

Natural Language Processing with Deep Learning

IFT6289, Winter 2022

Lecture 3: Language Modeling and RNN

Bang Liu

2 Lecture outline

1. Language Modeling
2. Neural Language Modeling
3. Recurrent Neural Networks (RNN)
4. RNN Variants

3 Certain Slides Adapted From or Referred To...

- **Stanford CS224n - Natural Language Processing with Deep Learning, Chris Manning**
 - Winter 2020: <http://web.stanford.edu/class/cs224n/>, lecture 6, 7
- https://www.simplilearn.com/tutorials/deep-learning-tutorial/rnn?source=sl_frs_nav_playlist_video_clicked
- <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Language Modeling



5 You Use Language Models Everyday



Q Covid|



Q covid **19**

Q covid **19 bc**

Q covid **symptoms**

Q covid **bc**

Q covid

Q covid **vaccine**

Q covid **19 quebec**

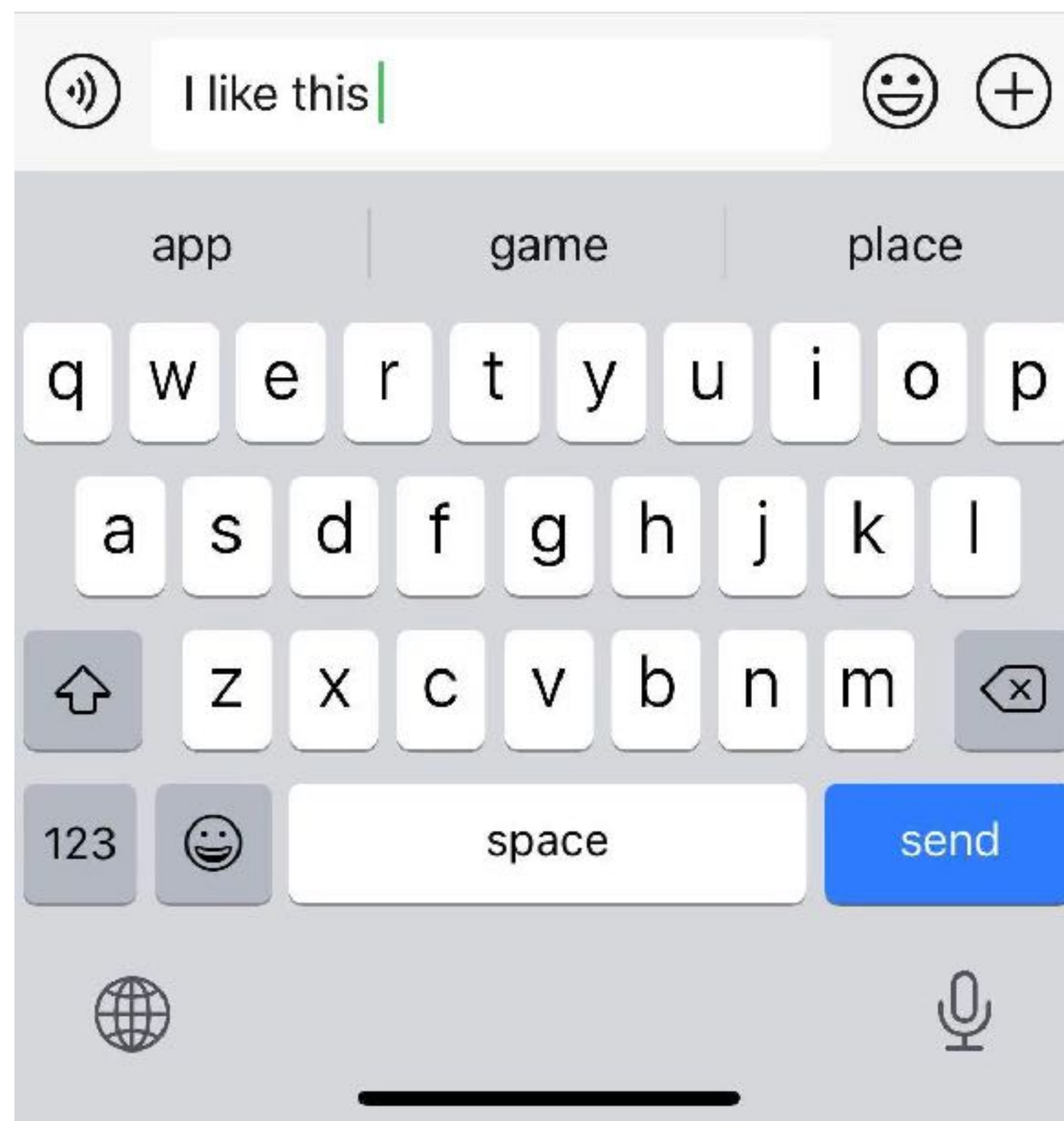
Q covid **cases in bc**

Q covid **19 symptoms**

Q covid **19 bc update**

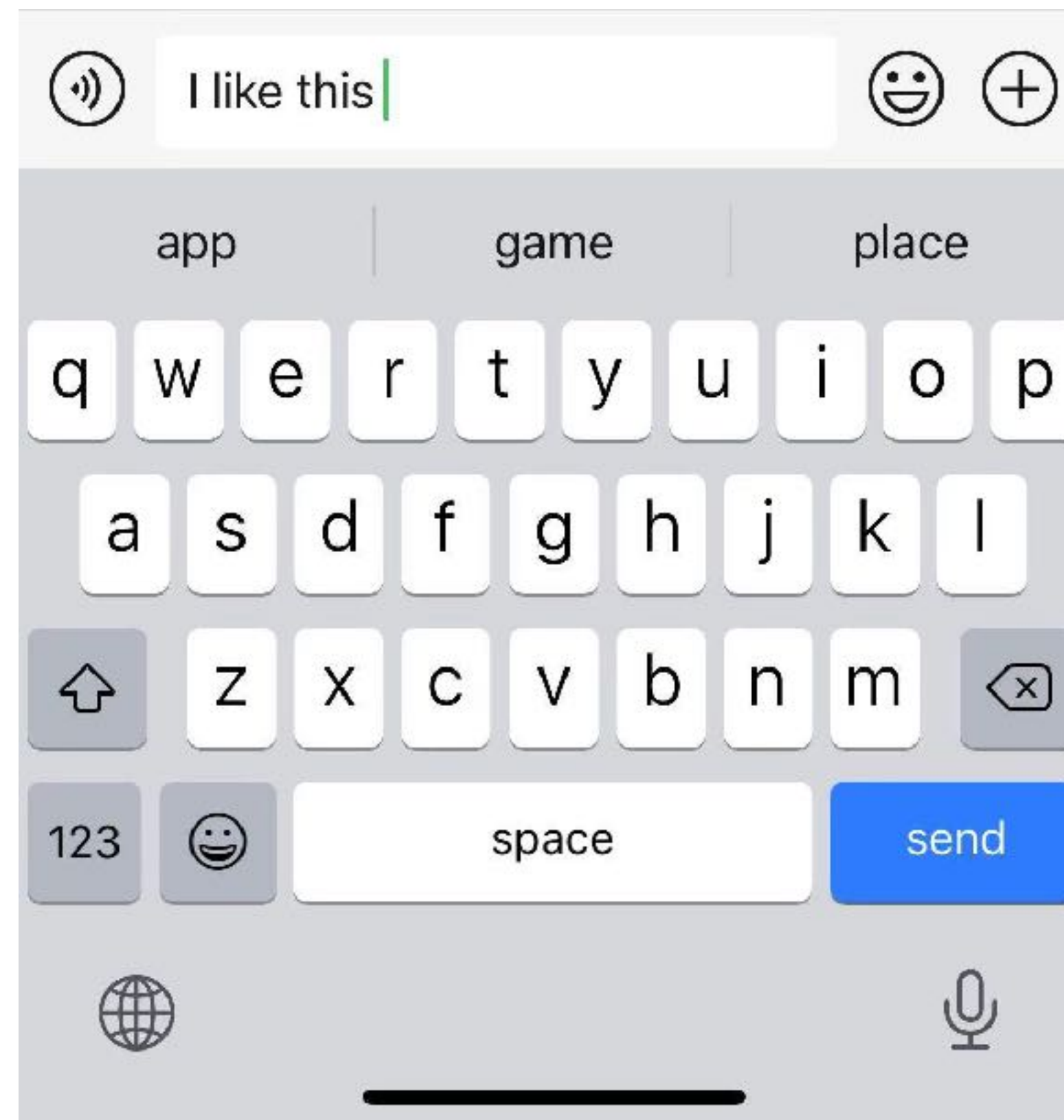
Report inappropriate predictions

6 You Use Language Models Everyday



7 What is Language Modeling

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



8 What is Language Modeling

- **Language Model (LM)** is a system that estimates the probability of a piece of text (a word sequence)

Given a sequence of words:

$$x^{(1)}, x^{(2)}, \dots, x^{(T)}$$

where $x^{(i)}$ can be any word in a vocabulary

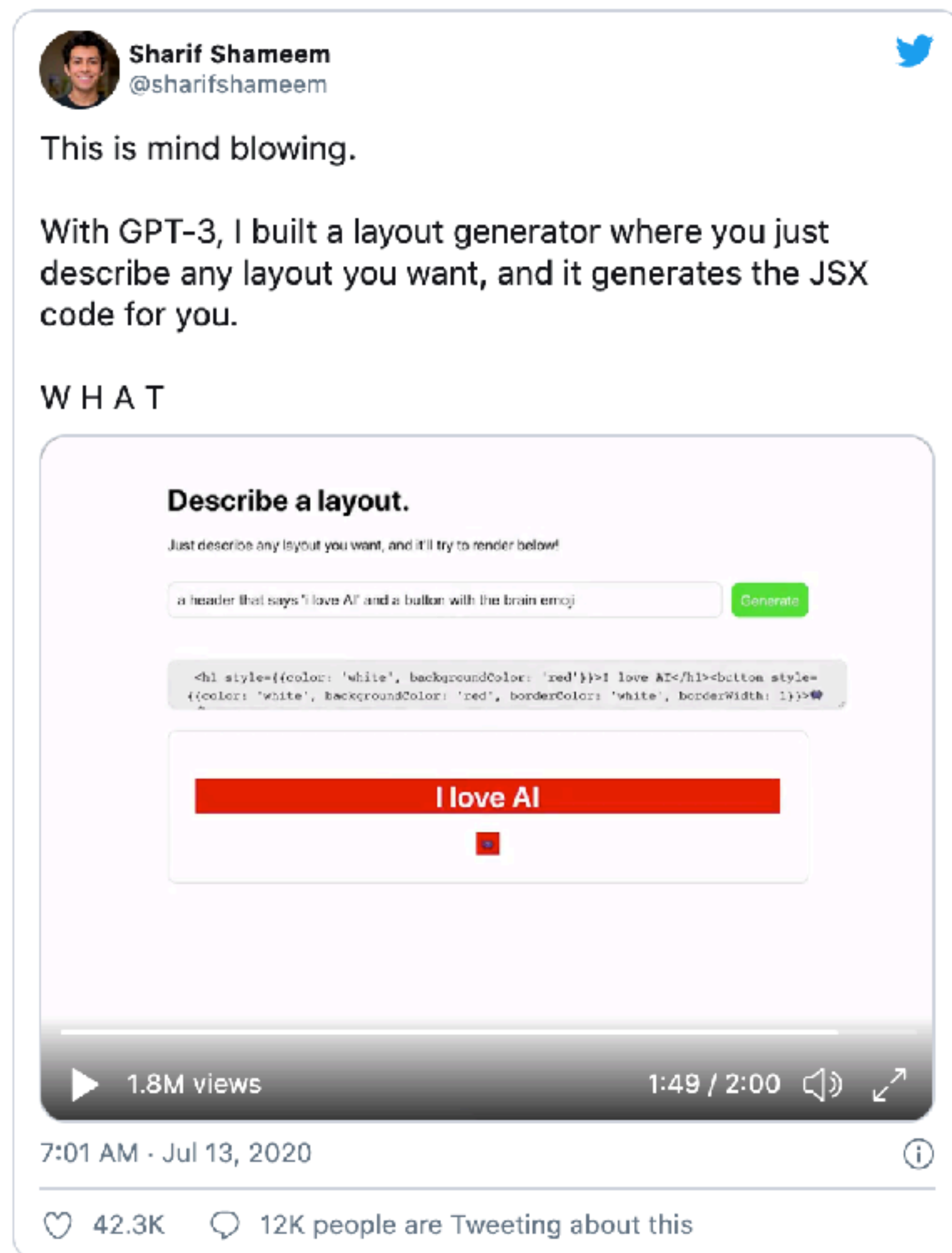
$$V = \{w_1, \dots, w_{|V|}\}$$

Estimate the probability:

$$\begin{aligned} P(x^{(1)}, x^{(2)}, \dots, x^{(T)}) &= P(x^{(1)}) \times P(x^{(2)} | x^{(1)}) \times \dots \times P(x^{(T)} | x^{(T-1)}, \dots, x^{(1)}) \\ &= \prod_{t=1}^T P(x^{(t)} | x^{(t-1)}, \dots, x^{(1)}) \end{aligned}$$

This is what a LM provides

9 Language Models are Changing Everything



Sharif Shameem
@sharifshameem

This is mind blowing.

With GPT-3, I built a layout generator where you just describe any layout you want, and it generates the JSX code for you.

W H A T

Describe a layout.
Just describe any layout you want, and it'll try to render below!

a header that says "I love AI" and a button with the brain emoji

```
<h1 style={{color: 'white', backgroundColor: 'red'}}>I love AI</h1><button style={{color: 'white', backgroundColor: 'red', borderColors: 'white', borderWidth: 1}}>🧠
```

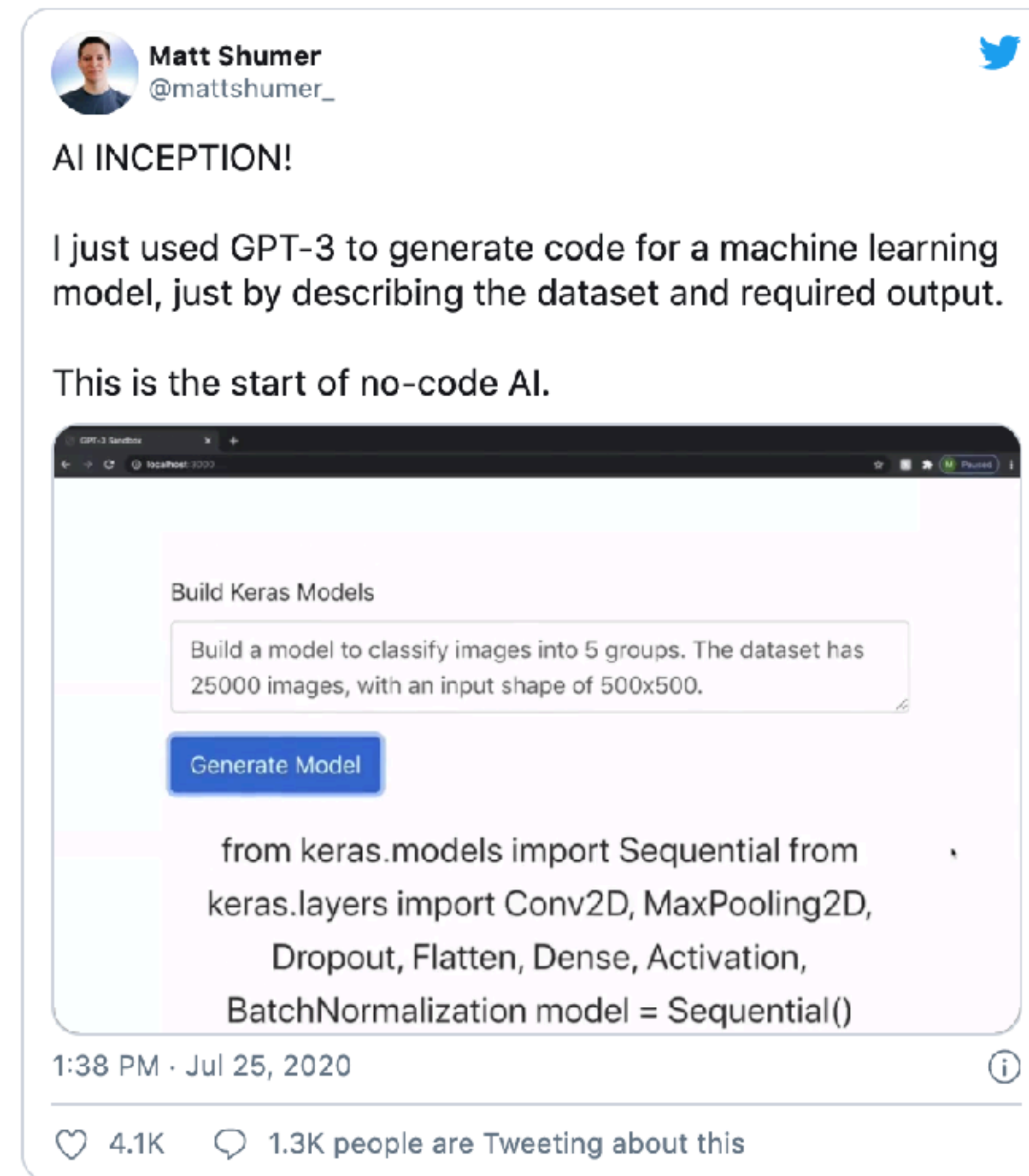
I love AI

1.8M views 1:49 / 2:00

7:01 AM · Jul 13, 2020

42.3K 12K people are Tweeting about this

Code



Matt Shumer
@mattshumer_

AI INCEPTION!

I just used GPT-3 to generate code for a machine learning model, just by describing the dataset and required output.

This is the start of no-code AI.

Build Keras Models

Build a model to classify images into 5 groups. The dataset has 25000 images, with an input shape of 500x500.

Generate Model

```
from keras.models import Sequential  
from keras.layers import Conv2D, MaxPooling2D,  
Dropout, Flatten, Dense, Activation,  
BatchNormalization  
model = Sequential()
```

1:38 PM · Jul 25, 2020

4.1K 1.3K people are Tweeting about this

Machine Learning Code

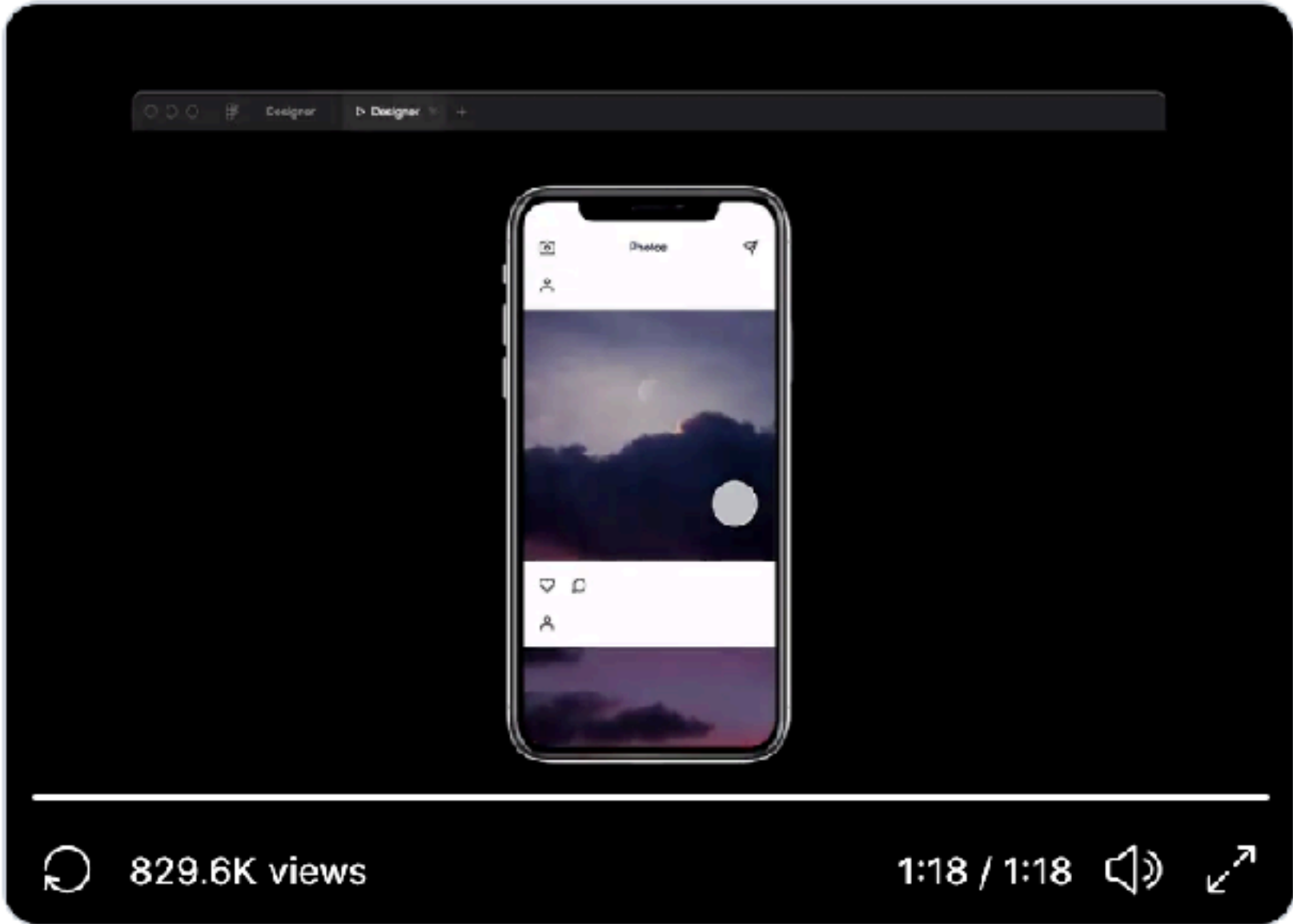
10 Language Models are Changing Everything

Jordan Singer
@jsngr

This changes everything. 🤖

With GPT-3, I built a Figma plugin to design for you.

I call it "Designer"



829.6K views 1:18 / 1:18

8:31 AM · Jul 18, 2020

13.9K 3.6K people are Tweeting about this

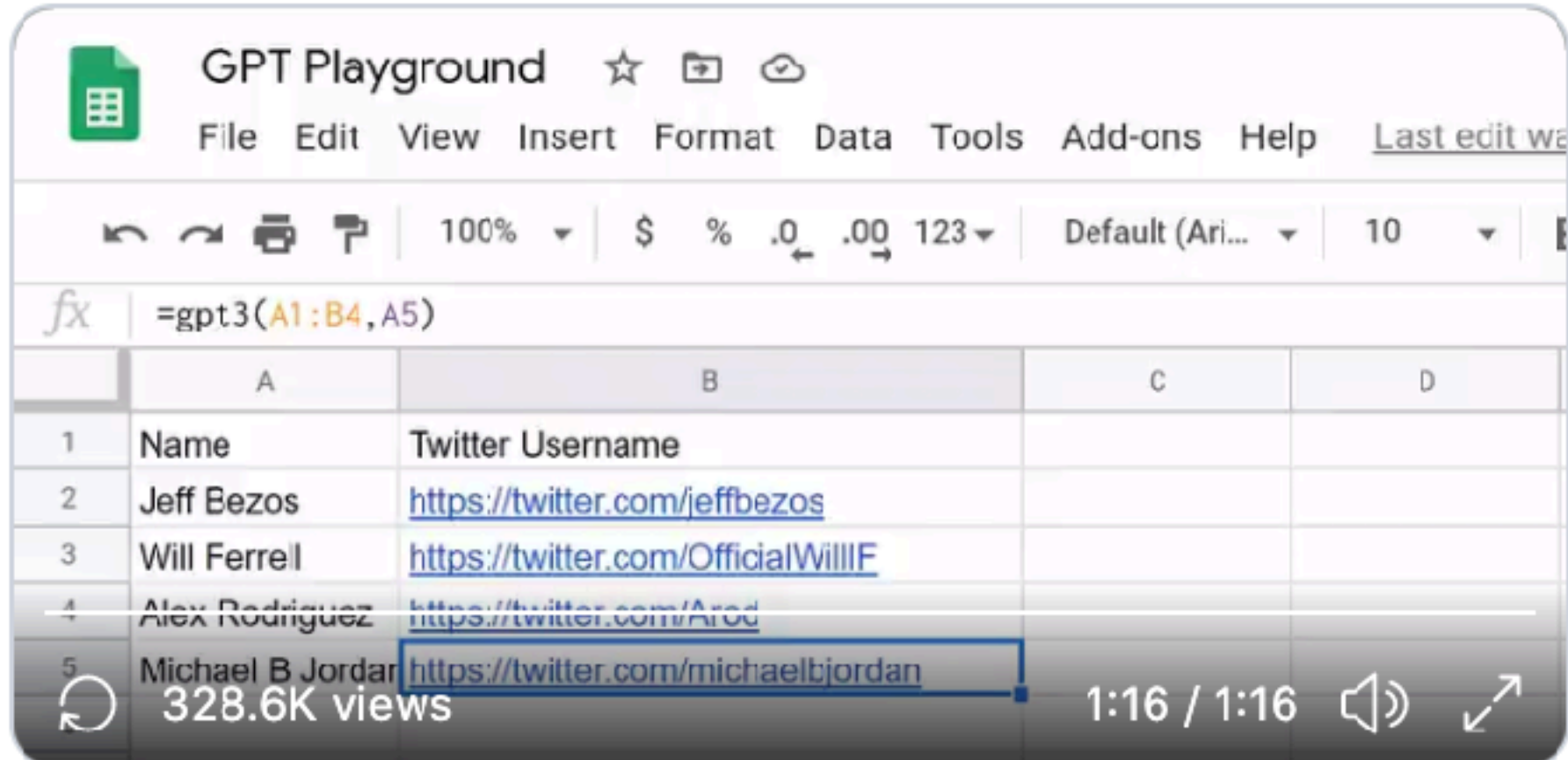
Design

Paul Katsen
@pavtalk

=GPT3()... the spreadsheet function to rule them all.

Impressed with how well it pattern matches from a few examples.

The same function looked up state populations, peoples' twitter usernames and employers, and did some math.



	A	B	C	D
1	Name	Twitter Username		
2	Jeff Bezos	https://twitter.com/jeffbezos		
3	Will Ferrel	https://twitter.com/OfficialWillIF		
4	Alex Rodriguez	https://twitter.com/Arod		
5	Michael B Jordan	https://twitter.com/michaeltjordan		

328.6K views 1:16 / 1:16

8:06 PM · Jul 20, 2020

9.4K 2K people are Tweeting about this

Spreadsheets

**How to Learn a
Language Model?
N-gram Language
Models**

12 N-gram Language Models

- ⦿ A **n-gram** is a chunk of n consecutive words.
Given a piece of text: “the students opened their ...”
 - **uni**grams: “the”, “students”, “opened”, “their”
 - **bi**grams: “the students”, “students opened”, “opened their”
 - **tri**grams: “the students opened”, “students opened their”
 - **4**-grams: “the students opened their”
- ⦿ **Idea (before deep learning)**: Collect statistics about how frequent different n-grams are, and use these to predict next word.

13 N-gram Language Models

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

$$\begin{aligned} P(x^{(t+1)} | x^{(t)}, \dots, x^{(1)}) &= P(x^{(t+1)} | \boxed{x^{(t)}, \dots, x^{(t-n+2)}}) && \text{(assumption)} \\ &= \frac{P(x^{(t+1)}, x^{(t)}, \dots, x^{(t-n+2)})}{P(x^{(t)}, \dots, x^{(t-n+2)})} && \text{(conditional probability)} \\ &\approx \frac{\text{count}(x^{(t+1)}, x^{(t)}, \dots, x^{(t-n+2)})}{\text{count}(x^{(t)}, \dots, x^{(t-n+2)})} && \text{(statistical approximation)} \end{aligned}$$

- We get these n-gram and (n-1)-gram probabilities by **counting** them in some large corpus of text!

14 N-gram Language Models: Example

- Suppose we are learning a 4-gram Language Model

~~as the proctor started the clock, the students~~ opened their _____

$$P(w|\text{students opened their}) = \frac{\text{count}(\text{students opened their } w)}{\text{count}(\text{students opened their})}$$

- For example, suppose that in the corpus:
 - “students opened their” occurred 1000 times
 - If “students opened their books” occurred 400 times

$$P(\text{books}|\text{students opened their}) = 0.4$$

- “students opened their exams” occurred 100 times

$$P(\text{exams}|\text{students opened their}) = 0.1$$

Shall we discard the context “proctor”?

15 Sparsity Problem with N-gram Language Models

Sparsity Problem 1

Problem: What if “students opened their w ” never occurred in data? Then w has probability 0!

(Partial) Solution: Add small δ to the count for every $w \in V$. This is called *smoothing*.

$$P(w|\text{students opened their}) = \frac{\text{count}(\text{students opened their } w)}{\text{count}(\text{students opened their})}$$

Sparsity Problem 2

Problem: What if “students opened their” never occurred in data? Then we can’t calculate probability for any w !

(Partial) Solution: Just condition on “opened their” instead. This is called *backoff*.

Note: Increasing n makes sparsity problems worse. Typically we can’t have n bigger than 5.

16 Storage Problem with N-gram Language Models

Storage: Need to store count for all n -grams you saw in the corpus.

$$P(\mathbf{w}|\text{students opened their}) = \frac{\text{count}(\text{students opened their } \mathbf{w})}{\text{count}(\text{students opened their})}$$

Increasing n or increasing corpus increases model size!

17 N-gram Language Models in Practice

You can build a simple trigram Language Model over a 1.7 million word corpus (Reuters) in a few seconds on your laptop*

Business and financial news

today the _____

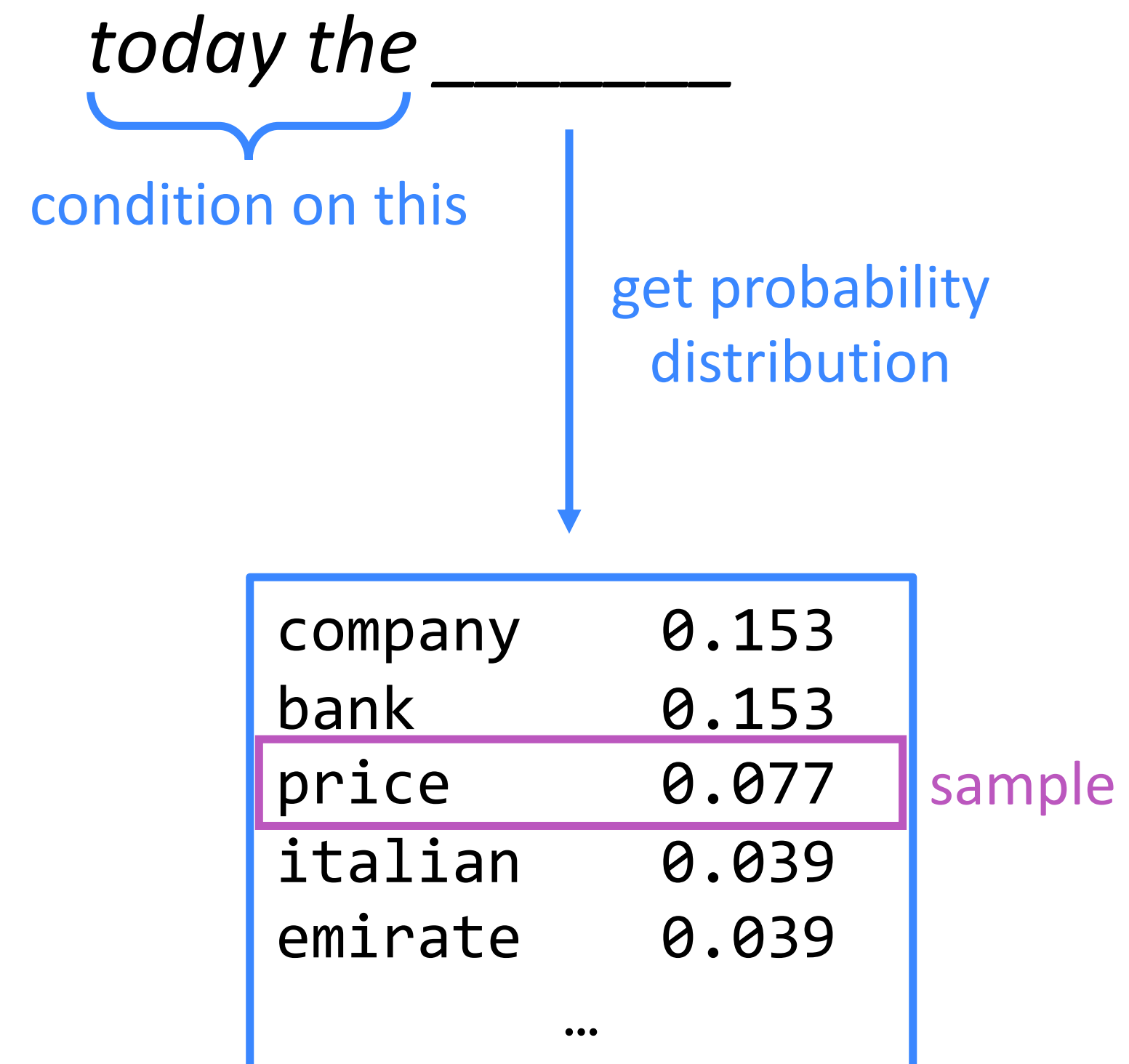
get probability
distribution

company	0.153
bank	0.153
price	0.077
italian	0.039
emirate	0.039
...	

Sparsity problem:
not much granularity
in the probability
distribution

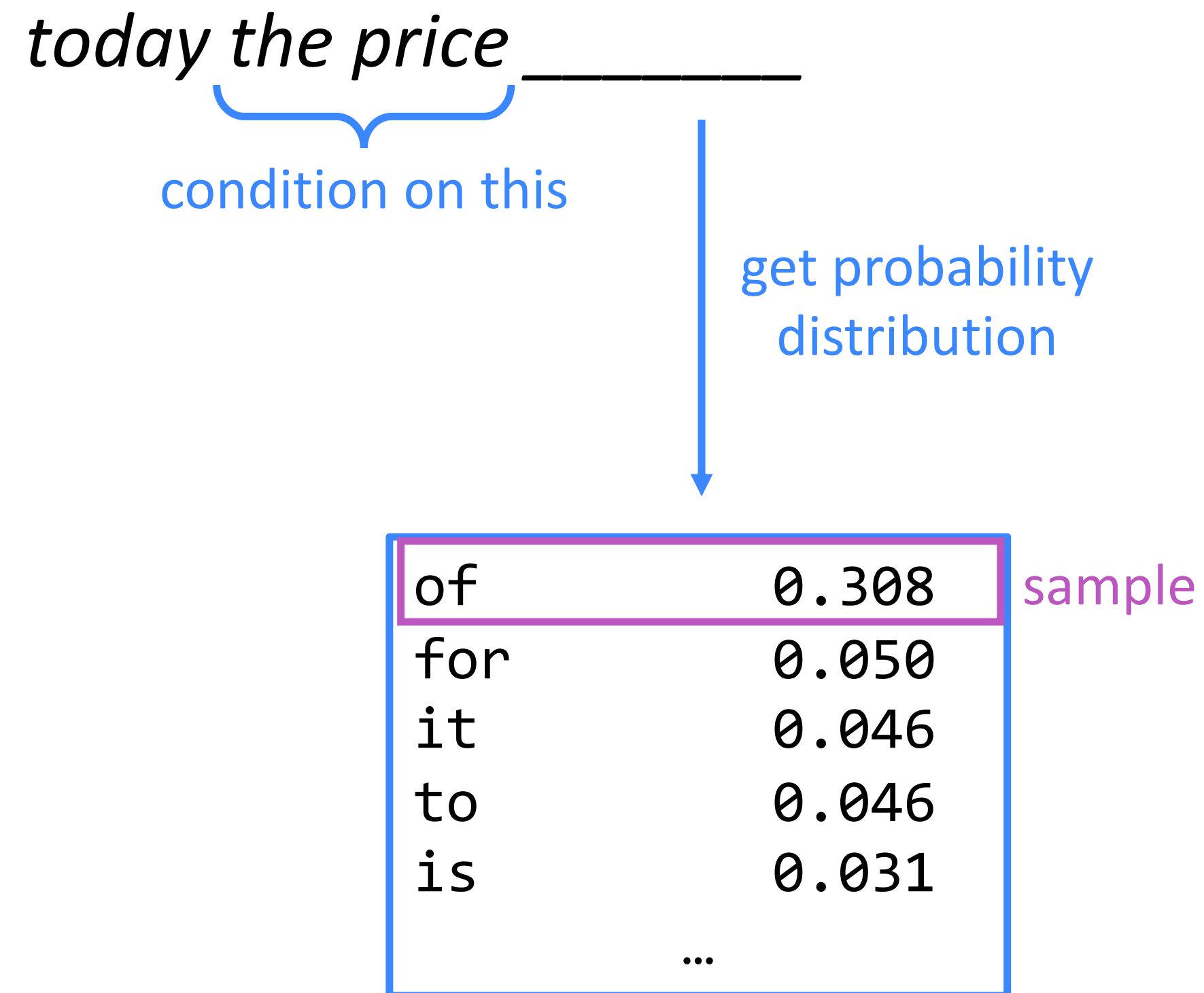
18 Generating Text with a N-gram Language Model

- You can also use a Language Model to generate text.



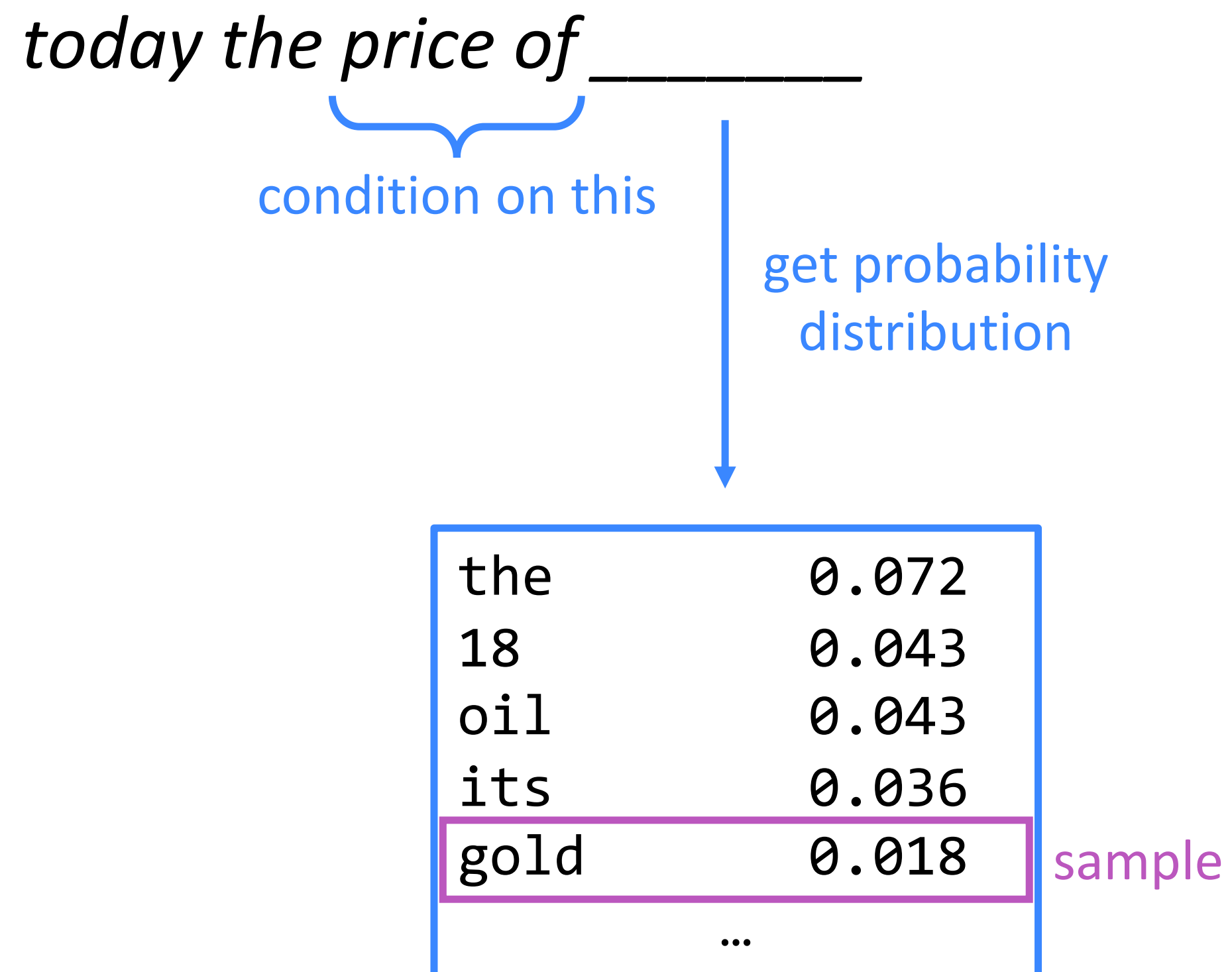
19 Generating Text with a N-gram Language Model

You can also use a Language Model to **generate text**.



20 Generating Text with a N-gram Language Model

You can also use a Language Model to *generate text*.



21 Generating Text with a N-gram Language Model

You can also use a Language Model to **generate text**.

today the price of gold _____

22 Generating Text with a N-gram Language Model

You can also use a Language Model to **generate text**.

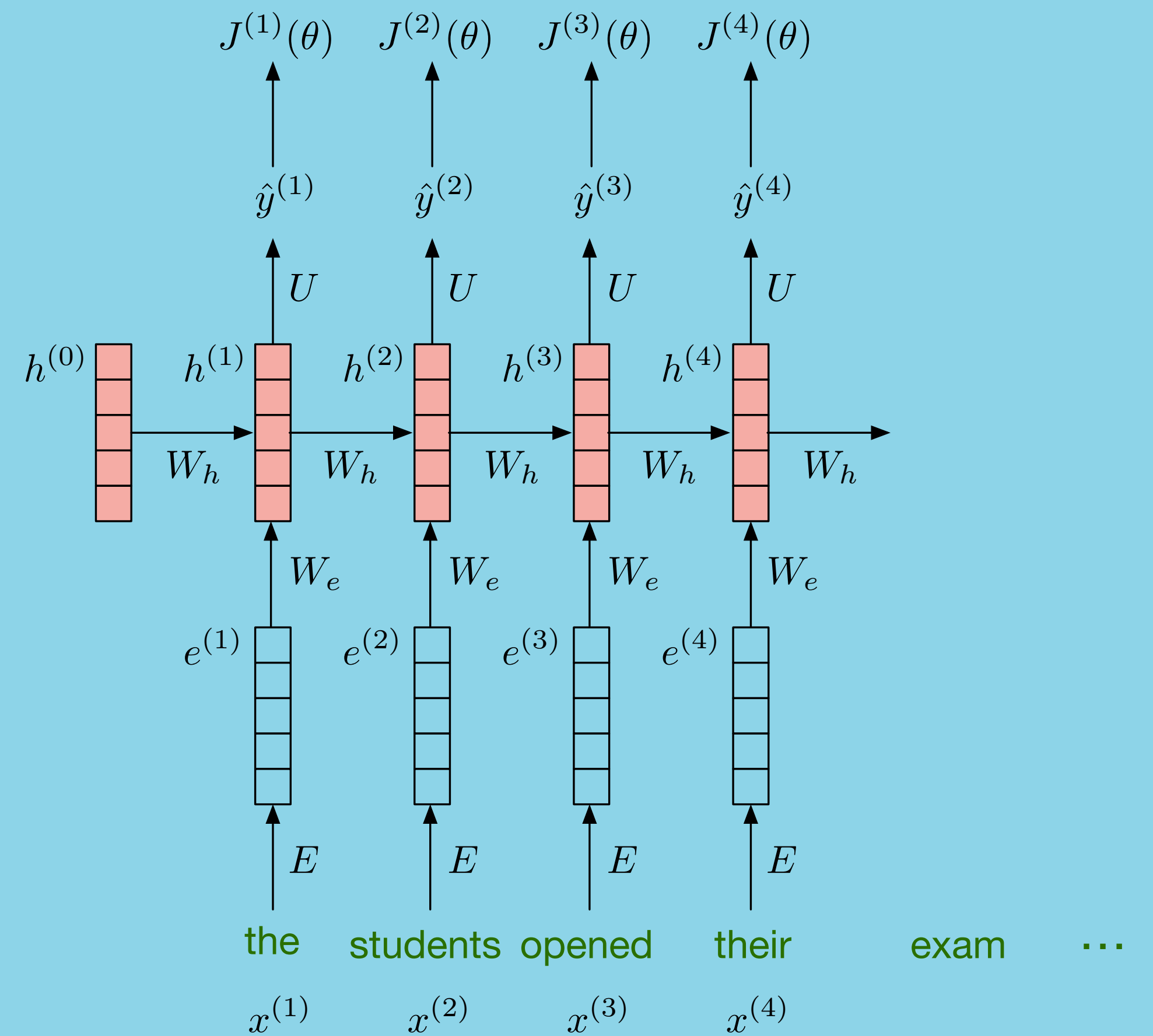
today the price of gold per ton , while production of shoe lasts and shoe industry , the bank intervened just after it considered and rejected an imf demand to rebuild depleted european stocks , sept 30 end primary 76 cts a share .

Surprisingly grammatical!

...but **incoherent**. We need to consider more than three words at a time if we want to model language well.

But increasing n worsens sparsity problem,
and increases model size...

Neural Language Modeling



24 A Fixed-Window Neural Language Model

- Estimate $P(x^{(t+1)} | x^{(t)}, \dots, x^{(t-n+2)})$ not from **count**, but from **neural network prediction**

output distribution

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

hidden layer

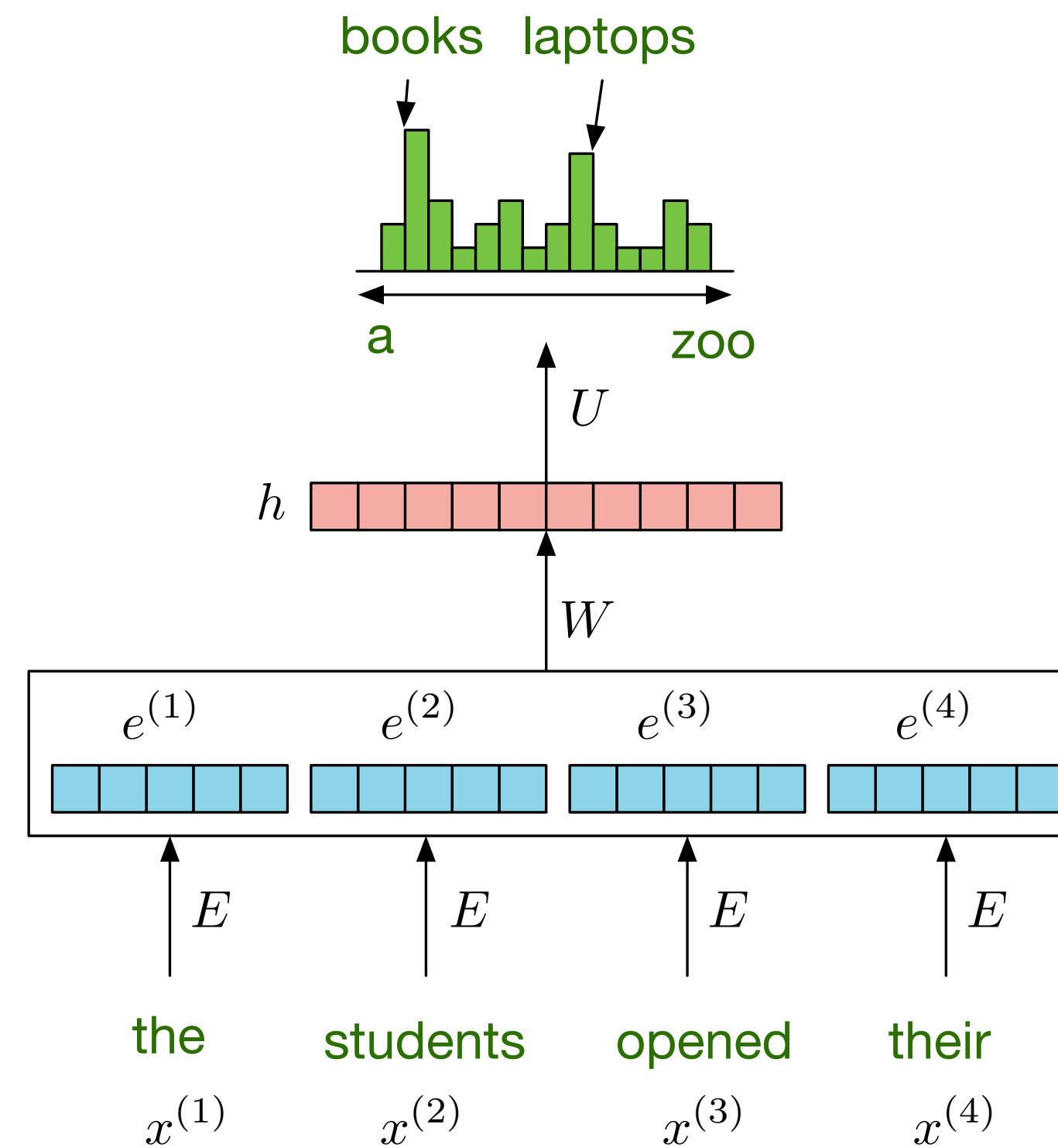
$$\mathbf{h} = f(\mathbf{W}\mathbf{e} + \mathbf{b}_1)$$

concatenated word embeddings

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

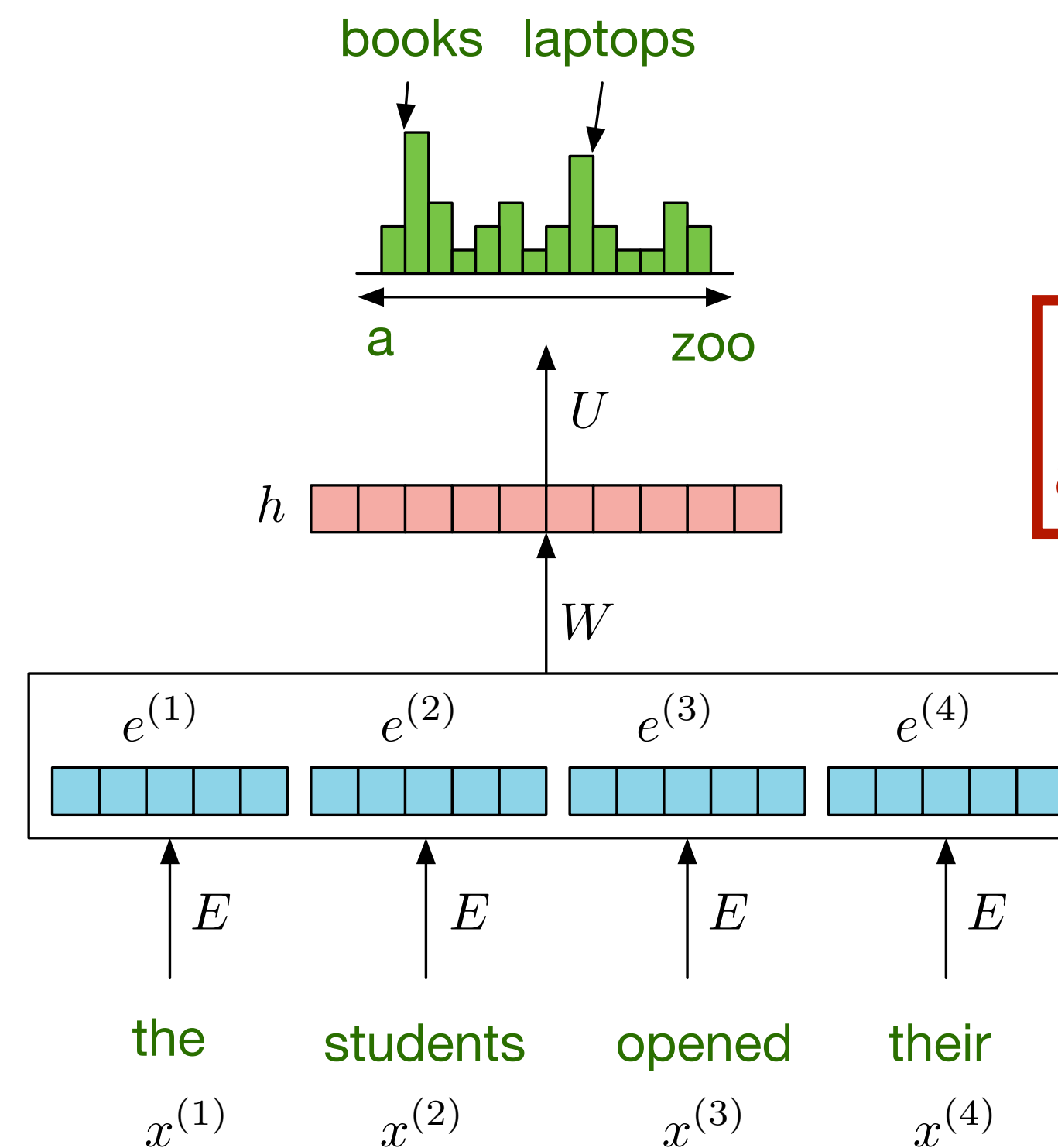
words / one-hot vectors

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



25 A Fixed-Window Neural Language Model

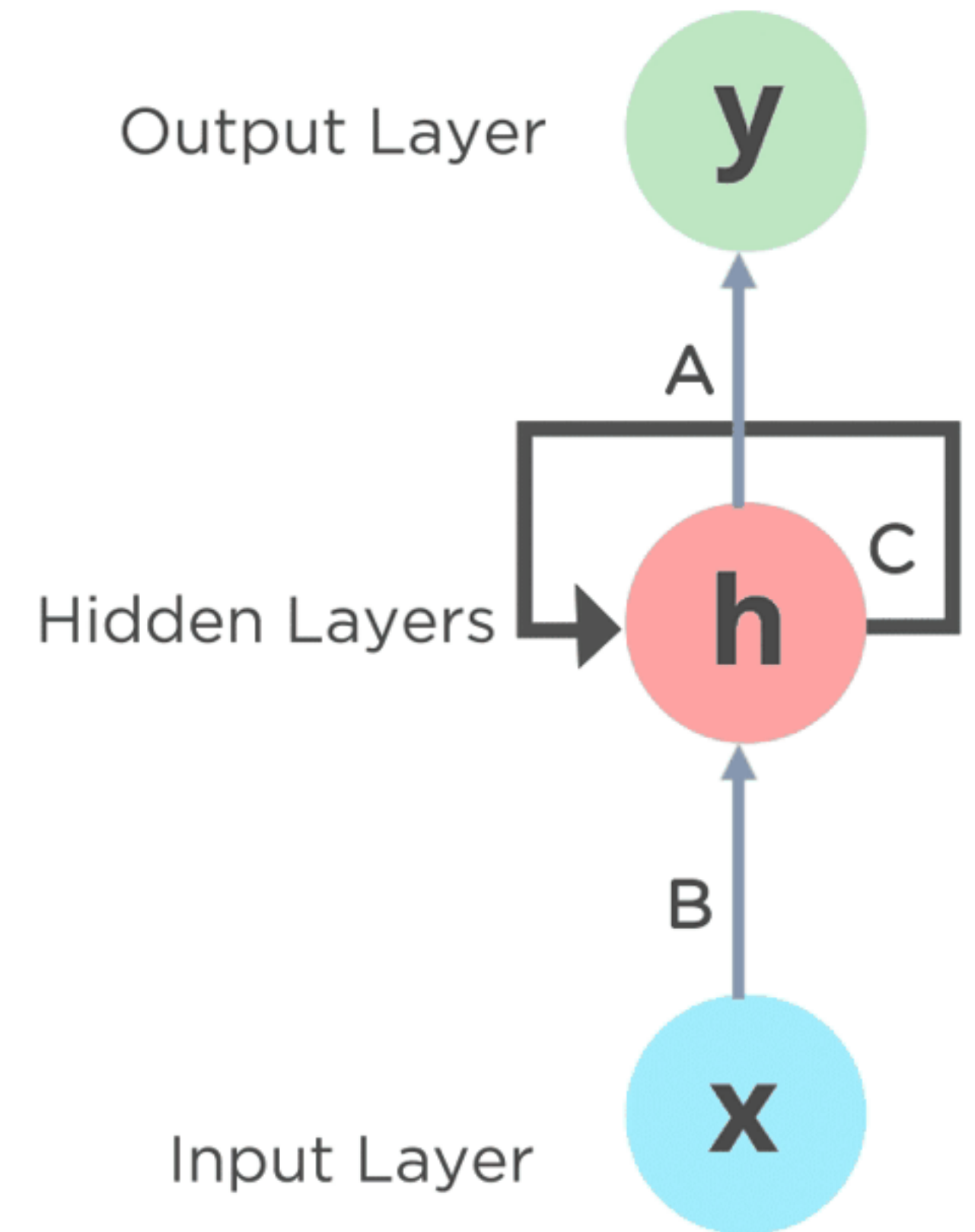
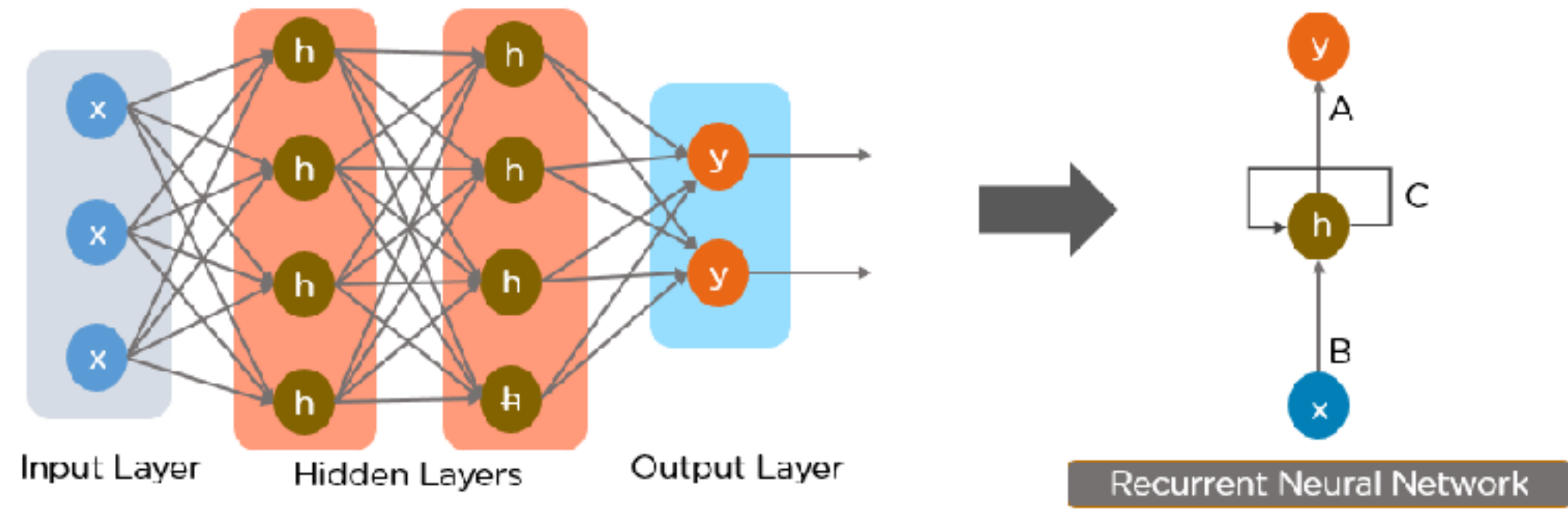
- **Improvements** to n-gram LM
 - No sparsity problem
 - Don't need to store all observed n-grams
- **Remaining problems**
 - Fixed window is too small
 - Enlarging window enlarges W
 - Words in different positions are multiplied by different W , **no symmetry** in how the inputs are processed.



How to process any length input?

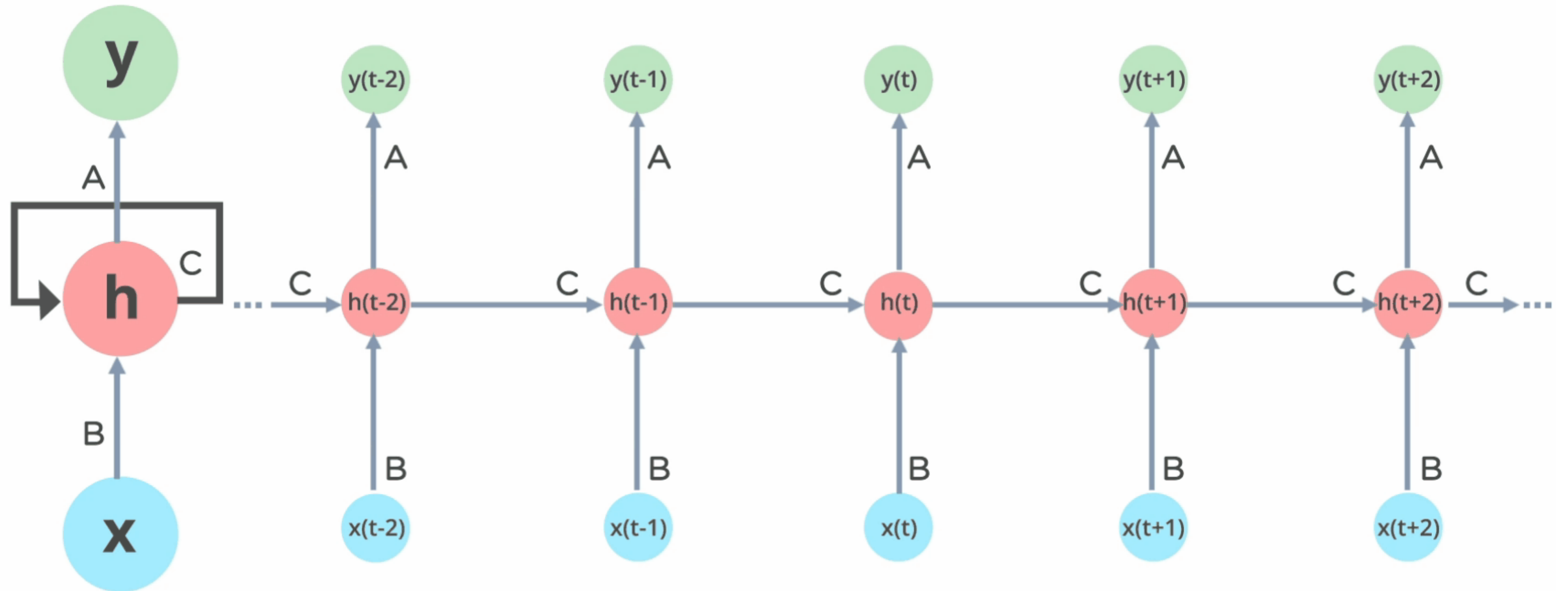
Recurrent Neural Networks

27 Recurrent Neural Networks (RNN)



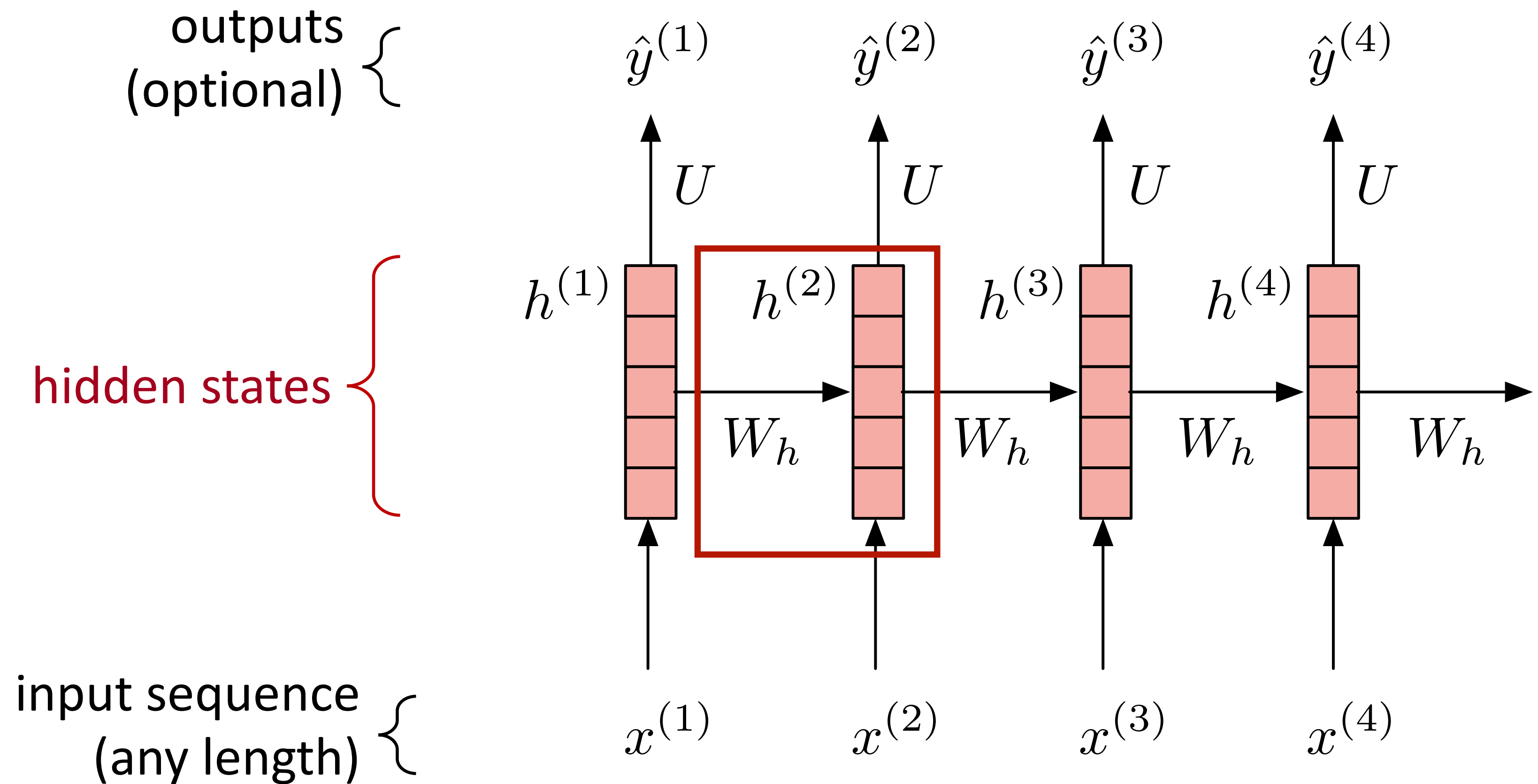
A, B and C are the parameters

28 Recurrent Neural Networks (RNN)



Recurrent Neural Networks (RNN)

- RNN is a family of neural architectures



Core idea: Apply the same weights W repeatedly

A Simple RNN Language Model

output distribution

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

hidden states

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

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

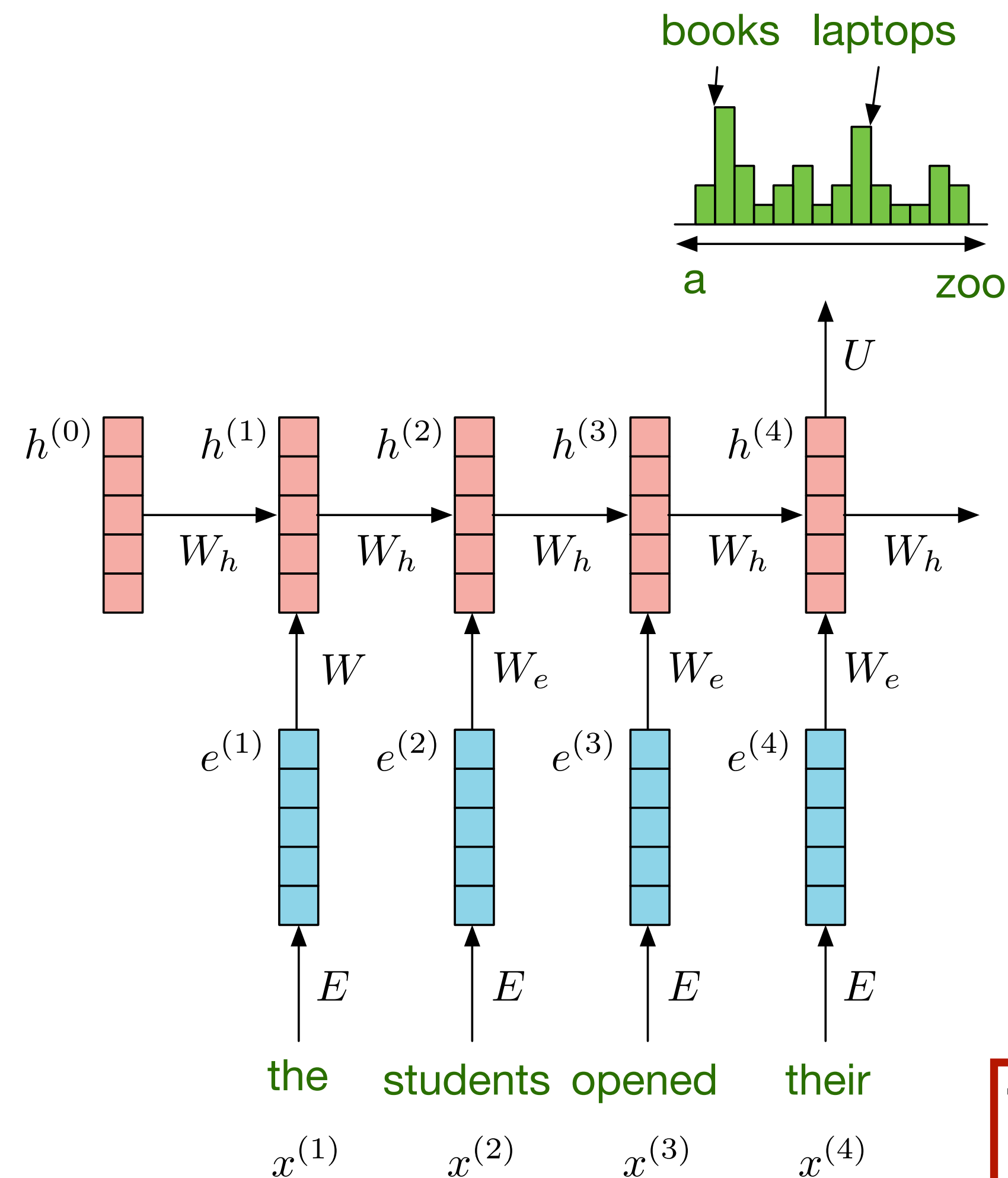
word embeddings

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

words / one-hot vectors

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

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



The input sequence can be any length

RNN Language Models

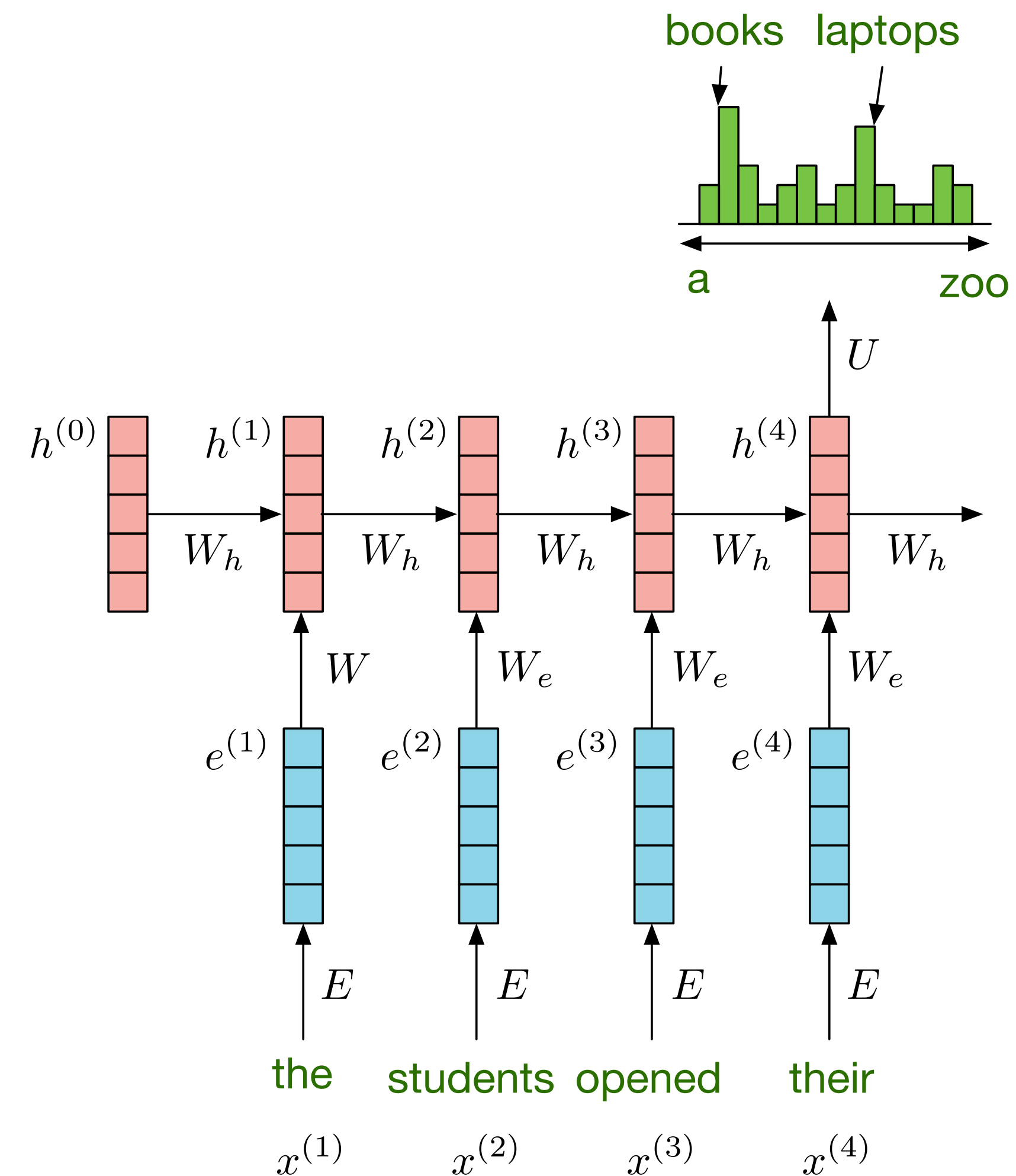
⦿ RNN advantages

- Can process any length input
- Computation in step t can (in theory) use information from many steps back
- Model size doesn't increase for longer inputs
- Same weights applied on each time step, so there is symmetry on how inputs are processed

⦿ RNN problems

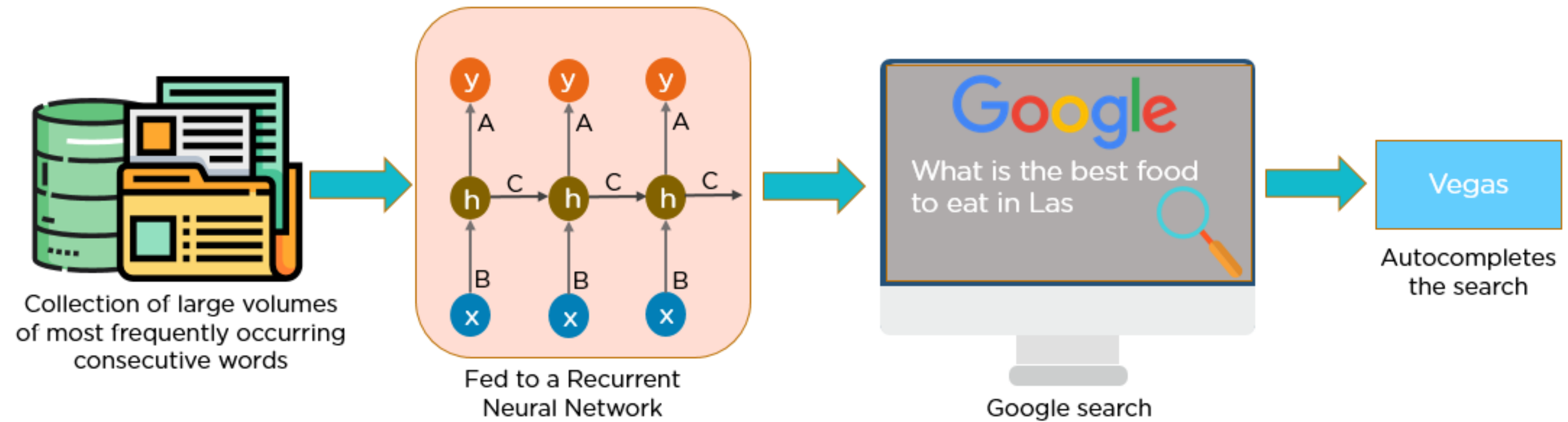
- RNN computations are slow
- In practice, it is difficult to access information from many steps back

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



Training Recurrent Neural Networks

33 RNN Language Models



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 $\hat{\mathbf{y}}^{(t)}$ **for every step t** .
- i.e. predict probability dist of *every word*, given words so far

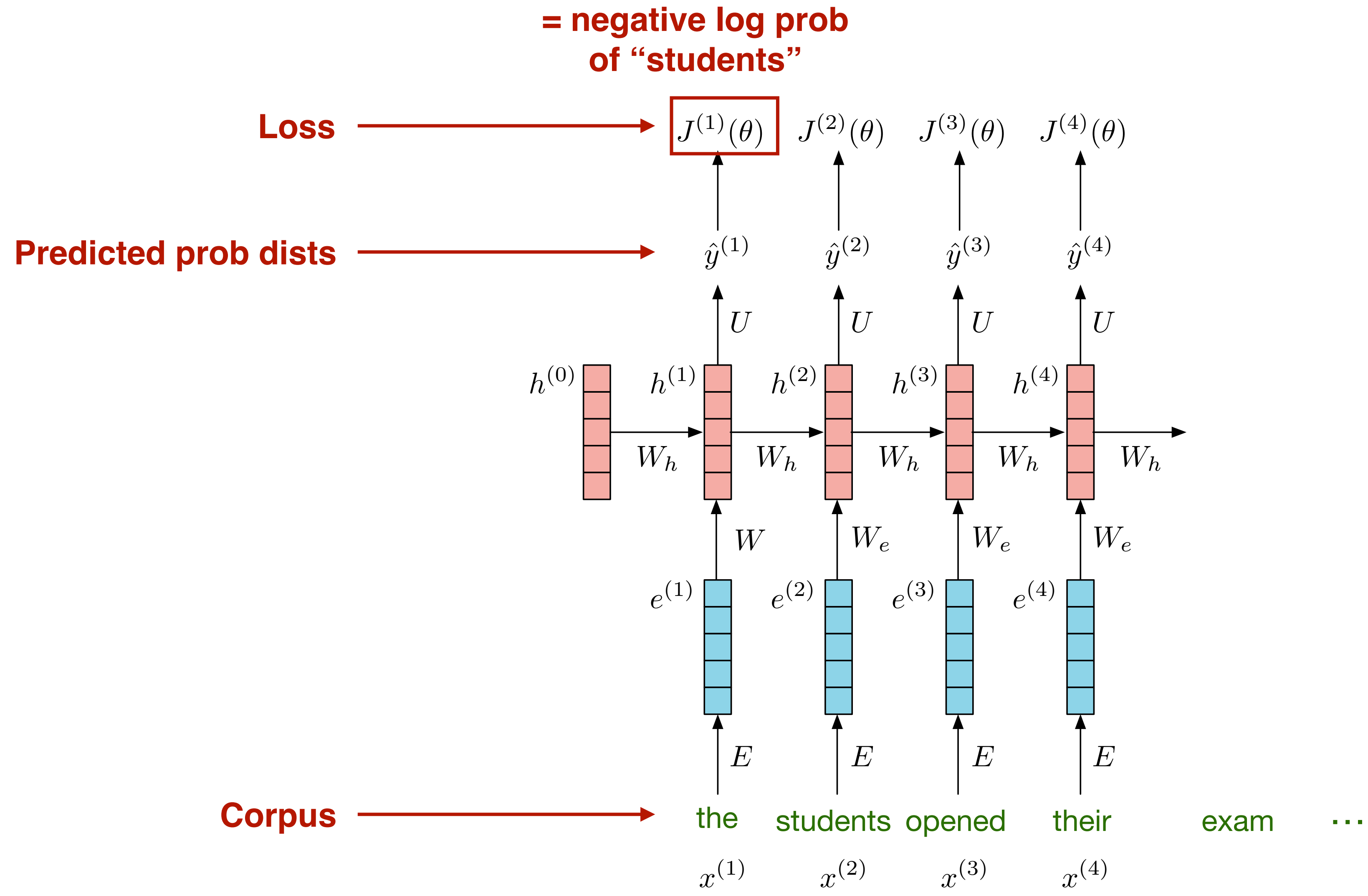
Loss function on step t is **cross-entropy** between predicted probability distribution $\hat{\mathbf{y}}^{(t)}$, and the true next word $\mathbf{y}^{(t)}$ (one-hot for $x^{(t+1)}$):

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

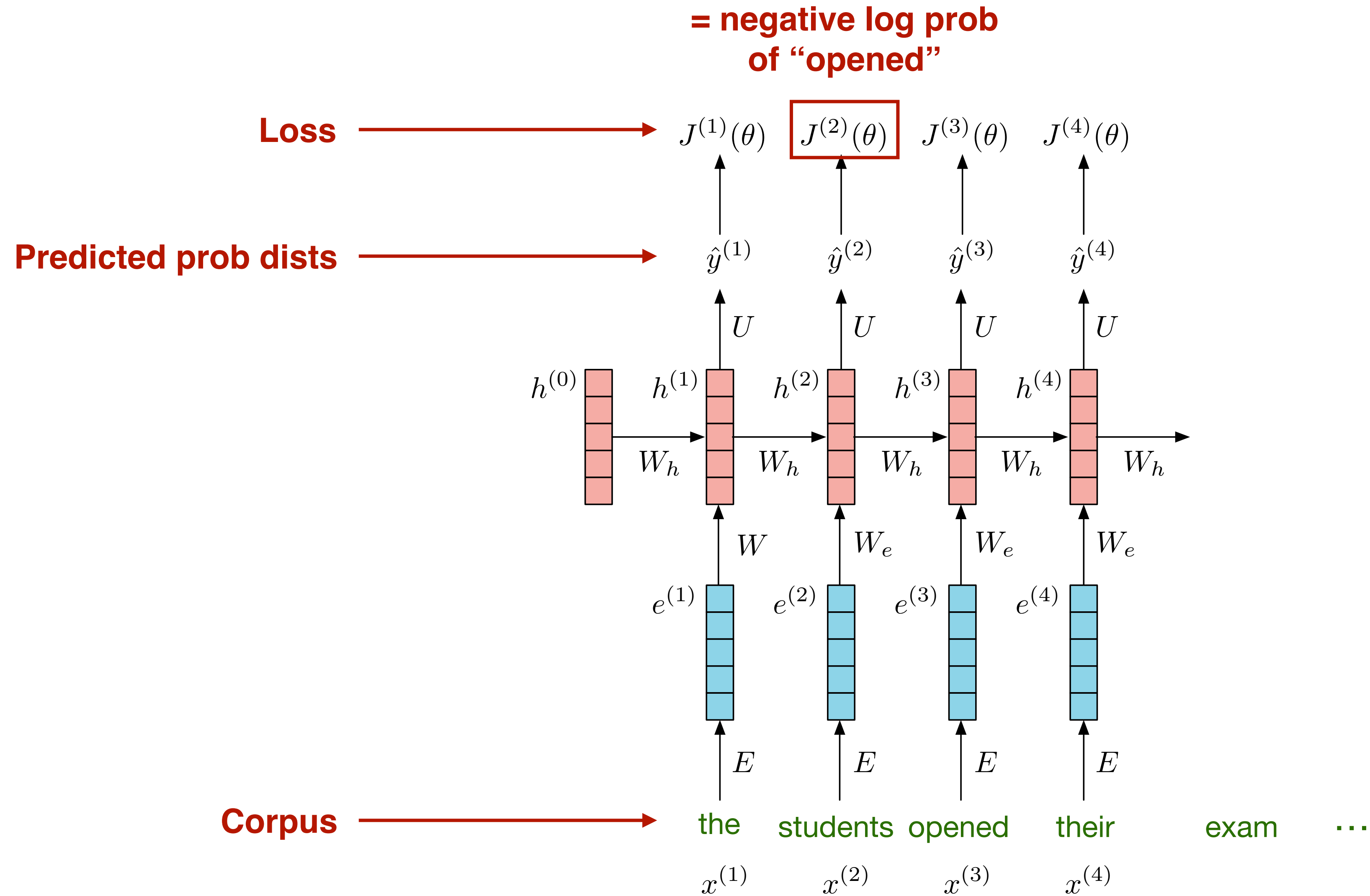
Average this to get **overall loss** for entire training set:

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

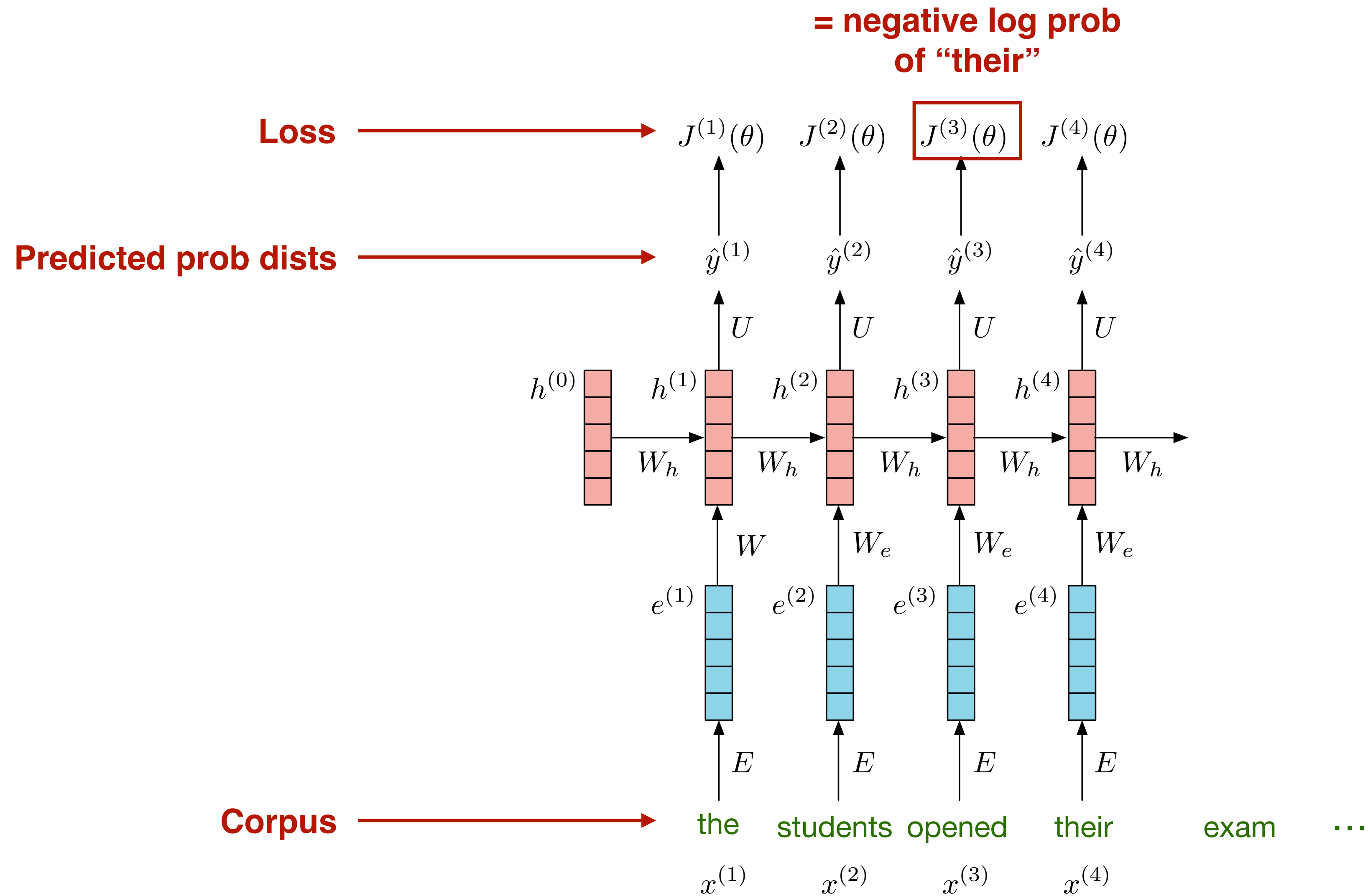
Training an RNN Language Model



Training an RNN Language Model

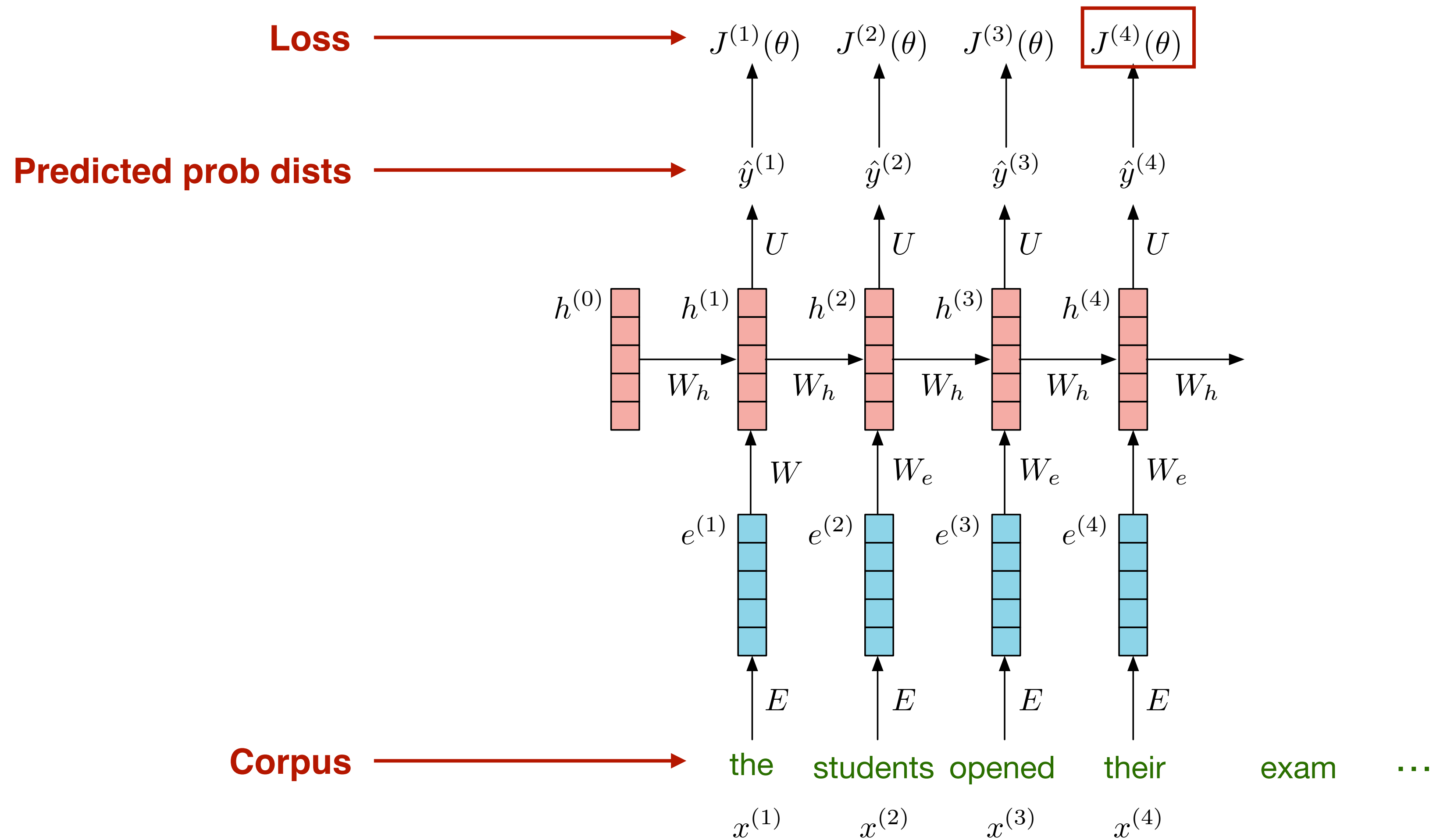


Training an RNN Language Model



Training an RNN Language Model

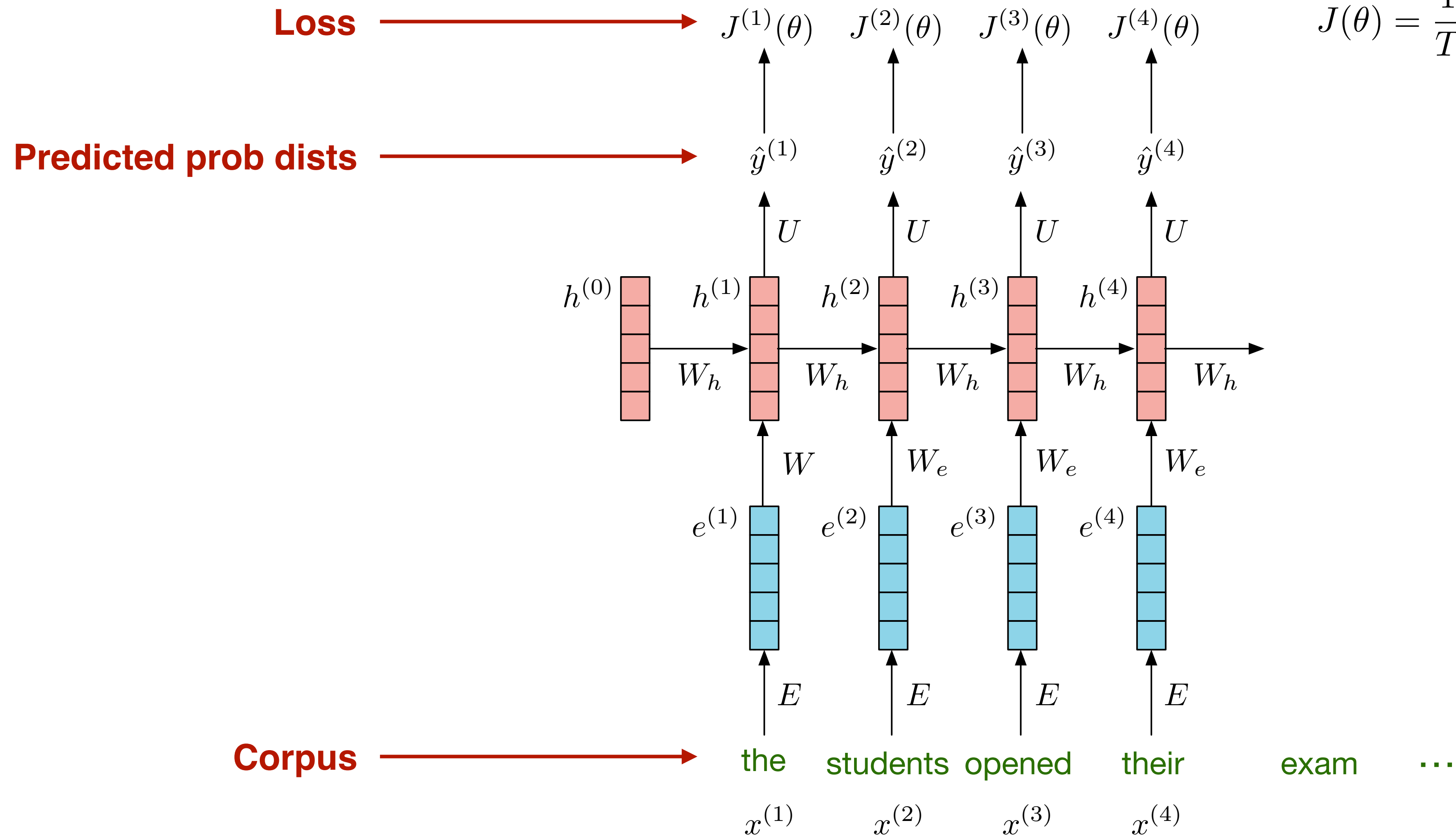
= negative log prob of "exam"



Training an RNN Language Model

“Teacher forcing”

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



40 Training an RNN Language Model

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

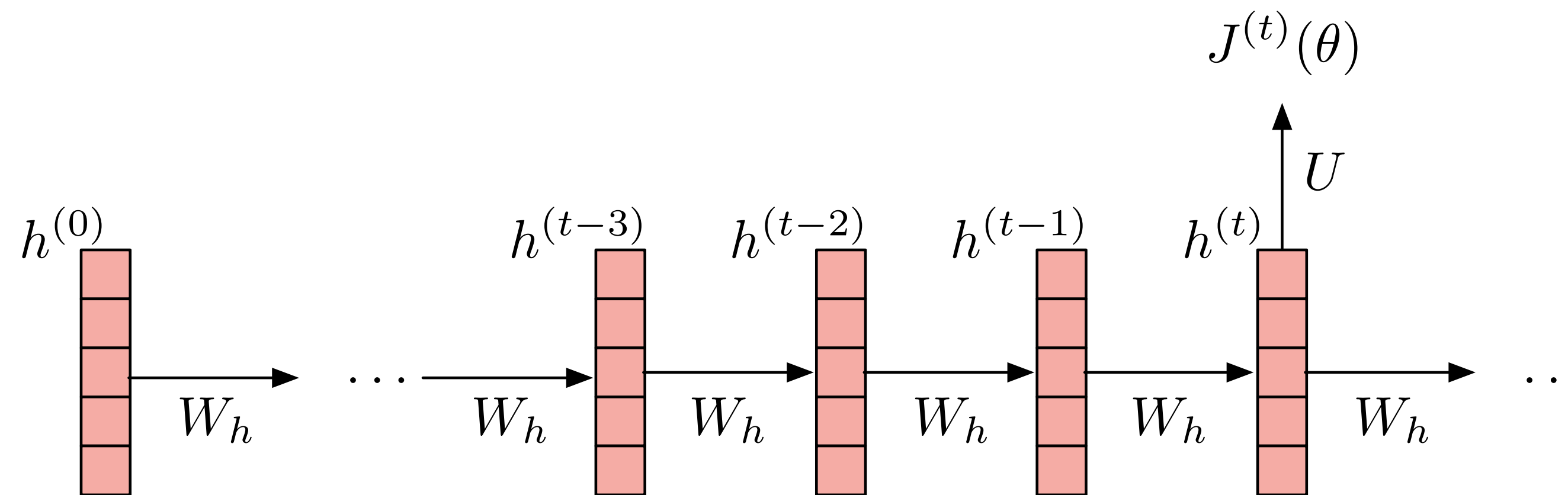
$$J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta)$$

In practice, consider $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}$ 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.

41 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”

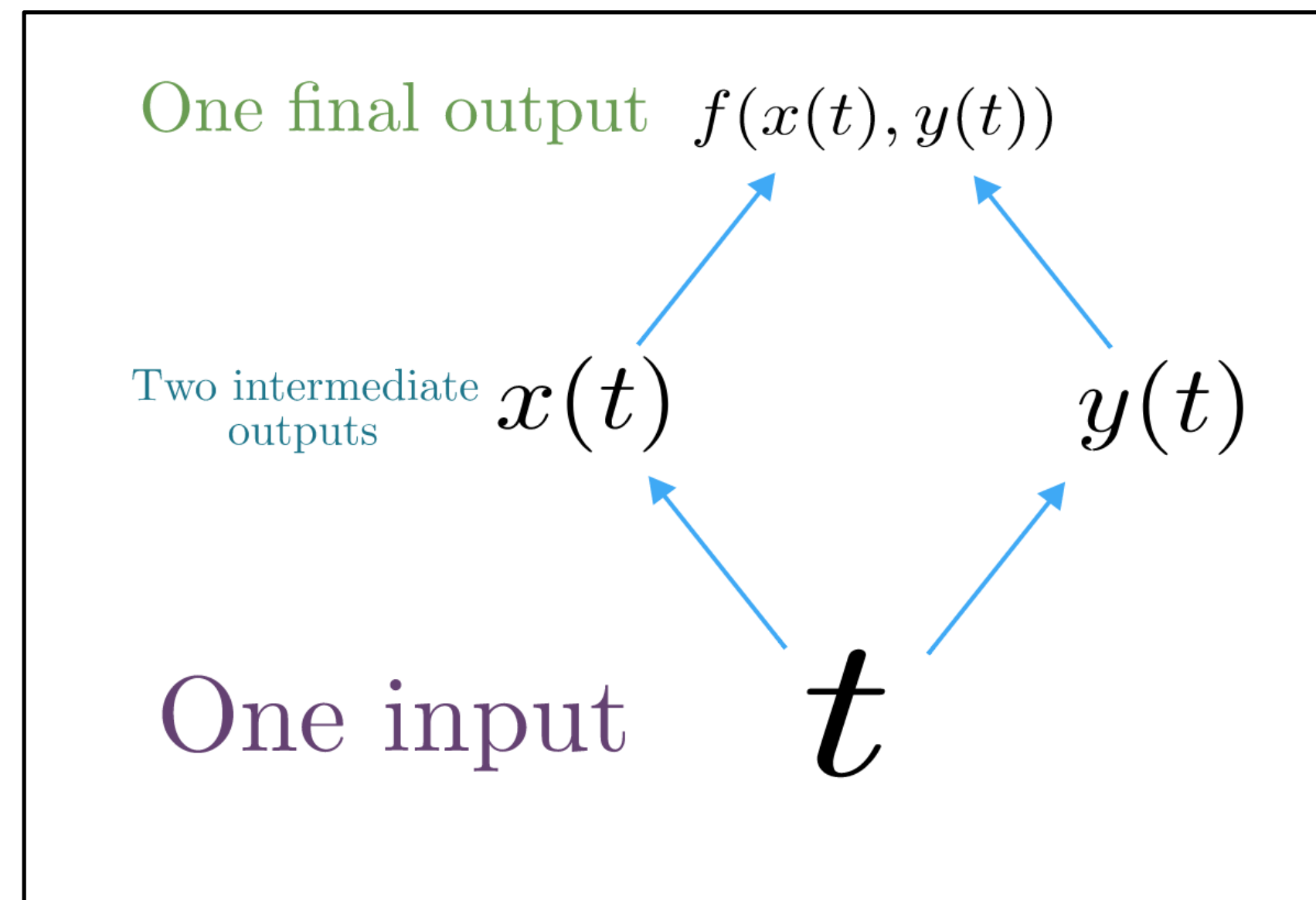
Why?

Multivariable Chain Rule

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

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

Derivative of composition function



Source:

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

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:

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

Derivative of composition function

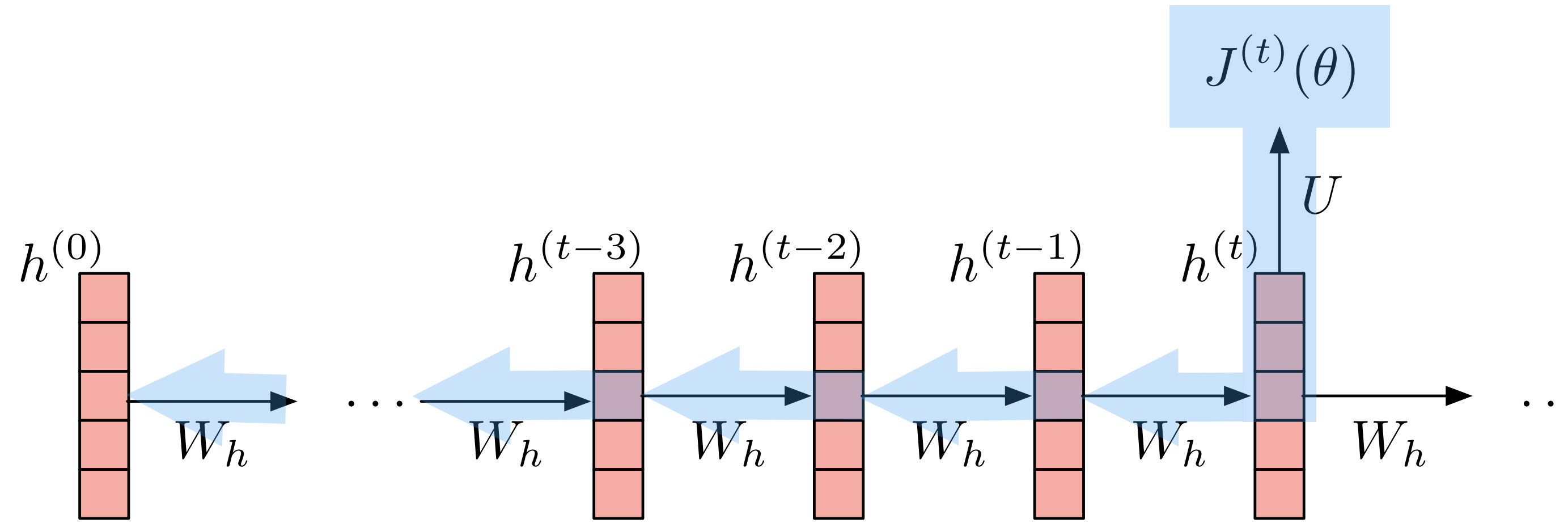
In our example:

Apply the multivariable chain rule:

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

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

44 Backpropagation for RNNs



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]

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

Question: How do we calculate this?

RNN Applications

46 Generating Text with a RNN Language Model

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>

47 Generating Text with a RNN Language Model

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>

48 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

Inverse probability of corpus, according to Language Model

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!

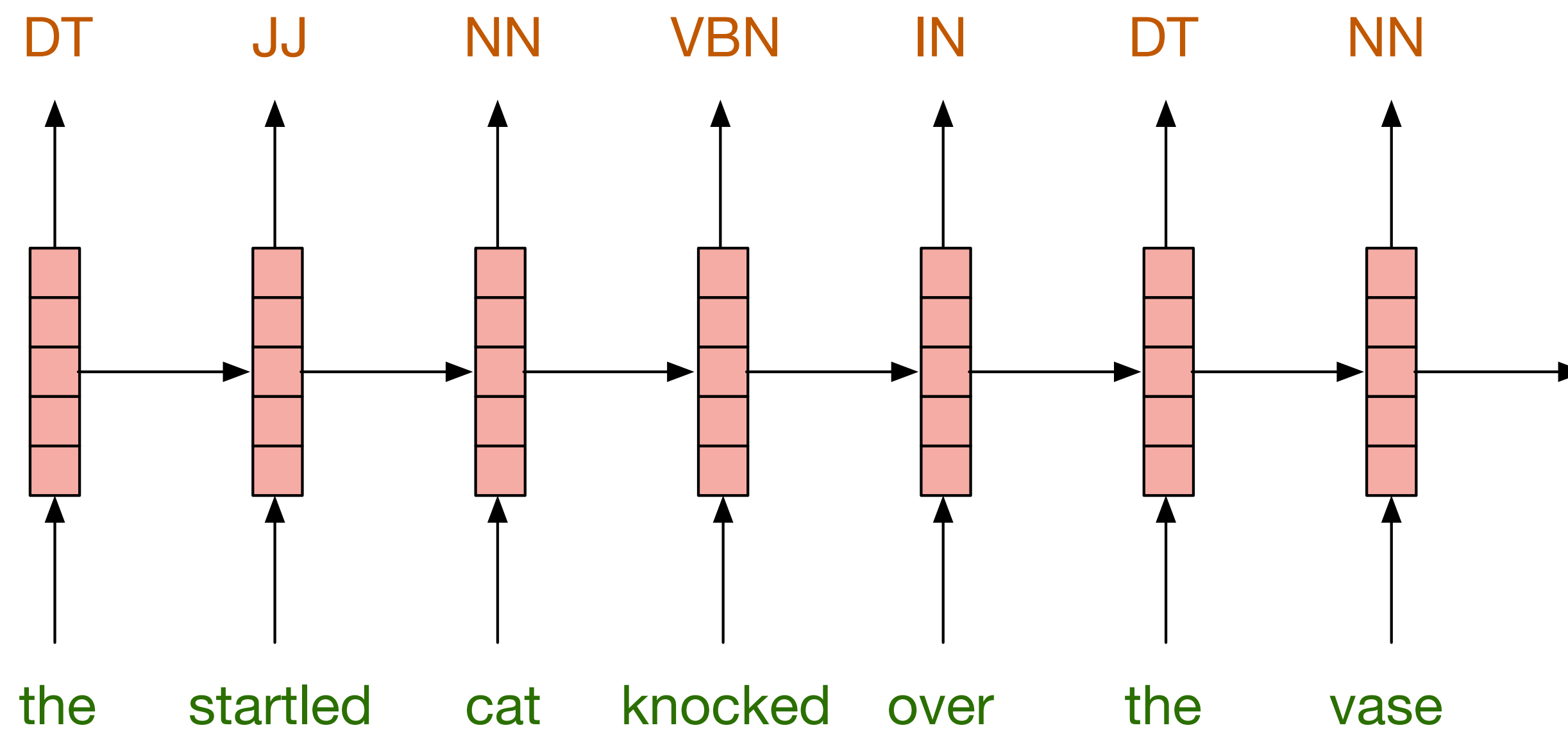
49 Why Should We Care about Language Modeling

Language Modeling is a **benchmark task** that helps us **measure our progress** on understanding language

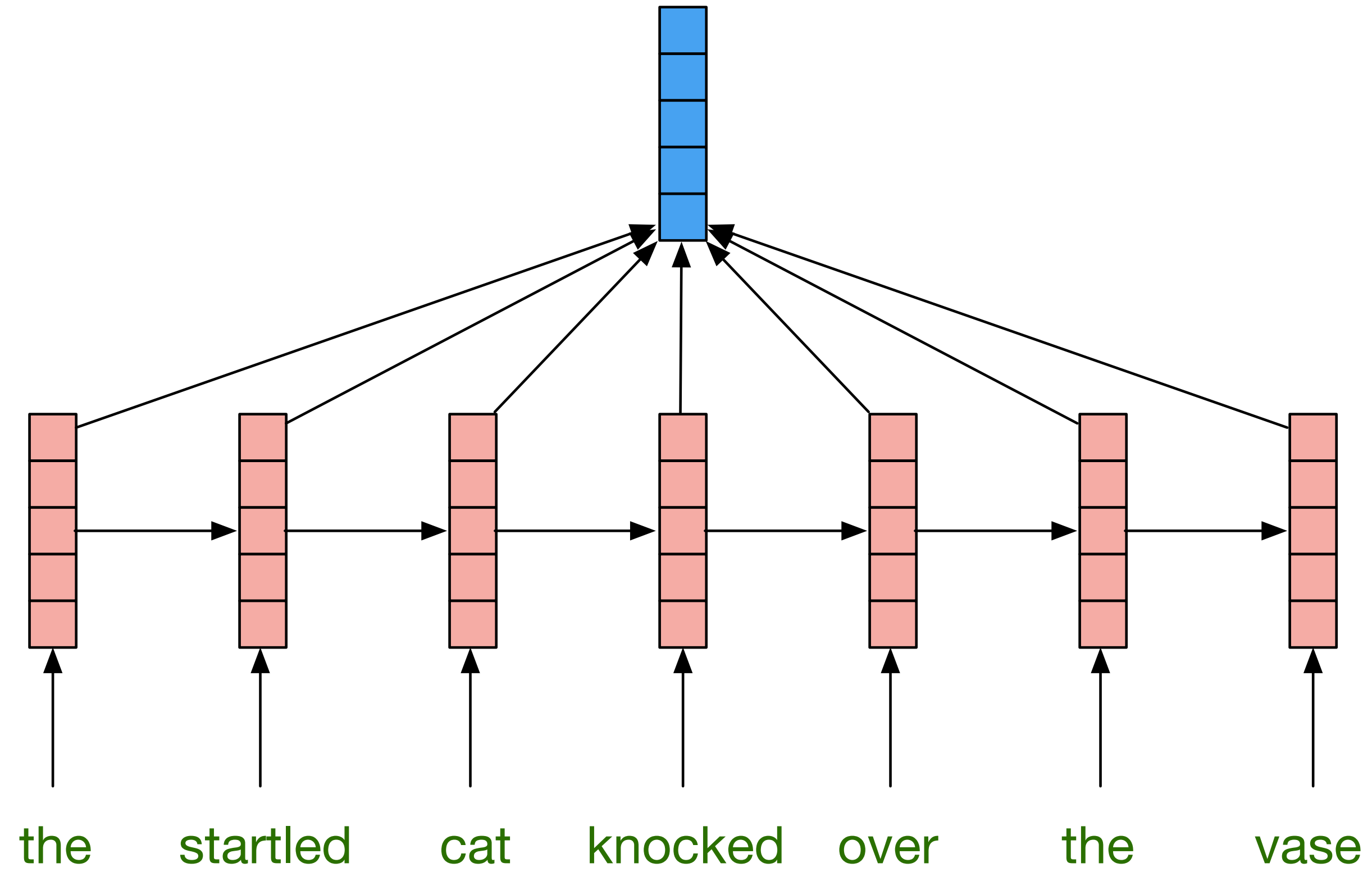
Language Modeling is a **subcomponent** of many NLP tasks, especially those involving **generating text** or **estimating the probability of text**:

- Predictive typing
- Speech recognition
- Handwriting recognition
- Spelling/grammar correction
- Authorship identification
- Machine translation
- Summarization
- Dialogue
- etc.

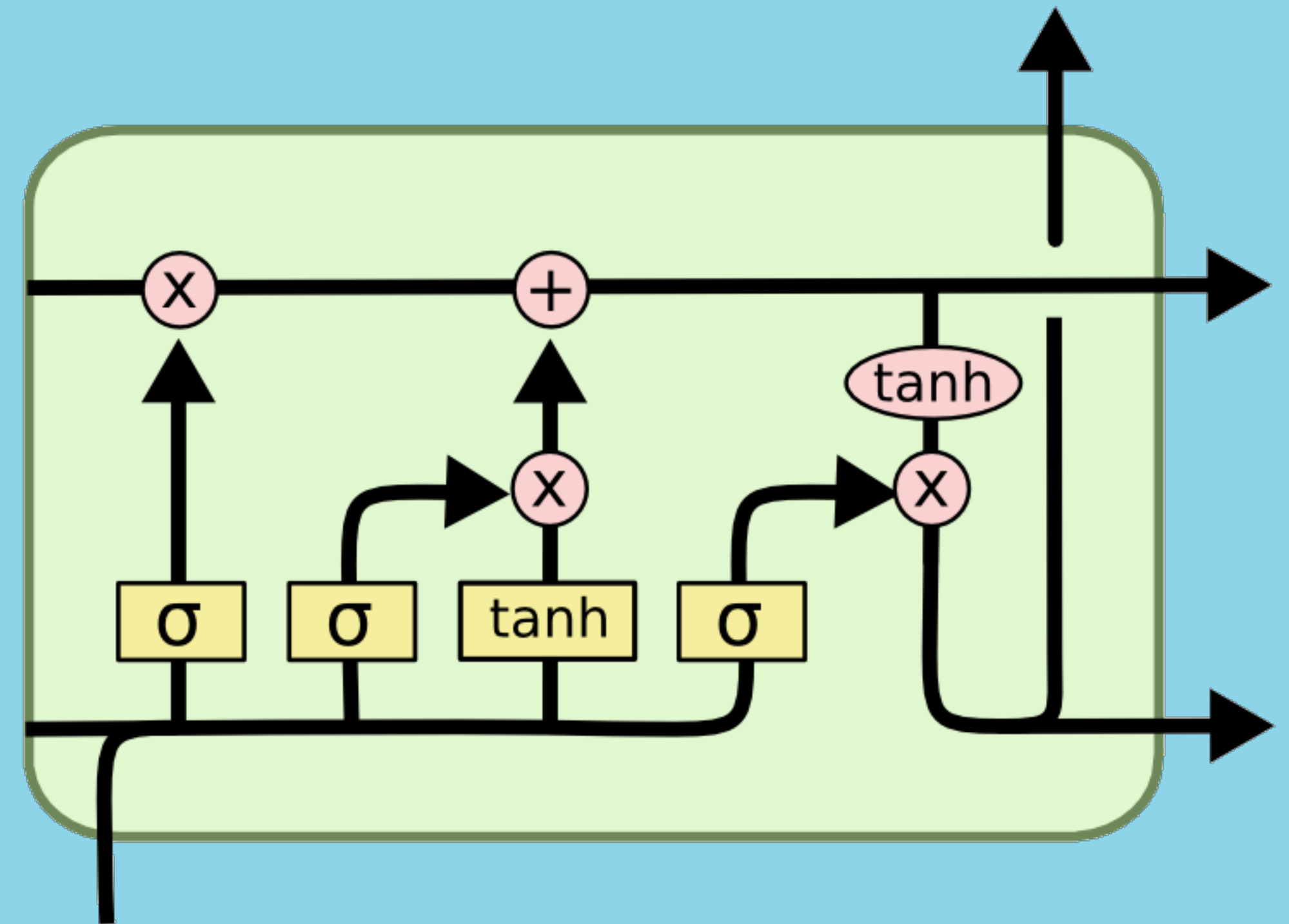
RNNs Can Be Used for Tagging



51 RNNs Can Be Used as an Encoder Module



RNN Variants



Get RNNs to Work

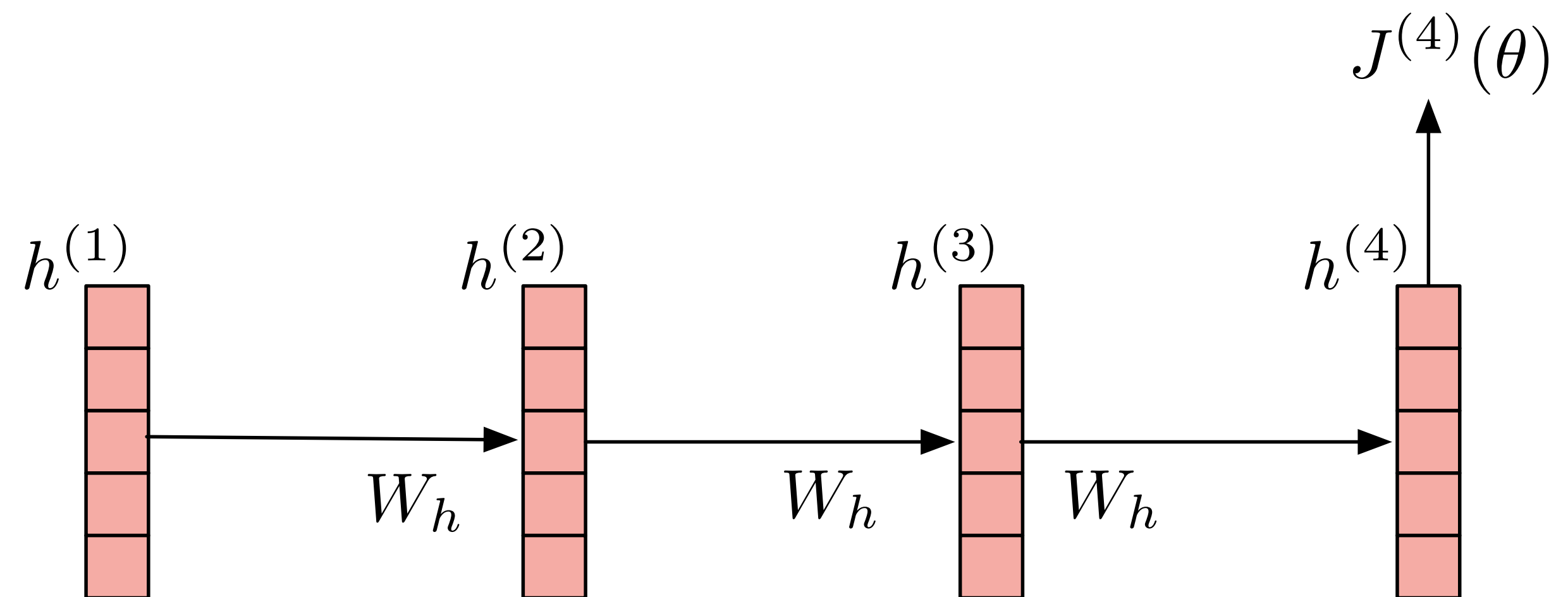
- Vanishing gradient problem
- Two new types of RNN: LSTM and GRU
- Other fixes for vanishing (or exploding) gradient:
 - Gradient clipping
 - Skip connections
- More fancy RNN variants:
 - Bidirectional RNNs
 - Multi-layer RNNs

motivates

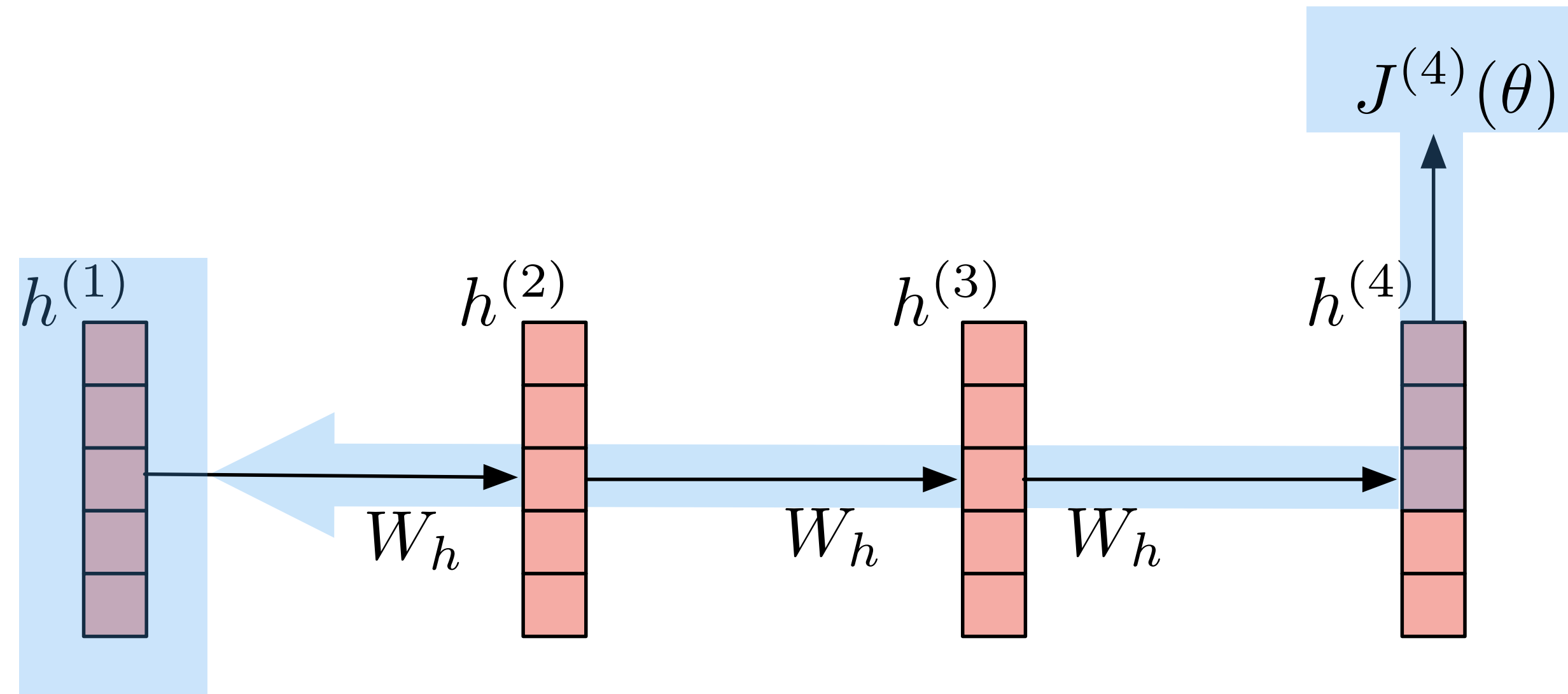
Lots of important definitions today!

Problems with RNNs: Vanishing Gradient and Exploding Gradient Problem

Vanishing Gradient Intuition

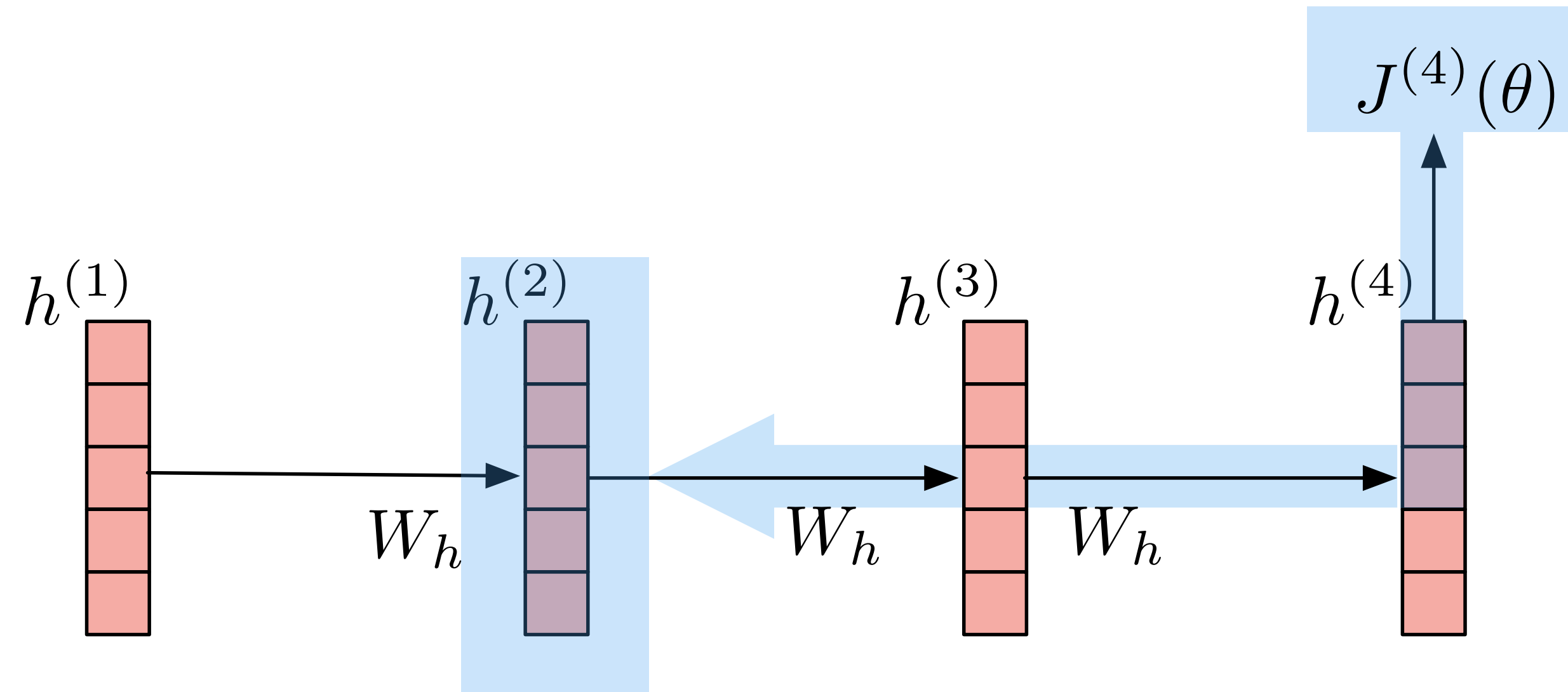


Vanishing Gradient Intuition



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

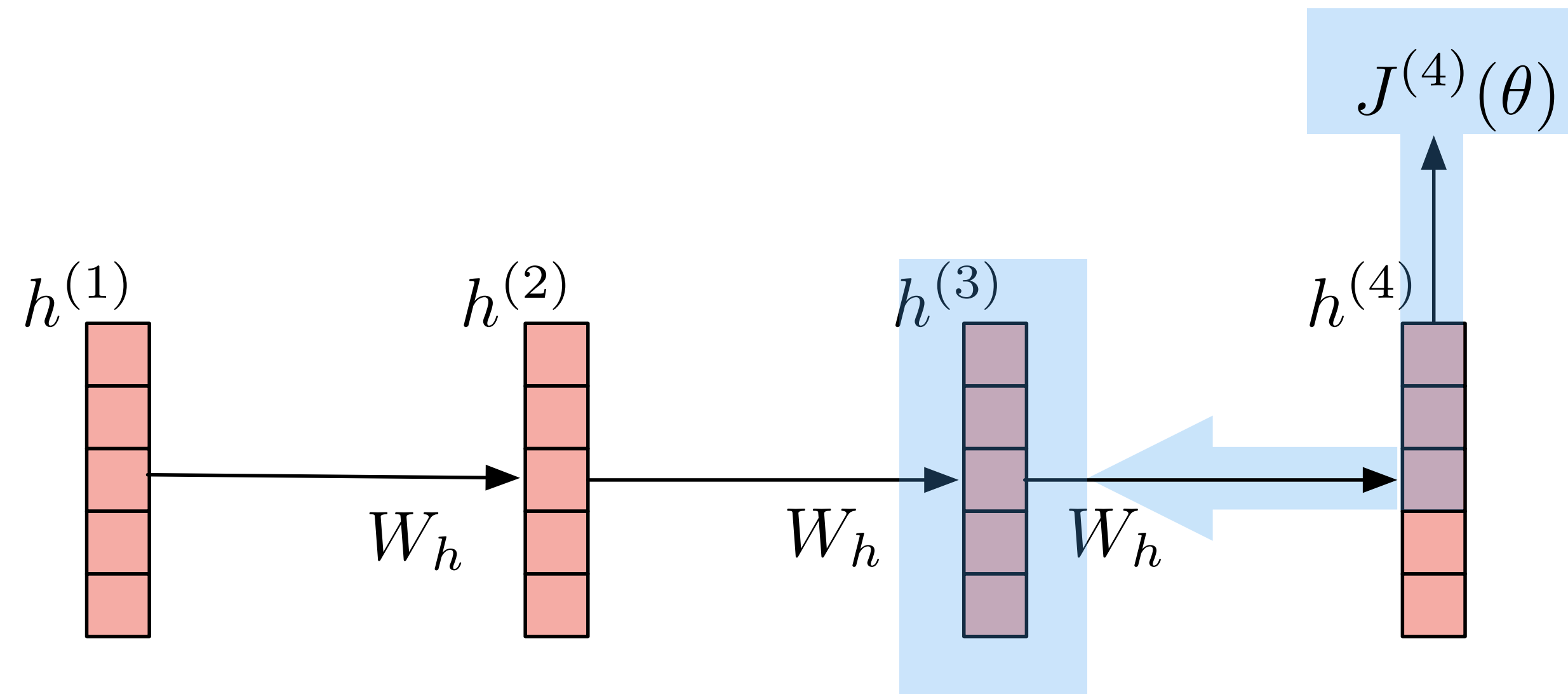
57 Vanishing Gradient Intuition



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

Chain Rule!

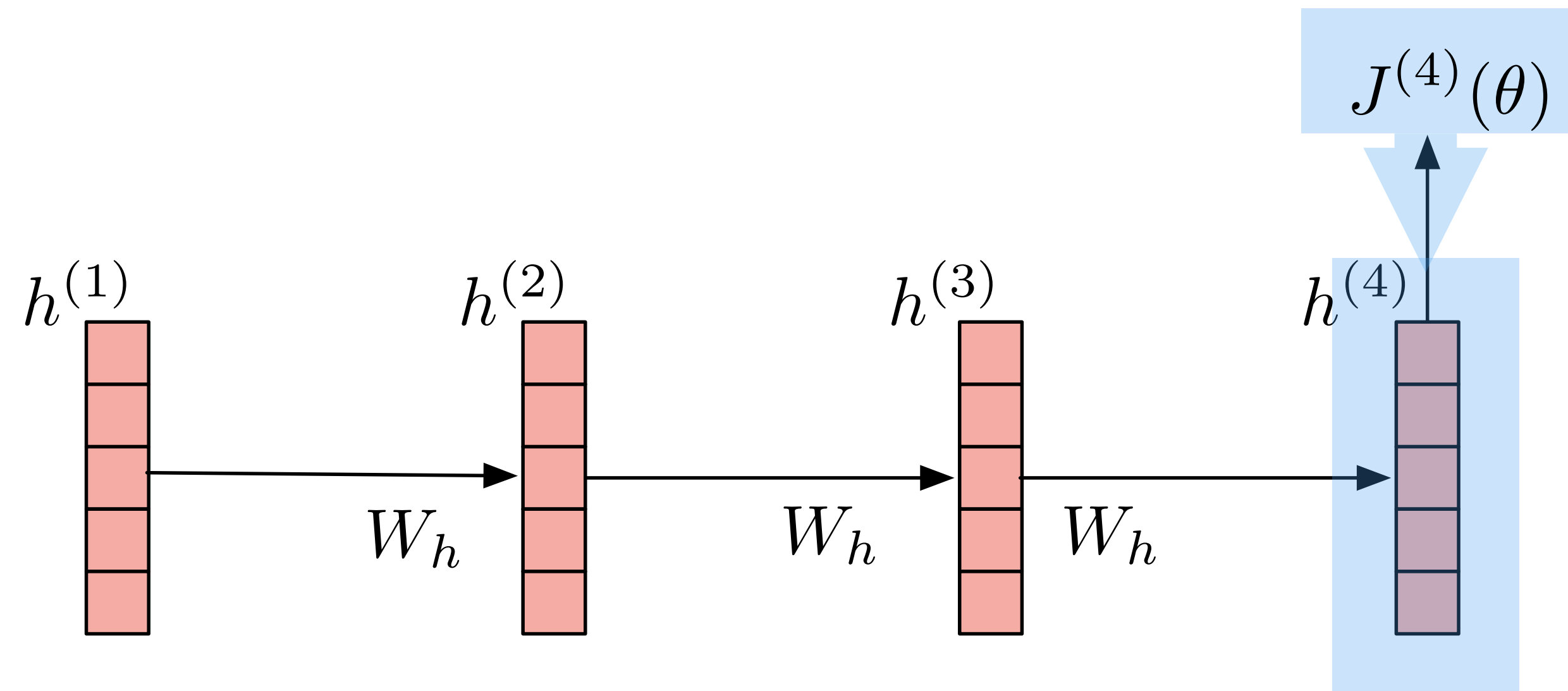
Vanishing Gradient Intuition



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

Chain Rule!

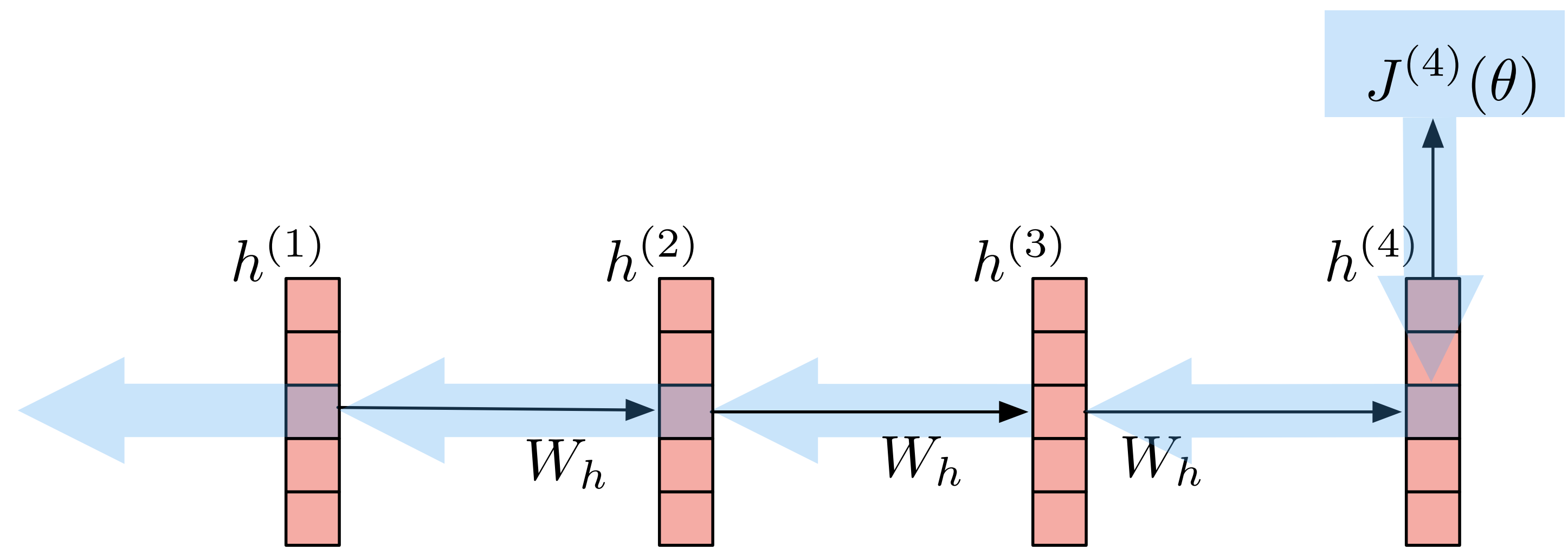
Vanishing Gradient Intuition



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

Chain Rule!

Vanishing Gradient Intuition



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

What happens if these are small?

Vanishing gradient problem:
the gradient signal gets smaller and smaller as it backpropagates further

61 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, $\sigma(x) = x$?

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

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 $\ell = i - j$

$$\begin{aligned} \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(j)}} &= \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \prod_{j < t \leq i} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} && \text{(chain rule)} \\ &= \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \prod_{j < t \leq i} \mathbf{W}_h = \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \mathbf{W}_h^\ell && \text{(value of } \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} \text{)} \end{aligned}$$

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

62 Vanishing Gradient Proof Sketch (Linear Case)

What's wrong with \mathbf{W}_h^ℓ ?

Consider if the eigenvalues of \mathbf{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 $\frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \mathbf{W}_h^\ell$ using the eigenvectors of \mathbf{W}_h as a basis:

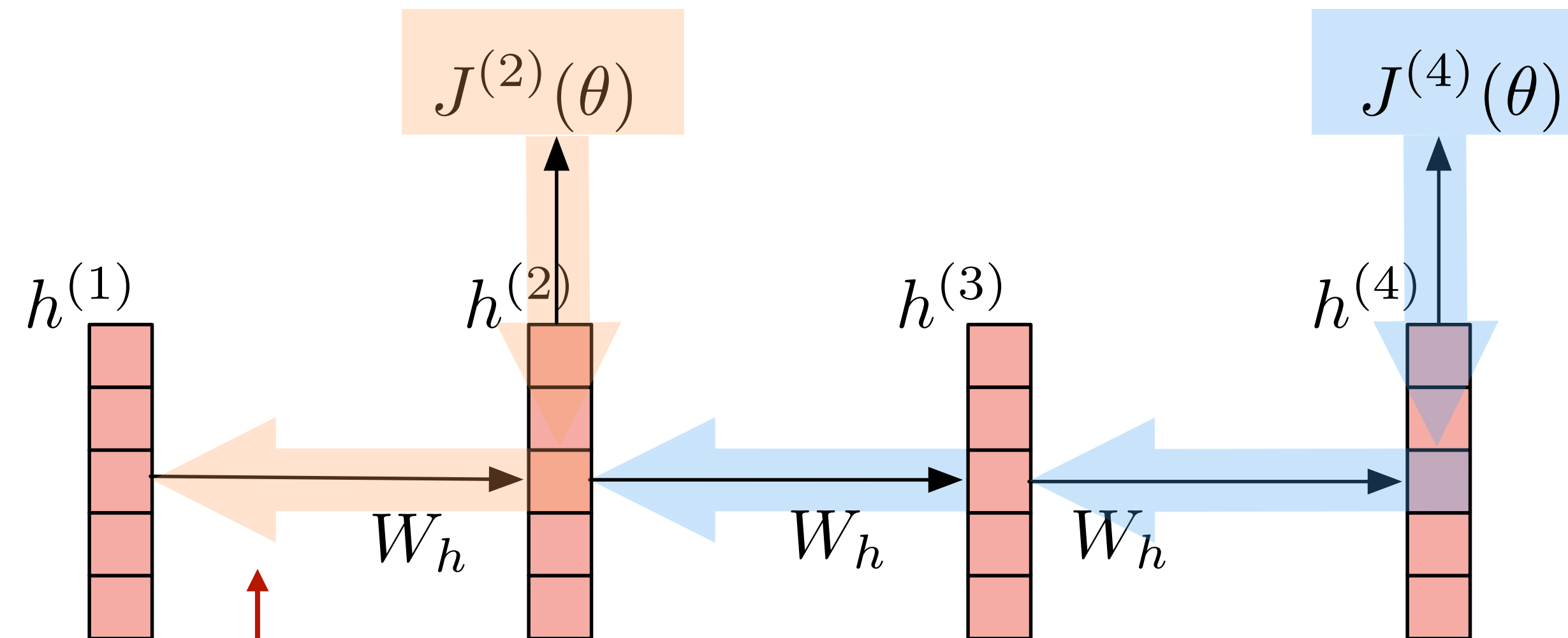
$$\frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \mathbf{W}_h^\ell = \sum_{i=1}^n c_i \lambda_i^\ell \mathbf{q}_i \approx \mathbf{0} \text{ (for large } \ell \text{)}$$

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

63 Why is Vanishing Gradient a Problem?



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

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

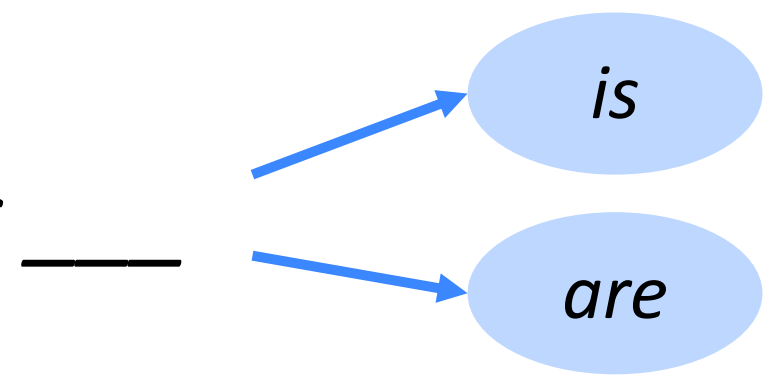


64 Why is Vanishing Gradient a Problem?

- Another explanation: Gradient can be viewed as a measure of *the effect of the past on the future*
- If the gradient becomes vanishingly small over longer distances (step t to step $t+n$), then we can't tell whether:
 1. There's **no dependency** between step t and $t+n$ in the data
 2. We have **wrong parameters** to capture the true dependency between t and $t+n$

65 Effect of Vanishing Gradient on RNN-LM

- **LM task:** *When she tried to print her tickets, she found that the printer was out of toner. She went to the stationery store to buy more toner. It was very overpriced. After installing the toner into the printer, she finally printed her _____*
- To learn from this training example, the RNN-LM needs to **model the dependency** between “*tickets*” on the 7th step and the target word “*tickets*” at the end.
- But if gradient is small, the model **can't learn this dependency**
 - So the model is **unable to predict similar long-distance dependencies** at test time

66 Effect of Vanishing Gradient on RNN-LM

- **LM task:** *The writer of the books ____* 
- **Correct answer:** *The writer of the books is planning a sequel*
- **Syntactic recency:** *The writer of the books is* (correct) 
- **Sequential recency:** *The writer of the books are* (incorrect) 
- Vanishing gradient problems may bias RNN-LMs towards learning from **sequential recency**, so they make this type of error more often than we'd like. [Linzen et al 2016]

67 Why is Exploding Gradient a problem?

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

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

- This can cause **bad updates**: we take too large a step and reach a 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)

68 Gradient Clipping - Solution for Exploding Gradient

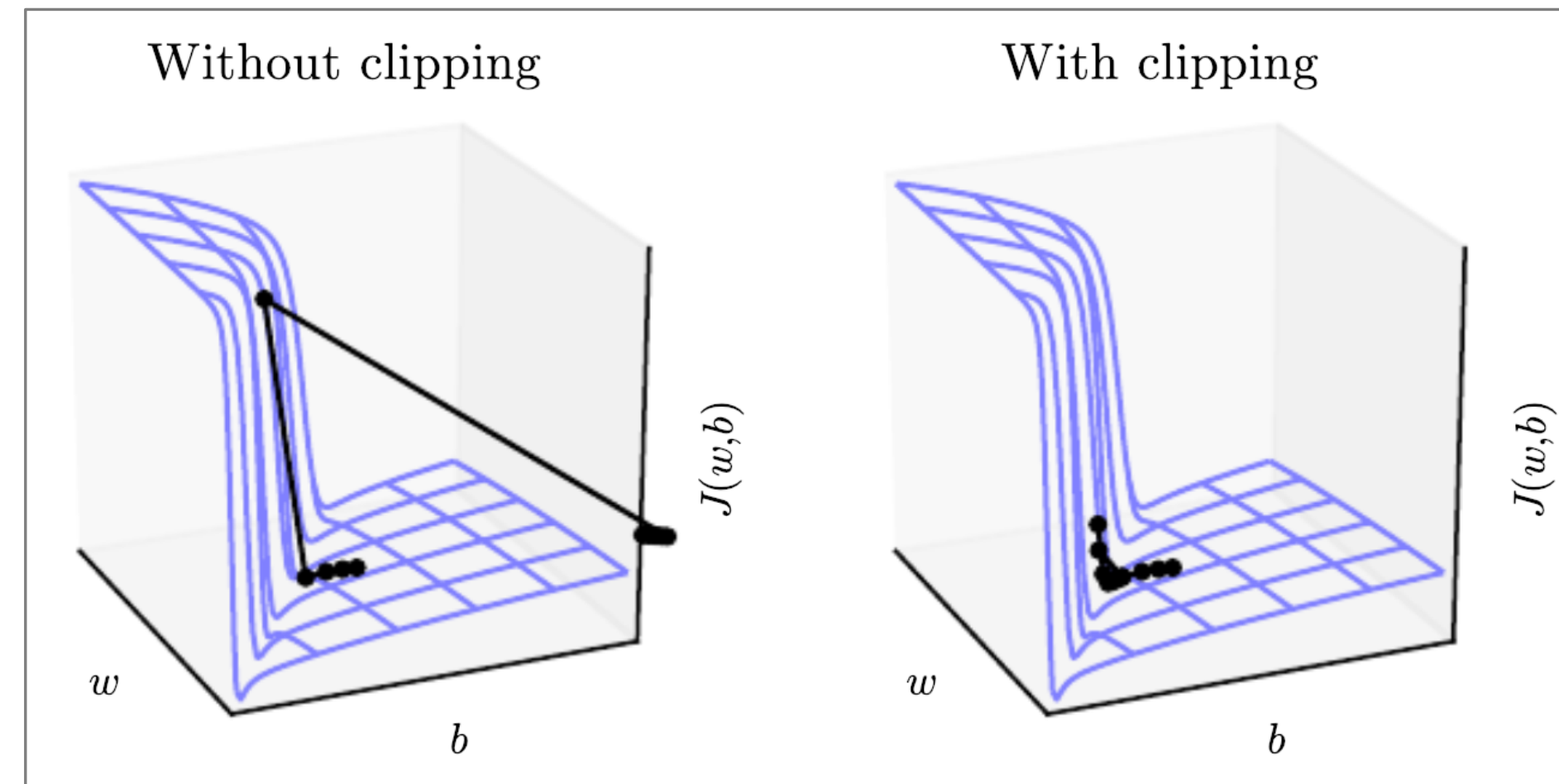
Gradient clipping: if the norm of the gradient is greater than some threshold, scale it down before applying SGD update

Algorithm 1 Pseudo-code for norm clipping

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

Intuition: take a step in the same direction, but a smaller step

69 Gradient Clipping - Solution for Exploding Gradient



- This shows the loss surface of a simple RNN (hidden state is a scalar not a vector)
- The “cliff” is dangerous because it has steep gradient
- On the left, gradient descent takes **two very big steps** due to steep gradient, resulting in climbing the cliff then shooting off to the right (both **bad updates**)
- On the right, gradient clipping reduces the size of those steps, so effect is **less drastic**

70 How to Fix 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 a RNN with separate memory?

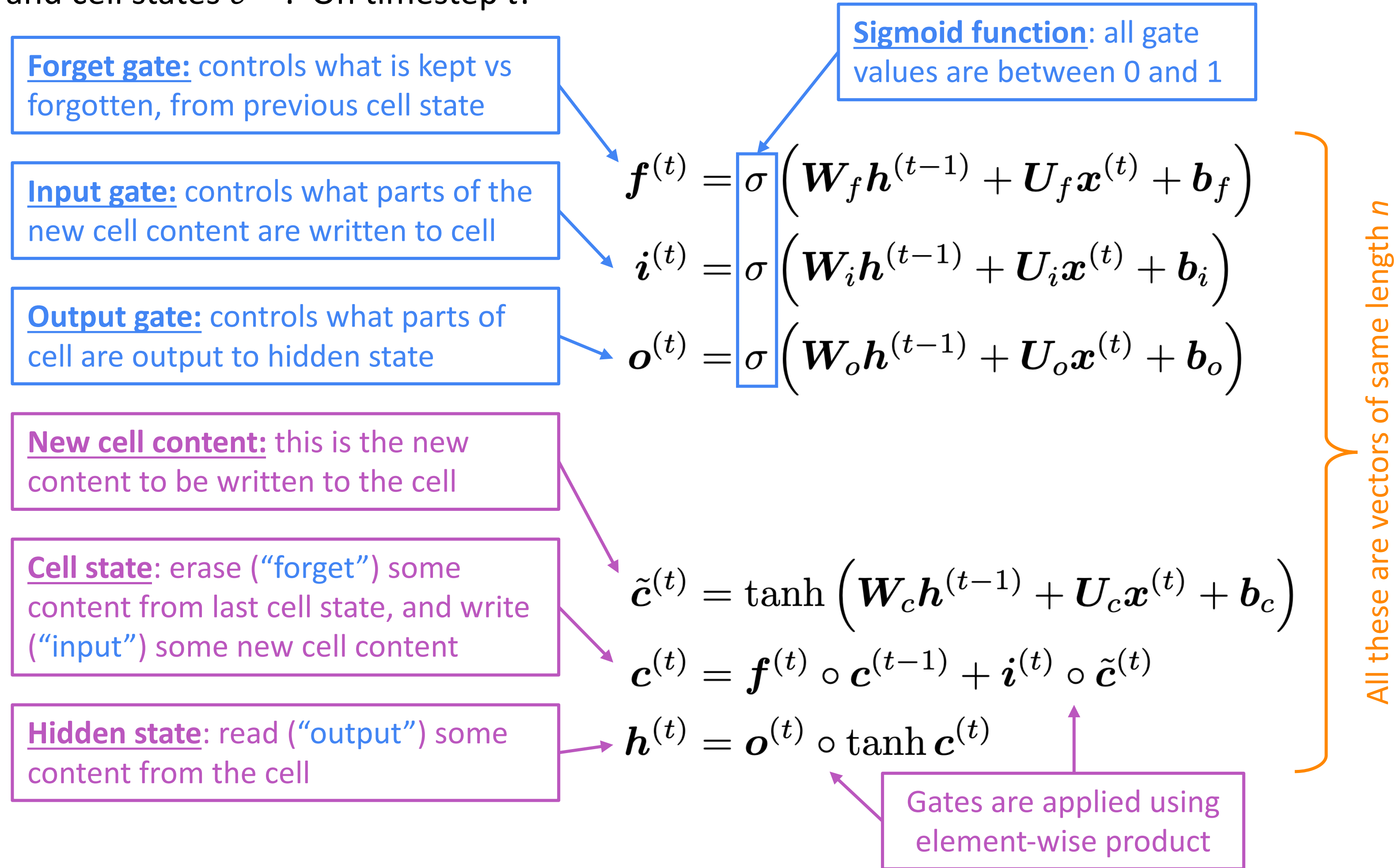
LSTM and GRU

72 Long Short-Term Memory (LSTM)

- A type of RNN proposed by Hochreiter and Schmidhuber in 1997 as a solution to the vanishing gradients problem.
- 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 **erase**, **write** and **read** information from the cell
- The selection of which information is erased/written/read is controlled by three corresponding **gates**
 - The gates are also vectors 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

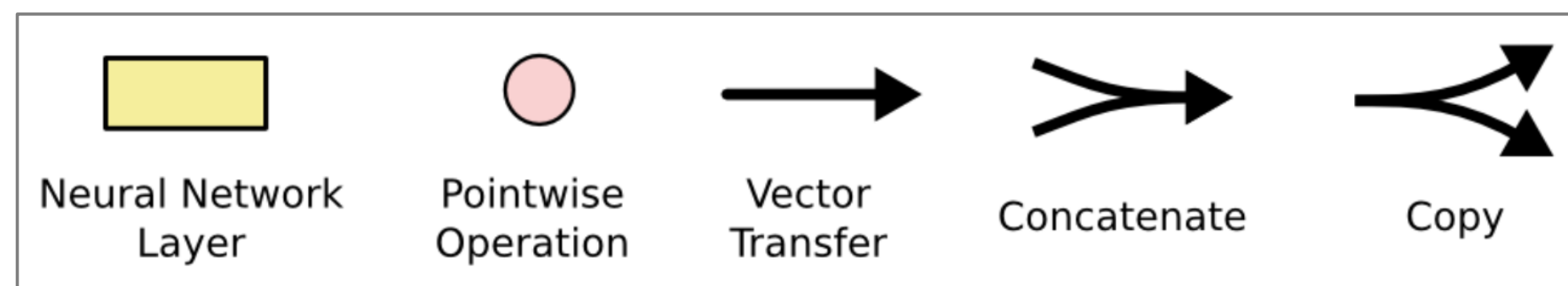
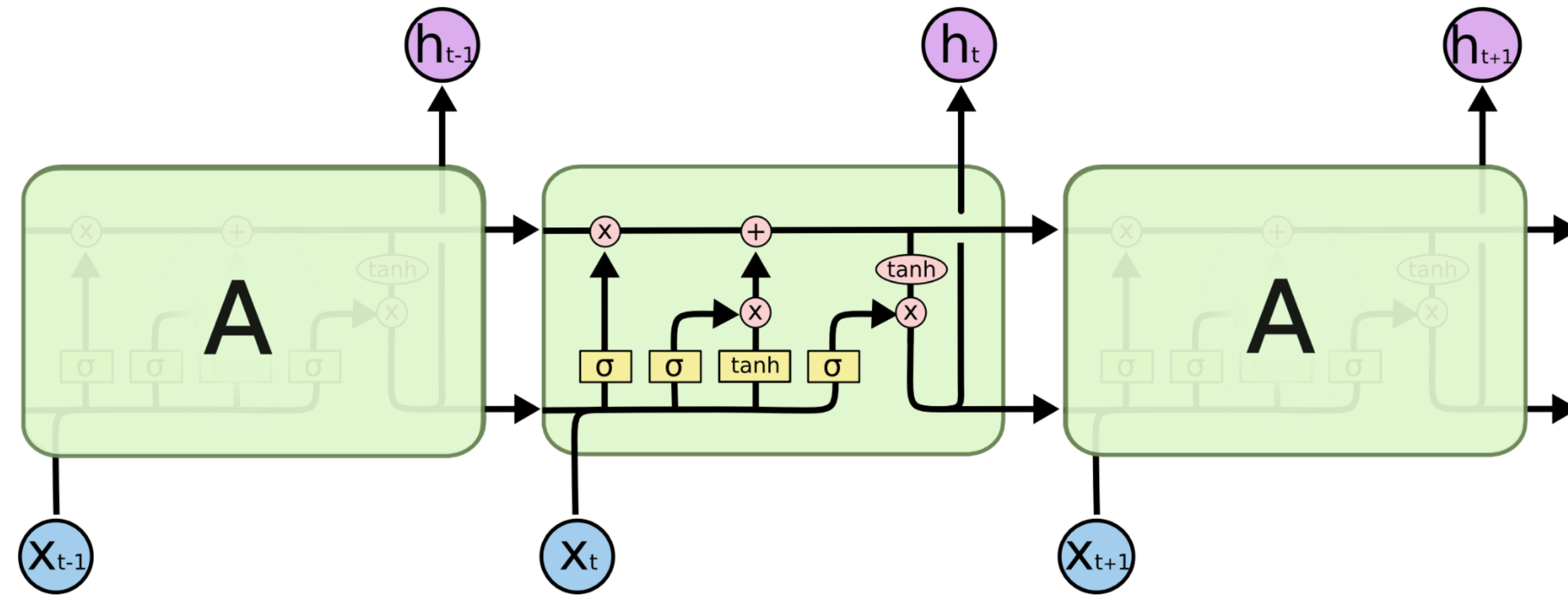
73 Long Short-Term Memory (LSTM)

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 :



74 Long Short-Term Memory (LSTM)

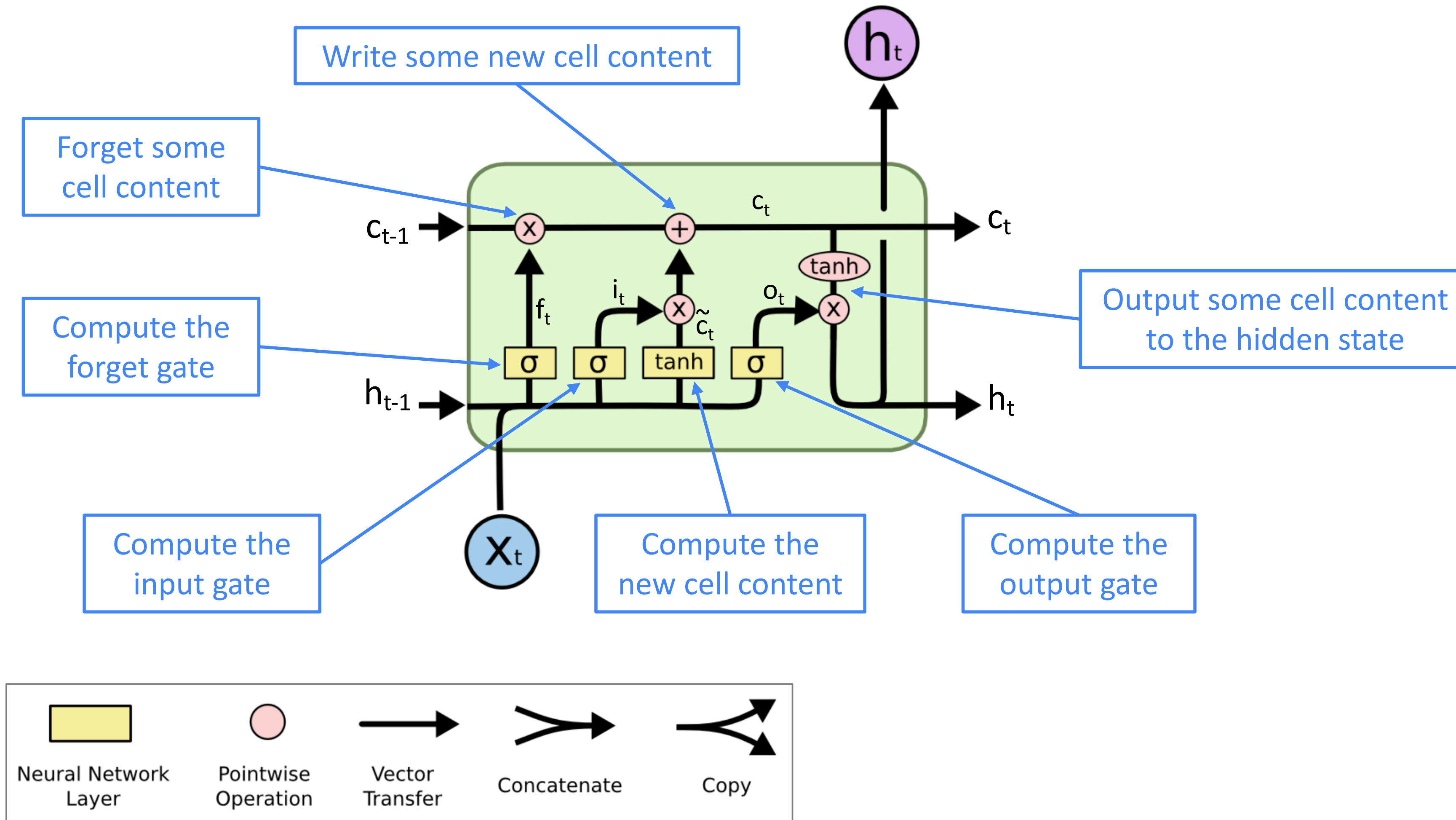
You can think of the LSTM equations visually like this:



Source: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

75 Long Short-Term Memory (LSTM)

You can think of the LSTM equations visually like this:



Source: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

76 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.
 - By contrast, it's harder for vanilla RNN to learn a recurrent weight matrix W_h that preserves info in 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

77 LSTMs: Real-World Success

- In 2013-2015, LSTMs started achieving state-of-the-art results
 - Successful tasks include: handwriting recognition, speech recognition, machine translation, parsing, image captioning
 - LSTM became the dominant approach
- Now (2020), other approaches (e.g. Transformers) have become more dominant for certain tasks.
 - For example in WMT (a MT conference + competition):
 - In WMT 2016, the summary report contains "RNN" 44 times
 - In WMT 2018, the report contains "RNN" 9 times and "Transformer" 63 times
 - In WMT 2019: "RNN" 7 times, "Transformer" 105 times

Source: "Findings of the 2016 Conference on Machine Translation (WMT16)", Bojar et al. 2016, <http://www.statmt.org/wmt16/pdf/W16-2301.pdf>

Source: "Findings of the 2018 Conference on Machine Translation (WMT18)", Bojar et al. 2018, <http://www.statmt.org/wmt18/pdf/WMT028.pdf>

Source: "Findings of the 2019 Conference on Machine Translation (WMT19)", Barrault et al. 2019, <http://www.statmt.org/wmt18/pdf/WMT028.pdf>

Gated Recurrent Unites (GRU)

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

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

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

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

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

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

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

$$\mathbf{h}^{(t)} = (1 - \mathbf{u}^{(t)}) \circ \mathbf{h}^{(t-1)} + \mathbf{u}^{(t)} \circ \tilde{\mathbf{h}}^{(t)}$$

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

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)

- 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.

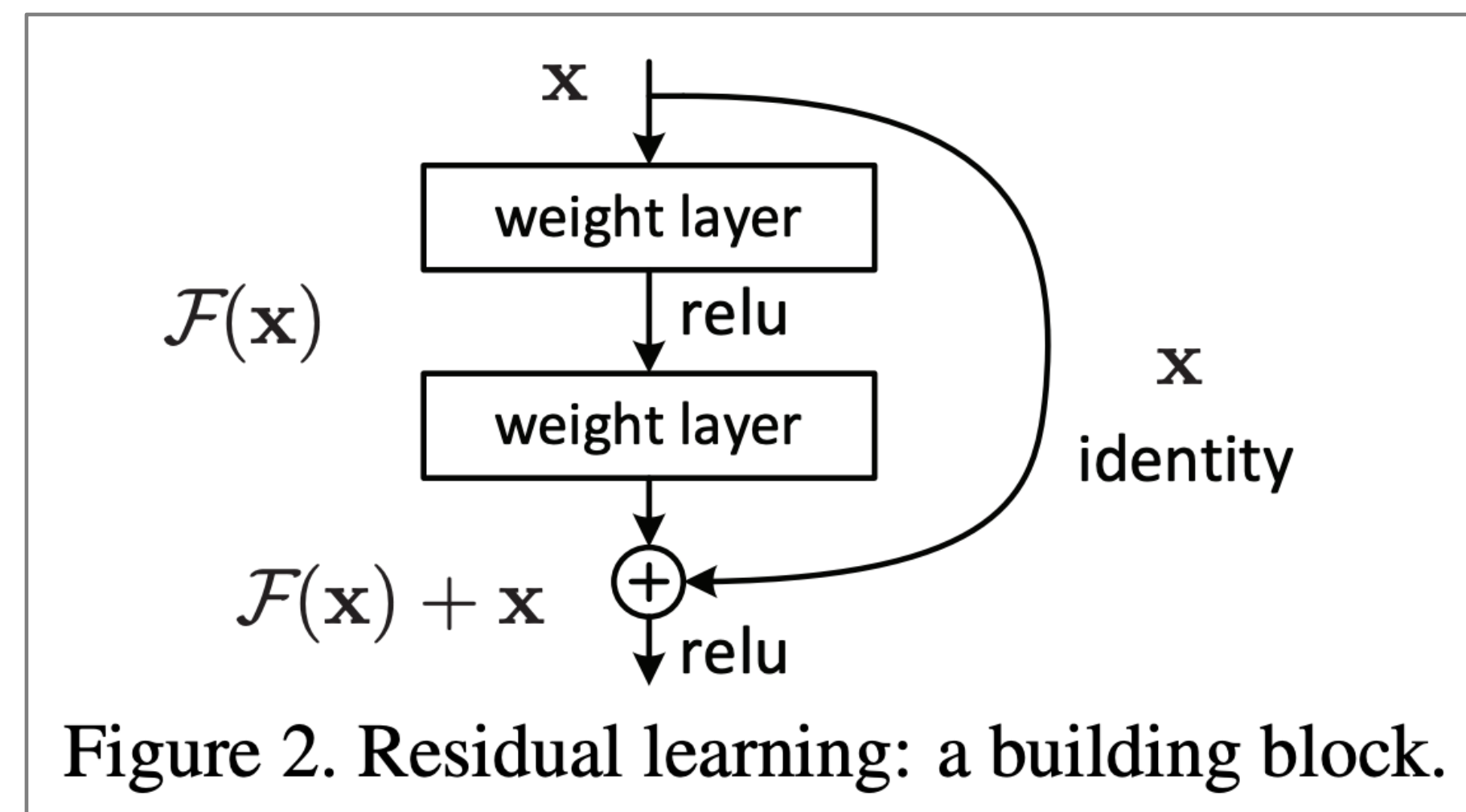
80 Is Vanishing/Exploding Gradient just a RNN problem?

- No! It can be a problem for all neural architectures (including **feed-forward** and **convolutional**), especially **deep** ones.
 - Due to chain rule / choice of nonlinearity function, gradient can become vanishingly small as it backpropagates
 - Thus lower layers are learnt very slowly (hard to train)
 - Solution: lots of new deep feedforward/convolutional architectures that **add more direct connections** (thus allowing the gradient to flow)

81 Residual Connections

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

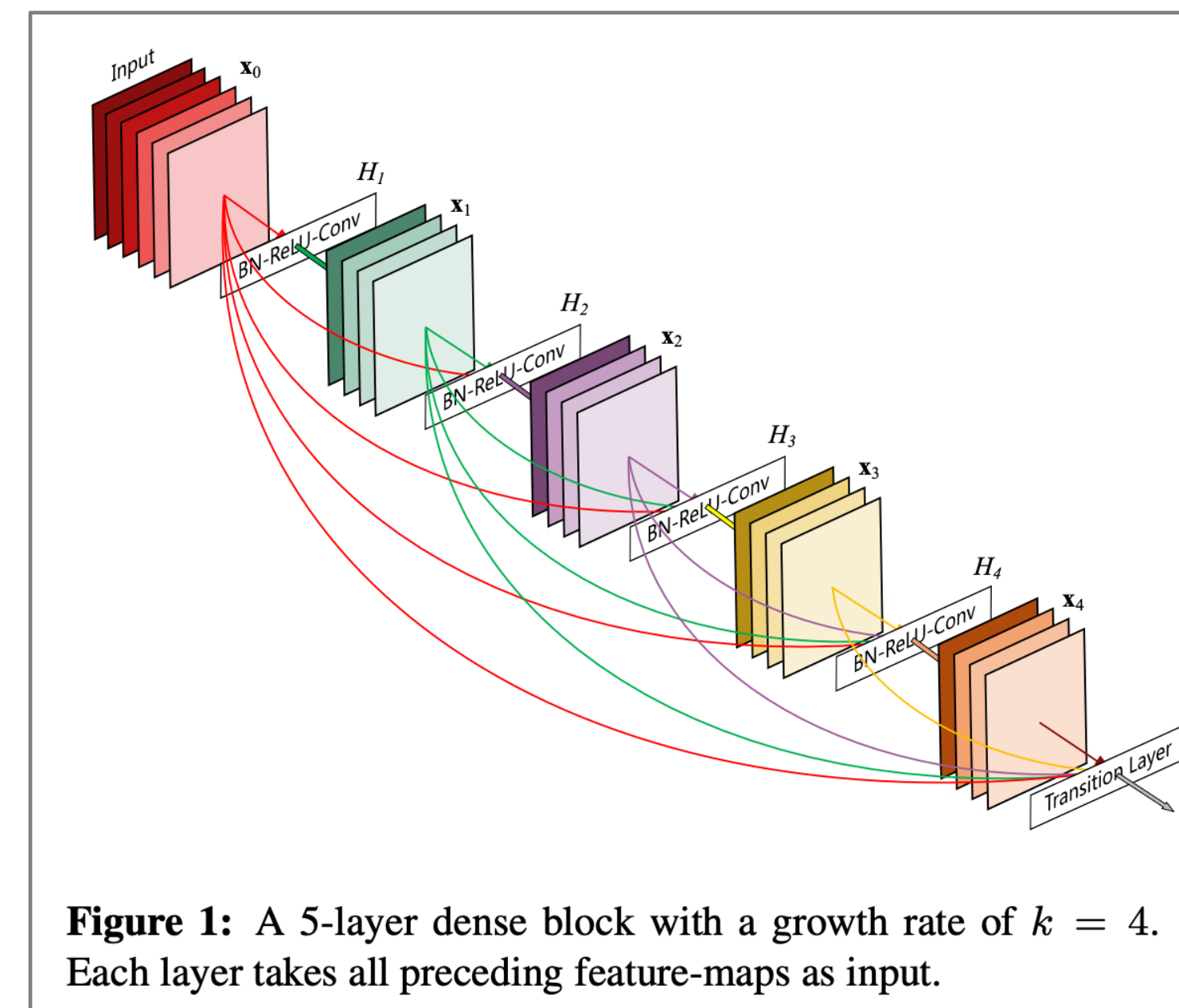


"Deep Residual Learning for Image Recognition", He et al, 2015. <https://arxiv.org/pdf/1512.03385.pdf>

Dense Connections

For example:

- Dense connections aka “DenseNet”
- Directly connect each layer to all future layers!



“Densely Connected Convolutional Networks”, Huang et al, 2017. <https://arxiv.org/pdf/1608.06993.pdf>

83 Highway Connections

For example:

- 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

“Highway Networks”, Srivastava et al, 2015. <https://arxiv.org/pdf/1505.00387.pdf>

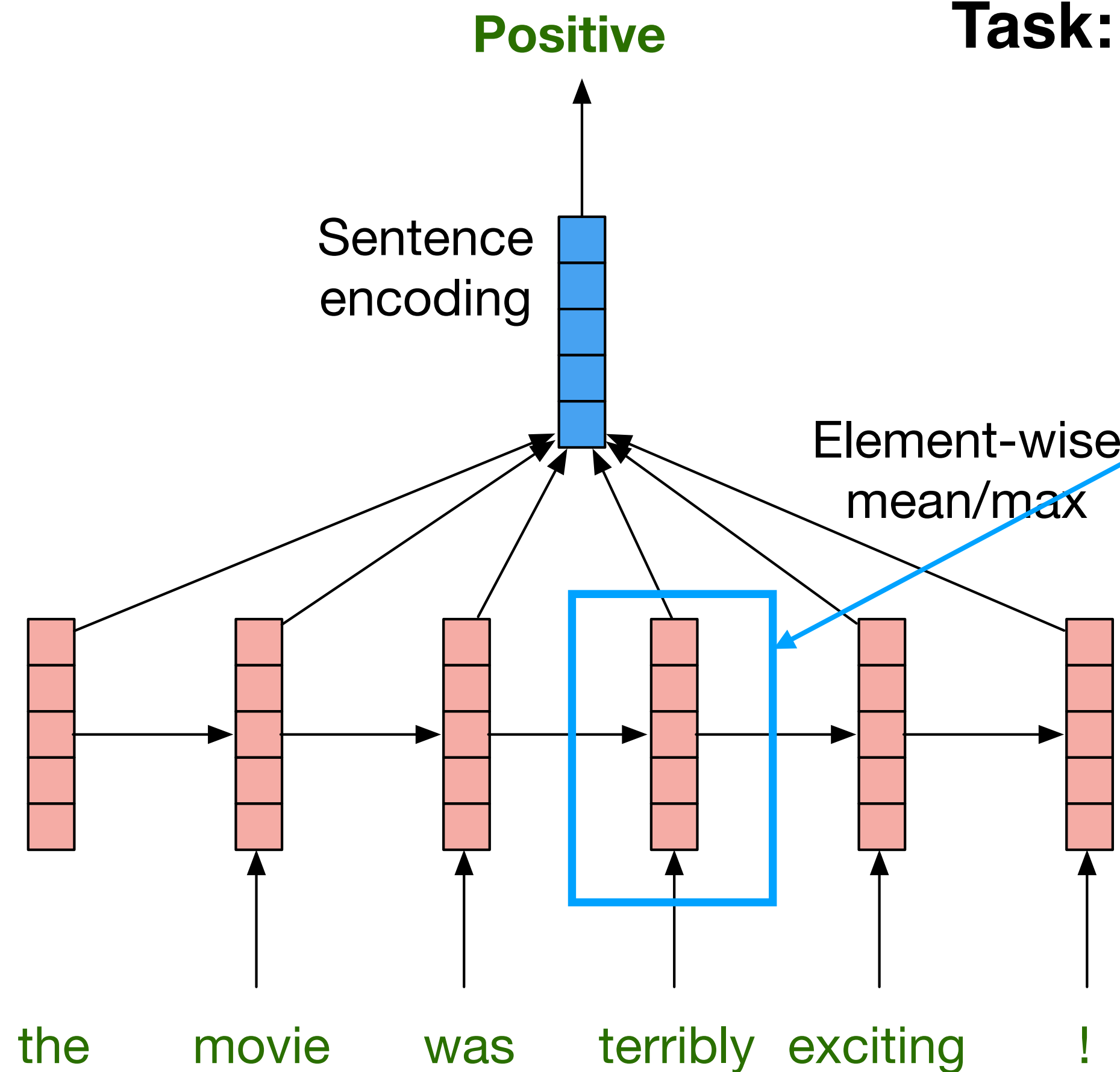
84 Is Vanishing/Exploding Gradient just a RNN problem?

- No! It can be a problem for all neural architectures (including **feed-forward** and **convolutional**), especially **deep** ones.
 - Due to chain rule / choice of nonlinearity function, gradient can become vanishingly small as it backpropagates
 - Thus lower layers are learnt very slowly (hard to train)
 - Solution: lots of new deep feedforward/convolutional architectures that **add more direct connections** (thus allowing the gradient to flow)
- 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]

Bi-RNN

Bidirectional RNNs: Motivation

Task: sentiment classification



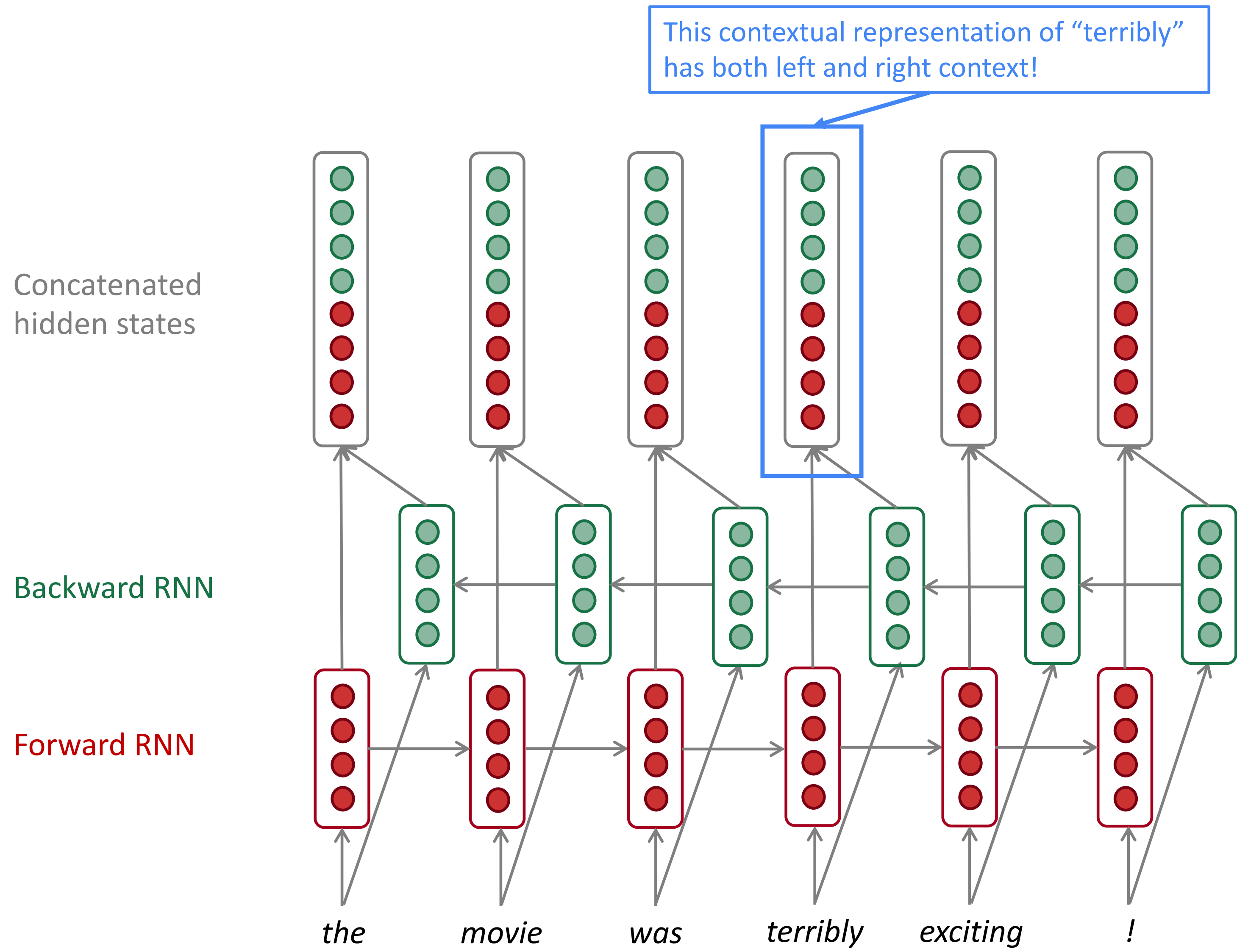
We can regard this hidden state as a representation of the word "terribly" in the context of this sentence. We call this a *contextual representation*.

These contextual representations only contain information about the *left* context (e.g. "the movie was").

What about right context?

In this example, "exciting" is in the right context and this modifies the meaning of "terribly" (from negative to positive)

Bidirectional RNNs



Bidirectional RNNs

On timestep t :

This is a general notation to mean “compute one forward step of the RNN” – it could be a vanilla, LSTM or GRU computation.

Forward RNN $\vec{h}^{(t)} = \text{RNN}_{\text{FW}}(\vec{h}^{(t-1)}, \mathbf{x}^{(t)})$

Backward RNN $\overleftarrow{h}^{(t)} = \text{RNN}_{\text{BW}}(\overleftarrow{h}^{(t+1)}, \mathbf{x}^{(t)})$

Generally, these two RNNs have separate weights

Concatenated hidden states $\mathbf{h}^{(t)} = [\vec{h}^{(t)}; \overleftarrow{h}^{(t)}]$

We regard this as “the hidden state” of a bidirectional RNN. This is what we pass on to the next parts of the network.

Bidirectional RNNs

- Note: bidirectional RNNs are only applicable if you have access to the **entire input sequence**.
 - They are **not** applicable to Language Modeling, because in LM you *only* have left context available.
- If you do have entire input sequence (e.g. any kind of encoding), **bidirectionality is powerful** (you should use it by default).
- For example, **BERT (Bidirectional Encoder Representations from Transformers)** is a powerful pretrained contextual representation system **built on bidirectionality**.
 - You will learn more about BERT later in the course!

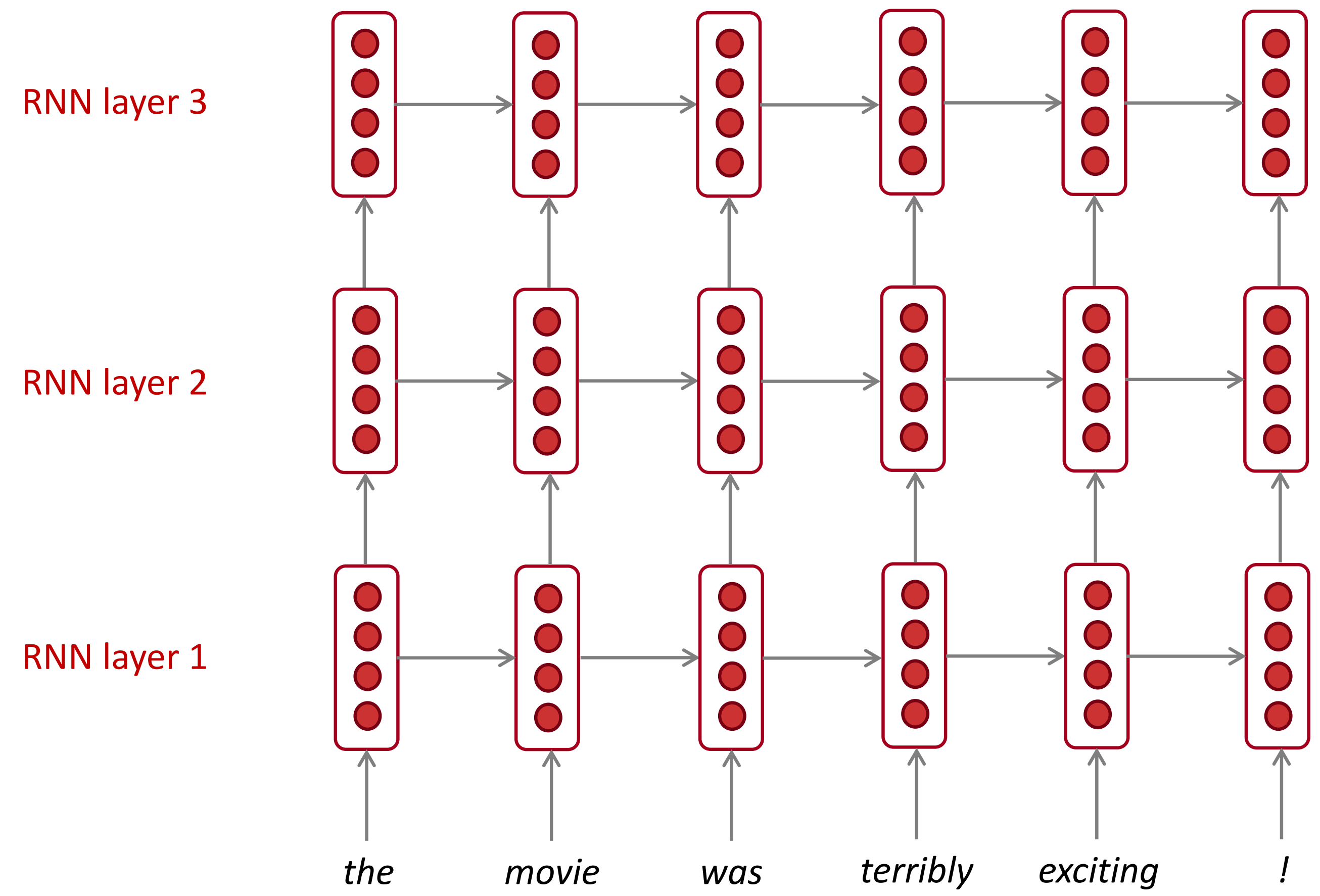
Deep-RNN

Multi-layer RNNs

- RNNs are already “deep” on one dimension (they unroll over many timesteps)
- We can also make them “deep” in another dimension by **applying multiple RNNs** – this is a multi-layer RNN.
- This allows the network to compute **more complex representations**
 - The **lower RNNs** should compute **lower-level features** and the **higher RNNs** should compute **higher-level features**.
- Multi-layer RNNs are also called ***stacked RNNs***.

Multi-layer RNNs

The hidden states from RNN layer i are the inputs to RNN layer $i+1$



Multi-layer RNNs in Practice

- High-performing RNNs are often multi-layer (but aren't as deep as convolutional or feed-forward networks)
- For example: In a 2017 paper, Britz et al find that for Neural Machine Translation, 2 to 4 layers is best for the encoder RNN, and 4 layers is best for the decoder RNN
 - However, skip-connections/dense-connections are needed to train deeper RNNs (e.g. 8 layers)
- Transformer-based networks (e.g. BERT) are frequently deeper, like 12 or 24 layers.
 - You will learn about Transformers later; they have a lot of skipping-like connections

◉ Suggested readings

- <https://web.stanford.edu/~jurafsky/slp3/3.pdf>
- <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>
- <https://www.deeplearningbook.org/contents/rnn.html>
- <http://norvig.com/chomsky.html>
- <https://arxiv.org/pdf/1211.5063.pdf>
- <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>
- <https://arxiv.org/pdf/1503.04069.pdf>

Next lecture: Word Meaning and Word Embedding

Thanks! Q&A

Bang Liu

Email: bang.liu@umontreal.ca

Homepage: <http://www-labs.iro.umontreal.ca/~liubang/>