

# Recurrent Neural Networks and Transformers

Shan-Hung Wu

*shwu@cs.nthu.edu.tw*

Department of Computer Science,  
National Tsing Hua University, Taiwan

Machine Learning

# Outline

## 1 RNNs

- Vanilla RNNs
- Design Alternatives

## 2 RNN Training

- Backprop through Time (BPTT)
- Optimization Techniques
- Optimization-Friendly Models & LSTM
- Parallelism & Teacher Forcing

## 3 RNNs with Attention Mechanism

- Attention for Image Captioning
- Attention for Neural Machine Translation (NMT)

## 4 Transformers

- Attention is All You Need
- Pretrained Language Models
- More Applications

## 5 Subword Tokenization

# Outline

## ① RNNs

- Vanilla RNNs
- Design Alternatives

## ② RNN Training

- Backprop through Time (BPTT)
- Optimization Techniques
- Optimization-Friendly Models & LSTM
- Parallelism & Teacher Forcing

## ③ RNNs with Attention Mechanism

- Attention for Image Captioning
- Attention for Neural Machine Translation (NMT)

## ④ Transformers

- Attention is All You Need
- Pretrained Language Models
- More Applications

## ⑤ Subword Tokenization

# Outline

## ① RNNs

- Vanilla RNNs
- Design Alternatives

## ② RNN Training

- Backprop through Time (BPTT)
- Optimization Techniques
- Optimization-Friendly Models & LSTM
- Parallelism & Teacher Forcing

## ③ RNNs with Attention Mechanism

- Attention for Image Captioning
- Attention for Neural Machine Translation (NMT)

## ④ Transformers

- Attention is All You Need
- Pretrained Language Models
- More Applications

## ⑤ Subword Tokenization

# Sequential Data

- So far, we assume that data points  $(x, y)$ 's in a dataset are i.i.d

# Sequential Data

- So far, we assume that data points  $(x,y)$ 's in a dataset are i.i.d
- Does **not** hold in many applications

# Sequential Data

- So far, we assume that data points  $(x,y)$ 's in a dataset are i.i.d
- Does **not** hold in many applications
- **Sequential data**: data points come in order and successive points may be dependent, e.g.,
  - Letters in a word
  - Words in a sentence/document
  - Phonemes in a spoken word utterance
  - Page clicks in a Web session
  - Frames in a video, etc.

# Sequential Data

- So far, we assume that data points  $(x,y)$ 's in a dataset are i.i.d
- Does **not** hold in many applications
- **Sequential data**: data points come in order and successive points may be dependent, e.g.,
  - Letters in a word
  - Words in a sentence/document
  - Phonemes in a spoken word utterance
  - Page clicks in a Web session
  - Frames in a video, etc.
- Dataset:  $\mathbf{X} = \{\mathbf{X}^{(n)}\}_n \in \mathbb{R}^{N \times (D,K) \times T}$

# Sequential Data

- So far, we assume that data points  $(x, y)$ 's in a dataset are i.i.d
- Does **not** hold in many applications
- **Sequential data**: data points come in order and successive points may be dependent, e.g.,
  - Letters in a word
  - Words in a sentence/document
  - Phonemes in a spoken word utterance
  - Page clicks in a Web session
  - Frames in a video, etc.
- Dataset:  $\mathbf{X} = \{X^{(n)}\}_n \in \mathbb{R}^{N \times (D, K) \times T}$ 
  - $X^{(n)} = \{(x^{(n,t)}, y^{(n,t)})\}_t$  a **sequence**, where the superscript  $n$  can be omitted for simplicity

# Sequential Data

- So far, we assume that data points  $(x, y)$ 's in a dataset are i.i.d
- Does **not** hold in many applications
- **Sequential data:** data points come in order and successive points may be dependent, e.g.,
  - Letters in a word
  - Words in a sentence/document
  - Phonemes in a spoken word utterance
  - Page clicks in a Web session
  - Frames in a video, etc.
- Dataset:  $\mathbf{X} = \{X^{(n)}\}_n \in \mathbb{R}^{N \times (D, K) \times T}$ 
  - $X^{(n)} = \{(x^{(n,t)}, y^{(n,t)})\}_t$  a **sequence**, where the superscript  $n$  can be omitted for simplicity
  - $T$  is called the **horizon** and may be different between  $x^{(n)}$  and  $y^{(n)}$  and across data points  $n$ 's

# Sequence Modeling I

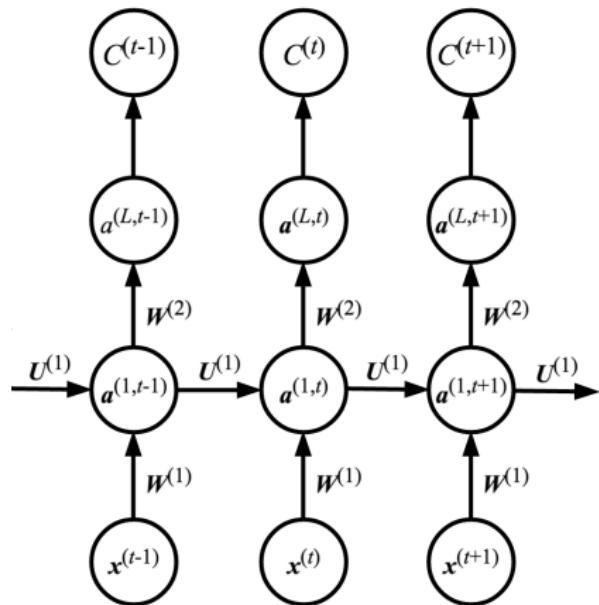
- How to model sequential data?

# Sequence Modeling I

- How to model sequential data?

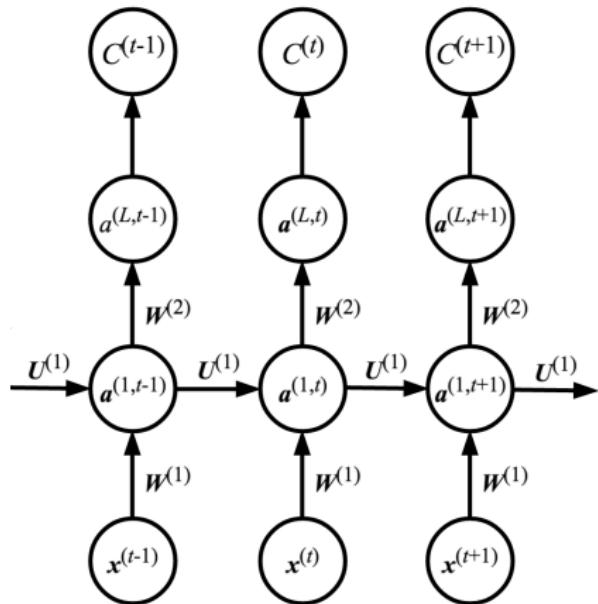
- **Recurrent neural networks**

(vanilla RNNs):



# Sequence Modeling I

- How to model sequential data?
- **Recurrent neural networks**  
(vanilla RNNs):
- $y^{(t)}$  depends on  $x^{(1)}, \dots, x^{(t)}$

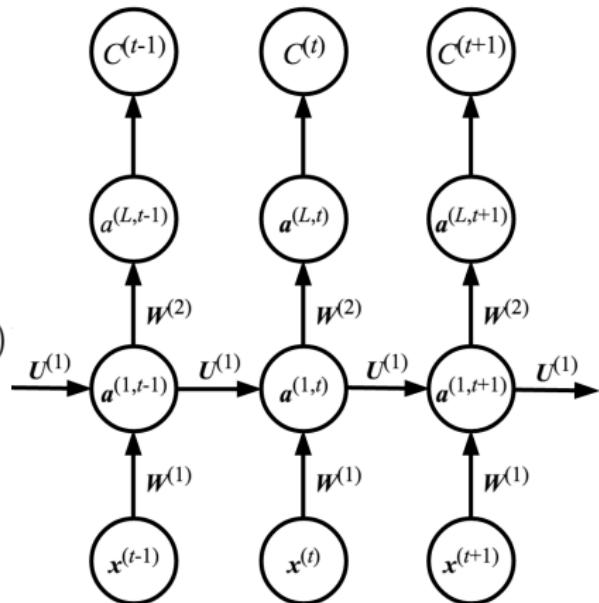


# Sequence Modeling I

- How to model sequential data?
- **Recurrent neural networks**  
(vanilla RNNs):
- $y^{(t)}$  depends on  $x^{(1)}, \dots, x^{(t)}$
- Output  $a^{(L,t)}$  depends on hidden activations:

$$\begin{aligned}\mathbf{a}^{(k,t)} &= \text{act}(\mathbf{z}^{(k,t)}) \\ &= \text{act}(\mathbf{U}^{(k)} \mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)} \mathbf{a}^{(k-1,t)})\end{aligned}$$

- Bias term omitted

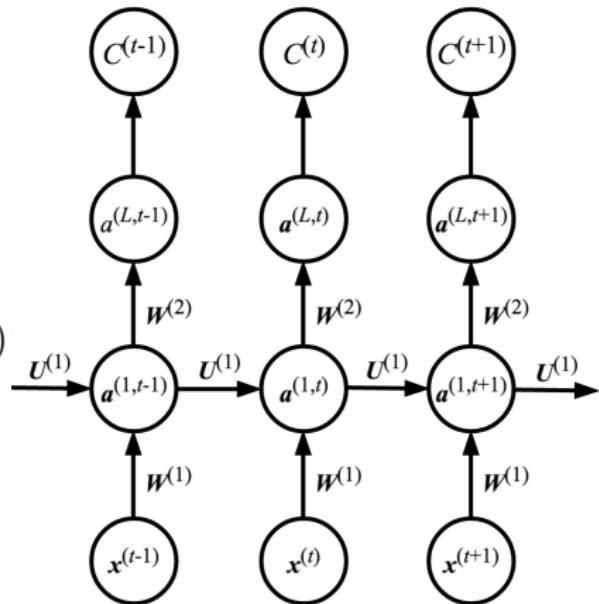


# Sequence Modeling I

- How to model sequential data?
- **Recurrent neural networks**  
(vanilla RNNs):
- $\mathbf{y}^{(t)}$  depends on  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}$
- Output  $\mathbf{a}^{(L,t)}$  depends on hidden activations:

$$\begin{aligned}\mathbf{a}^{(k,t)} &= \text{act}(\mathbf{z}^{(k,t)}) \\ &= \text{act}(\mathbf{U}^{(k)} \mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)} \mathbf{a}^{(k-1,t)})\end{aligned}$$

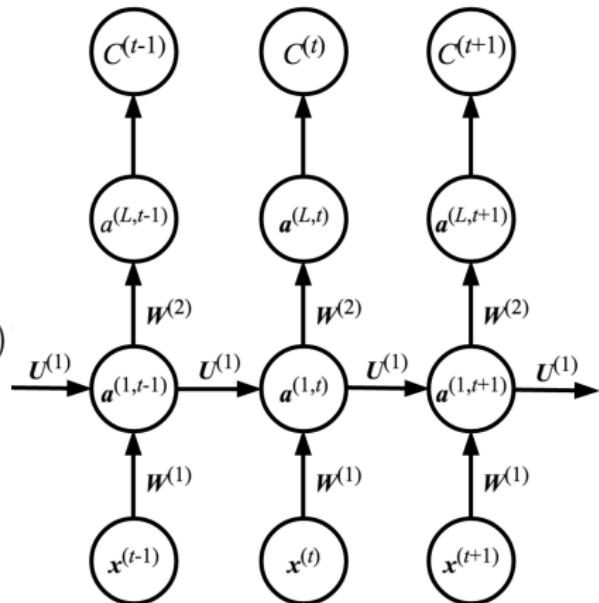
- Bias term omitted
- $\mathbf{a}^{(\cdot,t)}$  summarizes  $\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}$ 
  - Earlier points are less important



# Sequence Modeling I

- How to model sequential data?
- **Recurrent neural networks**  
(vanilla RNNs):
- $\mathbf{y}^{(t)}$  depends on  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}$
- Output  $\mathbf{a}^{(L,t)}$  depends on hidden activations:

$$\begin{aligned}\mathbf{a}^{(k,t)} &= \text{act}(\mathbf{z}^{(k,t)}) \\ &= \text{act}(\mathbf{U}^{(k)} \mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)} \mathbf{a}^{(k-1,t)})\end{aligned}$$

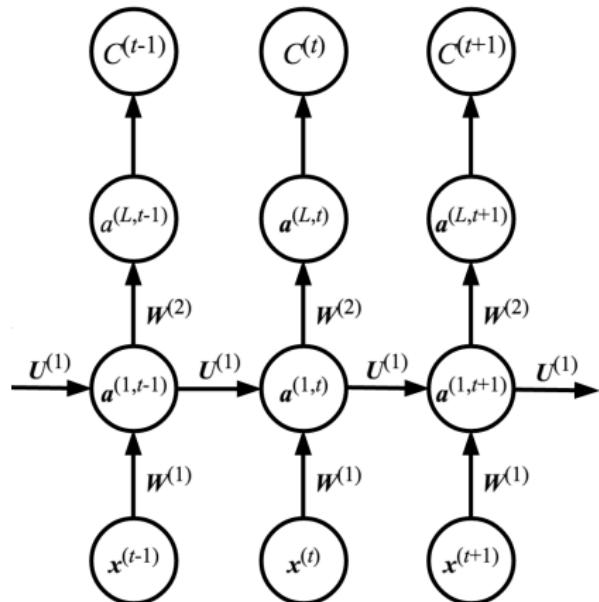


- Bias term omitted
- $\mathbf{a}^{(\cdot,t)}$  summarizes  $\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}$ 
  - Earlier points are less important
- $\mathbf{a}^{(\cdot,t)}$ 's at deeper layers give more abstract summarizations

# Sequence Modeling II

$$\begin{aligned}\mathbf{a}^{(k,t)} &= \text{act}(\mathbf{z}^{(k,t)}) \\ &= \text{act}(\mathbf{U}^{(k)} \mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)} \mathbf{a}^{(k-1,t)})\end{aligned}$$

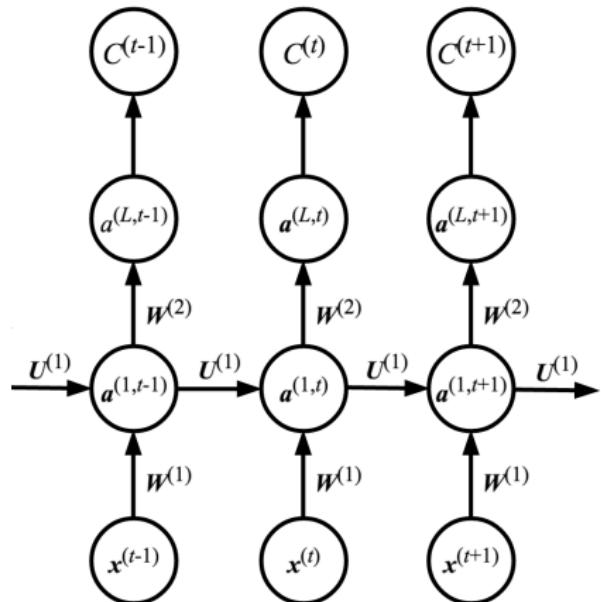
- Weights are *shared* across time instances



# Sequence Modeling II

$$\begin{aligned}\mathbf{a}^{(k,t)} &= \text{act}(\mathbf{z}^{(k,t)}) \\ &= \text{act}(\mathbf{U}^{(k)} \mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)} \mathbf{a}^{(k-1,t)})\end{aligned}$$

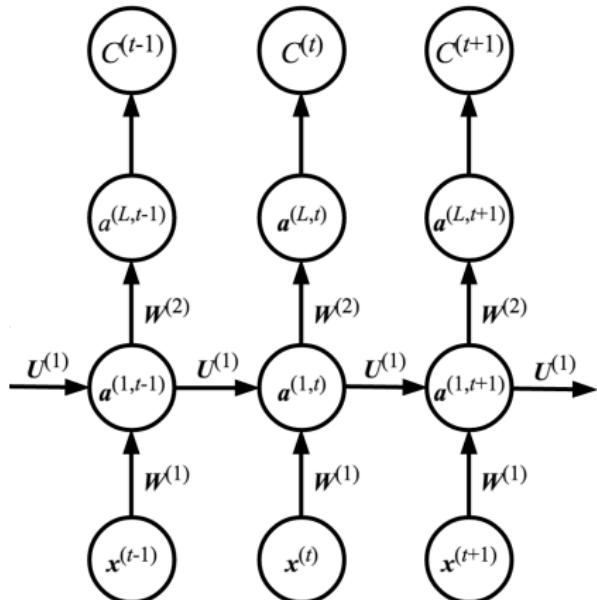
- Weights are *shared* across time instances
- Assumes that the “transition functions” are time invariant



# Sequence Modeling II

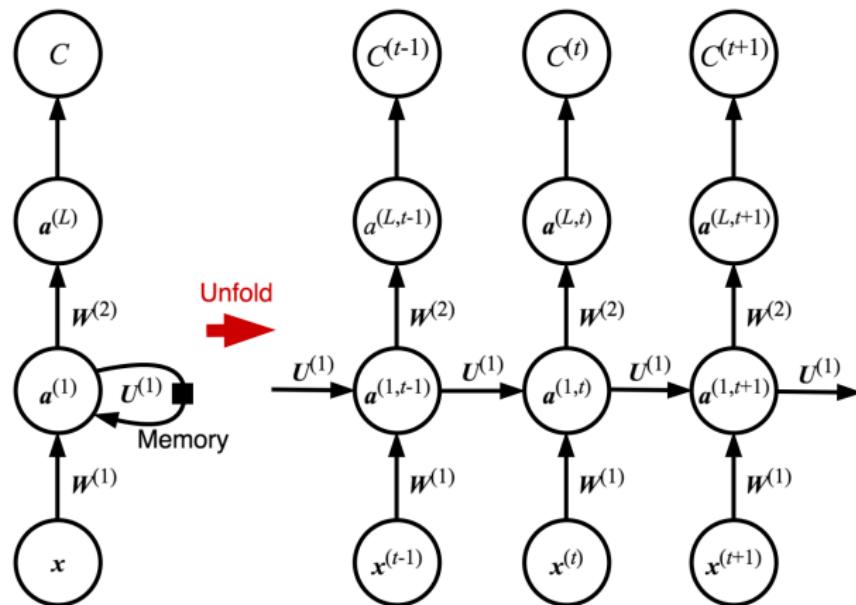
$$\begin{aligned}\mathbf{a}^{(k,t)} &= \text{act}(\mathbf{z}^{(k,t)}) \\ &= \text{act}(\mathbf{U}^{(k)} \mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)} \mathbf{a}^{(k-1,t)})\end{aligned}$$

- Weights are *shared* across time instances
- Assumes that the “transition functions” are time invariant
- Our goal is to learn  $\mathbf{U}^{(k)}$ ’s and  $\mathbf{W}^{(k)}$ ’s for  $k = 1, \dots, L$



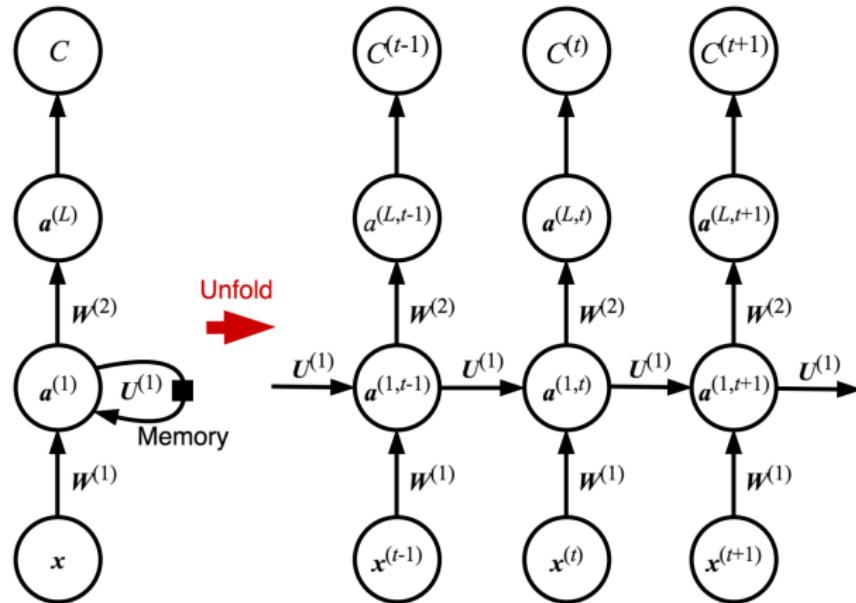
# RNNs have Memory

- The computational graph of an RNN can be *folded* in time



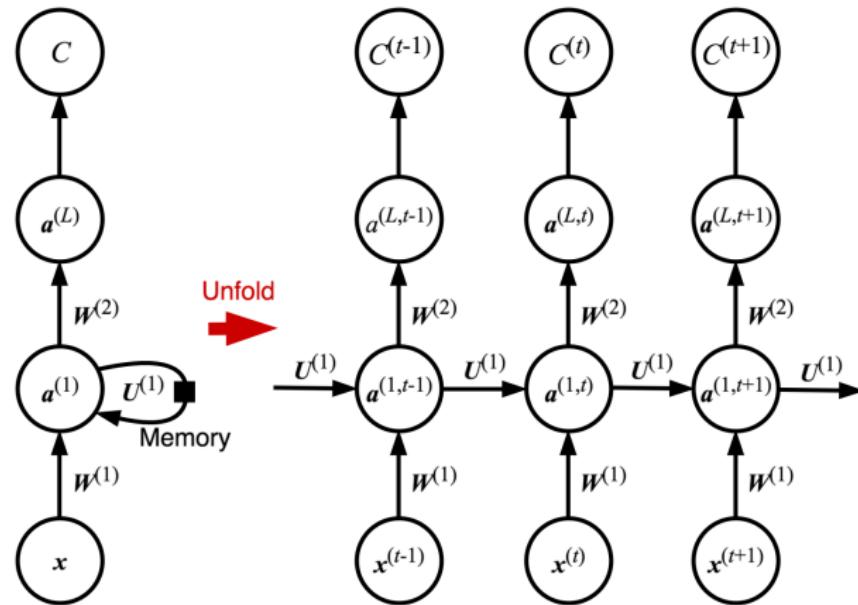
# RNNs have Memory

- The computational graph of an RNN can be *folded* in time
- Black squares denotes *memory* access



# Output Layer (1/2)

- With multi-class  $\mathbf{y}^{(t)}$ ,
  - $\mathbf{a}^{(L,t)}$  represents the probability of each class
  - $C^{(t)}$  is cross entropy
- How to obtain  $\hat{\mathbf{y}}^{(t)}$  from  $\mathbf{a}^{(L,t)}$  at inference time?



# Output Layer (2/2)

- ***Output sampling*** for multi-class tasks:
  - Greedy: sample  $\hat{\mathbf{y}}^{(t)}$  from  $\mathbf{a}^{(L,t)}$
  - Bean search: sample  $\hat{\mathbf{y}}^{(t)}$  from the most probable paths of the join distribution  $(\mathbf{a}^{(L,t)}, \mathbf{a}^{(L,t-1)}, \dots, \mathbf{a}^{(L,t-b)})$ , where  $b$  is bean size
  - Noisy: sample  $\hat{\mathbf{y}}^{(t)}$  based on  $\mathbf{a}^{(L,t)} + \text{noise}$

# Output Layer (2/2)

- **Output sampling** for multi-class tasks:
  - Greedy: sample  $\hat{y}^{(t)}$  from  $\mathbf{a}^{(L,t)}$
  - Bean search: sample  $\hat{y}^{(t)}$  from the most probable paths of the join distribution  $(\mathbf{a}^{(L,t)}, \mathbf{a}^{(L,t-1)}, \dots, \mathbf{a}^{(L,t-b)})$ , where  $b$  is bean size
  - Noisy: sample  $\hat{y}^{(t)}$  based on  $\mathbf{a}^{(L,t)} + \text{noise}$
- Problem: out of vocabulary or high dimensional  $\mathbf{a}^{(L,t)}$

The image shows the word "HELLO" in large blue capital letters at the center. Surrounding it are various international greetings in different colors and fonts:

- Top left: 你好 (red)
- Top middle: HALLO (green)
- Top right: 안녕 (purple)
- Middle left: CIAO (yellow)
- Middle center: HOLA (orange)
- Middle right: নমস্তে (brown)
- Bottom left: こんにちは (light green)
- Bottom center: привет (pink)
- Bottom right: OLÁ (red)
- Bottom left: BONJOUR (purple)
- Bottom center: مرحبا (blue)

- Solution?

# Output Layer (2/2)

- **Output sampling** for multi-class tasks:
  - Greedy: sample  $\hat{y}^{(t)}$  from  $\mathbf{a}^{(L,t)}$
  - Bean search: sample  $\hat{y}^{(t)}$  from the most probable paths of the join distribution  $(\mathbf{a}^{(L,t)}, \mathbf{a}^{(L,t-1)}, \dots, \mathbf{a}^{(L,t-b)})$ , where  $b$  is bean size
  - Noisy: sample  $\hat{y}^{(t)}$  based on  $\mathbf{a}^{(L,t)} + \text{noise}$
- Problem: out of vocabulary or high dimensional  $\mathbf{a}^{(L,t)}$

你好 HALLO 안녕  
HOLA نমস্তে  
CIAO  
ঘোড়া  
こんにちは  
Bonjour مرحبا  
привет  
OLÁ

- Solution? **Subword tokenization** (to be discussed later)

# RNNs vs CNNs for Sequential Data

- On processing a sequence of length  $T$  at each layer with
  - $D$ -dimensional point input and output
  - $F$  = the CNN filter/kernel size
  - #CNN filters =  $D$

	#Weights	Computation	Autoregressive	Point Distance
CNN	$O(FD^2)$	$O(TFD^2)$	No	$O(\frac{T}{F})$
RNN	$O(D^2)$	$O(TD^2)$	Yes	$O(T)$

# Outline

## ① RNNs

- Vanilla RNNs
- Design Alternatives

## ② RNN Training

- Backprop through Time (BPTT)
- Optimization Techniques
- Optimization-Friendly Models & LSTM
- Parallelism & Teacher Forcing

## ③ RNNs with Attention Mechanism

- Attention for Image Captioning
- Attention for Neural Machine Translation (NMT)

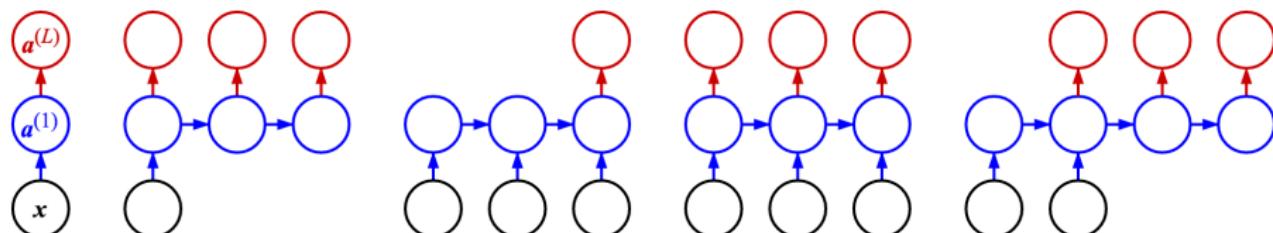
## ④ Transformers

- Attention is All You Need
- Pretrained Language Models
- More Applications

## ⑤ Subword Tokenization

# Input and Output

- $x^{(t)}$ 's and  $y^{(t)}$ 's do **not** need to have one-to-one correspondence:



NN

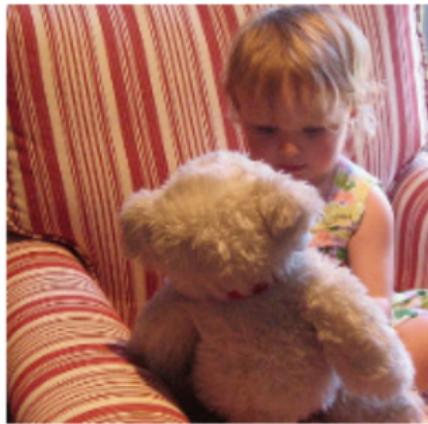
One to Many

Many to One

Many to Many  
(Synced)

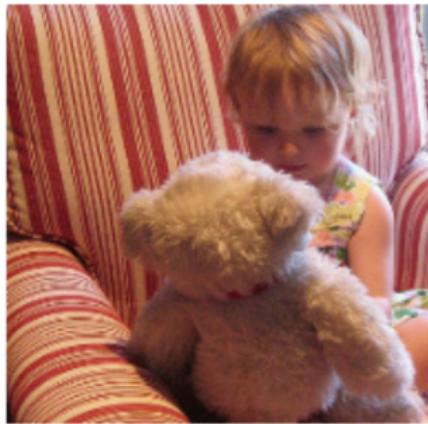
Many to Many  
(Unsynced)

# One2Many: Image Captioning

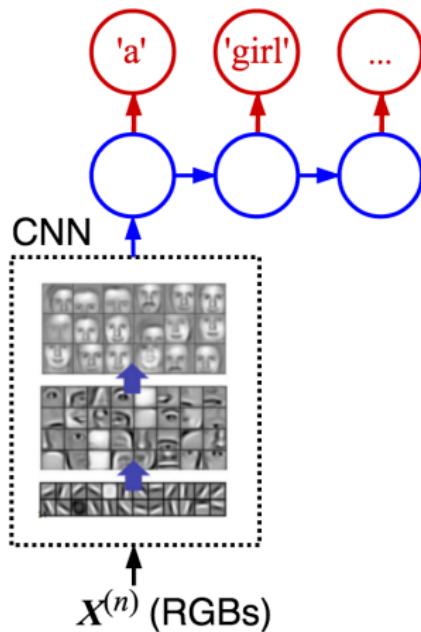


“A little girl sitting on a bed with a  
teddy bear.”

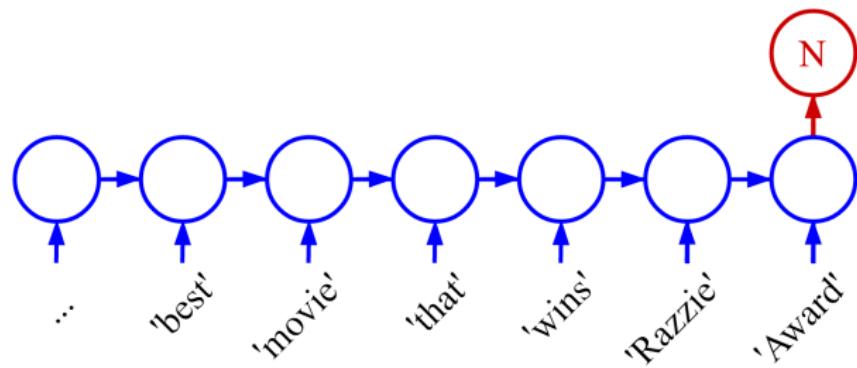
# One2Many: Image Captioning



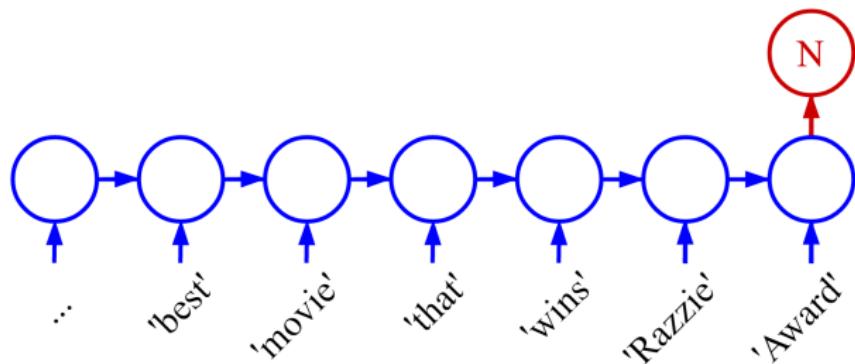
"A little girl sitting on a bed with a teddy bear."



# Many2One: Sentiment Analysis



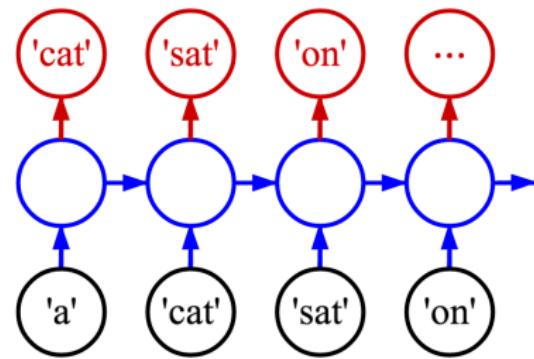
# Many2One: Sentiment Analysis



- A single word (e.g., "Razzie") can negate the entire input sentence

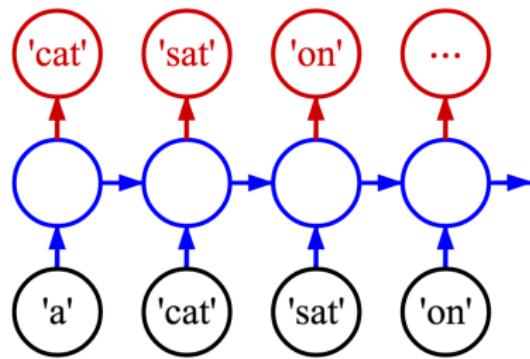
# Many2Many (Synced): Language Modeling

- **Language modeling**: predicting the next/nearby word based on the context



# Many2Many (Synced): Language Modeling

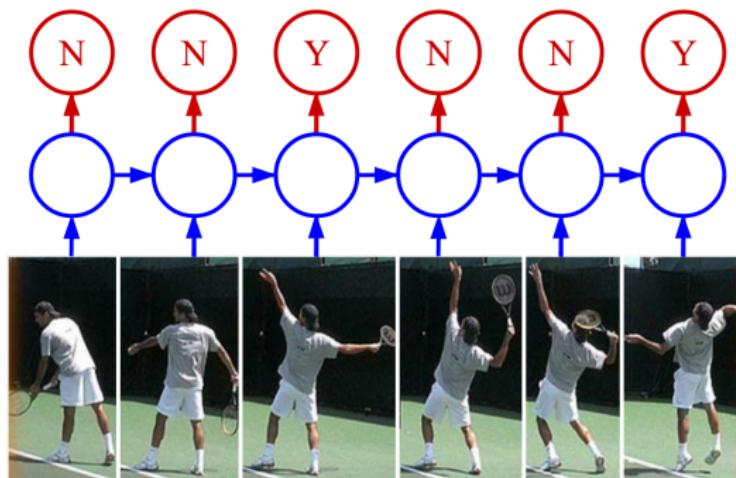
- **Language modeling**: predicting the next/nearby word based on the context



- Latent representations of RNN provide the context

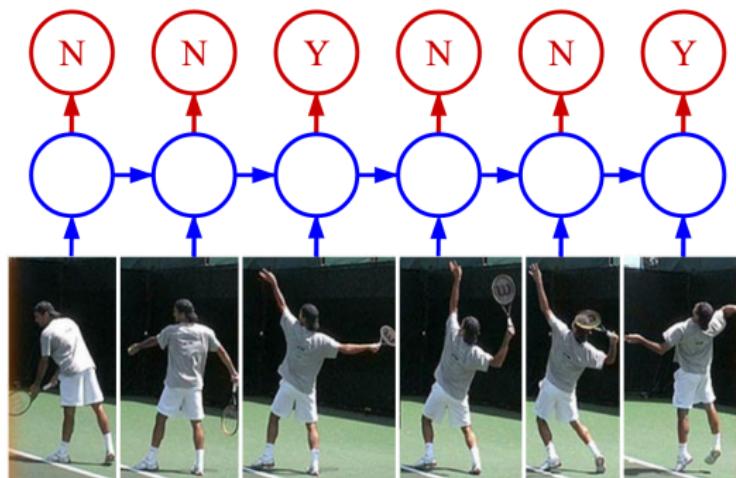
# Many2Many (Synced): Video Keyframe Tagging

- Video frame annotation:



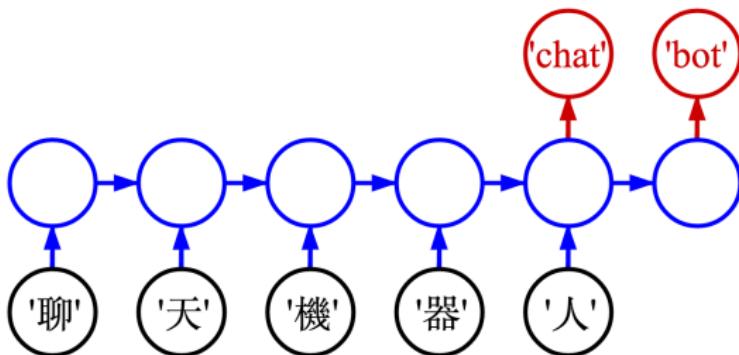
# Many2Many (Synced): Video Keyframe Tagging

- Video frame annotation:



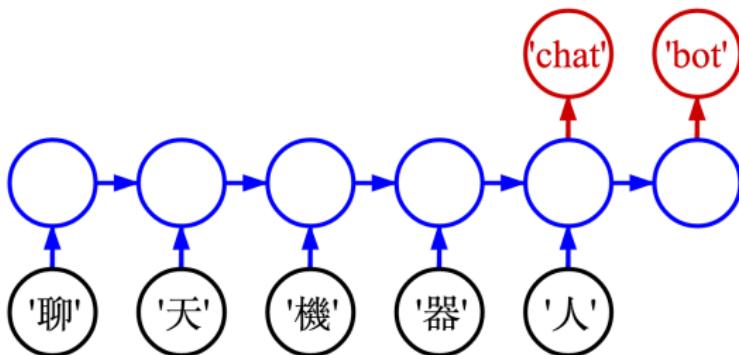
- Latent representations summarize “what’s going on”

# Many2Many (Unsynced): Machine Translation



- Latent representations support *encoding* first, and then *decoding*
- RNN learns the structure difference

# Many2Many (Unsynced): Machine Translation



- Latent representations support *encoding* first, and then *decoding*
- RNN learns the structure difference
- Also called *sequence to sequence* learning
  - Also used in other applications, e.g., chat bots

# Bidirectional RNNs

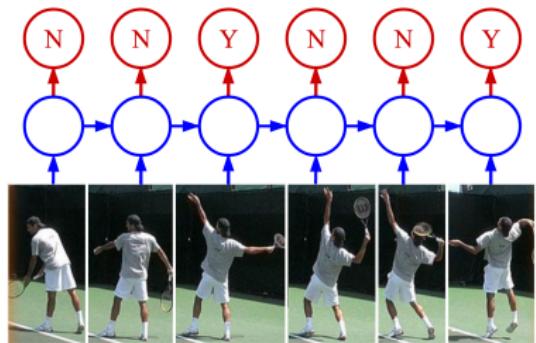
$$\mathbf{a}^{(k,t)} = \text{act}(\mathbf{U}^{(k)} \mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)} \mathbf{a}^{(k-1,t)})$$

- Each  $\mathbf{a}^{(\cdot,t)}$  summarizes  
 $\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}$

# Bidirectional RNNs

$$\mathbf{a}^{(k,t)} = \text{act}(\mathbf{U}^{(k)} \mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)} \mathbf{a}^{(k-1,t)})$$

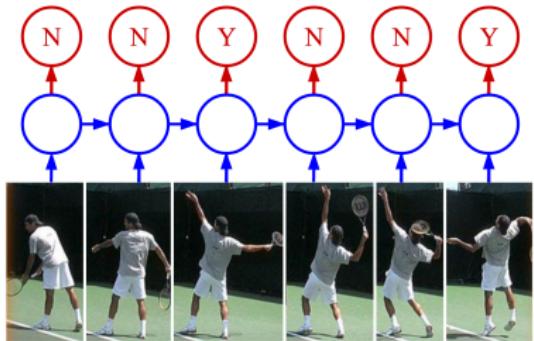
- Each  $\mathbf{a}^{(\cdot,t)}$  summarizes  $\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}$
- A prediction is made by considering previous frames



# Bidirectional RNNs

$$\mathbf{a}^{(k,t)} = \text{act}(\mathbf{U}^{(k)} \mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)} \mathbf{a}^{(k-1,t)})$$

- Each  $\mathbf{a}^{(\cdot,t)}$  summarizes  $\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}$
- A prediction is made by considering previous frames
- Can we take **future** frames into account?

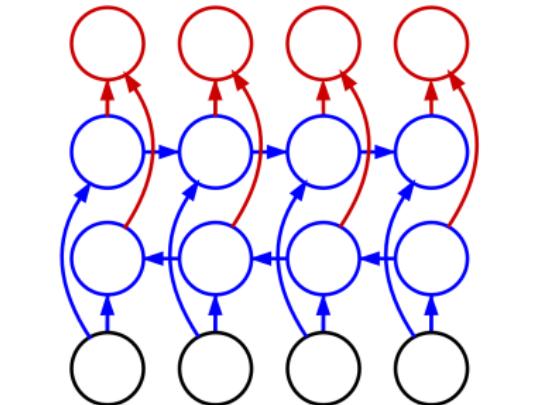
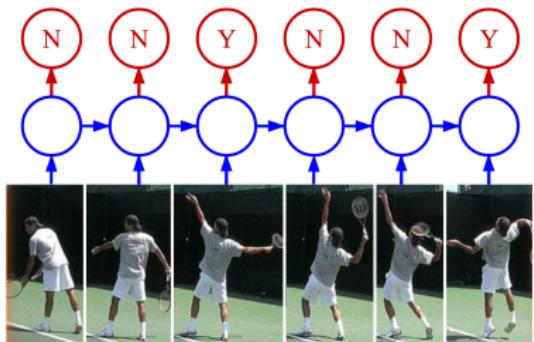


# Bidirectional RNNs

$$\mathbf{a}^{(k,t)} = \text{act}(\mathbf{U}^{(k)} \mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)} \mathbf{a}^{(k-1,t)})$$

- Each  $\mathbf{a}^{(\cdot,t)}$  summarizes  $\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}$
- A prediction is made by considering previous frames
- Can we take *future* frames into account?
- **Bidirectional RNNs**: output  $\mathbf{a}^{(L,t)}$  depends on both  $\mathbf{a}^{(k,t)}$ 's and  $\tilde{\mathbf{a}}^{(k,t)}$ 's

$$\tilde{\mathbf{a}}^{(k,t)} = \text{act}(\tilde{\mathbf{U}}^{(k)} \tilde{\mathbf{a}}^{(k,t+1)} + \tilde{\mathbf{W}}^{(k)} \tilde{\mathbf{a}}^{(k-1,t)})$$



# Recursive RNNs I

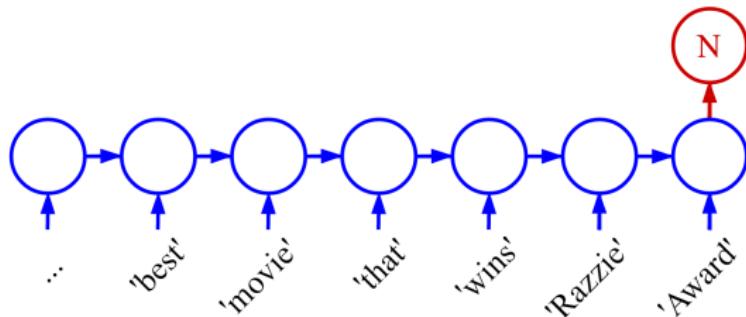
$$\mathbf{a}^{(k,t)} = \text{act}(\mathbf{U}^{(k)}\mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)}\mathbf{a}^{(k-1,t)})$$

- The transition of hidden representation is invariant in time
  - Earlier input/representation is less important to current  $\mathbf{a}^{(\cdot,t)}$

# Recursive RNNs I

$$\mathbf{a}^{(k,t)} = \text{act}(\mathbf{U}^{(k)} \mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)} \mathbf{a}^{(k-1,t)})$$

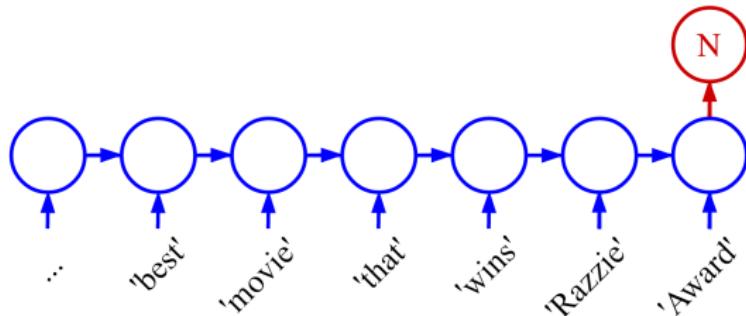
- The transition of hidden representation is invariant in time
  - Earlier input/representation is less important to current  $\mathbf{a}^{(\cdot,t)}$
- “Razzie” has less effect if it is far away from the prediction



# Recursive RNNs I

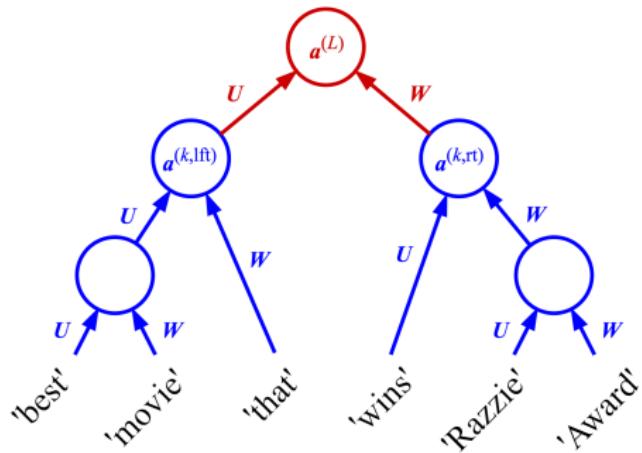
$$\mathbf{a}^{(k,t)} = \text{act}(\mathbf{U}^{(k)}\mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)}\mathbf{a}^{(k-1,t)})$$

- The transition of hidden representation is invariant in time
  - Earlier input/representation is less important to current  $\mathbf{a}^{(\cdot,t)}$
- “Razzie” has less effect if it is far away from the prediction
- In some applications, transitions are invariant in terms of other concepts



# Recursive RNNs II

- In natural language processing (NLP), we can parse the input sentence  $X^{(n)}$  into a tree
  - Following grammar rules

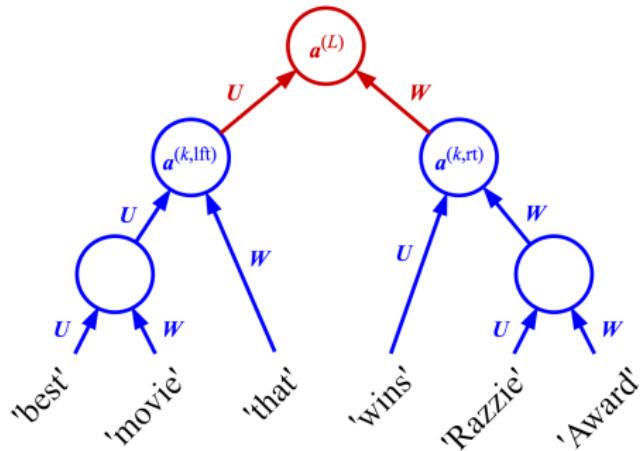


# Recursive RNNs II

- In natural language processing (NLP), we can parse the input sentence  $X^{(n)}$  into a tree
  - Following grammar rules
- **Recursive RNNs:** “subtree merges” are invariant

$$\mathbf{a}^{(k,t)} = \text{act}(\mathbf{U}\mathbf{a}^{(k-1, \text{ left })} + \mathbf{W}\mathbf{a}^{(k-1, \text{ right })})$$

- $\mathbf{U}$  and  $\mathbf{W}$  are shared recursively in subtrees

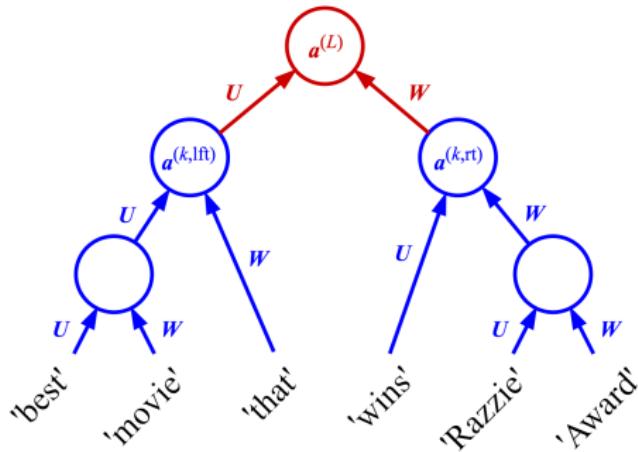


# Recursive RNNs II

- In natural language processing (NLP), we can parse the input sentence  $X^{(n)}$  into a tree
  - Following grammar rules
- **Recursive RNNs:** “subtree merges” are invariant

$$\mathbf{a}^{(k,t)} = \text{act}(\mathbf{U}\mathbf{a}^{(k-1, \text{ left })} + \mathbf{W}\mathbf{a}^{(k-1, \text{ right })})$$

- $\mathbf{U}$  and  $\mathbf{W}$  are shared recursively in subtrees
- Given sentence length  $T$ ,  $\mathbf{a}^{(L)}$  and  $\mathbf{a}^{(1,\cdot)}$  can be  $O(\log T)$  away in the best case



# Outline

## 1 RNNs

- Vanilla RNNs
- Design Alternatives

## 2 RNN Training

- Backprop through Time (BPTT)
- Optimization Techniques
- Optimization-Friendly Models & LSTM
- Parallelism & Teacher Forcing

## 3 RNNs with Attention Mechanism

- Attention for Image Captioning
- Attention for Neural Machine Translation (NMT)

## 4 Transformers

- Attention is All You Need
- Pretrained Language Models
- More Applications

## 5 Subword Tokenization

# Cost Function of Vanilla RNNs

- Parameters to learn:  $\Theta = \{\mathbf{W}^{(k)}, \mathbf{U}^{(k)}\}_k$  (bias terms omitted)
- Maximum likelihood:

$$\begin{aligned}& \arg \min_{\Theta} C(\Theta) \\&= \arg \min_{\Theta} -\log P(\mathbf{X} | \Theta) \\&= \arg \min_{\Theta} -\sum_{n,t} \log P(y^{(n,t)} | \mathbf{x}^{(n,t)}, \dots, \mathbf{x}^{(n,1)}, \Theta) \\&= \arg \min_{\Theta} -\sum_{n,t} C^{(n,t)}(\Theta)\end{aligned}$$

- $y^{(t)}$  depends only on  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}$

# Cost Function of Vanilla RNNs

- Parameters to learn:  $\Theta = \{\mathbf{W}^{(k)}, \mathbf{U}^{(k)}\}_k$  (bias terms omitted)
- Maximum likelihood:

$$\begin{aligned}& \arg \min_{\Theta} C(\Theta) \\&= \arg \min_{\Theta} -\log P(\mathbf{X} | \Theta) \\&= \arg \min_{\Theta} -\sum_{n,t} \log P(\mathbf{y}^{(n,t)} | \mathbf{x}^{(n,t)}, \dots, \mathbf{x}^{(n,1)}, \Theta) \\&= \arg \min_{\Theta} -\sum_{n,t} C^{(n,t)}(\Theta)\end{aligned}$$

- $\mathbf{y}^{(t)}$  depends only on  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}$
- For example, in binary classification:
- Assuming  $P(\mathbf{y}^{(n,t)} = 1 | \mathbf{x}^{(n,t)}, \dots, \mathbf{x}^{(n,1)}) \sim \text{Bernoulli}(\rho^{(t)})$ , we have

$$C^{(n,t)}(\Theta) = (\mathbf{a}^{(L,t)})^{\mathbf{y}^{(n,t)}} (1 - \mathbf{a}^{(L,t)})^{(1-\mathbf{y}^{(n,t)})}$$

- $\mathbf{a}^{(L,t)} = \rho^{(t)}$  are based on  $\mathbf{a}^{(\cdot,t)}$ 's, which summarize  $\mathbf{x}^{(n,t)}, \dots, \mathbf{x}^{(n,1)}$

# Outline

## 1 RNNs

- Vanilla RNNs
- Design Alternatives

## 2 RNN Training

- Backprop through Time (BPTT)
- Optimization Techniques
- Optimization-Friendly Models & LSTM
- Parallelism & Teacher Forcing

## 3 RNNs with Attention Mechanism

- Attention for Image Captioning
- Attention for Neural Machine Translation (NMT)

## 4 Transformers

- Attention is All You Need
- Pretrained Language Models
- More Applications

## 5 Subword Tokenization

# SGD-based Training

- RNN optimization problem can be solved using SGD:

$$\Theta^{(s+1)} \leftarrow \Theta^{(s)} - \eta \nabla_{\Theta} \sum_{n,t} C^{(n,t)}(\Theta^{(s)})$$

- Let  $c^{(n,t)} = C^{(n,t)}(\Theta^{(s)})$ , our goal is to evaluate  $\frac{\partial c^{(n,t)}}{\partial U_{i,j}^{(k)}}$  and  $\frac{\partial c^{(n,t)}}{\partial W_{i,j}^{(k)}}$
- Evaluation of  $\frac{\partial c^{(n,t)}}{\partial W_{i,j}^{(k)}}$  is similar to that in DNNs and omitted
- We focus on:

$$\frac{\partial c^{(n,t)}}{\partial U_{i,j}^{(k)}} = \frac{\partial c^{(n,t)}}{\partial z_j^{(k,t)}} \cdot \frac{\partial z_j^{(k,t)}}{\partial U_{i,j}^{(k)}} = \delta_j^{(k,t)} \frac{\partial z_j^{(k,t)}}{\partial U_{i,j}^{(k)}}$$

# Forward Pass through Time

- The second term:  $\frac{\partial z_j^{(k,t)}}{\partial U_{i,j}^{(k)}}$

# Forward Pass through Time

- The second term:  $\frac{\partial z_j^{(k,t)}}{\partial U_{i,j}^{(k)}}$
- We have  $z_j^{(k,t)} = \sum_i W_{i,j}^{(k)} a_i^{(k-1,t)} + \sum_i U_{i,j}^{(k)} a_i^{(k,t-1)}$  and

$$\frac{\partial z_j^{(k,t)}}{\partial U_{i,j}^{(k)}} = a_i^{(k,t-1)}$$

# Forward Pass through Time

- The second term:  $\frac{\partial z_j^{(k,t)}}{\partial U_{i,j}^{(k)}}$
- We have  $z_j^{(k,t)} = \sum_i W_{i,j}^{(k)} a_i^{(k-1,t)} + \sum_i U_{i,j}^{(k)} a_i^{(k,t-1)}$  and

$$\frac{\partial z_j^{(k,t)}}{\partial U_{i,j}^{(k)}} = a_i^{(k,t-1)}$$

- We can get all second terms starting from the most shallow layer and **earliest time**

# Backward Pass through Time I

- The first term (error signal):  $\delta_j^{(k,t)} = \frac{\partial c^{(n,t)}}{\partial z_j^{(k,t)}}$
- We have

$$\delta_j^{(k,t)} = \frac{\partial c^{(n,t)}}{\partial z_j^{(k,t)}} = \frac{\partial c^{(n,t)}}{\partial a_j^{(k,t)}} \cdot \frac{\partial a_j^{(k,t)}}{\partial z_j^{(k,t)}}$$

# Backward Pass through Time I

- The first term (error signal):  $\delta_j^{(k,t)} = \frac{\partial c^{(n,t)}}{\partial z_j^{(k,t)}}$

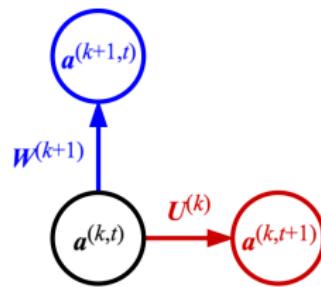
- We have

$$\delta_j^{(k,t)} = \frac{\partial c^{(n,t)}}{\partial z_j^{(k,t)}} = \frac{\partial c^{(n,t)}}{\partial a_j^{(k,t)}} \cdot \frac{\partial a_j^{(k,t)}}{\partial z_j^{(k,t)}} = \frac{\partial c^{(n,t)}}{\partial a_j^{(k,t)}} \cdot \text{act}'(z_j^{(k,t)})$$

# Backward Pass through Time I

- The first term (error signal):  $\delta_j^{(k,t)} = \frac{\partial c^{(n,t)}}{\partial z_j^{(k,t)}}$
- We have

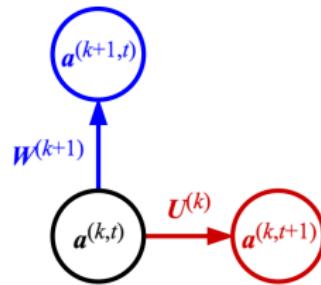
$$\begin{aligned}\delta_j^{(k,t)} &= \frac{\partial c^{(n,t)}}{\partial z_j^{(k,t)}} = \frac{\partial c^{(n,t)}}{\partial a_j^{(k,t)}} \cdot \frac{\partial a_j^{(k,t)}}{\partial z_j^{(k,t)}} = \frac{\partial c^{(n,t)}}{\partial a_j^{(k,t)}} \cdot \text{act}'(z_j^{(k,t)}) \\ &= \left( \sum_s \frac{\partial c^{(n,t)}}{\partial z_s^{(k+1,t)}} \cdot \frac{\partial z_s^{(k+1,t)}}{\partial a_j^{(k,t)}} + \sum_s \frac{\partial c^{(n,t)}}{\partial z_s^{(k,t+1)}} \cdot \frac{\partial z_s^{(k,t+1)}}{\partial a_j^{(k,t)}} \right) \text{act}'(z_j^{(k,t)})\end{aligned}$$



# Backward Pass through Time I

- The first term (error signal):  $\delta_j^{(k,t)} = \frac{\partial c^{(n,t)}}{\partial z_j^{(k,t)}}$
- We have

$$\begin{aligned}\delta_j^{(k,t)} &= \frac{\partial c^{(n,t)}}{\partial z_j^{(k,t)}} = \frac{\partial c^{(n,t)}}{\partial a_j^{(k,t)}} \cdot \frac{\partial a_j^{(k,t)}}{\partial z_j^{(k,t)}} = \frac{\partial c^{(n,t)}}{\partial a_j^{(k,t)}} \cdot \text{act}'(z_j^{(k,t)}) \\ &= \left( \sum_s \frac{\partial c^{(n,t)}}{\partial z_s^{(k+1,t)}} \cdot \frac{\partial z_s^{(k+1,t)}}{\partial a_j^{(k,t)}} + \sum_s \frac{\partial c^{(n,t)}}{\partial z_s^{(k,t+1)}} \cdot \frac{\partial z_s^{(k,t+1)}}{\partial a_j^{(k,t)}} \right) \text{act}'(z_j^{(k,t)}) \\ &= \left( \sum_s \delta_s^{(k+1,t)} W_{j,s}^{(k+1)} + \sum_s \delta_s^{(k,t+1)} U_{j,s}^{(k)} \right) \text{act}'(z_j^{(k,t)})\end{aligned}$$



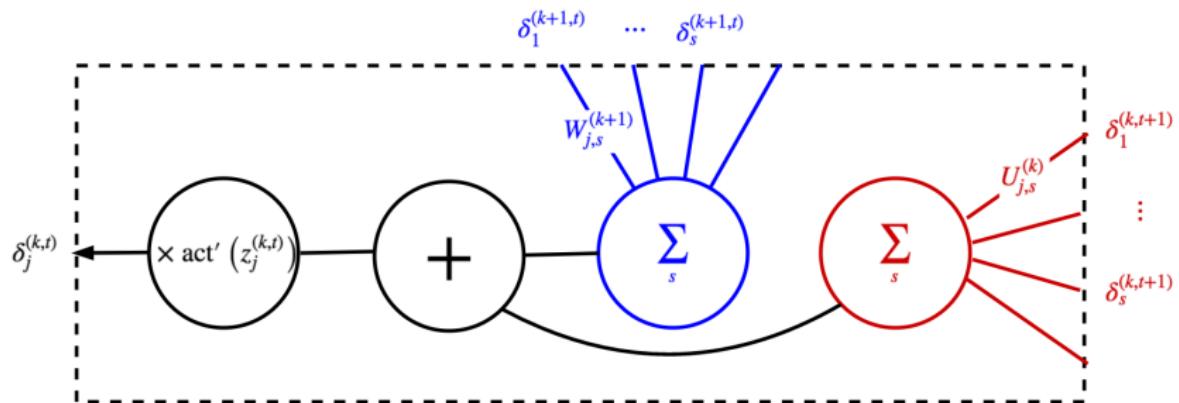
# Backward Pass through Time II

$$\delta_j^{(k,t)} = \left( \sum_s \delta_s^{(k+1,t)} W_{j,s}^{(k+1)} + \sum_s \delta_s^{(k,t+1)} U_{j,s}^{(k)} \right) \text{act}'(z_j^{(k,t)})$$

# Backward Pass through Time II

$$\delta_j^{(k,t)} = \left( \sum_s \delta_s^{(k+1,t)} W_{j,s}^{(k+1)} + \sum_s \delta_s^{(k,t+1)} U_{j,s}^{(k)} \right) \text{act}'(z_j^{(k,t)})$$

- We can evaluate all  $\delta_j^{(k,t)}$ 's starting from the deepest layer and **latest time**



# Backprop through Time (BPTT)

- So far, we have discussed how to compute the gradients of  $U_{i,j}^{(k)}$ 's (and  $W_{i,j}^{(k)}$ 's) for a single loss  $c^{(n,t)} = C^{(n,t)}(\Theta^{(s)})$  at specific time

# Backprop through Time (BPTT)

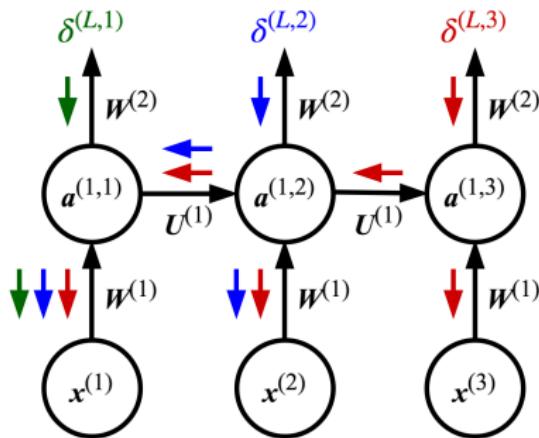
- So far, we have discussed how to compute the gradients of  $U_{i,j}^{(k)}$ 's (and  $W_{i,j}^{(k)}$ 's) for a single loss  $c^{(n,t)} = C^{(n,t)}(\Theta^{(s)})$  at specific time  $t$
- However, for each sequence  $n$ , we have multiple  $c^{(n,t)}$ 's at different  $t$ 's
  - Their gradient need to be summed

# Backprop through Time (BPTT)

- So far, we have discussed how to compute the gradients of  $U_{i,j}^{(k)}$ 's (and  $W_{i,j}^{(k)}$ 's) for a single loss  $c^{(n,t)} = C^{(n,t)}(\Theta^{(s)})$  at specific time
- However, for each sequence  $n$ , we have multiple  $c^{(n,t)}$ 's at different  $t$ 's
  - Their gradient need to be summed
- For different  $c^{(n,\cdot)}$ 's, the forward pass can be shared

# Backprop through Time (BPTT)

- So far, we have discussed how to compute the gradients of  $U_{i,j}^{(k)}$ 's (and  $W_{i,j}^{(k)}$ 's) for a single loss  $c^{(n,t)} = C^{(n,t)}(\Theta^{(s)})$  at specific time  $t$
- However, for each sequence  $n$ , we have multiple  $c^{(n,t)}$ 's at different  $t$ 's
  - Their gradient need to be summed
- For different  $c^{(n,\cdot)}$ 's, the forward pass can be shared
- **BPTT**: single forward pass, **multiple** backward passes



# Outline

## ① RNNs

- Vanilla RNNs
- Design Alternatives

## ② RNN Training

- Backprop through Time (BPTT)
- Optimization Techniques
- Optimization-Friendly Models & LSTM
- Parallelism & Teacher Forcing

## ③ RNNs with Attention Mechanism

- Attention for Image Captioning
- Attention for Neural Machine Translation (NMT)

## ④ Transformers

- Attention is All You Need
- Pretrained Language Models
- More Applications

## ⑤ Subword Tokenization

# Long-Term Dependencies

- Forward pass:

$$z^{(k,t)} \leftarrow U^{(k)\top} a^{(k,t-1)} + W^{(k)\top} a^{(k-1,t)}$$

$$a^{(k,t)} \leftarrow \text{act}(z^{(k,t)})$$

- Backward pass:

$$\delta^{(k,t)} \leftarrow \text{act}'(z^{(k,t)}) \odot (U^{(k)} \delta^{(k,t+1)} + W^{(k+1)} \delta^{(k+1,t)})$$

# Long-Term Dependencies

- Forward pass:

$$z^{(k,t)} \leftarrow \mathbf{U}^{(k)\top} \mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)\top} \mathbf{a}^{(k-1,t)}$$

$$\mathbf{a}^{(k,t)} \leftarrow \text{act}(z^{(k,t)})$$

- Backward pass:

$$\delta^{(k,t)} \leftarrow \text{act}'(z^{(k,t)}) \odot (\mathbf{U}^{(k)} \delta^{(k,t+1)} + \mathbf{W}^{(k+1)} \delta^{(k+1,t)})$$

- The dependency between  $\mathbf{a}^{(k,i)}$  (resp.  $\delta^{(k,i)}$ ) and  $\mathbf{a}^{(k,j)}$  (resp.  $\delta^{(k,j)}$ ) are maintained by  $(\mathbf{U}^{(k)})^{(j-i)}$ 
  - Ignoring activation function and depth, we have  $\mathbf{a}^{(k,j)} = (\mathbf{U}^{(k)\top})^{j-i} \mathbf{a}^{(k,i)}$  and  $\delta^{(k,i)} = (\mathbf{U}^{(k)})^{j-i} \delta^{(k,j)}$

# Long-Term Dependencies

- Forward pass:

$$z^{(k,t)} \leftarrow \mathbf{U}^{(k)\top} \mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)\top} \mathbf{a}^{(k-1,t)}$$

$$\mathbf{a}^{(k,t)} \leftarrow \text{act}(z^{(k,t)})$$

- Backward pass:

$$\delta^{(k,t)} \leftarrow \text{act}'(z^{(k,t)}) \odot (\mathbf{U}^{(k)} \delta^{(k,t+1)} + \mathbf{W}^{(k+1)} \delta^{(k+1,t)})$$

- The dependency between  $\mathbf{a}^{(k,i)}$  (resp.  $\delta^{(k,i)}$ ) and  $\mathbf{a}^{(k,j)}$  (resp.  $\delta^{(k,j)}$ ) are maintained by  $(\mathbf{U}^{(k)})^{(j-i)}$ 
  - Ignoring activation function and depth, we have  $\mathbf{a}^{(k,j)} = (\mathbf{U}^{(k)\top})^{j-i} \mathbf{a}^{(k,i)}$  and  $\delta^{(k,i)} = (\mathbf{U}^{(k)})^{j-i} \delta^{(k,j)}$
  - If  $\mathbf{a}^{(k,i)}$  and  $\mathbf{a}^{(k,j)}$  are far away in time, their long-term dependency causes optimization problems

# Exploding/Vanishing Gradient Problem

- Ignoring activation function and depth:

$$\mathbf{a}^{(k,j)} = (\mathbf{U}^{(k)\top})^{(j-i)} \mathbf{a}^{(k,i)} \text{ and } \boldsymbol{\delta}^{(k,i)} = (\mathbf{U}^{(k)})^{(j-i)} \boldsymbol{\delta}^{(k,j)}$$

- Given eigendecomposition:  $\mathbf{U}^{(k)} = \mathbf{Q} \begin{bmatrix} \lambda_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_{D^{(k)}} \end{bmatrix} \mathbf{Q}^\top$

# Exploding/Vanishing Gradient Problem

- Ignoring activation function and depth:

$$\mathbf{a}^{(k,j)} = (\mathbf{U}^{(k)\top})^{(j-i)} \mathbf{a}^{(k,i)} \text{ and } \boldsymbol{\delta}^{(k,i)} = (\mathbf{U}^{(k)})^{(j-i)} \boldsymbol{\delta}^{(k,j)}$$

- Given eigendecomposition:  $\mathbf{U}^{(k)} = \mathbf{Q} \begin{bmatrix} \lambda_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_{D^{(k)}} \end{bmatrix} \mathbf{Q}^\top$

- We have:  $(\mathbf{U}^{(k)})^{j-i} = \mathbf{Q} \begin{bmatrix} \lambda_1^{j-i} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_{D^{(k)}}^{j-i} \end{bmatrix} \mathbf{Q}^\top$

# Exploding/Vanishing Gradient Problem

- Ignoring activation function and depth:

$$\mathbf{a}^{(k,j)} = (\mathbf{U}^{(k)\top})^{(j-i)} \mathbf{a}^{(k,i)} \text{ and } \boldsymbol{\delta}^{(k,i)} = (\mathbf{U}^{(k)})^{(j-i)} \boldsymbol{\delta}^{(k,j)}$$

- Given eigendecomposition:  $\mathbf{U}^{(k)} = \mathbf{Q} \begin{bmatrix} \lambda_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_{D^{(k)}} \end{bmatrix} \mathbf{Q}^\top$
- We have:  $(\mathbf{U}^{(k)})^{j-i} = \mathbf{Q} \begin{bmatrix} \lambda_1^{j-i} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_{D^{(k)}}^{j-i} \end{bmatrix} \mathbf{Q}^\top$
- When  $j - i$  is large,  $\lambda_s^{j-i}$  is either very large or small
  - $\lambda_s = 1.01 \Rightarrow \lambda_s^{1000} \approx 20,000$
  - $\lambda_s = 0.99 \Rightarrow \lambda_s^{1000} \approx 0$

# Exploding/Vanishing Gradient Problem

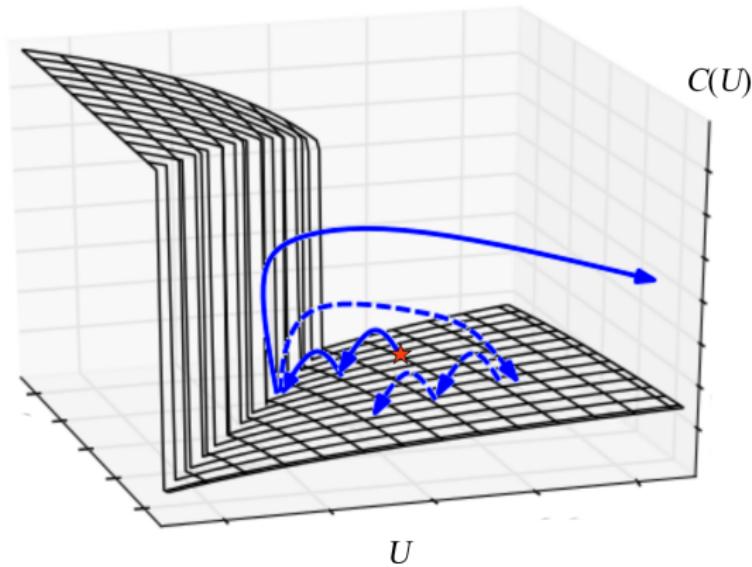
- Ignoring activation function and depth:

$$\mathbf{a}^{(k,j)} = (\mathbf{U}^{(k)\top})^{(j-i)} \mathbf{a}^{(k,i)} \text{ and } \boldsymbol{\delta}^{(k,i)} = (\mathbf{U}^{(k)})^{(j-i)} \boldsymbol{\delta}^{(k,j)}$$

- Given eigendecomposition:  $\mathbf{U}^{(k)} = \mathbf{Q} \begin{bmatrix} \lambda_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_{D^{(k)}} \end{bmatrix} \mathbf{Q}^\top$
- We have:  $(\mathbf{U}^{(k)})^{j-i} = \mathbf{Q} \begin{bmatrix} \lambda_1^{j-i} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_{D^{(k)}}^{j-i} \end{bmatrix} \mathbf{Q}^\top$
- When  $j - i$  is large,  $\lambda_s^{j-i}$  is either very large or small
  - $\lambda_s = 1.01 \Rightarrow \lambda_s^{1000} \approx 20,000$
  - $\lambda_s = 0.99 \Rightarrow \lambda_s^{1000} \approx 0$
- Exploding or vanishing gradients!

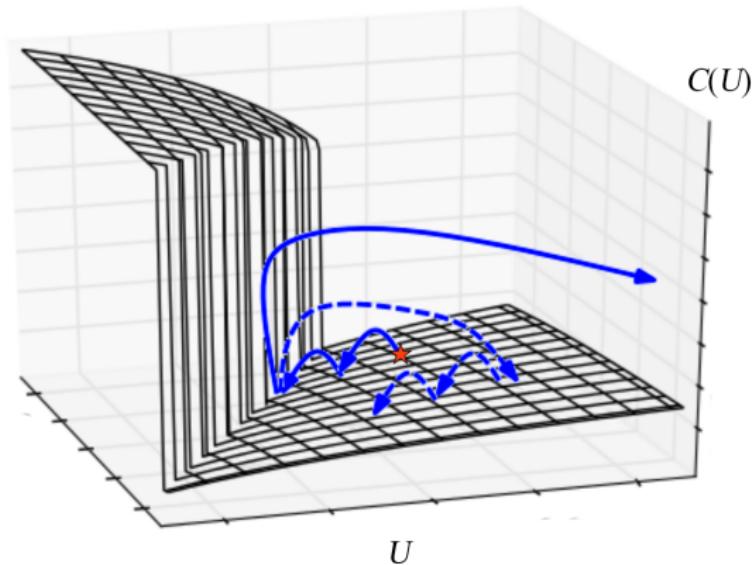
# Cost Surface

- The cost surface of  $C$  is either very flat or steep
- Hard for gradient-based optimization



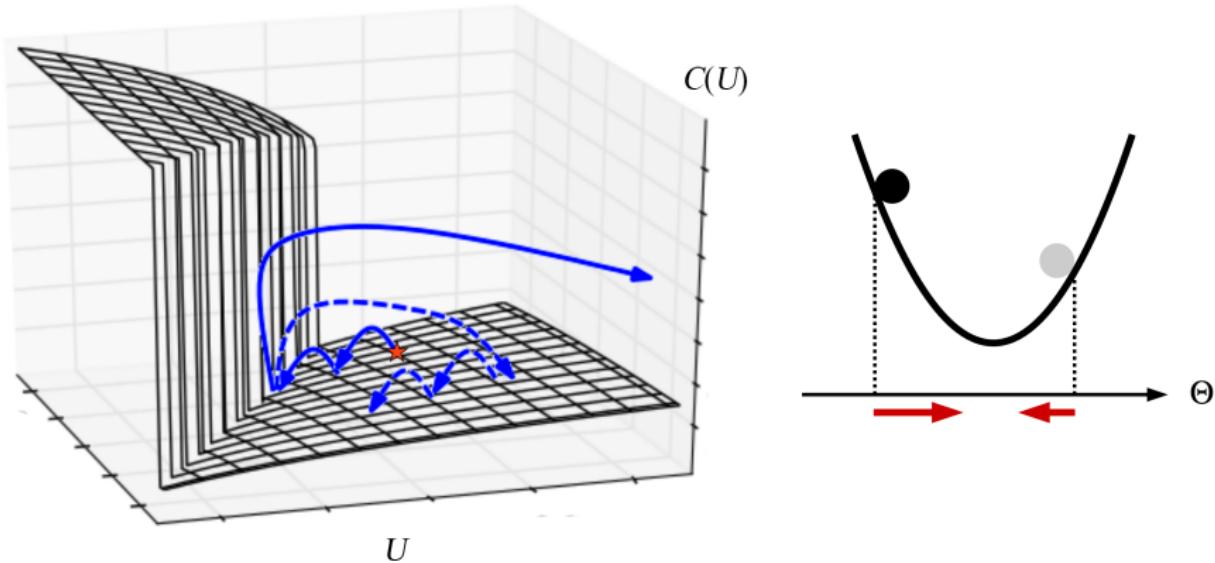
# Cost Surface

- The cost surface of  $C$  is either very flat or steep
- Hard for gradient-based optimization
- Optimization techniques?



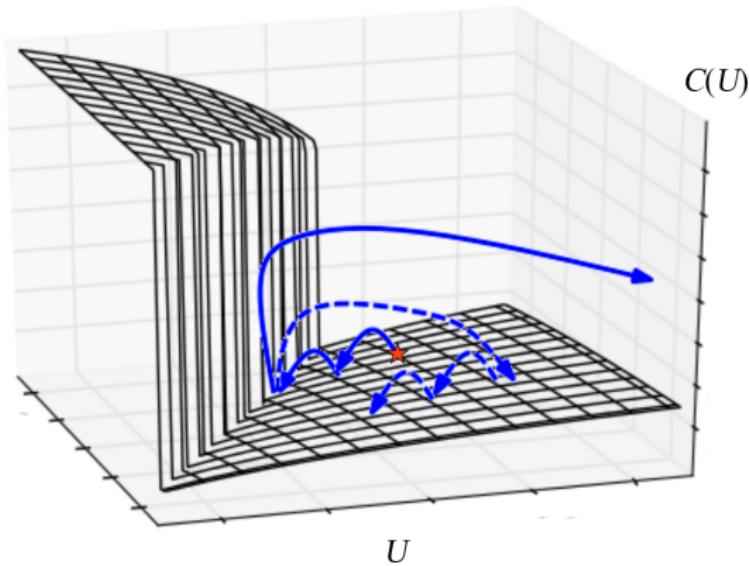
# Nesterov Momentum

- Use Nesterov momentum to “brake” before hitting the wall



# Gradient Clipping

- A simple way is to avoid the exploding gradient problem is to *clip* a gradient if it exceeds a predefined threshold
- Very effective in practice

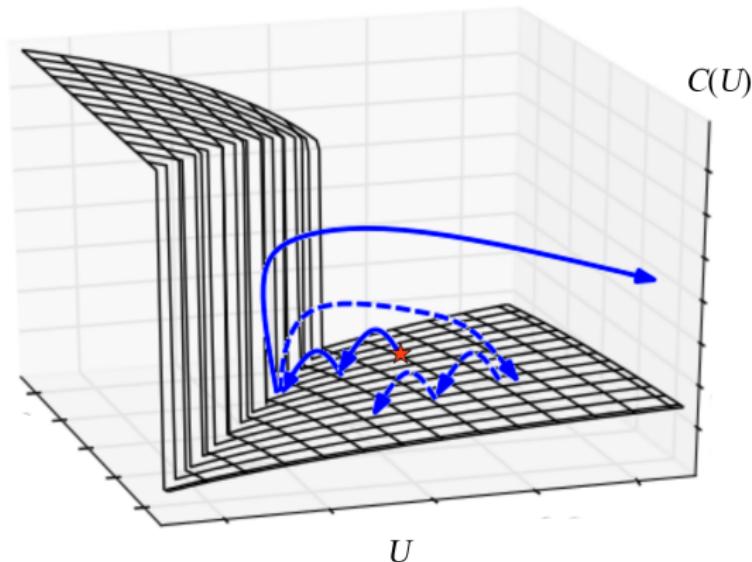


# RMS Prop

- Adaptive learning rate based on statistics of recent gradients:

$$\mathbf{r}^{(t+1)} \leftarrow \lambda \mathbf{r}^{(t)} + (1 - \lambda) \mathbf{g}^{(t)} \odot \mathbf{g}^{(t)}$$

$$\Theta^{(t+1)} \leftarrow \Theta^{(t)} - \frac{\eta}{\sqrt{\mathbf{r}^{(t+1)}}} \odot \mathbf{g}^{(t)}$$



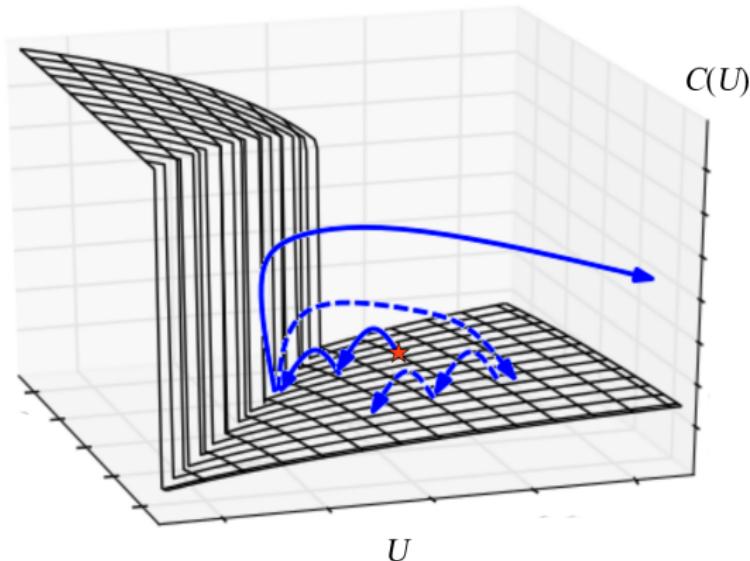
# RMS Prop

- Adaptive learning rate based on statistics of recent gradients:

$$\mathbf{r}^{(t+1)} \leftarrow \lambda \mathbf{r}^{(t)} + (1 - \lambda) \mathbf{g}^{(t)} \odot \mathbf{g}^{(t)}$$

$$\Theta^{(t+1)} \leftarrow \Theta^{(t)} - \frac{\eta}{\sqrt{\mathbf{r}^{(t+1)}}} \odot \mathbf{g}^{(t)}$$

- Reduce  $\lambda$



# Outline

## ① RNNs

- Vanilla RNNs
- Design Alternatives

## ② RNN Training

- Backprop through Time (BPTT)
- Optimization Techniques
- Optimization-Friendly Models & LSTM
- Parallelism & Teacher Forcing

## ③ RNNs with Attention Mechanism

- Attention for Image Captioning
- Attention for Neural Machine Translation (NMT)

## ④ Transformers

- Attention is All You Need
- Pretrained Language Models
- More Applications

## ⑤ Subword Tokenization

# Sigmoid Activation Function

- Traditionally, the hidden units of RNNs use the sigmoid activation functions
  - E.g.,  $\tanh(\cdot)$  the *hyperbolic tangent*

# Sigmoid Activation Function

- Traditionally, the hidden units of RNNs use the sigmoid activation functions
  - E.g.,  $\tanh(\cdot)$  the *hyperbolic tangent*
- Why sigmoid activation function?

# Sigmoid Activation Function

- Traditionally, the hidden units of RNNs use the sigmoid activation functions
  - E.g.,  $\tanh(\cdot)$  the *hyperbolic tangent*
- Why sigmoid activation function?
- Forward pass:

$$z^{(k,t)} \leftarrow \mathbf{U}^{(k)\top} \mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)\top} \mathbf{a}^{(k-1,t)}$$

$$\mathbf{a}^{(k,t)} \leftarrow \text{act}(z^{(k,t)})$$

Backward pass:

$$\delta^{(k,t)} \leftarrow \text{act}'(z^{(k,t)}) \odot (\mathbf{U}^{(k)} \delta^{(k,t+1)} + \mathbf{W}^{(k+1)} \delta^{(k+1,t)})$$

# Sigmoid Activation Function

- Traditionally, the hidden units of RNNs use the sigmoid activation functions
  - E.g.,  $\tanh(\cdot)$  the *hyperbolic tangent*
- Why sigmoid activation function?
- Forward pass:

$$z^{(k,t)} \leftarrow \mathbf{U}^{(k)\top} \mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)\top} \mathbf{a}^{(k-1,t)}$$

$$\mathbf{a}^{(k,t)} \leftarrow \text{act}(z^{(k,t)})$$

Backward pass:

$$\delta^{(k,t)} \leftarrow \text{act}'(z^{(k,t)}) \odot (\mathbf{U}^{(k)} \delta^{(k,t+1)} + \mathbf{W}^{(k+1)} \delta^{(k+1,t)})$$

- Sigmoid activation:
  - Bounded range in the forward pass
  - $\text{act}'(\cdot) < 1$  in the backward pass

# Sigmoid Activation Function

- Traditionally, the hidden units of RNNs use the sigmoid activation functions
  - E.g.,  $\tanh(\cdot)$  the *hyperbolic tangent*
- Why sigmoid activation function?
- Forward pass:

$$z^{(k,t)} \leftarrow \mathbf{U}^{(k)\top} \mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)\top} \mathbf{a}^{(k-1,t)}$$

$$\mathbf{a}^{(k,t)} \leftarrow \text{act}(z^{(k,t)})$$

Backward pass:

$$\delta^{(k,t)} \leftarrow \text{act}'(z^{(k,t)}) \odot (\mathbf{U}^{(k)} \delta^{(k,t+1)} + \mathbf{W}^{(k+1)} \delta^{(k+1,t)})$$

- Sigmoid activation:
  - Bounded range in the forward pass
  - $\text{act}'(\cdot) < 1$  in the backward pass
- Mitigates the exploding (but not vanishing) gradient problem for  $\mathbf{U}^{(k)}$ 's

# Sigmoid Activation Function

- Traditionally, the hidden units of RNNs use the sigmoid activation functions
  - E.g.,  $\tanh(\cdot)$  the *hyperbolic tangent*
- Why sigmoid activation function?
- Forward pass:

$$z^{(k,t)} \leftarrow \mathbf{U}^{(k)\top} \mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)\top} \mathbf{a}^{(k-1,t)}$$

$$\mathbf{a}^{(k,t)} \leftarrow \text{act}(z^{(k,t)})$$

Backward pass:

$$\delta^{(k,t)} \leftarrow \text{act}'(z^{(k,t)}) \odot (\mathbf{U}^{(k)} \delta^{(k,t+1)} + \mathbf{W}^{(k+1)} \delta^{(k+1,t)})$$

- Sigmoid activation:
  - Bounded range in the forward pass
  - $\text{act}'(\cdot) < 1$  in the backward pass
- Mitigates the exploding (but not vanishing) gradient problem for  $\mathbf{U}^{(k)}$ 's
- Introduces vanishing gradients of  $\mathbf{W}^{(k)}$ 's

# Learning Unitary $U^{(k)}$ 's

- Long-term dependency:  $(U^{(k)})^{j-i} = Q \begin{bmatrix} \lambda_1^{j-i} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_{D^{(k)}}^{j-i} \end{bmatrix} Q^\top$

# Learning Unitary $U^{(k)}$ 's

- Long-term dependency:  $(U^{(k)})^{j-i} = Q \begin{bmatrix} \lambda_1^{j-i} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_{D^{(k)}}^{j-i} \end{bmatrix} Q^\top$
- Why not make  $U^{(k)}$ 's unitary (i.e.,  $\lambda_s = 1$  for all  $s$ )?

# Learning Unitary $U^{(k)}$ 's

$$\begin{bmatrix} \lambda_1^{j-i} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_{D^{(k)}}^{j-i} \end{bmatrix} \mathbf{Q}^\top$$

- Long-term dependency:  $(\mathbf{U}^{(k)})^{j-i} = \mathbf{Q} \begin{bmatrix} \lambda_1^{j-i} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_{D^{(k)}}^{j-i} \end{bmatrix} \mathbf{Q}^\top$
- Why not make  $\mathbf{U}^{(k)}$ 's unitary (i.e.,  $\lambda_s = 1$  for all  $s$ )?
- Hinton et al. [10] propose IRNN:
  - **Initializes  $\mathbf{U}^{(k)} = \mathbf{I}$**  in SGD
  - Uses **ReLU** hidden units

# Learning Unitary $U^{(k)}$ 's

- Long-term dependency:  $(U^{(k)})^{j-i} = Q \begin{bmatrix} \lambda_1^{j-i} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_{D^{(k)}}^{j-i} \end{bmatrix} Q^\top$
- Why not make  $U^{(k)}$ 's unitary (i.e.,  $\lambda_s = 1$  for all  $s$ )?
- Hinton et al. [10] propose IRNN:
  - **Initializes  $U^{(k)} = I$**  in SGD
  - Uses **ReLU** hidden units
  - Empirically, requires a very small learning rate (e.g.,  $10^{-8}$ ) to work well
  - Simple, but very slow

# Learning Unitary $U^{(k)}$ 's

- Long-term dependency:  $(U^{(k)})^{j-i} = Q \begin{bmatrix} \lambda_1^{j-i} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_{D^{(k)}}^{j-i} \end{bmatrix} Q^\top$
- Why not make  $U^{(k)}$ 's unitary (i.e.,  $\lambda_s = 1$  for all  $s$ )?
- Hinton et al. [10] propose IRNN:
  - **Initializes  $U^{(k)} = I$**  in SGD
  - Uses **ReLU** hidden units
  - Empirically, requires a very small learning rate (e.g.,  $10^{-8}$ ) to work well
  - Simple, but very slow
- Krueger et al. [7] add a term  $\sum_{k,t} (\|a^{(k,t)}\| - \|a^{(k,t-1)}\|)^2$  to IRNN cost to stabilize the norms of  $a^{(k,t)}$ 's in time

# Learning Unitary $U^{(k)}$ 's

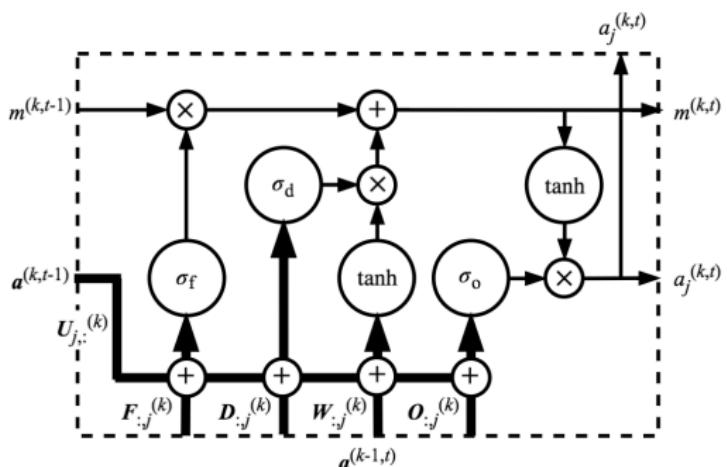
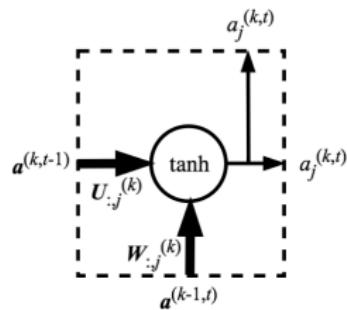
- Long-term dependency:  $(U^{(k)})^{j-i} = Q \begin{bmatrix} \lambda_1^{j-i} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_{D^{(k)}}^{j-i} \end{bmatrix} Q^\top$
- Why not make  $U^{(k)}$ 's unitary (i.e.,  $\lambda_s = 1$  for all  $s$ )?
- Hinton et al. [10] propose IRNN:
  - **Initializes  $U^{(k)} = I$**  in SGD
  - Uses **ReLU** hidden units
  - Empirically, requires a very small learning rate (e.g.,  $10^{-8}$ ) to work well
  - Simple, but very slow
- Krueger et al. [7] add a term  $\sum_{k,t} (\|a^{(k,t)}\| - \|a^{(k,t-1)}\|)^2$  to IRNN cost to stabilize the norms of  $a^{(k,t)}$ 's in time
- Bengio et al. [1] propose uRNN that learns unitary  $U^{(k)}$ 's explicitly

# Long Short-Term Memory (LSTM)

- Idea: to create *shortcut* in each neuron for the error signals to flow backward more smoothly

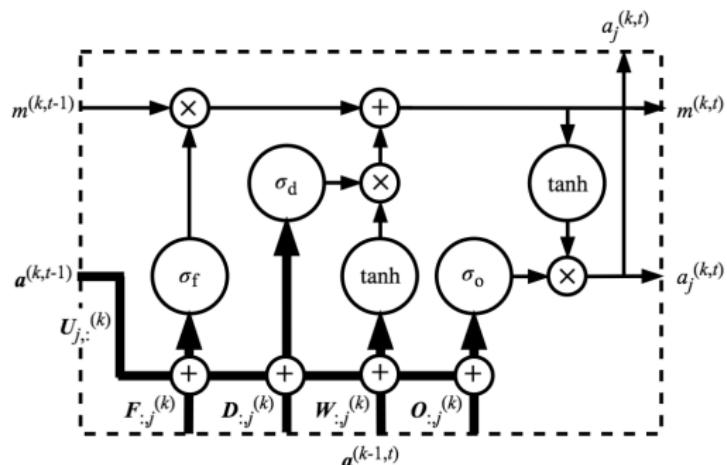
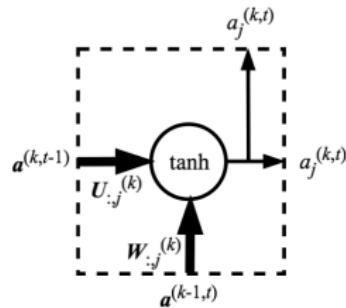
# Long Short-Term Memory (LSTM)

- Idea: to create **shortcut** in each neuron for the error signals to flow backward more smoothly
- The  $j$ -th LSTM unit at depth  $k$  and time  $t$ :



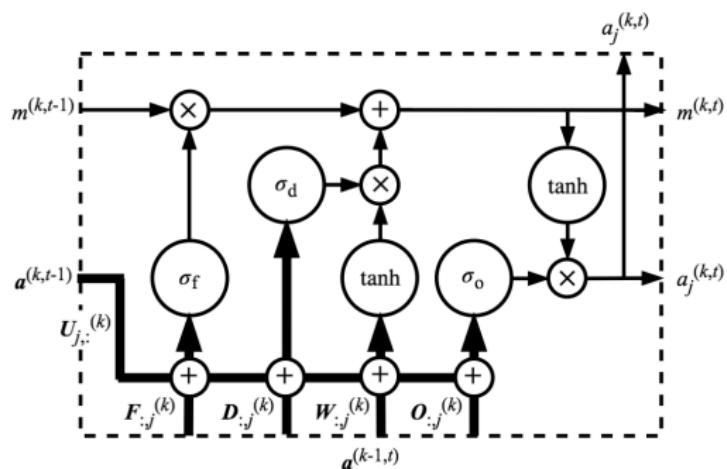
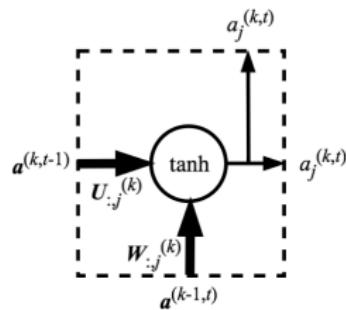
# Long Short-Term Memory (LSTM)

- Idea: to create **shortcut** in each neuron for the error signals to flow backward more smoothly
- The  $j$ -th LSTM unit at depth  $k$  and time  $t$ :
  - First tanh activation to be added into the **memory cell**  $m^{(k,t)}$
  - Second tanh activation as the final activation  $a_j^{(k,t)}$



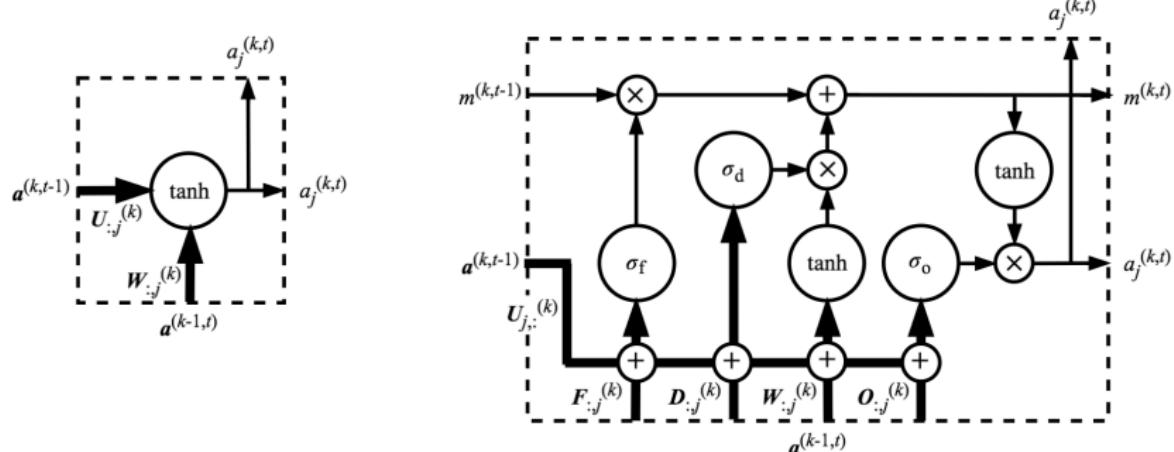
# Long Short-Term Memory (LSTM)

- Idea: to create **shortcut** in each neuron for the error signals to flow backward more smoothly
- The  $j$ -th LSTM unit at depth  $k$  and time  $t$ :
  - First tanh activation to be added into the **memory cell**  $m^{(k,t)}$
  - Second tanh activation as the final activation  $a_j^{(k,t)}$
  - Forget gate  $\sigma_f$ : whether to forget  $m^{(k,t-1)}$



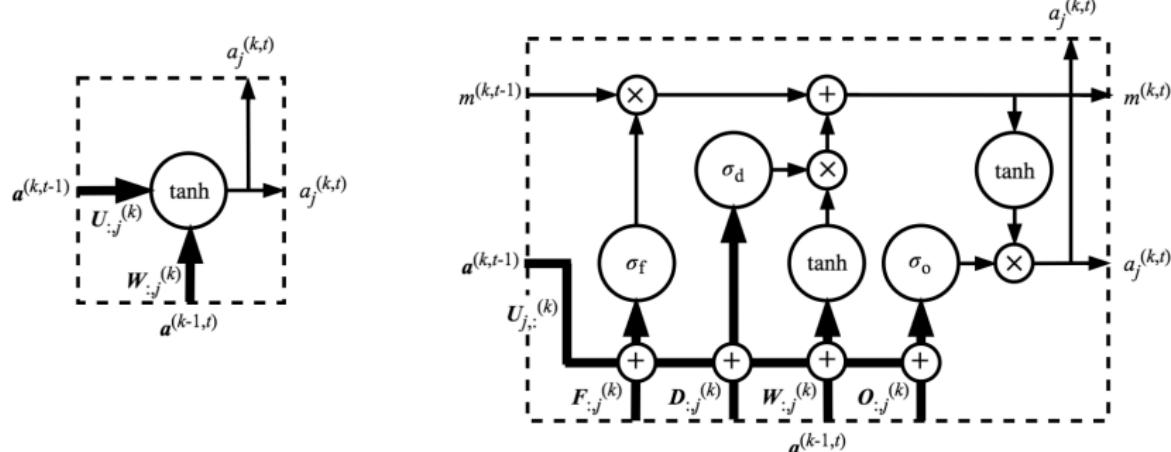
# Long Short-Term Memory (LSTM)

- Idea: to create **shortcut** in each neuron for the error signals to flow backward more smoothly
- The  $j$ -th LSTM unit at depth  $k$  and time  $t$ :
  - First tanh activation to be added into the **memory cell**  $m^{(k,t)}$
  - Second tanh activation as the final activation  $a_j^{(k,t)}$
  - Forget gate  $\sigma_f$ : whether to forget  $m^{(k,t-1)}$
  - Input gate  $\sigma_d$ : whether to store the first activation into  $m^{(k,t)}$



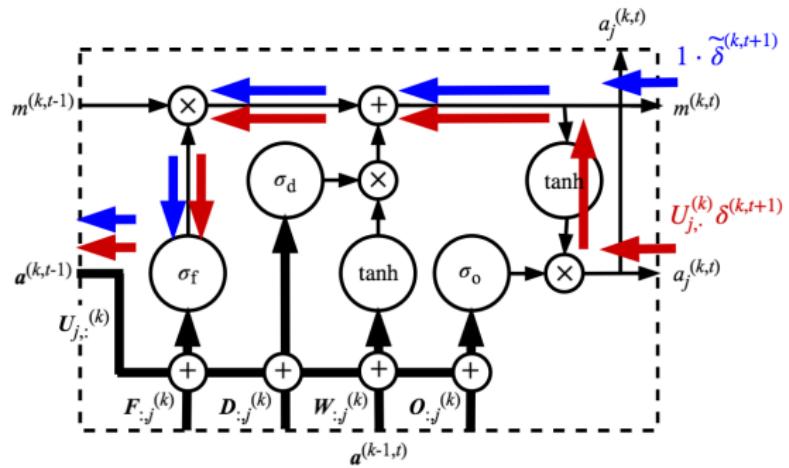
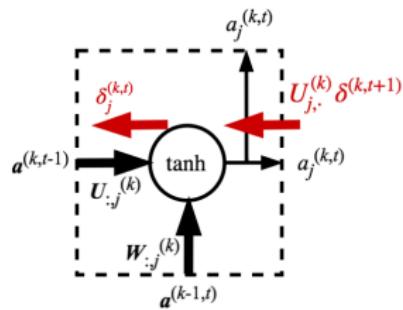
# Long Short-Term Memory (LSTM)

- Idea: to create **shortcut** in each neuron for the error signals to flow backward more smoothly
- The  $j$ -th LSTM unit at depth  $k$  and time  $t$ :
  - First tanh activation to be added into the **memory cell**  $m^{(k,t)}$
  - Second tanh activation as the final activation  $a_j^{(k,t)}$
  - Forget gate  $\sigma_f$ : whether to forget  $m^{(k,t-1)}$
  - Input gate  $\sigma_d$ : whether to store the first activation into  $m^{(k,t)}$
  - Output gate  $\sigma_o$ : whether to output the second activation



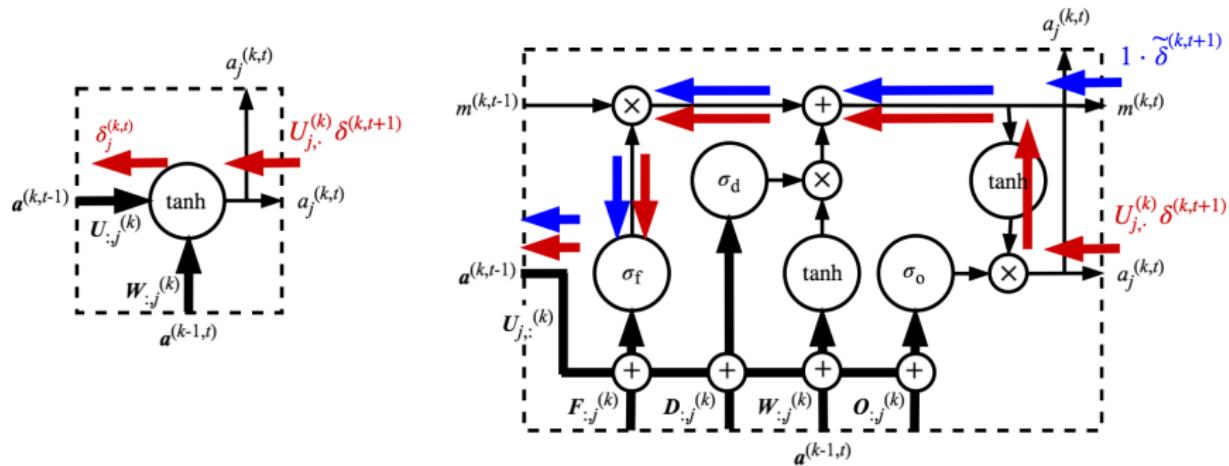
# Error Signals

- Error signals now have a second path
  - If the forget gate is open, error signals won't decay (**blue arrows**)



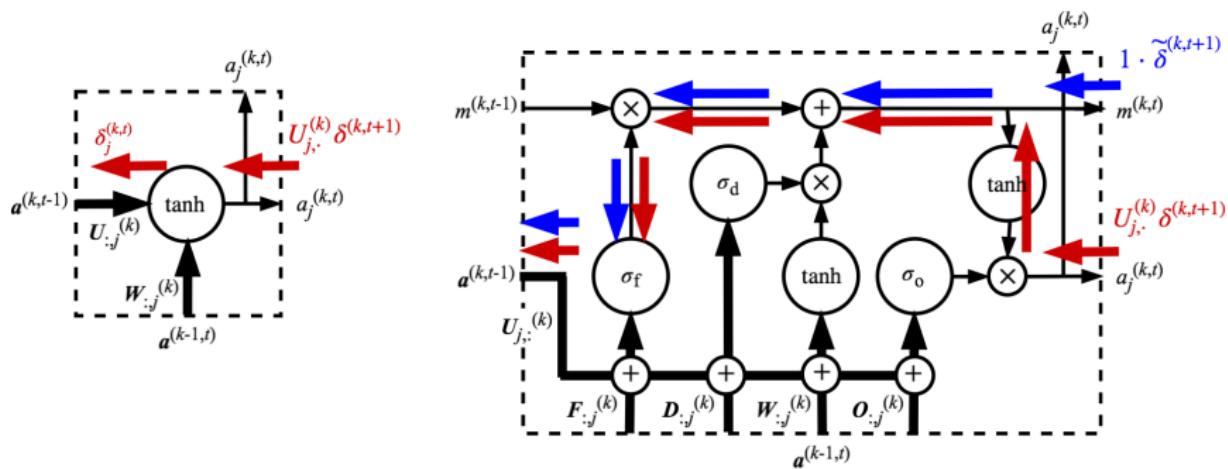
# Error Signals

- Error signals now have a second path
  - If the forget gate is open, error signals won't decay (**blue arrows**)
  - Avoids the vanishing gradients (but not exploding ones)

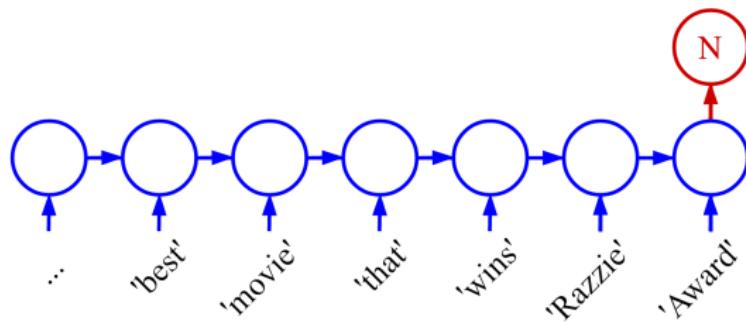


# Error Signals

- Error signals now have a second path
  - If the forget gate is open, error signals won't decay (**blue arrows**)
  - Avoids the vanishing gradients (but not exploding ones)
- When NN decides to close the forget gate, the vanishing gradient problem is irrelevant
- In practice, LSTM + gradient clipping works well together

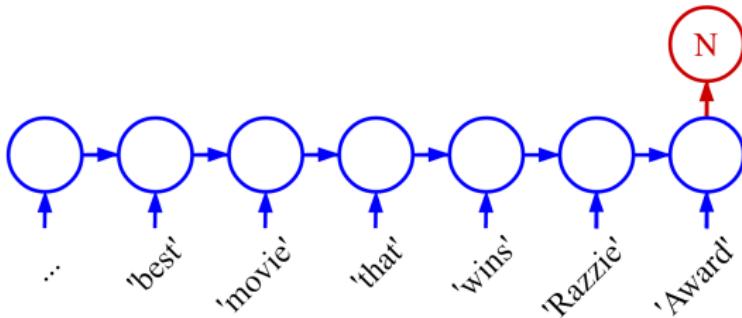


# Dynamic Representations for Sentiment Analysis



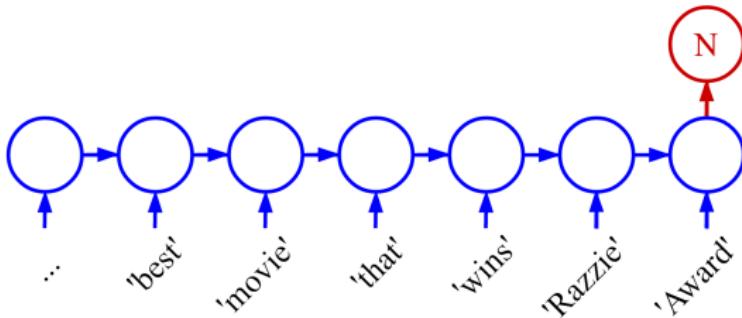
- LSTM units learn dynamic representations

# Dynamic Representations for Sentiment Analysis



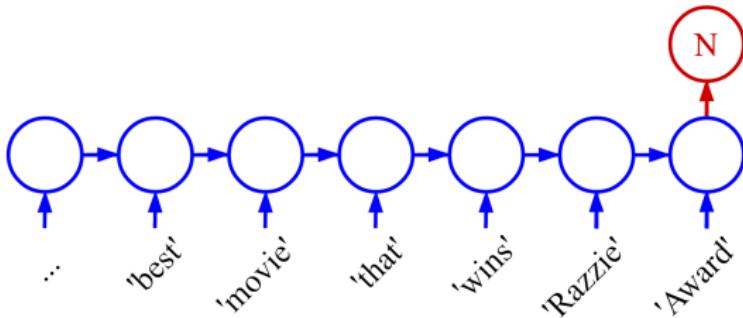
- LSTM units learn dynamic representations
- Closing forget gate for “Razzie”
  - To correct previous summarization and/or shift focus

# Dynamic Representations for Sentiment Analysis



- LSTM units learn dynamic representations
- Closing forget gate for “Razzie”
  - To correct previous summarization and/or shift focus
- Closing input gate for “movie”
  - Not to learn, to keep the same summarization

# Dynamic Representations for Sentiment Analysis



- LSTM units learn dynamic representations
- Closing forget gate for “Razzie”
  - To correct previous summarization and/or shift focus
- Closing input gate for “movie”
  - Not to learn, to keep the same summarization
- Closing output gate for “that”
  - To let the next neuron decide the activation/gate values by its own

# Dynamic Representations for Language Models

# Dynamic Representations for Language Models

## • Neuron activations for language modeling [6]

▶ Interactive Tool

Cell sensitive to position in line:

```
The sole importance of the crossing of the Berezina lies in the fact  
that it plainly and indubitably proved the fallacy of all the plans for  
cutting off the enemy's retreat and the soundness of the only possible  
line of action--the one Kutuzov and the general mass of the army  
demanded--namely, simply to follow the enemy up. The French crowd fled  
at a continually increasing speed and all its energy was directed to  
reaching its goal. It fled like a wounded animal and it was impossible  
to block its path. This was shown not so much by the arrangements it  
made for crossing as by what took place at the bridges. When the bridges  
broke down, unarmed soldiers, people from Moscow and women with children  
who were with the French transport, all--carried on by vis inertiae--  
pressed forward into boats and into the ice-covered water and did not,  
surrender.
```

Cell that turns on inside quotes:

```
"You mean to imply that I have nothing to eat out of.... On the  
contrary, I can supply you with everything even if you want to give  
dinner parties," warmly replied Chichagov, who tried by every word he  
spoke to prove his own rectitude and therefore imagined Kutuzov to be  
animated by the same desire.
```

```
Kutuzov, shrugging his shoulders, replied with his subtle penetrating  
smile: "I meant merely to say what I said."
```

Cell that robustly activates inside if statements:

```
static int __dequeue_signal(struct sigpending *pending, sigset_t *mask,  
    siginfo_t *info)  
{  
    int sig = next_signal(pending, mask);  
    if (sig) {  
        if (current->notifier) {  
            if (sigismember(current->notifier_mask, sig)) {  
                if (!!(current->notifier)(current->notifier_data)) {  
                    clear_thread_flag(TIF_SIGPENDING);  
                    return 0;  
                }  
            }  
        }  
        collect_signal(sig, pending, info);  
    }  
    return sig;  
}
```

# Learning Process of LSTMs

- Output at epoch 100:

*“... tyntd-iafhatawiaoihrdemot lytdws e ,tfti, astai f ogoh...”*

# Learning Process of LSTMs

- Output at epoch 100:

“... *tyntd-iafhatawiaoihrdemot lytdws e ,tfti, astai f ogoh...*”

- At 300 (spaces and periods):

“... *Phe lism thond hon at. MeiDimorotion in ther...*”

# Learning Process of LSTMs

- Output at epoch 100:

*“... tyntd-iafhatawiaoihrdemot lytdws e ,tfti, astai f ogoh...”*

- At 300 (spaces and periods):

*“... Phe lism thond hon at. MeiDimorotion in ther...”*

- At 500 (common words “we,” “he,” etc.):

*“... we counter. He stutn co des. His stanted out one...”*

# Learning Process of LSTMs

- Output at epoch 100:

*“... tyntd-iafhatawiaoihrdemot lytdws e ,tfti, astai f ogoh...”*

- At 300 (spaces and periods):

*“... Phe lism thond hon at. MeiDimorotion in ther...”*

- At 500 (common words “we,” “he,” etc.):

*“... we counter. He stutn co des. His stanted out one...”*

- At 700 (English-like structure):

*“... Aftair fall unsuch that the hall for Prince Velzonski's...”*

# Learning Process of LSTMs

- Output at epoch 100:

*“... tyntd-iafhatawiaoihrdemot lytdws e ,tfti, astai f ogoh...”*

- At 300 (spaces and periods):

*“... Phe lism thond hon at. MeiDimorotion in ther...”*

- At 500 (common words “we,” “he,” etc.):

*“... we counter. He stutn co des. His stanted out one...”*

- At 700 (English-like structure):

*“... Aftair fall unsuch that the hall for Prince Velzonski's...”*

- At 1200 (quotations and longer words):

*“... “Kite vouch!” he repeated by her door...”*

# Learning Process of LSTMs

- Output at epoch 100:

*“... tyntd-iafhatawiaoihrdemot lytdws e ,tfti, astai f ogoh...”*

- At 300 (spaces and periods):

*“... Phe lism thond hon at. MeiDimorotion in ther...”*

- At 500 (common words “we,” “he,” etc.):

*“... we counter. He stutn co des. His stanted out one...”*

- At 700 (English-like structure):

*“... Aftair fall unsuch that the hall for Prince Velzonski's...”*

- At 1200 (quotations and longer words):

*“... “Kite vouch!” he repeated by her door...”*

- At 2000 (topics and longer-term dependencies):

*“... “Why do what that day,” replied Natasha, ...”*

# Outline

## ① RNNs

- Vanilla RNNs
- Design Alternatives

## ② RNN Training

- Backprop through Time (BPTT)
- Optimization Techniques
- Optimization-Friendly Models & LSTM
- Parallelism & Teacher Forcing

## ③ RNNs with Attention Mechanism

- Attention for Image Captioning
- Attention for Neural Machine Translation (NMT)

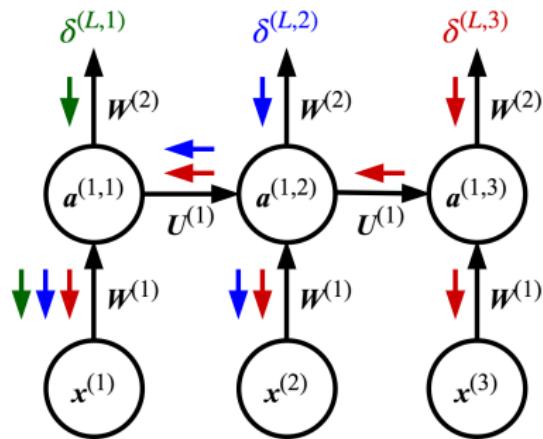
## ④ Transformers

- Attention is All You Need
- Pretrained Language Models
- More Applications

## ⑤ Subword Tokenization

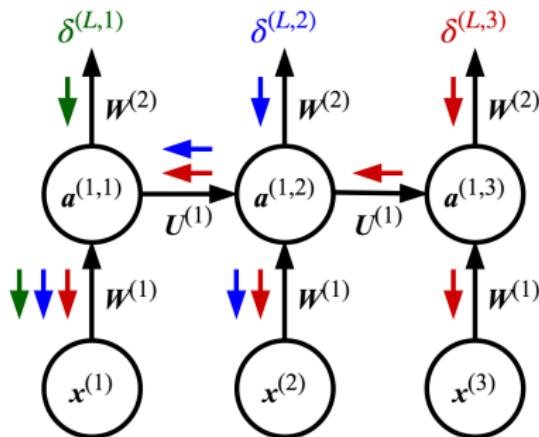
# Parallelism

- A forward/backward pass through time in BPTT cannot be parallelized



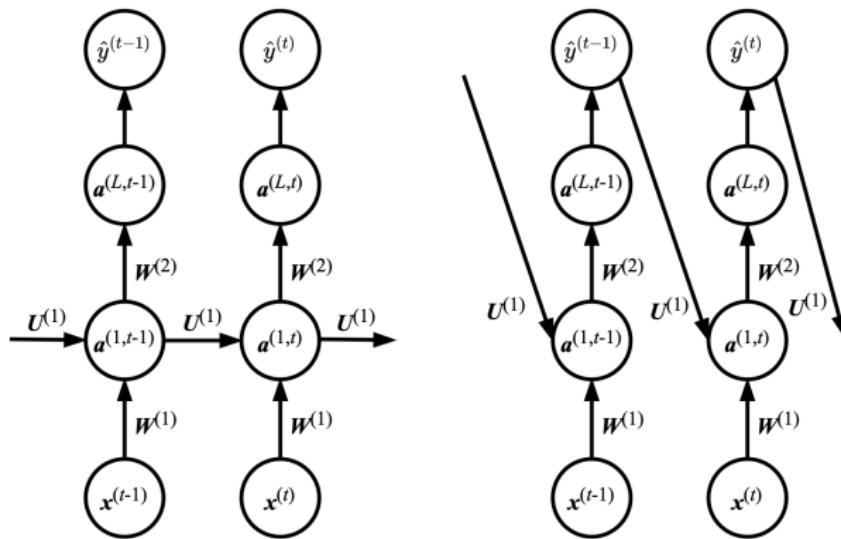
# Parallelism

- A forward/backward pass through time in BPTT cannot be parallelized
- The **hidden-to-hidden** recurrent connections in a vanilla RNN create dependency between
  - $a^{(k,t)}$ 's in forward pass
  - $\delta^{(k,t)}$ 's in backward pass



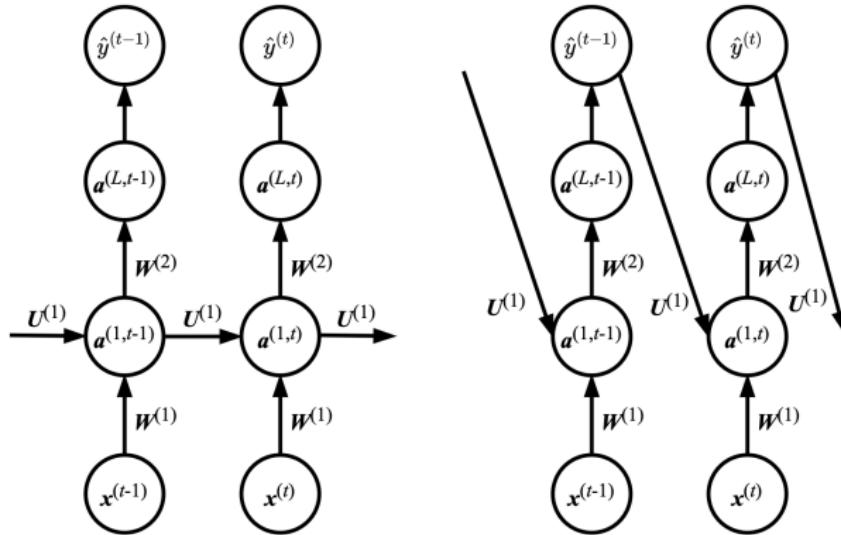
# Output Recurrence and Teacher Forcing

- **Teacher forcing:** replace hidden-to-hidden recurrence with output-to-hidden or output-to-input recurrence



# Output Recurrence and Teacher Forcing

- **Teacher forcing:** replace hidden-to-hidden recurrence with output-to-hidden or output-to-input recurrence
- At training time, use **correct labels  $y^{(\cdot)}$ 's** to train the model
  - So, the forward/backward pass through time can be parallelized
- At test time, switch back to using model output  $\hat{y}^{(\cdot)}$ 's



# Cost: Exposure Bias

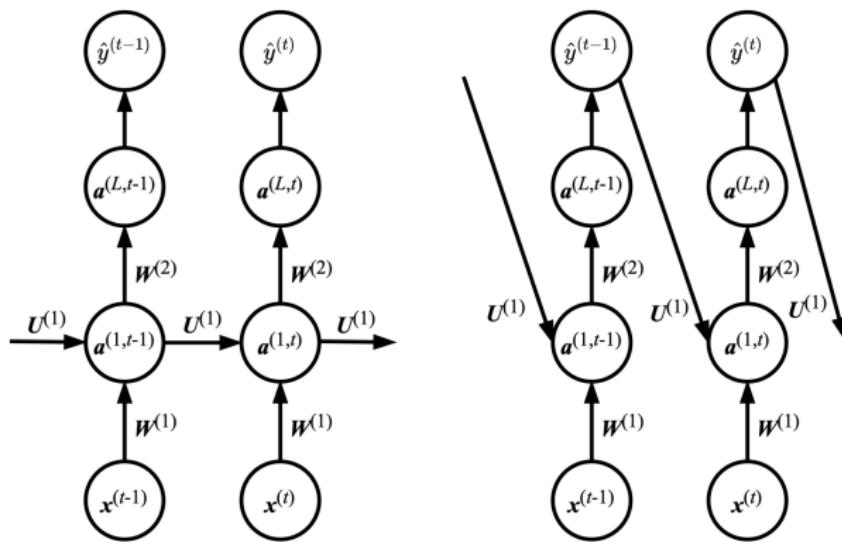
- Mismatch between  $y^{(\cdot)}$ 's and  $\hat{y}^{(\cdot)}$ 's hurts RNN performance
- Solution?

# Cost: Exposure Bias

- Mismatch between  $y^{(\cdot)}$ 's and  $\hat{y}^{(\cdot)}$ 's hurts RNN performance
- Solution? Scheduled sampling
- At training time,
  - ① Use  $y^{(\cdot)}$ 's initially
  - ② Gradually mix in  $\hat{y}^{(\cdot)}$ 's later

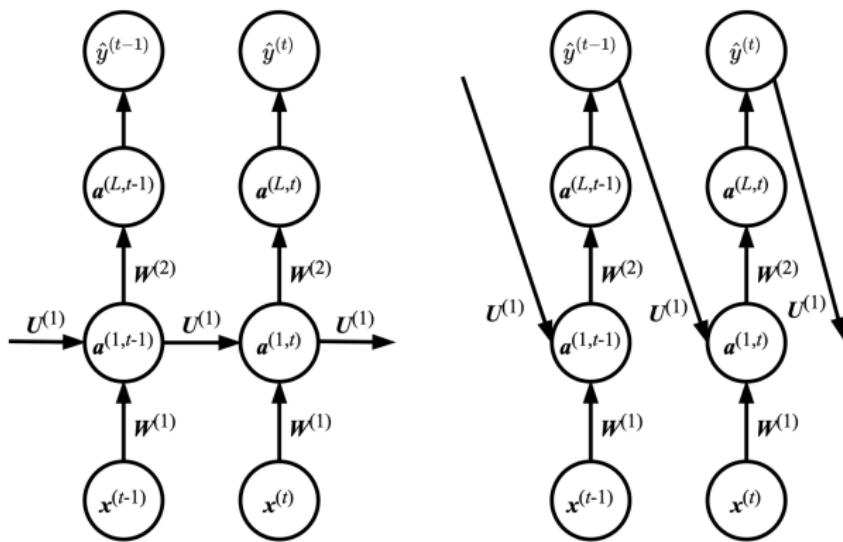
# Cost: Reduced Expressiveness

- The vanilla RNNs are universal in the sense that they can simulate Turing machines [18]



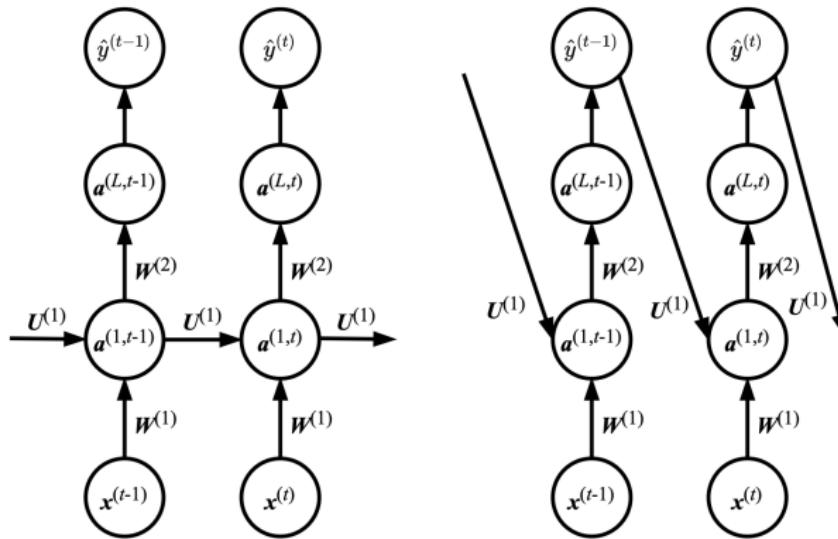
# Cost: Reduced Expressiveness

- The vanilla RNNs are universal in the sense that they can simulate Turing machines [18]
- Output-recurrent RNNs **cannot** simulate Turing machines and are strictly less powerful



# Cost: Reduced Expressiveness

- The vanilla RNNs are universal in the sense that they can simulate Turing machines [18]
- Output-recurrent RNNs **cannot** simulate Turing machines and are strictly less powerful
- The output  $a^{(L,\cdot)}$ 's are explicitly trained to match training targets
  - Cannot capture all required information in the past to predict the future



# Outline

## ① RNNs

- Vanilla RNNs
- Design Alternatives

## ② RNN Training

- Backprop through Time (BPTT)
- Optimization Techniques
- Optimization-Friendly Models & LSTM
- Parallelism & Teacher Forcing

## ③ RNNs with Attention Mechanism

- Attention for Image Captioning
- Attention for Neural Machine Translation (NMT)

## ④ Transformers

- Attention is All You Need
- Pretrained Language Models
- More Applications

## ⑤ Subword Tokenization

# Outline

## ① RNNs

- Vanilla RNNs
- Design Alternatives

## ② RNN Training

- Backprop through Time (BPTT)
- Optimization Techniques
- Optimization-Friendly Models & LSTM
- Parallelism & Teacher Forcing

## ③ RNNs with Attention Mechanism

- Attention for Image Captioning
- Attention for Neural Machine Translation (NMT)

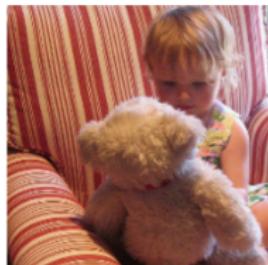
## ④ Transformers

- Attention is All You Need
- Pretrained Language Models
- More Applications

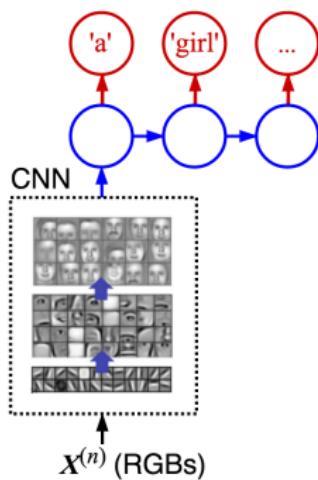
## ⑤ Subword Tokenization

# Limited Representation Size

- In some RNNs, a hidden representation  $\mathbf{a}^{(\cdot,t)}$  needs to support:
  - Current prediction  $\mathbf{a}^{(L,t)}$ , and
  - All** future predictions  $\mathbf{a}^{(L,t+1)}, \dots, \mathbf{a}^{(L,T)}$

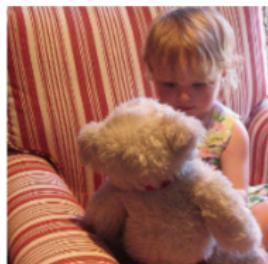


“A little girl sitting on a bed  
with a teddy bear.”

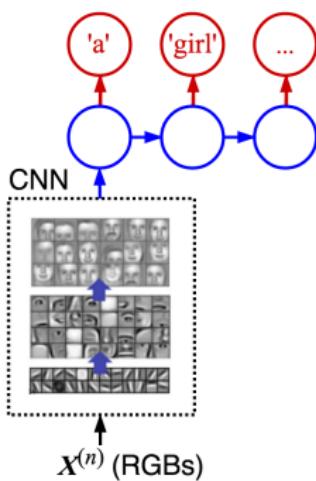


# Limited Representation Size

- In some RNNs, a hidden representation  $\mathbf{a}^{(\cdot,t)}$  needs to support:
  - Current prediction  $\mathbf{a}^{(L,t)}$ , and
  - **All** future predictions  $\mathbf{a}^{(L,t+1)}, \dots, \mathbf{a}^{(L,T)}$
- The fixed-size  $\mathbf{a}^{(t)}$  faces a trade-off between:
  - Representing face features for current prediction ("girl")
  - Representing other features for future predictions ("bear")



"A little girl sitting on a bed  
with a teddy bear."

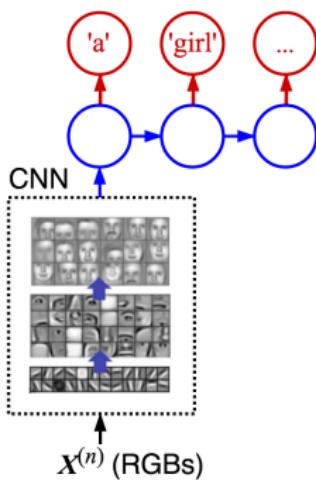


# Limited Representation Size

- In some RNNs, a hidden representation  $\mathbf{a}^{(\cdot,t)}$  needs to support:
  - Current prediction  $\mathbf{a}^{(L,t)}$ , and
  - **All** future predictions  $\mathbf{a}^{(L,t+1)}, \dots, \mathbf{a}^{(L,T)}$
- The fixed-size  $\mathbf{a}^{(t)}$  faces a trade-off between:
  - Representing face features for current prediction ("girl")
  - Representing other features for future predictions ("bear")
- Can we ease the job of  $\mathbf{a}^{(\cdot,t)}$ ?

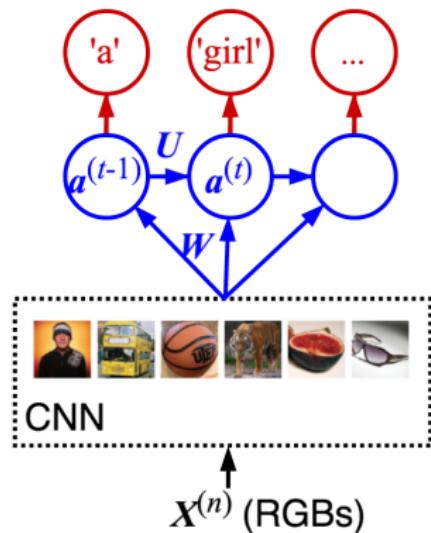


"A little girl sitting on a bed  
with a teddy bear."



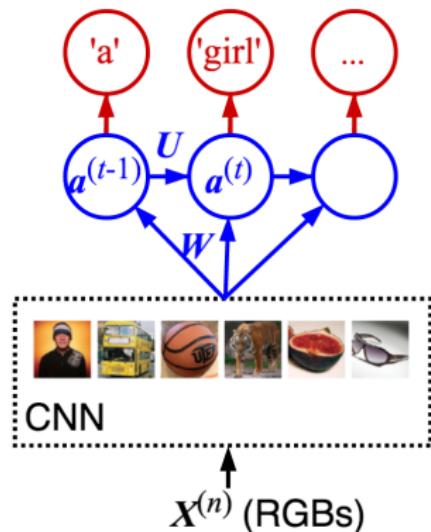
# Input Recurrence

- **Input recurrence**: to feed the entire input  $\mathbf{X}$  to all  $\mathbf{a}^{(1,t)}$ 's,  $\forall t$



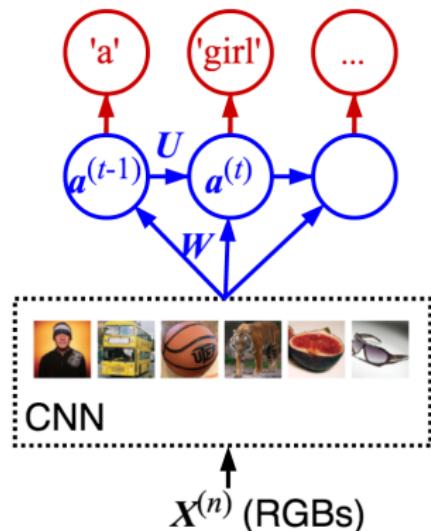
# Input Recurrence

- **Input recurrence**: to feed the entire input  $\mathbf{X}$  to all  $\mathbf{a}^{(1,t)}$ 's,  $\forall t$
- Now,  $\mathbf{a}^{(\cdot,t)}$  only needs to:
  - Support current prediction  $\mathbf{a}^{(L,t)}$ , and
  - Provide context to the **next** representation  $\mathbf{a}^{(\cdot,t+1)}$



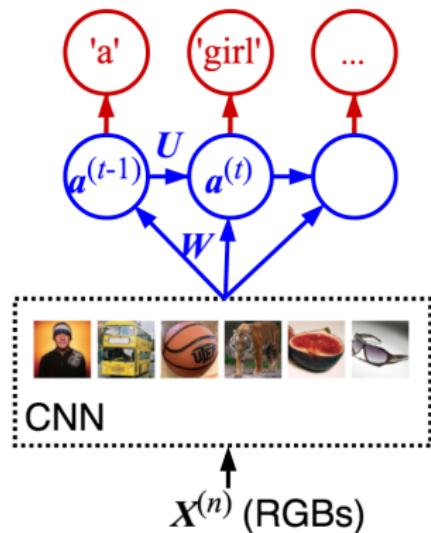
# Input Recurrence

- **Input recurrence**: to feed the entire input  $\mathbf{X}$  to all  $\mathbf{a}^{(1,t)}$ 's,  $\forall t$
- Now,  $\mathbf{a}^{(\cdot,t)}$  only needs to:
  - Support current prediction  $\mathbf{a}^{(L,t)}$ , and
  - Provide context to the **next** representation  $\mathbf{a}^{(\cdot,t+1)}$
- E.g., when predicting “girl,”  $\mathbf{a}^{(\cdot,t)}$  may pay attention to only few face-related images features of current input  $X^{(n)}$



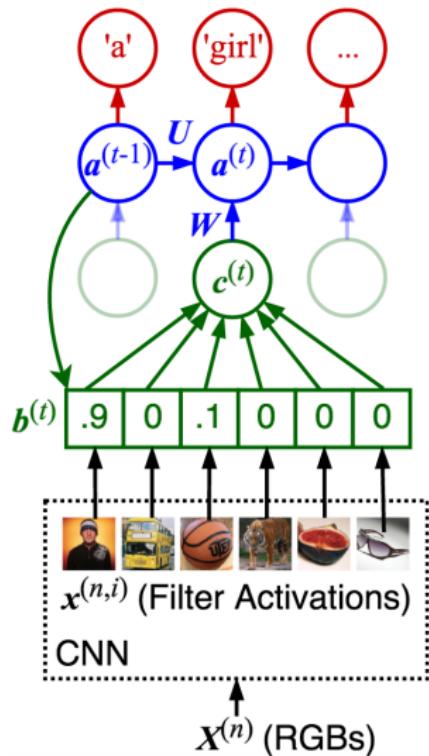
# Input Recurrence

- **Input recurrence**: to feed the entire input  $\mathbf{X}$  to all  $\mathbf{a}^{(1,t)}$ 's,  $\forall t$
- Now,  $\mathbf{a}^{(\cdot,t)}$  only needs to:
  - Support current prediction  $\mathbf{a}^{(L,t)}$ , and
  - Provide context to the **next** representation  $\mathbf{a}^{(\cdot,t+1)}$
- E.g., when predicting “girl,”  $\mathbf{a}^{(\cdot,t)}$  may pay attention to only few face-related images features of current input  $X^{(n)}$
- Why not model the attention explicitly?
  - So we can see where  $\mathbf{a}^{(\cdot,t)}$  is “looking at”



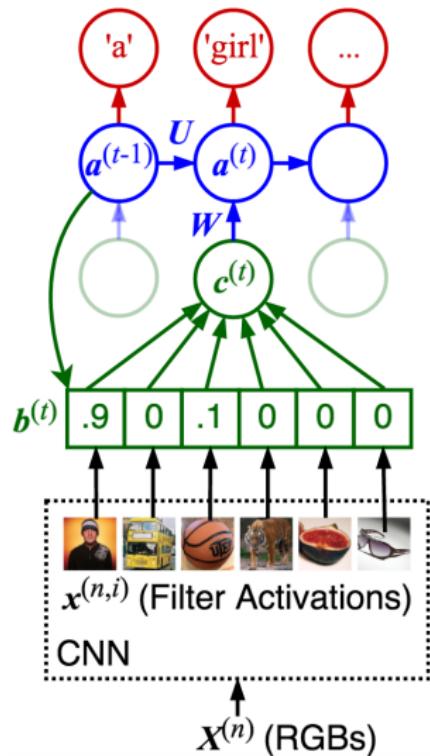
# Attention Mechanism

- Assumes that the input  $\mathbf{X} = \{\mathbf{x}^{(i)}\}_i$  can be broken into “parts”
  - E.g., with CNN,  $\mathbf{x}^{(i)}$  could be the activation values of a filter



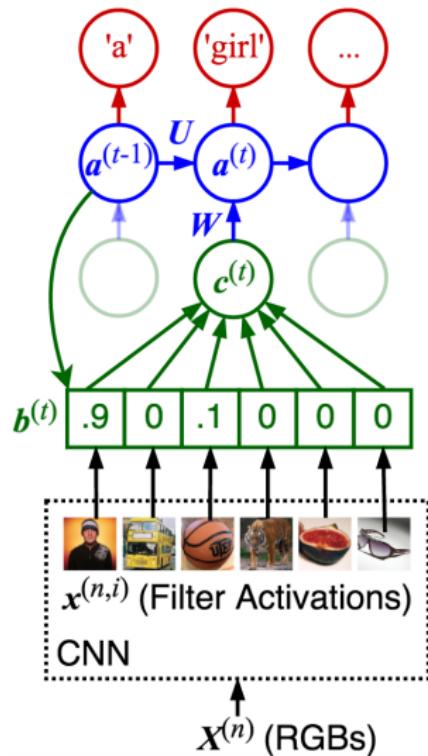
# Attention Mechanism

- Assumes that the input  $\mathbf{X} = \{\mathbf{x}^{(i)}\}_i$  can be broken into “parts”
  - E.g., with CNN,  $\mathbf{x}^{(i)}$  could be the activation values of a filter
- For each timestamp  $t$ :
  - Obtain from  $\mathbf{a}^{(L-1,t-1)}$  an **attention vector**  $\mathbf{b}^{(t)}$ ,  $\sum b_i^{(t)} = 1$



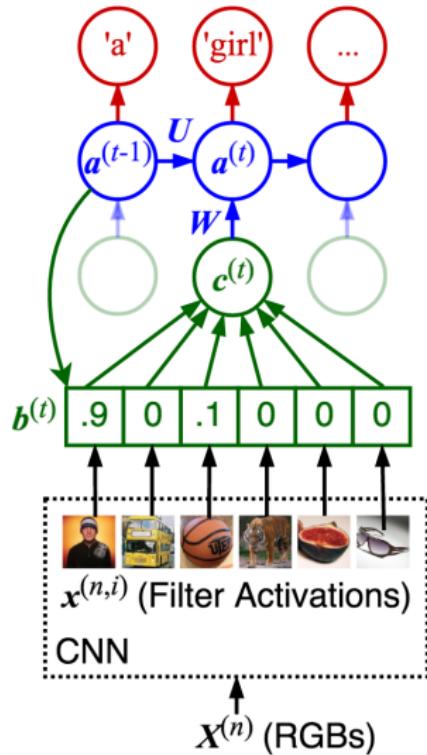
# Attention Mechanism

- Assumes that the input  $\mathbf{X} = \{\mathbf{x}^{(i)}\}_i$  can be broken into “parts”
  - E.g., with CNN,  $\mathbf{x}^{(i)}$  could be the activation values of a filter
- For each timestamp  $t$ :
  - Obtain from  $\mathbf{a}^{(L-1,t-1)}$  an **attention vector**  $\mathbf{b}^{(t)}$ ,  $\sum b_i^{(t)} = 1$
  - Feed  $\mathbf{a}^{(1,t)}$  with the weighted input  $\mathbf{c}^{(t)} = \sum_i b_i^{(t)} \mathbf{x}^{(i)}$



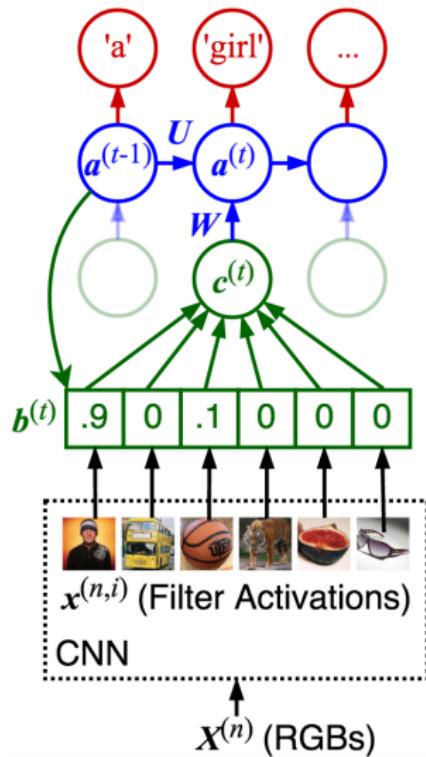
# Attention Mechanism

- Assumes that the input  $\mathbf{X} = \{\mathbf{x}^{(i)}\}_i$  can be broken into “parts”
  - E.g., with CNN,  $\mathbf{x}^{(i)}$  could be the activation values of a filter
- For each timestamp  $t$ :
  - Obtain from  $\mathbf{a}^{(L-1,t-1)}$  an **attention vector**  $\mathbf{b}^{(t)}$ ,  $\sum b_i^{(t)} = 1$
  - Feed  $\mathbf{a}^{(1,t)}$  with the weighted input  $\mathbf{c}^{(t)} = \sum_i b_i^{(t)} \mathbf{x}^{(i)}$
- Reduces size of  $\mathbf{W}$  from  $O(|\mathbf{X}| \cdot |\mathbf{a}^{(\cdot,t)}|)$  to  $O(|\mathbf{x}^{(i)}| \cdot |\mathbf{a}^{(\cdot,t)}|)$



# Attention Mechanism

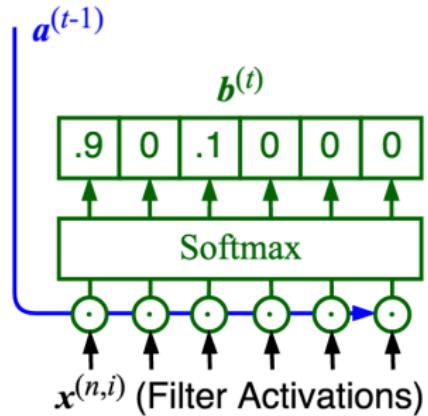
- Assumes that the input  $\mathbf{X} = \{\mathbf{x}^{(i)}\}_i$  can be broken into “parts”
  - E.g., with CNN,  $\mathbf{x}^{(i)}$  could be the activation values of a filter
- For each timestamp  $t$ :
  - Obtain from  $\mathbf{a}^{(L-1,t-1)}$  an **attention vector**  $\mathbf{b}^{(t)}$ ,  $\sum b_i^{(t)} = 1$
  - Feed  $\mathbf{a}^{(1,t)}$  with the weighted input  $\mathbf{c}^{(t)} = \sum_i b_i^{(t)} \mathbf{x}^{(i)}$
- Reduces size of  $\mathbf{W}$  from  $O(|\mathbf{X}| \cdot |\mathbf{a}^{(\cdot,t)}|)$  to  $O(|\mathbf{x}^{(i)}| \cdot |\mathbf{a}^{(\cdot,t)}|)$
- How to obtain  $\mathbf{b}^{(t)}$ ?



# Computing Attention Vector

- ① Use  $\mathbf{a}^{(L-1,t-1)}$  as a “query” to get a match score for each input part by using, e.g., a simple NN [2, 20]:

$$z_i = \text{act}(\mathbf{p}^\top \mathbf{a}^{(L-1,t-1)} + \mathbf{q}^\top \mathbf{x}^{(i)} + r)$$



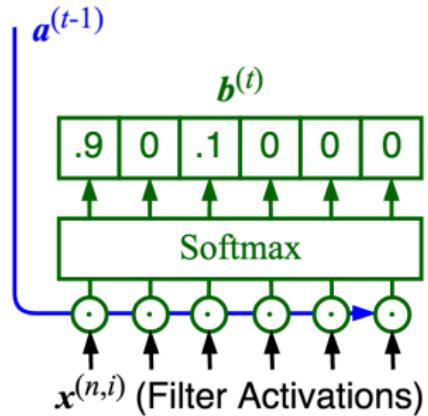
# Computing Attention Vector

- ① Use  $\mathbf{a}^{(L-1,t-1)}$  as a “query” to get a match score for each input part by using, e.g., a simple NN [2, 20]:

$$z_i = \text{act}(\mathbf{p}^\top \mathbf{a}^{(L-1,t-1)} + \mathbf{q}^\top \mathbf{x}^{(i)} + r)$$

- ② Normalize and concentrate on few larger scores by:

$$b_i = \text{softmax}(z)_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$



# Computing Attention Vector

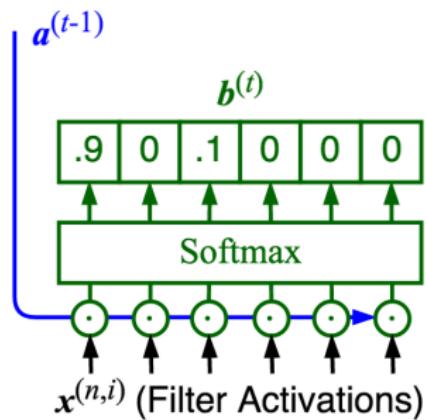
- ① Use  $\mathbf{a}^{(L-1,t-1)}$  as a “query” to get a match score for each input part by using, e.g., a simple NN [2, 20]:

$$z_i = \text{act}(\mathbf{p}^\top \mathbf{a}^{(L-1,t-1)} + \mathbf{q}^\top \mathbf{x}^{(i)} + r)$$

- Jointly trained with the main RNN

- ② Normalize and concentrate on few larger scores by:

$$b_i = \text{softmax}(z)_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$



# Computing Attention Vector

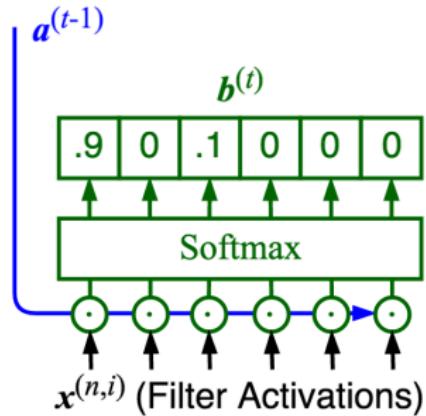
- ① Use  $\mathbf{a}^{(L-1,t-1)}$  as a “query” to get a match score for each input part by using, e.g., a simple NN [2, 20]:

$$z_i = \text{act}(\mathbf{p}^\top \mathbf{a}^{(L-1,t-1)} + \mathbf{q}^\top \mathbf{x}^{(i)} + r)$$

- Jointly trained with the main RNN
- $\mathbf{p}$ ,  $\mathbf{q}$ , and  $r$  are shared by different  $i$ 's and  $t$ 's (weight tying)

- ② Normalize and concentrate on few larger scores by:

$$b_i = \text{softmax}(z)_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$



# Visualizing Attention



sitting(0.29)



AI(0.99)



little(0.47)



girl(0.35)



on(0.23)



a(0.23)



bed(0.40)



with(0.27)



a(0.15)



teddy(0.31)



bear(0.24)

- How to draw a mask?

# Visualizing Attention



sitting(0.29)



on(0.23)



a(0.23)



bed(0.40)



a(0.15)



teddy(0.31)



- How to draw a mask? Threat  $\mathbf{c}^{(t)} = \sum_i \mathbf{b}_i^{(t)} \mathbf{x}^{(i)}$  as image and enlarge it

# Outline

## ① RNNs

- Vanilla RNNs
- Design Alternatives

## ② RNN Training

- Backprop through Time (BPTT)
- Optimization Techniques
- Optimization-Friendly Models & LSTM
- Parallelism & Teacher Forcing

## ③ RNNs with Attention Mechanism

- Attention for Image Captioning
- Attention for Neural Machine Translation (NMT)

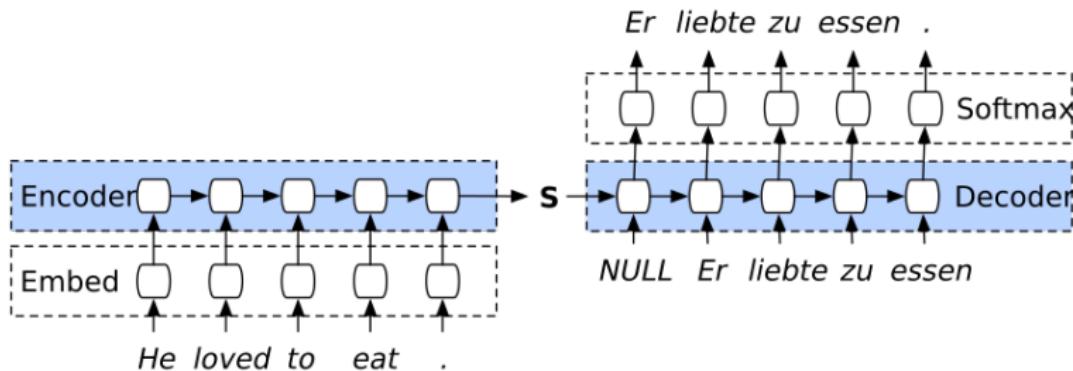
## ④ Transformers

- Attention is All You Need
- Pretrained Language Models
- More Applications

## ⑤ Subword Tokenization

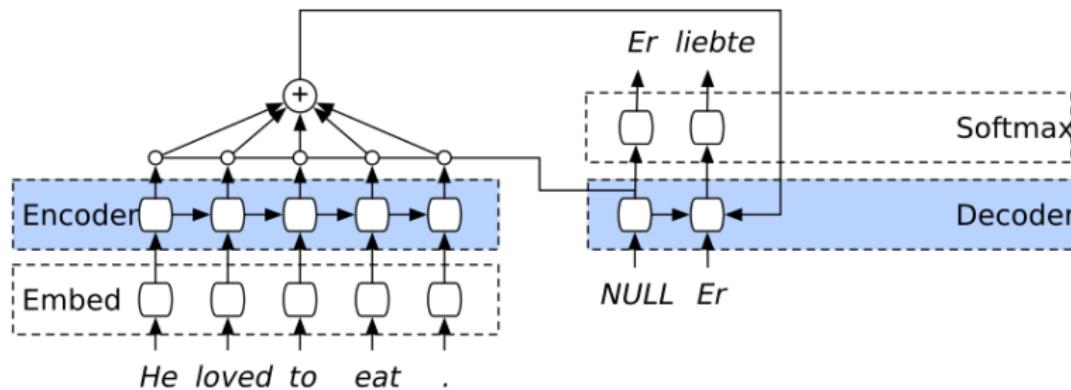
# GoogleNMTv1: Encoder-Decoder

- LSTM-based RNN
- Hidden state  $S$  encodes an entire input sequence
- Then supplied to the decoder



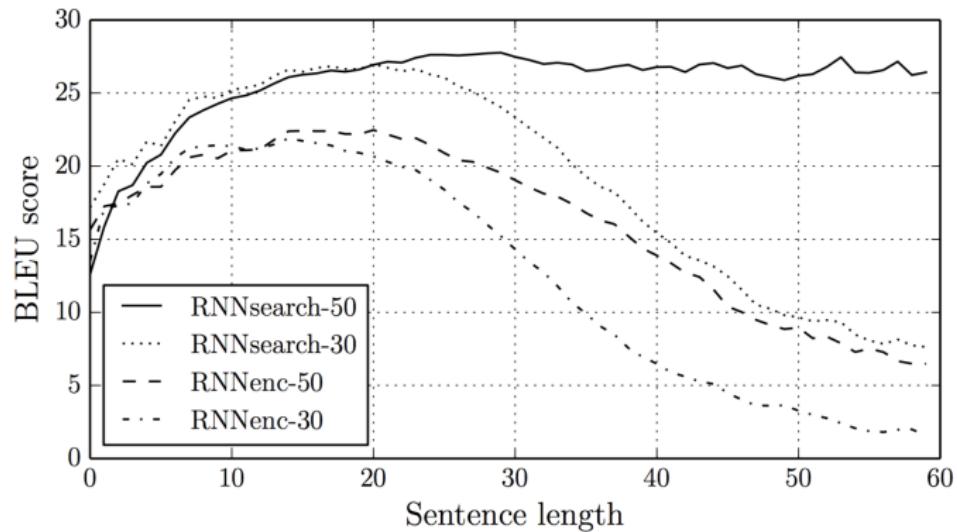
# GoogleNMTv2: Attention-based Encoder-Decoder

- No feed-forward connection between the encoder and decoder
- Instead, uses previous output as query to get attention and next output
  - Allows for retrieving different parts of input sentence depending on decoding context



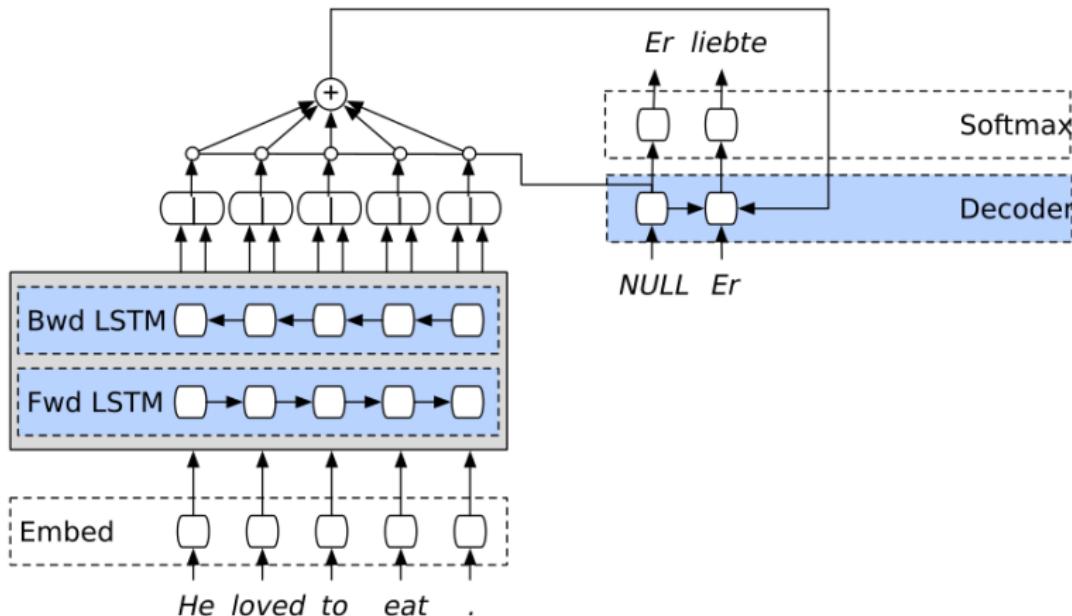
# Long Sequences

- Attention-based model generates long sequences better



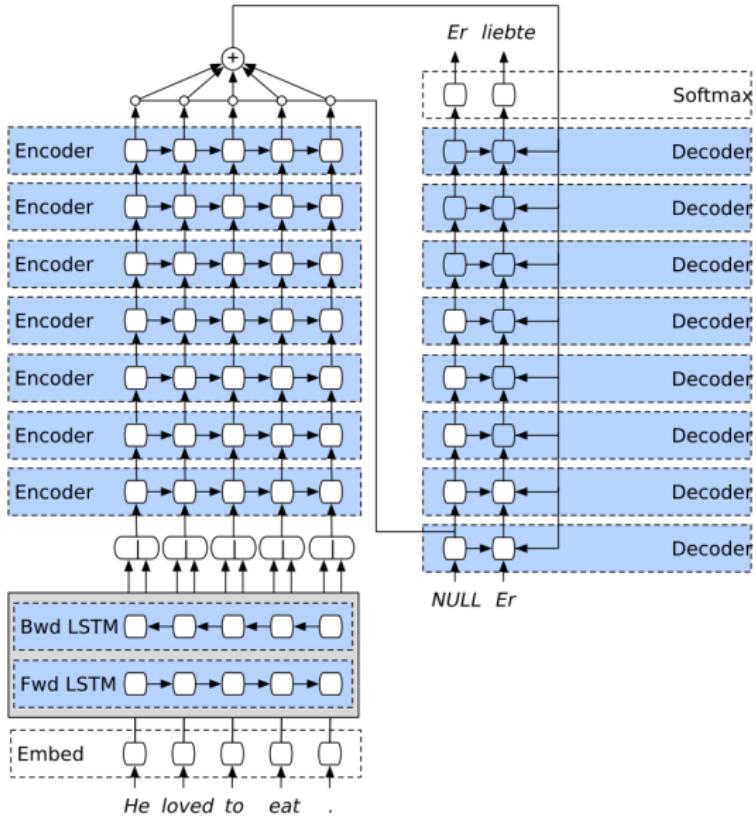
# GoogleNMTv3: Bidirectional Encoder Layer

- Takes into account future words when summarizing input sequence
- Better determines the meaning/context



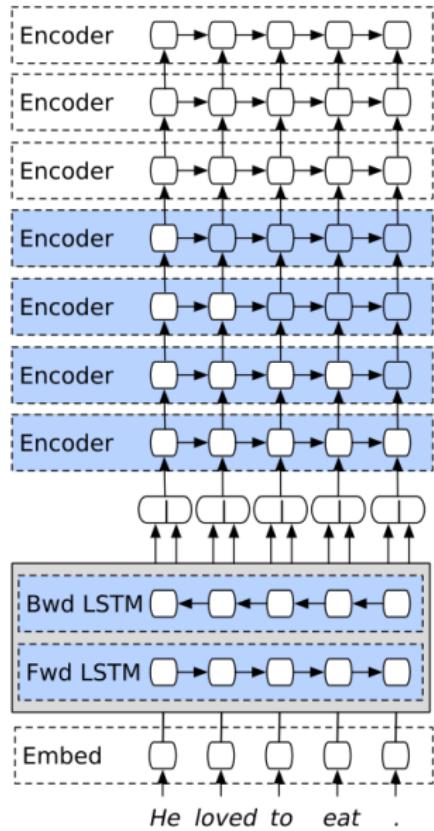
# GoogleNMTv4: Going Deep

- Encoder:
  - 1 bi-directional layer
  - 7 uni-directional layers
- Decoder:
  - 8 uni-directional layers
  - Lowest** decoder layer for querying attention

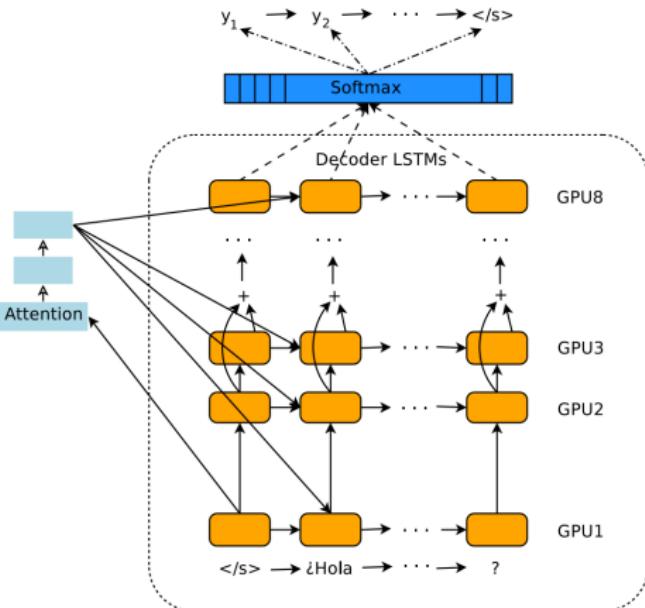
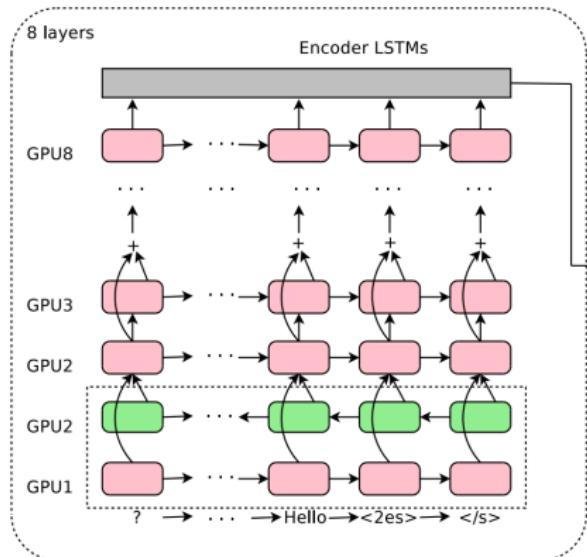


# GoogleNMTv5: Parallelization

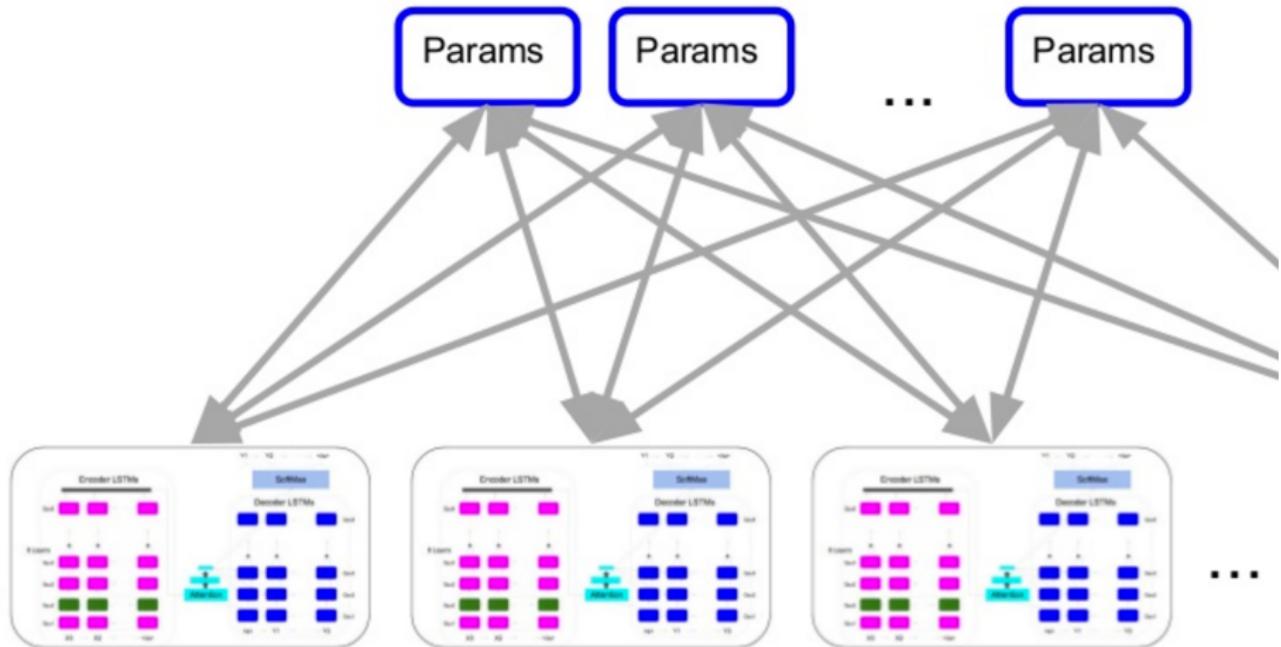
- Each layer trained by a GPU
- Forward pass:
  - Encoder: 7 uni-directional encoder layers trained in a pipeline
  - Decoder: pipeline starts as soon as encoder layers are ready
- Backward pass:
  - Teacher forcing



# Model Parallelism (1 Machine, 8 GPUs)

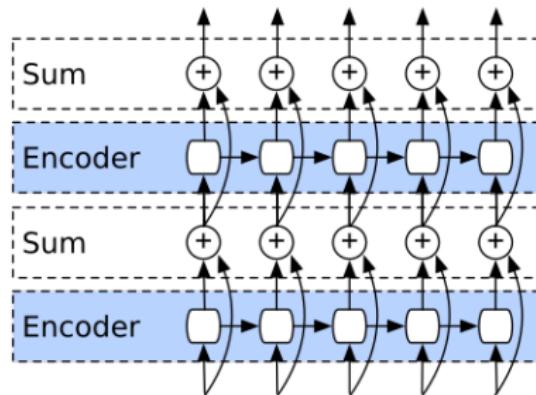


# Data Parallelism (Multiple Param Servers)



# GoogleNMTv6: Residuals

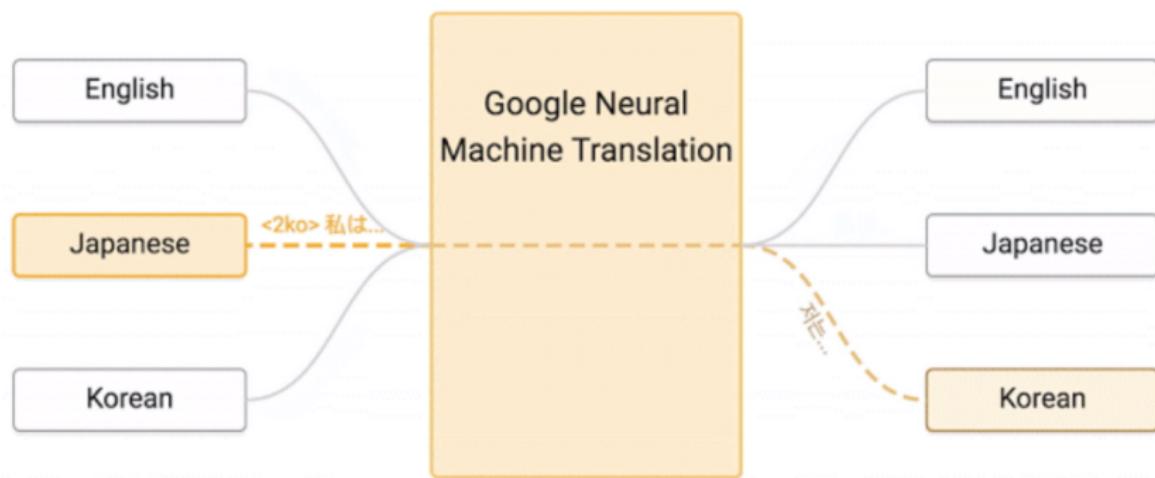
- Upper layer learns the *delta* function to the lower one
- Easier to train a deep NN



# GoogleNMTv7: Multilingual & Zero-Shot Translation

- Training: input  $x$  augmented with task identifier (language pair)
  - E.g., (eng, jp), (kr, en), (jp, en)
- Inference: **unknown** language-pair identifier
  - E.g., (kr, jp)

Zero-shot



# Outline

## ① RNNs

- Vanilla RNNs
- Design Alternatives

## ② RNN Training

- Backprop through Time (BPTT)
- Optimization Techniques
- Optimization-Friendly Models & LSTM
- Parallelism & Teacher Forcing

## ③ RNNs with Attention Mechanism

- Attention for Image Captioning
- Attention for Neural Machine Translation (NMT)

## ④ Transformers

- Attention is All You Need
- Pretrained Language Models
- More Applications

## ⑤ Subword Tokenization

# How Humans Process Text?

研表究明  
漢字的序順並不一定能影響讀  
比如當你看完這句話後  
才發這現裡的字全都是亂的

Aoccdrnig to a rscheearch at an Elingsh uinervtisy,  
it deosn't mttaer in waht oredr the ltteers in a wrod  
are, the olny iprmoetnt tihg is taht frist and lsat  
ltteer is at the rghit pclae. The rset can be a toatl  
mses and you can stil raed it wouthit porbelm. Tih  
is bcuseae we do not raed ervey lteter by itslef but  
the wrod as a whole.

# How Humans Process Text?

研表究明  
漢字的序順並不一定能影響讀  
比如當你看完這句話後  
才發這現裡的字全都是亂的

Aoccdrnig to a rscheearch at an Elingsh uinervtisy,  
it deosn't mttaer in waht oredr the ltteers in a wrod  
are, the olny iprmoetnt tihng is taht frist and lsat  
ltteer is at the rghit pclae. The rset can be a toatl  
mses and you can stil raed it wouthit porbelm. Tihs  
is bcuseae we do not raed ervey lteter by itslef but  
the wrod as a whole.

- *Not entirely sequential* as in RNNs
- *Not recursively based on local patterns* as in CNNs
- Any other “better” architectures?

# Outline

## ① RNNs

- Vanilla RNNs
- Design Alternatives

## ② RNN Training

- Backprop through Time (BPTT)
- Optimization Techniques
- Optimization-Friendly Models & LSTM
- Parallelism & Teacher Forcing

## ③ RNNs with Attention Mechanism

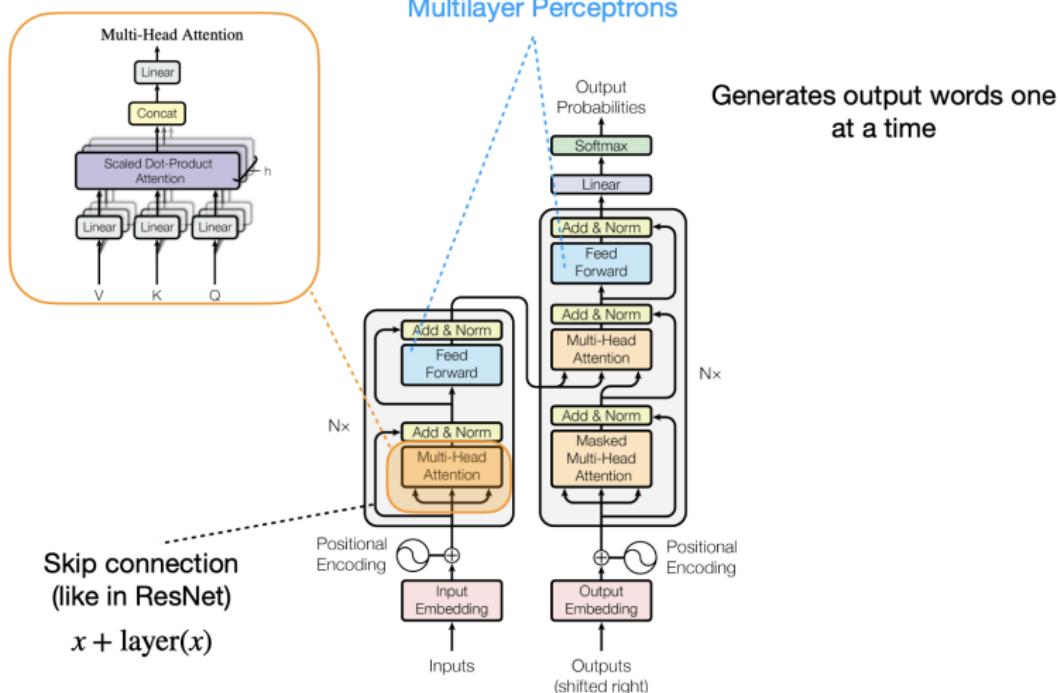
- Attention for Image Captioning
- Attention for Neural Machine Translation (NMT)

## ④ Transformers

- Attention is All You Need
- Pretrained Language Models
- More Applications

## ⑤ Subword Tokenization

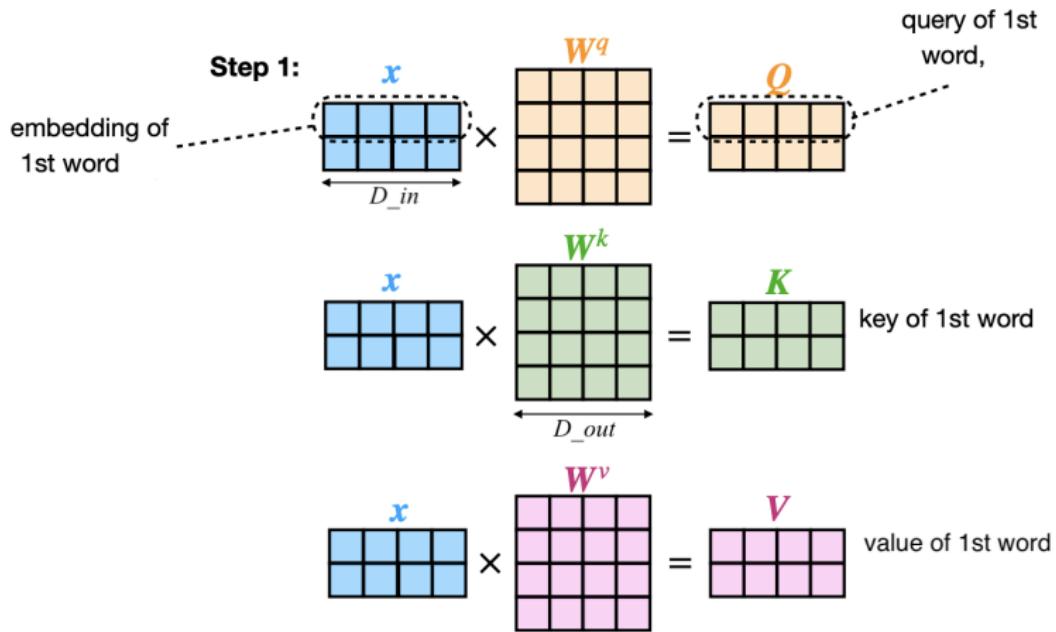
# Attention is All You Need [19]



- RNN layers are replaced with **self-** and **cross-attention** blocks

# Self-Attention

- Weights to learn at each layer:  $\mathbf{W}_{\text{query}}, \mathbf{W}_{\text{key}}, \mathbf{W}_{\text{value}} \in \mathbb{R}^{D \times D}$
- Batch, non-autoregressive processing of sequences



# Self-Attention

**Step 2:**  $Q \times K^T = QK^T$

key of 1st word,  $k_1$

relationship between 1st & 2nd word

$Q$

$K^T$

$QK^T$

**Step 3:**  $QK^T / \sqrt{D} = QK^T / \sqrt{D}$

$QK^T$

$/ \sqrt{D}$

$QK^T / \sqrt{D}$

**Step 4:** softmax  $(QK^T / \sqrt{D}) =$

attention of 1st word on values

softmax  $(QK^T / \sqrt{D})$

1st output value of 1st word

**Step 5:**  $\text{softmax}(QK^T / \sqrt{D}) \times V = A$

$V$

$A$

# Self-Attention

Step 2:

$$\begin{matrix} Q \\ \text{---} \\ \boxed{\text{---}} \end{matrix} \times \begin{matrix} K^T \\ \text{---} \\ \boxed{\text{---}} \end{matrix} = \begin{matrix} \text{key of 1st word, } k_1 \\ \text{---} \\ \boxed{\text{---}} \end{matrix}$$

relationship between  
1st & 2nd word

Step 3:

$$\begin{matrix} QK^T \\ \text{---} \\ \boxed{\text{---}} \end{matrix} / \sqrt{D} = \begin{matrix} QK^T / \sqrt{D} \\ \text{---} \\ \boxed{\text{---}} \end{matrix}$$

Step 4: softmax  $(QK^T / \sqrt{D}) = \begin{matrix} \text{---} \\ \boxed{\text{---}} \end{matrix}$

- Input:  $A^{(l-1)} \in \mathbb{R}^{T \times D}$
- Output:  $A^{(l)} \in \mathbb{R}^{T \times D}$ 
  - **Context** augmented
  - E.g, “**train** a model” vs. “get on a **train**”

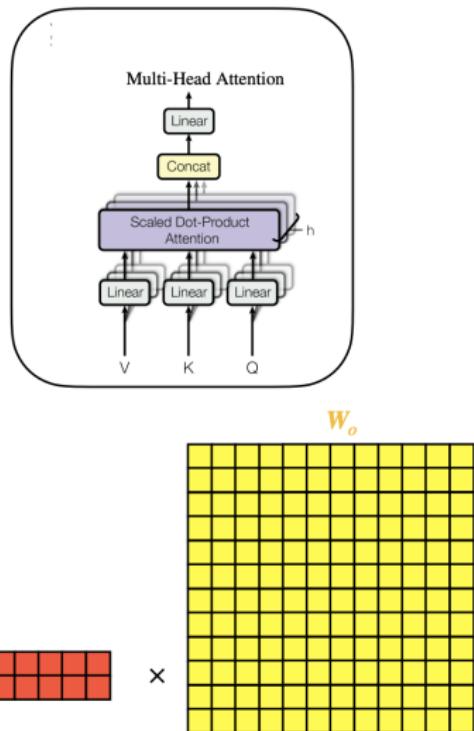
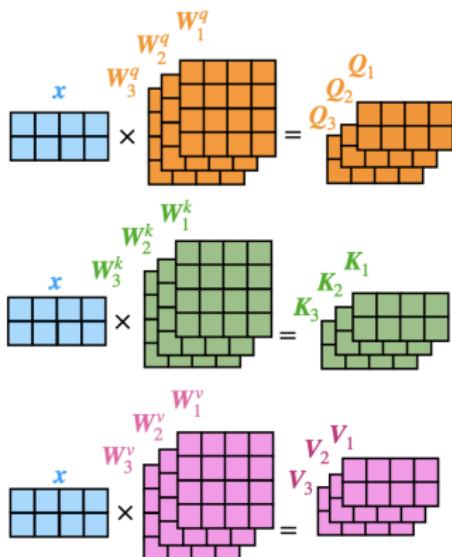
attention of 1st word on values

softmax  $(QK^T / \sqrt{D})$

Step 5:

$$\begin{matrix} \text{---} \\ \boxed{\text{---}} \end{matrix} \times \begin{matrix} V \\ \text{---} \\ \boxed{\text{---}} \end{matrix} = \begin{matrix} \text{1st output value of 1st word} \\ \text{---} \\ \boxed{A} \end{matrix}$$

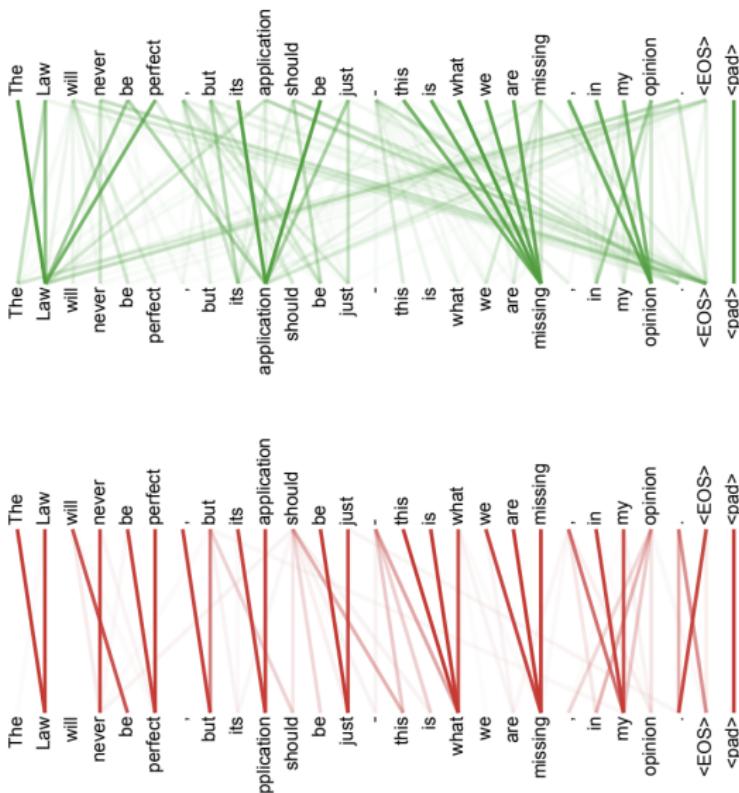
# Multi-head Self-Attention



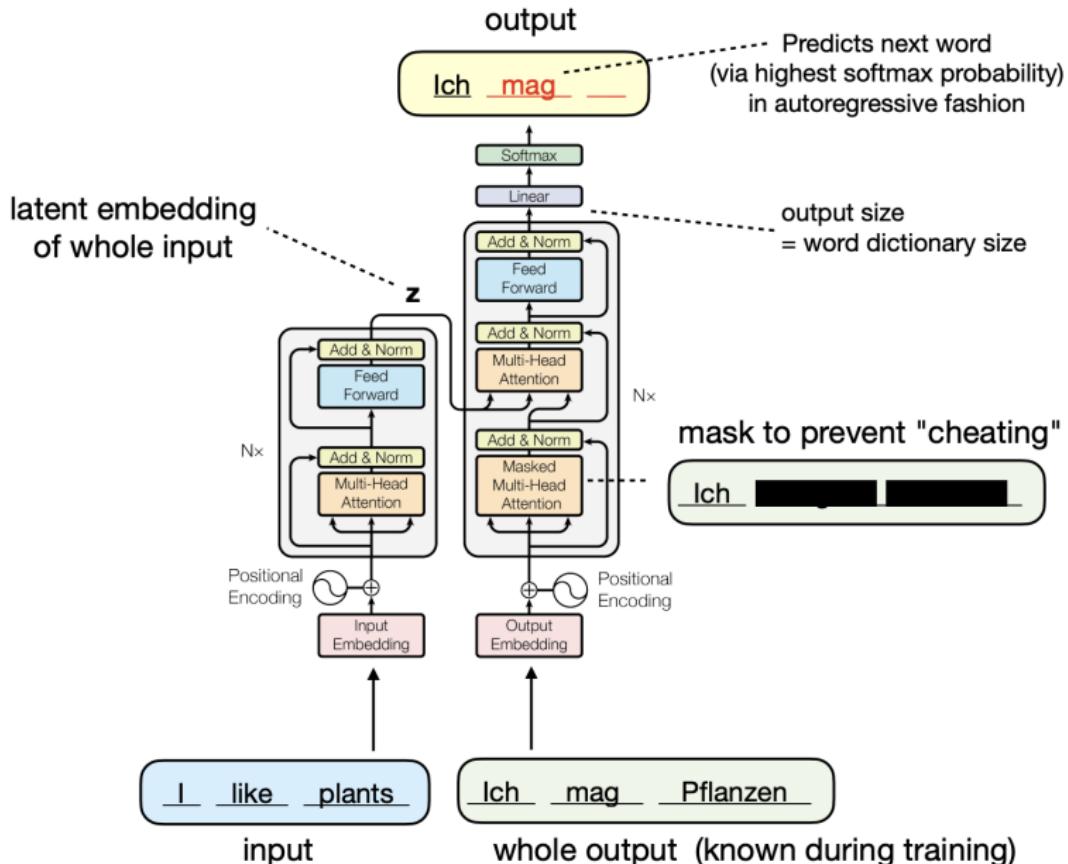
- Given  $H$  heads, we have  $W_{\text{query}}^{(h)}, W_{\text{key}}^{(h)}, W_{\text{value}}^{(h)} \in \mathbb{R}^{D \times \frac{D}{H}}$  and  $A \in \mathbb{R}^{T \times D}$

# Same Sequence, Different Context Augmentations

Step 4: softmax  $(QK^T / \sqrt{D}) = \boxed{\quad}$

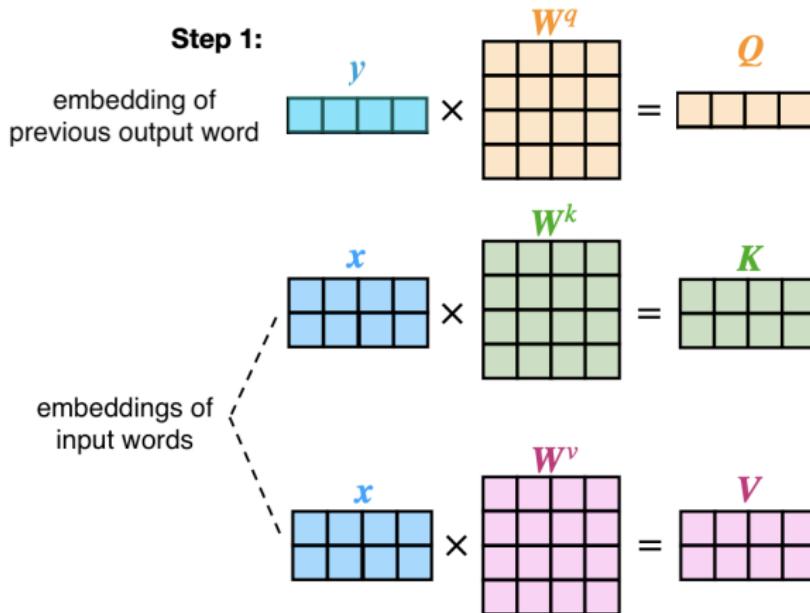


# Masked Self-Attention for Decoding Output

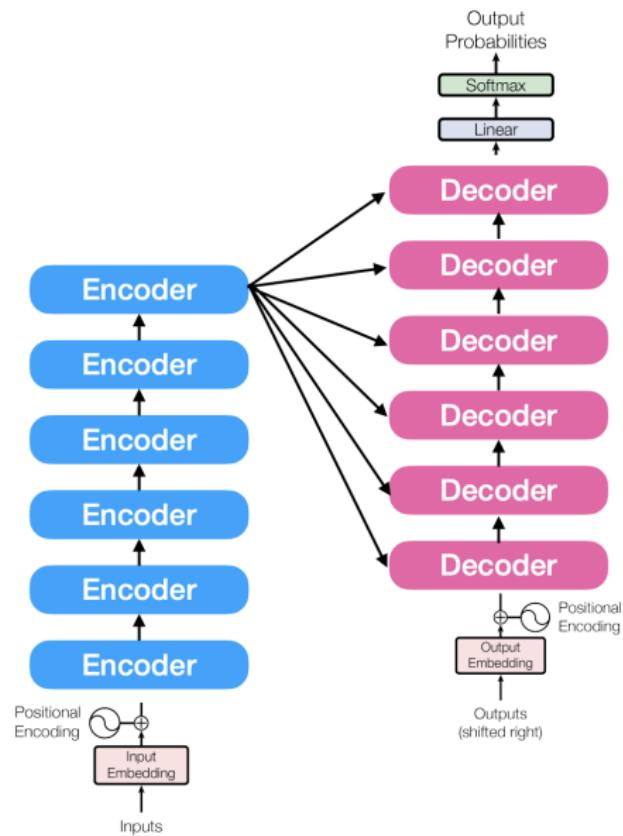


# Cross-Attention

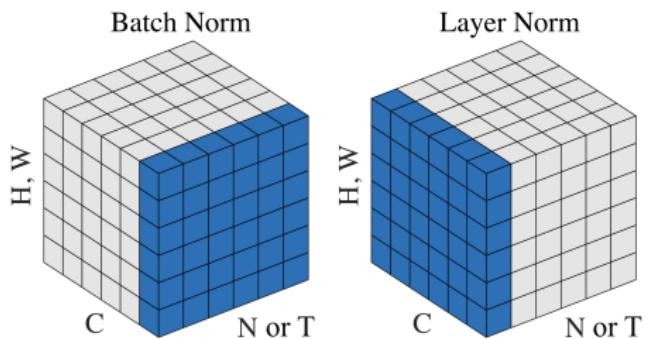
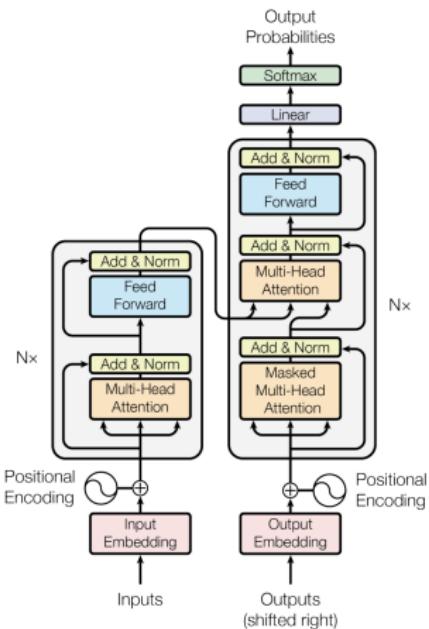
- Users previous output as query (autoregressive)
- Input as context



# Going Deep

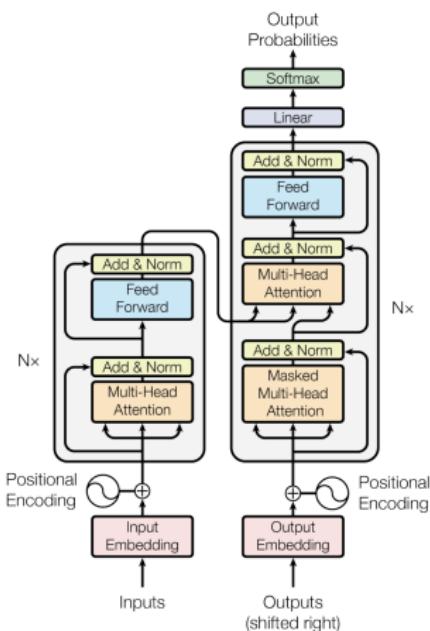


# Layer Normalization



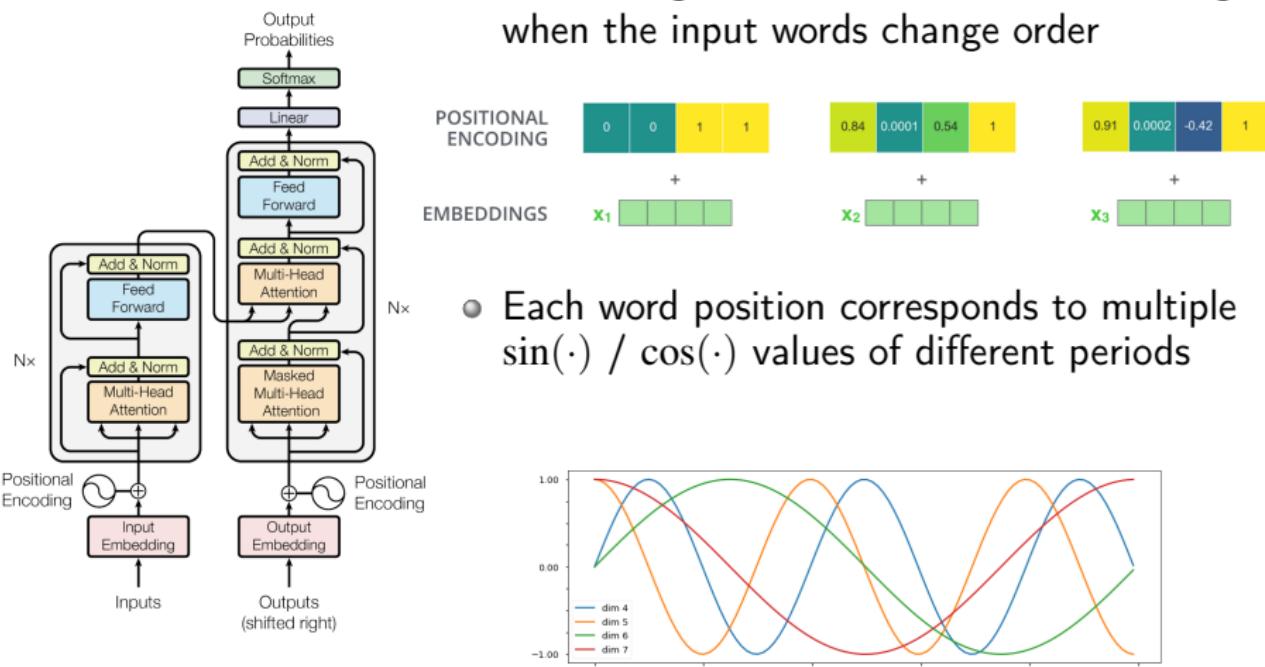
# Positional Encoding

- “John loves Mary”  $\neq$  “Mary loves John”
- So far, augmented context does **not** change when the input words change order



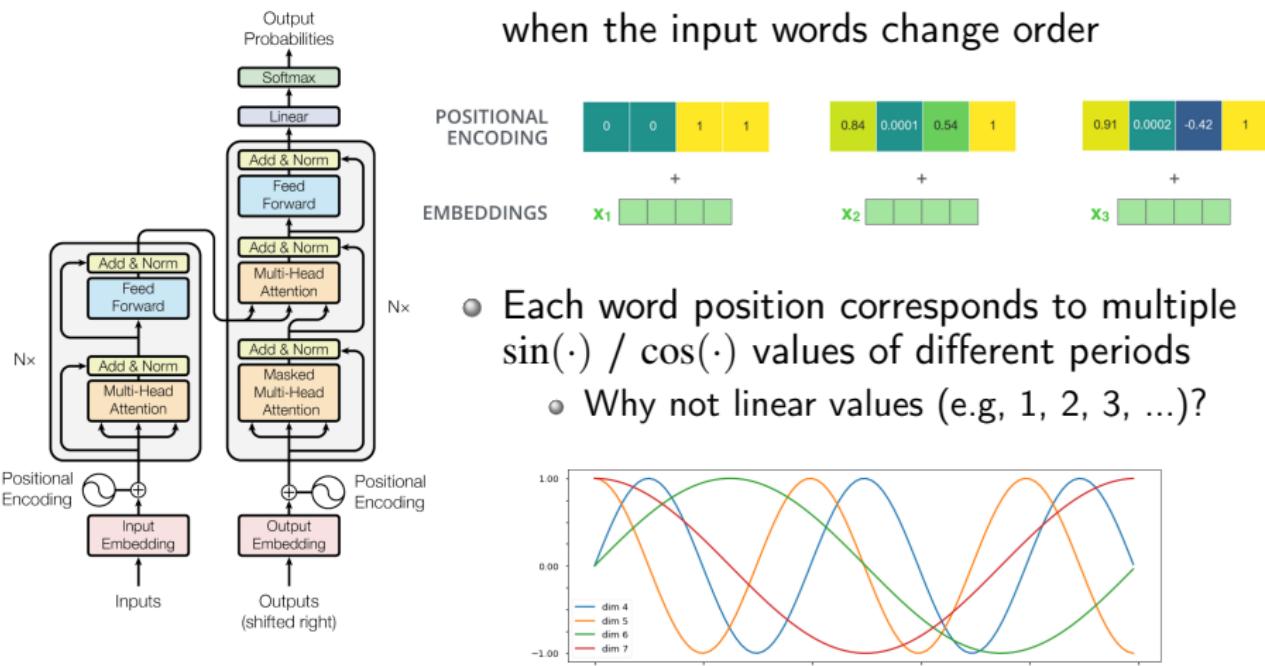
# Positional Encoding

- “John loves Mary”  $\neq$  “Mary loves John”
- So far, augmented context does **not** change when the input words change order



# Positional Encoding

- “John loves Mary”  $\neq$  “Mary loves John”
- So far, augmented context does **not** change when the input words change order



# RNNs vs. CNNs vs. Transformers

- On processing a sequence of length  $T$  at each layer with
  - $D$ -dimensional point input and output
  - $F$  = the CNN filter/kernel size
  - #CNN filters =  $D$
  - #attention heads =  $H$
  - Query, key, and value size =  $\frac{D}{H}$

	#Weights	Computation	Auto-reg.	Point Dist.
CNN	$O(FD^2)$	$O(TFD^2)$	No	$O(\frac{T}{F})$
RNN	$O(D^2)$	$O(TD^2)$	Yes	$O(T)$
Self-attention	$O(D^2)$	$O(TD^2 + T^2D)$	No	$O(1)$

# Outline

## ① RNNs

- Vanilla RNNs
- Design Alternatives

## ② RNN Training

- Backprop through Time (BPTT)
- Optimization Techniques
- Optimization-Friendly Models & LSTM
- Parallelism & Teacher Forcing

## ③ RNNs with Attention Mechanism

- Attention for Image Captioning
- Attention for Neural Machine Translation (NMT)

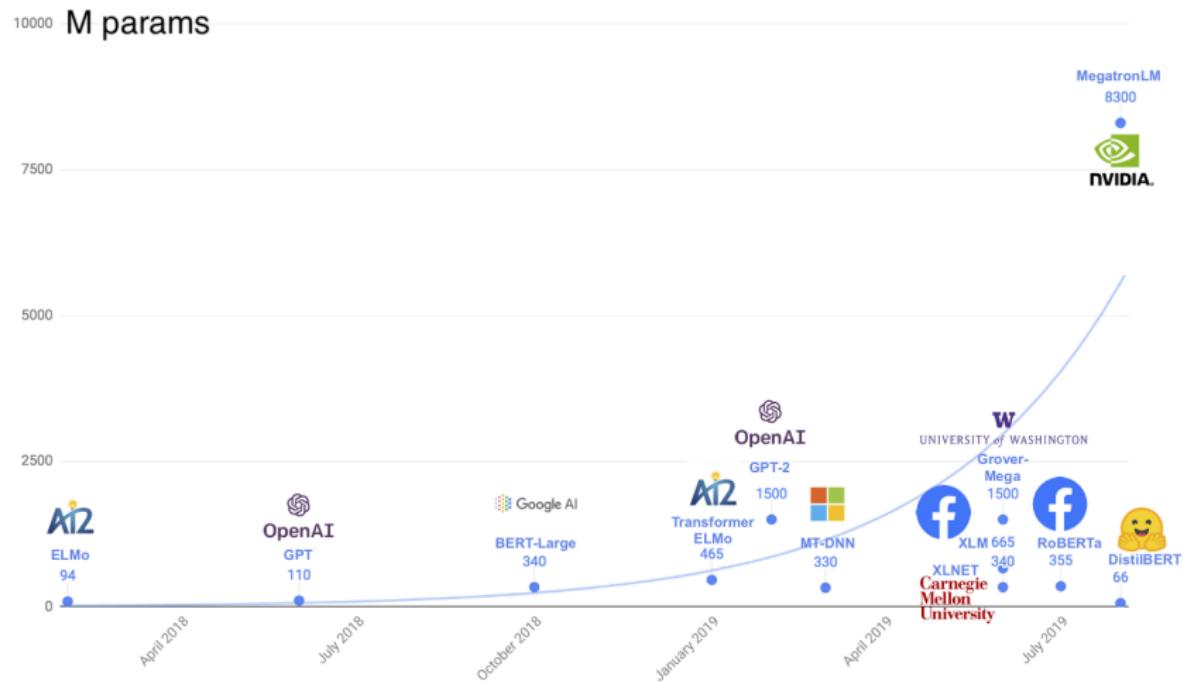
## ④ Transformers

- Attention is All You Need
- Pretrained Language Models
- More Applications

## ⑤ Subword Tokenization

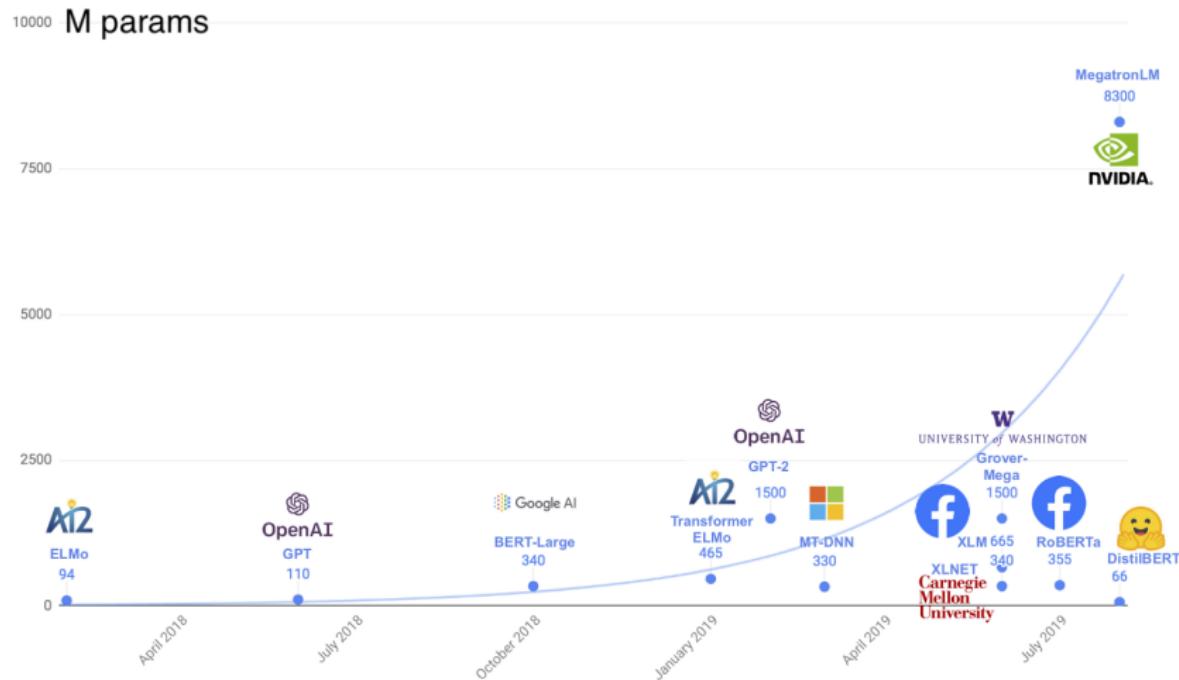
# Recent Growth of Language Transformers

- All based on *self-supervised pre-training*



# Recent Growth of Language Transformers

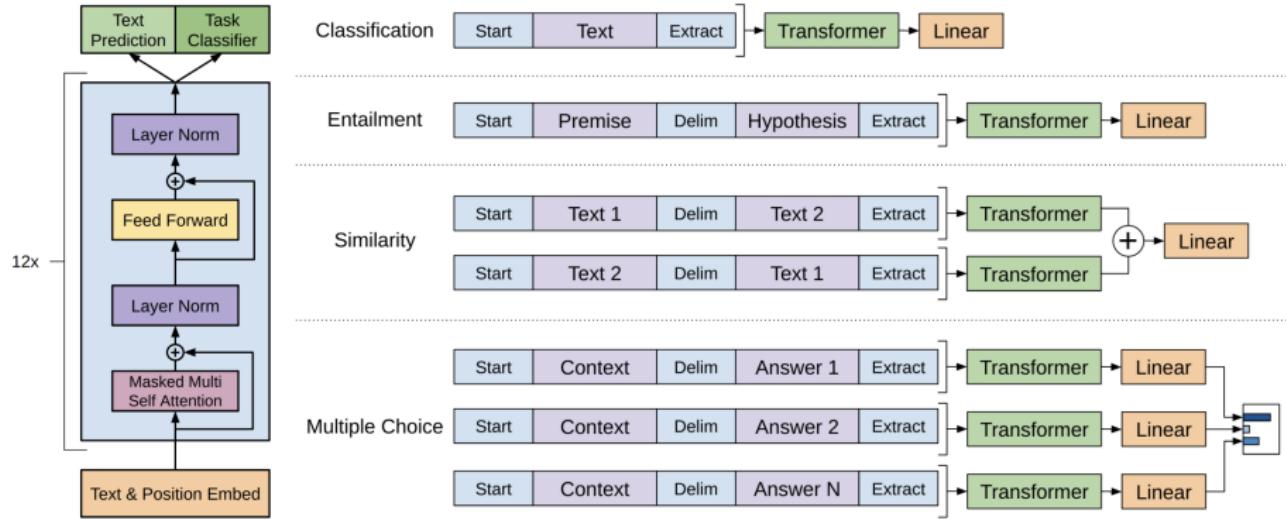
- All based on *self-supervised pre-training*
- Costs [17]: \$2.5k–\$50k (110M params), \$10k–\$200k (340M), \$80k–\$1.6m (1.5B)



# GPTv1 [12]

- 2-step training process:
  - ① Generative pre-training on unlabeled data
    - to predict next word in a sentence (language model)
  - ② Discriminative fine-tuning on labeled data
    - Here, fine-tuning could also mean training a new model based on extracted features
- Architecture based on the **decoder** of original transformer (“Attention is All You Need”)
  - Autoregressive

# Architecture & Fine-tuning Tasks



- Fine-tuning tasks: whether to minimize language model (LM) loss or not?

# Ablation Study

Table 5: Analysis of various model ablations on different tasks. Avg. score is a unweighted average of all the results. (*mc*= Mathews correlation, *acc*=Accuracy, *pc*=Pearson correlation)

Method	Avg. Score	CoLA (mc)	SST2 (acc)	MRPC (F1)	STS-B (pc)	QQP (F1)	MNLI (acc)	QNLI (acc)	RTE (acc)
Transformer w/ aux LM (full)	74.7	45.4	91.3	82.3	82.0	<b>70.3</b>	<b>81.8</b>	<b>88.1</b>	<b>56.0</b>
Transformer w/o pre-training	59.9	18.9	84.0	79.4	30.9	65.5	75.7	71.2	53.8
Transformer w/o aux LM	<b>75.0</b>	<b>47.9</b>	<b>92.0</b>	<b>84.9</b>	<b>83.2</b>	69.8	81.1	86.9	54.4
LSTM w/ aux LM	69.1	30.3	90.5	83.2	71.8	68.1	73.7	81.1	54.6

# BERT [4]

- Architecture based on the *encoder* of original transformer (“Attention is All You Need”)
  - Non-autoregressive
- Pre-training tasks:
  - Masked language model (“Cloze” task)
    - “A quick brown [MASK] jumps over the lazy dog” → “fox” 11%, “ant” 5%, ...
  - Next sentence prediction
    - “[MASK] go to store [SEP] to buy a [MASK] of milk” → True 93%, False 7%

# Pre-training & Fine-tuning Tasks

- Special input: [CLS] token

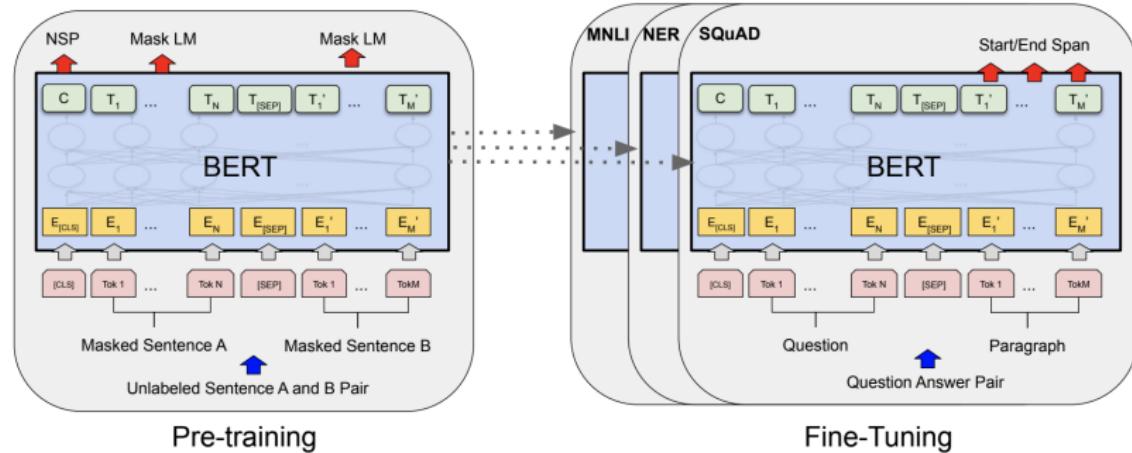


Figure 1: Overall pre-training and fine-tuning procedures for BERT. Apart from output layers, the same architectures are used in both pre-training and fine-tuning. The same pre-trained model parameters are used to initialize models for different down-stream tasks. During fine-tuning, all parameters are fine-tuned. [CLS] is a special symbol added in front of every input example, and [SEP] is a special separator token (e.g. separating questions/answers).

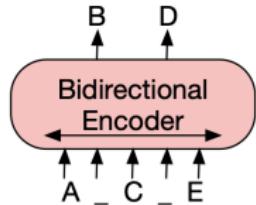
# BART [11]

- BERT: bidirectional, non-autoregressive (autoencoder like)
  - Good for embedding tasks
- GPT: unidirectional, autoregressive
  - Good for generative tasks

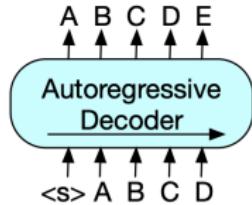
# BART [11]

- BERT: bidirectional, non-autoregressive (autoencoder like)
  - Good for embedding tasks
- GPT: unidirectional, autoregressive
  - Good for generative tasks
- BART: why not combine both?

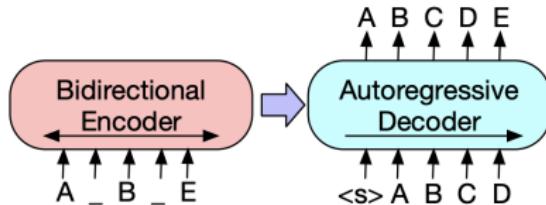
# Architecture



(a) BERT: Random tokens are replaced with masks, and the document is encoded bidirectionally. Missing tokens are predicted independently, so BERT cannot easily be used for generation.

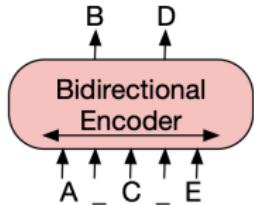


(b) GPT: Tokens are predicted auto-regressively, meaning GPT can be used for generation. However words can only condition on leftward context, so it cannot learn bidirectional interactions.

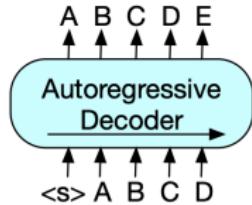


(c) BART: Inputs to the encoder need not be aligned with decoder outputs, allowing arbitrary noise transformations. Here, a document has been corrupted by replacing spans of text with mask symbols. The corrupted document (left) is encoded with a bidirectional model, and then the likelihood of the original document (right) is calculated with an autoregressive decoder. For fine-tuning, an uncorrupted document is input to both the encoder and decoder, and we use representations from the final hidden state of the decoder.

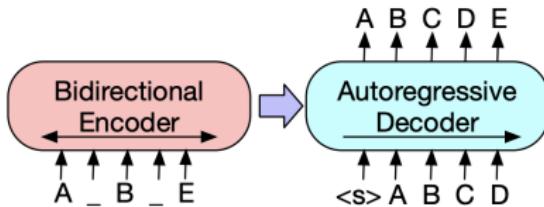
# Architecture



(a) BERT: Random tokens are replaced with masks, and the document is encoded bidirectionally. Missing tokens are predicted independently, so BERT cannot easily be used for generation.



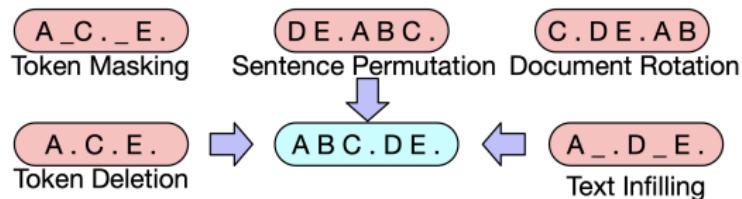
(b) GPT: Tokens are predicted auto-regressively, meaning GPT can be used for generation. However words can only condition on leftward context, so it cannot learn bidirectional interactions.



(c) BART: Inputs to the encoder need not be aligned with decoder outputs, allowing arbitrary noise transformations. Here, a document has been corrupted by replacing spans of text with mask symbols. The corrupted document (left) is encoded with a bidirectional model, and then the likelihood of the original document (right) is calculated with an autoregressive decoder. For fine-tuning, an uncorrupted document is input to both the encoder and decoder, and we use representations from the final hidden state of the decoder.

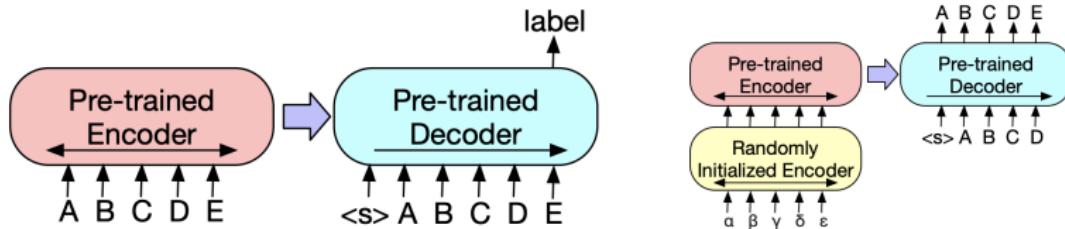
- Isn't it "Attention is All You Need?"

# Denoising S2S Pre-training



- Similar to denoising autoencoder

# Fine-tuning Tasks



(a) To use BART for classification problems, the same input is fed into the encoder and decoder, and the representation from the final output is used.

(b) For machine translation, we learn a small additional encoder that replaces the word embeddings in BART. The new encoder can use a disjoint vocabulary.

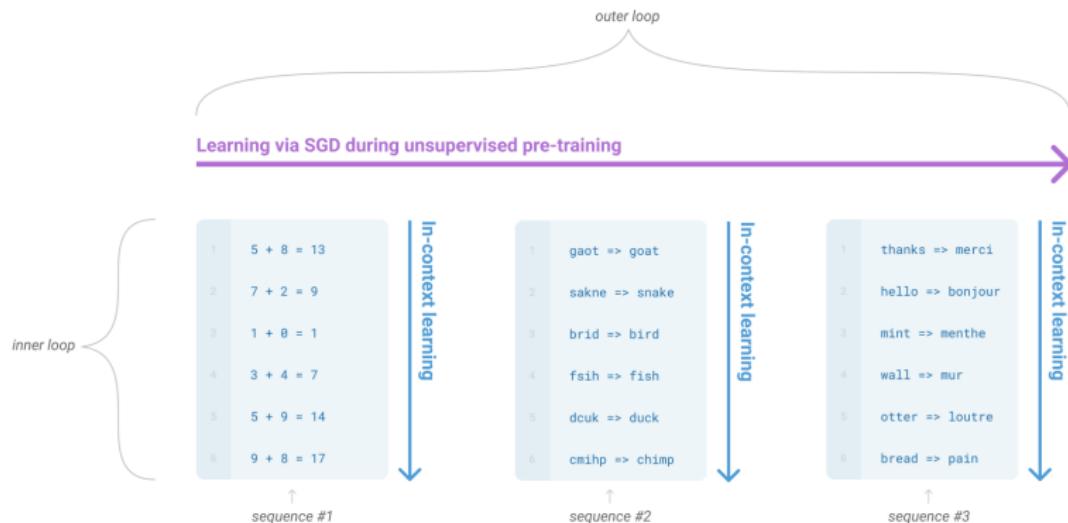
# GPTv2 [13] & GPTv3 [3]

- Still autoregressive
- But focus on *transfer* and *few-shot learning* scenarios
  - No need for fine-tuning [▶ Demo 1](#) and [▶ Demo 2](#)
- How?

# GPTv2 [13] & GPTv3 [3]

- Still autoregressive
- But focus on *transfer* and *few-shot learning* scenarios
  - No need for fine-tuning [▶ Demo 1](#) and [▶ Demo 2](#)
- How? *Multitask Language Model*
- Pre-training: to maximize  $\text{Pr}(\text{output}|\text{input}, \text{task})$  in multiple LM tasks
  - Translation:  $\text{Pr}(\text{french text}|\text{en text}, \text{translation})$
  - Question answering:  $\text{Pr}(\text{answer}|\text{question}, \text{qa})$
  - Reading comprehension:  $\text{Pr}(\text{answer}|\text{document}, \text{question}, \text{reading})$

# In-context Learning



**Figure 1.1: Language model meta-learning.** During unsupervised pre-training, a language model develops a broad set of skills and pattern recognition abilities. It then uses these abilities at inference time to rapidly adapt to or recognize the desired task. We use the term “in-context learning” to describe the inner loop of this process, which occurs within the forward-pass upon each sequence. The sequences in this diagram are not intended to be representative of the data a model would see during pre-training, but are intended to show that there are sometimes repeated sub-tasks embedded within a single sequence.

# Zero/Few Shot Learning vs. Fine-tuning

## Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



---

## One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.



---

## Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

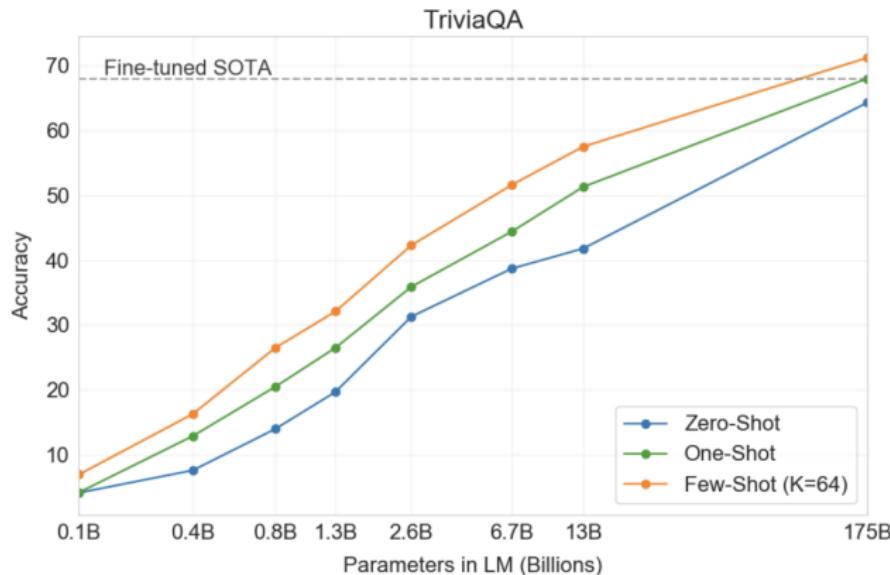


## Fine-tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.



# Performance vs. Model Size



**Figure 3.3:** On TriviaQA GPT3's performance grows smoothly with model size, suggesting that language models continue to absorb knowledge as their capacity increases. One-shot and few-shot performance make significant gains over zero-shot behavior, matching and exceeding the performance of the SOTA fine-tuned open-domain model, RAG [LPP<sup>+</sup>20]

# Outline

## ① RNNs

- Vanilla RNNs
- Design Alternatives

## ② RNN Training

- Backprop through Time (BPTT)
- Optimization Techniques
- Optimization-Friendly Models & LSTM
- Parallelism & Teacher Forcing

## ③ RNNs with Attention Mechanism

- Attention for Image Captioning
- Attention for Neural Machine Translation (NMT)

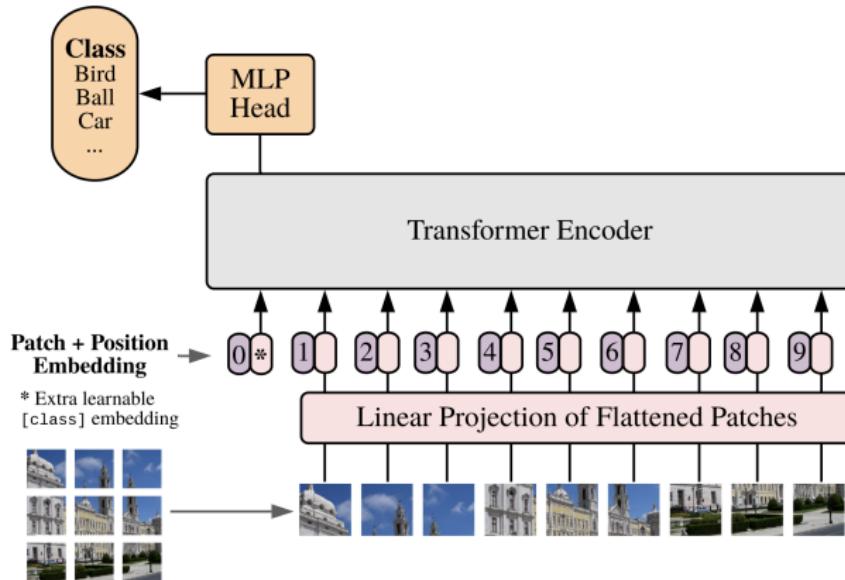
## ④ Transformers

- Attention is All You Need
- Pretrained Language Models
- More Applications

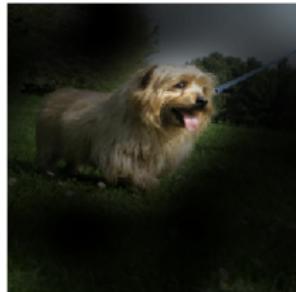
## ⑤ Subword Tokenization

# Vision Transformers (ViT) [5]

- Splits an image into  $16 \times 16$  patches (words)
- Like BERT, uses [CLS] input word to get class predictions



# Attention



# DALL-E [14]

- “An astronaut riding a horse in a photorealistic style”

[Live Demo](#)



# Outline

## ① RNNs

- Vanilla RNNs
- Design Alternatives

## ② RNN Training

- Backprop through Time (BPTT)
- Optimization Techniques
- Optimization-Friendly Models & LSTM
- Parallelism & Teacher Forcing

## ③ RNNs with Attention Mechanism

- Attention for Image Captioning
- Attention for Neural Machine Translation (NMT)

## ④ Transformers

- Attention is All You Need
- Pretrained Language Models
- More Applications

## ⑤ Subword Tokenization

# Sequence Tokenization

- Word-level
  - High input/output dimension ( $D$ )
  - Out-of-vocabulary (OOV) problem
  - “old”  $\neq$  “older”  $\neq$  “oldest”
- Char-level
  - Long sequence ( $T$ )

# Sequence Tokenization

- Word-level
  - High input/output dimension ( $D$ )
  - Out-of-vocabulary (OOV) problem
  - “old”  $\neq$  “older”  $\neq$  “oldest”
- Char-level
  - Long sequence ( $T$ )
- Subword-level?
  - How to deal with OOV problem?
  - How to support different languages?

# Byte Pair Encoding (BPE) [16]

- Given a training set of words

- Uses Unicode bytes as base symbols  $\{s^{(1)}, s^{(2)}, \dots\}$
- Merge two symbol  $s^{(i)}$  and  $s^{(j)}$  having highest  $\text{Pr}(s^{(i)})$  and  $\text{Pr}(s^{(j)})$
- Repeat step 2 until target #symbols is met

	Vocabulary	Encoded Sentence
Initialization	[‘a’, ‘c’, ‘b’, ‘e’, ‘i’, ‘</w>’, ‘k’, ‘m’, ‘o’, ‘n’, ‘p’, ‘s’, ‘r’, ‘u’, ‘t’, ‘v’, ‘x’]	vi e t n a m </w> t a k e s </w> me a s u r e s </w> t o </w> b o o s t </w> r i c e </w> e x p o r t s </w>
After merge operation 1	[‘a’, ‘c’, ‘b’, ‘e’, ‘p’, ‘</w>’, ‘k’, ‘m’, ‘o’, ‘n’, ‘i’, ‘s’, ‘r’, ‘u’, ‘t’, ‘v’, ‘x’, ‘s</w>’]	vi e t n a m </w> t a k e s </w> me a s u r e s </w> t o </w> b o o s t </w> r i c e </w> e x p o r t s </w>
After merge operations 10	[‘</w>’, ‘vi’, ‘as’, ‘es</w>’, ‘s</w>’, ‘nam’, ‘to’, ‘ri’, ‘t</w>’, ‘ort’, ‘a’, ‘c’, ‘b’, ‘e’, ‘k’, ‘m’, ‘o’, ‘p’, ‘s’, ‘r’, ‘u’, ‘t’, ‘x’]	vi e t n a m </w> t a k e s </w> m e a s u r e s </w> t o </w> b o o s t </w> r i c e </w> e x p o r t s </w>
After merge operations 34	[‘takes</w>’, ‘measures</w>’, ‘exports</w>’, ‘boost</w>’, ‘rice</w>’, ‘vietnam</w>’, ‘to</w>’]	vietnam </w> takes </w> mea sures </w> to </w> boost </w> rice </w> exports </w>

# Other Variants

- BPE is used by GPT
- WordPiece [15]:
  - Merge  $s^{(i)}$  and  $s^{(j)}$  having highest  $\frac{\Pr(s^{(i)}, s^{(j)})}{\Pr(s^{(i)}) \Pr(s^{(j)})}$
  - Used by BERT
- Unigram Language Model [8]
  - Top-down, probabilistic
- SentencePiece [9]
  - Treat space “\_” as symbols to support different languages
  - Merge algorithm: BPE or Unigram Language Model
  - Used by ALBERT, XLNet, Marian, T5

# Reference I

- [1] Martin Arjovsky, Amar Shah, and Yoshua Bengio.  
Unitary evolution recurrent neural networks.  
*arXiv preprint arXiv:1511.06464*, 2015.
- [2] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio.  
Neural machine translation by jointly learning to align and translate.  
*arXiv preprint arXiv:1409.0473*, 2014.
- [3] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al.  
Language models are few-shot learners.  
*Advances in neural information processing systems*, 33:1877–1901, 2020.

## Reference II

- [4] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding.  
*arXiv preprint arXiv:1810.04805*, 2018.
- [5] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al.  
An image is worth 16x16 words: Transformers for image recognition at scale.  
*arXiv preprint arXiv:2010.11929*, 2020.
- [6] Andrej Karpathy, Justin Johnson, and Li Fei-Fei.  
Visualizing and understanding recurrent networks.  
*arXiv preprint arXiv:1506.02078*, 2015.

# Reference III

- [7] David Krueger and Roland Memisevic.  
Regularizing rnns by stabilizing activations.  
*arXiv preprint arXiv:1511.08400*, 2015.
- [8] Taku Kudo.  
Subword regularization: Improving neural network translation models with multiple subword candidates.  
*arXiv preprint arXiv:1804.10959*, 2018.
- [9] Taku Kudo and John Richardson.  
Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing.  
*arXiv preprint arXiv:1808.06226*, 2018.
- [10] Quoc V Le, Navdeep Jaitly, and Geoffrey E Hinton.  
A simple way to initialize recurrent networks of rectified linear units.  
*arXiv preprint arXiv:1504.00941*, 2015.

## Reference IV

- [11] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer.  
Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension.  
*arXiv preprint arXiv:1910.13461*, 2019.
- [12] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al.  
Improving language understanding by generative pre-training.  
*OpenAI blog*, 2018.
- [13] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al.  
Language models are unsupervised multitask learners.  
*OpenAI blog*, 2019.

# Reference V

- [14] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen.  
Hierarchical text-conditional image generation with clip latents.  
*arXiv preprint arXiv:2204.06125*, 2022.
- [15] Mike Schuster and Kaisuke Nakajima.  
Japanese and korean voice search.  
In *2012 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 5149–5152. IEEE, 2012.
- [16] Rico Sennrich, Barry Haddow, and Alexandra Birch.  
Neural machine translation of rare words with subword units.  
*arXiv preprint arXiv:1508.07909*, 2015.
- [17] Or Sharir, Barak Peleg, and Yoav Shoham.  
The cost of training nlp models: A concise overview.  
*arXiv preprint arXiv:2004.08900*, 2020.

## Reference VI

- [18] Eduardo D Sontag.  
On the computational power of neural nets'.  
*JOURNAL OF COMPUTER AND SYSTEM SCIENCES*, 50:132–150,  
1995.
- [19] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion  
Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin.  
Attention is all you need.  
*Advances in neural information processing systems*, 30, 2017.
- [20] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville,  
Ruslan Salakhutdinov, Richard S Zemel, and Yoshua Bengio.  
Show, attend and tell: Neural image caption generation with visual  
attention.  
*arXiv preprint arXiv:1502.03044*, 2(3):5, 2015.