

CSCE 638 Natural Language Processing Foundation and Techniques

Lecture 5: Convolutional Neural Network and Recurrent Neural Network

Kuan-Hao Huang

Spring 2025



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

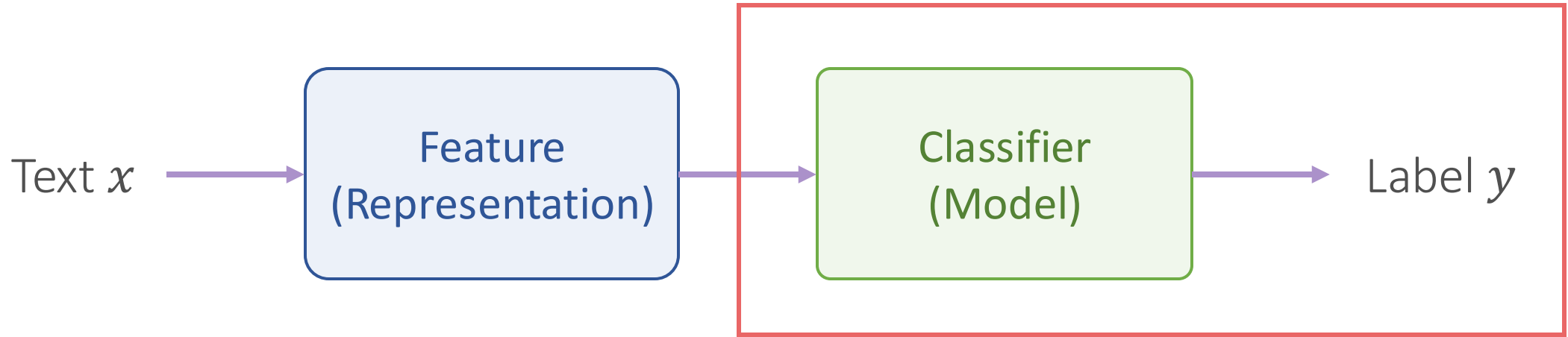
Assignment 0

- https://khhuang.me/CSCE638-S25/assignments/assignment0_0122.pdf
- Due: 1/29/2025 11:59pm
- Submit a .zip file to Canvas
 - [submission.pdf](#) for the writing section
 - [submission.py](#) and [submission.ipynb](#) for the coding section
- For questions
 - Discuss on Canvas
 - Send an email to csce638-ta-25s@list.tamu.edu

Lecture Plan

- Convolutional Neural Network
- Recurrent Neural Network
 - Long Short-Term Memory
 - Gated Recurrent Units

Recap: A General Framework for Text Classification

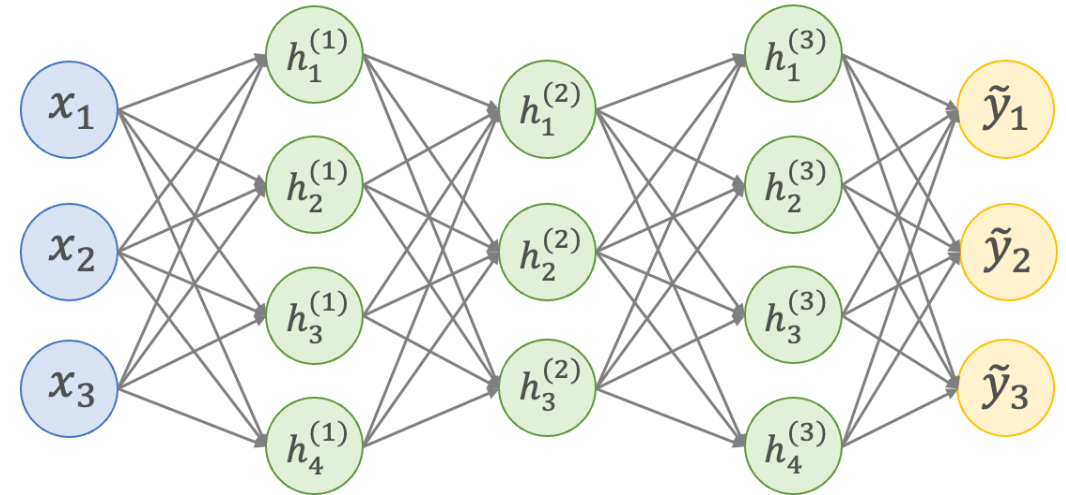


- Teach the model how to **make prediction y**
- Logistic regression, neural networks, CNN, RNN, LSTM, Transformers

A Simple Approach: Averaged Embeddings + DNN

	Alice	treats	Bob	well	
Dimension 1	0.7	2.7	-0.1	-5.7	-0.6
Dimension 2	8.6	-3.9	6.7	-9.8	0.4
Dimension 3	-2.4	-5.6	1.5	-1.6	-1.6
Dimension 4	2.3	1.1	2.0	-1.0	1.1

Any problems?

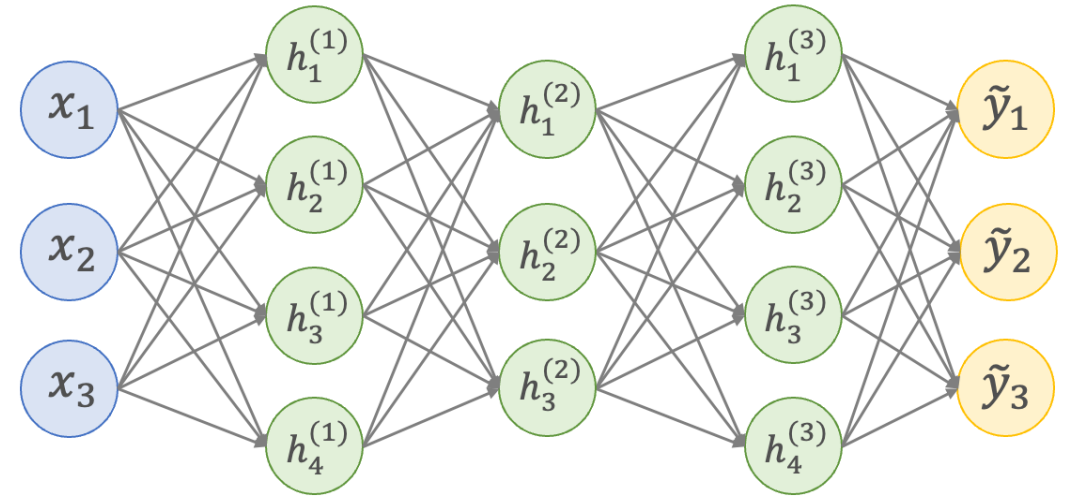


$$\mathcal{L}_{CE}(y, \tilde{y}) = - \sum_{c=0}^C y_c \log P(y = c | \mathbf{x})$$

A Simple Approach: Concatenated Embeddings + DNN

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

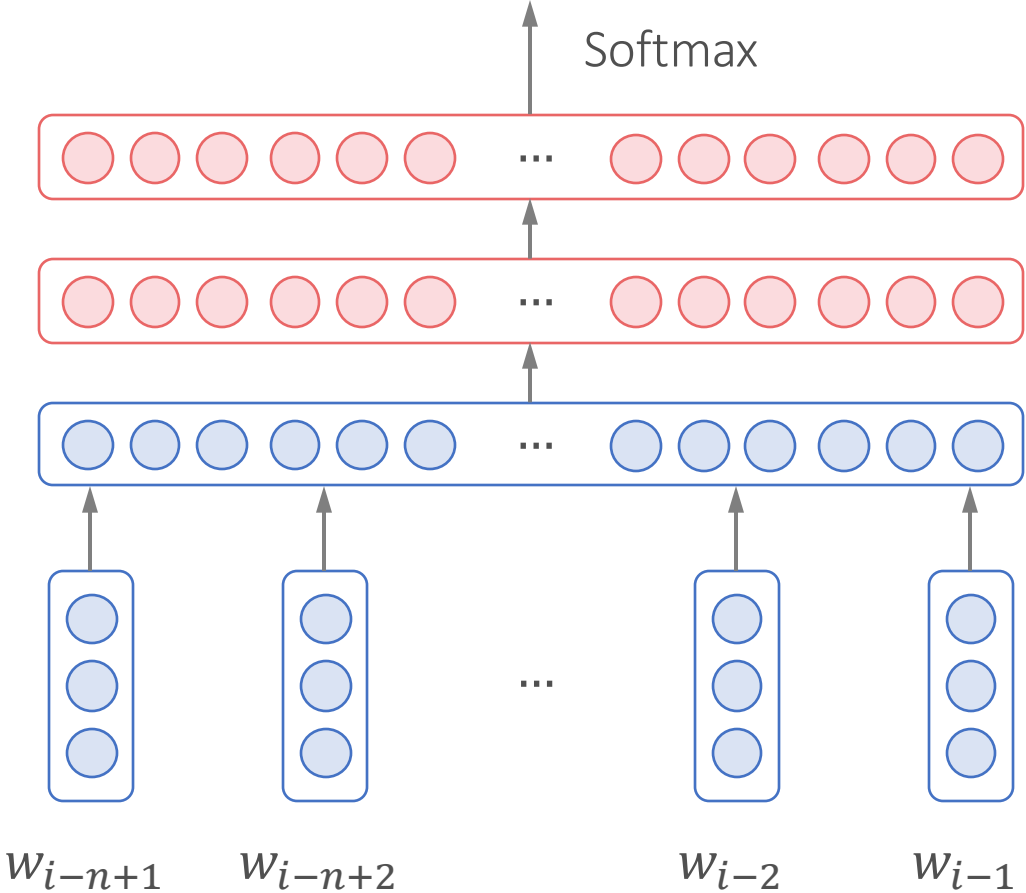
Any problems?



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

Recap: Neural Language Models

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



Training corpus

- I like apples
- I love oranges

Test set

- I love apples
- I like oranges

Challenges

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

Solution 1: Capture Local Order Information

Bob likes Alice very much

Unigram

{Bob, likes, Alice, very, much}

Bigram

{Bob likes, likes Alice, Alice very, very much}

Trigram

{Bob likes Alice, likes Alice very, Alice very much}

4-gram

{Bob likes Alice very, likes Alice very much}

We can infer global order information from local order information

Convolutional Neural Network (CNN)

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

Convolutional Neural Network (CNN)

Learnable Weight (Filter)
Filter Size = 3

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



Convolutional Neural Network (CNN)

Learnable Weight (Filter)
Filter Size = 3

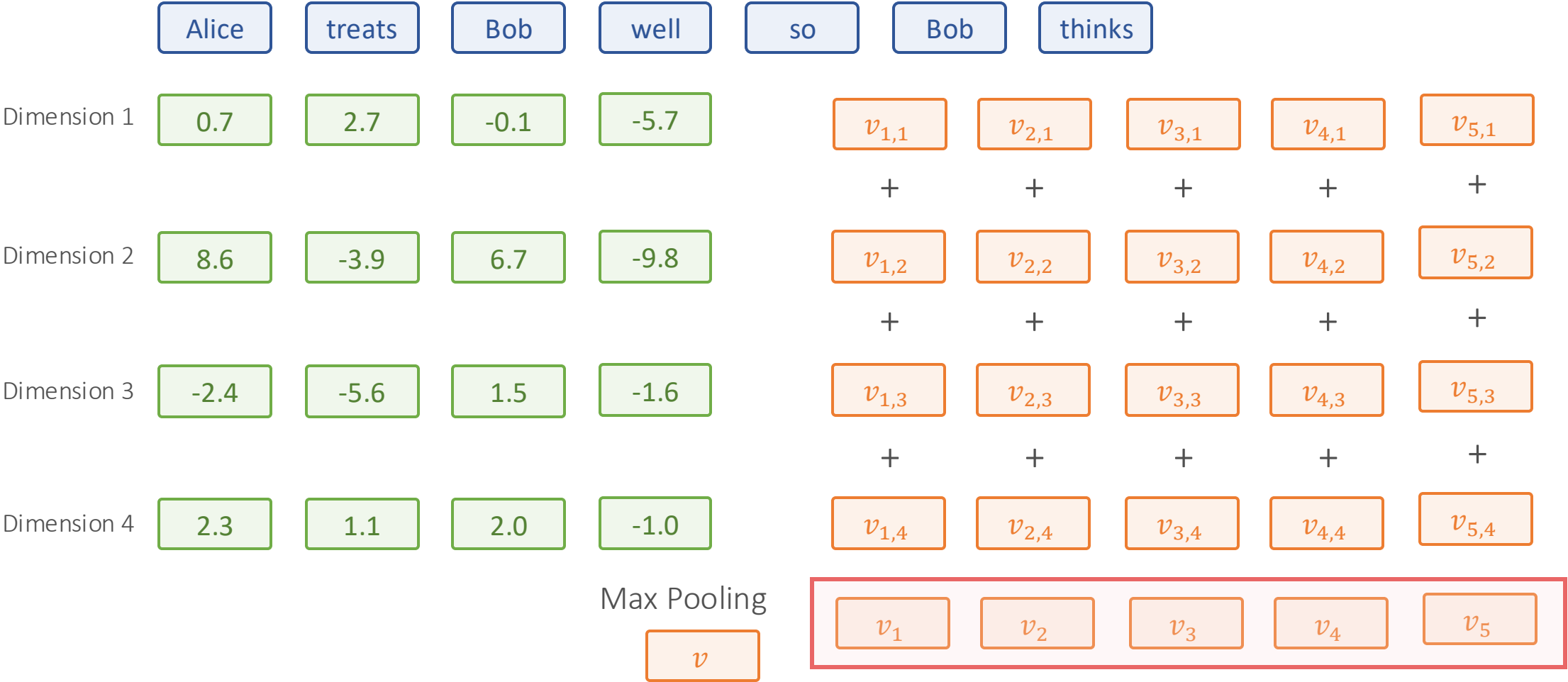
$$W = \begin{bmatrix} W_{1,1} & W_{1,2} & W_{1,3} \\ \dots & \dots & \dots \\ W_{4,1} & W_{4,2} & W_{4,3} \end{bmatrix}$$



Convolutional Neural Network (CNN)

Learnable Weight (Filter)
Filter Size = 3

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



Convolutional Neural Network (CNN)

Learnable Weight (Filter)
Filter Size = 3

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

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

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

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



Convolutional Neural Network (CNN)

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

Filter Size = 2 $W = \begin{bmatrix} W_{1,1} & W_{1,2} \\ \dots & \dots \\ W_{4,1} & W_{4,2} \end{bmatrix}$

Filter Size = 4 $W = \begin{bmatrix} W_{1,1} & \dots & W_{1,4} \\ \dots & \dots & \dots \\ W_{4,1} & \dots & W_{4,4} \end{bmatrix}$

Alice treats Bob well

Dimension 1

0.7	2.7	-0.1	-5.7
-----	-----	------	------

$W_{1,1} * 0.7 + W_{1,2} * 2.7$

Dimension 2

8.6	-3.9	6.7	-9.8
-----	------	-----	------

$W_{2,1} * 8.6 + W_{2,2} * -3.9 + W_{2,3} * 6.7 + W_{2,4} * -9.8$

Dimension 3

-2.4	-5.6	1.5	-1.6
------	------	-----	------

Dimension 4

2.3	1.1	2.0	-1.0
-----	-----	-----	------

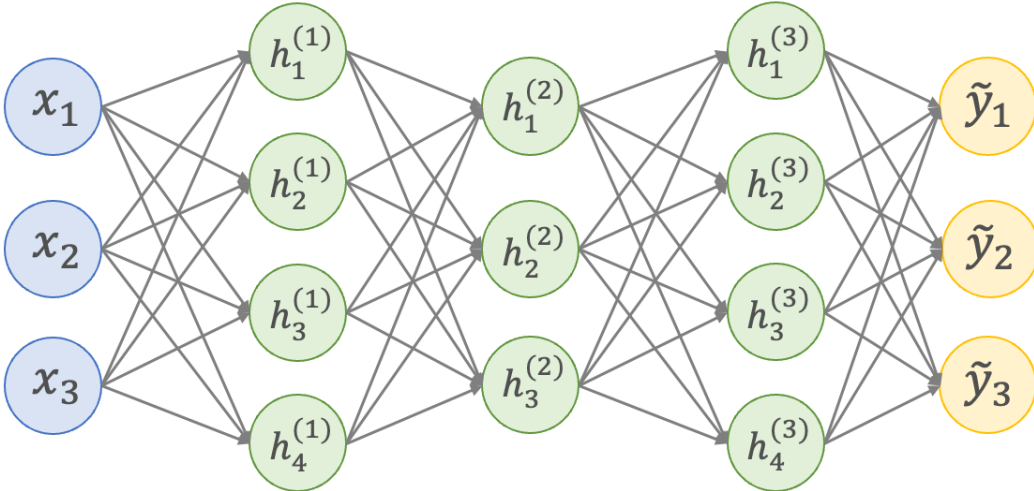
v

v

v

Convolutional Neural Network (CNN)

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



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

From Single Layer to Multiple Layers

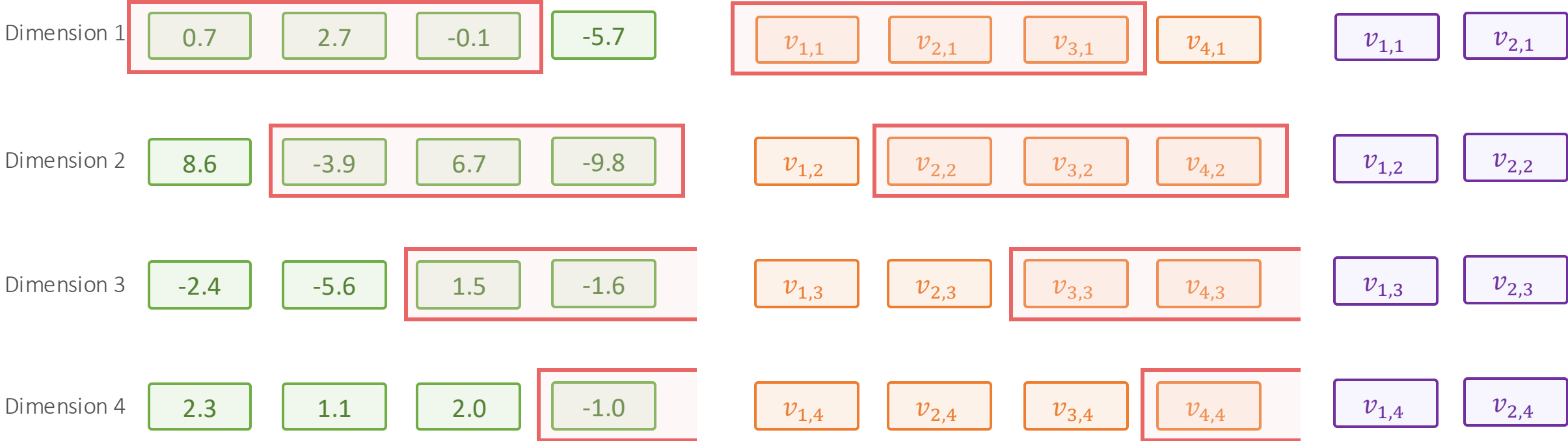
Learnable Weight (Layer 1)
Filter Size = 3

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

Learnable Weight (Layer 2)
Filter Size = 3

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

Alice treats Bob well so Bob thinks



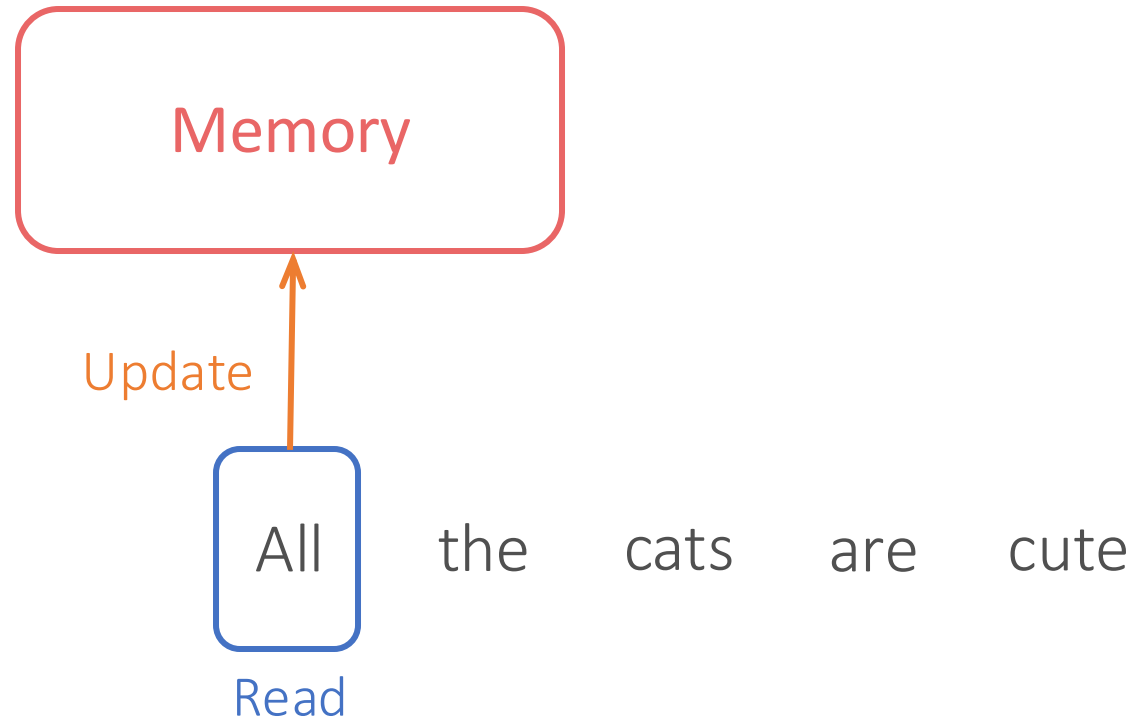
Capture high-order or hierarchical information

Convolutional Neural Network (CNN)

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

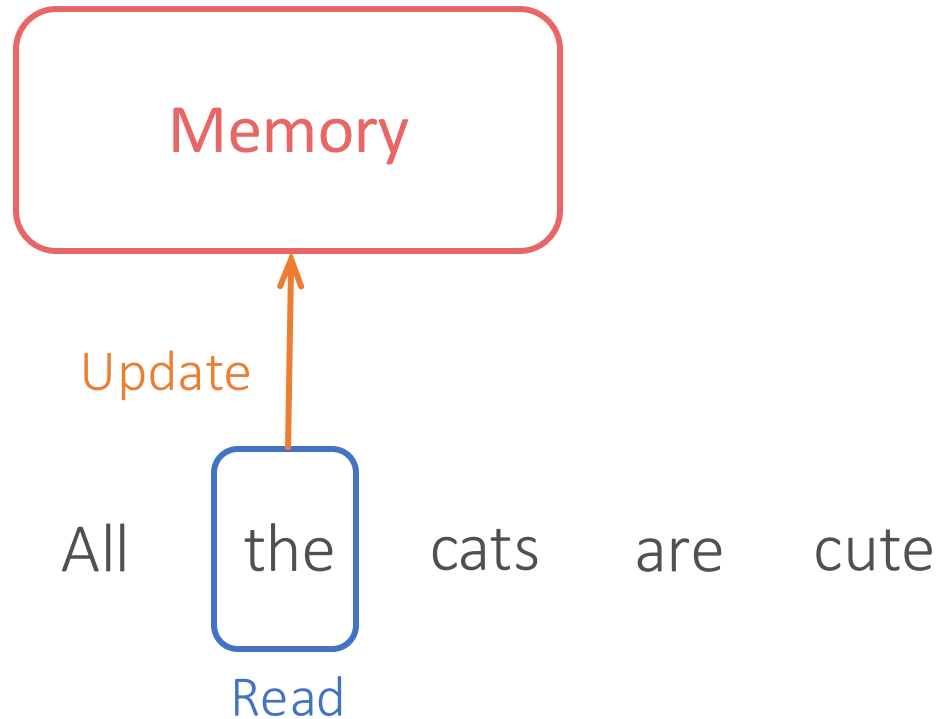
Recurrent Neural Network (RNN)

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



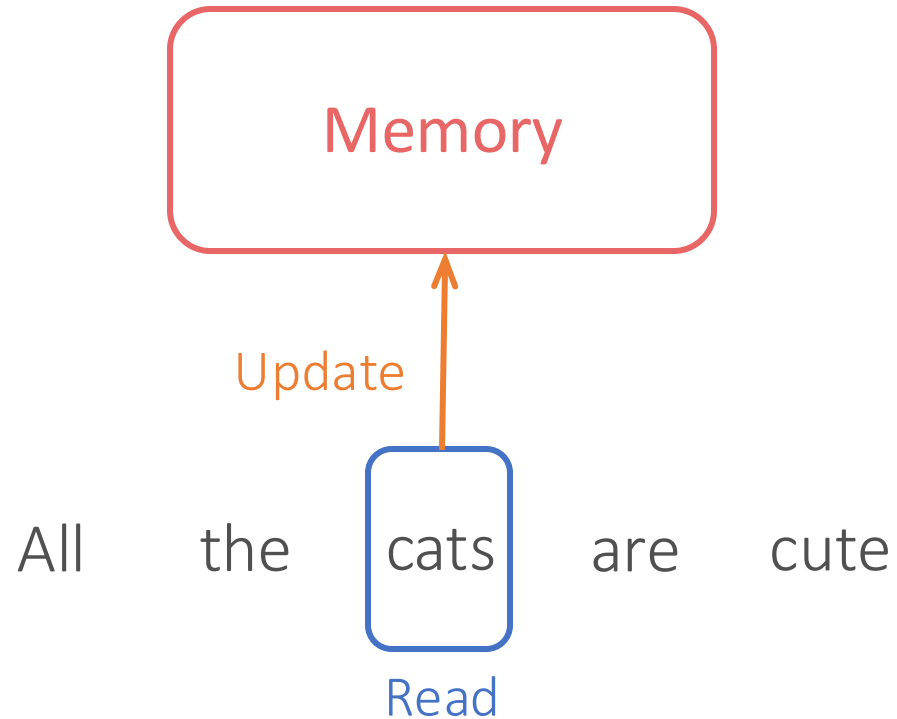
Recurrent Neural Network (RNN)

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



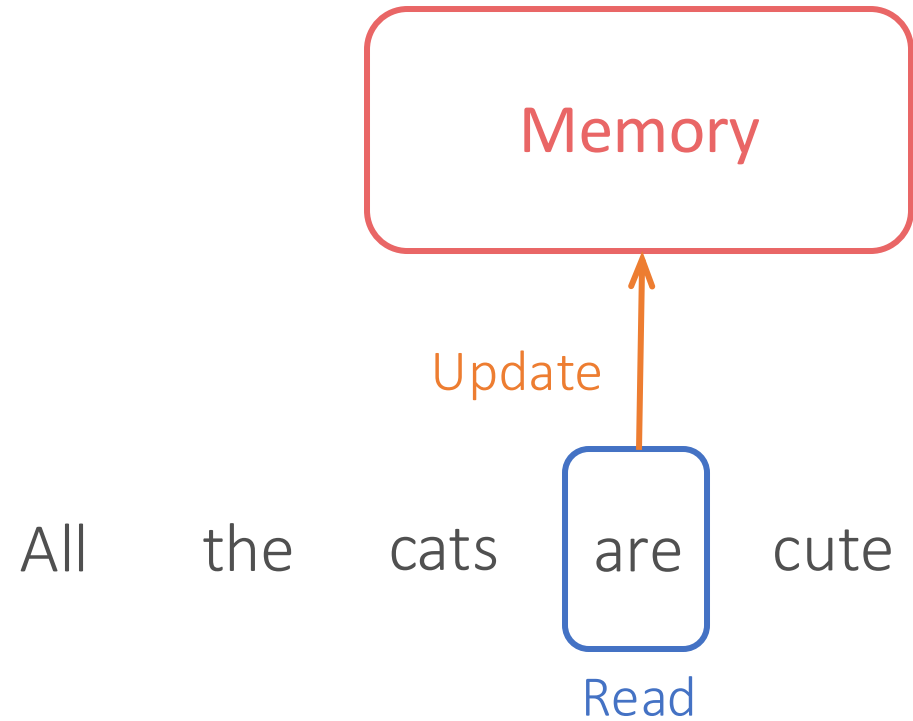
Recurrent Neural Network (RNN)

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



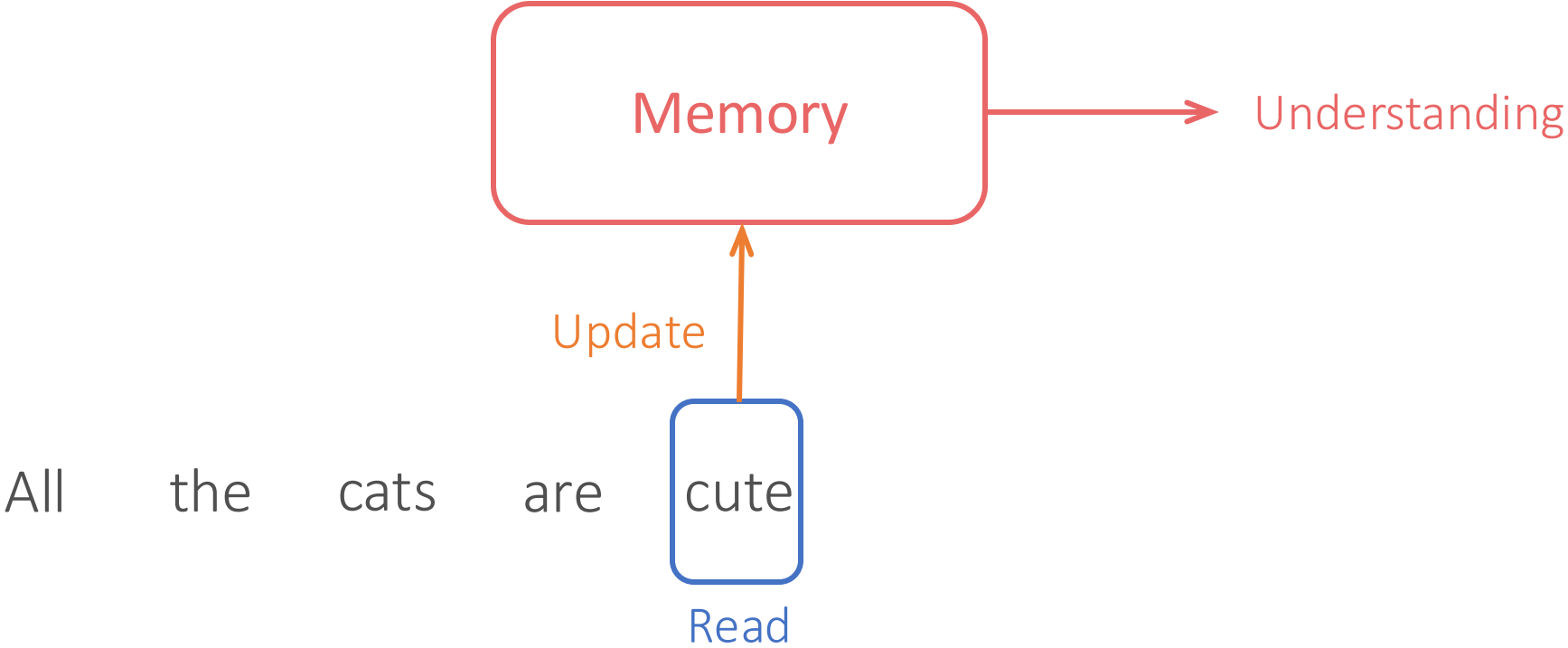
Recurrent Neural Network (RNN)

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



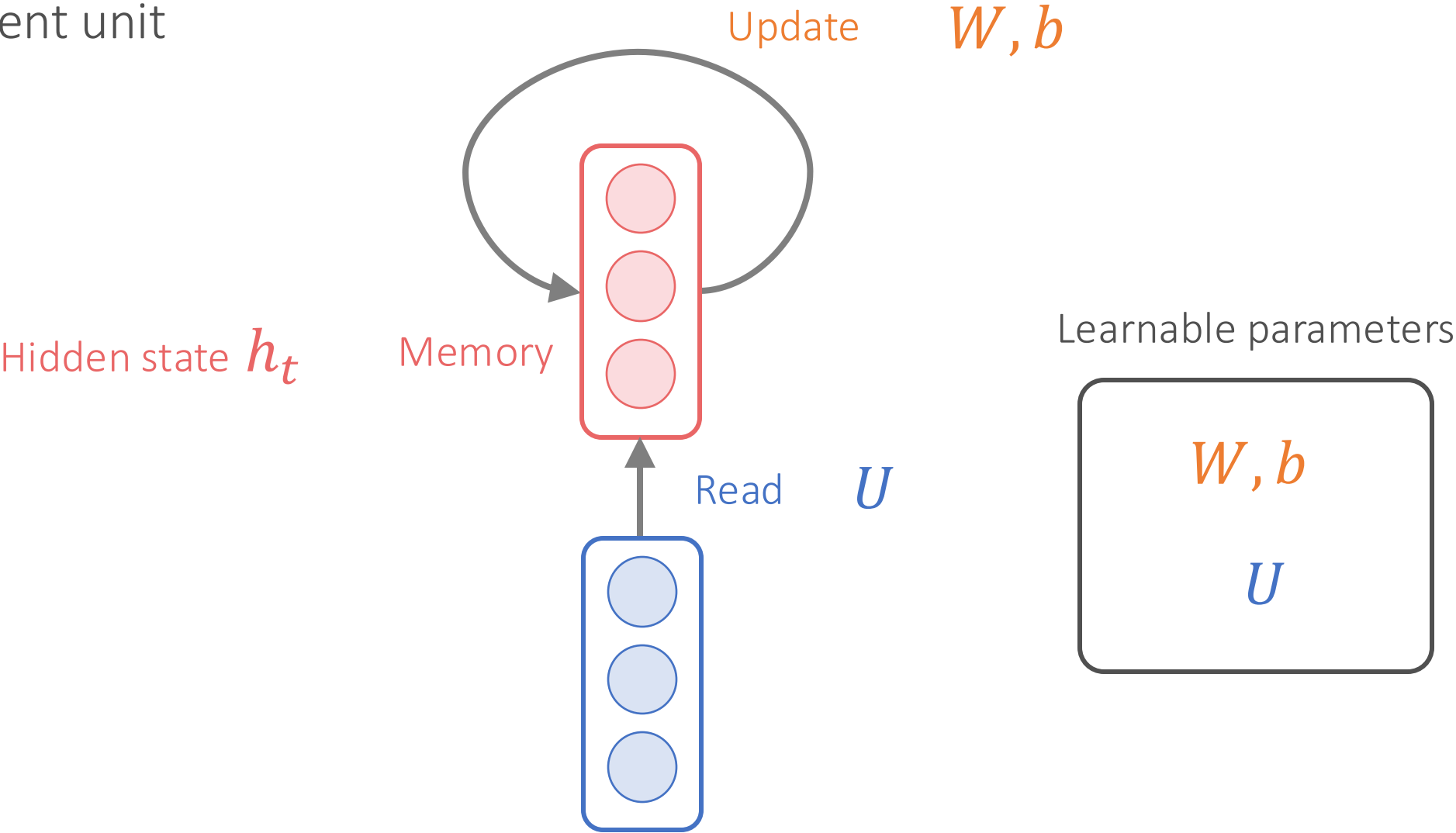
Recurrent Neural Network (RNN)

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

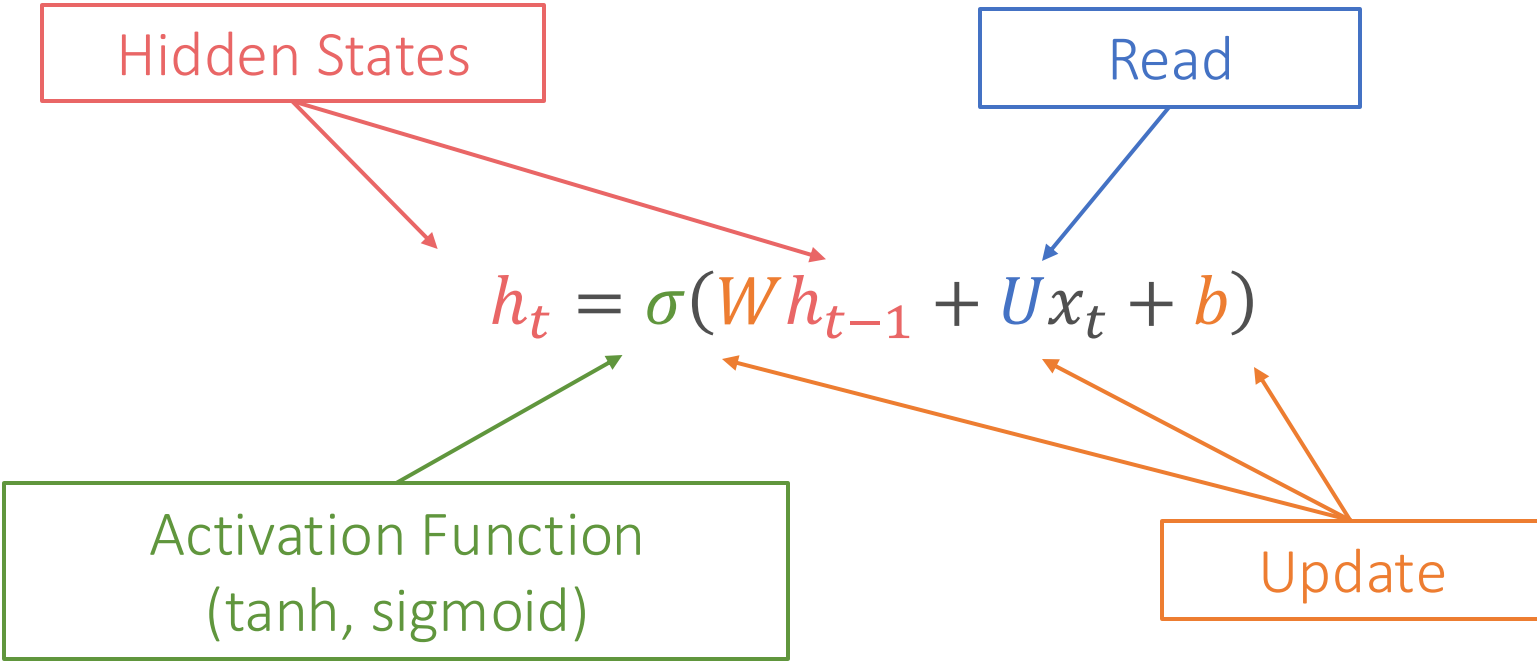


Recurrent Neural Network (RNN)

- Recurrent unit

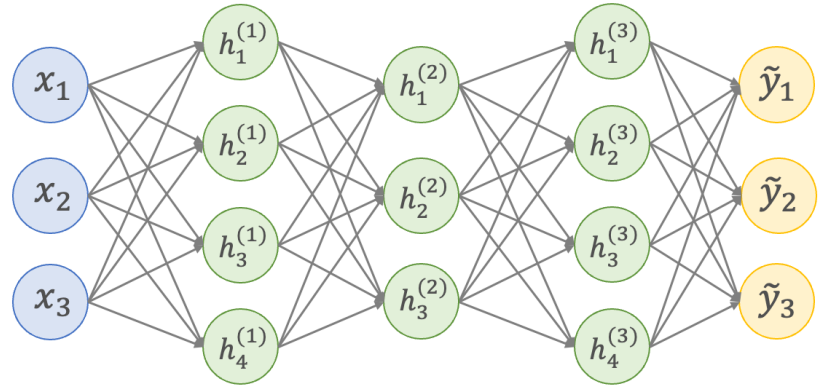
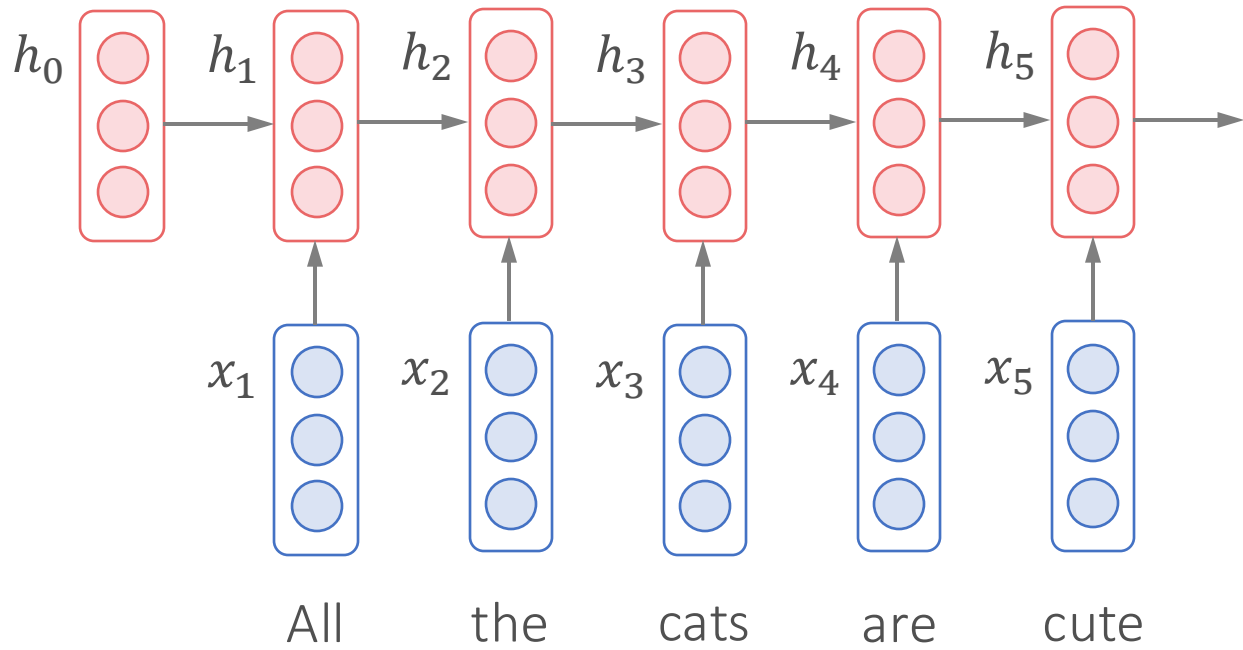


Recurrent Neural Network (RNN)



Recurrent Neural Network (RNN)

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

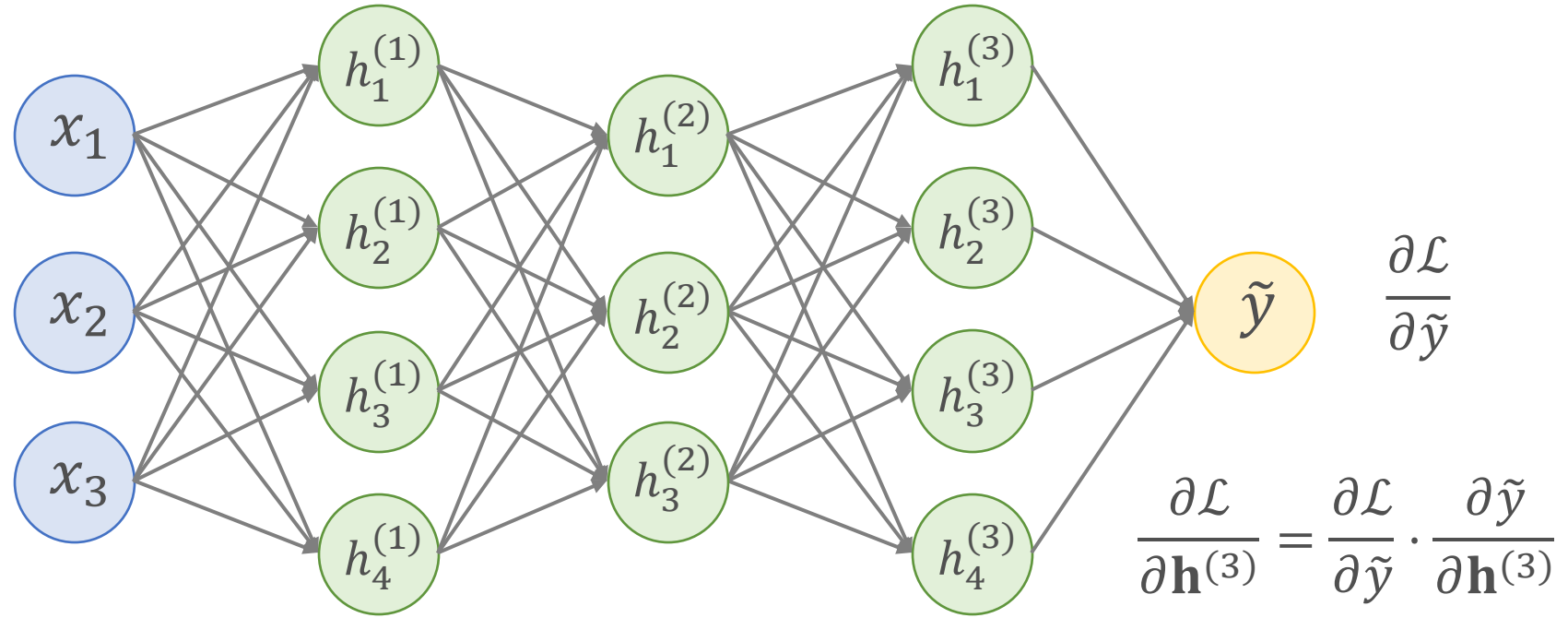


$$\mathcal{L}_{CE}(y, \tilde{y}) = - \sum_{c=0}^C y_c \log P(y = c | \mathbf{x})$$

Recurrent Neural Network (RNN)

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

Back-Propagation



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

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

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

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

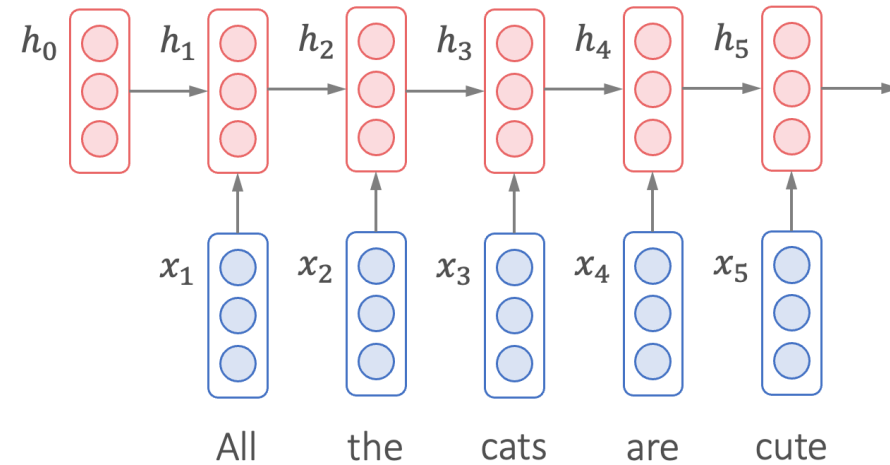
Vanishing Gradient Problem

$$h_2 = \sigma(W h_1 + U x_2 + b)$$

$$h_3 = \sigma(W h_2 + U x_3 + b)$$

$$h_4 = \sigma(W h_3 + U x_4 + b)$$

$$h_5 = \sigma(W h_4 + U x_5 + b)$$



$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial W} = & \frac{\partial \mathcal{L}}{\partial h_5} \cdot \frac{\partial h_5}{\partial W} + \frac{\partial \mathcal{L}}{\partial h_5} \cdot \frac{\partial h_5}{\partial h_4} \cdot \frac{\partial h_4}{\partial W} + \frac{\partial \mathcal{L}}{\partial h_5} \cdot \frac{\partial h_5}{\partial h_4} \cdot \frac{\partial h_4}{\partial h_3} \cdot \frac{\partial h_3}{\partial W} \\ & + \frac{\partial \mathcal{L}}{\partial h_5} \cdot \frac{\partial h_5}{\partial h_4} \cdot \frac{\partial h_4}{\partial h_3} \cdot \frac{\partial h_3}{\partial h_2} \cdot \frac{\partial h_2}{\partial W} + \frac{\partial \mathcal{L}}{\partial h_5} \cdot \frac{\partial h_5}{\partial h_4} \cdot \frac{\partial h_4}{\partial h_3} \cdot \frac{\partial h_3}{\partial h_2} \cdot \frac{\partial h_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial W} \end{aligned}$$

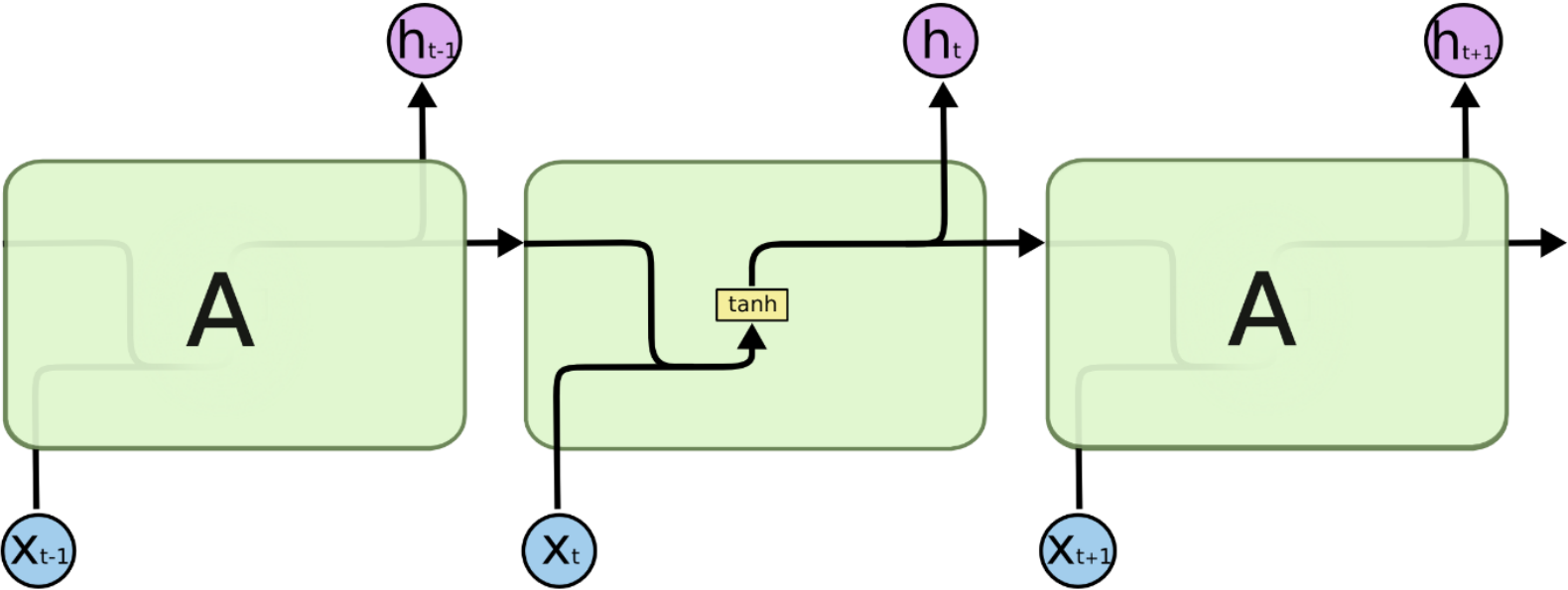
When these are small, the gradient signal gets smaller and smaller as it back-propagates further

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

Long Short-Term Memory (LSTM)

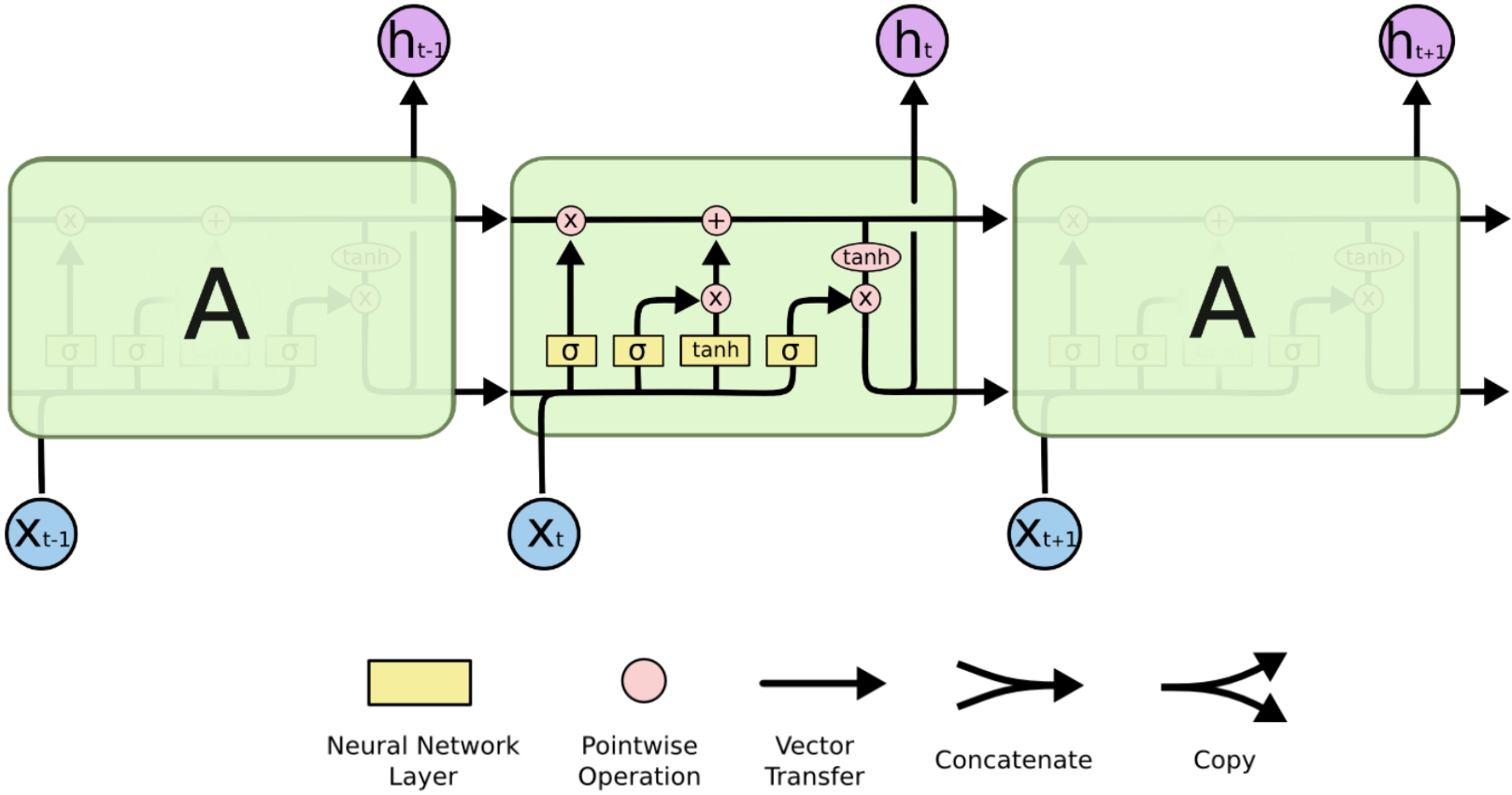
- Short-term memory: hidden state h_t
- Long-term memory: cell state c_t
- Key idea
 - Turn **multiplication** into **addition** (partially reduce gradient vanishing)
 - use **gates** to control how much information to add/erase

Recurrent Neural Network (RNN)

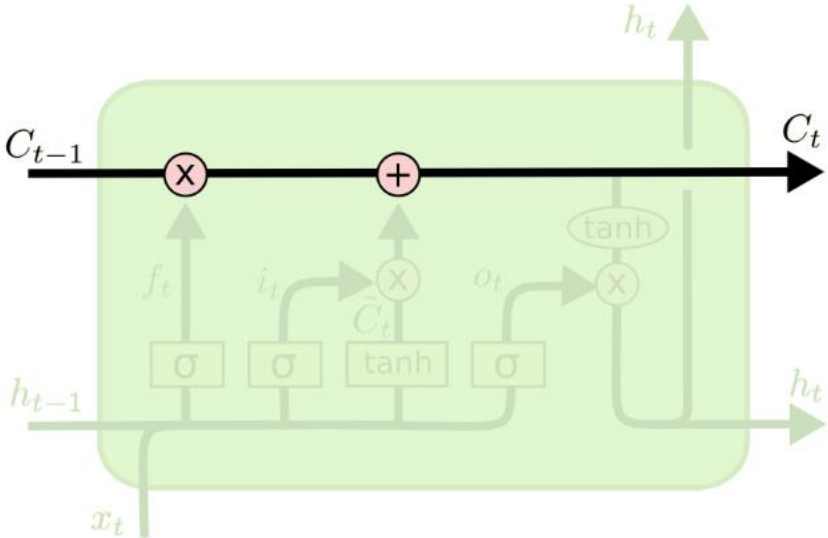


$$h_t = \sigma(W h_{t-1} + U x_t + b)$$

Long Short-Term Memory (LSTM)

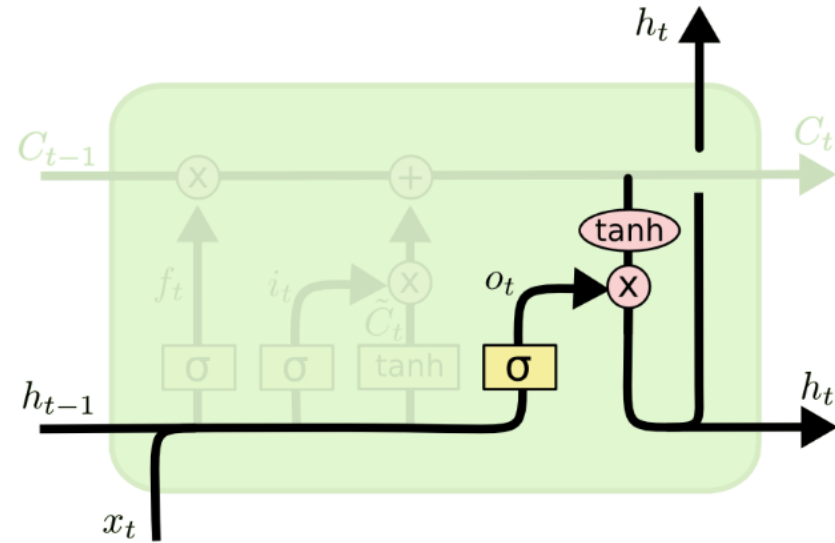


Long Short-Term Memory (LSTM)



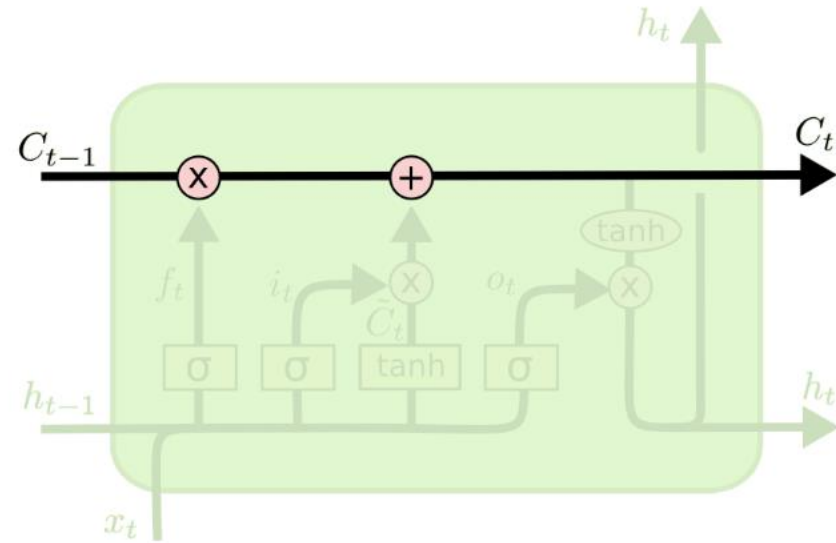
The cell state stores **long-term information**

Long Short-Term Memory (LSTM)



The hidden state stores **short-term information**

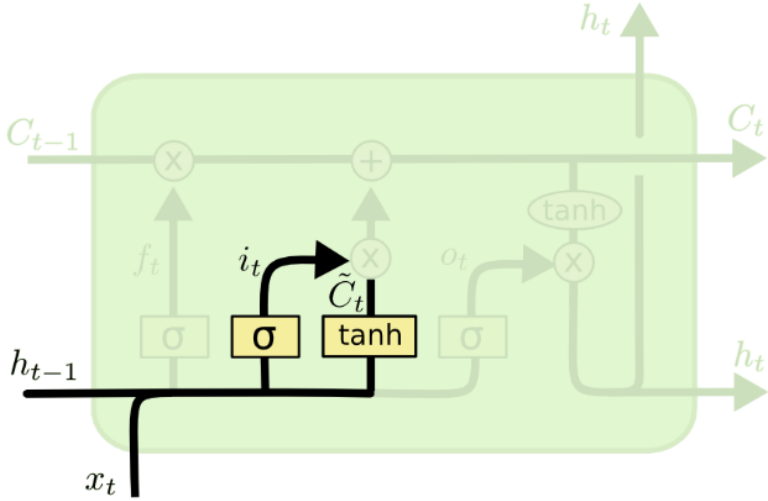
Long Short-Term Memory (LSTM)



The cell state stores **long-term information**

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

Long Short-Term Memory (LSTM)



Sigmoid function: gate values are between 0 and 1

Input gate

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

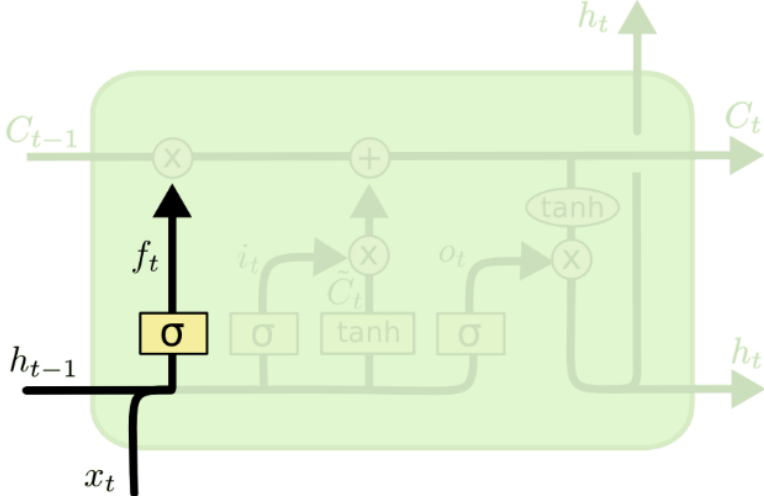
How much we should write

New information

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

What we should write

Long Short-Term Memory (LSTM)



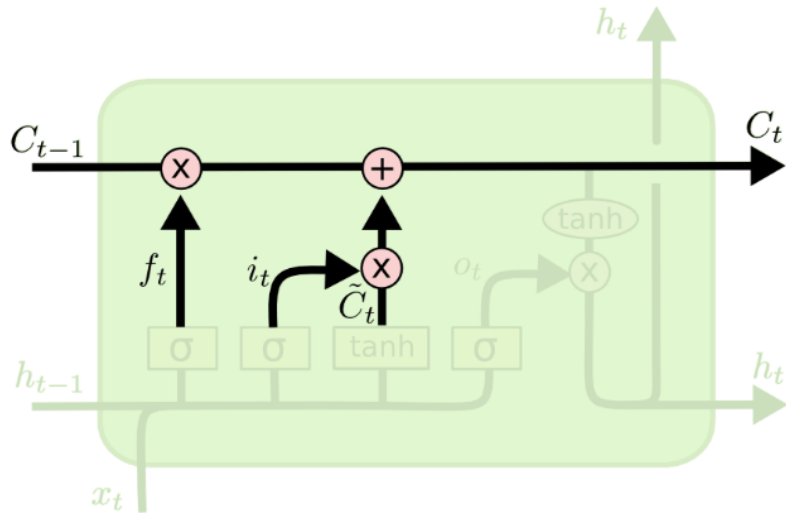
Sigmoid function: gate values are between 0 and 1

Forget gate

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

How much we should **erase**

Long Short-Term Memory (LSTM)



What we should write

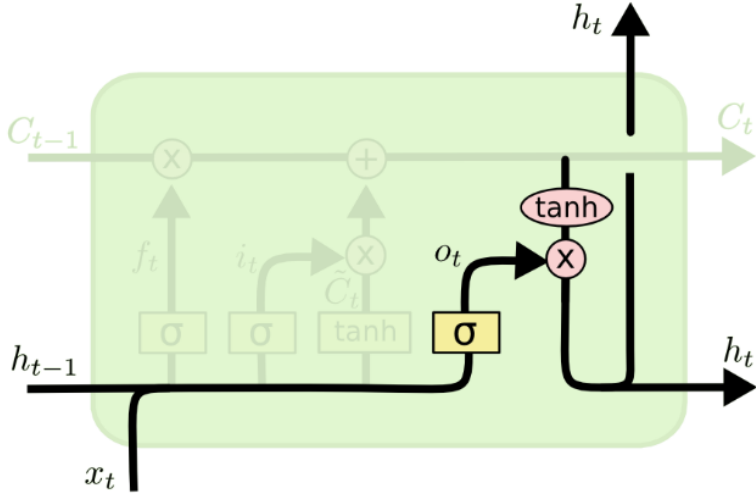
Update cell state

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

How much we should erase

How much we should write

Long Short-Term Memory (LSTM)



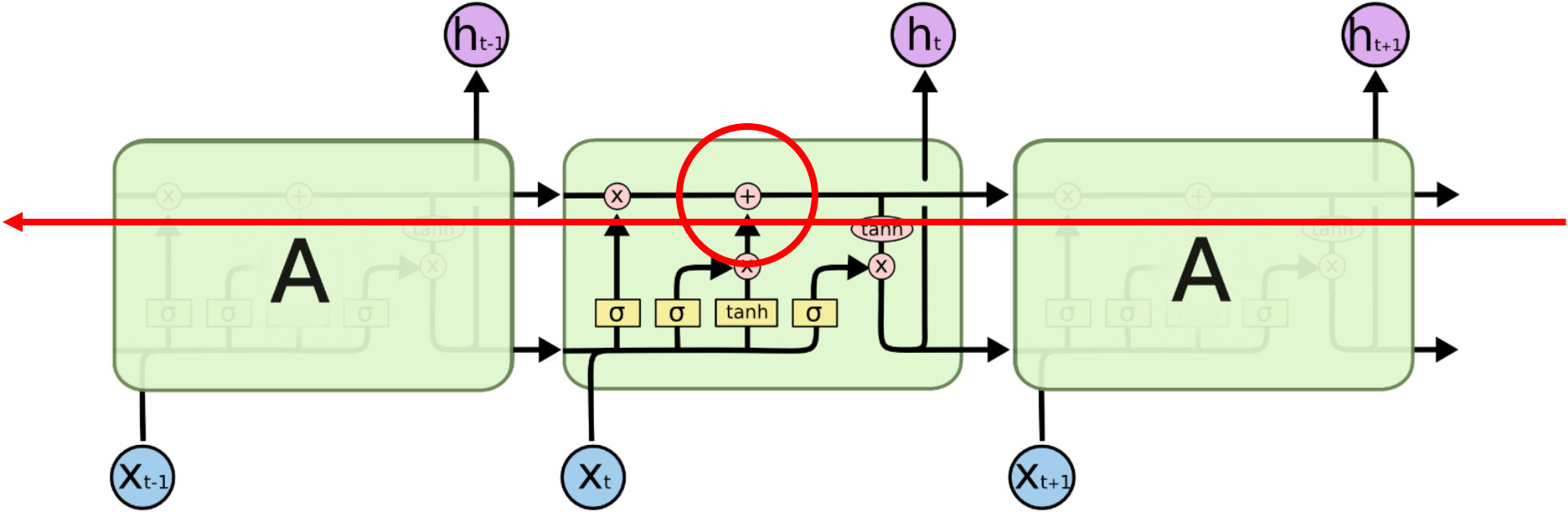
$$o_t = \sigma(W^{(o)}h_{t-1} + U^{(o)}x_t + b^{(o)})$$

Update hidden state

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

Long Short-Term Memory (LSTM)

Uninterrupted gradient flow

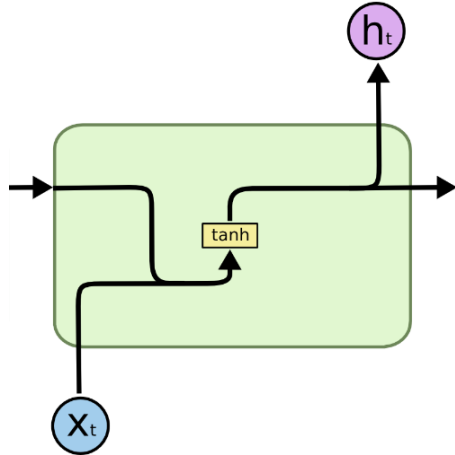


The addition is the key

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

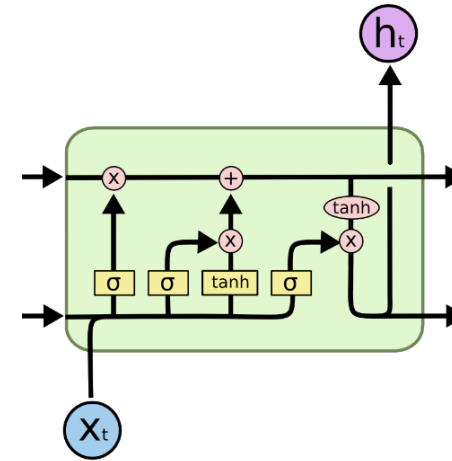
RNN vs. LSTM

RNN



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

LSTM



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

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

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

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

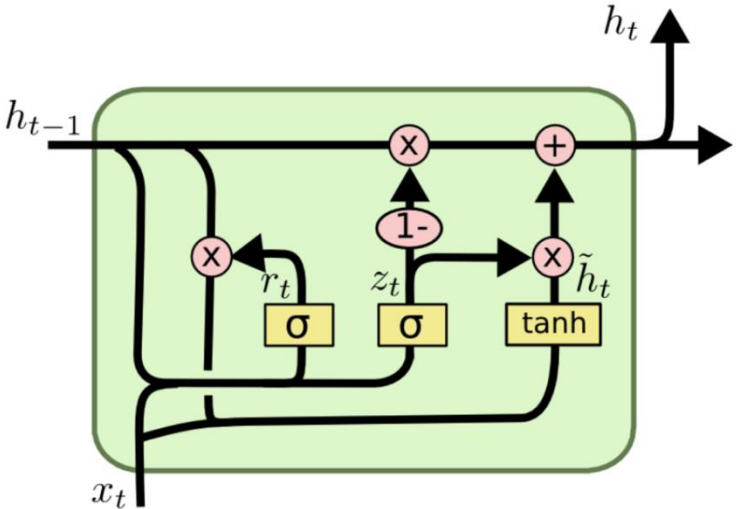
$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

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

Gated Recurrent Units (GRU)

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

Gated Recurrent Units (GRU)



Reset gate

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

Update gate

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

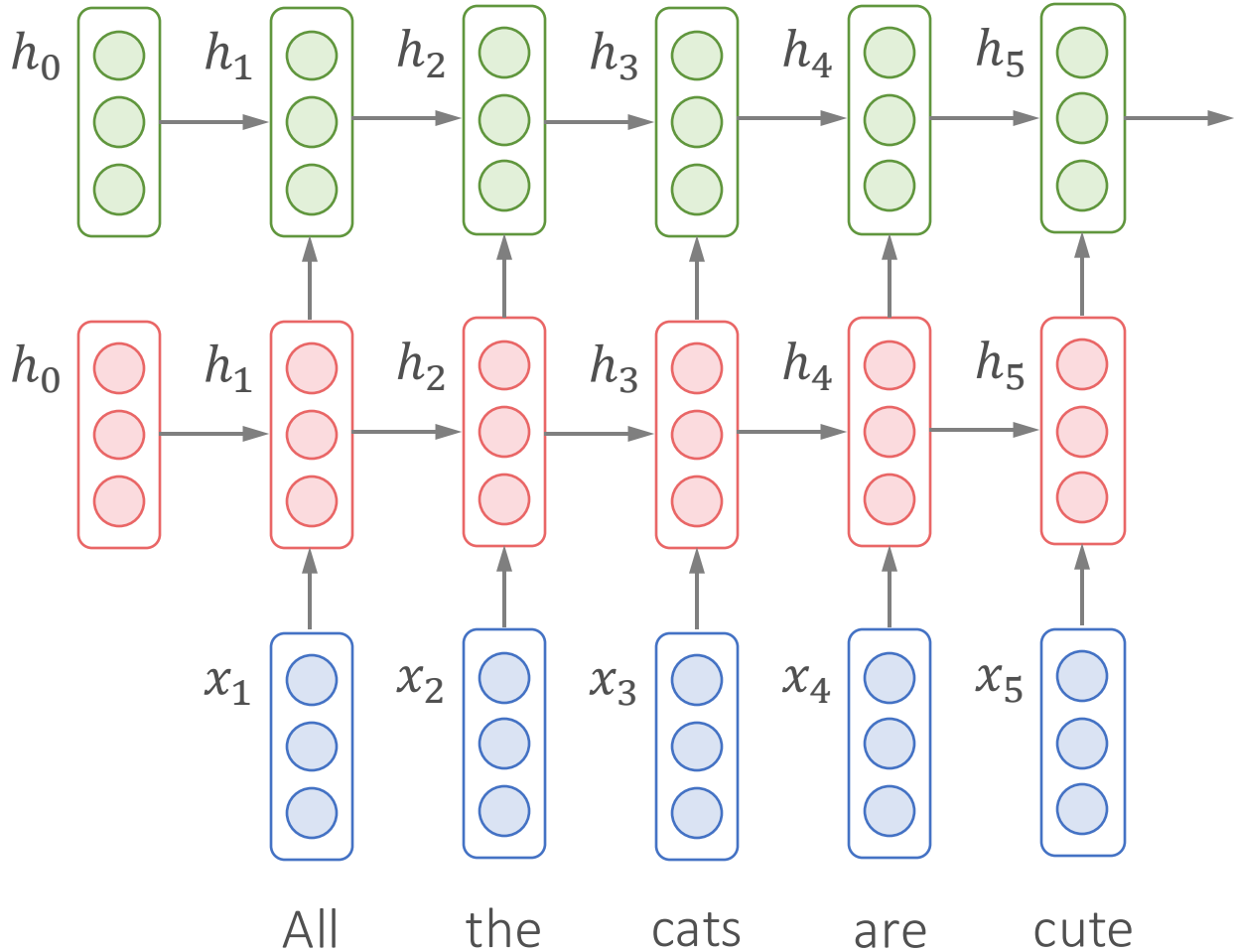
New hidden state

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

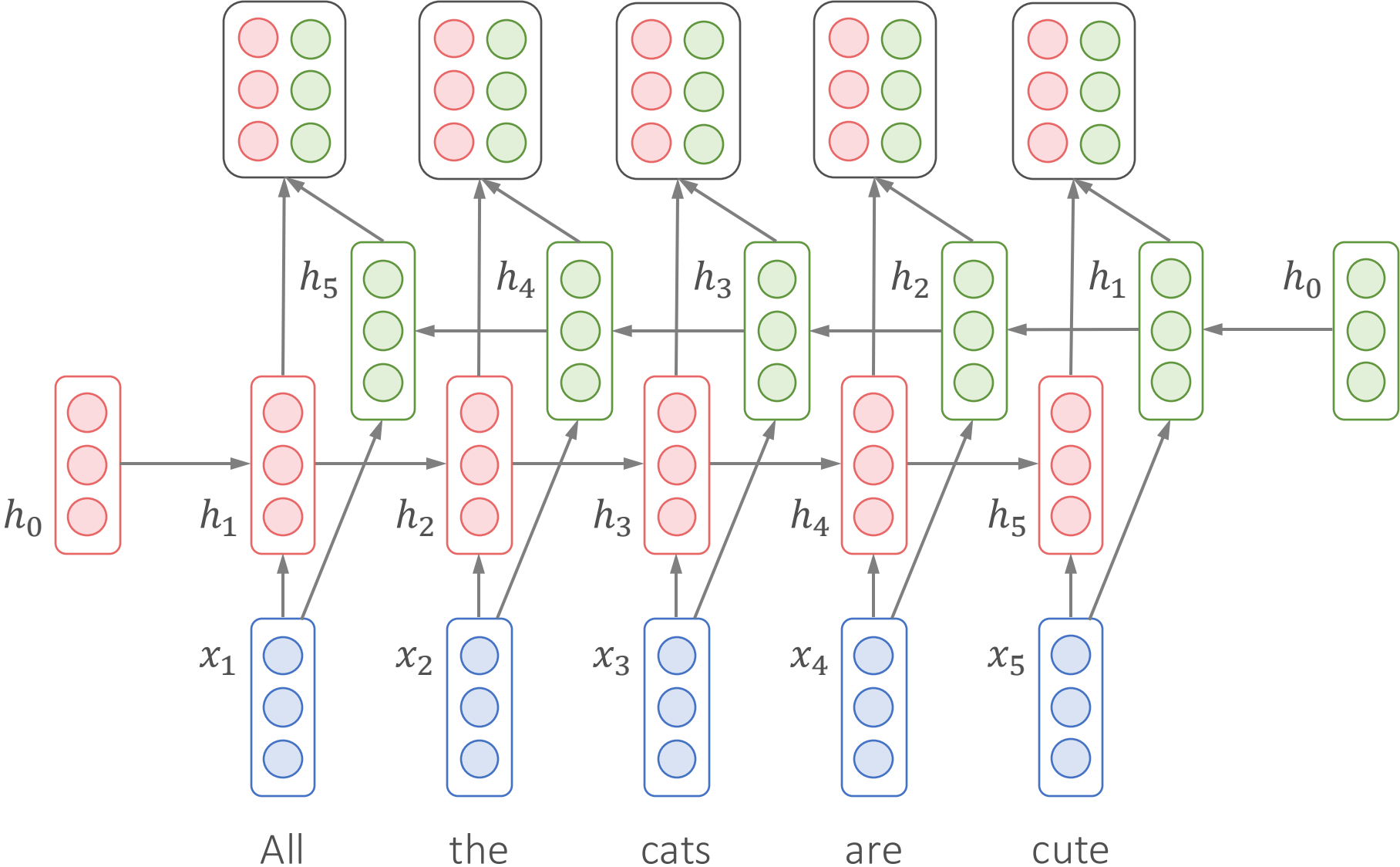
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Merge input and forget gate

Multi-Layer RNN (Stacked RNN)



Bidirectional RNN



Lecture Plan

- Convolutional Neural Network
- Recurrent Neural Network
 - Long Short-Term Memory
 - Gated Recurrent Units