

CSCI 1470/2470
Spring 2024

Ritambhara Singh

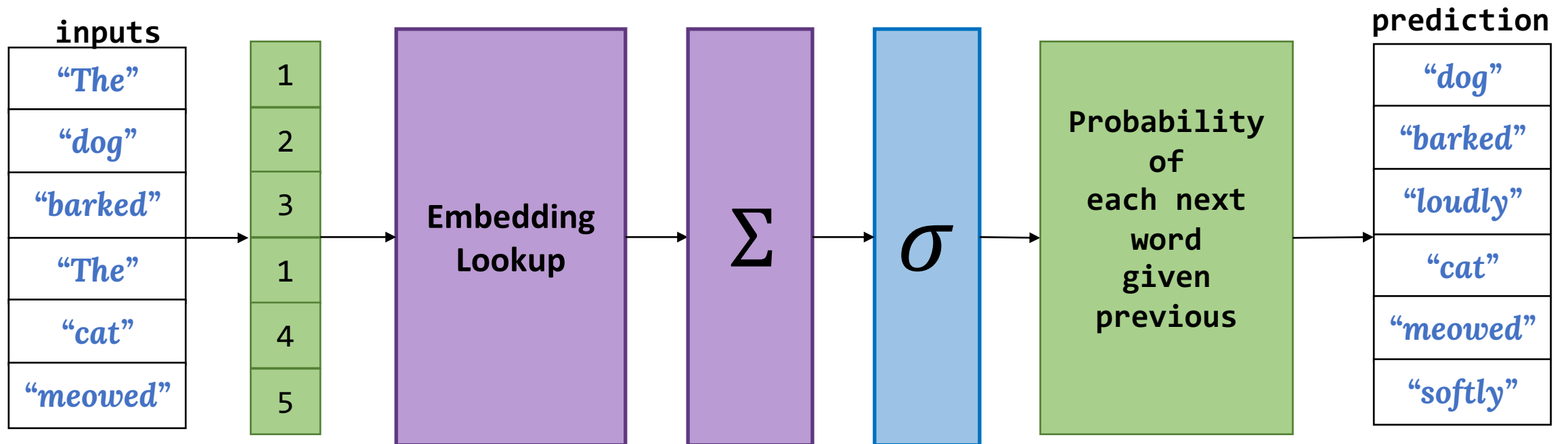
March 01, 2024
Friday

Recurrent Neural Networks

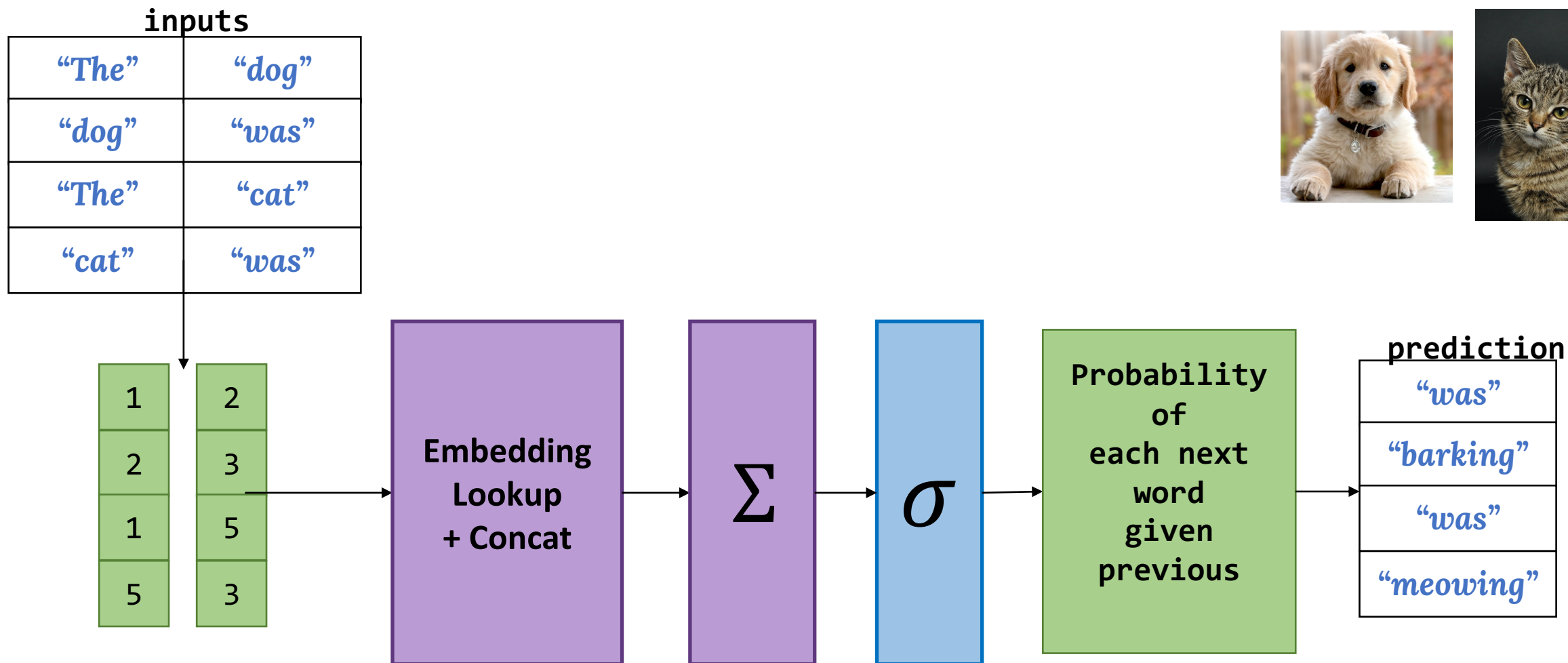
Deep Learning

Deep Learning Days – May 6 and 7, 2024
Project teams announced!
Please complete labs!

Review: Bigram Language Model Architecture



Review: Complete Trigram Language Model

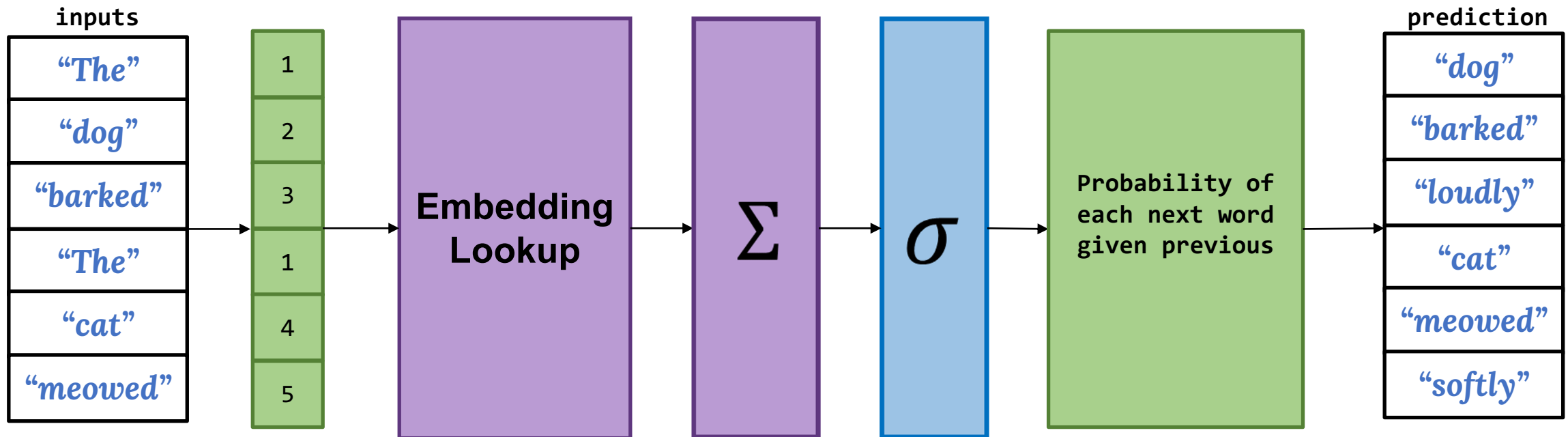


Limitations of the N-gram model

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

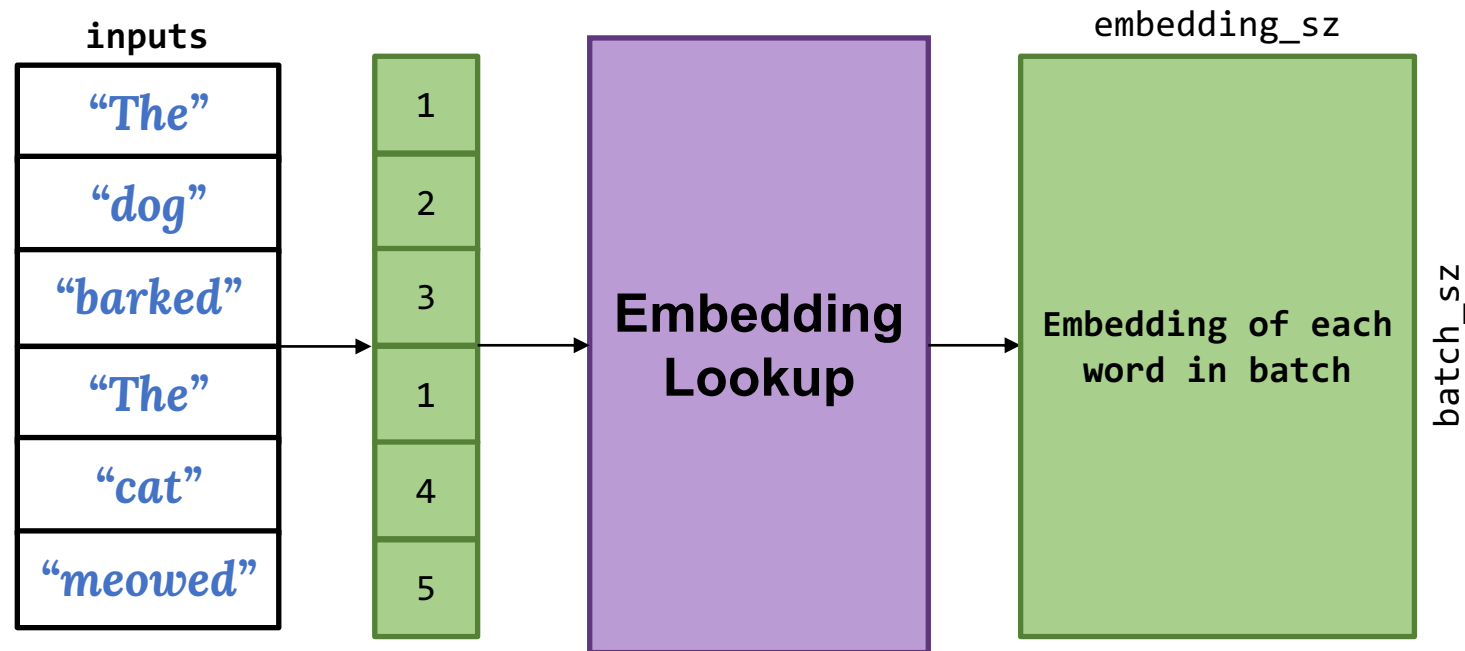
Size of Feed Forward bigram Model

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



Size of Feed Forward bigram Model

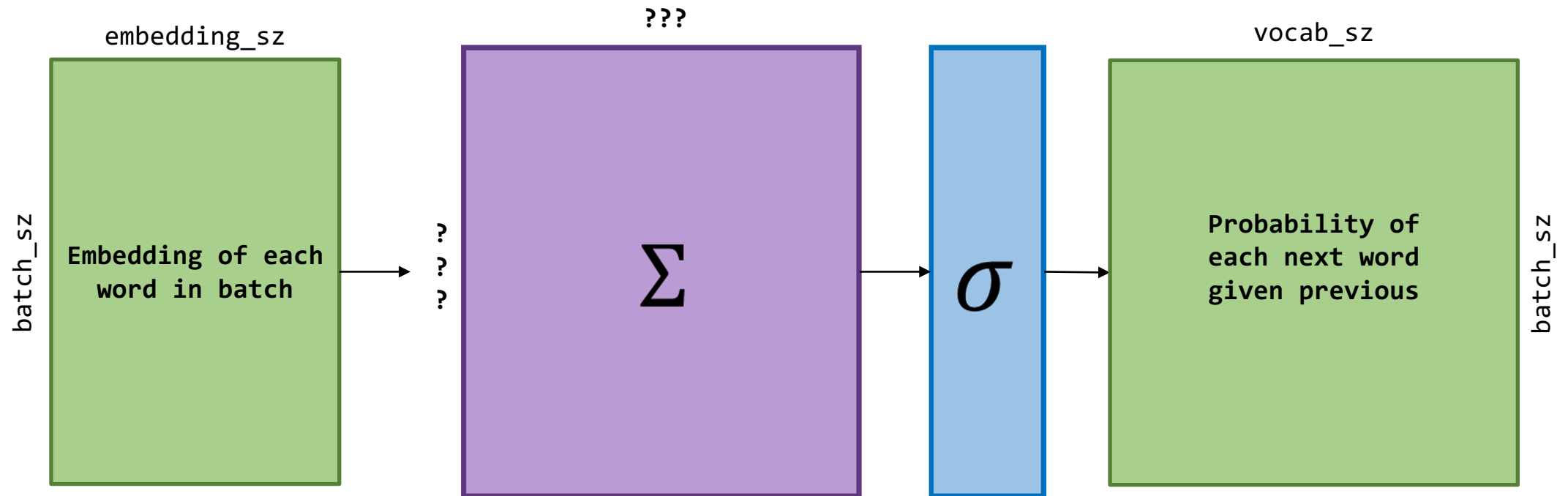
To perform embedding lookup on our entire batch, we just need one embedding matrix of size: (vocab_sz, embedding_sz)



Size of Feed Forward bigram Model

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

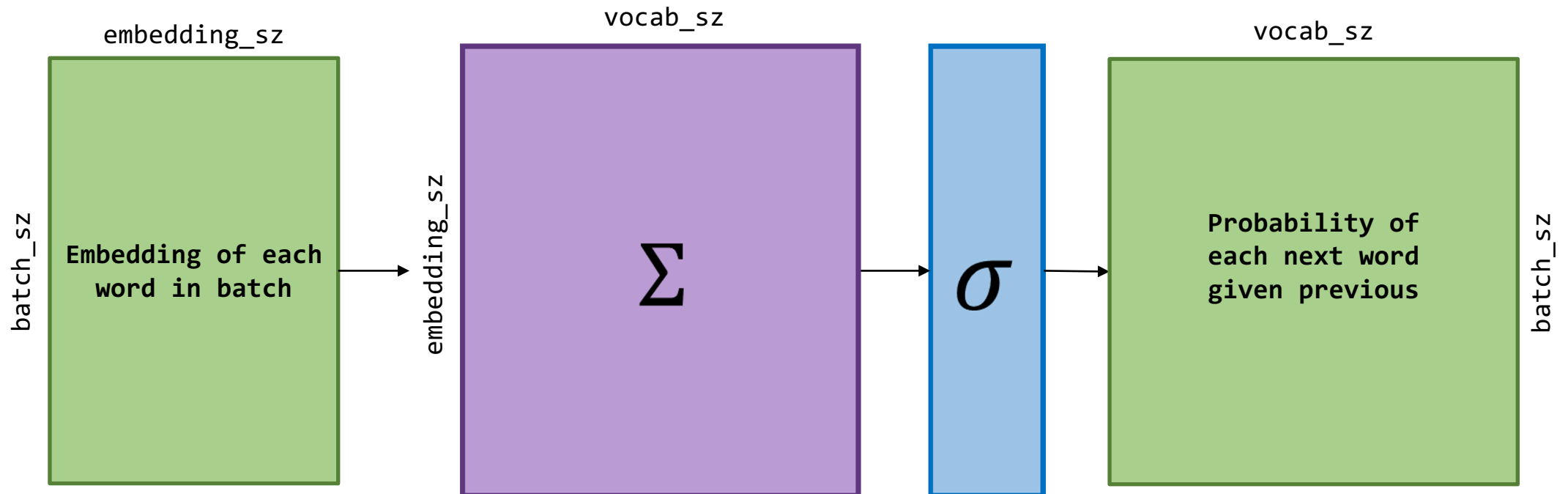
$$(\text{batch_sz}, \text{embedding_sz}) \times (???, ???) \rightarrow (\text{batch_sz}, \text{vocab_sz})$$



Size of Feed Forward bigram Model

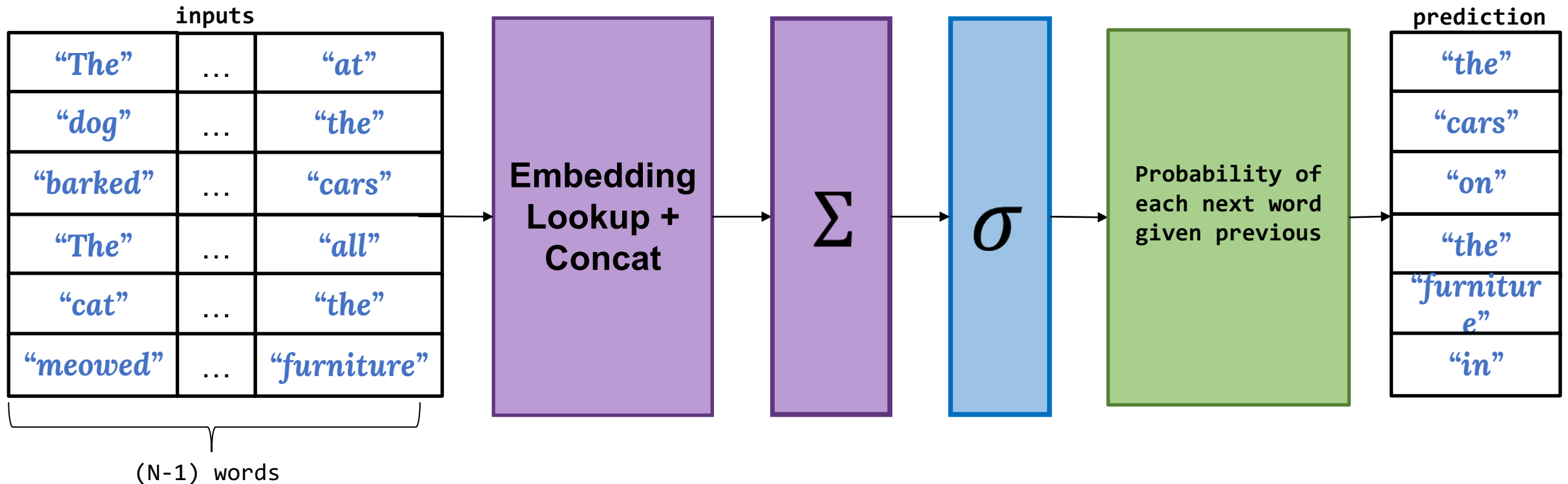
What size do we need the linear layer to be in order to map:

$$(\text{batch_sz}, \text{embedding_sz}) \times (???, ???) \rightarrow (\text{batch_sz}, \text{vocab_sz})$$



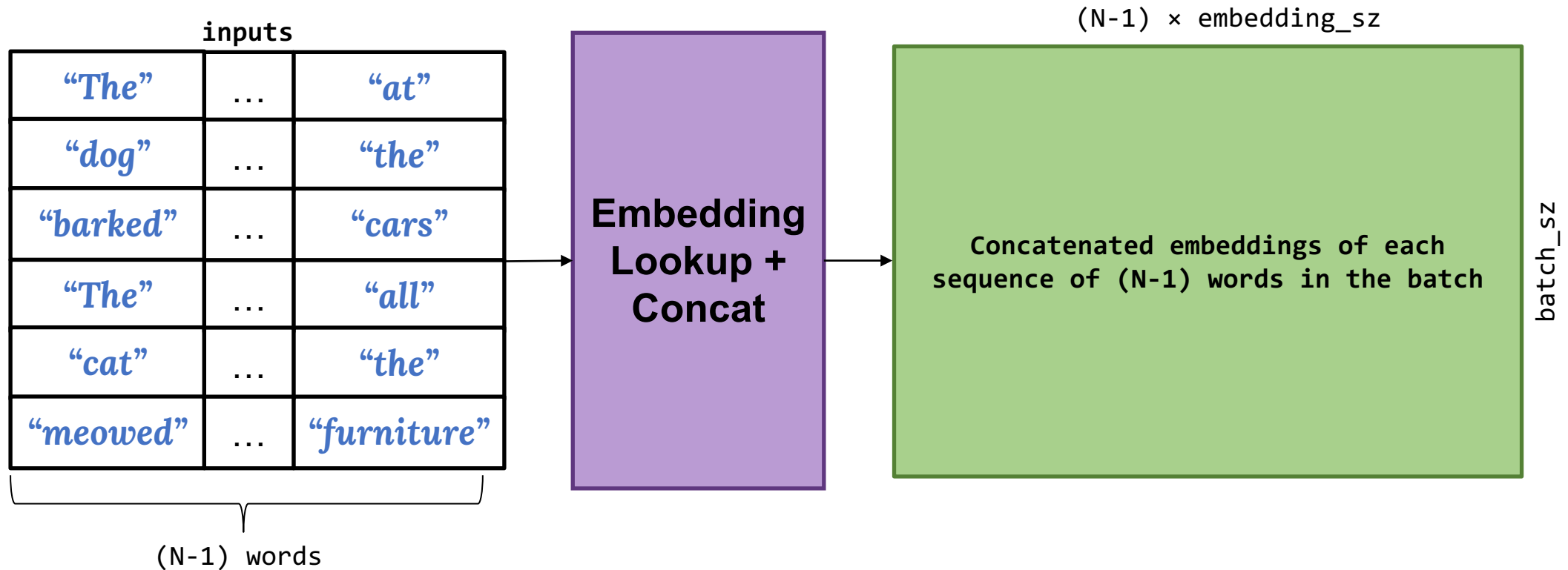
Size of Feed Forward N-gram Model

So what happens in the N-gram case?



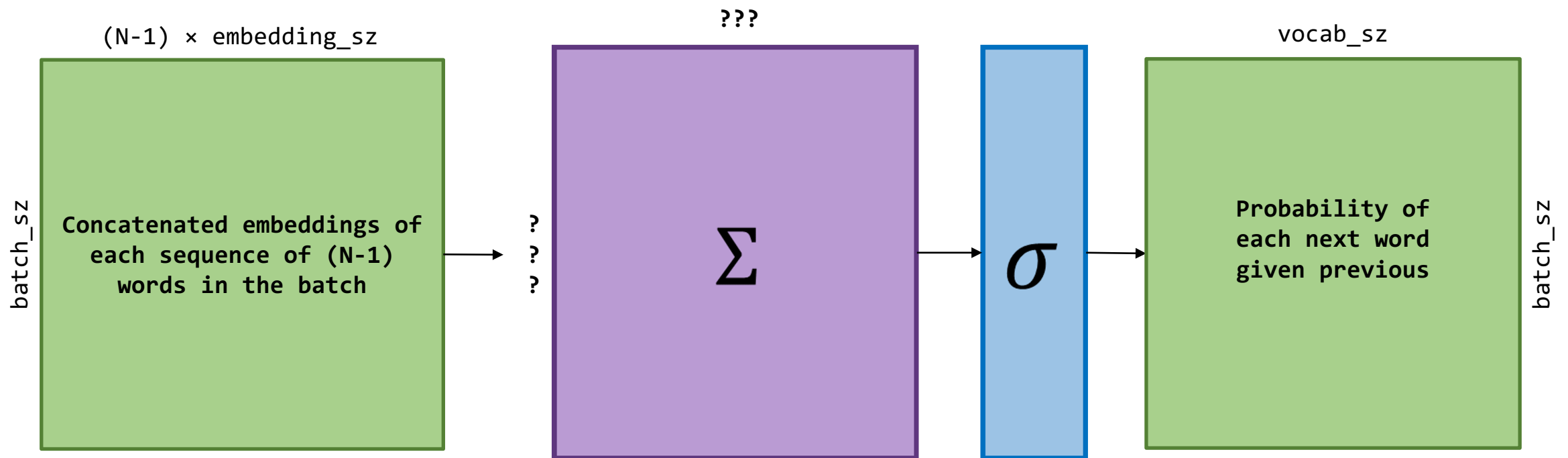
Size of Feed Forward N-gram Model

Embedding lookup + Concatenation still requires only one embedding matrix of size: $(\text{vocab_sz}, \text{embedding_sz})$



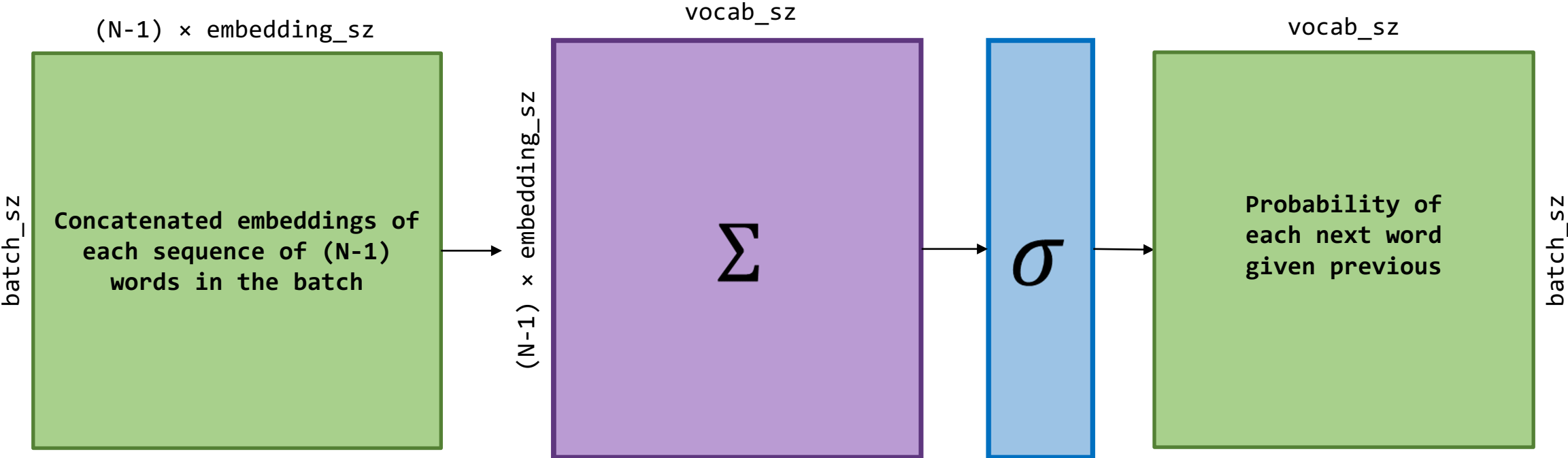
Size of Feed Forward N-gram Model

But what happens to our feed forward layer?



Size of Feed Forward N-gram Model

It needs to be size: $((N-1) \times \text{embedding_sz}, \text{vocab_sz})$
 For every word, we add $(\text{embedding_sz} \times \text{vocab_sz})$ more weights!



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.

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.

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.”



(“The”, “dog”) → “barked”



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.”

(“at”, “one”, “of”, “the”) → ???



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

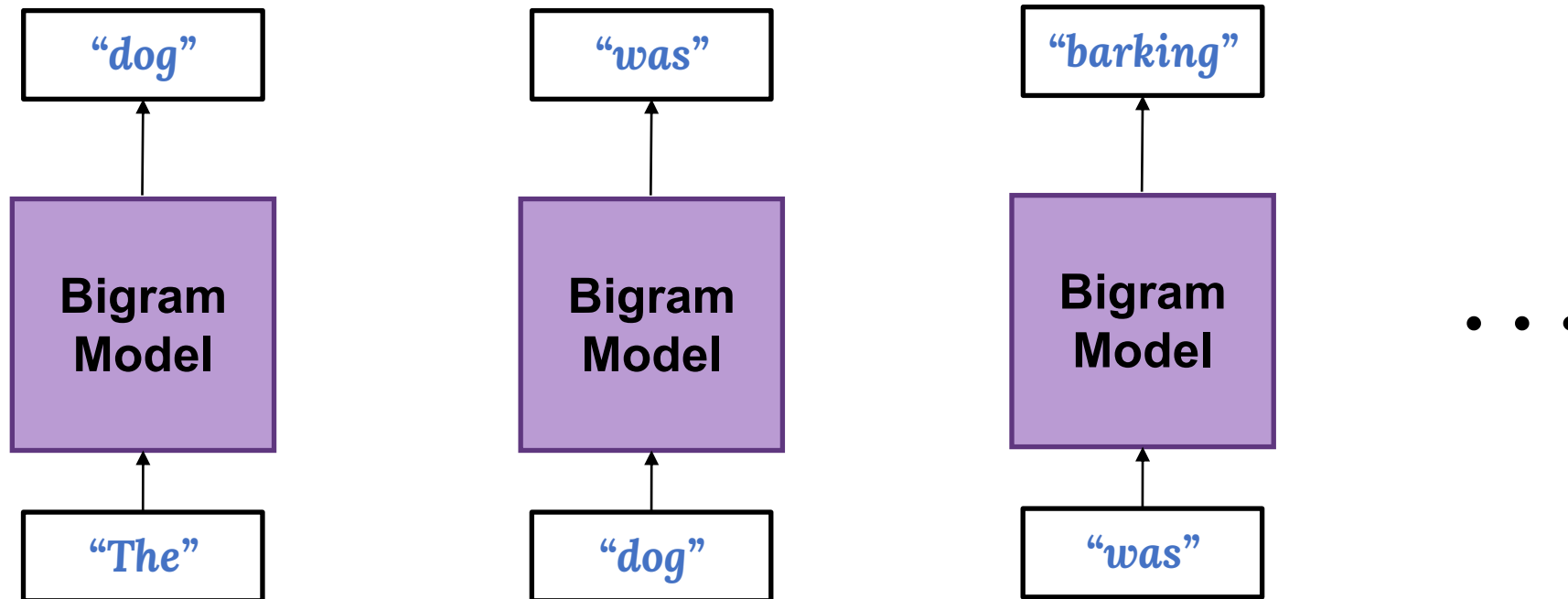
What problems do we run into using Feed Forward N-gram models?

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

New Approach

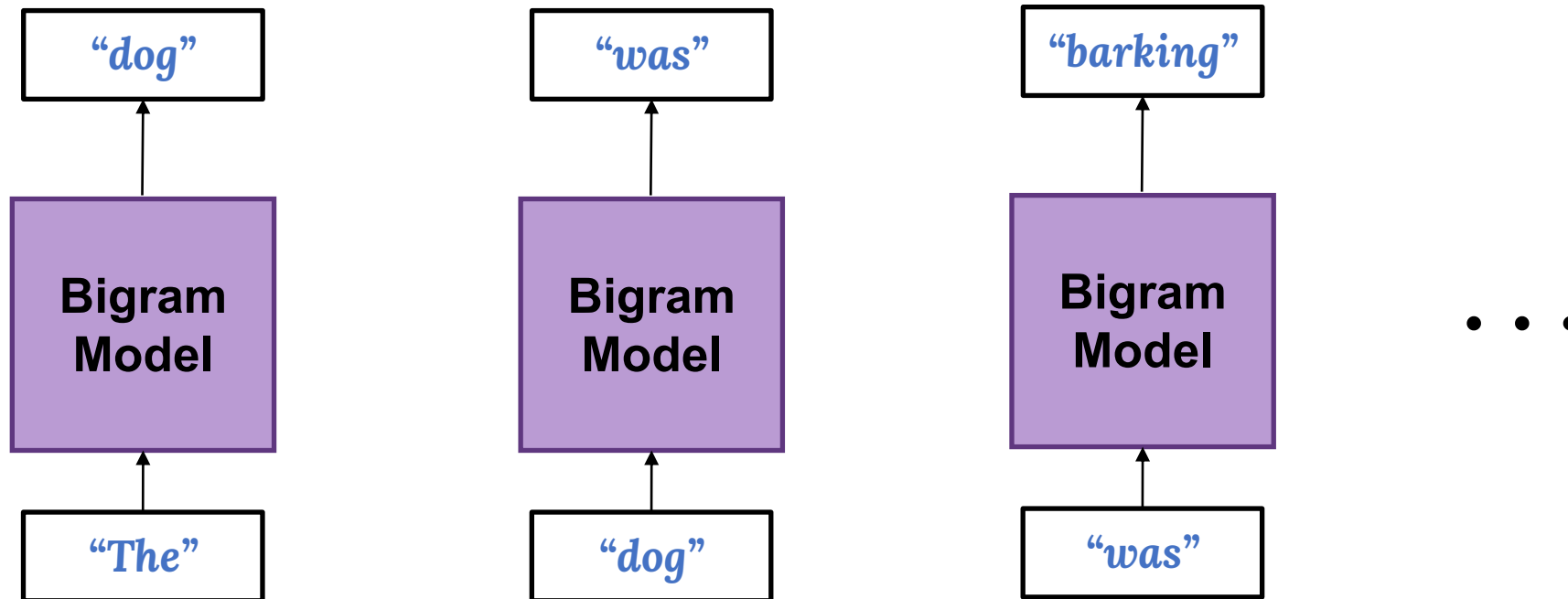
Let's revisit the bigram model and see several iterations of prediction using a bigram model:



New Approach

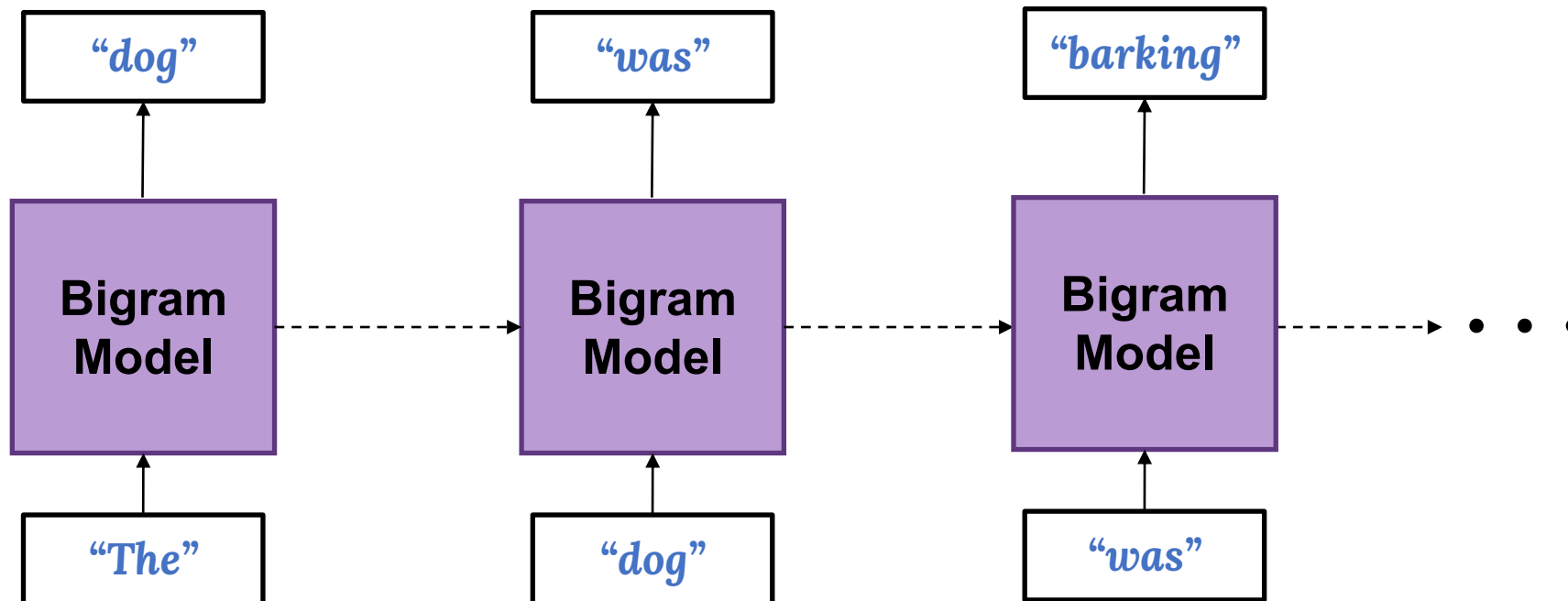
Any ideas?

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



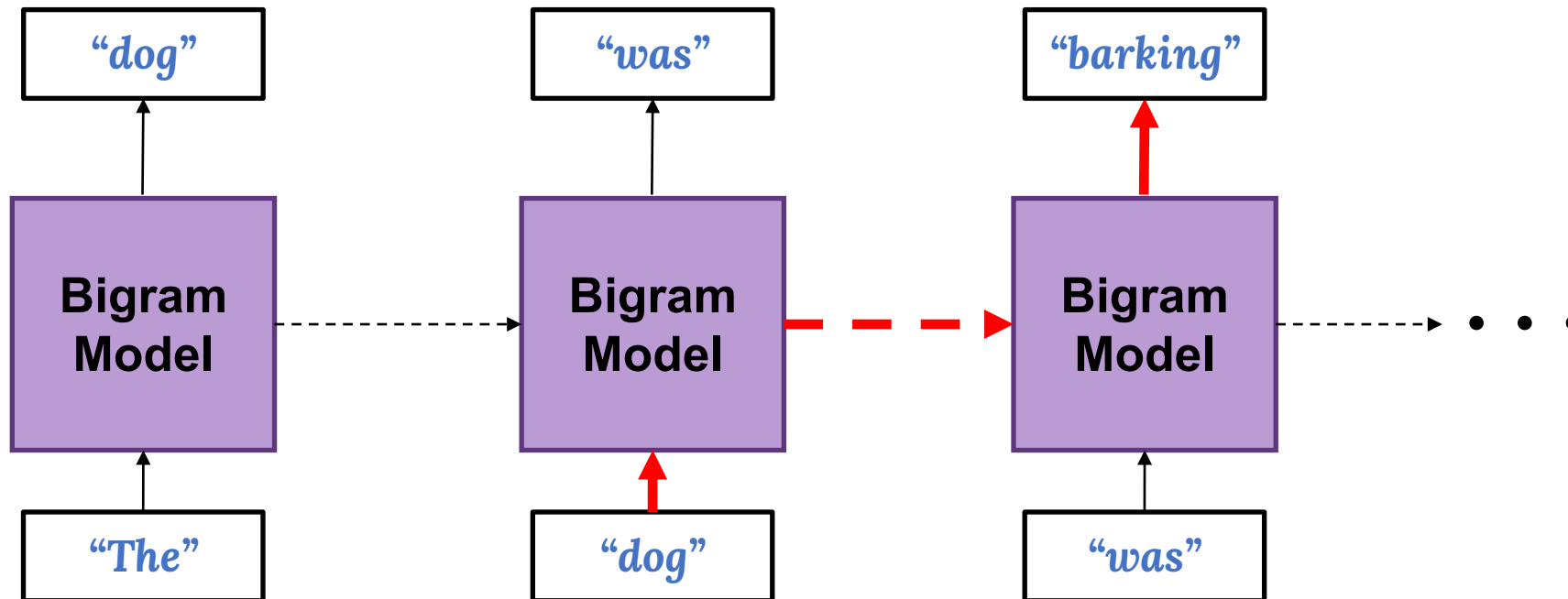
New Approach

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



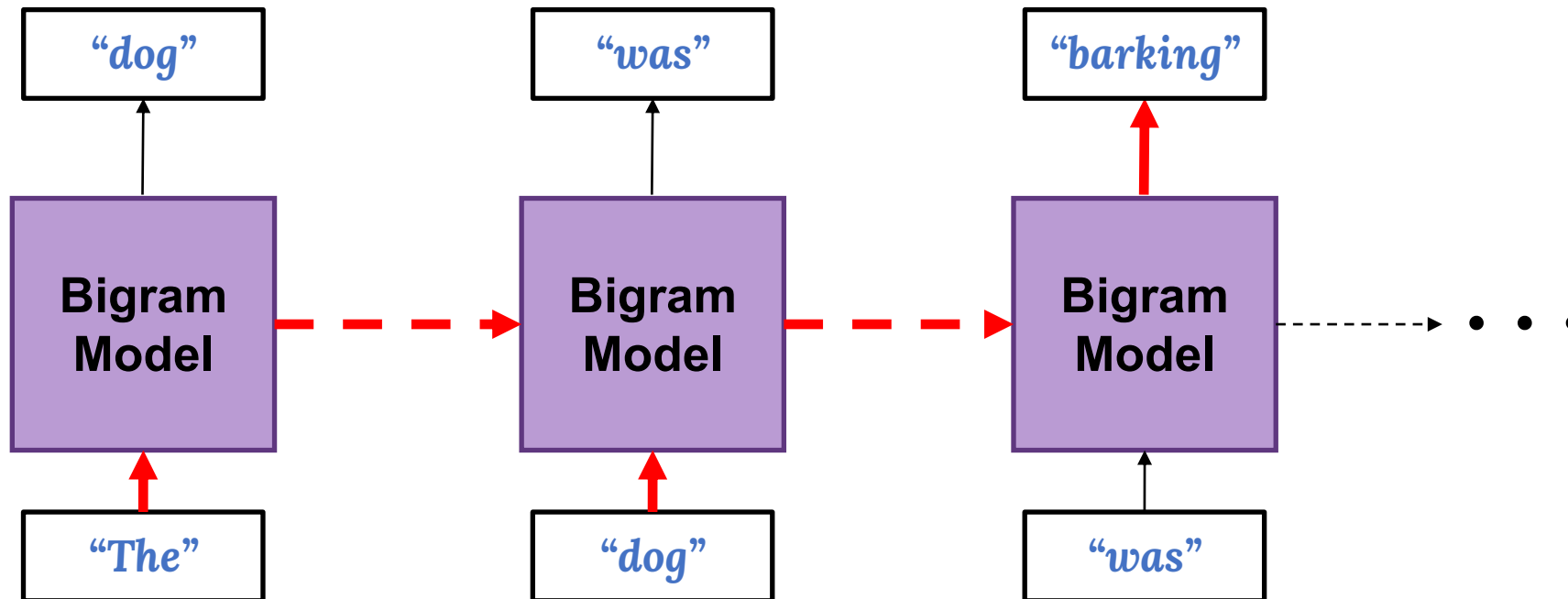
New Approach

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



New Approach

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

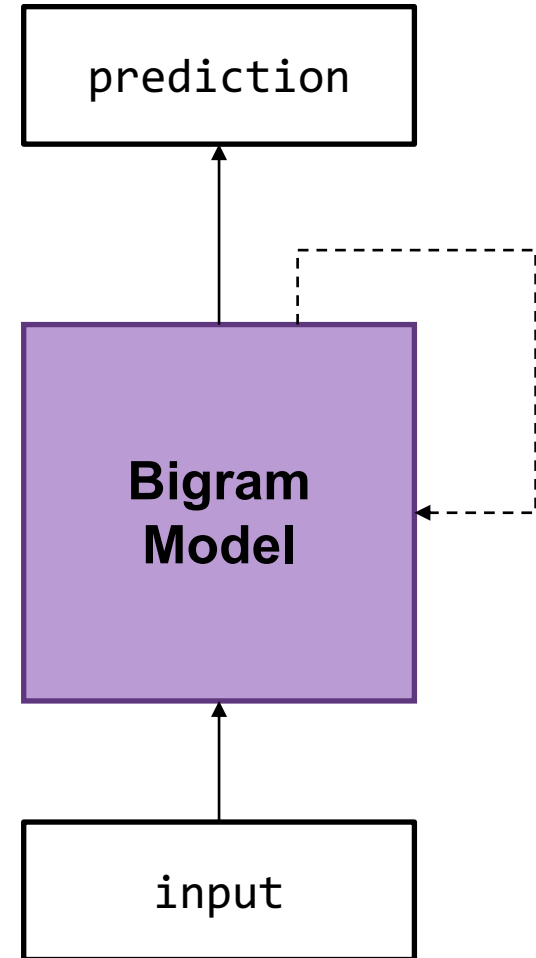


New Approach

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

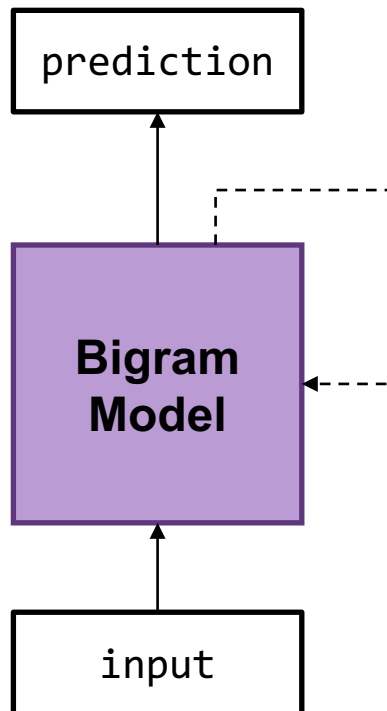
We call this connection a *recurrent* connection.

We call the previous representation the “unrolled” representation.

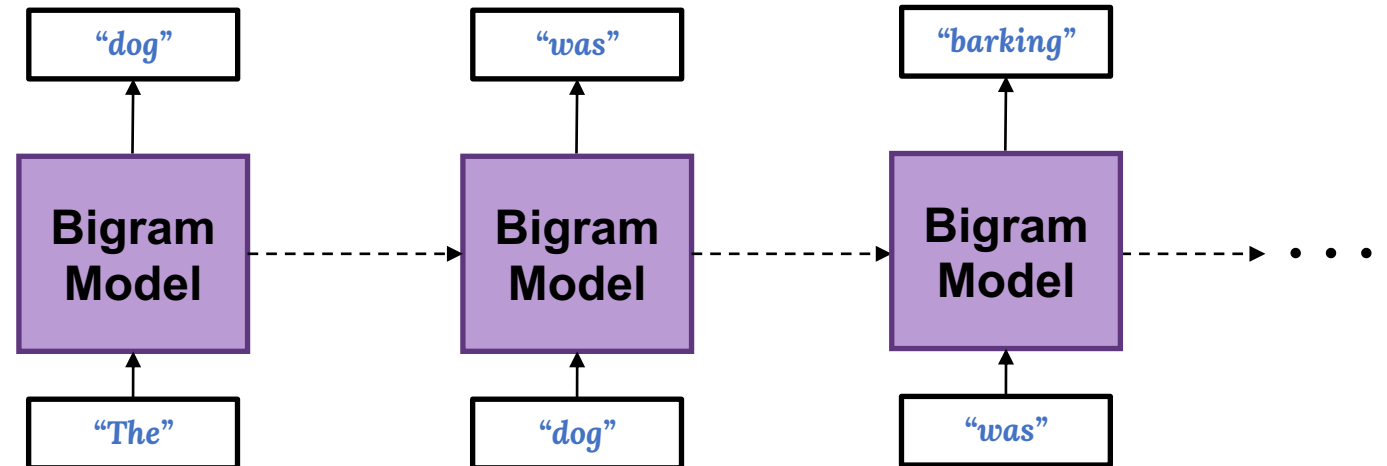


Different views of recurrent models

Recurrent view



Unrolled view



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.

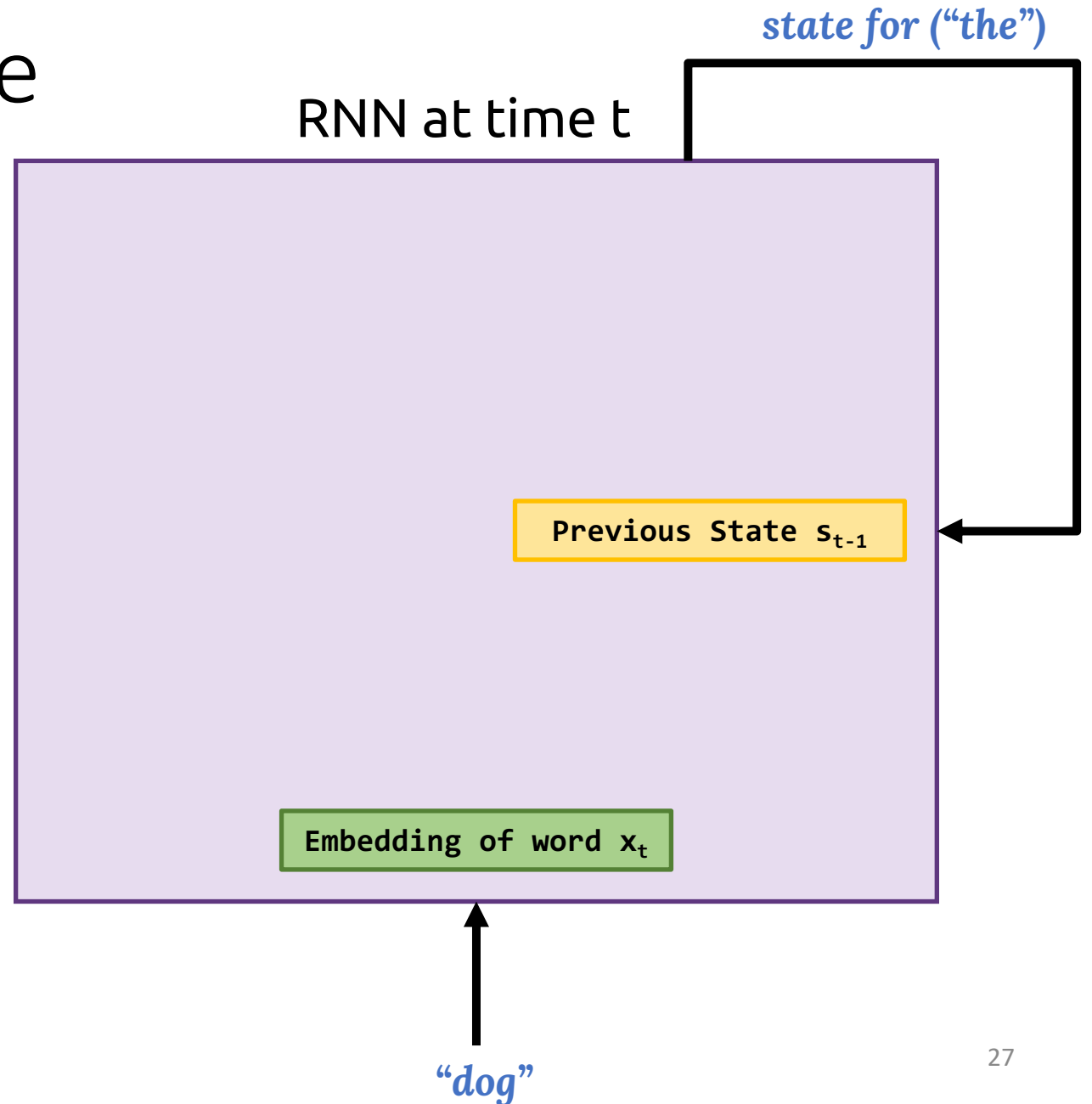
They can be used to process sequence data with relatively low model complexity when compared to feed forward models.

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!

RNN Cell Architecture

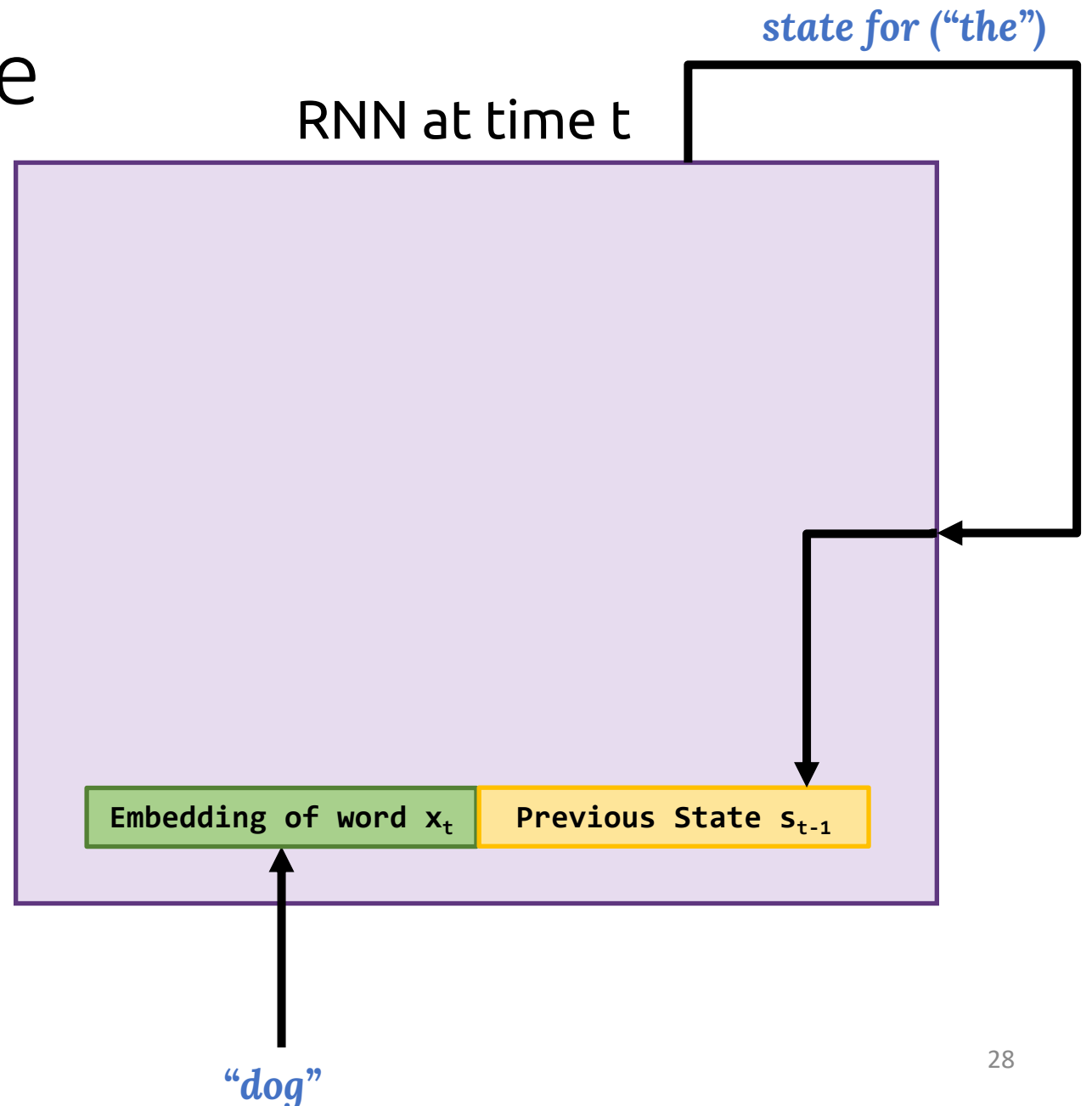
At each step of our RNN, we will get an input word, and a state vector from the previous cell.



RNN Cell Architecture

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.

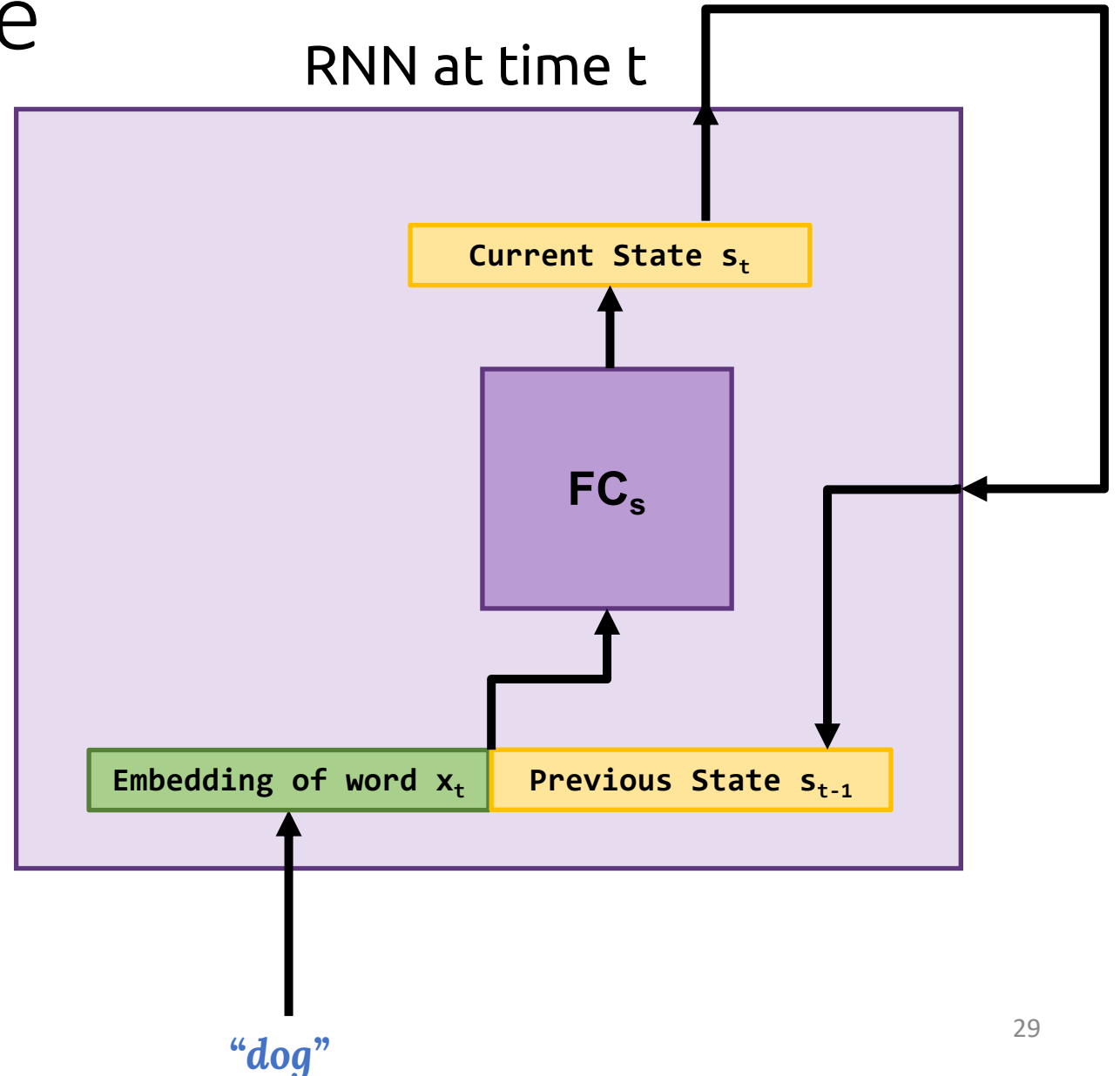


RNN Cell Architecture

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



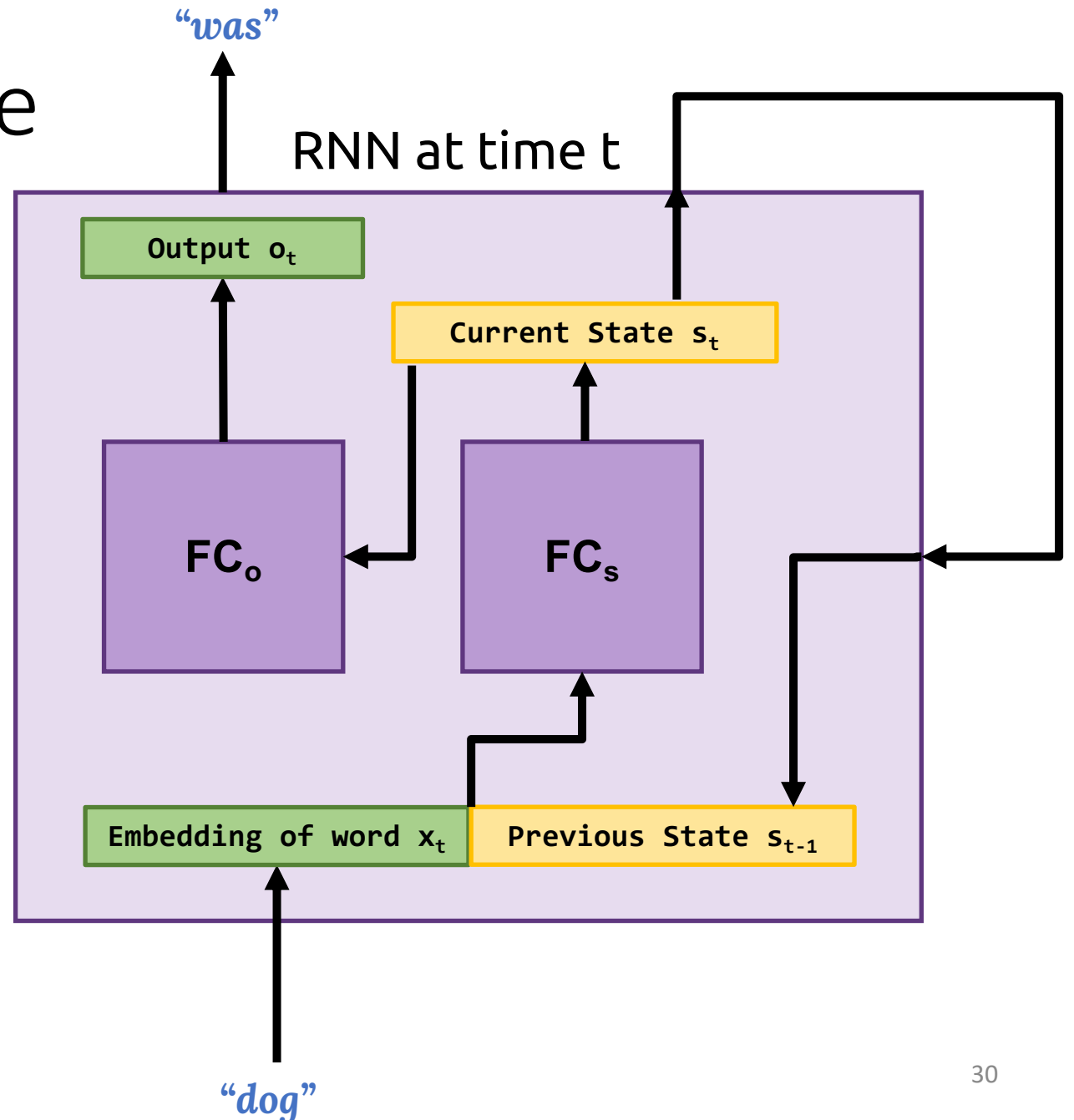
RNN Cell Architecture

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.

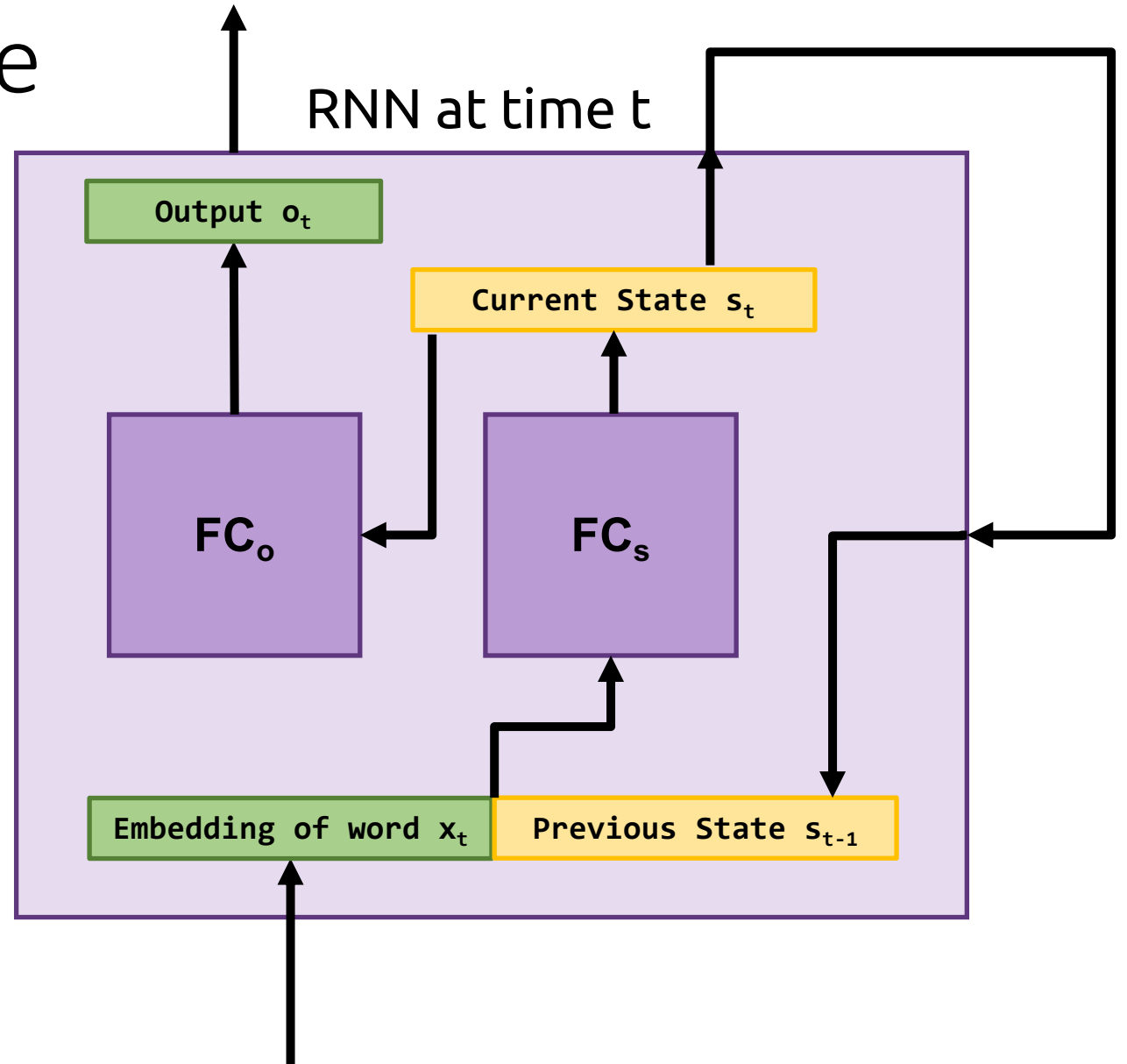


RNN Cell Architecture

We can represent the RNN in with the following equations:

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

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



RNN Cell Architecture

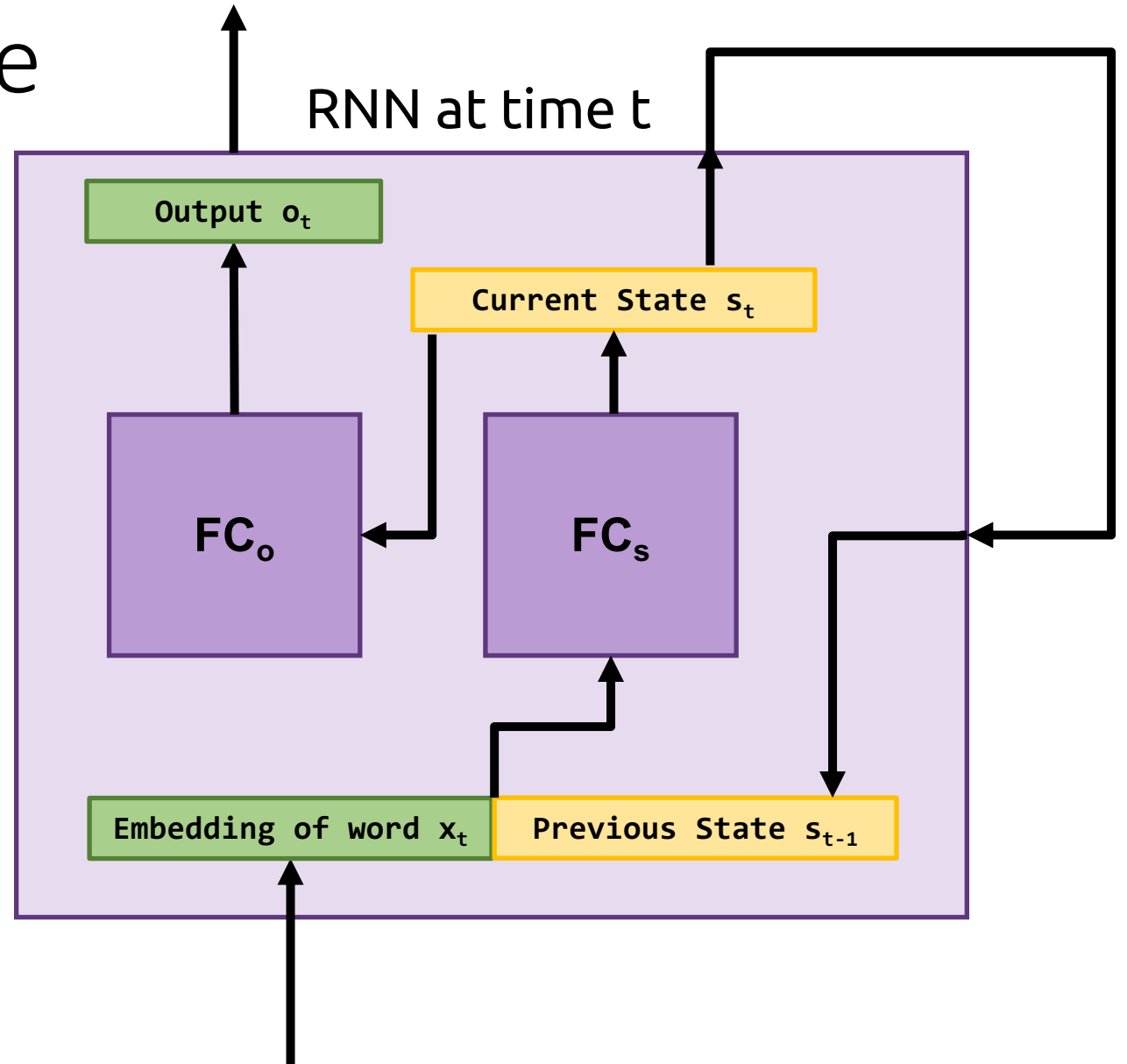
We can represent the RNN in with the following equations:

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

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

**Nonlinear activations
(e.g. sigmoid, tanh)**

Any questions?



RNN Cell Architecture

We can represent the RNN in with the following equations:

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

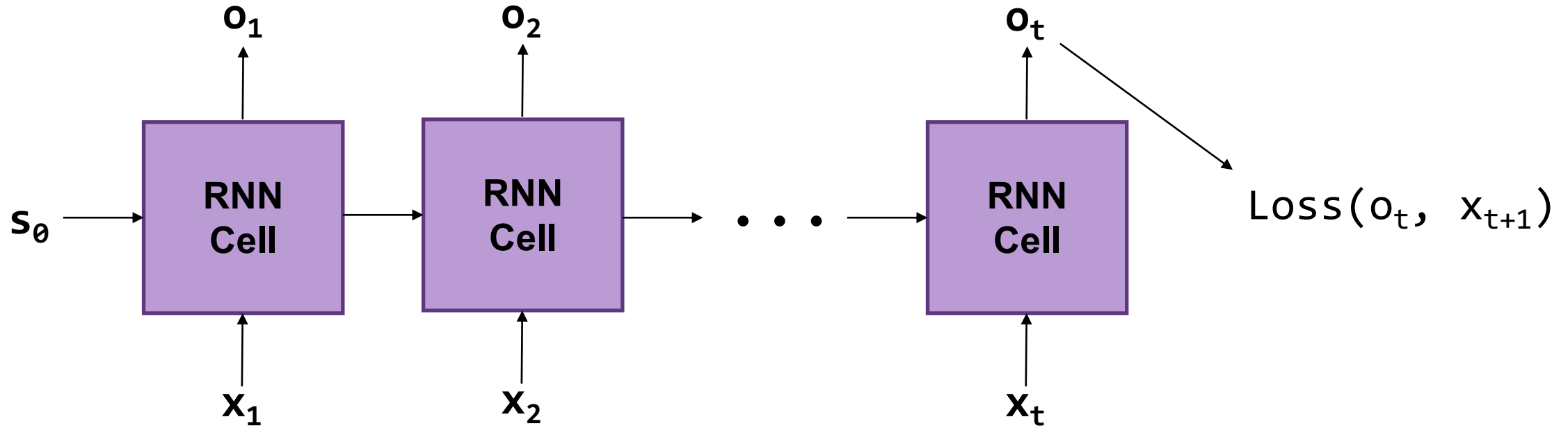
$$o_t = \sigma(s_t W_o + b_o)$$

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”)

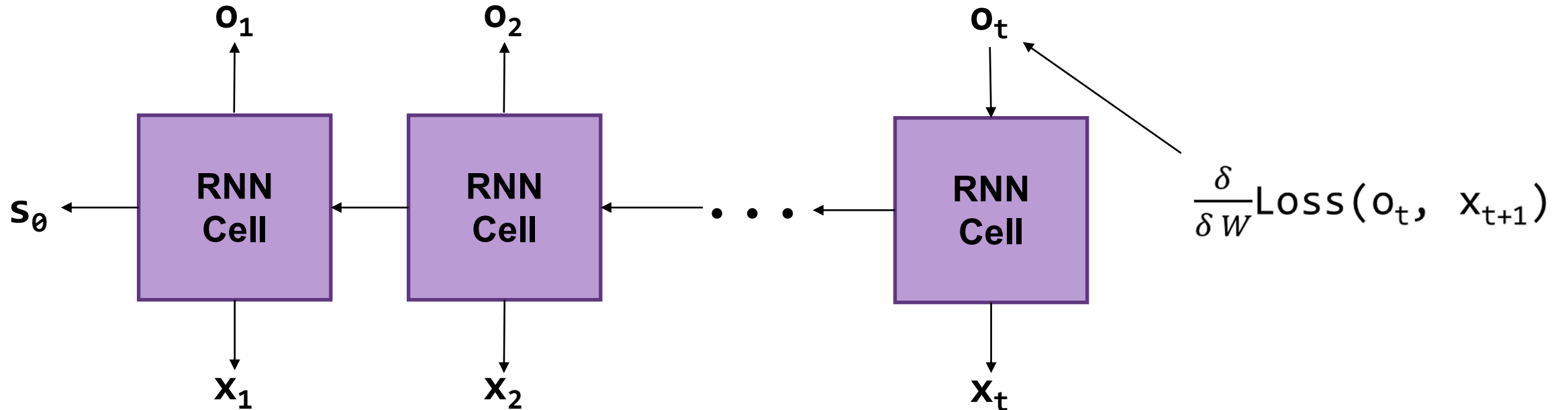
Training RNNs

We can calculate the cross entropy loss just as before since for any sequence of input words (x_1, x_2, \dots, x_t), we know the true next word x_{t+1}



Training RNNs

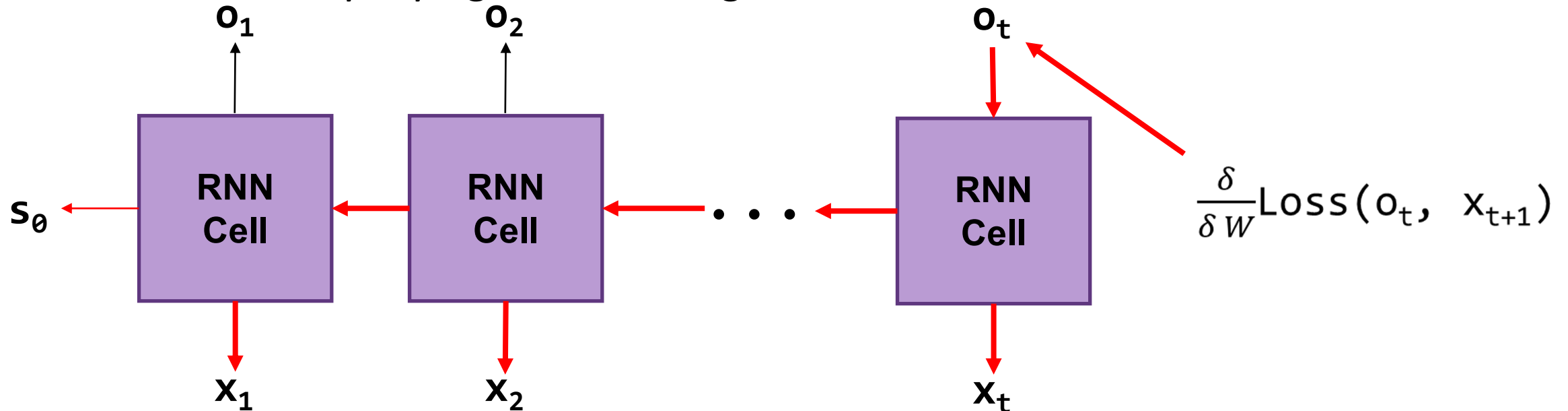
But what happens when we differentiate the loss and backpropagate?



Training RNNs

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

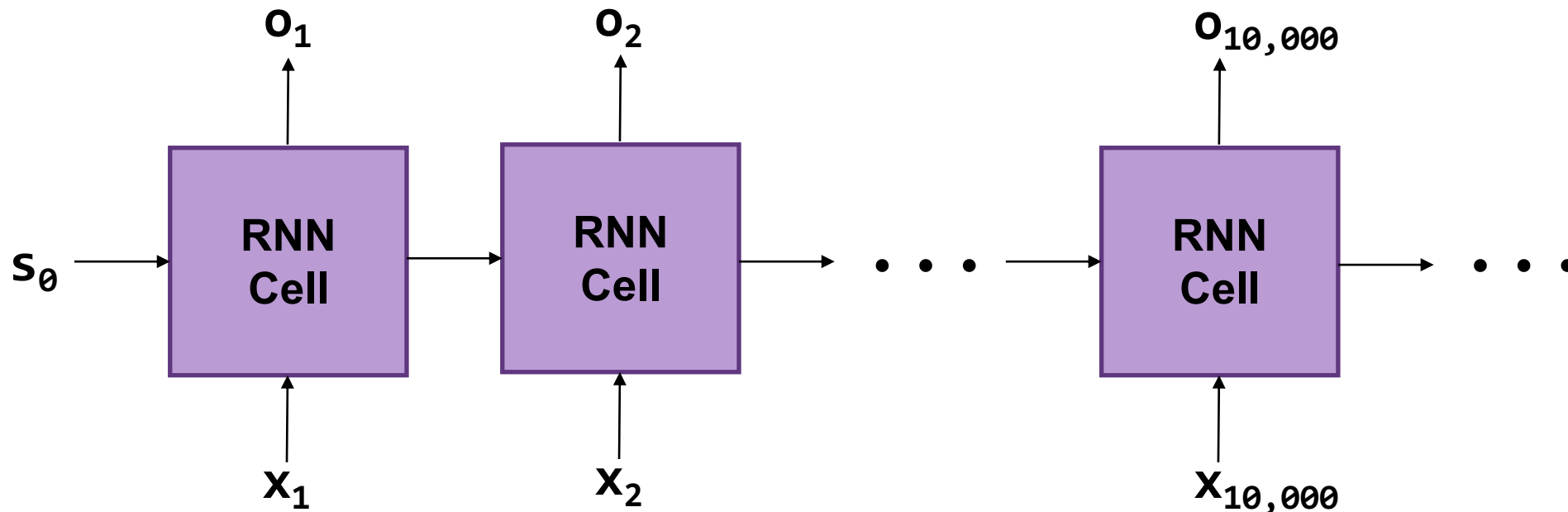
We call this *backpropagation through time*.



Training RNNs

But at what point do we stop and calculate the loss/update?

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



Training RNNs

Solution: We define a new hyperparameter called `window_sz`.

We now chop our corpus into sequences of words of size `window_sz`

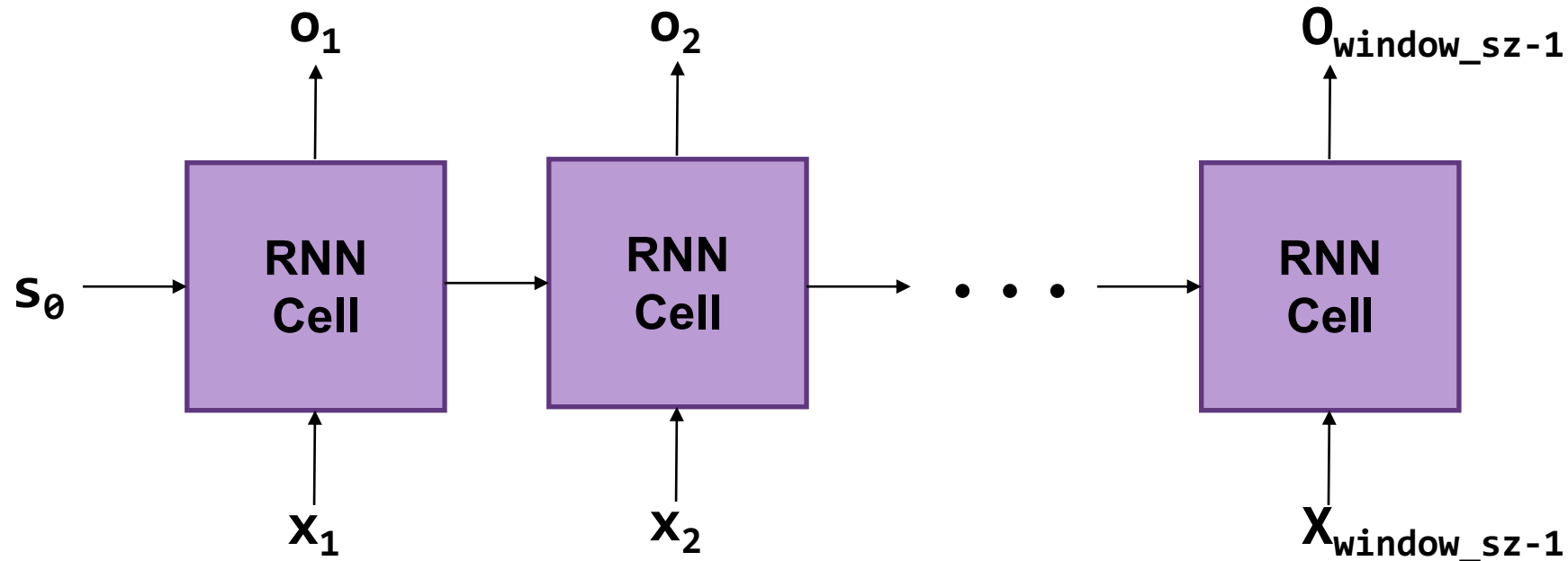
The new shape of our data should be:

`(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.

Training RNNs

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

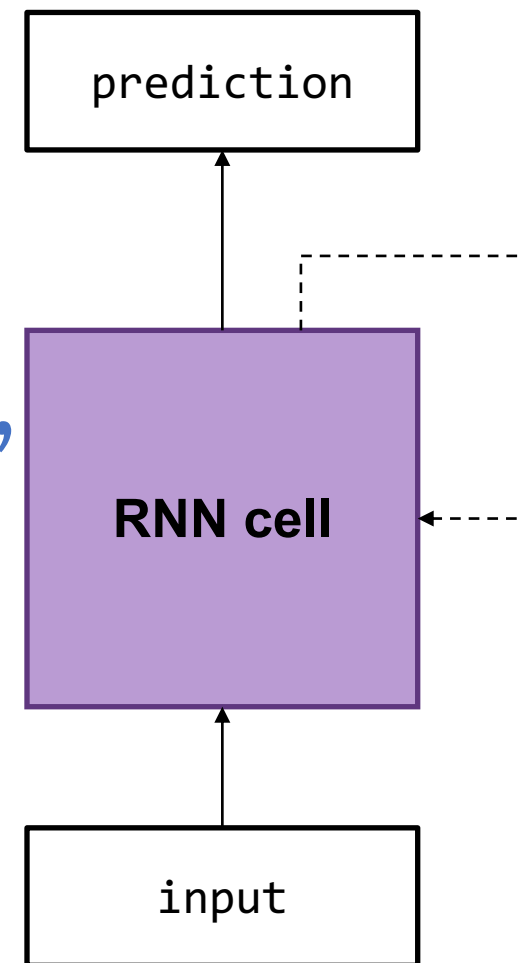
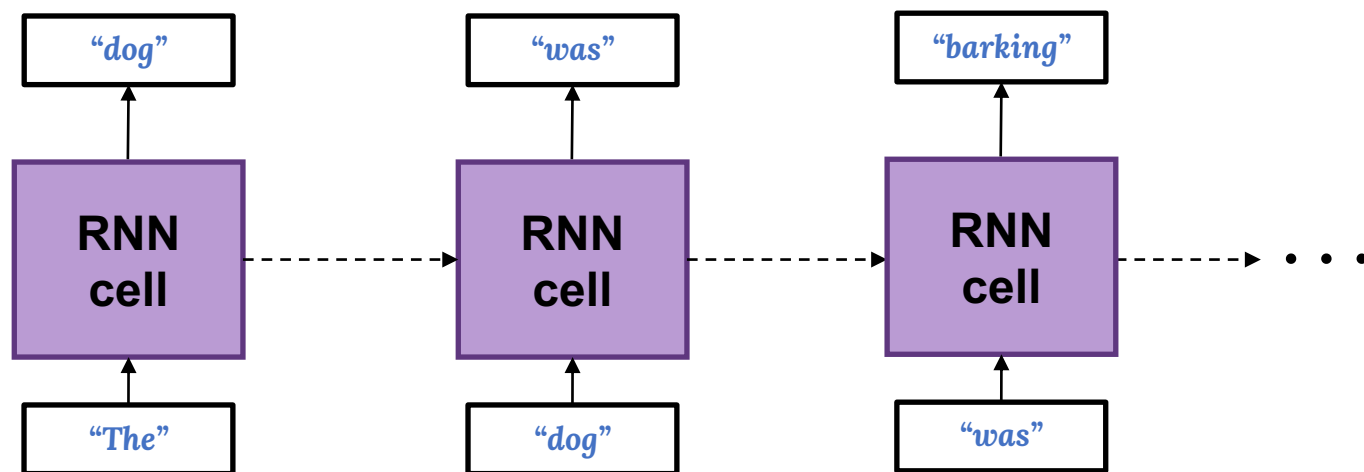




Does RNN fix the limitations of the N-gram model?

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

“The dog was barking at one of the cats.”



RNNs in Tensorflow

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 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)
```

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)
```

The size of our output vectors



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)
```

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



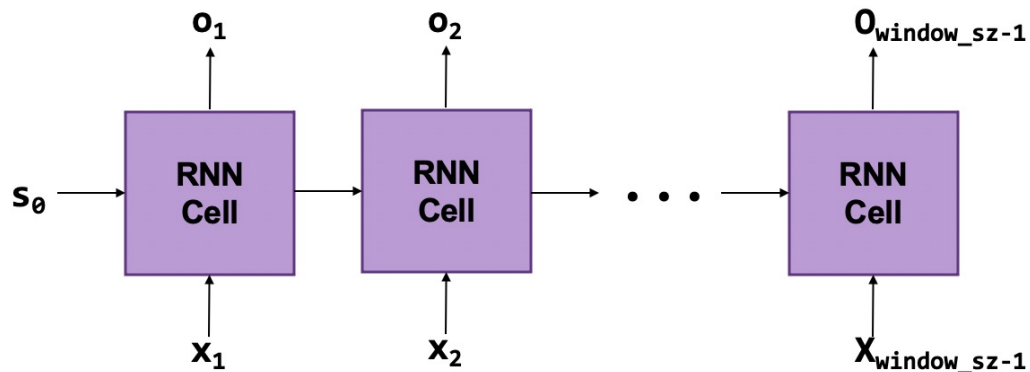
RNNs in Tensorflow

Any intuition why we would want `return_sequences` to be `TRUE`?

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 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)
```

Usage:

```
RNN = SimpleRNN(10) # RNN with 10-dimensional output vectors
```

```
Final_output = RNN(inputs) # inputs: a [batch_sz, seq_length, embedding_sz] tensor
```

Any questions?



RNNs in Tensorflow

```
inputs = np.random.random([32, 10, 8]).astype(np.float32)
simple_rnn = tf.keras.layers.SimpleRNN(4)
```

```
output = simple_rnn(inputs)
```

inputs: a [batch_sz,
seq_length, embedding_sz]
tensor

Join at menti.com | use code **8372 9615**

```
simple_rnn = tf.keras.layers.SimpleRNN(4,  
return_sequences=True)
```

```
whole_sequence_output = simple_rnn(inputs)
```

What is the size of
(a) output
(b) whole_sequence_output?

RNNs are a marked improvement
over previous language models
we've seen

But what are the implications
when language models get really
good?

Like really, really, *really* good

GPT-3

GPT-3, explained: This new language AI is uncanny, funny — and a big deal

Computers are getting closer to passing the Turing Test.

By Kelsey Piper | Aug 13, 2020, 9:50am EDT

[f](#) [🐦](#) [📄](#) SHARE



OpenAI co-founder and chair Greg Brockman, OpenAI co-founder and CEO Sam Altman, and TechCrunch news editor Frederic Lardinois during TechCrunch Disrupt San Francisco 2019. | Steve Jennings/Getty Images for TechCrunch

GPT-3

Background:

- Attention- and transformer-based model (We'll be discussing the details of these models next week in class, so stay tuned!)
- First Described by OpenAI in May 2020

ChatGPT

Background:

- Attention- and transformer-based model + **Reinforcement Learning**
(We'll cover this at the end of the course!)
- Launched in November 2022

These models are dangerous. Why?

Example1: Impacts of fake news and bot networks

- What if fake news didn't even have to be written by humans?
- **Could this significantly undermine trust in written sources?**



You

Write a news piece about Pepper the penguin who was granted knighthood



ChatGPT

Title: Pepper the Penguin Knighted for Exceptional Service to Conservation

Date: February 28, 2024

By: [Your Name], Staff Writer

In a heartwarming ceremony that captivated the world, Pepper the Penguin has been granted knighthood in recognition of his extraordinary dedication to conservation efforts and raising awareness about the plight of his species. The conferral of this prestigious honor marks a historic moment in the realm of wildlife preservation.

Pepper, a beloved resident of Antarctica, has long been a symbol of resilience and determination in the face of environmental challenges. His tireless advocacy for penguin conservation and his charismatic presence have endeared him to people around the globe, inspiring countless individuals to take action to protect endangered species and their habitats.

Other Limitations

Example 2: “While we’ve made efforts to make the model refuse inappropriate requests, it will sometimes respond to harmful instructions or exhibit biased behavior.”

The Washington Post
Democracy Dies in Darkness

TECH Help Desk Artificial Intelligence Internet Culture Space Tech Policy

ChatGPT leans liberal, research shows

Chatbots are ingrained with political biases picked up from their training data – which in most cases is unfiltered text from the web

 By [Gerrit De Vynck](#)

Updated August 16, 2023 at 7:21 p.m. EDT | Published August 16, 2023 at 7:01 p.m. EDT

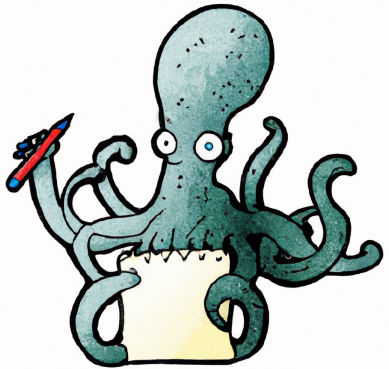
Other Limitations

Example 3: “ChatGPT sometimes writes plausible-sounding but incorrect or nonsensical answers”

OpenAI will not make GPT-3/ChatGPT open source— instead, they have released an API and limits the number of people allowed to try the actual model.

Recap

Limitations of N-gram models



RNNs

Size of weights dependent on N

Limited Flexibility

Recurrent connection can help

RNN cell architecture

Backprop through time

RNNs in Tensorflow

prediction

Bigram Model

input

