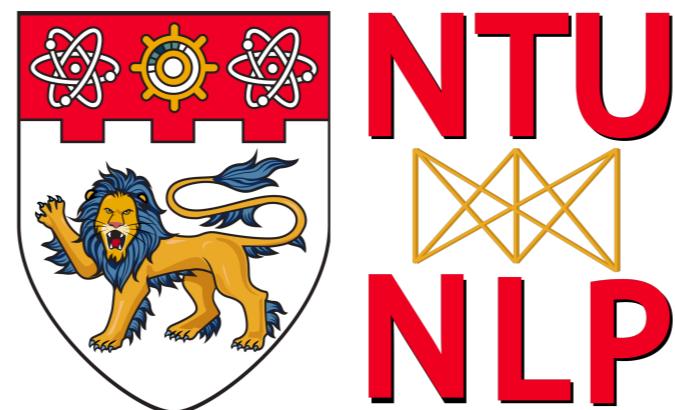


# Deep Learning for Natural Language Processing

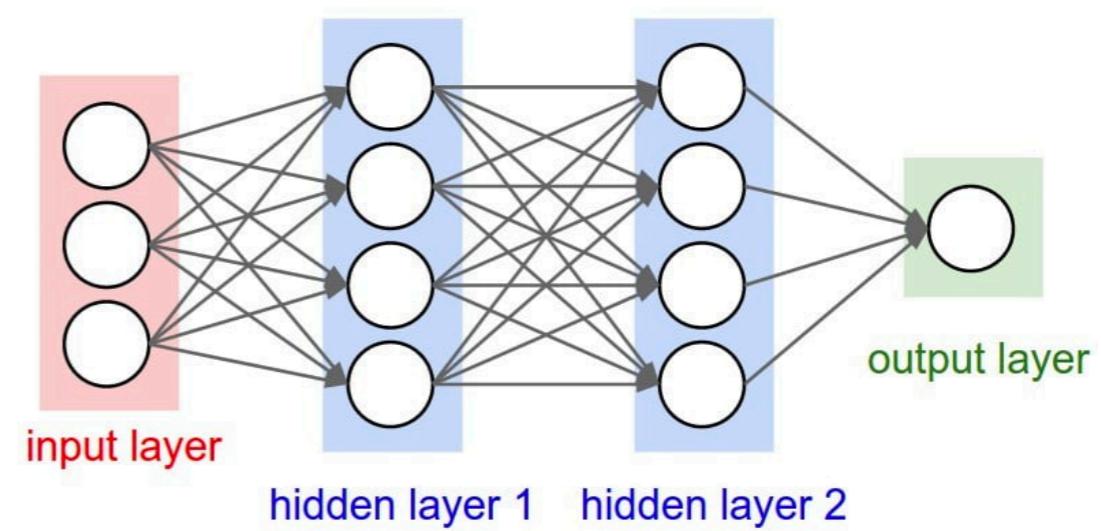
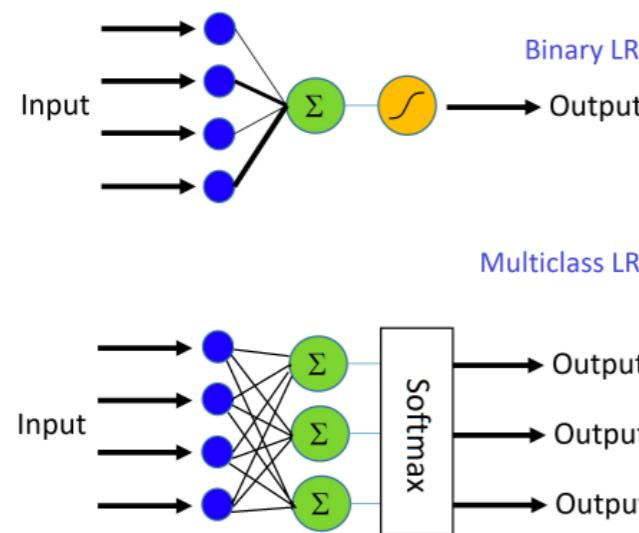
Shafiq Joty



Lecture 6: Recurrent Neural Nets

# Things Covered So Far

- Lecture 2: Linear Models (Linear & Logistic Regression)
  - Training with gradient descent and SGD
- Lecture 3: Extended linear models to MLP/FNN
  - Training with SGD (+ Backprop)



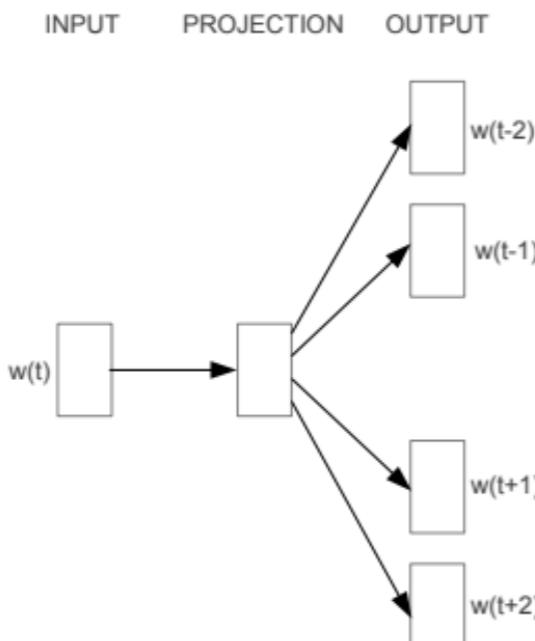
# Things Covered So Far

- Lecture 4: Used simple FNN as
  - A. Word2Vec Models
  - B. Cross-lingual Mapping

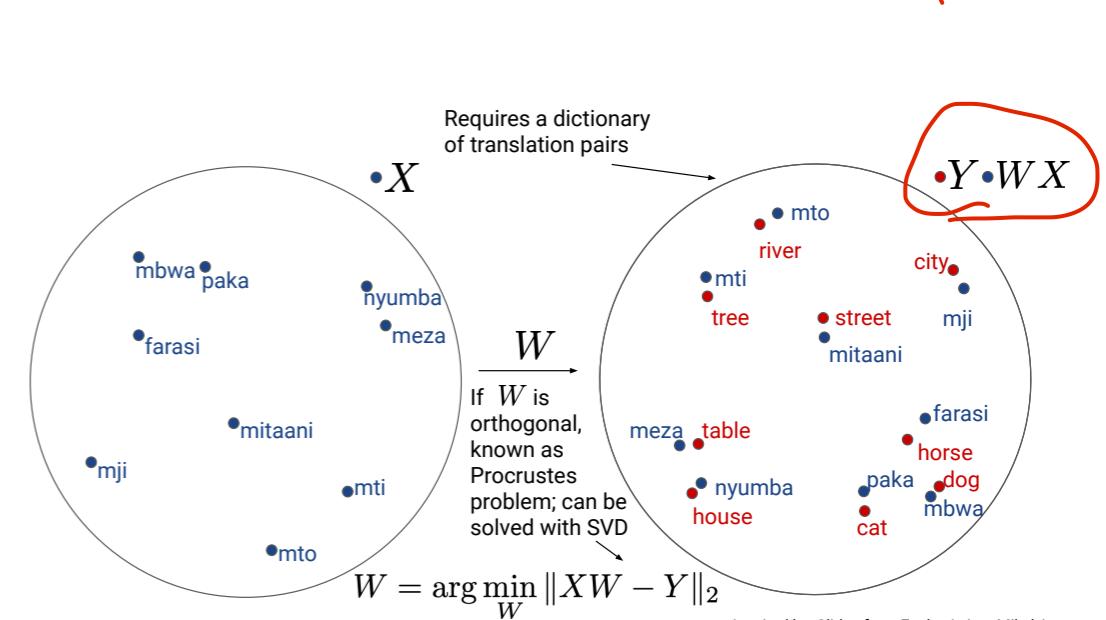
each word has a static embedding

learn mapping for these Embedding

?

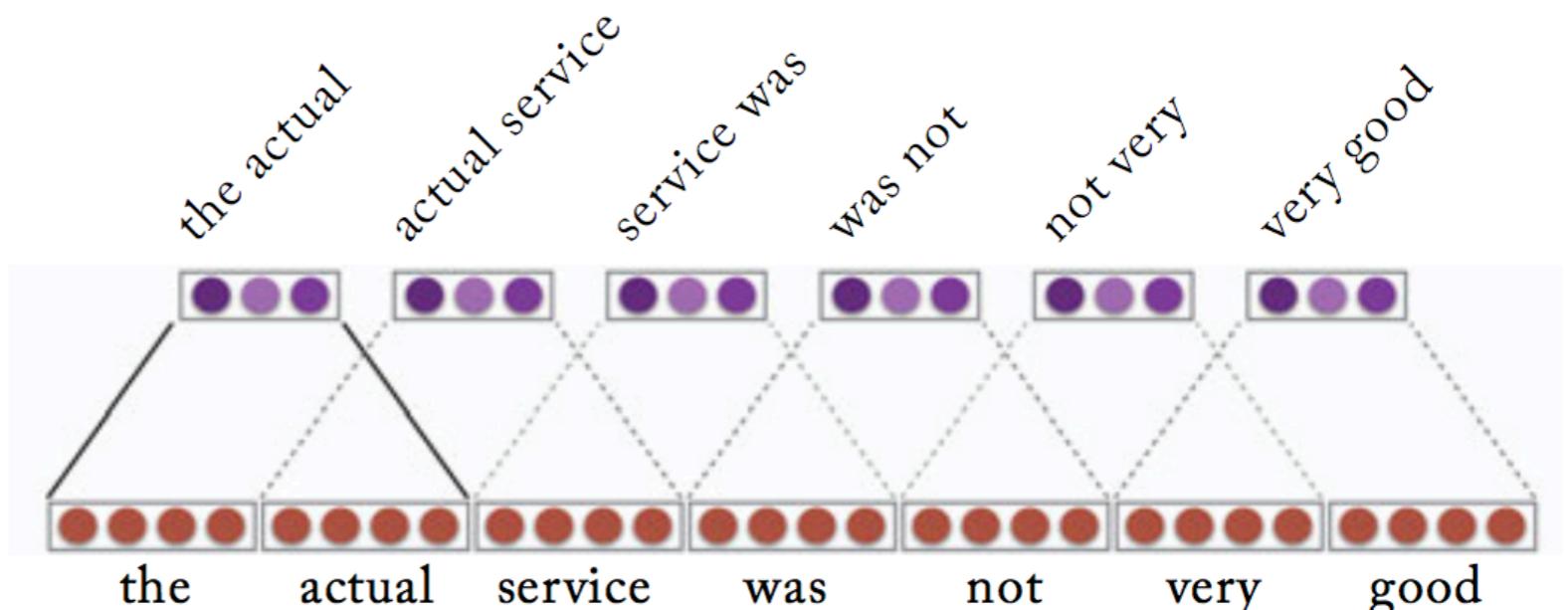
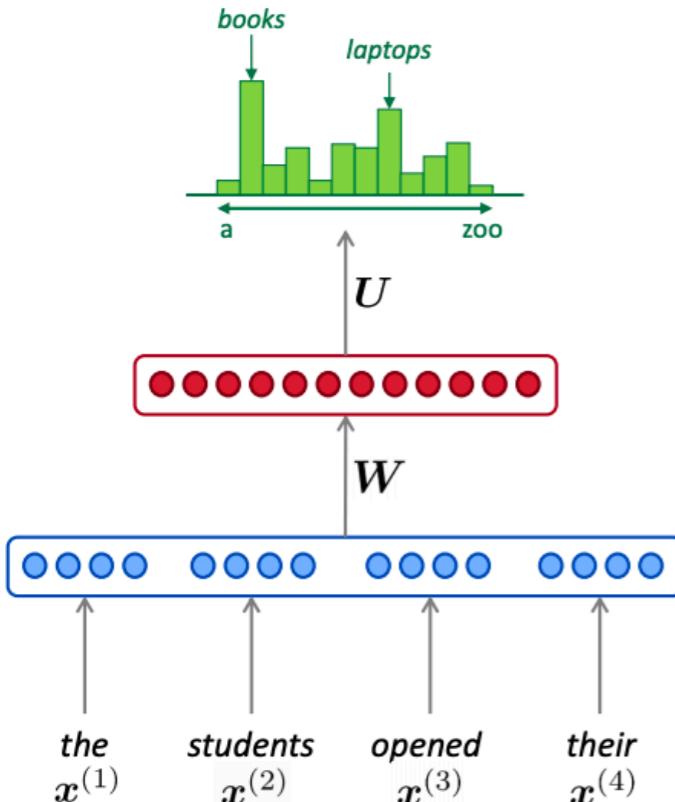


Skip-gram

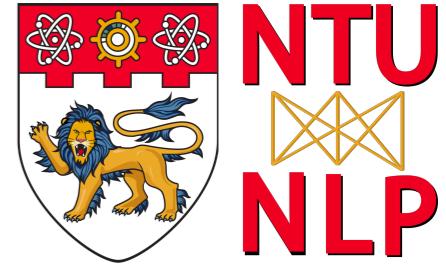


# Things Covered So Far

- Lecture 5: Window-based approach & Conv. Nets
  - Language modelling with window-based approach
  - Conv. Nets for encoding & classification



# Where we are



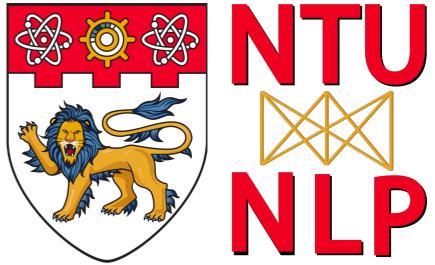
## Models/Algorithms

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

## NLP tasks/applications

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

# Lecture Plan

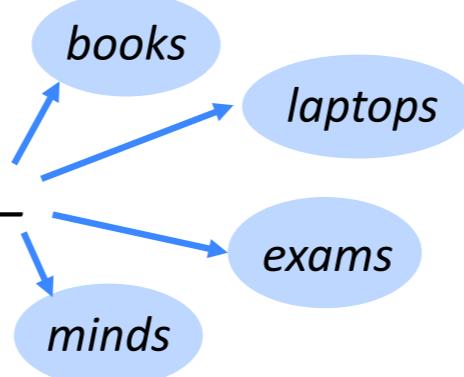


- Language Modeling
- Recurrent Neural Networks
- Backpropagation through time
- Text Generation with RNNs
- Sequence Tagging (NER, POS) with RNNs
- Sequence Classification with RNNs
- Bi-directional and stacked RNNs
- Gated RNNs (GRU, LSTM)

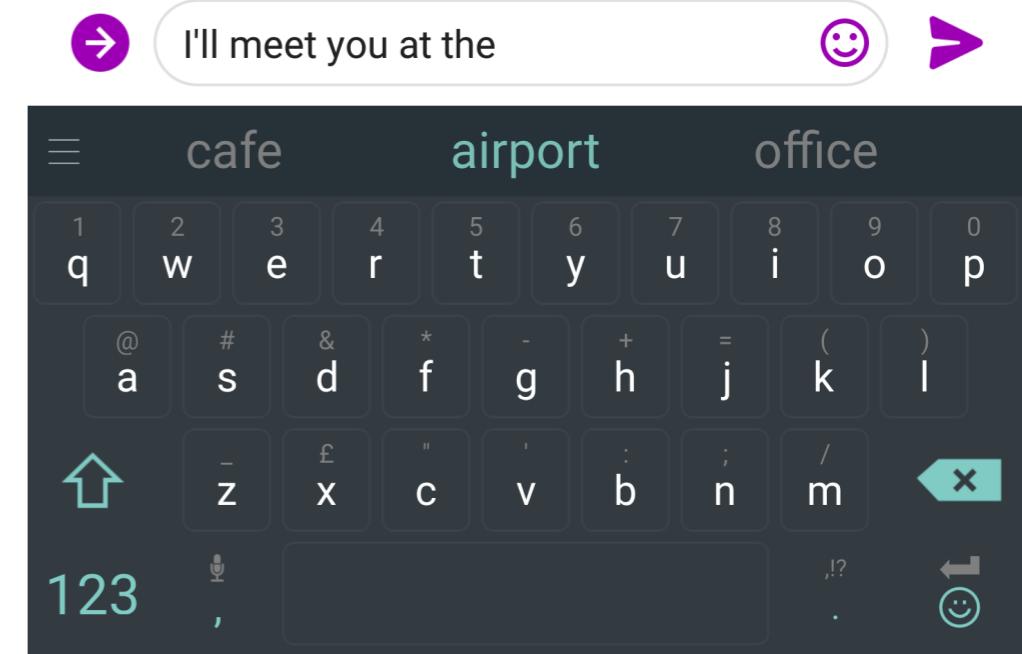
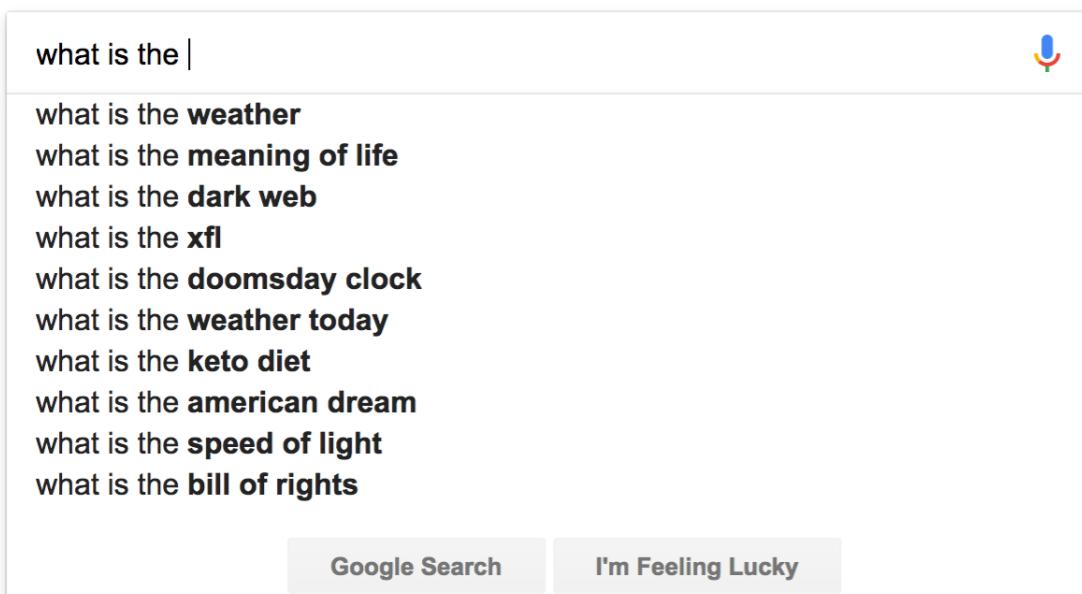
# Language Modelling

A language model takes a list of words (history/context), and attempts to predict the word that follows them

*the students opened their* \_\_\_\_\_



# Google



# Language Modelling

A language model takes a list of words (history/context), and attempts to predict the word that follows them

More formally: given a sequence of words  $x^{(1)}, x^{(2)}, \dots, x^{(t)}$ , compute the probability distribution of the next word  $x^{(t+1)}$ :

$$P(x^{(t+1)} | x^{(t)}, \dots, x^{(1)})$$

where  $x^{(t+1)}$  can be any word in the vocabulary  $V = \{w_1, \dots, w_{|V|}\}$

# Language Modeling

Language Models (LM) also assigns a probability to a piece of text.

For example, if we have some text  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}$ , then the probability of this text (according to the Language Model) is:

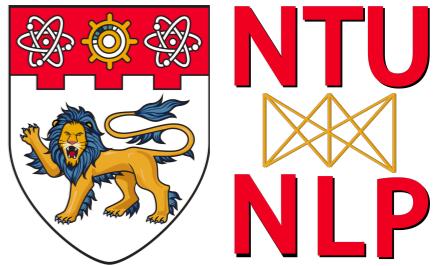
$$P(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}) = P(\mathbf{x}^{(1)}) \times P(\mathbf{x}^{(2)} | \mathbf{x}^{(1)}) \times \dots \times P(\mathbf{x}^{(T)} | \mathbf{x}^{(T-1)}, \dots, \mathbf{x}^{(1)})$$

$$= \prod_{t=1}^T P(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(1)})$$

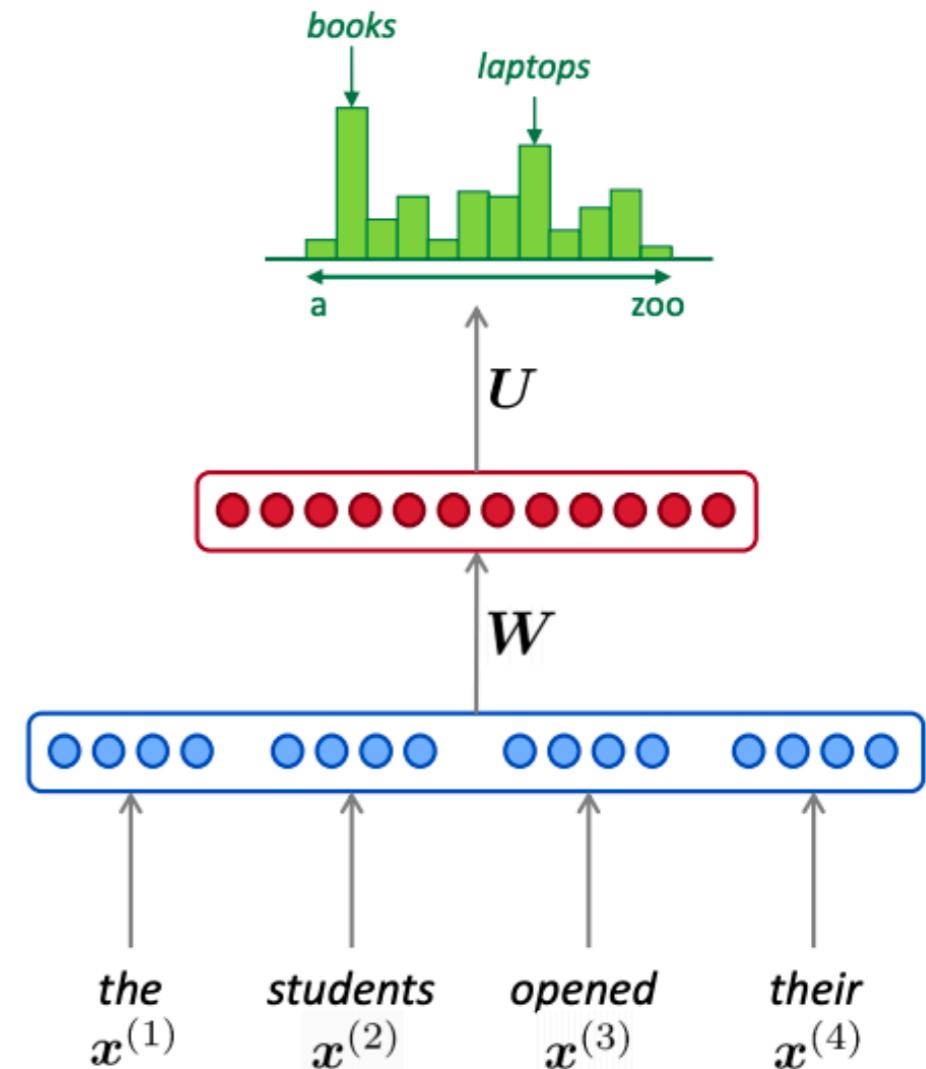
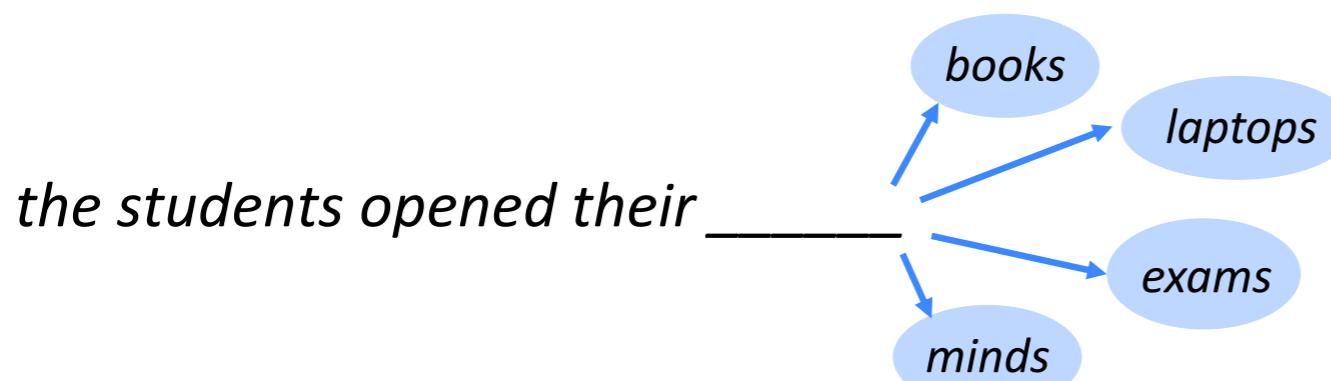


This is what our LM provides

# Language Modeling with Window Based Approach (Revisit)



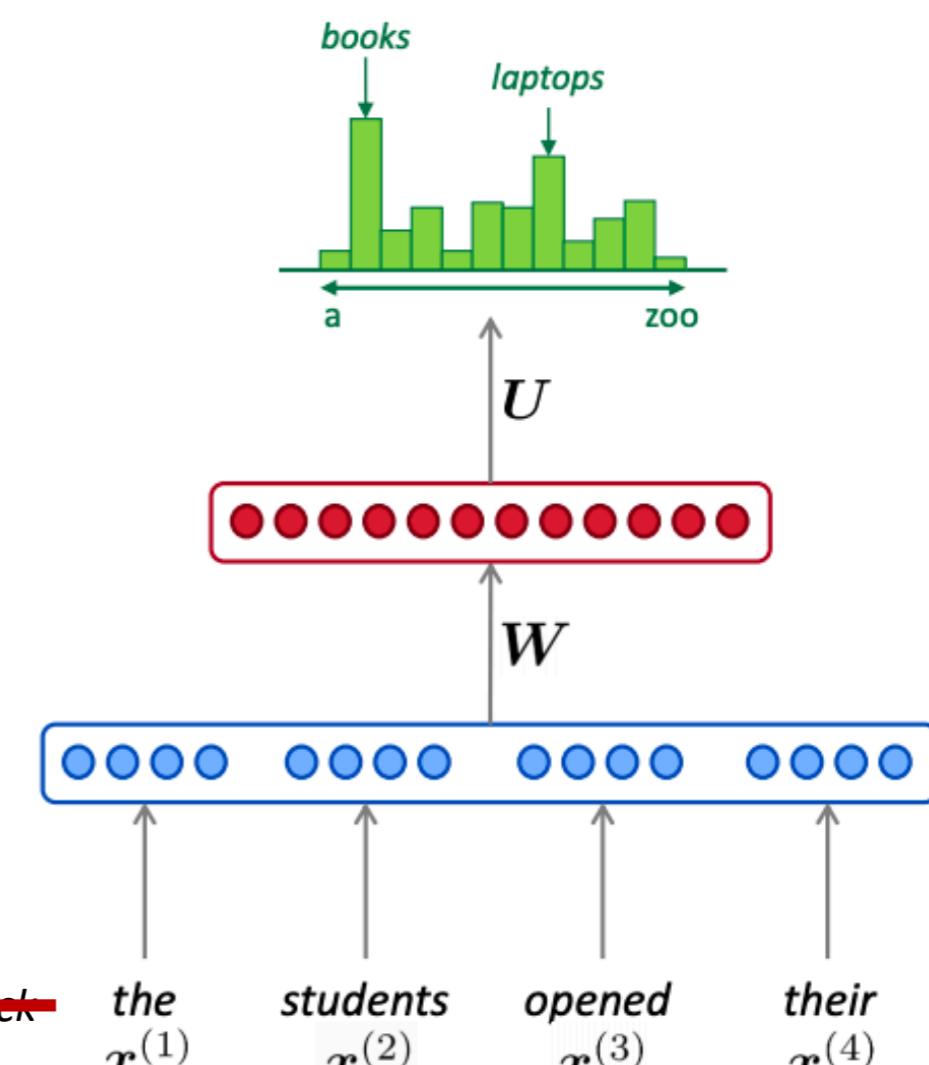
A language model takes a list of words (history/context), and attempts to predict the word that follows them



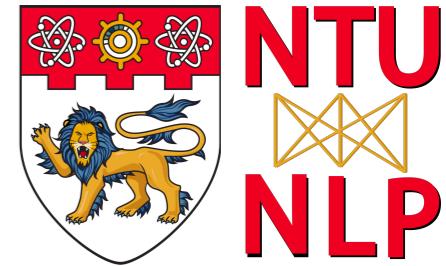
# Language Modeling with Window Based Approach (Revisit)

A language model takes a list of words (history/context), and attempts to predict the word that follows them

the students opened their \_\_\_\_\_  
here, we don't consider the entire context.  
if so our net architecture ~~was~~ the proctor started the clock ~~discard~~  
would be huge and in this case  
we should increase the volume of dataset in case overfitting.



# Language Modeling with Window Based Approach (Revisit)

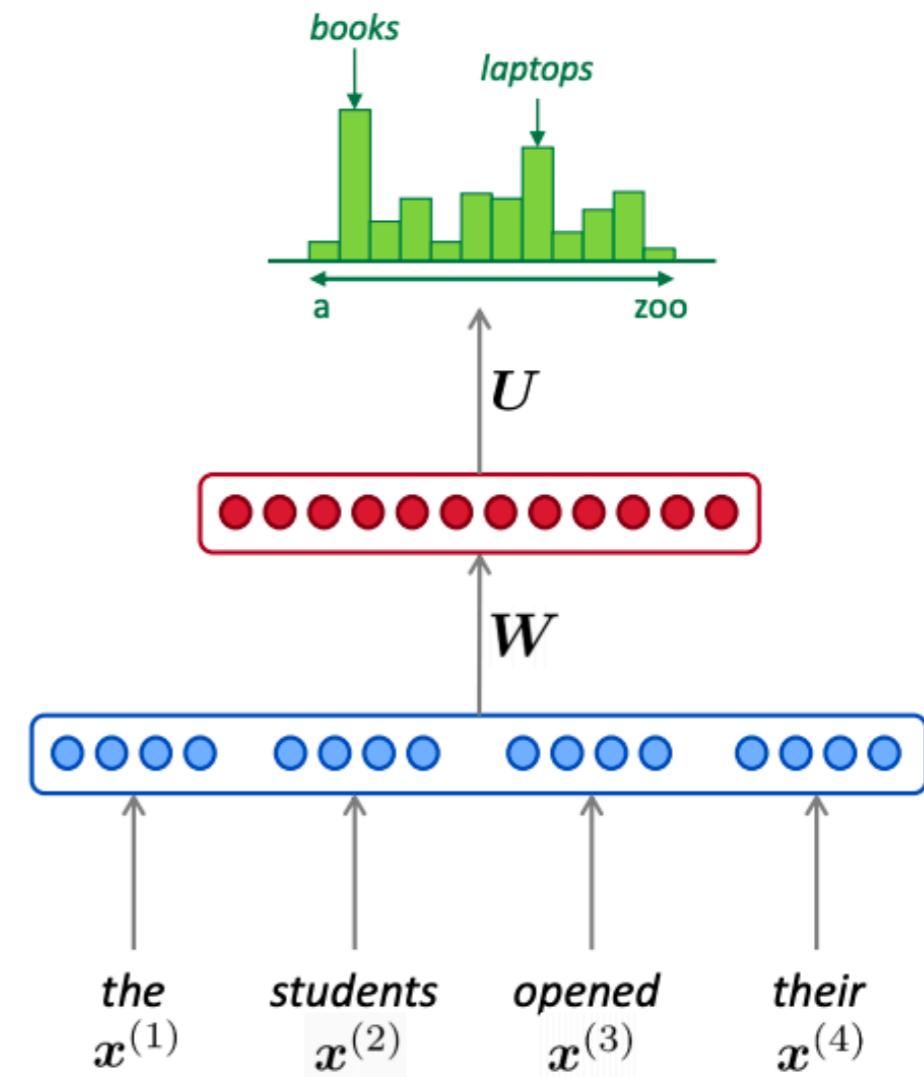


Improvements over n-gram LM:

- No sparsity problem
- Don't need to store the n-grams

Remaining **problems**:

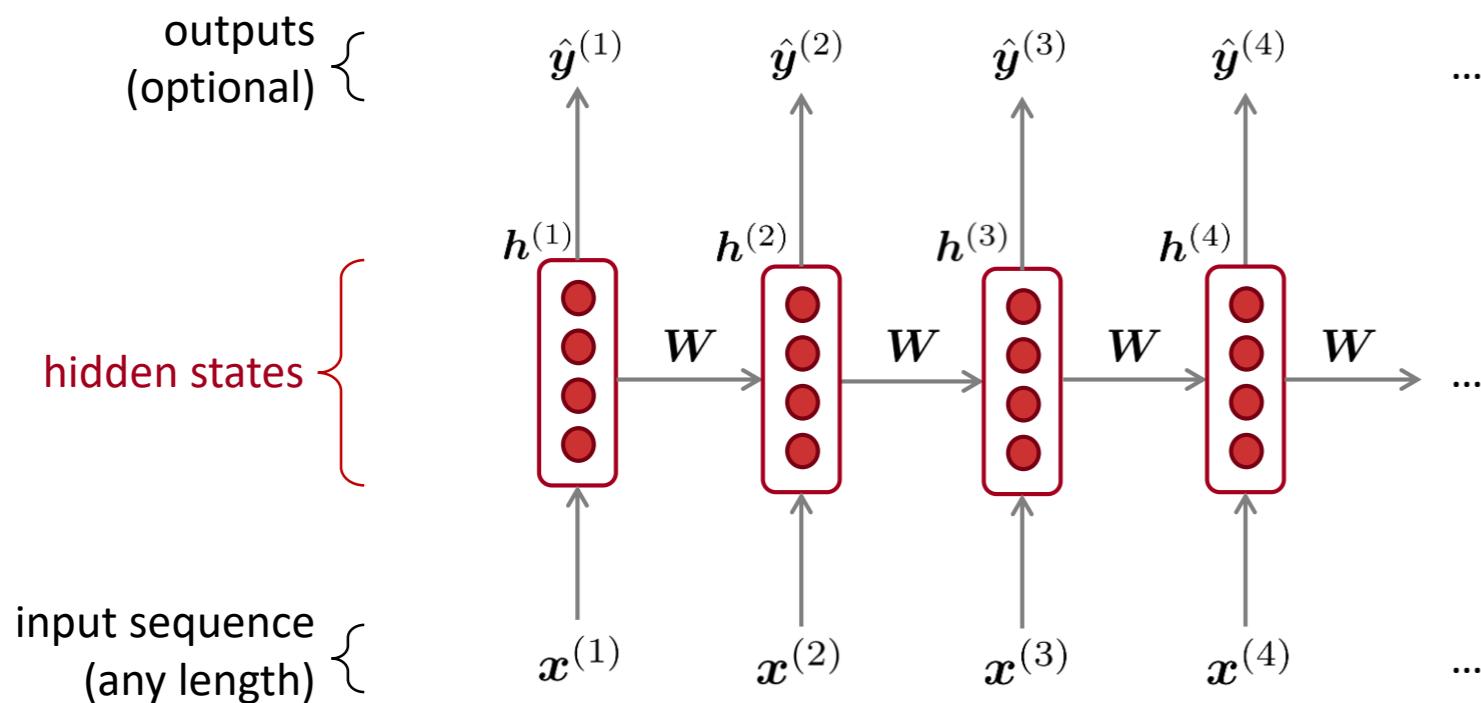
- Fixed window is **too small**
- Window can never be large enough! Enlarging window enlarges W
- **Each input vector is multiplied by completely different weights. No symmetry** in how the inputs are processed.
- Convolution? **Not suitable for variable length sequences**



We need a neural architecture that can process any length input

# Recurrent Neural Networks

Main idea: Apply the same weights repeatedly (recurrently)



Notice:

- Output at each step is optional
- Recurrence is done with hidden states

# A RNN Language Model

$$\hat{\mathbf{y}}^{(4)} = P(\mathbf{x}^{(5)} | \text{the students opened their})$$

output distribution

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{U}\mathbf{h}^{(t)} + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$

hidden states

$$\mathbf{h}^{(t)} = \sigma(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_e \mathbf{e}^{(t)} + \mathbf{b}_1)$$

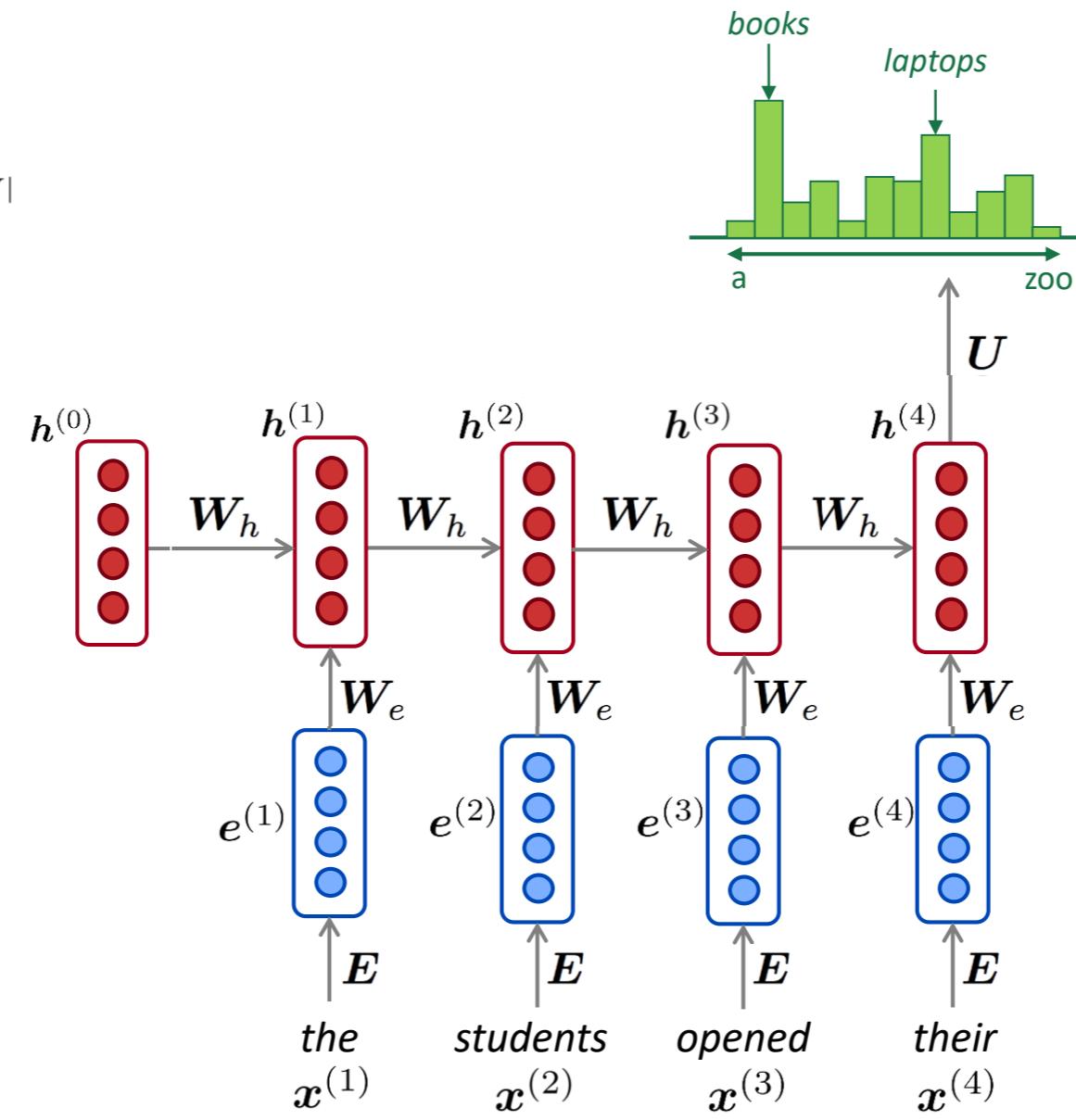
$\mathbf{h}^{(0)}$  is the initial hidden state

word embeddings

$$\mathbf{e}^{(t)} = \mathbf{E} \mathbf{x}^{(t)}$$

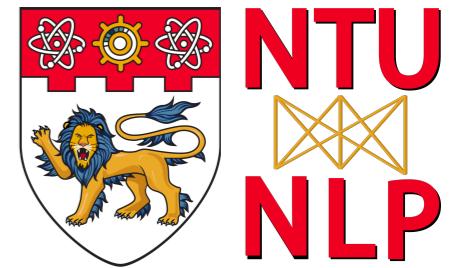
words / one-hot vectors

$$\mathbf{x}^{(t)} \in \mathbb{R}^{|V|}$$



the input sequence can be of arbitrary length

# Language Modeling with RNNs

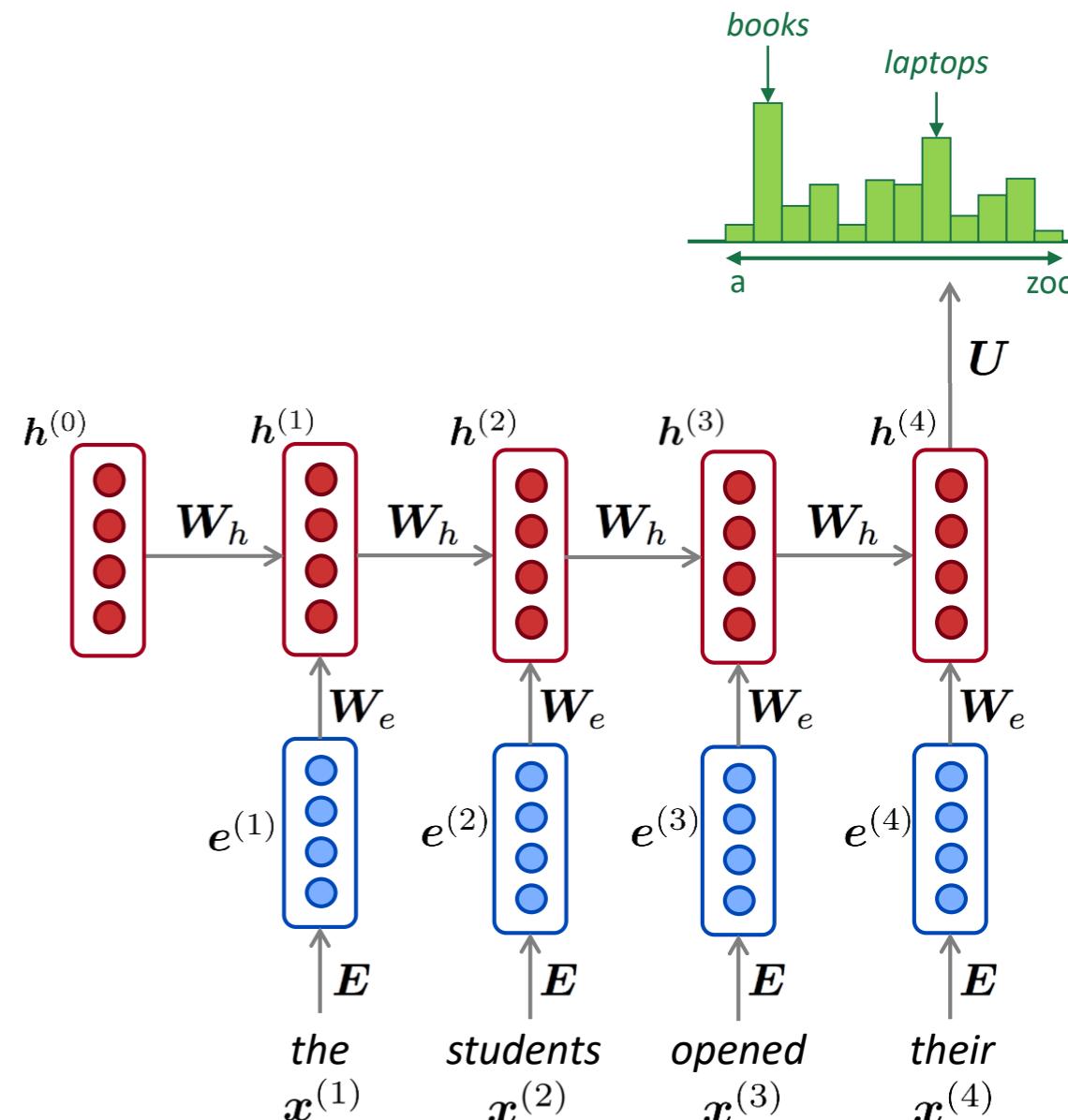


## Advantages:

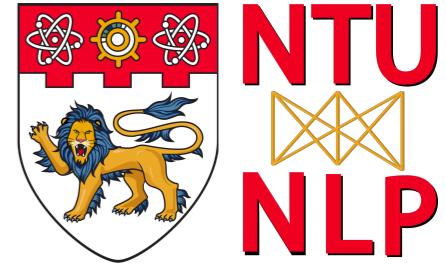
- Can process **any length** input
- Computation at step  $t$  can (in theory) use information from **many steps back**
- Model size doesn't increase for longer input
- Same weights applied on every timestep, so there is **symmetry** in how inputs are processed.

## Problems:

- Recurrent computation is **slow**
- In practice, **difficult to access** information from many steps back  
(remembers only recent states)



# Training a RNN-LM



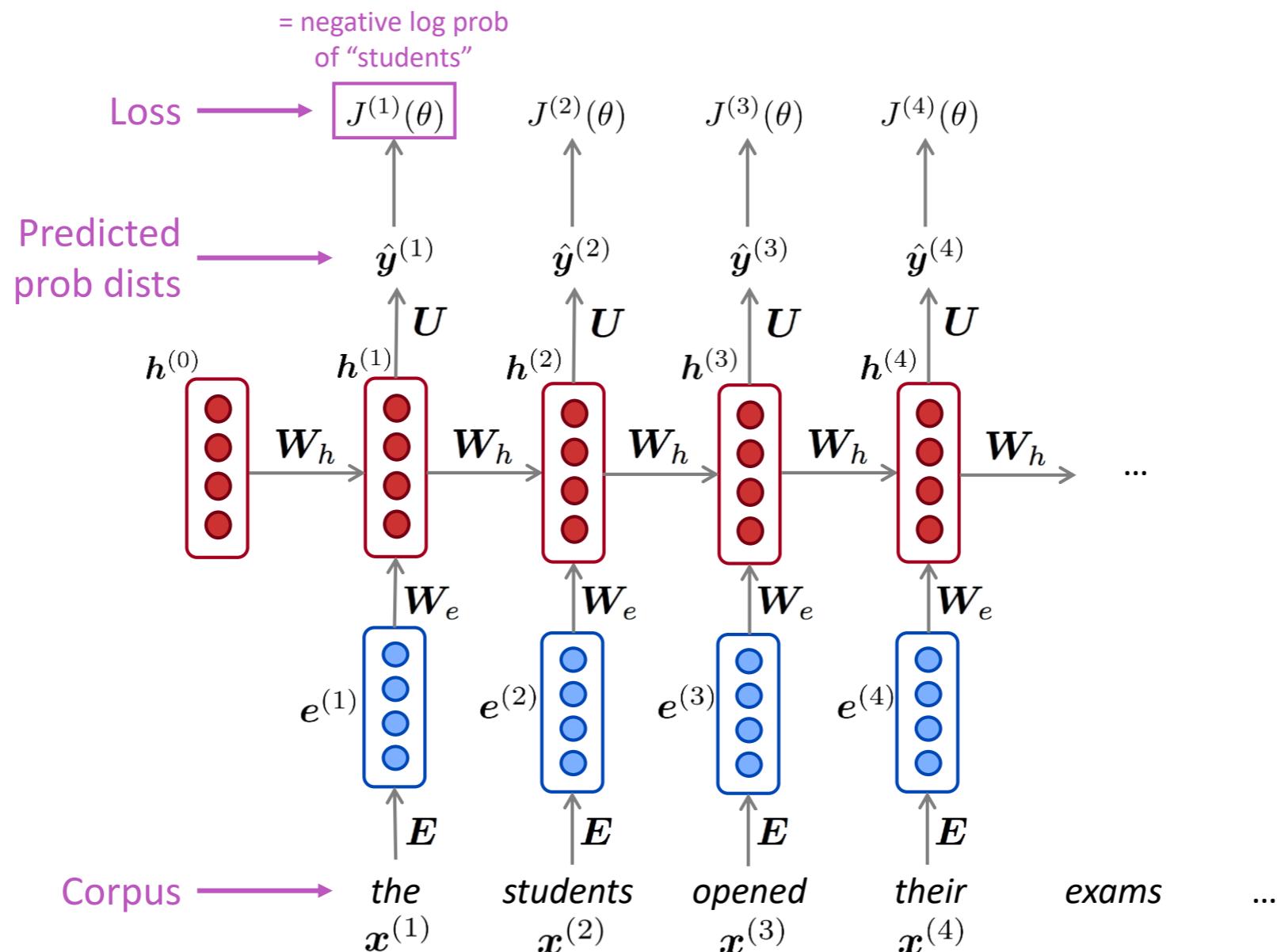
- Get a **big corpus of text** which is a sequence of words  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}$
- Feed into RNN-LM; compute output distribution  $\hat{\mathbf{y}}^{(t)}$  **for every step  $t$ .**
  - i.e. predict probability dist of *every word*, given words so far
- **Loss function** on step  $t$  is **cross-entropy** between predicted probability distribution  $\hat{\mathbf{y}}^{(t)}$ , and the true next word  $\mathbf{y}^{(t)}$  (one-hot for  $\mathbf{x}^{(t+1)}$ ):

$$J^{(t)}(\theta) = CE(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = - \sum_{w \in V} \mathbf{y}_w^{(t)} \log \hat{\mathbf{y}}_w^{(t)} = - \log \hat{\mathbf{y}}_{\mathbf{x}_{t+1}}^{(t)}$$

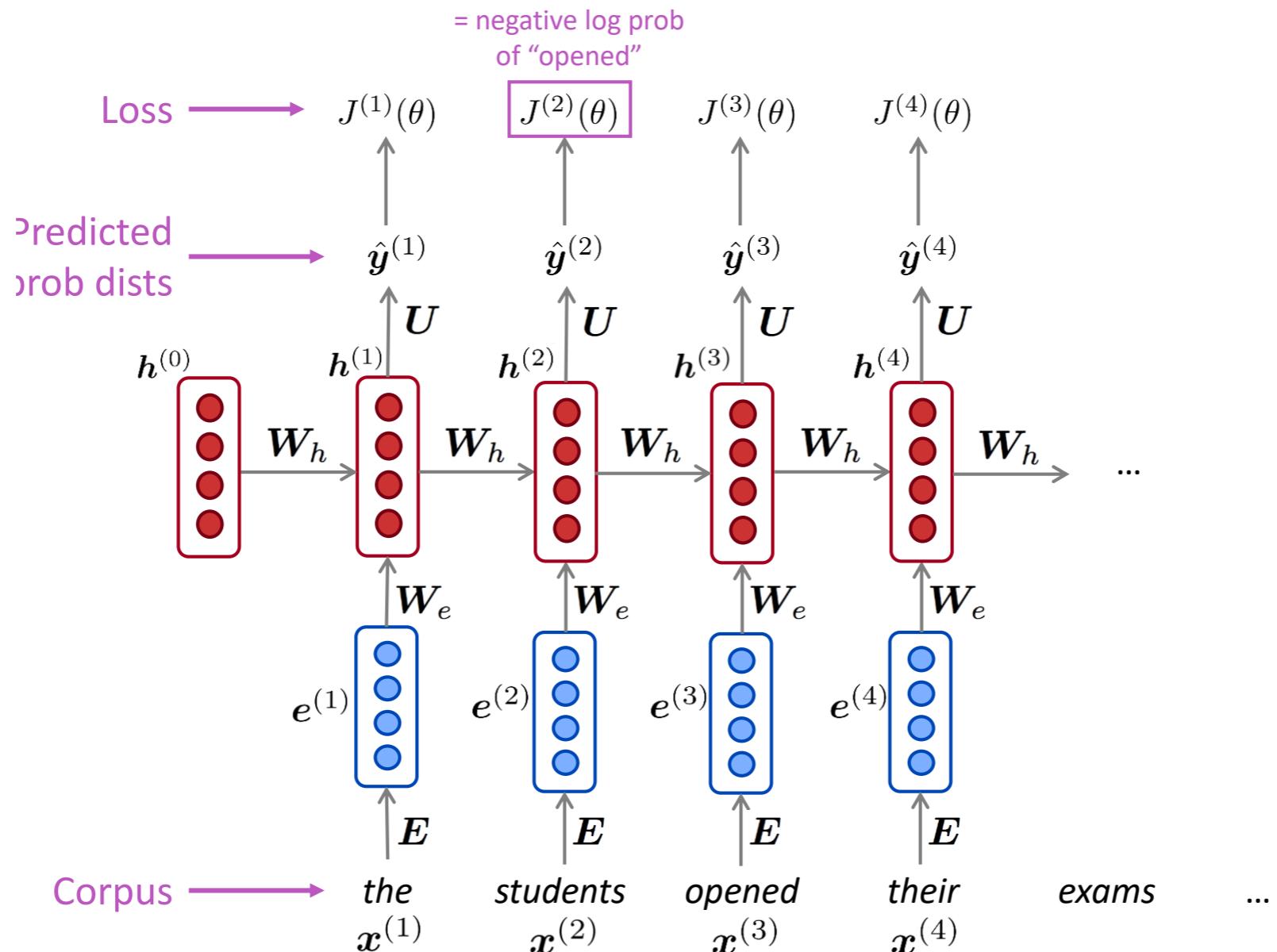
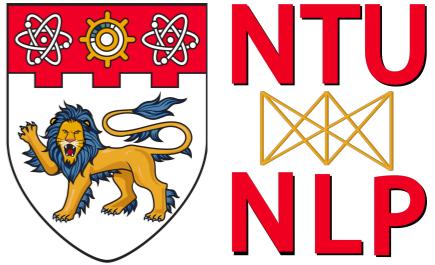
- Average this to get **overall loss** for entire training set:

$$J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta) = \frac{1}{T} \sum_{t=1}^T - \log \hat{\mathbf{y}}_{\mathbf{x}_{t+1}}^{(t)}$$

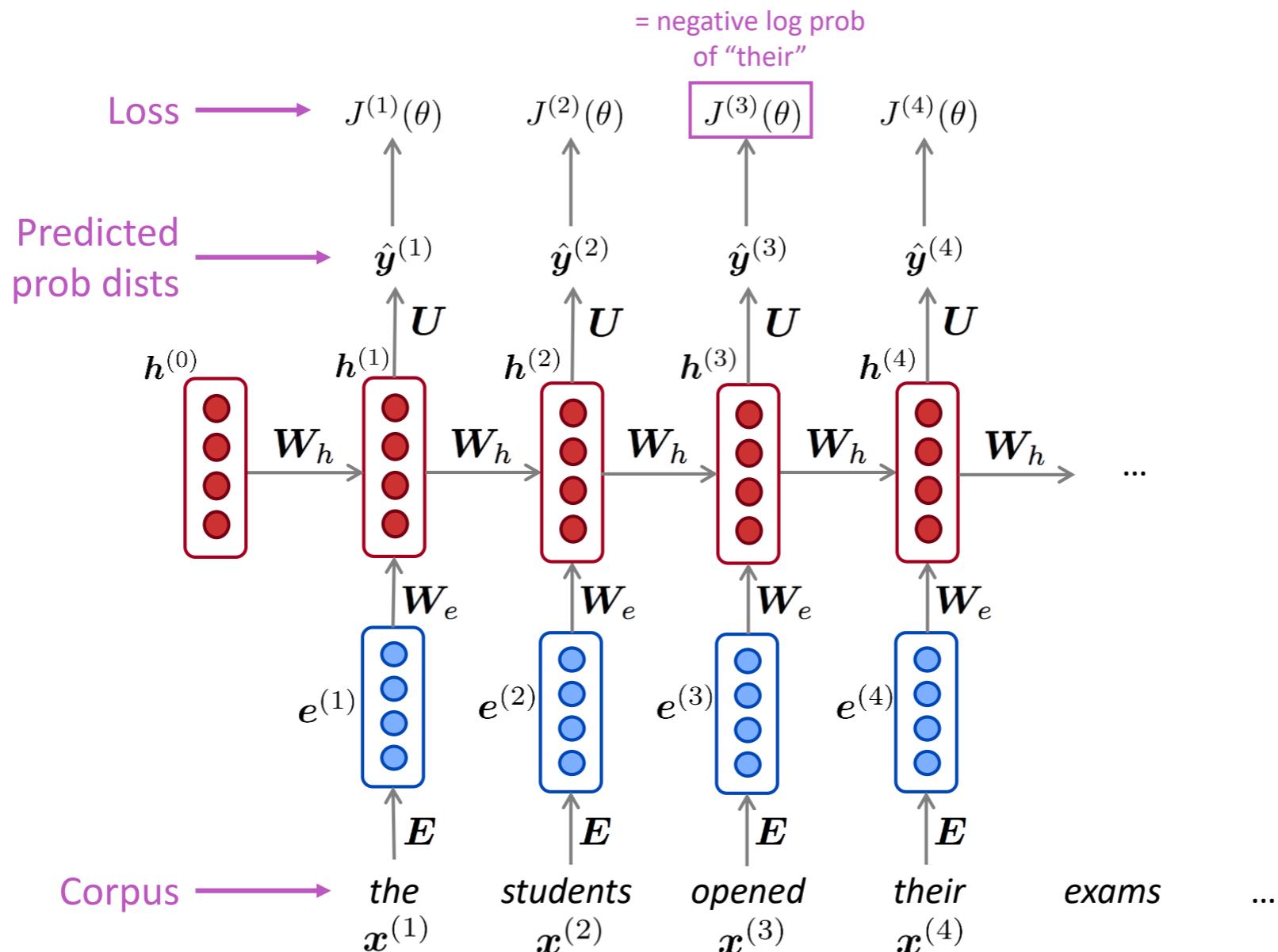
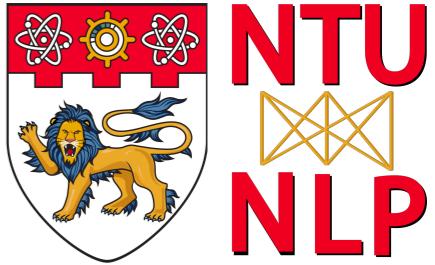
# Training a RNN-LM



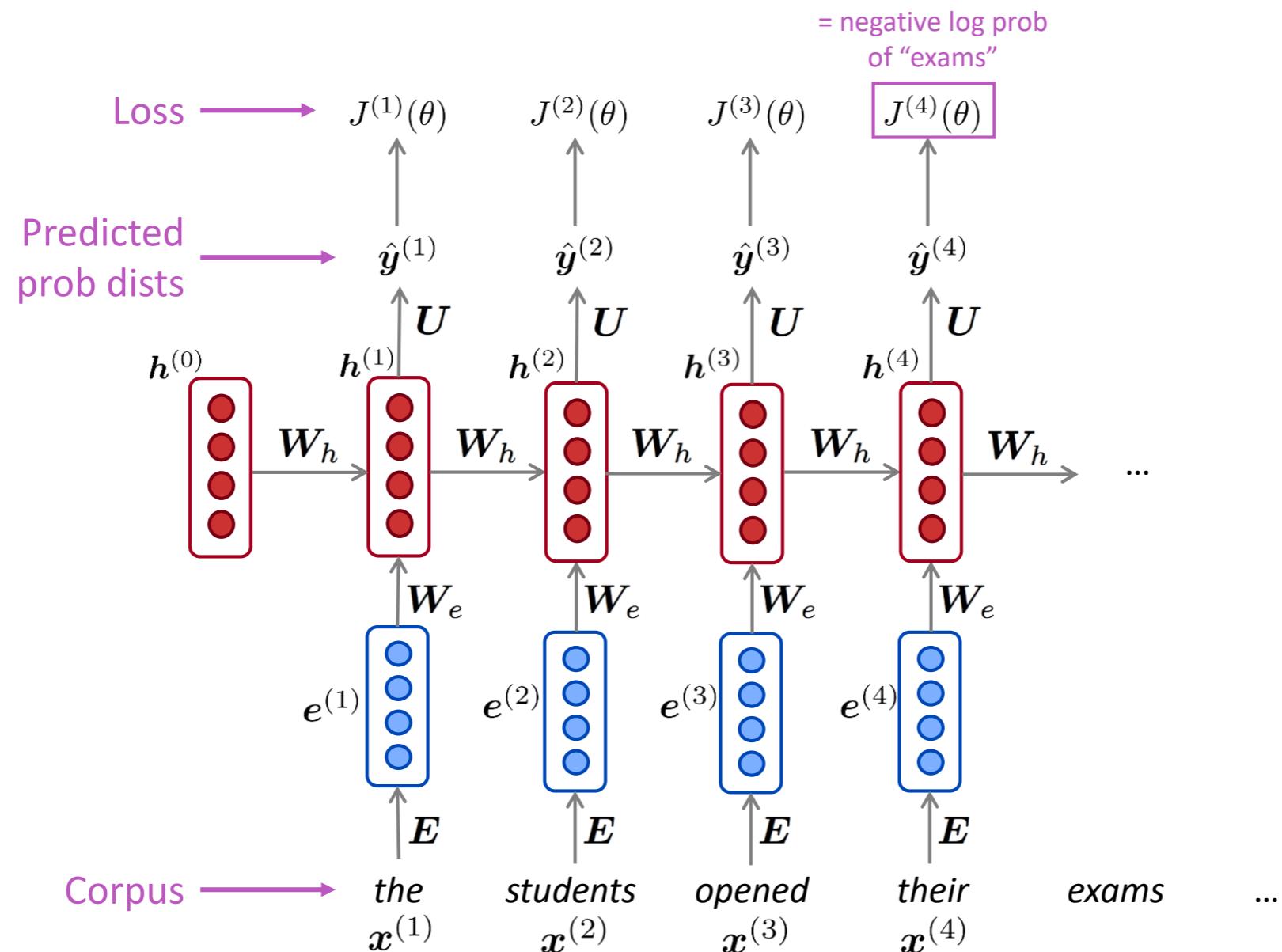
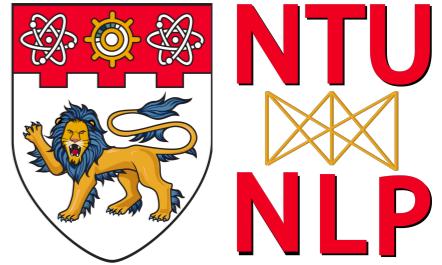
# Training a RNN-LM



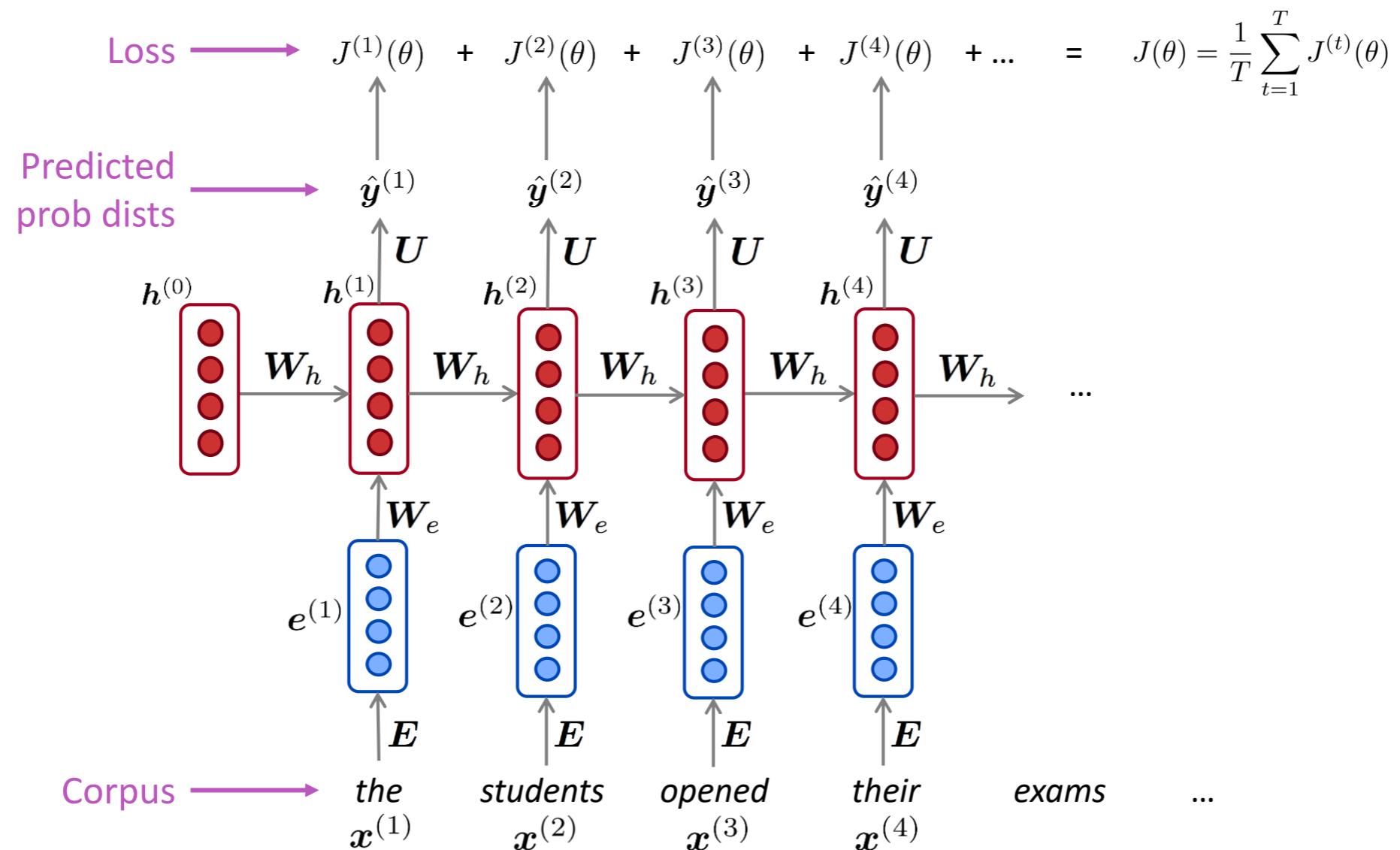
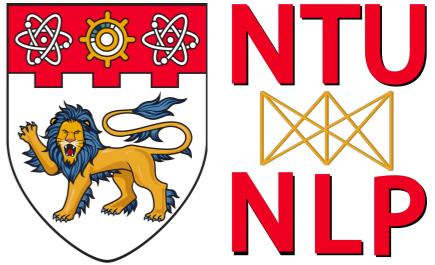
# Training a RNN-LM



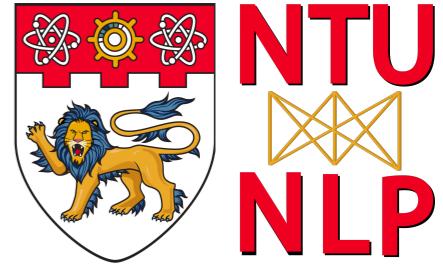
# Training a RNN-LM



# Training a RNN-LM

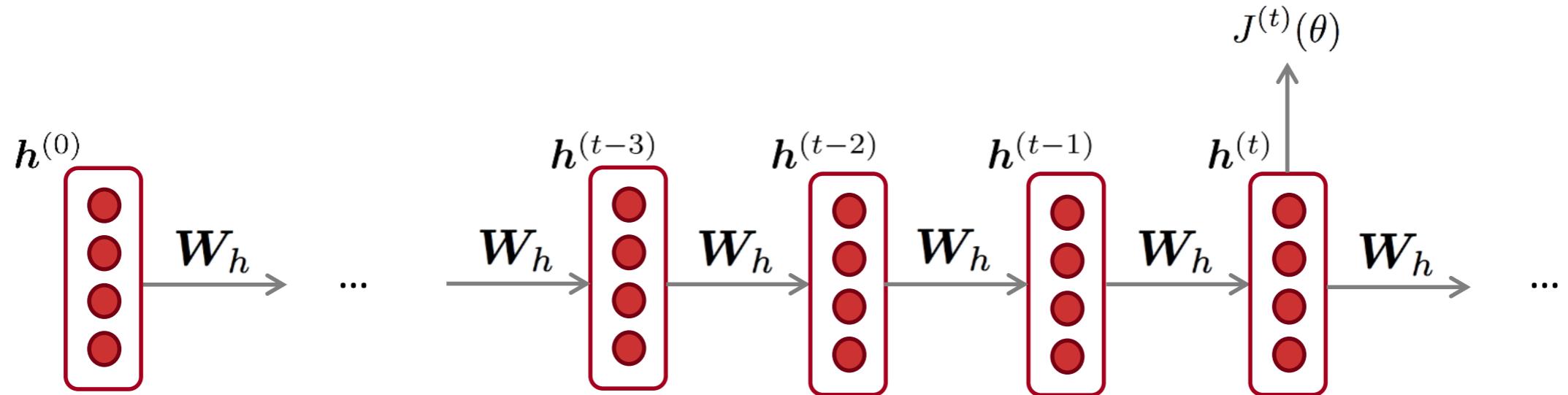
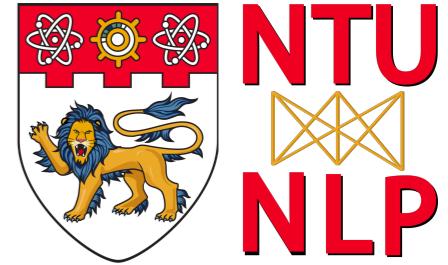


# Training a RNN-LM



- However: Computing loss and gradients across the entire corpus is way too expensive (no parallelisation)
- In practice, consider as a sentence (or a fixed length sequence)
- In SGD, we compute loss and gradients for batches and update, and then repeat

# Backpropagation for RNNs



The gradient w.r.t. a repeated weight is the sum of the gradient w.r.t. each time it appears

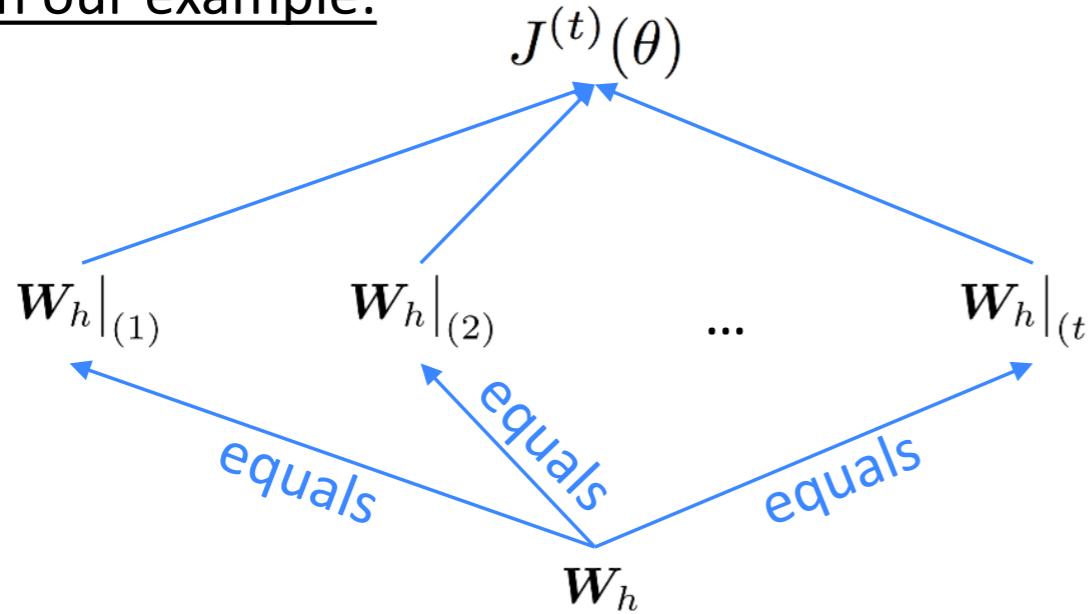
$$\frac{\partial J^{(t)}}{\partial W_h} = \sum_{i=1}^t \left. \frac{\partial J^{(t)}}{\partial W_h} \right|_{(i)}$$

# Backpropagation for RNNs

- Given a multivariable function  $f(x, y)$ , and two single variable functions  $x(t)$  and  $y(t)$ , here's what the multivariable chain rule says:

$$\underbrace{\frac{d}{dt} f(\textcolor{teal}{x}(t), \textcolor{red}{y}(t))}_{\text{Derivative of composition function}} = \frac{\partial f}{\partial \textcolor{teal}{x}} \frac{d\textcolor{teal}{x}}{dt} + \frac{\partial f}{\partial \textcolor{red}{y}} \frac{d\textcolor{red}{y}}{dt}$$

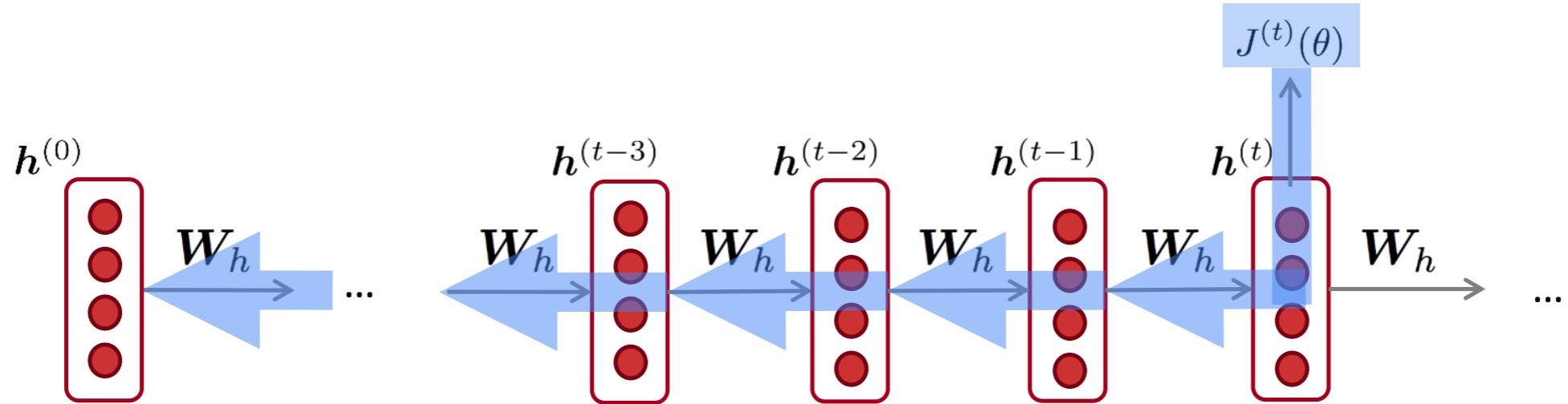
In our example:



Apply the multivariable chain rule:

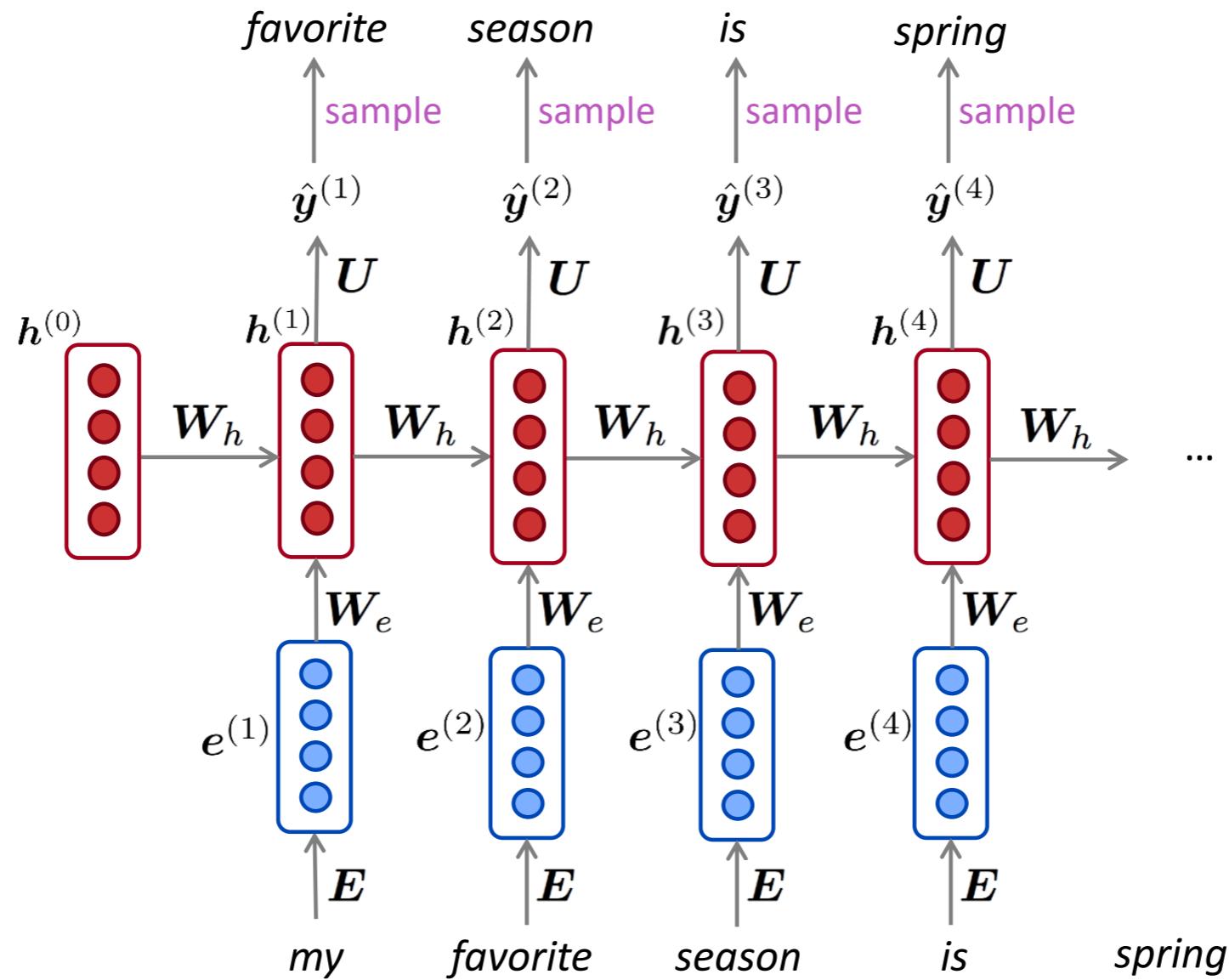
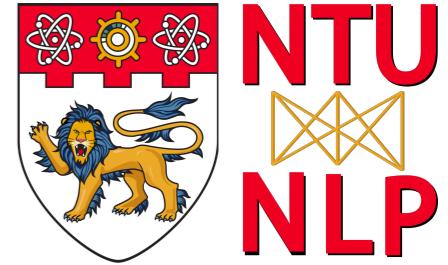
$$\begin{aligned} \frac{\partial J^{(t)}}{\partial \mathbf{W}_h} &= \sum_{i=1}^t \frac{\partial J^{(t)}}{\partial \mathbf{W}_h} \Big|_{(i)} \frac{\partial \mathbf{W}_h|_{(i)}}{\partial \mathbf{W}_h} \\ &= \sum_{i=1}^t \frac{\partial J^{(t)}}{\partial \mathbf{W}_h} \Big|_{(i)} \end{aligned}$$

# Backpropagation for RNNs



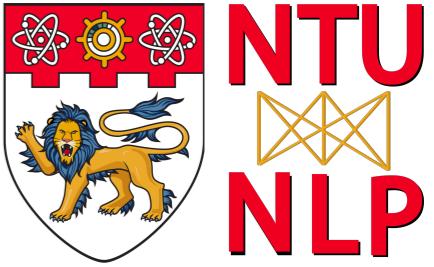
Backpropagate over timesteps  $i=t, \dots, 0$ , summing gradients as you go.  
 This algorithm is called “**backpropagation through time**”

# Generating Text with RNN-LM



Just like a n-gram LM, you can use an RNN-LM to **generate text by repeated sampling**. Sampled output is next step's input.

# Fun with RNN Generation



Char-RNN trained on Obama Speeches



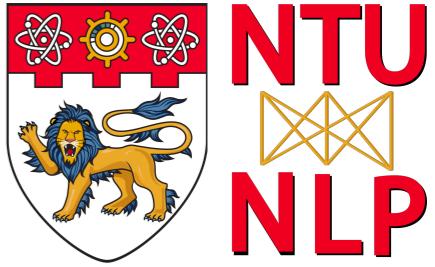
## SEED: Jobs

*Good afternoon. God bless you.*

*The United States will step up to the cost of a new challenges of the American people that will share the fact that we created the problem. They were attacked and so that they have to say that all the task of the final days of war that I will not be able to get this done. The promise of the men and women who were still going to take out the fact that the American people have fought to make sure that they have to be able to protect our part. It was a chance to stand together to completely look for the commitment to borrow from the American people. And the fact is the men and women in uniform and the millions of our country with the law system that we should be a strong stretches of the forces that we can afford to increase our spirit of the American people and the leadership of our country who are on the Internet of American lives.*

*Thank you very much. God bless you, and God bless the United States of America.*

# Fun with RNN Generation



Char-RNN trained on Obama Speeches

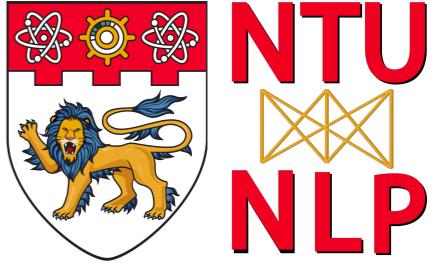
## SEED: Democracy

*Good morning. One of the borders will be able to continue to be here today. We have to say that the partnership was a partnership with the American people and the street continually progress that is a process and distant lasting peace and support that they were supporting the work of concern in the world. They were in the streets and communities that could have to provide steps to the people of the United States and Afghanistan. In the streets – the final decade of the country that will include the people of the United States of America. Now, humanitarian crisis has already rightly achieved the first American future in the same financial crisis that they can find reason to invest in the world.*

*Thank you very much. God bless you. God bless you. Thank you.*



# Fun with RNN Generation



Char-RNN trained on Harry Potter



## Part 1

*“The Malfoys!” said Hermione.*

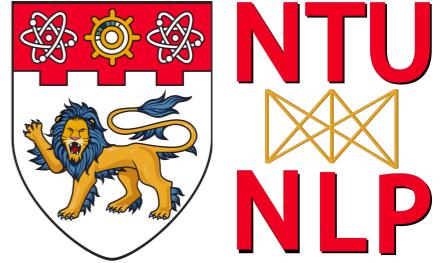
*Harry was watching him. He looked like Madame Maxime. When she strode up the wrong staircase to visit himself.*

*“I’m afraid I’ve definitely been suspended from power, no chance – indeed?” said Snape. He put his head back behind them and read groups as they crossed a corner and fluttered down onto their ink lamp, and picked up his spoon. The doorbell rang. It was a lot cleaner down in London. Hermione yelled. The party must be thrown by Krum, of course.*

*Harry collected fingers once more, with Malfoy. “Why, didn’t she never tell me. ...” She vanished. And then, Ron, Harry noticed, was nearly right.*

*“Now, be off,” said Sirius, “I can’t trace a new voice.”*

# Fun with RNN Generation



## Char-RNN trained on Food Recipe

"MMMMM----- Recipe via Meal-Master (tm) v8.05

Title: CARAMEL CORN GARLIC BEEF

Categories: Soups, Desserts

Yield: 10 Servings

2 tb Parmesan cheese, ground

1/4 ts Ground cloves

-- diced

1 ts Cayenne pepper

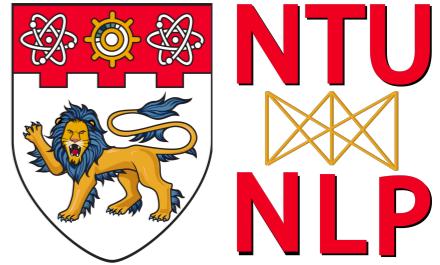
Cook it with the batter. Set aside to cool. Remove the peanut oil in a small saucepan and pour into the margarine until they are soft. Stir in a mixer (dough). Add the chestnuts, beaten egg whites, oil, and salt and brown sugar and sugar; stir onto the boquly brown it.

The recipe from an oiled by fried and can. Beans, by Judil Cookbook, Source: Pintore, October, by Chocolates, Breammons of Jozen, Empt.com



And many more ...

# Evaluating Language Models



We use **perplexity** as the standard evaluation metric for LMs

$$\text{perplexity} = \prod_{t=1}^T \left( \frac{1}{P_{\text{LM}}(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})} \right)^{1/T}$$

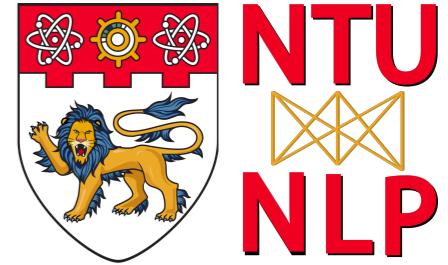
Inverse probability of corpus, according to Language Model

Normalized by  
number of words

This is equal to the **exponential of the cross-entropy loss**

$$= \prod_{t=1}^T \left( \frac{1}{\hat{y}_{\mathbf{x}_{t+1}}^{(t)}} \right)^{1/T} = \exp \left( \frac{1}{T} \sum_{t=1}^T -\log \hat{y}_{\mathbf{x}_{t+1}}^{(t)} \right) = \exp(J(\theta))$$

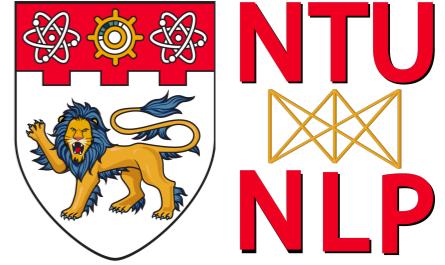
# Language Modeling Results



RNNs improve results quite significantly

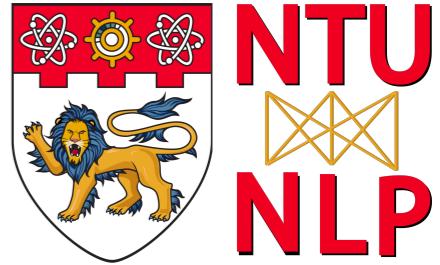
Model	Perplexity
Interpolated Kneser-Ney 5-gram (Chelba et al., 2013)	67.6
RNN-1024 + MaxEnt 9-gram (Chelba et al., 2013)	51.3
RNN-2048 + BlackOut sampling (Ji et al., 2015)	68.3
Sparse Non-negative Matrix factorization (Shazeer et al., 2015)	52.9
LSTM-2048 (Jozefowicz et al., 2016)	43.7
2-layer LSTM-8192 (Jozefowicz et al., 2016)	30
<b>Ours small</b> (LSTM-2048)	43.9
<b>Ours large</b> (2-layer LSTM-2048)	39.8

# Why Language Modeling is Important



- A **benchmark** task to track our progress on understanding language
  - An important **component** of many NLP tasks, especially those involving **generating text** or **estimating the probability** of a text
    - Speech recognition
    - Spelling/grammar correction
    - Authorship identification
    - Machine translation
    - Summarization
    - Dialogue
    - etc.
- (BERT)
- to learn contextual embedding in unsupervised way.*

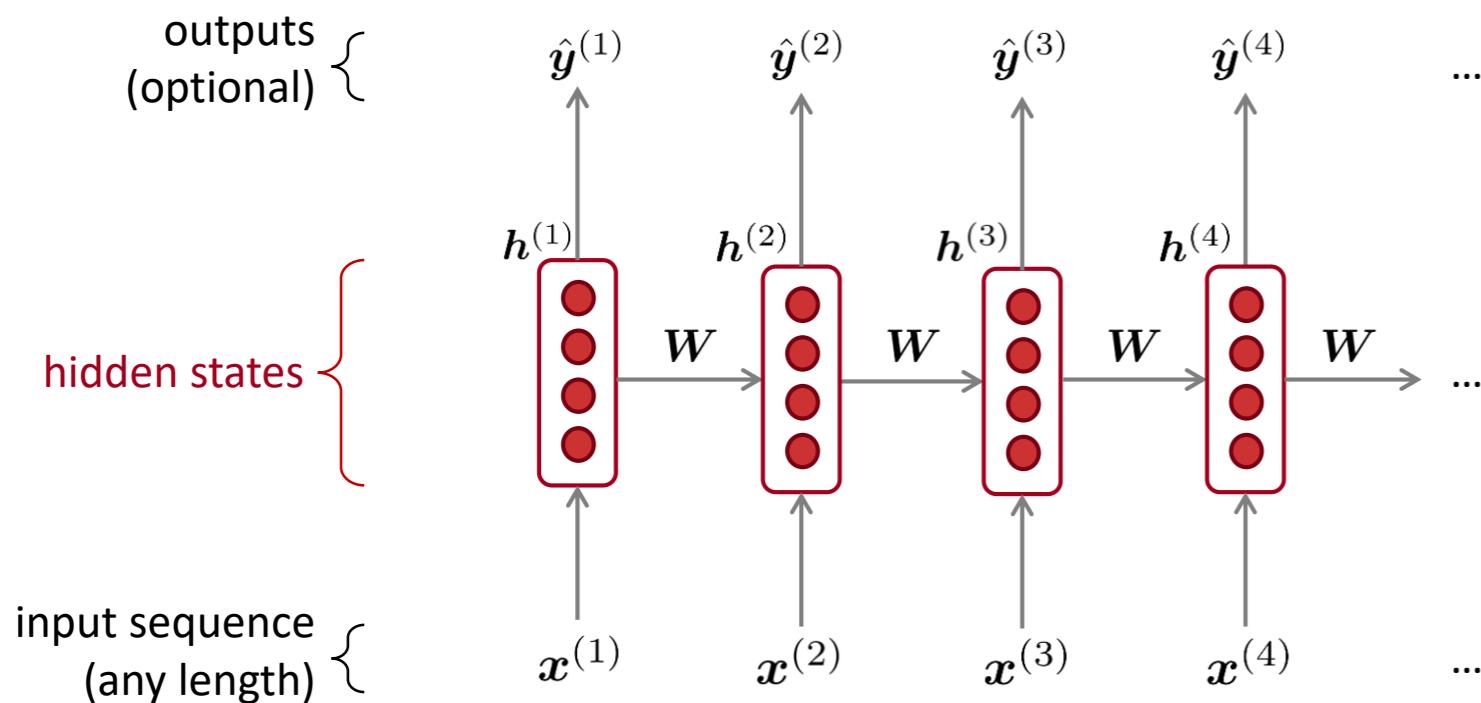
# Lecture Plan



- Language Modeling
  - Recurrent Neural Networks
  - Backpropagation through time
  - Text Generation with RNNs
  - Sequence Tagging (NER, POS) with RNNs
  - Sequence Classification with RNNs
  - Bi-directional and stacked RNNs
  - Gated RNNs (GRU, LSTM)
- not just Language Modeling*
- general model can be used in many other purposes*

# Recurrent Neural Networks (Revisit)

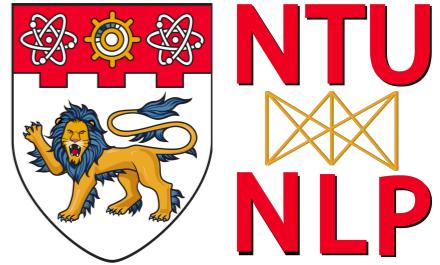
Main idea: Apply the same weights repeatedly (recurrently)



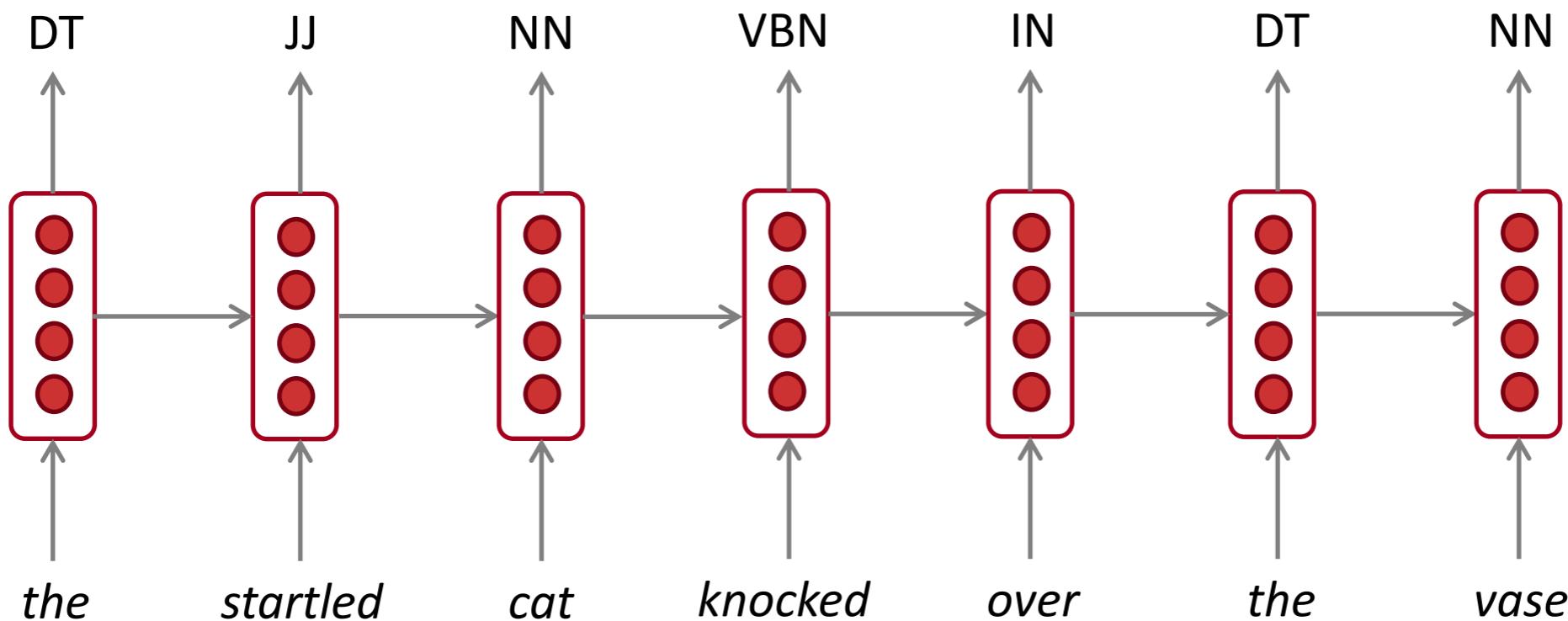
Notice:

- Output at each step is optional
- Recurrence is done with hidden states

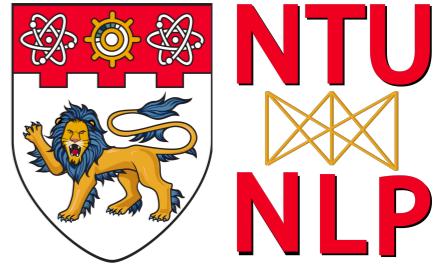
# Sequence Tagging with RNNs



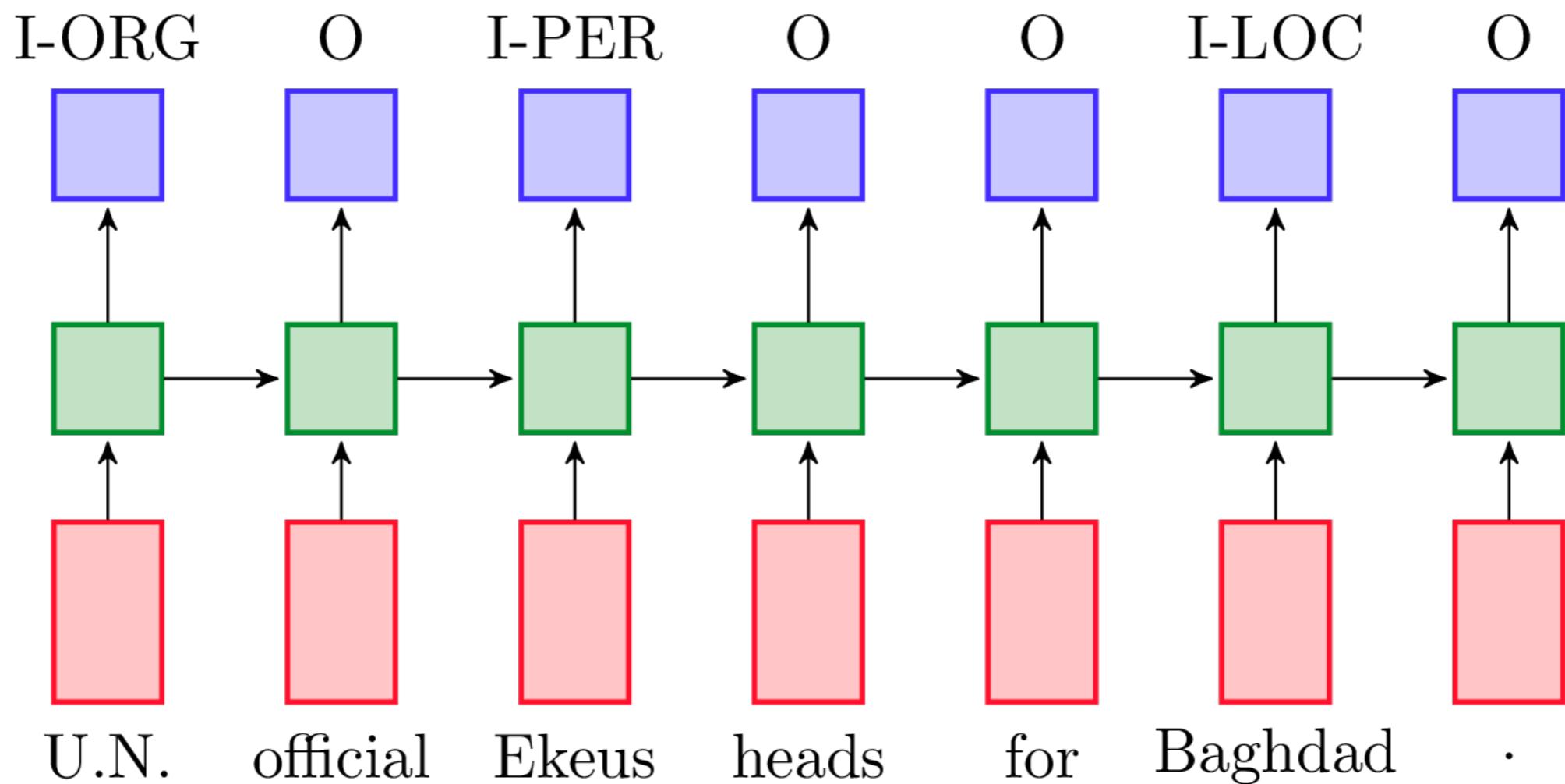
- POS Tagging:



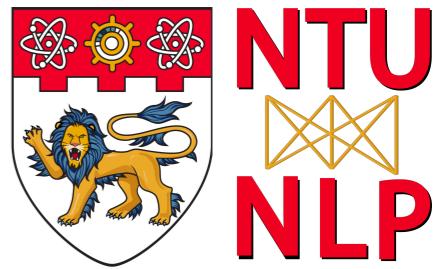
# Sequence Tagging with RNNs



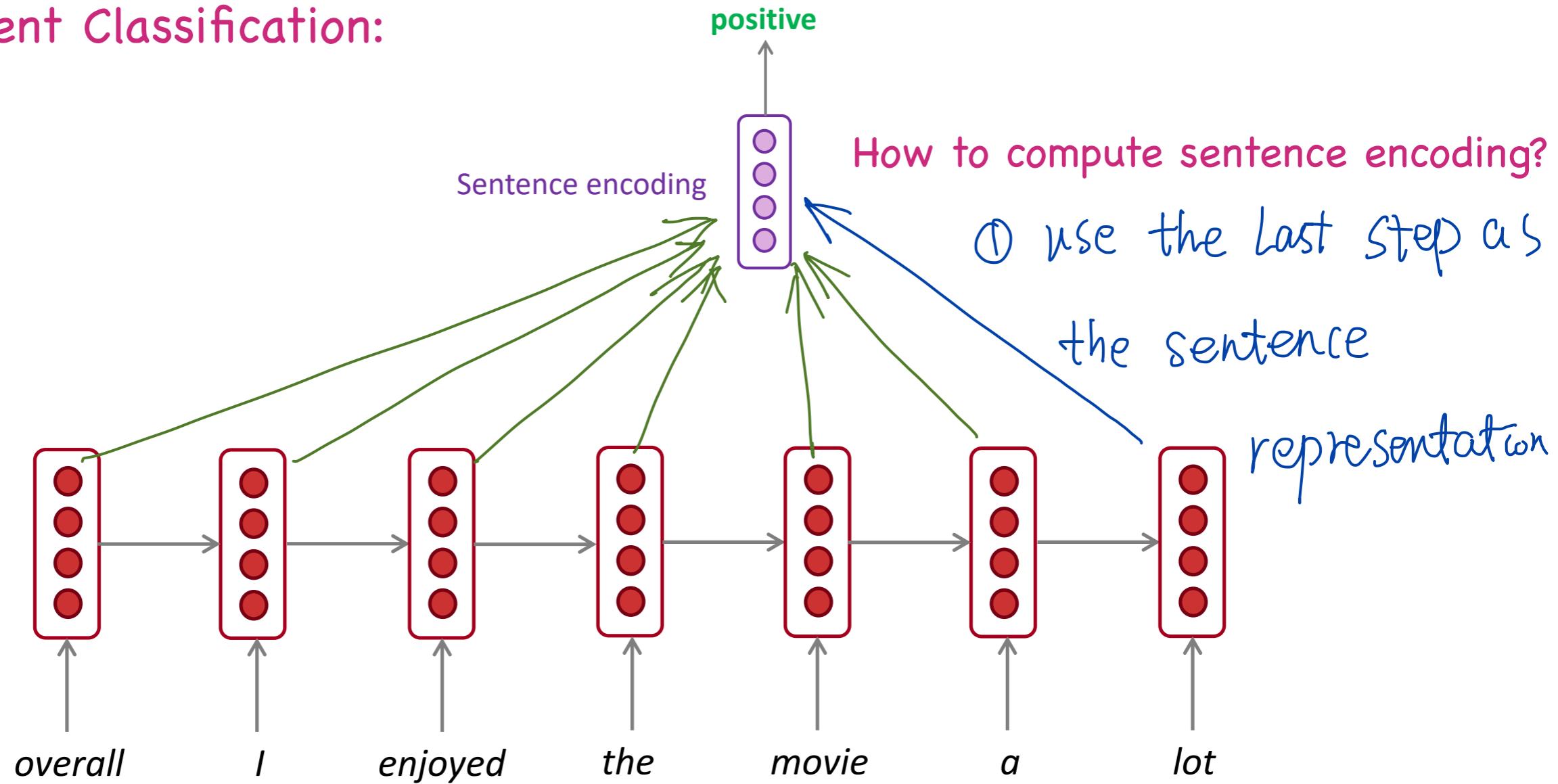
- NER Tagging:



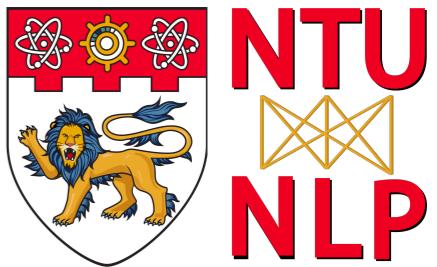
# Sequence Classification with RNNs



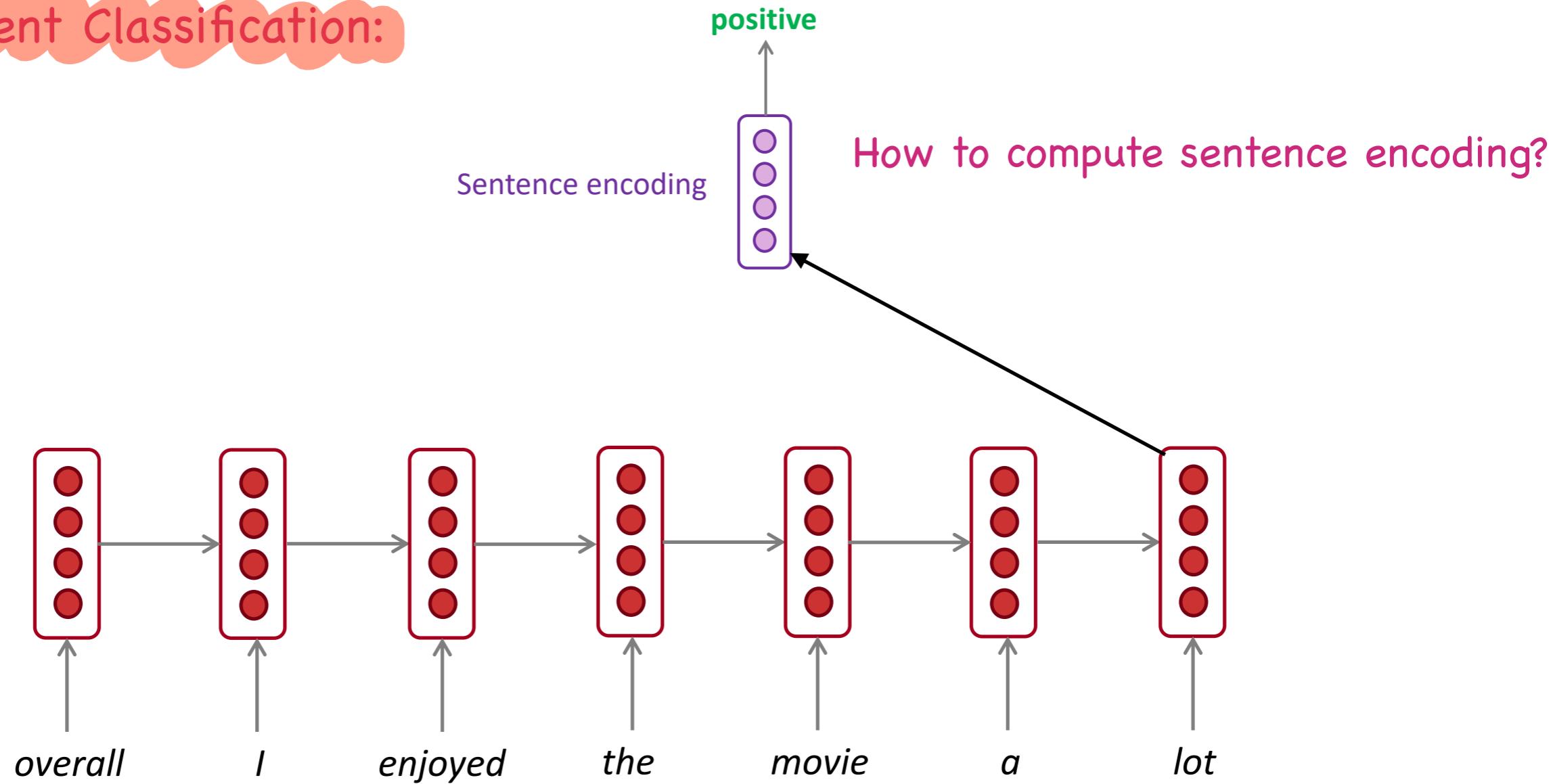
## ● Sentiment Classification:



# Sequence Classification with RNNs



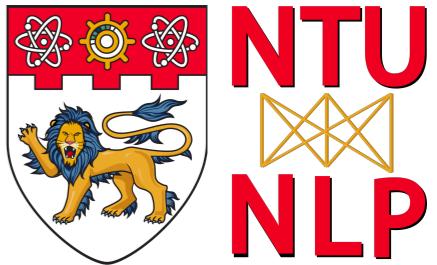
- Sentiment Classification:



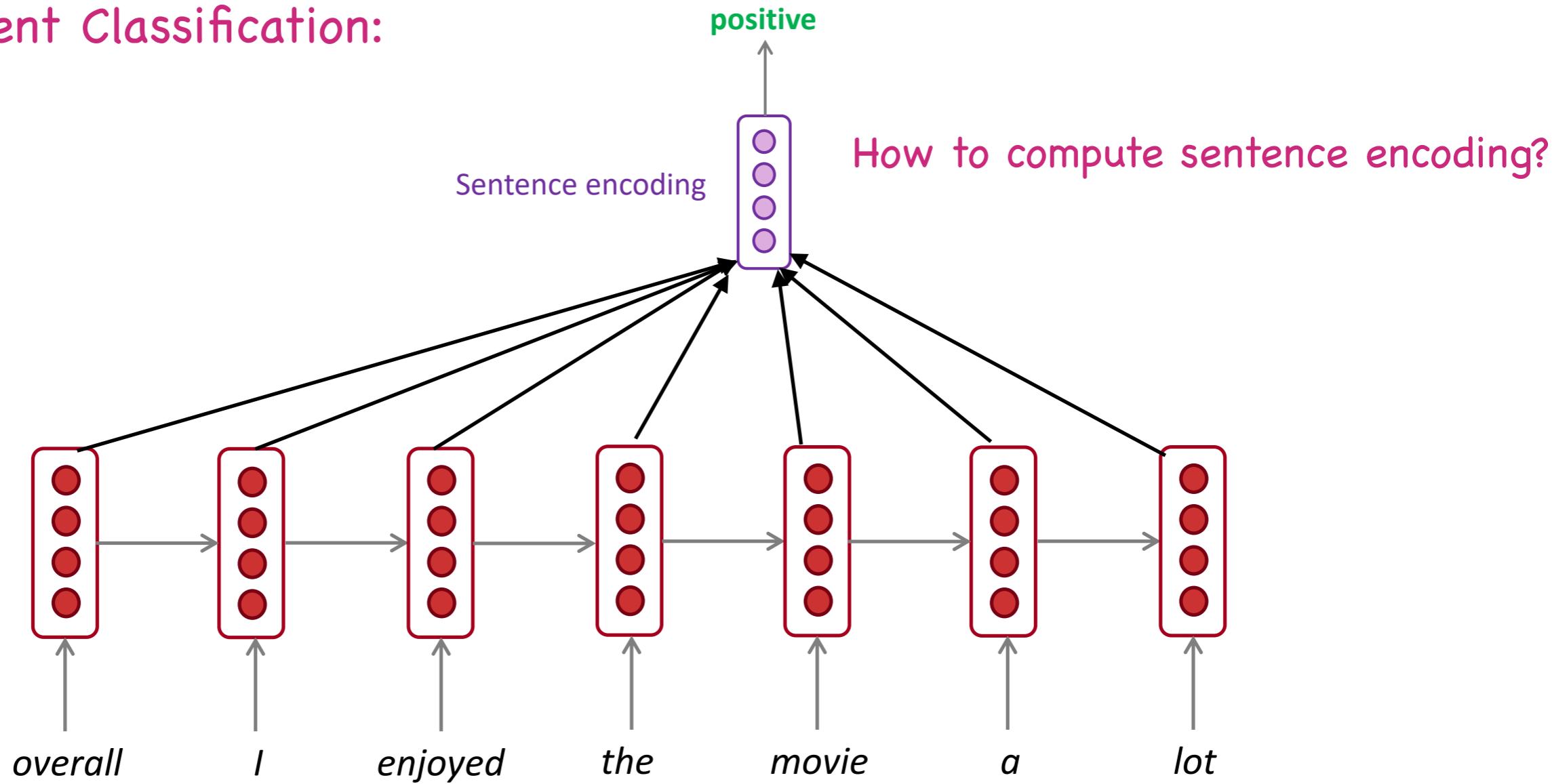
First choice :

Use the final hidden state as the summary of the sentence

# Sequence Classification with RNNs

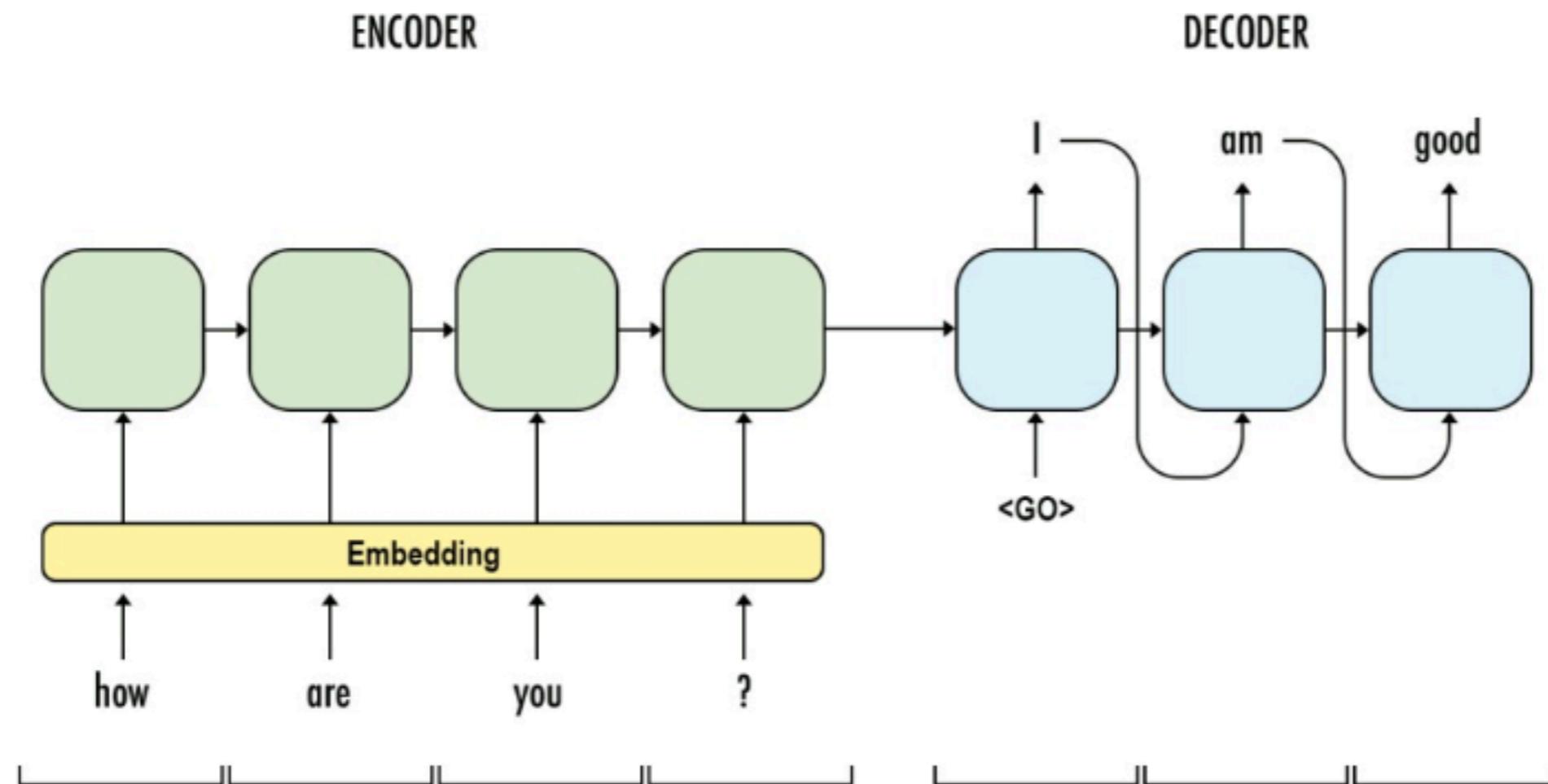
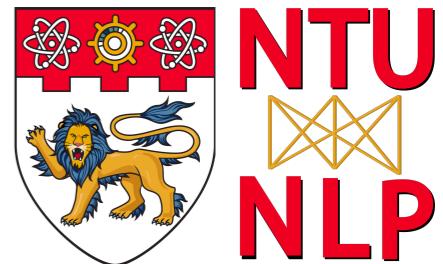


## ● Sentiment Classification:



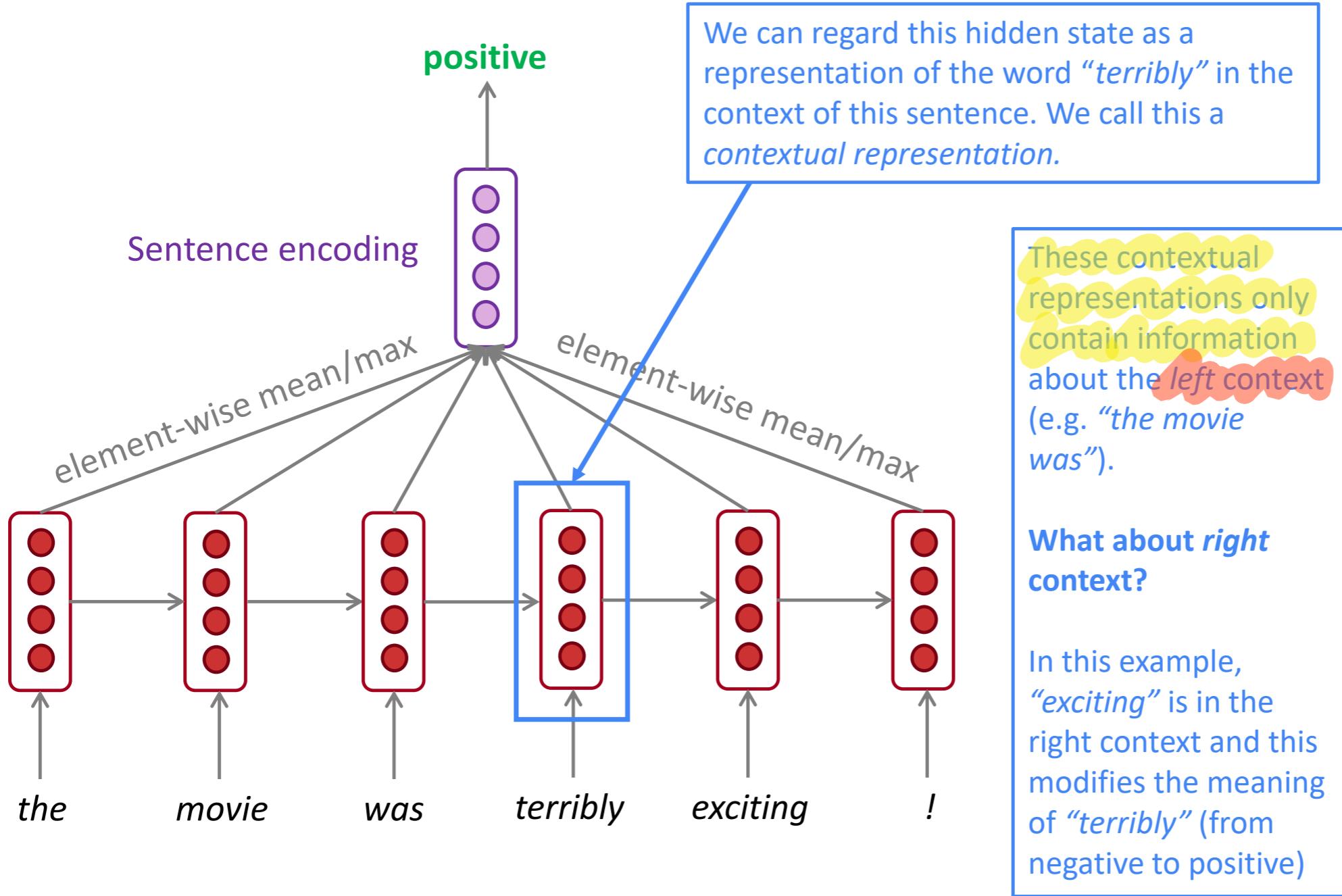
2<sup>nd</sup> choice = Use max/average pooling as the summary of the sentence

# RNNs as General Purpose Encoders/Decoders

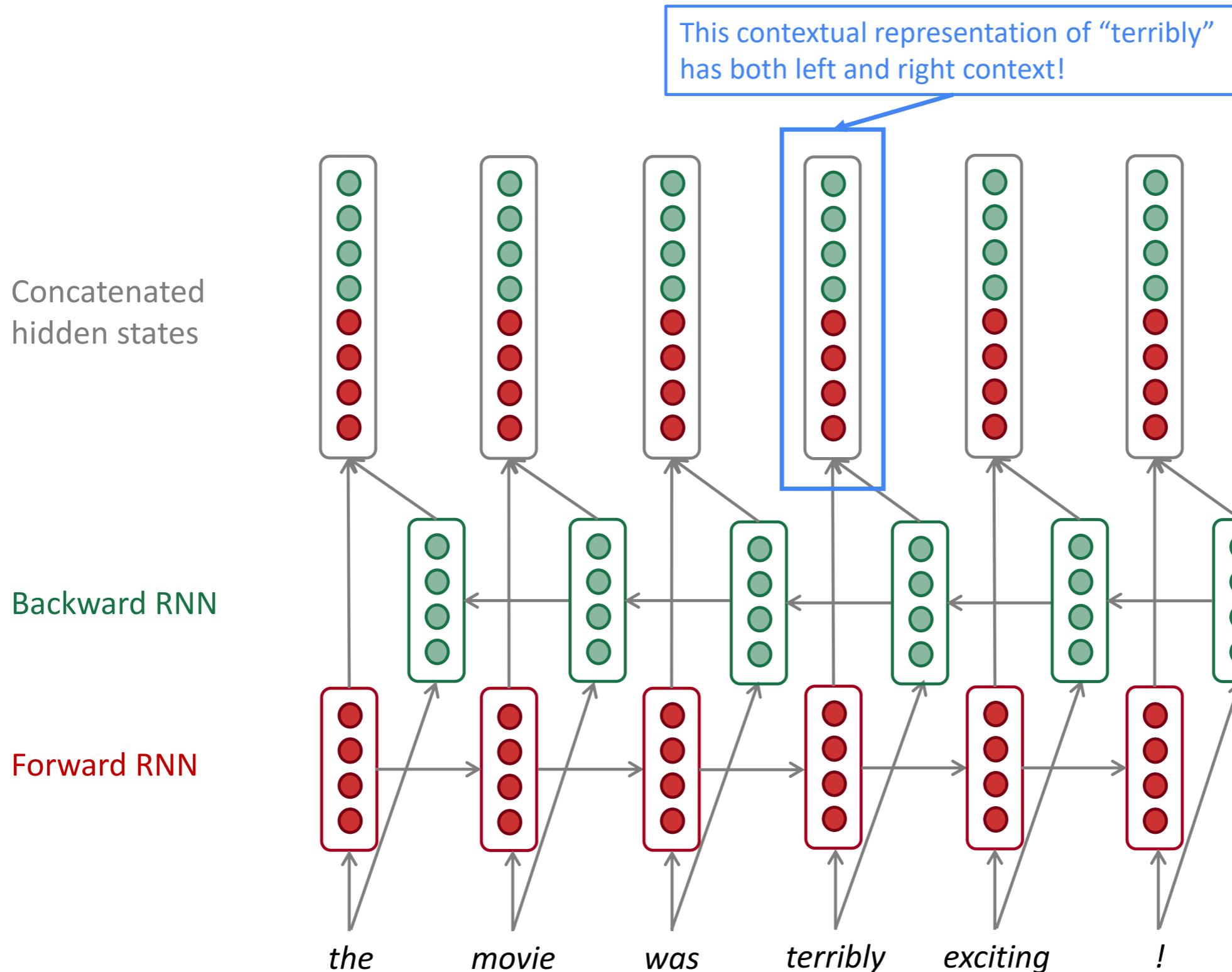


- We will see more details on this when we cover Seq2Seq model and Neural Machine Translation

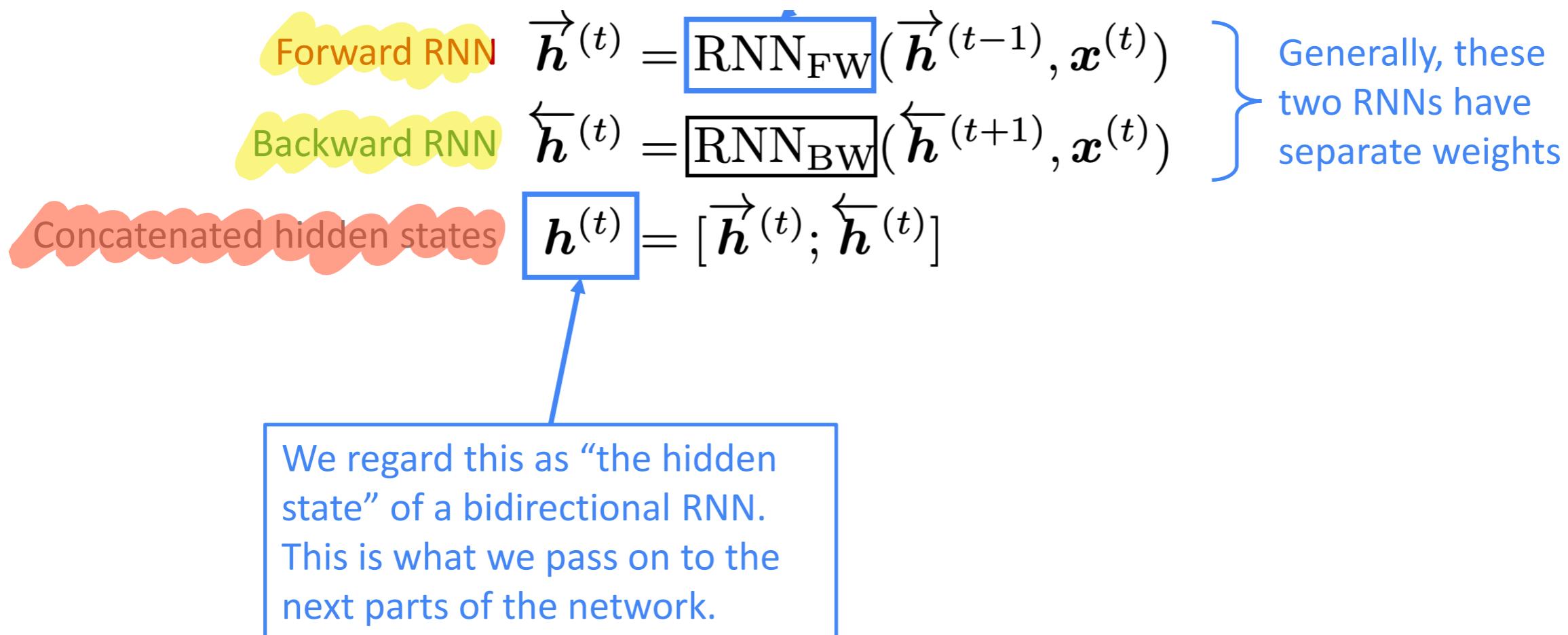
# Bi-directional RNN: Motivation



# Bi-directional RNN



# Bi-directional RNN



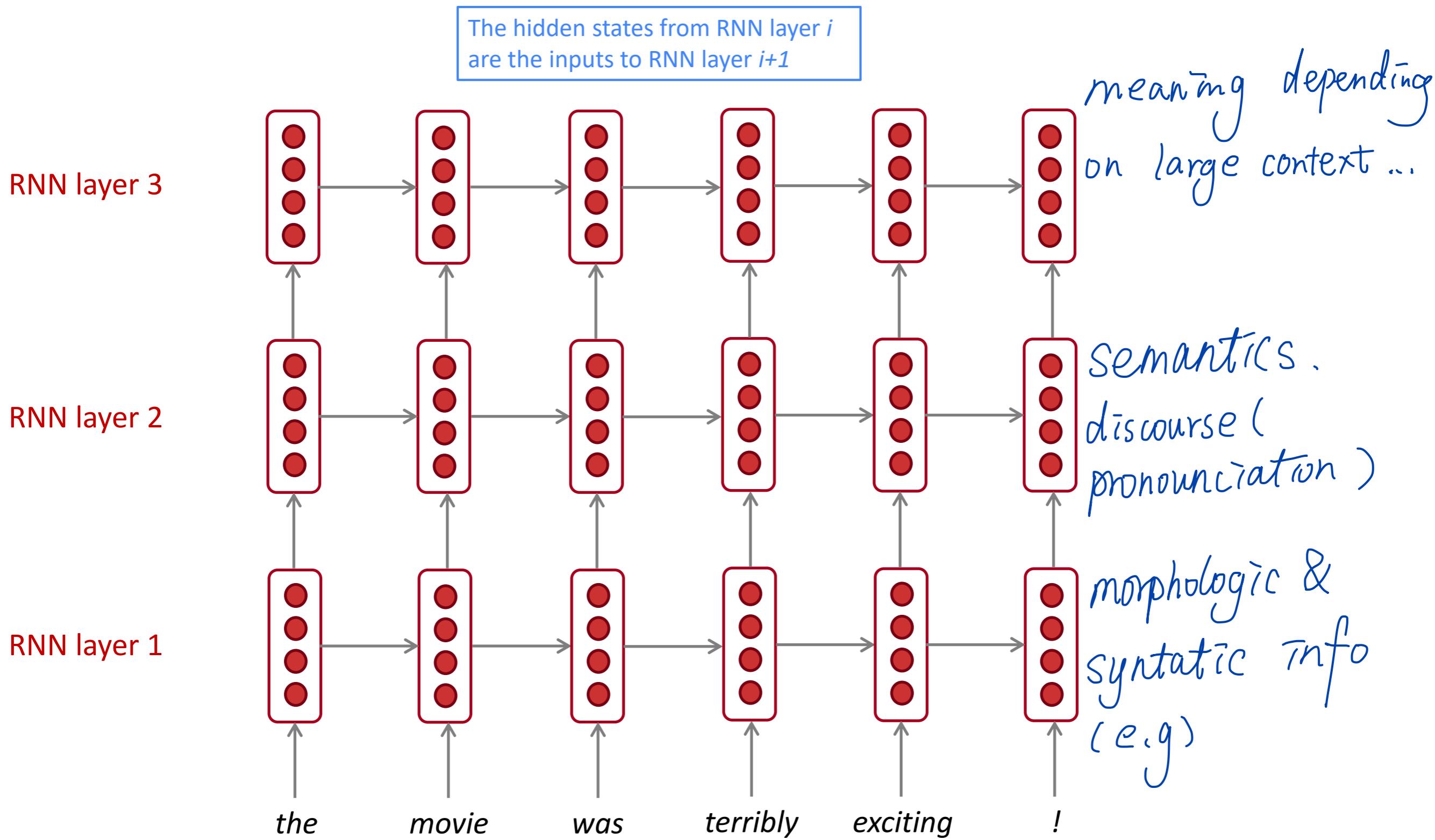
# Bi-directional RNN

- Bidirectional RNNs are only applicable when you have access to the entire input sequence. They are not applicable to Language Modeling as the future tokens are not accessible.
- If you do have entire input sequence, bidirectionality is powerful for encoding (you should use it by default).
- For example, BERT (Bidirectional Encoder Representations from Transformers) is a powerful pretrained contextual representation system built on bidirectionality (we'll learn more about BERT later)

# Stacked RNNs

- RNNs are already “deep” in one dimension (time)
- We can also make them “deep” in another dimension by putting an RNN on top of another - this is a multi-layer RNN.
- This allows the network to compute more complex representations
  - The lower RNNs should compute lower-level features and the higher RNNs should compute higher-level features.
- Multi-layer RNNs are also called stacked RNNs.

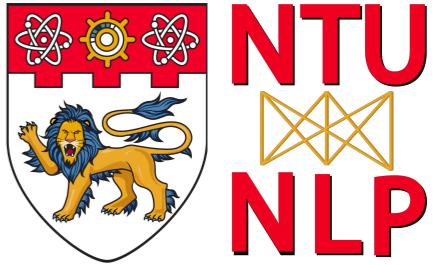
# Stacked RNNs



# Stacked RNNs

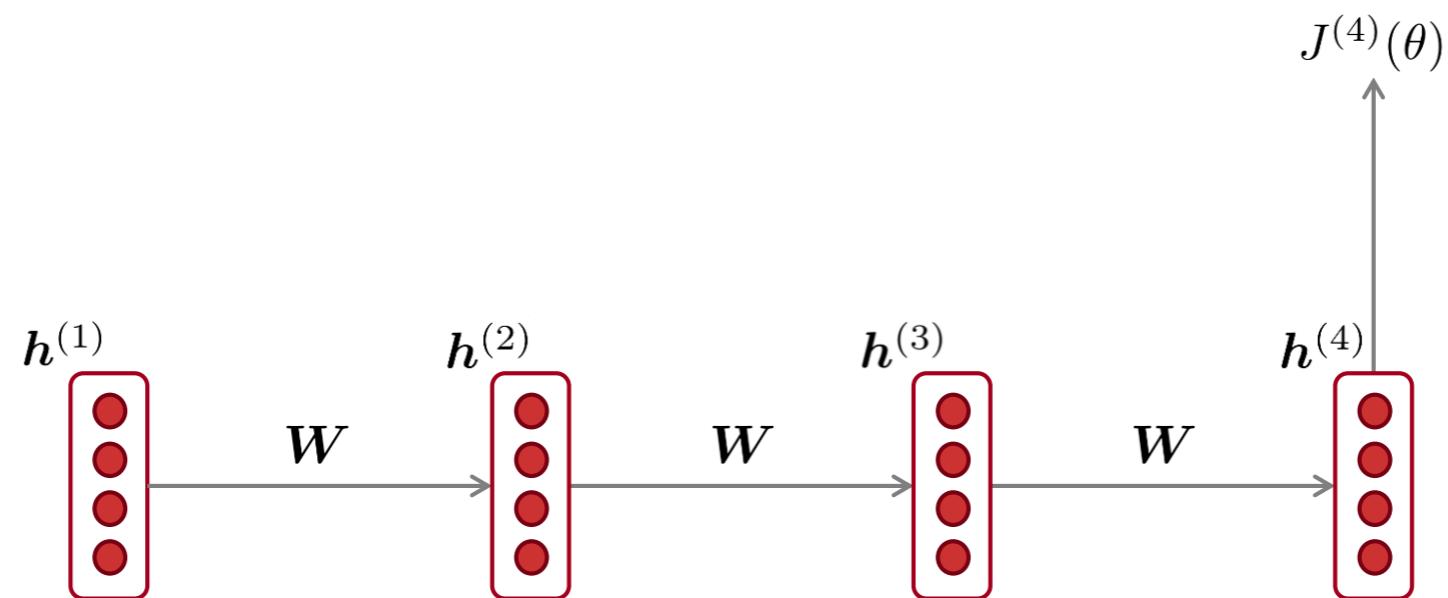
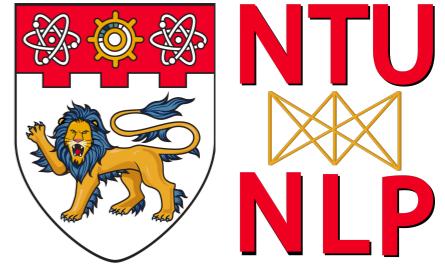
- High-performing RNNs are often multi-layer (but aren't as deep as convolutional or feed-forward networks) already deep in time dimension .
- For example: In a 2017 paper, Britz et al find that for Neural Machine Translation, 2 to 4 layers is best for the encoder RNN, and 4 layers is best for the decoder RNN
  - However, skip-connections/dense-connections are needed to train deeper RNNs (e.g. 8 layers)
  - Transformer-based networks (e.g. BERT) can be up to 24 layers

# Lecture Plan

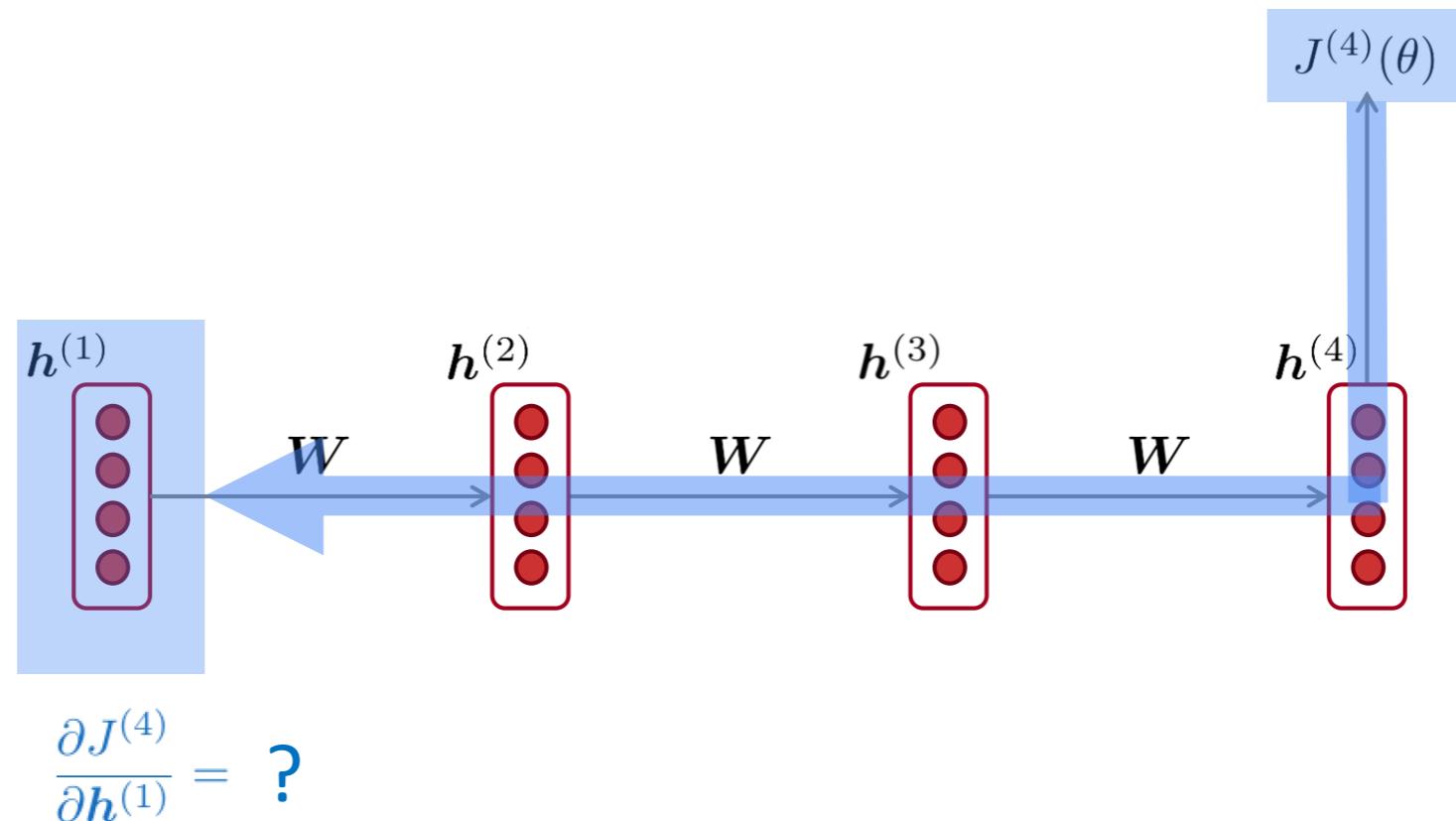
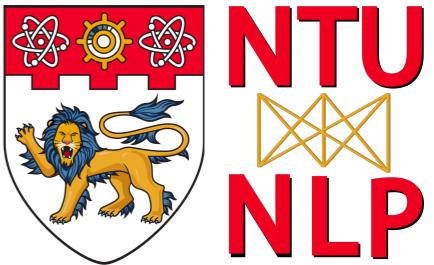


- Language Modeling
- Recurrent Neural Networks
- Backpropagation through time
- Text Generation with RNNs
- Sequence Tagging (NER, POS) with RNNs
- Sequence Classification with RNNs
- Bi-directional and stacked RNNs
- Gated RNNs (GRU, LSTM)

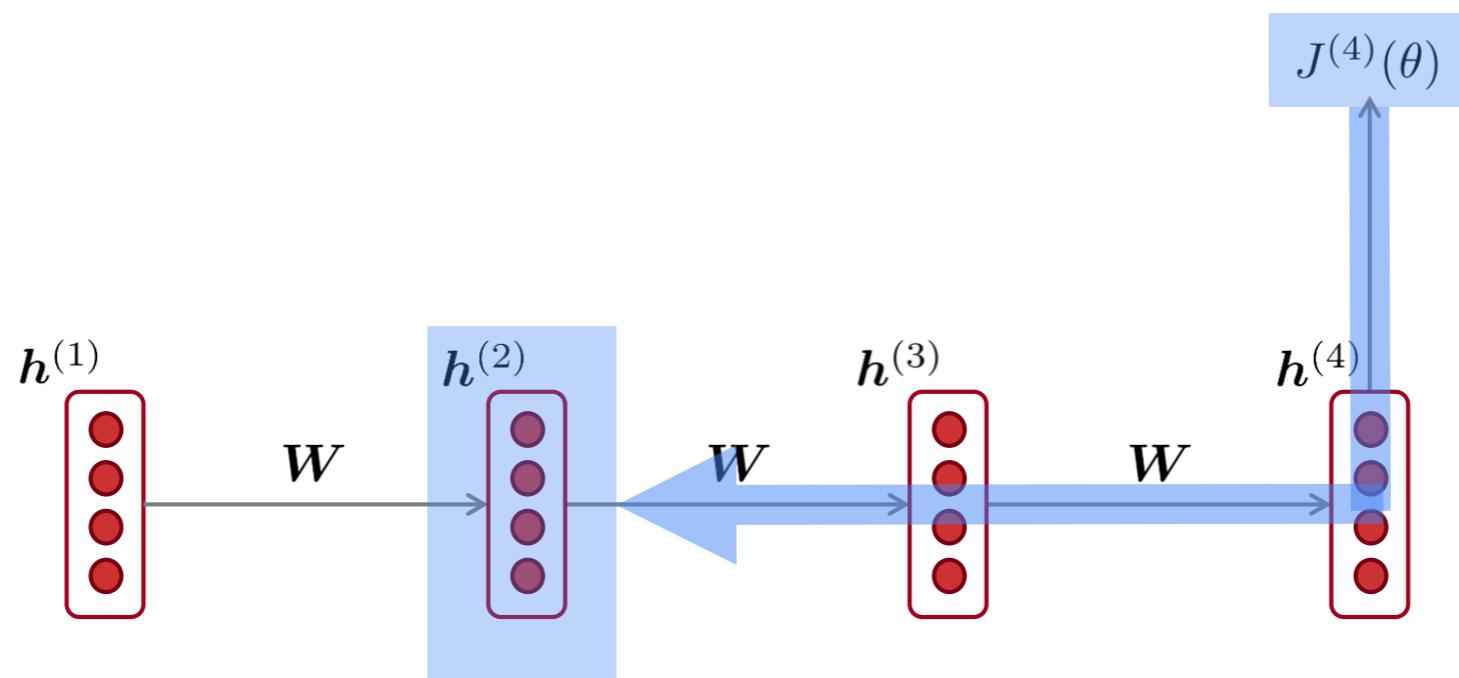
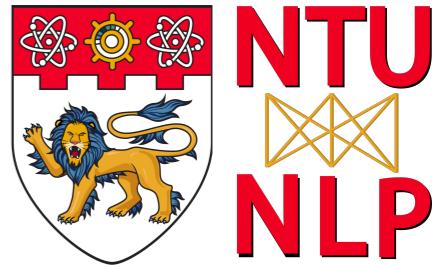
# Vanishing Gradient Problem with Vanilla RNNs



# Vanishing Gradient Problem with Vanilla RNNs



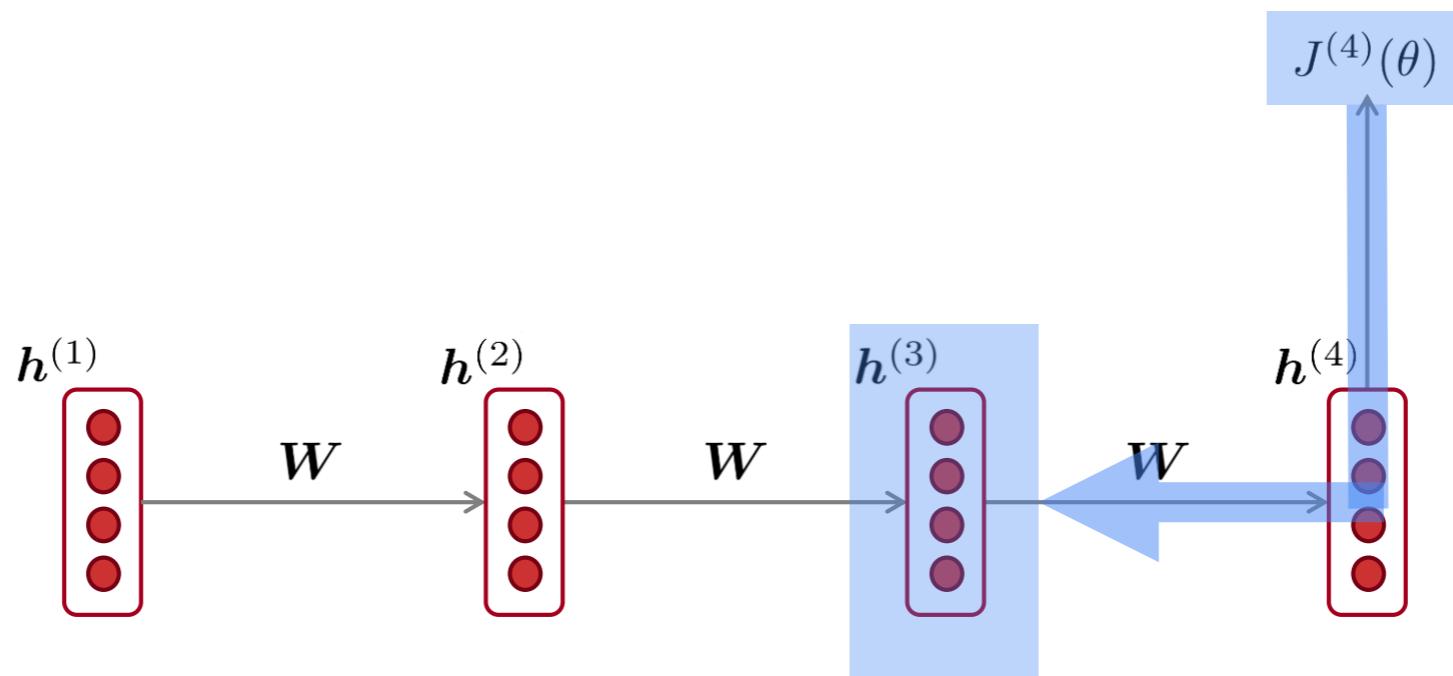
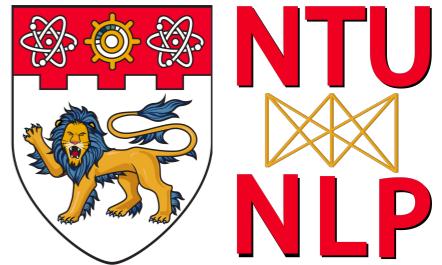
# Vanishing Gradient Problem with Vanilla RNNs



$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = \frac{\partial h^{(2)}}{\partial h^{(1)}} \times \frac{\partial J^{(4)}}{\partial h^{(2)}}$$

chain rule!

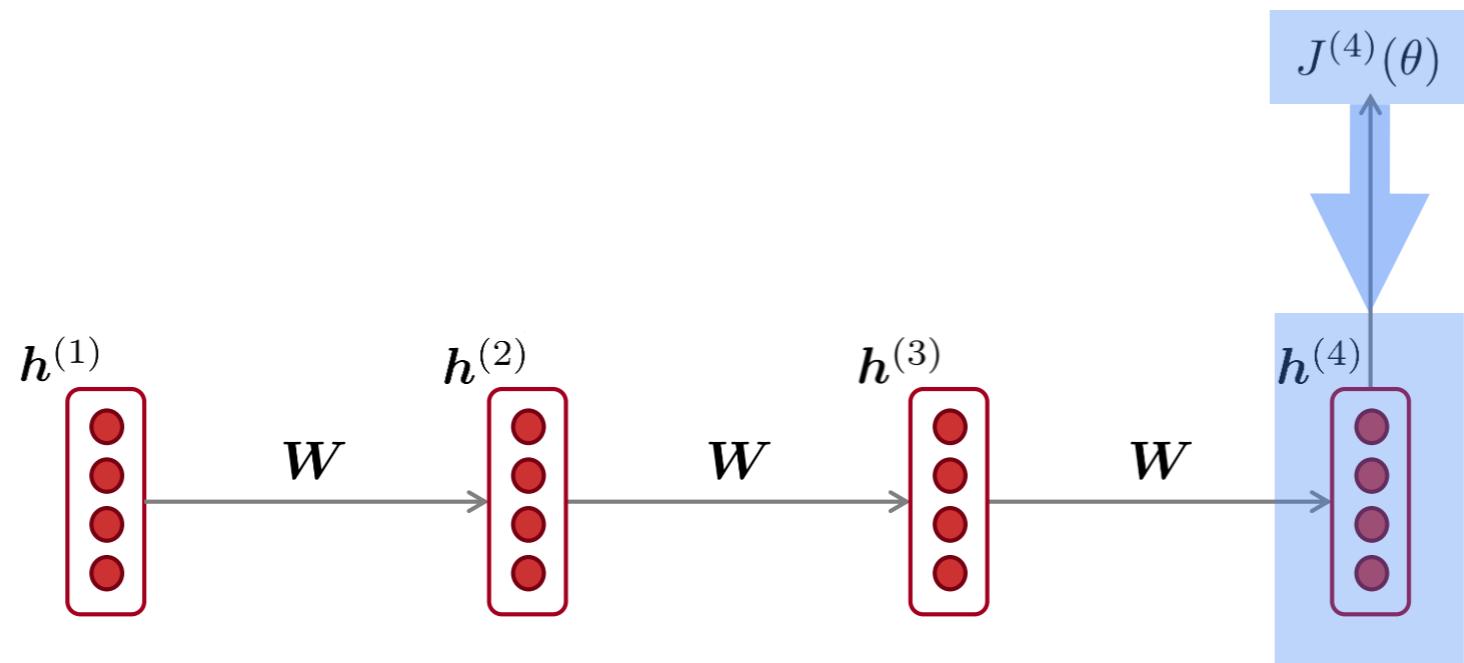
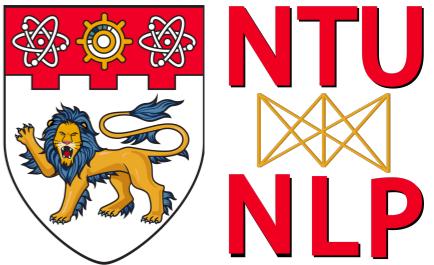
# Vanishing Gradient Problem with Vanilla RNNs



$$\frac{\partial J^{(4)}}{\partial \mathbf{h}^{(1)}} = \frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{h}^{(1)}} \times \frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{h}^{(2)}} \times \frac{\partial J^{(4)}}{\partial \mathbf{h}^{(3)}}$$

chain rule!

# Vanishing Gradient Problem with Vanilla RNNs



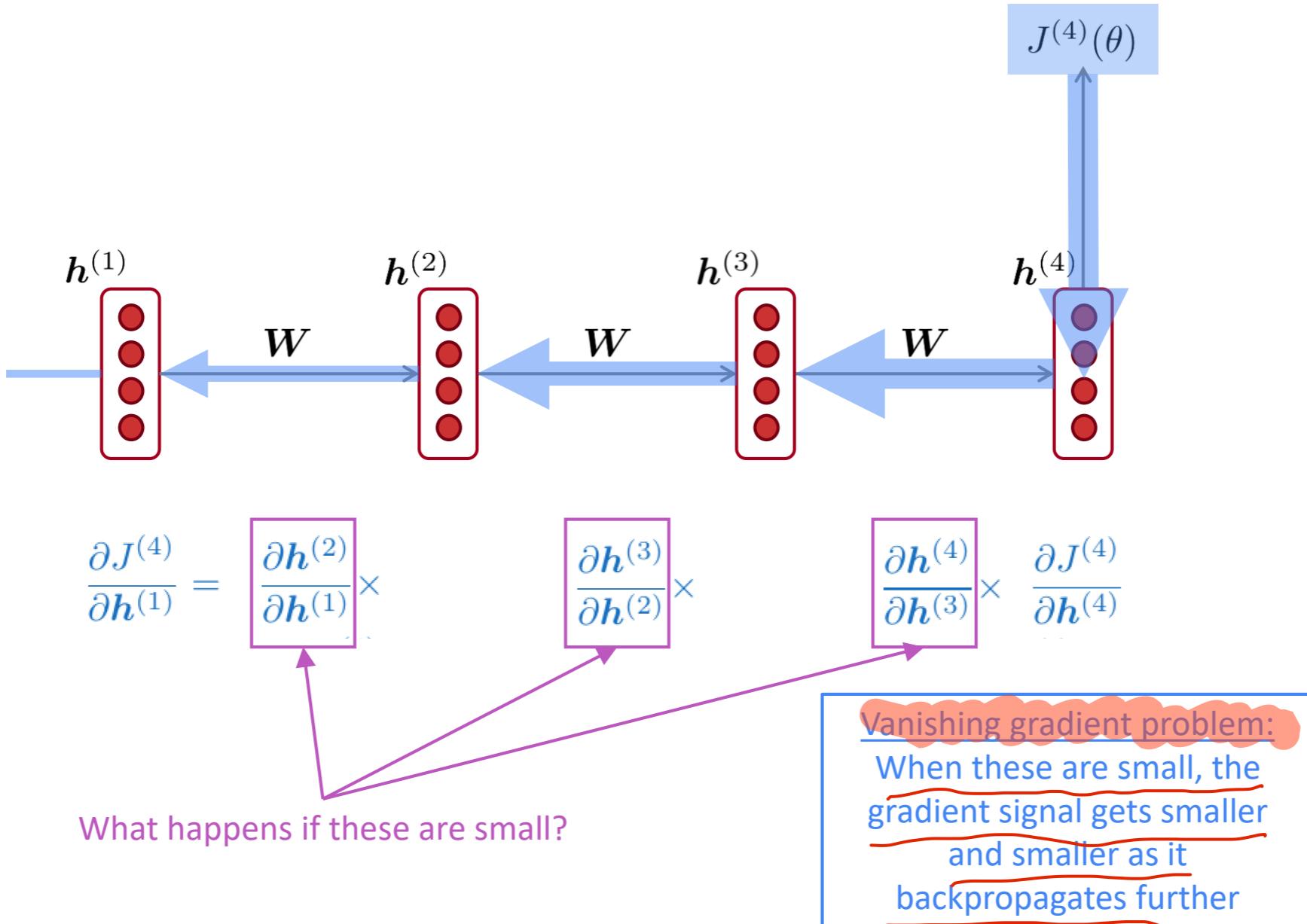
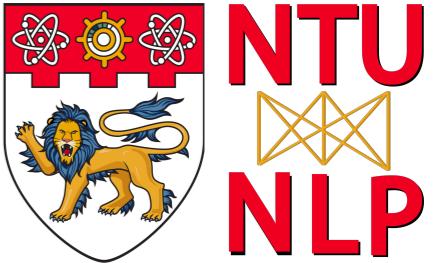
$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = \frac{\partial h^{(2)}}{\partial h^{(1)}} \times$$

$$\frac{\partial h^{(3)}}{\partial h^{(2)}} \times$$

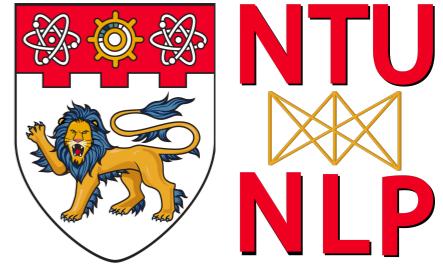
$$\frac{\partial h^{(4)}}{\partial h^{(3)}} \times \frac{\partial J^{(4)}}{\partial h^{(4)}}$$

chain rule!

# Vanishing Gradient Problem with Vanilla RNNs



# Vanishing Gradient Problem with Vanilla RNNs



- A hidden layer in RNN:  $\mathbf{h}^{(t)} = \sigma(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)} + \mathbf{b}_1)$
- Taking the derivative:  $\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} = \text{diag}\left(\sigma'\left(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)} + \mathbf{b}_1\right)\right) \mathbf{W}_h$
- The gradient of the loss on step  $i$  with respect to the hidden state on some previous step  $j$ :

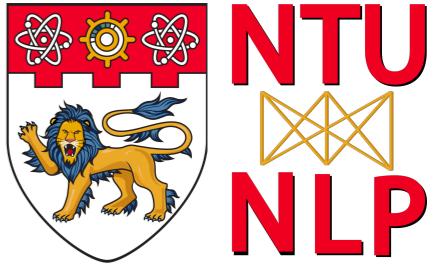
$$\frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(j)}} = \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \prod_{j < t \leq i} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} \quad (\text{chain rule})$$

$$= \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \boxed{\mathbf{W}_h^{(i-j)}} \prod_{j < t \leq i} \text{diag}\left(\sigma'\left(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)} + \mathbf{b}_1\right)\right) \quad (\text{value of } \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}})$$



If  $\mathbf{W}_h$  is small, then this term gets vanishingly small as  $i$  and  $j$  get further apart

# Vanishing Gradient Problem with Vanilla RNNs

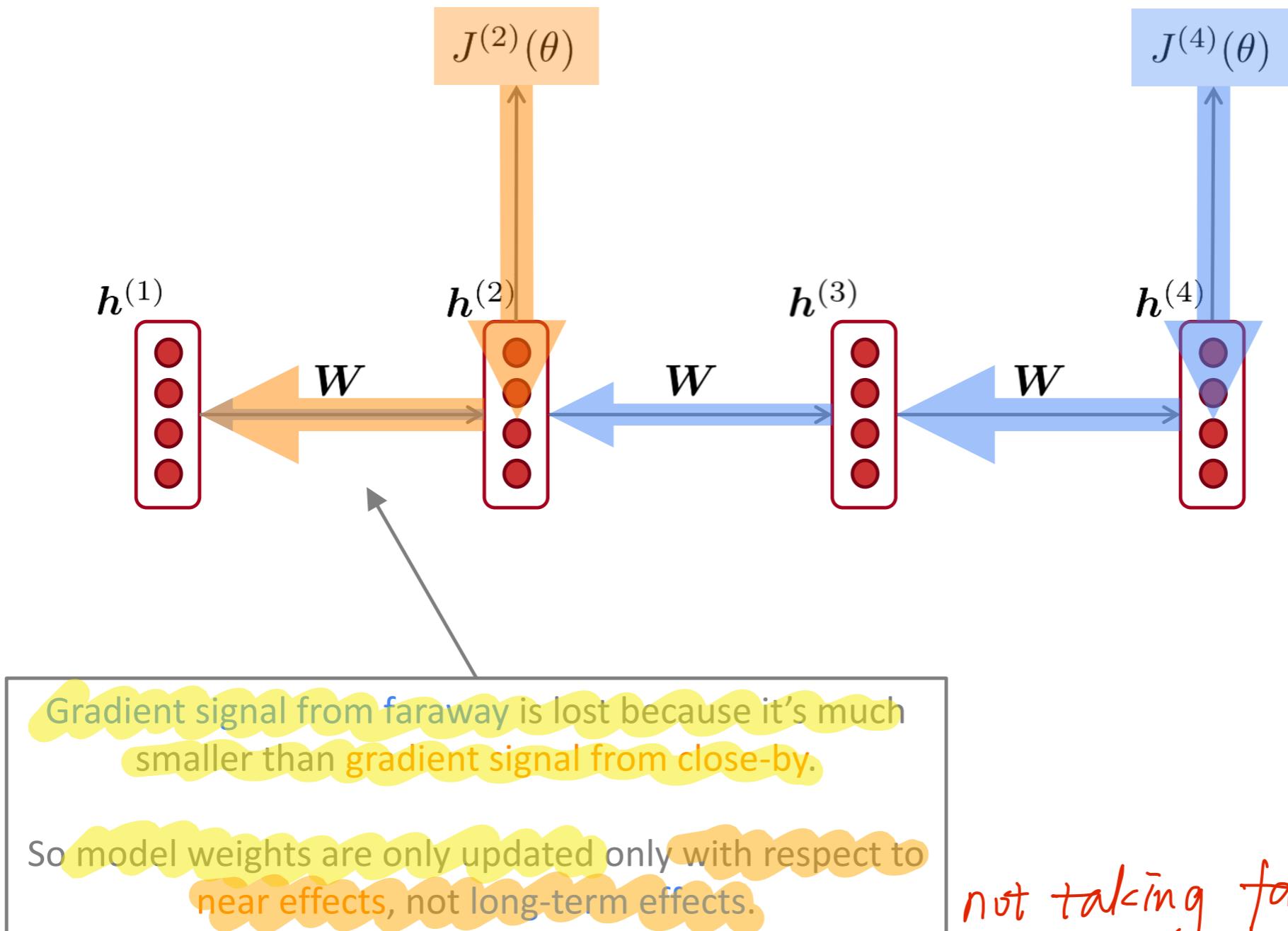


- With the Matrix L2 norm:

$$\left\| \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(j)}} \right\| \leq \left\| \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \right\| \|\mathbf{W}_h\|^{(i-j)} \prod_{j < t \leq i} \left\| \text{diag} \left( \sigma' \left( \mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)} + \mathbf{b}_1 \right) \right) \right\|$$

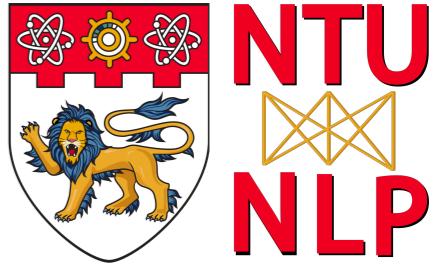
- Pascanu et al showed that if the largest eigenvalue of  $\mathbf{W}_h$  is less than 1 (bound is 1 for sigmoid nonlinearity), then the gradient  $\left\| \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(j)}} \right\|$  will shrink exponentially
- Similarly a largest eigenvalue  $>1$  results in exploding gradients

# Why Vanishing Gradient is a Problem



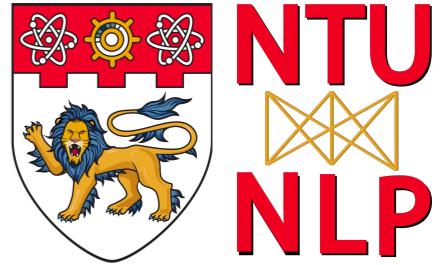
not taking far weights into account. fails to capture the defined answers. (supervised)

# Why Vanishing Gradient is a Problem



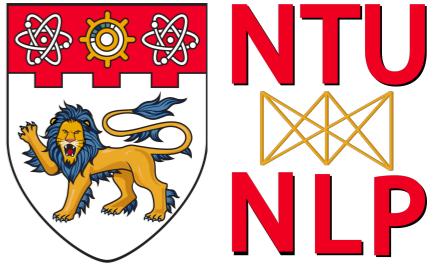
- Gradient can be viewed as a measure of the effect of the past on the future.
- If the gradient becomes vanishingly small over longer distances (step  $t$  to step  $t+n$ ), then we can't tell whether:
  - There's **no dependency** between step  $t$  and  $t+n$  in the data or
  - We have **wrong parameters** to capture the true dependency between  $t$  and  $t+n$

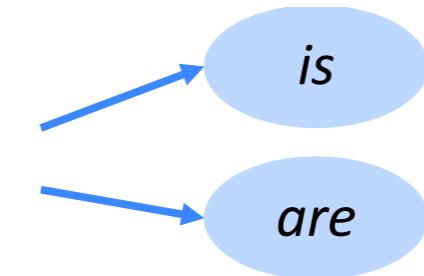
# Why Vanishing Gradient is a Problem



- **LM task:** When she tried to print her tickets, she found that the printer was out of toner. She went to the stationery store to buy more toner. It was very overpriced. After installing the toner into the printer, she finally printed her \_\_\_\_\_
- To learn from this training example, the RNN-LM needs to **model the dependency** between “tickets” on the 7<sup>th</sup> step and the target word “tickets” at the end.
- But if gradient is small, the model **can't learn this dependency**

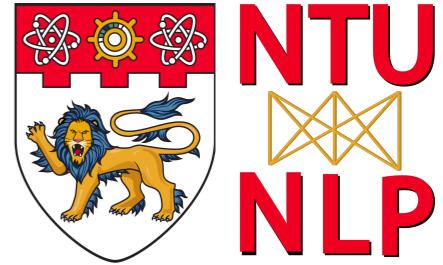
# Why Vanishing Gradient is a Problem



- LM task: *The writer of the books*    
- Correct answer: *The writer of the books is planning a sequel*
- Syntactic recency: *The writer of the books is* (correct)
- Sequential recency: *The writer of the books are* (incorrect)

Due to vanishing gradient, RNN-LMs are better at learning from sequential recency than syntactic recency, so they make this type of error more often [Linzen et al 2016]

# Why Exploding Gradient is a Problem



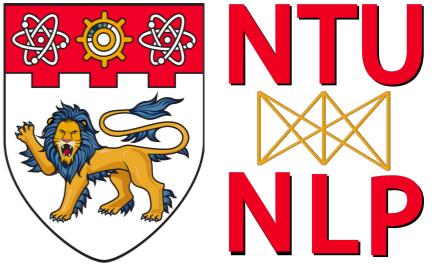
- If the gradient becomes too big, then the SGD update step becomes too big:

$$\theta^{new} = \theta^{old} - \underbrace{\alpha \nabla_{\theta} J(\theta)}_{\text{gradient}}$$

learning rate

- This can cause **bad updates**: we take too large a step and reach a bad parameter configuration (with large loss)
- In the worst case, this will result in **Inf** or **Nan** in your network (then you have to restart training from an earlier checkpoint)

# Solution for Exploding Gradient



Gradient clipping: if the norm of the gradient is greater than some threshold, scale it down before applying SGD update

---

## Algorithm 1 Pseudo-code for norm clipping

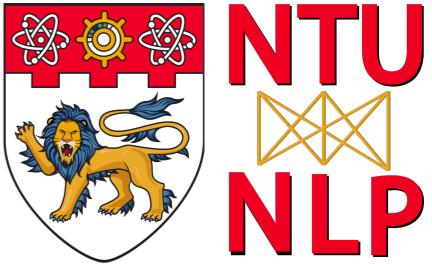
---

```
hat{g} ← ∂E / ∂θ
if ||hat{g}|| ≥ threshold then
    hat{g} ← threshold / ||hat{g}|| * hat{g}
end if
```

---

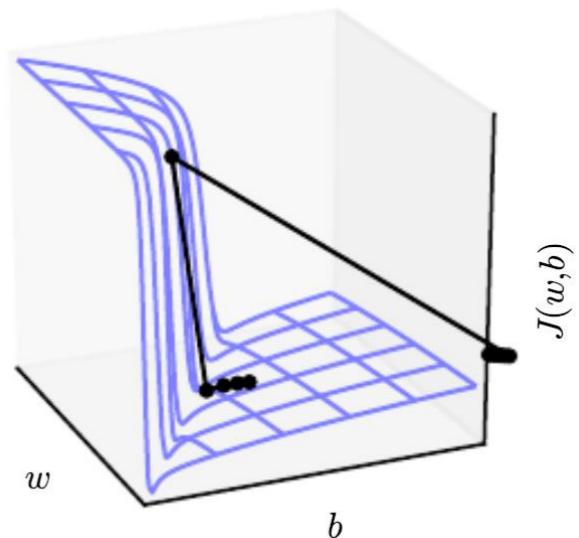
- Intuition: take a step in the same direction, but a smaller step

# Solution for Exploding Gradient

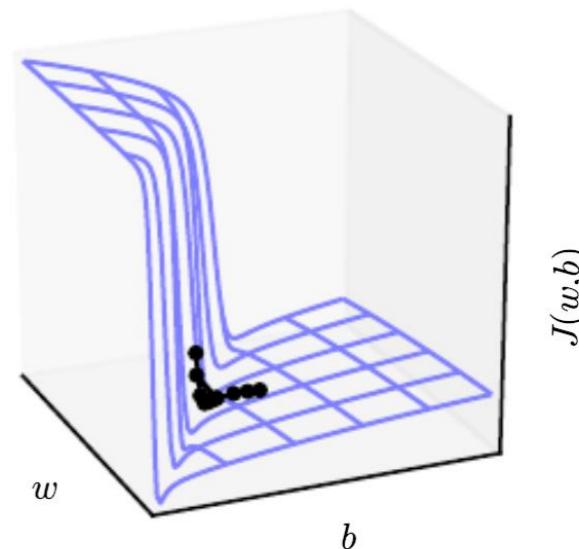


- This shows the loss surface of a simple RNN (hidden state is a scalar not a vector)
- The “**cliff**” is dangerous because it has steep gradient
- In the Fig. at the top, gradient descent takes **two very big steps** due to steep gradient, resulting in climbing the cliff then shooting off to the right (both **bad updates**)
- At the bottom, gradient clipping reduces the size of those steps, so effect is **less drastic**

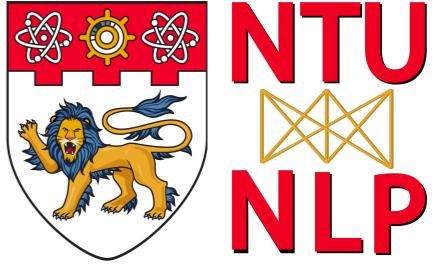
Without clipping



With clipping



# Tackling Vanishing Gradient

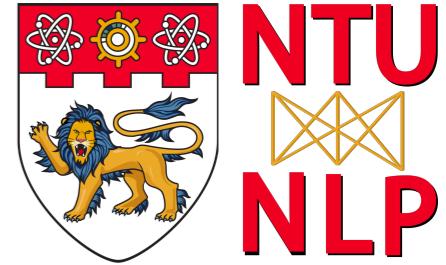


- The main problem is that it's too difficult for the vanilla RNNs to learn to preserve information over many timesteps. The "same" hidden state is constantly updated

$$\mathbf{h}^{(t)} = \sigma \left( \mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)} + \mathbf{b} \right)$$

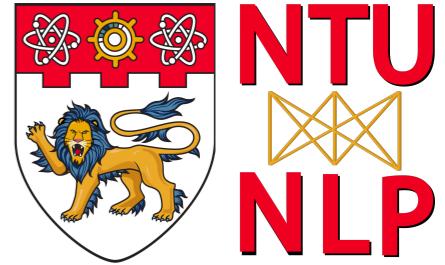
- How about a RNN with separate memory?

# Long Short-Term Memory (LSTM)



- On step  $t$ , there is a hidden state  $h^{(t)}$  and a cell state  $c^{(t)}$ 
  - Both are vectors length  $n$  neurons in  $h$
  - The cell stores long-term information
  - The LSTM can erase, write and read information from the cell
- The selection of which information is erased/written/read is controlled by three corresponding gates
  - The gates are also vectors length  $n$
  - On each timestep, each element of the gates can be open (1), closed (0), or somewhere in-between.
  - The gates are dynamic: their value is computed based on the current context previous text and current input

# Long Short-Term Memory (LSTM)



- Given current input and previous hidden state, we compute the gates

**Forget gate:** controls what is kept vs forgotten, from previous cell state

**Input gate:** controls what parts of the new cell content are written to cell

**Output gate:** controls what parts of cell are output to hidden state

**Sigmoid function:** all gate values are between 0 and 1

$$f^{(t)} = \sigma(W_f h^{(t-1)} + U_f x^{(t)} + b_f)$$
$$i^{(t)} = \sigma(W_i h^{(t-1)} + U_i x^{(t)} + b_i)$$
$$o^{(t)} = \sigma(W_o h^{(t-1)} + U_o x^{(t)} + b_o)$$

- Now we compute how much to forget, update and output

**New cell content:** this is the new content to be written to the cell

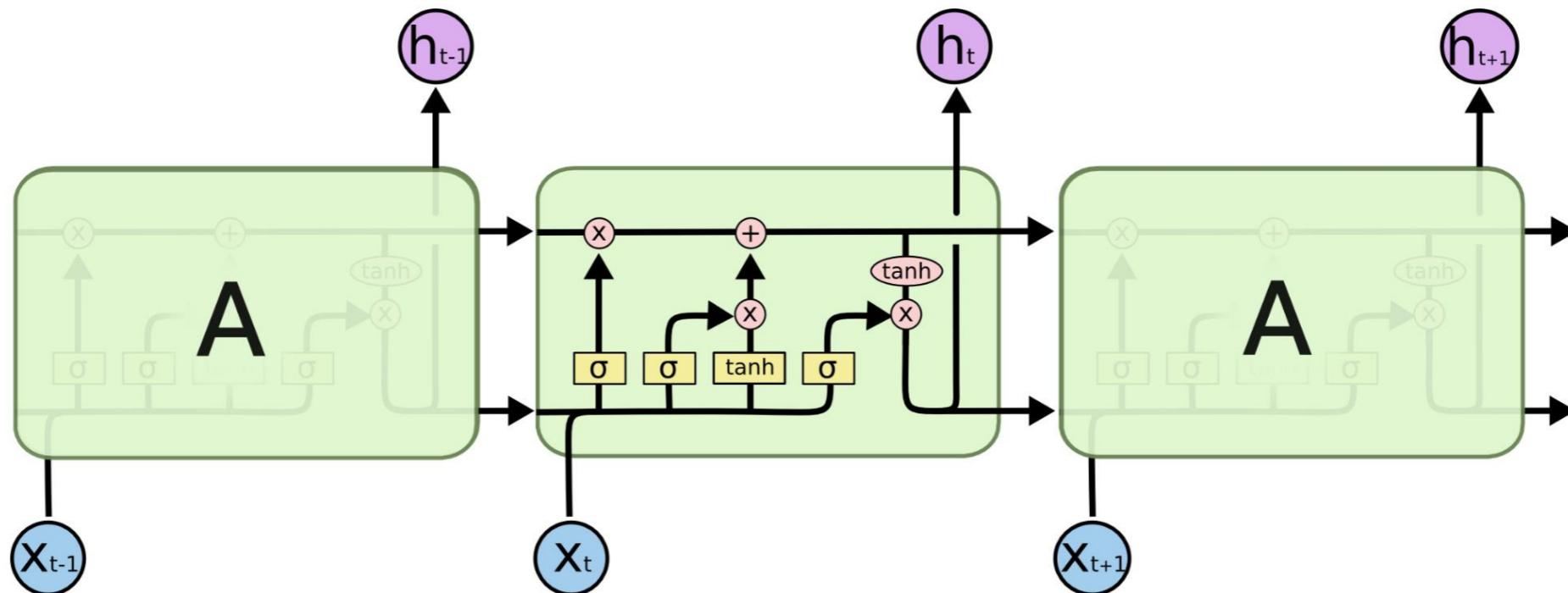
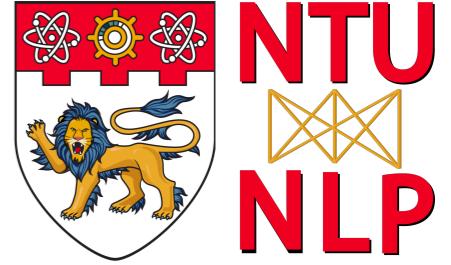
**Cell state:** erase ("forget") some content from last cell state, and write ("input") some new cell content

**Hidden state:** read ("output") some content from the cell

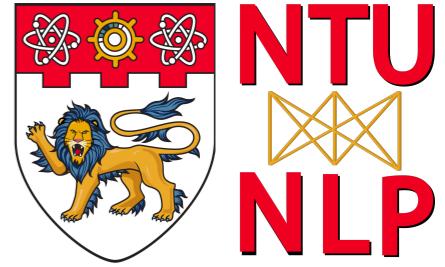
$$\tilde{c}^{(t)} = \tanh(W_c h^{(t-1)} + U_c x^{(t)} + b_c)$$
$$c^{(t)} = f^{(t)} \circ c^{(t-1)} + i^{(t)} \circ \tilde{c}^{(t)}$$
$$h^{(t)} = o^{(t)} \circ \tanh c^{(t)}$$

Gates are applied using element-wise product

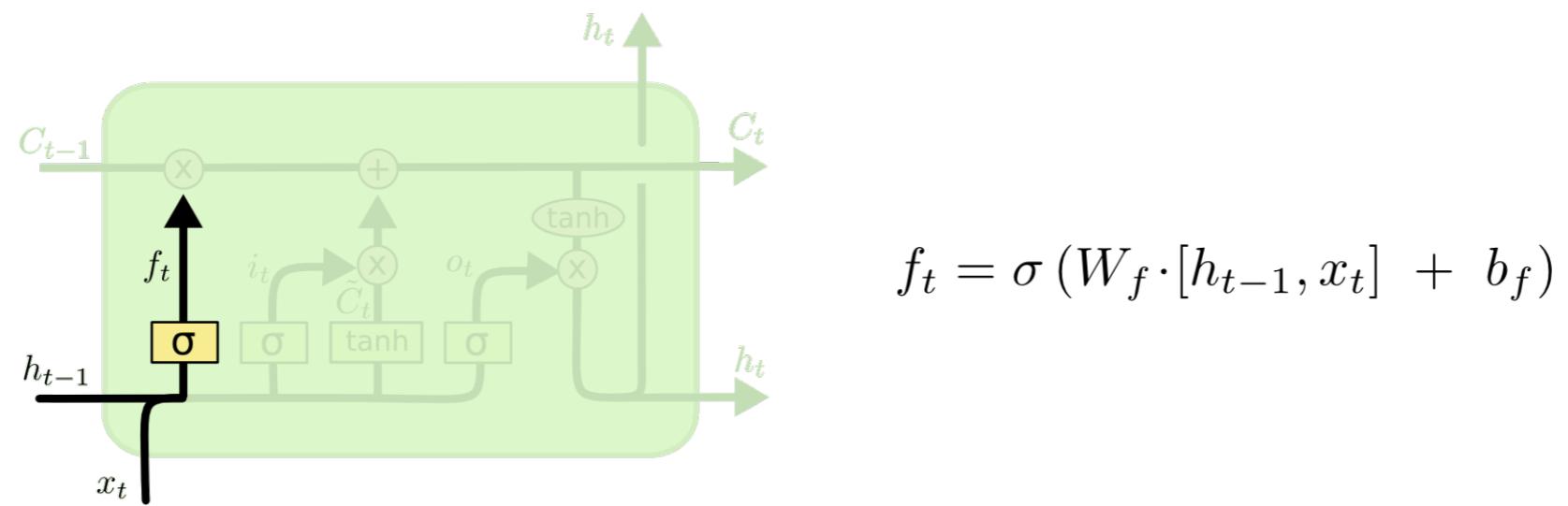
# Long Short-Term Memory (LSTM)



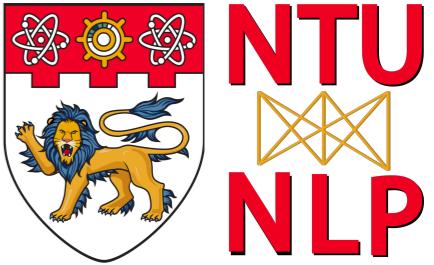
# Long Short-Term Memory (LSTM)



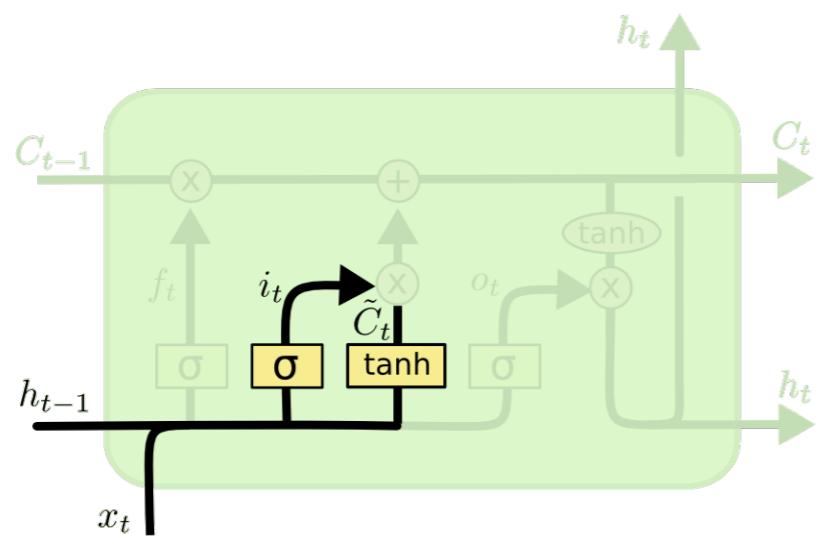
- Compute forget gate



# Long Short-Term Memory (LSTM)



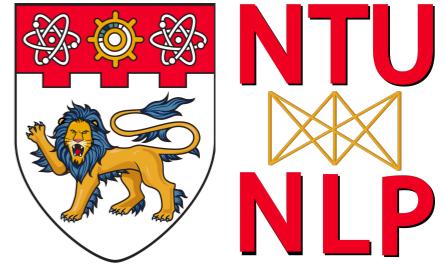
- Compute current cell state and input gate



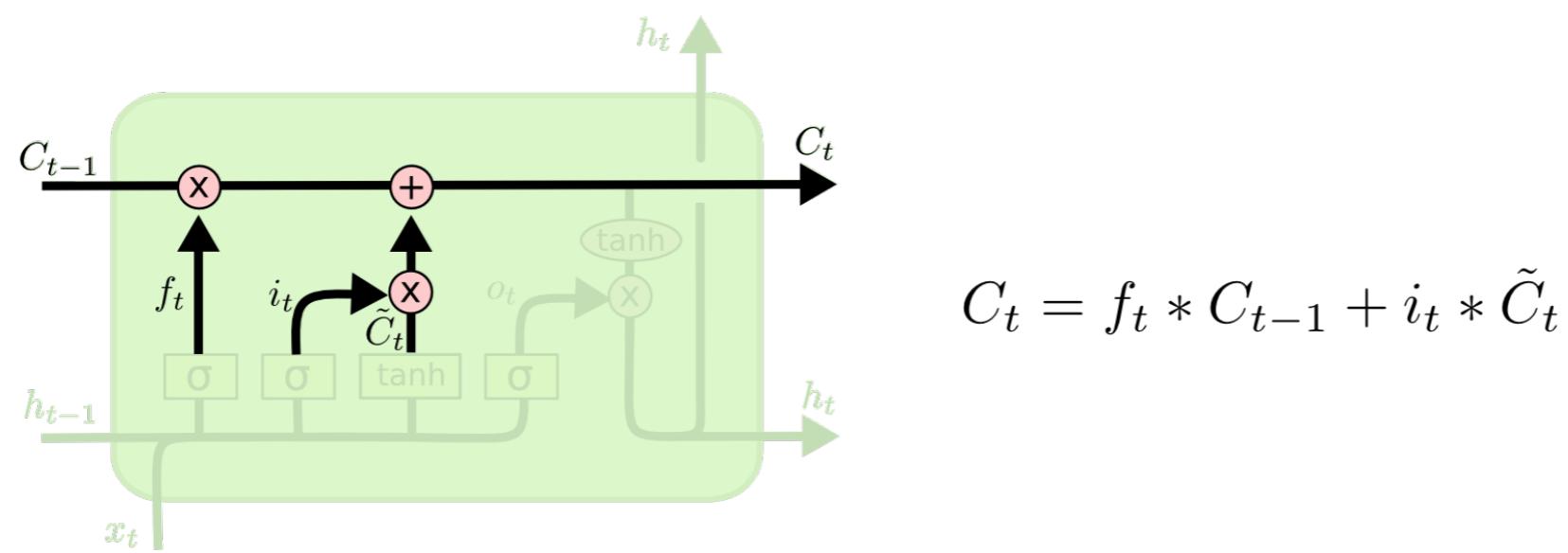
$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

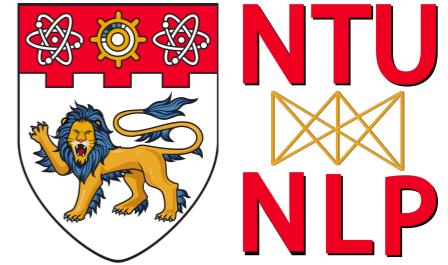
# Long Short-Term Memory (LSTM)



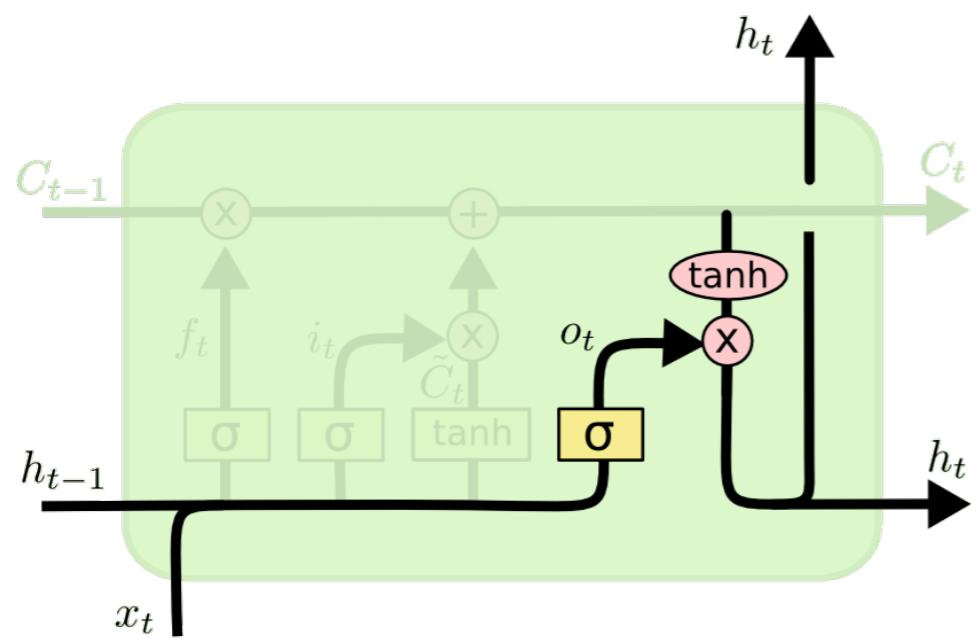
- Forget (erase) something and write (update) something new



# Long Short-Term Memory (LSTM)



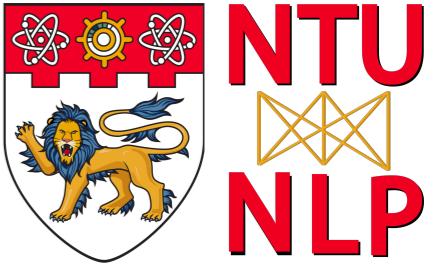
- Compute output gate and output state



$$o_t = \sigma (W_o [ h_{t-1}, x_t ] + b_o)$$

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

# How LSTM Remembers



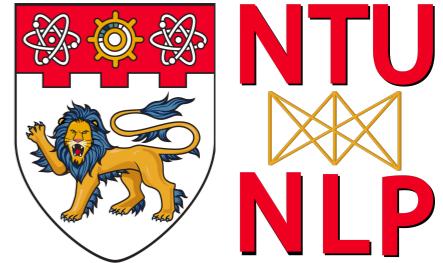
- The LSTM architecture makes it easier for the RNN to preserve information over many time steps.
  - If the forget gate is set to remember everything on every time step, then the info in the cell is preserved indefinitely.
  - In practice, LSTM uses adaptive gates, which are also learned (i.e., using their respective parameters)
- LSTM however doesn't guarantee that there is no vanishing/exploding gradient, but it does provide an easier way for the model to learn long-distance dependencies.

# Success of LSTMs

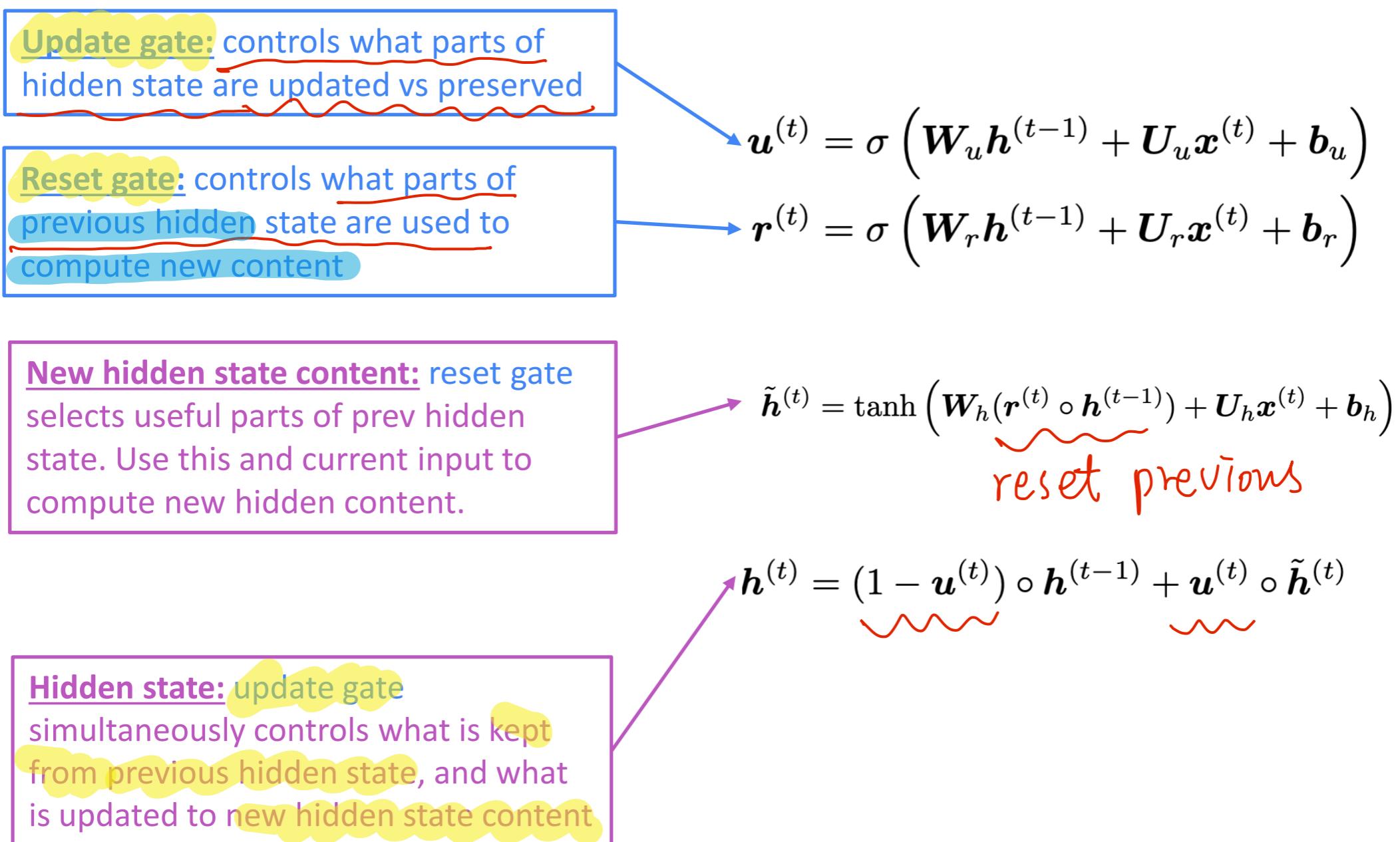
- In 2013-2015, LSTMs started achieving state-of-the-art results
  - Successful tasks: handwriting recognition, speech recognition, machine translation, parsing, image captioning.
  - Became the dominant approach
- In recent years (2017-2020), other approaches (e.g. Transformers) have become more dominant for certain tasks. For example in WMT (a MT conference + competition):
  - WMT 2016, the summary report contains “RNN” 44 times
  - In WMT 2018, the report contains “RNN” 9 times and “Transformer” 63 times

# Gated Recurrent Unit (GRU)

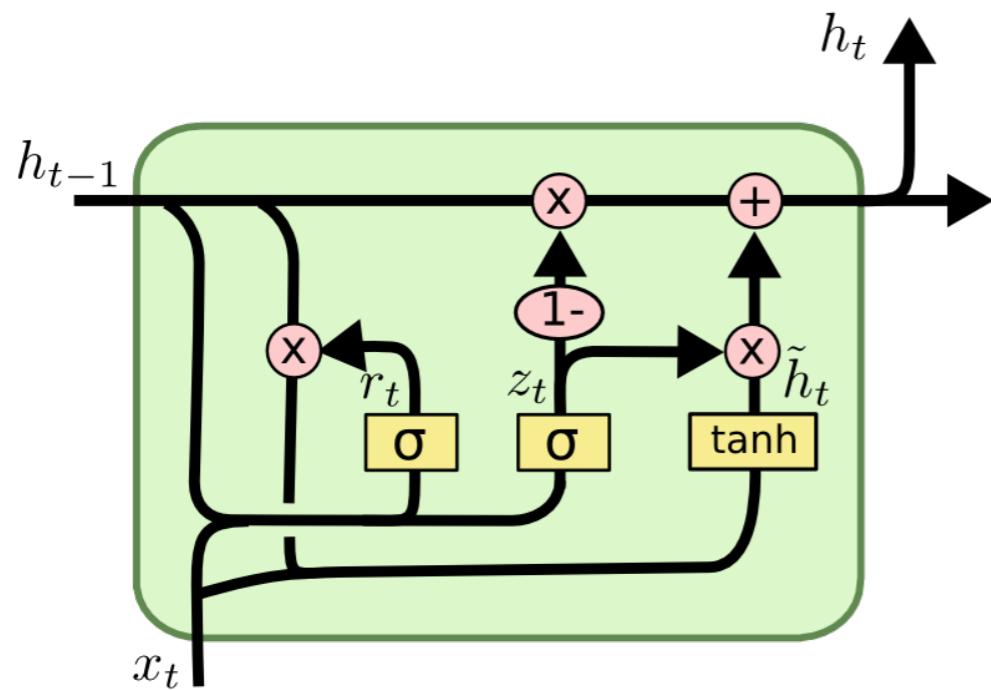
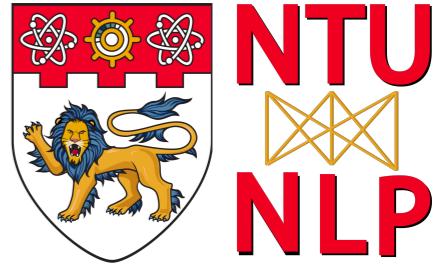
A simpler alternative to LSTM



- No cell state, remembers through two gates



# Gated Recurrent Unit (GRU)



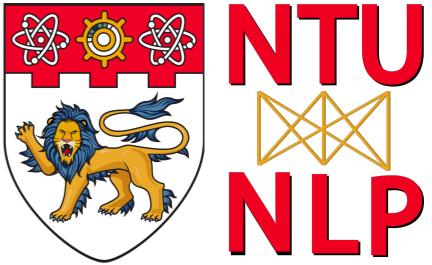
$$z_t = \sigma (W_z \cdot [h_{t-1}, x_t])$$

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

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

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

# LSTM vs. GRU

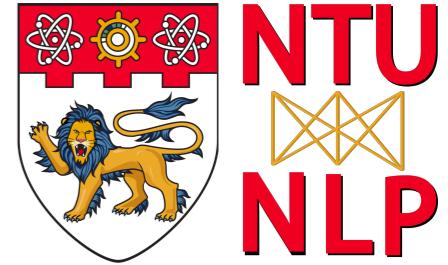


- GRU is quicker to compute and has fewer parameters
- No conclusive evidence that one consistently performs better than the other
- LSTM is a good default choice (especially if your data has particularly long dependencies, or you have lots of training data)
- Rule of thumb: start with LSTM, but switch to GRU if you want something more efficient

# Vanishing/Exploding Gradient in Other Architectures

- Vanishing/Exploding gradient problems are **not specific** to only RNNs
- Can happen in all neural **deep** architectures (including **feed-forward** and **convolutional**; RNNs are **deep in time**).
- Due to **chain rule** and/or choice of **nonlinearity function**, **gradient can become vanishingly small as it backpropagates**.
- Thus lower layers are learnt very slowly (hard to train)
- Solution: lots of new deep feedforward/convolutional architectures **add more direct connections** (thus allowing the gradient to flow)

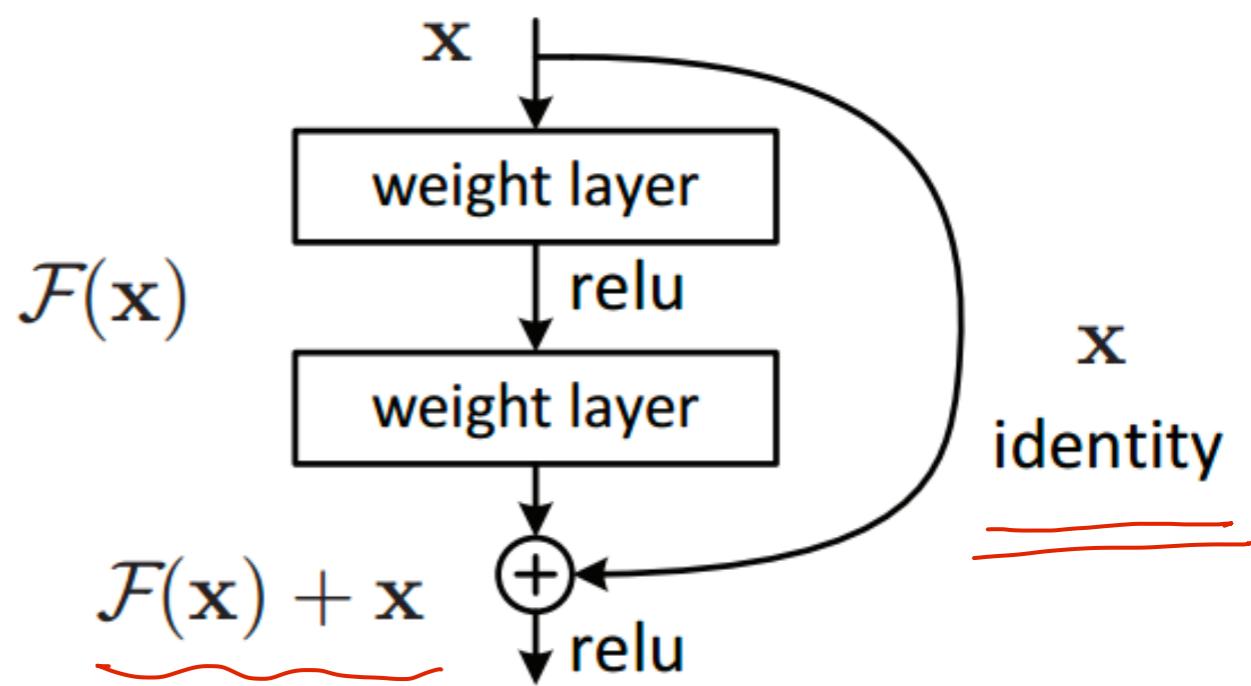
# Direct Connections: Residual Links



- Solution: lots of new deep feedforward/convolutional architectures **add more direct connections** (thus allowing the gradient to flow)

- Residual connections aka "**ResNet**"
- Also known as **skip-connections**
- The identity connection preserves information by default
- This makes deep networks much easier to train

$$H(x) = F(x) + x$$
$$H(x) - x = F(x)$$



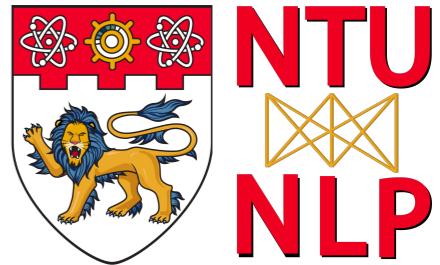
before going to final NonLinear layer

Why?

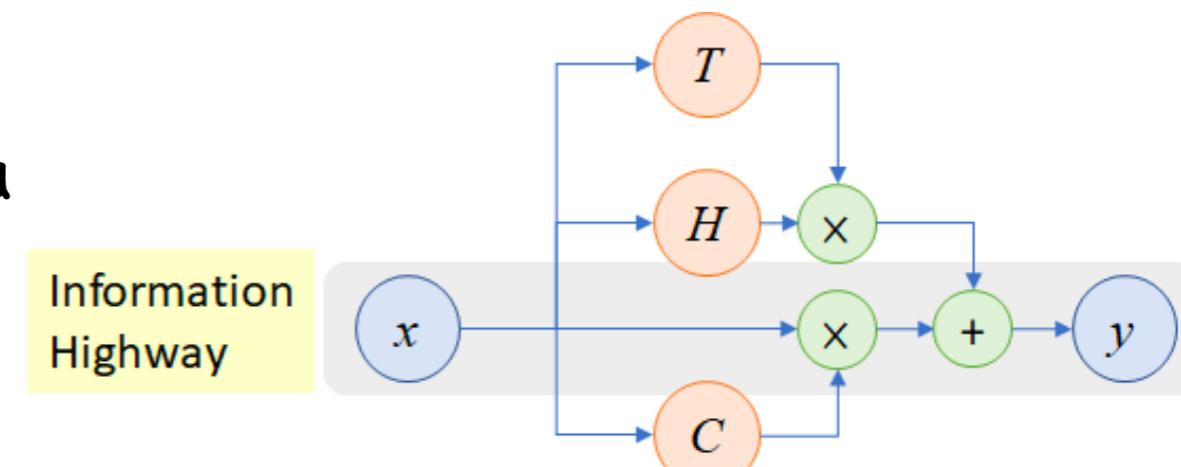
<https://arxiv.org/pdf/1512.03385.pdf>

Source: Stanford 224n

# Direct Connections: Highway Connections



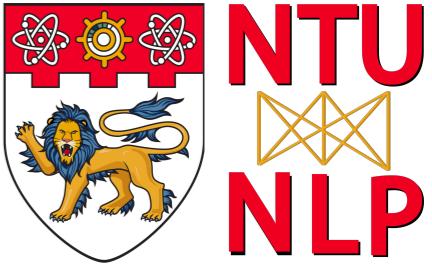
- Solution: lots of new deep feedforward/convolutional architectures **add more direct connections** (thus allowing the gradient to flow)
- Similar to residual connections, but the identity connection vs the regular (transformation) layer is controlled by a **dynamic gate**
- Inspired by LSTM's gating but applied to other architectures.
- Two non-linear transforms **T** and **C** are introduced: **T is the Transform Gate** and **C is the Carry Gate;  $C = 1 - T$**
- $H(x, W_H)$  is any hidden layer with  $x$  as input and  $W_H$  as weights



$$y = H(x, W_H) \cdot T(x, W_T) + x \cdot C(x, W_C).$$

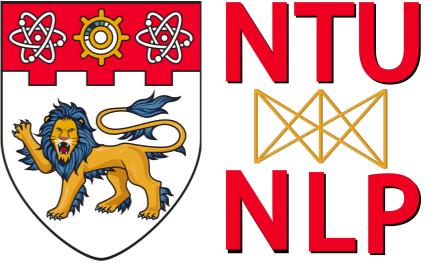
$$y = H(x, W_H) \cdot T(x, W_T) + x \cdot (1 - T(x, W_T)).$$

# Final Remark

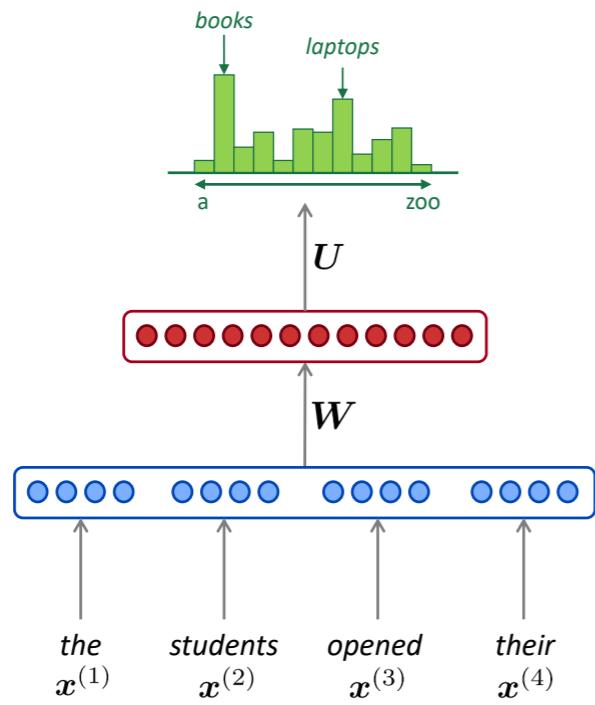


- Though vanishing/exploding gradients are a general problem, RNNs are particularly unstable due to the repeated multiplication by the same weight matrix [Bengio et al, 1994]

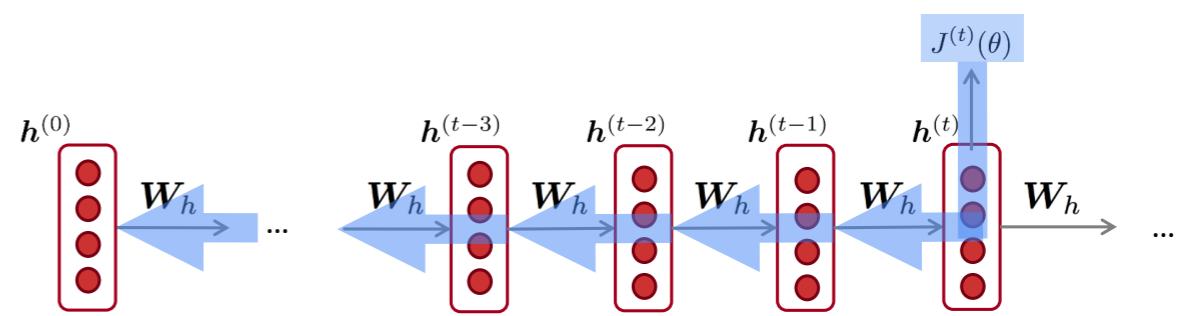
# Summary



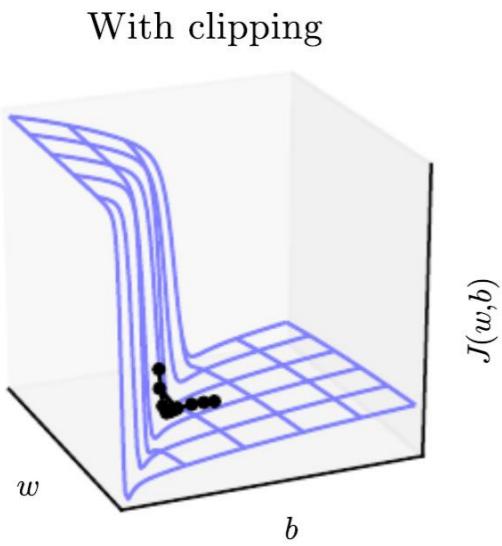
## RNN-LM



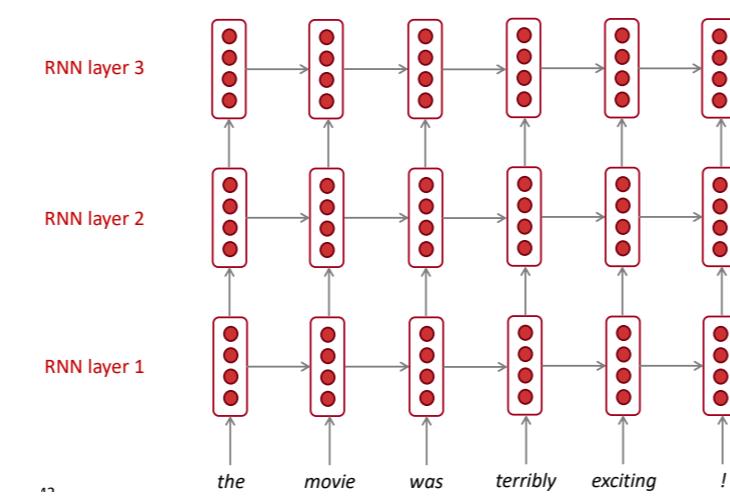
## BPTT



## Gradient clipping for exploding gradient



## Multi-layer RNNs



## LSTM/GRU

