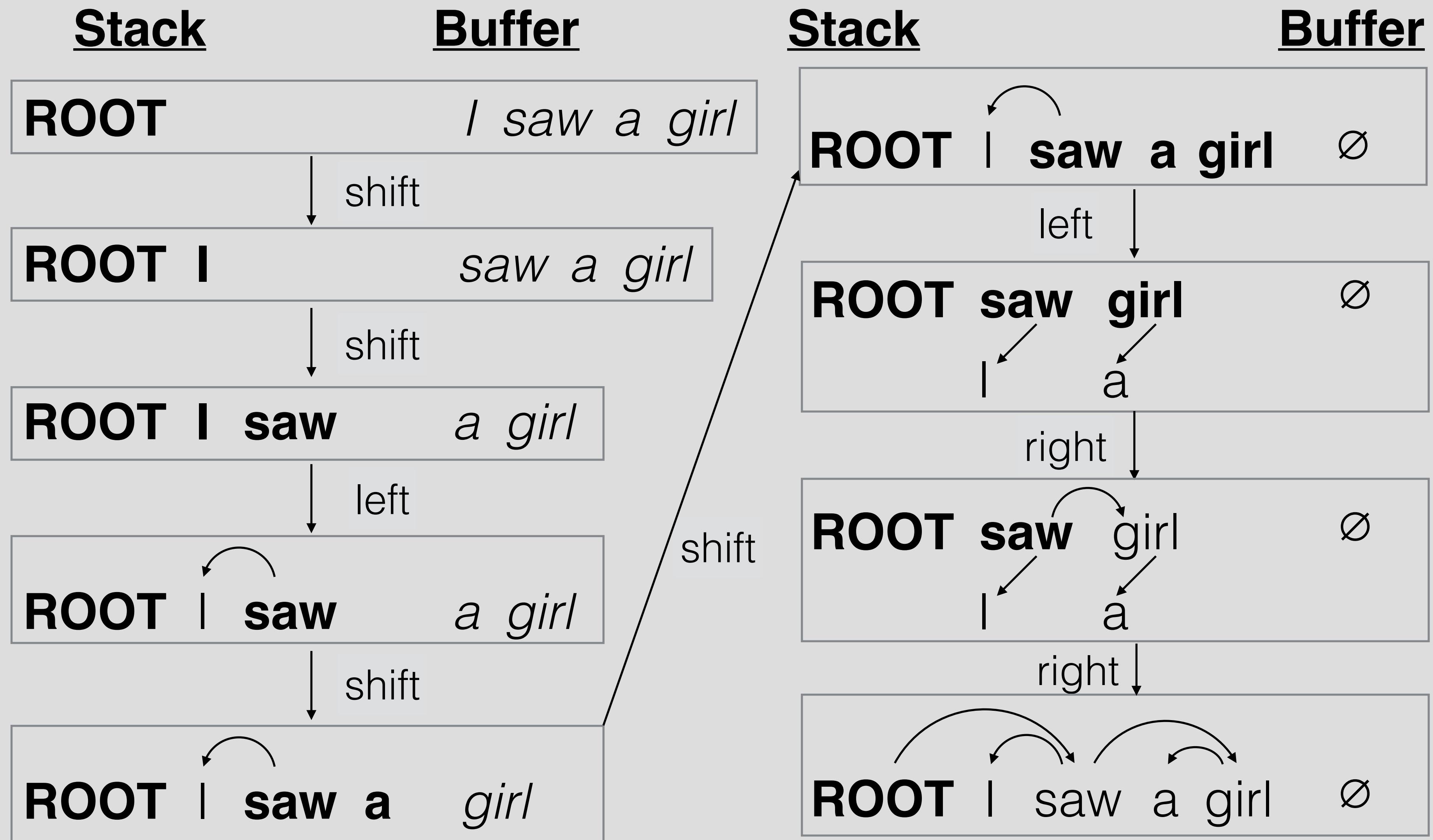
# Shift Reduce Example



# Shift Reduce Example



# Shift Reduce Example

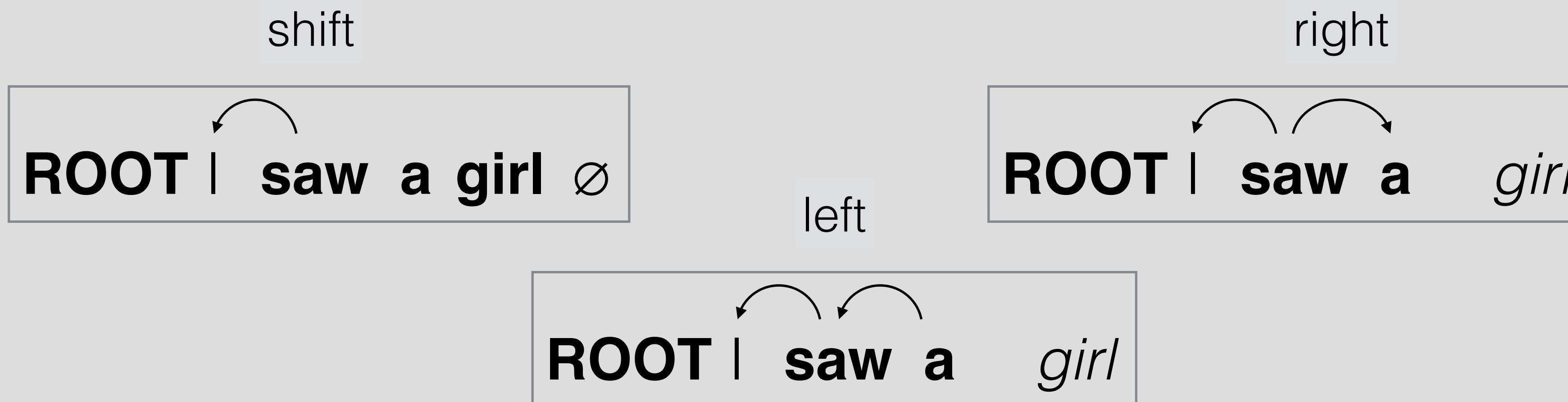


# Classification for Shift-reduce

- Given a **configuration**

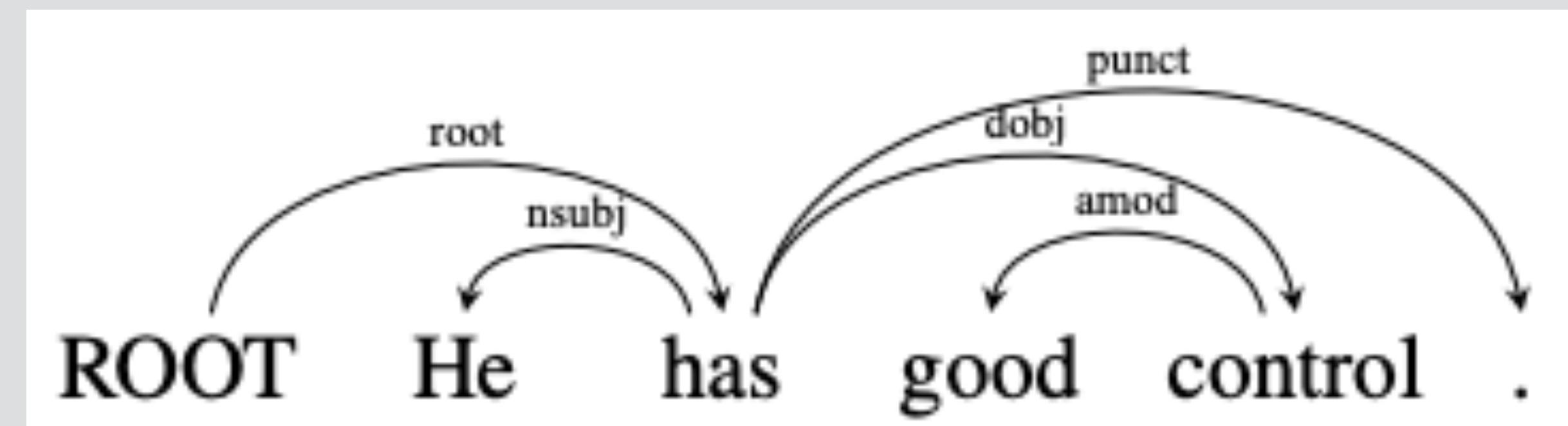


- Which **action** do we choose?



# What Actions to Choose?

- Building a multi-class classifier
  - $\text{num\_classes} = 2 \times \text{num\_arc\_types}$  (left/right reduce with a certain arc type) + 1 (shift)



# What Actions to Choose?

- Building a multi-class classifier
  - $\text{num\_classes} = 2 \times \text{num\_arc\_types}$  (left/right reduce with a certain arc type) + 1 (shift)
- What features?

$s_i$ : the i-th top element on the stack;

$w$ : word token;

$t$ : POS tag;

$b_i$ : the i-th element on the buffer.

---

### Single-word features (9)

$s_1.w; s_1.t; s_1.wt; s_2.w; s_2.t;$   
 $s_2.wt; b_1.w; b_1.t; b_1.wt$

---

### Word-pair features (8)

$s_1.wt \circ s_2.wt; s_1.wt \circ s_2.w; s_1.wts_2.t;$   
 $s_1.w \circ s_2.wt; s_1.t \circ s_2.wt; s_1.w \circ s_2.w$   
 $s_1.t \circ s_2.t; s_1.t \circ b_1.t$

---

### Three-word feaures (8)

$s_2.t \circ s_1.t \circ b_1.t; s_2.t \circ s_1.t \circ lc_1(s_1).t;$   
 $s_2.t \circ s_1.t \circ rc_1(s_1).t; s_2.t \circ s_1.t \circ lc_1(s_2).t;$   
 $s_2.t \circ s_1.t \circ rc_1(s_2).t; s_2.t \circ s_1.w \circ rc_1(s_2).t;$   
 $s_2.t \circ s_1.w \circ lc_1(s_1).t; s_2.t \circ s_1.w \circ b_1.t$

---

# What Actions to Choose?

- Building a multi-class classifier
  - $\text{num\_classes} = 2 \times \text{num\_arc\_types}$  (left/right reduce with a certain arc type) + 1 (shift)
- What features?
  - Issues with conventional feature-engineering approaches:  
Sparsity, incompleteness (in terms of feature space), expensive feature computation (as the parsing runs dynamically and you have to look up the feature in a huge feature table)

# What Actions to Choose?

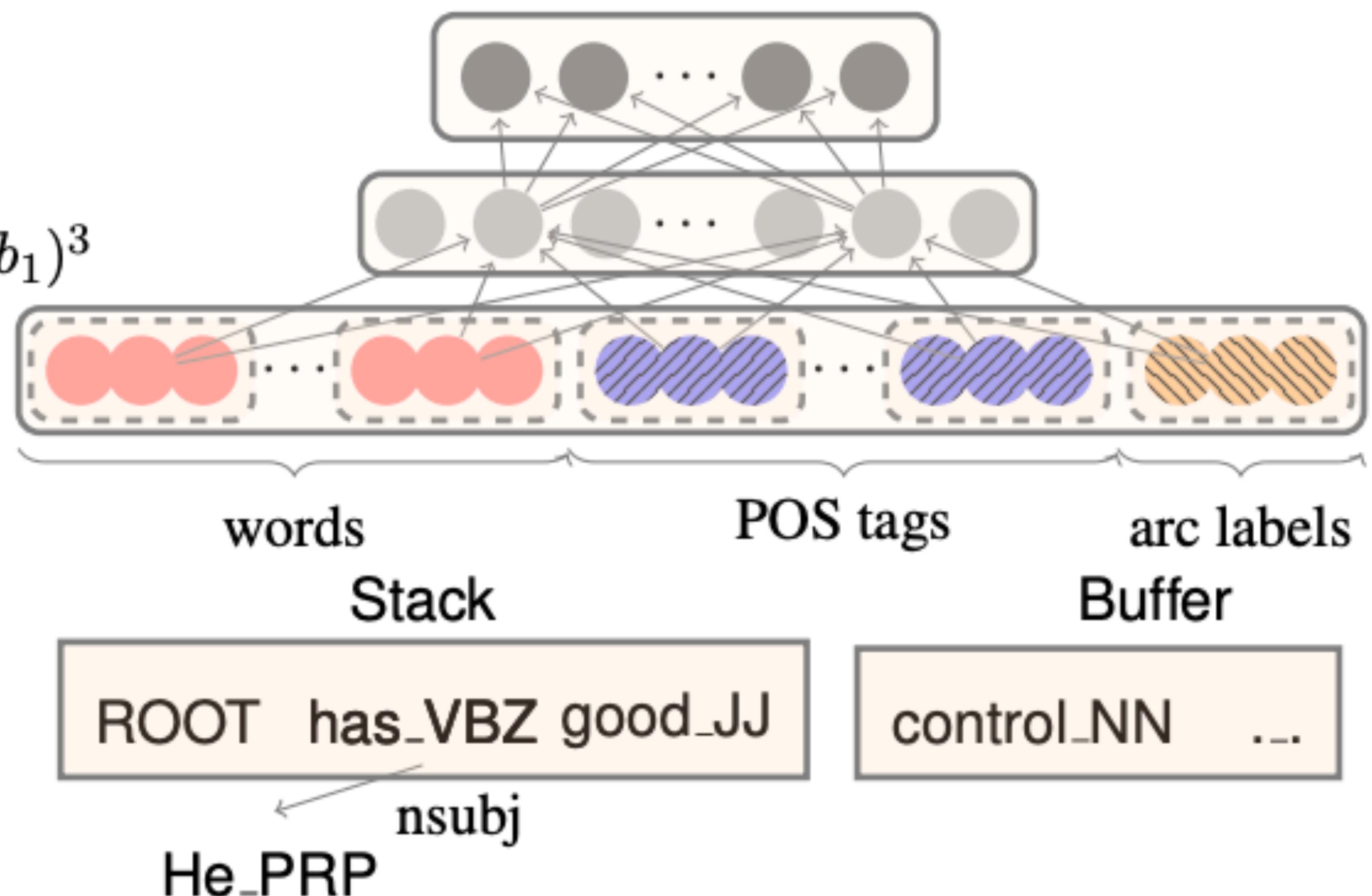
**Softmax layer:**

$$p = \text{softmax}(W_2 h)$$

**Hidden layer:**

$$h = (W_1^w x^w + W_1^t x^t + W_1^l x^l + b_1)^3$$

**Input layer:**  $[x^w, x^t, x^l]$

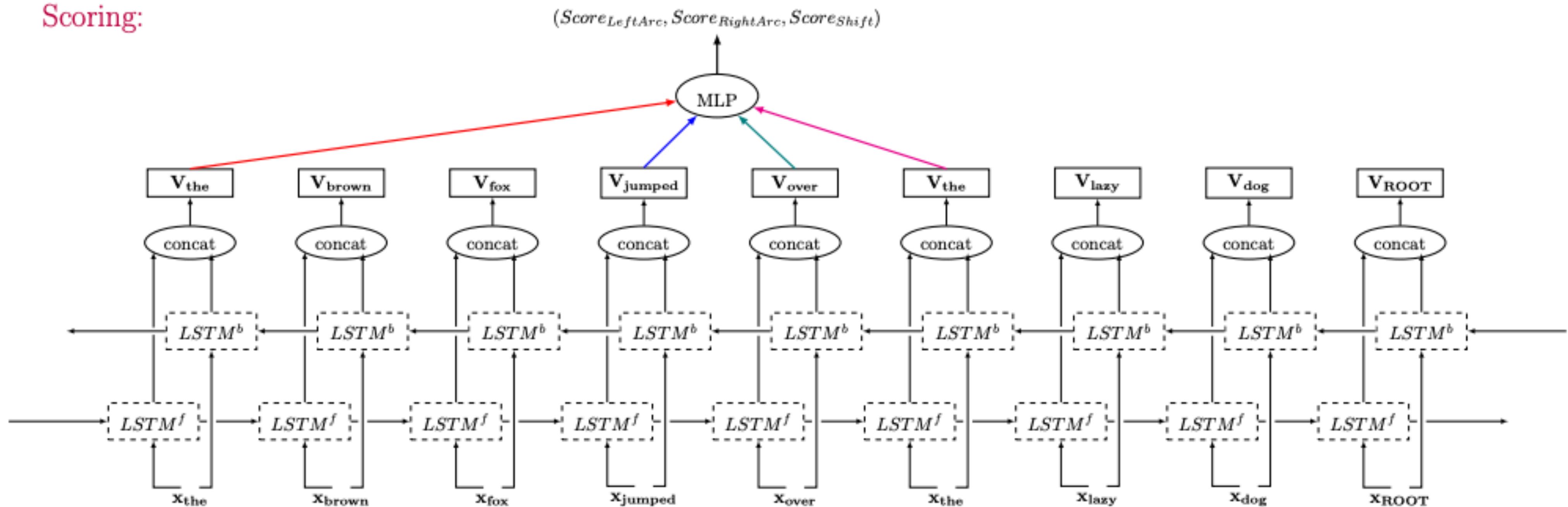


# What Actions to Choose?

Configuration:



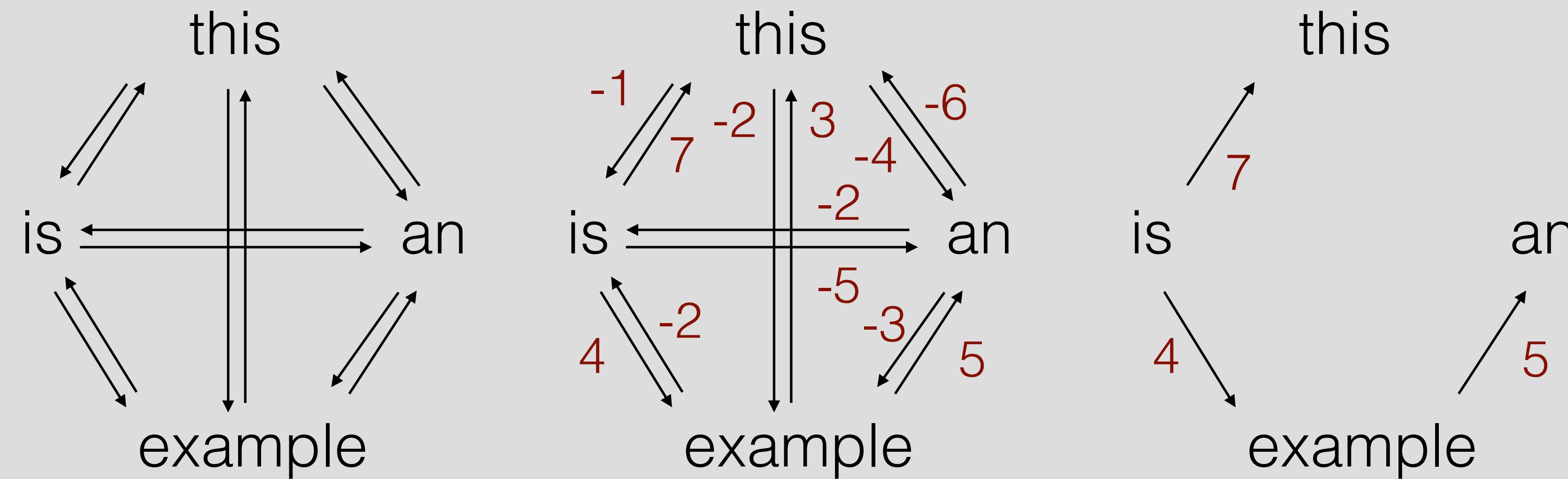
Scoring:



# Graph-based Parsing

# Graph-based Dependency Parsing

- Express sentence as fully connected directed graph
- Score each *edge* independently
- Find *maximal spanning tree*



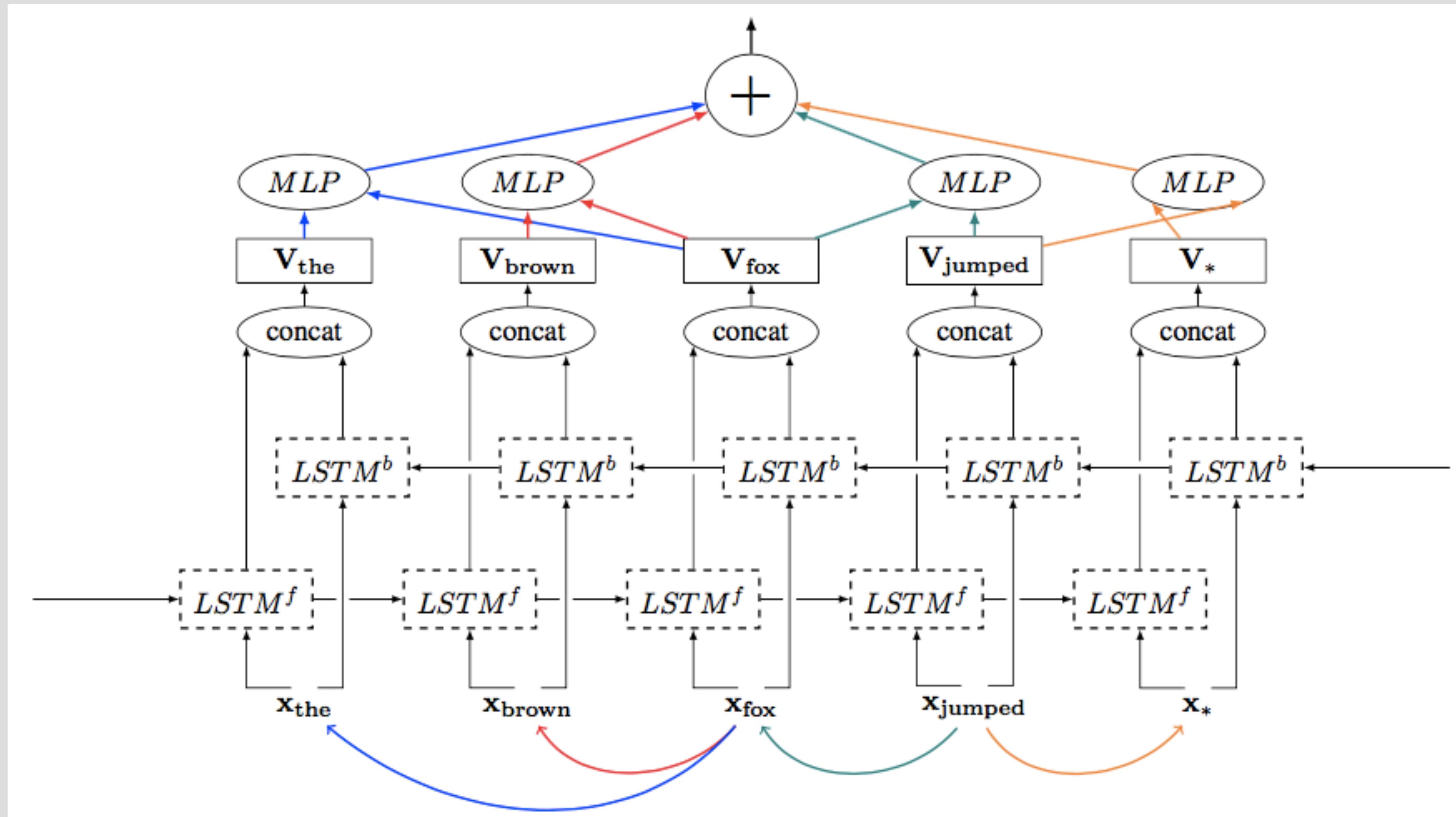
# How to Score the Edge?

- Words in sentence  $x$ , tree  $T$  is a collection of directed edges ( $\text{parent}(i)$ ,  $i$ ) for each word  $i$ 
  - Parsing = identify  $\text{parent}(i)$  for each word
  - Each word has exactly one parent
- Log-linear CRF (discriminative):  $P(T|x) = \exp\left(\sum_i w^\top f(i, \text{parent}(i), x)\right)$
- Example of a feature =  $I[\text{head}=\text{to} \ \& \ \text{modifier}=\text{house}]$



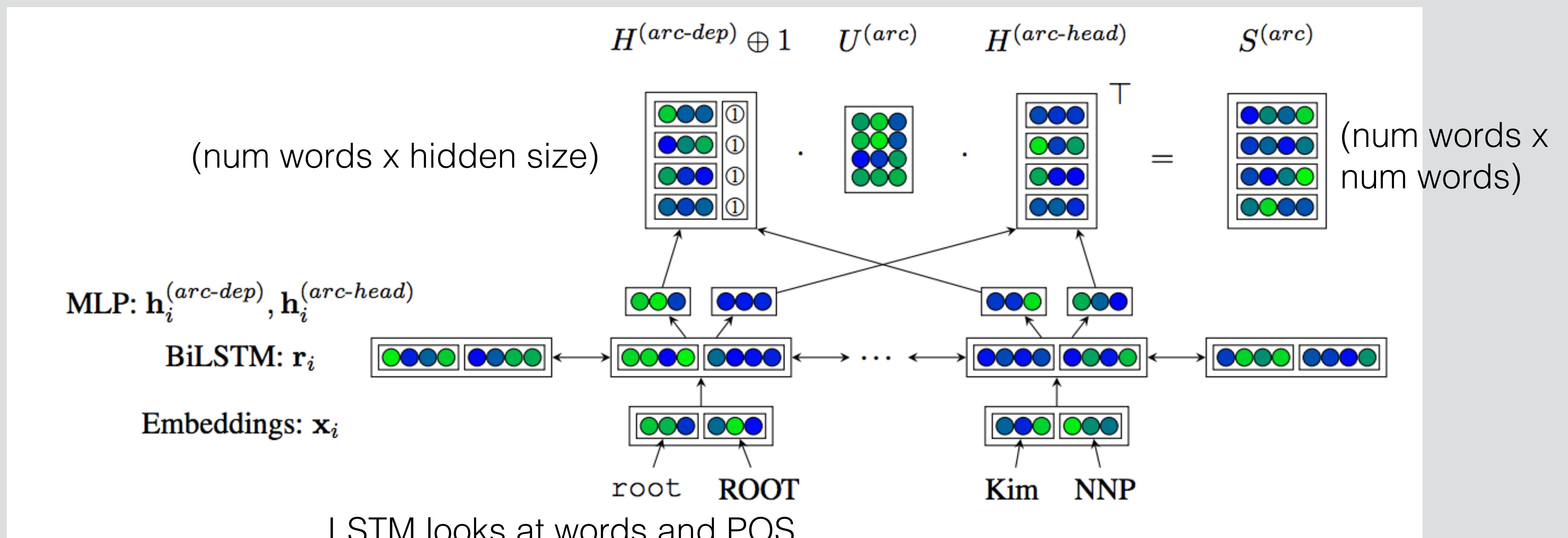
# Sequence Model Feature Extractors

(Kipperwasser and Goldberg 2016)



# Biaffine Neural Parsing

- Neural CRFs for dependency parsing: let  $c$  = LSTM embedding of  $i$ ,  $p$  = LSTM embedding of  $\text{parent}(i)$ .  $\text{score}(i, \text{parent}(i), x) = p^T U c$



# Graph-based vs. Transition Based

- **Transition-based**

- + Easily condition on infinite tree context (structured prediction)
- - Greedy search algorithm causes short-term mistakes

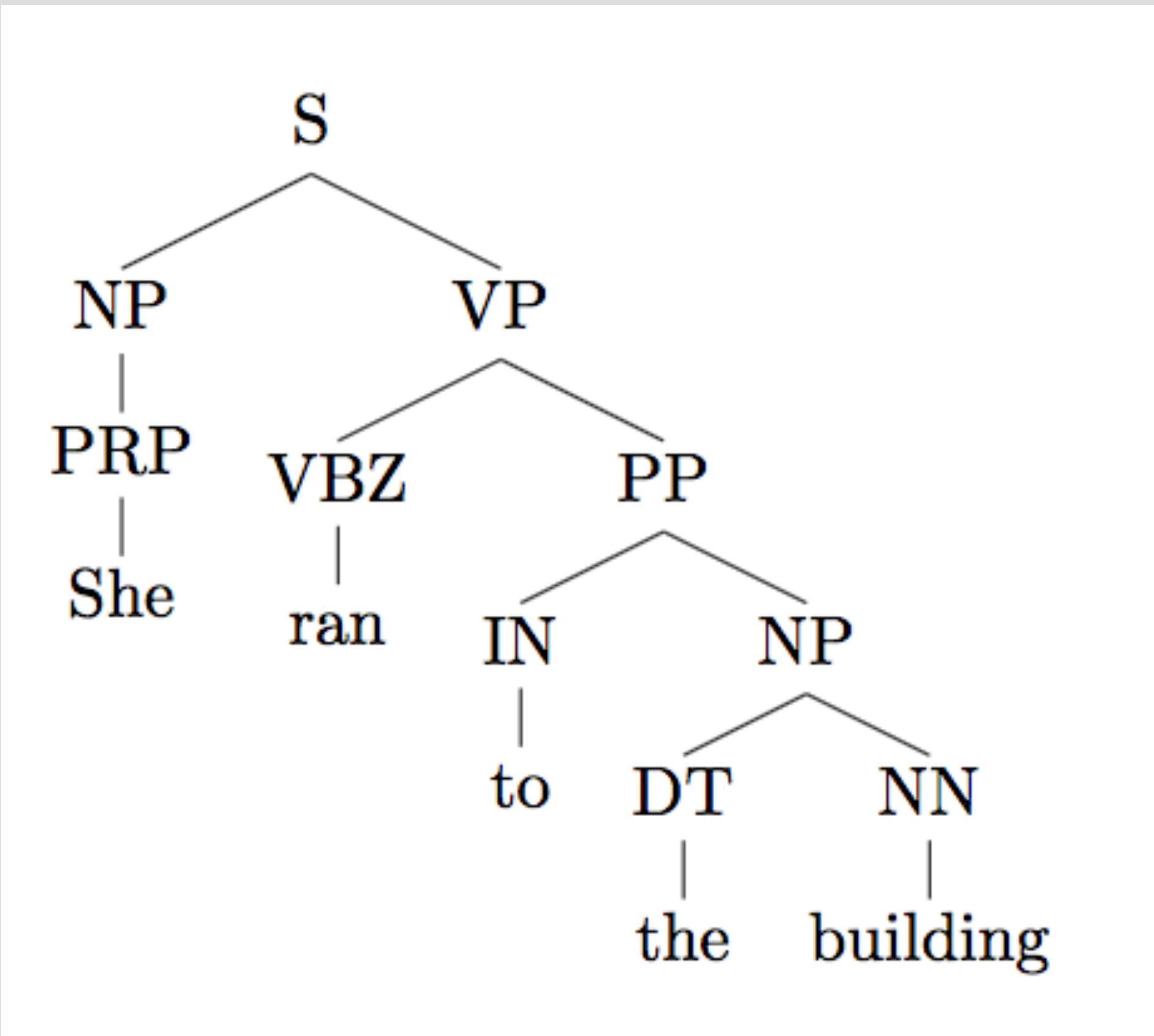
- **Graph-based**

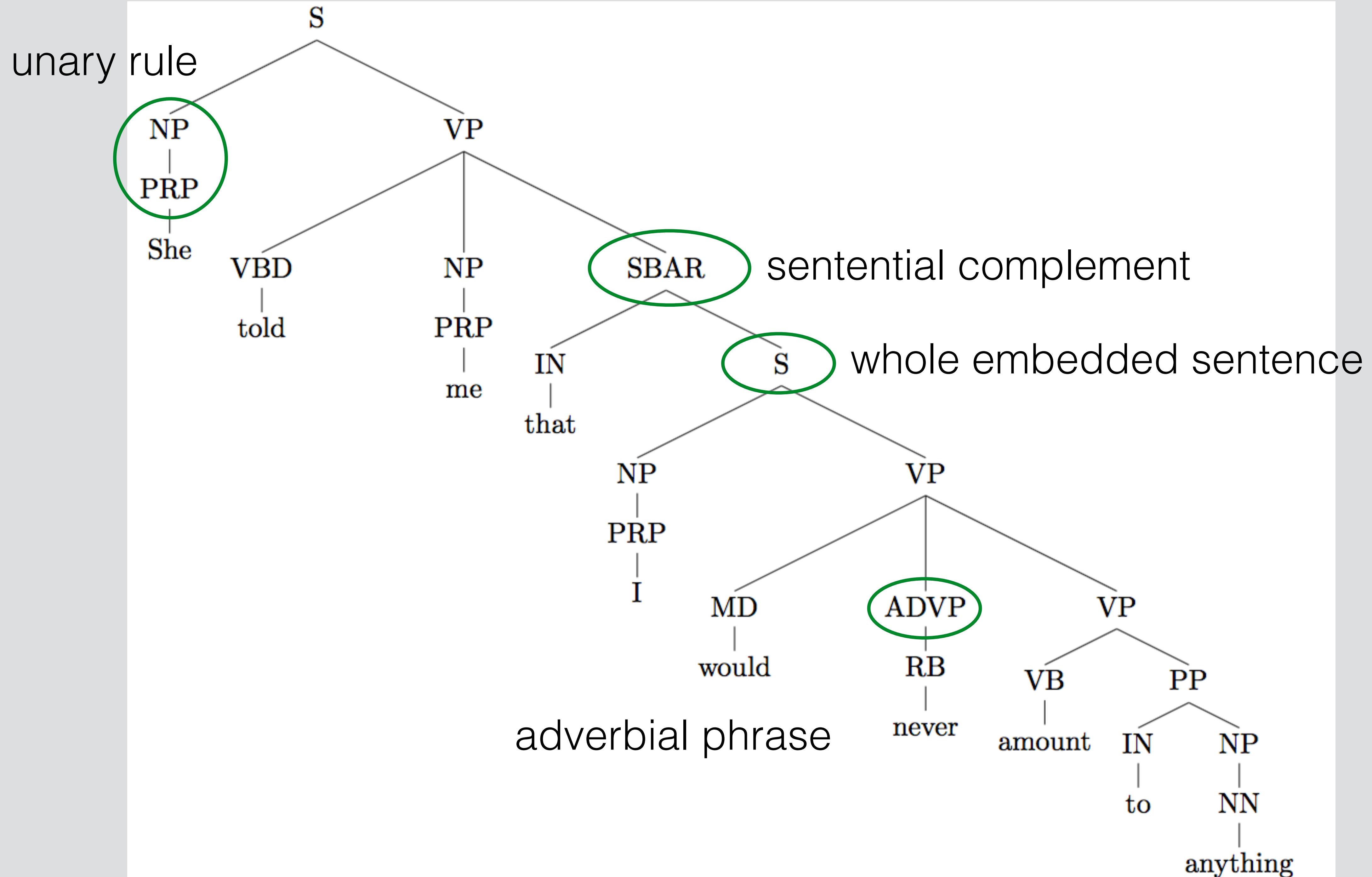
- + Can find exact best global solution via DP algorithm
- - Have to make local independence assumptions

# Syntactic Parsing 2: Constituency Parsing

# Constituency Parsing

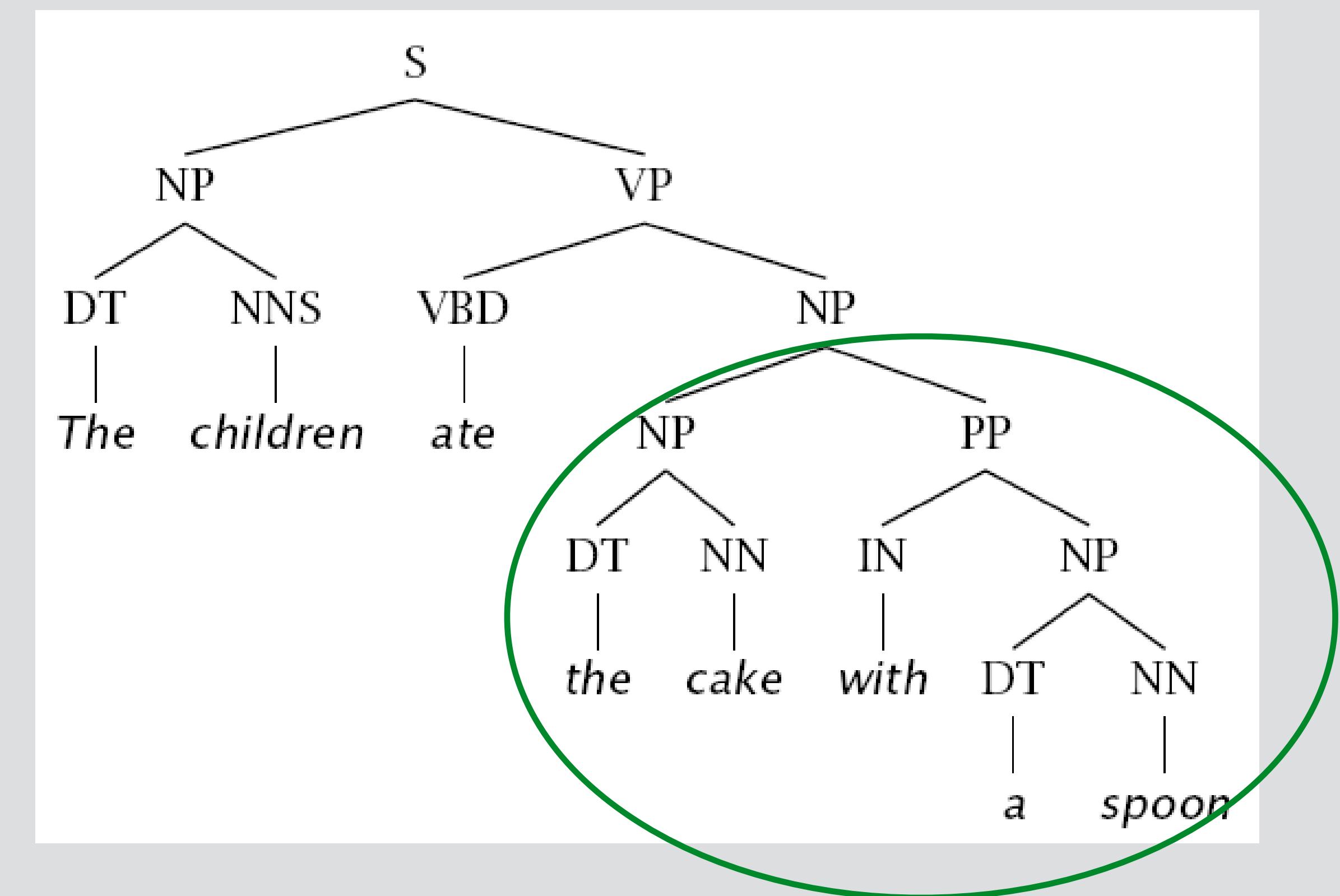
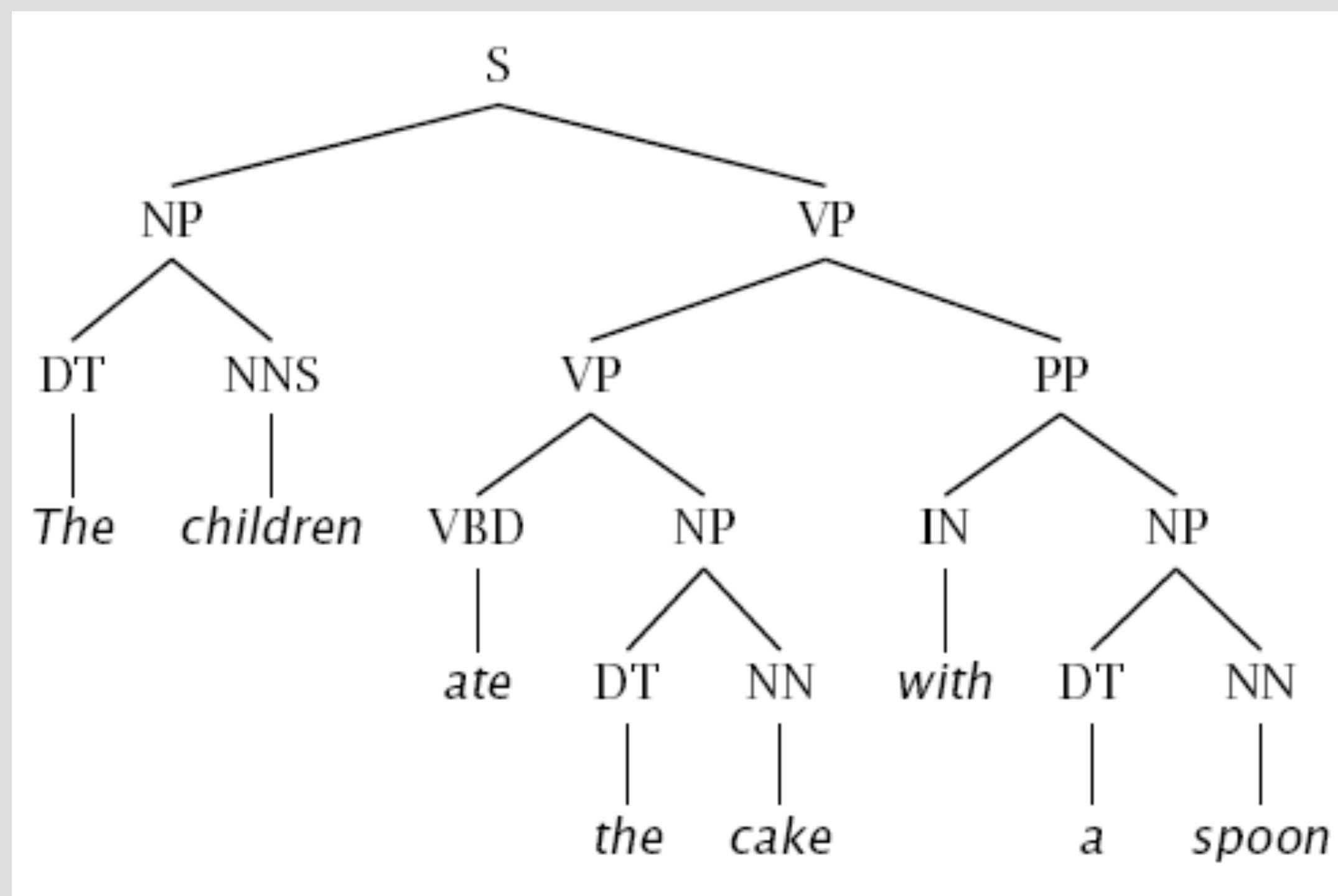
- Tree-structured syntactic analyses of sentences
- Common things: noun phrases, verb phrases, prepositional phrases
- Bottom layer is POS tags
- Examples will be in English. Constituency makes sense for a lot of languages but not all





# Challenges

- PP (Prepositional Phrase) attachment



same parse as “the cake with some icing”

# Challenges

- Modifier scope:
  - e.g., plastic cup holder
- Complement structure:
  - e.g., The students complained to the professor that they didn't understand
- Coordination scope:
  - e.g., The man picked up his hammer and saw (*verb or noun?*)

# Context-Free Grammars, CKY

# CFGs and PCFGs

## Grammar (CFG)

$\text{ROOT} \rightarrow S$

$S \rightarrow NP\ VP$

$NP \rightarrow DT\ NN$

$NP \rightarrow NN\ NNS$

$NP \rightarrow NP\ PP$

$VP \rightarrow VBP\ NP$

$VP \rightarrow VBP\ NP\ PP$

$PP \rightarrow IN\ NP$

## Lexicon

$NN \rightarrow \text{interest}$

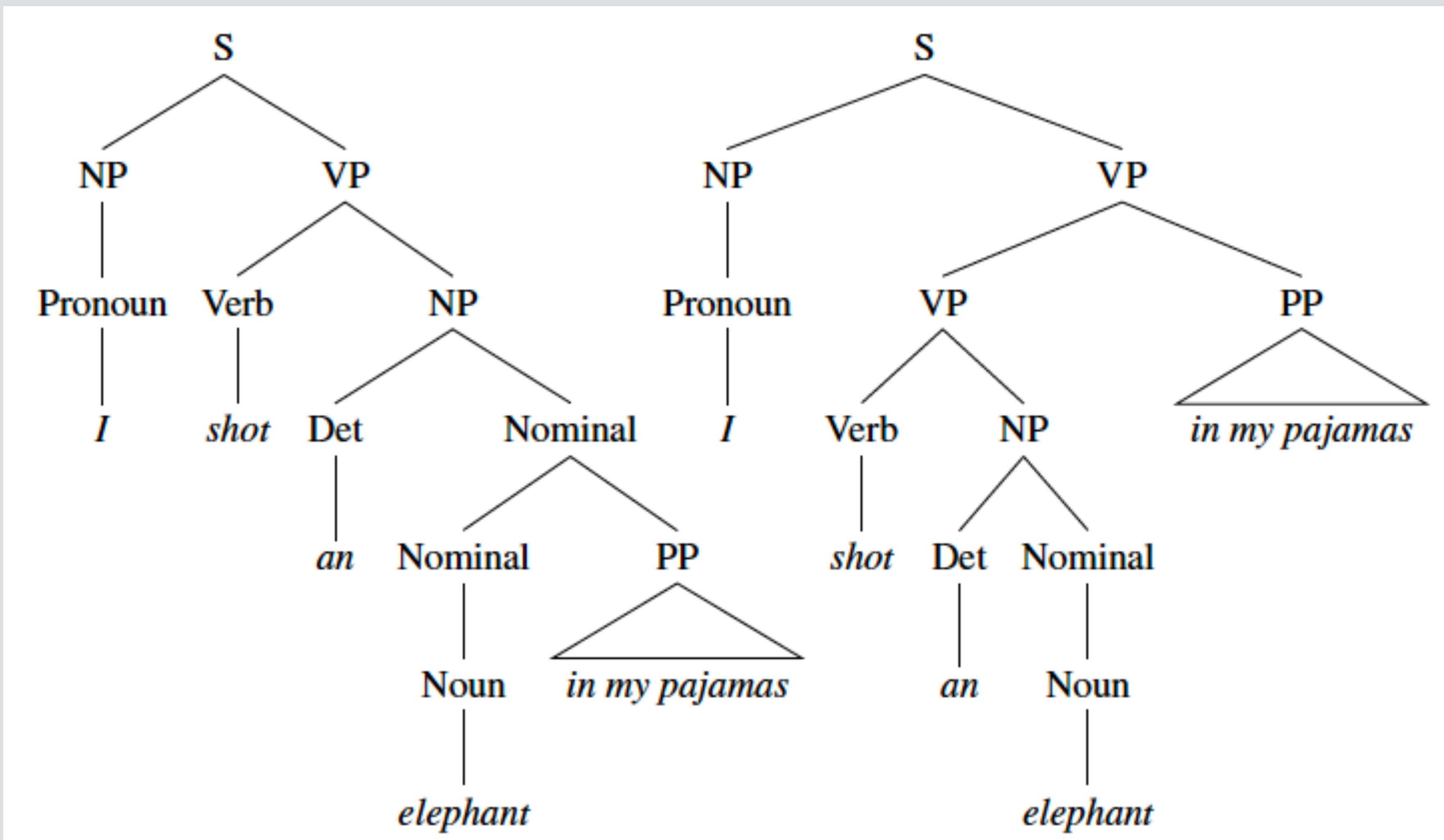
$NNS \rightarrow \text{raises}$

$VBP \rightarrow \text{interest}$

$VBZ \rightarrow \text{raises}$

- **Context-free grammar (CFG):** 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

Ambiguity under the same set of grammar and lexicon, e.g.,



Which one is more likely? Modeling the “probability”!

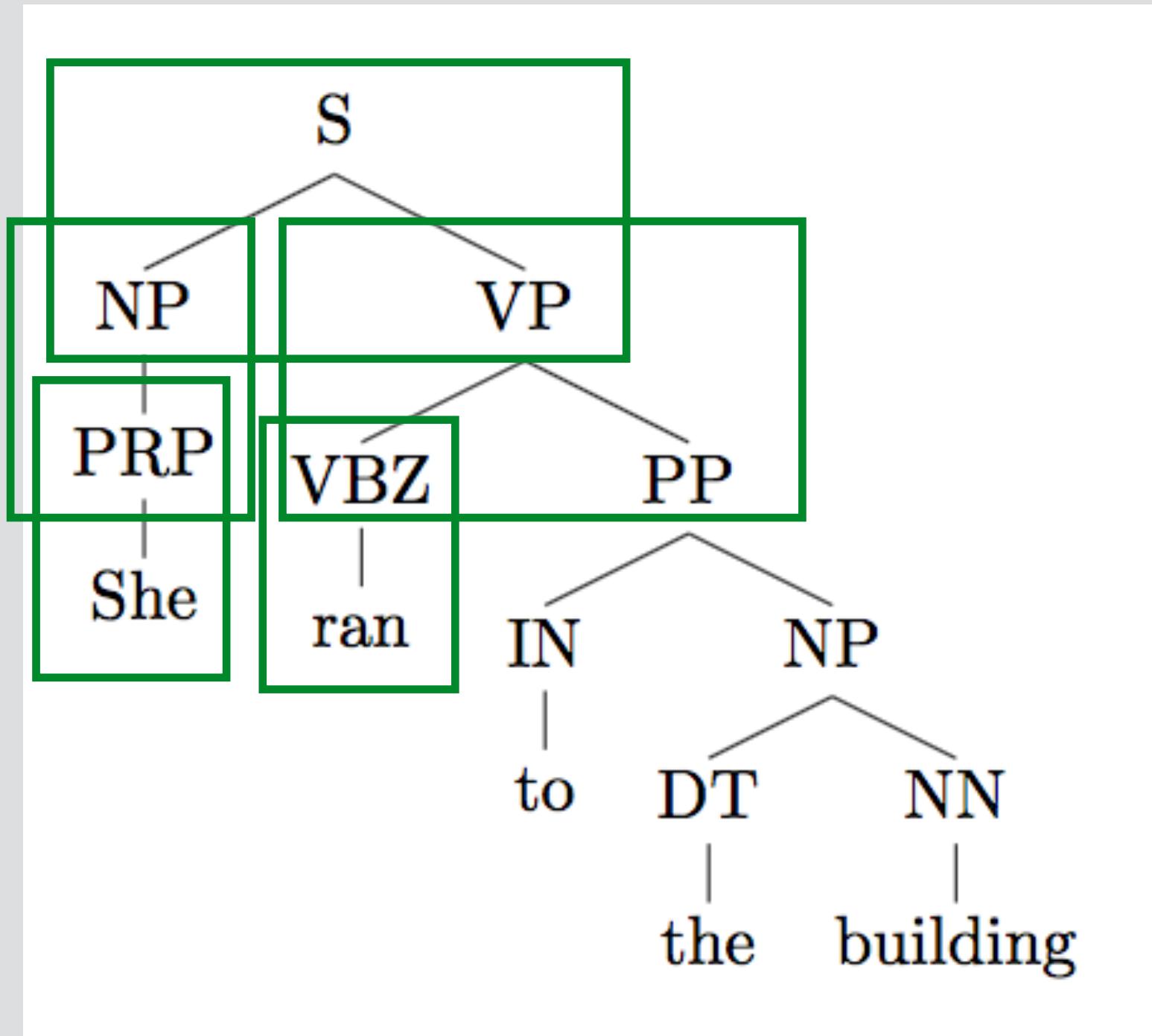
# CFGs and PCFGs

Grammar (CFG)				Lexicon	
ROOT → S	1.0	NP → NP PP	0.3	NN → interest	1.0
S → NP VP	1.0	VP → VBP NP	0.7	NNS → raises	1.0
NP → DT NN	0.2	VP → VBP NP PP	0.3	VBP → interest	1.0
NP → NN NNS	0.5	PP → IN NP	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

# PCFGs

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



$$\begin{aligned} P(T) = & P(S \rightarrow NP \ VP) \times P(NP \rightarrow PRP) \\ & \times P(PRP \rightarrow She) \times P(VP \rightarrow VBZ \ PP) \times P(VBZ \rightarrow ran) \\ & \times P(PP \rightarrow IN \ NP) \times P(IN \rightarrow to) \times P(NP \rightarrow DT \ NN) \\ & \times P(DT \rightarrow the) \times P(NN \rightarrow building) \end{aligned}$$

**Learning** PCFG: how do I know the probability for each derivation?

**Using** PCFG: how do I apply PCFG to parse a new sentence into a constituency parse tree?

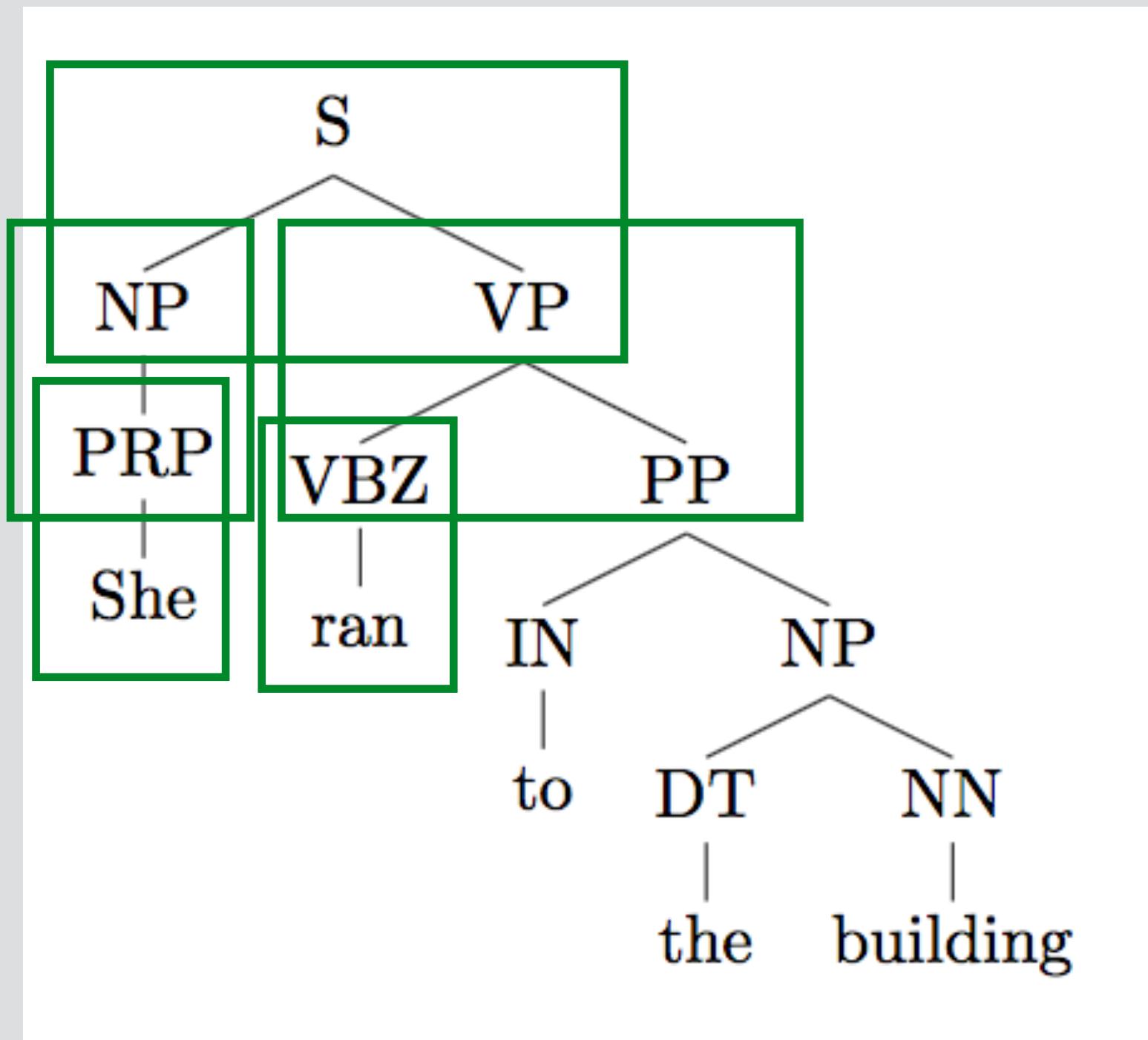
# Learning a PCFG

- Given a set of example trees (from a treebank), the underlying CFG can simply be all production rules seen in the corpus
- Learning goal: estimate the probability of each production rule
  - e.g.,  $P(NP \rightarrow PRP)$  or  $P(PRP | NP)$  describes the probability of applying the rule “NP → PRP” when the current (parent node) is “NP”
- Maximum likelihood estimation (similar as HMMs’)

$$P(\alpha - > \beta) = \frac{Count(\alpha - > \beta)}{Count(\alpha)}$$

# PCFGs

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



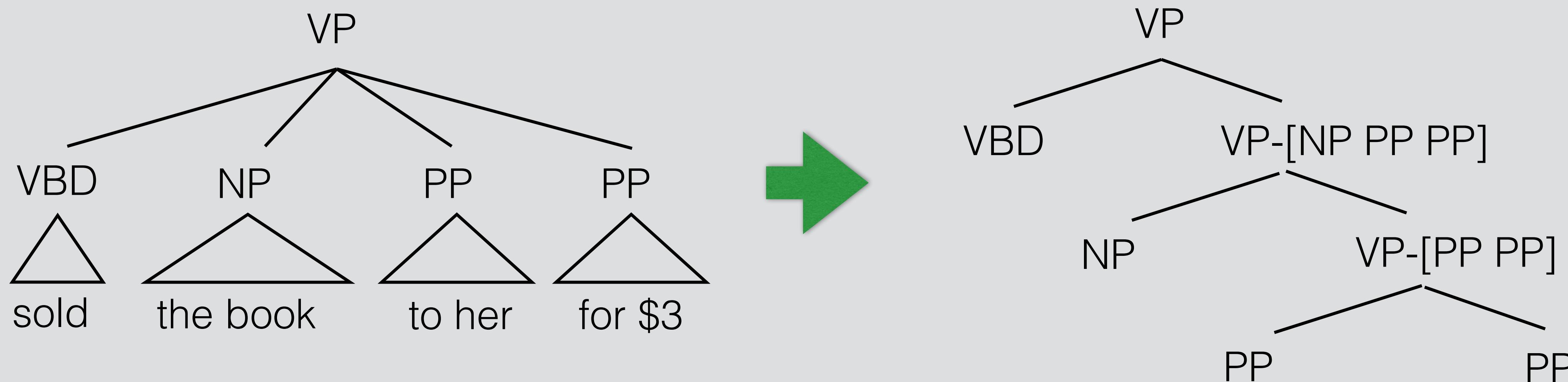
$$\begin{aligned} P(T) = & P(S \rightarrow NPVP) \times P(NP \rightarrow PRP) \\ & \times P(PRP \rightarrow She) \times P(VP \rightarrow VBZ\ PP) \times P(VBZ \rightarrow ran) \\ & \times P(PP \rightarrow INNP) \times P(IN \rightarrow to) \times P(NP \rightarrow DT\ NN) \\ & \times P(DT \rightarrow the) \times P(NN \rightarrow building) \end{aligned}$$

**Learning** PCFG: MLE by counting + normalization, similar to learning HMMs

**Using** PCFG:  $\text{argmax}_T P(T)$

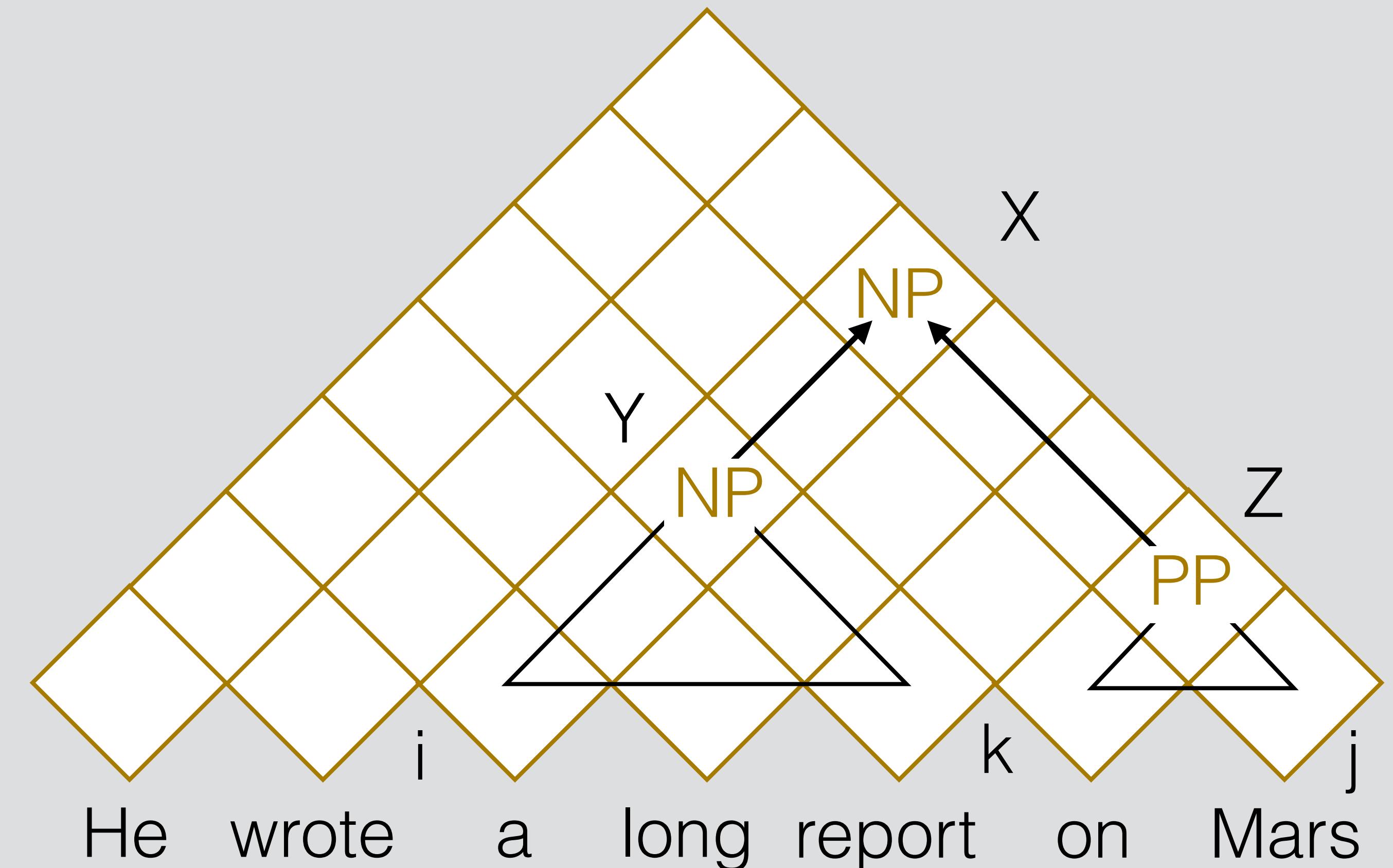
# Binarization

- To parse efficiently, we assume PCFGs in **Chomsky Normal Form (CNF)**, i.e., having either two (nonterminal) subtree children or one (terminal) leaf value.



# Inference: CKY

- Find  $\text{argmax}_T P(T)$
- Dynamic programming: chart maintains the *best* way of building **symbol X** over **span (i, j)**
- CKY is like Viterbi

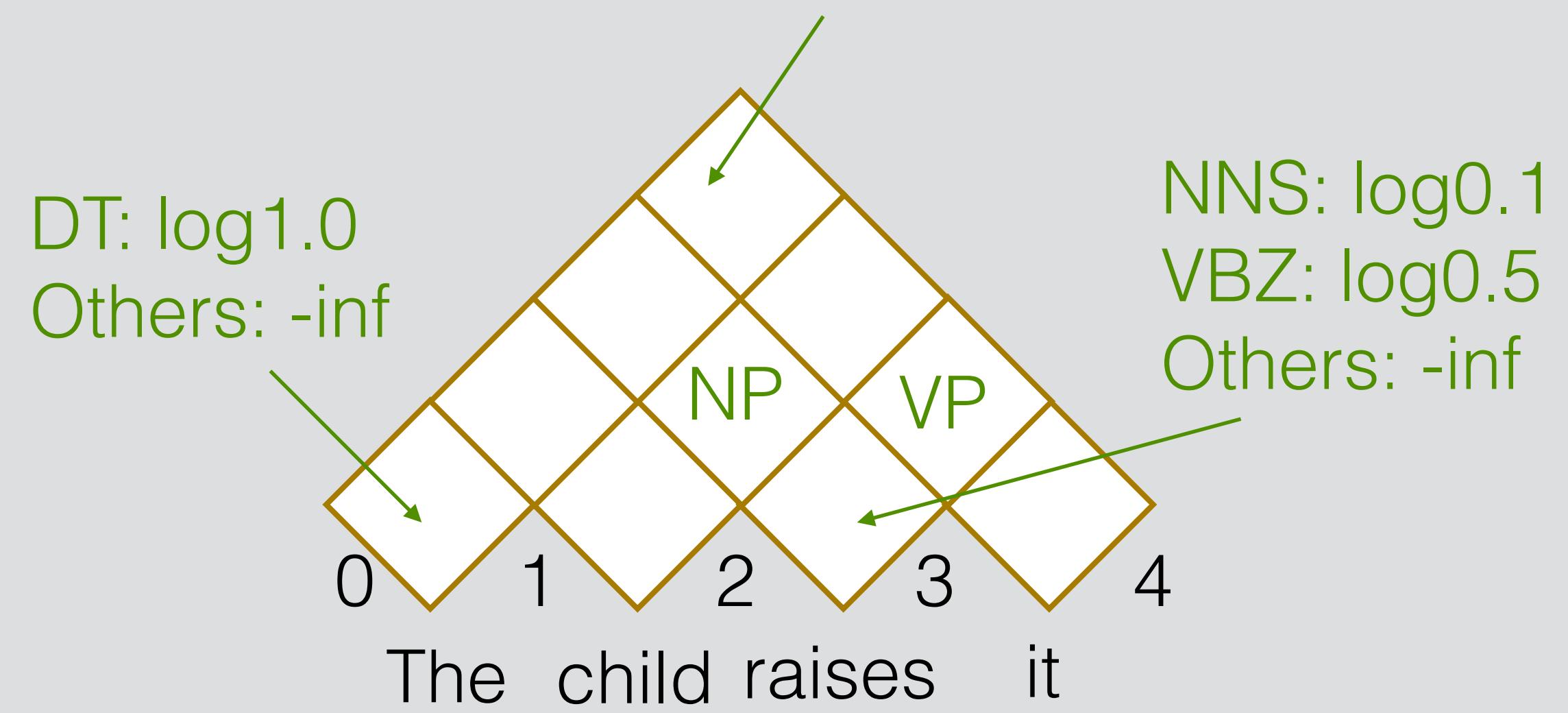


# CKY

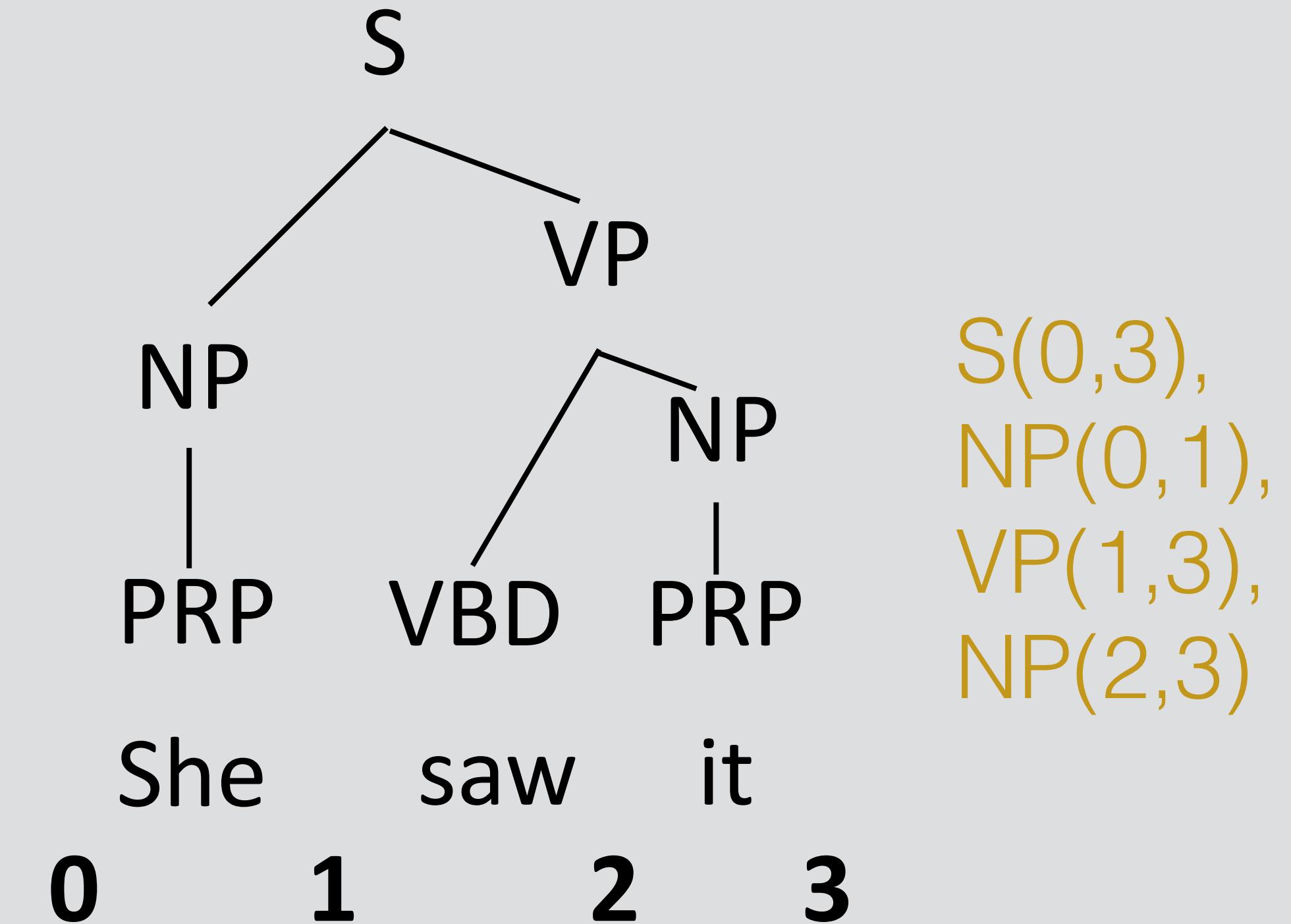
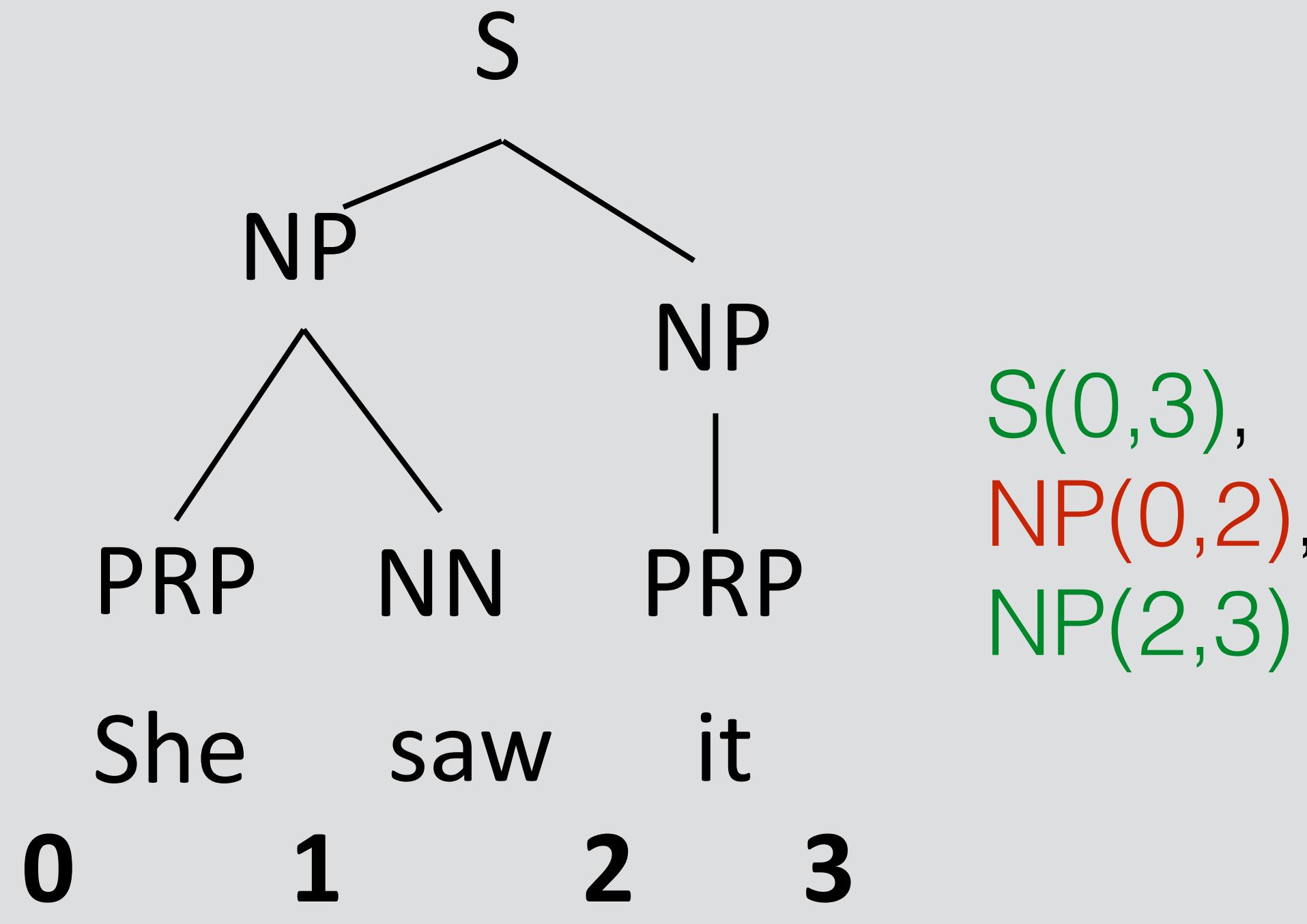
- Chart:  $T[i,j,X] = \text{best score for } X \text{ (a POS tag) over } (i, j)$
- Base:  $T[i,i+1,X] = \log P(X \rightarrow w_i)$
- Loop over all split points  $k$ , apply rules  $X \rightarrow Y Z$  to build  $X$  in every possible way
- Recurrence:  

$$T[i,j,X] = \max_{r: X \rightarrow X_1 X_2} \max_k T[i,k,X_1] + T[k,j,X_2] + \log P(X \rightarrow X_1 X_2)$$
- Runtime:  $O(n^3G)$   $G = \text{grammar constant}$

$T[0,4,X]$  - how many ways of parsing in the end & how likely is each of them?



# Parser Evaluation



- Precision: number of correct brackets / num pred brackets =  $2/3$
- Recall: number of correct brackets / num of gold brackets =  $2/4$
- F1: harmonic mean of precision and recall =  $(1/2 * ((2/4)^{-1} + (2/3)^{-1}))^{-1}$   
 $= 0.57$

# Results

- Standard dataset for English: Penn Treebank (Marcus et al., 1993)
  - Evaluation: F1 over labeled constituents of the sentence
- Vanilla PCFG: ~75 F1
- Best performance: ~96 F1

<https://github.com/sebastianruder/NLP-progress>

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

# Dependency vs. Constituency

- Dependency is often more useful in practice (models predicate argument structure)
- Slightly different representational choices:
  - PP attachment is better modeled under dependency
  - Coordination is better modeled under constituency
- Dependency parsers are easier to build: no “grammar engineering”, no unaries, easier to get structured discriminative models working well
- Dependency parsers are usually faster
- Dependencies are more universal cross-lingually