# ROB 498/599: Deep Learning for Robot Perception (DeepRob)

Lecture 16: Sequences (RNNs, Seq2Seq)

03/17/2025

https://deeprob.org/w25/

ROBOTICS

# Today

- Feedback and Recap (5min)
- Processing Sequences (40min)
  - RNN
  - LSTM
  - Seq2Seq
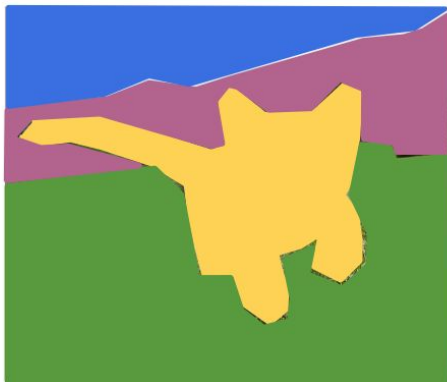- Midterm Review (~20min)
- Summary and Takeaways (5min)
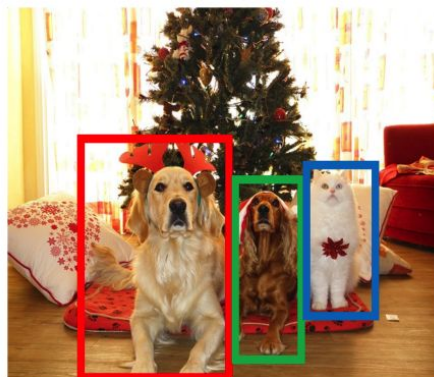
# So far…



**Classification**

CAT

No spatial extent

**Semantic Segmentation**

GRASS, CAT, TREE, SKY

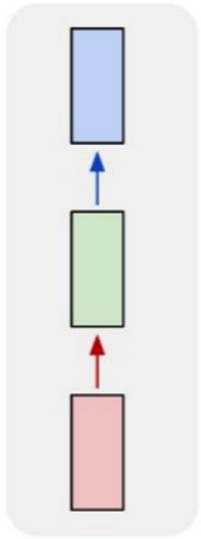No objects, just pixels

**Object Detection**

DOG, DOG, CAT

**Instance Segmentation**

DOG, DOG, CAT

Multiple Objects
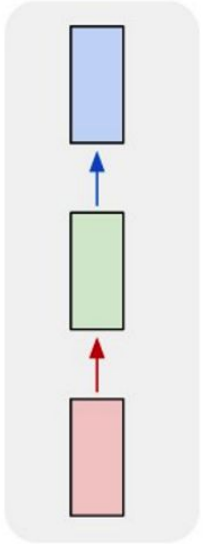
ROBOTICS

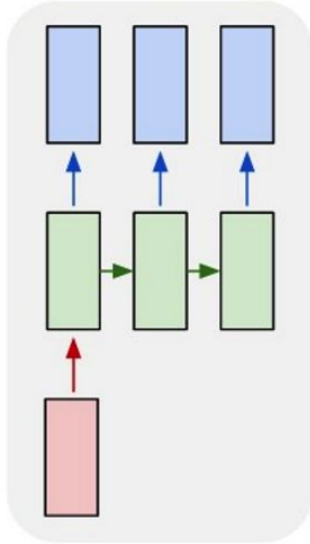# Recurrent Neural Networks



one to one

e.g. **Image classification**
Image -> Label

# Recurrent Neural Networks
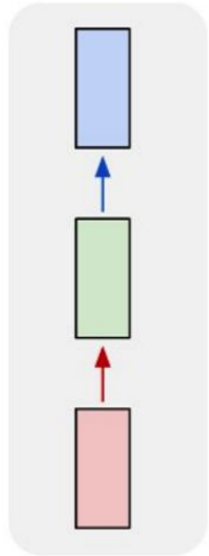


one to one    one to many

e.g. **Image Captioning**:
Image -> sequence of words

# Recurrent Neural Networks



one to one  one to many  many to one

e.g. **Video classification:**
Sequence of images -> label

ROBOTICS

# Recurrent Neural Networks



one to one     one to many     many to one     many to many

e.g. **Machine Translation:**
Sequence of words -> Sequence of words

ROBOTICS

# Recurrent Neural Networks



e.g. **Per-frame video classification:**
Sequence of images -> Sequence of labels

# 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 **x** by applying a **recurrence formula** at every time step:

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

new state — $h_t$

some function with parameters W — $f_W$

old state — $h_{t-1}$

input vector at some time step — $x_t$

y

RNN

x

ROBOTICS

# (Vanilla) Recurrent Neural Networks



The state consists of a single *"hidden"* vector **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/0BWPC

Q1: why re-use (shared) weights?

# (Vanilla) Recurrent Neural Networks: Gradient Flow



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

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

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

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994  https://ieeexplore.ieee.org/document/279181
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013 https://proceedings.mlr.press/v28/pascanu13.html

ROBOTICS

# (Vanilla) RNN: Gradient Flow

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



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

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

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

M | ROBOTICS

# (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** $\longrightarrow$ **Change RNN architecture!**

# Long Short Term Memory (LSTM)

*Long-range dependency*

**Vanilla RNN**

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

**LSTM**

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

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

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

Two vectors at each timestep:
Cell state: $c_t \in \mathbb{R}^H$
Hidden state: $h_t \in \mathbb{R}^H$

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

ROBOTICS

# 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

ROBOTICS

# 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

Input vector (x)

| x |
| h |

W

Previous hidden state (h)

sigmoid → i

sigmoid → f

sigmoid → o

tanh → g

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

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

4h x 2h    2h        4h            4*h
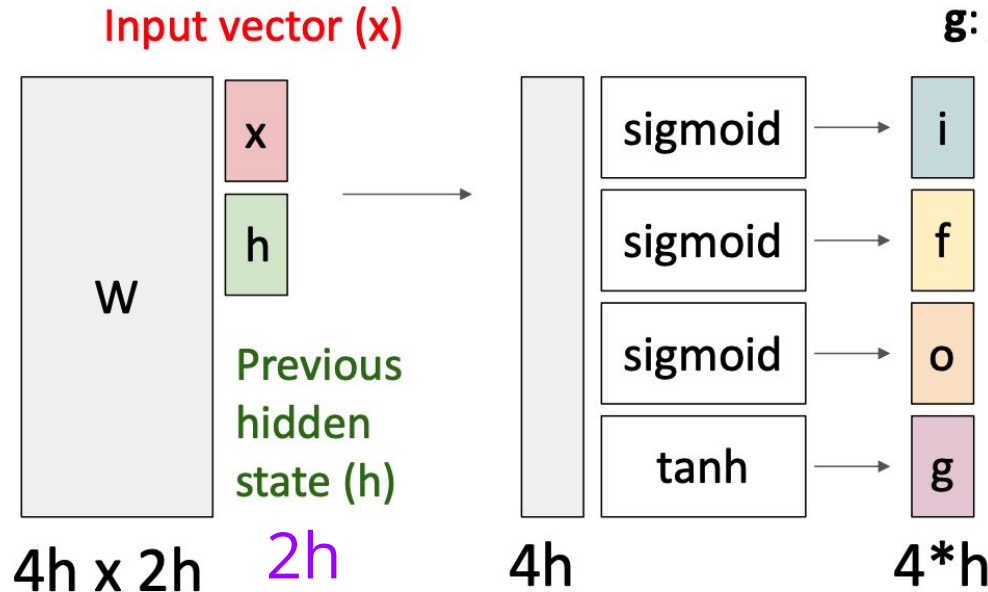
# Long Short Term Memory (LSTM)



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

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

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

M | ROBOTICS

# 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:

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

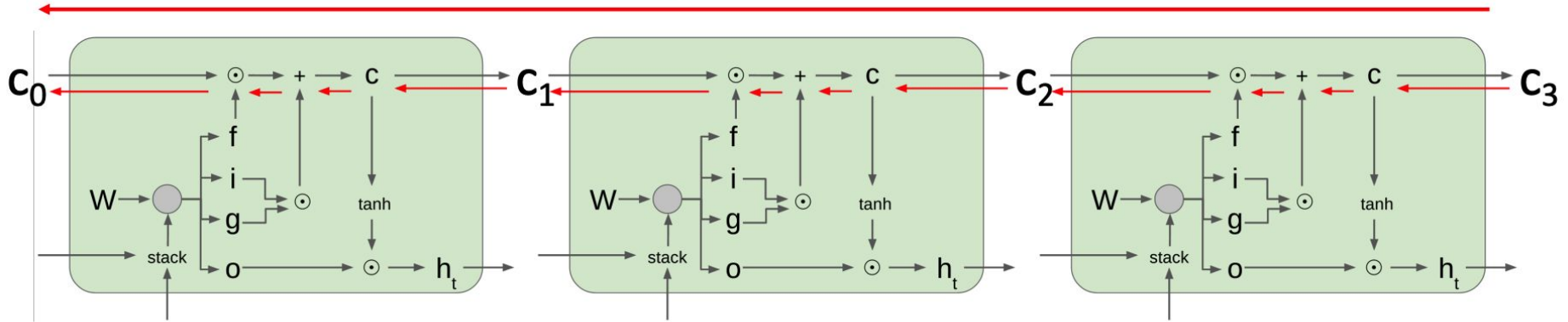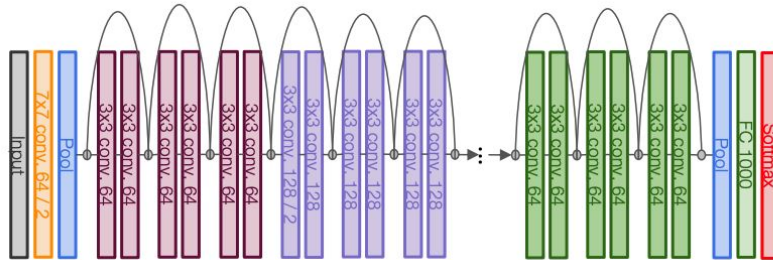$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$

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



time

# 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.

time

# Vanilla RNN - 112 Lines of Python Code

Minimal character-level language model with a Vanilla Recurrent Neural Network, in Python/numpy
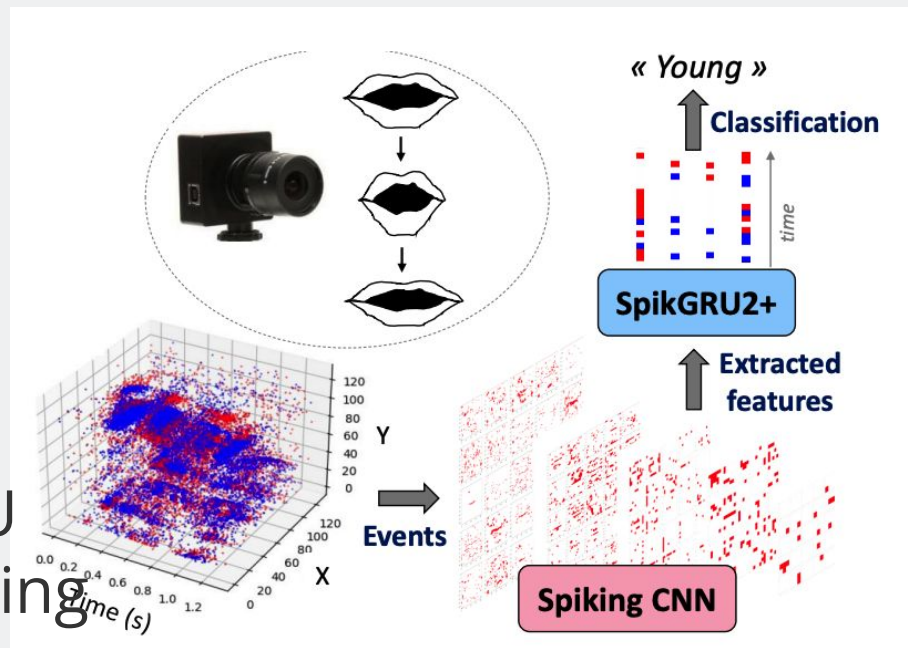
```python
min-char-rnn.py
1  """
2  Minimal character-level Vanilla RNN model. Written by Andrej Karpathy (@karpathy)
3  BSD License
4  """
5  import numpy as np
6
7  # data I/O
8  data = open('input.txt', 'r').read() # should be simple plain text file
9  chars = list(set(data))
10 data_size, vocab_size = 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/d4dee56686 7f8291f086

# Gated Recurrent Unit (GRU)

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

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

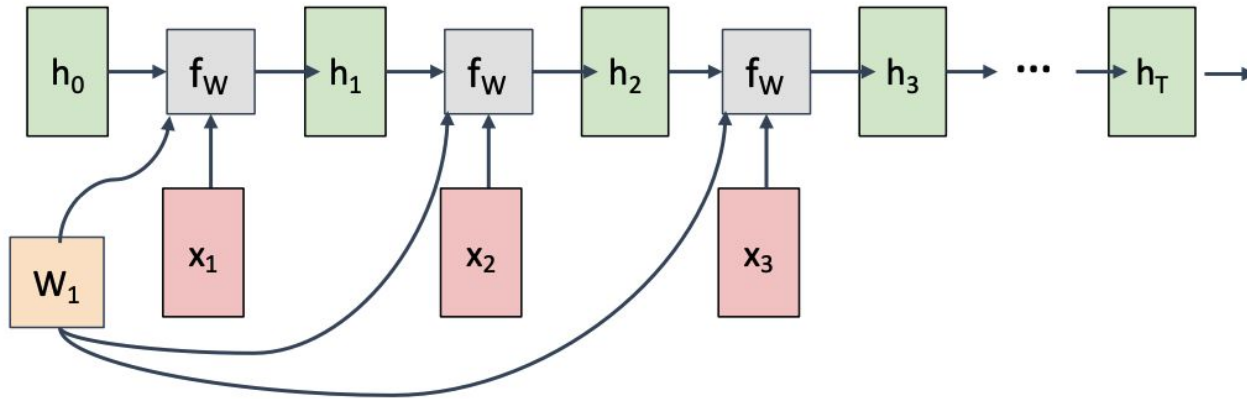https://openaccess.thecvf.com/content/CVPR2024W/EVW/papers/Dampfhoffer_Neuromorphic_Lip-Reading_with_Signed_Spiking_Gated_Recurrent_Units_CVPRW_2024_paper.pdf
Spike GPT https://arxiv.org/pdf/2302.13939

ROBOTICS

# Seq2Seq: Sequence to Sequence



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

https://proceedings.neurips.cc/paper_files/paper/2014/file/5a18e133cbf9f257297f410bb7eca942-Paper.pdf

# Seq2Seq: Sequence to Sequence

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

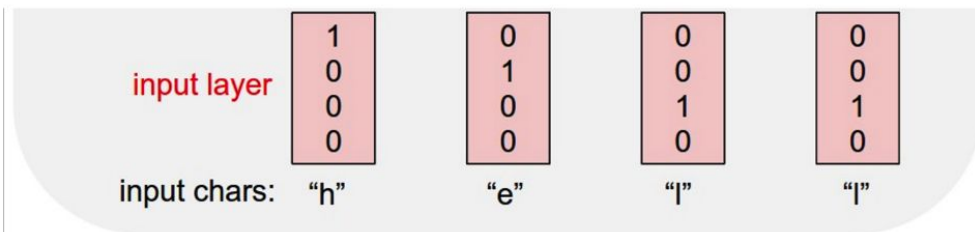**Many to one**: Encode input sequence in a single vector

ROBOTICS

# Example: Language Modeling

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

Training sequence: "hello"

Vocabulary: [h, e, l, o]



input layer

| | 1 | 0 | 0 | 0 |
|---|---|---|---|---|
| | 0 | 1 | 0 | 0 |
| | 0 | 0 | 1 | 1 |
| | 0 | 0 | 0 | 0 |

input chars:  "h"    "e"    "l"    "l"

# 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

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

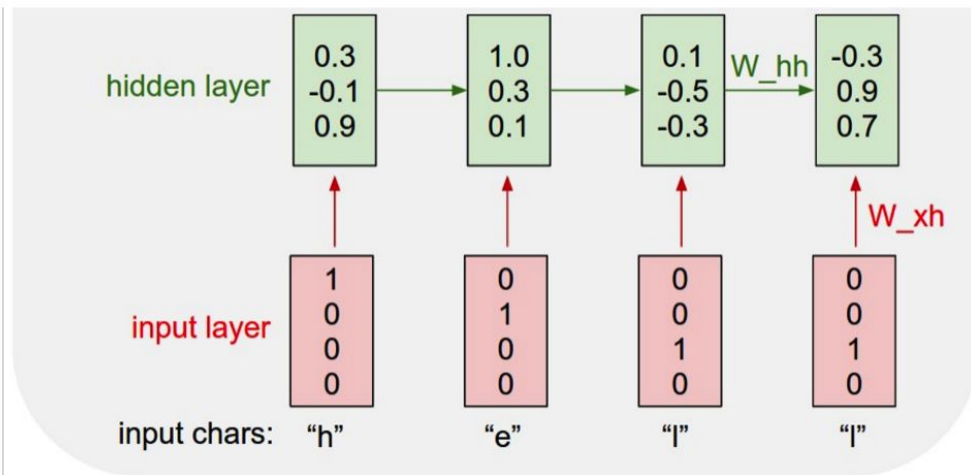Training sequence: "hello"

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"

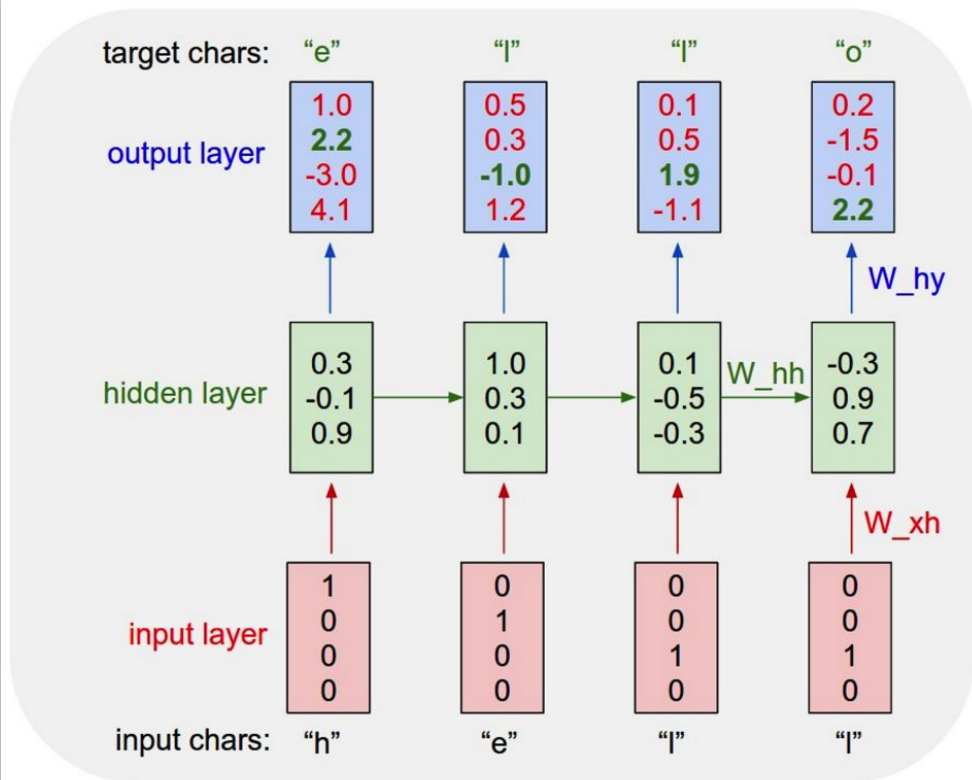Vocabulary: [h, e, l, o]

# 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

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

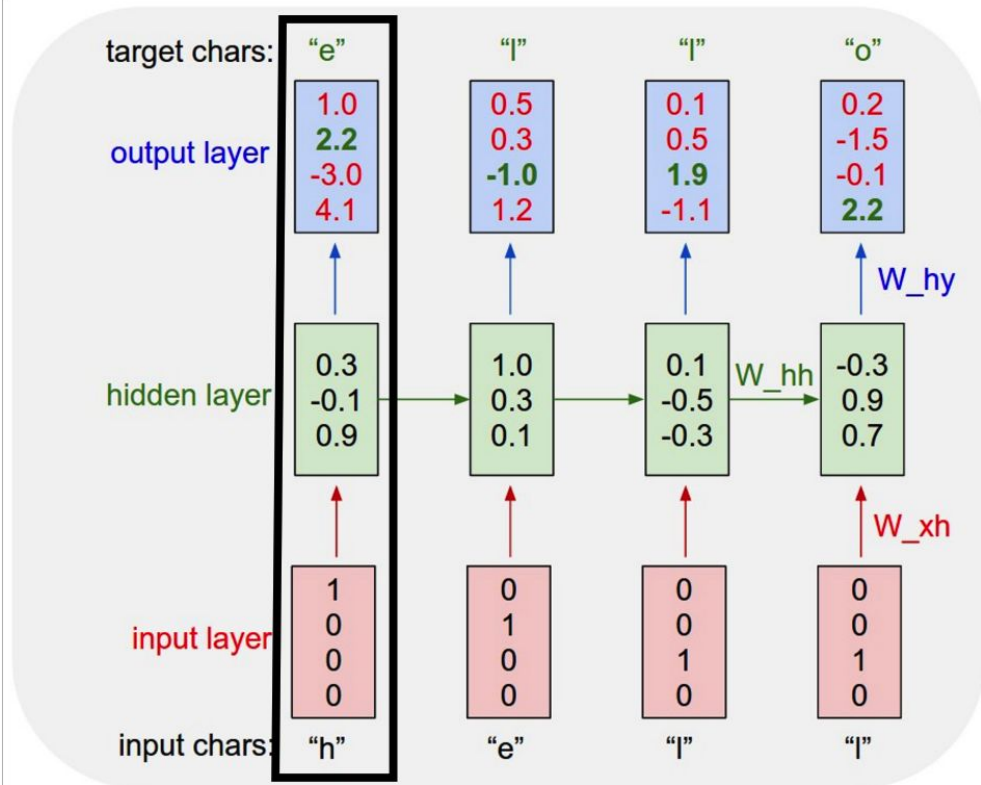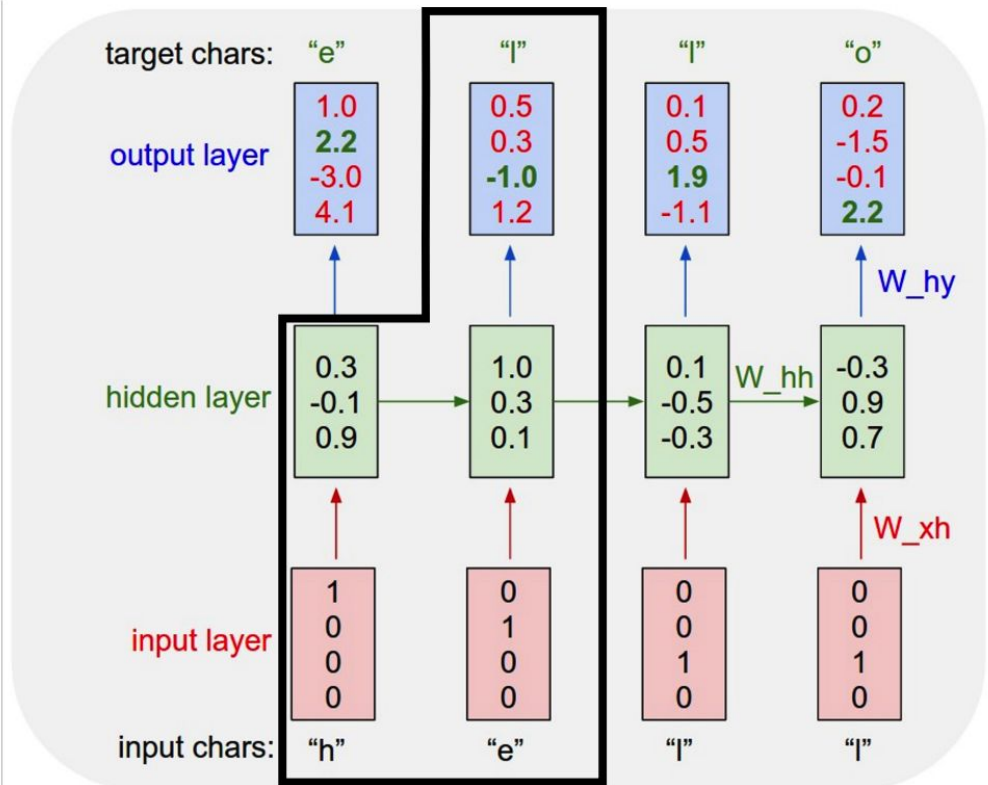Training sequence: "hello"

Vocabulary: [h, e, l, o]

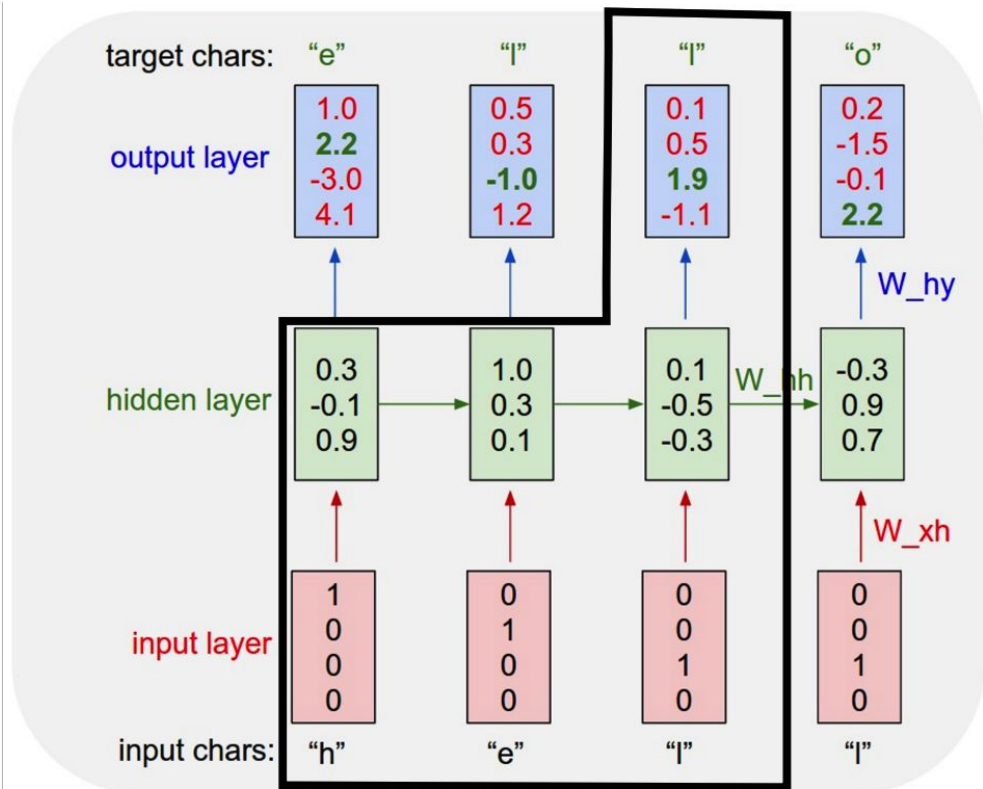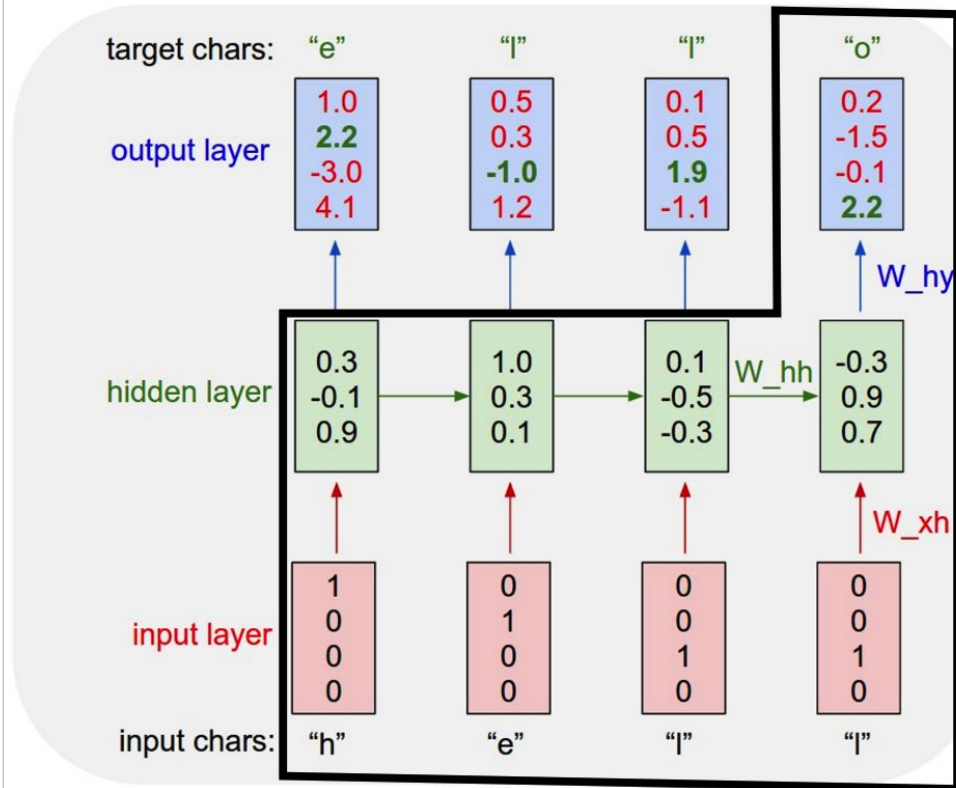# Example: Language Modeling

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]

Sample

Softmax

"e"

| |
|---|
| .03 |
| .13 |
| .00 |
| .84 |

output layer

| |
|---|
| 1.0 |
| 2.2 |
| -3.0 |
| 4.1 |

hidden layer

| |
|---|
| 0.3 |
| -0.1 |
| 0.9 |

input layer

| |
|---|
| 1 |
| 0 |
| 0 |
| 0 |

input chars:   "h"

ROBOTICS

# 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]

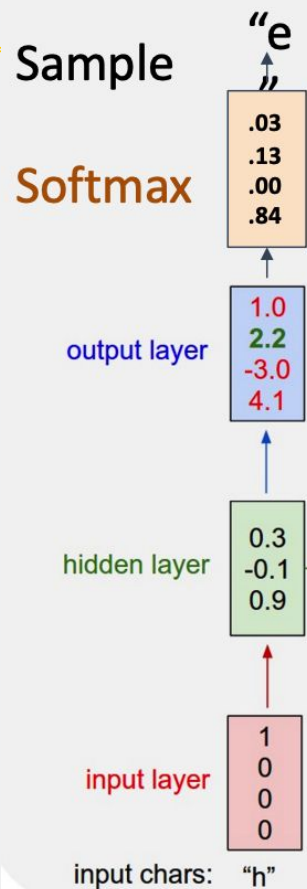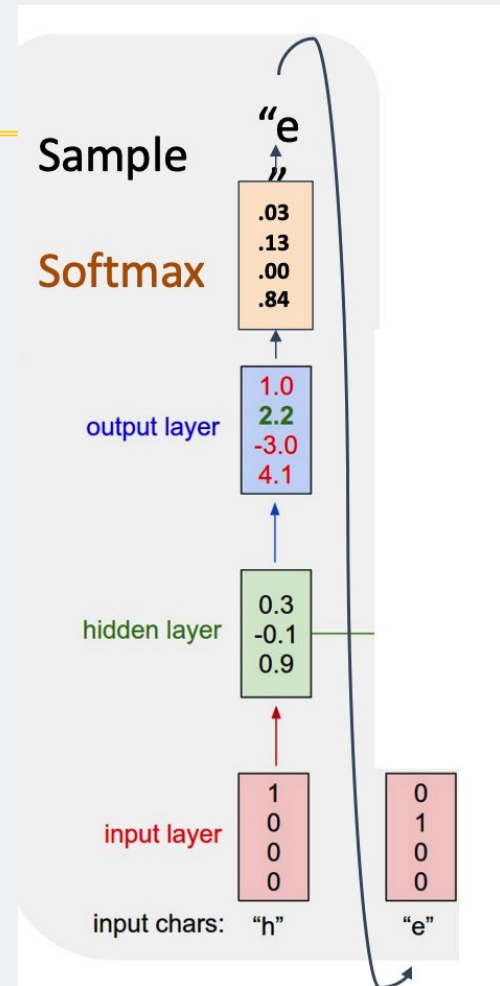# Example: Language Modeling

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

Training sequence: "hello"

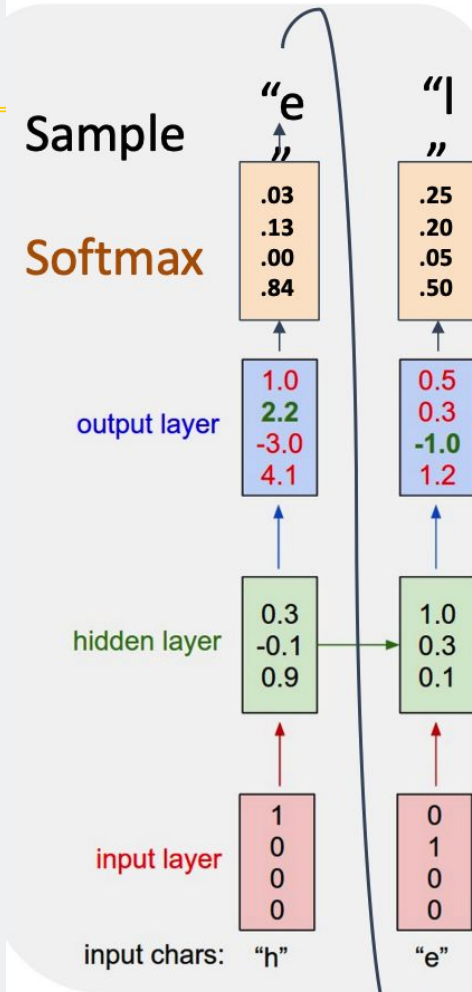Vocabulary: [h, e, l, o]

# 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

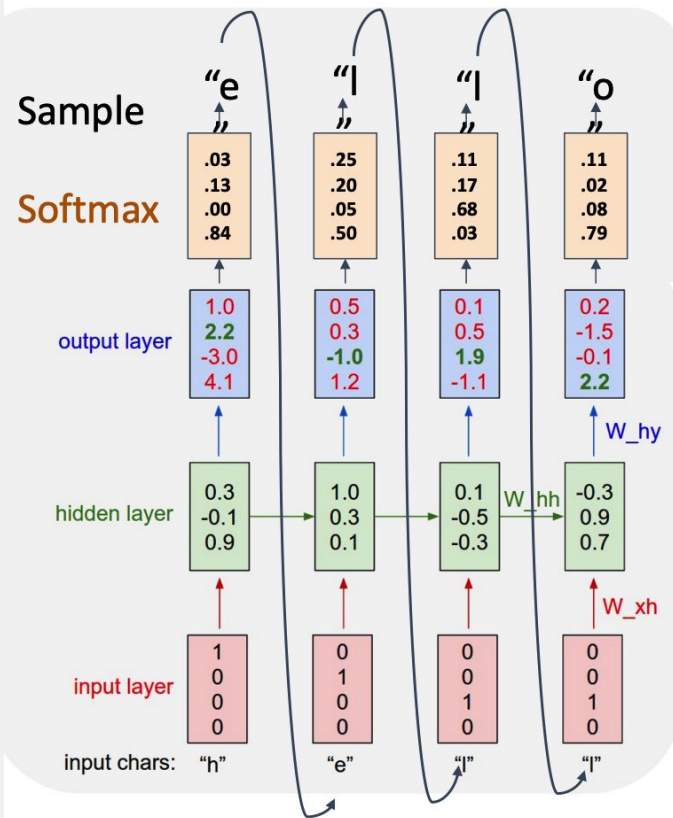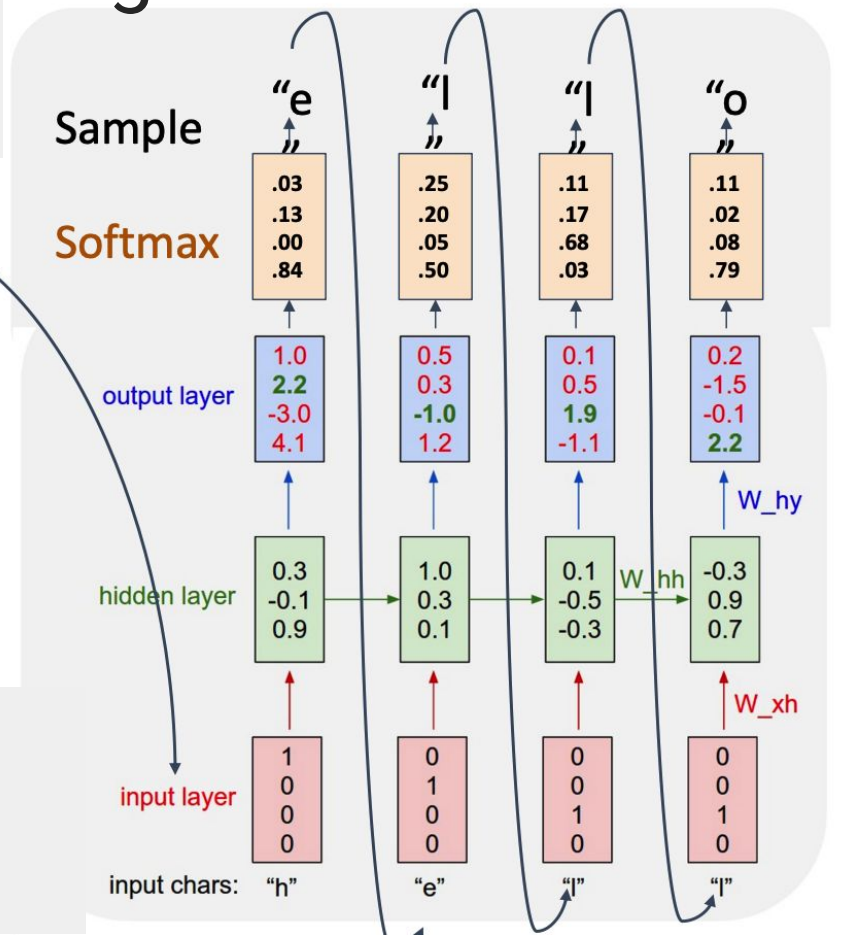So far: encode inputs
as **one-hot-vector**

$$[w_{11} \; w_{12} \; w_{13} \; w_{14}] \; [1] \qquad [w_{11}]$$
$$[w_{21} \; w_{22} \; w_{23} \; w_{14}] \; [0] \; = \; [w_{21}]$$
$$[w_{31} \; w_{32} \; w_{33} \; w_{14}] \; [0] \qquad [w_{31}]$$
$$[0]$$

# Example: Language Modeling

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

$$[w_{11} \ w_{12} \ w_{13} \ w_{14}] \ [1] \qquad [w_{11}]$$
$$[w_{21} \ w_{22} \ w_{23} \ w_{14}] \ [0] \ = \ [w_{21}]$$
$$[w_{31} \ w_{32} \ w_{33} \ w_{14}] \ [0] \qquad [w_{31}]$$
$$[0]$$

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!

Loss

# 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

# MidTerm, P4, Final Project

# P4 Reminder

https://deeprob.org/w25/projects/project4/
Due March 30, 2025

Two parts:
1. PoseCNN (see Lecture 13 for more hints)
2. Vision Transformer (this week)

*Start NOW!!!*

ROBOTICS

# Final Project Teams Assigned

https://docs.google.com/spreadsheets/d/1FjWAjJ8p26xZmZaqsW4Iew8H4iKe0FA78Q30eZ38g7A/edit?usp=sharing

Next TO-DOs:  (final project total - 23% grade)
April 1st, 5-min poster "lightning talk", 5% grade

18% grade

April 22nd, final project showcase @FRB atrium
April 28th, final project (report, code, video/website) DUE

# Final Project Teams Assigned

Reminder: <span style="color:magenta">Group Collaboration</span> Policy (refer to Course Information Document)

- **"I participated and contributed to team discussions on each problem, and I attest to the integrity of each solution. Our team met as a group on [DATE(S)]. "**
- **"Contribution of Authors: [Team member A] did [Task XXX]; [Team members B and C] did [Task YYY]; [Team members A, B and C] did [ZZZ]. [All authors] [gave feedback on the software development, contributed to writing the report/making the demo presentation, and approved the final version for submission.]"** (*Modify the texts in brackets according to your specific team situation and member contribution. Ideally, each member/subset of members contributed to something unique, and all authors contributed to giving feedback and writing/making the final report/demo/presentations and approving the final version for submission.)

⟹ Set and <u>assign smaller goals to each person</u> early!!!

<span style="color:magenta">GenAI</span>: Permitted with disclosure (specify prompt, platform, results). Suggest <u>verify</u> and <u>edit</u> on top of GenAI results (if using).

https://genai.umich.edu/

**M | ROBOTICS**