



# WOMEN IN AI PANEL

MCGILL AI SOCIETY X MCWICS

MAR 14 | 5:30-7:30PM | MCENG 304

PROF. PINEAU | L. ASRI | A. DURAND | N. ROSTAMZADEH



# **COMP 551 - Applied Machine Learning**

## **Lecture 16 – Recurrent neural networks**

---

William L. Hamilton

(with slides and content from Joelle Pineau and Ryan Lowe)

\* Unless otherwise noted, all material posted for this course are copyright of the instructor, and cannot be reused or reposted without the instructor's written permission.

# Midterm

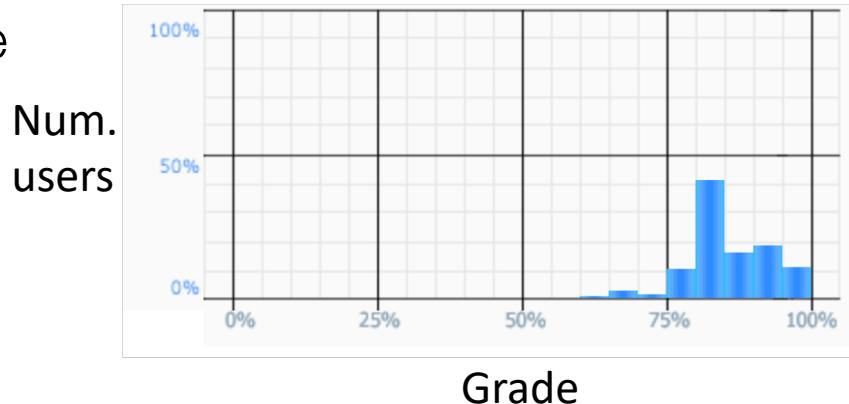
---

- I released practice questions. (See the announcement on MyCourses).
  - Without solutions: [https://www.cs.mcgill.ca/~wlh/comp551/files/midterm\\_practice\\_nosols.pdf](https://www.cs.mcgill.ca/~wlh/comp551/files/midterm_practice_nosols.pdf)
  - With solutions: [https://www.cs.mcgill.ca/~wlh/comp551/files/midterm\\_practice\\_sols.pdf](https://www.cs.mcgill.ca/~wlh/comp551/files/midterm_practice_sols.pdf)
- The midterm is 6-8pm on March 26<sup>th</sup>.
  - We will send room assignments soon!
- The midterm itself will be 100 minutes.
- Clarifications:
  - Allowed to bring one double-sided **handwritten** page of notes (A4 or letter size).
  - Allowed to bring a non-programmable calculator.
- There will be long and short answer (like the practice questions) as well as multiple choice (like the quizzes).

# MiniProject 2

---

- Grades are out!
- The class did well: 86% average



- We released examples of excellent reports (see the MyCourses announcement) to help everyone improve.

# MiniProject 2

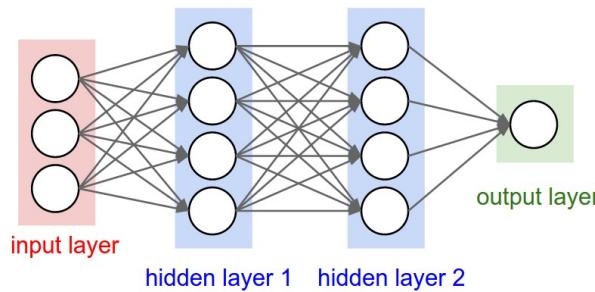
---

- Top-performing groups used:
  - Stacking ensembles
  - Naïve Bayes SVM (NB-SVM, <https://www.aclweb.org/anthology/P12-2018>)
  - Pre-trained deep learning models (e.g., BERT)
- The professor's benchmark was the NB-SVM model with unigrams, bigrams, and NLTK tokenization.

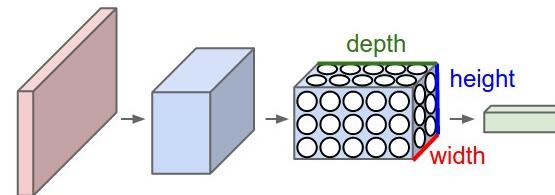
# Last time: Convolutional Neural Networks

---

Feedforward network



Convolutional neural network (CNN)



- **CNN characteristics:**

- Input is usually a 3D tensor: 2D image  $\times$  3 colours
- Each layer transforms an input 3D tensor to an output 3D tensor using a differentiable function.

# Tip #1: Dropout regularization

---

- **Goal:** Learn model that generalizes well, robust to variability.
- **Method:** Independently set each hidden unit activity to zero with probability  $p$  (usually  $p=0.5$  works best).
- **Effect:** Can greatly reduce overfitting.



# Tip #2: Batch normalization

---

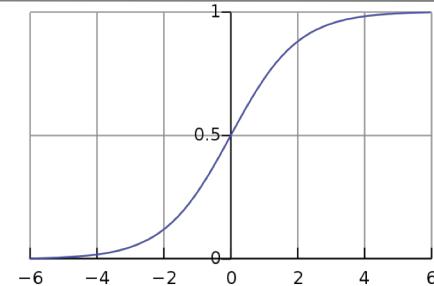
- Idea: Feature scaling/normalization makes gradient descent easier.
  - We already apply this at the input layer; extend to other layers.
  - Use empirical batch statistics to choose re-scaling parameters.
- For each mini-batch of data, at each layer  $k$  of the network:
  - Compute empirical mean and var independently for each dimension
  - Normalize each input: 
$$\hat{x}^{(k)} = \frac{x^{(k)} - E[x^{(k)}]}{\sqrt{VAR[x^{(k)}]}}$$
  - Output has tunable parameters  $(\gamma, \beta)$  for each layer: 
$$y^k = \gamma^k \cdot \hat{x}^{(k)} + \beta^k$$
- Effect: More stable gradient estimates, especially for deep networks.

# Tip #3: softmax vs. sigmoid

- We have discussed the **sigmoid function** for binary classification

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

- Maps a real-valued scalar to probability.
- Combine with cross-entropy loss to train binary classification models.



- It can be extended to the **softmax function** for multiclass classification

$$\text{softmax}(\mathbf{z}, k) = \frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}}$$

- Maps a real-valued K-dimensional vector to a K-way categorical distribution.
- Combine with cross-entropy loss to train K-way classification models

# Tip #3: softmax vs. sigmoid

---

$$P(\hat{y} = k) = \text{softmax}(\mathbf{z}, k) = \frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}}$$

Probability that output belongs to class  $k$ .

K-dimensional real-valued vector.

Index of a particular dimension in  $\mathbf{z}$  (i.e.,  $k \in \{1, \dots, K\}$ )

$k$ -th entry of  $\mathbf{z}$ .

Sum over all entries of  $\mathbf{z}$  (exponentiated)

# Tip #3: softmax vs. sigmoid

---

- Sigmoid is essentially a special case of the softmax with only two classes:

$$\sigma(z) = \frac{1}{1 + e^{-z}} = \frac{e^z}{1 + e^z} = \frac{e^z}{e^z + e^0}$$

- Cross entropy can be generalized to the multiclass/softmax setting:

$$\begin{aligned}\text{softmax-cross-entropy}(y, \mathbf{z}) &= -\log(P(\hat{y} = y)) \\ &= -\log(\text{softmax}(\mathbf{z}, y))\end{aligned}$$

# Tip #3: softmax vs. sigmoid

---

- Multi-class classification (K-dimensional softmax) is not the same as multi-output classification (K-independent sigmoid functions).
- Use softmax/multiclass when every point belongs to one of K-classes.
  - E.g., standard MNIST digit classification
  - More common scenario.
- Use K independent sigmoids when there are multiple labels for each point.
  - E.g., suppose we are detecting shapes in an images that contain multiple shapes

# Major paradigms for deep learning

---

- **Deep neural networks:** The model should be interpreted as a computation graph.
  - **Supervised training:** E.g., feedforward neural networks.
  - **Unsupervised training (later in the course):** E.g., autoencoders.
- Special architectures for different problem domains.
  - Computer vision => Convolutional neural nets.
  - Text and speech => Recurrent neural nets.

# Major paradigms for deep learning

---

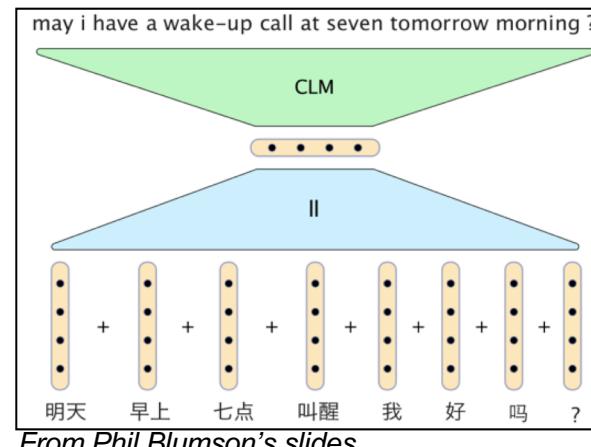
- Deep neural networks: The model should be interpreted as a computation graph.
  - Supervised training: E.g., feedforward neural networks.
  - Unsupervised training (later in the course): E.g., autoencoders.
- Special architectures for different problem domains.
  - Computer vision => Convolutional neural nets.
  - Text and speech => Recurrent neural nets.

# Neural models for sequences

- Many datasets contain **sequences** of data (e.g. time-series, text)
- How could we process sequences with a **feed-forward neural network**?
  1. Take vectors representing the last N timesteps and **concatenate** (join) them
  2. Take vectors representing the last N timesteps and **average** them

Example: machine translation

- Input: sequence of words
- Output: sequence of words



# Neural models for sequences

---

- Problem: these approaches **don't exploit the sequential nature of the data!!**
- Also, they can only consider information from a fixed-size context window
- **Temporal information is very important in sequences!!**
- E.g. machine translation:

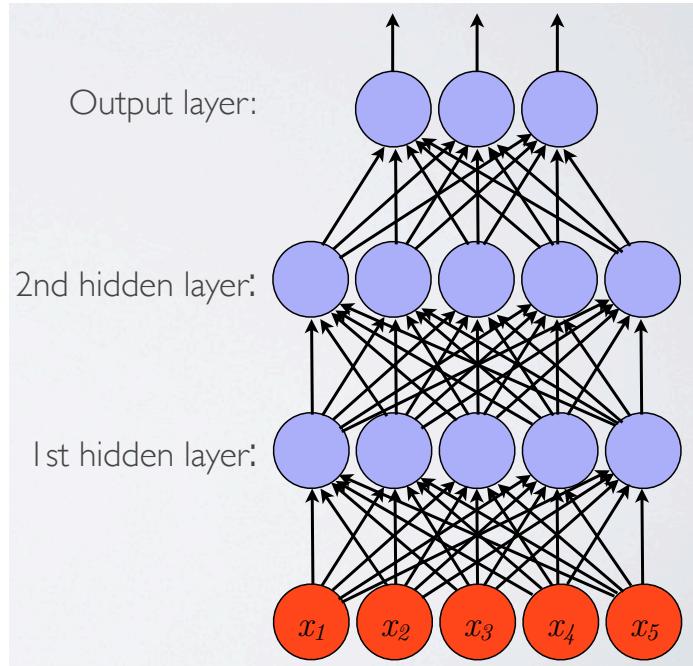
“John hit Steve on the head with a bat”

$\neq$  “Steve hit John on the bat with a head”

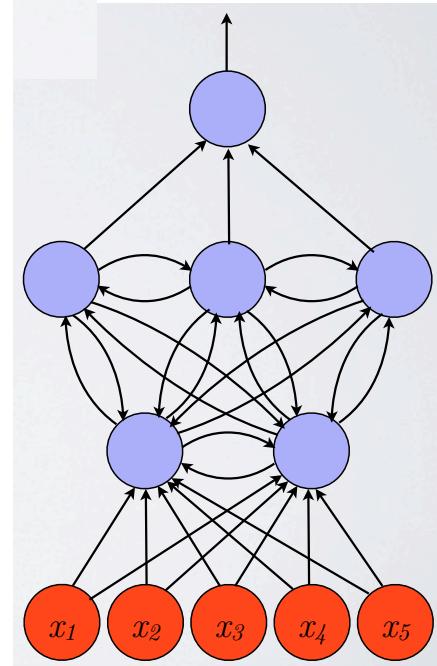
$\neq$  “Bat hit a with head on the John Steve”

# Recurrent Neural Networks (RNNs)

Feed-forward neural net

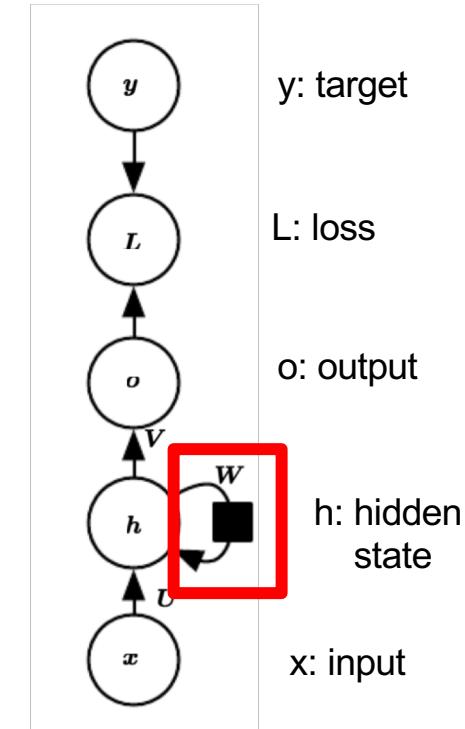


Add cycles in network



# Recurrent Neural Networks (RNNs)

- What kind of cycles?
- Cycles with a **time delay**
- Box means that the information is sent at the next time step (no infinite loops)

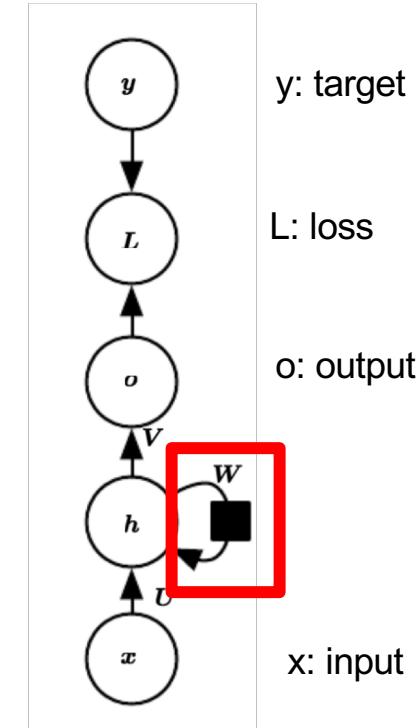


# Recurrent Neural Networks (RNNs)

- What does this allow us to do?
- Can view RNNs as having a **hidden state  $h_t$**  that changes over time.
- $h_t$  represents “useful information” from past inputs.
- A standard/simple RNN:

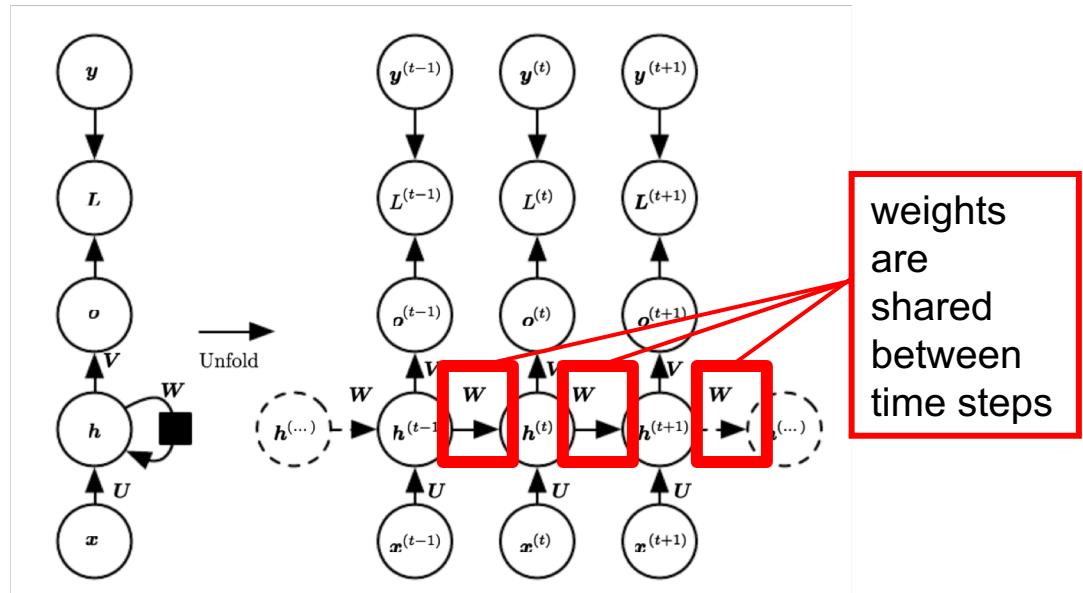
$$h_t = \sigma(\mathbf{W}h_{t-1} + \mathbf{U}x_t + \mathbf{b})$$

$$\mathbf{o}_t = \phi(\mathbf{V}h_t + \mathbf{c})$$



# Recurrent Neural Networks (RNNs)

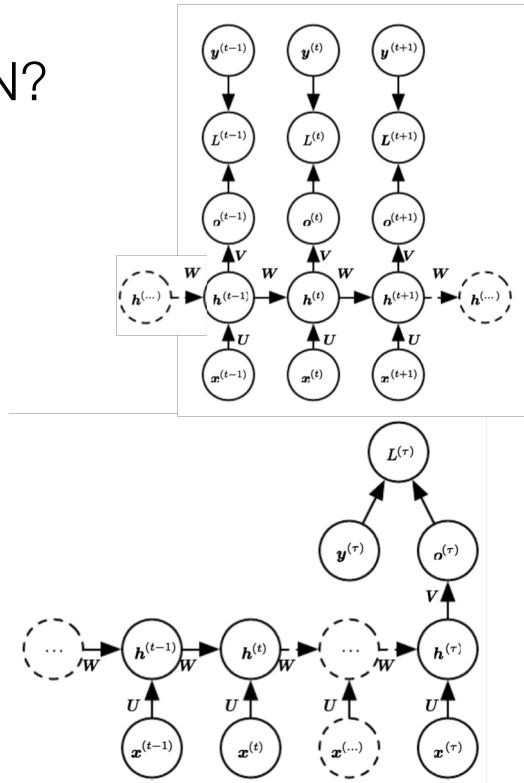
- Can **unroll** the RNN over time to form an acyclic graph.
- RNN = special kind of feed-forward network



$$\text{E.g., } \mathbf{h}_3 = \sigma(\mathbf{W}(\sigma(\mathbf{W}\sigma(\mathbf{W}\mathbf{h}_0 + \mathbf{U}\mathbf{x}_1 + \mathbf{b}) + \mathbf{U}\mathbf{x}_2\mathbf{b}) + \mathbf{U}\mathbf{x}_3 + \mathbf{b})$$

# Kinds of output

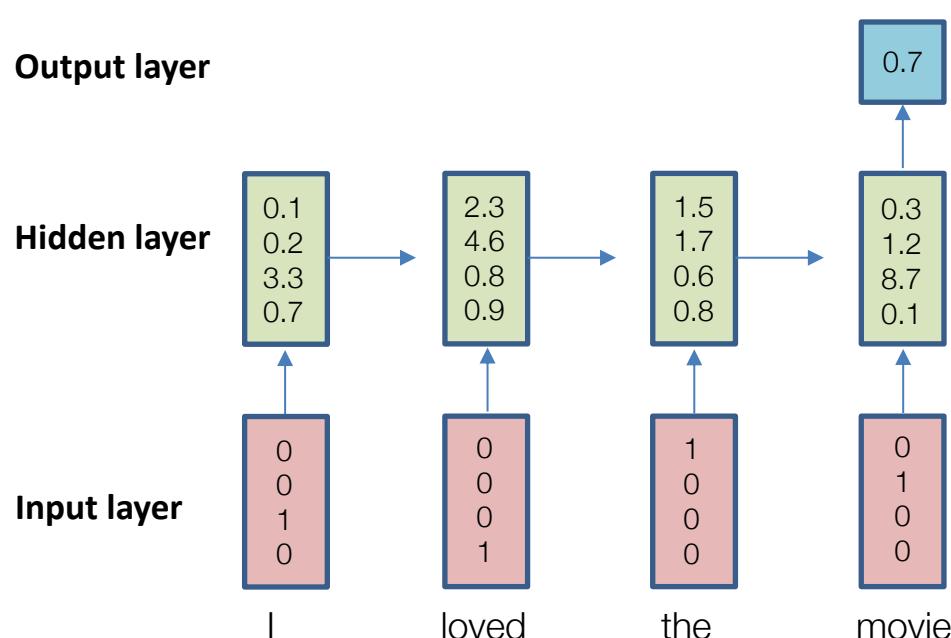
- How do we specify the target output of an RNN?
- Many ways! Two main ones:
  - 1) Can specify one target **at the end of the sequence**
    - Ex: sentiment classification
  - 2) Can specify a target **at each time step**
    - Ex: generating language



# Sentiment classification

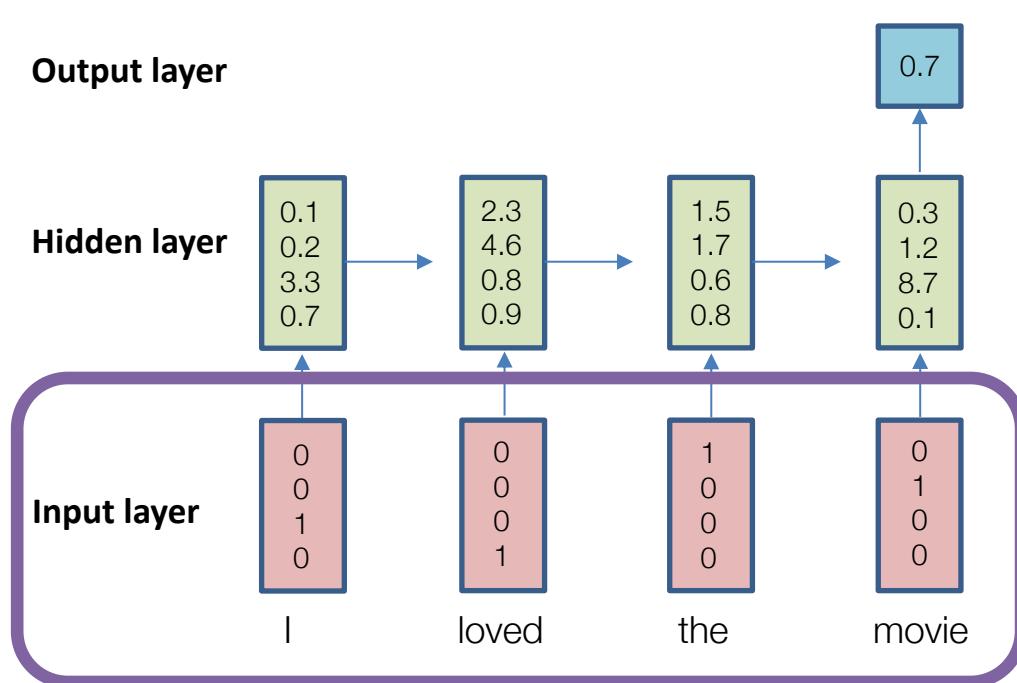
---

- **Input:** Sequence of words
  - E.g., a review, a tweet, a news article
- **Output:** A single value
  - E.g., indicating the probability that the text has a positive sentiment



# Sentiment classification

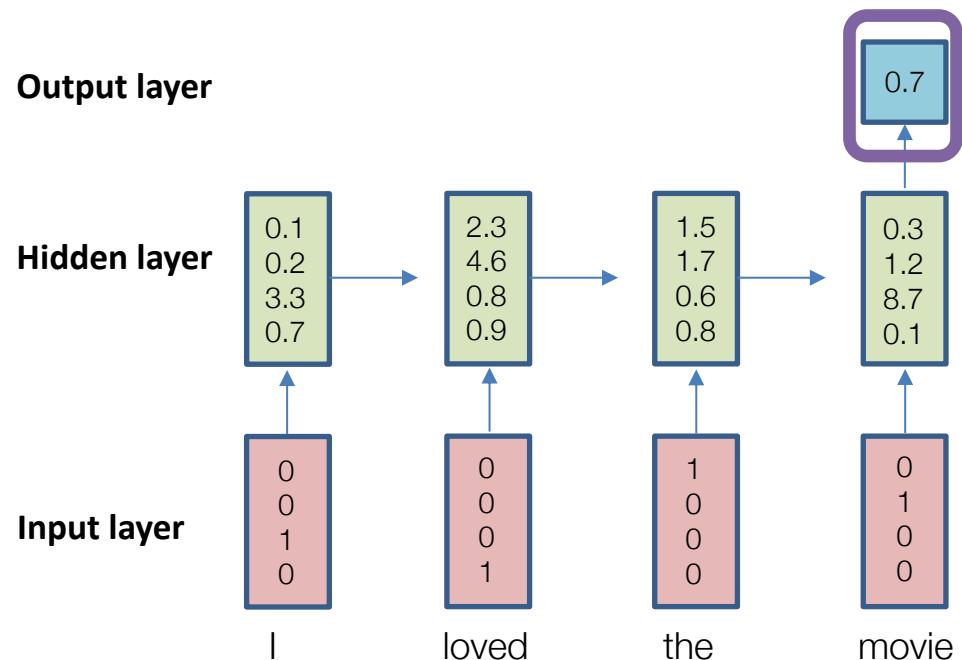
- **Input:** Sequence of words
  - E.g., review, tweet, or news article
- **Output:** A single value
  - E.g., indicating the probability that the text has a positive sentiment
- Words can be encoded as “one-hot” vectors.



# Sentiment classification

---

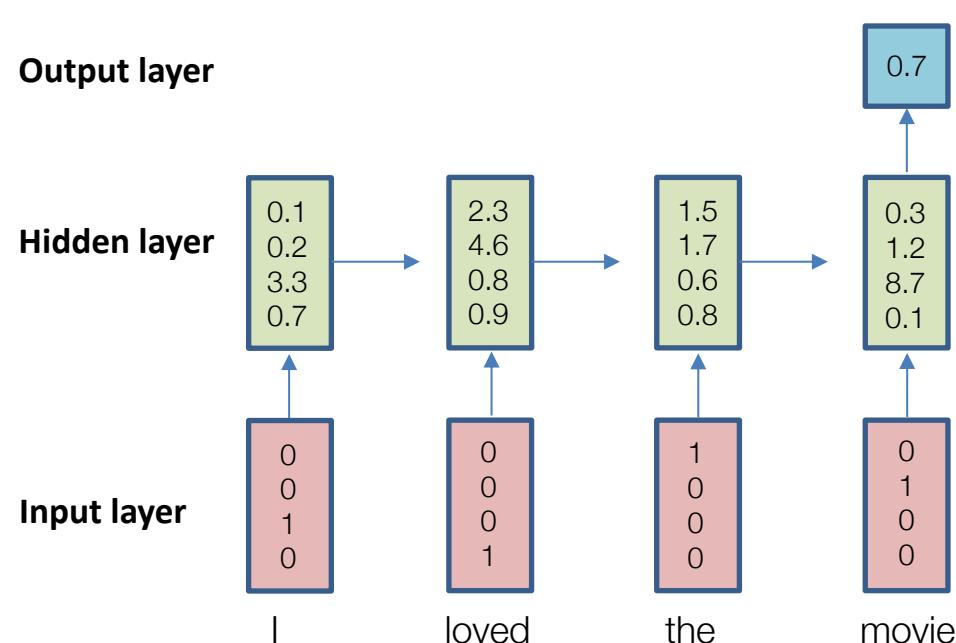
- **Input:** Sequence of words
  - E.g., review, tweet, or news article
- **Output:** A single value
  - E.g., indicating the probability that the text has a positive sentiment
- There is only output at the end of the sequence



# Sentiment classification

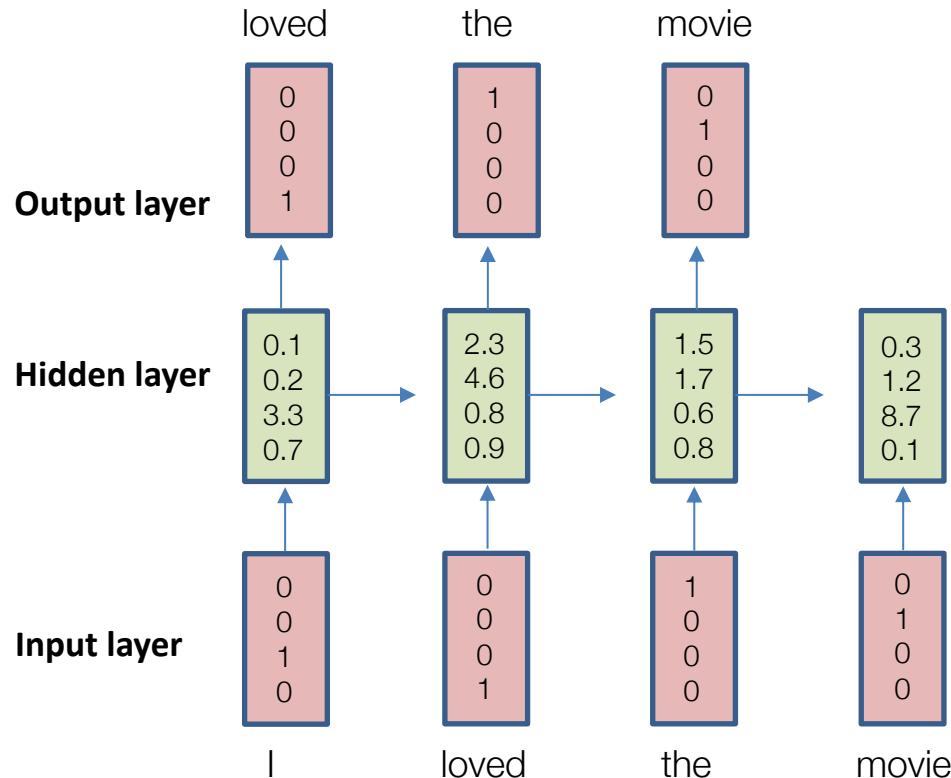
---

- **Input:** Sequence of words
  - E.g., review, tweet, or news article
- **Output:** A single value
  - E.g., indicating the probability that the text has a positive sentiment
- **Classic example of a “sequence classification” task.**



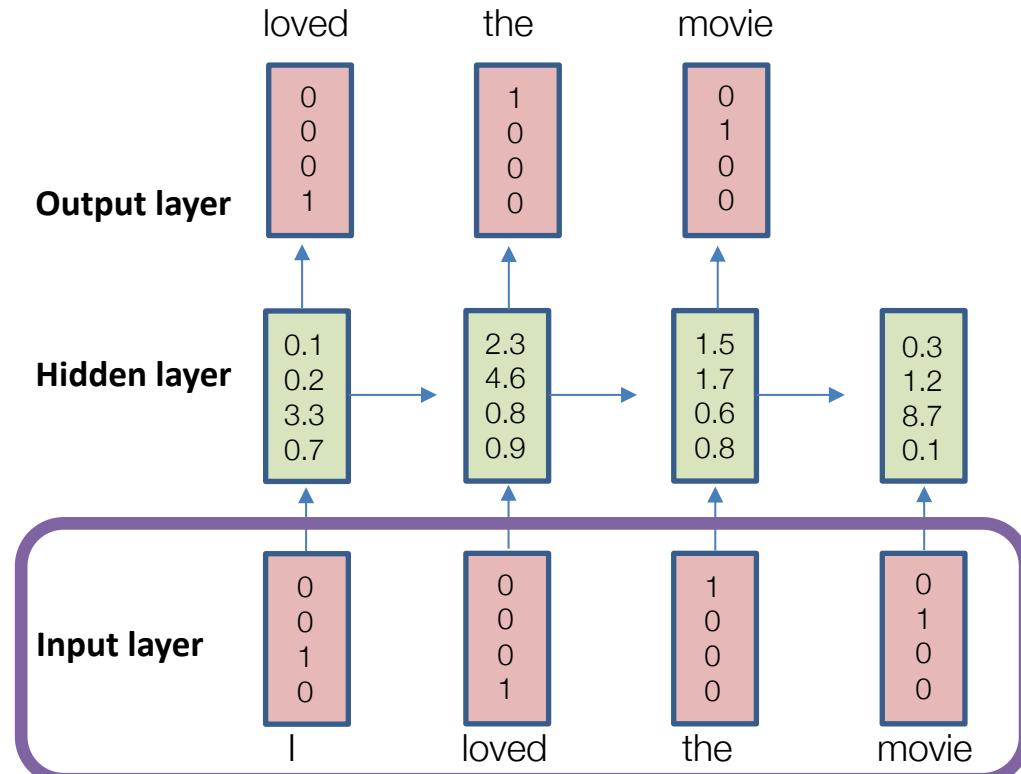
# Language modeling

- **Input:** Sequence of words
  - E.g., review, tweet, or news article
- **Output:** Sequence of words
  - E.g., predicting the next word that will occur in a sentence.



# Language modeling

- **Input:** Sequence of words
  - E.g., review, tweet, or news article
- **Output:** Sequence of words
  - E.g., predicting the next word that will occur in a sentence.
- Same input representation as sentiment classification

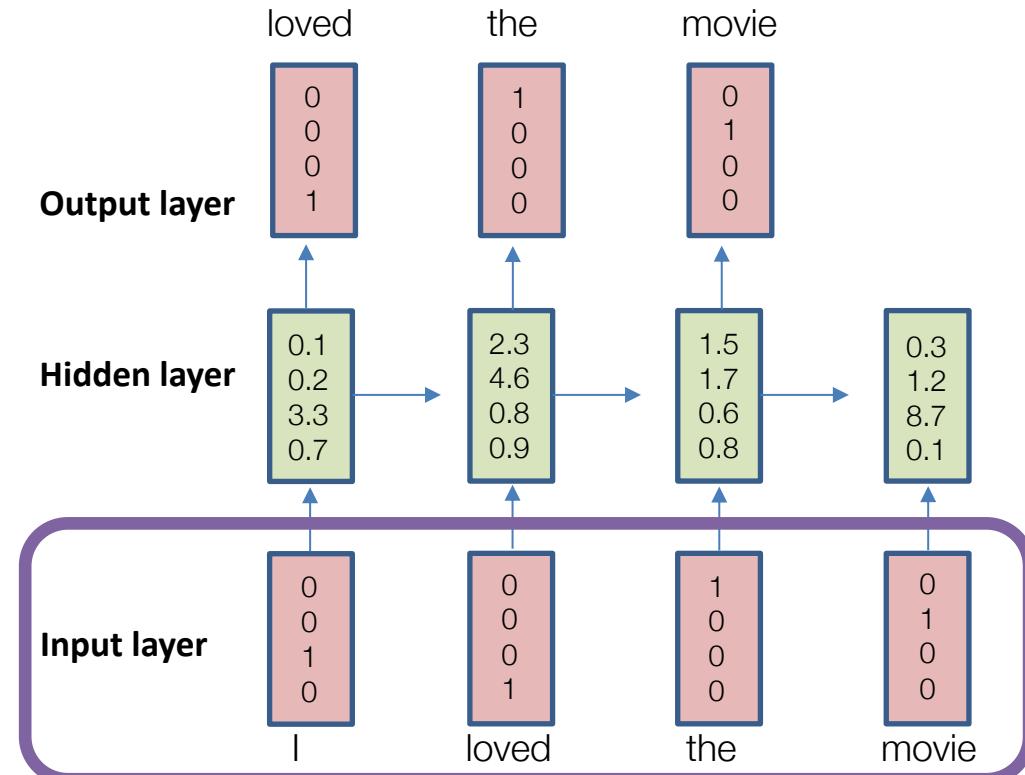


# Aside: Word embeddings

- Multiplication of one-hot vectors by a weight matrix is equivalent to mapping each word to a column of the weight matrix.

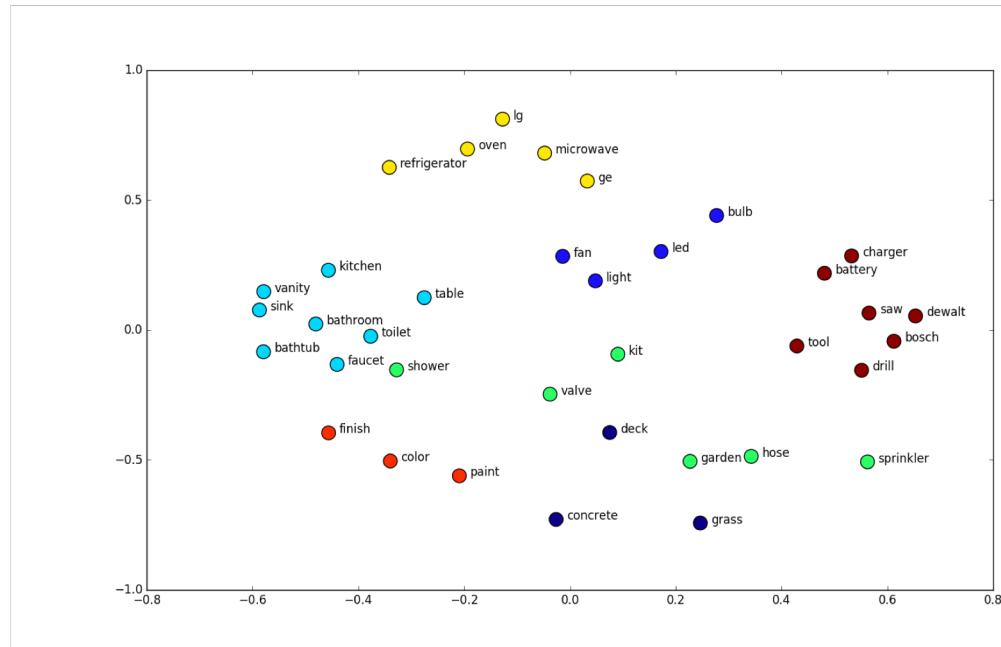
$$\mathbf{h}_t = \sigma(\mathbf{W}\mathbf{o}_{t-1} + \mathbf{U}\mathbf{x}_t + \mathbf{b})$$

- We call these word embeddings.
- Usually implemented via sparse “embedding lookup” rather than basic matrix multiplication.



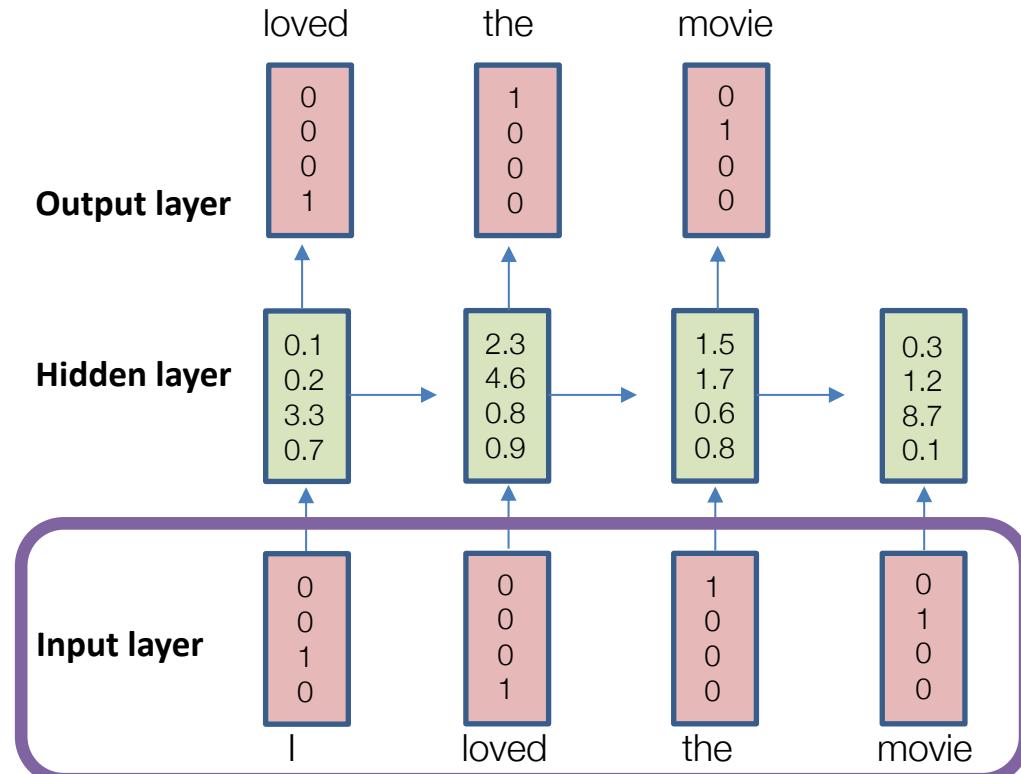
# Aside: Word embeddings

- Word embeddings are essentially  $d$ -dimensional vector representations of words.
- There are techniques to **pretrain** these word embeddings using large corpora.
- To be discussed in more detail in Lectures 18/19.



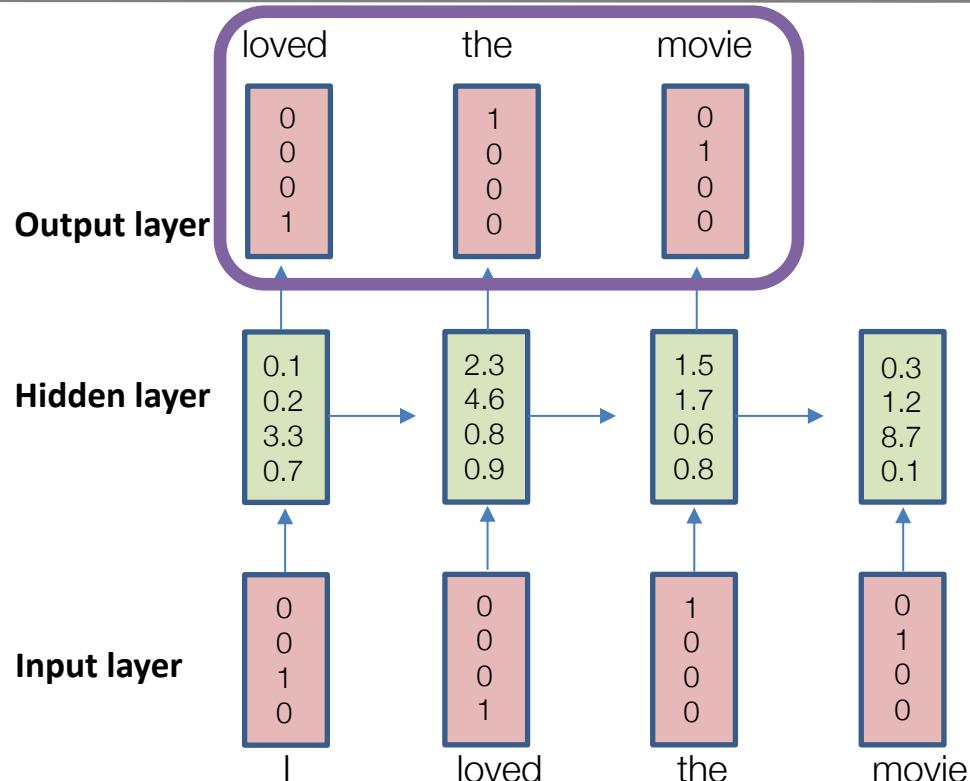
# Language modeling

- **Input:** Sequence of words
  - E.g., review, tweet, or news article
- **Output:** Sequence of words
  - E.g., predicting the next word that will occur in a sentence.
- Same input representation as sentiment classification



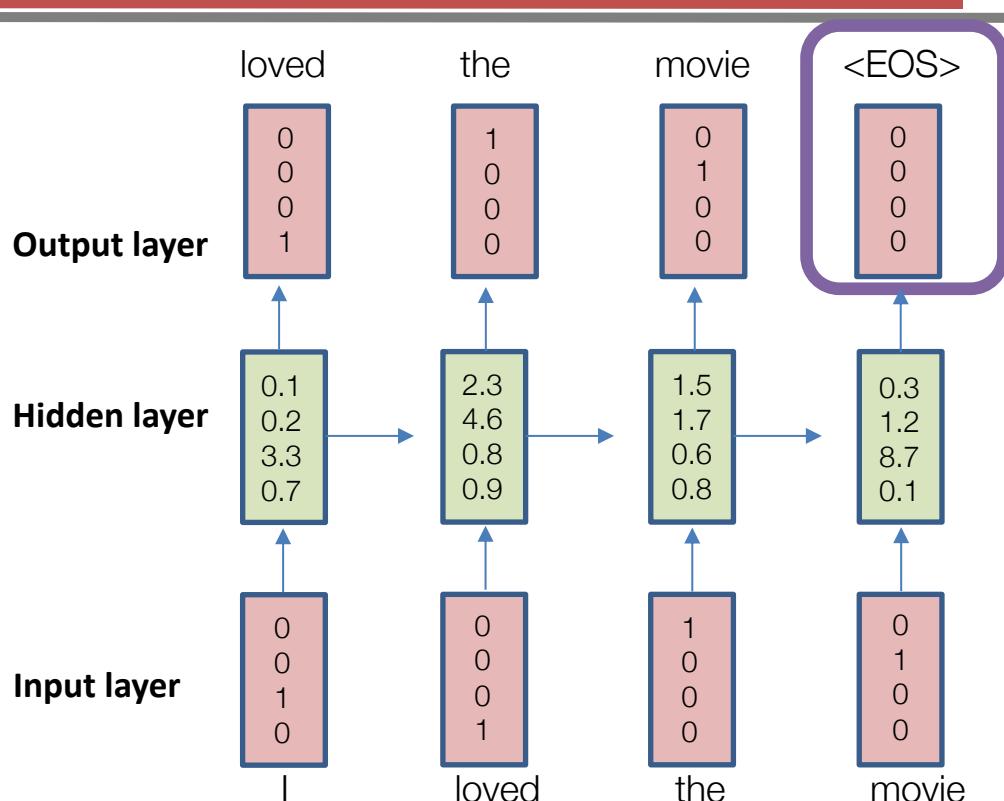
# Language modeling

- **Input:** Sequence of words
  - E.g., review, tweet, or news article
- **Output:** Sequence of words
  - E.g., predicting the next word that will occur in a sentence.
- **But now we have an output at each time-step!**
- Predicting the next word is a multiclass prediction problem (each word is a class), so use a softmax!



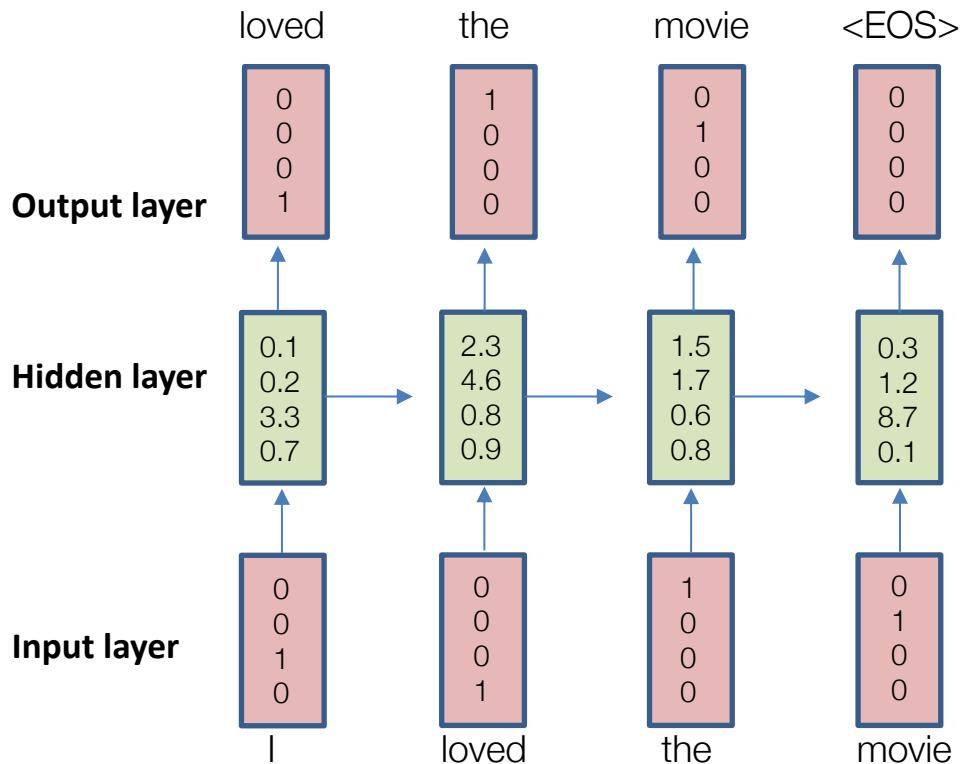
# Language modeling

- **Input:** Sequence of words
  - E.g., review, tweet, or news article
- **Output:** Sequence of words
  - E.g., predicting the next word that will occur in a sentence.
- Usually add an “end of sentence” token



# Language modeling

- **Input:** Sequence of words
  - E.g., review, tweet, or news article
- **Output:** Sequence of words
  - E.g., predicting the next word that will occur in a sentence.
- Classic “sequence modeling” task.
- Language modelling is the “backbone” of NLP.
  - Useful for machine translation, dialogue systems, automated captioning, etc.



# Training RNNs

---

- How can we train RNNs?
- Same as feed-forward networks: train with backpropagation on unrolled computation graph!

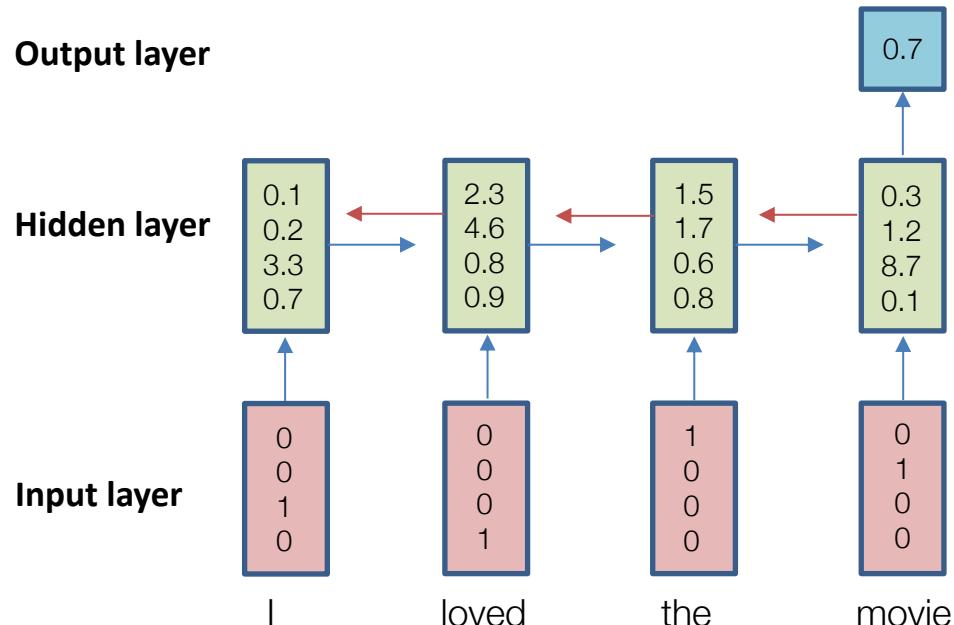
# Training RNNs

---

- How can we train RNNs?
  - Same as feed-forward networks: train with backpropagation on unrolled computation graph!
- 
- This is called **backpropagation through time** (BPTT)
  - Same derivation as regular backprop (use chain rule)

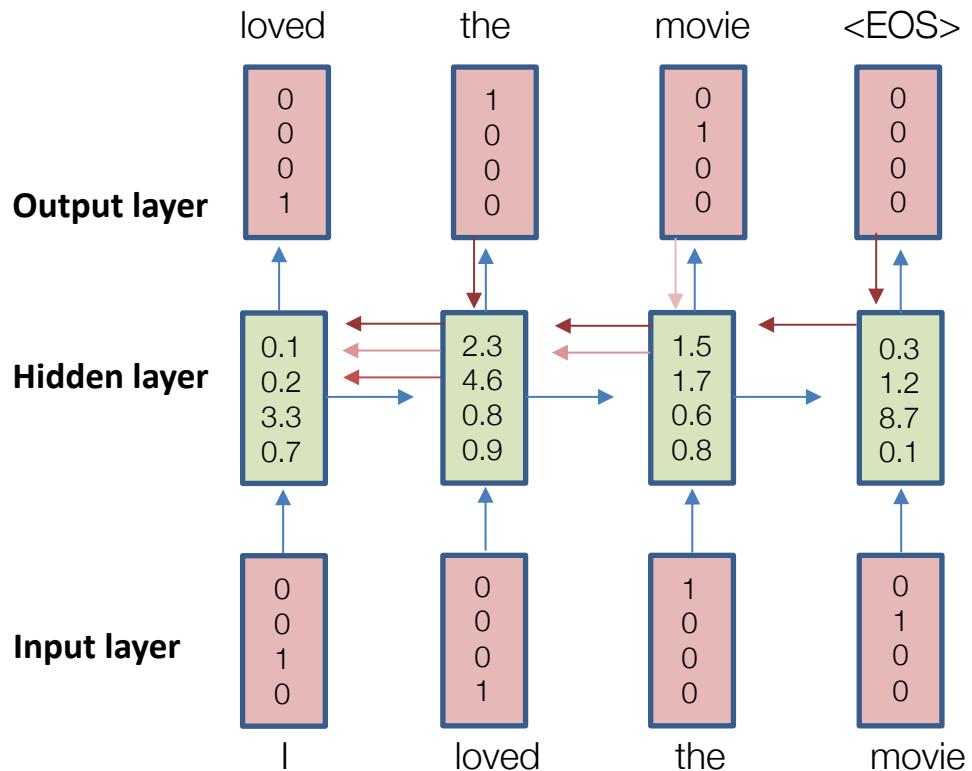
# Backpropagation through time

- BPTT is straightforward for sequence classification.
- Gradient flows from the final prediction back through all the layers.



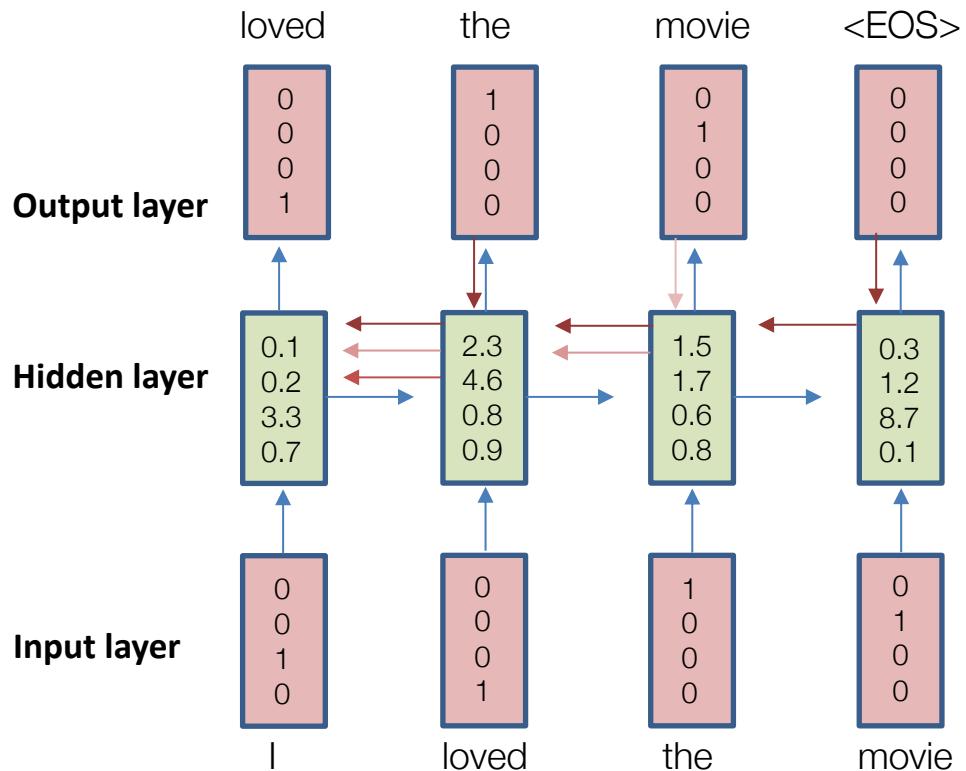
# Backpropagation through time

- BPTT is less straightforward for language modeling.
- Gradient flows from the prediction at each time-step to the preceding time-steps.



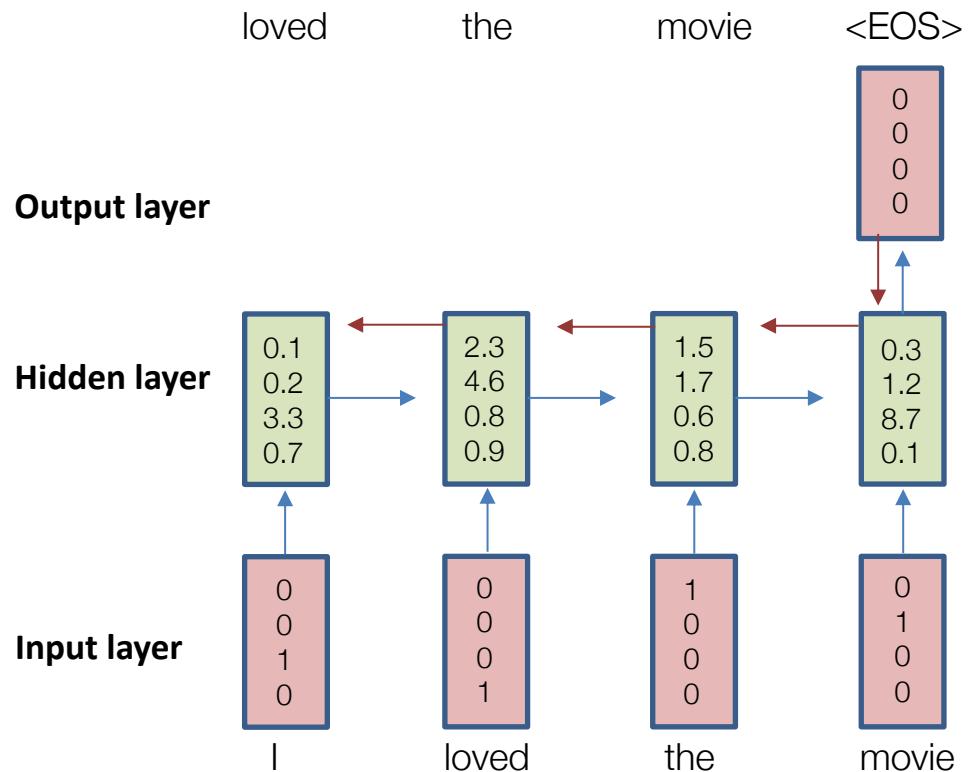
# Backpropagation through time

- BPTT is less straightforward for language modeling.
- Conceptually, we can think that we are jointly training on three sequence classification tasks.



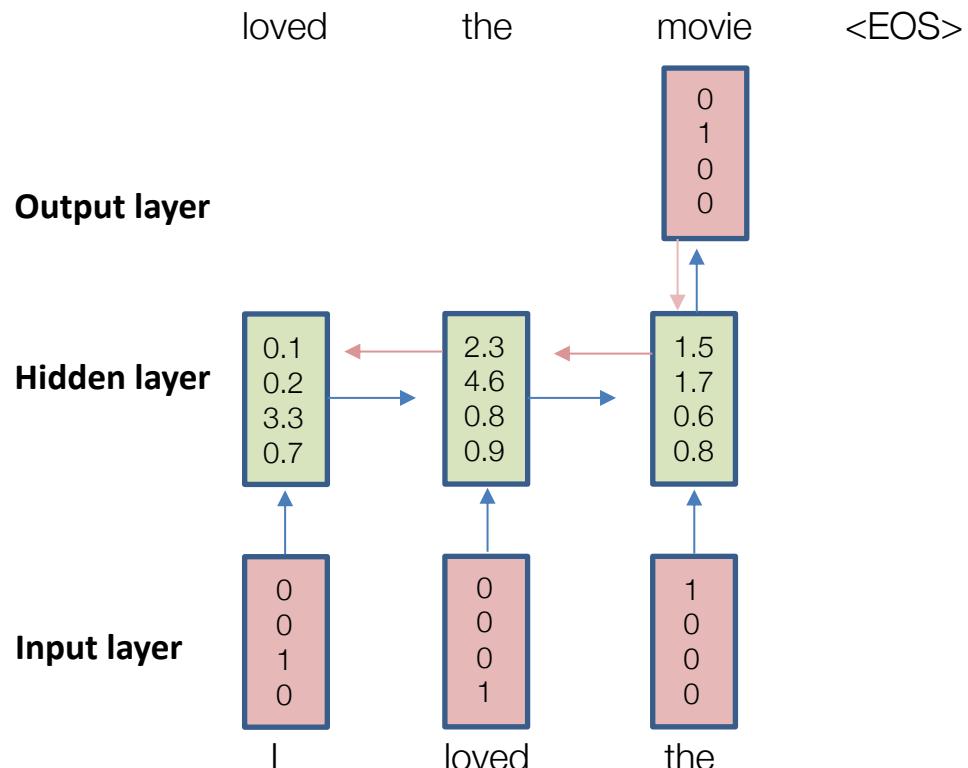
# Backpropagation through time

- BPTT is less straightforward for language modeling.
- Conceptually, we can think that we are jointly training on three sequence classification tasks.



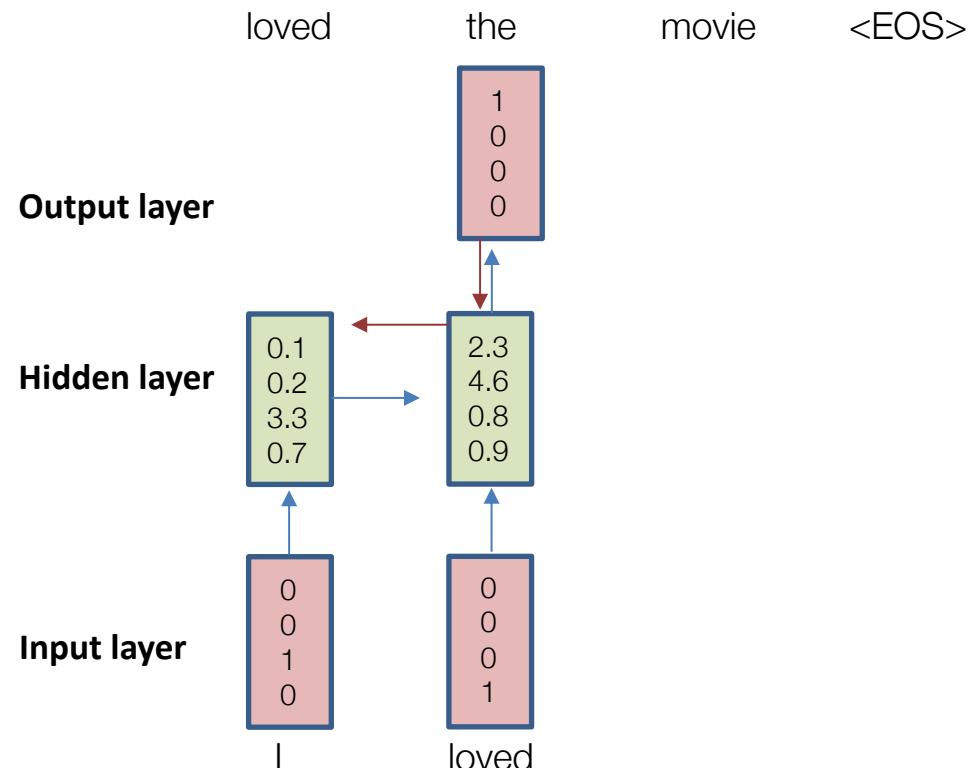
# Backpropagation through time

- BPTT is less straightforward for language modeling.
- Conceptually, we can think that we are jointly training on three sequence classification tasks.



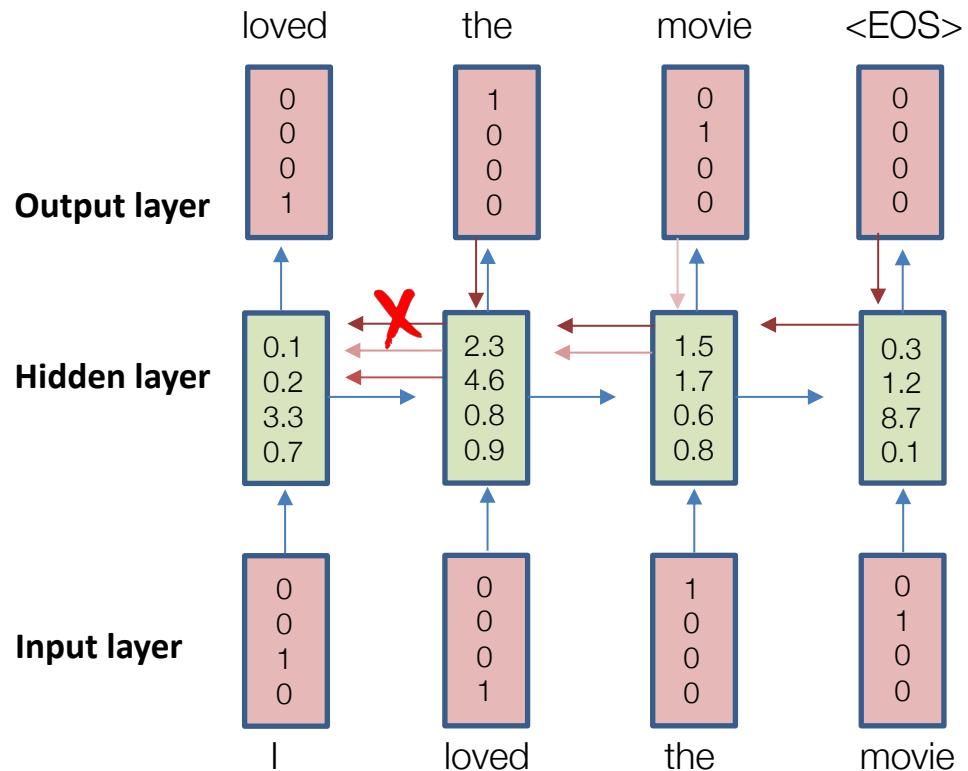
# Backpropagation through time

- BPTT is less straightforward for language modeling.
- Conceptually, we can think that we are jointly training on three sequence classification tasks.



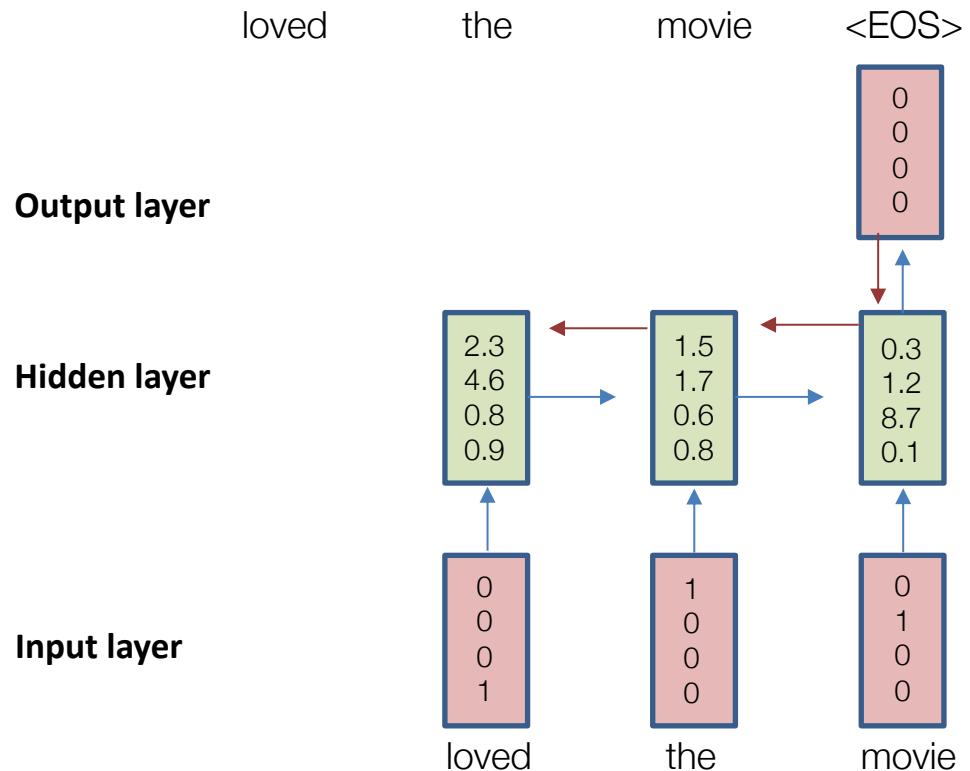
# Backpropagation through time

- BPTT is less straightforward for language modeling.
- For very long sequence/language modeling tasks, sometimes we “truncate” the gradient flow.



# Backpropagation through time

- BPTT is less straightforward for language modeling.
- For very long sequence/language modeling tasks, sometimes we “truncate” the gradient flow.
  - I.e., only the last K time-steps are used to train the prediction.

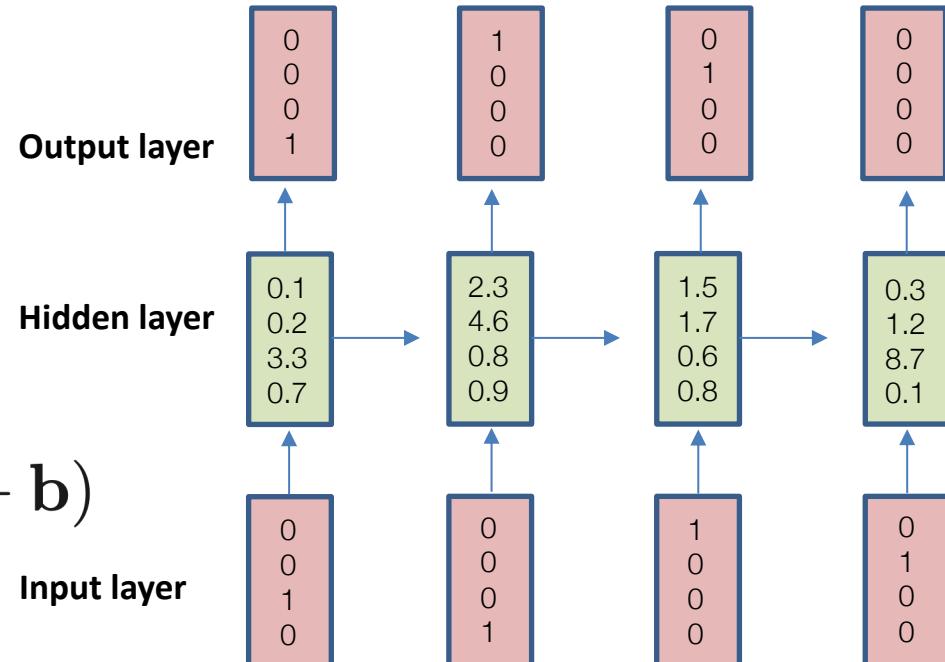


# There are many ways to add recurrence

- So far, we have been considering a standard/simple RNN.

- Recurrence is between the hidden states:

$$\mathbf{h}_t = \sigma(\mathbf{W}\mathbf{h}_{t-1} + \mathbf{U}\mathbf{x}_t + \mathbf{b})$$



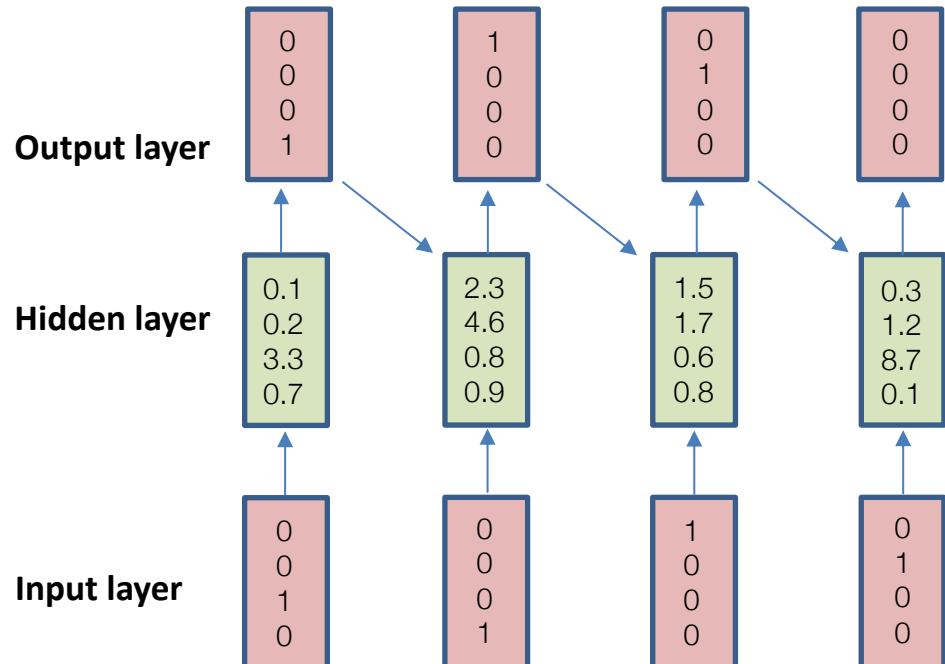
- This is called the **Elman RNN**.

# There are many ways to add recurrence

- But there are other options!
- E.g., recurrence based on the output:

$$\mathbf{h}_t = \sigma(\mathbf{W}\mathbf{o}_{t-1} + \mathbf{U}\mathbf{x}_t + \mathbf{b})$$

- This is called the **Jordan RNN**.



# There are many ways to add recurrence

- Q: Which is better?

- Elman RNN:

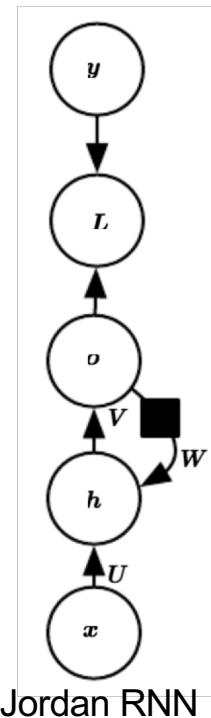
$$\mathbf{h}_t = \sigma(\mathbf{W}\mathbf{h}_{t-1} + \mathbf{U}\mathbf{x}_t + \mathbf{b})$$

$$\mathbf{o}_t = \phi(\mathbf{V}\mathbf{h}_t + \mathbf{c})$$

- Jordan RNN:

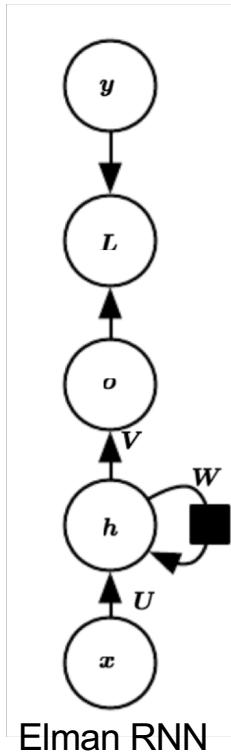
$$\mathbf{h}_t = \sigma(\mathbf{W}\mathbf{o}_{t-1} + \mathbf{U}\mathbf{x}_t + \mathbf{b})$$

$$\mathbf{o}_t = \phi(\mathbf{V}\mathbf{h}_t + \mathbf{c})$$



Jordan RNN

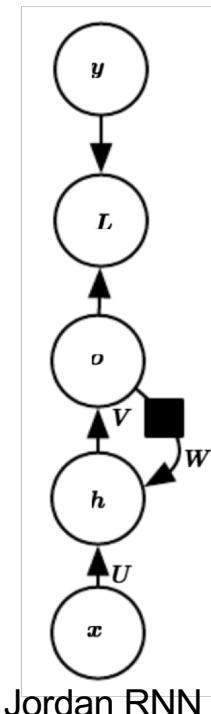
vs.



Elman RNN

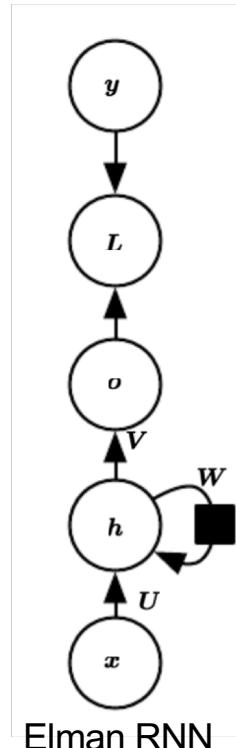
# There are many ways to add recurrence

- Q: Which is better?
- A: Elman RNN. Usually output  $o$  is constrained in some way, and may be missing some important info from the past.



Jordan RNN

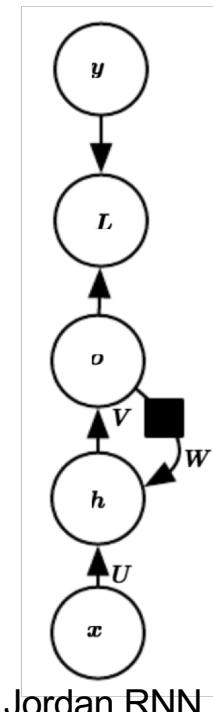
vs.



Elman RNN

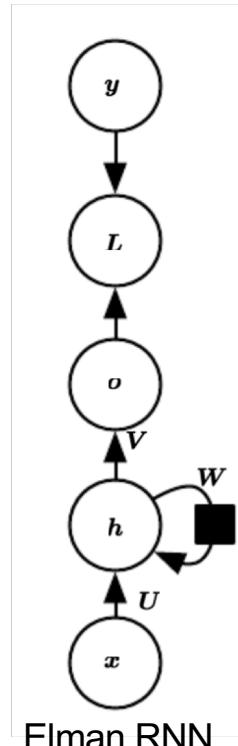
# There are many ways to add recurrence

- Q: Which is better?
- A: Elman RNN. Usually output  $o$  is constrained in some way, and may be missing some important info from the past.
- We can also add both types of recurrence at once!



Jordan RNN

vs.



Elman RNN

# We can also stack recurrent layers

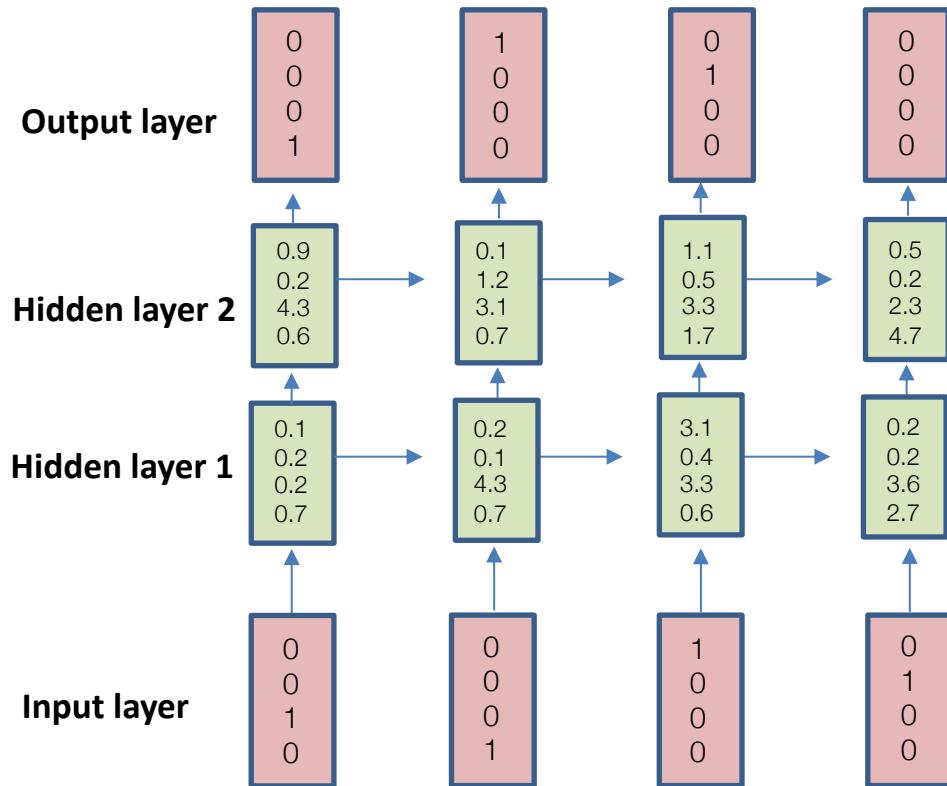
- I.e., RNNs can be “deep”

- E.g.,

$$\mathbf{h}_t^{(1)} = \sigma(\mathbf{W}^{(1)}\mathbf{h}_{t-1}^{(1)} + \mathbf{U}^{(1)}\mathbf{x}_t + \mathbf{b}^{(1)})$$

$$\mathbf{h}_t^{(2)} = \sigma(\mathbf{W}^{(2)}\mathbf{h}_{t-1}^{(2)} + \mathbf{U}^{(2)}\mathbf{h}_t^{(1)} + \mathbf{b}^{(2)})$$

- Again, there are many options for how to connect the deep recurrent layers...



# Beyond Elman and Jordan RNNs

---

- Elman and Jordan RNNs are relatively straightforward.
- But in practice they are very hard to train!
- **Issue:** Multiplying by the same  $\mathbf{W}$  matrix over and over is very unstable...

$$\mathbf{h}_t = \sigma(\mathbf{W}\mathbf{h}_{t-1} + \mathbf{U}\mathbf{x}_t + \mathbf{b})$$

$$\text{E.g., } \mathbf{h}_3 = \sigma(\mathbf{W}(\sigma(\mathbf{W}\sigma(\mathbf{W}\mathbf{h}_0 + \mathbf{U}\mathbf{x}_1 + \mathbf{b}) + \mathbf{U}\mathbf{x}_2\mathbf{b}) + \mathbf{U}\mathbf{x}_3 + \mathbf{b})$$

- There are recurrent architectures that fix this! (Next lecture).