

Natural Language Processing

AI51701/CSE71001

Lecture 9

10/10/2023

Instructor: Taehwan Kim

Announcements

- We have a substitute class on Oct. 17 at 4:00pm

Announcements

❑ Final project proposal

- Presentation on 10/12: ~6 min. for presentation and 1 min. for QnA

	10/12
1	Jihyoung Jang, Eunchan Lee, Juhyeong Lee
2	Jaewon Yang, Sungho Jeon
3	Gaurav Saha, Youngbin Ki, Eldor Fozilov, Rogelio Ruzcko Tobias
4	Dahee Kim, Song Kim
5	Sangjune Park, Jongsung Lee
6	Namwoo Kwon, Gawon Choi
7	Jinwoo Lee, Seongouk Kim, Meraj Mammadov
8	Taegyeong Lee, Taesoo Kim, Sergey Pyatkovskiy
9	Jaemu Heo, Hyunmin Song
10	Ryskul Tagmanova

Announcements

- ❑ Final project proposal
 - Presentation on 10/12: 6 min. for presentation and 1 min. for QnA
 - Written proposal: 3–4 pages (in latex, *EMNLP 2023* format)
 - Submit **both written proposal and presentation slides** via BlackBoard (Assignments>Final project proposal),
due Oct. 11, 11:59pm

Language Modeling + RNNs

Language Modeling

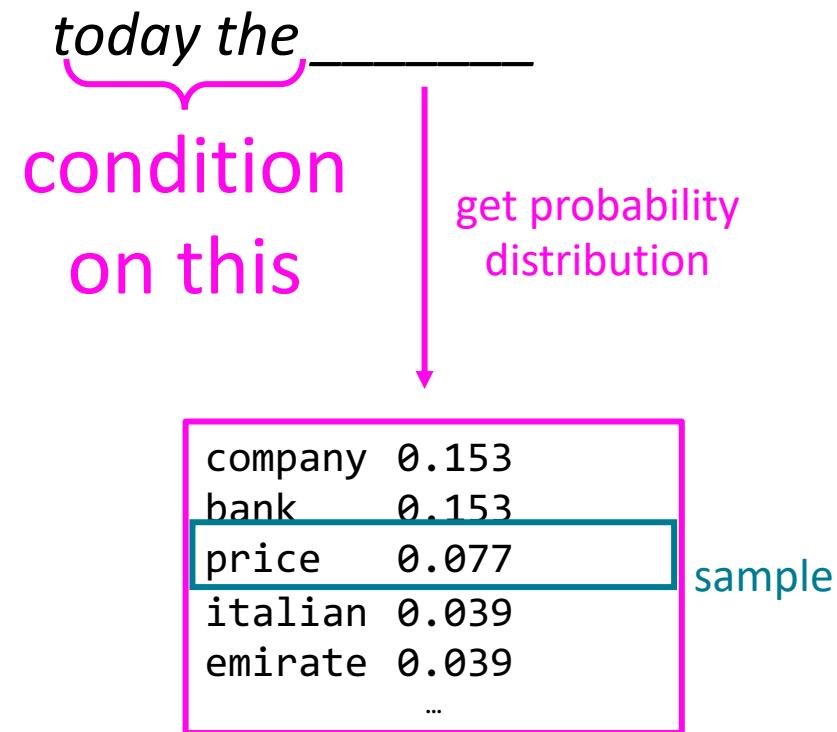
- ❑ **Language Modeling** is the task of predicting what word comes next

$$P(w_k \mid w_1, w_2, \dots, w_{k-1})$$

- ❑ A system that compute the probability distribution of the next word is called a **Language Model**

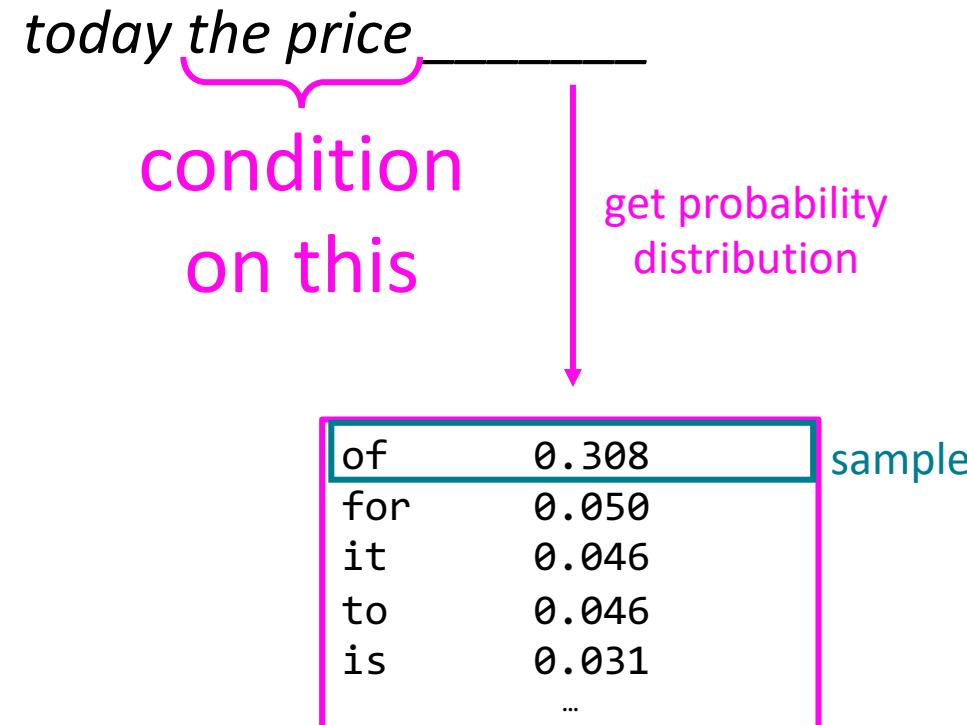
Generating text with a n-gram Language Model

- You can use a Language Model to generate text



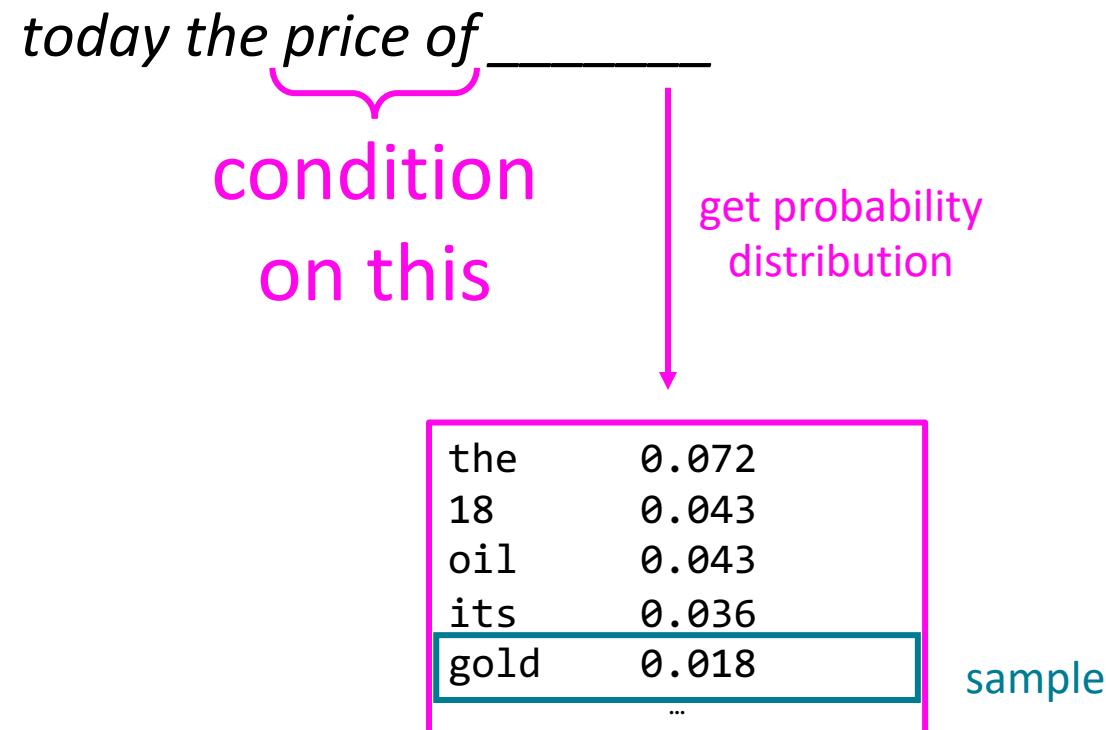
Generating text with a n-gram Language Model

- You can use a Language Model to generate text



Generating text with a n-gram Language Model

- You can use a Language Model to generate text



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

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

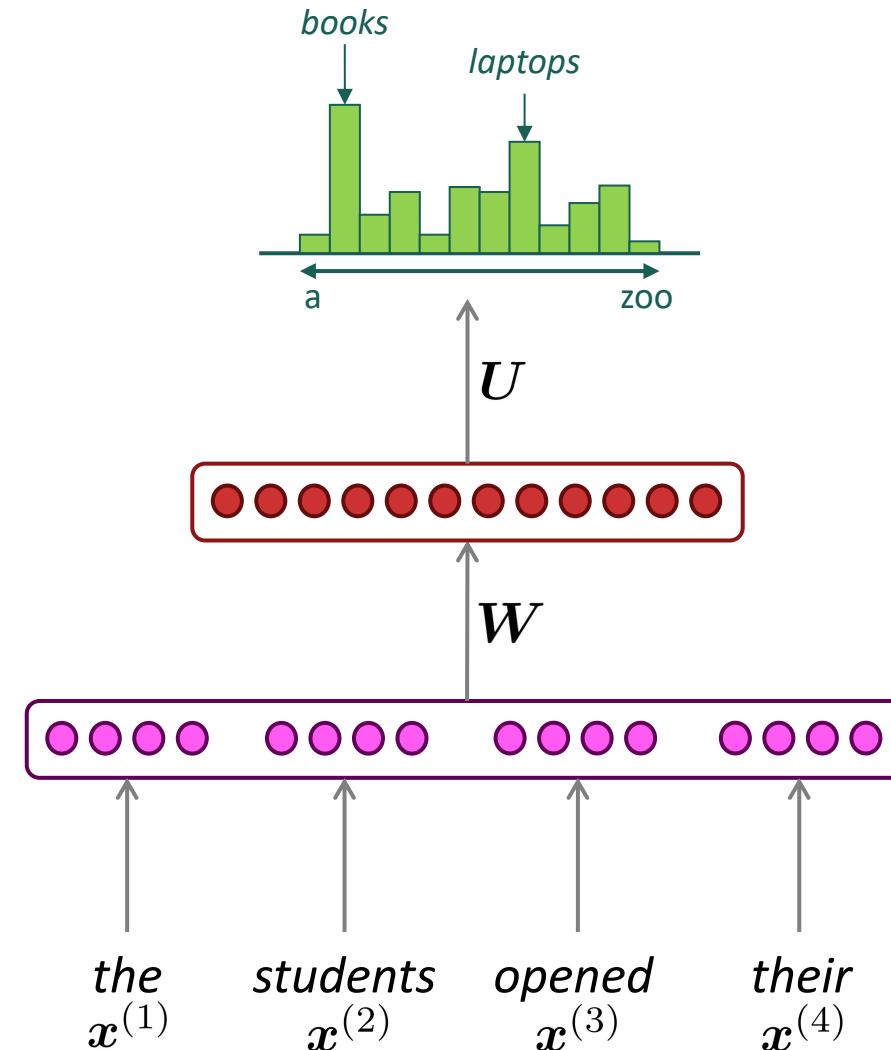
hidden layer

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

concatenated word embeddings

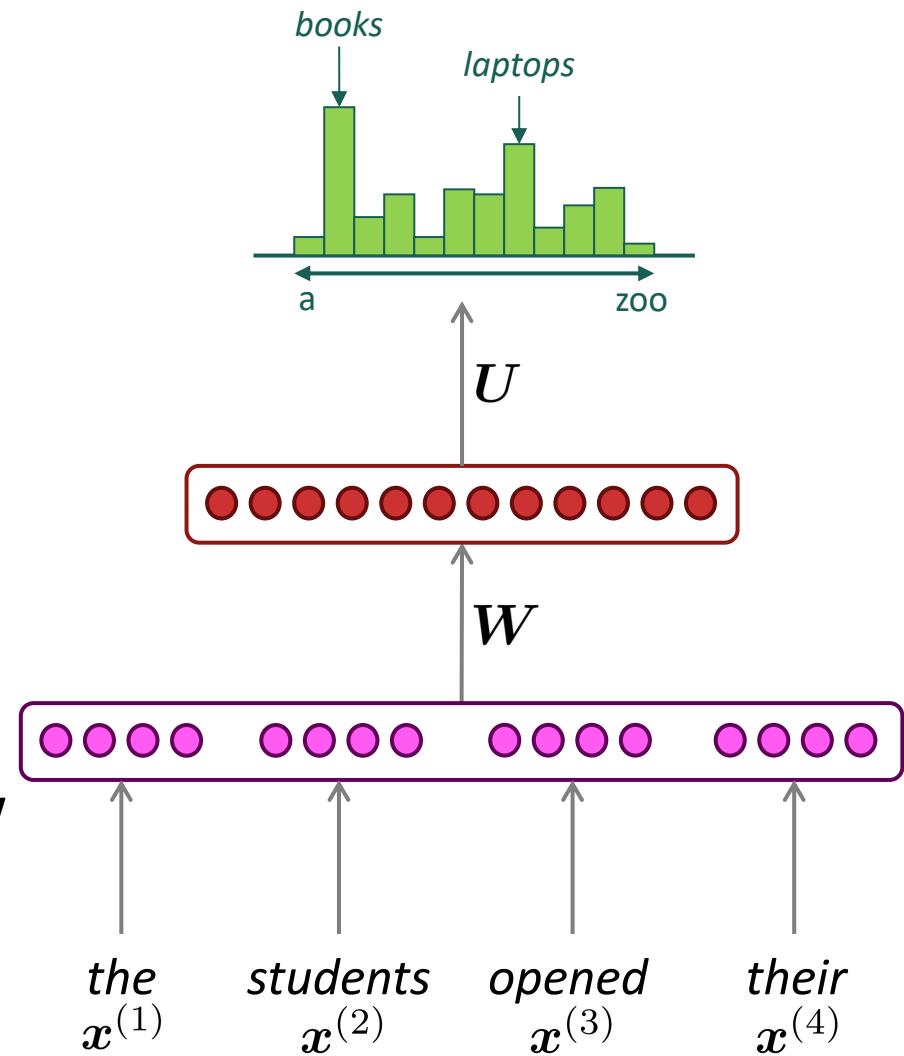
$$\mathbf{e} = [\mathbf{e}^{(1)}; \mathbf{e}^{(2)}; \mathbf{e}^{(3)}; \mathbf{e}^{(4)}]$$

words / one-hot vectors
 $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)}, \mathbf{x}^{(4)}$



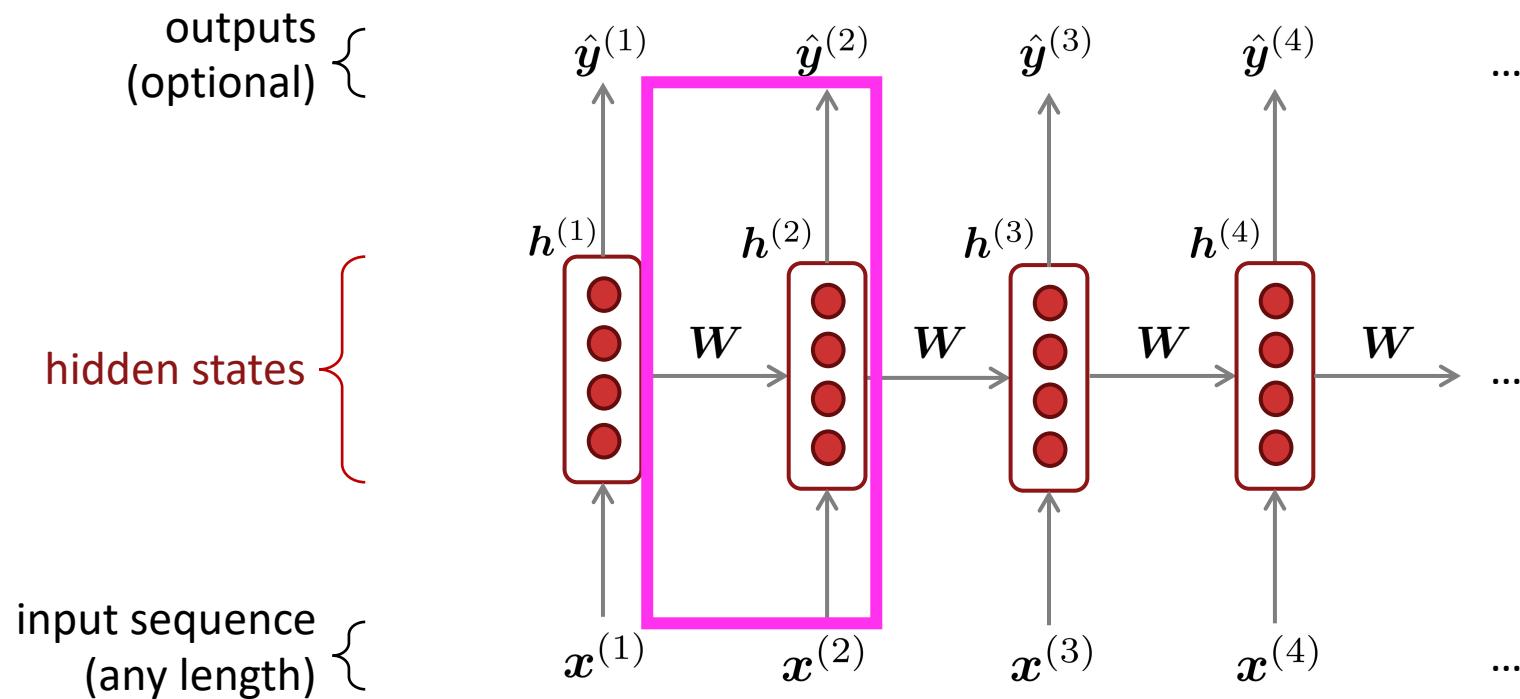
A fixed-window neural Language Model

- 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!
 - $x(0)$ and $x(1)$ 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!**



Recurrent Neural Networks (RNN)

- ❑ A family of neural architectures
- ❑ Core idea: Apply the same weights W repeatedly



A Simple RNN Language Model

output distribution

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

hidden states

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

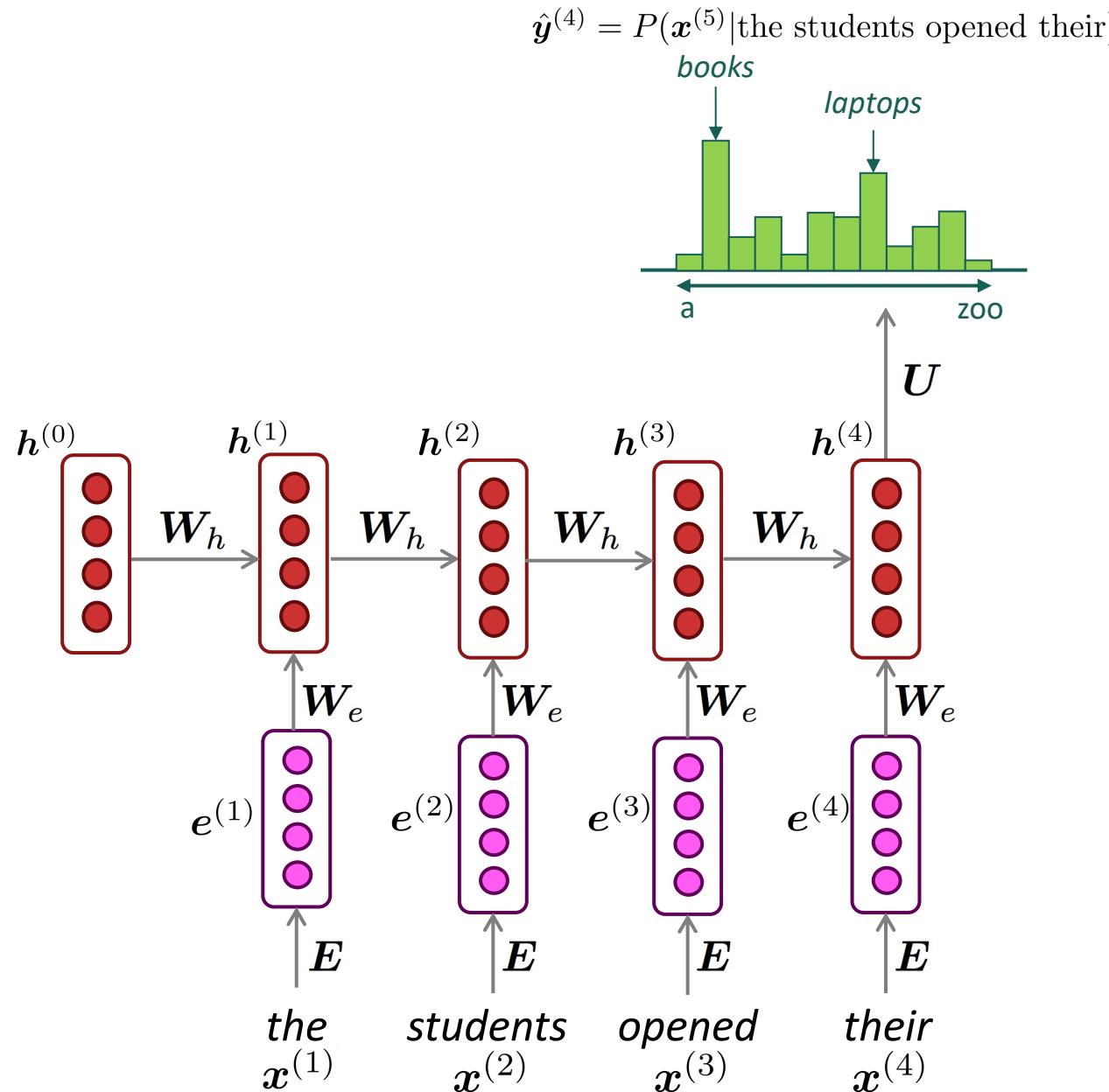
$\mathbf{h}^{(0)}$ is the initial hidden state

word embeddings

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

words / one-hot vectors

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



RNN Language Models

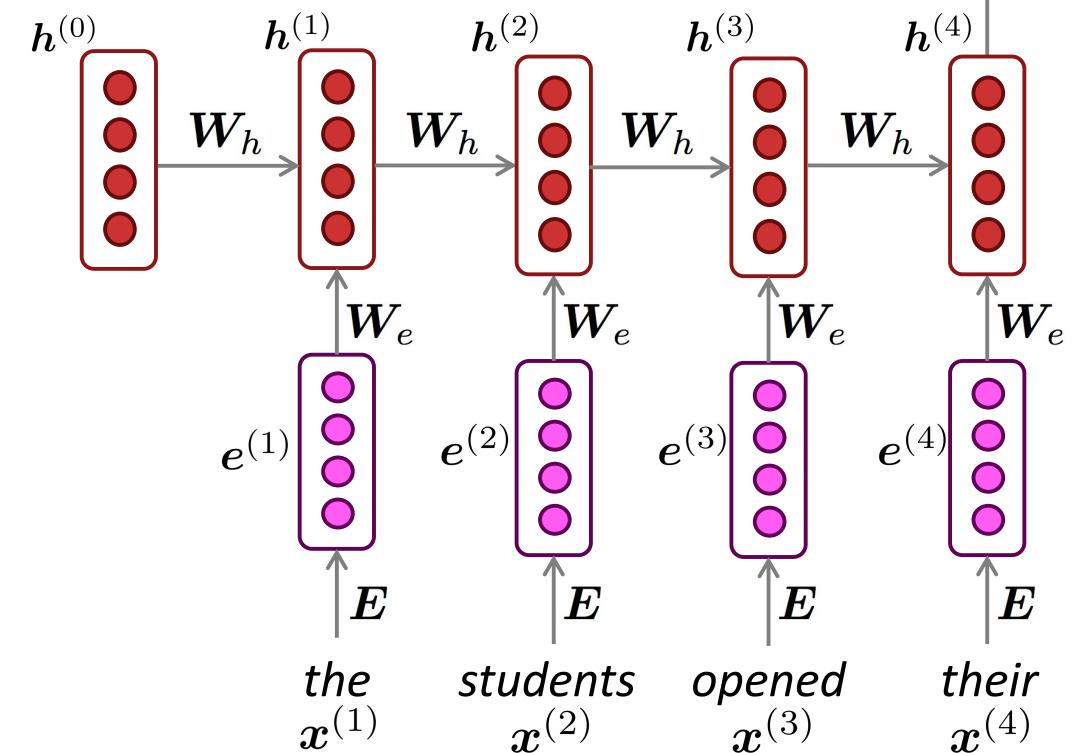
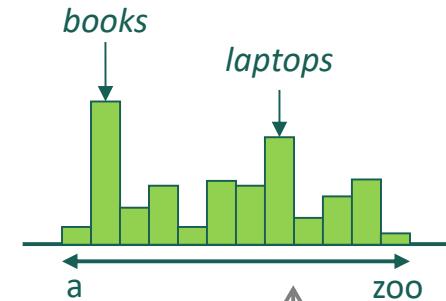
□ RNN Advantages:

- 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

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



Training an RNN Language Model

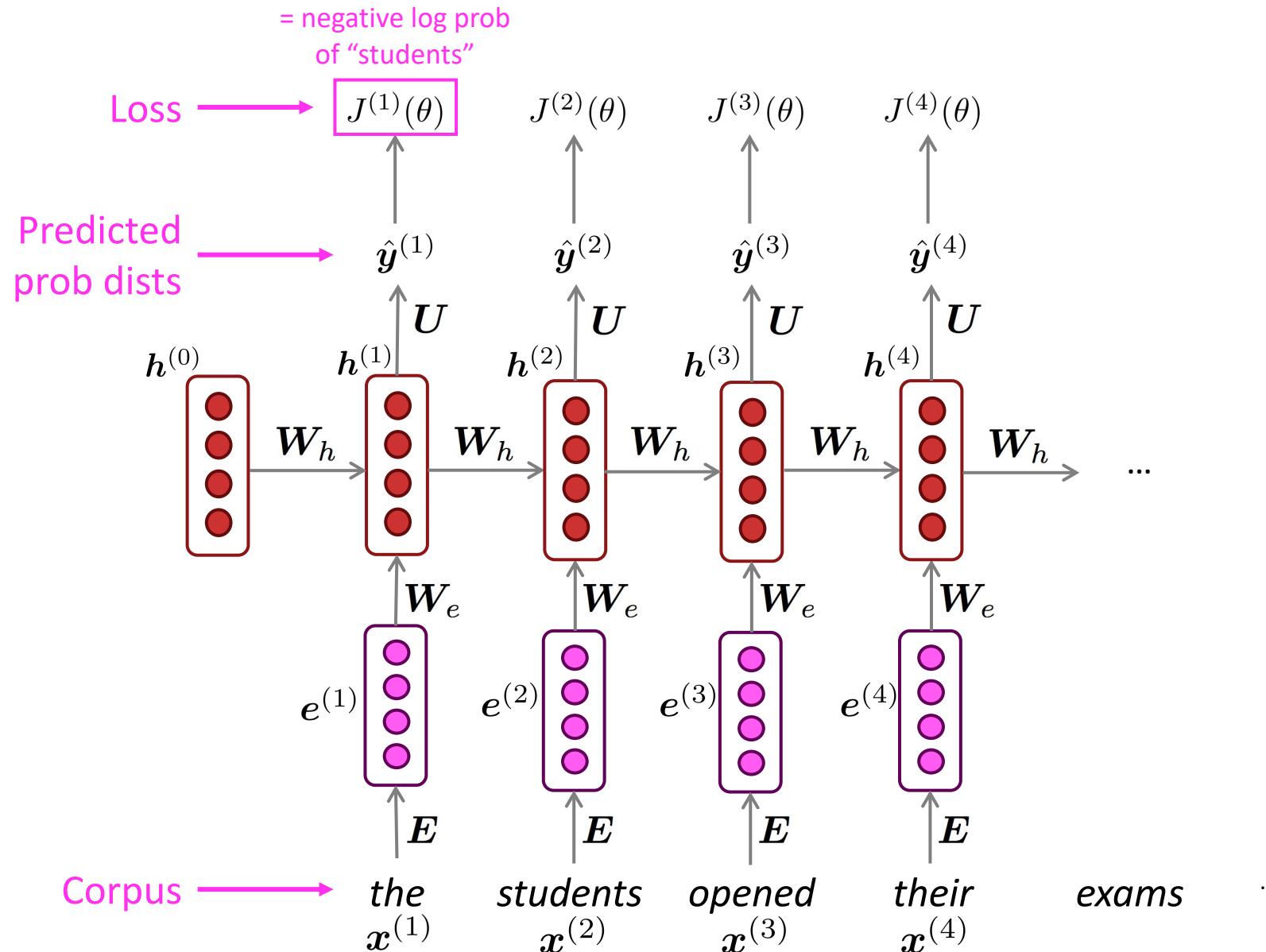
- Get a **big corpus of text** which is a sequence of words $x^{(1)}, \dots, x^{(T)}$
- Feed into RNN-LM; compute output distribution $\hat{y}^{(t)}$ for every step t
 - i.e. predict probability distribution of *every word*, given words so far
- **Loss function** on step t is **cross-entropy** between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

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

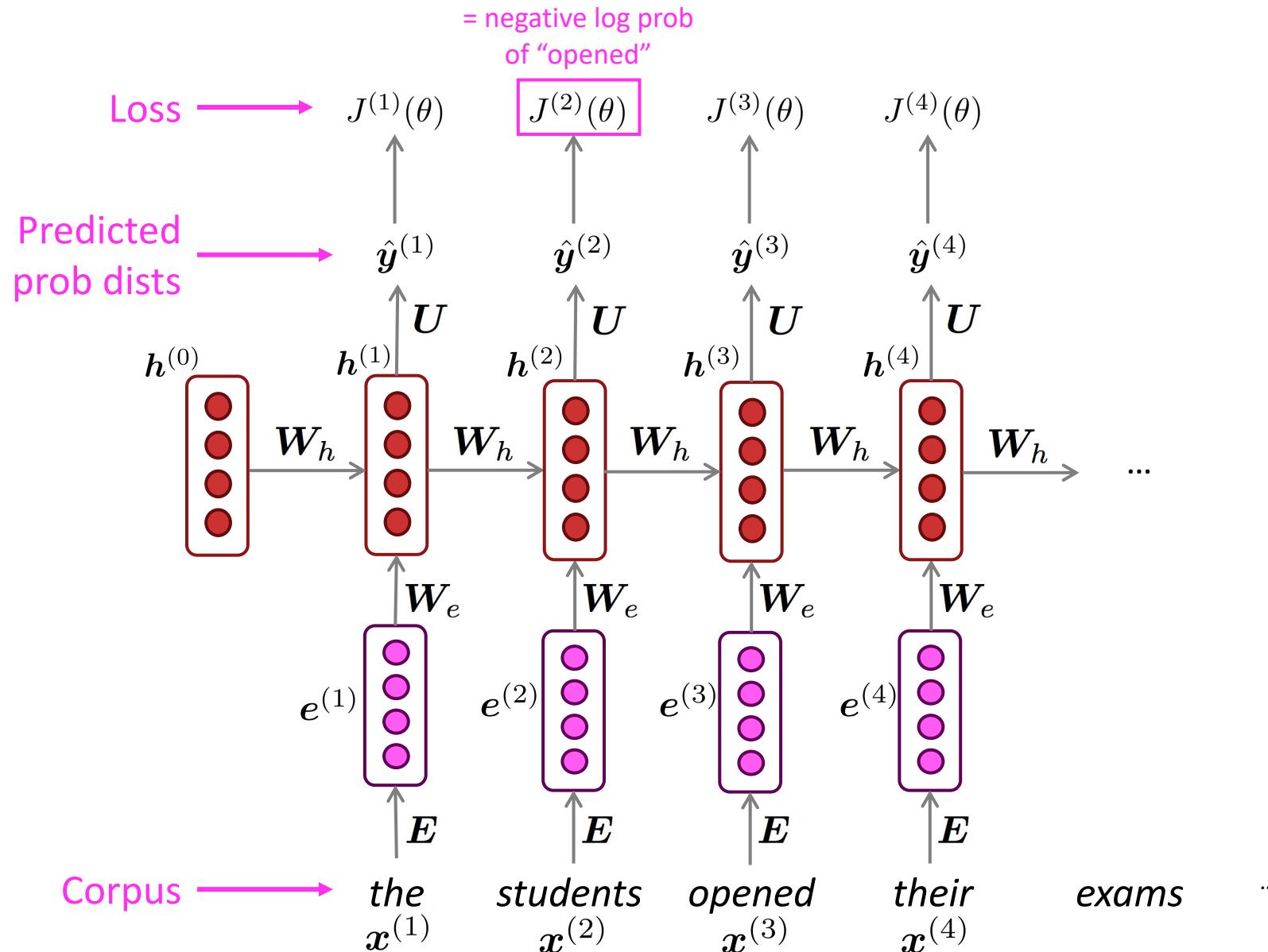
- Average this to get **overall loss** for entire training set:

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

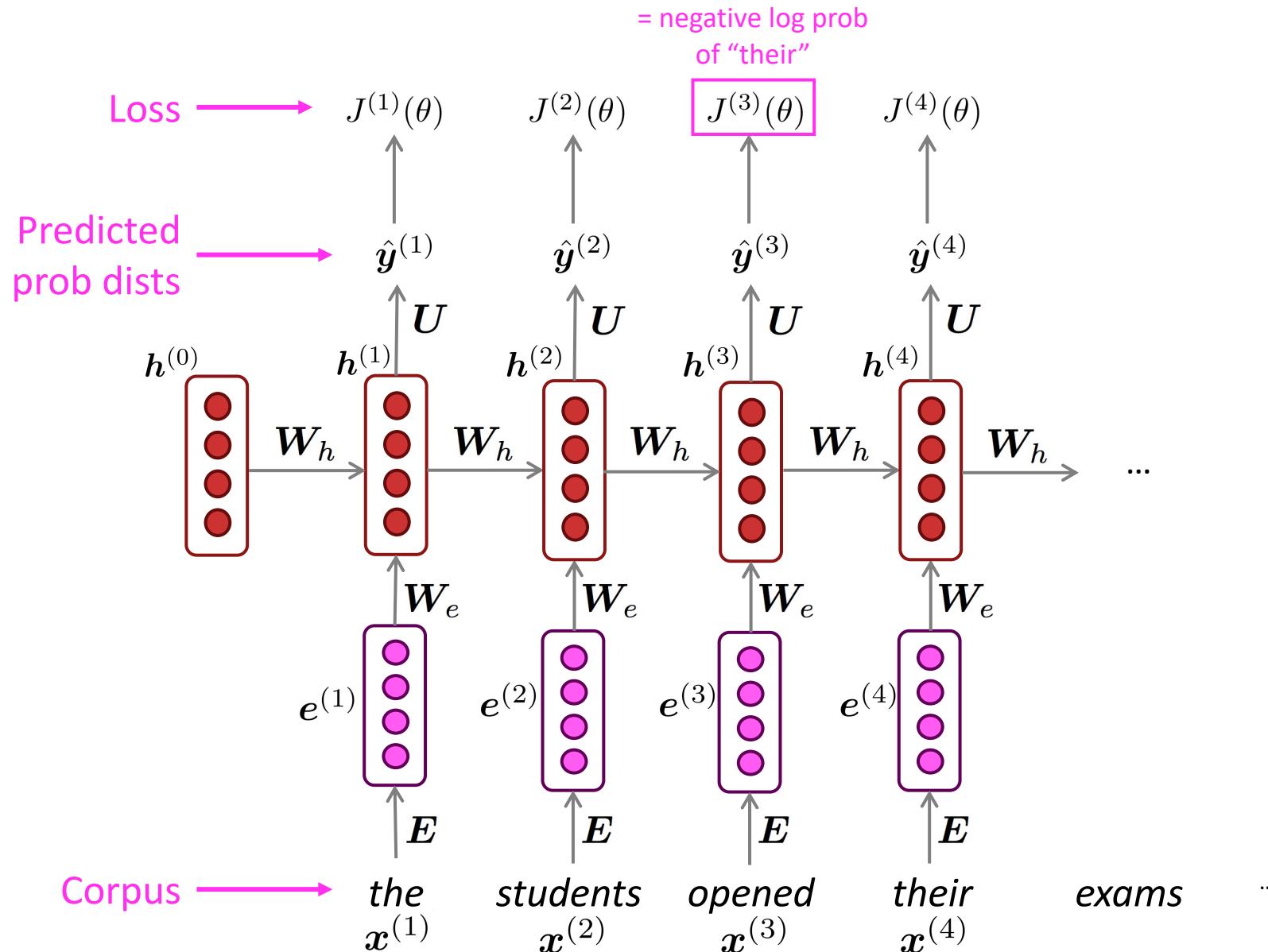
Training an RNN Language Model



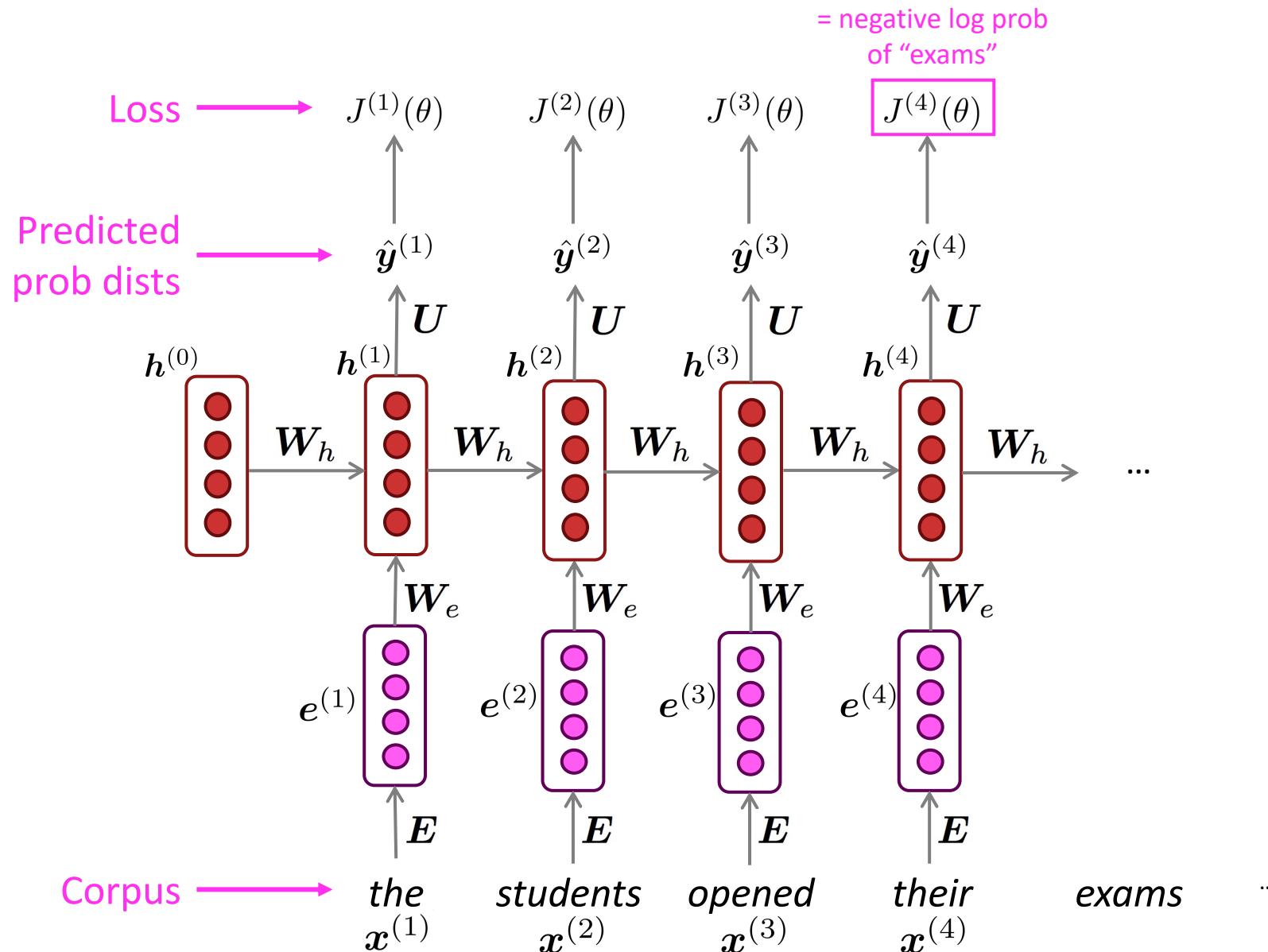
Training an RNN Language Model



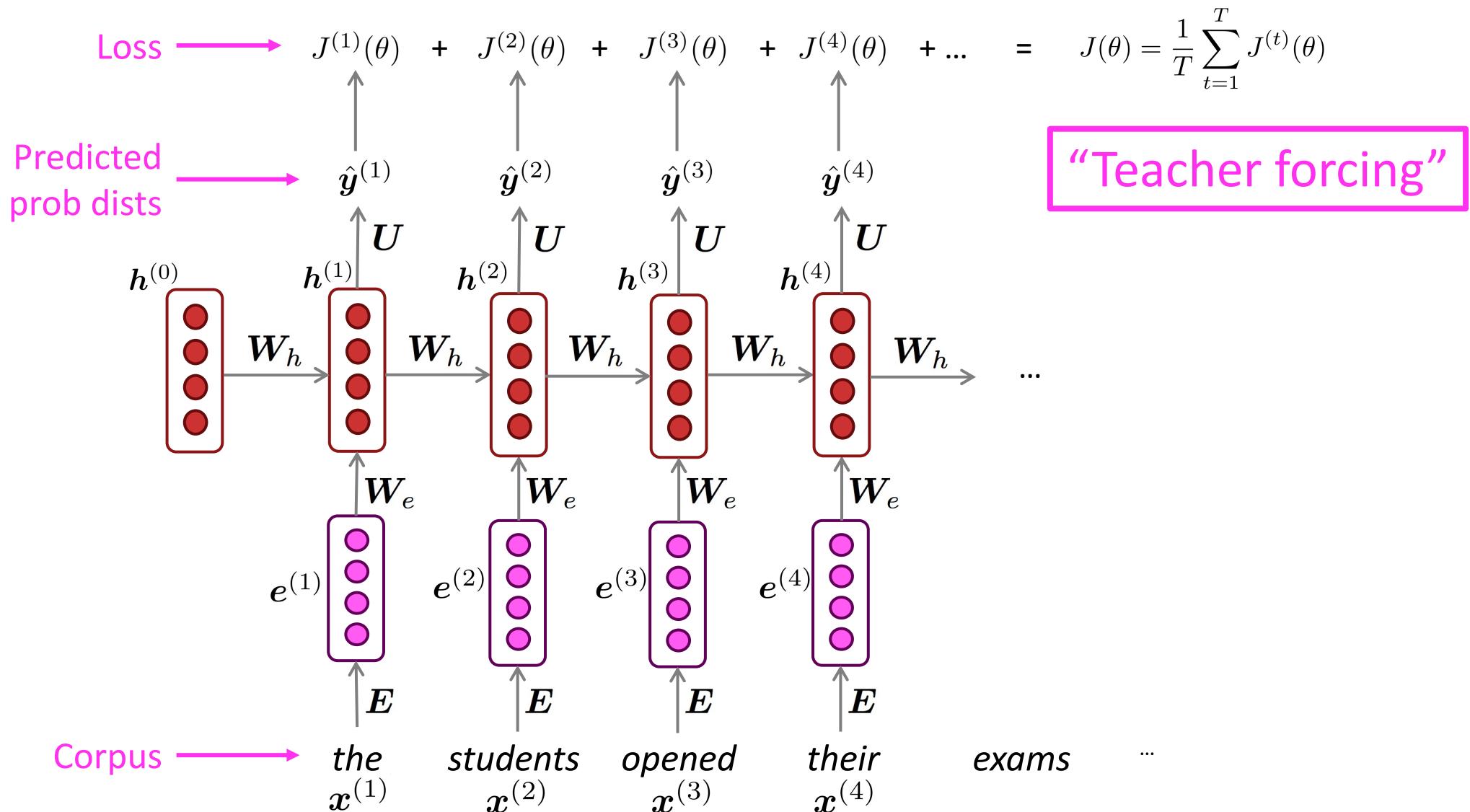
Training an RNN Language Model



Training an RNN Language Model



Training an RNN Language Model



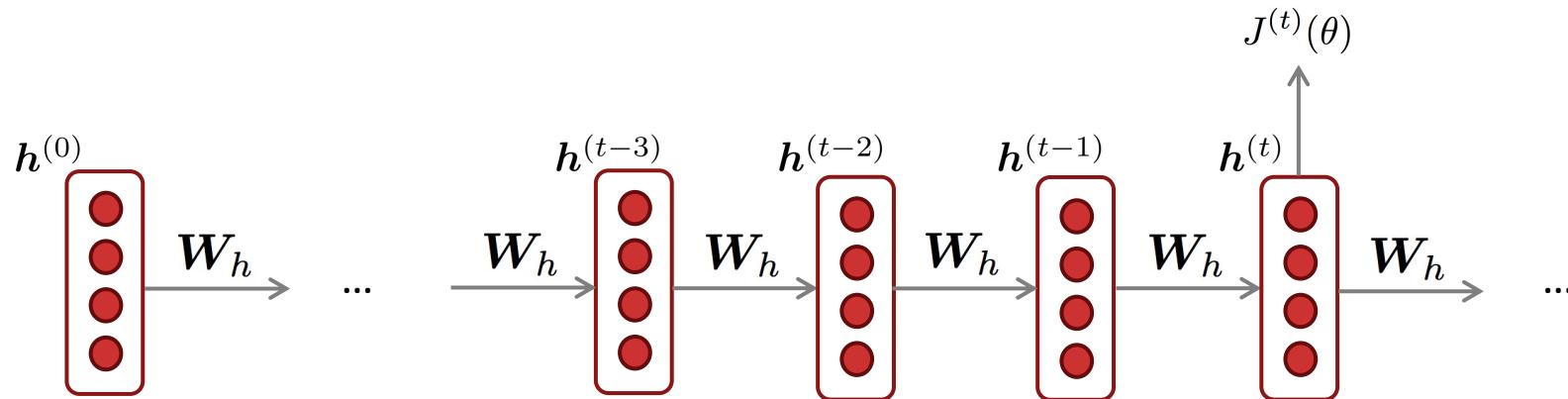
Training an RNN Language Model

- ❑ 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
- ❑ Compute loss $J(\theta)$ for a sentence (actually, a batch of sentences), compute gradients and update weights. Repeat.

Backpropagation for RNNs



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

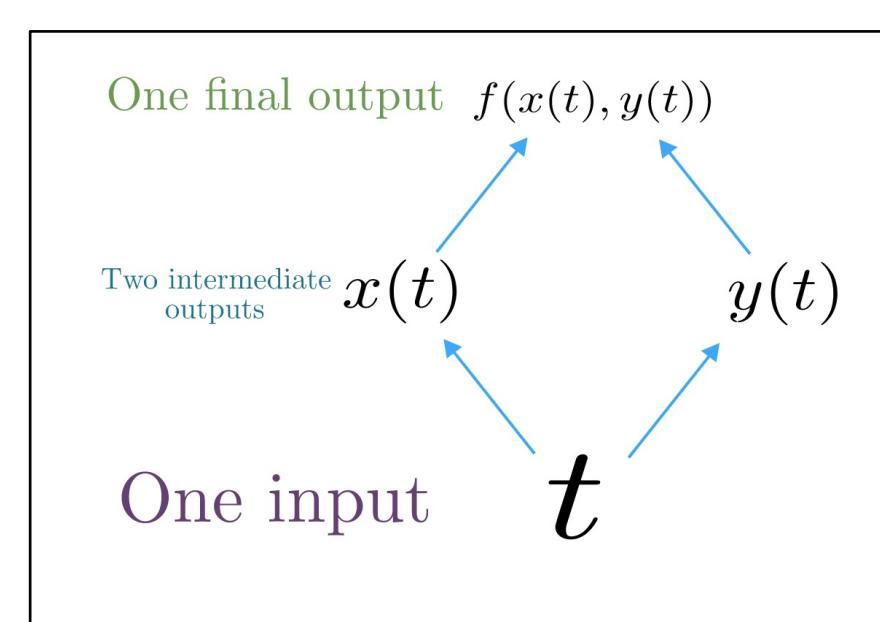
❑ **Answer:** $\frac{\partial J^{(t)}}{\partial W_h} = \sum_{i=1}^t \frac{\partial J^{(t)}}{\partial W_h} \Big|_{(i)}$

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

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:

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

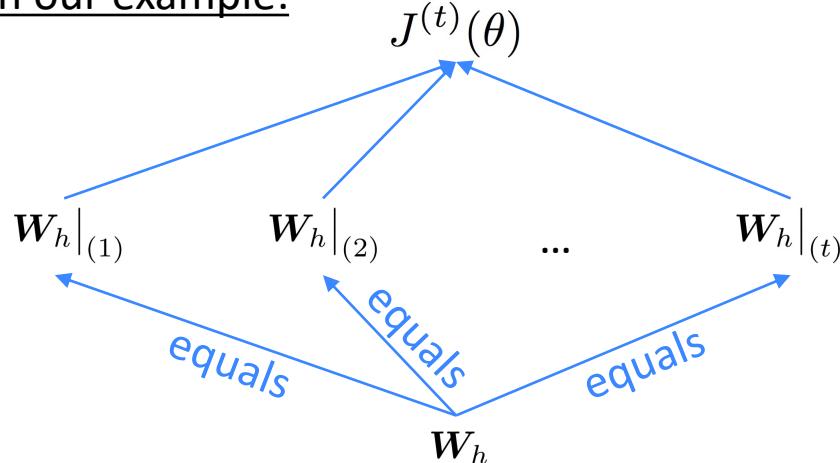


Backpropagation for RNNs: Proof sketch

- 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:

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

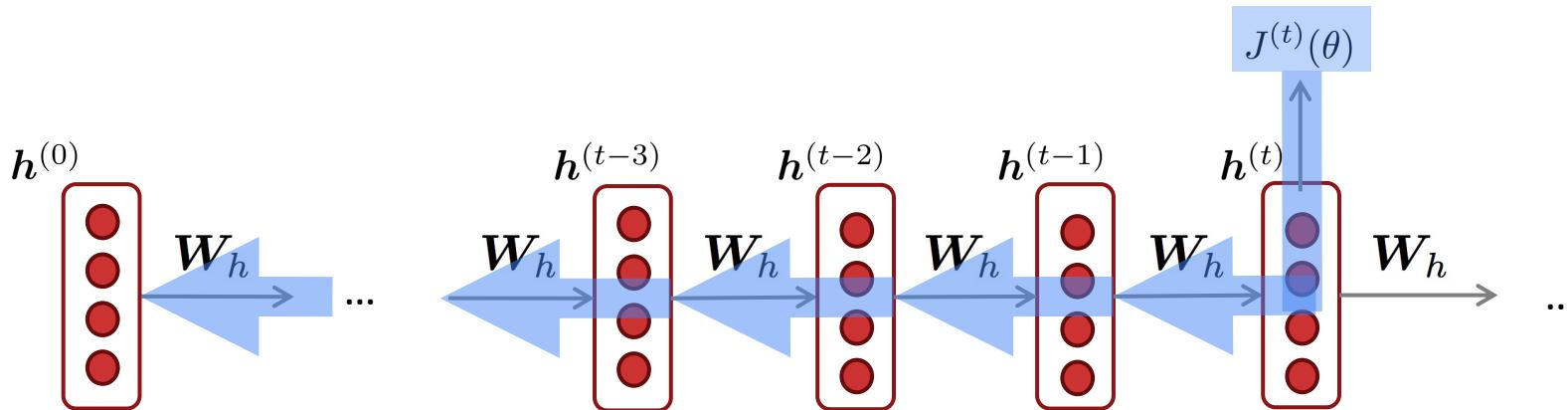
In our example:



Apply the multivariable chain rule:

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

Backpropagation for RNNs



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

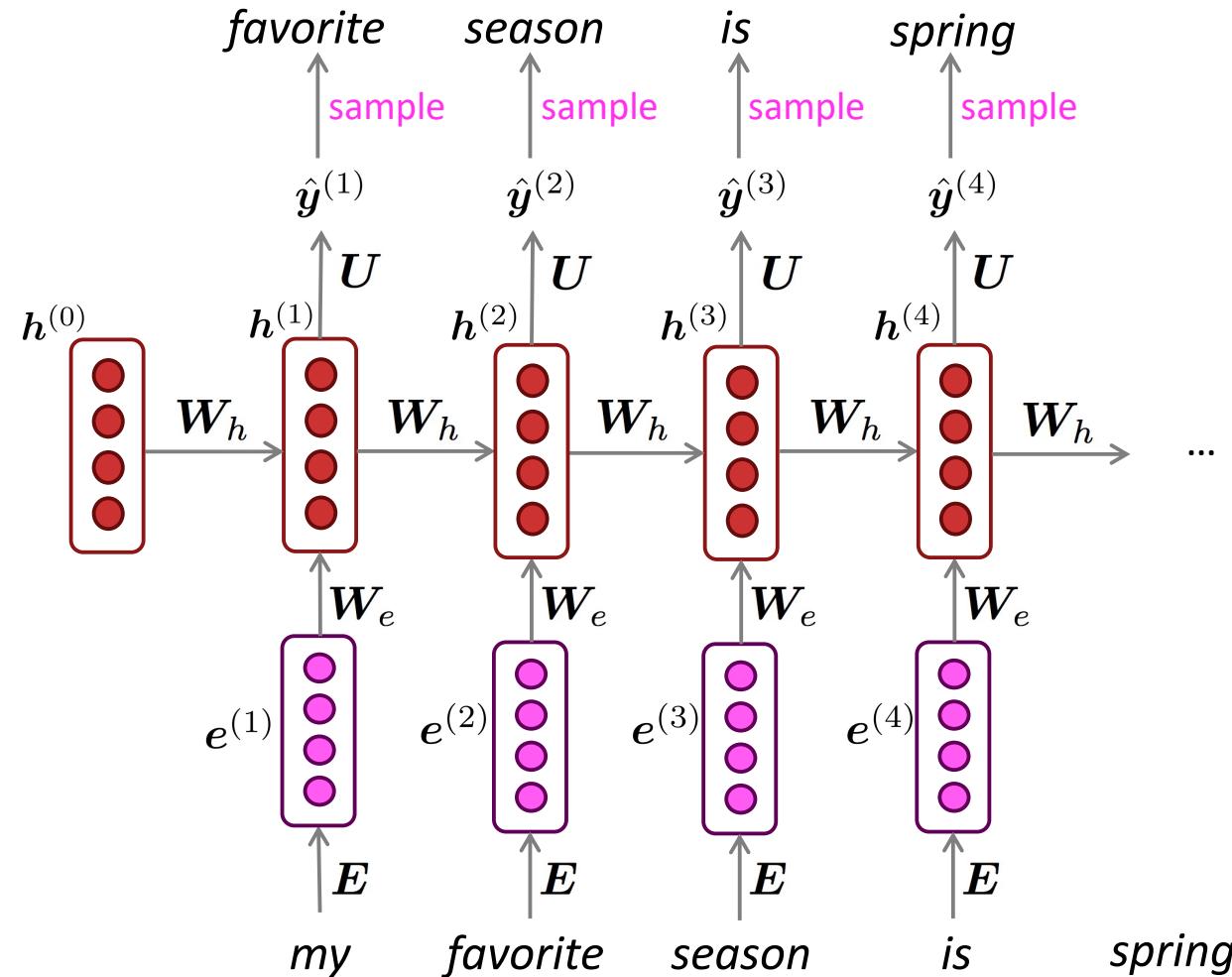
Question: How do we calculate this?

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

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 becomes next step's input.



Generating text with a RNN Language Model

□ 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.

Generating text with a RNN Language Model

❑ 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 *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.

Generating text with a RNN Language Model

❑ 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 *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>

Evaluating Language Models

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

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

Inverse probability of corpus, according to Language Model

Normalized by
number of words

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

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

Lower perplexity is better!

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?

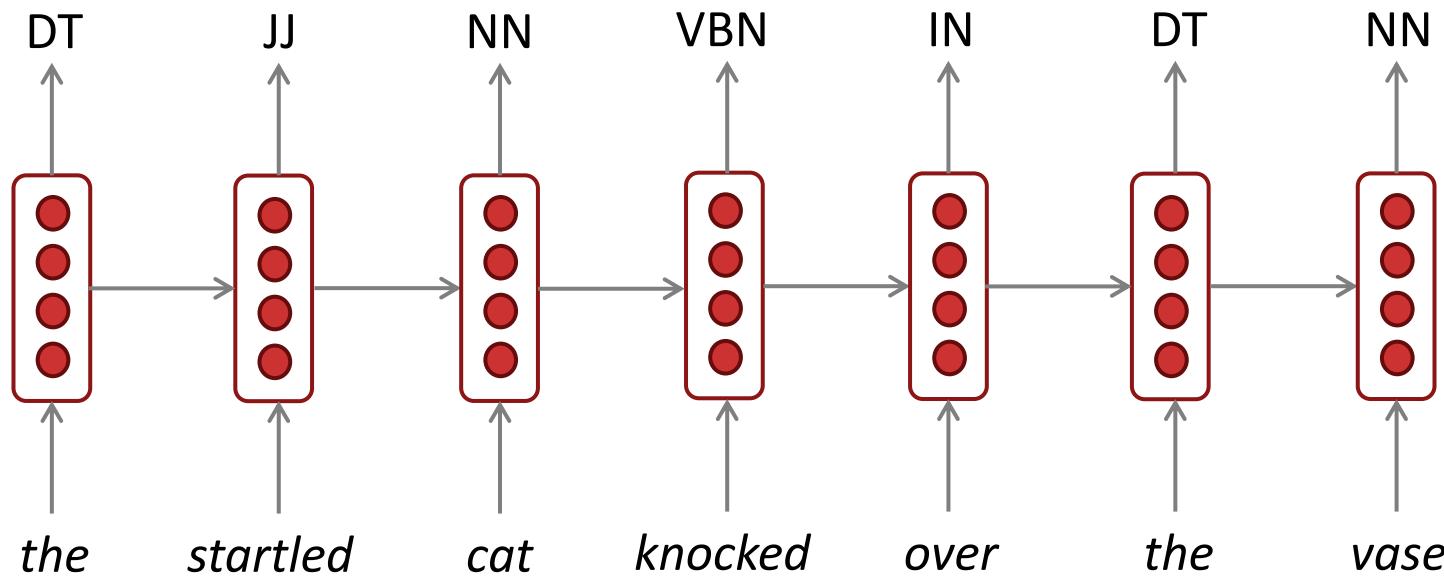
- ❑ Language Modeling is a **benchmark task** that helps us **measure our progress** on understanding language
- ❑ Language Modeling is a **subcomponent** of many NLP tasks, especially those involving **generating text** or **estimating the probability of text**:
 - Predictive typing
 - Speech recognition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

Summary

- ❑ 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!

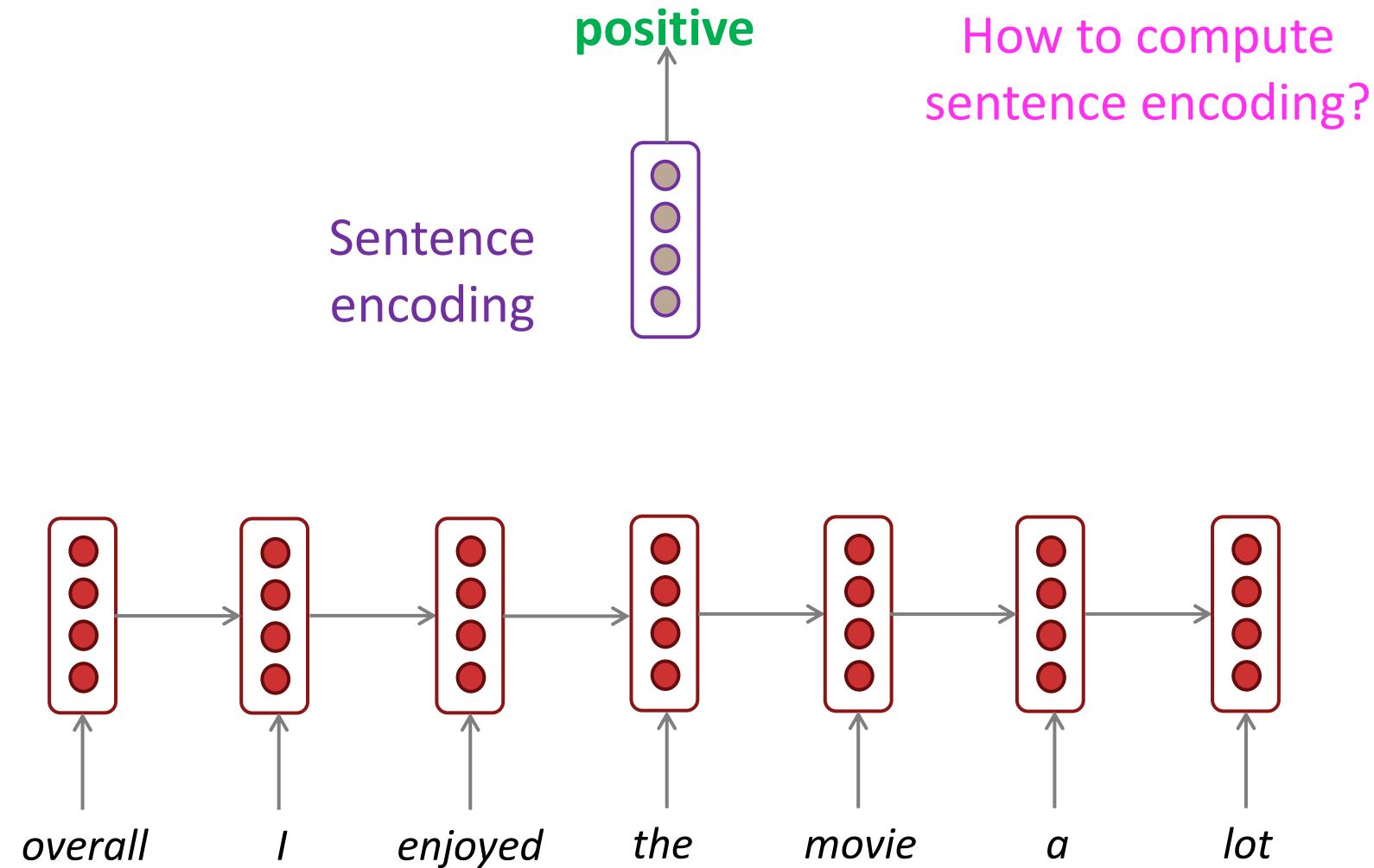
RNNs can be used for tagging

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



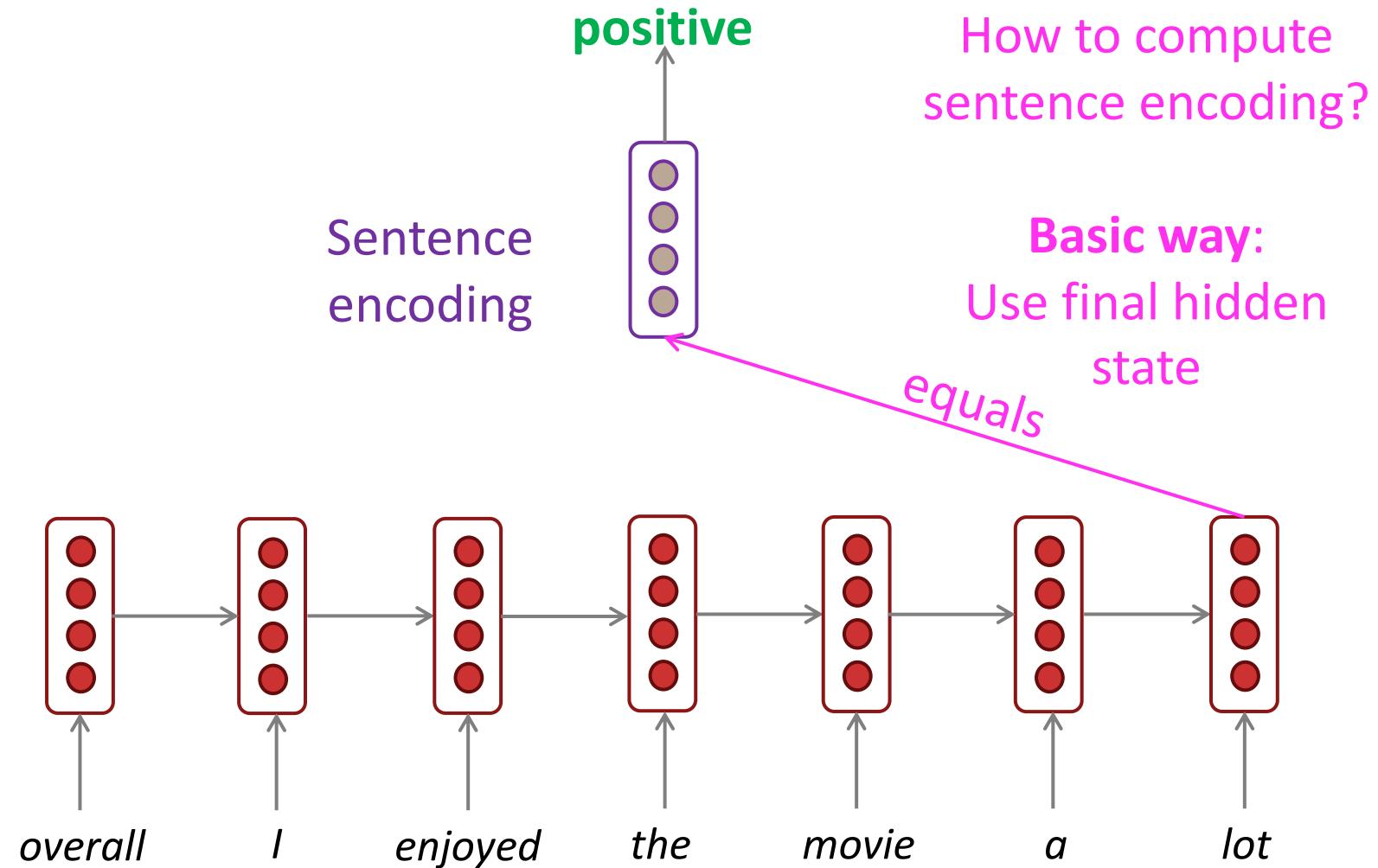
RNNs can be used for sentence classification

- e.g., sentiment classification



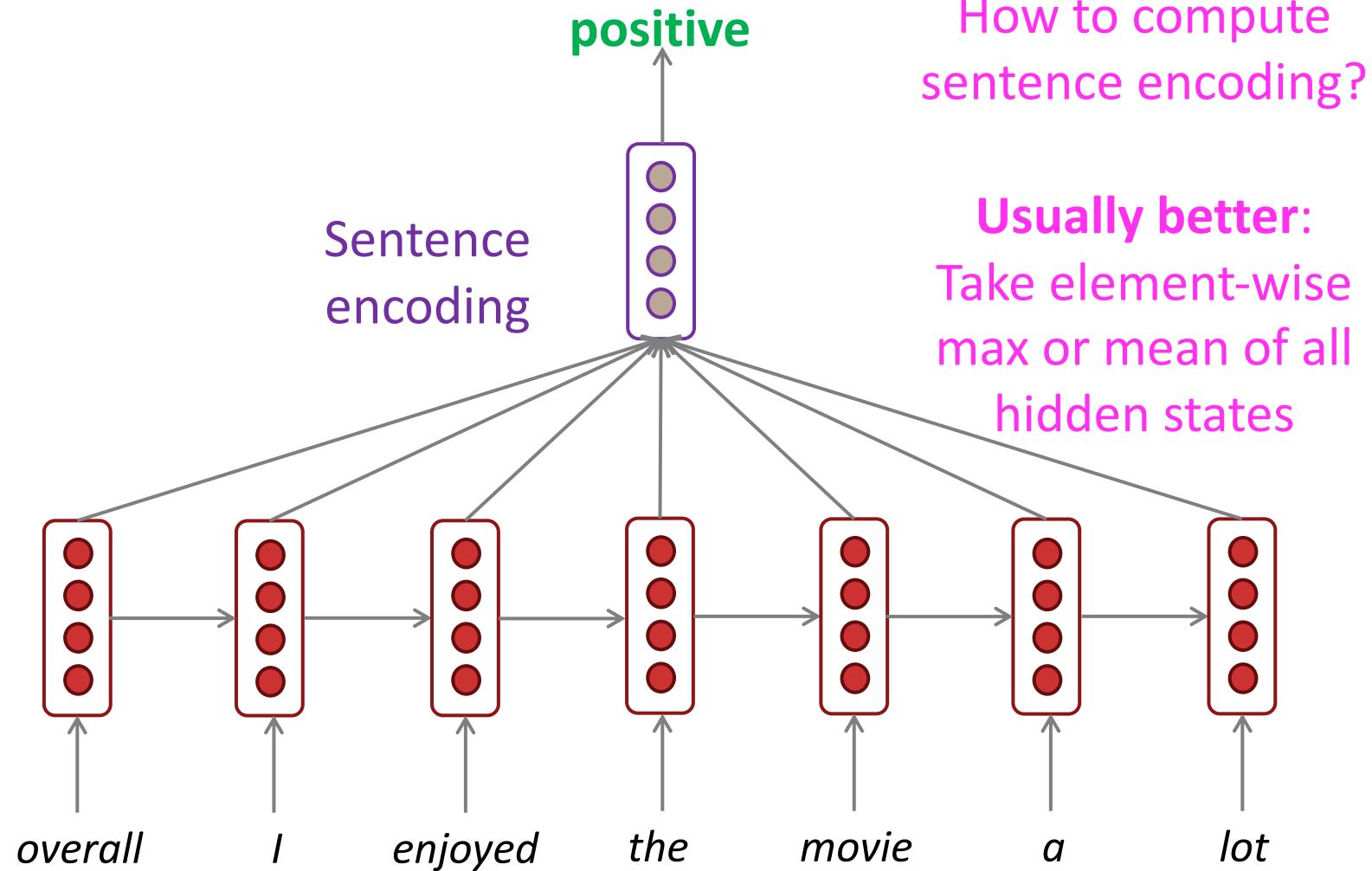
RNNs can be used for sentence classification

- e.g., sentiment classification



RNNs can be used for sentence classification

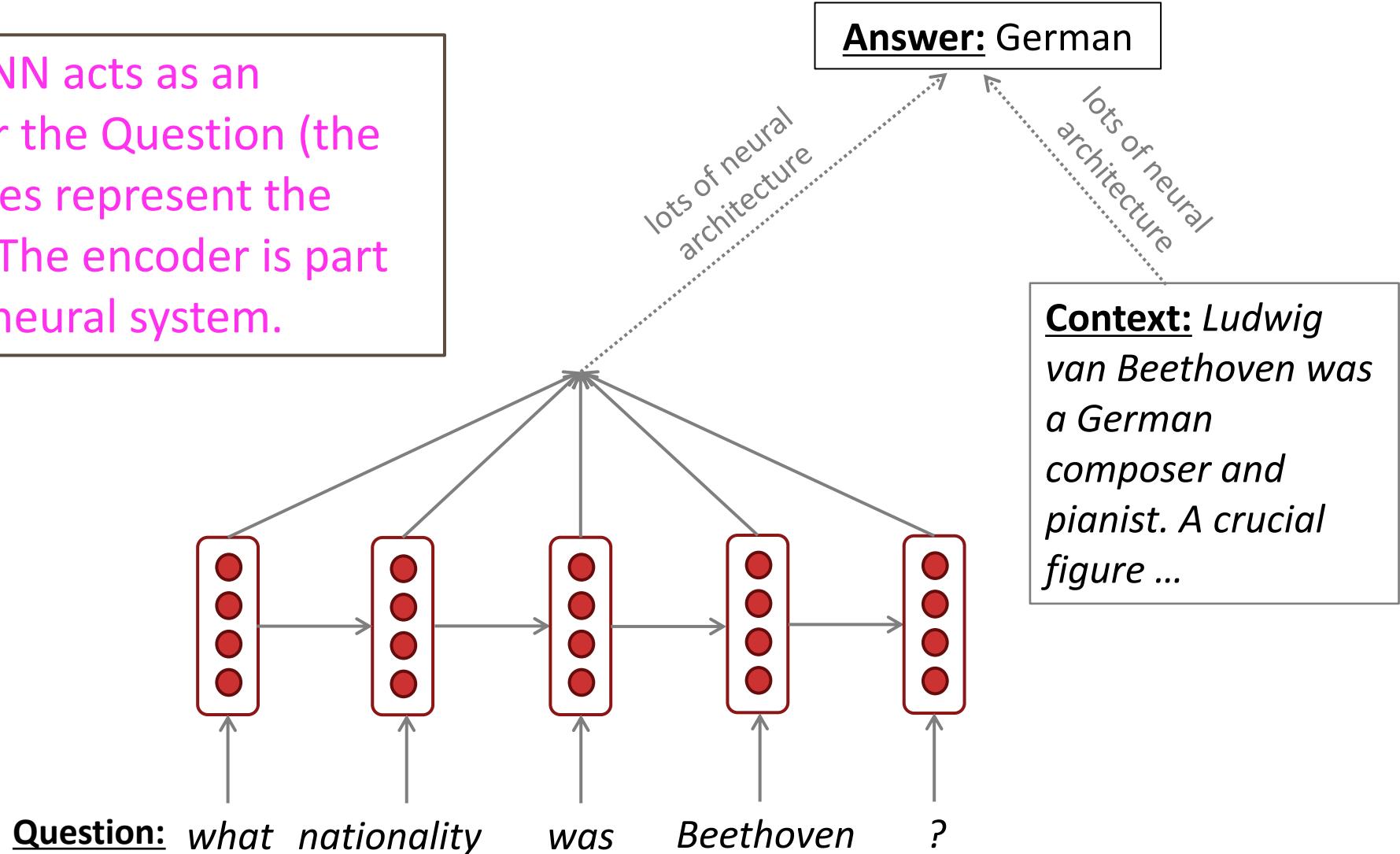
- e.g., sentiment classification



RNNs can be used as an encoder module

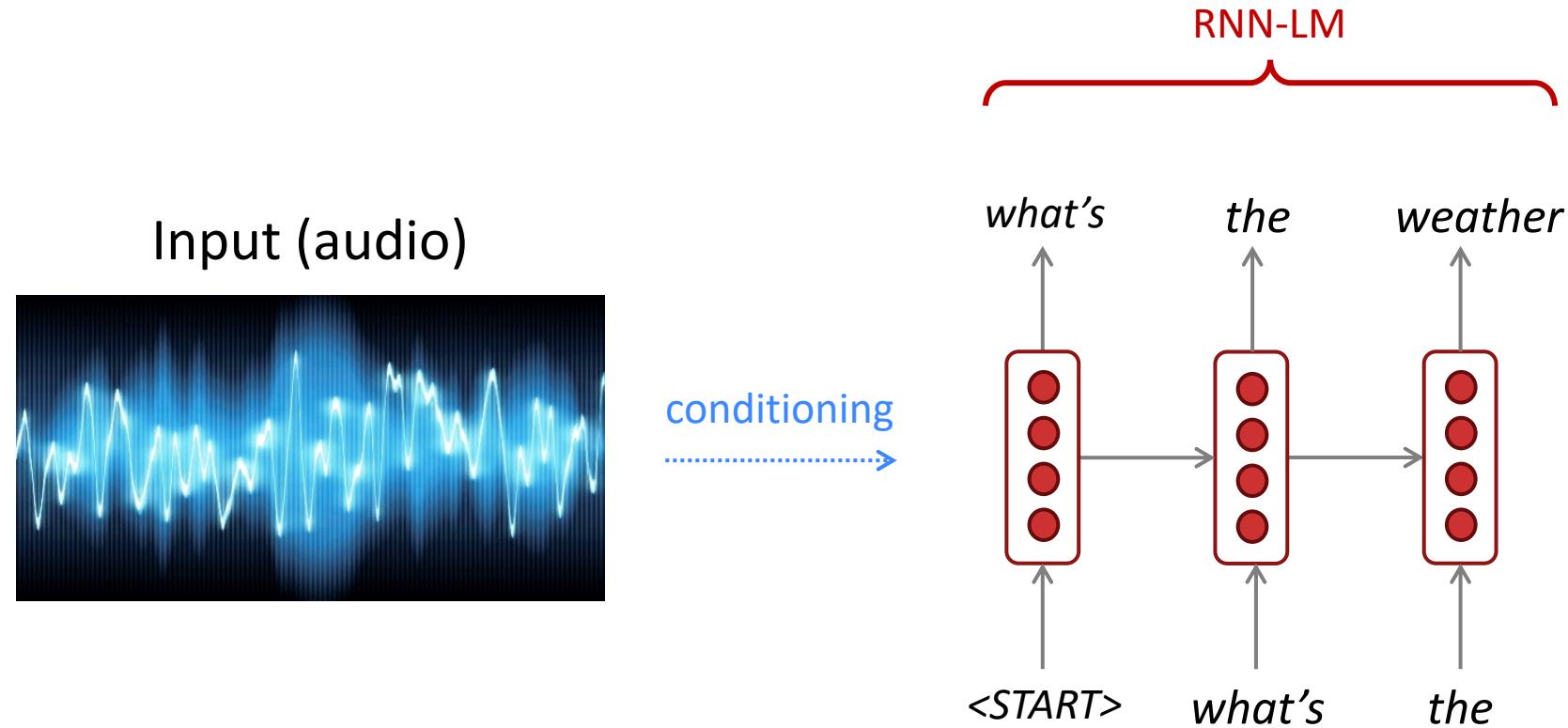
- e.g., **question answering**, machine translation, many other tasks!

Here the RNN acts as an **encoder** for the Question (the hidden states represent the Question). The encoder is part of a larger neural system.



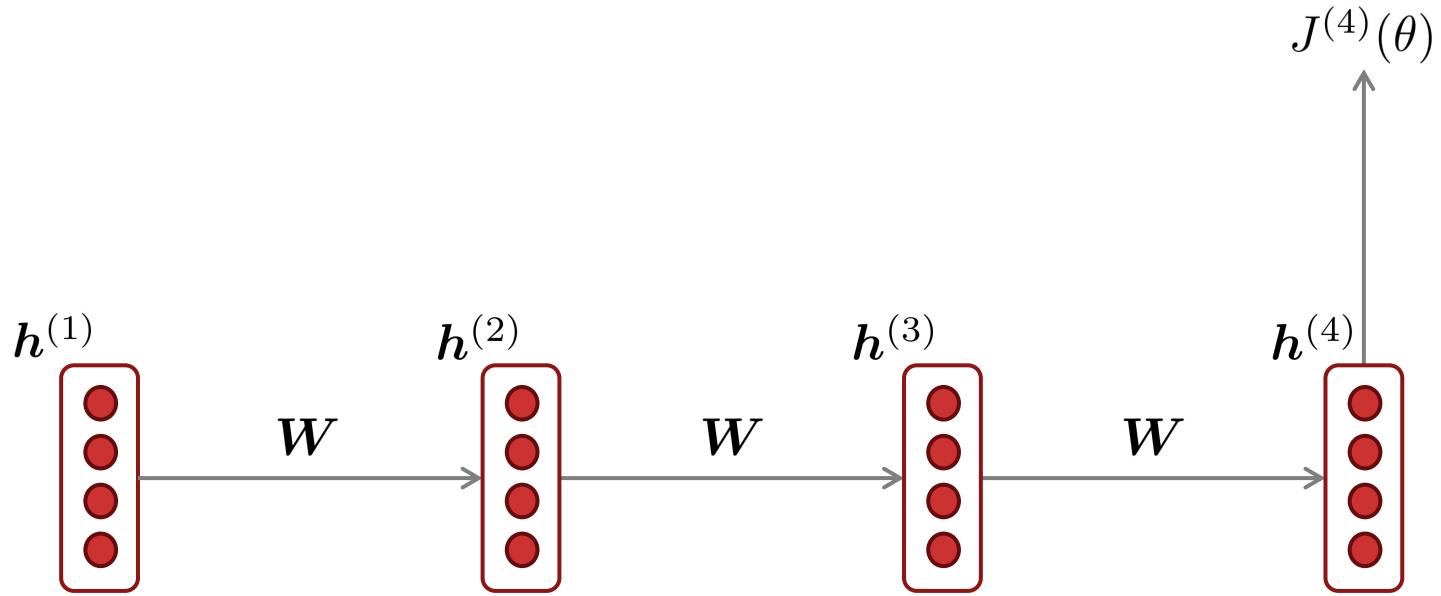
RNN-LMs can be used to generate text

- e.g., speech recognition, machine translation, summarization

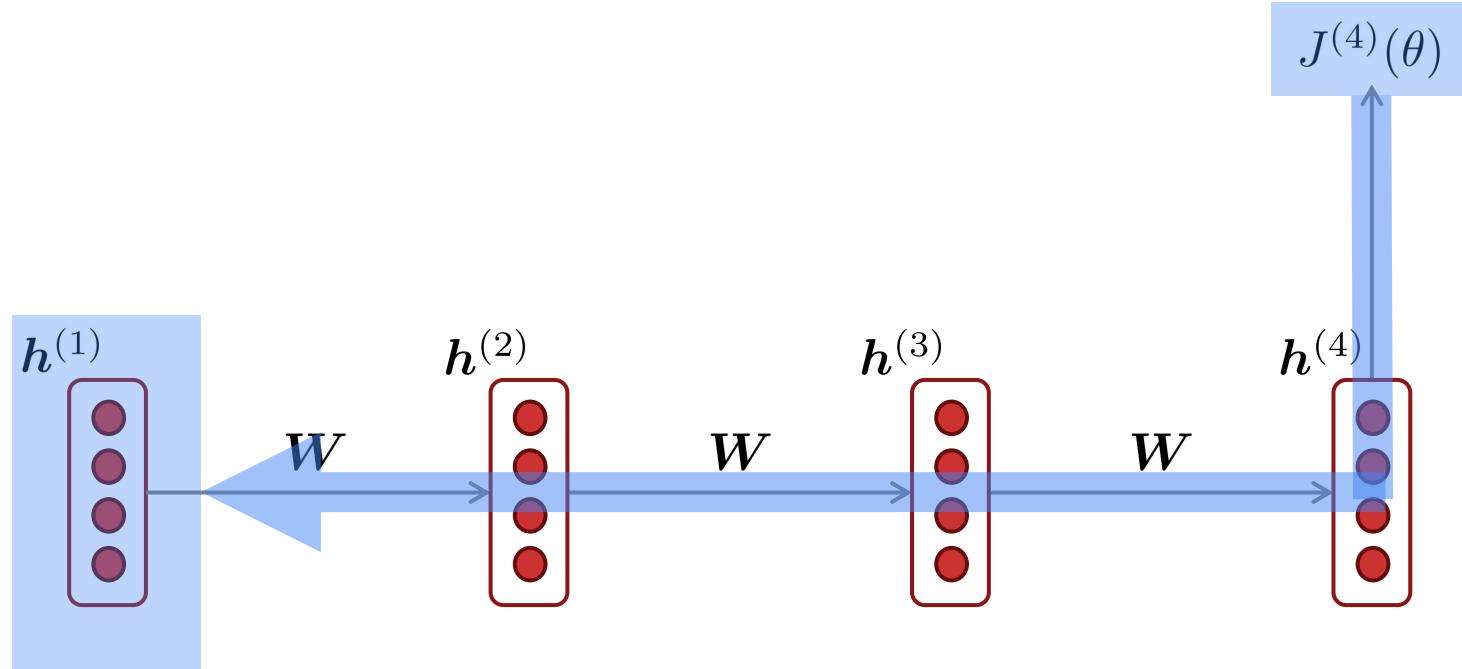


This is an example of a *conditional language model*.
We'll see Machine Translation in much more detail later.

Problems with RNNs: Vanishing and Exploding Gradients

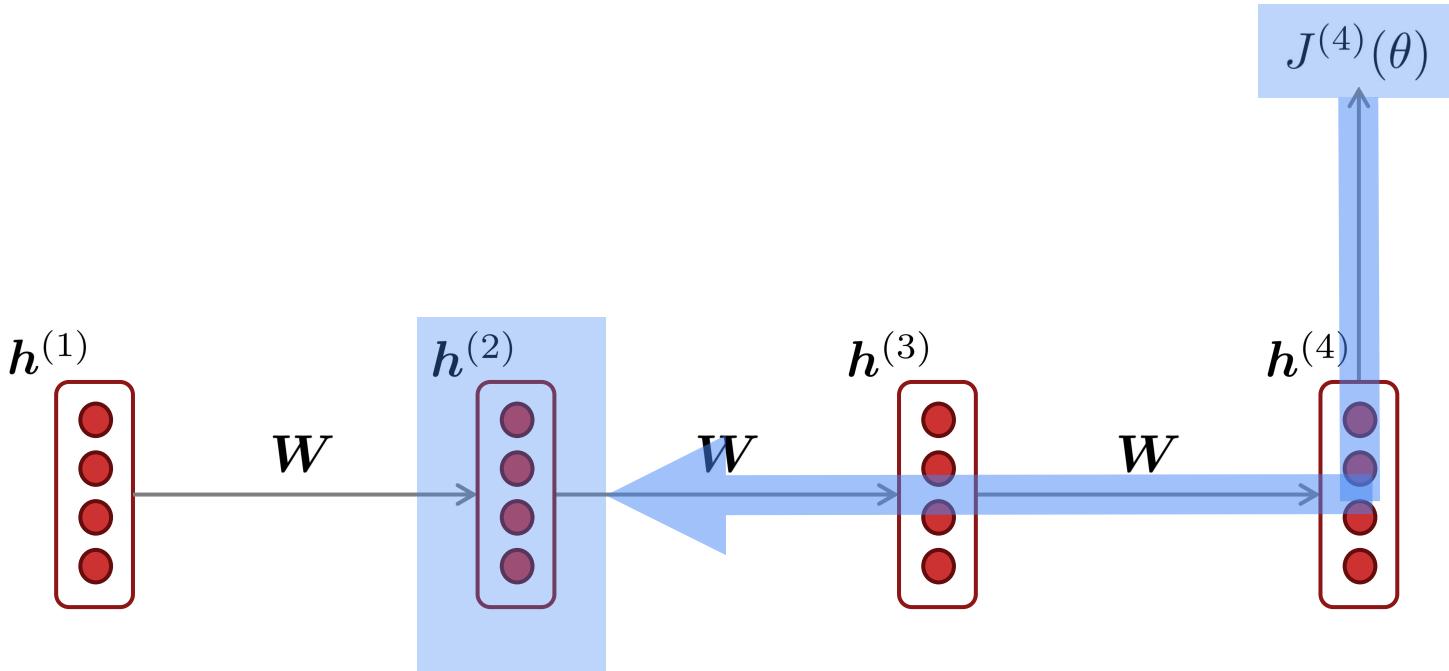


Vanishing gradient intuition



$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = ?$$

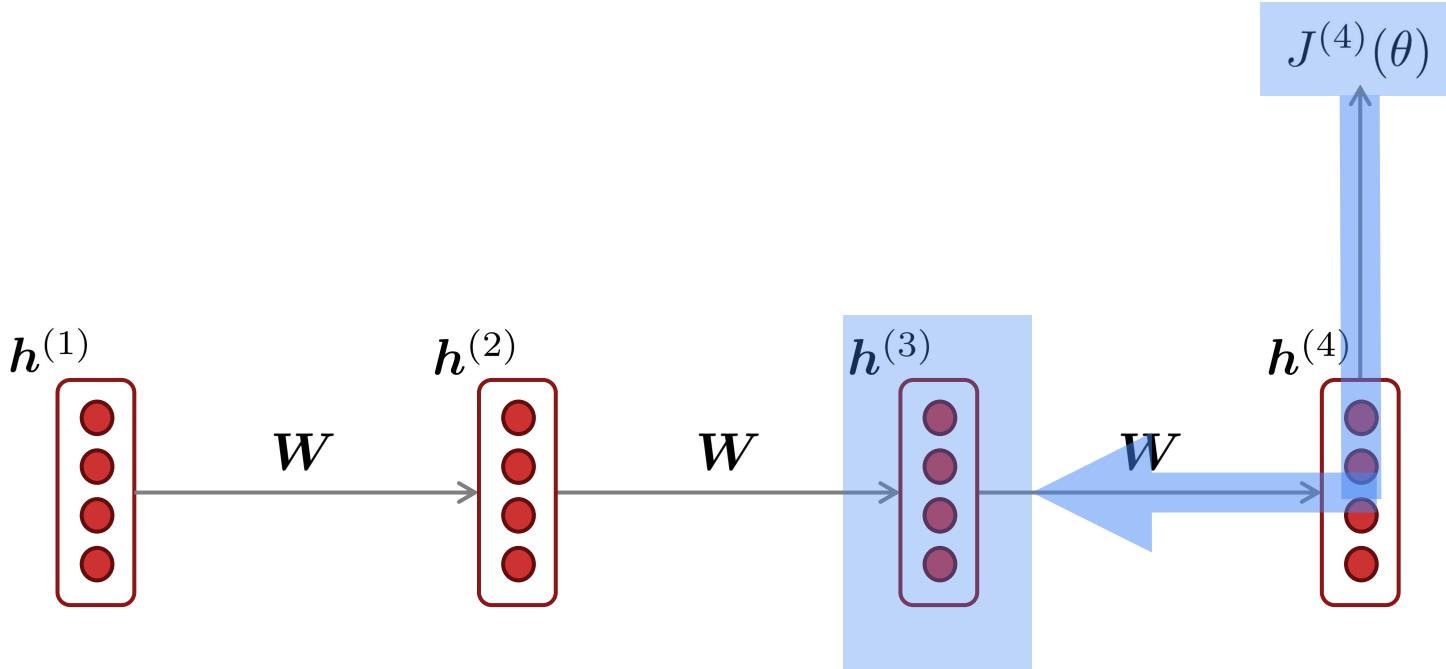
Vanishing gradient intuition



$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = \frac{\partial h^{(2)}}{\partial h^{(1)}} \times \frac{\partial J^{(4)}}{\partial h^{(2)}}$$

chain rule!

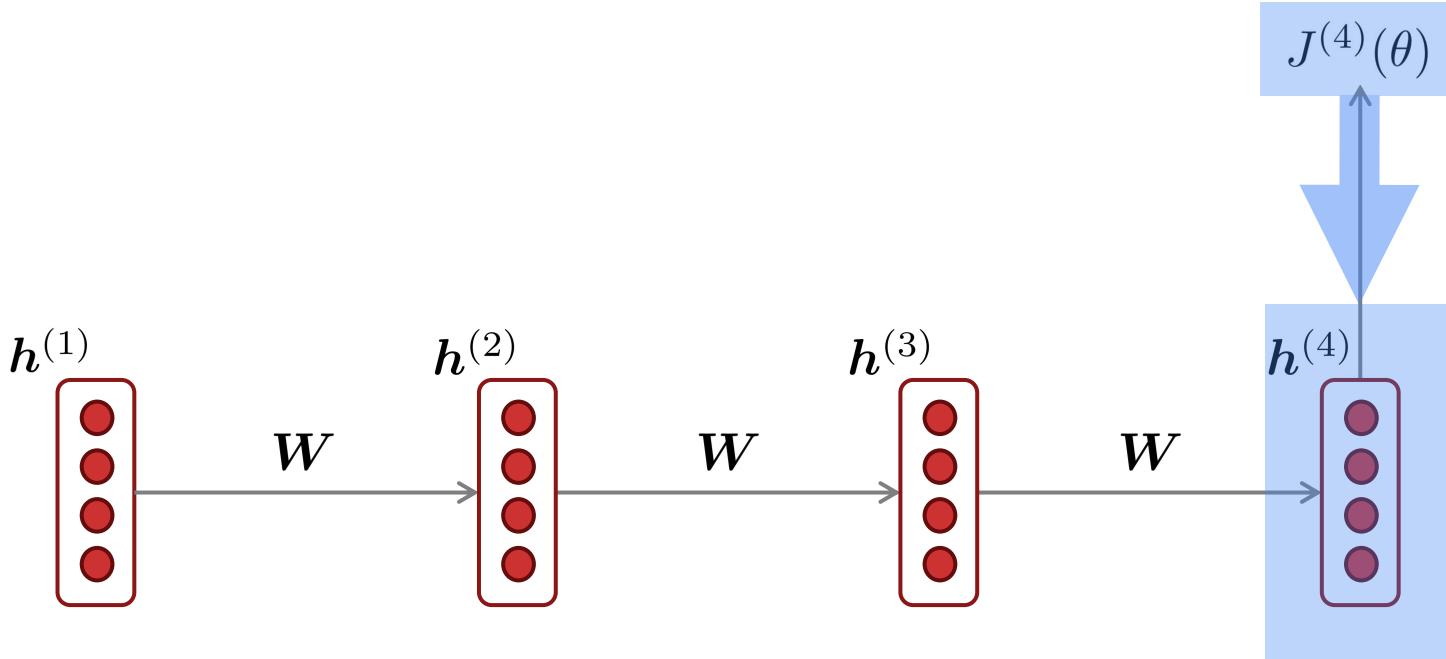
Vanishing gradient intuition



$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = \frac{\partial h^{(2)}}{\partial h^{(1)}} \times \frac{\partial h^{(3)}}{\partial h^{(2)}} \times \frac{\partial J^{(4)}}{\partial h^{(3)}}$$

chain rule!

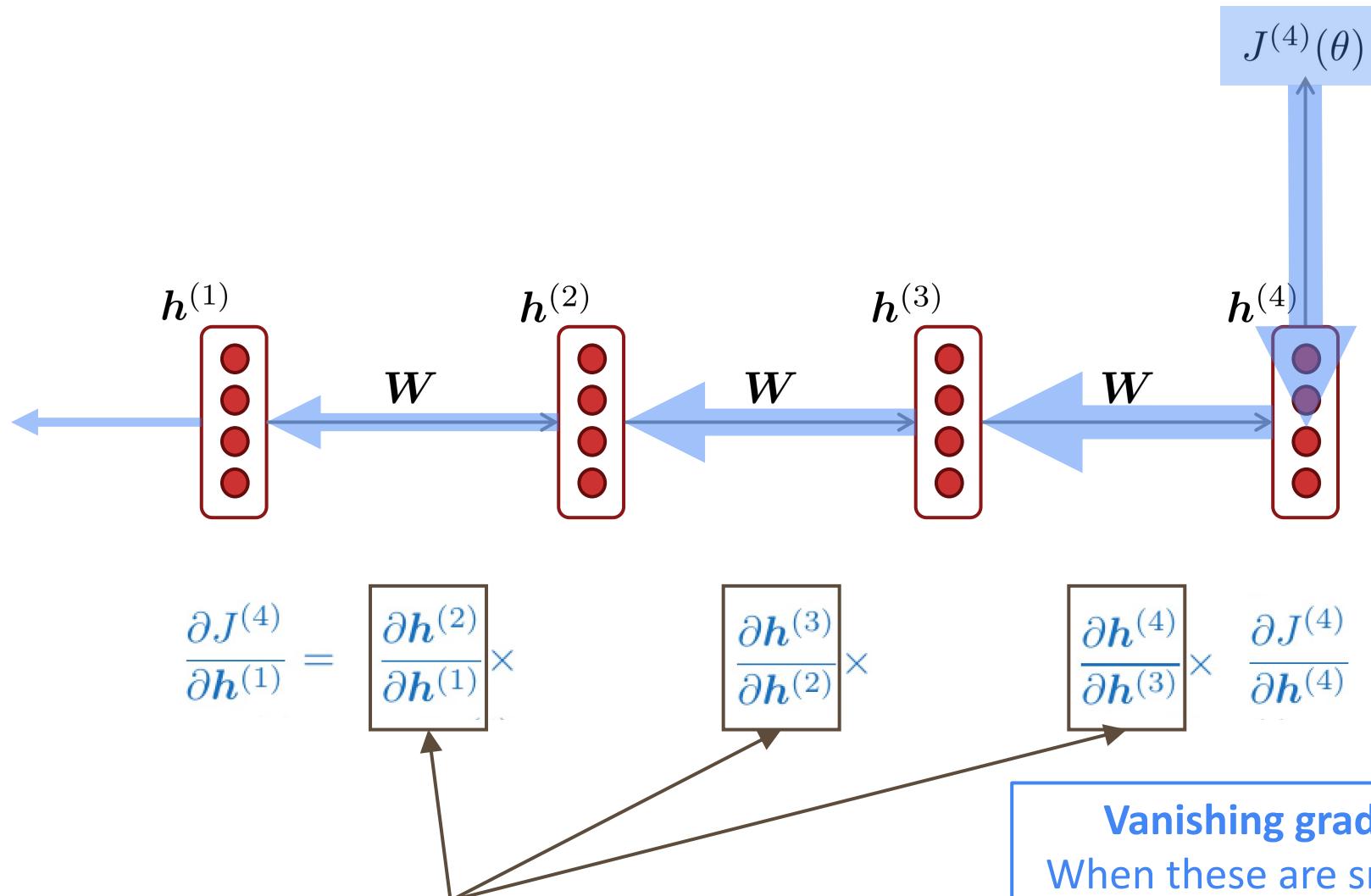
Vanishing gradient intuition



$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = \frac{\partial h^{(2)}}{\partial h^{(1)}} \times \frac{\partial h^{(3)}}{\partial h^{(2)}} \times \frac{\partial h^{(4)}}{\partial h^{(3)}} \times \frac{\partial J^{(4)}}{\partial h^{(4)}}$$

chain rule!

Vanishing gradient intuition



What happens if these are small?

Vanishing gradient problem:
When these are small, the gradient signal gets smaller and smaller as it backpropagates further

Vanishing gradient proof sketch (linear case)

□ Recall:

$$\mathbf{h}^{(t)} = \sigma \left(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)} + \mathbf{b}_1 \right)$$

□ What if σ were the identity function, $\sigma(x) = x$

$$\begin{aligned} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} &= \text{diag} \left(\sigma' \left(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)} + \mathbf{b}_1 \right) \right) \mathbf{W}_h && \text{(chain rule)} \\ &= \mathbf{I} \quad \mathbf{W}_h = \mathbf{W}_h \end{aligned}$$

□ Consider the gradient of the loss $J^{(i)}(\theta)$ on step i , with respect to the hidden state $\mathbf{h}^{(j)}$ on some previous step j . Let $\ell = i - j$

$$\frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(j)}} = \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \prod_{j < t \leq i} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} \quad \text{(chain rule)}$$

$$= \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \prod_{j < t \leq i} \mathbf{W}_h = \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \boxed{\mathbf{W}_h^\ell}$$

↑
(value of $\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}}$)

If \mathbf{W}_h is “small”, then this term gets exponentially problematic as ℓ becomes large

Source: “On the difficulty of training recurrent neural networks”, Pascanu et al, 2013

Vanishing gradient proof sketch (linear case)

- ❑ What's wrong with W_h^ℓ ?
- ❑ Consider if the eigenvalues of W_h are all less than 1

$$\lambda_1, \lambda_2, \dots, \lambda_n < 1$$

$\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n$ (eigenvectors)

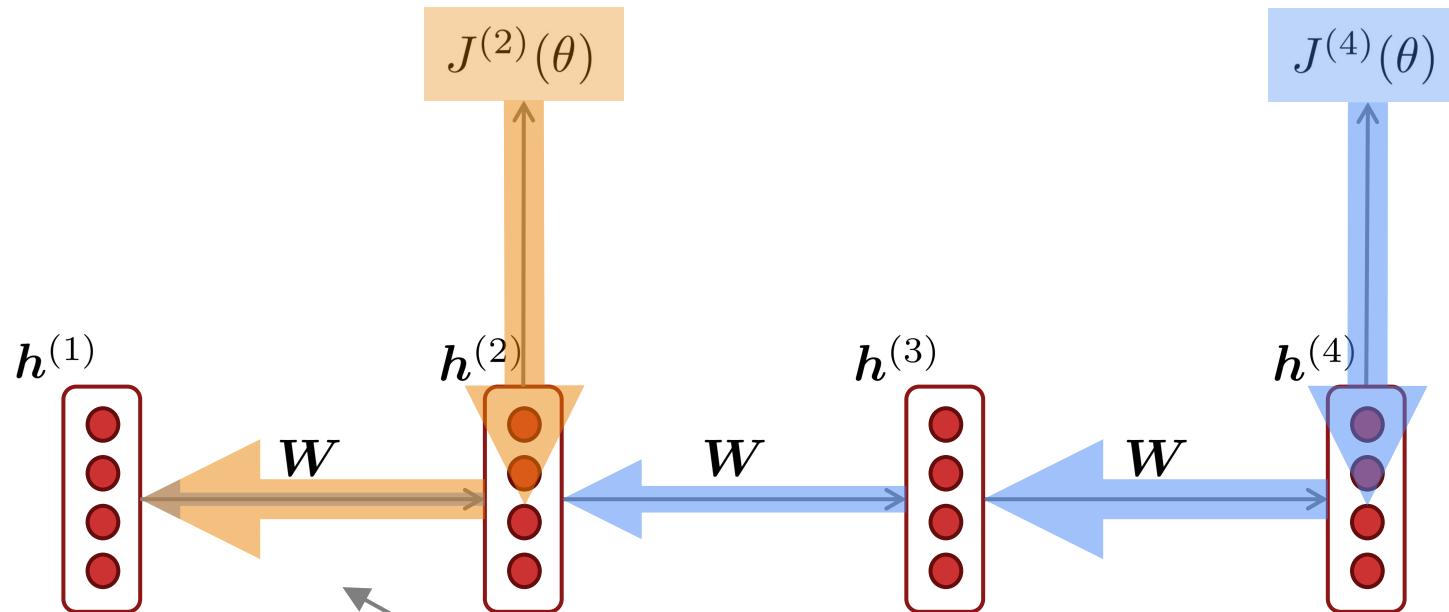
- ❑ We can write $\frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}}$ W_h^ℓ using the eigenvectors of W_h as a basis:

$$\frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} W_h^\ell = \sum_{i=1}^n c_i \boxed{\lambda_i^\ell} \mathbf{q}_i \approx \mathbf{0} \text{ (for large } \ell\text{)}$$

Approaches 0 as ℓ grows, so gradient vanishes

- ❑ What about nonlinear activations σ (i.e., what we use?)
 - ❑ Pretty much the same thing, except the proof requires $\lambda_i < \gamma$ for some γ dependent on dimensionality and σ

Why is vanishing gradient a problem?



Gradient signal from far away is lost because it's much smaller than gradient signal from close-by.

So, model weights are updated only with respect to near effects, not long-term effects.

Why is exploding gradient a problem?

- ❑ If the gradient becomes too big, then the SGD update step becomes too big:

$$\theta^{new} = \theta^{old} - \underbrace{\alpha \nabla_{\theta} J(\theta)}_{\text{gradient}}$$

learning rate

- ❑ This can cause **bad updates**: we take too large a step and reach a weird and bad parameter configuration (with large loss)
- ❑ In the worst case, this will result in **Inf** or **NaN** in your network (then you have to restart training from an earlier checkpoint)

Gradient clipping: solution for exploding gradient

- ❑ **Gradient clipping:** if the norm of the gradient is greater than some threshold, scale it down before applying SGD update

Algorithm 1 Pseudo-code for norm clipping

```
 $\hat{\mathbf{g}} \leftarrow \frac{\partial \mathcal{E}}{\partial \theta}$ 
if  $\|\hat{\mathbf{g}}\| \geq \text{threshold}$  then
     $\hat{\mathbf{g}} \leftarrow \frac{\text{threshold}}{\|\hat{\mathbf{g}}\|} \hat{\mathbf{g}}$ 
end if
```

- ❑ **Intuition:** take a step in the same direction, but a smaller step
- ❑ In practice, remembering to clip gradients is important, but exploding gradients are an easy problem to solve

How to fix the vanishing gradient problem?

- ❑ The main problem is that *it's too difficult for the RNN to learn to preserve information over many timesteps.*
- ❑ In a vanilla RNN, the hidden state is constantly being **rewritten**

$$\mathbf{h}^{(t)} = \sigma \left(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)} + \mathbf{b} \right)$$

- ❑ How about an RNN with separate **memory** which is added to?

Long Short-Term Memory RNNs (LSTMs)

- A type of RNN proposed by Hochreiter and Schmidhuber in 1997 as a solution to the vanishing gradients problem.
 - Everyone cites that paper but really a crucial part of the modern LSTM is from Gers et al. (2000)
- On step t, there is a **hidden state** $h^{(t)}$ and a **cell state** $c^{(t)}$
 - Both are vectors length n
 - The cell stores **long-term information**
 - The LSTM can **read**, **erase**, and **write** information from the cell
 - The cell becomes conceptually rather like RAM in a computer

“Long short-term memory”, Hochreiter and Schmidhuber, 1997. <https://www.bioinf.jku.at/publications/older/2604.pdf>

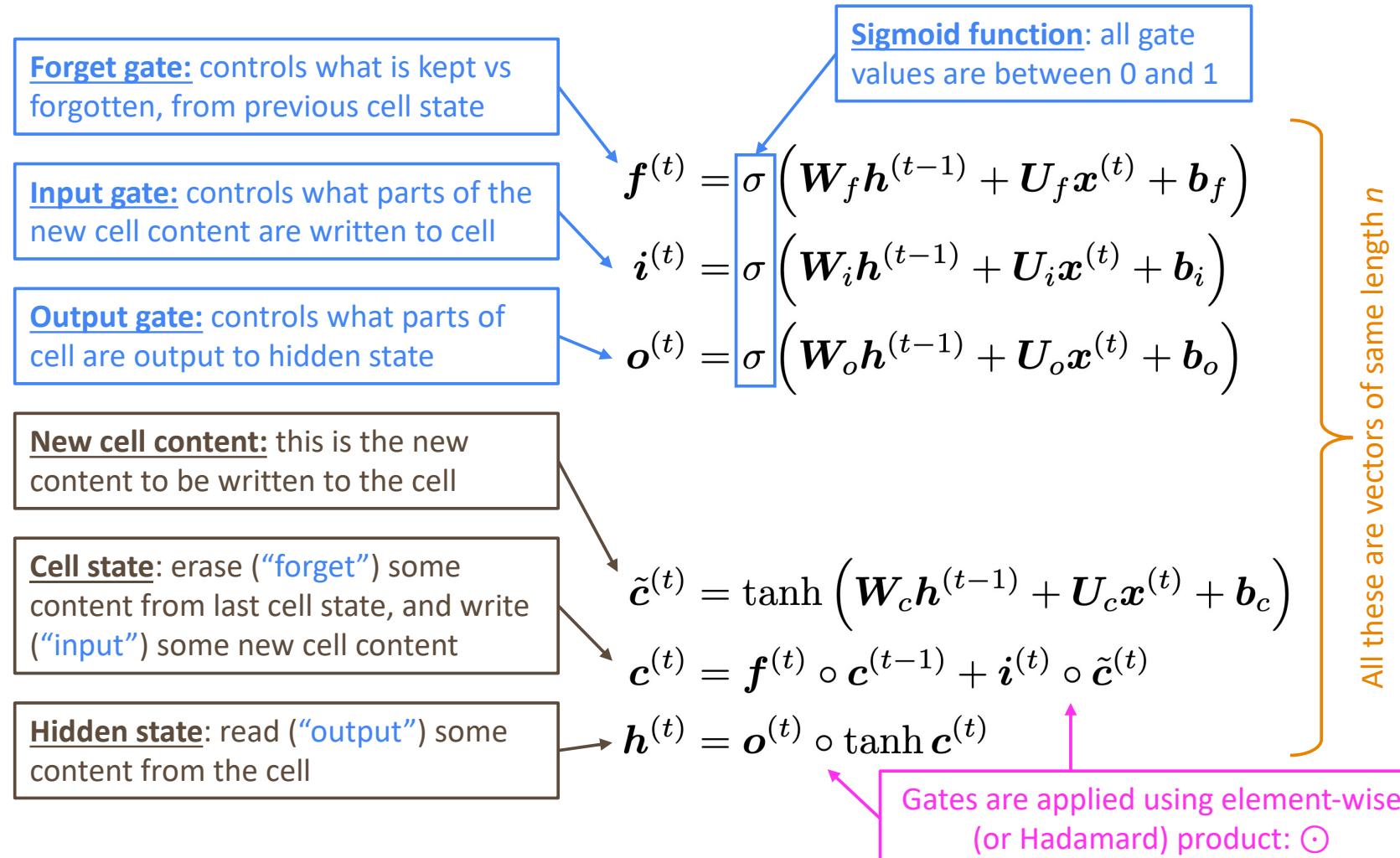
“Learning to Forget: Continual Prediction with LSTM”, Gers, Schmidhuber, and Cummins, 2000. <https://dl.acm.org/doi/10.1162/089976600300015015>

Long Short-Term Memory RNNs (LSTMs)

- On step t , there is a **hidden state** $h^{(t)}$ *and* a **cell state** $c^{(t)}$
 - Both are vectors length n
 - The cell stores long-term information
 - The LSTM can read, erase, and write information from the cell
 - The cell becomes conceptually rather like RAM in a computer
- The selection of which information is erased/written/read is controlled by three corresponding gates
 - The gates are also vectors of length n
 - On each timestep, each element of the gates can be **open** (1), **closed** (0), or somewhere in-between
 - The gates are **dynamic**: their value is computed based on the current context

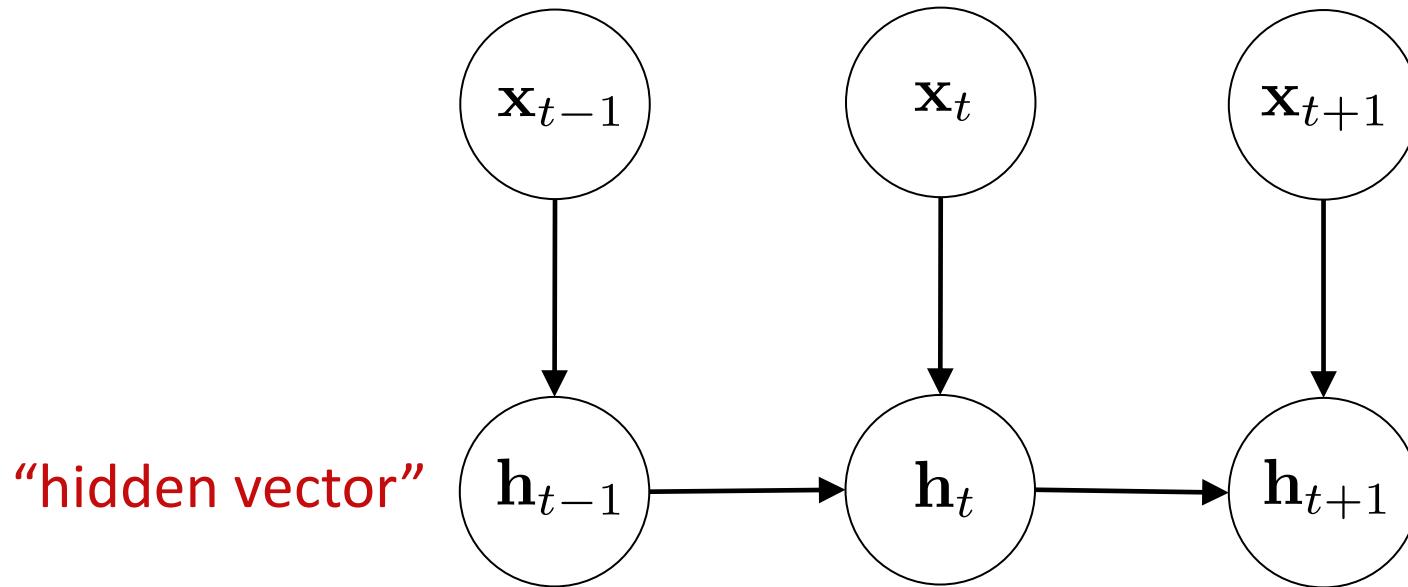
Long Short-Term Memory RNNs (LSTMs)

- We have a sequence of inputs $x^{(t)}$, and we will compute a sequence of hidden states $h^{(t)}$ and cell states $c^{(t)}$. On timestep t :



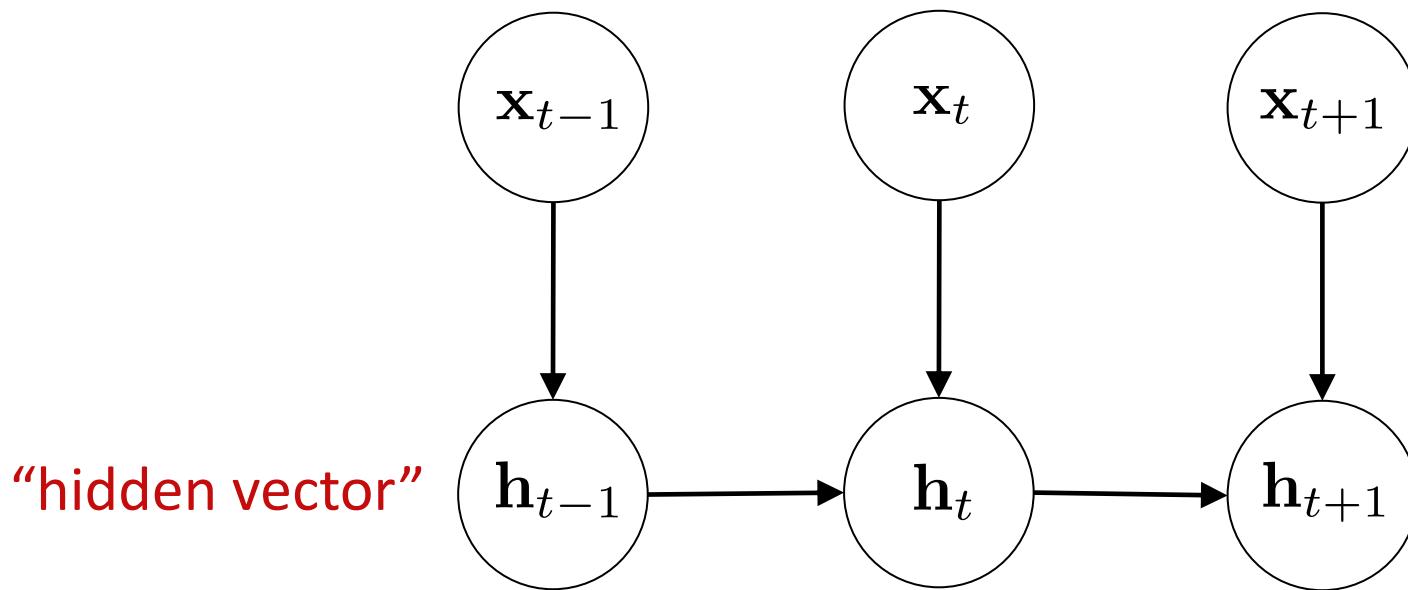
Recurrent Neural Networks

- ❑ Input is a sequence:

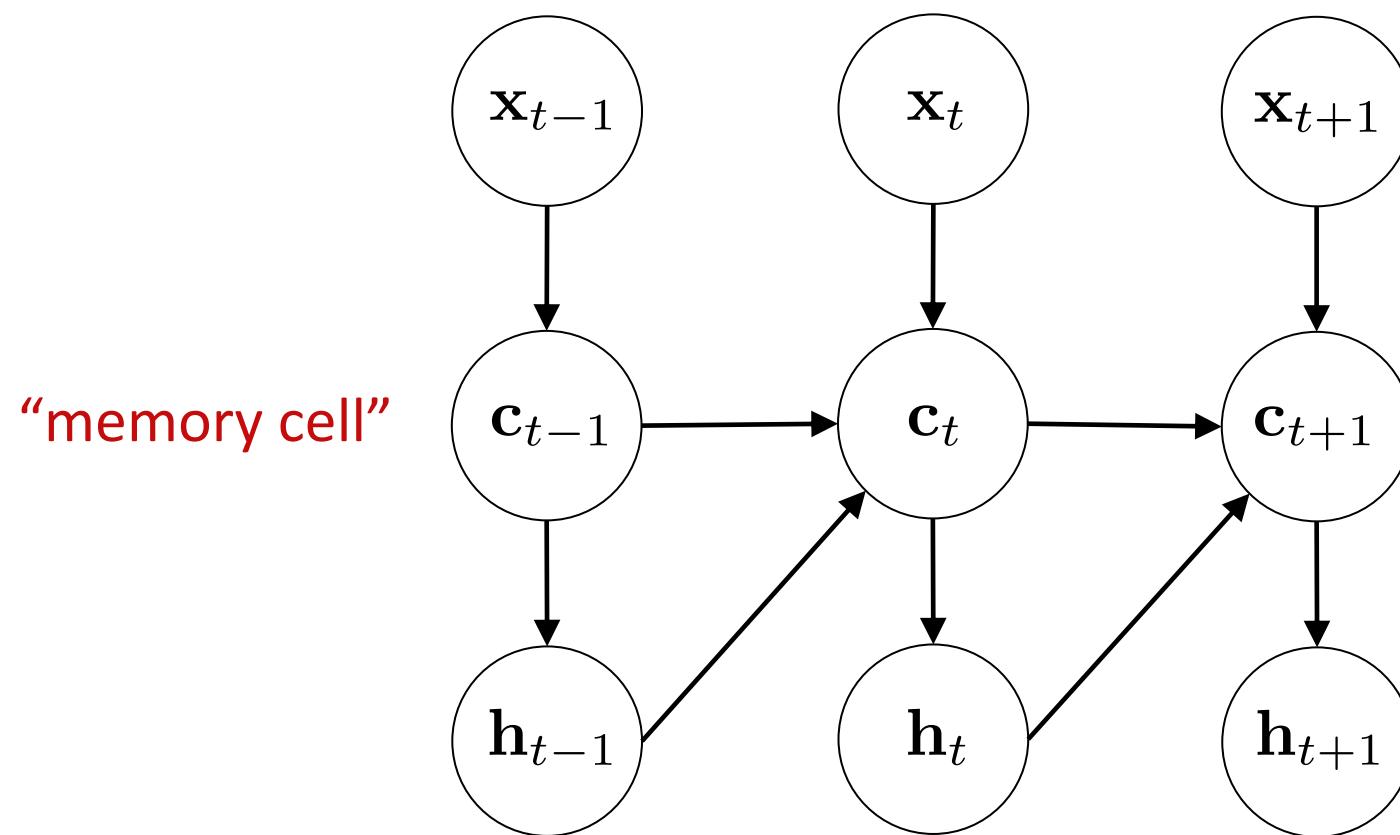


Recurrent Neural Networks

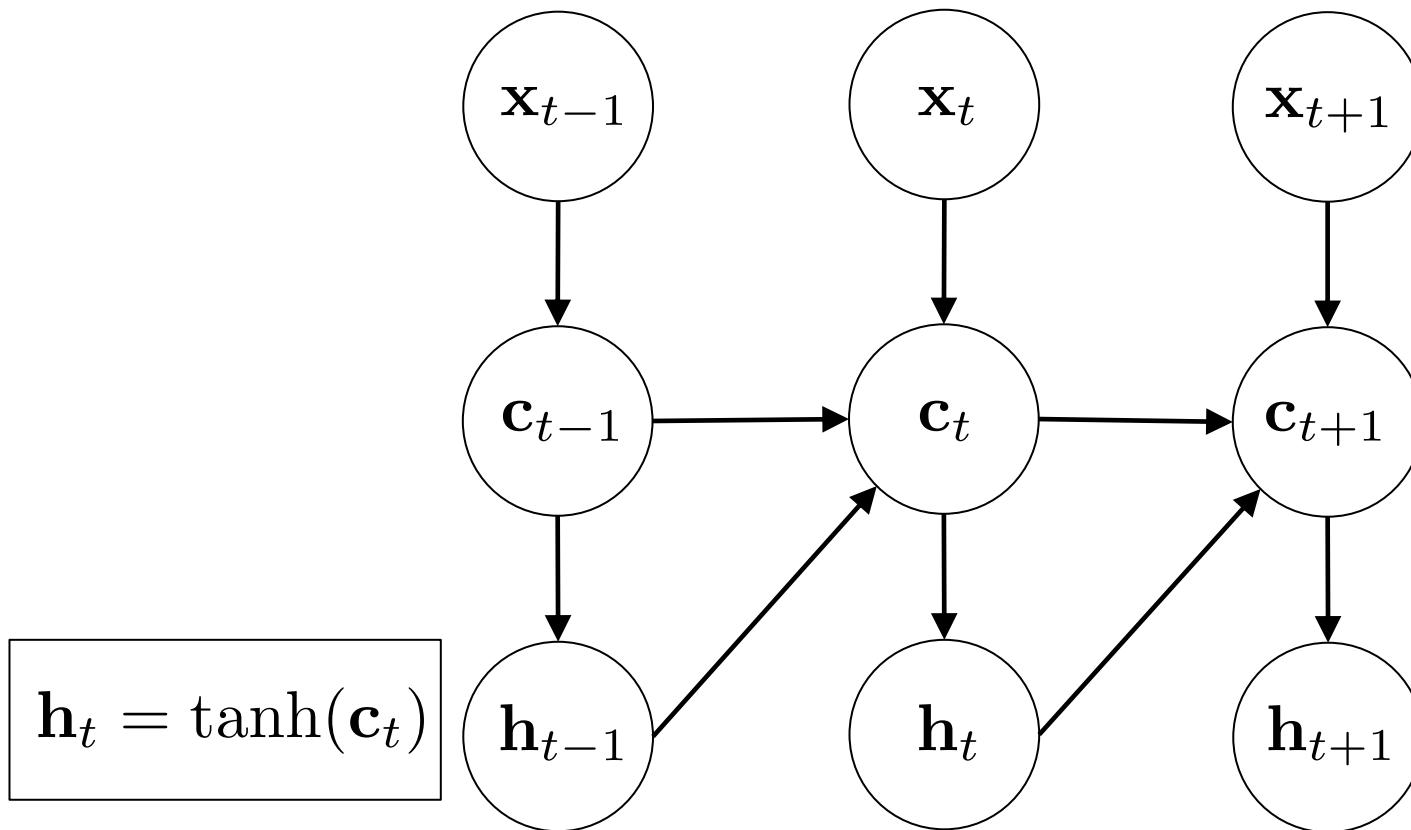
$$\mathbf{h}_t = \tanh \left(\mathbf{W}^{(x)} \mathbf{x}_t + \mathbf{W}^{(h)} \mathbf{h}_{t-1} + \mathbf{b} \right)$$



Long Short-Term Memory Networks (gateless)

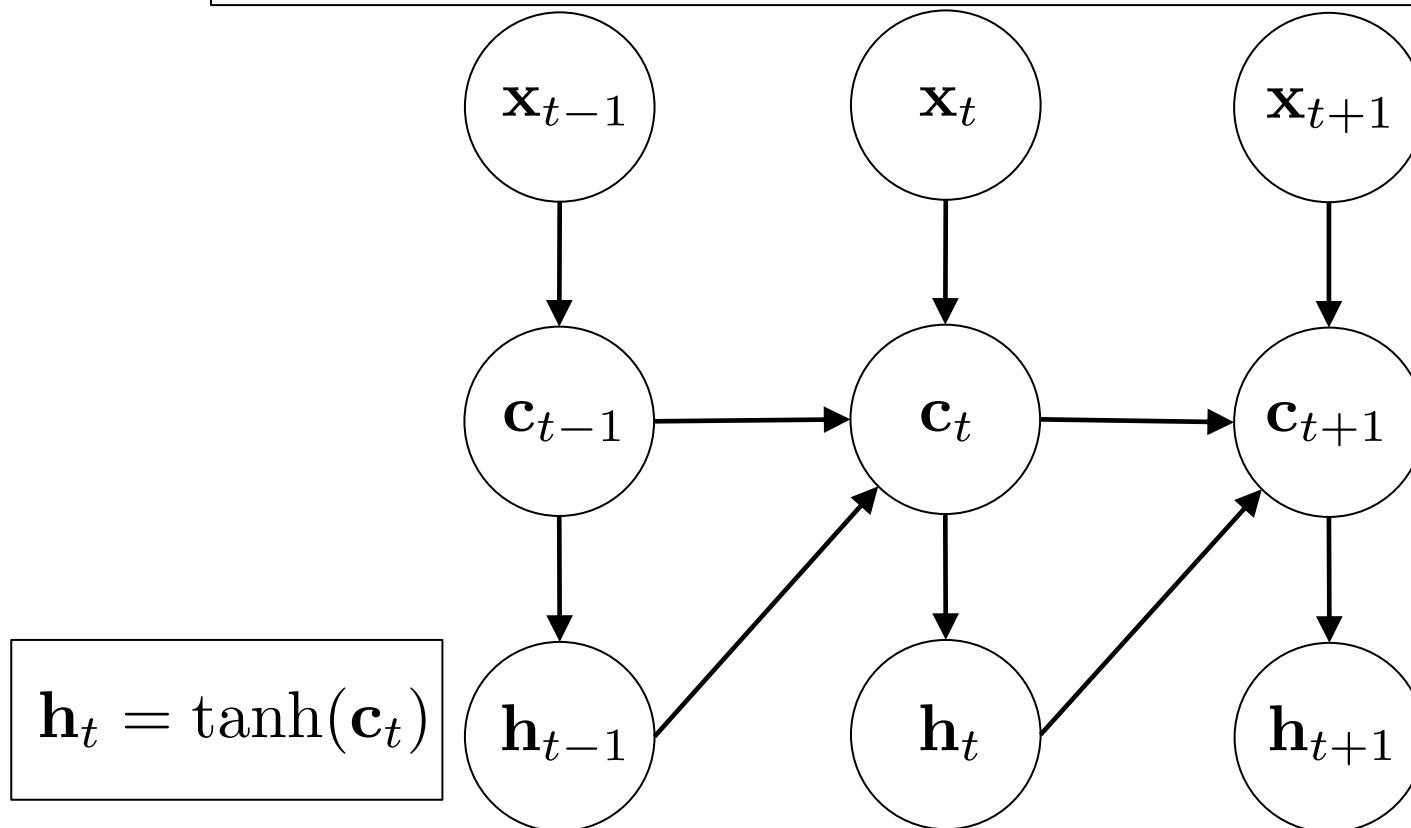


Long Short-Term Memory Networks (gateless)



Long Short-Term Memory Networks (gateless)

$$\mathbf{c}_t = \mathbf{c}_{t-1} + \tanh \left(\mathbf{W}^{(xc)} \mathbf{x}_t + \mathbf{W}^{(hc)} \mathbf{h}_{t-1} + \mathbf{b}^{(c)} \right)$$



Long Short-Term Memory Networks (gateless)

$$\mathbf{c}_t = \mathbf{c}_{t-1} + \tanh \left(\mathbf{W}^{(xc)} \mathbf{x}_t + \mathbf{W}^{(hc)} \mathbf{h}_{t-1} + \mathbf{b}^{(c)} \right)$$

Experiment: text classification

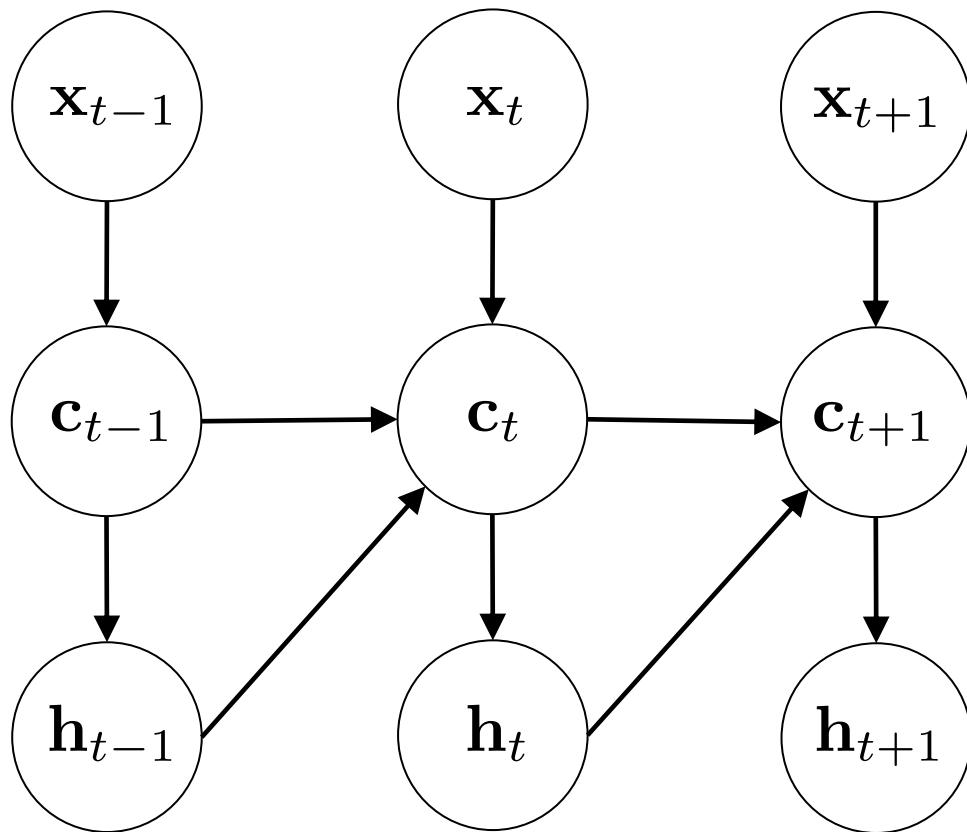
- Stanford Sentiment Treebank
 - binary classification (positive/negative)
- 25-dim word vectors
- 50-dim cell/hidden vectors
- classification layer on **final** hidden vector
- AdaGrad, 10 epochs, mini-batch size 10
- early stopping on dev set

accuracy

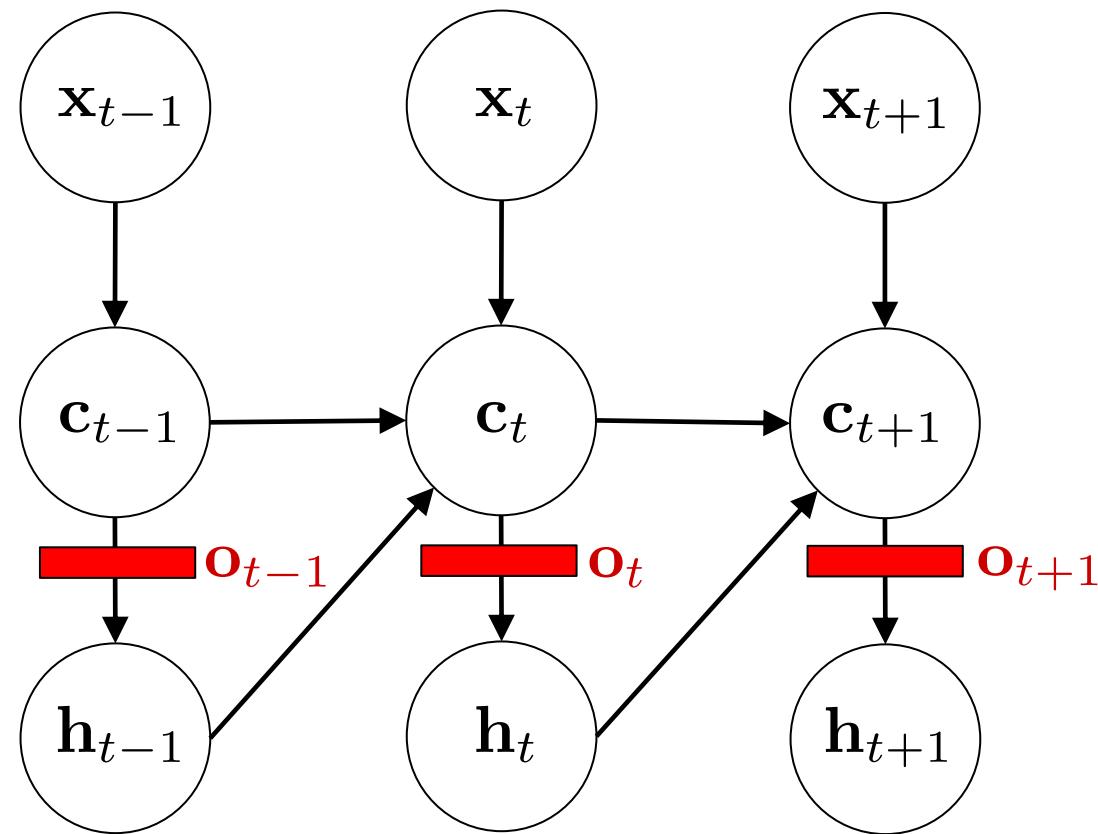
80.6



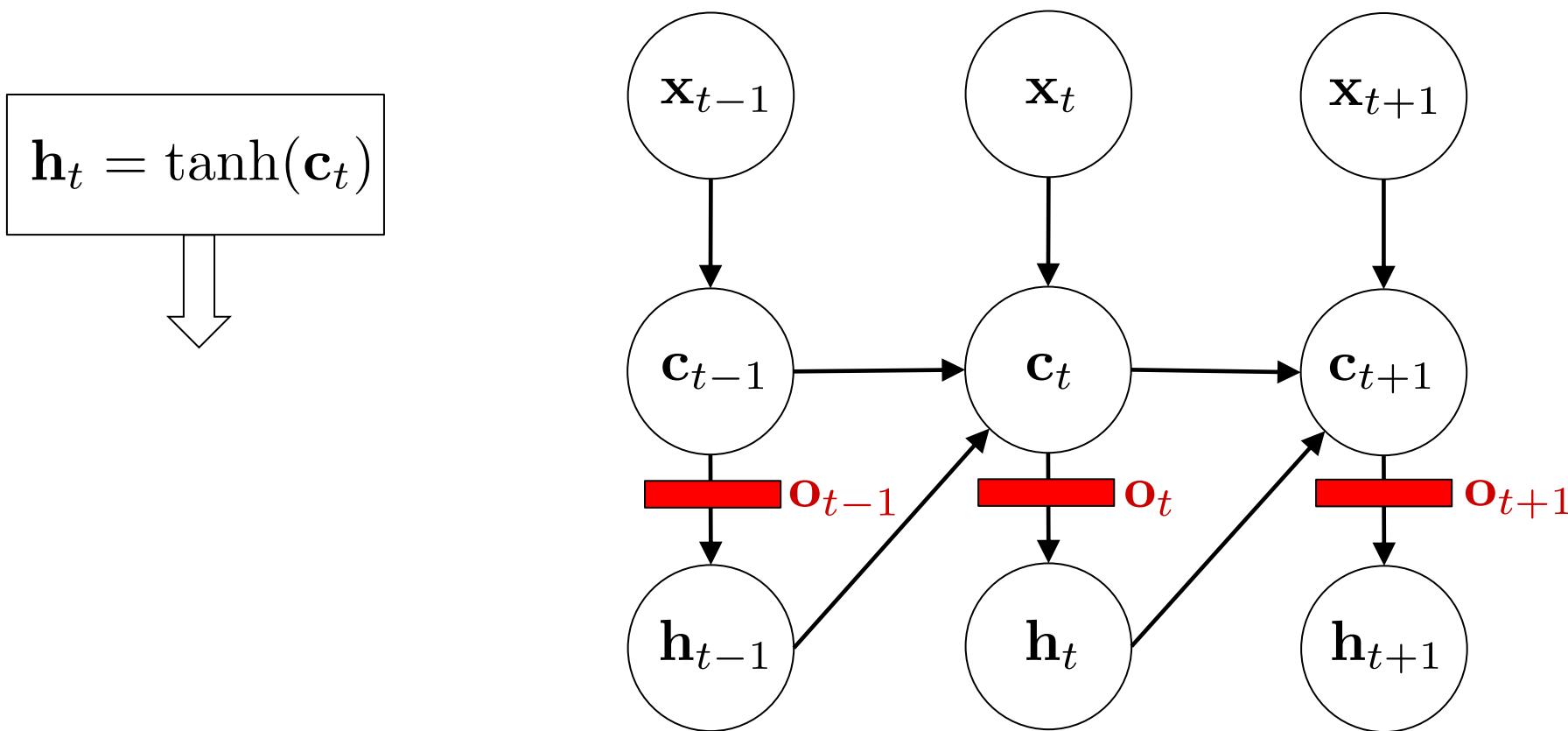
Adding Output Gates



Adding Output Gates

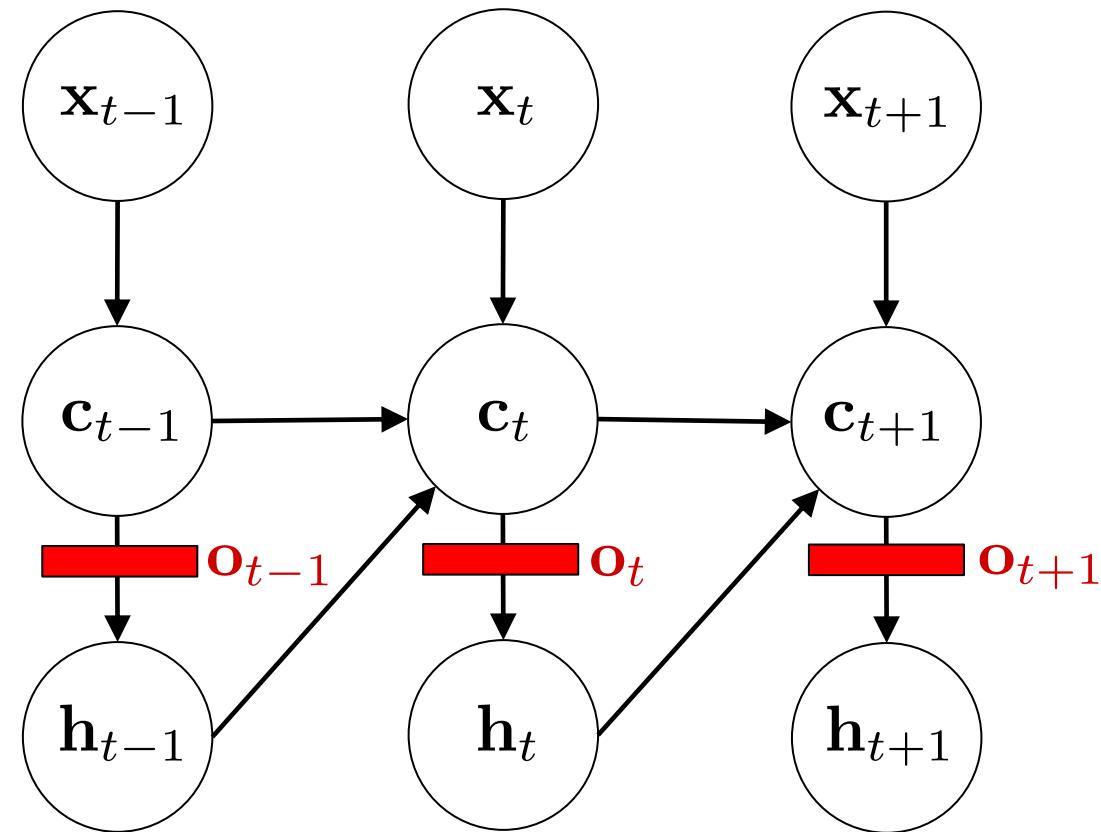


Adding Output Gates



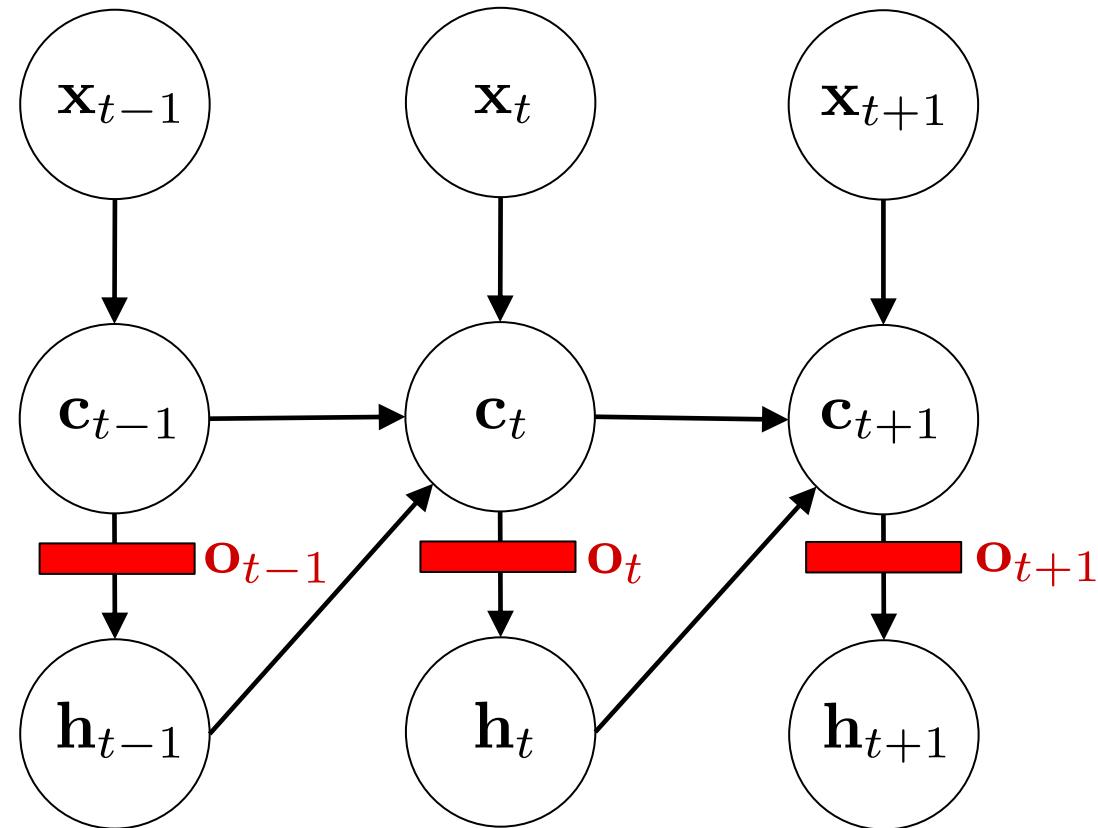
Adding Output Gates

$$\begin{array}{l} \boxed{\mathbf{h}_t = \tanh(\mathbf{c}_t)} \\ \downarrow \\ \boxed{\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)} \end{array}$$



Adding Output Gates

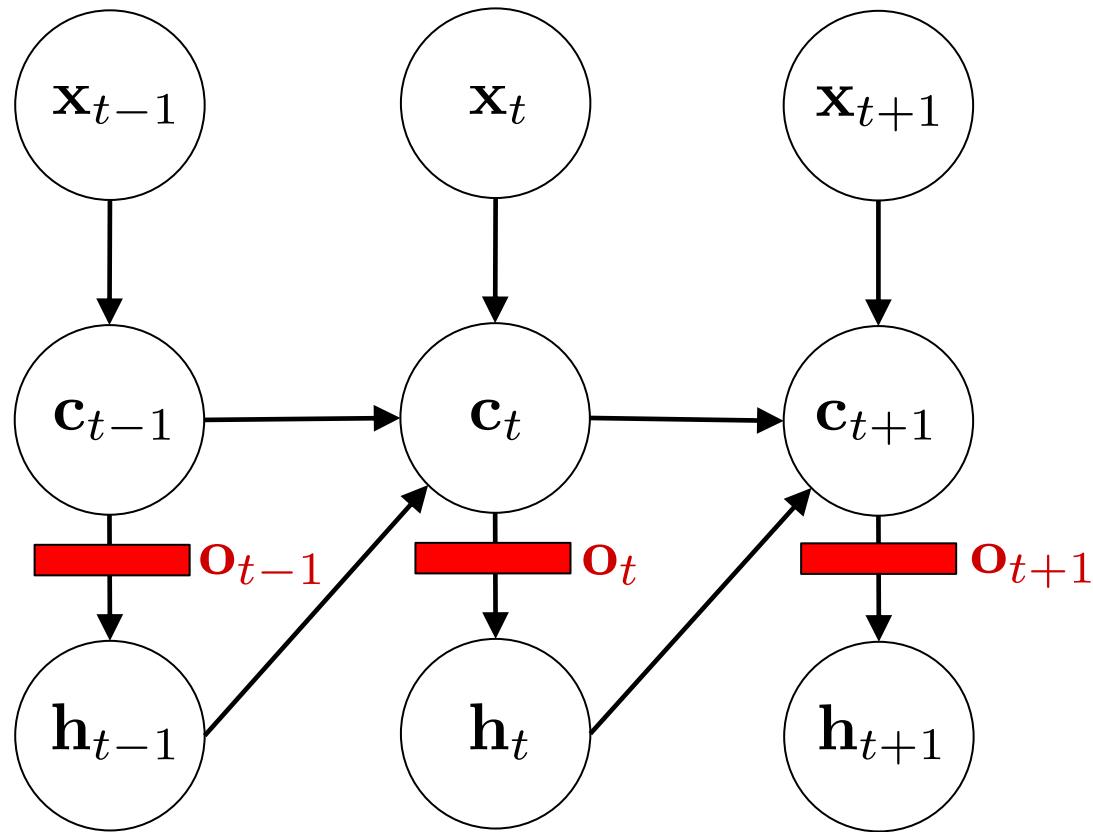
$$\begin{array}{l} \boxed{\mathbf{h}_t = \tanh(\mathbf{c}_t)} \\ \downarrow \\ \boxed{\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)} \\ \text{this is pointwise multiplication!} \\ \mathbf{o}_t \text{ is a vector} \end{array}$$



Adding Output Gates

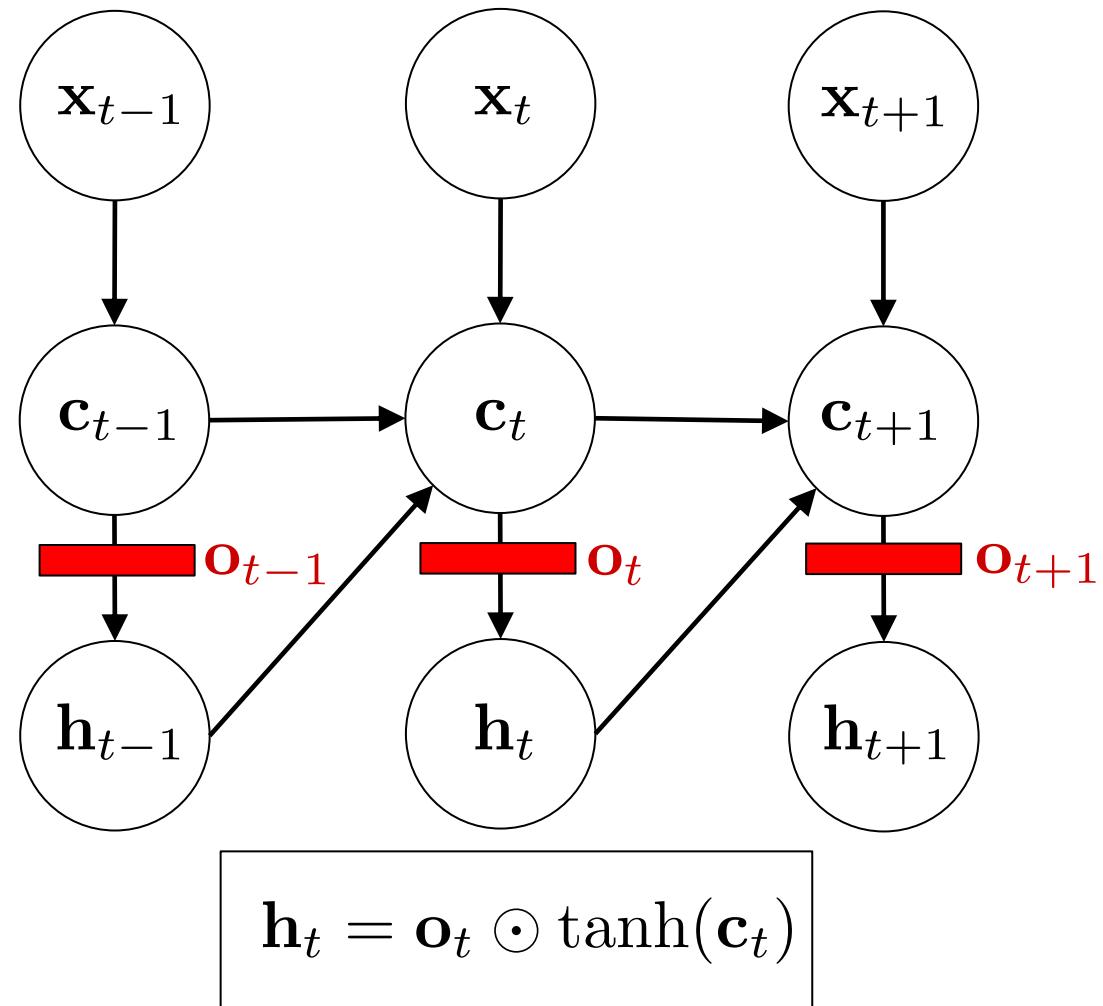
$$\begin{array}{l} \boxed{\mathbf{h}_t = \tanh(\mathbf{c}_t)} \\ \downarrow \\ \boxed{\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)} \end{array}$$

output gate affects how much “information” is transmitted from cell vector to hidden vector



Adding Output Gates

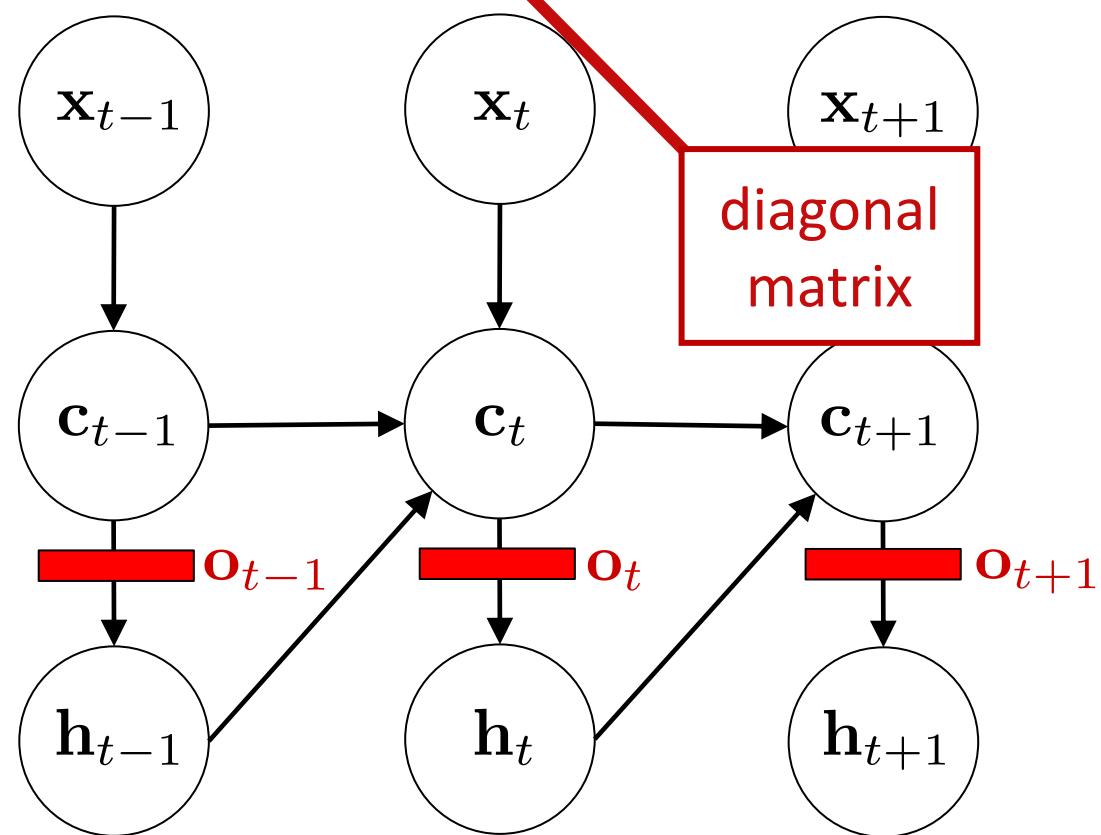
$$\mathbf{o}_t = \sigma \left(\mathbf{W}^{(xo)} \mathbf{x}_t + \mathbf{W}^{(ho)} \mathbf{h}_{t-1} + \mathbf{W}^{(co)} \mathbf{c}_t + \mathbf{b}^{(o)} \right)$$



Adding Output Gates

$$\mathbf{o}_t = \sigma \left(\mathbf{W}^{(xo)} \mathbf{x}_t + \mathbf{W}^{(ho)} \mathbf{h}_{t-1} + \mathbf{W}^{(co)} \mathbf{c}_t + \mathbf{b}^{(o)} \right)$$

logistic sigmoid, so output ranges from 0 to 1

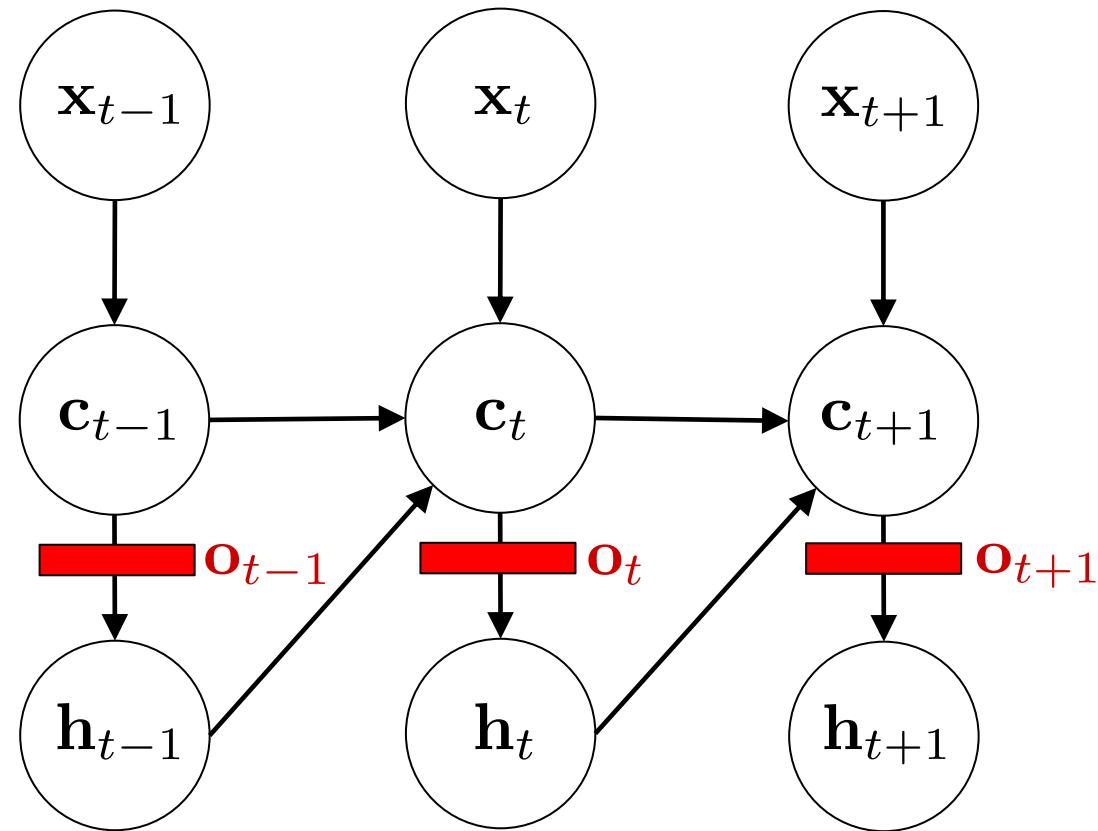


$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)$$

Adding Output Gates

$$\mathbf{o}_t = \sigma \left(\mathbf{W}^{(xo)} \mathbf{x}_t + \mathbf{W}^{(ho)} \mathbf{h}_{t-1} + \mathbf{W}^{(co)} \mathbf{c}_t + \mathbf{b}^{(o)} \right)$$

output gate is a function
of current observation,
previous hidden vector,
and current cell vector



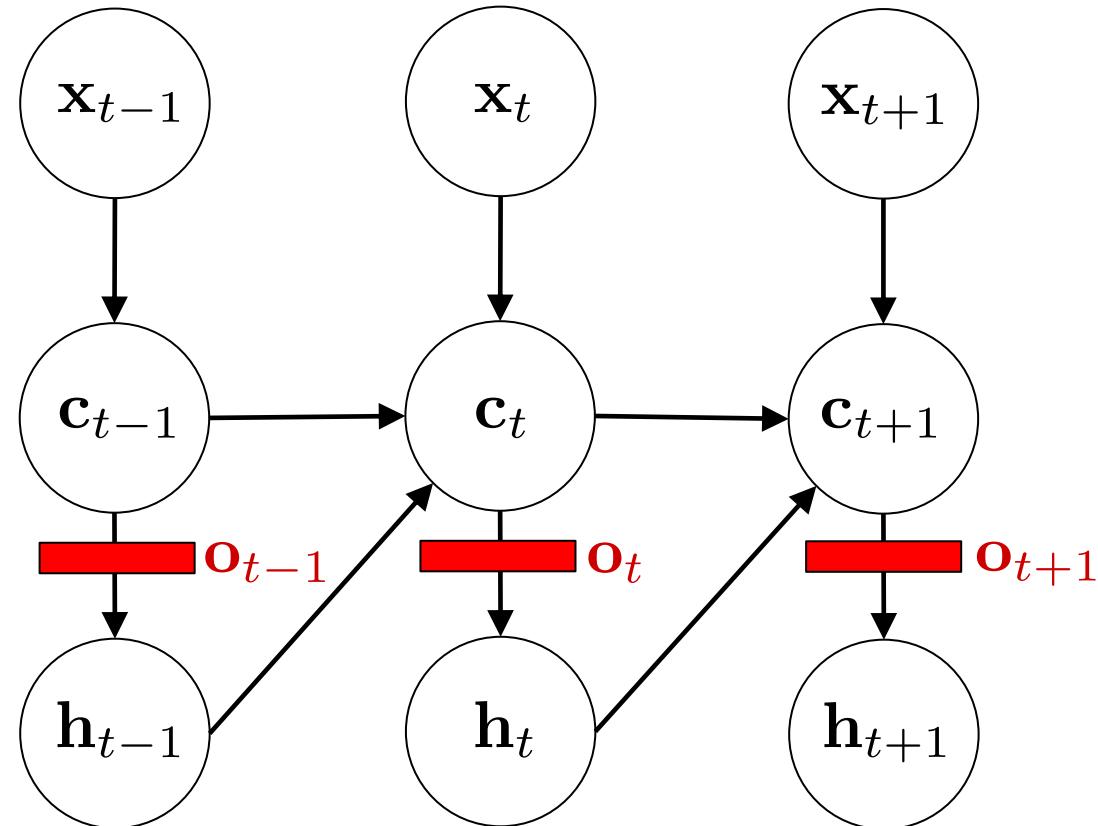
$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)$$

Adding Output Gates

$$\mathbf{o}_t = \sigma \left(\mathbf{W}^{(xo)} \mathbf{x}_t + \mathbf{W}^{(ho)} \mathbf{h}_{t-1} + \mathbf{W}^{(co)} \mathbf{c}_t + \mathbf{b}^{(o)} \right)$$

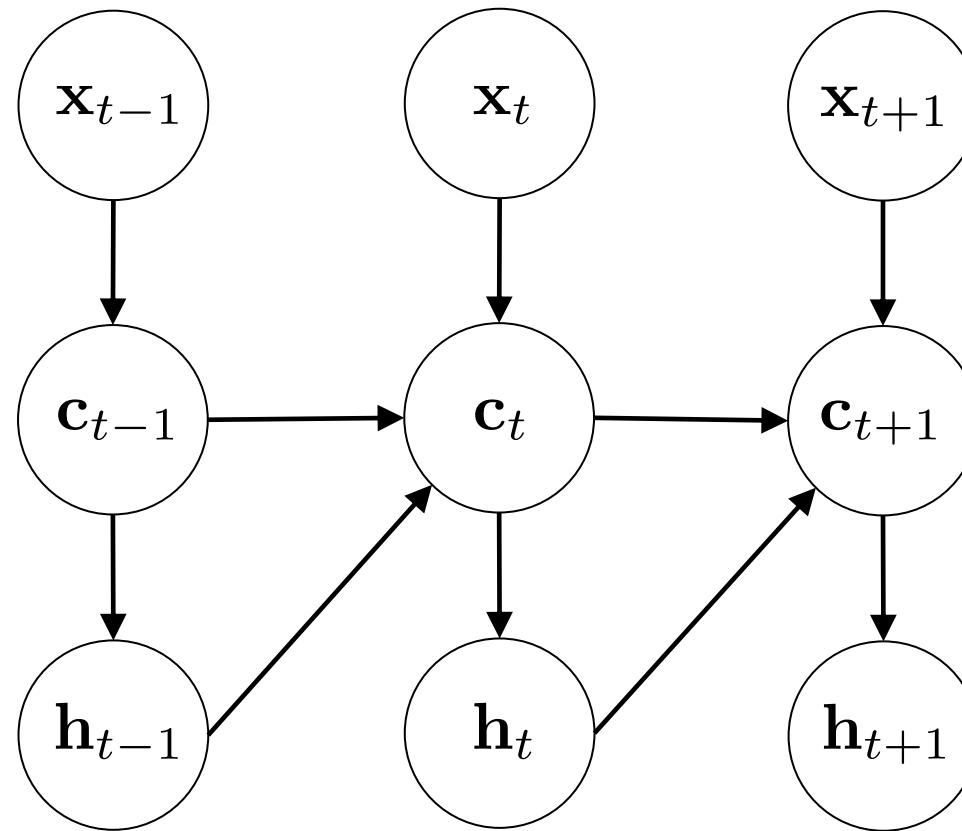
output gate is a function
of current observation,
previous hidden vector,
and current cell vector

	acc.
gateless	80.6
output gates	81.9

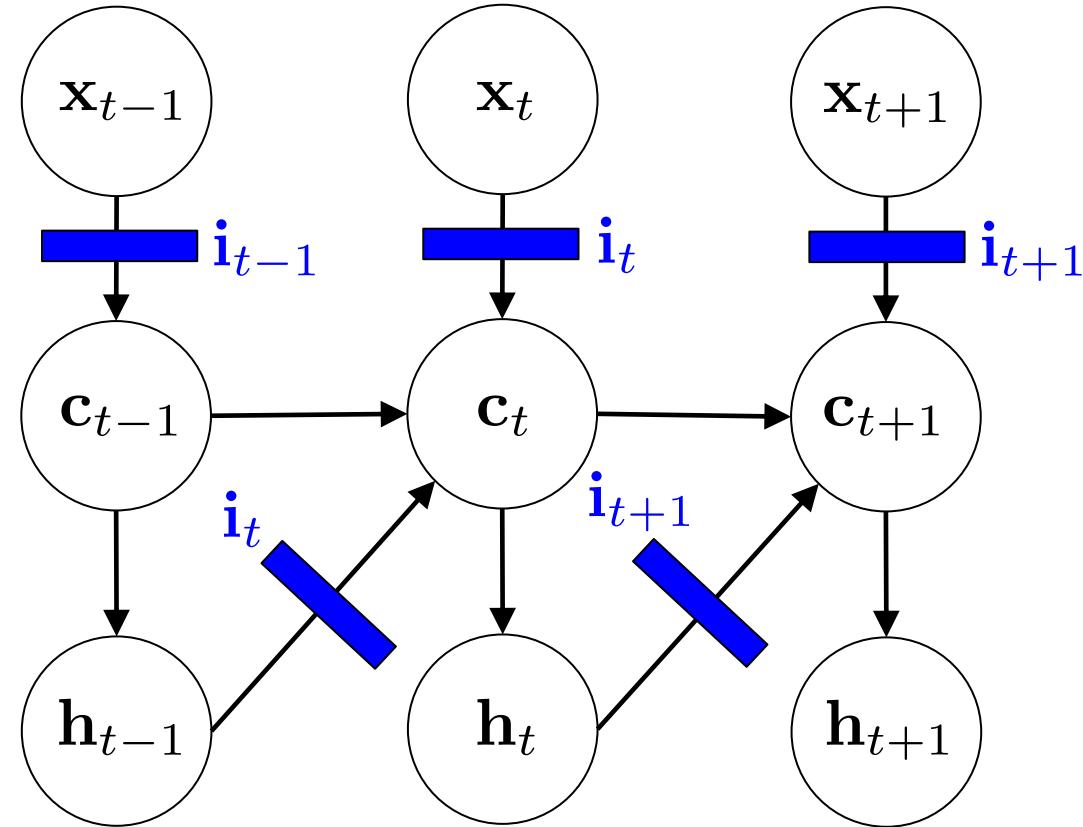


$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)$$

Adding Input Gates



Adding Input Gates

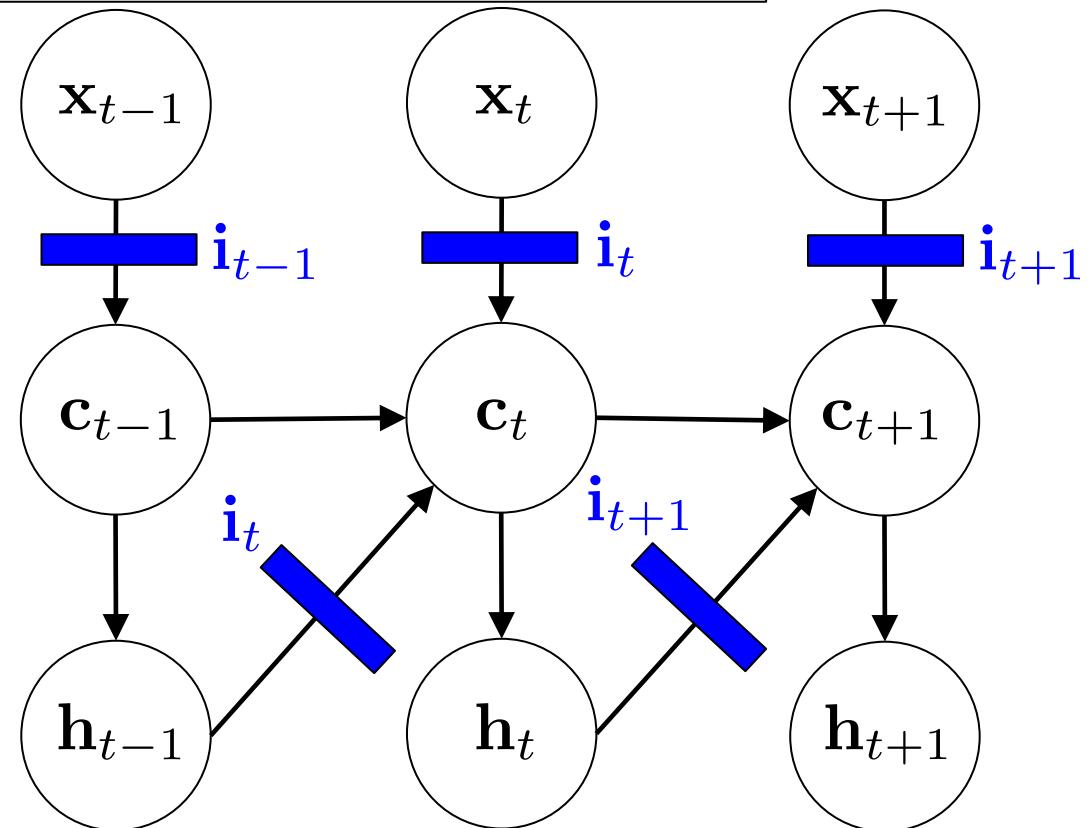


Adding Input Gates

$$\mathbf{c}_t = \mathbf{c}_{t-1} + \tanh(\mathbf{W}^{(xc)} \mathbf{x}_t + \mathbf{W}^{(hc)} \mathbf{h}_{t-1} + \mathbf{b}^{(c)})$$

$$\mathbf{c}_t = \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tanh(\mathbf{W}^{(xc)} \mathbf{x}_t + \mathbf{W}^{(hc)} \mathbf{h}_{t-1} + \mathbf{b}^{(c)})$$

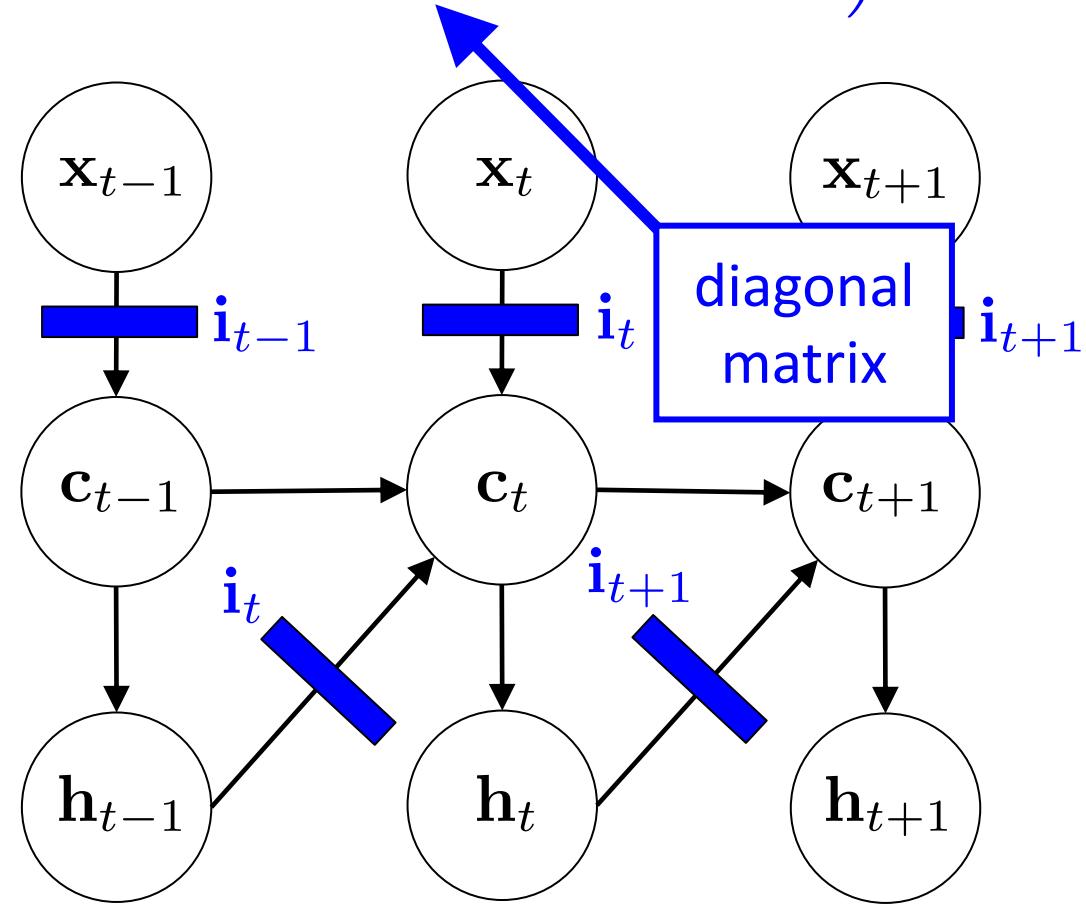
input gate controls how much cell is affected by current observation and previous hidden vector



Adding Input Gates

$$\mathbf{i}_t = \sigma \left(\mathbf{W}^{(xi)} \mathbf{x}_t + \mathbf{W}^{(hi)} \mathbf{h}_{t-1} + \mathbf{W}^{(ci)} \mathbf{c}_{t-1} + \mathbf{b}^{(i)} \right)$$

input gate is a function of current observation, previous hidden vector, and previous cell vector

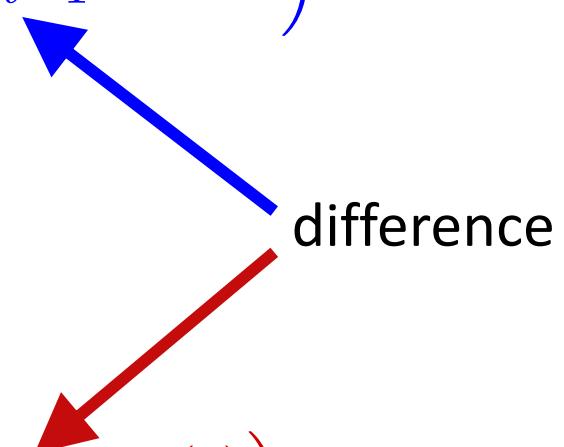


Input Gates

$$\mathbf{i}_t = \sigma \left(\mathbf{W}^{(xi)} \mathbf{x}_t + \mathbf{W}^{(hi)} \mathbf{h}_{t-1} + \mathbf{W}^{(ci)} \mathbf{c}_{t-1} + \mathbf{b}^{(i)} \right)$$

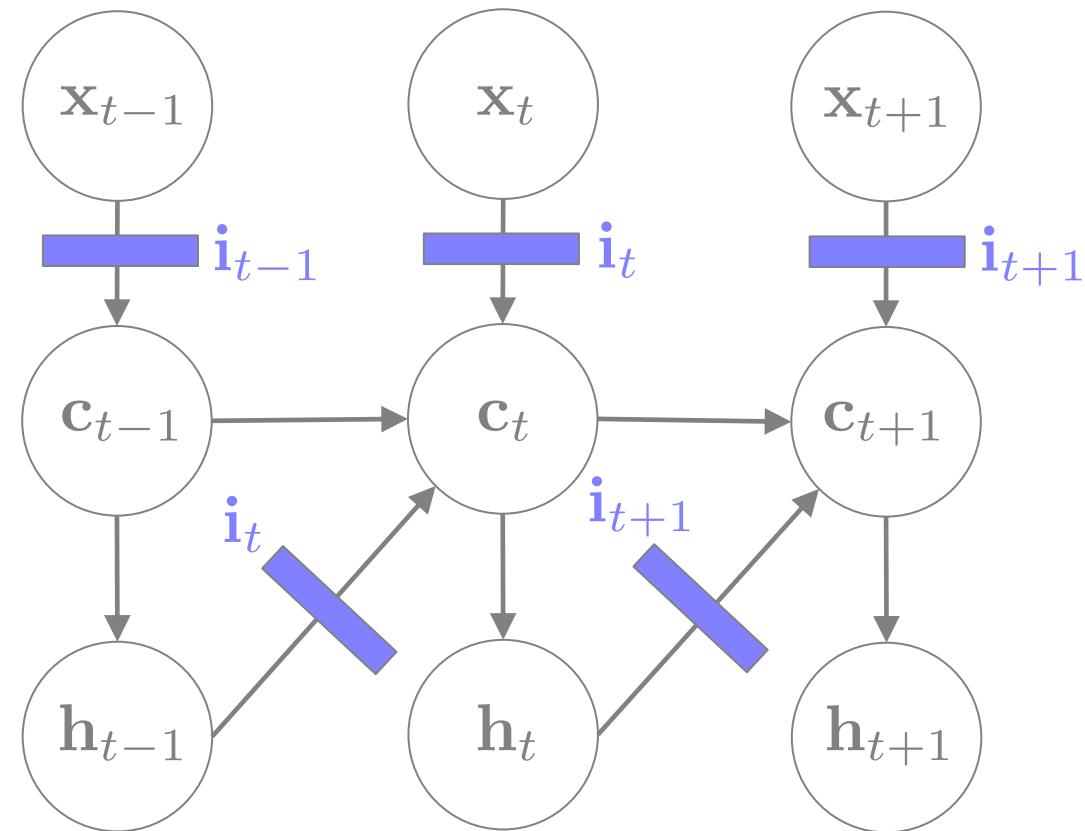
Output Gates

$$\mathbf{o}_t = \sigma \left(\mathbf{W}^{(xo)} \mathbf{x}_t + \mathbf{W}^{(ho)} \mathbf{h}_{t-1} + \mathbf{W}^{(co)} \mathbf{c}_t + \mathbf{b}^{(o)} \right)$$



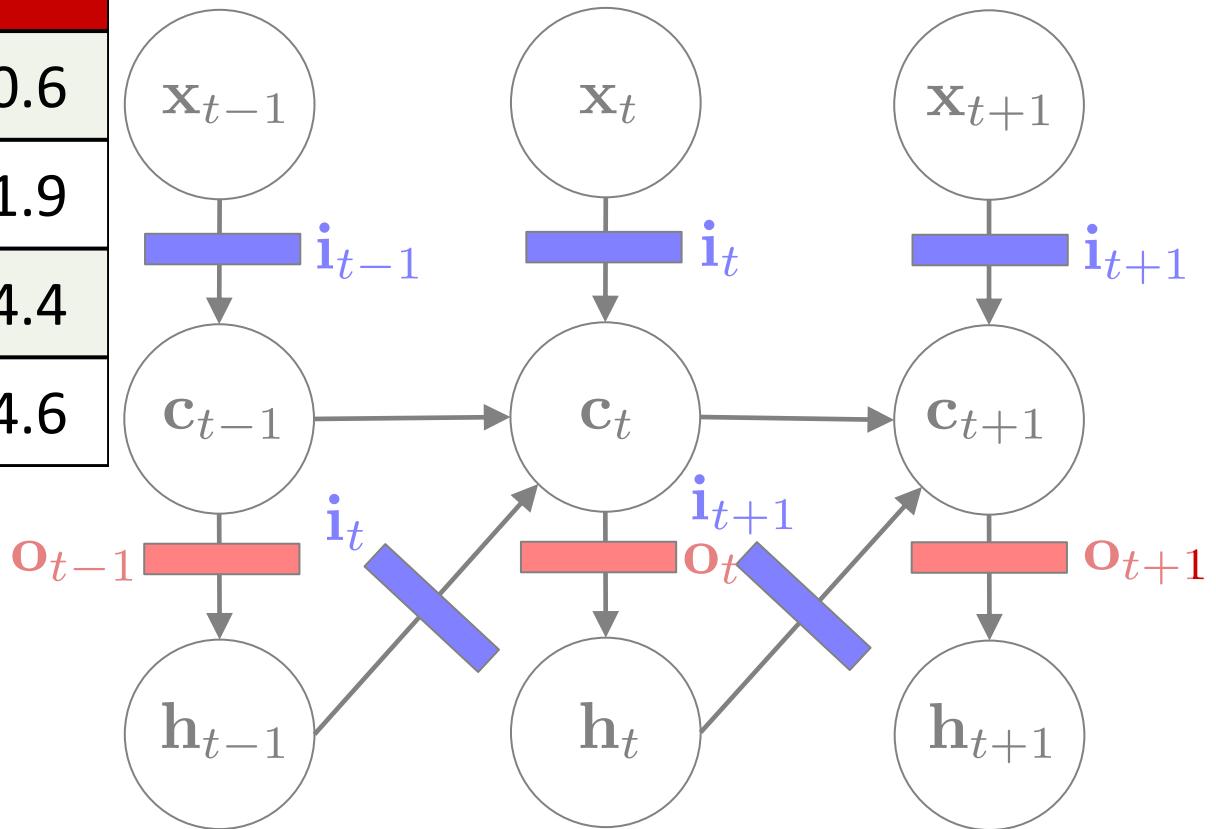
Input Gates

	acc.
gateless	80.6
output gates	81.9
input gates	84.4

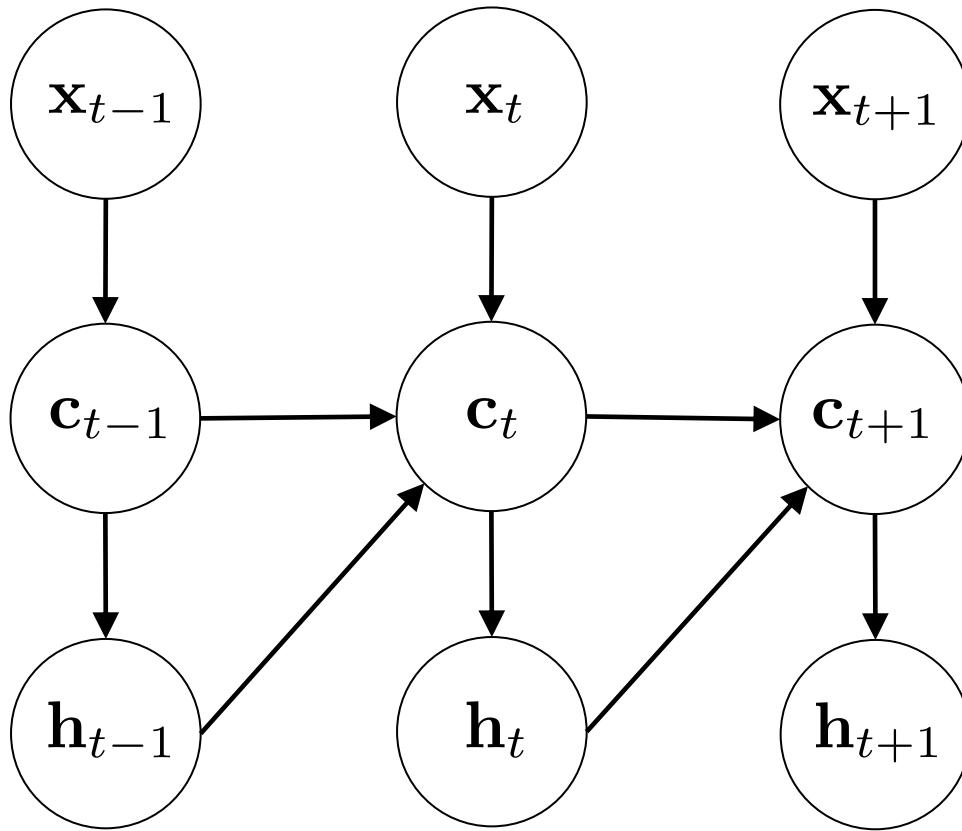


Input Gates

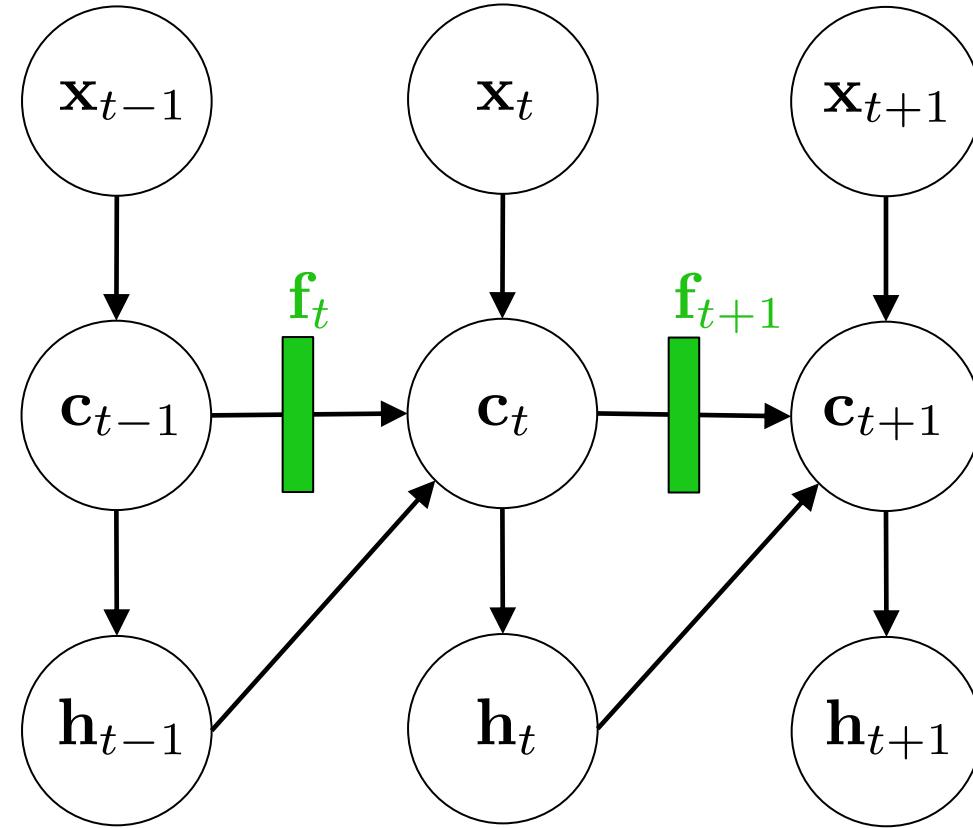
	acc.
gateless	80.6
output gates	81.9
input gates	84.4
input & output gates	84.6



Adding Forget Gates



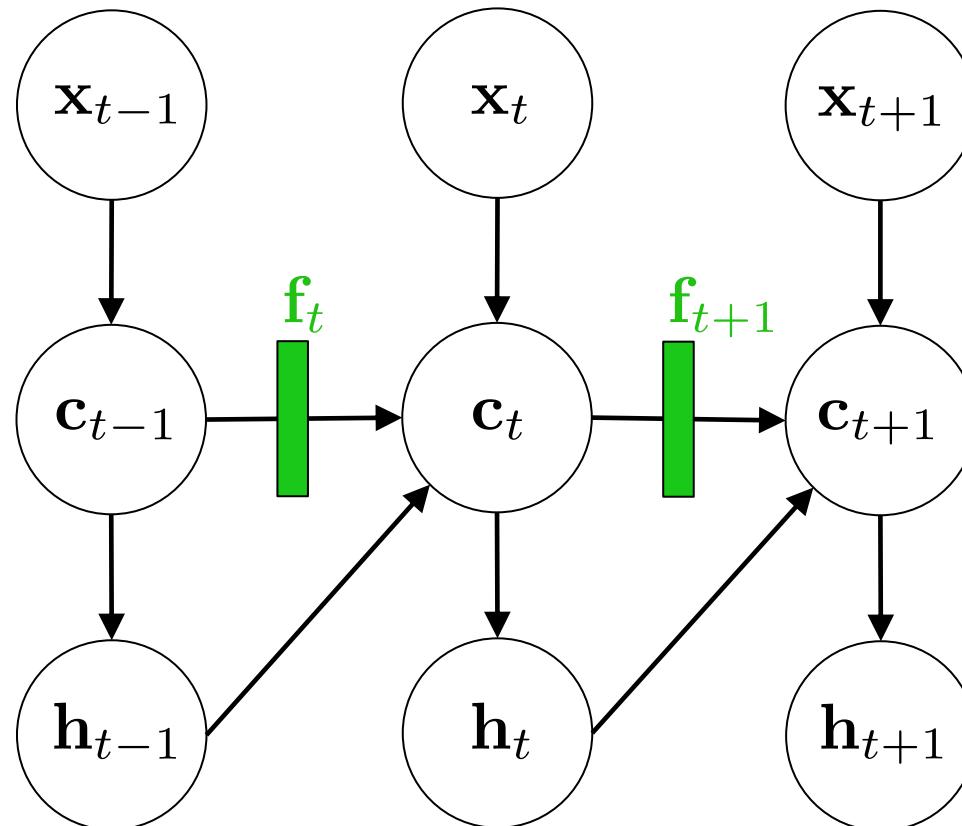
Adding Forget Gates



Adding Forget Gates

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \tanh \left(\mathbf{W}^{(xc)} \mathbf{x}_t + \mathbf{W}^{(hc)} \mathbf{h}_{t-1} + \mathbf{b}^{(c)} \right)$$

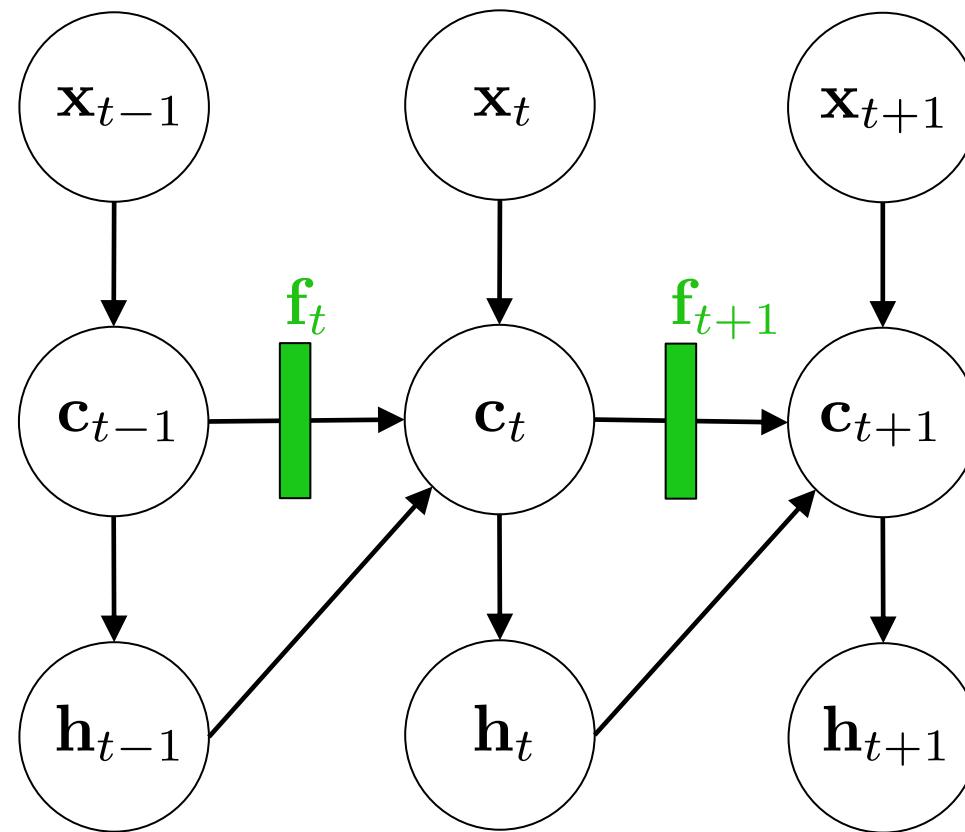
forget gate controls how much “information” is kept from the previous cell vector



Adding Forget Gates

$$\mathbf{f}_t = \sigma \left(\mathbf{W}^{(xf)} \mathbf{x}_t + \mathbf{W}^{(hf)} \mathbf{h}_{t-1} + \mathbf{W}^{(cf)} \mathbf{c}_{t-1} + \mathbf{b}^{(f)} \right)$$

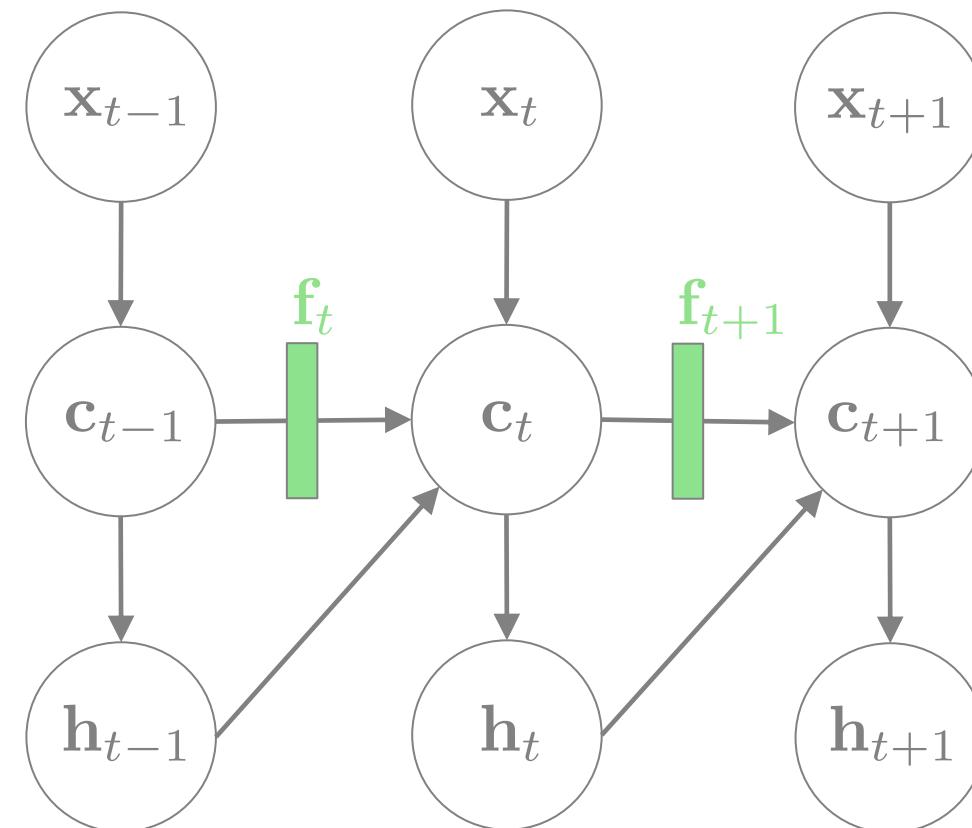
forget gate depends on
current observation,
previous hidden vector,
and previous cell vector



Adding Forget Gates

$$\mathbf{f}_t = \sigma \left(\mathbf{W}^{(xf)} \mathbf{x}_t + \mathbf{W}^{(hf)} \mathbf{h}_{t-1} + \mathbf{W}^{(cf)} \mathbf{c}_{t-1} + \mathbf{b}^{(f)} \right)$$

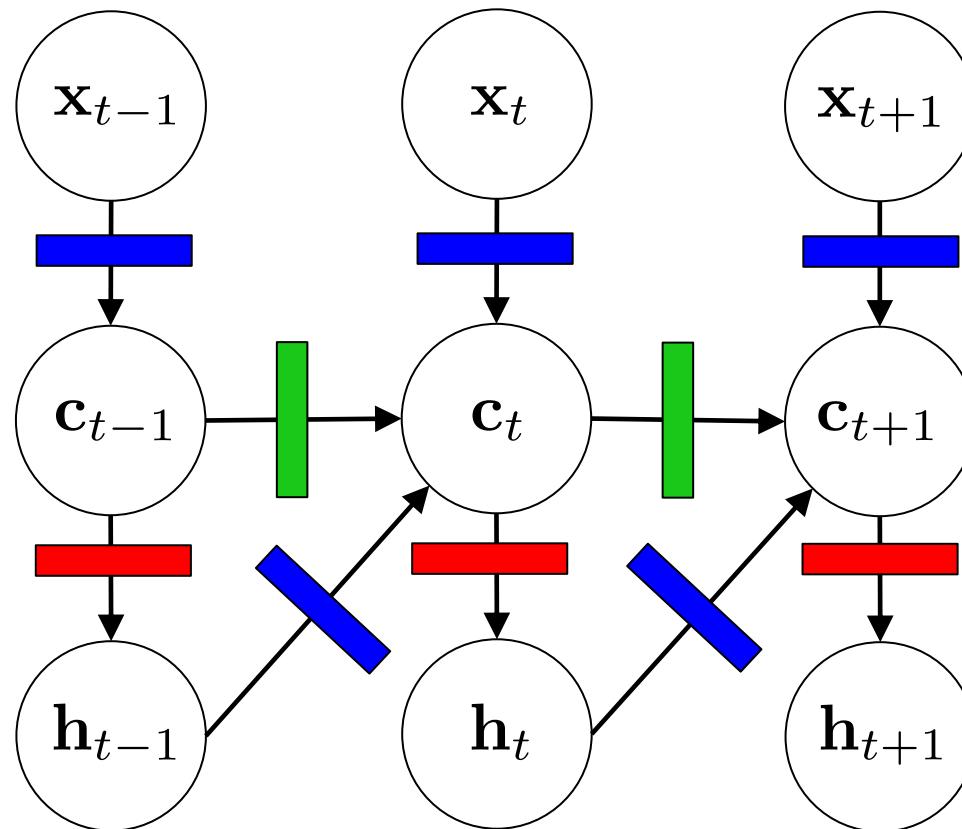
	acc.
gateless	80.6
output gates	81.9
input gates	84.4
forget gates	82.1



All Gates

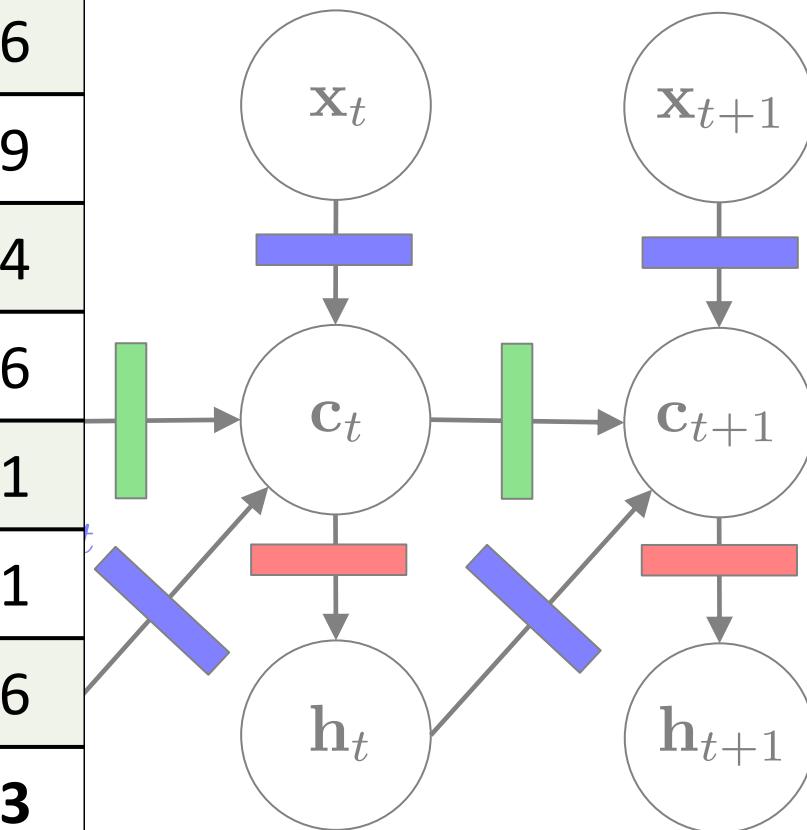
$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tanh \left(\mathbf{W}^{(xc)} \mathbf{x}_t + \mathbf{W}^{(hc)} \mathbf{h}_{t-1} + \mathbf{b}^{(c)} \right)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)$$



All Gates

	acc.
gateless	80.6
output gates	81.9
input gates	84.4
input & output gates	84.6
forget gates	82.1
input & forget gates	84.1
forget & output gates	82.6
input, forget, output gates	85.3



Gated Recurrent Units (GRU)

- ❑ Proposed by Cho et al. in 2014 as a simpler alternative to the LSTM.
- ❑ On each timestep t we have input $x^{(t)}$ and hidden state $h^{(t)}$ no cell state)

Gated Recurrent Units (GRU)

Update gate: controls what parts of hidden state are updated vs preserved

Reset gate: controls what parts of previous hidden state are used to compute new content

New hidden state content: reset gate selects useful parts of prev hidden state. Use this and current input to compute new hidden content.

Hidden state: update gate simultaneously controls what is kept from previous hidden state, and what is updated to new hidden state content

$$\begin{aligned} \mathbf{u}^{(t)} &= \sigma \left(\mathbf{W}_u \mathbf{h}^{(t-1)} + \mathbf{U}_u \mathbf{x}^{(t)} + \mathbf{b}_u \right) \\ \mathbf{r}^{(t)} &= \sigma \left(\mathbf{W}_r \mathbf{h}^{(t-1)} + \mathbf{U}_r \mathbf{x}^{(t)} + \mathbf{b}_r \right) \end{aligned}$$

$$\begin{aligned} \tilde{\mathbf{h}}^{(t)} &= \tanh \left(\mathbf{W}_h (\mathbf{r}^{(t)} \circ \mathbf{h}^{(t-1)}) + \mathbf{U}_h \mathbf{x}^{(t)} + \mathbf{b}_h \right) \\ \mathbf{h}^{(t)} &= (1 - \mathbf{u}^{(t)}) \circ \mathbf{h}^{(t-1)} + \mathbf{u}^{(t)} \circ \tilde{\mathbf{h}}^{(t)} \end{aligned}$$

How does this solve vanishing gradient?
Like LSTM, GRU makes it easier to retain info long-term (e.g., by setting update gate to 0)

LSTM vs GRU

- ❑ Researchers have proposed many gated RNN variants, but LSTM and GRU are the most widely-used
- ❑ **Rule of thumb:** LSTM is a **good default choice** (especially if your data has particularly long dependencies, or you have lots of training data); Switch to **GRUs** for **speed** and **fewer parameters**.

How does LSTM solve vanishing gradients?

- The LSTM architecture makes it **easier** for the RNN to **preserve information over many timesteps**
 - e.g., if the forget gate is set to 1 for a cell dimension and the input gate set to 0, then the information of that cell is preserved indefinitely.
 - In contrast, it's harder for a vanilla RNN to learn a recurrent weight matrix W_h that preserves info in the hidden state
- LSTM doesn't *guarantee* that there is no vanishing/exploding gradient, but it does provide an easier way for the model to learn long-distance dependencies

LSTMs: real-world success

- In 2013–2015, LSTMs started achieving state-of-the-art results
 - Successful tasks include handwriting recognition, speech recognition, machine translation
 - LSTMs became the **dominant approach** for most NLP tasks
- Now (2019–2022), other approaches (e.g., **Transformers**) have become dominant for many tasks
 - For example, in **WMT** (a Machine Translation conference + competition):
 - In WMT 2014, there were 0 neural machine translation systems (!)
 - In WMT 2016, the summary report contains “**RNN**” 44 times (and these systems won)
 - In WMT 2019: “**RNN**” 7 times, “**Transformer**” 105 times

Is vanishing/exploding gradient just a RNN problem?

- ❑ No! It can be a problem for all neural architectures (including **feed-forward** and **convolutional**), especially **very deep** ones.
 - Due to chain rule/choice of nonlinearity function, gradient can become vanishingly small as it backpropagates
 - Thus, lower layers are learned very slowly (hard to train)
- ❑ Solution: lots of new deep feedforward/convolutional architectures **add more direct connections** (thus allowing the gradient to flow)
 - For example:
 - **Residual connections** aka “ResNet”
 - Also known as **skip-connections**
 - The **identity connection preserves information** by default
 - This makes **deep networks much easier to train**

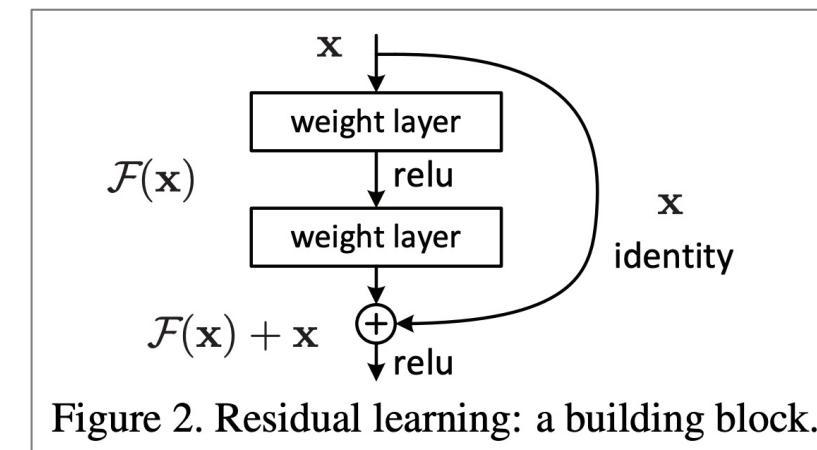


Figure 2. Residual learning: a building block.

Is vanishing/exploding gradient just a RNN problem?

❑ Other methods:

- Dense connections aka “DenseNet”
 - Directly connect each layer to all future layers!
- Highway connections aka “HighwayNet”
 - Similar to residual connections, but the identity connection vs the transformation layer is controlled by a dynamic gate
 - Inspired by LSTMs, but applied to deep feedforward/convolutional networks

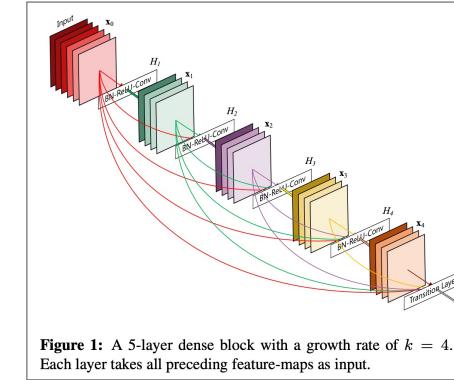
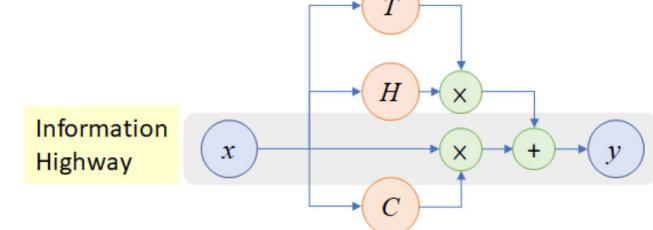


Figure 1: A 5-layer dense block with a growth rate of $k = 4$. Each layer takes all preceding feature-maps as input.



- ## ❑ Conclusion:
- Though vanishing/exploding gradients are a general problem, RNNs are particularly unstable due to the repeated multiplication by the same weight matrix [Bengio et al, 1994]