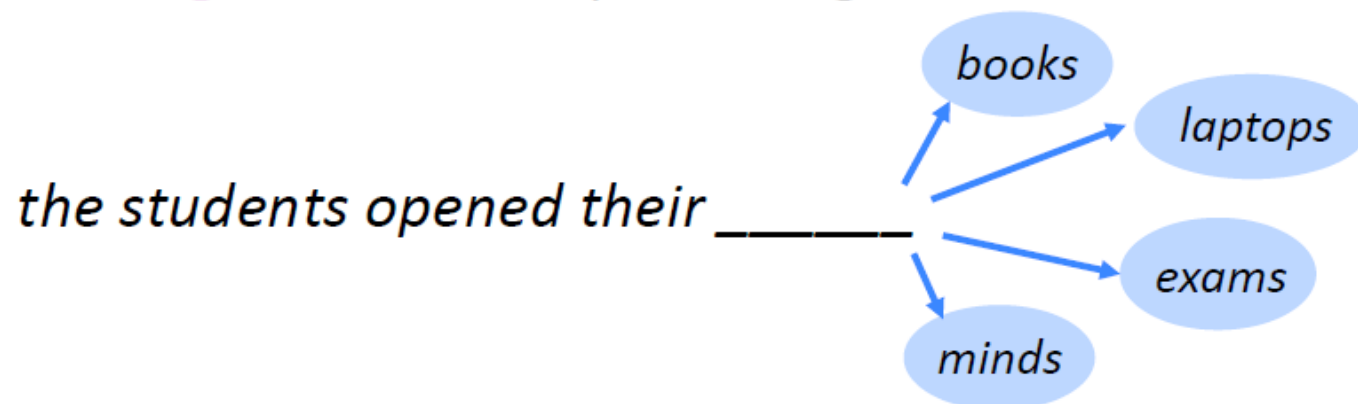# Recurrent Neural Networks (RNN)

# Language Modeling

- **Language Modeling** is the task of predicting what word comes next.

*the students opened their _____*

- books
- laptops
- exams
- minds

- More formally: given a sequence of words $x^{(1)}, x^{(2)}, \ldots, x^{(t)}$, compute the probability distribution of the next word $x^{(t+1)}$ :

$$P(x^{(t+1)} \mid x^{(t)}, \ldots, x^{(1)})$$

where $x^{(t+1)}$ can be any word in the vocabulary $V = \{w_1, \ldots, w_{|V|}\}$

- A system that does this is called a **Language Model**.

# n-gram Language Models

- First we make a simplifying assumption: $x^{(t+1)}$ depends only on the preceding *n-1* words.

$$\underbrace{P(x^{(t+1)}|x^{(t)},\ldots,x^{(1)}) = P(x^{(t+1)}|\overbrace{x^{(t)},\ldots,x^{(t-n+2)}}^{\text{n-1 words}})} \qquad \text{(assumption)}$$

prob of a n-gram

prob of a (n-1)-gram

$$= \frac{P(x^{(t+1)},x^{(t)},\ldots,x^{(t-n+2)})}{P(x^{(t)},\ldots,x^{(t-n+2)})} \qquad \begin{array}{l}\text{(definition of} \\ \text{conditional prob)}\end{array}$$

- **Question:** How do we get these *n*-gram and (*n*-1)-gram probabilities?
- **Answer:** By counting them in some large corpus of text!

$$\approx \frac{\text{count}(x^{(t+1)},x^{(t)},\ldots,x^{(t-n+2)})}{\text{count}(x^{(t)},\ldots,x^{(t-n+2)})} \qquad \begin{array}{l}\text{(statistical} \\ \text{approximation)}\end{array}$$

# A fixed-window neural language model



output distribution

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

hidden layer

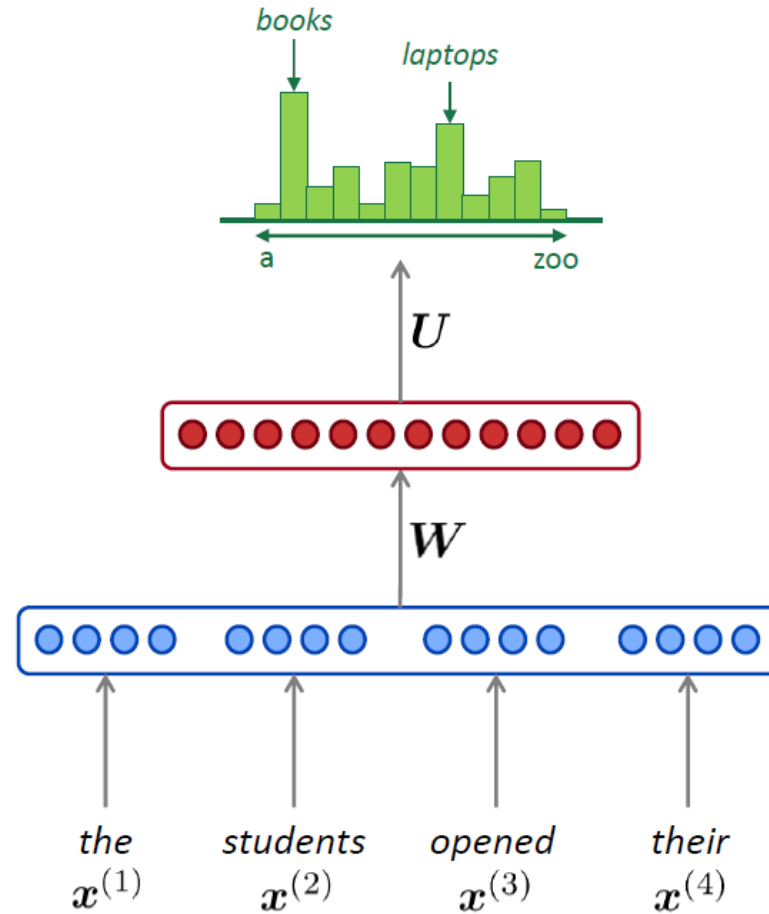$$h = f(\boldsymbol{W}e + \boldsymbol{b}_1)$$

concatenated word embeddings

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

words / one-hot vectors

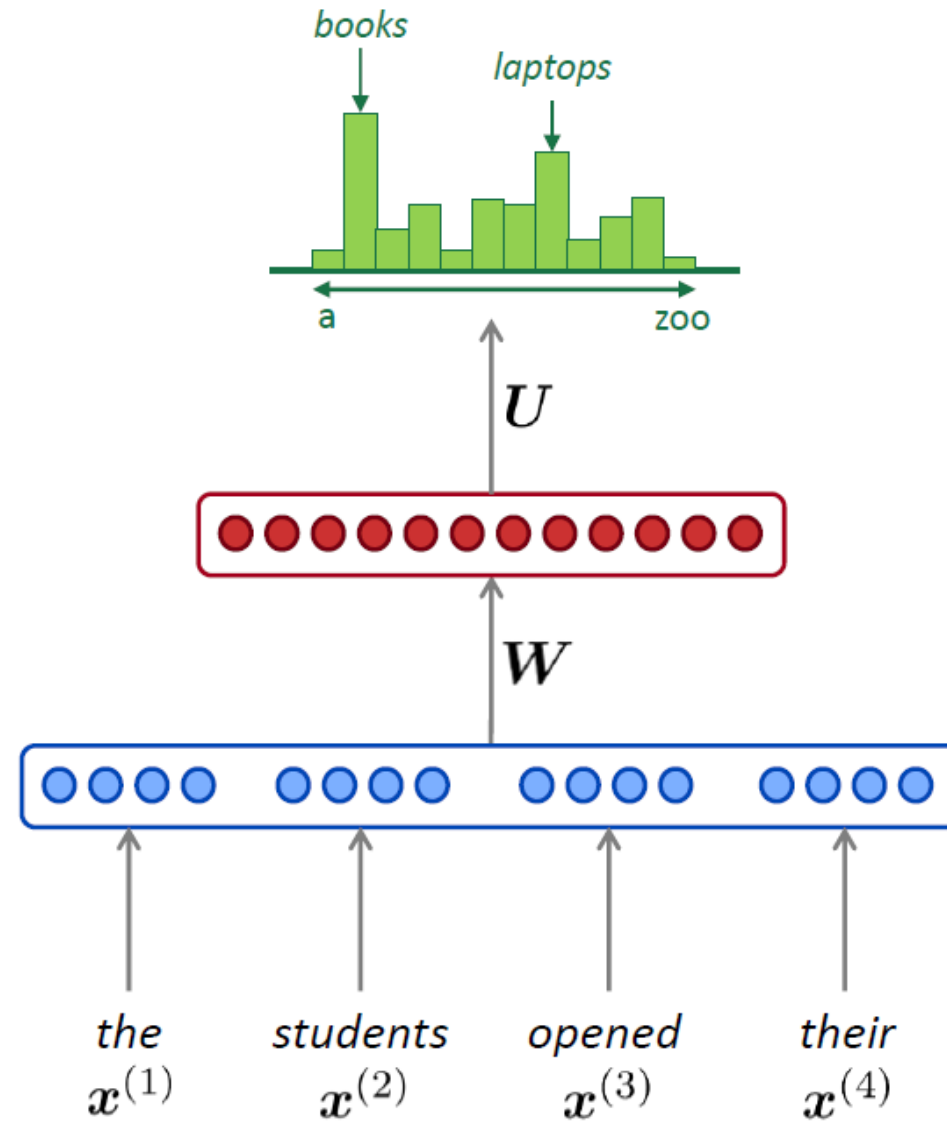$$x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$$

# 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.
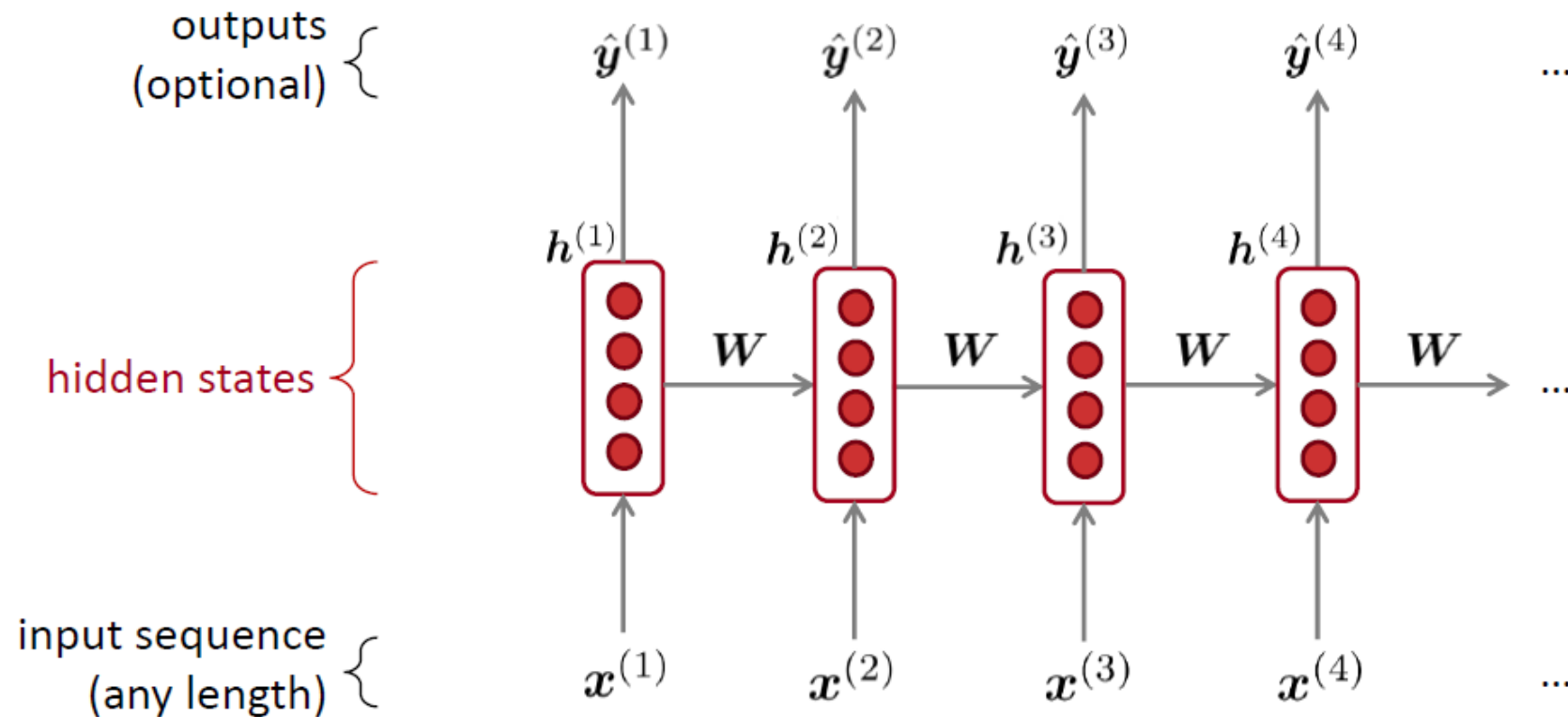
We need a neural architecture that can process *any length input*

books

laptops

a                    zoo

$U$

$W$

the
$x^{(1)}$

students
$x^{(2)}$

opened
$x^{(3)}$

their
$x^{(4)}$

# Recurrent Neural Networks (RNN)

A family of neural architectures

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

outputs (optional) {

$\hat{y}^{(1)}$     $\hat{y}^{(2)}$     $\hat{y}^{(3)}$     $\hat{y}^{(4)}$     ...

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

hidden states {

$W$    $W$    $W$    $W$    ...

input sequence (any length) {

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

# A RNN Language Model

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

books

laptops

output distribution

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

a                    zoo

$U$

hidden states

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

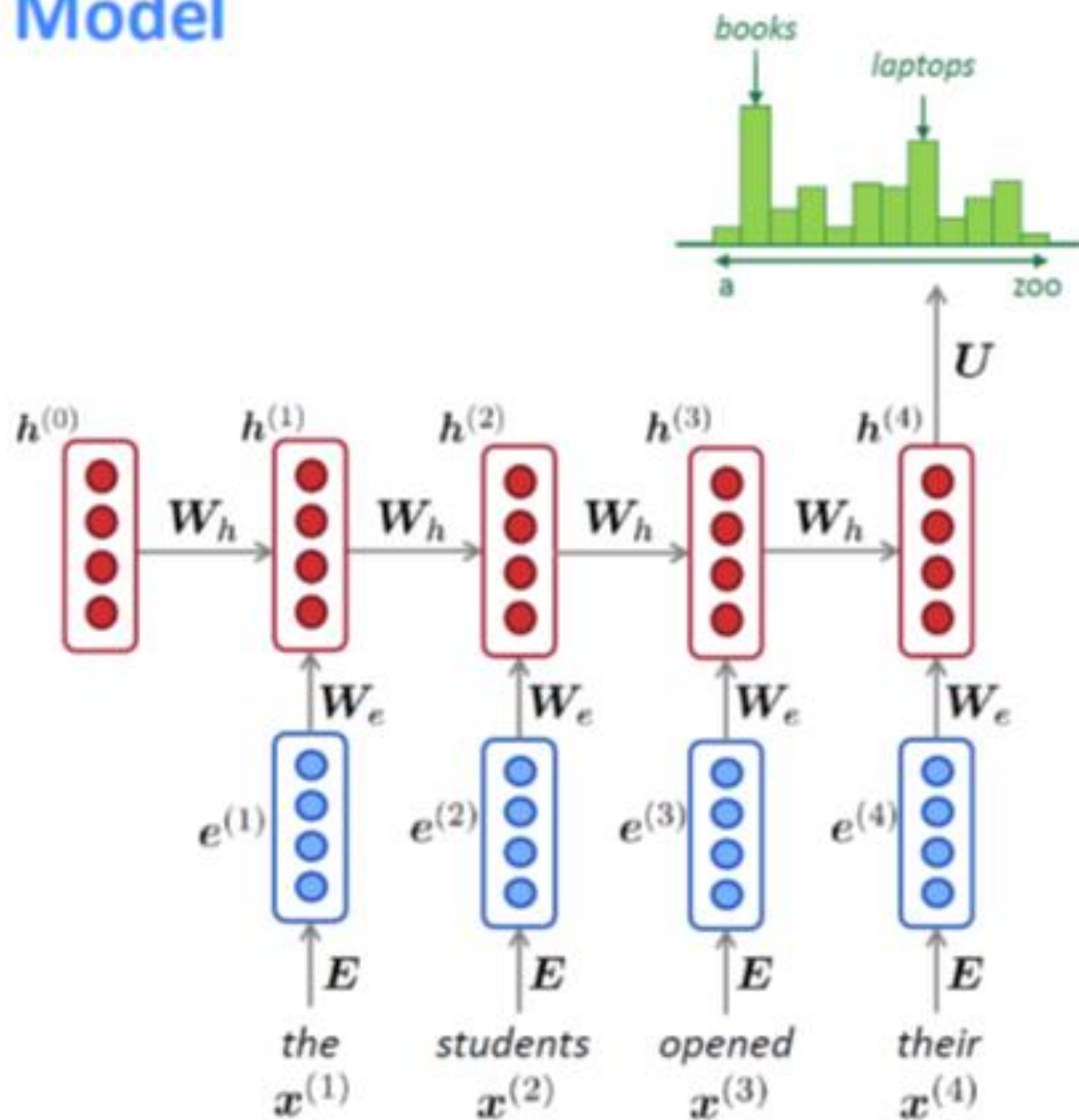$h^{(0)}$ is the initial hidden state

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

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

word embeddings

$$e^{(t)} = Ex^{(t)}$$

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

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

words / one-hot vectors

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

the          students       opened        their

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

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
- 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

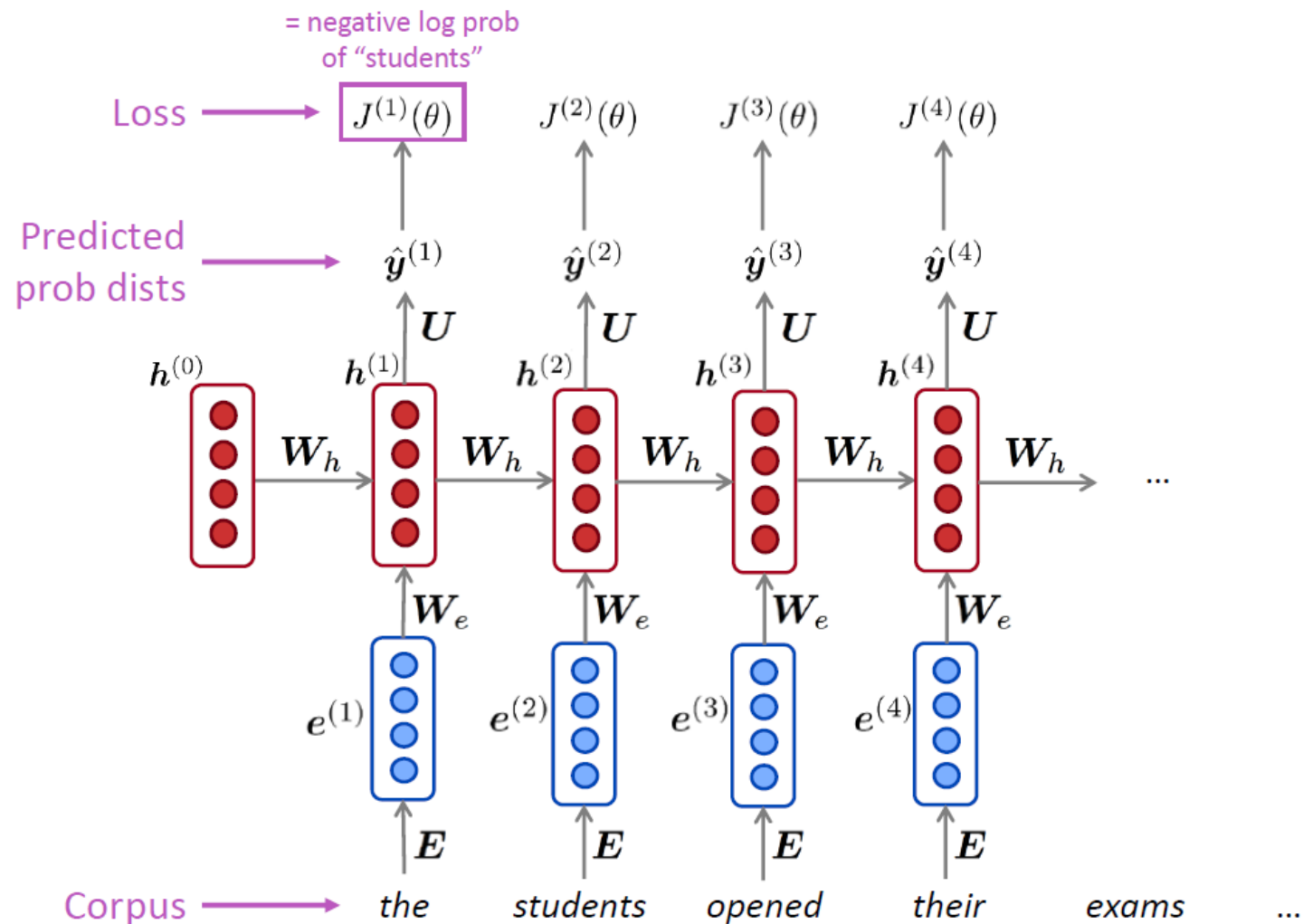More on these later in the course

# Training the RNN 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 $x^{(t+1)}$):

$$J^{(t)}(\theta) = CE(y^{(t)}, \hat{y}^{(t)}) = -\sum_{w \in V} y_w^{(t)} \log \hat{y}_w^{(t)} = -\log \hat{y}_{x_{t+1}}^{(t)}$$
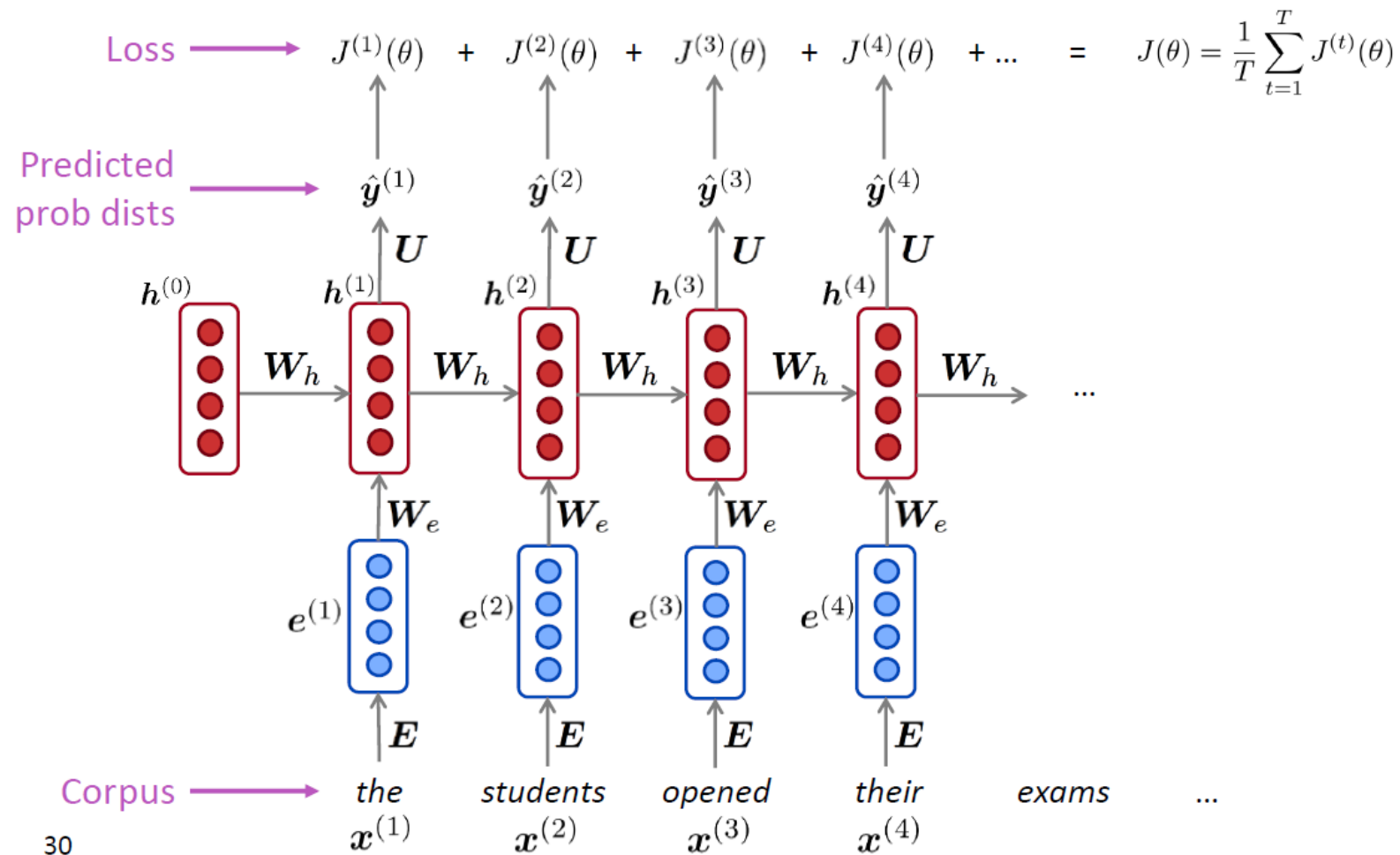
- Average this to get overall loss for entire training set:

$$J(\theta) = \frac{1}{T} \sum_{t=1}^{T} J^{(t)}(\theta) = \frac{1}{T} \sum_{t=1}^{T} -\log \hat{y}_{x_{t+1}}^{(t)}$$
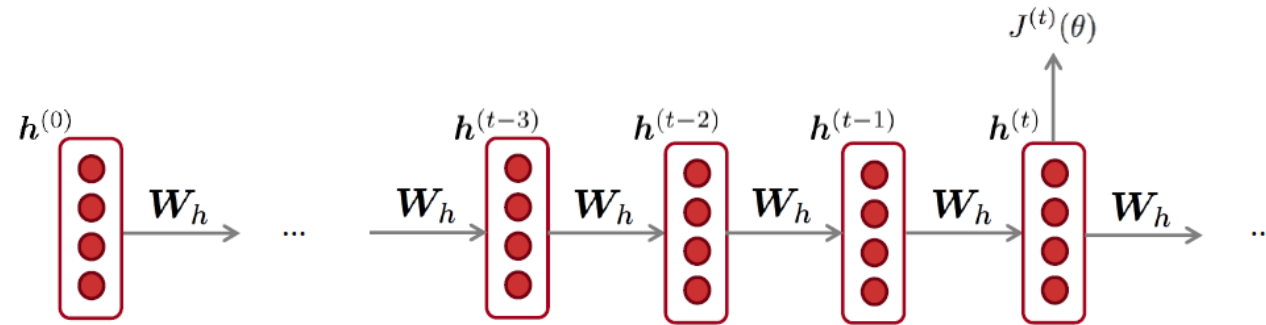
# Training the RNN model

# Training the RNN model



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

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

$\boldsymbol{U}$ $\boldsymbol{U}$ $\boldsymbol{U}$ $\boldsymbol{U}$

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

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

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

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

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

Corpus ⟶ the students opened their exams ...

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

30

# Backpropagation

## Backpropagation for RNNs



**Question:** What's the derivative of $J^{(t)}(\theta)$ w.r.t. the repeated weight matrix $W_h$ ?

**Answer:** 
$$\frac{\partial J^{(t)}}{\partial W_h} = \sum_{i=1}^{t} \left. \frac{\partial J^{(t)}}{\partial W_h} \right|_{(i)}$$

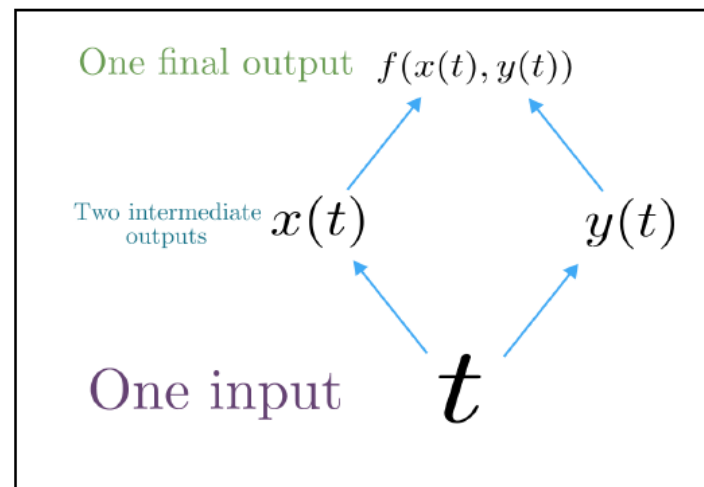> "The gradient w.r.t. a repeated weight is the sum of the gradient w.r.t. each time it appears"

**Why?**

# MCR

## Multivariable Chain Rule

- Given a multivariable function $f(x, y)$, and two single variable functions $x(t)$ and $y(t)$, here's what the multivariable chain rule says:

$$\underbrace{\frac{d}{dt} f(x(t), y(t))}_{\text{Derivative of composition function}} = \frac{\partial f}{\partial x}\frac{dx}{dt} + \frac{\partial f}{\partial y}\frac{dy}{dt}$$
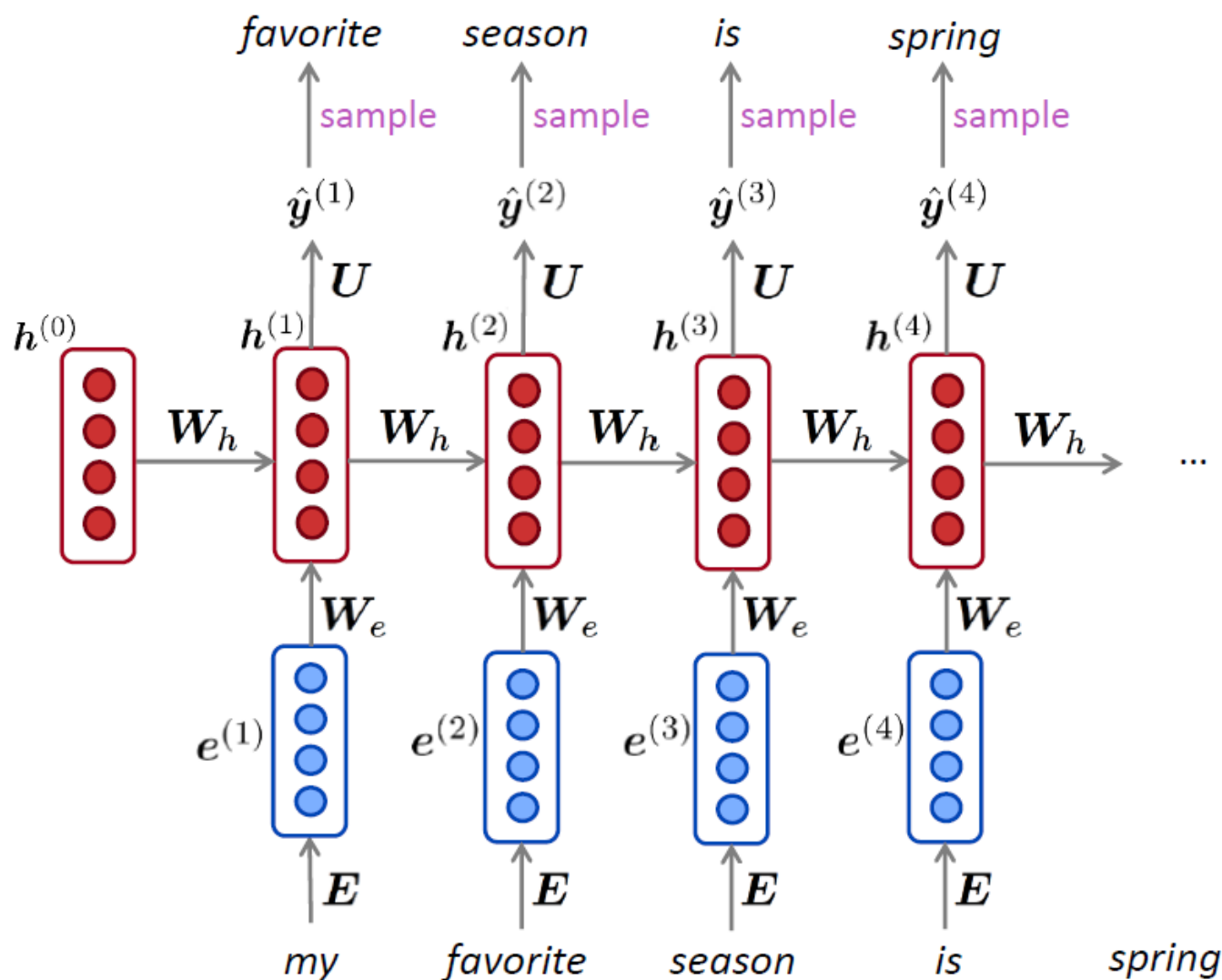
One final output $f(x(t), y(t))$

Two intermediate outputs $x(t)$ $y(t)$
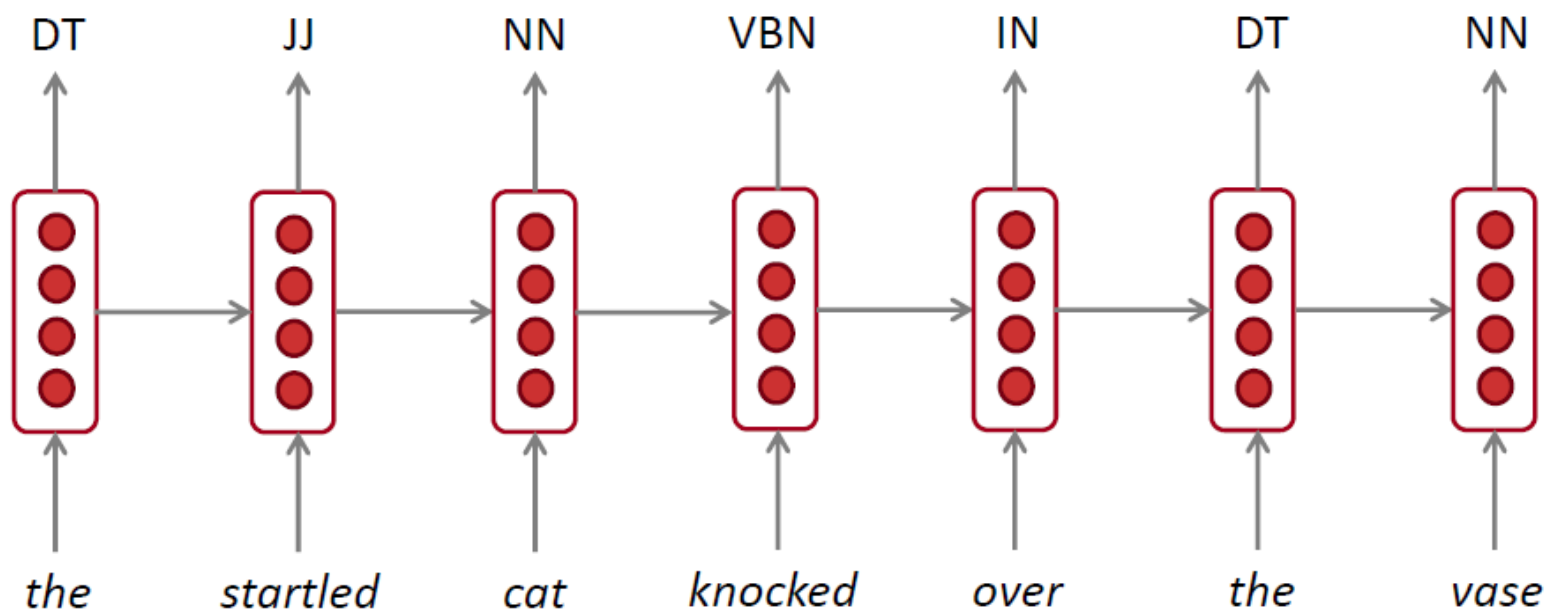
One input $t$

# Lets have a real example

☺

# 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 is next step's input.

# RNNs can be used for tagging

e.g. part-of-speech tagging, named entity recognition

# RNNs can be used for sentence classification

e.g. sentiment classification

positive

How to compute sentence encoding?

Sentence encoding



overall     I     enjoyed     the     movie     a     lot