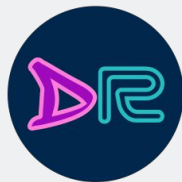


# ROB 430/599.430: Deep Learning for Robot Perception and Manipulation (DeepRob)

---

Lecture 16: Sequence Prediction (RNNs, Seq2Seq, etc.);  
Attention and Transformers

03/16/2026



<https://deeprob.org/w26/>

# Today

---

\*P4 released\*

- Feedback and Recap (5min)
- Processing Sequences (1hr10min)
  - RNN
  - LSTM
  - Seq2Seq
  - Attention
  - Transformer
  - ViT (if time allows)
- Summary and Takeaways (5min)

# So far...

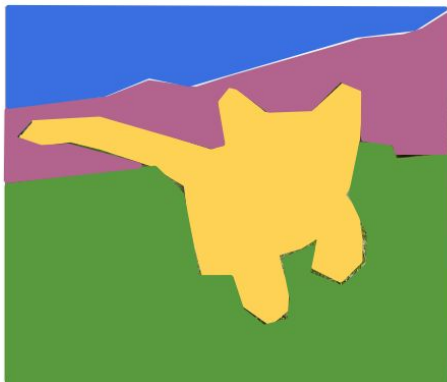
## Classification



**CAT**

No spatial extent

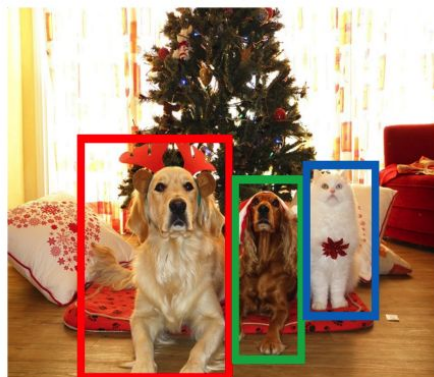
## Semantic Segmentation



**GRASS, CAT, TREE, SKY**

No objects, just pixels

## Object Detection



**DOG, DOG, CAT**

Multiple Objects

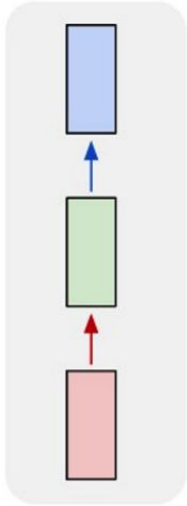
## Instance Segmentation



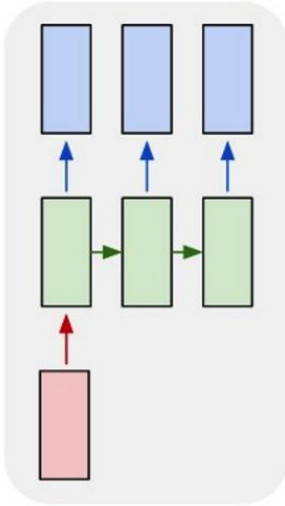
**DOG, DOG, CAT**

# Recurrent Neural Networks

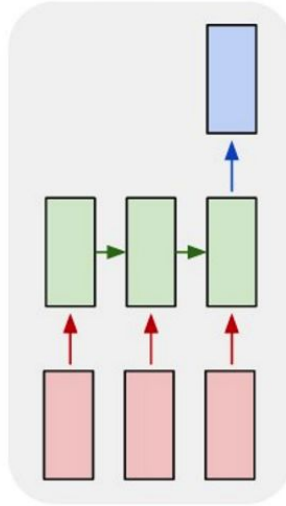
(1)  
one to one



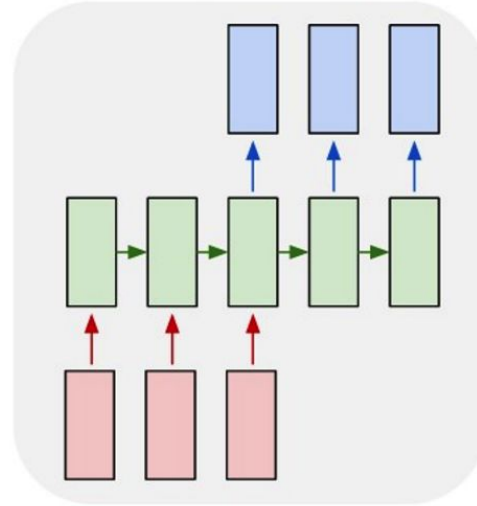
(2)  
one to many



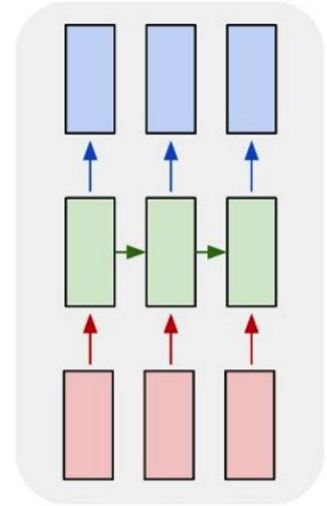
(3)  
many to one



(4)  
many to many



(5)  
many to many



# Aha Slides (In-class participation)

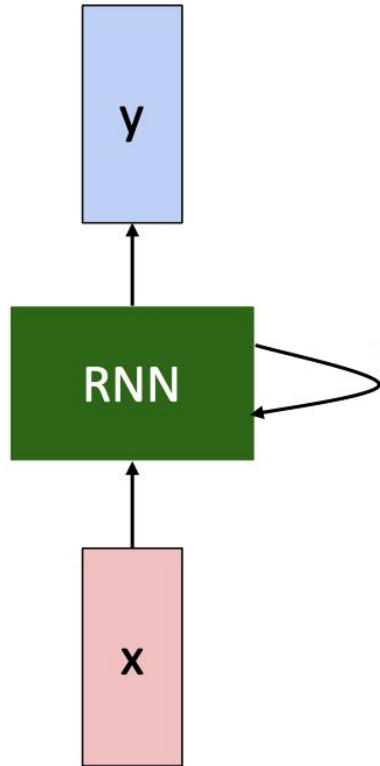
<https://ahaslides.com/ATZ7W>



Q1: what are some examples of applications?

# Recurrent Neural Networks

---

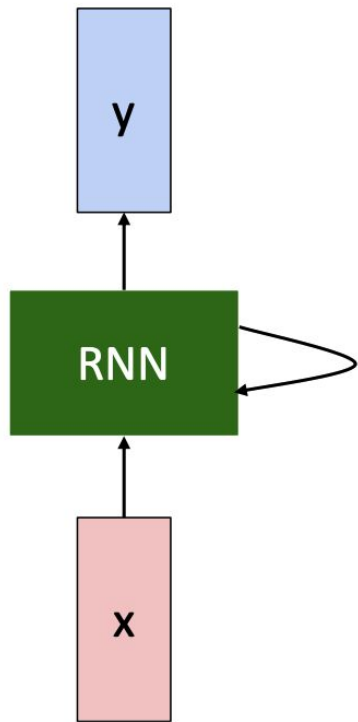


Key idea: RNNs have an “internal state” that is updated as a sequence is processed

# Recurrent Neural Networks

---

We can process a sequence of vectors  $\mathbf{x}$  by applying a **recurrence formula** at every time step:



$$h_t = f_W(h_{t-1}, x_t)$$

new state

old state

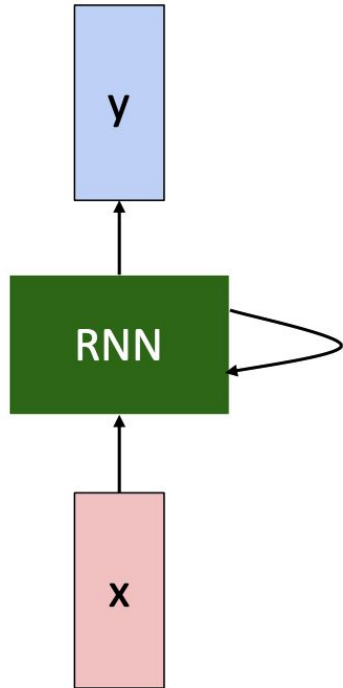
input vector at  
some time step

some function  
with parameters  $W$

# (Vanilla) Recurrent Neural Networks

The state consists of a single “hidden” vector  $\mathbf{h}$ :

$$h_t = f_W(h_{t-1}, x_t)$$



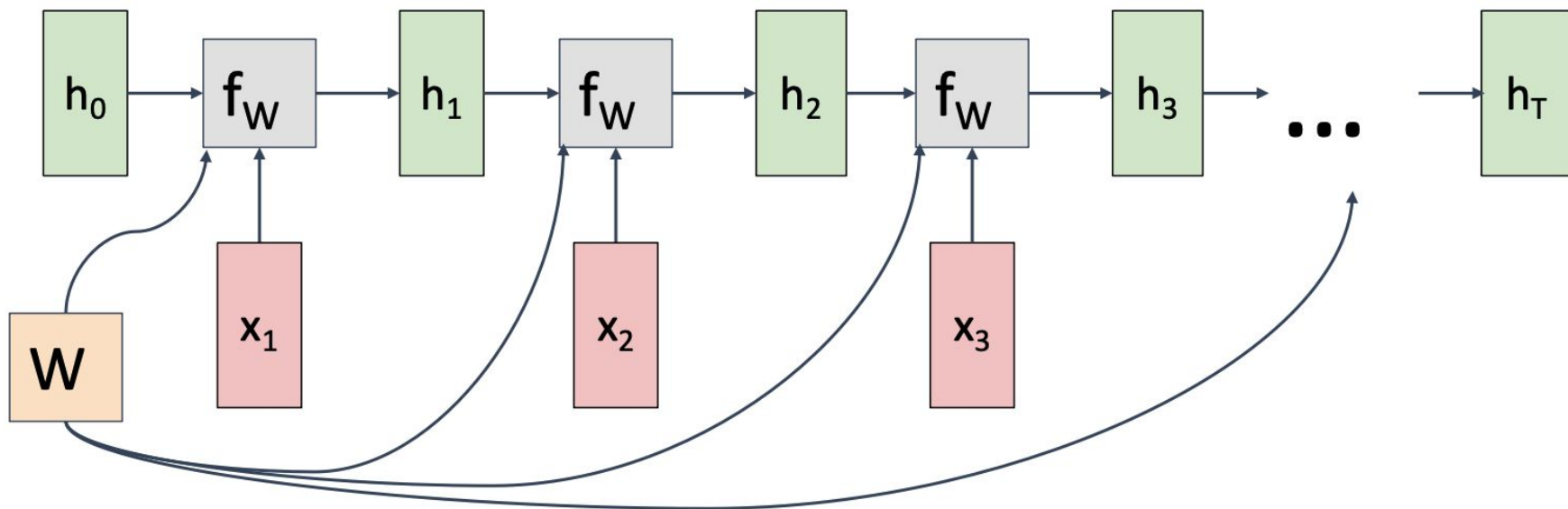
$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$$

$$y_t = W_{hy}h_t + b_y$$

Sometimes called a “Vanilla RNN” or an “Elman RNN” after Prof. Jeffrey Elman

# Recurrent Neural Networks: computational graph

Re-use the same weight matrix at every time-step



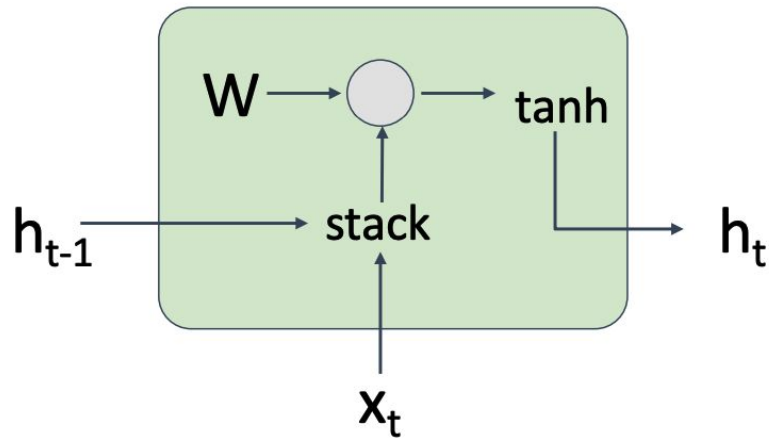
# Aha Slides (In-class participation)

<https://ahaslides.com/ATZ7W>



Q2: why re-use (shared) weights?

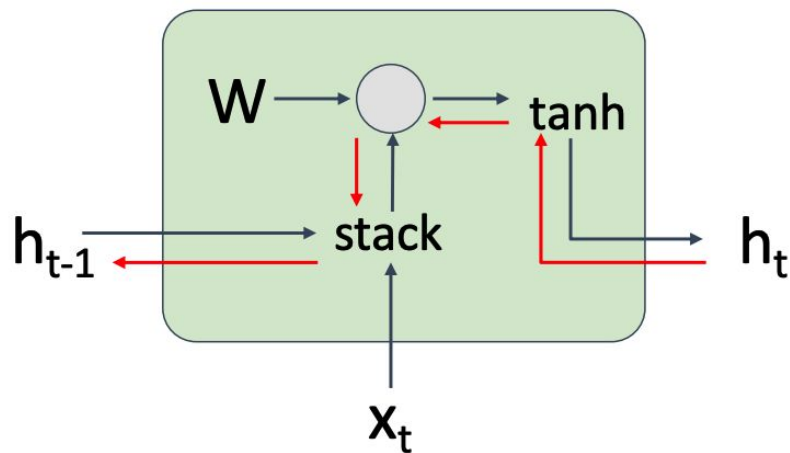
# (Vanilla) Recurrent Neural Networks: Gradient Flow



$$\begin{aligned}h_t &= \tanh(W_{hh}h_{t-1} + W_{hx}x_t + b_h) \\ &= \tanh\left((W_{hh} \quad W_{hx}) \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h\right) \\ &= \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h\right)\end{aligned}$$

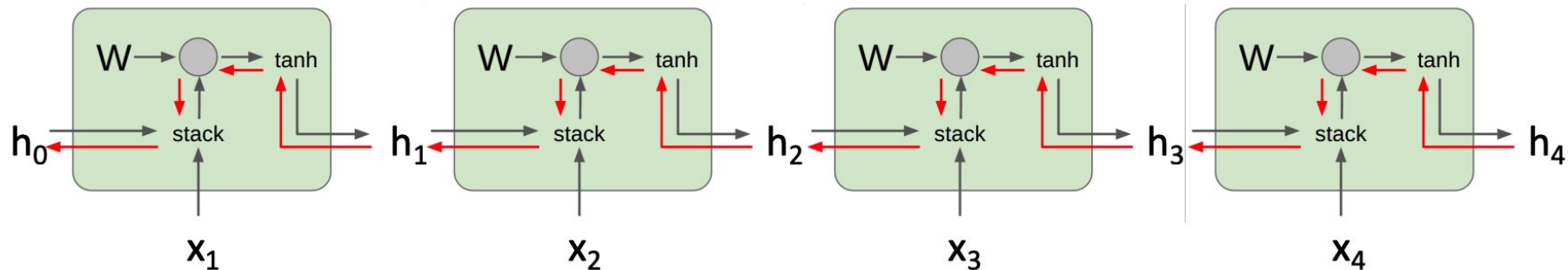
# (Vanilla) RNN: Gradient Flow

Backpropagation from  $h_t$  to  $h_{t-1}$  multiplies by  $W$  (actually  $W_{hh}^T$ )



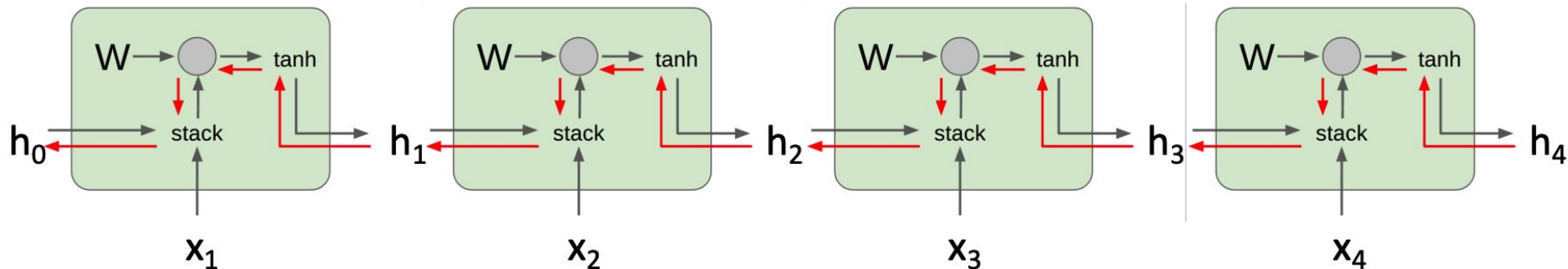
$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{hx}x_t + b_h) \\ &= \tanh\left((W_{hh} \quad W_{hx}) \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h\right) \\ &= \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h\right) \end{aligned}$$

# (Vanilla) RNN: Gradient Flow



Computing gradient of  $h_0$  involves many factors of  $W$  (and repeated  $\tanh$ )

# (Vanilla) RNN: Gradient Flow

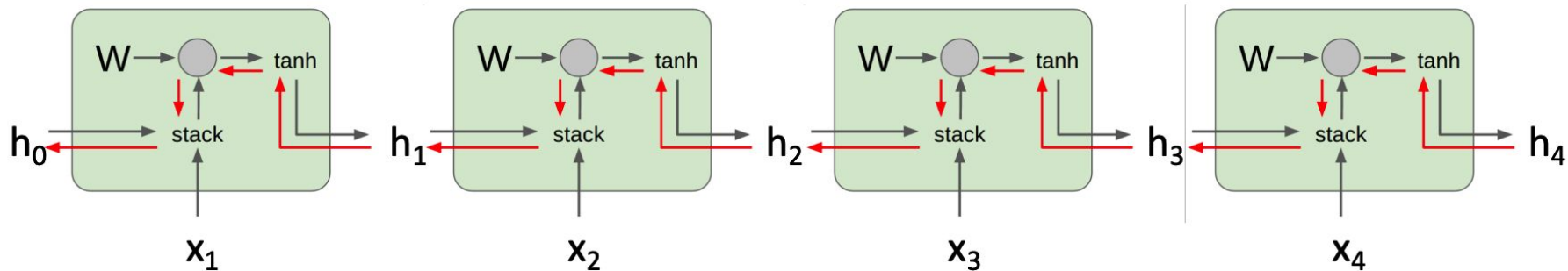


Computing gradient of  $h_0$  involves many factors of  $W$  (and **repeated** tanh)

Largest singular value  $> 1$ :  
**Exploding gradients**

Largest singular value  $< 1$ :  
**Vanishing gradients**

# (Vanilla) RNN: Gradient Flow



Computing gradient of  $h_0$  involves many factors of  $W$  (and repeated tanh)

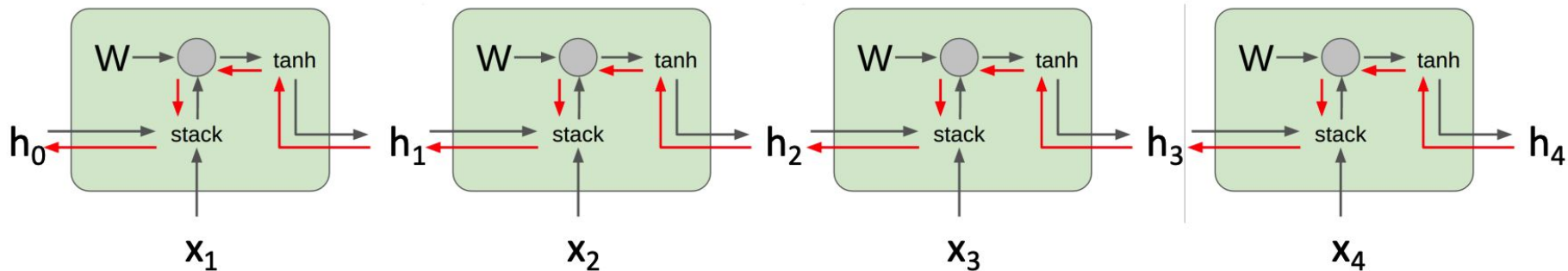
Largest singular value  $> 1$ :  
**Exploding gradients**

Largest singular value  $< 1$ :  
**Vanishing gradients**

**Gradient clipping:** Scale gradient if its norm is too big

```
grad_norm = np.sum(grad * grad)
if grad_norm > threshold:
    grad *= (threshold / grad_norm)
```

# (Vanilla) RNN: Gradient Flow



Computing gradient of  $h_0$  involves many factors of  $W$  (and repeated tanh)

Largest singular value  $> 1$ :  
**Exploding gradients**

Largest singular value  $< 1$ :  
**Vanishing gradients**

→ **Change RNN architecture!**

# Long Short Term Memory (LSTM)

\*Long-range dependency\*

## Vanilla RNN

$$h_t = \tanh \left( W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h \right)$$

Two vectors at each timestep:

Cell state:  $c_t \in \mathbb{R}^H$

Hidden state:  $h_t \in \mathbb{R}^H$

## LSTM

$$\begin{pmatrix} i_t \\ f_t \\ o_t \\ g_t \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} \left( W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h \right)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$

$$h_t = o_t \odot \tanh(c_t)$$

Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation 1997 <https://www.bioinf.jku.at/publications/older/2604.pdf>

# Long Short Term Memory (LSTM)

Compute four “gates” per timestep:

Input gate:  $i_t \in \mathbb{R}^H$

Forget gate:  $f_t \in \mathbb{R}^H$

Output gate:  $o_t \in \mathbb{R}^H$

“Gate?” gate:  $g_t \in \mathbb{R}^H$

## LSTM

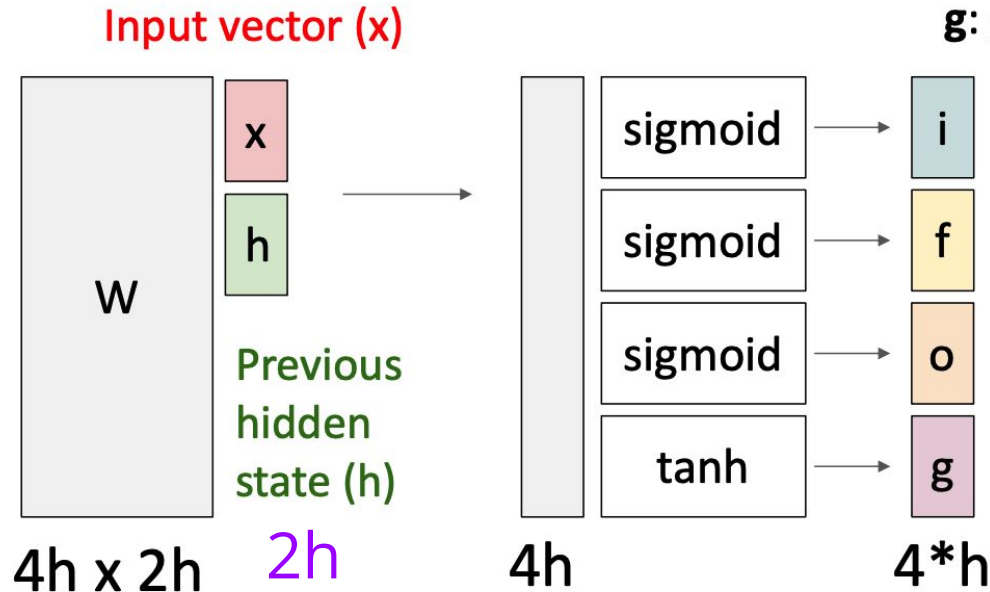
$$\begin{pmatrix} i_t \\ f_t \\ o_t \\ g_t \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} \left( W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h \right)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$

$$h_t = o_t \odot \tanh(c_t)$$

Hochreiter and Schmidhuber, “Long Short Term Memory”, Neural  
Computation 1997 <https://www.bioinf.jku.at/publications/older/2604.pdf>

# Long Short Term Memory (LSTM)

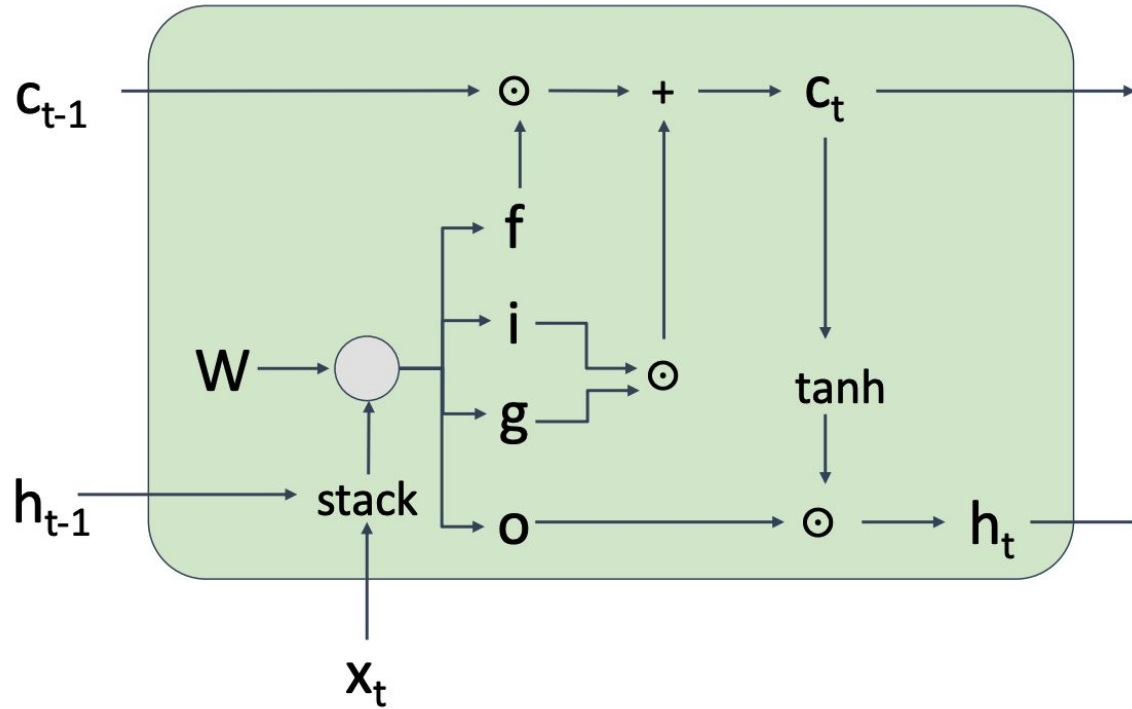


- i**: Input gate, whether to write to cell
- f**: Forget gate, Whether to erase cell
- o**: Output gate, How much to reveal cell
- g**: Gate gate (?), How much to write to cell

$$\begin{pmatrix} i_t \\ f_t \\ o_t \\ g_t \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} \left( W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h \right)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$
$$h_t = o_t \odot \tanh(c_t)$$

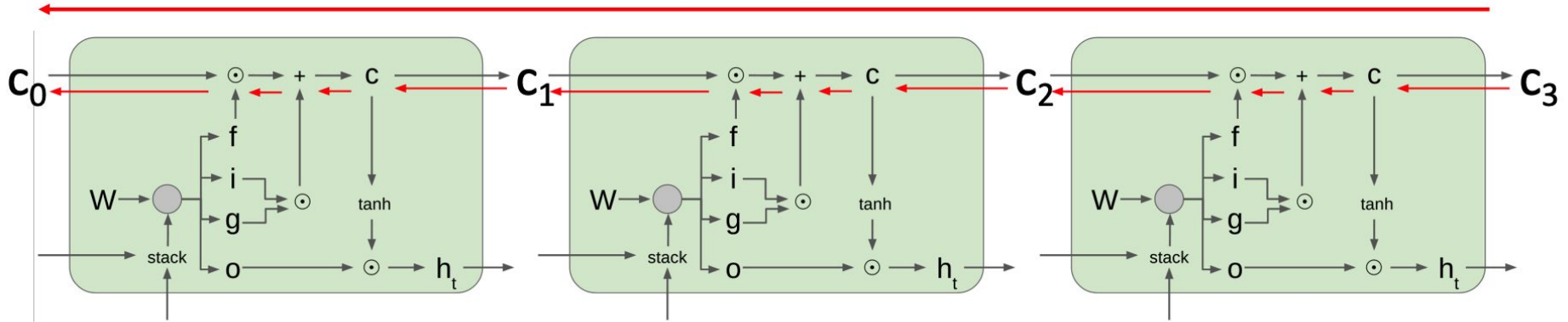
# Long Short Term Memory (LSTM)



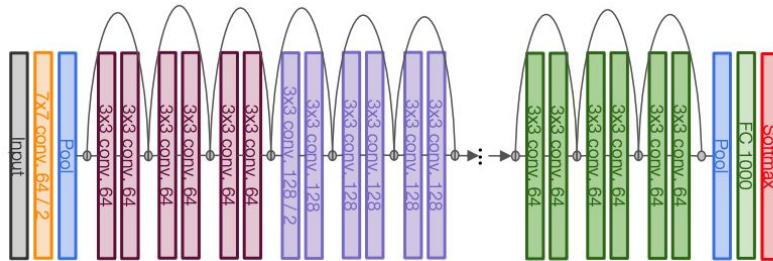
$$\begin{pmatrix} i_t \\ f_t \\ o_t \\ g_t \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} \left( W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h \right)$$
$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$
$$h_t = o_t \odot \tanh(c_t)$$

# Long Short Term Memory (LSTM)

Uninterrupted gradient flow!



Similar to ResNet!



“Highway network”

<https://static.googleusercontent.com/media/research.google.com/en/pubs/archive/46171.pdf>

<https://arxiv.org/pdf/1505.00387>

# Single-Layer RNN

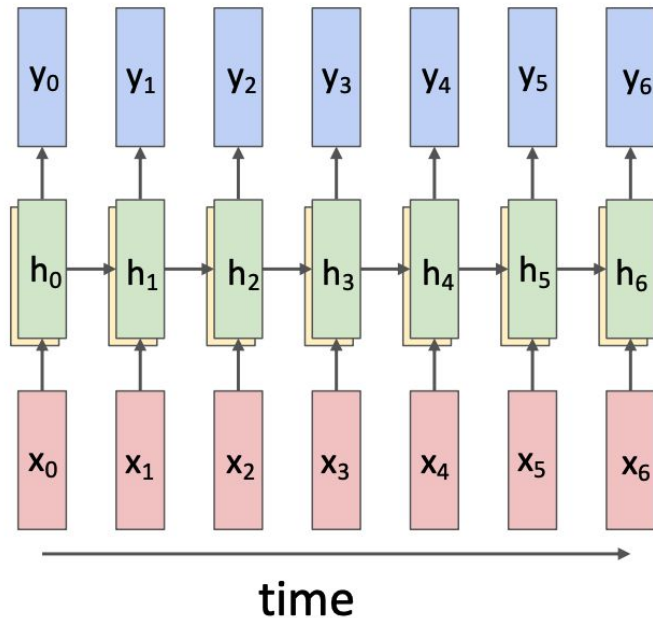
$$h_t = \tanh \left( W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h \right)$$

LSTM:

$$\begin{pmatrix} i_t \\ f_t \\ o_t \\ g_t \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} \left( W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h \right)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$

$$h_t = o_t \odot \tanh(c_t)$$



# Multi-Layer RNN

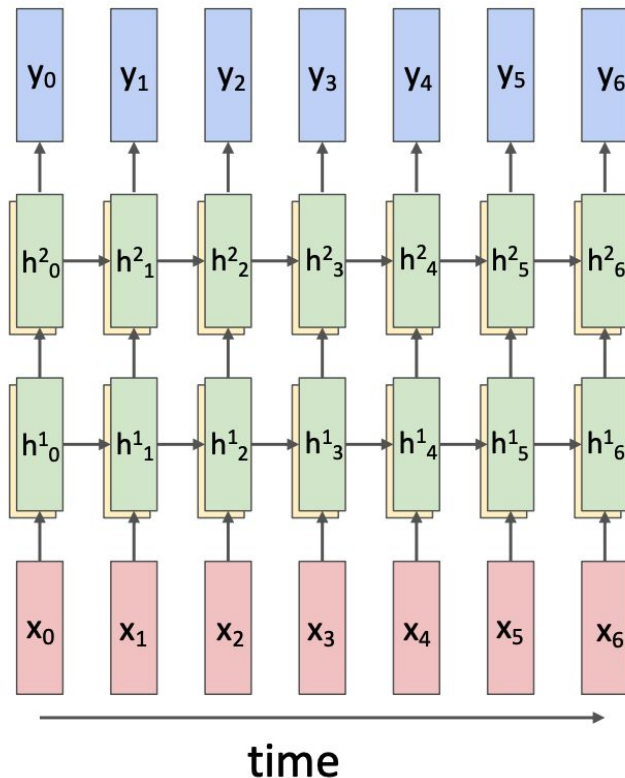
$$h_t^\ell = \tanh \left( W \begin{pmatrix} h_{t-1}^\ell \\ h_t^{\ell-1} \end{pmatrix} + b_h^\ell \right)$$

LSTM:

$$\begin{pmatrix} i_t^\ell \\ f_t^\ell \\ o_t^\ell \\ g_t^\ell \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} \left( W \begin{pmatrix} h_{t-1}^\ell \\ h_t^{\ell-1} \end{pmatrix} + b_h^\ell \right)$$
$$c_t^\ell = f_t^\ell \odot c_{t-1}^\ell + i_t^\ell \odot g_t^\ell$$
$$h_t^\ell = o_t^\ell \odot \tanh(c_t^\ell)$$

depth

**Two-layer RNN:** Pass hidden states from one RNN as inputs to another RNN



Add Three-layer, etc.

# Vanilla RNN - 112 Lines of Python Code

---

Minimal character-level language model with a Vanilla Recurrent Neural Network, in Python/numpy

min-char-rnn.py

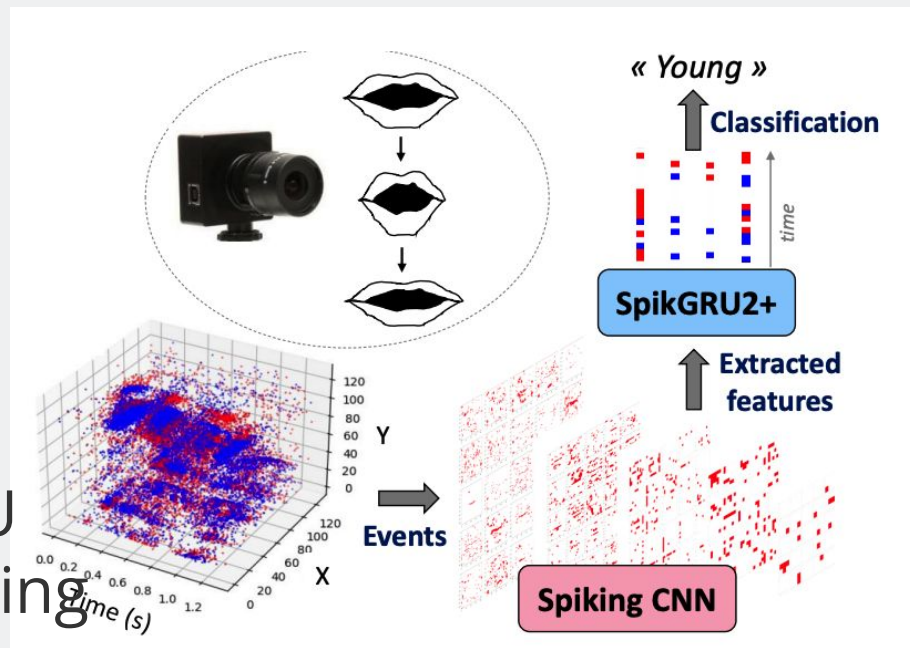
```
1 """
2 Minimal character-level Vanilla RNN model. Written by Andrej Karpathy (@karpathy)
3 BSD License
4 """
5 import numpy as np
6
7 # data I/O
8 data = open('input.txt', 'r').read() # should be simple plain text file
9 chars = list(set(data))
10 data_size, vocab_size = len(data), len(chars)
11 print 'data has %d characters, %d unique.' % (data_size, vocab_size)
12 char_to_ix = { ch:i for i,ch in enumerate(chars) }
13 ix_to_char = { i:ch for i,ch in enumerate(chars) }
14
15 # hyperparameters
16 hidden_size = 100 # size of hidden layer of neurons
17 seq_length = 25 # number of steps to unroll the RNN for
18 learning_rate = 1e-1
19
20 # model parameters
21 Wxh = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
22 Whh = np.random.randn(hidden_size, hidden_size)*0.01 # hidden to hidden
23 Why = np.random.randn(vocab_size, hidden_size)*0.01 # hidden to output
24 bh = np.zeros((hidden_size, 1)) # hidden bias
25 by = np.zeros((vocab_size, 1)) # output bias
26
27 def lossFun(inputs, targets, hprev):
28     """
29     inputs, targets are both list of integers.
30     hprev is Hx1 array of initial hidden state
31     returns the loss, gradients on model parameters, and last hidden state
32     """
33     xs, hs, ys, ps = {}, {}, {}, {}
34     hs[-1] = np.copy(hprev)
35     loss = 0
```

<https://gist.github.com/karpathy/d4dee566867f8291f086>

# Gated Recurrent Unit (GRU)

$$\begin{aligned}r_t &= \sigma(W_{xr}x_t + W_{hr}h_{t-1} + b_r) \\z_t &= \sigma(W_{xz}x_t + W_{hz}h_{t-1} + b_z) \\ \tilde{h}_t &= \tanh(W_{xh}x_t + W_{hh}(r_t \odot h_{t-1}) + b_h) \\h_t &= z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t\end{aligned}$$

E.g., CVPR 2024  
Spiking Neural Networks + GRU  
for speech recognition/lip reading



Cho et al "Learning phrase representations using RNN encoder-decoder for statistical machine translation", 2014

<https://arxiv.org/abs/1406.1078>

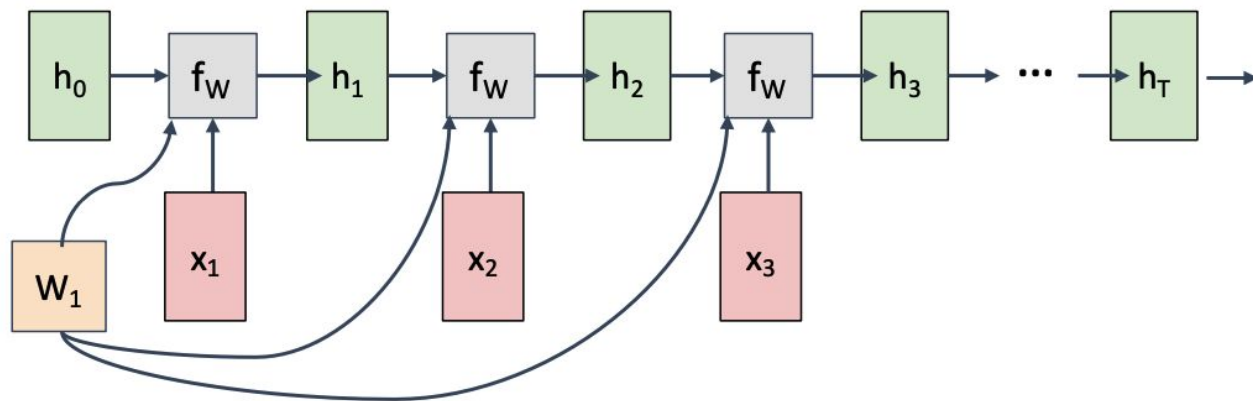
[https://openaccess.thecvf.com/content/CVPR2024W/EVW/papers/Dampfhofer\\_Neuromorphic\\_Lip-Reading\\_with\\_Signed\\_Spiking\\_Gated\\_Recurrent\\_Units\\_CVPRW\\_2024\\_paper.pdf](https://openaccess.thecvf.com/content/CVPR2024W/EVW/papers/Dampfhofer_Neuromorphic_Lip-Reading_with_Signed_Spiking_Gated_Recurrent_Units_CVPRW_2024_paper.pdf)

Spike GPT <https://arxiv.org/pdf/2302.13939>

# Seq2Seq: Sequence to Sequence

---

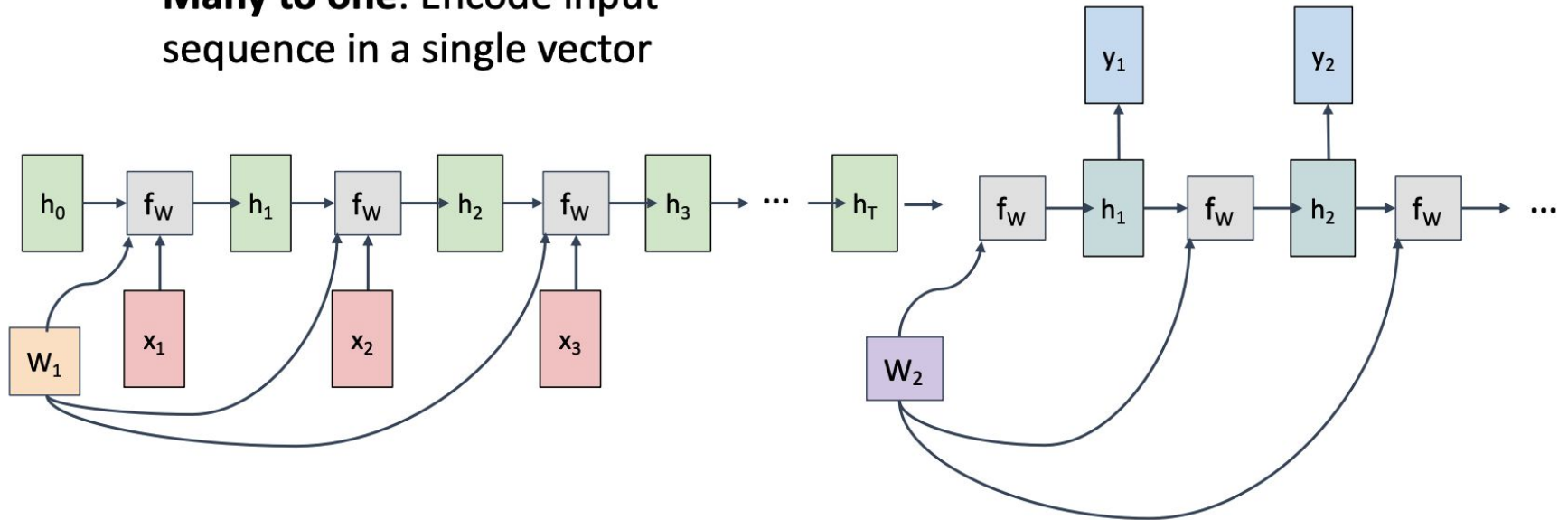
**Many to one:** Encode input sequence in a single vector



# Seq2Seq: Sequence to Sequence

**One to many:** Produce output sequence from single input vector

**Many to one:** Encode input sequence in a single vector



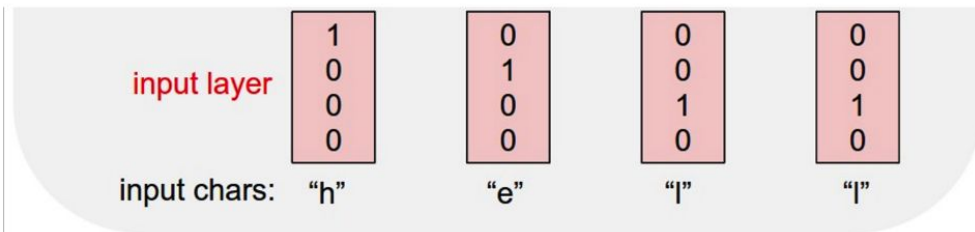
# Example: Language Modeling

---

Given characters 1, 2, ..., t-1,  
model predicts character t

Training sequence: "hello"

Vocabulary: [h, e, l, o]



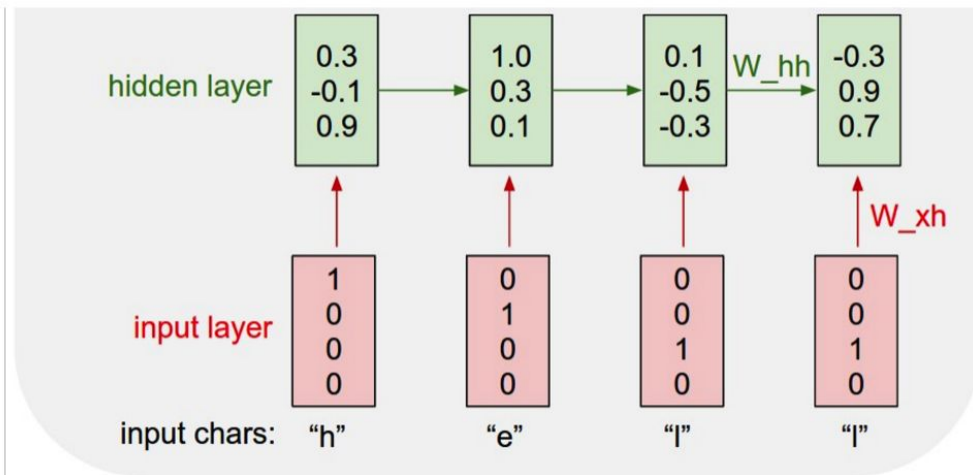
# Example: Language Modeling

Given characters 1, 2, ..., t-1,  
model predicts character t

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

Training sequence: "hello"

Vocabulary: [h, e, l, o]



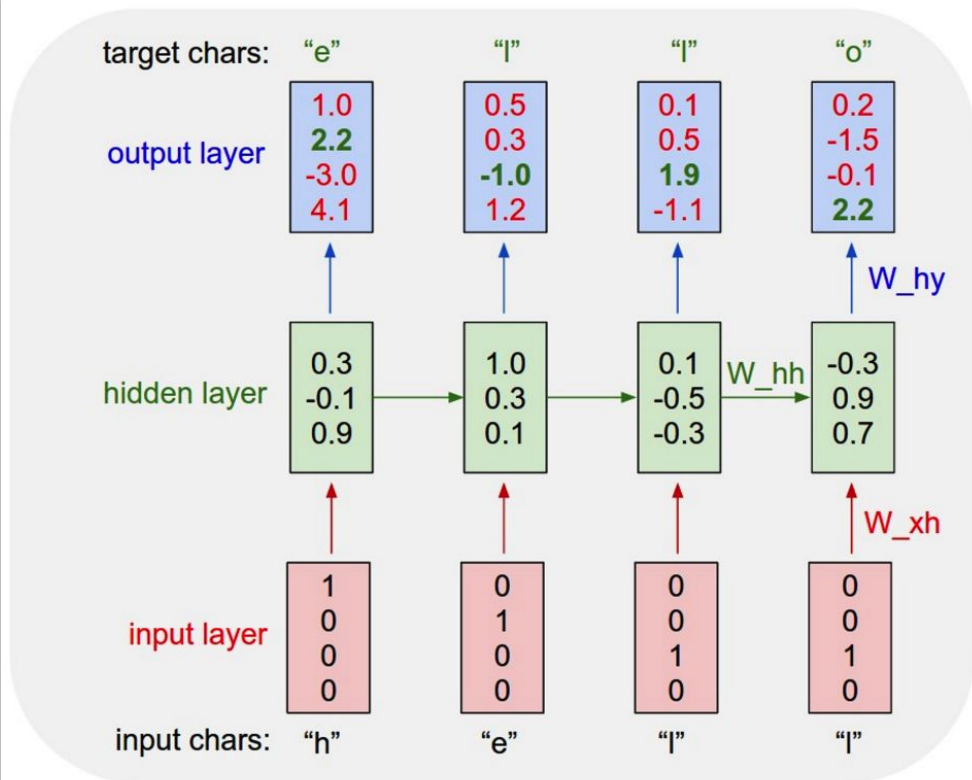
# Example: Language Modeling

Given characters 1, 2, ..., t-1,  
model predicts character t

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

Training sequence: "hello"

Vocabulary: [h, e, l, o]



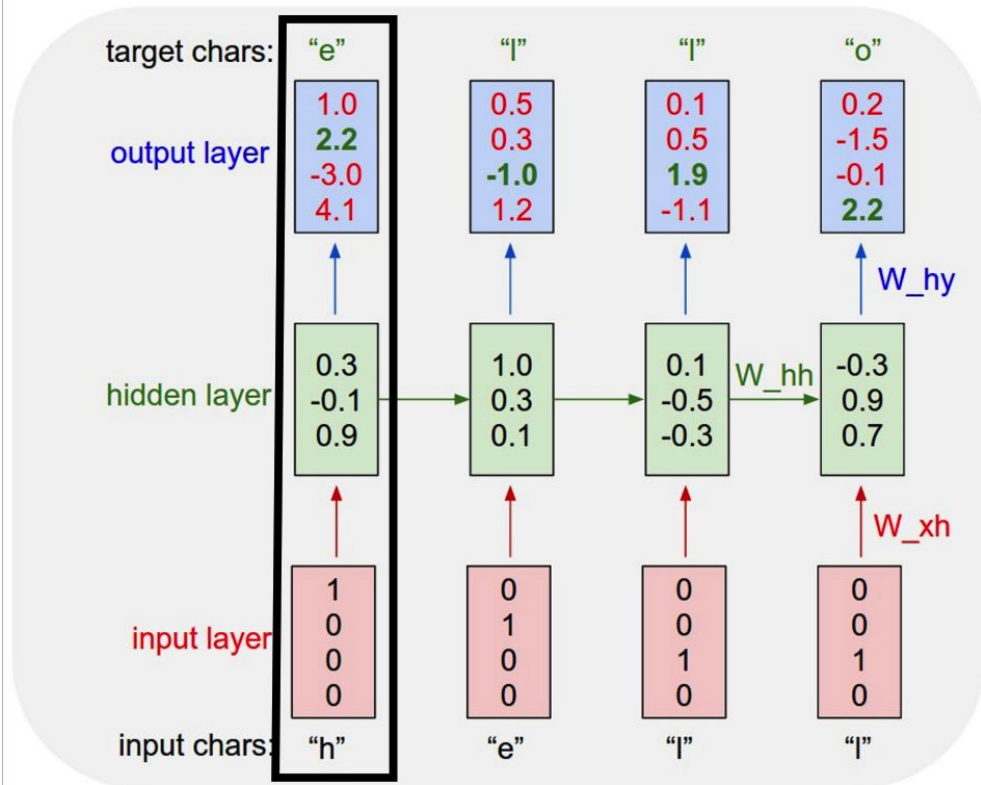
# Example: Language Modeling

Given characters 1, 2, ..., t-1,  
model predicts character t

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

Training sequence: "hello"

Vocabulary: [h, e, l, o]



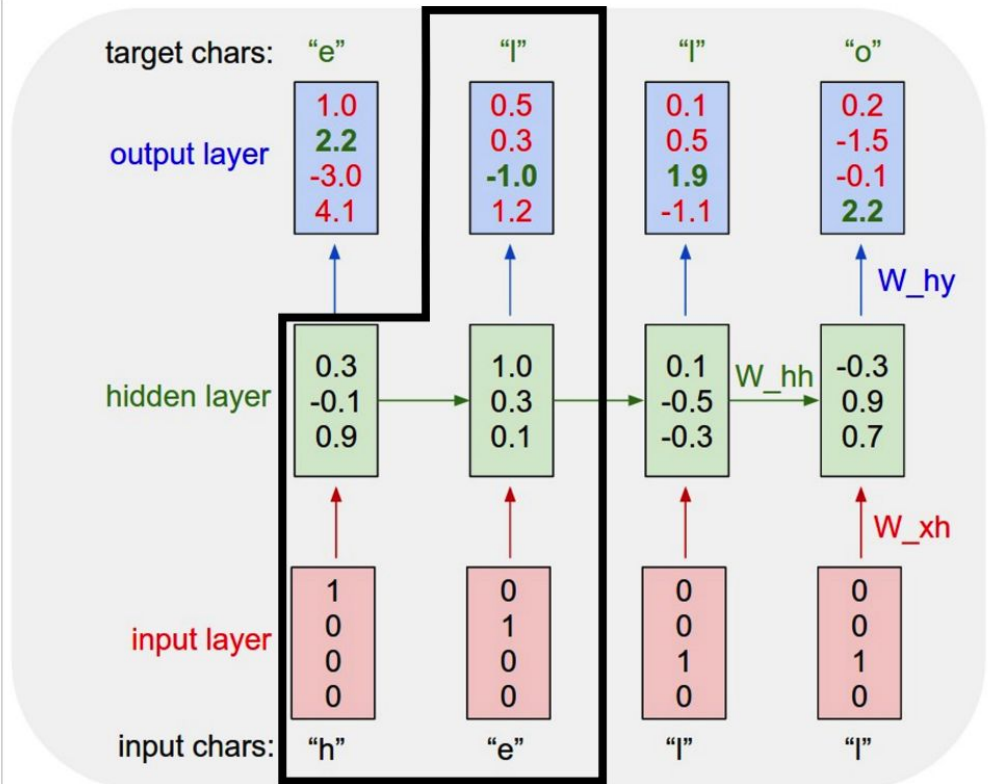
# Example: Language Modeling

Given characters 1, 2, ..., t-1,  
model predicts character t

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

Training sequence: "hello"

Vocabulary: [h, e, l, o]



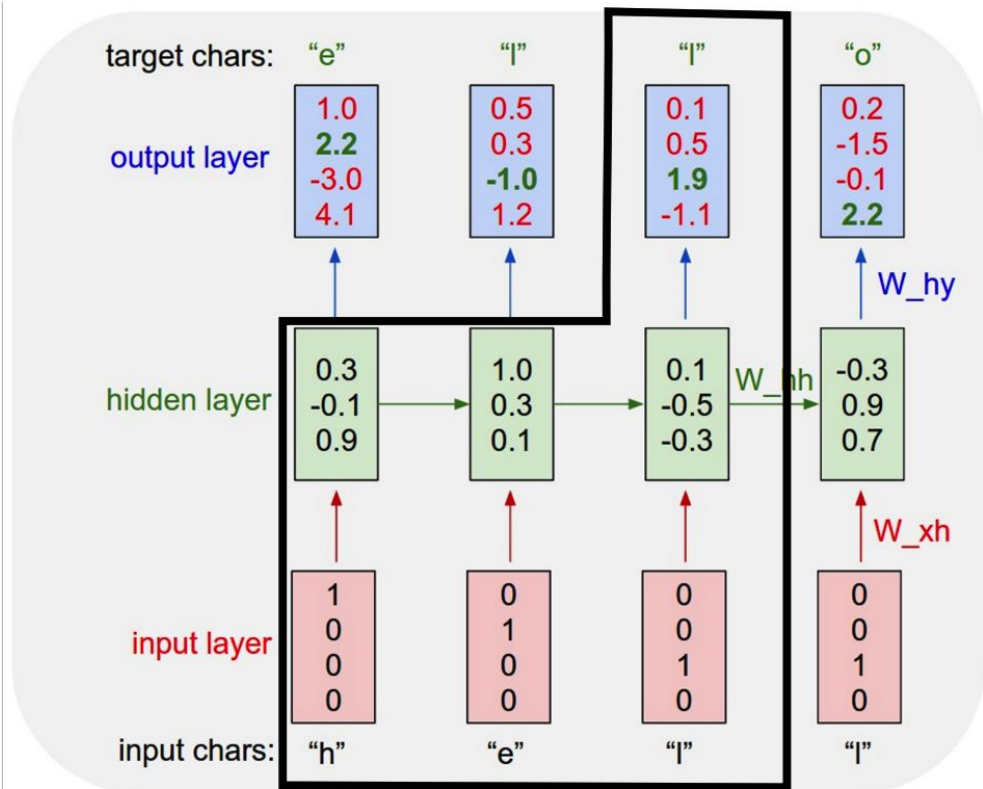
# Example: Language Modeling

Given characters 1, 2, ..., t-1,  
model predicts character t

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

Training sequence: "hello"

Vocabulary: [h, e, l, o]



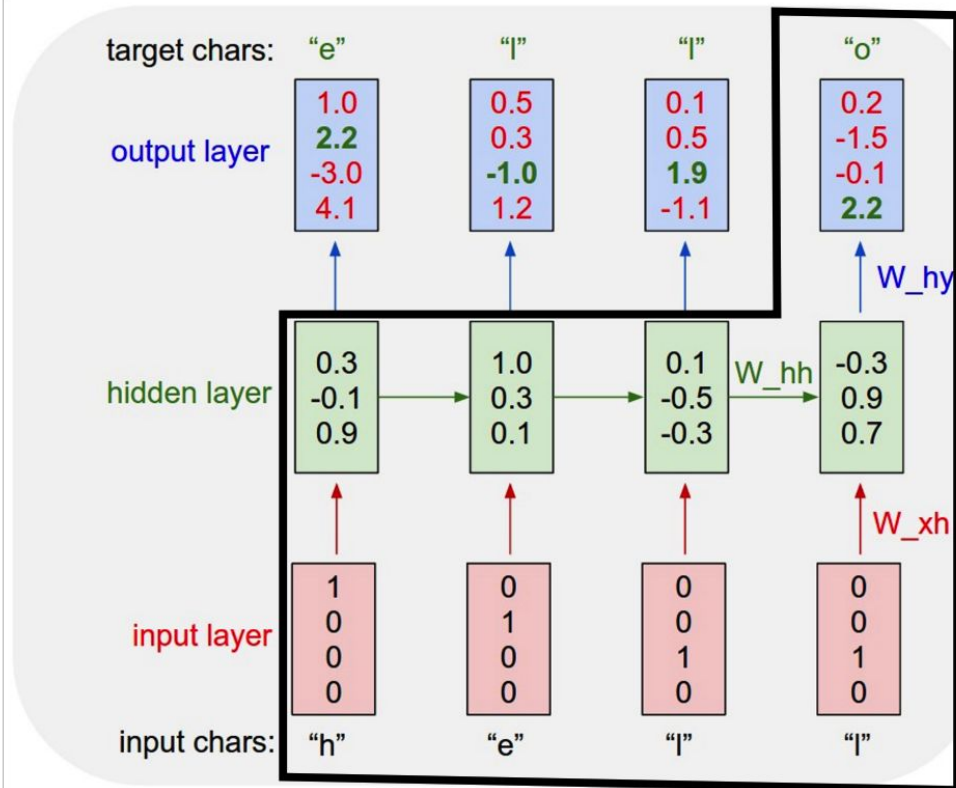
# Example: Language Modeling

Given characters 1, 2, ..., t-1,  
model predicts character t

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

Training sequence: "hello"

Vocabulary: [h, e, l, o]

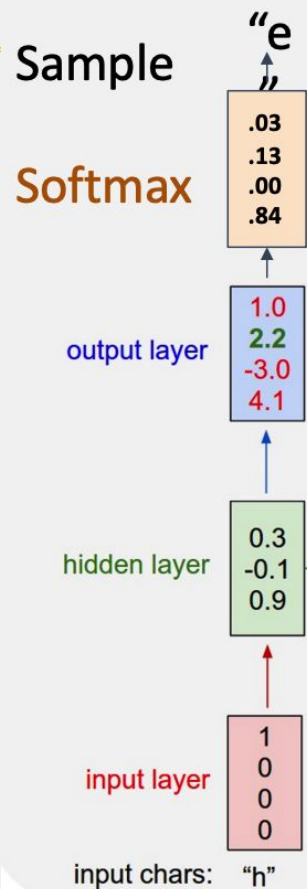


# Example: Language Modeling

At **test-time**, **generate** new text: sample characters one at a time, feed back to model

Training sequence: "hello"

Vocabulary: [h, e, l, o]

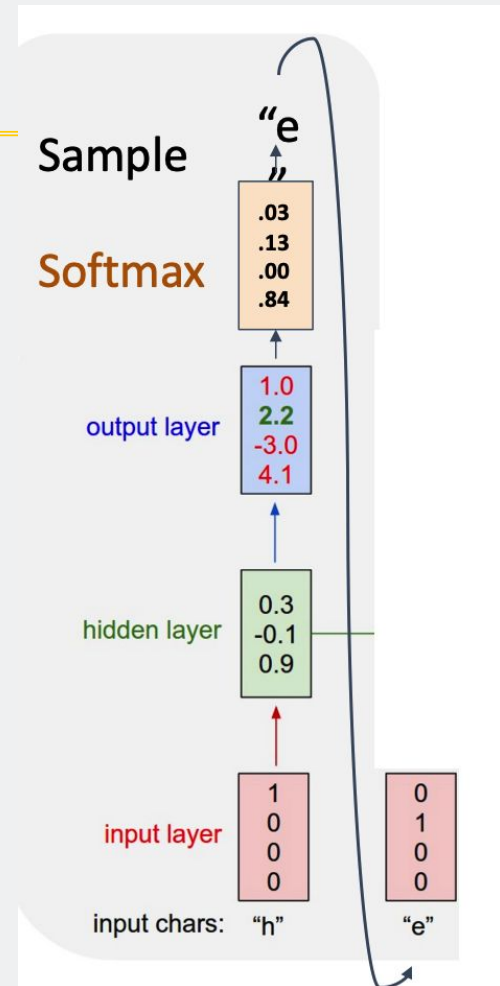


# Example: Language Modeling

At test-time, **generate** new text: sample characters one at a time, feed back to model

Training sequence: "hello"

Vocabulary: [h, e, l, o]

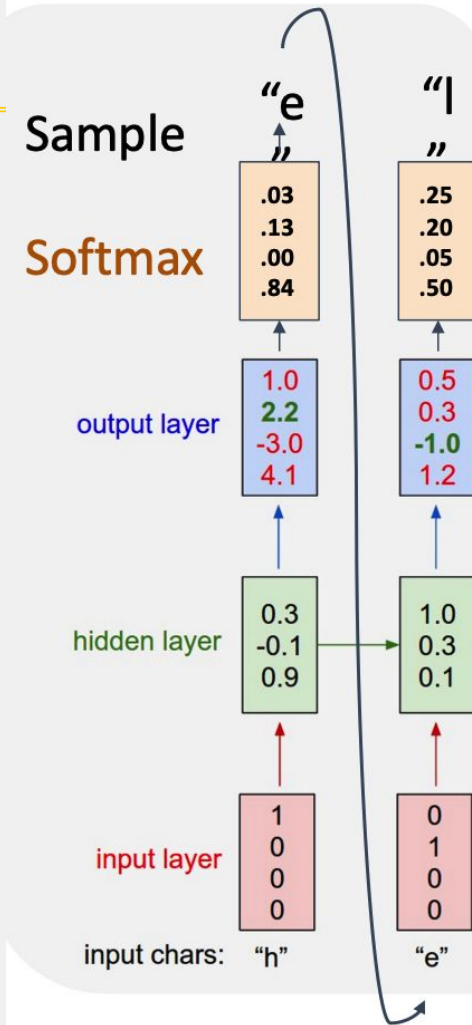


# Example: Language Modeling

At test-time, **generate** new text: sample characters one at a time, feed back to model

Training sequence: "hello"

Vocabulary: [h, e, l, o]

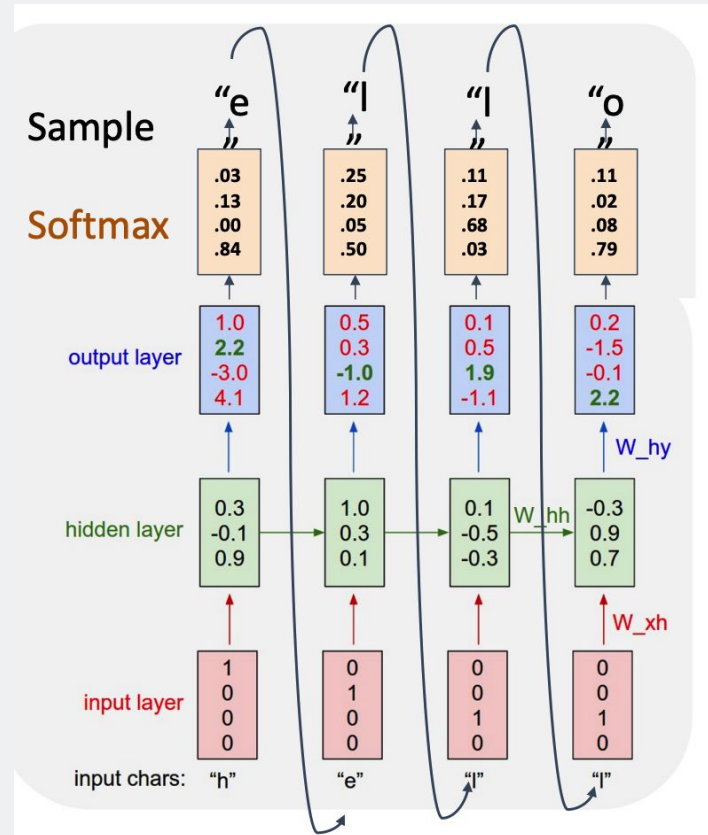


# Example: Language Modeling

At test-time, **generate** new text: sample characters one at a time, feed back to model

Training sequence: "hello"

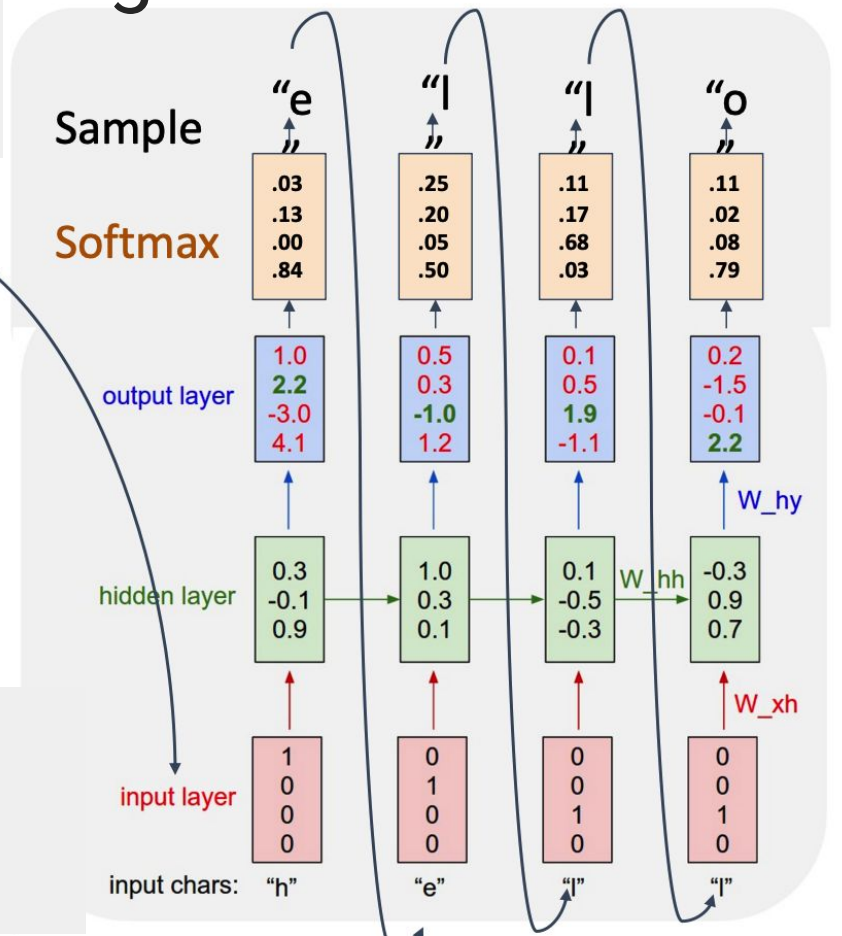
Vocabulary: [h, e, l, o]



# Example: Language Modeling

So far: encode inputs as **one-hot-vector**

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \\ w_{31} & w_{32} & w_{33} & w_{34} \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} w_{11} \\ w_{21} \\ w_{31} \end{bmatrix}$$



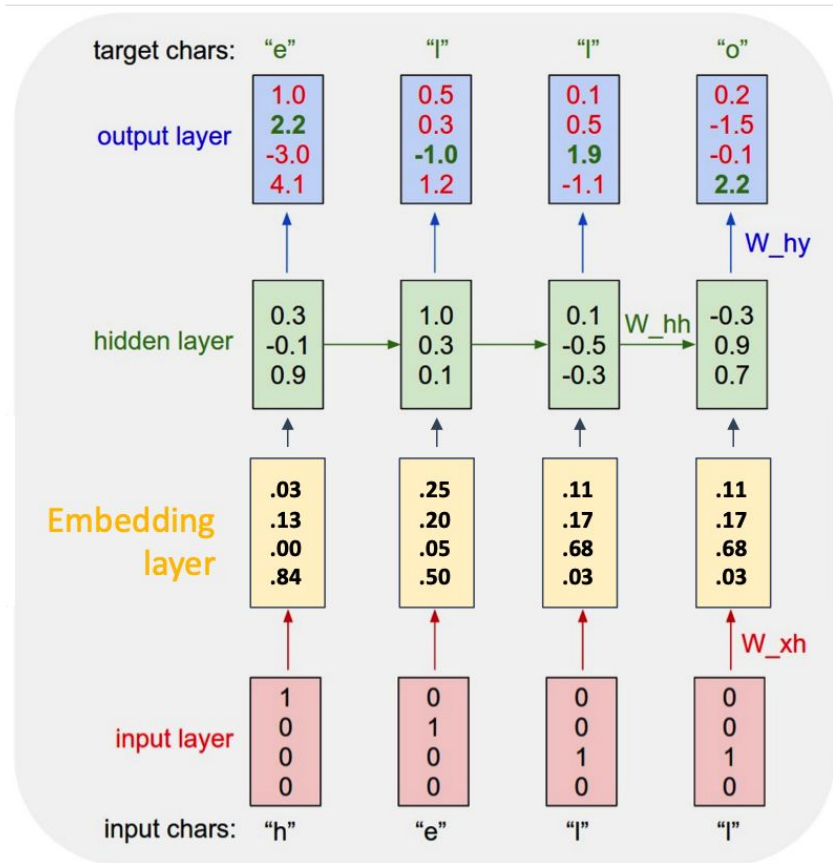
# Example: Language Modeling

\*NN likes dense real-valued vectors

So far: encode inputs  
as **one-hot-vector**

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} w_{11} \\ w_{21} \\ w_{31} \end{bmatrix}$$

Matrix multiply with a one-hot vector just extracts a column from the weight matrix. Often extract this into a separate **embedding layer**

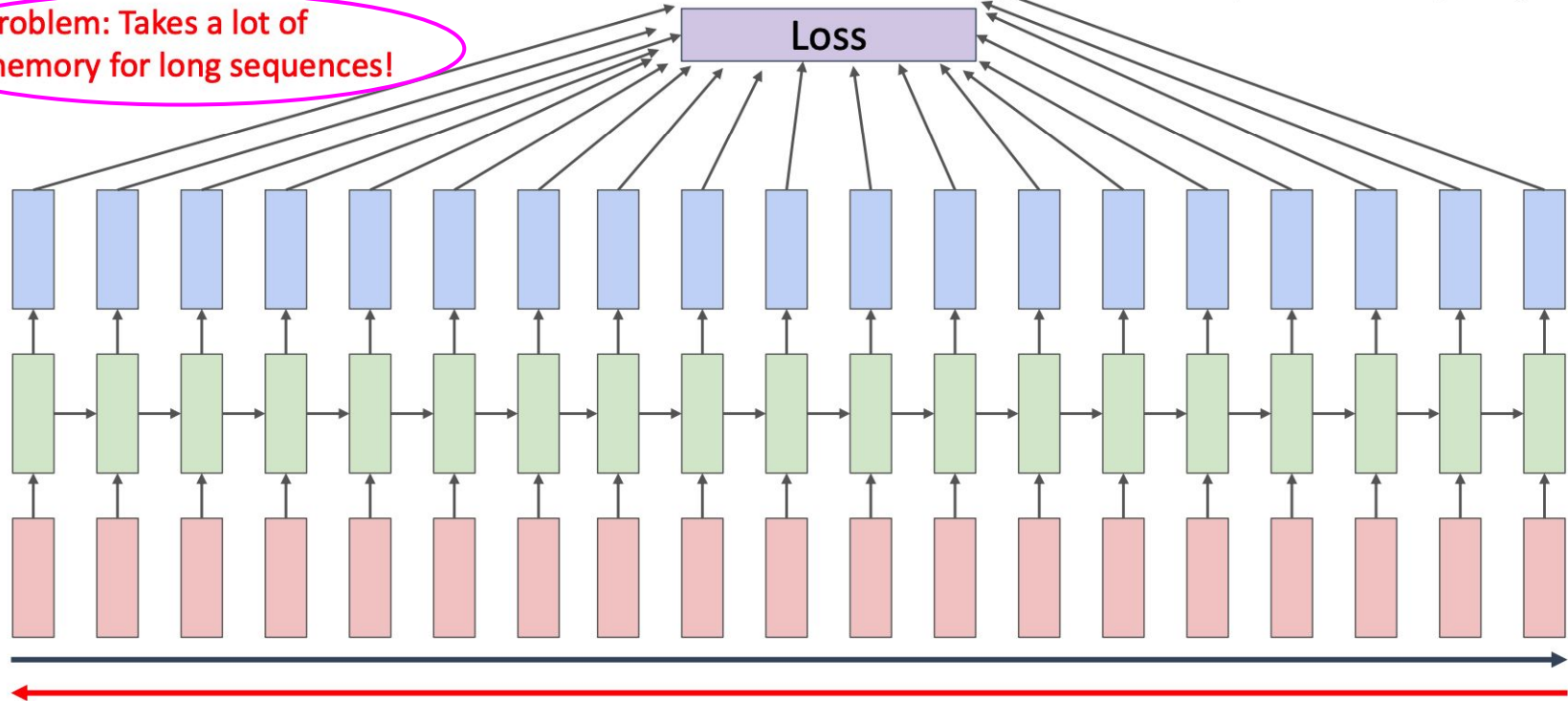


# Example: Language Modeling (Backprop)

## Backpropagation Through Time

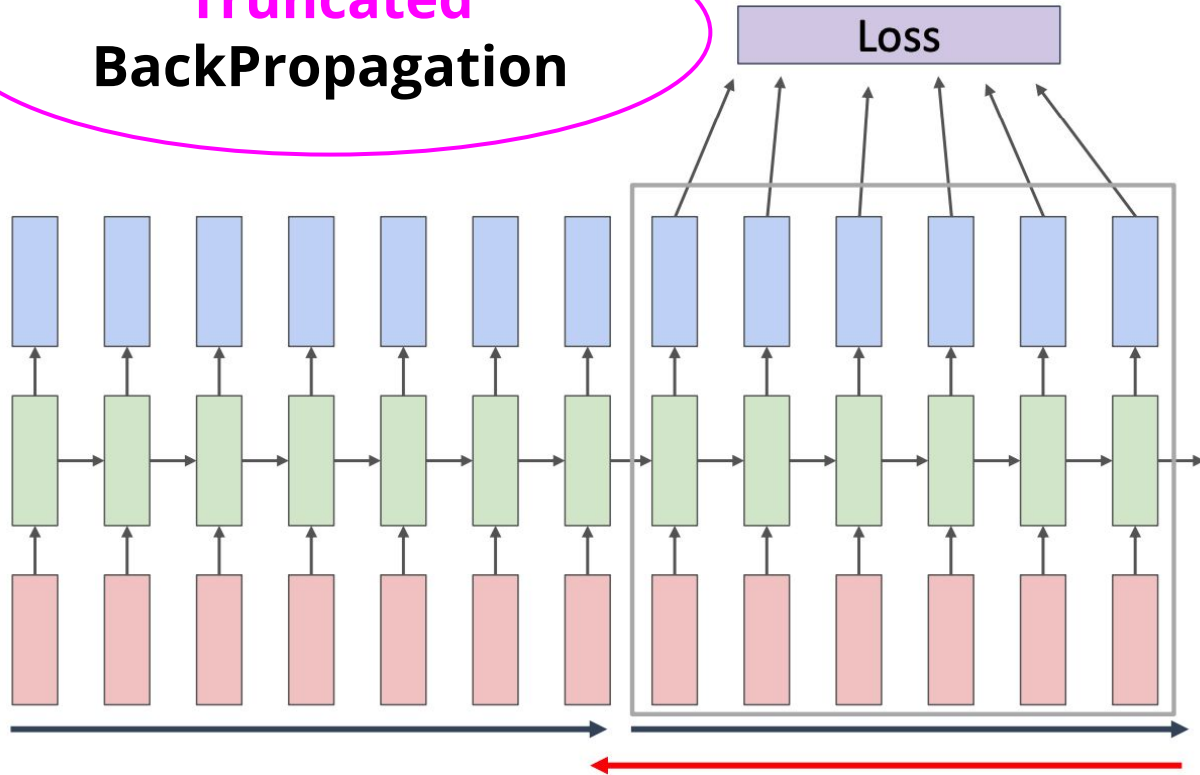
Forward through entire sequence to compute loss, then backward through entire sequence to compute gradient

Problem: Takes a lot of memory for long sequences!



# Example: Language Modeling (Backprop)

## Truncated BackPropagation



- Run forward and backward through **chunks of the sequence** instead of whole sequence
- Carry hidden states forward in time forever, but **only backpropagate for some smaller number of steps**

# Limitations of RNNs

---

- Modeling **long-range dependencies** limited by **vanishing gradient**
- **Computational and memory efficiency**, especially for long sequences
- **Parallelization** of layers that depend on **sequential** information

# Modeling Long-Range Dependencies

## One Example: Language Translation

**Input:** Sequence  $x_1, \dots, x_T$

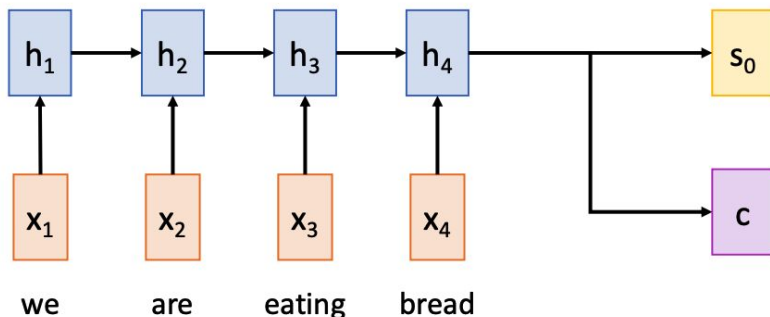
**Output:** Sequence  $y_1, \dots, y_{T'}$

**Encoder:**  $h_t = f_W(x_t, h_{t-1})$

From final hidden state predict:

**Initial decoder state**  $s_0$

**Context vector**  $c$  (often  $c=h_T$ )



Recall: Seq2Seq + RNN

# Modeling Long-Range Dependencies

## One Example: Language Translation

**Input:** Sequence  $x_1, \dots, x_T$

**Output:** Sequence  $y_1, \dots, y_T$

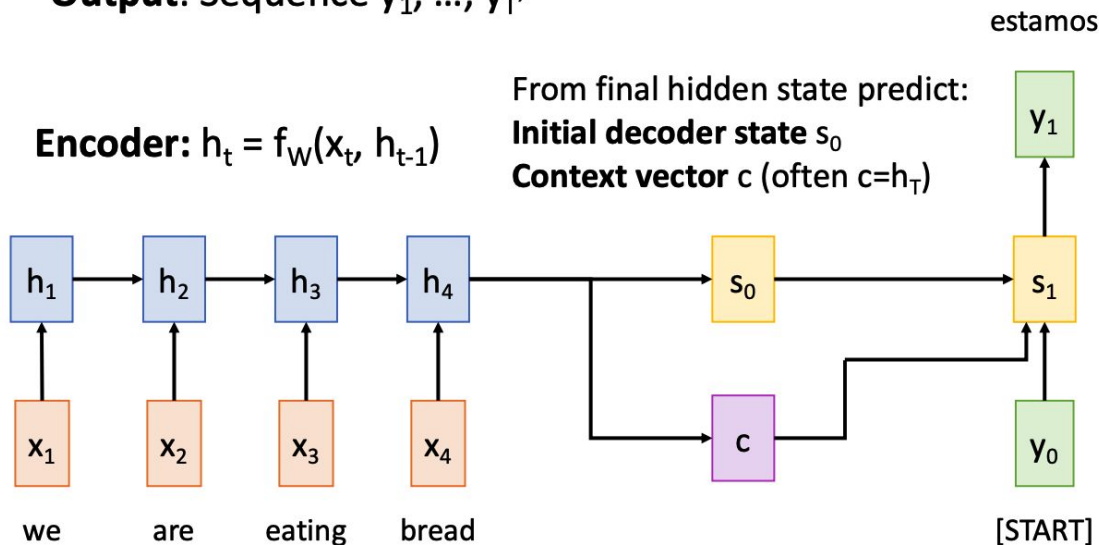
**Decoder:**  $s_t = g_U(y_{t-1}, s_{t-1}, c)$

**Encoder:**  $h_t = f_W(x_t, h_{t-1})$

From final hidden state predict:

**Initial decoder state**  $s_0$

**Context vector**  $c$  (often  $c=h_T$ )



# Modeling Long-Range Dependencies

## One Example: Language Translation

**Input:** Sequence  $x_1, \dots, x_T$

**Output:** Sequence  $y_1, \dots, y_{T'}$

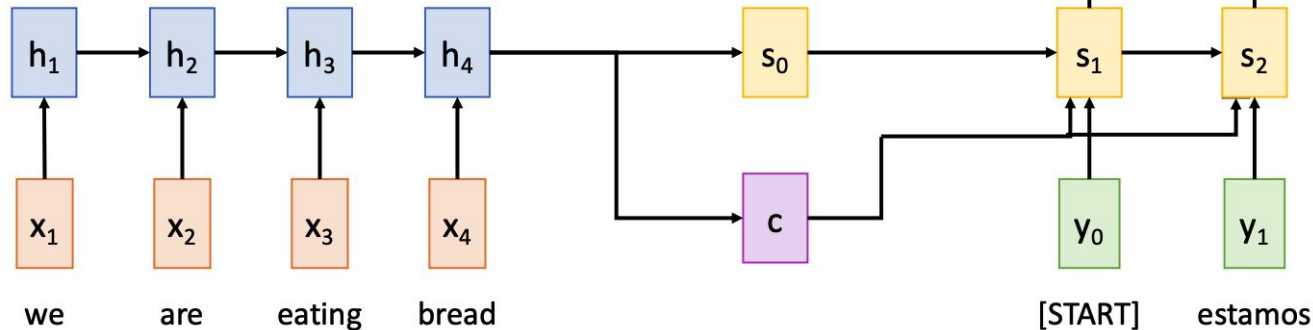
**Decoder:**  $s_t = g_U(y_{t-1}, s_{t-1}, c)$

**Encoder:**  $h_t = f_W(x_t, h_{t-1})$

From final hidden state predict:

**Initial decoder state**  $s_0$

**Context vector**  $c$  (often  $c=h_T$ )



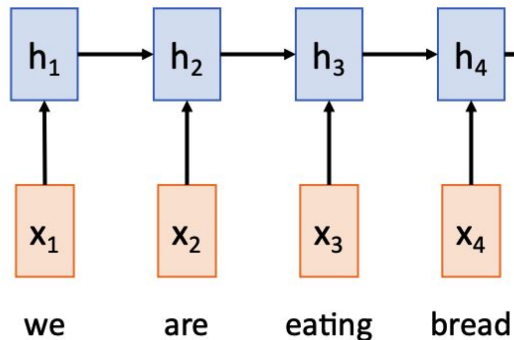
# Modeling Long-Range Dependencies

## One Example: Language Translation

**Input:** Sequence  $x_1, \dots, x_T$

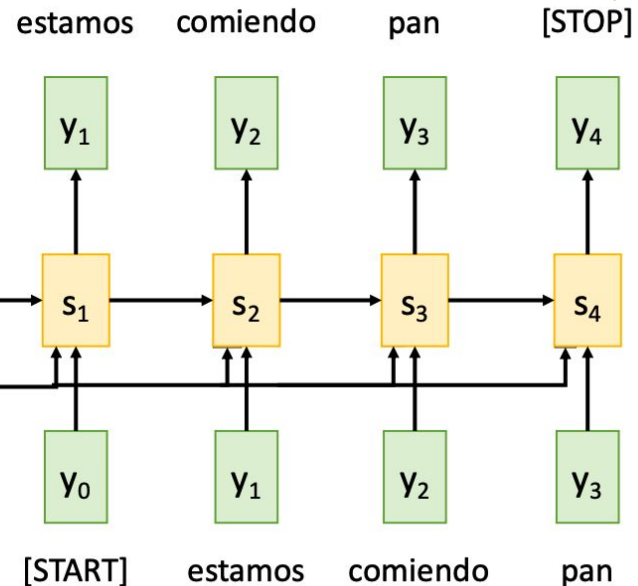
**Output:** Sequence  $y_1, \dots, y_{T'}$

**Encoder:**  $h_t = f_W(x_t, h_{t-1})$



From final hidden state predict:  
**Initial decoder state**  $s_0$   
**Context vector**  $c$  (often  $c=h_T$ )

**Decoder:**  $s_t = g_U(y_{t-1}, s_{t-1}, c)$



# Modeling Long-Range Dependencies

## One Example: Language Translation

**Input:** Sequence  $x_1, \dots, x_T$

**Output:** Sequence  $y_1, \dots, y_T$

**Encoder:**  $h_t = f_W(x_t, h_{t-1})$

From final hidden state predict:

**Initial decoder state**  $s_0$

**Context vector**  $c$  (often  $c=h_T$ )

**Decoder:**  $s_t = g_U(y_{t-1}, s_{t-1}, c)$

estamos comiendo pan [STOP]

$y_1$   $y_2$   $y_3$   $y_4$

Problem: Input sequence bottlenecked through  
fixed-sized vector. What if  $T=1000$ ?

we are eating bread

[START] estamos comiendo pan

# Modeling Long-Range Dependencies

## One Example: Language Translation

**Input:** Sequence  $x_1, \dots, x_T$

**Output:** Sequence  $y_1, \dots, y_T$

**Encoder:**  $h_t = f_W(x_t, h_{t-1})$

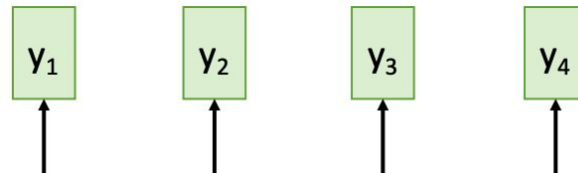
From final hidden state predict:

**Initial decoder state**  $s_0$

**Context vector**  $c$  (often  $c=h_T$ )

**Decoder:**  $s_t = g_U(y_{t-1}, s_{t-1}, c)$

estamos comiendo pan [STOP]



**Idea:** Instead of than summarizing entire input sequence into one context vector  $c$ , let's have **decoder compute its own context vector**

we are eating bread

[START] estamos comiendo pan

# Modeling Long-Range Dependencies

## One Example: Language Translation

Use a different context vector in each timestep of decoder

- Input sequence not bottlenecked through single vector
- At each timestep of decoder, context vector “looks at” different parts of the input sequence

