# ITCS 4101: Introduction to Natural Language Processing

N-gram Language Models

Neural Language Models

Recurrent Neural Networks (RNNs)

Razvan C. Bunescu
Comptuer Science @ CCI
University of North Carolina at Charlotte
https://webpages.charlotte.edu/rbunescu
*razvan.bunescu@charlotte.edu*

# Language Modeling (LM)

- **Causal** Language Modeling:
  - **Predict the next word in a sequence**:
    - AI systems use machine _____

      eat?

      learning?

      frogs?

      …
    - **The LM estimates P(word | word$_{-1}$, word$_{-2}$, ...)**
      - we want P(learning | machine, use) >> P(about | machine, eat).

  - **Decoder** neural architectures are widely used to train LMs:
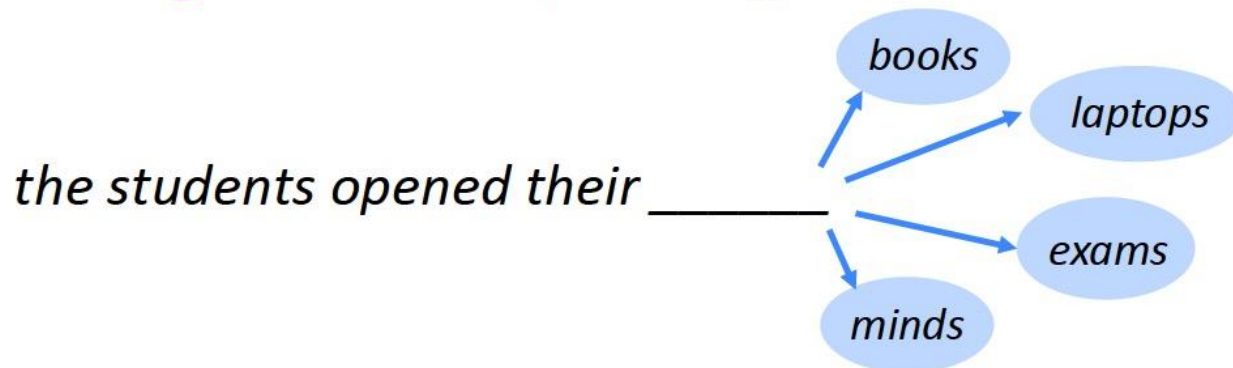    - GPT, Gemini, Llama, Mixtral, Claude, …

# Language Modeling (LM)

- **Masked** Language Modeling:
  - **Predict the most likely word in a context**:
    - AI systems use machine _____ models for language understanding .
      - eat?
      - learning?
      - frogs?
      - …
    - The LM estimates P(word | $word_{-1}$, $word_{-2}$, ...; $word_1$, $word_2$, ...)
      - we want P(learning | machine, use; models, for )
        - $>>$ P(frogs | machine, use; models, for).

  - **Encoder** neural architectures are used to train masked LMs.
    - BERT, RoBERTa, …

# Language Modeling

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

books

laptops

*the students opened their* _____

exams

minds

- More formally: given a sequence of words $x^{(1)}, x^{(2)}, \ldots, x^{(t)}$, compute the probability distribution of the next word $x^{(t+1)}$ :

$$P(x^{(t+1)} \mid x^{(t)}, \ldots, x^{(1)})$$

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

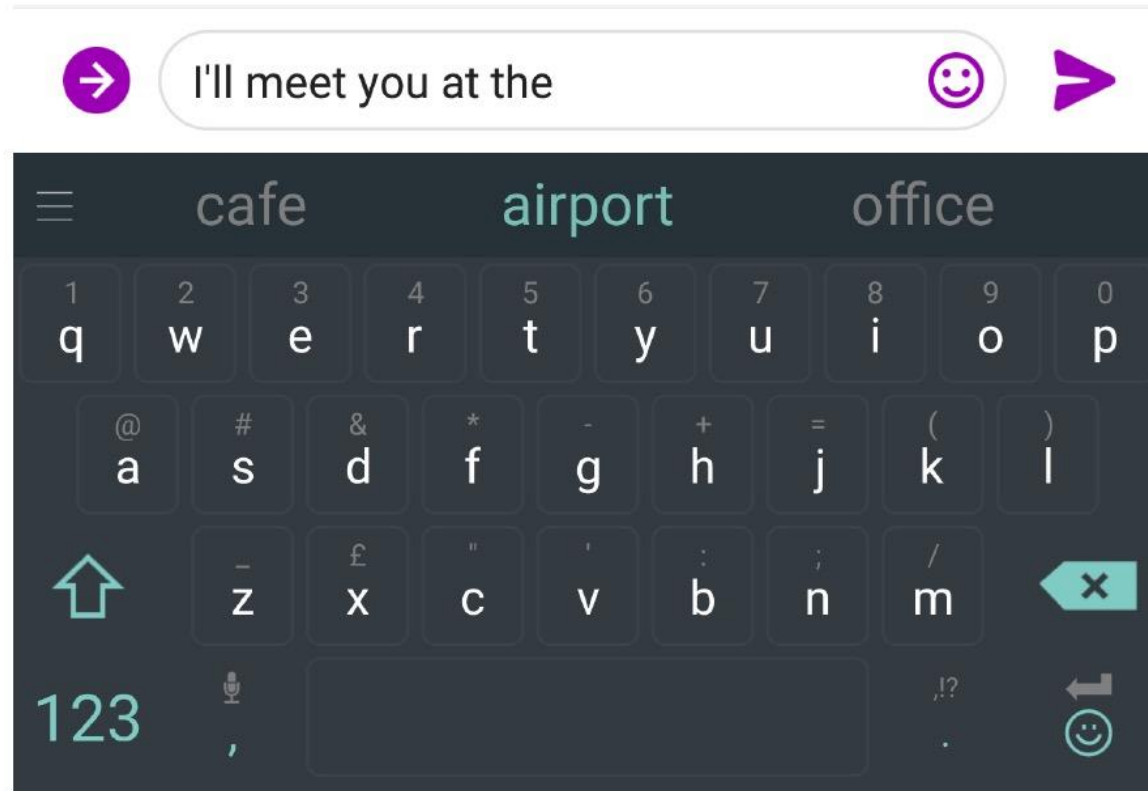- A system that does this is called a **Language Model**.

# Language Modeling

- You can also think of a Language Model as a system that assigns probability to a piece of text.

- For example, if we have some text $\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(T)}$, then the probability of this text (according to the Language Model) is:

$$P(\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(T)}) = P(\boldsymbol{x}^{(1)}) \times P(\boldsymbol{x}^{(2)} \mid \boldsymbol{x}^{(1)}) \times \cdots \times P(\boldsymbol{x}^{(T)} \mid \boldsymbol{x}^{(T-1)}, \ldots, \boldsymbol{x}^{(1)})$$

$$= \prod_{t=1}^{T} P(\boldsymbol{x}^{(t)} \mid \boldsymbol{x}^{(t-1)}, \ldots, \boldsymbol{x}^{(1)})$$

This is what our LM provides

# You use Language Models every day!

# You use Language Models every day!

# n-gram Language Models

*the students opened their* _____

- **Question**: How to learn a Language Model?
- **Answer** (pre- Deep Learning): learn a *n*-gram Language Model!

- Definition: A *n*-gram is a chunk of *n* consecutive words.
  - unigrams: "the", "students", "opened", "their"
  - bigrams: "the students", "students opened", "opened their"
  - trigrams: "the students opened", "students opened their"
  - 4-grams: "the students opened their"

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

# n-gram Language Models

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

*n*-1 words

$$P(x^{(t+1)}|x^{(t)}, \ldots, x^{(1)}) = P(x^{(t+1)}|x^{(t)}, \ldots, x^{(t-n+2)}) \qquad \text{(assumption)}$$

prob of a n-gram

prob of a (n-1)-gram

$$= \frac{P(x^{(t+1)}, x^{(t)}, \ldots, x^{(t-n+2)})}{P(x^{(t)}, \ldots, x^{(t-n+2)})} \qquad \begin{array}{l}\text{(definition of}\\ \text{conditional prob)}\end{array}$$

- **Question:** How do we get these *n*-gram and (*n*-1)-gram probabilities?

- **Answer:** By counting them in some large corpus of text!

$$\approx \frac{\text{count}(x^{(t+1)}, x^{(t)}, \ldots, x^{(t-n+2)})}{\text{count}(x^{(t)}, \ldots, x^{(t-n+2)})} \qquad \begin{array}{l}\text{(statistical}\\ \text{approximation)}\end{array}$$

8

# n-gram Language Models: Example

Suppose we are learning a 4-gram Language Model.

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

discard

condition on this

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

For example, suppose that in the corpus:

- "students opened their" occurred 1000 times
- "students opened their books" occurred 400 times
  - → P(books | students opened their) = 0.4
- "students opened their exams" occurred 100 times
  - → P(exams | students opened their) = 0.1

Should we have discarded the "proctor" context?

# Sparsity Problems 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 $\delta$ 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.

# Storage Problems with n-gram Language Models

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

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

Increasing *n* or increasing corpus increases model size!

# 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

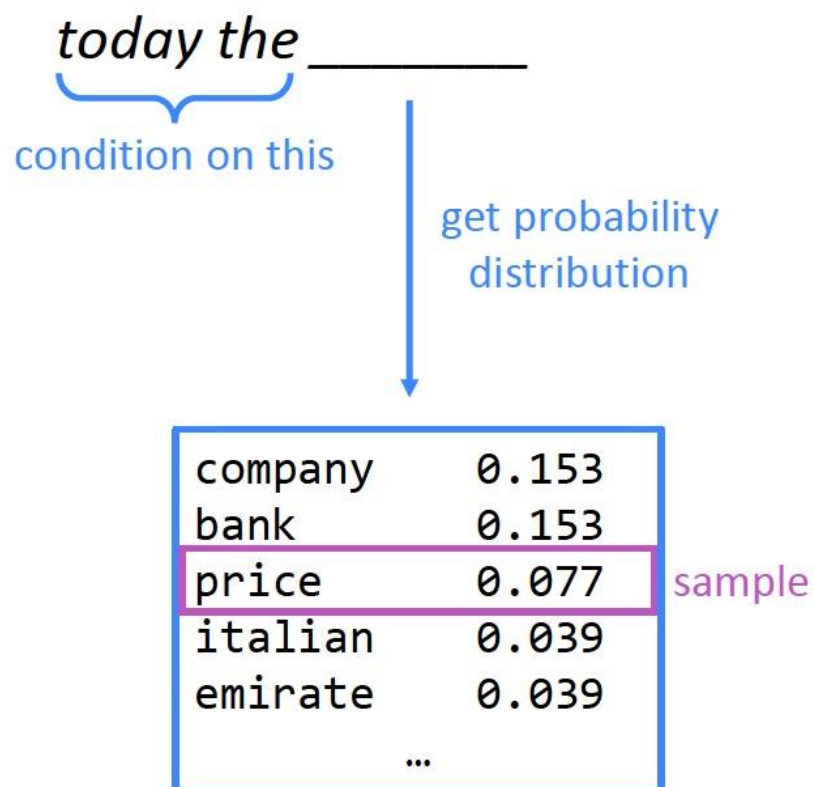  | | |
  |---|---|
  | company | 0.153 |
  | bank | 0.153 |
  | price | 0.077 |
  | italian | 0.039 |
  | emirate | 0.039 |
  | ... | |

  **Sparsity problem**:
  not much granularity
  in the probability
  distribution
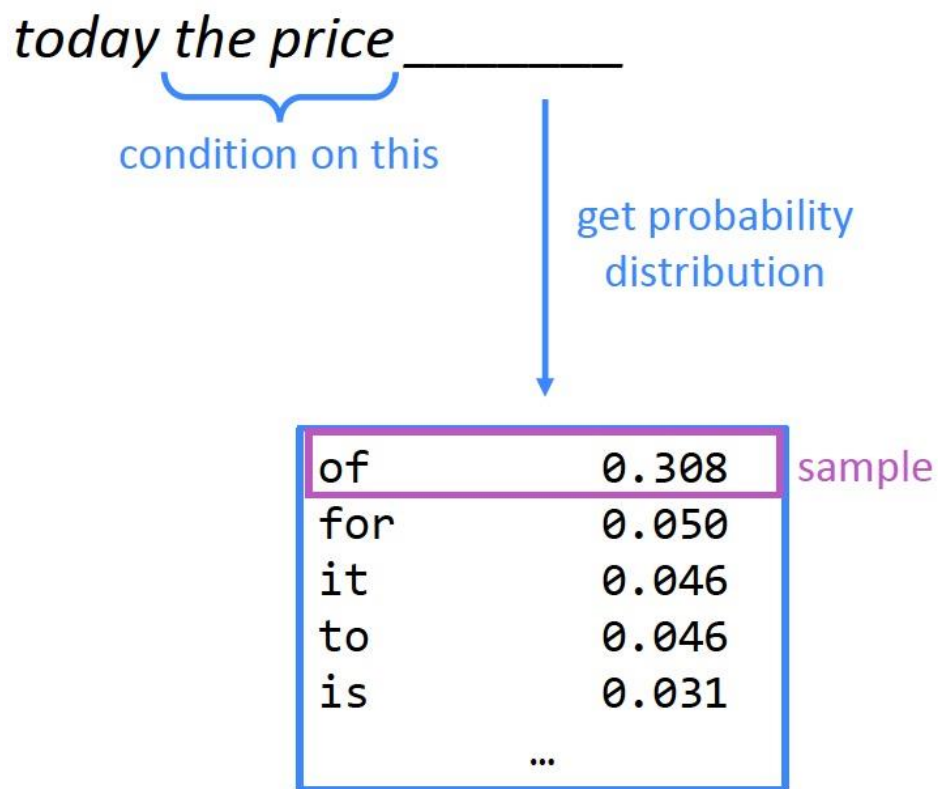
  Otherwise, seems reasonable!

  **\* Try for yourself:** https://nlpforhackers.io/language-models/

# Generating text with a n-gram Language Model

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

*today the* _____

condition on this

get probability
distribution

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

sample

# Generating text with a n-gram Language Model

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

*today the price* _____

condition on this

get probability
distribution

| | |
|---|---|
| of | 0.308 |
| for | 0.050 |
| it | 0.046 |
| to | 0.046 |
| is | 0.031 |
| ... | |

sample

# Generating text with a n-gram Language Model

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

*today the price of* _____

condition on this

get probability
distribution

| | |
|------|-------|
| the | 0.072 |
| 18 | 0.043 |
| oil | 0.043 |
| its | 0.036 |
| gold | 0.018 |
| ... | |

sample

# Generating text with a n-gram Language Model

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

*today the price of gold _____*

# 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!

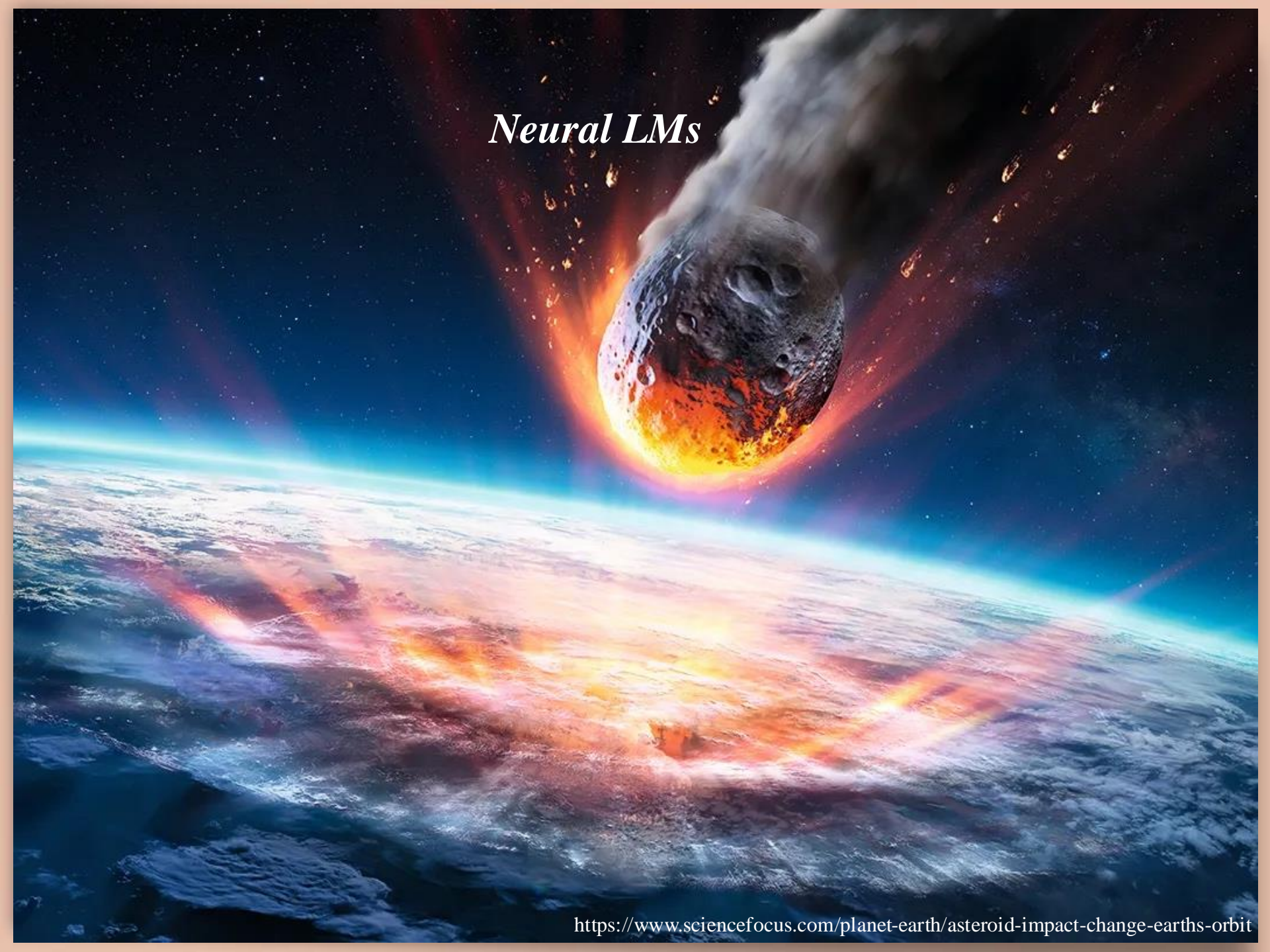…but **incoherent.** We need to consider more than three words at a time if we want to model language well.

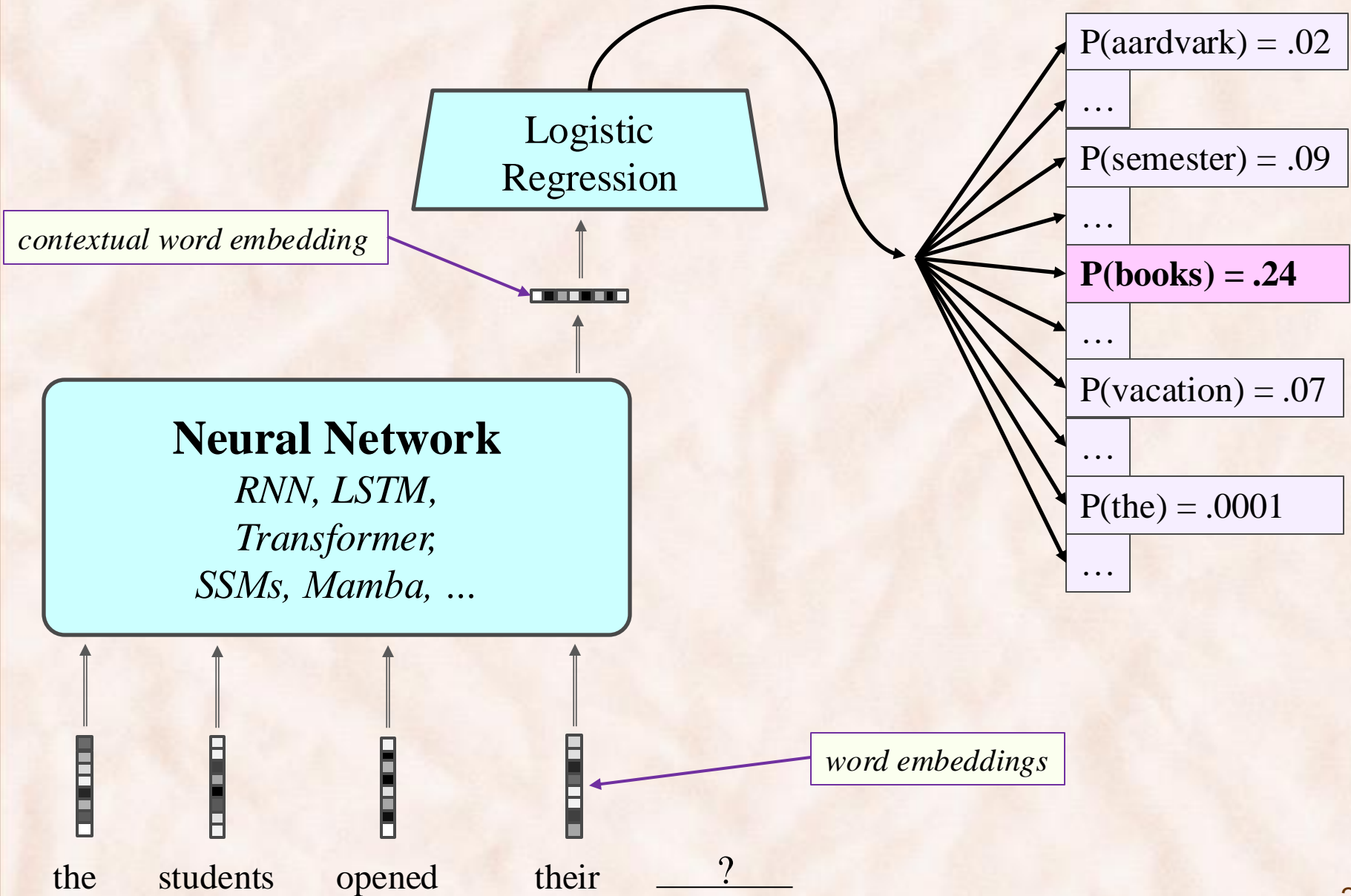But increasing *n* worsens sparsity problem, and increases model size…

# N-gram Language Models are Shallow Learners

- N-gram LMs exhibit extremely little "understanding" of language due to their **weak generalization**.

- N-gram LMs do not generalize to:
  - Unseen words:
    - *spinach* and *kale* are considered completely different words.
    - sim(*spinach*, *kale*) = sim(*spinach*, *laptop*) = 0.

  - Unseen combinations of known words.

  - Longer context.

*Neural LMs*

https://www.sciencefocus.com/planet-earth/asteroid-impact-change-earths-orbit

Logistic Regression

*contextual word embedding*

**Neural Network**
*RNN, LSTM, Transformer, SSMs, Mamba, …*

P(aardvark) = .02

…

P(semester) = .09

…

**P(books) = .24**

…

P(vacation) = .07

…

P(the) = .0001

…

*word embeddings*

the    students    opened    their    ____?____

21

# What can be learned from predicting missing words?

Asheville is a city located in the state of _____.

# What can be learned from predicting missing words?

I took ___ dog out for a walk .

# What can be learned from predicting missing words?

Upon exiting the restaurant, the man realized ___ left ___ phone at the table .

# What can be learned from predicting missing words?

I stopped by the grocery store to buy bread, blueberry pie, milk, and ____ .

# What can be learned from predicting missing words?

Overall , the value I got from the two hours watching it was

the sum total of the popcorn and the drink .

The movie was _____ .

# What can be learned from predicting missing words?

Andrei was eating in the kitchen .

Roxby joined him for breakfast .

After a while , Andrei went to the living room to watch TV .

Once she was done with breakfast , Roxby left the _____ .

# What can be learned from predicting missing words?

Dan and Tom go to a restaurant for dinner . Dan leaves his coat on the chair , then goes to the bathroom . While Dan was gone , Tom hangs Dan 's coat on the coat rack . When Dan comes back , he thinks his coat is on the ___ .

# What can be learned from predicting missing words?

I have been thinking of the sequence that goes 1, 2, 4, 7, 11, 16, 22, _____

# What can be learned from predicting missing words?

Theorem: √2 is an irrational number .

Proof: ‾‾‾‾

Theorem: √2 is an irrational number .

Proof: Suppose that √2 were a ‾‾‾‾

Theorem: √2 is an irrational number .

Proof: Suppose that √2 were a rational number , so by definition
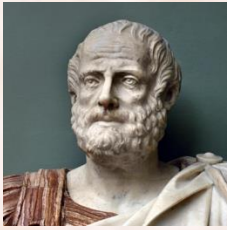
√2 = a ‾‾‾‾

# Word Embeddings

- Neural LMs use as input **word embeddings**:
  - **Vectors** that are **short**, e.g. 256, 512, 768, or 1024, and **dense**, e.g. most if not all entries are non-zero.
    - Very unlike *one-hot encodings,* which are long and sparse.

- We want word embeddings to represent **word meanings**:
  - What do we mean by "word meaning"?
  - How do we train word embeddings to represent word meaning?
    - How do we determine their quality?

# What do words mean?

- Classical approach uses a **dictionary** and the **context**.

- Depending on the context, a word is used to refer to a particular *concept*, i.e. *sense*, i.e. *word meaning*:
  - The word "pepper" has multiple meanings, as listed in the dictionary:
    - sense 1: spice from pepper plant
    - sense 2: the pepper plant itself
    - sense 3: another similar plant (Jamaican pepper)
    - sense 4: another plant with peppercorns (California pepper)
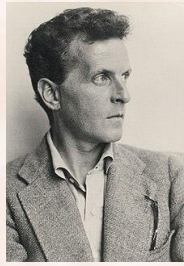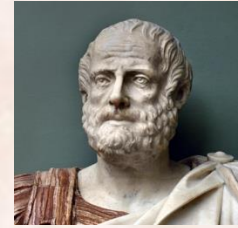    - sense 5: capsicum (i.e. chili, paprika, bell pepper, etc)

# Aristotle: What is a concept?

- Classical theory of concepts:
  - The meaning of a word is a concept that is defined by **necessary** and **sufficient** conditions.
    - The following necessary conditions, jointly, are sufficient for an object $x$ to be a **square**:
      - $x$ has (exactly) four sides
      - each of $x$'s sides is straight
      - $x$ is a closed figure
      - $x$ lies in a plane
      - each of $x$'s sides is equal in length to each of the others
      - each of $x$'s interior angles is equal to the others (right angles)
      - the sides of $x$ are joined at their ends

# Wittgenstein to Aristotle:
## *But what is a **game***?

- Philosophical Investigations (1945, # 66):

  Don't say "there must be something common, or they would not be called `games'"—but *look and see* whether there is anything common to all"

  – Is it amusing? Is there competition? Is there long-term strategy?

  – Is skill required? Must luck play a role?

  – Are there cards? Is there a ball?

  – …

- Family Resemblance theory (Rosch and Mervis):

  – Each item has at least one, and probably several, elements in common with one or more items, but no, or few, elements are common to all items.

| Game 1 | Game 2 | Game 3 | Game 4 |
| --- | --- | --- | --- |
| ABC | BCD | ACD | ABD |

# Distributional Semantics Idea

- Wittgenstein (1945):

  **The meaning of a word is its use in the language.**

- Harris (1954):

  **Words that occur in the same contexts tend to have similar meanings.**

- Firth (1935, 1957):

  **The complete meaning of a word is always contextual, and no study of meaning apart from a complete context can be taken seriously.**

  **You shall know a word by the company it keeps.**

# Distributional Semantics

- The **meaning** of a word is **determined by** the words that appear nearby, i.e. its **context**.
    - **Words that appear in the same contexts tend to have similar meanings.**
    - One of the most successful ideas of modern statistical NLP!

# What does recent English borrowing *ongchoi* mean?

- Suppose you see these sentences:
  - Ong choi is delicious **sautéed with garlic**
  - Ong choi is superb **over rice**
  - Ong choi **leaves** with **salty** sauces

- And you've also seen these:
  - …spinach **sautéed with garlic over rice**
  - Chard stems and **leaves** are **delicious**
  - Collard greens and other **salty** leafy greens

- Conclusion:
  - Ongchoi is a leafy green like spinach, chard, or collard greens

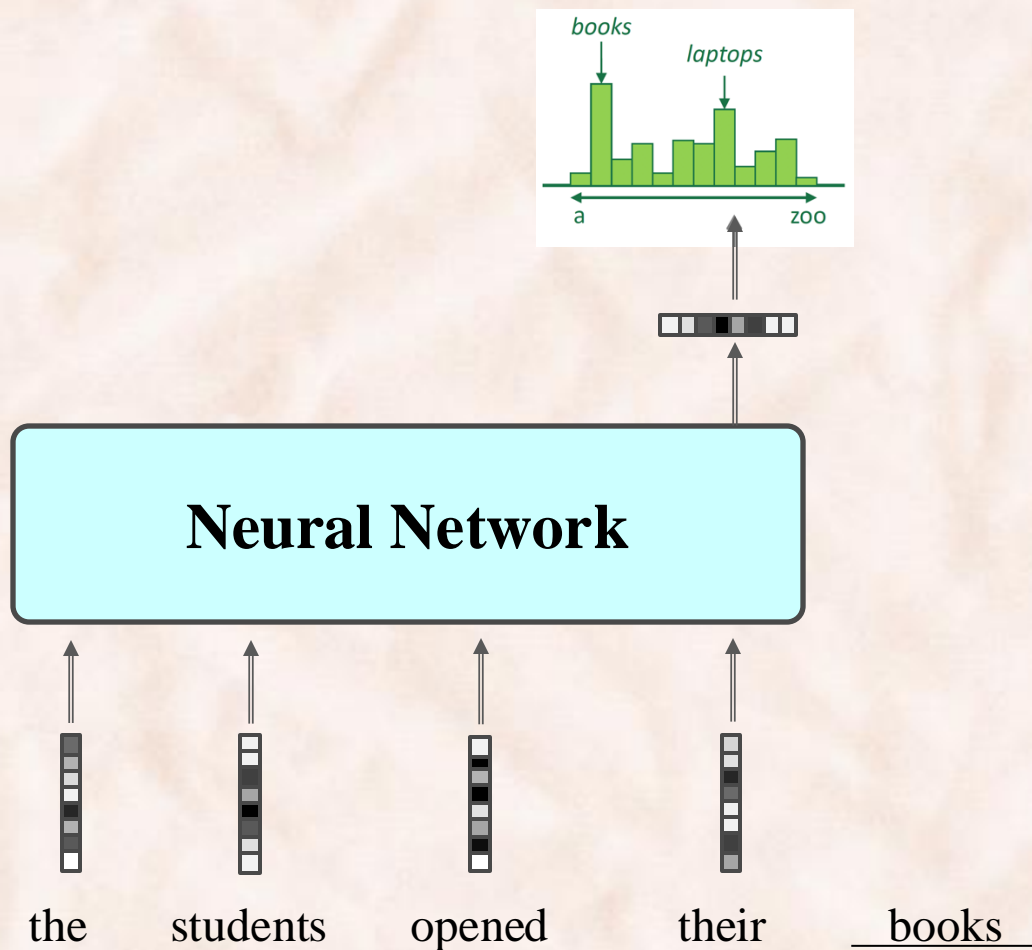# Ongchoi: *Ipomoea aquatica "Water Spinach"*

空心菜
*kangkong*
rau muống
…



Yamaguchi, Wikimedia Commons, public domain

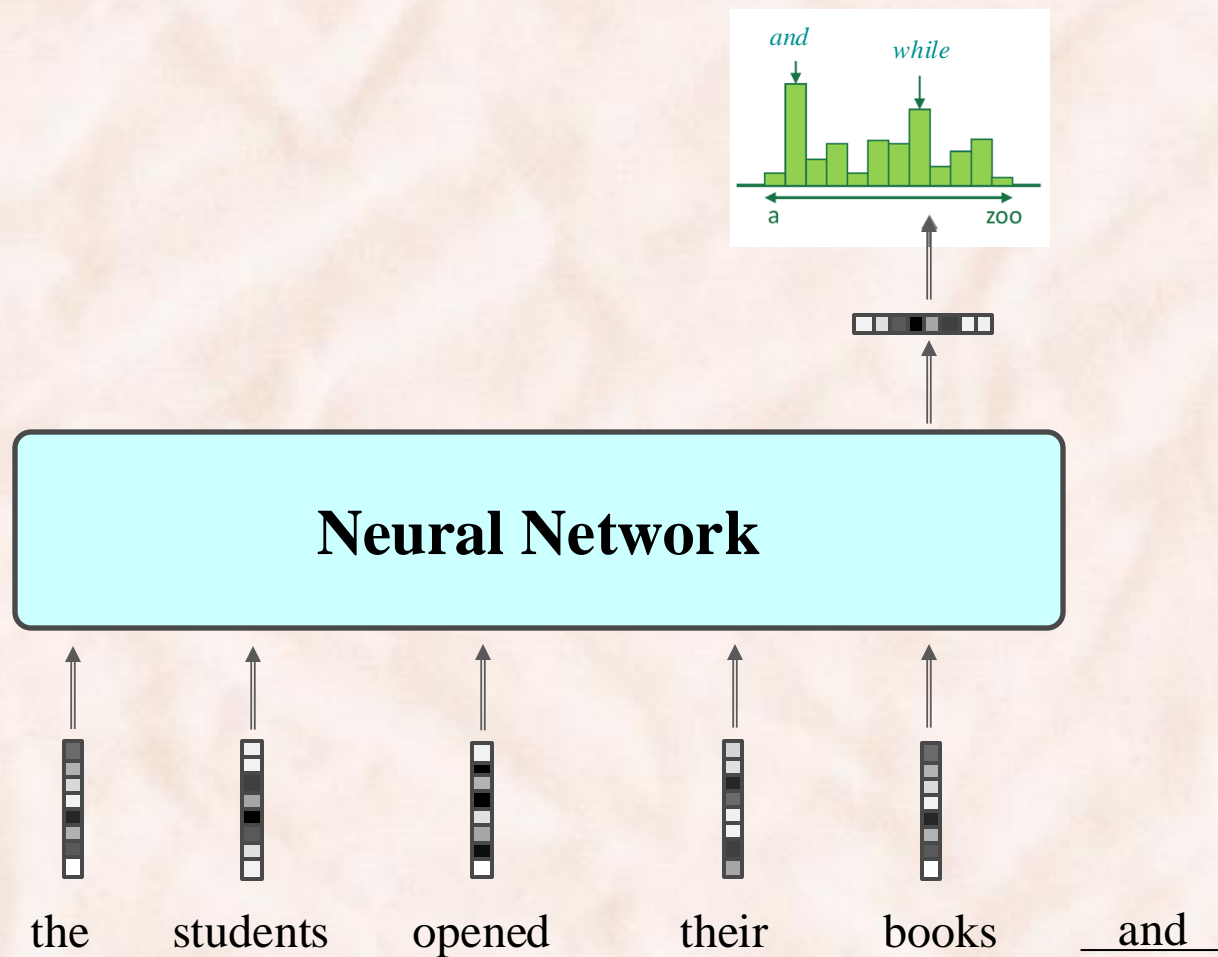# Distributional Semantics and Neural LMs

- Learn word embeddings by training a neural network to **predict a missing word** given words in the **context**.
    - This is a special type of *reconstructing the input* idea used in other modalities, such as computer vision (see *autoencoders*).

- **Causal LMs** trained using the distributional hypothesis:
    - Context: words so far.
    - Missing word: next word.

- **Masked LMs** trained using the distributional hypothesis:
    - Context: words to the left and to the right of the center word.
    - Missign word: word in the center.
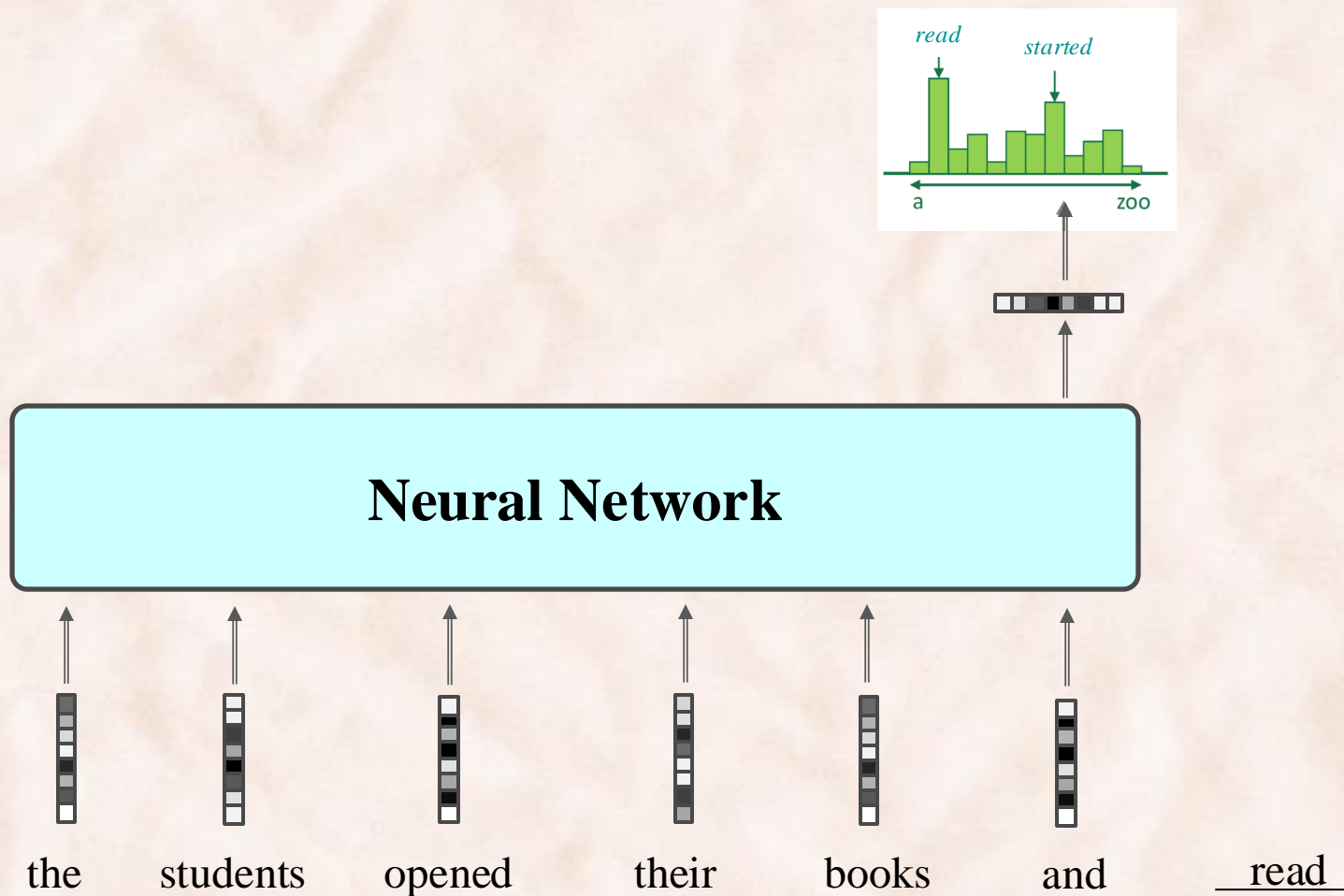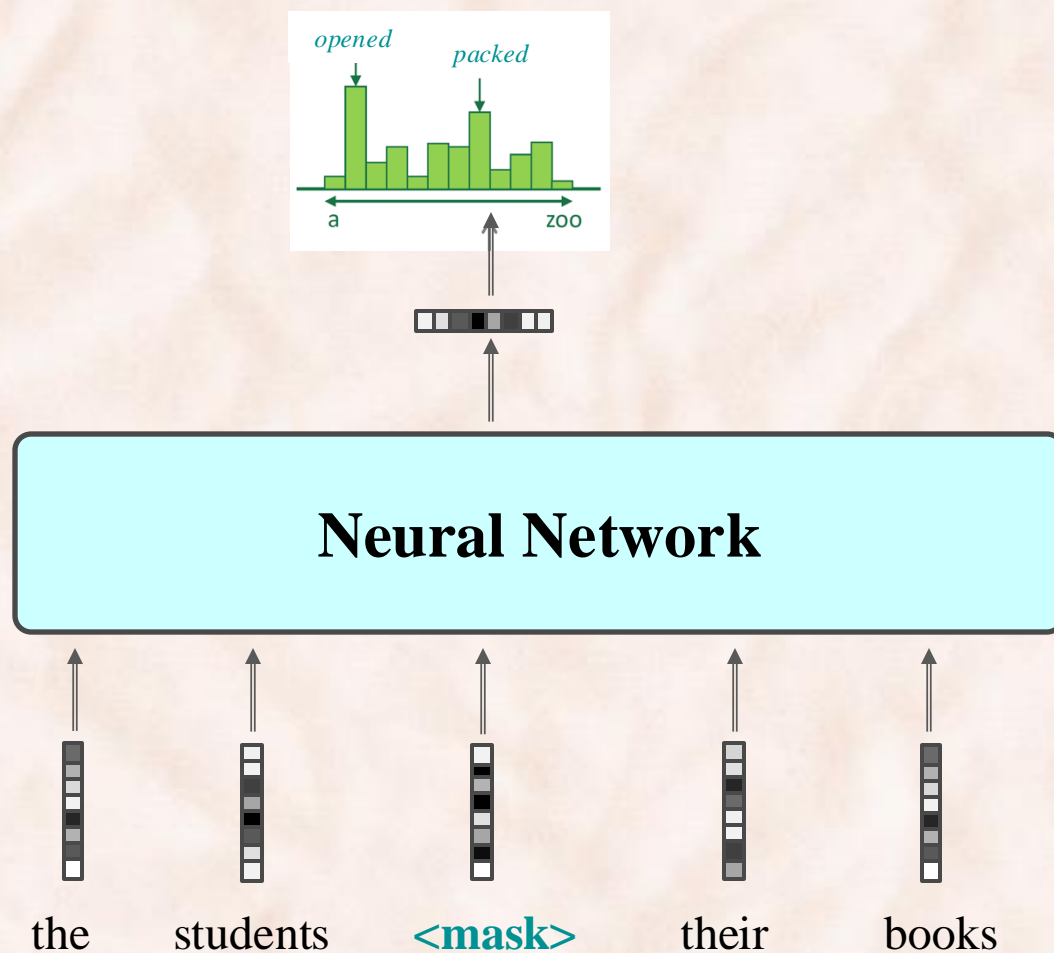
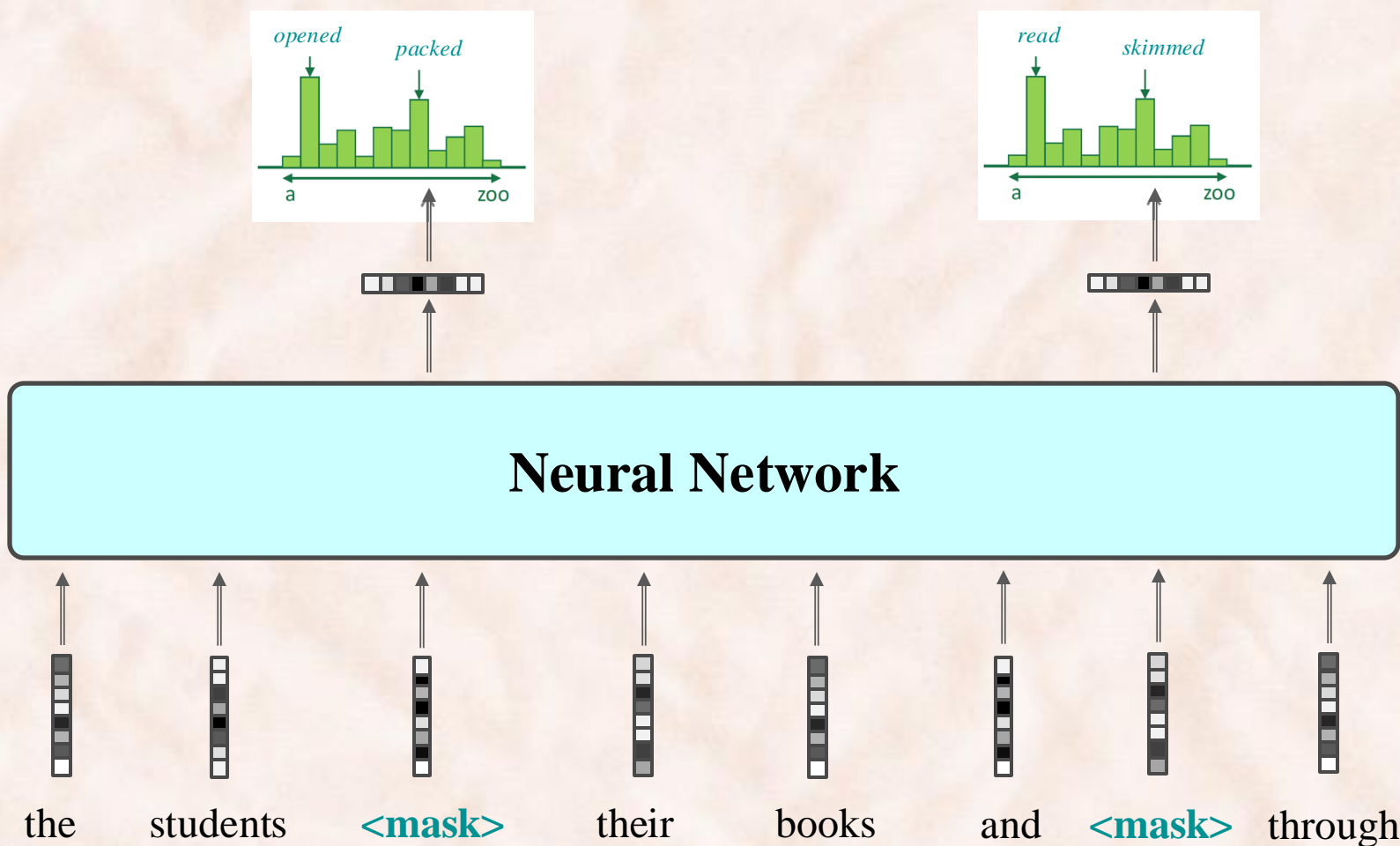# Neural Language Modeling: **Decoders**



**Neural Network**

the      students     opened     their     __books__

# Neural Language Modeling: **Decoders**



**Neural Network**

the    students    opened    their    books    __and__

# Neural Language Modeling: **Decoders**



**Neural Network**

the      students      opened      their      books      and      __read__

# Neural Language Modeling: **Encoders**



**Neural Network**

the    students    **<mask>**    their    books

# Neural Language Modeling: **Encoders**



**Neural Network**

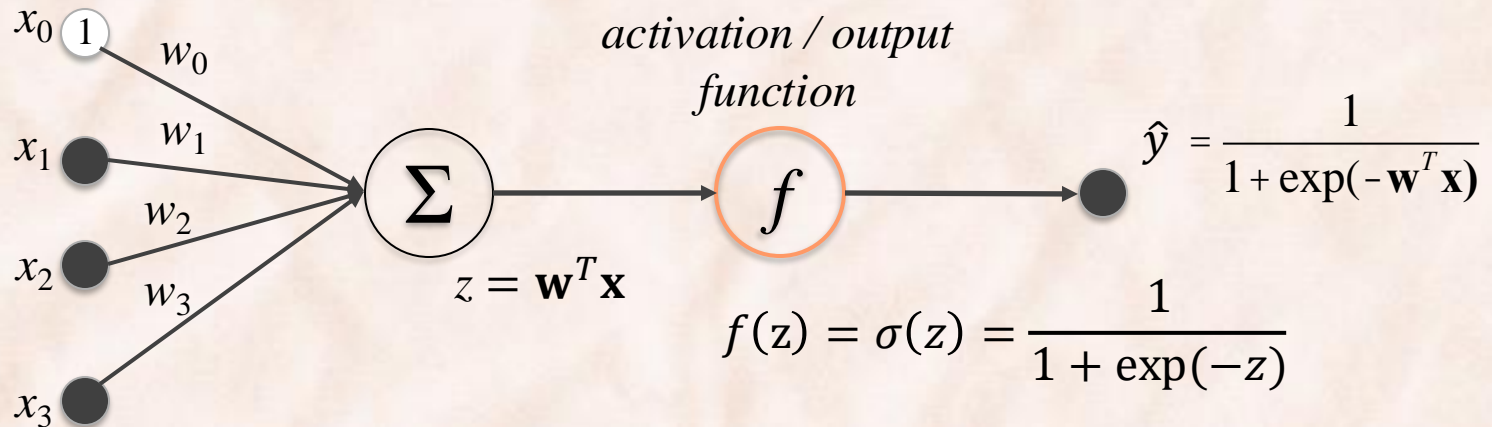the    students    **&lt;mask&gt;**    their    books    and    **&lt;mask&gt;**    through

# Start: Interlude on Neural Networks

- Required reading:
    - [Chapter 7 on Neural Networks](#) from the J&M textbook.
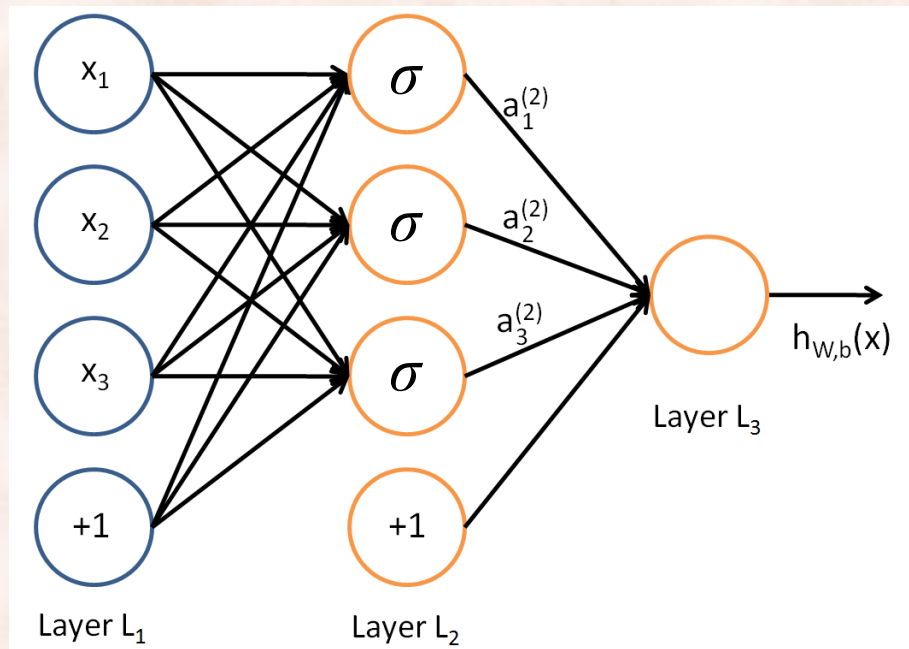        - Sections on Training are optional.

# Logistic Neuron = Logistic Regression



$x_0$ 1
$w_0$
$x_1$
$w_1$
$x_2$
$w_2$
$x_3$
$w_3$

$$z = \mathbf{w}^T\mathbf{x}$$

*activation / output function*

$f$

$$f(z) = \sigma(z) = \frac{1}{1 + \exp(-z)}$$

$$\hat{y} = \frac{1}{1 + \exp(-\mathbf{w}^T\mathbf{x})}$$

- Algebraic interpretation:
  - The output of the neuron is a **linear combination** of inputs from other neurons, **rescaled by** the synaptic **weights**.
    - weights $w_i$ correspond to the synaptic weights (activating or inhibiting).
    - summation corresponds to combination of signals in the soma.
  - It is often transformed through a monotonic **activation function**.

# Neural Network Model

- Put together many neurons in layers, such that the output of a neuron on layer $l$ can be the input of another neuron on layer $l + 1$:
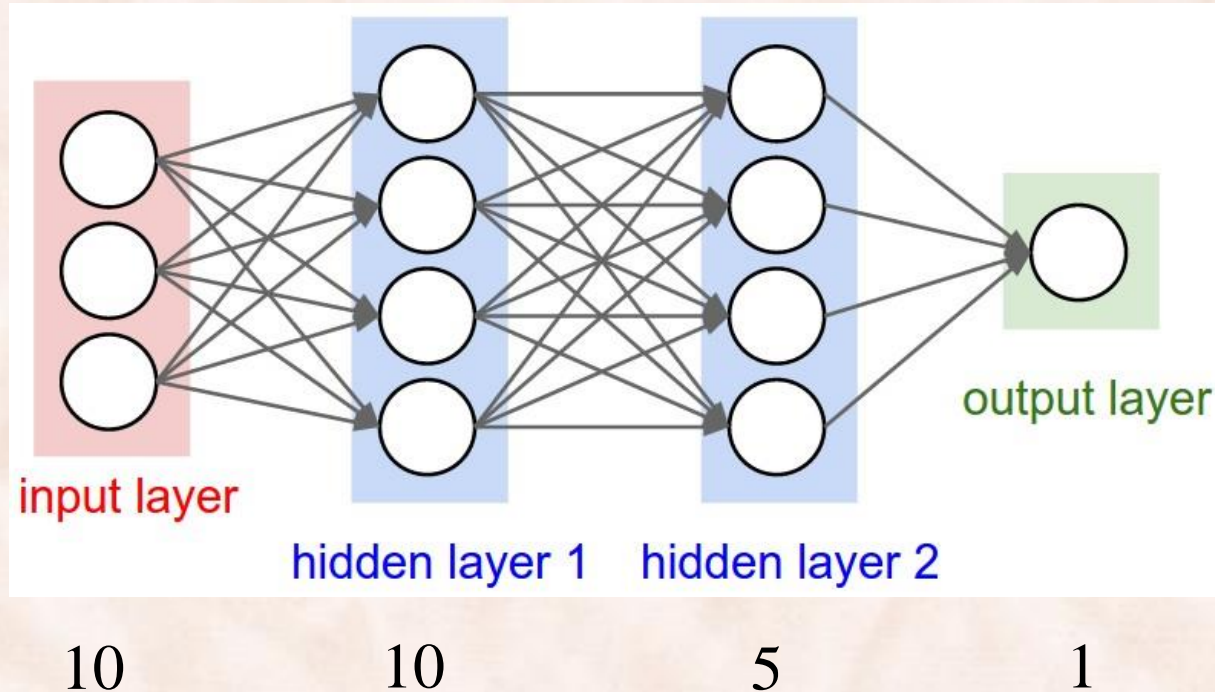


***input layer***          ***hidden layer***          ***output layer***
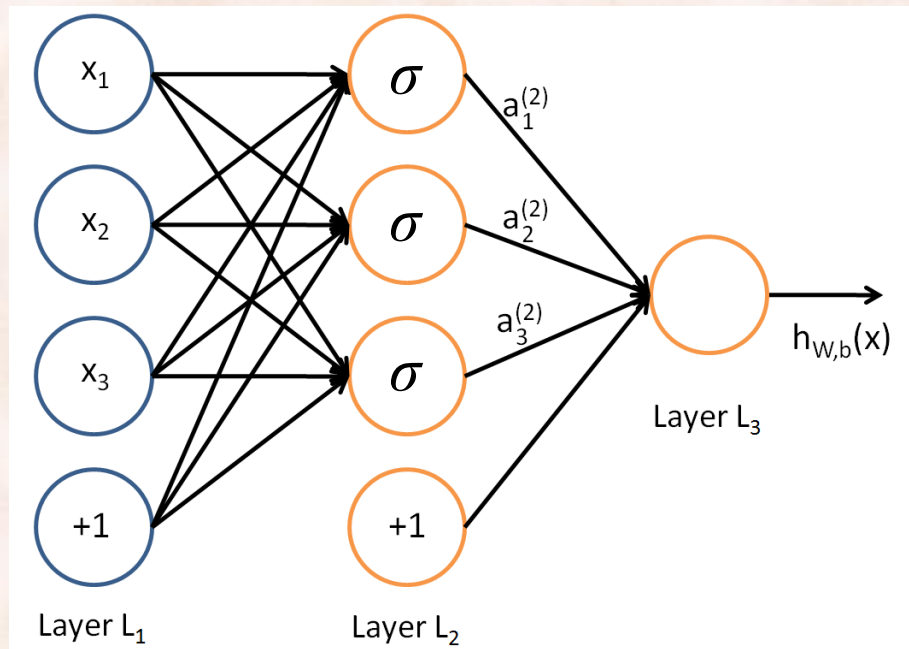
# Feed-Forward Neural Networks



| 10 | 10 | 5 | 1 |

1. For each neuron in hidden layer 1, we need $10 + 1 = 11$ params. For the 10 neurons on hidden layer 1, we need in total $10 * 11 = $ **110 params**.
2. For the 5 neurons on hidden layer 2, we need $5 * 11 = $ **55 params**.
3. For the output neurons, we need $5 + 1 = $ **6 params**.

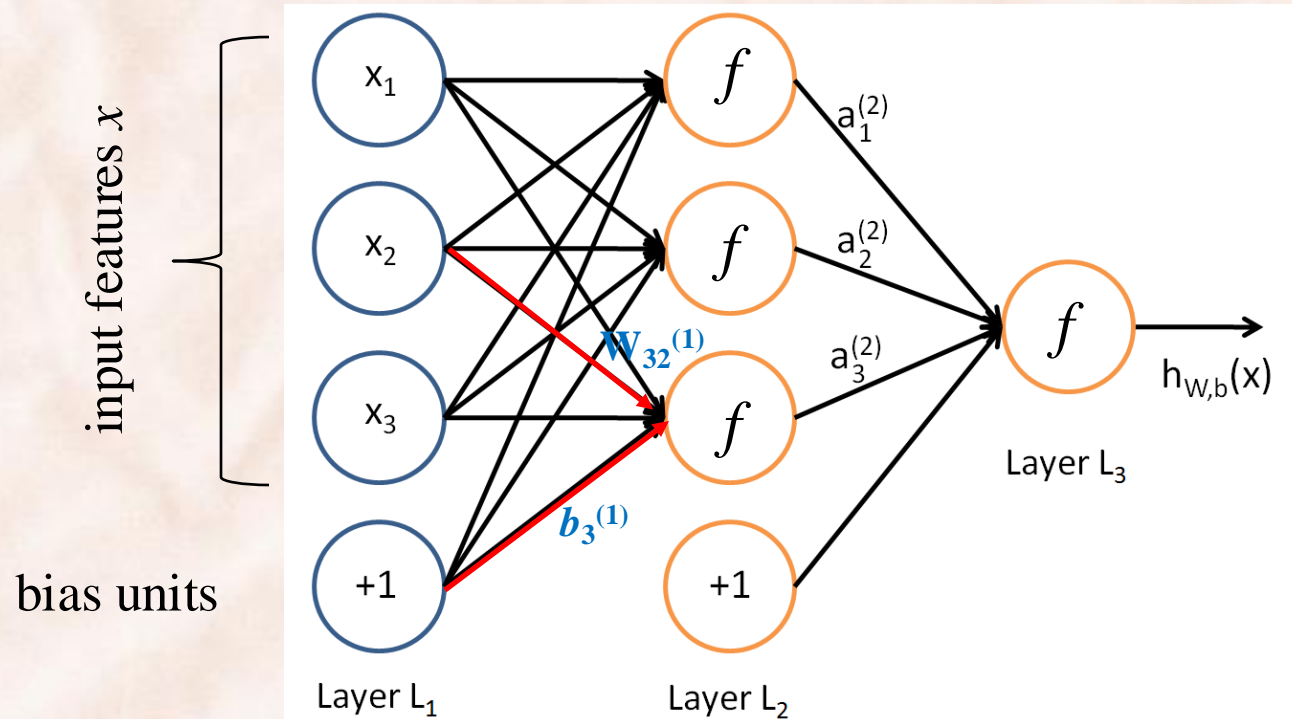# Neural Network Model

- Put together many neurons in layers, such that the output of a neuron can be the input of another:



***input layer***        ***hidden layer***        ***output layer***

- $n_l = 3$ is the number of **layers**.
  - $L_1$ is the input layer, $L_3$ is the output layer
- $(\mathbf{W}, \mathbf{b}) = (W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)})$ are the parameters:
  - $W^{(l)}_{ij}$ is the **weight** of the connection between unit $j$ in layer $l$ and unit $i$ in layer $l + 1$.
  - $b^{(l)}_i$ is the **bias** associated unit unit $i$ in layer $l + 1$.
- $a^{(l)}_i$ is the **activation** of unit i in layer $l$, e.g. $a^{(1)}_i = x_i$ and $a^{(3)}_1 = h_{W,b}(x)$.

# Inference: Forward Propagation

- The activations in the hidden layer are:

$$a_1^{(2)} = f(W_{11}^{(1)}x_1 + W_{12}^{(1)}x_2 + W_{13}^{(1)}x_3 + b_1^{(1)})$$
$$a_2^{(2)} = f(W_{21}^{(1)}x_1 + W_{22}^{(1)}x_2 + W_{23}^{(1)}x_3 + b_2^{(1)})$$
$$a_3^{(2)} = f(W_{31}^{(1)}x_1 + W_{32}^{(1)}x_2 + W_{33}^{(1)}x_3 + b_3^{(1)})$$

- The activations in the output layer are:

$$h_{W,b}(x) = a_1^{(3)} = f(W_{11}^{(2)}a_1^{(2)} + W_{12}^{(2)}a_2^{(2)} + W_{13}^{(2)}a_3^{(2)} + b_1^{(2)})$$

- Compressed notation:

$$a_i^{(l)} = f(z_i^{(l)}) \text{ where } z_i^{(2)} = \sum_{j=1}^{n} W_{ij}^{(1)}x_j + b_i^{(1)}$$

# Forward Propagation

- Forward propagation (unrolled):

$$a_1^{(2)} = f(W_{11}^{(1)}x_1 + W_{12}^{(1)}x_2 + W_{13}^{(1)}x_3 + b_1^{(1)})$$

$$a_2^{(2)} = f(W_{21}^{(1)}x_1 + W_{22}^{(1)}x_2 + W_{23}^{(1)}x_3 + b_2^{(1)})$$

$$a_3^{(2)} = f(W_{31}^{(1)}x_1 + W_{32}^{(1)}x_2 + W_{33}^{(1)}x_3 + b_3^{(1)})$$

$$h_{W,b}(x) = a_1^{(3)} = f(W_{11}^{(2)}a_1^{(2)} + W_{12}^{(2)}a_2^{(2)} + W_{13}^{(2)}a_3^{(2)} + b_1^{(2)})$$

- Forward propagation (compressed):

$$z^{(2)} = W^{(1)}x + b^{(1)}$$

$$a^{(2)} = f(z^{(2)})$$

$$z^{(3)} = W^{(2)}a^{(2)} + b^{(2)}$$

$$h_{W,b}(x) = a^{(3)} = f(z^{(3)})$$

- Element-wise application:

$$f(\mathbf{z}) = [f(z_1), f(z_2), f(z_3)]$$

# Forward Propagation

- Forward propagation (compressed):

$$z^{(2)} = W^{(1)}x + b^{(1)}$$
$$a^{(2)} = f(z^{(2)})$$
$$z^{(3)} = W^{(2)}a^{(2)} + b^{(2)}$$
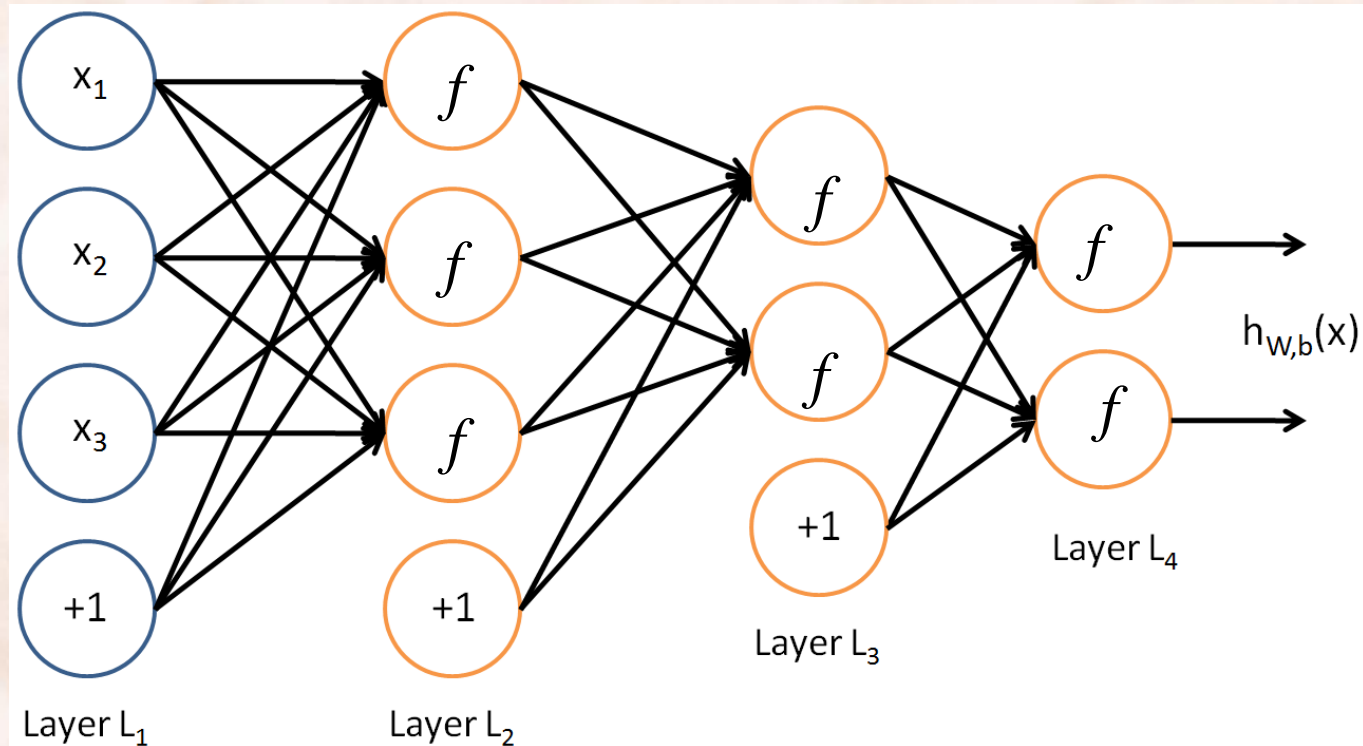$$h_{W,b}(x) = a^{(3)} = f(z^{(3)})$$

- Composed of two *forward propagation steps*:

$$z^{(l+1)} = W^{(l)}a^{(l)} + b^{(l)}$$
$$a^{(l+1)} = f(z^{(l+1)})$$

# Multiple Hidden Units, Multiple Outputs

- Write down the forward propagation steps for:
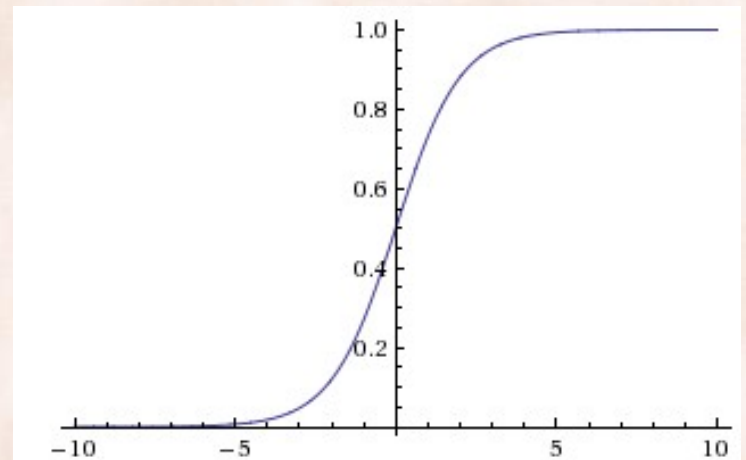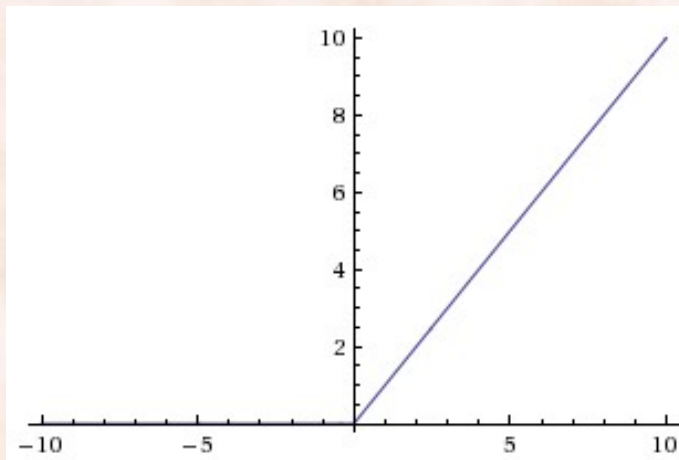
# ReLU and Generalizations

- It has become more common to use piecewise linear activation functions for hidden units instead of σ:

  – **ReLU**: the rectifier activation $g(z) = \max\{0, z\}$.

  – **Absolute value ReLU**: $g(z) = |z|$.

  – **Maxout**: $g(a_1, ..., a_k) = \max\{a_1, ..., a_k\}$.

    - needs k weight vectors instead of 1.

  – **Leaky ReLU**: $g(a) = \max\{0, a\} + \alpha \min(0, a)$.

$\Rightarrow$ the network computes a *piecewise linear function* (up to the output activation function).
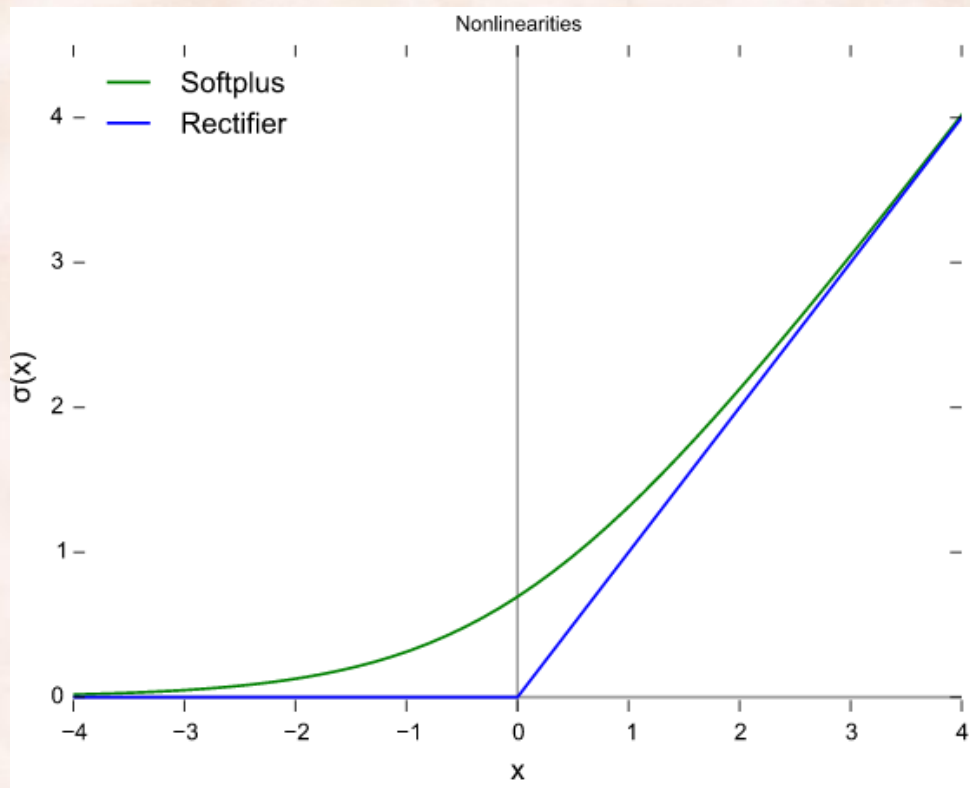
# ReLU vs. Sigmoid and Tanh

- Sigmoid and Tanh saturate for values not close to 0:
  - "kill" gradients, bad behavior for gradient-based learning.
- ReLU does not saturate for values > 0:
  - greatly accelerates learning, fast implementation.
  - fragile during training and can "die", due to 0 gradient:
    - initialize all $b$'s to a small, positive value, e.g. 0.1.

# ReLU vs. Softplus

- Softplus $g(z) = \ln(1+e^z)$ is a smooth version of the rectifier.
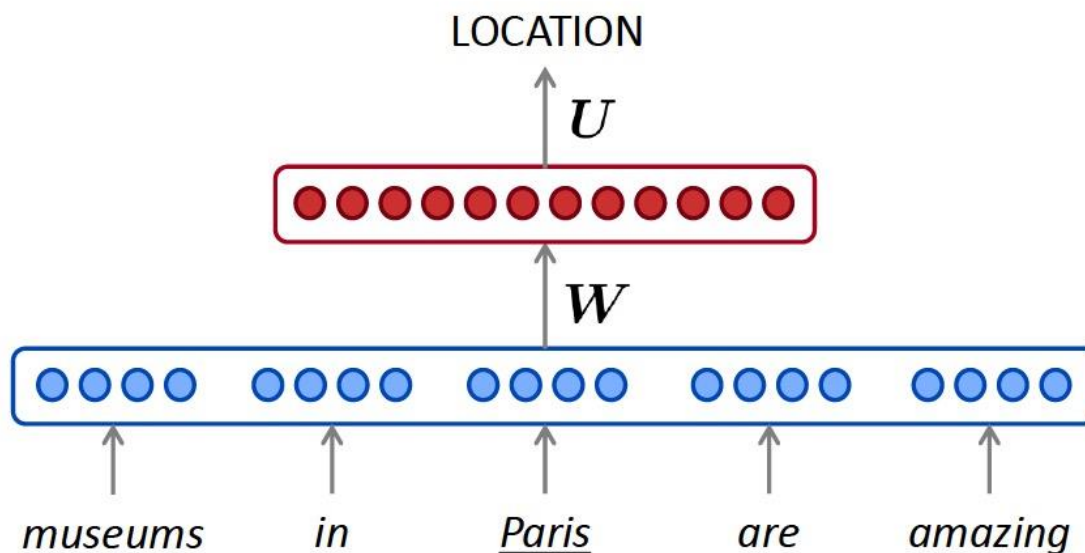  - Saturates less than ReLU, yet ReLU still does better [Glorot, 2011].

# End: Interlude on Neural Networks

- Required reading:
    - [Chapter 7 on Neural Networks](#) from the J&M textbook.
        - Sections on Training are optional.

# How to build a *neural* Language Model?

- Recall the Language Modeling task:
  - Input: sequence of words $x^{(1)}, x^{(2)}, \ldots, x^{(t)}$
  - Output: prob dist of the next word $P(x^{(t+1)} \mid x^{(t)}, \ldots, x^{(1)})$

- How about a window-based neural model?
  - We can apply this to Named Entity Recognition:

LOCATION

$U$

$W$

*museums*   *in*   *Paris*   *are*   *amazing*

# A fixed-window neural Language Model

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

discard

fixed window

# A fixed-window neural Language Model

output distribution

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

hidden layer

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

concatenated word embeddings

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

words / one-hot vectors

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

books

laptops

a          zoo

$\boldsymbol{U}$

$\boldsymbol{W}$

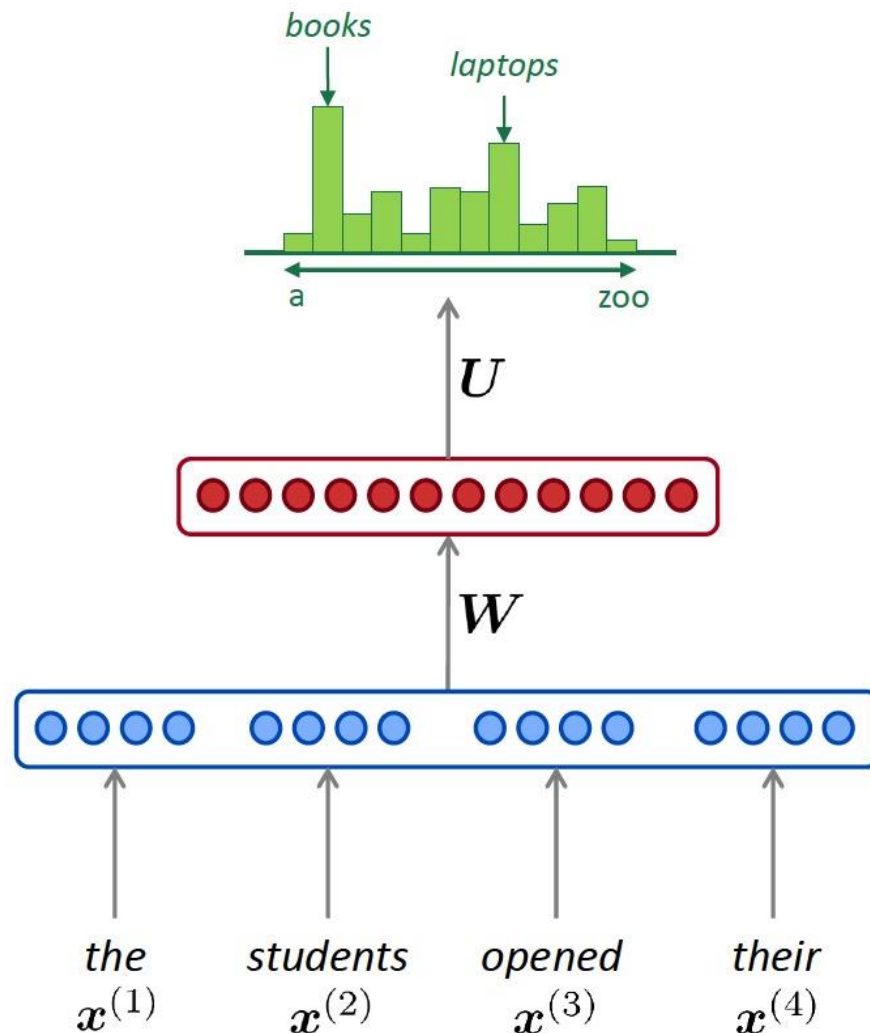| the | students | opened | their |
| $\boldsymbol{x}^{(1)}$ | $\boldsymbol{x}^{(2)}$ | $\boldsymbol{x}^{(3)}$ | $\boldsymbol{x}^{(4)}$ |

20

# A fixed-window neural Language Model

**Improvements** over *n*-gram LM:
- No sparsity problem
- Don't need to store all observed *n*-grams

Remaining **problems**:
- Fixed window is too small
- Enlarging window enlarges $W$
- Window can never be large enough!
- $x^{(1)}$ and $x^{(2)}$ are multiplied by completely different weights in $W$. No symmetry in how the inputs are processed.

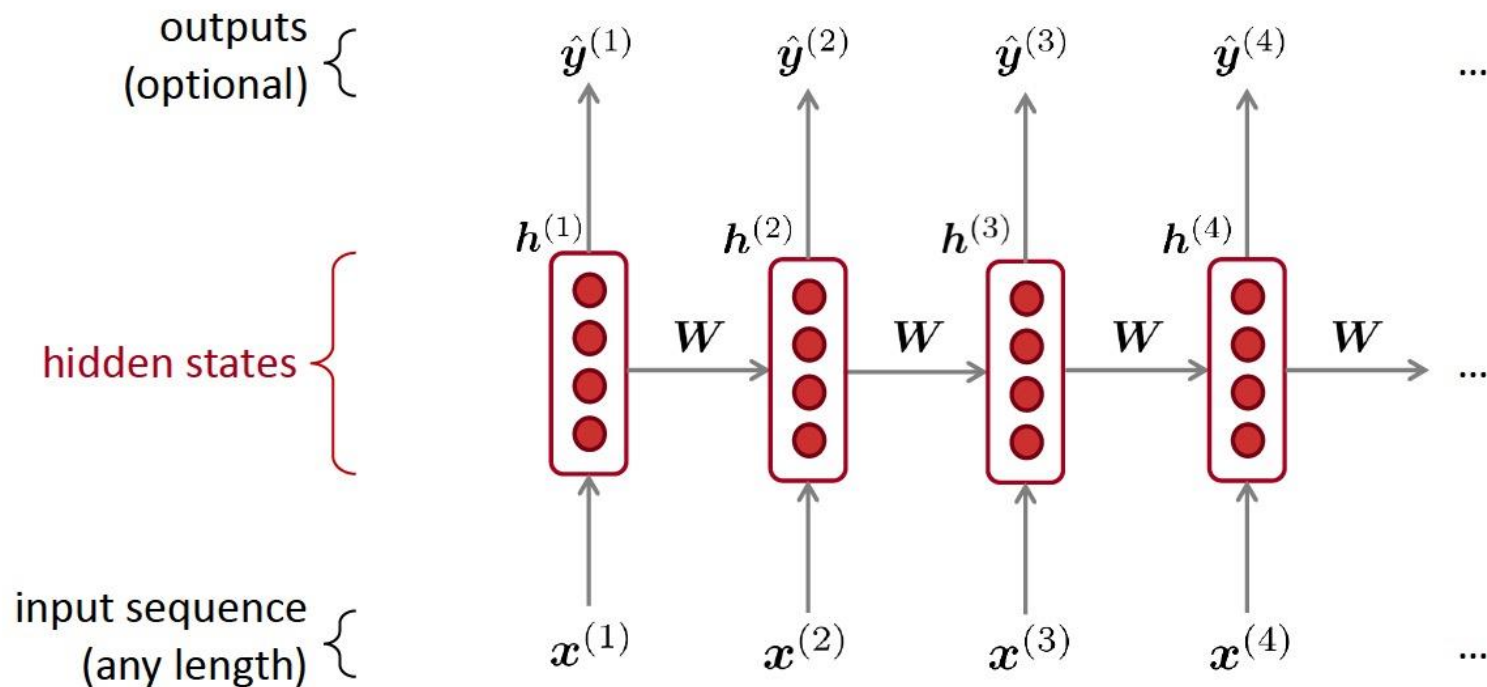We need a neural architecture that can process *any length input*

books

laptops

a                    zoo

$U$

$W$

the $x^{(1)}$    students $x^{(2)}$    opened $x^{(3)}$    their $x^{(4)}$

# Recurrent Neural Networks (RNN)

A family of neural architectures

**Slides from the CS224N at Stanford**



outputs (optional) $\{$  $\hat{y}^{(1)}$  $\hat{y}^{(2)}$  $\hat{y}^{(3)}$  $\hat{y}^{(4)}$  ...

$h^{(1)}$  $h^{(2)}$  $h^{(3)}$  $h^{(4)}$

hidden states $\{$  $W$  $W$  $W$  $W$  ...

input sequence (any length) $\{$  $x^{(1)}$  $x^{(2)}$  $x^{(3)}$  $x^{(4)}$  ...

22

# A RNN Language Model

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

books  laptops

a  zoo

**output distribution**

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

$\boldsymbol{U}$

$\boldsymbol{h}^{(0)}$  $\boldsymbol{h}^{(1)}$  $\boldsymbol{h}^{(2)}$  $\boldsymbol{h}^{(3)}$  $\boldsymbol{h}^{(4)}$

**hidden states**

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

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

$\boldsymbol{W}_h$  $\boldsymbol{W}_h$  $\boldsymbol{W}_h$  $\boldsymbol{W}_h$

$\boldsymbol{W}_e$  $\boldsymbol{W}_e$  $\boldsymbol{W}_e$  $\boldsymbol{W}_e$

**word embeddings**

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

$\boldsymbol{e}^{(1)}$  $\boldsymbol{e}^{(2)}$  $\boldsymbol{e}^{(3)}$  $\boldsymbol{e}^{(4)}$

$\boldsymbol{E}$  $\boldsymbol{E}$  $\boldsymbol{E}$  $\boldsymbol{E}$

**words / one-hot vectors**

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

the  students  opened  their
$\boldsymbol{x}^{(1)}$  $\boldsymbol{x}^{(2)}$  $\boldsymbol{x}^{(3)}$  $\boldsymbol{x}^{(4)}$

23

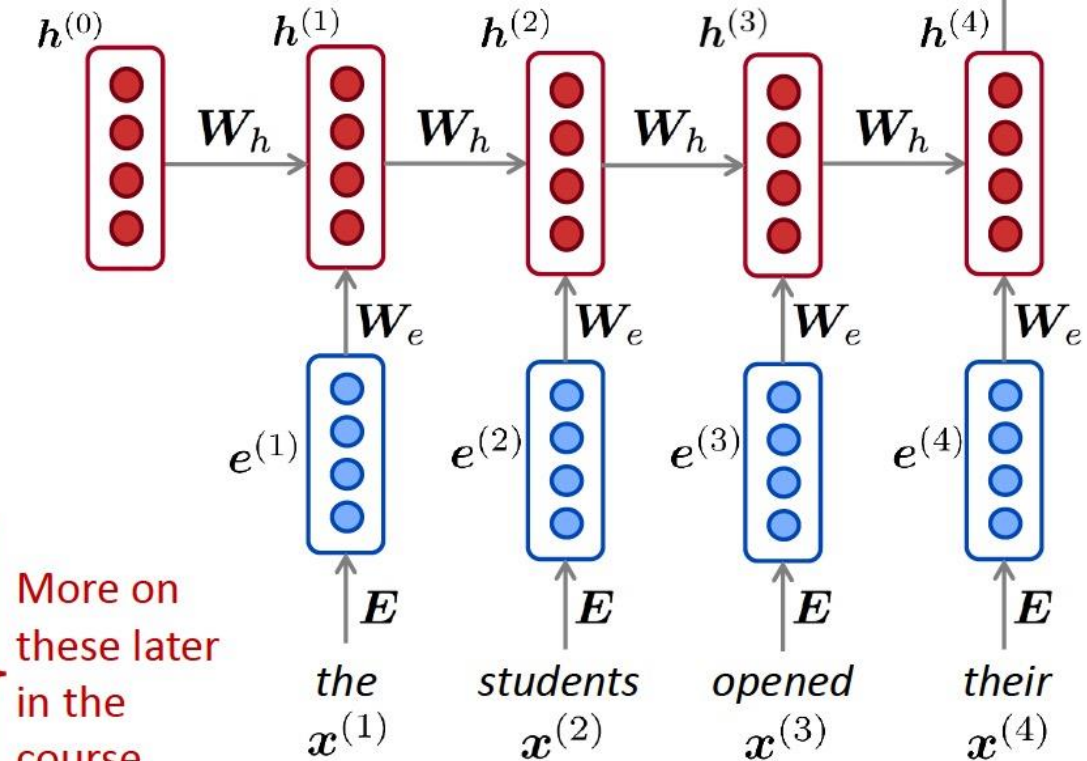**Note**: this input sequence could be much longer, but this slide doesn't have space!

# A RNN Language Model

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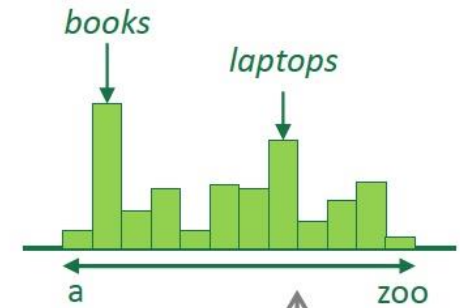- Can process any length input
- Computation for step $t$ can (in theory) use information from many steps back
- Model size doesn't increase for longer input
- Same weights applied on every timestep, so there is symmetry in how inputs are processed.

RNN **Disadvantages**:
- Recurrent computation is slow
- In practice, difficult to access information from many steps back

More on these later in the course

24



$$\hat{y}^{(4)} = P(x^{(5)}|\text{the students opened their})$$
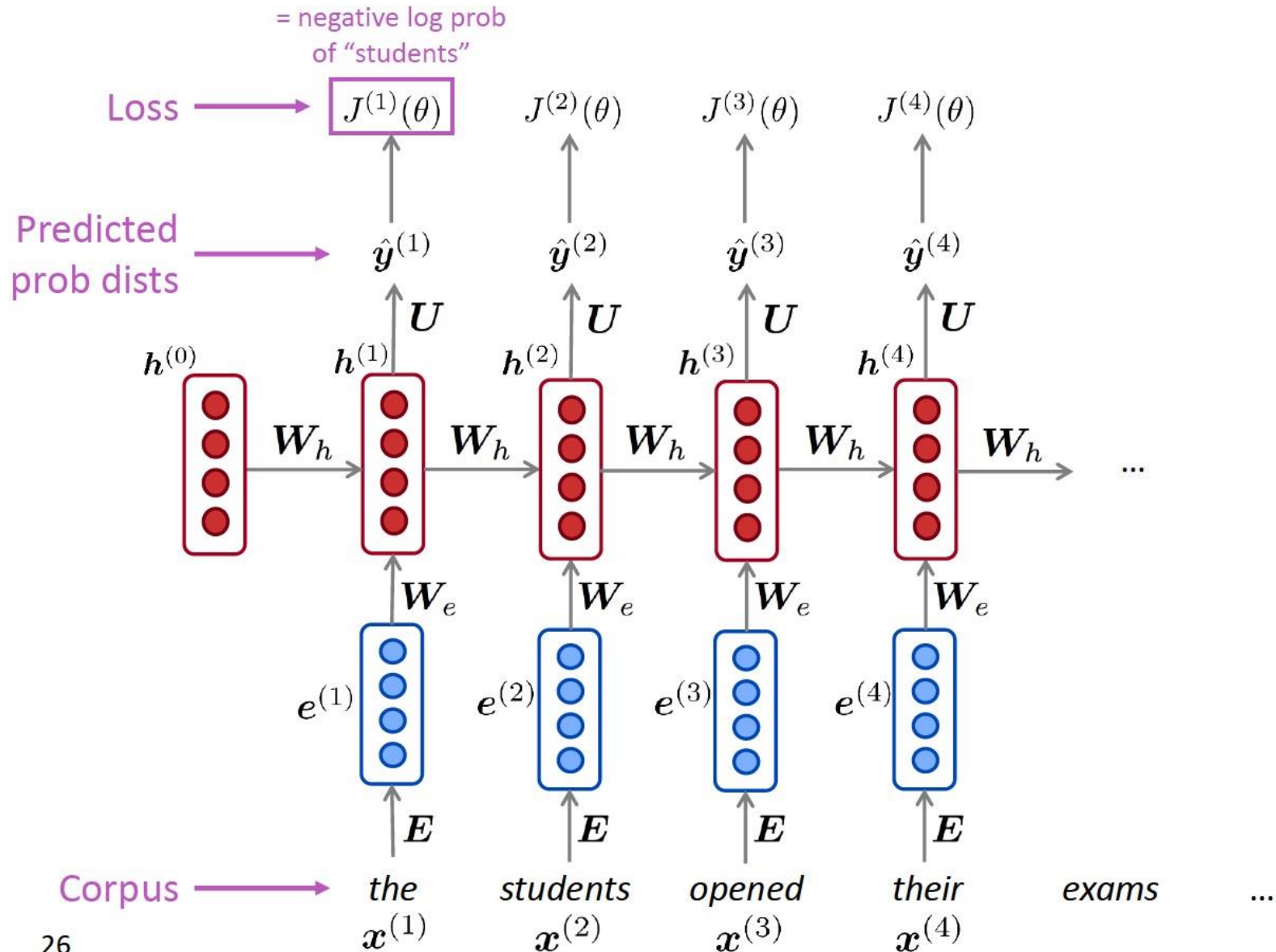
# Start Supplemental Material: Training RNNs

# Training a RNN Language Model

- Get a big corpus of text which is a sequence of words $x^{(1)}, \ldots, x^{(T)}$
- Feed into RNN-LM; compute output distribution $\hat{y}^{(t)}$ for *every step t*.
  - i.e. predict probability dist of *every word*, given words so far

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

$$J^{(t)}(\theta) = CE(y^{(t)}, \hat{y}^{(t)}) = - \sum_{w \in V} y_w^{(t)} \log \hat{y}_w^{(t)} = - \log \hat{y}_{x_{t+1}}^{(t)}$$
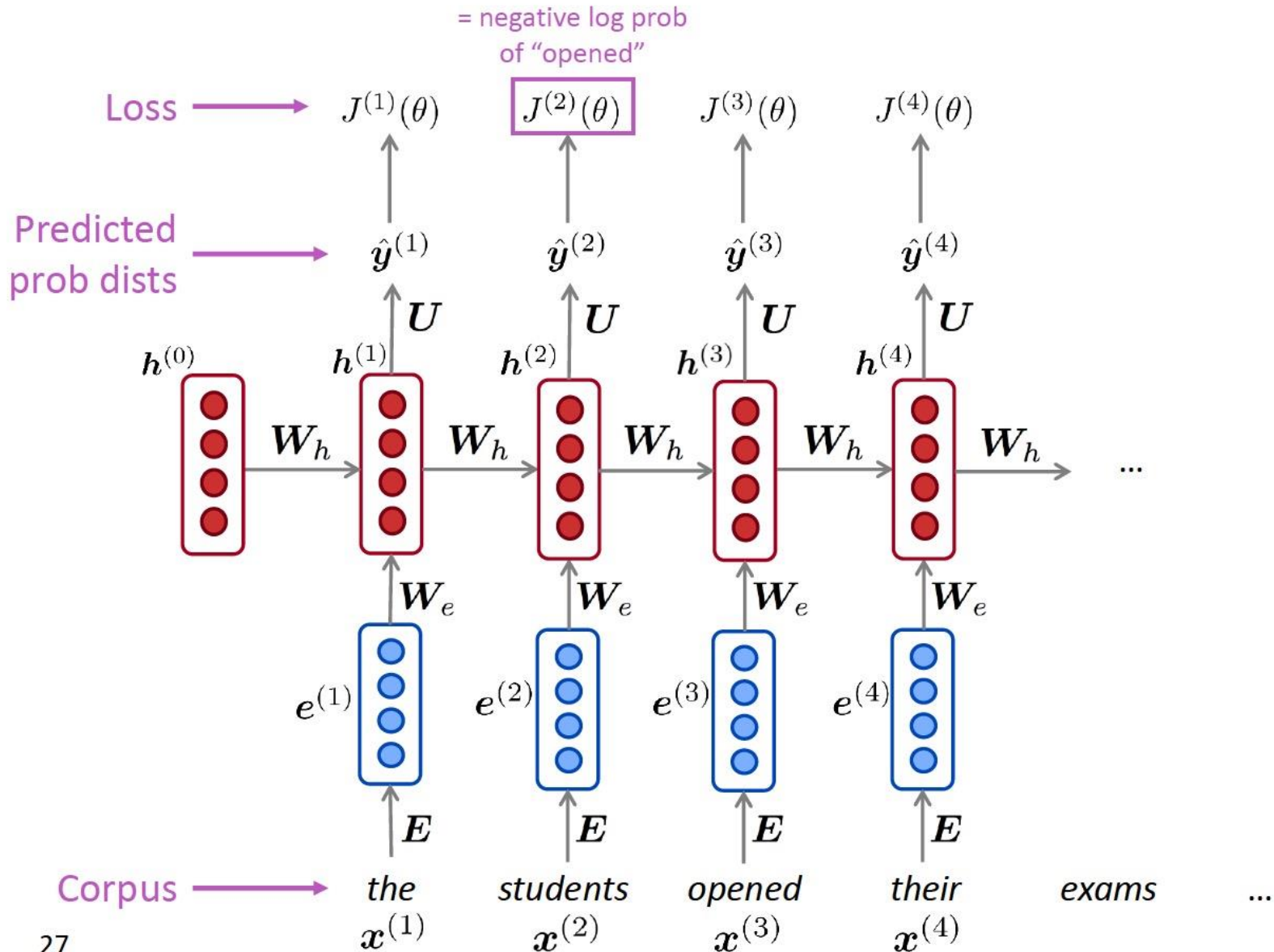
- Average this to get overall loss for entire training set:

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

25

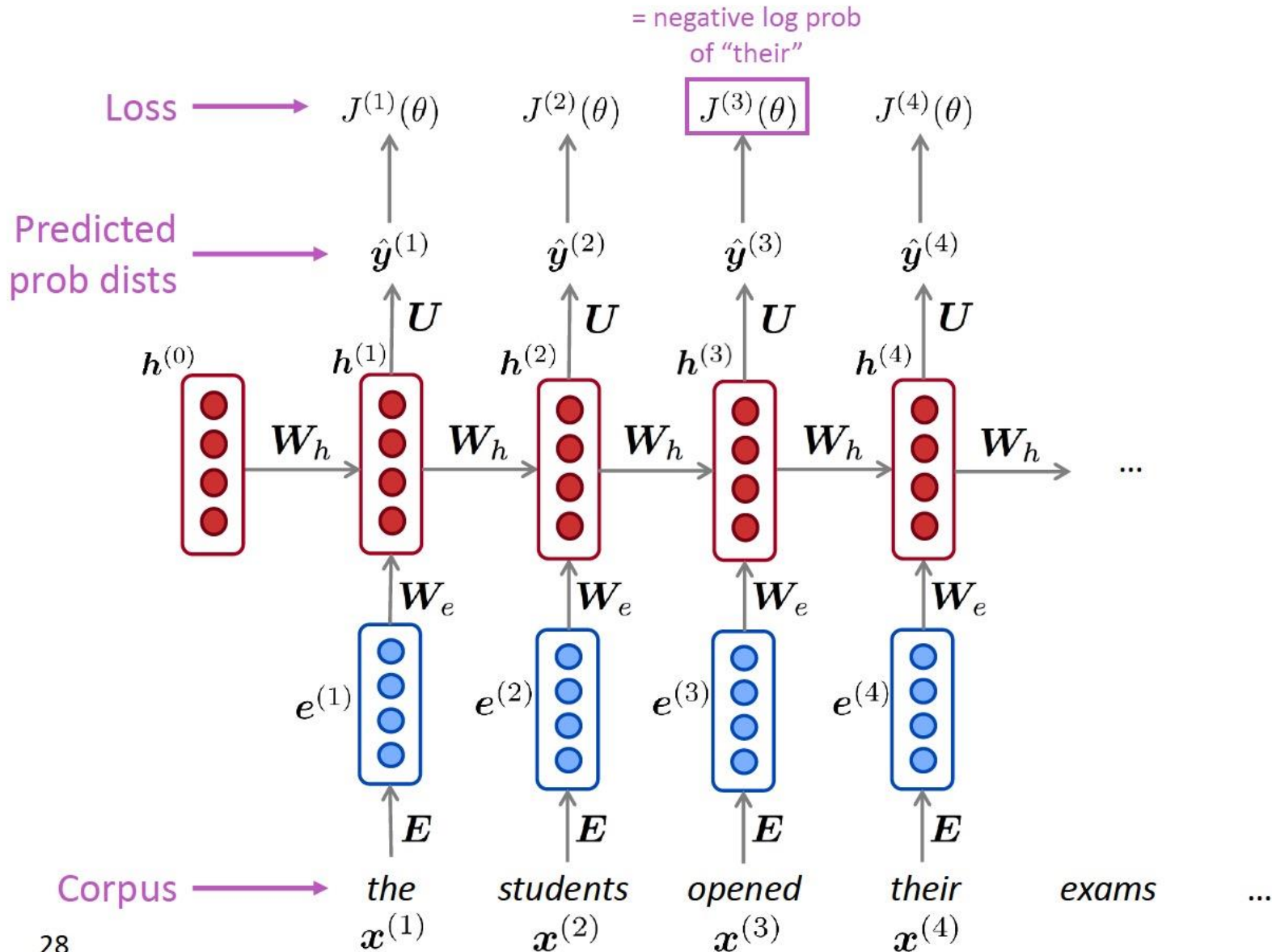# Training a RNN Language Model



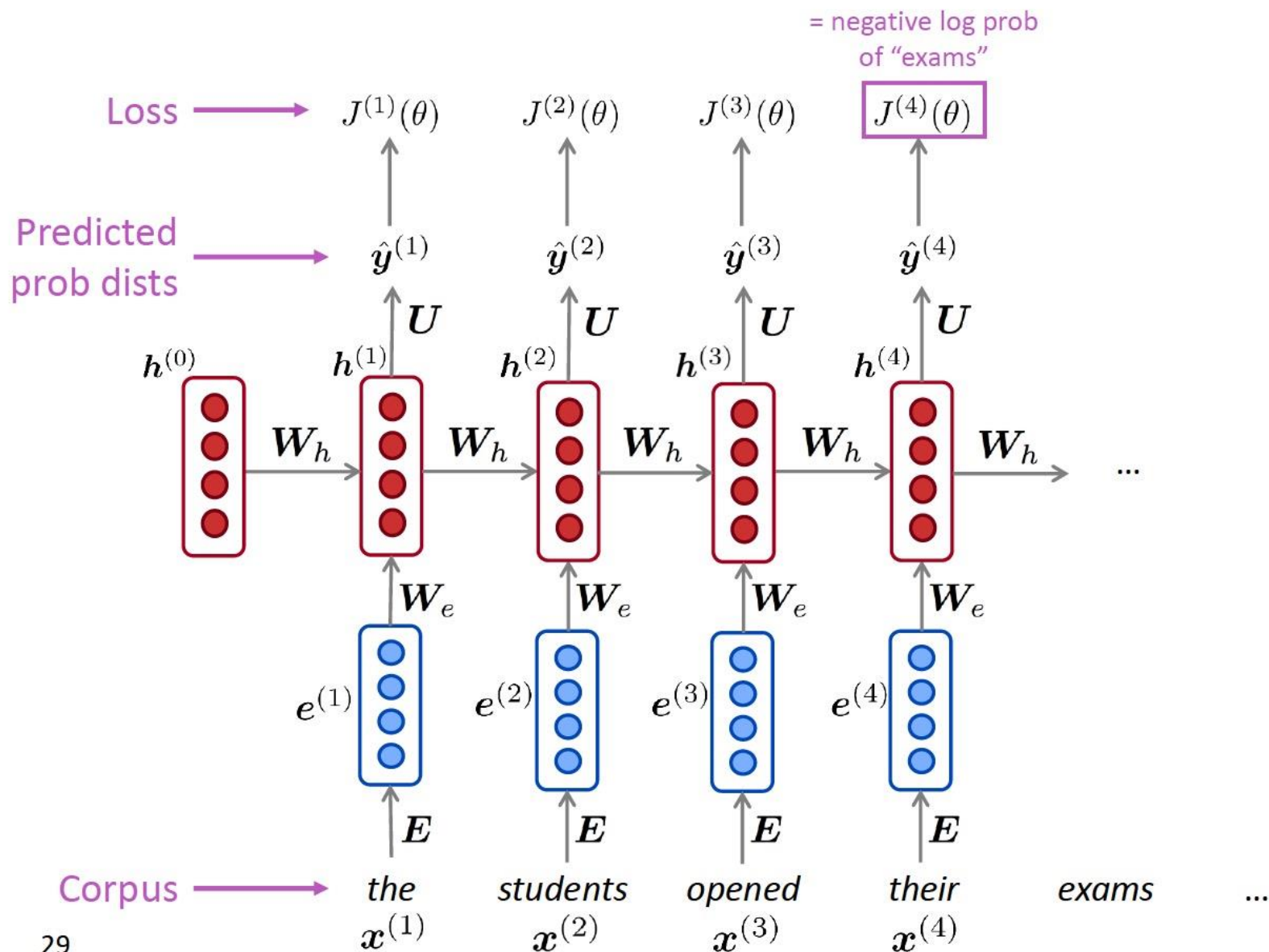= negative log prob of "students"

Loss $\longrightarrow$ $J^{(1)}(\theta)$ $\quad J^{(2)}(\theta)$ $\quad J^{(3)}(\theta)$ $\quad J^{(4)}(\theta)$

Predicted prob dists $\longrightarrow$ $\hat{y}^{(1)}$ $\quad \hat{y}^{(2)}$ $\quad \hat{y}^{(3)}$ $\quad \hat{y}^{(4)}$

$U \quad U \quad U \quad U$

$h^{(0)} \quad h^{(1)} \quad h^{(2)} \quad h^{(3)} \quad h^{(4)}$

$W_h \quad W_h \quad W_h \quad W_h \quad W_h \quad \dots$

$W_e \quad W_e \quad W_e \quad W_e$

$e^{(1)} \quad e^{(2)} \quad e^{(3)} \quad e^{(4)}$

$E \quad E \quad E \quad E$

Corpus $\longrightarrow$ *the* *students* *opened* *their* *exams* *...*

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

26

# Training a RNN Language Model



Loss → $J^{(1)}(\theta)$   $J^{(2)}(\theta)$ = negative log prob of "opened"   $J^{(3)}(\theta)$   $J^{(4)}(\theta)$

Predicted prob dists → $\hat{y}^{(1)}$   $\hat{y}^{(2)}$   $\hat{y}^{(3)}$   $\hat{y}^{(4)}$

$U$

$h^{(0)}$   $h^{(1)}$   $h^{(2)}$   $h^{(3)}$   $h^{(4)}$

$W_h$

$W_e$

$e^{(1)}$   $e^{(2)}$   $e^{(3)}$   $e^{(4)}$

$E$

Corpus → the $x^{(1)}$   students $x^{(2)}$   opened $x^{(3)}$   their $x^{(4)}$   exams   ...

27

# Training a RNN Language Model



Loss $\longrightarrow$ $J^{(1)}(\theta)$ $\quad J^{(2)}(\theta)$ $\quad J^{(3)}(\theta)$ $\quad J^{(4)}(\theta)$

= negative log prob of "their"

Predicted prob dists $\longrightarrow$ $\hat{y}^{(1)}$ $\quad \hat{y}^{(2)}$ $\quad \hat{y}^{(3)}$ $\quad \hat{y}^{(4)}$

$U$

$h^{(0)}$ $\quad h^{(1)}$ $\quad h^{(2)}$ $\quad h^{(3)}$ $\quad h^{(4)}$

$W_h$

$W_e$

$e^{(1)}$ $\quad e^{(2)}$ $\quad e^{(3)}$ $\quad e^{(4)}$

$E$

Corpus $\longrightarrow$ the $\quad$ students $\quad$ opened $\quad$ their $\quad$ exams $\quad$ ...

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

28

# Training a RNN Language Model



= negative log prob of "exams"

Loss ⟶ $J^{(1)}(\theta)$  $J^{(2)}(\theta)$  $J^{(3)}(\theta)$  $\boxed{J^{(4)}(\theta)}$

Predicted prob dists ⟶ $\hat{\boldsymbol{y}}^{(1)}$  $\hat{\boldsymbol{y}}^{(2)}$  $\hat{\boldsymbol{y}}^{(3)}$  $\hat{\boldsymbol{y}}^{(4)}$

$\boldsymbol{U}$

$\boldsymbol{h}^{(0)}$  $\boldsymbol{h}^{(1)}$  $\boldsymbol{h}^{(2)}$  $\boldsymbol{h}^{(3)}$  $\boldsymbol{h}^{(4)}$

$\boldsymbol{W}_h$  $\boldsymbol{W}_h$  $\boldsymbol{W}_h$  $\boldsymbol{W}_h$  $\boldsymbol{W}_h$  ...

$\boldsymbol{W}_e$

$\boldsymbol{e}^{(1)}$  $\boldsymbol{e}^{(2)}$  $\boldsymbol{e}^{(3)}$  $\boldsymbol{e}^{(4)}$

$\boldsymbol{E}$

Corpus ⟶ the $\boldsymbol{x}^{(1)}$  students $\boldsymbol{x}^{(2)}$  opened $\boldsymbol{x}^{(3)}$  their $\boldsymbol{x}^{(4)}$  exams  ...

29

# Training a RNN Language Model

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



Predicted prob dists $\longrightarrow$ $\hat{\boldsymbol{y}}^{(1)}$ $\hat{\boldsymbol{y}}^{(2)}$ $\hat{\boldsymbol{y}}^{(3)}$ $\hat{\boldsymbol{y}}^{(4)}$

$\boldsymbol{U}$ $\boldsymbol{U}$ $\boldsymbol{U}$ $\boldsymbol{U}$

$\boldsymbol{h}^{(0)}$ $\boldsymbol{h}^{(1)}$ $\boldsymbol{h}^{(2)}$ $\boldsymbol{h}^{(3)}$ $\boldsymbol{h}^{(4)}$

$\boldsymbol{W}_h$ $\boldsymbol{W}_h$ $\boldsymbol{W}_h$ $\boldsymbol{W}_h$ $\boldsymbol{W}_h$ ...

$\boldsymbol{W}_e$ $\boldsymbol{W}_e$ $\boldsymbol{W}_e$ $\boldsymbol{W}_e$

$\boldsymbol{e}^{(1)}$ $\boldsymbol{e}^{(2)}$ $\boldsymbol{e}^{(3)}$ $\boldsymbol{e}^{(4)}$

$\boldsymbol{E}$ $\boldsymbol{E}$ $\boldsymbol{E}$ $\boldsymbol{E}$

Corpus $\longrightarrow$ *the* *students* *opened* *their* *exams* ...
$\boldsymbol{x}^{(1)}$ $\boldsymbol{x}^{(2)}$ $\boldsymbol{x}^{(3)}$ $\boldsymbol{x}^{(4)}$

30

# Training a RNN Language Model

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

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

- In practice, consider $x^{(1)}, \ldots, x^{(T)}$ 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.

# Backpropagation for RNNs



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

**Answer:** $\dfrac{\partial J^{(t)}}{\partial \boldsymbol{W_h}} = \displaystyle\sum_{i=1}^{t} \dfrac{\partial J^{(t)}}{\partial \boldsymbol{W_h}}\bigg|_{(i)}$

"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:

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

One final output $f(x(t), y(t))$

Two intermediate outputs $x(t)$ $y(t)$

One input $t$

# 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:

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

In our example:



Apply the multivariable chain rule:

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

$= 1$

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

34

# Backpropagation for RNNs



$$\frac{\partial J^{(t)}}{\partial \boldsymbol{W_h}} = \sum_{i=1}^{t} \frac{\partial J^{(t)}}{\partial \boldsymbol{W_h}}\bigg|_{(i)}$$

**Question:** How do we calculate this?

**Answer:** Backpropagate over timesteps $i=t,\ldots,0$, summing gradients as you go.
This algorithm is called **"backpropagation through time"**

# Generating text with a RNN Language Model

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

# End Supplemental Material: Training RNNs

# Generating text with a RNN Language Model

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



*The United States will step up to the cost of a new challenges of the American people that will share the fact that we created the problem. They were attacked and so that they have to say that all the task of the final days of war that I will not be able to get this done.*

37

# Generating text with a RNN Language Model

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



"Sorry," Harry shouted, panicking—"I'll leave those brooms in London, are they?"

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

# Generating text with a RNN Language Model

- Let's have some fun!

- You can train a RNN-LM on any kind of text, then generate text in that style.

- RNN-LM trained on recipes:

```
        Title: CHOCOLATE RANCH BARBECUE
  Categories: Game, Casseroles, Cookies, Cookies
       Yield: 6 Servings

       2 tb Parmesan cheese -- chopped
       1 c  Coconut milk
       3    Eggs, beaten
```
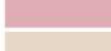
```
Place each pasta over layers of lumps. Shape mixture into the moderate oven and simmer
until firm. Serve hot in bodied fresh, mustard, orange and cheese.

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

**Source:** https://gist.github.com/nylki/1efbaa36635956d35bcc

39

# Generating text with a RNN Language Model

- Let's have some fun!

- You can train a RNN-LM on any kind of text, then generate text in that style.

- RNN-LM trained on paint color names:

| | |
|---|---|
| Ghasty Pink 231 137 165 | Sand Dan 201 172 143 |
| Power Gray 151 124 112 | Grade Bat 48 94 83 |
| Navel Tan 199 173 140 | Light Of Blast 175 150 147 |
| Bock Coe White 221 215 236 | Grass Bat 176 99 108 |
| Horble Gray 178 181 196 | Sindis Poop 204 205 194 |
| Homestar Brown 133 104 85 | Dope 219 209 179 |
| Snader Brown 144 106 74 | Testing 156 101 106 |
| Golder Craam 237 217 177 | Stoner Blue 152 165 159 |
| Hurky White 232 223 215 | Burble Simp 226 181 132 |
| Burf Pink 223 173 179 | Stanky Bean 197 162 171 |
| Rose Hork 230 215 198 | Turdly 190 164 116 |

This is an example of a character-level RNN-LM (predicts what character comes next)

40

# Evaluating Language Models

- The standard evaluation metric for Language Models is perplexity.

$$\text{perplexity} = \prod_{t=1}^{T} \left( \frac{1}{P_{\text{LM}}(\boldsymbol{x}^{(t+1)} \mid \boldsymbol{x}^{(t)}, \ldots, \boldsymbol{x}^{(1)})} \right)^{1/T}$$

Normalized by number of words

Inverse probability of corpus, according to Language Model

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

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

**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)

42

# Why should we care about Language Modeling?
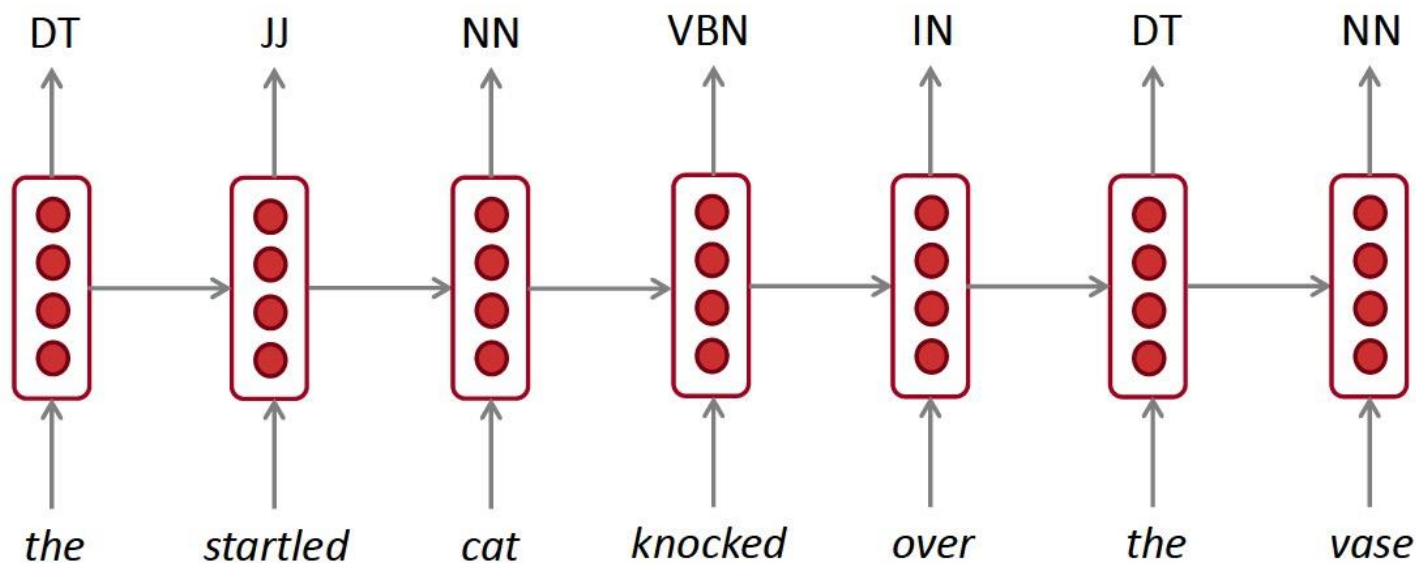
- Language Modeling is a benchmark task that helps us measure our progress on understanding language

- 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.

# Recap

- **Language Model**: A system that predicts the next word

- **Recurrent Neural Network**: A family of neural networks that:
  - Take sequential input of any length
  - Apply the same weights on each step
  - Can optionally produce output on each step

- Recurrent Neural Network ≠ Language Model

- We've shown that RNNs are a great way to build a LM.

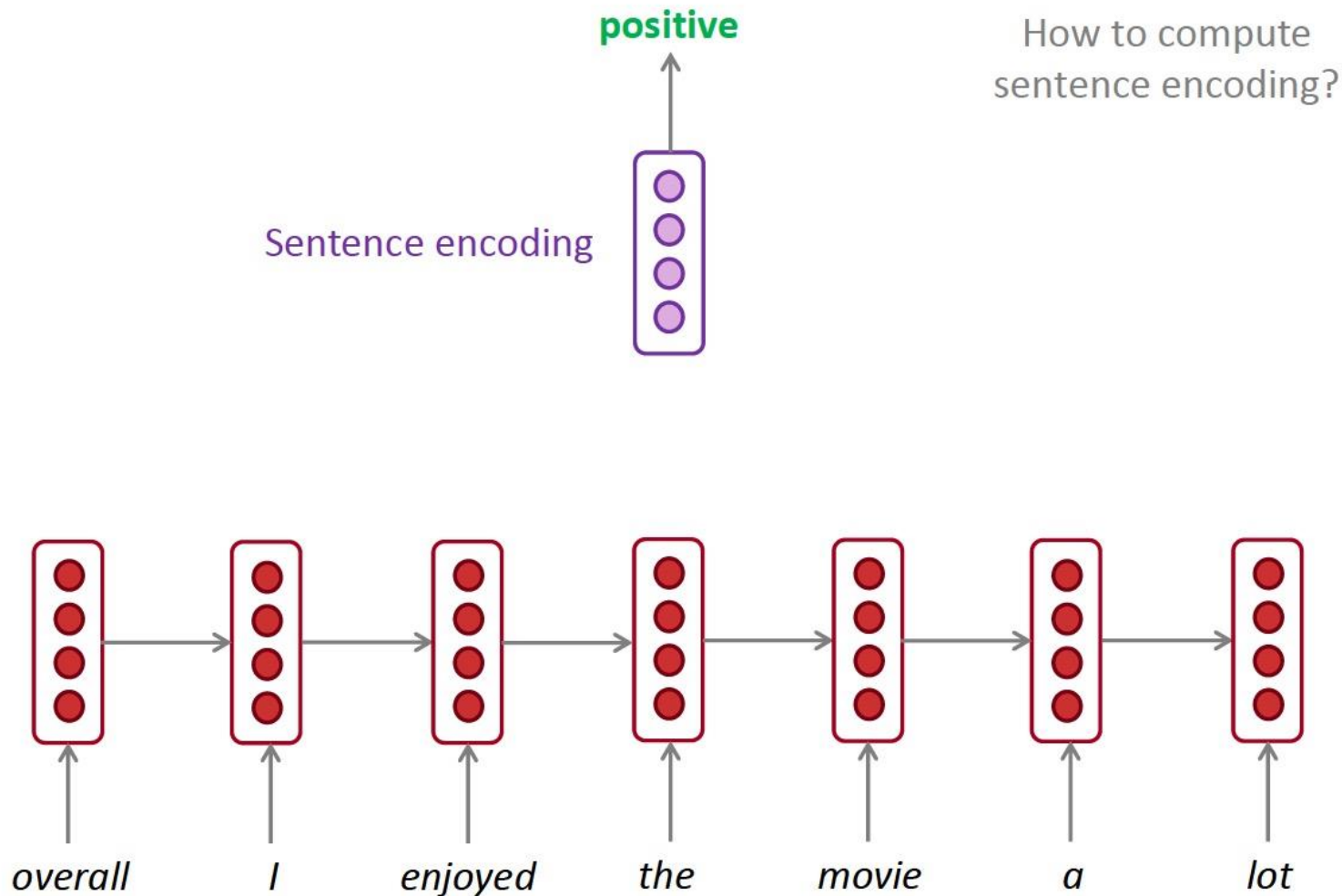- But RNNs are useful for much more!

# RNNs can be used for tagging

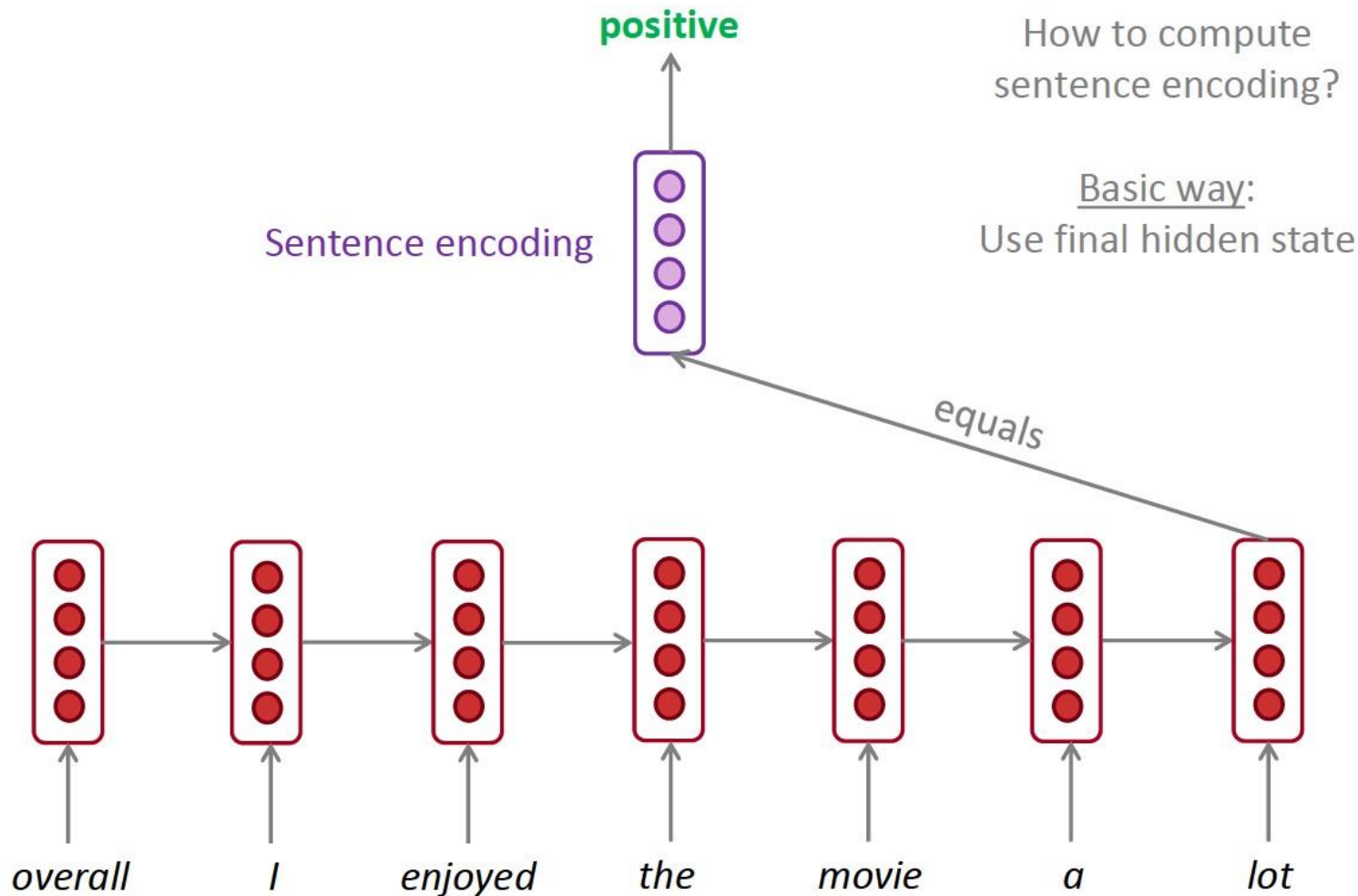e.g. part-of-speech tagging, named entity recognition

# RNNs can be used for sentence classification
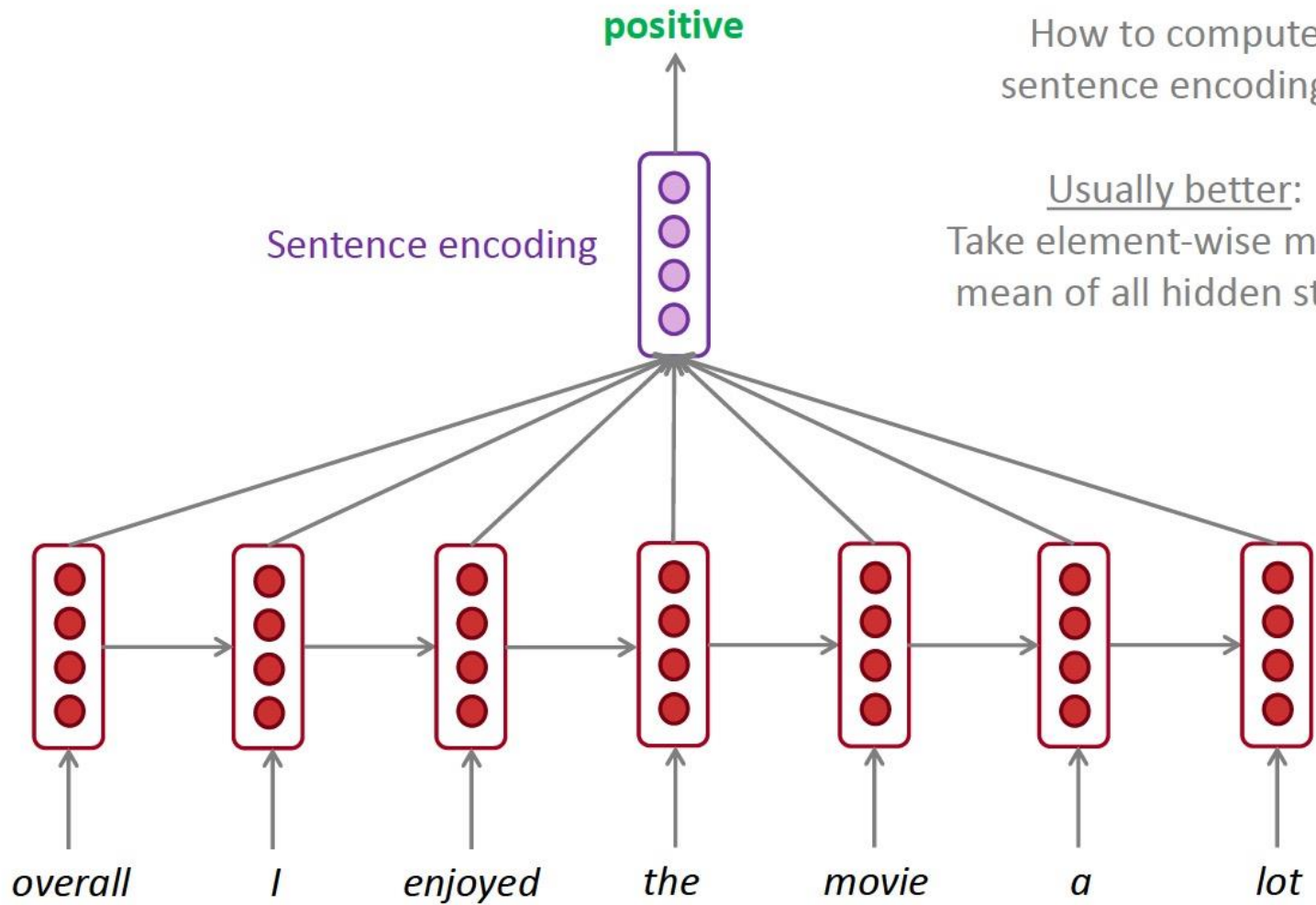
e.g. sentiment classification

**positive**

How to compute
sentence encoding?

Sentence encoding

overall    I    enjoyed    the    movie    a    lot

# RNNs can be used for sentence classification

e.g. sentiment classification

**positive**

Sentence encoding

How to compute sentence encoding?

Basic way:
Use final hidden state

*equals*

overall    I    enjoyed    the    movie    a    lot

# RNNs can be used for sentence classification

e.g. sentiment classification

**positive**

How to compute
sentence encoding?

Sentence encoding

Usually better:
Take element-wise max or
mean of all hidden states

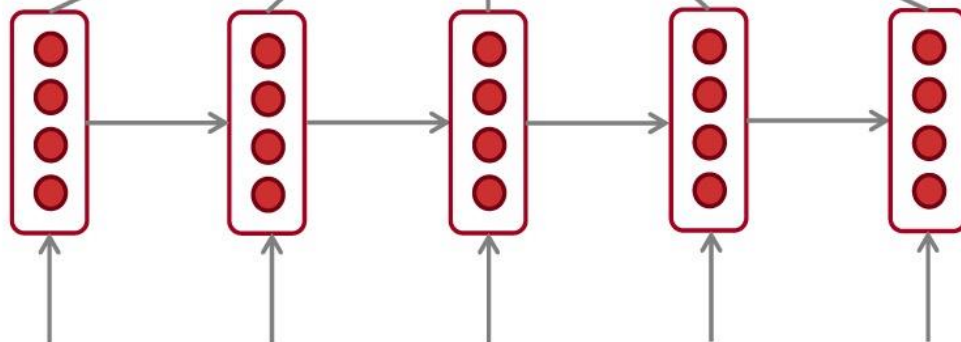*overall*  *I*  *enjoyed*  *the*  *movie*  *a*  *lot*

# RNNs can be used as an encoder module

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

**Answer:** German

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.

*lots of neural architecture*

*lots of neural architecture*

**Context:** *Ludwig van Beethoven was a German composer and pianist. A crucial figure ...*



**Question:** *what    nationality    was    Beethoven    ?*

49

# RNN-LMs can be used to generate text

e.g. speech recognition, machine translation, summarization

# A note on terminology

RNN described in this lecture = "vanilla RNN"

**Next lecture:** You will learn about other RNN flavors

like GRU      and LSTM      and multi-layer RNNs

**By the end of the course:** You will understand phrases like
*"stacked bidirectional LSTM with residual connections and self-attention"*
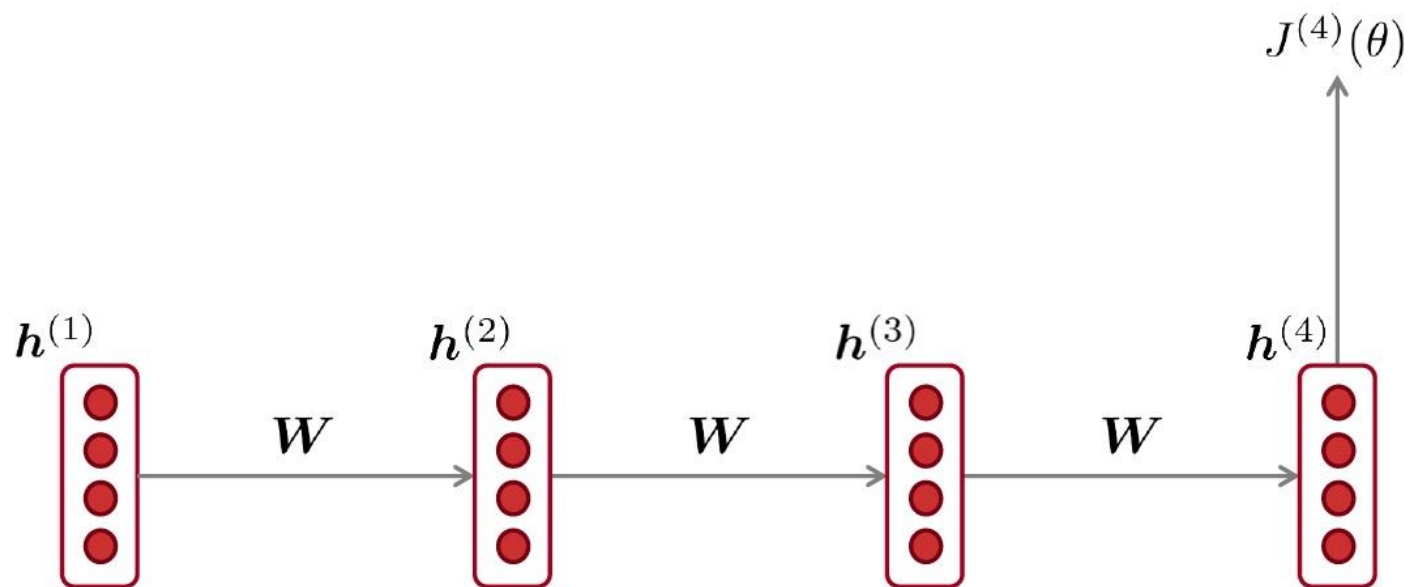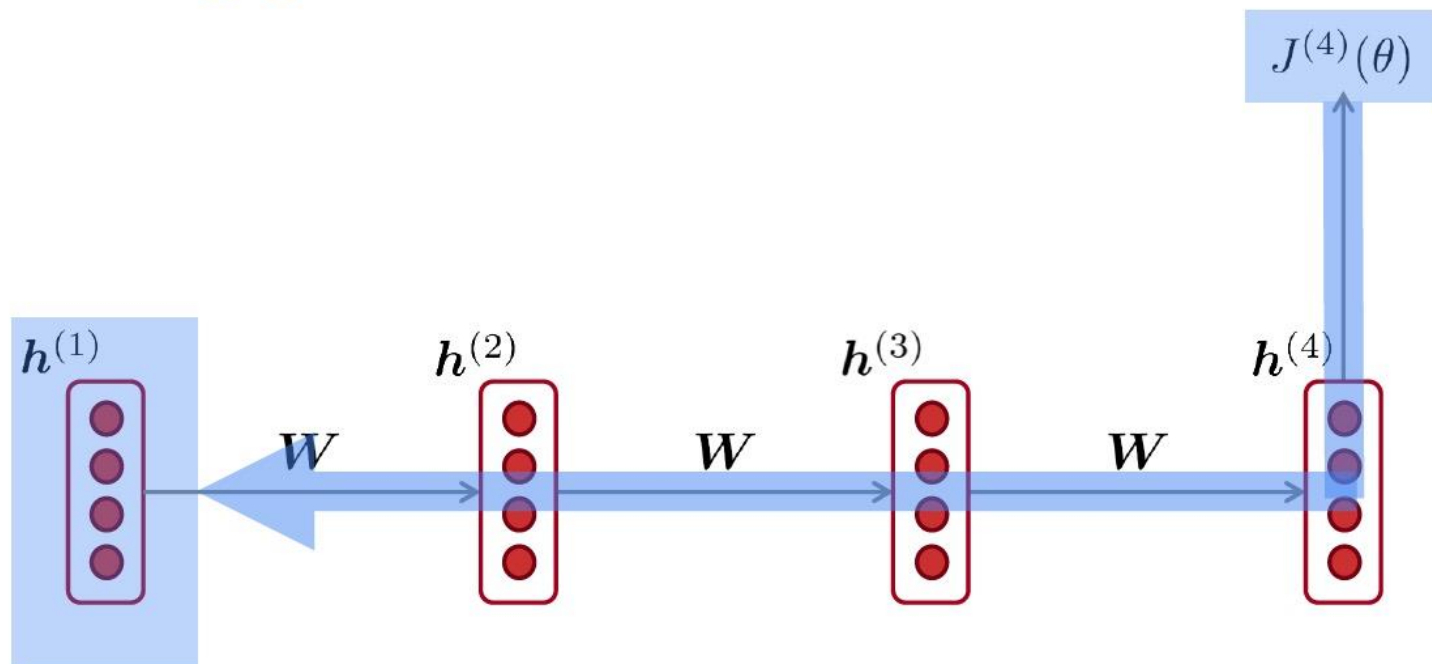
51

- **Problems** with RNNs!
  - Vanishing gradients

    **motivates**

- **Fancy RNN** variants!
  - LSTM
  - GRU
  - multi-layer
  - bidirectional

# Vanishing gradient intuition

# Vanishing gradient intuition



$J^{(4)}(\theta)$

$\boldsymbol{h}^{(1)}$    $\boldsymbol{h}^{(2)}$    $\boldsymbol{h}^{(3)}$    $\boldsymbol{h}^{(4)}$

$\boldsymbol{W}$    $\boldsymbol{W}$    $\boldsymbol{W}$

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

# Vanishing gradient intuition



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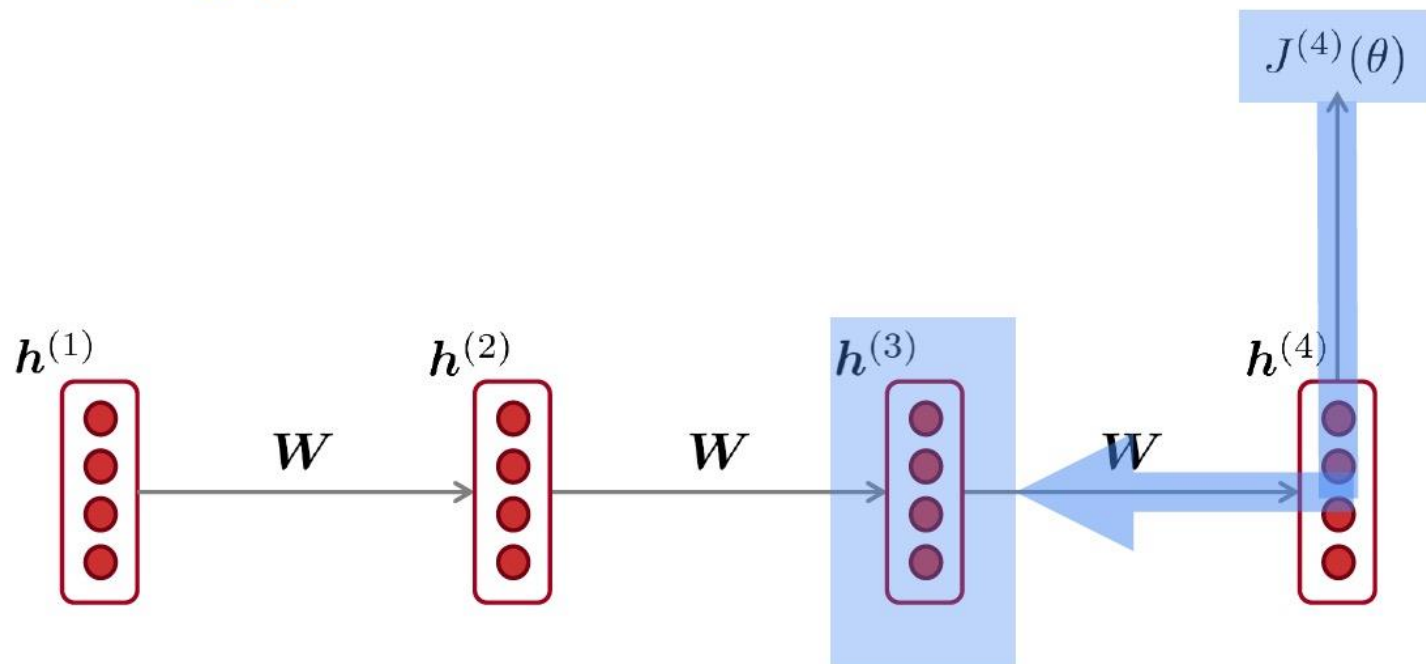
chain rule!

8

# Vanishing gradient intuition



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

chain rule!

9

# Vanishing gradient intuition



$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = \quad \frac{\partial h^{(2)}}{\partial h^{(1)}} \times \qquad\qquad \frac{\partial h^{(3)}}{\partial h^{(2)}} \times \qquad\qquad \frac{\partial h^{(4)}}{\partial h^{(3)}} \times \quad \frac{\partial J^{(4)}}{\partial h^{(4)}}$$
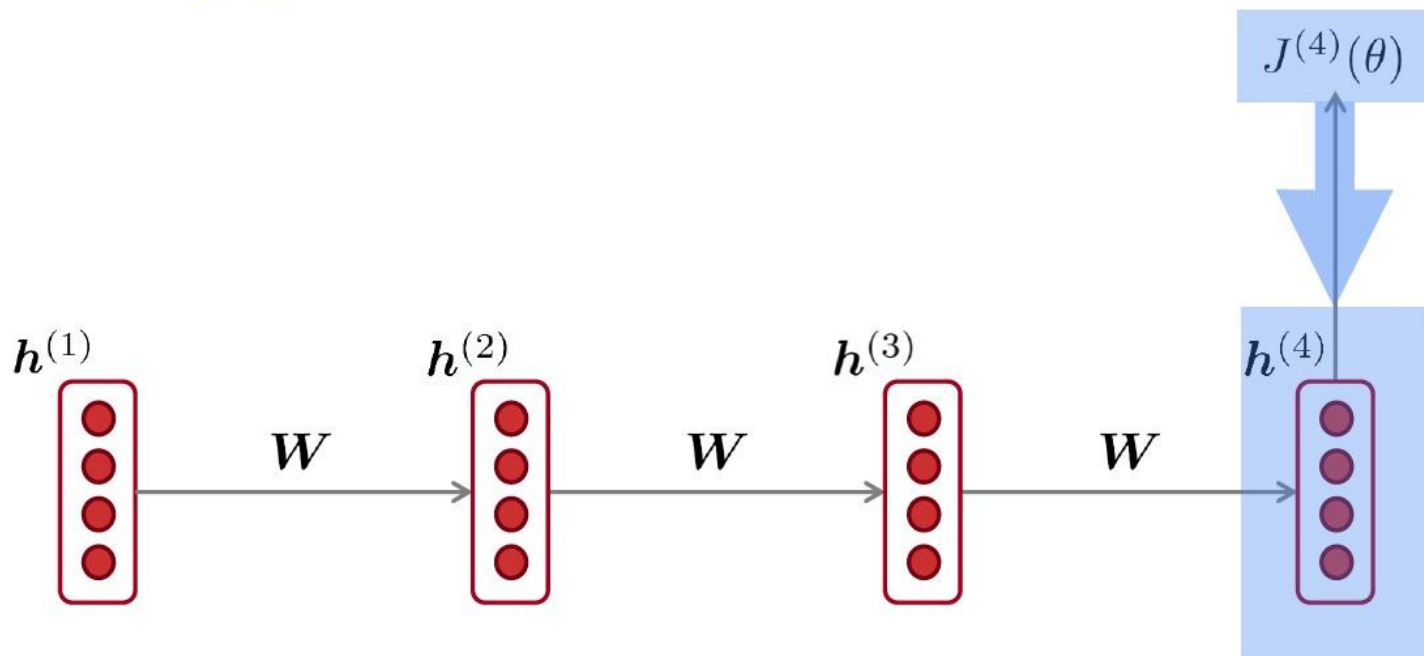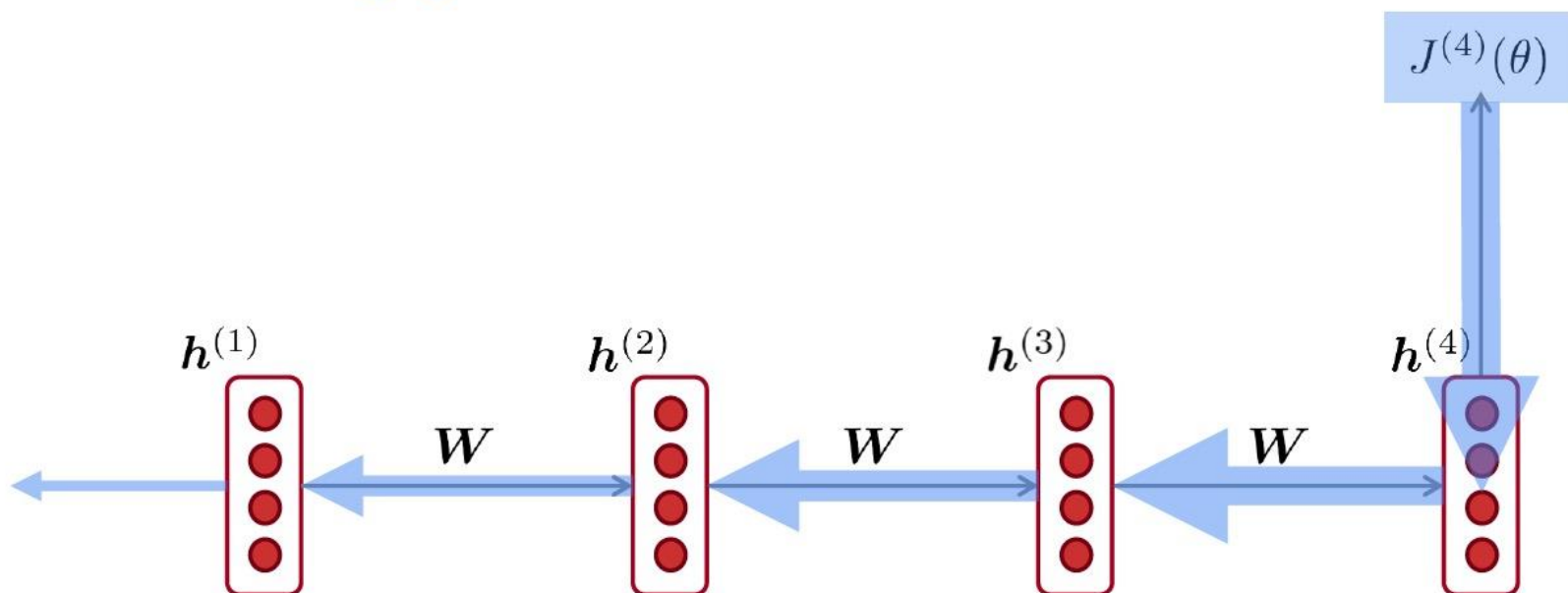
chain rule!

# Vanishing gradient intuition



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

What happens if these are small?

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

11

# Why is vanishing gradient a problem?



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

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

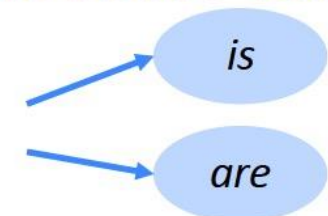# Why is vanishing gradient a problem?

- 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 $t+n$

# 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 $7^{th}$ 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

# Effect of vanishing gradient on RNN-LM

- **LM task:** *The writer of the books ___*

  *is*

  *are*

- **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)

- Due to vanishing gradient, RNN-LMs are better at learning from sequential recency than syntactic recency, so they make this type of error more often than we'd like [Linzen et al 2016]

"Assessing the Ability of LSTMs to Learn Syntax-Sensitive Dependencies", Linzen et al, 2016. https://arxiv.org/pdf/1611.01368.pdf

# 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

$$h^{(t)} = \sigma \left( W_h h^{(t-1)} + W_x x^{(t)} + b \right)$$

- How about a RNN with separate memory?

# 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

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

# Long Short-Term Memory (LSTM)

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

**Forget gate:** controls what is kept vs forgotten, from previous cell state

**Input gate:** controls what parts of the new cell content are written to cell

**Output gate:** controls what parts of cell are output to hidden state

**Sigmoid function:** all gate values are between 0 and 1

$$\boldsymbol{f}^{(t)} = \sigma\left(\boldsymbol{W}_f \boldsymbol{h}^{(t-1)} + \boldsymbol{U}_f \boldsymbol{x}^{(t)} + \boldsymbol{b}_f\right)$$

$$\boldsymbol{i}^{(t)} = \sigma\left(\boldsymbol{W}_i \boldsymbol{h}^{(t-1)} + \boldsymbol{U}_i \boldsymbol{x}^{(t)} + \boldsymbol{b}_i\right)$$

$$\boldsymbol{o}^{(t)} = \sigma\left(\boldsymbol{W}_o \boldsymbol{h}^{(t-1)} + \boldsymbol{U}_o \boldsymbol{x}^{(t)} + \boldsymbol{b}_o\right)$$

**New cell content:** this is the new content to be written to the cell

**Cell state:** erase ("forget") some content from last cell state, and write ("input") some new cell content

**Hidden state:** read ("output") some content from the cell

$$\tilde{\boldsymbol{c}}^{(t)} = \tanh\left(\boldsymbol{W}_c \boldsymbol{h}^{(t-1)} + \boldsymbol{U}_c \boldsymbol{x}^{(t)} + \boldsymbol{b}_c\right)$$

$$\boldsymbol{c}^{(t)} = \boldsymbol{f}^{(t)} \circ \boldsymbol{c}^{(t-1)} + \boldsymbol{i}^{(t)} \circ \tilde{\boldsymbol{c}}^{(t)}$$
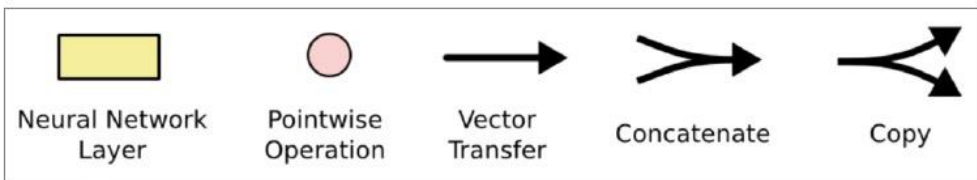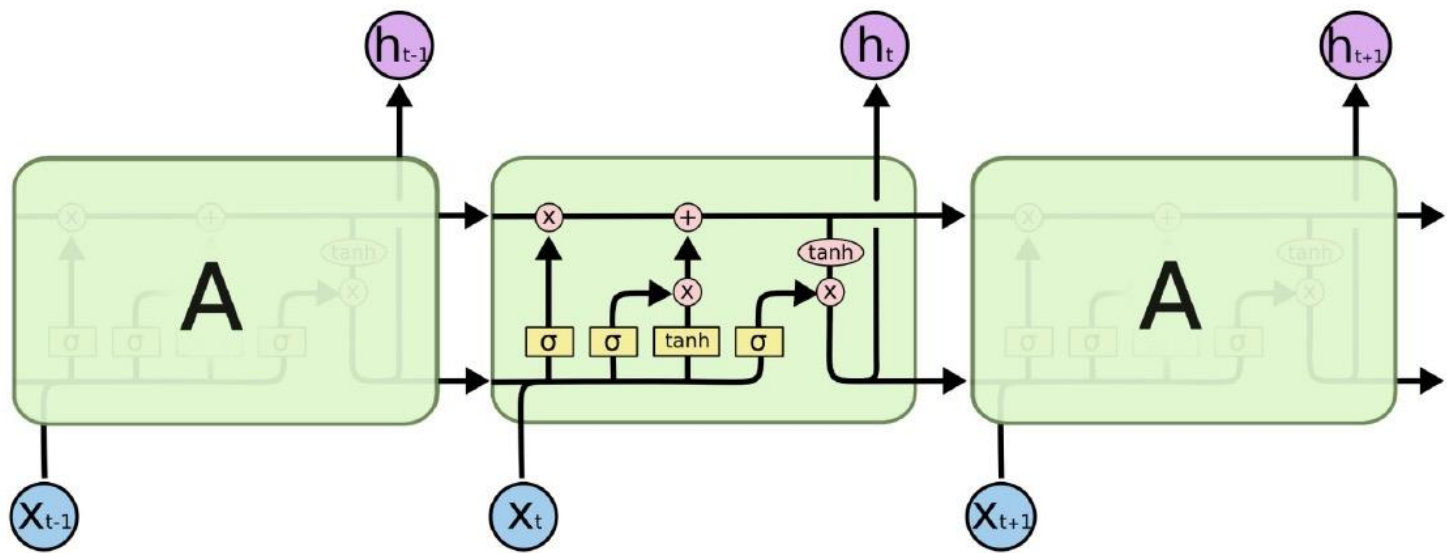
$$\boldsymbol{h}^{(t)} = \boldsymbol{o}^{(t)} \circ \tanh \boldsymbol{c}^{(t)}$$

All these are vectors of same length $n$

Gates are applied using element-wise product

23

# Long Short-Term Memory (LSTM)

You can think of the LSTM equations visually like this:



| | | | | |
|---|---|---|---|---|
| Neural Network Layer | Pointwise Operation | Vector Transfer | Concatenate | Copy |

24

# Long Short-Term Memory (LSTM)

You can think of the LSTM equations visually like this:

# 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 remember everything on every timestep, then the info in the 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

# 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 (2019), 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

**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

# 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 $\boldsymbol{x}^{(t)}$ and hidden state $\boldsymbol{h}^{(t)}$ (no cell state).

**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

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

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

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

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

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

**Hidden state:** update gate simultaneously controls what is kept from previous hidden state, and what is updated to new hidden state content
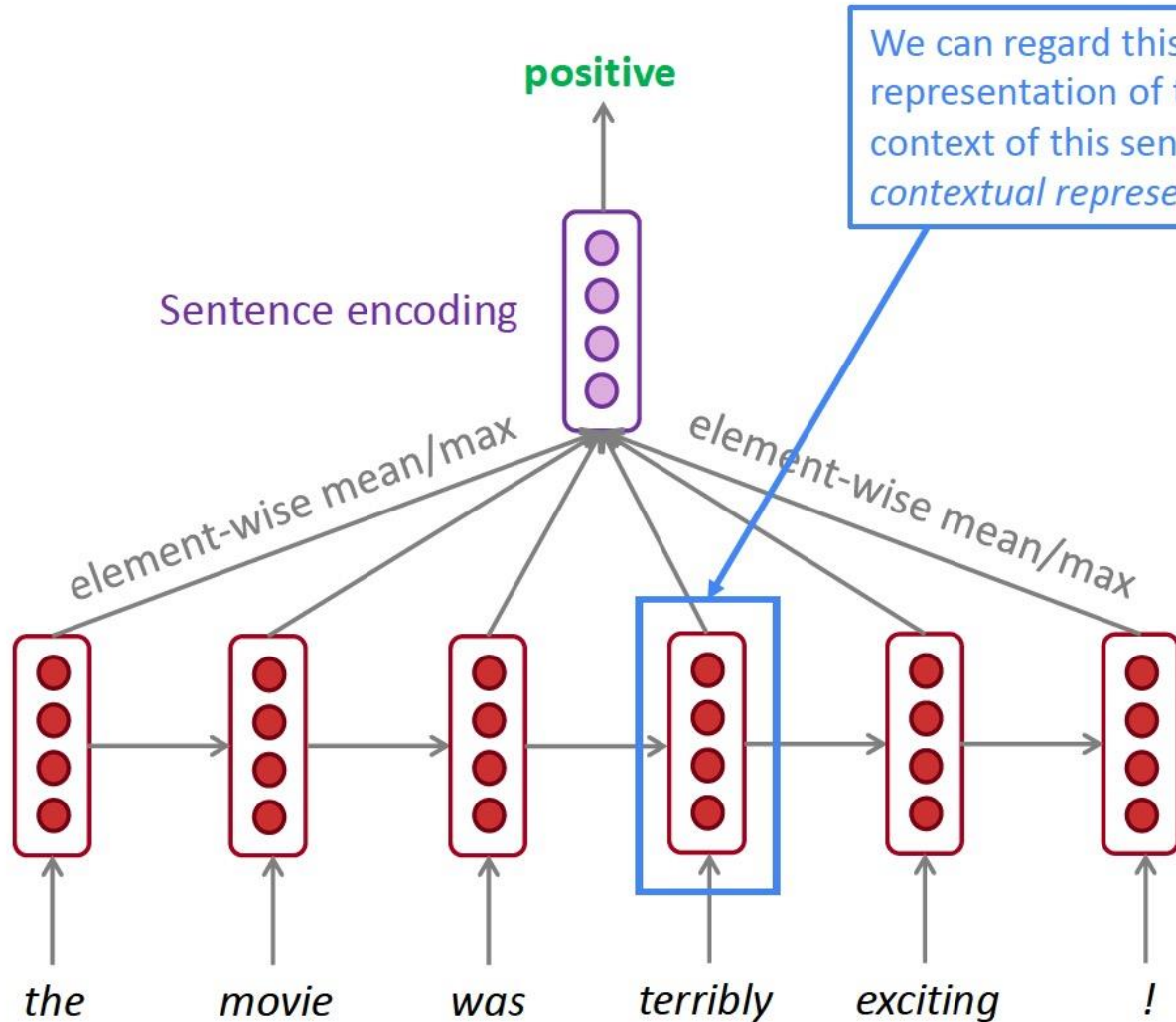
**How does this solve vanishing gradient?**
Like LSTM, GRU makes it easier to retain info long-term (e.g. by setting update gate to 0)

"Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation", Cho et al. 2014, https://arxiv.org/pdf/1406.1078v3.pdf

# LSTM vs GRU

- Researchers have proposed many gated RNN variants, but LSTM and GRU are the most widely-used

- The biggest difference is that GRU is quicker to compute and has fewer parameters

- There is no conclusive evidence that one consistently performs better than the other

- LSTM is a good default choice (especially if your data has particularly long dependencies, or you have lots of training data)

- Rule of thumb: start with LSTM, but switch to GRU if you want something more efficient

# Bidirectional RNNs: motivation

Task: Sentiment Classification



positive

Sentence encoding

element-wise mean/max

element-wise mean/max

the    movie    was    terribly    exciting    !

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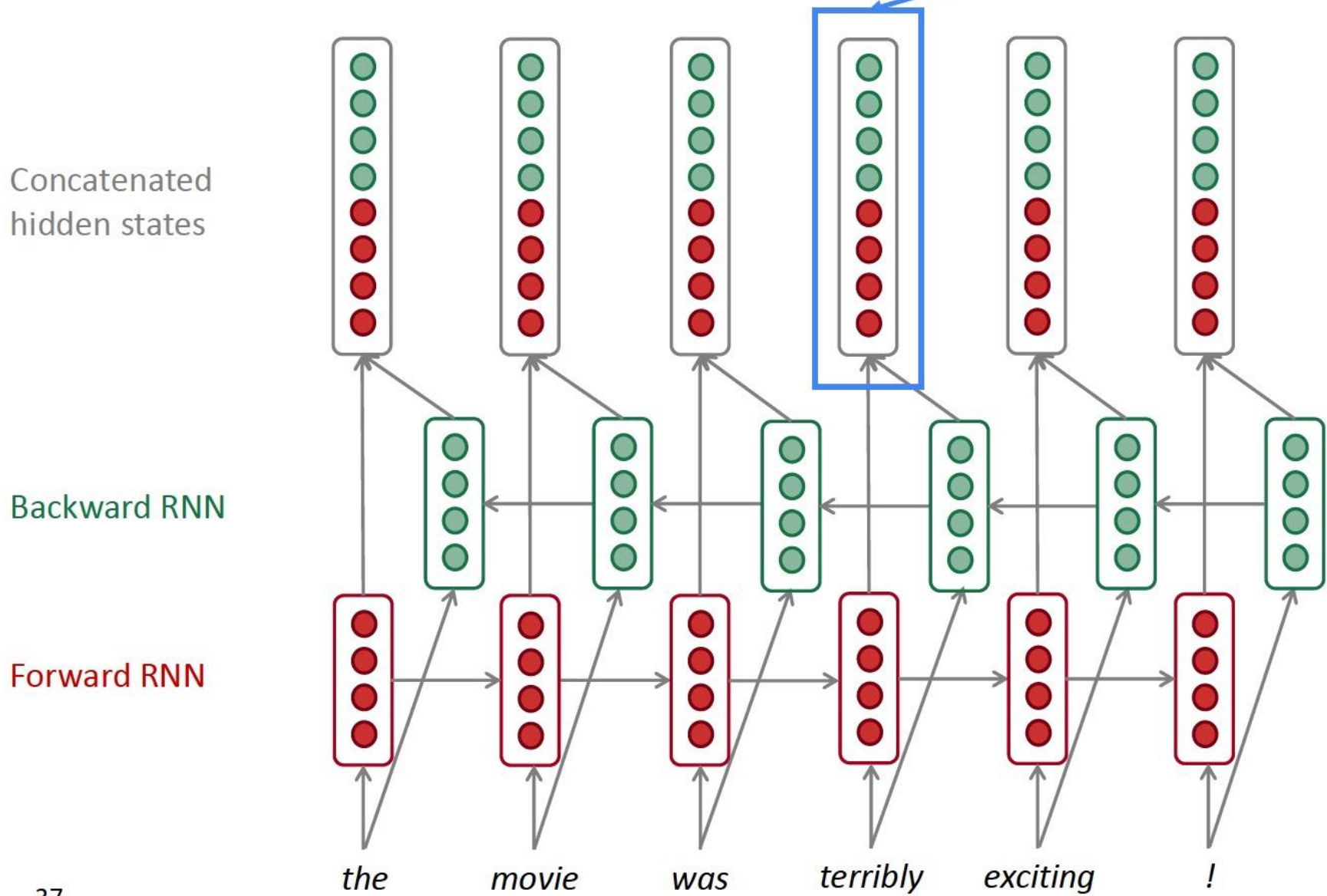
These contextual representations only contain information about the *left* context (e.g. *"the movie was"*).

**What about *right* context?**

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

36

# Bidirectional RNNs

This contextual representation of "terribly" has both left and right context!

Concatenated hidden states

Backward RNN

Forward RNN

the    movie    was    terribly    exciting    !

37

# 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 $\qquad \overrightarrow{\boldsymbol{h}}^{(t)} = \boxed{\text{RNN}_{\text{FW}}}(\overrightarrow{\boldsymbol{h}}^{(t-1)}, \boldsymbol{x}^{(t)})$

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

Generally, these two RNNs have separate weights

Concatenated hidden states $\qquad \boxed{\boldsymbol{h}^{(t)}} = [\overrightarrow{\boldsymbol{h}}^{(t)}; \overleftarrow{\boldsymbol{h}}^{(t)}]$

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

38

# Bidirectional RNNs: simplified diagram



The two-way arrows indicate bidirectionality and the depicted hidden states are assumed to be the concatenated forwards+backwards states.

# 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.
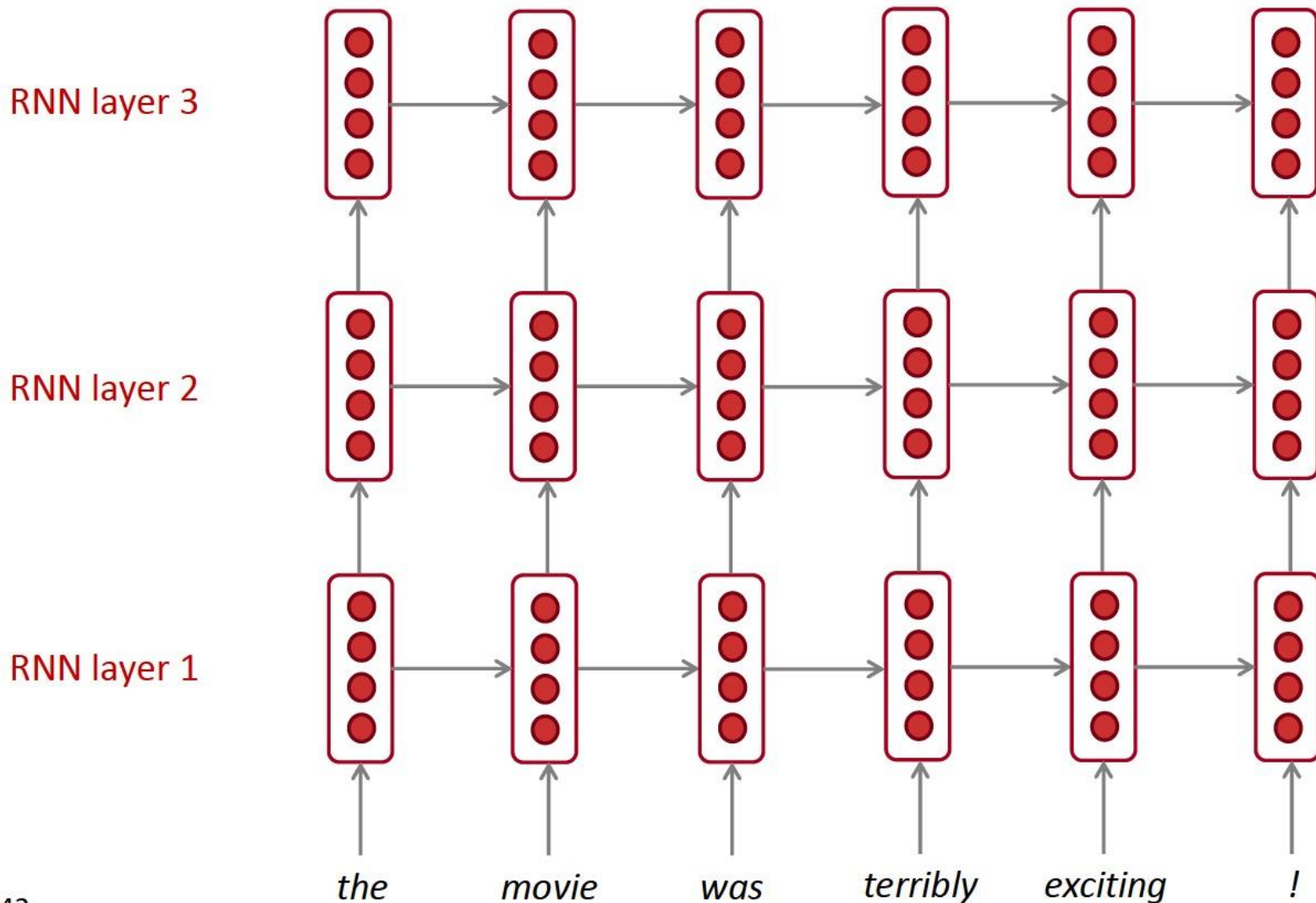  - You will learn more about BERT later in the course!

# 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

RNN layer 3

RNN layer 2

RNN layer 1

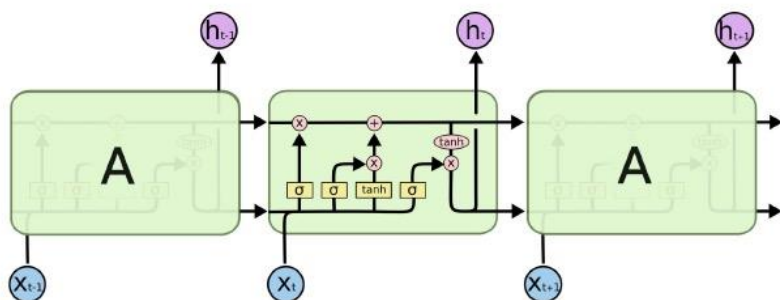the    movie    was    terribly    exciting    !

42
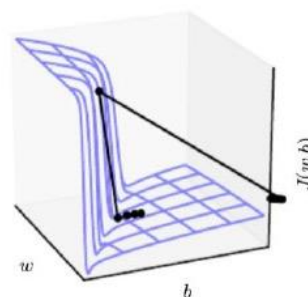
# Multi-layer RNNs in practice

- 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) can be up to 24 layers
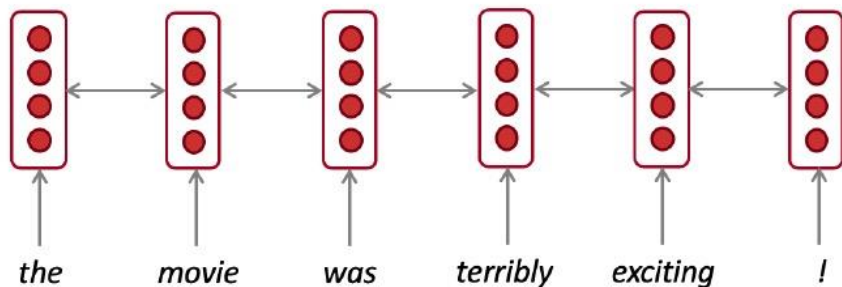  - You will learn about Transformers later; they have a lot of skipping-like connections

"Massive Exploration of Neural Machine Translation Architecutres", Britz et al, 2017. https://arxiv.org/pdf/1703.03906.pdf

# In summary

Lots of new information today! What are the practical takeaways?



1. LSTMs are powerful but GRUs are faster



2. Clip your gradients



*the*   *movie*   *was*   *terribly*   *exciting*   *!*

3. Use bidirectionality when possible



4. Multi-layer RNNs are powerful, but you might need skip/dense-connections if it's deep

44

# Required Readings

- [Chapter 7 on Neural Networks](#) from the J&M textbook.
  - Sections on Training are optional.
- [Chapter 8 on RNNs and LSTMs](#) from the J&M textbook.
  - Sections on Training are optional.