# EECS 487: Introduction to Natural Language Processing

Instructor: Prof. Lu Wang

Computer Science and Engineering
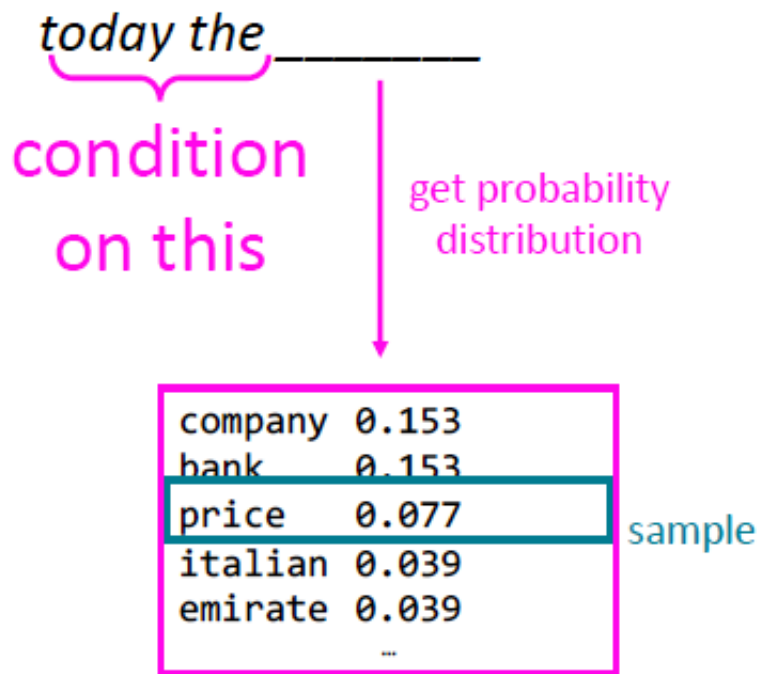
University of Michigan

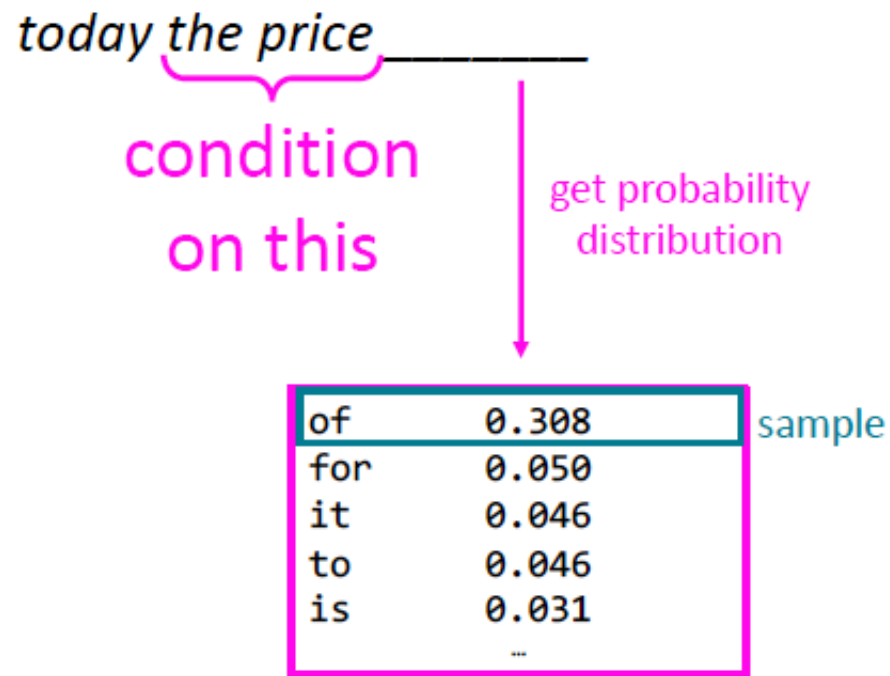Webpage: web.eecs.umich.edu/~wangluxy

1

# Outline

- Recurrent neural network (RNN) for language modeling

[Some slides are adopted from Stanford's cs224n]

# Generating Text with n-gram Language Models

today the _____

condition on this

get probability distribution

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

sample

# Generating Text with n-gram Language Models

today *the price* _____

condition on this

get probability distribution

| | |
|---|---|
| of | 0.308 |
| for | 0.050 |
| it | 0.046 |
| to | 0.046 |
| is | 0.031 |
| ... | |

sample

# Generating Text with n-gram Language Models

*today the price of* _____

condition on this

get probability distribution

| | |
|-----|-------|
| the | 0.072 |
| 18 | 0.043 |
| oil | 0.043 |
| its | 0.036 |
| gold | 0.018 |
| ... | |

sample

# How to build a neural language model?

- Recall the Language Modeling task:
  - Input: sequence of words $x^{(1)}, x^{(2)}, \ldots, x^{(t)}$
  - Output: prob dist of the next word $P(x^{(t+1)} \mid x^{(t)}, \ldots, x^{(1)})$



Predicting the next word?

$U$

$W$

museums    in    Paris    are    amazing

# A fixed-window neural Language Model

as      the      proctor  started  the      clock

**discard**

the      students  opened    their  _____

*fixed window*

# A fixed-window neural Language Model

output distribution
$$\hat{y} = \mathrm{softmax}(\boldsymbol{U}\boldsymbol{h} + \boldsymbol{b}_2) \in \mathbb{R}^{|V|}$$

hidden layer
$$\boldsymbol{h} = f(\boldsymbol{W}\boldsymbol{e} + \boldsymbol{b}_1)$$

concatenated word embeddings
$$\boldsymbol{e} = [\boldsymbol{e}^{(1)}; \boldsymbol{e}^{(2)}; \boldsymbol{e}^{(3)}; \boldsymbol{e}^{(4)}]$$

words / one-hot vectors
$$\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \boldsymbol{x}^{(3)}, \boldsymbol{x}^{(4)}$$
[1, 0, 0, …] [0, 0, 1, …]  …

books

laptops

a          zoo

$U$

$W$

the          students     opened      their
$x^{(1)}$      $x^{(2)}$        $x^{(3)}$        $x^{(4)}$

8

# A fixed-window neural Language Model

**Improvements** over *n*-gram LM:
- No sparsity problem
- Don't need to store all observed *n*-grams



$U$

$W$

the $x^{(1)}$   students $x^{(2)}$   opened $x^{(3)}$   their $x^{(4)}$
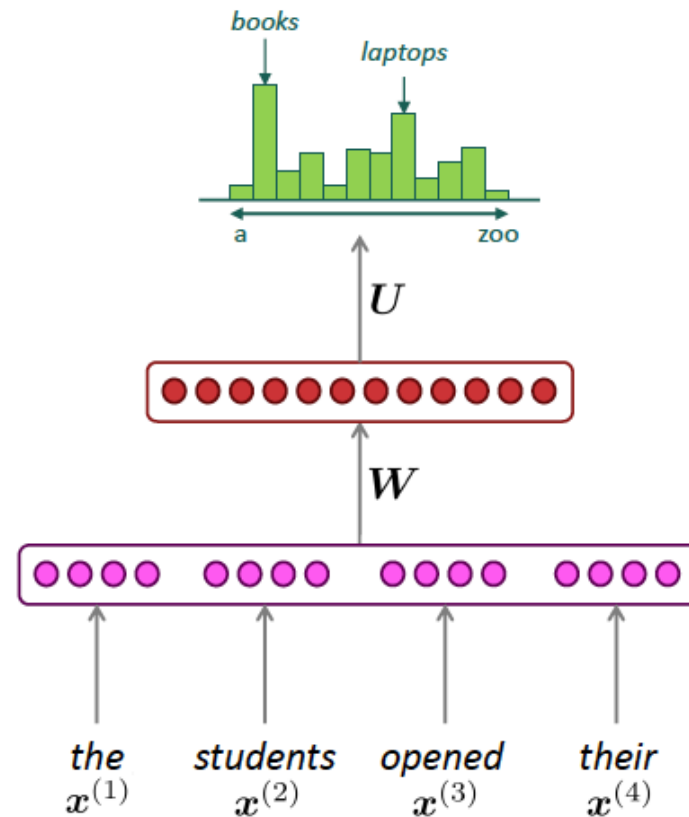
# A fixed-window neural Language Model

**Improvements** over *n*-gram LM:
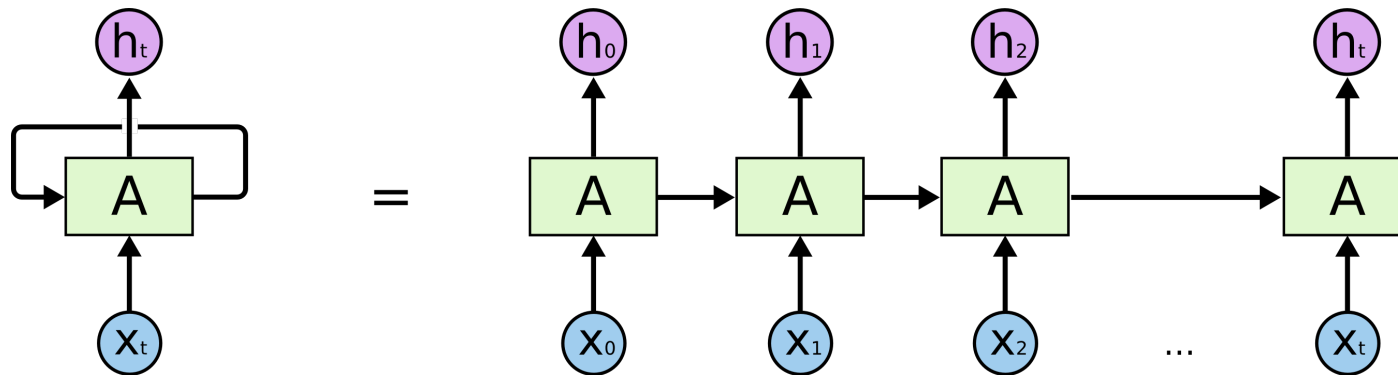- No sparsity problem
- Don't need to store all observed *n*-grams

Remaining **problems**:
- Fixed window is too small
- Enlarging window enlarges $W$
- Window can never be large enough!
- $x^{(1)}$ and $x^{(2)}$ are multiplied by completely different weights in $W$. No symmetry in how the inputs are processed.



$U$

$W$

the $x^{(1)}$  students $x^{(2)}$  opened $x^{(3)}$  their $x^{(4)}$
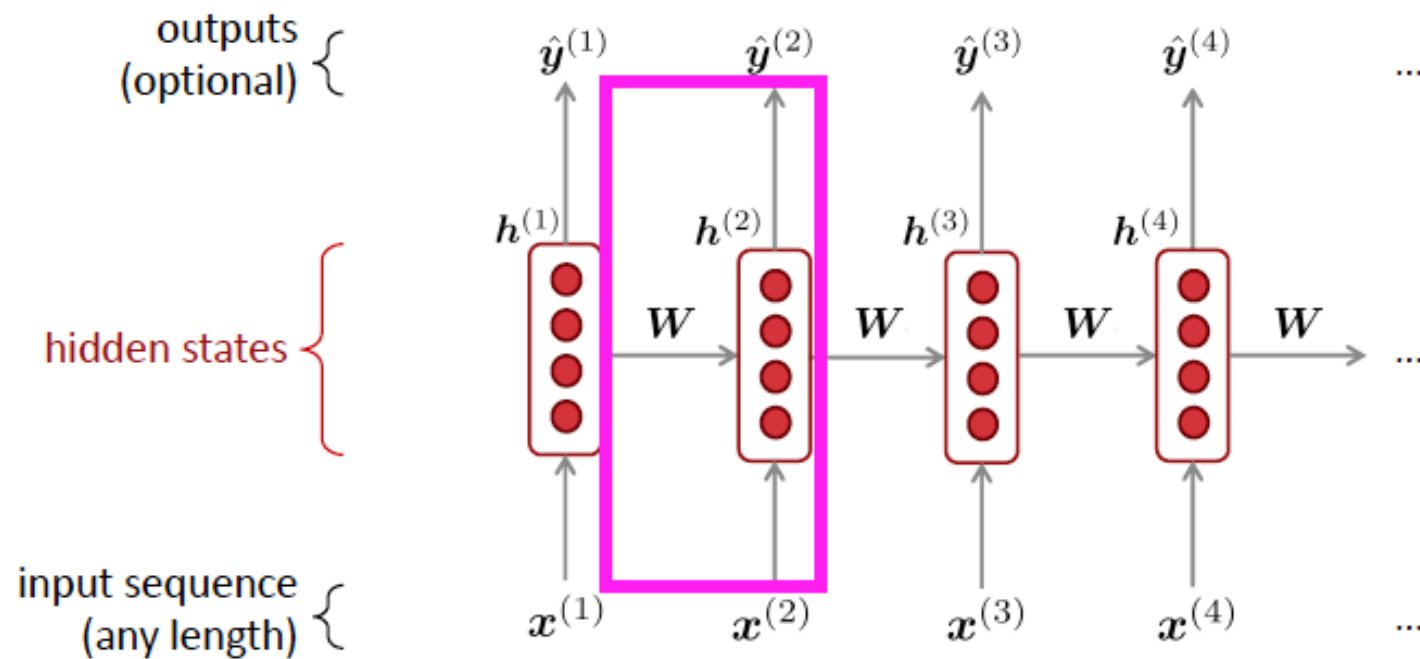
# Long Distance Dependencies

- It is very difficult to train NNs to retain information over many time steps

- This makes it very difficult to handle long-distance dependencies.

- E.g. Jane walked into the room. John walked in too. It was late in the day. Jane said hi to _?_

# Recurrent Neural Networks (RNN)

- Core idea: Apply the same weights $W$ repeatedly

# A Simple RNN Language Model

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

output distribution

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

hidden states

$$\boldsymbol{h}^{(t)} = \sigma\left(\boldsymbol{W}_h\boldsymbol{h}^{(t-1)} + \boldsymbol{W}_e\boldsymbol{e}^{(t)} + \boldsymbol{b}_1\right)$$

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

word embeddings

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

words / one-hot vectors

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



books

laptops

a                    zoo

$\boldsymbol{U}$

$\boldsymbol{h}^{(0)}$  $\boldsymbol{h}^{(1)}$  $\boldsymbol{h}^{(2)}$  $\boldsymbol{h}^{(3)}$  $\boldsymbol{h}^{(4)}$

$\boldsymbol{W}_h$  $\boldsymbol{W}_h$  $\boldsymbol{W}_h$  $\boldsymbol{W}_h$

$\boldsymbol{W}_e$  $\boldsymbol{W}_e$  $\boldsymbol{W}_e$  $\boldsymbol{W}_e$

$\boldsymbol{e}^{(1)}$  $\boldsymbol{e}^{(2)}$  $\boldsymbol{e}^{(3)}$  $\boldsymbol{e}^{(4)}$

$\boldsymbol{E}$  $\boldsymbol{E}$  $\boldsymbol{E}$  $\boldsymbol{E}$

the        students      opened        their
$\boldsymbol{x}^{(1)}$   $\boldsymbol{x}^{(2)}$   $\boldsymbol{x}^{(3)}$   $\boldsymbol{x}^{(4)}$

13

# Pros and Cons

RNN **Advantages**:
- Can process any length input
- Computation for step *t* can (in theory) use information from many steps back
- Model size doesn't increase for longer input context
- Same weights applied on every timestep, so there is symmetry in how inputs are processed.

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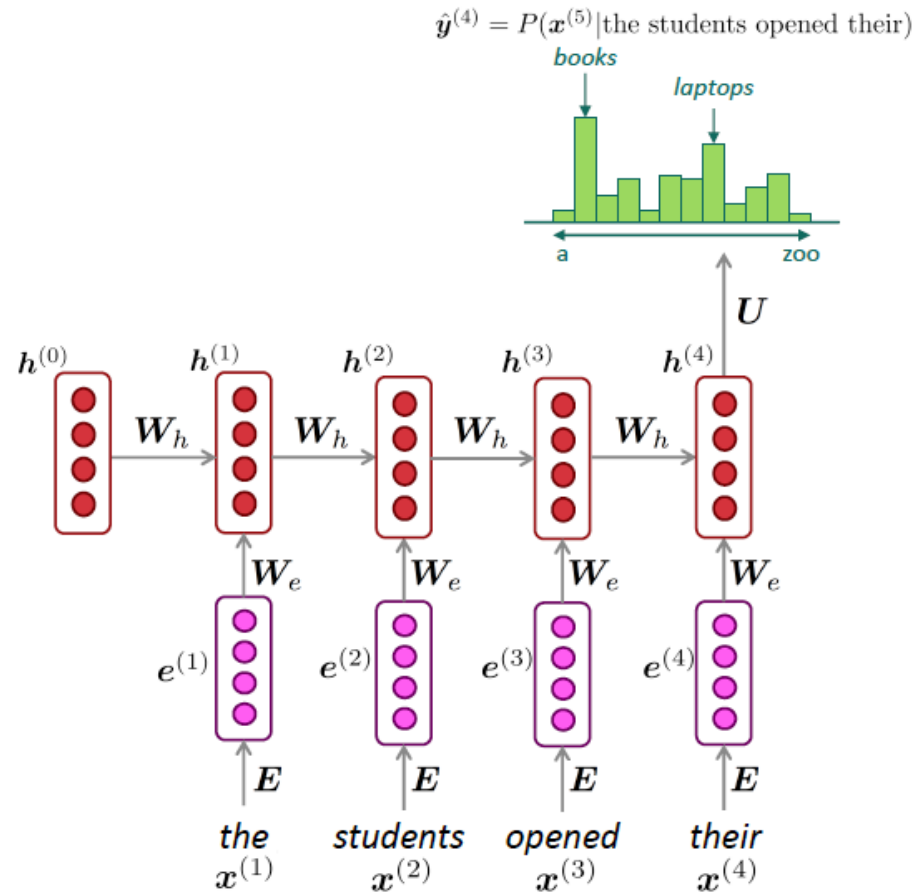
# Pros and Cons

**RNN Advantages:**
- Can process any length input
- Computation for step *t* can (in theory) use information from many steps back
- Model size doesn't increase for longer input context
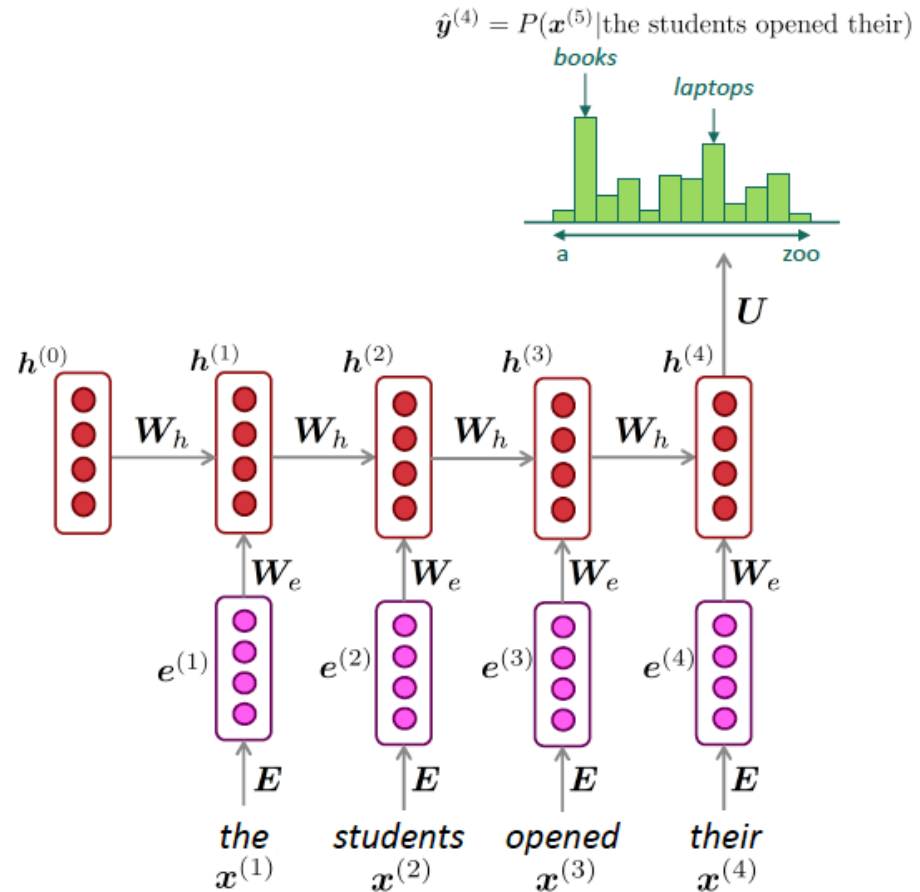- Same weights applied on every timestep, so there is symmetry in how inputs are processed.

**RNN Disadvantages:**
- Recurrent computation is slow
- In practice, difficult to access information from many steps back

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

*books*

*laptops*

*a*     *zoo*

$U$

$h^{(0)}$   $h^{(1)}$   $h^{(2)}$   $h^{(3)}$   $h^{(4)}$

$W_h$   $W_h$   $W_h$   $W_h$

$W_e$   $W_e$   $W_e$   $W_e$

$e^{(1)}$   $e^{(2)}$   $e^{(3)}$   $e^{(4)}$

$E$   $E$   $E$   $E$

*the*   *students*   *opened*   *their*
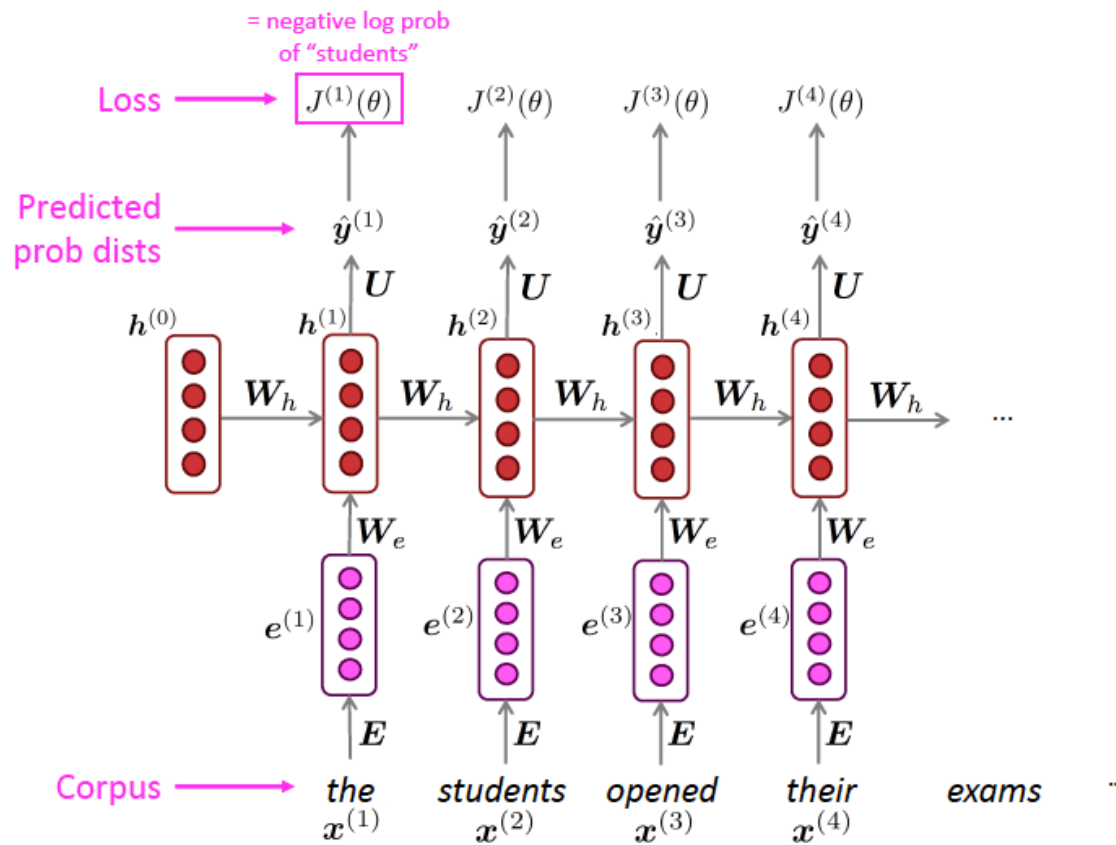
$x^{(1)}$   $x^{(2)}$   $x^{(3)}$   $x^{(4)}$

15

# Training an RNN Language Model

- Get a big corpus of text which is a sequence of words $x^{(1)}, \ldots, x^{(T)}$
- Feed into RNN-LM; compute output distribution $\hat{y}^{(t)}$ 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 $\hat{y}^{(t)}$, and the true next word $y^{(t)}$ (one-hot for $x^{(t+1)}$):

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

# Training an RNN Language Model

# Training an RNN Language Model



= negative log prob
of "opened"

Loss $\longrightarrow$ $J^{(1)}(\theta)$ $\boxed{J^{(2)}(\theta)}$ $J^{(3)}(\theta)$ $J^{(4)}(\theta)$

Predicted $\longrightarrow$ $\hat{y}^{(1)}$ $\hat{y}^{(2)}$ $\hat{y}^{(3)}$ $\hat{y}^{(4)}$
prob dists

Corpus $\longrightarrow$ the students opened their exams ...
$x^{(1)}$ $x^{(2)}$ $x^{(3)}$ $x^{(4)}$

# Training an RNN Language Model

# Training an RNN Language Model



= negative log prob of "exams"

Loss ⟶ $J^{(1)}(\theta)$     $J^{(2)}(\theta)$     $J^{(3)}(\theta)$     $\boxed{J^{(4)}(\theta)}$

Predicted prob dists ⟶ $\hat{y}^{(1)}$     $\hat{y}^{(2)}$     $\hat{y}^{(3)}$     $\hat{y}^{(4)}$

$U$   $U$   $U$   $U$

$h^{(0)}$   $h^{(1)}$   $h^{(2)}$   $h^{(3)}$   $h^{(4)}$

$W_h$   $W_h$   $W_h$   $W_h$   $W_h$   ...

$W_e$   $W_e$   $W_e$   $W_e$

$e^{(1)}$   $e^{(2)}$   $e^{(3)}$   $e^{(4)}$

$E$   $E$   $E$   $E$

Corpus ⟶ the    students    opened    their    exams   ...

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

# Training an RNN Language Model

Loss $\longrightarrow$ $J^{(1)}(\theta)$ + $J^{(2)}(\theta)$ + $J^{(3)}(\theta)$ + $J^{(4)}(\theta)$ + ... = $J(\theta) = \frac{1}{T}\sum_{t=1}^{T} J^{(t)}(\theta)$

Predicted prob dists $\longrightarrow$ $\hat{y}^{(1)}$ $\hat{y}^{(2)}$ $\hat{y}^{(3)}$ $\hat{y}^{(4)}$

$U$ $U$ $U$ $U$

$h^{(0)}$ $h^{(1)}$ $h^{(2)}$ $h^{(3)}$ $h^{(4)}$

$W_h$ $W_h$ $W_h$ $W_h$ $W_h$ ...

$W_e$ $W_e$ $W_e$ $W_e$

$e^{(1)}$ $e^{(2)}$ $e^{(3)}$ $e^{(4)}$

$E$ $E$ $E$ $E$

Corpus $\longrightarrow$ the   students   opened   their   exams   ...

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