# CSCE 689: Special Topics in Trustworthy NLP

## Lecture 4: Convolutional Neural Network, Recurrent Neural Network

Kuan-Hao Huang

khhuang@tamu.edu

Kuan-Hao Huang

khhuang@tamu.edu

(Some slides adapted from Chris Manning, Abigail See, Karthik Narasimhan, and Danqi Chen)

# LaTeX Assignment

- LaTeX Assignment (1%)
- Due: Sep 11, 11:59pm

---

**Your Name**
Your UID and email

### Overview

This assignment is designed to give you practice with LATEX, which you are expected to use for your literature review, project proposal, and final report in this course.

### Instructions

For this assignment, you will create a PDF containing your answers using LATEX. If this is your first time working with LATEX, we recommend starting with this short tutorial, which covers the basic features you will need for this course. Please use the Association for Computational Linguistics LATEX template (link), a template widely used in major NLP conferences. We suggest using Overleaf as your online editor, since it automatically manages packages for you.

By default, the template is set to *review mode*. To switch to *final mode*, change:

```
Review Mode (Default)
\usepackage[review]{acl}
```

to:

```
Final Mode
\usepackage[final]{acl}
```

This allows you to display the author information. Be sure to include your name, UID, and email.

The following sections contain questions on some commonly used LATEX commands. There are a total of 100 points for this assignment. Please answer each question in a separate *section*, and submit the final PDF generated using LATEX.

### 1 Including Equations [20pts]

Typeset the following expression using LATEX:

$$\frac{\partial \mathcal{L}_{\text{total}}}{\partial \mathbf{w}_j} = -\frac{1}{m} \sum_{i=1}^{m} \big( y_i - \sigma(z_i) \big) \cdot \mathbf{x}_{i,j}$$

### 2 Including Images [20pts]

Select a picture of a cat and include it with a caption. The figure below is provided as an example.

Figure 1: This is a cute cat!

### 3 Including Tables [20pts]

Create a table that displays your name, UIN, and email. You can follow the example below as a template.

| Name | Kuan-Hao Huang |
|-------|------------------|
| UIN | 123456789 |
| Email | khhuang@tamu.edu |

Table 1: Example table.

### 4 Including Lists [20pts]

Create a list that displays your name, UIN, and email. You can follow the example below as a template.

- Name: Kuan-Hao Huang
- UIN: 123456789
- Email: khhuang@tamu.edu

### 5 Including Citations [20pts]

Use *BibTex* to include the following paper: Paper 1 (Vaswani et al., 2017) and Paper 2 (Devlin et al., 2019). You can learn more about *BibTex* here.
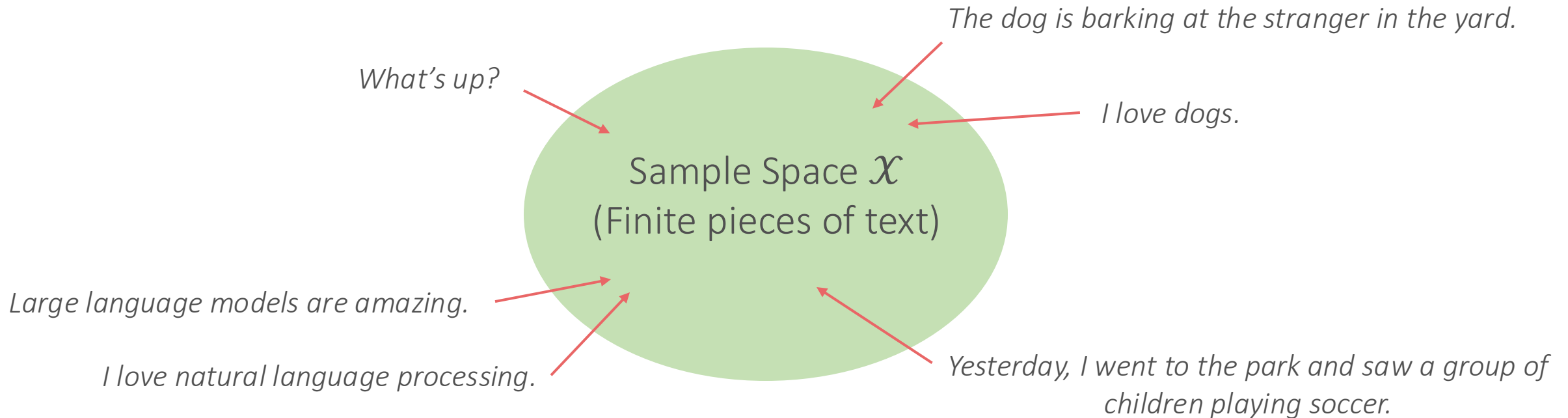
# Topic Sign-Up

- Sign-up: https://tinyurl.com/2p9mr2wa
  - Log in with TAMU account
  - Due: Sep 10 before lecture

| | | | Put Preference with Topic IDs | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Team** | **Member 1** | **Member 2 (optional)** | **Preference 1** | **Preference 2** | **Preference 3** | **Preference 4** | **Preference 5** | **Preference 6** | **Preference 7** | **Preference 8** | **Preference 9** |
| Example | First_name Last_name | First_name Last_name | 4 | 10 | 1 | 7 | 3 | 9 | 8 | 6 | 12 |
| 1 | Kunal Jain | | 3 | 12 | 8 | 6 | 11 | 10 | 9 | | |
| 2 | Muhan Gao | | 5 | 6 | 4 | 9 | 10 | 7 | 8 | 11 | 12 |
| 3 | Serhii Honcharenko | | 9 | 10 | 4 | 7 | 6 | 1 | | | |
| 4 | Oscar Chew | | 9 | 3 | 11 | 7 | 8 | 10 | 4 | 2 | 1 |
| 5 | Junggeun Do | | 8 | 7 | 11 | 10 | 2 | 6 | 12 | 11 | 5 |
| 6 | Jiongran Wang | | 3 | 9 | 10 | 6 | 4 | 2 | 7 | 5 | 1 |
| 7 | Sicong Liang | | 10 | 9 | 3 | 13 | 7 | | | | |
| 8 | Kowsalya Balamuralei Umamaheswari | | 3 | 10 | 12 | 8 | | | | | |
| 9 | Yi Wen | | 12 | 4 | 7 | 9 | 6 | 8 | | | |
| 10 | Quang Nguyen | | 6 | 3 | 4 | 12 | 2 | 1 | 8 | 9 | 7 |
| 11 | Aaron Xu | | 4 | 8 | 7 | 3 | 9 | 6 | 11 | 2 | 1 |
| 12 | Bhaskar Ruthvik Bikkina | | 7 | 11 | 12 | 9 | 10 | 5 | 3 | 4 | 6 |
| 13 | | | | | | | | | | | |
| 14 | | | | | | | | | | | |
| 15 | | | | | | | | | | | |
| 16 | | | | | | | | | | | |
| 17 | | | | | | | | | | | |

# Recap: Language Models

- Learn the probability distribution over texts $x = [w_1, w_2, \ldots, w_l] \in \mathcal{X}$

$$P(x) = P(w_1, w_2, \ldots, w_l)$$

*The dog is barking at the stranger in the yard.*

*What's up?*

*I love dogs.*

Sample Space $\mathcal{X}$
(Finite pieces of text)

*Large language models are amazing.*

*I love natural language processing.*

*Yesterday, I went to the park and saw a group of children playing soccer.*

# Recap: Auto-Regressive Language Models

$$P(w_1, w_2, w_3, \dots, w_l) = P(w_1)P(w_2, w_3, \dots, w_l | w_1)$$

$$= P(w_1)P(w_2|w_1)(w_3, \dots, w_l | w_1, w_2)$$

$$= P(w_1)P(w_2|w_1)P(w_3|w_1, w_2)(w_4, \dots, w_l | w_1, w_2, w_3)$$

$$= \prod_{i=1}^{l} P(w_i | w_1, w_2, \dots, w_{i-1})$$

$$P(\textit{She likes to go hiking}) = P(\textit{She}) \cdot P(\text{likes}|\text{She}) \cdot P(\text{to}|\text{She likes})$$

$$\cdot P(\text{go}|\text{She likes to}) \cdot P(\text{hiking}|\text{She likes to go})$$

# Recap: N-Gram Language Models

Assumption: $P(w_i | w_1, w_2, \dots, w_{i-1}) \approx P(w_i | w_{i-n+1}, \dots, w_{i-1})$

$$P(w_1, w_2, w_3, \dots, w_l) \approx \prod_i P(w_i | w_{i-n+1}, \dots, w_{i-1})$$
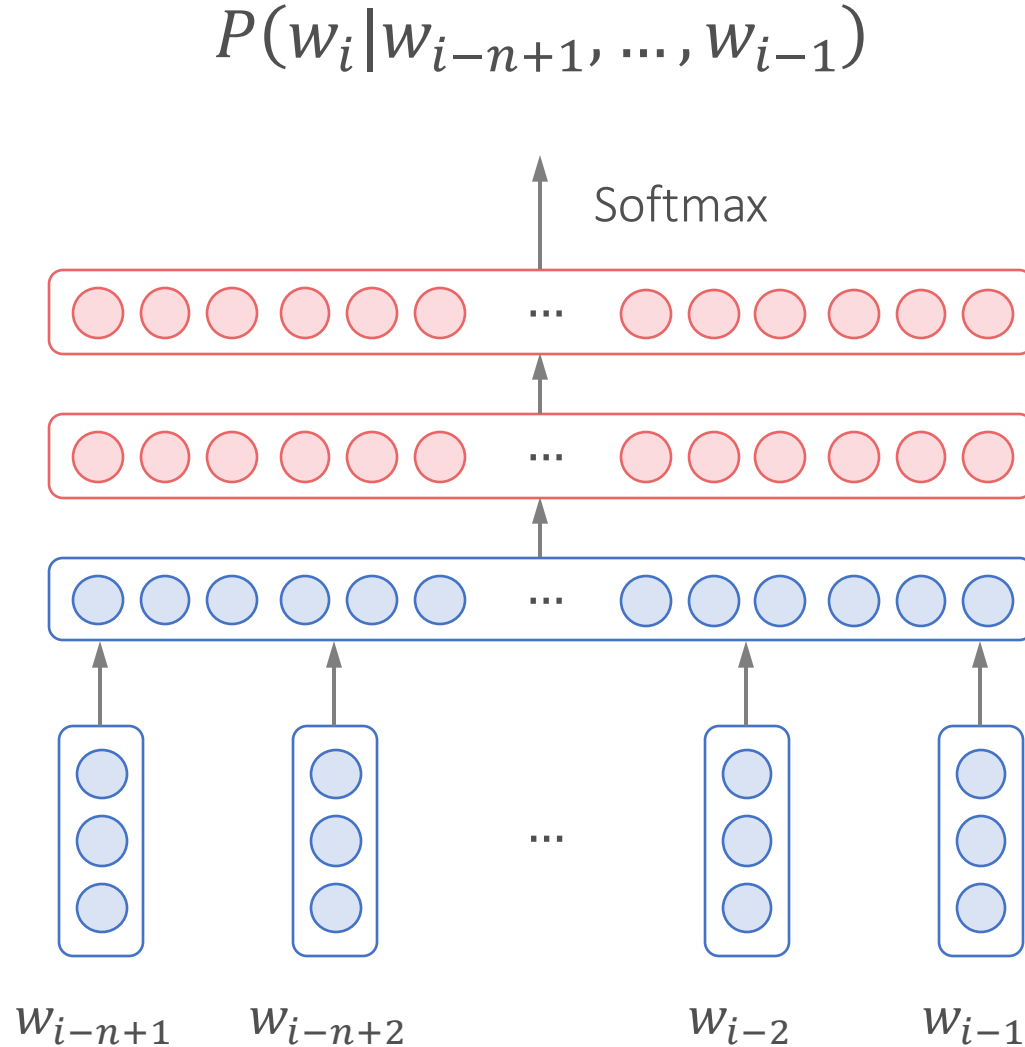
The prediction of the next token depends on the previous **n** tokens

A count-based solution: Collect training corpus and count

$$P(w_i | w_{i-n+1}, \dots, w_{i-1}) = \frac{C_{train}(w_{i-n+1}, \dots, w_{i-1}, w_i)}{C_{train}(w_{i-n+1}, \dots, w_{i-1})}$$

# Recap: What Can Language Models Do?

- Score texts

$P$(*The dog is barking at the stranger in the yard.*)  → High

$P$(*Cats upon they chairs sleeping their dreams fall.*)  → Low

- Generate texts

$$\tilde{x} \sim P(\mathcal{X})$$

Compute perplexity

Sample from word distribution

# N-Gram Language Models

Assumption: $P(w_i|w_1, w_2, \ldots, w_{i-1}) \approx P(w_i|w_{i-n+1}, \ldots, w_{i-1})$

$$P(w_1, w_2, w_3, \ldots, w_l) \approx \prod_i P(w_i|w_{i-n+1}, \ldots, w_{i-1})$$

The prediction of the next token depends on the previous **n** tokens

A count-based solution: Collect training corpus and count

$$P(w_i|w_{i-n+1}, \ldots, w_{i-1}) = \frac{C_{train}(w_{i-n+1}, \ldots, w_{i-1}, w_i)}{C_{train}(w_{i-n+1}, \ldots, w_{i-1})}$$

Any Problems?

# Unseen Patterns in Training Corpus

- Not all n-grams in the test set will be observed in training corpus
- Training corpus
  - I like apples
  - I love oranges
- Test set
  - I love apples
  - I like oranges
- This problem becomes severe when n is large

# Laplace Smoothing

- Handle sparsity by making sure all probabilities are non-zero in our model
  - Just add $\alpha$ to all counts and renormalize

Bigram language model before smoothing

$$P(w_i|w_{i-1}) = \frac{C_{train}(w_{i-1}, w_i)}{C_{train}(w_{i-1})}$$

Bigram language model after smoothing

$$P(w_i|w_{i-1}) = \frac{C_{train}(w_{i-1}, w_i) + \alpha}{C_{train}(w_{i-1}) + \alpha|\mathcal{V}|}$$

# Linear Interpolation

$$\hat{P}(w_i|w_{i-1}, w_{i-2}) = \lambda_1 P(w_i|w_{i-1}, w_{i-2}) \qquad \text{Trigram}$$

$$+ \lambda_2 P(w_i|w_{i-1}) \qquad \text{Bigram}$$

$$+ \lambda_3 P(w_i) \qquad \text{Unigram}$$

$$\sum_i \lambda_i = 1$$

Strong empirical performance!

# N-Gram Language Models

Assumption: $P(w_i|w_1, w_2, \ldots, w_{i-1}) \approx P(w_i|w_{i-n+1}, \ldots, w_{i-1})$

$$P(w_1, w_2, w_3, \ldots, w_l) \approx \prod_i P(w_i|w_{i-n+1}, \ldots, w_{i-1})$$

The prediction of the next token depends on the previous **n** tokens

A count-based solution: Collect training corpus and count

$$P(w_i|w_{i-n+1}, \ldots, w_{i-1}) = \frac{C_{train}(w_{i-n+1}, \ldots, w_{i-1}, w_i)}{C_{train}(w_{i-n+1}, \ldots, w_{i-1})}$$

Can we compute the probability in a different way?

# Neural Language Models

$$P(w_i | w_{i-n+1}, \ldots, w_{i-1})$$



Softmax

Training corpus
- I like apples
- I love oranges

Test set
- I love apples
- I like oranges

$w_{i-n+1}$  $w_{i-n+2}$  $\ldots$  $w_{i-2}$  $w_{i-1}$

# Auto-Regressive Language Models

$$P(w_1, w_2, w_3, \dots, w_l) = P(w_1)P(w_2, w_3, \dots, w_l|w_1)$$

$$= P(w_1)P(w_2|w_1)(w_3, \dots, w_l|w_1, w_2)$$

$$= P(w_1)P(w_2|w_1)P(w_3|w_1, w_2)(w_4, \dots, w_l|w_1, w_2, w_3)$$

$$= \prod_{i=1}^{l} P(w_i|w_1, w_2, \dots, w_{i-1})$$

$P(\textit{She likes to go hiking}) = P(\textit{She}) \cdot P(\text{likes}|\text{She}) \cdot P(\text{to}|\text{She likes})$

$\cdot P(\text{go}|\text{She likes to}) \cdot P(\text{hiking}|\text{She likes to go})$

# Recap: A General Framework for Text Classification

Text $x$ → **Feature (Representation)** → **Classifier (Model)** → Label $y$

- Teach the model how to make prediction $y$
- Logistic regression, neural networks, CNN, RNN, LSTM, Transformers

# Recap: Logistic Regression

- Logistic Regression for multiclass classification

Feature Vector $\mathbf{x} = [x_1, x_2, x_3, \ldots, x_d]$     Label $y = 0, 1, \ldots, C-1$

Weight Vectors $\mathbf{w}_c = [w_{c,1}, w_{c,2}, w_{c,3}, \ldots, w_{c,d}]$     Bias $b_c$

Learnable Parameters

$$z_c = \mathbf{w}_c \cdot \mathbf{x} + b_c$$

$$P(y = c | \mathbf{x}) = \mathrm{softmax}(z_c)$$

$$\mathrm{softmax}(z_c) = \frac{e^{z_c}}{\sum_t e^{z_t}}$$

Softmax Function

$$\mathrm{Prediction} = \arg\max_c P(y = c | \mathbf{x})$$

# Recap: Neural Networks



Prediction $= \underset{c}{\arg \max} \, \tilde{y}_c$

Multiclass Cross Entropy Loss

$$\mathcal{L}_{CE}(y, \tilde{y}) = -\sum_{c=0}^{C} y_c \log P(y = c \mid \mathbf{x})$$

# A Simple Approach: Averaged Embeddings + DNN



| | Alice | treats | Bob | well |
|---|---|---|---|---|
| Dimension 1 | 0.7 | 2.7 | -0.1 | -5.7 |
| Dimension 2 | 8.6 | -3.9 | 6.7 | -9.8 |
| Dimension 3 | -2.4 | -5.6 | 1.5 | -1.6 |
| Dimension 4 | 2.3 | 1.1 | 2.0 | -1.0 |

-0.6
0.4
-1.6
1.1

Any problems?

$$\mathcal{L}_{CE}(y, \tilde{y}) = -\sum_{c=0}^{C} y_c \log P(y = c | \mathbf{x})$$

# A Simple Approach: Concatenated Embeddings + DNN

| | Alice | treats | Bob | well |
|---|---|---|---|---|
| Dimension 1 | 0.7 | 2.7 | -0.1 | -5.7 |
| Dimension 2 | 8.6 | -3.9 | 6.7 | -9.8 |
| Dimension 3 | -2.4 | -5.6 | 1.5 | -1.6 |
| Dimension 4 | 2.3 | 1.1 | 2.0 | -1.0 |

0.7
8.6
-2.4
2.3
2.7
-3.9
-5.6
1.1
...
-5.7
-9.8
-1.6
-1.0

Any problems?

$$\mathcal{L}_{CE}(y, \tilde{y}) = -\sum_{c=0}^{C} y_c \log P(y = c \mid \mathbf{x})$$

# Challenges

- Averaged Embeddings
  - Lose order information
- Concatenated Embeddings
  - Cannot handle various lengths

# Solution: Capture Local Order Information

*Bob likes Alice very much*

| | |
|---|---|
| **Unigram** | *{Bob, likes, Alice, very, much}* |
| **Bigram** | *{Bob likes, likes Alice, Alice very, very much}* |
| **Trigram** | *{Bob likes Alice, likes Alice very, Alice very much}* |
| **4-gram** | *{Bob likes Alice very, likes Alice very much}* |

We can infer global order information from local order information

# Convolutional Neural Network (CNN)

- Capture local features (N-grams)
  - Filters (Kernels)
- Hierarchical feature learning
  - Multiple layers

# Convolutional Neural Network (CNN)

Learnable Weight (Filter)
Filter Size = 3

$$\mathbf{W} = \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} \\ \dots & \dots & \dots \\ w_{4,1} & w_{4,2} & w_{4,3} \end{bmatrix}$$
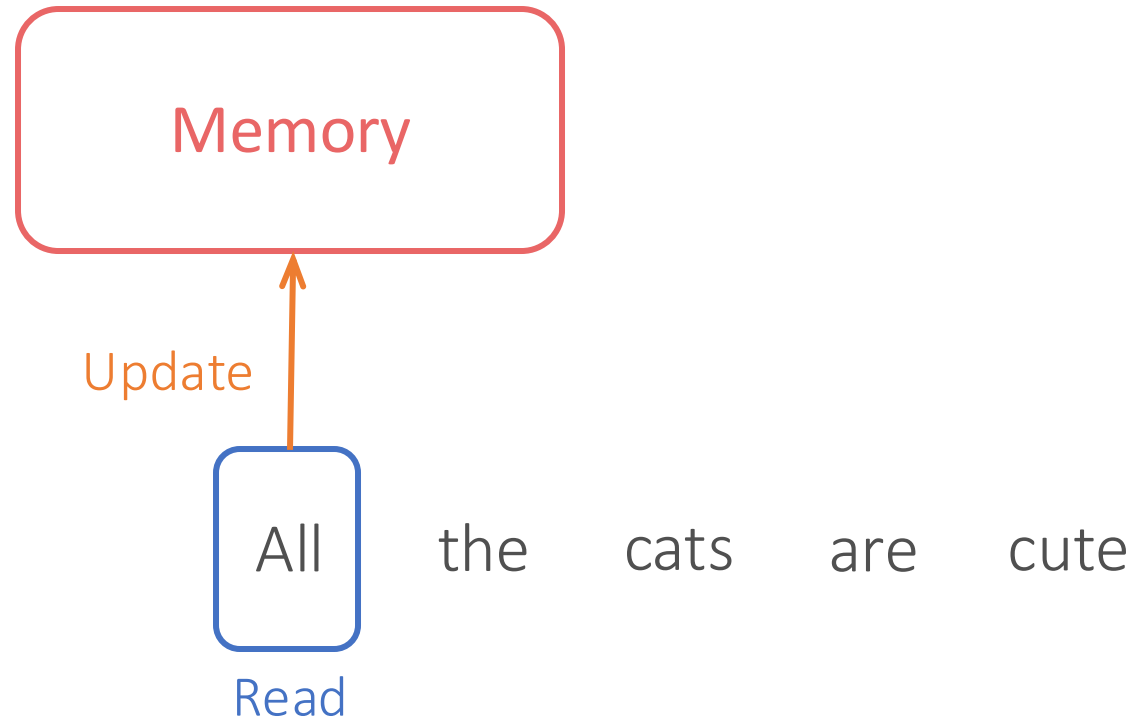
| Alice | treats | Bob | well |
|-------|--------|-----|------|

Dimension 1 | 0.7 | 2.7 | -0.1 | -5.7

$w_{1,1}*0.7 \ + \ w_{1,2}*2.7 + \ w_{1,3}*-0.1$

Dimension 2 | 8.6 | -3.9 | 6.7 | -9.8

$w_{2,1}*8.6 \ + \ w_{2,2}*-3.9 + \ w_{2,3}*6.7$

Dimension 3 | -2.4 | -5.6 | 1.5 | -1.6

$w_{3,1}*-2.4 + \ w_{3,2}*-5.6 + \ w_{3,3}*1.5$

Dimension 4 | 2.3 | 1.1 | 2.0 | -1.0

$w_{4,1}*2.3 \ + \ w_{4,2}*1.1 \ + \ w_{4,3}*2.0$

$v_{1,1}$

$+$

$v_{1,2}$

$+$

$v_{1,3}$

$+$

$v_{1,4}$

$v_1$

# Convolutional Neural Network (CNN)

Learnable Weight (Filter)
Filter Size = 3

$$\mathbf{W} = \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} \\ \dots & \dots & \dots \\ w_{4,1} & w_{4,2} & w_{4,3} \end{bmatrix}$$



Alice | treats | Bob | well

Dimension 1: 0.7 | 2.7 | -0.1 | -5.7
$w_{1,1}*2.7 + w_{1,2}*-0.1 + w_{1,3}*-5.7$

Dimension 2: 8.6 | -3.9 | 6.7 | -9.8
$w_{2,1}*-3.9 + w_{2,2}*6.7 + w_{2,3}*-9.8$

Dimension 3: -2.4 | -5.6 | 1.5 | -1.6
$w_{3,1}*-5.6 + w_{3,2}*1.5 + w_{3,3}*-1.6$

Dimension 4: 2.3 | 1.1 | 2.0 | -1.0
$w_{4,1}*1.1 + w_{4,2}*2.0 + w_{4,3}*-1.0$

$v_{1,1}$ | $v_{2,1}$
+ | +
$v_{1,2}$ | $v_{2,2}$
+ | +
$v_{1,3}$ | $v_{2,3}$
+ | +
$v_{1,4}$ | $v_{2,4}$

$v_1$ | $v_2$

# Convolutional Neural Network (CNN)

Learnable Weight (Filter)
Filter Size = 3

$$W = \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} \\ \dots & \dots & \dots \\ w_{4,1} & w_{4,2} & w_{4,3} \end{bmatrix}$$

| | Alice | treats | Bob | well | | so | Bob | thinks |
|---|---|---|---|---|---|---|---|---|

| Dimension 1 | 0.7 | 2.7 | -0.1 | -5.7 | | $v_{1,1}$ | $v_{2,1}$ | $v_{3,1}$ | $v_{4,1}$ | $v_{5,1}$ |
| | | | | | | + | + | + | + | + |
| Dimension 2 | 8.6 | -3.9 | 6.7 | -9.8 | | $v_{1,2}$ | $v_{2,2}$ | $v_{3,2}$ | $v_{4,2}$ | $v_{5,2}$ |
| | | | | | | + | + | + | + | + |
| Dimension 3 | -2.4 | -5.6 | 1.5 | -1.6 | | $v_{1,3}$ | $v_{2,3}$ | $v_{3,3}$ | $v_{4,3}$ | $v_{5,3}$ |
| | | | | | | + | + | + | + | + |
| Dimension 4 | 2.3 | 1.1 | 2.0 | -1.0 | | $v_{1,4}$ | $v_{2,4}$ | $v_{3,4}$ | $v_{4,4}$ | $v_{5,4}$ |

Max Pooling

$v$

| $v_1$ | $v_2$ | $v_3$ | $v_4$ | $v_5$ |

# Convolutional Neural Network (CNN)

Learnable Weight (Filter)
Filter Size = 3

$$\mathbf{W} = \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} \\ \cdots & \cdots & \cdots \\ w_{4,1} & w_{4,2} & w_{4,3} \end{bmatrix} \quad \mathbf{W} = \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} \\ \cdots & \cdots & \cdots \\ w_{4,1} & w_{4,2} & w_{4,3} \end{bmatrix} \quad \mathbf{W} = \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} \\ \cdots & \cdots & \cdots \\ w_{4,1} & w_{4,2} & w_{4,3} \end{bmatrix}$$

| | Alice | treats | Bob | well |
|---|---|---|---|---|
| Dimension 1 | 0.7 | 2.7 | -0.1 | -5.7 |
| Dimension 2 | 8.6 | -3.9 | 6.7 | -9.8 |
| Dimension 3 | -2.4 | -5.6 | 1.5 | -1.6 |
| Dimension 4 | 2.3 | 1.1 | 2.0 | -1.0 |

$v$   $v$   $v$

# Convolutional Neural Network (CNN)

Filter Size = 3 $\quad \mathbf{W} = \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} \\ \cdots & \cdots & \cdots \\ w_{4,1} & w_{4,2} & w_{4,3} \end{bmatrix}$ Filter Size = 2 $\quad \mathbf{W} = \begin{bmatrix} w_{1,1} & w_{1,2} \\ \cdots & \cdots \\ w_{4,1} & w_{4,2} \end{bmatrix}$ Filter Size = 4 $\quad \mathbf{W} = \begin{bmatrix} w_{1,1} & \cdots & w_{1,4} \\ \cdots & \cdots & \cdots \\ w_{4,1} & \cdots & w_{4,4} \end{bmatrix}$

| Alice | treats | Bob | well |

Dimension 1: 0.7 | 2.7 | -0.1 | -5.7

$w_{1,1}*0.7 \; + \; w_{1,2}*2.7$

Dimension 2: 8.6 | -3.9 | 6.7 | -9.8

$w_{2,1}*8.6 \; + \; w_{2,2}*-3.9 + w_{2,3}*6.7 + w_{2,4}*-9.8$

$v$ $\qquad v \qquad v$

Dimension 3: -2.4 | -5.6 | 1.5 | -1.6

Dimension 4: 2.3 | 1.1 | 2.0 | -1.0

# Convolutional Neural Network (CNN)

| | Alice | treats | Bob | well |
|---|---|---|---|---|
| Dimension 1 | 0.7 | 2.7 | -0.1 | -5.7 |
| Dimension 2 | 8.6 | -3.9 | 6.7 | -9.8 |
| Dimension 3 | -2.4 | -5.6 | 1.5 | -1.6 |
| Dimension 4 | 2.3 | 1.1 | 2.0 | -1.0 |

$v$ $v$ $v$ $v$ $v$

$$\mathcal{L}_{CE}(y, \tilde{y}) = -\sum_{c=0}^{C} y_c \log P(y = c \mid \mathbf{x})$$

# From Single Layer to Multiple Layers

Learnable Weight (Layer 1)
Filter Size = 3

$$\mathbf{W} = \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} \\ \dots & \dots & \dots \\ w_{4,1} & w_{4,2} & w_{4,3} \end{bmatrix}$$

Learnable Weight (Layer 2)
Filter Size = 3

$$\mathbf{W} = \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} \\ \dots & \dots & \dots \\ w_{4,1} & w_{4,2} & w_{4,3} \end{bmatrix}$$

| Alice | treats | Bob | well | so | Bob | thinks |

Dimension 1: 0.7 | 2.7 | -0.1 | -5.7 | $v_{1,1}$ | $v_{2,1}$ | $v_{3,1}$ | $v_{4,1}$ | $v_{1,1}$ | $v_{2,1}$

Dimension 2: 8.6 | -3.9 | 6.7 | -9.8 | $v_{1,2}$ | $v_{2,2}$ | $v_{3,2}$ | $v_{4,2}$ | $v_{1,2}$ | $v_{2,2}$

Dimension 3: -2.4 | -5.6 | 1.5 | -1.6 | $v_{1,3}$ | $v_{2,3}$ | $v_{3,3}$ | $v_{4,3}$ | $v_{1,3}$ | $v_{2,3}$

Dimension 4: 2.3 | 1.1 | 2.0 | -1.0 | $v_{1,4}$ | $v_{2,4}$ | $v_{3,4}$ | $v_{4,4}$ | $v_{1,4}$ | $v_{2,4}$

Capture high-order or hierarchical information

# Convolutional Neural Network (CNN)

- Capture local features (N-grams)
  - Filters (Kernels)
- Hierarchical feature learning
  - Multiple layers
- The whole process is still not similar to how human read texts
- Can we model reading texts in a sequential way?

# Recurrent Neural Network (RNN)

- Read texts sequentially like a human with memory
  - Read → update memory



Memory

Update

All    the    cats    are    cute

Read

# Recurrent Neural Network (RNN)

- Read texts sequentially like a human with memory
  - Read → update memory



All   the   cats   are   cute

# Recurrent Neural Network (RNN)

- Read texts sequentially like a human with memory
  - Read → update memory



Memory

Update

All    the    cats    are    cute

Read

# Recurrent Neural Network (RNN)

- Read texts sequentially like a human with memory
  - Read → update memory

# Recurrent Neural Network (RNN)

- Read texts sequentially like a human with memory
  - Read → update memory



Memory

Understanding

Update

All    the    cats    are    cute

Read

# Recurrent Neural Network (RNN)

- Recurrent unit

Update $\quad W, b$

Hidden state $h_t$

Memory

Read $\quad U$

Learnable parameters

$$W, b$$

$$U$$

# Recurrent Neural Network (RNN)

Hidden States

Read

$$h_t = \sigma(W h_{t-1} + U x_t + b)$$

Activation Function
(tanh, sigmoid)

Update

# Recurrent Neural Network (RNN)

$$h_t = \sigma(Wh_{t-1} + Ux_t + b)$$



$$\mathcal{L}_{CE}(y, \tilde{y}) = -\sum_{c=0}^{C} y_c \log P(y = c | \mathbf{x})$$

# Recurrent Neural Network (RNN)

- Advantages
  - Can process any length input
  - Model size doesn't increase for longer input context
  - Computation for step t can (in theory) use information from many steps back
- Disadvantages
  - Recurrent computation is slow
  - In practice, difficult to access information from many steps back
  - Vanishing gradient

# Back-Propagation



$$\frac{\partial \mathcal{L}}{\partial \tilde{y}}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{h}^{(3)}} = \frac{\partial \mathcal{L}}{\partial \tilde{y}} \cdot \frac{\partial \tilde{y}}{\partial \mathbf{h}^{(3)}}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(1)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{h}^{(1)}} \cdot \frac{\partial \mathbf{h}^{(1)}}{\partial \mathbf{W}^{(1)}}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{h}^{(2)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{h}^{(3)}} \cdot \frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{h}^{(2)}}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}^{(1)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{h}^{(1)}} \cdot \frac{\partial \mathbf{h}^{(1)}}{\partial \mathbf{b}^{(1)}}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{h}^{(1)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{h}^{(2)}} \cdot \frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{h}^{(1)}}$$

# Vanishing Gradient Problem

$$h_2 = \sigma(Wh_1 + Ux_2 + b)$$

$$h_3 = \sigma(Wh_2 + Ux_3 + b)$$

$$h_4 = \sigma(Wh_3 + Ux_4 + b)$$

$$h_5 = \sigma(Wh_4 + Ux_3 + b)$$



$$\frac{\partial \mathcal{L}}{\partial W} = \frac{\partial \mathcal{L}}{\partial h_5} \cdot \frac{\partial h_5}{\partial W} + \frac{\partial \mathcal{L}}{\partial h_5} \cdot \boxed{\frac{\partial h_5}{\partial h_4}} \cdot \frac{\partial h_4}{\partial W} + \frac{\partial \mathcal{L}}{\partial h_5} \cdot \boxed{\frac{\partial h_5}{\partial h_4} \cdot \frac{\partial h_4}{\partial h_3}} \cdot \frac{\partial h_3}{\partial W}$$

$$+ \frac{\partial \mathcal{L}}{\partial h_5} \cdot \boxed{\frac{\partial h_5}{\partial h_4} \cdot \frac{\partial h_4}{\partial h_3} \cdot \frac{\partial h_3}{\partial h_2}} \cdot \frac{\partial h_2}{\partial W} + \frac{\partial \mathcal{L}}{\partial h_5} \cdot \boxed{\frac{\partial h_5}{\partial h_4} \cdot \frac{\partial h_4}{\partial h_3} \cdot \frac{\partial h_3}{\partial h_2} \cdot \frac{\partial h_2}{\partial h_1}} \cdot \frac{\partial h_1}{\partial W}$$

When these are small, the gradient signal gets smaller and smaller as it back-propagates further

Model weights are updated only with respect to short-term effect rather than long-term effect

# Long Short-Term Memory (LSTM)

- Short-term memory: hidden state $h_t$

- Long-term memory: cell state $c_t$

- Key idea

  - Turn <span style="color:red">multiplication</span> into <span style="color:green">addition</span> (partially reduce gradient vanishing)

  - use <span style="color:green">gates</span> to control how much information to add/erase

# Recurrent Neural Network (RNN)



$$h_t = \sigma(W h_{t-1} + U x_t + b)$$

# Long Short-Term Memory (LSTM)



Neural Network Layer

Pointwise Operation

Vector Transfer

Concatenate

Copy

# Long Short-Term Memory (LSTM)



The cell state stores long-term information

# Long Short-Term Memory (LSTM)



The hidden state stores short-term information

https://colah.github.io/posts/2015-08-Understanding-LSTMs/

# Long Short-Term Memory (LSTM)



The cell state stores long-term information

Whenever reading a word, we will write/forget information to the cell state

# Long Short-Term Memory (LSTM)



Sigmoid function: gate values are between 0 and 1

Input gate

$$i_t = \sigma\big(W^{(i)}h_{t-1} + U^{(i)}x_t + b^{(i)}\big)$$

How much we should write

New information

$$\tilde{C}_t = \tanh\big(W^{(C)}h_{t-1} + U^{(C)}x_t + b^{(C)}\big)$$

What we should write

# Long Short-Term Memory (LSTM)



Sigmoid function: gate values are between 0 and 1

Forget gate

$$f_t = \sigma\big(W^{(f)}h_{t-1} + U^{(f)}x_t + b^{(f)}\big)$$

How much we should erase

# Long Short-Term Memory (LSTM)



Update cell state $\qquad C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$

What we should write

How much we should erase

How much we should write

https://colah.github.io/posts/2015-08-Understanding-LSTMs/

# Long Short-Term Memory (LSTM)



$$o_t = \sigma\big(W^{(o)}h_{t-1} + U^{(o)}x_t + b^{(o)}\big)$$

Update hidden state

$$h_t = o_t * \tanh(C_t)$$

# Long Short-Term Memory (LSTM)

Uninterrupted gradient flow



The addition is the key

LSTM does not guarantee that there is no vanishing gradient
but it does provide an easier way to learn long-distance dependencies

https://colah.github.io/posts/2015-08-Understanding-LSTMs/

# RNN vs. LSTM

RNN



LSTM



$$h_t = \sigma(W h_{t-1} + U x_t + b)$$

$$i_t = \sigma\big(W^{(i)} h_{t-1} + U^{(i)} x_t + b^{(i)}\big)$$

$$f_t = \sigma\big(W^{(f)} h_{t-1} + U^{(f)} x_t + b^{(f)}\big)$$

$$o_t = \sigma\big(W^{(o)} h_{t-1} + U^{(o)} x_t + b^{(o)}\big)$$

$$\tilde{C}_t = \tanh\big(W^{(C)} h_{t-1} + U^{(C)} x_t + b^{(C)}\big)$$

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

$$h_t = o_t * \tanh(C_t)$$

https://colah.github.io/posts/2015-08-Understanding-LSTMs/

# Gated Recurrent Units (GRU)

- Simplify 3 gates to 2 gates
  - Reset gate and update gate
- No explicit cell state
- More training-efficient

# Gated Recurrent Units (GRU)



Reset gate

$$r_t = \sigma\big(W^{(r)}h_{t-1} + U^{(r)}x_t + b^{(r)}\big)$$

Update gate

$$z_t = \tanh\big(W^{(z)}h_{t-1} + U^{(z)}x_t + b^{(z)}\big)$$

$$\tilde{h}_t = \tanh(W(r_t * h_{t-1}) + Ux_t + b)$$

New hidden state

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

Merge input and forget gate

# Multi-Layer RNN (Stacked RNN)

# Bidirectional RNN

# RNN is Flexible

- Can be used for both classification and generation
  - Encoder
  - Decoder
  - Encoder-decoder

# RNN as Sentence-Level Encoder



Sentence Embedding

$h_0$ $h_1$ $h_2$ $h_3$ $h_4$ $h_5$

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$

All   the   cats   are   cute

# RNN as Token-Level Encoder

Token Embeddings



The embeddings are contextualized

# Part-of-Speech (POS) Tagging

*You*    *can*    *close*    *the*    *door*

PRP    MD    VB    DT    NN

*It*    *is*    *close*    *to*    *the*    *door*

PRP    VBZ    JJ    IN    DT    NN

It's a structed prediction problem

# POS Tagging with Word Embeddings

# POS Tagging with Sequential Labeling

# Named Entity Recognition

John went to New York City to visit Kuan-Hao Huang

Entity        Entity        Entity

**BIO Sequence**

| John | went | to | New | York | City | to | visit | Kuan-Hao | Huang |
|------|------|------|------|------|------|------|------|------|------|
| B-Entity | Other | Other | B-Entity | I-Entity | I-Entity | Other | Other | B-Entity | I-Entity |

B-Entity: Begin of an entity span, I-Entity: Inside of an entity span

It's a structed prediction problem

# Named Entity Recognition as Sequential Labeling

# Sequential Labeling

- A sequence of dependent classification



$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^{N} \mathcal{L}_{CE}^{(i)}$$

# RNN as Decoder (Generator)

- RNN Language Modeling
  - Generation is a sequence of word classification

$$P(w_1, w_2, w_3, \ldots, w_l) = P(w_1)P(w_2, w_3, \ldots, w_l | w_1)$$



$P(w_i | w_{i-n+1}, \ldots, w_{i-1})$

Softmax

$w_{i-n+1} \quad w_{i-n+2} \qquad w_{i-2} \quad w_{i-1}$

Neural language models
with context window

$$= P(w_1)P(w_2|w_1)(w_3, \ldots, w_l | w_1, w_2)$$

$$= P(w_1)P(w_2|w_1)P(w_3|w_1, w_2)(w_4, \ldots, w_l | w_1, w_2, w_3)$$

$$= \prod_{i=1}^{l} P(w_i | w_1, w_2, \ldots, w_{i-1})$$

# RNN as Decoder (Generator)

- RNN Language Modeling
  - Generation is a sequence of word classification

# RNN as Decoder (Generator)

- RNN Language Modeling
  - Generation is a sequence of word classification



$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^{N} \mathcal{L}_{CE}^{(i)}$$

# Encoder vs. Decoder

- Encoder
  - Focus more on representations and understanding
- Decoder
  - Focus on generation

# Encoder



All    the    cats    are    cute

# Encoder

# Decoder

# Sequence-to-Sequence Models (Seq2Seq)

- When we need understanding and generation at the same time

# Sequence-to-Sequence Tasks

# Translation

- Translate English to Chinese

Classification over the whole vocabulary

我

$h_0$ $h_1$ $h_2$ $h_3$ $h_4$ $h_5$ $s_1$

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$ $y_1$

I    like    cats    a    lot    <bos>    Begin-of-Sentence Token

# Translation

- Translate English to Chinese



End-of-Sentence Token

我　很　喜　歡　貓　<eos>

$h_0$　$h_1$　$h_2$　$h_3$　$h_4$　$h_5$　$s_1$　$s_2$　$s_3$　$s_4$　$s_5$　$s_6$

$x_1$　$x_2$　$x_3$　$x_4$　$x_5$　$y_1$　$y_2$　$y_3$　$y_4$　$y_5$　$y_6$

I　like　cats　a　lot　<bos>　我　很　喜　歡　貓

# Sequence-to-Sequence Model Loss

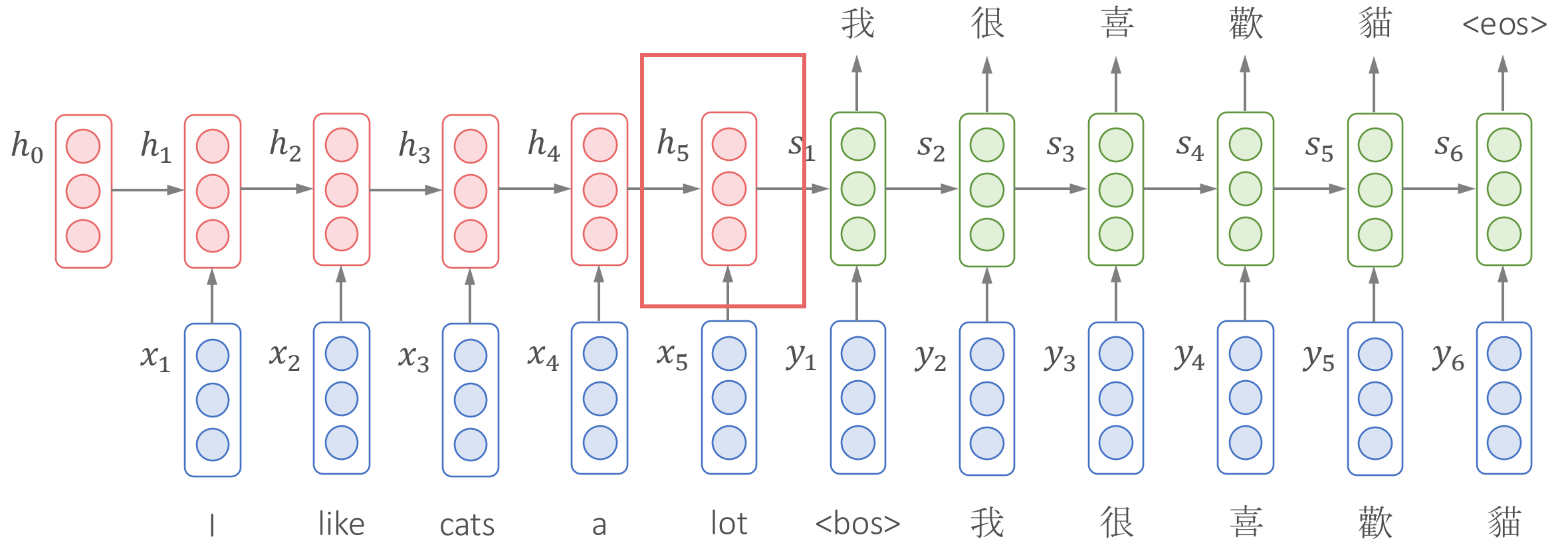$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^{N} \mathcal{L}_{CE}^{(i)}$$

# Decoder-Only Models vs. Seq2Seq Models

- Decoder-only models with prompting
  - Continue writing
- Seq2Seq models
  - Encode first, then generate
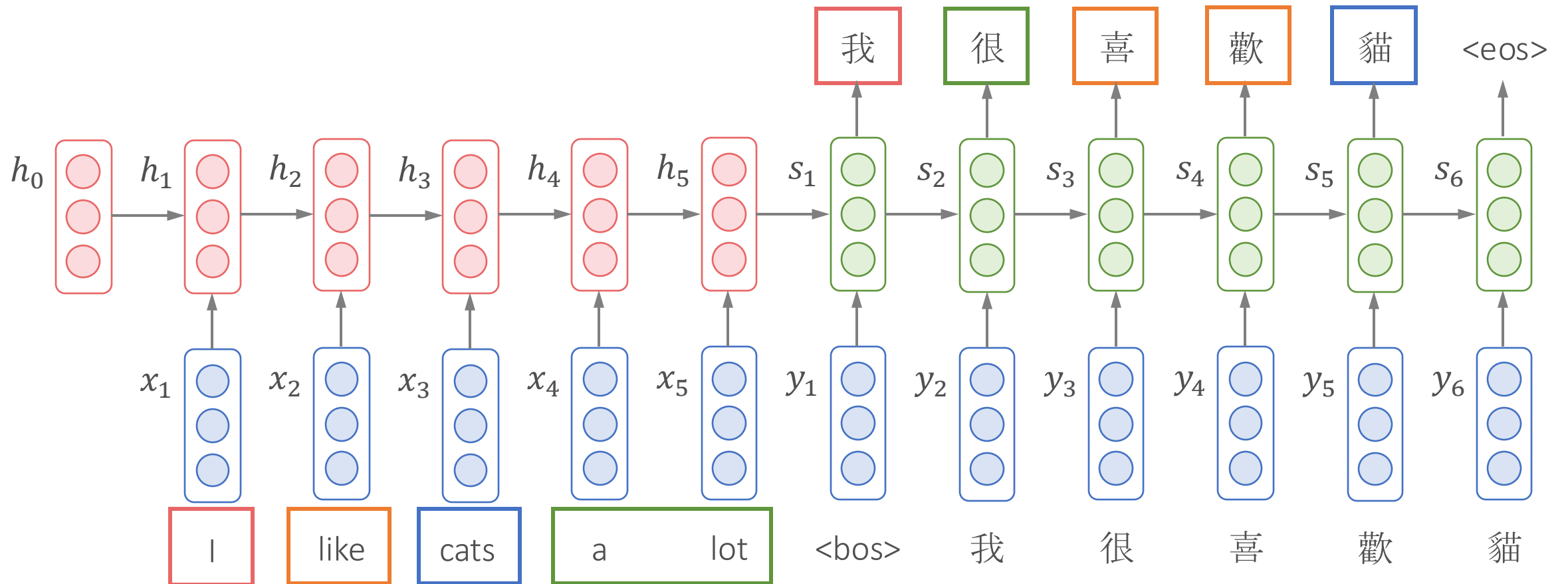- The difference becomes larger when we talk about Transformers!

# Seq2Seq: Bottleneck

- A single vector needs to capture all the information about source sentence
- Longer sequences can still lead to vanishing gradients

# Focus on A Particular Part When Decoding

- Each token classification requires different part of information from source sentence

# Next: Attention

- Attention provides a solution to the bottleneck problem
- Key idea: At each time step during decoding, focus on a particular part of source sentence