

# Part-of-Speech Tagging: Neural Methods

Katharina Kann — CSCI/LING5832

# Word Classes: Parts of Speech

- Traditional parts of speech:
  - Noun, verb, adjective, preposition, adverb, article, interjection, pronoun, conjunction, etc.

# Word Classes: Parts of Speech

- Traditional parts of speech:
  - Noun, verb, adjective, preposition, adverb, article, interjection, pronoun, conjunction, etc.
  - Lots of names for this notion: Part of speech, lexical category, word class, morphological class, lexical tag...

# Word Classes: Parts of Speech

- Traditional parts of speech:
  - Noun, verb, adjective, preposition, adverb, article, interjection, pronoun, conjunction, etc.
  - Lots of names for this notion: Part of speech, lexical category, word class, morphological class, lexical tag...
  - Lots of debate within linguistics about the number, nature, and universality of these categories

# Word Classes: Parts of Speech

- Sources of evidence
  - Morphological evidence
    - walk, walking, walked, walks
    - probably a verb

# Word Classes: Parts of Speech

- Sources of evidence
  - Morphological evidence
    - walk, walking, walked, walks
    - probably a verb
  - Distributional evidence
    - The crash, A crash, Two crashes,  
The big crash...
    - probably a noun

# Penn Treebank POS Tagset

Tag	Description	Example	Tag	Description	Example
CC	coordin. conjunction	<i>and, but, or</i>	SYM	symbol	<i>+, %, &amp;</i>
CD	cardinal number	<i>one, two, three</i>	TO	“to”	<i>to</i>
DT	determiner	<i>a, the</i>	UH	interjection	<i>ah, oops</i>
EX	existential ‘there’	<i>there</i>	VB	verb, base form	<i>eat</i>
FW	foreign word	<i>mea culpa</i>	VBD	verb, past tense	<i>ate</i>
IN	preposition/sub-conj	<i>of, in, by</i>	VBG	verb, gerund	<i>eating</i>
JJ	adjective	<i>yellow</i>	VBN	verb, past participle	<i>eaten</i>
JJR	adj., comparative	<i>bigger</i>	VBP	verb, non-3sg pres	<i>eat</i>
JJS	adj., superlative	<i>wildest</i>	VBZ	verb, 3sg pres	<i>eats</i>
LS	list item marker	<i>1, 2, One</i>	WDT	wh-determiner	<i>which, that</i>
MD	modal	<i>can, should</i>	WP	wh-pronoun	<i>what, who</i>
NN	noun, sing. or mass	<i>llama</i>	WP\$	possessive wh-	<i>whose</i>
NNS	noun, plural	<i>llamas</i>	WRB	wh-adverb	<i>how, where</i>
NNP	proper noun, singular	<i>IBM</i>	\$	dollar sign	<i>\$</i>
NNPS	proper noun, plural	<i>Carolinas</i>	#	pound sign	<i>#</i>
PDT	predeterminer	<i>all, both</i>	“	left quote	<i>‘ or “</i>
POS	possessive ending	<i>’s</i>	”	right quote	<i>’ or ”</i>
PRP	personal pronoun	<i>I, you, he</i>	(	left parenthesis	<i>[, (, {, &lt;</i>
PRP\$	possessive pronoun	<i>your, one’s</i>	)	right parenthesis	<i>], ), }, &gt;</i>
RB	adverb	<i>quickly, never</i>	,	comma	<i>,</i>
RBR	adverb, comparative	<i>faster</i>	.	sentence-final punc	<i>. ! ?</i>
RBS	adverb, superlative	<i>fastest</i>	:	mid-sentence punc	<i>: ; ... --</i>
RP	particle	<i>up, off</i>			

# POS Tagging

- The process of assigning a part of speech or lexical class marker to each word in a text.
- Often a useful first step in an NLP pipeline.



# POS Tagging

- The process of assigning a part of speech or lexical class marker to each word in a text.
- Often a useful first step in an NLP pipeline.
- Fast and accurate taggers are widely available for many languages.
  - (Now we will learn how to make one!)

# POS Tagging

- The process of assigning a part of speech or lexical class marker to each word in a text.
- This is our first example of a **sequence labeling** task:
  - Assigning a category label to each element of a sequence.

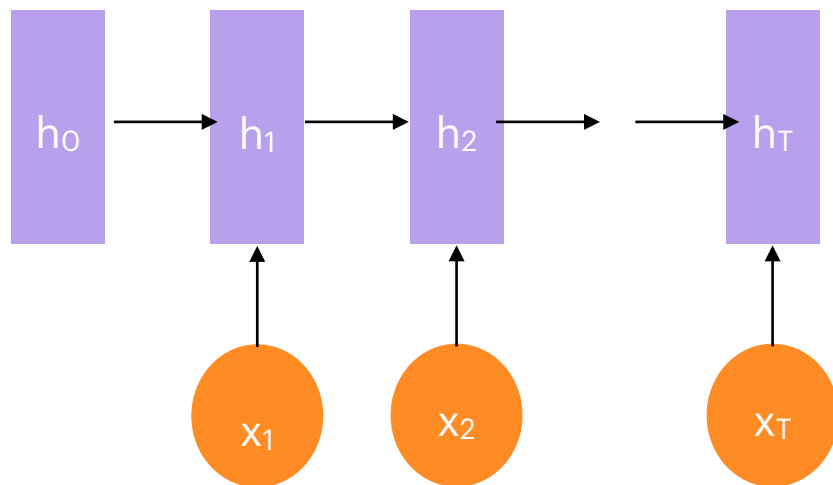
# Recap from Lecture 07: More on RNNs

# Recurrent Neural Networks

Input:  $x_1, x_2, \dots, x_T$

$$h_i = \text{sigmoid}(W_h h_{i-1} + W_x x_i)$$

$$y = \text{argmax}(\text{softmax}(W_y h_T))$$

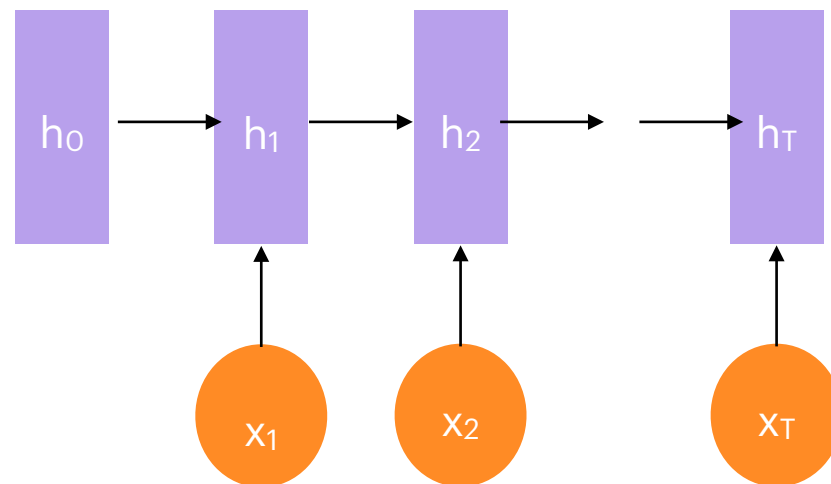


# Recurrent Neural Networks

Input:  $x_1, x_2, \dots, x_T$

$$h_i = \text{sigmoid}(W_h h_{i-1} + W_x x_i)$$

$$y = \text{argmax}(\text{softmax}(W_y h_T))$$



# Problems with Recurrent Neural Networks

- During backpropagation, since RNNs multiply the same matrix over and over, gradients vanish or explode.

# Problems with Recurrent Neural Networks

- During backpropagation, since RNNs multiply the same matrix over and over, gradients vanish or explode.
- Cannot capture long-term dependencies (or learn much) with uninformative gradients!

# Problems with Recurrent Neural Networks

- During backpropagation, since RNNs multiply the same matrix over and over, gradients vanish or explode.
- Cannot capture long-term dependencies (or learn much) with uninformative gradients!
- History gets forgotten over time.



# Problems with Recurrent Neural Networks

The flights the airline was cancelling were full.

# Problems with Recurrent Neural Networks

The flights the  airline  was cancelling were full.

# Problems with Recurrent Neural Networks

The flights the airline was cancelling were full.

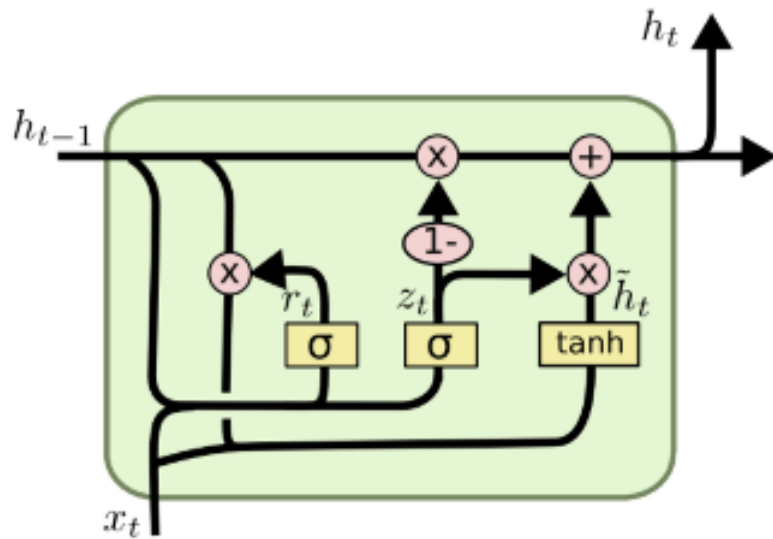
# Solutions

- **Gradient Clipping:** Normalize the gradient of a parameter vector when its L2 norm reaches a certain value.

# Solutions

- **Gradient Clipping**: Normalize the gradient of a parameter vector when its L2 norm reaches a certain value.
- Explicit **memory cell** with which a model can keep important long-term information.

# Gated Recurrent Units (GRUs)



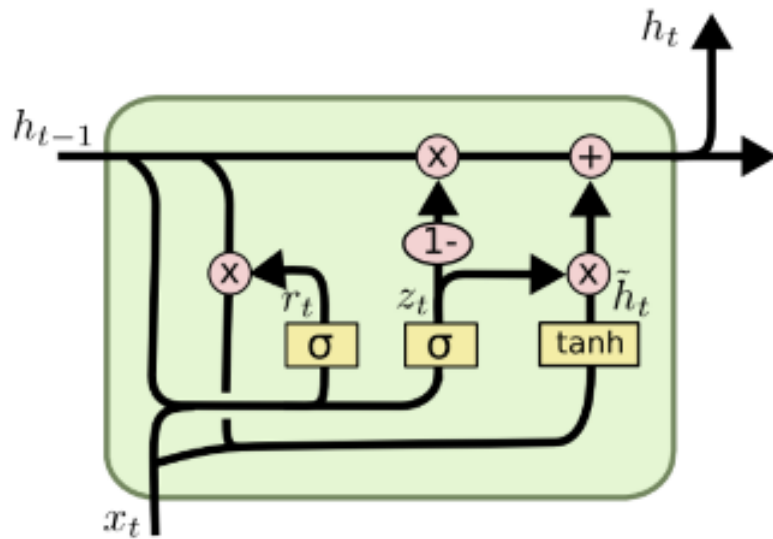
$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

# Gated Recurrent Units (GRUs)



$$z_t = \sigma (W_z \cdot [h_{t-1}, x_t])$$

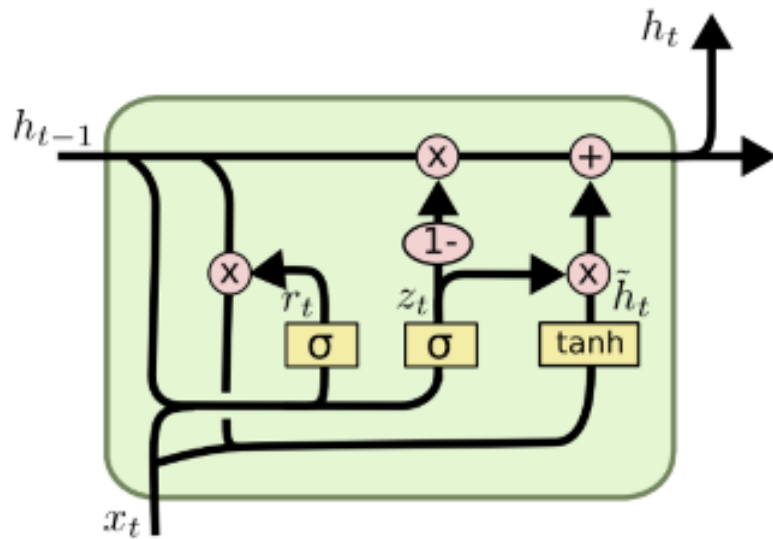
$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh (W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

**Update gate:** controls how much of the previous hidden state to update

# Gated Recurrent Units (GRUs)



$$z_t = \sigma (W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh (W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

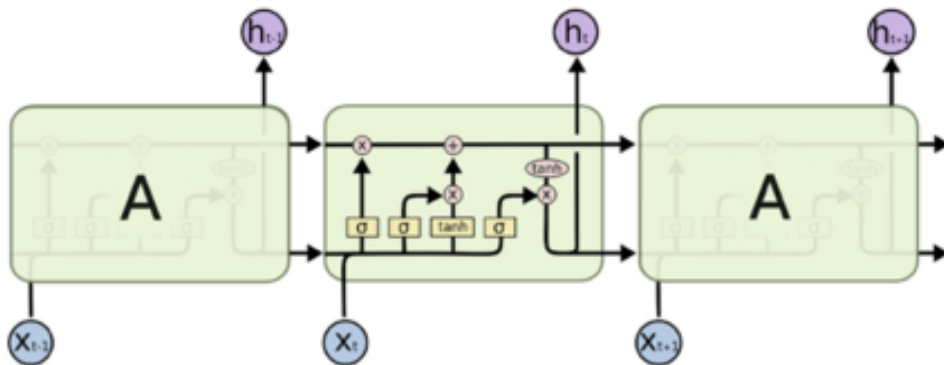
**Reset gate:** controls how much of the previous hidden state to keep



# Long Short-Term Memory Networks (LSTMs)

- First proposed by Hochreiter and Schmidhuber (1997).
- Similar to the GRU, but some more parameters.
- Performance is usually extremely similar.

# Long Short-Term Memory Networks (LSTMs)



The repeating module in an LSTM contains four interacting layers.

$$g_t = \tanh(U_g h_{t-1} + W_g x_t)$$

$$i_t = \sigma(U_i h_{t-1} + W_i x_t)$$

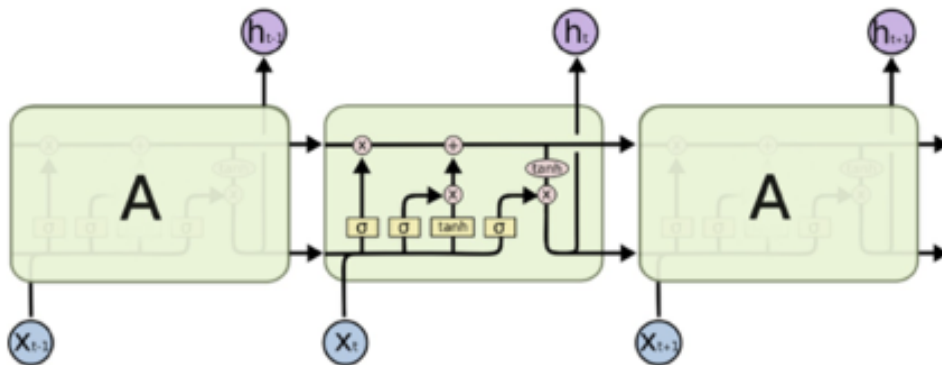
$$f_t = \sigma(U_f h_{t-1} + W_f x_t)$$

$$o_t = \sigma(U_o h_{t-1} + W_o x_t)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$

$$h_t = o_t \odot \tanh(c_t)$$

# Long Short-Term Memory Networks (LSTMs)

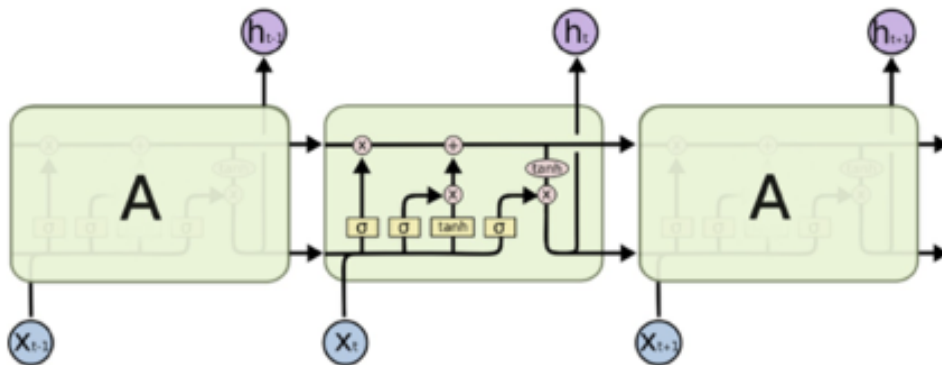


The repeating module in an LSTM contains four interacting layers.

$$\begin{aligned}
 g_t &= \tanh(U_g h_{t-1} + W_g x_t) \\
 i_t &= \sigma(U_i h_{t-1} + W_i x_t) \\
 f_t &= \sigma(U_f h_{t-1} + W_f x_t) \\
 o_t &= \sigma(U_o h_{t-1} + W_o x_t) \\
 c_t &= f_t \odot c_{t-1} + i_t \odot g_t \\
 h_t &= o_t \odot \tanh(c_t)
 \end{aligned}$$

**Add gate**

# Long Short-Term Memory Networks (LSTMs)



The repeating module in an LSTM contains four interacting layers.

$$g_t = \tanh(U_g h_{t-1} + W_g x_t)$$

$$i_t = \sigma(U_i h_{t-1} + W_i x_t)$$

$$f_t = \sigma(U_f h_{t-1} + W_f x_t)$$

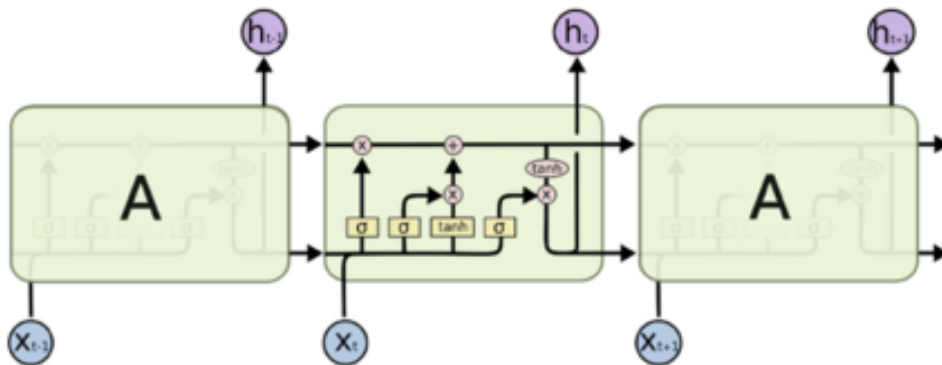
$$o_t = \sigma(U_o h_{t-1} + W_o x_t)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$

$$h_t = o_t \odot \tanh(c_t)$$

**Forget gate**

# Long Short-Term Memory Networks (LSTMs)



The repeating module in an LSTM contains four interacting layers.

$$g_t = \tanh(U_g h_{t-1} + W_g x_t)$$

$$i_t = \sigma(U_i h_{t-1} + W_i x_t)$$

$$f_t = \sigma(U_f h_{t-1} + W_f x_t)$$

$$o_t = \sigma(U_o h_{t-1} + W_o x_t)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$

$$h_t = o_t \odot \tanh(c_t)$$

**Output gate**

# How Do We Get a Sentence Representation?

- Apply the RNN to all words in a sequence.
- Use the last hidden vector as the representation of your sentence.
- The sentence representation is often called  $h$ .

# Bidirectional Recurrent Neural Networks

- Most prominent one: BiLSTM
- Idea: Read the input from both sides.
  - Beginning isn't forgotten!
- For the final representation:
  - Concatenate the last hidden vectors of the two individual RNNs.
  - (What is the size of the final sentence representation computed by a bidirectional RNN?)

# In-Class Exercise (Exercise 3 in Lecture 07)

- Assume that we want to classify the sentiment of tweets into *positive*, *neutral*, and *negative*.
- Design a neural network that can learn this task, given a large-enough training set.
- Which preprocessing steps do we need?
- Write down all dimensions explicitly. What is your output dimension?
- (Keep your architecture safe for later!)



# In-Class Exercise (Exercise 3 in Lecture 07)

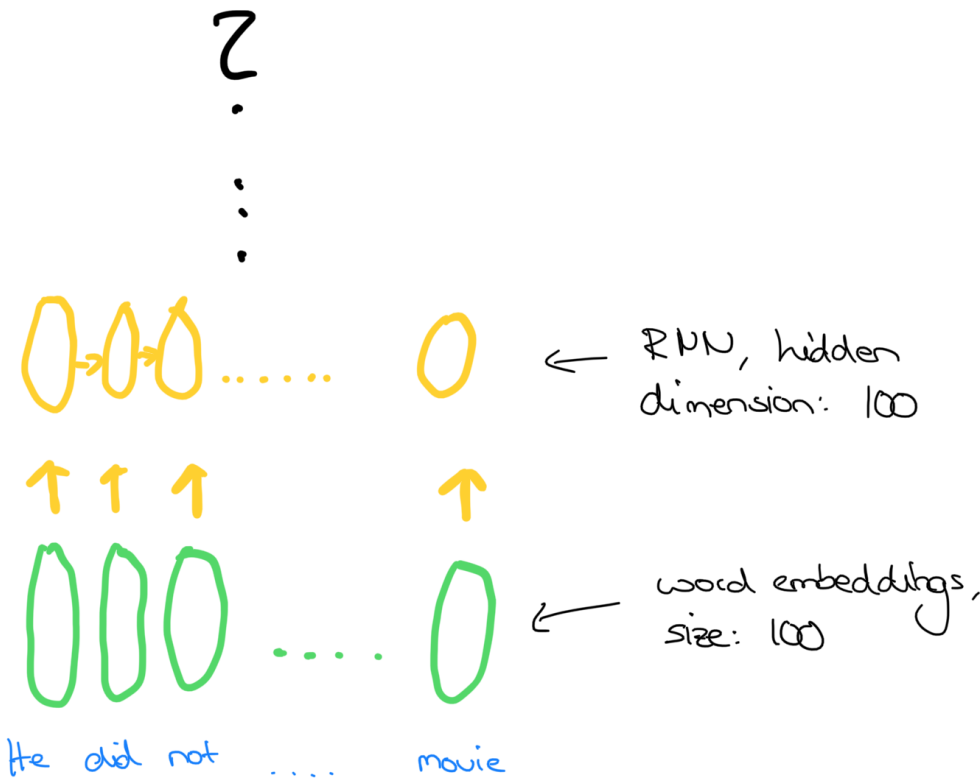
- 

Diagram illustrating a Recurrent Neural Network (RNN) architecture for sentiment classification. The input layer consists of word embeddings (green ovals) for the sentence "He did not ... movie". The hidden layer consists of RNN states (yellow ovals) connected sequentially. The output layer is a single RNN state (yellow oval) representing the sentiment classification. Labels include "word embeddings, size: 100" and "RNN, hidden dimension: 100".
  - classify the sentiment as *positive*, *neutral*, or *negative*.
  - How can we learn this from a training set?
  - What do we need?
  - Explicitly. What
- (Keep your architecture safe for later!)

# In-Class Exercise (Exercise 3 in Lecture 07)

- 

Input: word embeddings, size: 10

1st layer: RNN, hidden dimension: 20

...



He did not

....



movie

word embeddings,  
size: 100

- 

- 

- (Keep your architecture safe for later!)

# RNNs for Sequence Labeling

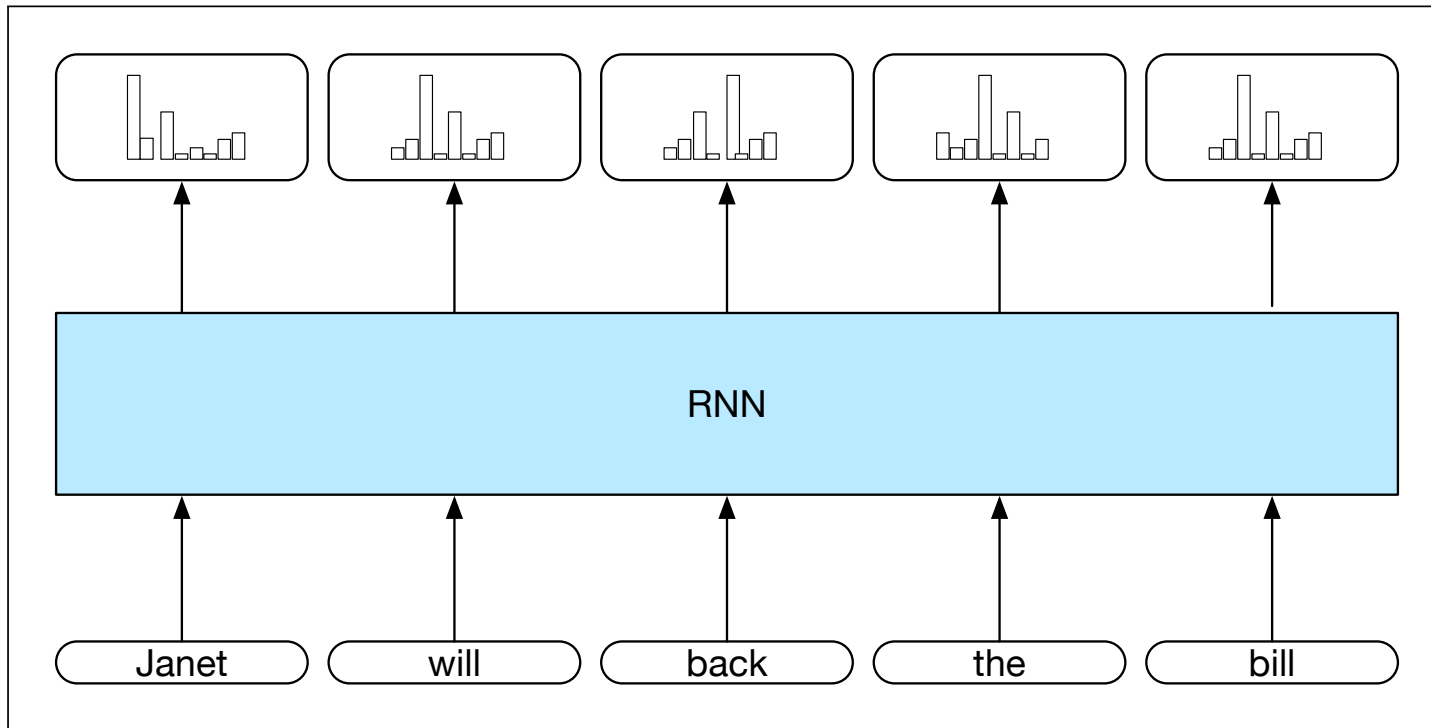
# Sequence Labeling with RNNs

- For sequence labeling, we don't exclusively use the last hidden state of the RNN
  - (Why?)

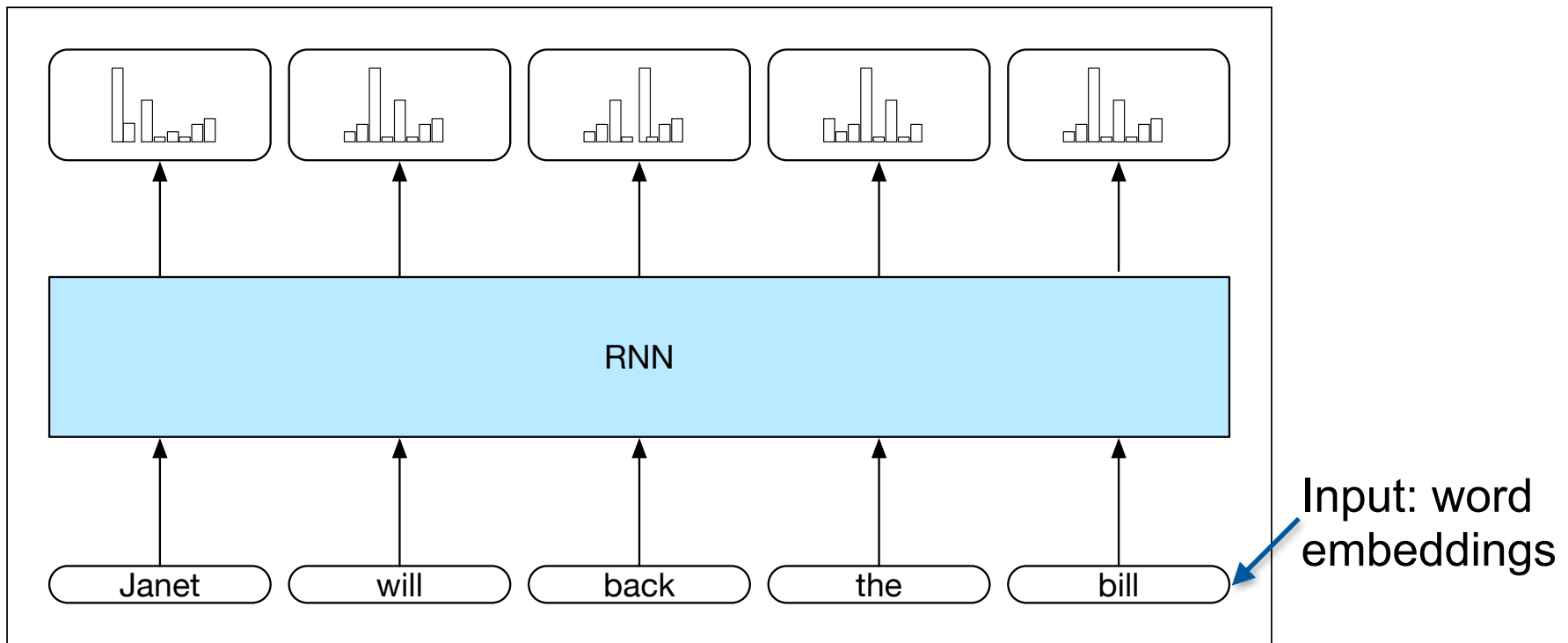
# Sequence Labeling with RNNs

- For sequence labeling, we don't exclusively use the last hidden state of the RNN
  - (Why?)
- Instead, we assign a class to each hidden state
  - Use a feedforward layer and a sigmoid/softmax function

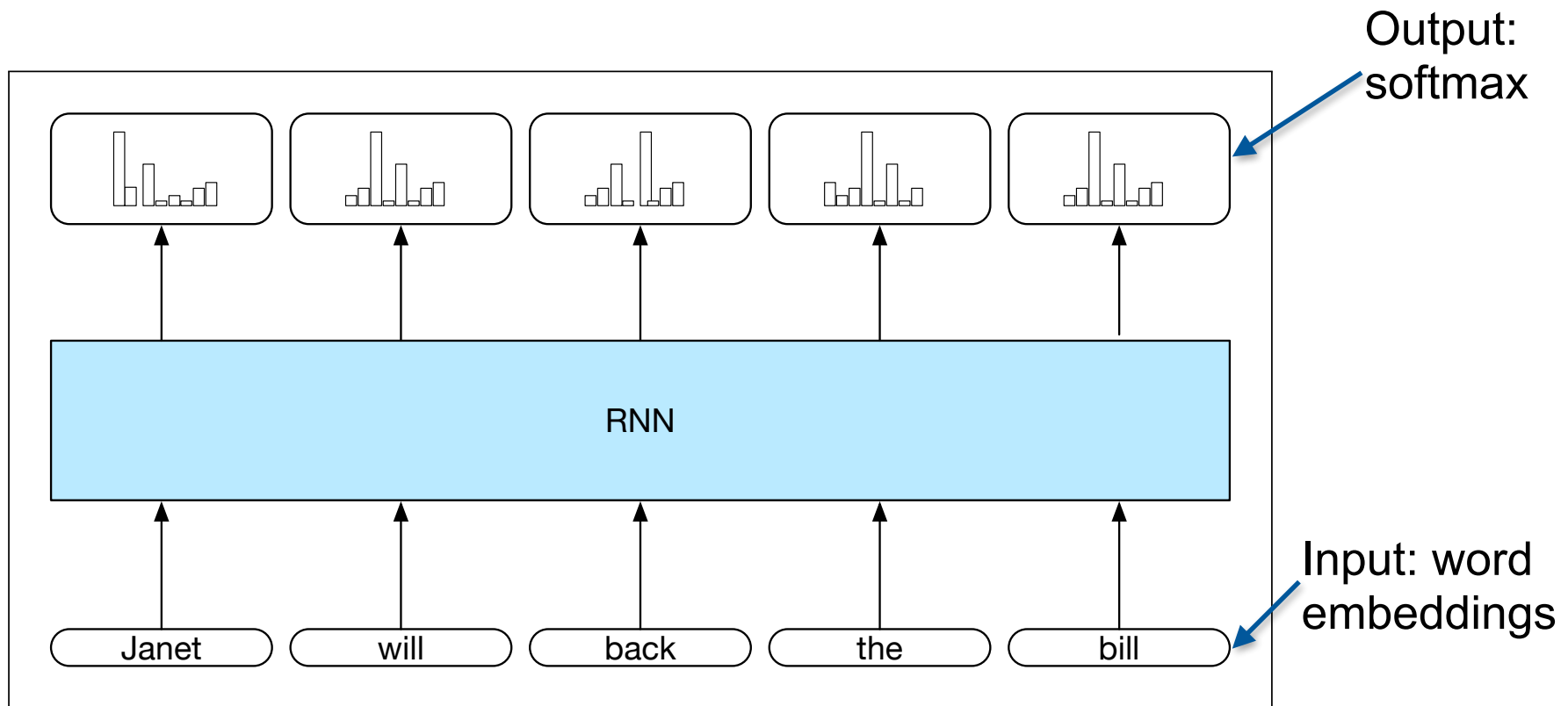
# Sequence Labeling with RNNs



# Sequence Labeling with RNNs



# Sequence Labeling with RNNs





# Example: Sequence Labeling with RNNs

Input:  $x_1, x_2, \dots, x_T$

$$h_i = \text{sigmoid}(W_h h_{i-1} + W_x x_i)$$

$$y_i = \text{argmax}(\text{softmax}(W_y h_i))$$

# Example: Sequence Labeling with RNNs

Input:  $x_1, x_2, \dots, x_T$

$$h_i = \text{sigmoid}(W_h h_{i-1} + W_x x_i)$$

$$y_i = \text{argmax}(\text{softmax}(W_y h_i))$$

- (What would be the formula to predict the label for sequence labeling?)

# Example: Sequence Labeling with RNNs

Input:  $x_1, x_2, \dots, x_T$

$$h_i = \text{sigmoid}(W_h h_{i-1} + W_x x_i)$$

$$y_i = \text{argmax}(\text{softmax}(W_y h_i))$$

- **we like books**
- $\text{we} = [1, 2]$ ,  $\text{like} = [0.5, 0.1]$ ,  $\text{books} = [0.2, 1]$
- $W_h = [[0.4]]$ ,  $W_x = [[0.1, 0.14]]$ ,  $W_y = [[0.3], [0.2], [0.01]]$ ,  $h_0 = [0]$
- Compute the class of “we”!

# Example: Sequence Labeling with RNNs

Input:  $x_1, x_2, \dots, x_T$

$$h_i = \text{sigmoid}(W_h h_{i-1} + W_x x_i)$$

$$y_i = \text{argmax}(\text{softmax}(W_y h_i))$$

- **we like books**
- $\text{we} = [1, 2]$ ,  $\text{like} = [0.5, 0.1]$ ,  $\text{books} = [0.2, 1]$
- $W_h = [[0.4]]$ ,  $W_x = [[0.1, 0.14]]$ ,  $W_y = [[0.3], [0.2], [0.01]]$ ,  $h_0 = [0]$
- Compute the class of “we”!
- $h_1 = \text{sigmoid}([0] + [0.1 + 0.28]) = \text{sigmoid}([0.38]) = 0.59$
- $y_1 = \text{argmax}(\text{softmax}([0.3 \cdot 0.59, 0.2 \cdot 0.59, 0.01 \cdot 0.59])) = \text{argmax}(\text{softmax}([1.77, 1.18, 0.0059])) = \text{argmax}([0.58, 0.32, 0.1])$

# Example: Sequence Labeling with RNNs

Input:  $x_1, x_2, \dots, x_T$

$$h_i = \text{sigmoid}(W_h h_{i-1} + W_x x_i)$$

$$y_i = \text{argmax}(\text{softmax}(W_y h_i))$$

- **we like books**

- $\text{we} = [1, 2]$ ,  $\text{like} = [0.5, 0.1]$ ,  $\text{books} = [0.2, 1]$
- $W_h = [[0.4]]$ ,  $W_x = [[0.1, 0.14]]$ ,  $W_y = [[0.3], [0.2], [0.01]]$ ,  $h_0 = [0]$
- Compute the class of “we”!
- $h_1 = \text{sigmoid}([0] + [0.1 + 0.28]) = \text{sigmoid}([0.38]) = 0.59$
- $y_1 = \text{argmax}(\text{softmax}([0.3 \cdot 0.59, 0.2 \cdot 0.59, 0.01 \cdot 0.59])) = \text{argmax}(\text{softmax}([1.77, 1.18, 0.0059])) = \text{argmax}([0.58, 0.32, 0.1])$

# In-Class Exercise 2

- **we like books**

Input:  $x_1, x_2, \dots, x_T$

$$h_i = \text{sigmoid}(W_h h_{i-1} + W_x x_i)$$

$$y_i = \text{argmax}(\text{softmax}(W_y h_i))$$

- $\text{we} = [1, 2]$ ,  $\text{like} = [0.5, 0.1]$ ,  $\text{books} = [0.2, 1]$
- $W_h = [[0.4]]$ ,  $W_x = [[0.1, 0.14]]$ ,  $W_y = [[0.3], [0.2], [0.01]]$ ,  $h_0 = [0]$
- $h_1 = \text{sigmoid}([0] + [0.1 + 0.28]) = \text{sigmoid}([0.38]) = 0.59$
- $y_1 = \text{argmax}(\text{softmax}([0.3 \cdot 0.59, 0.2 \cdot 0.59, 0.01 \cdot 0.59])) = \text{argmax}(\text{softmax}([1.77, 1.18, 0.0059])) = \text{argmax}([0.58, 0.32, 0.1])$
- Compute the classes of “like” and “books”!

# Hierarchical LSTMs

# Characters and POS Tagging

- Until now, we considered only word-level information for (neural) POS tagging
- Character information can be beneficial!



# Characters and POS Tagging

- Until now, we considered only word-level information for (neural) POS tagging
- Character information can be beneficial!
  - employer**er**, runner**er**, worker**er**, ...

# Characters and POS Tagging

- Until now, we considered only word-level information for (neural) POS tagging
- Character information can be beneficial!
  - employer**er**, runner**er**, worker**er**, ...
  - happiness**ness**, sadness**ness**, richness**ness**...

# Characters and POS Tagging

- Until now, we considered only word-level information for (neural) POS tagging
- Character information can be beneficial!
  - employ**er**, runner**er**, worker**er**, ...
  - happi**ness**, sad**ness**, rich**ness**...
  - help**ful**, joy**ful**, cheer**ful**, ...

# Characters and POS Tagging

- We want to use this information for POS tagging
- How can we do that?

# Characters and POS Tagging

- We want to use this information for POS tagging
- How can we do that?
  - RNN over characters

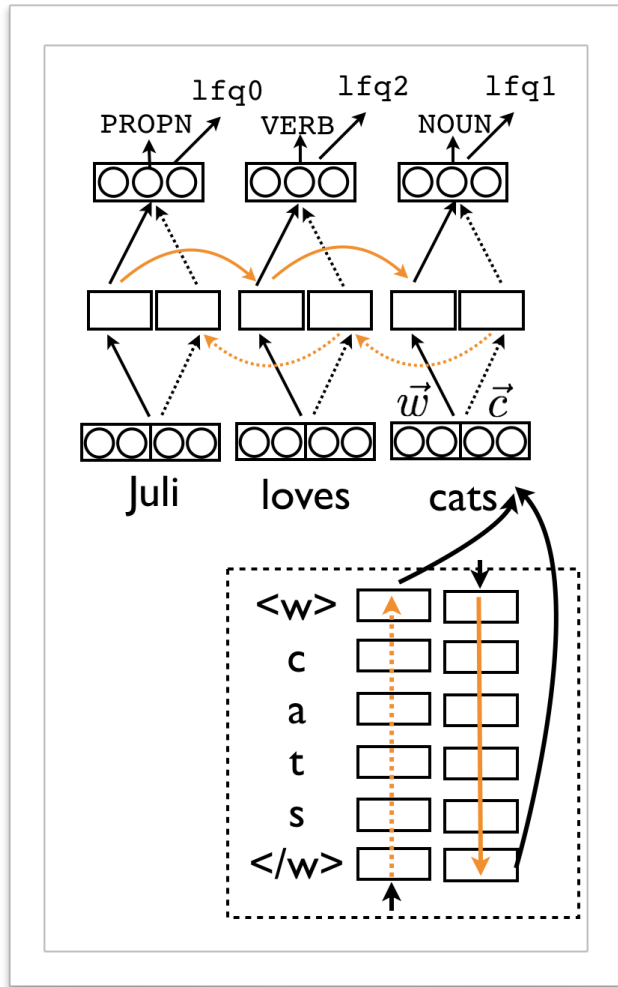
# Characters and POS Tagging

- We want to use this information for POS tagging
- How can we do that?
  - RNN over characters
  - CNN over characters

# Characters and POS Tagging

- We want to use this information for POS tagging
- How can we do that?
  - RNN over characters
  - CNN over characters
  - Two RNNs: one over characters and one over words!

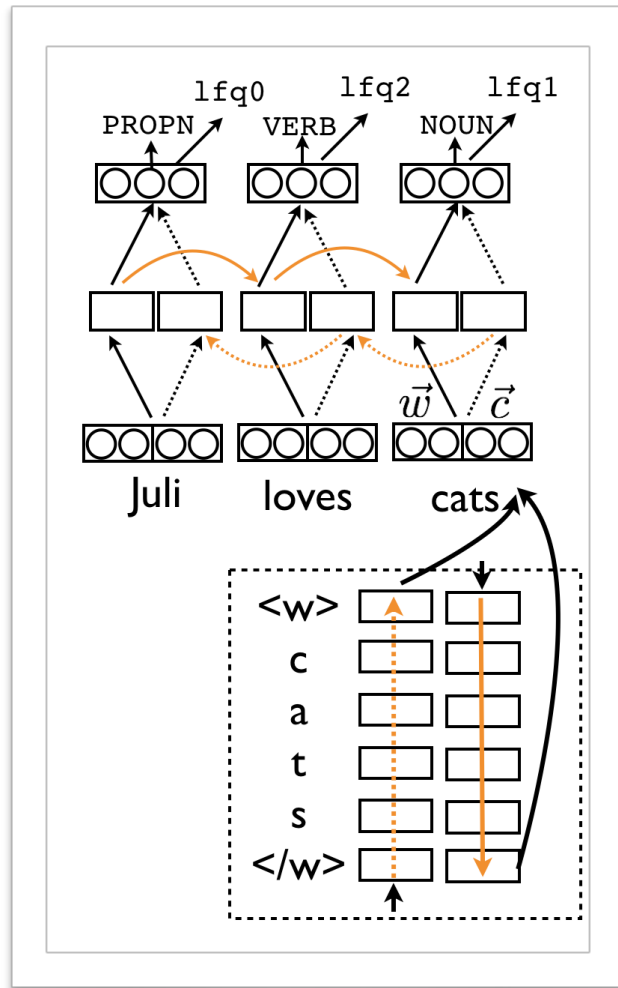
# Hierarchical RNN for POS Tagging



From [Plank et al., 2016](#)



# Hierarchical RNN for POS Tagging



- Concatenate last hidden states of the character-level RNN with the word embeddings!

# Other Sequence Labeling Tasks in NLP

# Named Entity Recognition

- Named entities are spans in a text that correspond to names of people, places or organizations:
  - Barack Obama
  - CNN
  - Donald Trump Jr.
- Task: Find named entities in a given text
  - This required finding **spans** which correspond to named entities

# IOB Encoding

- Label token with
  - B: beginning of span of interest
  - I: inside of span of interest
  - O: “outside”, not part of the span of interest

# Named Entity Recognition

(9.13) *United cancelled the flight from Denver to San Francisco.*

B        O                O O        O        B                O B        I

J&M, Ch.9

# Named Entity Recognition

(9.13) *United cancelled the flight from Denver to San Francisco.*

B O O O O B O B I

J&M, Ch.9

(9.14) *United cancelled the flight from Denver to San Francisco.*

B-ORG O O O B-LOC O B-LOC I-LOC

J&M, Ch.9

- We can encode more detail!

# More Sequence Labeling Tasks

- Morphological tagging
- Language identification for code-switched text
- Constituency parsing
- ...

# More Sequence Labeling Tasks

- Morphological tagging
- Language identification for code-switched text
- Constituency parsing
- ...

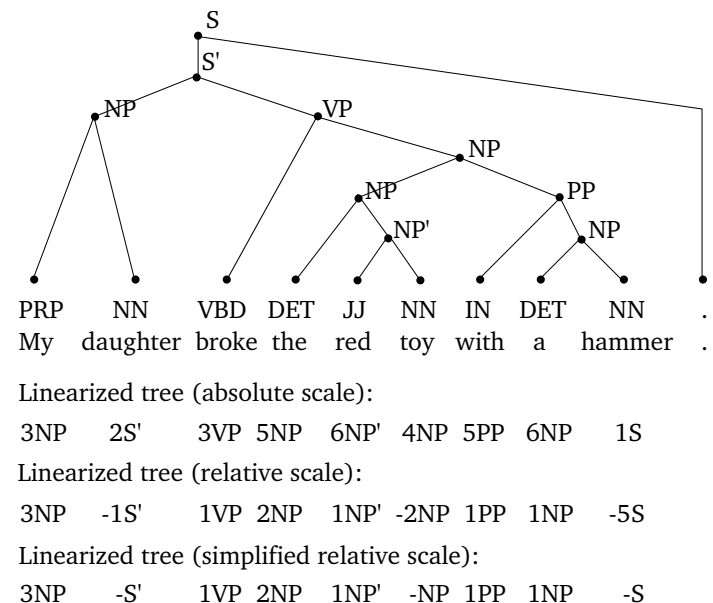


Figure 2: An example of a binarized constituency tree, linearized both applying absolute and relative scales.



# Wrapping up

- Discussed today:
  - Recap: RNNs
  - RNNs for sequence labeling tasks
  - Hierarchical RNNs
- On Wednesday: Training neural networks