

# Cornell Bowers C-IS

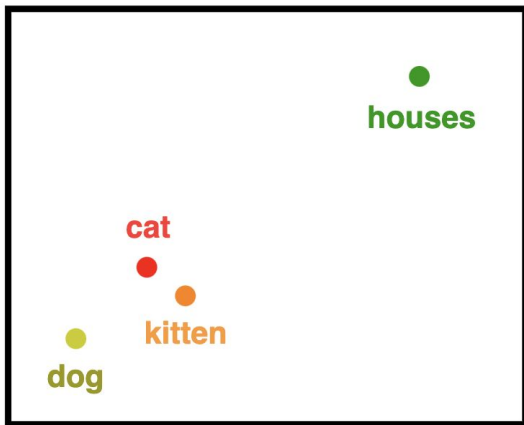
College of Computing and Information Science

# Deep Learning

Week 03: RNN/LSTM

# Why Do We Need Word Embeddings?

- Numerical Input
- Shows Similarity and Distance



	<i>living being</i>	<i>feline</i>	<i>human</i>	<i>gender</i>	<i>royalty</i>	<i>verb</i>	<i>plural</i>
<b>cat</b>	0.6	0.9	0.1	0.4	-0.7	-0.3	-0.2
<b>kitten</b>	0.5	0.8	-0.1	0.2	-0.6	-0.5	-0.1
<b>dog</b>	0.7	-0.1	0.4	0.3	-0.4	-0.1	-0.3
<b>houses</b>	-0.8	-0.4	-0.5	0.1	-0.9	0.3	0.8
<b>man</b>	0.6	-0.2	0.8	0.9	-0.1	-0.9	-0.7
<b>woman</b>	0.7	0.3	0.9	-0.7	0.1	-0.5	-0.4
<b>king</b>	0.5	-0.4	0.7	0.8	0.9	-0.7	-0.6
<b>queen</b>	0.8	-0.1	0.8	-0.9	0.8	-0.5	-0.9

embedding using features of words

## Word2Vec

- We want vectors for words so that the context of a word can suggest the vector of this word, and vice versa
- Idea: **Similar words appear in similar contexts**

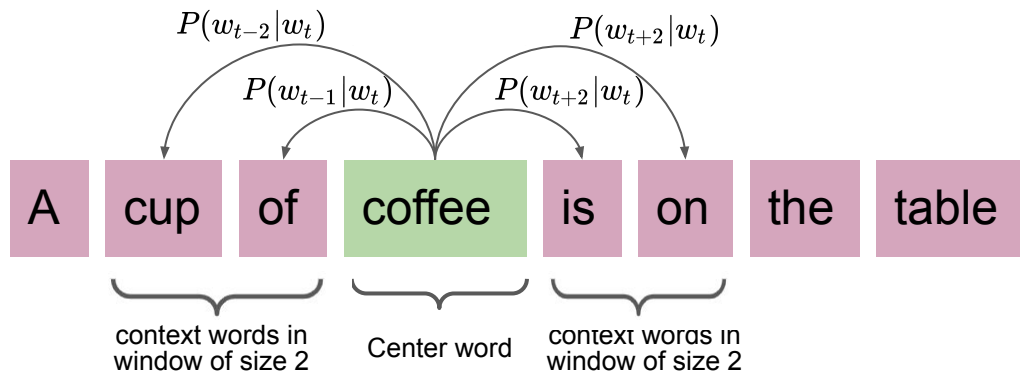
A cup of **coffee** is on the table.

**Coffee** helps me focus.

Espresso is my favorite type of **coffee**.

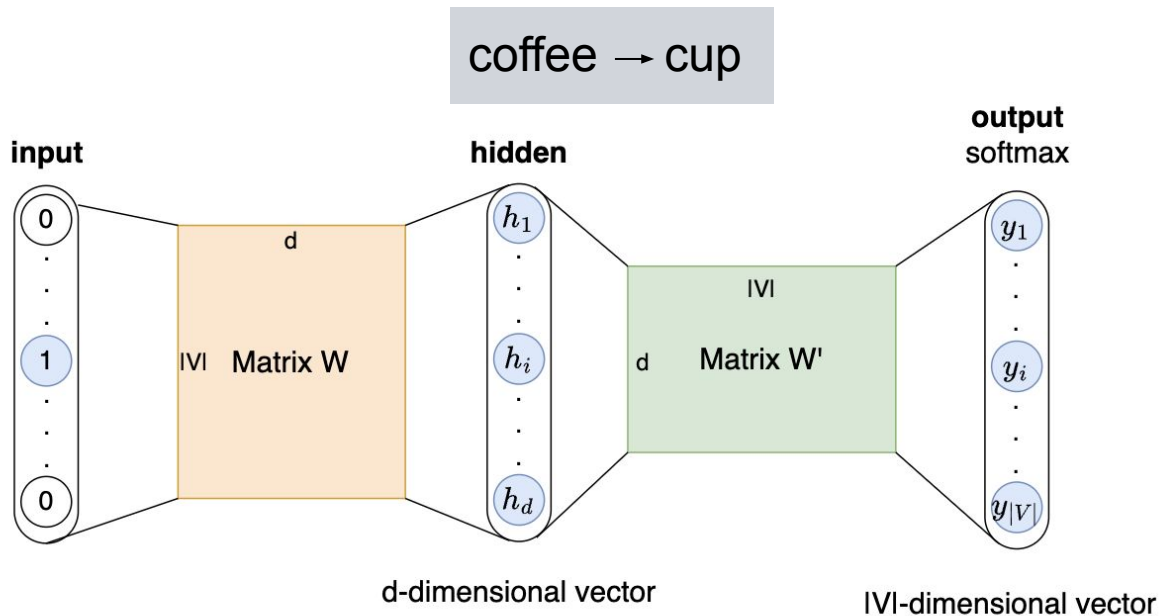
# Word2Vec - Training

SkipGram - Predict context from target



# Word2Vec Architecture - SkipGram

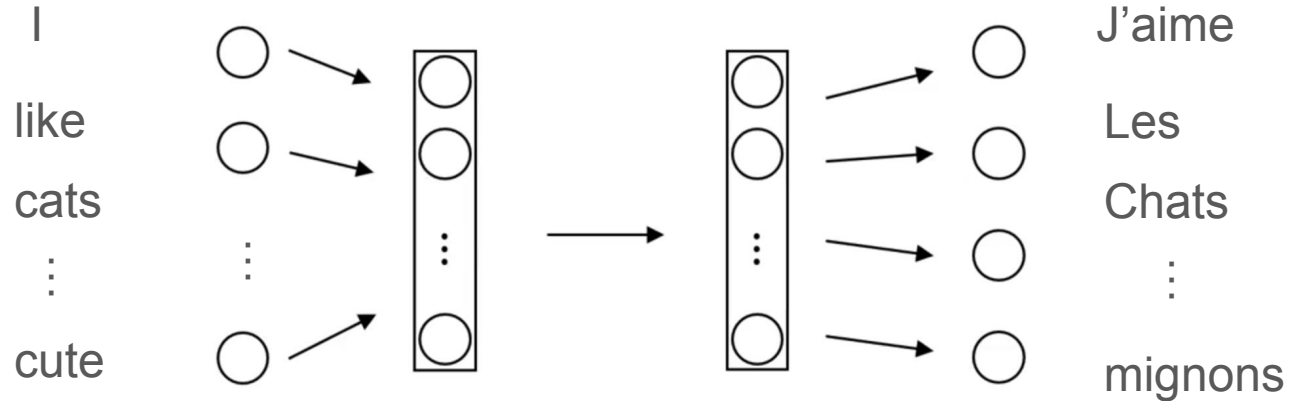
Predict every target word from each context word!



# Big Question: How to model sequences of words?

I	0.7	-0.1	0.4	0.3	-0.4	-0.1	-0.3	$\mathbf{v}_1 \in \mathbb{R}^d$
like	0.6	-0.2	0.8	0.9	-0.1	-0.9	-0.7	$\mathbf{v}_2 \in \mathbb{R}^d$
cats	0.6	0.9	0.1	0.4	-0.7	-0.3	-0.2	
because	0.5	0.8	-0.1	0.2	-0.6	-0.5	-0.1	$\vdots$
they	-0.8	-0.4	-0.5	0.1	-0.9	0.3	0.8	
look	0.5	-0.4	0.7	0.8	0.9	-0.7	-0.6	
cute	0.8	-0.1	0.8	-0.9	0.8	-0.5	-0.9	$\mathbf{v}_T \in \mathbb{R}^d$

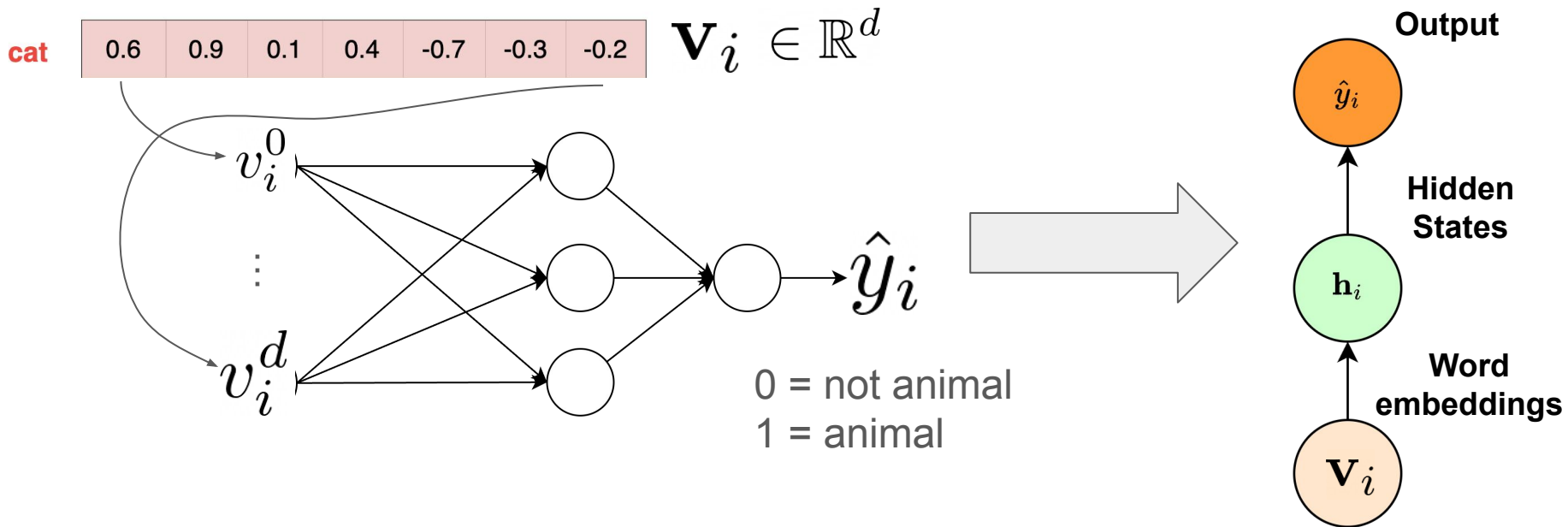
# How to use word vectors with neural networks?



- Inputs and outputs don't have fixed lengths
- Weights are not shared

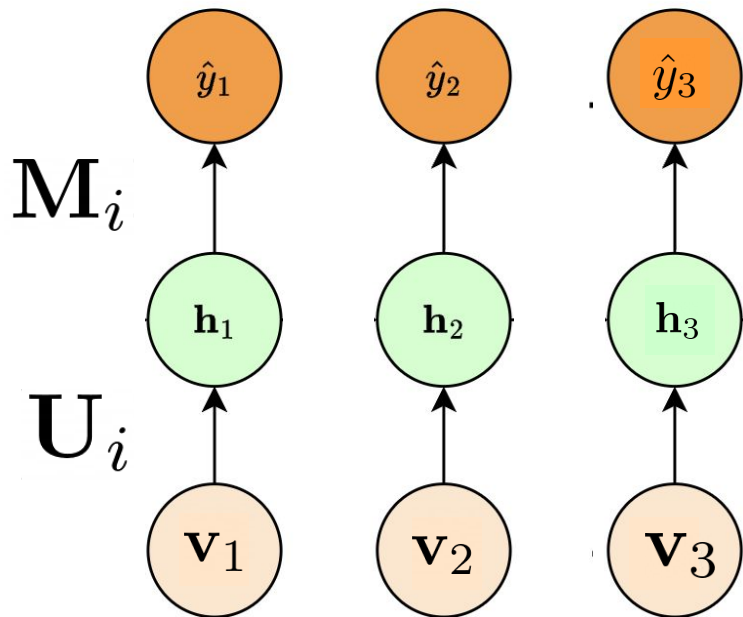
# Let's simplify!

What if we have a single word and a single output?





# Towards RNNs



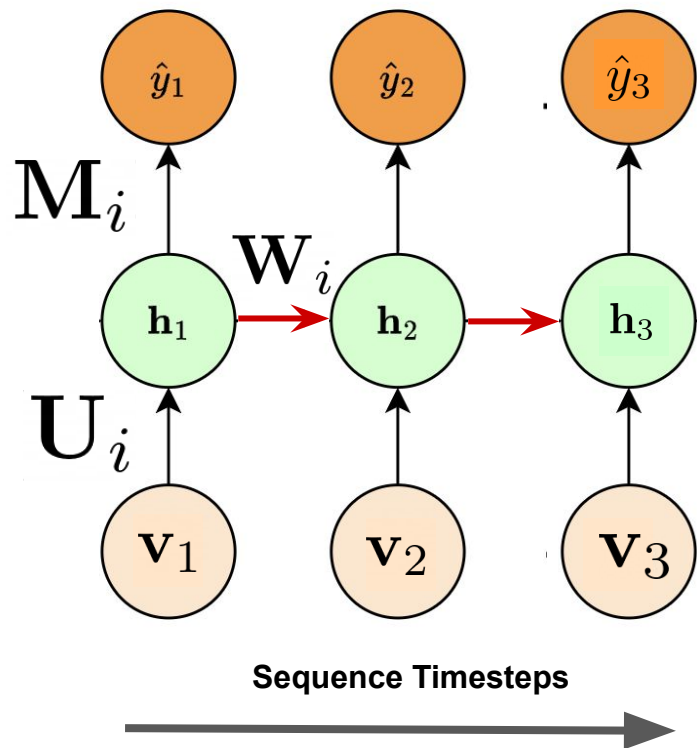
Output

$$\hat{y}_i = \mathbf{M}_i \mathbf{h}_i$$

Hidden State

$$\mathbf{h}_i = \sigma(\mathbf{U}_i \mathbf{v}_i)$$

## Towards RNNs



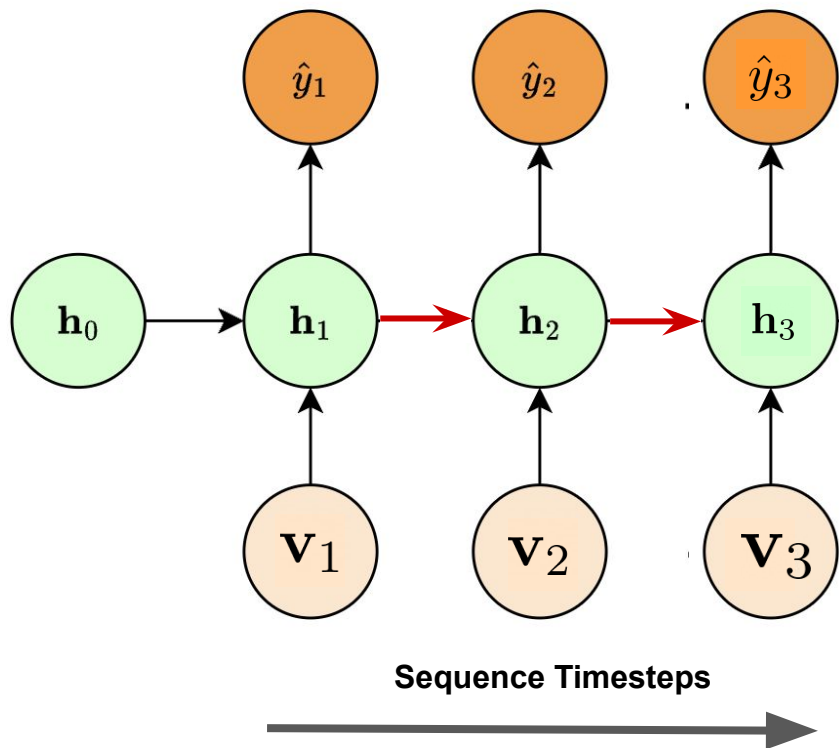
Output

$$\hat{y}_i = \mathbf{M}_i \mathbf{h}_i$$

Hidden State

$$\mathbf{h}_i = \sigma(\mathbf{U}_i \mathbf{v}_i + \mathbf{W}_i \mathbf{h}_{i-1})$$

## Towards RNNs



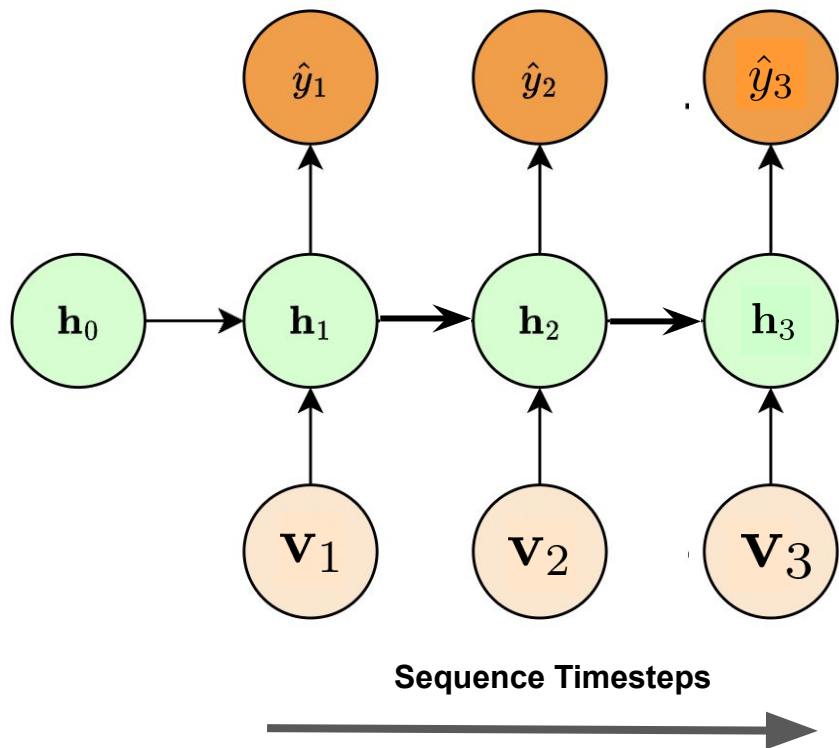
Output

$$\hat{y}_i = \mathbf{M}_i \mathbf{h}_i$$

Hidden State

$$\mathbf{h}_i = \sigma(\mathbf{U}_i \mathbf{v}_i + \mathbf{W}_i \mathbf{h}_{i-1})$$

What's the issue with this setup?



- Too many parameters if we have a long sequence!
- Longer sequence parameters will not receive many updates

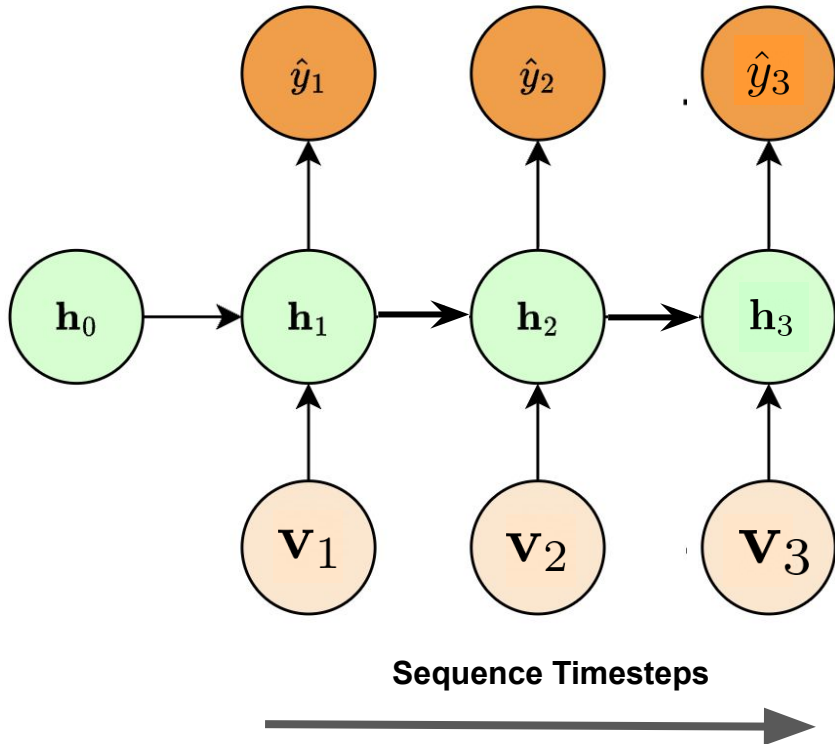
Output

$$\hat{y}_i = \mathbf{M}_i \mathbf{h}_i$$

Hidden State

$$\mathbf{h}_i = \sigma(\mathbf{U}_i \mathbf{v}_i + \mathbf{W}_i \mathbf{h}_{i-1})$$

# Recurrent neural network (RNN)



Use the same parameters across different timesteps.

Output

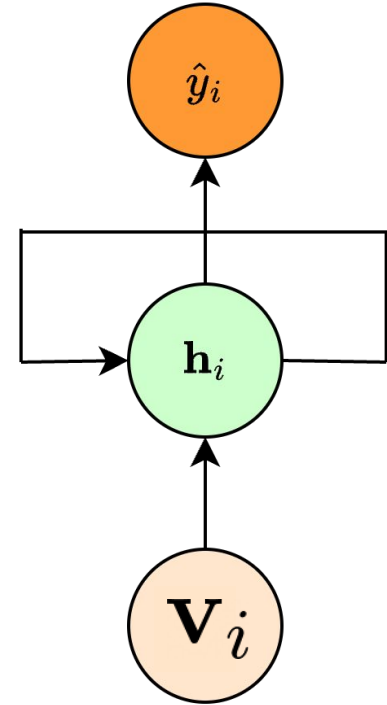
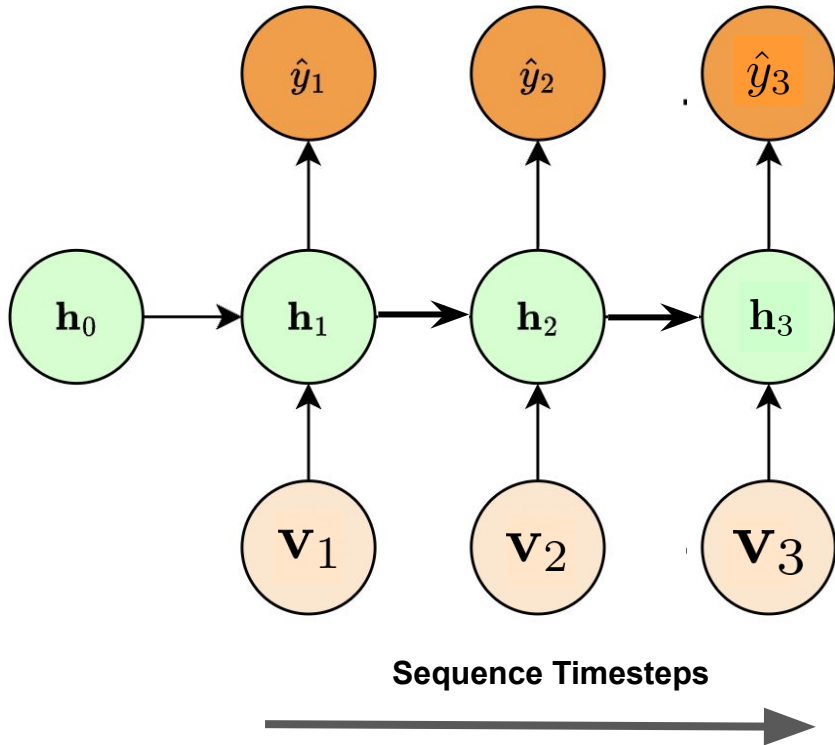
$$\hat{y}_i = \mathbf{M} \mathbf{h}_i$$

Hidden State

$$\mathbf{h}_i = \sigma(\mathbf{U} \mathbf{v}_i + \mathbf{W} \mathbf{h}_{i-1})$$

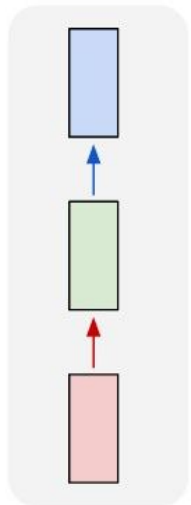
# Recurrent neural network (RNN)

Use the same parameters across different timesteps.

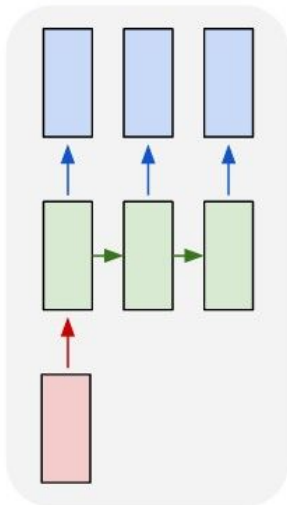


# Discuss: Which tasks can you perform with RNNs? Can you find an example of each task?

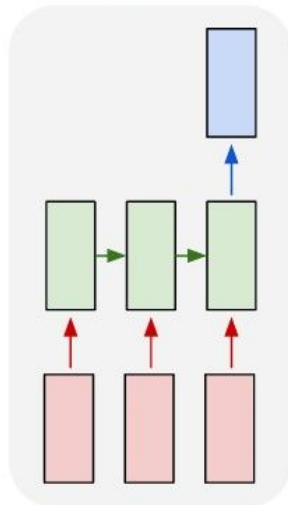
one to one



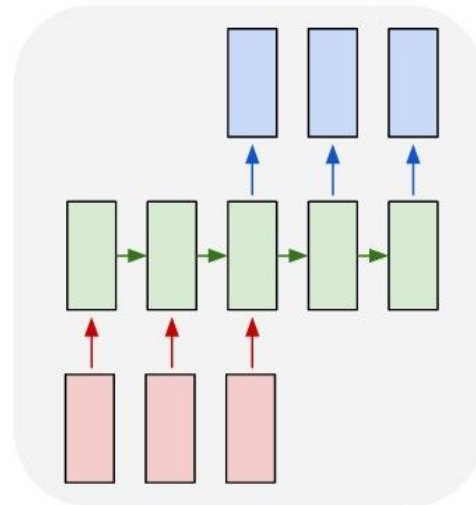
one to many



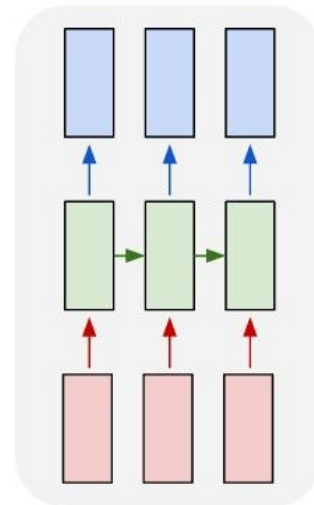
many to one



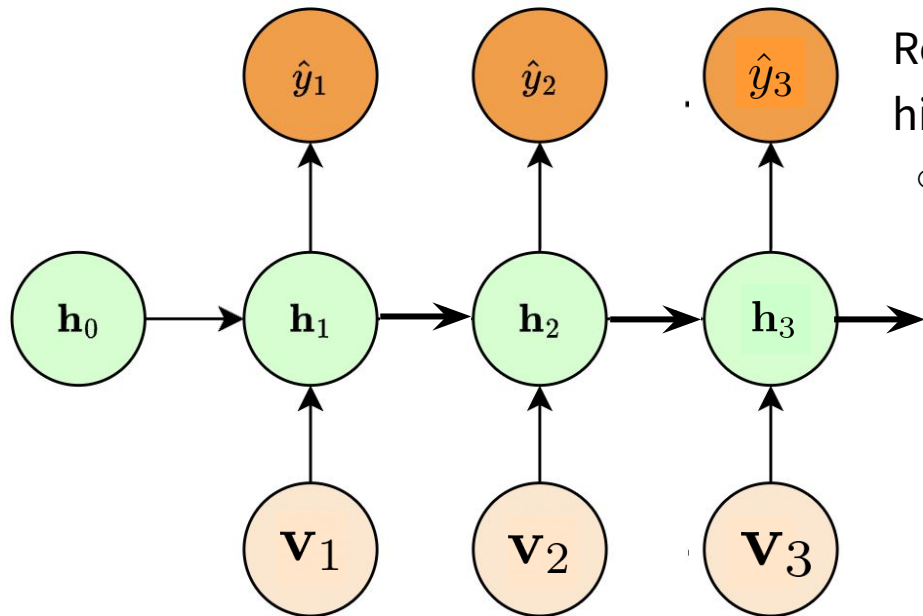
many to many



many to many

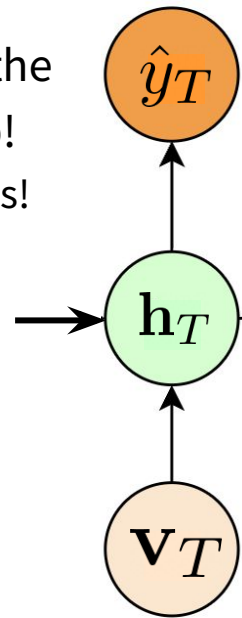


# RNN: Issues under **Looooooooong** Context



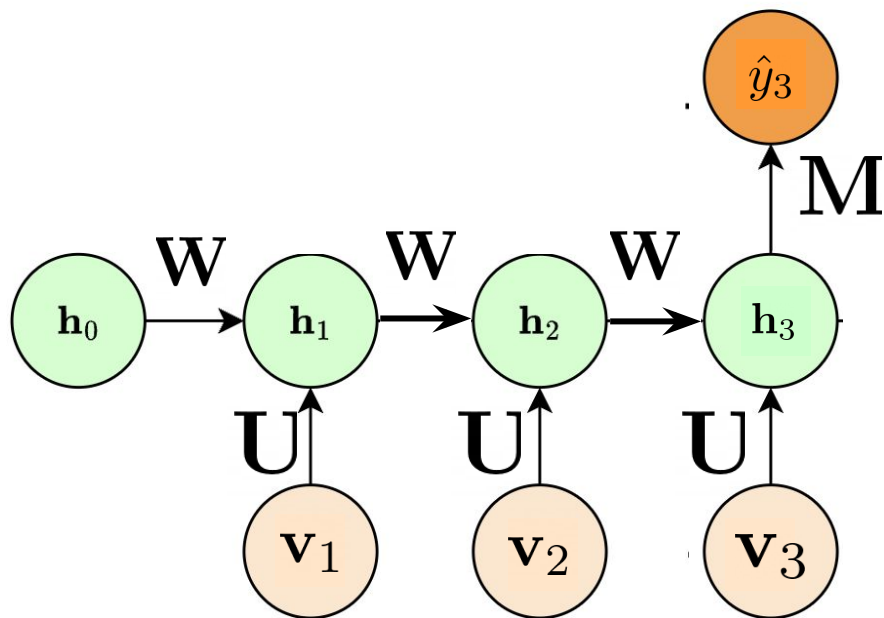
Recurrent forward will **rewrite** the hidden states on every timestep!

- What will happen? Let's discuss!

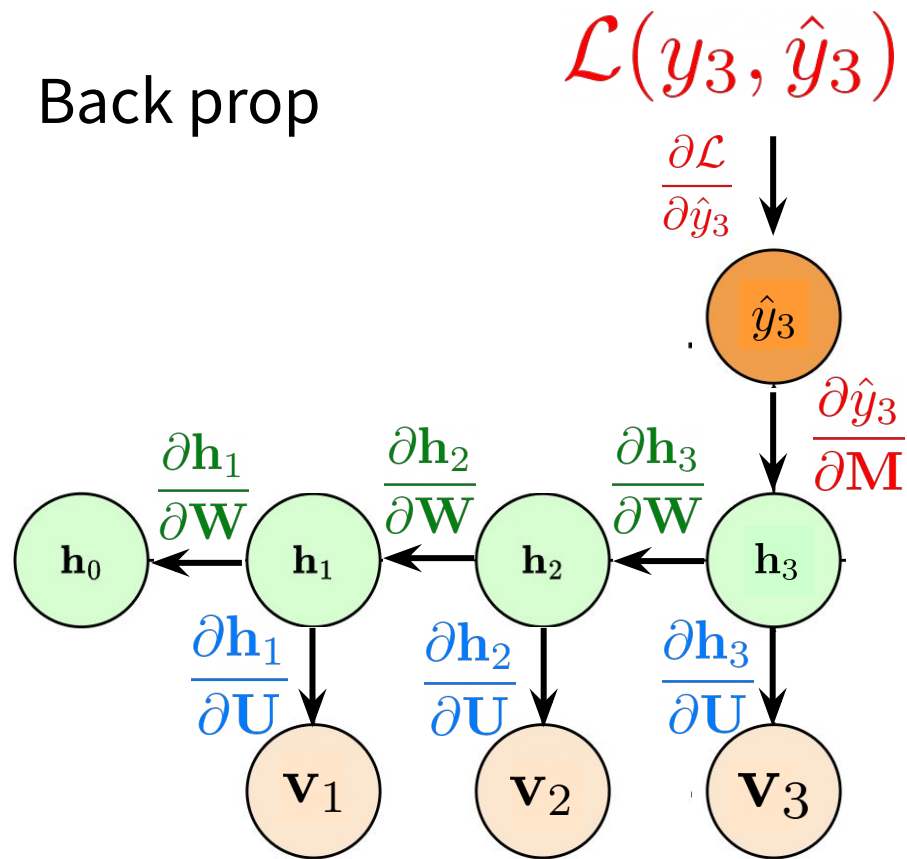




Forward prop

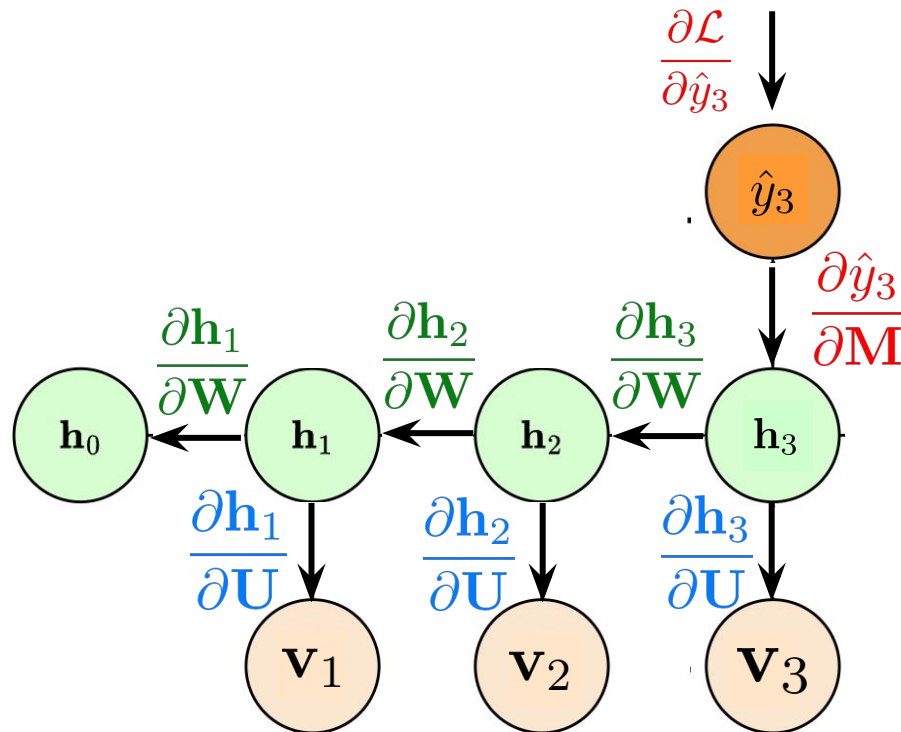


Back prop



## Back prop through time

$$\mathcal{L}(y_3, \hat{y}_3)$$



- Unfold a recurrent neural network in time
- Gradients are accumulated across all time steps by applying the chain rule
- Propagate gradients backwards through time steps

# Back prop through time

$$\frac{\partial \mathcal{L}}{\partial \mathbf{M}} = \frac{\partial \mathcal{L}}{\partial \hat{y}_3} \cdot \frac{\partial \hat{y}_3}{\partial \mathbf{M}}$$

Gradients wrt  $\mathbf{W}$  from last time step:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \frac{\partial \mathcal{L}}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial \mathbf{h}_3} \frac{\partial \mathbf{h}_3}{\partial \mathbf{W}}$$

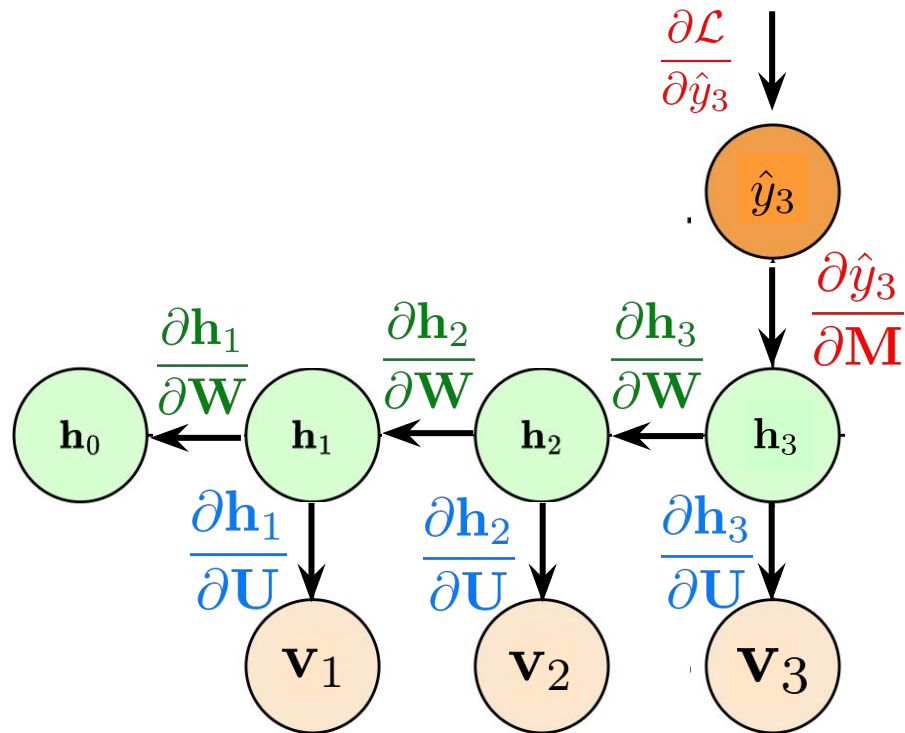
Gradients wrt  $\mathbf{W}$  from time step 2:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \frac{\partial \mathcal{L}}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial \mathbf{h}_3} \frac{\partial \mathbf{h}_3}{\partial \mathbf{h}_2} \frac{\partial \mathbf{h}_2}{\partial \mathbf{W}}$$

Gradients wrt  $\mathbf{W}$  from time step 1:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \frac{\partial \mathcal{L}}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial \mathbf{h}_3} \frac{\partial \mathbf{h}_3}{\partial \mathbf{h}_2} \frac{\partial \mathbf{h}_2}{\partial \mathbf{h}_1} \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}}$$

$$\mathcal{L}(y_3, \hat{y}_3)$$



$T = 3$ 

What is the general form of  $\frac{\partial \mathcal{L}}{\partial \mathbf{W}}$  with  $T = T$ , at time step  $t$ ?

Gradients wrt  $\mathbf{W}$  from time step 3:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \frac{\partial \mathcal{L}}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial \mathbf{h}_3} \frac{\partial \mathbf{h}_3}{\partial \mathbf{W}}$$

Gradients wrt  $\mathbf{W}$  from time step 2:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \frac{\partial \mathcal{L}}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial \mathbf{h}_3} \frac{\partial \mathbf{h}_3}{\partial \mathbf{h}_2} \frac{\partial \mathbf{h}_2}{\partial \mathbf{W}}$$

Gradients wrt  $\mathbf{W}$  from time step 1:

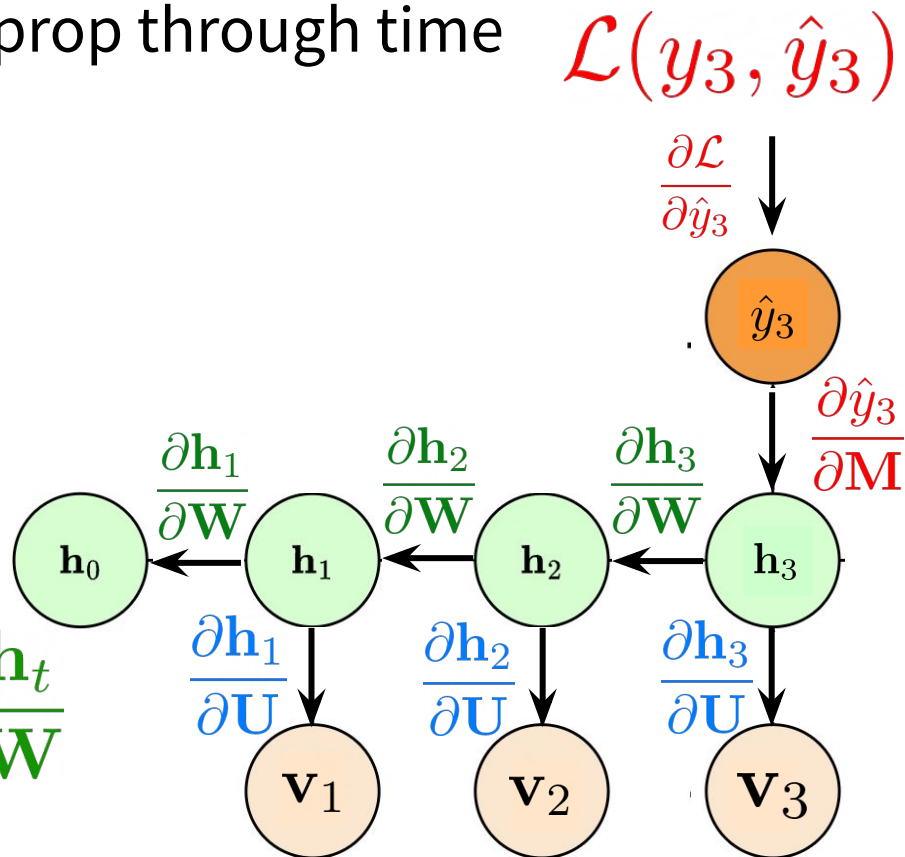
$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \frac{\partial \mathcal{L}}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial \mathbf{h}_3} \frac{\partial \mathbf{h}_3}{\partial \mathbf{h}_2} \frac{\partial \mathbf{h}_2}{\partial \mathbf{h}_1} \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}}$$

## Back prop through time

Gradients wrt W from time step t:

Each timestamp contributes to the gradient!  
Summing over all timestamps:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \sum_{t=1}^T \frac{\partial \mathcal{L}}{\partial \hat{y}_T} \frac{\partial \hat{y}_T}{\partial \mathbf{h}_T} \left( \prod_{i=t}^{T-1} \frac{\partial \mathbf{h}_{i+1}}{\partial \mathbf{h}_i} \right) \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}}$$



RNN: Issues under **Loooooong** Context

$$\frac{\partial \mathcal{L}(\hat{y}_T)}{\partial \mathbf{h}_1} = \frac{\partial \mathcal{L}(\hat{y}_T)}{\partial \mathbf{h}_T} \prod_{1 < t \leq T} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}}$$

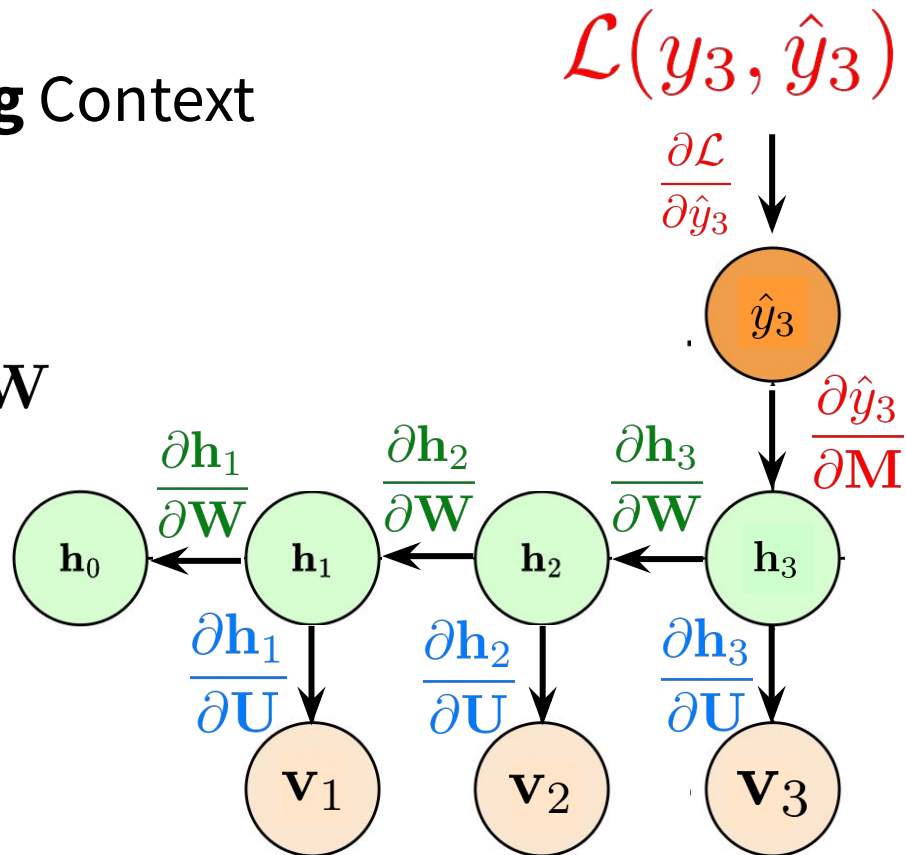
$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} = \text{diag}(\sigma'(\mathbf{W}\mathbf{h}_{t-1} + \mathbf{U}\mathbf{v}_t))\mathbf{W}$$

- **Vanishing gradients: grad to 0**

If  $\left\| \frac{\partial \mathbf{h}_i}{\partial \mathbf{h}_{i-1}} \right\| < 1$  and  $T$  is large,  $\left\| \frac{\partial \mathcal{L}(\hat{y}_T)}{\partial \mathbf{h}_1} \right\| \rightarrow 0$ .

- **Exploding gradients: grad to *inf***

If  $\left\| \frac{\partial \mathbf{h}_i}{\partial \mathbf{h}_{i-1}} \right\| > 1$  and  $T$  is large,  $\left\| \frac{\partial \mathcal{L}(\hat{y}_T)}{\partial \mathbf{h}_1} \right\| \rightarrow \text{inf}$ .

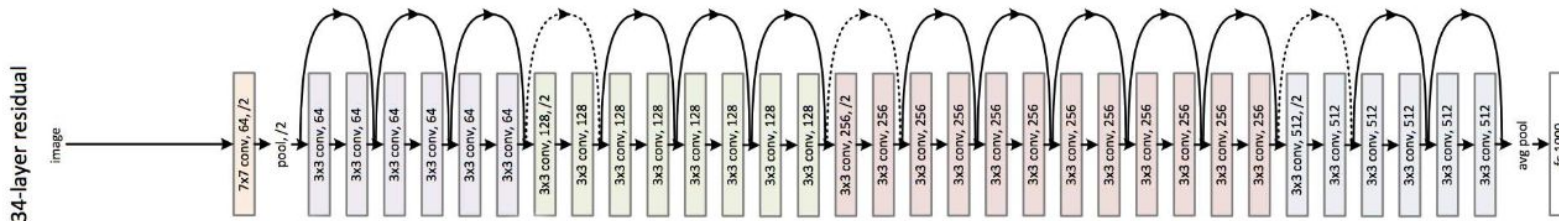


# Recall: ResNet

## “Plain” Network

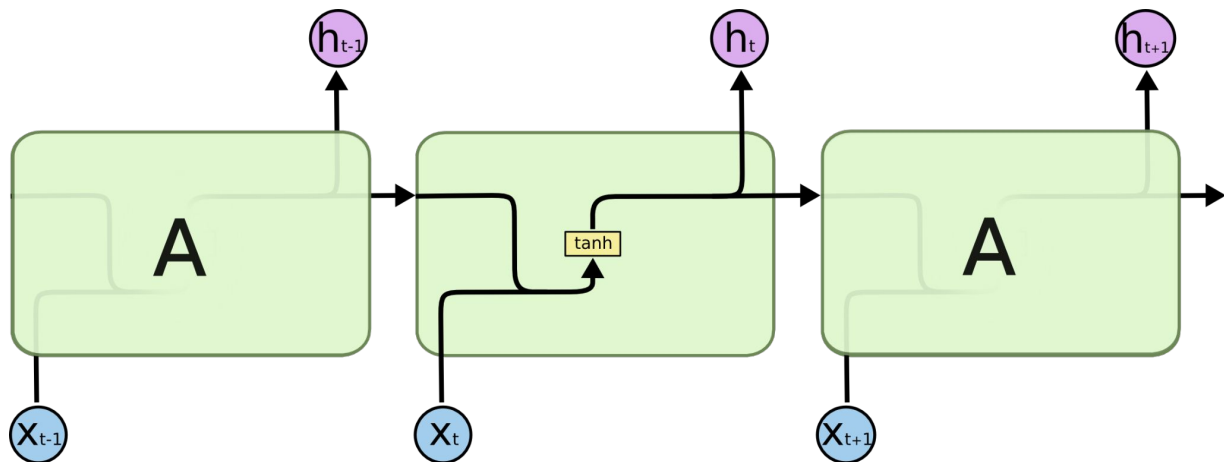


## ResNet

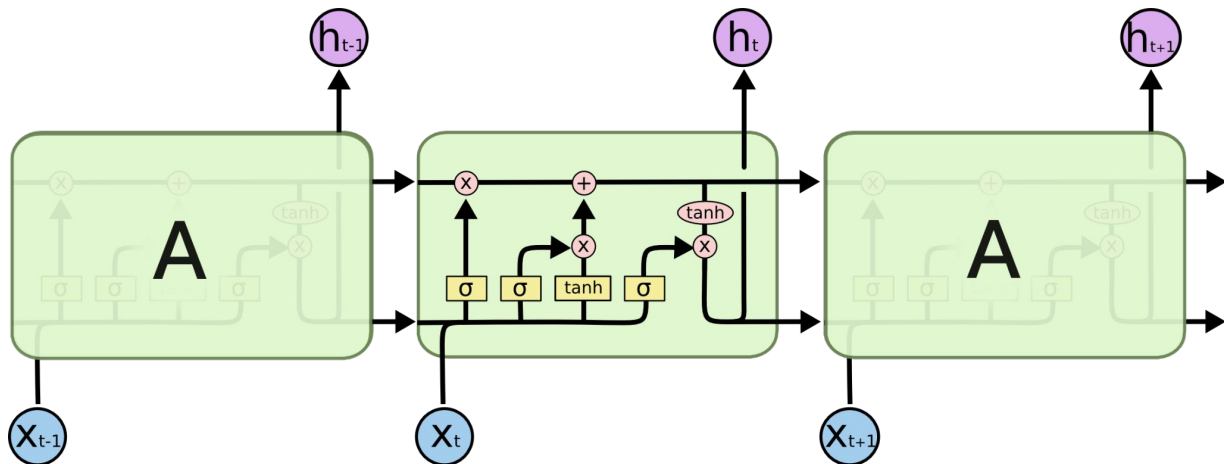


[He, Kaiming, et al. "Deep residual learning for image recognition." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.]

RNN



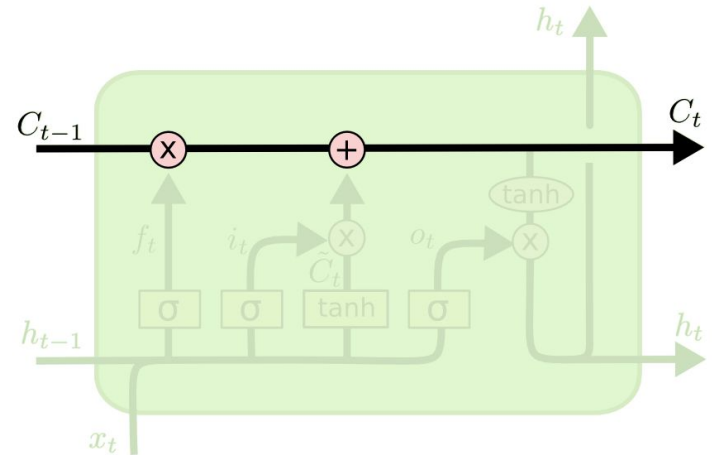
Long-short Term  
Memory (LSTM)





# Long-short Term Memory (LSTM)

- Main idea: add a “cell” state that allows information to flow easily
  - Similar to residual connections
  - No repeated matrix multiplications!

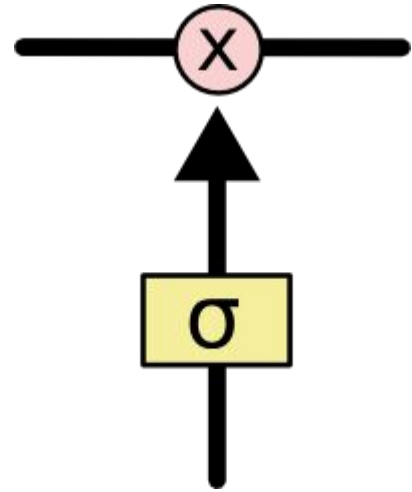
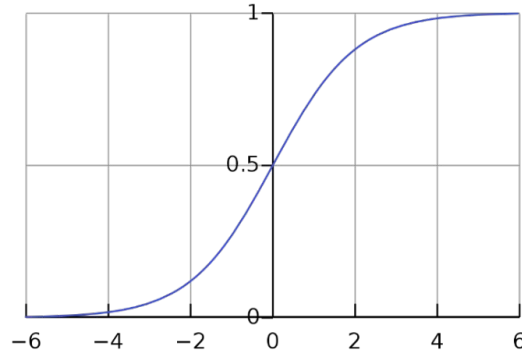


# LSTMs- Gates

- Control the flow of information with “gates”
  - Element-wise product with the output of a sigmoid activation

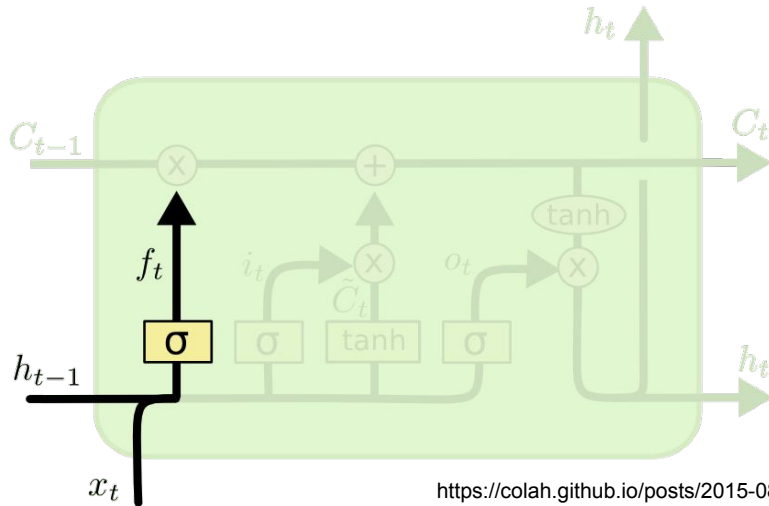
Sigmoid  
Function

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



# LSTMs- Forget Gate

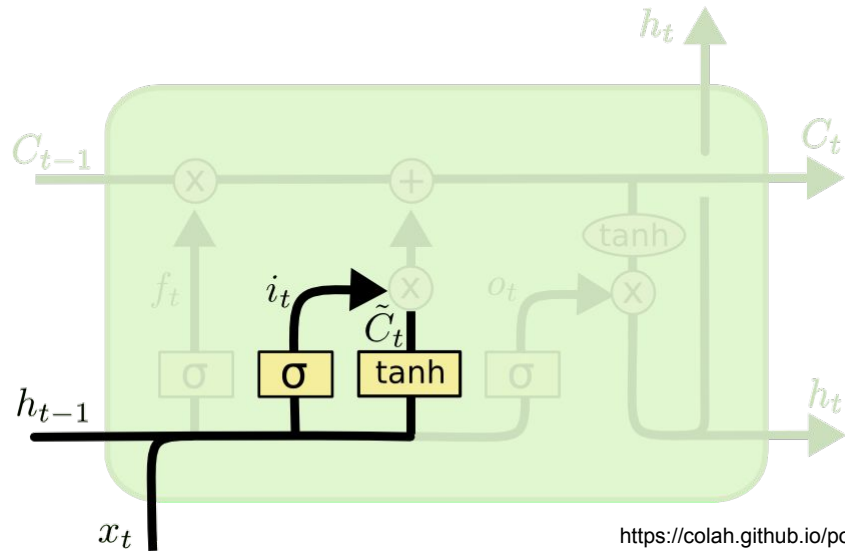
- Forget gate- function of current input and previous hidden state
- Controls what should be remembered in the cell state



$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

# LSTMs- Input Gate

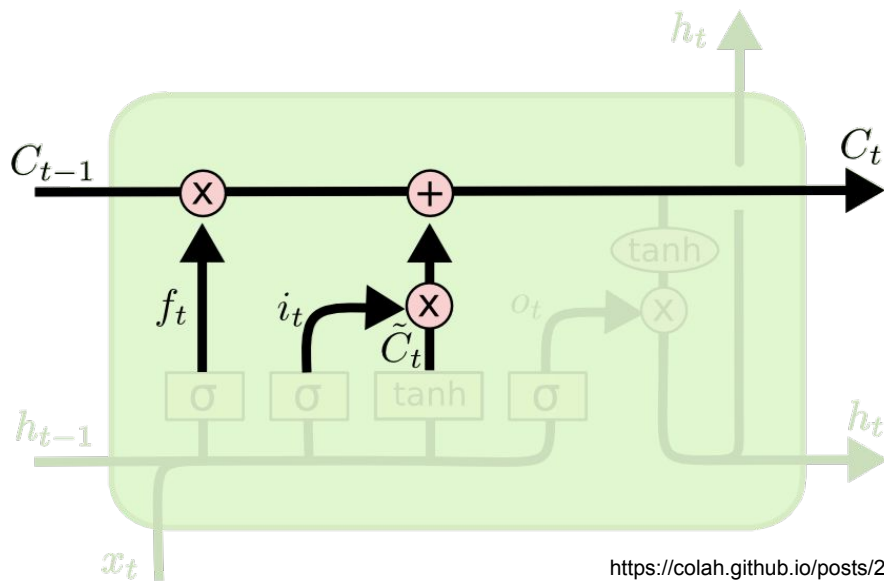
- Input gate- function of current input and previous hidden state
- Decides what information to write to the cell state



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

# LSTM- Cell Update

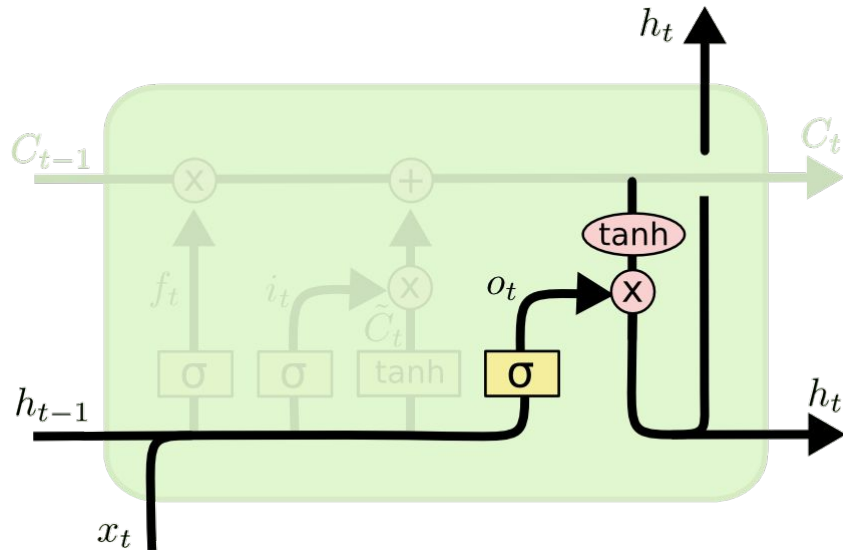
- Forget irrelevant information
- Add new information from the current token



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

# LSTM- Output Gate

- Output gate- function of current input and previous hidden state
- Controls flow of information from the cell state to the hidden state
- Given some weight matrix  $W_o$ , how do we write to  $o_t$  and  $h_t$ ?



# LSTMs

- Add a cell state to store information
  - Gradient flows along the cell state
- Update cell state with parameterized gating functions
- Performs better with long sequences

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

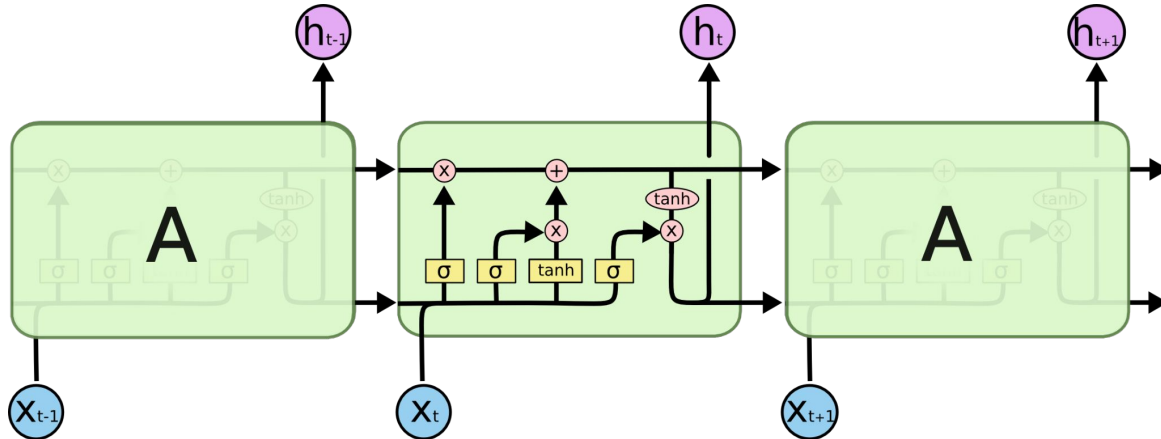
$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

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

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

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

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



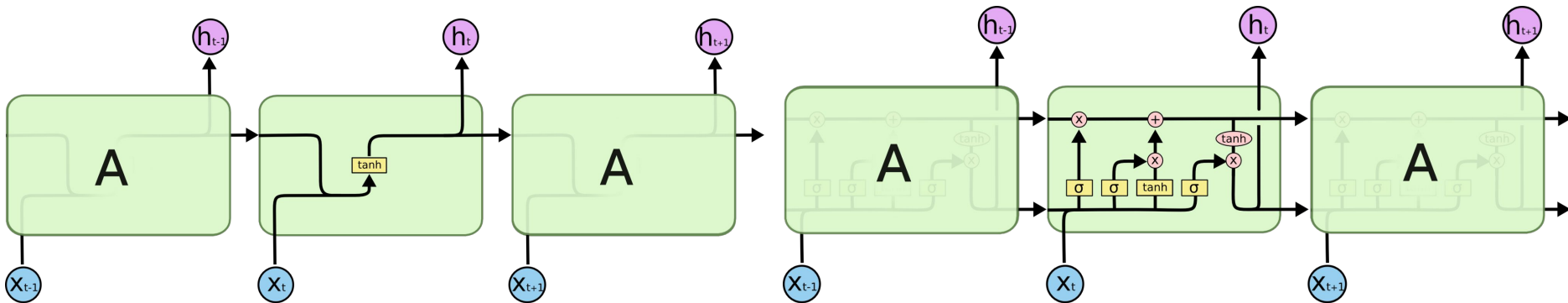
# RNN vs. LSTM

- **RNN**

- Can be applied to variable-length sequences
- Share parameters across time
- Hard to train!

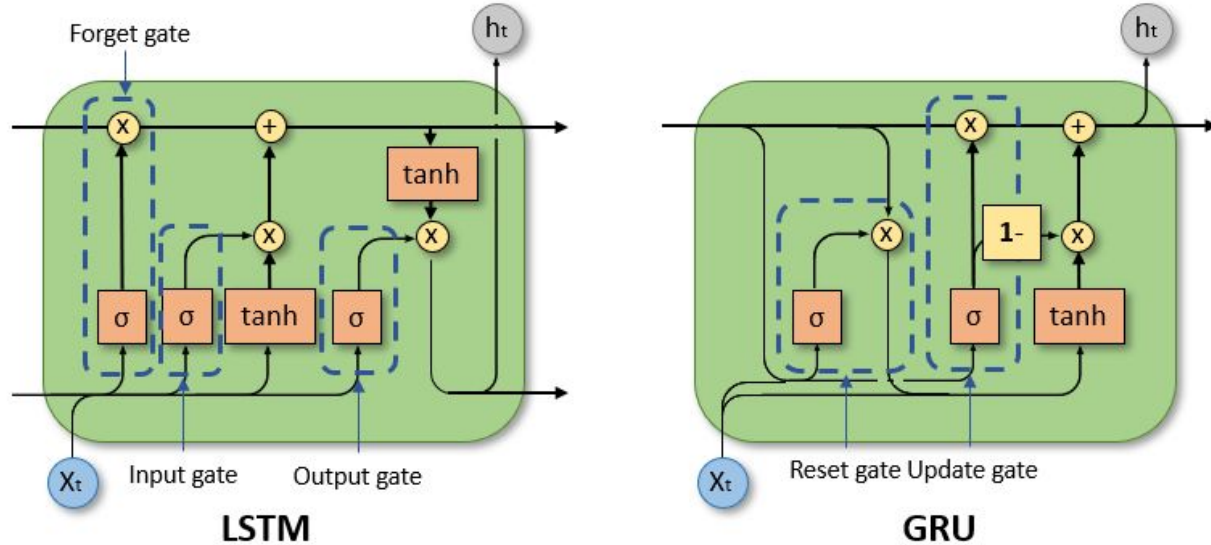
- **LSTM**

- Mitigates the vanishing gradient problem with the cell state
- Better for long sequences



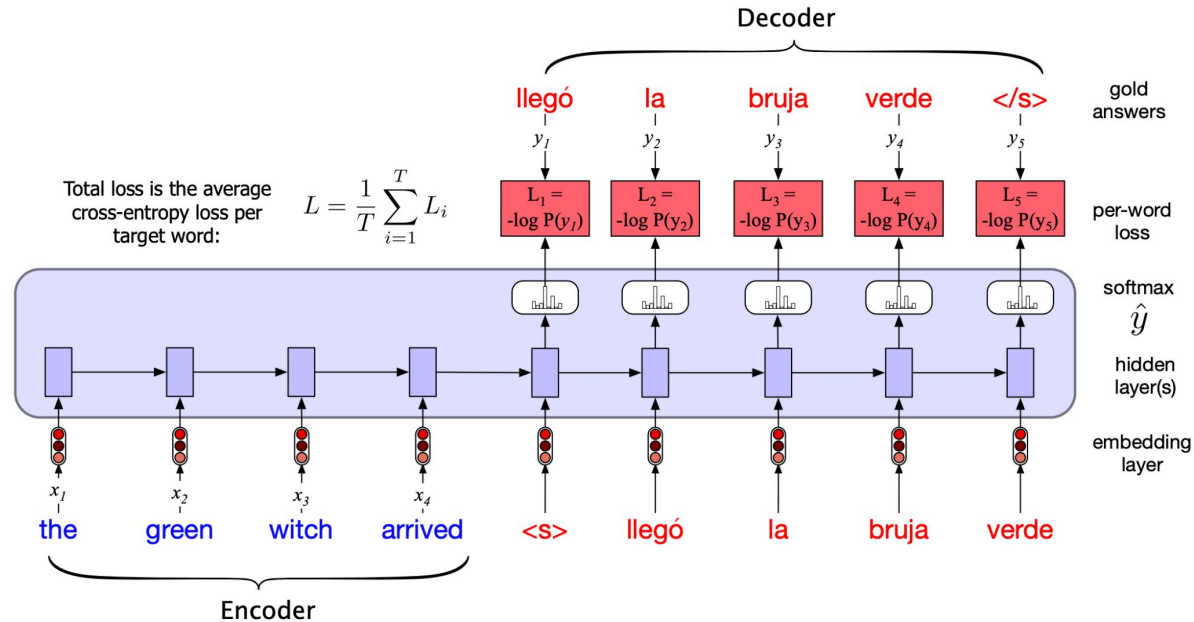


# Gated recurrent units (GRUs)



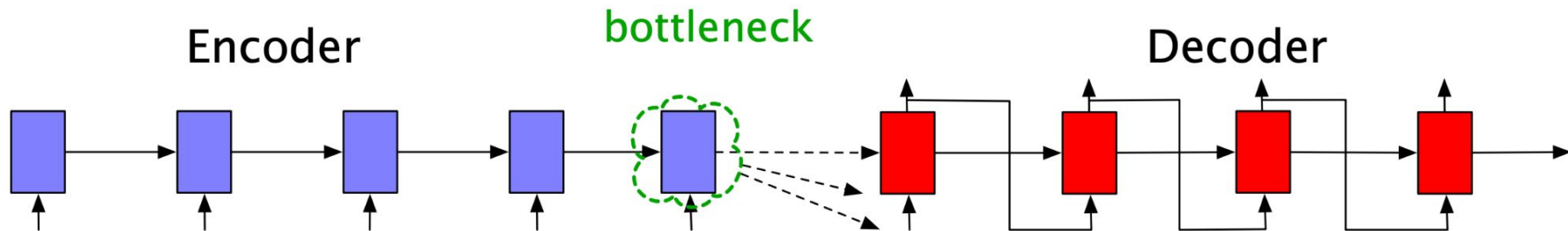
# Sequence-to-Sequence Generation

- Map some input sequence to a target sequence
- Applications:
  - Machine translation
  - News summarization
  - ChatGPT!



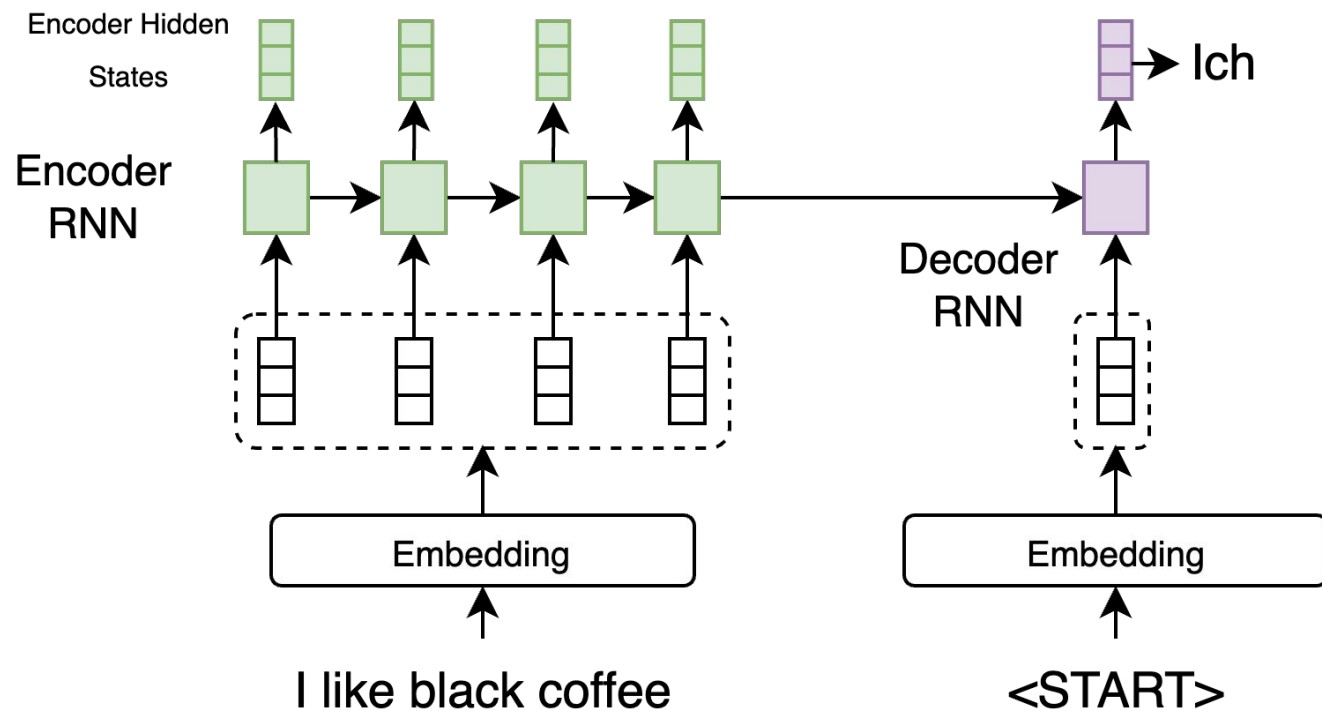
# Bottleneck Problem

- All the information about the source sequence must be stored in a single vector
  - How to translate a long paragraph?
  - How to summarize long articles?



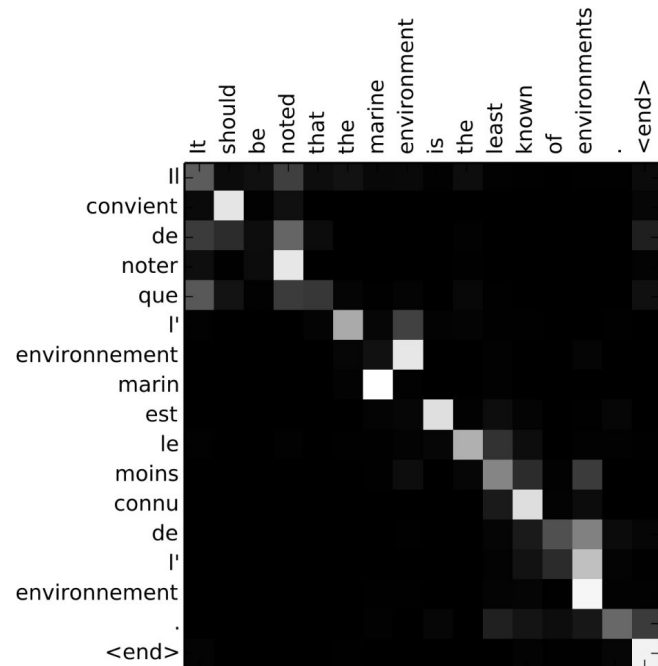
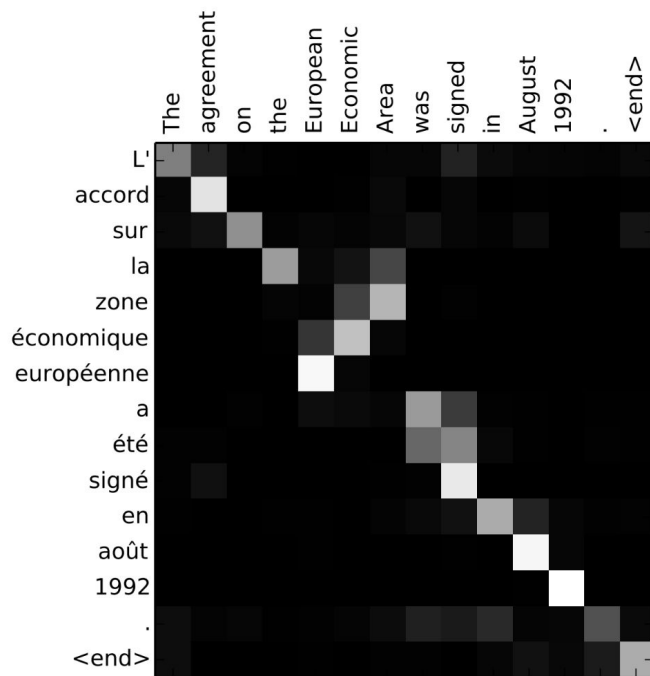
# RNN for Machine Translation

Would be nice if we could “look back” at previous hidden states



# Visualizing Attention

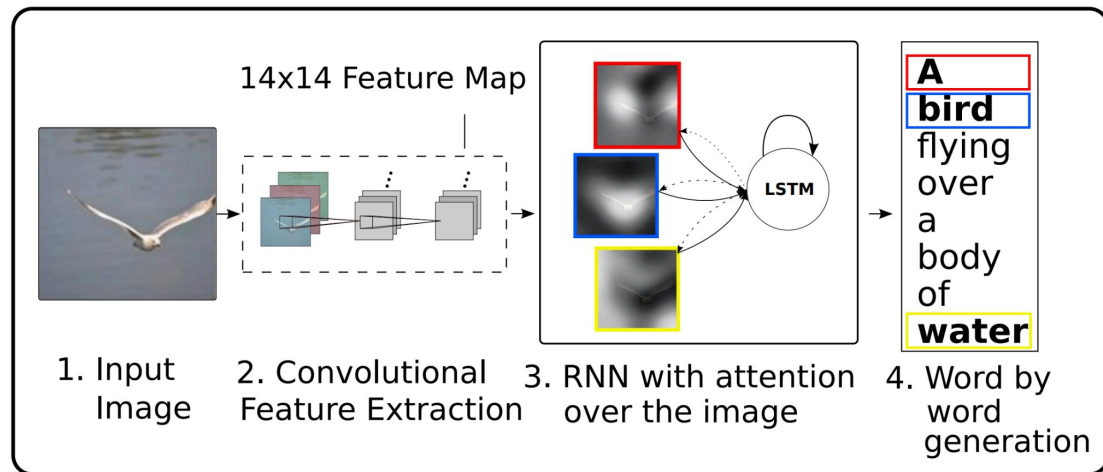
- Plot attention weights to see where the model is “looking”
  - Learns language alignment for translation!



# Attention Application- Image Captioning!

- Extract image features with a CNN
- Use an LSTM with attention to generate image captions

*Figure 1.* Our model learns a words/image alignment. The visualized attentional maps (3) are explained in section 3.1 & 5.4



# Visualize Attention Weights

- Learns to focus on relevant regions of the image

Figure 3. Examples of attending to the correct object (*white* indicates the attended regions, *underlines* indicated the corresponding word)



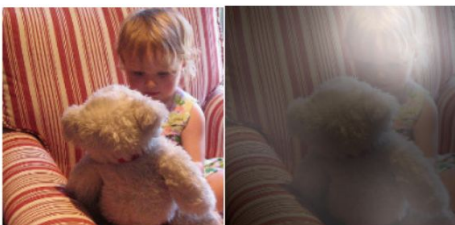
A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



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



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

# Recap

- RNNs can be applied to arbitrary length sequences
  - Run into vanishing/exploding gradient problems
- LSTMs add a cell state to RNNs to improve gradient flow
  - Better a handling long sequences
- Attention can look back at past feature vectors!
  - Scales better to long sequences
  - Can incorporate image features
  - Many, many more applications!