



**BOSTON UNIVERSITY**  
**MACHINE INTELLIGENCE**  
**COMMUNITY**

# Sequential Models

RNN, LSTM

Presenters: Julius, Curtis



# Background

Conditional Probability and Sequential Data

# Conditional Probability

1. Conditional Probability is the probability something happens given another event happened.
  - a.  $P(\text{passing} \mid \text{studying}) > P(\text{passing} \mid \text{gaming})$
2. Conditional independence/dependance
  - a. Conditional independence:  $P(x) = P(x|y)$
  - b. Conditional dependance:  $P(x) \neq P(x|y)$



# Sequential Data

1. Data is **sequential** if the probability of one data point is dependant on previous data

a.  $P(x_t) \neq P(x_t | x_{t-1}, \dots, x_0)$

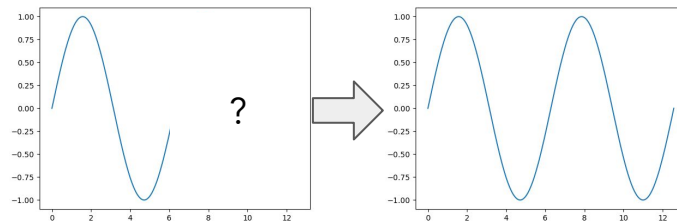
2. Example

a. NLP

i. The color of the bus is \_\_\_\_\_.

ii.  $P(\text{yellow} | \text{"The color of the bus is"}) > P(\text{yellow})$

b. Sine wave



# Sequential models

1. Sequential models predict the next data point given previous data points

2. Model Examples

- a. Recurrent Neural Networks (RNN)
- b. LSTM
- c. GRU

$$P(x_t | x_{t-1}, \dots, x_0)$$

3. Applications

- a. Natural Language Processing
- b. Audio Signal Processing
- c. Reinforcement Learning





# Recurrent Neural Networks

Combining new inputs with previous data

# Recursive Functions

1. Our sequential data is  $x_1, x_2, \dots, x_t$
2. We wish to map each input,  $x$ , with a prediction,  $y_1, y_2, \dots, y_t$

We could learn a function that maps all of the previous inputs to one output, but that would require a function that takes a different number of inputs for each step of the sequence.

$$y_t = f(x_1, \dots, x_t)$$

Instead, we want a function that takes only one state that summarizes all information necessary to make a prediction.

$$y_t = f(s_t)$$



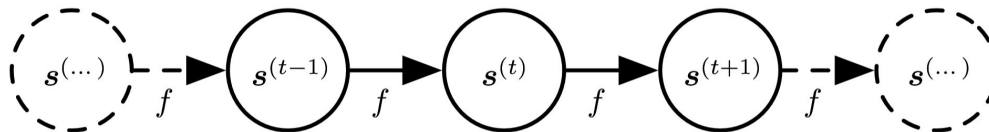
# Parameterized Recursive Functions

1. A recursive function takes an input as a result from itself
  - a. We can add additional parameterization with the parameter theta
  - b. This function could be a feed forward neural network

$$\underbrace{s^{(t)}}_{\text{State at step t}} = f\left(\underbrace{s^{(t-1)}}_{\text{State at step t-1}}; \theta\right)$$

State at step t

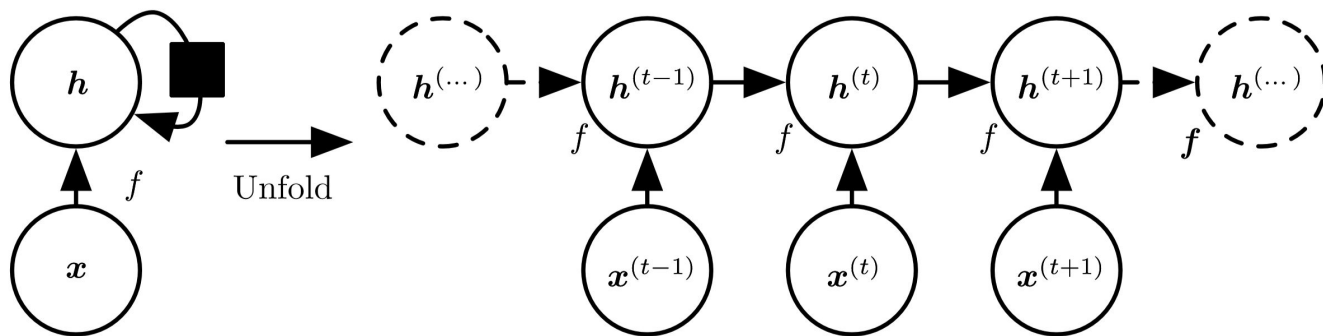
State at step t-1



A state is equal to the function output of the previous state



# Unfolding

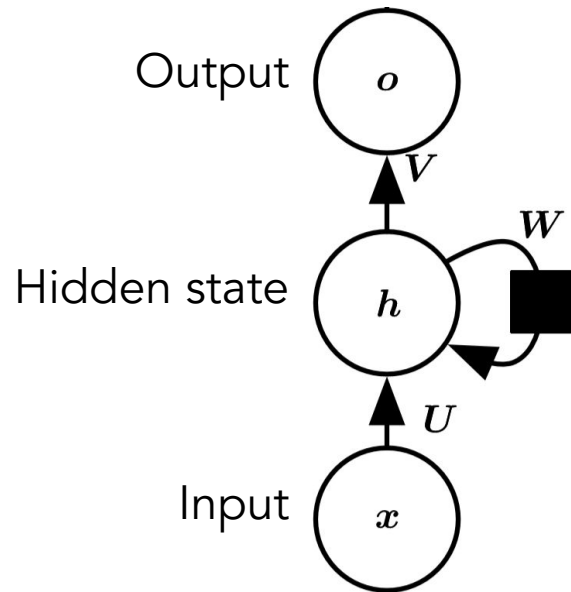


# Recurrent Neural Network Goals

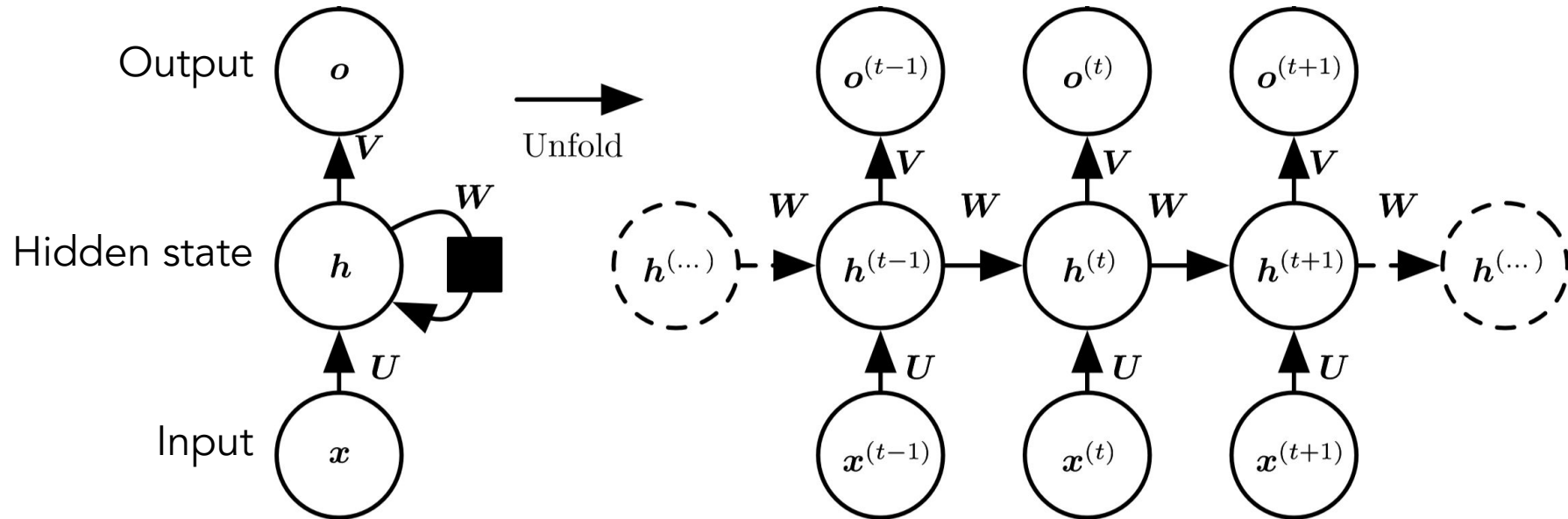
1. Summarize: Our network should take the previous state,  $h_{t-1}$ , and add the new input,  $x_t$ , to get the current state,  $h_t$ .
  - a. This acts to summarize all previous inputs,  $x_{1..t}$ , into one state vector,  $h_t$ .
2. Predict: Our network should predict the output,  $o_t$ , from the current state,  $h_t$ .
  - a. These predictions are also noted as  $\hat{y}_t$ .
3. Train: The network should then compare the output,  $o_t$ , with the target value,  $y_t$ , and generate a loss,  $L_t$ .
  - a. We then train the network through backpropagation using the gradient of the loss with respect to our weight matrices.

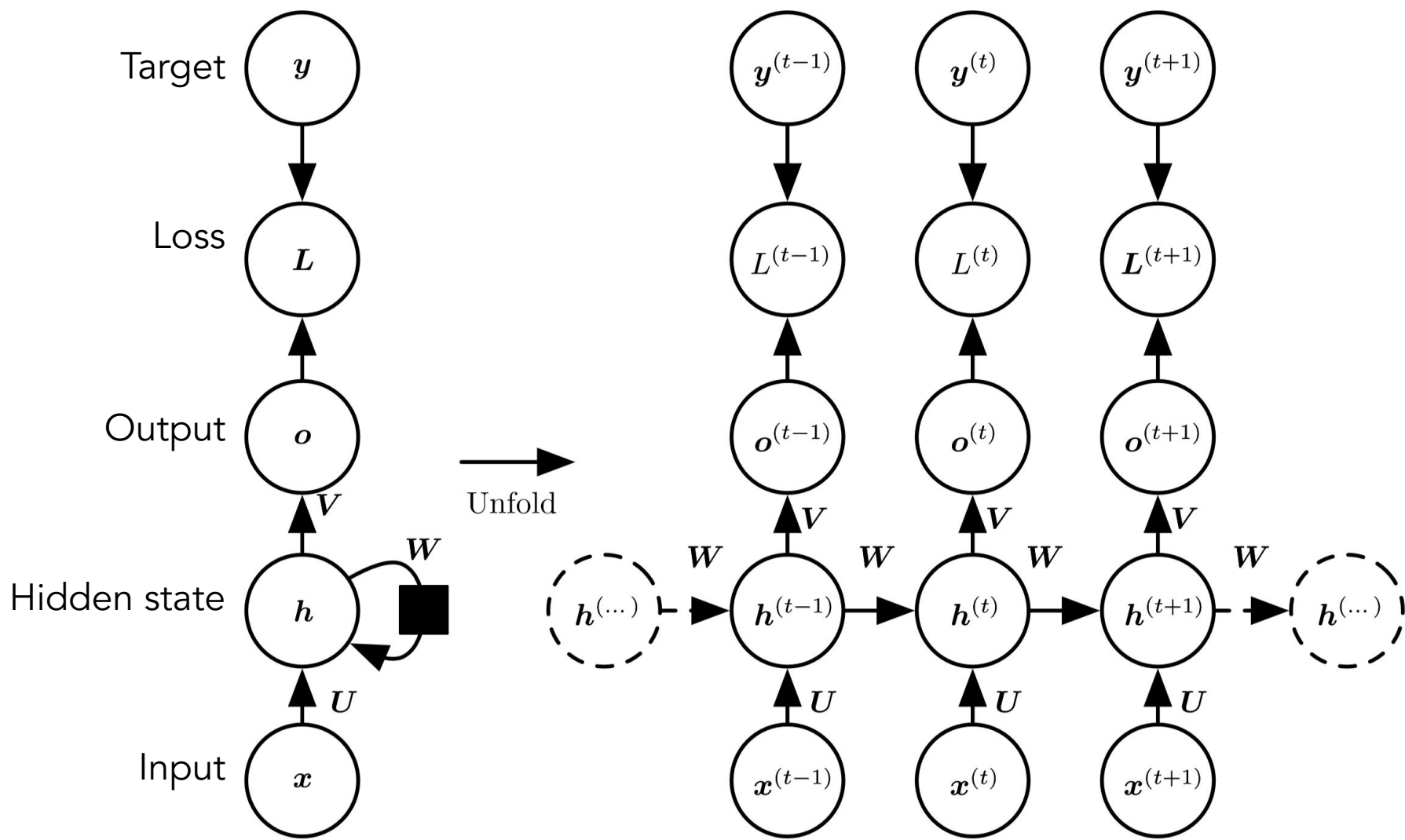


# RNN Cell



# RNN Graph





# RNN Details

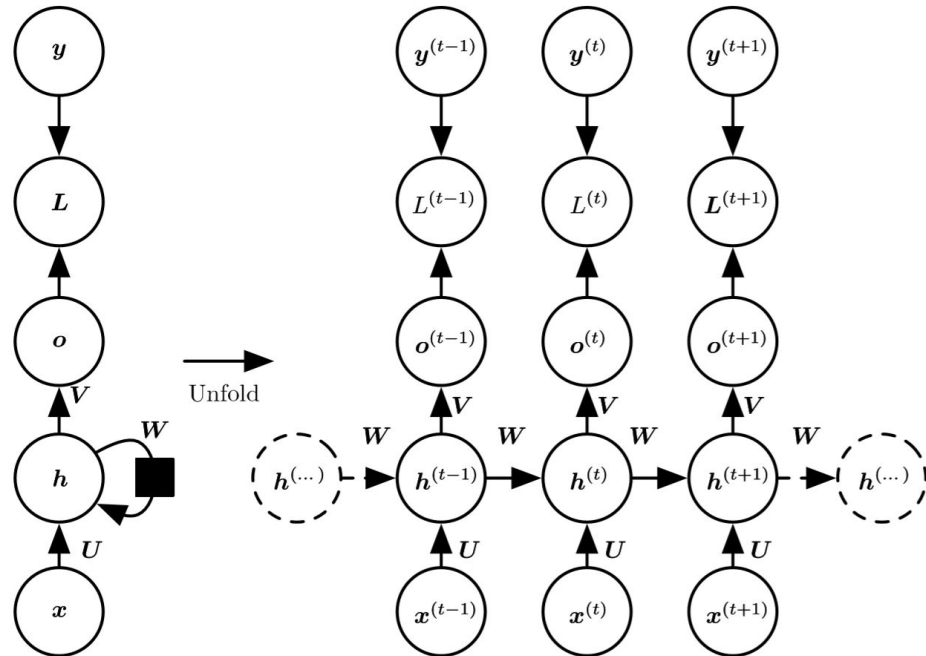
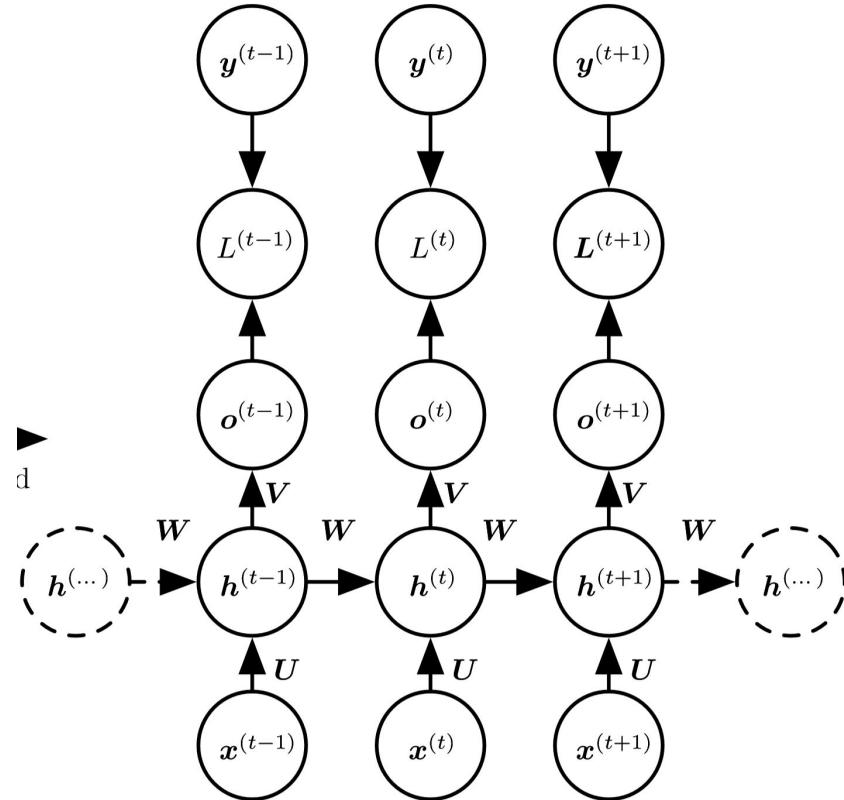


Figure 10.3: The computational graph to compute the training loss of a recurrent network that maps an input sequence of  $x$  values to a corresponding sequence of output  $o$  values. A loss  $L$  measures how far each  $o$  is from the corresponding training target  $y$ . When using softmax outputs, we assume  $o$  is the unnormalized log probabilities. The loss  $L$  internally computes  $\hat{y} = \text{softmax}(o)$  and compares this to the target  $y$ . The RNN has input to hidden connections parametrized by a weight matrix  $U$ , hidden-to-hidden recurrent connections parametrized by a weight matrix  $W$ , and hidden-to-output connections parametrized by a weight matrix  $V$ . Equation 10.8 defines forward propagation in this model. (Left) The RNN and its loss drawn with recurrent connections. (Right) The same seen as a time-unfolded computational graph, where each node is now associated with one particular time instance.

# RNN Feedforward Explanation

1.  $h_0$  is initialized
2.  $x_1$  is combined with  $h_0$  to get  $h_1$
3.  $o_1$  is obtained from  $h_1$
4.  $y_{\text{hat}}_1$  probabilities is obtained from  $o_1$
5.  $L_1$  is computed from  $y_1, y_{\text{hat}}_1$  with softmax
6. Process repeats for each  $t$
- 7.



# RNN Feedforward

Affine  $\mathbf{a}^{(t)} = \underbrace{\mathbf{b} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)}}_{\text{Affine function (hidden network)}},$

Hidden state  $\mathbf{h}^{(t)} = \underbrace{\tanh(\mathbf{a}^{(t)})}_{\text{Activation function}},$

Output  $\mathbf{o}^{(t)} = \underbrace{\mathbf{c} + \mathbf{V}\mathbf{h}^{(t)}}_{\text{Output network}},$



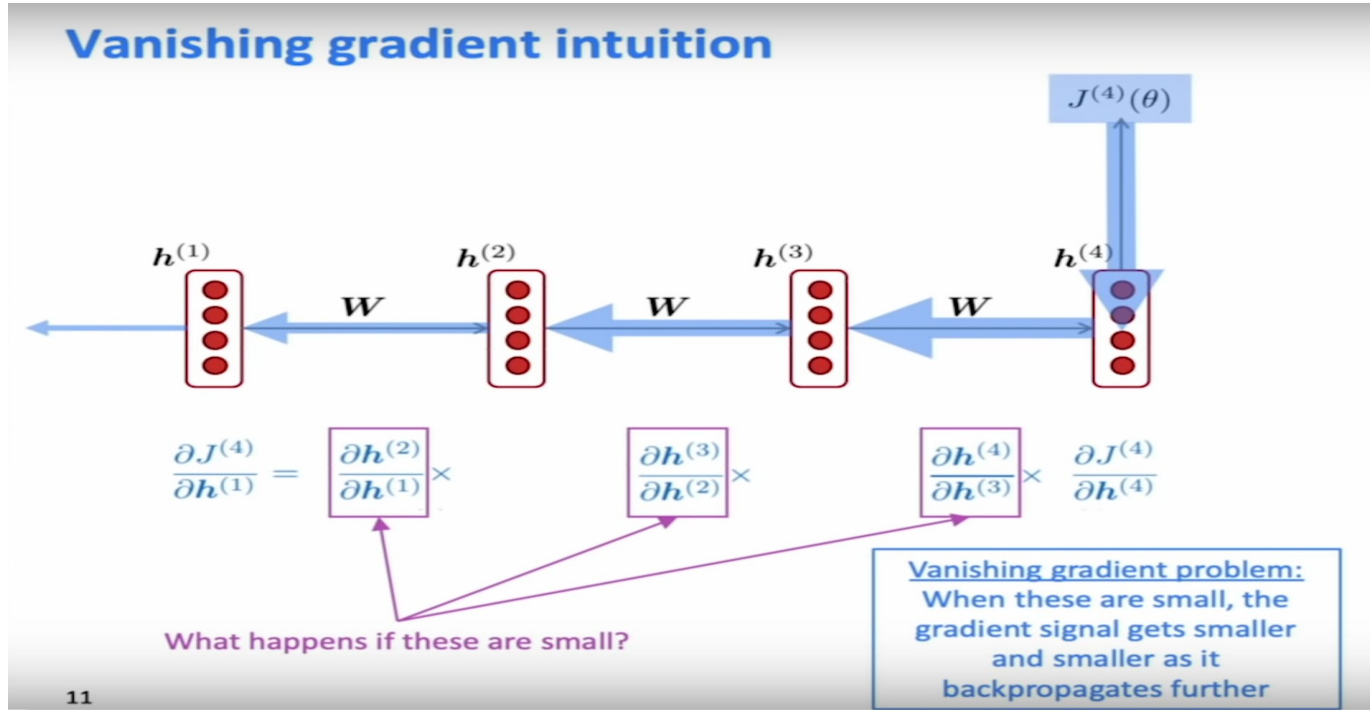


# RNN Feedforward

Predictions  $\hat{\mathbf{y}}^{(t)} = \underbrace{\text{softmax}(\mathbf{o}^{(t)})}_{\text{Prediction probabilities}},$

Loss 
$$\begin{aligned} & L\left(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)}\}, \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)}\}\right) \\ &= \sum_t L^{(t)} \\ &= - \sum_t \log p_{\text{model}}\left(\mathbf{y}^{(t)} \mid \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}\}\right), \end{aligned}$$

# Backpropagation Through Time



11



# Backpropagation Through Time

Gradient formulas for reference:

$$\nabla_{\mathbf{h}^{(\tau)}} L = \mathbf{V}^\top \nabla_{\mathbf{o}^{(\tau)}} L. \quad (10.19)$$

We can then iterate backward in time to back-propagate gradients through time, from  $t = \tau - 1$  down to  $t = 1$ , noting that  $\mathbf{h}^{(t)}$  (for  $t < \tau$ ) has as descendents both  $\mathbf{o}^{(t)}$  and  $\mathbf{h}^{(t+1)}$ . Its gradient is thus given by

$$\nabla_{\mathbf{h}^{(t)}} L = \left( \frac{\partial \mathbf{h}^{(t+1)}}{\partial \mathbf{h}^{(t)}} \right)^\top (\nabla_{\mathbf{h}^{(t+1)}} L) + \left( \frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{h}^{(t)}} \right)^\top (\nabla_{\mathbf{o}^{(t)}} L) \quad (10.20)$$

$$= \mathbf{W}^\top (\nabla_{\mathbf{h}^{(t+1)}} L) \text{diag} \left( 1 - (\mathbf{h}^{(t+1)})^2 \right) + \mathbf{V}^\top (\nabla_{\mathbf{o}^{(t)}} L), \quad (10.21)$$

$$\nabla_{\mathbf{c}} L = \sum_t \left( \frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{c}} \right)^\top \nabla_{\mathbf{o}^{(t)}} L = \sum_t \nabla_{\mathbf{o}^{(t)}} L, \quad (10.22)$$

$$\nabla_{\mathbf{b}} L = \sum_t \left( \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{b}^{(t)}} \right)^\top \nabla_{\mathbf{h}^{(t)}} L = \sum_t \text{diag} \left( 1 - (\mathbf{h}^{(t)})^2 \right) \nabla_{\mathbf{h}^{(t)}} L, \quad (10.23)$$

$$\nabla_{\mathbf{V}} L = \sum_t \sum_i \left( \frac{\partial L}{\partial o_i^{(t)}} \right) \nabla_{\mathbf{V} o_i^{(t)}} = \sum_t (\nabla_{\mathbf{o}^{(t)}} L) \mathbf{h}^{(t)\top}, \quad (10.24)$$

375

---

## CHAPTER 10

---

$$\nabla_{\mathbf{W}} L = \sum_t \sum_i \left( \frac{\partial L}{\partial h_i^{(t)}} \right) \nabla_{\mathbf{W}^{(t)} h_i^{(t)}} \quad (10.25)$$

$$= \sum_t \text{diag} \left( 1 - (\mathbf{h}^{(t)})^2 \right) (\nabla_{\mathbf{h}^{(t)}} L) \mathbf{h}^{(t-1)\top}, \quad (10.26)$$

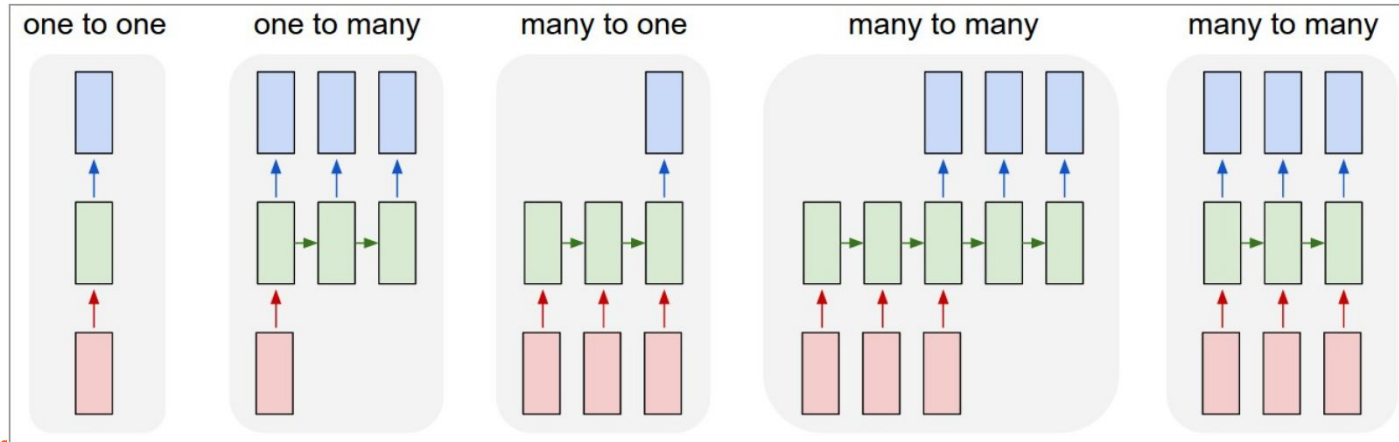
$$\nabla_{\mathbf{U}} L = \sum_t \sum_i \left( \frac{\partial L}{\partial h_i^{(t)}} \right) \nabla_{\mathbf{U}^{(t)} h_i^{(t)}} \quad (10.27)$$

$$= \sum_t \text{diag} \left( 1 - (\mathbf{h}^{(t)})^2 \right) (\nabla_{\mathbf{h}^{(t)}} L) \mathbf{x}^{(t)\top}, \quad (10.28)$$



# Types of RNN problems

1. One-one: (given an image, is this a bird or not)
2. One-many: (given an image, give a description of this image)
3. Many-one: (given video, classify genre; predict sentiment of sentence)
4. Many-many: (given video, predict label for each frame; translate a sentence)



# Deep RNNs

1. Why not just make RNNs more deep? Doesn't solve...
  - a. Vanishing gradient problem
  - b. The network "forgets" information in the hidden state over long time steps
2. Adding Convolutions
  - a. Use CNNs to have a "sliding window" over time steps
  - b. Takes advantage of extracting information from features close in time with less computational overhead
3. Adding Attention
  - a. Makes the network focus on particular parts of the input that is relevant to prediction





# Long short-term memory (LSTM)

Brief introduction to LSTM's

# Issues with vanilla RNN's

1. Vanishing Gradients
2. Difficult to retain long term memory within our hidden state  $h$

# Solution

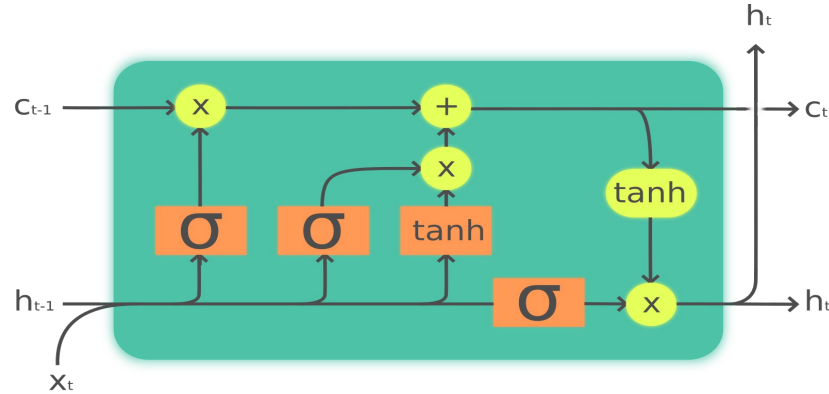
1. Long Short Term Memory cells which maintain a hidden state, and a cell state.
2. RNN is able to read, write and erase information from cell state through gates.
  - a. Gates are nearly open or nearly closed, open meaning information can pass through, closed meaning information cannot pass through.





# Vanilla LSTM Cell

1.



Legend:

Layer



Pointwise op



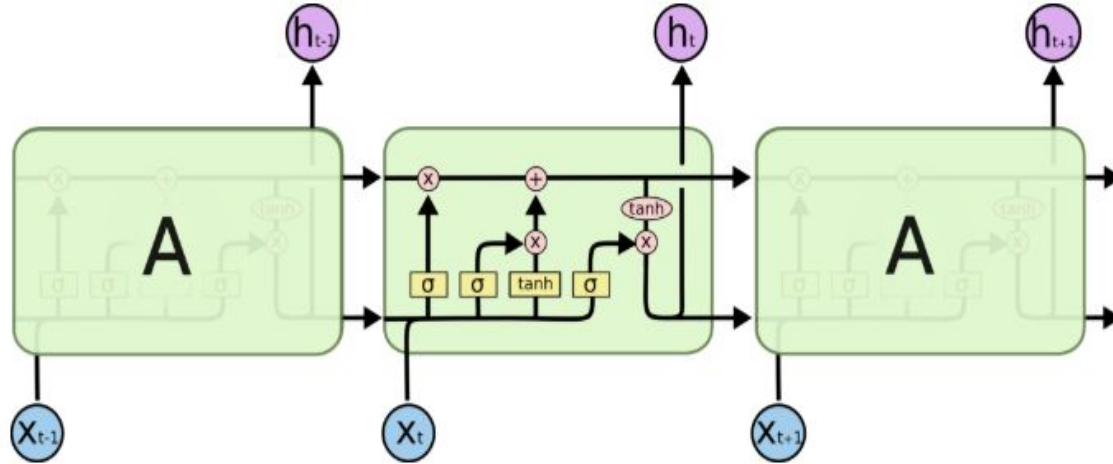
Copy



[https://en.wikipedia.org/wiki/Long\\_short-term\\_memory#/media/File:The\\_LSTM\\_cell.png](https://en.wikipedia.org/wiki/Long_short-term_memory#/media/File:The_LSTM_cell.png)

# Multiple LSTM Cells

1.

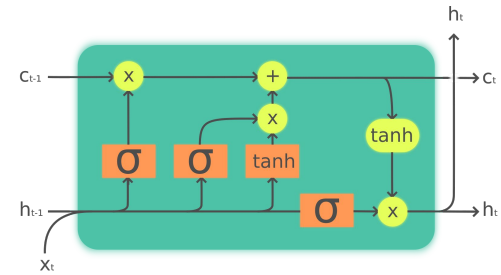


The repeating module in an LSTM contains four interacting layers.

<https://hackernoon.com/understanding-architecture-of-lstm-cell-from-scratch-with-code-8da40f0b71f4>

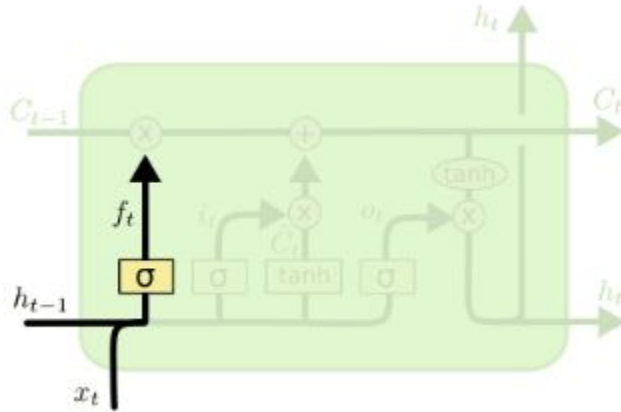
# Vanilla LSTM Cell

1. Composed of several smaller neural networks
2. Each layer/nn have their own weights
3. Layers typically end in sigmoid or tanh
  - a. Sigmoid squashes inputs between 0 and 1 deciding what information should be forgotten.
  - b. Tanh rescales inputs between -1 and 1.



$H :=$  Previous Output,  $X :=$  new input,  $C :=$  Cell State

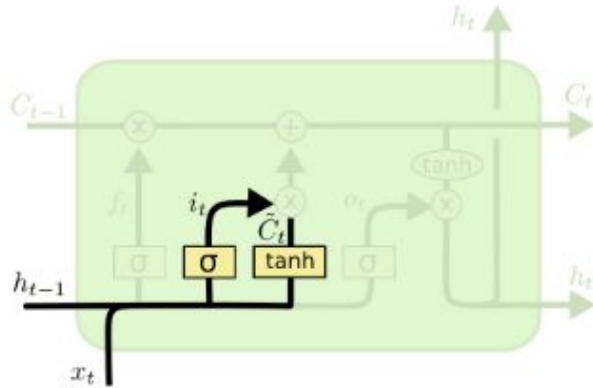
# Intuition: Forget Gate



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

- Based on the previous output and current input, decides what information is relevant by squashing the combined feature vector  $[h_{t-1}, x_t]$  with a sigmoid function (squashes inputs between 0 and 1, 0 being forgotten input).
- Uses element-wise multiplication to effectively forget unimportant cell state features based on the relevancy matrix  $f_t$

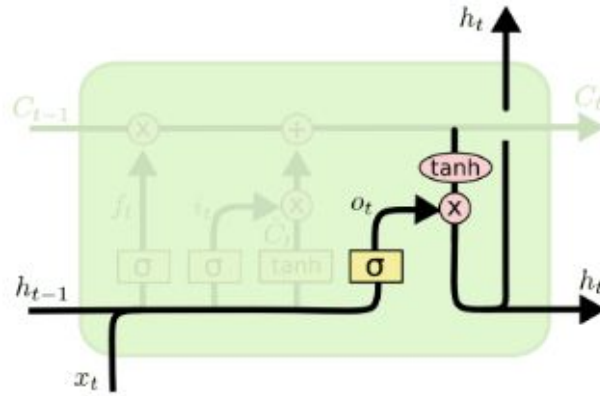
# Intuition: Input Gate



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

- Adds in new information to the cell state. Note, cell state acts as long term memory within a lstm cell.
- Does this by determining what is relevant like in the forget gate ( $i_t$ ) and element wise multiplying that by rescaled inputs ( $\tilde{C}_t$ ) then adding it to the current cell state ( $C_t$ ).

# Intuition: Output Gate



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

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

- Uses the existing cell state to decide what information will be propagated onwards to the next step/what the general output will be.
- Once again decides what's relevant (sigmoid) then forgets/reinforces relevant information through a dot product.

# Conclusion

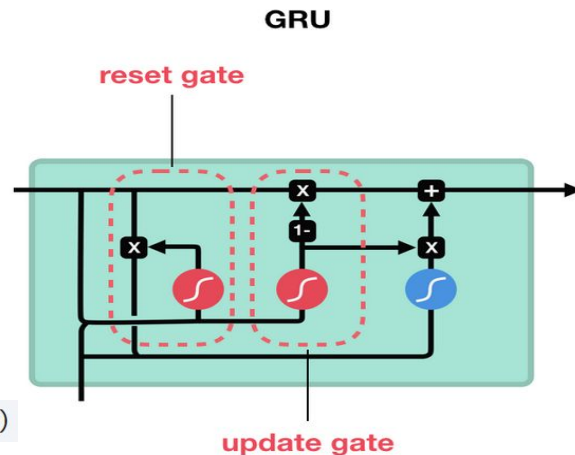
1. LSTM's allow long term memory to be stored for longer in comparison to vanilla RNN's
2. LSTM's do not avoid vanishing gradients entirely

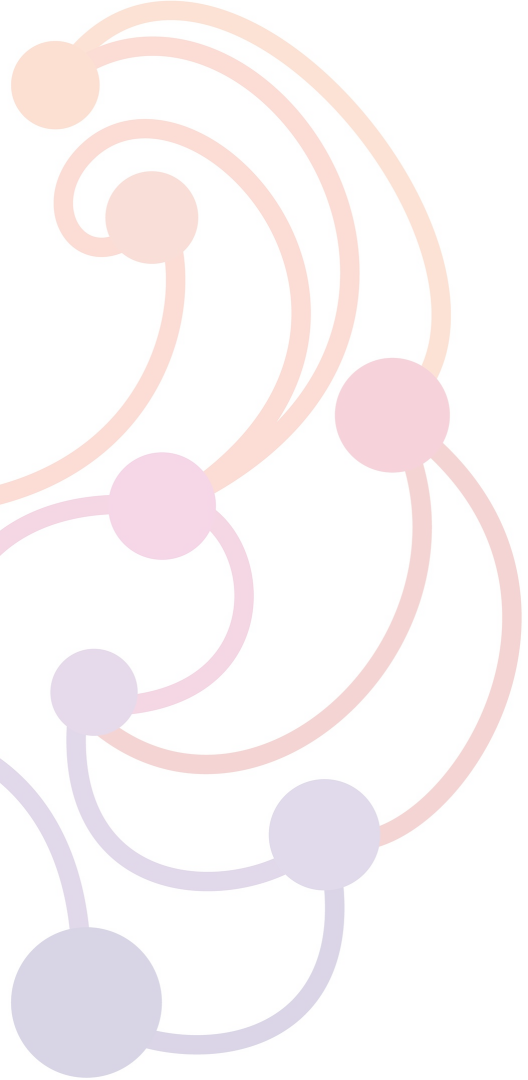
GRUs act similarly to LSTM's with fewer weights and fewer gates.

Remember you can just do :)

```
CLASS torch.nn.LSTM(*args, **kwargs)
```

```
CLASS torch.nn.GRU(*args, **kwargs)
```





Thank you for coming!



# References

Goodfellow, Ian, et al. *Deep Learning*. MIT Press, 2017. <http://www.deeplearningbook.org>