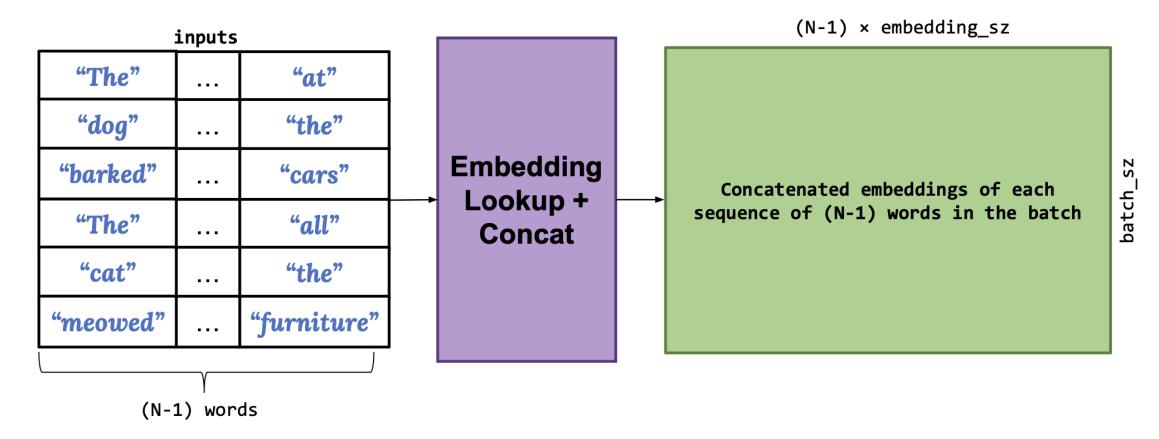


Size of Feed Forward N-gram Model

Embedding lookup + Concatenation still requires only one embedding matrix of size: (vocab_sz, embedding_sz)



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

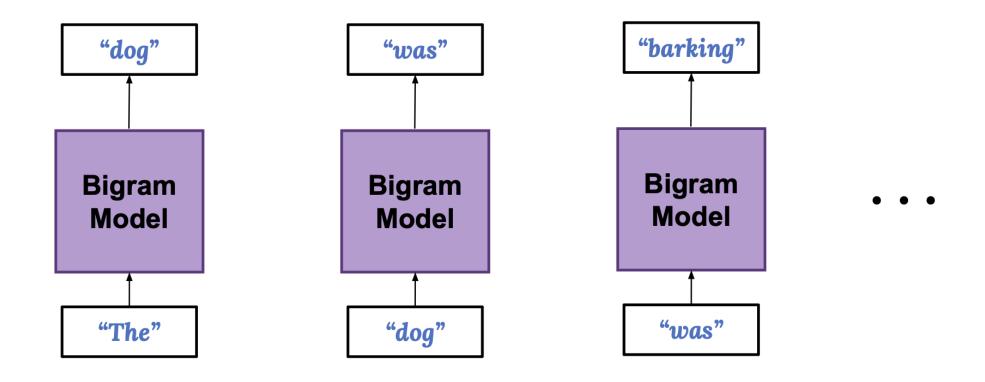
What problems do we run into using Feed Forward N-gram models?

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.

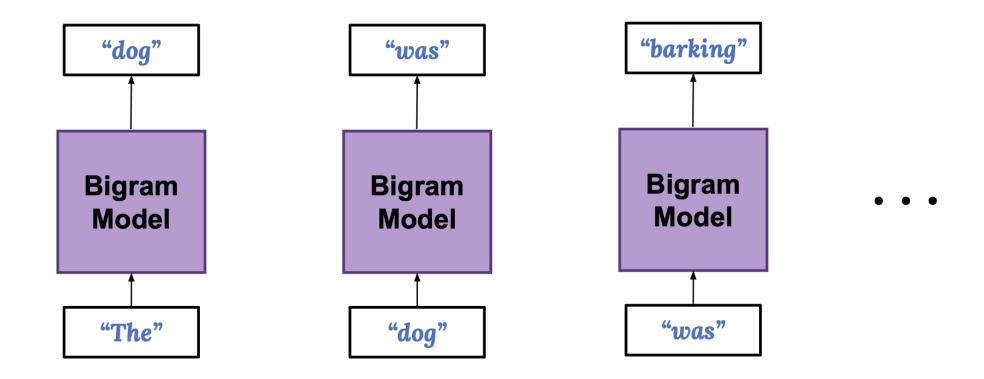
Let's revisit the bigram model and see several iterations of prediction using a bigram model:



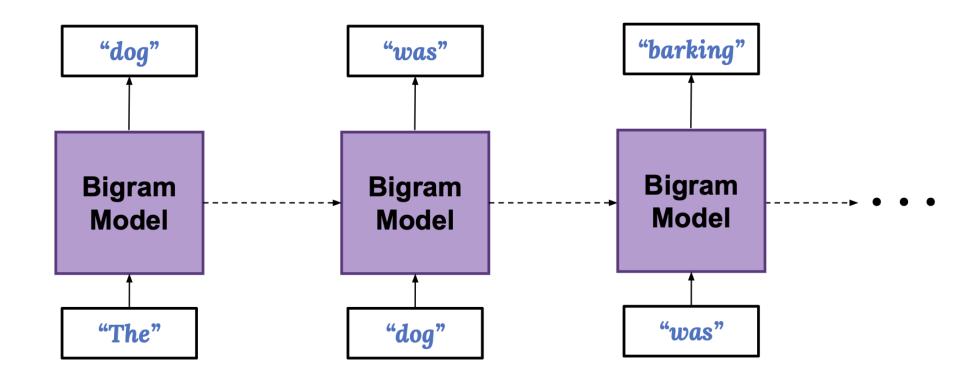
Any ideas?

New Approach

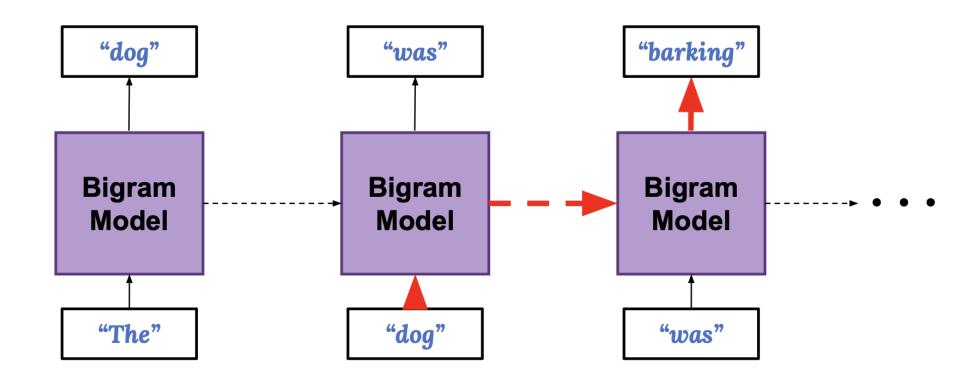
Ideally, we would like to be able to keep "memory" of what words occurred in the past.



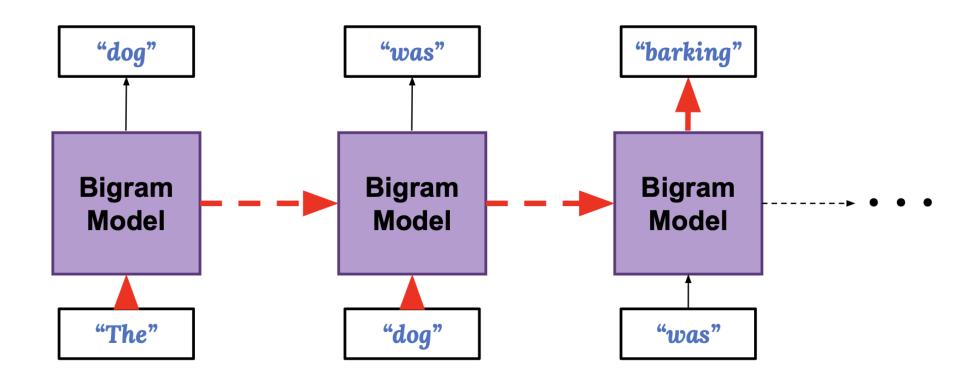
What if we sequentially passed information from our previous bigram block into our next block?



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



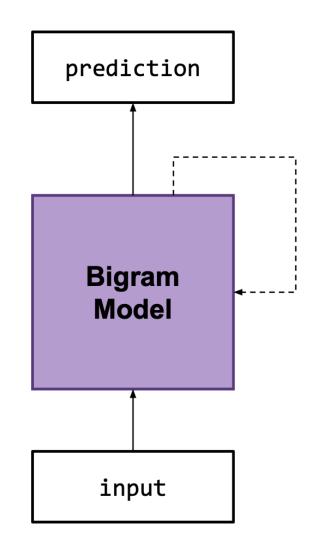
In fact, we even have a way of knowing that "The" was observed!



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.

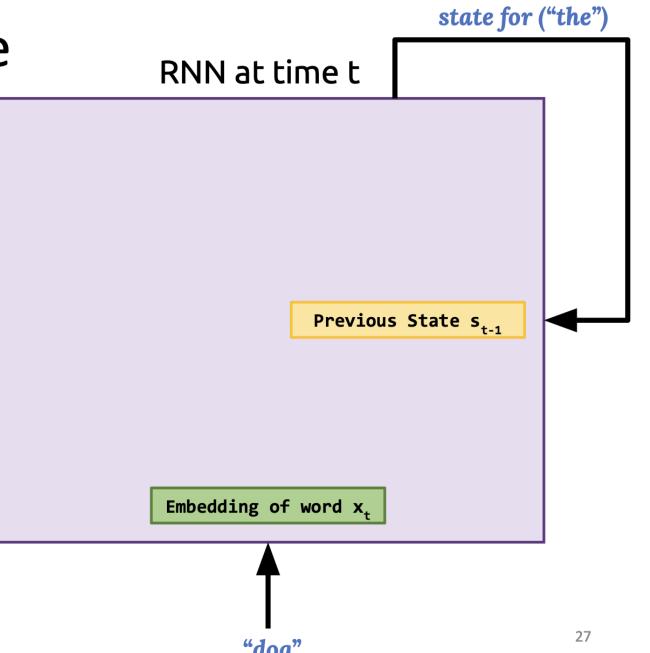
They pass previous *state* information from previous computations to the next.

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.

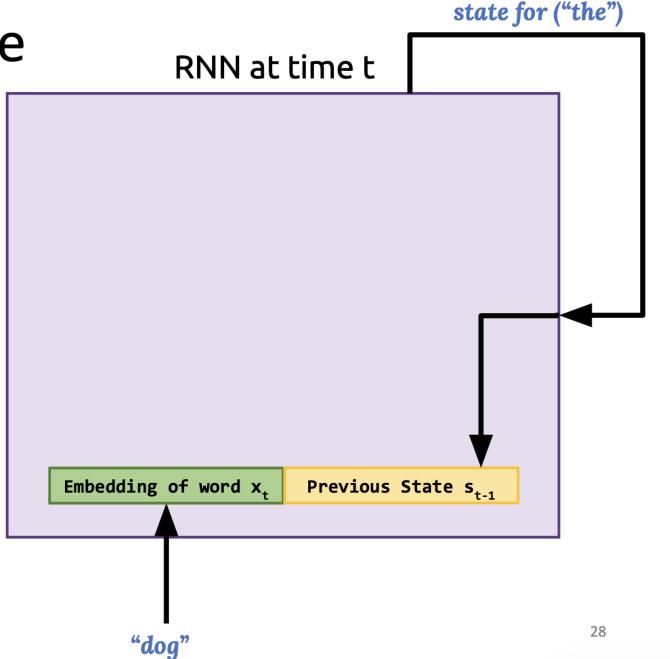
Let's see how we can build one!

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.

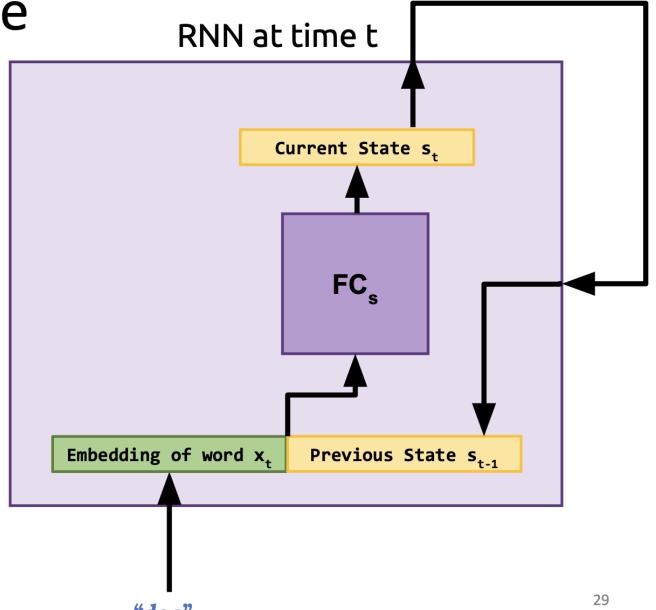
We then concatenate the embedding and state vectors.



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

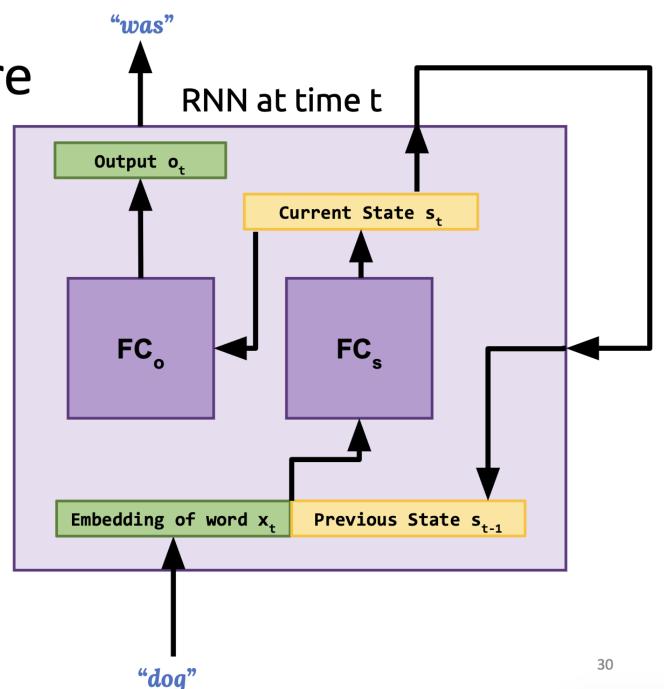


At each step of our RNN, we will get an input word, and a state vector from the previous cell.

We then concatenate the embedding and state vectors.

We use a fully connected layer to compute the next state

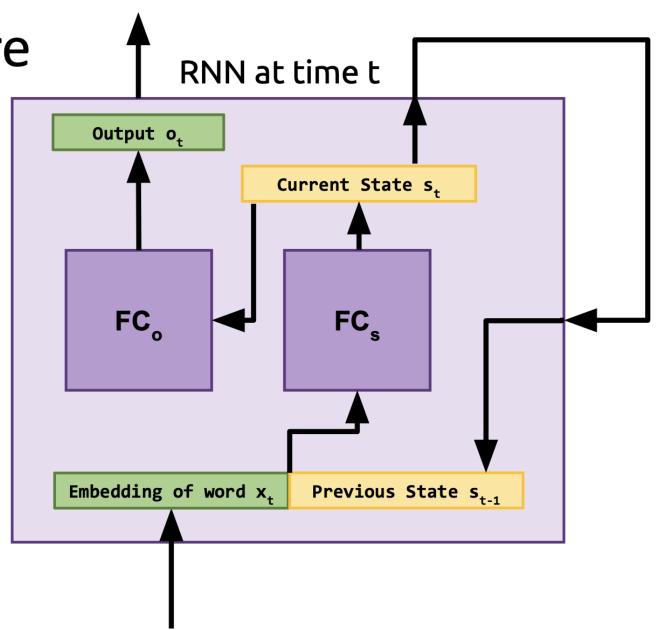
We use another connected layer to get the output.



We can represent the RNN in with the following equations:

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

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



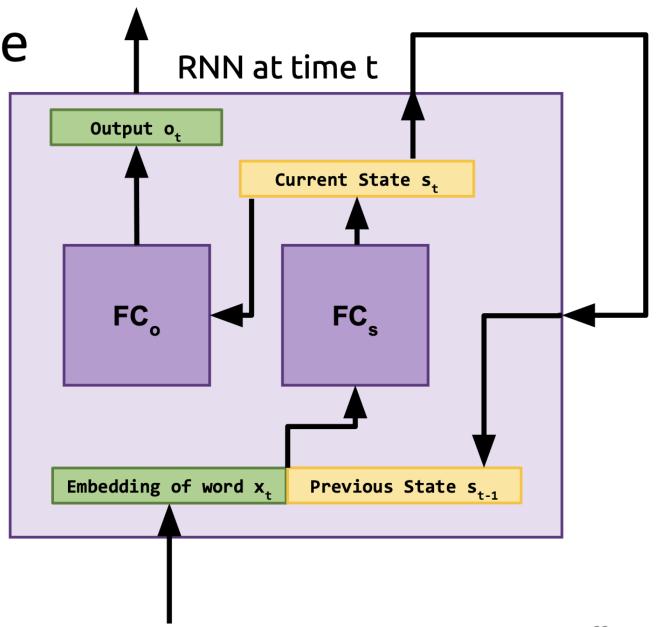
We can represent the RNN in with the following equations:

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

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

Nonlinear activations (e.g. sigmoid, tanh)

Any questions?



We can represent the RNN in with the following equations:

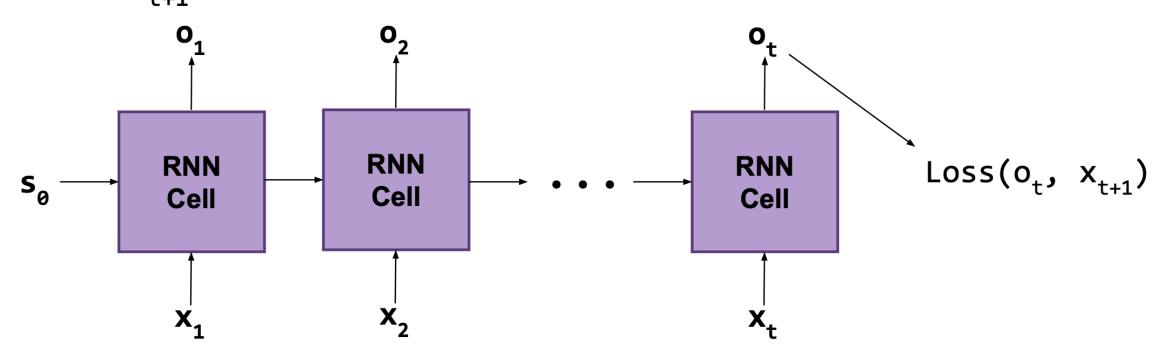
$$s_t = \rho \big((e_t, s_{t-1}) W_r + b_r \big)$$

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

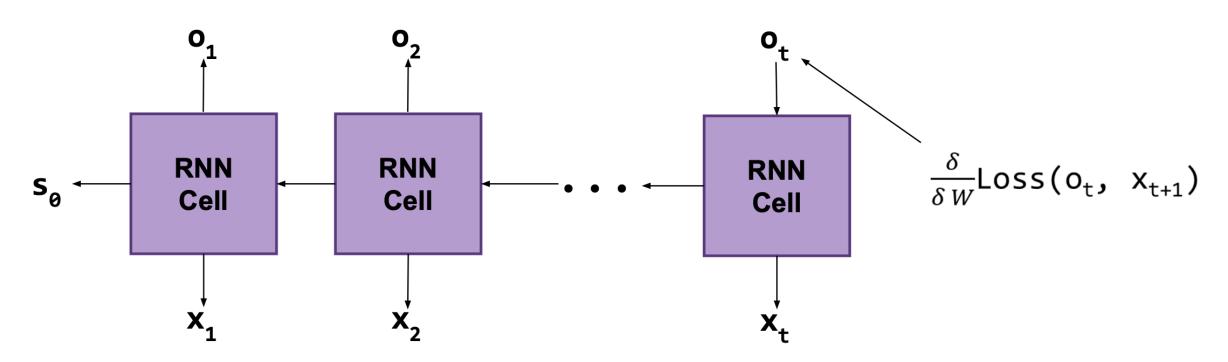
This brings up an immediate question: what is s_0 ?

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}

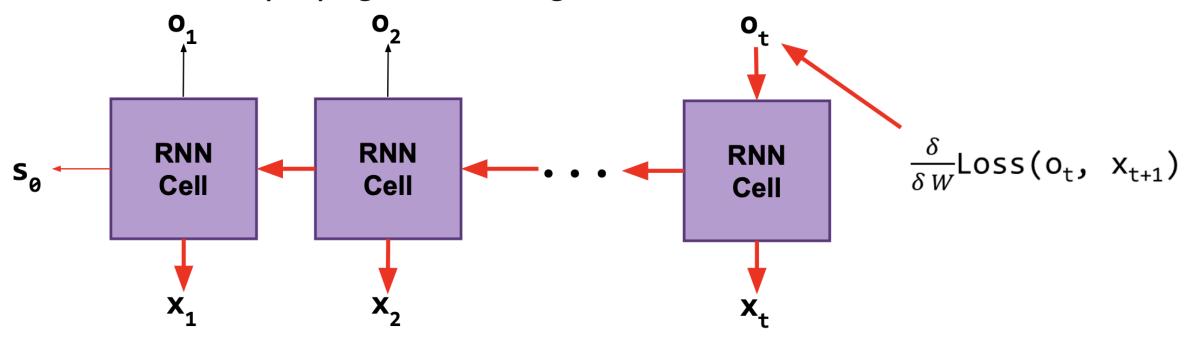


But what happens when we differentiate the loss and backpropagate?

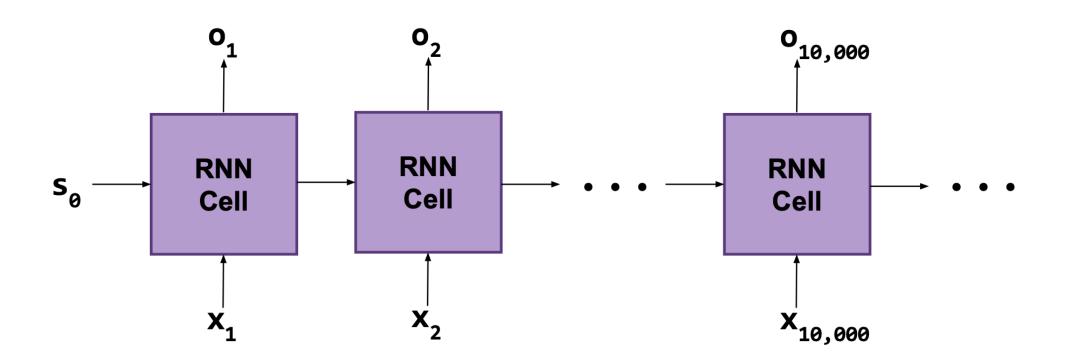


Not only do our gradients for o_t depend on x_t , but also on all of the previous inputs.

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.

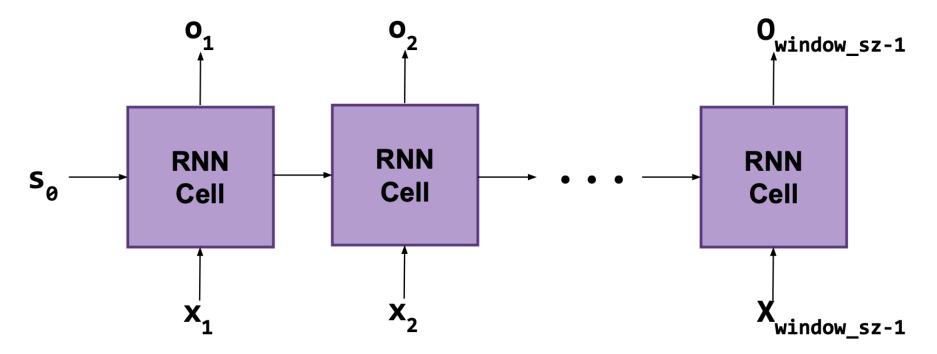
We now chop our corpus into sequences of words of size window_sz

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

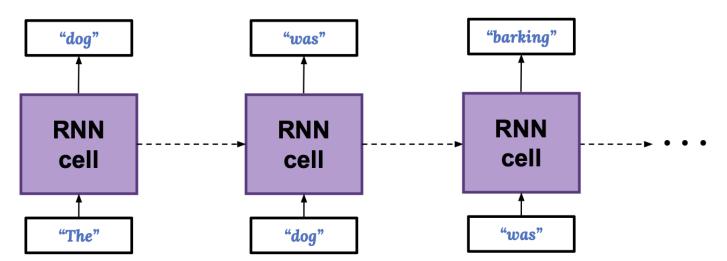


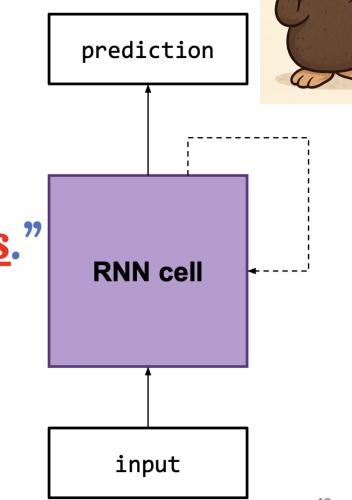
Any questions?

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

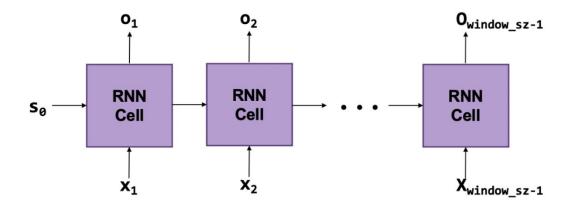
Any intuition why we would want return_sequences to be TRUE?

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
```

RNN

"The dog that my family had when I was a child had a fluffy

Want: "tail"



RNN Weaknesses

But....RNNs are not very good at remembering things *far* in the past.

Issue #1: Vanishing Gradients

Issue #2: What should we remember?



RNN Weaknesses

"The dog that my family had when I was a child had a fluffy

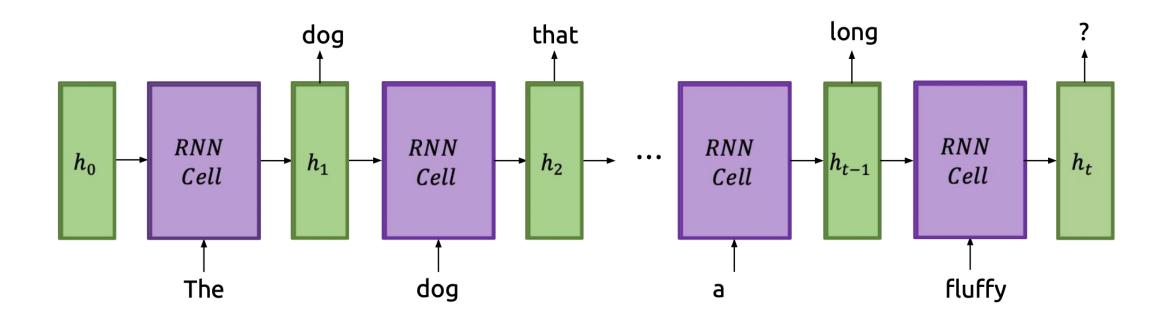
To predict "tail" RNN needs to remember the subject of the sentence
 "dog"

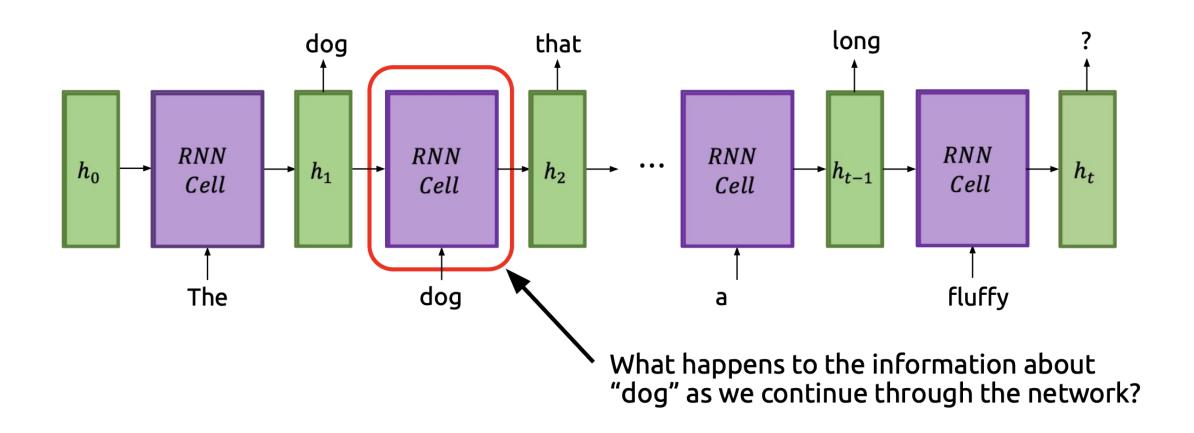
RNN Weaknesses

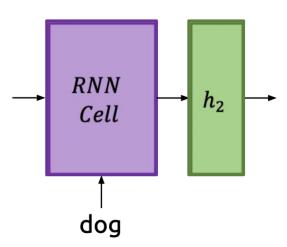
"The dog that my family had when I was a child had a fluffy

- To predict "tail" RNN needs to remember the subject of the sentence
 "dog"
- "dog" and predicted word are separated by 12 words
 - On the outer limit of what a vanilla RNN would be able to remember.

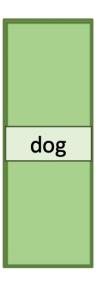
An Illustrative Example:



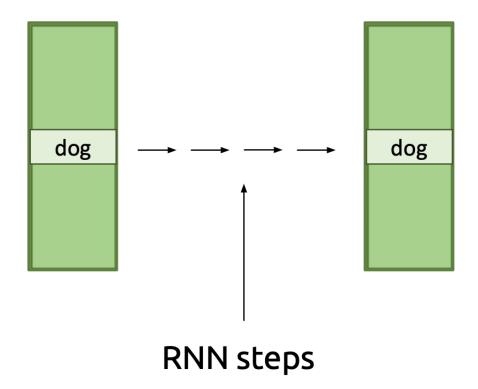




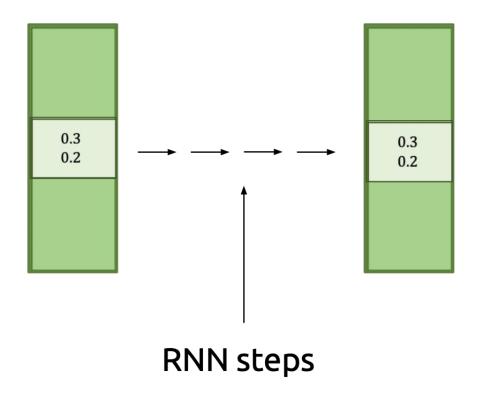
Can imagine that the information about "dog" is stored in some part of the RNN's hidden state vector



Through all subsequent RNN steps, we want "dog" to stay the same



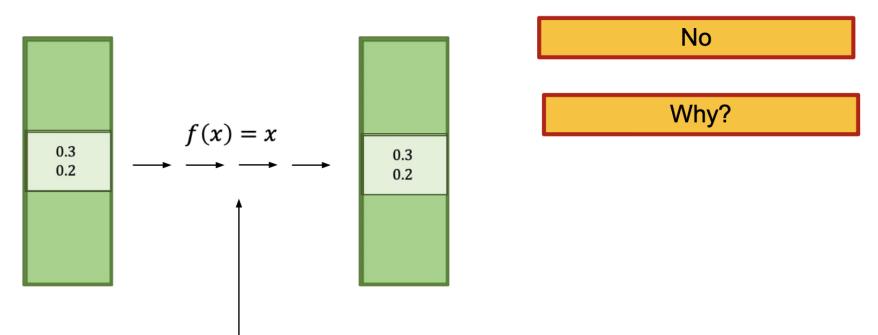
If we think of "dog" as just a few entries in the vector...



...to preserve "dog", we need to compute the identity function over

RNN steps

the part of the vector that stores it



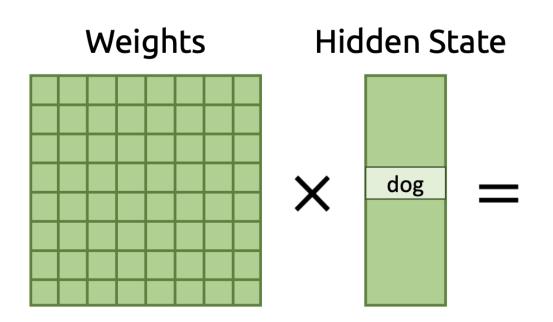
But will that happen?

RNN update

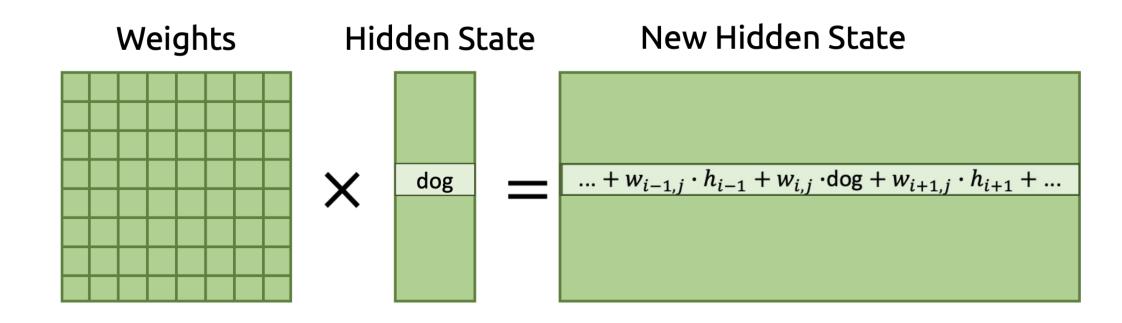
$$h_t = \rho((e_t, h_{t-1})W_r + b_r)$$

The hidden state goes through a fully connected layer!

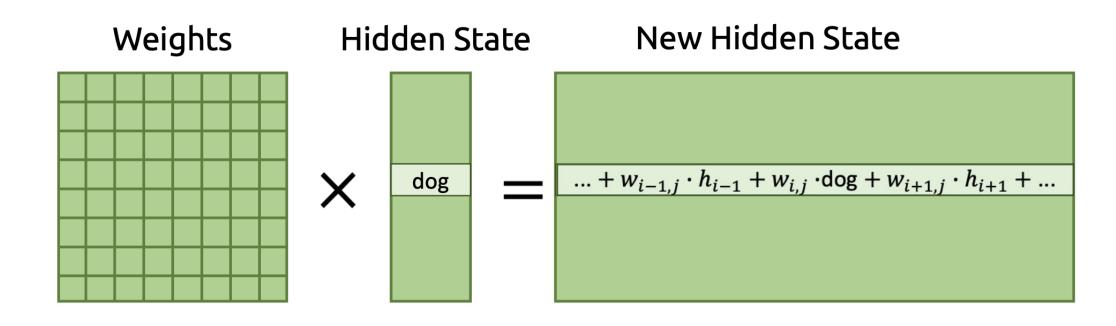
 What will happen to our dog after we multiply our weights by our hidden state?



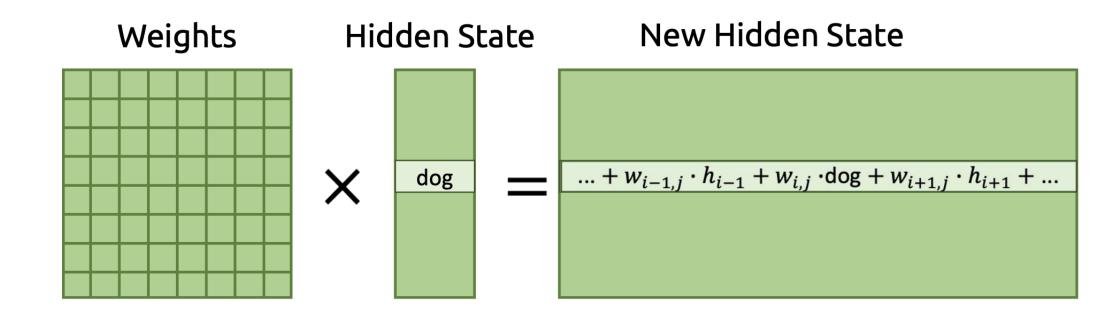
 What will happen to our dog after we multiply our weights by our hidden state?



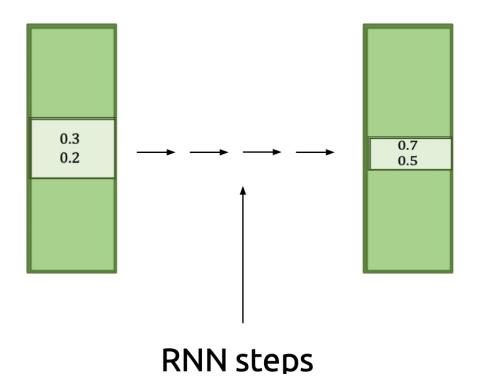
Dog gets lost in all the other information!



Dog gets lost in all the other information!

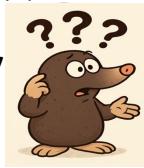


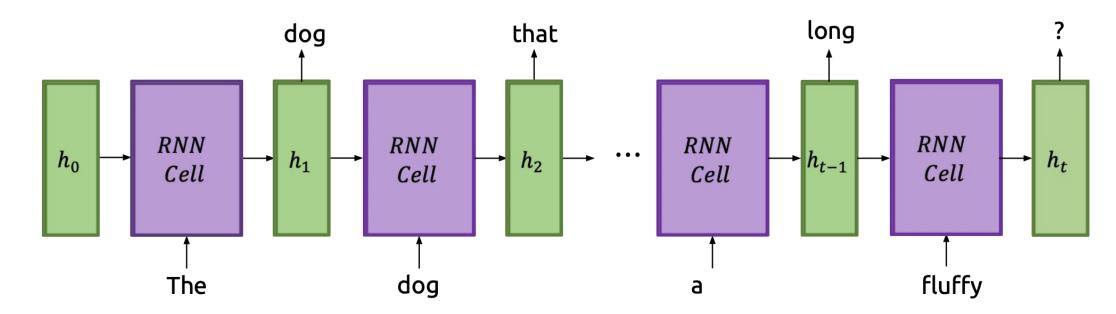
 "dog" in hidden state gets combined and mixed with rest of hidden state



RNN forgets about the dog after a certain time

RNNs cannot learn "long term" dependency



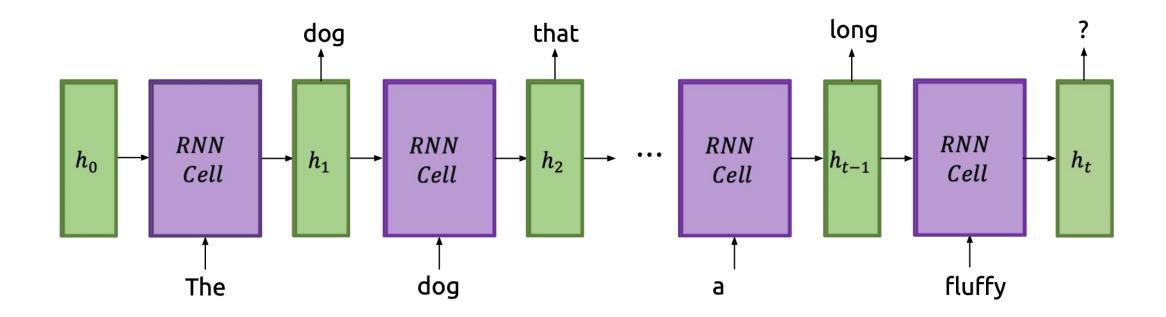


We need new way to update hidden state!

How?

An analogy to human (or computer) memory:

- RNN hidden state → "short term memory/RAM"
 - Like how you lose contents of RAM if you shut down a computer...
 - ...or how human short-term memory fades after time

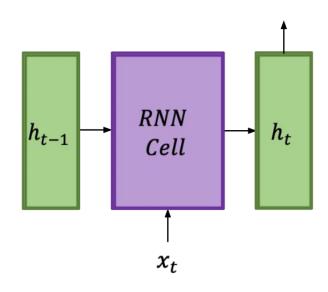


An analogy to human (or computer) memory:

- RNN hidden state → "short term memory/RAM"
 - Like how you lose contents of RAM if you shut down a computer...
 - ...or how human short-term memory fades after time
- What we want → "long term memory/disk"
 - Some state representing knowledge that persists
 - Like how contents of disk persist across shut-downs...
 - · ...or how sleep consolidates human memory into long-term memory
- <u>Long</u> Short Term Memory (LSTM)
 - "Short-term memory that persists over time"
 - i.e. "hidden states that remember information for longer"

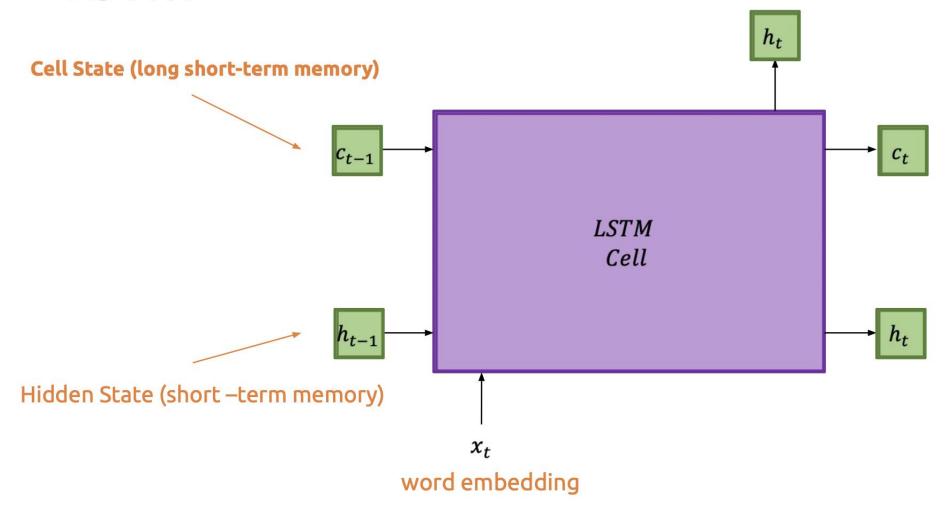
What is different?

Vanilla RNN



LSTM h_t **LSTM** Cell x_t

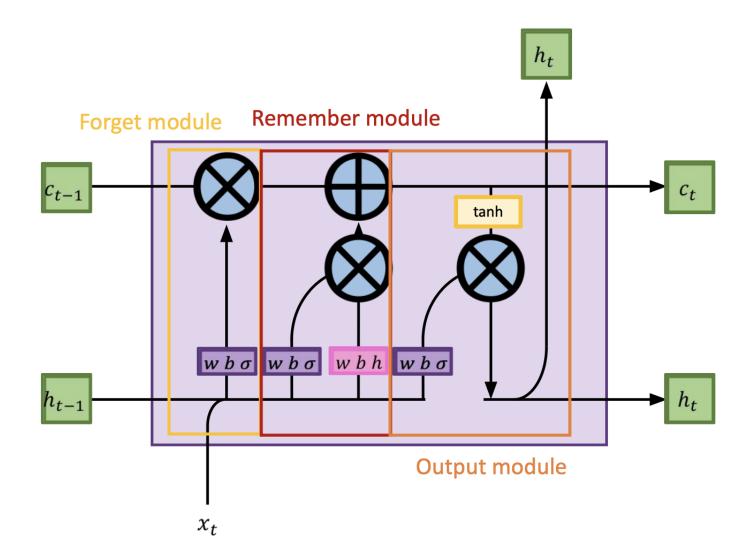
LSTM



How an LSTM works

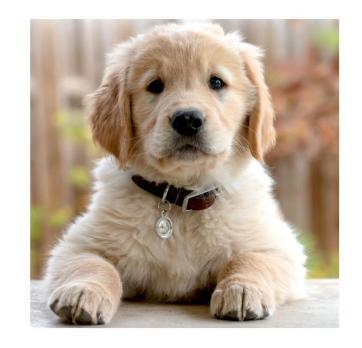
- An LSTM consists of 3 major modules:
 - Forget module
 - Remember module
 - Output module

The Complete LSTM



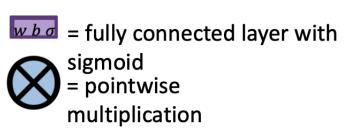
Say we just predicted "tail" in "My dog has a fluffy _____."

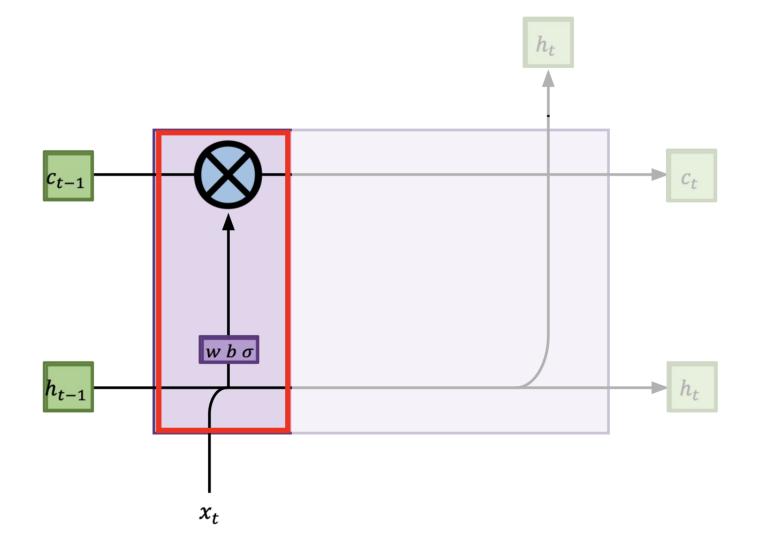
Next set of words: "I love my dog"



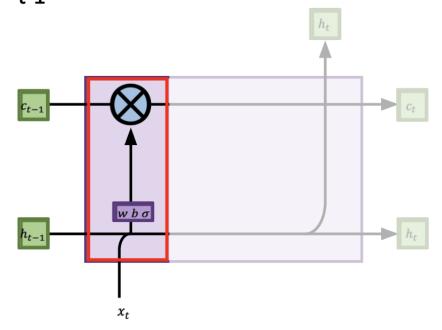
- Model no longer needs to know about "dog"
- Ready to **delete** information about subject





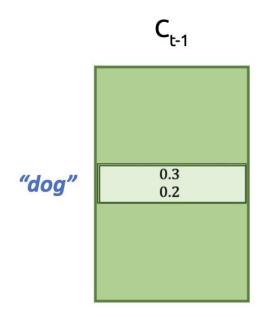


- Filters out what gets allowed into the LSTM cell from the last state
 - Example: If it's remembering gender pronouns, and a new subject is seen, it will forget the old gender pronouns
- Either lets parts of C_{t-1} pass through or not



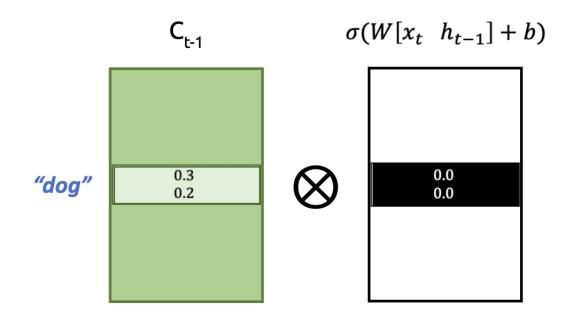
Forgetting information

- Use pointwise multiplication by a mask vector to forget information
 - What do we want to forget from last cell state?



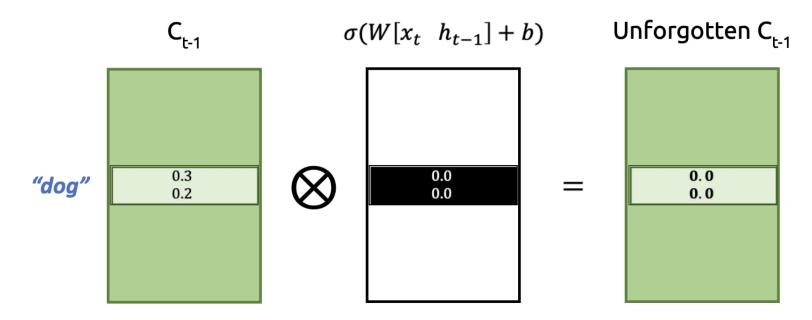
Forgetting information

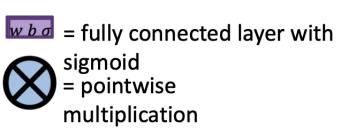
- Use pointwise multiplication by a mask vector to forget information
 - What do we want to forget from last cell state?
 - Output of fully connected + sigmoid is what we want to forget

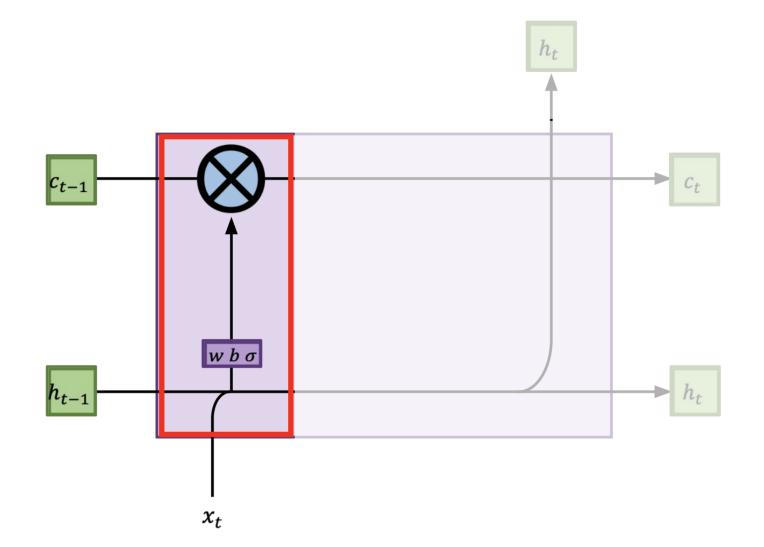


Forgetting information

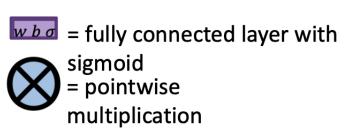
- Use pointwise multiplication by a mask vector to forget information
 - What do we want to forget from last cell state?
 - Output of fully connected + sigmoid is what we want to forget
 - "Zeros out" a part of the cell state
 - Pointwise multiplication by a learned mask vector is known as gating

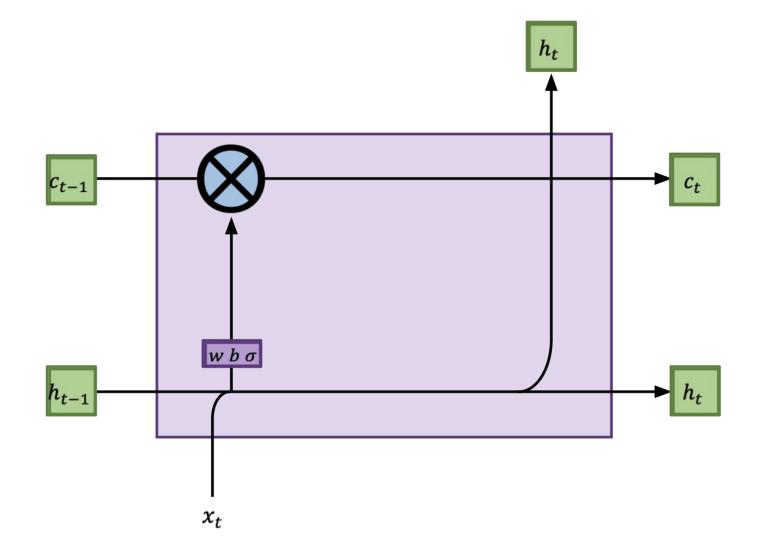






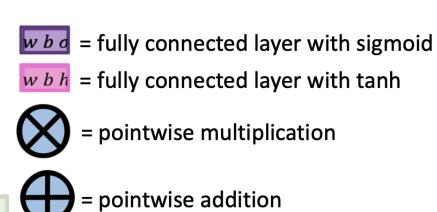
What's next?

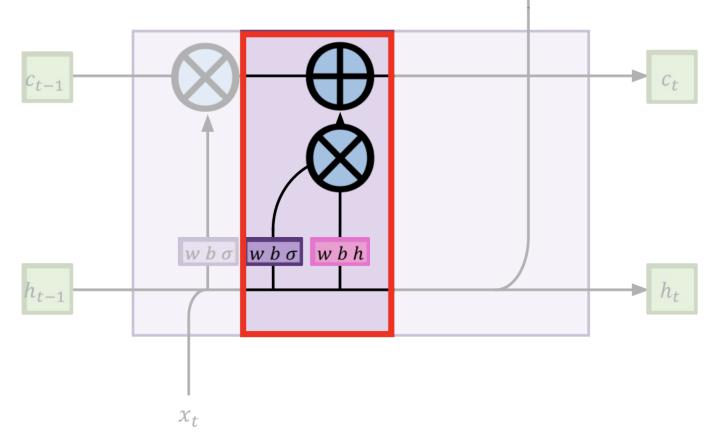




Remember Module

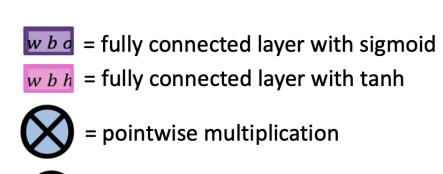
 We can save information that we want to remember by adding it into "empty" slots in the cell state



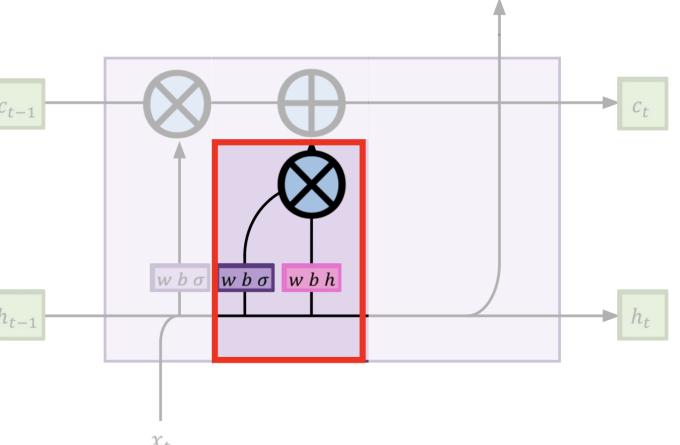


Remember Module

 First: use gating to decide what to remember

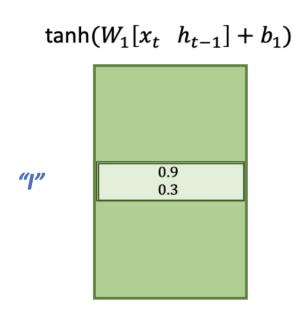


= pointwise addition



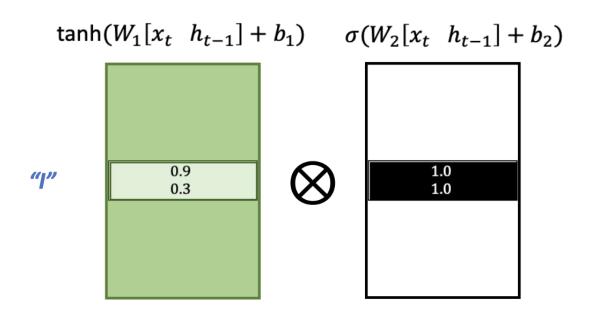
Gating for 'selective memory'

 A fully-connected + tanh on [input, memory] computes some new memory



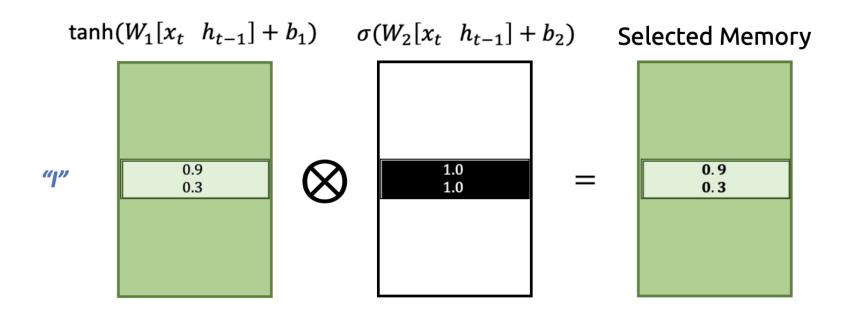
Gating for 'selective memory'

- A fully-connected + tanh on [input, memory] computes some new memory
- We gate this memory to decide what bits of it we want to remember long-term in the cell state



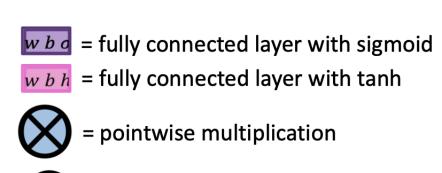
Gating for 'selective memory'

- A fully-connected + tanh on [input, memory] computes some new memory
- We gate this memory to decide what bits of it we want to remember long-term in the cell state

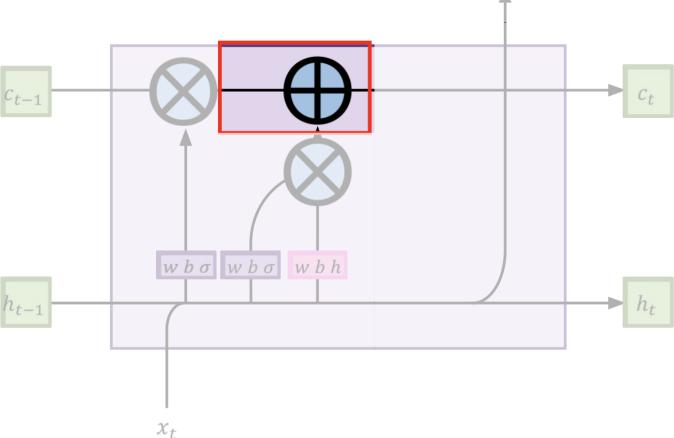


Remember Module

 Then: we add this selective memory into the cell state



= pointwise addition

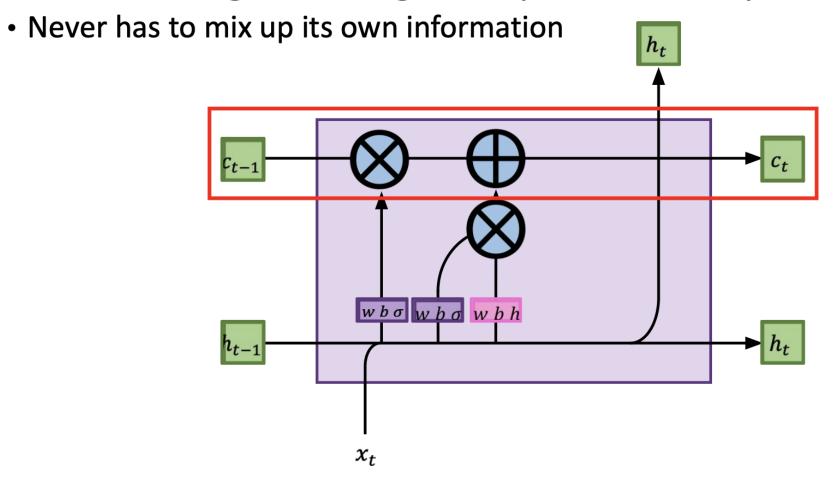


Remembering information

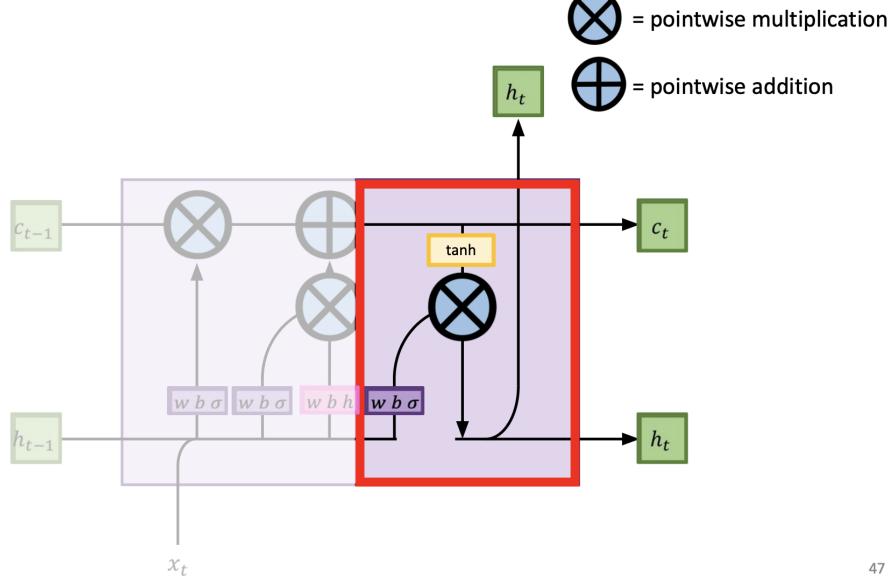
 Add what we <u>didn't forget</u> to what we <u>did remember</u> Unforgotten C₁₋₁ **Selected Memory** 0.0 0.9 0.9 + 0.0 0.3 0.3

Why does this solve our problem?

Cell state never goes through a fully connected layer!



Output Module



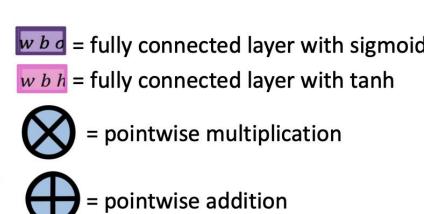
wbo = fully connected layer with sigmoid

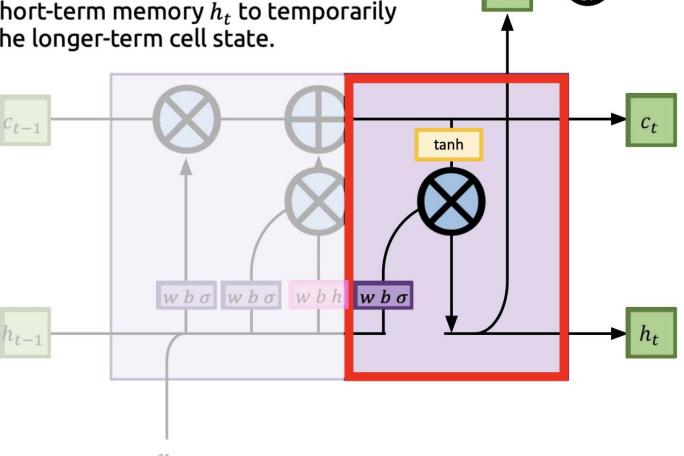
w b h = fully connected layer with tanh

Output Module

Same structure as the remember module

• Provides path for short-term memory h_t to temporarily acquire info from the longer-term cell state.





The Complete LSTM



$$i_{t} = \sigma(W_{i}h_{t-1} + U_{i}x_{t} + b_{i})$$

$$f_{t} = \sigma(W_{f}h_{t-1} + U_{f}x_{t} + b_{f})$$

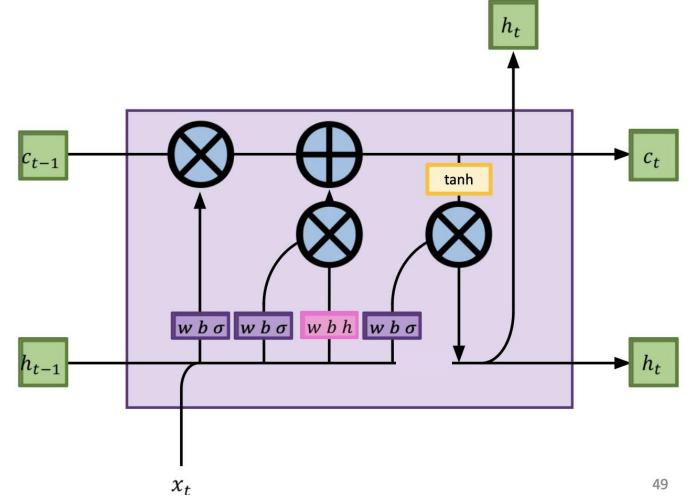
$$o_{t} = \sigma(W_{o}h_{t-1} + U_{o}x_{t} + b_{o})$$

$$\widetilde{c}_{t} = tanh(Wh_{t-1} + Ux_{t} + b)$$

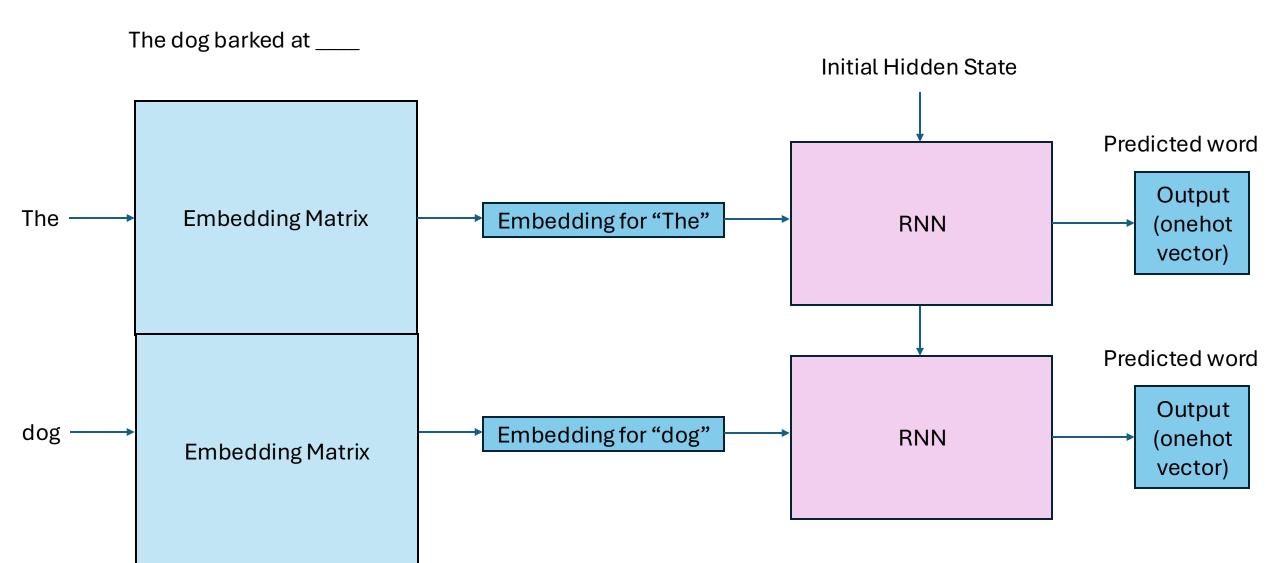
$$c_{t} = f_{t} \circ c_{t-1} + i_{t} \circ \widetilde{c}_{t}$$

$$h_{t} = o_{t} \circ tanh(c_{t})$$

$$y_{t} = h_{t}$$



Overview of RNN Sequence Prediction



When to compute loss

We have predictions and ground truths at every step, why not compute the loss after every word and backprop?

Should your model be penalized equally for incorrect predictions?

- 1) The ____
- 2) The dog barked at ____

Well... What task are you actually training it for?

Recap of LSTMs

Two "Outputs": Cell and Hidden states

- Hidden state h_t is used for output (and short term memory)
- Cell State used for long term memory

Forget Module:

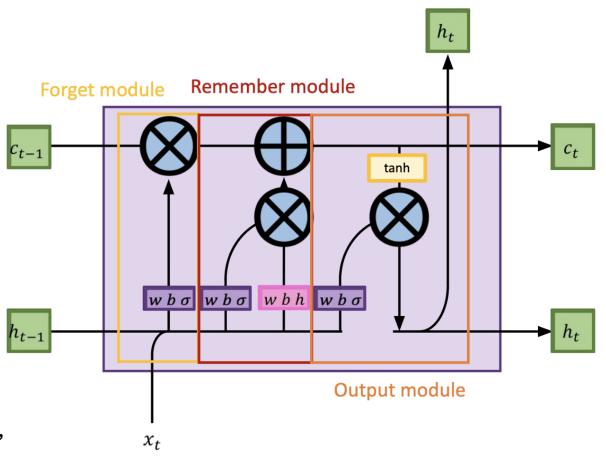
 Multiply cell state by numbers between 0 and 1 (0 forgets information, 1 keeps it)

Remember Module:

 Adds information from short term memory/input to long term memory

Output Module:

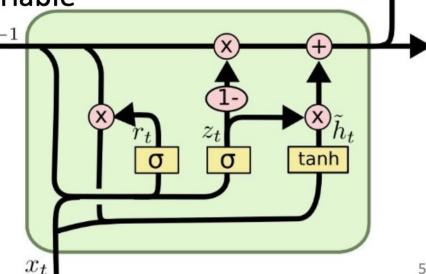
- Combines input (x_t) , short term memory (h_t) , and long term memory (c_t) .
- Produces output
- Output is passed along as short term memory



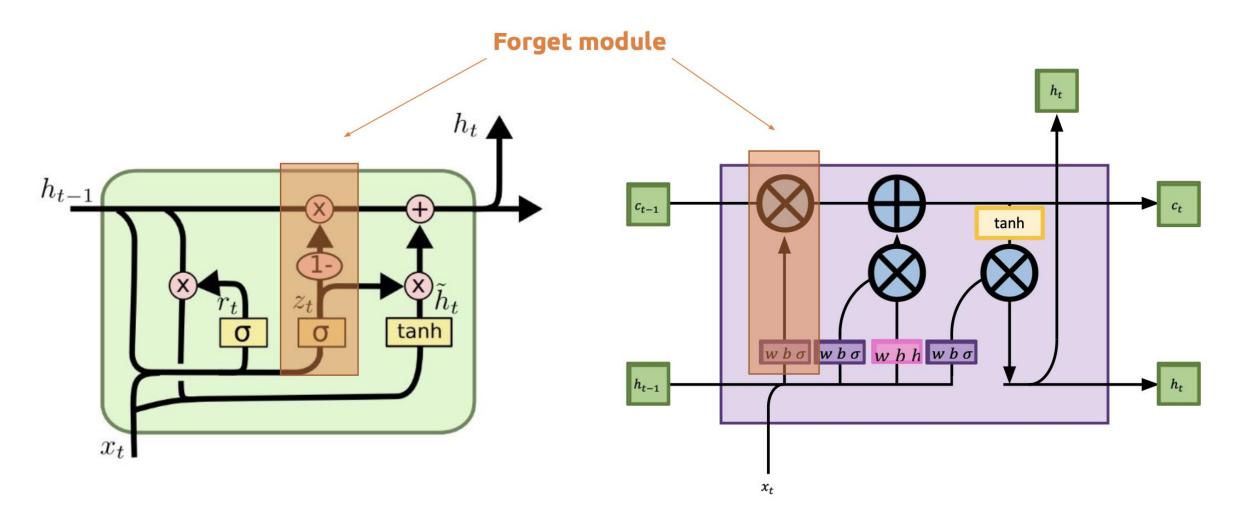
GRU

- Gated Recurrent Unit
- In practice, similar performance and may train faster
 - Removes cell state, computationally more efficient and less complex
- In theory, weaker than LSTMs since it <u>cannot unboundedly count</u>
 - · Counting: track increment or decrement of variable
 - e.g. Validate brackets in code

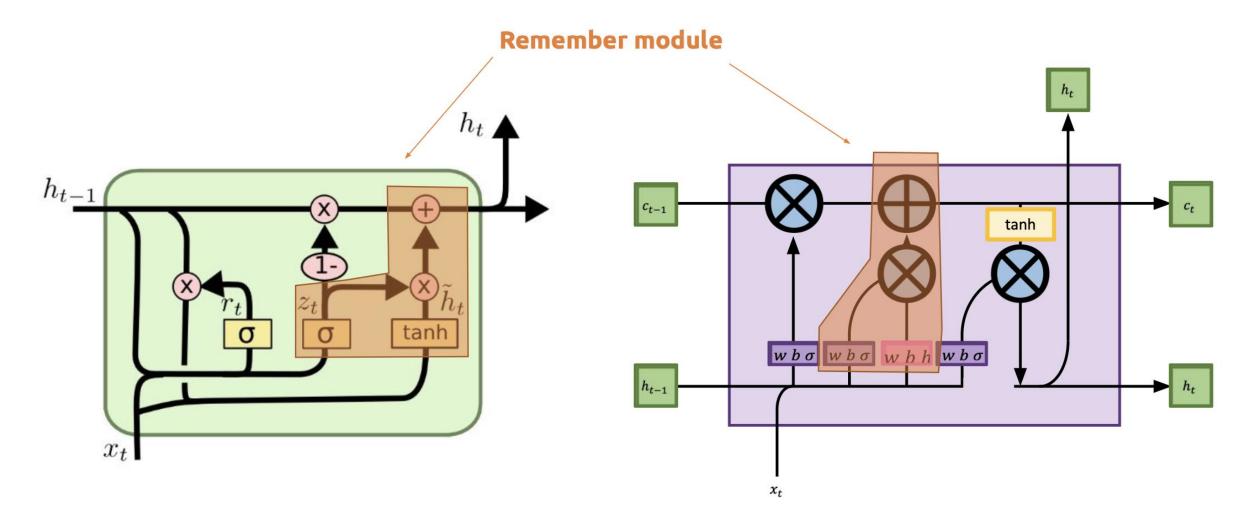
Requires counting brackets & nesting levels



GRU vs LSTM



GRU vs LSTM



GRU vs LSTM

