

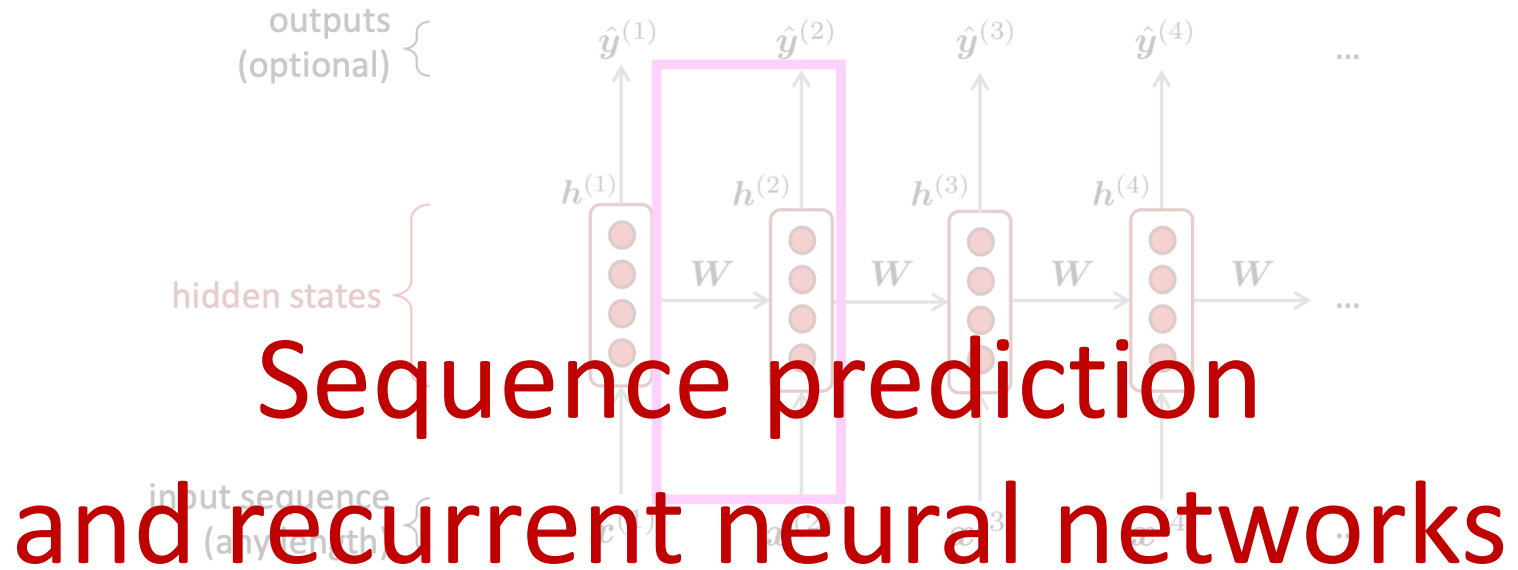
CSCI 567: Machine Learning

Vatsal Sharan
Spring 2024

Lecture 8, March 8

Administrivia

- HW3 due in less than 3 weeks
- No office hours next week due to spring break
- Project proposals due today on Gradescope & Google form
- Today's plan:
 - Sequential prediction, Markov models, recurrent neural networks, attention & Transformers




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 ?



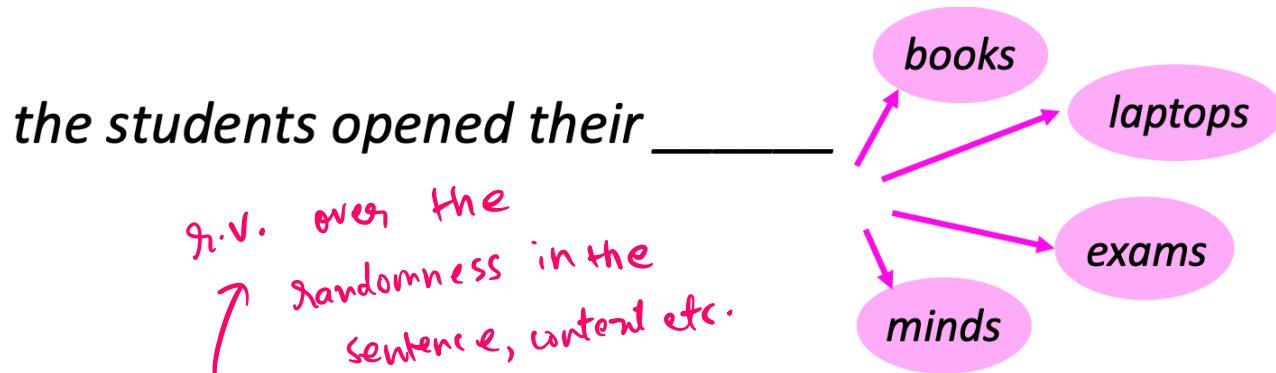
Examples:

- 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**.

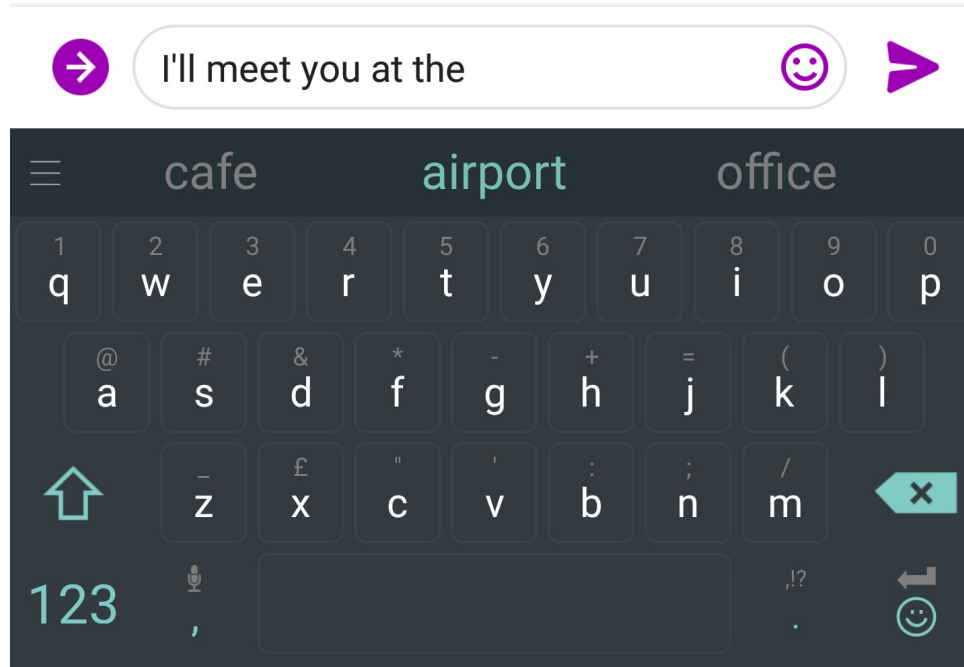
Language modelling

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:


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

You use Language Models every day!



You use Language Models every day!



what is the | 

- what is the **weather**
- what is the **meaning of life**
- what is the **dark web**
- what is the **xfi**
- what is the **doomsday clock**
- what is the **weather today**
- what is the **keto diet**
- what is the **american dream**
- what is the **speed of light**
- what is the **bill of rights**

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 $\{1, \dots, S\}$
- $P(X_{t+1} = s' | X_t = s) = a_{s,s'}$, known as **transition probability**

the size of dictionary of all possible words

- $P(X_1 = s) = \pi_s$ → *initial probability*

- $(\{\pi_s\}, \{a_{s,s'}\}) = (\boldsymbol{\pi}, \mathbf{A})$ are **parameters of the model** (s, s') entry of \mathbf{A} is $a_{s,s'}$

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

Markov model: examples

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

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

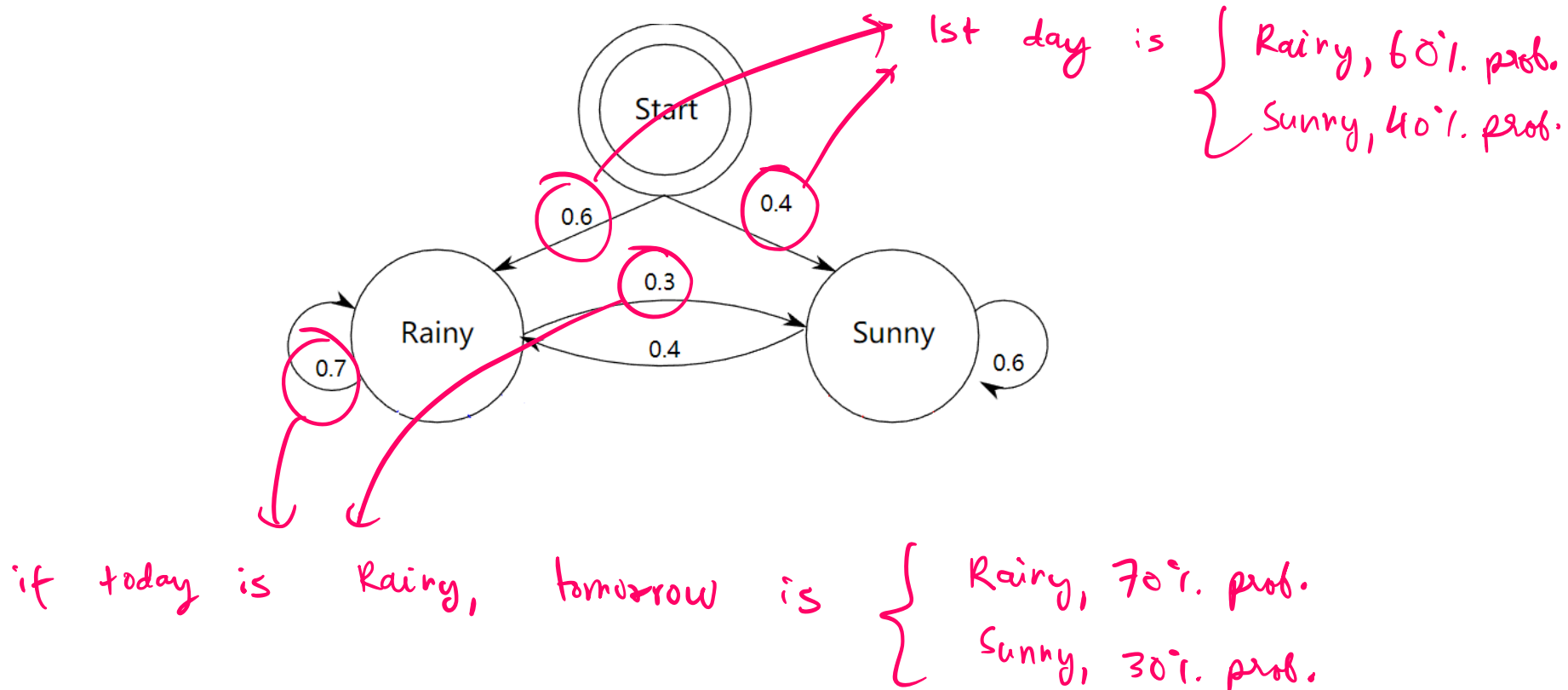
is an example of the transition probability.

- 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**



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* (π, \mathbf{A}) ?

Learning Markov models: MLE

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

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

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

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

Prob. 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

s, s' entry is a_{s, s'}

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

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,

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

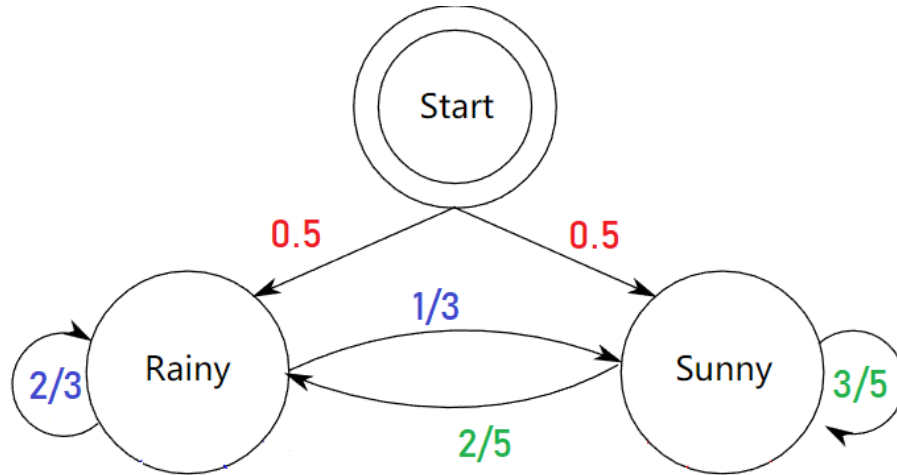
The initial state distribution follows similarly,

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

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):

$$P(X_{t+1} \mid X_t, \dots, X_1) = P(X_{t+1} \mid \overbrace{X_t, X_{t-1}, \dots, X_{t-n+2}}^{\text{previous } n \text{ observations}}) \quad (n\text{-th order Markov assumption})$$

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

$$\begin{aligned} P(X_{t+1} = x_{t+1} \mid X_t = x_t, \dots, X_1 = x_1) &= P(X_{t+1} = x_{t+1} \mid X_t = x_t, X_{t-1} = x_{t-1}, \dots, X_{t-n+2} = x_{t-n+2}) \\ &= \frac{P(X_{t+1} = x_{t+1}, X_t = x_t, X_{t-1} = x_{t-1}, \dots, X_{t-n+2} = x_{t-n+2})}{P(X_t = x_t, X_{t-1} = x_{t-1}, \dots, X_{t-n+2} = x_{t-n+2})} \\ &\approx \frac{\text{count}(x_{t-n+2}, \dots, x_{t-1}, x_t, x_{t+1}) \text{ in the data}}{\text{count}(x_{t-n+2}, \dots, x_{t-1}, x_t) \text{ in the data}} \end{aligned}$$

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 _____

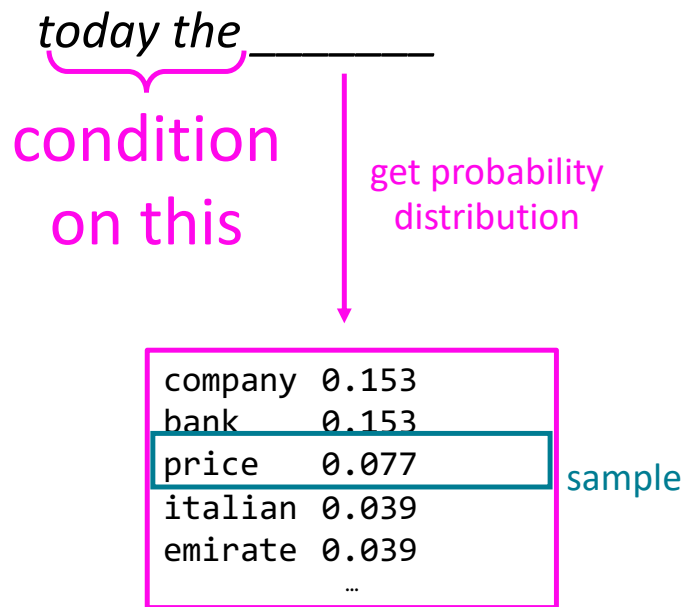
get probability distribution

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.

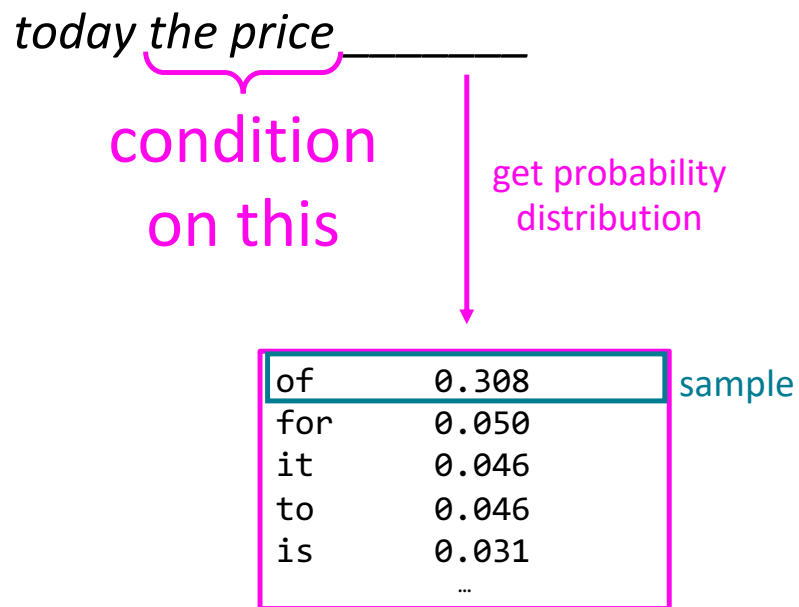
Generating text with a n-gram Language Model

You can also use a Language Model to generate text



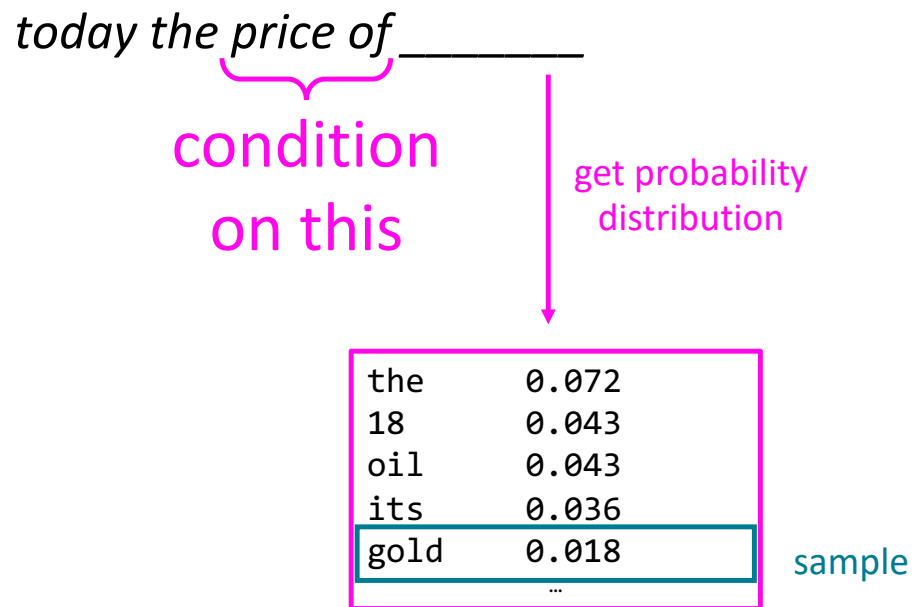
Generating text with a n-gram Language Model

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

w_1 w_2 w_3 \neq

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 $x^{(1)}, x^{(2)}, \dots, x^{(t)}$

- Output: prob dist of the next word $P(x^{(t+1)} | x^{(t)}, \dots, x^{(1)})$

$x^{(t)}$ refer to both

g.v. & the value
it takes

- How about a **window-based neural model**?

Word embeddings/vectors

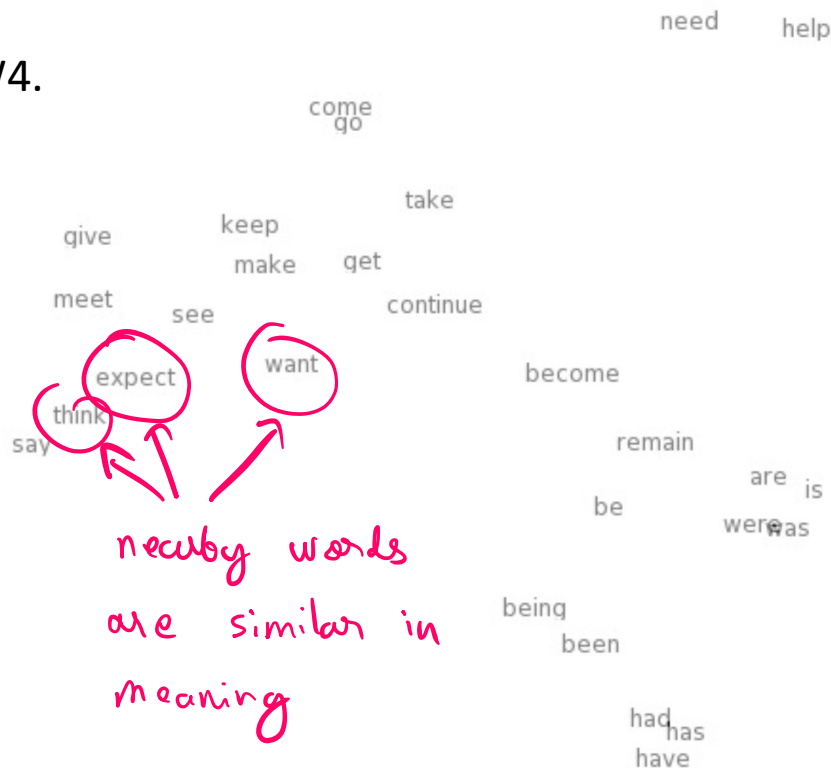
A word embedding is a (dense) mapping from words, to vector representations of the words.

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.

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

10-dim



A fixed-window neural Language Model

Same as in HW3 architecture

output distribution

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

hidden layer

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

f : non-linearity (ReLU)

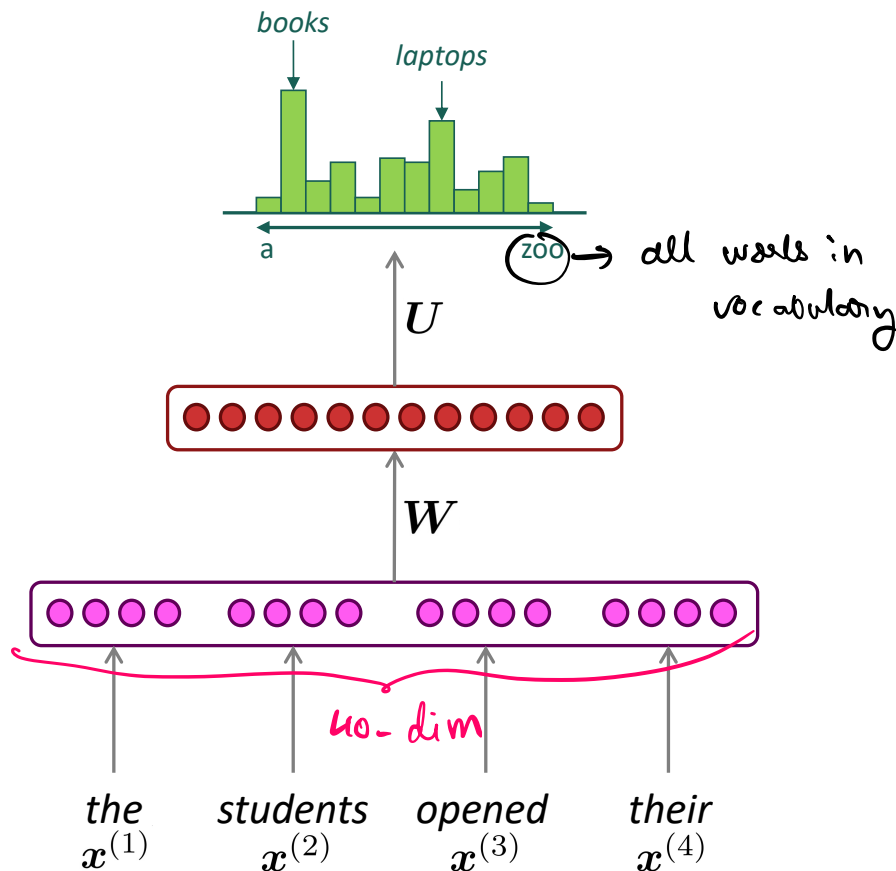
concatenated word embeddings

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

suppose each is 10-dim

words / one-hot vectors

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



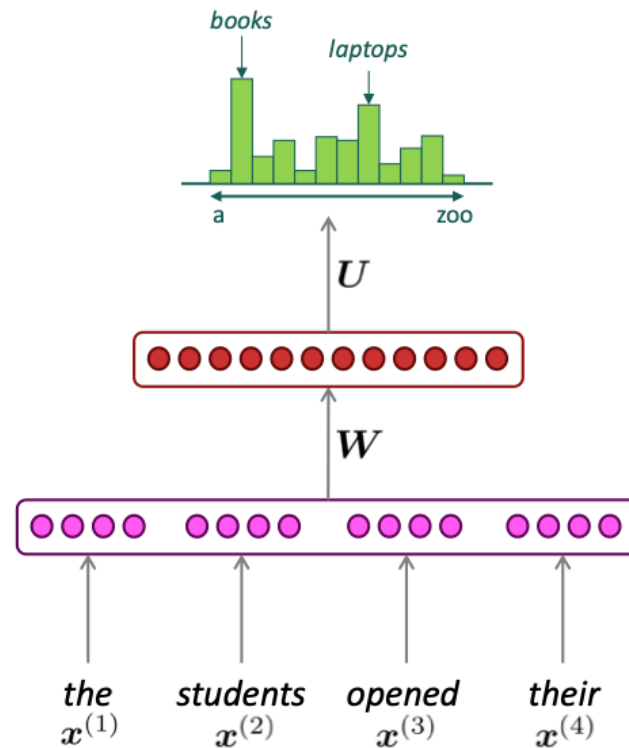
The problem with this architecture

- 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?*

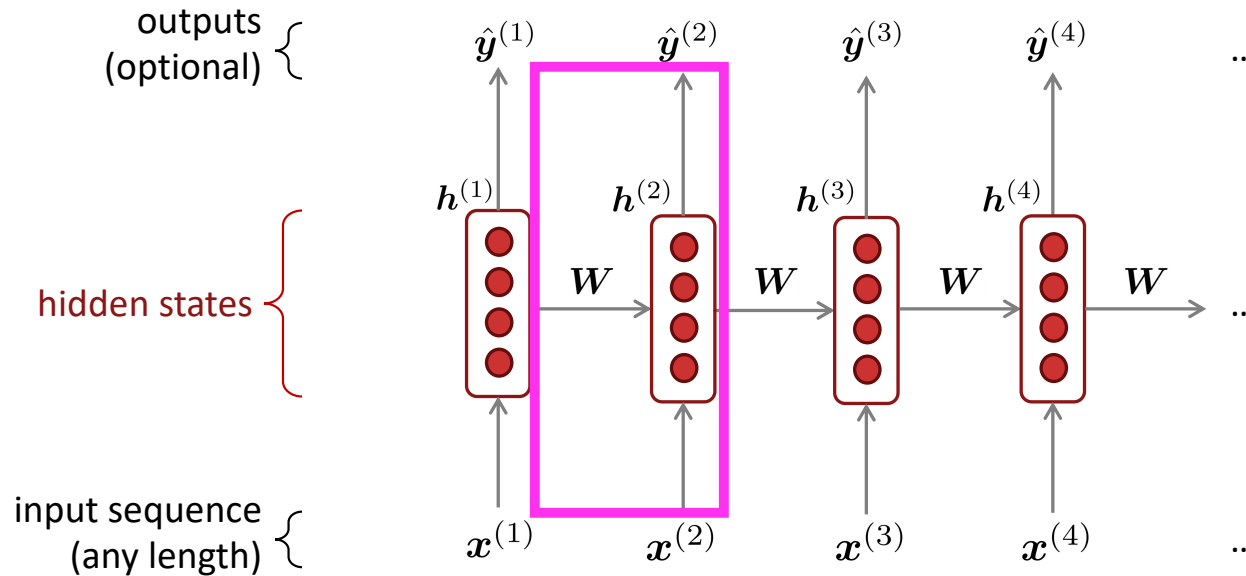


Recurrent Neural Networks (RNN)

A family of neural architectures

Core idea: Apply the same weights W repeatedly

↓
similar to filters

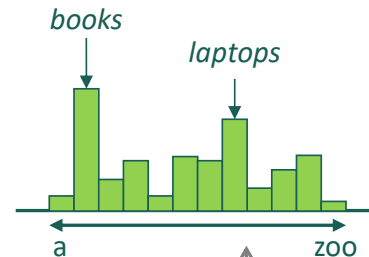


A Simple RNN Language Model

output distribution

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

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



hidden states

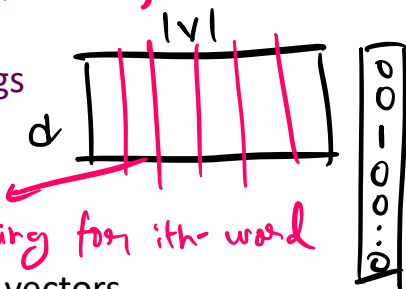
$$\mathbf{h}^{(t)} = \sigma(W_h \mathbf{h}^{(t-1)} + W_e \mathbf{e}^{(t)} + \mathbf{b}_1)$$

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

σ : Activation (ReLU)

word embeddings

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

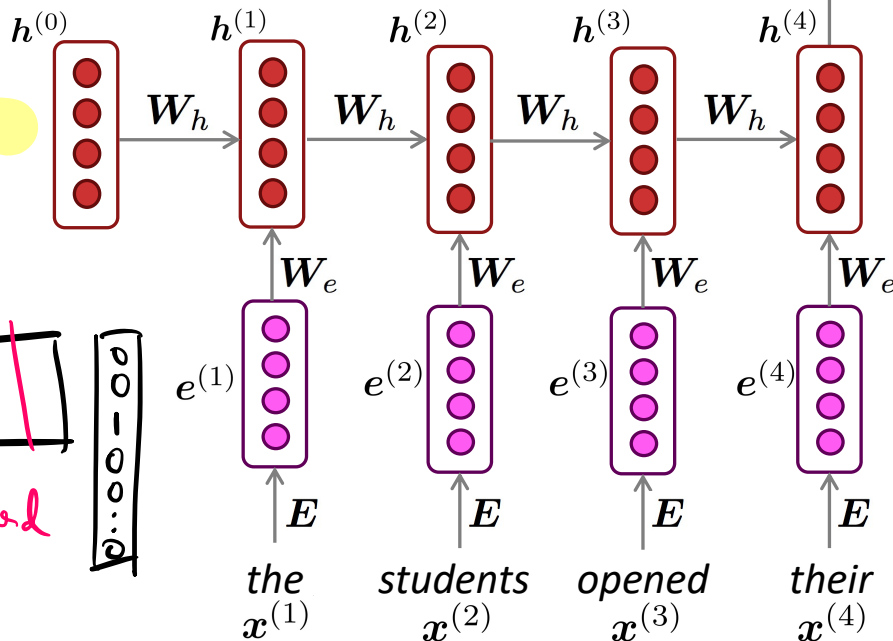


i -th column is embedding for i -th word

words / one-hot vectors

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

Note: this input sequence could be much longer now!



Training an RNN Language Model

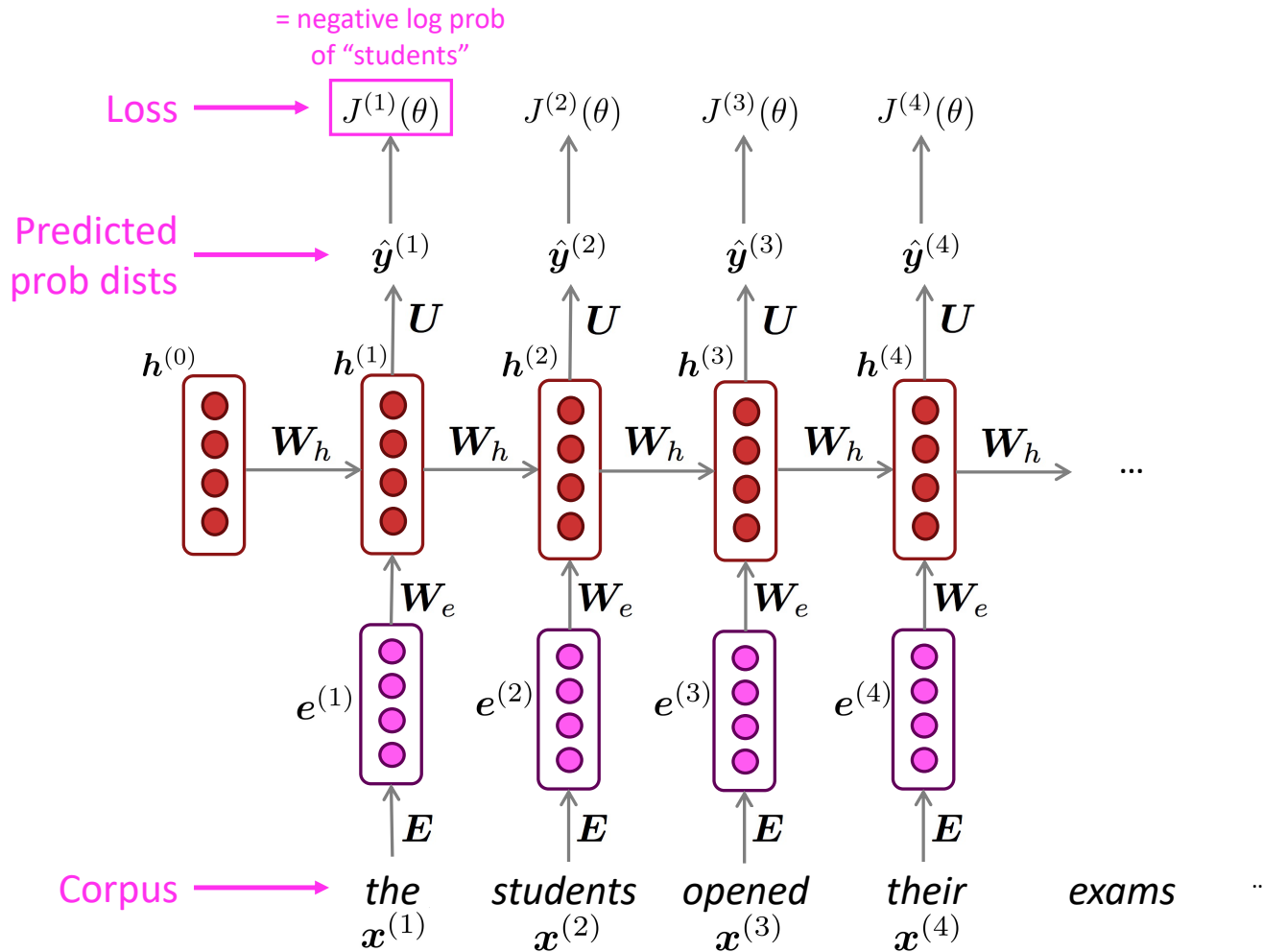
- Get a **big corpus of text** which is a sequence of words $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}$
- Feed into RNN-LM; compute output distribution $\hat{\mathbf{y}}^{(t)}$ **for every step t** .
 - i.e. predict probability dist of *every word*, given words so far
- **Loss function** on step t is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

$$J^{(t)}(\theta) = CE(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = - \sum_{w \in V} \mathbf{y}_w^{(t)} \log \hat{\mathbf{y}}_w^{(t)} = - \log \hat{\mathbf{y}}_{\mathbf{x}_{t+1}}^{(t)}$$

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

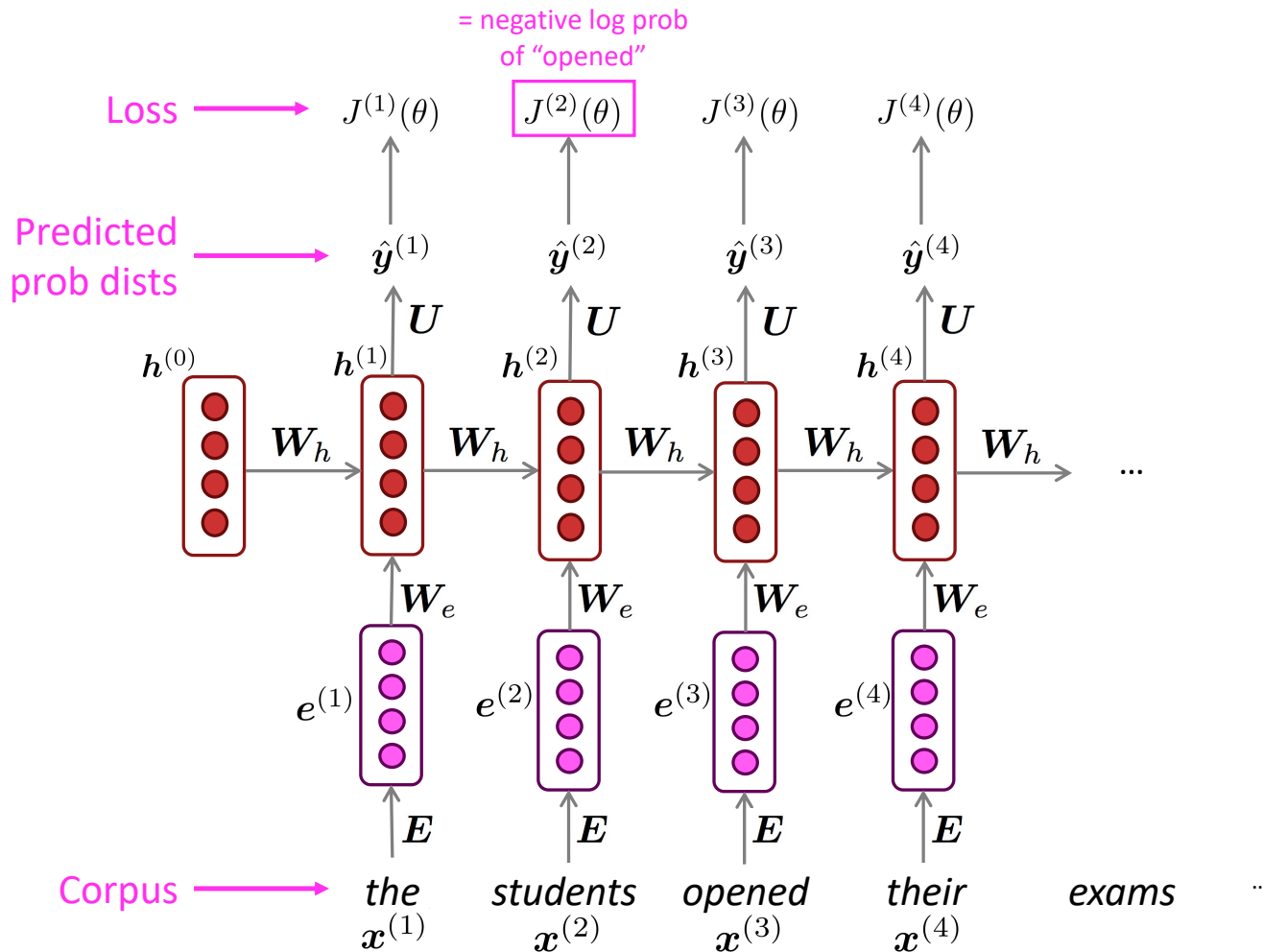
$$J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta) = \frac{1}{T} \sum_{t=1}^T - \log \hat{\mathbf{y}}_{\mathbf{x}_{t+1}}^{(t)}$$

Training an RNN Language Model



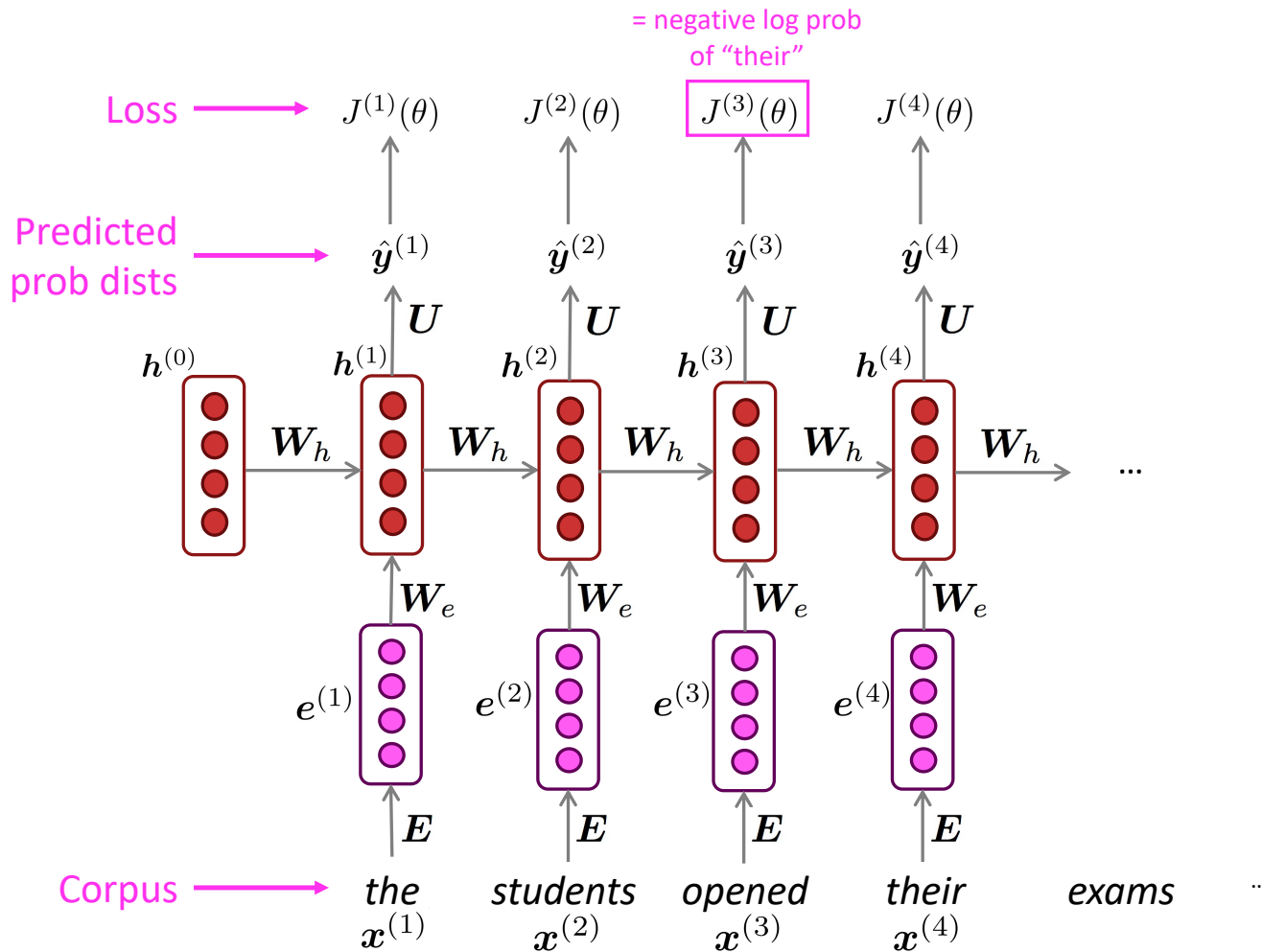
Slide adapted from CS224n by Chris Manning (Lecture 5)

Training an RNN Language Model

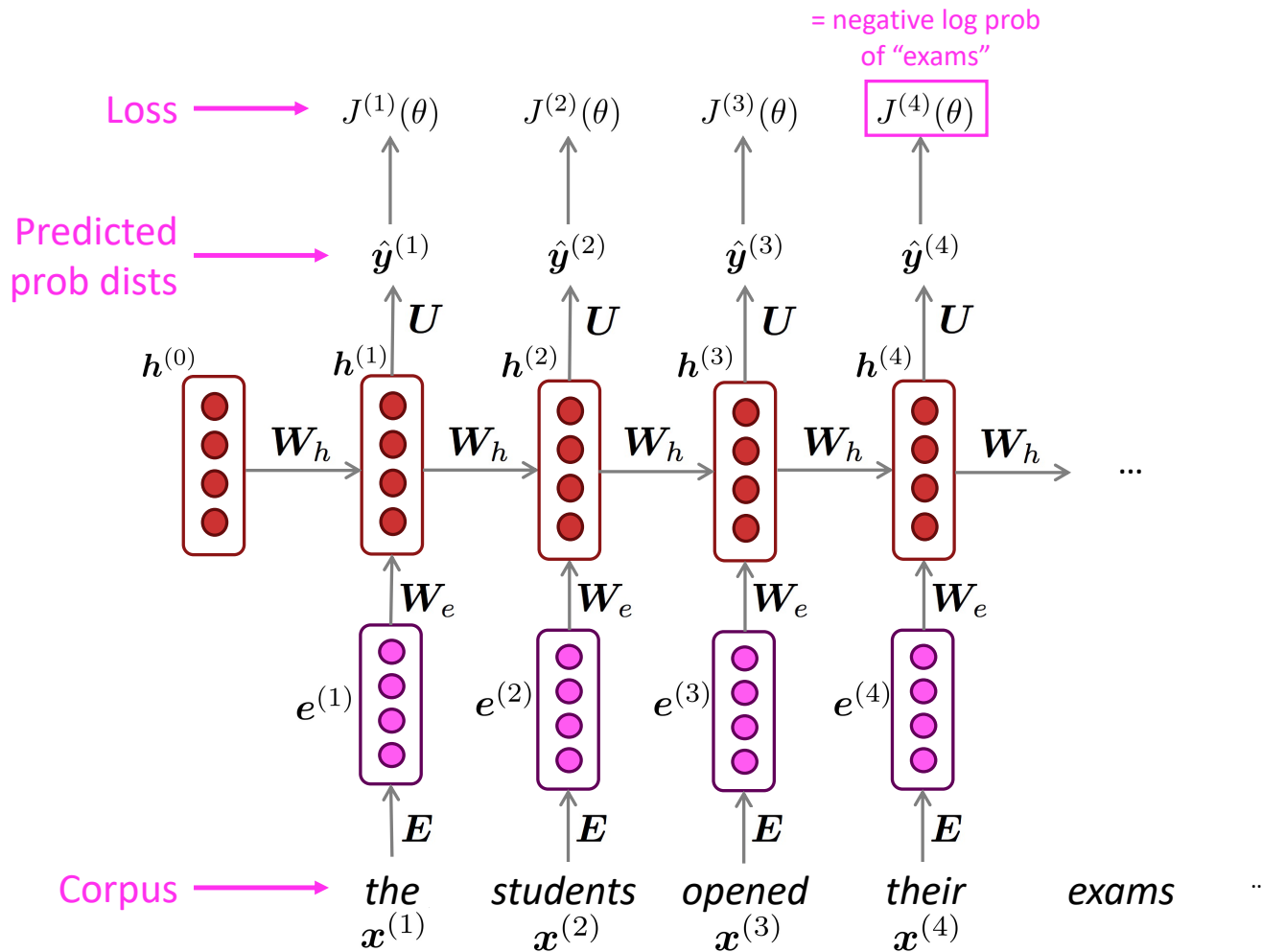


Slide adapted from CS224n by Chris Manning (Lecture 5)

Training an RNN Language Model



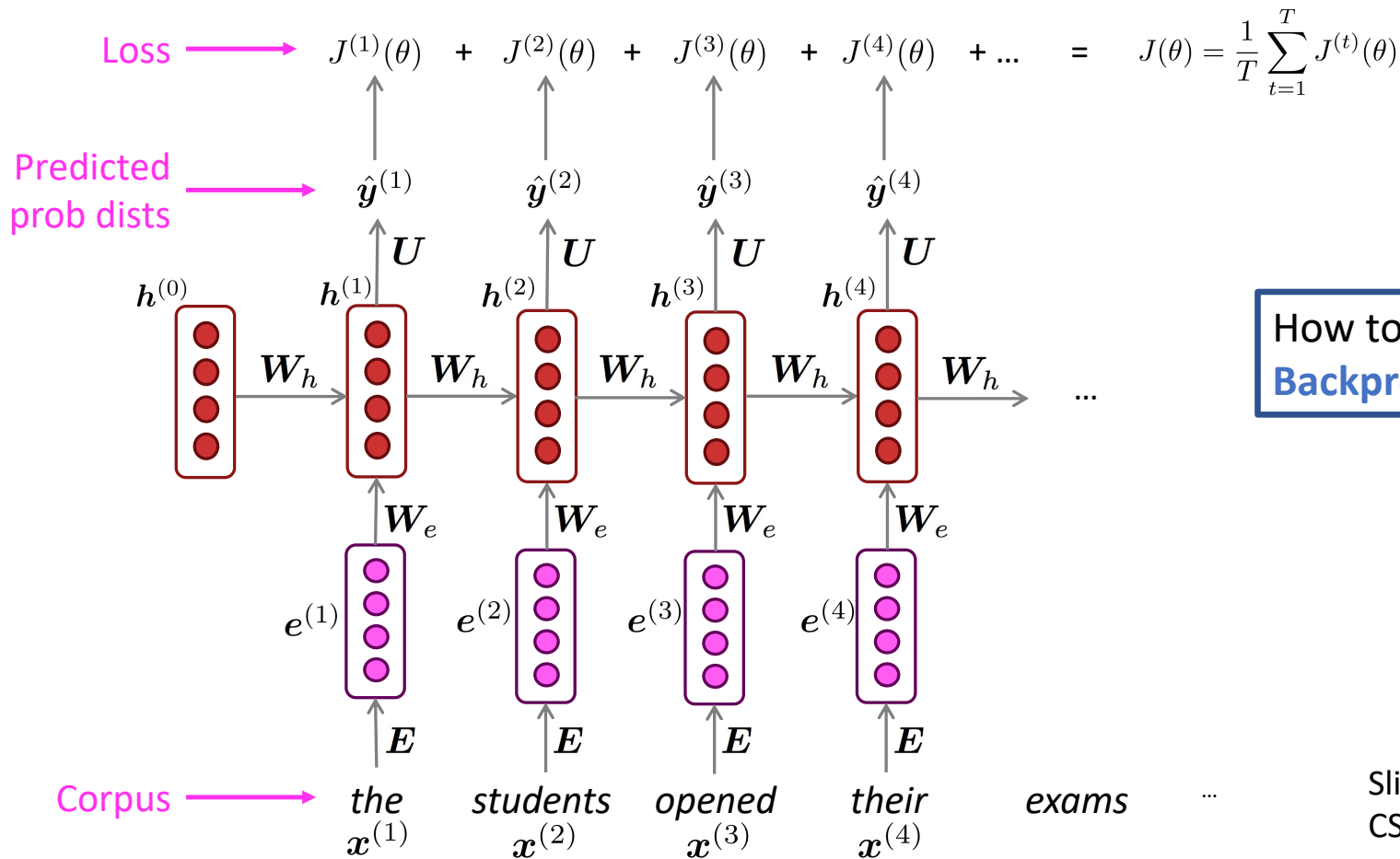
Training an RNN Language Model



Slide adapted from CS224n by Chris Manning (Lecture 5)

Training an RNN Language Model

“Teacher forcing”

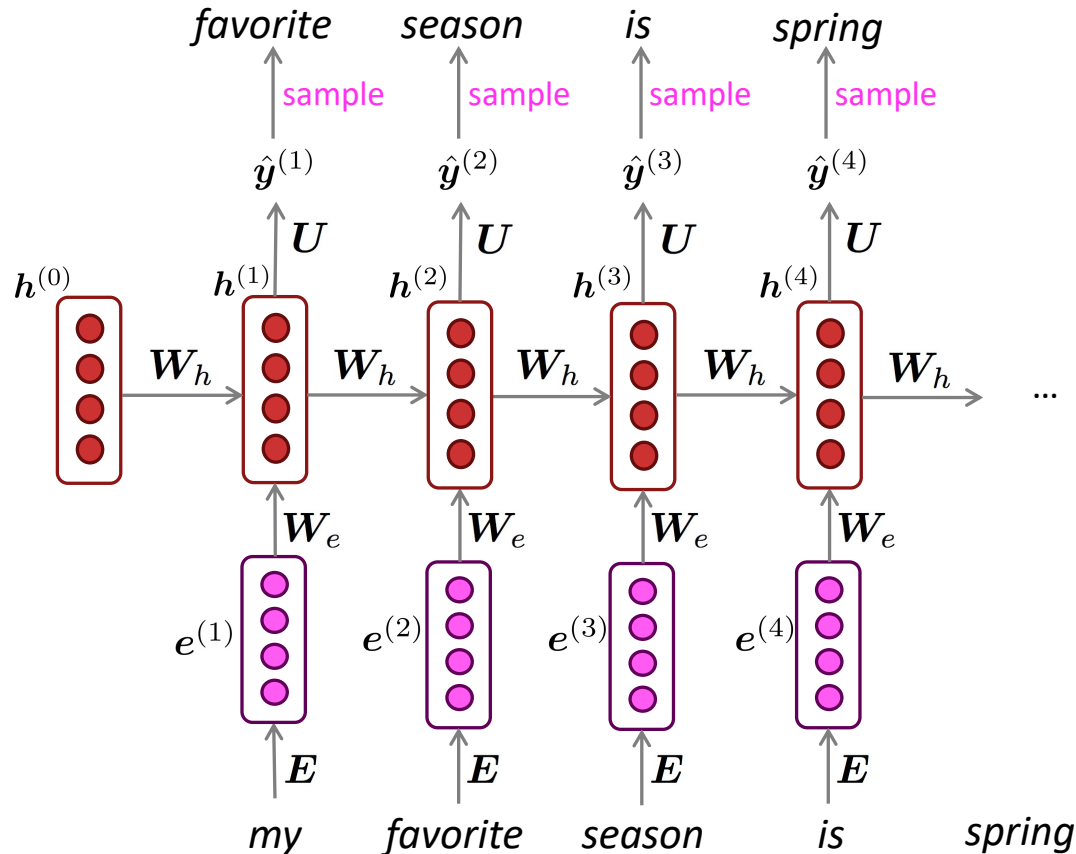


How to train this?
Backprop + SGD

Slide adapted from CS224n by Chris Manning (Lecture 5)

Generating text with a RNN Language Model

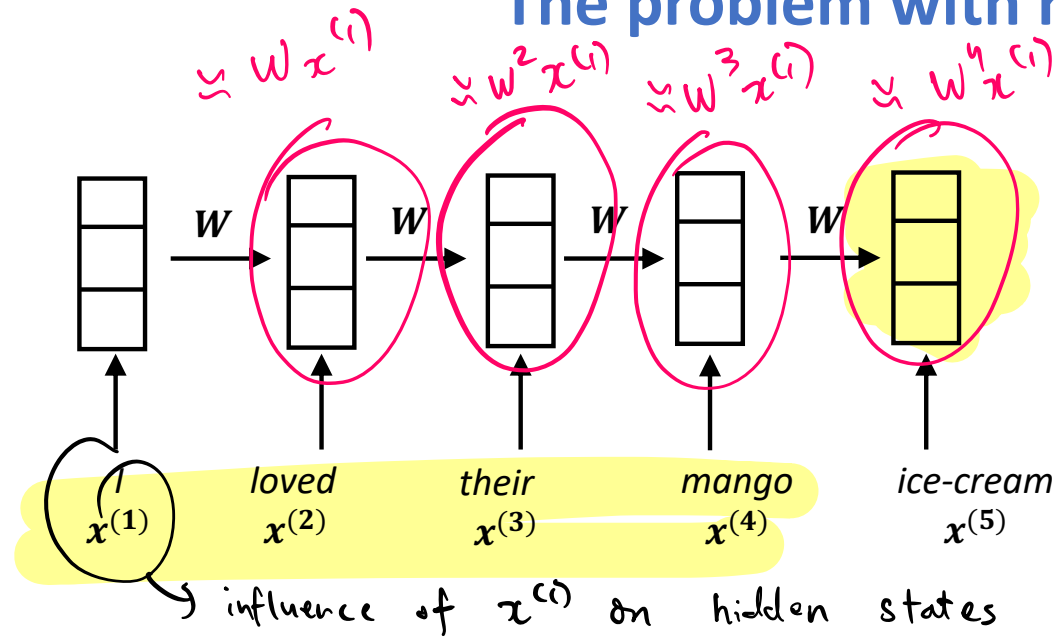
Just like a n-gram Language Model, you can use a RNN Language Model to **generate text** by **repeated sampling**. Sampled output becomes next step's input.





Transformers

The problem with recurrence



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

VSC
(connection :))

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

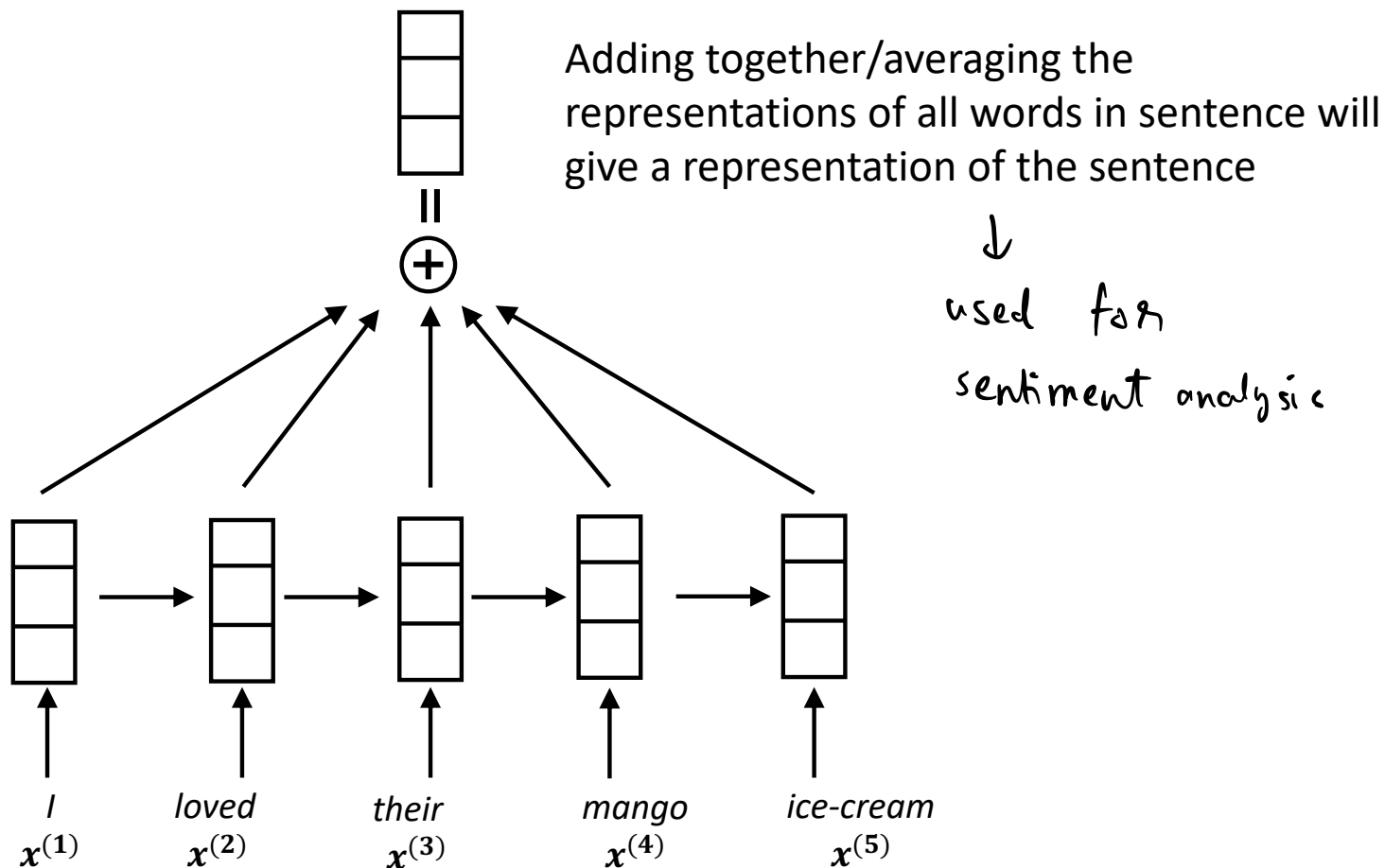
Lukasz Kaiser*
Google Brain
lukaszkaizer@google.com

Illia Polosukhin* †
illia.polosukhin@gmail.com

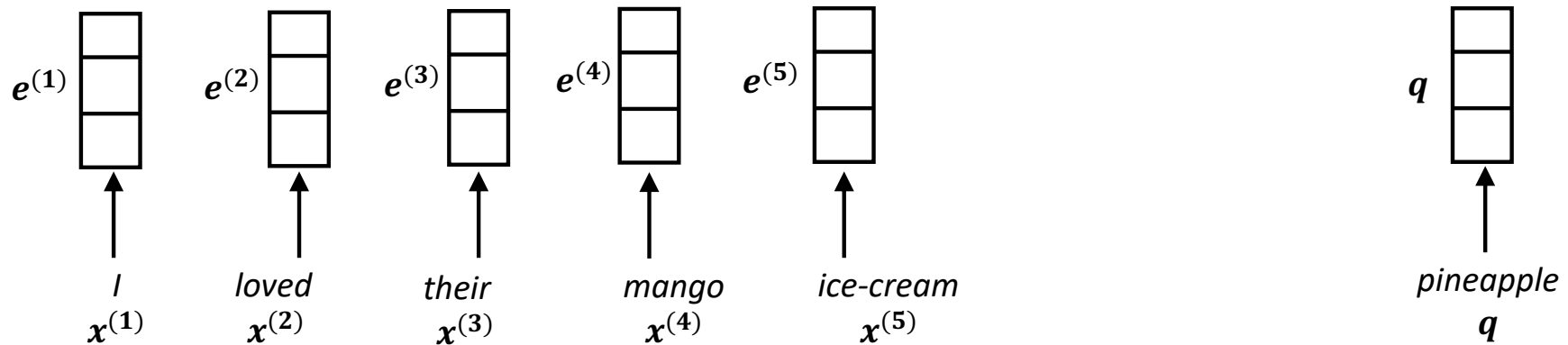
Abstract

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.

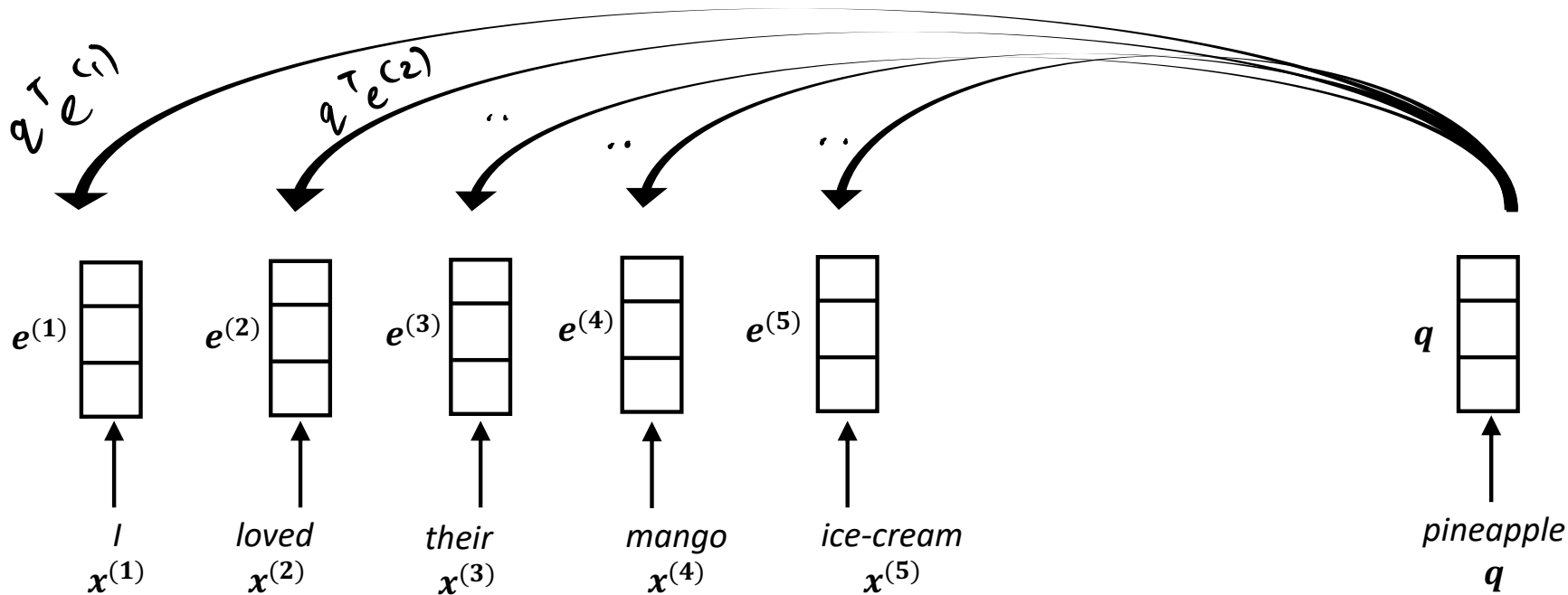
Starting point: Averaging word representations



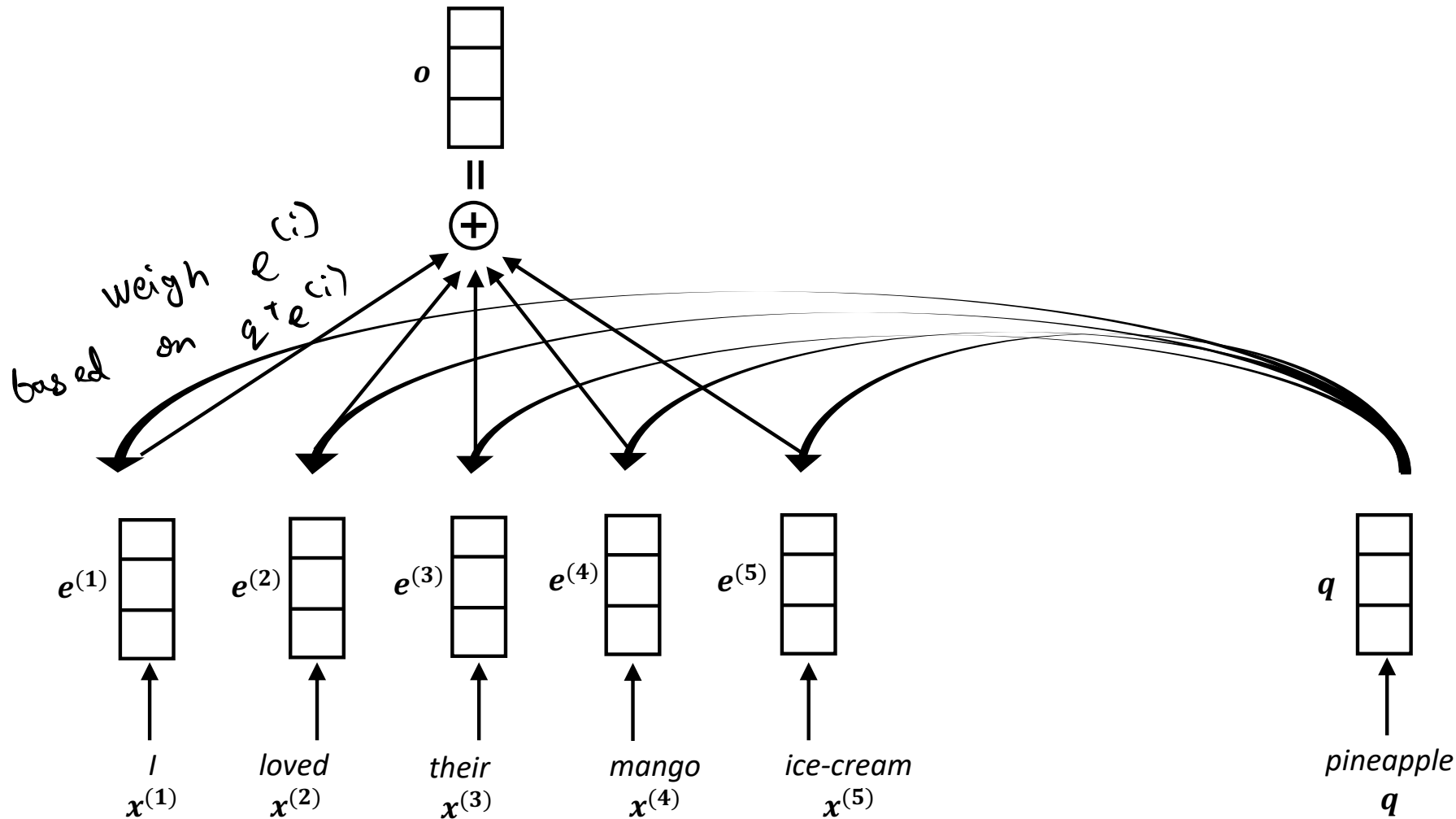
Attention: Weighted averaging



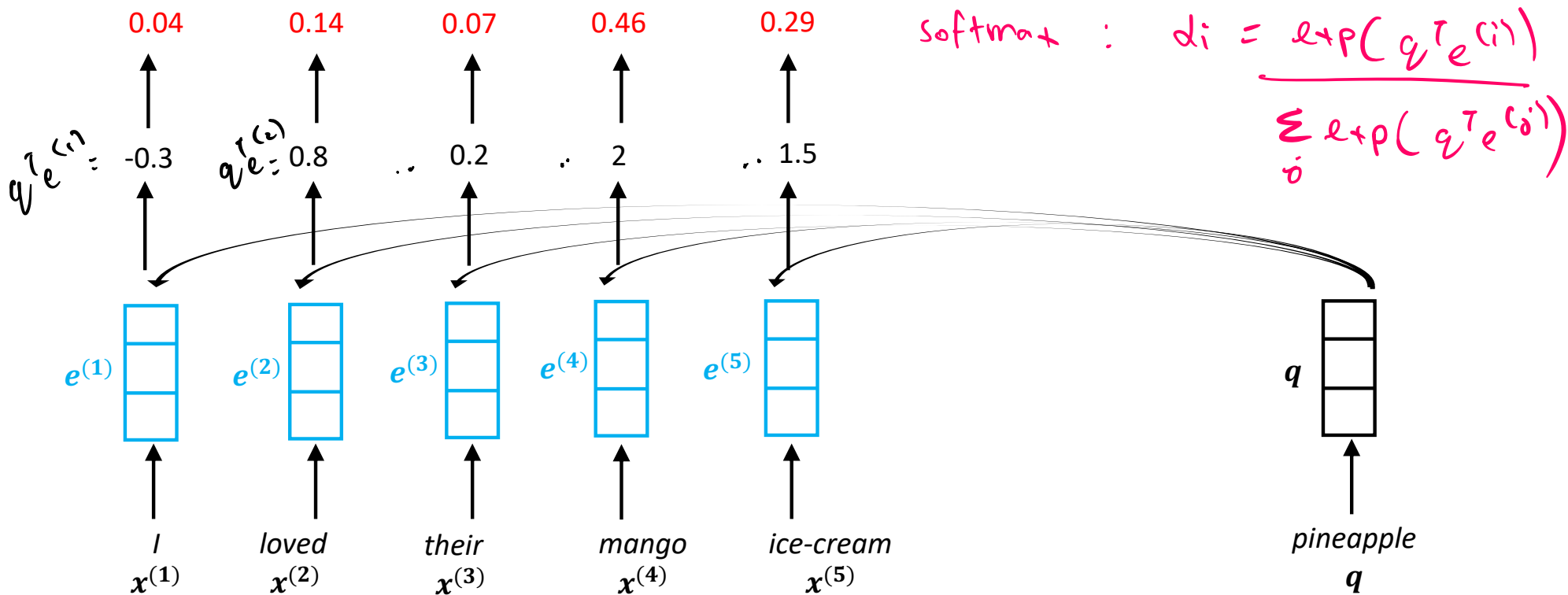
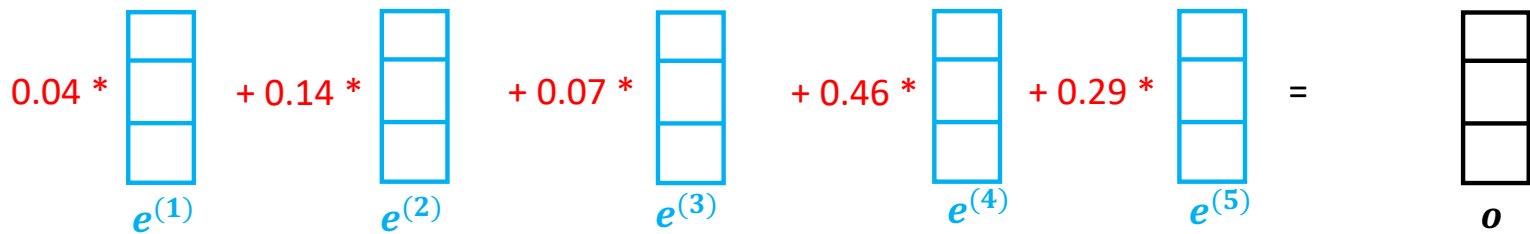
Attention: Weighted averaging



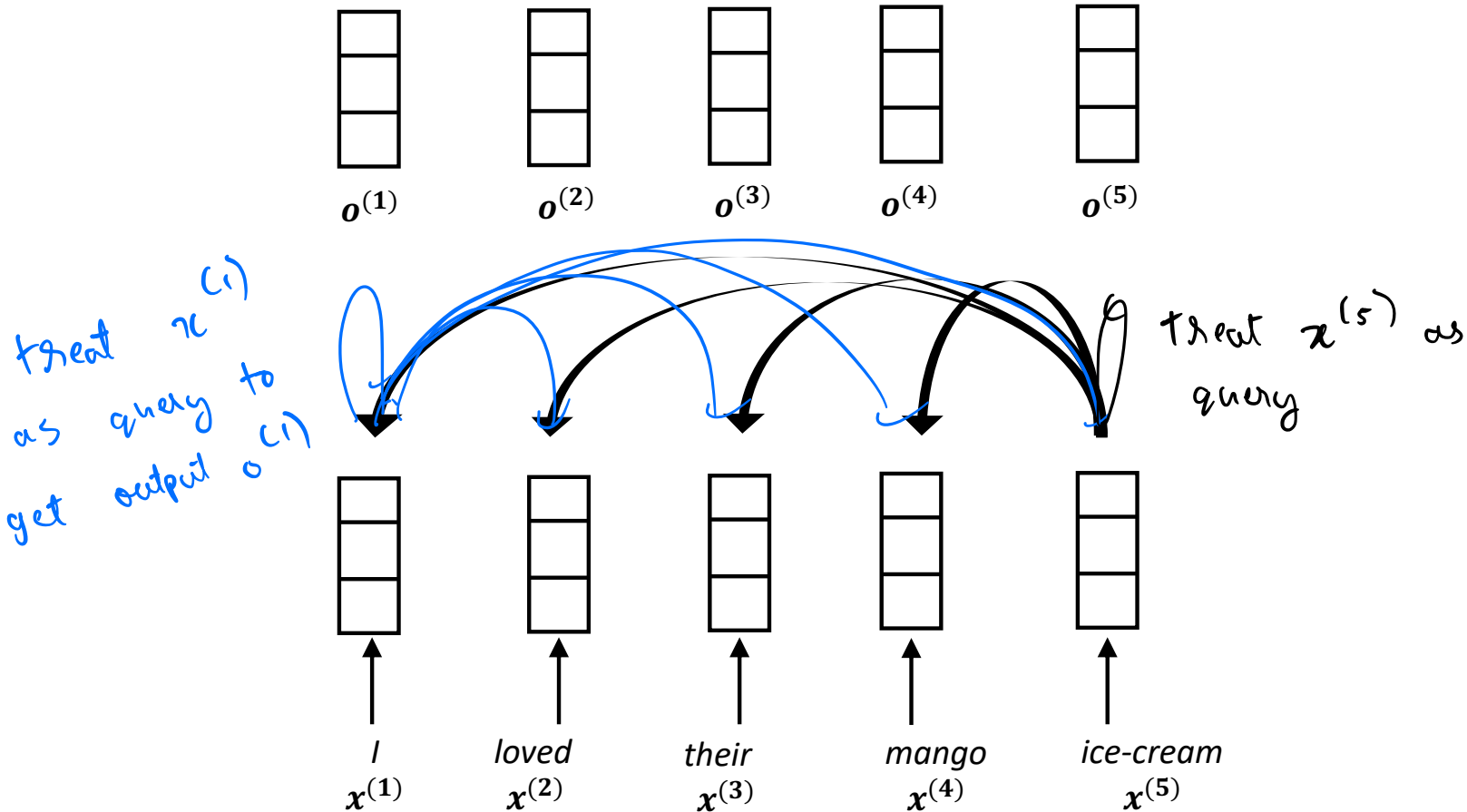
Attention: Weighted averaging



Attention: Weighted averaging

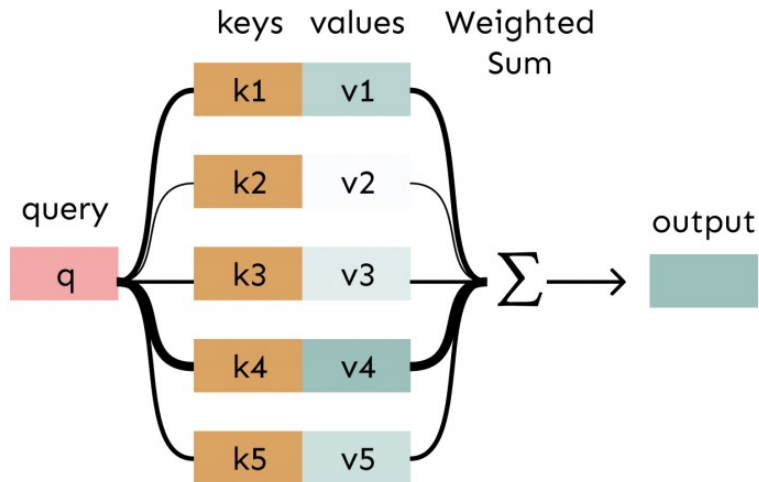


Self-attention

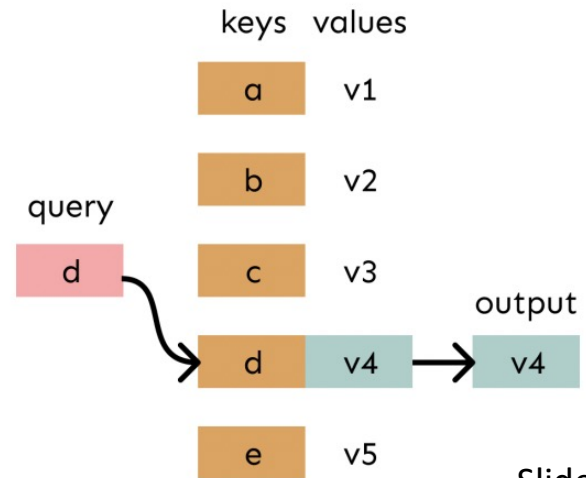


Attention as soft lookup

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



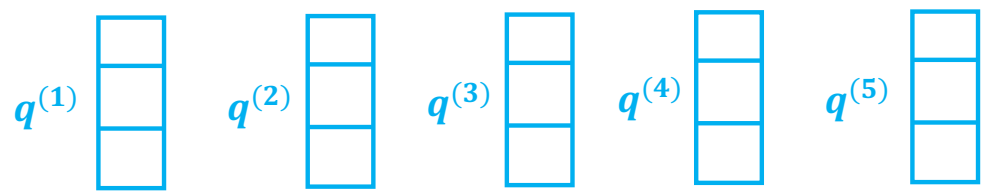
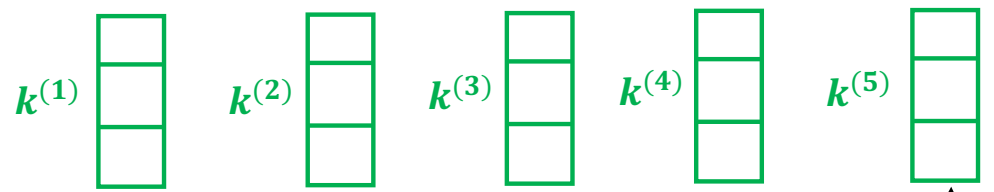
Lookup: find query in database, return value corresponding to its key



Self-attention

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

$(q^{(5)})^T k^{(s)}$
0.12
0.28
0.09
0.23
0.28
 $(q^{(5)})^T k^{(s)}$



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

→ softmax

How to get k, q, v ?

$$q^{(s)} = Q x^{(s)}$$

$$k^{(s)} = K x^{(s)}$$

$$v^{(s)} = V x^{(s)}$$

$$Q, K, V \in \mathbb{R}^{d \times d}$$

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

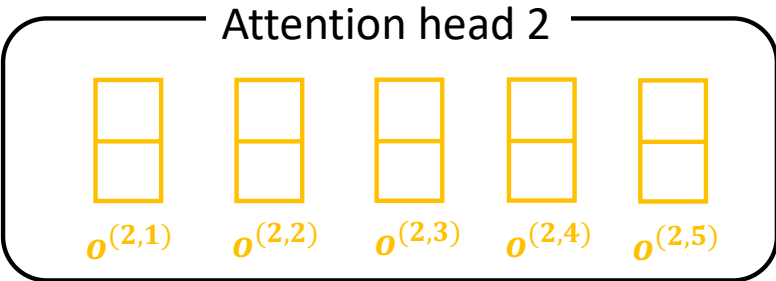
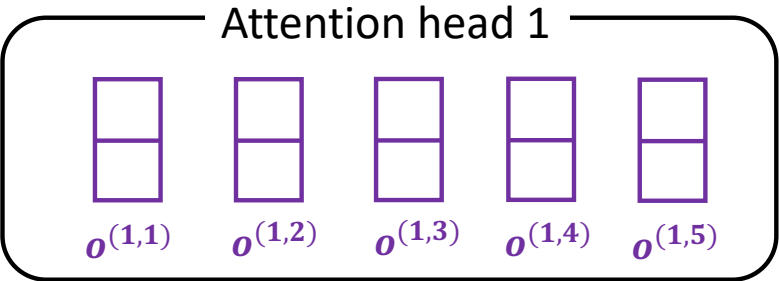
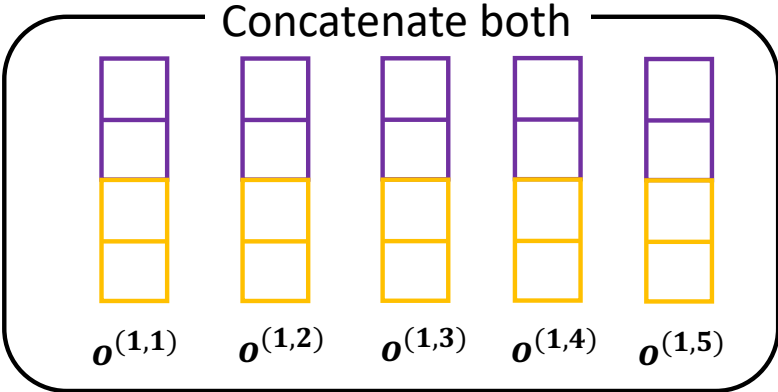
$$\alpha_{ij} = \mathbf{q}_i^\top \mathbf{k}_j$$

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

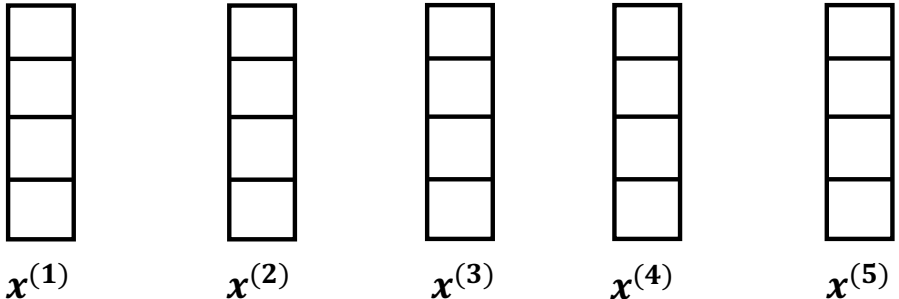
3. Compute output for each word as weighted sum of values

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

Multi headed self-attention



apply attention layer with $Q, K, V \in \mathbb{R}^{(d/2)+d}$



Multi headed self-attention

- Input: List of vectors $\mathbf{x}_1, \dots, \mathbf{x}_T$, each of size d
- Output: List of vectors $\mathbf{h}_1, \dots, \mathbf{h}_T$, each of size d
- Formula: For each head i :
 - Compute self attention output using $\mathbf{Q}_i, \mathbf{K}_i, \mathbf{V}_i$
 - Finally, concatenate results for all heads
- Parameters:
 - For each head i , parameter matrices $\mathbf{Q}_i, \mathbf{K}_i, \mathbf{V}_i$ of size $d_{\text{attn}} \times d$
 - # of heads n is hyperparameter, $d_{\text{attn}} = d/n$

What do attention heads learn?

