# Computational Linguistics LT3233

Jixing Li

Lecture 12: Recurrent Neural Network

Slides adapted from Chris Manning

# Lecture plan

- Recap: Language model with FFNN
- RNN
- Short break (15 mins)
- Hands-on exercises

# Language model

**Language Model (LM): A system that predicts the next word**

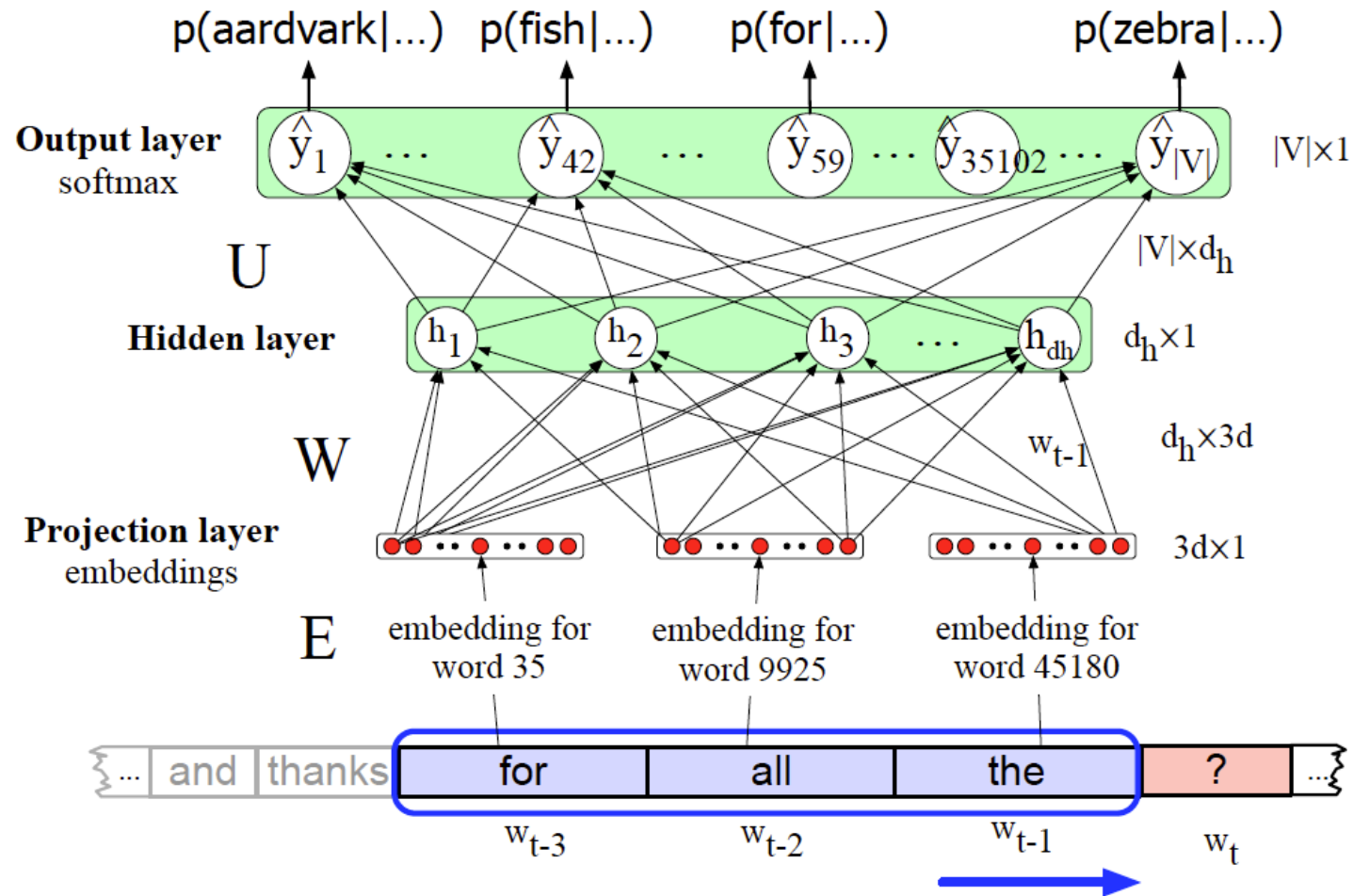**Goal:** Compute the probability of a sentence or sequence of words

**P(***a pile of shaving cream***)**     **or**     **P(***cream|a pile of shaving***)**

**P(***W***) = P(***w_1, w_2, w_3, ... w_n***)**           **P(***w_n|w_1, w_2, w_3, ... w_{n-1}***)**

**N-gram language model:**

$$\mathbf{P}(w_i|w_{i-1}) = \frac{\mathbf{Count}(w_{i-1}, w_i)}{\mathbf{Count}(w_{i-1})}$$

# Simple FFNN language model

Predicting next word $w_t$ given prior words $w_{t-1}$, $w_{t-2}$, $w_{t-3}$, ... using **sliding windows (of fixed length)**
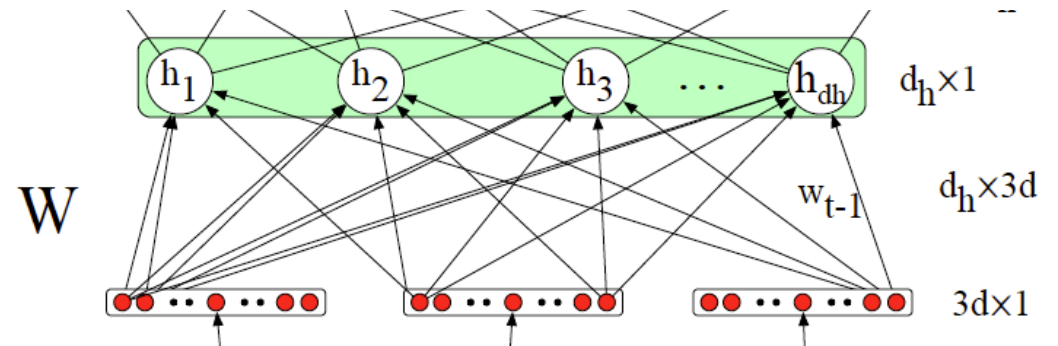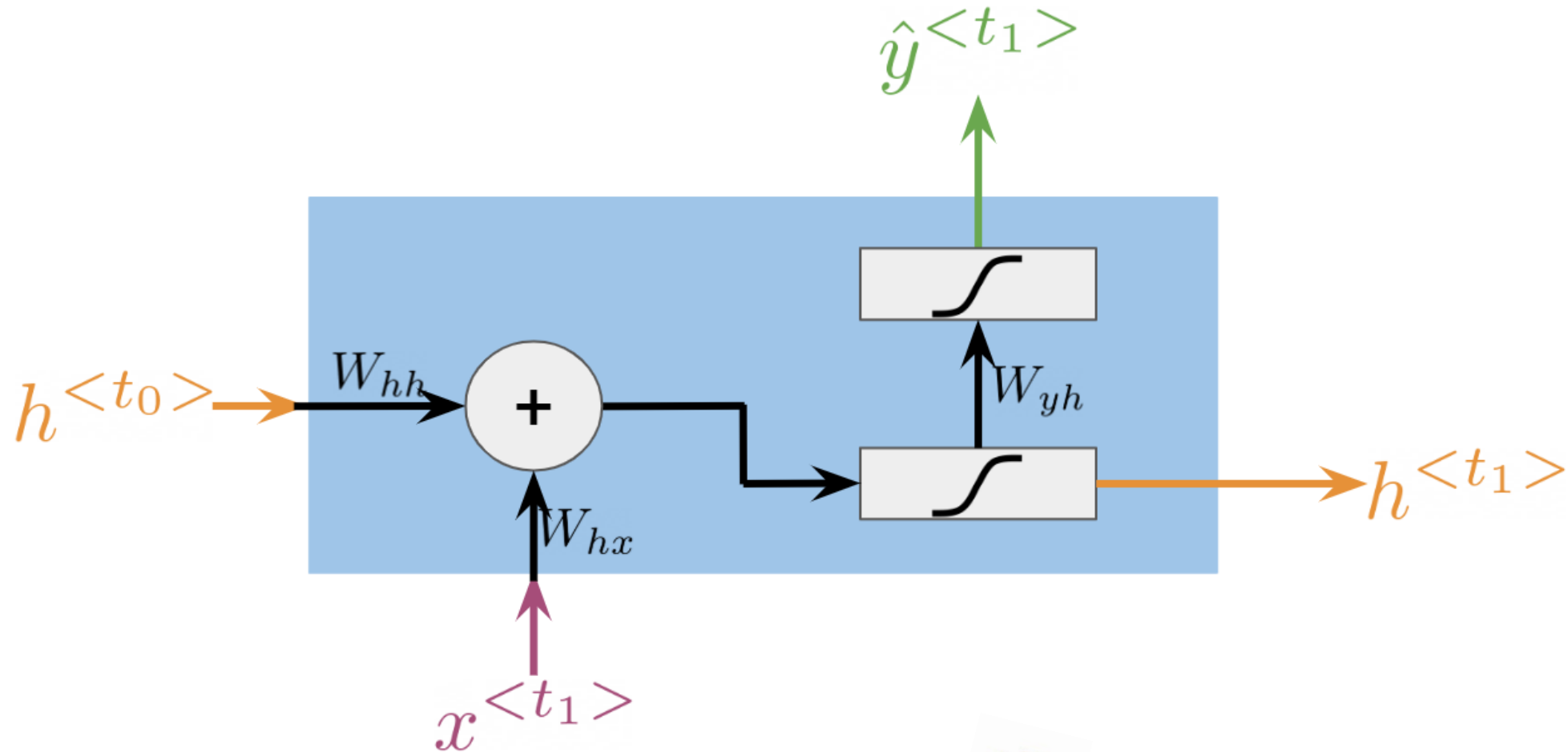
# FFNN LMs

Improvements over n-gram LMs:

- No sparsity problem
- Don't need to store all observed n-grams
- Embeddings can generalize and predict unseen words

Remaining problems:

- Fixed window is **too small**
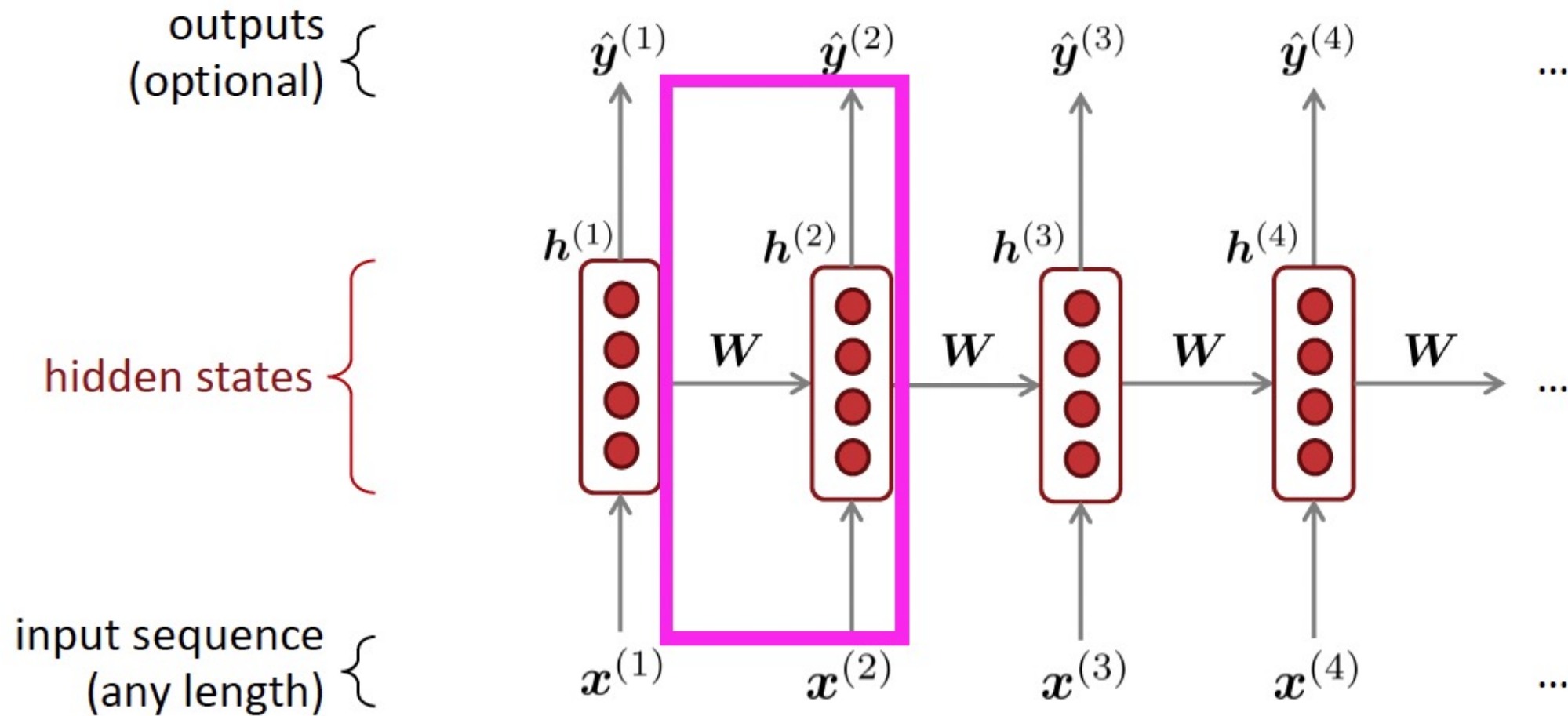- Enlarging window enlarges $W$

# Recurrent Neural Networks (RNN)


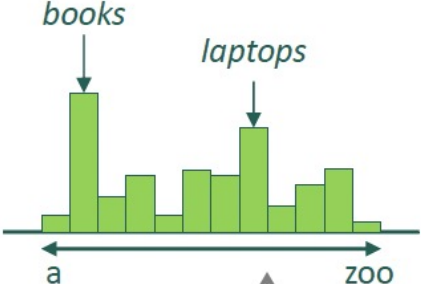
**simple/vanilla/Elman RNN**

# RNN basic structure



© Jixing Li

# RNN language model

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

output distribution

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

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

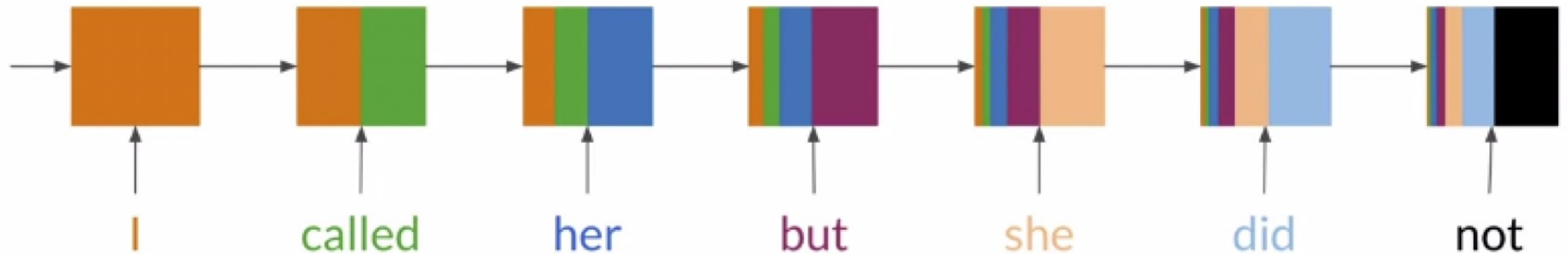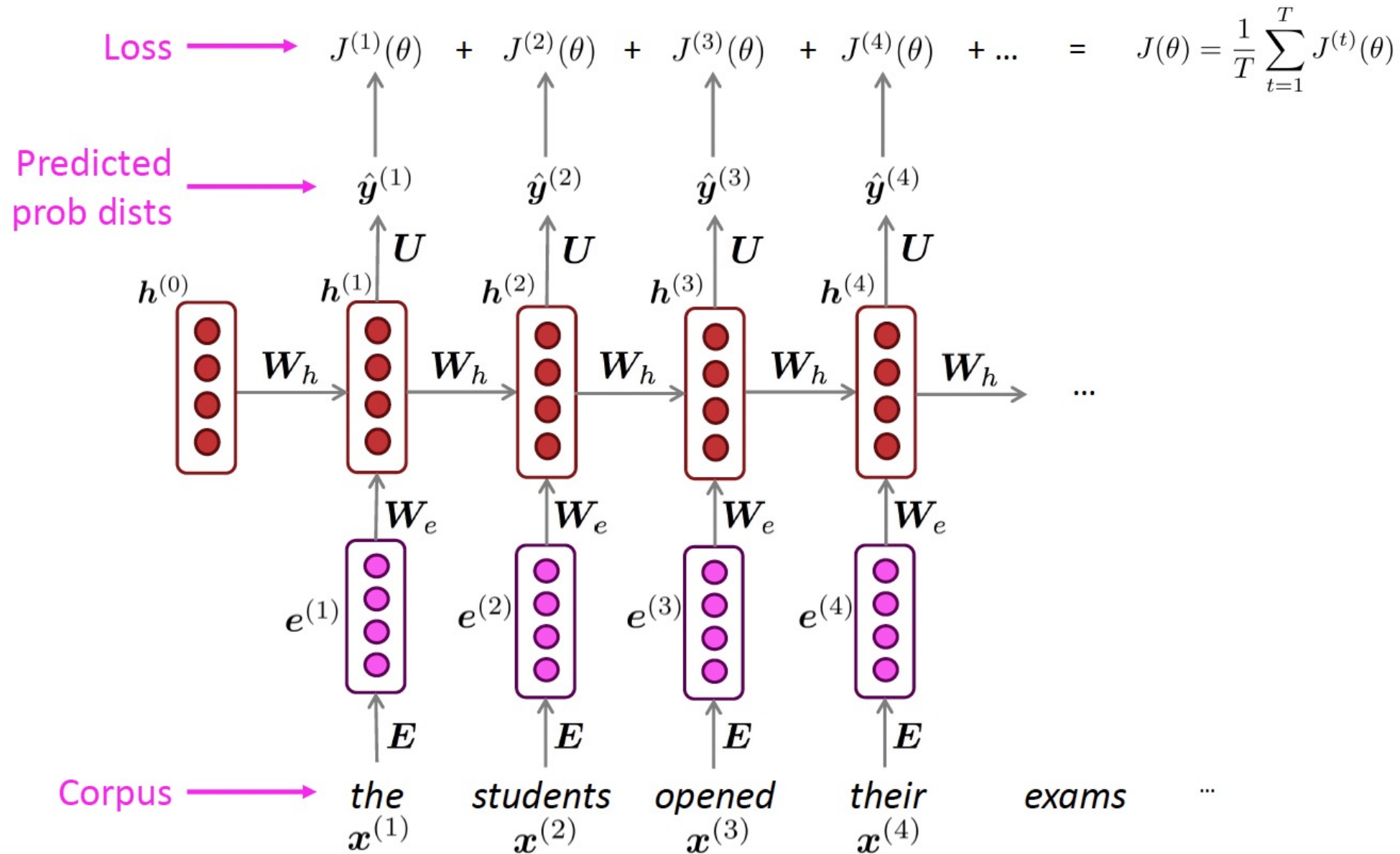$\boldsymbol{x}^{(t)} \in \mathbb{R}^{|V|}$
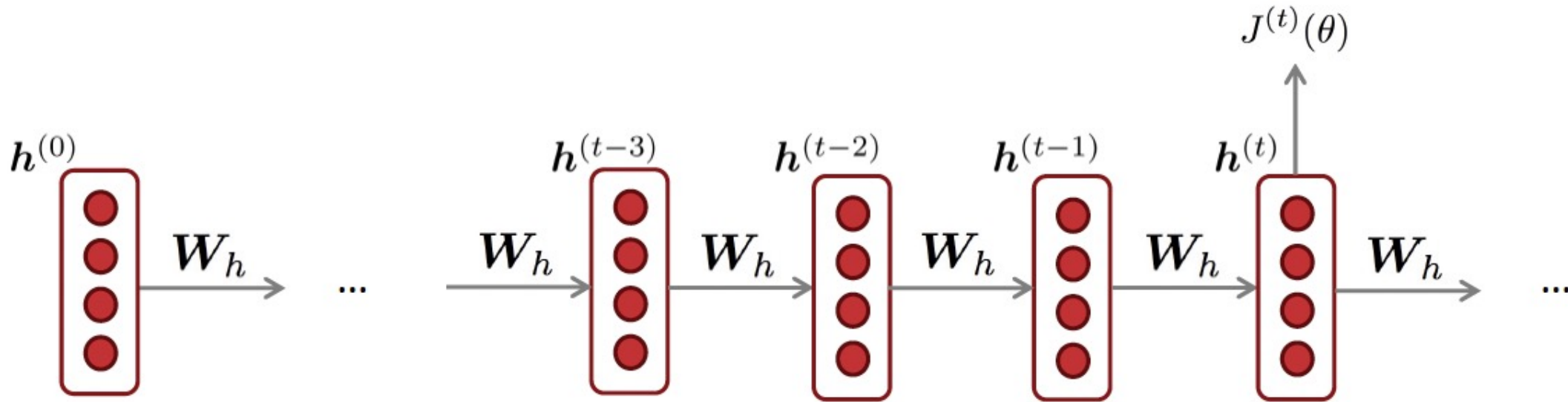
# RNN language model

**Advantages:**

• Can process any length input

• Computation for step $t$ can (in theory) use information from many steps back

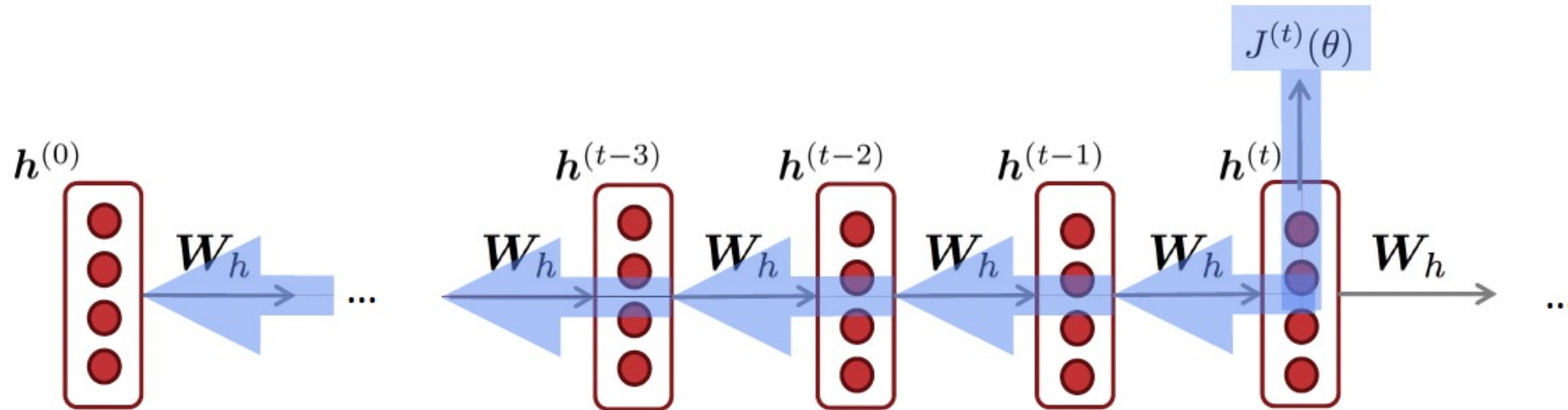• Model size doesn't increase for longer input context: Same weights applied on every timestep



I    called    her    but    she    did    not

# Training an RNN LM

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

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 ...

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

© Jixing Li

# Backpropagation for RNNs



The derivative of $J^{(t)}(\theta)$ w.r.t. the repeated weight matrix $\boldsymbol{W}_h$ is the sum of the gradient w.r.t. each time it appears

$$\frac{\partial J^{(t)}}{\partial \boldsymbol{W}_h} = \sum_{i=1}^{t} \frac{\partial J^{(t)}}{\partial \boldsymbol{W}_h}\bigg|_{(i)}$$
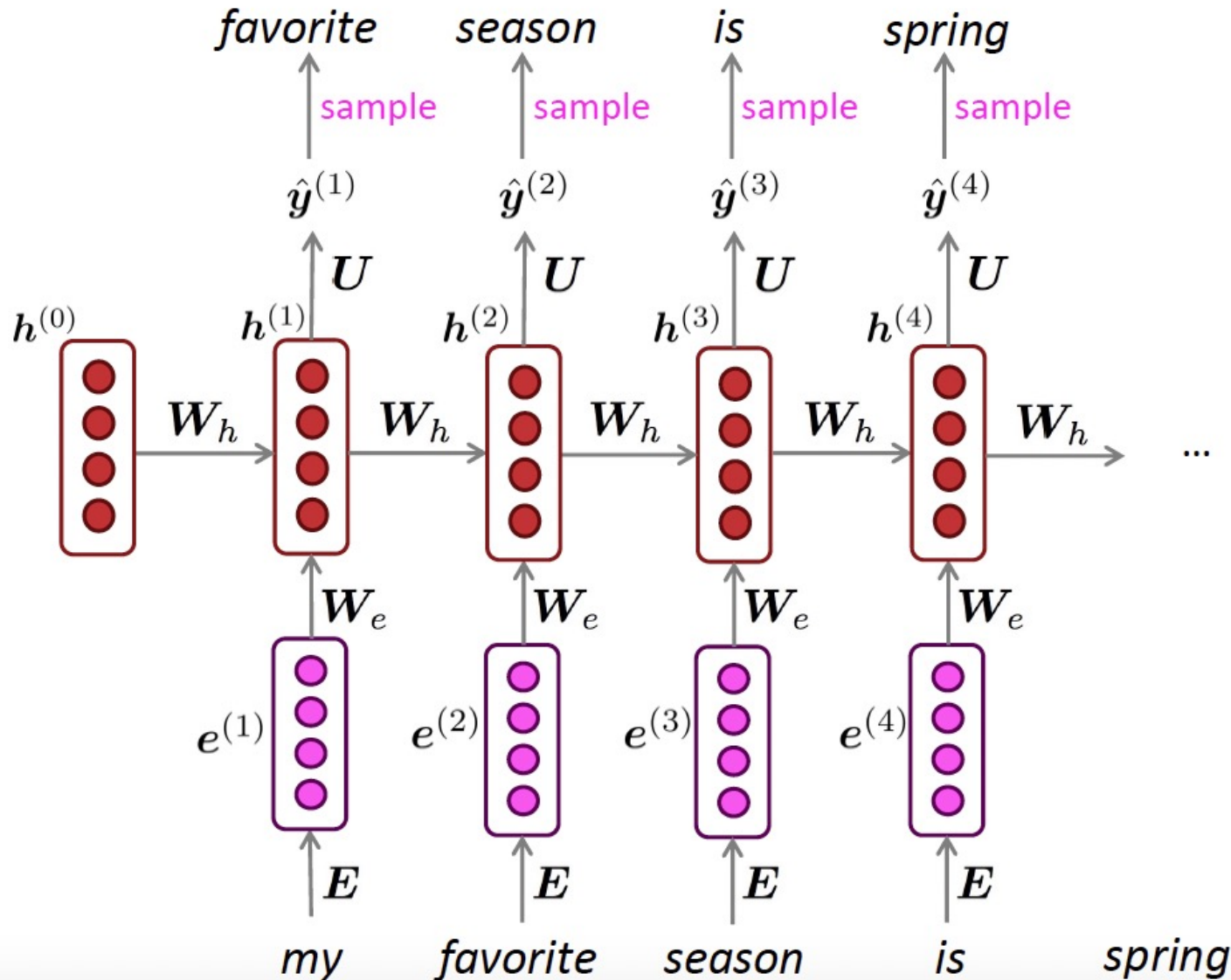
© Jixing Li

# Backpropagation through time



Backpropagate over timesteps $i=t,...,0$, summing gradients as you go → **"backpropagation through time"**
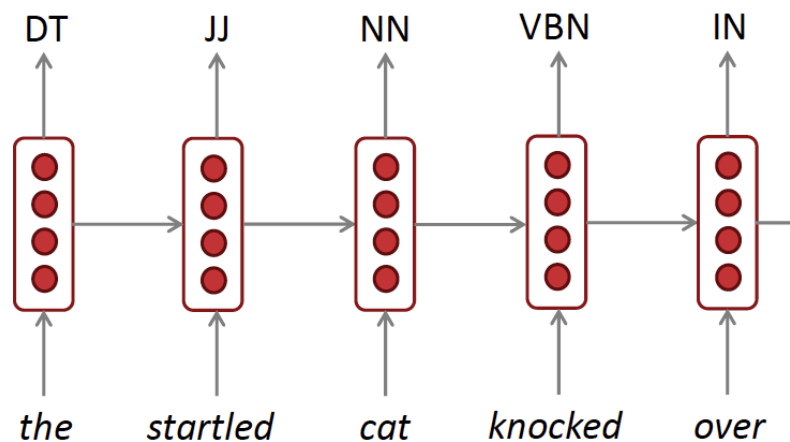
In practice, often "truncated" after ~20 timesteps for training efficiency reasons
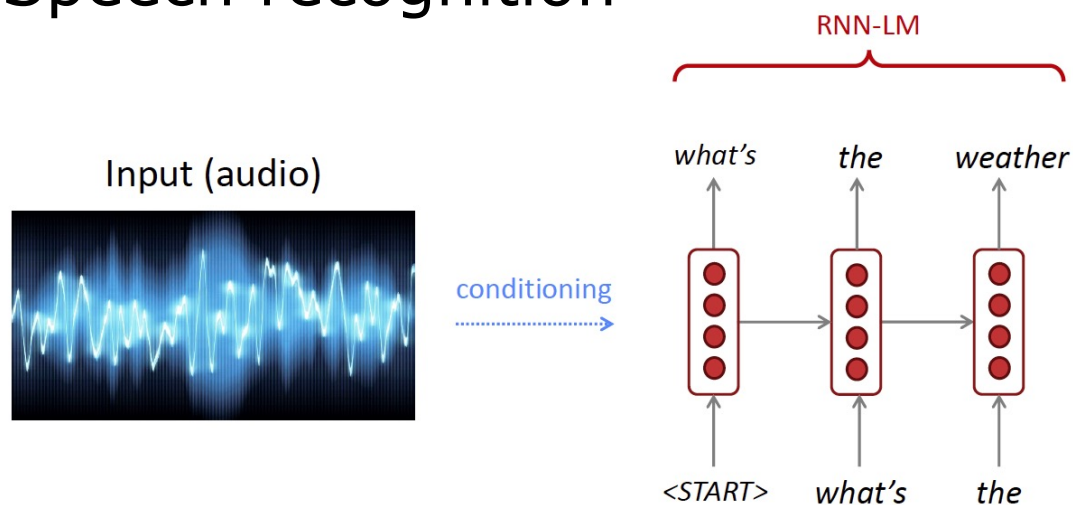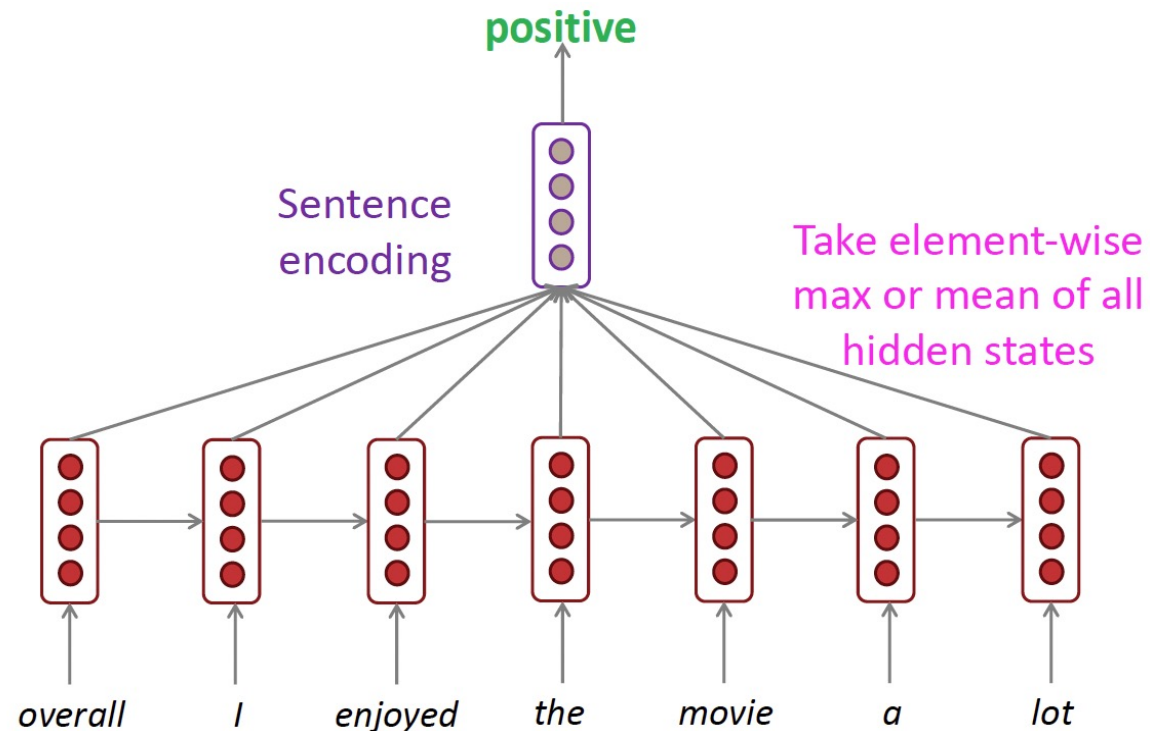
# Generating text with an RNN LM

# RNN Applications

## POS tagging



## Speech recognition



## Sentiment analysis



© Jixing Li

# To do

- Optional reading: **SLP** Ch9
- Do HW9