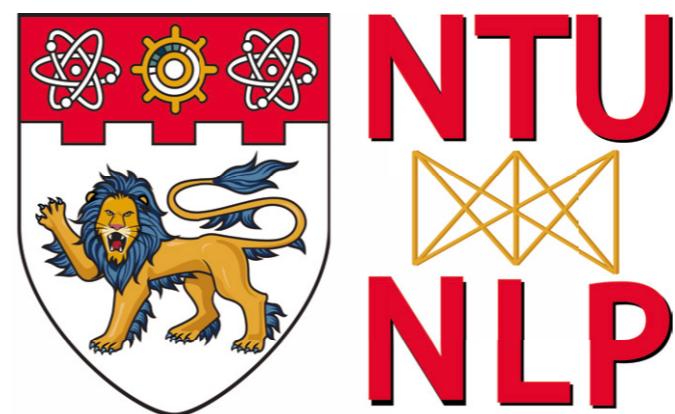


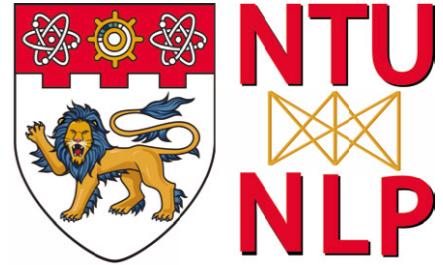
Deep Learning for Natural Language Processing

Shafiq Joty



Lecture 7: Recursive Neural Nets & Parsing

Where we are



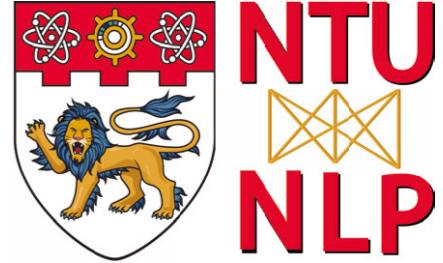
Models/Algorithms

- Linear models
- Feed-forward Neural Nets (FNN)
- Window-based methods
- Convolutional Nets
- Recurrent Neural Nets
- Recursive Neural Nets

NLP tasks/applications

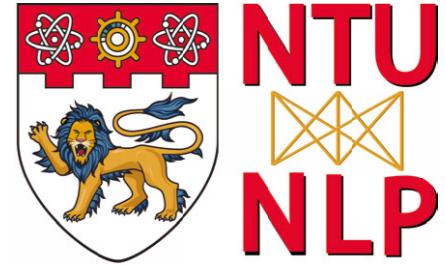
- Word meaning
- Language modelling
- Sequence tagging
- Sequence encoding
- Parsing
- Hierarchical encoding

Today's Plan

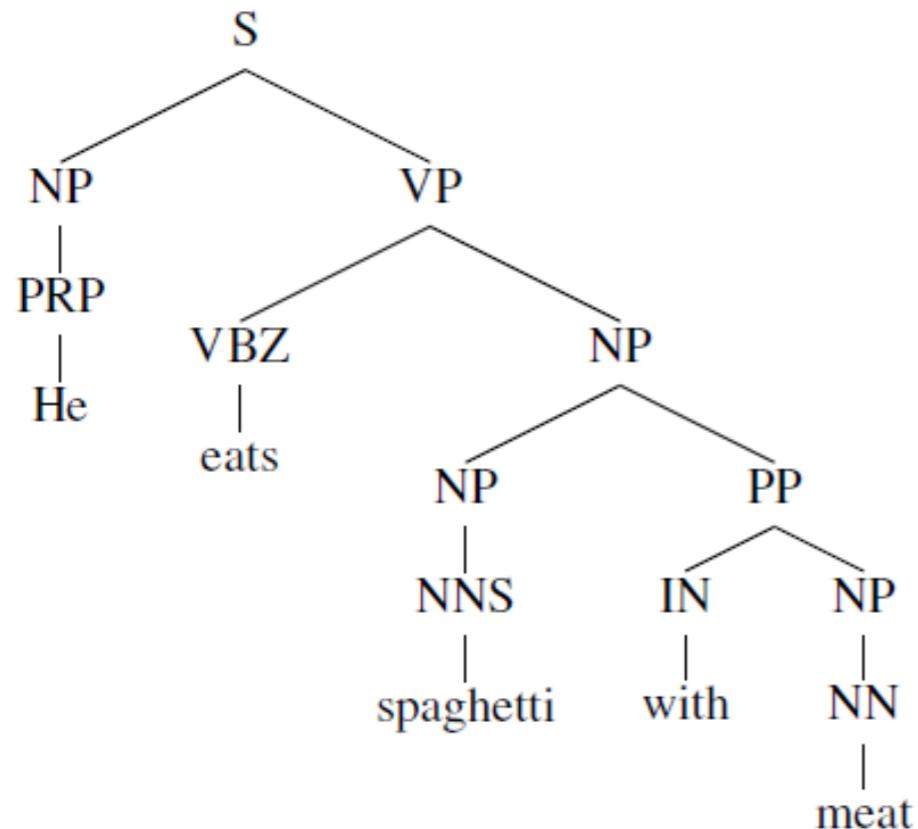


- Syntactic structures
 - Constituency & Dependency structures
- Need for syntactic structures
- Are languages hierarchical?
- Parsing with Recursive Neural Networks
- Backpropagation through tree
- Semantic composition with Tree RNNs
- Recent Parsers

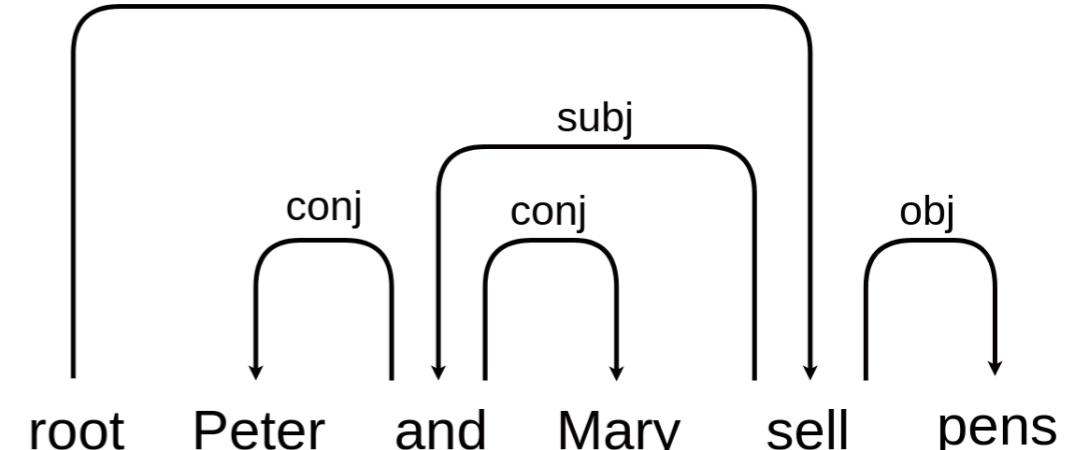
Two Views of Syntactic Structures



Two types of syntactic structures of a sentence

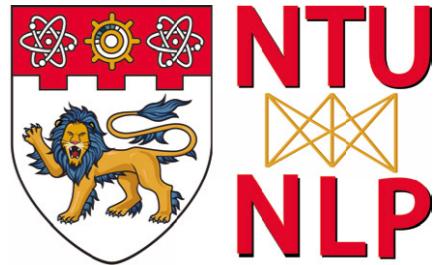


Constituency/phrase structure

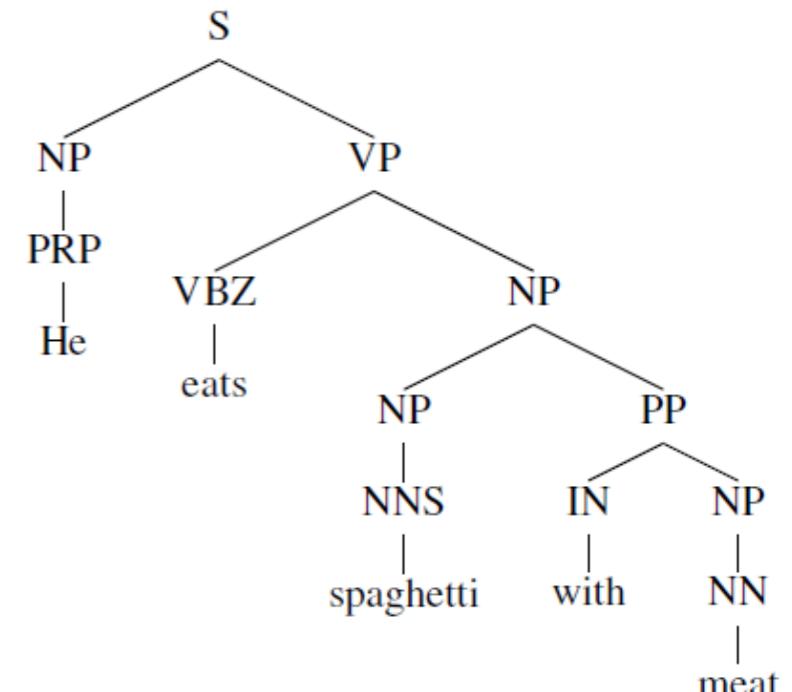


Dependency structure

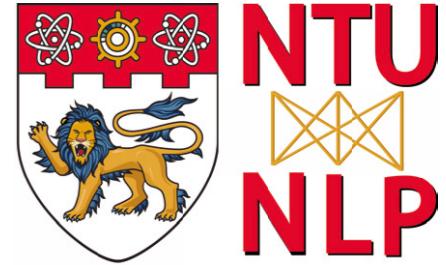
Constituency Structure



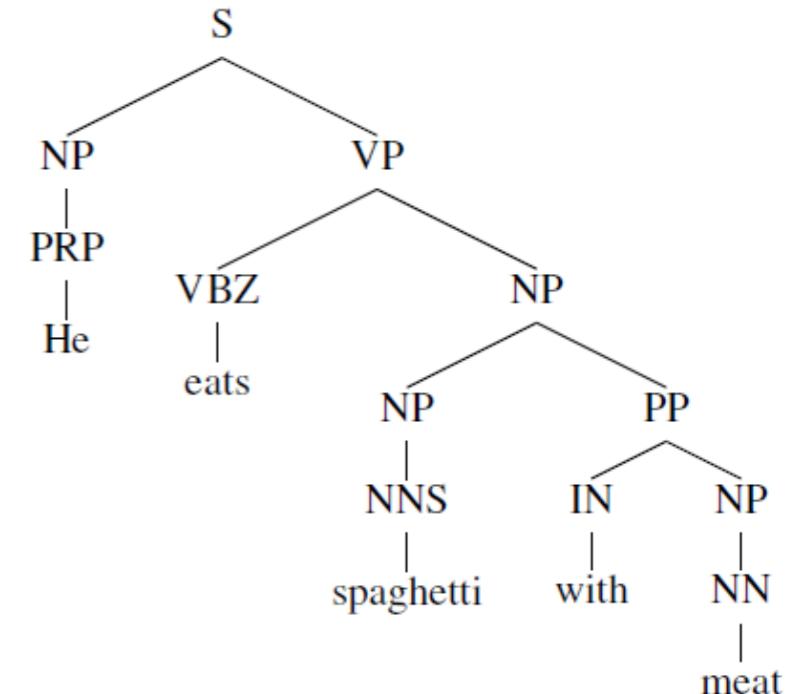
- Phrase structure organizes words into nested constituents
- Atomic units are words. E.g., He, eats, spaghetti, with, meat
- Words combine into phrases: with meat,
- Phrases combine into larger phrases (recursively): spaghetti with meat



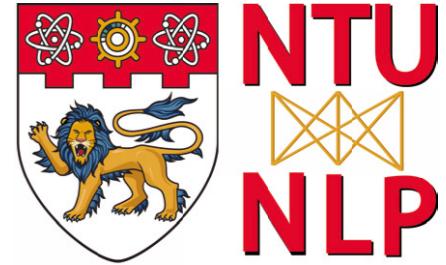
Constituency Structure



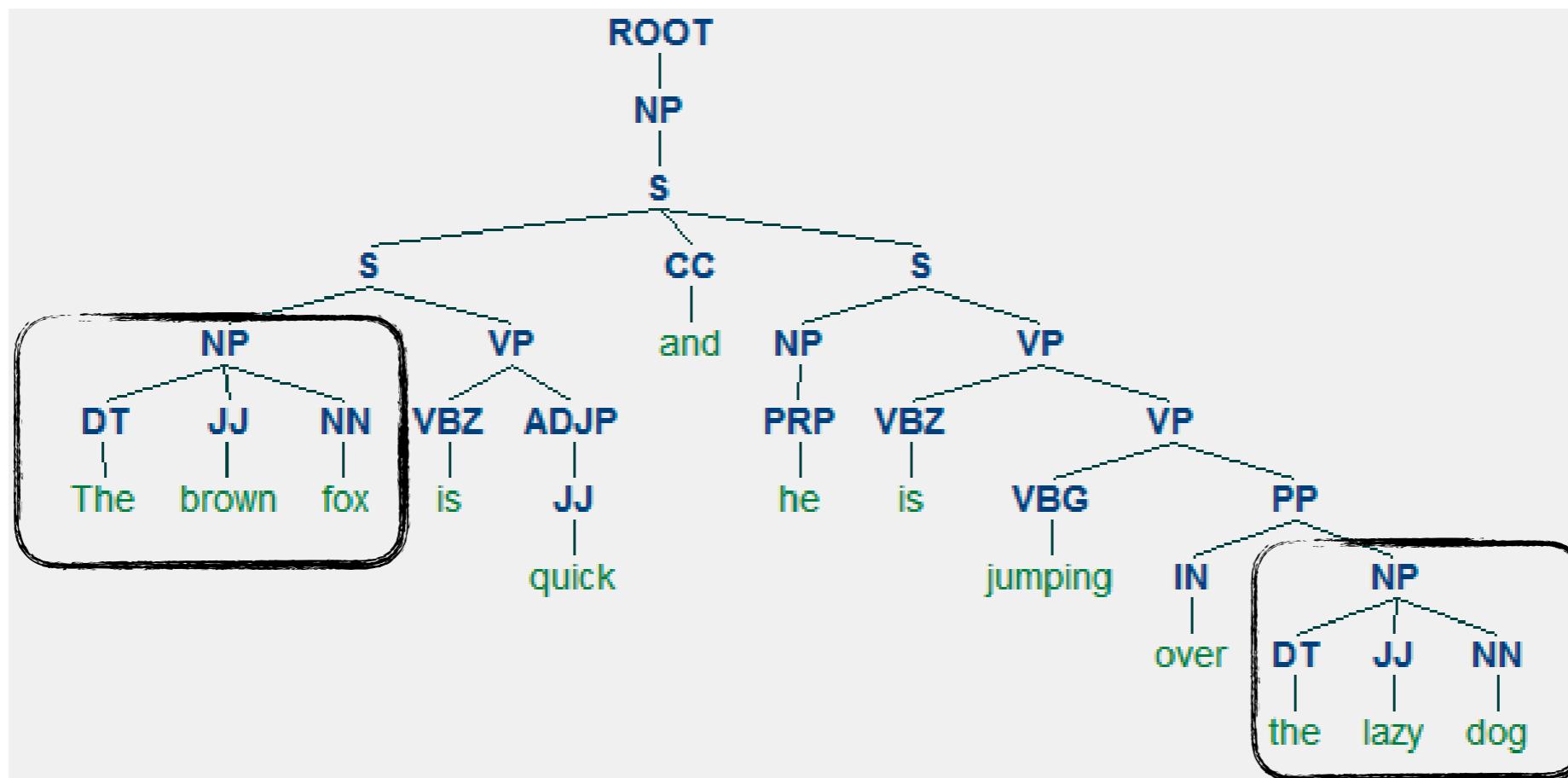
- Phrase structure organizes words into nested constituents
 - Can represent the process (grammar) with CFG rules
- Atomic units are words: words are assigned a category (part-of-speech). E.g., He-PRP, eats-VBZ, with-IN, meat-NN
- Words combine into phrases: with phrasal categories. E.g., PP → IN NP (with meat)
- Phrases combine into bigger phrases (recursively) with phrasal categories. E.g., NP → NP PP (spaghetti with meat)



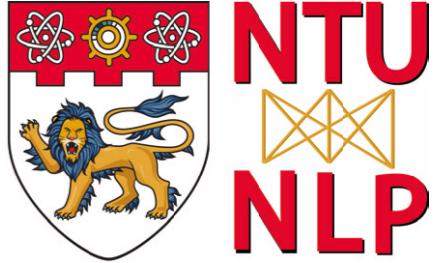
Constituency Structure



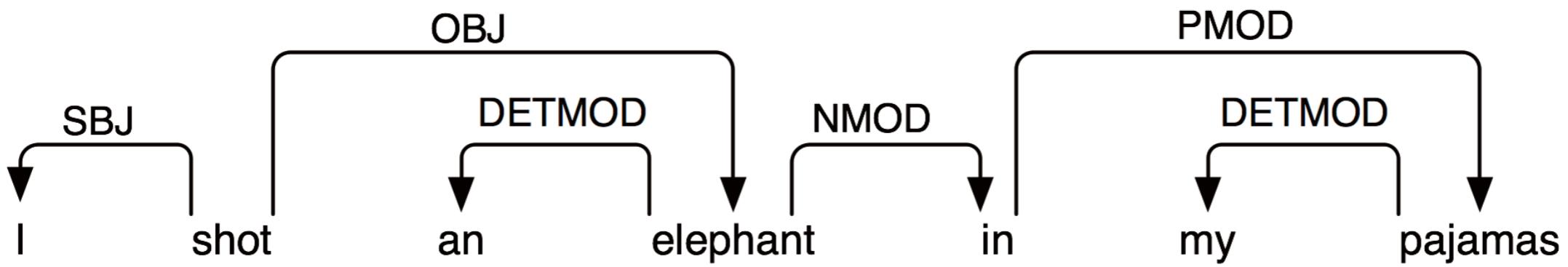
- Phrase structure organizes words into nested constituents
 - Can represent the process (grammar) with CFG rules
 - We can describe a language (e.g., English) using the CFG rules
 - Dominant structure in linguistics



Dependency Structure

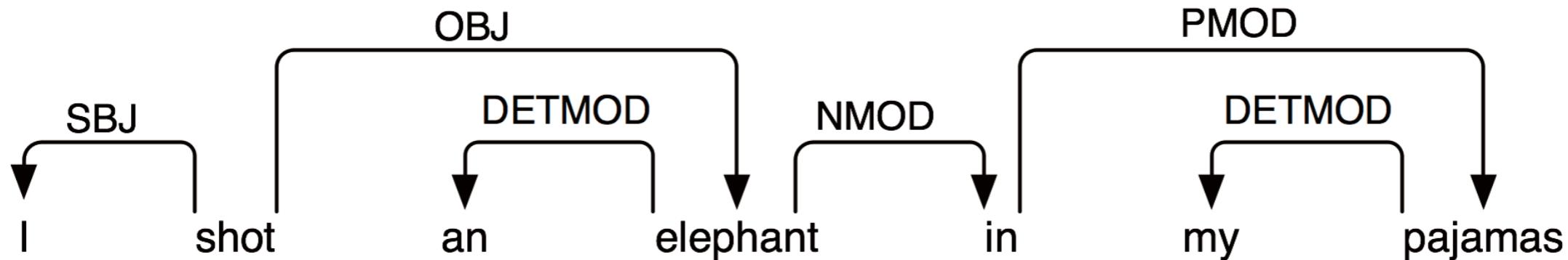


- Dependency structure shows which words depend on (modify or are arguments of) which other words
- Dependency is a binary asymmetric relation that holds between a **head** and its **dependents**.



Directed Acyclic Graph (Tree)

Dependency Structure



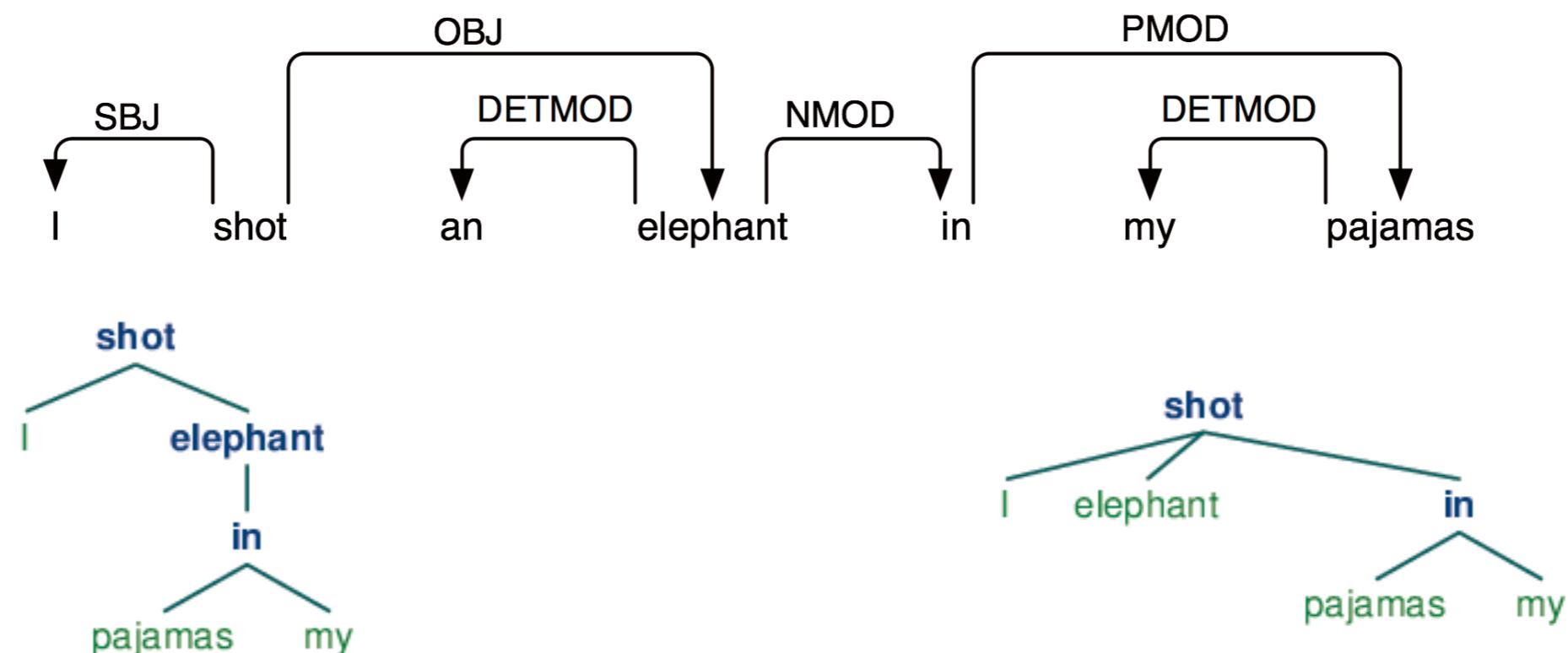
- Arrows point from heads to their dependents.
- Labels indicate the grammatical function of the dependent **as subject, object or modifier**
- Dependency grammar is used to **directly express grammatical functions as a type of dependency**.

```

... 'shot' -> 'I' | 'elephant' | 'in'
... 'elephant' -> 'an' | 'in'
... 'in' -> 'pajamas'
... 'pajamas' -> 'my'
  
```

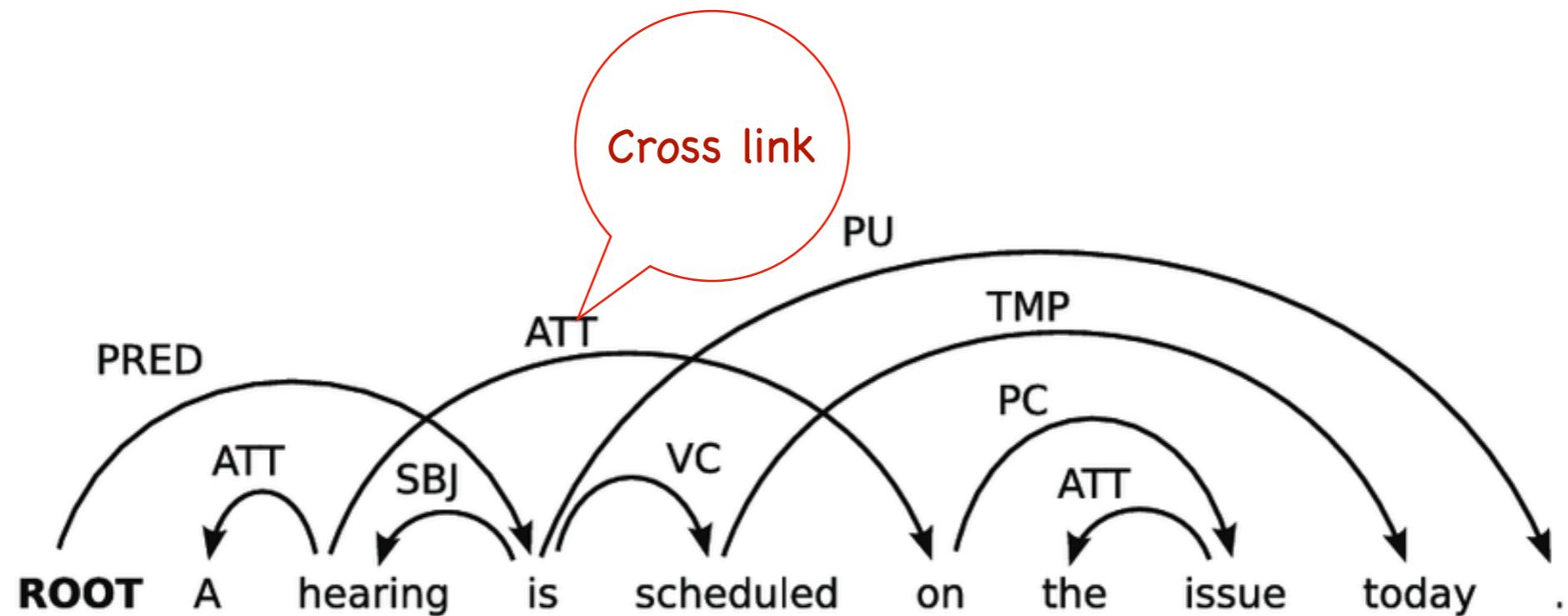
Dependency Structure

- A dependency graph is **projective** if, when the words are written in linear order, the edges can be drawn above the words without crossing.
- This is equivalent to saying that a word and all its descendants (dependents and dependents of its dependents, etc.) form a contiguous sequence of words within the sentence.
- We can parse most English sentences using a projective parser



Dependency Structure

- Non-projective dependency structure



Why Syntactic Structures?

- To be able to interpret language correctly.
- Humans communicate complex ideas by composing words together into bigger units to convey complex meanings.
- We interpret the meaning of larger text units – entities, descriptive terms, facts, arguments, stories – by **semantic composition** of smaller elements
- We need to know what is connected to what to be able to understand bigger things from knowing about smaller parts

A person on a snowboard jumps into the air

Why Syntactic Structures?

- PP attachment ambiguity

San Jose cops kill man with knife
Text Paper Close
Translate Listen

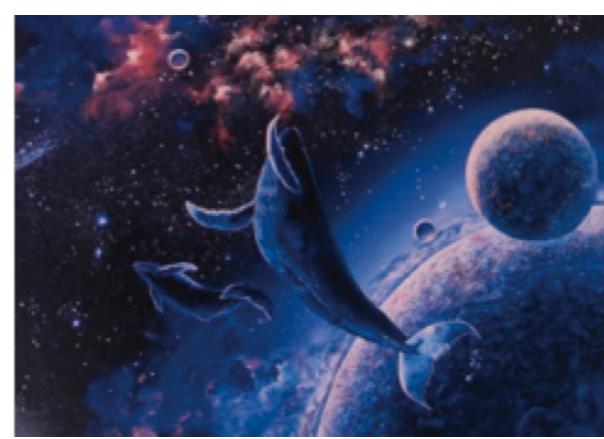
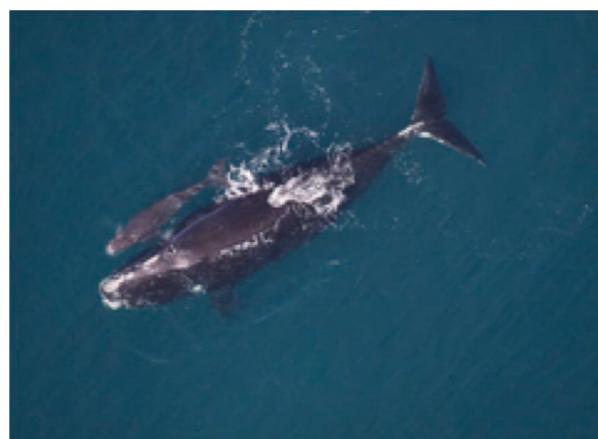
San Jose cops kill man with knife

Why Syntactic Structures?

- PP attachment ambiguity



Scientists count whales from space



Why Syntactic Structures?

- PP attachment ambiguity
- A key parsing decision is how we ‘attach’ various constituents
 - PPs, adverbial or participial phrases, infinitives, coordinations

The board approved [its acquisition] [by Royal Trustco Ltd.]

[of Toronto]

[for \$27 a share]

[at its monthly meeting].

- # of possible structures: Catalan numbers: $C_n = (2n)!/[(n+1)!n!]$
- An exponentially growing series, which arises in many tree-like contexts:
- E.g., the number of possible triangulations of a polygon with $n+2$ sides

Why Syntactic Structures?

- Coordination scope ambiguity

Small rats and mice can squeeze into holes or cracks in the wall.

Shuttle veteran and longtime NASA executive Fred Gregory appointed to board

Why Syntactic Structures?

- Adjectival modifier ambiguity

Students get their **first** lab experience

Why Syntactic Structures?

● Verb Phrase attachment



theguardian

home > world > americas asia ≡ all

Rio de Janeiro

Mutilated body washes up
on Rio beach to be used for
Olympics beach volleyball

Are Languages Hierarchical?

The Faculty of Language: What Is It, Who Has It, and How Did It Evolve?

Marc D. Hauser^{1,*}, Noam Chomsky², W. Tecumseh Fitch¹

⁺ See all authors and affiliations

Science 22 Nov 2002:
Vol. 298, Issue 5598, pp. 1569-1579
DOI: 10.1126/science.298.5598.1569

Article

Figures & Data

Info & Metrics

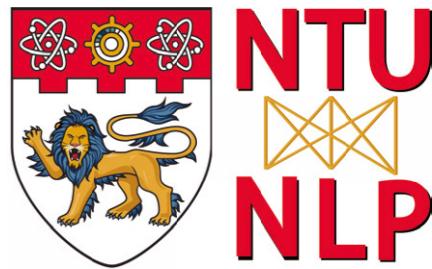
eLetters

 **PDF**

Abstract

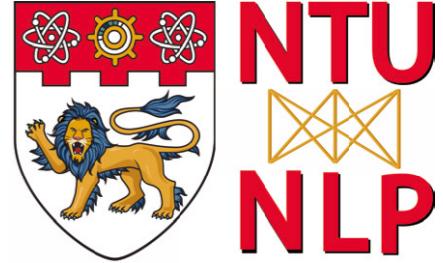
We argue that an understanding of the faculty of language requires substantial interdisciplinary cooperation. We suggest how current developments in linguistics can be profitably wedded to work in evolutionary biology, anthropology, psychology, and neuroscience. We submit that a distinction should be made between the faculty of language in the broad sense (FLB) and in the narrow sense (FLN). FLB includes a sensory-motor system, a conceptual-intentional system, and the computational mechanisms for recursion, providing the capacity to generate an infinite range of expressions from a finite set of elements. We hypothesize that FLN only includes recursion and is the only uniquely human component of the faculty of language. We further argue that FLN may have evolved for reasons other than language, hence comparative studies might look for evidence of such computations outside of the domain of communication (for example, number, navigation, and social relations).

Are Languages Hierarchical?



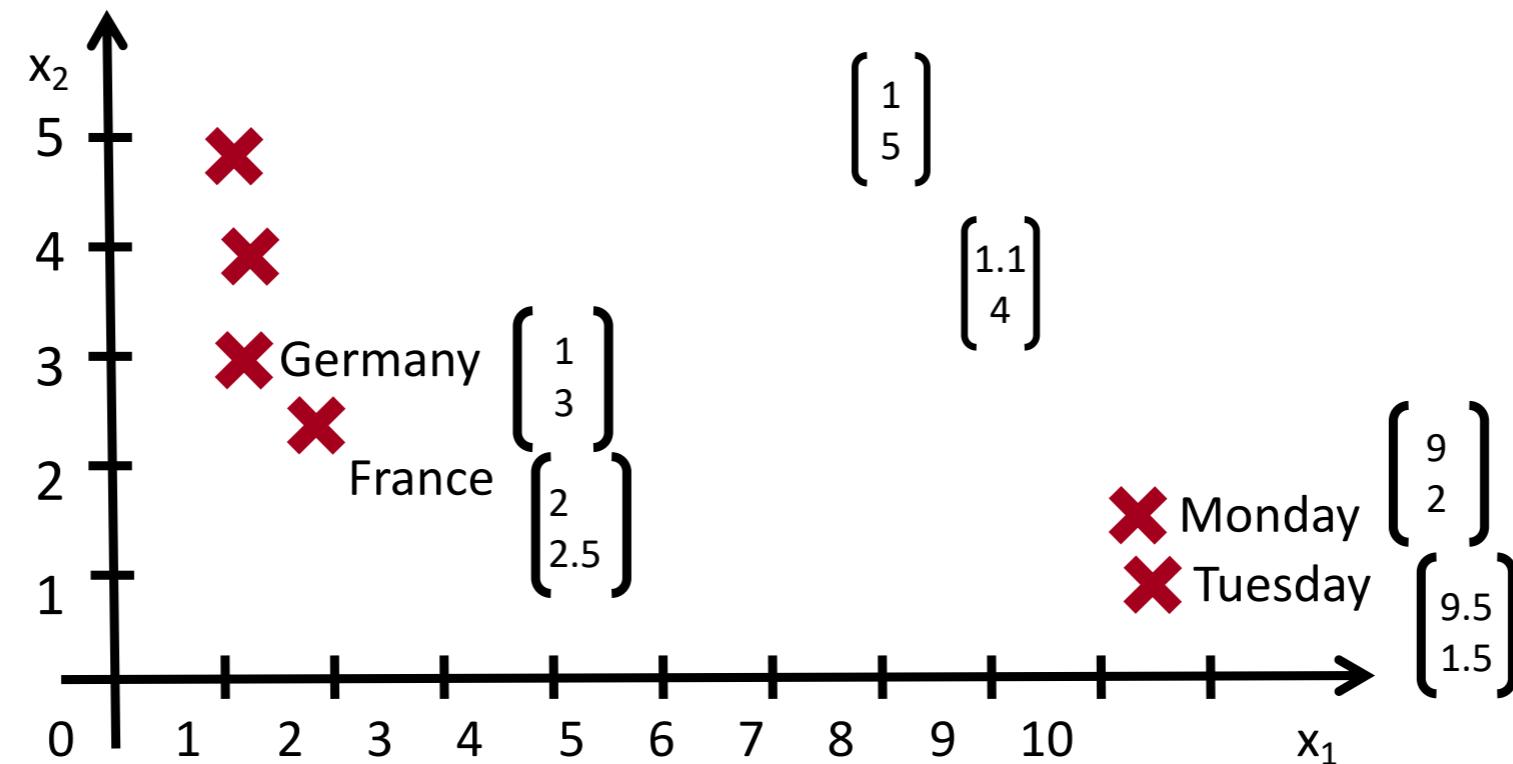
- Cognitively somewhat debatable (need to head to infinity)
- But: recursion is natural for describing language
 - [The person standing next to [the man from [the company that purchased [the firm that you used to work at]]]]
 - NP containing a NP containing a NP
- It's a very powerful prior for language structure

Today's Plan



- Syntactic structures
 - Constituency & Dependency structures
- Need for syntactic structures
- Are languages hierarchical?
- Parsing with Recursive Neural Networks
- Backpropagation through tree
- Semantic composition with Tree RNNs
- Recent Parsers

Building on Word Vectors

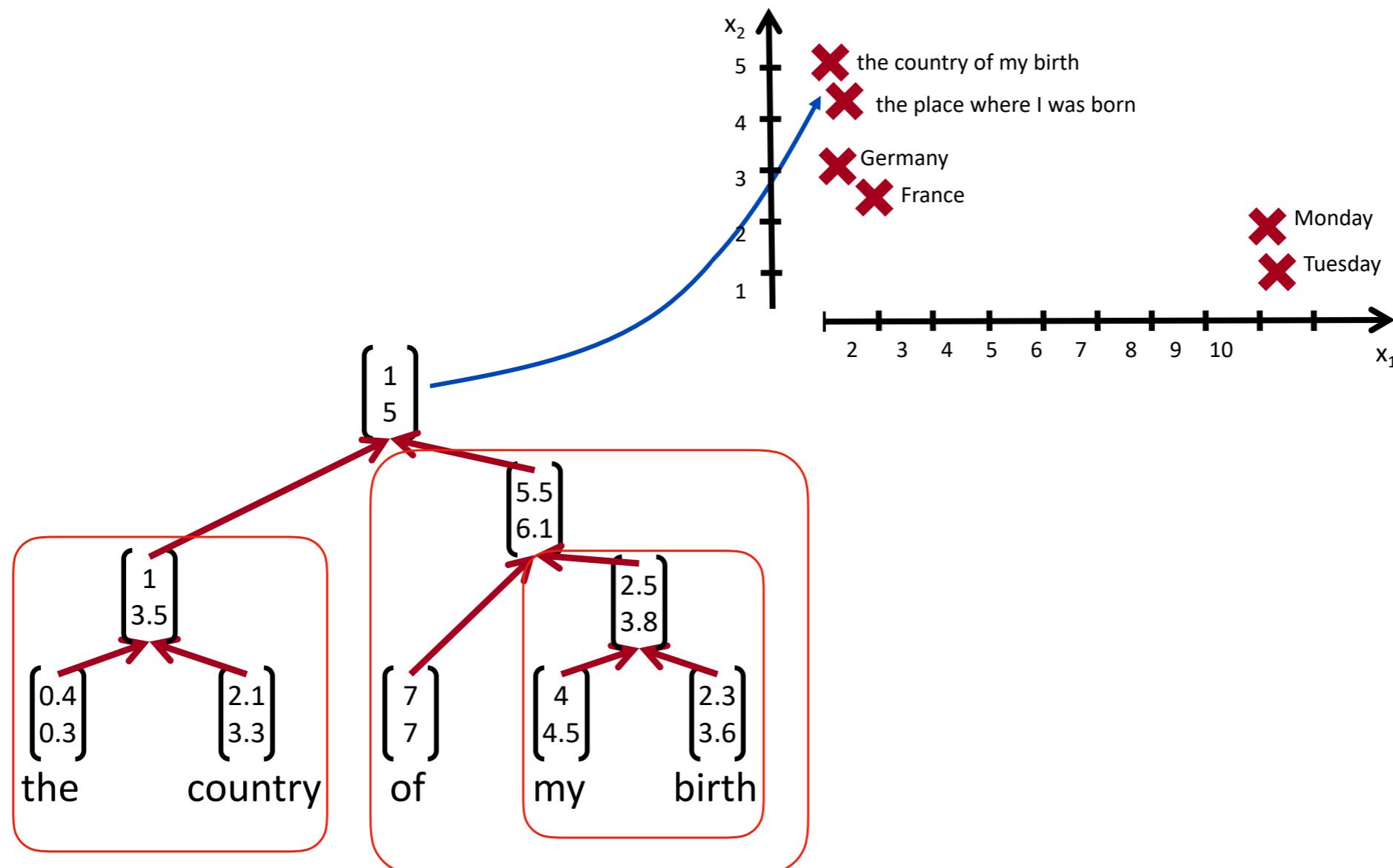


- How can we represent the meaning of longer phrases
 - map them into the same vector space!

My birth place
the place where I was born
The place of my birth

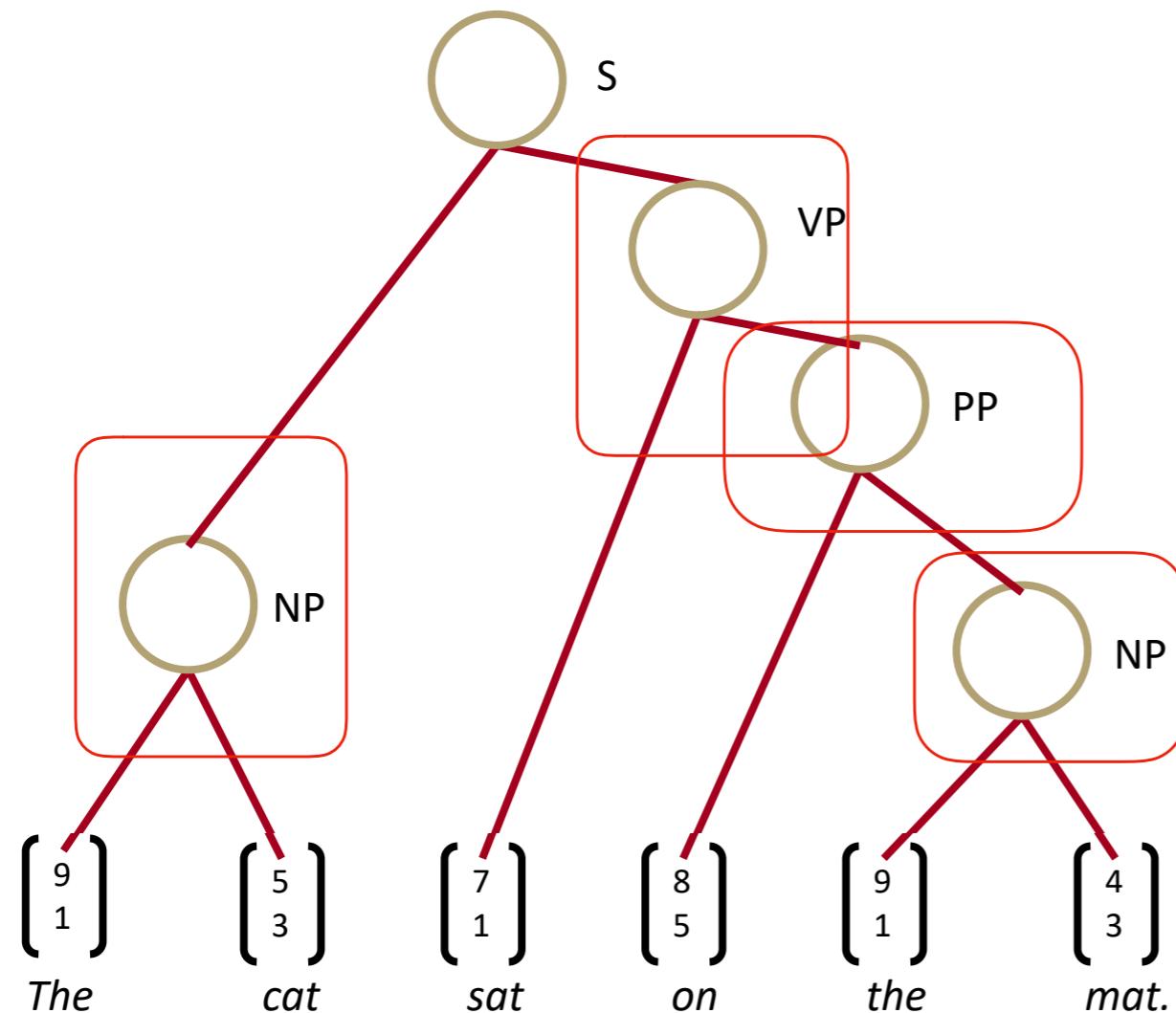
How Should We Map Phrases

- Map them into the same vector space!



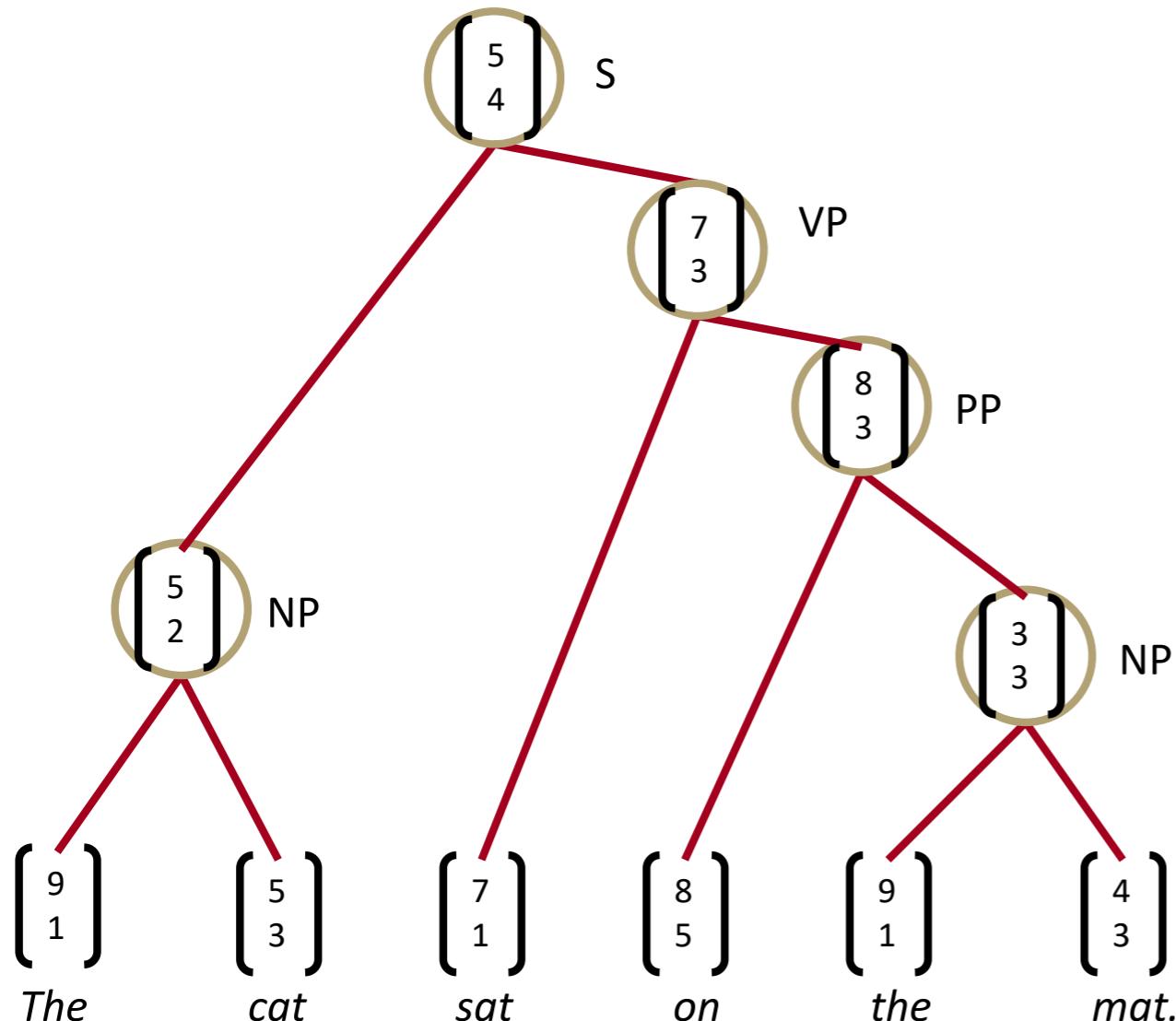
The meaning (vector) of a sentence is determined by
(1) the meanings of its words and (2) the rules that combine them.

Constituency Parsing



Jointly learn **parse trees**
and compositional vector
representations

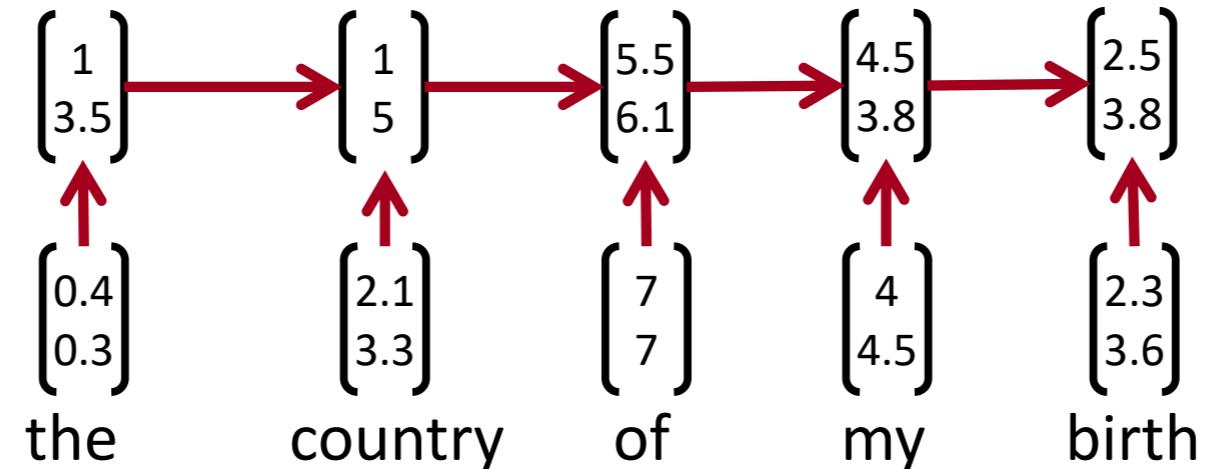
Constituency Parsing



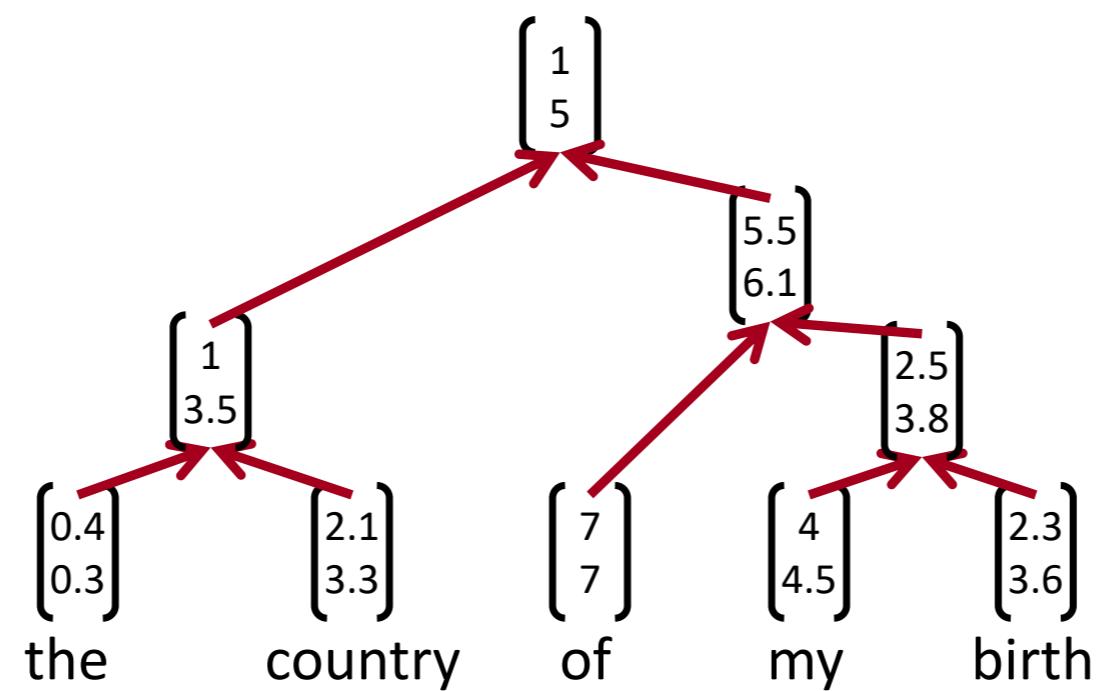
Jointly learn parse trees
and **compositional vector**
representations

Recursive vs. Recurrent

Recurrent

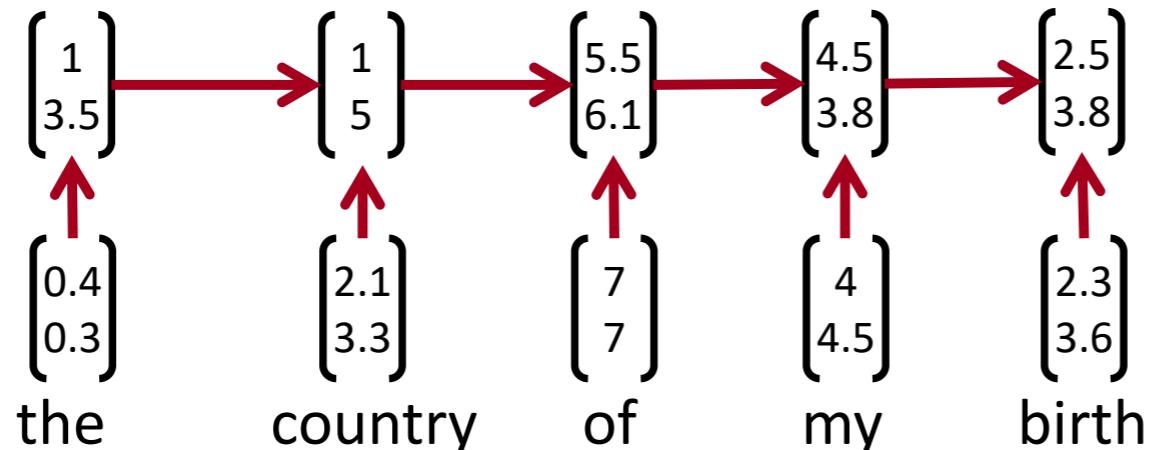


Recursive



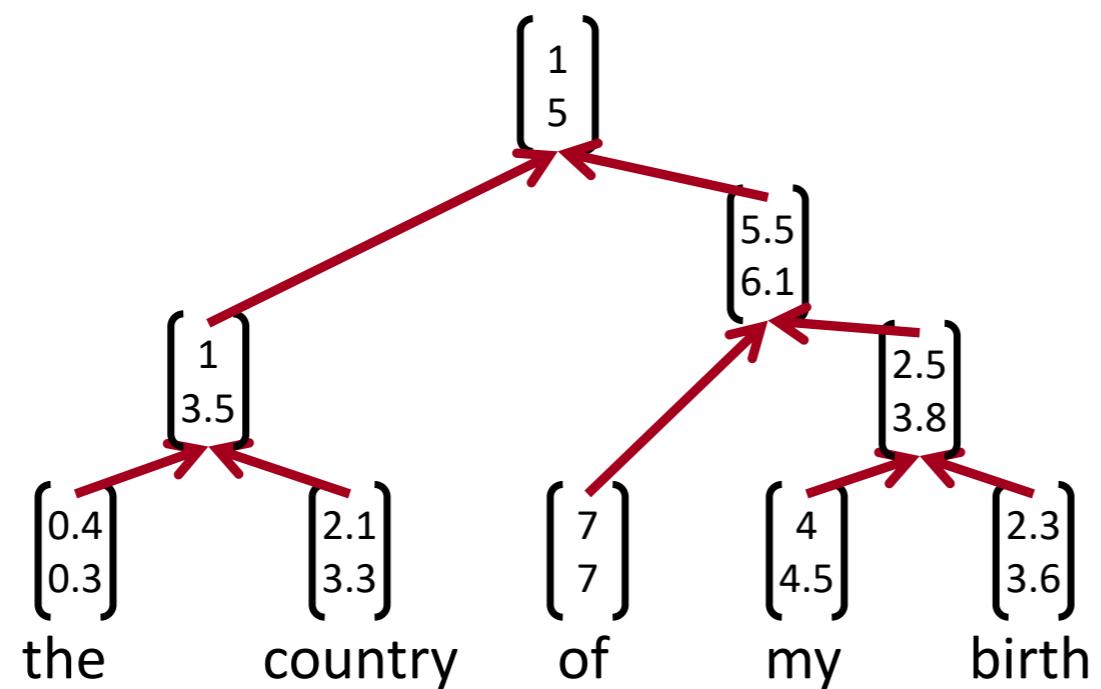
Recursive vs. Recurrent

Recurrent



Cannot capture phrases without prefix context and often capture too much of last words in final vector

Recursive



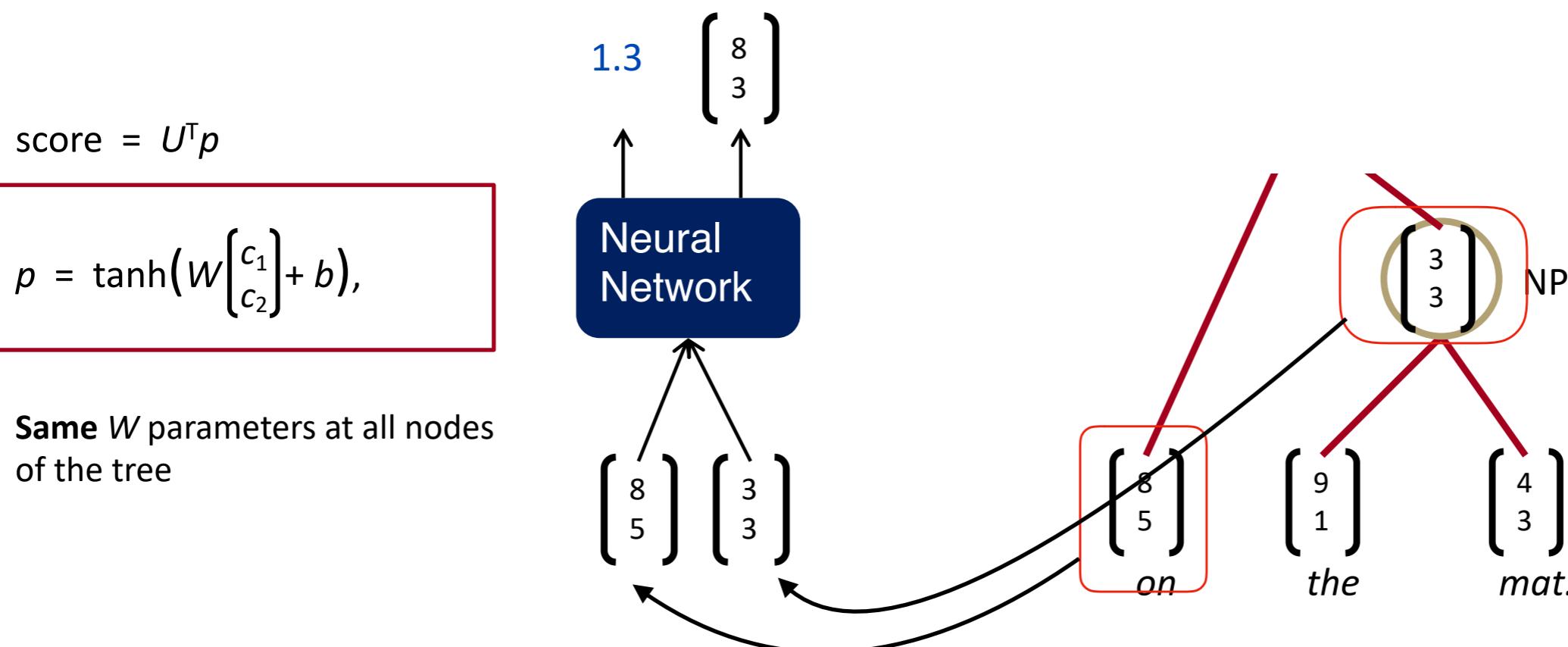
Require a tree structure to get phrase representations.

Recursive Nets for Structure Prediction

Inputs: two candidate children's representations

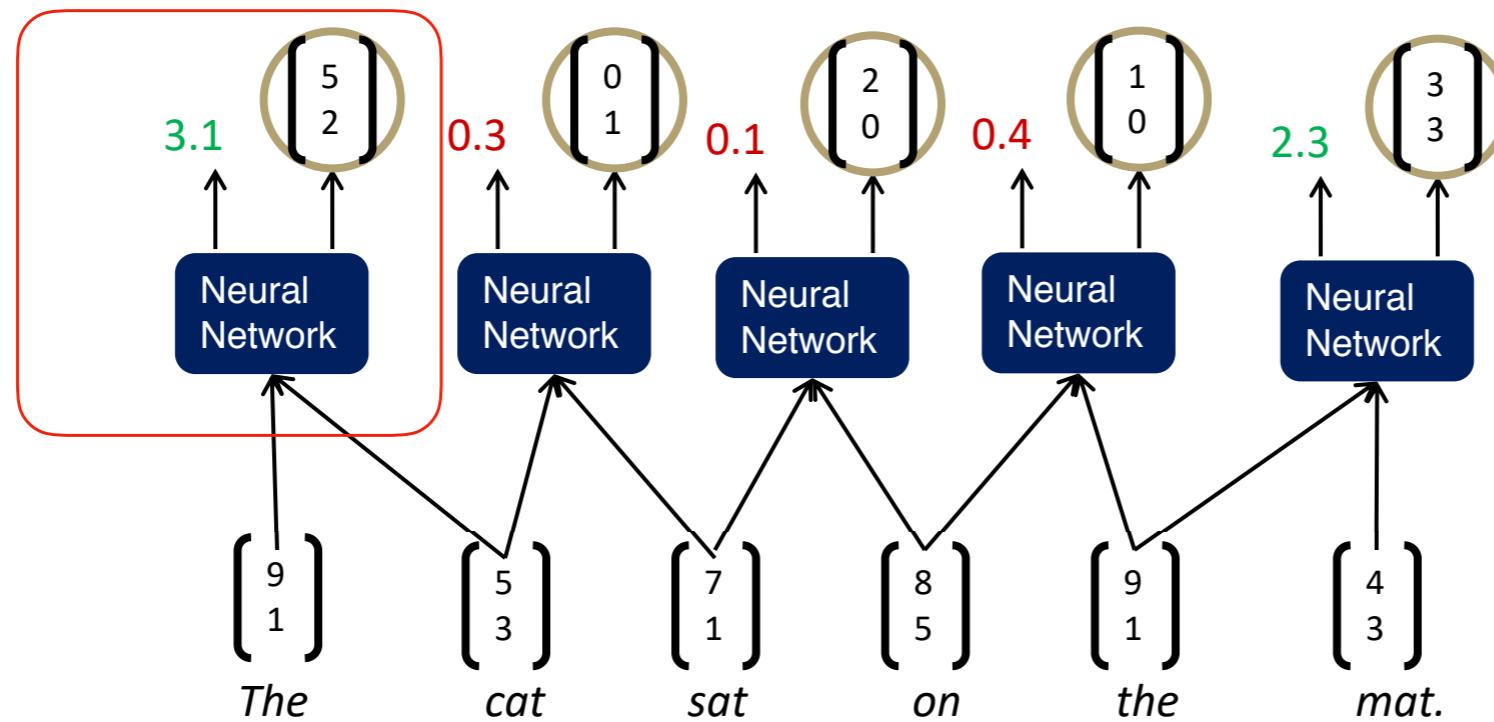
Outputs:

1. The semantic representation if the two nodes are merged.
2. Score of how plausible the new node would be.



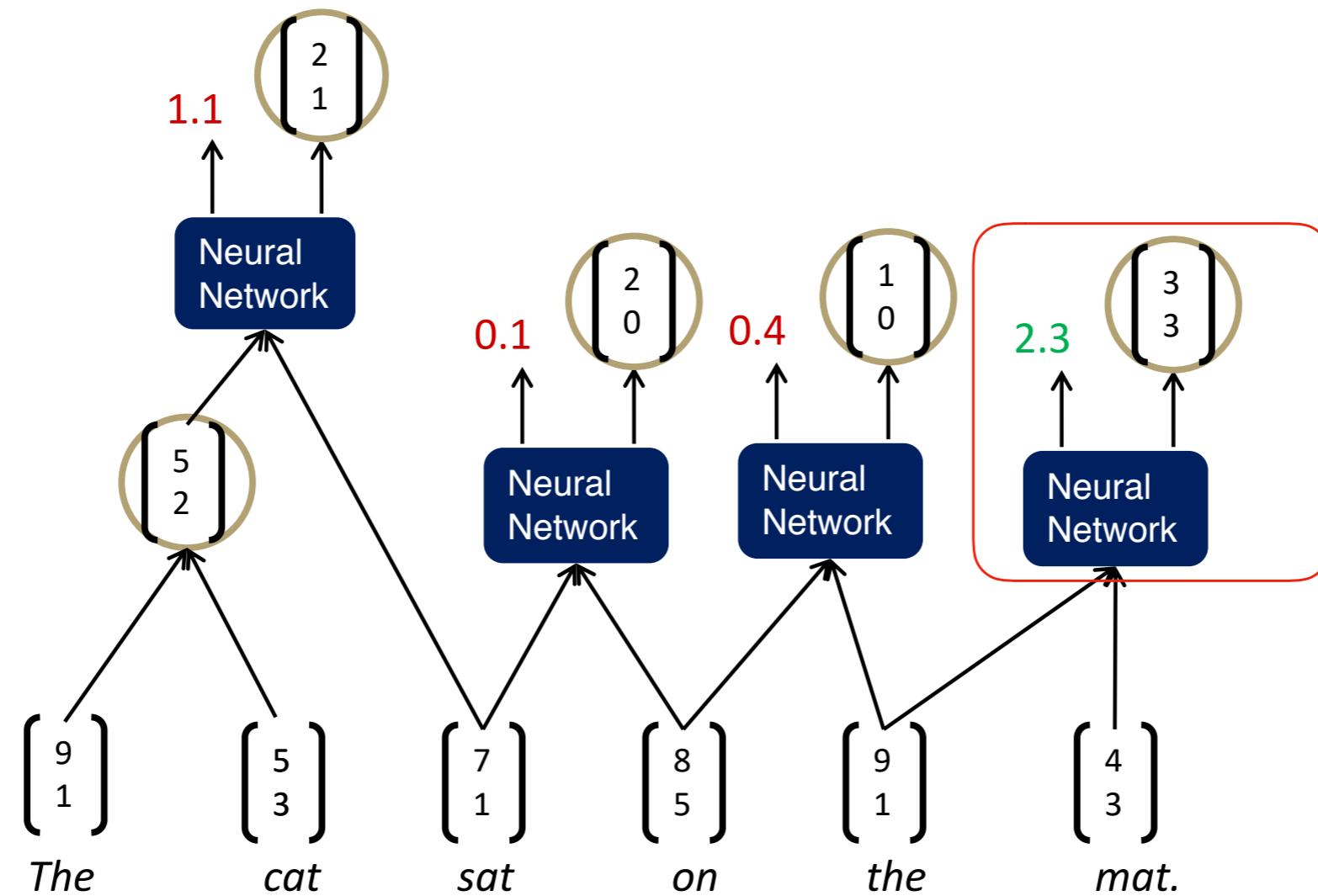
Parsing a Sentence

Greedy Approach



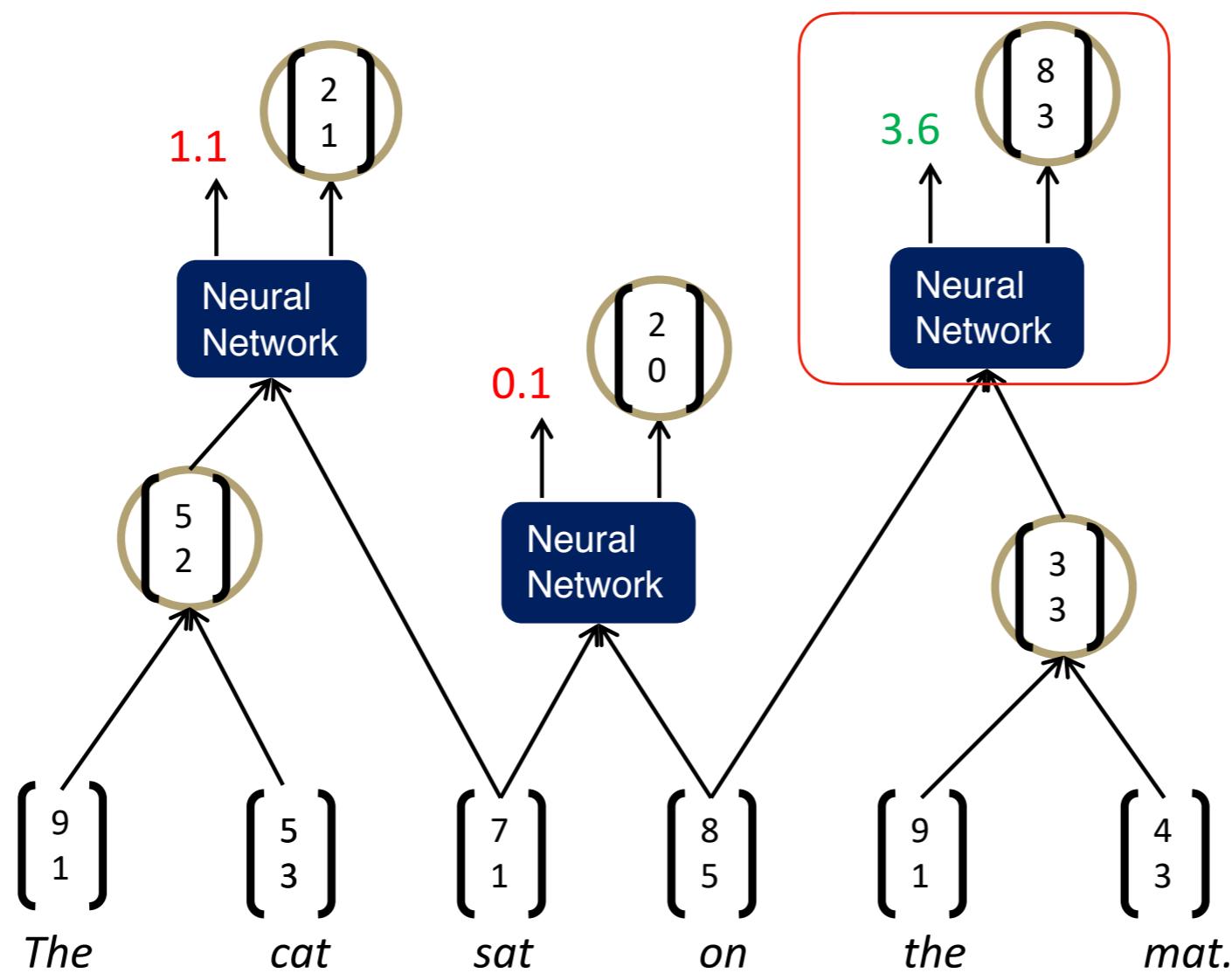
Parsing a Sentence

Greedy Approach



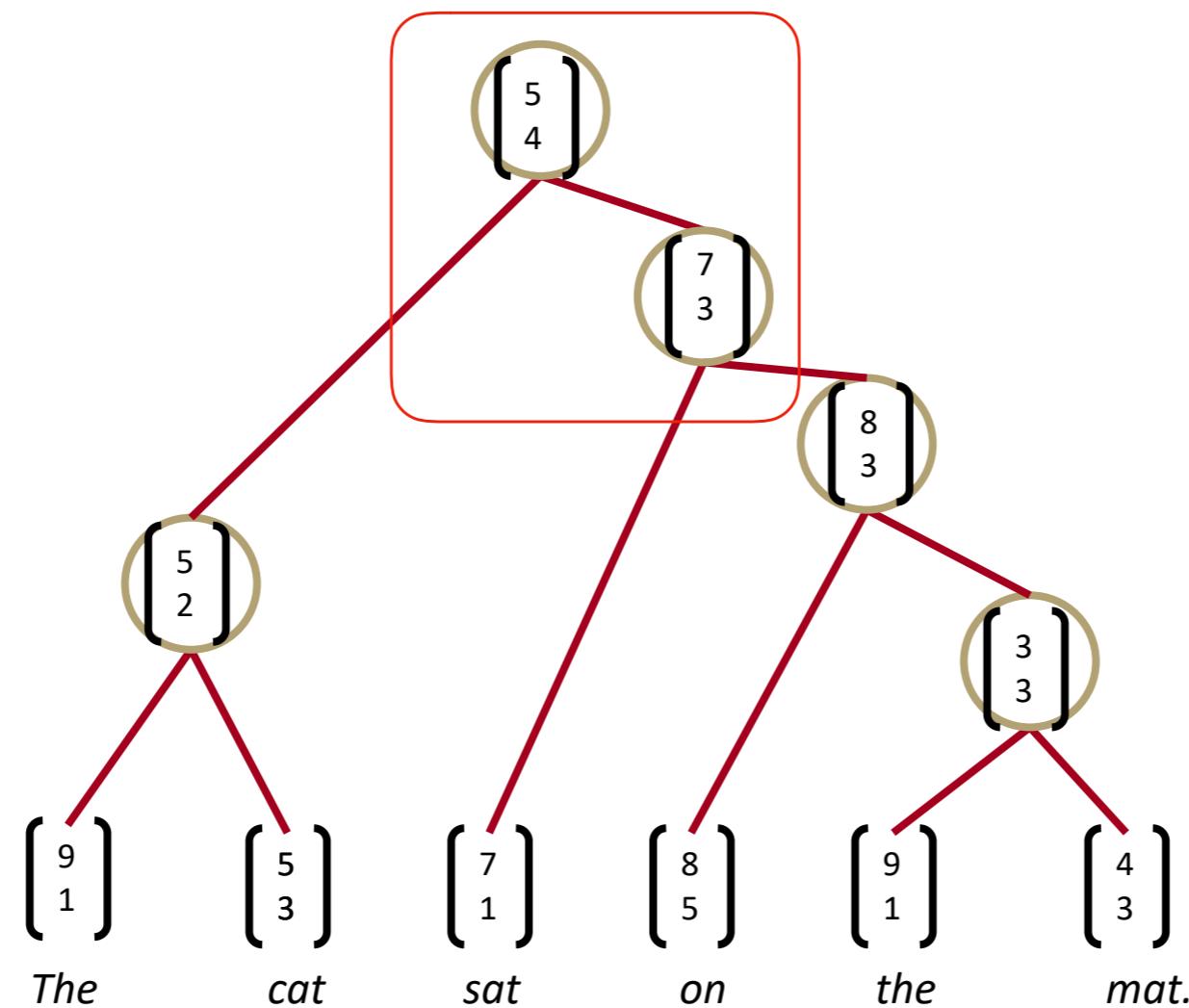
Parsing a Sentence

Greedy Approach



Parsing a Sentence

Greedy Approach



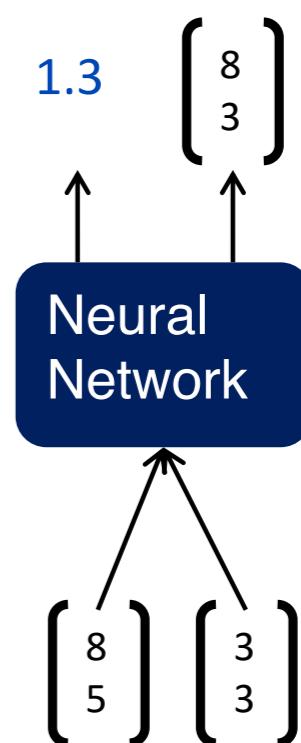
Parsing a Sentence

Max-margin Training

- The score of a tree is computed by the sum of the parsing decision scores at each node:

$$s(x, y) = \sum_{n \in nodes(y)} s_n$$

- x is sentence; y is parse tree



Parsing a Sentence

Max-margin Training

- The score of a tree is computed by the sum of the parsing decision scores at each node:

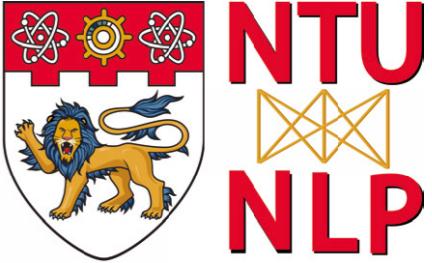
$$J = \sum_i s(x_i, y_i) - \max_{y \in A(x_i)} (s(x_i, y) + \Delta(y, y_i))$$

Beam search with chart
(instead of Greedy in previous slides)

penalizes all incorrect decisions

- Cannot do CKY (Dynamic Programming)

Today's Plan



- Syntactic structures
 - Constituency & Dependency structures
- Need for syntactic structures
- Are languages hierarchical?
- Parsing with Recursive Neural Networks
- Backpropagation through tree
- Semantic composition with Tree RNNs
- Recent Parsers

Backpropagation through tree

Same backpropagation as before

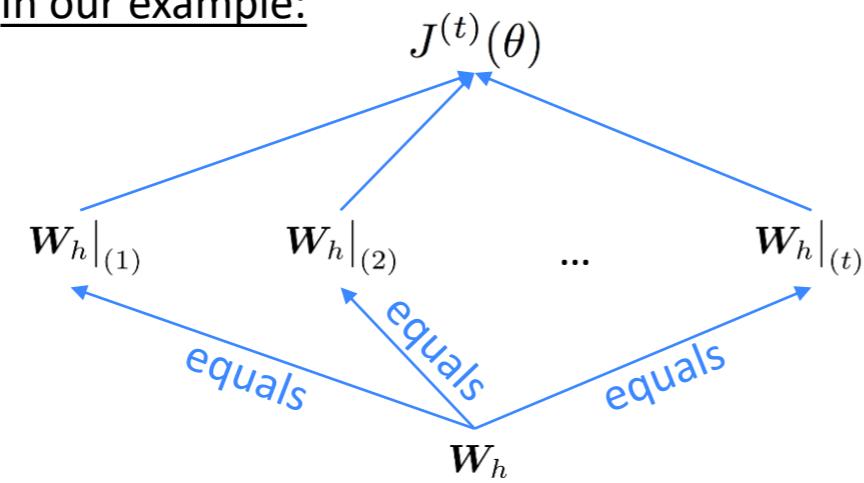
- Sum derivatives of W from all nodes (like RNN)
- Split derivatives at each node (for tree).
- Add error messages from parent + node itself

Backpropagation through tree

Same backpropagation as before

- Sum derivatives of \mathbf{W} from all nodes (like RNN)
- Split derivatives at each node (for tree).
- Add error messages from parent + node itself

In our example:



Apply the multivariable chain rule:

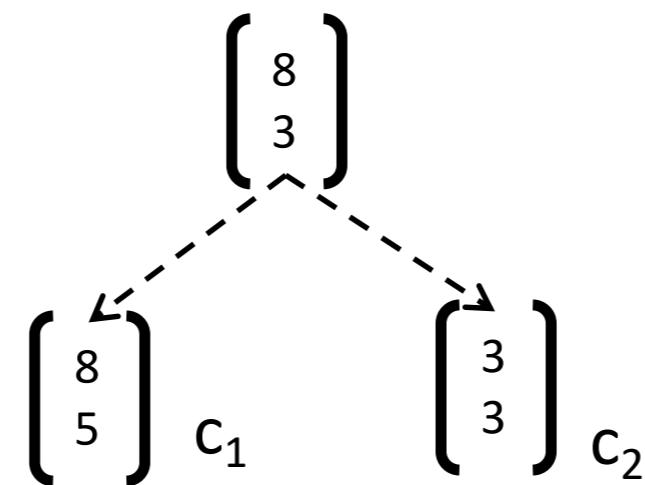
$$\begin{aligned} \frac{\partial J^{(t)}}{\partial \mathbf{W}_h} &= \sum_{i=1}^t \frac{\partial J^{(t)}}{\partial \mathbf{W}_h} \Big|_{(i)} \frac{\partial \mathbf{W}_h|_{(i)}}{\partial \mathbf{W}_h} \\ &= \sum_{i=1}^t \frac{\partial J^{(t)}}{\partial \mathbf{W}_h} \Big|_{(i)} \end{aligned}$$

= 1

Backpropagation through tree

Same backpropagation as before

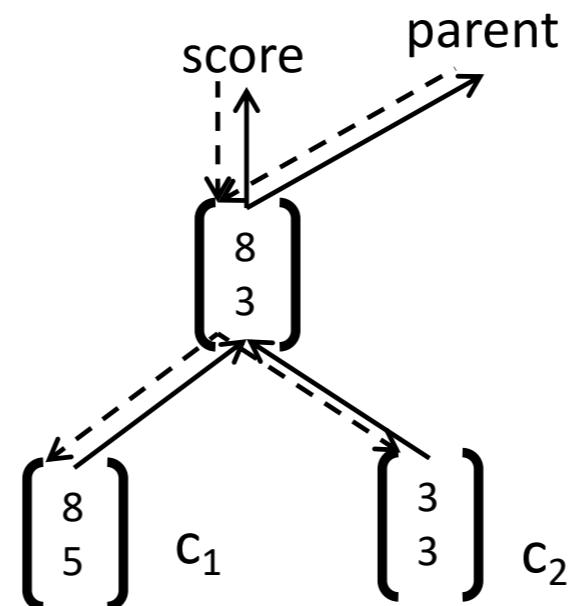
- Sum derivatives of W from all nodes (like RNN)
- **Split derivatives at each node (for tree).**
- Add error messages from parent + node itself



Backpropagation through tree

Same backpropagation as before

- Sum derivatives of W from all nodes (like RNN)
- **Split derivatives at each node (for tree).**
- Add error messages from parent + node itself



Python Code

```
def forwardProp(self,node):
    # Recursion
    ...

    # This node's hidden activation
    node.h = np.dot(self.W,np.hstack([node.left.h, node.right.h])) + self.b
    # Relu
    node.h[node.h<0] = 0

    # Softmax
    node.probs = np.dot(self.Ws,node.h) + self.bs
    node.probs -= np.max(node.probs)
    node.probs = np.exp(node.probs)
    node.probs = node.probs/np.sum(node.probs)
```

Python Code

```

def backProp(self, node, error=None):
    # Softmax grad
    deltas = node.probs
    deltas[node.label] -= 1.0
    self.dWs += np.outer(deltas, node.h)
    self.dbs += deltas
    deltas = np.dot(self.Ws.T, deltas)

    # Add deltas from above
    if error is not None:
        deltas += error

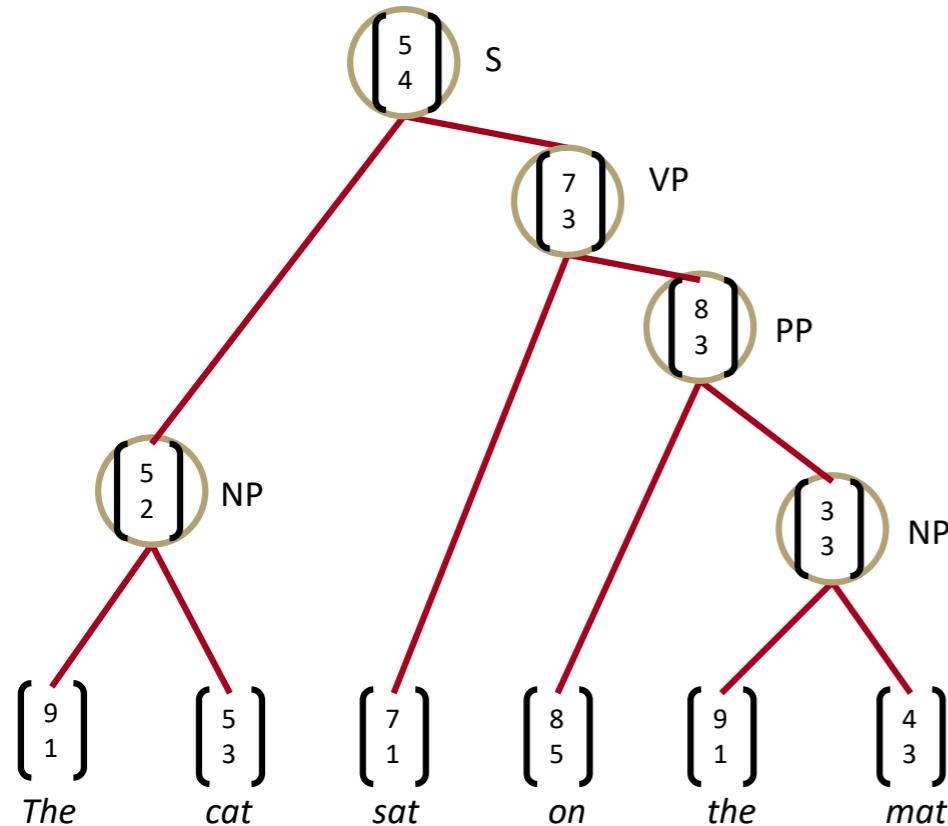
# f'(z) now:
deltas *= (node.h != 0)

# Update word vectors if leaf node:
if node.isLeaf:
    self.dL[node.word] += deltas
    return

# Recursively backprop
if not node.isLeaf:
    self.dW += np.outer(deltas, np.hstack([node.left.h, node.right.h]))
    self.db += deltas
    # Error signal to children
    deltas = np.dot(self.W.T, deltas)
    self.backProp(node.left, deltas[:self.hiddenDim])
    self.backProp(node.right, deltas[self.hiddenDim:])

```

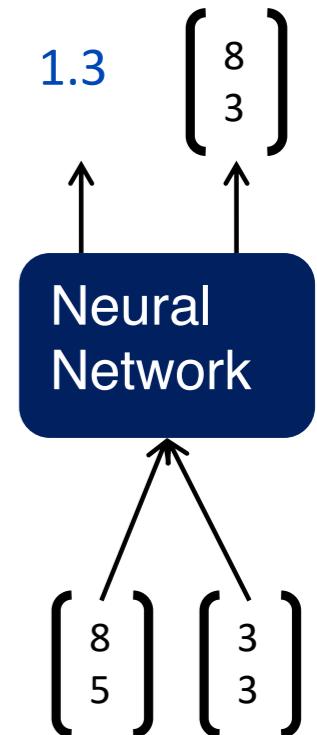
Limitations of Basic Tree-RNNs



$$\text{score} = U^T p$$

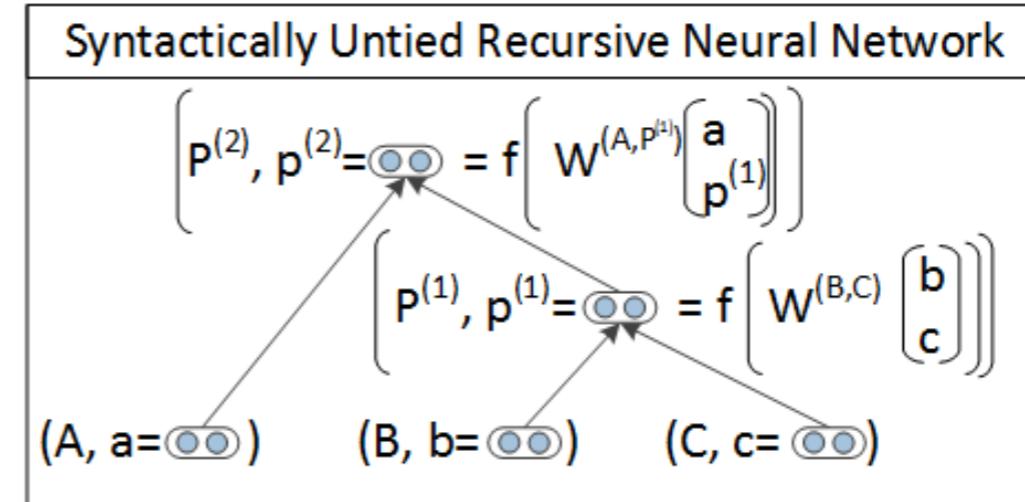
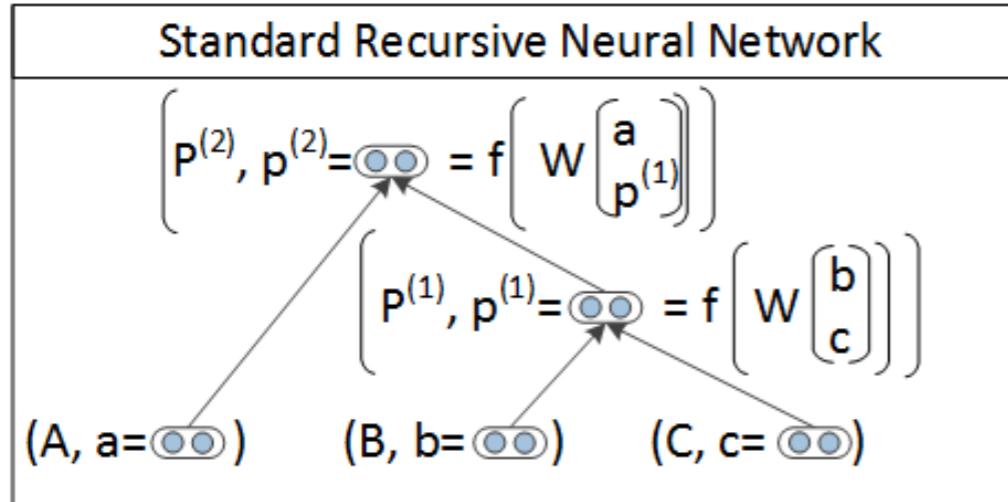
$$p = \tanh\left(W \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} + b\right),$$

Same W parameters at all nodes of the tree



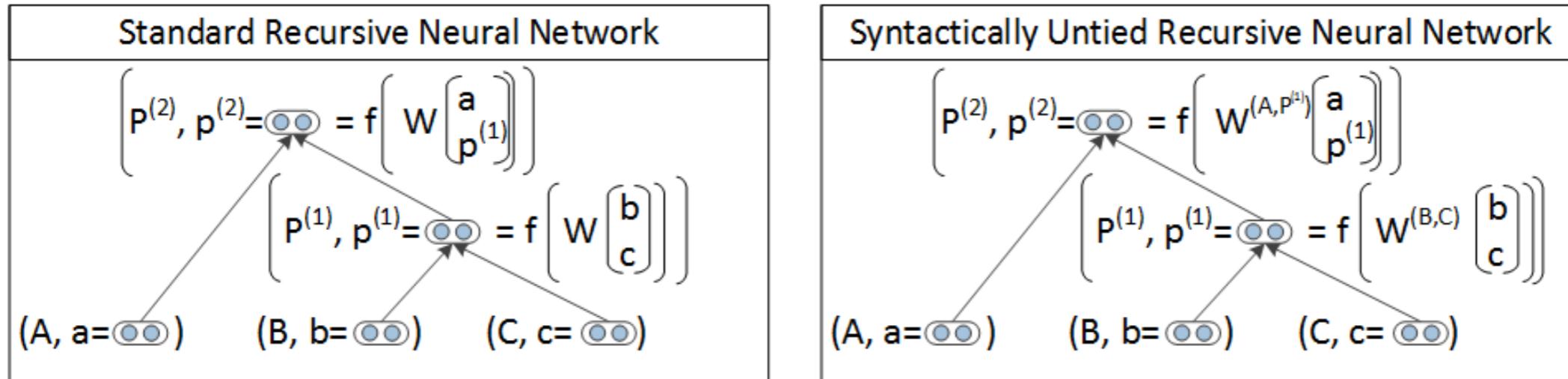
- Single weight matrix TreeRNN could capture some phenomena but not adequate for more complex, higher order composition and parsing long sentences
- There is no real interaction between the input words
- The composition function is the same for all syntactic categories, punctuation, etc

Version 2: Syntactically-Untied RNN



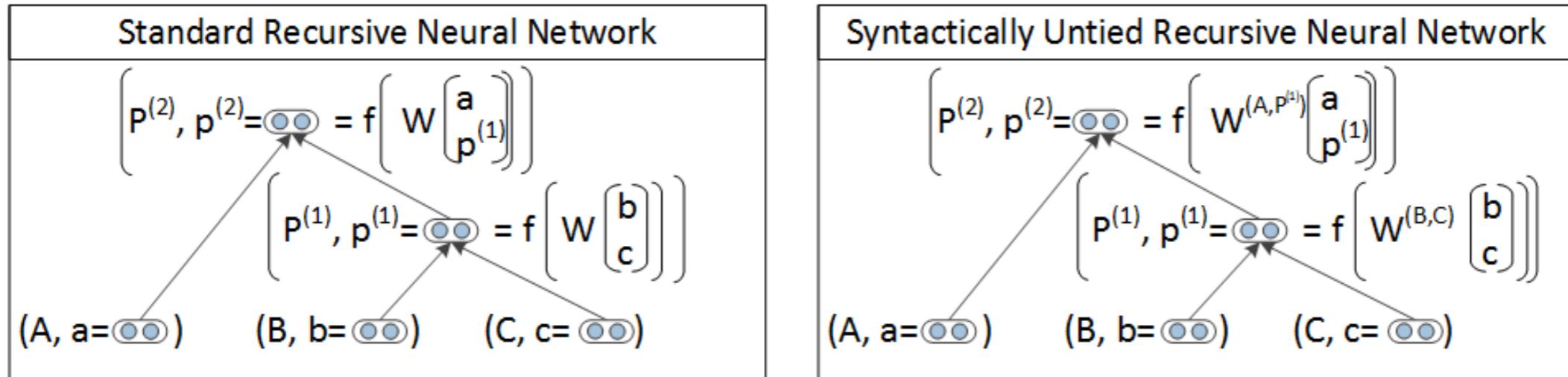
- A symbolic Context-Free Grammar (CFG) backbone is adequate for basic syntactic structure
- We use the discrete syntactic categories of the children to choose the composition matrix
- A TreeRNN can do better with different composition matrix for different syntactic environments
- The result gives us a better semantics

Version 2: Syntactically-Untied RNN



- **Problem:** Speed. Every candidate score in beam search needs a matrix-vector product.
- **Solution:** Compute score only for a subset of trees coming from a simpler, faster model (PCFG)
 - Prunes very unlikely candidates for speed
 - Provides coarse syntactic categories of the children for each beam candidate
- Compositional Vector Grammar = PCFG + TreeRNN

Version 2: Syntactically-Untied RNN



- Compositional Vector Grammar = PCFG + TreeRNN
- Scores at each node computed by combination of PCFG and SU-RNN:

$$s(p^{(1)}) = (v^{(B,C)})^T p^{(1)} + \log P(P_1 \rightarrow B \ C)$$

- Factoring discrete and continuous parsing in one model:

$$\begin{aligned} P((P_1, p_1) \rightarrow (B, b)(C, c)) \\ = P(p_1 \rightarrow b \ c | P_1 \rightarrow B \ C) P(P_1 \rightarrow B \ C) \end{aligned}$$

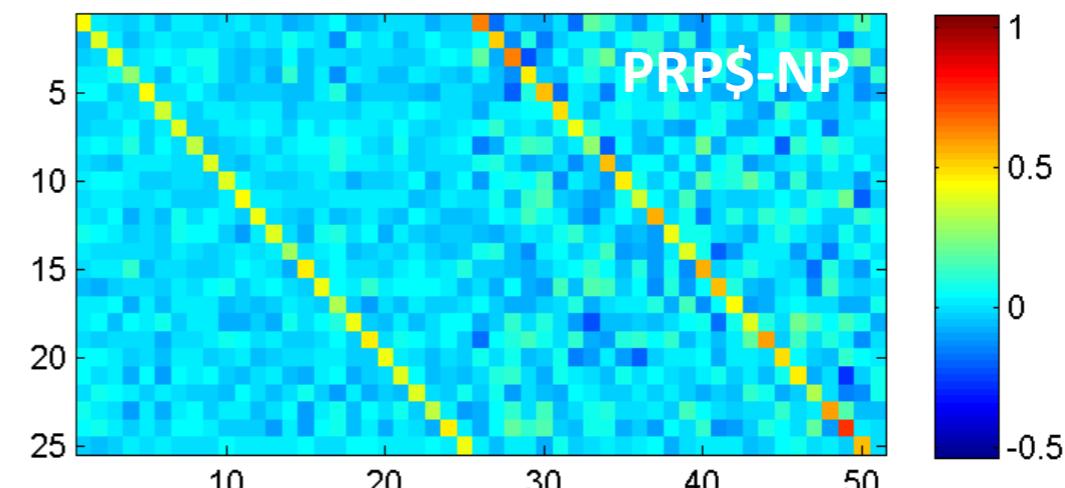
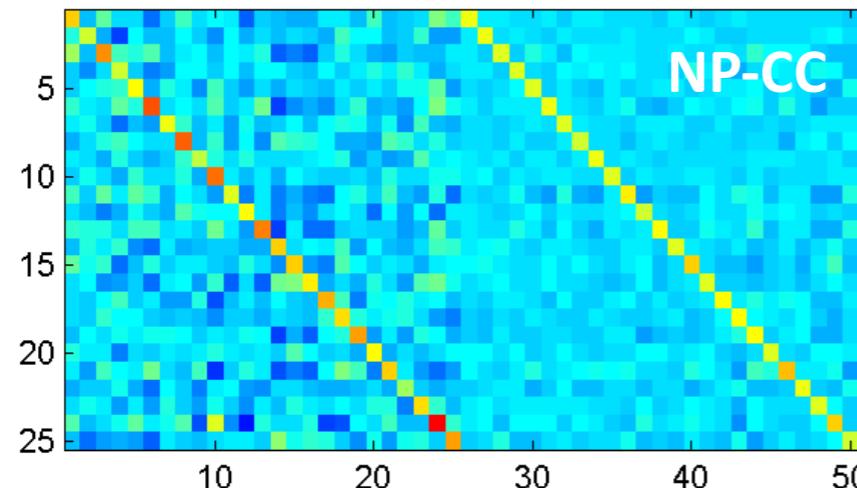
Parsing Results

Parser	Test, All Sentences
Stanford PCFG, (Klein and Manning, 2003a)	85.5
Stanford Factored (Klein and Manning, 2003b)	86.6
Factored PCFGs (Hall and Klein, 2012)	89.4
Collins (Collins, 1997)	87.7
SSN (Henderson, 2004)	89.4
Berkeley Parser (Petrov and Klein, 2007)	90.1
CVG (RNN) (Socher et al., ACL 2013)	85.0
CVG (SU-RNN) (Socher et al., ACL 2013)	90.4
Charniak - Self Trained (McClosky et al. 2006)	91.0
Charniak - Self Trained-ReRanked (McClosky et al. 2006)	92.1

- Standard WSJ split, labeled F1
- Based on simple PCFG with fewer states
- Fast pruning of search space, few matrix-vector products
- 3.8% higher F1

Analysis of Learned Representations

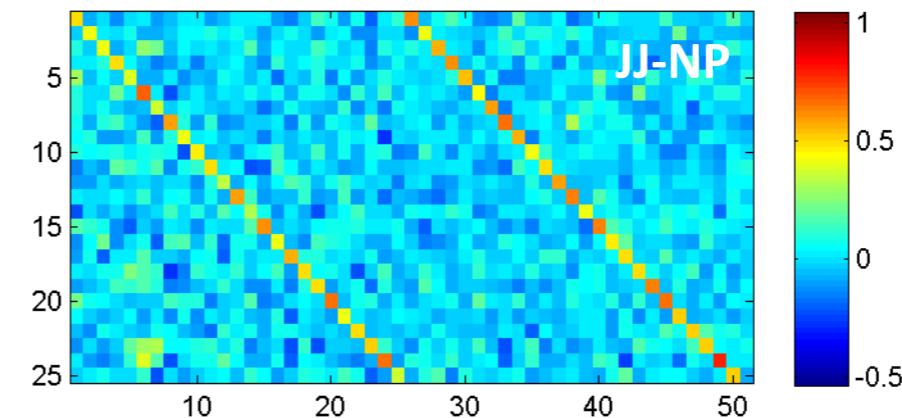
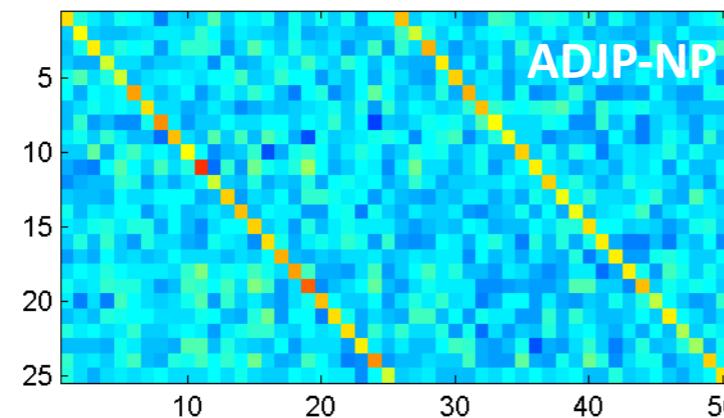
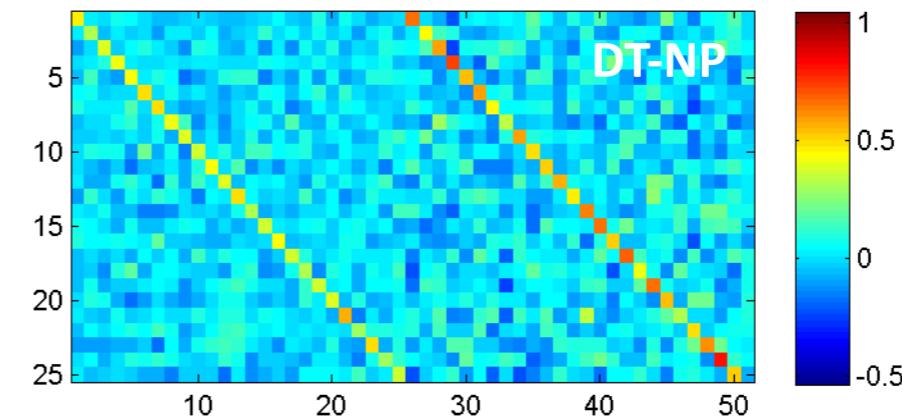
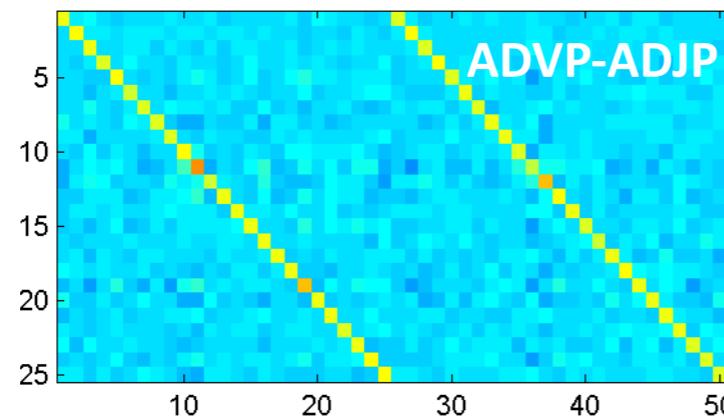
Initialization $W^{(\cdot)} = 0.5[I_{n \times n} I_{n \times n} 0_{n \times 1}] + \epsilon$



- Learns soft notion of head words

Analysis of Learned Representations

Initialization $W^{(\cdot)} = 0.5[I_{n \times n} I_{n \times n} 0_{n \times 1}] + \epsilon$



- Learns soft notion of head words

Analysis of Learned Representations

All the figures are adjusted for seasonal variations

1. All the numbers are adjusted for seasonal fluctuations
2. All the figures are adjusted to remove usual seasonal patterns

Knight-Ridder wouldn't comment on the offer

1. Harsco declined to say what country placed the order
2. Coastal wouldn't disclose the terms

Sales grew almost 7% to \$UNK m. from \$UNK m.

1. Sales rose more than 7% to \$94.9 m. from \$88.3 m.
2. Sales surged 40% to UNK b. yen from UNK b.

Version 3: Compositionality Through Recursive Matrix-Vector Spaces

So far:

$$p = \tanh(w \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} + b)$$

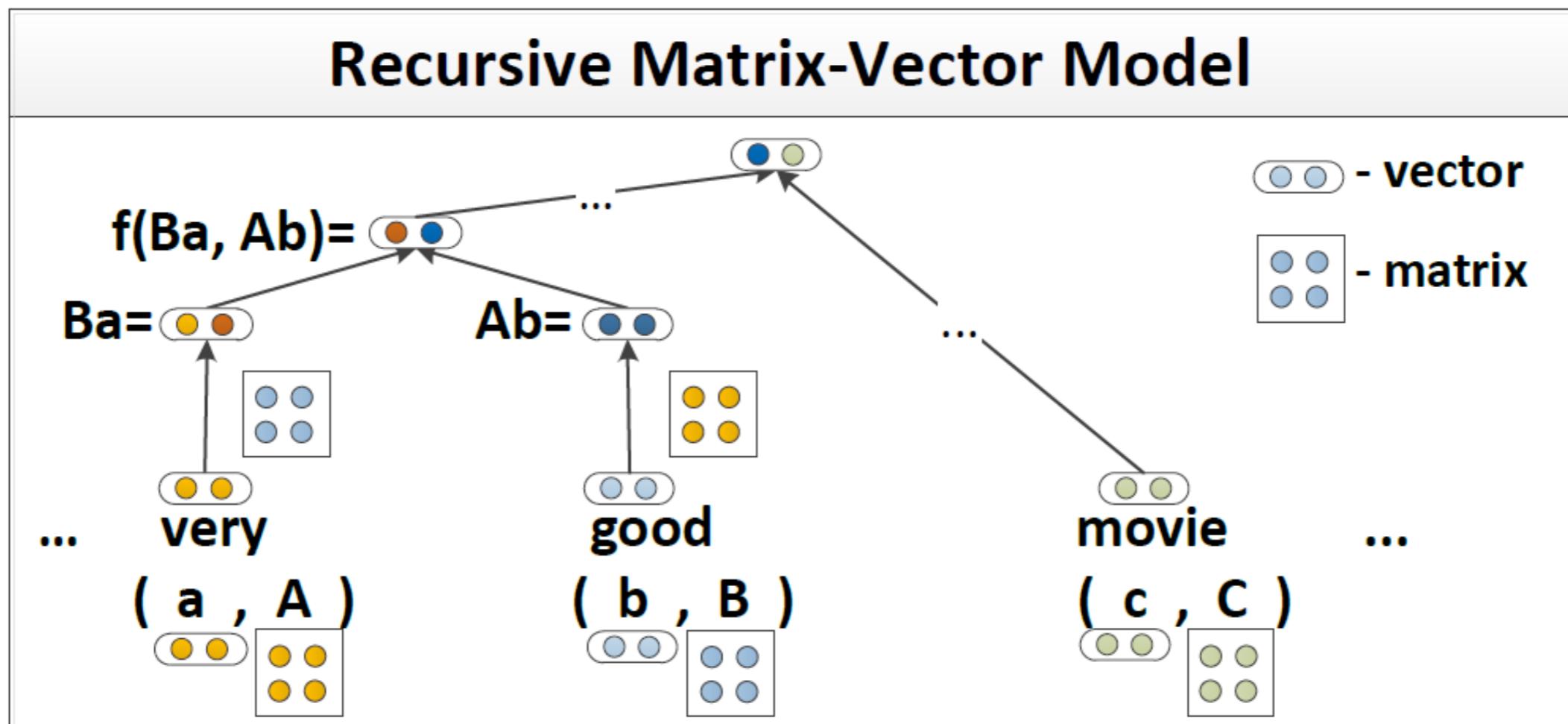
We untied the weights in Version 2 to make the composition function more powerful

But what if words act mostly as an operator, e.g. “very” in
very good

Version 3: Compositionality Through Recursive Matrix-Vector Spaces

$$p = \tanh\left(W \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} + b\right)$$

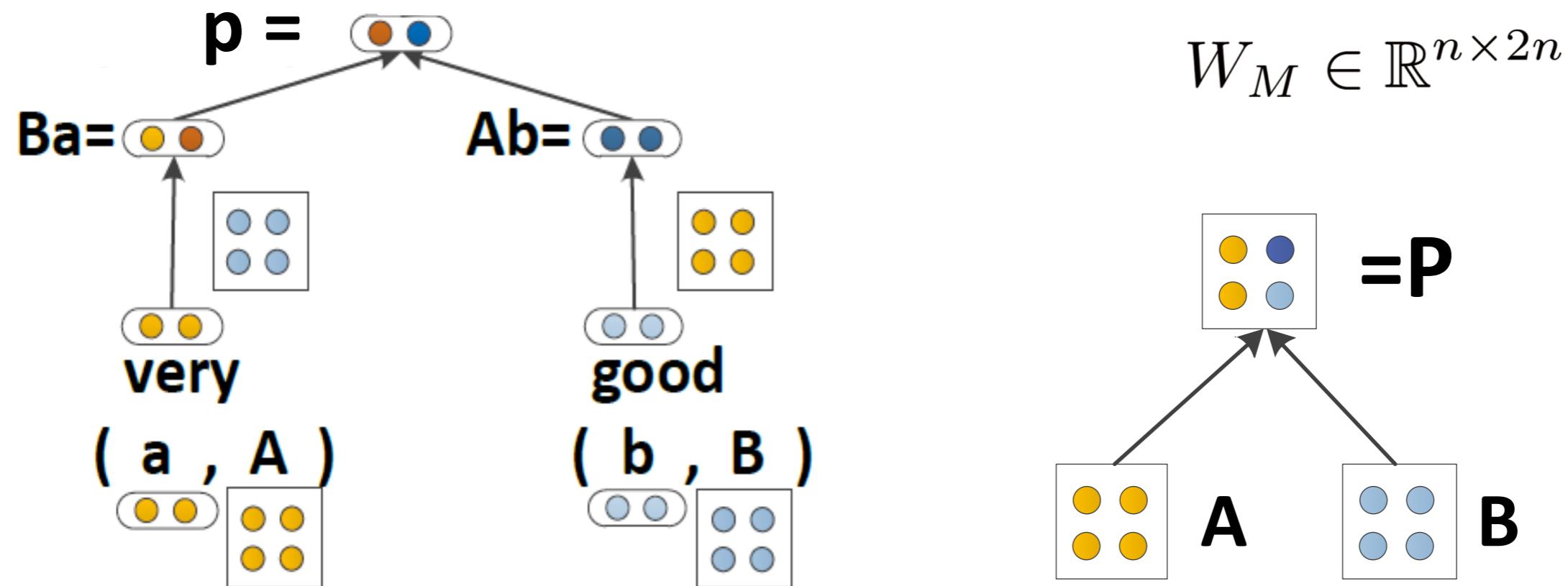
$$p = \tanh\left(W \begin{bmatrix} c_2 c_1 \\ c_1 c_2 \end{bmatrix} + b\right)$$



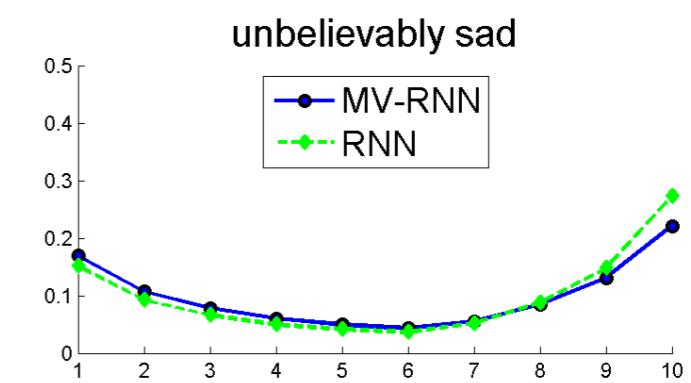
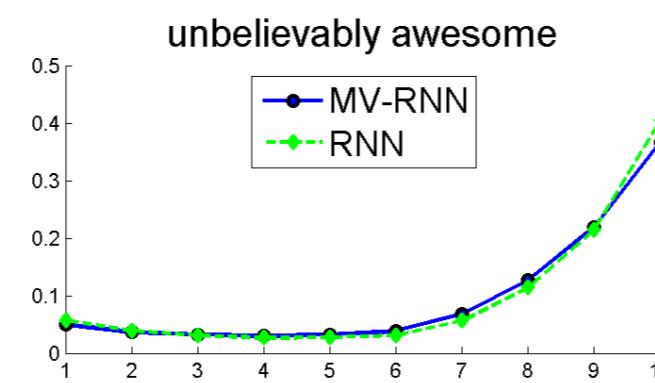
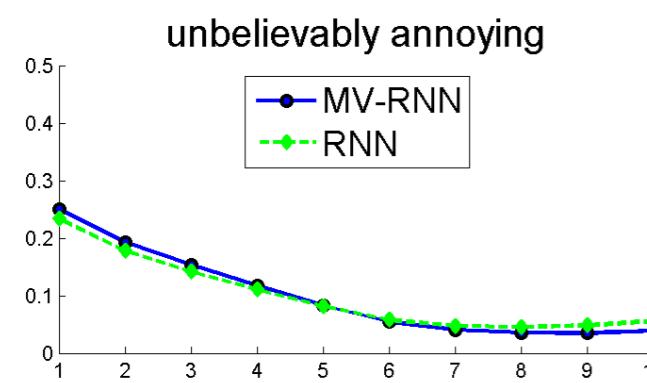
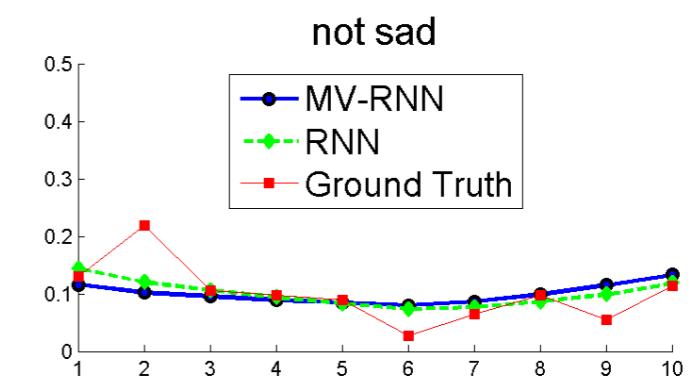
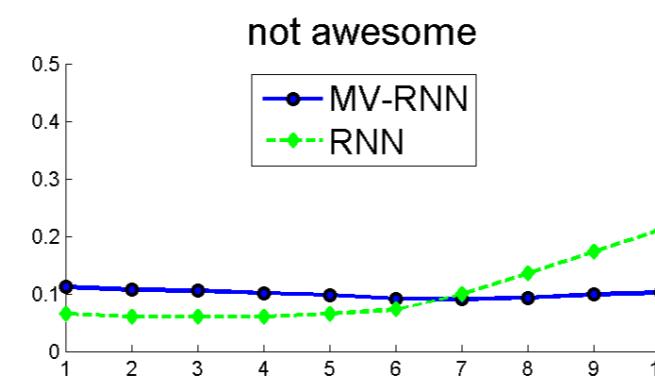
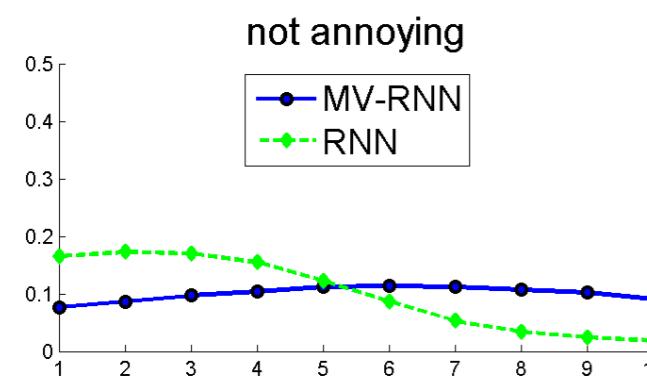
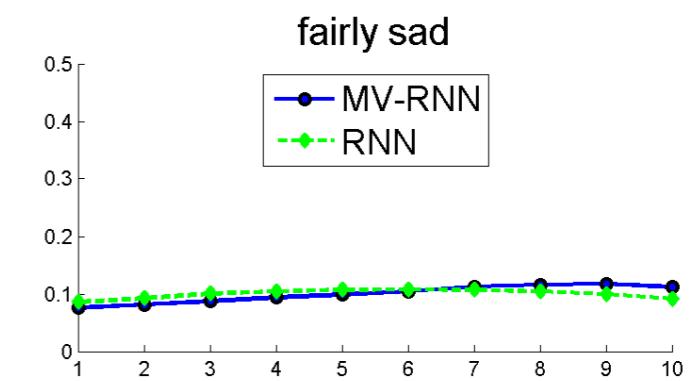
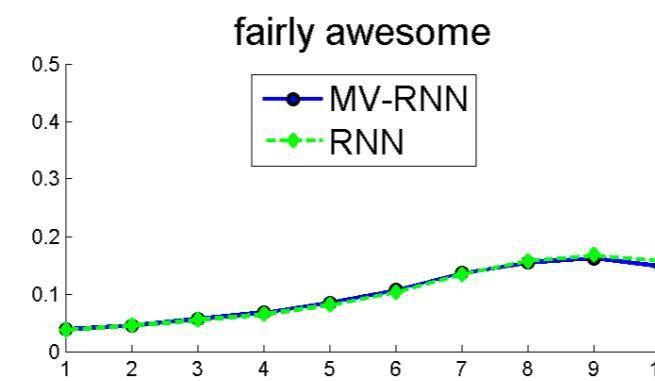
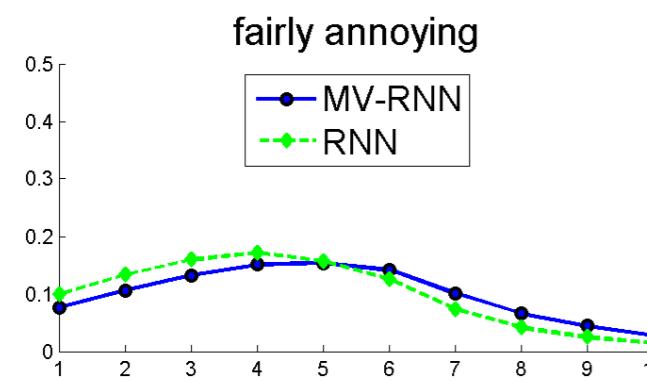
Version 3: Compositionality Through Recursive Matrix-Vector Spaces

$$p = f \left(W \begin{bmatrix} Ba \\ Ab \end{bmatrix} \right)$$

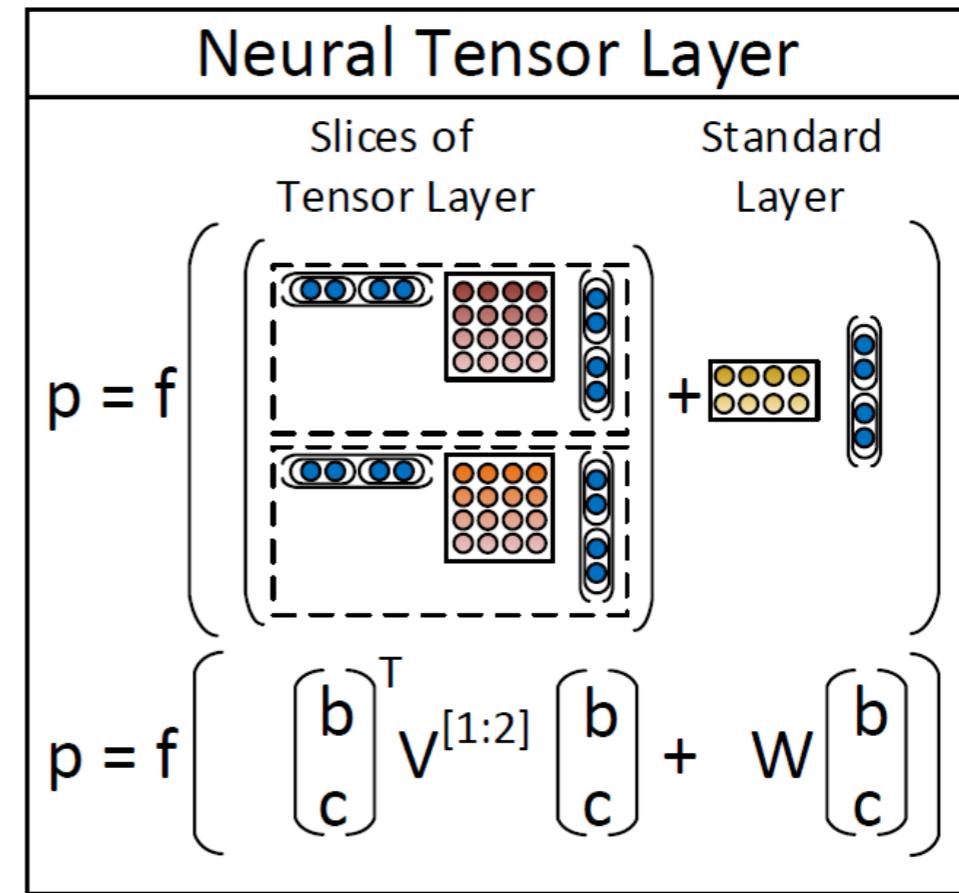
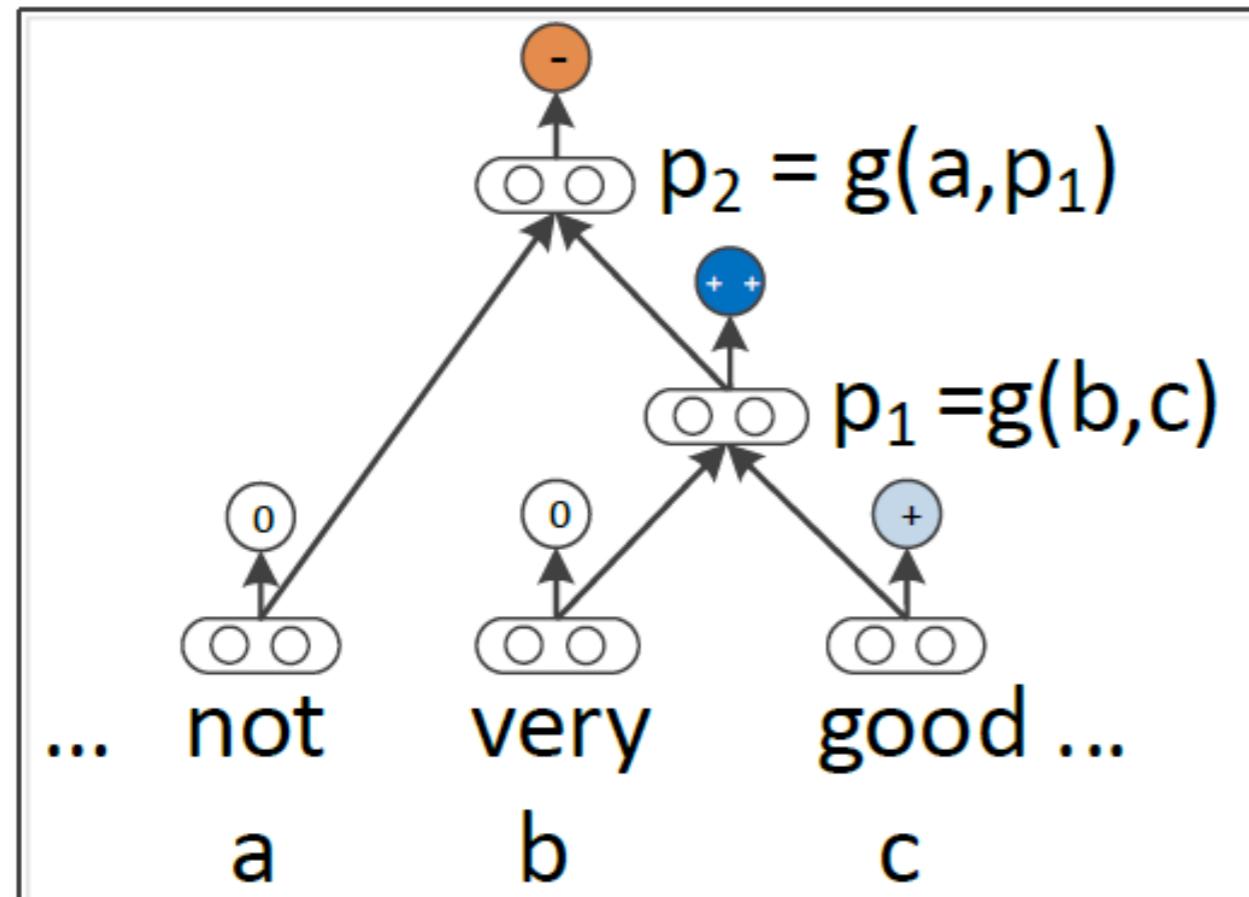
$$P = g(A, B) = W_M \begin{bmatrix} A \\ B \end{bmatrix}$$



Predicting Sentiment Distributions



Version 4: Recursive Neural Tensor Network



- Less parameters than MV-RNN
- Allows the two word or phrase vectors to interact multiplicatively

Sentiment Analysis

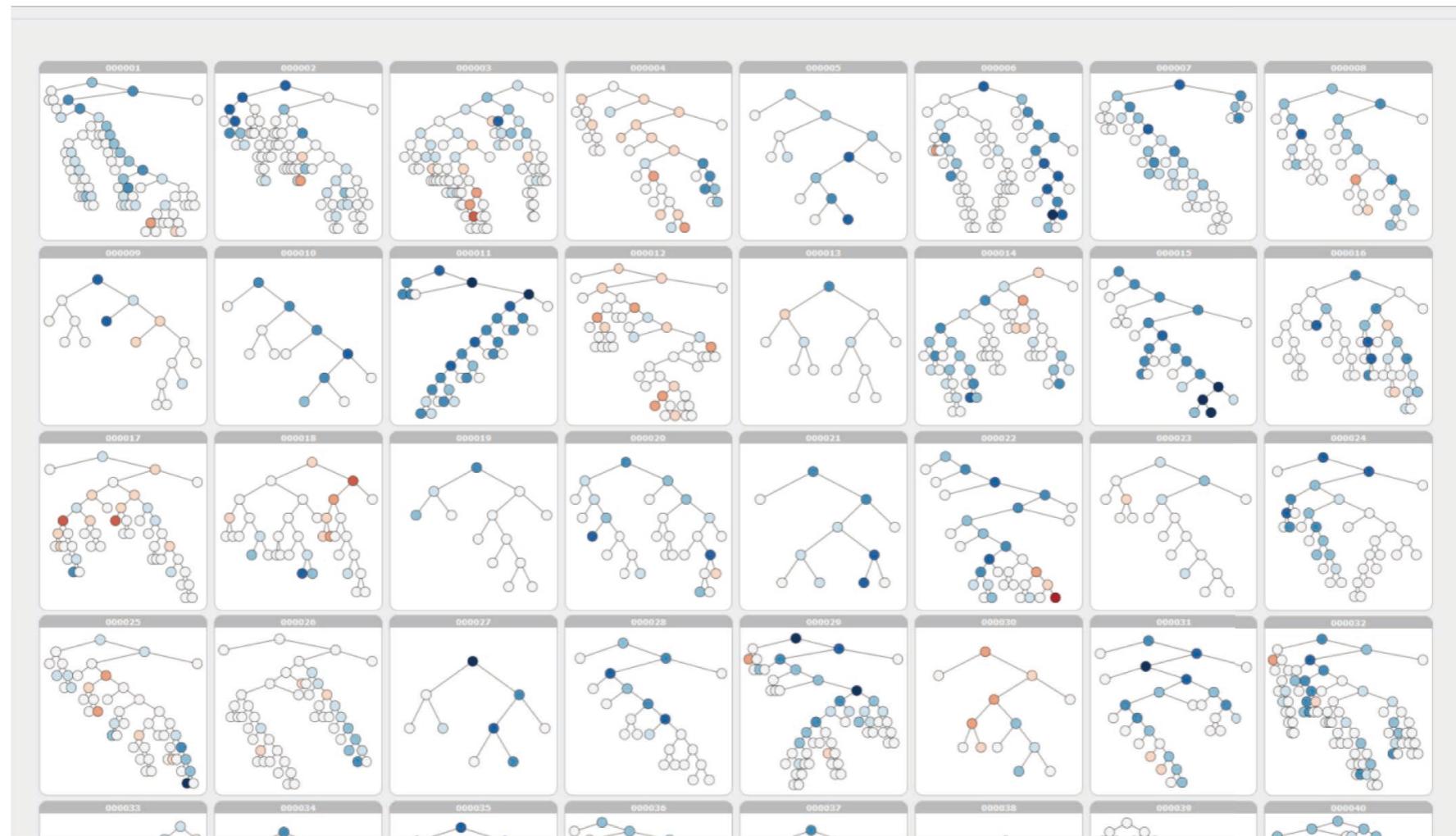
- Is the tone of a piece of text positive, negative, or neutral?

I can't remember the last time I saw a movie that contained as many genres as 'Parasite'. The movie starts out almost like an 'Ocean's Eleven' heist film and then expands into a comedy, mystery, thriller, drama, romance, crime and even horror film. It really did have everything and it was strikingly good at all of them too.

Source: IMDB

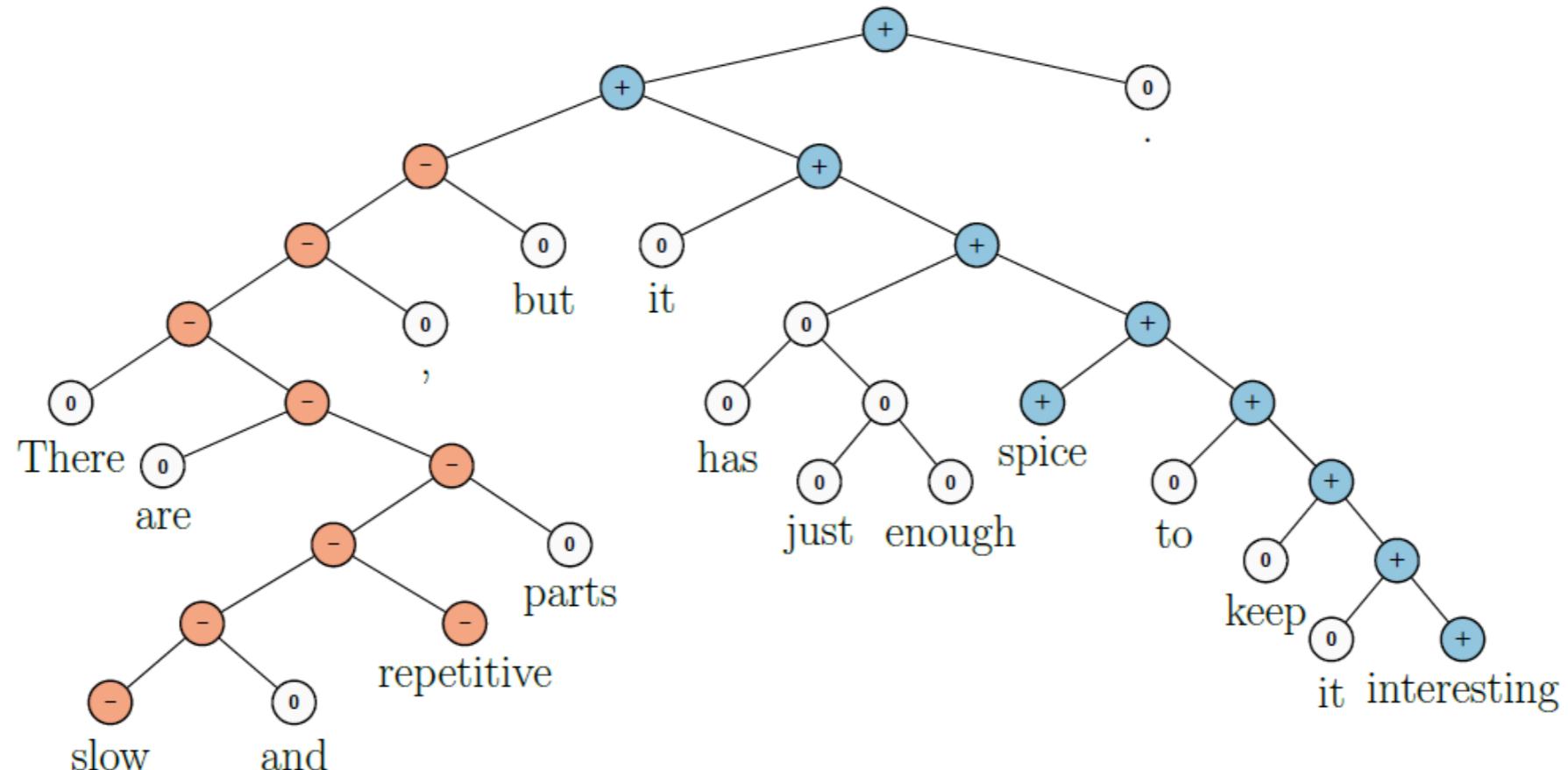
Sentiment Analysis: Stanford Sentiment Treebank

- 215,154 phrases labeled in 11,855 sentences
- Can actually train and test compositions



<https://nlp.stanford.edu/sentiment/index.html>

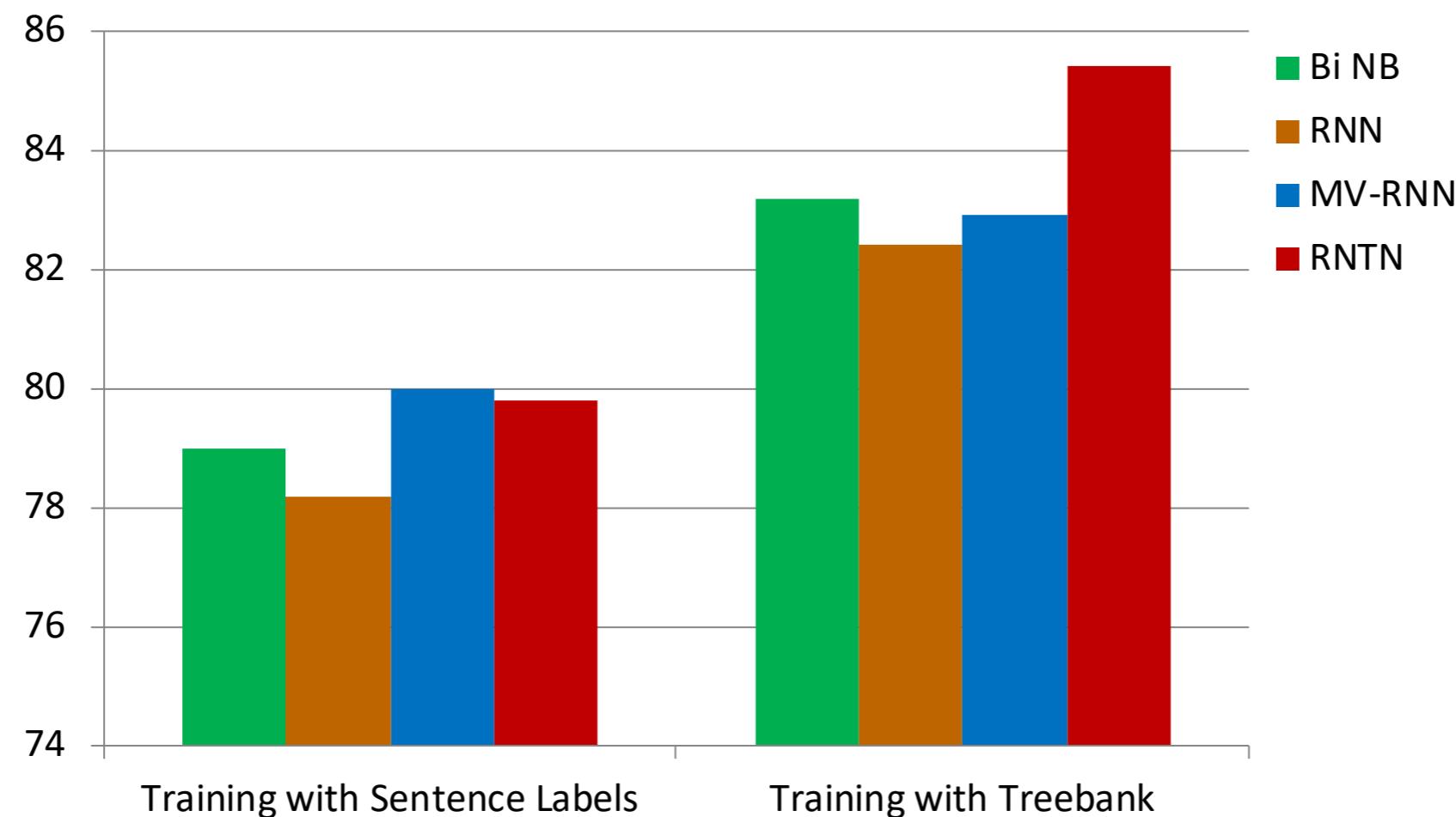
Sentiment Analysis: Stanford Sentiment Treebank



<https://nlp.stanford.edu/sentiment/index.html>

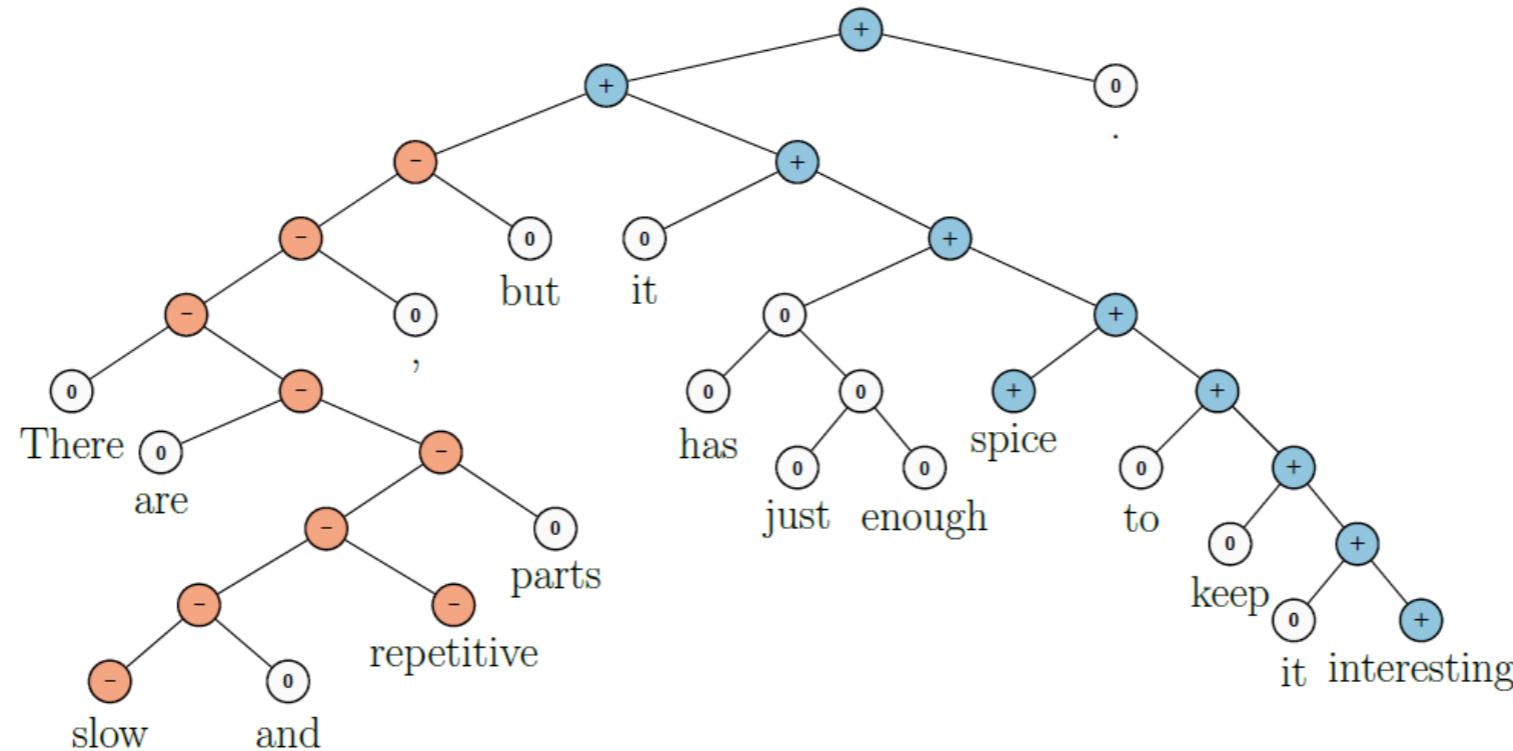
Sentiment Analysis

Classifying Sentences: Accuracy improves to 85.4



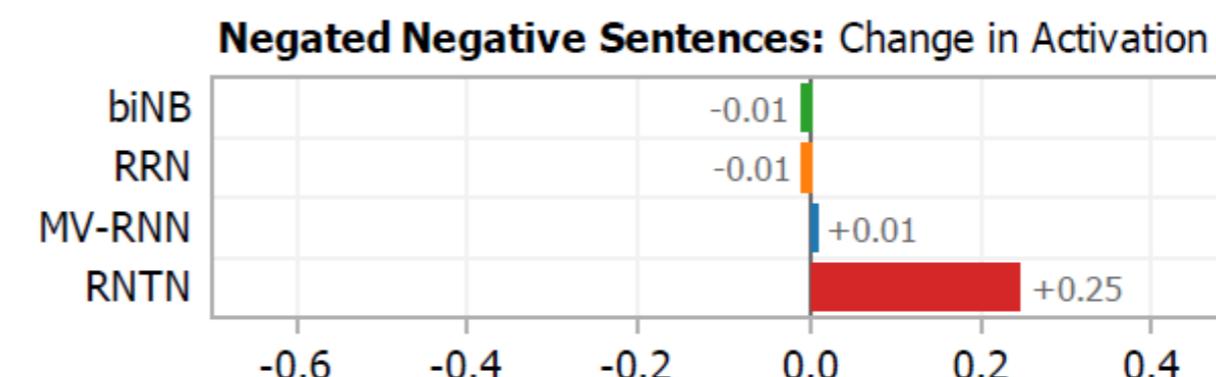
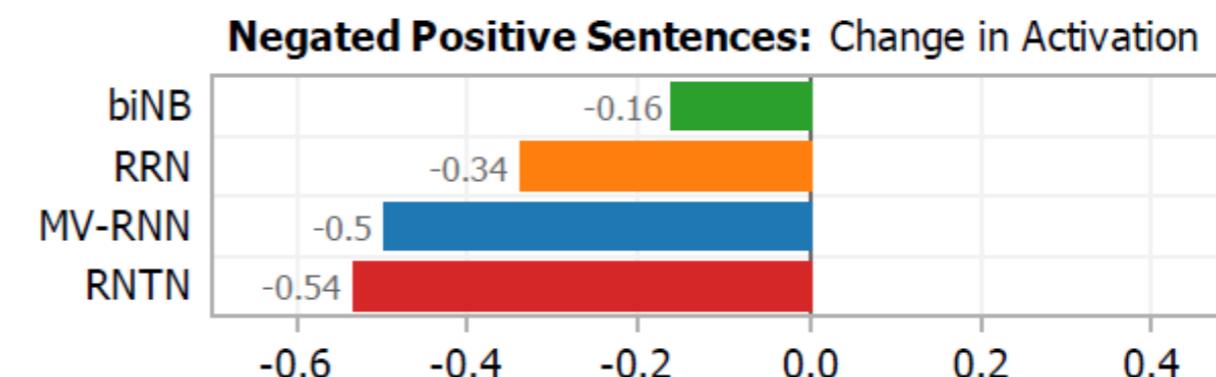
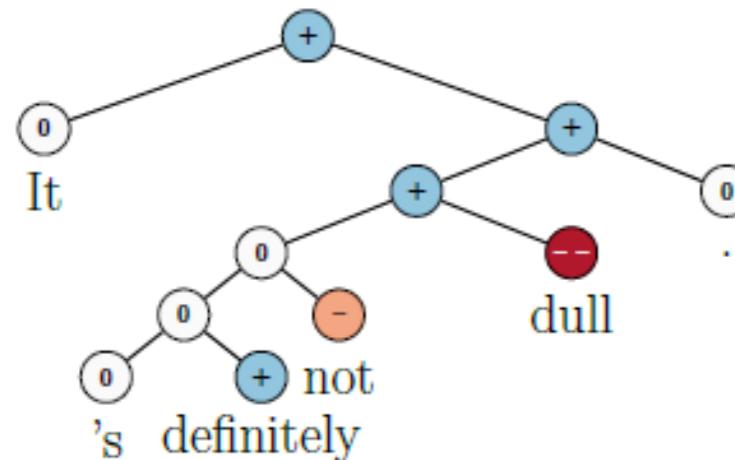
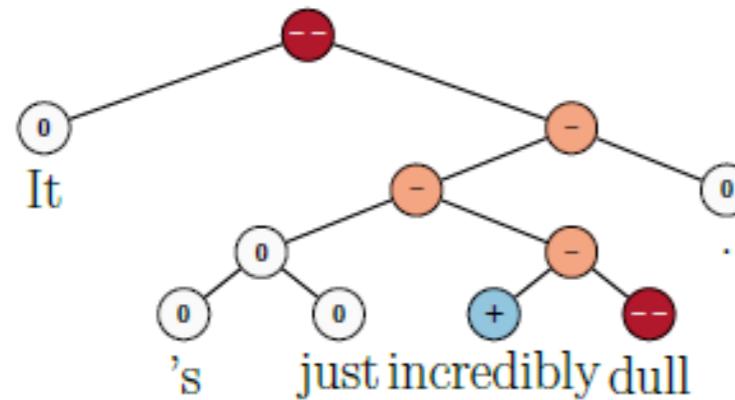
Positive/Negative

Sentiment Analysis



- RNTN can capture constructions like X but Y
- RNTN accuracy of 72%, compared to MV-RNN (65%), bi-word NB (58%) and RNN (54%)

Sentiment Analysis



When negating negatives, positive activation should increase!

Version 5: Tree LSTM

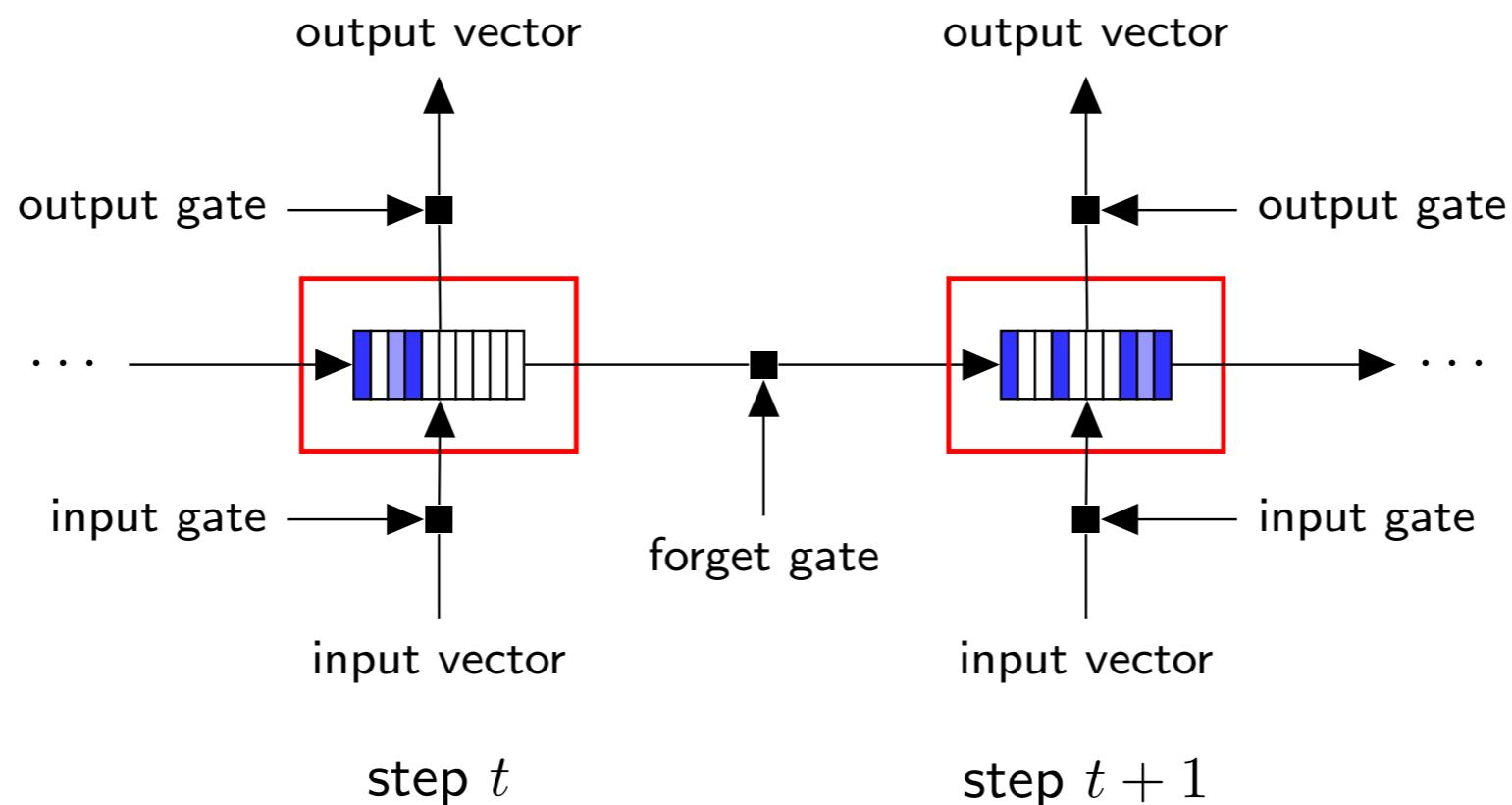
- ① Generalization of sequential LSTM to tree structures.
- ② Similar to recursive networks, it uses the **tree structure** to recursively **encode** node representations.
- ③ Generalizes to any number of children (not just binary).

● Goal:

To represent the meaning of a sentence as a point in a (high-dimensional, continuous) vector space in a way that accurately handles semantic composition and sentence meaning

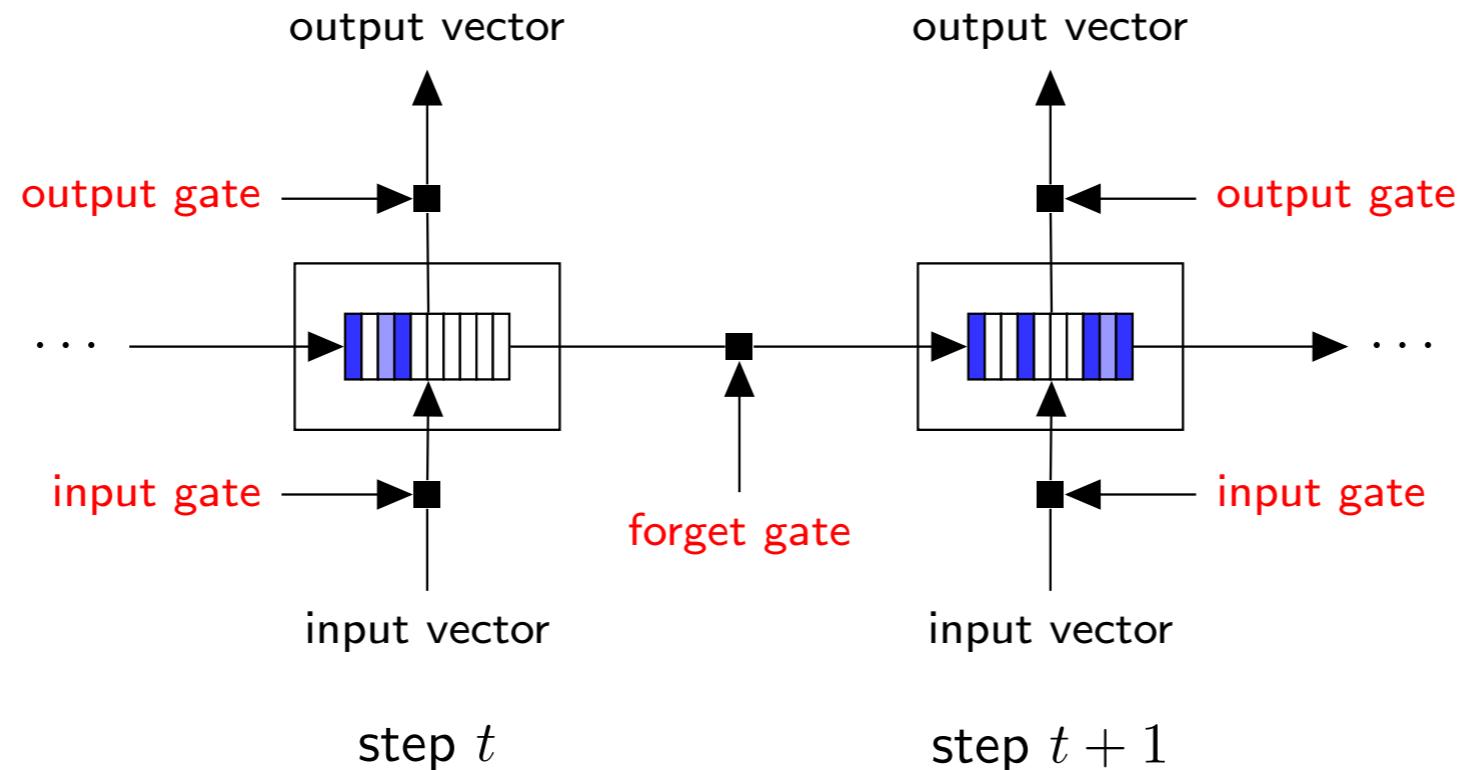
Sequential LSTM

ϕ is the composition function to be learned from data.



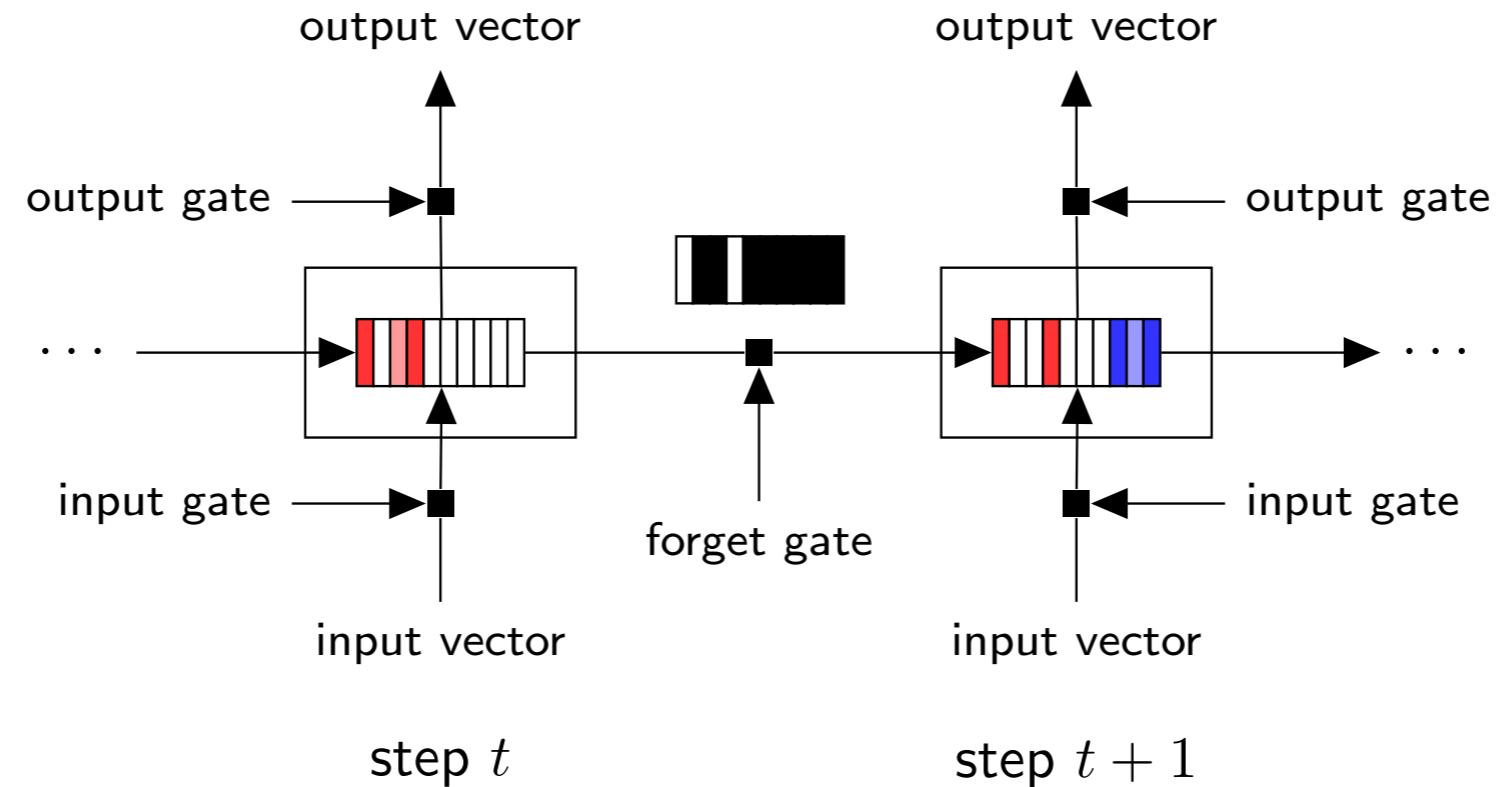
- ① LSTM is a particular parameterization of ϕ .
- ② **Memory cell**: a vector representing the inputs seen so far. Its state can be preserved over many time steps.

Sequential LSTM



- ① **input/output/forget** gates are vectors in $[0, 1]^d$, used as **soft mask** (i.e., by element-wise multiplication) for selective memory read/write/information-propagation.
- ② Each gate is **parameterized and conditioned** on the current input x_t and previous hidden state \mathbf{h}_t (*output vector* in Fig.).

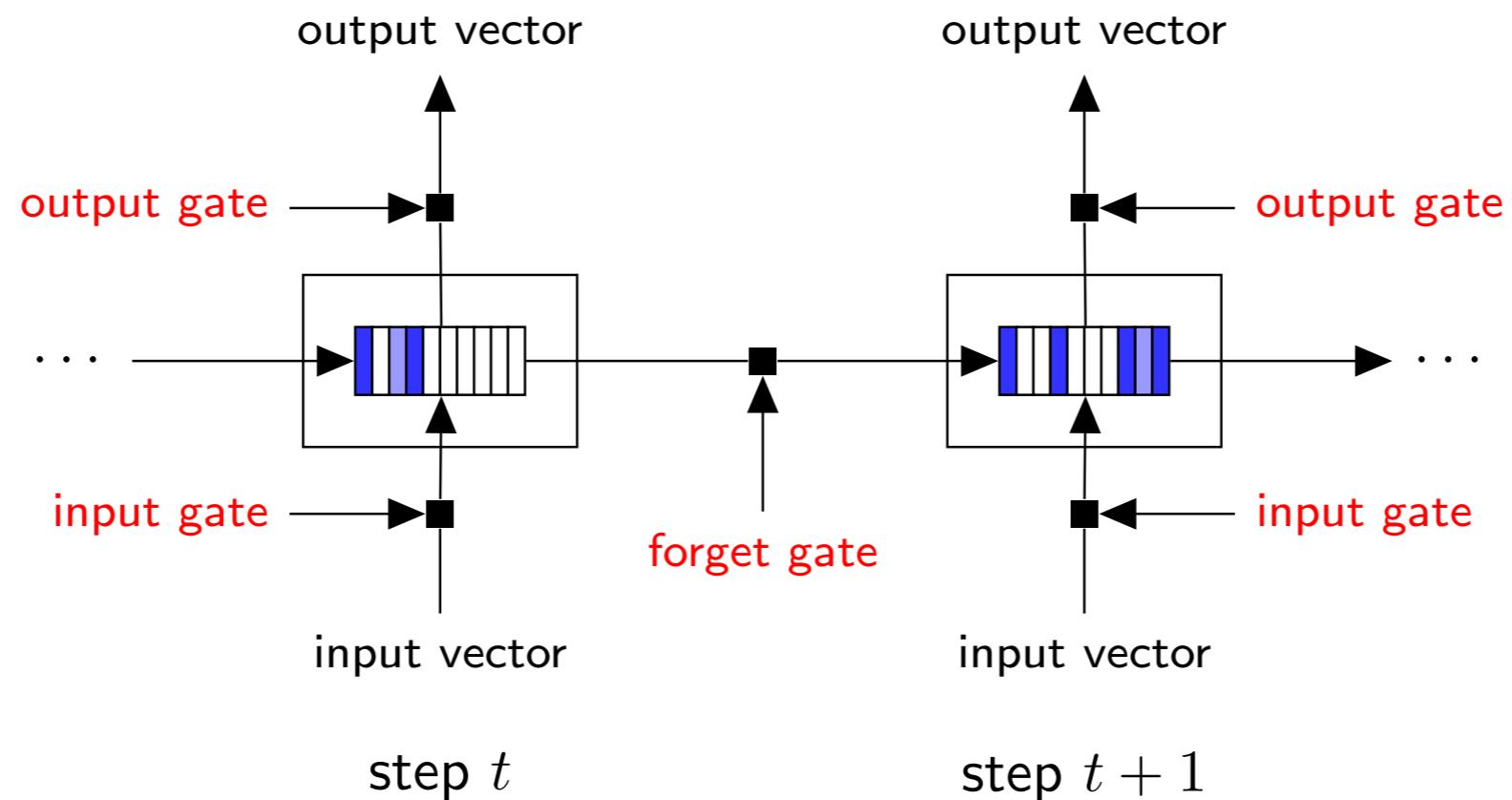
Sequential LSTM



Starting at state $t - 1$:

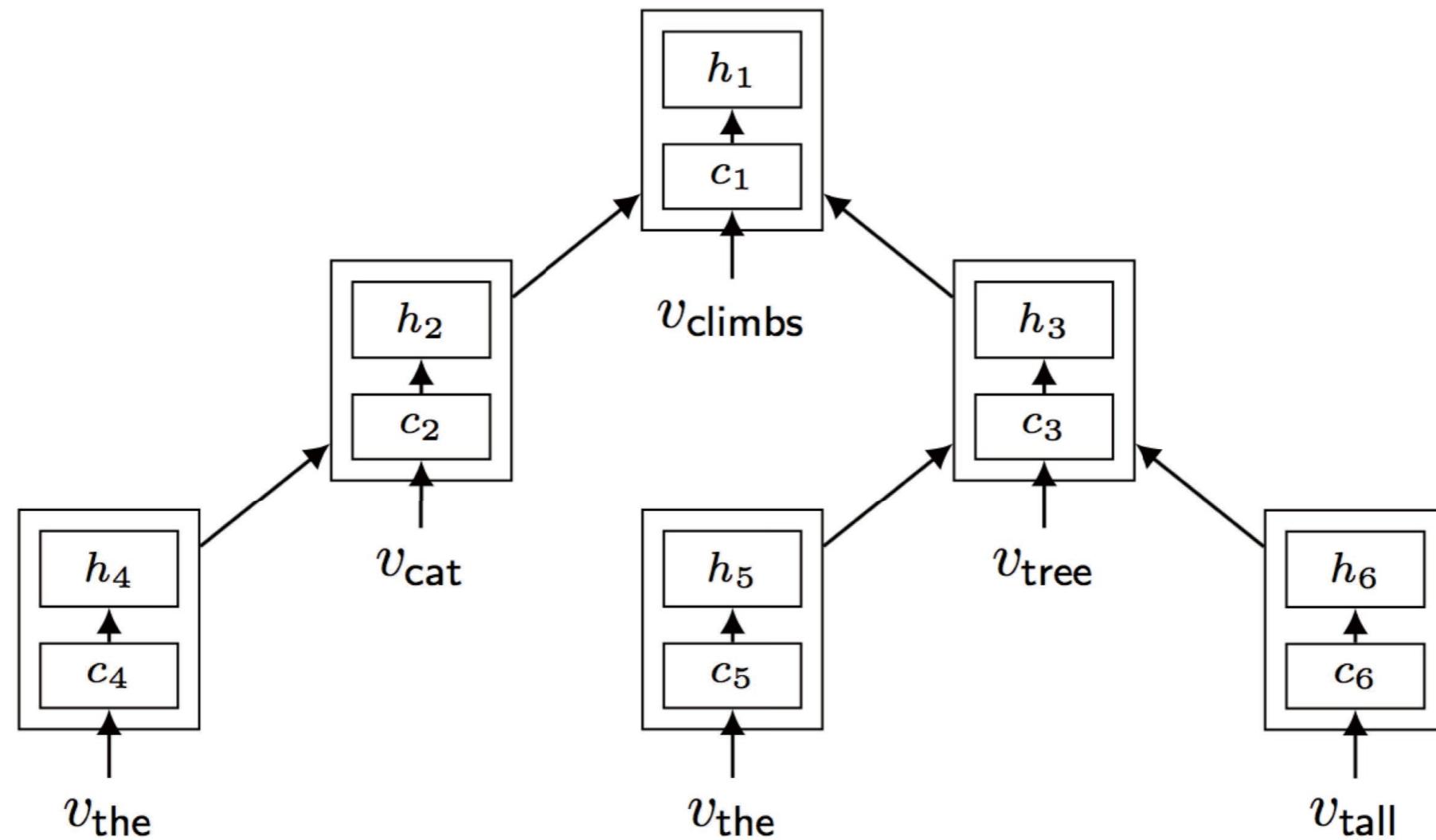
- ① Predict gates ($\mathbf{i}_t, \mathbf{f}_t, \mathbf{o}_t$) and intermediate state ($\bar{\mathbf{c}}_t$) from \mathbf{x}_t and \mathbf{h}_{t-1} .
- ② Update cell state: $\mathbf{c}_t = \mathbf{i}_t \otimes \bar{\mathbf{c}}_t + \mathbf{f}_t \otimes \mathbf{c}_{t-1}$.
- ③ Compute hidden state (output vector): $\mathbf{h}_t = \mathbf{o}_t \otimes \mathbf{c}_t$

Generalising LSTM

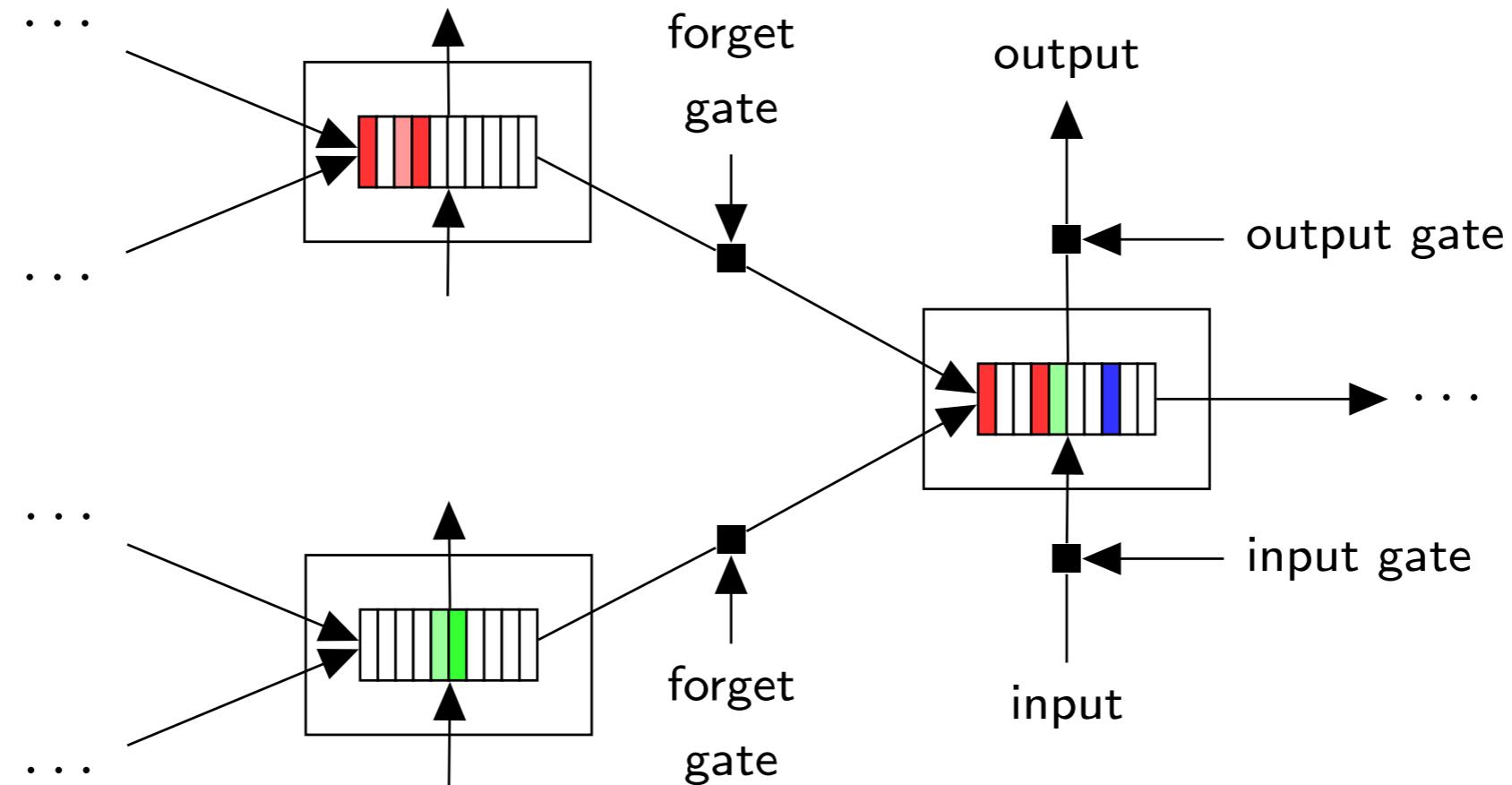


- ① Sequential LSTM receives only one child at a time.
- ② Generalize it to receive **multiple children**.

Tree LSTM

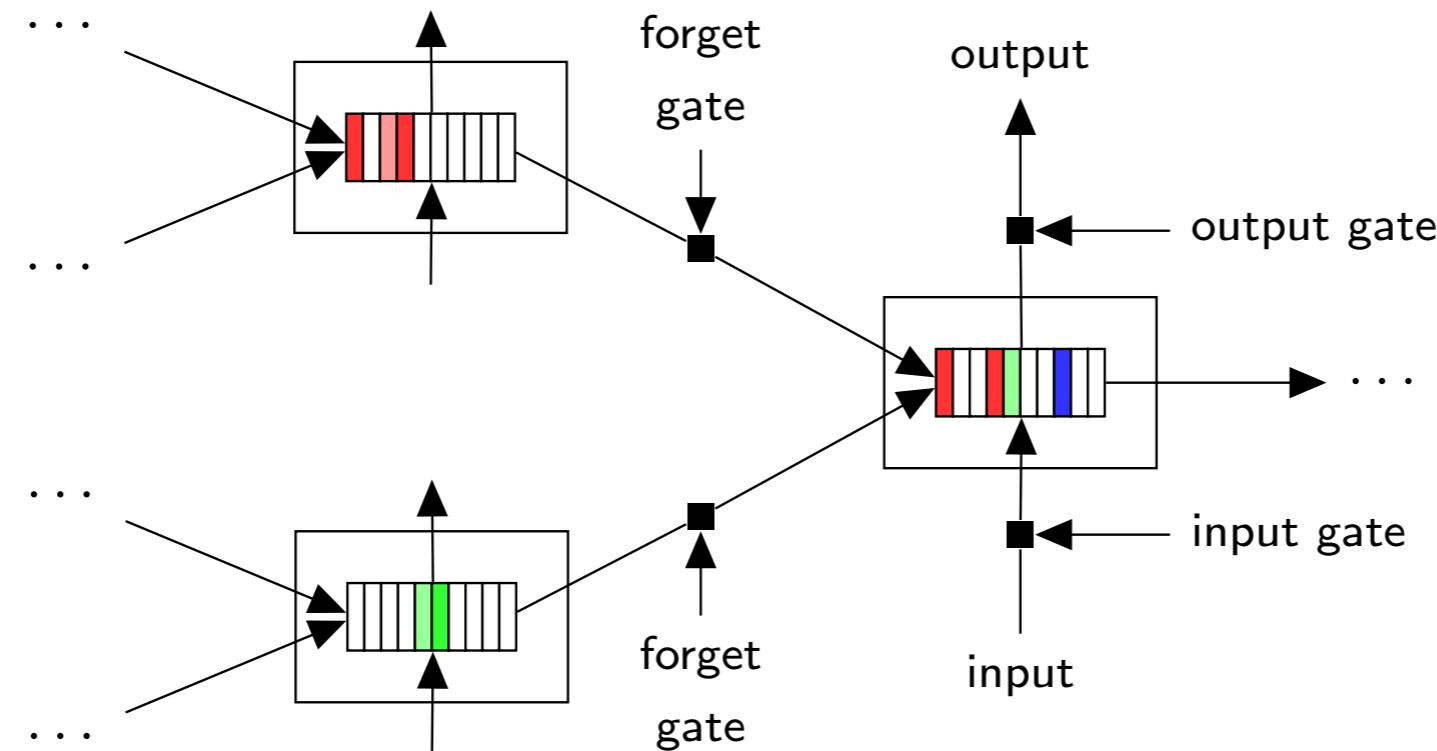


Tree LSTM



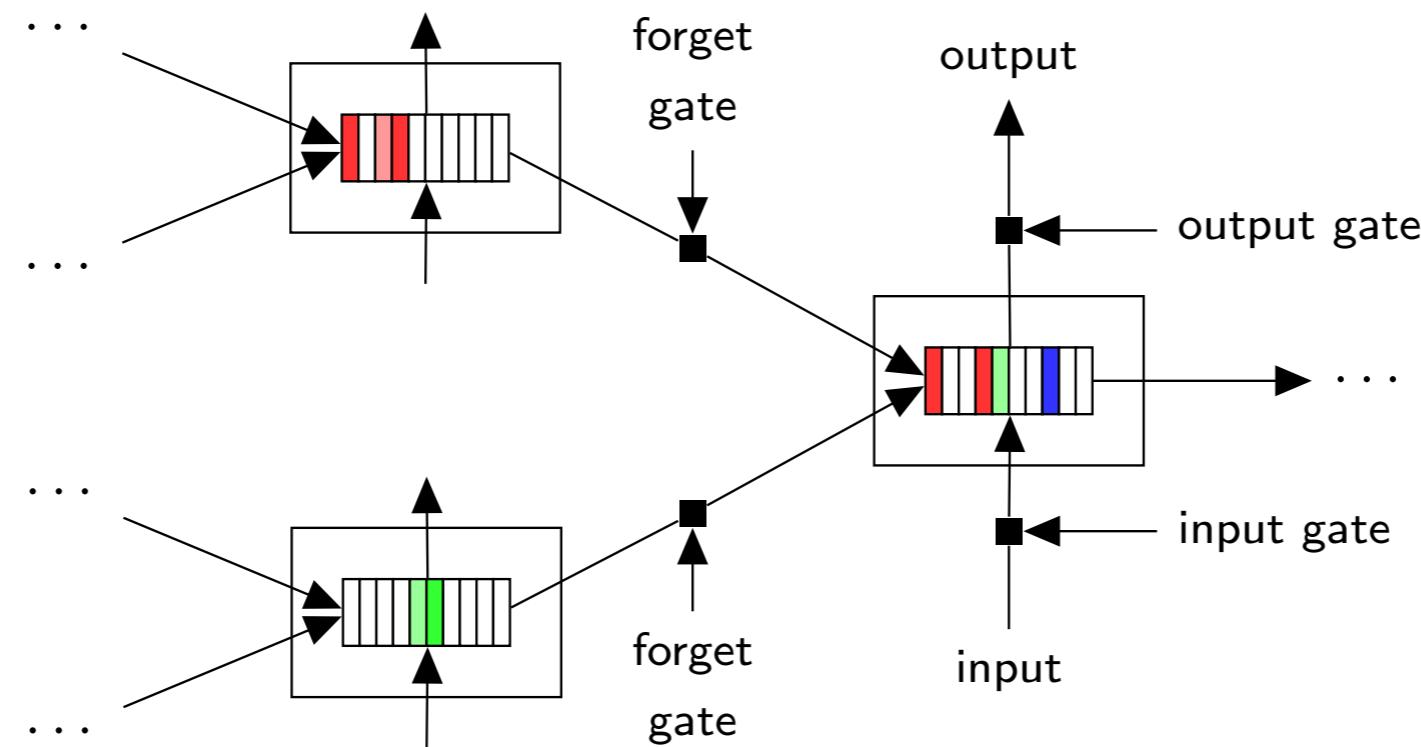
- ① Allows arbitrary **branching factor**.
- ② Each branch has a **separate forget gate**.

Tree-Sum LSTM



- ① **Sum up** all child vectors for computing gates and cell states.
- ② Forget gate parameters are **shared across branches**.
- ③ Supports trees with variable branch numbers.

N-ary Tree LSTM



- ① Each branch uses its own parameters for computing gates and cells.
- ② Supports trees with a **max branching factors**.

Evaluation (Sentiment Analysis)

Method	5-class	Binary
RNTN (Socher et al., 2013)	45.7	85.4
Paragraph-Vec (Le & Mikolov, 2014)	48.7	87.8
Convolutional NN (Kim 2014)	47.4	88.1
Epic (Hall et al., 2014)	49.6	—
DRNN (Irsoy & Cardie, 2014)	49.8	86.6
LSTM	46.4	84.9
Bidirectional LSTM	49.1	87.5
Constituency Tree-LSTM	51.0	88.0

Evaluation (Semantic Relatedness)

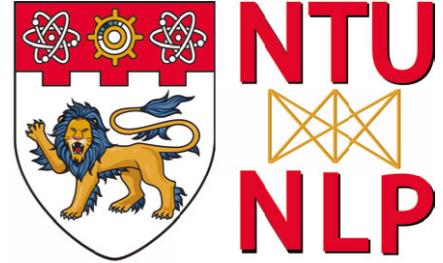
“the snowboarder is leaping over white snow” ? “a person who is practicing snowboarding jumps into the air”

- **Task:** Predict the semantic relatedness of sentence pairs
- **Dataset:** SICK from SemEval 2014, Task 1
- **Supervision:** Human-annotated relatedness scores $y \in [1, 5]$.
- **Model:** Sentence representation with Tree-LSTM on dependency parses. Similarity predicted by NN regressor given representations at root nodes

Evaluation (Semantic Relatedness)

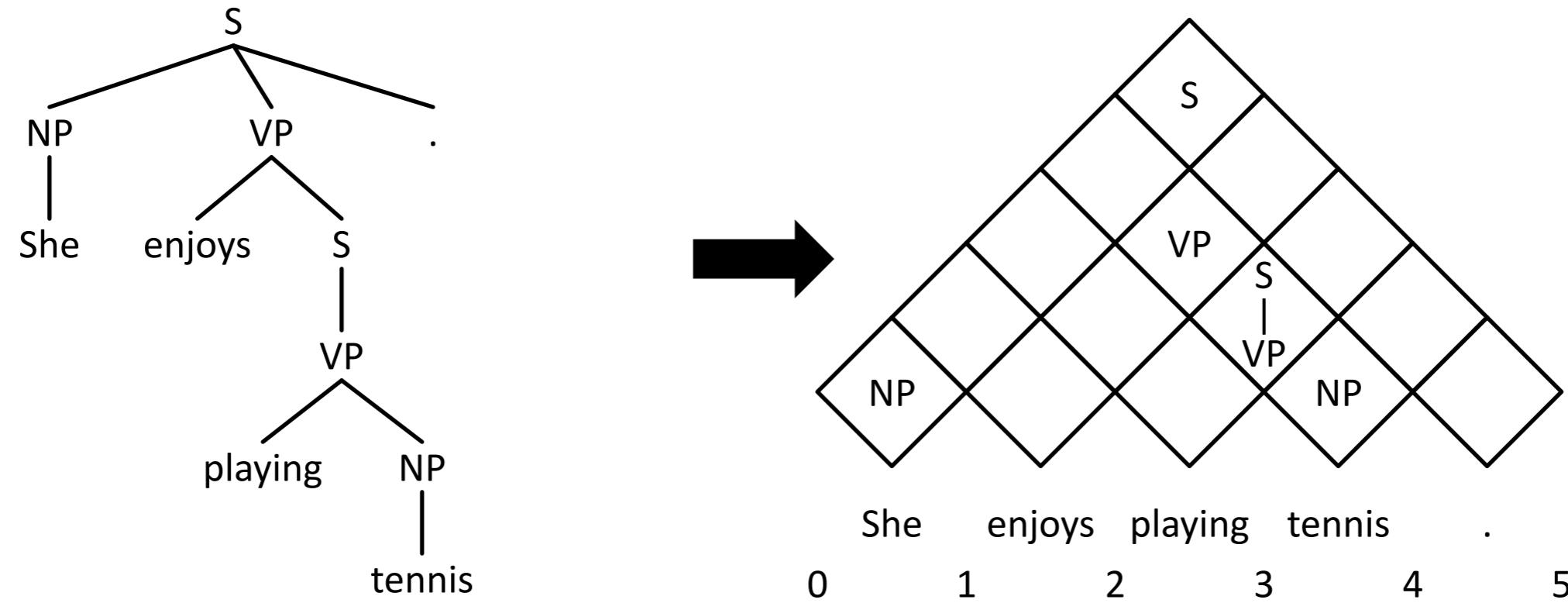
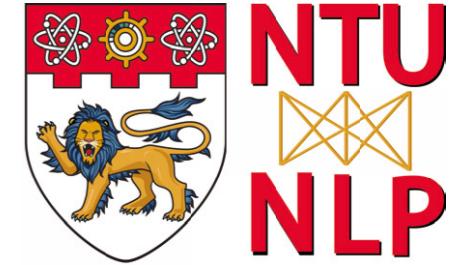
Method	Pearson's r
Word vector average	0.758
Meaning Factory (Bjerva et al., 2014)	0.827
ECNU (Zhao et al., 2014)	0.841
LSTM	0.853
Bidirectional LSTM	0.857
Dependency Tree-LSTM	0.868

Today's Plan



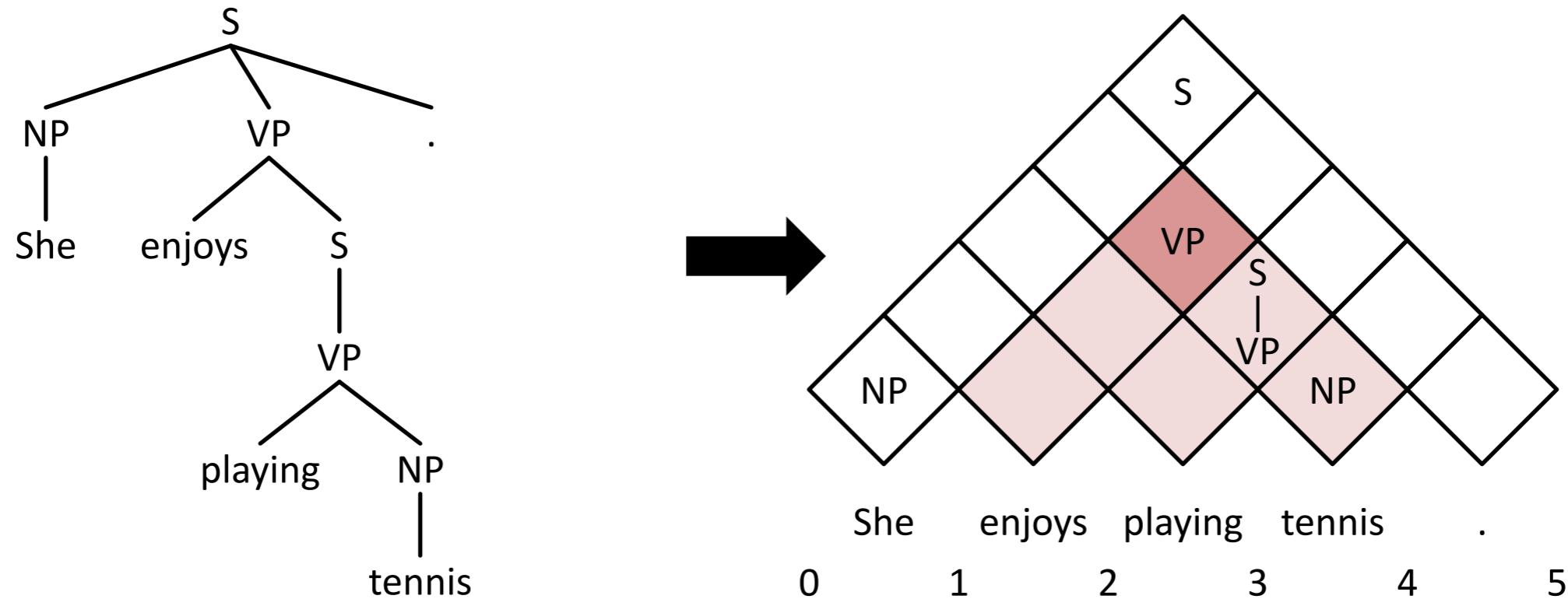
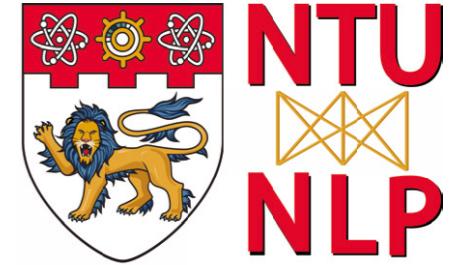
- Syntactic structures
 - Constituency & Dependency structures
- Need for syntactic structures
- Are languages hierarchical?
- Parsing with Recursive Neural Networks
- Backpropagation through tree
- Semantic composition with Tree RNNs
- Recent Parsers

Chart Parsing



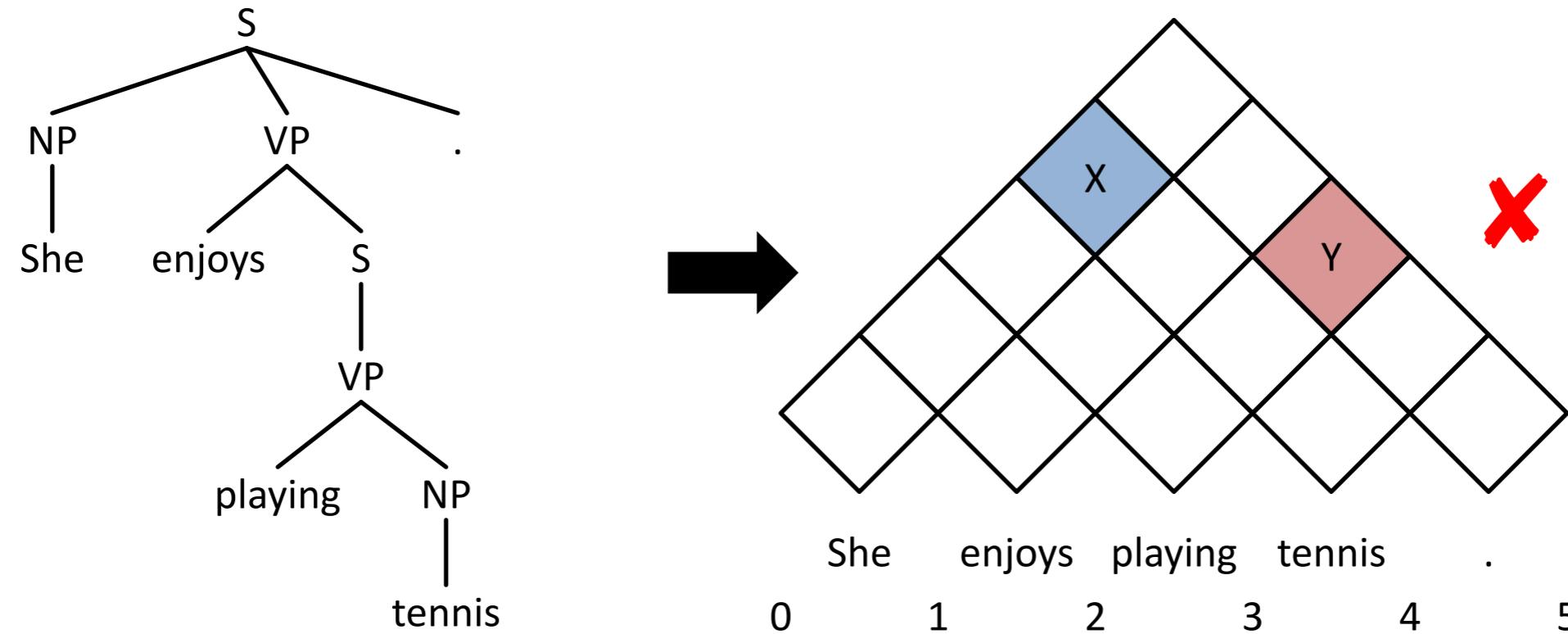
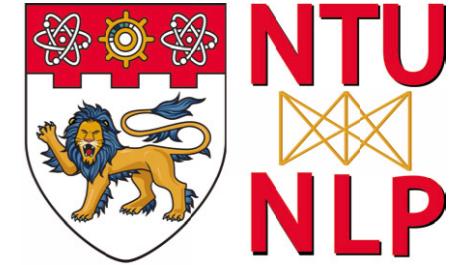
Constituency parses can be represented by a chart

Chart Parsing



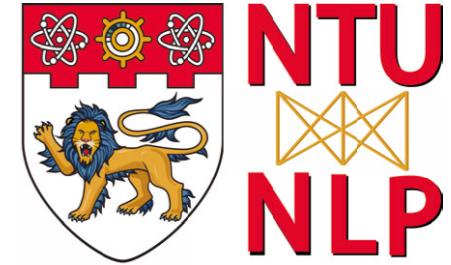
Constituency parses can be represented by a chart

Chart Parsing



Needs tree constraints

Use of CFG



$S(\text{decided}) \rightarrow NP(\text{workers}) VP(\text{decided})$

[Collins (1999)]

$S \rightarrow NP^S \quad VP^S$ [Klein and Manning (2003)]

$S \rightarrow NP \quad VP$ [Hall et al. (2014)]

S **Modern Parsers**

Modern Parsers don't use any CFGs; uses only labels

Scoring in Parsing

$\text{score}(S \rightarrow \text{NP}^S \text{ VP}^S)$

[Klein and Manning (2003)]

$\text{score}(i, k, j, S \rightarrow \text{NP VP})$

[Hall et al. (2014)]

$\text{score}(i, j, S)$ and $\text{score}_{\text{action}}(i, k, j)$

[Cross and Huang (2016)]

Shift-Reduce

$\text{score}(i, j, S)$

Modern Parsers

Decoding in Parsing

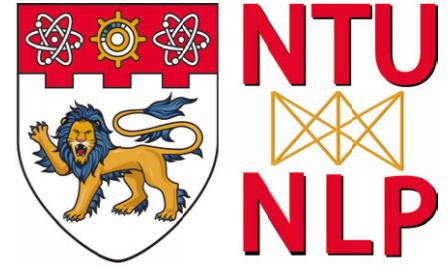


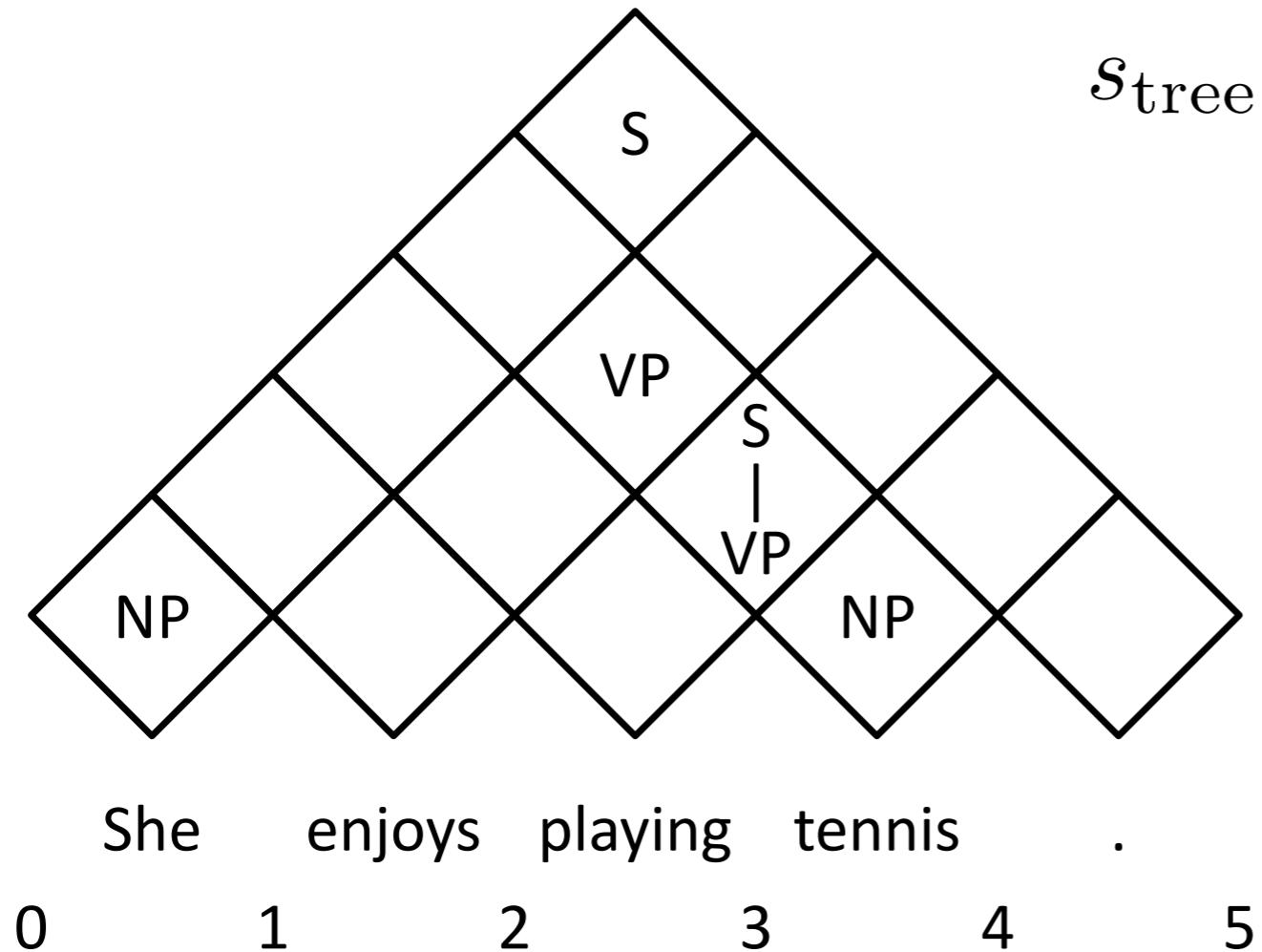
Chart-based

Globally optimal, $O(n^3)$ time complexity

Transition-based

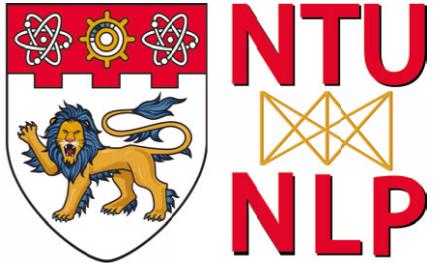
Greedy, $O(n)$ or $O(n^2)$ time complexity

Scoring in Modern Parsers



$$\begin{aligned}
 s_{\text{tree}}(T) &= \sum_{(\ell, (i, j)) \in T} s(i, j, \ell) \\
 &= s(0, 5, S) \\
 &\quad + s(0, 1, \text{NP}) \\
 &\quad + s(1, 4, \text{VP}) \\
 &\quad + s(2, 4, \text{S-VP}) \\
 &\quad + s(3, 4, \text{NP})
 \end{aligned}$$

Dynamic Programming (CKY)



Base case:

$$s_{\text{best}}(i, i + 1) = \max_{\ell} [s(i, i + 1, \ell)]$$

Pick best label

General case:

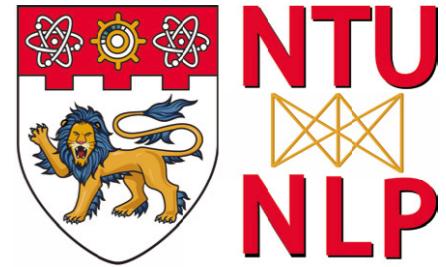
$$s_{\text{best}}(i, j) = \max_{\ell} [s(i, j, \ell)] + \max_k [s_{\text{best}}(i, k) + s_{\text{best}}(k, j)]$$

Pick best label

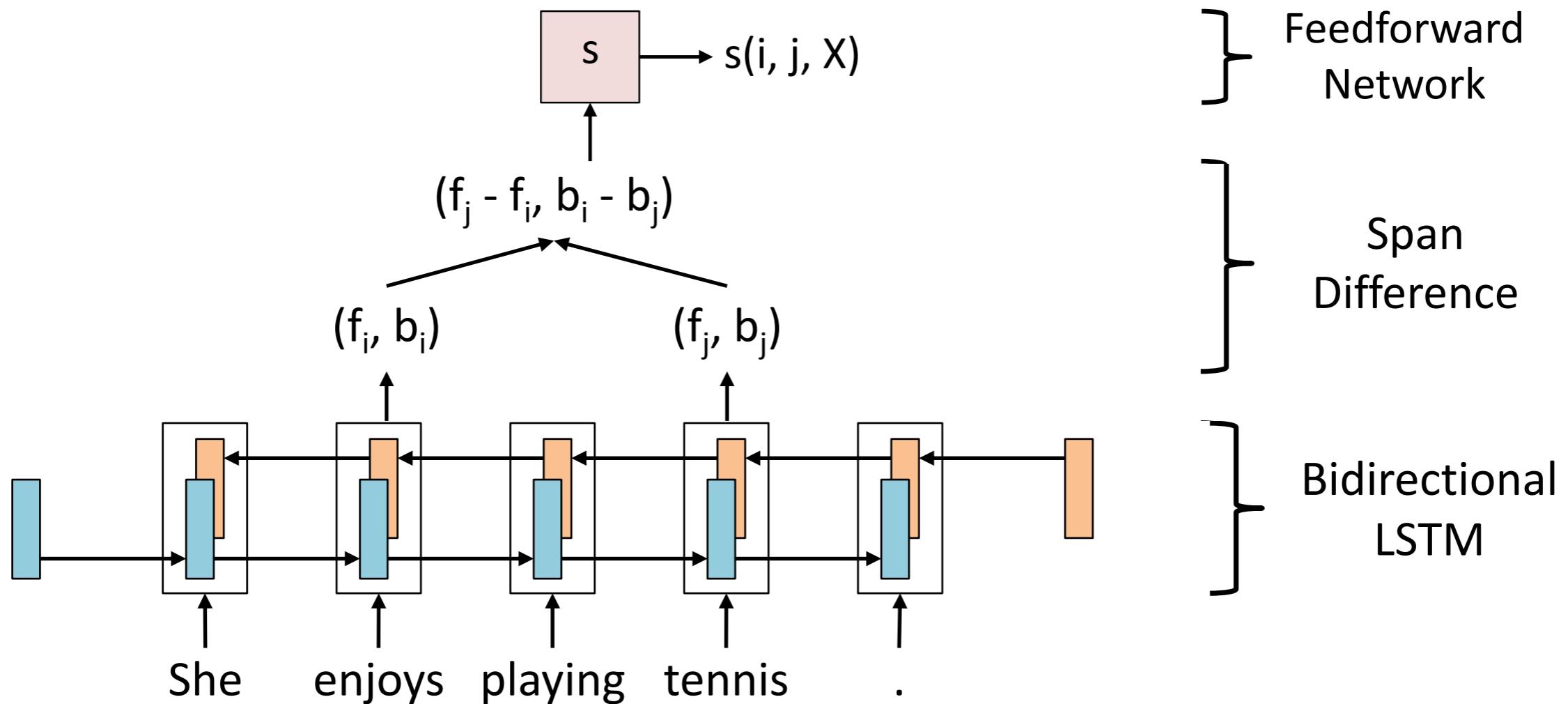
Bottom Up

Pick best split point

Scoring Functions

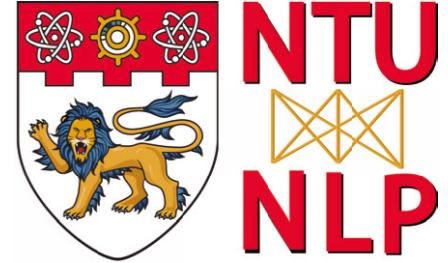


Using Sequential B-LSTM as the Sentence Encoder (Stern et al, 2017)



[Inspired by Cross and Huang (2016)]

Max-Margin Training



(Stern et al, 2017)

Want $s_{\text{tree}}(T^*) > s_{\text{tree}}(T)$ for all $T \neq T^*$

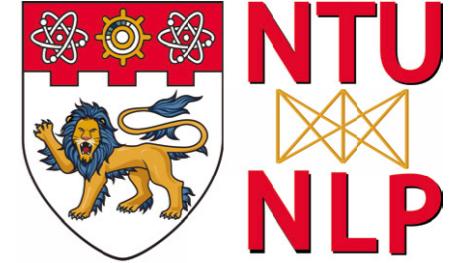
Require larger margin for higher loss:

$$s_{\text{tree}}(T^*) \geq \Delta(T, T^*) + s_{\text{tree}}(T)$$

Use hinge penalty function:

$$\max \left(0, \Delta(\hat{T}, T^*) - s_{\text{tree}}(T^*) + s_{\text{tree}}(\hat{T}) \right)$$

Max-Margin Training



(Stern et al, 2017)

Use loss-augmented decoding during training:

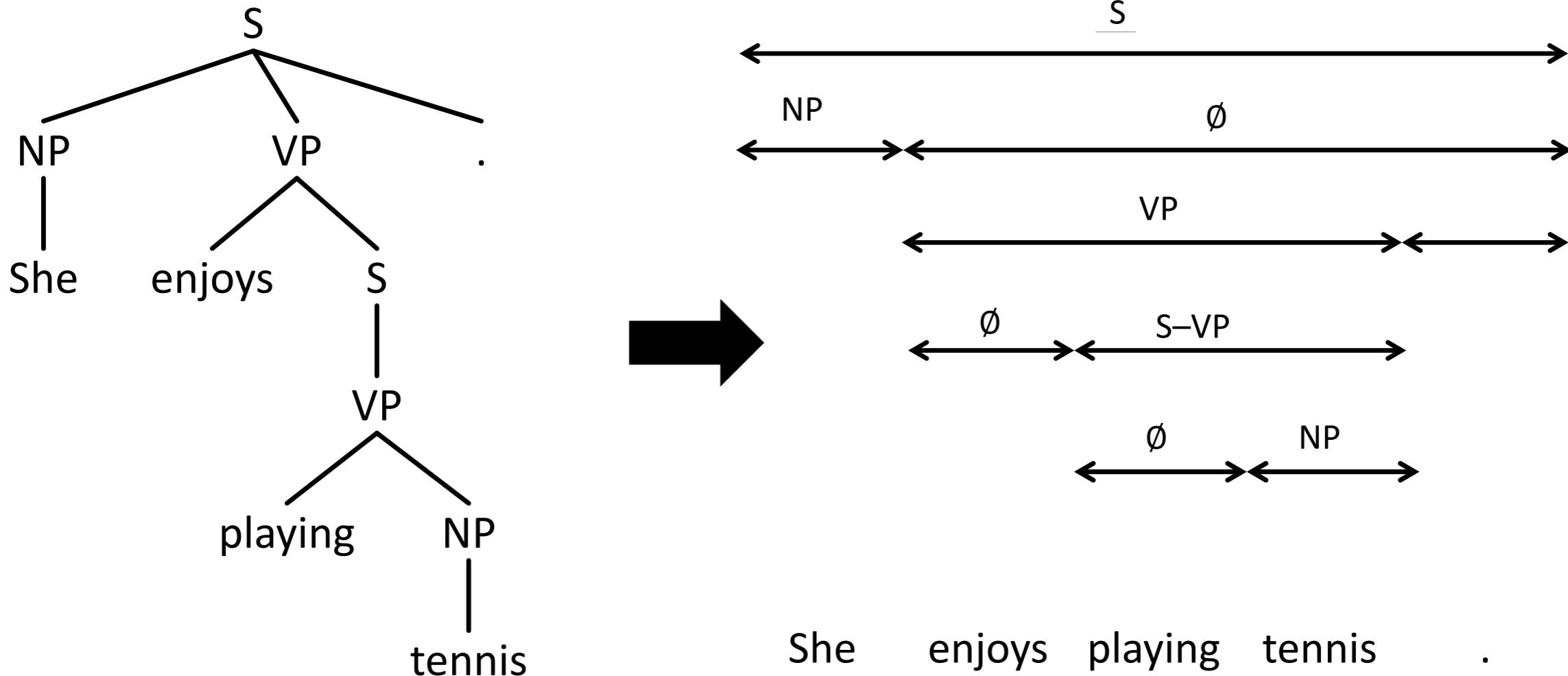
$$\hat{T} = \max_T [\Delta(T, T^*) + s_{\text{tree}}(T)]$$

Loss-augmented decoding for Hamming loss:

Replace $s(i, j, \ell)$ with $s(i, j, \ell) + 1(\ell \neq \ell_{ij}^*)$

Use Dynamic Programming (CKY) to solve it

Top-Down Greedy Parsing



Top-Down Greedy Parsing

$$\hat{\ell} = \operatorname{argmax}_{\ell} [s_{\text{label}}(i, j, \ell)]$$

$$\hat{k} = \operatorname{argmax}_k [s_{\text{span}}(i, k) + s_{\text{span}}(k, j)]$$

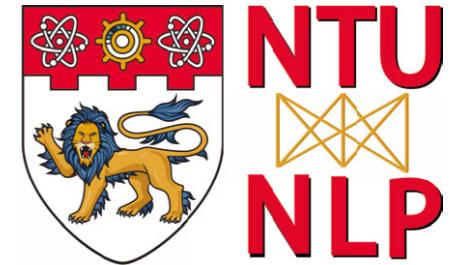
Training

Margin constraint for each decision:

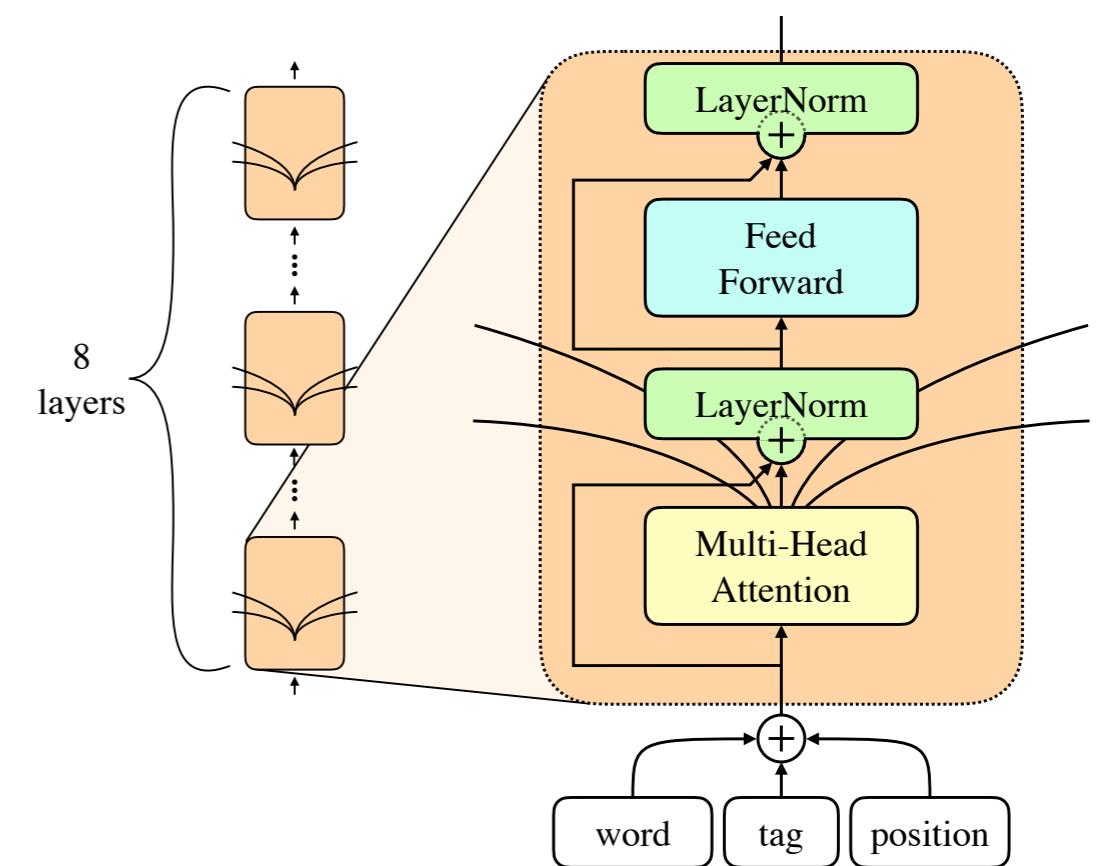
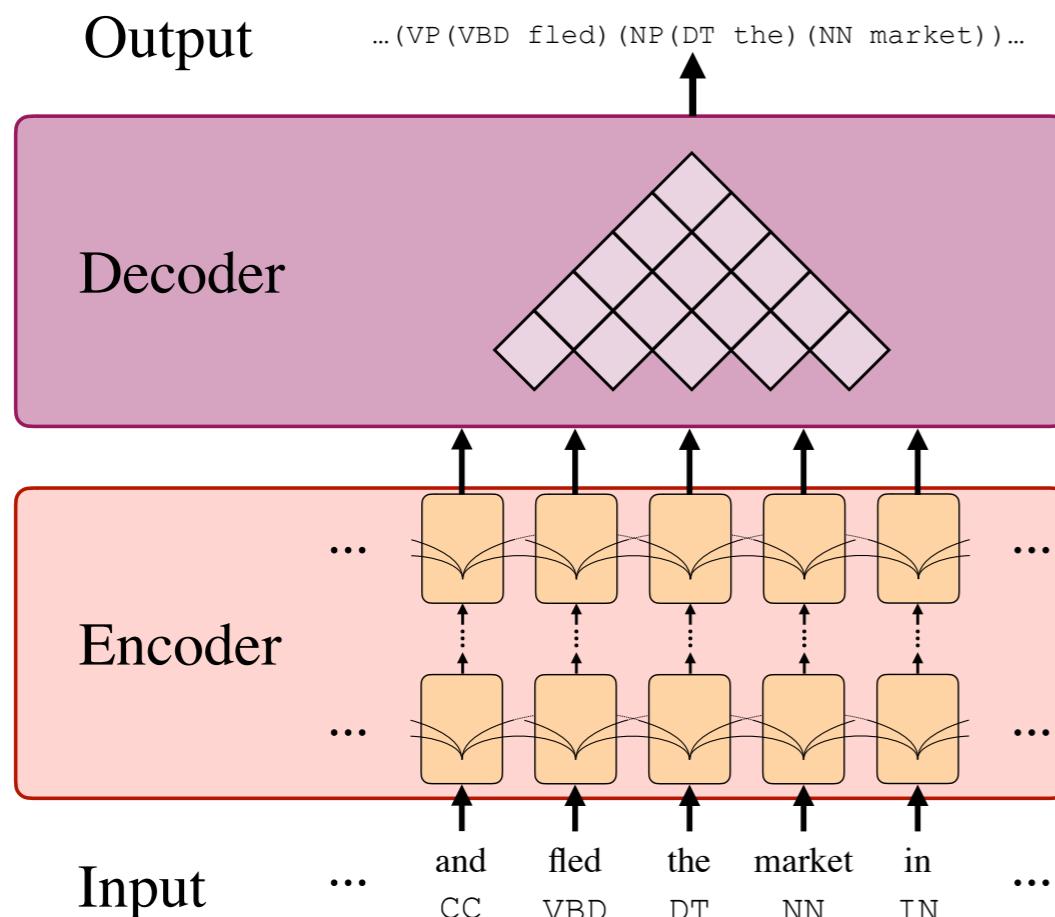
$$\text{score(gold)} \geq 1 + \text{score(other)}$$

- Similar to CKY, but doesn't consider the best sub-trees underneath
- Neural encoder is powerful enough.

Self-Attentive Encoder



Better Encoder

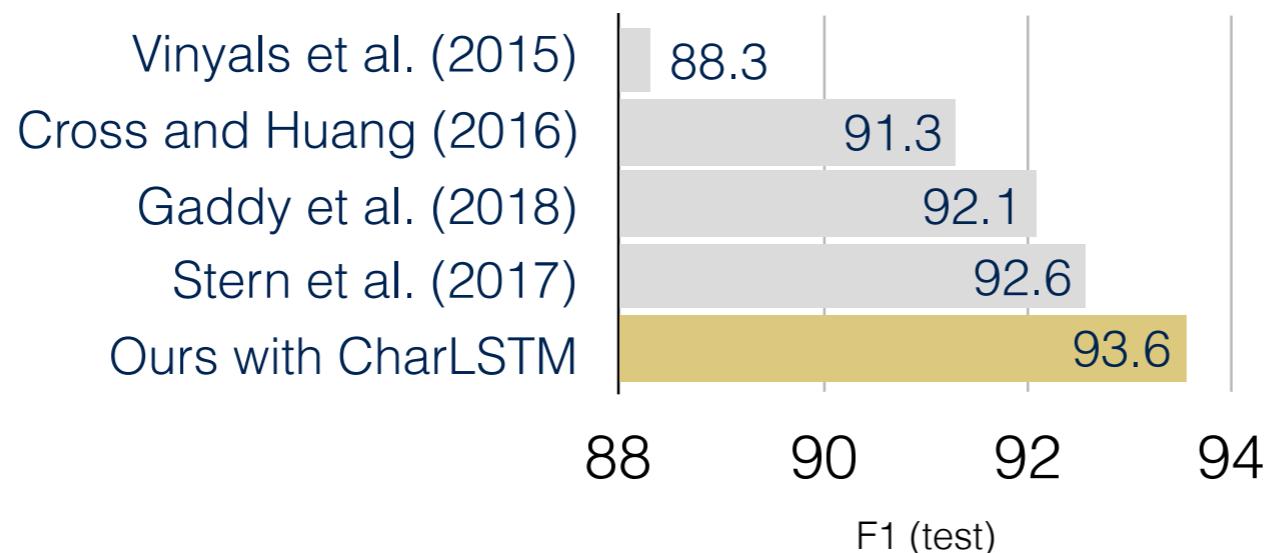


- Scoring and decoding are same as Stern et al, (2017)

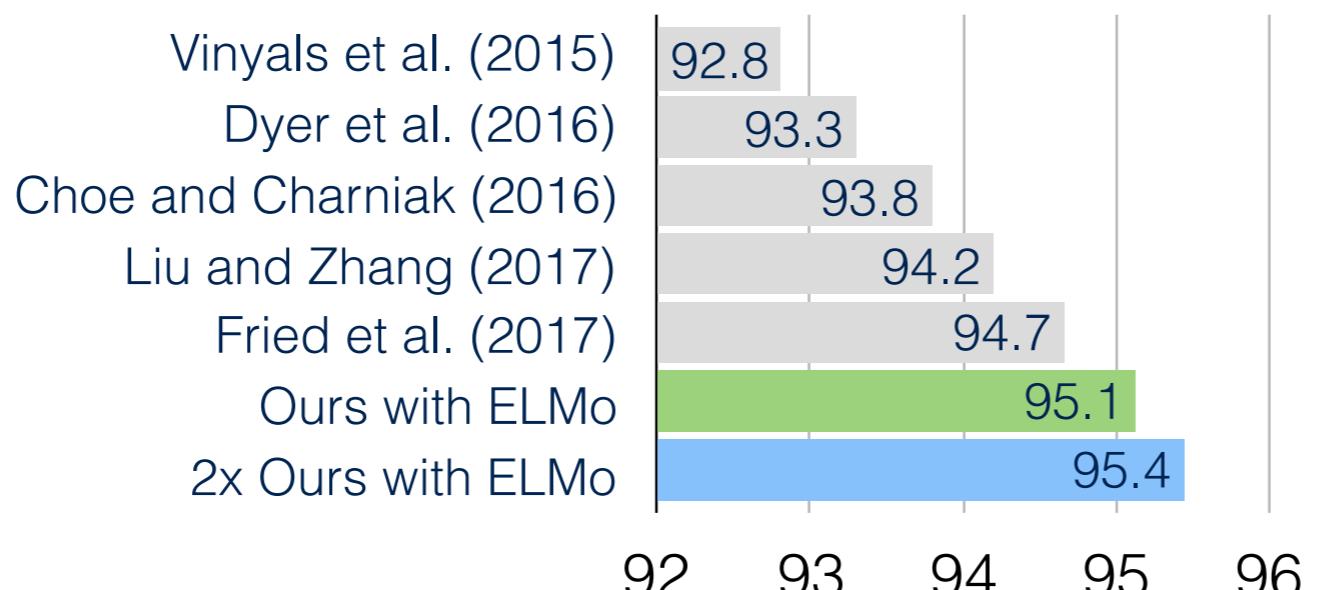
Results

English

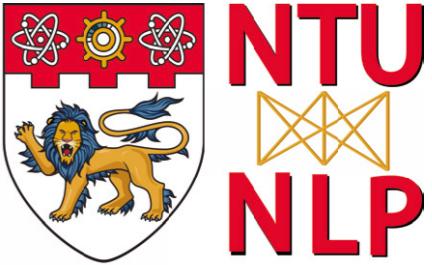
Single Model, WSJ Only



Multi-Model / External

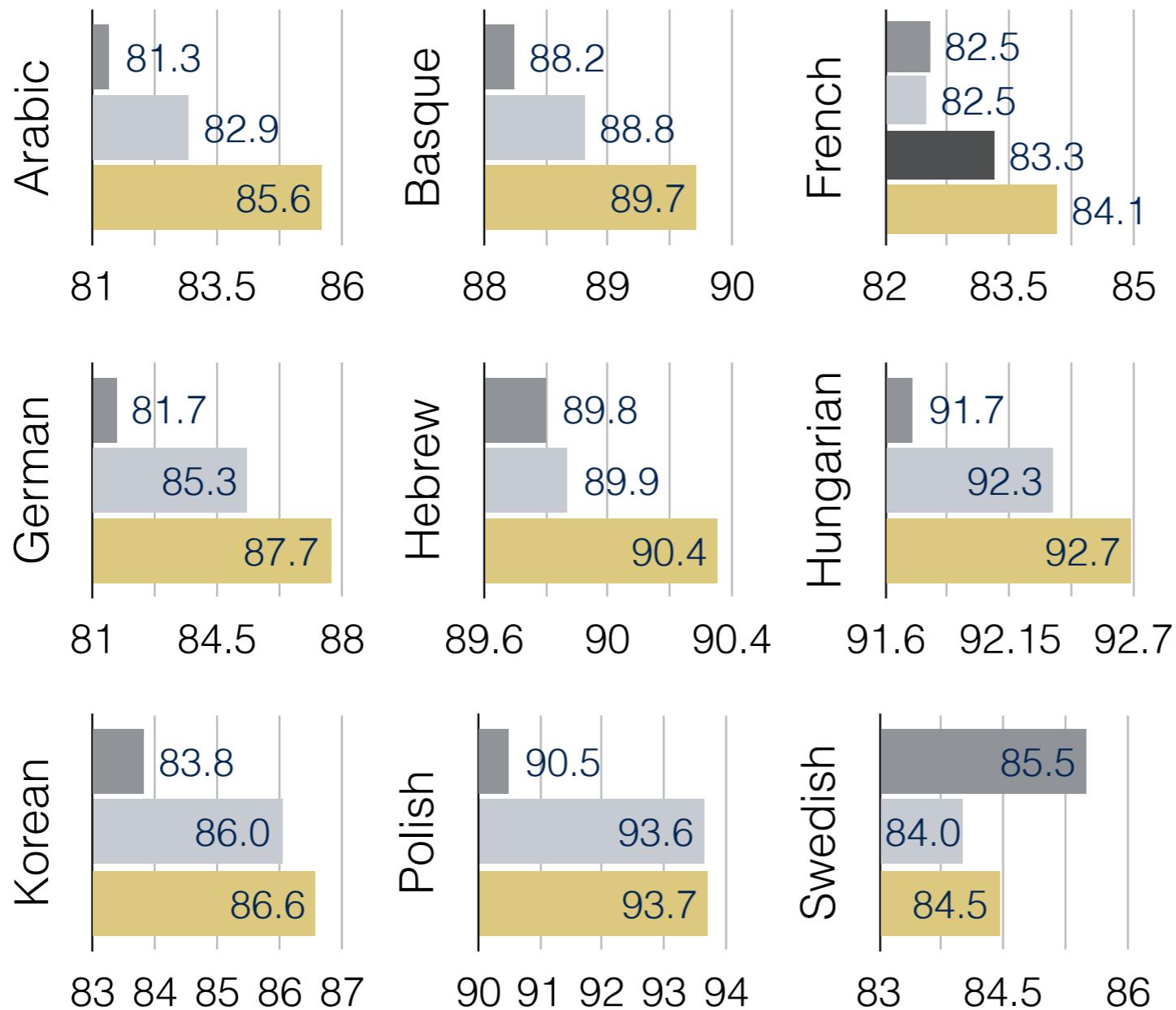


Results

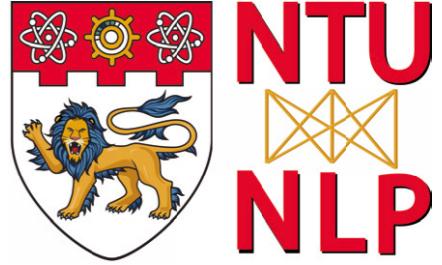


Multilingual (SPMRL)

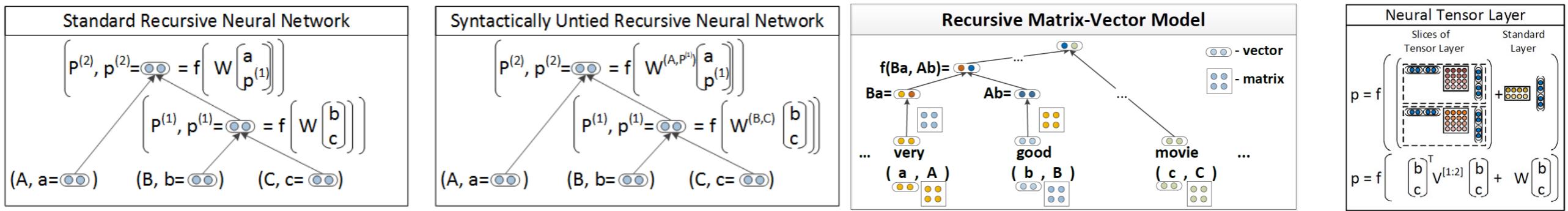
■ Björkelund et al. (2014) ■ Coavoux and Crabbé (2017)
■ Cross and Huang (2016) ■ Ours with CharLSTM



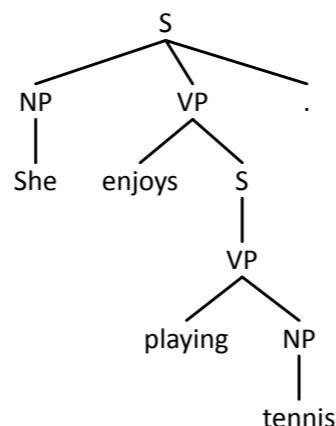
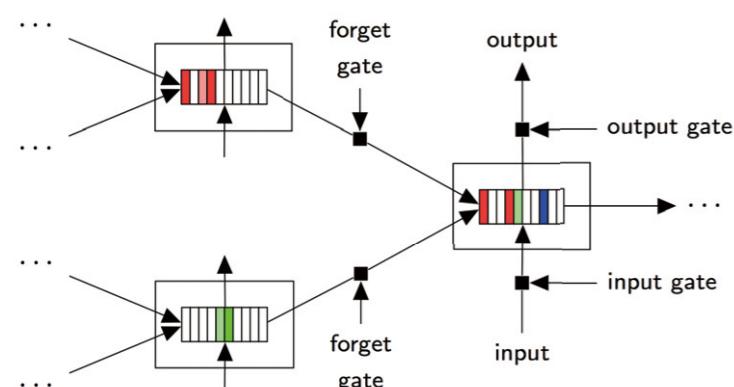
Summary



Tree-RNNs



Tree-LSTM



Modern Parsers

