



Carnegie Mellon University

Language  
Technologies  
Institute

# 11-411/11-611 Natural Language Processing

## Neural Language Models

---

David R. Mortensen

October 4, 2022

Language Technologies Institute

# Learning Objectives

**At the end of this lecture, you should be able to do the following things:**

- Generalize the notion of language models beyond the ngram models that we have discussed thus far
- State four advantages of neural LMs over ngram LMs
- State one disadvantage of NLMs (runtime efficiency)
- Apply your knowledge of neural nets (feedforward neural nets, RNNs,

LSTMs) to language modeling architectures

- Explain how a feedforward neural language model is applied and trained
- Explain the relationship between embeddings and LMs
- Explain the advantages of RNNs over feedforward neural nets for language modeling
- Explain the advantage of LSTMs over RNNs
- Describe masked language models

## Units—Building Blocks of Neural Networks

---

## The Fundamentals

The fundamental equation that describes a unit of a neural network should look very familiar:

$$z = b + \sum_i w_i x_i \quad (1)$$

Which we will represent as

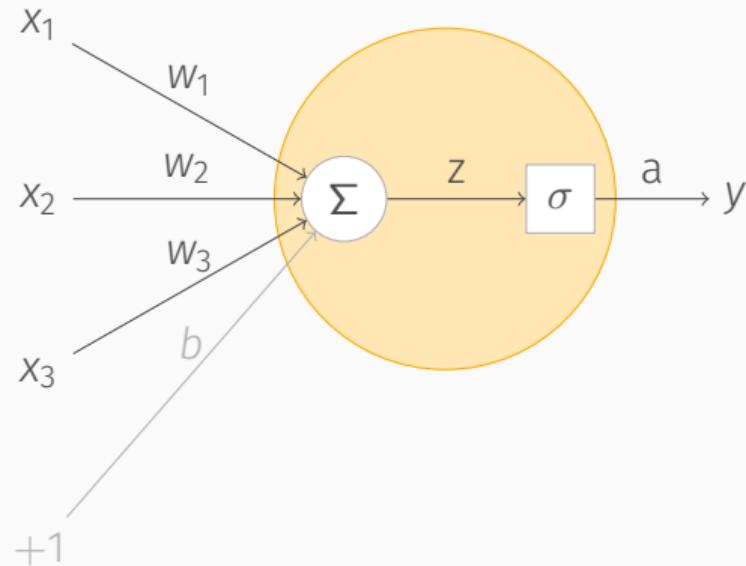
$$z = \mathbf{w} \cdot \mathbf{x} + b \quad (2)$$

But we do not use  $z$  directly. Instead, we pass it through a non-linear function, like the sigmoid function:

$$y = \sigma(z) = \frac{1}{1 + e^{-z}} \quad (3)$$

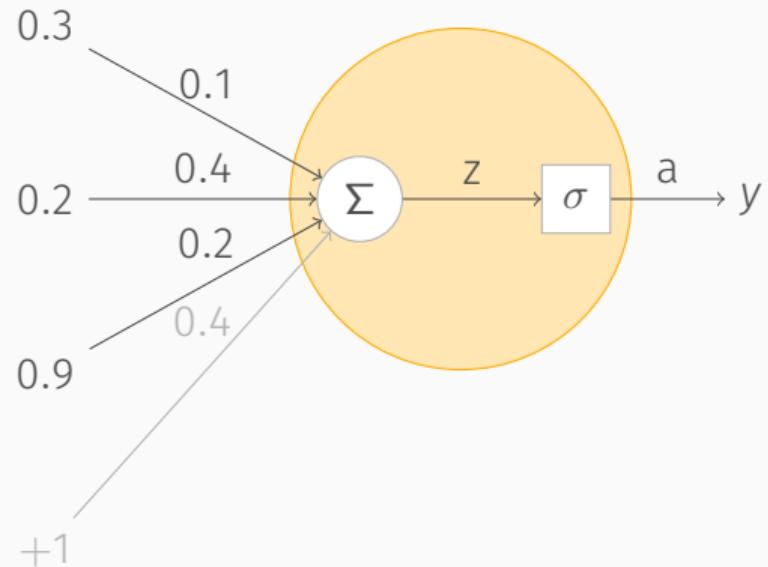
(which has some nice properties even though, in practice, we will prefer other functions like tanh and ReLU).

## A Unit Illustrated

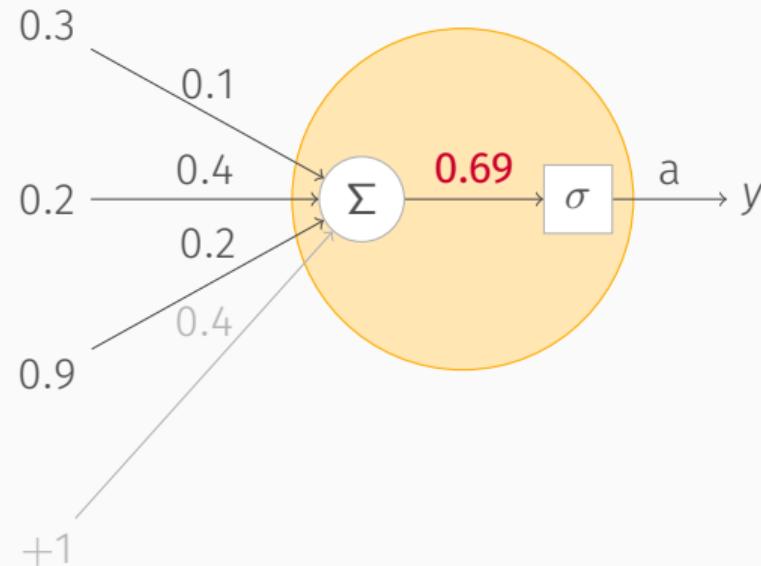


Take, for example, a scenario in which our unit has the weights [0.1, 0.4, 0.2] and the bias term 0.4 and the input vector  $x$  has the values [0.3, 0.2, 0.9].

## Filling in the Input Values and Weights

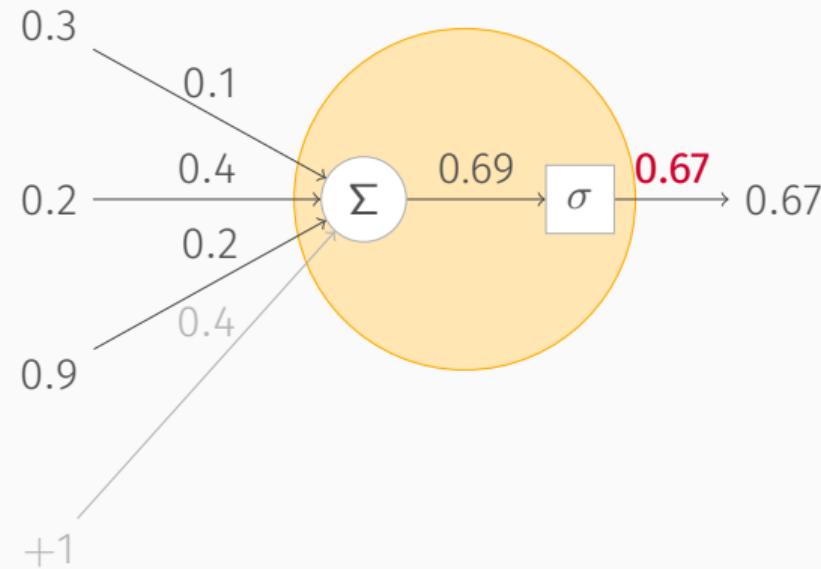


## Multiplying the Input Values and Weights and Summing Them (with the Bias Term)



$$z = x_1w_1 + x_2w_2 + x_3w_3 + b = 0.1(0.3) + 0.4(0.2) + 0.2(0.9) + 0.4 = 0.69 \quad (4)$$

## Applying the Activation Function (Sigmoid)



$$y = \sigma(0.69) = \frac{1}{1 + e^{-0.69}} = 0.67 \quad (5)$$

## Non-Linear Activation Functions

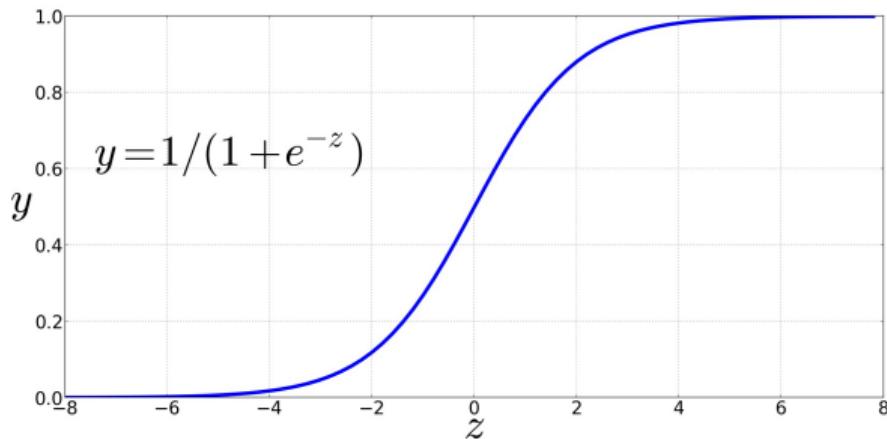
---

# Non-Linear Activation Functions

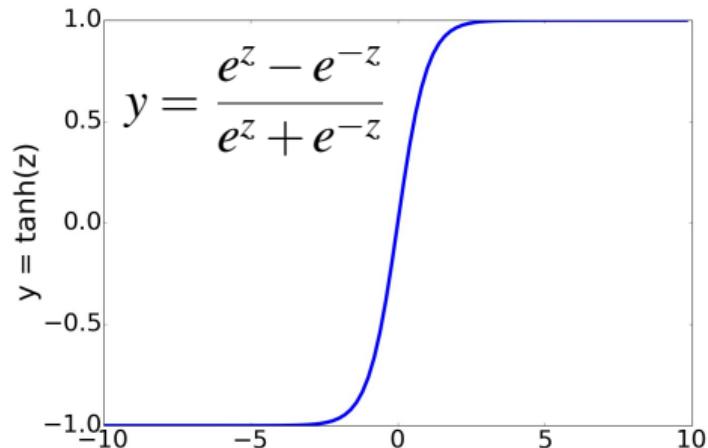
We're already seen the sigmoid for logistic regression:

Sigmoid

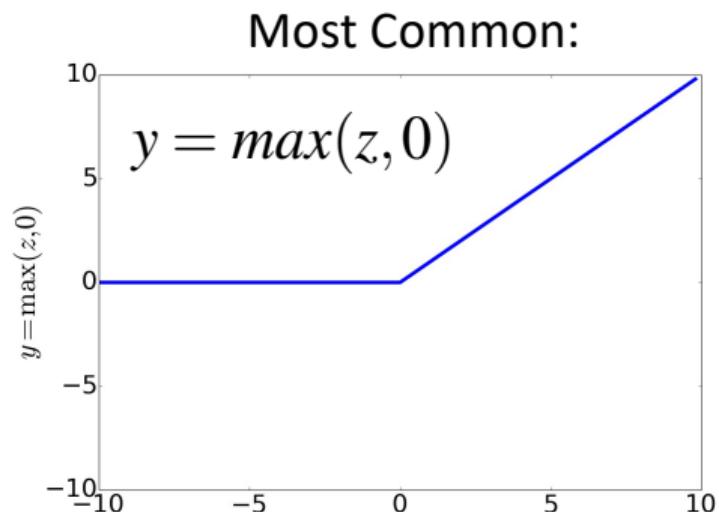
$$y = \sigma(z) = \frac{1}{1 + e^{-z}}$$



# Non-Linear Activation Functions besides sigmoid



tanh



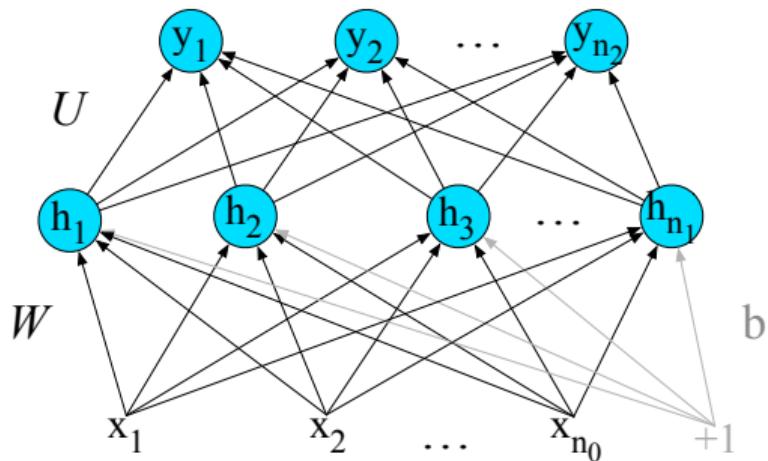
ReLU  
Rectified Linear Unit

## Feedforward Neural Networks

---

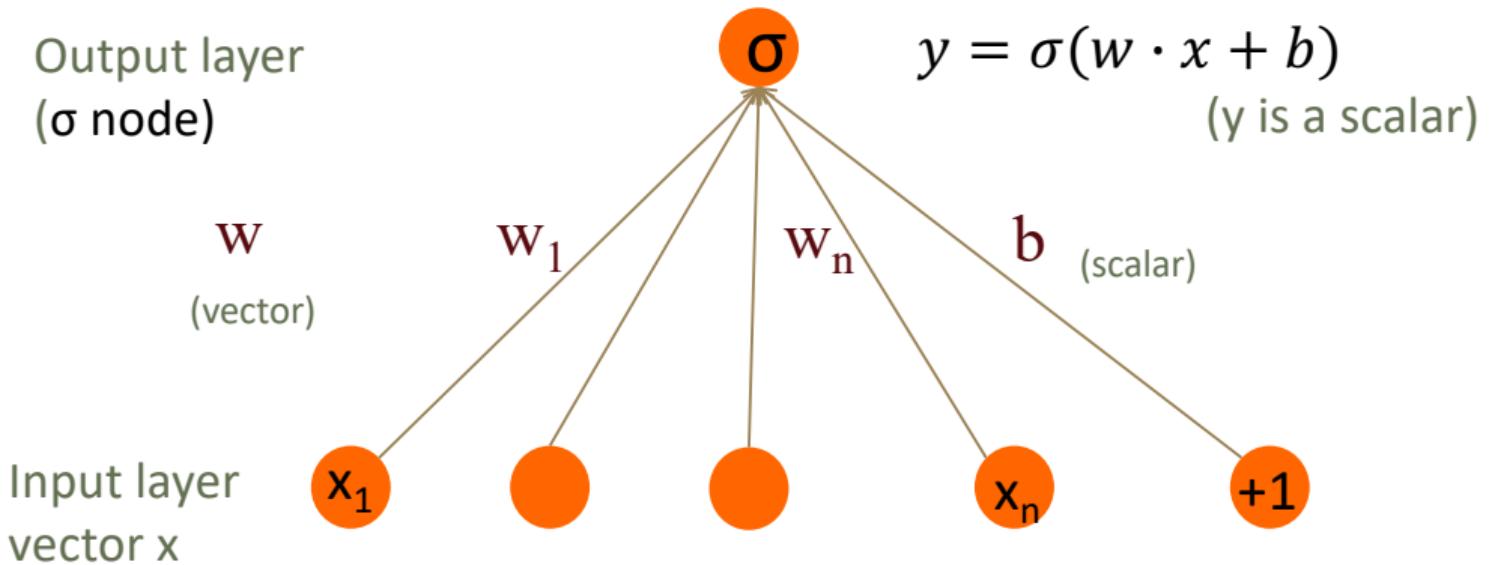
# Feedforward Neural Networks

Can also be called **multi-layer perceptrons** (or **MLPs**) for historical reasons



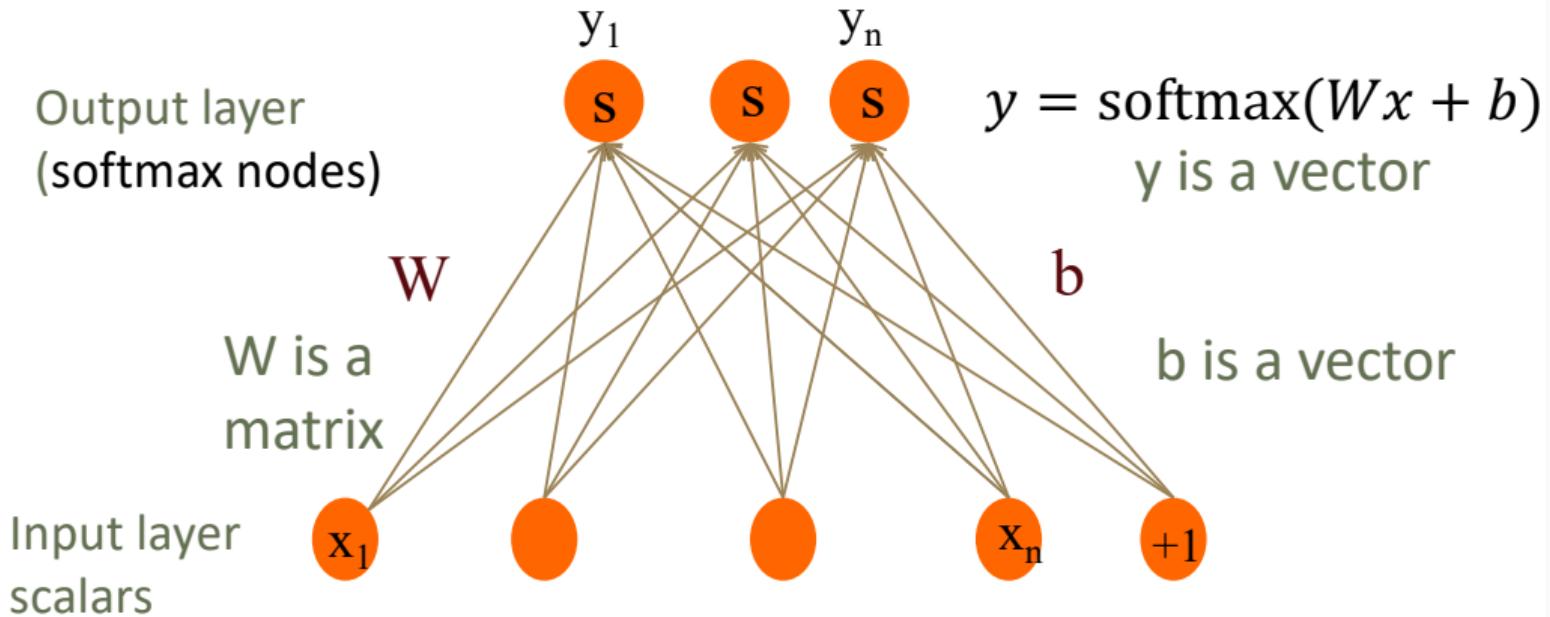
# Binary Logistic Regression as a 1-layer Network

(we don't count the input layer in counting layers!)

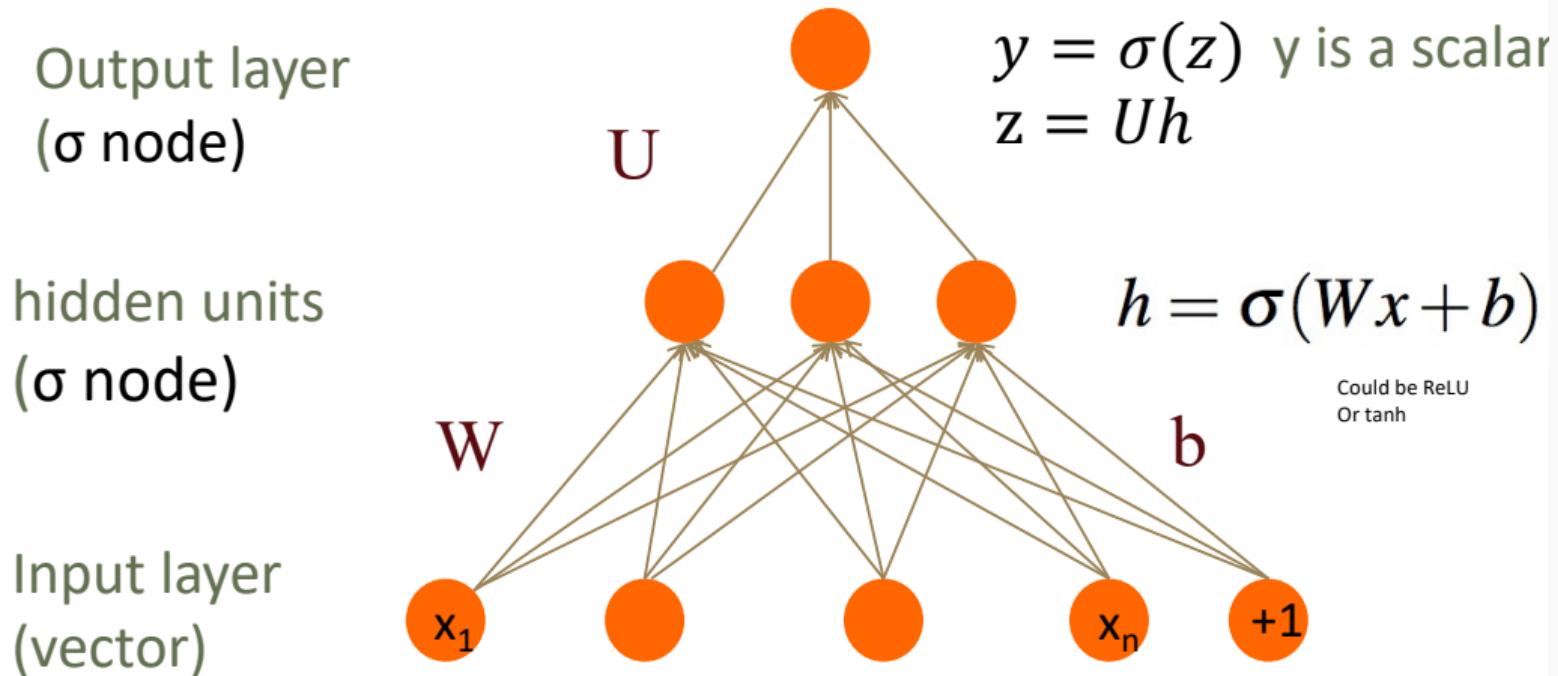


# Multinomial Logistic Regression as a 1-layer Network

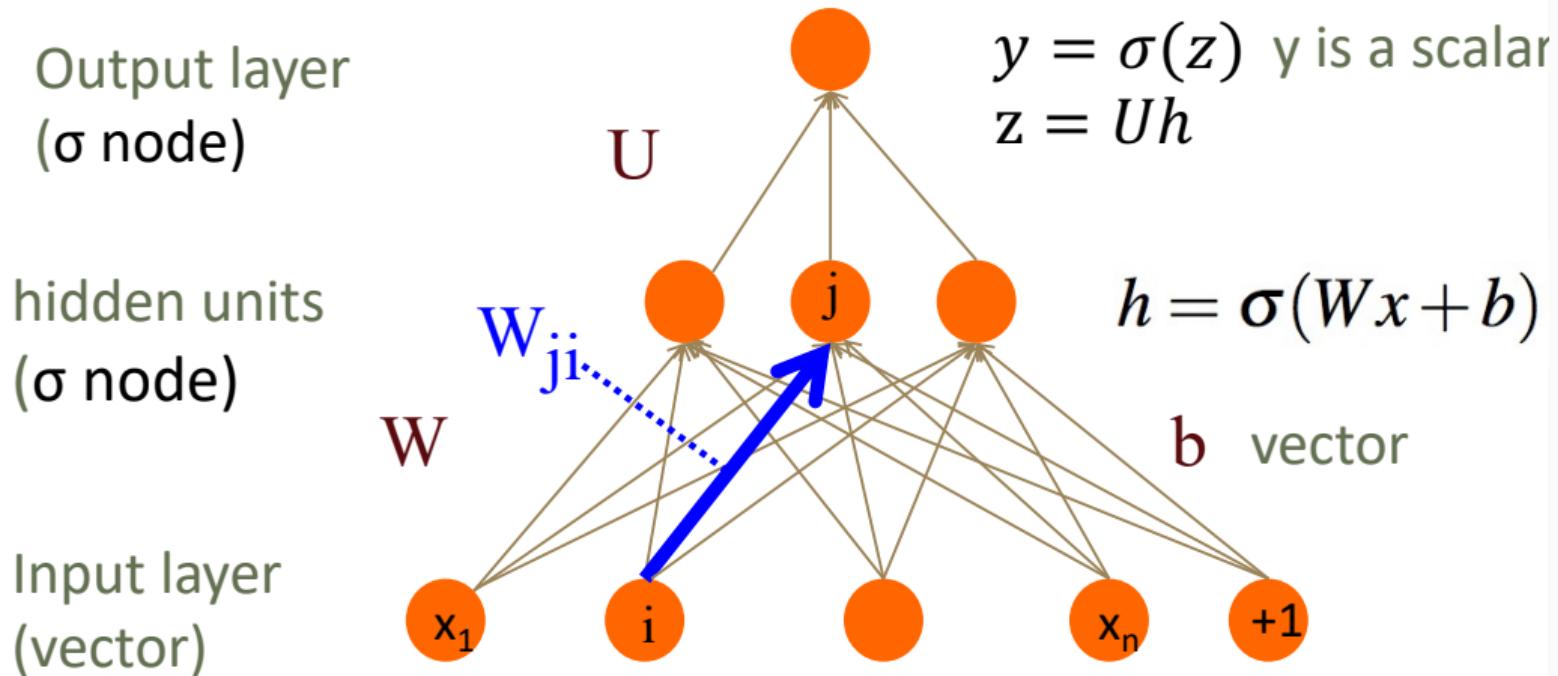
Fully connected single layer network



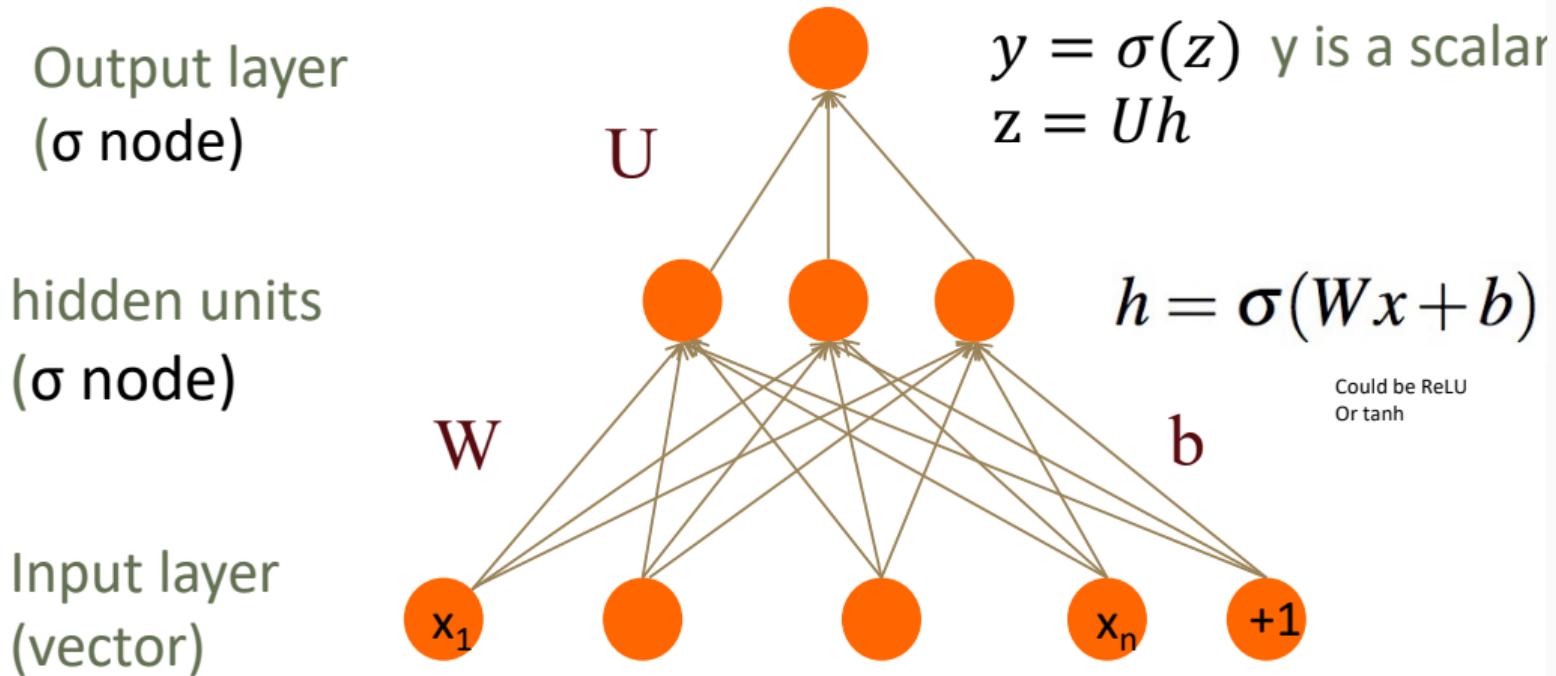
# Two-Layer Network with scalar output



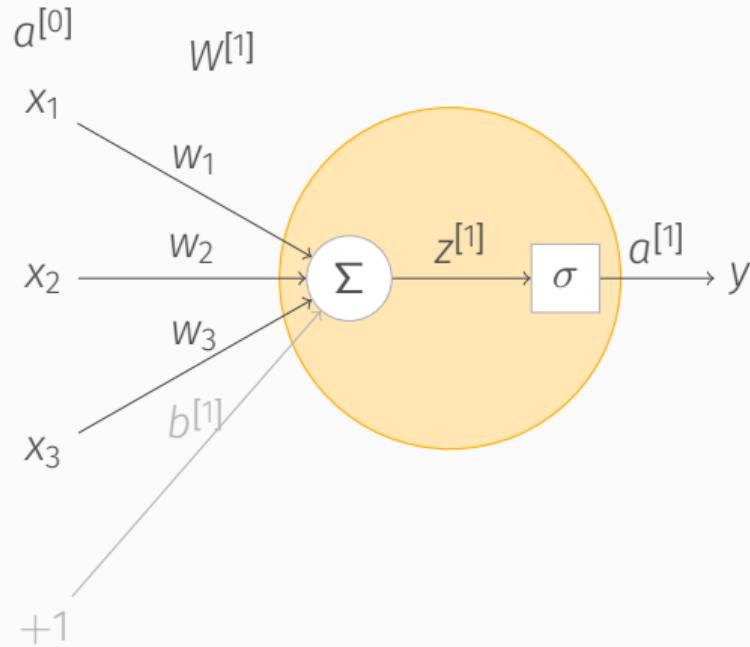
# Two-Layer Network with scalar output



# Two-Layer Network with scalar output

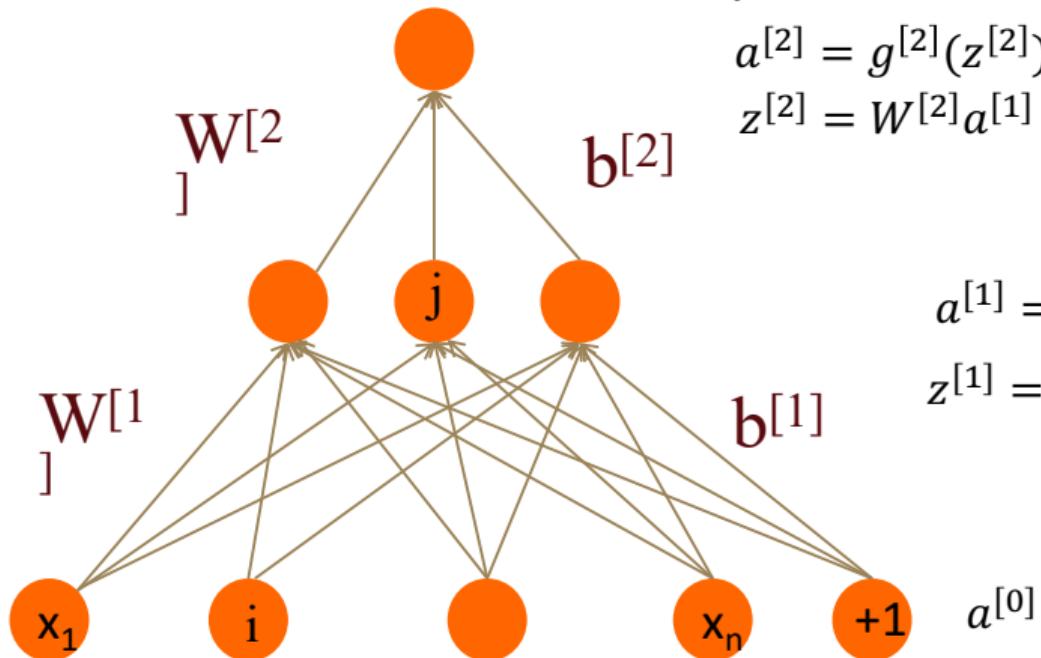


## A Forward Pass in Terms of Multi-Layer Notation



```
for  $i \in 1..n$  do  
   $z^{[i]} \leftarrow W^{[i]}a^{[i-1]} + b^{[i]}$   
   $a^{[i]} \leftarrow g^{[i]}(z^{[i]})$   
 $\hat{y} \leftarrow a^{[n]}$ 
```

# Multi-layer Notation



$$y = a^{[2]}$$

$$a^{[2]} = g^{[2]}(z^{[2]}) \quad \text{sigmoid or softmax}$$

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[1]} = g^{[1]}(z^{[1]}) \quad \text{ReLU}$$

$$z^{[1]} = W^{[1]}a^{[0]} + b^{[1]}$$

$$a^{[0]}$$

# Replacing the bias unit

Instead of:

$$x = x_1, x_2, \dots, x_{n_0}$$

$$h = \sigma(Wx + b)$$

$$h_j = \sigma \left( \sum_{i=1}^{n_0} W_{ji} x_i + b_j \right)$$

We'll do this:

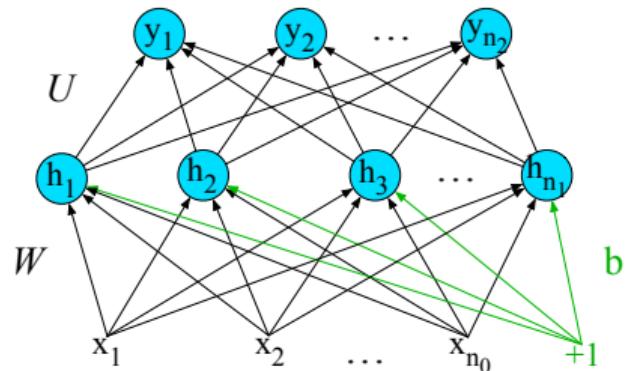
$$x = x_0, x_1, x_2, \dots, x_{n_0}$$

$$h = \sigma(Wx)$$

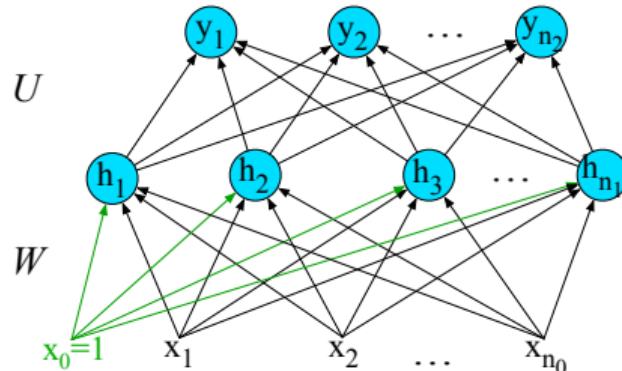
$$\sigma \left( \sum_{i=0}^{n_0} W_{ji} x_i \right)$$

# Replacing the bias unit

Instead of:



We'll do this:



# Feedforward Neural Nets as Classifiers

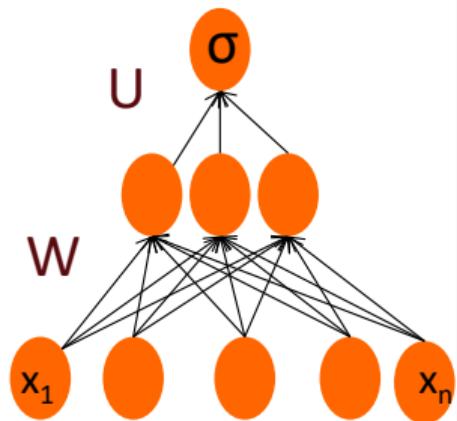
---

# Classification: Sentiment Analysis

We could do exactly what we did with logistic regression

Input layer are binary features as before

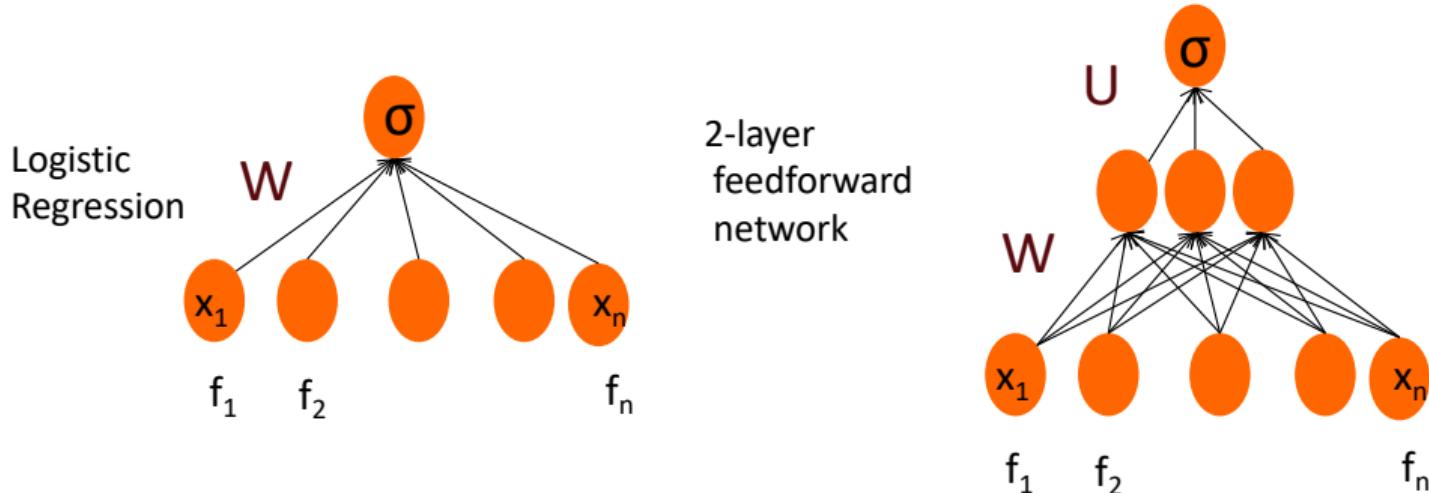
Output layer is 0 or 1



# Sentiment Features

Var	Definition
$x_1$	$\text{count}(\text{positive lexicon}) \in \text{doc}$
$x_2$	$\text{count}(\text{negative lexicon}) \in \text{doc}$ )
$x_3$	$\begin{cases} 1 & \text{if “no”} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$
$x_4$	$\text{count}(1\text{st and 2nd pronouns} \in \text{doc})$
$x_5$	$\begin{cases} 1 & \text{if “!”} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$
$x_6$	$\log(\text{word count of doc})$

# Feedforward nets for simple classification



Just adding a hidden layer to logistic regression

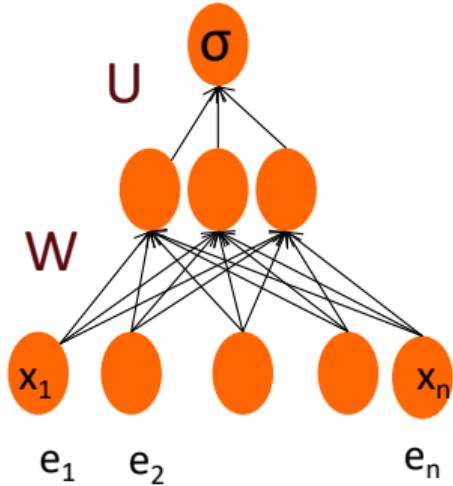
- allows the network to use non-linear interactions between features
- which may (or may not) improve performance.

## Even better: representation learning

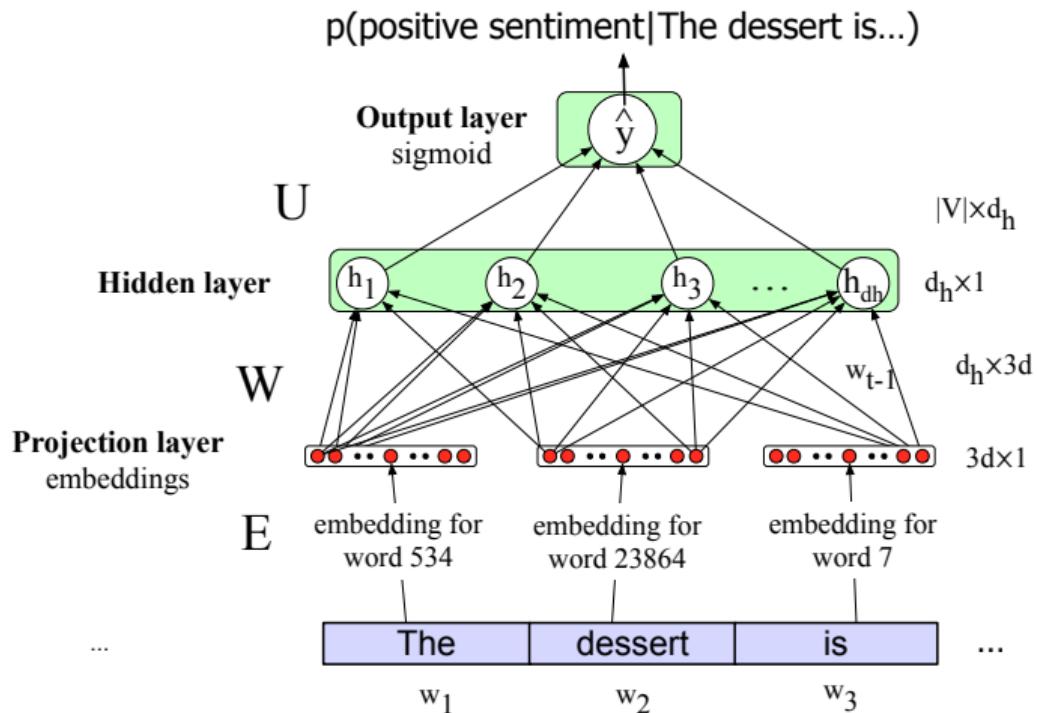
The real power of deep learning comes from the ability to **learn** features from the data

Instead of using hand-built human-engineered features for classification

Use learned representations like embeddings!



# Neural Net Classification with embeddings as input features!



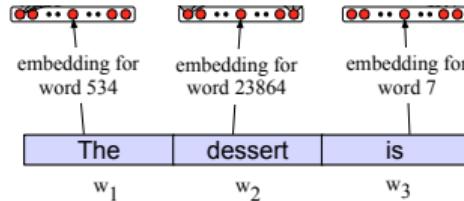
# Issue: texts come in different sizes

This assumes a fixed size length (3)!

Kind of unrealistic.

Some simple solutions (more sophisticated solutions later)

1. Make the input the length of the longest review
  - If shorter then pad with zero embeddings
  - Truncate if you get longer reviews at test time
2. Create a single "sentence embedding" (the same dimensionality as a word) to represent all the words
  - Take the mean of all the word embeddings
  - Take the element-wise max of all the word embeddings
    - For each dimension, pick the max value from all words

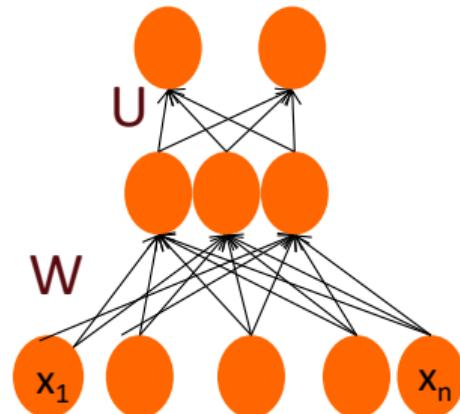


# Reminder: Multiclass Outputs

What if you have more than two output classes?

- Add more output units (one for each class)
- And use a “softmax layer”

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}} \quad 1 \leq i \leq D$$



# Introducing Feedforward Language Models

---

## We Start with Feedforward LMs because They Are Simple

- To start, we will look at language models built on feedforward neural networks
- In practice, it is not common to use such LMs
- Instead, RNNs or the Transformer are used as the basis for most neural language models
- However, feedforward neural nets are easy to understand and will allow us to discuss the basic concepts involved

## Defining Feedforward LMs

**Input** A representation at time  $t$  of some number of previous words  
 $(w_{t-1}, w_{t-2}, \dots)$

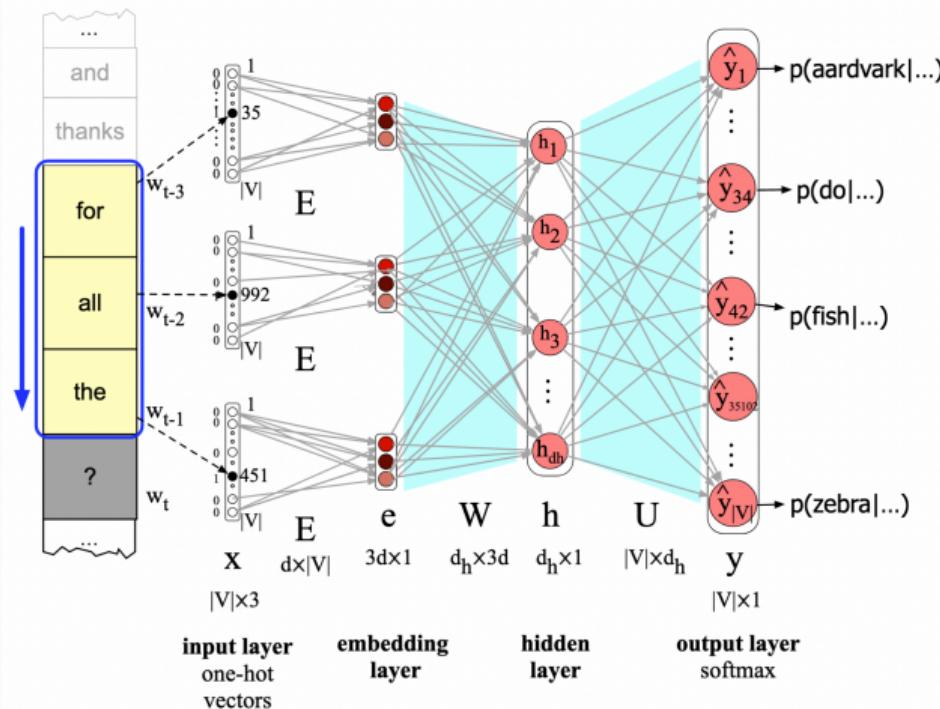
**Output** A probability distribution over possible next words

Such a language model is like an n-gram language model in that it approximates the probability of a word given the **entire** prior context  $P(w_t|w_1 : t - 1)$  by approximating based on the  $N$  previous words:

$$P(w_t|w_1, \dots, w_{t-1}) \approx P(w_t|w_{t-N+1}, \dots, w_{t-1})$$

### The Markov Assumption

# Calculating a Probability Distribution over $w_t$ Using a Simple Feedforward LM



1. Take  $N$  context words at each timestep, converting each to a  $d$ -dimensional EMBEDDING and concatenating them together (yielding a  $1 \times Nd$  unit input layer)
2. Multiply by  $W$  Multiply these units by a weight matrix  $W$
3. Apply activation function element-wise to produce a hidden layer  $h$
4. Multiply by  $U$ , another weight matrix
5. Apply softmax, predicting at each node  $i$  the probability that the new word  $w_i$  will be the vocabulary word  $V_i$

# Why Neural LMs work better than N-gram LMs

## Training data:

We've seen: I have to make sure that the cat gets fed.

Never seen: dog gets fed

## Test data:

I forgot to make sure that the dog gets \_\_

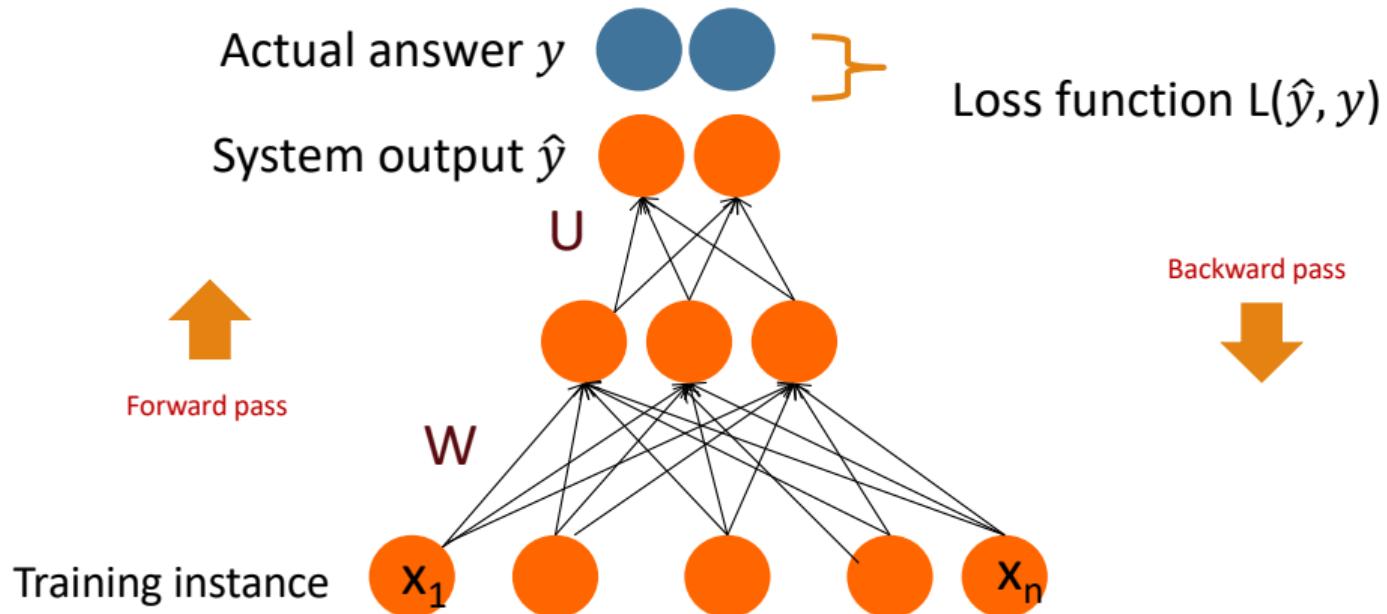
N-gram LM can't predict "fed"!

Neural LM can use similarity of "cat" and "dog" embeddings to generalize and predict "fed" after dog

## Training Neural Models

---

# Intuition: training a 2-layer Network



# The Intuition Behind Training a 2-Layer Network

For every training tuple  $(x, y)$

1. Run **forward** computation to find the estimate  $\hat{y}$
2. Run **backward** computation to update weights
  - For every output node
    - Compute the loss  $L$  between true  $y$  and estimated  $\hat{y}$
    - For every weight  $w$  from the hidden layer to the output layer
    - Update the weight
  - For every hidden node
    - Assess how much blame it deserves for the current answer
    - From every weight  $w$  from the input layer to the hidden layer
    - Update the weight

## This Brings Us to Computation Graphs

For training, we need the derivative of the loss with respect to each weight in every layer of the network

- But the loss is computed only at the very end of the network!
- Solution: **error backpropagation** (Rumelhart, Hinton, Williams, 1986)
- Backprop is a special case of **backward differentiation** which relies on **computation graphs**.

# Computation Graphs

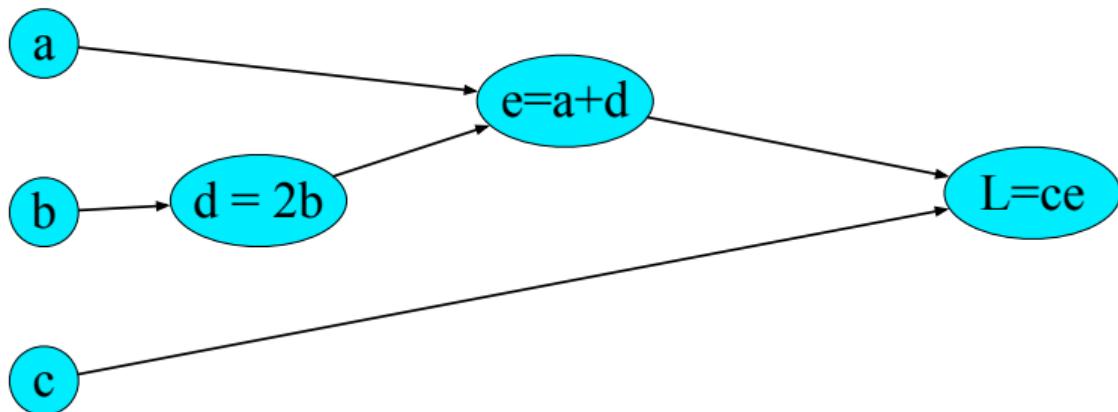
A computation graph represents the process of computing a mathematical expression

Example:  $L(a,b,c) = c(a+2b)$

$$d = 2 * b$$

Computations:  $e = a + d$

$$L = c * e$$



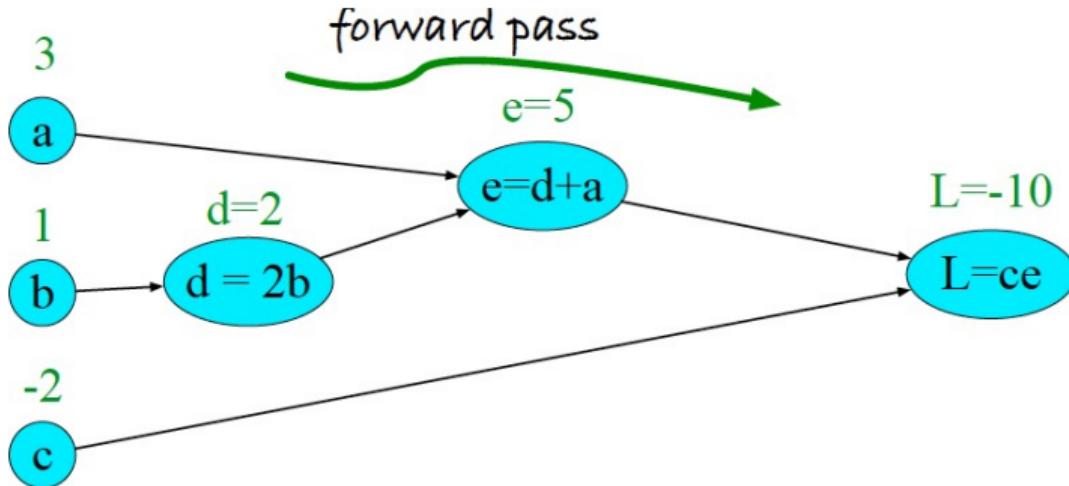
Example:  $L(a, b, c) = c(a + 2b)$

$$d = 2 * b$$

Computations:

$$e = a + d$$

$$L = c * e$$



## Backwards differentiation in computation graphs

The importance of the computation graph comes from the backward pass

This is used to compute the derivatives that we'll need for the weight update.

Example  $L(a,b,c) = c(a + 2b)$

$$d = 2 * b$$

$$e = a + d$$

$$L = c * e$$

We want:  $\frac{\partial L}{\partial a}$ ,  $\frac{\partial L}{\partial b}$ , and  $\frac{\partial L}{\partial c}$

The derivative  $\frac{\partial L}{\partial a}$ , tells us how much a small change in  $a$  affects  $L$ .

# The chain rule

Computing the derivative of a composite function:

$$f(x) = u(v(x))$$

$$\frac{df}{dx} = \frac{du}{dv} \cdot \frac{dv}{dx}$$

$$f(x) = u(v(w(x)))$$

$$\frac{df}{dx} = \frac{du}{dv} \cdot \frac{dv}{dw} \cdot \frac{dw}{dx}$$

Example  $L(a, b, c) = c(a + 2b)$

$$d = 2 * b$$

$$e = a + d$$

$$L = c * e$$

$$\frac{\partial L}{\partial c} = e$$

$$\frac{\partial L}{\partial a} = \frac{\partial L}{\partial e} \frac{\partial e}{\partial a}$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial e} \frac{\partial e}{\partial d} \frac{\partial d}{\partial b}$$

## Example

$$L(a, b, c) = c(a + 2b)$$

$$d = 2 * b$$

$$e = a + d$$

$$L = c * e$$

$$\begin{aligned}\frac{\partial L}{\partial a} &= \frac{\partial L}{\partial e} \frac{\partial e}{\partial a} \\ \frac{\partial L}{\partial b} &= \frac{\partial L}{\partial e} \frac{\partial e}{\partial d} \frac{\partial d}{\partial b}\end{aligned}$$

$$L = ce : \quad \frac{\partial L}{\partial e} = c, \frac{\partial L}{\partial c} = e$$

$$e = a + d : \quad \frac{\partial e}{\partial a} = 1, \frac{\partial e}{\partial d} = 1$$

$$d = 2b : \quad \frac{\partial d}{\partial b} = 2$$

## Example

$$a=3$$

a

$$b=1$$

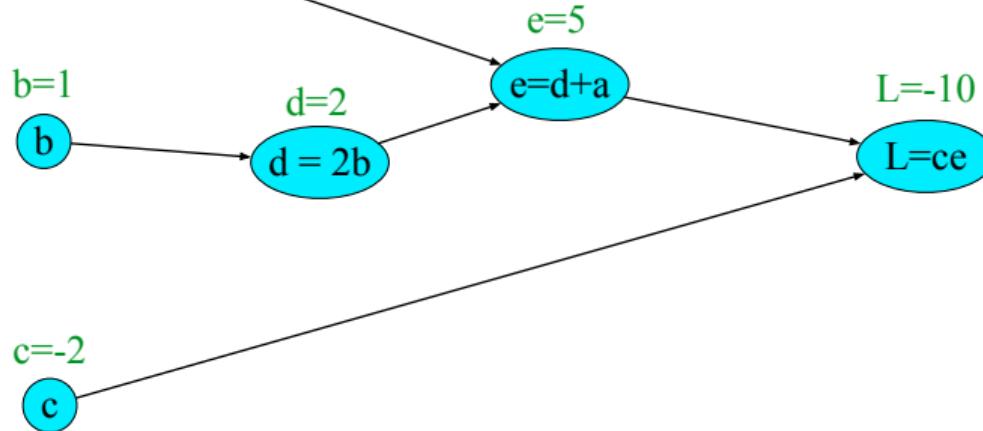
b

$$c=-2$$

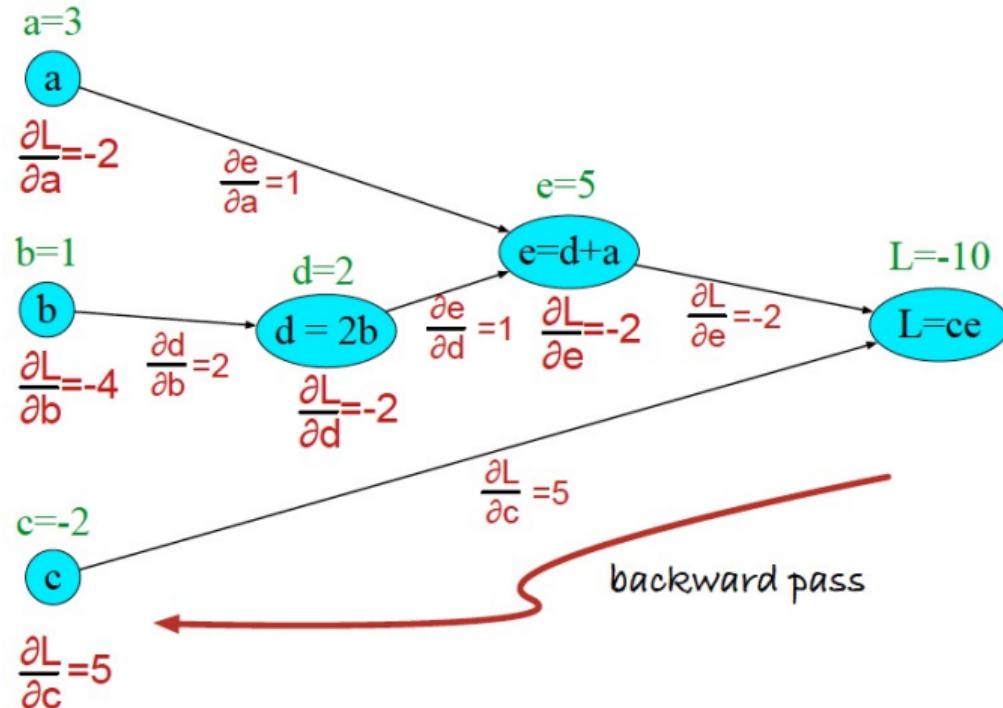
c

$$\frac{\partial L}{\partial a} = \frac{\partial L}{\partial e} \frac{\partial e}{\partial a}$$
$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial e} \frac{\partial e}{\partial d} \frac{\partial d}{\partial b}$$

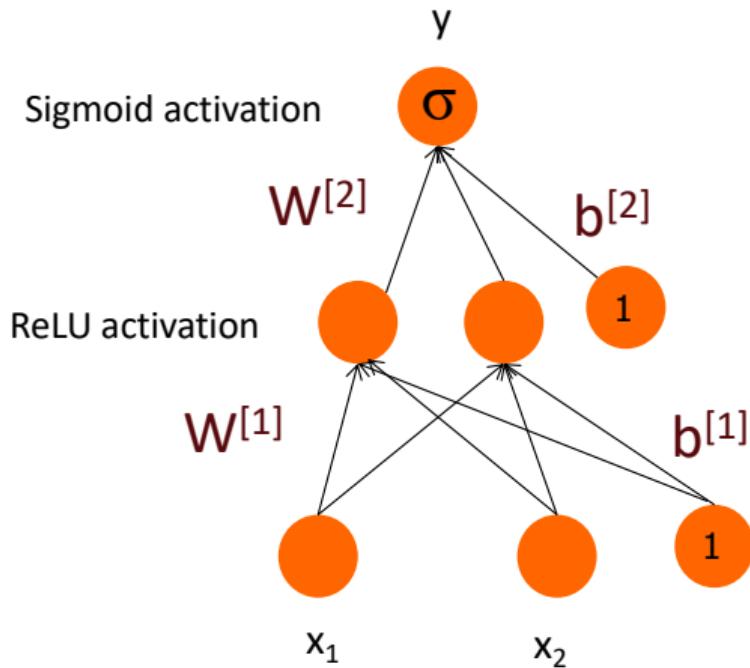
$$L = ce : \quad \frac{\partial L}{\partial e} = c, \frac{\partial L}{\partial c} = e$$
$$e = a + d : \quad \frac{\partial e}{\partial a} = 1, \frac{\partial e}{\partial d} = 1$$
$$d = 2b : \quad \frac{\partial d}{\partial b} = 2$$



# Example



# Backward differentiation on a two layer network



$$z^{[1]} = W^{[1]} \mathbf{x} + b^{[1]}$$

$$a^{[1]} = \text{ReLU}(z^{[1]})$$

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = \sigma(z^{[2]})$$

$$\hat{y} = a^{[2]}$$

## Backward differentiation on a two layer network

$$z^{[1]} = W^{[1]} \mathbf{x} + b^{[1]}$$

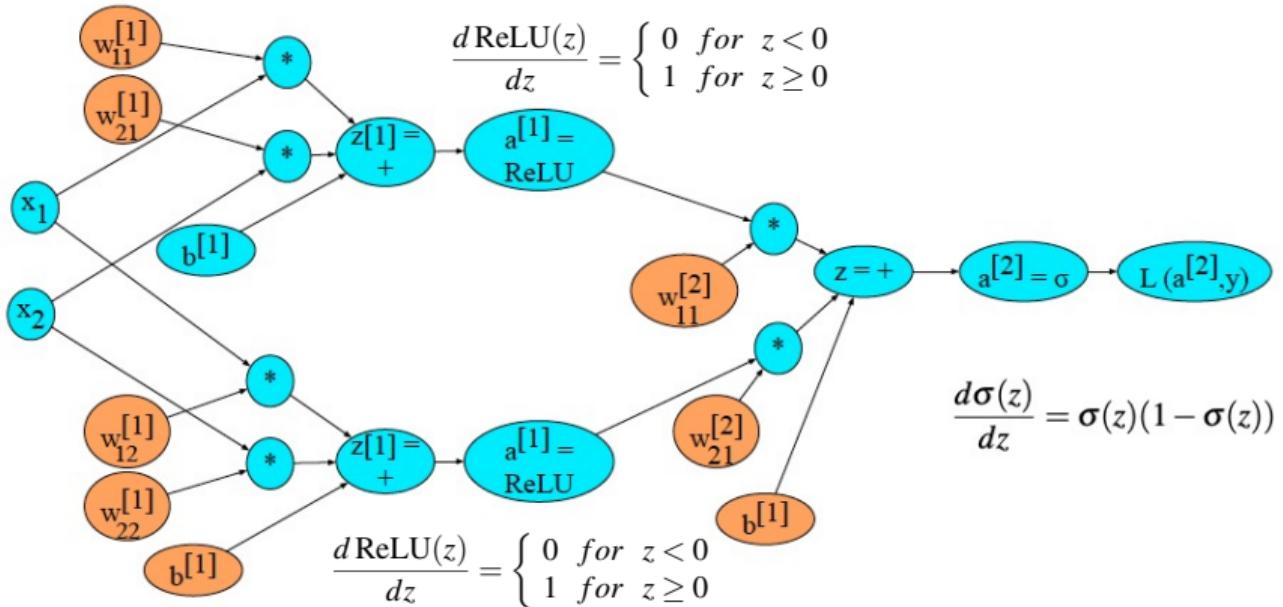
$$a^{[1]} = \text{ReLU}(z^{[1]}) \quad \frac{d \text{ReLU}(z)}{dz} = \begin{cases} 0 & \text{for } z < 0 \\ 1 & \text{for } z \geq 0 \end{cases}$$

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = \sigma(z^{[2]}) \quad \frac{d\sigma(z)}{dz} = \sigma(z)(1 - \sigma(z))$$

$$\hat{y} = a^{[2]}$$

# Backward differentiation on a 2-layer network



Starting off the backward pass:  $\frac{\partial L}{\partial z}$   
 (I'll write  $a$  for  $a^{[2]}$  and  $z$  for  $z^{[2]}$ )

$$\begin{aligned} z^{[1]} &= W^{[1]} \mathbf{x} + b^{[1]} \\ a^{[1]} &= \text{ReLU}(z^{[1]}) \\ z^{[2]} &= W^{[2]} a^{[1]} + b^{[2]} \\ a^{[2]} &= \sigma(z^{[2]}) \\ \hat{y} &= a^{[2]} \end{aligned}$$

$$L(\hat{y}, y) = -(y \log(\hat{y}) + (1 - y) \log(1 - \hat{y}))$$

$$L(a, y) = -(y \log a + (1 - y) \log(1 - a))$$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \frac{\partial a}{\partial z}$$

$$\begin{aligned} \frac{\partial L}{\partial a} &= - \left( \left( y \frac{\partial \log(a)}{\partial a} \right) + (1 - y) \frac{\partial \log(1 - a)}{\partial a} \right) \\ &= - \left( \left( y \frac{1}{a} \right) + (1 - y) \frac{1}{1 - a} (-1) \right) = - \left( \frac{y}{a} + \frac{y - 1}{1 - a} \right) \end{aligned}$$

$$\frac{\partial a}{\partial z} = a(1 - a) \quad \frac{\partial L}{\partial z} = - \left( \frac{y}{a} + \frac{y - 1}{1 - a} \right) a(1 - a) = a - y$$

# Summary

For training, we need the derivative of the loss with respect to weights in early layers of the network

- But loss is computed only at the very end of the network!

Solution: **backward differentiation**

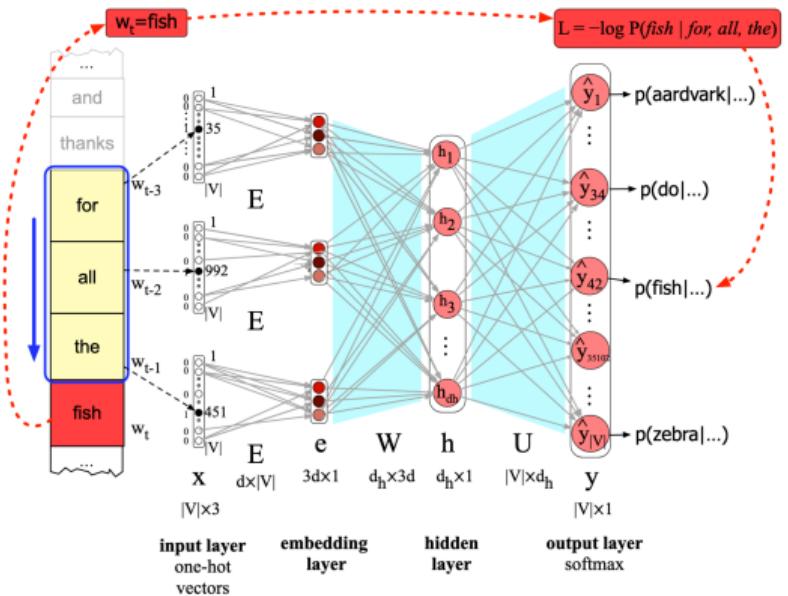
Given a computation graph and the derivatives of all the functions in it we can automatically compute the derivative of the loss with respect to these early weights.

## Back to the Narrative: Training (Feedforward) Neural Language Models

---

# Training a Feedforward LM: Forward Pass

1. **Select three embeddings from E** take words  $w_{t-3}$ ,  $w_{t-2}$ , and  $w_{t-1}$ , create three ONE-HOT VECTORS for them, and multiply each by embedding matrix  $E$ , yielding the PROJECTION LAYER. Concatenate the three embeddings.
2. **Multiply by W** and add  $b$  (the bias term). Pass the result through the ReLU (or other) activation function to get hidden layer  $h$
3. **Multiply by U**  $h$  is multiplied by  $U$
4. **Apply softmax** resulting in a layer where each node  $i$  in the output layer estimates the probability  $P(w_t = i | w_{t-1}, w_{t-2}, w_{t-3})$



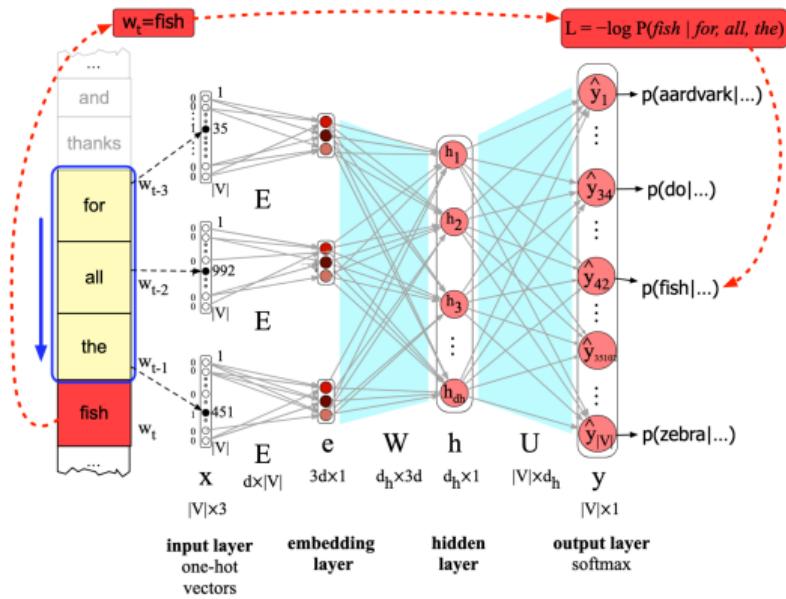
# Training a Feedforward LM: Backward Pass

1. Optional: freeze embedding layer  
(appropriate for some tasks)
2. Define loss (e.g., **cross-entropy loss**)

$$L_{CE}(\hat{y}, y) = -\log \hat{y}_c$$

(in our case, the negative log of the output probability corresponding to the correct class,  $c$ )

3. Set parameters  $\theta = E, W, U, b$  via **gradient descent**, using **error backpropagation** on the **computation graph** to compute the gradient



## More Detail on the Backward Pass

Concatenate all sentences in the training set (starting with random weights) and move through this text iteratively, predicting each word  $w_t$ . Use cross entropy (that is, negative log likelihood) at each  $w_t$ :

$$L_{CE}(\hat{y}, y) = -\log \hat{y}_i \quad (6)$$

where  $i$  is the correct class.  $\hat{y}_i$  means the probability that the model assigns the correct next word  $w_t$ :

$$L_{CE} = -\log p(w_t | w_{t-1}, \dots, w_{t-n+1}) \quad (7)$$

So the parameter update for the stochastic gradient descent for this loss from set  $s$  to  $s + 1$  is:

$$\theta^{s+1} = \theta^s - \eta \frac{\partial [-\log p(w_t | w_{t-1}, \dots, w_{t-n+1})]}{\partial \theta} \quad (8)$$

Computing this gradient and backpropagating through  $\theta = \mathbf{E}, \mathbf{W}, \mathbf{U}, \mathbf{b}$  is trivial and is left as an exercise for PyTorch.

## There are Two Benefits to Training a Neural Language Model

If you train a neural language model, you get two things:

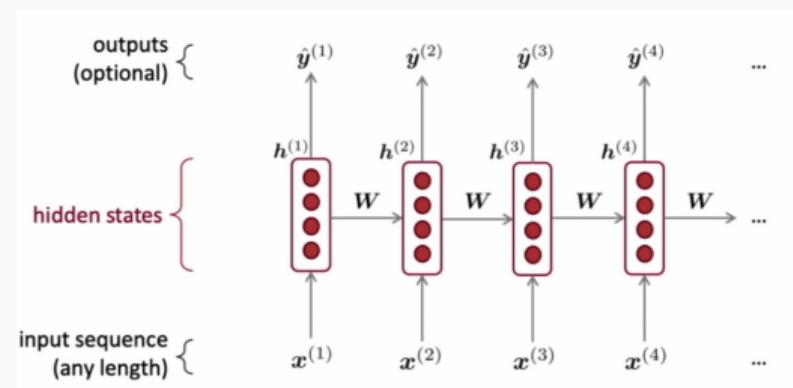
1. An algorithm that will allow you to predict the next word in a sequence
2. A set of embeddings  $\mathbf{E}$  that can be used to represent words in other tasks  
(assuming that you did not freeze the embedding layer)

# Language Models with Recurrent Neural Networks

---

# Refresher: RNNs

- RNNs are neural networks with a recurrent structure
  - Each step takes an input, emits an output and pass internal state on to the next step
  - Hidden state multiplied by weight matrix  $W$  at each step
- Good for modeling sequences



# An RNN Language Model

output distribution

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

hidden states

$$h^{(t)} = \sigma(W_h h^{(t-1)} + W_e e^{(t)} + b_1)$$

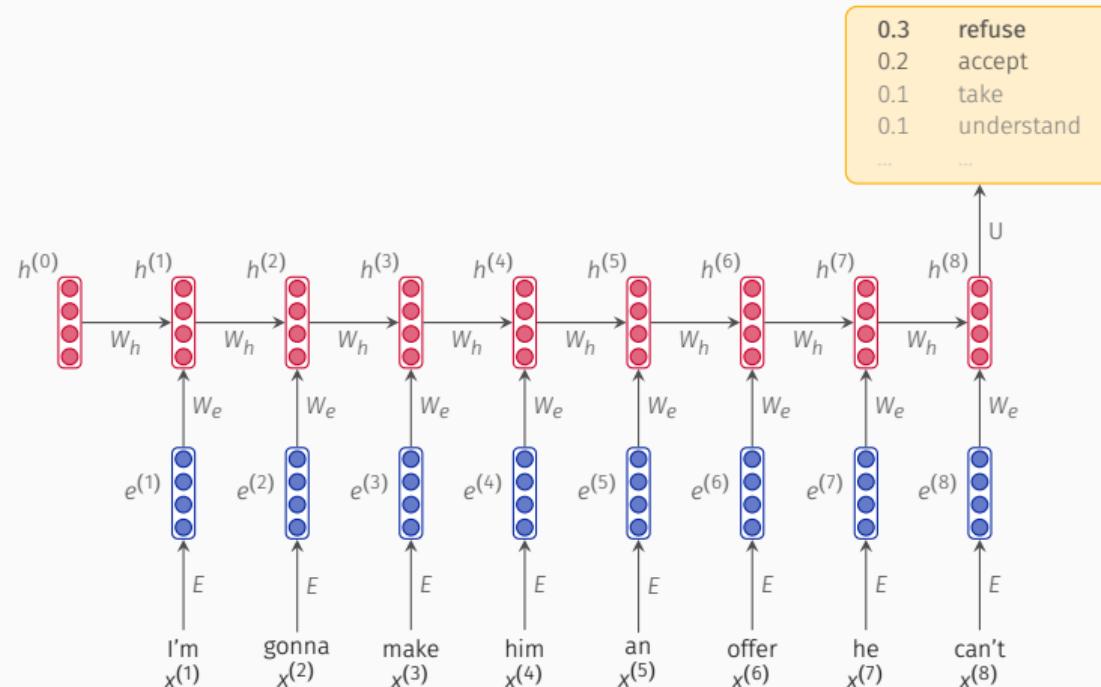
$h^{(0)}$  is the initial hidden state

word embeddings

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

one-hot vectors

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



## Training an RNN Language Model

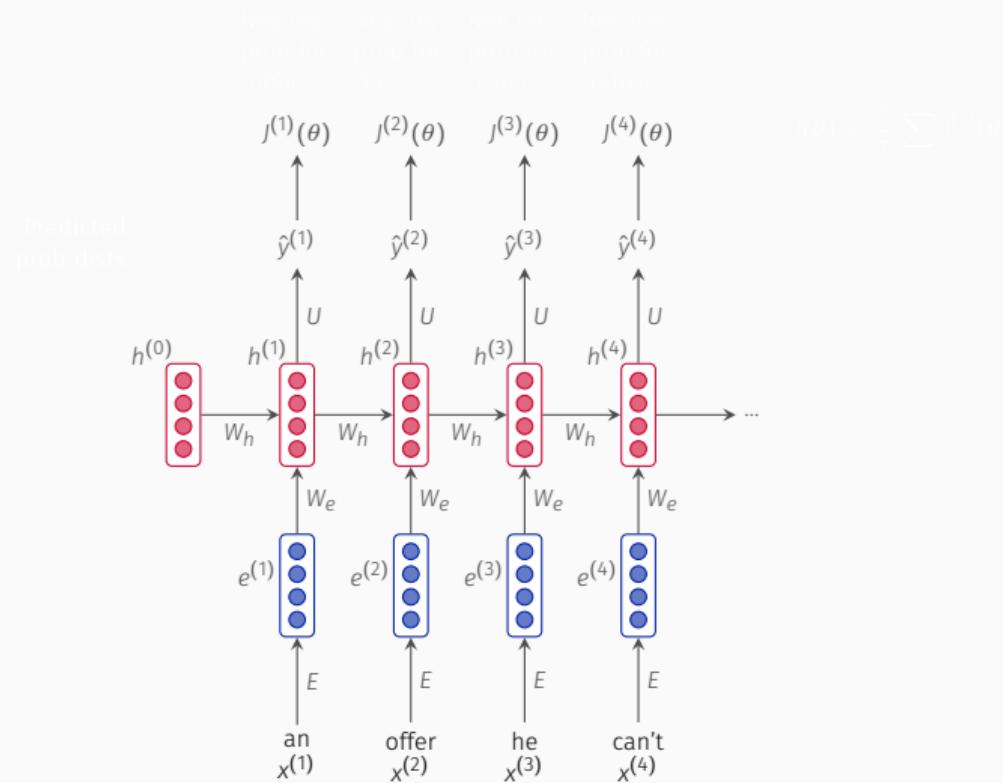
- Get a big corpus of text, which is a sequence of words  $x^{(1)}, \dots, x^{(T)}$
- Feed it into the RNN-LM, computing output distribution  $(t)$  for every step  $t$ .
- Loss function on step  $t$  is **cross-entropy** between the predicted probability distribution  $\hat{y}^{(t)}$  and the true next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

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

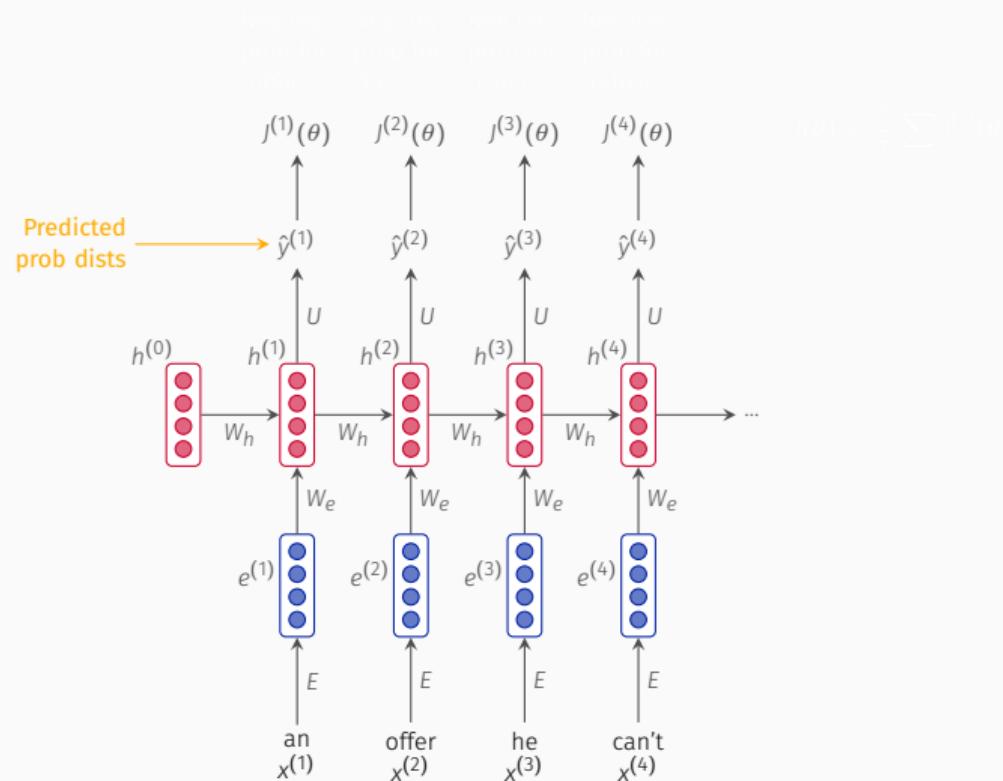
- Average this to get overall loss for the entire training set:

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

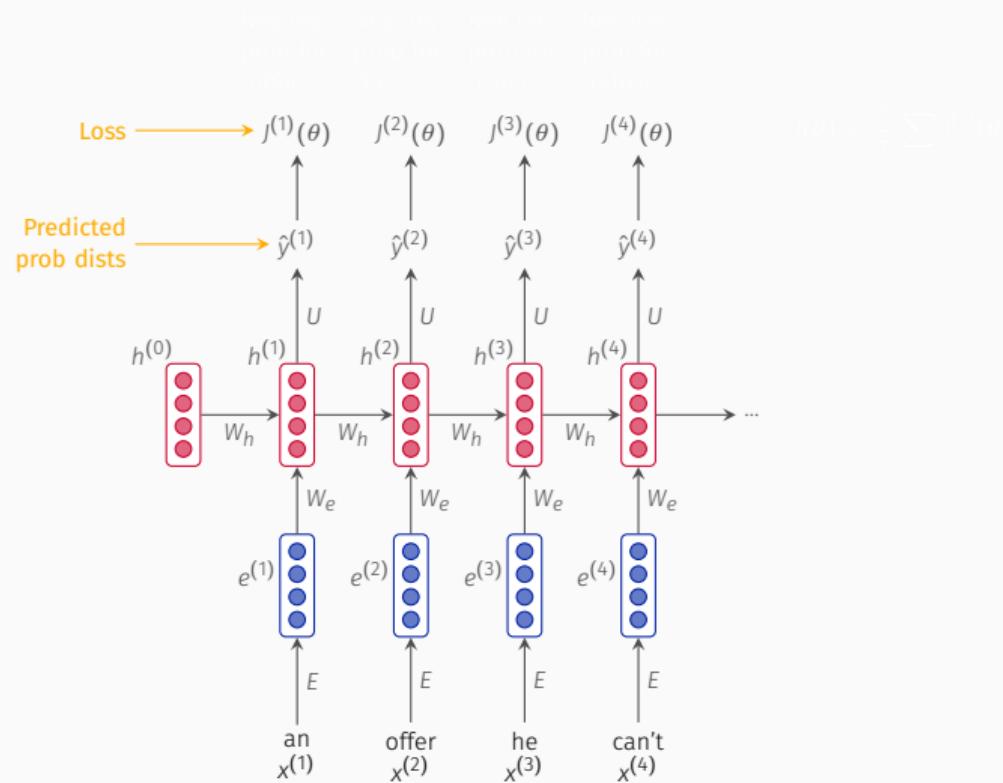
# Training an RNN Language Model



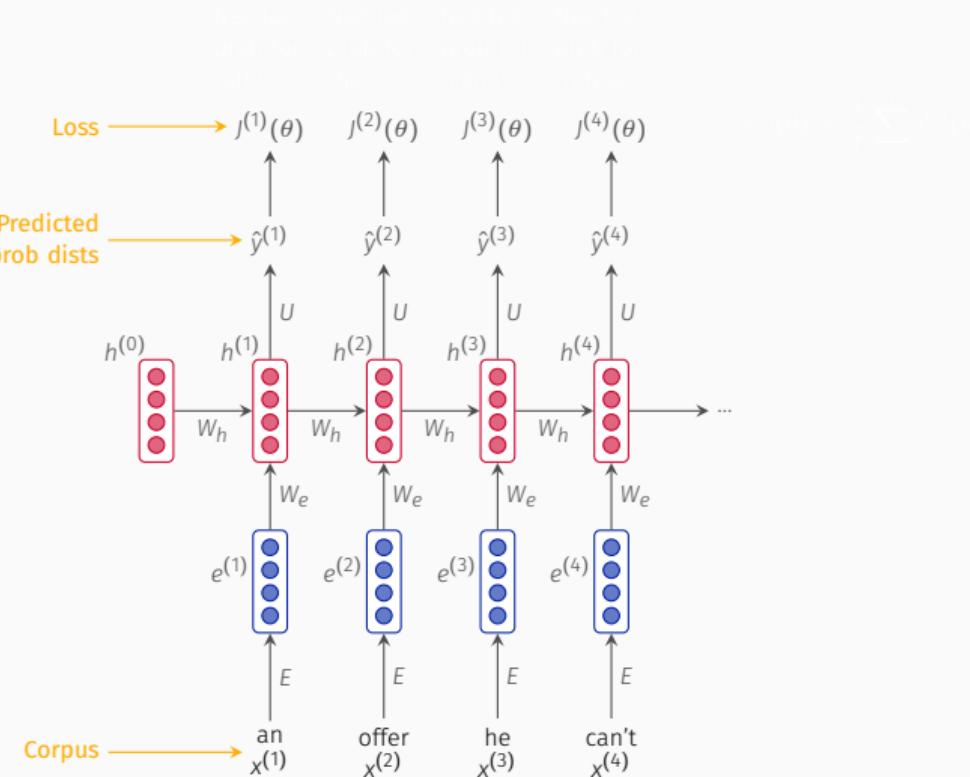
# Training an RNN Language Model



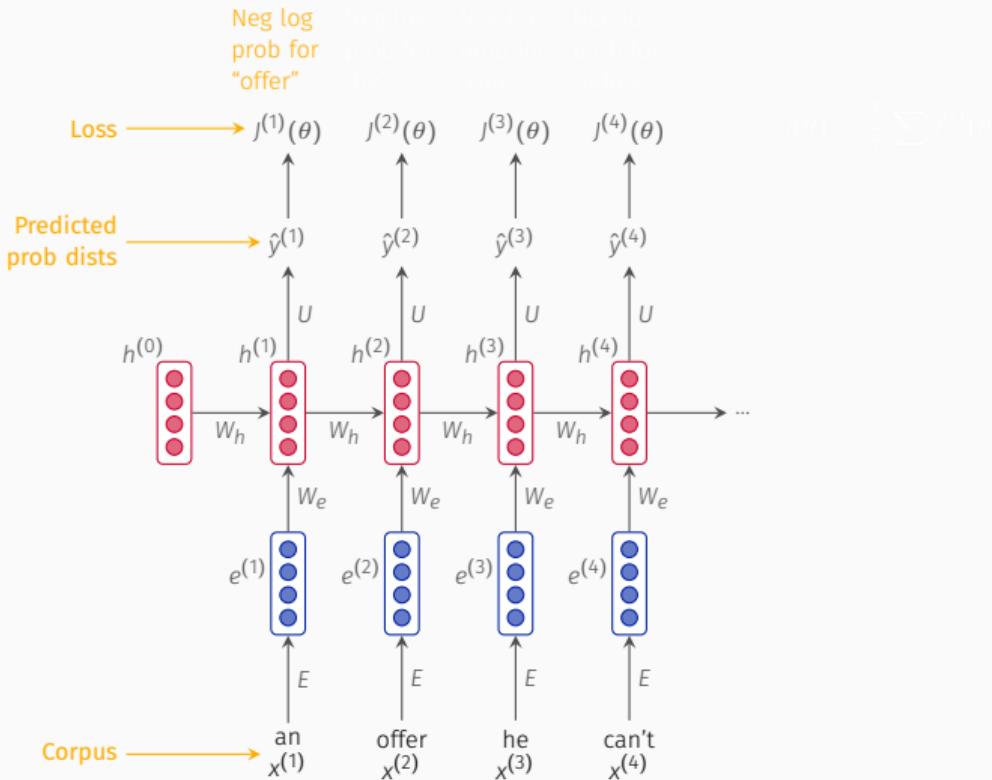
# Training an RNN Language Model



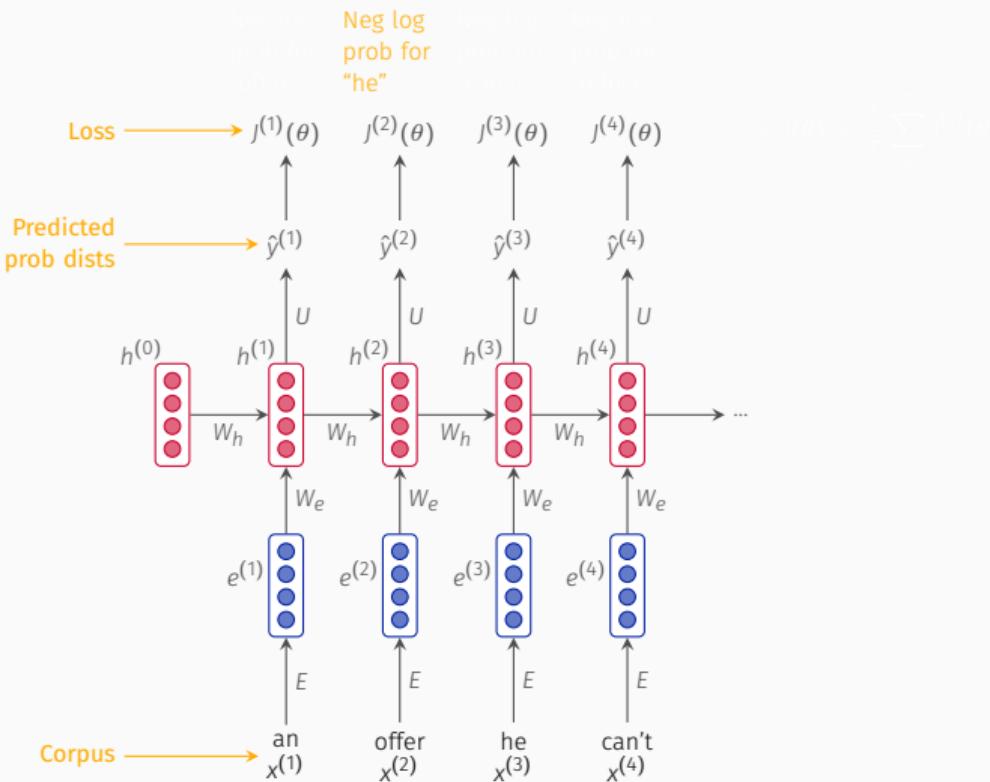
# Training an RNN Language Model



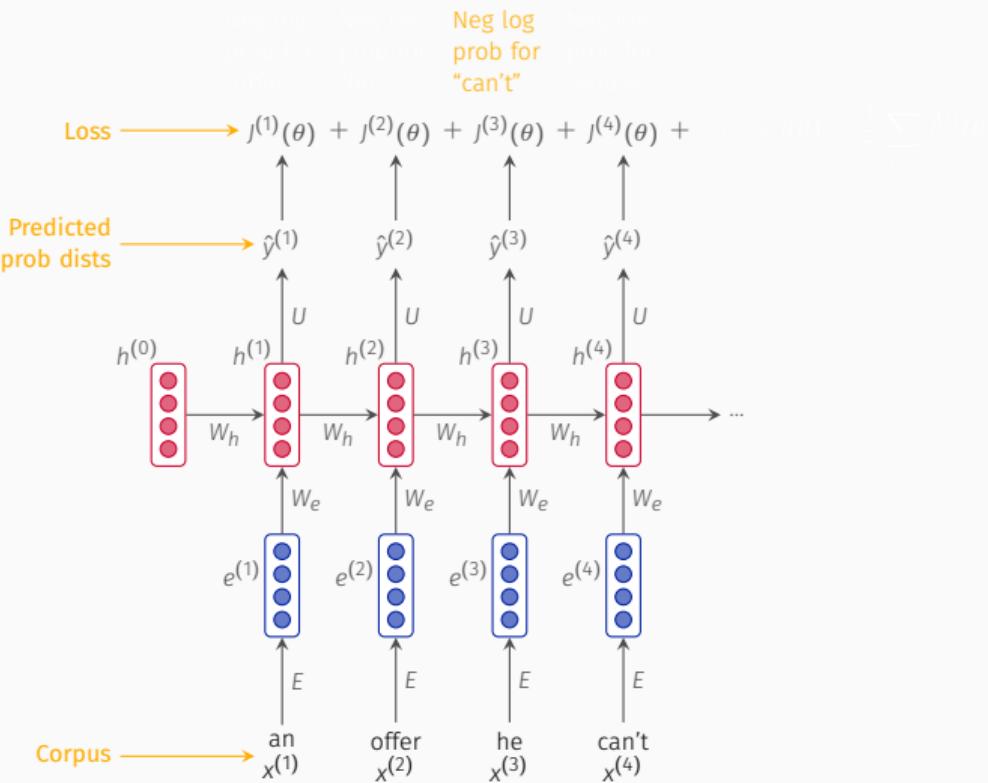
# Training an RNN Language Model



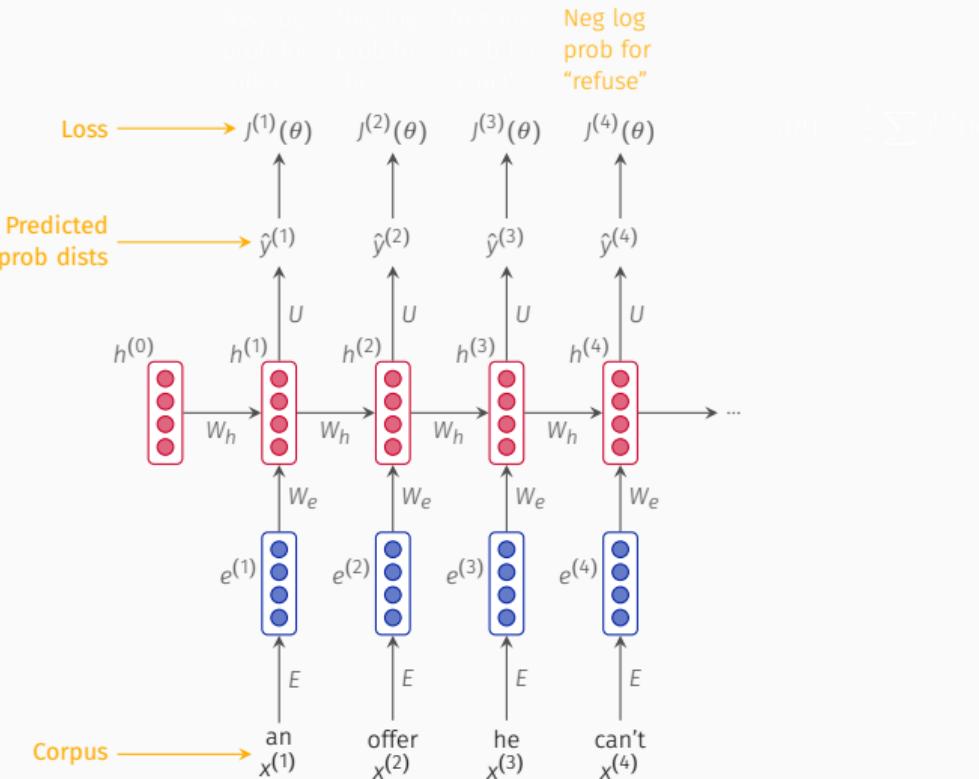
# Training an RNN Language Model



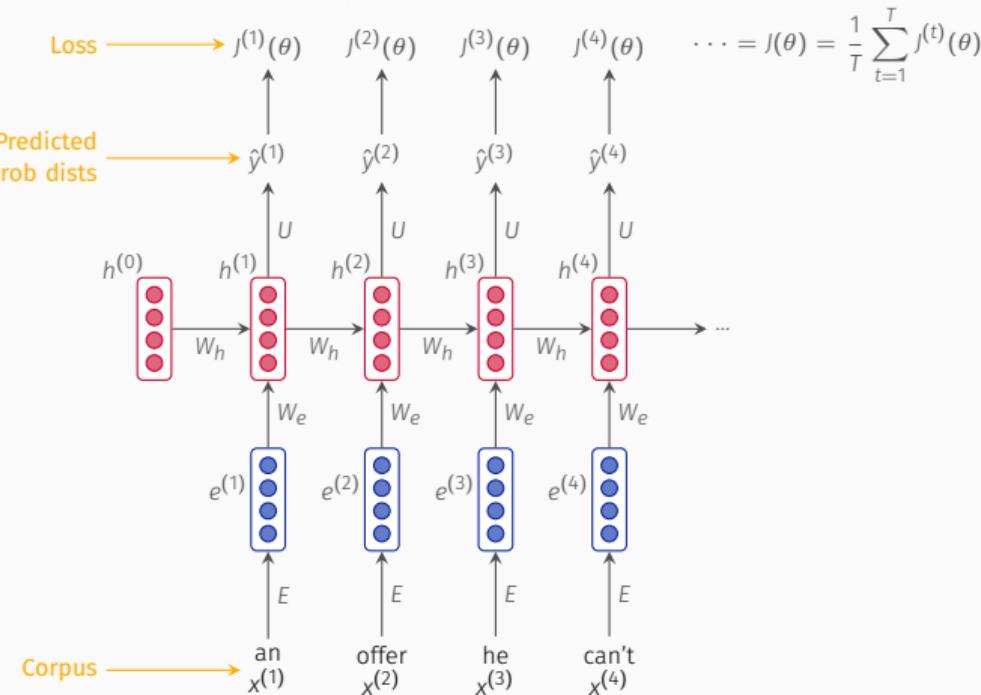
# Training an RNN Language Model



# Training an RNN Language Model



# Training an RNN Language Model



## Computing Loss and Gradients in Practice

- In principle, we could compute loss and gradients across the whole corpus  $(x^{(1)}, \dots, x^{(T)})$  but that would be incredibly expensive!

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

- Instead, we usually treat  $x^{(1)}, \dots, x^{(T)}$  as a document, or even a sentence
- This works much better with **Stochastic Gradient Descent**, which lets us compute loss and gradients for little chunks and update as we go.
- Actually, we do this in batches: compute  $J(\theta)$  for a batch of sentences; update weights; repeat.

## We Will Skip the Details of Backpropogation in RNNs for Now

- The fact that training RNNs involves backpropagation over timesteps, summing as you go, means that it (the **backpropagation through time** algorithm) is a bit more complicated than backpropagation in feedforward neural networks.
- We will skip these details for now, but you will want to learn them if you are doing serious work with RNNs.

## Advantages of RNN Language Models

- Input can be of an arbitrary length
- Computation can use information from many steps back (in principle)
- Longer inputs do not mean larger model sizes
- Same weights applied at every time step—**symmetry**

## Disadvantages of RNN LMs

- Recurrent computation is **slow**
- In practice, it is difficult to access information from many steps back (cf. the VANISHING GRADIENT PROBLEM)

## LSTMs Address (but Do not Solve) the Vanishing and Exploding Gradient Problems

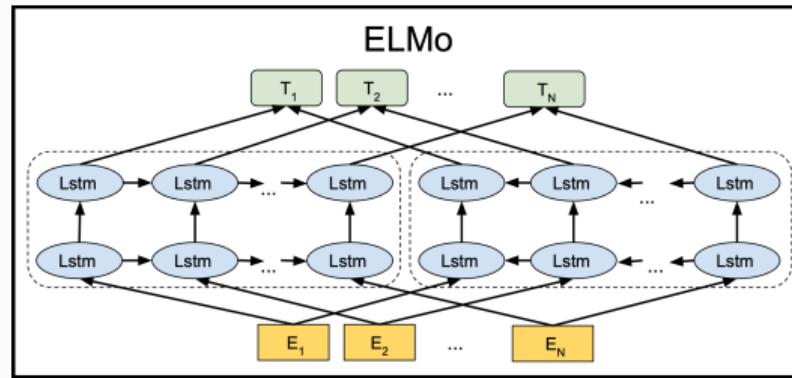
- LSTMs: Long Short Term Memory
- Process data sequentially, but keep hidden state through time
- Still subject, at some level, to vanishing gradients, but to a lesser degree than traditional RNNs
- Widely used in language modeling

# Prehistory of Pretrained Language Models

---

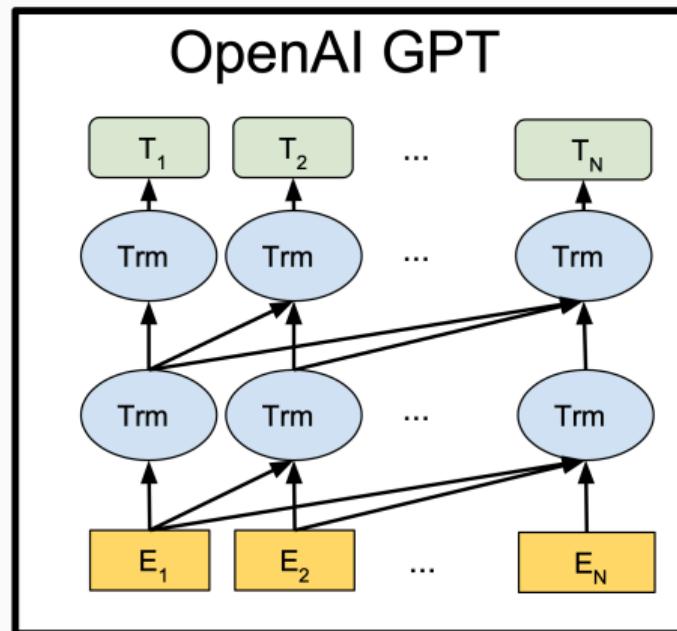
## ELMo

“ELMo is a deep contextualized word representation that models both (1) complex characteristics of word use (e.g., syntax and semantics), and (2) how these uses vary across linguistic contexts (i.e., to model polysemy). These word vectors are learned functions of the internal states of a deep bidirectional language model (biLM), which is pre-trained on a large text corpus. They can be easily added to existing models and significantly improve[d] the state of the art across a broad range of challenging NLP problems, including question answering, textual entailment and sentiment analysis.”



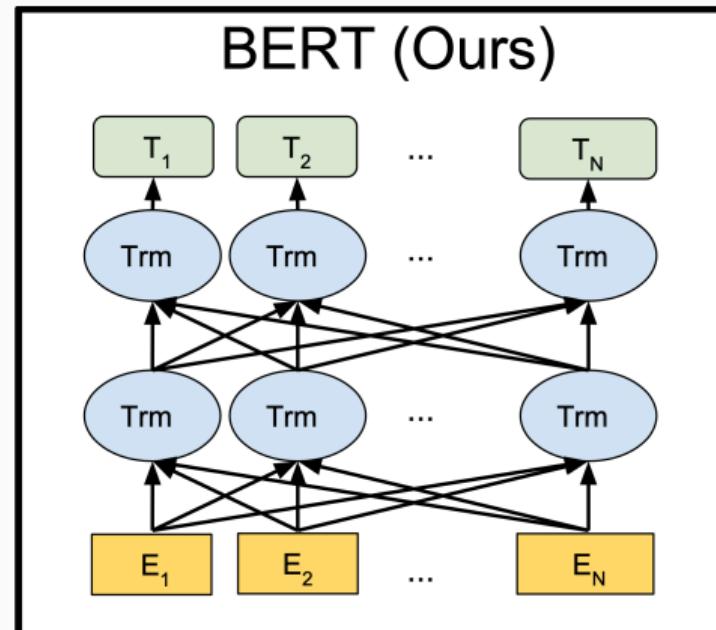
# GPT (Generative Pre-Trained Transformer)

The original GPT, released in 2018, combined transformers and unsupervised pre-training (ingredients that would also be central to BERT, GPT-2, and GPT-3). A broad view of its architecture is given below:



# BERT

BERT innovated beyond GPT in being bi-directional (like ELMo, only more so):



BERT also introduced a new training regimen.

## Transformers and BERT

---

## GPT and BERT Rely upon Transformers

- The ELMo architecture is based around LSTMs
- BERT and GPT incorporate a newer architectural idea: the Transformer
- Transformers have revolutionized many NLP tasks

## Transformers Improved upon RNNs and CNNs

- Google introduced Transformers in 2017
- At that time, most neural NLP models were based on:
  - RNNs
  - CNNs
- These were good
- For many tasks, Transformers were better

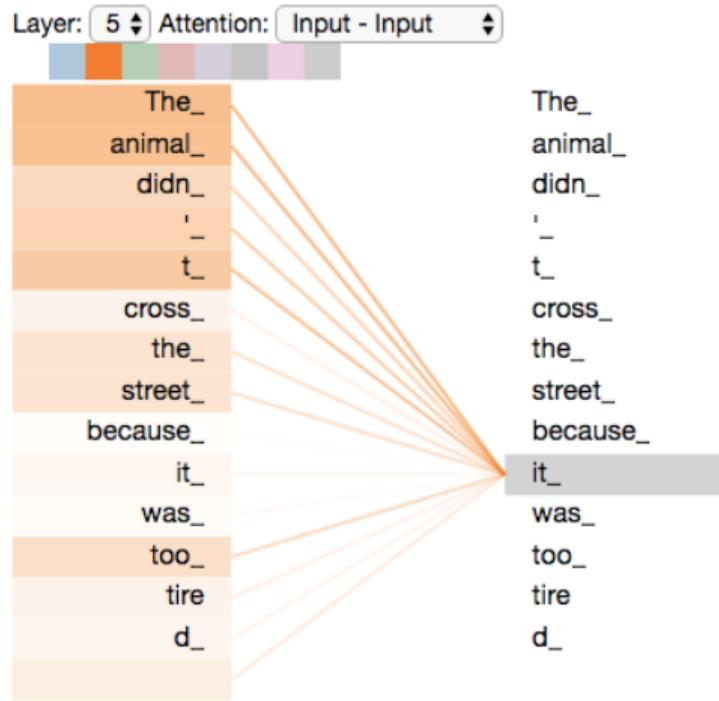
# Attention: All You Need?

## Attention:

- Each output element is connected to each input element
- The weightings between them are calculated dynamically, based on their connection
- Attention has long been important in NLP models for machine translation, etc.

Attention is fundamental to Transformers.

# Okay, Self-Attention Is All You Need!



Take the sentence: “**The animal didn’t cross the street because it was too tired**”. What is the antecedent of *it*?

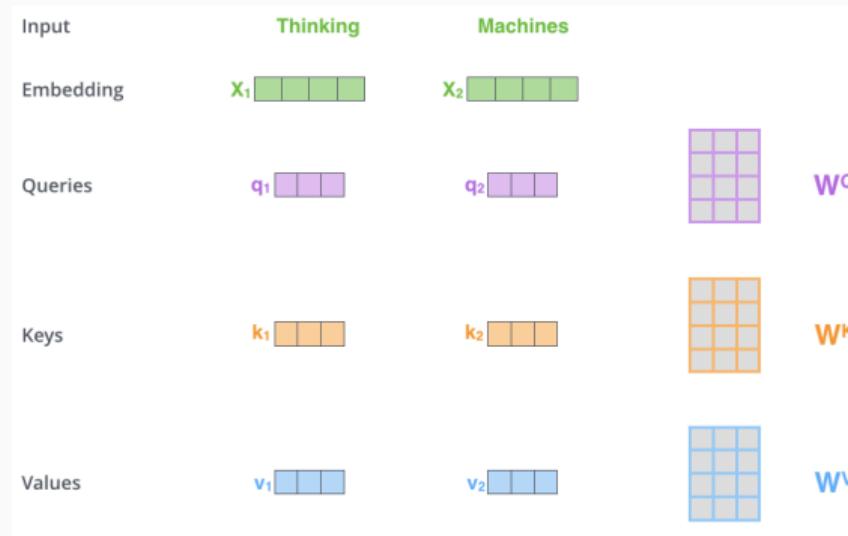
Self-attention allows the model to “attend” to all of the other positions and to process each position (including *the* and *animal*) to help it better encode the pronoun *it*.

You can compare this to the hidden state in an RNN—it conveys information about other words in the sequence to the position one is currently processing.

Transformers rely on self-attention.

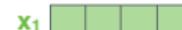
# A More Detailed Look at Self-Attention

Start with embeddings for our words (*thinking* and *machines*), and query, key, and value vectors obtained by multiplying the embedding by the weight vectors  $W^Q$ ,  $W^K$ , and  $W^V$ .



## Take the Dot-Product of Query and Keys

Here we take the dot product of the query vector for *thinking* ( $q_1$ ) and the key vectors for *thinking* ( $k_1$ ) and *machines* ( $k_2$ ):

Input	Thinking	Machines
Embedding	$x_1$ 	$x_2$ 
Queries	$q_1$ 	$q_2$ 
Keys	$k_1$ 	$k_2$ 
Values	$v_1$ 	$v_2$ 
Score	$q_1 \cdot k_1 = 112$	$q_1 \cdot k_2 = 96$

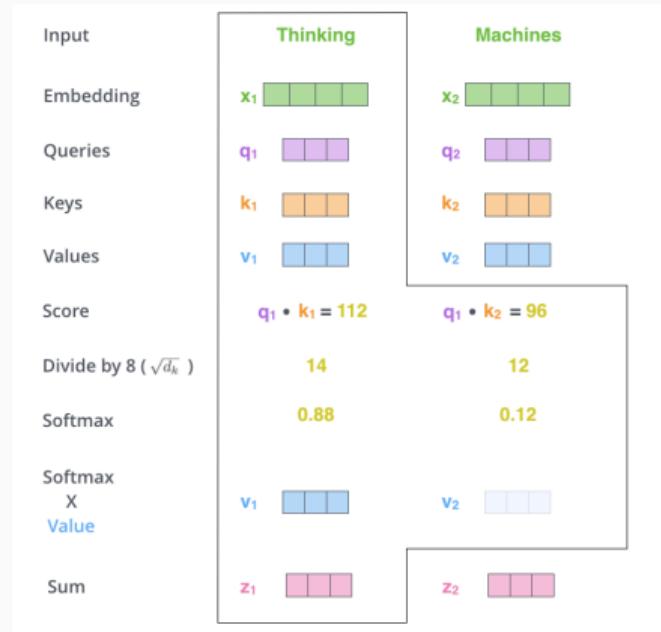
# Divide and Take Softmax

Divide the result by ( $\sqrt{d_k}$ ) and take the softmax (so that we get probabilities corresponding to each word):

Input	Thinking		Machines	
Embedding	$x_1$	[4 green boxes]	$x_2$	[4 green boxes]
Queries	$q_1$	[3 purple boxes]	$q_2$	[3 purple boxes]
Keys	$k_1$	[3 orange boxes]	$k_2$	[3 orange boxes]
Values	$v_1$	[3 blue boxes]	$v_2$	[3 blue boxes]
Score	$q_1 \cdot k_1 = 112$		$q_1 \cdot k_2 = 96$	
Divide by 8 ( $\sqrt{d_k}$ )	14		12	
Softmax	0.88		0.12	

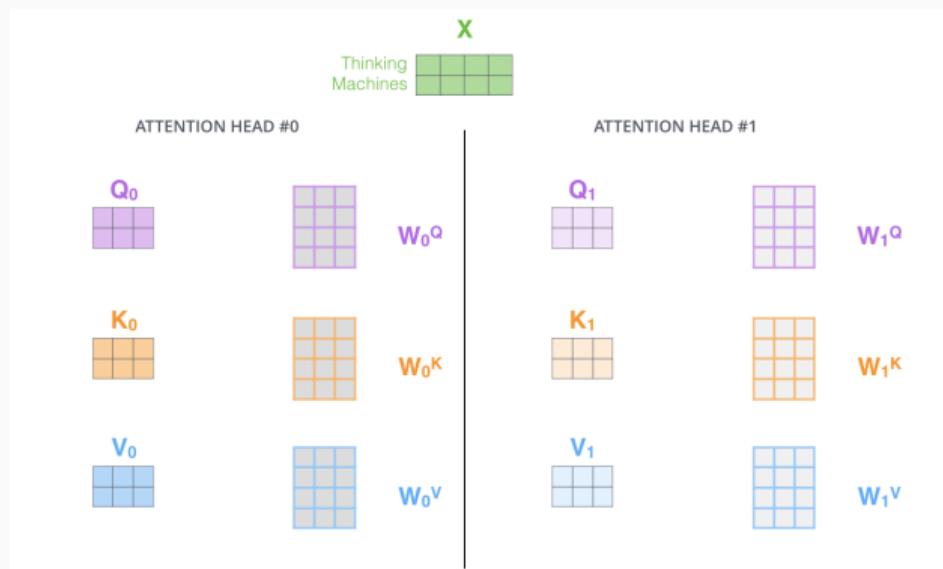
# Multiply and Sum

**Multiply** each value vector by the softmax score (drown out irrelevant words by multiplying them by very small numbers). **Sum** the weighted value vectors.



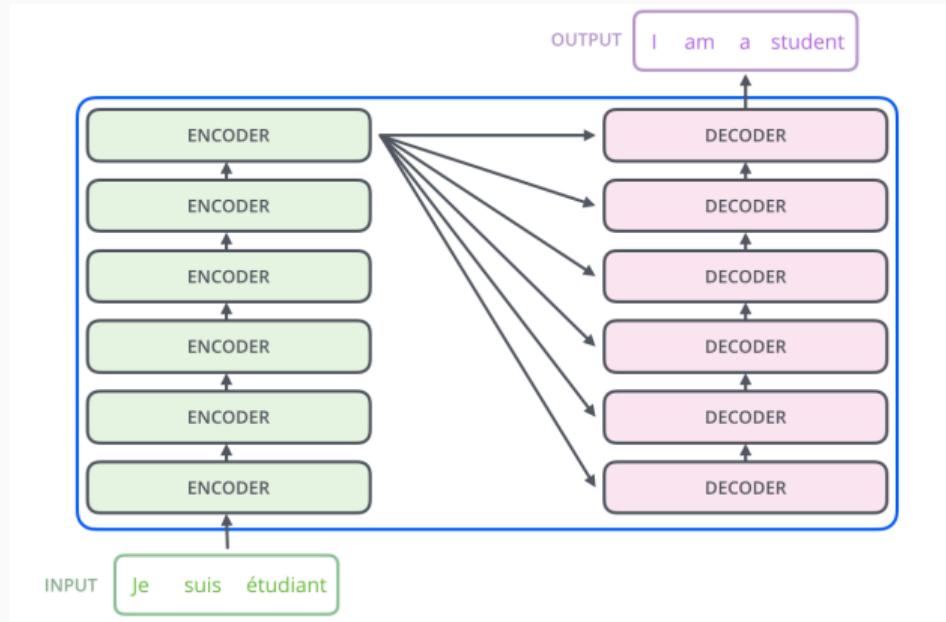
# Multi-Headed Attention Expands Transformer Models' Ability to Focus on Different Positions

Maintain distinct weight matrices for each attention head—distinct representational subspaces:



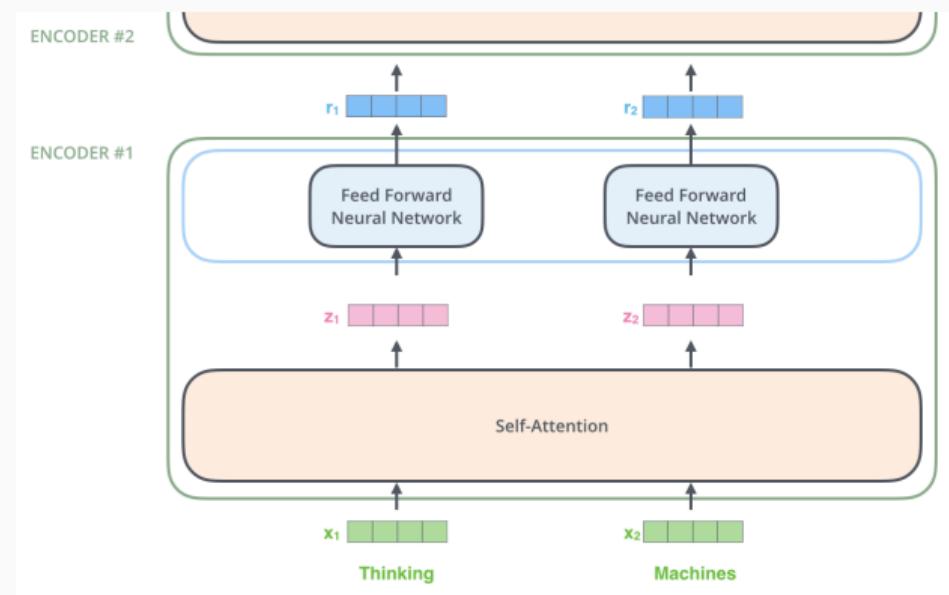
# Transformers: Stacks of Encoders and Decoders

A transformer is a stack of ~6 encoders and decoders. The encoders are identical in structure but do not share weights.



# The Structure of the Encoders

Each encoder takes input vectors. For the first encoder, these are the embeddings of the input words. For other encoders, these are the outputs of the preceding encoders. The inputs pass through a self attention layer, then through a feed forward neural net.

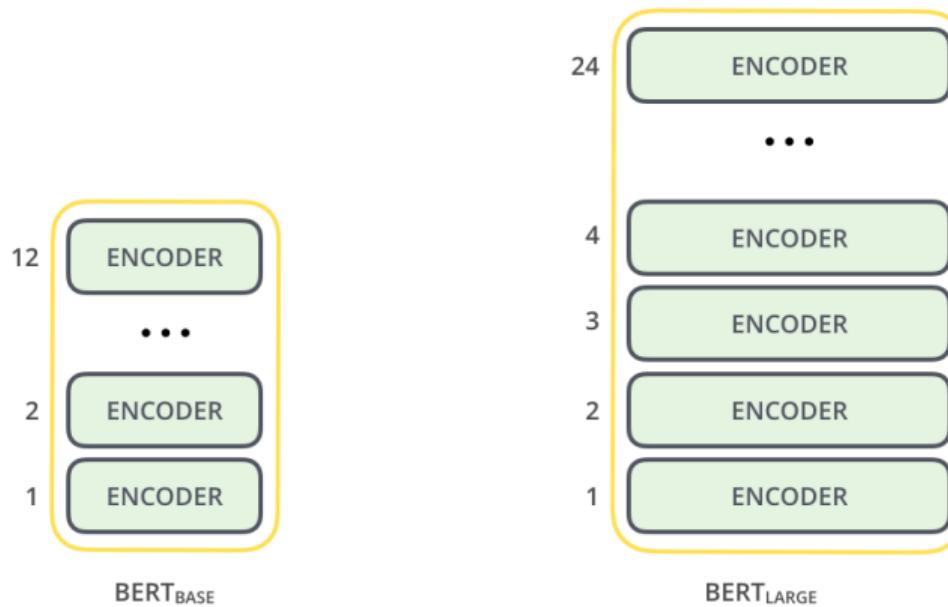


# Bidirectional Encoder Representations from Transformers (BERT)

---

Because BERT is a language model, it only needs  
**encoders** so I won't talk about decoders here.

# BERT Is Basically a Trained Transformer Encoder Stack



## BERT Involves Many Uninteresting Numbers

BERT Comes in Two Sizes: BASE and LARGE. The parameters of these models are as follows:

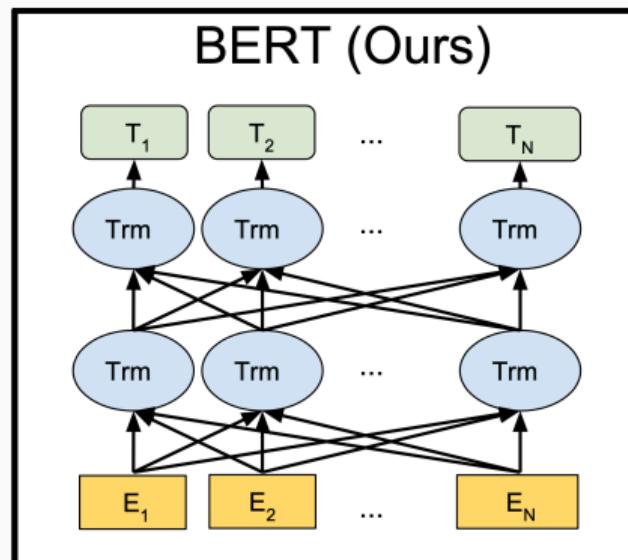
---

	BERT <sub>BASE</sub>	BERT <sub>LARGE</sub>
encoder layers	12	24
hidden units	768	1024
attention heads	12	16

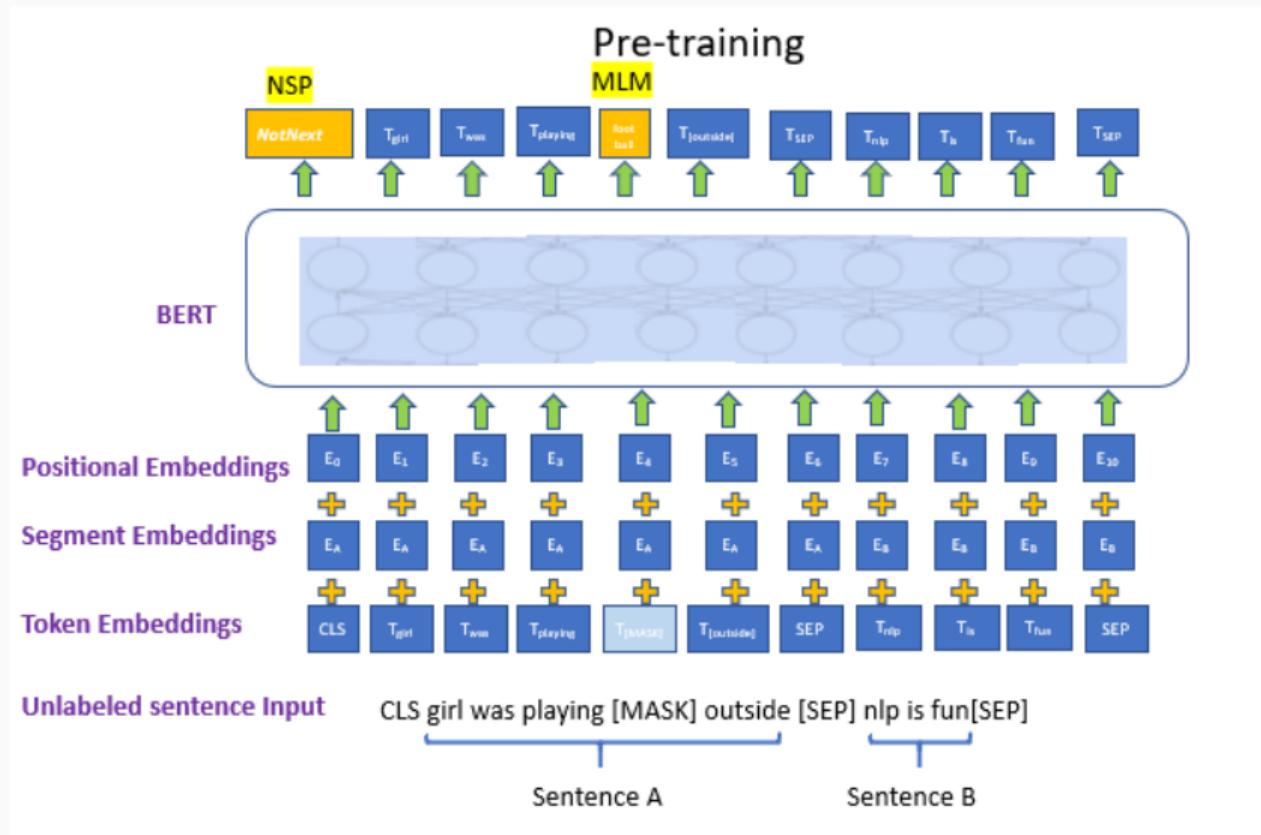
---

# BERT Is Bidirectional (or, Rather, Non-Directional)

- Traditional LMs are unidirectional—given some sequence of words  $w_1, w_2, \dots, w_t$ , they try to predict  $w_{t+1}$ .
- BERT uses the whole sequence to predict words in that sequence.



# Pretraining BERT: MLM and NSP



## The Cloze Task

- The **cloze** task comes from psycholinguistics (the branch of linguistics and cognitive science that uses experimental methods to study how language works in human brains).
- It is a fill-in-the-blank task:

**He drove the yellow \_\_\_\_\_ into the front of our house.**

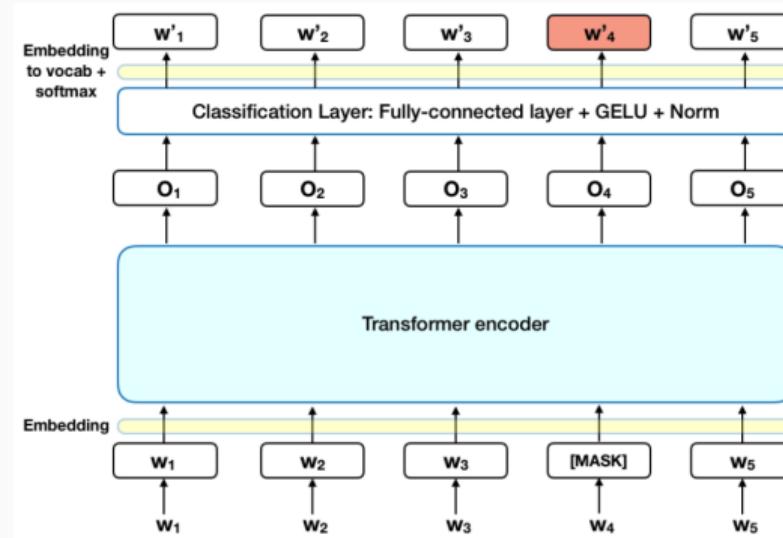
- Subjects are presented with these frames and asked to fill in the missing words
- This allows experimenters to assess what a speaker understands about grammar, semantics, etc.
- According to the original BERT paper, this task provided the inspiration for BERT's masked language modeling (MLM) training task.
- But compare various kinds of denoising algorithms.

## BERT is Trained to Perform Masked Language Modeling

- Since, in BERT's innards, everything is literally connected to everything, each word would **indirectly see itself**
- Enter masking!
- In training, random words are concealed from BERT using the [MASK] token
- BERT is trained to guess these masked words
- Take the sentence: “the girl was playing outside.”
  - Suppose that a [MASK] token is inserted before *outside*, thus masking it
  - We then have “the girl was playing [MASK] *outside*.”
  - The model now cannot “see” *outside*; instead it is trained to generate it based on context
  - This is part of why BERT can be so good at representing words according to the context in which they occur (not just the distribution of identically-spelled tokens)

# How BERT Is Trained to Perform Masked Language Modeling

1. Add classification layer
2. Multiply output vectors by embedding layer → vocabulary dimension
3. Calculate probabilities using softmax



## GELU: A Note

GELU (from the last slide) is a Gaussian Error Linear Unit:

$$\text{GELU}(x) = xP(X \leq x) = x\Phi(x) = x \cdot \frac{1}{2} \left[ 1 + \operatorname{erf}\left(\frac{x}{\sqrt{2}}\right) \right]$$

where  $\Phi(x)$  is the standard Gaussian cumulative distribution function.

GELU is an activation function that is similar, in some respects, to ReLU. Unlike ReLU, GELU non-linearly weights inputs by their percentile (ReLU gates inputs by their sign).

## BERT is also Trained to Perform Next Sentence Prediction

- Suppose you are given two pairs of sentences:
  1. Good pair
    - 1.1 Pickles are delicious.
    - 1.2 However, cucumbers make me burp.
  2. Bad pair
    - 2.1 Pickles are delicious.
    - 2.2 The NASDAC was up 10 points today amid heavy trading.
- For each pair, you (a human) can tell whether the first and the second are likely to occur in sequence
- BERT is trained to do the same
- As a result, BERT learns patterns that exist above the level of the sentence

## Embeddings from BERT

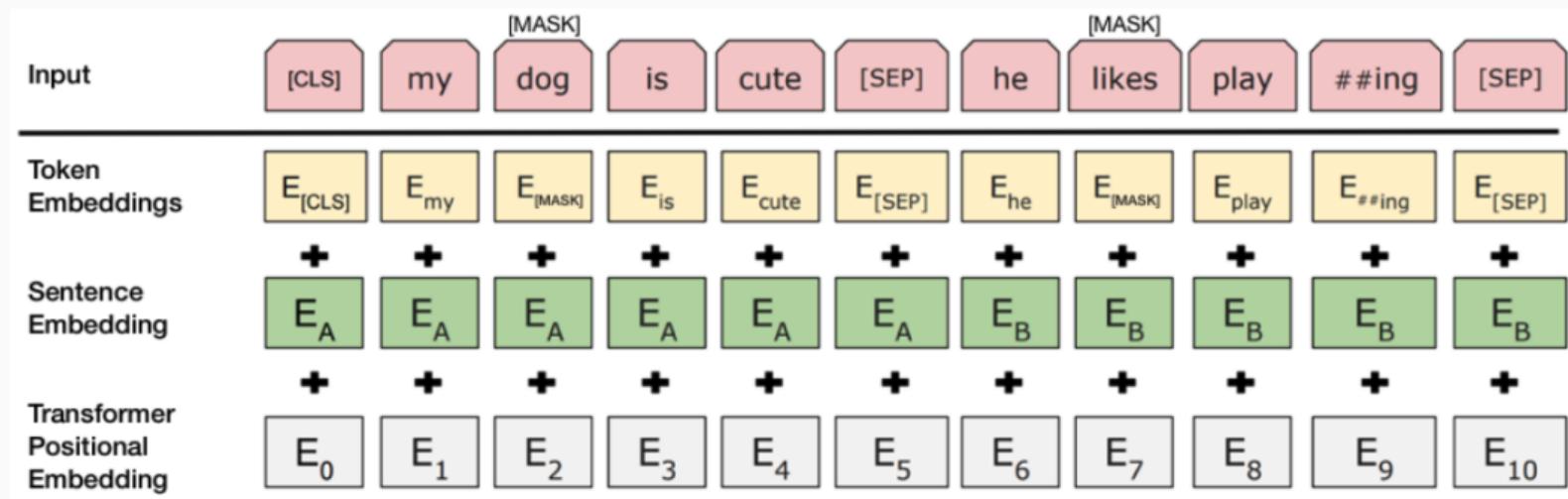
---

## BERT Has Its Own Tokenization Scheme

- BERT has its own tokenizer and tokenization scheme
- You must use it
- BERT has a few special tokens:
  - [CLS] “Classify” (occurs at the beginning of inputs)
  - [SEP] “Separator” (occurs between sentences)
  - [MASK] “Mask” (indicates masked words—training only)
  - [PAD] “Padding” (used when, for example, padding batching sequences of different lengths)
  - [UNK] “Unknown” (used when a token cannot be converted into an ID)

# Preparing Input for BERT

BERT needs inputs to have the follow structure:

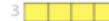
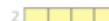


Getting this is simpler than you might think.

## How Do You Get Embeddings from BERT?

- $\text{BERT}_{\text{BASE}}$  has 12 layers
- The output of each base is an embedding
- Choose one of these, or some combination

# To Concatenate or to Sum?

What is the best contextualized embedding for “Help” in that context? For named-entity recognition task CoNLL-2003 NER		
		Dev F1 Score
12 	First Layer	Embedding 
• • •	Last Hidden Layer	12 
7 		12  +
6 	Sum All 12 Layers	2  + 1  = 
5 		95.5
4 	Second-to-Last Hidden Layer	11 
3 		95.6
2 	Sum Last Four Hidden	12  +
1 		11  +
		10  +
		9  = 
Help	Concat Last Four Hidden	96.1

# Questions?