

Text Generation with Language Model

NHÓM 6

TRẦN TRUNG ANH – 18520473

NGUYỄN ANH KHOA – 18520923

VÕ QUỐC AN – 18520440

Giới thiệu

- ❑ Text generation: dự đoán từ tiếp theo dựa trên những từ đã có trước.
- ❑ Mục tiêu: giống nhất có thể với ngôn ngữ tự nhiên.

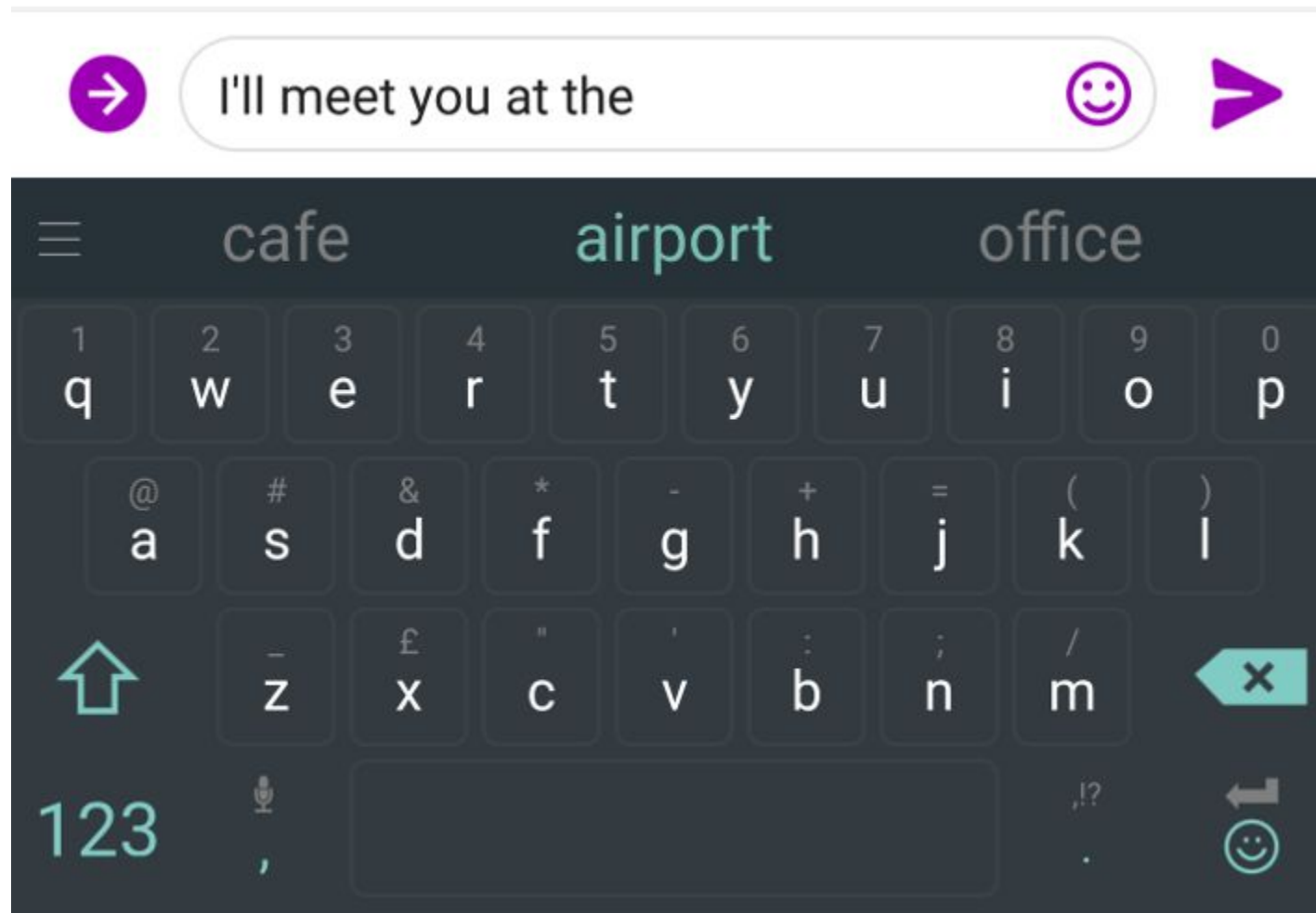
Ứng dụng

- ❑ Auto-complete.
- ❑ Hệ thống hỏi - đáp.
- ❑ Chatbot.

Ứng dụng



Ứng dụng



Phương pháp xây dựng

- ☐ Chuẩn bị dữ liệu.
- ☐ Xây dựng mô hình ngôn ngữ.
- ☐ Sử dụng mô hình để dự đoán từ tiếp theo.

Phương pháp xây dựng

Dữ liệu từ [Kaggle](#)

New York Times Articles

Số câu sử dụng: 6011

A Quote Disproved

Hot Stuff Turns Cold

At the Top

Years of Bizarre Behavior

Will the Court Stand Up to Mr. Trump?

Hope in Arizona

How Human Rights Groups Failed on Economic Equality

An American Tragedy in Nashville

We Don't Need No Education

Paul Relents

Making Change

Long Overlooked

Phương pháp xây dựng

Tiền xử lý dữ liệu:

Chuyển viết thường

Loại bỏ dấu câu

Loại bỏ kí tự đặc biệt

...

```
['a quote disproved',  
'hot stuff turns cold',  
'at the top',  
'years of bizarre behavior',  
'will the court stand up to mr trump',  
'hope in arizona',  
'how human rights groups failed on economic equality',  
'an american tragedy in nashville',  
'we dont need no education',  
'paul relents',  
'making change',  
'long overlooked']
```


Mô hình N-gram

the students opened their _____

N-gram là một chuỗi các từ liên tiếp nhau

- **uni**grams: “the”, “students”, “opened”, “their”
- **bi**grams: “the students”, “students opened”, “opened their”
- **tri**grams: “the students opened”, “students opened their”
- **4**-grams: “the students opened their”

Mô hình N-gram

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

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

Mô hình N-gram

Smoothing

- Laplace:
$$P(\boldsymbol{w}|\text{students opened their}) = \frac{\text{count}(\text{students opened their } \boldsymbol{w}) + 1}{\text{count}(\text{students opened their}) + |V|}$$
- Backoff: nếu không thể tính n -gram thì thay bằng $n-1, n-2, \dots$ gram

Mô hình N-gram

Dự đoán từ tiếp theo

In the class, the students open their _____

books: 0.153

bags: 0.056

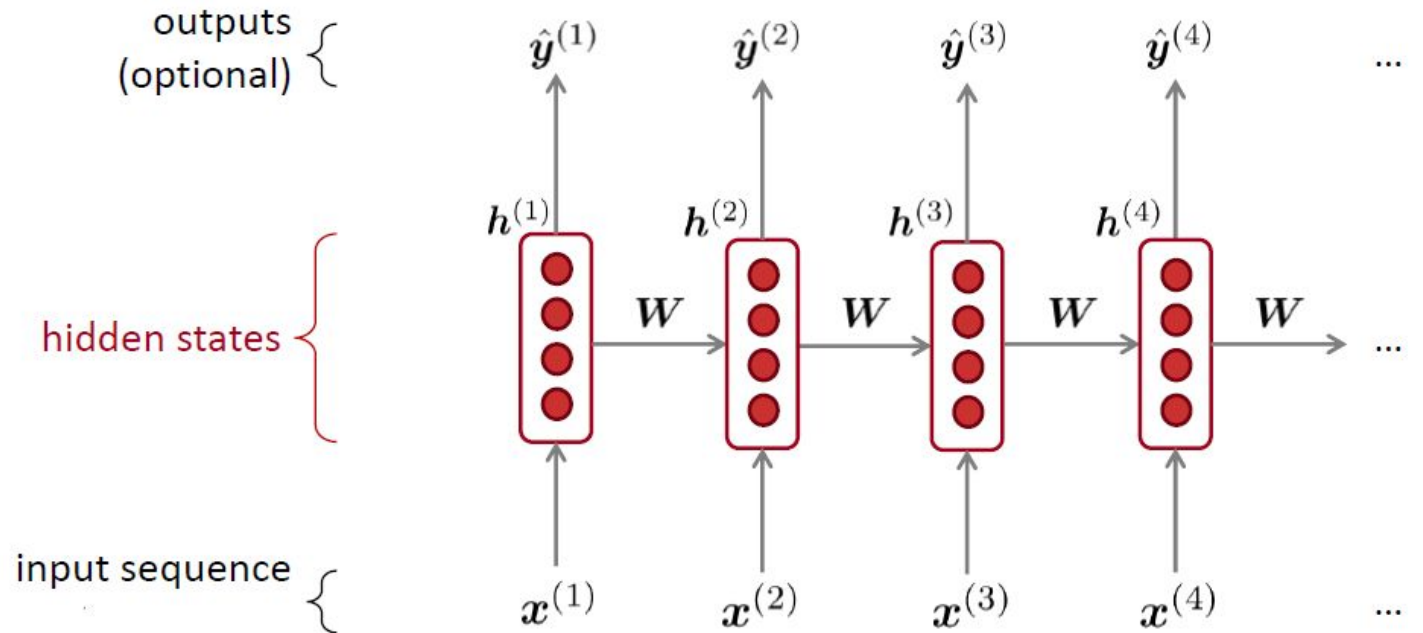
hands: 0.0032

brains: 0.00001

Recurrent Neural Networks (RNN)

$$\mathbf{h}^{(t)} = g_h(W_I \mathbf{x}^{(t)} + W_R \mathbf{h}^{(t-1)} + \mathbf{b}_h)$$

$$\mathbf{y}^{(t)} = g_y(W_y \mathbf{h}^{(t)} + \mathbf{b}_y)$$



Mô hình ngôn ngữ RNN

output distribution

$$\hat{y}^{(t)} = \text{softmax} \left(U h^{(t)} + b_2 \right) \in \mathbb{R}^{|V|}$$

hidden states

$$h^{(t)} = \sigma \left(W_h h^{(t-1)} + W_e e^{(t)} + b_1 \right)$$

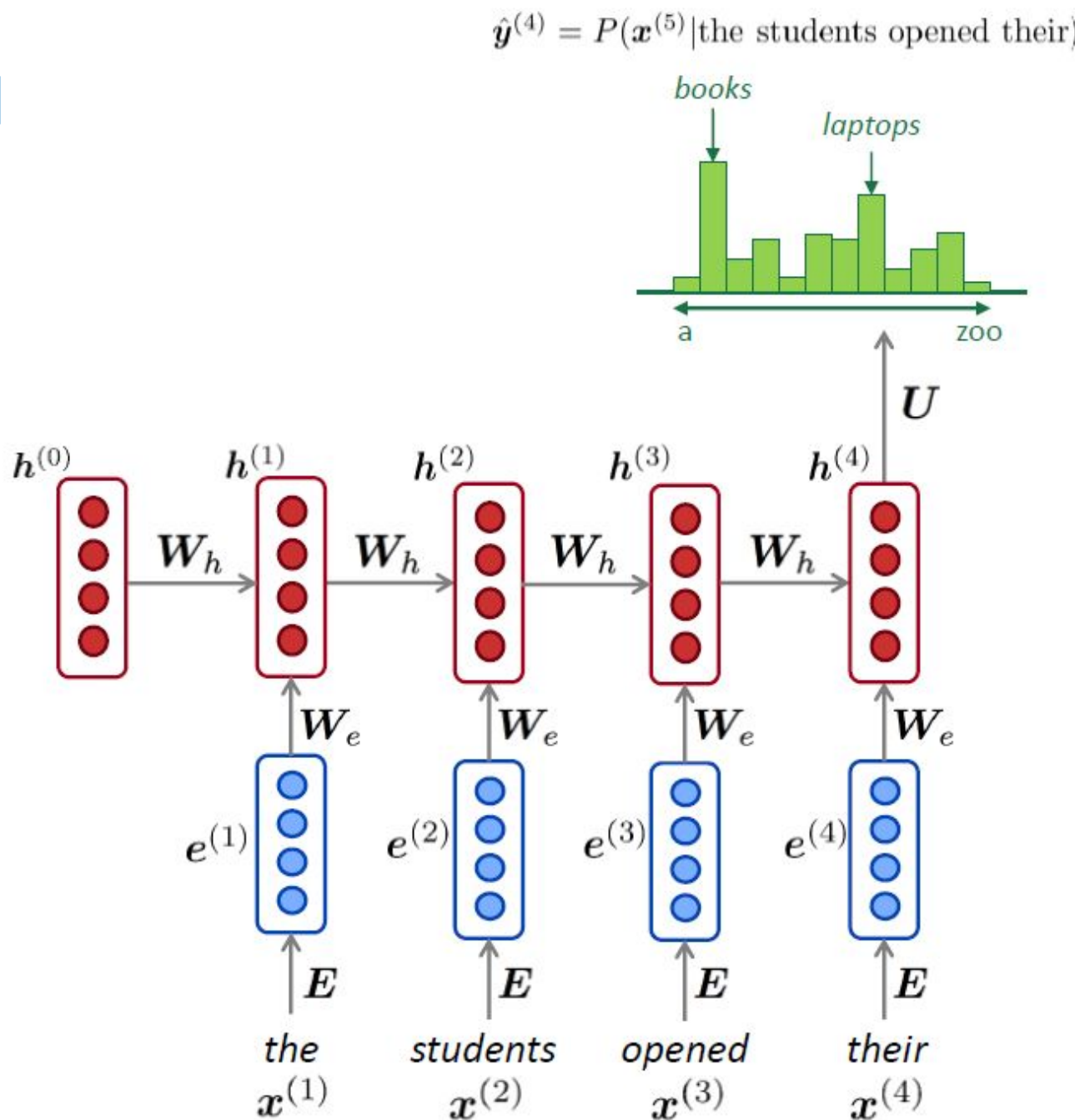
$h^{(0)}$ is the initial hidden state

word embeddings

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

words / one-hot vectors

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



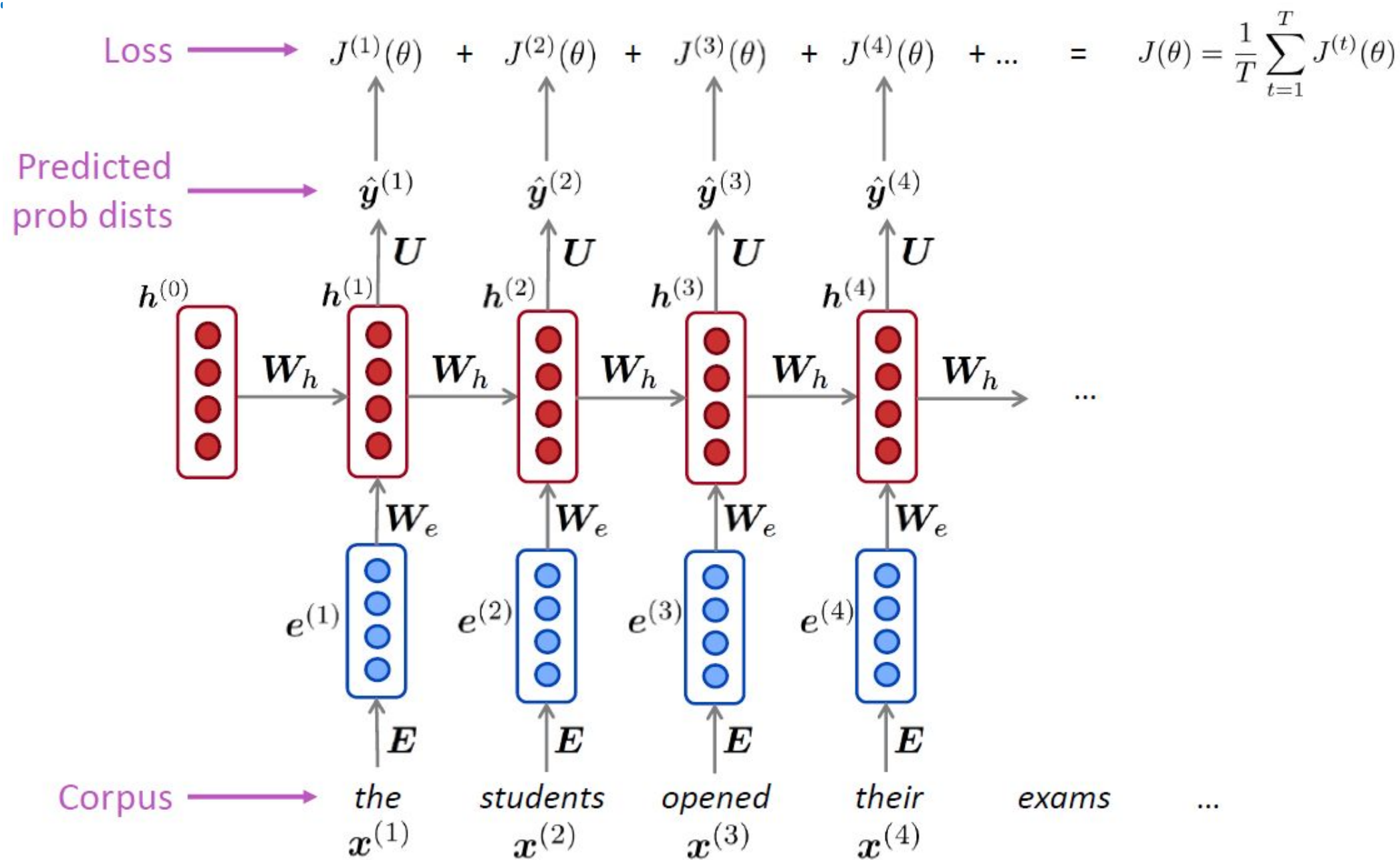
Huấn luyện RNN

- Input: one-hot của mỗi từ trong bộ dữ liệu
- Tính phân phối xác suất của tất cả các từ trong dữ liệu
- Loss function: cross-entropy

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

- Tính trung bình cho toàn tập dữ liệu

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



Vấn đề:

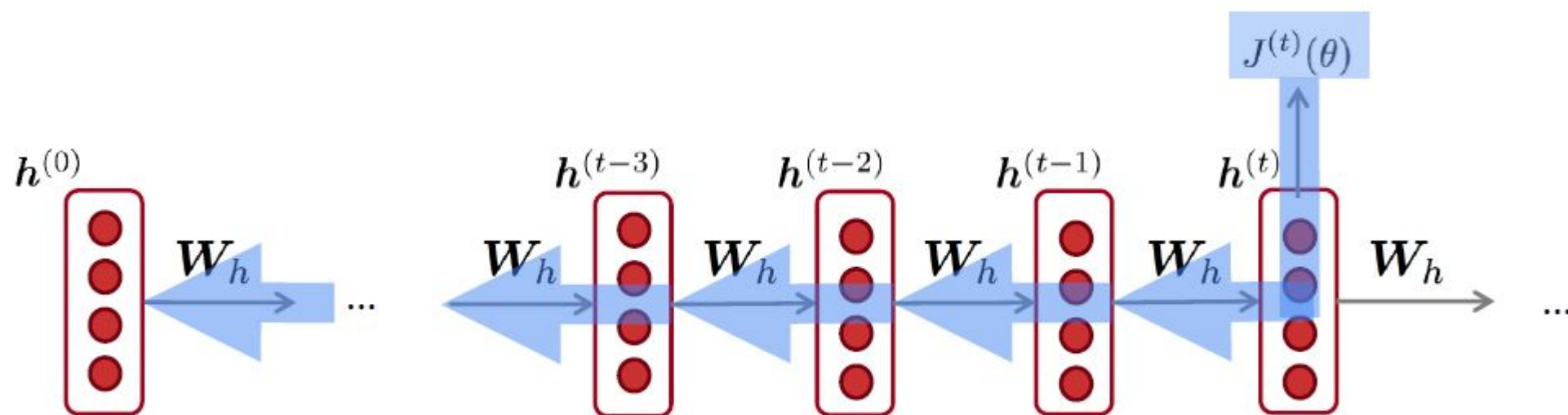
- Chi phí lớn để tính hàm mất mát và đạo hàm của toàn bộ tập dữ liệu

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

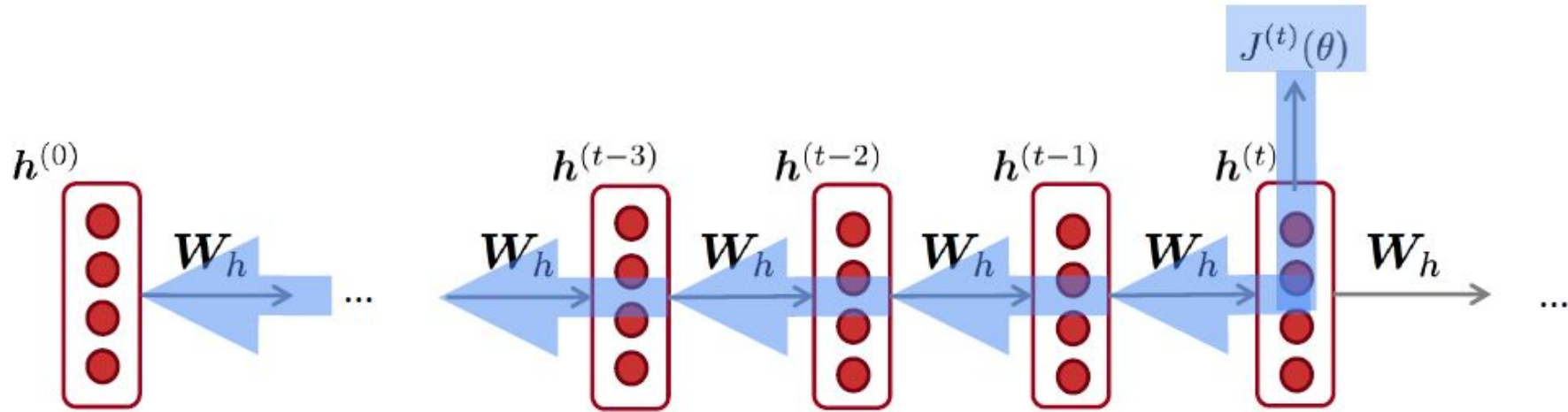
Giải pháp:

- Stochastic Gradient Descent

Backpropagation



Backpropagation through time



Tính đạo hàm theo từng
timestep $i = t, \dots, 0$
 \Rightarrow Backpropagation through time

$$\frac{\partial J^{(t)}}{\partial W_h} = \sum_{i=1}^t \frac{\partial J^{(t)}}{\partial W_h} \Big|_{(i)}$$

THANKS FOR LISTENING