

TTIC 31210:
Advanced Natural Language Processing

Kevin Gimpel
Spring 2017

Lecture 15:
Structured Prediction

- No class Monday May 29 (Memorial Day)
- Final class is Wednesday May 31

- Assignment 3 has been posted, due Thursday June 1
- Final project report due Friday, June 9

Modeling, Inference, Learning

inference: solve argmax

modeling: define score function

$$\operatorname{classify}(x, \theta) = \operatorname{argmax}_y \operatorname{score}(x, y, \theta)$$

learning: choose θ

Structured Prediction:

size of output space is exponential in size of input
or is unbounded (e.g., machine translation)
(we can't just enumerate all possible outputs)

- 2 categories of structured prediction:
score-based and search-based

Score-Based Structured Prediction

- focus on defining the score function of the structured input/output pair:

$$\text{score}(\boldsymbol{x}, \boldsymbol{y}, \boldsymbol{\theta})$$

- cleanly separates score function, inference algorithm, and training loss

Modeling in Score-Based SP

- define score as a sum or product over “parts” of the structured input/output pair:

$$\text{score}(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) = \sum_{\langle \mathbf{x}_r, \mathbf{y}_r \rangle \in \text{parts}(\mathbf{x}, \mathbf{y})} \text{score}(\mathbf{x}_r, \mathbf{y}_r, \boldsymbol{\theta})$$

$$\text{score}(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) = \prod_{\langle \mathbf{x}_r, \mathbf{y}_r \rangle \in \text{parts}(\mathbf{x}, \mathbf{y})} \text{score}(\mathbf{x}_r, \mathbf{y}_r, \boldsymbol{\theta})$$

Parts Functions in Score-Based SP

- for an HMM:

$$\text{parts}_{\text{HMM}}(\mathbf{x}, \mathbf{y}) = \{\langle x_t, y_t \rangle\}_{t=1}^T \cup \{\langle \emptyset, y_{t-1:t} \rangle\}_{t=1}^T$$

- each word-label pair forms a part, and each label bigram forms a part

Parts Functions in Score-Based SP

- for a linear chain CRF:

$$\text{parts}_{\text{chainCRF}}(\mathbf{x}, \mathbf{y}) = \{\langle \mathbf{x}, y_{t-1:t} \rangle\}_{t=1}^T$$

- each label bigram forms a part (each of which includes entire input!)

Parts Functions in Score-Based SP

- for a PCFG:

$$\text{parts}_{\text{PCFG}}(\mathbf{x}, \mathbf{y}) = \left(\bigcup_{(y \rightarrow x) \in \mathbf{y}} \langle x, y \rangle \right) \cup \left(\bigcup_{(y \rightarrow y_1, y_2) \in \mathbf{y}} \langle \emptyset, \langle y, y_1, y_2 \rangle \rangle \right)$$

- each context-free grammar rule forms a part

Parts Functions in Score-Based SP

- for an arc-factored dependency parser:

$$\text{parts}_{\text{arcdep}}(\mathbf{x}, \mathbf{y}) = \{\langle \mathbf{x}, y_t \rangle\}_{t=1}^T$$

where y_t is the index of the parent of x_t in \mathbf{y}

- each dependency arc forms a part

Inference in Score-Based SP

- inference algorithms are defined based on decomposition of score into parts

$$\text{score}(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) = \sum_{\langle \mathbf{x}_r, \mathbf{y}_r \rangle \in \text{parts}(\mathbf{x}, \mathbf{y})} \text{score}(\mathbf{x}_r, \mathbf{y}_r, \boldsymbol{\theta})$$

- smaller parts = easier to define efficient exact inference algorithms

Inference Algorithms for Score-Based SP

- exact inference algorithms are often based on **dynamic programming**

Dynamic Programming (DP)

- a class of algorithms that break problems into smaller pieces and reuse solutions for pieces
 - applicable if problem has certain properties (**optimal substructure** and **overlapping sub-problems**)
- in NLP, we use DP to iterate over exponentially-large output spaces in polynomial time
 - Viterbi and forward/backward for HMMs
 - CKY for PCFGs
 - Eisner algorithm for (arc-factored) dependency parsing

Viterbi Algorithm

- recursive equations + memoization:

base case:

returns score of sequence starting with label y for first word



$$V(1, y) = \text{score}(x_1, y) \text{score}(\emptyset, \langle, \langle s \rangle, y \rangle)$$

$$V(t, y) = \max_{y' \in \mathcal{L}} (\text{score}(x_t, y) \text{score}(\emptyset, \langle y', y \rangle) V(t - 1, y'))$$



recursive case:

computes score of max-scoring label sequence that ends with label y at position t

final value is in: $V(|\mathbf{x}| + 1, \langle /s \rangle)$

Viterbi Algorithm

- space and time complexity?
- can be read off from the recursive equations:

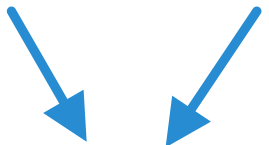
space complexity:

size of memoization table, which is # of unique indices of recursive equations

length of
sentence

*

number
of labels


$$V(t, y) = \max_{y' \in \mathcal{L}} (\text{score}(x_t, y) \text{ score}(\emptyset, \langle y', y \rangle) V(t - 1, y'))$$

so, space complexity is $O(|x| |L|)$

Viterbi Algorithm

- space and **time** complexity?
- can be read off from the recursive equations:

time complexity:

size of memoization table * complexity of computing each entry

length of sentence * number of labels * each entry requires iterating through the labels

$$V(t, y) = \max_{y' \in \mathcal{L}} (\text{score}(x_t, y) \text{ score}(\emptyset, \langle y', y \rangle) V(t - 1, y'))$$

so, time complexity is $O(|x| |\mathcal{L}| |\mathcal{L}|) = O(|x| |\mathcal{L}|^2)$

Feature Locality

- **feature locality**: how big are the parts?
- for efficient inference (w/ DP or other methods), we need to be mindful of this
- parts can be arbitrarily big in terms of input, but not in terms of *output*!
- HMM parts are small in both the input and output (only two pieces at a time)

Learning with Score-Based SP: Empirical Risk Minimization

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \sum_{\langle \mathbf{x}, \mathbf{y} \rangle \in \mathcal{D}} \operatorname{cost}(\mathbf{y}, \operatorname{predict}(\mathbf{x}, \theta))$$

$$\operatorname{predict}(\mathbf{x}, \theta) = \underset{\mathbf{y}}{\operatorname{argmax}} \operatorname{score}(\mathbf{x}, \mathbf{y}, \theta)$$

Cost Functions

- **cost function**: how different are these two structures?

$$\text{cost} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_{\geq 0}$$

- typically used to compare predicted structure to gold standard
- should reflect evaluation metric for task

- usual conventions: $\text{cost}(\mathbf{y}, \mathbf{y}) = 0$

$$\text{cost}(\mathbf{y}, \mathbf{y}') = \text{cost}(\mathbf{y}', \mathbf{y})$$

Cost Functions

- for classification, we used:

$$\text{cost}(y, y') = \mathbb{I}[y \neq y']$$

- how about for sequences?

- “Hamming cost”: $\text{cost}(\mathbf{y}, \mathbf{y}') = \sum_{t=1}^{|\mathbf{y}|} \mathbb{I}[y_t \neq y'_t]$

- “0-1 cost”: $\text{cost}(\mathbf{y}, \mathbf{y}') = \mathbb{I}[\mathbf{y} \neq \mathbf{y}']$

Empirical Risk Minimization

$$\hat{\theta} = \operatorname{argmin}_{\theta} \sum_{\langle \mathbf{x}, \mathbf{y} \rangle \in \mathcal{D}} \operatorname{cost}(\mathbf{y}, \operatorname{predict}(\mathbf{x}, \theta))$$

$$\operatorname{predict}(\mathbf{x}, \theta) = \operatorname{argmax}_{\mathbf{y}} \operatorname{score}(\mathbf{x}, \mathbf{y}, \theta)$$

- this is intractable so we typically minimize a **surrogate loss function** instead

Loss Functions for Score-Based SP

name	loss	where used
cost ("0-1")	$\text{cost}(\mathbf{y}, \text{predict}(\mathbf{x}, \boldsymbol{\theta}))$	
percep- tron		
hinge		

Loss Functions for Score-Based SP

name	loss	where used
cost ("0-1")	$\text{cost}(\mathbf{y}, \text{predict}(\mathbf{x}, \boldsymbol{\theta}))$	MERT (Och, 2003)
percep- tron	$-\text{score}(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) + \max_{\mathbf{y}'} \text{score}(\mathbf{x}, \mathbf{y}', \boldsymbol{\theta})$	structured perceptron (Collins, 2002)
hinge	$-\text{score}(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) + \max_{\mathbf{y}'} (\text{score}(\mathbf{x}, \mathbf{y}', \boldsymbol{\theta}) + \text{cost}(\mathbf{y}, \mathbf{y}'))$	structured SVMs (Taskar et al., <i>inter alia</i>)

Loss Functions for Score-Based SP

name	loss	where used
cost ("0-1")	$\text{cost}(\mathbf{y}, \text{predict}(\mathbf{x}, \boldsymbol{\theta}))$	MERT (Och, 2003)
percep- tron	$-\text{score}(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) + \max_{\mathbf{y}'} \text{score}(\mathbf{x}, \mathbf{y}', \boldsymbol{\theta})$	structured perceptron (Collins, 2002)
hinge	$-\text{score}(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) + \max_{\mathbf{y}'} (\text{score}(\mathbf{x}, \mathbf{y}', \boldsymbol{\theta}) + \text{cost}(\mathbf{y}, \mathbf{y}'))$	structured SVMs (Taskar et al., <i>inter alia</i>)
log	$-\text{score}(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) + \log \sum_{\mathbf{y}'} \exp \{ \text{score}(\mathbf{x}, \mathbf{y}', \boldsymbol{\theta}) \}$	CRFs (Lafferty et al., 2001)

Loss Functions for Score-Based SP

name	loss	where used
cost ("0-1")	$\text{cost}(\mathbf{y}, \text{predict}(\mathbf{x}, \boldsymbol{\theta}))$	MERT (Och, 2003)
percep- tron	$-\text{score}(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) + \max_{\mathbf{y}'} \text{score}(\mathbf{x}, \mathbf{y}', \boldsymbol{\theta})$	structured perceptron (Collins, 2002)
hinge	$-\text{score}(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) + \max_{\mathbf{y}'} (\text{score}(\mathbf{x}, \mathbf{y}', \boldsymbol{\theta}) + \text{cost}(\mathbf{y}, \mathbf{y}'))$	structured SVMs (Taskar et al., <i>inter alia</i>)
log	$-\text{score}(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) + \log \sum_{\mathbf{y}'} \exp \{ \text{score}(\mathbf{x}, \mathbf{y}', \boldsymbol{\theta}) \}$	CRFs (Lafferty et al., 2001)
softma- x- margin	$-\text{score}(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) + \log \sum_{\mathbf{y}'} \exp \{ \text{score}(\mathbf{x}, \mathbf{y}', \boldsymbol{\theta}) + \text{cost}(\mathbf{y}, \mathbf{y}') \}$	Povey et al. (2008), Gimpel & Smith (2010)

Perceptron



Conditional Likelihood



Max-Margin

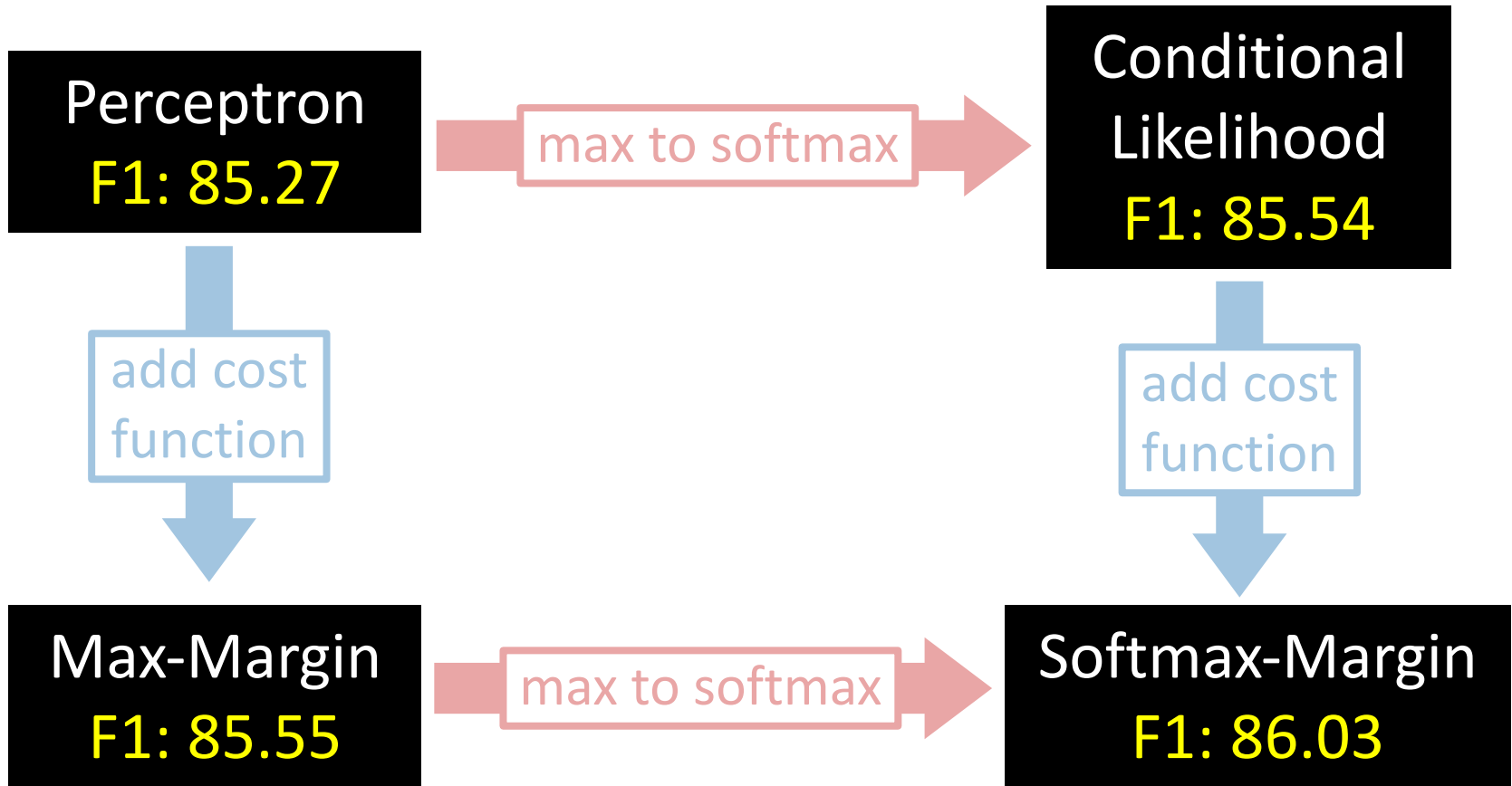


Softmax-Margin

$$-\text{score}(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) + \log \sum_{\mathbf{y}'} \exp \{ \text{score}(\mathbf{x}, \mathbf{y}', \boldsymbol{\theta}) + \text{cost}(\mathbf{y}, \mathbf{y}') \}$$

Results: Named Entity Recognition

(Gimpel & Smith, 2010)



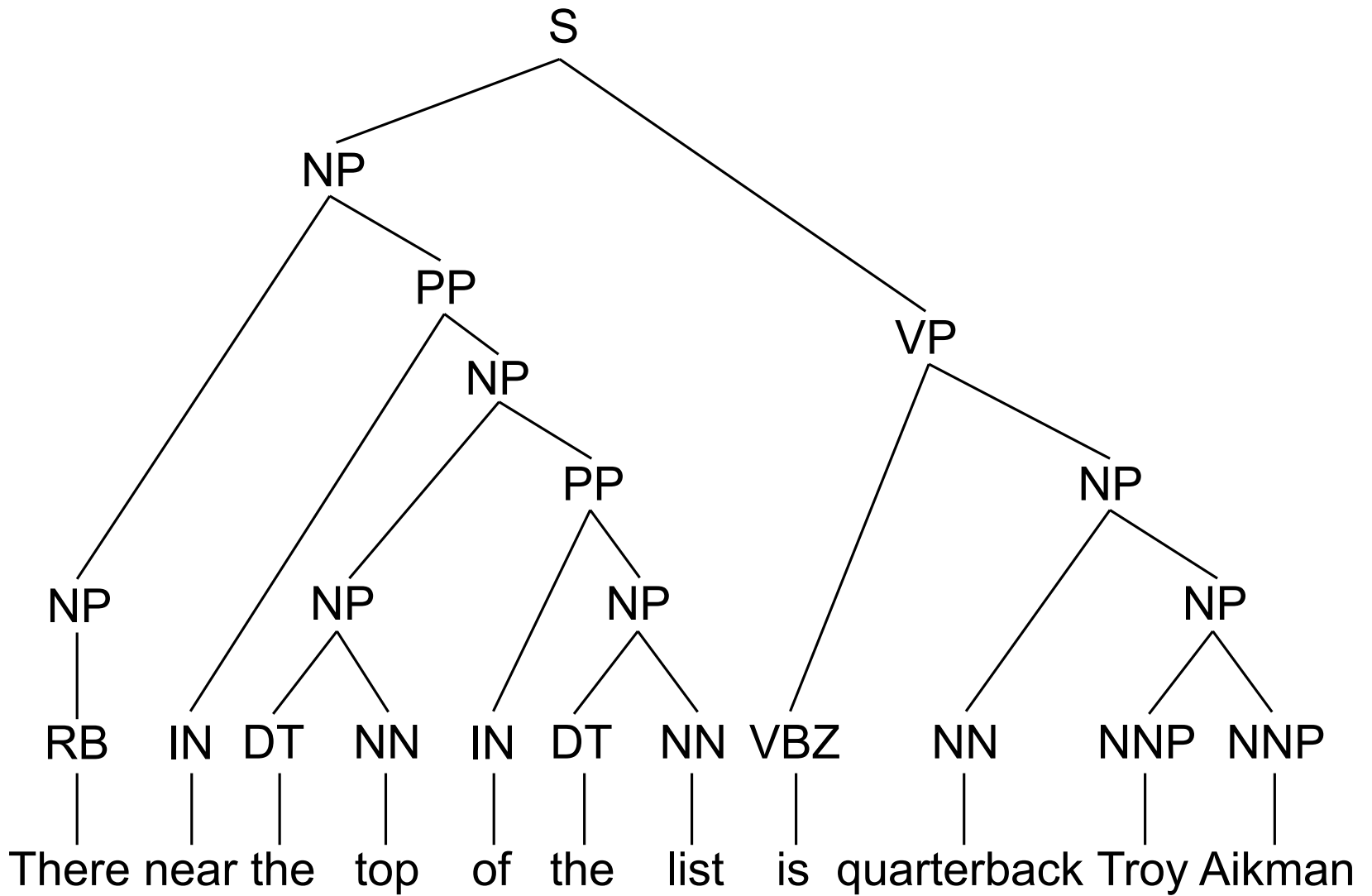
Inference Algorithms for Score-Based SP

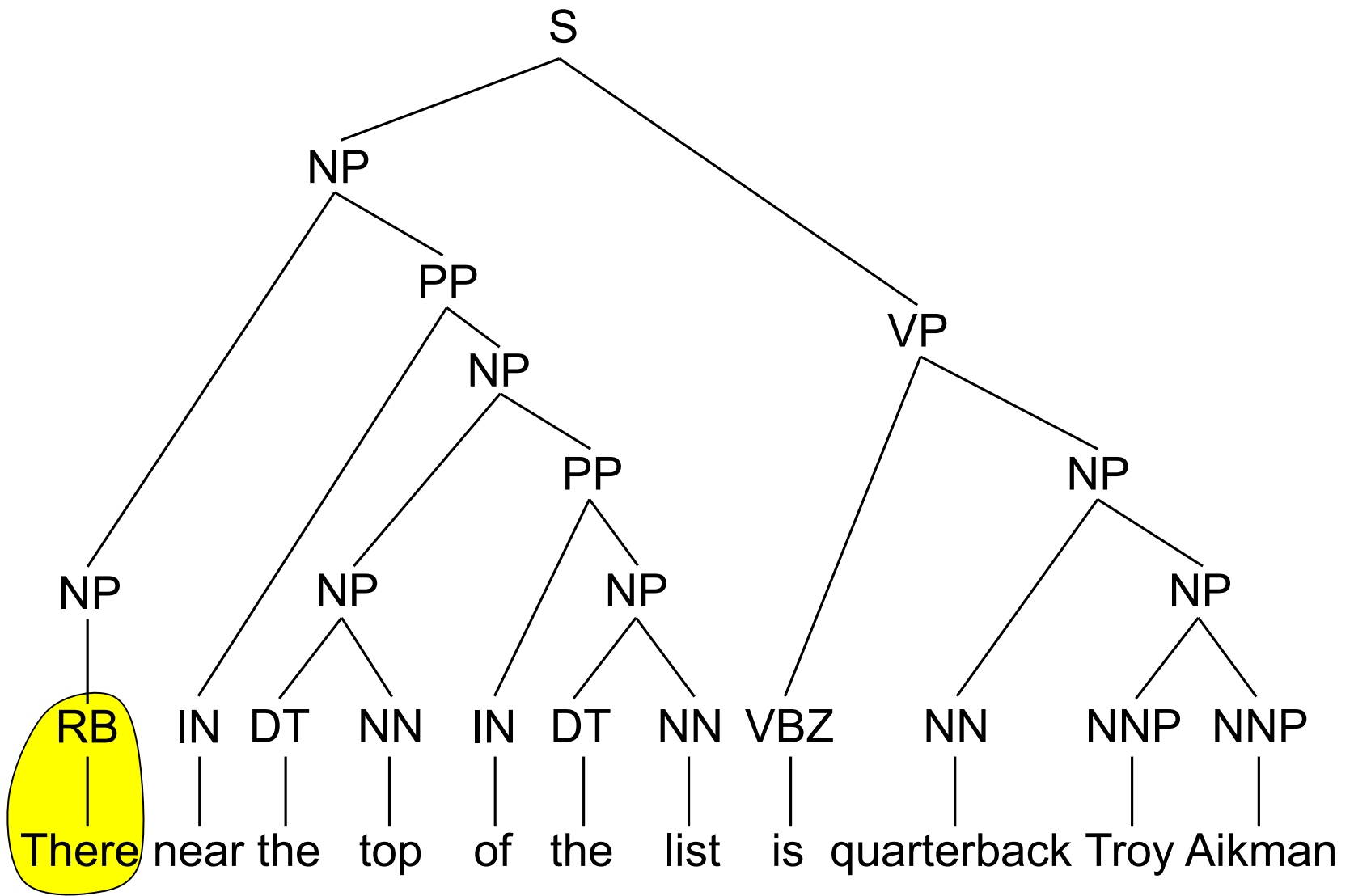
- dynamic programming
 - exact, but parts must be small for efficiency
- dynamic programming + “cube pruning”
 - permits approximate incorporation of large parts (“non-local features”) while still using dynamic program backbone
- integer linear programming

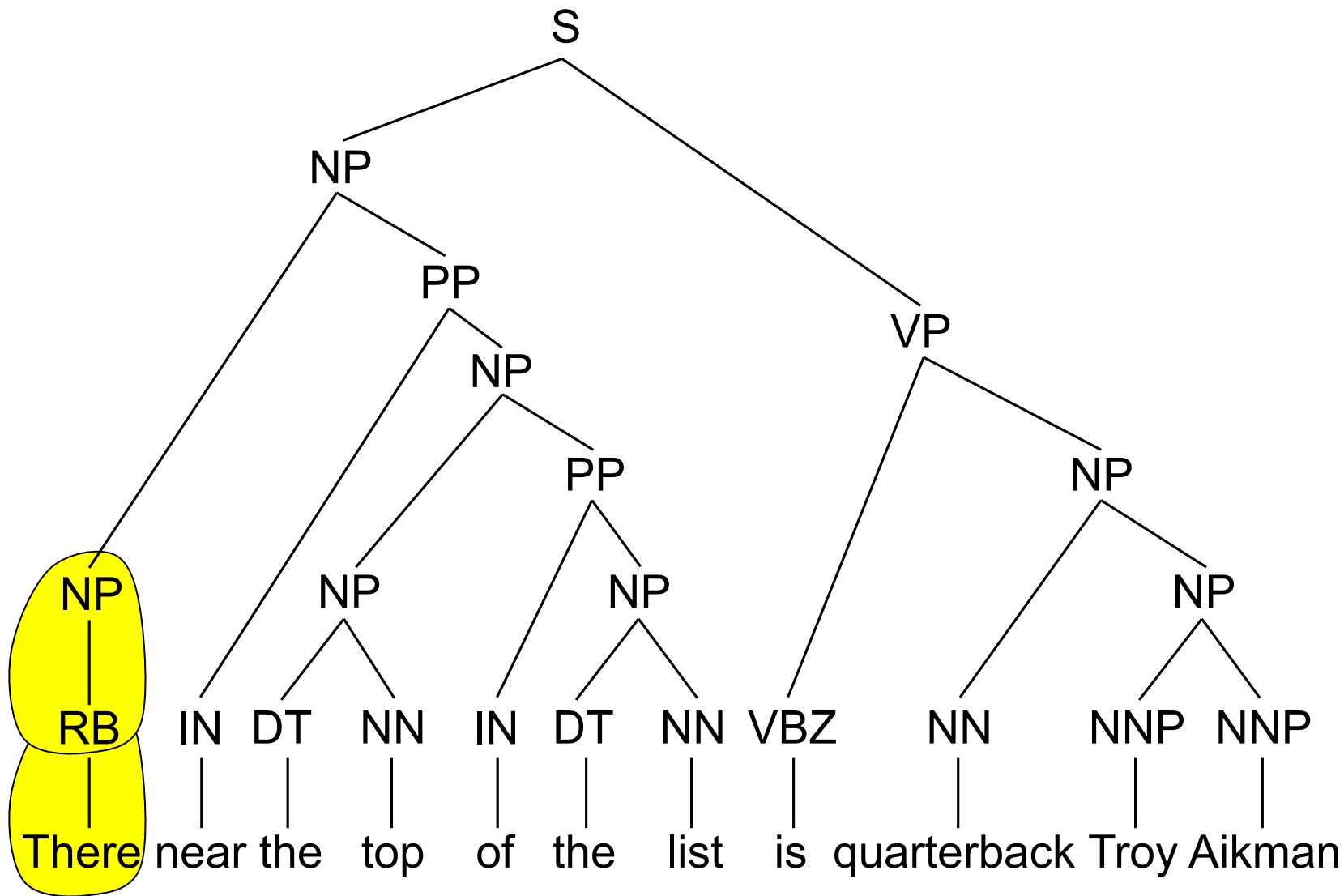
Cube Pruning

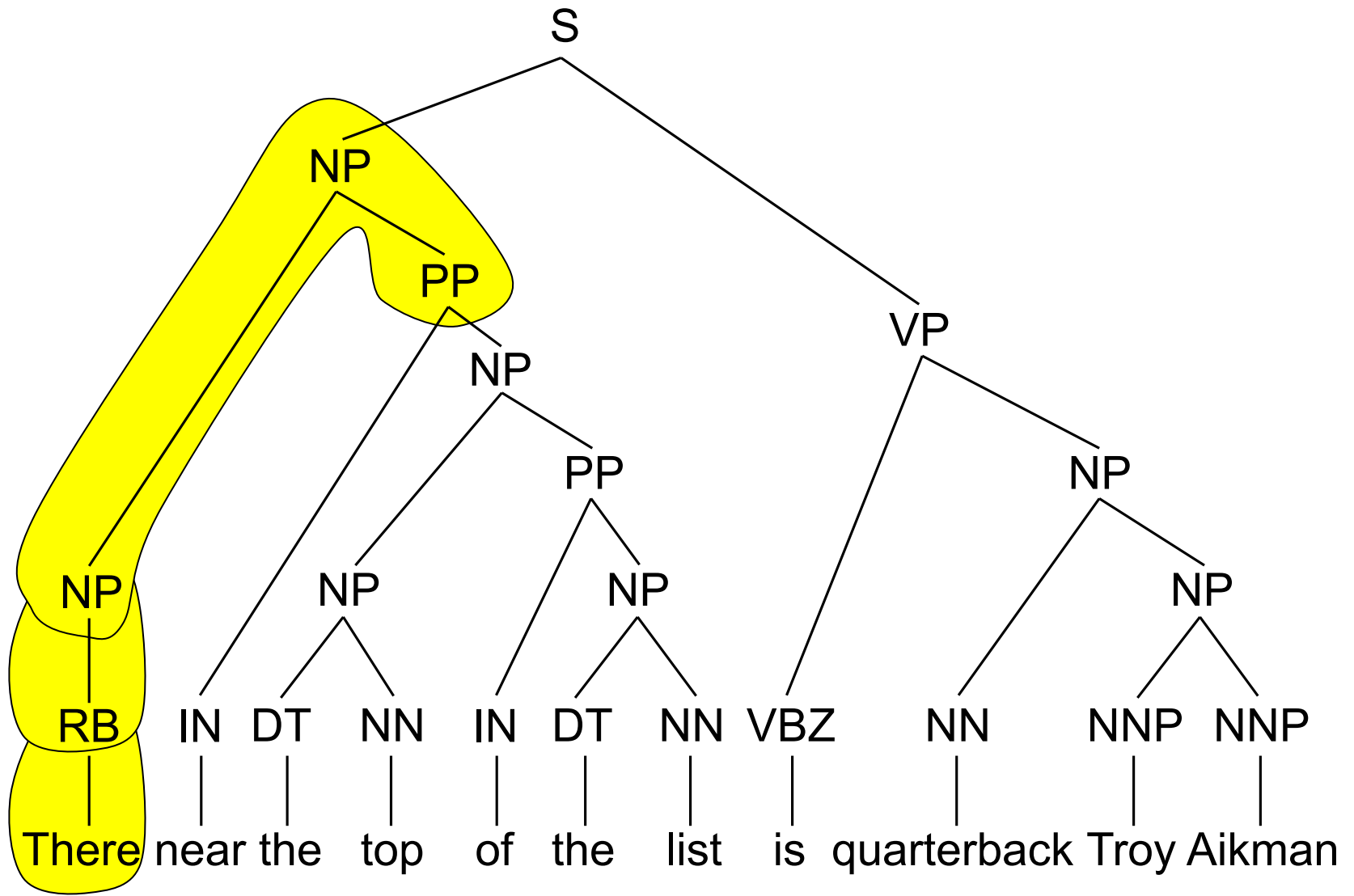
(Chiang, 2007; Huang & Chiang, 2007)

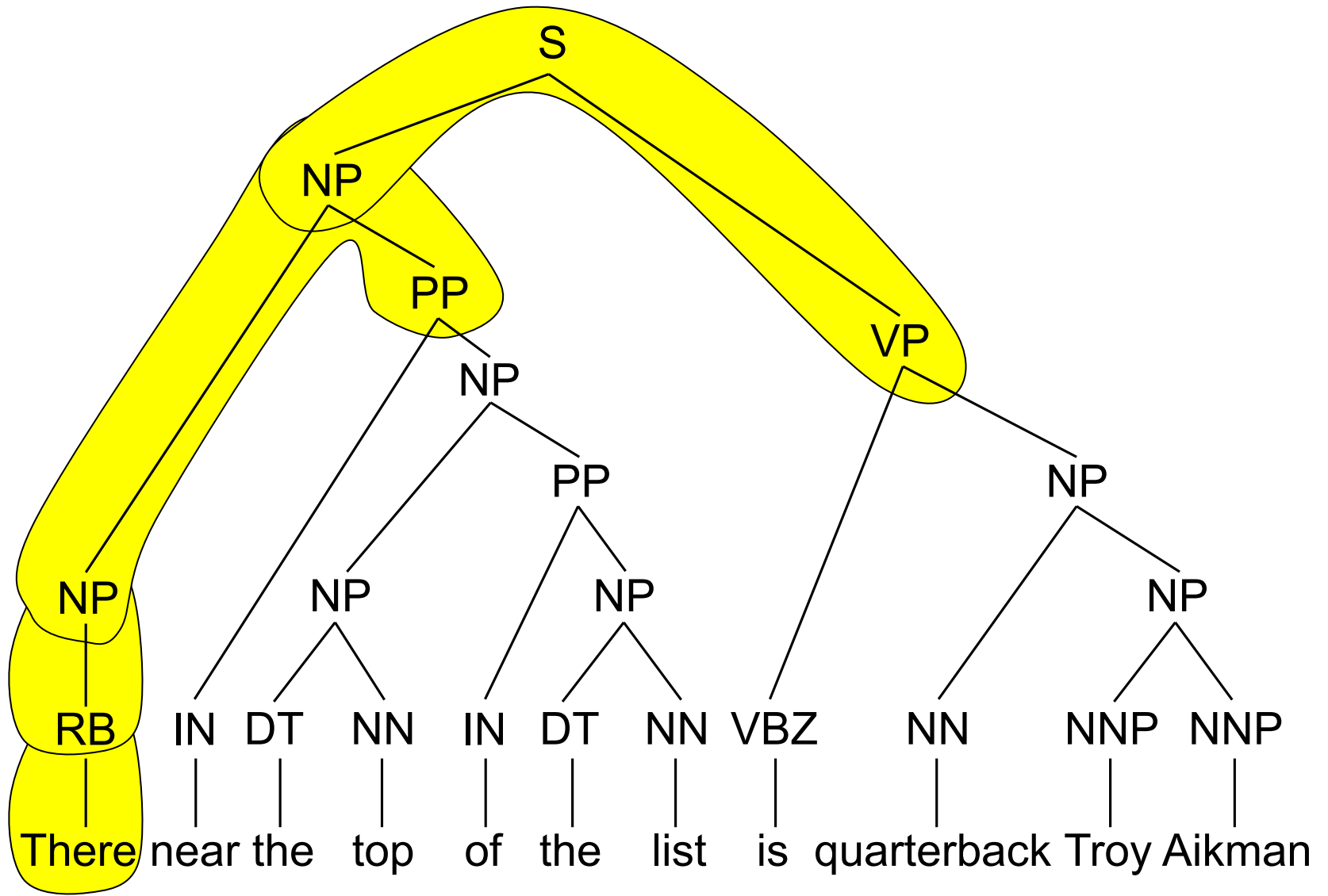
- Modification to dynamic programming algorithms for decoding to use non-local features approximately
- Keeps a k -best list of derivations for each item
- Applies non-local feature functions on these derivations when defining new items





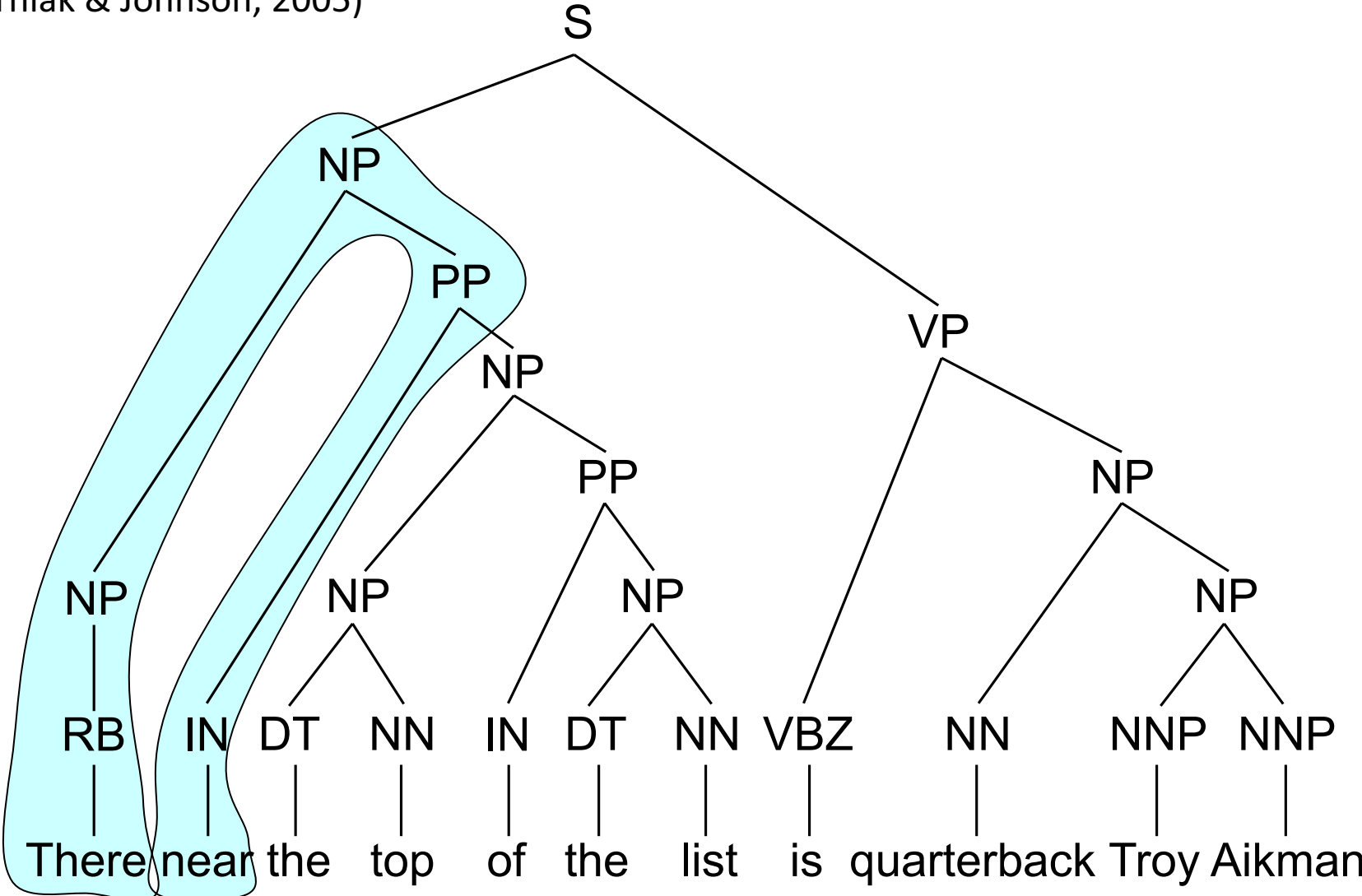






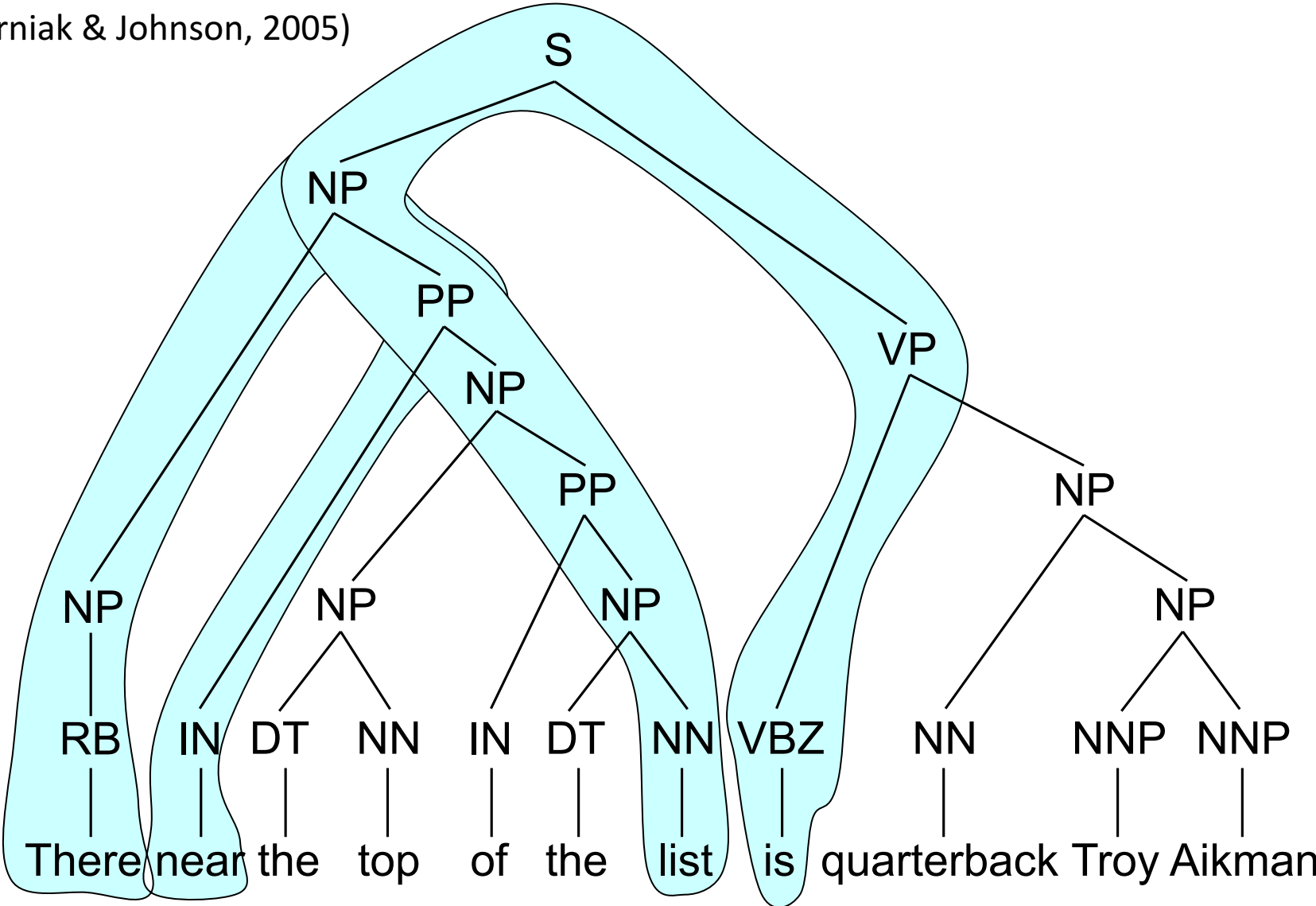
“NGramTree” feature

(Charniak & Johnson, 2005)

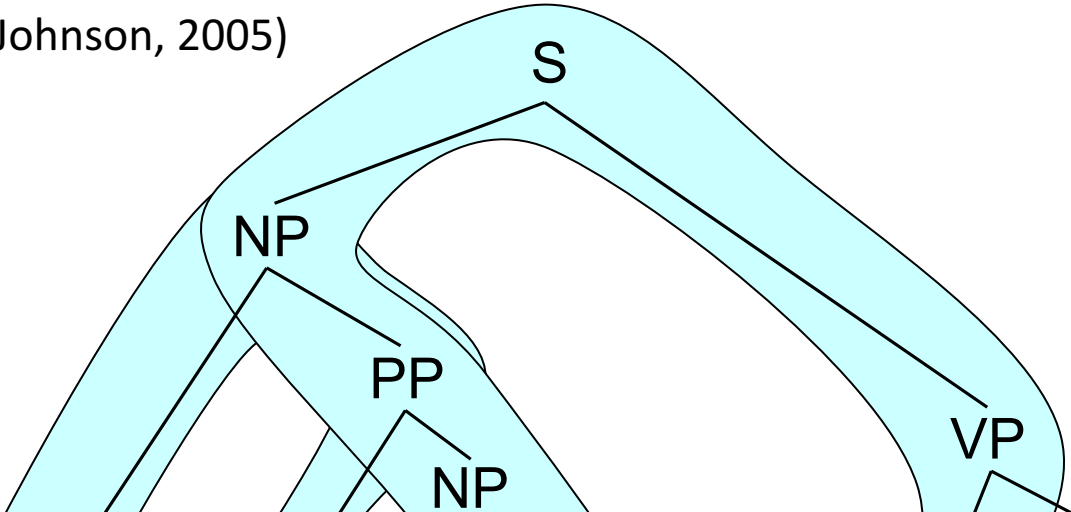


“NGramTree” feature

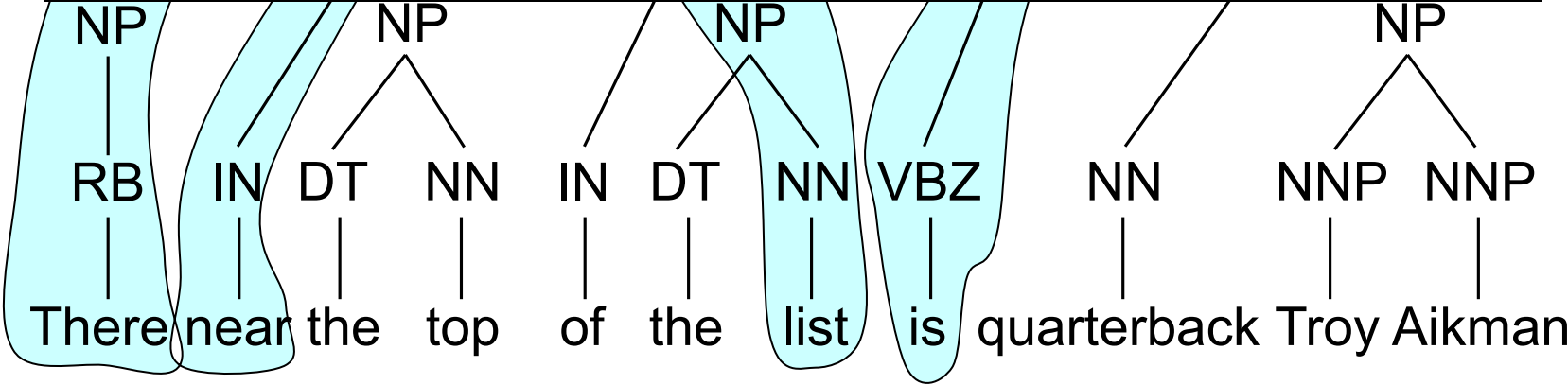
(Charniak & Johnson, 2005)



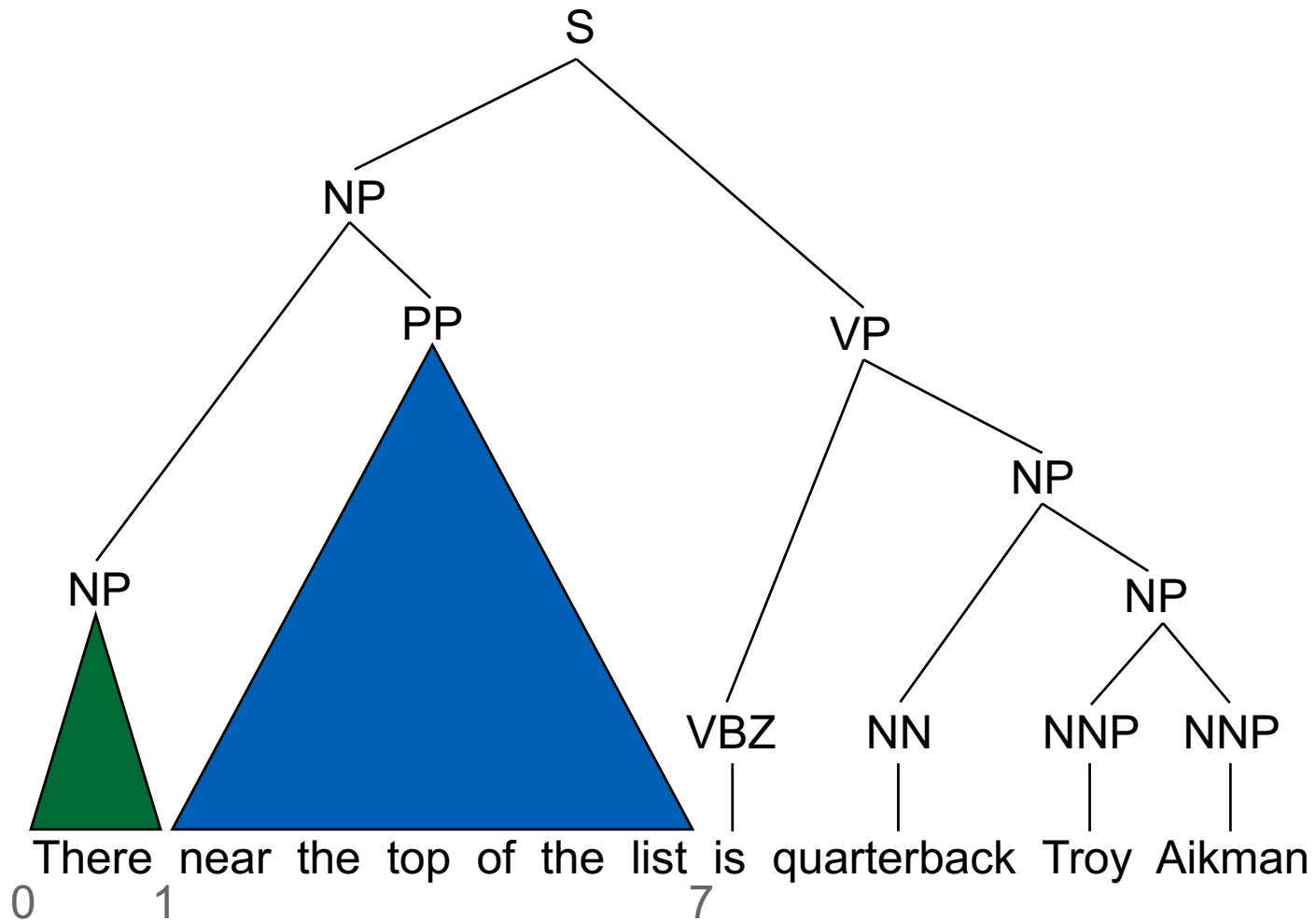
“NGramTree” feature
(Charniak & Johnson, 2005)



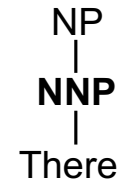
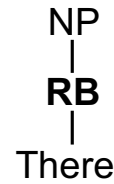
Non-local features break dynamic programming!



$$C_{NP,0,7} = C_{NP,0,1} \times C_{PP,1,7} \times \lambda_{NP \rightarrow NP PP}$$

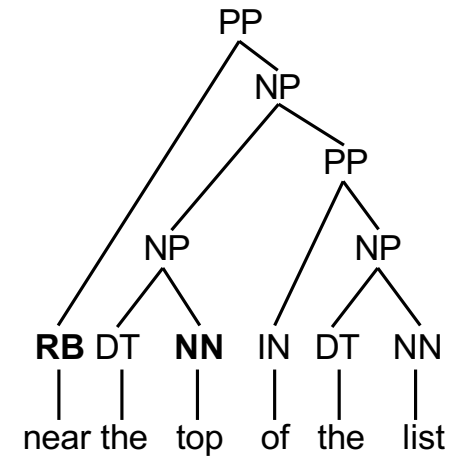
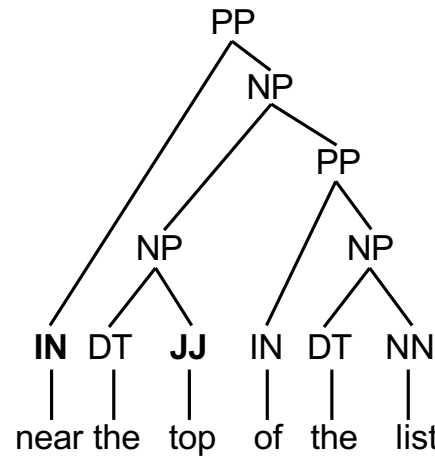
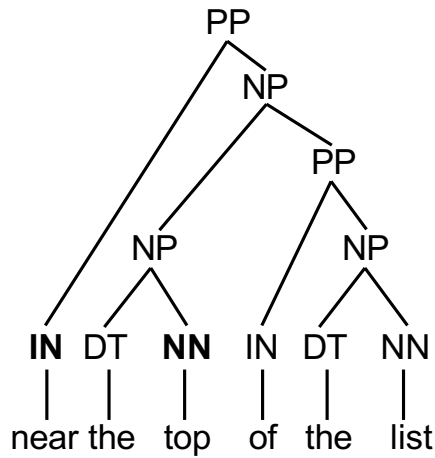


$$C_{NP,0,7} = C_{NP,0,1} \times C_{PP,1,7} \times \lambda_{NP \rightarrow NP PP}$$



$$C_{NP,0,1} =$$

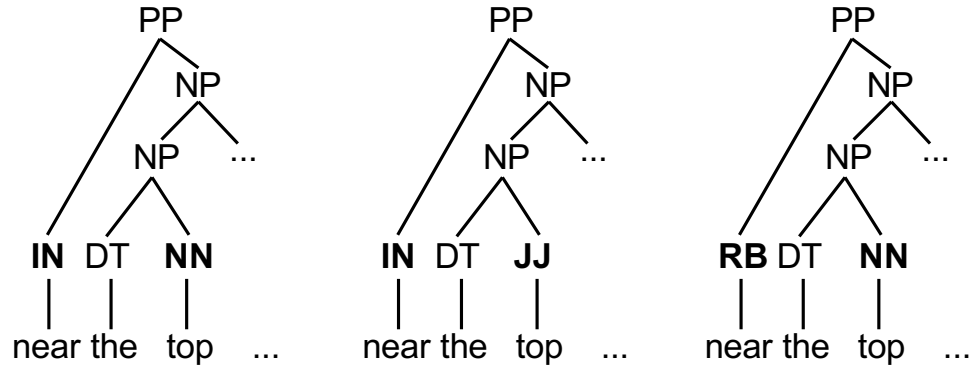
0.4	0.3	0.02
-----	-----	------



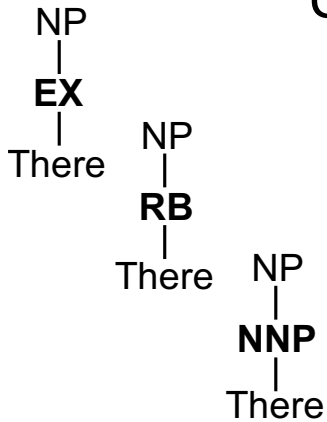
$$C_{PP,1,7} =$$

0.2	0.1	0.05
-----	-----	------

$$C_{NP,0,7} = C_{NP,0,1} \times C_{PP,1,7} \times \lambda_{NP \rightarrow NP PP}$$

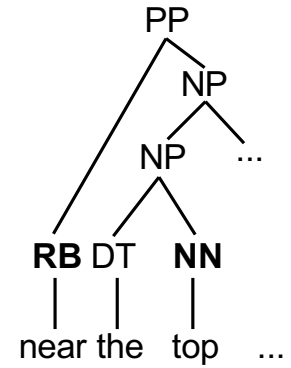
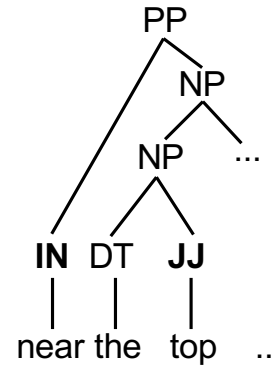
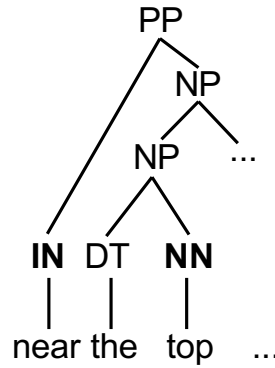


	$C_{PP,1,7}$	0.2	0.1	0.05
$C_{NP,0,1}$	0.4	0.08	0.04	0.02
	0.3	0.06	0.03	0.015
	0.02	0.004	0.002	0.001



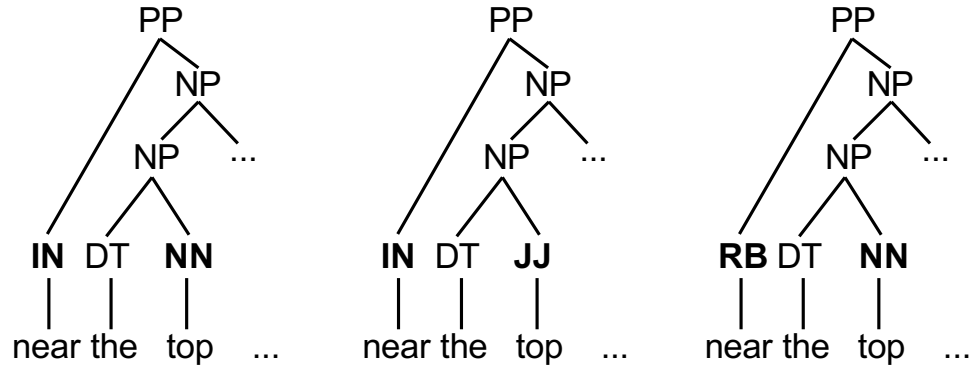
$$C_{NP,0,7} = C_{NP,0,1} \times C_{PP,1,7} \times \lambda_{NP \rightarrow NP PP}$$

$$\lambda_{NP \rightarrow NP PP} = 0.5$$

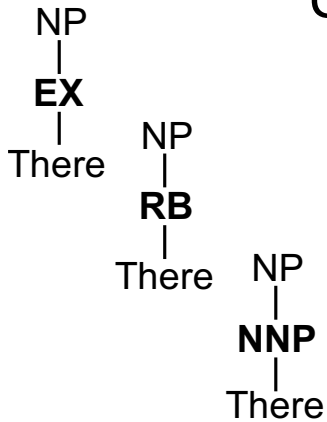


		$C_{PP,1,7}$		
$C_{NP,0,1}$		0.2	0.1	0.05
NP EX There	NP RB There	0.4	0.08 × 0.5	0.02 × 0.5
	NP NNP There	0.3	0.06 × 0.5	0.015 × 0.5
		0.02	0.004 × 0.5	0.001 × 0.5

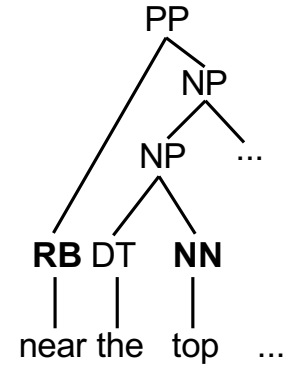
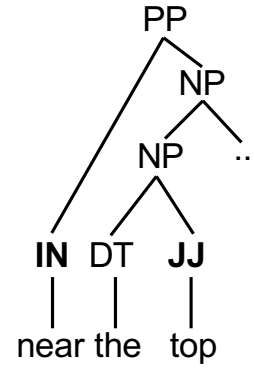
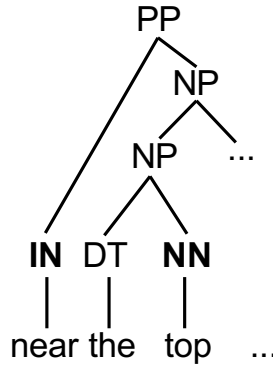
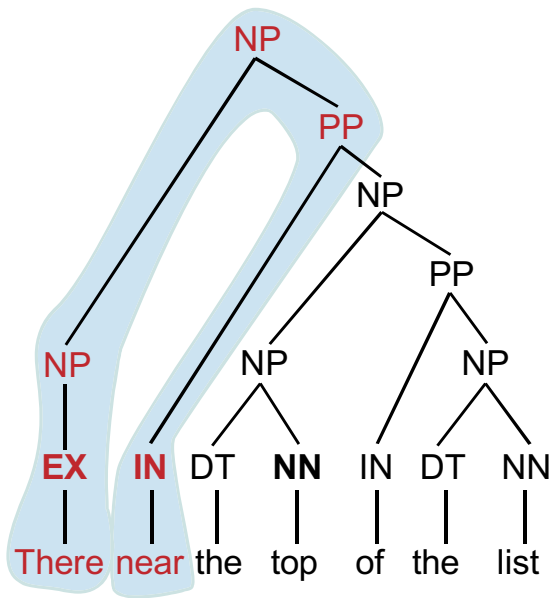
$$C_{NP,0,7} = C_{NP,0,1} \times C_{PP,1,7} \times \lambda_{NP \rightarrow NP PP}$$



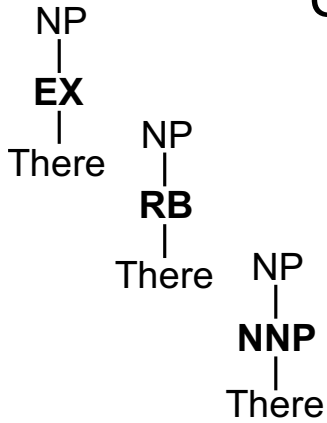
	$C_{PP,1,7}$	0.2	0.1	0.05
$C_{NP,0,1}$	0.4	0.04	0.02	0.01
	0.3	0.03	0.015	0.0075
	0.02	0.002	0.001	0.0005



$$\lambda_{\text{There EX NP NP PP IN near}} = 0.2$$



	$C_{PP,1,7}$	0.2	0.1	0.05
$C_{NP,0,1}$	0.4	0.04 × 0.2	0.02 × 0.2	0.01
	0.3	0.03	0.015	0.0075
	0.02	0.002	0.001	0.0005



$$\lambda_{\text{There EX NP NP PP IN near}} = 0.2$$

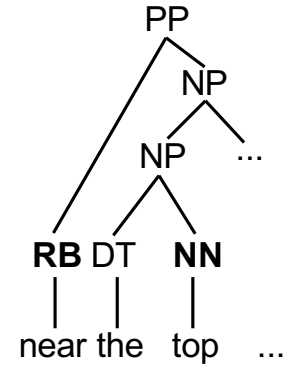
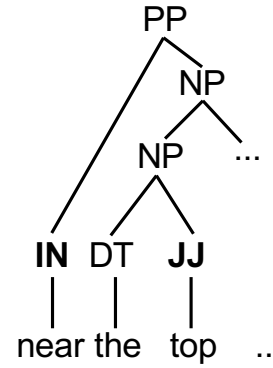
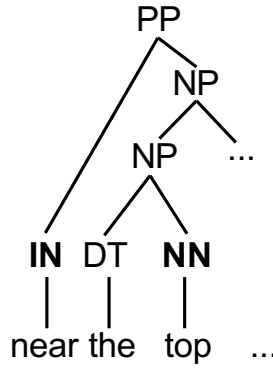
$$\lambda_{\text{There RB NP NP PP IN near}} = 0.6$$

$$\lambda_{\text{There NNP NP NP PP IN near}} = 0.1$$

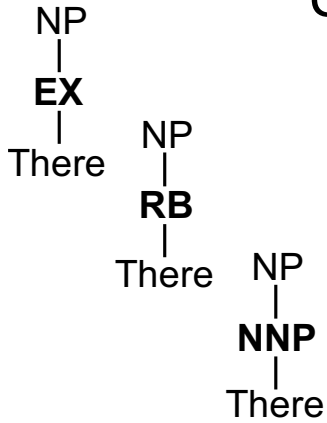
$$\lambda_{\text{There EX NP NP PP RB near}} = 0.1$$

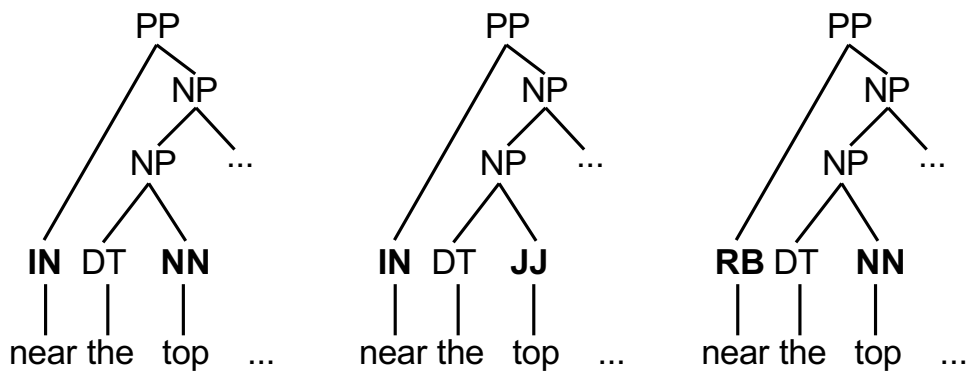
$$\lambda_{\text{There RB NP NP PP RB near}} = 0.4$$

$$\lambda_{\text{There NNP NP NP PP RB near}} = 0.2$$

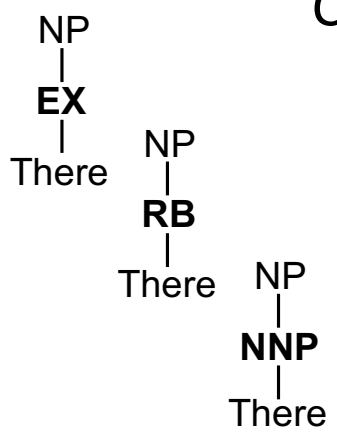


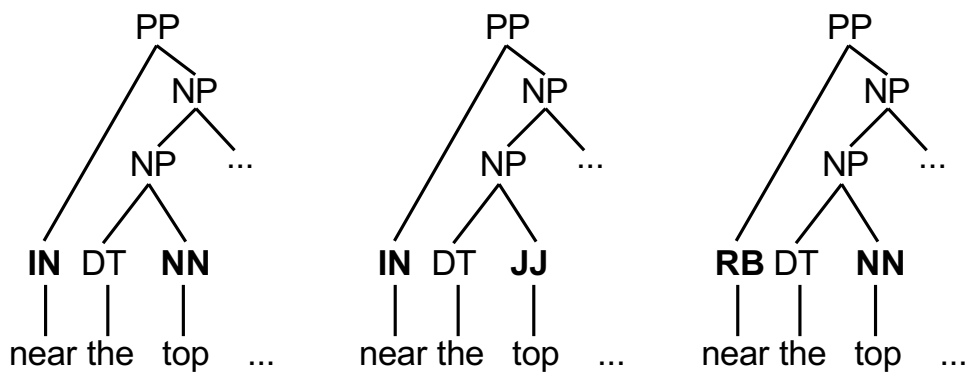
	$C_{NP,0,1}$	$C_{PP,1,7}$			
		0.2	0.1	0.05	
0.4	0.04 × 0.2	0.02 × 0.2	0.01 × 0.1		
0.3	0.03 × 0.6	0.015 × 0.6	0.0075 × 0.4		
0.02	0.002 × 0.1	0.001 × 0.1	0.0005 × 0.2		



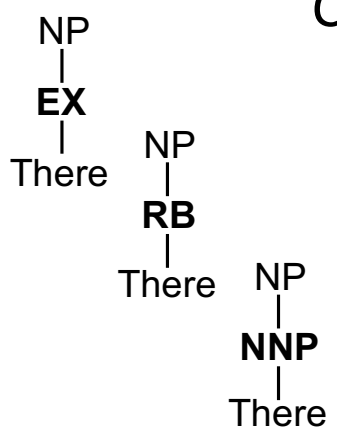


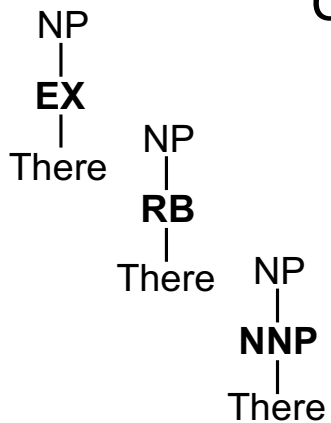
	$C_{PP,1,7}$	0.2	0.1	0.05
$C_{NP,0,1}$	0.4	0.008	0.004	0.001
	0.3	0.018	0.009	0.003
	0.02	0.0002	0.0001	0.0001



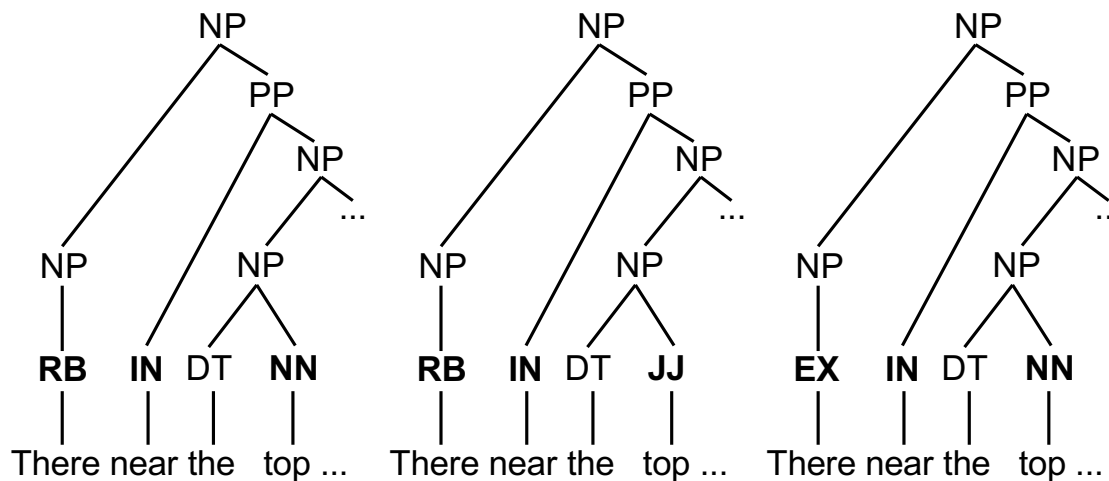


	$C_{PP,1,7}$	0.2	0.1	0.05
$C_{NP,0,1}$	0.4	0.008	0.004	0.001
	0.3	0.018	0.009	0.003
	0.02	0.0002	0.0001	0.0001





	$C_{PP,1,7}$	0.2	0.1	0.05
$C_{NP,0,1}$	0.4	0.008	0.004	0.001
	0.3	0.018	0.009	0.003
	0.02	0.0002	0.0001	0.0001



$C_{NP,0,7}$	0.018	0.009	0.008
--------------	-------	-------	-------

Clarification

- Cube pruning does not actually expand all k^2 proofs as this example showed
- It uses a fast approximation that only looks at $O(k)$ proofs

Integer Linear Programming

- (on board)