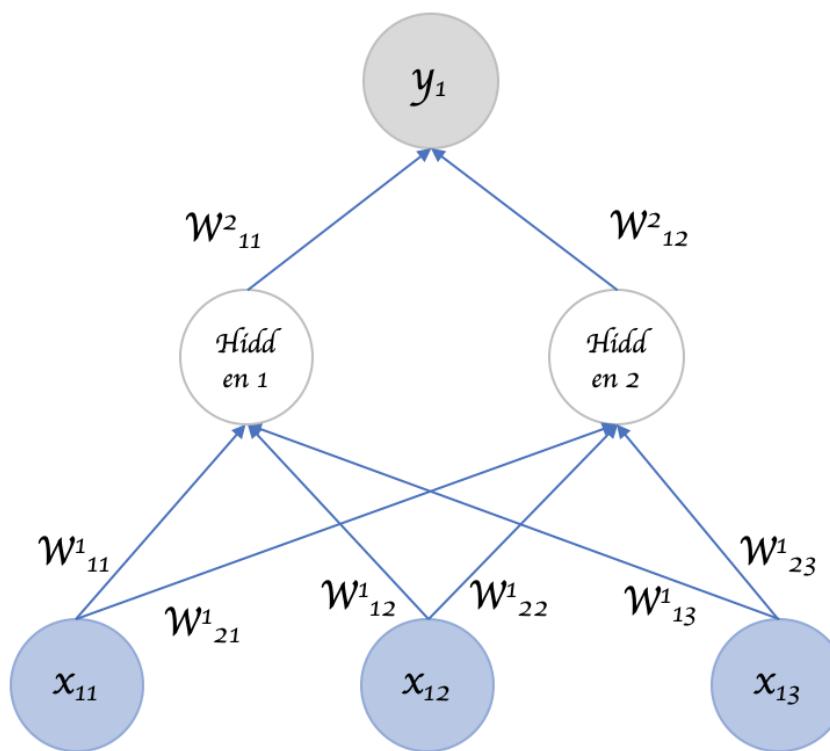
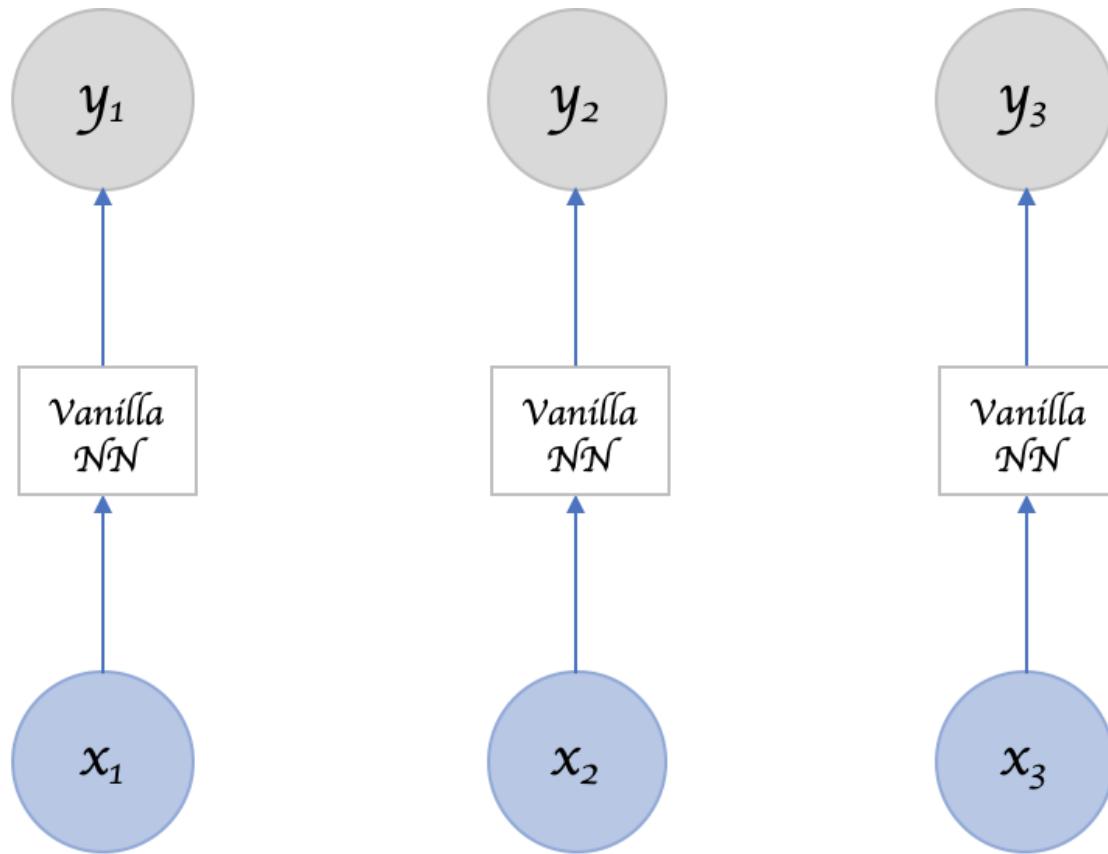


Introduction to Recurrent Neural Network

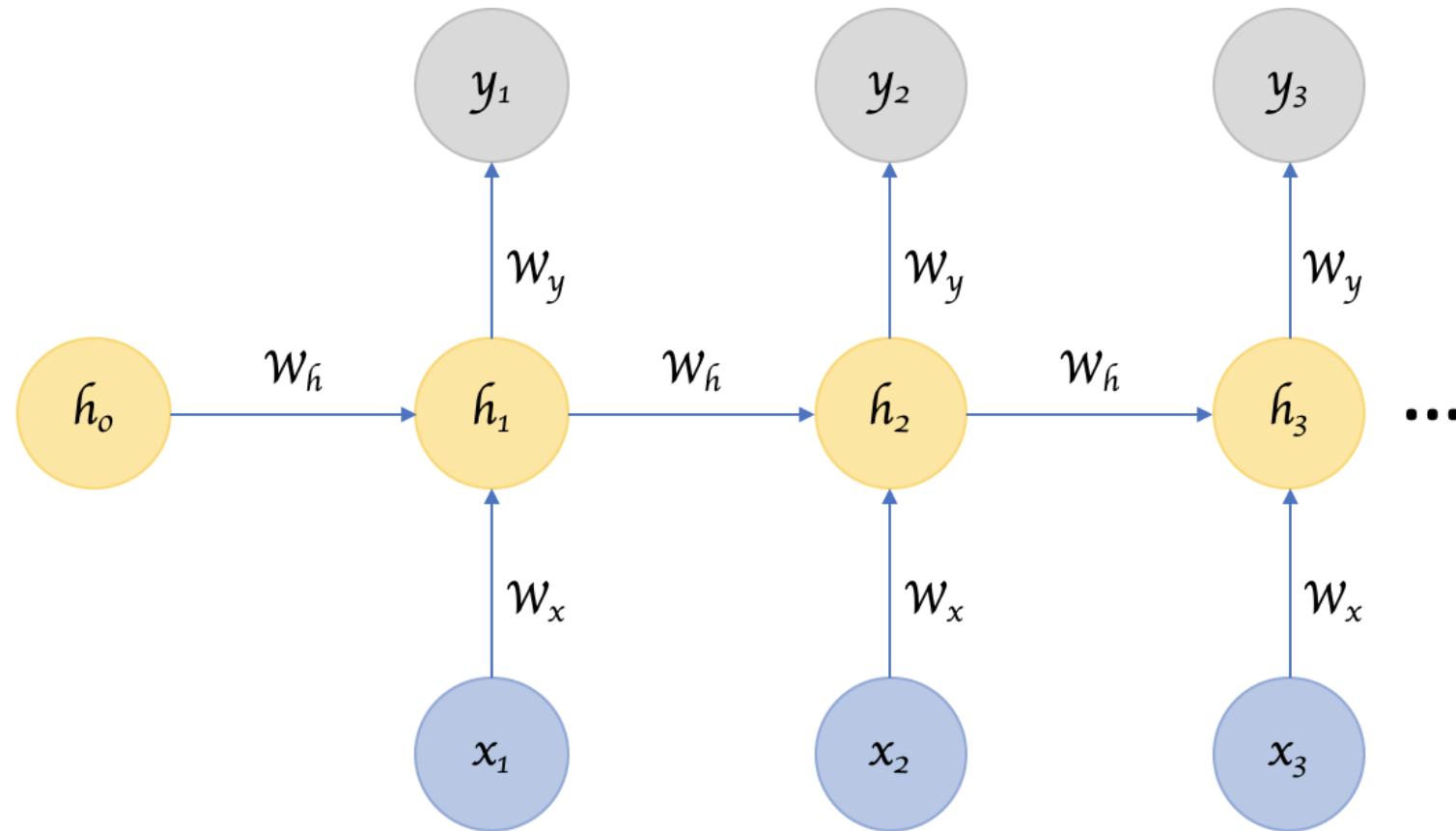
A simple neural network with 1 hidden layer



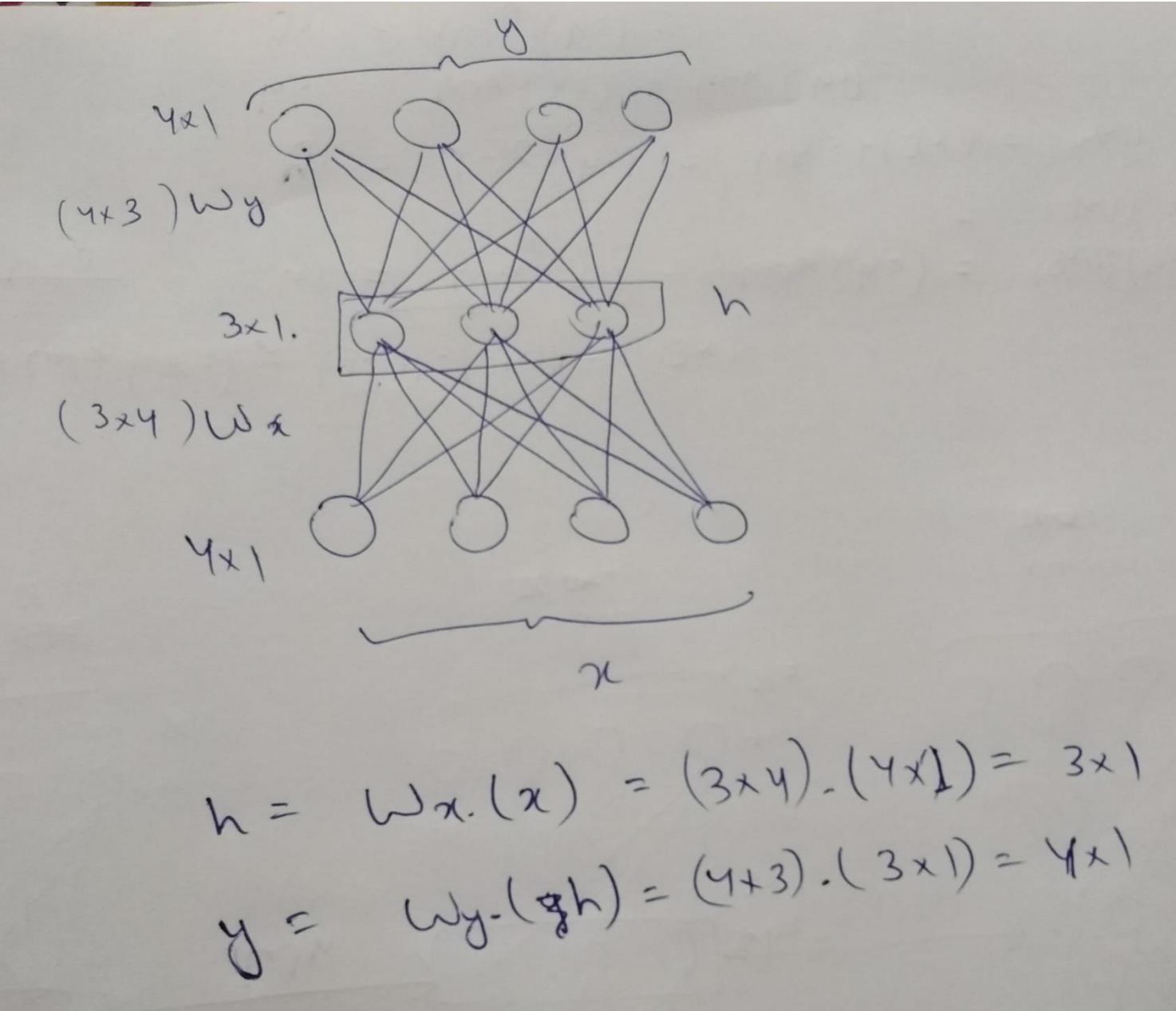
A simple neural network with 1 hidden layer



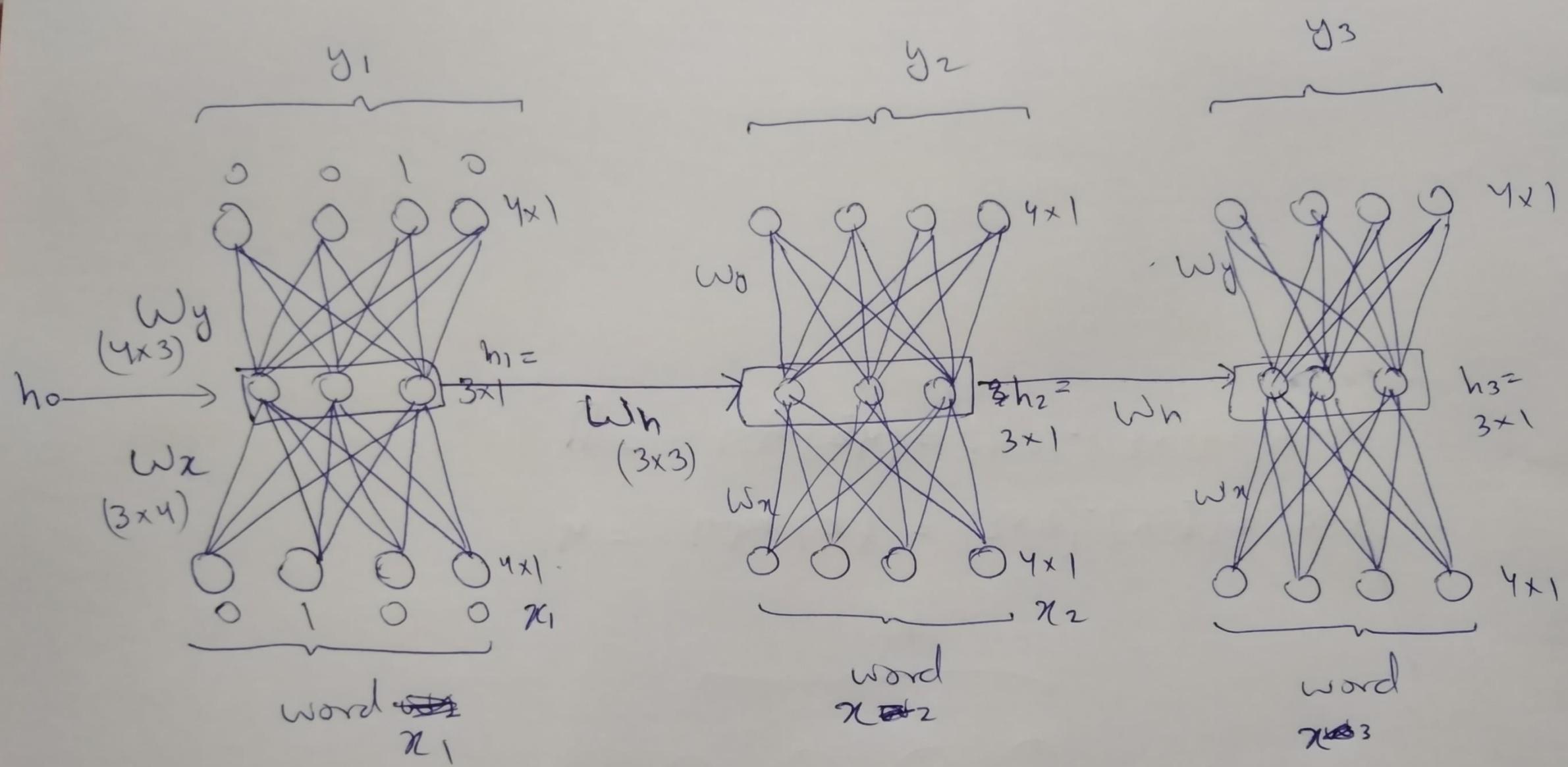
A recurrent neural network with 1 hidden layer



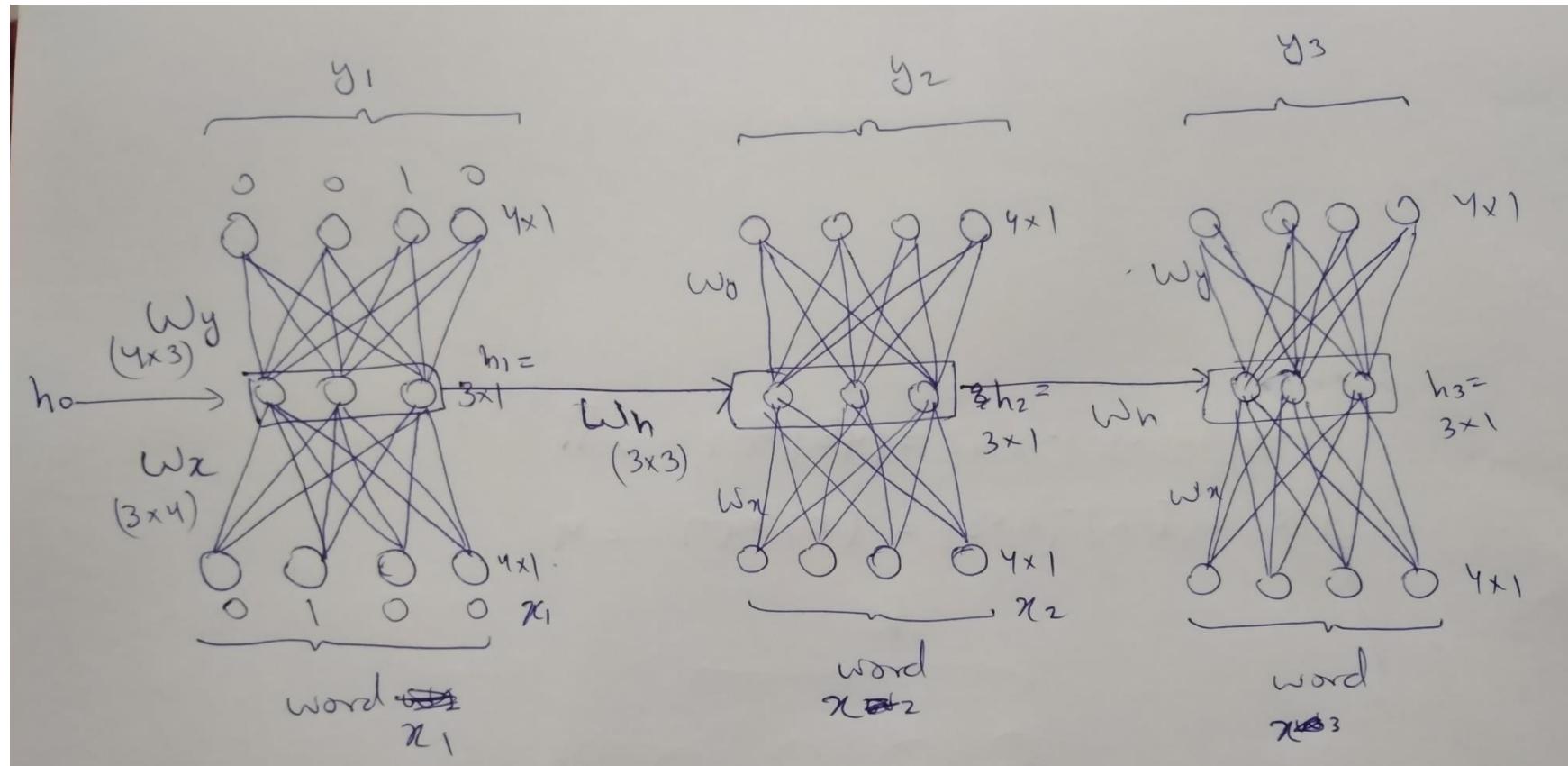
A simple neural network with 1 hidden layer



A simple recurrent neural network with 1 hidden layer



A simple recurrent neural network with 1 hidden layer



$$v = w_x x_1 + w_h h_0 + \dots$$

$$\begin{aligned} h_1 &= w_x(x_1) + w_h(h_0) = (3 \times 4)(4 \times 1) + (3 \times 3)(3 \times 1) \\ &= (3 \times 1) + (3 \times 1) = 3 \times 1 \end{aligned}$$

$$y_1 = w_y(h_1) = (4 \times 3)(3 \times 1) = 4 \times 1$$

$$h_2 = w_x(x_2) + w_h(h_1)$$

$$y_2 = (w_y)(h_2) = -$$

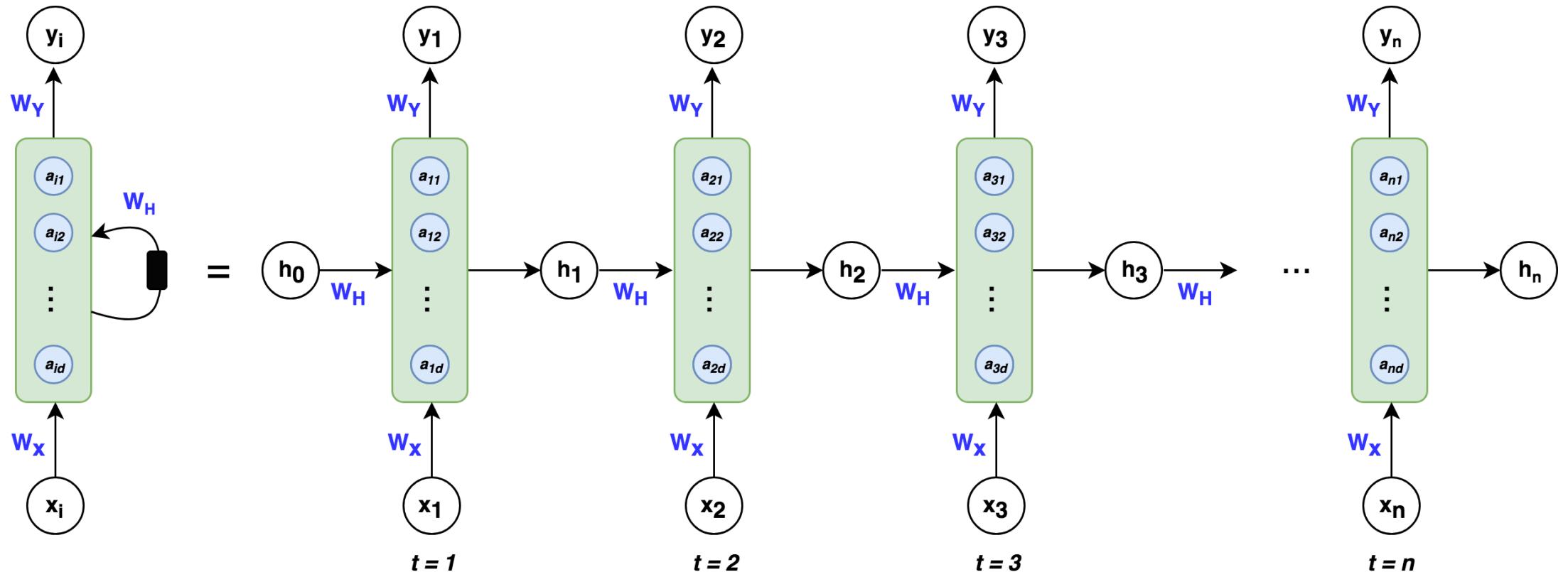
Applications of RNN

Sequence Data

- Text (sequence of words or sequence of characters)
- Stock prices (sequence of stock prices daily)
- DNA sequence data
- Consumer purchase history
- Web surfing history
- Many more

RNN Equations for Forward Pass

Left: Shorthand notation often used for RNNs, **Right:** Unfolded notation for RNNs



One-Hot Encoding of the word “dogs”

$$\mathbf{d} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

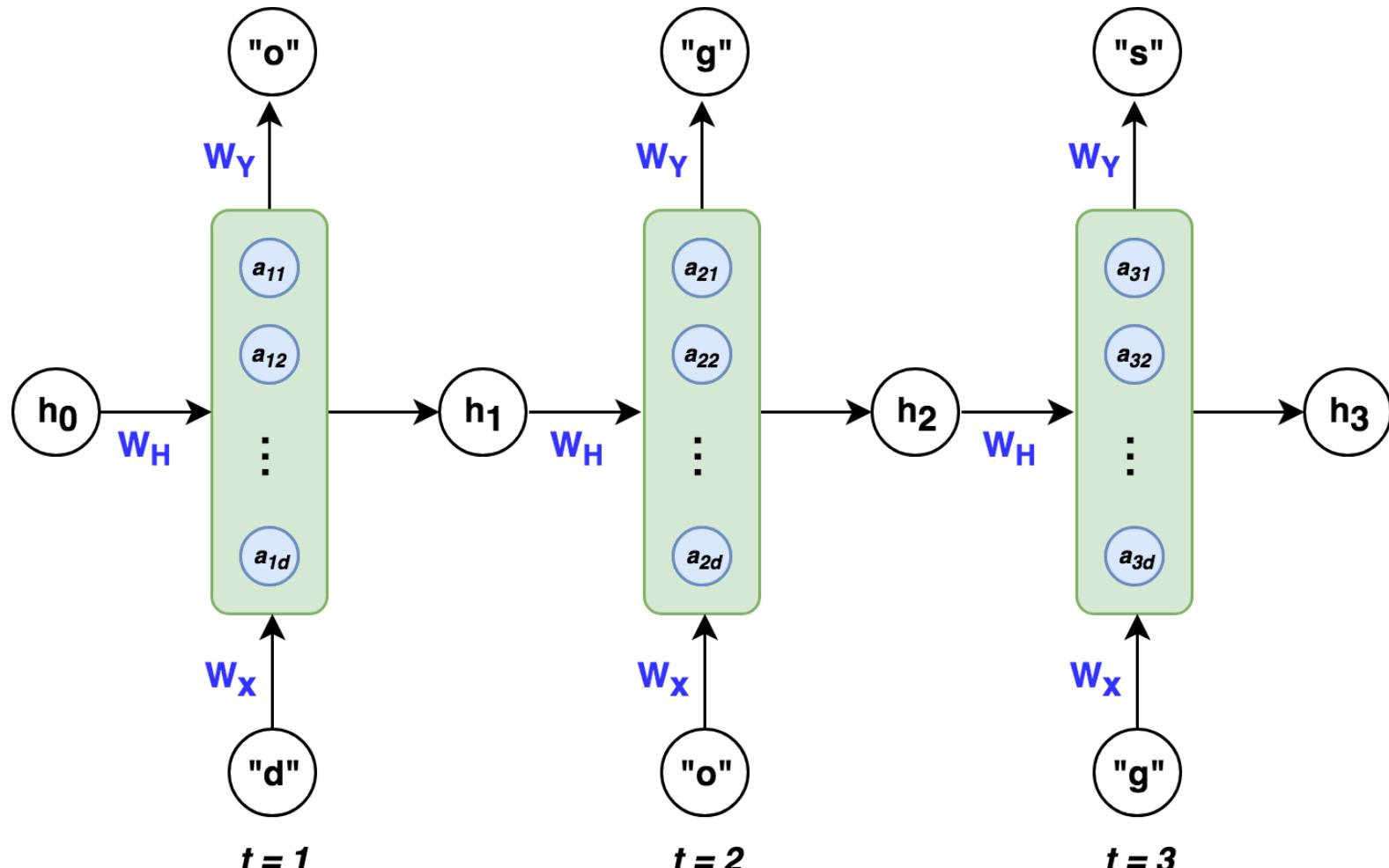
$$\mathbf{o} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\mathbf{g} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

$$\mathbf{s} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

An Example

RNN architecture predicting the letter “s” in “dogs”



RNN Equations

$$a_t = W_H h_{t-1} + W_X X_t$$

Hidden Nodes

**Output
From Hidden
State**

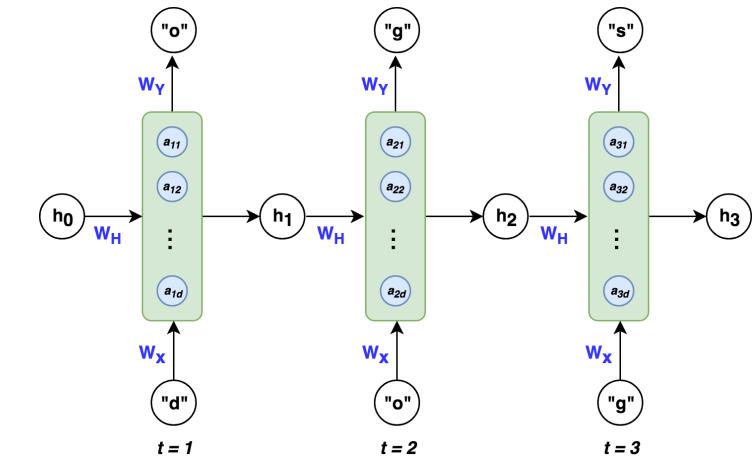
$$\rightarrow h_t = \tanh(a_t)$$

Activation Function

Hidden State

**Prediction at
time t**

$$\rightarrow y_t = \text{softmax}(W_Y h_t)$$



RNN Equations

$$a_t = W_H h_{t-1} + W_X X_t$$

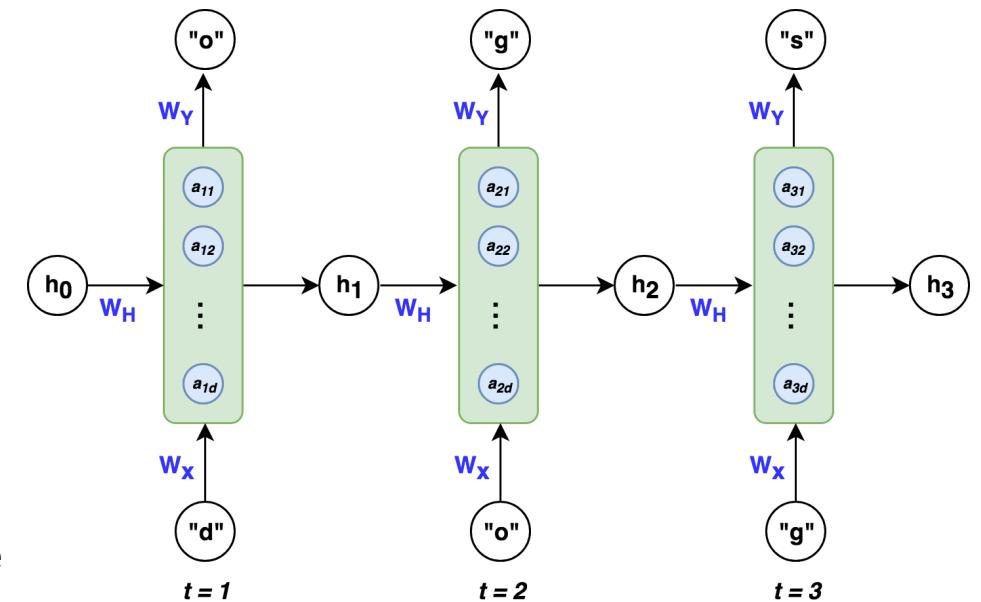
Where a_t is a vector of dimension $(k \times 1)$, W_H is a matrix of dimension $(d \times d)$, h_{t-1} is a vector of dimension $(d \times 1)$, W_X is a matrix of dimension $(d \times k)$, and X_t is a vector of dimension $(k \times 1)$.

$$h_t = \tanh(a_t)$$

Where h_t is a vector of dimension $(k \times 1)$, a_t is a vector of dimension $(k \times 1)$, and \tanh is the hyperbolic tangent function.

$$y_t = \text{softmax}(W_Y h_t)$$

Where y_t is a vector of dimension $(k \times 1)$, W_Y is a matrix of dimension $(k \times d)$, and softmax is the softmax function.



Where

k is the dimension of the input vector x_i

d is the number of hidden nodes

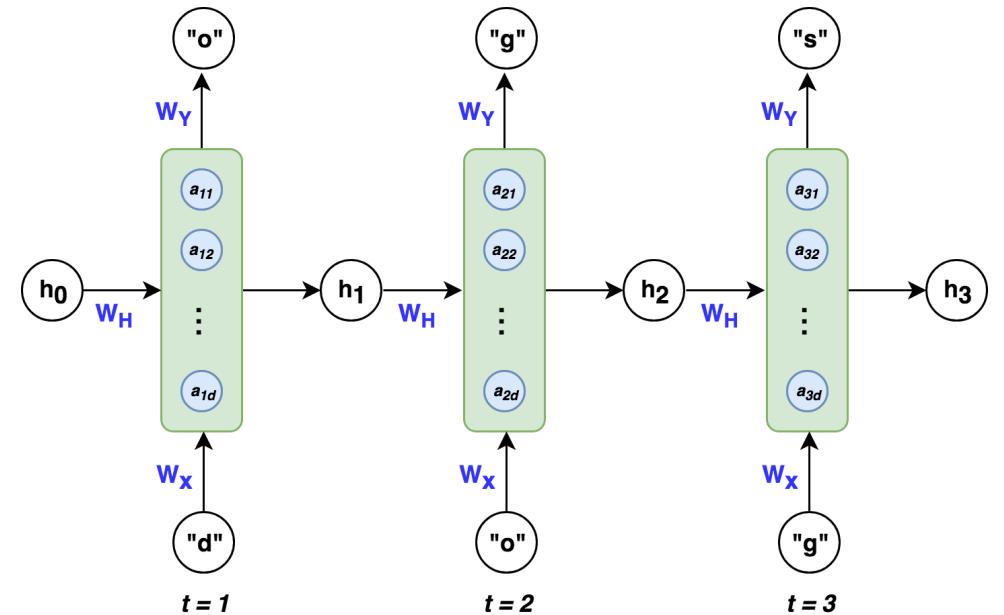
Matrix Dimensions

$$a_t = W_H h_{t-1} + W_X X_t$$

(3x3) (3x1) (4x1)

(3x1) $\longrightarrow h_t = \tanh(a_t)$ (3x4) (4x3)

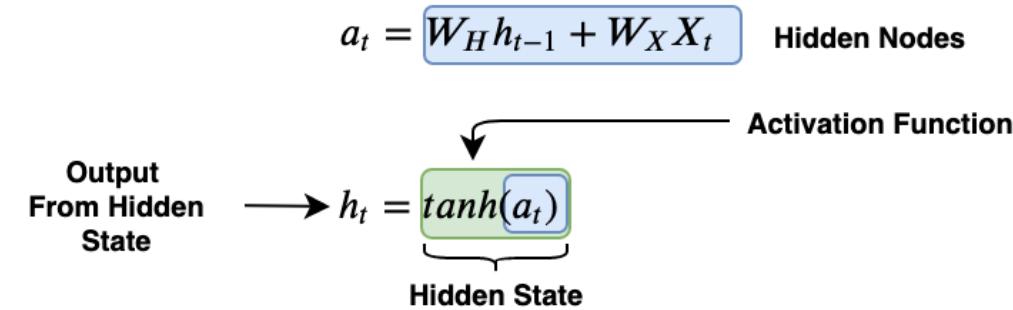
(4x1) $\longrightarrow y_t = \text{softmax}(W_Y h_t)$



Forward Propagation

At t = 1

$$a_1 = \begin{pmatrix} W_{H,11} & W_{H,12} & W_{H,13} \\ W_{H,21} & W_{H,22} & W_{H,23} \\ W_{H,31} & W_{H,32} & W_{H,33} \end{pmatrix} \begin{pmatrix} h_{0,1} \\ h_{0,2} \\ h_{0,3} \end{pmatrix} + \begin{pmatrix} W_{X,11} & W_{X,12} & W_{X,13} & W_{X,14} \\ W_{X,21} & W_{X,22} & W_{X,23} & W_{X,24} \\ W_{X,31} & W_{X,32} & W_{X,33} & W_{X,34} \end{pmatrix} \begin{pmatrix} x_{1,1} \\ x_{1,2} \\ x_{1,3} \\ x_{1,4} \end{pmatrix} = \begin{pmatrix} a_{1,1} \\ a_{1,2} \\ a_{1,3} \end{pmatrix}$$



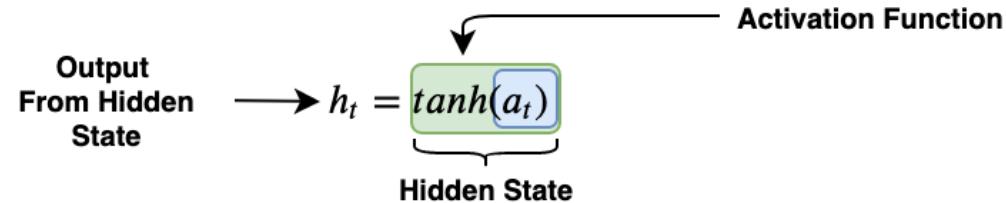
$$\text{Prediction at time t} \rightarrow y_t = \text{softmax}(W_Y h_t)$$

$$h_1 = \tanh\left(\begin{pmatrix} a_{1,1} \\ a_{1,2} \\ a_{1,3} \end{pmatrix}\right) = \begin{pmatrix} h_{1,1} \\ h_{1,2} \\ h_{1,3} \end{pmatrix}$$

$$y_1 = \text{softmax}\left(\begin{pmatrix} W_{Y,11} & W_{Y,12} & W_{Y,13} \\ W_{Y,21} & W_{Y,22} & W_{Y,23} \\ W_{Y,31} & W_{Y,32} & W_{Y,33} \\ W_{Y,41} & W_{Y,42} & W_{Y,43} \end{pmatrix} \begin{pmatrix} h_{1,1} \\ h_{1,2} \\ h_{1,3} \end{pmatrix}\right) = \begin{pmatrix} y_{1,1} \\ y_{1,2} \\ y_{1,3} \\ y_{1,4} \end{pmatrix}$$

Forward Propagation

$$a_t = W_H h_{t-1} + W_X X_t \quad \text{Hidden Nodes}$$



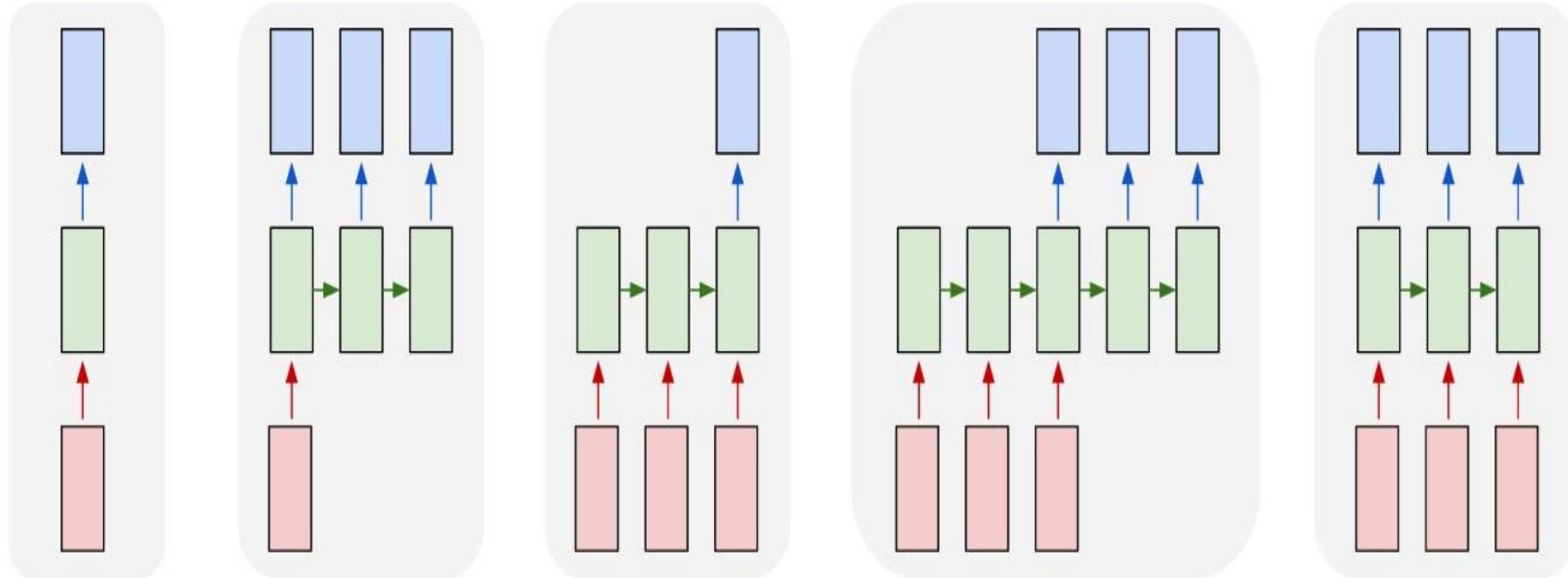
$$\text{Prediction at time t} \longrightarrow y_t = \text{softmax}(W_Y h_t)$$

$$a_1 = \begin{pmatrix} 0.1 & 0.5 & 0.1 \\ 0.5 & 0.9 & 0.3 \\ 0.3 & 0.2 & 0.1 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} + \begin{pmatrix} 0.6 & 0.8 & 0.4 & 0.8 \\ 0.2 & 0.2 & 0.8 & 0.7 \\ 0.9 & 0.8 & 0.1 & 0.2 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0.6 \\ 0.2 \\ 0.9 \end{pmatrix}$$

$$h_1 = \tanh\left(\begin{pmatrix} 0.6 \\ 0.2 \\ 0.9 \end{pmatrix}\right) = \begin{pmatrix} 0.54 \\ 0.20 \\ 0.72 \end{pmatrix}$$

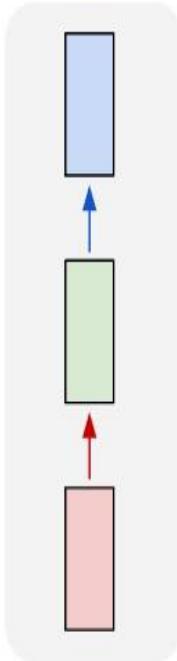
$$y_1 = \text{softmax}\left(\begin{pmatrix} 0.9 & 0.8 & 0.3 \\ 0.2 & 0.3 & 0.4 \\ 0.6 & 0.9 & 0.1 \\ 0.5 & 0.0 & 0.3 \end{pmatrix} \begin{pmatrix} 0.54 \\ 0.20 \\ 0.72 \end{pmatrix}\right) = \begin{pmatrix} 0.32 \\ 0.21 \\ 0.24 \\ 0.22 \end{pmatrix}$$

one to one one to many many to one many to many many to many

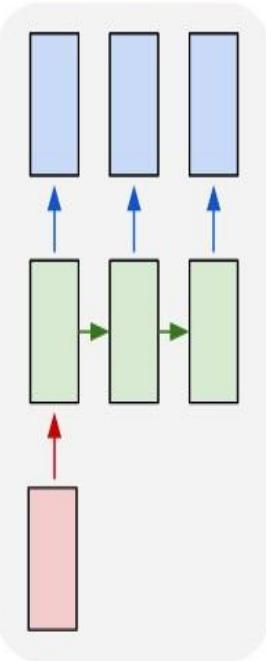


Vanilla Neural Networks

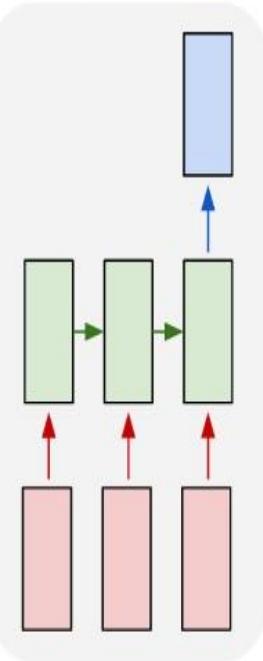
one to one



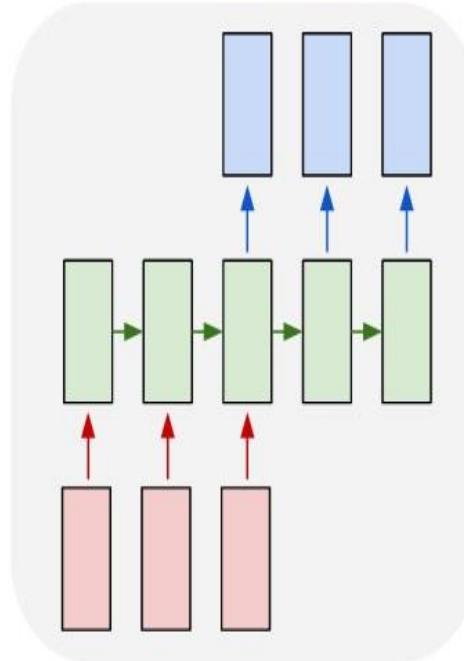
one to many



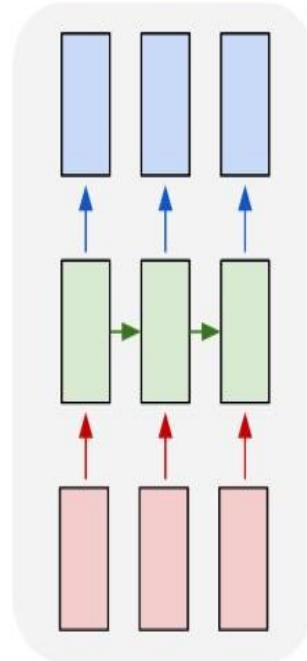
many to one



many to many

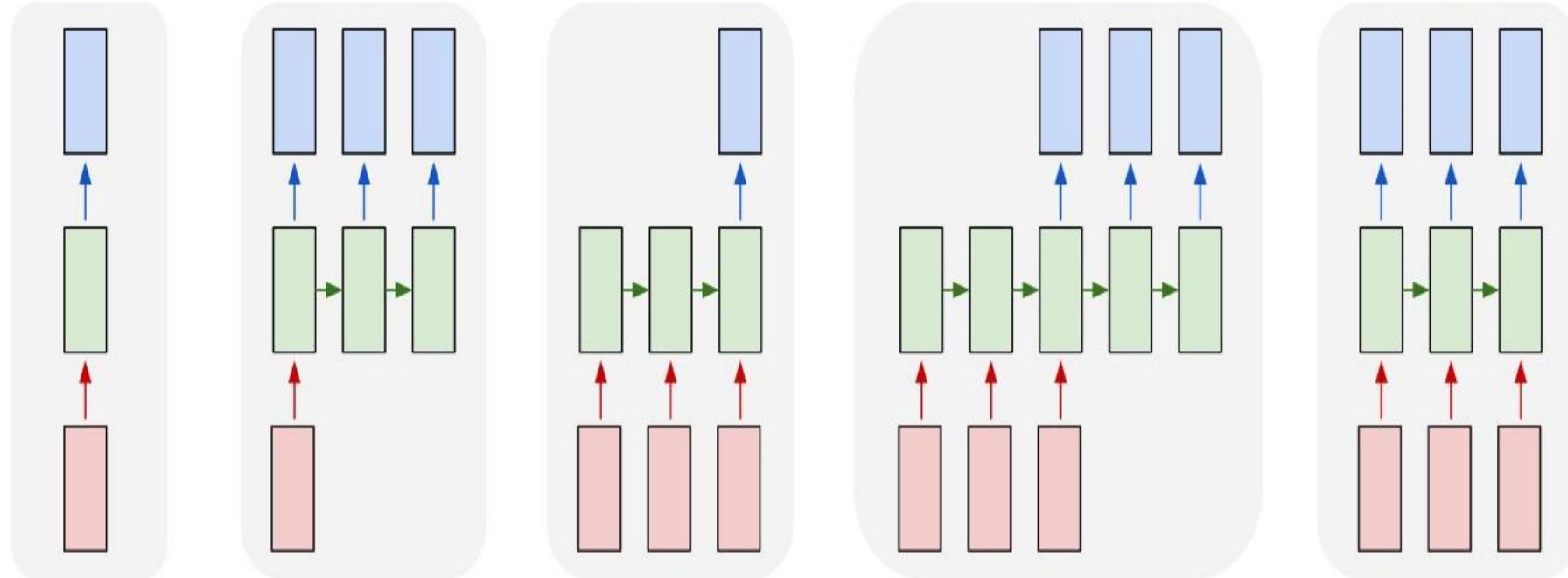


many to many



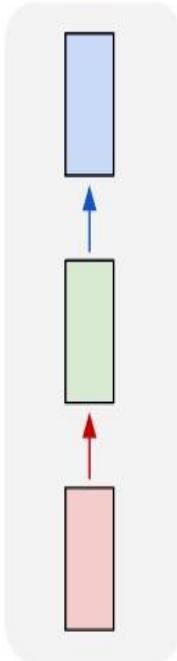
e.g. **Image Captioning**
image -> sequence of words

one to one one to many many to one many to many many to many

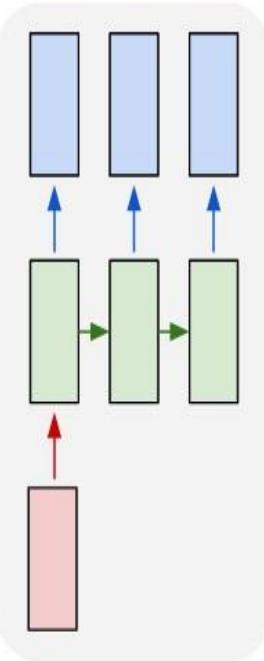


e.g. **Sentiment Classification**
sequence of words -> sentiment

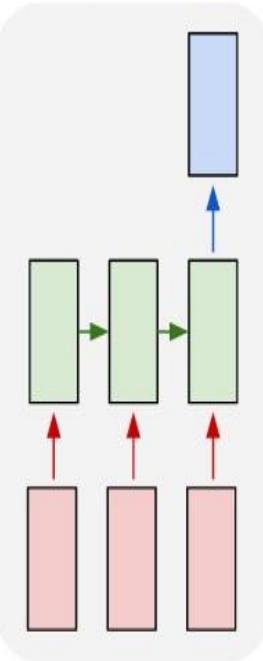
one to one



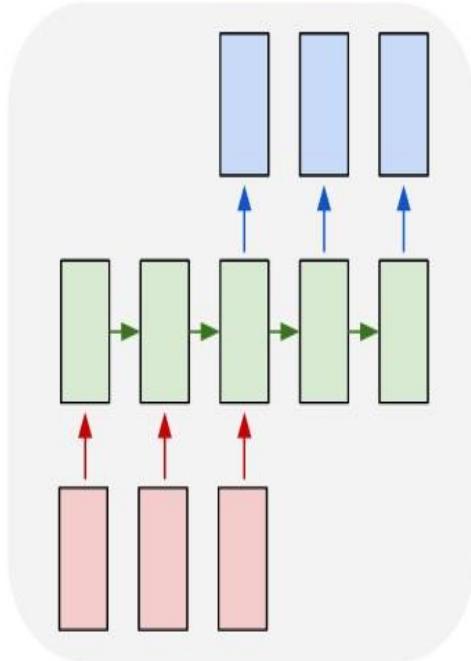
one to many



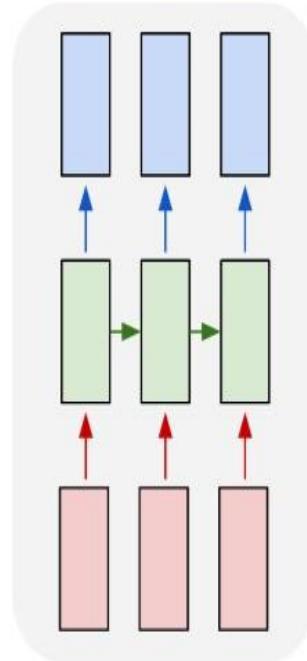
many to one



many to many

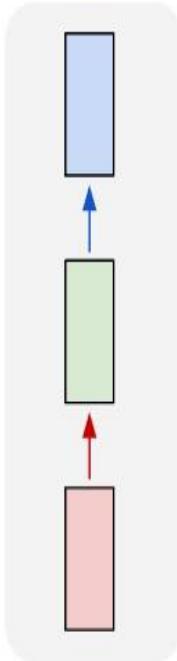


many to many

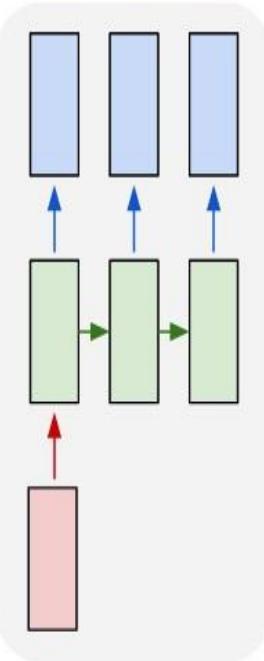


e.g. **Machine Translation**
seq of words \rightarrow seq of words

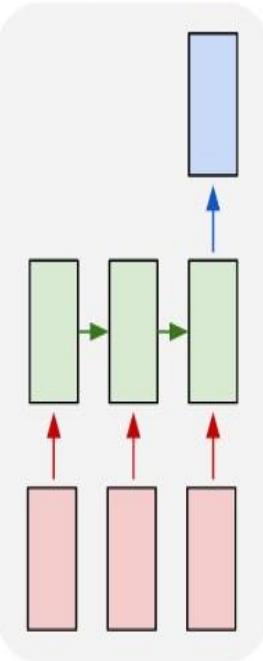
one to one



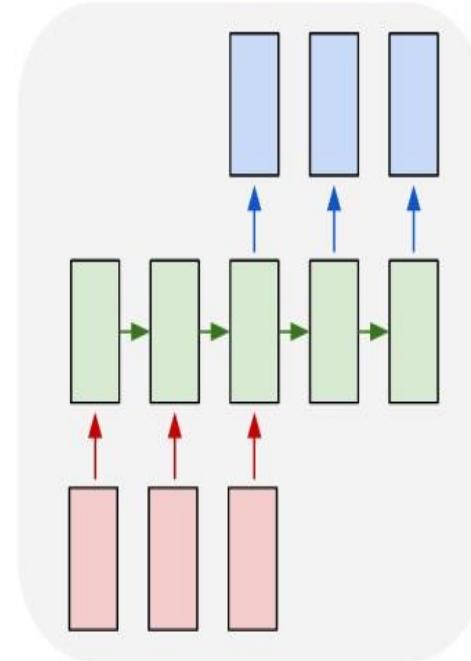
one to many



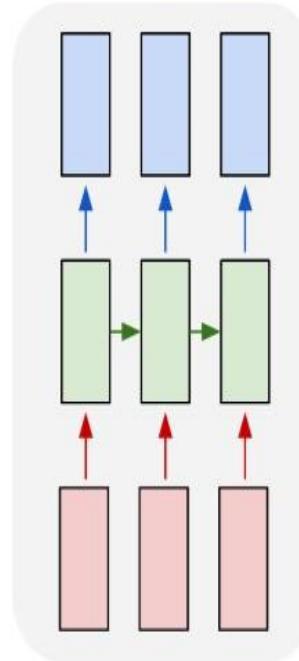
many to one



many to many



many to many



e.g. Video classification on frame level

A RNN Language Model

output distribution

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

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

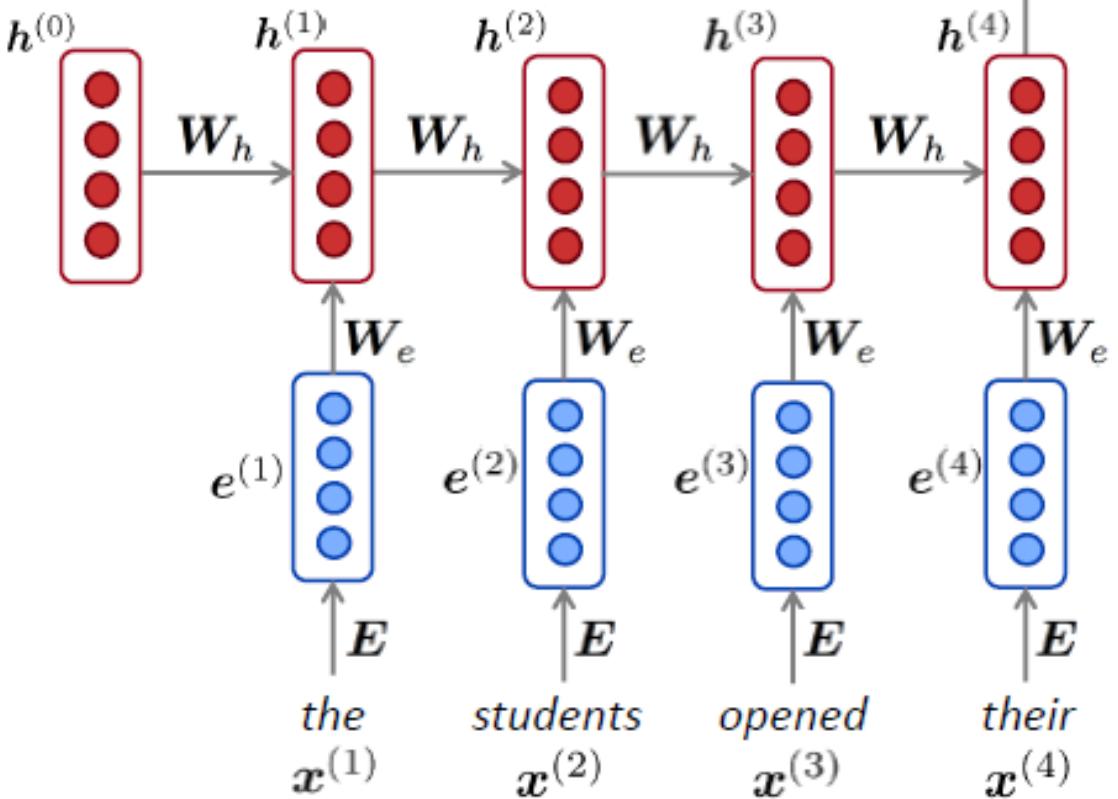
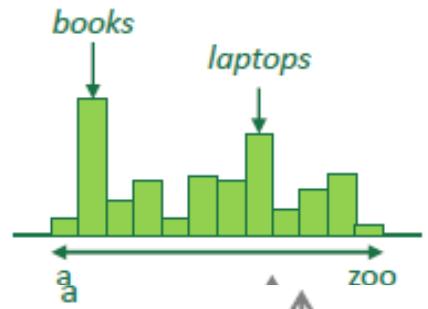
word embeddings

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

words / one-hot vectors

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

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



Note: this input sequence could be much longer, but this slide doesn't have space!

A RNN Language Model

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

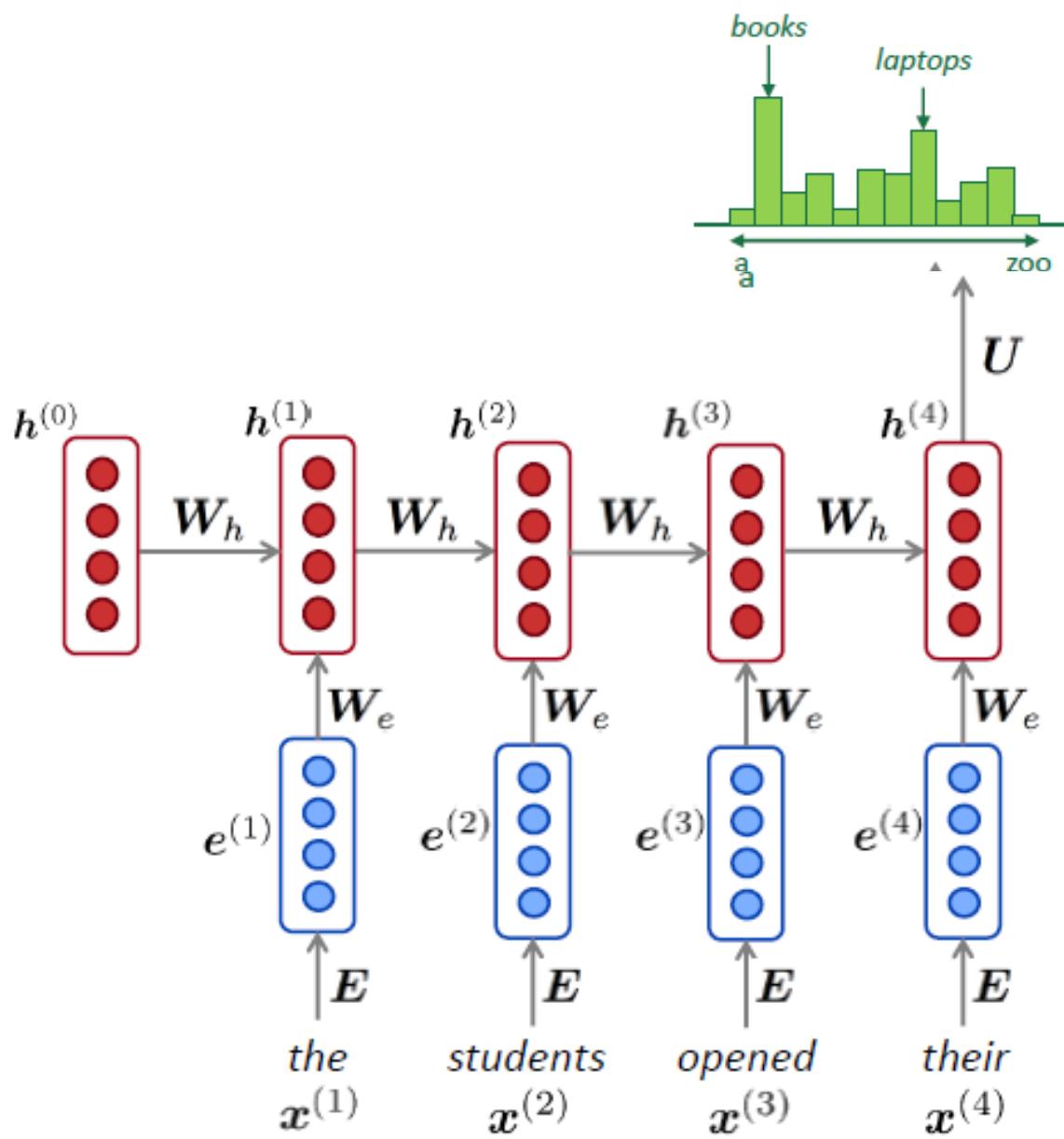
RNN Advantages:

- 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



Task

- Draw RNN architecture diagram and write equations along with dimensions of all layers and weight matrices for following language model:
- input sequence of length 5 (lets say 5 words).
- Hidden layer units are 4.
- $V = \text{vocabulary} = 10$

For reference revisit slides of lecture. Submit handwritten solution (snapshot) on Google classroom under assignment.

RNN for Language Modeling (Loss Function)

Cross Entropy Loss

Output for Binary Labels

$$y = [1, 0]$$

Cross Entropy Loss

$$= \underbrace{y \log \hat{y}}_{y=1} + \underbrace{(1-y) \log (1-\hat{y})}_{y=0}$$

Cross Entropy Loss for Language Modeling

Output Vector = $y_w = \begin{bmatrix} 0 & 0 & \dots & 0 & 1 & 0 & \dots & 0 \end{bmatrix}$
 $w_1 \quad w_2 \quad \dots \quad w_{i+1} \quad w_v$

$\hat{y}_w = \begin{bmatrix} 0.01 & 0.02 & 0 & \dots & 0.3 & 0 & \dots & 0.1 \end{bmatrix}$
 $w_1 \quad w_2 \quad \dots \quad w_{i+1} \quad \dots \quad w_v$

Cross Entropy Loss

$$y_{w_1} \log \hat{y}_{w_1} + \dots + y_{w_{i+1}} \log \hat{y}_{w_{i+1}} + \dots + y_{w_v} \log \hat{y}_{w_v}$$

since all y_w are 0 except y_{w_i} so

$$\text{Loss} = -y_{w_{i+1}} \log \hat{y}_{w_{i+1}}$$

Training a RNN Language Model (Loss Function)

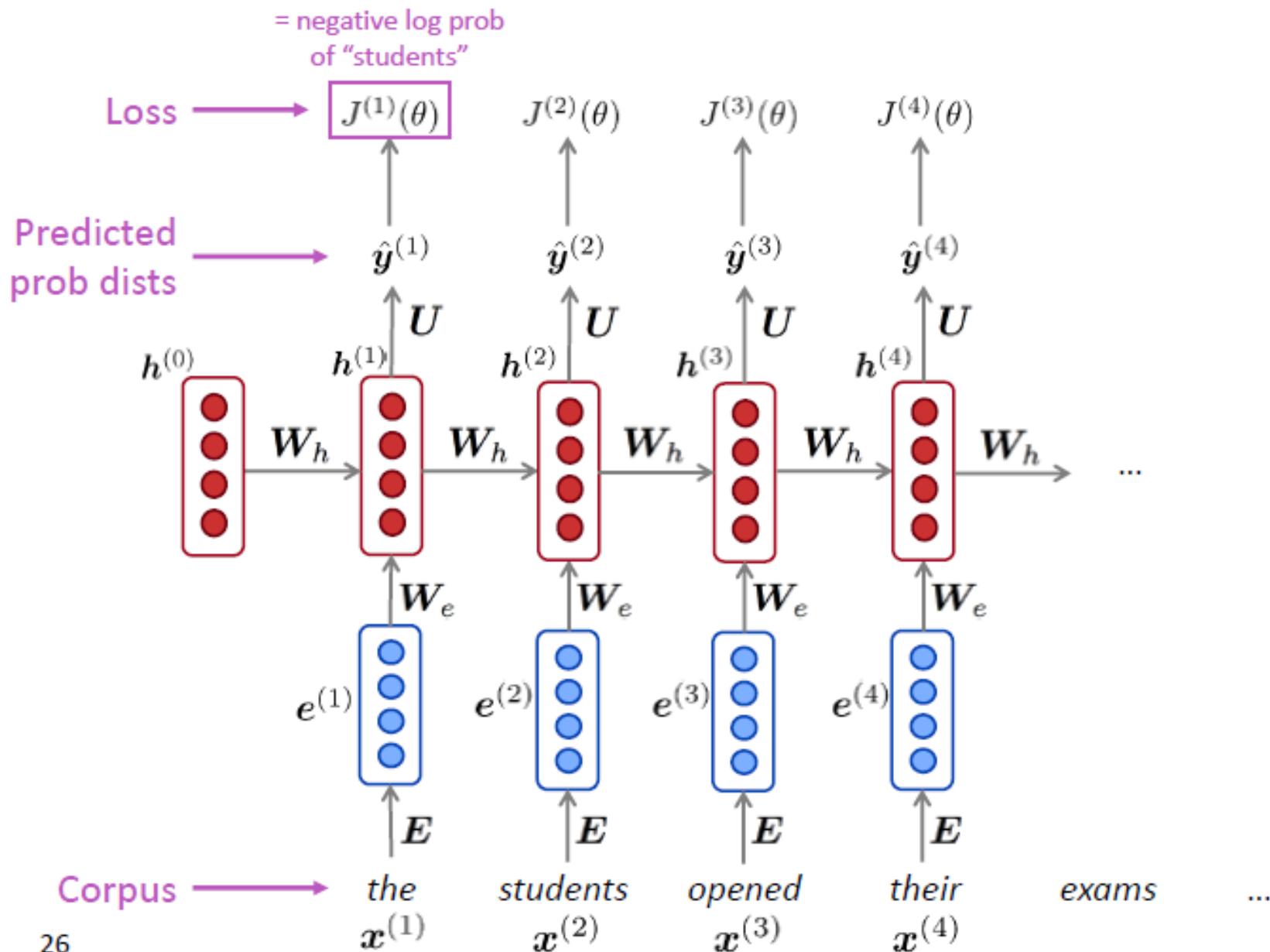
- Get a **big corpus of text** which is a sequence of words $x^{(1)}, \dots, x^{(T)}$
- Feed into RNN-LM; compute output distribution $\hat{y}^{(t)}$ **for every step t .**
 - i.e. predict probability dist of *every word*, given words so far
- Loss function on step t is **cross-entropy** between predicted probability distribution $\hat{y}^{(t)}$, and the true next word $y^{(t)}$ (one-hot for $x^{(t+1)}$):

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

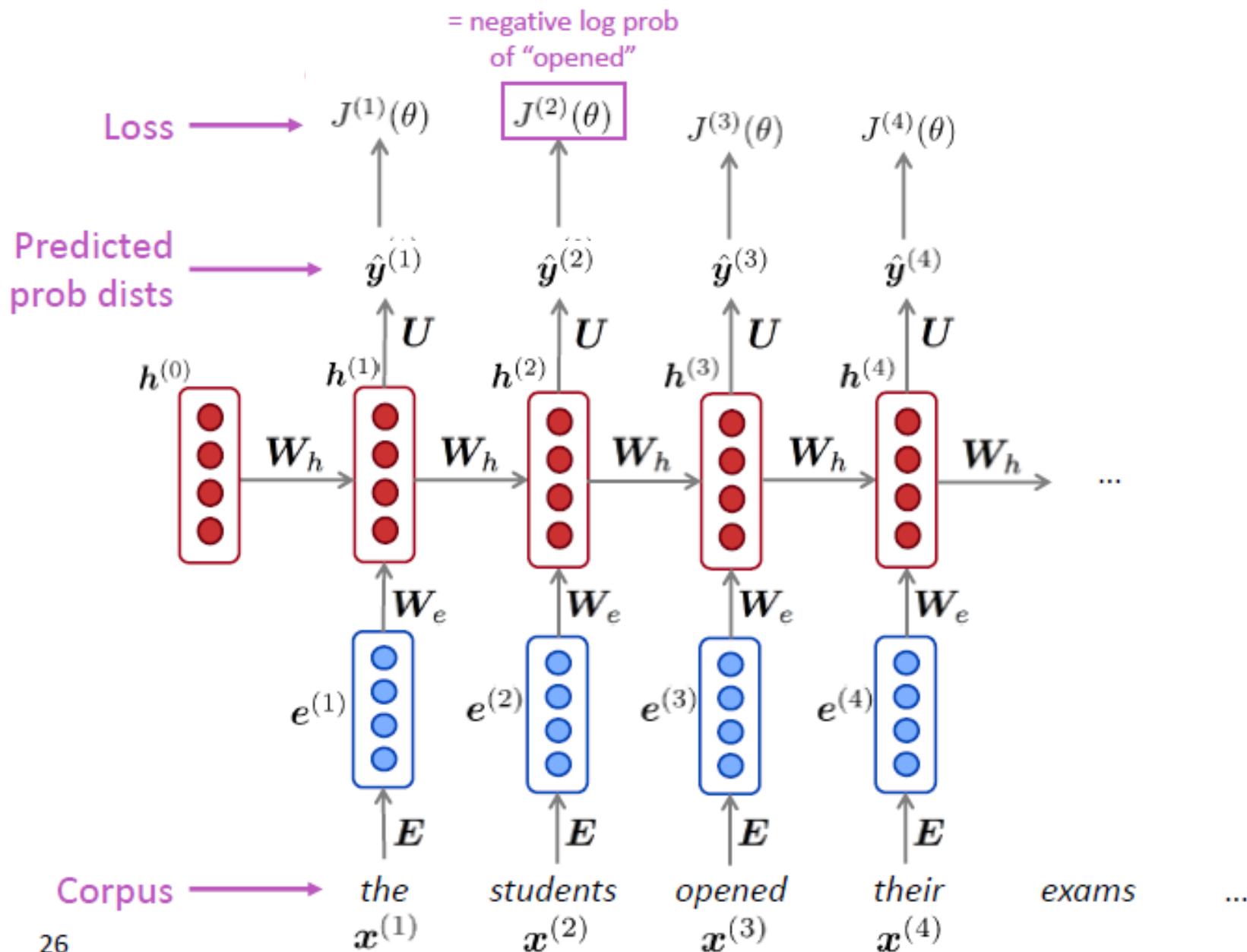
- Average this to get **overall loss** for entire training set:

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

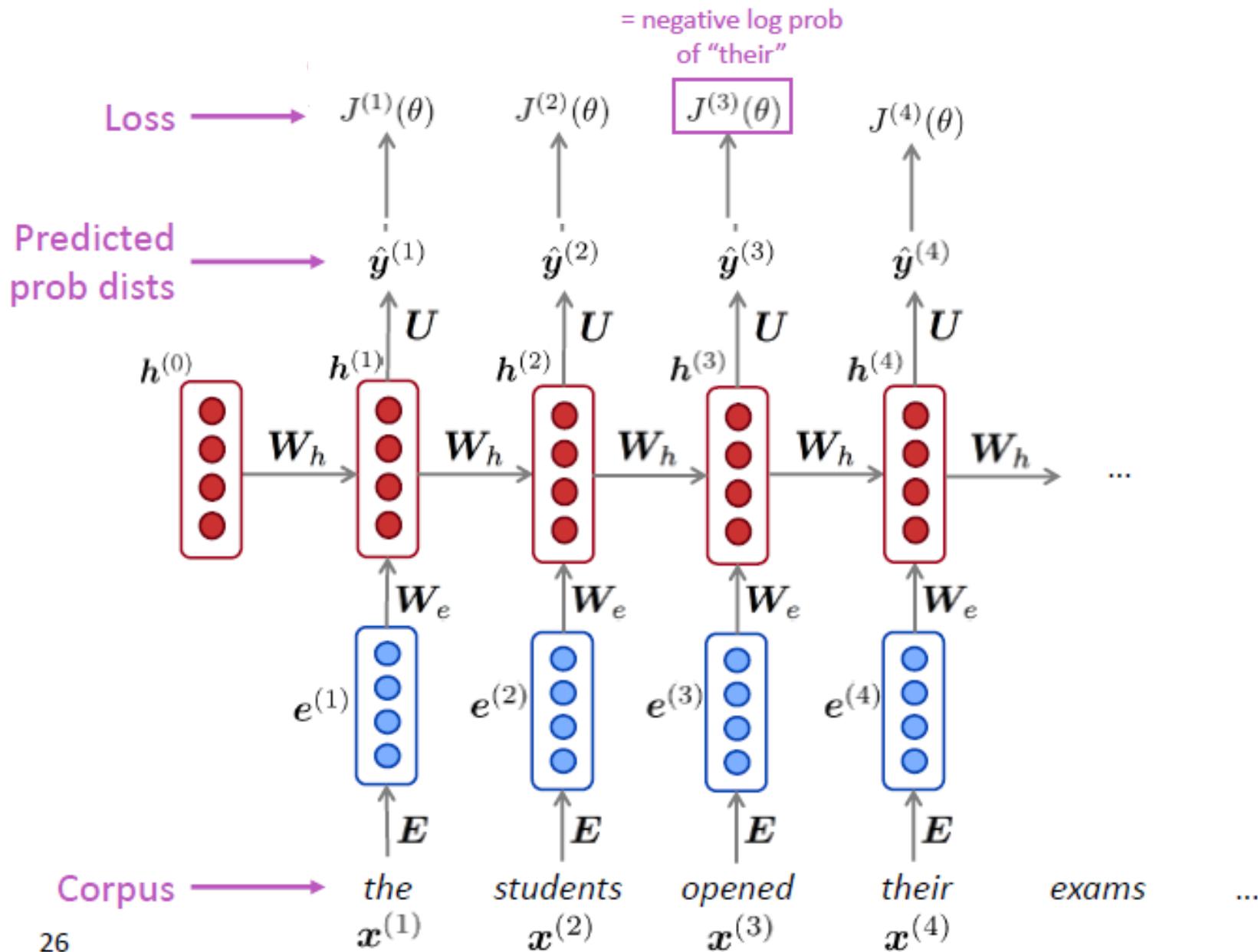
Training a RNN Language Model (Loss Function)



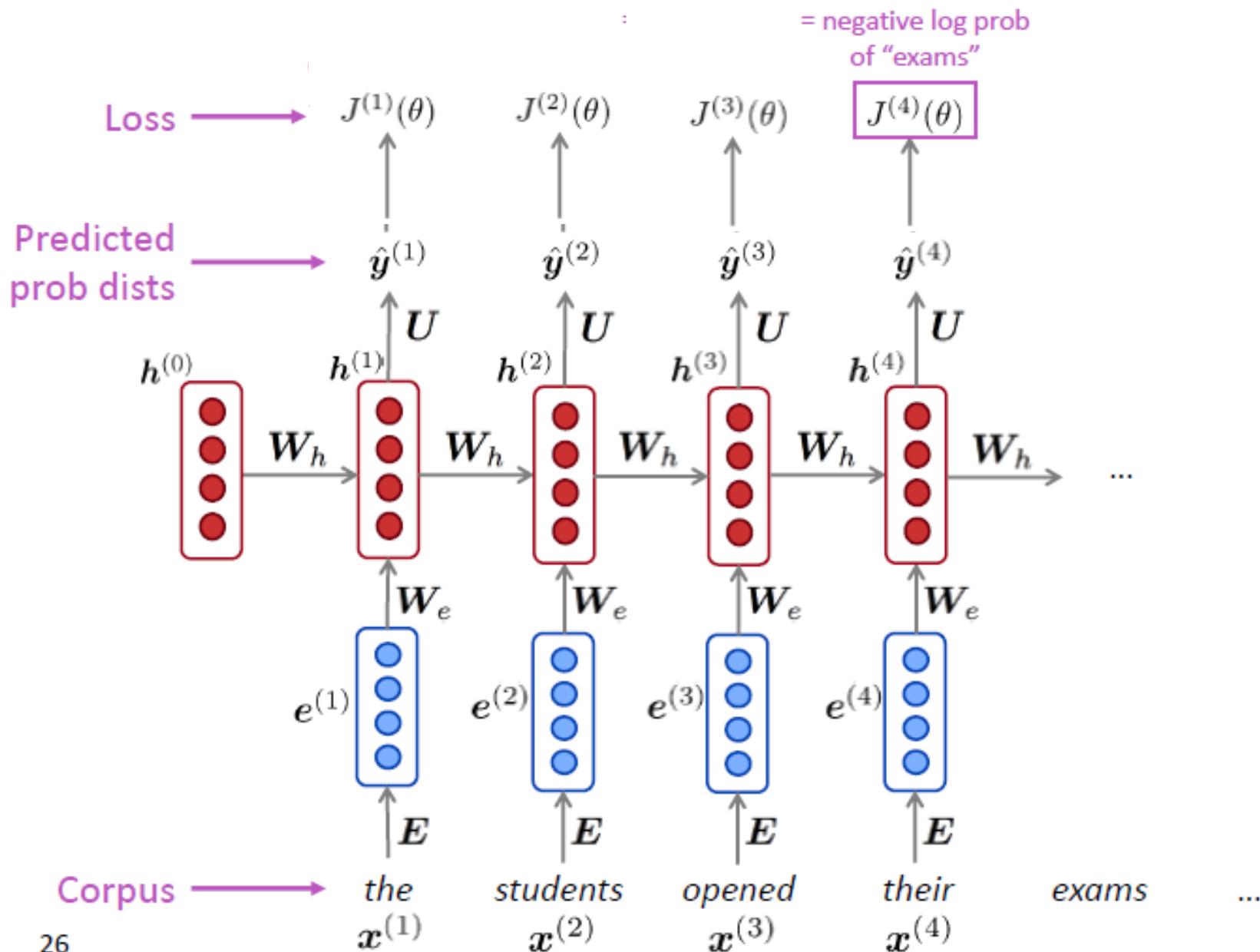
Training a RNN Language Model (Loss Function)



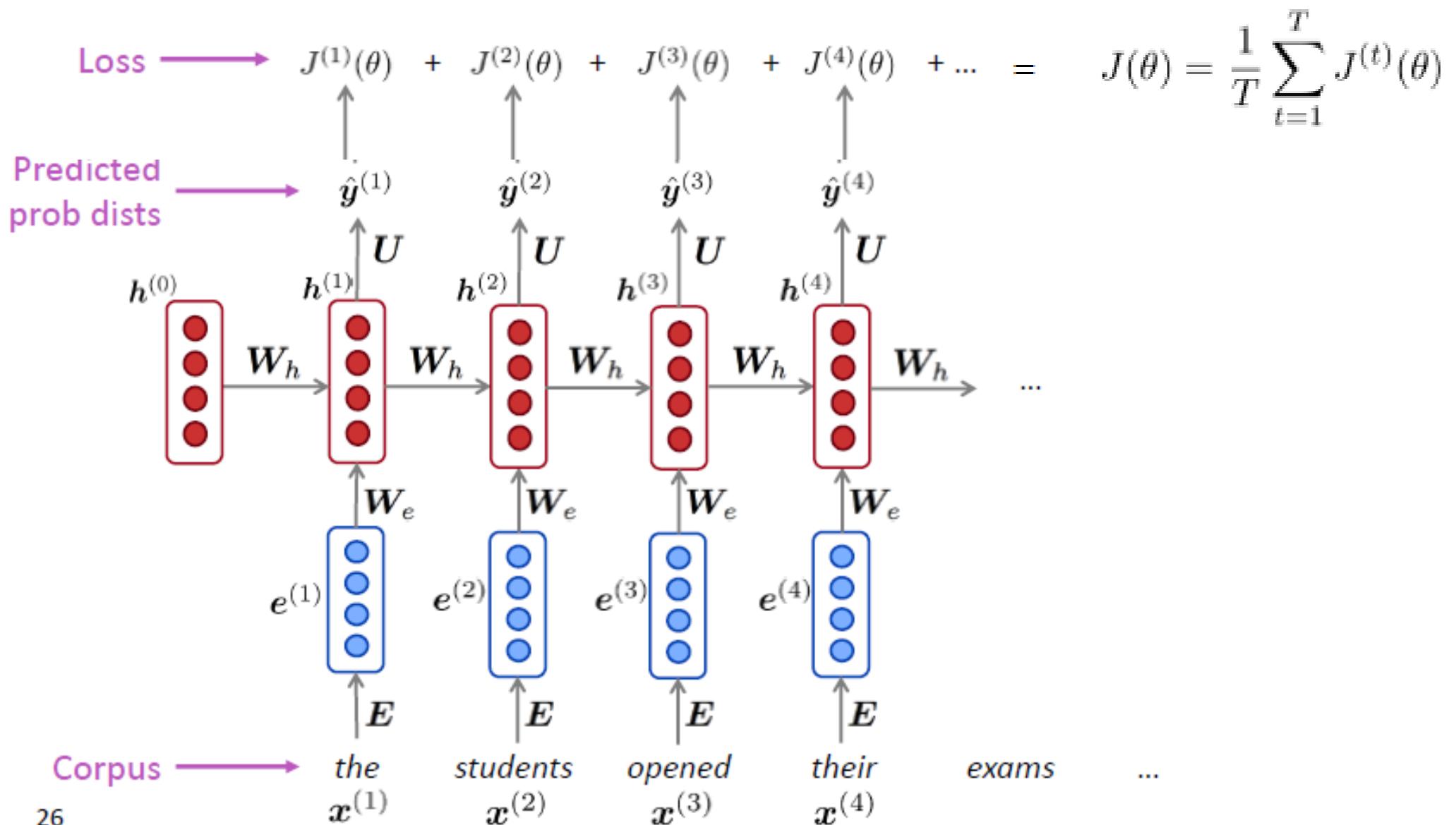
Training a RNN Language Model (Loss Function)



Training a RNN Language Model (Loss Function)



Training a RNN Language Model (Loss Function)



Training a RNN Language Model (Loss Function)

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

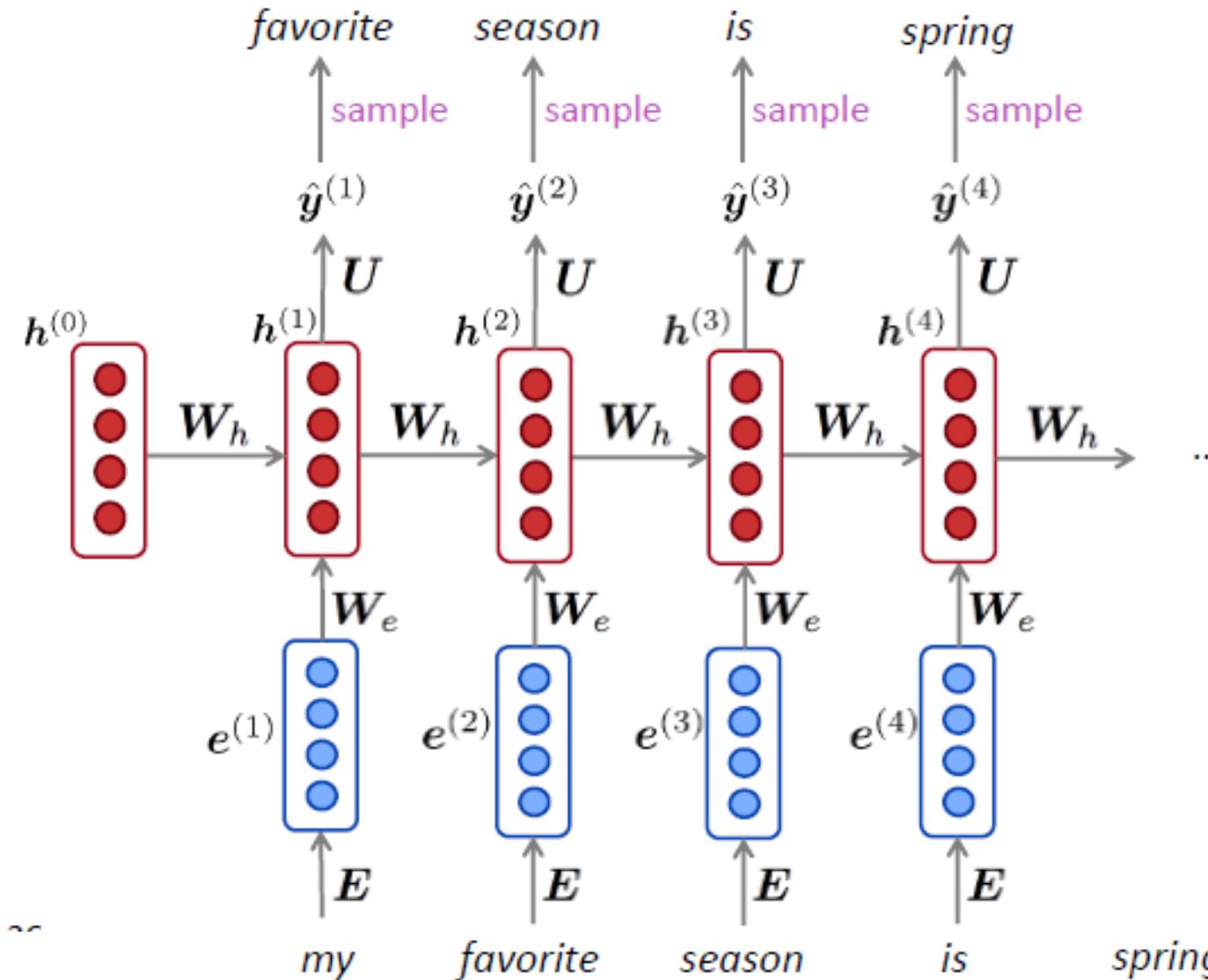
$$J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta)$$

- 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.

Generating Text with RNN

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.



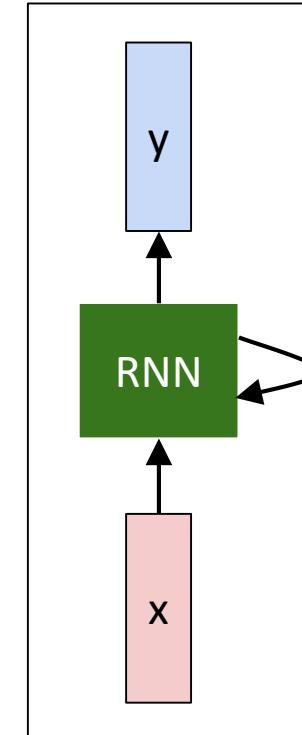
Sampling

Character	Probability
a	0.09
b	0.31
c	0.30
d	0.25
e	0.05

A Character based RNN Language Model

Vocabulary:
[h,e,l,o]

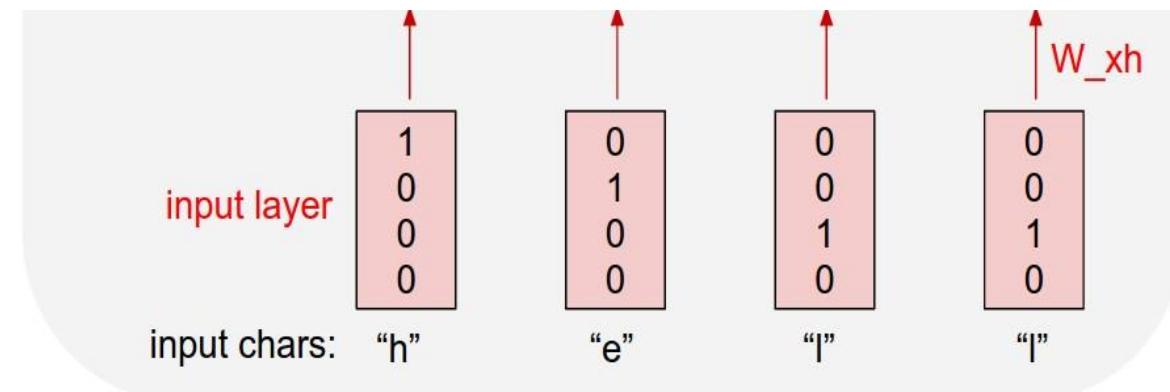
Example training
sequence:
“hello”



A Character based RNN Language Model

Vocabulary:
[h,e,l,o]

Example training
sequence:
“hello”

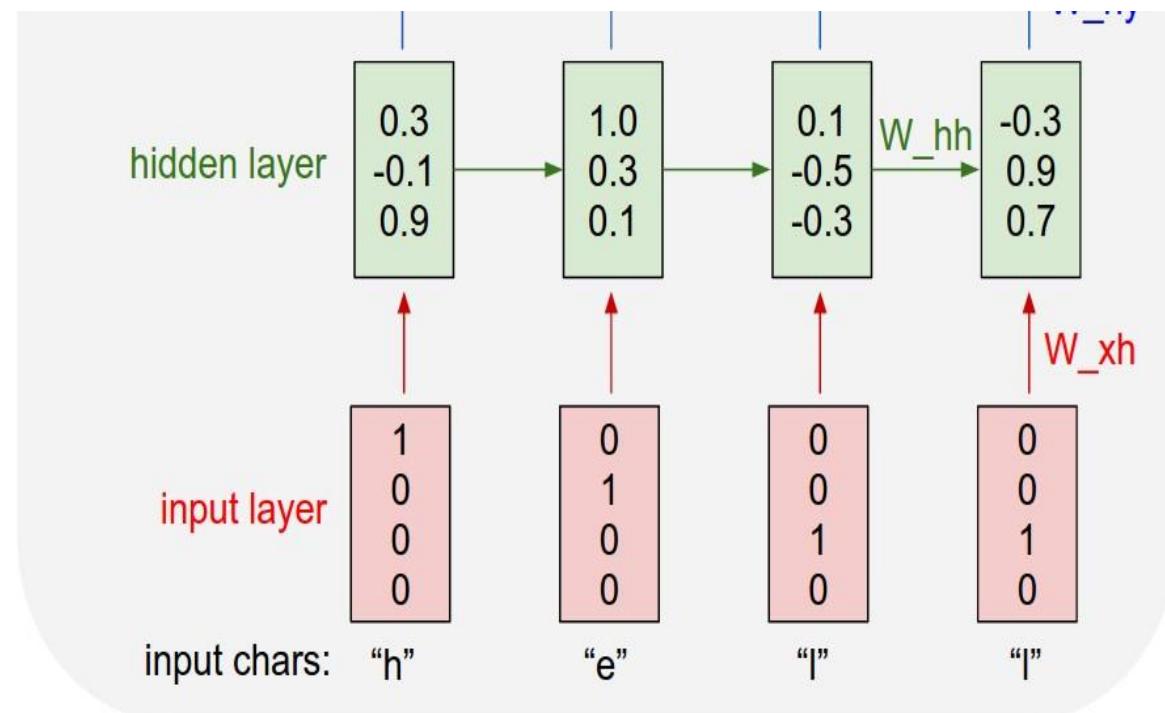


A Character based RNN Language Model

Vocabulary:
[h,e,l,o]

Example training
sequence:
“hello”

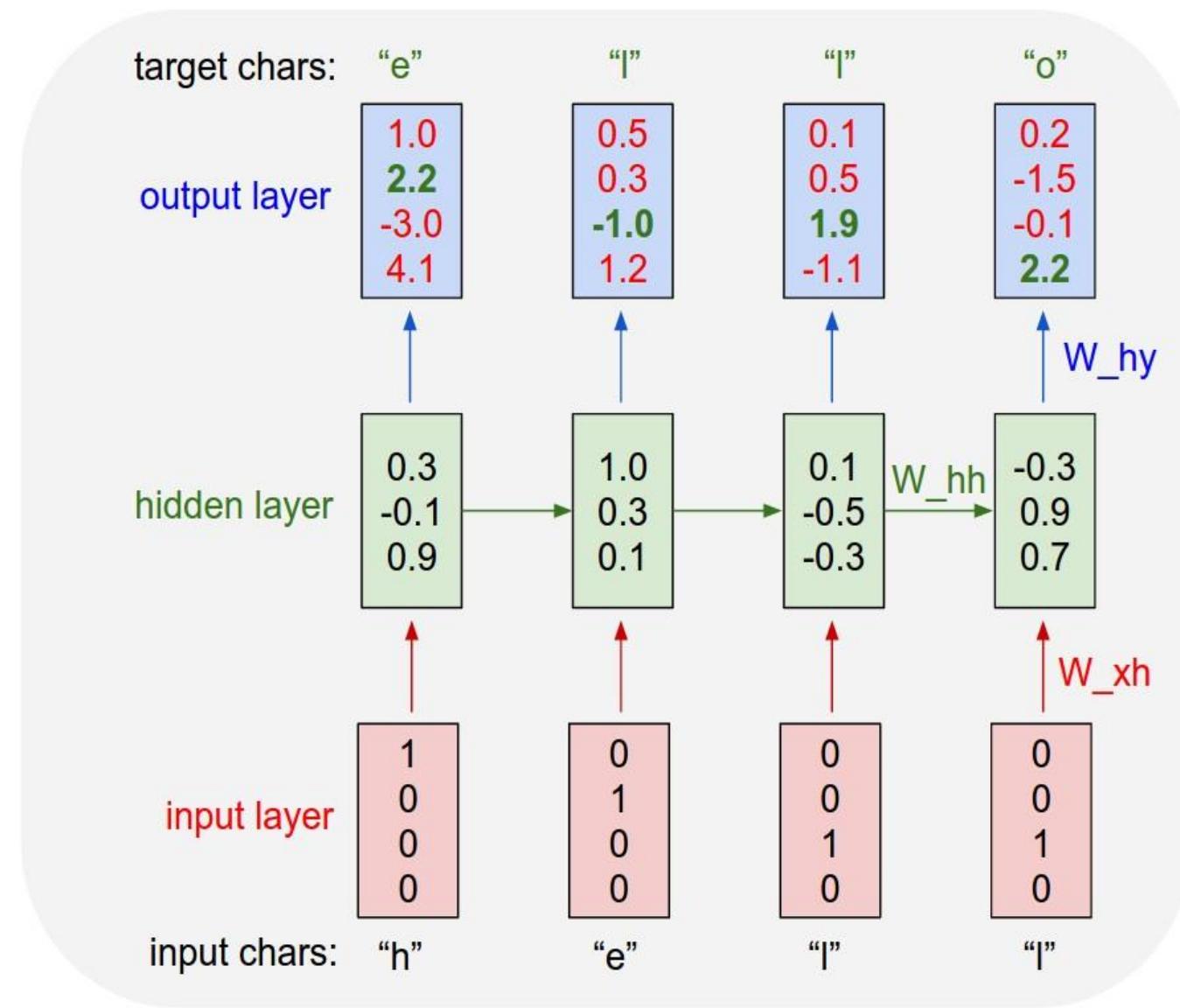
$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$



A Character based RNN Language Model

Vocabulary:
[h,e,l,o]

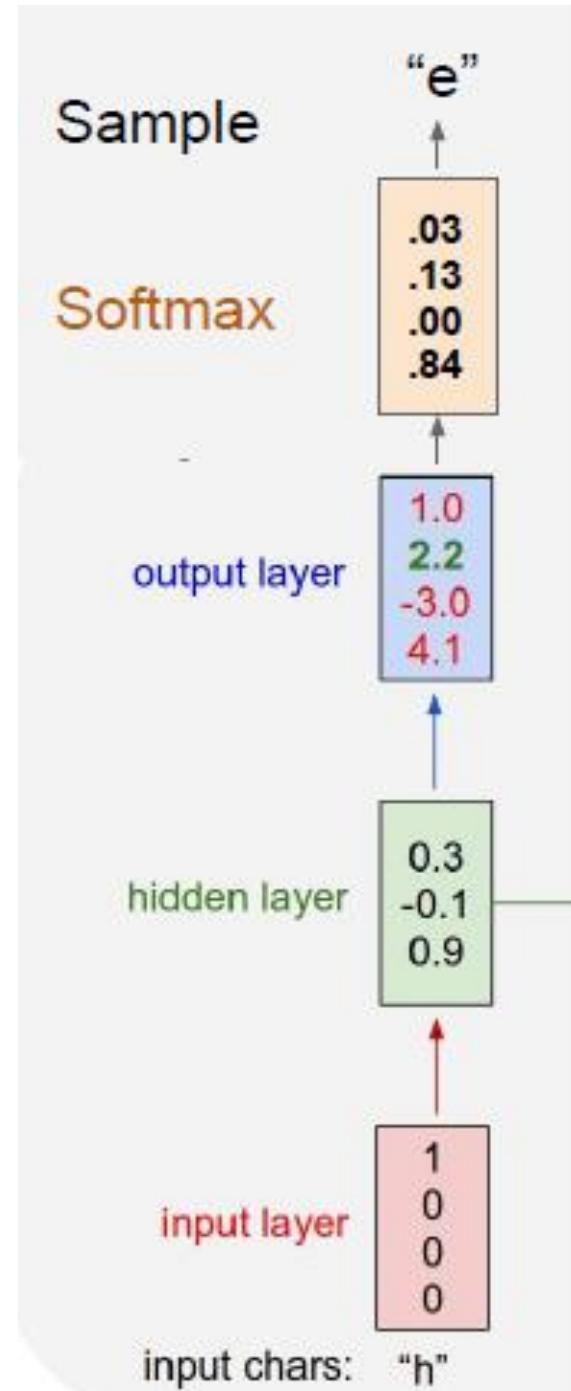
Example training
sequence:
“hello”



Example: Character-level Language Model Sampling

Vocabulary:
[h,e,l,o]

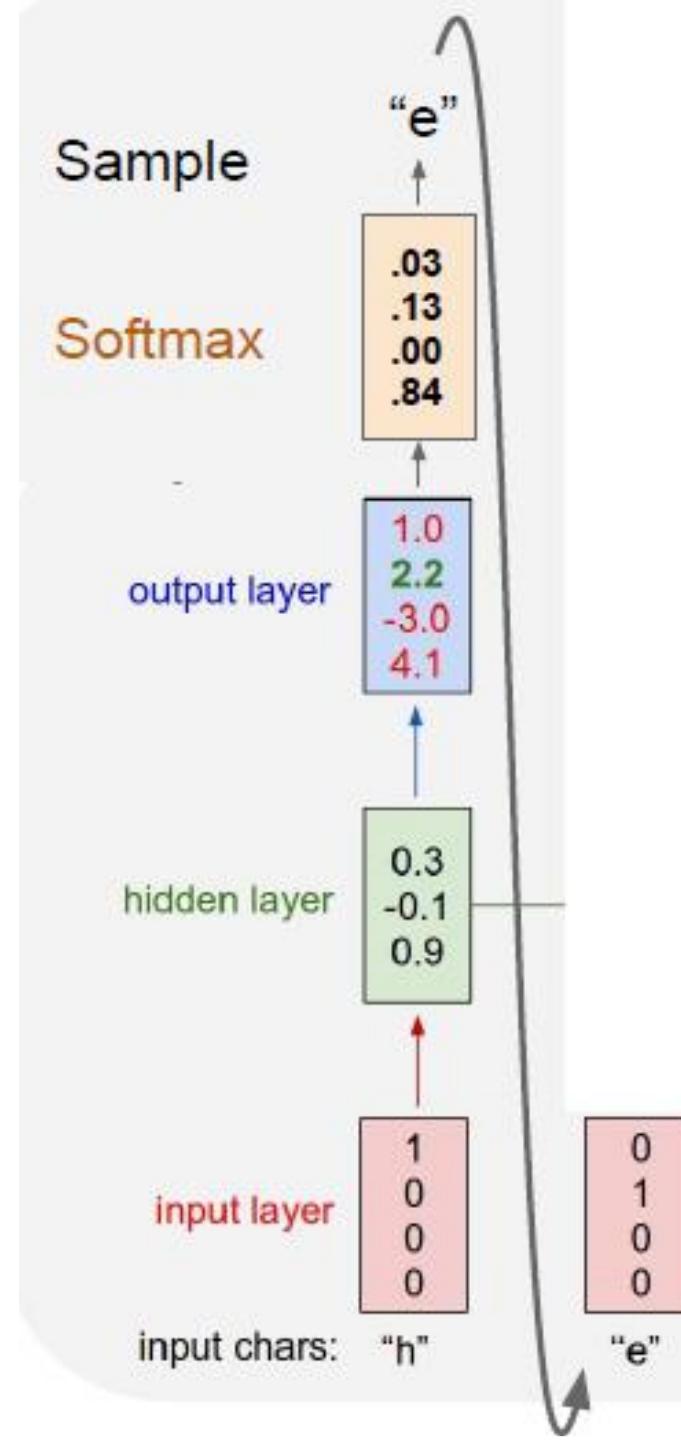
At test-time sample
characters one at a time,
feed back to model



Example: Character-level Language Model Sampling

Vocabulary:
[h,e,l,o]

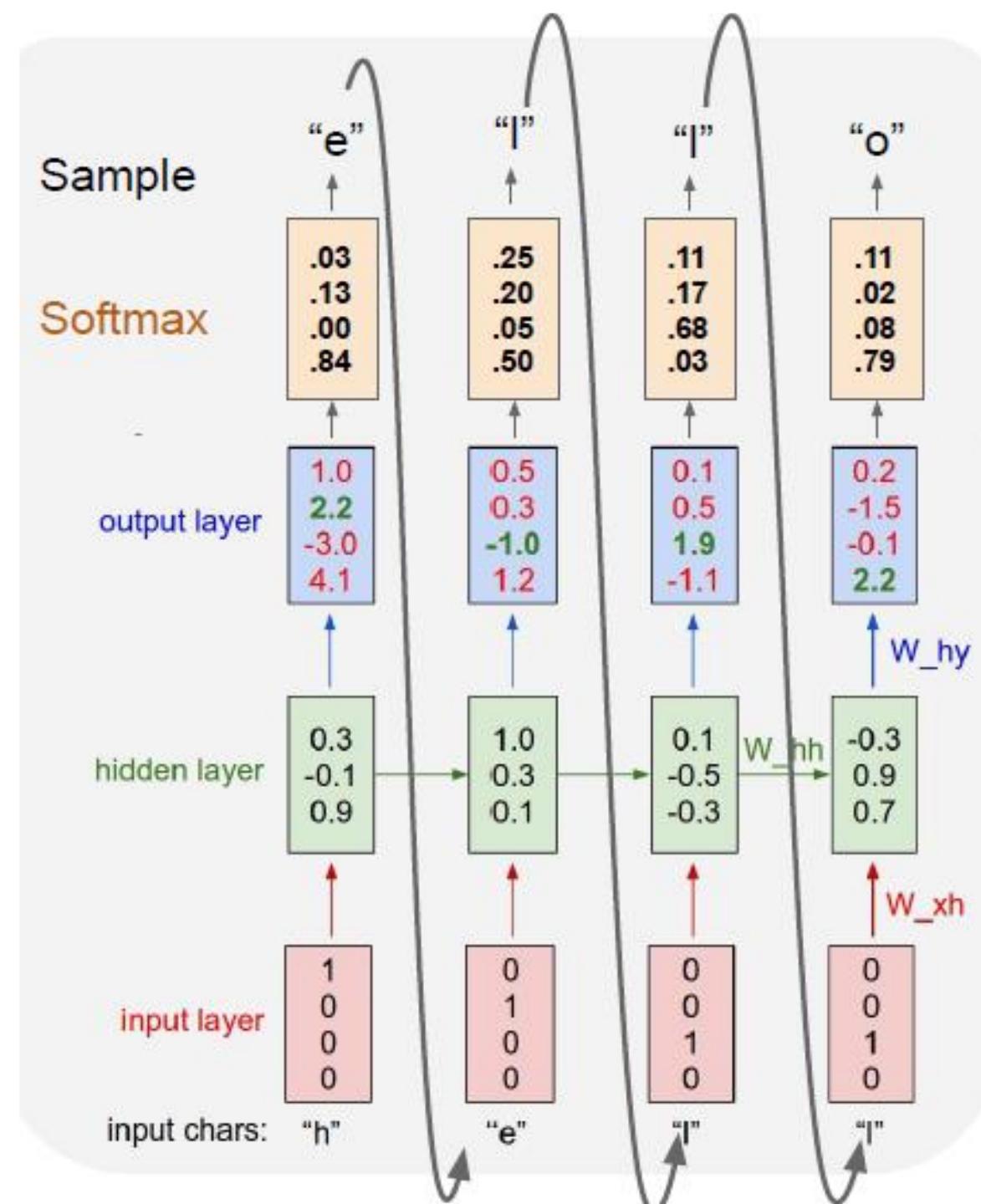
At test-time sample
characters one at a time,
feed back to model



Example: Character-level Language Model Sampling

Vocabulary:
[h,e,l,o]

At test-time sample
characters one at a time,
feed back to model



A Simple RNN code based on RNN equations is available at following link

min-char-rnn.py gist: 112 lines of Python

```
1  #!/usr/bin/python
2  #
3  # Minimal character-level Vanilla RNN model. Written by Andrej Karpathy (@karpathy)
4  #
5  # BSD License
6  #
7  # import numpy as np
8  #
9  #
10 # data size
11 data = open('input.txt', 'r').read() # should be simple plain text file
12 chars = list(data)
13 data_size, vocab_size = len(chars), len(chars)
14 print 'data has %d characters, %d unique.' % (data_size, vocab_size)
15 char_to_ix = { ch:i for i,ch in enumerate(chars) }
16 ix_to_char = { i:ch for i,ch in enumerate(chars) }
17
18 # hyperparameters
19 hidden_size = 100 # size of hidden layer of neurons
20 seq_length = 25 # number of steps to unroll the RNN for
21 learning_rate = 1e-1
22
23 # model parameters
24 Wxh = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
25 Whh = np.random.randn(hidden_size, hidden_size)*0.01 # hidden to hidden
26 why = np.random.randn(vocab_size, hidden_size)*0.01 # hidden to output
27 bh = np.zeros((hidden_size, 1)) # hidden bias
28 by = np.zeros((vocab_size, 1)) # output bias
29
30 def lossFun(inputs, targets, hprev):
31     """ 
32     inputs,targets are both list of integers.
33     hprev is ini array of initial hidden state
34     returns the loss, gradients in model parameters, and last hidden state
35     """
36     xx, hs, ys, ps = [], [], [], []
37     hprev = np.copy(hprev)
38     loss = 0
39     n = 0
40     for t in range(len(inputs)):
41         x = np.zeros((vocab_size, 1)) # encode in 1-of-K representation
42         x[inputs[t]] = 1
43         hs[t] = np.tanh(np.dot(Whh, hs[t-1]) + Wh) # hidden state
44         ys[t] = np.dot(why, hs[t]) + by # unnormalized log probabilities for next char
45         ps[t] = np.exp(ys[t]) / np.sum(np.exp(ys[t])) # probabilities for next char
46         loss += -np.log(ps[t][targets[t], 0]) # softmax (cross-entropy loss)
47         # backward pass: compute gradients going backwards
48         dWxh, dWhh, dWhy = np.zeros_like(Wxh), np.zeros_like(Whh), np.zeros_like(Why)
49         dWxh, dWhh, dWhy = np.zeros_like(Whh), np.zeros_like(by)
50         dnext = np.zeros_like(hs[0])
51         for t in reversed(range(len(inputs)))�
52             dy = np.copy(ps[t])
53             dy[targets[t]] -= 1 # backprop into y
54             dby += np.dot(dy, hs[t])
55             dy *= by
56             dh = np.dot(Why, T) * dy + dnext # backprop into h
57             dWxh += (x * dy).T
58             dWhh += (hs[t-1] * dy).T
59             dWhy += (1 - hs[t] * hs[t]).T
60             dnext = np.dot(Whh, T) * dy
61         for param in [dWxh, dWhh, dWhy, dWxh, dWhh, dWhy]:
62             np.clip(param, -5, 5, out=param) # clip to mitigate exploding gradients
63     return loss, dWxh, dWhh, dWhy, dWxh, dWhh, dWhy, hprev, n
```

```
64 def sample(h, seed_ix, n):
65     """ 
66     sample a sequence of integers from the model
67     h is memory state, seed_ix is seed letter for first time step
68     """
69     x = np.zeros((vocab_size, 1))
70     s[seed_ix] = 1
71     ixes = []
72     for t in range(n):
73         h = np.tanh(np.dot(Whh, h) + np.dot(Whh, h) + bh)
74         y = np.dot(Why, h) + by
75         p = np.exp(y) / np.sum(np.exp(y))
76         ix = np.random.choice(range(vocab_size), p=p, random())
77         x[ix] = 1
78         ixes.append(ix)
79     return ixes
80
81 n, p = 0, 0
82 nh0, nh1, why0 = np.zeros_like(Whh), np.zeros_like(Whh), np.zeros_like(Why)
83 nh1, nh2 = np.zeros_like(bh), np.zeros_like(by) # memory variables for adapted
84 smooth_loss = np.laplace(n/vocab_size)*seq_length # loss at iteration n
85 while True:
86     # process inputs (we're sweeping from left to right in steps seq_length long)
87     if n+seq_length > len(data) or n == 99:
88         hprev = np.zeros((hidden_size, 1)) # reset memmory
89         p = 0 # p = pr from start of run
90         inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
91         targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
92
93     # sample from the model now and then
94     if n % 100 == 0:
95         sample_ix = sample(hprev, inputs[0], 500)
96         txt = ''.join(ix_to_char[ix] for ix in sample_ix)
97         print '----\n' + txt + '\n----'
98
99     # forward seq_length characters through the net and fetch gradient
100    loss, dWxh, dWhh, dWhy, dWxh, dWhh, dWhy = lossFun(inputs, targets, hprev)
101    smooth_loss = smooth_loss * 0.999 + loss * 0.001
102    if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_loss) # print progress
103
104    # perform parameter update with Adam
105    for param, dparam, new in zip([Wxh, Whh, Why, bh, by],
106                                [dWxh, dWhh, dWhy, dbh, dby],
107                                [nh0, nh1, nh2, nh0, nh1]):
108        new += dparam ** 2
109        param += -learning_rate * dparam / np.sqrt(new + 1e-8) * adgrad_update
110
111    p += seq_length # move data pointer
112    n += 1 # iteration counter
```

(<https://gist.github.com/karpathy/d4dee566867f8291f086>)

(<https://gist.github.com/karpathy/d4dee566867f8291f086>)

Generating text with a RNN Language Model

- Let's have some fun!
- You can train a 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.

Source: <https://medium.com/@samim/obama-rnn-machine-generated-political-speeches-c8abd18a2ea0>

Generating text with a RNN Language Model

- Let's have some fun!
- You can train a RNN-LM on any kind of text, then generate text in that style.
- RNN-LM trained on *Harry Potter*:



“Sorry,” Harry shouted, panicking—“I’ll leave those brooms in London, are they?”

“No idea,” said Nearly Headless Nick, casting low close by Cedric, carrying the last bit of treacle Charms, from Harry’s shoulder, and to answer him the common room perched upon it, four arms held a shining knob from when the spider hadn’t felt it seemed. He reached the teams too.

Source: <https://medium.com/deep-writing/harry-potter-written-by-artificial-intelligence-8a9431803da6>

Generating text with a RNN Language Model

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

Title: CHOCOLATE RANCH BARBECUE

Categories: Game, Casseroles, Cookies, Cookies

Yield: 6 Servings

2 tb Parmesan cheese --- chopped

1 c Coconut milk

3 Eggs, beaten

Place each pasta over layers of lumps. Shape mixture into the moderate oven and simmer until firm. Serve hot in bodied fresh, mustard, orange and cheese.

Combine the cheese and salt together the dough in a large skillet; add the ingredients and stir in the chocolate and pepper.



Source: <https://gist.github.com/nylki/1efbaa36635956d35bcc>

RNN Trained on some novel

at first:

tyntd-iafhatawiaoihrdemot lytdws e ,fti, astai f ogoh eoase rrranbyne 'nhthnee e
plia tkldrgd t o idoe ns,smtt h ne etie h,hregtrs nigtike,aoaenngs lbg

train more

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

train more

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 offer.

train more

"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.

PANDARUS:

Alas, I think he shall be come approached and the day
When little strain would be attain'd into being never fed,
And who is but a chain and subjects of his death,
I should not sleep.

Second Senator:

They are away this miseries, produced upon my soul,
Breaking and strongly should be buried, when I perish
The earth and thoughts of many states.

DUKE VINCENTIO:

Well, your wit is in the care of side and that.

Second Lord:

They would be ruled after this chamber, and
my fair nues begun out of the fact, to be conveyed,
Whose noble souls I'll have the heart of the wars.

Clown:

Come, sir, I will make did behold your worship.

VIOLA:

I'll drink it.

VIOLA:

Why, Salisbury must find his flesh and thought
That which I am not aps, not a man and in fire,
To show the reining of the raven and the wars
To grace my hand reproach within, and not a fair are hand,
That Caesar and my goodly father's world;
When I was heaven of presence and our fleets,
We spare with hours, but cut thy council I am great,
Murdered and by thy master's ready there
My power to give thee but so much as hell:
Some service in the noble bondman here,
Would show him to her wine.

KING LEAR:

O, if you were a feeble sight, the courtesy of your law,
Your sight and several breath, will wear the gods
With his heads, and my hands are wonder'd at the deeds,
So drop upon your lordship's head, and your opinion
Shall be against your honour.



The Stacks Project

[home](#) [about](#) [tags explained](#) [tag lookup](#) [browse](#) [search](#) [bibliography](#) [recent comments](#) [blog](#) [add slogans](#)

Browse chapters

Part	Chapter	online	TeX source	view pdf
Preliminaries				
	1. Introduction	online	tex	pdf
	2. Conventions	online	tex	pdf
	3. Set Theory	online	tex	pdf
	4. Categories	online	tex	pdf
	5. Topology	online	tex	pdf
	6. Sheaves on Spaces	online	tex	pdf
	7. Sites and Sheaves	online	tex	pdf
	8. Stacks	online	tex	pdf
	9. Fields	online	tex	pdf
	10. Commutative Algebra	online	tex	pdf

Parts

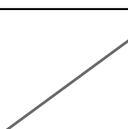
1. [Preliminaries](#)
2. [Schemes](#)
3. [Topics in Scheme Theory](#)
4. [Algebraic Spaces](#)
5. [Topics in Geometry](#)
6. [Deformation Theory](#)
7. [Algebraic Stacks](#)
8. [Miscellany](#)

Statistics

The Stacks project now consists of

- o 455910 lines of code
- o 14221 tags (56 inactive tags)
- o 2366 sections

Latex source



Algebra Equations and proofs using RNN

For $\bigoplus_{n=1,\dots,m} \mathcal{L}_{m,n} = 0$, hence we can find a closed subset \mathcal{H} in \mathcal{H} and any sets \mathcal{F} on X , U is a closed immersion of S , then $U \rightarrow T$ is a separated algebraic space.

Proof. Proof of (1). It also start we get

$$S = \text{Spec}(R) = U \times_X U \times_X U$$

and the comparicoly in the fibre product covering we have to prove the lemma generated by $\coprod Z \times_U U \rightarrow V$. Consider the maps M along the set of points Sch_{fppf} and $U \rightarrow U$ is the fibre category of S in U in Section ?? and the fact that any U affine, see Morphisms, Lemma ???. Hence we obtain a scheme S and any open subset $W \subset U$ in $\text{Sh}(G)$ such that $\text{Spec}(R') \rightarrow S$ is smooth or an

$$U = \bigcup U_i \times_{S_i} U_i$$

which has a nonzero morphism we may assume that f_i is of finite presentation over S . We claim that $\mathcal{O}_{X,x}$ is a scheme where $x, x', s'' \in S'$ such that $\mathcal{O}_{X,x'} \rightarrow \mathcal{O}'_{X',x'}$ is separated. By Algebra, Lemma ?? we can define a map of complexes $\text{GL}_S(x'/S'')$ and we win. \square

To prove study we see that $\mathcal{F}|_U$ is a covering of \mathcal{X}' , and \mathcal{T}_i is an object of $\mathcal{F}_{X/S}$ for $i > 0$ and \mathcal{F}_p exists and let \mathcal{F}_i be a presheaf of \mathcal{O}_X -modules on \mathcal{C} as a \mathcal{F} -module. In particular $\mathcal{F} = U/\mathcal{F}$ we have to show that

$$\widetilde{M}^\bullet = \mathcal{T}^\bullet \otimes_{\text{Spec}(k)} \mathcal{O}_{S,s} - i_X^{-1} \mathcal{F}$$

is a unique morphism of algebraic stacks. Note that

$$\text{Arrows} = (\text{Sch}/S)_{fppf}^{\text{opp}}, (\text{Sch}/S)_{fppf}$$

and

$$V = \Gamma(S, \mathcal{O}) \longrightarrow (U, \text{Spec}(A))$$

is an open subset of X . Thus U is affine. This is a continuous map of X is the inverse, the groupoid scheme S .

Proof. See discussion of sheaves of sets. \square

The result for prove any open covering follows from the less of Example ???. It may replace S by $X_{\text{spaces},\text{étale}}$ which gives an open subspace of X and T equal to S_{Zar} , see Descent, Lemma ???. Namely, by Lemma ?? we see that R is geometrically regular over S .

Lemma 0.1. Assume (3) and (3) by the construction in the description.

Suppose $X = \lim |X|$ (by the formal open covering X and a single map $\underline{\text{Proj}}_X(\mathcal{A}) = \text{Spec}(B)$ over U compatible with the complex

$$\text{Set}(\mathcal{A}) = \Gamma(X, \mathcal{O}_{X,\mathcal{O}_X}).$$

When in this case of to show that $\mathcal{Q} \rightarrow \mathcal{C}_{Z/X}$ is stable under the following result in the second conditions of (1), and (3). This finishes the proof. By Definition ?? (without element is when the closed subschemes are catenary. If T is surjective we may assume that T is connected with residue fields of S . Moreover there exists a closed subspace $Z \subset X$ of X where U in X' is proper (some defining as a closed subset of the uniqueness it suffices to check the fact that the following theorem

(1) f is locally of finite type. Since $S = \text{Spec}(R)$ and $Y = \text{Spec}(R)$.

Proof. This is form all sheaves of sheaves on X . But given a scheme U and a surjective étale morphism $U \rightarrow X$. Let $U \cap U = \coprod_{i=1,\dots,n} U_i$ be the scheme X over S at the schemes $X_i \rightarrow X$ and $U = \lim_i X_i$. \square

The following lemma surjective restrocomposes of this implies that $\mathcal{F}_{x_0} = \mathcal{F}_{x_0} = \mathcal{F}_{\mathcal{X},\dots,0}$.

Lemma 0.2. Let X be a locally Noetherian scheme over S , $E = \mathcal{F}_{X/S}$. Set $\mathcal{I} = \mathcal{J}_1 \subset \mathcal{I}'_n$. Since $\mathcal{I}'_n \subset \mathcal{I}^n$ are nonzero over $i_0 \leq p$ is a subset of $\mathcal{J}_{n,0} \circ \overline{A}_2$ works.

Lemma 0.3. In Situation ???. Hence we may assume $q' = 0$.

Proof. We will use the property we see that p is the next functor (??). On the other hand, by Lemma ?? we see that

$$D(\mathcal{O}_{X'}) = \mathcal{O}_X(D)$$

where K is an F -algebra where δ_{n+1} is a scheme over S . \square

Algebra Equations and proofs using RNN

Proof. Omitted. \square

Lemma 0.1. Let \mathcal{C} be a set of the construction.

Let \mathcal{C} be a gerber covering. Let \mathcal{F} be a quasi-coherent sheaves of \mathcal{O} -modules. We have to show that

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

Proof. This is an algebraic space with the composition of sheaves \mathcal{F} on $X_{\text{étale}}$ we have

$$\mathcal{O}_X(\mathcal{F}) = \{\text{morph}_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

where \mathcal{G} defines an isomorphism $\mathcal{F} \rightarrow \mathcal{G}$ of \mathcal{O} -modules. \square

Lemma 0.2. This is an integer \mathcal{Z} is injective.

Proof. See Spaces, Lemma ??.

Lemma 0.3. Let S be a scheme. Let X be a scheme and X is an affine open covering. Let $\mathcal{U} \subset \mathcal{X}$ be a canonical and locally of finite type. Let X be a scheme. Let X be a scheme which is equal to the formal complex.

The following to the construction of the lemma follows.

Let X be a scheme. Let X be a scheme covering. Let

$$b : X \rightarrow Y' \rightarrow Y \rightarrow Y \rightarrow Y' \times_X Y \rightarrow X.$$

be a morphism of algebraic spaces over S and Y .

Proof. Let X be a nonzero scheme of X . Let X be an algebraic space. Let \mathcal{F} be a quasi-coherent sheaf of \mathcal{O}_X -modules. The following are equivalent

- (1) \mathcal{F} is an algebraic space over S .
- (2) If X is an affine open covering.

Consider a common structure on X and X the functor $\mathcal{O}_X(U)$ which is locally of finite type. \square

This since $\mathcal{F} \in \mathcal{F}$ and $x \in \mathcal{G}$ the diagram

$$\begin{array}{ccccc}
 S & \xrightarrow{\quad} & & & \\
 \downarrow & & & & \\
 \xi & \longrightarrow & \mathcal{O}_{X'} & & \\
 \text{gor}_s & & \uparrow & \searrow & \\
 & & =\alpha' & \longrightarrow & \\
 & & \uparrow & & \\
 & & =\alpha' & \longrightarrow & \alpha \\
 & & & & \\
 \text{Spec}(K_\psi) & & \text{Mor}_{\text{Sets}} & & d(\mathcal{O}_{X_{/\kappa}}, \mathcal{G}) \\
 & & & & \\
 & & X & &
 \end{array}$$

is a limit. Then \mathcal{G} is a finite type and assume S is a flat and \mathcal{F} and \mathcal{G} is a finite type f_* . This is of finite type diagrams, and

- the composition of \mathcal{G} is a regular sequence,
- $\mathcal{O}_{X'}$ is a sheaf of rings.

Proof. We have see that $X = \text{Spec}(R)$ and \mathcal{F} is a finite type representable by algebraic space. The property \mathcal{F} is a finite morphism of algebraic stacks. Then the cohomology of X is an open neighbourhood of U . \square

Proof. This is clear that \mathcal{G} is a finite presentation, see Lemmas ??.

A reduced above we conclude that U is an open covering of \mathcal{C} . The functor \mathcal{F} is a “field”

$$\mathcal{O}_{X,x} \longrightarrow \mathcal{F}_{\overline{x}} \dashleftarrow (\mathcal{O}_{X_{\text{étale}}}) \longrightarrow \mathcal{O}_{X_\ell}^{-1} \mathcal{O}_{X_\lambda}(\mathcal{O}_{X_\eta}^{\overline{y}})$$

is an isomorphism of covering of \mathcal{O}_{X_ℓ} . If \mathcal{F} is the unique element of \mathcal{F} such that X is an isomorphism.

The property \mathcal{F} is a disjoint union of Proposition ?? and we can filtered set of presentations of a scheme \mathcal{O}_X -algebra with \mathcal{F} are opens of finite type over S .

If \mathcal{F} is a scheme theoretic image points. \square

If \mathcal{F} is a finite direct sum \mathcal{O}_{X_λ} is a closed immersion, see Lemma ?? . This is a sequence of \mathcal{F} is a similar morphism.

C Code generation using RNN

This repository Search Explore Gist Blog Help karpathy + ⌂ ⚙ ⌂

torvalds / linux Watch 3,711 Star 23,054 Fork 9,141

Linux kernel source tree

520,037 commits 1 branch 420 releases 5,039 contributors

branch: master linux / +

Merge branch 'drm-fixes' of git://people.freedesktop.org/~airlied/linux ...
torvalds authored 9 hours ago latest commit 4b1706927d

Category	Commit Message	Date
Documentation	Merge git://git.kernel.org/pub/scm/linux/kernel/git/nab/target-pending	6 days ago
arch	Merge branch 'x86-urgent-for-linus' of git://git.kernel.org/pub/scm/l...	a day ago
block	block: discard bdi_unregister() in favour of bdi_destroy()	9 days ago
crypto	Merge git://git.kernel.org/pub/scm/linux/kernel/git/herbert/crypto-2.6	10 days ago
drivers	Merge branch 'drm-fixes' of git://people.freedesktop.org/~airlied/linux	9 hours ago
firmware	firmware/hex2fw.c: restore missing default in switch statement	2 months ago
fs	vfs: read file_handle only once in handle_to_path	4 days ago
include	Merge branch 'perf-urgent-for-linus' of git://git.kernel.org/pub/scm/...	a day ago
init	init: fix regression by supporting devices with major:minor:offset fo...	a month ago
iosc	iosc: kernels now bound to virtio-unit kernel context due to driver ref...	a month ago

Code 74 Pull requests

Pulse Graphs

HTTPS clone URL <https://github.com/torvalds/linux>
You can clone with [HTTPS](#), [SSH](#), or [Subversion](#).

Clone in Desktop Download ZIP

Generated C code

```
static void do_command(struct seq_file *m, void *v)
{
    int column = 32 << (cmd[2] & 0x80);
    if (state)
        cmd = (int)(int_state ^ (in_8(&ch->ch_flags) & Cmd) ? 2 : 1);
    else
        seq = 1;
    for (i = 0; i < 16; i++) {
        if (k & (1 << 1))
            pipe = (in_use & UMXTHREAD_UNCCA) +
                ((count & 0x00000000fffff8) & 0x0000000f) << 8;
        if (count == 0)
            sub(pid, ppc_md.kexec_handle, 0x20000000);
        pipe_set_bytes(i, 0);
    }
    /* Free our user pages pointer to place camera if all dash */
    subsystem_info = &of_changes[PAGE_SIZE];
    rek_controls(offset, idx, &offset);
    /* Now we want to deliberately put it to device */
    control_check_polarity(&context, val, 0);
    for (i = 0; i < COUNTER; i++)
        seq_puts(s, "policy ");
}
```

Generated C code

```
/*
 * Copyright (c) 2006-2010, Intel Mobile Communications. All rights reserved.
 *
 * This program is free software; you can redistribute it and/or modify it
 * under the terms of the GNU General Public License version 2 as published by
 * the Free Software Foundation.
 *
 * This program is distributed in the hope that it will be useful,
 * but WITHOUT ANY WARRANTY; without even the implied warranty of
 * MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
 *
 * GNU General Public License for more details.
 *
 * You should have received a copy of the GNU General Public License
 * along with this program; if not, write to the Free Software Foundation,
 * Inc., 675 Mass Ave, Cambridge, MA 02139, USA.
 */

#include <linux/kexec.h>
#include <linux/errno.h>
#include <linux/io.h>
#include <linux/platform_device.h>
#include <linux/multi.h>
#include <linux/ckevent.h>

#include <asm/io.h>
#include <asm/prom.h>
#include <asm/e820.h>
#include <asm/system_info.h>
#include <asm/seteew.h>
#include <asm/pgproto.h>
```

Generated C code

```
#include <asm/io.h>
#include <asm/prom.h>
#include <asm/e820.h>
#include <asm/system_info.h>
#include <asm/seteew.h>
#include <asm/pgproto.h>

#define REG_PG    vesa_slot_addr_pack
#define PFM_NOCOMP AFSR(0, load)
#define STACK_DDR(type)      (func)

#define SWAP_ALLOCATE(nr)      (e)
#define emulate_sigs() arch_get_unaligned_child()
#define access_rw(TST)  asm volatile("movd %%esp, %0, %3" : : "r" (0)); \
if (__type & DO_READ)

static void stat_PC_SEC __read_mostly offsetof(struct seq_argsqueue, \
pC>[1]);

static void
os_prefix(unsigned long sys)
{
#ifdef CONFIG_PREEMPT
    PUT_PARAM_RAID(2, sel) = get_state_state();
    set_pid_sum((unsigned long)state, current_state_str(),
                (unsigned long)-1->lr_full; low;
}

```

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

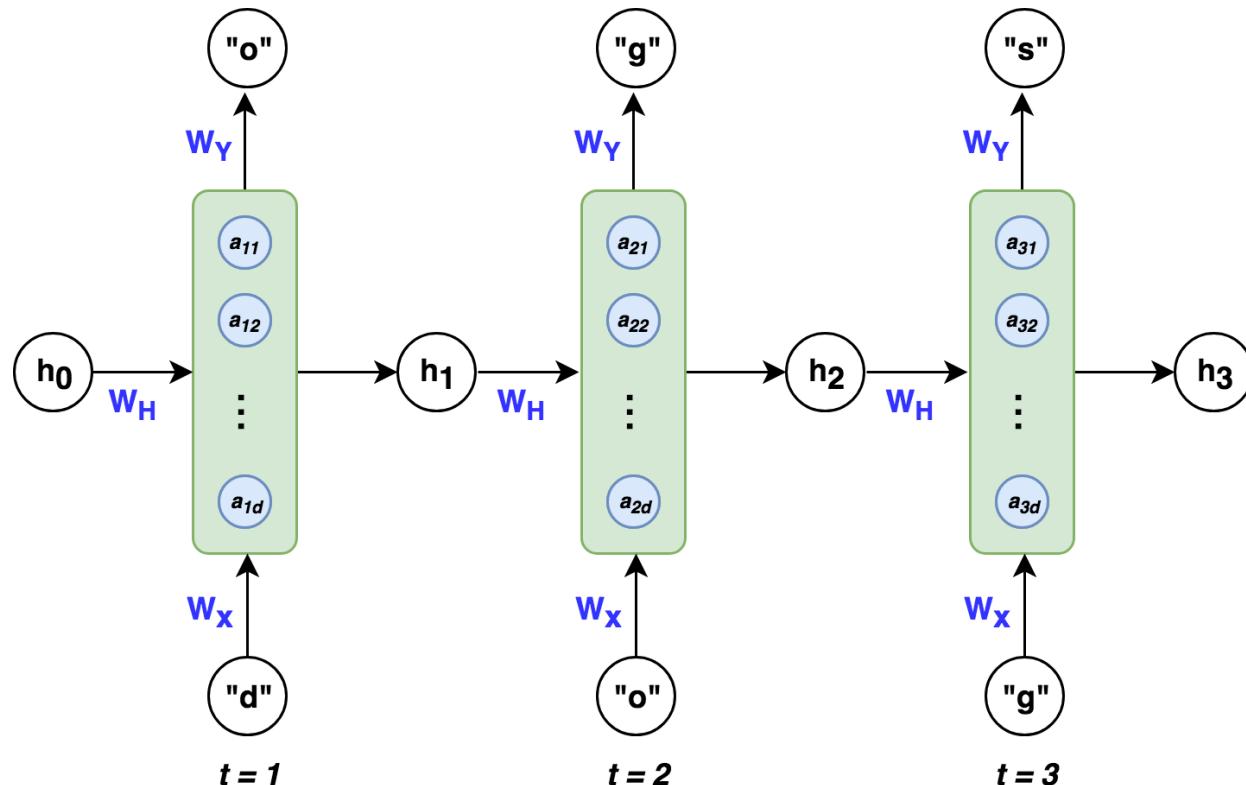
Perplexity improves
(lower is better)

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.

RNN Back Propagation

RNN Forward Pass



Output
From Hidden
State

Prediction at
time t

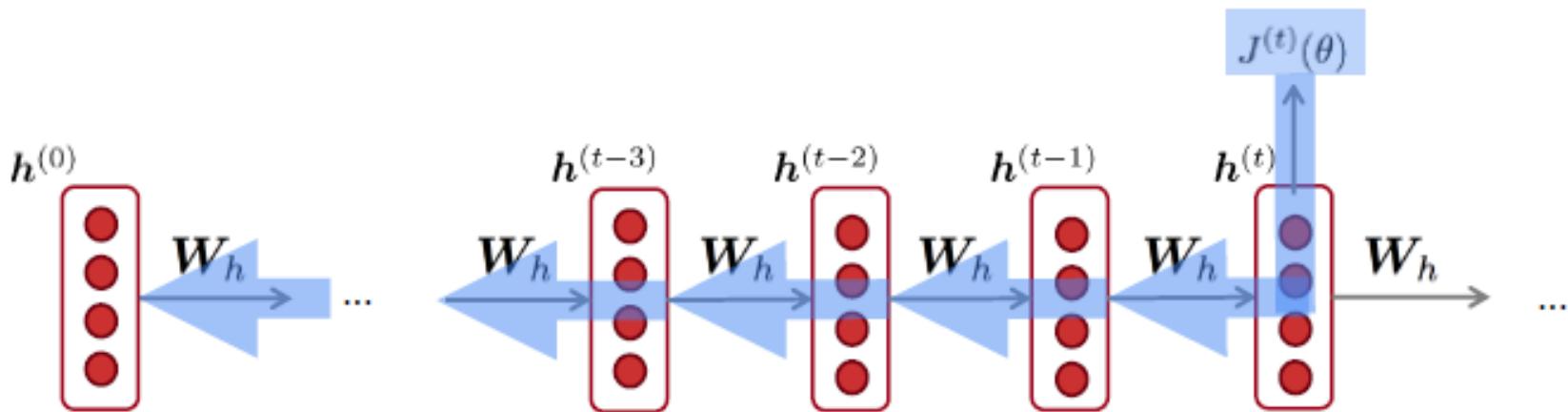
$$a_t = W_H h_{t-1} + W_X X_t \quad \text{Hidden Nodes}$$

Activation Function

$$\rightarrow h_t = \underbrace{\tanh(a_t)}_{\text{Hidden State}}$$

$$\rightarrow y_t = \text{softmax}(W_Y h_t)$$

Backpropagation for RNNs



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

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”

Backpropagation for RNNs

- 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:

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

Derivative of composition function

RNN can model long term dependency

The brown and black dog, which was playing with the cat, was a german shepherd.

(x_2)

(x_4)

(x_5)

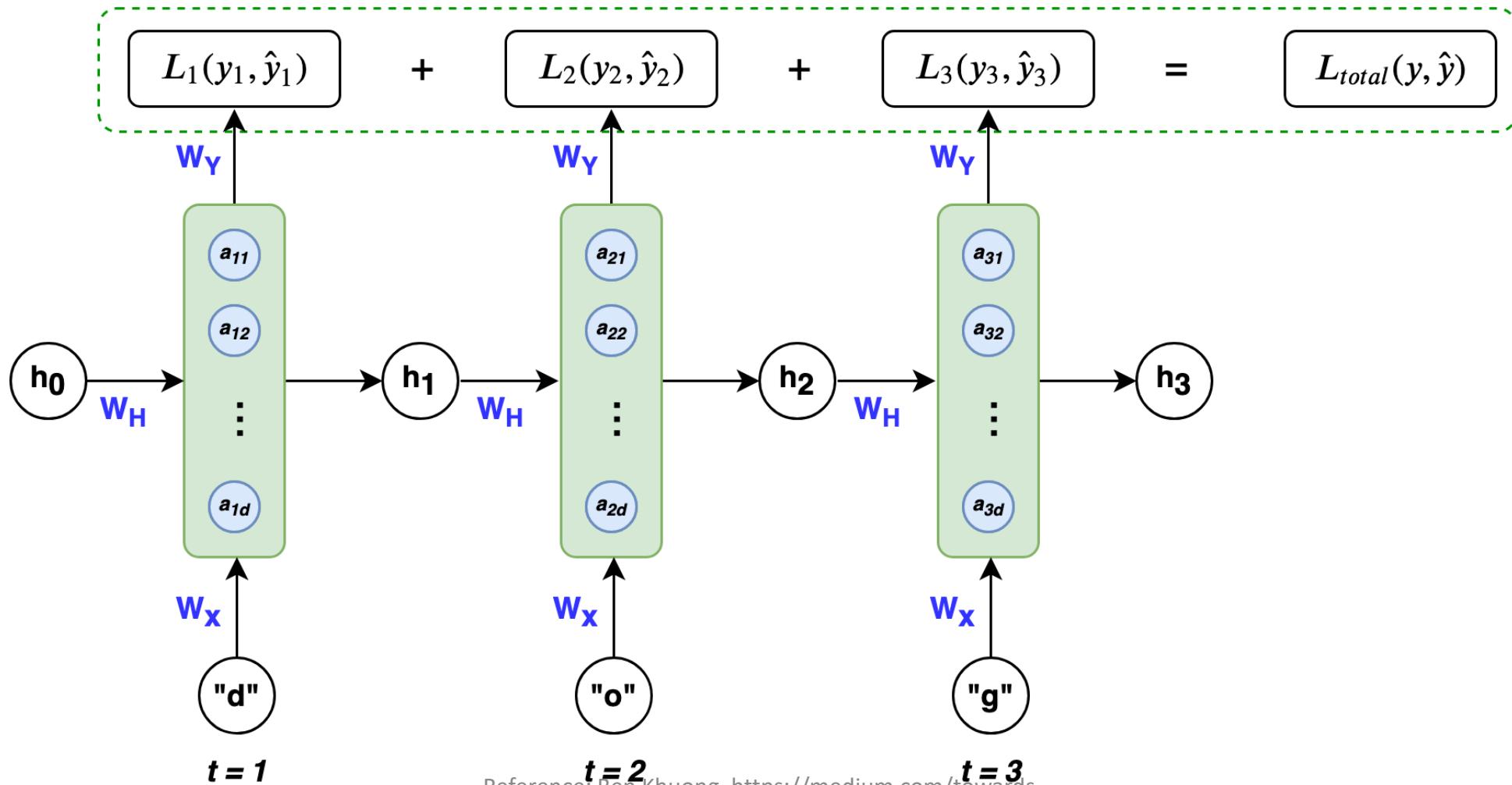
(x_{14})

(x_{15})

Backpropagation through time

1. Initialize weight matrices Wx , Wy , Wh randomly
2. Forward propagation to compute predictions
3. Compute the loss
4. Backpropagation to compute gradients
5. Update weights based on gradients
6. Repeat steps 2–5

Total Loss



Reference: Ben Khuong, <https://medium.com/towards-artificial-intelligence/whirlwind-tour-of-rnns-a11effb7808f>

Weight Update

$$W_i := W_i - \eta \frac{\partial L_{total}(y, \hat{y})}{\partial W_i}$$

where i = x, y, and h as a shorthand for the 3 weight matrices

The main idea is chain rule and to account for the loss at each time step

Function dependencies with respect to W_Y

$$L_t = -y_t \log(\hat{y}_t) \rightarrow \hat{y}_t = \text{softmax}(z_t) \rightarrow z_t = W_Y h_t$$

Chain rule with respect to W_Y

$$\frac{\partial L_t}{\partial W_Y} = \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial z_t} \frac{\partial z_t}{\partial W_Y}$$

Gradient for W_Y

$$L_{total}(y, \hat{y}) = - \sum_{t=1}^n y_t \log(\hat{y}_t)$$

$$\frac{\partial L_{total}}{\partial W_Y} = \frac{\partial L_1}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial z_1} \frac{\partial z_1}{\partial W_Y} + \frac{\partial L_2}{\partial \hat{y}_2} \frac{\partial \hat{y}_2}{\partial z_2} \frac{\partial z_2}{\partial W_Y} + \frac{\partial L_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial z_3} \frac{\partial z_3}{\partial W_Y} = \sum_{t=1}^n \frac{\partial L_t}{\partial W_Y}$$

Function Dependencies with respect to W_X

$$L_t = -y_t \log(\hat{y}_t) \rightarrow \hat{y}_t = \text{softmax}(z_t) \rightarrow z_t = W_Y h_t \rightarrow h_t = \tanh(W_H h_{t-1} + W_X X_t) \quad (1)$$

Chain rule with respect to W_x

$\frac{\partial L_t}{\partial W_x} = \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial z_t} \frac{\partial z_t}{\partial h_t} \frac{\partial h_t}{\partial W_x}$ but note, that within h_t , h_{t-1} also contains W_x , thus we need to recursively chain rule h_{t-1} until we reach h_0

$$\frac{\partial L_t}{\partial W_x} = \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial z_t} \frac{\partial z_t}{\partial h_t} \frac{\partial h_t}{\partial W_x} + \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial z_t} \frac{\partial z_t}{\partial h_t} \frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial W_x} + \dots + \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial z_t} \frac{\partial z_t}{\partial h_t} \frac{\partial h_t}{\partial h_{t-n}} \frac{\partial h_{t-n}}{\partial W_x} = \sum_{k=0}^n \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial z_t} \frac{\partial z_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial W_x} \quad (3)$$

Gradient for W_X

$$\frac{\partial L_{total}}{\partial W_X} = \sum_{t=1}^n \sum_{k=0}^n \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial z_t} \frac{\partial z_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial W_x} \quad (4)$$

Task 1

- Write partial derivative equations using chain rule for W_h . They will be similar to equations for W_x on previous slide. Submit on Google classroom

Reading

- Chapter 9, Speech and Language Processing. Daniel Jurafsky & James H. Martin. Third edition

<https://web.stanford.edu/~jurafsky/slp3/9.pdf>