CSCI 1470

Eric Ewing

Friday,
3/7/25

# Deep Learning

Day 18: RNNs and LSTMs

# Some Game Theory

# Other Updates

- Will try to handle point deductions over the weekend

- For those who didn't come forward:
  - I will (eventually) look at **all** submissions and look for known AI patterns and send cases I'm certain of.
  - This won't be quick… I'll prioritize seniors and masters students first

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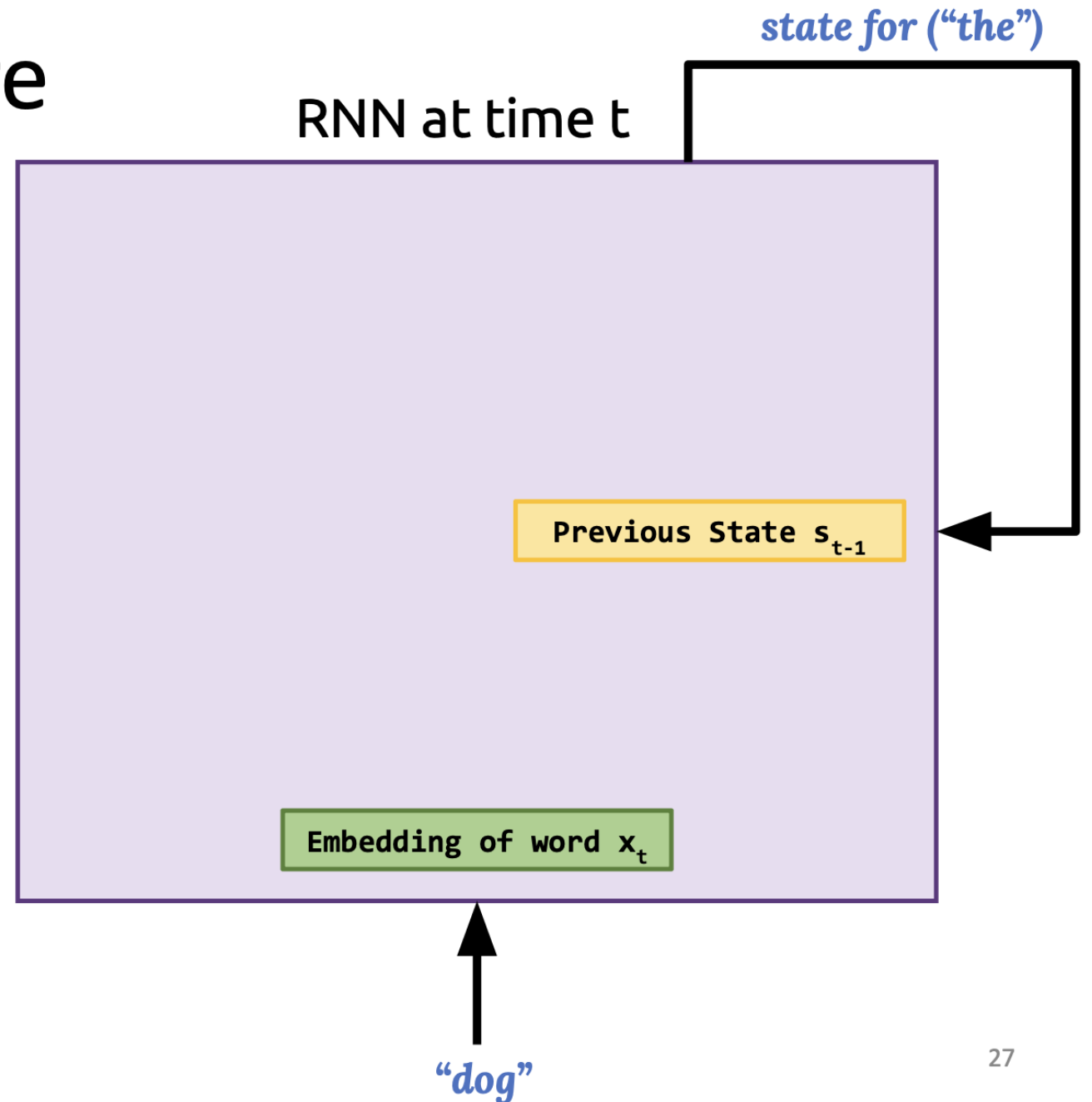
The block of computation that feeds its own output into its input is called the *RNN cell.*

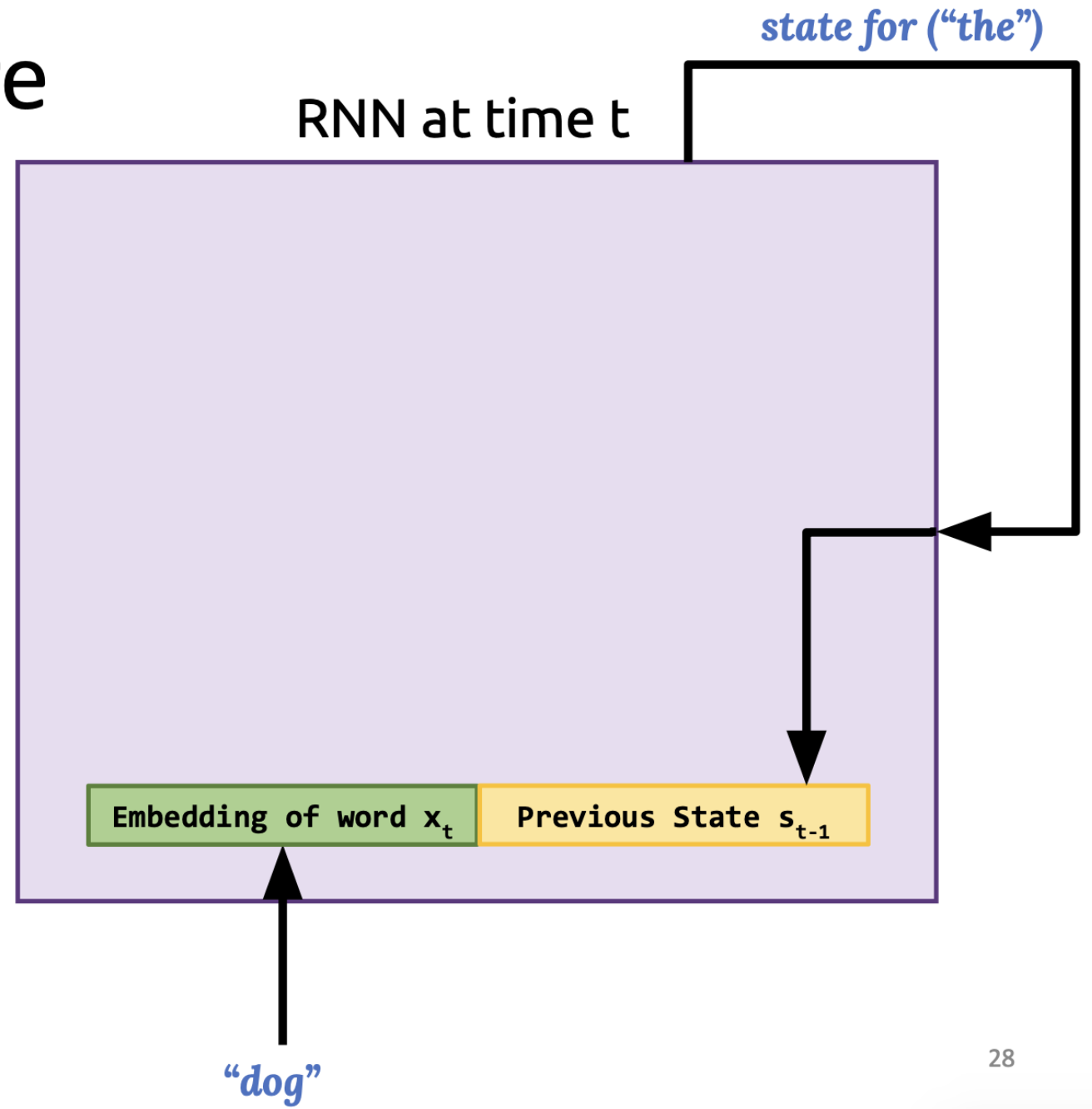Let's see how we can build one!

# RNN Cell Architecture

RNN at time t

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

Previous State $s_{t-1}$

Embedding of word $x_t$

*"dog"*

27

# RNN Cell Architecture

RNN at time t

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.

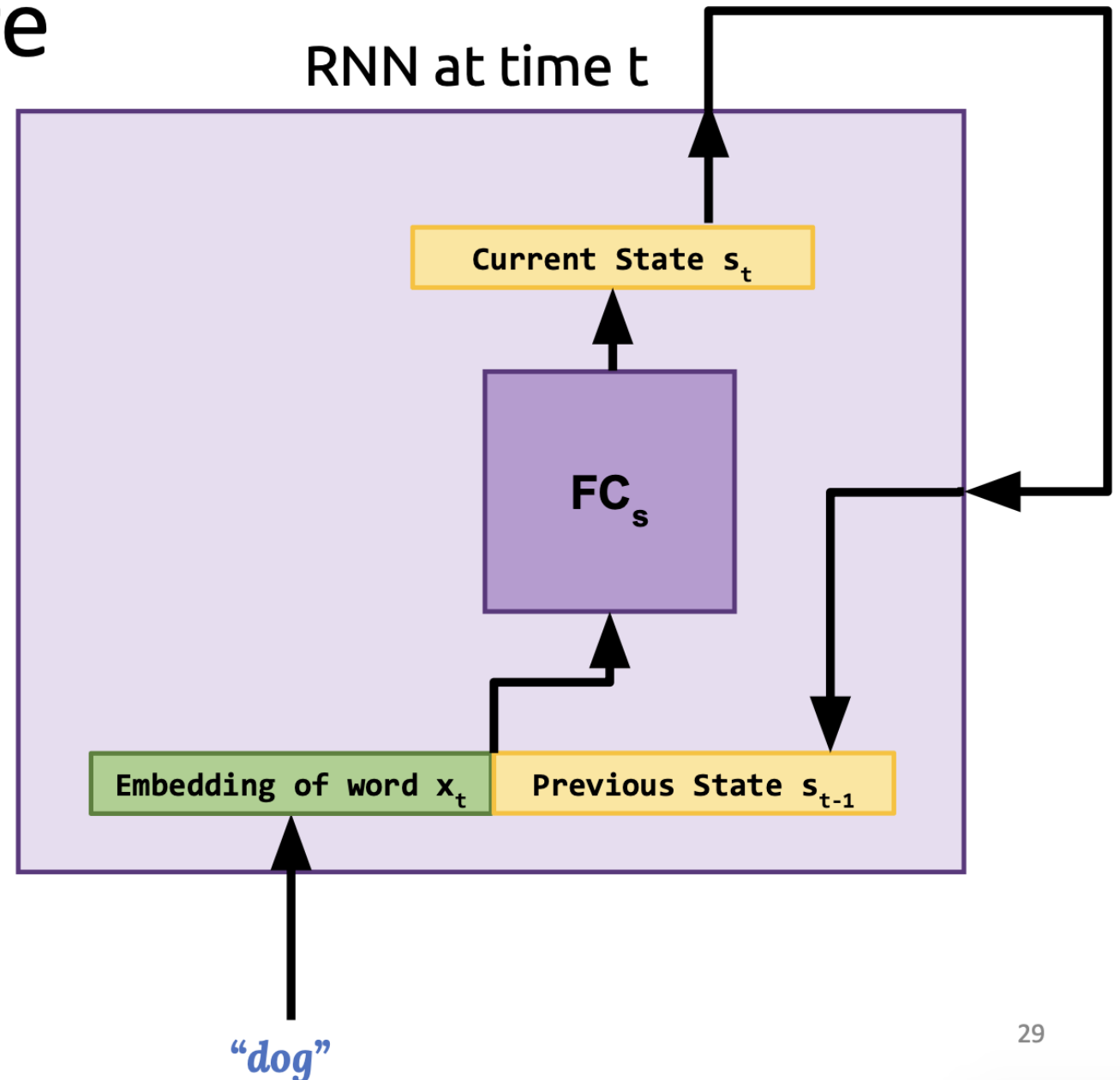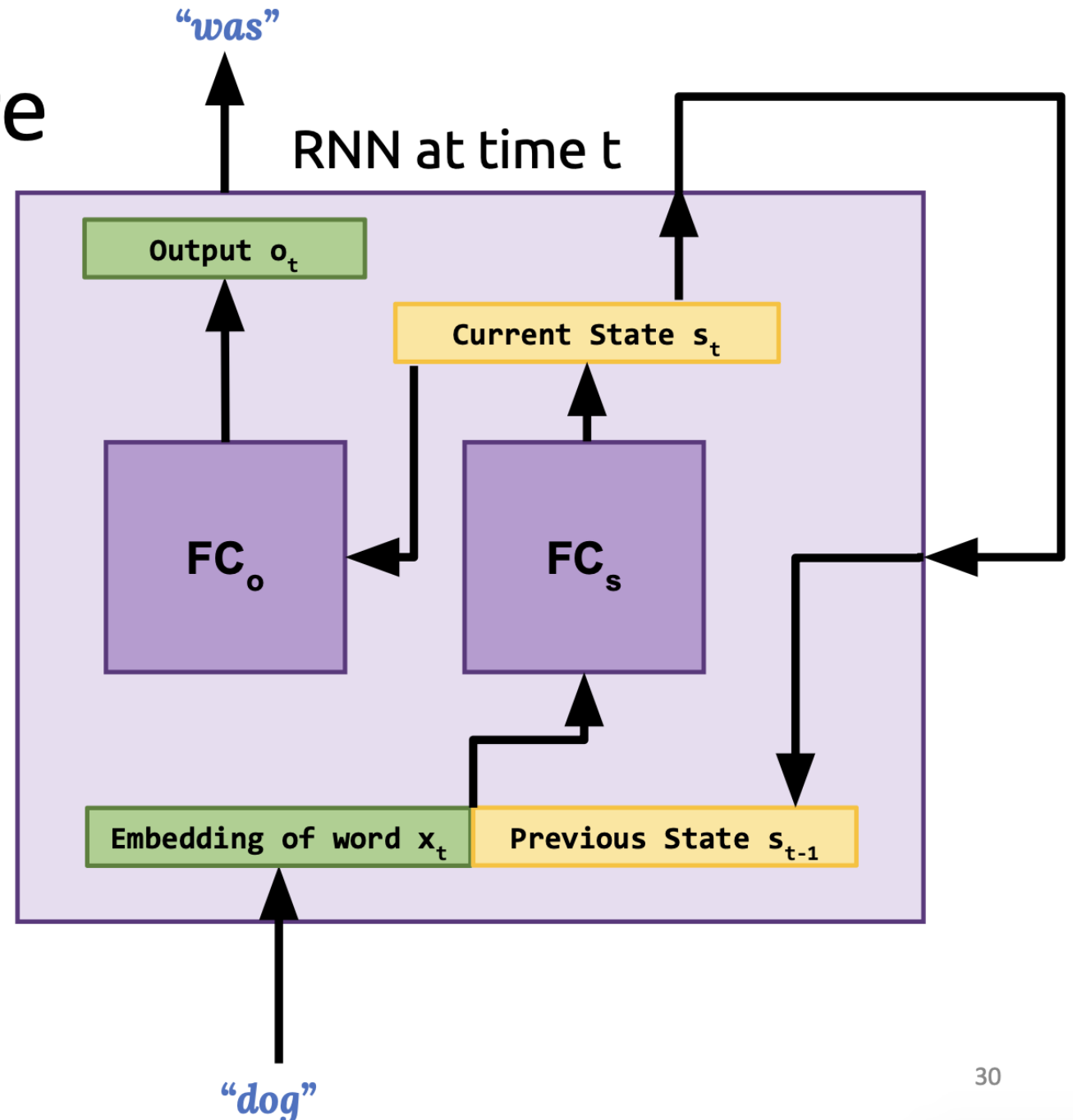| Embedding of word $x_t$ | Previous State $s_{t-1}$ |
|---|---|

"dog"

28

# 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 at time t

Current State $s_t$

$FC_s$

Embedding of word $x_t$ | Previous State $s_{t-1}$

*"dog"*

29

# 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

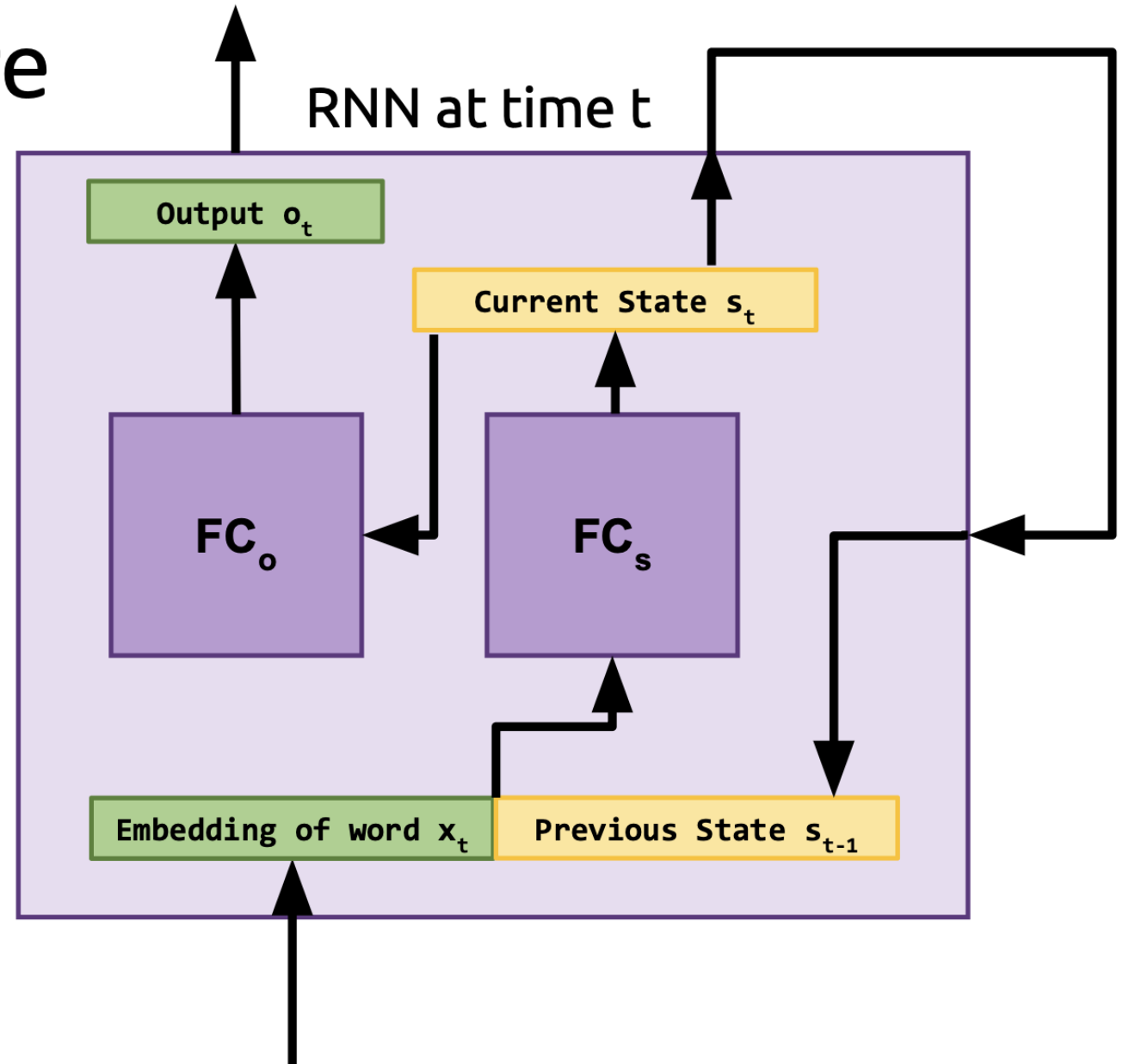We use another connected layer to get the output.

*"was"*

RNN at time t

Output $o_t$

Current State $s_t$

$FC_o$

$FC_s$

Embedding of word $x_t$

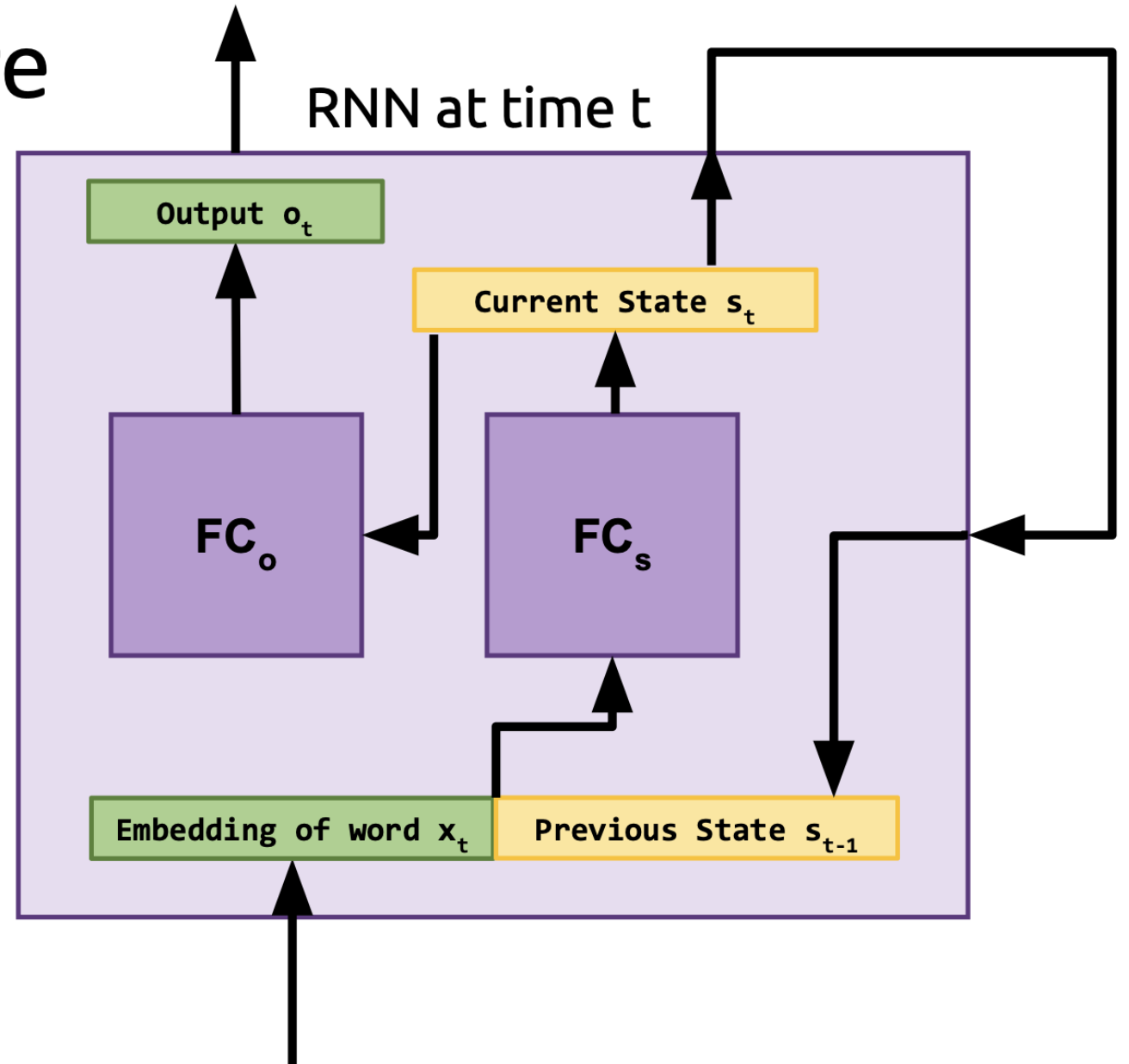Previous State $s_{t-1}$

*"dog"*

30

# RNN Cell Architecture

We can represent the RNN in with the following equations:

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

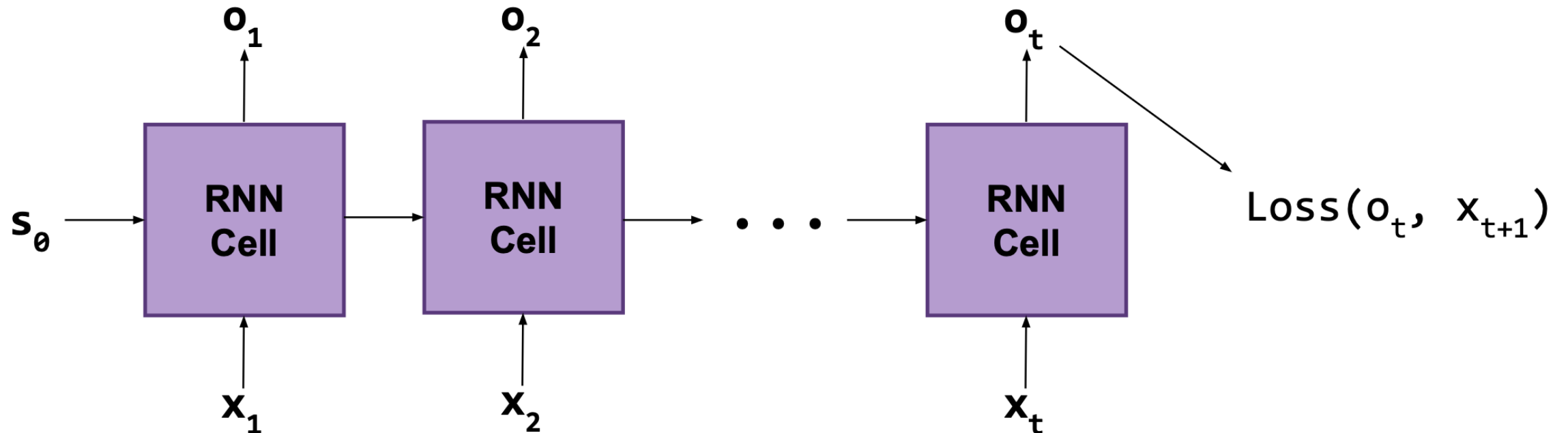$$o_t = \sigma(s_t W_o + b_o)$$

RNN at time t

Output $o_t$

Current State $s_t$

FC$_o$

FC$_s$

Embedding of word $x_t$

Previous State $s_{t-1}$

31

# RNN Cell Architecture

We can represent the RNN in with the following equations:

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

$$o_t = \boxed{\sigma}(s_t W_o + b_o)$$

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

Any questions?

RNN at time t

Output $o_t$

Current State $s_t$

FC$_o$

FC$_s$

Embedding of word $x_t$

Previous State $s_{t-1}$

# RNN Cell Architecture

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")

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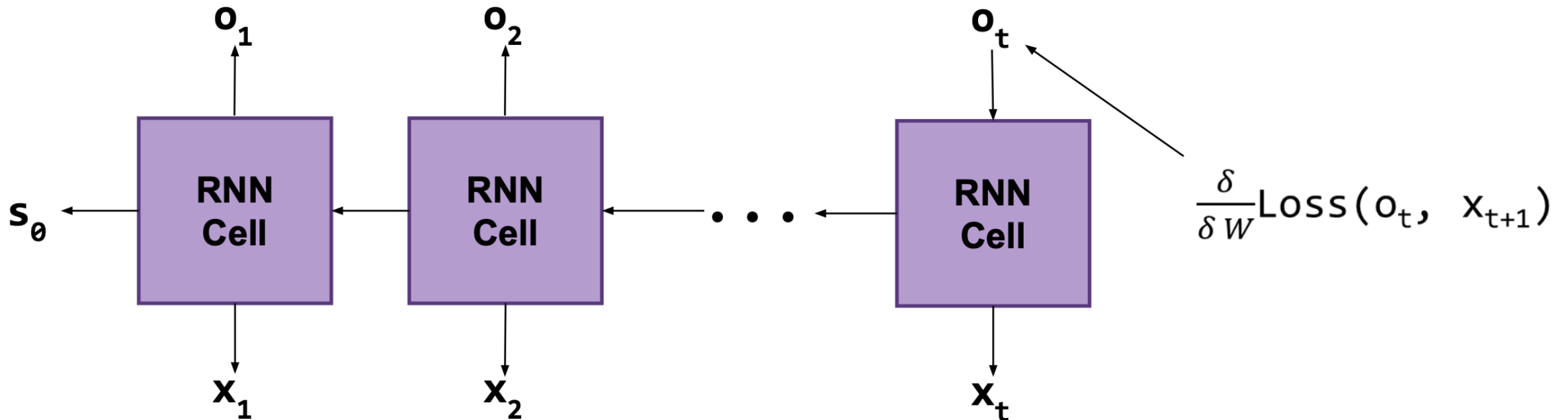
We can calculate the cross entropy loss just as before since for any sequence of input words $(x_1, x_2, \ldots, x_t)$, we know the true next word $x_{t+1}$
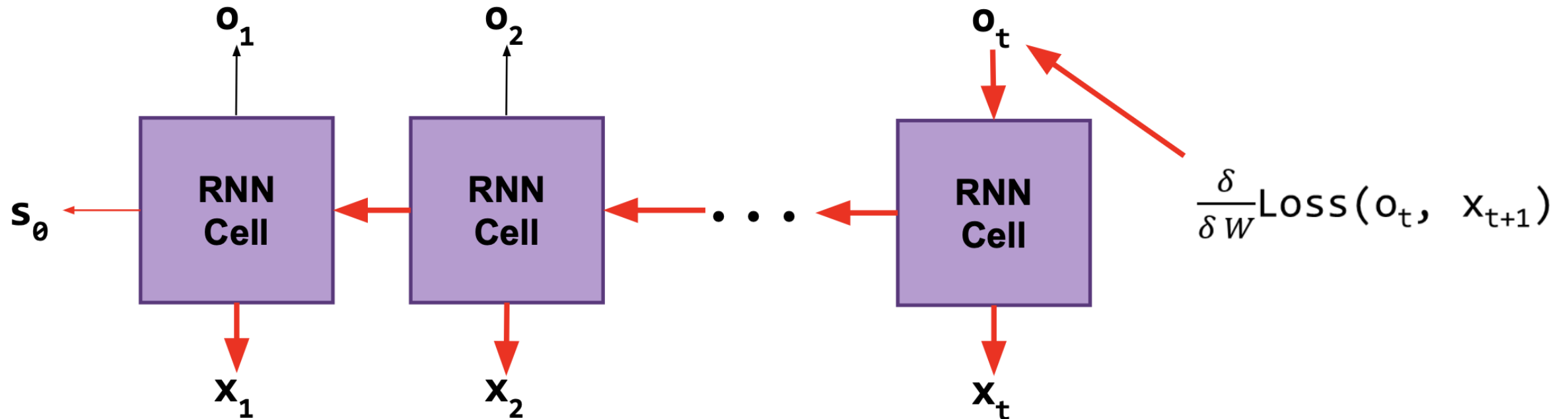
# Training RNNs

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.
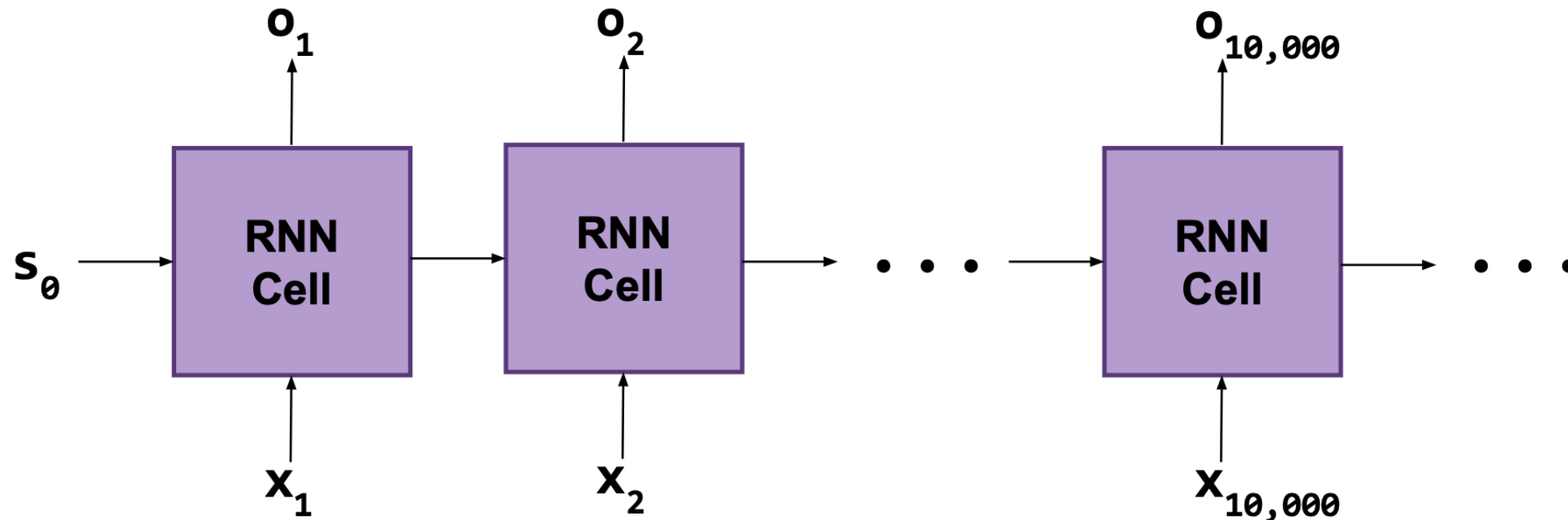
We call this *backpropagation through time.*



$$\frac{\delta}{\delta W} \text{Loss}(o_t, \; x_{t+1})$$

# Training RNNs

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`.

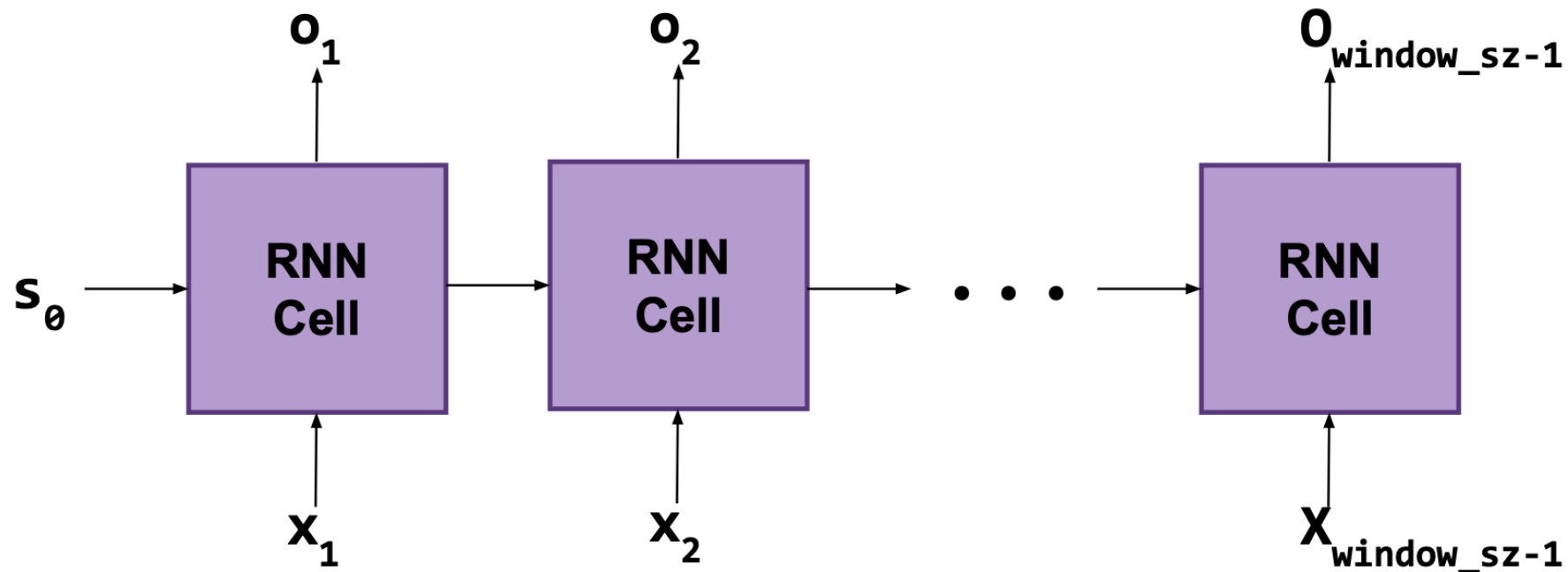We now chop our corpus into sequences of words of size `window_sz`

The new shape of our data should be:

$$\texttt{(batch\_sz, window\_sz, embedding\_sz)}$$

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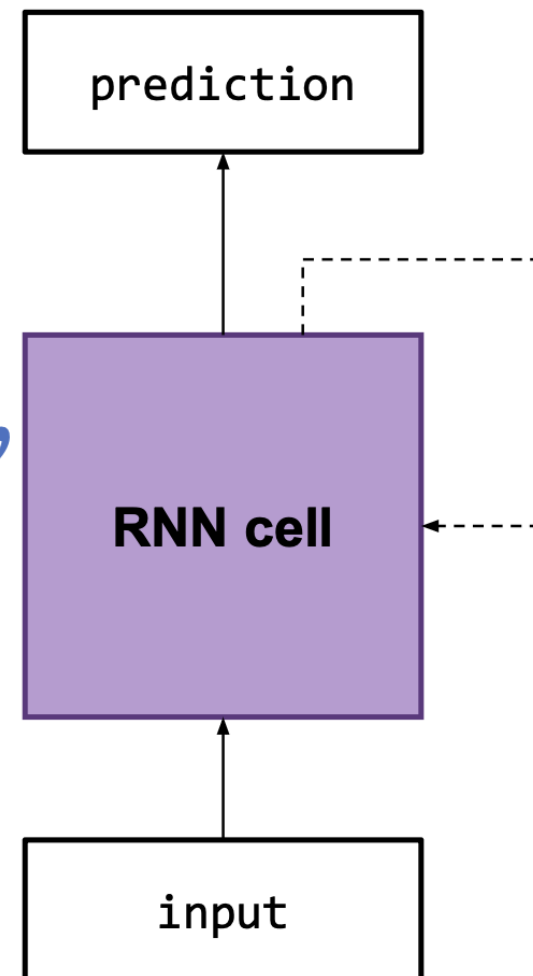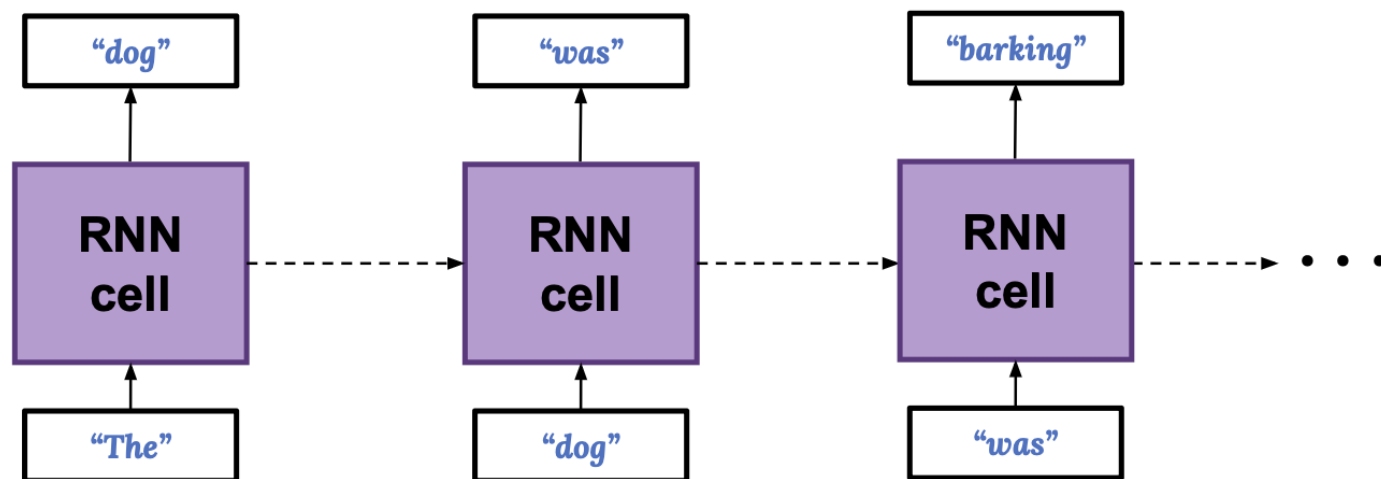
# Does RNN fix the limitations of the N-gram model?

1. Number of 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 <u>cats</u>."*



40

# 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

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)

45

# 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
```

# 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.

# 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

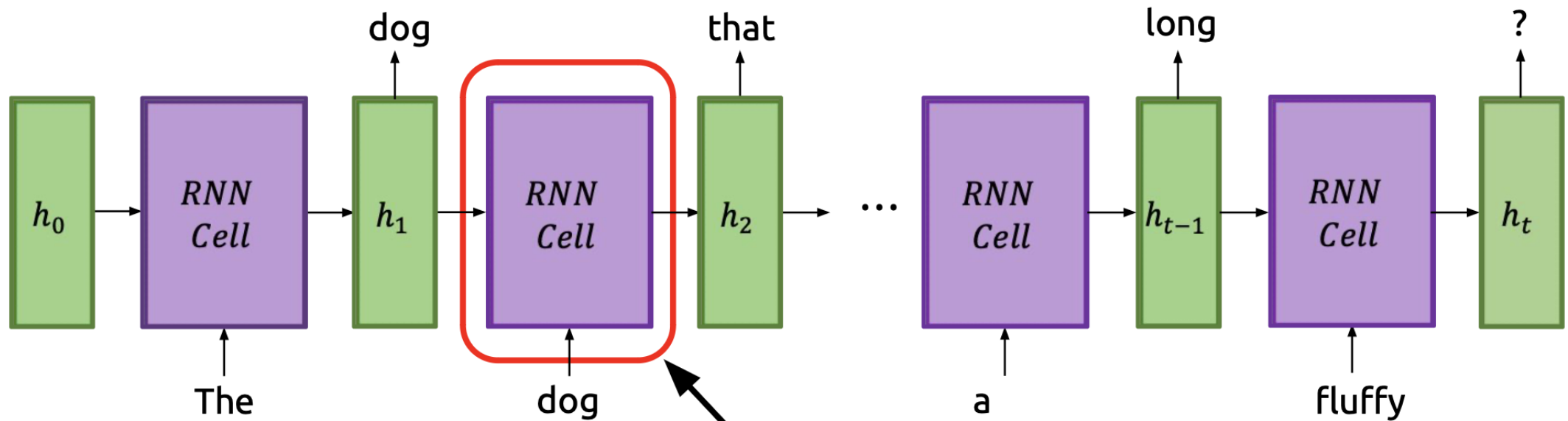*"The dog that my family had when I was a child had a fluffy \_\_\_\_."*

- 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:

# An Illustrative Example:



What happens to the information about "dog" as we continue through the network?
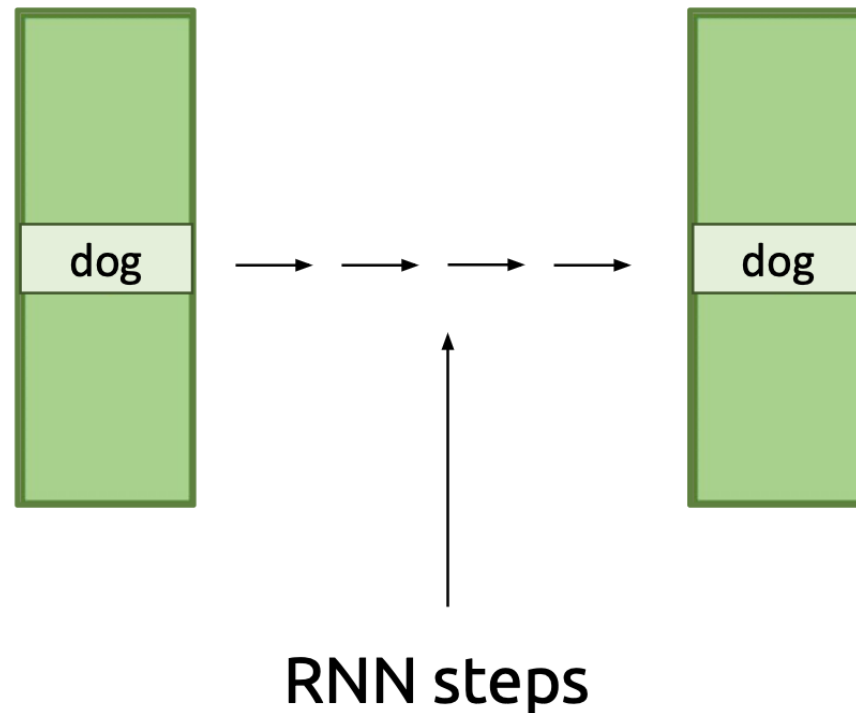
# An Illustrative Example:

# An Illustrative Example:

Can imagine that the information about *"dog"* is stored in some part of the RNN's hidden state vector
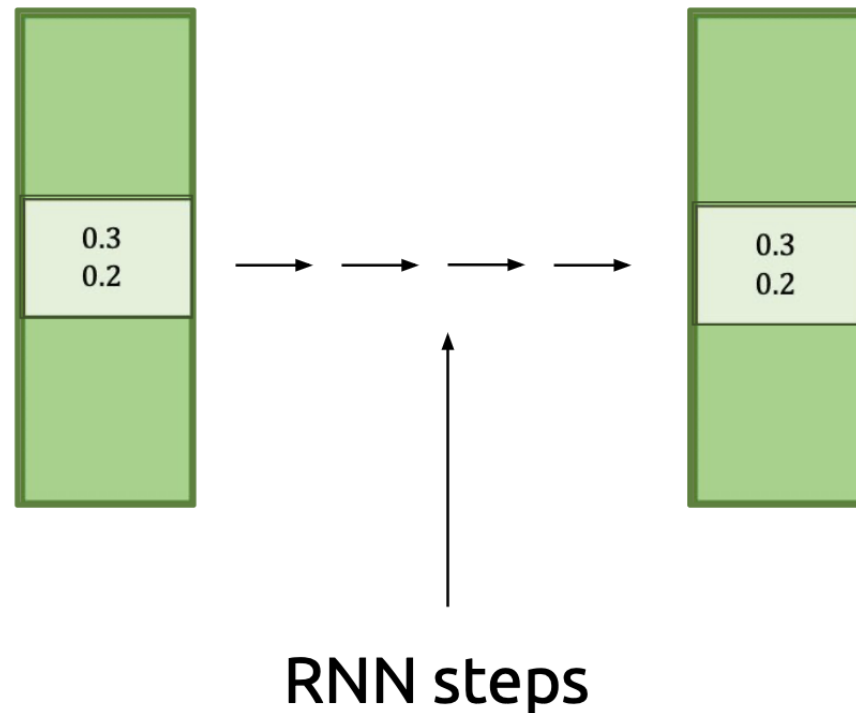
# An Illustrative Example:

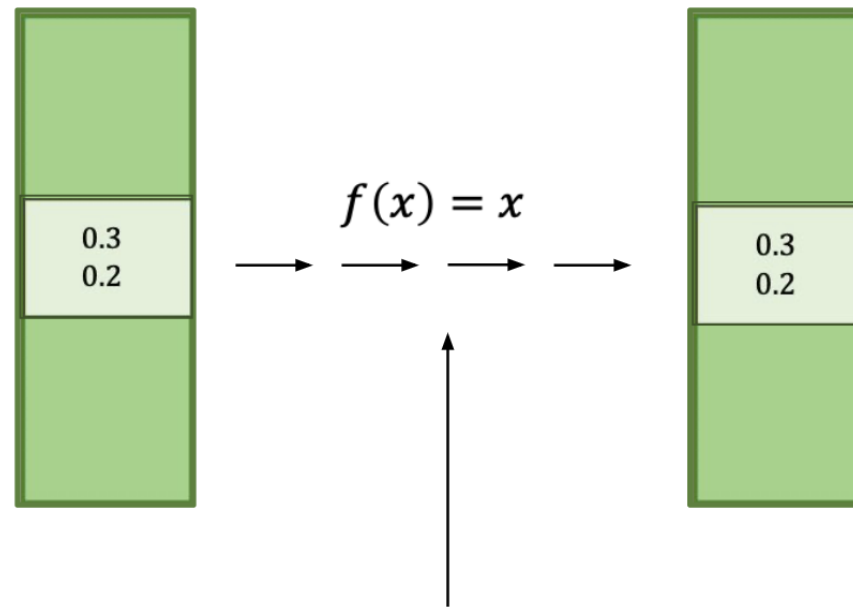Through all subsequent RNN steps, we want *"dog"* to stay the same

# An Illustrative Example:

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



RNN steps

# An Illustrative Example:

...to preserve *"dog"* , we need to compute the ***identity function*** over the part of the vector that stores it

But will that happen?

No

Why?

$$f(x) = x$$

0.3
0.2

0.3
0.2

RNN steps

# How does this affect the hidden state?

RNN update

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

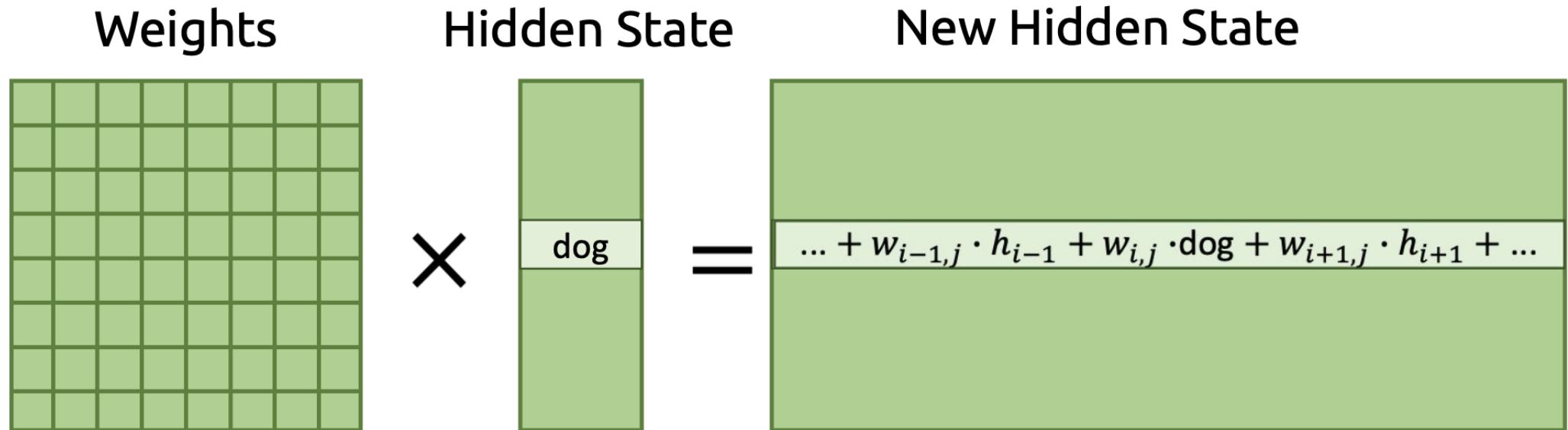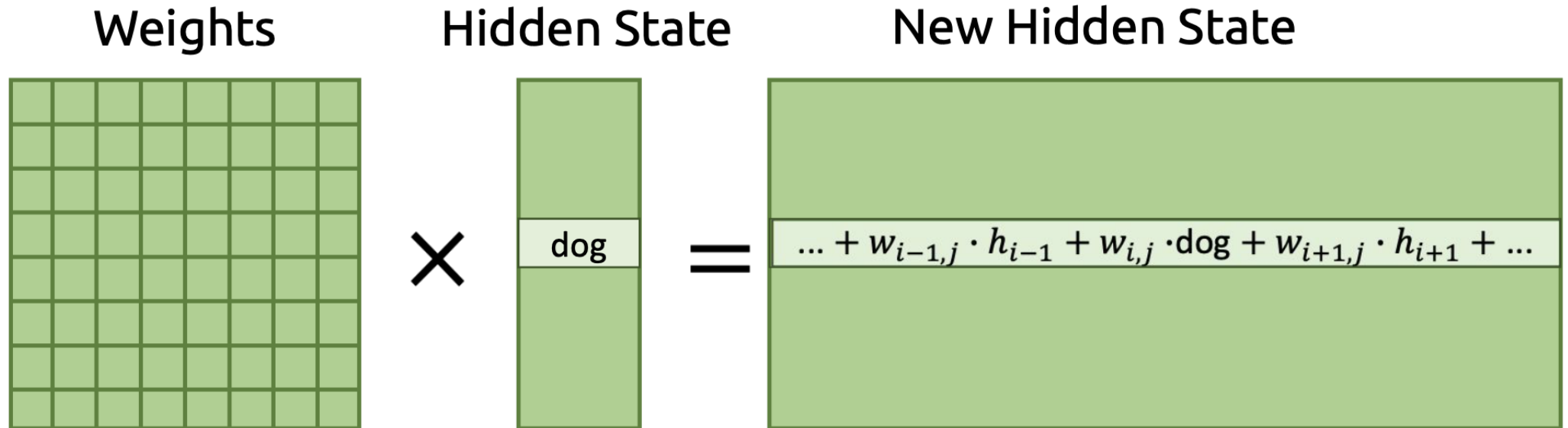The hidden state goes through a fully connected layer!

# How does this affect the hidden state?

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

# How does this affect the hidden state?

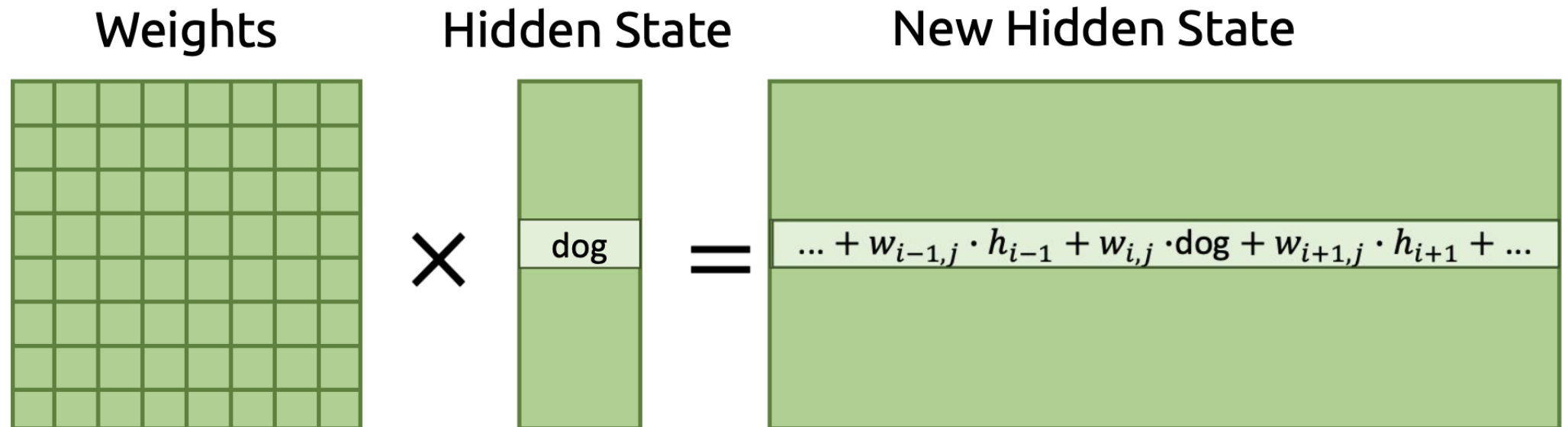- What will happen to our dog after we multiply our weights by our hidden state?



Weights $\times$ Hidden State $=$ New Hidden State

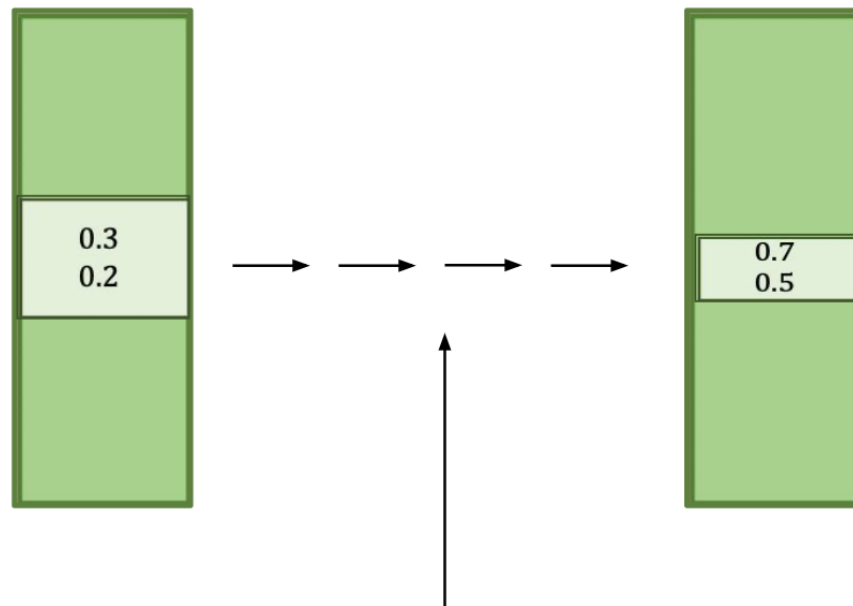$$\dots + w_{i-1,j} \cdot h_{i-1} + w_{i,j} \cdot \text{dog} + w_{i+1,j} \cdot h_{i+1} + \dots$$

# How does this affect the hidden state?

Dog gets lost in all the other information!



Weights $\times$ Hidden State = New Hidden State

$$... + w_{i-1,j} \cdot h_{i-1} + w_{i,j} \cdot dog + w_{i+1,j} \cdot h_{i+1} + ...$$

# How does this affect the hidden state?

Dog gets lost in all the other information!



Weights     Hidden State     New Hidden State

$$\ldots + w_{i-1,j} \cdot h_{i-1} + w_{i,j} \cdot \text{dog} + w_{i+1,j} \cdot h_{i+1} + \ldots$$

# How does this affect the hidden state?

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



RNN steps

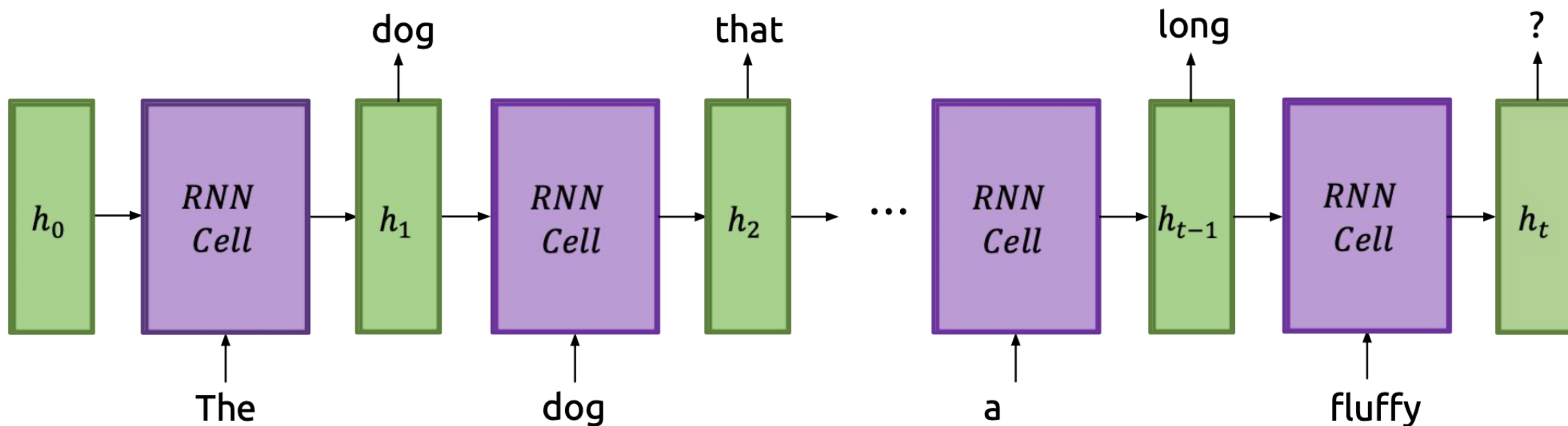RNN forgets about the dog after a certain time ☹
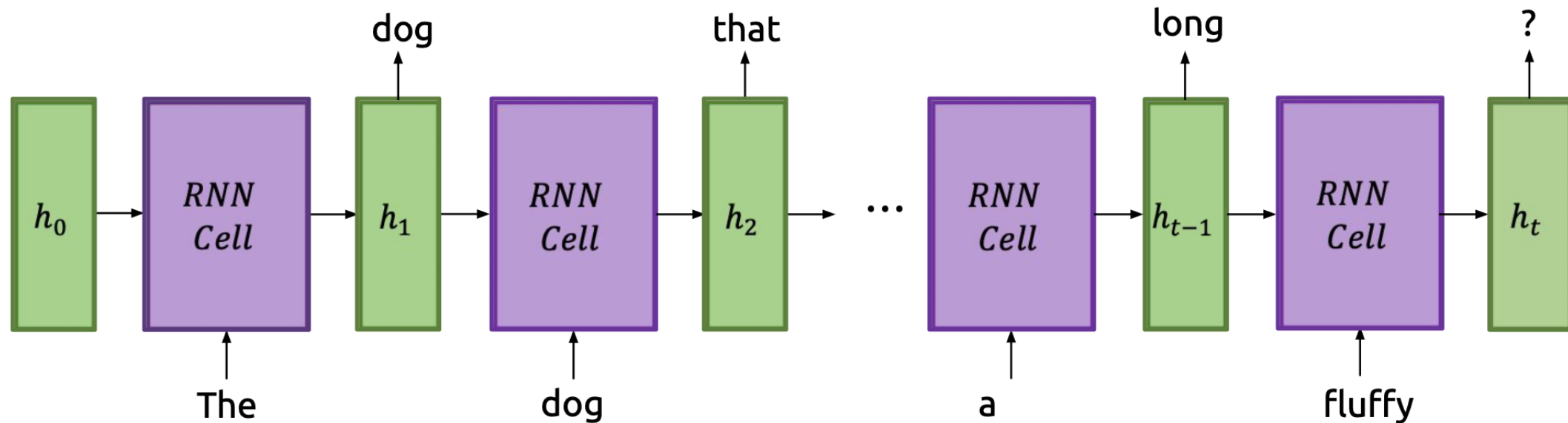
# RNNs cannot learn "long term" dependency



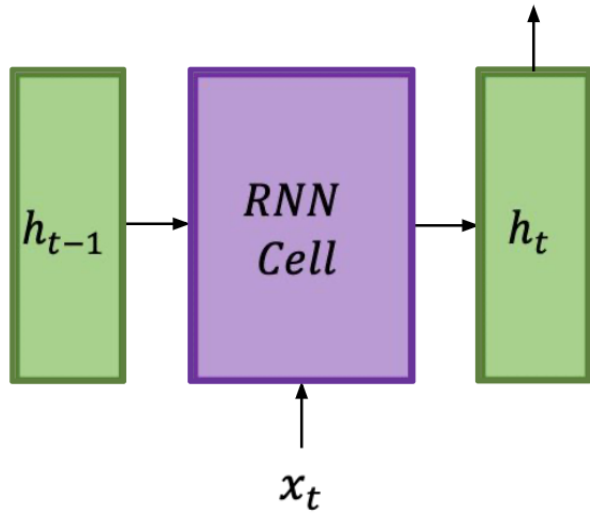We need new way to update hidden state!

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

- **_Long_** Short Term Memory (LSTM)
  - "Short-term memory that persists over time"
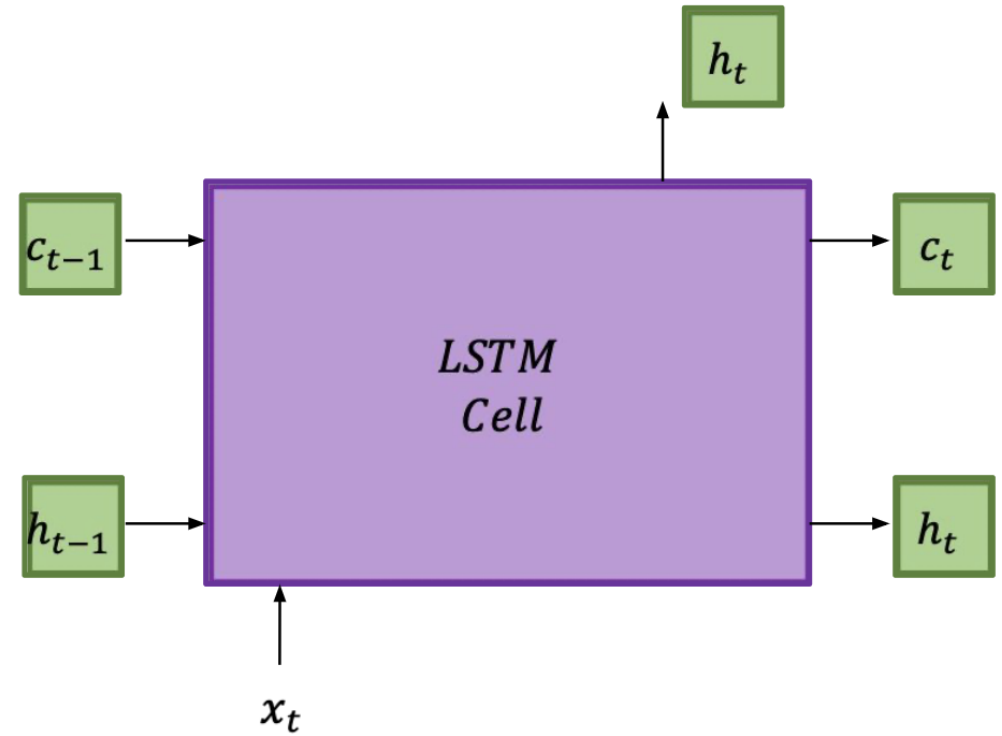  - i.e. "hidden states that remember information for longer"
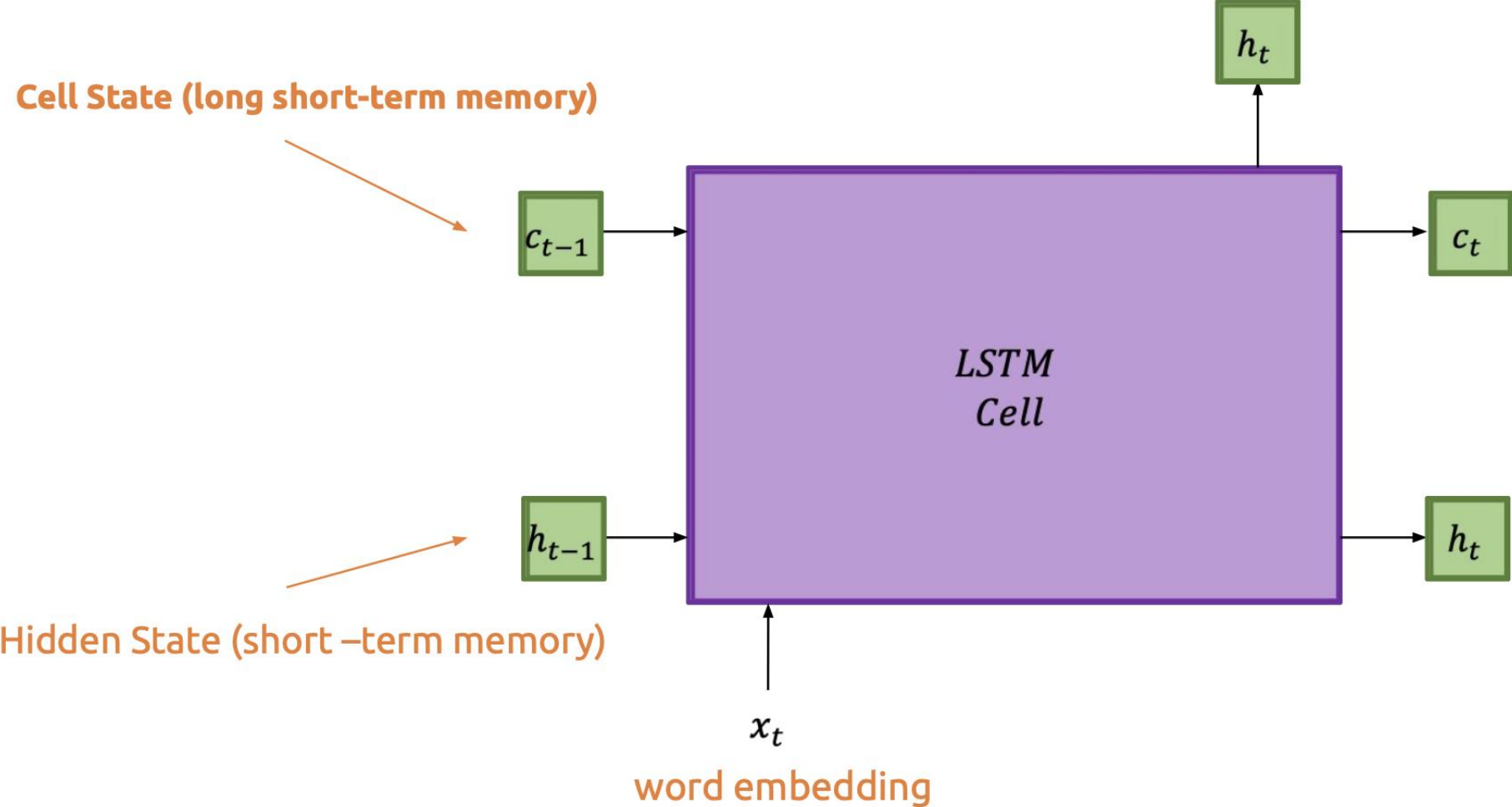
**What is different?**

**Vanilla RNN**



**LSTM**

# LSTM

Cell State (long short-term memory)

Hidden State (short –term memory)



$h_t$

$c_{t-1}$

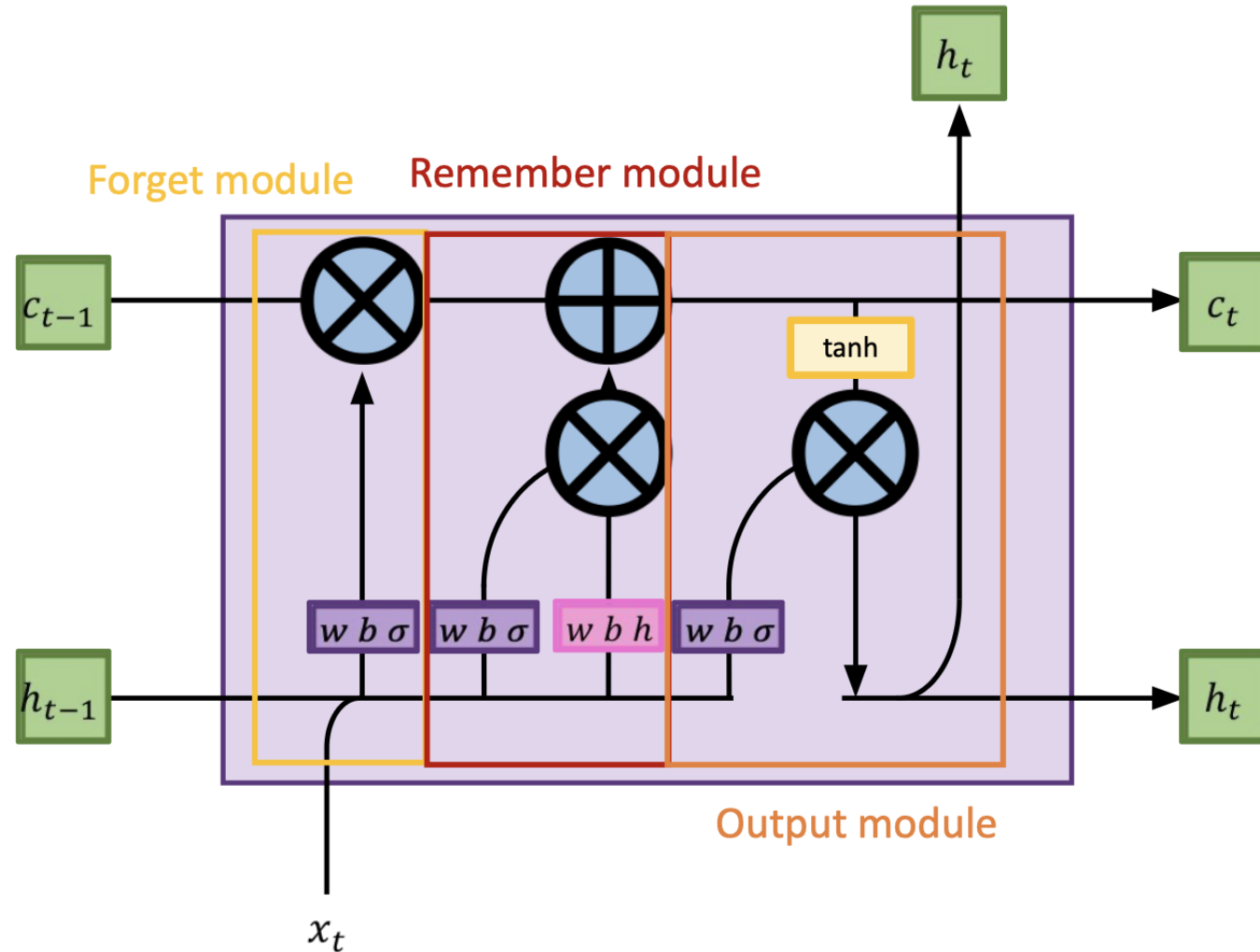$h_{t-1}$

LSTM
Cell

$c_t$

$h_t$

$x_t$

word embedding

# How an LSTM works

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

# The Complete LSTM

# Forget Module

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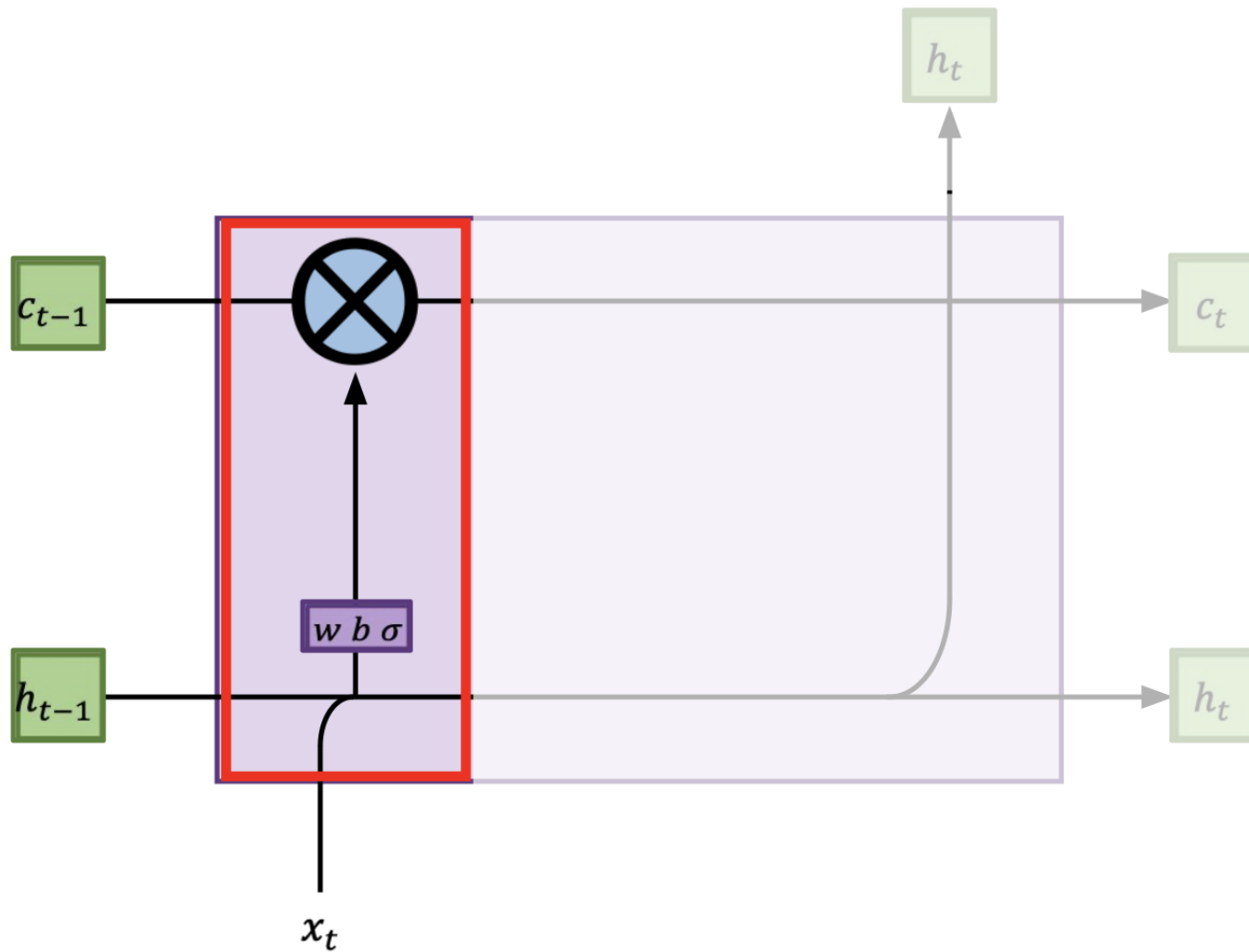


Next set of words: *"I love my dog"*

# Forget Module

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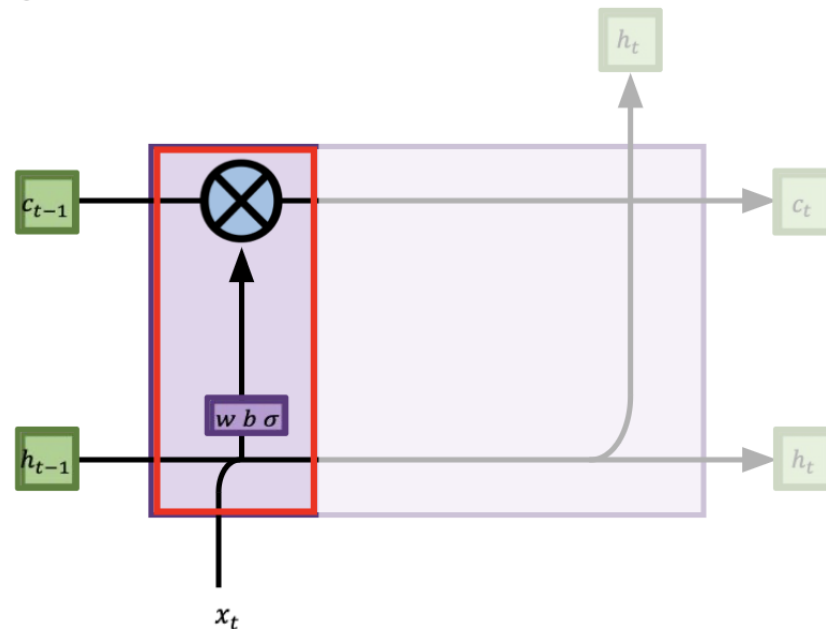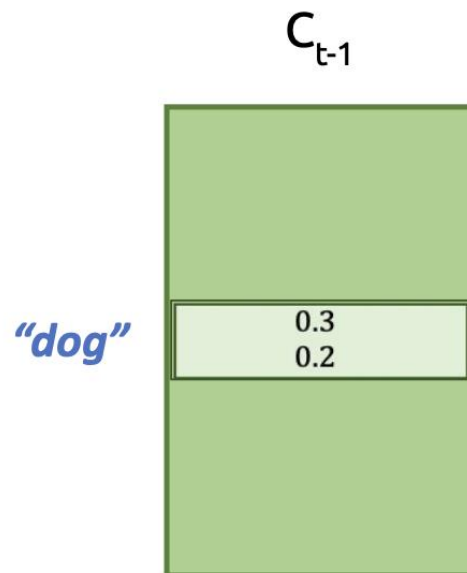
# Forget Module

# Forget Module

- 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
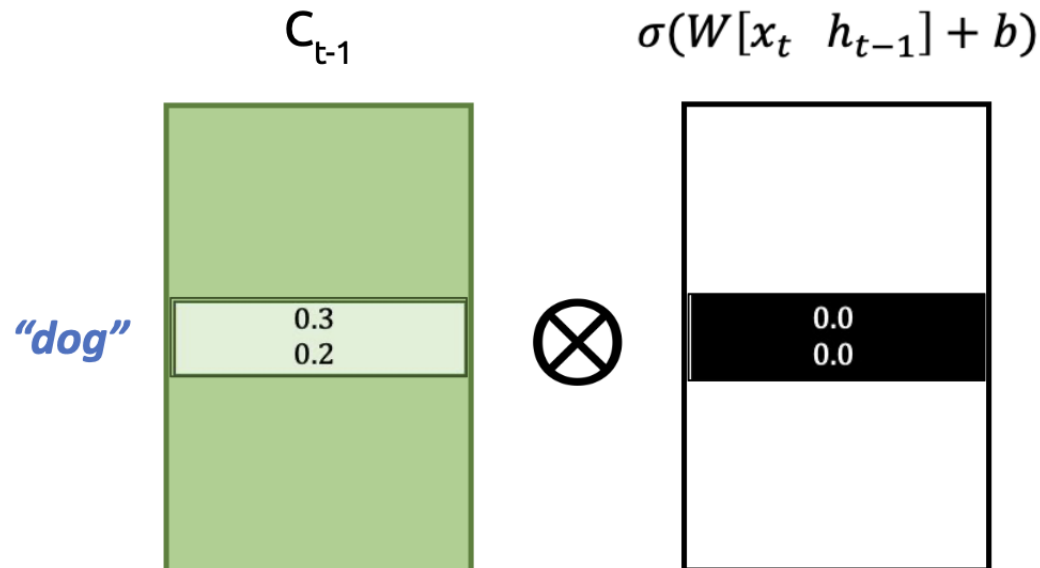
- Either lets parts of $C_{t-1}$ pass through or not

# Forgetting information

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

$C_{t-1}$

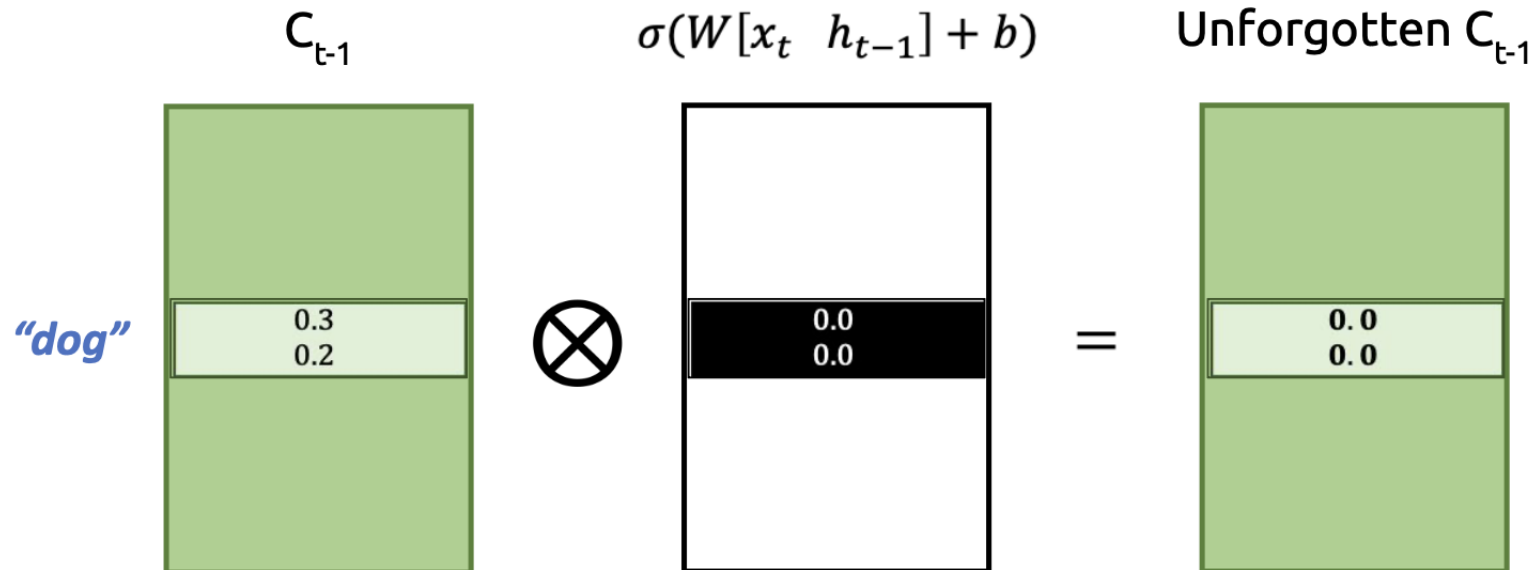**"dog"**

| |
|---|
| 0.3 |
| 0.2 |

33

# 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

$$C_{t-1} \qquad\qquad \sigma(W[x_t \quad h_{t-1}] + b)$$

*"dog"*

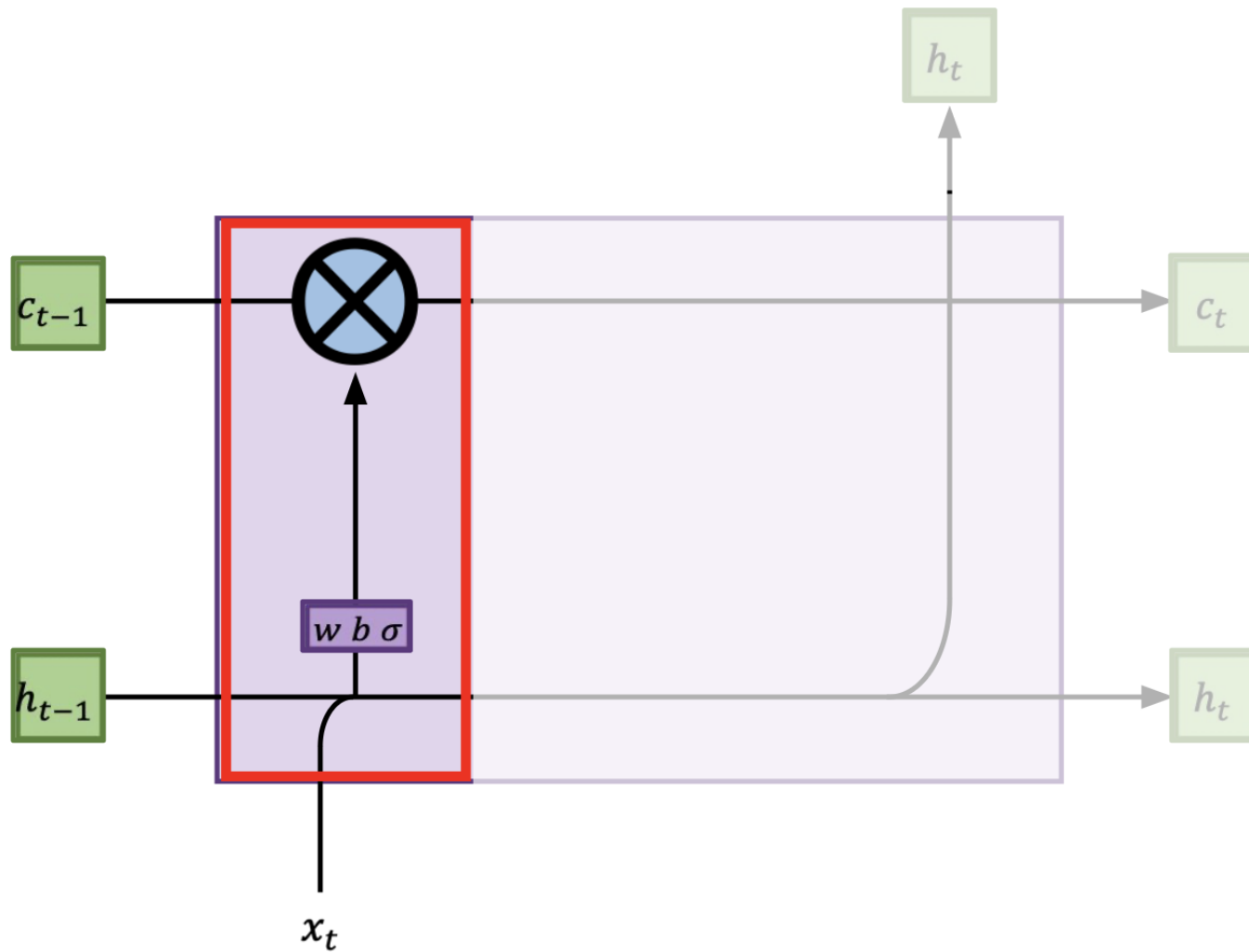| | |
|---|---|
| 0.3 | 0.0 |
| 0.2 | 0.0 |

# 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
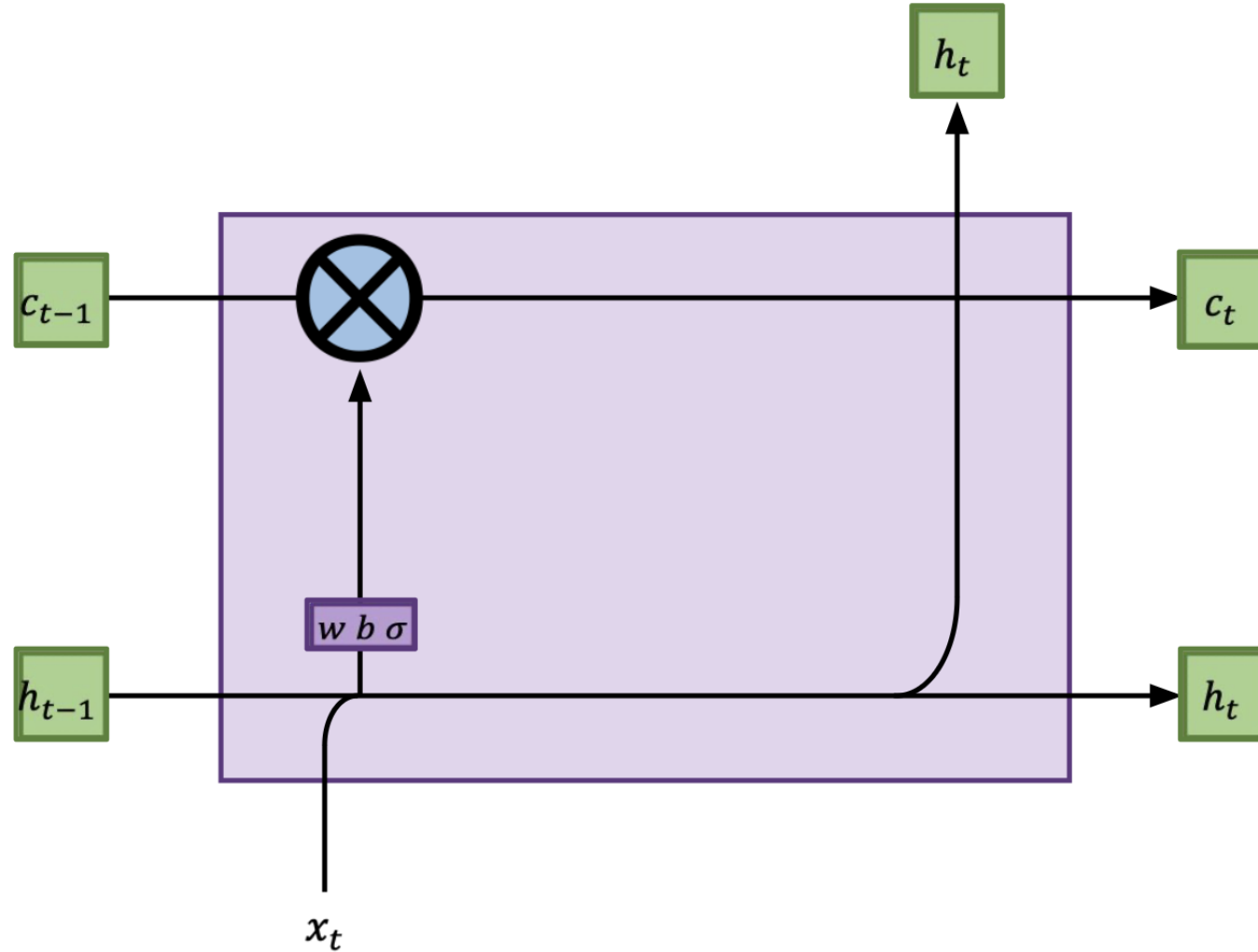  - Pointwise multiplication by a learned mask vector is known as *gating*

$C_{t-1}$ $\qquad\qquad$ $\sigma(W[x_t \quad h_{t-1}] + b)$ $\qquad$ Unforgotten $C_{t-1}$



*"dog"* | 0.3 0.2 | $\otimes$ | 0.0 0.0 | $=$ | 0.0 0.0

# Forget Module
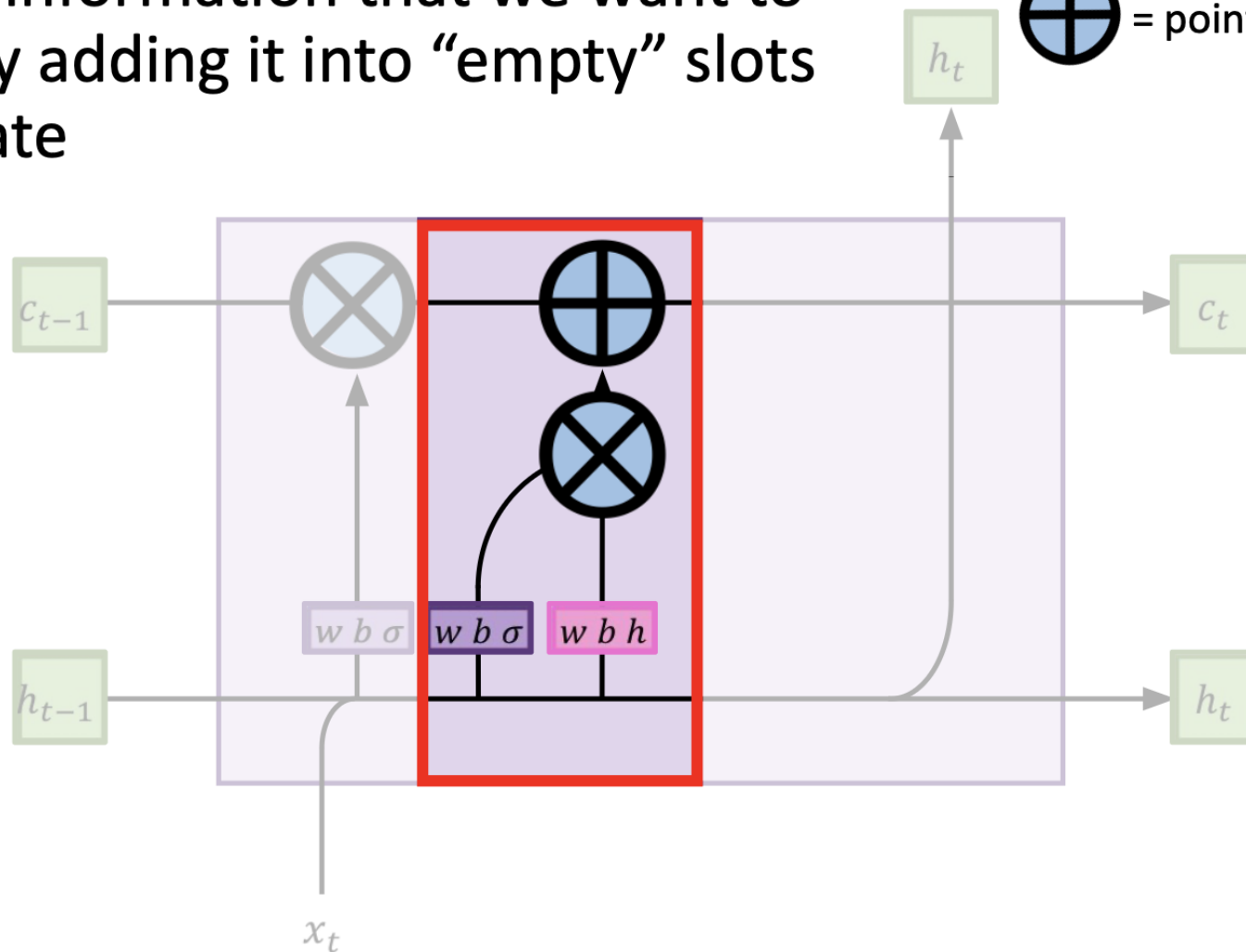
# Remember Module

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



$w\ b\ \sigma$ = fully connected layer with sigmoid

$w\ b\ h$ = fully connected layer with tanh

$\otimes$ = pointwise multiplication
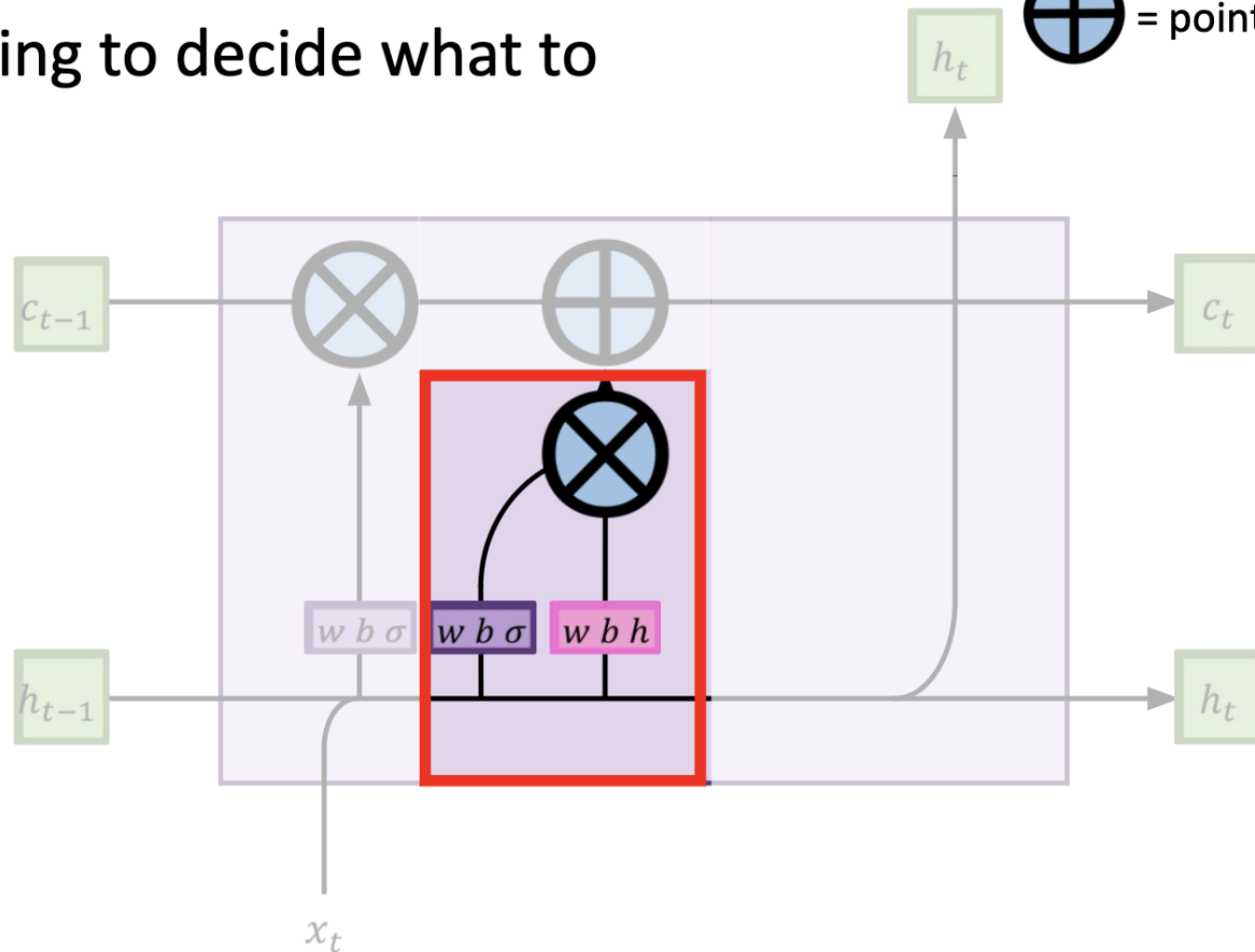
$\oplus$ = pointwise addition

# Remember Module

- First: use gating to decide what to remember

# Gating for 'selective memory'

- A fully-connected + tanh on [input, memory] computes some new memory
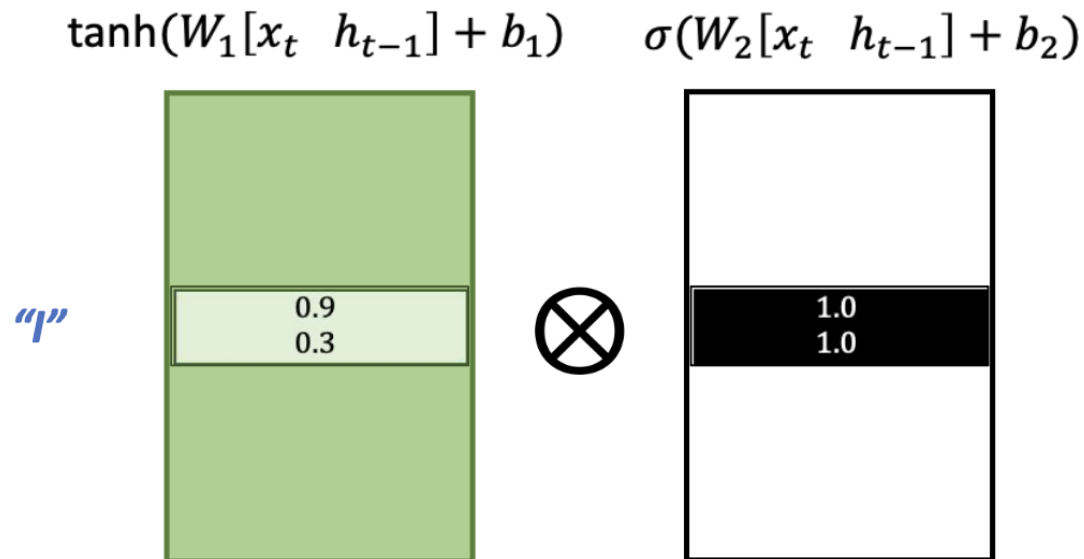
$$\tanh(W_1[x_t \quad h_{t-1}] + b_1)$$

*"I"*

| |
|---|
| 0.9 |
| 0.3 |

# 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

$$\tanh(W_1[x_t \quad h_{t-1}] + b_1) \qquad \sigma(W_2[x_t \quad h_{t-1}] + b_2)$$



"I"   | 0.9 |  ⊗  | 1.0 |
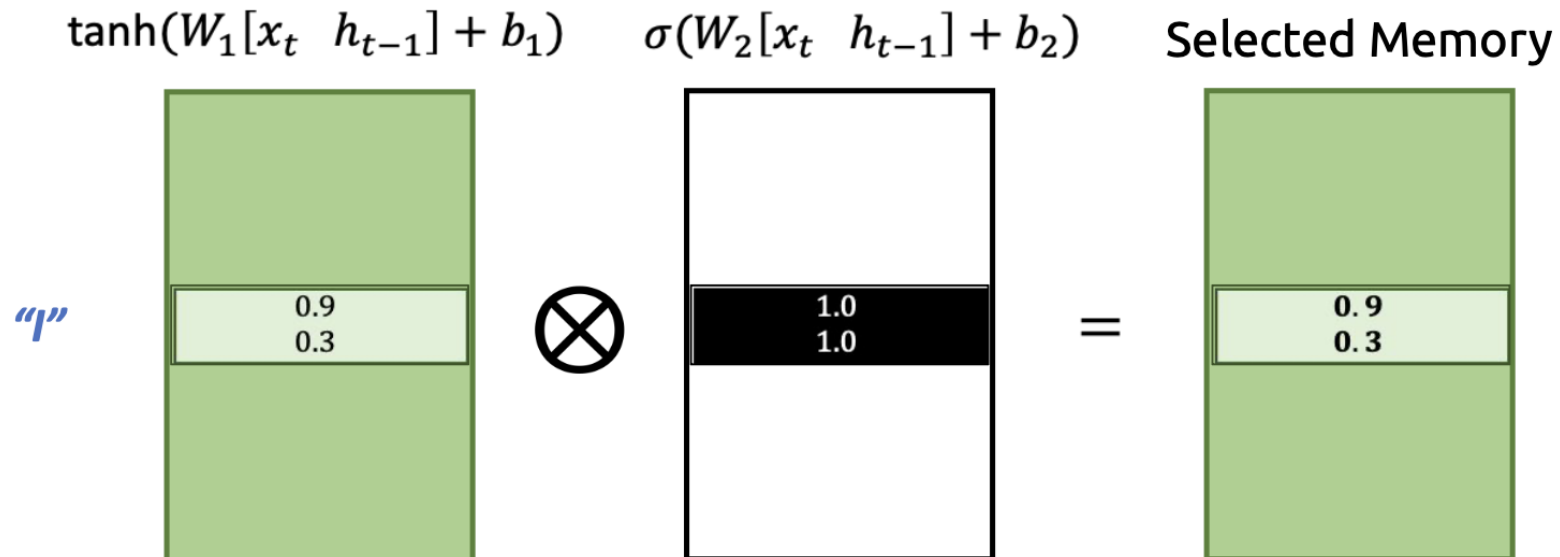      | 0.3 |     | 1.0 |

41

# 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

$$\tanh(W_1[x_t \quad h_{t-1}] + b_1) \qquad \sigma(W_2[x_t \quad h_{t-1}] + b_2) \qquad \text{Selected Memory}$$



"I"

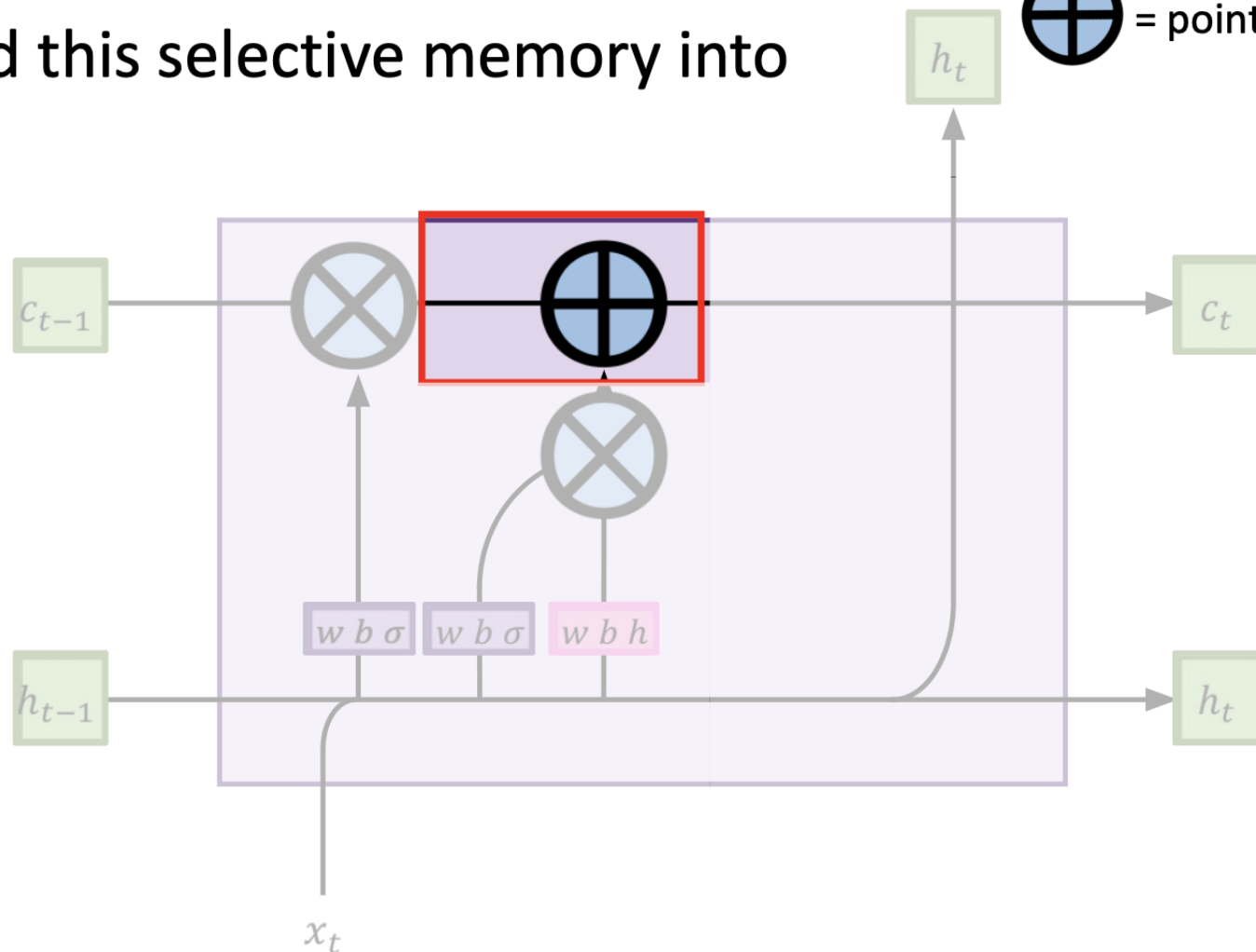| 0.9 | $\otimes$ | 1.0 | = | 0.9 |
| 0.3 | | 1.0 | | 0.3 |

# Remember Module

- Then: we add this selective memory into the cell state

$w\ b\ \sigma$ = fully connected layer with sigmoid

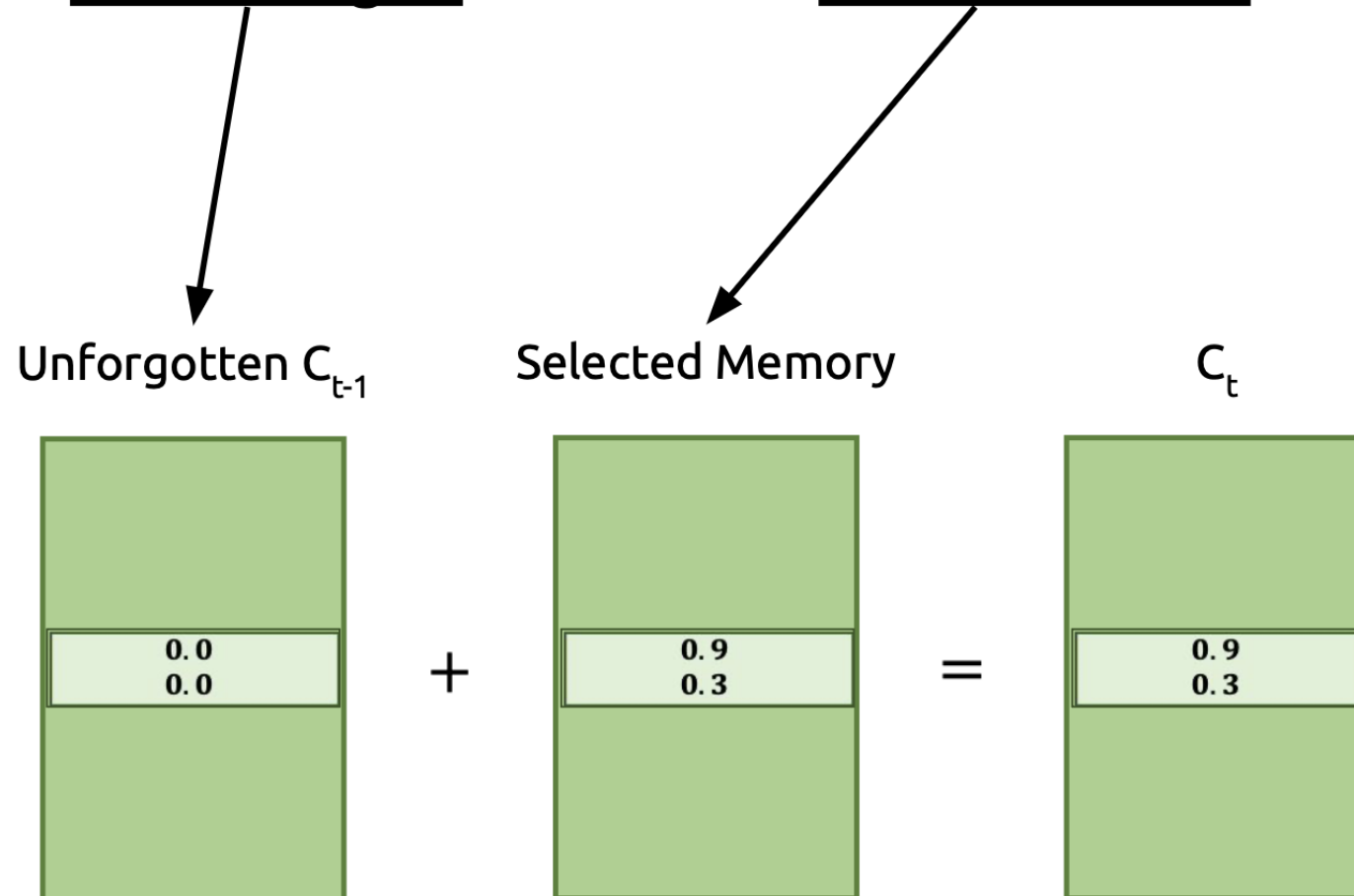$w\ b\ h$ = fully connected layer with tanh

⊗ = pointwise multiplication

⊕ = pointwise addition

$h_t$

$c_{t-1}$

$c_t$

$w\ b\ \sigma$ $w\ b\ \sigma$ $w\ b\ h$

$h_{t-1}$

$h_t$

$x_t$

44

# Remembering information

- Add what we <u>didn't forget</u> to what we <u>did remember</u>

Unforgotten $C_{t-1}$       Selected Memory       $C_t$

| 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!
  - Never has to mix up its own information

# Output Module



legend:
- $w\,b\,\sigma$ = fully connected layer with sigmoid
- $w\,b\,h$ = fully connected layer with tanh
- $\otimes$ = pointwise multiplication
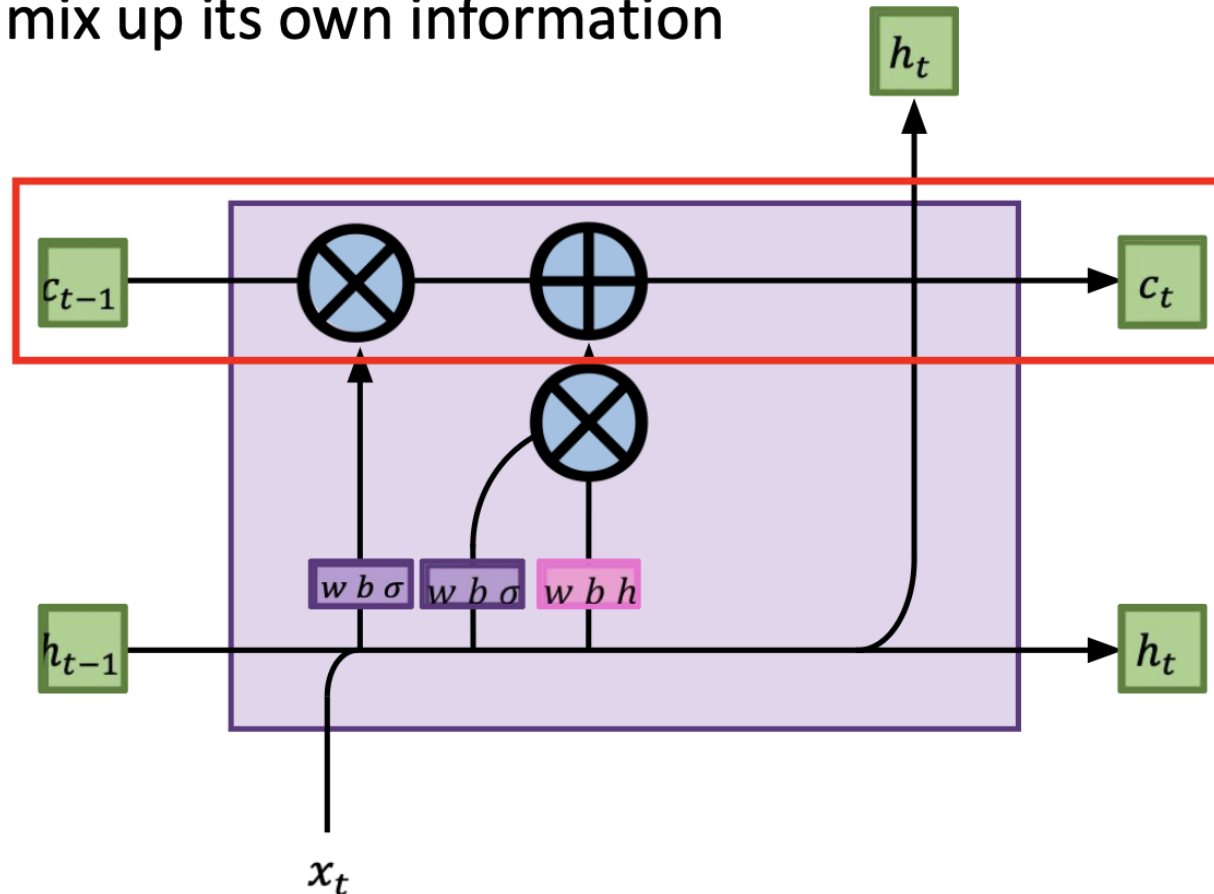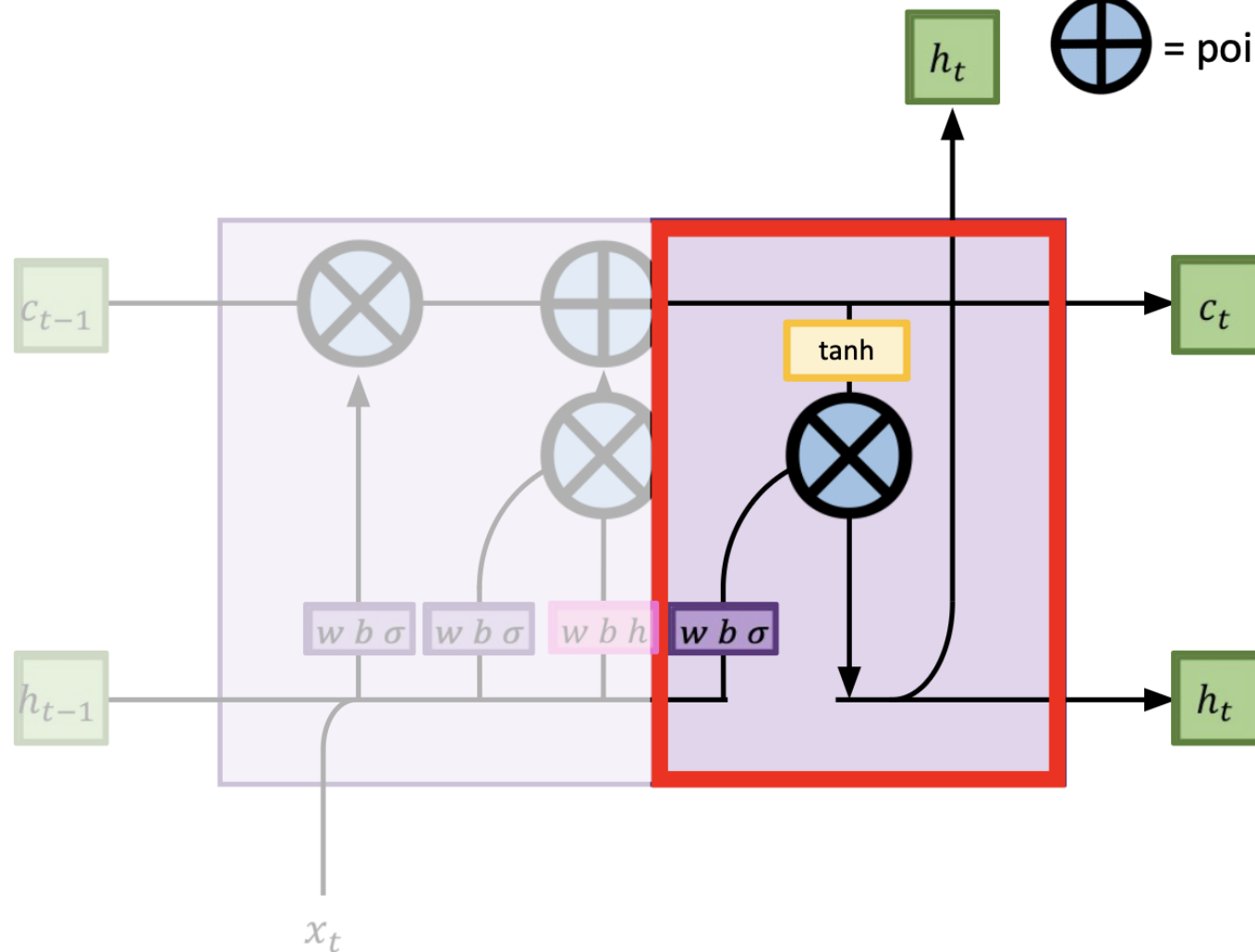- $\oplus$ = pointwise addition

$h_t$

tanh

$c_{t-1}$   $c_t$

$h_{t-1}$   $h_t$

$w\,b\,\sigma$   $w\,b\,\sigma$   $w\,b\,h$   $\boldsymbol{w\,b\,\sigma}$
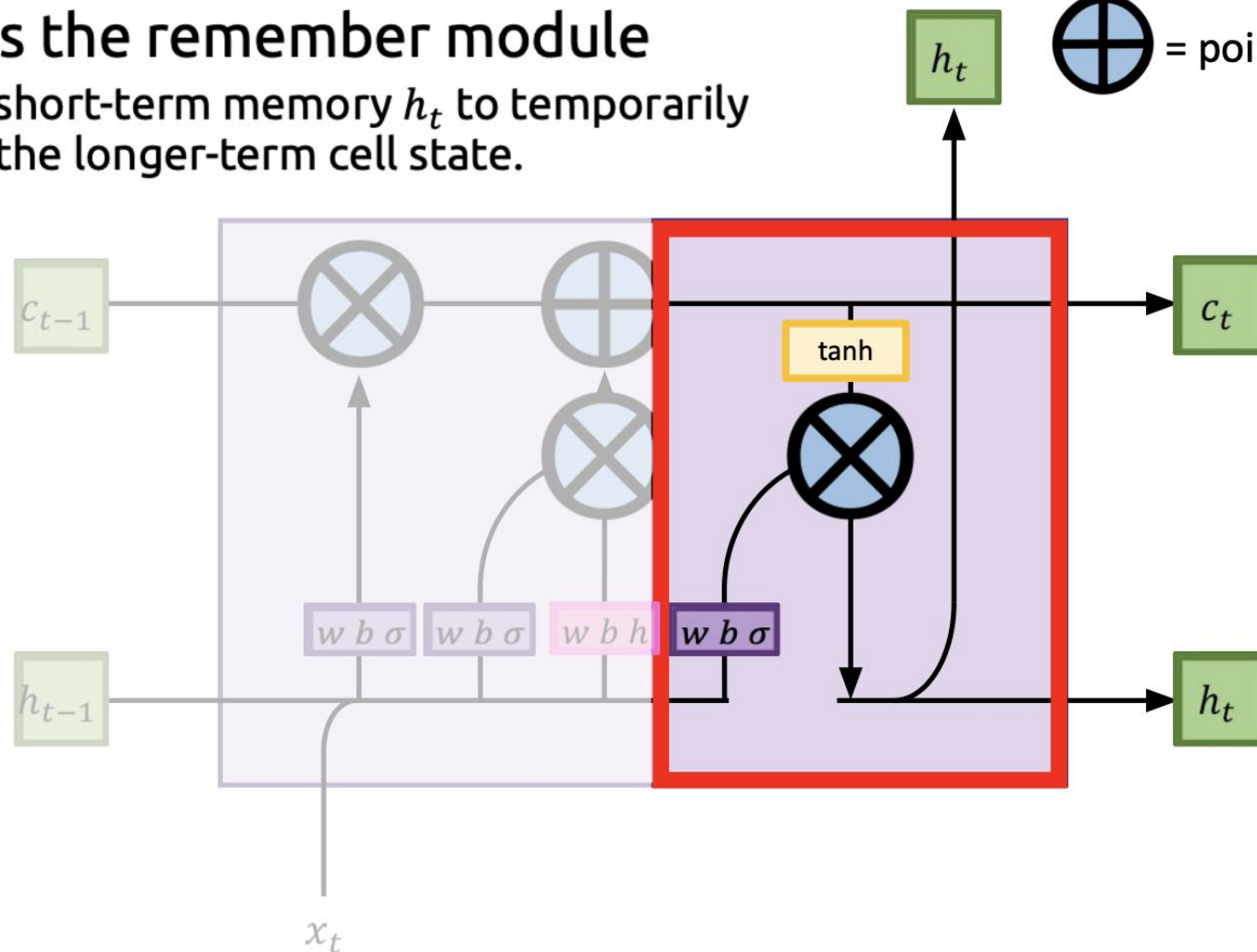
$x_t$

# 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.



$\boxed{w\ b\ \sigma}$ = fully connected layer with sigmoid

$\boxed{w\ b\ h}$ = fully connected layer with tanh

$\otimes$ = pointwise multiplication

$\oplus$ = pointwise addition

# The Complete LSTM

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i)$$
$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f)$$
$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o)$$
$$\widetilde{c}_t = tanh(W h_{t-1} + U x_t + b)$$
$$c_t = f_t \circ c_{t-1} + i_t \circ \widetilde{c}_t$$
$$h_t = o_t \circ tanh(c_t)$$
$$y_t = h_t$$