

# CST8507: NATURAL LANGUAGE PROCESSING

## **WEEK#5** **INTRODUCTION TO LANGUAGE MODEL**

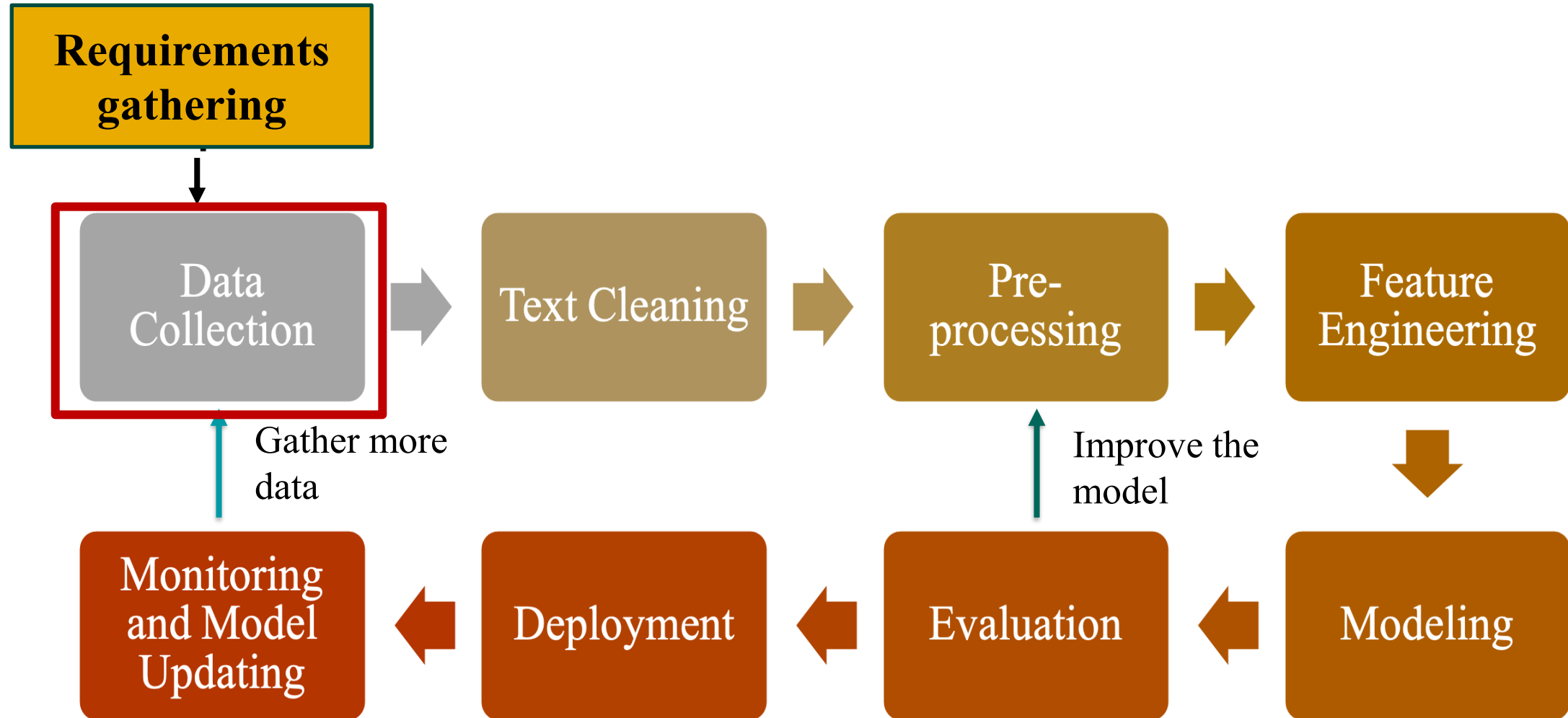
DEVELOPED BY  
HALA OWN, PH.D.

# Lesson Agenda

- Lab 3
- Text Collection(overview)
- Language Model
  - N-gram
  - NN Language model
    - Recurrent Neural Networks RNN
    - LSTMs



# NLP Development Life Cycle



# Data generated in one minute on various social platforms

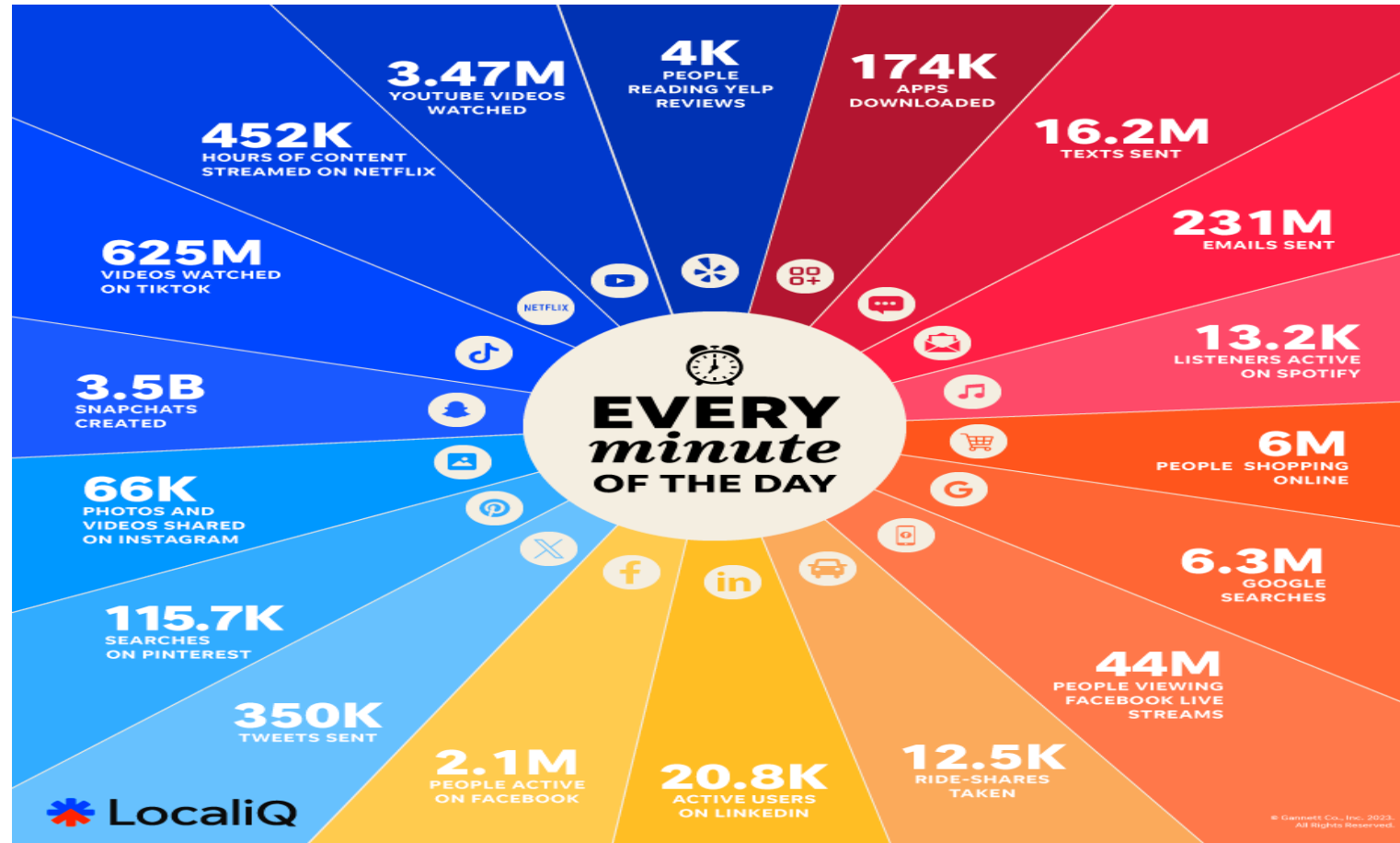
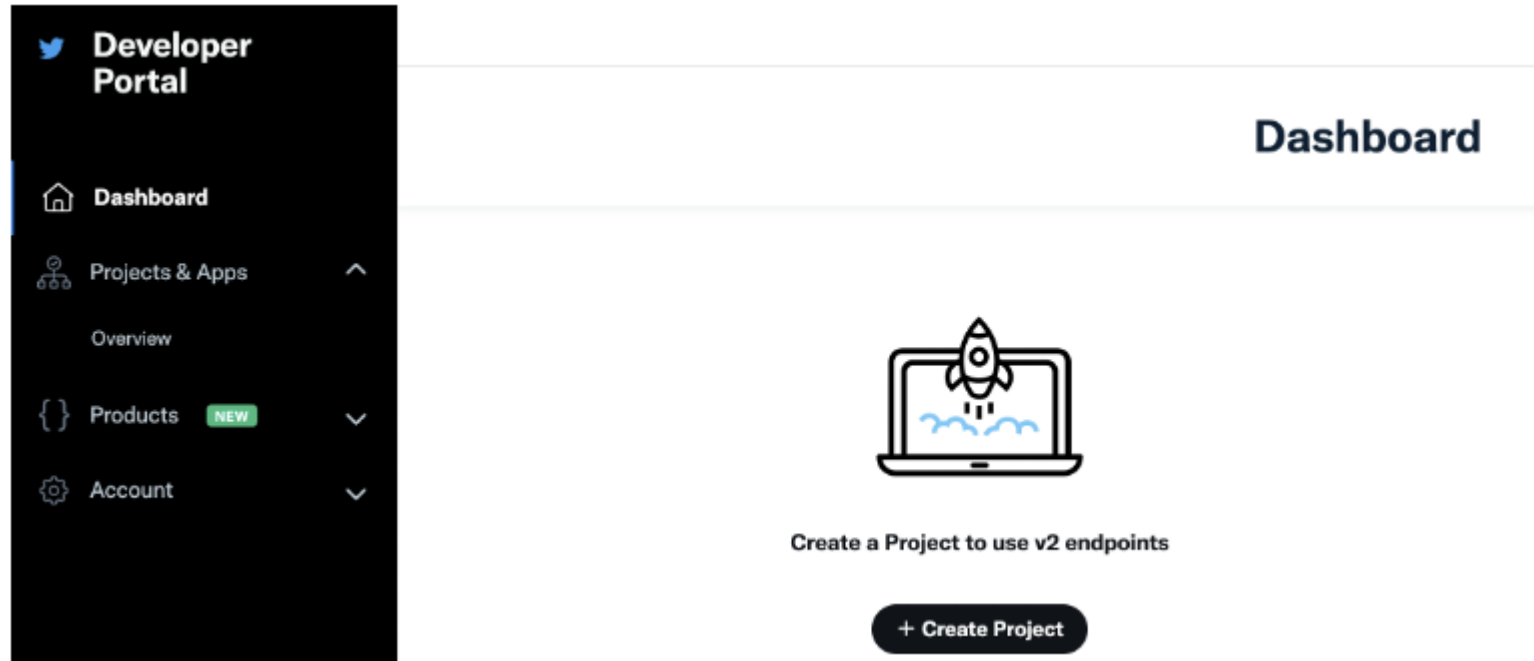


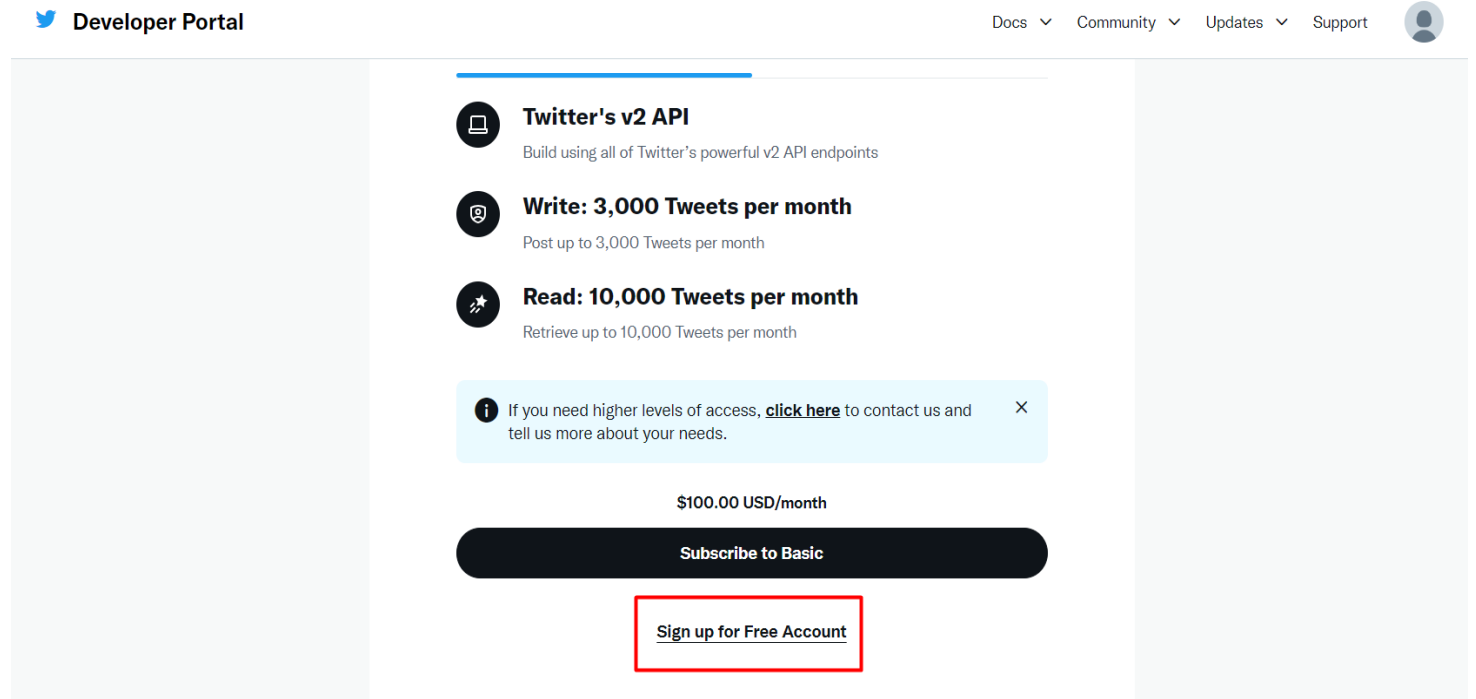
Image source: [HTTPS://localiq.com/blog/what-happens-in-an-internet-minute/](https://localiq.com/blog/what-happens-in-an-internet-minute/)

# Text Collection

- Tweet Collecting
  - X API



# Create X Developer Account



<https://help.rssground.com/articles/233141-how-to-create-x-twitter-developer-app>



# Web Scraping: Extraction of data from a website

Python libraries are widely used for parsing HTML:

- 1. BeautifulSoup:** A popular library for parsing HTML and XML documents. It simplifies extracting data from web pages and has an active community with detailed documentation.
- 2. lxml:** Known for its speed, lxml is one of the fastest parsing libraries available. It receives regular updates, with the latest released in July 2023.
- 3. html5lib:** A pure-Python library designed to conform to the WHATWG (**Web Hypertext Application Technology Working Group**) HTML specification, ensuring compatibility with major web browsers.



# Demo

- Inclass code





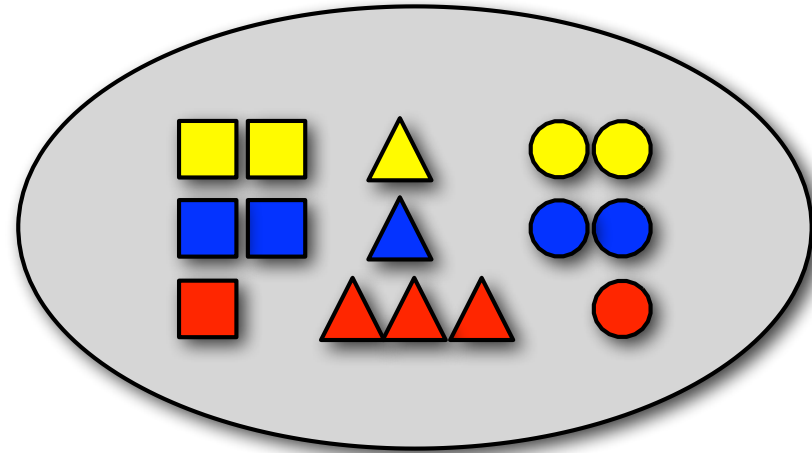
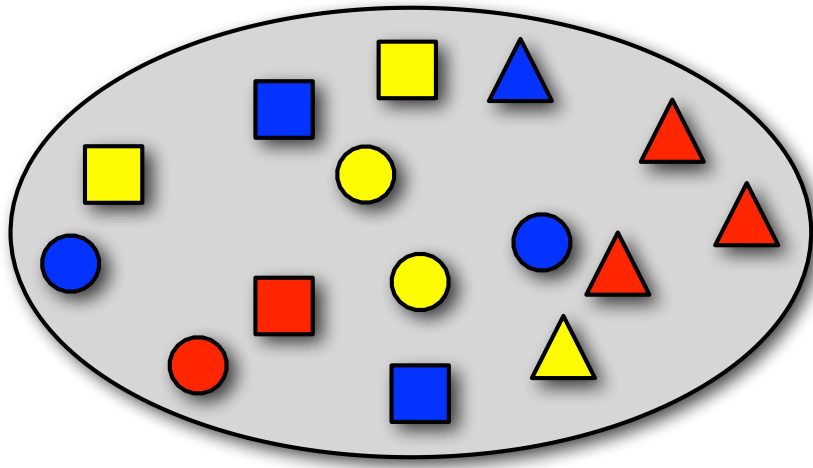
Reminder

# PROBABILITY THEORY



# Basic Probability Theory: Sampling with replacement

Pick a random shape, then put it back in the bag.



$$P(\blacksquare) = 2/15$$

$$P(\text{blue}) = 5/15$$

$$P(\text{blue} | \square) = 2/5$$

$$P(\blacksquare) = 1/15$$

$$P(\text{red}) = 5/15$$

$$P(\square) = 5/15$$

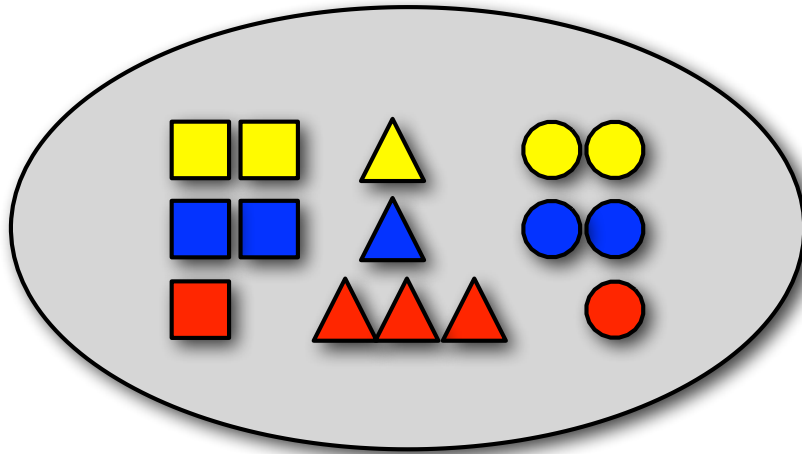
$$P(\blacksquare \text{ or } \blacktriangle) = 2/15$$

$$P(\triangle | \text{red}) = 3/5$$



# Sampling with replacement

Pick a random shape, then put it back in the bag.  
What **sequence of shapes** will you draw?



$$\begin{aligned}P(\text{red circle, yellow triangle, blue triangle, blue square}) \\&= 1/15 \times 1/15 \times 1/15 \times 2/15 \\&= 2/50625\end{aligned}$$

$$\begin{aligned}P(\text{red triangle, yellow circle, blue circle, red triangle}) \\&= 3/15 \times 2/15 \times 2/15 \times 3/15 \\&= 36/50625\end{aligned}$$

$$P(\text{blue square}) = 2/15$$

$$P(\text{blue}) = 5/15$$

$$P(\text{blue} | \text{square}) = 2/5$$

$$P(\text{red square}) = 1/15$$

$$P(\text{red}) = 5/15$$

$$P(\text{square}) = 5/15$$

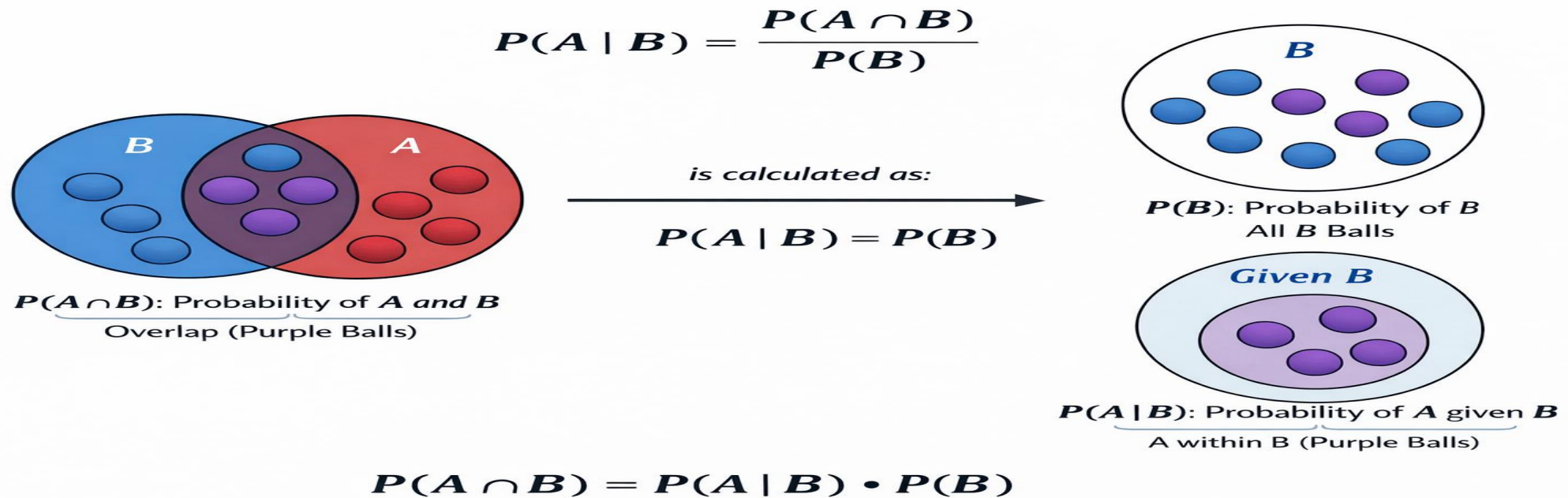
$$P(\text{red square or blue triangle}) = 2/15$$

$$P(\text{triangle} | \text{red}) = 3/5$$

# Conditional Probability

$$P(X|Y) = \frac{P(X, Y)}{P(Y)}$$

The **conditional probability of  $X$  given  $Y$** , Probability that one event occurs given that another event has already occurred.



# Chain Rule of Probability

The **chain rule** expresses a joint probability as a **product** of conditional probabilities.

For a sequence of events  $X_1, X_2, \dots, X_n$

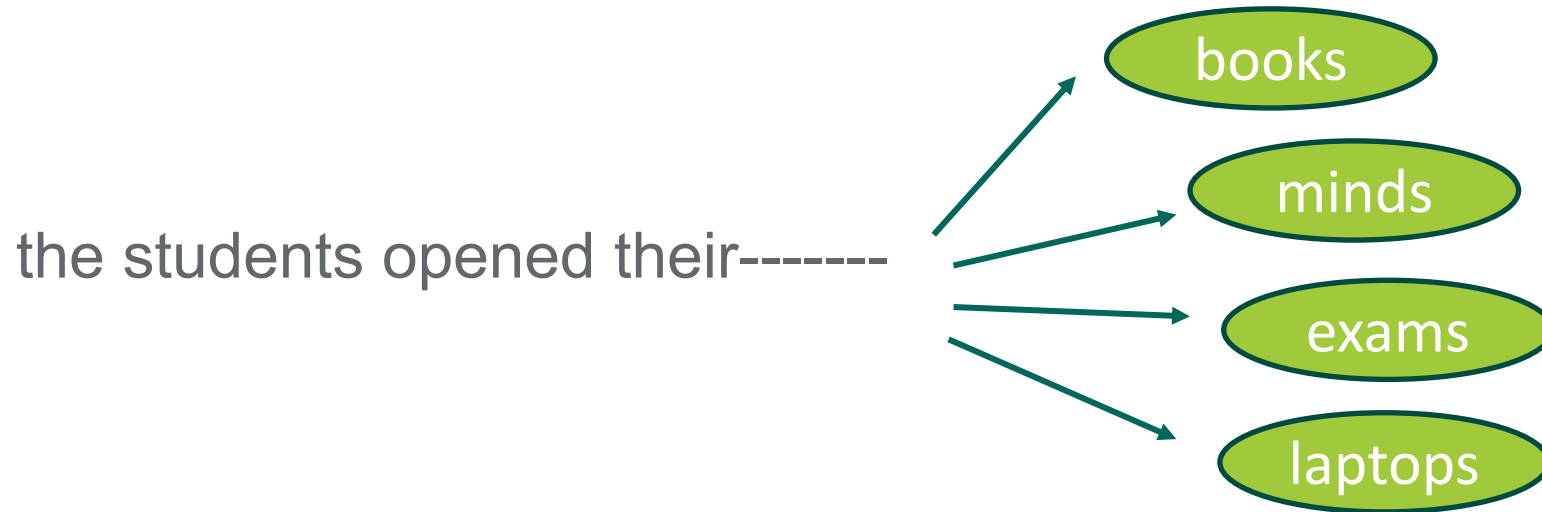
$$P(X_1, X_2, \dots, X_n) = P(X_1) P(X_2 \mid X_1) P(X_3 \mid X_1, X_2) \cdots P(X_n \mid X_1, \dots, X_{n-1})$$



# LANGUAGE MODELING



**Language Modeling:** the task of **predicting** what word comes next



Given a sequence of words  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(t)}$  compute the probability distribution of the next word  $\mathbf{x}^{(t+1)}$   
 $P(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})$

Where  $\mathbf{x}^{(t+1)}$  can be any word in the vocabulary  $V = \{w_1, \dots, w_{|V|}\}$

➤ A system that does this is called a **Language Model**.



# Popular Usages

I saw a cat

They All Saw a Cat

Book by Brendan Wenzel

i saw a cat in my dream

i saw a cat out there article

i saw a cat in my dream what does it mean

i saw a cat and it disappeared

i saw a cat in the garden

i saw a cat walking on the sidewalk

i saw a cat out there

i saw a cat in spanish

7.3K/s

Google Translate

SPANISH - detected

Si alguna vez hubo intención de representar la obra es tema de debate

ENGLISH

If there was ever an intention to represent the work, it is a matter of debate among scholars.

NEW TRANSLATION

Phonetic spelling mistakes

Fizix

 is an interesting 

sudgekt.

Physics

 is an interesting 

subject.

Felix	is an interesting	subject	0.03
Felix	is an interesting	dialect	0.005
Physics	is an interesting	dialect	0.0034
Physics	is an interesting	subject	0.01





# Goal of Language Modeling

learn patterns in text and predict the  
next word (or sequence of words)  
based on prior context.



# N-gram Language Modeling

This is Big Data AI Book

*Uni-Gram*

This	Is	Big	Data	AI	Book
------	----	-----	------	----	------

*Bi-Gram*

This is	Is Big	Big Data	Data AI	AI Book
---------	--------	----------	---------	---------

*Tri-Gram*

This is Big	Is Big Data	Big Data AI	Data AI Book
-------------	-------------	-------------	--------------

**IDEA:** Collect statistics about how frequent different n-grams are, and use these to predict next word.

# N-gram Language Modeling...

- For example, if we have sequence of tokens  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}$  then the probability to see these tokens in this order is:

## Using chain Rule

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

This is what our LM provides



# Language Modeling: n gram...

$n-1$  words

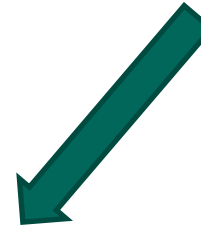
**Our assumption**  $P(x^{(t+1)} | x^{(t)}, \dots, x^{(1)}) = P(x^{(t+1)} | \overbrace{x^{(t)}, \dots, x^{(t-n+2)}}^{n-1 \text{ words}})$



Recall the definition of conditional probabilities

$$p(B|A) = P(A, B) / P(A)$$

$$P(A, B) = P(A) P(B|A)$$



$$\approx \frac{\text{count}(x^{(t+1)}, x^{(t)}, \dots, x^{(t-n+2)})}{\text{count}(x^{(t)}, \dots, x^{(t-n+2)})}$$



# n-gram Language Models: Example using 4- gram

~~as the proctor started the clock the~~ students opened their  
discard fixed window

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

For example, suppose that in the corpus:

- “students opened their” occurred 1000 times
- “students opened their books” occurred 400 times
  - $\rightarrow P(\text{books} \mid \text{students opened their}) = 0.4$
- “students opened their exams” occurred 100 times
  - $\rightarrow P(\text{exams} \mid \text{students opened their}) = 0.1$



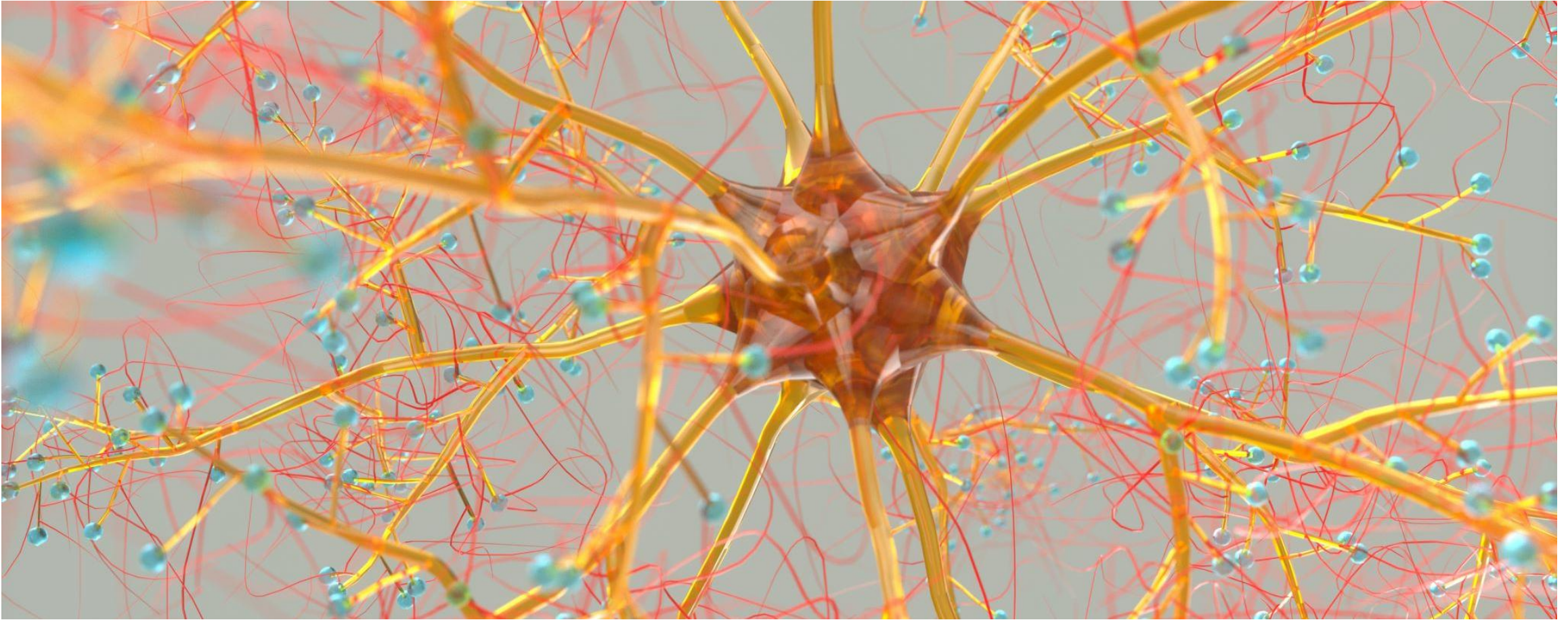
# N-grams : limitations and challenges

- Data Sparsity
- Computational Complexity
- Context Limitations

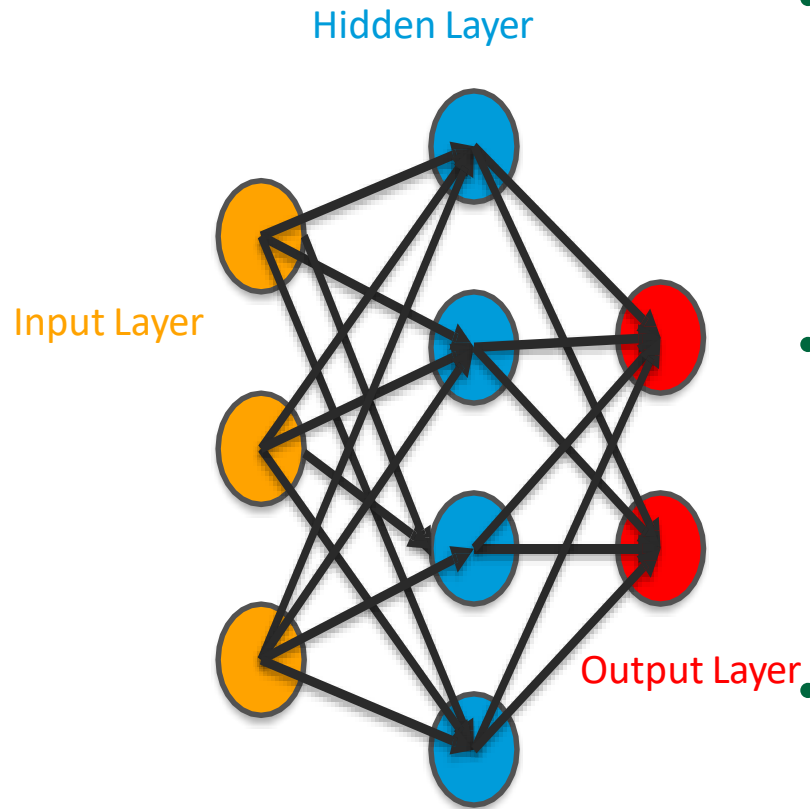




# Neural Network Based Language Models



# A Quick Review Of Neural Nets

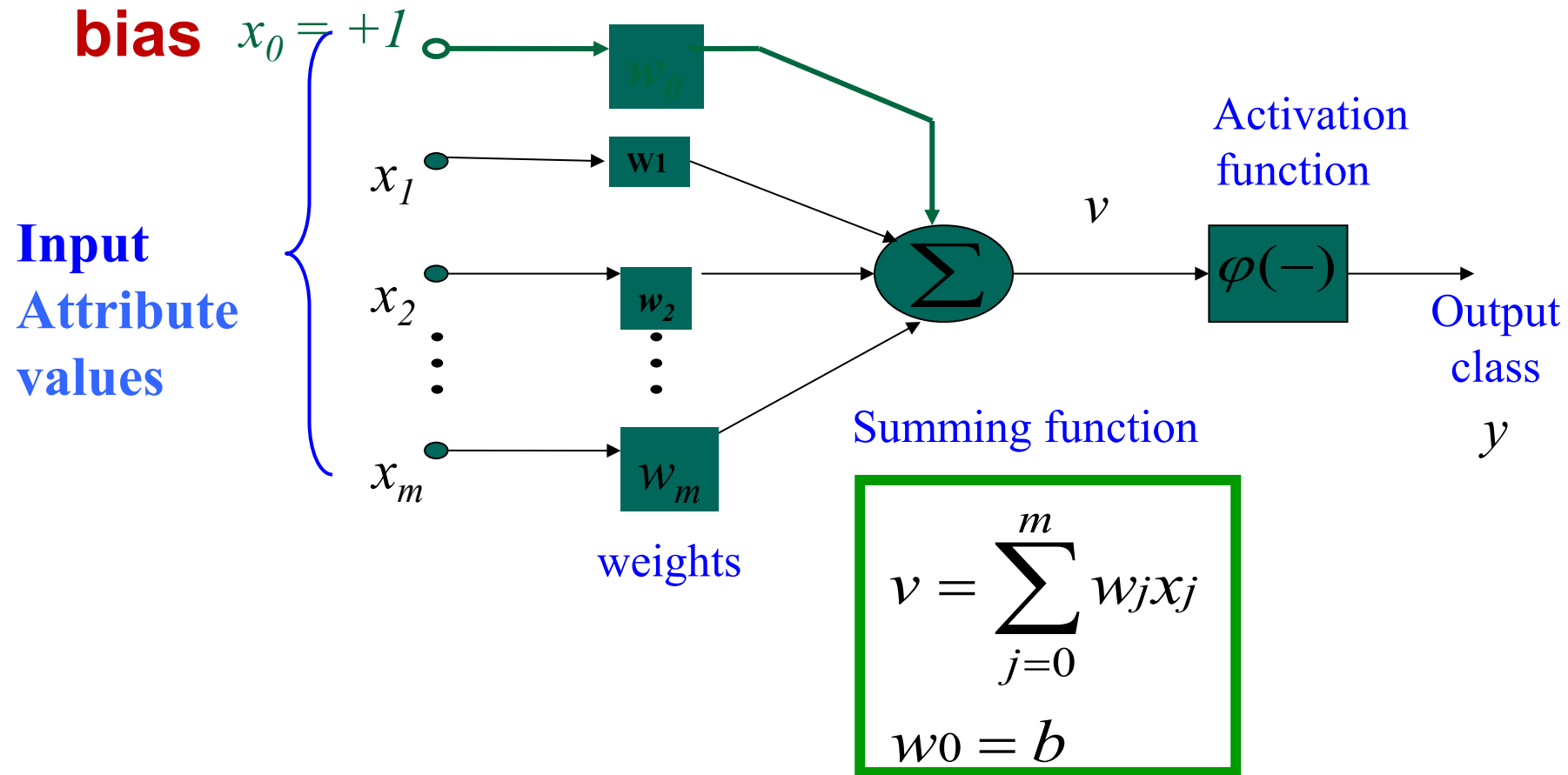


- **Input layer** is a set of features; each arrow represents a weight (float number) that tells us how much each input contributes to each following step.
- Each node in the hidden layer is some combination of all the inputs. **The hidden layer** acts as the 'input' for the output layer.
- Backpropagation allows us to adjust the weights to improve accuracy and find the 'correct' way to combine the inputs and hidden layers to get the best possible results.





# NN basic element: Perceptron or Neuron



# Language Model: Neural Nets

~~as the proctor started the clock~~      the students opened their      -----?-----

discard      fixed window

# Language Model: Neural Nets ...

output distribution

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

hidden layer

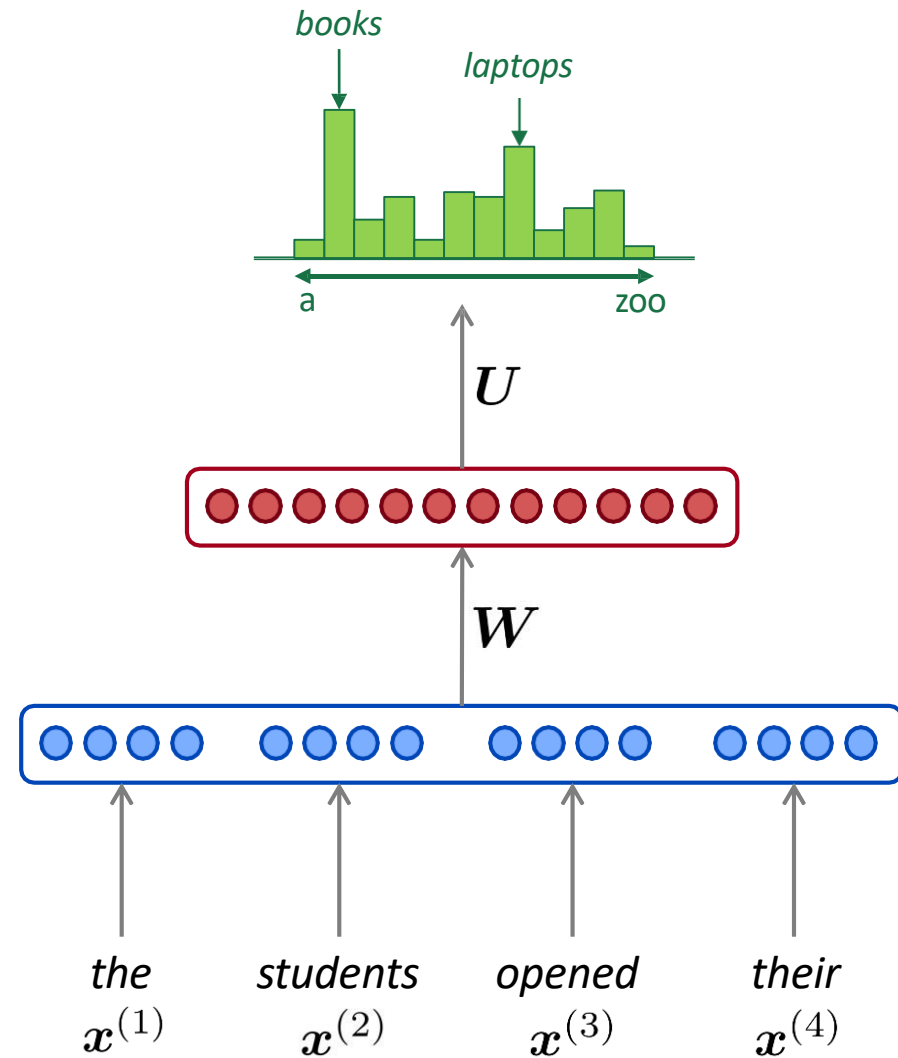
$$h = f(We + b_1)$$

concatenated word embeddings

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

words / one-hot vectors

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



# Feed Forward NN: Limitation...

~~as the proctor started the clock~~ the students opened their -----  
discard fixed window



# Feed Forward NN: Limitation

“The food was good, not bad at all”

“The food was bad, not good at all”



# Feed Forward NN: Limitation...

“Just watched the new movie. Loved it! 🍿 #entertained”

“The storyline was captivating, the characters were well-developed, and the cinematography was impressive. Overall, a fantastic movie night! 🎬 👍 #movienight #recommend”

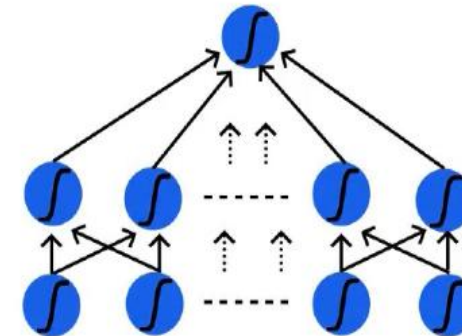
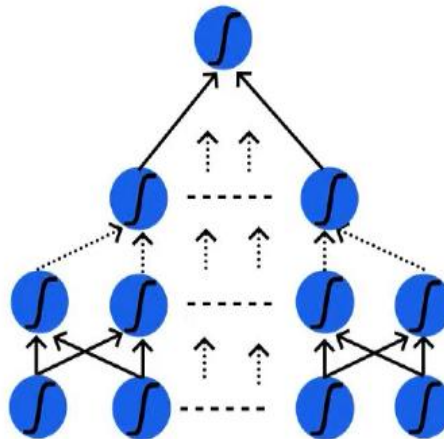
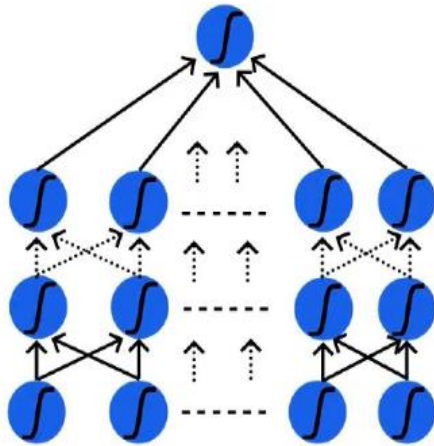
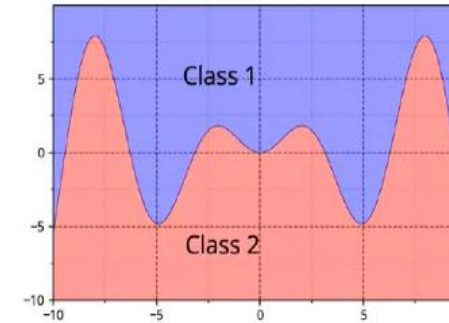
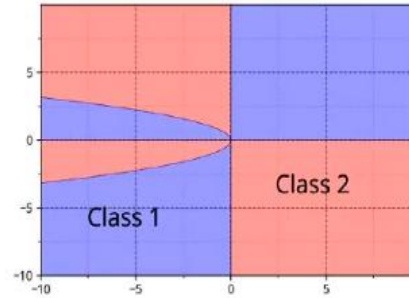
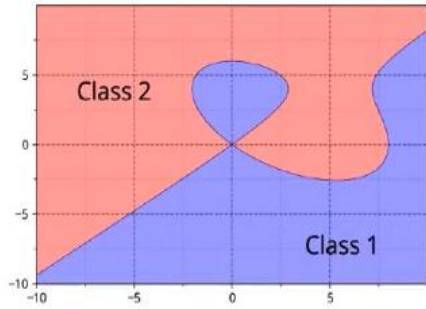


# Sequence Modeling: Motivations

- Handle **variable length** sequence data
- Track long **term dependency**
- Maintain information about **order**
- **Share information** across the sequence



# DNN: Universal Approximation Theorem (UAT)



proven by George Cybenko in 1989

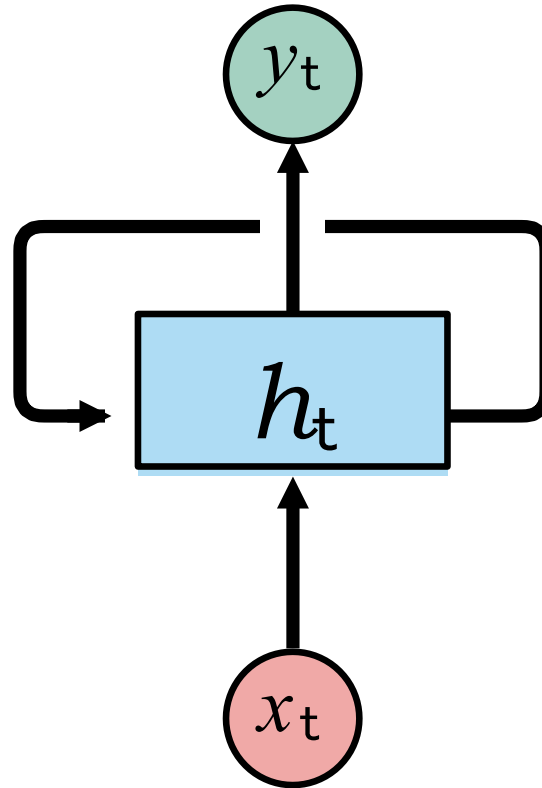




# Core idea of Recurrent Neural Networks (RNNs) RNNs

Stateful computation

---

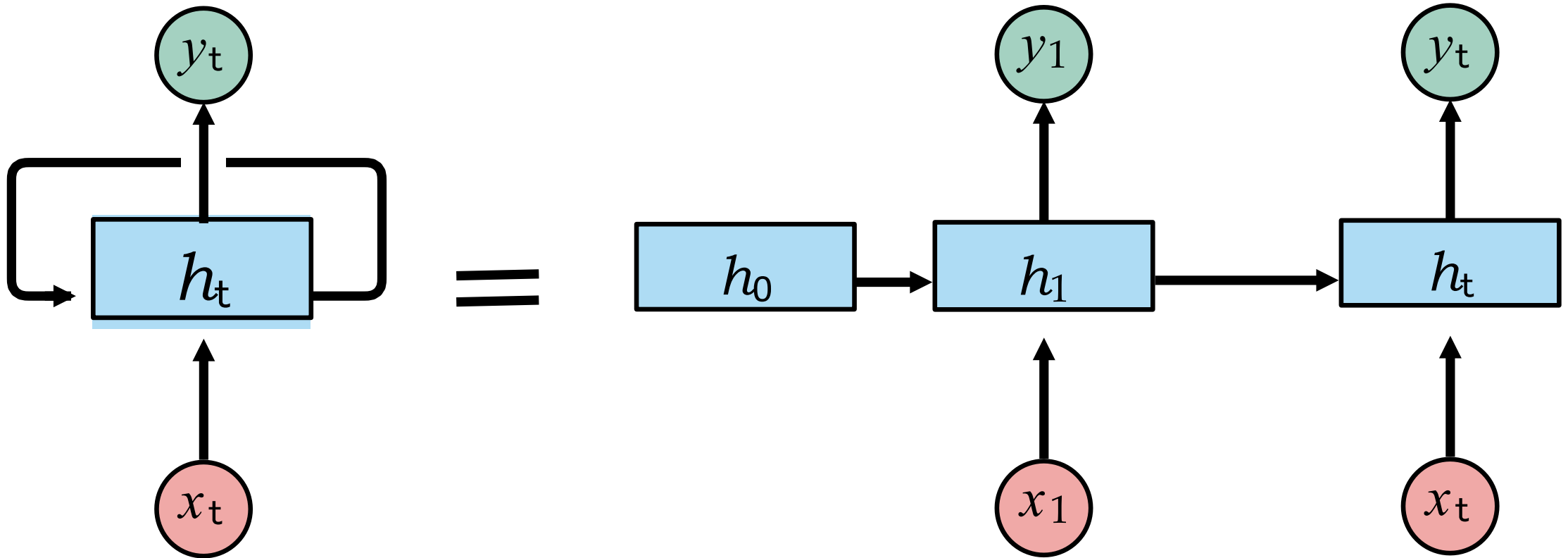


$$y_t, h_t = f(x_t, h_{t-1})$$



# Core idea of RNNs ...

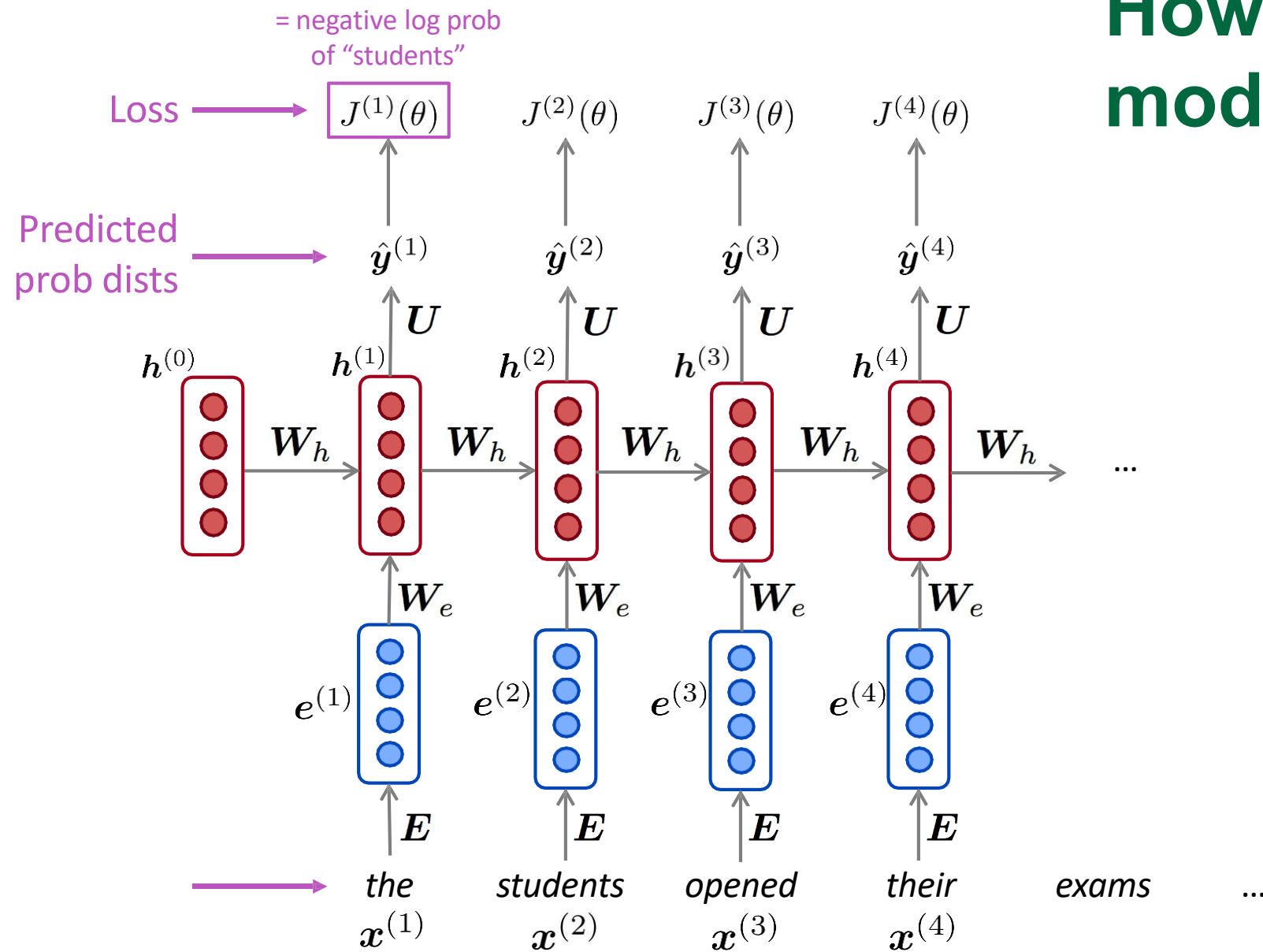
## Stateful computation



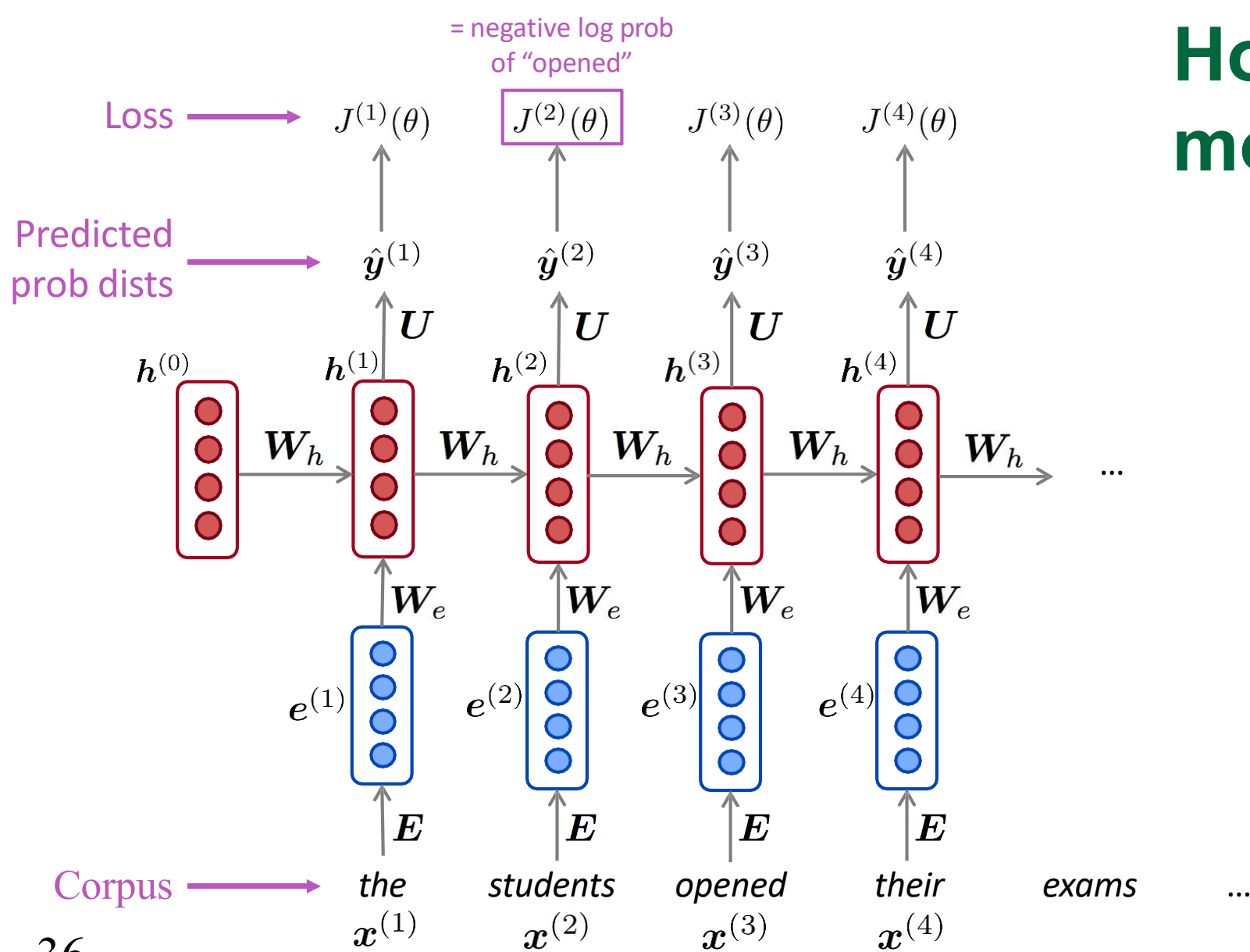
$$h_t = W_x x_t + W_h h_{t-1} + b$$



# How we train the model

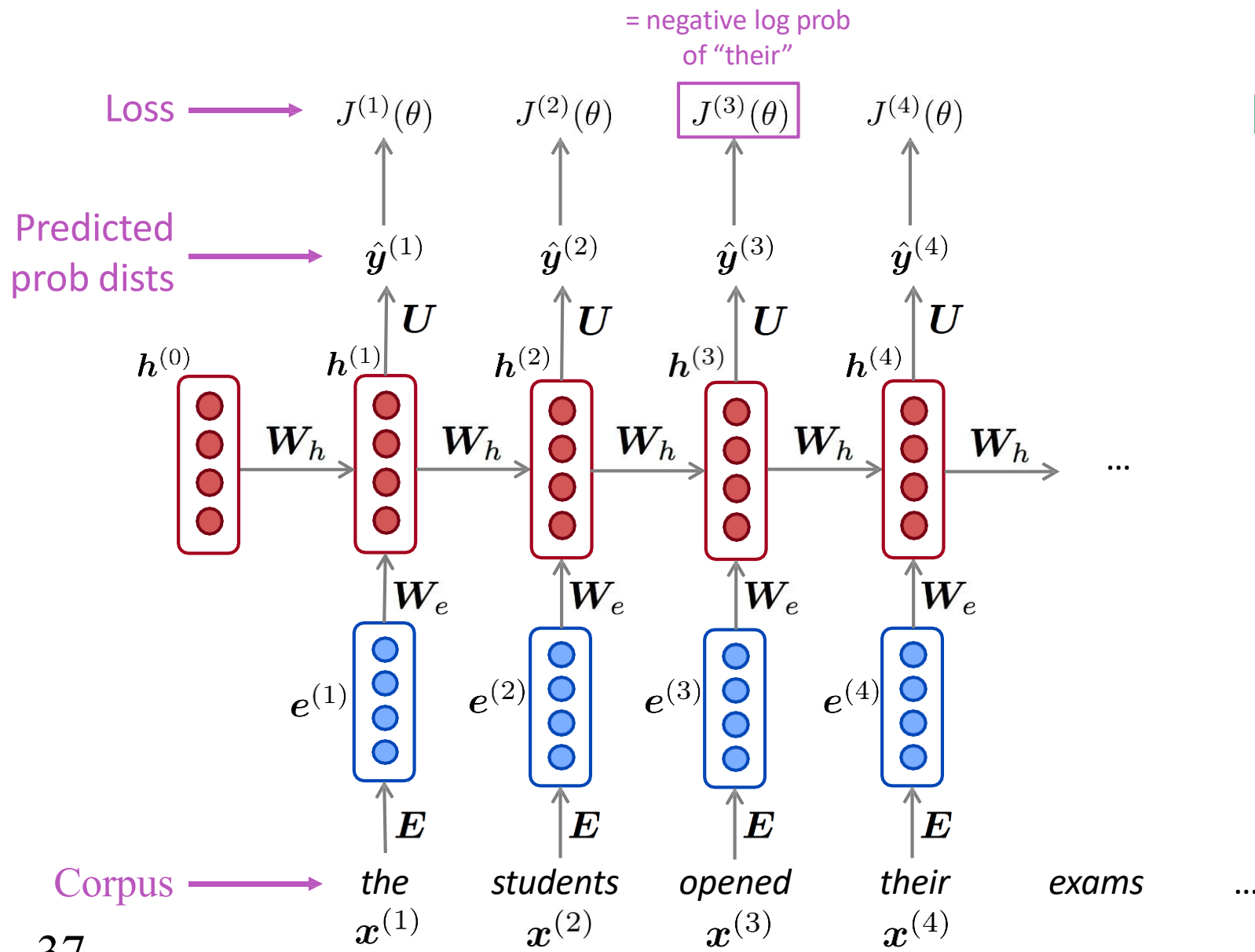


# How we train the model



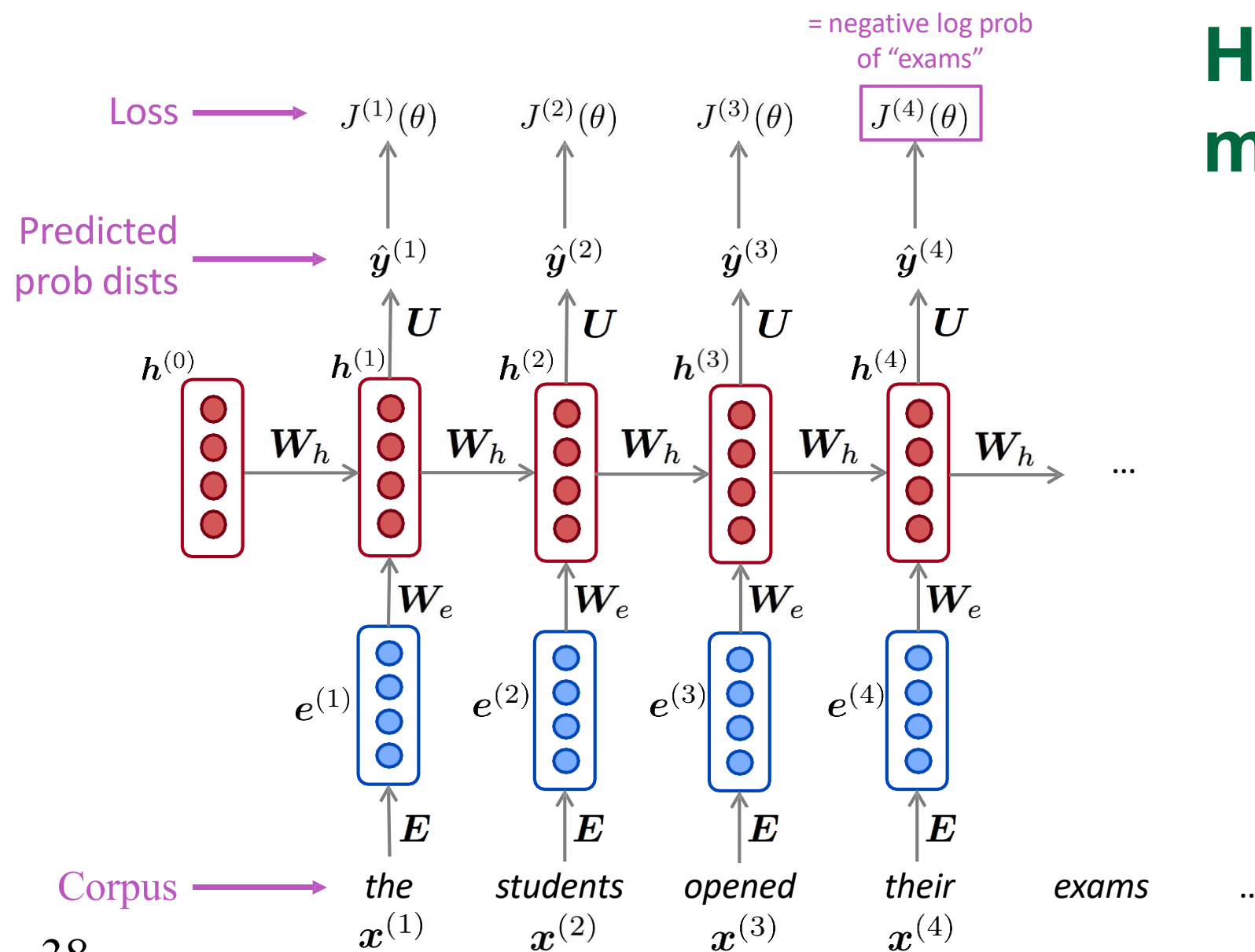
36

# How we train the model



37

# How we train the model

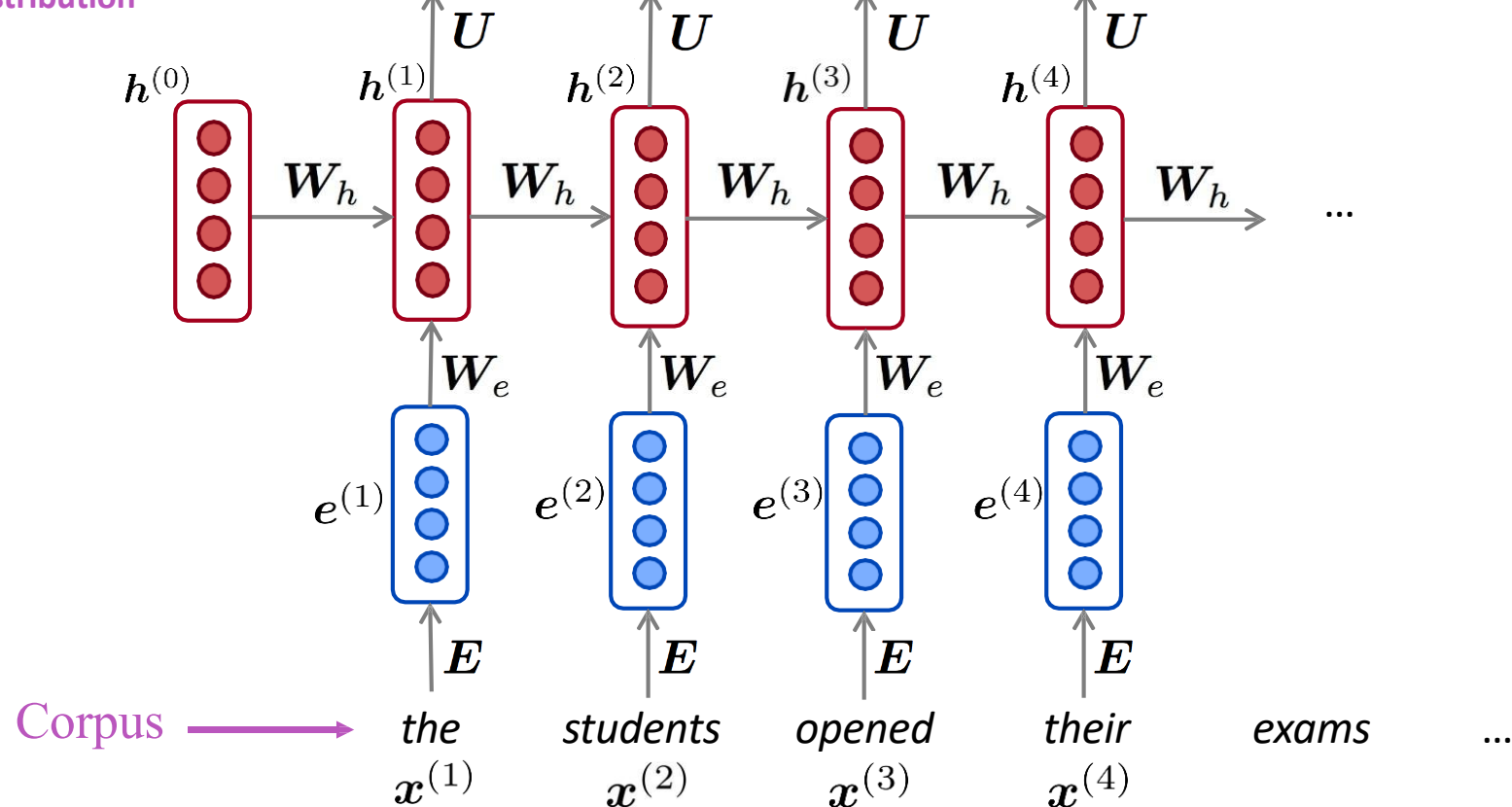


38

Loss  $\longrightarrow J^{(1)}(\theta) + J^{(2)}(\theta) + J^{(3)}(\theta) + J^{(4)}(\theta) + \dots = J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta)$

Predicted probability distribution  $\longrightarrow$

## How we train the model



# Language Model: RNN

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

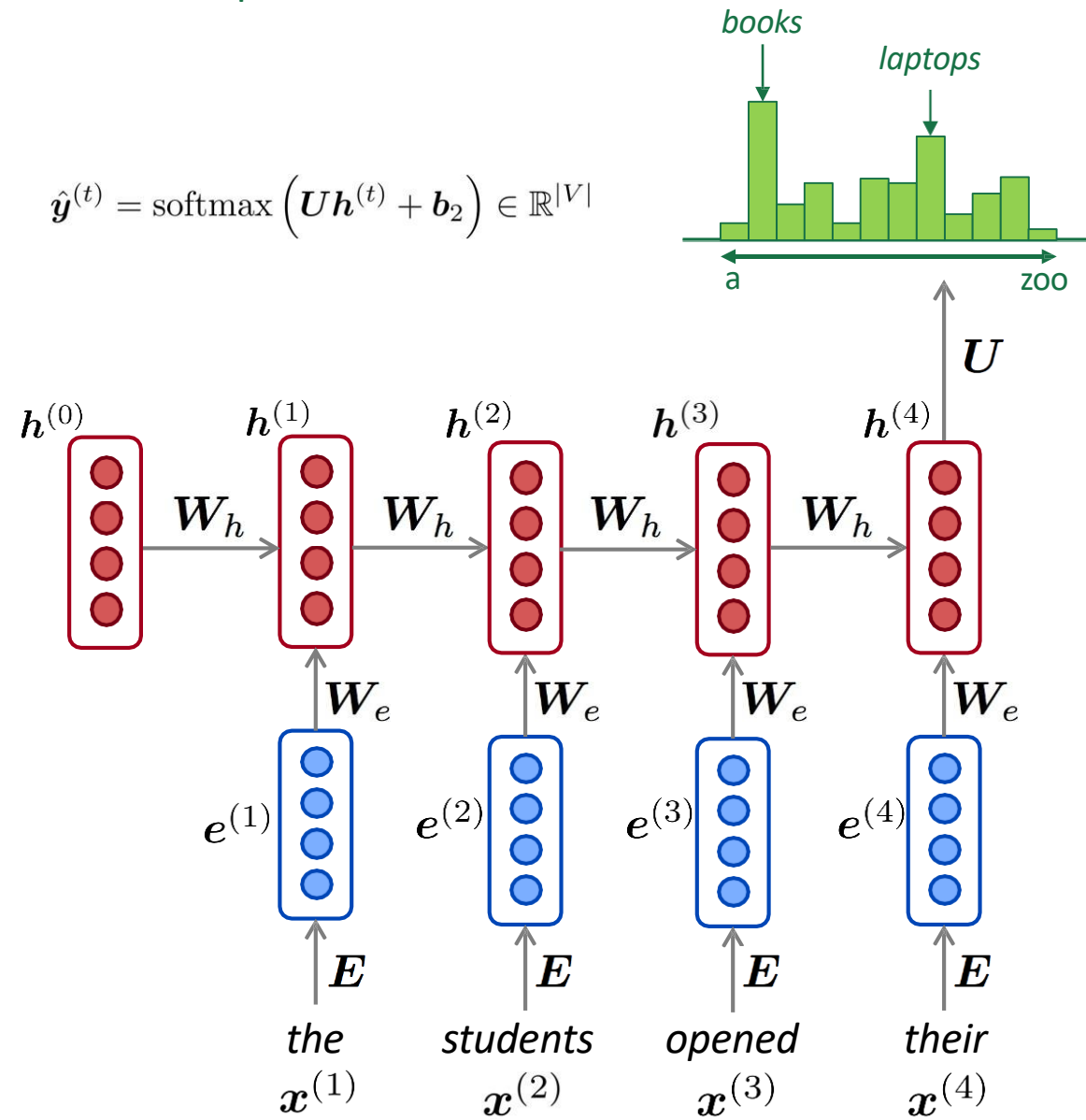
word embeddings

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

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

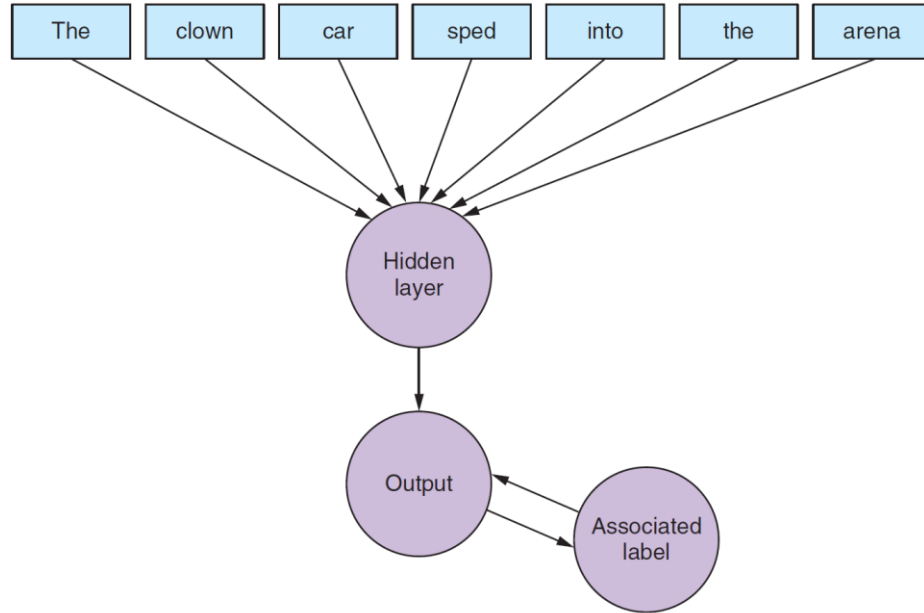
output distribution  $\hat{y}^{(4)} = P(x^{(5)} | \text{the students opened their})$

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

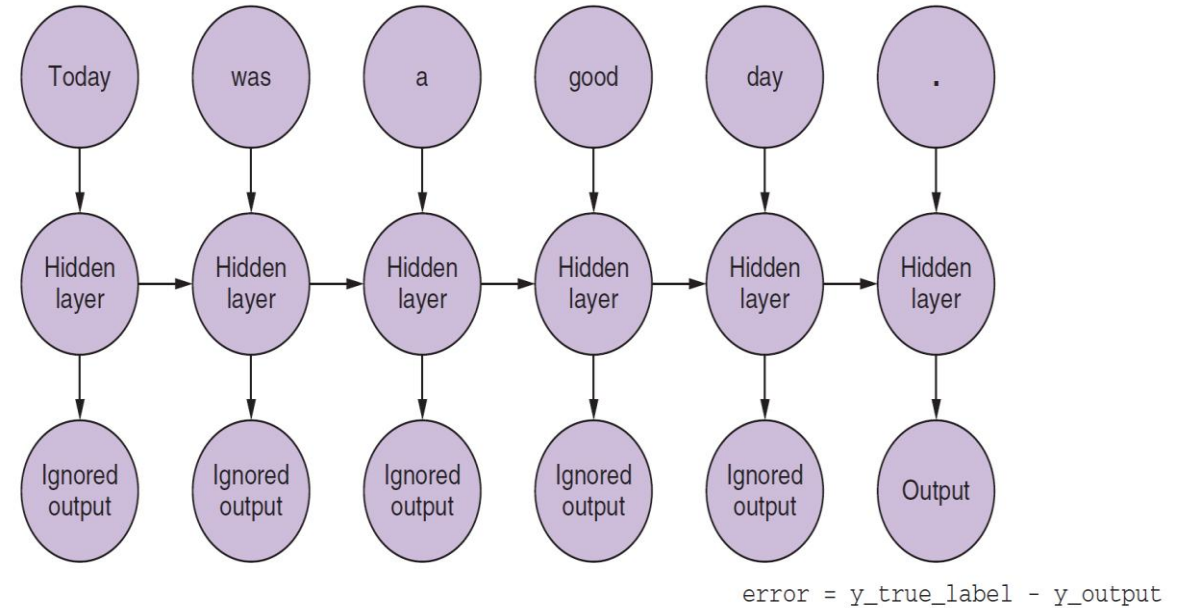




# Difference between NN and RNN



Traditional NN for LM



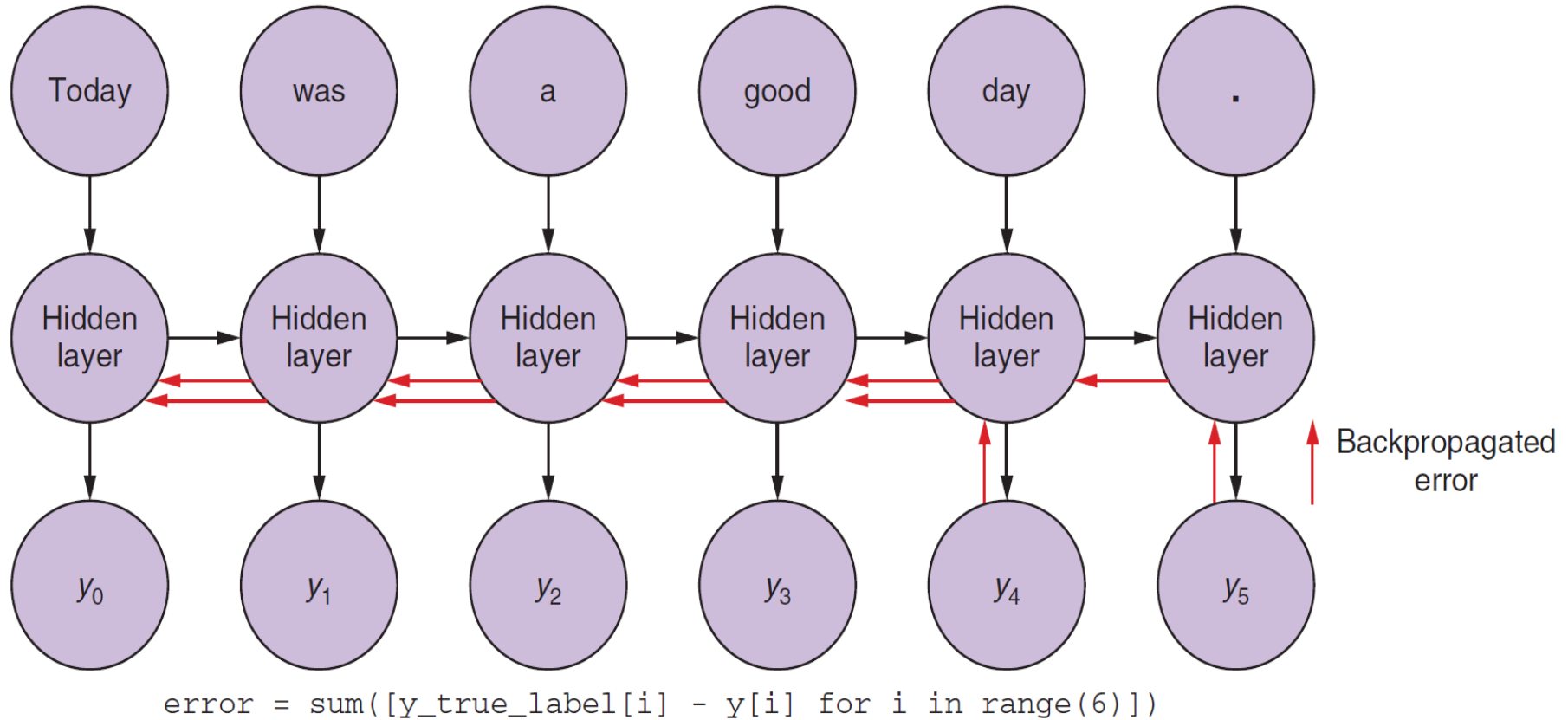
RNN for LM

# Fun With RNN Language Model

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



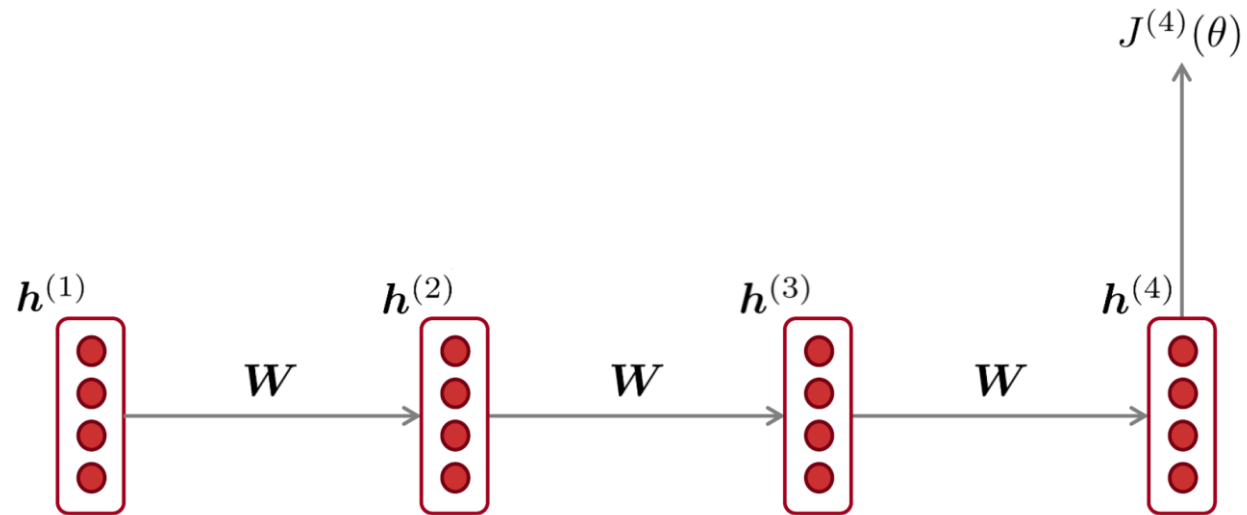
# Back Propagation in RNN



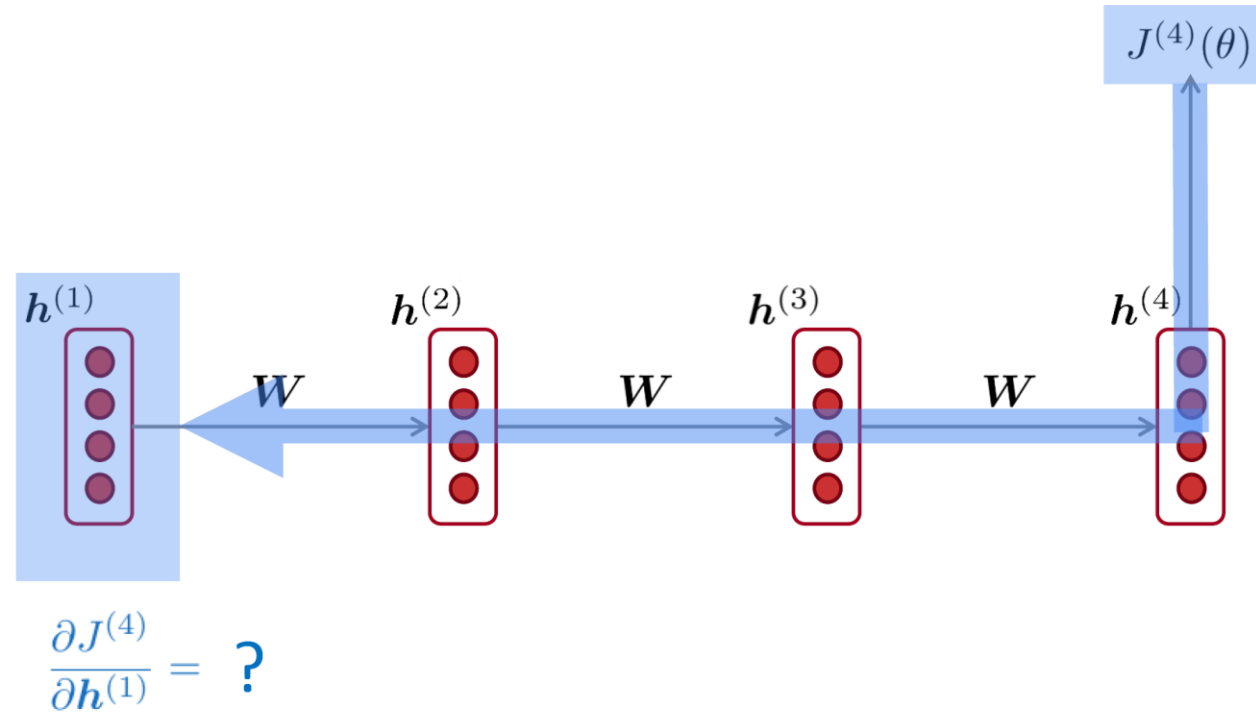
**Backpropagation Through Time (BPTT).**



# RNN Vanishing Gradient Intuition



# Vanishing gradient intuition

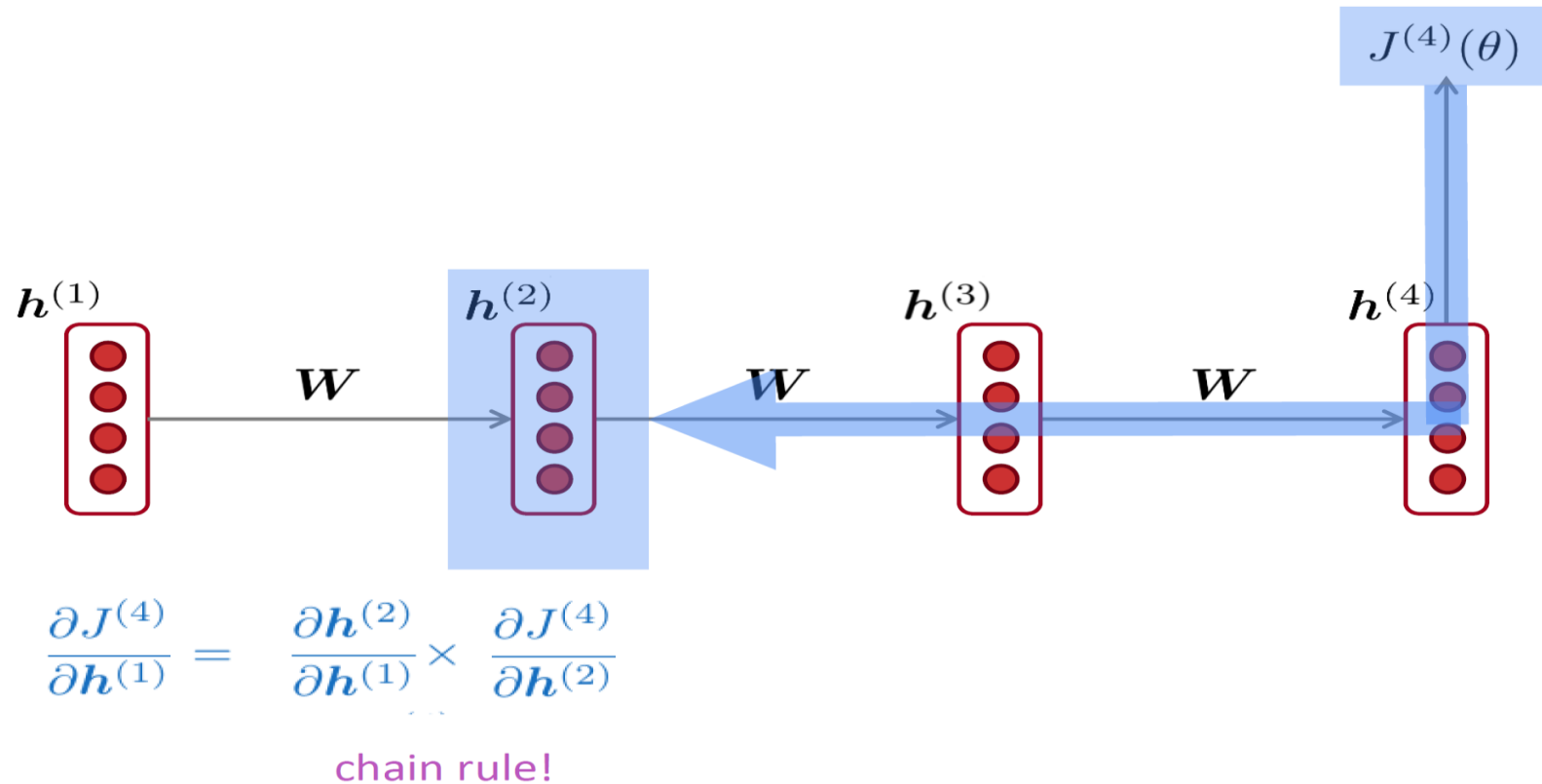


# Vanishing gradient intuition

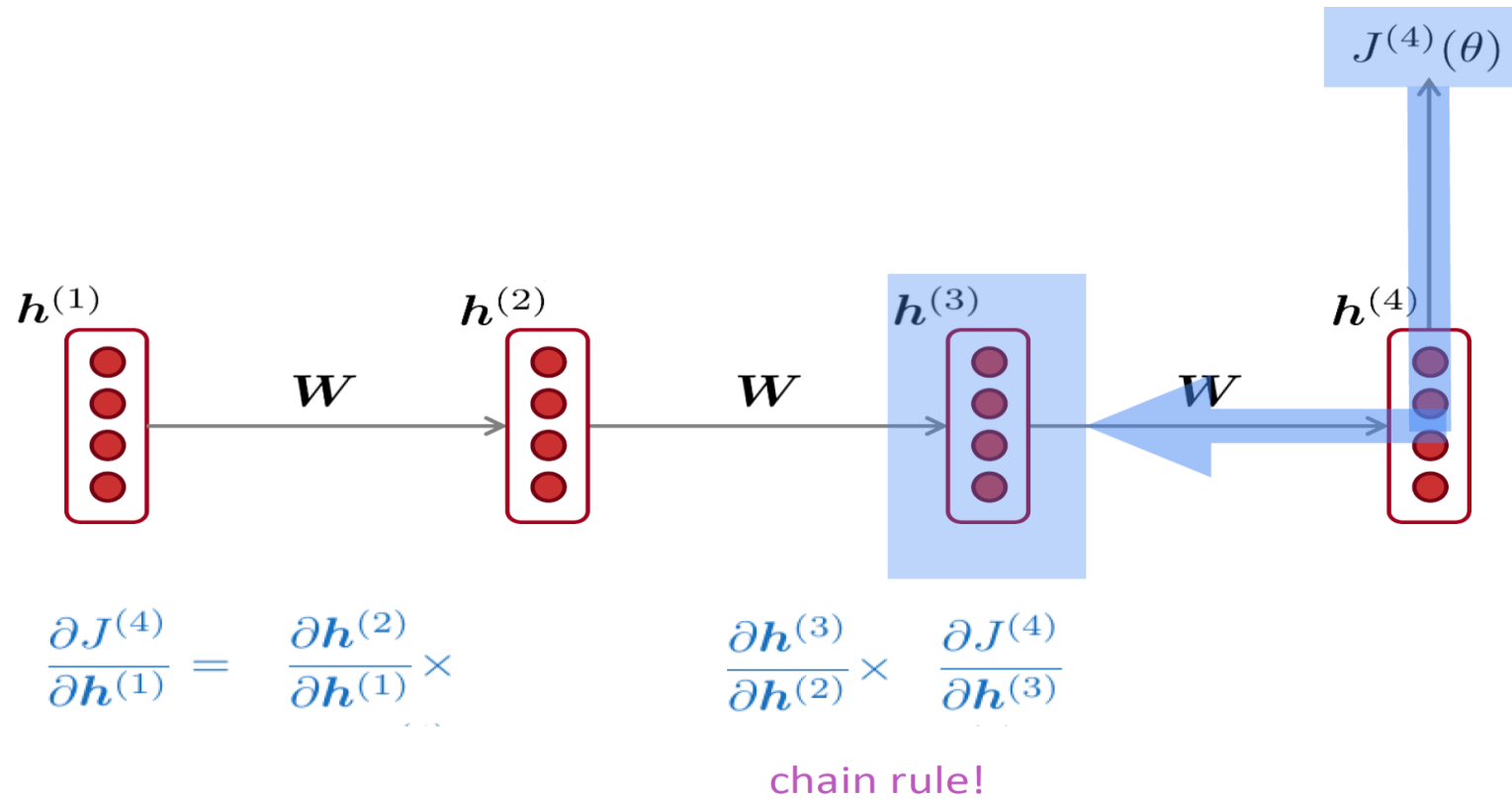
The Chain Rule

If  $y = f(u)$ , where  $u = g(x)$

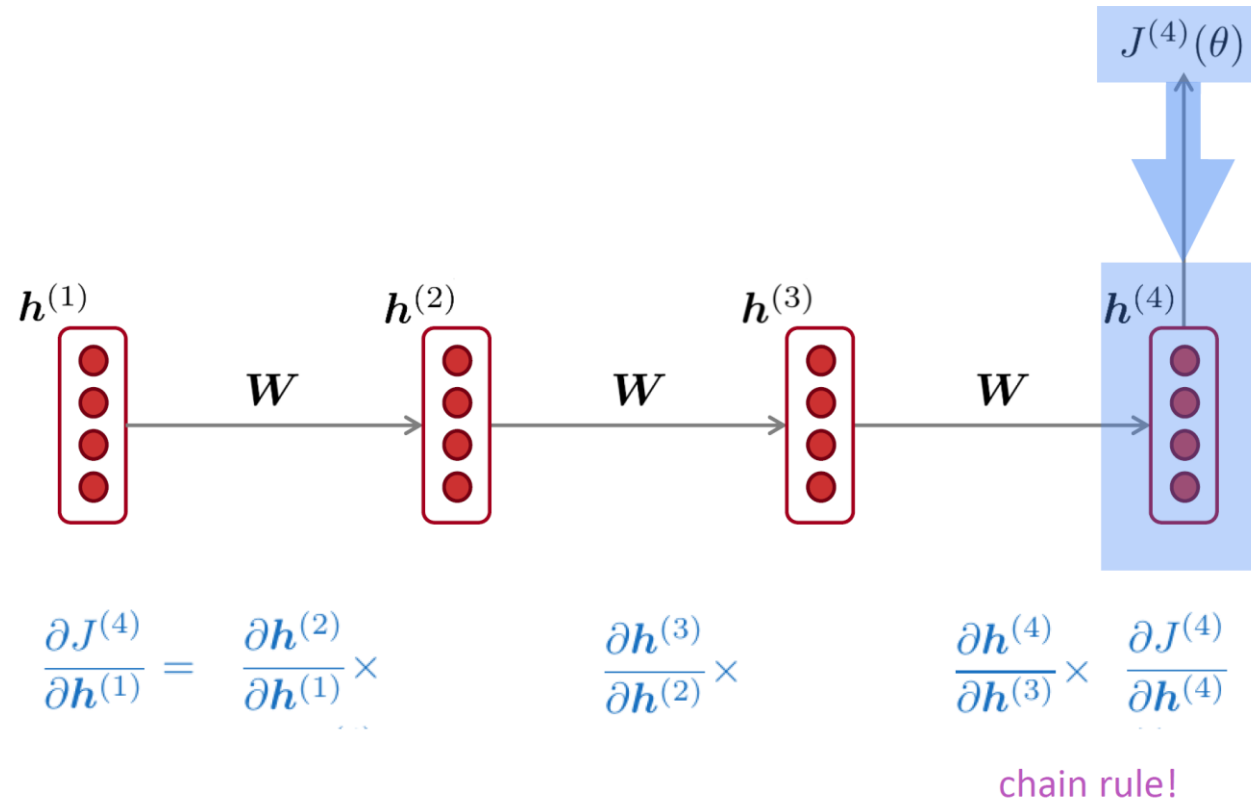
$$\frac{dy}{dx} = \frac{dy}{du} \cdot \frac{du}{dx}$$



# Vanishing gradient intuition

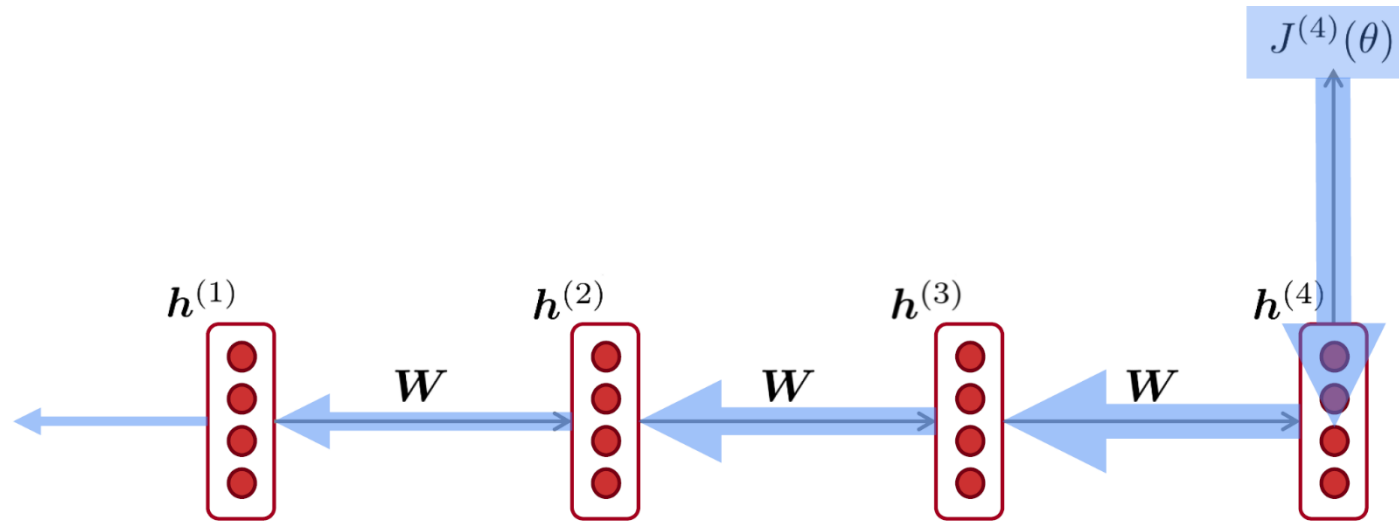


# Vanishing gradient intuition





# Vanishing gradient intuition



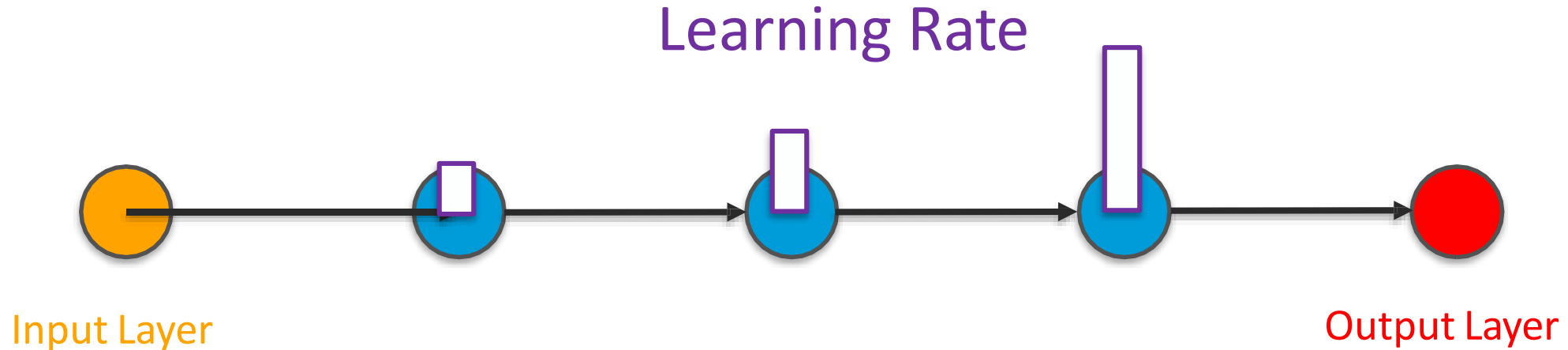
$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = \boxed{\frac{\partial h^{(2)}}{\partial h^{(1)}}} \times \boxed{\frac{\partial h^{(3)}}{\partial h^{(2)}}} \times \boxed{\frac{\partial h^{(4)}}{\partial h^{(3)}}} \times \frac{\partial J^{(4)}}{\partial h^{(4)}}$$

What happens if these are small?

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

# Why Vanishing Gradients is Problem

Vanishing gradients occur when the values of a gradient are too small and the model stops learning or takes way too long as a result



# Vanishing Gradients Problem...

## *Example*

*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-----*

RNN-LM needs to **model the dependency** between “tickets” on the 7<sup>th</sup> step and the target word “tickets” at the end



# Long Short-Term Memory (LSTM)

- **Hochreiter & Schmidhuber (1997)** solved the problem of getting an RNN to remember things for a long time.
  - At each timestep  $t$ , the LSTM maintains two key components:
    - **Hidden state** – captures short-term dependencies.
    - **Cell state** – acts as a memory unit, storing long-term information.



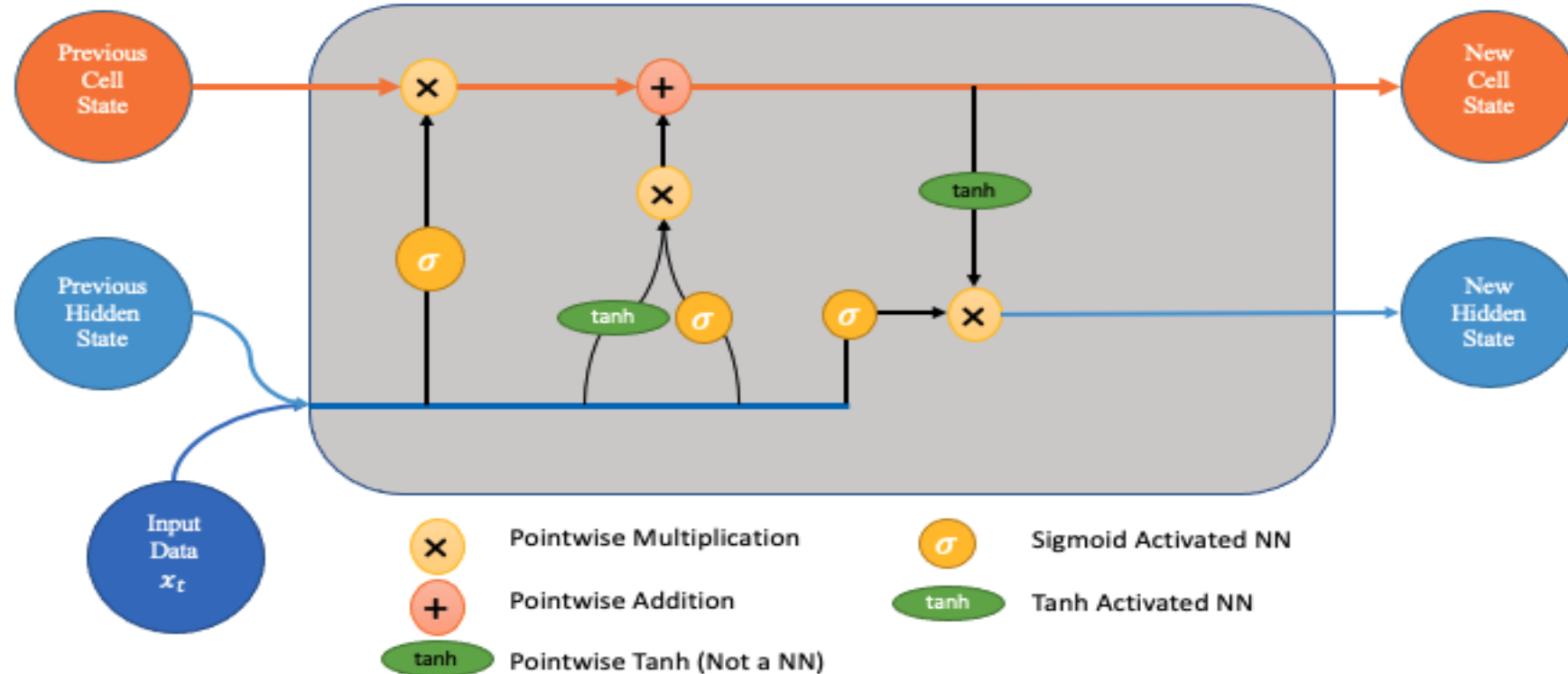
# Long Short-Term Memory (LSTM)

## Key Concepts:

- Unlike standard RNNs, LSTMs can **control the flow of information** through three specialized gates:
  - **Forget gate** – decides which information to erase.
  - **Input gate** – determines what new information should be stored.
  - **Output gate** – regulates what information is passed to the next timestep.
- Each gate is represented as a vector of size  $n$  and can take values between **0** (**closed**) and **1** (**open**) dynamically, based on the current context.



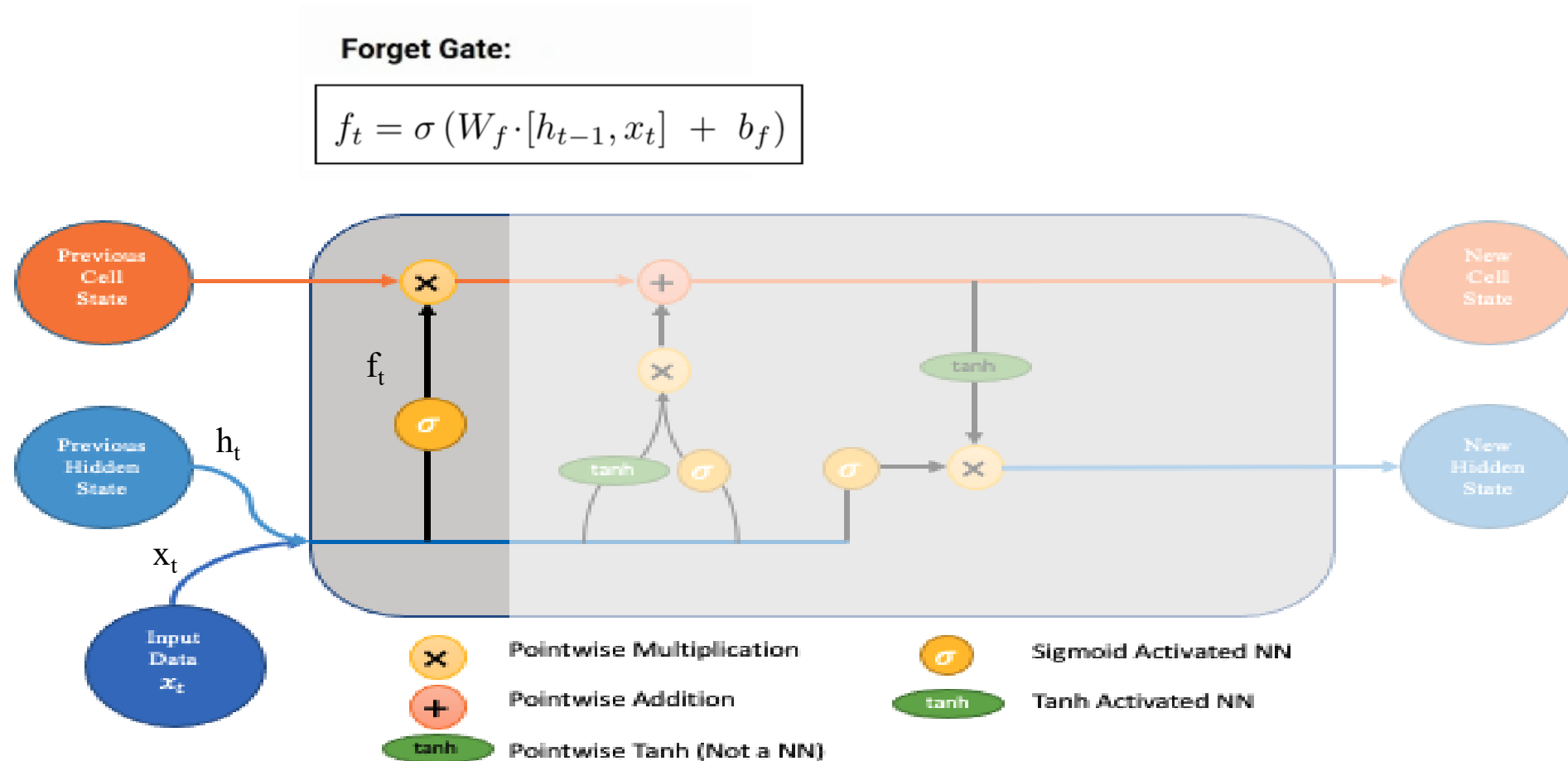
# Long Short -Term Memory (LSTM)



LSTM at time stamp T

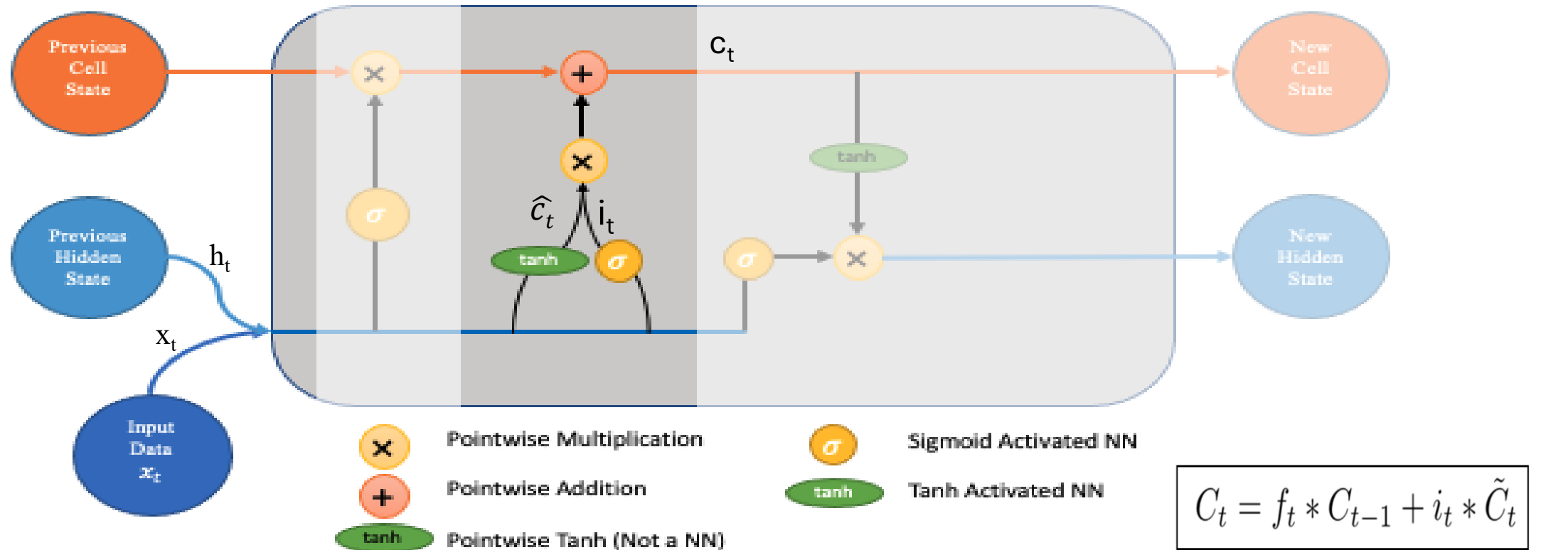
Image source: <https://towardsdatascience.com/lstm-networks-a-detailed-explanation-8fae6aefc7f9>

# Long Short Term Memory (LSTM): step 1



**Forget gate: decide what parts of old state to forget**

# Long Short Term Memory (LSTM):step2



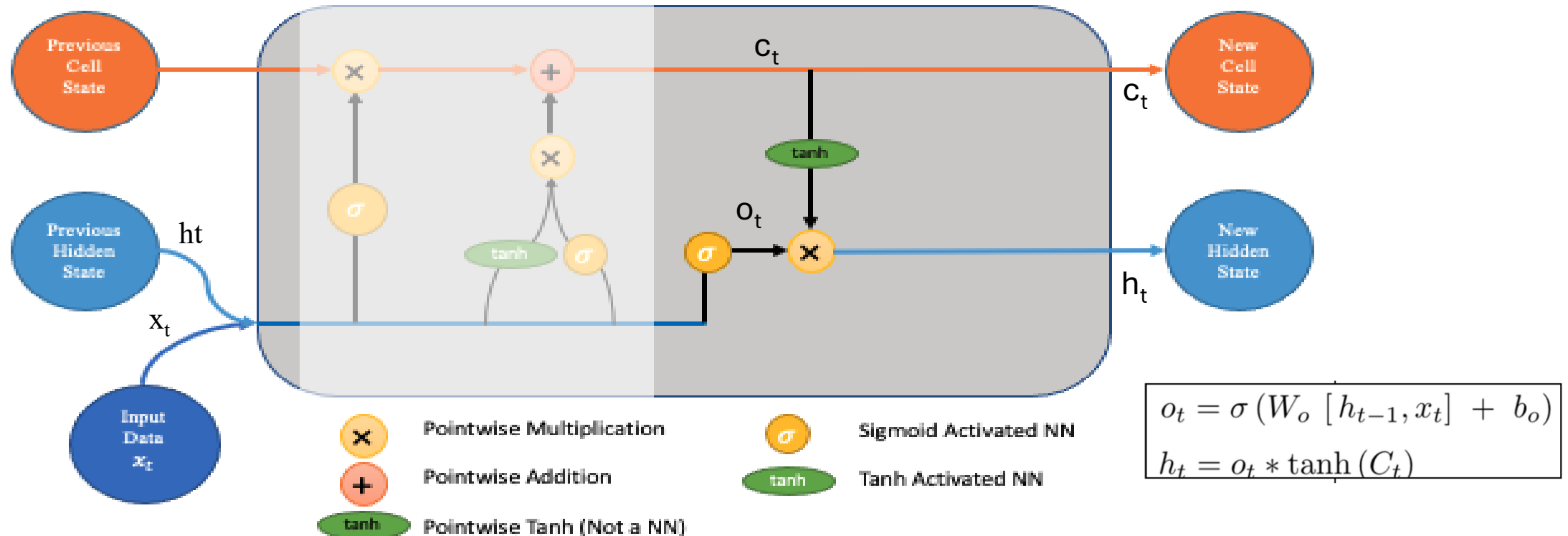
Input gate: decide how to update the cell state



# Long Short Term Memory (LSTM):step3

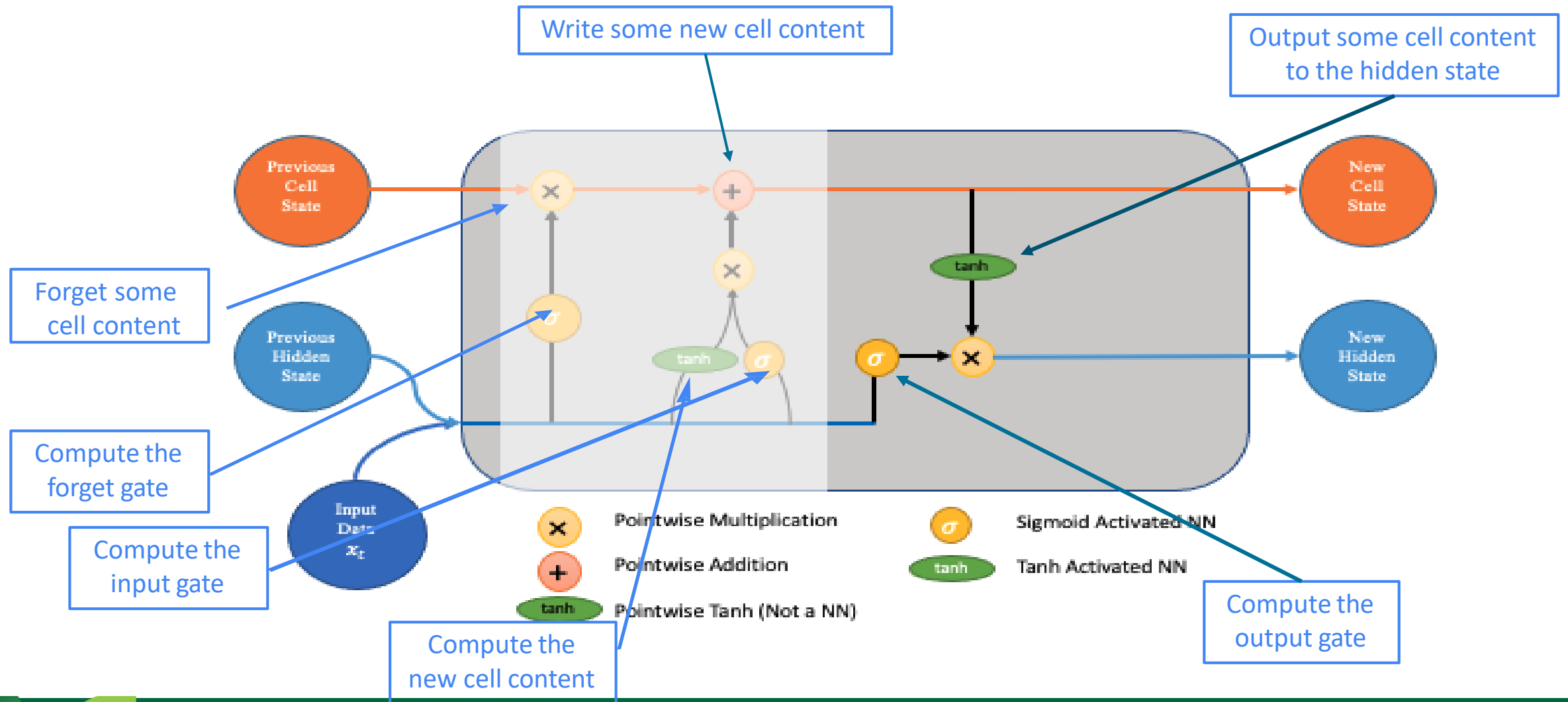
Output Gate:

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$



Finally, decide what to output as hidden state

# Long Short Term Memory (LSTM)



# LSTM Great resources

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



# Keras – Simplifying LSTMs in Python

Keras is a Python package that makes building and training TensorFlow neural networks really simple. We'll be working with the "Sequential" model which lets you add layers one at a time. As an example, let's see how to build a 1-layer LSTM model with 10 hidden nodes.

```
from keras.models import Sequential
from keras.layers import Dense, Activation, LSTM

model = Sequential()
model.add(LSTM(10, input_shape=(TIMESTEPS, FEATURE_LENGTH)))
model.add(Dense(NUMBER_OF_OUTPUT_NODES))
model.add(Activation('softmax'))
```



# 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

**Low perplexity** → the model predicts the text well

**High perplexity** → the text is unexpected for the model

**Perplexity (PPL)** measures **how confused** a language model is when predicting the next word in a sentence.



# Summary

- We introduced the concepts of recurrent neural networks and how it can be applied to language problems.
- RNNs can be trained with a straightforward extension of the backpropagation algorithm.
- How LSTM used for text generation
- Applications of LSTM for sequence-to-sequence modeling



# Q&A

