

# Final Project Logistics

DL Day (Poster Session) will take place on Thursday, Dec 11th

If you're capstoning, you don't need to fill out a form (for me)

Groups of 3-4 students, assigned to a mentor

# **Project Options**

**Core to all**: You need to train a neural network (i.e., using an LLMfor something new does not count if no training occurs)

- 1. Re-implement a paper, extend to different dataset, run new experiments, etc.
- 2. Try something new (i.e., ask a research question)
- 3. Al enabled: Use LLM coding tools to create a DL project. (higher expectations in terms of production and scope of project)

# **Project Options**

**Core to all**: You need to train a neural network (i.e., using an LLMfor something new does not count if no training occurs)

- 1. Re-implement a paper, extend to different dataset, run new experiments, etc.
- 2. Try something new (i.e., ask a research question)
- 3. Al enabled: Use LLM coding tools to create a DL project. (higher expectations in terms of production and scope of project)

# Recap

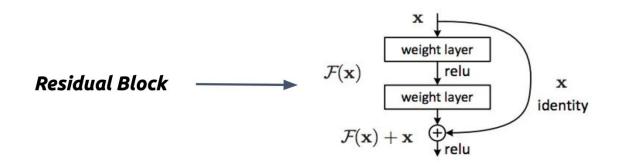
Convolutions provide spatial reasoning capabilities by **breaking symmetries** between inputs (not all pixels are connected to all hidden units).

This means that the **spatial layout** of the original image **matters** for convolutional neural
networks

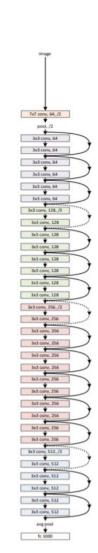
#### **More Complicated Networks**

#### ResNet:

Lots of layers, tons of learnable parameters Avoids Vanishing Gradient problem



K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. arXiv preprint arXiv:1512.03385, 2015.



# Recap

Convolutions provide spatial reasoning capabilities by **breaking symmetries** between inputs (not all pixels are connected to all hidden units).

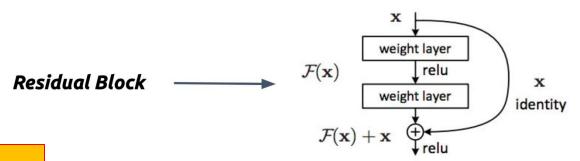
This means that the **spatial layout** of the original image **matters** for convolutional neural
networks

For a new data type (images), we introduced a new network architecture (convolutions) to fix an issue with MLPs.

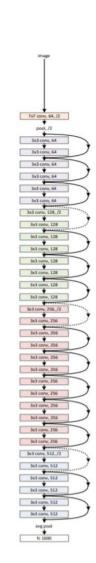
#### More Complicated Networks

#### ResNet:

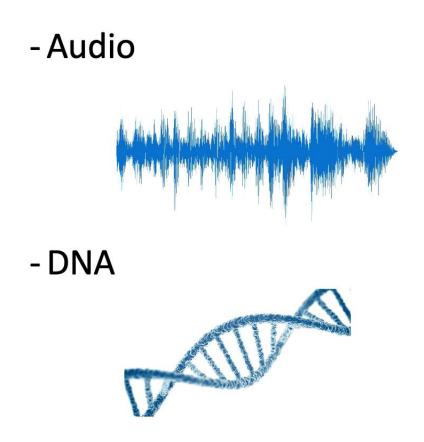
Lots of layers, tons of learnable parameters Avoids Vanishing Gradient problem



K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. arXiv preprint arXiv:1512.03385, 2015.



# New data type: sequences



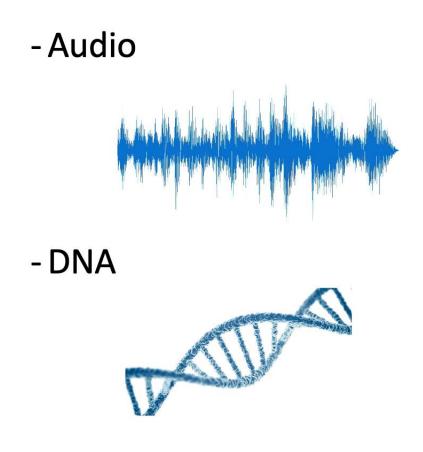




- Weather



# New data type: sequences







- Weather



What is the data property here that we could leverage?

# Natural Language

"language that has developed naturally in use"

# Natural Language

"language that has developed naturally in use"

Compare to constructed or formal language

```
-code: for i in range(50):
```

-math: 52 + 94 = 147

-logic: A ^ B -> C (if A and B, then C)

# Natural Language

In this class: **sequence of** *words* 

"They went to the grocery store and bought bread, peanut butter, and jam."



Input: X

I do not want sour cream in my burrito



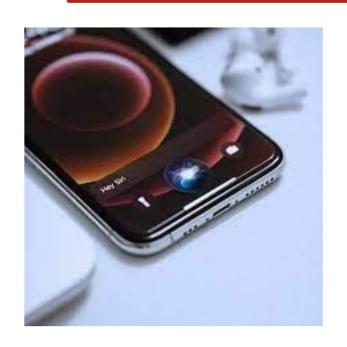
Function: f



No quiero crema agrea en mi burrito

Output: Y

Example of prediction?



Input: X

I do not want sour cream in my burrito



Function: f



No quiero crema agrea en mi burrito

Output: Y

Example of classification?

Input: X

"The story telling was erratic and, at times, slow"

"Loved the diverse cast of this movie" Output: Y

"Good review?"





Function: f



Example of prediction?

```
"They went to the grocery store and bought... bread?

milk?

rock?
```

**Generating artificial sentences:** Here each word is a discrete unit; predicting the next part of the sequence means predicting words

# Language models

Definition: Probability distribution over strings in a language.

Exponentially-many strings means each string has very low probability

Relative probabilities are meaningful:

P("they went to the store") >> P("butter dancing rock")

# Language models logic: leverage sentence structure

P(any sequence) is determined by P(the words in the sequence).

# Language models logic: leverage sentence structure

P(any sequence) is determined by P(the words in the sequence).

Said differently, we can represent a sequence as  $w_1, w_2, ... w_n$ , and

$$P(w_1, w_2, ... w_n) = P(w_1) * P(w_2|w_1) * P(w_3|w_1, w_2) * \cdots P(w_n|w_1 ... w_{n-1})$$

# Language models logic: leverage sentence structure

P(any sequence) is determined by P(the words in the sequence).

Said differently, we can represent a sequence as  $w_1, w_2, ... w_n$ , and

$$P(w_1, w_2, ... w_n) = P(w_1) * P(w_2|w_1) * P(w_3|w_1, w_2) * \cdots P(w_n|w_1 ... w_{n-1})$$

"The probability of a sentence is the product of the probabilities of each word given the previous words"

This is an application of the chain rule for probabilities

# Language models: weird & cool!

Model trained on the King James Bible, Structure and Interpretation of Computer Programs, and some of Eric S. Raymond's writings:

- The righteous shall inherit the land, and leave it for an inheritance unto the children of Gad according to the number of steps that is linear in b.
- And this I pray, that your love may abound yet more and more like a controlled use of shared memory.

(King James Programming)

https://kingjamesprogramming.tumblr.com/

Discriminitive models learn conditional probabilities P(Y|X=x)

Discriminitive models learn conditional probabilities P(Y|X=x)Given some features, what is the probability of a label?

Discriminitive models learn conditional probabilities P(Y|X=x)Given some features, what is the probability of a label?

Generative models learn joint probabilities P(X,Y)

Discriminitive models learn conditional probabilities P(Y|X=x)Given some features, what is the probability of a label?

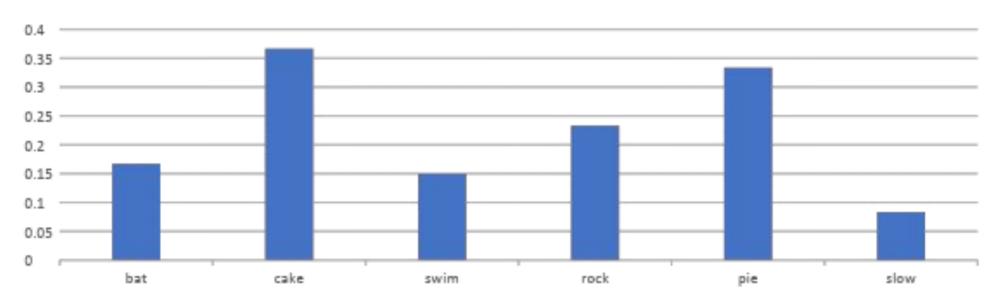
Generative models learn joint probabilities P(X,Y) models how a signal was generated

# Language models: the math

At each step, we look at a probability distribution for what the *next* word might be.

They went to the grocery store and bought ..



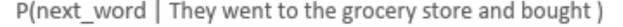


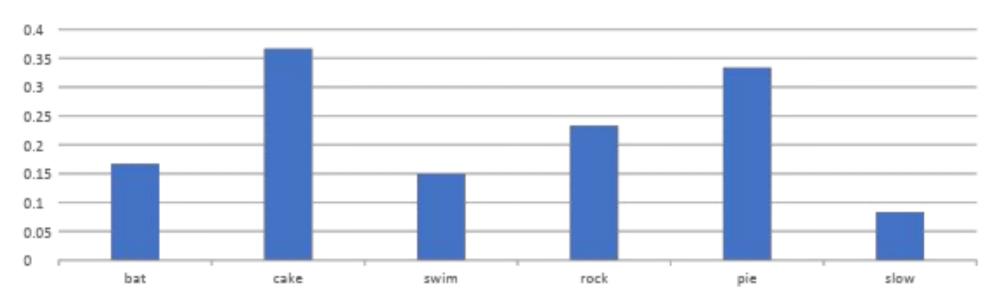
But first, how do we represent sentence?

# Language models: the math

At each step, we look at a probability distribution for what the *next* word might be.

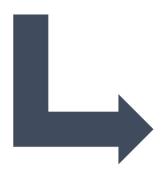
They went to the grocery store and bought ..





## Natural language: tokenization

"They went to the grocery store and bought bread, peanut butter, and jam."



```
["they", "went", "to", "the",
"grocery", "store", "and",
"bought", "bread", "peanut",
"butter", "and", "jam"]
```

# Natural language: tokenization

"They went to the grocery store and bought bread, peanut butter, and jam."

- Consistent casing
- Strip punctuation
- One word is one token
- -Split on spaces

```
["they", "went", "to", "the",
"grocery", "store", "and",
"bought", "bread", "peanut",
"butter", "and", "jam"]
```

# Aside: Tokenization itself can be challenging...

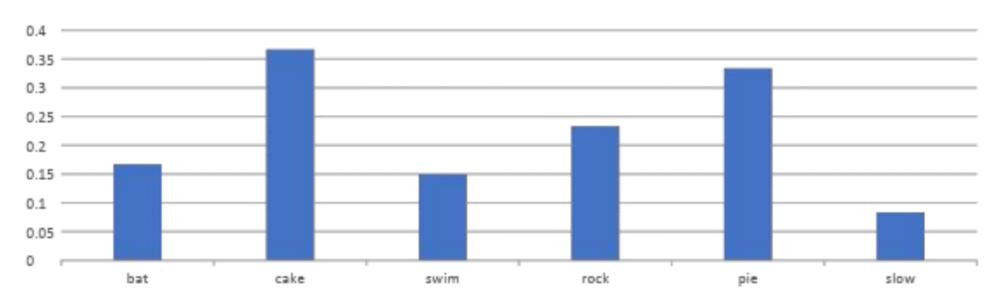
- A lot easier in English than other languages (e.g. Chinese)
  - Chinese is character-based; words & phrases have different character lengths
  - No spaces

# Language models: the math

At each step, we look at a probability distribution for what the *next* word might be.

They went to the grocery store and bought ..



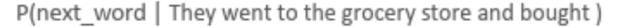


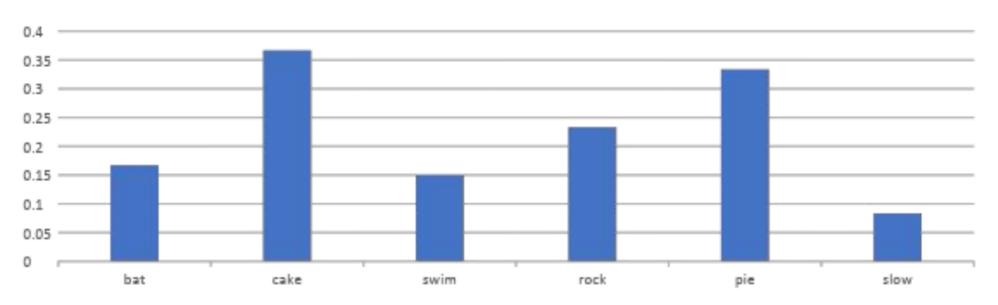
# Language models: the math

How do we know which words to calculate probabilities for?

At each step, we look at a probability distribution for what the *next* word might be.

They went to the grocery store and bought ..





# Vocabularies: Defining a finite set of words

Vocabularies: the set of all words "known" to the model

#### Why?

- We need a finite set of words in order to define a discrete distribution over it.

# Vocabularies: Defining a finite set of words

Vocabularies: the set of all words "known" to the model

#### Why?

- We need a finite set of words in order to define a discrete distribution over it.

#### How?

- Choose a hyperparameter vocab\_size for how many words the model should know
- Keep only the vocab\_size with most frequent words replace everything else with "UNK"

#### Vocabularies: how

- Original sentence:
  - "They galloped to the Ratty for dinner, and ate exactly seventy-three waffle fries and chocolate peamilk."

#### Vocabularies: how

- Original sentence:
  - "They galloped to the Ratty for dinner, and ate exactly seventy-three waffle fries and chocolate peamilk."
- Tokenized:

```
- ["they", "galloped", "to", "the", "ratty", "for", "dinner", "and", "ate", "exactly", "seventy-three", "waffle", "fries", "and", "chocolate", "peamilk"]
```

#### Vocabularies: how

- Original sentence:
  - "They galloped to the Ratty for dinner, and ate exactly seventy-three waffle fries and chocolate peamilk."
- Tokenized:

```
- ["they", "galloped", "to", "the", "ratty", "for", "dinner", "and", "ate", "exactly", "seventy-three", "waffle", "fries", "and", "chocolate", "peamilk"]
```

#### Vocabularies: how

- Original sentence:
  - "They galloped to the Ratty for dinner, and ate exactly seventy-three waffle fries and chocolate peamilk."
- Tokenized:

```
- ["they", "galloped", "to", "the", "ratty", "for", "dinner", "and", "ate", "exactly", "seventy-three", "waffle", "fries", "and", "chocolate", "peamilk"]
```

#### - UNKed:

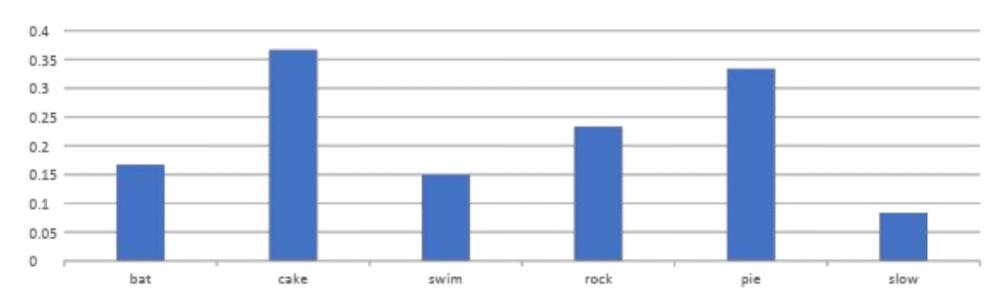
```
- ["they", "UNK", "to", "the", "UNK", "for", "dinner", "and", "ate", "exactly", "UNK", "waffle", "fries", "and", "chocolate", "UNK"]
```

#### Language models: the math

At each step, we look at a probability distribution for what the *next* word might be.

They went to the grocery store and bought ..





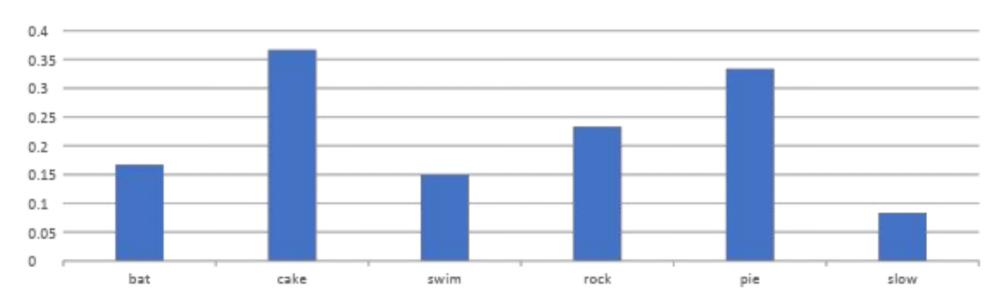
#### Language models: the math

How to calculate the probability for words in our vocabulary?

At each step, we look at a probability distribution for what the *next* word might be.

They went to the grocery store and bought ..





- Goal: predict next word given a preceding sequence
  - $P(word_n | word_1, word_2, ... word_{n-1}) = \frac{Count(word_1, word_2, ... word_{n-1}, word_n)}{Count(word_1, word_2, ... word_{n-1})}$

- Goal: predict next word given a preceding sequence
  - $-P(\boldsymbol{word_n}|\ word_1, word_2, ... word_{n-1}) = \frac{Count(word_1, word_2, ... word_{n-1}, \boldsymbol{word_n})}{Count(word_1, word_2, ... word_{n-1})}$
- Example task: predict the next word
  - he danced \_\_\_\_

- Goal: predict next word given a preceding sequence
  - $-P(\boldsymbol{word_n}|\ word_1, word_2, ... word_{n-1}) = \frac{Count(word_1, word_2, ... word_{n-1}, \boldsymbol{word_n})}{Count(word_1, word_2, ... word_{n-1})}$
- Example task: predict the next word
  - he danced \_\_\_\_
- Strategy: iterate through all words in vocabulary, and calculate

- Our training sentences were:
  - "She danced happily"
  - "They sang beautifully"
  - "He danced energetically"
  - "He sang happily"
  - "She danced gracefully"
- -"He danced \_ \_ \_ "
- "He danced happily"

 $\frac{Count(he \ danced < word >)}{Count(he \ danced)}$ 

- Our training sentences were:
  - "She danced happily"
  - "They sang beautifully"
  - "He danced energetically"
  - "He sang happily"
  - "She danced gracefully"
- -"He danced"
- "He danced happily" Has 0 probability

 $Count(he\ danced\ < word\ >)$ Count(he danced)

- Our training sentences were:
  - "She danced happily"
  - "They sang beautifully"
  - "He danced energetically"
  - "He sang happily"
  - "She danced gracefully"
- -"He danced"
- "He danced happily" Has 0 probability

 $Count(he\ danced\ < word\ >)$ Count(he danced)

Why doesn't this work?

- Our training sentences were:
  - "She danced happily"
  - "They sang beautifully"
  - "He danced energetically"
  - "He sang happily"
  - "She danced gracefully"
- -"He danced \_ \_ \_ "
- "He danced happily"

Has 0 probability

```
\frac{Count(he \ danced < word >)}{Count(he \ danced)}
```

Why doesn't this work?

This strategy depends on having instances of sentence prefixes.

### LM implementation: N-gram counting

Improvement: N-gram model – only look at N words at a time

## LM implementation: N-gram counting

Improvement: N-gram model – only look at N words at a time (in this case, bigrams look at 2 words at a time)

- -"She danced happily"
- -"They sang beautifully"
- -"He danced energetically"
- -"He sang happily"
- -"She danced gracefully"

### LM implementation: N-gram counting

Improvement: N-gram model – only look at N words at a time (in this case, bigrams look at 2 words at a time)

```
-"danced happily"
-"sang beautifully"
-"danced energetically"
-"sang happily"
-"danced gracefully"
```

"He danced happily" now has 1/3 probability!

But what if the answer was "He danced beautifully"?

# LM implementation

Problem: it's impossible for the training set to have *every possible valid* sequence of words!

Let's try to learn a better numerical representation

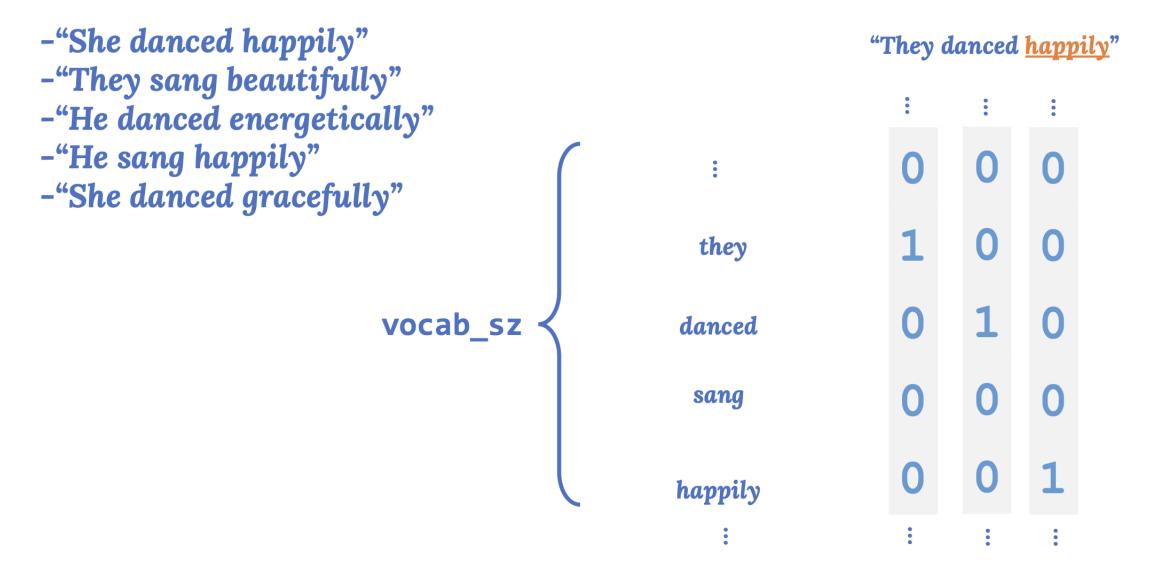
# LM implementation

Problem: it's impossible for the training set to have *every possible valid* sequence of words!

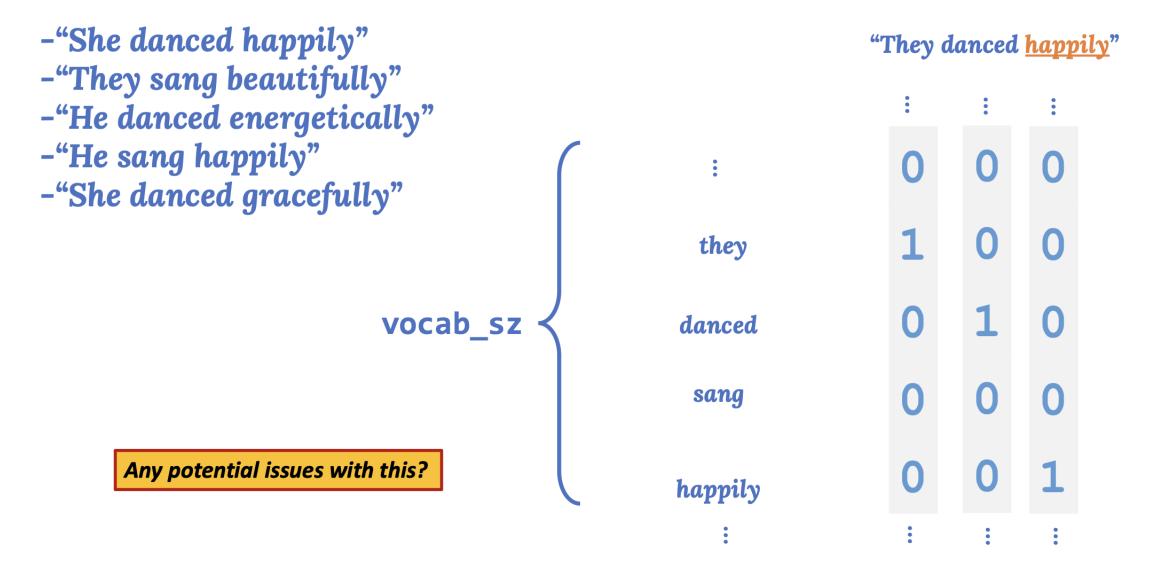
Let's try to learn a better numerical representation

What is the simplest thing you can think of?

# LM implementation: Simple approach



# LM implementation: Simple approach



# LM implementation

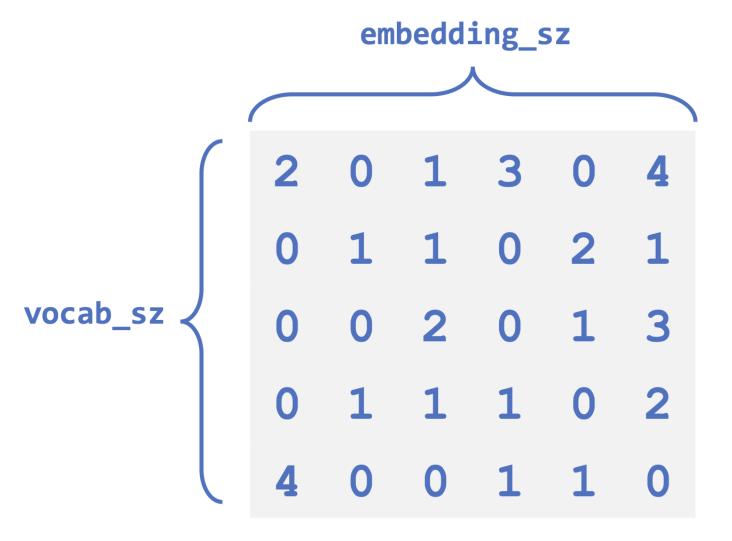
Problem: it's impossible for the training set to have every possible valid sequence of words!

Can we learn a better numerical representation which associates related words with one another?

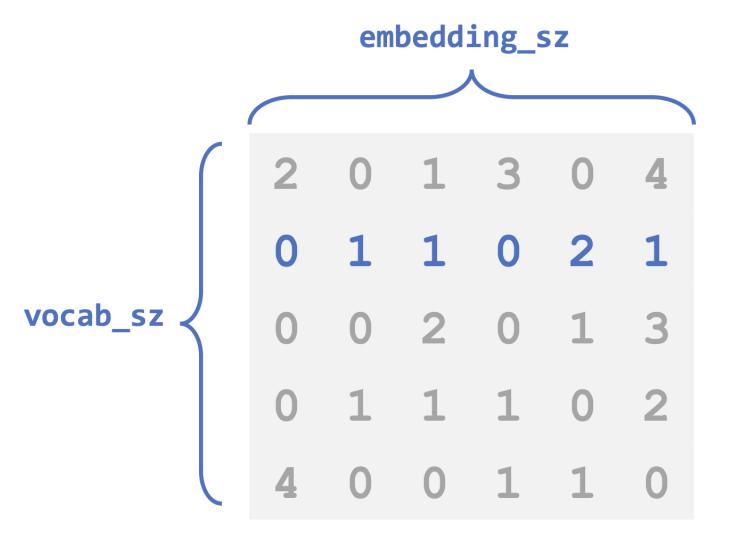




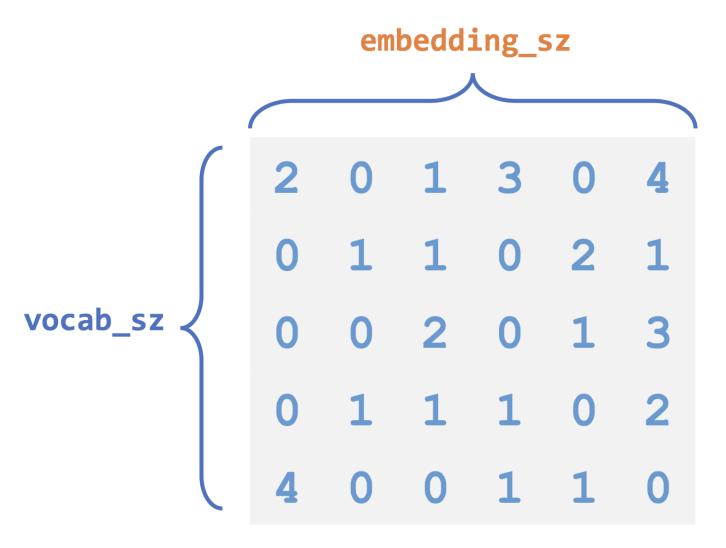




-2d matrix: vocab\_sz
x embedding\_sz



- -2d matrix: vocab\_sz
  x embedding\_sz
- -each word correspondsto an index, or word ID- hence the vocab\_szdimension



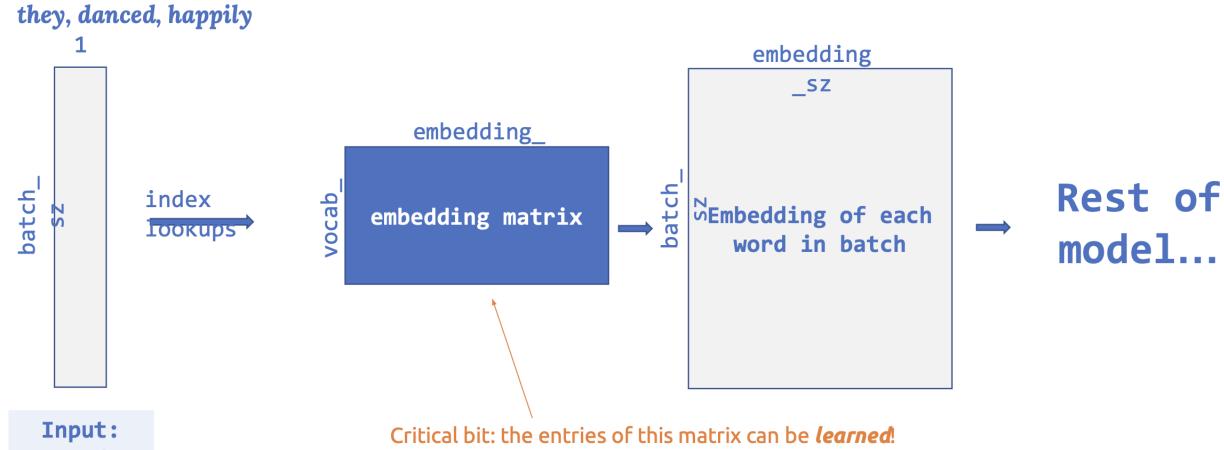
- -2d matrix: vocab\_sz
  x embedding\_sz
- -each word correspondsto an index, or word ID –hence the vocab\_szdimension
- embedding\_sz is a hyperparameter

## LM implementation: deep learning

Deep learning helps solve this!

We can learn an *embedding matrix* that associates *related* words with one another for solving a prediction task.

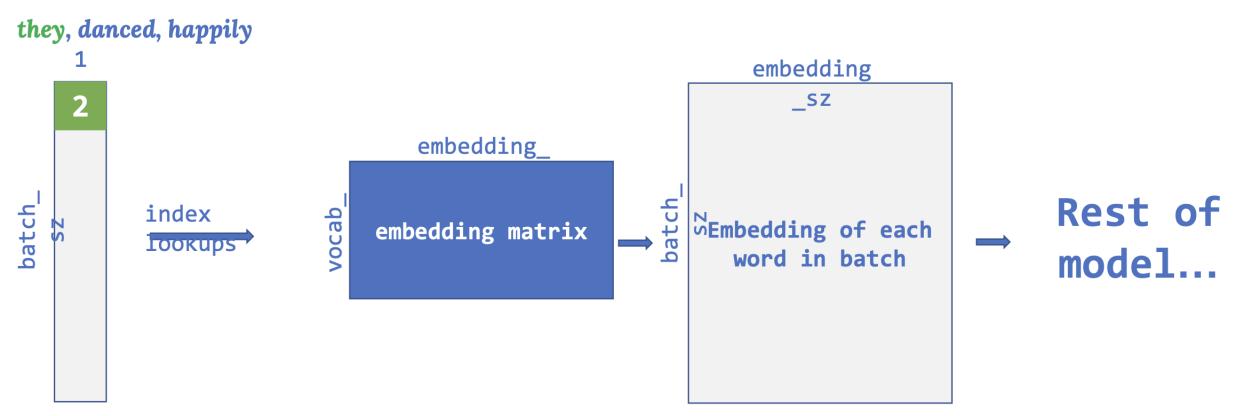
If you want to input a [batch of] words into a neural net, this is how:



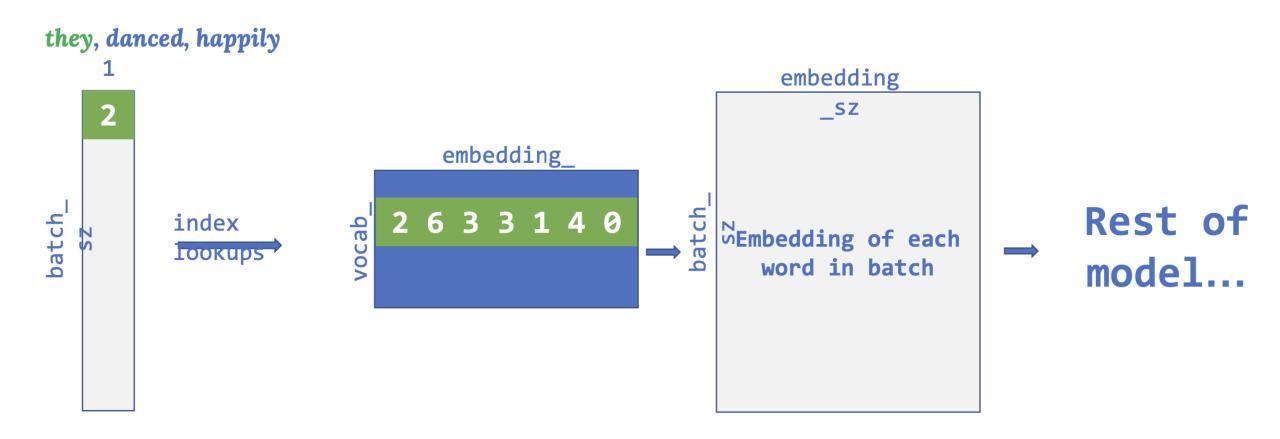
The network learns what word embeddings are most effective for performing its task

indicac

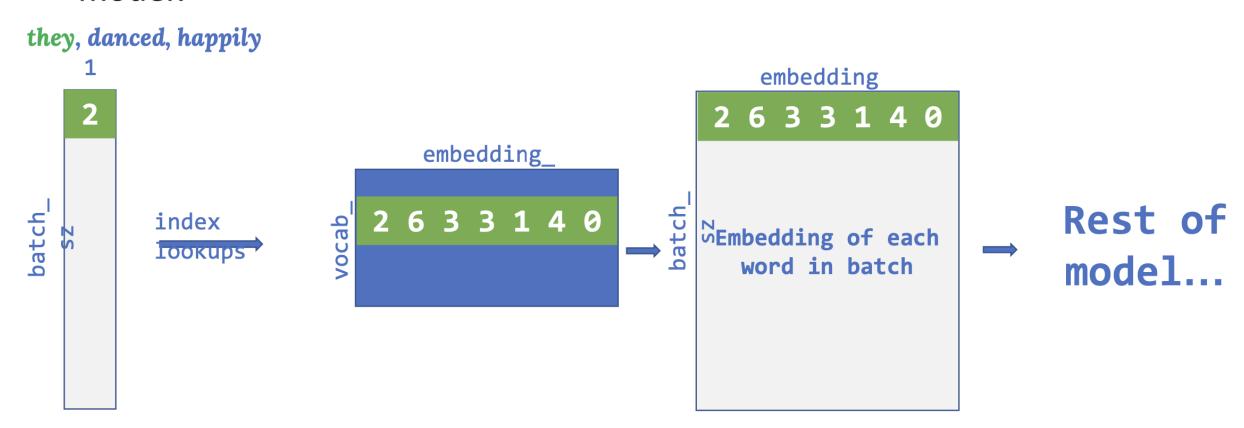
Let's look at the 0<sup>th</sup> word in this batch; its ID in the vocab is 2.



So we look at row 2 of the embedding matrix.



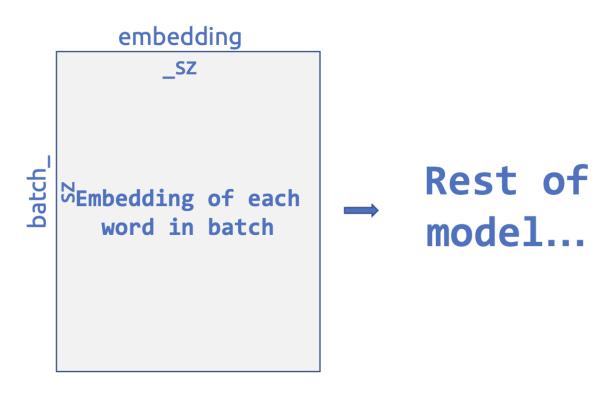
We can then pull out this embedding so we can use it in the rest of the model!



In tensorflow, we can use

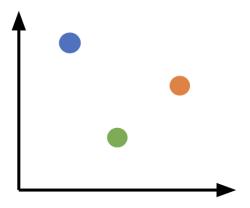
tf.nn.embedding\_lookup

which takes in an embedding matrix and a list of indices, and returns the embedding corresponding to each index.



 Each row in the matrix can be viewed as a vector in vector space

Example 2-D vector space:

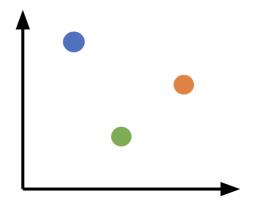


Vocab size: 3

Embed size: 2

- Each row in the matrix can be viewed as a vector in vector space
- "Embedding": We're embedding a non-Euclidian entity [a word] into Euclidian space

Example 2-D vector space:

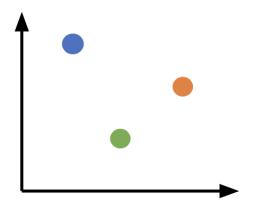


Vocab size: 3

Embed size: 2

- Each row in the matrix can be viewed as a vector in vector space
- "Embedding": We're embedding a non-Euclidian entity [a word] into Euclidian space
- Each row represents the "embedding" for a single word

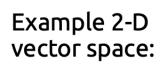
Example 2-D vector space:

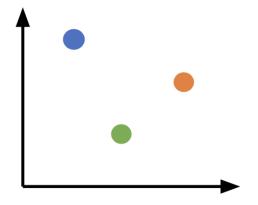


Vocab size: 3

Embed size: 2

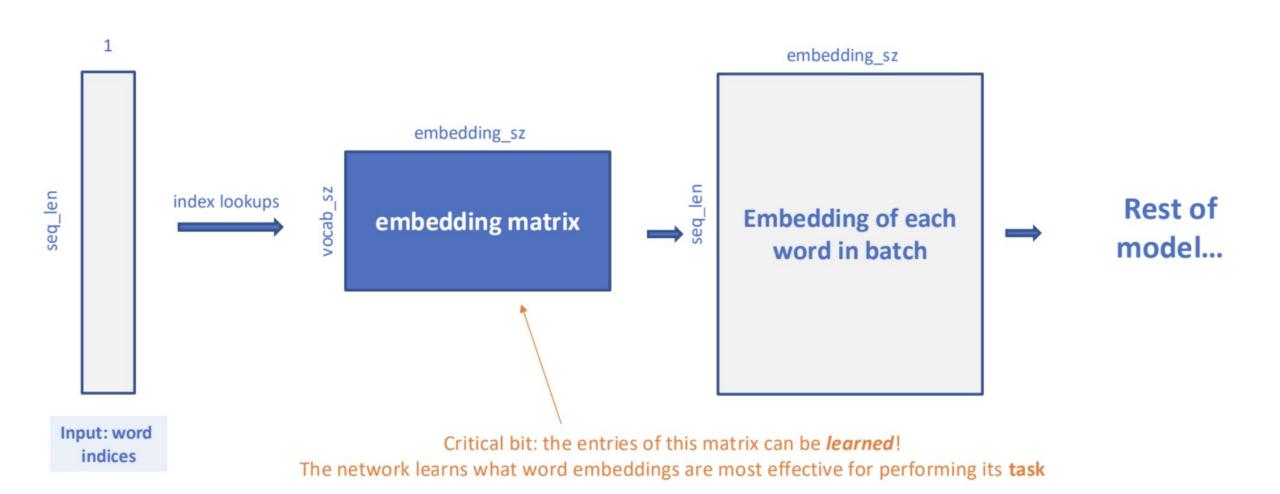
- Each row in the matrix can be viewed as a vector in vector space
- "Embedding": We're embedding a non-Euclidian entity [a word] into Euclidian space
- Each row represents the "embedding" for a single word
- This has pretty remarkable properties!





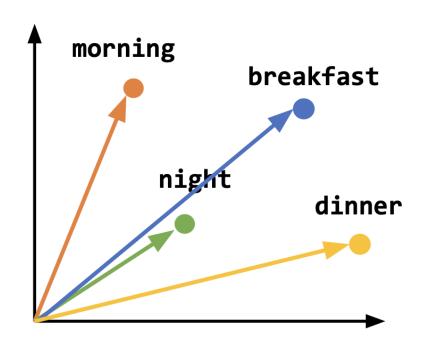
Vocab size: 3

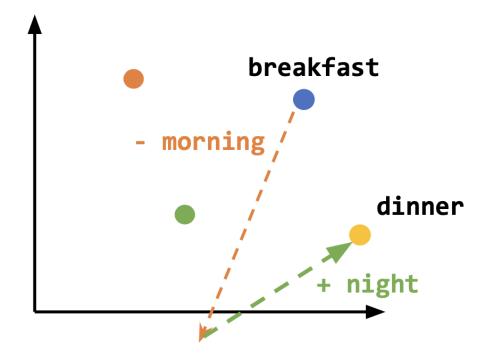
Embed size: 2



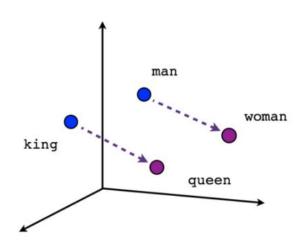
# Vector arithmetic in the embedding matrix

Demo here: https://turbomaze.github.io/word2vecjson/





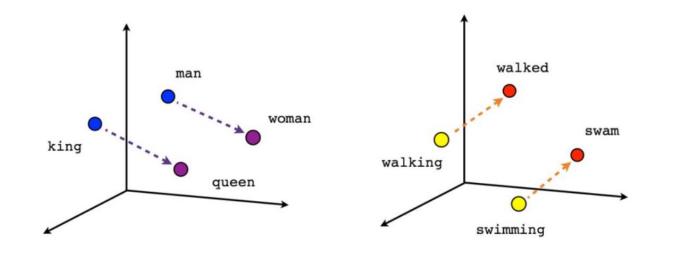
# More 'semantic directions' in embedding space



Male-Female

```
E(queen) - E(king) ≈
E(woman) - E(man)
```

# More 'semantic directions' in embedding space



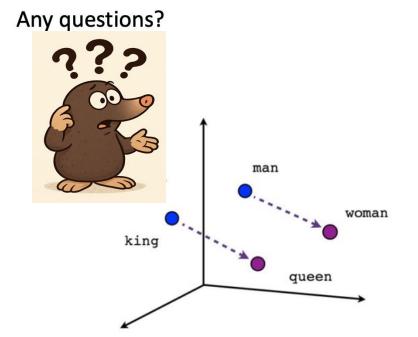
Male-Female

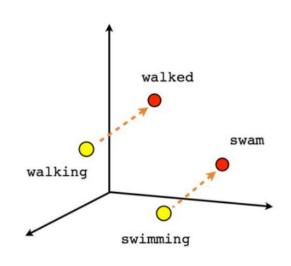
 $E(queen) - E(king) \approx$ E(woman) - E(man)

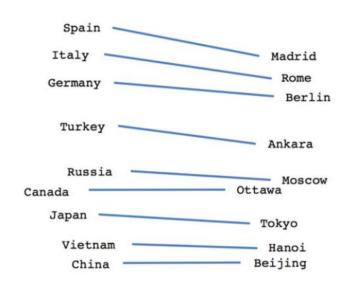
Verb tense

E(walked) - E(walking) ≈ E(swam) - E(swimming)

### More 'semantic directions' in embedding space







Male-Female

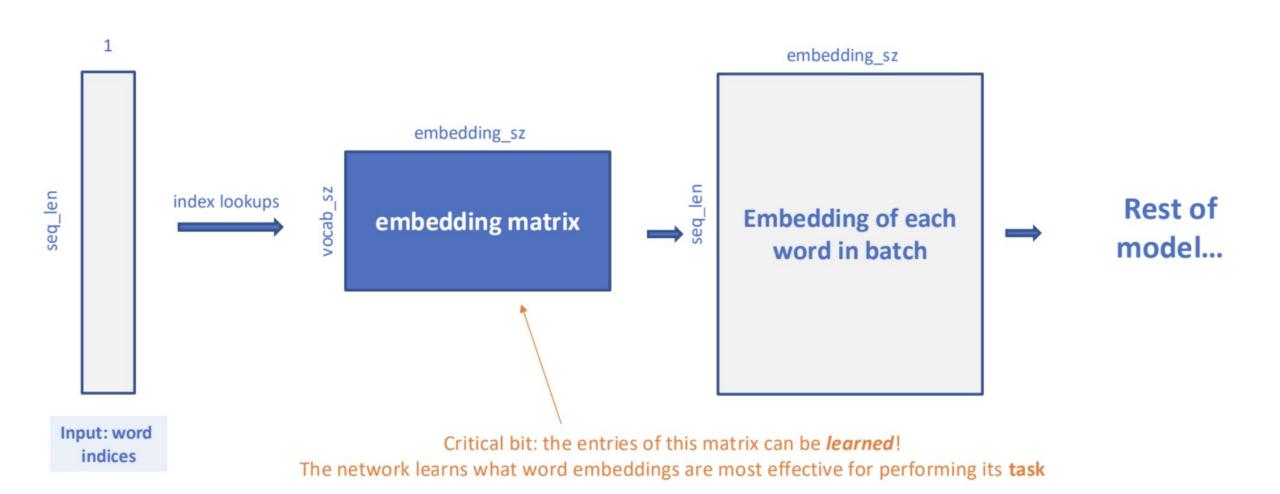
 $E(queen) - E(king) \approx$ 

E(woman) - E(man)

Verb tense

Country-Capital

# Using the Embedding Matrix in a Network

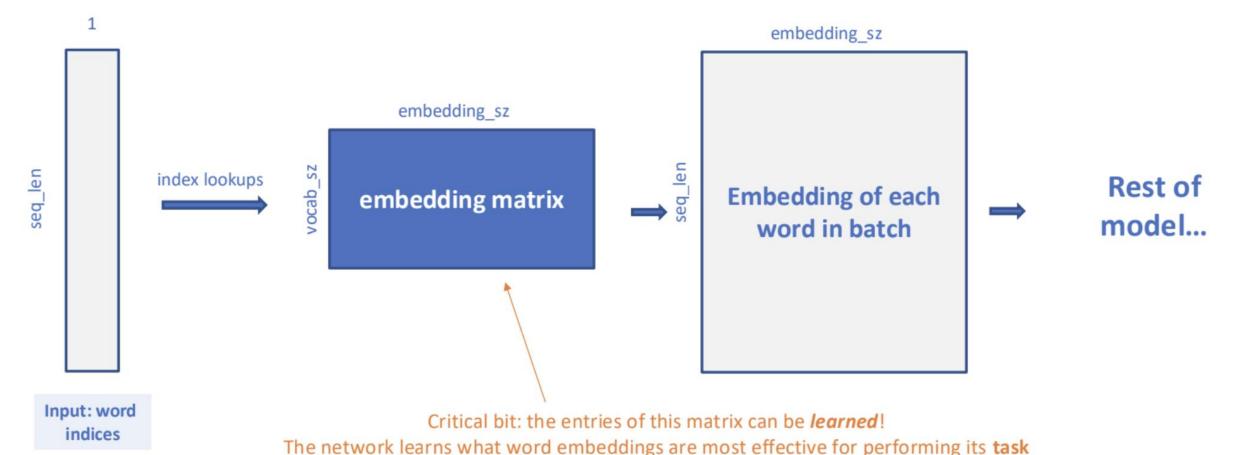


#### Using the Embedding Matrix in a Network

Say in the middle of training, the model sees:

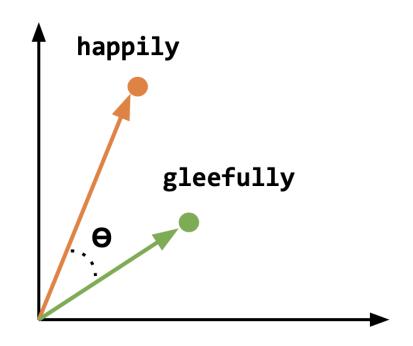
Then the model sees a lot of "danced gleefully"

P("happily"|"They Danced") = HighP("gleefully"|"They Danced") = Low



# Quantifying "similarity"

cosine similarity = 
$$\cos(\theta) = \frac{A \cdot B}{||A|| ||B||} = \frac{\sum_{i=1}^{n} A_i B_i}{\sqrt{\sum_{i=1}^{n} A_i^2} \sqrt{\sum_{i=1}^{n} B_i^2}}$$

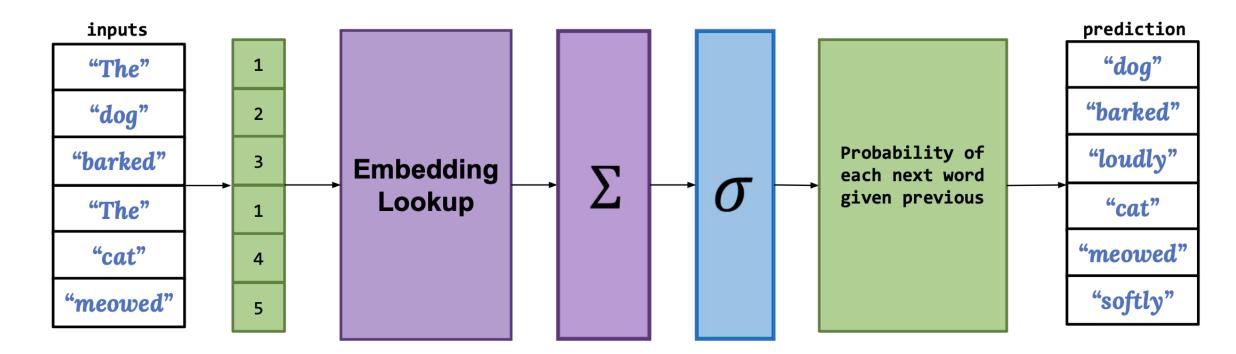


$$cos(0^\circ) = 1$$
  
 $cos(90^\circ) = -0.448$   
 $cos(180^\circ) = -0.598$ 

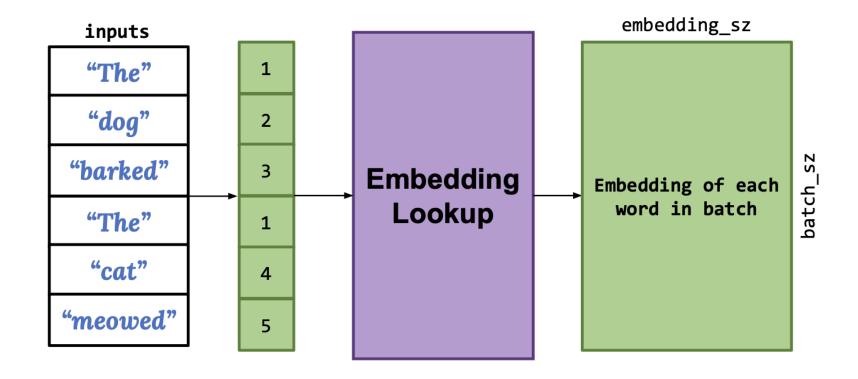
#### Limitations of the N-gram model

What issues do we run into using feed-forward N-gram models?

Let's look at bigram model and count the number of weights.

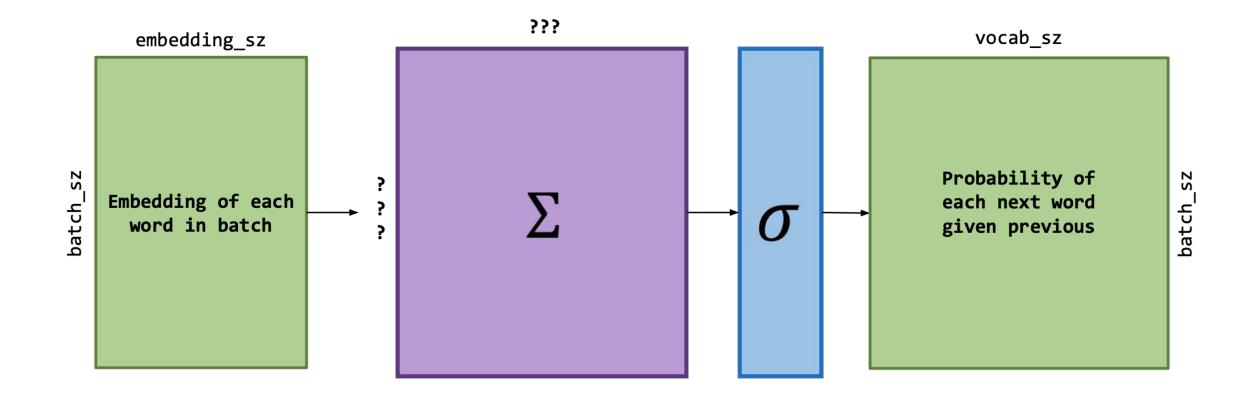


To preform embedding lookup on our entire batch, we just need one embedding matrix of size: (vocab\_sz, embedding\_sz)



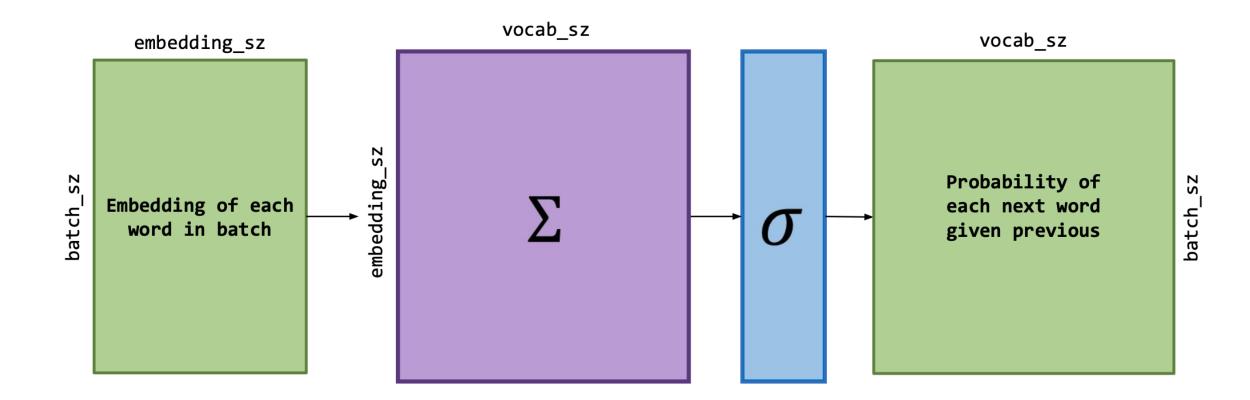
What size do we need the linear layer to be in order to map:

(batch\_sz, embedding\_sz)  $\times$  (???, ???)  $\rightarrow$  (batch\_sz, vocab\_sz)

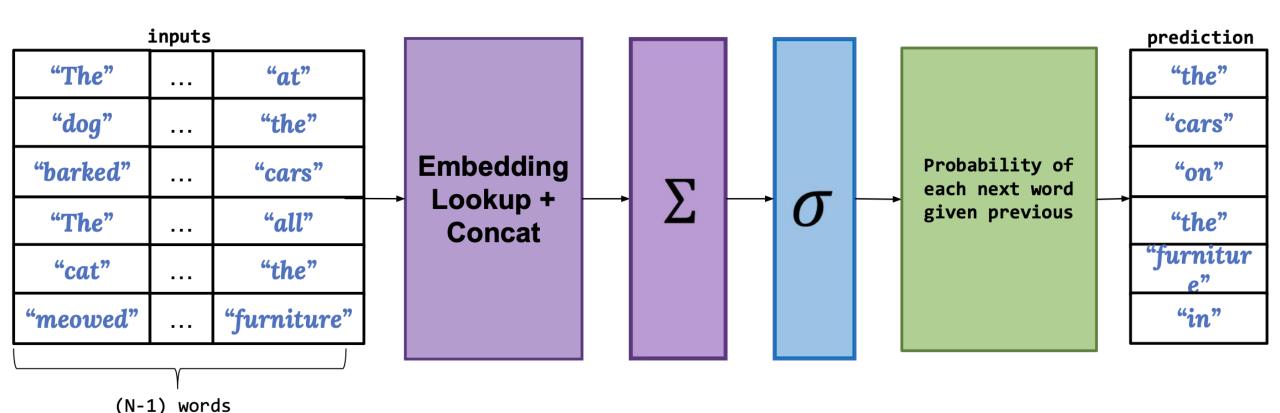


What size do we need the linear layer to be in order to map:

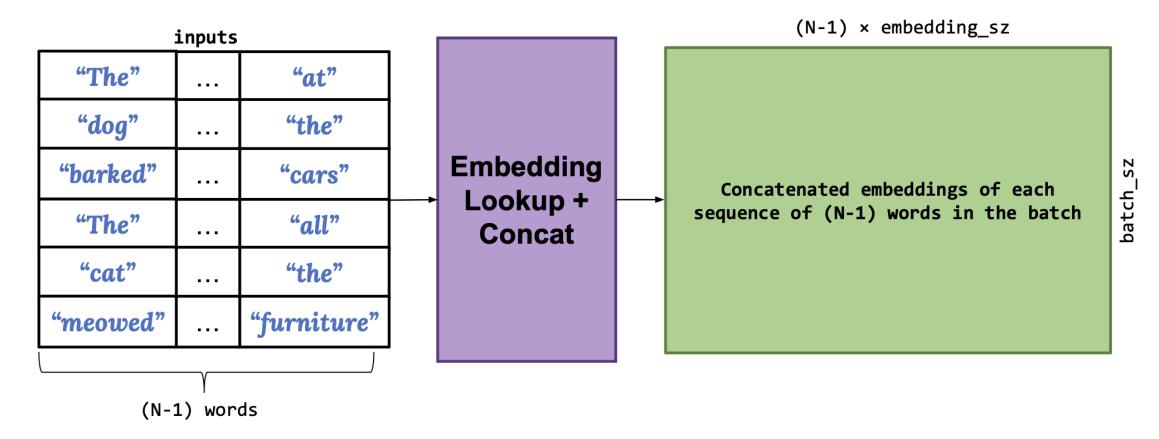
(batch\_sz, embedding\_sz)  $\times$  (???, ???)  $\rightarrow$  (batch\_sz, vocab\_sz)



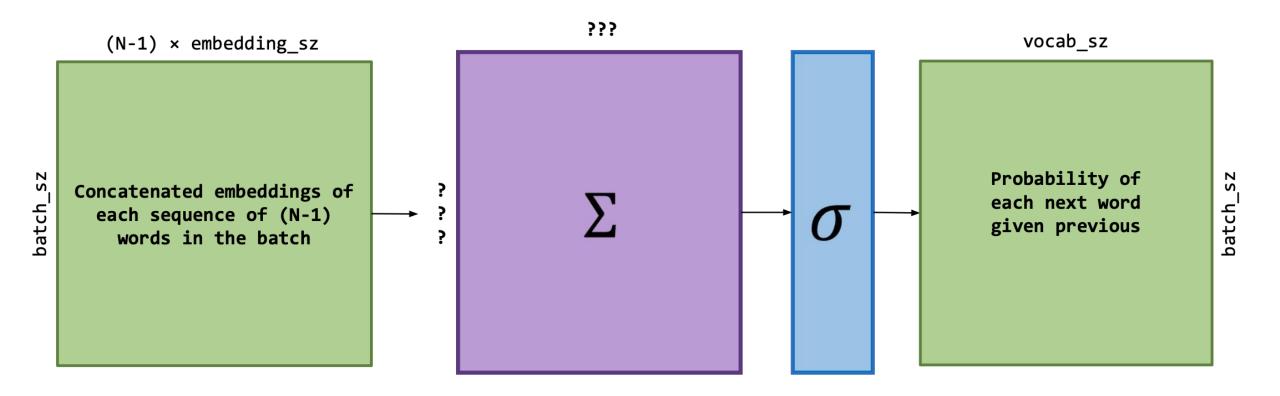
So what happens in the N-gram case?



Embedding lookup + Concatenation still requires only one embedding matrix of size: (vocab\_sz, embedding\_sz)



But what happens to our feed forward layer?



#### Limitations of the N-gram model

- What issues do we run into using feed-forward N-gram models?
  - As the size of **N** increases, the number of weights needed for the linear layer becomes far too large.

#### Limitations of the N-gram model

- What issues do we run into using feed-forward N-gram models?
  - As the size of **N** increases, the number of weights needed for the linear layer becomes far too large.
  - Using a fixed N creates problems with the flexibility of our model.

# Lack of Flexibility with N-grams

We would like for our language model to be more aware of context when deciding on how many words in the past to consider as "relevant".

For example, we can see that at some parts of the sentence below, smaller N-gram models should be sufficient to make predictions:



### "The dog barked at one of the cats."





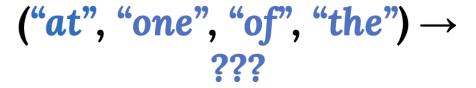
#### Lack of Flexibility with N-grams

We would like for our language model to be more aware of context when deciding on how many words in the past to consider as "relevant".

But when we look at other portions, common phrases and sequences of words may make it impossible to have any idea what should come next.



#### "The dog barked at one of the cats."





#### Lack of Flexibility with N-grams

We would like for our language model to be more aware of context when deciding on how many words in the past to consider as "relevant".

But when we look at other portions, common phrases and sequences of words may make it impossible to have any idea what should come next.

#### "The dog barked at one of the cats."

We want our model to recognize these patterns and dynamically adapt how it makes a prediction based on context.



# Limitations of the N-gram model

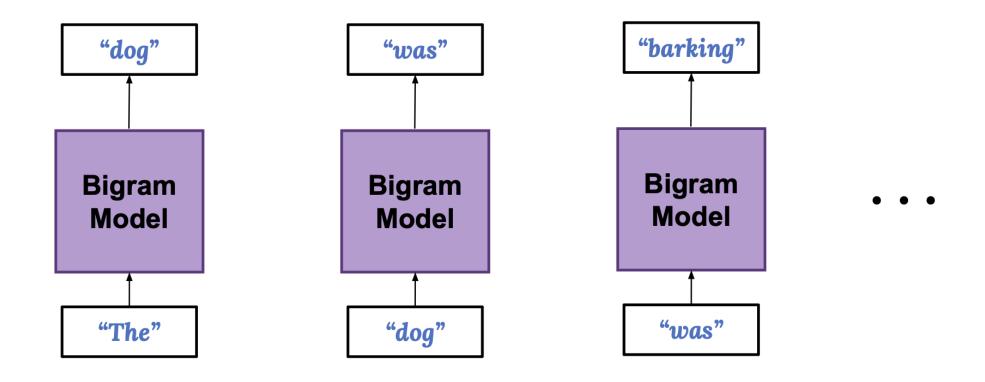
What problems do we run into using Feed Forward N-gram models?

1. As the size of N increases, the number of weights needed for the linear layer becomes far too large.

2. Using a fixed **N** creates problems with the flexibility of our model.

We need a solution that is both computationally cheap and more dynamic in terms of its memory of previously seen words.

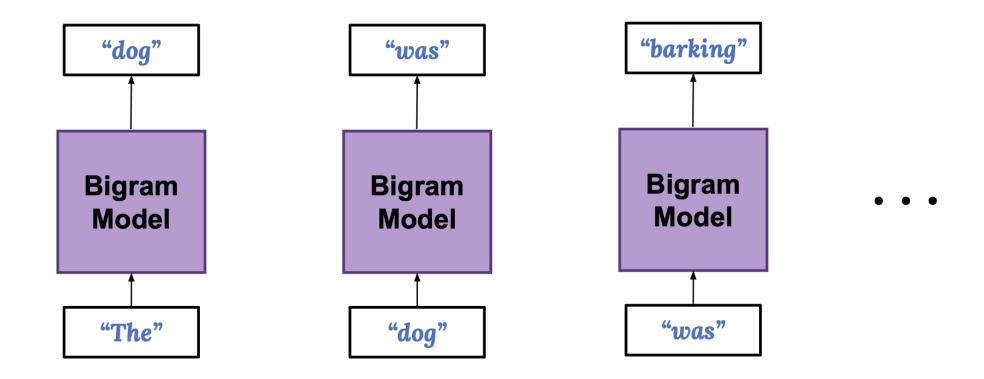
Let's revisit the bigram model and see several iterations of prediction using a bigram model:



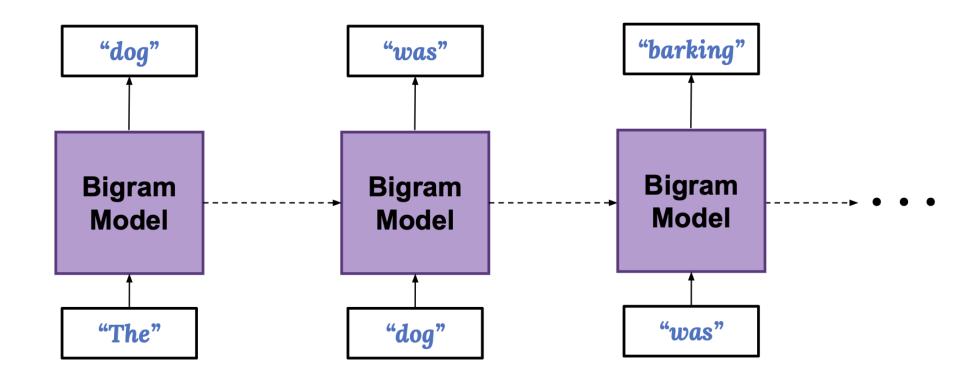
Any ideas?

#### New Approach

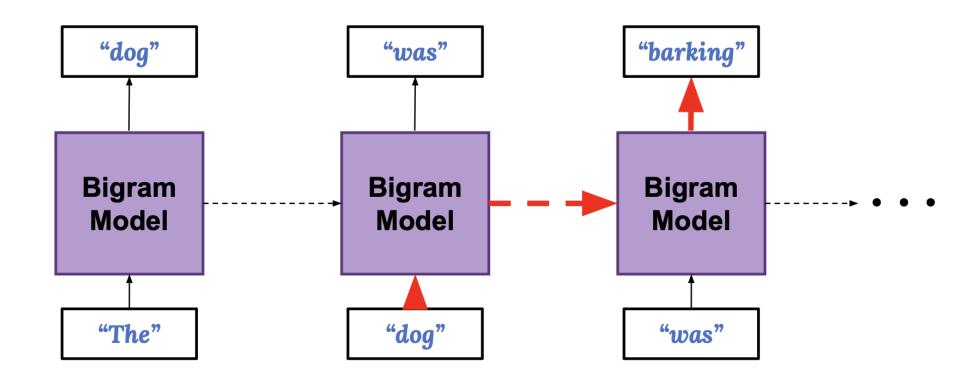
Ideally, we would like to be able to keep "memory" of what words occurred in the past.



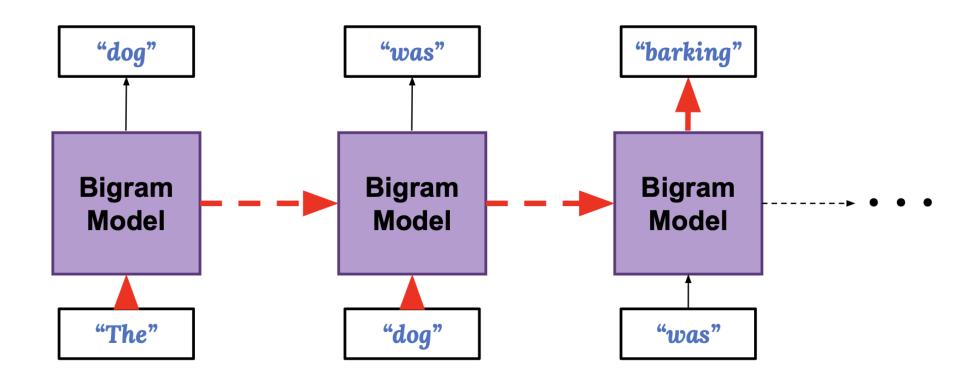
What if we sequentially passed information from our previous bigram block into our next block?



If we follow the information flow, we see that when predicting "barking", we have some way of knowing that "dog" was previously observed:



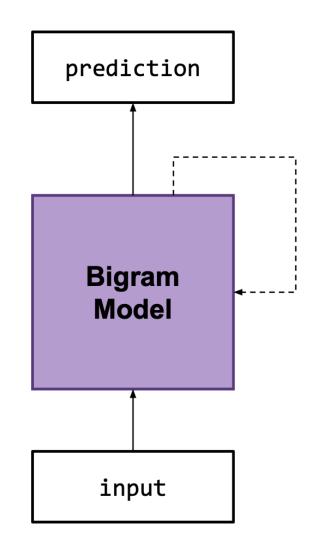
In fact, we even have a way of knowing that "The" was observed!



We can represent this relationship using only one bigram block and connection that feeds from the output of the model back into the input.

We call this connection a *recurrent* connection.

We call the previous representation the "unrolled" representation.



#### Recurrent Neural Network (RNN)

Recurrent Neural Networks are networks in the form of a directed cyclic graph.

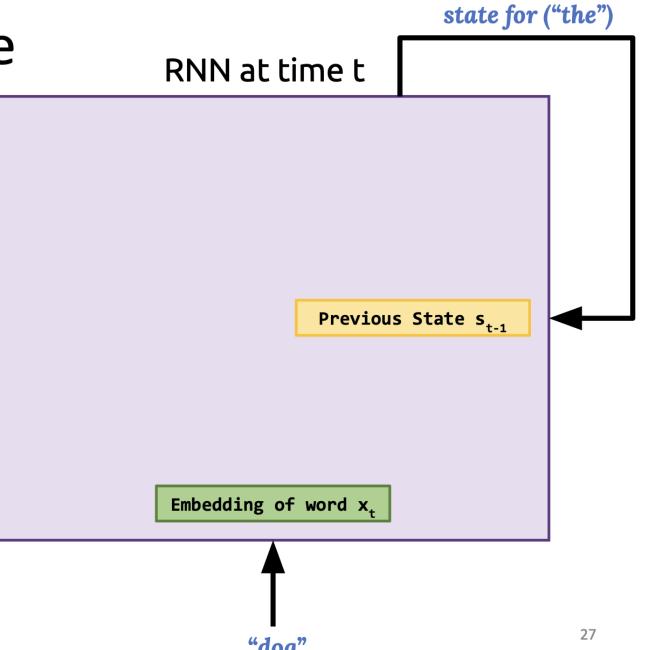
They pass previous *state* information from previous computations to the next.

They can be used to process sequence data with relatively low model complexity when compared to feed forward models.

The block of computation that feeds its own output into its input is called the RNN cell.

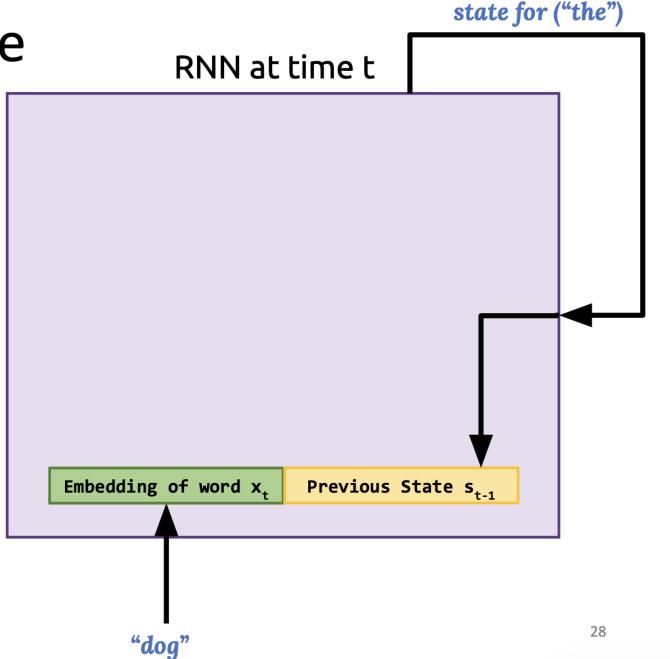
Let's see how we can build one!

At each step of our RNN, we will get an input word, and a state vector from the previous cell.



At each step of our RNN, we will get an input word, and a state vector from the previous cell.

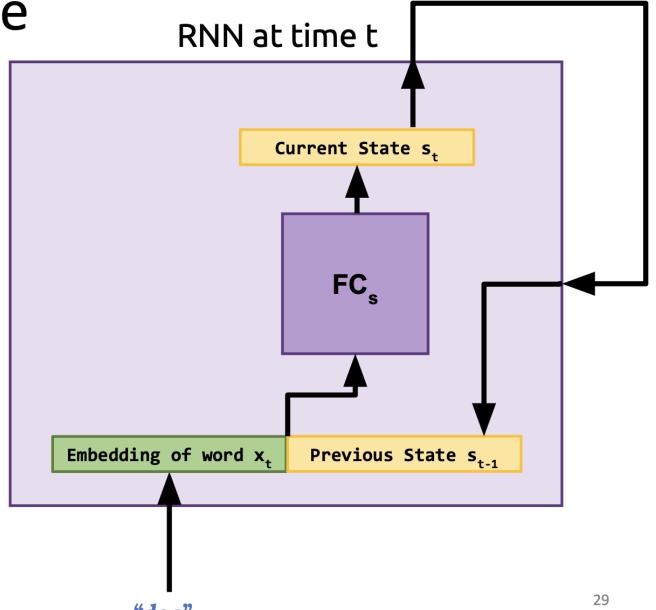
We then concatenate the embedding and state vectors.



At each step of our RNN, we will get an input word, and a state vector from the previous cell.

We then concatenate the embedding and state vectors.

We use a fully connected layer to compute the next state

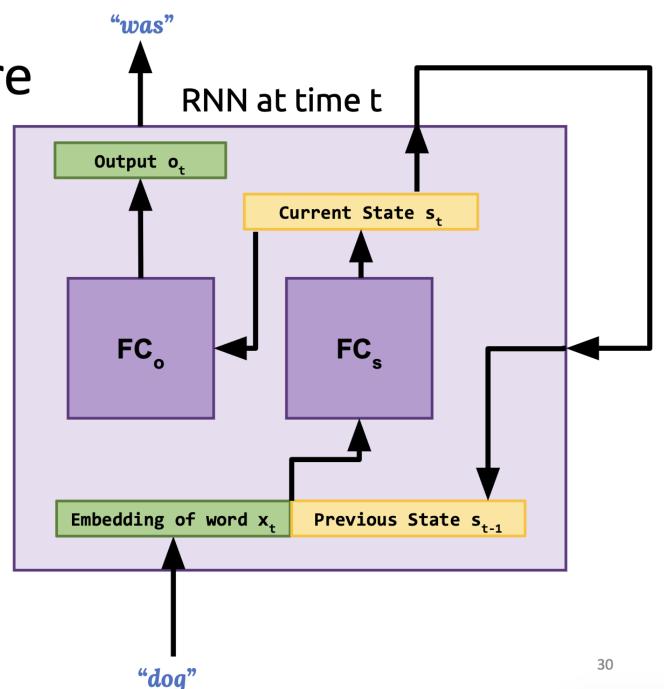


At each step of our RNN, we will get an input word, and a state vector from the previous cell.

We then concatenate the embedding and state vectors.

We use a fully connected layer to compute the next state

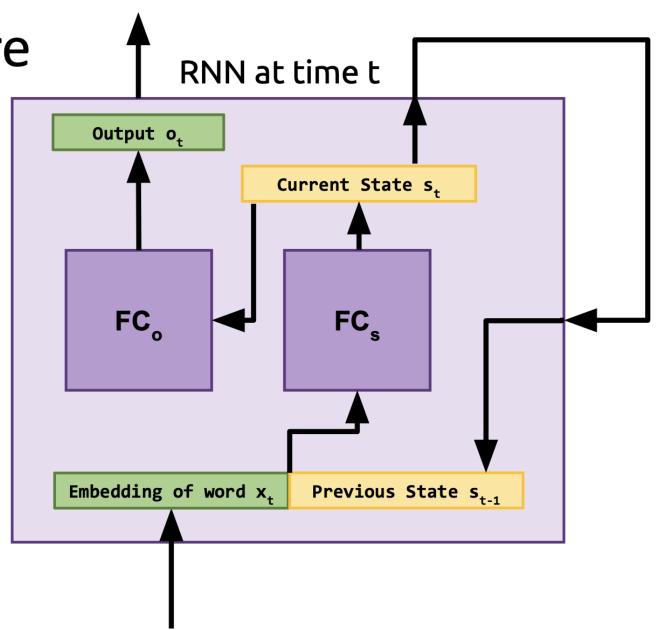
We use another connected layer to get the output.



We can represent the RNN in with the following equations:

$$s_t = \rho \big( (e_t, s_{t-1}) W_r + b_r \big)$$

$$o_t = \sigma(s_t W_o + b_o)$$



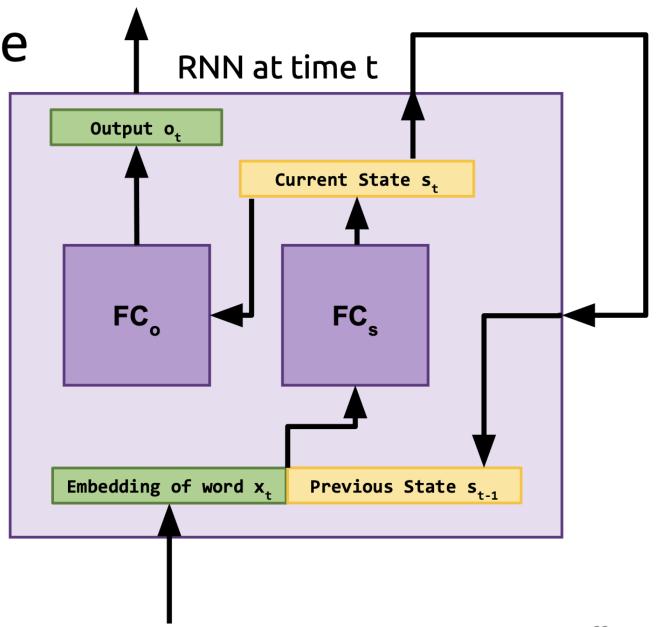
We can represent the RNN in with the following equations:

$$s_t = \rho ((e_t, s_{t-1})W_r + b_r)$$

$$o_t = \sigma(s_t W_o + b_o)$$

Nonlinear activations (e.g. sigmoid, tanh)

Any questions?



We can represent the RNN in with the following equations:

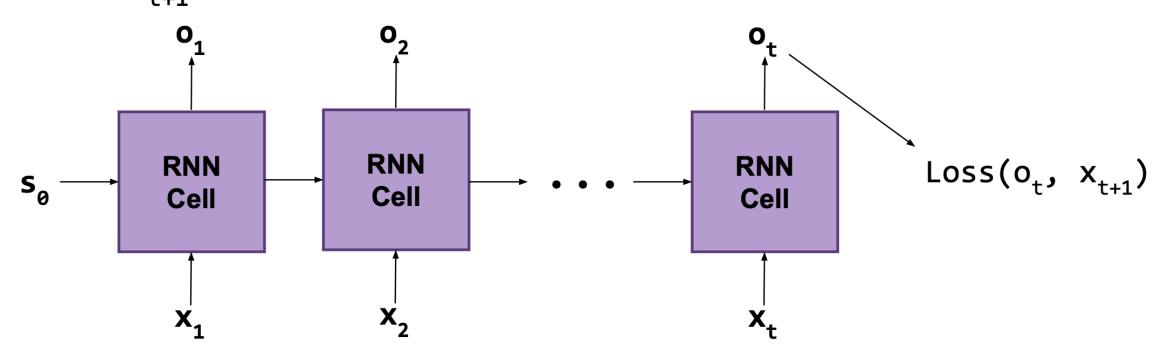
$$s_t = \rho \big( (e_t, s_{t-1}) W_r + b_r \big)$$

$$o_t = \sigma(s_t W_o + b_o)$$

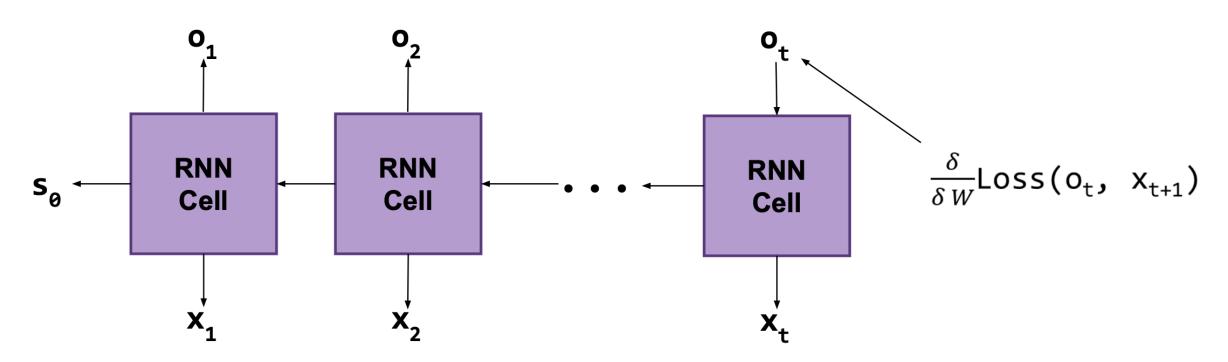
This brings up an immediate question: what is  $s_0$ ?

Typically, we initialize  $s_0$  to be a vector of zeros (i.e. "initially, there is no memory of any previous words")

We can calculate the cross entropy loss just as before since for any sequence of input words  $(x_1, x_2, ..., x_t)$ , we know the true next word  $x_{+11}$ 

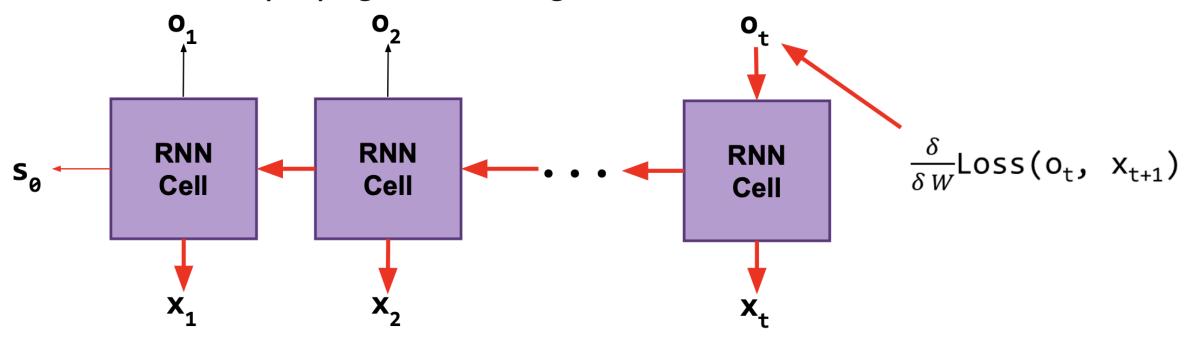


But what happens when we differentiate the loss and backpropagate?

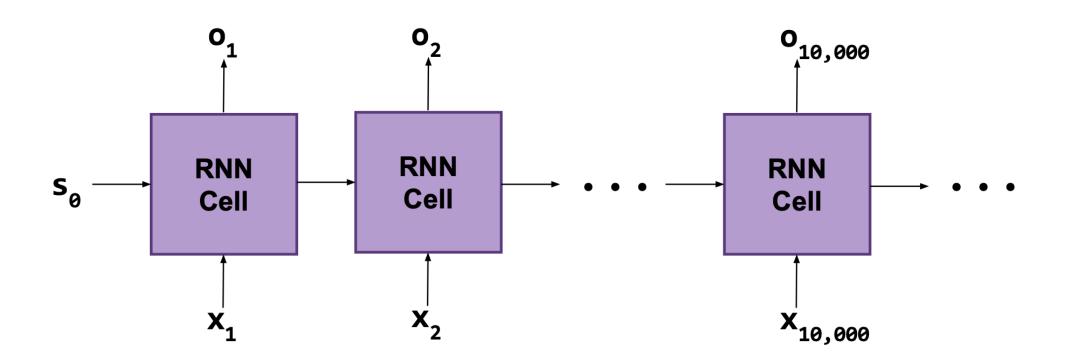


Not only do our gradients for  $o_t$  depend on  $x_t$ , but also on all of the previous inputs.

We call this backpropagation through time.



With this architecture, we can run the RNN cell for as many steps as we want, constantly accumulating memory in the state vector.



Solution: We define a new hyperparameter called window\_sz.

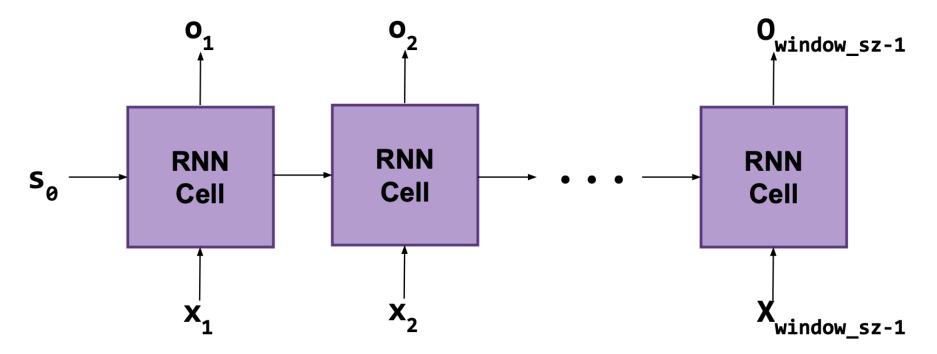
We now chop our corpus into sequences of words of size window\_sz

The new shape of our data should be:

(batch\_sz, window\_sz, embedding\_sz)

Each example in our batch is a "window" of window\_sz many words. Since each word is represented as an embedding\_sz, that is the last dimension of the data.

Now that every example is a window or words, we can run the RNN till the end of that window, and compute the loss for that specific window and update our weights

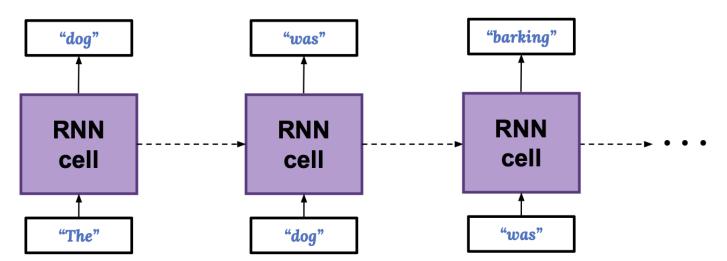


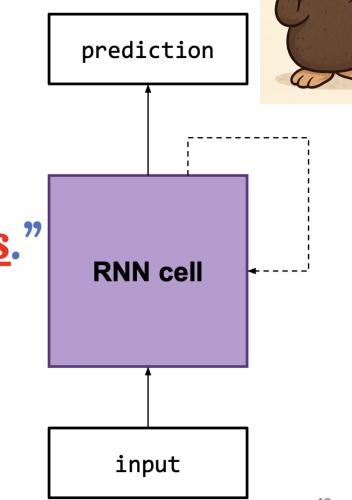
#### Any questions?

# Does RNN fix the limitations of the N-gram ??? model?

- Number of of weights not dependent on N
- State gives flexibility to choose context from near or far

# "The dog was barking at one of the cats."





RNNs can be built from scratch using Python for loops:

```
prev state = Zero vector
for i from 0 to window_sz:
  state and input = concat(inputs[i], prev state)
  current state = fc state(state and input)
  outputs[i] = fc output(current state)
  prev state = current state
return outputs
```

RNNs can be built from scratch using Python for loops.

There's also a handy built-in Keras recurrent layer:

tf.keras.layers.SimpleRNN(units, activation, return\_sequences)

RNNs can be built from scratch using Python for loops.

There's also a handy built-in Keras recurrent layer:

tf.keras.layers.SimpleRNN(units, activation, return\_sequences)

The size of our output vectors

RNNs can be built from scratch using Python for loops.

There's also a handy built-in Keras recurrent layer:

tf.keras.layers.SimpleRNN(units, activation, return\_sequences)

The activation function to be used in the FC layers inside of the RNN Cell

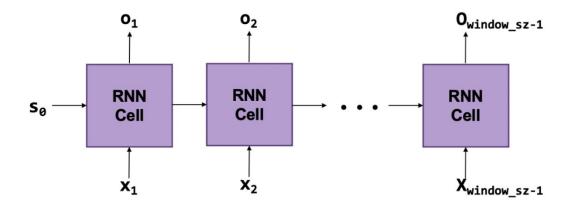
Any intuition why we would want return\_sequences to be TRUE?

#### RNNs in Tensorflow

RNNs can be built from scratch using Python for loops.

There's also a handy built-in Keras recurrent layer:

tf.keras.layers.SimpleRNN(units, activation, return\_sequences)



- If True: calling the RNN on an input sequence returns the whole sequence of outputs + final state output
- If False: calling the RNN on an input sequence returns just the final state output (Default)

RNNs can be built from scratch using Python for loops.

There's also a handy built-in Keras recurrent layer:

tf.keras.layers.SimpleRNN(units, activation, return\_sequences)

```
Usage:
```

```
RNN = SimpleRNN(10) # RNN with 10-dimensional output vectors
Final_output = RNN(inputs) # inputs: a [batch_sz, seq_length, embedding_sz] tensor
```