

Sequence Modeling: Recurrent and Recursive neural Nets

KISTI-UST
JUYEON YU

Contents

1. Introduction
2. Recurrent Neural Networks (RNN)
3. RNN Language Model
 - Backpropagation
 - Multivariable Chain Rule
4. Example of RNN Language Model
5. Understanding what's going on
6. Conclusion

01 | Introduction

A fixed-window neural Language Model

output distribution

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

hidden layer

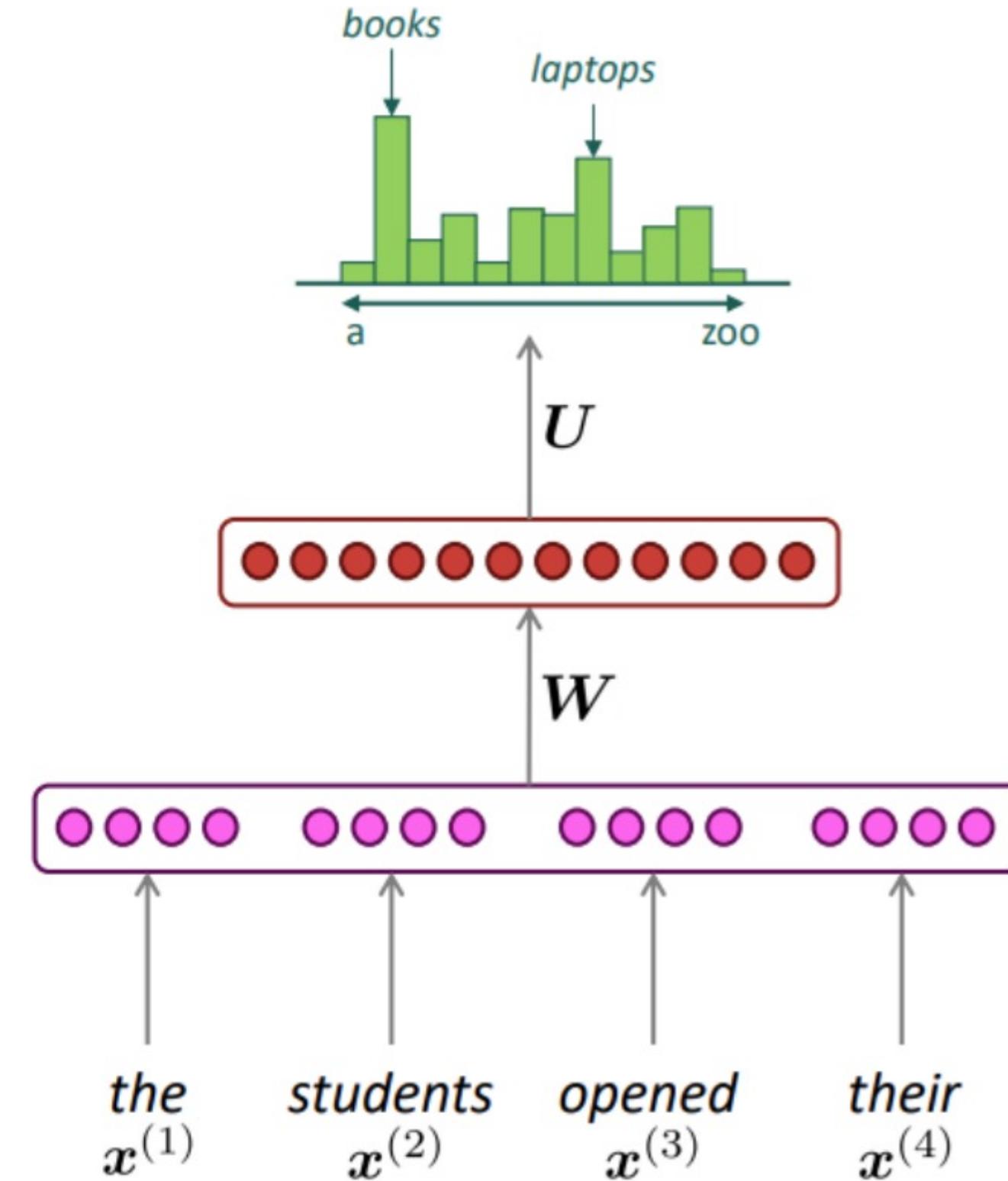
$$\mathbf{h} = f(\mathbf{W}\mathbf{e} + \mathbf{b}_1)$$

concatenated word embeddings

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

words / one-hot vectors

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



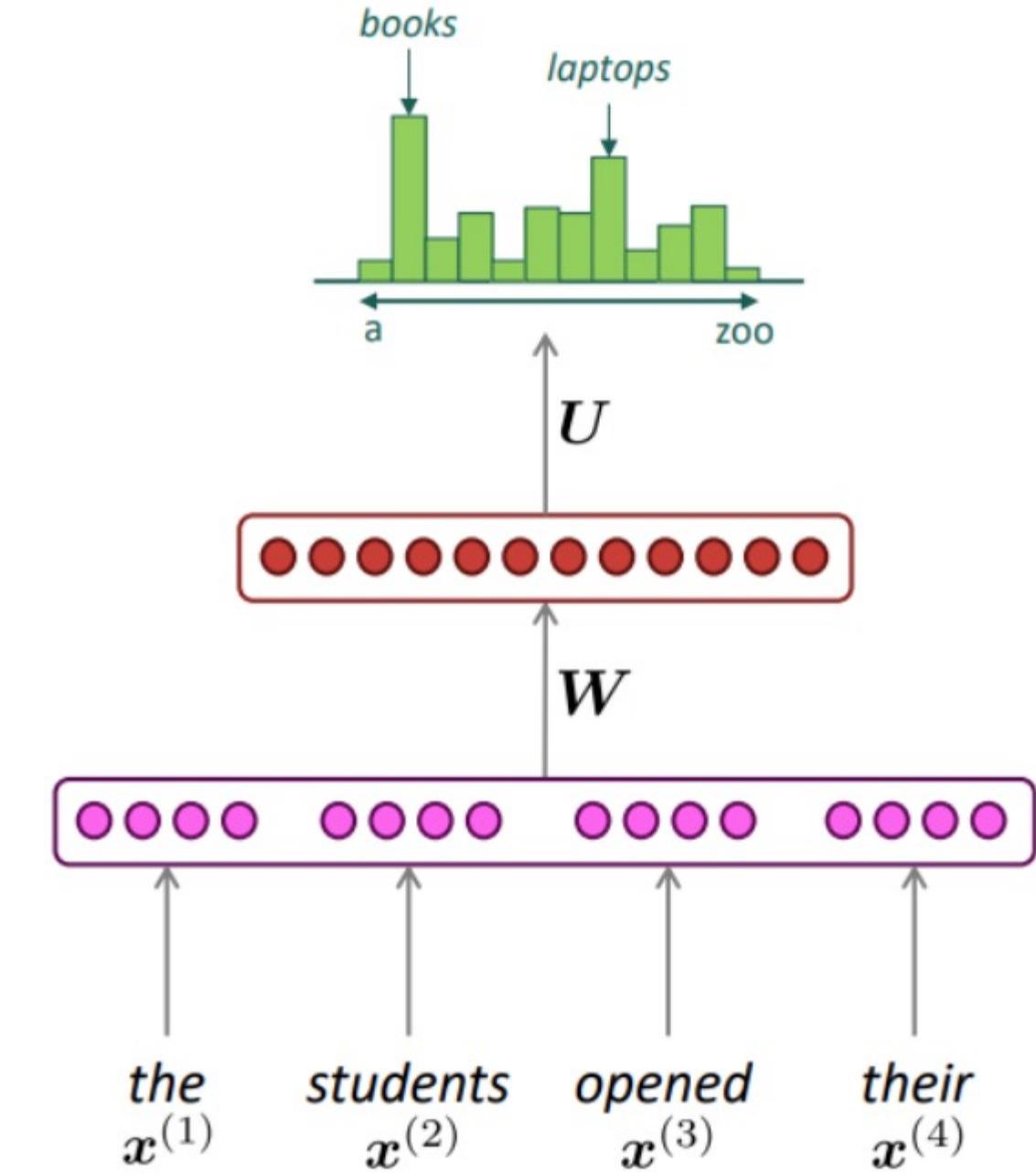
01 | Introduction

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

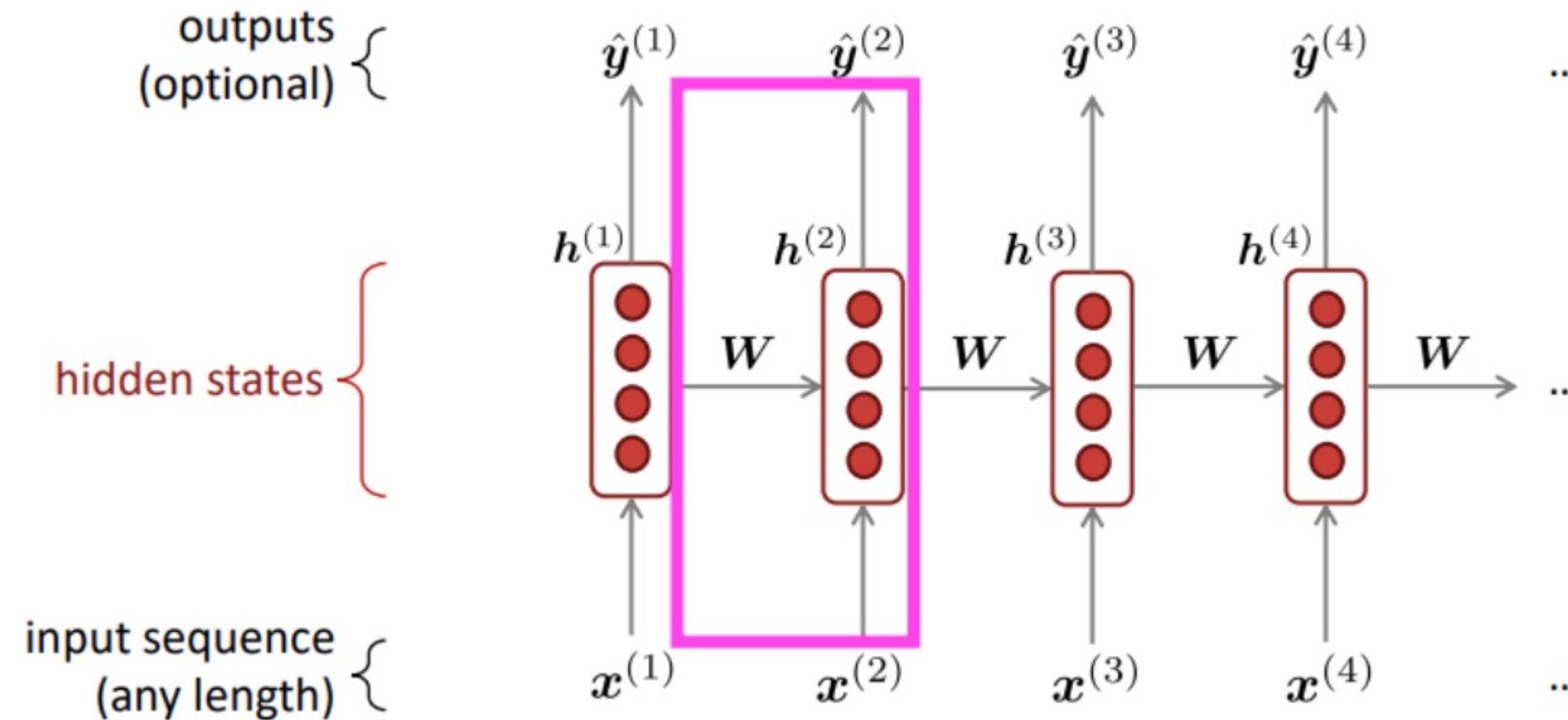


Y. Bengio, et al. (2000/2003): A Neural Probabilistic Language Model

02 | Recurrent Neural Networks

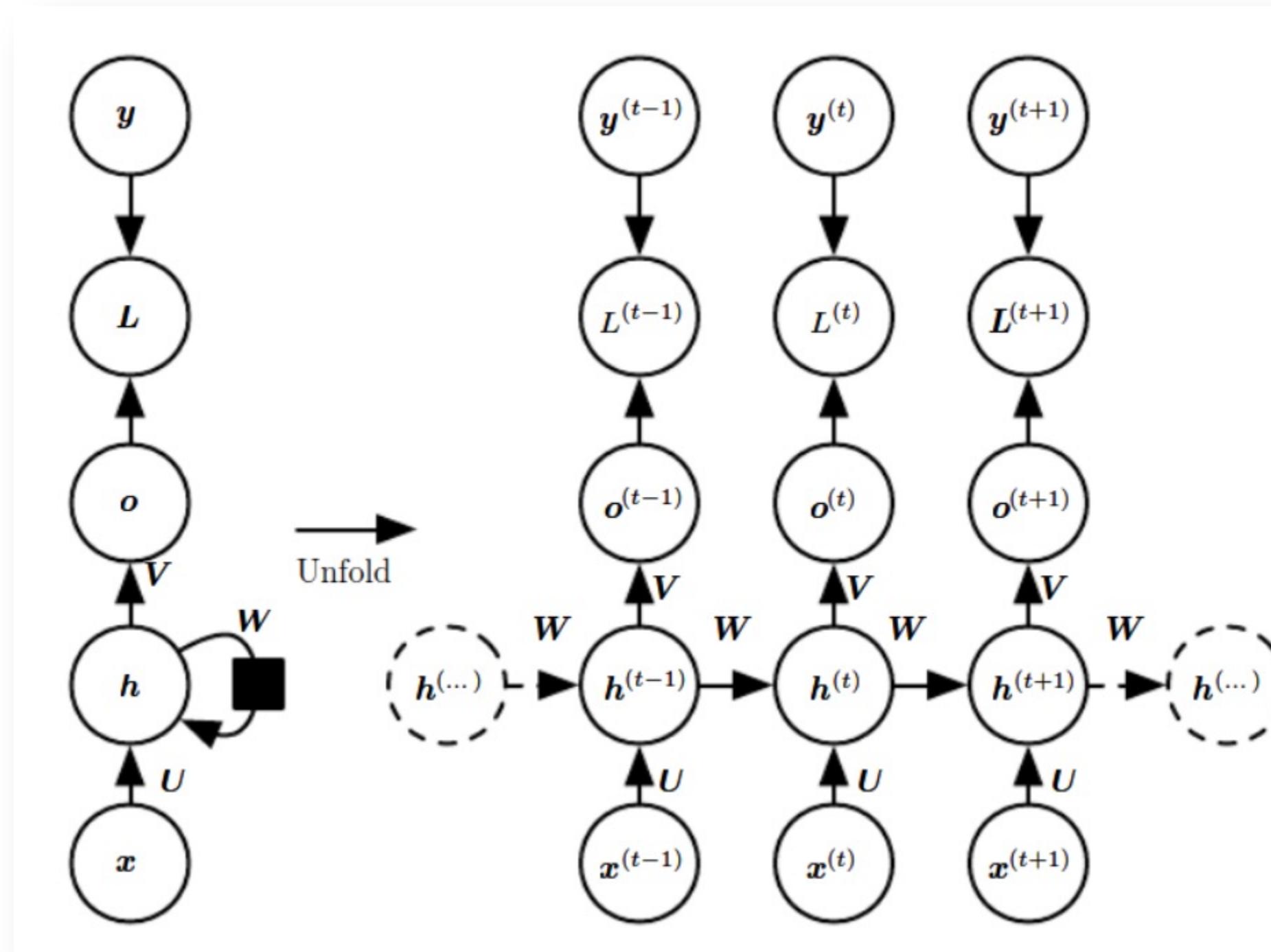
A family of neural architectures

- Core idea: Apply the same weights W repeatedly



02 | Recurrent Neural Networks

The computational graph



$$\mathbf{h}^{(t)} = f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}; \boldsymbol{\theta}), \quad (10.5)$$

$$\mathbf{h}^{(t)} = g^{(t)}(\mathbf{x}^{(t)}, \mathbf{x}^{(t-1)}, \mathbf{x}^{(t-2)}, \dots, \mathbf{x}^{(2)}, \mathbf{x}^{(1)}) \quad (10.6)$$

$$= f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}; \boldsymbol{\theta}). \quad (10.7)$$

$$\mathbf{a}^{(t)} = \mathbf{b} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)}, \quad (10.8)$$

$$\mathbf{h}^{(t)} = \tanh(\mathbf{a}^{(t)}), \quad (10.9)$$

$$\mathbf{o}^{(t)} = \mathbf{c} + \mathbf{V}\mathbf{h}^{(t)}, \quad (10.10)$$

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{o}^{(t)}), \quad (10.11)$$

- $U(Wxh)$: $p \times d$ input-hidden matrix
- $W(Whh)$: $p \times p$ hidden-hidden matrix
- $V(Why)$: $d \times p$ hidden-output matrix
- both $x(t)$ and $y(t)$ are d -dimensional for a lexicon(vocabulary) of size d .
- p regulates the complexity of the embedding.

02 | Recurrent Neural Networks

Advantages and Disadvantages

- 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

02 | Recurrent Neural Networks

Computation

```
rnn = RNN()  
y = rnn.step(x) # x is an input vector, y is the RNN's output vector
```

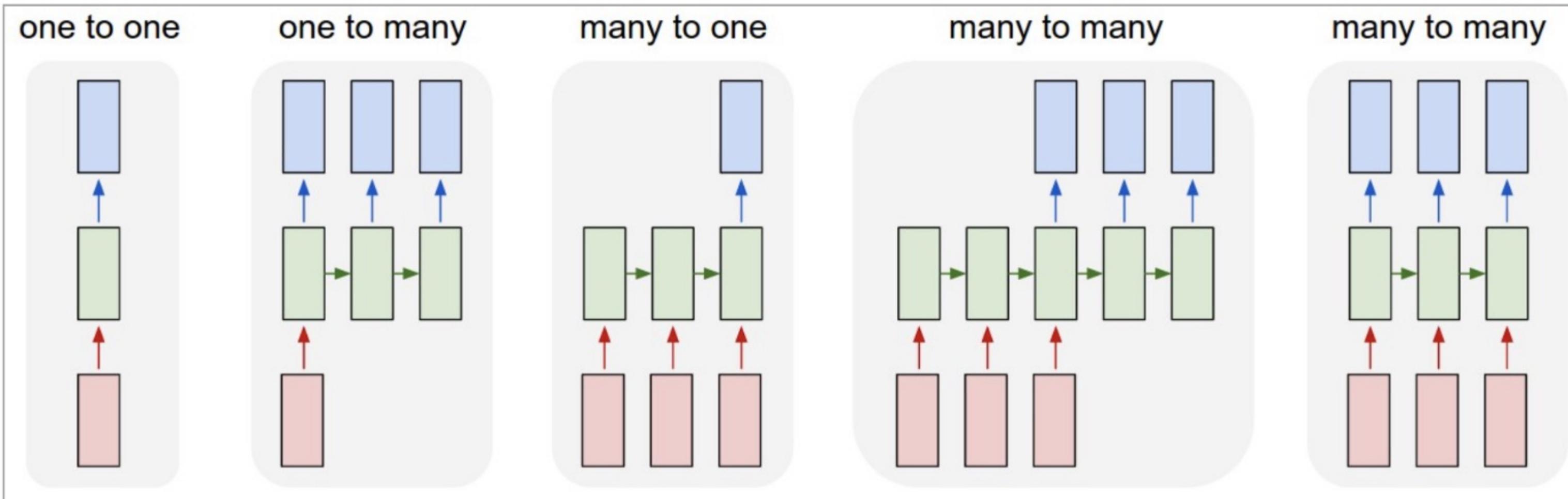
```
class RNN:  
    # ...  
    def step(self, x):  
        # update the hidden state  
        self.h = np.tanh(np.dot(self.W_hh, self.h) + np.dot(self.W_xh, x))  
        # compute the output vector  
        y = np.dot(self.W_hy, self.h)  
        return y
```

Going deep

```
y1 = rnn1.step(x)  
y = rnn2.step(y1)
```

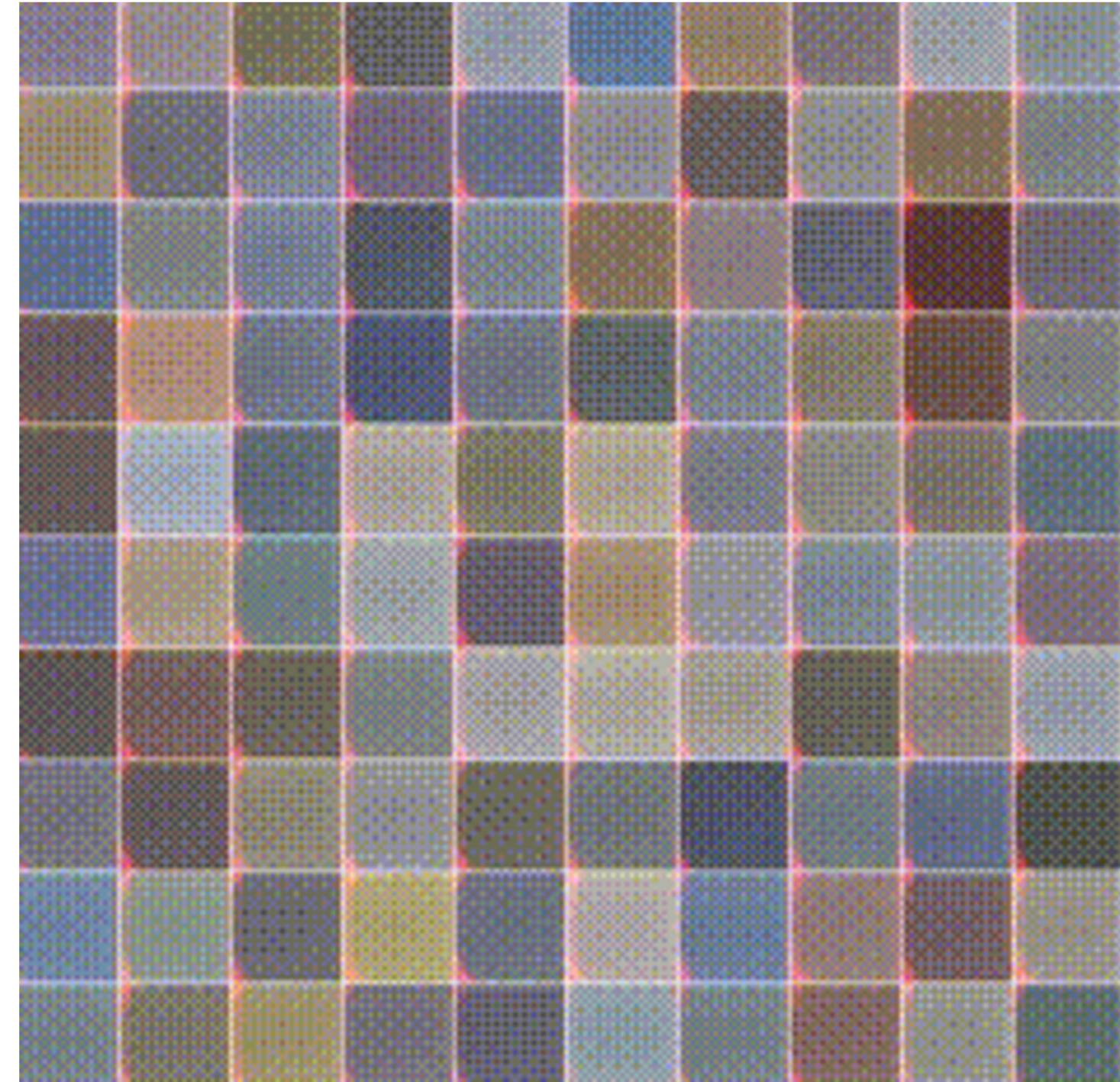
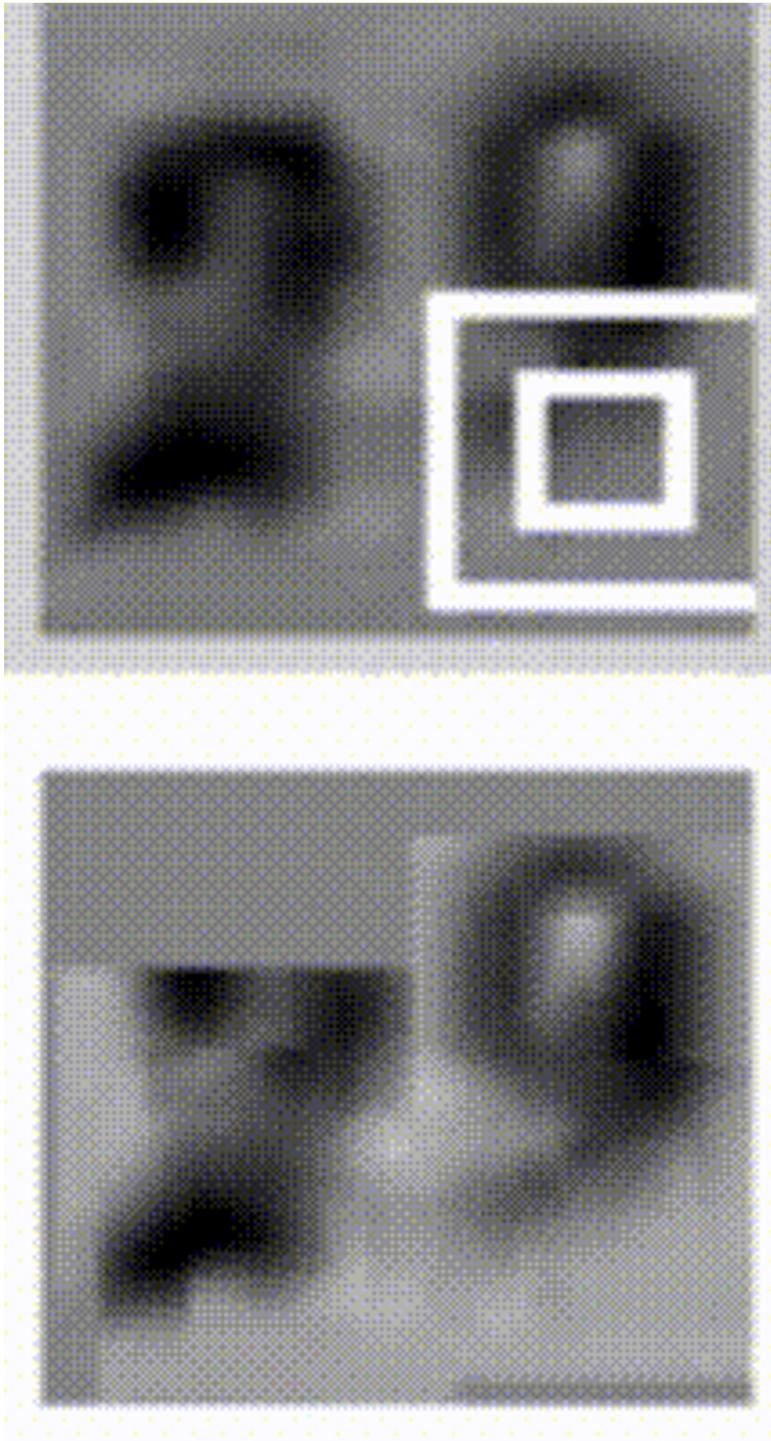
02 | Recurrent Neural Networks

The different variations of RNN with missing inputs and outputs



02 | Recurrent Neural Networks

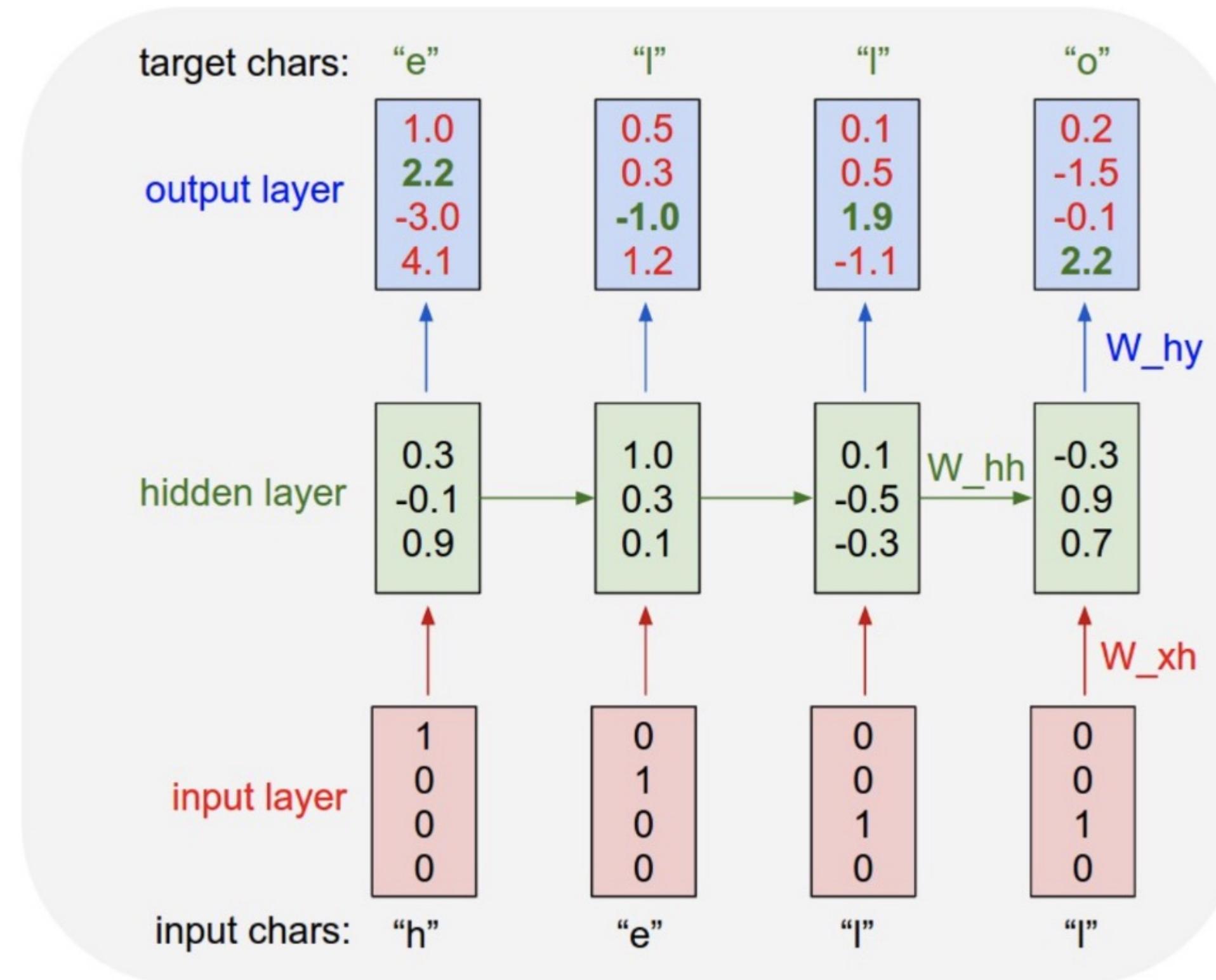
Sequential processing in absence of sequences.



Left: RNN learns to read house numbers. Right: RNN learns to paint house numbers.

03 | RNN Language Model

Character-Level Language Models



03 | RNN Language Model

Word-Level Language Models

output distribution

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

hidden states

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

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

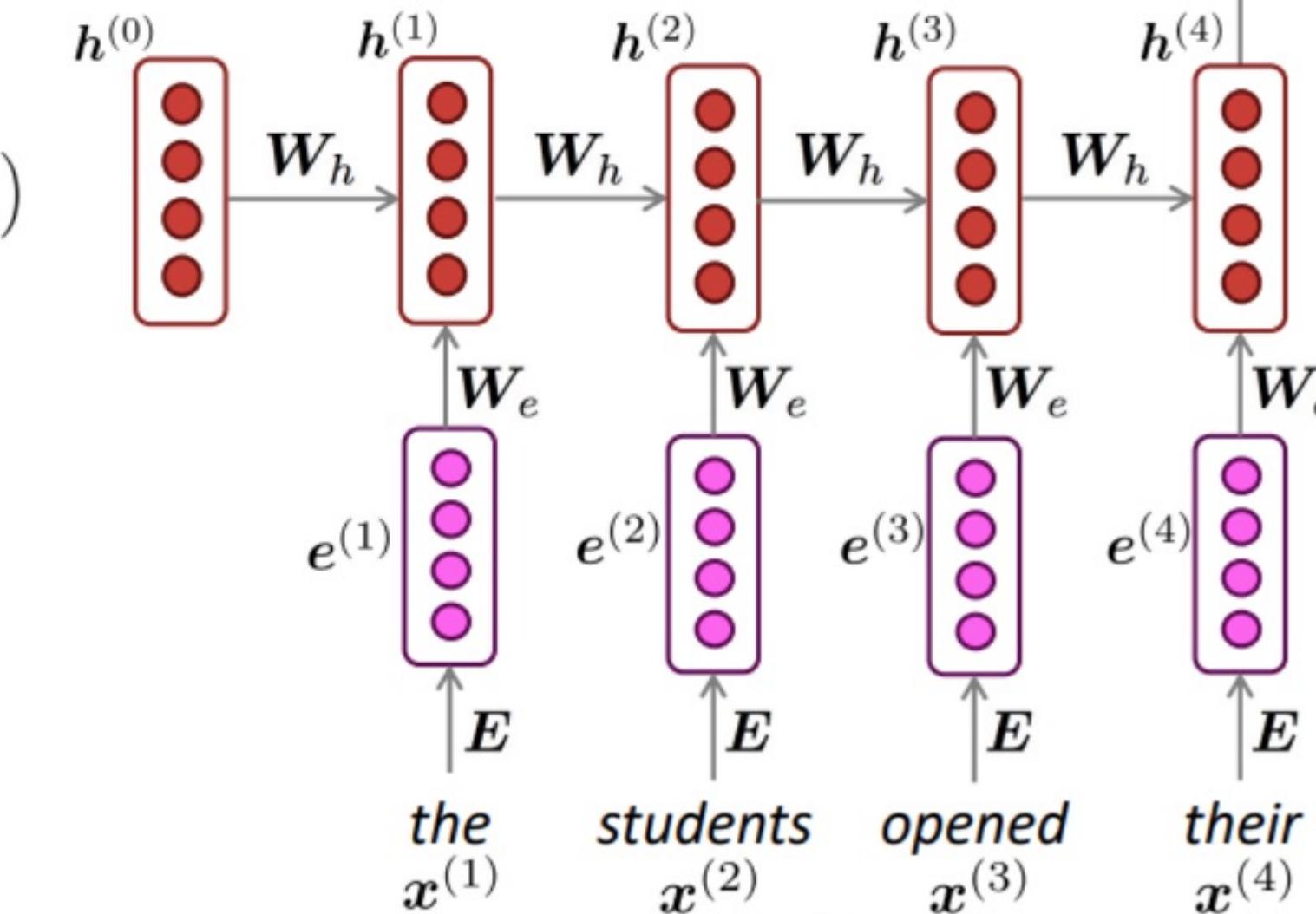
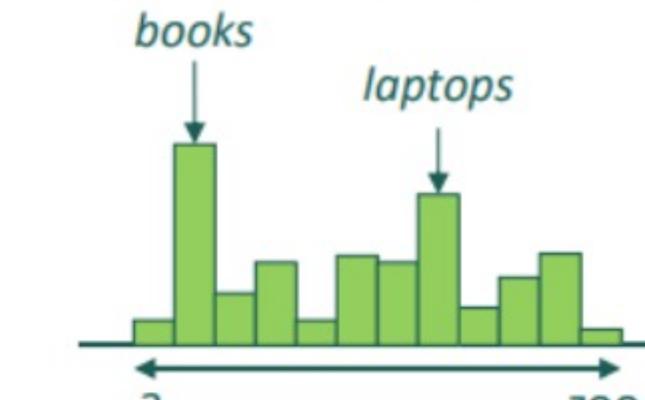
word embeddings

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

words / one-hot vectors

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

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



Note: this input sequence could be much longer now!

03 | RNN Language Model

Training an RNN Language Model

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

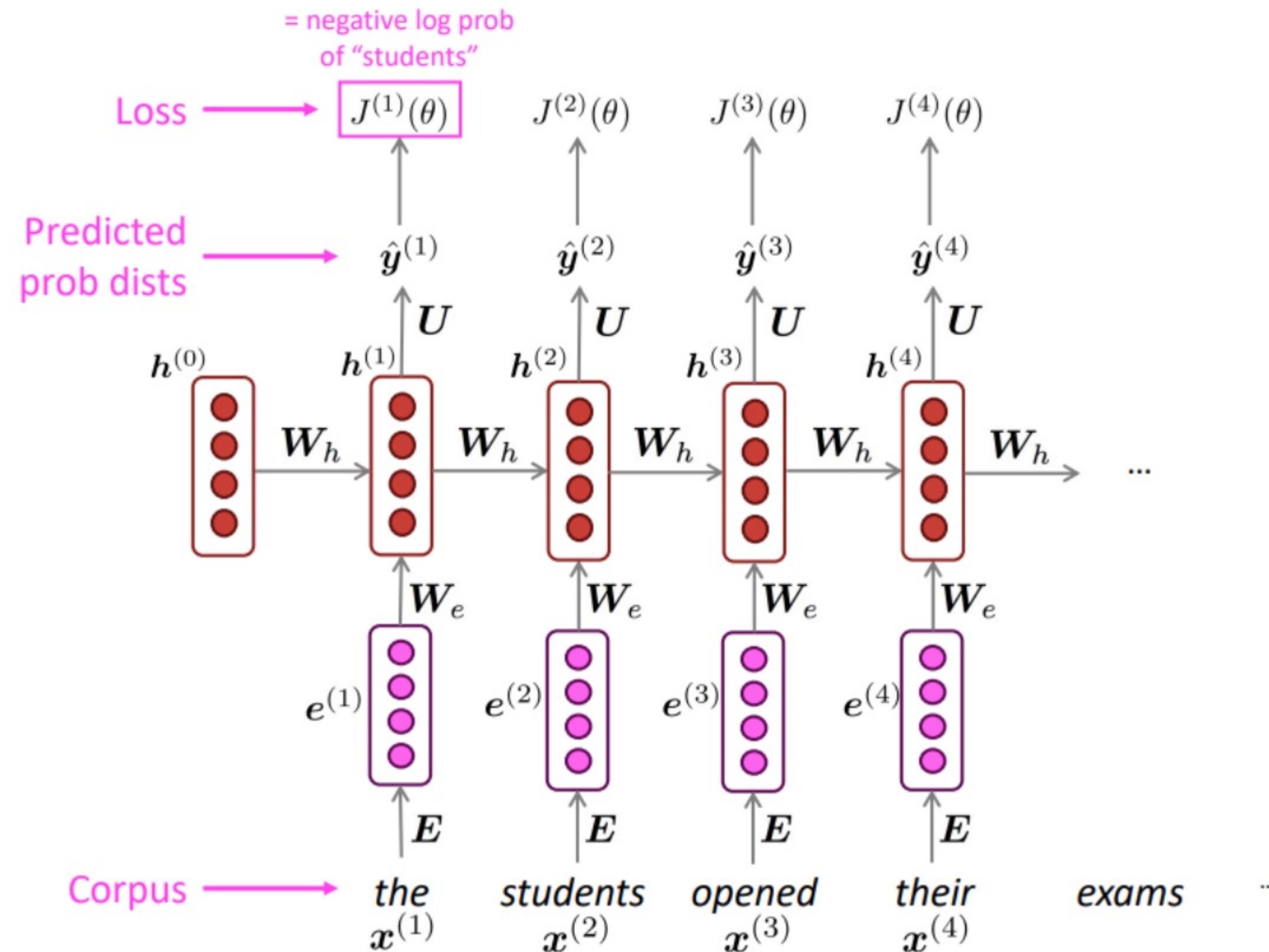
$$J^{(t)}(\theta) = CE(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = - \sum_{w \in V} \mathbf{y}_w^{(t)} \log \hat{\mathbf{y}}_w^{(t)} = - \log \hat{\mathbf{y}}_{\mathbf{x}_{t+1}}^{(t)}$$

- Average this to get **overall loss** for entire training set:

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

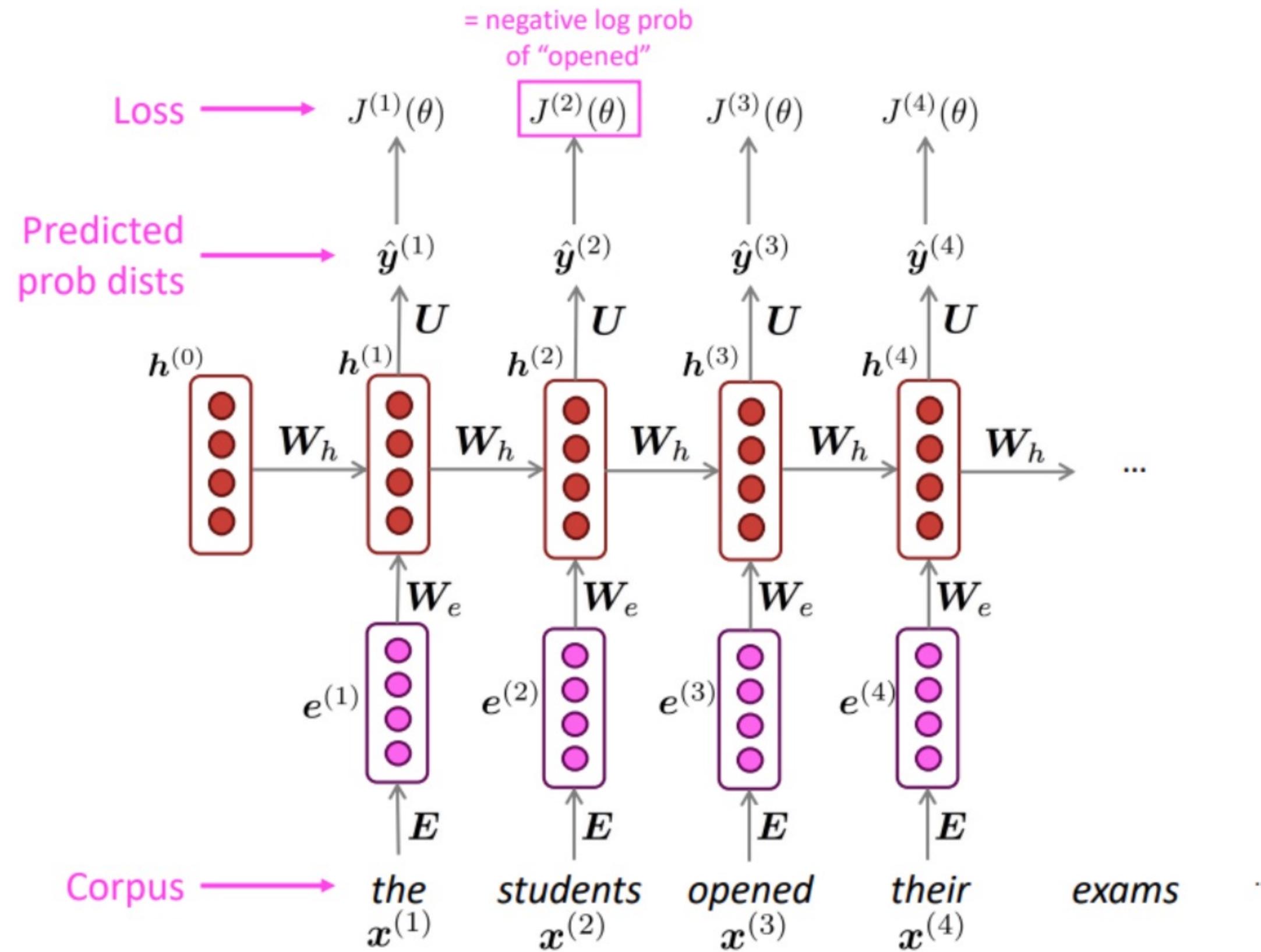
03 | RNN Language Model

Training an RNN Language Model



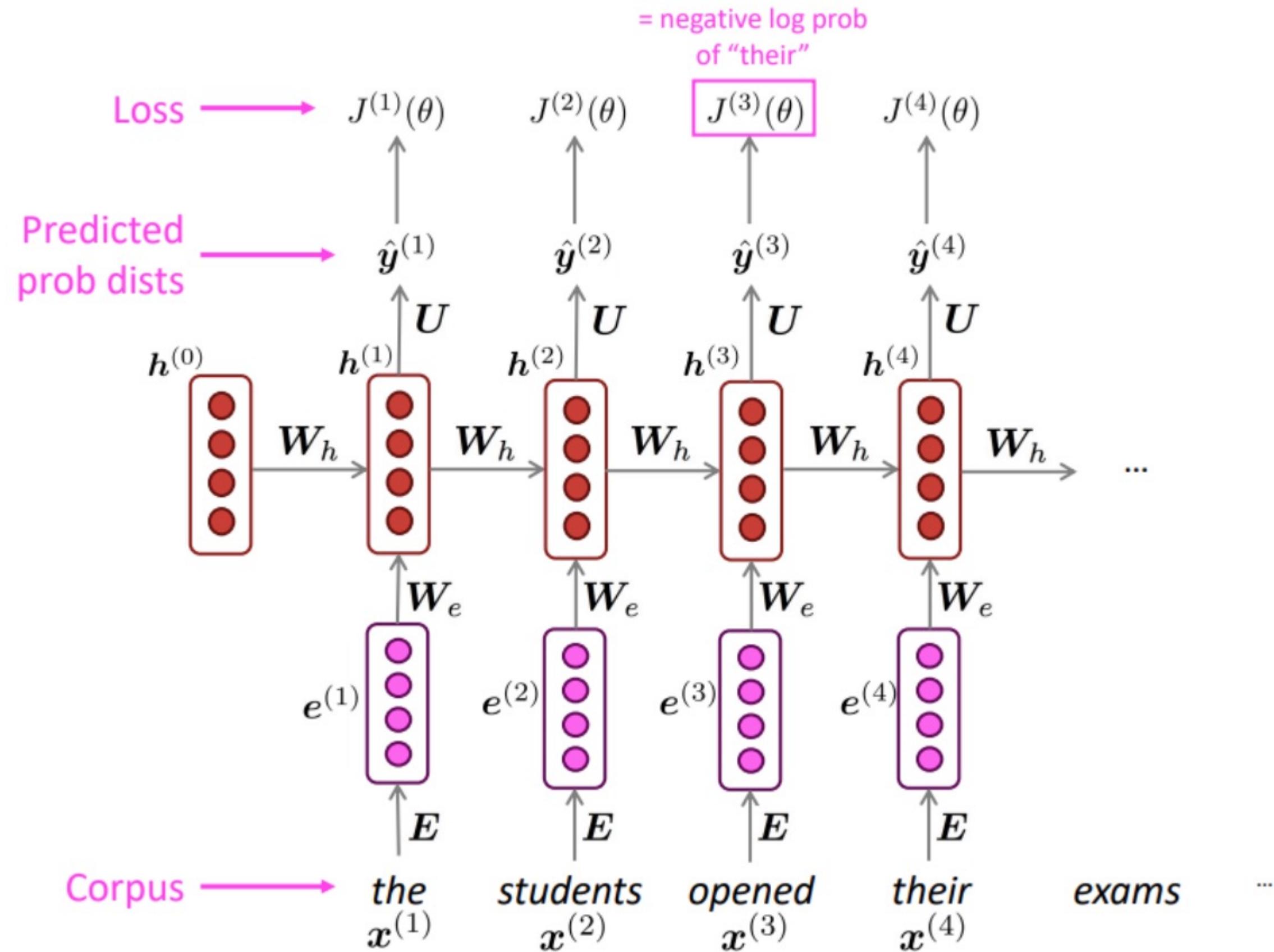
03 | RNN Language Model

Training an RNN Language Model



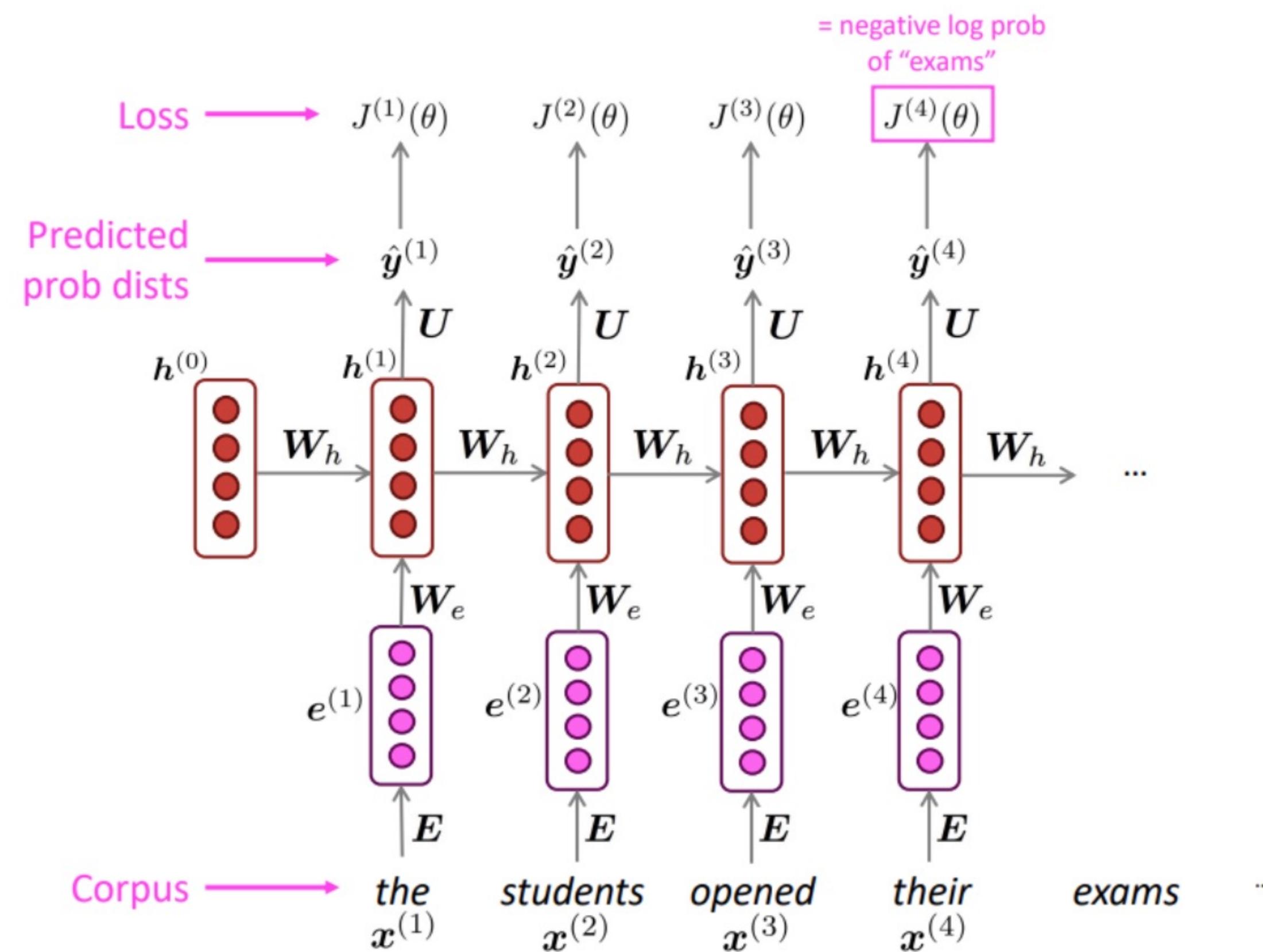
03 | RNN Language Model

Training an RNN Language Model



03 | RNN Language Model

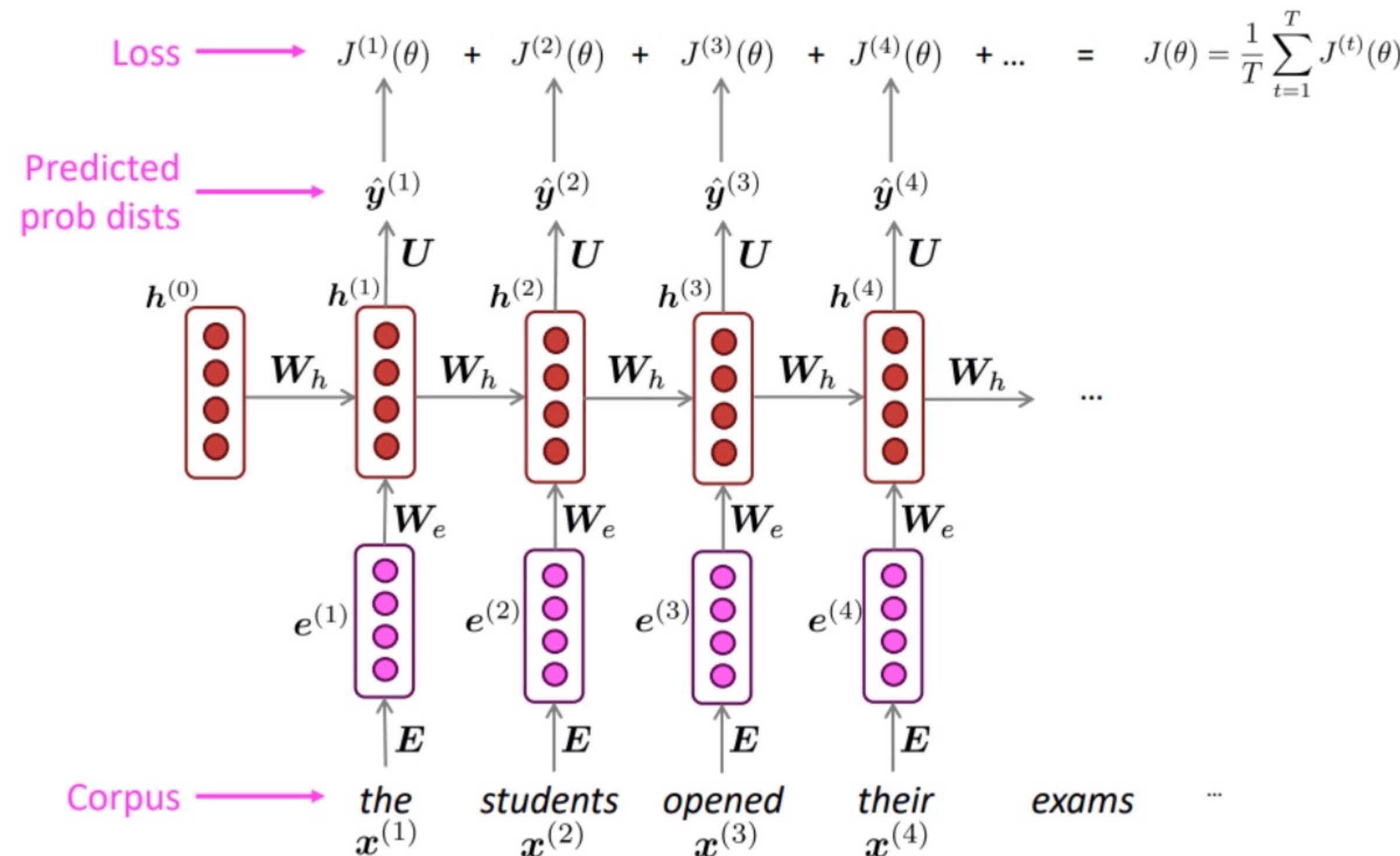
Training an RNN Language Model



03 | RNN Language Model

Training an RNN Language Model

“Teacher forcing”



03 | RNN Language Model

Training an RNN Language Model

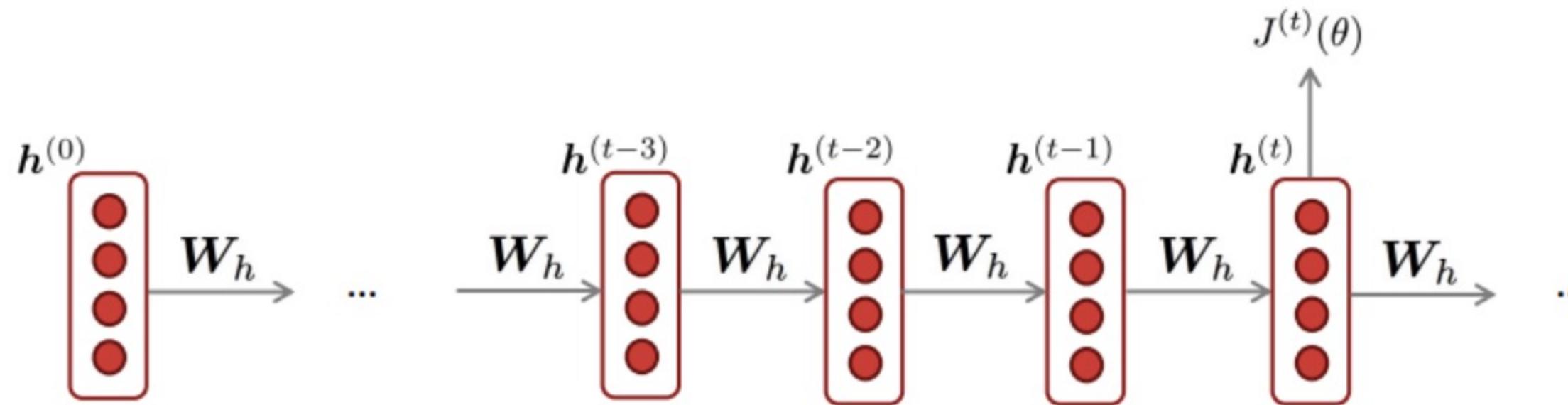
- However: Computing loss and gradients across **entire corpus** $x^{(1)}, \dots, x^{(T)}$ is **too expensive!**

$$J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta)$$

- In practice, consider $x^{(1)}, \dots, x^{(T)}$ as a **sentence (or a document)**
- Recall: **Stochastic Gradient Descent** allows us to compute loss and gradients for small chunk of data, and update.
- Compute loss $J(\theta)$ for a sentence (actually, a batch of sentences), compute gradients and update weights. Repeat.

03 | RNN Language Model

Backpropagation



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

Answer:
$$\frac{\partial J^{(t)}}{\partial \mathbf{W}_h} = \sum_{i=1}^t \frac{\partial J^{(t)}}{\partial \mathbf{W}_h} \Big|_{(i)}$$

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

Why?

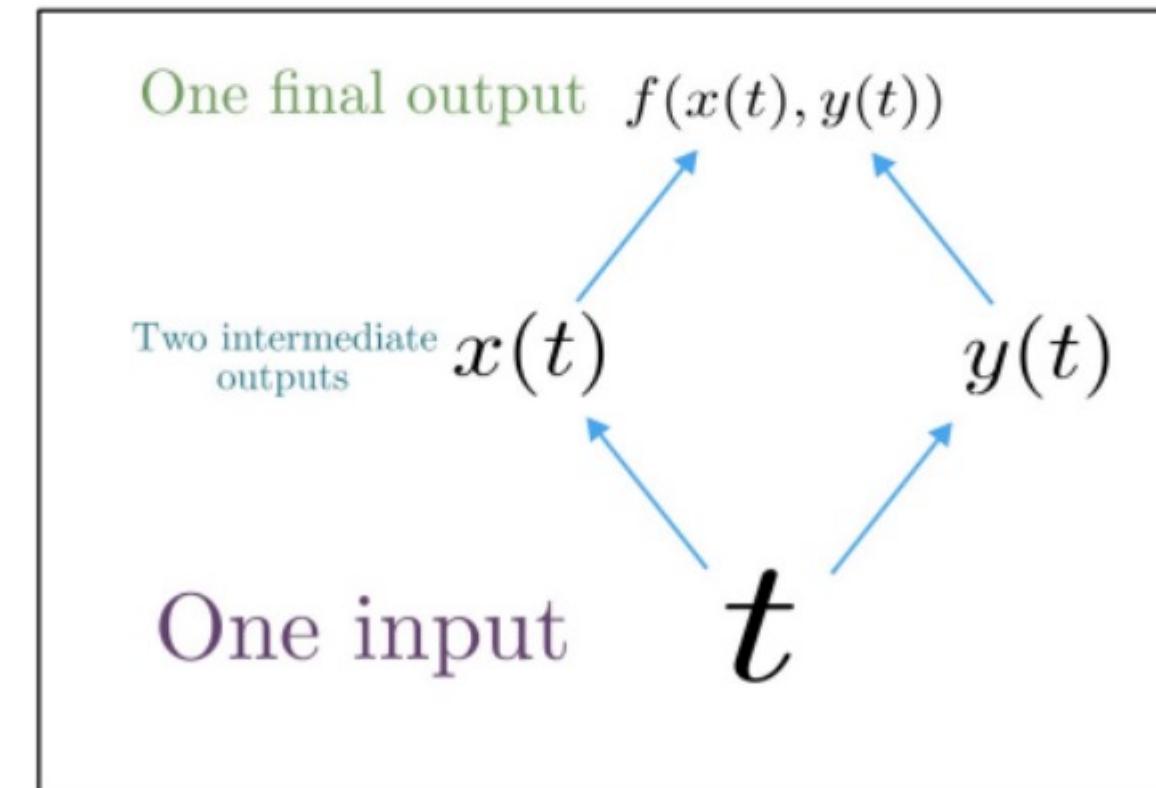
03 | RNN Language Model

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:

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

Derivative of composition function



Source:

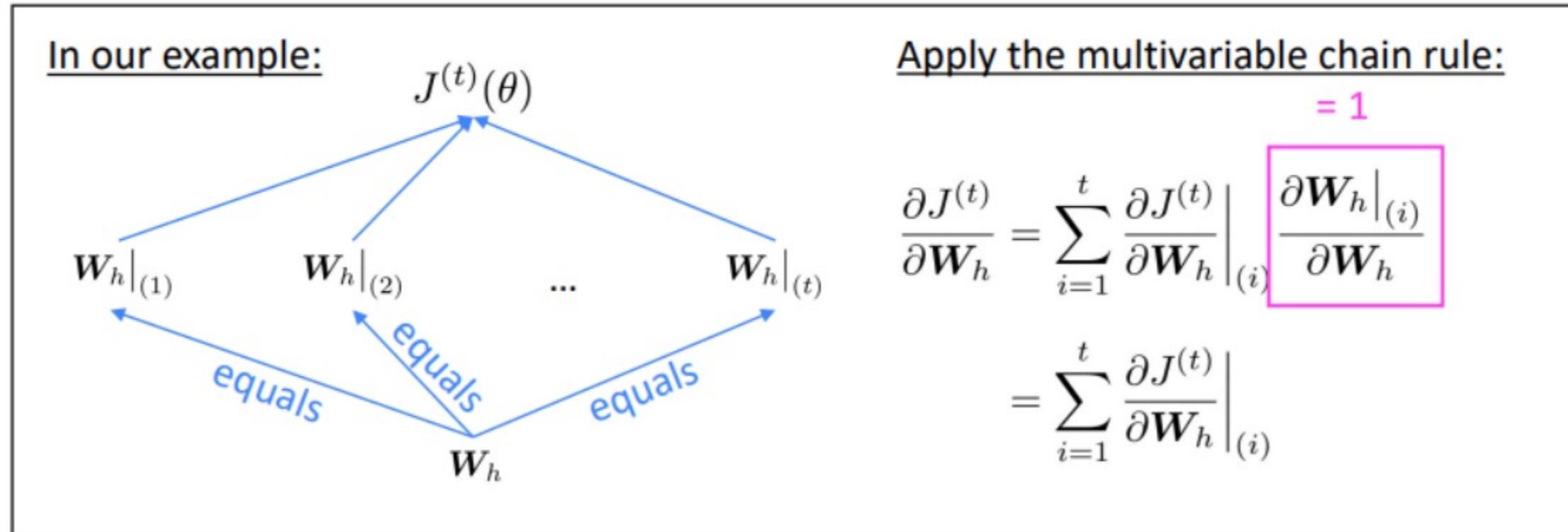
<https://www.khanacademy.org/math/multivariable-calculus/multivariable-derivatives/differentiating-vector-valued-functions/a/multivariable-chain-rule-simple-version>

03 | RNN Language Model

Backpropagation: Proof sketch

- 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:

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

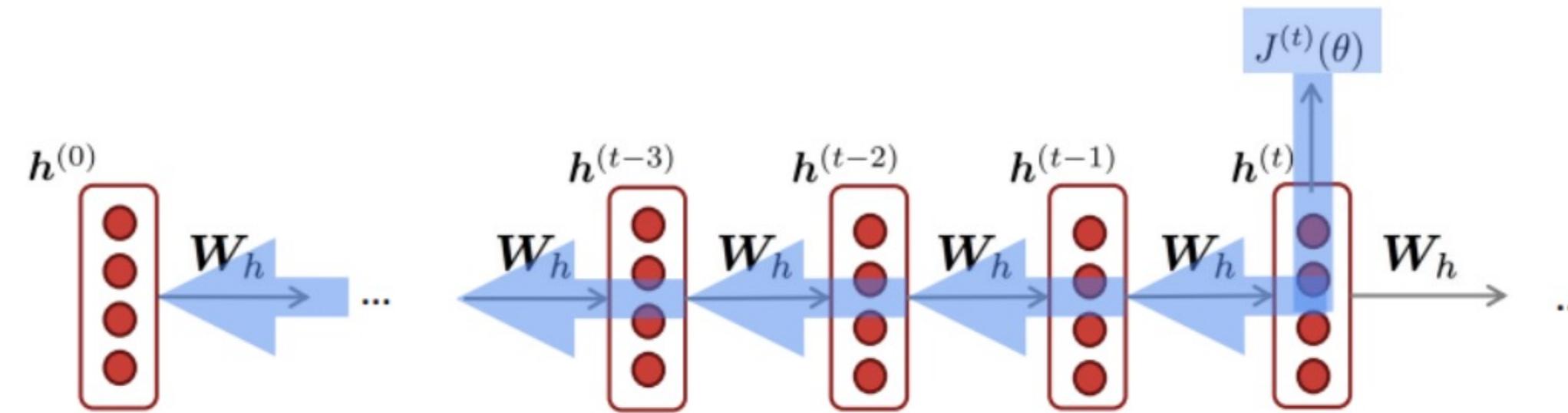


Source:

<https://www.khanacademy.org/math/multivariable-calculus/multivariable-derivatives/differentiating-vector-valued-functions/a/multivariable-chain-rule-simple-version>

03 | RNN Language Model

Backpropagation



$$\frac{\partial J^{(t)}}{\partial \mathbf{W}_h} = \left[\sum_{i=1}^t \frac{\partial J^{(t)}}{\partial \mathbf{W}_h} \right]_{(i)}$$

Question: How do we calculate this?

Answer: Backpropagate over timesteps $i=t, \dots, 0$, summing gradients as you go.
This algorithm is called “**backpropagation through time**” [Werbos, P.G., 1988, *Neural Networks 1*, and others]

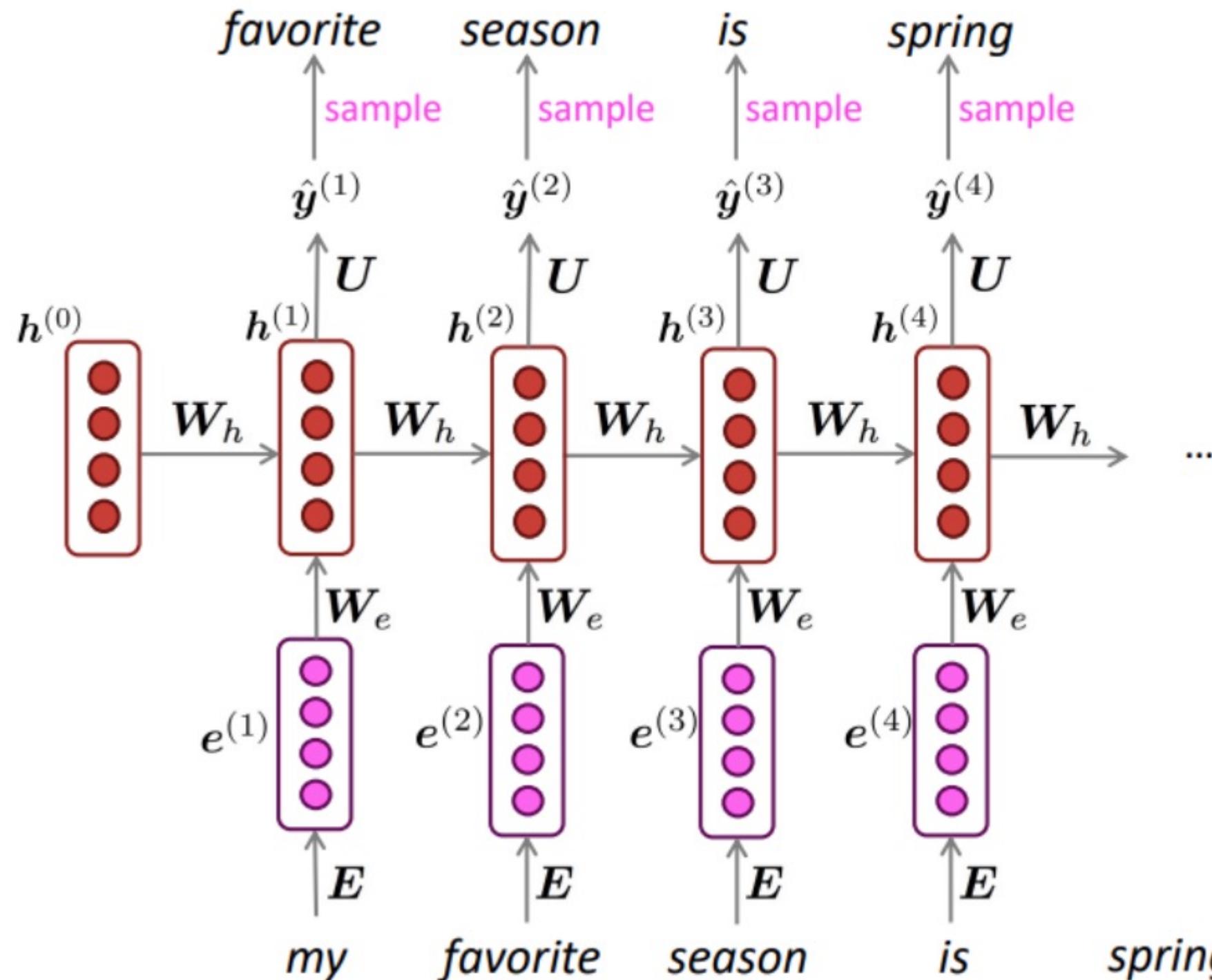
In practice, often “truncated” after ~ 20 timesteps for training efficiency reasons

04

Example of RNN Language Model

Generating text

Just like a n-gram Language Model, you can use a RNN Language Model to generate text by **repeated sampling**. Sampled output becomes next step's input.



- Let's have some fun!
 - You can train an RNN-LM on any kind of text, then generate text in that style.
 - RNN-LM trained on Obama speeches:



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.

04 | Example of RNN Language Model

Evaluating Language Models

- The standard **evaluation metric** for Language Models is **perplexity**.

$$\text{perplexity} = \prod_{t=1}^T \left(\frac{1}{P_{\text{LM}}(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})} \right)^{1/T}$$

Normalized by
number of words

- This is equal to the exponential of the cross-entropy loss $J(\theta)$:

$$= \prod_{t=1}^T \left(\frac{1}{\hat{\mathbf{y}}_{\mathbf{x}_{t+1}}^{(t)}} \right)^{1/T} = \exp \left(\frac{1}{T} \sum_{t=1}^T -\log \hat{\mathbf{y}}_{\mathbf{x}_{t+1}}^{(t)} \right) = \exp(J(\theta))$$

Lower perplexity is better!

04

Example of RNN Language Model

RNNs have greatly improved perplexity

n-gram model →

Increasingly complex RNNs

Model	Perplexity
Interpolated Kneser-Ney 5-gram (Chelba et al., 2013)	67.6
RNN-1024 + MaxEnt 9-gram (Chelba et al., 2013)	51.3
RNN-2048 + BlackOut sampling (Ji et al., 2015)	68.3
Sparse Non-negative Matrix factorization (Shazeer et al., 2015)	52.9
LSTM-2048 (Jozefowicz et al., 2016)	43.7
2-layer LSTM-8192 (Jozefowicz et al., 2016)	30
Ours small (LSTM-2048)	43.9
Ours large (2-layer LSTM-2048)	39.8

Perplexity improves
(lower is better)

Source: <https://research.fb.com/building-an-efficient-neural-language-model-over-a-billion-words/>

Why should we care about Language Modeling?

- Language Modeling is a **benchmark task** that helps us **measure our progress** on understanding language
- Language Modeling is a **subcomponent** of many NLP tasks, especially those involving **generating text** or **estimating the probability of text**:
 - Predictive typing
 - Speech recognition
 - Handwriting recognition
 - Spelling/grammar correction
 - Authorship identification
 - Machine translation
 - Summarization
 - Dialogue
 - etc.
- Language Modeling has been extended to cover everything else in NLP: **GPT-3 is an LM!**

04 | Example of RNN Language Model

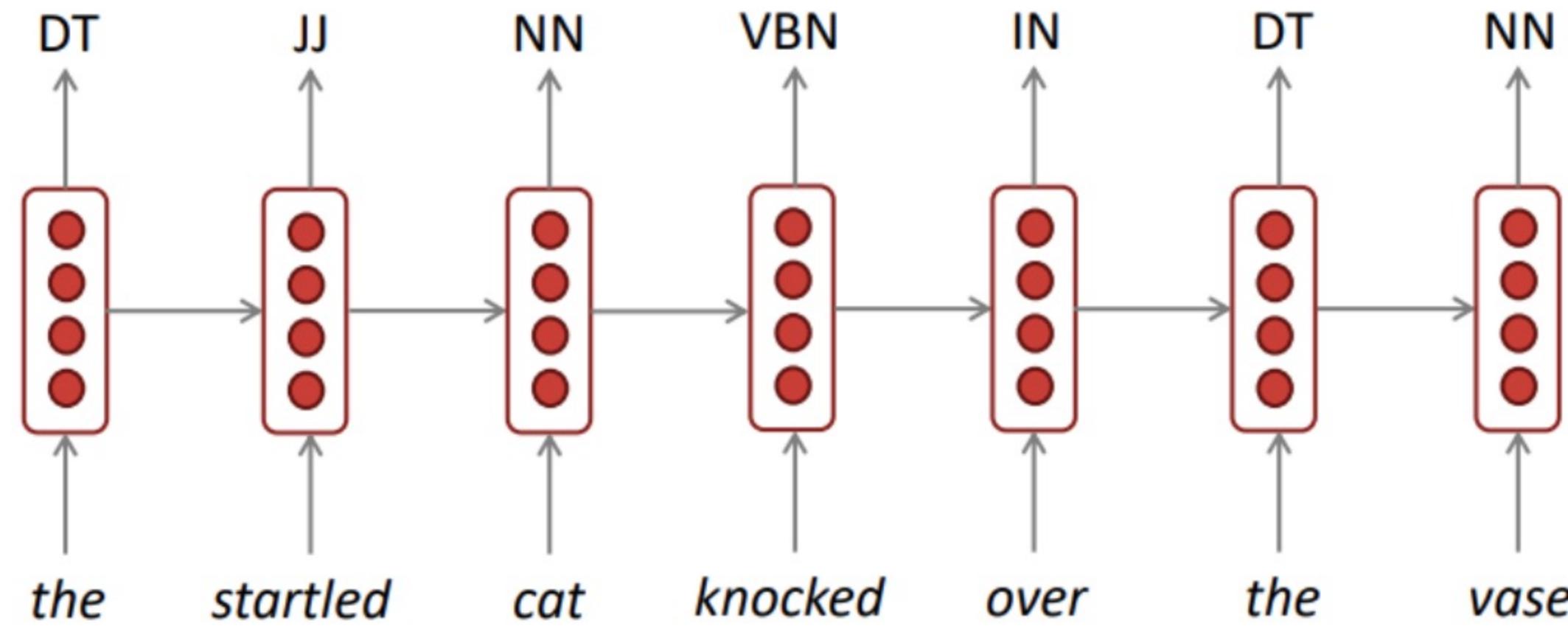
Recap

- **Language Model:** A system that predicts the next word
- **Recurrent Neural Network:** A family of neural networks that:
 - Take sequential input of any length
 - Apply the same weights on each step
 - Can optionally produce output on each step
- Recurrent Neural Network \neq Language Model
- We've shown that RNNs are a great way to build a LM.
- But RNNs are useful for much more!

04

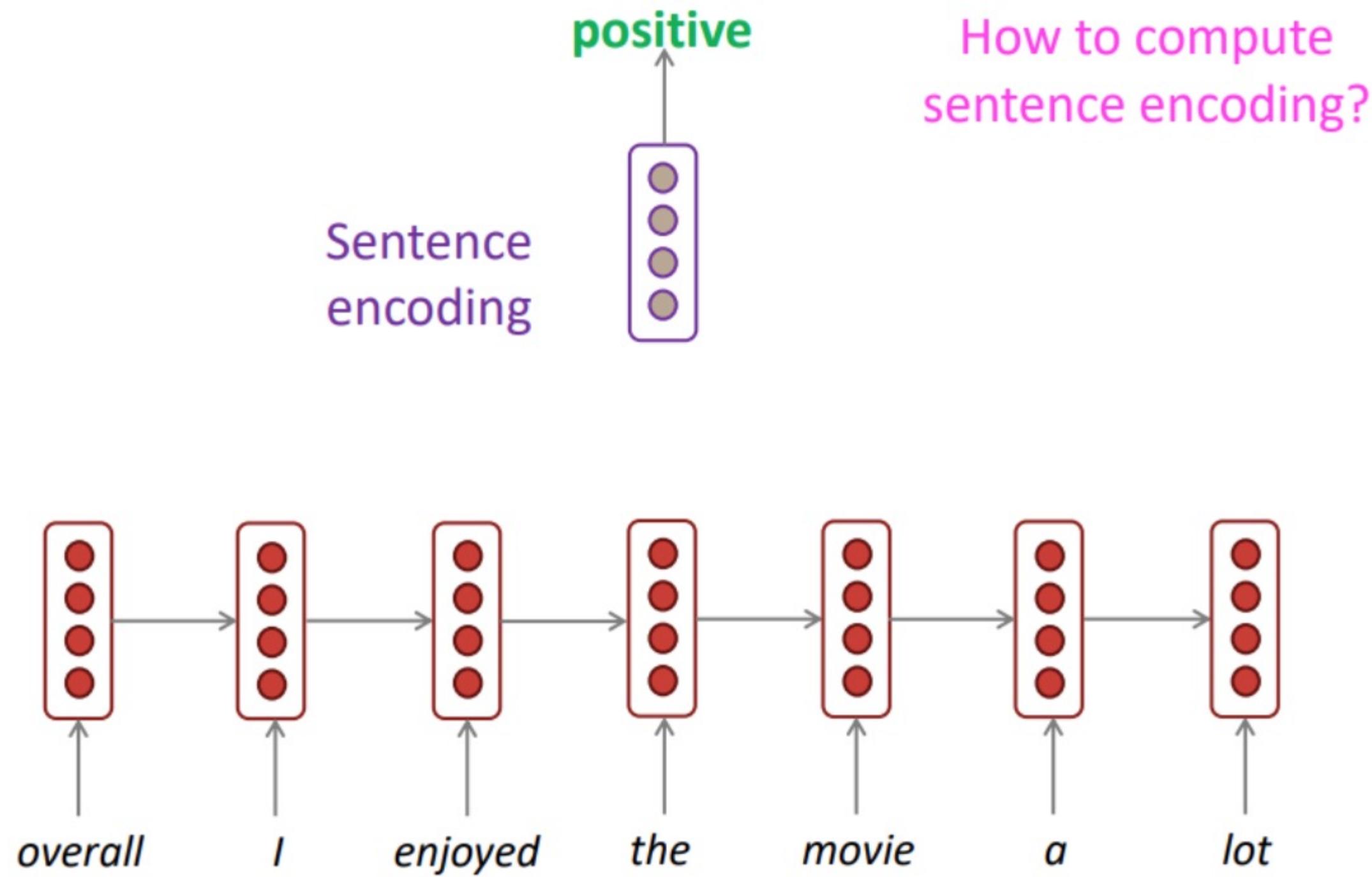
RNNs can be used for tagging

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



04 | RNNs can be used for tagging

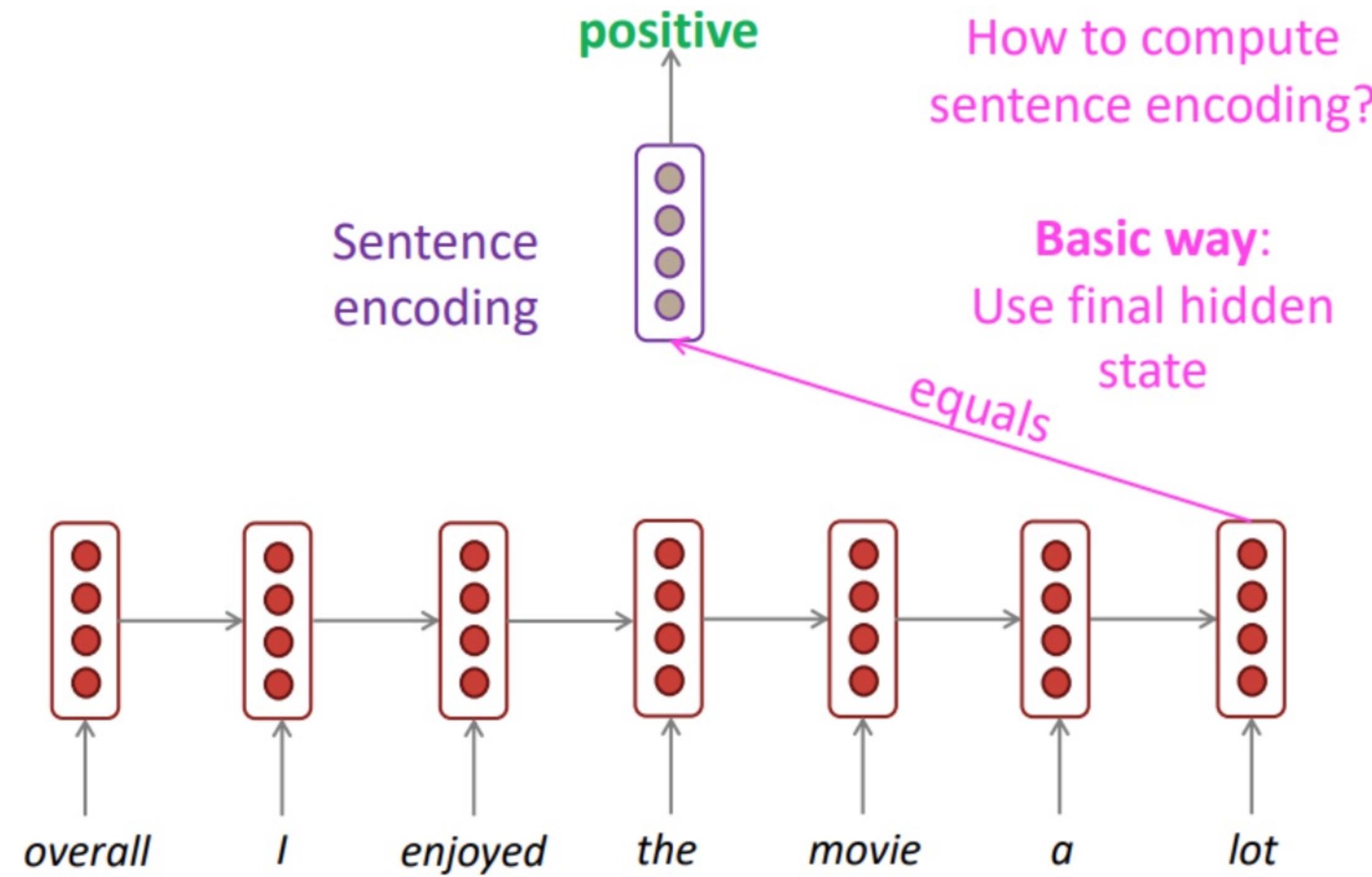
e.g., sentiment classification



04

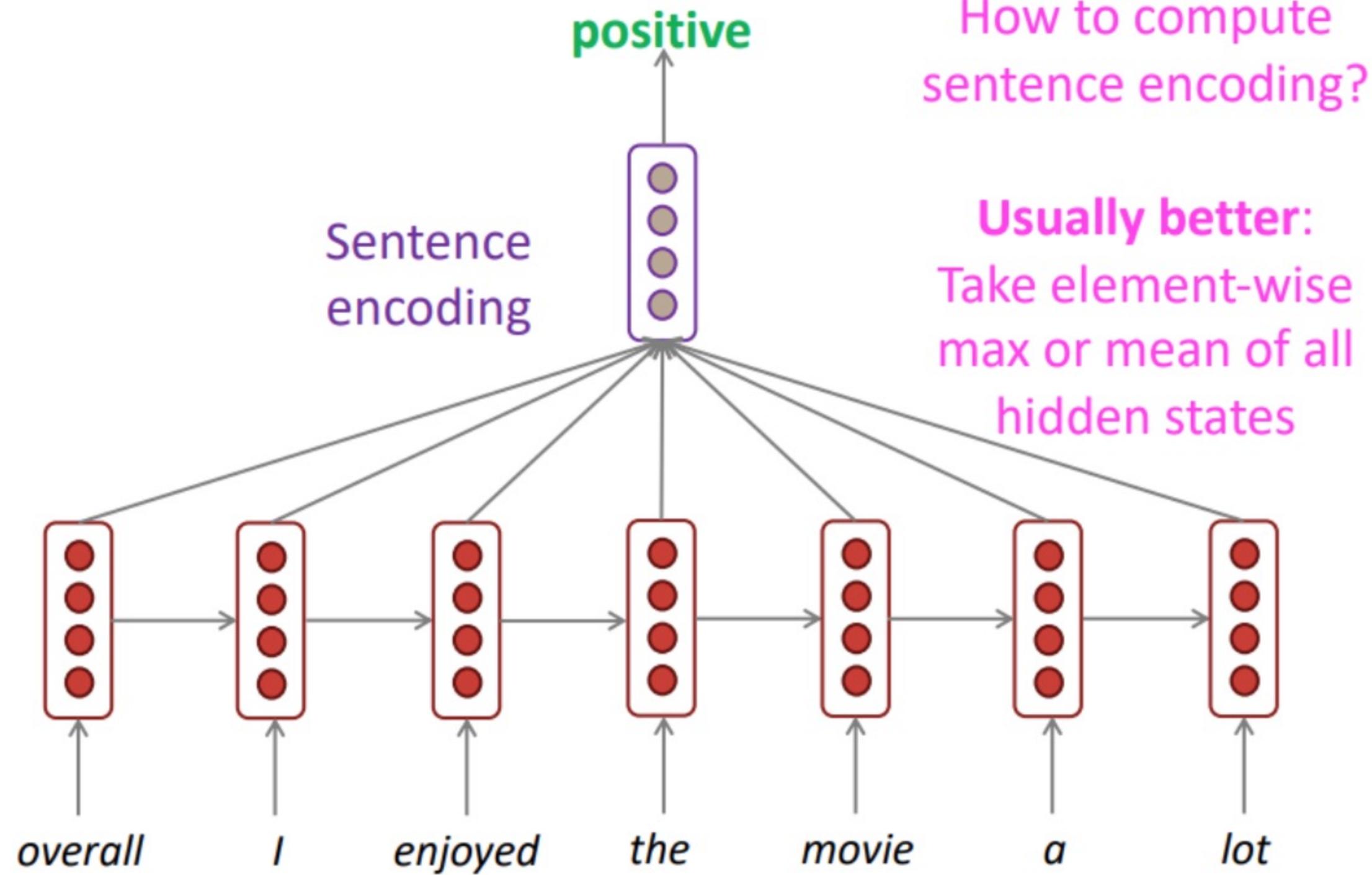
RNNs can be used for tagging

e.g., sentiment classification



04 | RNNs can be used for tagging

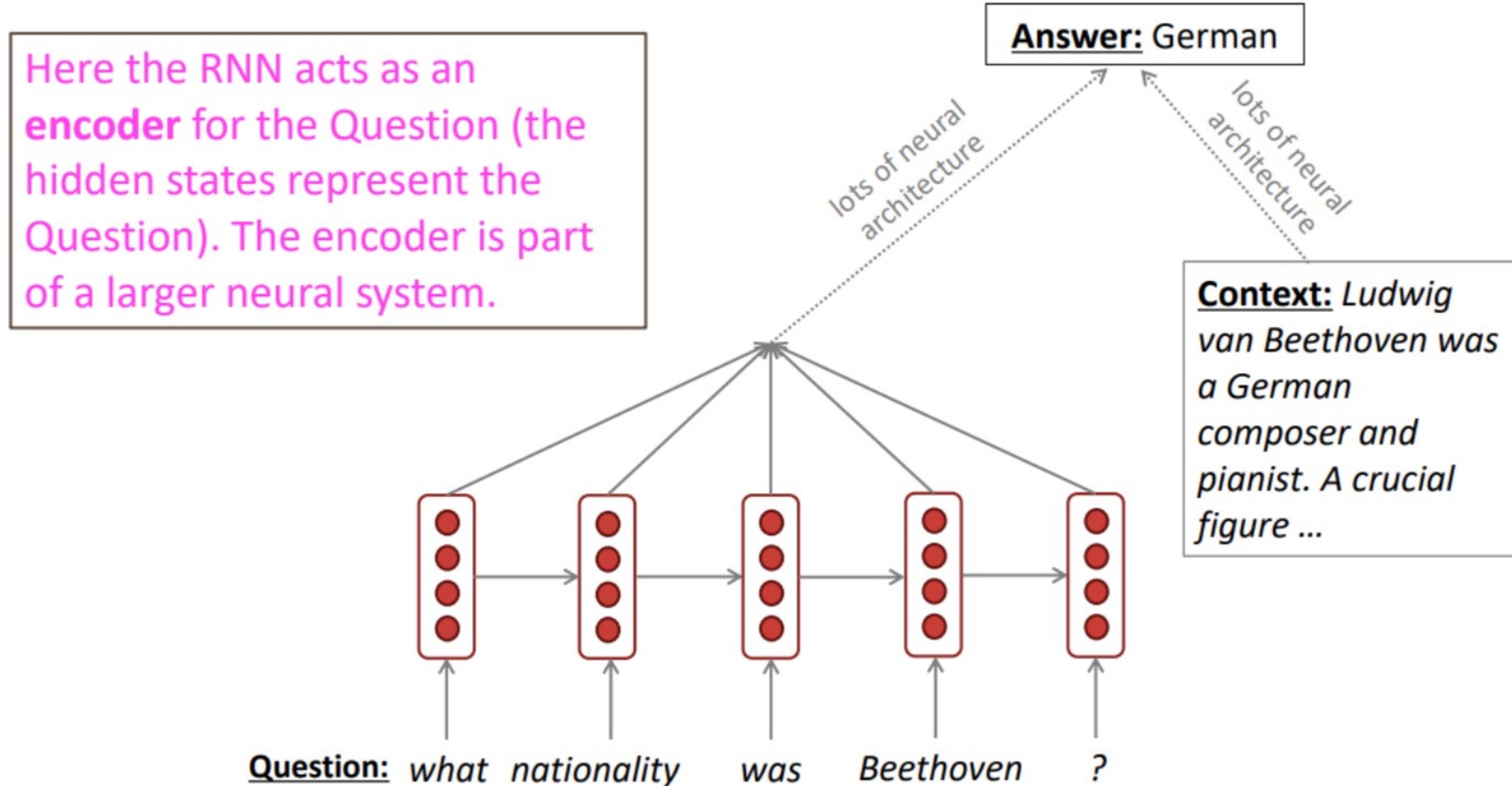
e.g., sentiment classification



04

RNNs can be used for tagging

e.g., question answering, machine translation, many other tasks!



05

Understanding what's going on

The evolution of samples while training

- Trained an LSTM of Leo Tolstoy's War and Peace and then generated samples every 100 iterations of training.

At iteration 100

```
tyntd-iafhatawiaoahrdemot lytdws e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e  
plia tklrgd t o idoe ns,smtt h ne etie h,hregtrs nigtike,aoaenns lng
```

At iteration 300

```
"Tmont thithey" fomesscerliund  
Keushey. Thom here  
sheulke, anmerenith ol sivh I lalterthend Bleipile shuwy fil on aseterlome  
coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."
```

At iteration 500

```
we counter. He stutn co des. His stanted out one ofler that concossions and was  
to gearang reay Jotrets and with fre colt oft paitt thin wall. Which das stimm
```

At iteration 700

```
Aftair fall unsuch that the hall for Prince Velzonski's that me of  
her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort  
how, and Gogition is so overelical and ofter.
```

At iteration 1200

```
"Kite vouch!" he repeated by her  
door. "But I would be done and quarts, feeling, then, son is people...."
```

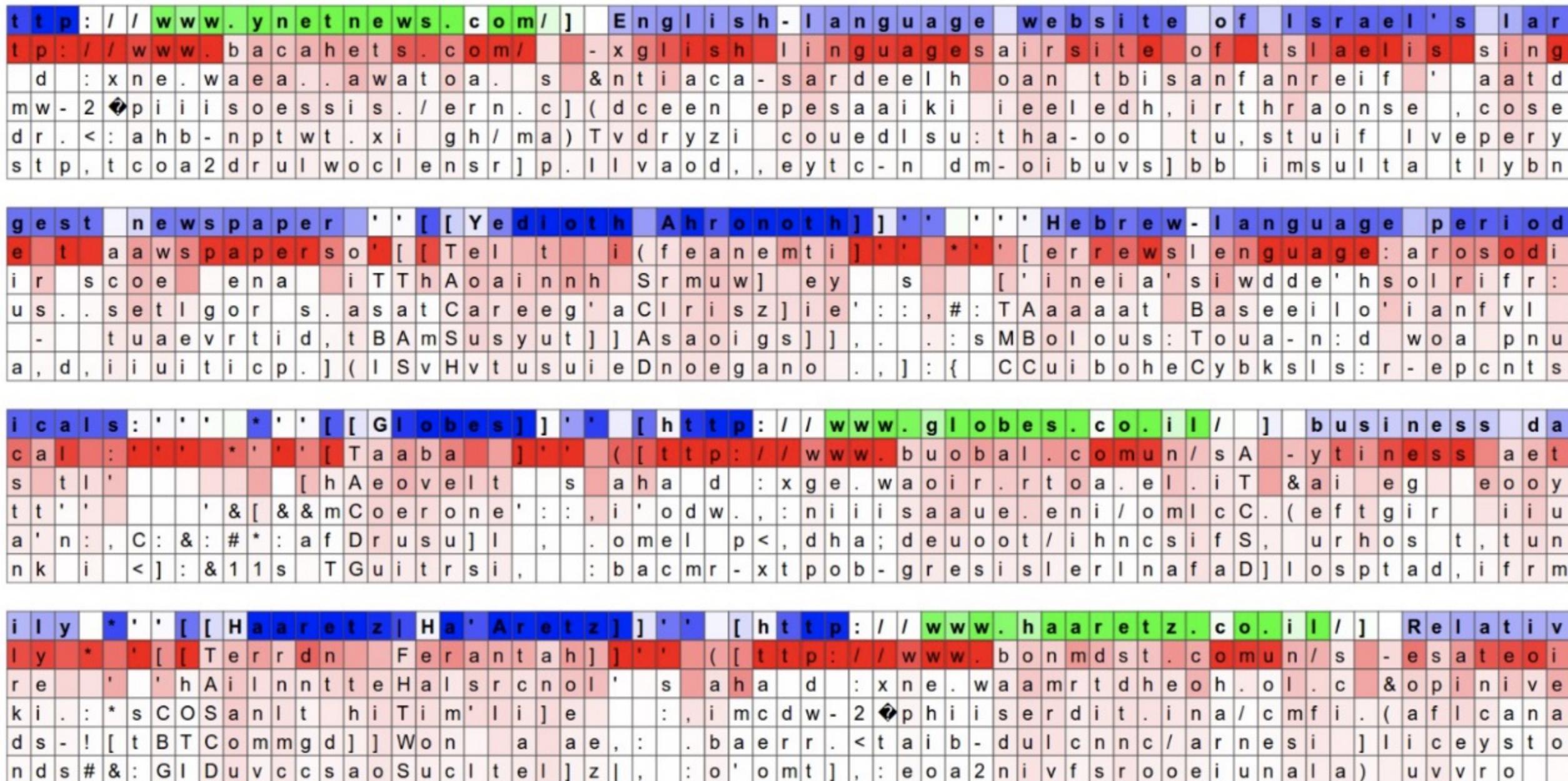
At iteration 2000

```
"Why do what that day," replied Natasha, and wishing to himself the fact the  
princess, Princess Mary was easier, fed in had oftened him.  
Pierre aking his soul came to the packs and drove up his father-in-law women.
```

05 | Understanding what's going on

Visualizing the predictions and the “neuron” firings in the RNN

- Feed a Wikipedia RNN model character data from the validation set (shown along the blue/green rows)
- visualize (in red) the top 5 guesses that the model assigns for the next character.
- The guesses are colored by their probability.
 - so dark red = judged as very likely
 - white = not very likely
- The input character sequence (blue/green) is colored based on the firing of a randomly chosen neuron in the hidden representation of the RNN.
 - green = very excited
 - blue = not very excited



The neuron highlighted in this image seems to get very excited about URLs and turns off outside of the URLs. The LSTM is likely using this neuron to remember if it is inside a URL or not.

05 | Understanding what's going on

Visualizing the predictions and the “neuron” firings in the RNN

Cell sensitive to position in line:

```
The sole importance of the crossing of the Berezina lies in the fact  
that it plainly and indubitably proved the fallacy of all the plans for  
cutting off the enemy's retreat and the soundness of the only possible  
line of action--the one Kutuzov and the general mass of the army  
demanded--namely, simply to follow the enemy up. The French crowd fled  
at a continually increasing speed and all its energy was directed to  
reaching its goal. It fled like a wounded animal and it was impossible  
to block its path. This was shown not so much by the arrangements it  
made for crossing as by what took place at the bridges. When the bridges  
broke down, unarmed soldiers, people from Moscow and women with children  
who were with the French transport, all--carried on by vis inertiae--  
pressed forward into boats and into the ice-covered water and did not,  
surrender.
```

Cell that turns on inside quotes:

```
"You mean to imply that I have nothing to eat out of.... On the  
contrary, I can supply you with everything even if you want to give  
dinner parties," warmly replied Chichagov, who tried by every word he  
spoke to prove his own rectitude and therefore imagined Kutuzov to be  
animated by the same desire.
```

```
Kutuzov, shrugging his shoulders, replied with his subtle penetrating  
smile: "I meant merely to say what I said."
```

Cell that robustly activates inside if statements:

```
static int __dequeue_signal(struct sigpending *pending, sigset_t *mask,  
    siginfo_t *info)  
{  
    int sig = next_signal(pending, mask);  
    if (sig) {  
        if (current->notifier) {  
            if (sigismember(current->notifier_mask, sig)) {  
                if (!!(current->notifier)(current->notifier_data)) {  
                    clear_thread_flag(TIF_SIGPENDING);  
                    return 0;  
                }  
            }  
        }  
        collect_signal(sig, pending, info);  
    }  
    return sig;  
}
```

A large portion of cells are not easily interpretable. Here is a typical example:

```
/* Unpack a filter field's string representation from user-space  
 * buffer. */  
char *audit_unpack_string(void **bufp, size_t *remain, size_t len)  
{  
    char *str;  
    if (!*bufp || (len == 0) || (len > *remain))  
        return ERR_PTR(-EINVAL);  
    /* of the currently implemented string fields, PATH_MAX  
     * defines the longest valid length.  
     */  
    if (len > PATH_MAX)  
        return ERR_PTR(-ENAMETOOLONG);  
    str = kmalloc(len + 1, GFP_KERNEL);  
    if (unlikely(!str))  
        return ERR_PTR(-ENOMEM);  
    memcpy(str, *bufp, len);  
    str[len] = 0;  
    *bufp += len;  
    *remain -= len;  
    return str;
```

Cell that turns on inside comments and quotes:

```
/* Duplicate LSM field information. The lsm_rule is opaque, so  
 * re-initialized. */  
static inline int audit_dupe_lsm_field(struct audit_field *df,  
    struct audit_field *sf)  
{  
    int ret = 0;  
    char *lsm_str;  
    /* our own copy of lsm_str */  
    lsm_str = kstrdup(sf->lsm_str, GFP_KERNEL);  
    if (unlikely(!lsm_str))  
        return -ENOMEM;  
    df->lsm_str = lsm_str;  
    /* our own (refreshed) copy of lsm_rule */  
    ret = security_audit_rule_init(df->type, df->op, df->lsm_str,  
        (void **)&df->lsm_rule);  
    /* Keep currently invalid fields around in case they  
     * become valid after a policy reload. */  
    if (ret == -EINVAL) {  
        pr_warn("audit rule for LSM '%s' is invalid\n",  
            df->lsm_str);  
        ret = 0;  
    }  
    return ret;  
}
```

Cell that is sensitive to the depth of an expression:

```
#ifdef CONFIG_AUDITSYSCALL  
static inline int audit_match_class_bits(int class, u32 *mask)  
{  
    int i;  
    if (classes[class]) {  
        for (i = 0; i < AUDIT_BITMASK_SIZE; i++)  
            if (mask[i] & classes[class][i])  
                return 0;  
    }  
    return 1;  
}
```

Cell that might be helpful in predicting a new line. Note that it only turns on for some ")":

```
char *audit_unpack_string(void **bufp, size_t *remain, si  
{  
    char *str;  
    if (!*bufp || (len == 0) || (len > *remain))  
        return ERR_PTR(-EINVAL);  
    /* of the currently implemented string fields, PATH_MAX  
     * defines the longest valid length.  
     */  
    if (len > PATH_MAX)  
        return ERR_PTR(-ENAMETOOLONG);  
    str = kmalloc(len + 1, GFP_KERNEL);  
    if (unlikely(!str))  
        return ERR_PTR(-ENOMEM);  
    memcpy(str, *bufp, len);  
    str[len] = 0;  
    *bufp += len;  
    *remain -= len;  
    return str;
```

06 | Conclusion

Terminology and a look forward

The RNN described in this lecture = **simple/vanilla/Elman RNN**

Next lecture: You will learn about other RNN flavors

like **GRU**



and **LSTM**



and multi-layer RNNs

Gated Recurrent Units

Long Short-Term Memory



References

1. <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>
2. <https://www.deeplearningbook.org/contents/rnn.html>
3. <http://web.stanford.edu/class/cs224n/slides/cs224n-2022-lecture05-rnnlm.pdf>
4. http://web.stanford.edu/class/cs224n/readings/cs224n-2019-notes05-LM_RNN.pdf
5. Neural Networks and Deep Learning_A Textbook

THANK YOU

Q & A