

# 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

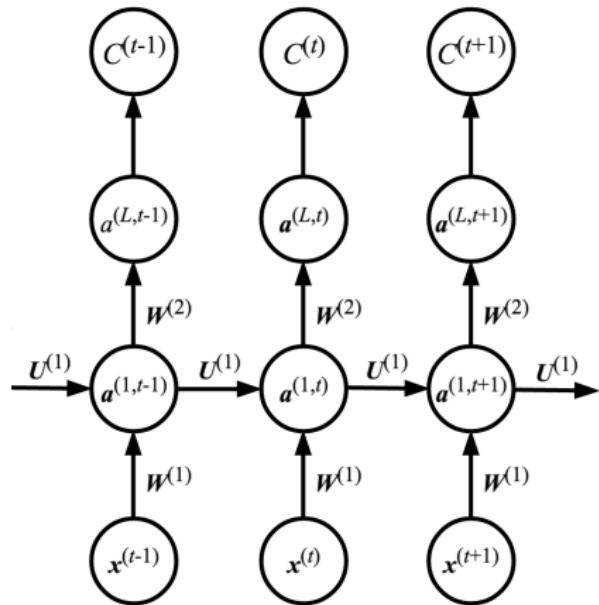
- How to model sequential data?

# Sequence Modeling I

- How to model sequential data?

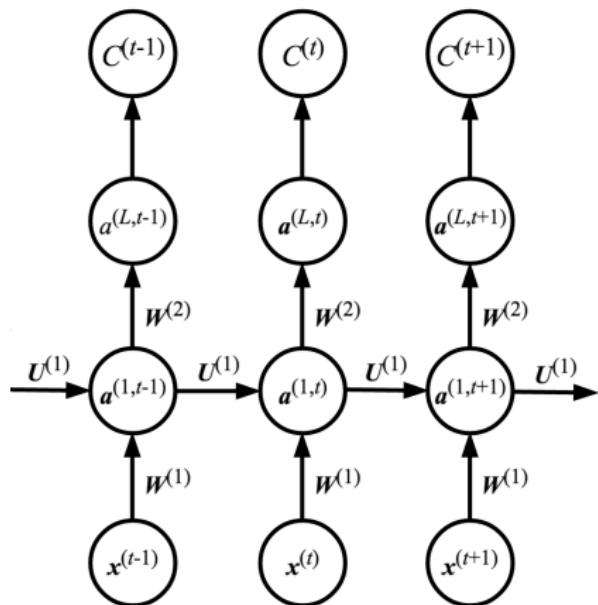
- **Recurrent neural networks**

(vanilla RNNs):



# Sequence Modeling I

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

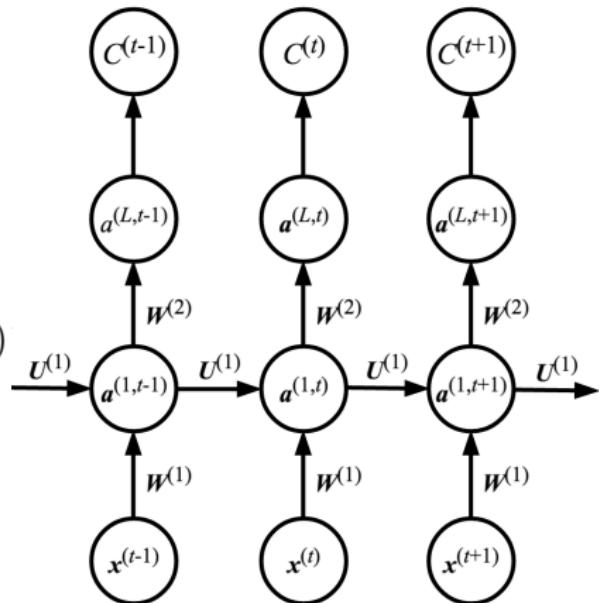


# Sequence Modeling I

- 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

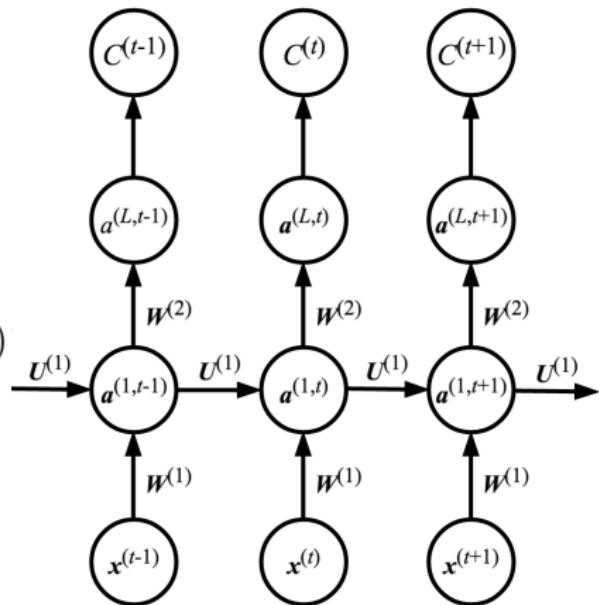


# Sequence Modeling I

- How to model sequential data?
- **Recurrent neural networks**  
(vanilla RNNs):
- $y^{(t)}$  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
- $\mathbf{a}^{(\cdot,t)}$  summarizes  $x^{(t)}, \dots, x^{(1)}$ 
  - Earlier points are less important



# Sequence Modeling I

- How to model sequential data?

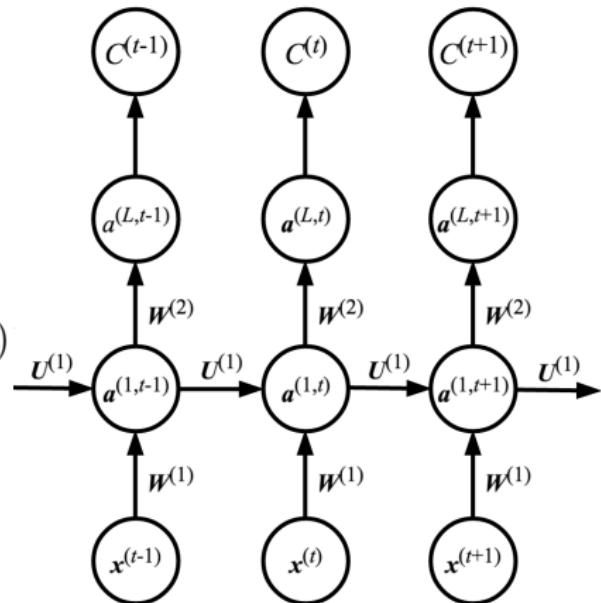
- **Recurrent neural networks**

(vanilla RNNs):

- $y^{(t)}$  depends on  $x^{(1)}, \dots, x^{(t)}$
- Output  $a^{(L,t)}$  depends on hidden activations:

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

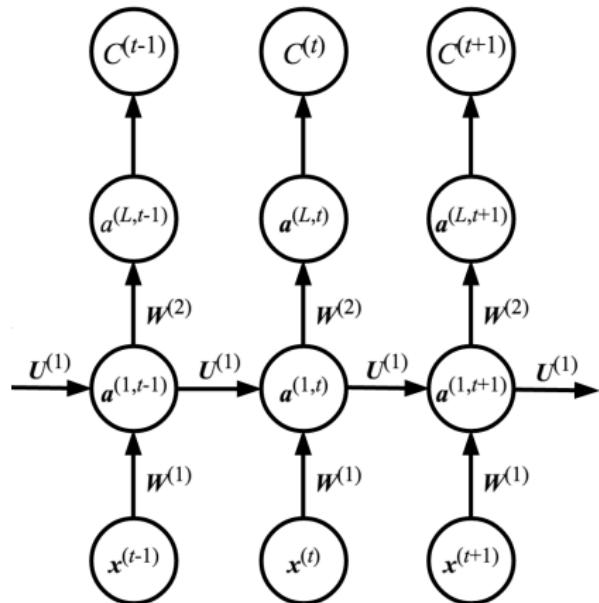
- Bias term omitted
- $\mathbf{a}^{(\cdot,t)}$  summarizes  $x^{(t)}, \dots, x^{(1)}$ 
  - Earlier points are less important
- $\mathbf{a}^{(\cdot,t)}$ 's at deeper layers give more abstract summarizations



# Sequence Modeling II

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

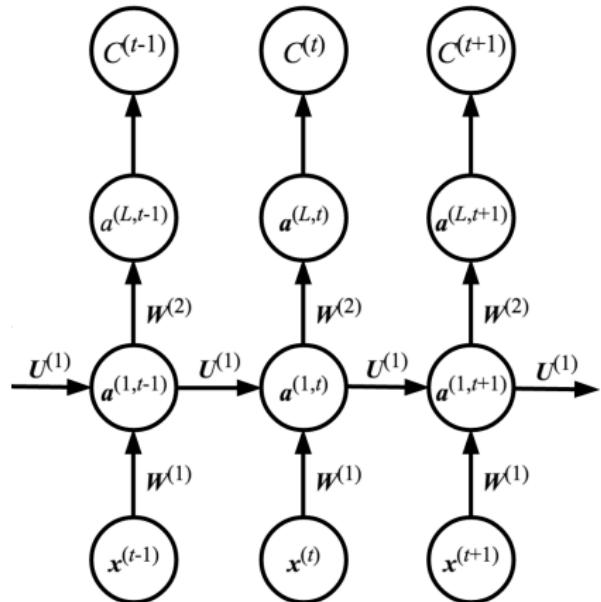
- Weights are *shared* across time instances



# Sequence Modeling II

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

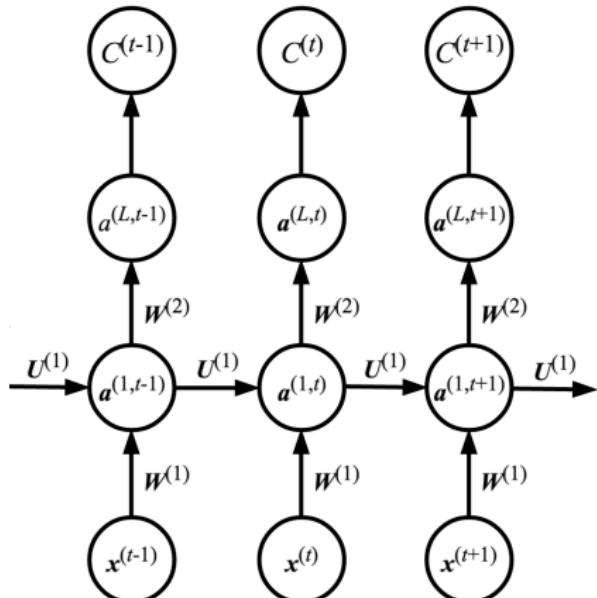
- Weights are *shared* across time instances
- Assumes that the “transition functions” are time invariant



# Sequence Modeling II

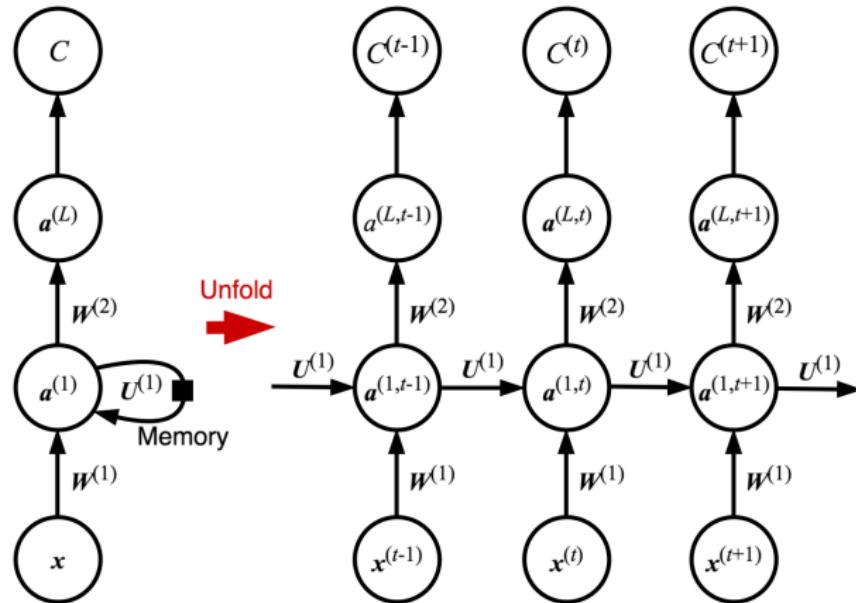
$$\begin{aligned}\mathbf{a}^{(k,t)} &= \text{act}(\mathbf{z}^{(k,t)}) \\ &= \text{act}(\mathbf{U}^{(k)} \mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)} \mathbf{a}^{(k-1,t)})\end{aligned}$$

- 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  $k = 1, \dots, L$



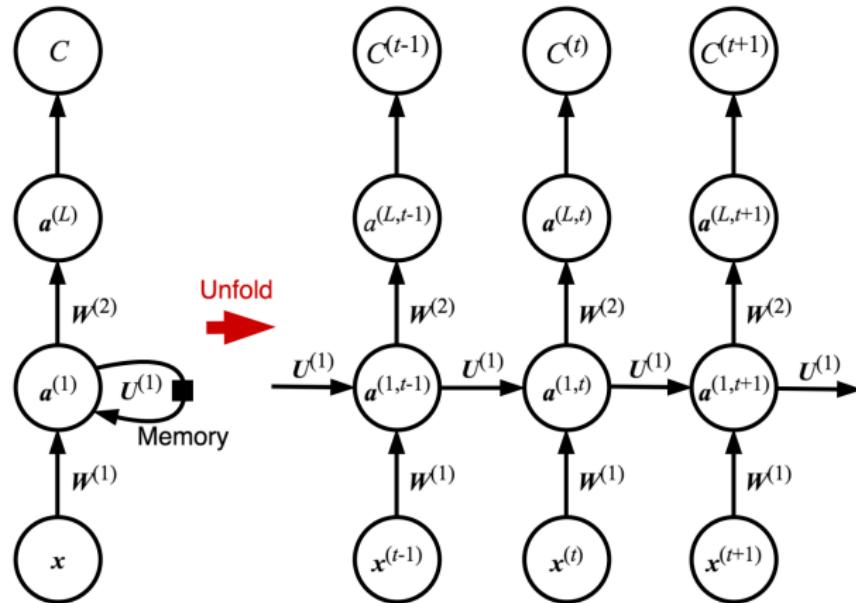
# RNNs have Memory

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



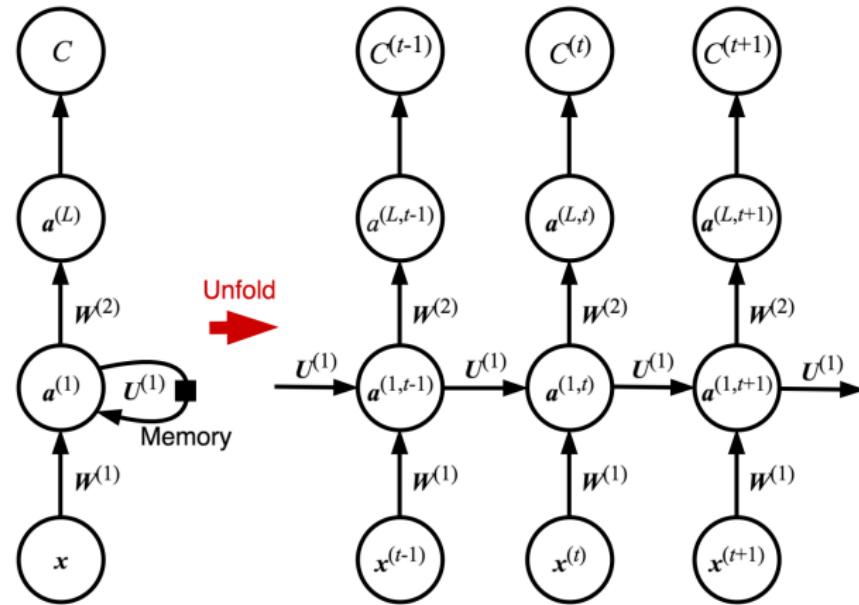
# RNNs have Memory

- 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  $\mathbf{a}^{(L,t)}$

你好 HALLO 안녕  
CIAO HOLA নমস্তে  
¡Hola! HELLO  
こんにちは プリベ特  
BONJOUR مرحباً OLÁ

- 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  
CIAO HOLA নমস্তে  
¡Hola! HELLO  
こんにちは プリベ特  
BONJOUR مرحباً 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	$O(\log_F T)$
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)

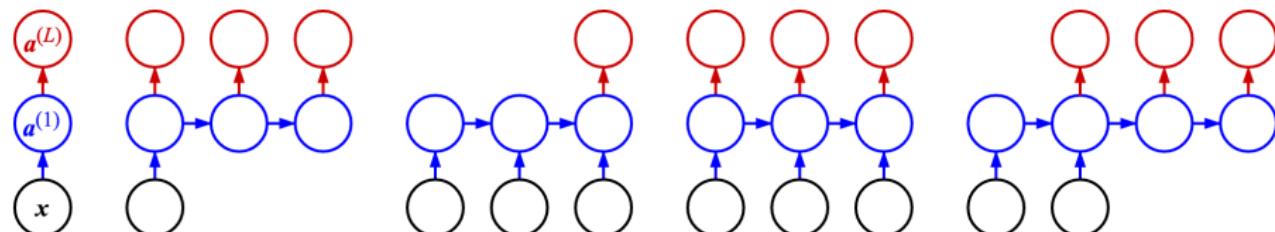
## ④ Transformers

- 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

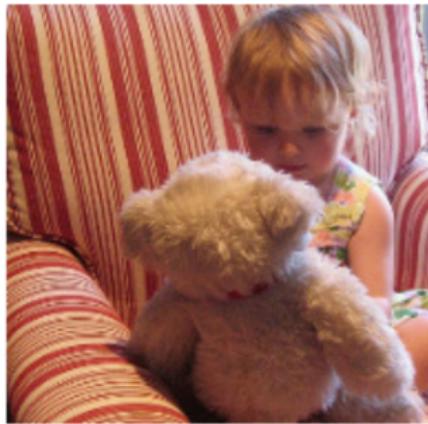
One to Many

Many to One

Many to Many  
(Synced)

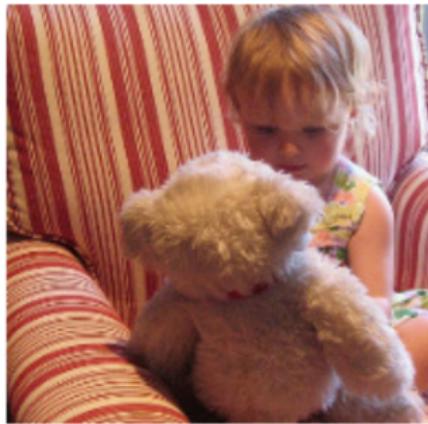
Many to Many  
(Unsynced)

# One2Many: Image Captioning

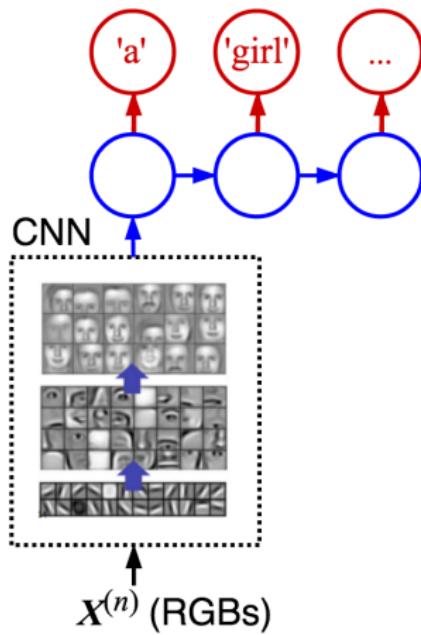


“A little girl sitting on a bed with a  
teddy bear.”

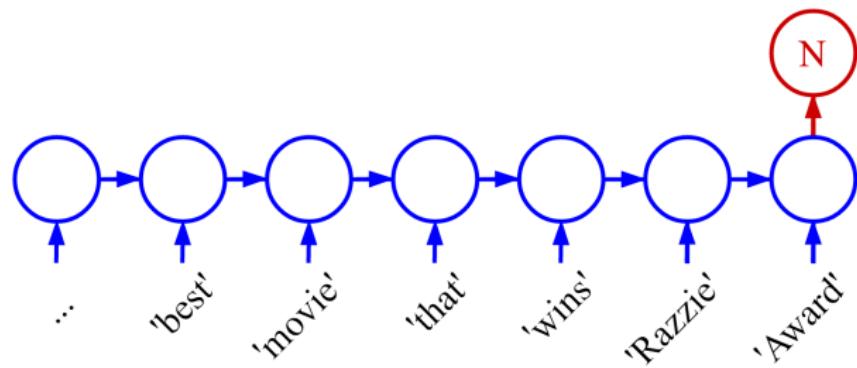
# One2Many: Image Captioning



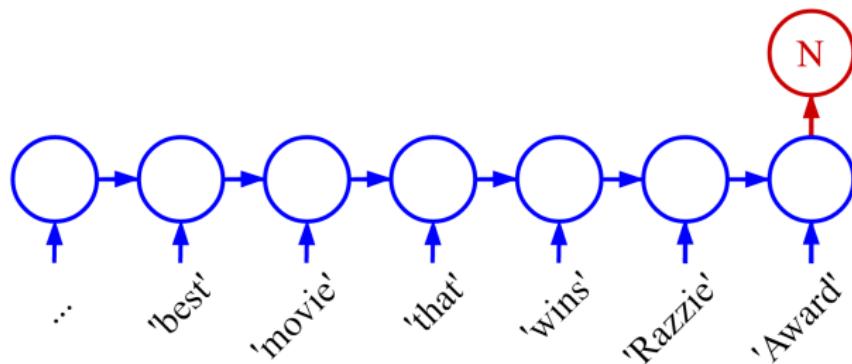
"A little girl sitting on a bed with a teddy bear."



# Many2One: Sentiment Analysis



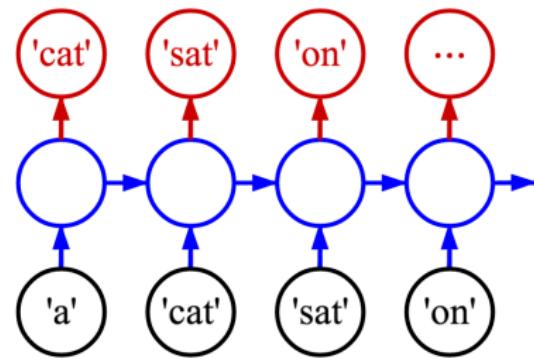
# Many2One: Sentiment Analysis



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

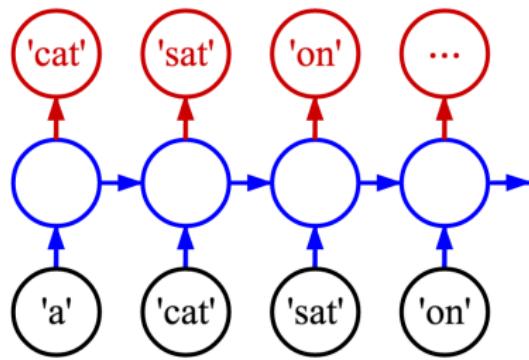
# Many2Many (Synced): Language Modeling

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



# Many2Many (Synced): Language Modeling

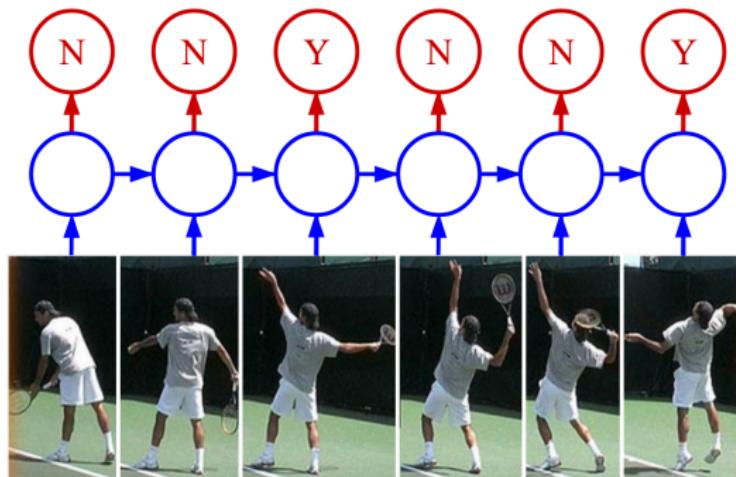
- **Language modeling**: predicting the next/nearby word based on the context



- Latent representations of RNN provide the context

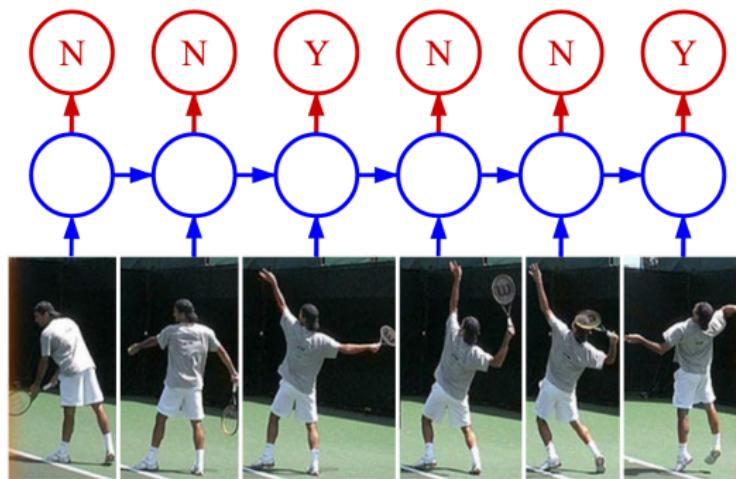
# Many2Many (Synced): Video Keyframe Tagging

- Video frame annotation:



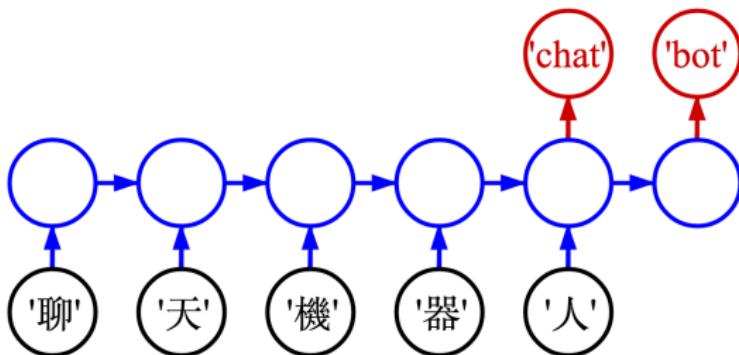
# Many2Many (Synced): Video Keyframe Tagging

- Video frame annotation:



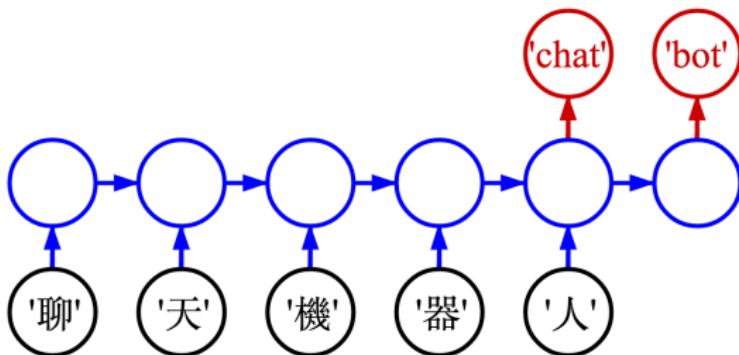
- Latent representations summarize “what’s going on”

# 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

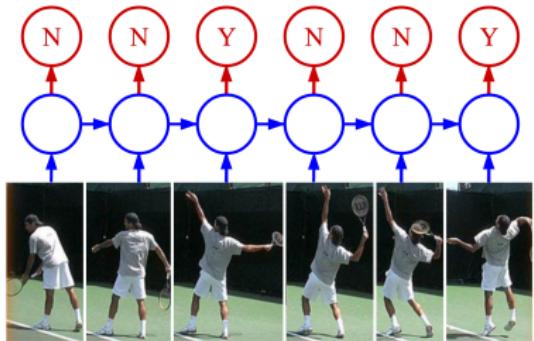
$$\mathbf{a}^{(k,t)} = \text{act}(\mathbf{U}^{(k)} \mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)} \mathbf{a}^{(k-1,t)})$$

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

# Bidirectional RNNs

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

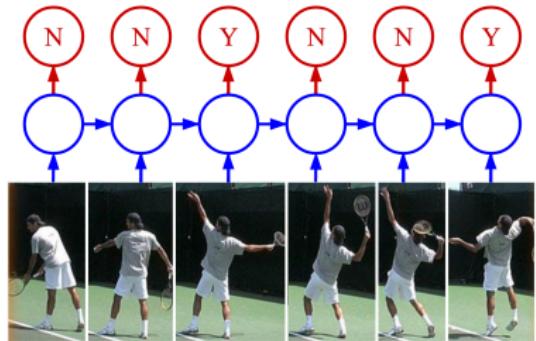
- Each  $\mathbf{a}^{(\cdot,t)}$  summarizes  $\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}$
- A prediction is made by considering previous frames



# 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?

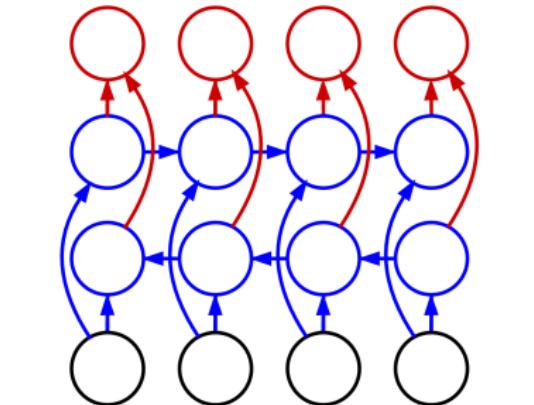
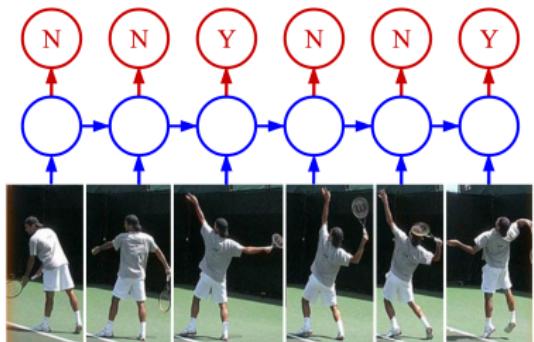


# Bidirectional RNNs

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

- 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

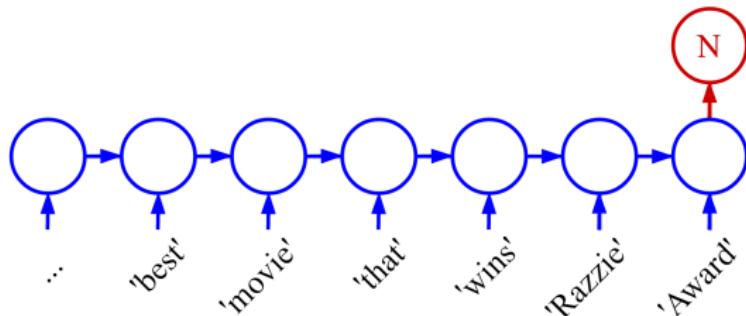
$$\mathbf{a}^{(k,t)} = \text{act}(\mathbf{U}^{(k)}\mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)}\mathbf{a}^{(k-1,t)})$$

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

# Recursive RNNs I

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

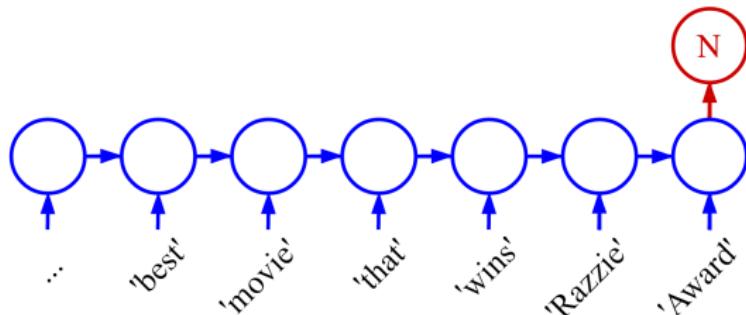
- The transition of hidden representation is invariant in time
  - Earlier input/representation is less important to current  $\mathbf{a}^{(k,t)}$
- “Razzie” has less effect if it is far away from the prediction



# Recursive RNNs I

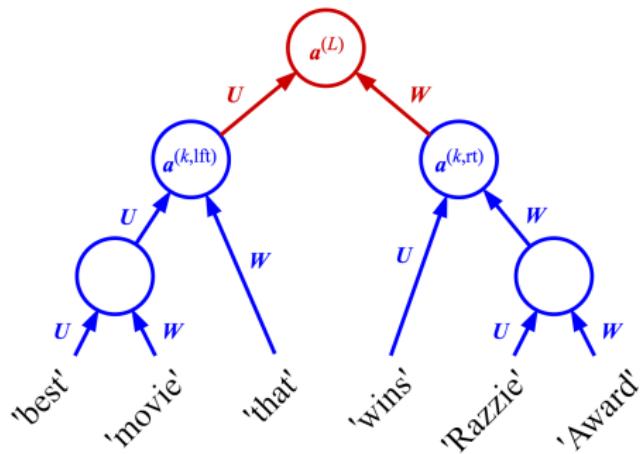
$$\mathbf{a}^{(k,t)} = \text{act}(\mathbf{U}^{(k)}\mathbf{a}^{(k,t-1)} + \mathbf{W}^{(k)}\mathbf{a}^{(k-1,t)})$$

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

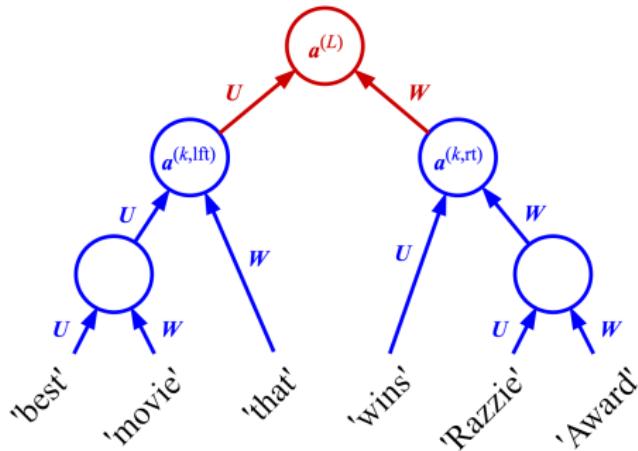


# Recursive RNNs II

- 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

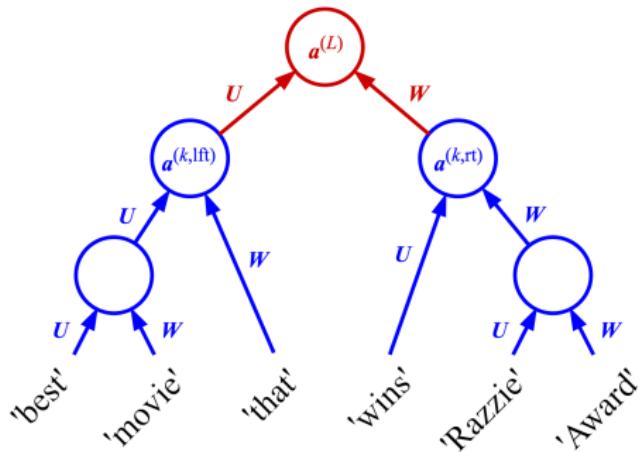


# Recursive RNNs II

- 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

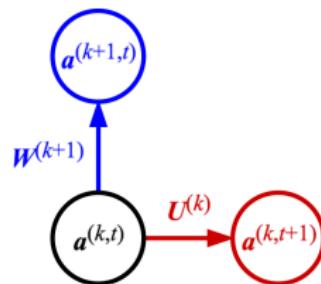
$$\delta_j^{(k,t)} = \frac{\partial c^{(n,t)}}{\partial z_j^{(k,t)}} = \frac{\partial c^{(n,t)}}{\partial a_j^{(k,t)}} \cdot \frac{\partial a_j^{(k,t)}}{\partial z_j^{(k,t)}} = \frac{\partial c^{(n,t)}}{\partial a_j^{(k,t)}} \cdot \text{act}'(z_j^{(k,t)})$$

# 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

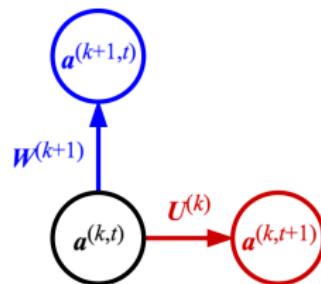
$$\begin{aligned}\delta_j^{(k,t)} &= \frac{\partial c^{(n,t)}}{\partial z_j^{(k,t)}} = \frac{\partial c^{(n,t)}}{\partial a_j^{(k,t)}} \cdot \frac{\partial a_j^{(k,t)}}{\partial z_j^{(k,t)}} = \frac{\partial c^{(n,t)}}{\partial a_j^{(k,t)}} \cdot \text{act}'(z_j^{(k,t)}) \\ &= \left( \sum_s \frac{\partial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



# 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

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



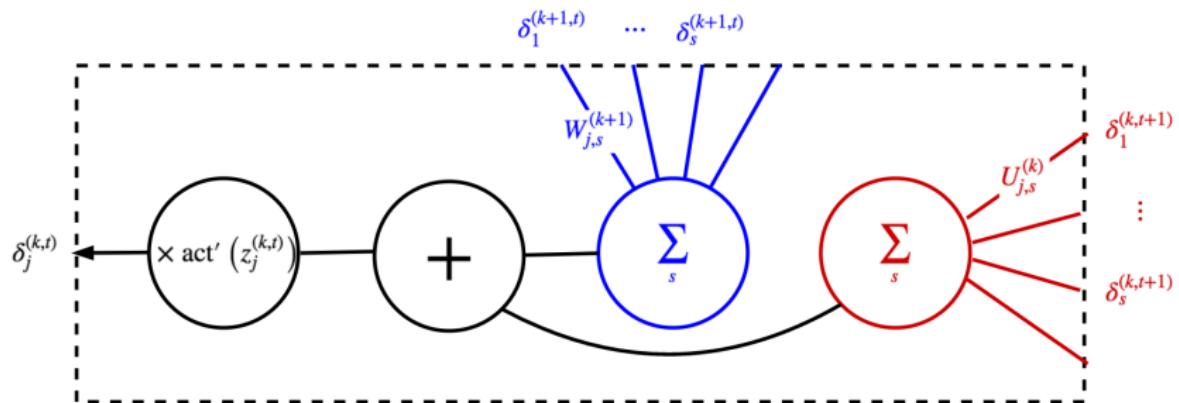
# Backward Pass through Time II

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

# 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)

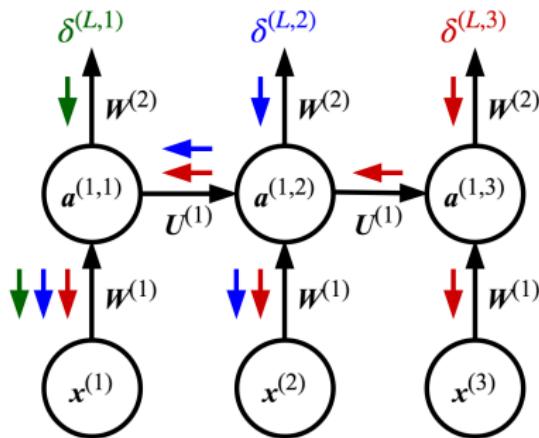
- So far, we have discussed how to compute the gradients of  $U_{i,j}^{(k)}$ 's (and  $W_{i,j}^{(k)}$ 's) for a single loss  $c^{(n,t)} = C^{(n,t)}(\Theta^{(s)})$  at specific time  $t$
- However, for each sequence  $n$ , we have multiple  $c^{(n,t)}$ 's at different  $t$ 's
  - Their gradient need to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

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

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

- Backward pass:

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

# 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

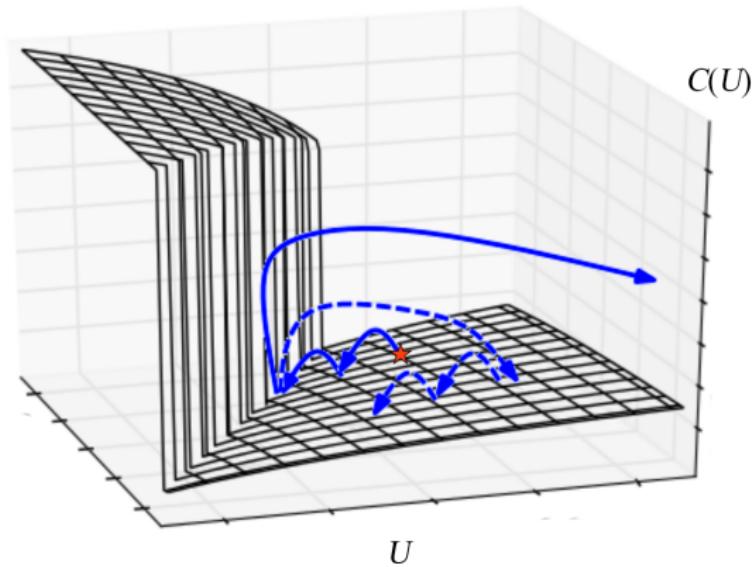
- Ignoring activation function and depth:

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

- 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!

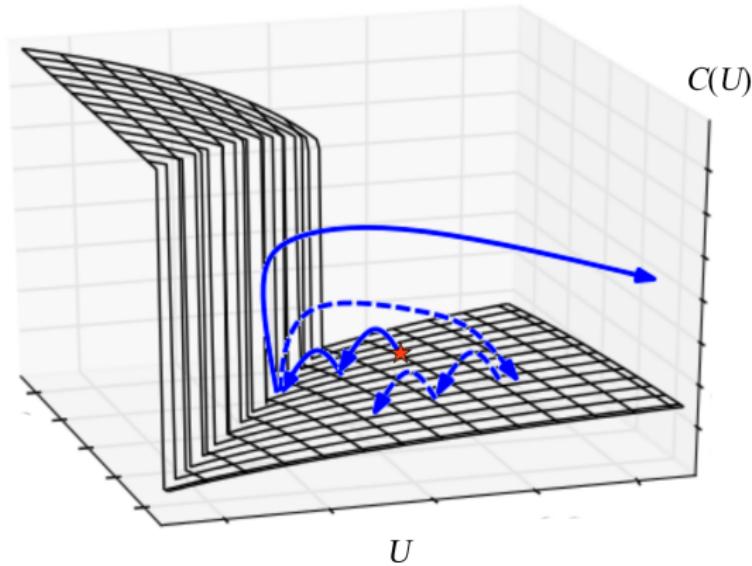
# Cost Surface

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



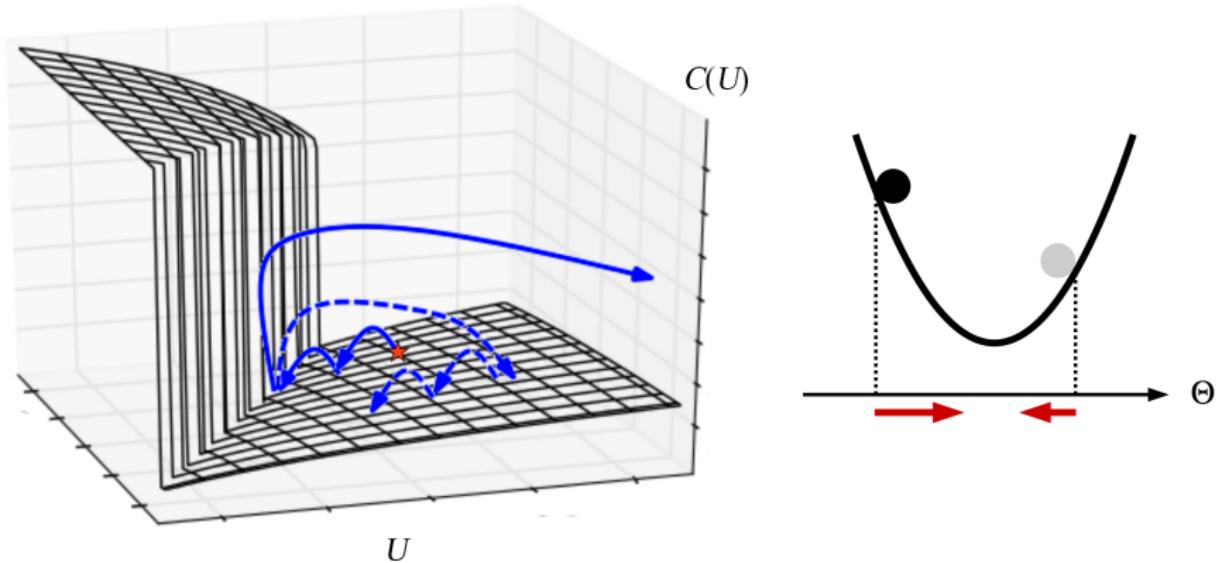
# Cost Surface

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



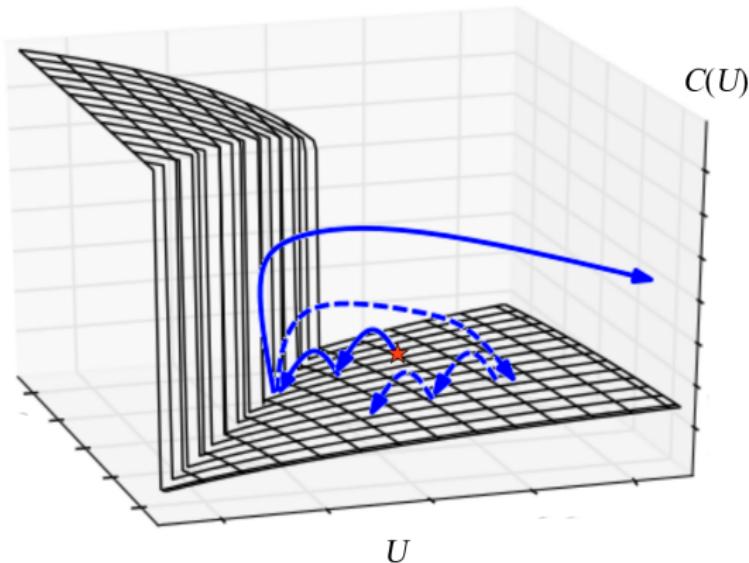
# Nesterov Momentum

- 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

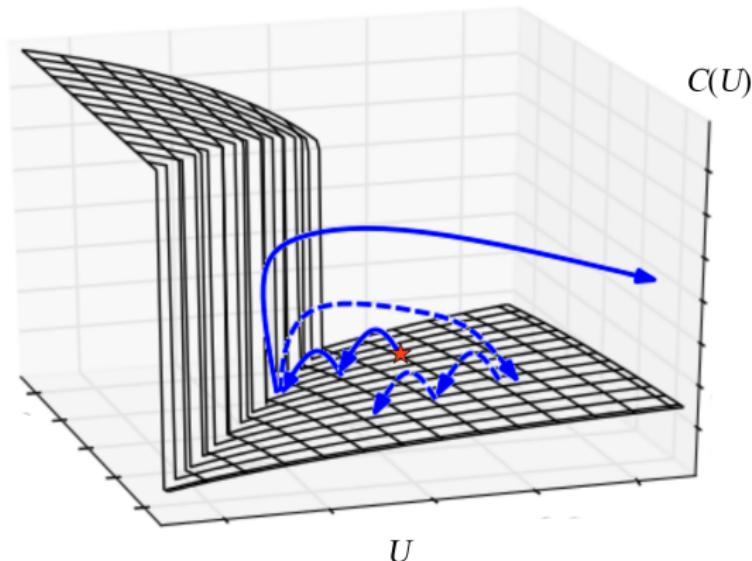


# RMS Prop

- Adaptive learning rate based on statistics of recent gradients:

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

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



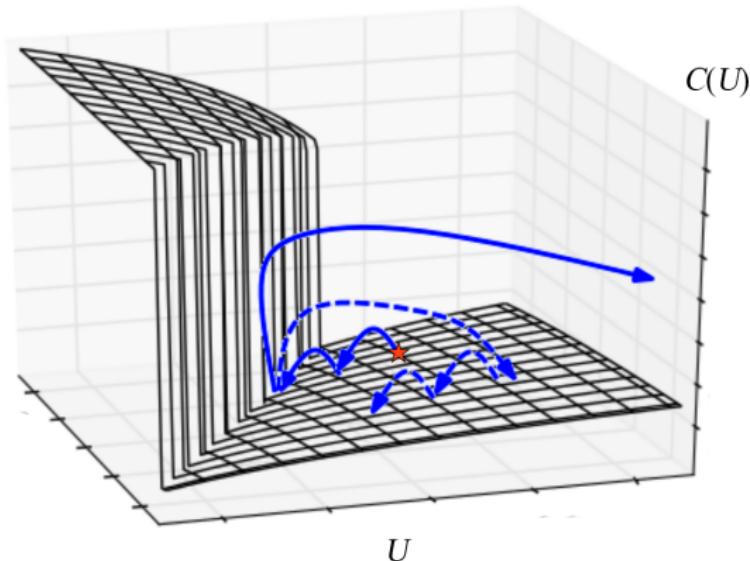
# RMS Prop

- 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

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

- 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. [5] propose IRNN:
  - **Initializes  $U^{(k)} = I$**  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. [5] 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

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

- 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. [5] propose IRNN:
  - **Initializes  $\mathbf{U}^{(k)} = \mathbf{I}$**  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. [4] add a term  $\sum_{k,t} (\|\mathbf{a}^{(k,t)}\| - \|\mathbf{a}^{(k,t-1)}\|)^2$  to IRNN cost to stabilize the norms of  $\mathbf{a}^{(k,t)}$ 's in time

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

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

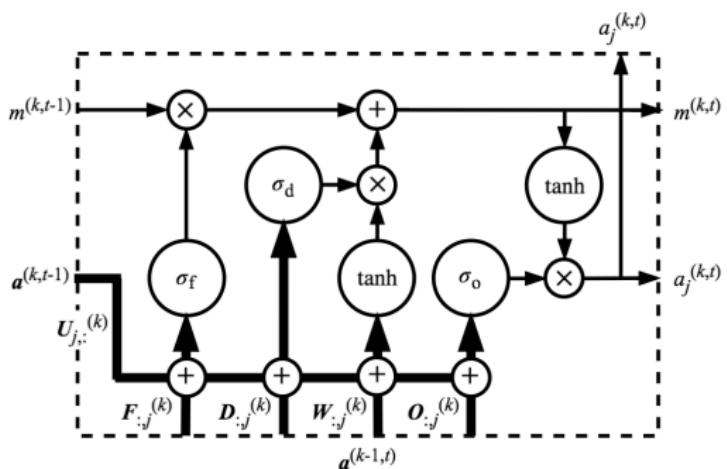
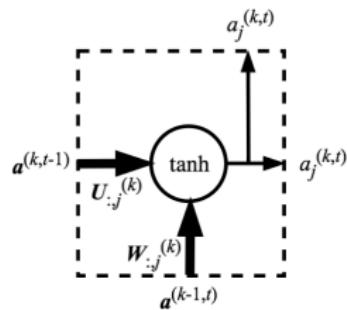
- Long-term dependency:  $(\mathbf{U}^{(k)})^{j-i} = \mathbf{Q} \begin{bmatrix} \lambda_1^{j-i} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_{D^{(k)}}^{j-i} \end{bmatrix} \mathbf{Q}^\top$
- Why not make  $\mathbf{U}^{(k)}$ 's unitary (i.e.,  $\lambda_s = 1$  for all  $s$ )?
- Hinton et al. [5] propose IRNN:
  - **Initializes**  $\mathbf{U}^{(k)} = \mathbf{I}$  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. [4] add a term  $\sum_{k,t} (\|\mathbf{a}^{(k,t)}\| - \|\mathbf{a}^{(k,t-1)}\|)^2$  to IRNN cost to stabilize the norms of  $\mathbf{a}^{(k,t)}$ 's in time
- Bengio et al. [1] propose uRNN that learns unitary  $\mathbf{U}^{(k)}$ 's explicitly

# Long Short-Term Memory (LSTM)

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

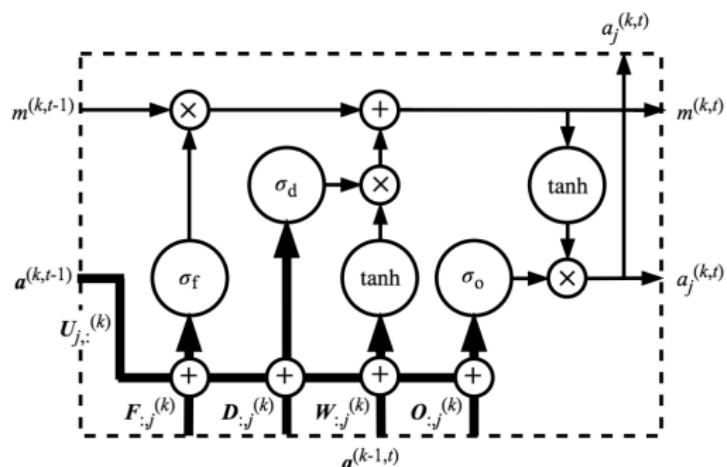
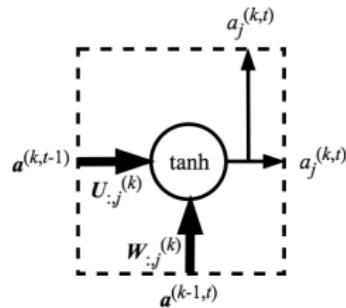
# Long Short-Term Memory (LSTM)

- 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  $t$ :



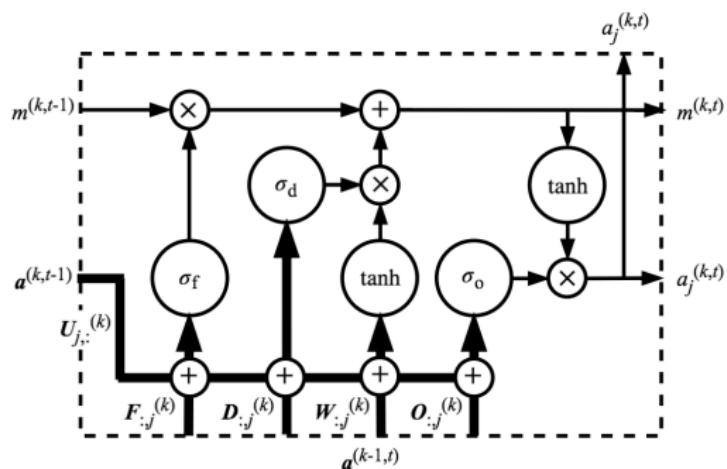
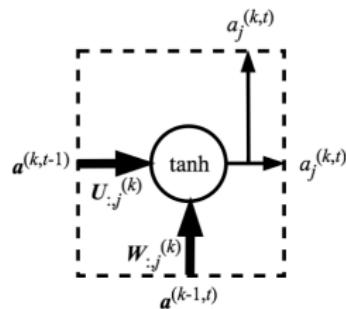
# Long Short-Term Memory (LSTM)

- 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  $a_j^{(k,t)}$



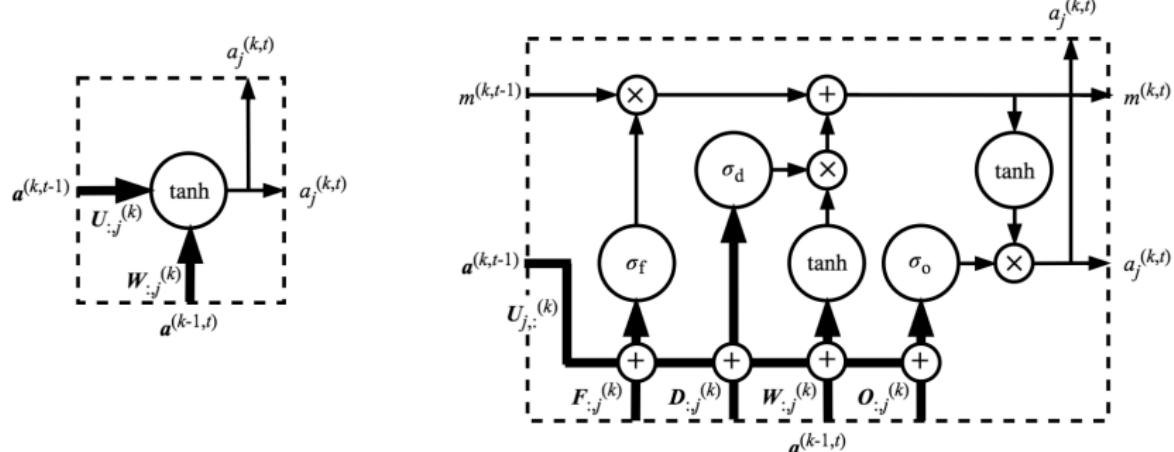
# Long Short-Term Memory (LSTM)

- 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  $m^{(k,t-1)}$



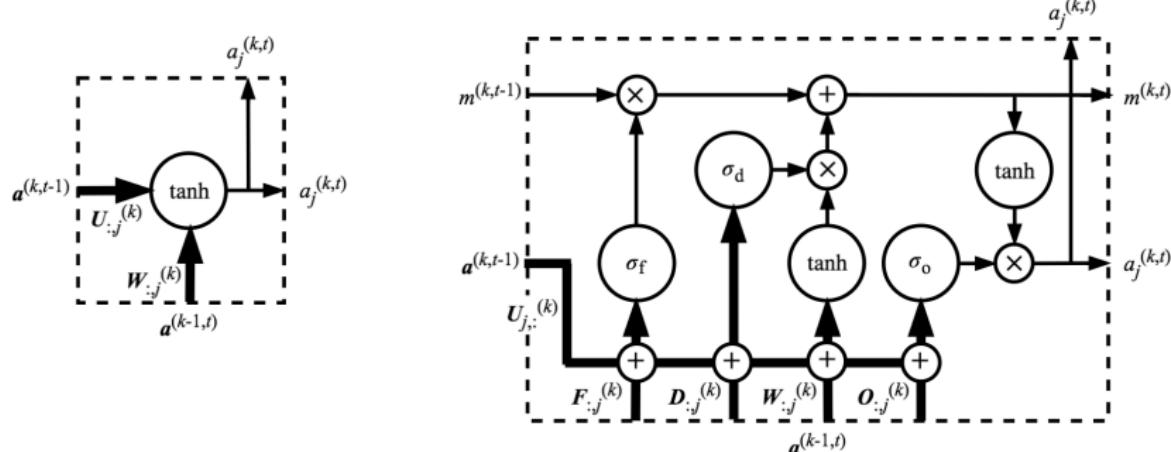
# Long Short-Term Memory (LSTM)

- 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  $m^{(k,t)}$



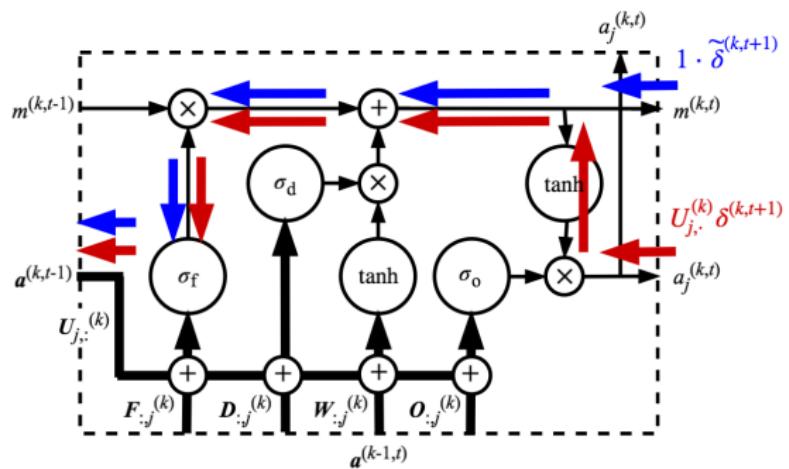
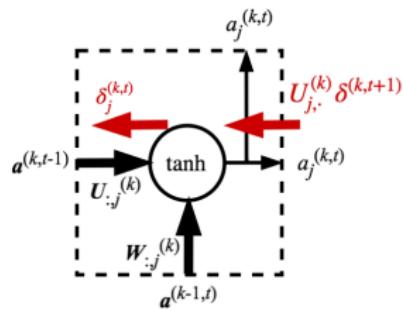
# Long Short-Term Memory (LSTM)

- 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



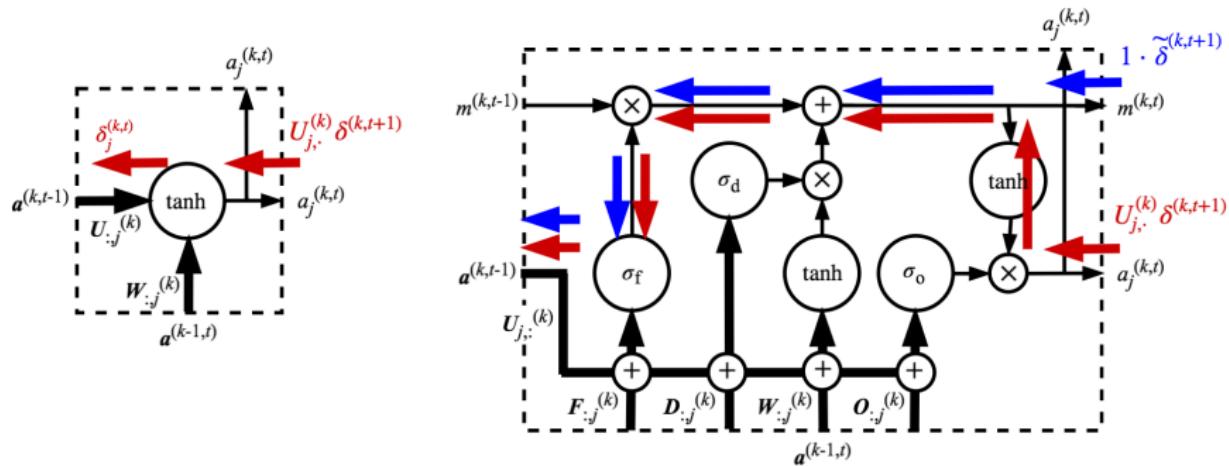
# Error Signals

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



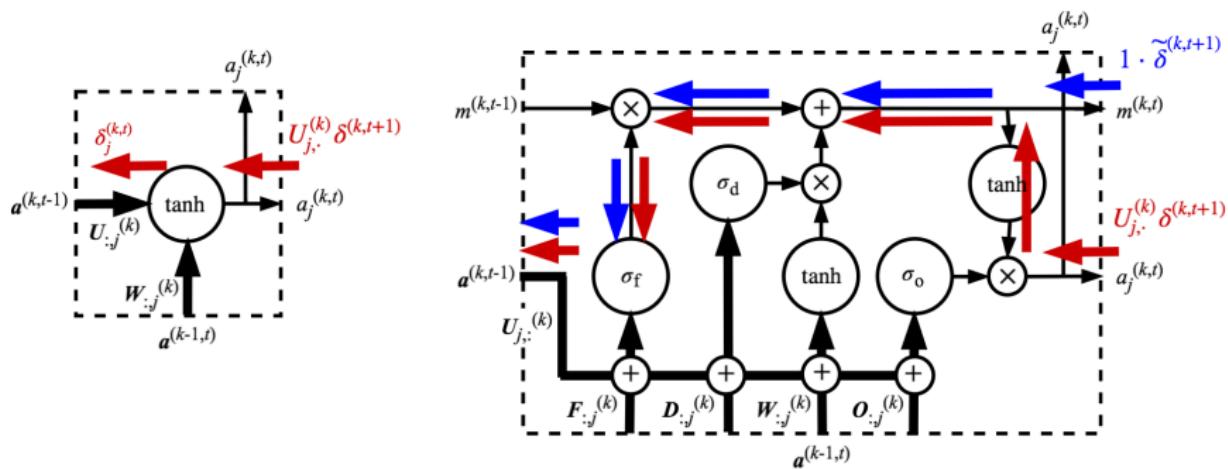
# Error Signals

- 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)

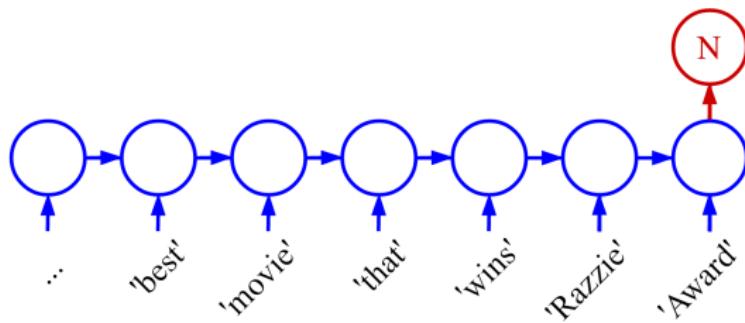


# Error Signals

- 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

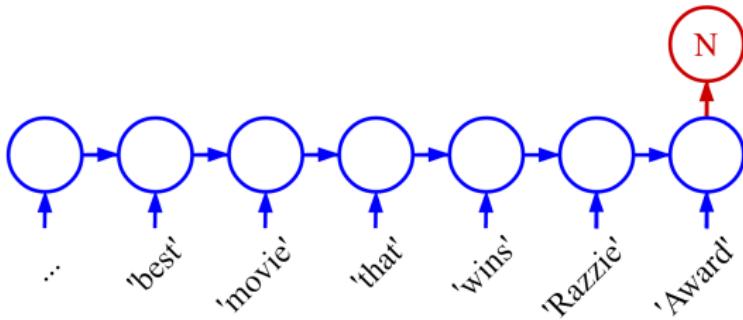


# Dynamic Representations for Sentiment Analysis



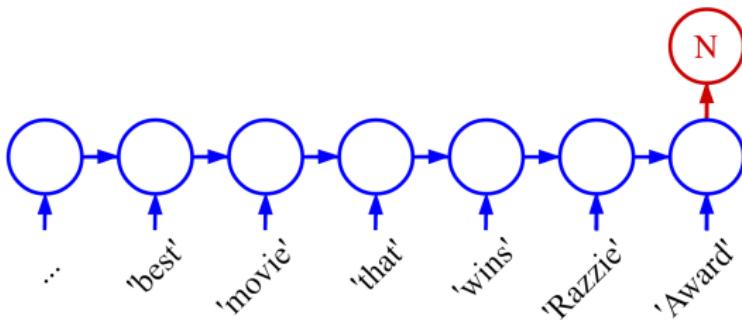
- LSTM units learn dynamic representations

# Dynamic Representations for Sentiment Analysis



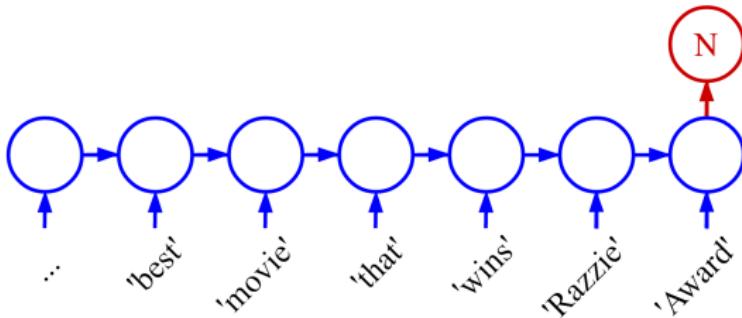
- LSTM units learn dynamic representations
- Closing forget gate for “Razzie”
  - To correct previous summarization and/or shift focus

# 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 [3] [▶ 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...*”

# 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)

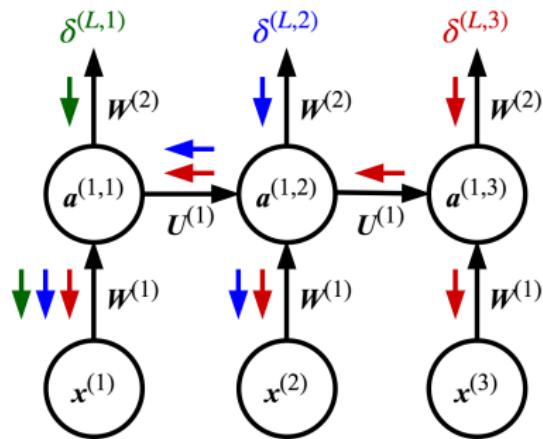
## ④ Transformers

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

## ⑤ Subword Tokenization

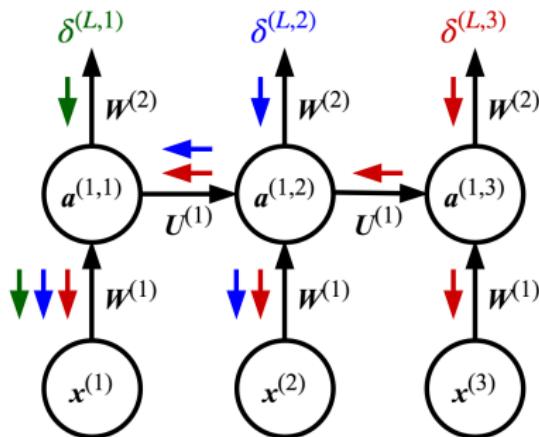
# Parallelism

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



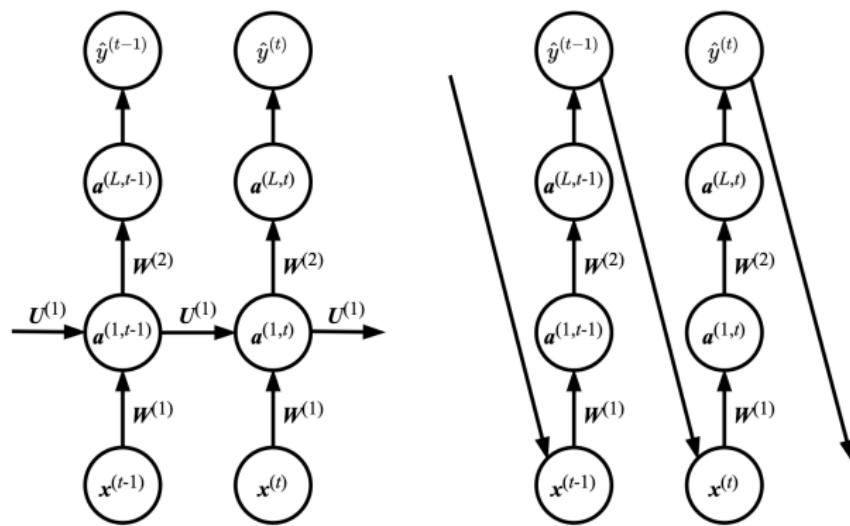
# Parallelism

- 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



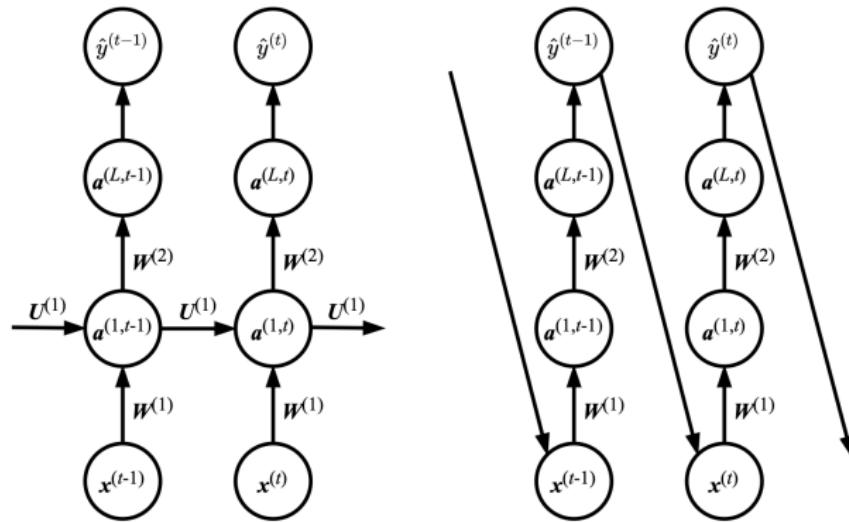
# Output Recurrence and Teacher Forcing

- **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



# Exposure Bias

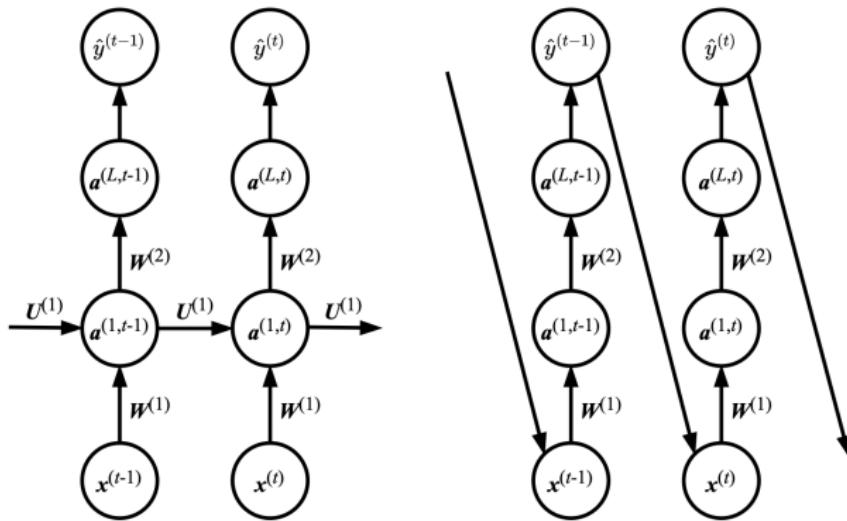
- Mismatch between  $y^{(\cdot)}$ 's and  $\hat{y}^{(\cdot)}$ 's hurts RNN performance
- Solution?

# 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

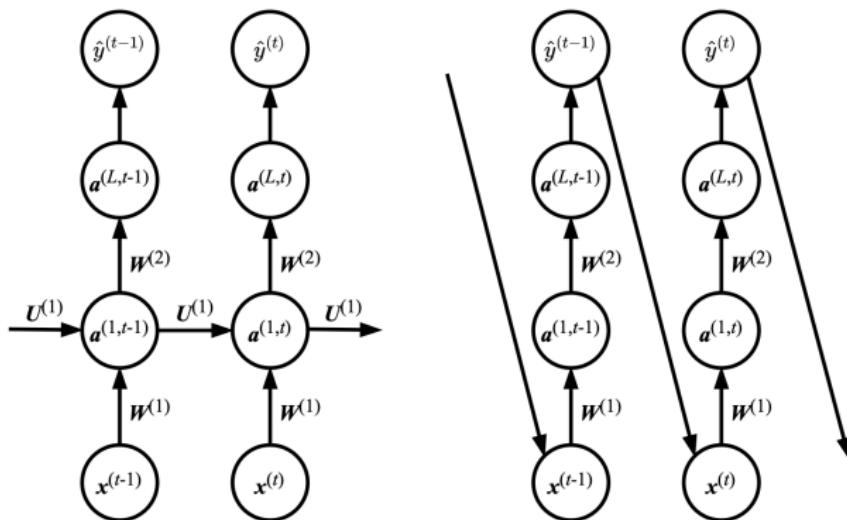
# Reduced Expressiveness

- The vanilla RNNs are universal in the sense that they can simulate Turing machines [6]



# Reduced Expressiveness

- The vanilla RNNs are universal in the sense that they can simulate Turing machines [6]
- Output-recurrent RNNs **cannot** simulate Turing machines and are strictly less powerful



# Reduced Expressiveness

- The vanilla RNNs are universal in the sense that they can simulate Turing machines [6]
- Output-recurrent RNNs **cannot** simulate Turing machines and are strictly less powerful
- The output  $\hat{y}^{(L,\cdot)}$ 's are explicitly trained to match training targets
  - Cannot capture all required information in the past to predict the future

