

344.063 KV Special Topic:

Natural Language Processing with Deep Learning

Language Modeling with Recurrent Neural Networks



Navid Rekab-saz

navid.rekabsaz@jku.at

Agenda

- Neural n -gram Language Model
- Recurrent Neural Networks
- Language Modeling with RNN

Probability

- Conditional probability, given two random variables X and Y :

$$P(Y|X)$$

- Probability distribution
 - For a **discrete** random variable Y with K states (classes)
 - $0 \leq P(Y_i) \leq 1$
 - $\sum_{i=1}^K P(Y_i) = 1$
 - E.g. with $K = 4$ states: [0.2 0.3 0.45 0.05]

Agenda

- **Neural n -gram Language Model**
- Recurrent Neural Networks
- Language Modeling with RNN

Language Modeling

- Language Modeling is the task of predicting a word/subword/character/etc. given a context:

$$P(v|\text{context})$$

- A Language Model can answer the questions like



$$P(v|the\ students\ opened\ their)$$

Language Modeling – formal definition

- Given a sequence of words $x^{(1)}, x^{(2)}, \dots, x^{(t)}$, a **language model** calculates the **probability distribution** of next word $x^{(t+1)}$ over all words in vocabulary

$$P(x^{(t+1)} | x^{(t)}, x^{(t-1)}, \dots, x^{(1)})$$

x is any word in the vocabulary $\mathbb{V} = \{v_1, v_2, \dots, v_N\}$

☞ Note: this is the definition of directed left-to-right language models.

N-gram

- Recall: a *n*-gram is a chunk of *n* consecutive words.

the students opened their _____

- unigrams: “*the*”, “*students*”, “*opened*”, “*their*”
- bigrams: “*the students*”, “*students opened*”, “*opened their*”
- trigrams: “*the students opened*”, “*students opened their*”
- 4-grams: “*the students opened their*”

N-gram Language Model



- **Markov assumption:** decision at time t depends only on the current state
- In n -gram Language Model: predicting $x^{(t+1)}$ depends on preceding $n-1$ words
- Without Markovian assumption:

$$P(x^{(t+1)} | x^{(t)}, x^{(t-1)}, \dots, x^{(1)})$$

- n -gram Language Model:

$$P(x^{(t+1)} | \underbrace{x^{(t)}, x^{(t-1)}, \dots, x^{(t-n+2)}}_{n-1 \text{ words}})$$

N-gram language modeling with neural networks

Recall

- The aim of a n -gram Language Model is to calculate:

$$P(x^{(t+1)} | x^{(t)}, \dots, x^{(t-n+2)})$$

- We can use a feed forward neural network to estimate this probability
- Immediate benefits:
 - Smooth probability estimation
 - Exploiting the semantic space of word embeddings (probably better generalization)

Neural n -gram LM – preparing training data

- Preparing training data for a neural 4-gram Language Model in the form of (context, next word), namely $(x^{(t-2)}x^{(t-1)}x^{(t)}, x^{(t+1)})$:

- For a given text corpus:

a fluffy cat sunbathes on the bank of river ...

- Training data items would be:

(`<bos> <bos> <bos>`, a)

(`<bos> <bos> a`, fluffy)

(`<bos> a fluffy`, cat)

(a fluffy cat, sunbathes)

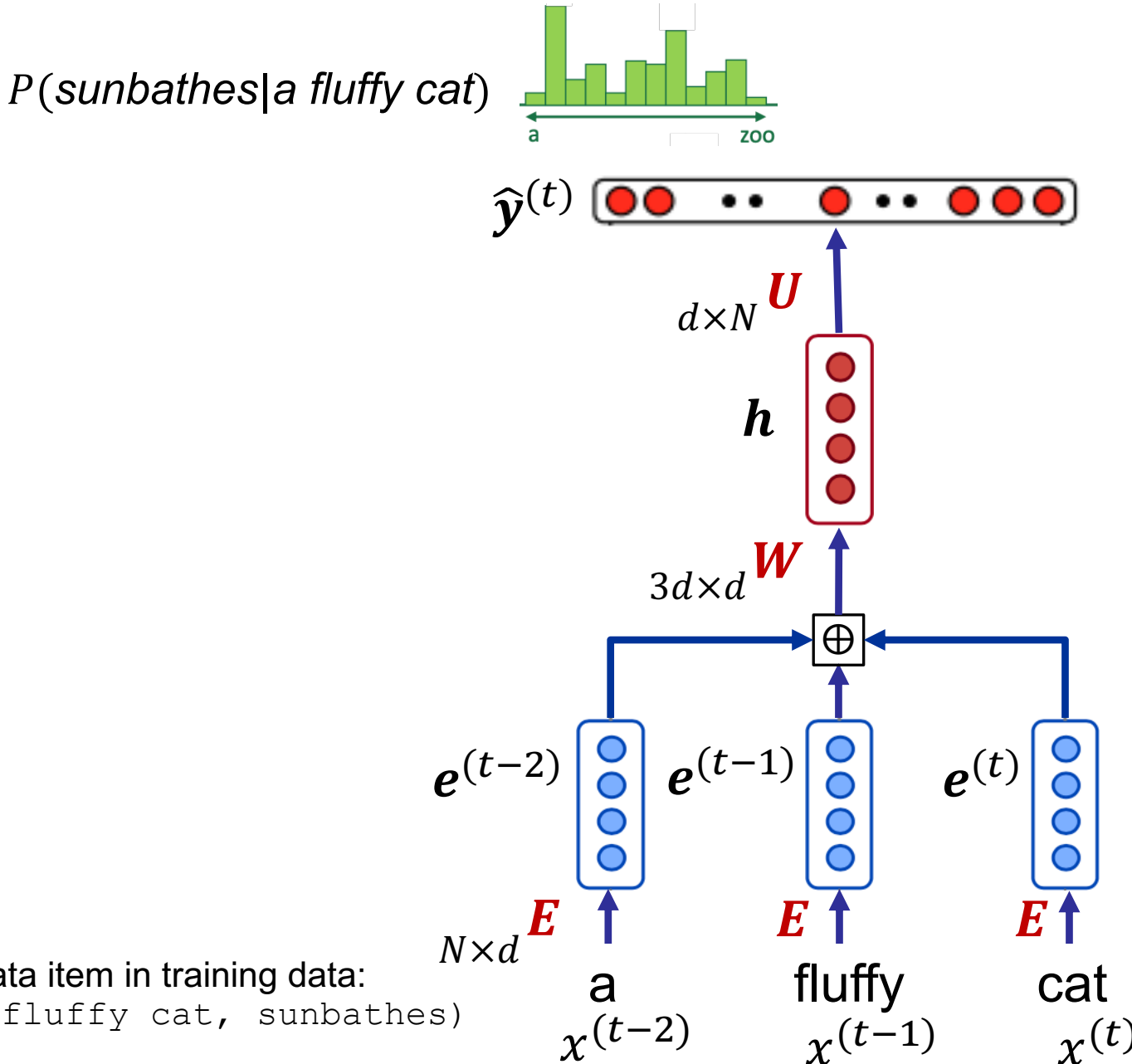
(fluffy cat sunbathes, on)

(cat sunbathes on, the)

(sunbathes on the, bank)

...

Neural n -gram Language Model – architecture



A data item in training data:
(a fluffy cat, sunbathes)

Formulation

Encoder

- From words to word embeddings:
 - One-hot vector of word $x^{(t)} \rightarrow \mathbf{x}^{(t)} \in \mathbb{R}^N$
 - Fetching word embedding $\rightarrow \mathbf{e}^{(t)} = \mathbf{x}^{(t)} \mathbf{E}$
 - In practice, $\mathbf{e}^{(t)}$ is achieved by fetching the vector of $x^{(t)}$ from \mathbf{E} (no need for $\mathbf{x}^{(t)}$)
- Concatenation of word embeddings: $\mathbf{e} = [\mathbf{e}^{(t-2)}, \mathbf{e}^{(t-1)}, \mathbf{e}^{(t)}]$
- Hidden layer: $\mathbf{h} = \tanh(\mathbf{W}\mathbf{e} + \mathbf{b})$

Decoder

- Predicted probabilities:
 - Predicted probability distribution:
$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{U}\mathbf{h} + \mathbf{b}) \in \mathbb{R}^N$$
 - Probability of any next word v at step t :

$$P(v | x^{(t)}, \dots, x^{(t-n+2)}) = \hat{y}_v^{(t)}$$

Model parameters

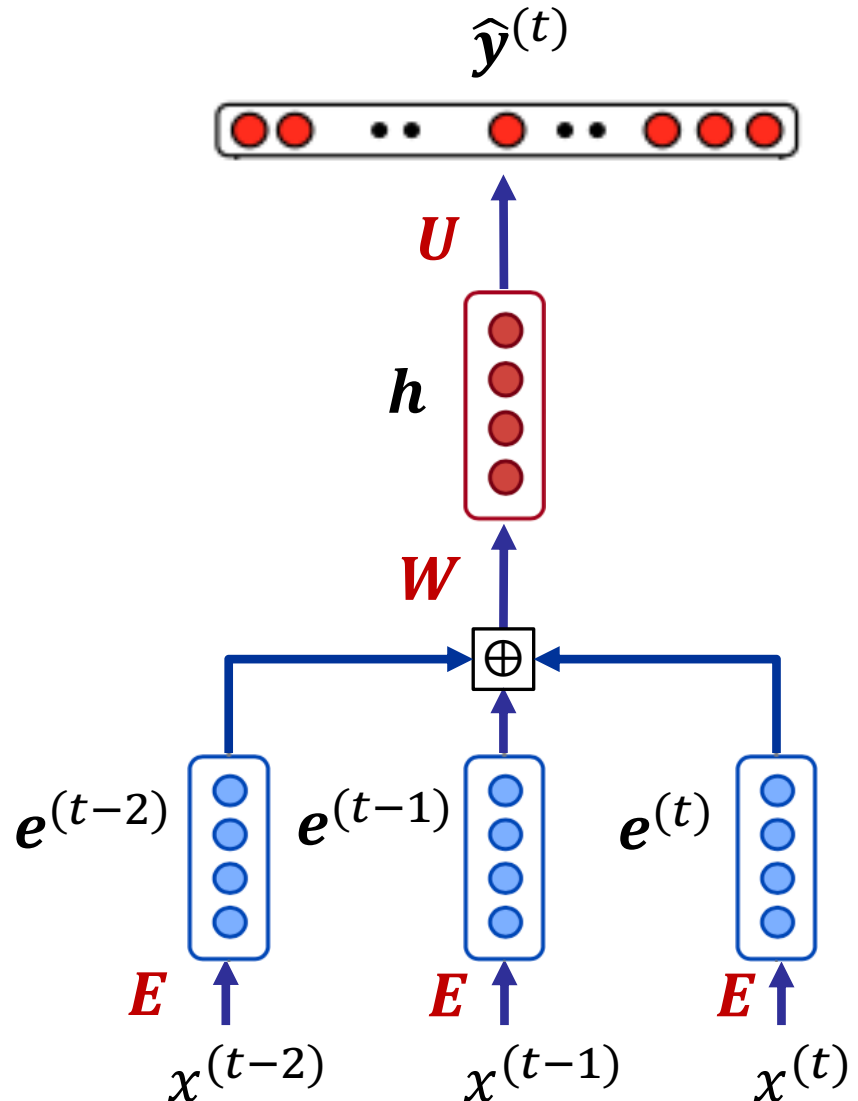
- $E \rightarrow N \times h$
- $W \rightarrow (n \times h) \times h$
- $U \rightarrow h \times N$

h embeddings dimension

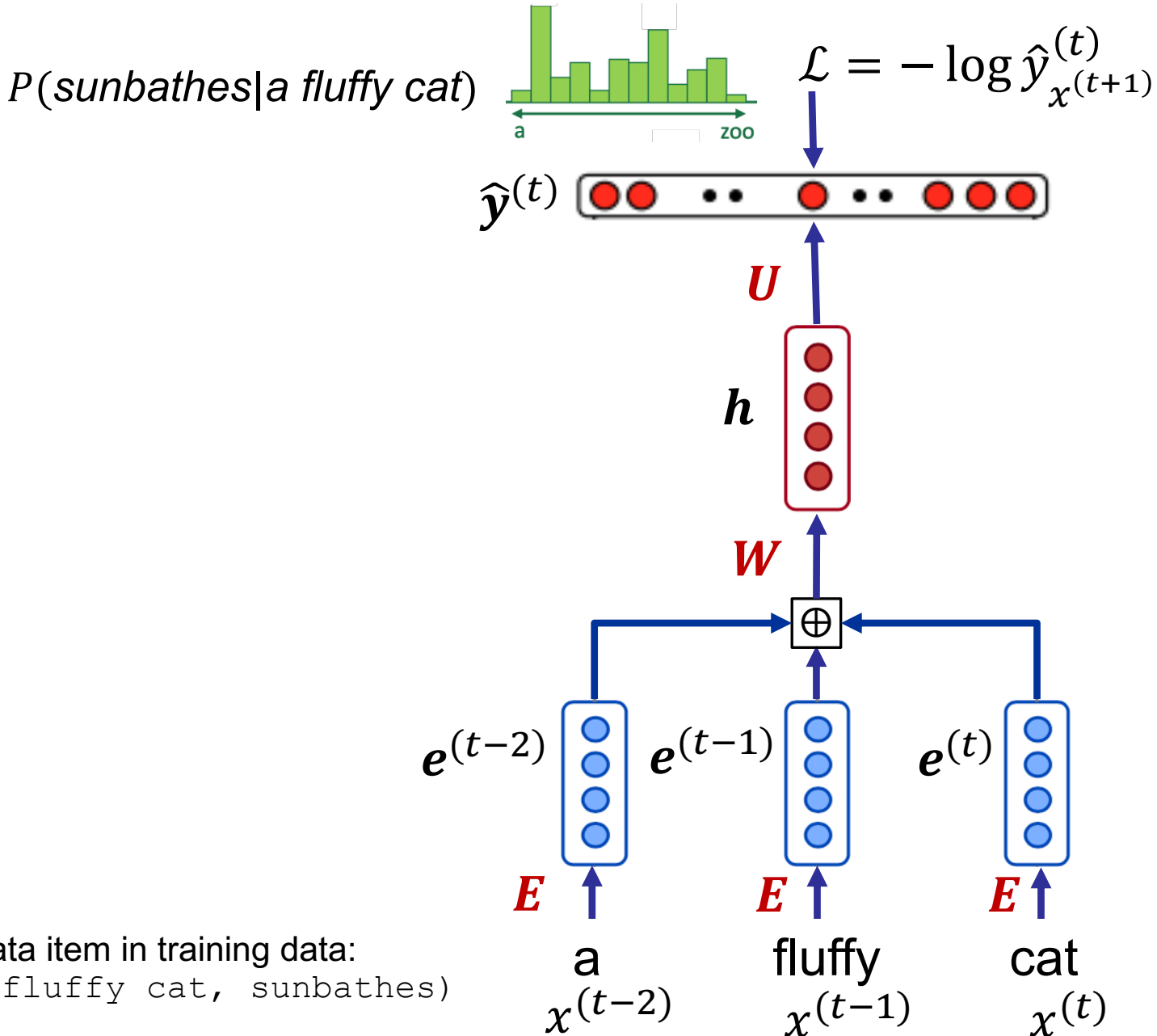
n number of preceding words

E is called **encoder embedding**

U is called **decoder embedding** or **output projection**



Loss function



A data item in training data:
(a fluffy cat, sunbathes)

Training procedure

- Start with a large text corpus: $x^{(1)}, \dots, x^{(T)}$
- For every **step** t predict the **output distribution** $\hat{y}^{(t)}$ given n previous words
- Loss function at t is Negative Log Likelihood of the **predicted probability** of the **word at** $x^{(t+1)}$

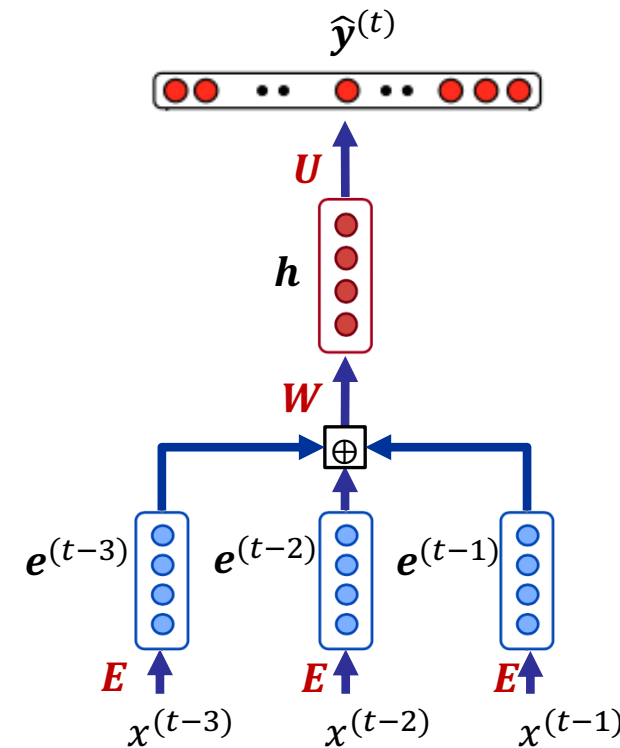
$$\mathcal{L}^{(t)} = -\log \hat{y}_{x^{(t+1)}}^{(t)}$$

- **Overall loss** is the average over all time steps:

$$\mathcal{L} = \frac{1}{T} \sum_{t=1}^T \mathcal{L}^{(t)}$$

Neural n -gram LMs – summary

- Neural n -gram LMs predict co-occurrence probabilities
- Neural n -gram LMs provide a smooth probability distribution
- At inference time, neural n -gram LMs require a forward pass



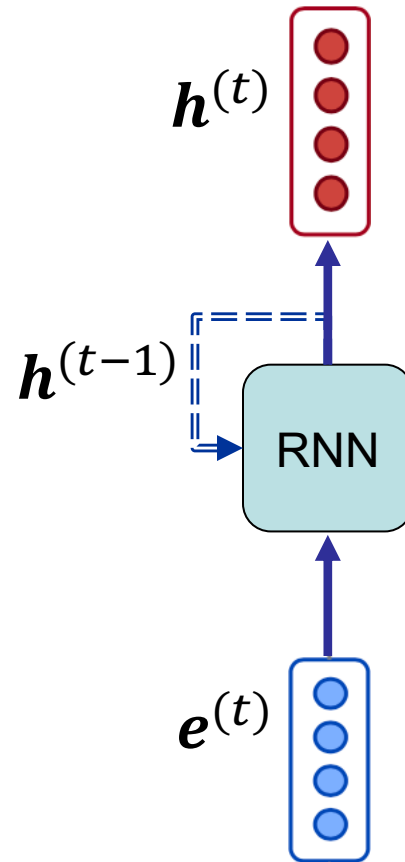
Agenda

- Neural n -gram Language Model
- **Recurrent Neural Networks**
- Language Modeling with RNN

Recurrent Neural Network

- Recurrent Neural Network (RNN) encodes/embeds a **sequential input of any size** into compositional embeddings
- A sequence can be ...
 - a stream of word/subword/character vectors
 - time series
 - etc.
- RNN models ...
 - capture dependencies through the sequence
 - apply the **same parameters** repeatedly
 - output a **final embedding** but also **intermediary embeddings** on each time step

Recurrent Neural Networks

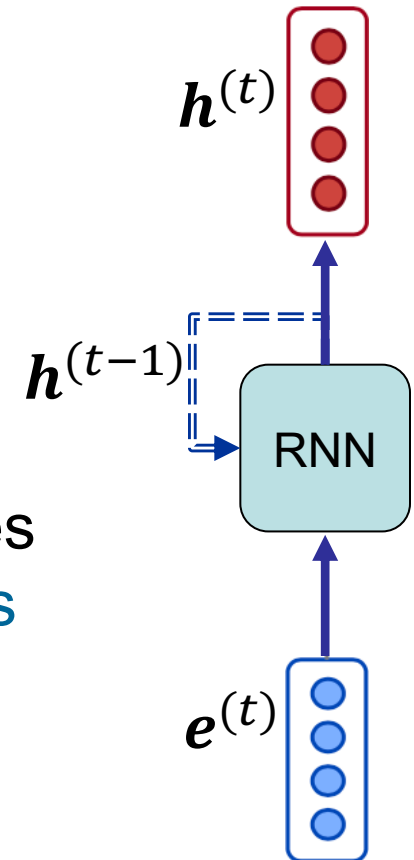


Recurrent Neural Networks

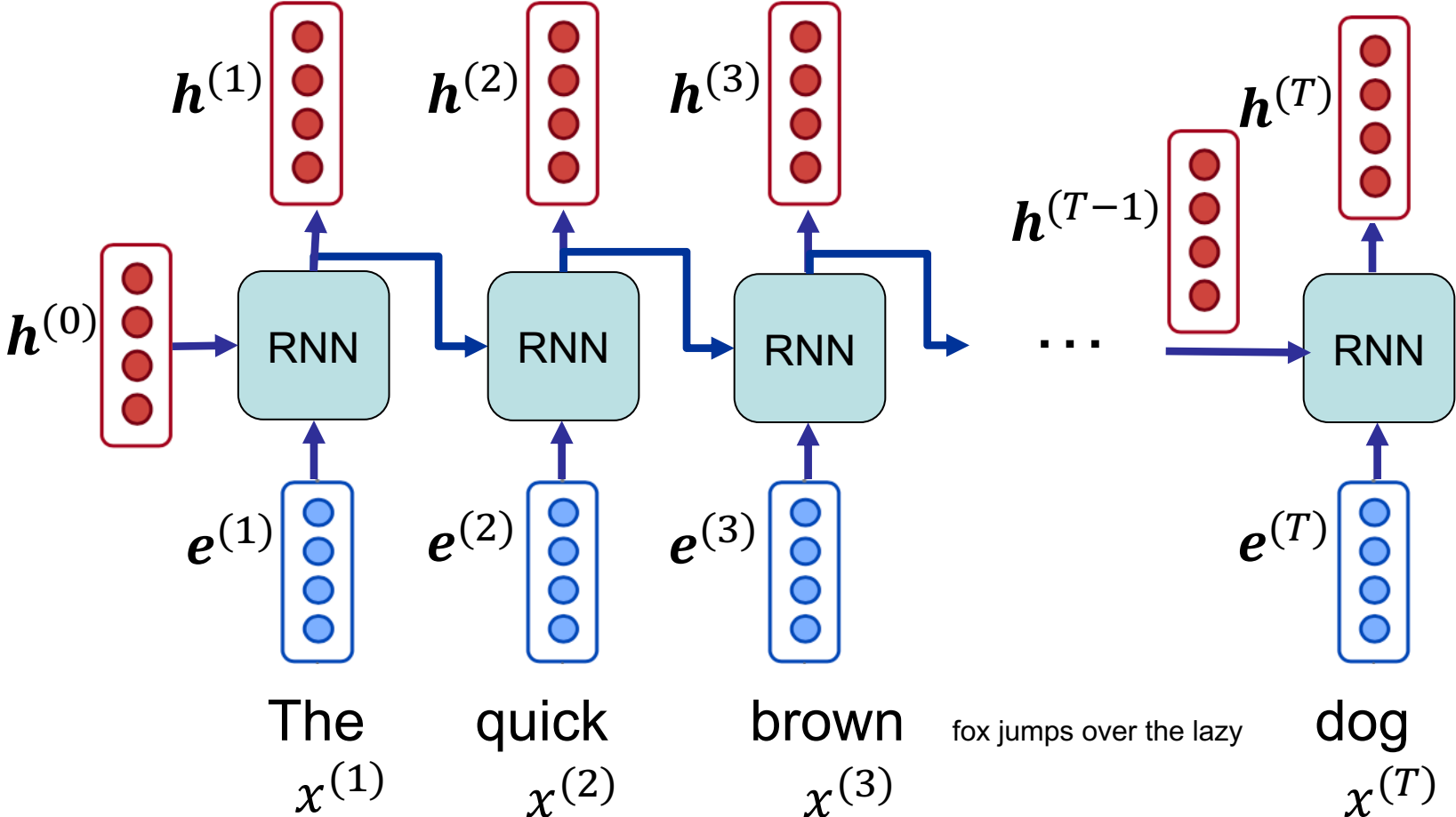
- Output $\mathbf{h}^{(t)}$ is a function of input $\mathbf{e}^{(t)}$ and the output of the previous time step $\mathbf{h}^{(t-1)}$

$$\mathbf{h}^{(t)} = \text{RNN}(\mathbf{h}^{(t-1)}, \mathbf{e}^{(t)})$$

- $\mathbf{h}^{(t)}$ is called **hidden state**
- With hidden state $\mathbf{h}^{(t-1)}$, the model accesses to a sort of **memory** from all **previous entities**



RNN – Unrolling



Standard (Elman) RNN

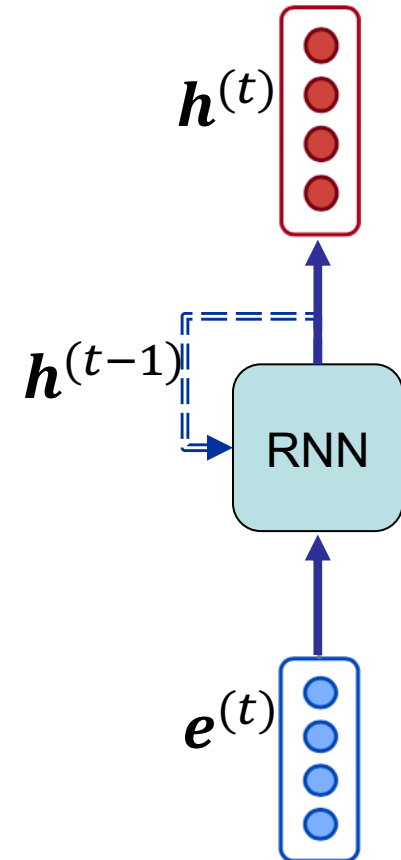
- General form of an RNN function

$$\mathbf{h}^{(t)} = \text{RNN}(\mathbf{h}^{(t-1)}, \mathbf{e}^{(t)})$$

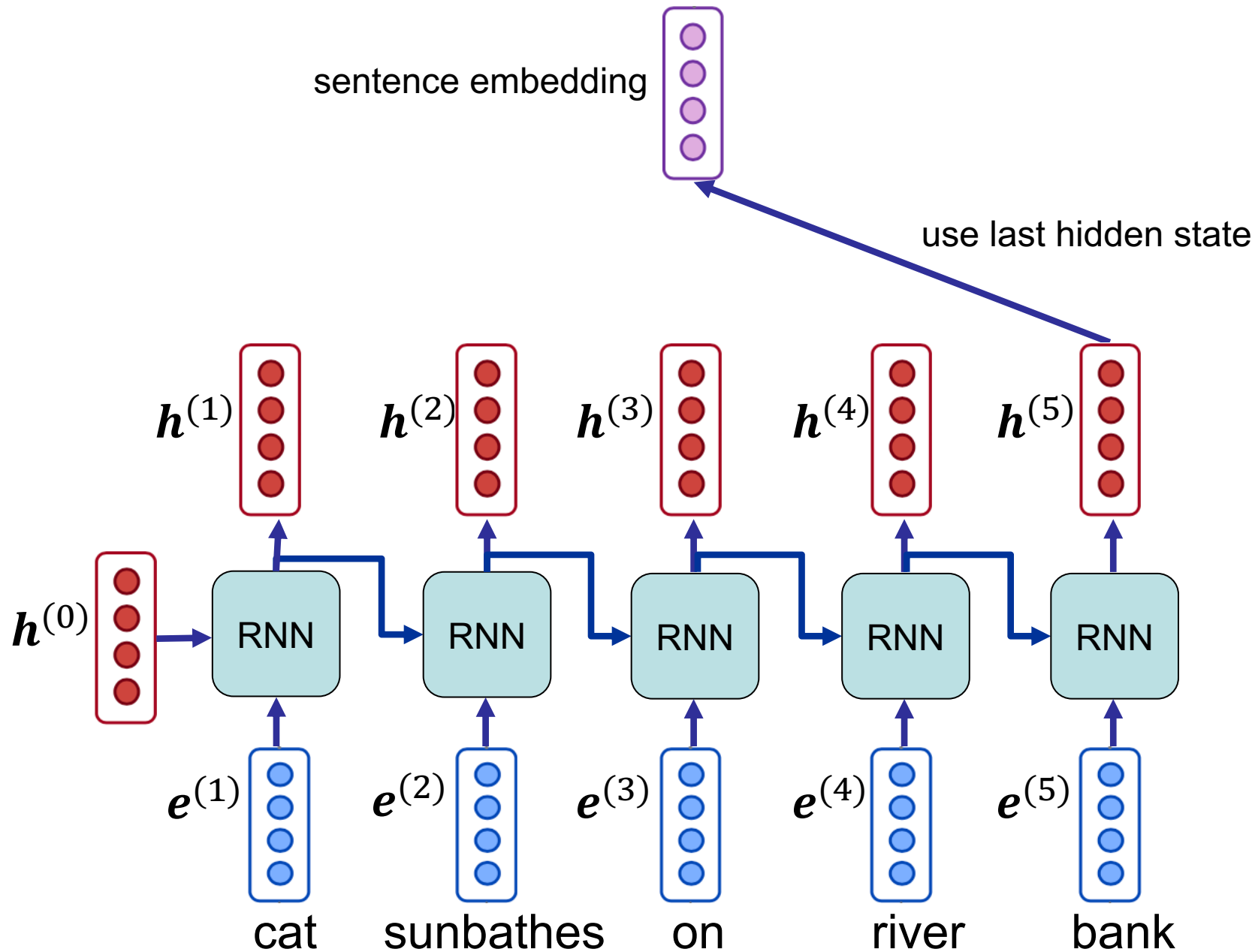
- Standard RNN:

- linear projection of the previous hidden state $\mathbf{h}^{(t-1)}$
- linear projection of input $\mathbf{e}^{(t)}$
- summing the projections and applying a non-linearity

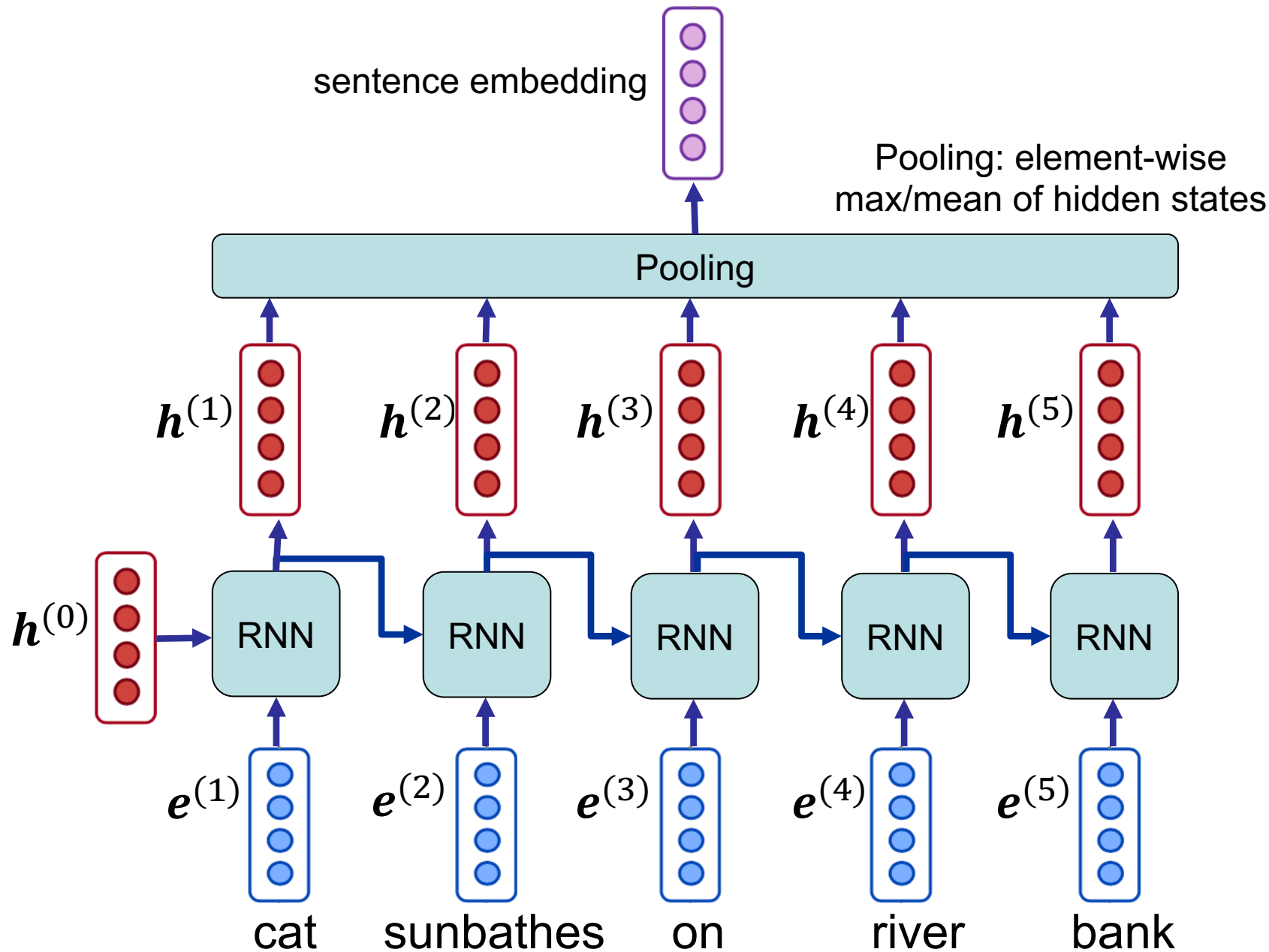
$$\mathbf{h}^{(t)} = \sigma(\mathbf{h}^{(t-1)} \mathbf{W}_h + \mathbf{e}^{(t)} \mathbf{W}_e + \mathbf{b})$$



RNN – Compositional embedding



RNN – Compositional embedding



Bidirectional RNNs

- Bidirectional RNN consists of **two RNNs**, one reads from the beginning to the end of sequence (**forward**), and the other reads from the end to the beginning (**backward**)

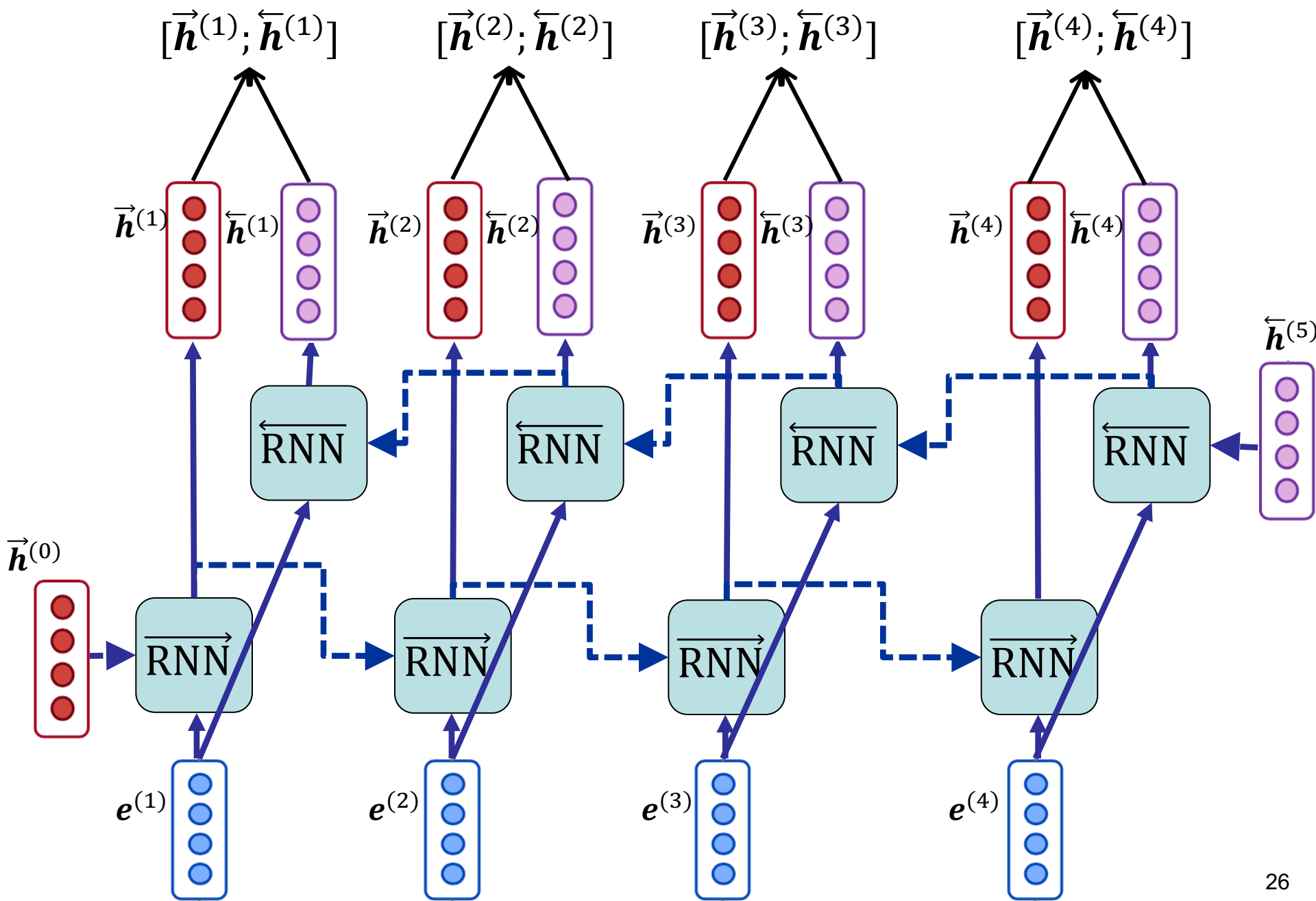
$$\vec{\mathbf{h}}^{(t)} = \overrightarrow{\text{RNN}} (\vec{\mathbf{h}}^{(t-1)}, \mathbf{e}^{(t)})$$

$$\overleftarrow{\mathbf{h}}^{(t)} = \overleftarrow{\text{RNN}} (\overleftarrow{\mathbf{h}}^{(t+1)}, \mathbf{e}^{(t)})$$

- Output at each time step is the **concatenation** of the outputs of both RNNs at that time step:

$$\mathbf{h}^{(t)} = [\vec{\mathbf{h}}^{(t)}; \overleftarrow{\mathbf{h}}^{(t)}]$$

- *To remember:* Using bidirectional RNN is only possible when the **entire sequence** is available



Agenda

- Neural n -gram Language Model
- Recurrent Neural Networks
- **Language Modeling with RNN**

Language Modeling – recap

- Given a sequence of words $x^{(1)}, x^{(2)}, \dots, x^{(t)}$, a **language model** calculates the **probability distribution** of next word $x^{(t+1)}$ over all words in vocabulary

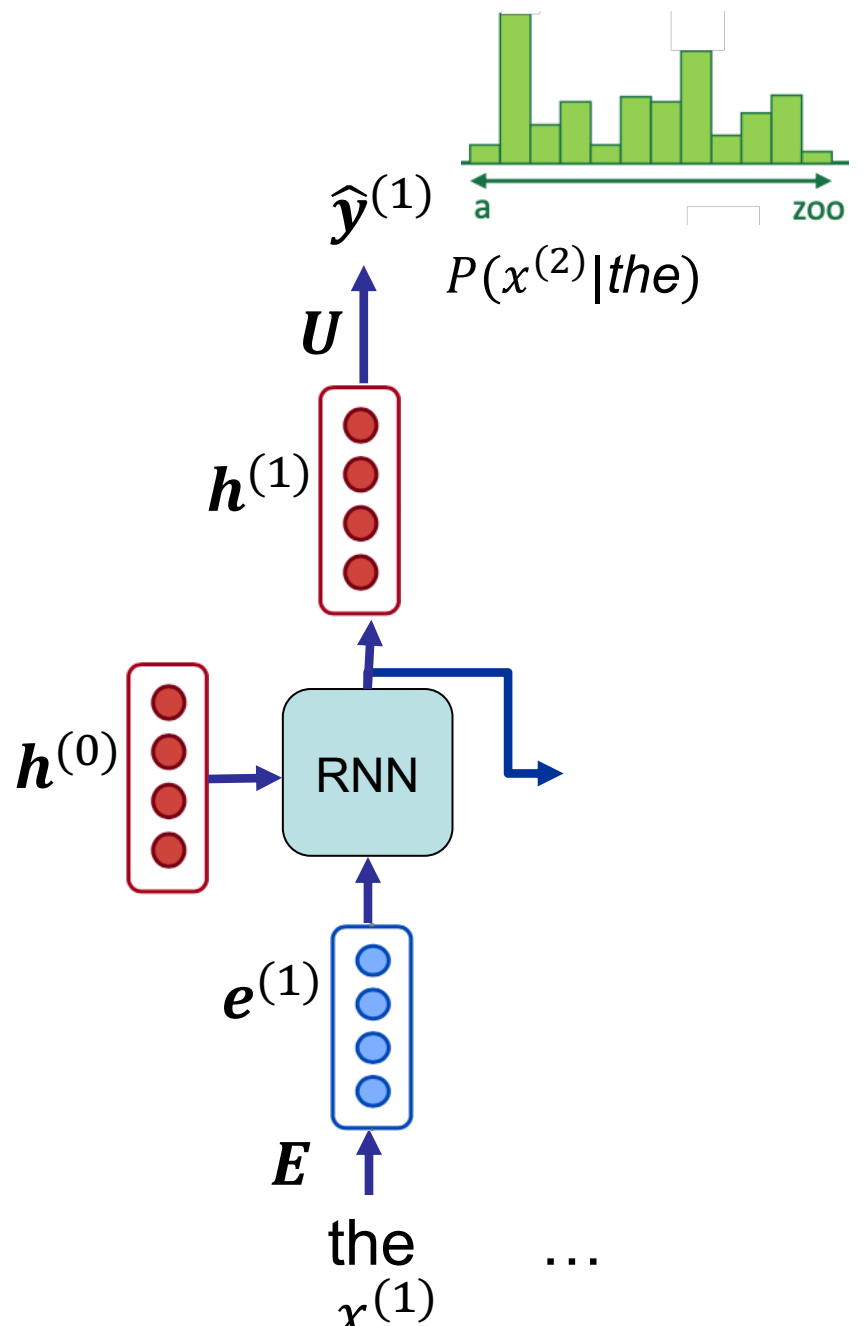
$$P(x^{(t+1)} | x^{(t)}, x^{(t-1)}, \dots, x^{(1)})$$

x is any word in the vocabulary $\mathbb{V} = \{v_1, v_2, \dots, v_N\}$

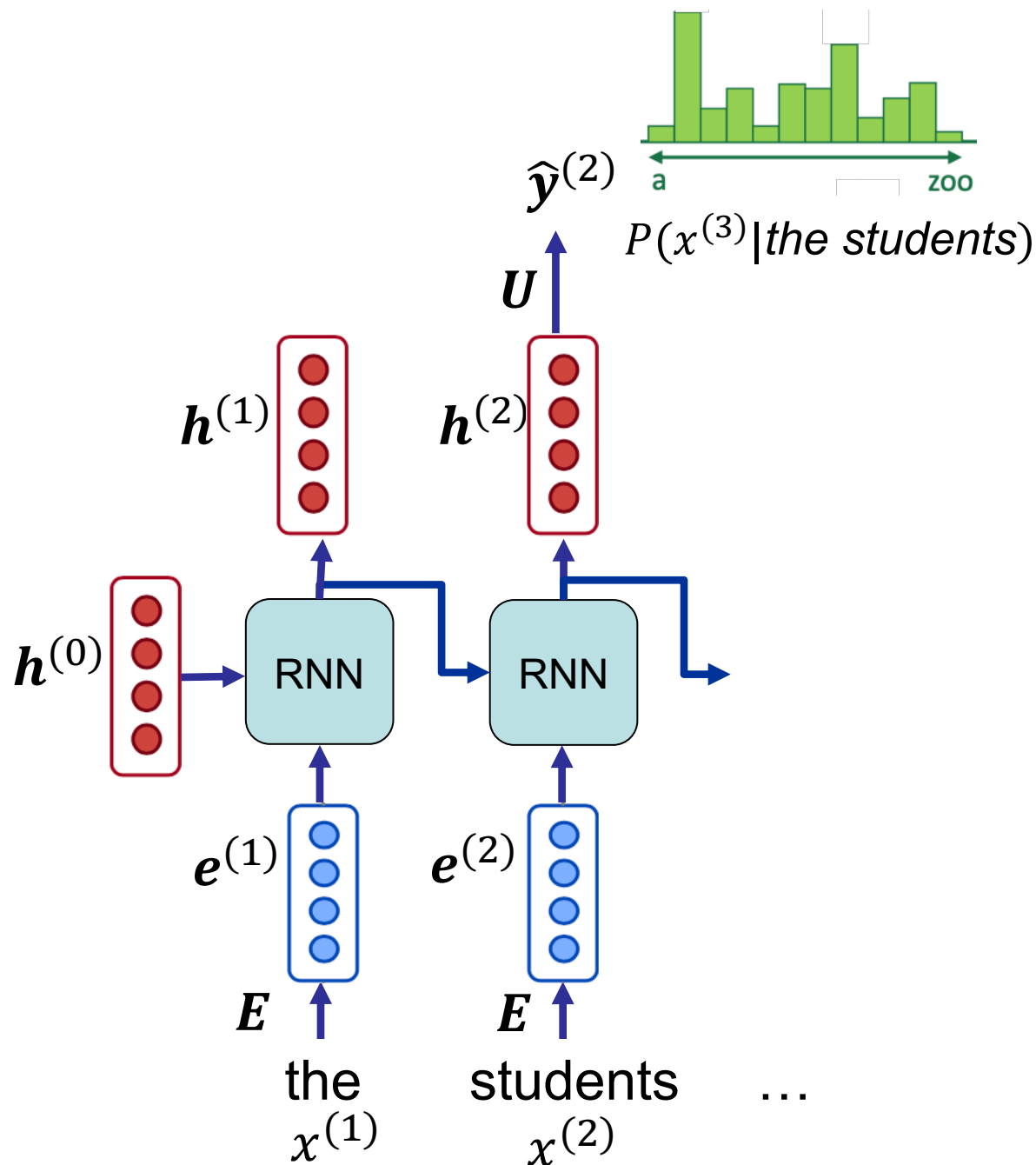


$$P(v | \text{the students opened their})$$

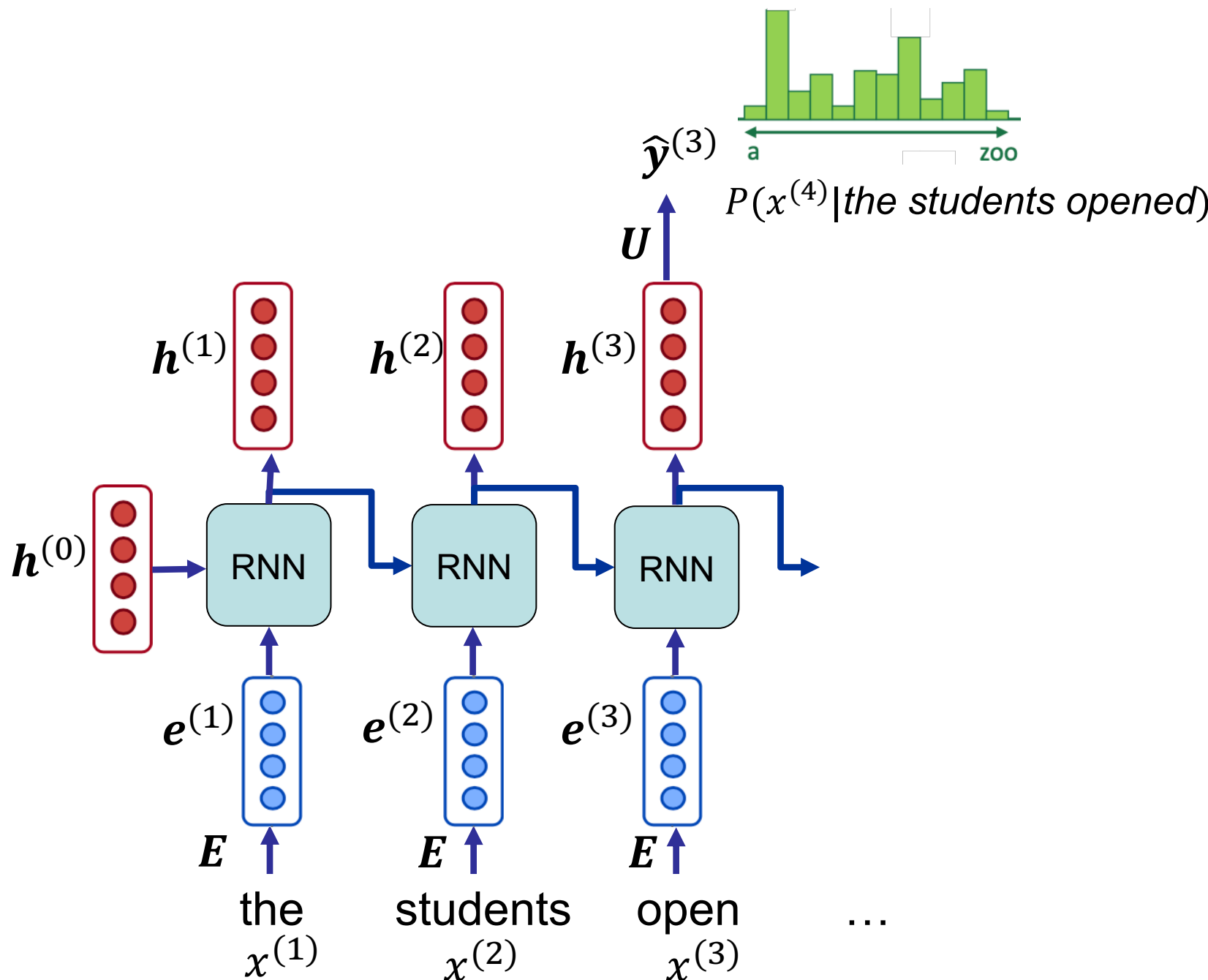
RNN Language Model



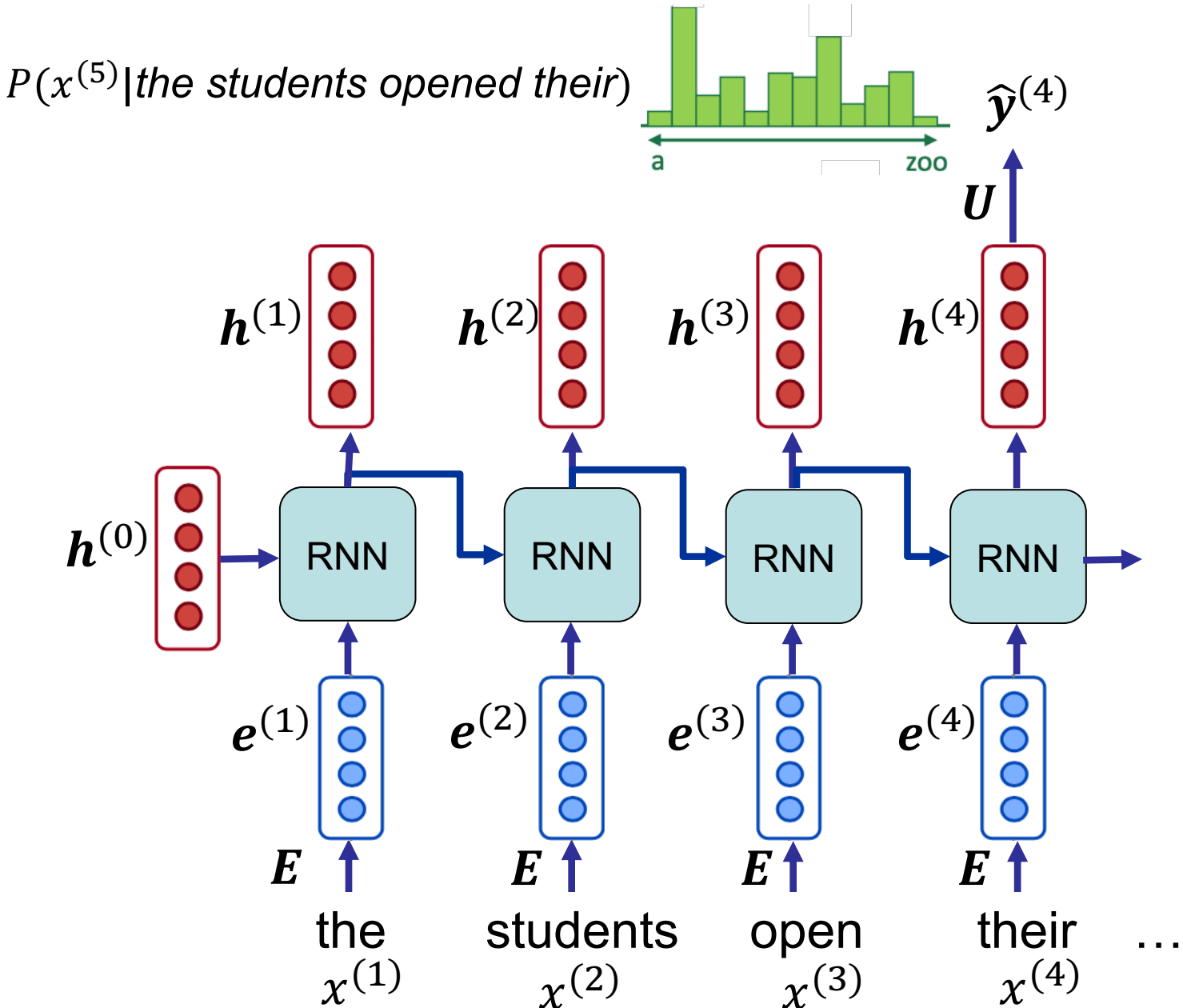
RNN Language Model



RNN Language Model



RNN Language Model



Formulation

Encoder

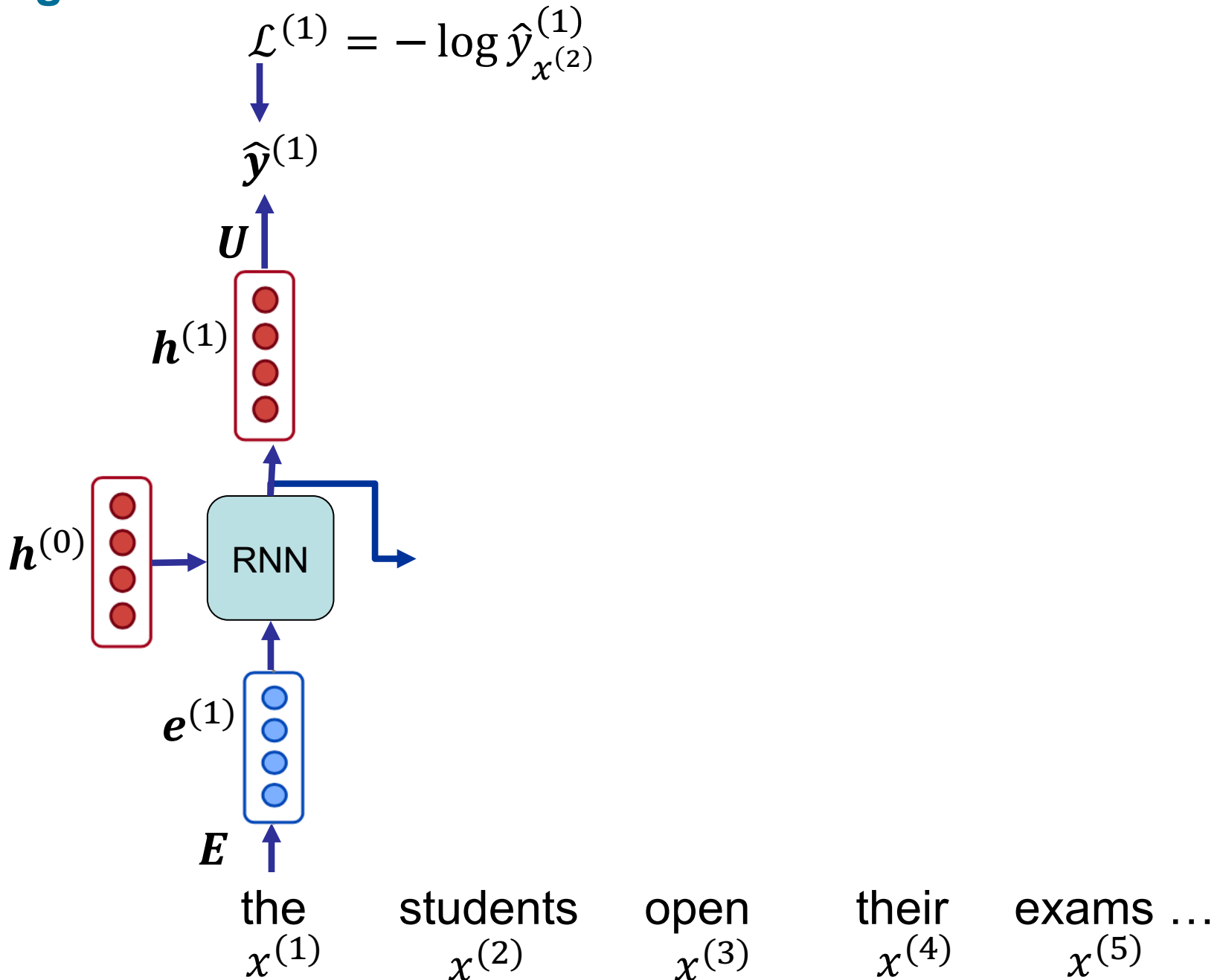
- From words to word embeddings:
 - One-hot vector of word $x^{(t)} \rightarrow \mathbf{x}^{(t)} \in \mathbb{R}^N$
 - Fetching word embedding $\rightarrow \mathbf{e}^{(t)} = \mathbf{x}^{(t)} \mathbf{E}$
 - In practice, $\mathbf{e}^{(t)}$ is achieved by fetching the vector of $x^{(t)}$ from \mathbf{E} (no need for $\mathbf{x}^{(t)}$)
- RNN: $\mathbf{h}^{(t)} = \text{RNN}(\mathbf{h}^{(t-1)}, \mathbf{e}^{(t)})$

Decoder

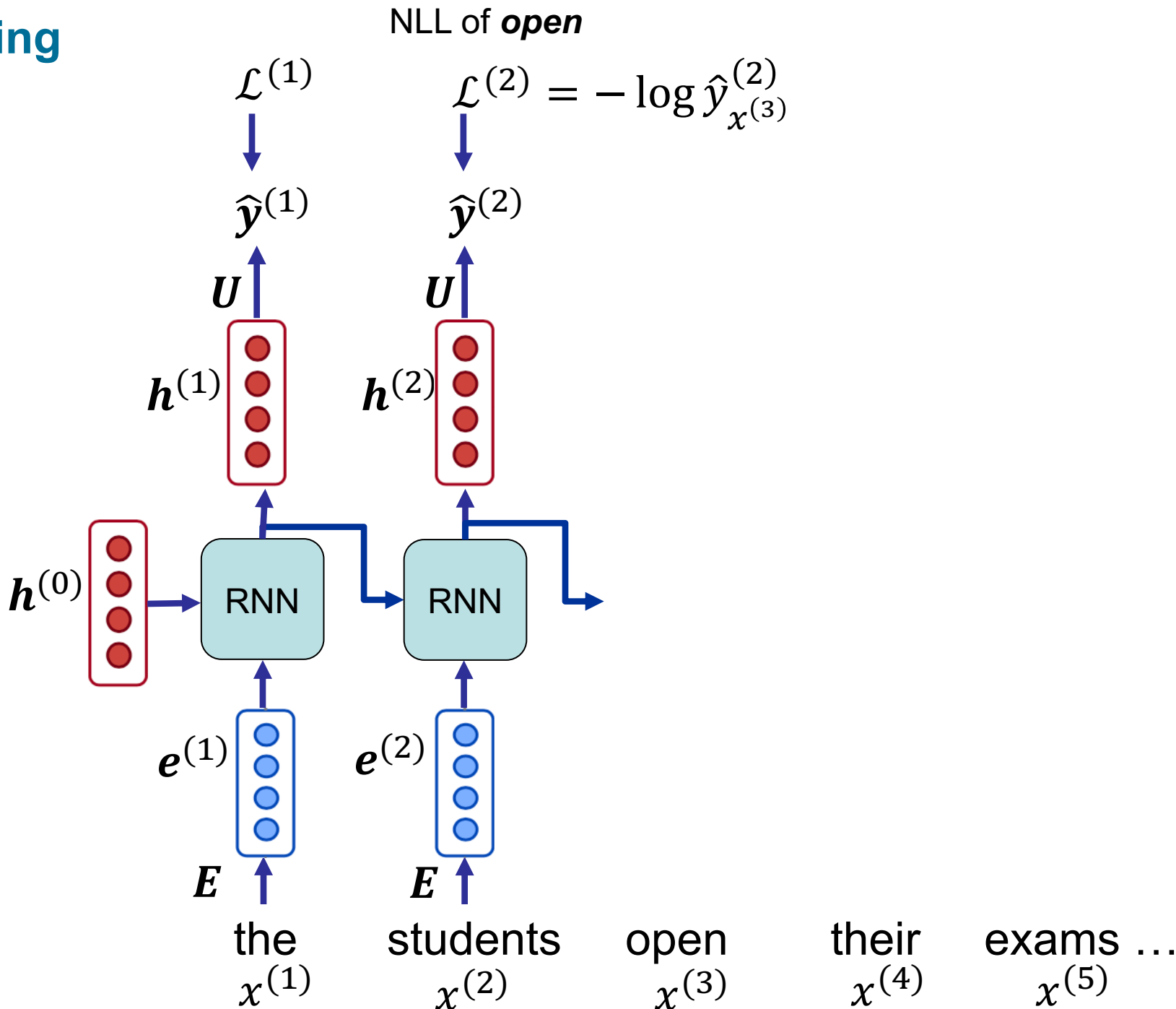
- Predicted probabilities
 - Predicted probability distribution:
$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{U}\mathbf{h}^{(t)} + \mathbf{b}) \in \mathbb{R}^N$$
 - Probability of any word v at step t :
$$P(v|x^{(t-1)}, \dots, x^{(1)}) = \hat{y}_v^{(t)}$$

Training

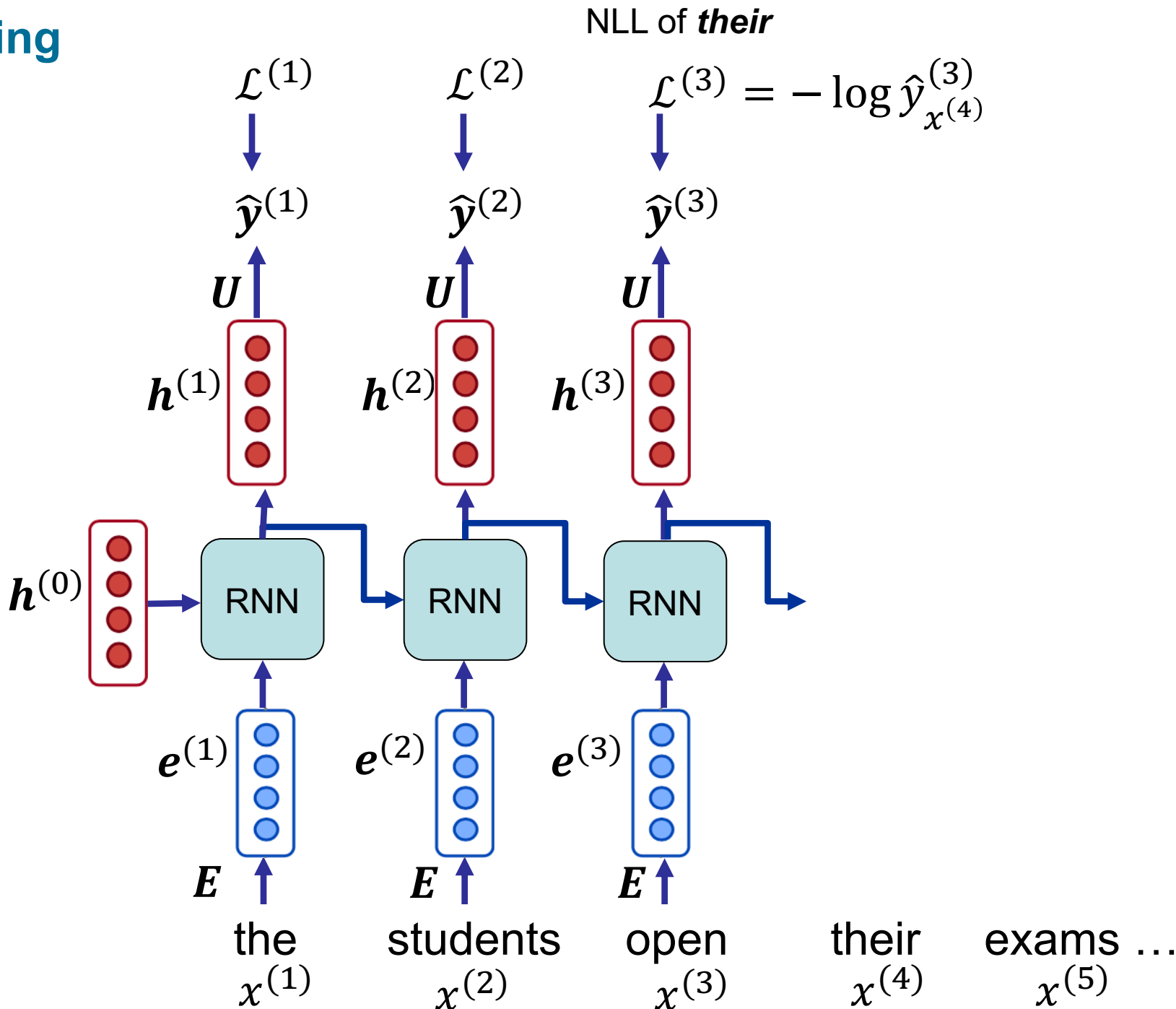
NLL of *students*



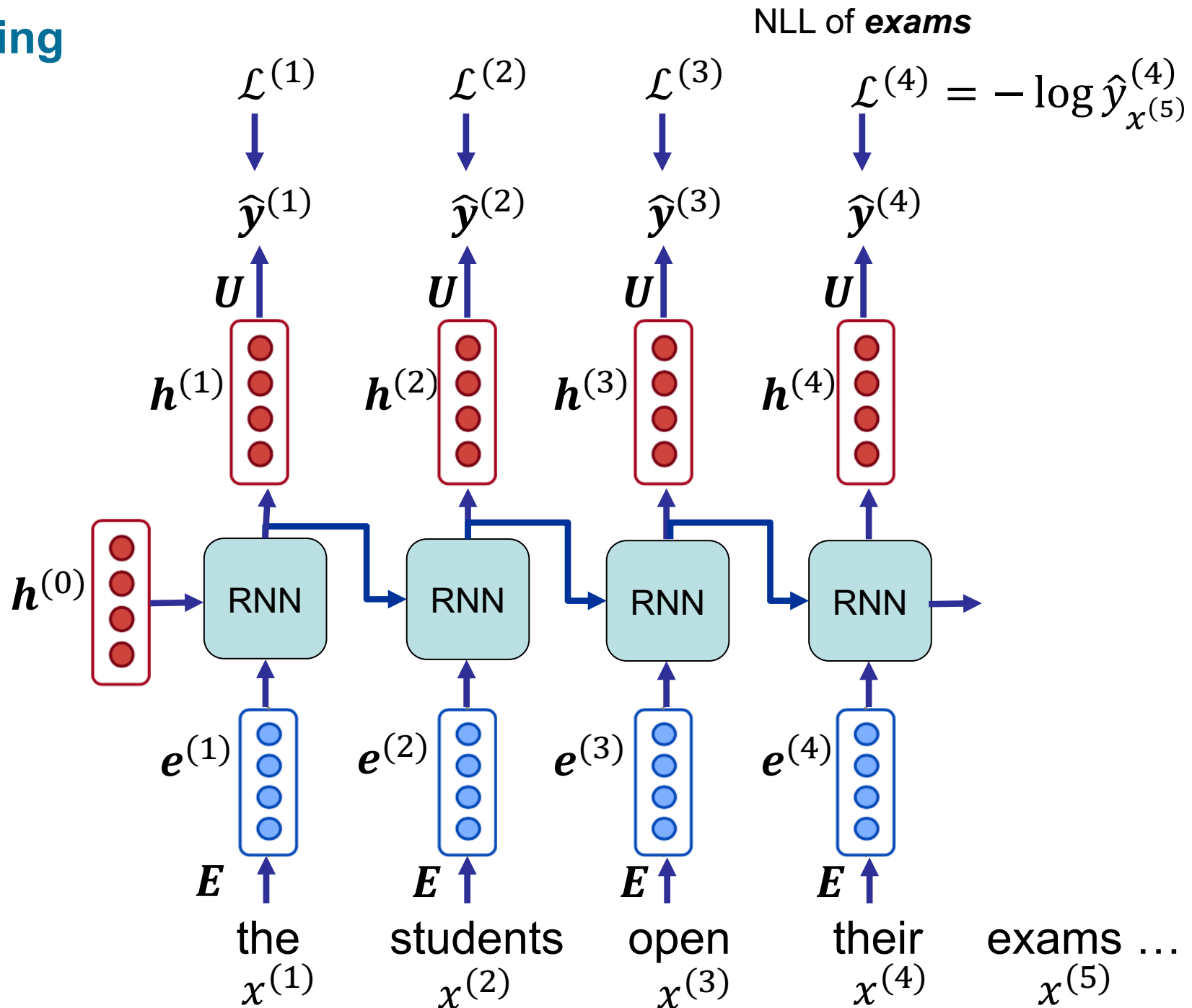
Training



Training



Training



Training an RNN Language Model

- Start with a large text corpus: $x^{(1)}, \dots, x^{(T)}$
- For every **step t** predict the **output distribution $\hat{y}^{(t)}$**
- Calculate the loss function: Negative Log Likelihood of the **predicted probability** of the true **next word $x^{(t+1)}$**

$$\mathcal{L}^{(t)} = -\log \hat{y}_{x^{(t+1)}}^{(t)} = -\log P(x^{(t+1)} | x^{(t)}, \dots, x^{(1)})$$

- **Overall loss** is the average of loss values over the entire training set:

$$\mathcal{L} = \frac{1}{T} \sum_{t=1}^T \mathcal{L}^{(t)}$$

Training RNN Language Model – Mini batches

- In practice, the overall loss is calculated not over whole the corpus, but over (mini) batches of length L :

$$\mathcal{L} = \frac{1}{L} \sum_{t=1}^L \mathcal{L}^{(t)}$$

- After calculating the \mathcal{L} for one batch, gradients are computed, and weights are updated (e.g. using SGD)

Training RNN Language Model – Data preparation

- In practice, every **forward pass** contains k batches. These batches are all trained **in parallel**
 - with **batch size** k , tensor of each forward pass has the shape: $[k, L]$
- To prepare the data in this form, the corpus is splitted into **sub-corpora**. Each sub-corpus contains the text for each row of forward tensors
- For example, for batch size $k = 2$, the corpus is splitted in middle, and the first forward-pass tensor looks like:

$$\text{batch size } k \rightarrow \begin{bmatrix} x^{(1)} & \dots & x^{(L)} \\ x^{(\tau+1)} & \dots & x^{(\tau+L)} \end{bmatrix}$$

$\tau = T/2$ is the overall length of the first sub-corpus

Training RNN Language Model – $h^{(0)}$

- At the beginning, the **initial hidden state** $h^{(0)}$ is set to vectors of **zeros** (no memory)
- To carry the memory of previous forward passes, at each forward pass, the initial hidden states are initialized with the values of the **last hidden states** of the **previous forward pass**

Example

- First forward pass: $h^{(0)}$ is set to zero values

$$\begin{bmatrix} x^{(1)} & \dots & x^{(L)} \\ x^{(\tau+1)} & \dots & x^{(\tau+L)} \end{bmatrix}$$

- Second forward pass: $h^{(0)}$ is set to the $h^{(L)}$ values in the first pass

$$\begin{bmatrix} x^{(L+1)} & \dots & x^{(2L)} \\ x^{(L+\tau+1)} & \dots & x^{(\tau+2L)} \end{bmatrix}$$

Training RNN Language Model – Parameters

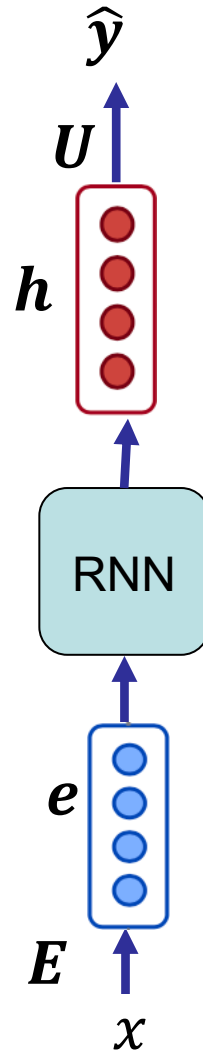
- Parameters in a vanilla RNN (bias terms discarded)

- $E \rightarrow N \times d$
- $U \rightarrow h \times N$
- $W_i \rightarrow d \times h$
- $W_h \rightarrow h \times h$

d : dimension of input embedding

h : dimension of hidden vectors

- Encoder and decoder embeddings have most of the parameters



Training RNN Language Model – Weight Tying

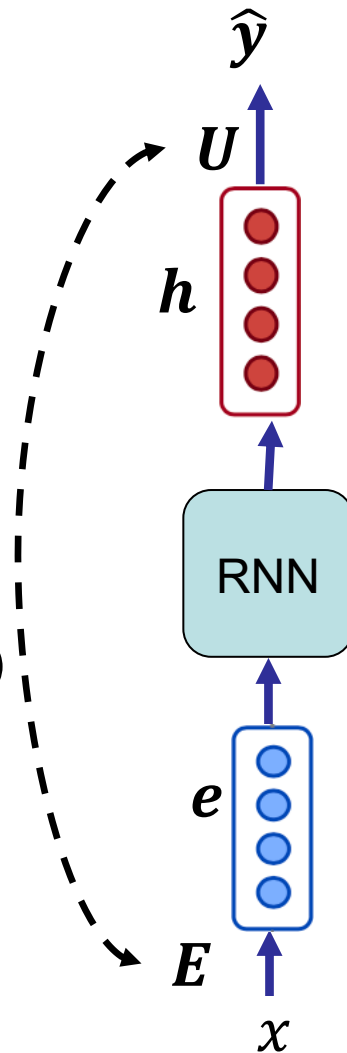
- Parameters in a vanilla RNN (bias terms discarded)

- $E \rightarrow N \times d$
- $U \rightarrow h \times N$
- $W_i \rightarrow d \times h$
- $W_h \rightarrow h \times h$

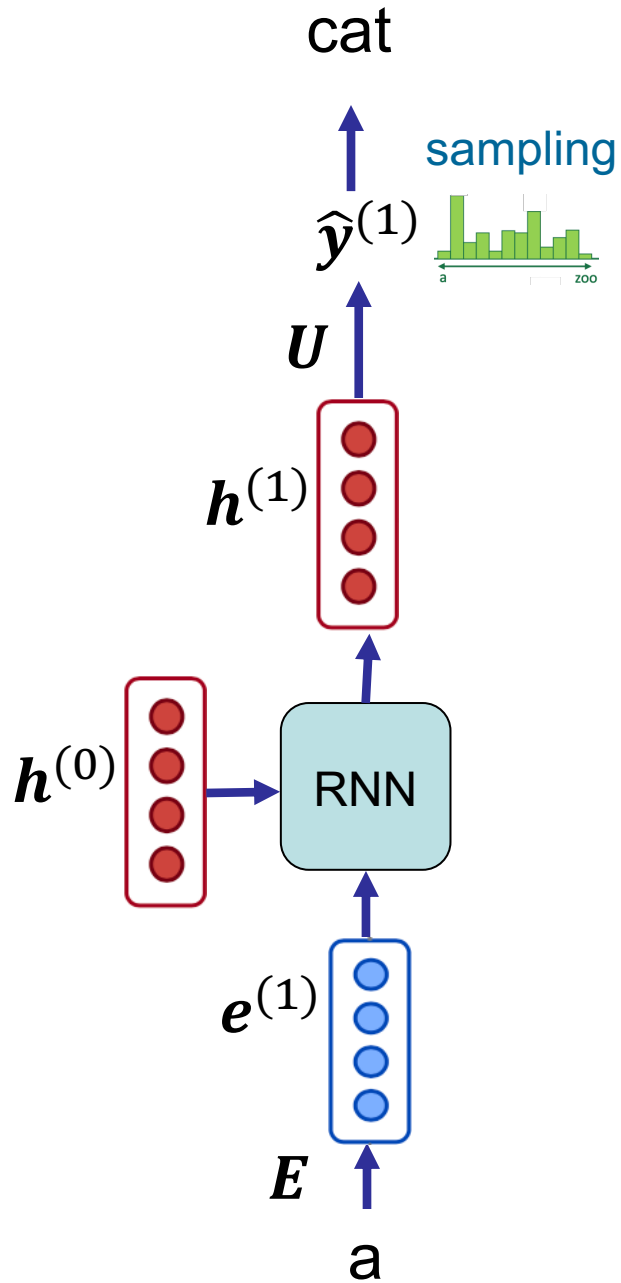
d : dimension of input embedding

h : dimension of hidden vectors

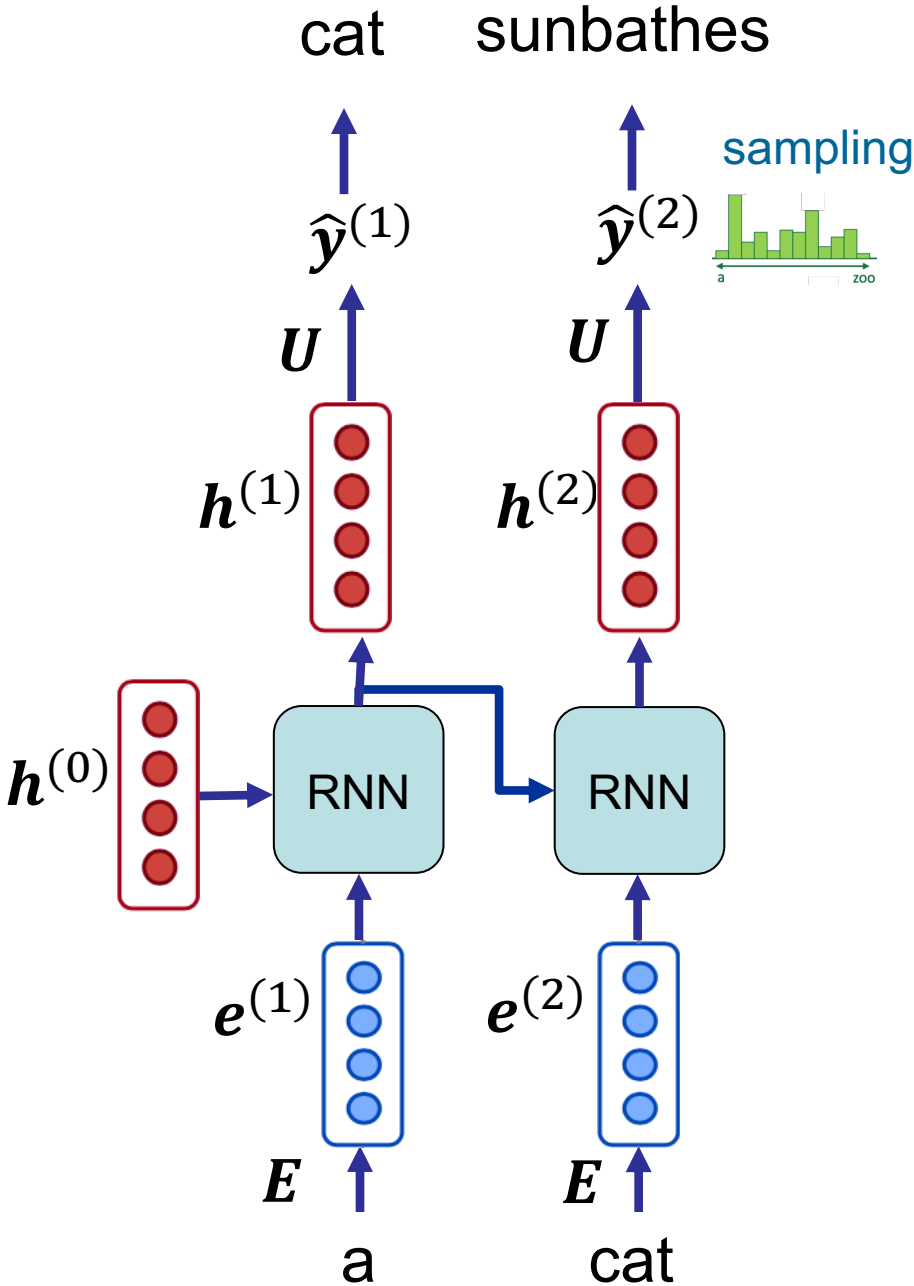
- Weight tying**: set the decoder parameters the same as encoder parameters (saving $N \times h$ decoding parameters)
 - In this case d must be equal to h
 - If $d \neq h$, usually a linear projects the output vector from h to d dimensions



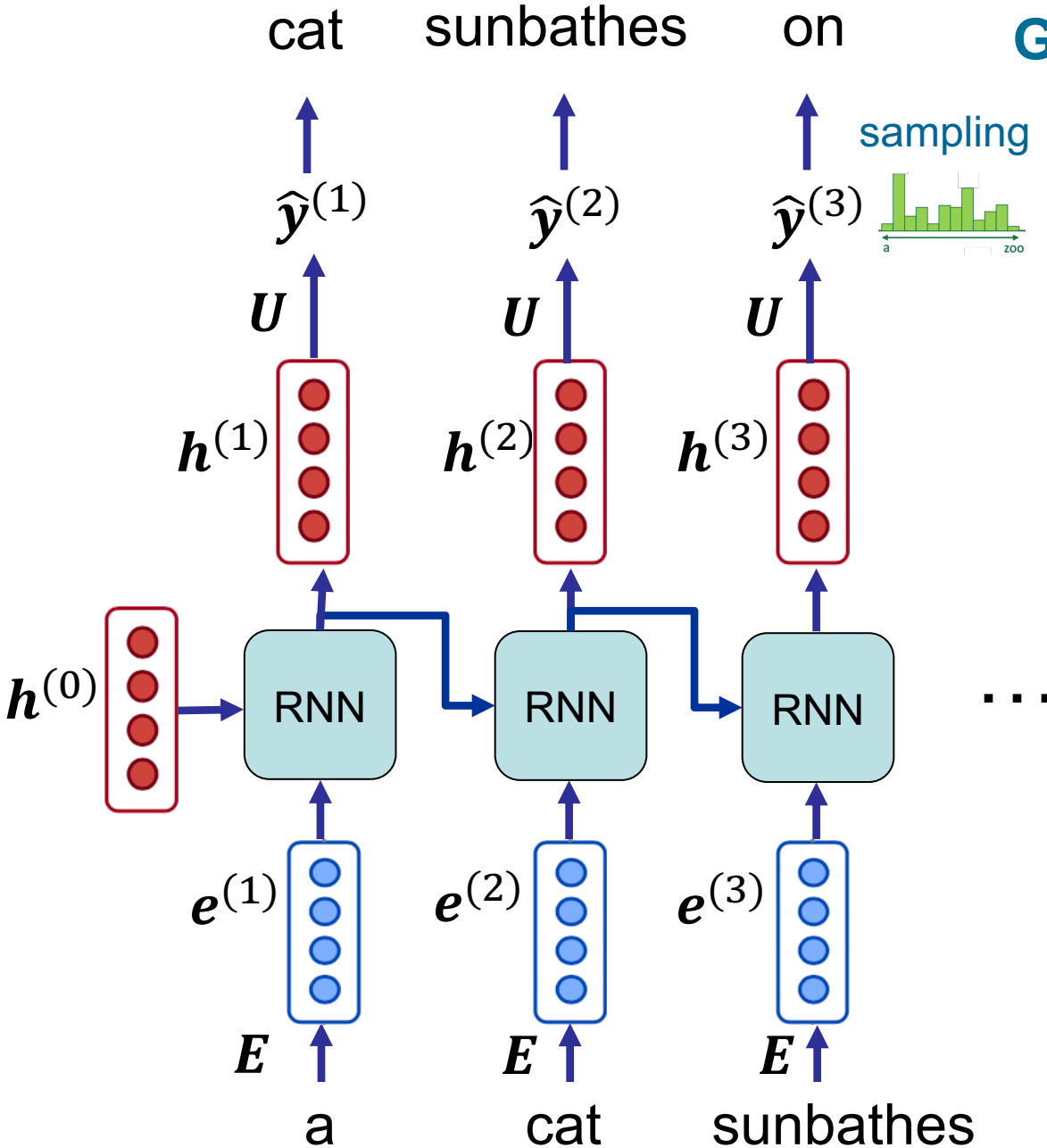
Generating text



Generating text



Generating text



Generating text with RNN Language Model

- Trained on Obama speeches



Jobs

The United States will step up to the cost of a new challenges of the American people that will share the fact that we created the problem. They were attacked and so that they have to say that all the task of the final days of war that I will not be able to get this done. The promise of the men and women who were still going to take out the fact that the American people have fought to make sure that they have to be able to protect our part. It was a chance to stand together to completely look for the commitment to borrow from the American people.

Generating text with RNN Language Model

- Trained on Trump speeches



make the country rich. it was terrible. but saudi arabia, they make a billion dollars a day. i was the king. i was the king. i was the smartest person in yemen, we have to get to business. i have to say, but he was an early starter. and we have to get to business. i have to say, donald, i can't believe it. it's so important. but this is what they're saying, dad, you're going to be really pro, growth, blah, blah. it's disgusting what's disgusting., and it was a 19 set washer and to go to japan. did you hear that character. we are going to have to think about it. but you know, i've been nice to me.

Summary

RNN for Language Modeling

- **Pros:**
 - RNN can process any length input
 - RNN can (in theory) use information from many steps back
 - Model size doesn't increase for longer input sequences

- **Cons:**
 - Recurrent computation is slow → (in its vanilla form) does not fully exploit the parallel computation capability of GPUs
 - Biased toward recent incidents → in practice, RNN has the tendency to use the information of recent steps (harder to access information from many steps back)