



Cornell Bowers C-IS

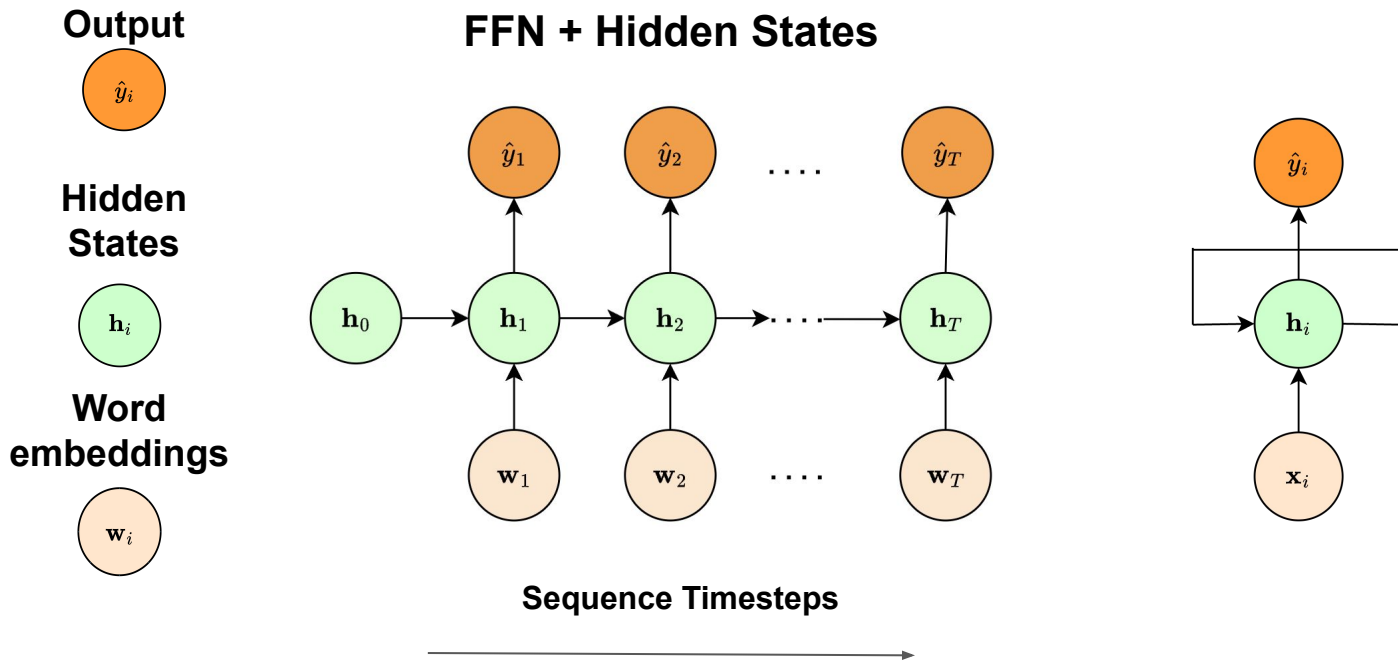
College of Computing and Information Science

Deep Learning

Week 02: LSTMs/Attention

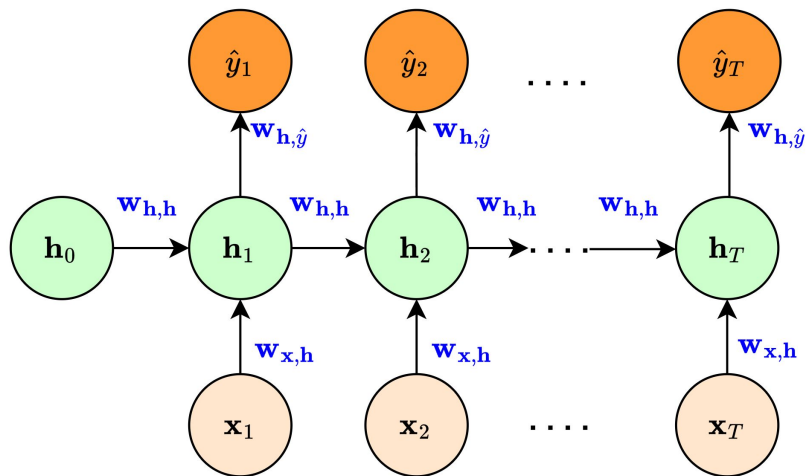
Big Question: How to model sequences of words?

Recurrent neural network (RNN)



RNN w/ parameter-sharing

Use the same parameters across different timesteps.



Hidden State

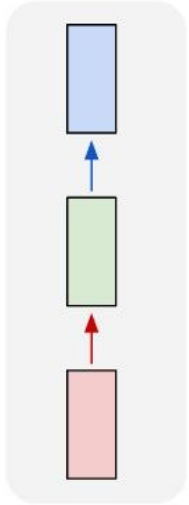
$$\begin{aligned} \mathbf{h}_i &= f(\mathbf{x}_i, \mathbf{h}_{i-1}) \\ &= \sigma(\mathbf{w}_{\mathbf{x},h} \mathbf{x}_i + \mathbf{w}_{\mathbf{h},h} \mathbf{h}_{i-1}) \end{aligned}$$

Output

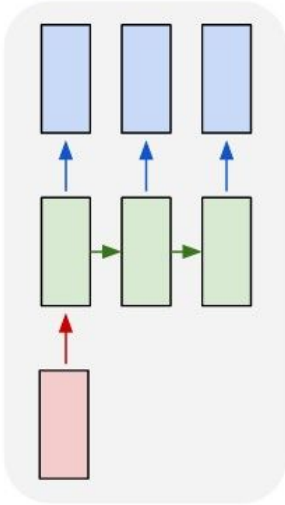
$$\hat{y}_i = \mathbf{w}_{\mathbf{h},\hat{y}} \mathbf{h}_i$$

Discuss: What types of tasks can you perform with RNNs?

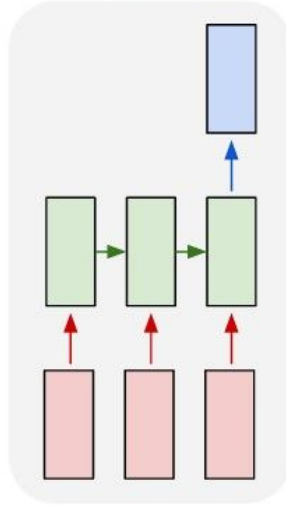
one to one



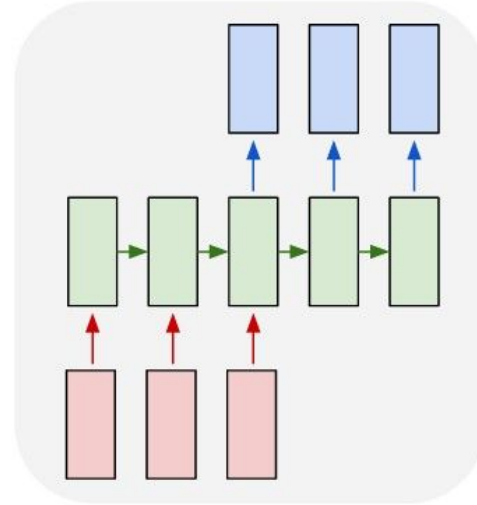
one to many



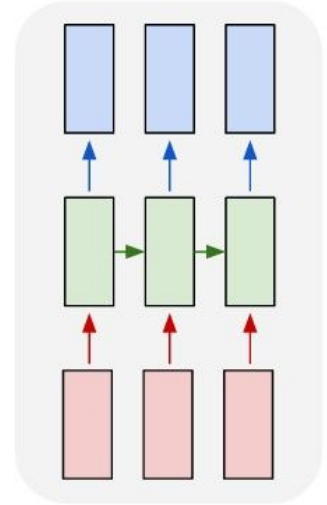
many to one



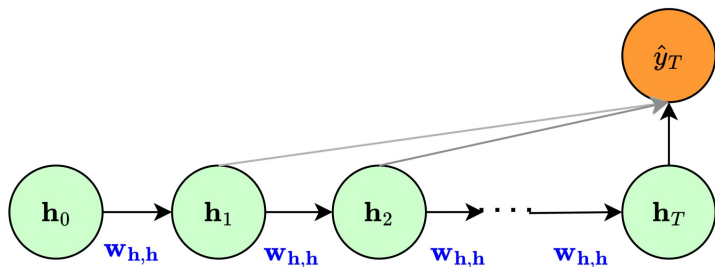
many to many



many to many



RNN: Issues under **Looooooooong** Context

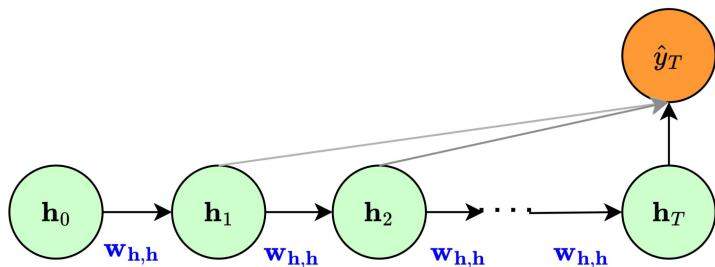


Sequence Timesteps

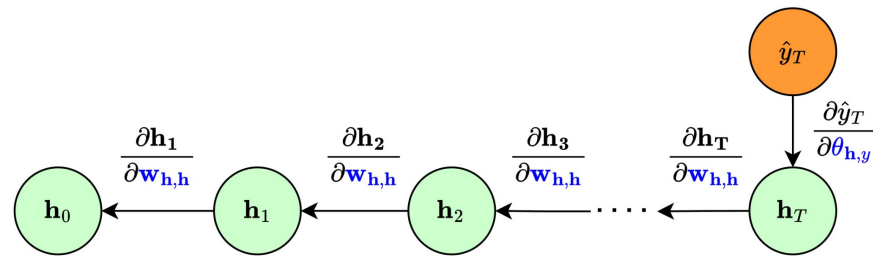


- Recurrent forward will **rewrite** the hidden states on every timestep!
 - What will happen? Let's discuss!

RNN: Issues under **Looooooooong** Context



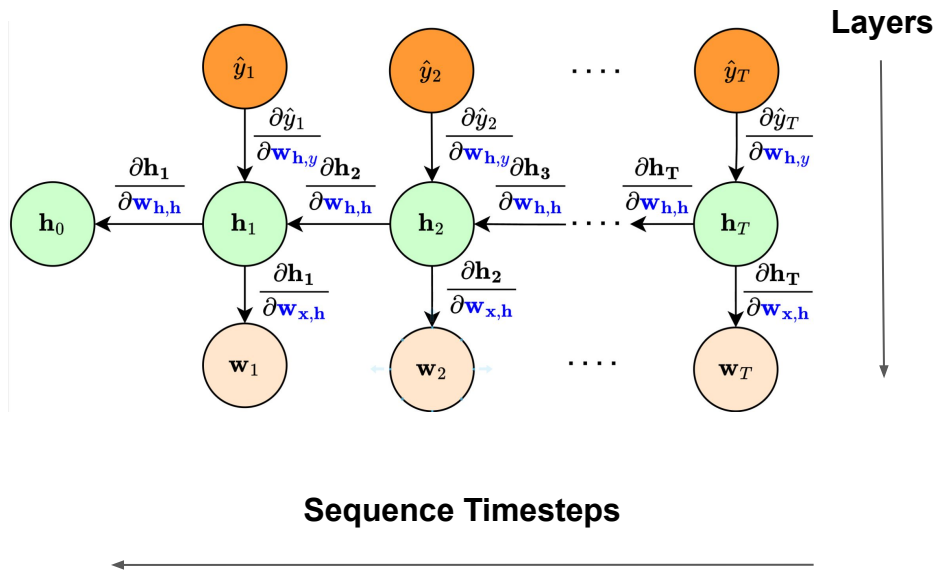
Sequence Timesteps



Sequence Timesteps



Backpropagation through the Time (BPTT)



- 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

Backpropagation through the Time (BPTT)

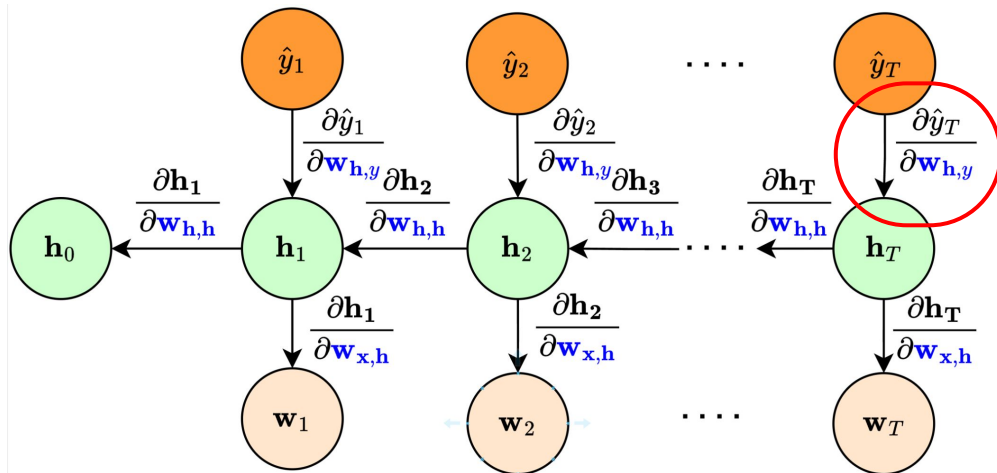
Hidden State

$$\begin{aligned} \mathbf{h}_i &= f(\mathbf{x}_i, \mathbf{h}_{i-1}) \\ &= \sigma(\mathbf{w}_{\mathbf{x},h}\mathbf{x} + \mathbf{w}_{\mathbf{h},h}\mathbf{h}_{i-1}) \end{aligned}$$

Output

$$\hat{y}_i = \mathbf{w}_{\mathbf{h},\hat{y}}\mathbf{h}_i$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}_{\mathbf{h},\hat{y}}} = \frac{\partial \mathcal{L}}{\partial \hat{y}_T} \cdot \frac{\partial \hat{y}_T}{\partial \mathbf{w}_{\mathbf{h},\hat{y}}}$$



Backpropagation through the Time (BPTT)

Assume we only compute the loss on the last time step

Last time step:

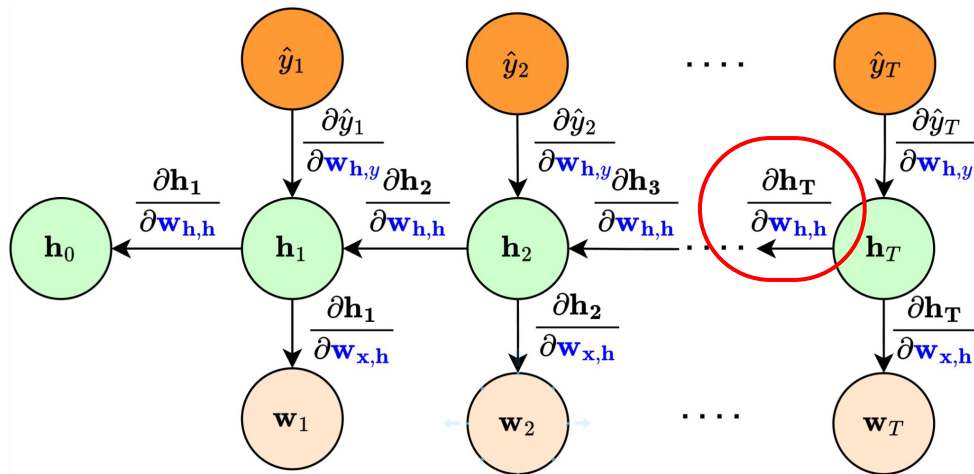
$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}_{\mathbf{h},\mathbf{h}}} = \frac{\partial \mathcal{L}}{\partial \hat{y}_T} \cdot \frac{\partial \hat{y}_T}{\partial \mathbf{h}_T} \cdot \frac{\partial \mathbf{h}_T}{\partial \mathbf{w}_{\mathbf{h},\mathbf{h}}}$$

Hidden State

$$\begin{aligned} \mathbf{h}_i &= f(\mathbf{x}_i, \mathbf{h}_{i-1}) \\ &= \sigma(\mathbf{w}_{\mathbf{x},\mathbf{h}}\mathbf{x} + \mathbf{w}_{\mathbf{h},\mathbf{h}}\mathbf{h}_{i-1}) \end{aligned}$$

Output

$$\hat{y}_i = \mathbf{w}_{\mathbf{h},\hat{y}}\mathbf{h}_i$$



Backpropagation through the Time (BPTT)

Assume we only compute the loss on the last time step

Last time step:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}_{h,h}} = \frac{\partial \mathcal{L}}{\partial \hat{y}_T} \cdot \frac{\partial \hat{y}_T}{\partial \mathbf{h}_T} \cdot \frac{\partial \mathbf{h}_T}{\partial \mathbf{w}_{h,h}}$$

T-1th time step:

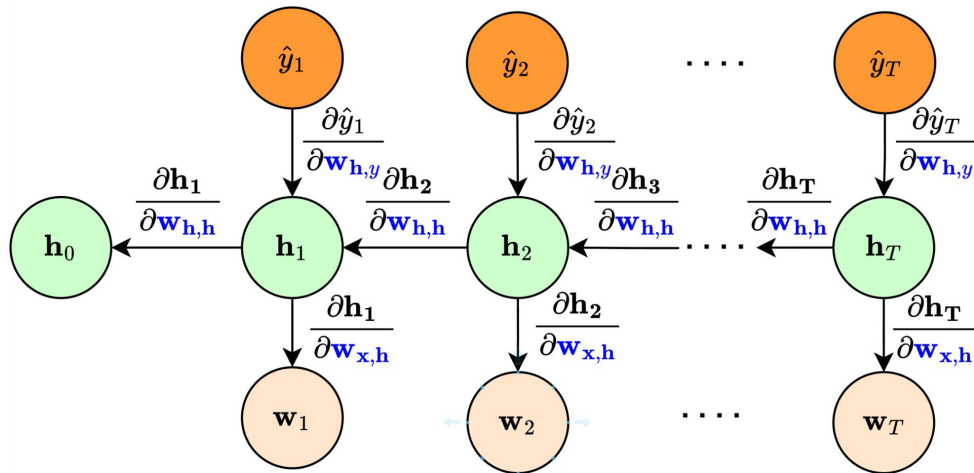
$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}_{h,h}} = \frac{\partial \mathcal{L}}{\partial \hat{y}_T} \cdot \frac{\partial \hat{y}_T}{\partial \mathbf{h}_T} \cdot \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_{T-1}} \cdot \frac{\partial \mathbf{h}_{T-1}}{\partial \mathbf{w}_{h,h}}$$

Hidden State

$$\begin{aligned} \mathbf{h}_i &= f(\mathbf{x}_i, \mathbf{h}_{i-1}) \\ &= \sigma(\mathbf{w}_{x,h} \mathbf{x} + \mathbf{w}_{h,h} \mathbf{h}_{i-1}) \end{aligned}$$

Output

$$\hat{y}_i = \mathbf{w}_{h,\hat{y}} \mathbf{h}_i$$



Backpropagation through the Time (BPTT)

Assume we only compute the loss on the last time step

Last time step:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}_{h,h}} = \frac{\partial \mathcal{L}}{\partial \hat{y}_T} \cdot \frac{\partial \hat{y}_T}{\partial \mathbf{h}_T} \cdot \frac{\partial \mathbf{h}_T}{\partial \mathbf{w}_{h,h}}$$

T-1th time step:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}_{h,h}} = \frac{\partial \mathcal{L}}{\partial \hat{y}_T} \cdot \frac{\partial \hat{y}_T}{\partial \mathbf{h}_T} \cdot \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_{T-1}} \cdot \frac{\partial \mathbf{h}_{T-1}}{\partial \mathbf{w}_{h,h}}$$

Generalizing and summing over all time steps:

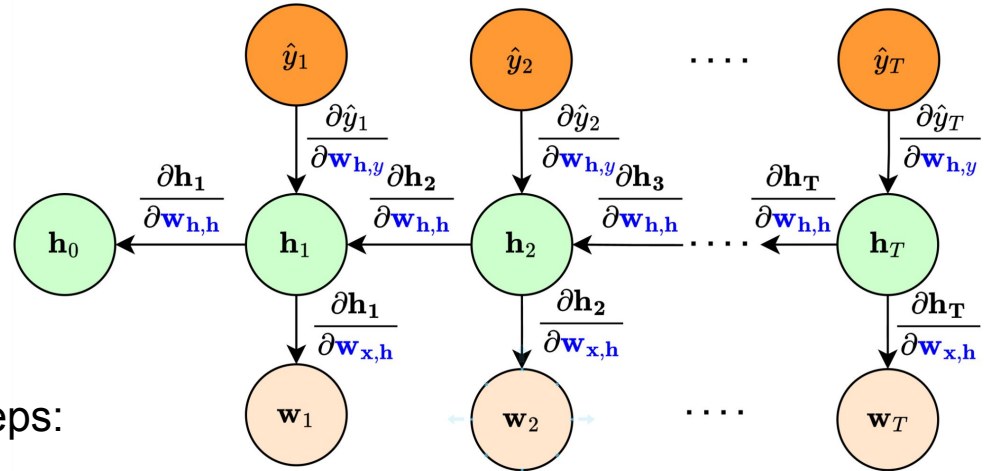
$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{w}_{h,h}} &= \sum_{k=1}^T \frac{\partial \mathcal{L}}{\partial \hat{y}_T} \cdot \frac{\partial \hat{y}_T}{\partial \mathbf{h}_T} \cdot \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_k} \cdot \frac{\partial \mathbf{h}_k}{\partial \mathbf{w}_{h,h}} \\ &= \sum_{k=1}^T \frac{\partial \mathcal{L}}{\partial \hat{y}_T} \cdot \frac{\partial \hat{y}_T}{\partial \mathbf{h}_T} \cdot \left(\prod_{j=k}^{T-1} \frac{\partial \mathbf{h}_{j+1}}{\partial \mathbf{h}_j} \right) \cdot \frac{\partial \mathbf{h}_k}{\partial \mathbf{w}_{h,h}} \end{aligned}$$

Hidden State

$$\begin{aligned} \mathbf{h}_i &= f(\mathbf{x}_i, \mathbf{h}_{i-1}) \\ &= \sigma(\mathbf{w}_{x,h} \mathbf{x} + \mathbf{w}_{h,h} \mathbf{h}_{i-1}) \end{aligned}$$

Output

$$\hat{y}_i = \mathbf{w}_{h,\hat{y}} \mathbf{h}_i$$



RNN: Issues under **Looooooooong** Context

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

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} = \text{diag} \left(\sigma' \left(\mathbf{w}_{\mathbf{h},\mathbf{h}} \mathbf{h}_{t-1} + \mathbf{w}_{\mathbf{h},\mathbf{x}} \mathbf{x}_t \right) \right) \mathbf{w}_{\mathbf{h},\mathbf{h}}$$

- **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***

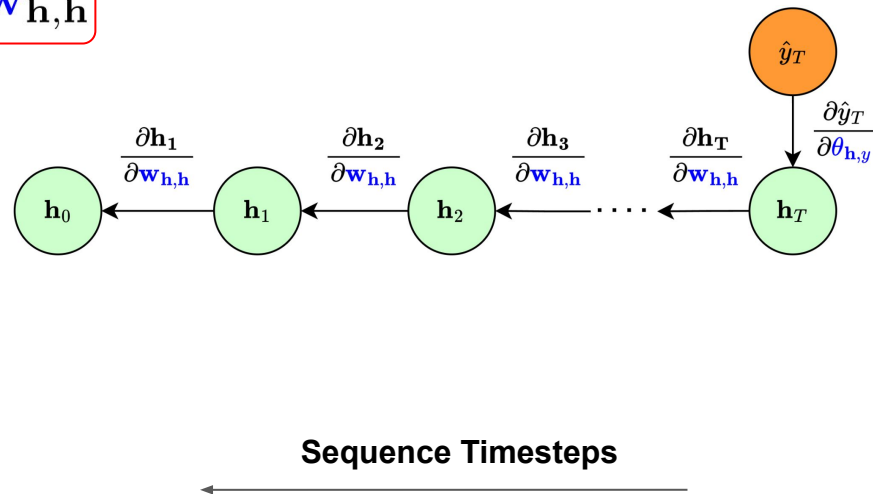
If $\left\| \frac{\partial \mathbf{h}_i}{\partial \mathbf{h}_{i-1}} \right\| > 1$ and T is large, $\left\| \frac{\partial \mathcal{L}(\hat{y}_T)}{\partial \mathbf{h}_1} \right\| \rightarrow \text{inf}$.

Hidden State

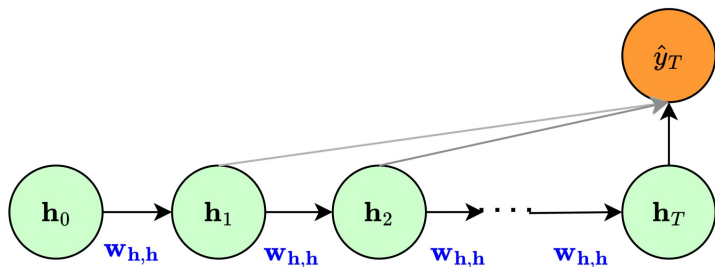
$$\begin{aligned} \mathbf{h}_i &= f(\mathbf{x}_i, \mathbf{h}_{i-1}) \\ &= \sigma(\mathbf{w}_{\mathbf{x},\mathbf{h}} \mathbf{x}_i + \mathbf{w}_{\mathbf{h},\mathbf{h}} \mathbf{h}_{i-1}) \end{aligned}$$

Output

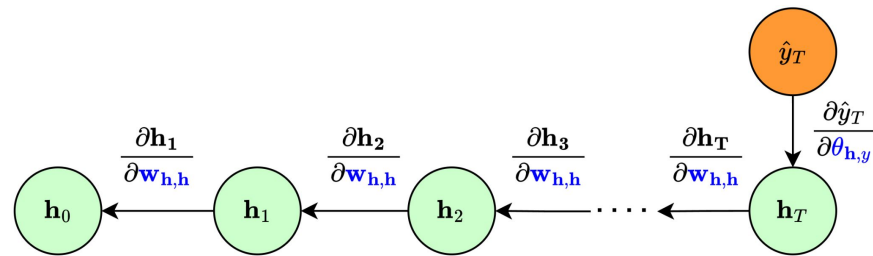
$$\hat{y}_i = \mathbf{w}_{\mathbf{h},\hat{y}} \mathbf{h}_i$$



RNN: Issues under **Looooooooong** Context



Sequence Timesteps

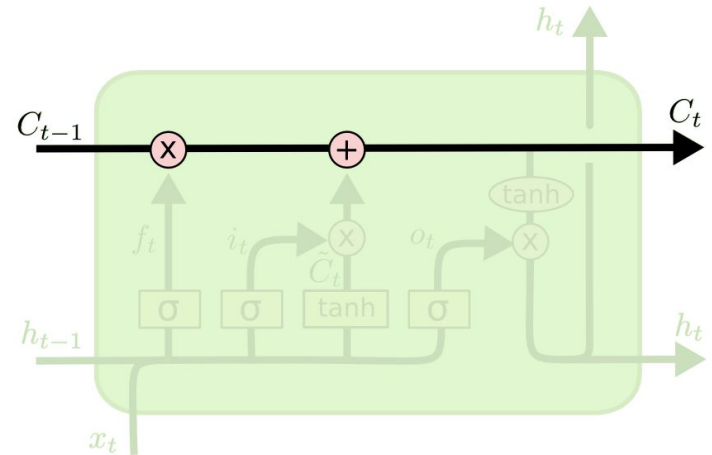


Sequence Timesteps



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!

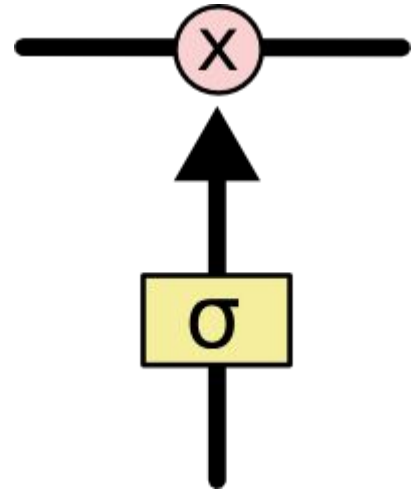
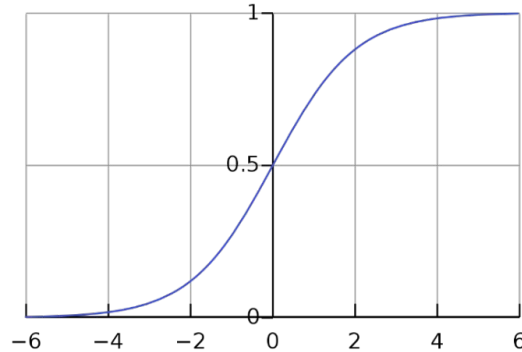


LSTMs- Gates

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

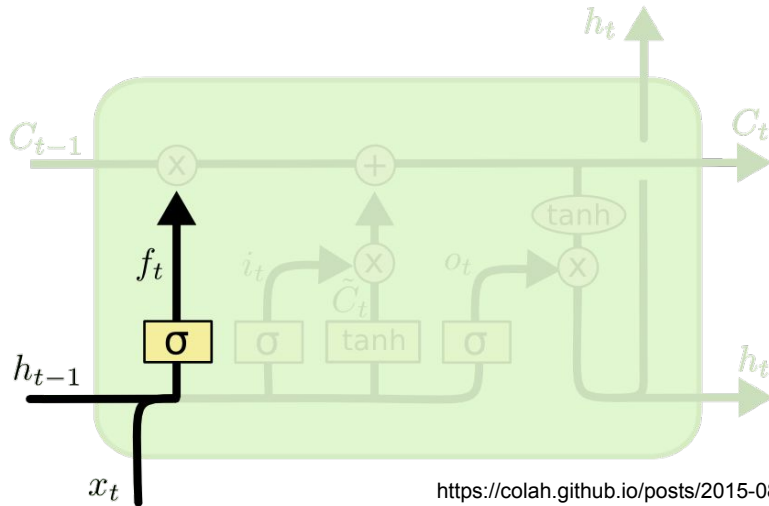
Sigmoid
Function

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



LSTMs- Forget Gate

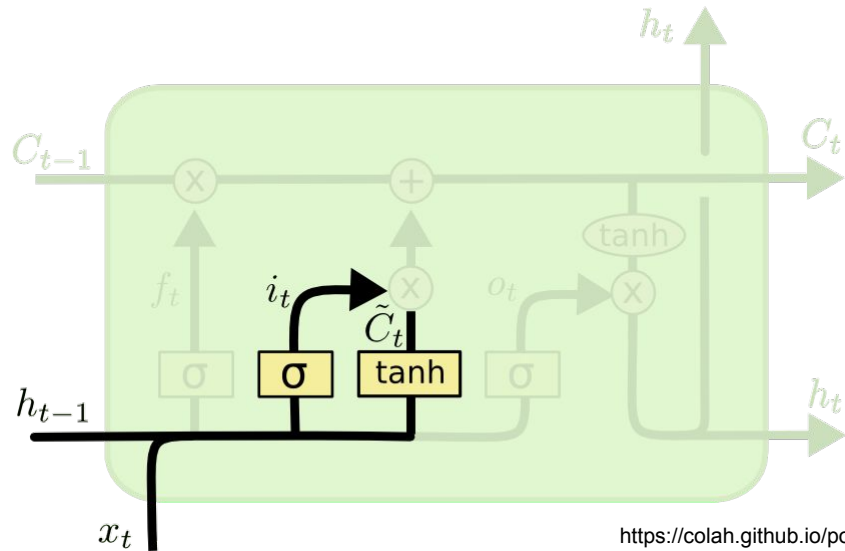
- Forget gate- function of current input and previous hidden state
- Controls what should be remembered in the cell state



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

LSTMs- Input Gate

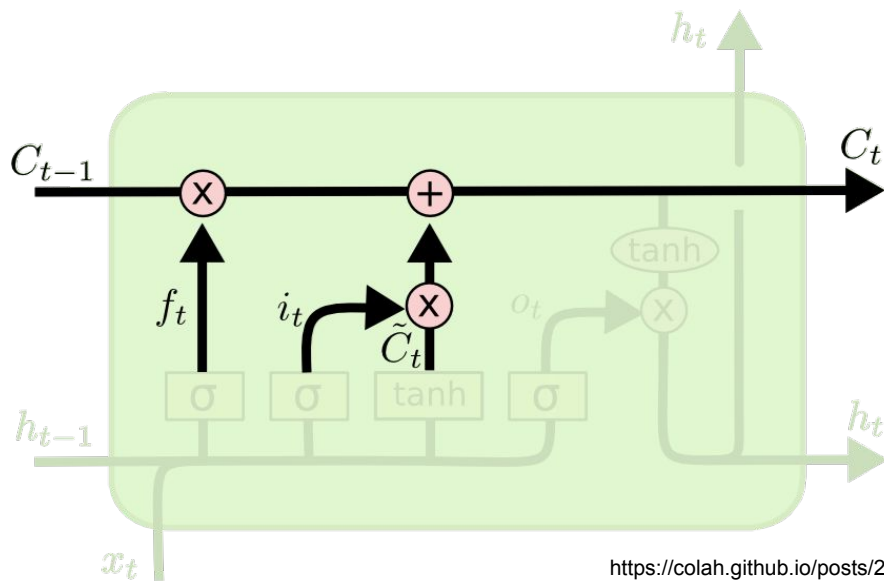
- Input gate- function of current input and previous hidden state
- Decides what information to write to the cell state



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

LSTM- Cell Update

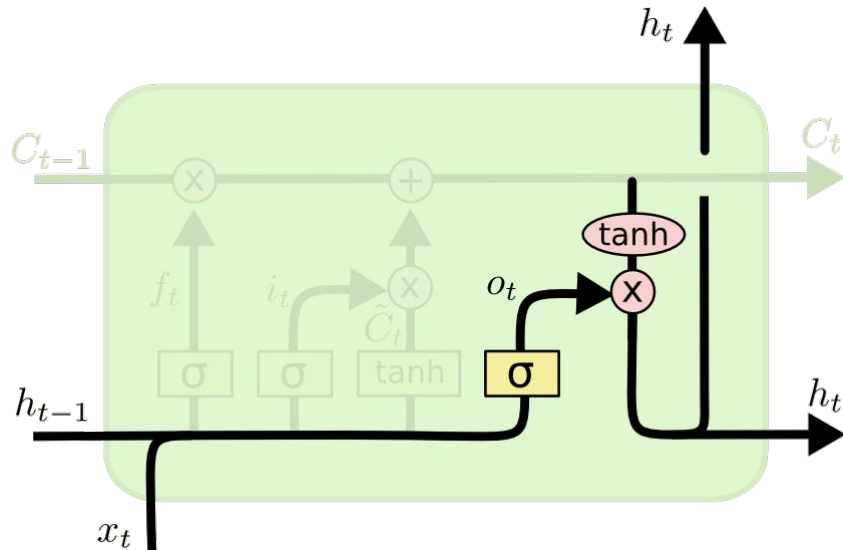
- Forget irrelevant information
- Add new information from the current token



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

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
- Discuss: Given some weight matrix W_o , how could we write the update?



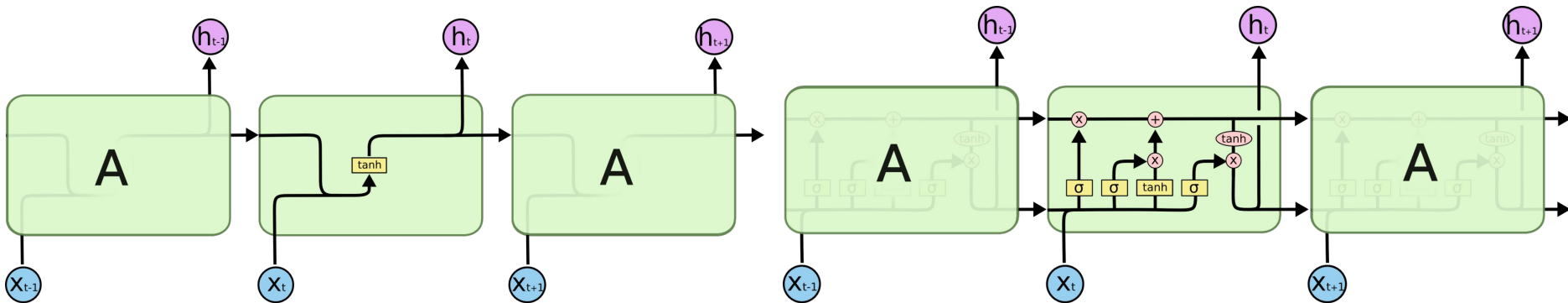
RNN vs. LSTM

- **RNN**

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

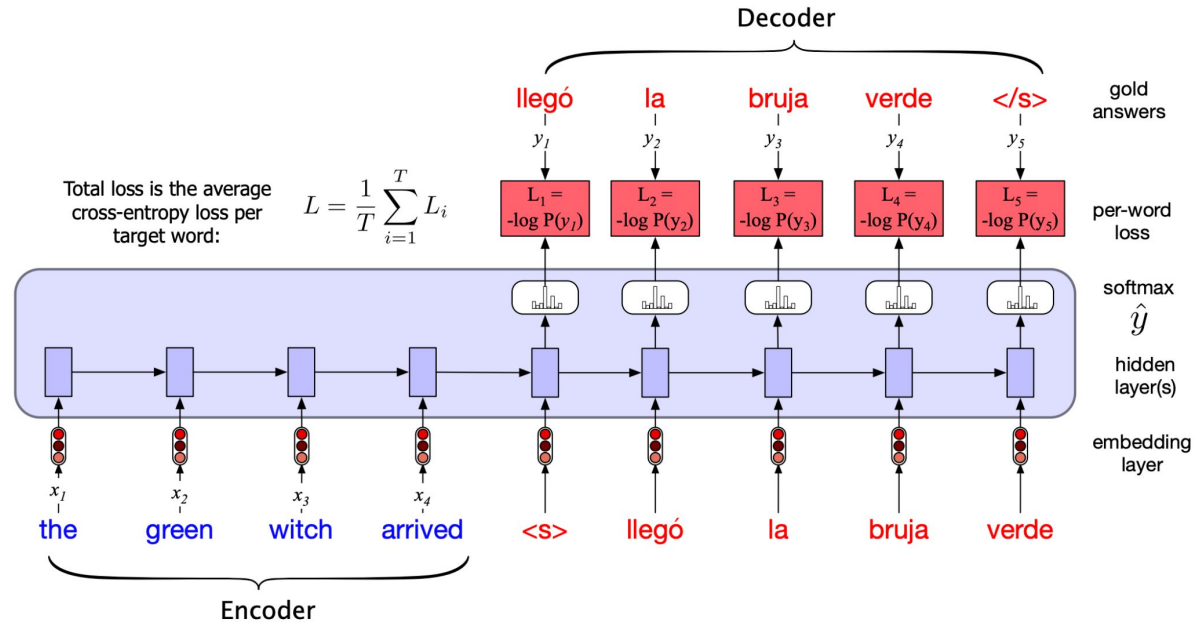
- **LSTM**

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

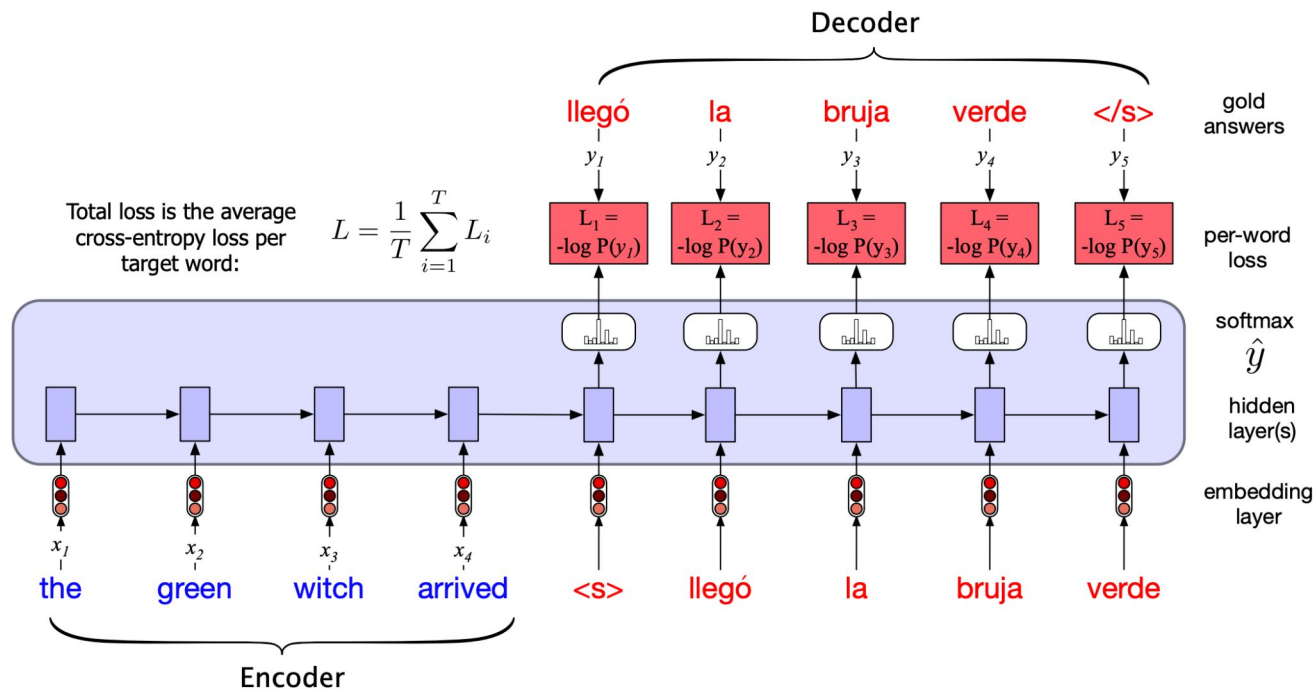


Sequence-to-Sequence Generation

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

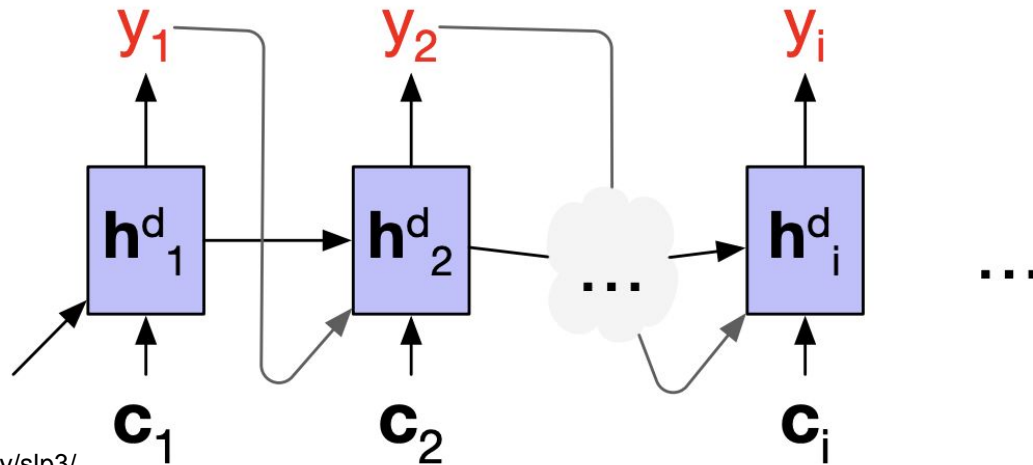


Discuss: Potential problems with Sequence to Sequence Models



Attention

- Attention gives the network a way to “look back” at all previous hidden states
 - Introduced to handle long source sentences in neural machine translation (NMT)



Attention Mechanism

Consists of 3 “general” steps:

Step 1: Compute score of each embedding/input

$$\text{score}(\mathbf{s}_{t-1}, \mathbf{h}_i)$$

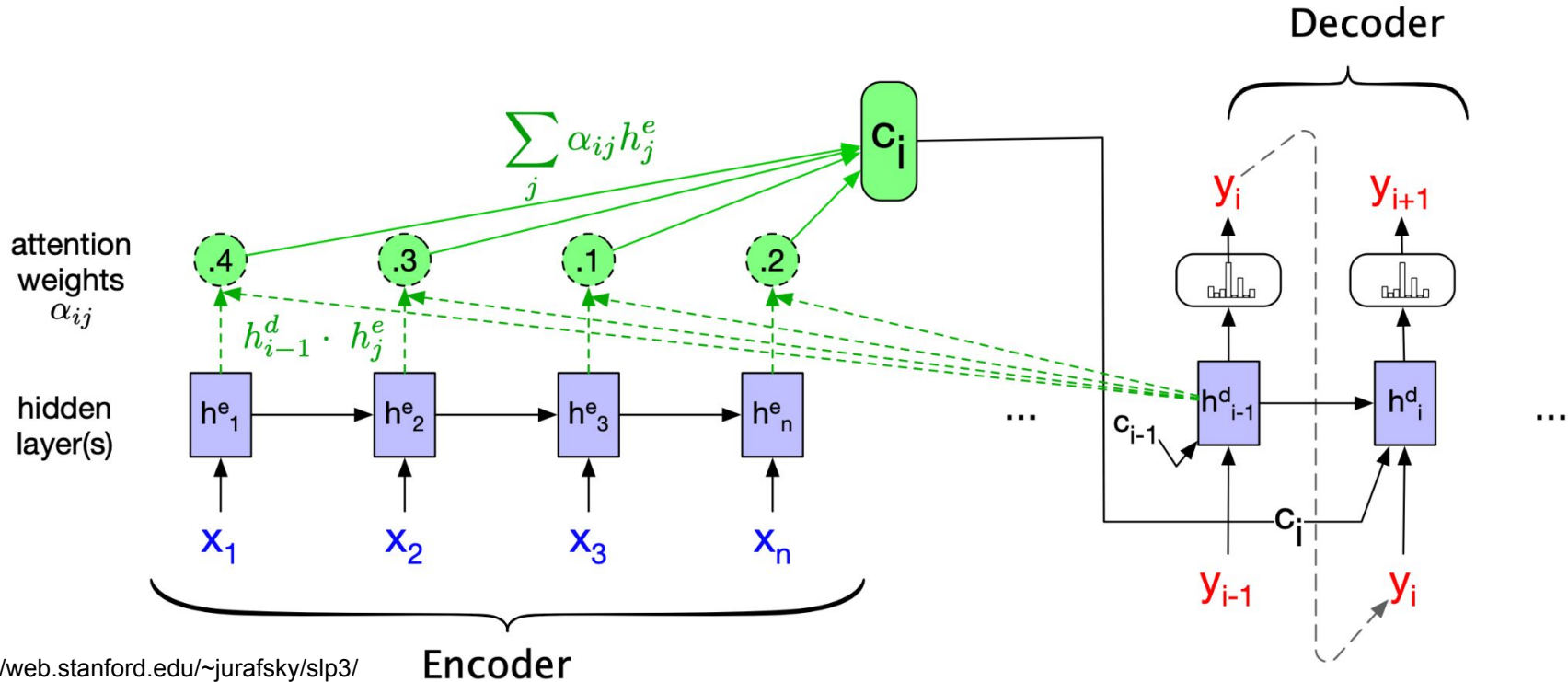
Step 2: Compute attention weights according to alignment with outputs (general attention) or inputs (self attention)

$$\begin{aligned}\alpha_{t,i} &= \text{align}(\mathbf{y}_t, \mathbf{x}_i) \\ &= \frac{\exp(\text{score}(\mathbf{s}_{t-1}, \mathbf{h}_i))}{\sum_{i'=1}^n \exp(\text{score}(\mathbf{s}_{t-1}, \mathbf{h}_{i'}))}\end{aligned}$$

Step 3: Compute the context vector, scaled according to attention weights

$$\mathbf{c}_t = \sum_{i=1}^n \alpha_{t,i} \mathbf{h}_i$$

Attention Mechanism



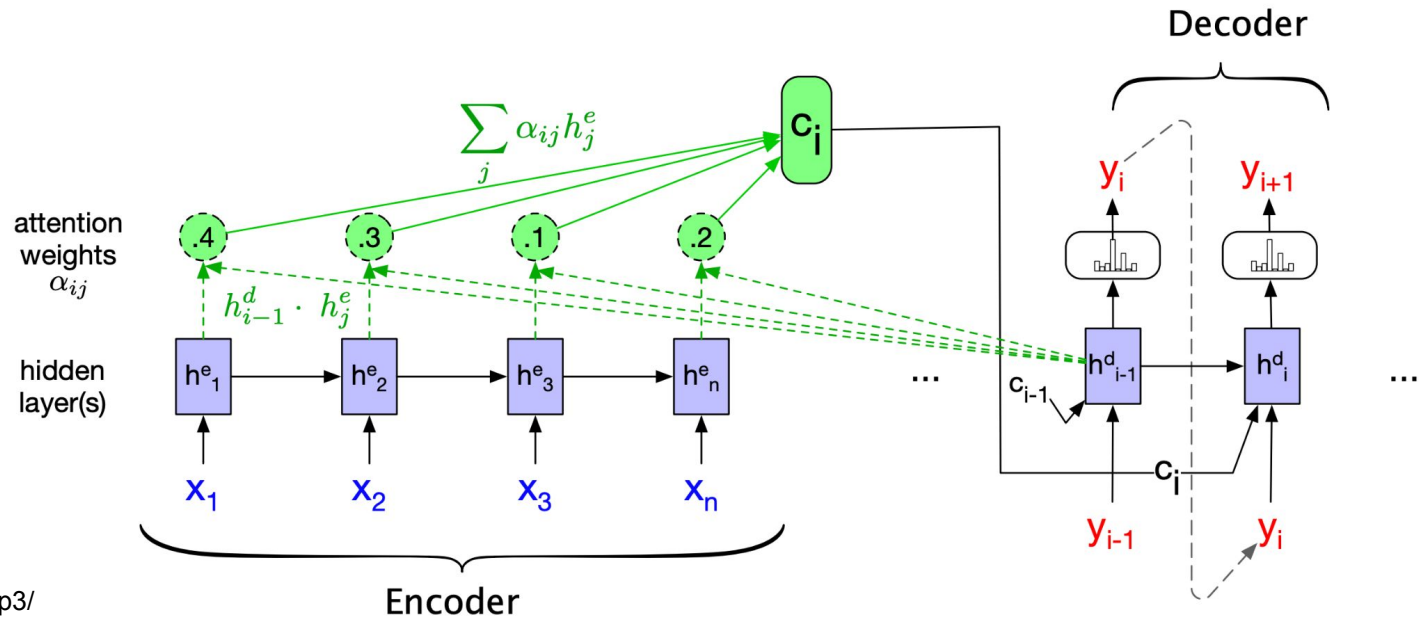
Popular Attention Formulations

- Different score functions have been introduced
 - In practice, the dot-product is simple and effective

Name	Alignment score function	Citation
Content-base attention	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \text{cosine}[\mathbf{s}_t, \mathbf{h}_i]$	<u>Graves2014</u>
Additive(*)	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \mathbf{v}_a^\top \tanh(\mathbf{W}_a[\mathbf{s}_{t-1}; \mathbf{h}_i])$	<u>Bahdanau2015</u>
Dot-Product	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \mathbf{s}_t^\top \mathbf{h}_i$	<u>Luong2015</u>
Scaled Dot-Product(^)	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \frac{\mathbf{s}_t^\top \mathbf{h}_i}{\sqrt{n}}$ Note: very similar to the dot-product attention except for a scaling factor; where n is the dimension of the source hidden state.	<u>Vaswani2017</u>

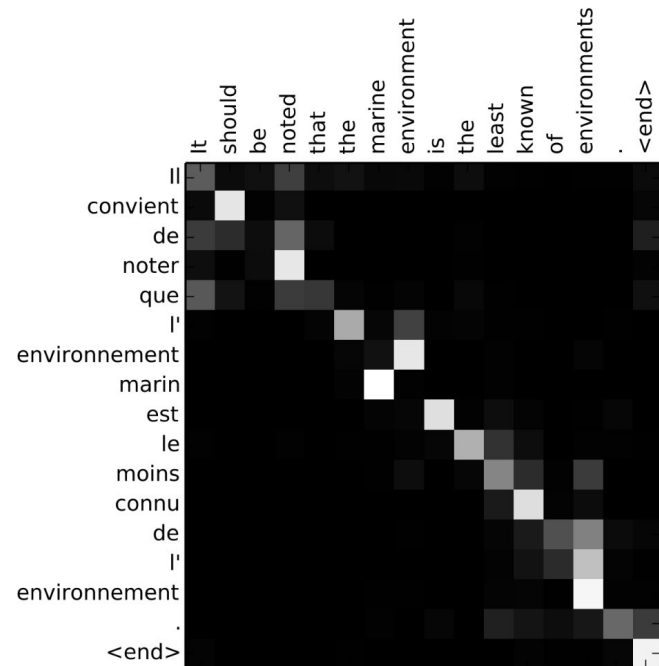
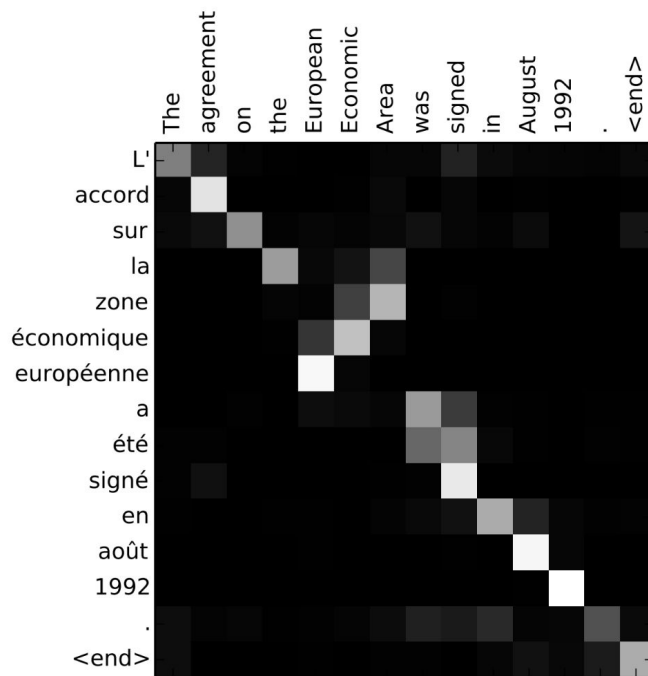
Attention

- Computes time-dependent weighted averages over previous vectors
- Can focus on different aspects of the past sequence



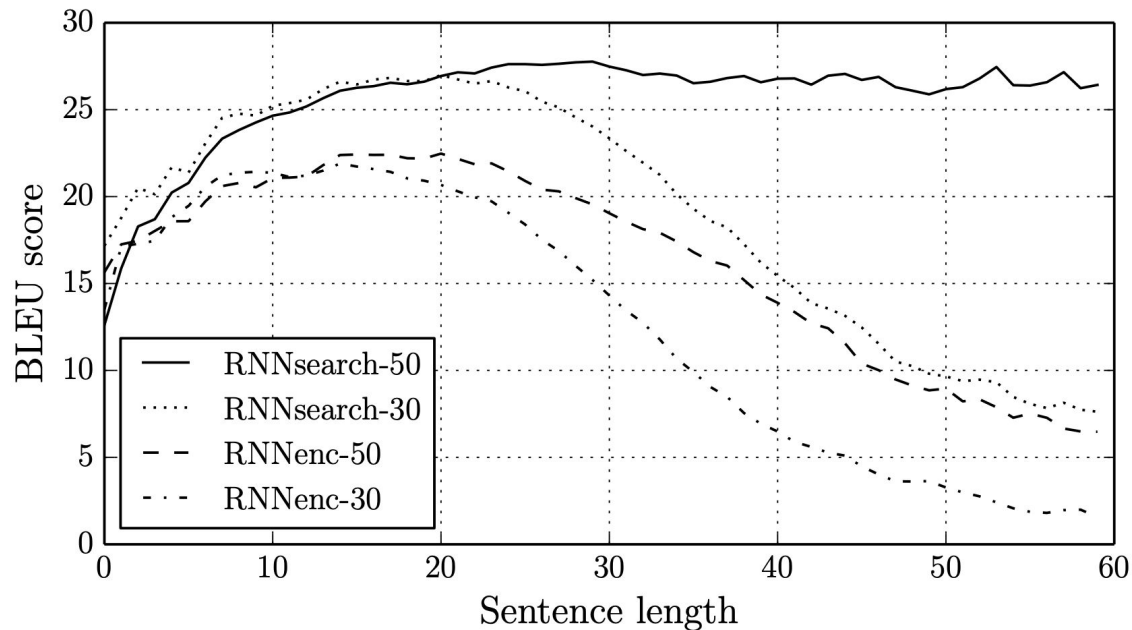
Visualizing Attention

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



Impact of Attention

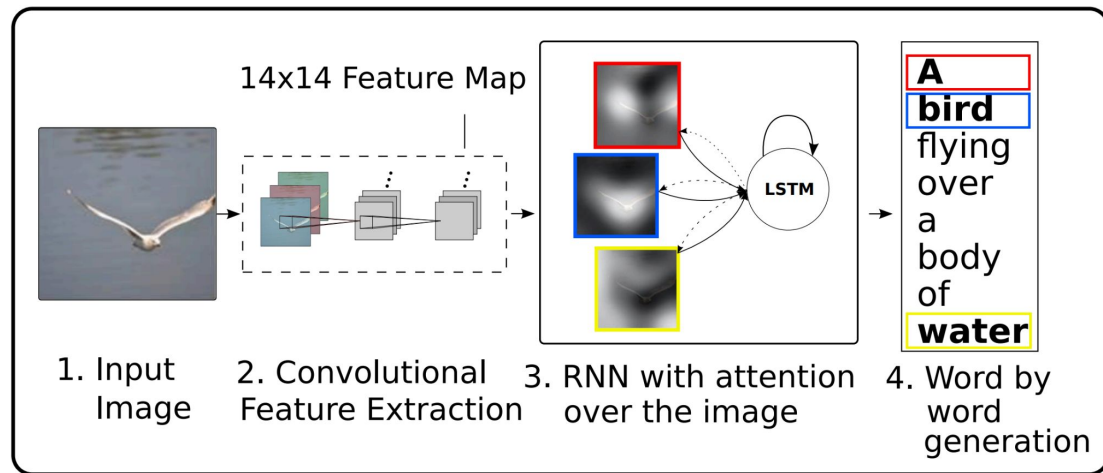
- Really helpful for long sequences
 - Helps solve bottleneck problem!



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)



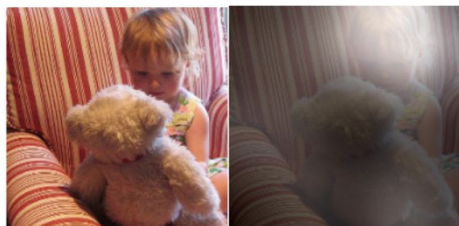
A woman is throwing a frisbee in a park.



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.



A(0.98)



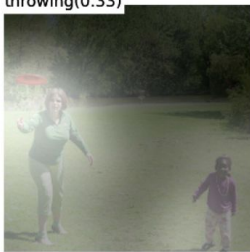
woman(0.54)



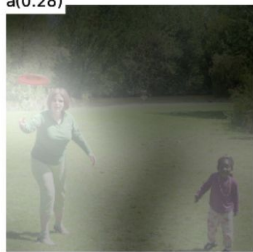
is(0.37)



throwing(0.33)



a(0.28)



frisbee(0.37)



in(0.21)



a(0.18)



park(0.35)



.(0.33)



(b) A woman is throwing a frisbee in a park.

Error Analysis!

Figure 5. Examples of mistakes where we can use attention to gain intuition into what the model saw.



A large white bird standing in a forest.



A woman holding a clock in her hand.



A man wearing a hat and
a hat on a skateboard.



A person is standing on a beach
with a surfboard.



A woman is sitting at a table
with a large pizza.



A man is talking on his cell phone
while another man watches.

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 at handling long sequences
- Attention can look back at past feature vectors!
 - Scales better to long sequences
 - Can incorporate image features
 - Many, many more applications!