

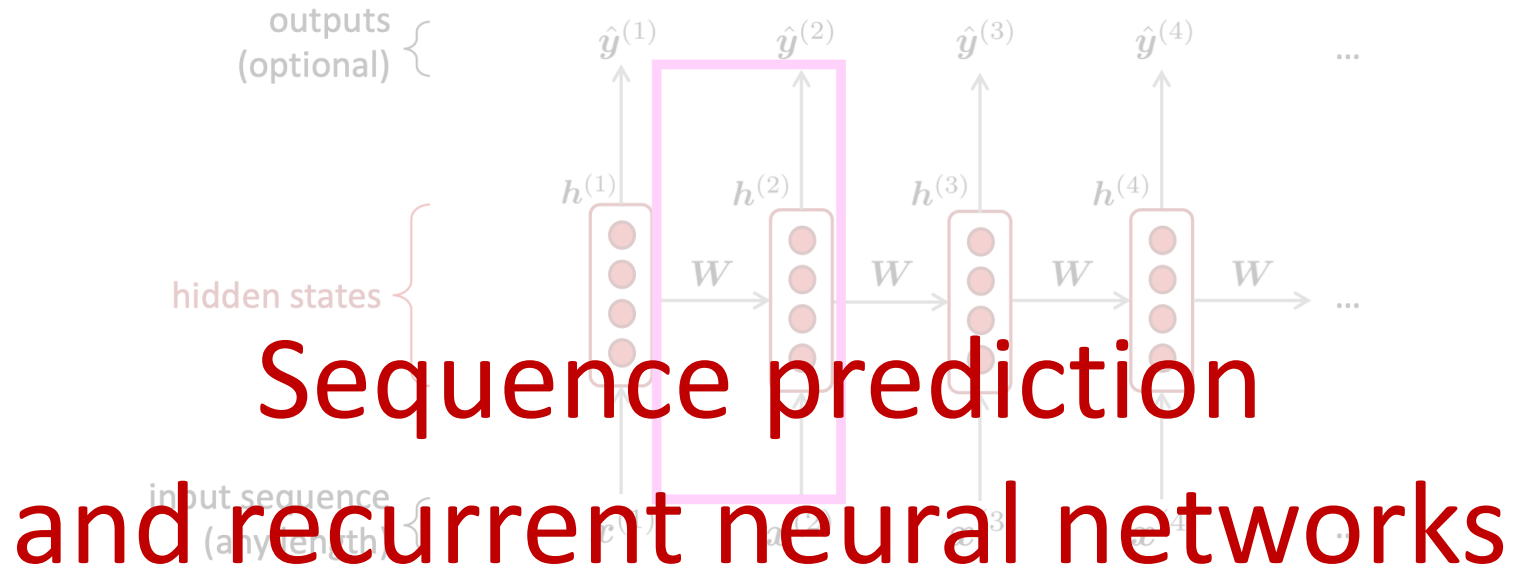
# CSCI 567: Machine Learning

Vatsal Sharan  
Spring 2026

Lecture 9, March 27

# Administrivia

- HW3 due in less than 1 week



# Acknowledgements

We borrow heavily from:

- Stanford's CS224n: <https://web.stanford.edu/class/cs224n/>

# Sequential prediction

Given observations  $x_1, x_2, \dots, x_{t-1}$  what is  $x_t$ ?

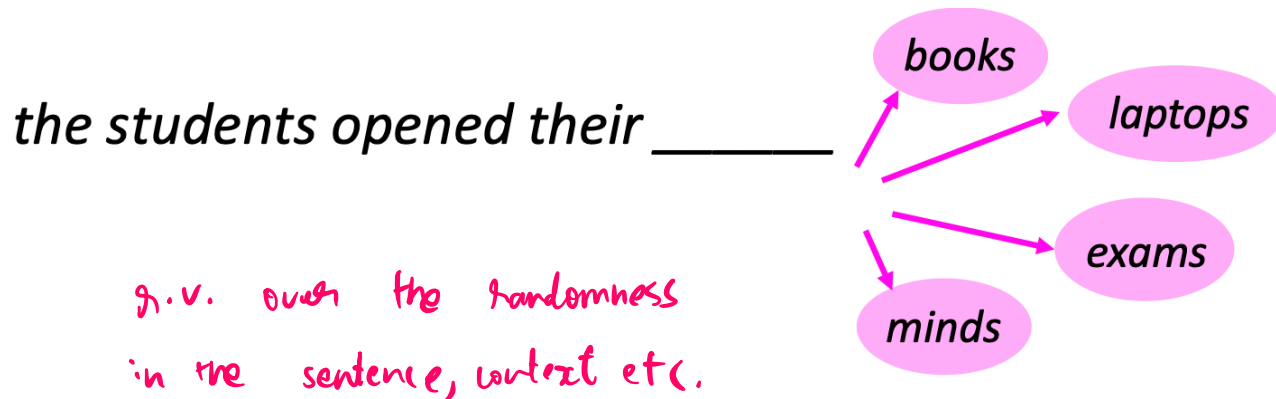
Examples:

- text or speech data
- stock market data
- weather data
- ...

In this lecture, we will mostly focus on text data (**language modelling**).

# Language modelling

**Language modelling** is the task of predicting what word comes next:



More formally, let  $X_i$  be the random variable for the  $i$ -th word in the sentence, and let  $x_i$  be the **value** taken by the random variable. Then the goal is to compute

$$P(X_{t+1} | X_t = x_t, \dots, X_1 = x_1).$$

A system that does this is known as a **Language Model**.

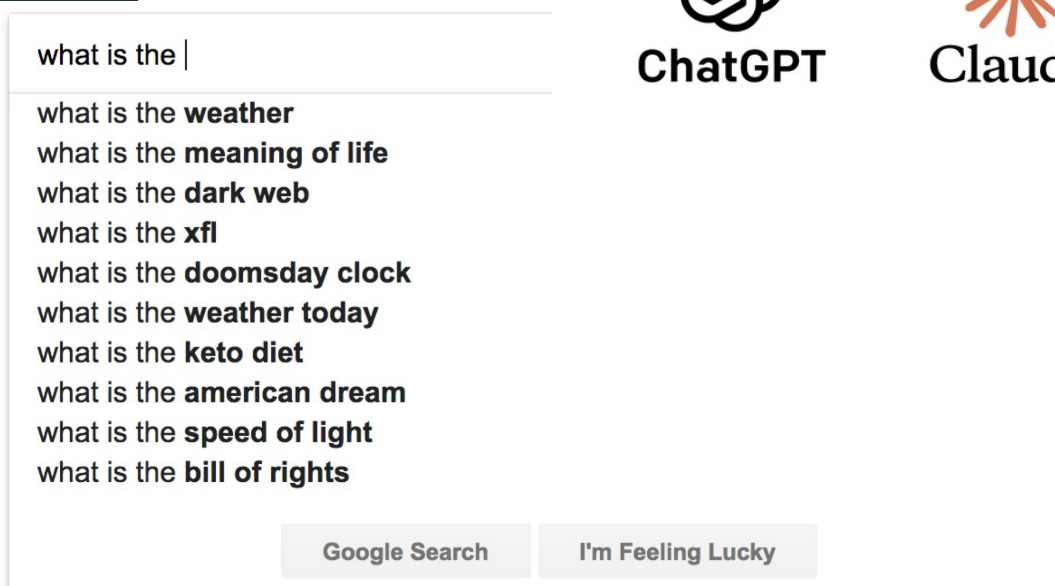
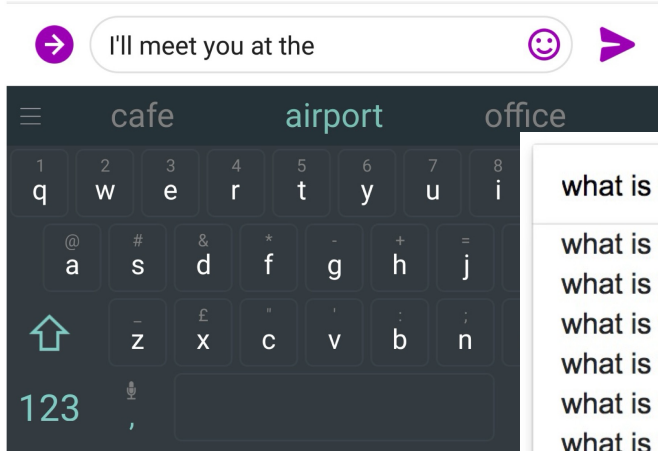
# Language modelling

We can also think of a Language Model as a system that *assigns a probability to a piece of text*.

For example, if we have some text  $x_1, \dots, x_T$ , then the probability of this text (according to the Language Model) is:

$$\begin{aligned} P(X_1 = x_1, \dots, X_T = x_T) &= P(X_1 = x_1) \times P(X_2 = x_2 | X_1 = x_1) \\ &\times \dots \times P(X_T = x_T | X_{T-1} = x_{T-1}, \dots, X_1 = x_1) \\ &= \prod_{t=1}^T P(X_t = x_t | X_{t-1} = x_{t-1}, \dots, X_1 = x_1). \end{aligned}$$

# Language modelling is used everywhere



# n-gram Language Models

*the students opened their \_\_\_\_\_*

- **Question:** How to learn a Language Model?
- **Answer** (pre- Deep Learning): learn an *n*-gram Language Model!
- **Definition:** An *n*-gram is a chunk of *n* consecutive words.
  - **unigrams:** “the”, “students”, “opened”, “their”
  - **bigrams:** “the students”, “students opened”, “opened their”
  - **trigrams:** “the students opened”, “students opened their”
  - **four-grams:** “the students opened their”
- **Idea:** Collect statistics about how frequent different n-grams are and use these to predict next word.

# $n$ -gram language model: A type of Markov model

A **Markov model or Markov chain** is a sequence of random variables with the **Markov property**: a sequence of random variables  $X_1, X_2, \dots$  s.t.

$$P(X_{t+1} | X_{1:t}) = P(X_{t+1} | X_t) \quad (\text{Markov property})$$

i.e. *the next state only depends on the most recent state* (notation  $X_{1:t}$  denotes the sequence  $X_1, \dots, X_t$ ). This is a *bigram model*.

We will consider the following setting:

- All  $X_t$ 's take value from the same **discrete** set  $\{1, \dots, S\}$  the size of the dictionary of all possible words
- $P(X_{t+1} = s' | X_t = s) = a_{s,s'}$ , known as **transition probability**
- $P(X_1 = s) = \pi_s$  initial probability
- $(\{\pi_s\}, \{a_{s,s'}\}) = (\boldsymbol{\pi}, \mathbf{A})$  are **parameters of the model**

$$A \in \mathbb{R}^{S \times S} \quad A = \begin{pmatrix} a_{s,s'} \end{pmatrix}$$

$$P(X_1, \dots, X_T) = P(X_1) \cdot P(X_2 | X_1) \cdot P(X_3 | X_2) \cdot \dots \cdot P(X_T | X_{T-1})$$

# Markov model: examples

- Example 1 (**Language model**)  
States  $[S]$  represent a dictionary of words,

$$a_{\text{ice,cream}} = P(X_{t+1} = \text{cream} \mid X_t = \text{ice})$$

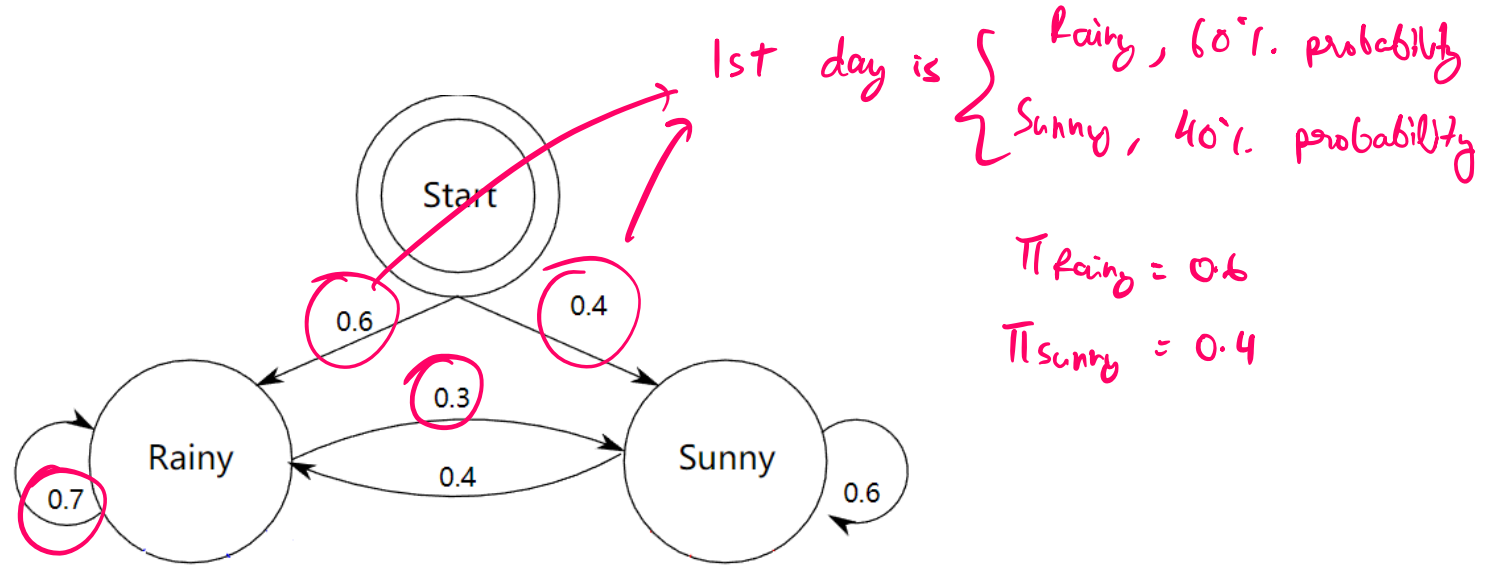
is an example of the transition probability.

- Example 2 (**Weather**)  
States  $[S]$  represent weather at each day

$$a_{\text{sunny,rainy}} = P(X_{t+1} = \text{rainy} \mid X_t = \text{sunny})$$

# Markov model: Graphical representation

A Markov model is nicely represented as a **directed graph**



if today is rainy, next day is  $\begin{cases} \text{rainy, } 70\% \text{ prob.} \\ \text{Sunny, } 30\% \text{ prob.} \end{cases}$

# Learning Markov models

Now suppose we have observed  $n$  **sequences of examples**:

- $x_{1,1}, \dots, x_{1,T}$  (rainy, sunny, ..., rainy)
- ...
- $x_{i,1}, \dots, x_{i,T}$
- ...
- $x_{n,1}, \dots, x_{n,T}$

where

- for simplicity we assume each sequence has the same **length  $T$**
- lower case  $x_{i,t}$  represents the **value** of the random variable  $X_{i,t}$

From these observations how do we **learn the model parameters**  $(\pi, \mathbf{A})$ ?

# Learning Markov models: MLE

Same story, find the **MLE**. The log-likelihood of a sequence  $x_1, \dots, x_T$  is

$$\begin{aligned} \ln P(X_{1:T} = x_{1:T}) \\ &= \sum_{t=1}^T \ln P(X_t = x_t \mid X_{1:t-1} = x_{1:t-1}) && \text{(always true)} \\ &= \sum_{t=1}^T \ln P(X_t = x_t \mid X_{t-1} = x_{t-1}) && \text{(Markov property)} \end{aligned}$$

$\ln P(X_1 = x_1)$   
 $= \ln \pi_{x_1}$

$\ln P(X_{1:T} = x_{1:T}) = \ln \pi_{x_1} + \sum_{t=2}^T \ln a_{x_{t-1}, x_t}$

↳  $\ln \pi_{x_1}$  is the probability of starting at  $x_1$ .  
↳  $\ln a_{x_{t-1}, x_t}$  is the probability of transitioning from  $x_{t-1} \rightarrow x_t$ .

$$= \sum_s \mathbb{I}[x_1 = s] \ln \pi_s + \sum_{s, s'} \left( \sum_{t=2}^T \mathbb{I}[x_{t-1} = s, x_t = s'] \right) \ln a_{s, s'}$$

This is over one sequence, can sum over all

# Learning Markov models: MLE

So MLE is

$$\operatorname{argmax}_{\pi, A} \sum_s (\text{\#initial states with value } s) \ln \pi_s + \sum_{s, s'} (\text{\#transitions from } s \text{ to } s') \ln a_{s, s'}$$

*if this is large for some s*

*this should be large for that s*

This is an optimization problem, and can be solved by hand (though we'll skip in class).  
The solution is:

$$\pi_s = \frac{\text{\#initial states with value } s}{\text{\#initial states}}$$
$$a_{s, s'} = \frac{\text{\#transitions from } s \text{ to } s'}{\text{\#transitions from } s \text{ to any state}}$$

# Learning Markov models: Another perspective

Let's first look at the transition probabilities. By the Markov assumption,

$$P(X_{t+1} = x_{t+1} \mid X_t = x_t, \dots, X_1 = x_1) = P(X_{t+1} = x_{t+1} \mid X_t = x_t)$$

Using the definition of conditional probability,

$$P(X_{t+1} = x_{t+1} \mid X_t = x_t) = \frac{P(X_{t+1} = x_{t+1}, X_t = x_t)}{P(X_t = x_t)}$$

We can estimate this using data,

$$\frac{P(X_{t+1} = x_{t+1}, X_t = x_t)}{P(X_t = x_t)} \approx \frac{\text{\#times } (x_t, x_{t+1}) \text{ appears} \quad / \quad \cancel{\text{\# sequences}}}{\text{\# times } (x_t) \text{ appears (and is not the last state)} \quad | \quad \cancel{\text{\# sequences}}}$$

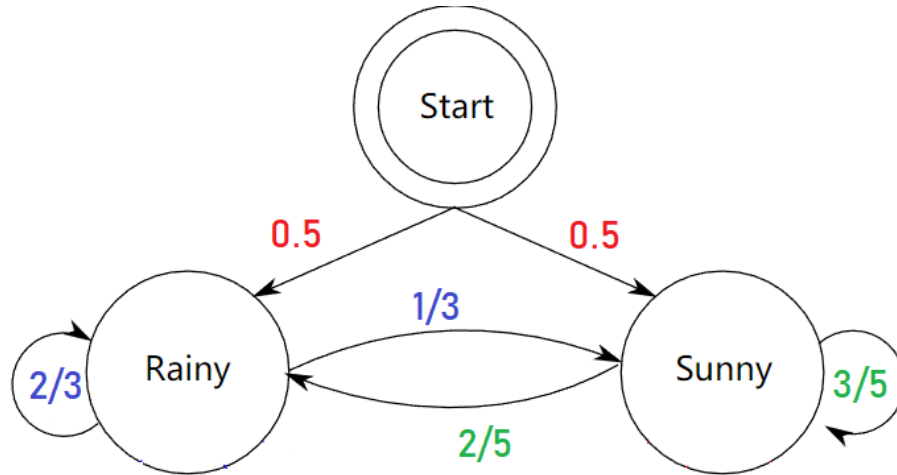
The initial state distribution follows similarly,

$$P(X_1 = s) \approx \frac{\text{\#times } s \text{ is first state}}{\text{\#sequences}}$$

# Learning Markov models: Example

Suppose we observed the following 2 sequences of length 5

- sunny, sunny, rainy, rainy, rainy
- rainy, sunny, sunny, sunny, rainy



# Higher-order Markov models

*Is the Markov assumption reasonable?* Not so in many cases, such as for language modeling.

Higher order Markov chains make it a bit more reasonable, e.g.

$$P(X_{t+1} \mid X_t, \dots, X_1) = P(X_{t+1} \mid X_t, X_{t-1}) \quad (\text{second-order Markov assumption})$$

i.e. the current word only depends on the last two words. This is a *trigram model*, since we need statistics of three words at a time to learn. In general, we can consider a  $n$ -th Markov model (or a  $(n + 1)$ -gram model):

$$P(X_{t+1} \mid X_t, \dots, X_1) = P(X_{t+1} \mid \overbrace{X_t, X_{t-1}, \dots, X_{t-n+1}}^{\text{previous } n \text{ observations}}) \quad (n\text{-th order Markov assumption})$$

Learning higher order Markov chains is similar, but more expensive.

$$\begin{aligned} P(X_{t+1} = x_{t+1} \mid X_t = x_t, \dots, X_1 = x_1) &= P(X_{t+1} = x_{t+1} \mid X_t = x_t, X_{t-1} = x_{t-1}, \dots, X_{t-n+1} = x_{t-n+1}) \\ &= \frac{P(X_{t+1} = x_{t+1}, X_t = x_t, X_{t-1} = x_{t-1}, \dots, X_{t-n+1} = x_{t-n+1})}{P(X_t = x_t, X_{t-1} = x_{t-1}, \dots, X_{t-n+1} = x_{t-n+1})} \\ &\approx \frac{\text{count}(x_{t-n+1}, \dots, x_{t-1}, x_t, x_{t+1}) \text{ in the data}}{\text{count}(x_{t-n+1}, \dots, x_{t-1}, x_t) \text{ in the data}} \end{aligned}$$

# n-gram Language Models: Example

Suppose we are learning a 4-gram Language Model.

~~as the proctor started the clock, the~~ students opened their \_\_\_\_\_  
discard condition on this

$$P(\mathbf{w} | \text{students opened their}) = \frac{\text{count}(\text{students opened their } \mathbf{w})}{\text{count}(\text{students opened their})}$$

For example, suppose that in the corpus:

- “students opened their” occurred 1000 times
  - “students opened their books” occurred 400 times
    - $\rightarrow P(\text{books} | \text{students opened their}) = 0.4$
  - “students opened their exams” occurred 100 times
    - $\rightarrow P(\text{exams} | \text{students opened their}) = 0.1$
- Should we have discarded the “proctor” context?

# n-gram Language Models in practice

- You can build a simple trigram Language Model over a 1.7 million word corpus (Reuters) in a few seconds on your laptop

Business and financial news

*today the* \_\_\_\_\_

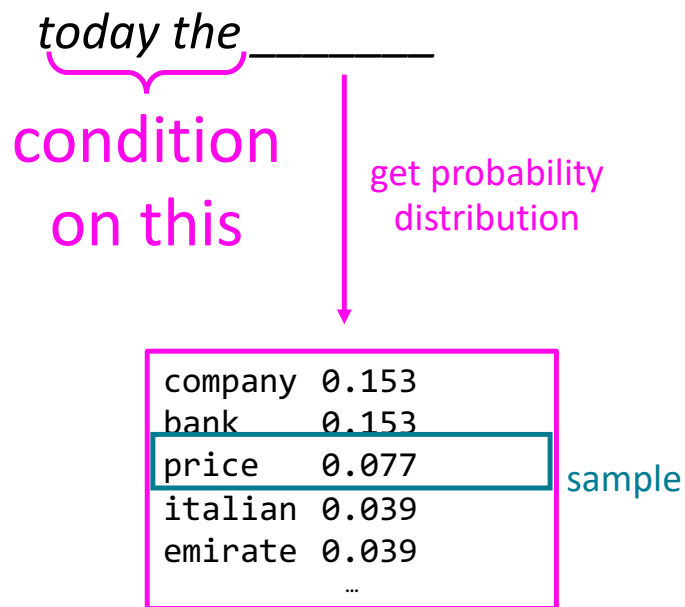
get probability  
distribution

company	0.153
bank	0.153
price	0.077
italian	0.039
emirate	0.039
...	

Notice that there isn't that much granularity in the distribution, because "*today the*" doesn't appear too often in corpus. Most two-grams won't appear too often.

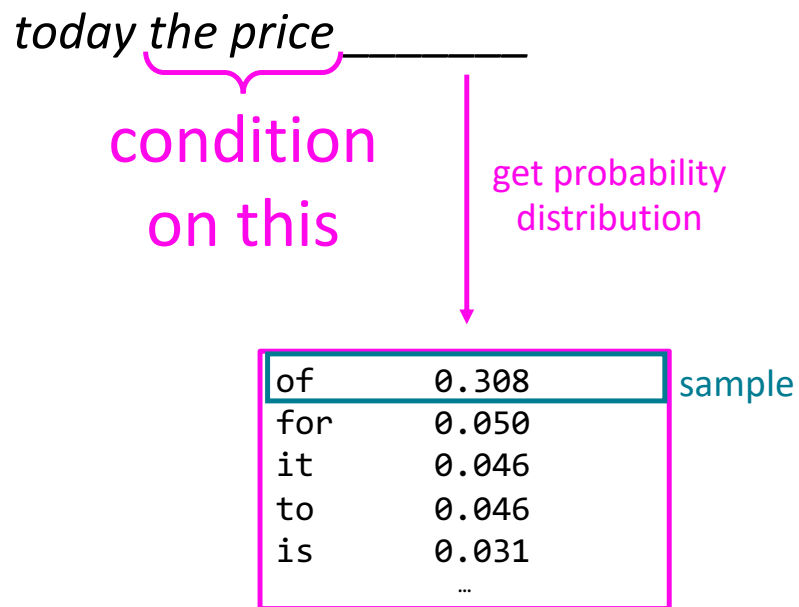
# Generating text with a n-gram Language Model

You can also use a Language Model to generate text



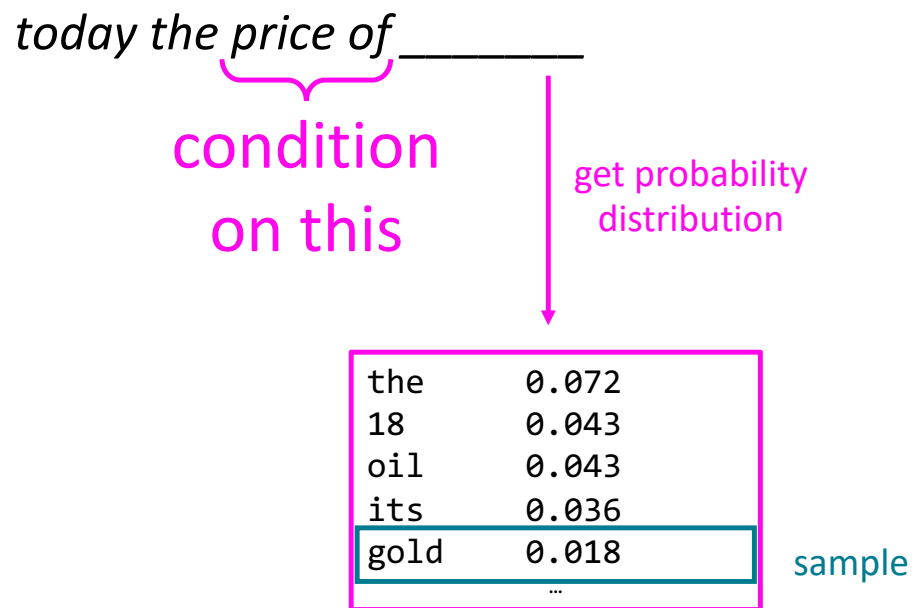
# Generating text with a n-gram Language Model

You can also use a Language Model to generate text



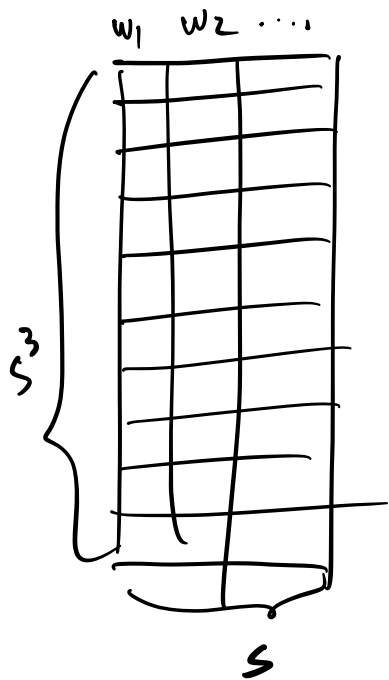
# Generating text with a n-gram Language Model

You can also use a Language Model to generate text



# Generating text with a n-gram Language Model

You can also use a Language Model to generate text



*today the price of gold per ton , while production of shoe lasts and shoe industry , the bank intervened just after it considered and rejected an imf demand to rebuild depleted european stocks , sept 30 end primary 76 cts a share .*

Surprisingly grammatical!

...but **incoherent**. We need to consider more than three words at a time if we want to model language well.

However, larger  $n$  increases model size and requires too much data to learn

# How to build a *neural* Language Model?

- Recall the Language Modeling task:
  - Input: sequence of words  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(t)}$
  - Output: prob dist of the next word  $P(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})$
- How about a **window-based neural model**?

Overloading notation,  $x^{(i)}$   
is overloaded to the r.v.  $X^{(i)}$   
↗ & its value  $x^{(i)}$ .



# Word embeddings/vectors

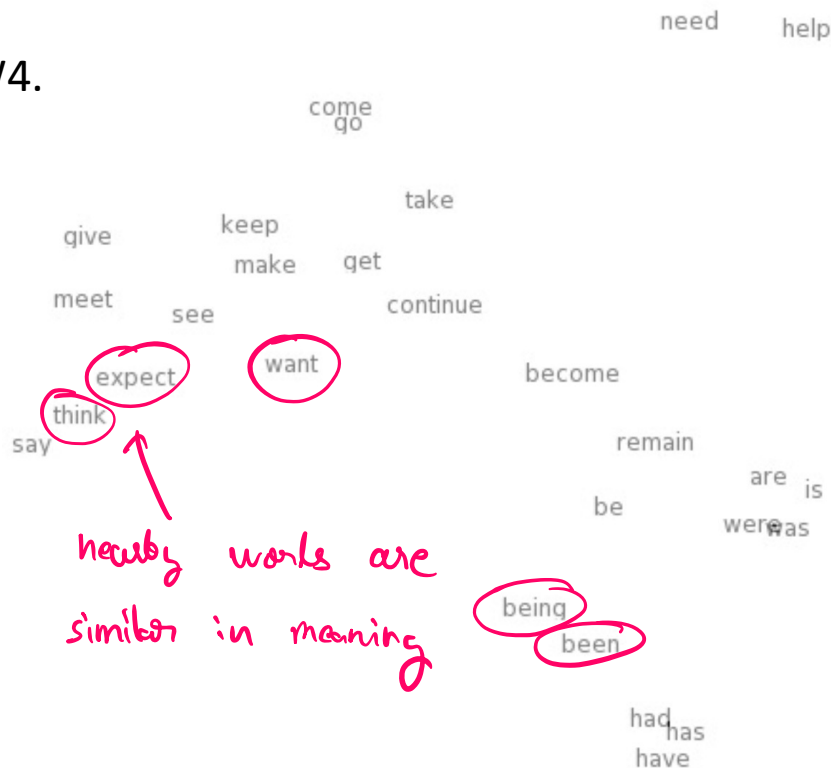
A word embedding is a (dense) mapping from words, to vector representations of the words.

Ideally, this mapping has the property that words similar in meaning have representations which are close to each other in the vector space.

You'll see a simple way to construct these in HW4.

*expect* =  $\begin{pmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \\ 0.487 \end{pmatrix}$

*10-dim embedding*



*nearby words are similar in meaning*

# A fixed-window neural Language Model

output distribution

$$\hat{y} = \text{softmax}(Uh + b_2) \in \mathbb{R}^{|V|}$$

using softmax to get distribution

hidden layer

$$h = f(We + b_1)$$

$f$ : non-linearity (e.g. ReLU)

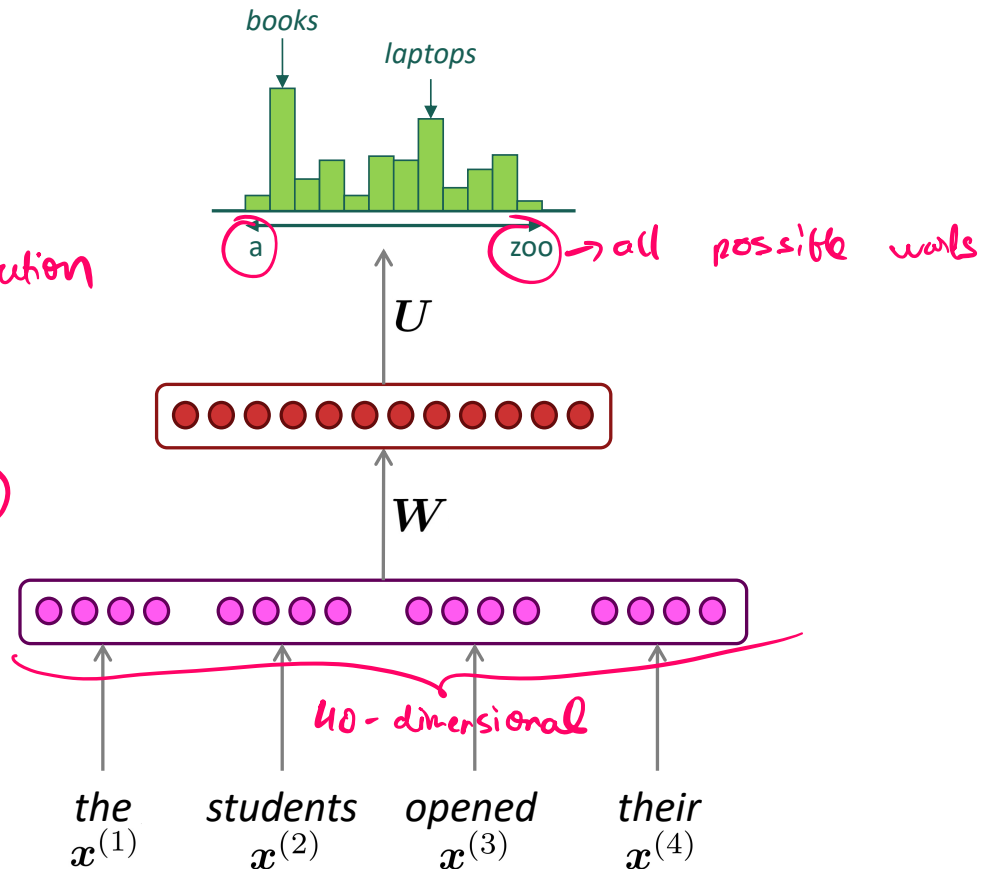
concatenated word embeddings

$$e = [e^{(1)}; e^{(2)}; e^{(3)}; e^{(4)}]$$

suppose each is 10-dim

words / one-hot vectors

$$x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$$



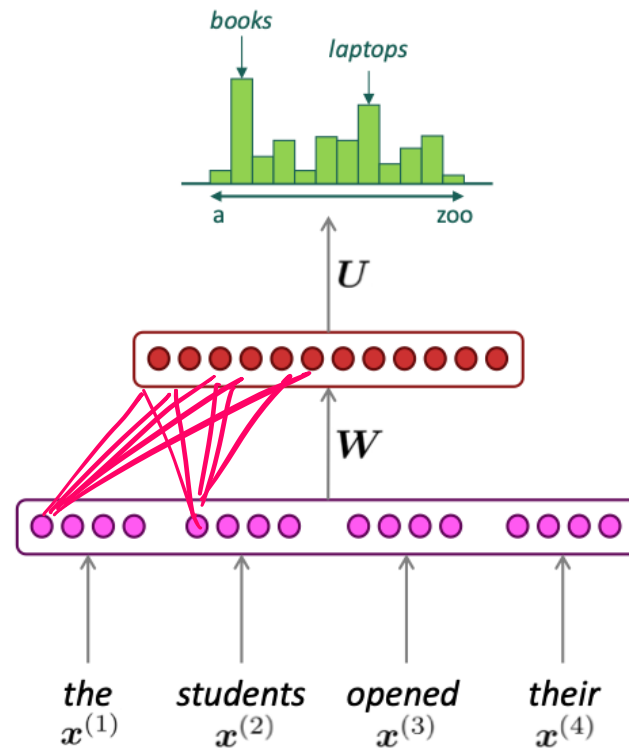
# The problem with this architecture

- Uses a fixed window, which can be too small.
- Enlarging this window will enlarge the size of the weight matrix  $W$ .
- The inputs  $x^{(1)}$  and  $x^{(2)}$  are multiplied by completely different weights in  $W$ .

**No symmetry in how inputs are processed!**

As with CNNs for images before, we need an architecture which has similar symmetries as the data.

In this case, *can we have an architecture that can process any input length?*

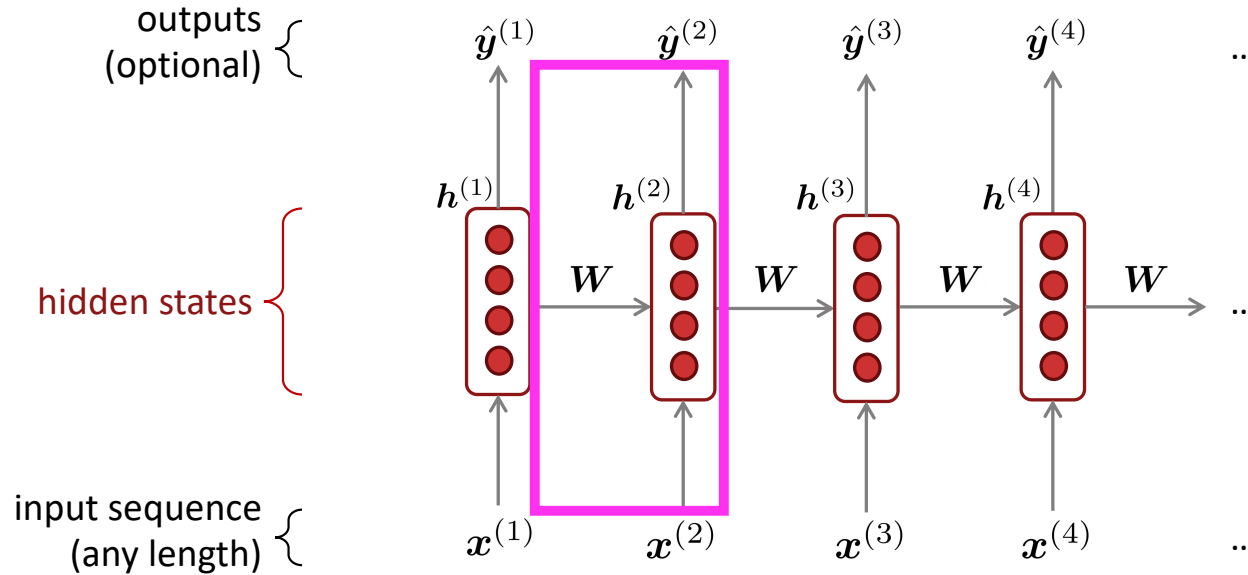


# Recurrent Neural Networks (RNN)

A family of neural architectures

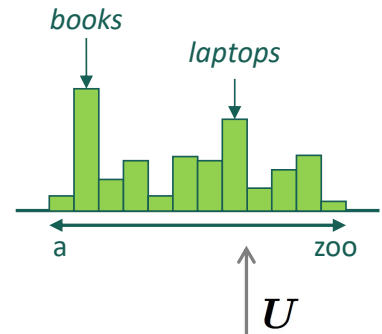
Core idea: Apply the same weights  $W$  repeatedly

similar to filters in CNNs



# A Simple RNN Language Model

$$\hat{y}^{(4)} = P(x^{(5)} | \text{the students opened their})$$



output distribution

$$\hat{y}^{(t)} = \text{softmax}(U\mathbf{h}^{(t)} + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$

linear prediction + softmax

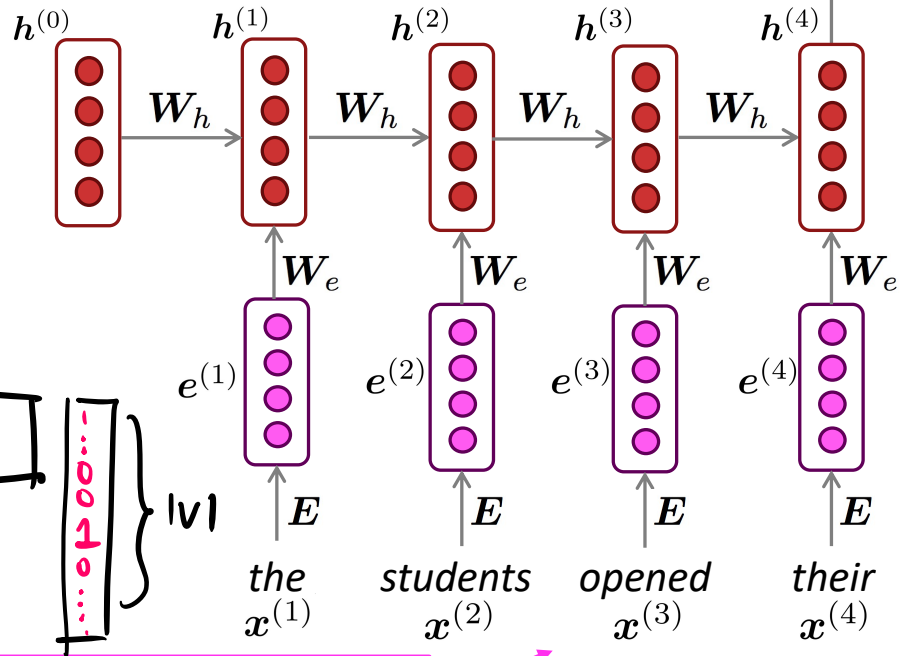
hidden states

$$\mathbf{h}^{(t)} = \sigma(W_h \mathbf{h}^{(t-1)} + W_e e^{(t)} + \mathbf{b}_1)$$

$\mathbf{h}^{(0)}$  is the initial hidden state

$\sigma$ : Activation (ReLU)

bias



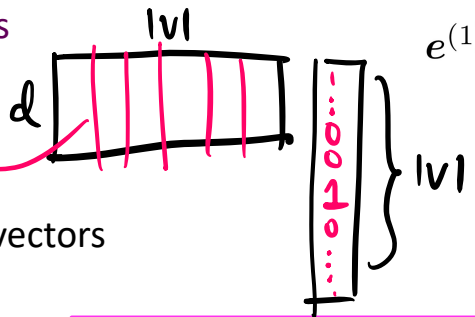
word embeddings

$$e^{(t)} = E\mathbf{x}^{(t)}$$

columns of  $E$  are word embeddings

words / one-hot vectors

$$\mathbf{x}^{(t)} \in \mathbb{R}^{|V|}$$

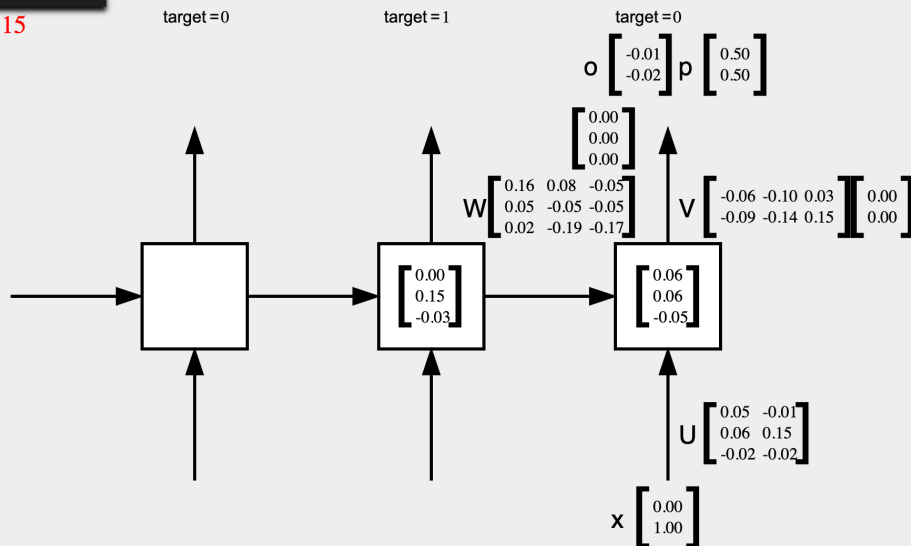


Note: this input sequence could be much longer now!

# Visualizing RNNs

```
#Initialize weights and biases
input_weights_U = np.random.randn(hidden_size, vocab_size) * 0.01
hidden_weights_W = np.random.randn(hidden_size, hidden_size) * 0.01
hidden_bias = np.zeros((hidden_size, 1))
output_weights_V = np.random.randn(vocab_size, hidden_size) * 0.01
output_bias = np.zeros((vocab_size, 1))
#Forward pass
xs, hidden_states, outputs, probabilities = {}, {}, {}, {}
loss = 0
hidden_states[-1] = np.copy(hidden_state_prev)
for t in range(len(inputs)):
    # one-hot-encoding the input character
    xs[t] = np.zeros((vocab_size,1))
    character = inputs[t]
    xs[t][character] = 1
    target = targets[t]
    # Compute hidden state
    hidden_states[t] = np.tanh(input_weights_U @ xs[t] + hidden_weights_W @ hidden_states[t-1] + hidden_bias)
    # Compute output and probabilities
    outputs[t] = output_weights_V @ hidden_states[t] + output_bias
    probabilities[t] = np.exp(outputs[t]) / np.sum(np.exp(outputs[t]))
    #Compute cross-entropy loss
    loss += -np.log(probabilities[t][target,0])
```

loss = 2.08315



<https://joshvarty.github.io/VisualizingRNNs/>

# Training an RNN Language Model

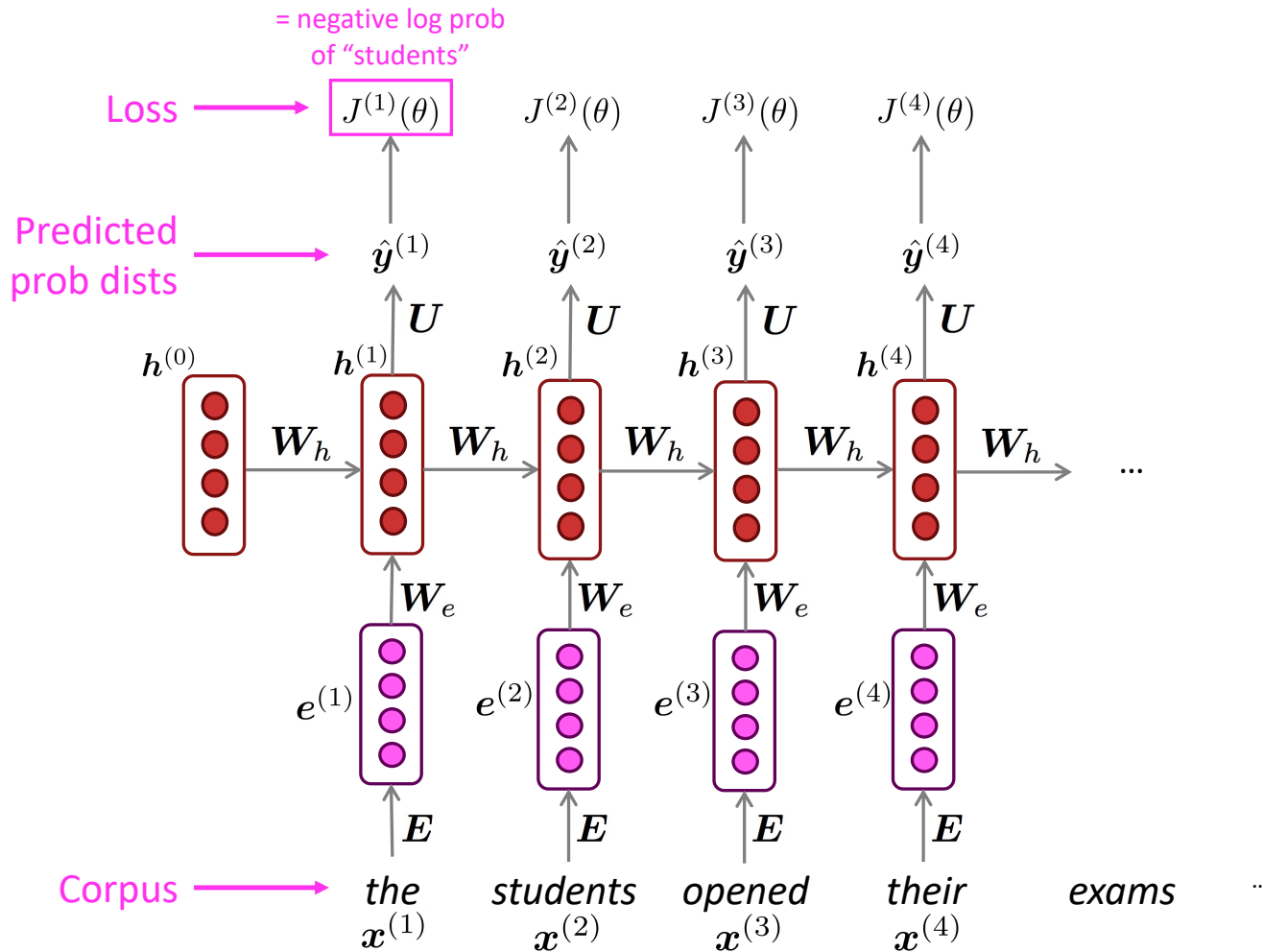
- Get a **big corpus of text** which is a sequence of words  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}$
- Feed into RNN-LM; compute output distribution  $\hat{\mathbf{y}}^{(t)}$  **for every step  $t$** .
  - i.e. predict probability distribution of *every word*, given words  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}$  so far
  - Suppose the parameters of the RNN-LM are  $\theta$ , and our goal is to learn  $\theta$
- **Loss function**  $J^{(t)}(\theta)$  on step  $t$  is **cross-entropy** between predicted probability distribution  $\hat{\mathbf{y}}^{(t)}$  and the true next word, which is  $\mathbf{x}^{(t+1)}$ :

$$\begin{aligned} J^{(t)}(\theta) &= \text{CrossEntropy}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) \\ &= -\log(\text{Predicted probability of next word being } \mathbf{x}^{(t+1)}) \\ &= -\log \hat{y}_{\mathbf{x}^{(t+1)}}^{(t)} \end{aligned}$$

- Average this to get **overall loss** for entire training set:

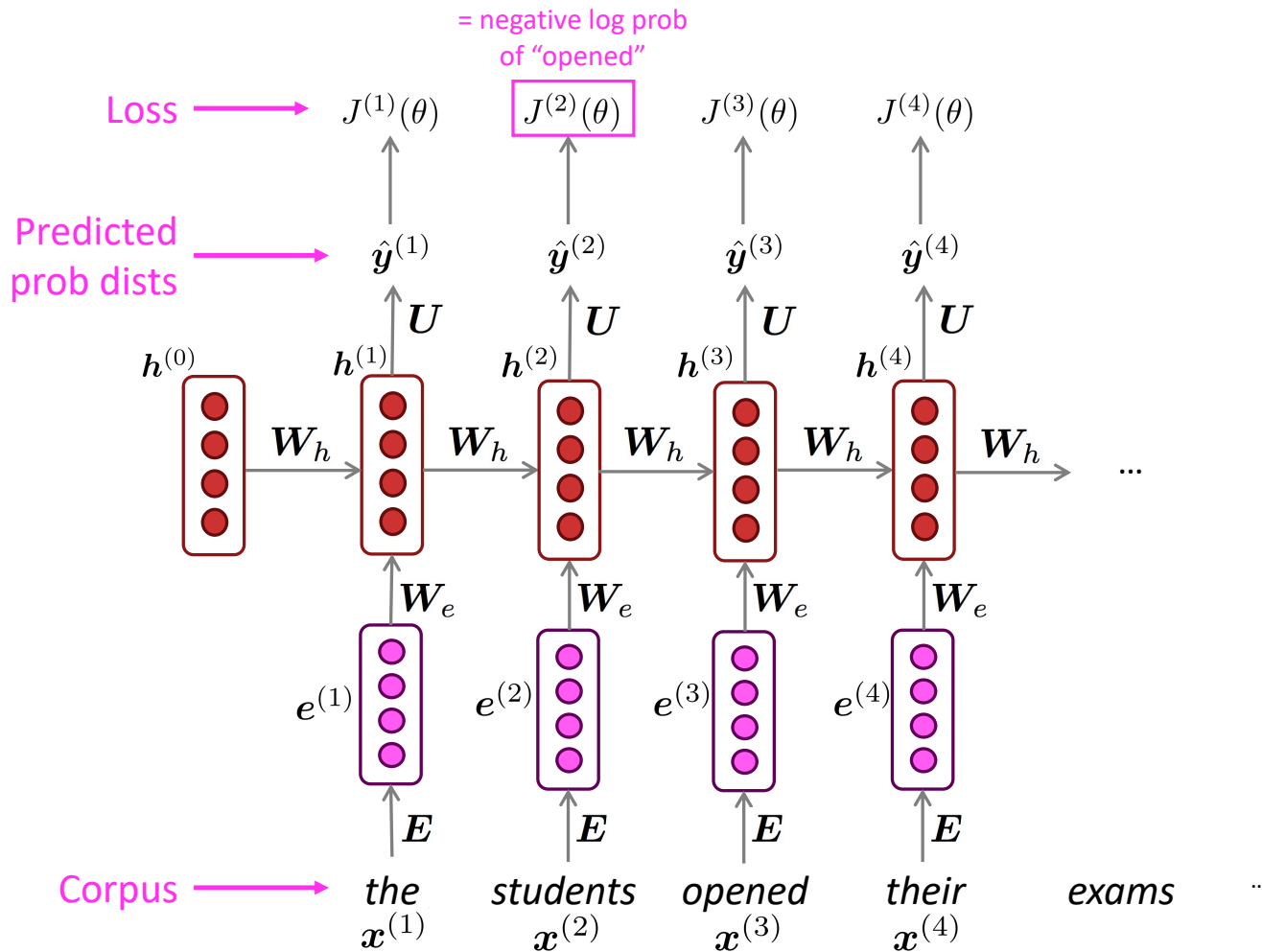
$$J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta) = \frac{1}{T} \sum_{t=1}^T -\log \hat{y}_{\mathbf{x}^{(t+1)}}^{(t)}$$

# Training an RNN Language Model



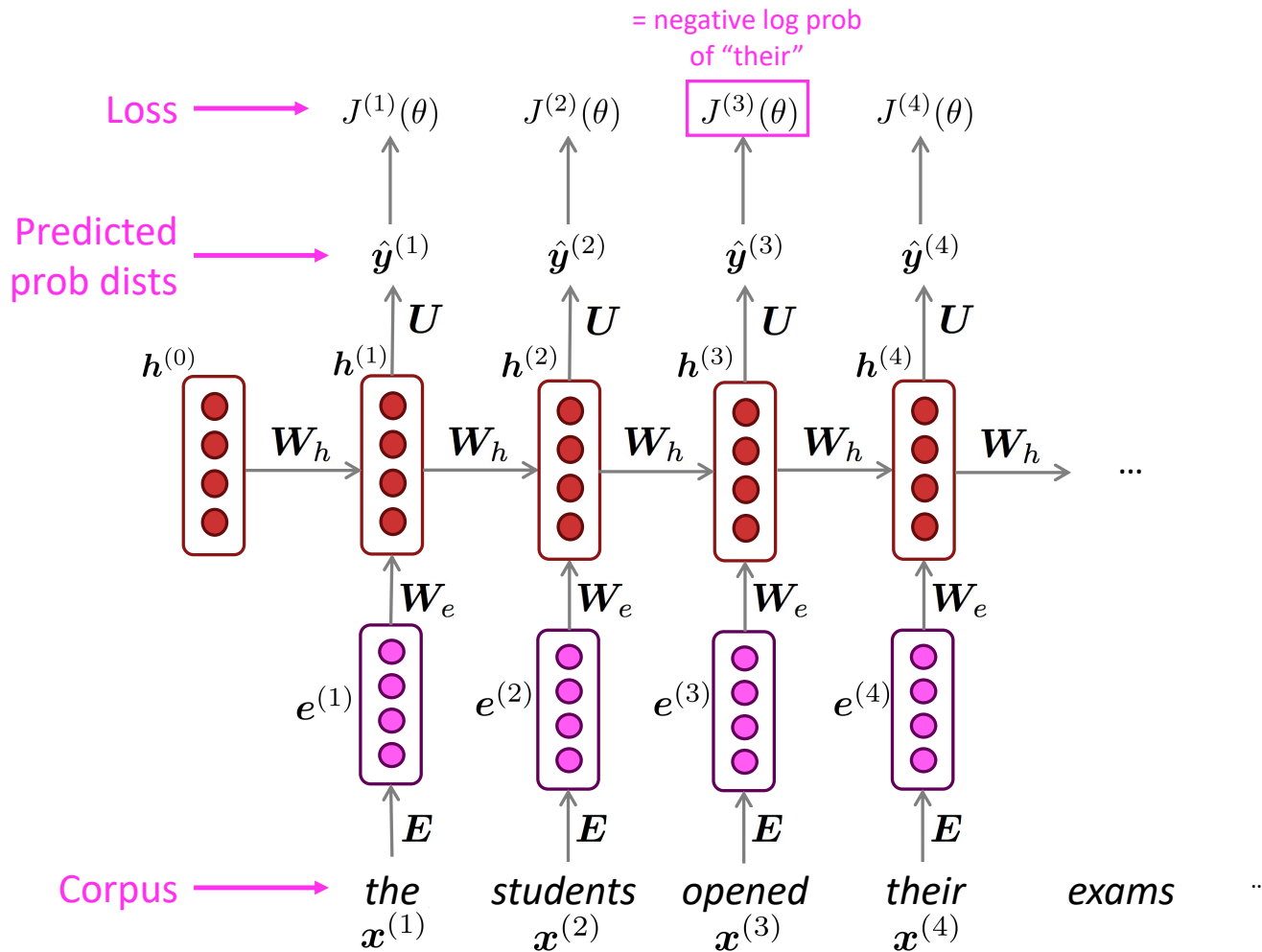
Slide adapted from CS224n by Chris Manning (Lecture 5)

# Training an RNN Language Model

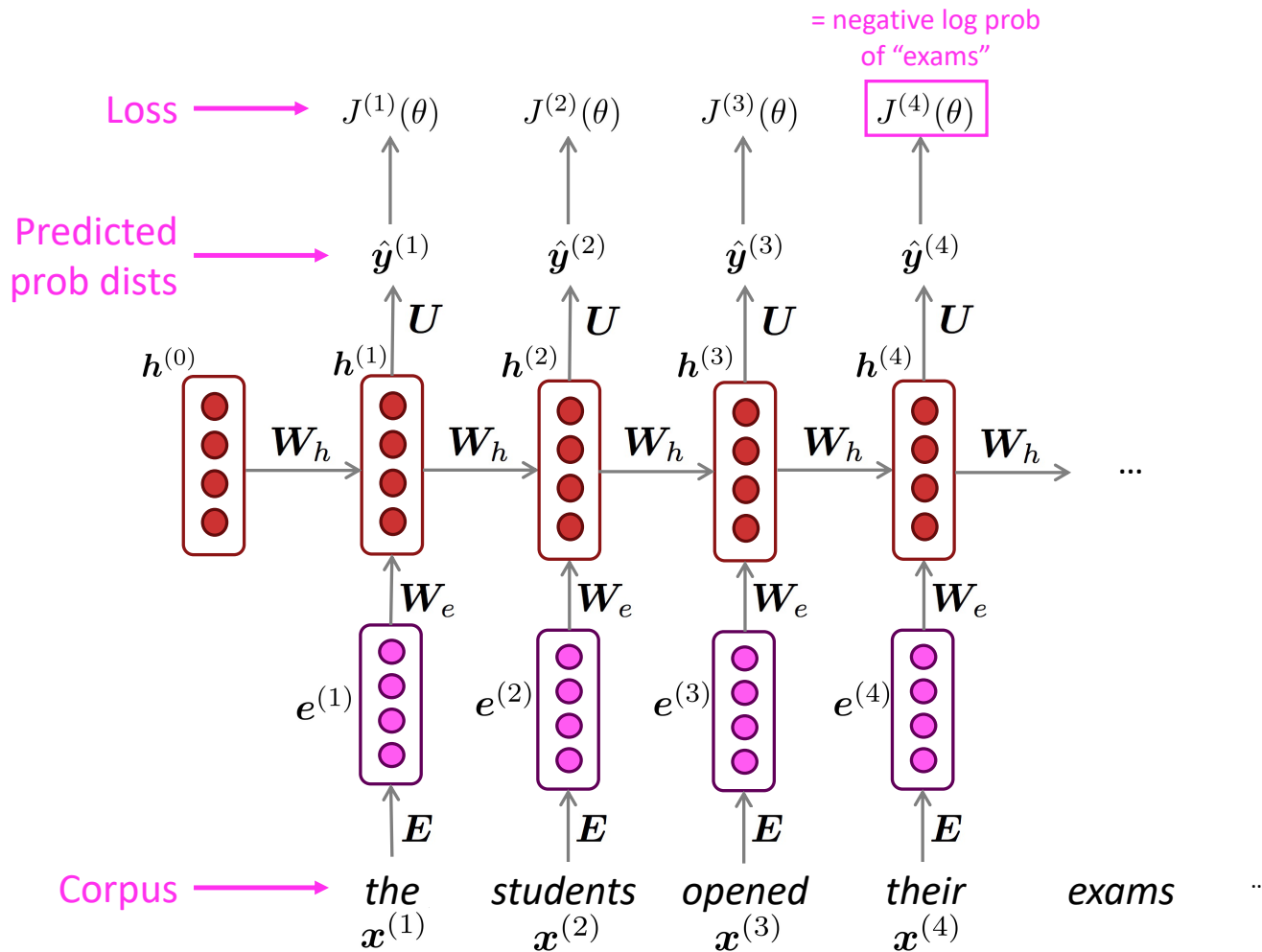


Slide adapted from CS224n by Chris Manning (Lecture 5)

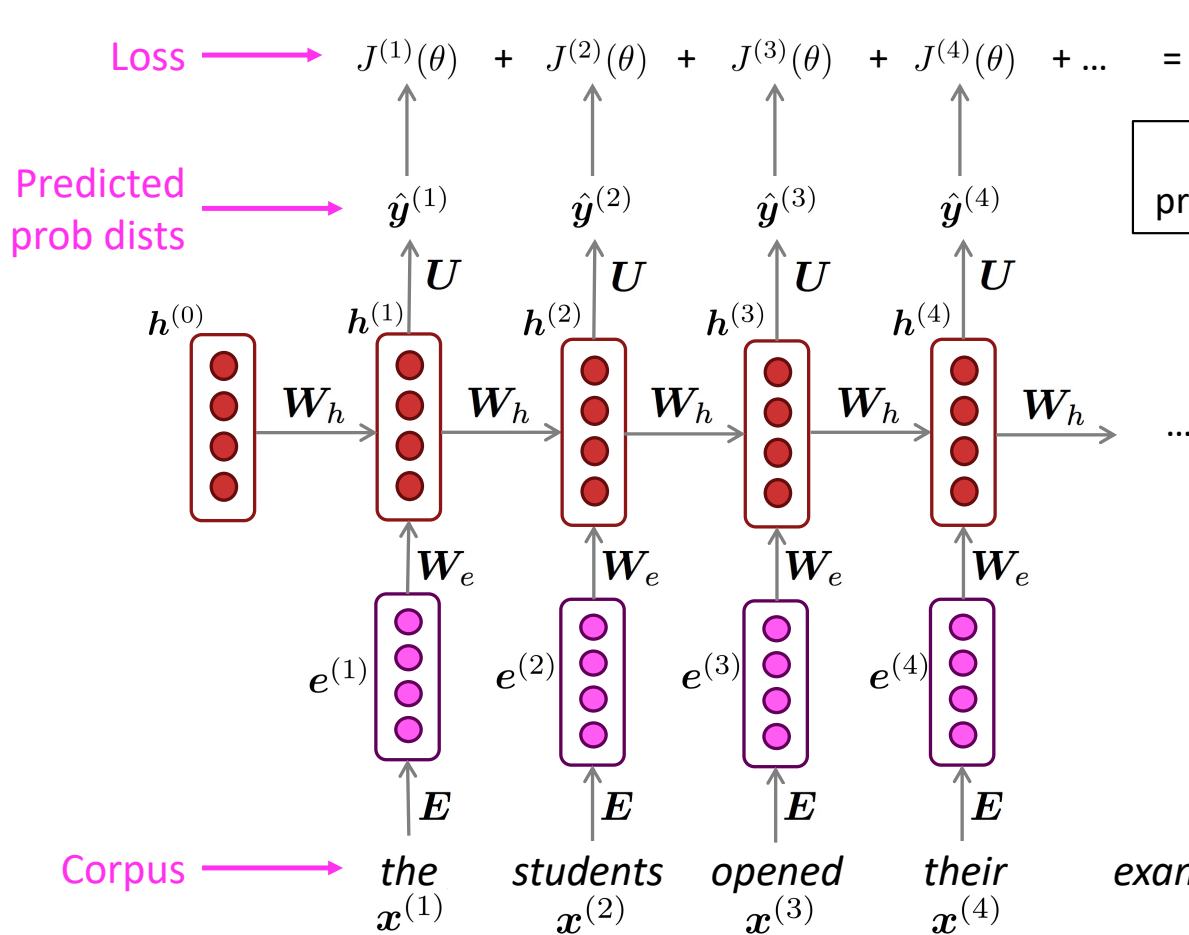
# Training an RNN Language Model



# Training an RNN Language Model



# Training an RNN Language Model



“Teacher forcing”

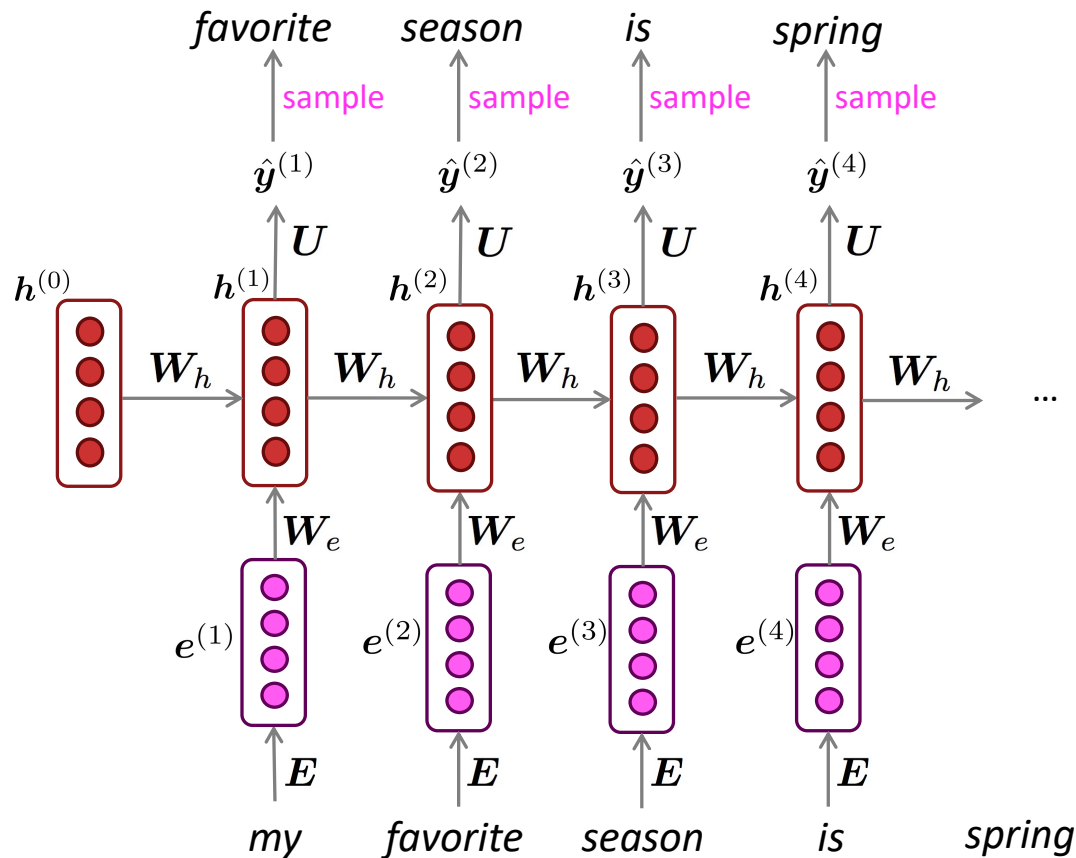
$$J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta)$$

Use true next word (instead of predicted next word) as input to model

How to train this?  
Backprop + SGD

# Generating text with a RNN Language Model

Just like a n-gram Language Model, you can use a RNN Language Model to **generate text** by **repeated sampling**. Sampled output becomes next step's input.



Slide adapted from  
CS224n by Chris  
Manning (Lecture 5)

# Example: A surprisingly effective character-level RNN

## The Unreasonable Effectiveness of Recurrent Neural Networks

May 21, 2015

VIOLA:

Why, Salisbury must find his flesh and thought  
That which I am not apt, not a man and in fire,  
To show the reining of the raven and the wars  
To grace my hand reproach within, and not a fair are hand,  
That Caesar and my goodly father's world;  
When I was heaven of presence and our fleets,  
We spare with hours, but cut thy council I am great,  
Murdered and by thy master's ready there  
My power to give thee but so much as hell:  
Some service in the noble bondman here,  
Would show him to her wine.

KING LEAR:

O, if you were a feeble sight, the courtesy of your law,  
Your sight and several breath, will wear the gods  
With his heads, and my hands are wonder'd at the deeds,  
So drop upon your lordship's head, and your opinion  
Shall be against your honour.

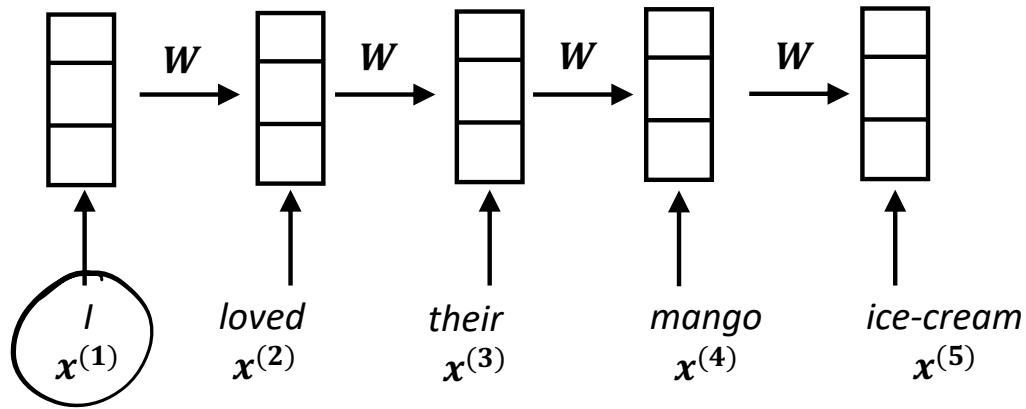
<https://karpathy.github.io/2015/05/21/rnn-effectiveness/>



Transformers

# The problem with recurrence

$$\approx Wx^{(1)} \quad \approx W^2x^{(1)} \quad \approx W^3x^{(1)} \quad \approx W^4x^{(1)}$$



1. Must always compress all necessary information into one hidden state representation
2. Cannot capture long-range dependencies in input (“vanishing gradients problem”)

Inputs from sufficiently far away do not contribute to hidden state representation:

Suppose 
$$W = \begin{pmatrix} 0.8 & 0.2 \\ -0.6 & 0.9 \end{pmatrix}$$

Then 
$$W^5 = \begin{pmatrix} -0.31 & 0.35 \\ -1.06 & -0.13 \end{pmatrix}, \quad W^{10} = \begin{pmatrix} -0.28 & -0.16 \\ 0.47 & -0.36 \end{pmatrix}, \quad W^{50} = \begin{pmatrix} 0.01 & 0.00 \\ -0.01 & 0.01 \end{pmatrix}$$

# A solution: Attention

---

## Attention Is All You Need

---

**Ashish Vaswani\***  
Google Brain  
avaswani@google.com

**Noam Shazeer\***  
Google Brain  
noam@google.com

**Niki Parmar\***  
Google Research  
nikip@google.com

**Jakob Uszkoreit\***  
Google Research  
usz@google.com

**Llion Jones\***  
Google Research  
llion@google.com

**Aidan N. Gomez\* †**  
University of Toronto  
aidan@cs.toronto.edu

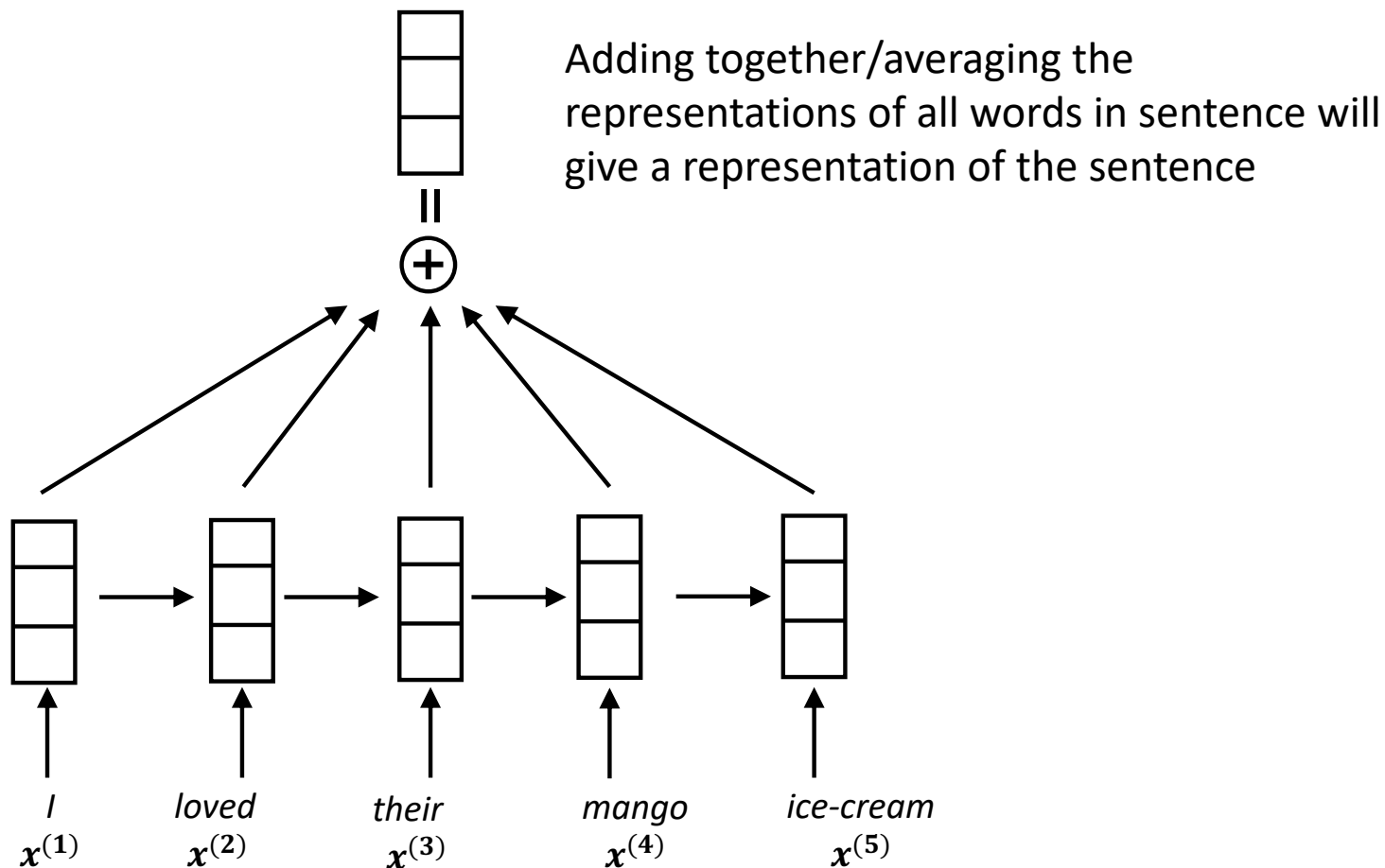
**Lukasz Kaiser\***  
Google Brain  
lukaszkaizer@google.com

**Illia Polosukhin\* ‡**  
illia.polosukhin@gmail.com

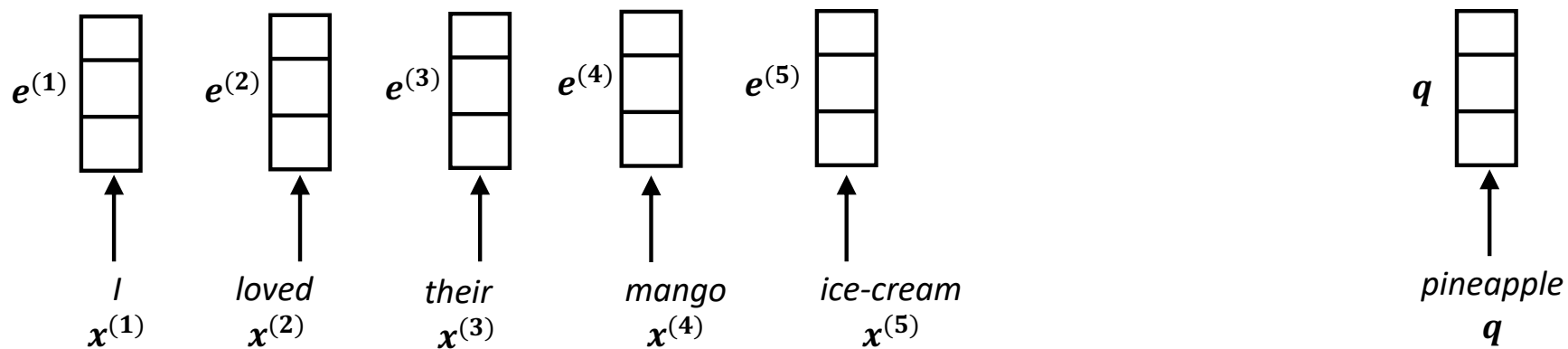
### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.0 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature.

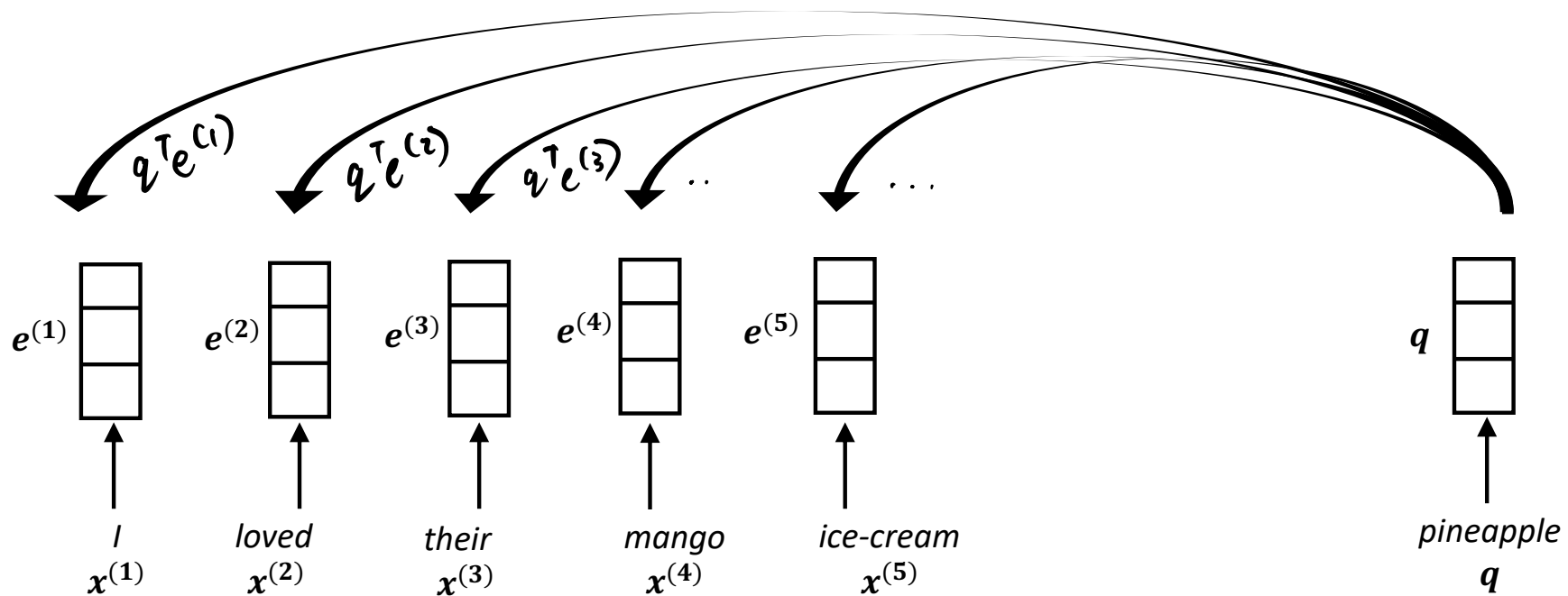
# Starting point: Averaging word representations



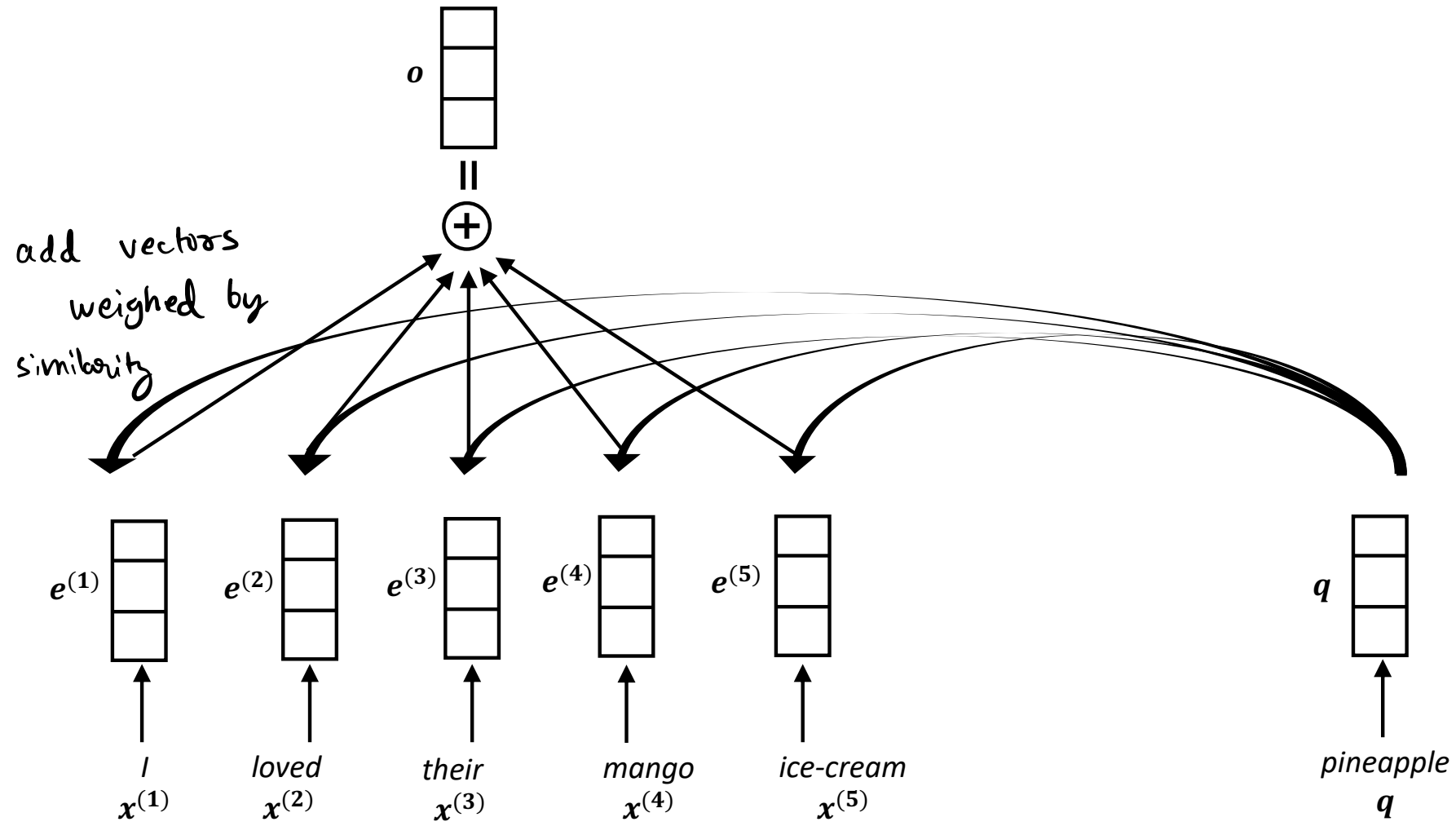
# Attention: Weighted averaging



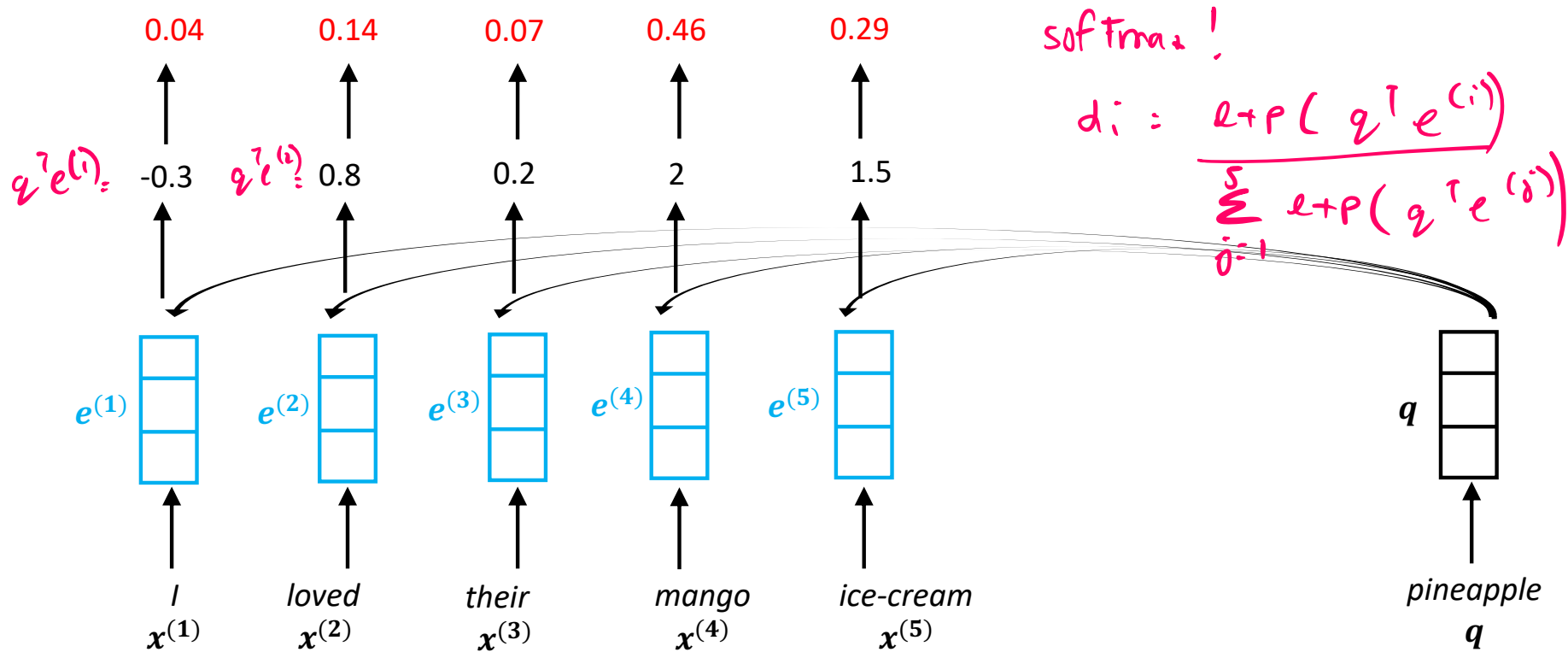
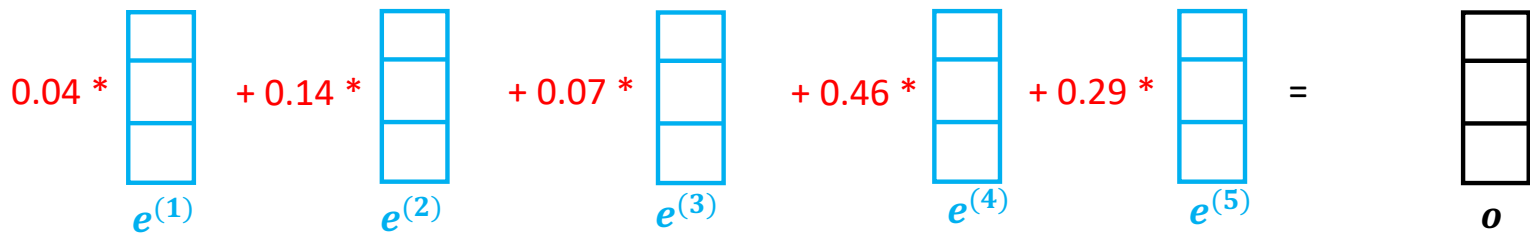
# Attention: Weighted averaging



# Attention: Weighted averaging

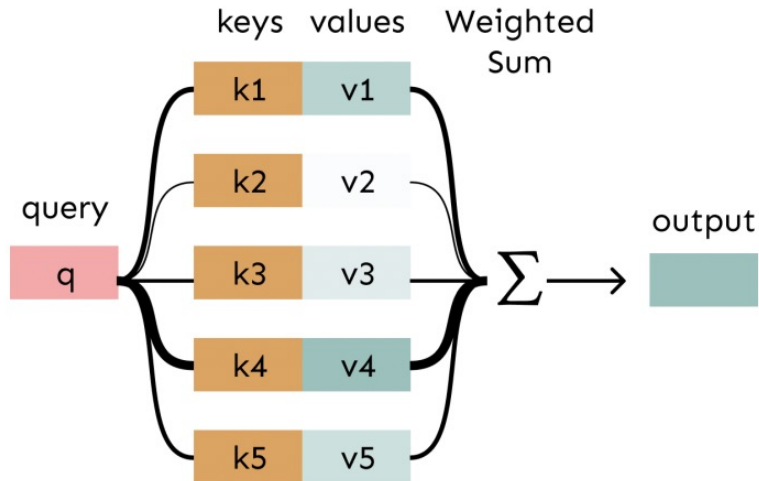


# Attention: Weighted averaging

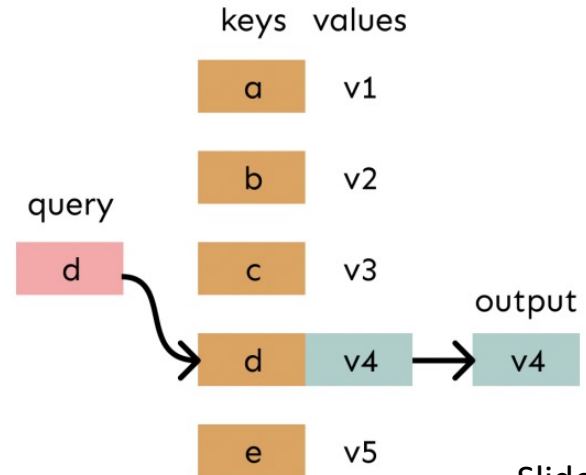


# Attention as soft lookup

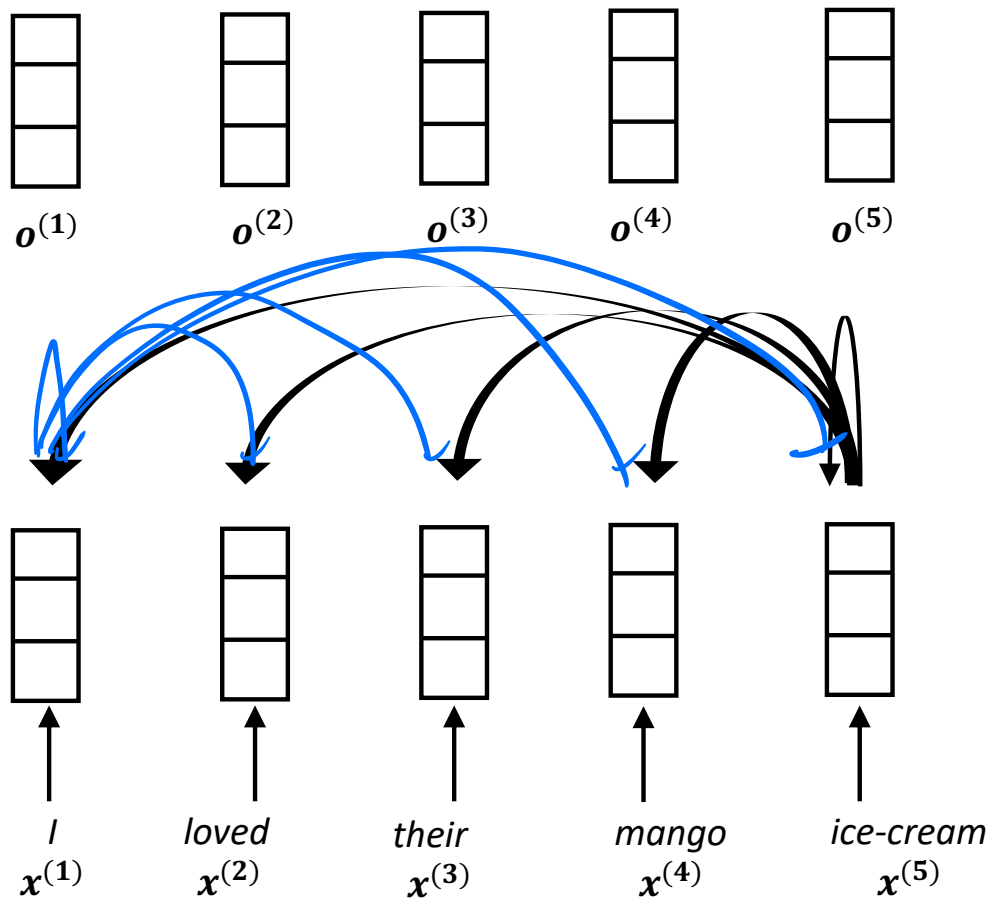
**Attention:** match query  $q$  to keys  $k_1, k_2, \dots, k_5$  to get weights between 0 and 1. Sum up values corresponding to each key with respective weight



**Lookup:** find query in database, return value corresponding to its key



# Self-attention

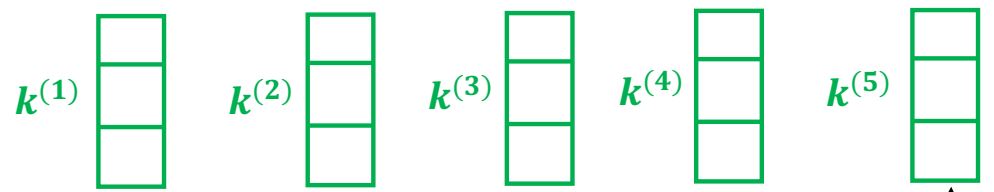


# Self-attention

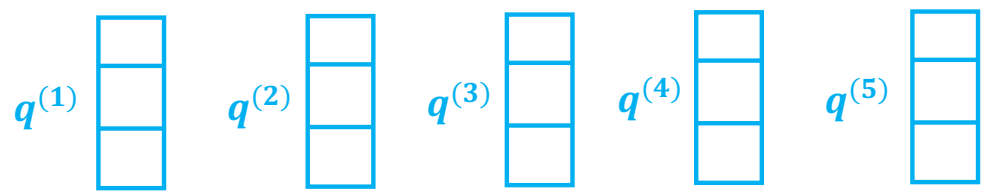
$$0.12 * v^{(1)} + 0.28 * v^{(2)} + 0.09 * v^{(3)} + 0.23 * v^{(4)} + 0.28 * v^{(5)} = o^{(5)}$$

0.12      0.28      0.09      0.23      0.28       $\rightarrow \text{softmax}$

$(q^{(5)})^T k^{(i)}$       0.1      1      -0.2      0.8      1      =  $(q^{(5)})^T k^{(5)}$



How to get  $k, q, v$ ?



$q^{(s)} = Q x^{(s)}$   
 $k^{(s)} = K x^{(s)}$   
 $v^{(s)} = V x^{(s)}$

$I$       loved      their      mango      ice-cream  
 $x^{(1)}$        $x^{(2)}$        $x^{(3)}$        $x^{(4)}$        $x^{(5)}$

# Self-attention in matrix form

1. Transform each word embedding with weight matrices  $\mathbf{Q}$ ,  $\mathbf{K}$ ,  $\mathbf{V}$ , each in  $\mathbb{R}^{d \times d}$

$$\mathbf{q}_i = \mathbf{Q} \mathbf{x}_i \quad (\text{queries})$$

$$\mathbf{k}_i = \mathbf{K} \mathbf{x}_i \quad (\text{keys})$$

$$\mathbf{v}_i = \mathbf{V} \mathbf{x}_i \quad (\text{values})$$

2. Compute pairwise similarities between keys and queries; normalize with softmax

$$\alpha_{ij} = \mathbf{q}_i^\top \mathbf{k}_j$$

$$w_{ij} = \frac{\exp(\alpha_{ij})}{\sum_{j'} \exp(\alpha_{ij'})}$$

3. Compute output for each word as weighted sum of values

$$\mathbf{o}_i = \sum_j w_{ij} \mathbf{v}_j$$