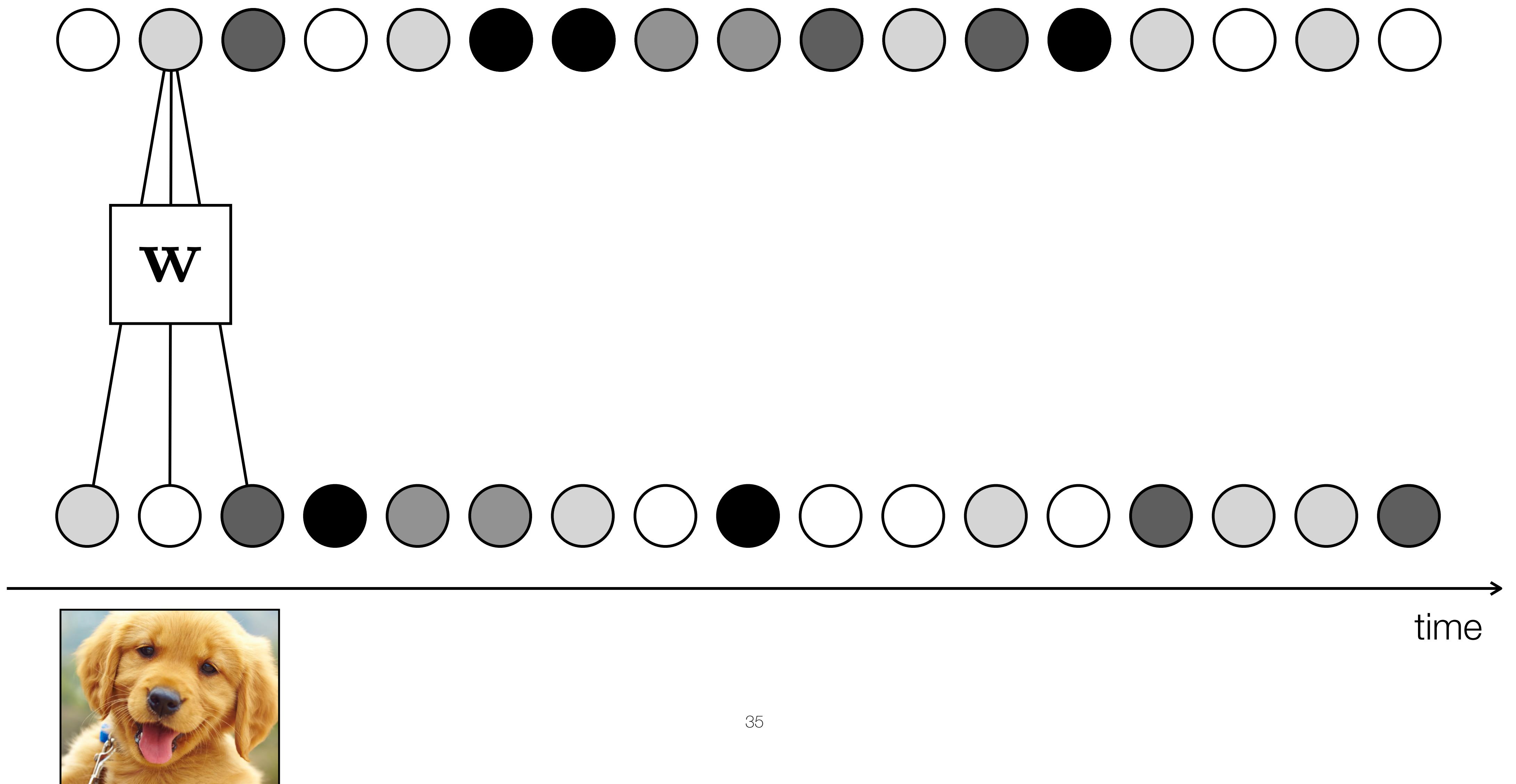
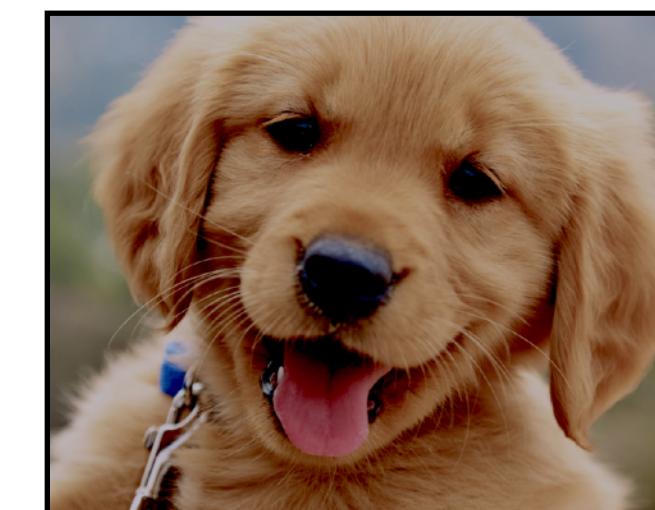
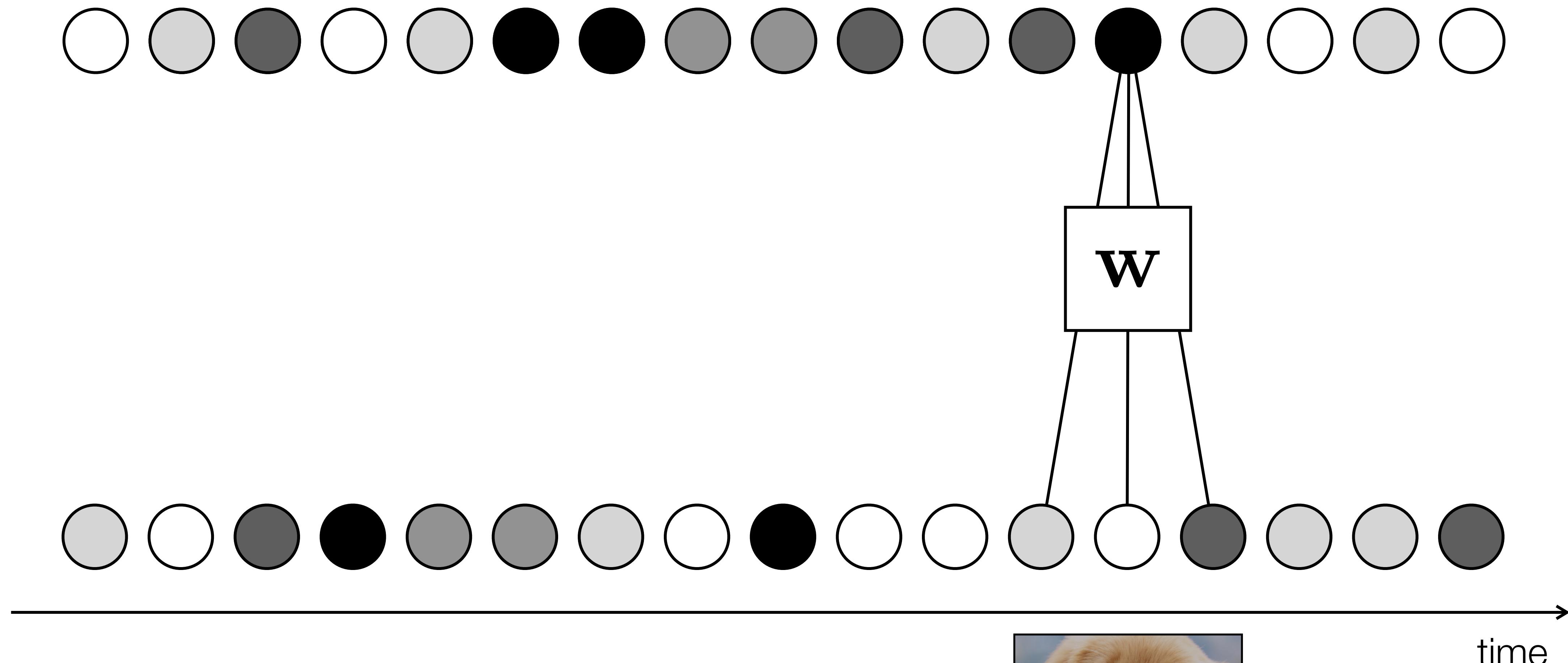


Lecture 12: Recurrent Models & Language

Rufus

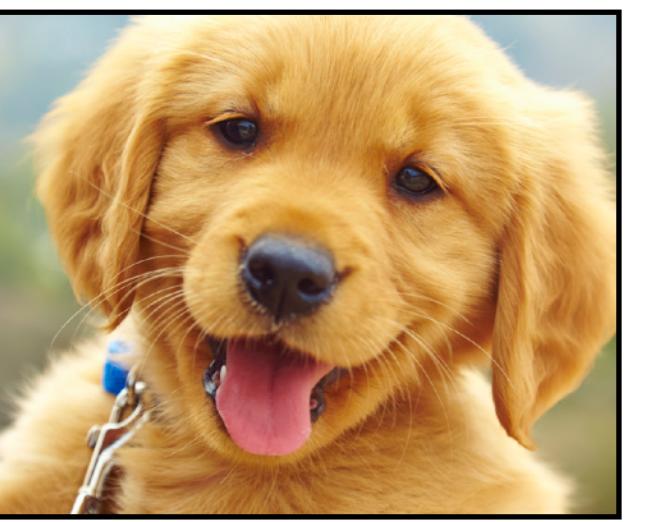
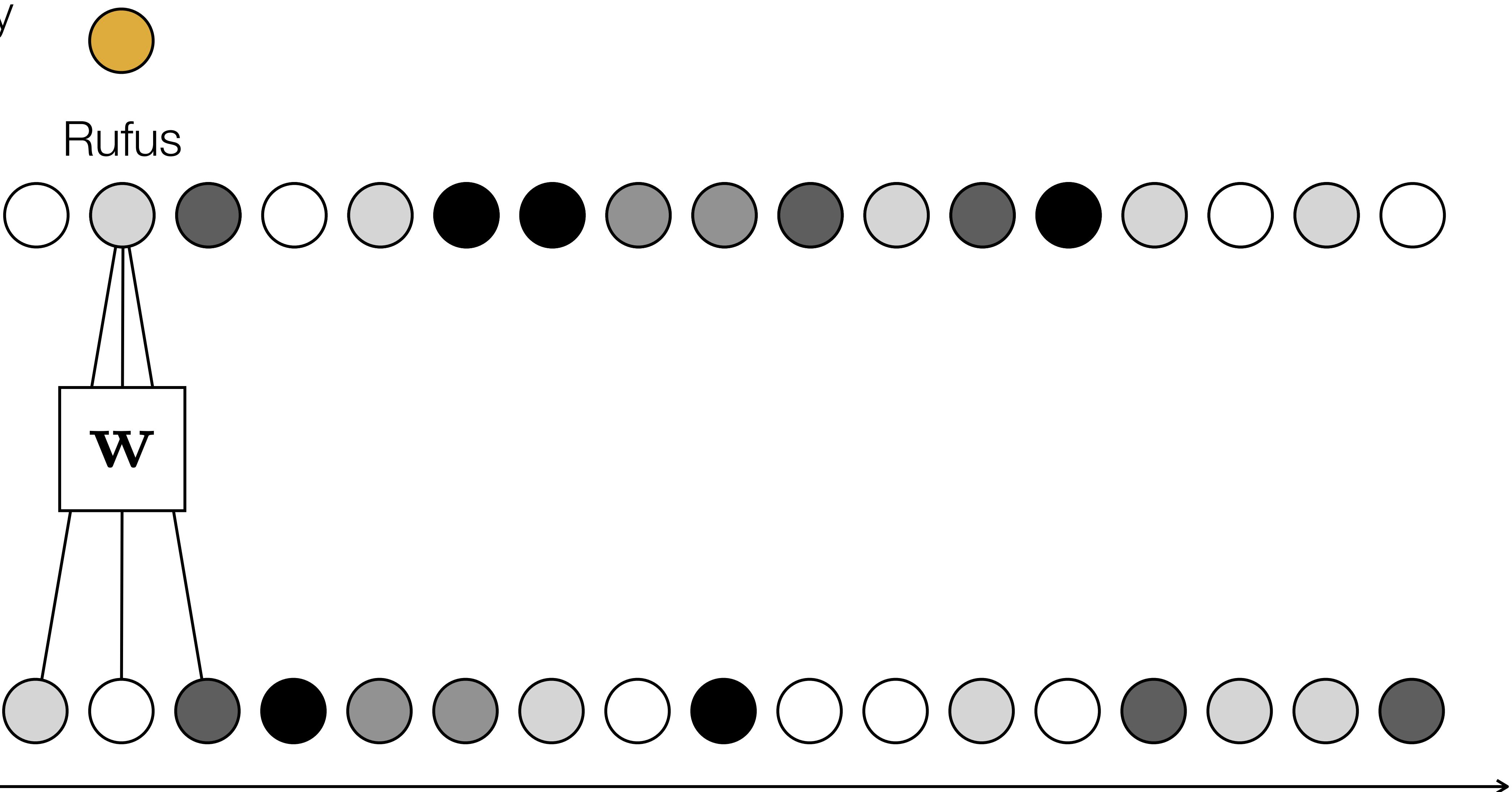


Douglas

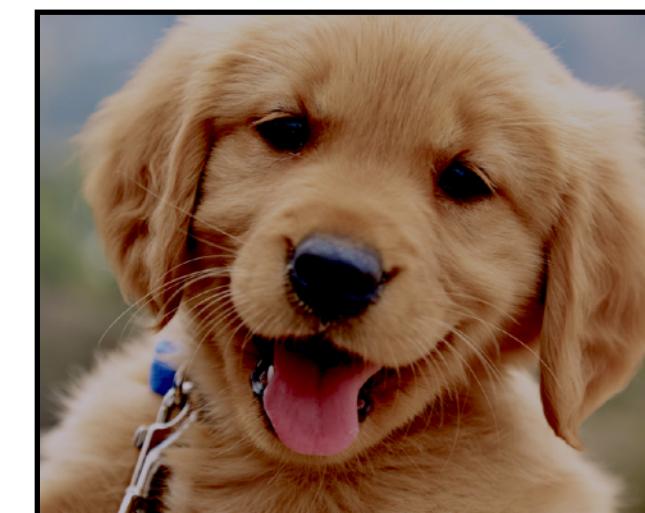
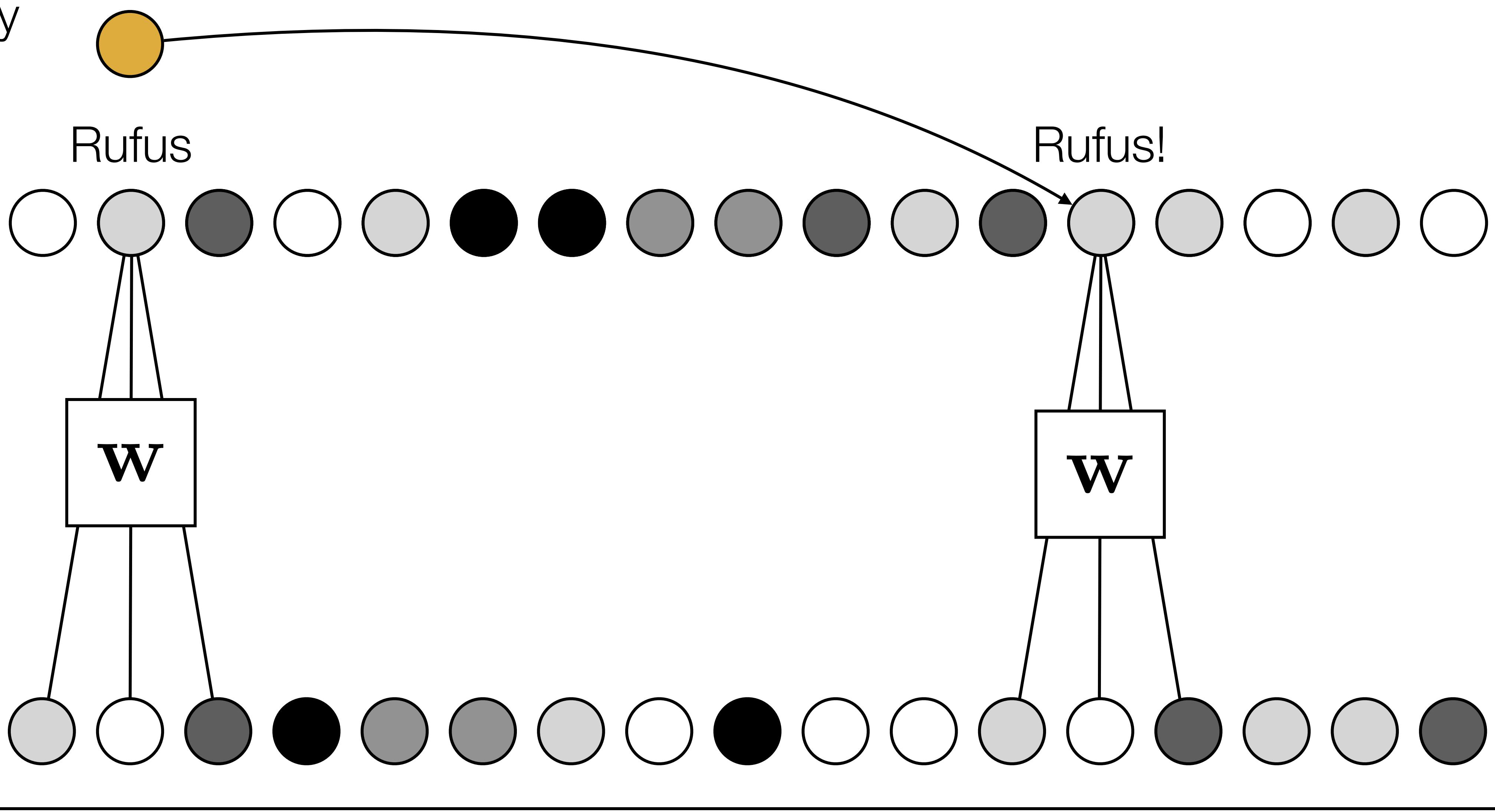


Memory
unit

Rufus



Memory
unit

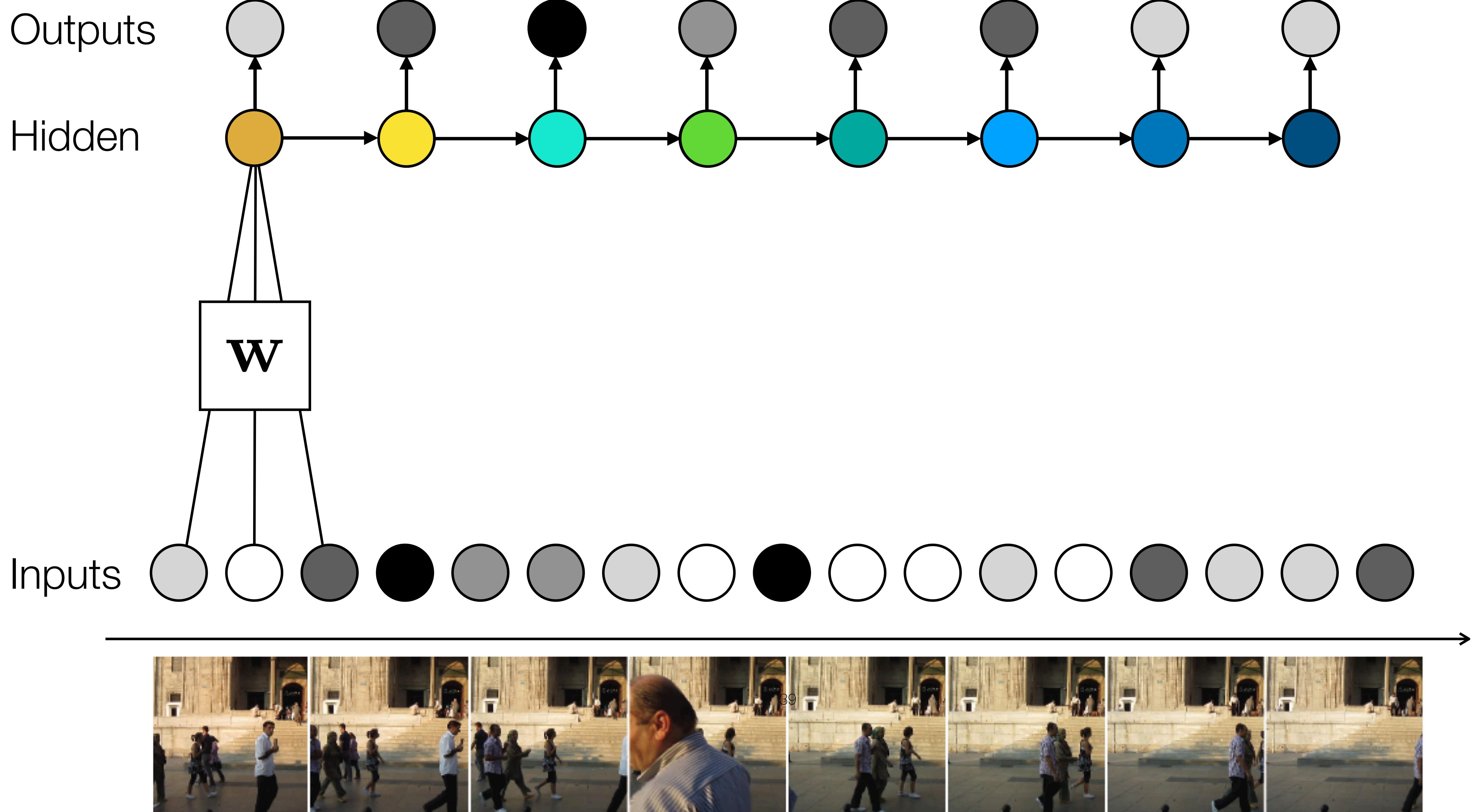


38

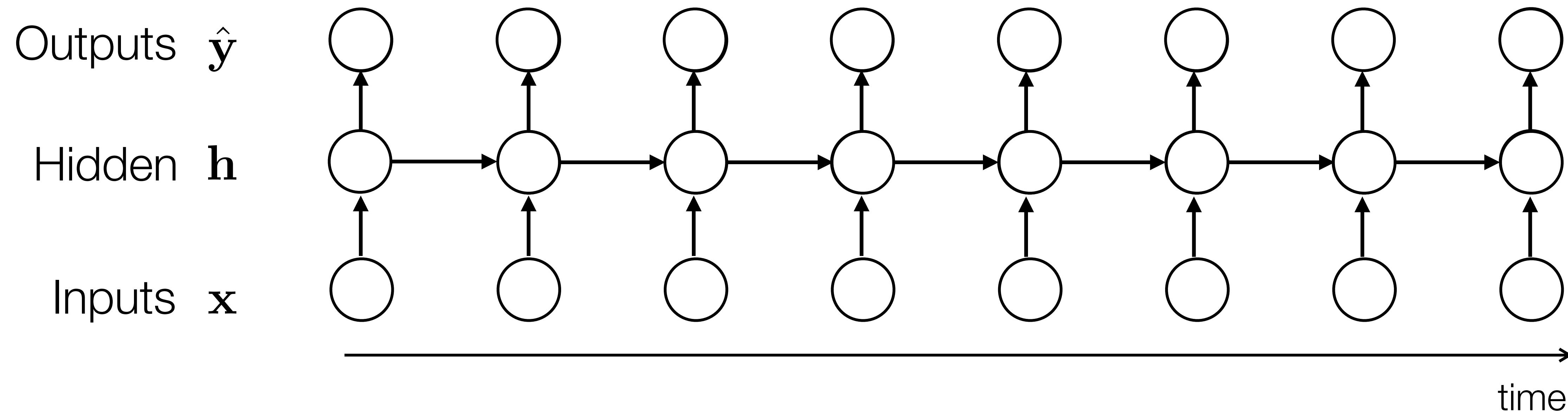
time

Source: Torralba, Freeman, Isola

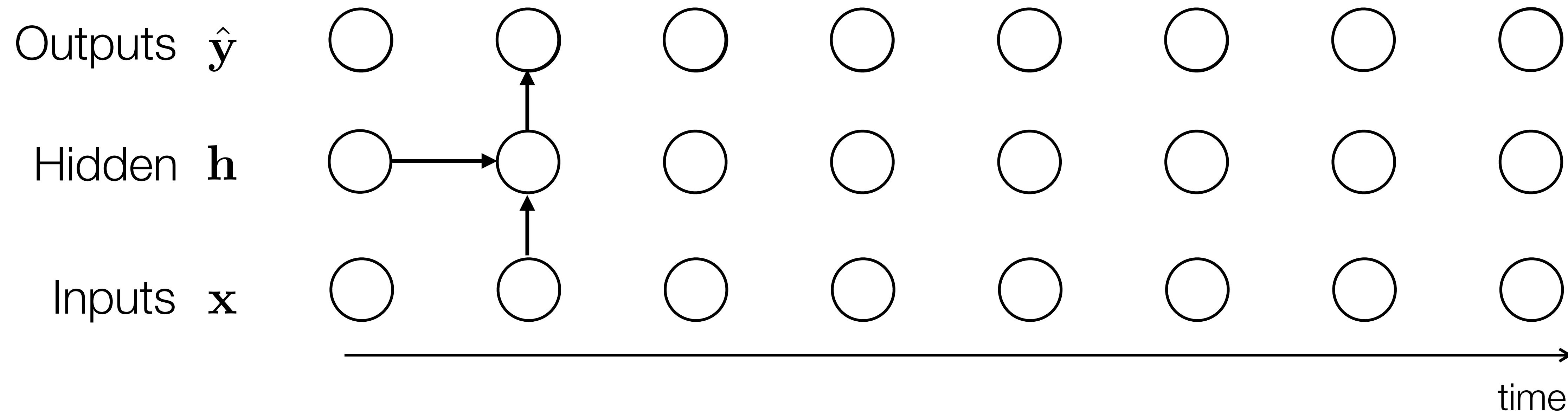
Recurrent Neural Networks (RNNs)



Recurrent Neural Networks (RNNs)



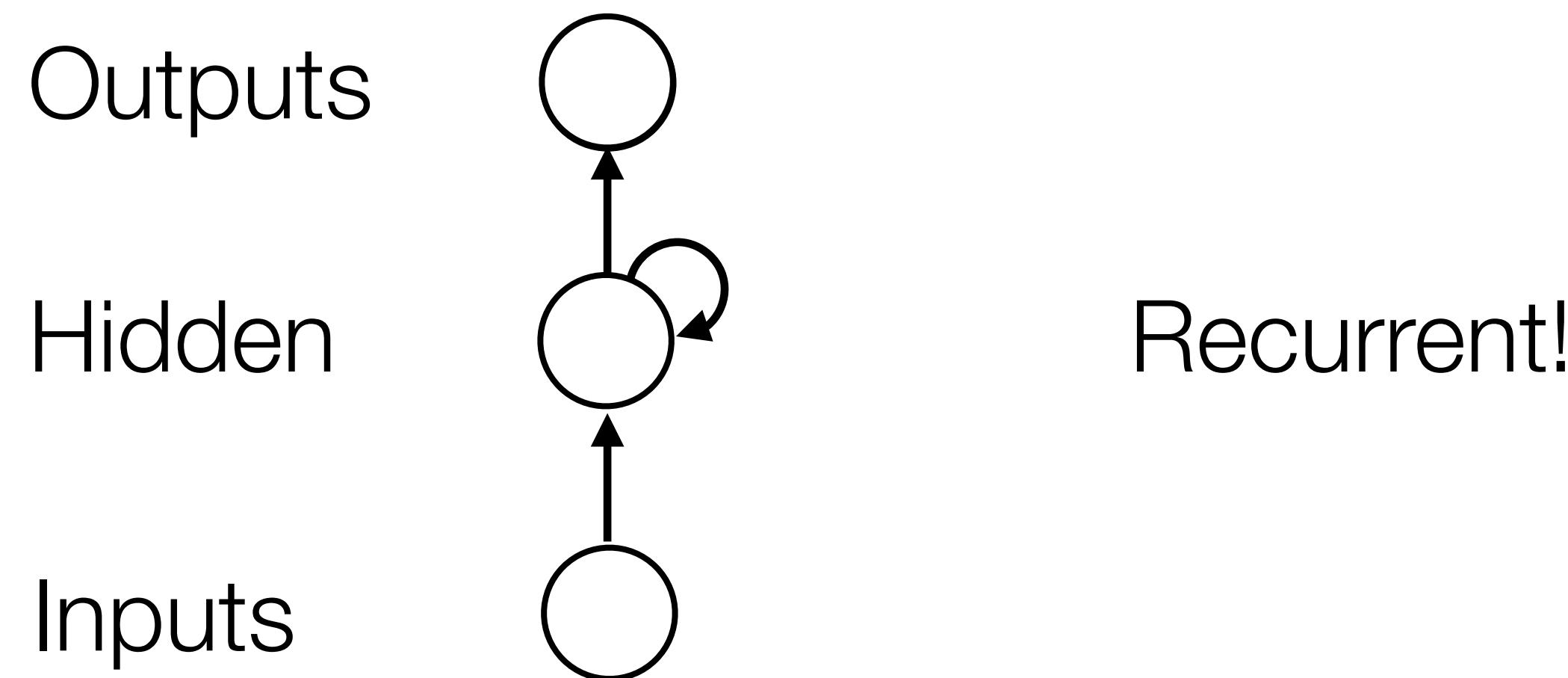
Recurrent Neural Networks (RNNs)



$$\mathbf{h}^{(t)} = f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)})$$

$$\hat{\mathbf{y}}^{(t)} = g(\mathbf{h}^{(t)})$$

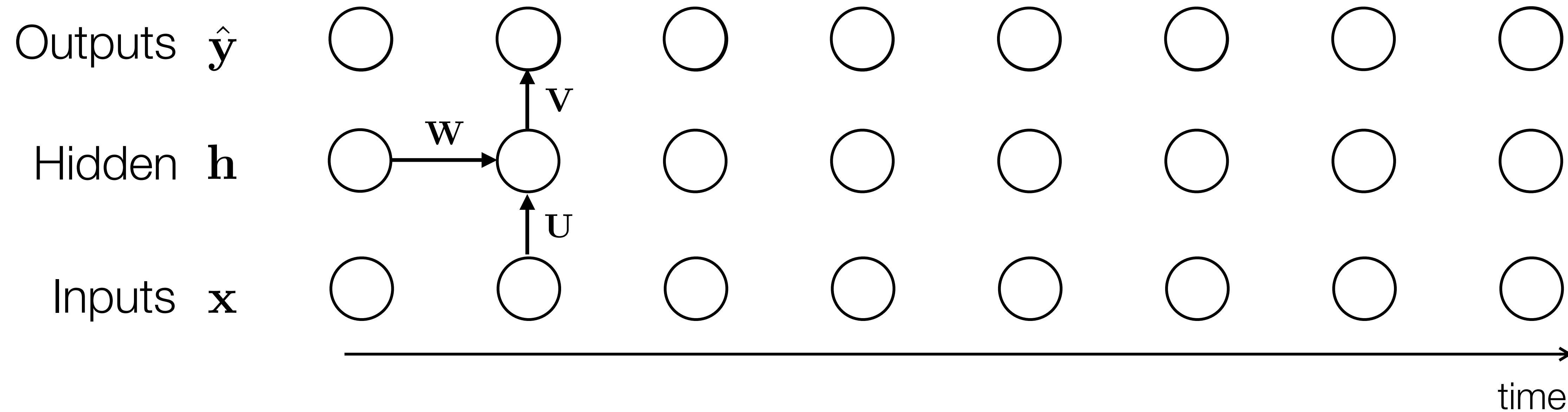
Recurrent Neural Networks (RNNs)



$$\mathbf{h}^{(t)} = f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)})$$

$$\mathbf{y}^{(t)} = g(\mathbf{h}^{(t)})$$

Recurrent Neural Networks (RNNs)



$$\mathbf{a}^{(t)} = \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)} + \mathbf{b}$$

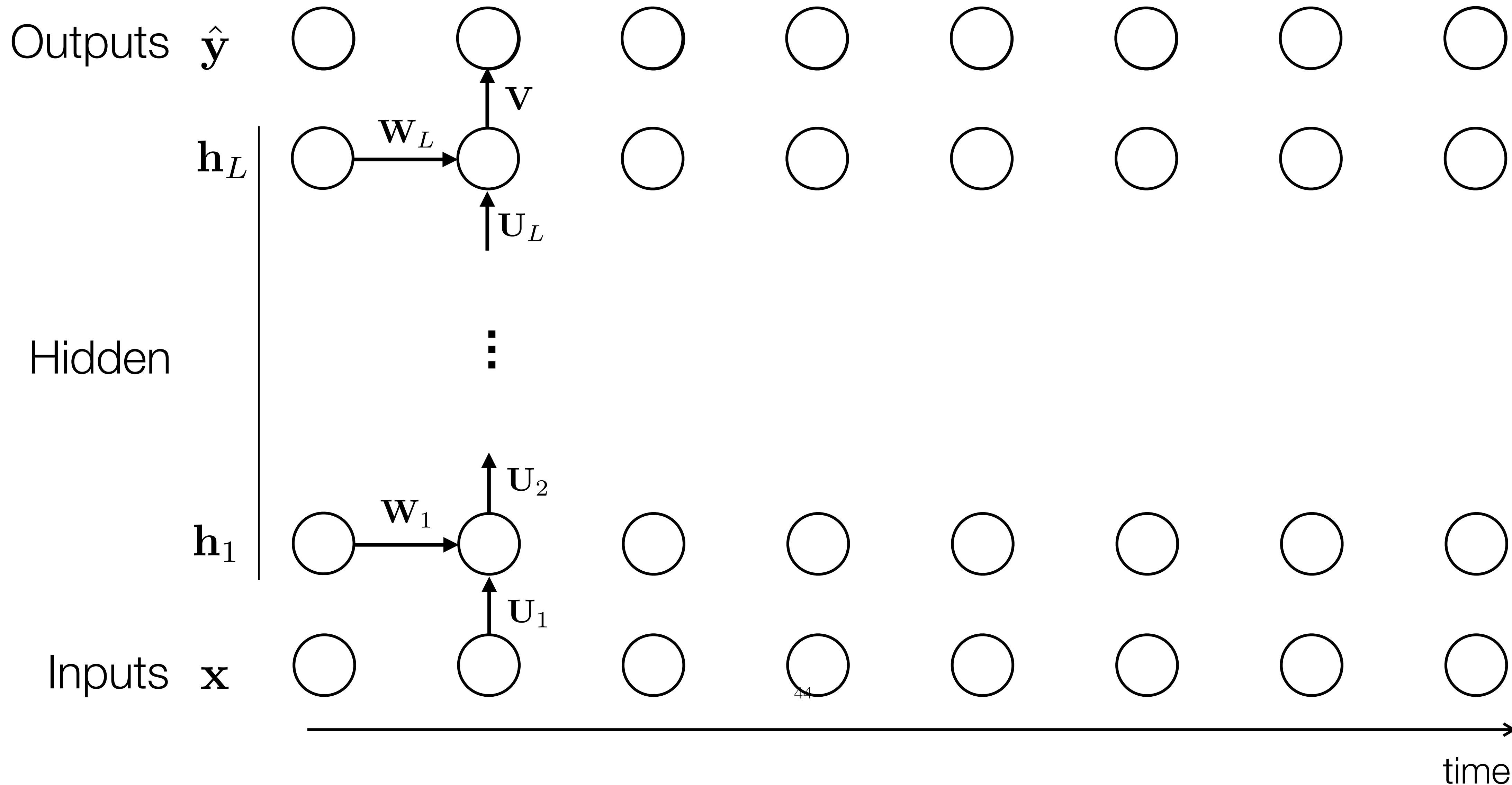
$$\mathbf{h}^{(t)} = \tanh(\mathbf{a}^{(t)})$$

$$\mathbf{o}^{(t)} = \mathbf{V}\mathbf{h}^{(t)} + \mathbf{c}$$

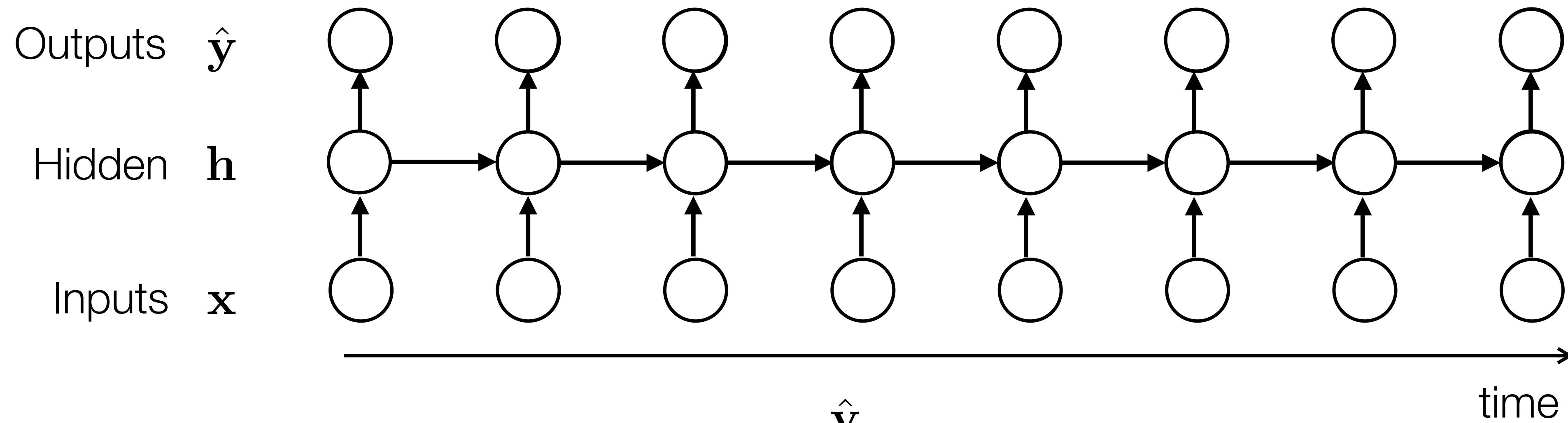
43

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{o}^{(t)})$$

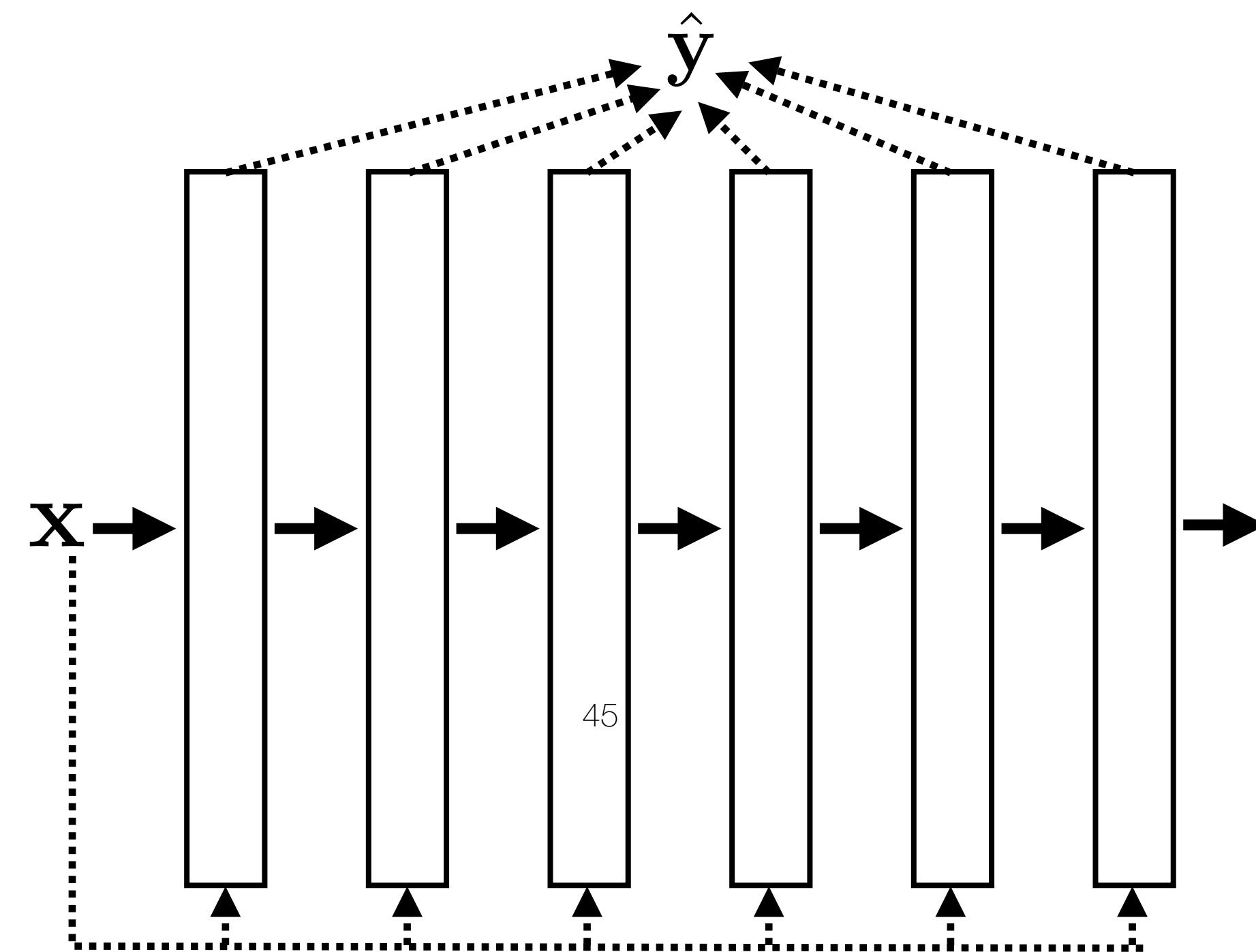
Deep Recurrent Neural Networks (RNNs)



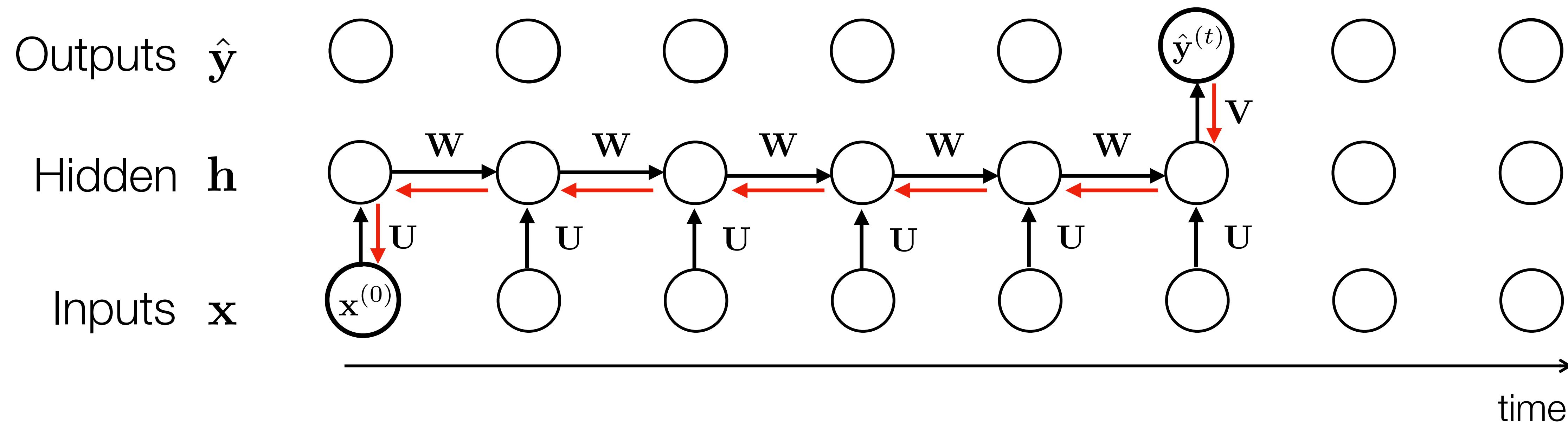
Unrolling an RNN



Equivalent to
"unrolled"
network with
lots of layers

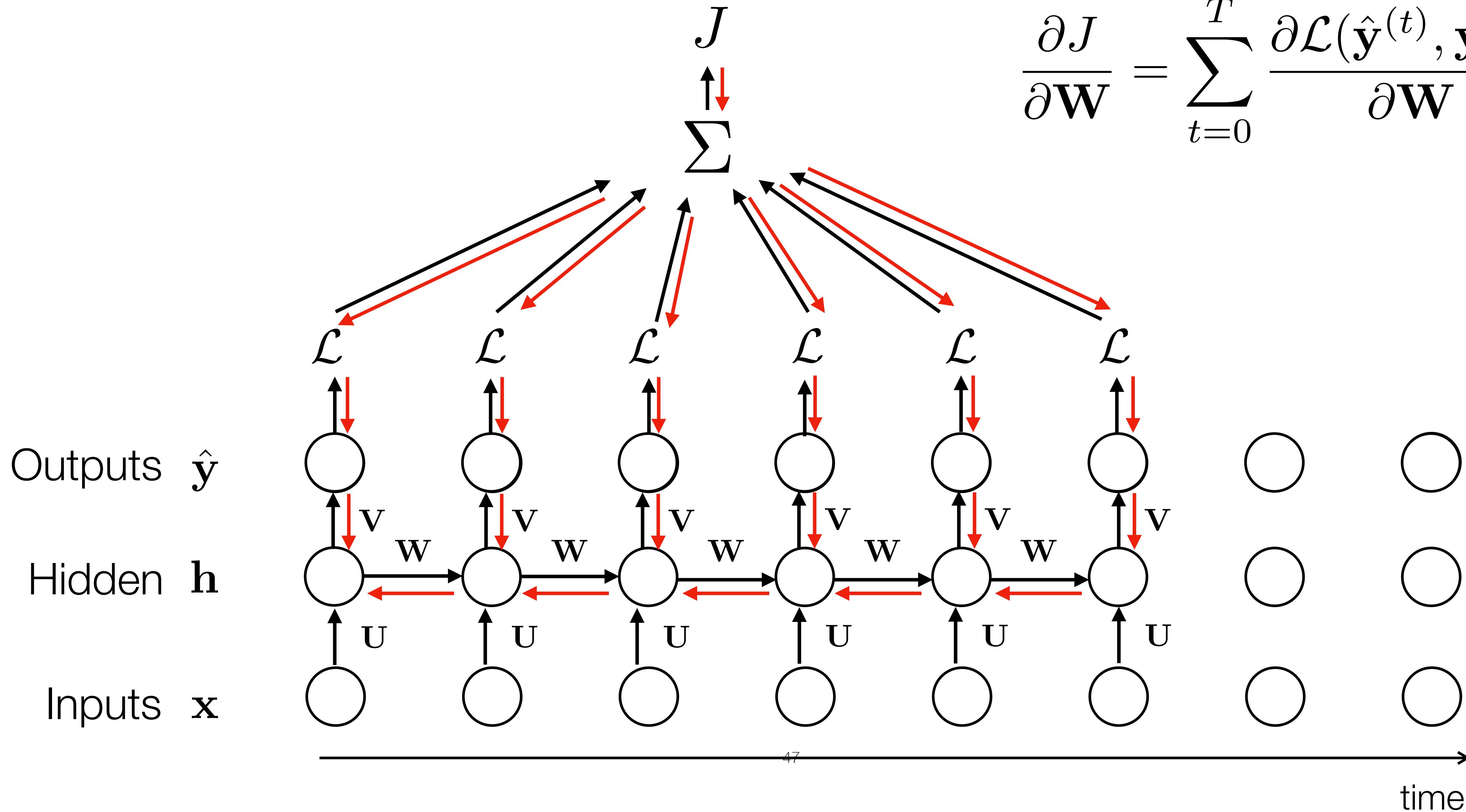


Backprop through time

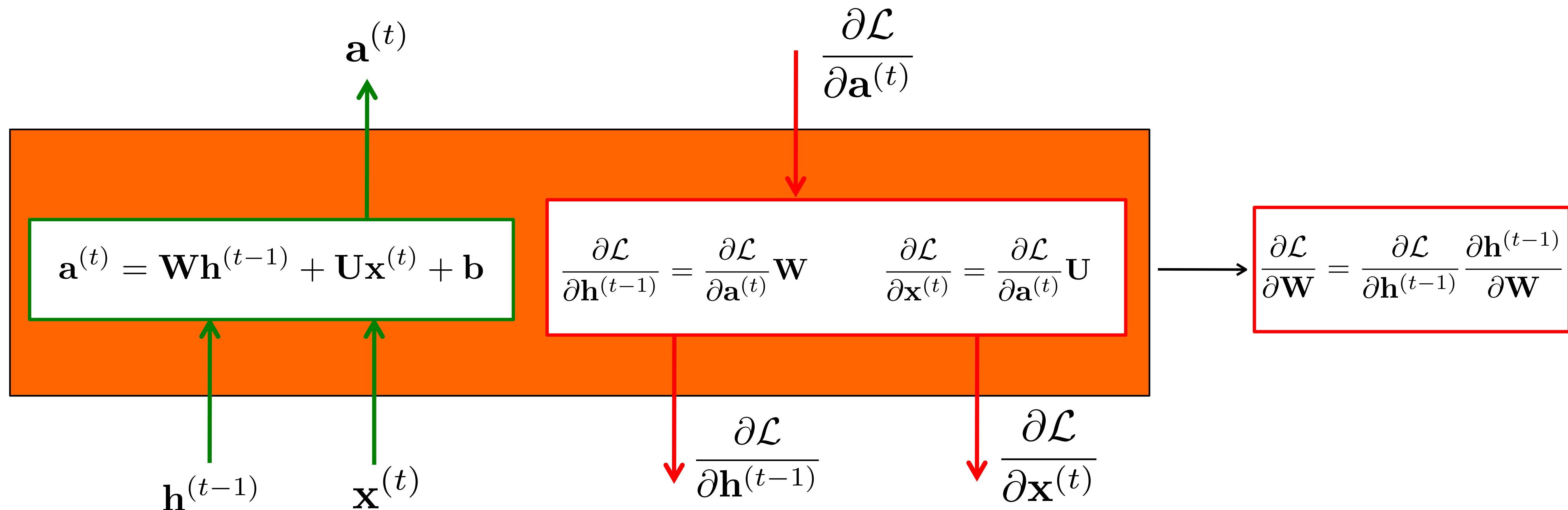


$$\frac{\partial \hat{y}^{(t)}}{\partial \mathbf{x}^{(0)}} = \frac{\partial \hat{y}^{(t)}}{\partial \mathbf{h}^{(t)}} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} \cdots \frac{\partial \mathbf{h}^{(1)}}{\partial \mathbf{h}^{(0)}} \frac{\partial \mathbf{h}^{(0)}}{\partial \mathbf{x}^{(0)}}$$

$$\frac{\partial J}{\partial \mathbf{W}} = \sum_{t=0}^T \frac{\partial \mathcal{L}(\hat{\mathbf{y}}^{(t)}, \mathbf{y}^{(t)})}{\partial \mathbf{W}}$$

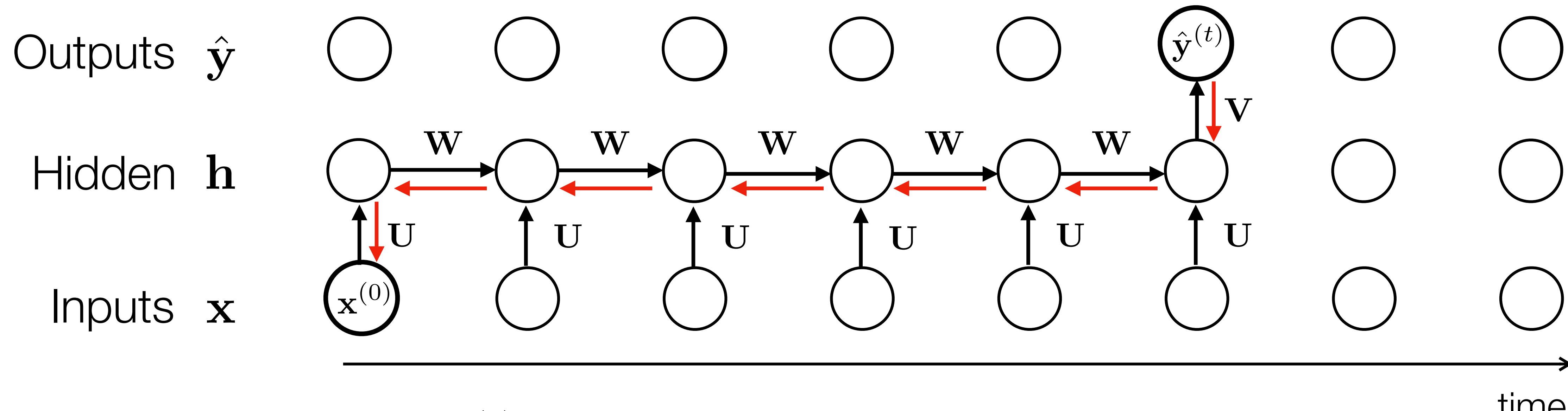


Recurrent linear layer



$$\frac{\partial J}{\partial \mathbf{W}} = \sum_{t=0}^T \frac{\partial \mathcal{L}(\hat{\mathbf{y}}^{(t)}, \mathbf{y}^{(t)})}{\partial \mathbf{W}}$$

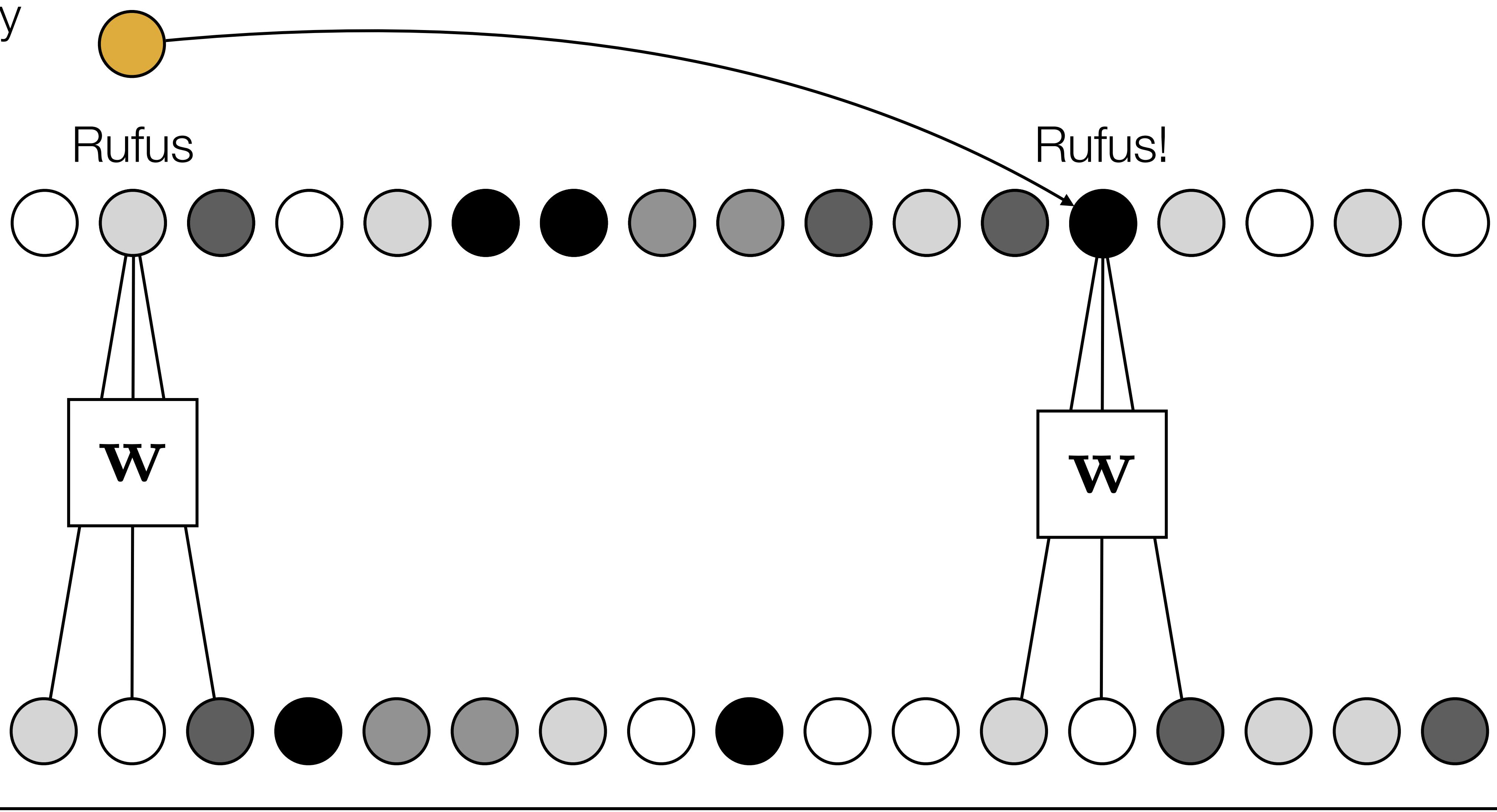
The problem of long-range dependences



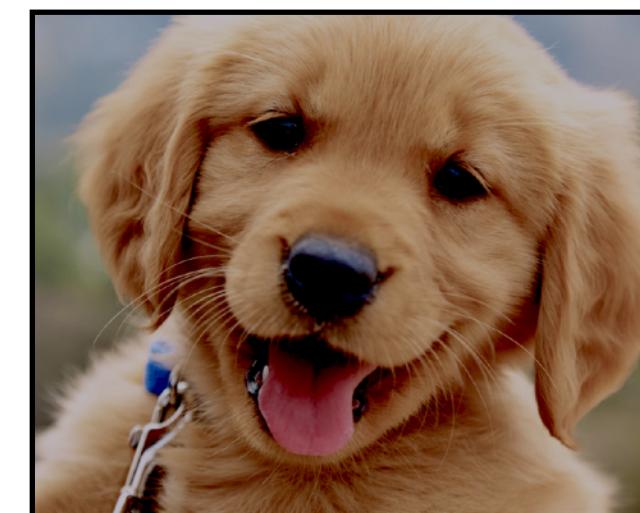
$$\frac{\partial \hat{y}^{(t)}}{\partial x^{(0)}} = \frac{\partial \hat{y}^{(t)}}{\partial h^{(t)}} \frac{\partial h^{(t)}}{\partial h^{(t-1)}} \cdots \frac{\partial h^{(1)}}{\partial h^{(0)}} \frac{\partial h^{(0)}}{\partial x^{(0)}}$$

- Capturing long-range dependences requires propagating information through a long chain of dependences.
- Old observations are forgotten
- Stochastic gradients become high variance (noisy), and gradients may **vanish or explode**

Memory
unit



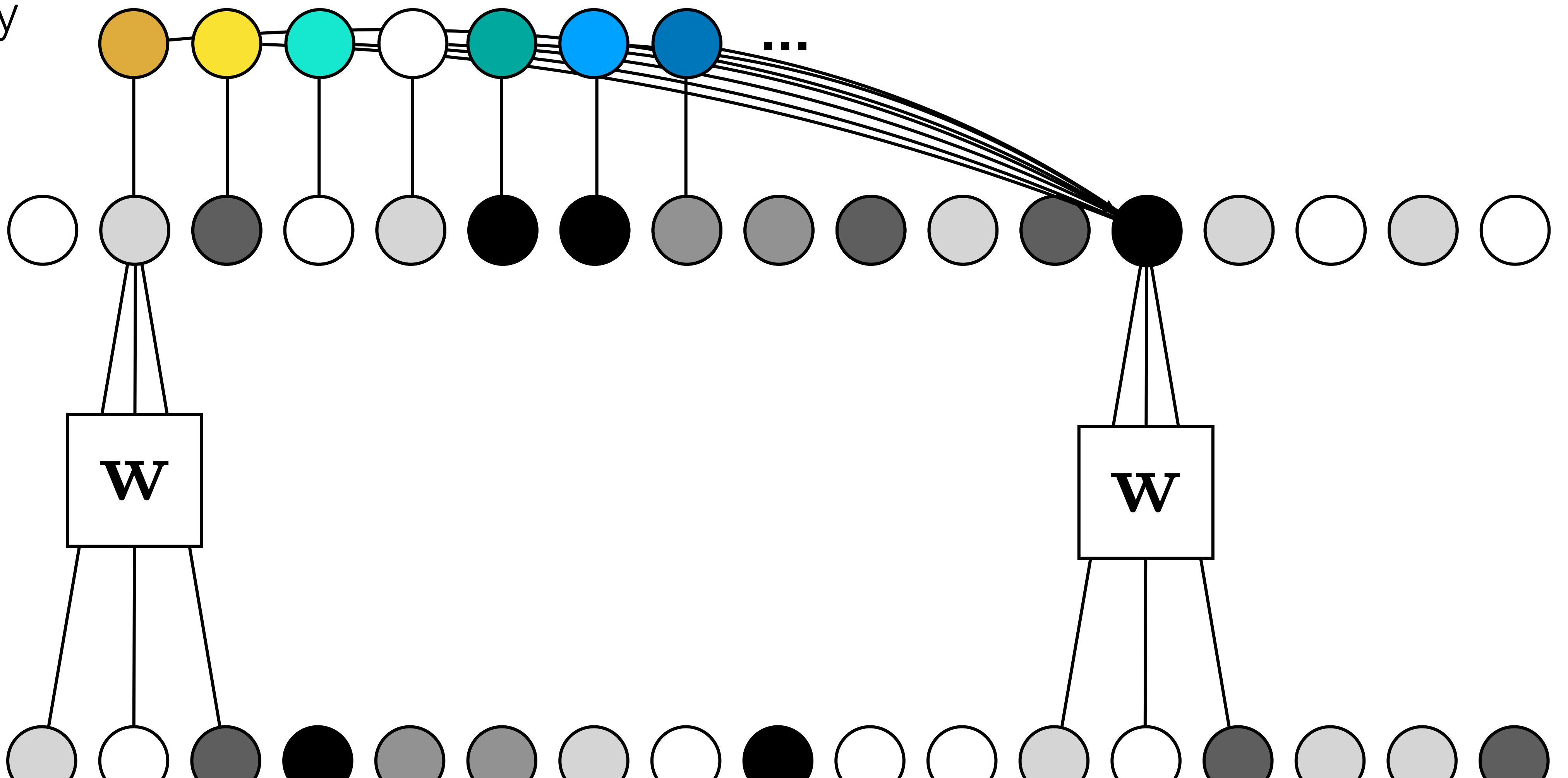
50



time

Source: Torralba, Freeman, Isola

Memory
units



51

time

Source: Torralba, Freeman, Isola

The problem of long-range dependences

Why not remember everything?

- Memory size grows with t
- This kind of memory is **nonparametric**: there is no finite set of parameters we can use to model it
- Markov-like assumption: future state only dependent on preceding hidden state. Very easy to forget things.
- By putting the right info in to the hidden state, RNNs can model dependencies that are arbitrarily⁶² far apart

LSTMs

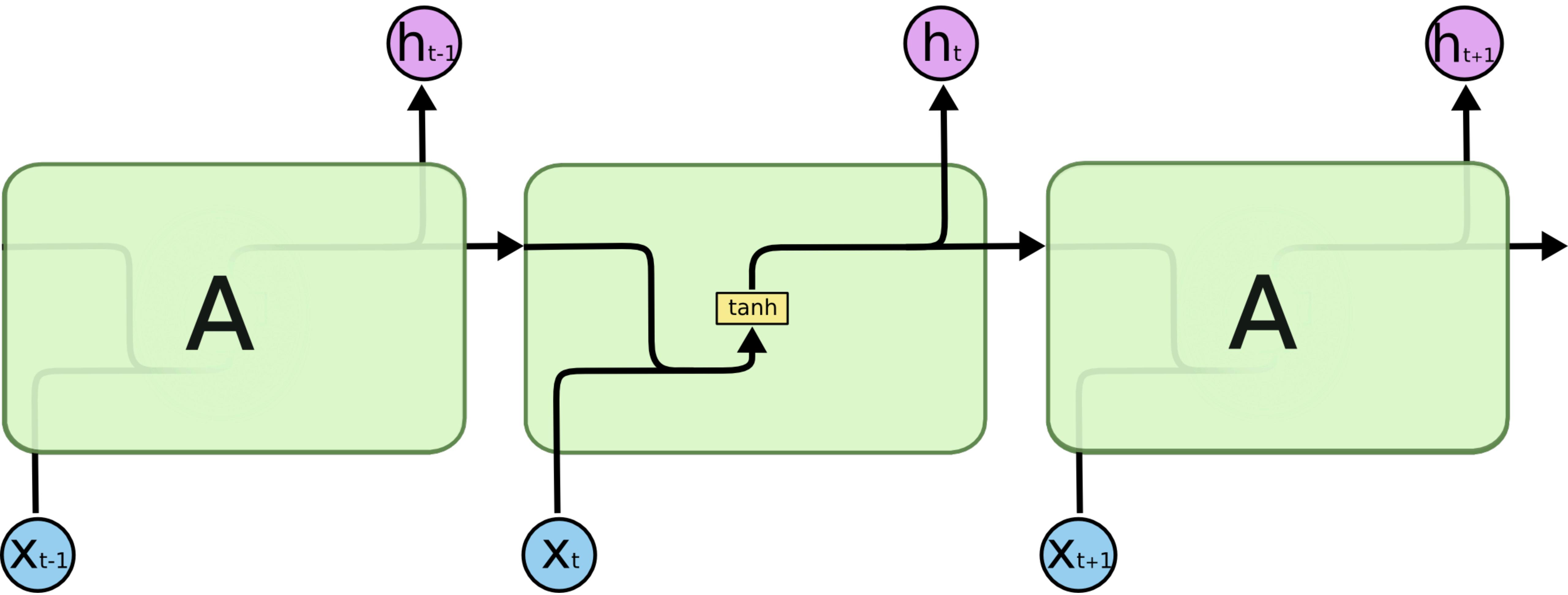
Long Short Term Memory

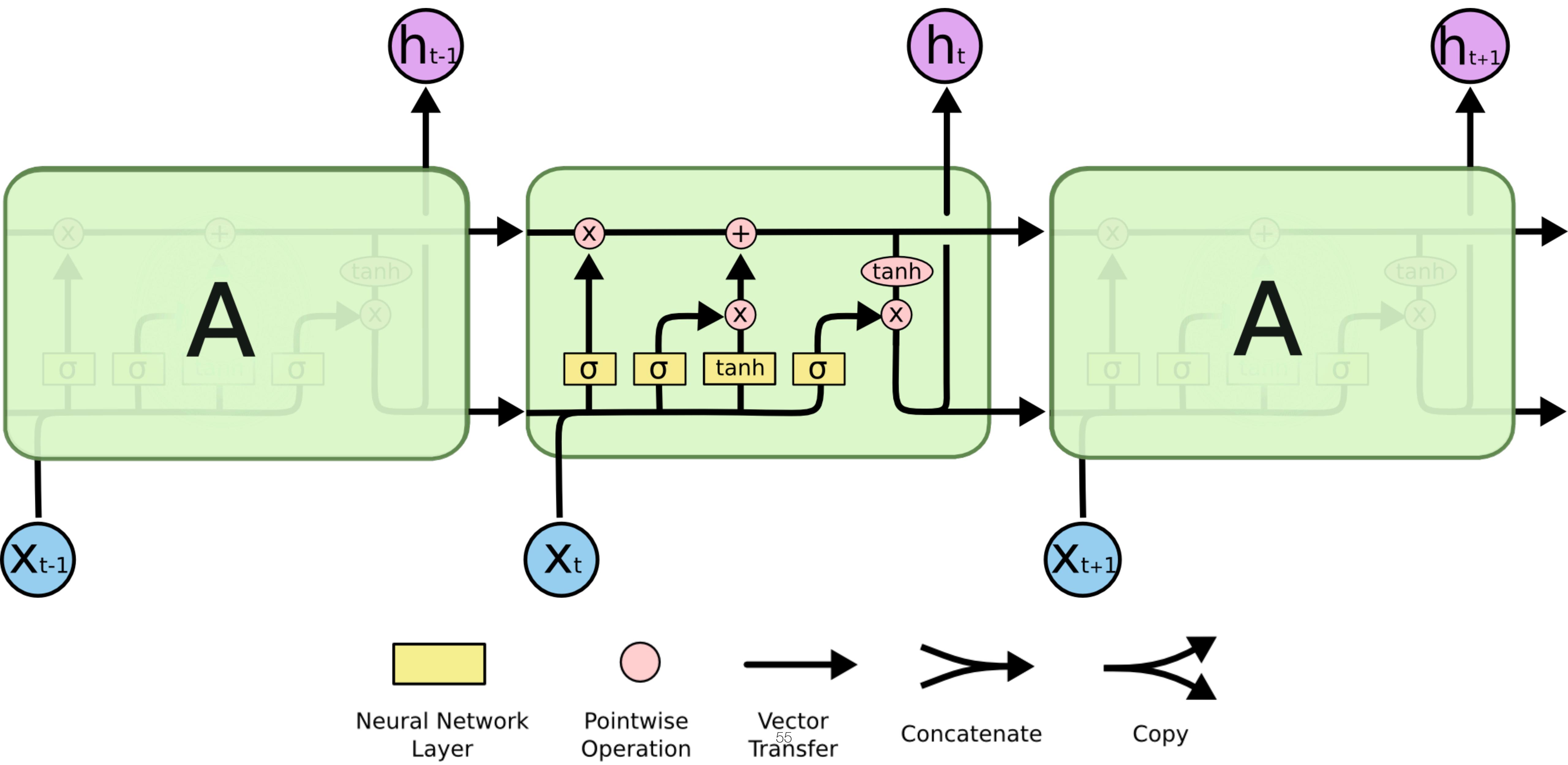
A special kind of RNN designed to avoid forgetting.

Related to ResNets inductive bias is that state transition is an identity function.

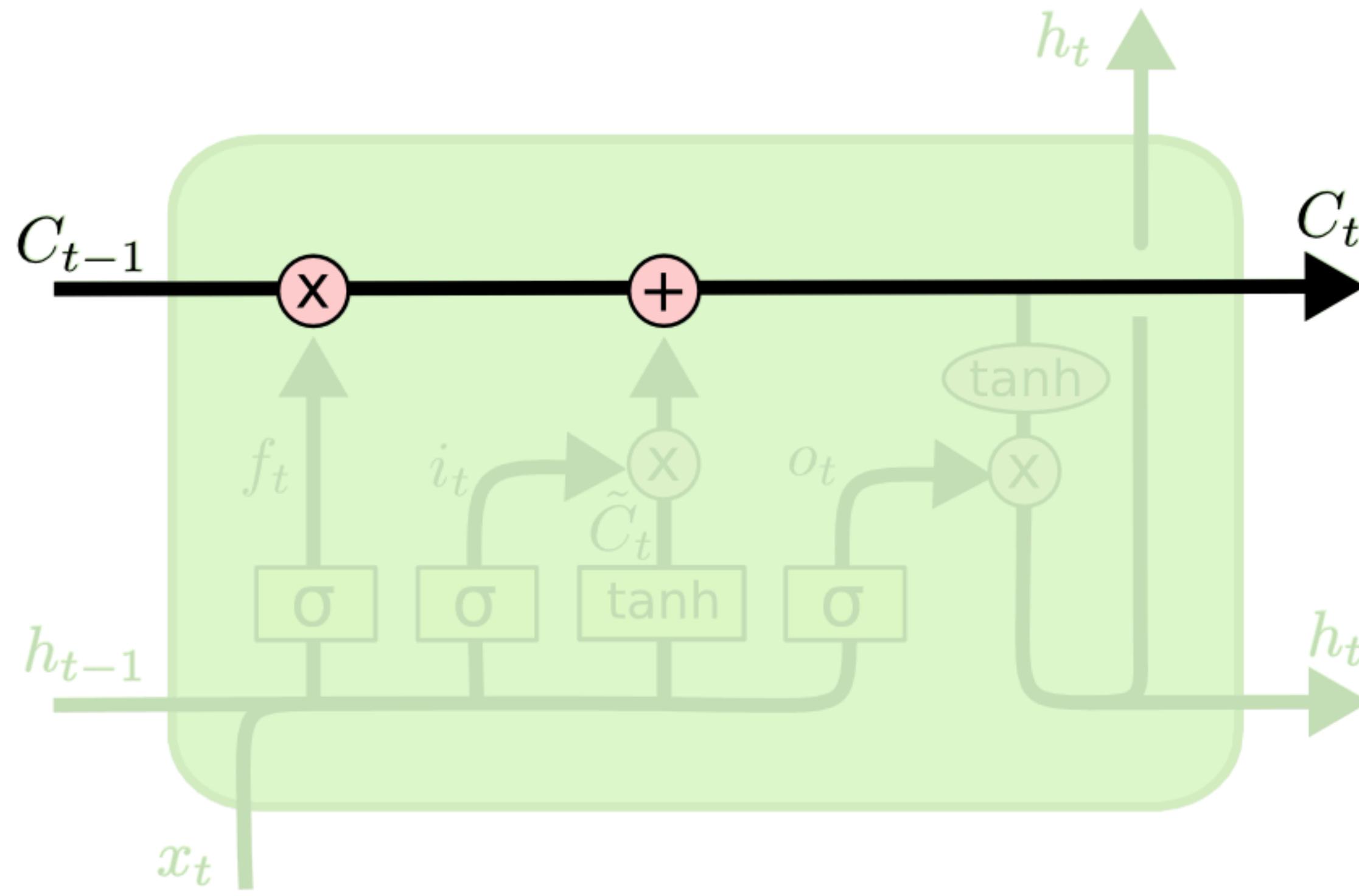
This way the default behavior is not to forget an old state. Instead of forgetting by default, the network has to *learn to forget*.

Bit of a complex design. Works well but simpler methods like Gated Recurrent Unit (GRU) are competitive [Jozefowicz et al. 2015].

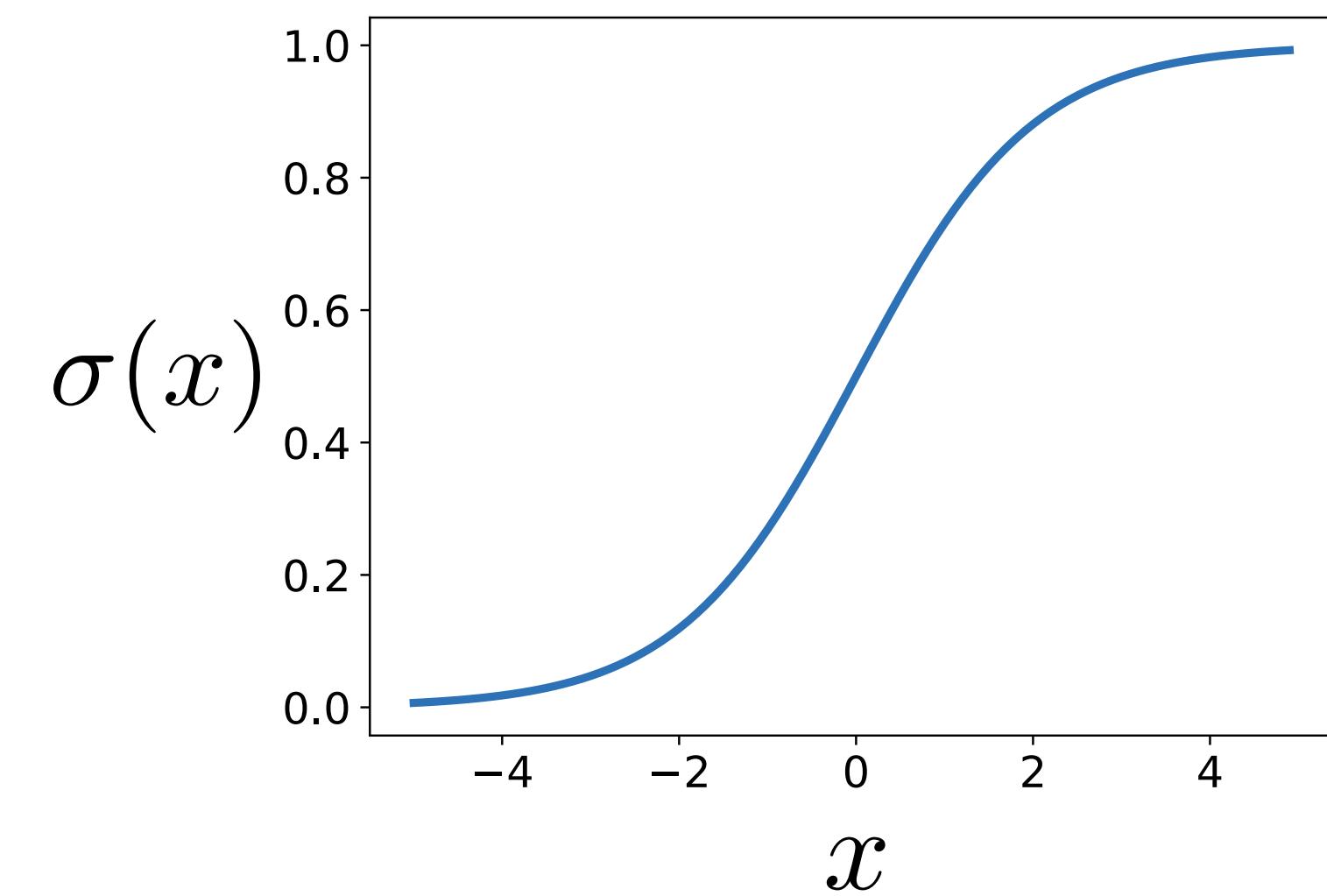
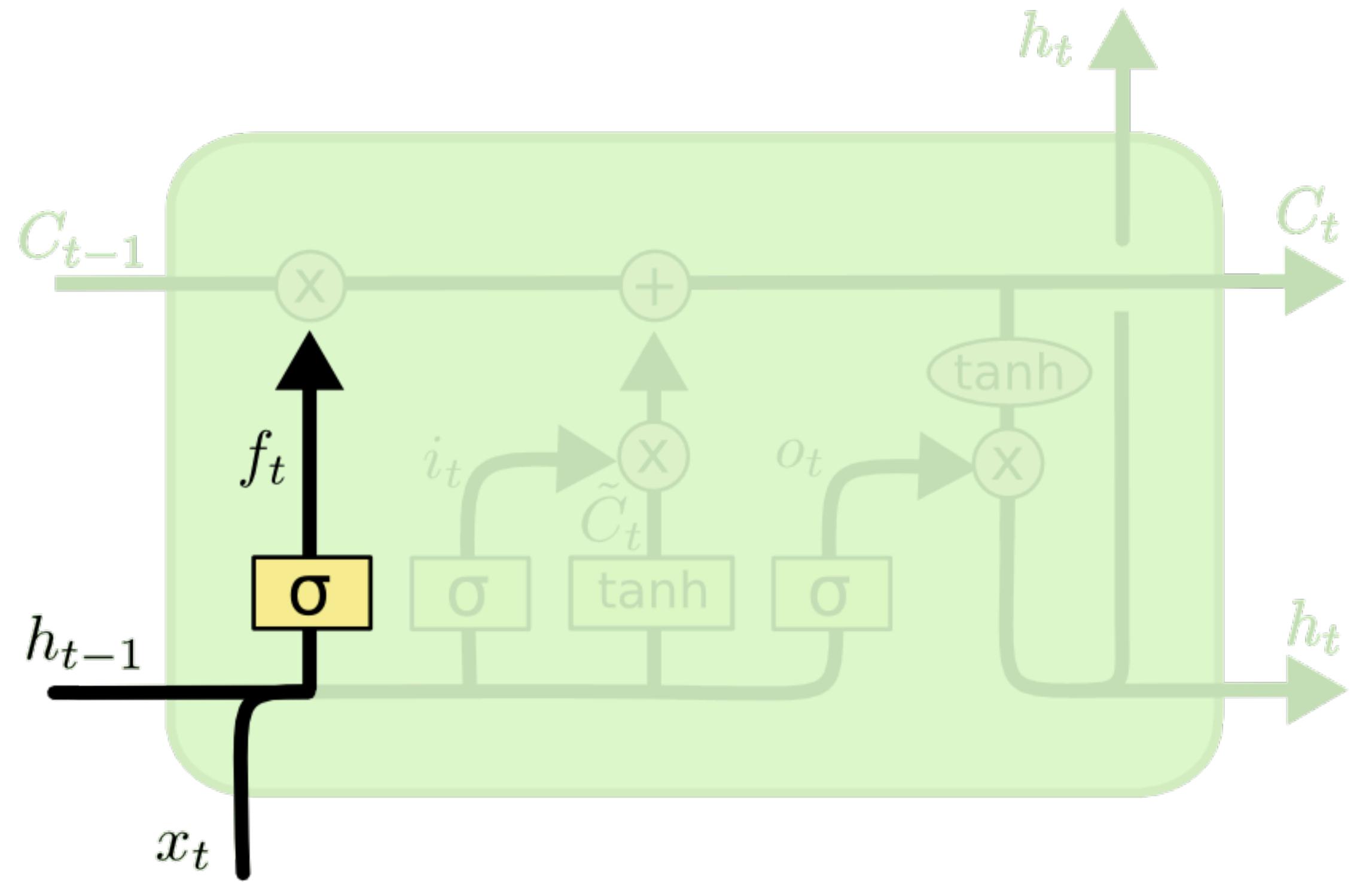




[Slide derived from Chris Olah: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>]



$C_t = \text{Cell state}$

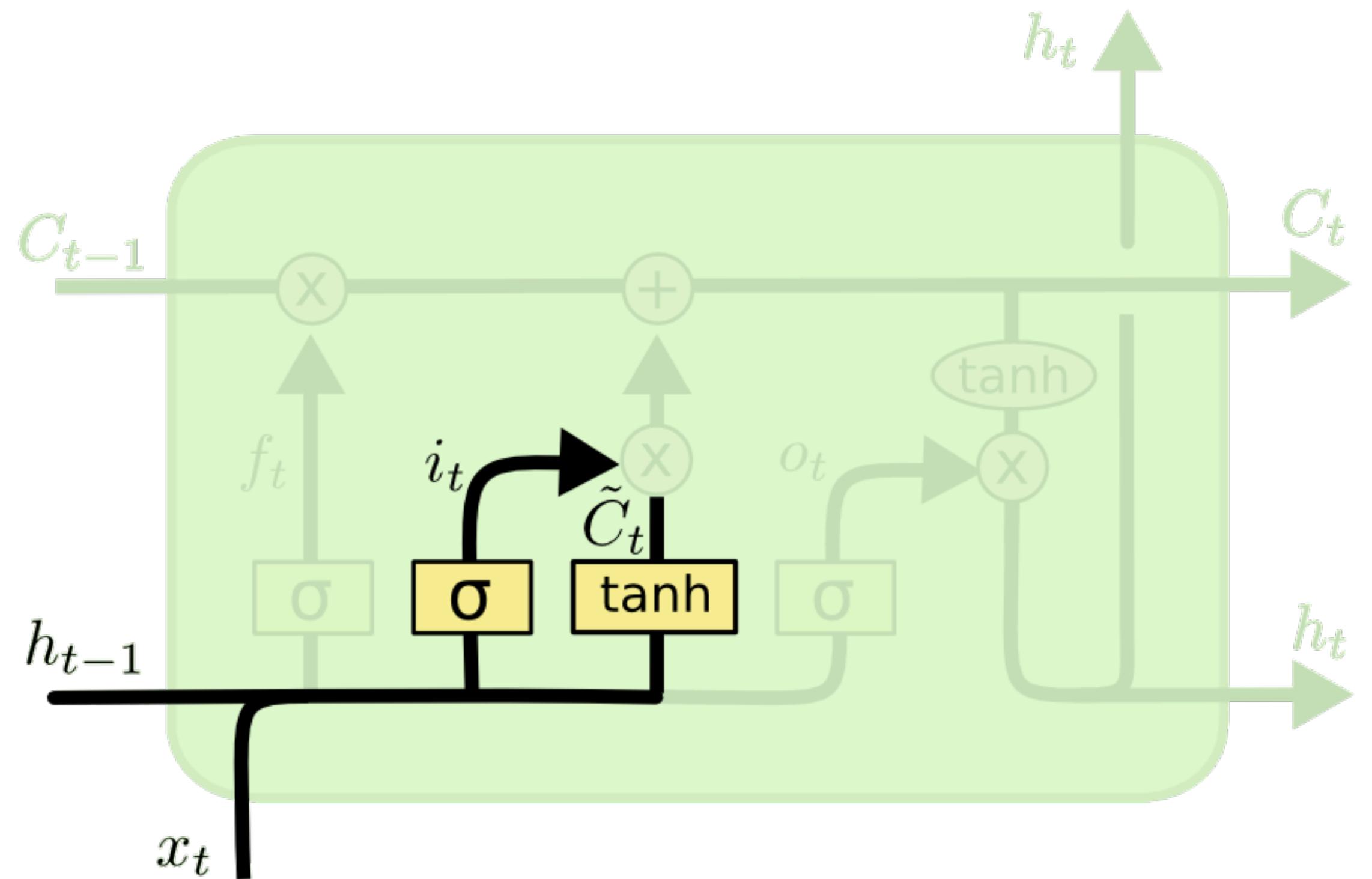


$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

Decide what information to throw away from the cell state.

Each element of cell state is multiplied by ~ 1 (remember) or ~ 0 (forget).

[Slide derived from Chris Olah: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>]



which indices to write to

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

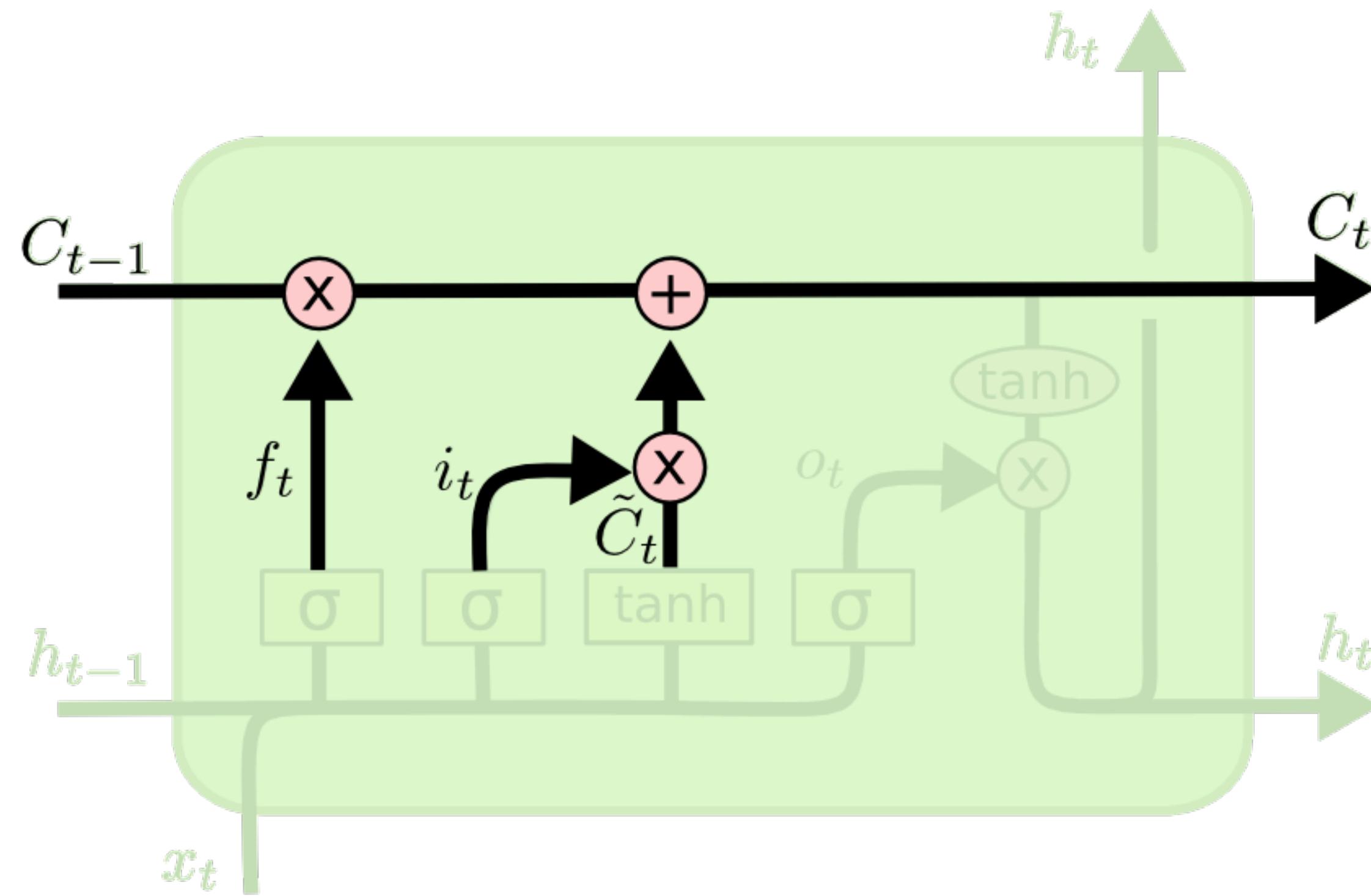
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

what to write to those indices

Decide what new information to add to the cell state.

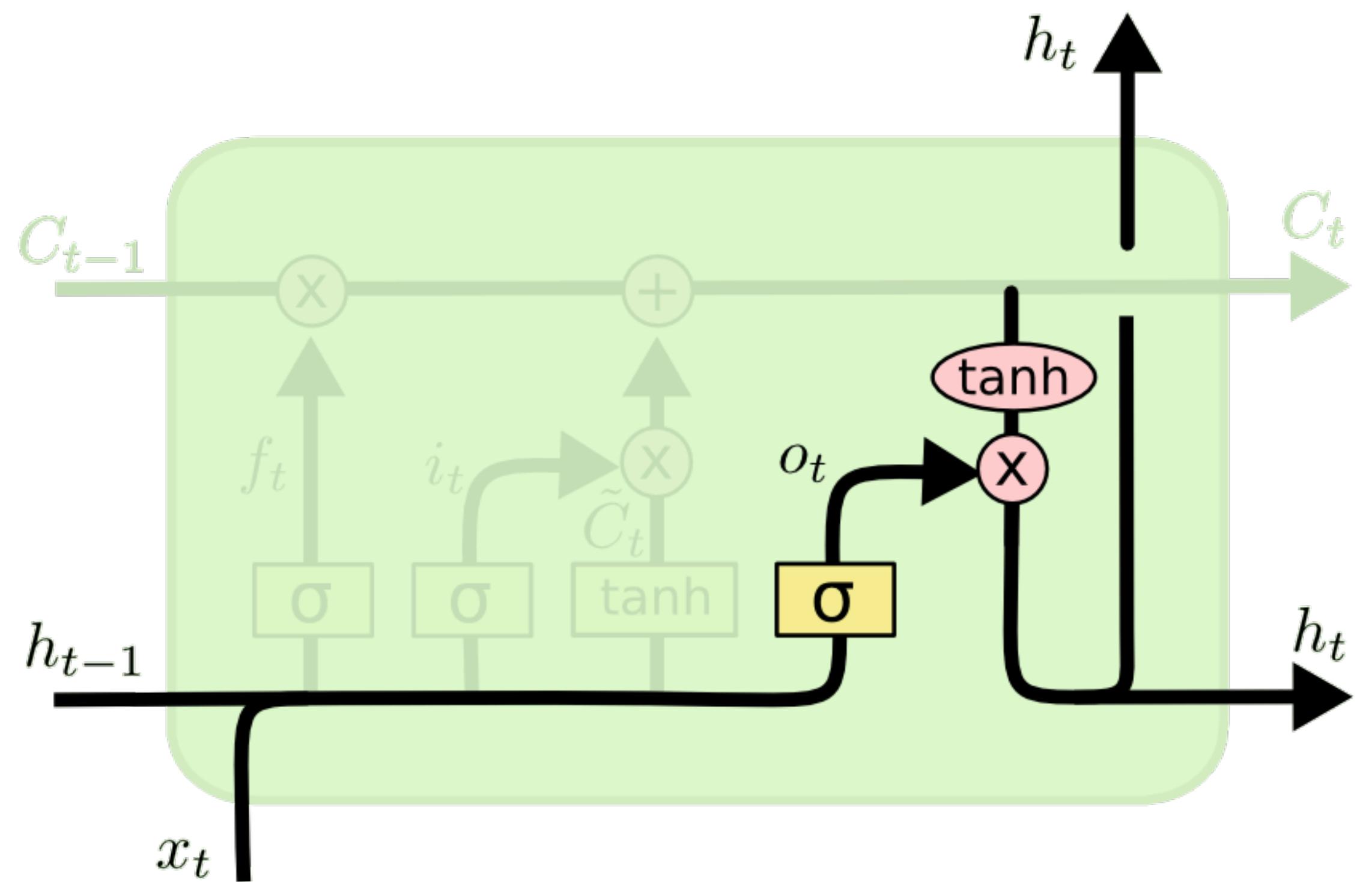
58

[Slide derived from Chris Olah: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>]



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Forget selected old information, write selected new information.



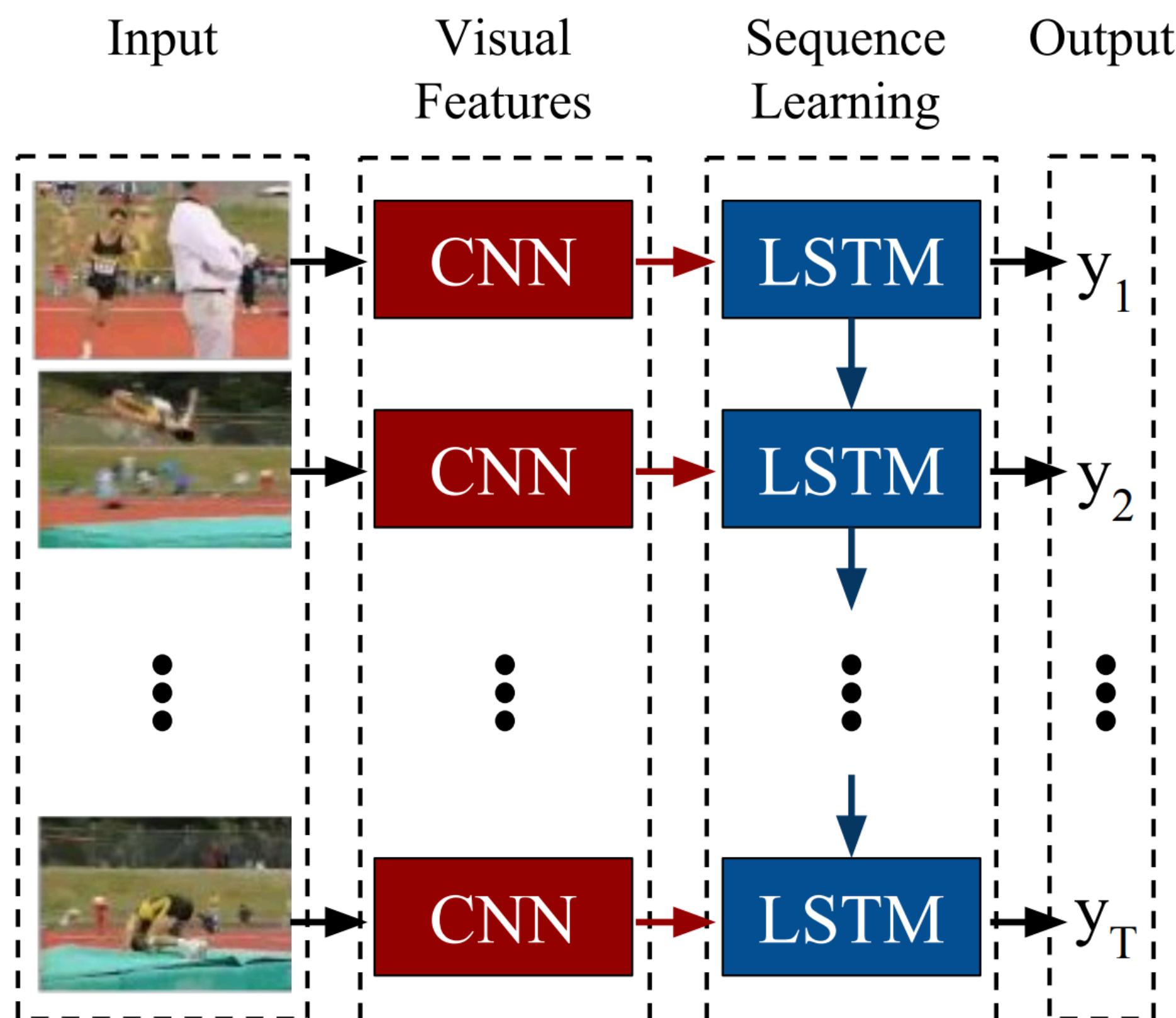
$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh (C_t)$$

After having updated the cell state's information, decide what to output.

[Slide derived from Chris Olah: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>]

Some uses for LSTMs



Activity Recognition
Sequences in the Input

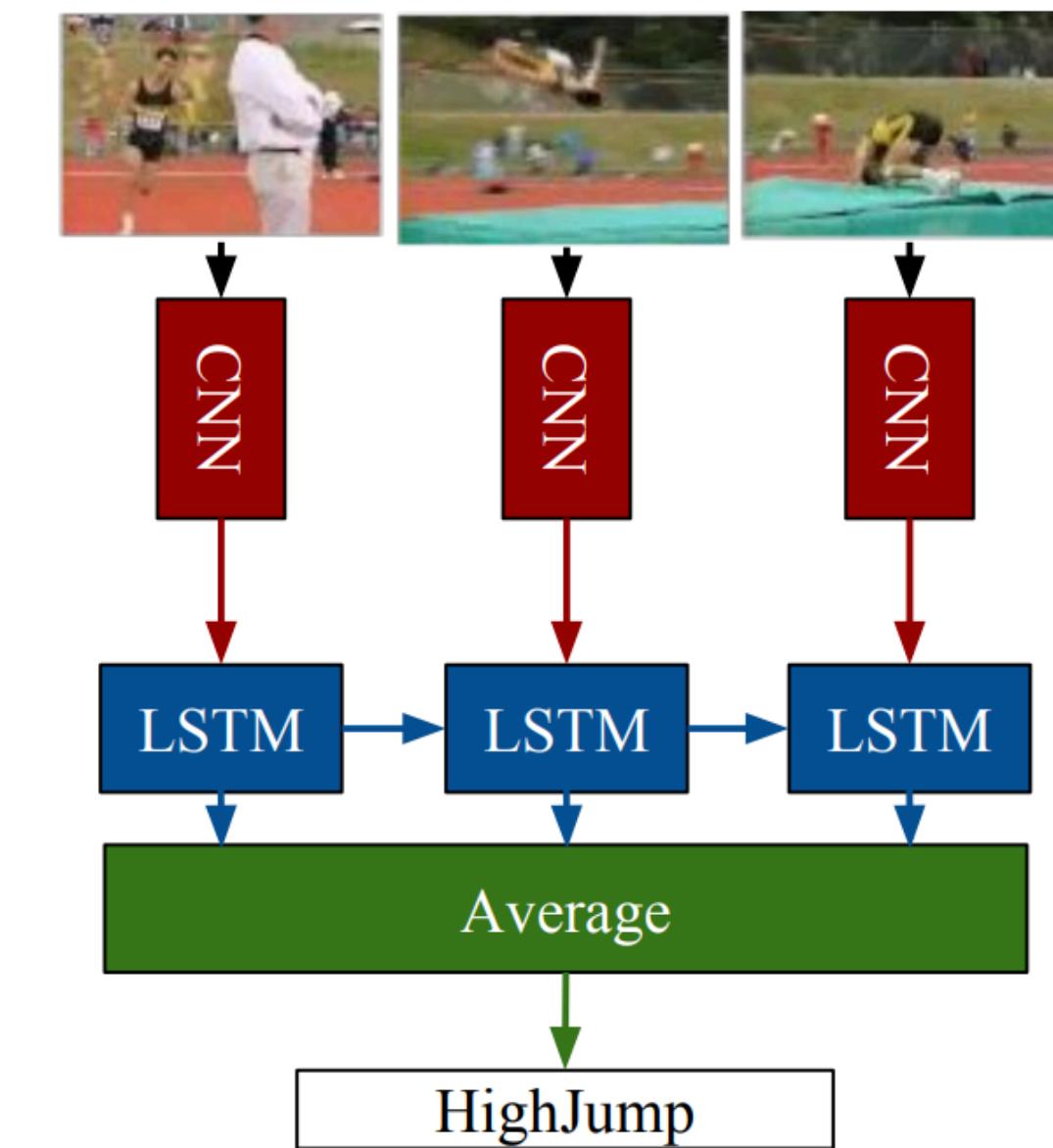
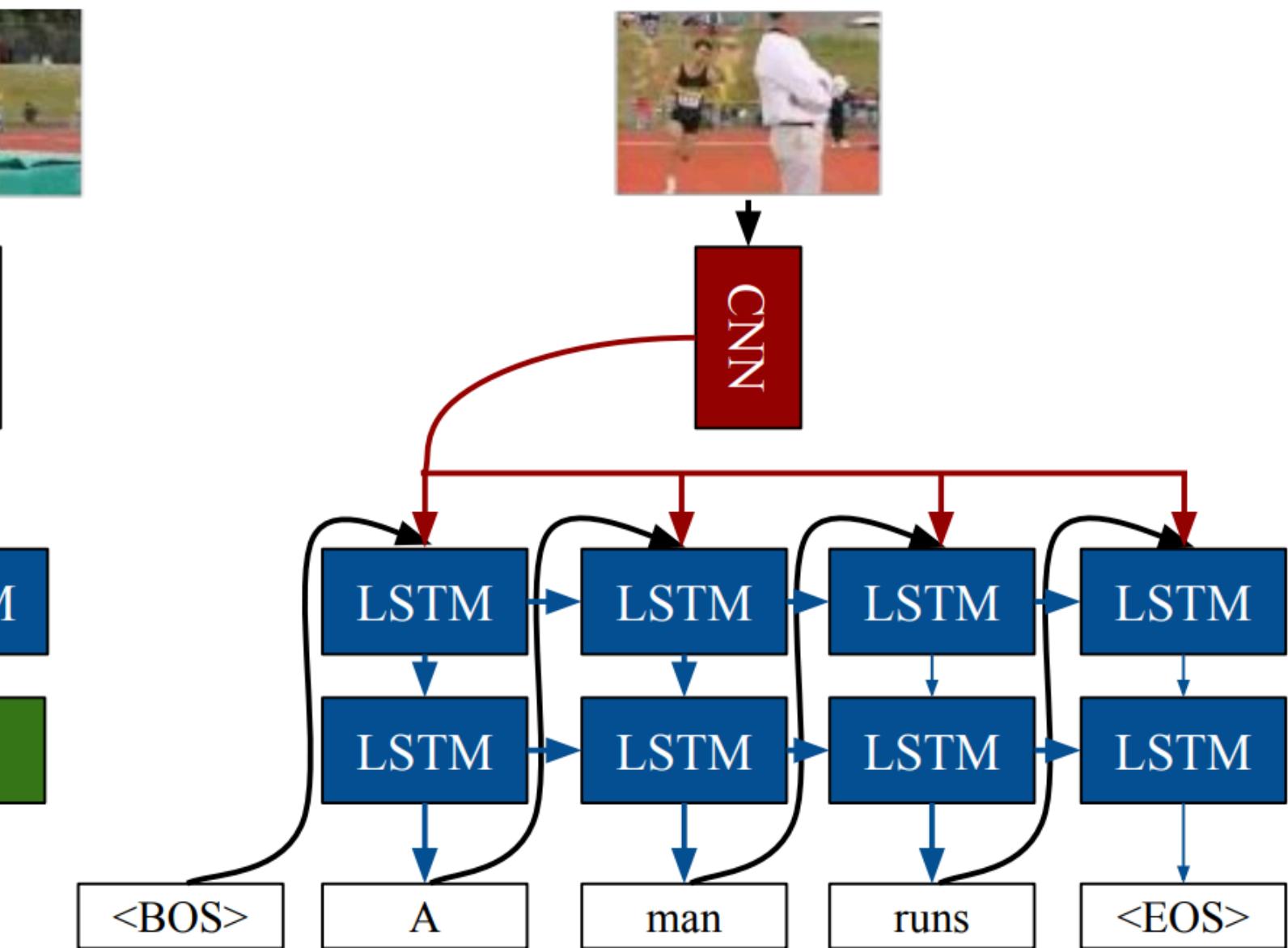


Image Captioning
Sequences in the Output



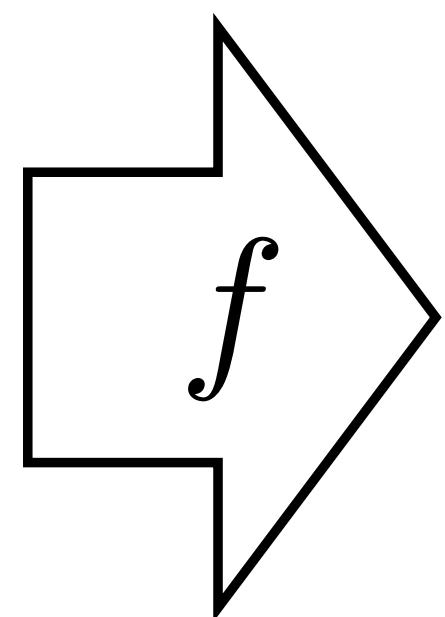
The problem of long-range dependences

Other methods exist that do directly link old “memories” (observations or hidden states) to future predictions:

- Temporal convolutions (but temporal horizon often short)
- Attention (see <https://arxiv.org/abs/1706.03762>)
- Memory networks (see <https://arxiv.org/abs/1410.3916>)

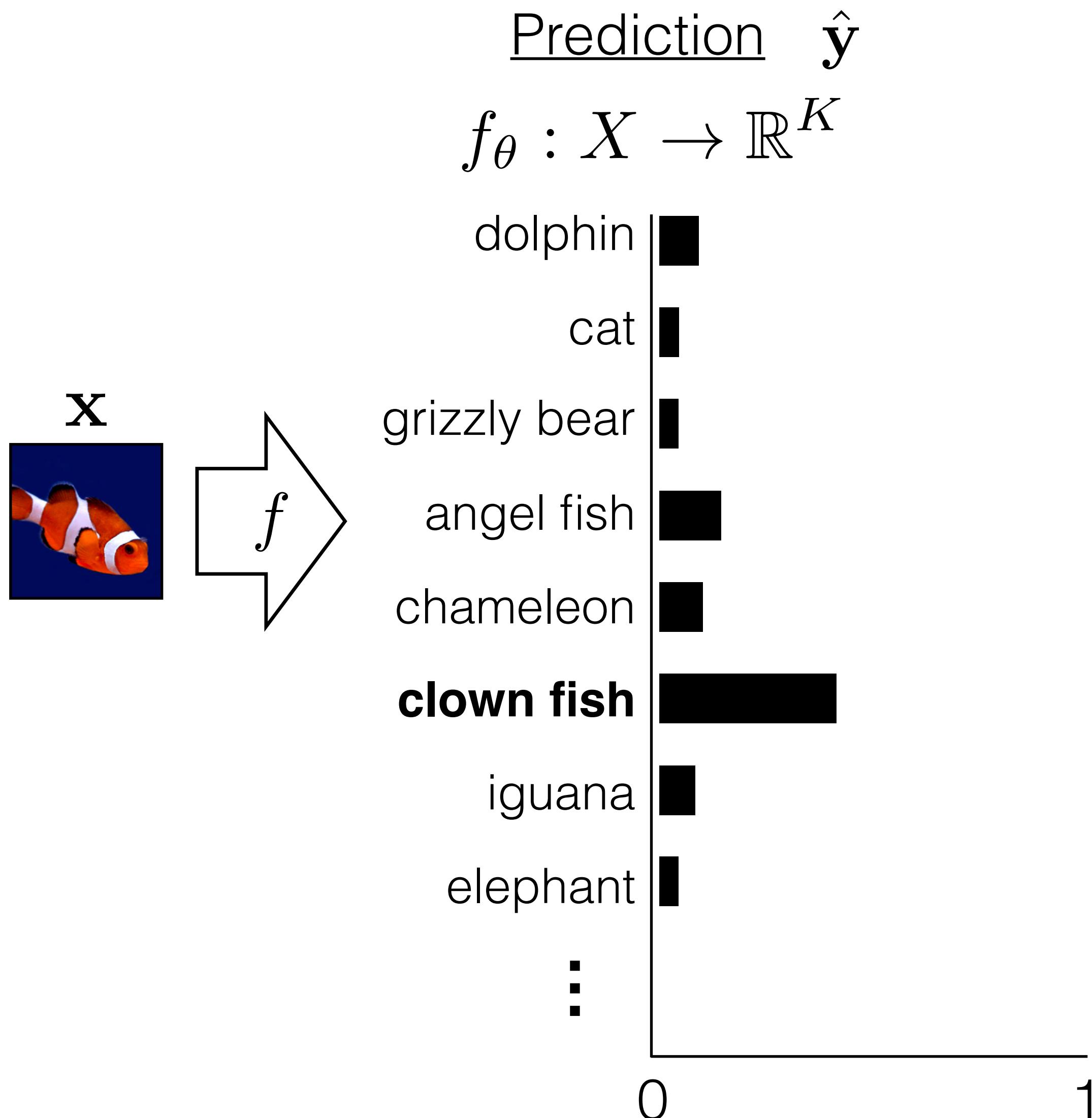
Language

Image captioning



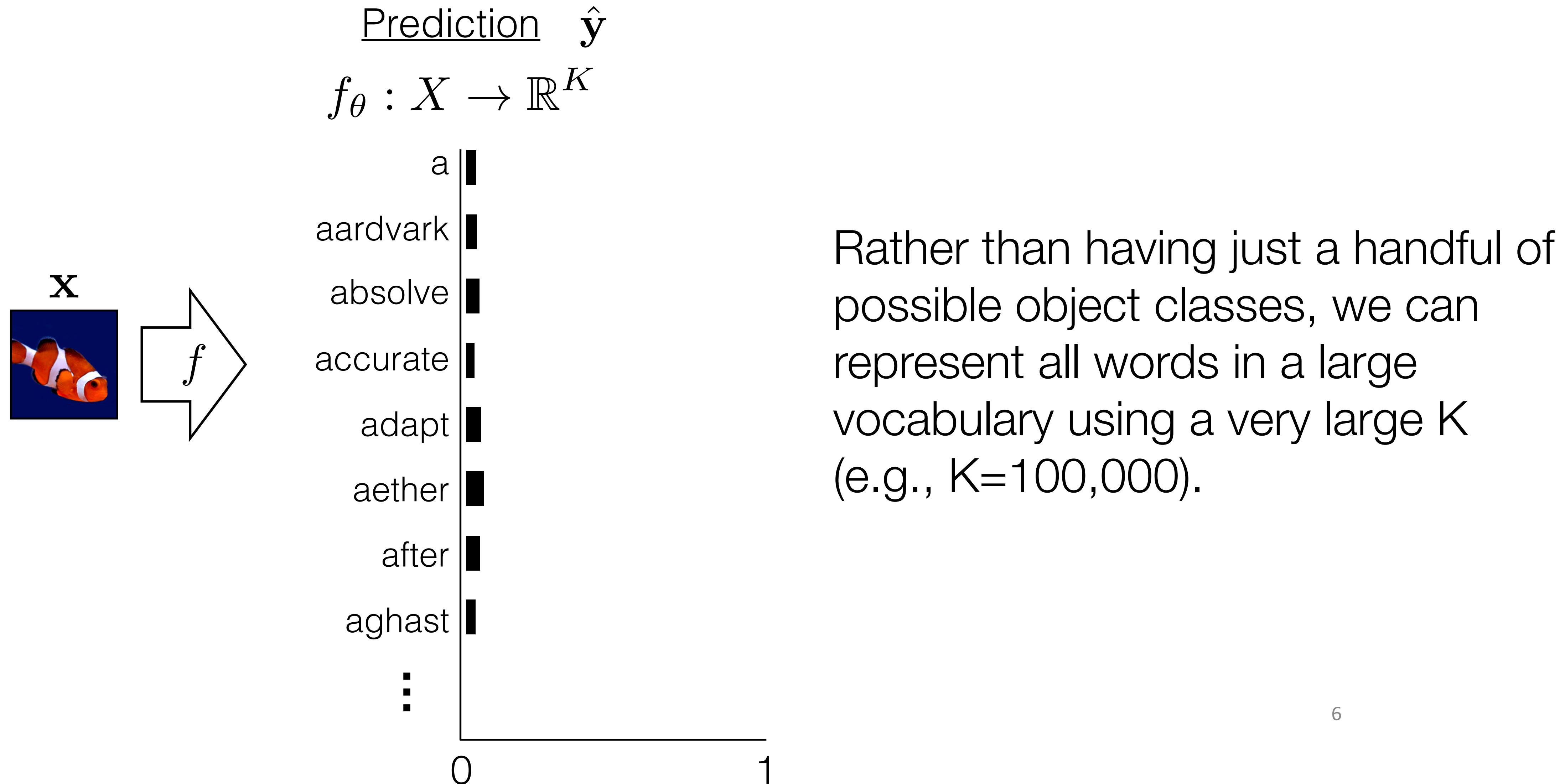
“A flock of birds against
a gray sky”

How to represent words as numbers?

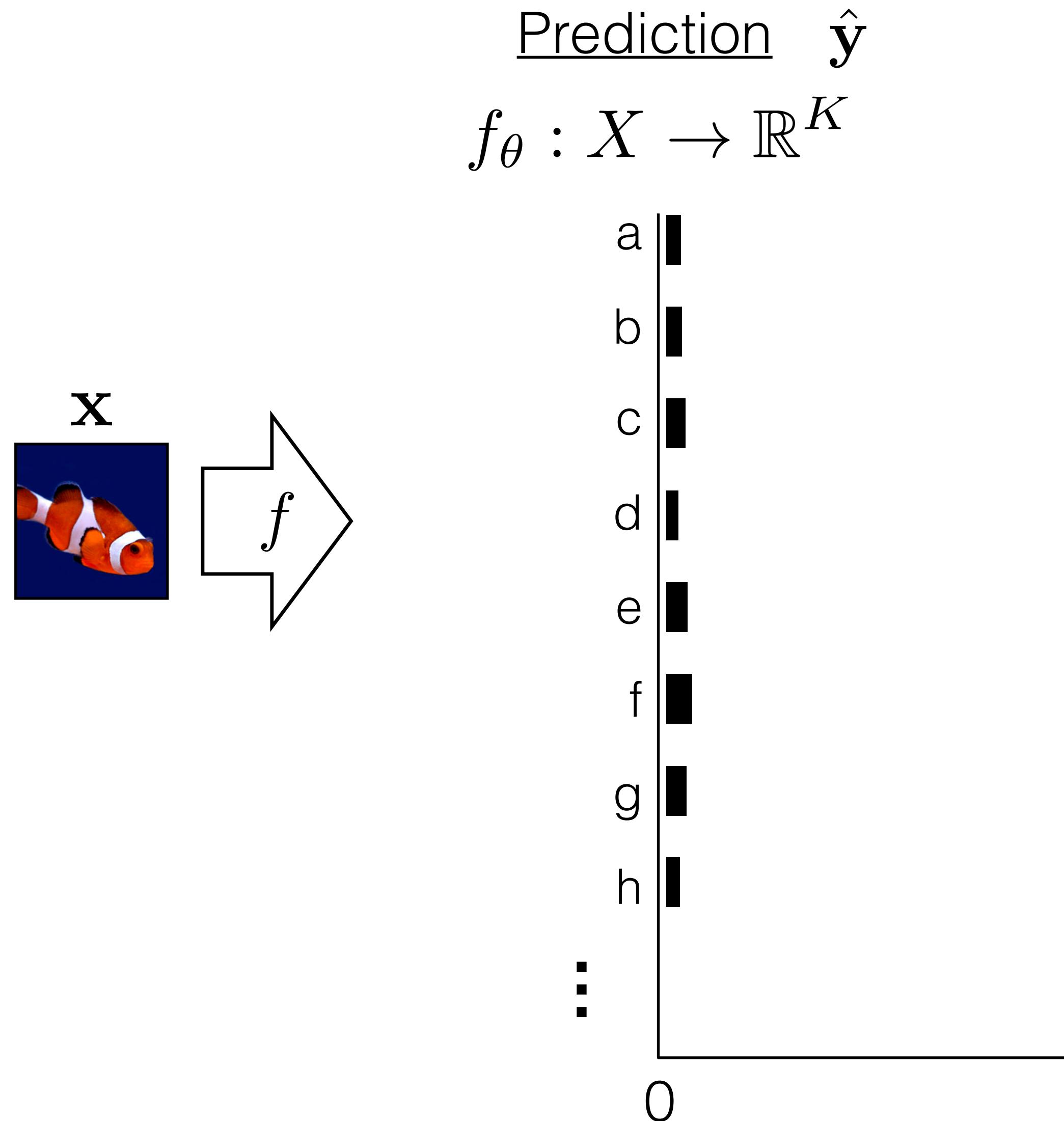


5

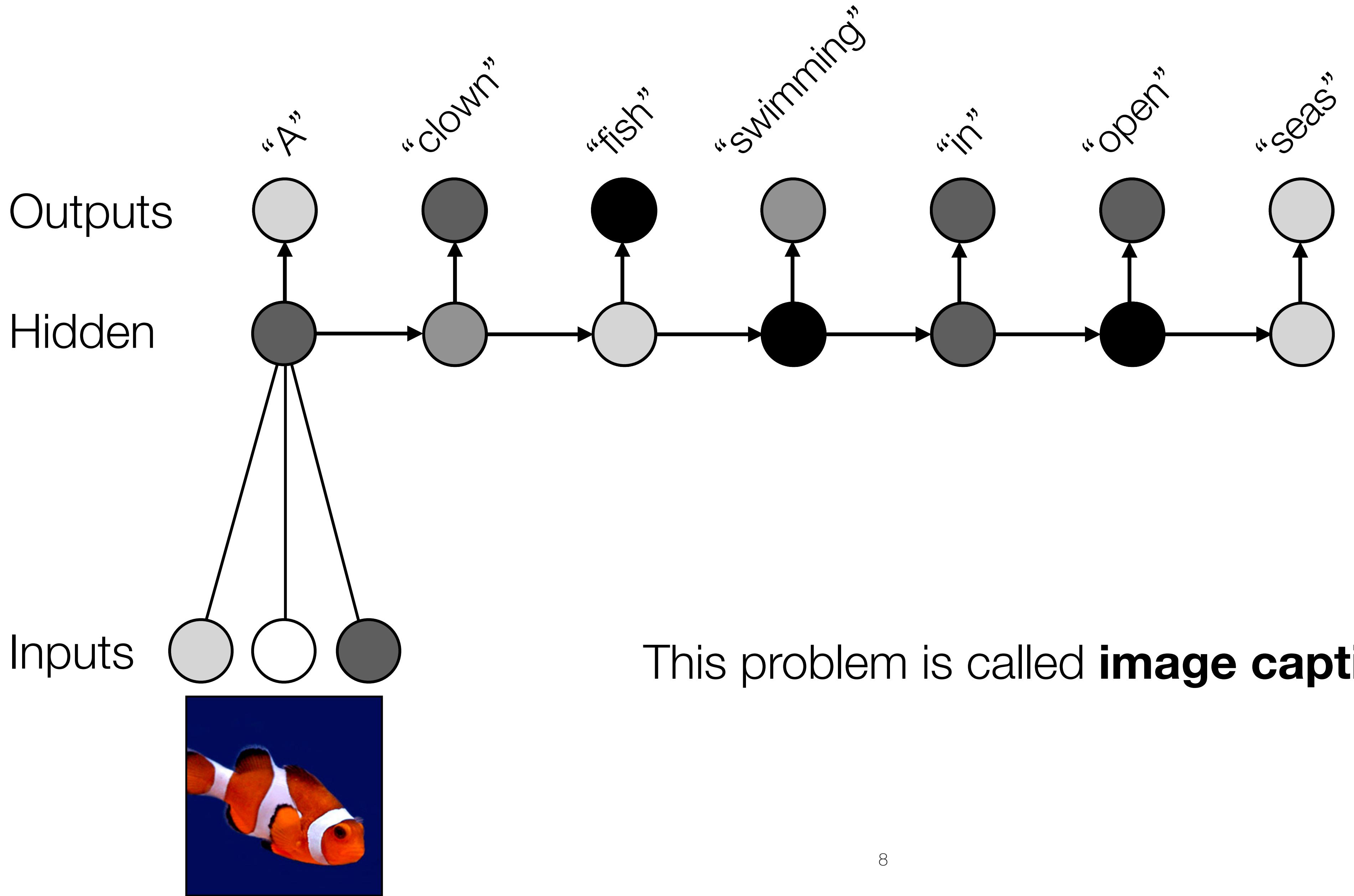
How to represent words as numbers?



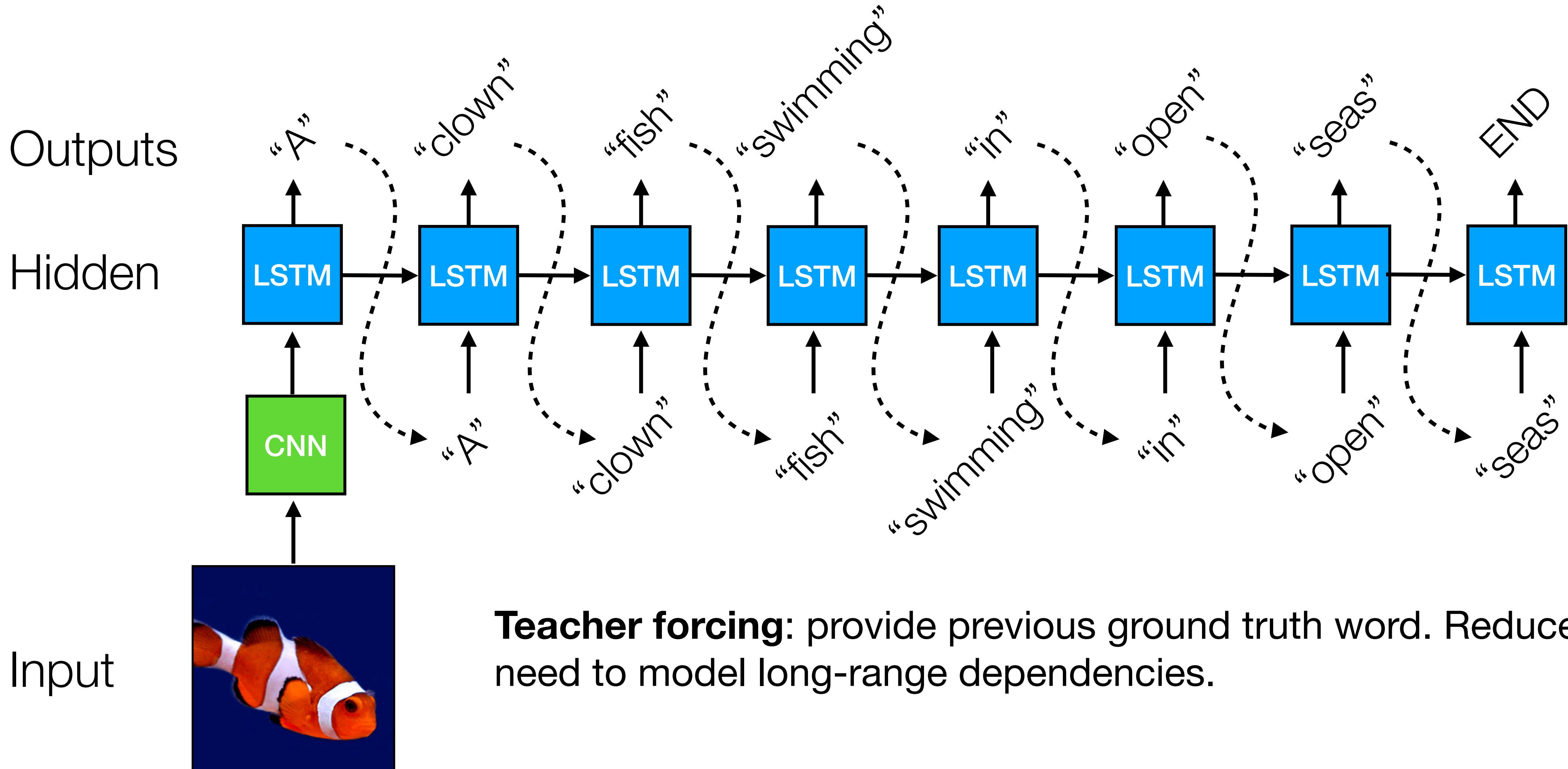
How to represent words as numbers?



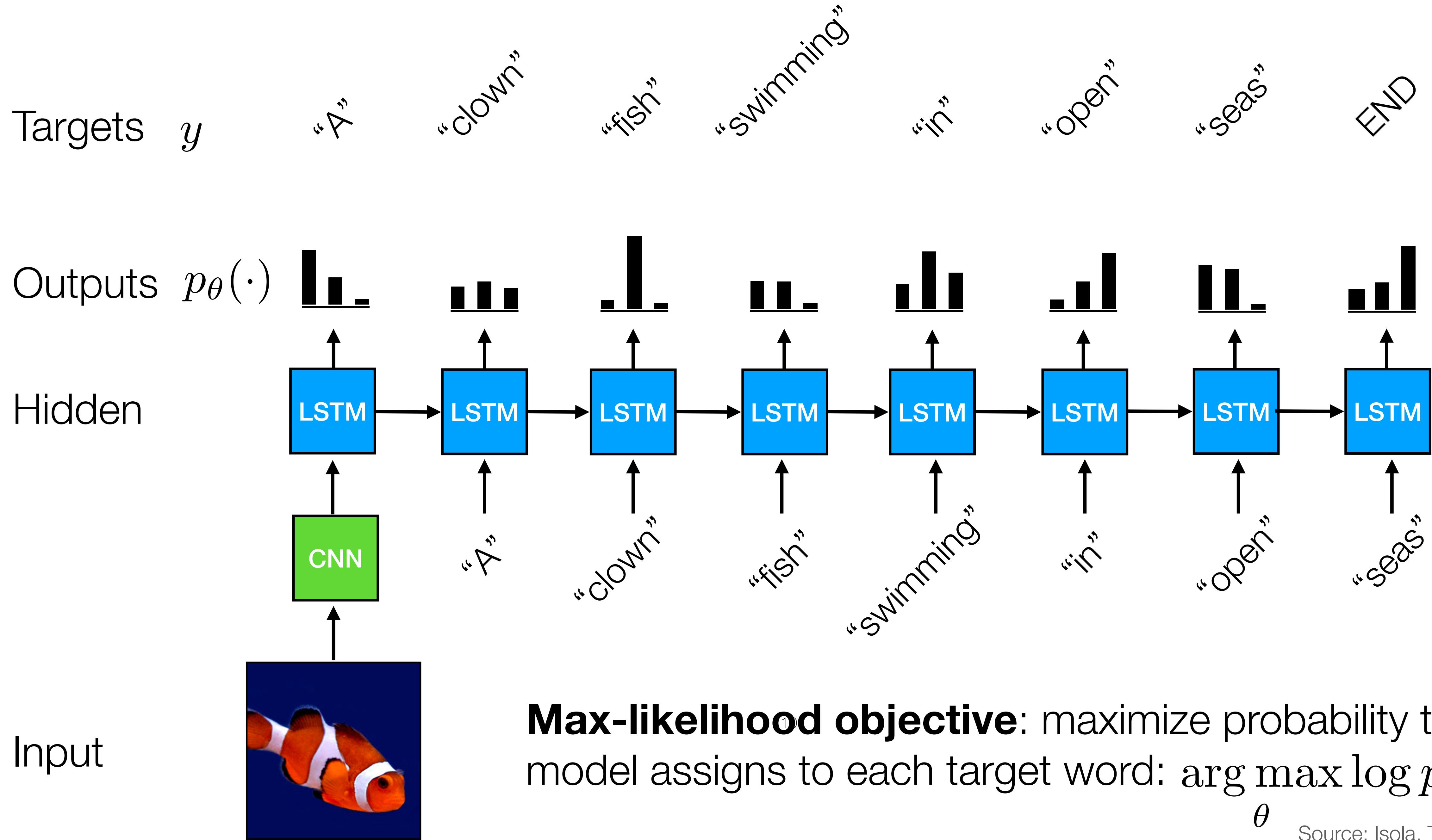
Or, represent each character as a class (e.g., $K=26$ for English letters), and represent words as a sequence of characters.



This problem is called **image captioning**.



Teacher forcing: provide previous ground truth word. Reduces need to model long-range dependencies.



Max-likelihood objective: maximize probability the model assigns to each target word: $\arg \max_{\theta} \log p_{\theta}(y)$

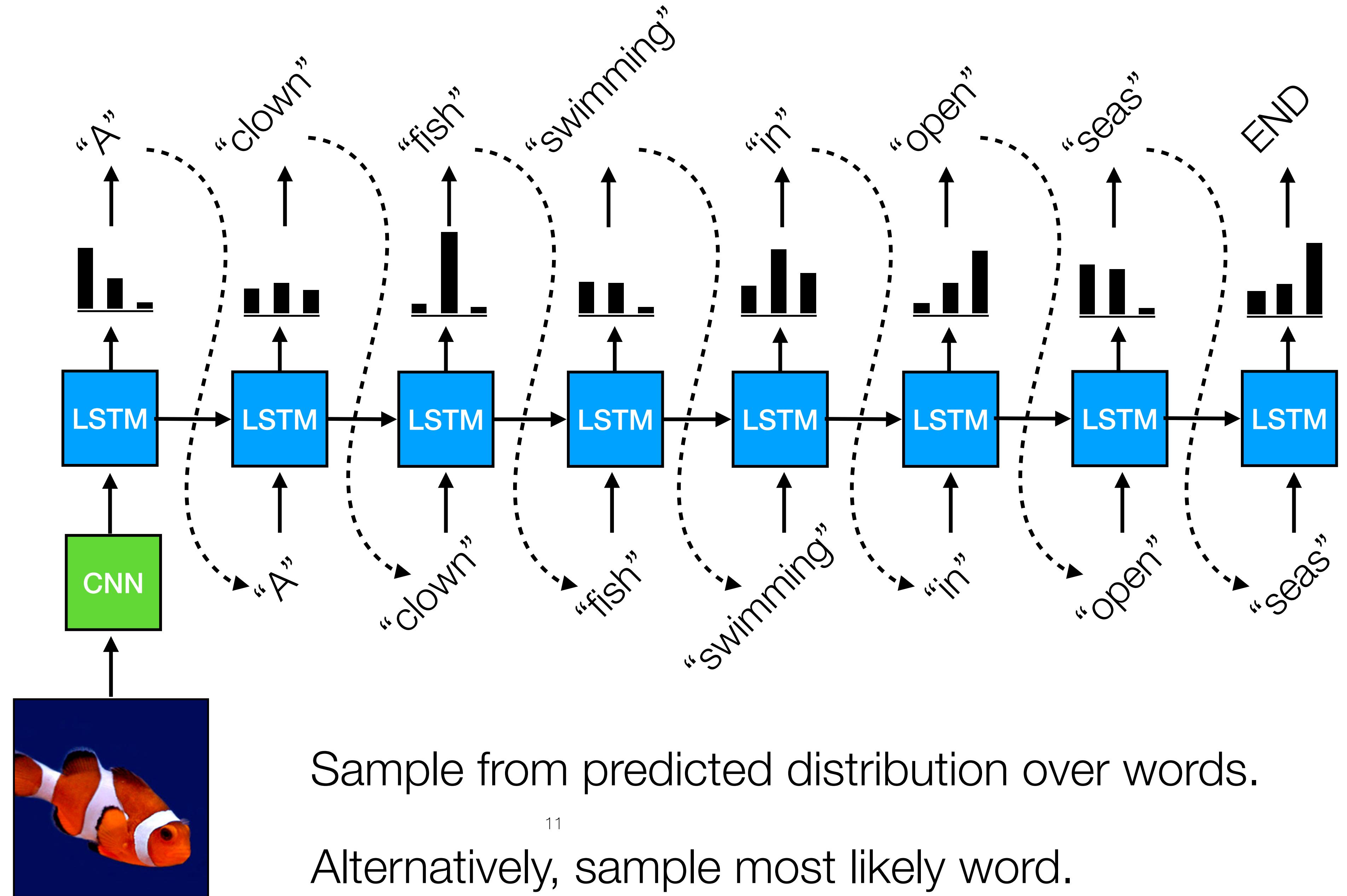
Testing

Samples

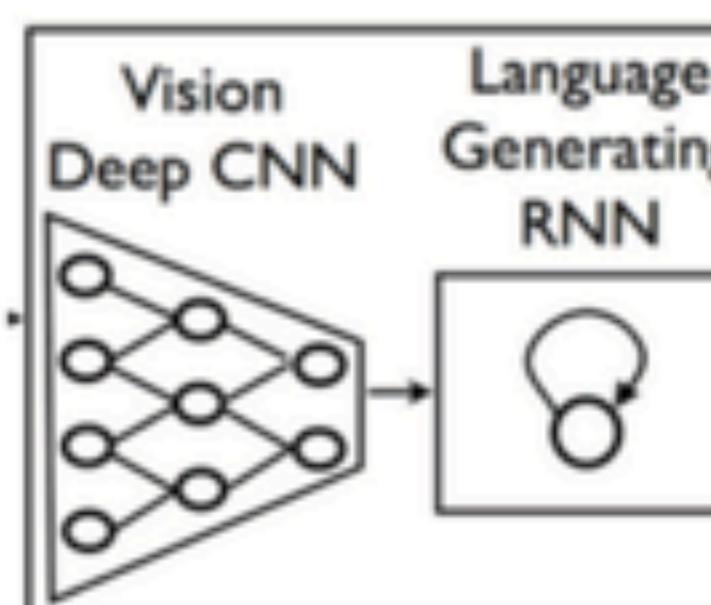
Outputs $p_\theta(\cdot)$

Hidden

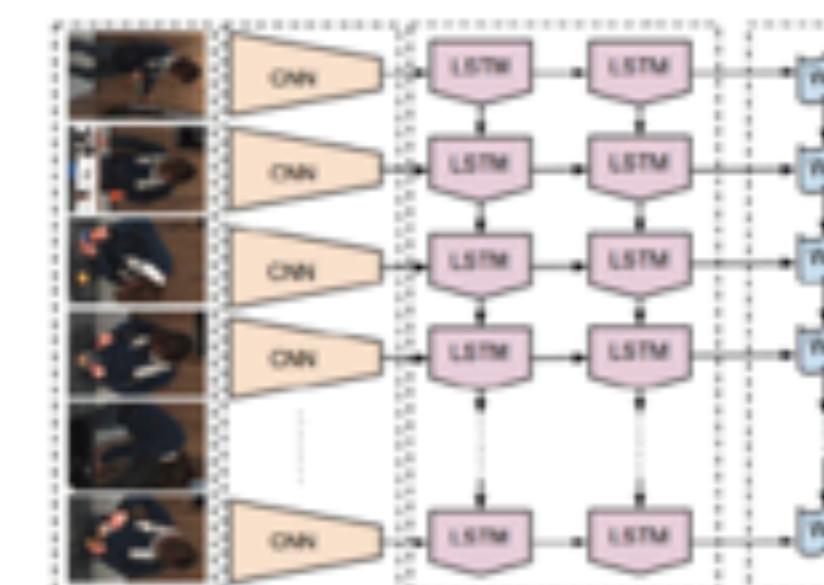
Input



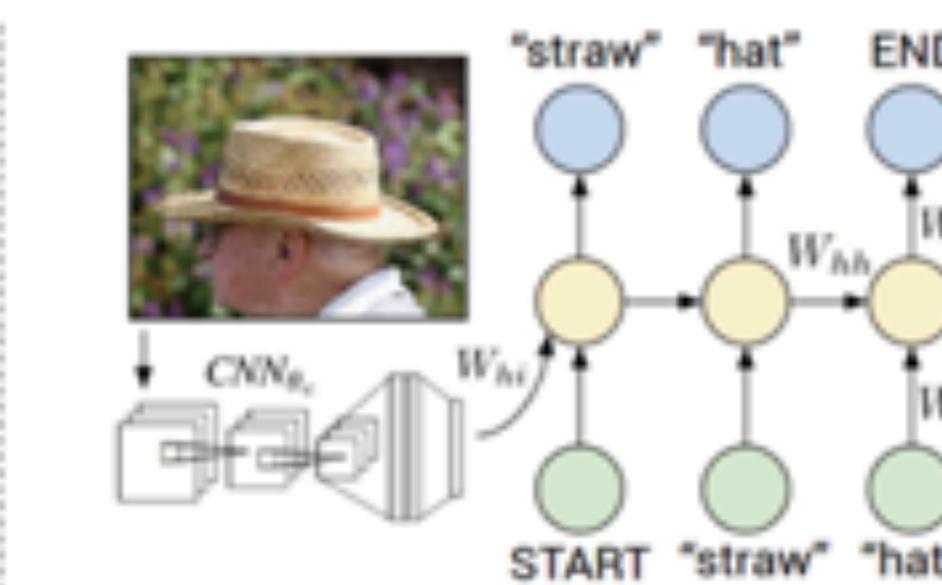
It was very popular a few years ago



Vinyals et al., 2015



Donahue et al., 2015



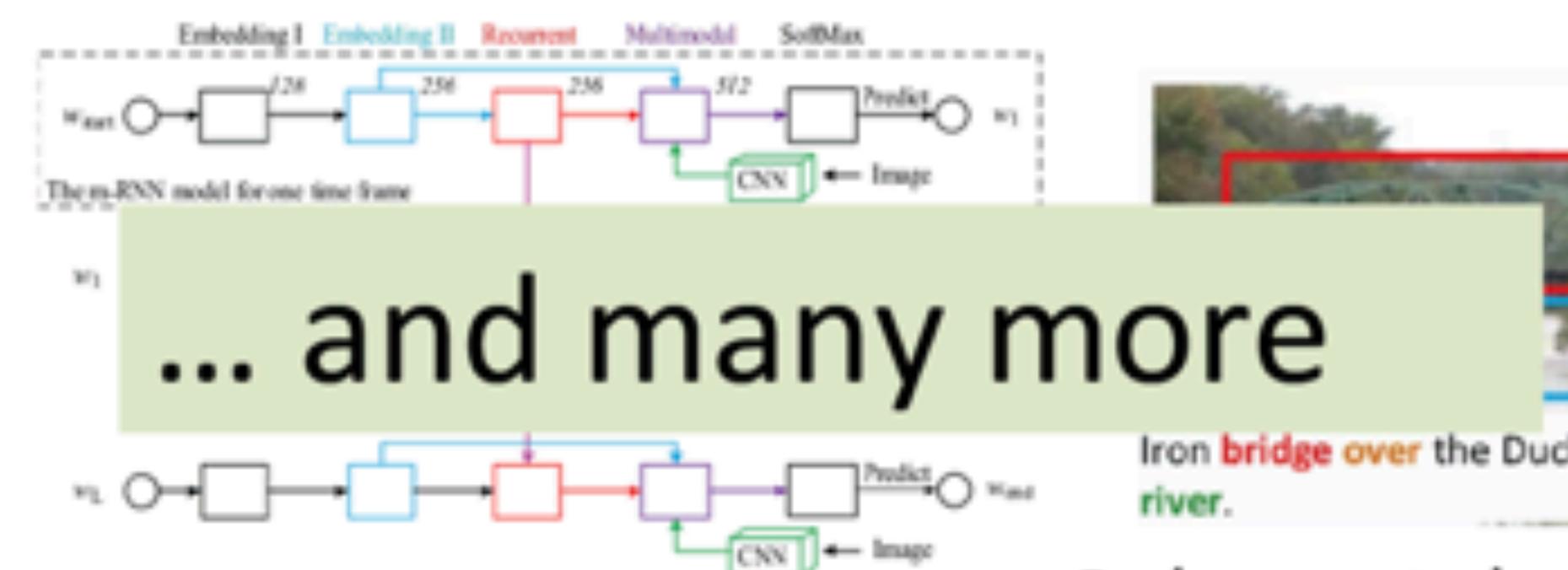
Karpathy and Fei-Fei, 2015



Hodosh et al., 2013



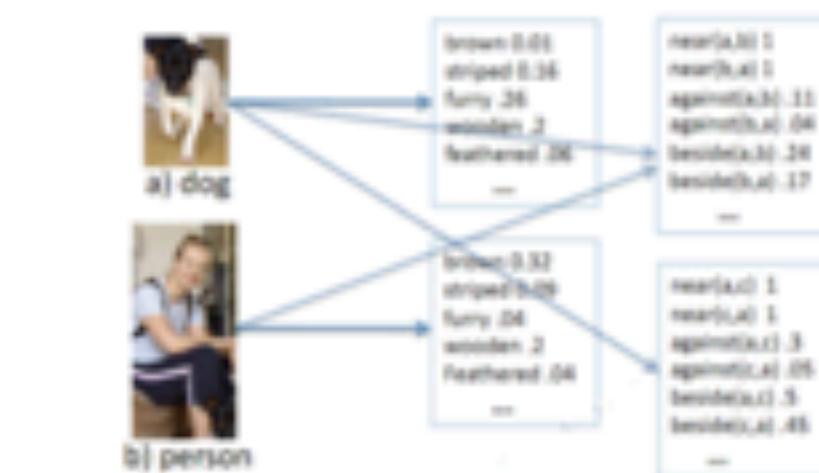
Fang et al., 2015



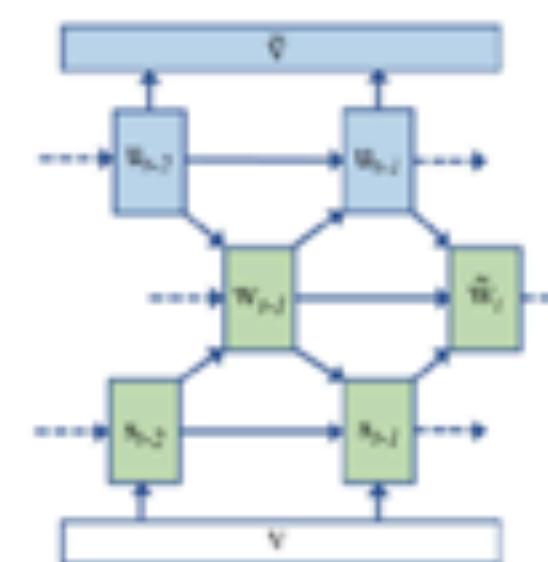
Mao et al., 2015



Ordonez et al., 2011



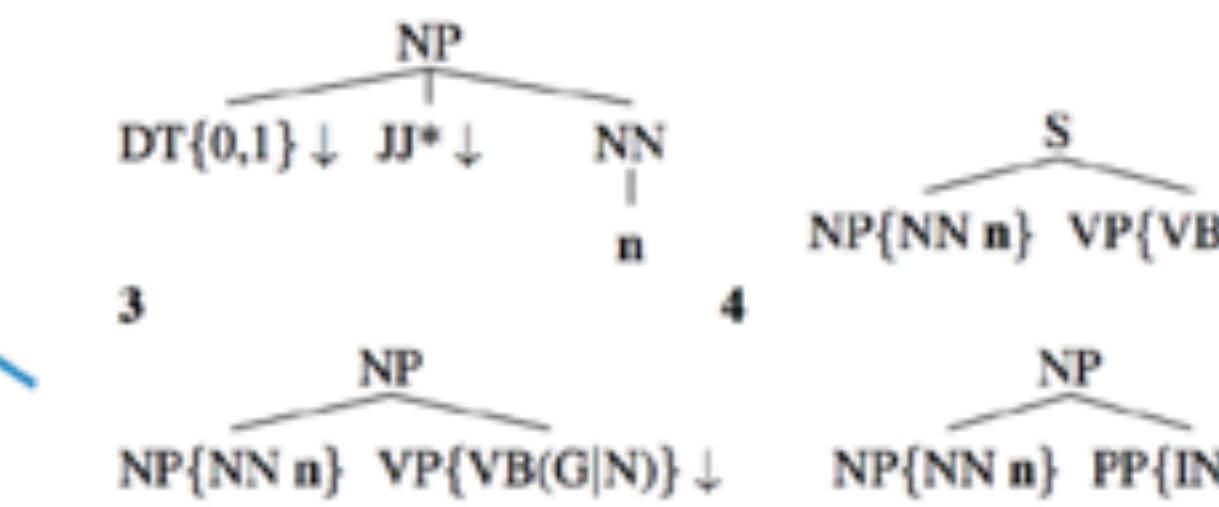
Kulkarni et al., 2011



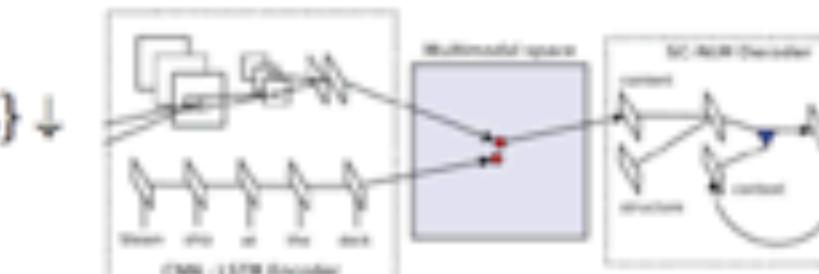
Chen and Zitnick, 2015



Farhadi et al., 2010



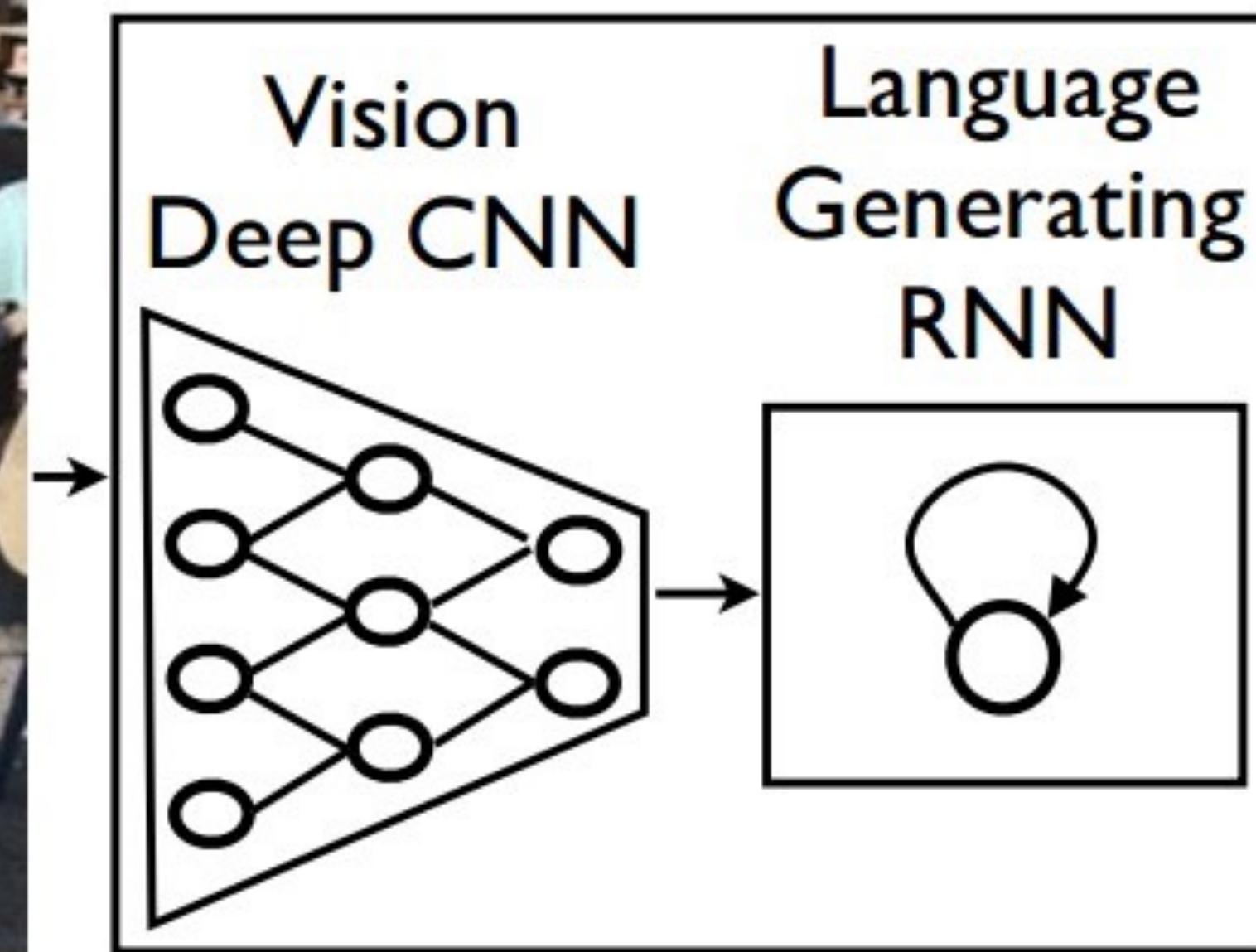
Mitchell et al., 2012



Kiros et al., 2015

Show and Tell: A Neural Image Caption Generator

[Vinyals et. al., CVPR 2015]

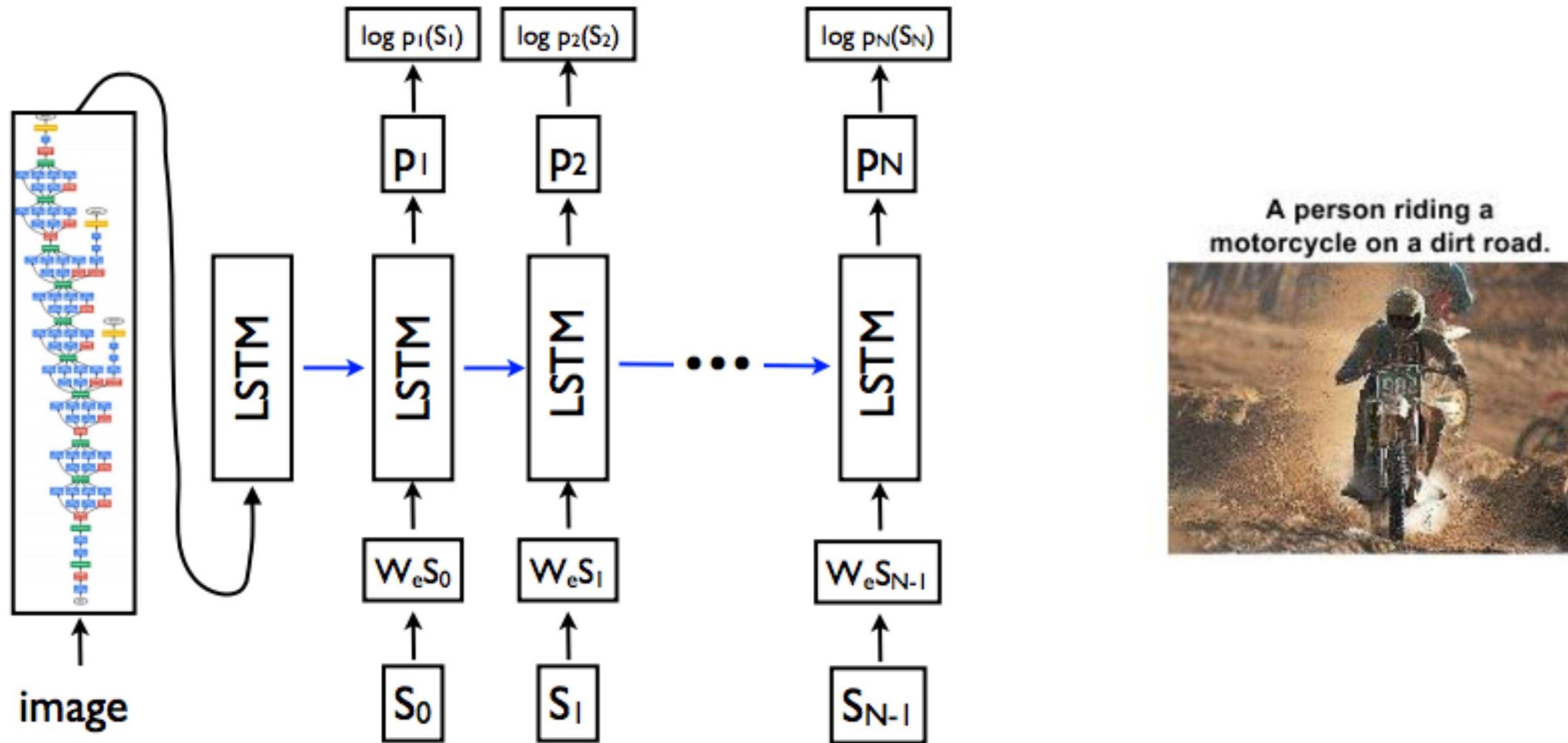


A group of people shopping at an outdoor market.

There are many vegetables at the fruit stand.

Show and Tell: A Neural Image Caption Generator

[Vinyals et. al., CVPR 2015]



A person riding a motorcycle on a dirt road.



A group of young people playing a game of frisbee.



A herd of elephants walking across a dry grass field.



Two dogs play in the grass.



Two hockey players are fighting over the puck.



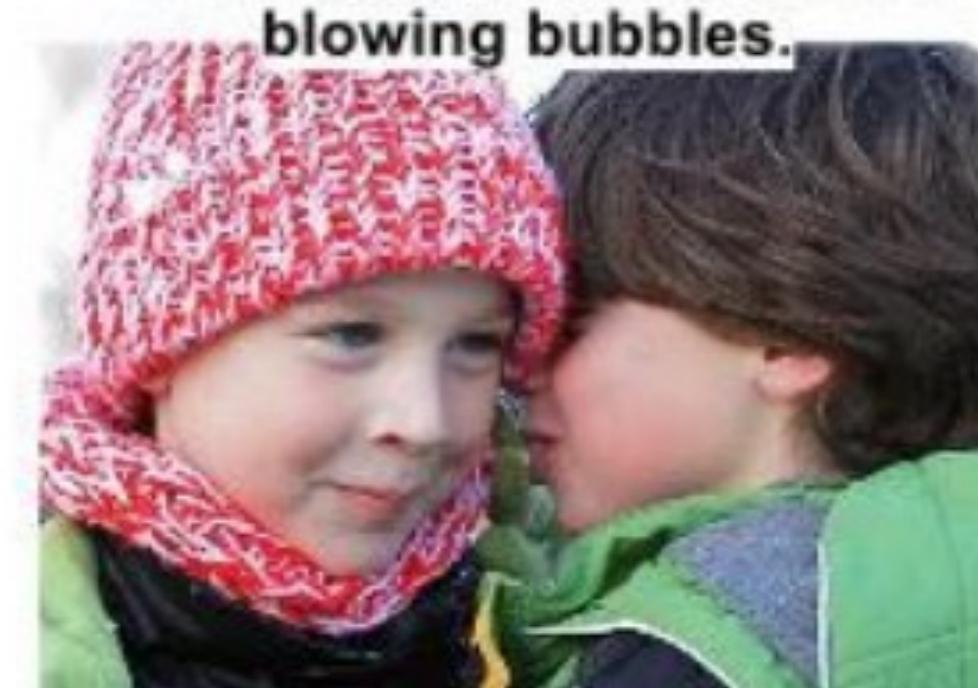
A close up of a cat laying on a couch.



A skateboarder does a trick on a ramp.



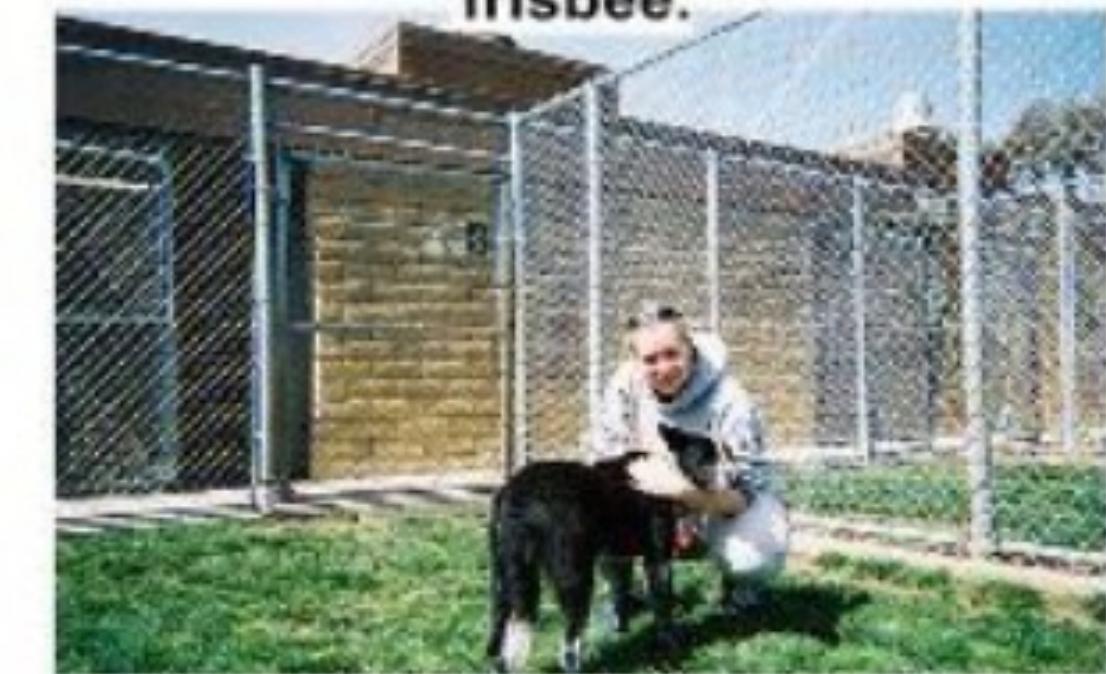
A little girl in a pink hat is blowing bubbles.



A red motorcycle parked on the side of the road.



A dog is jumping to catch a frisbee.



A refrigerator filled with lots of food and drinks.



A yellow school bus parked in a parking lot.

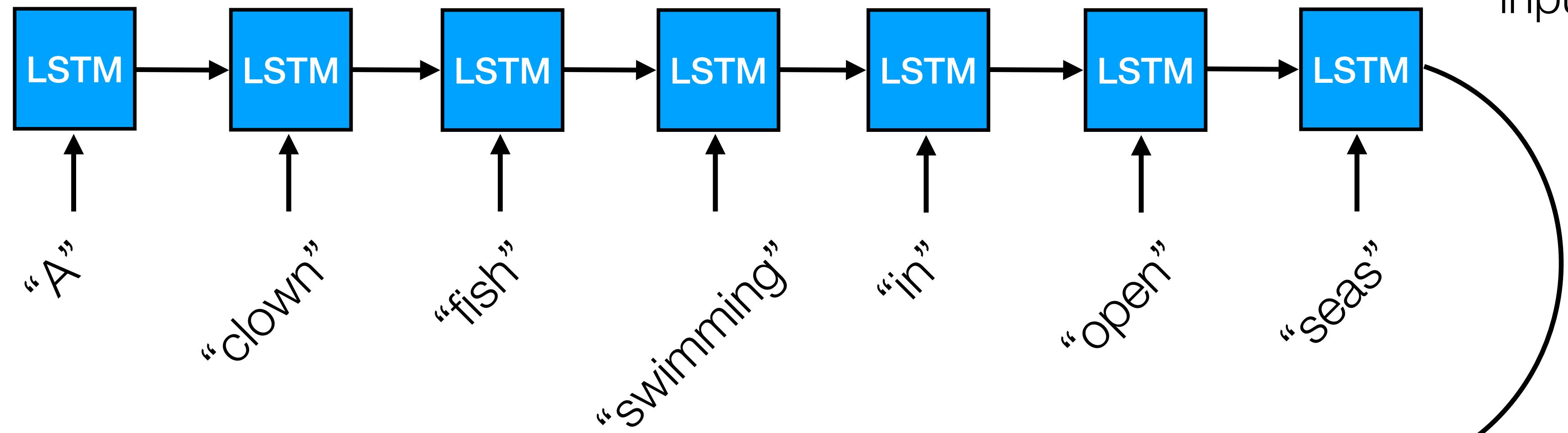


Shortcomings of recurrent models

- The recurrent state needs to **remember** a lot
- Instead of remembering: look at the input data! This idea is often implemented using **attention**.
- Example: “sequence to sequence” language translation

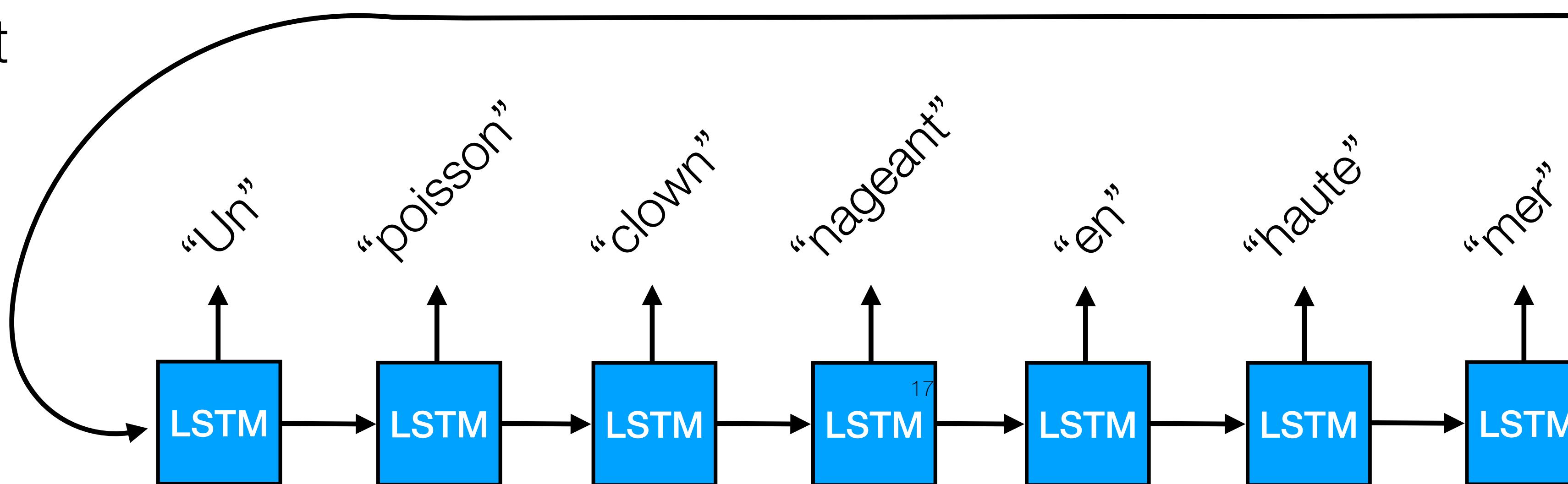
Example: Seq2Seq Translation

Input



Hidden state
remembers whole
input sentence!

Output

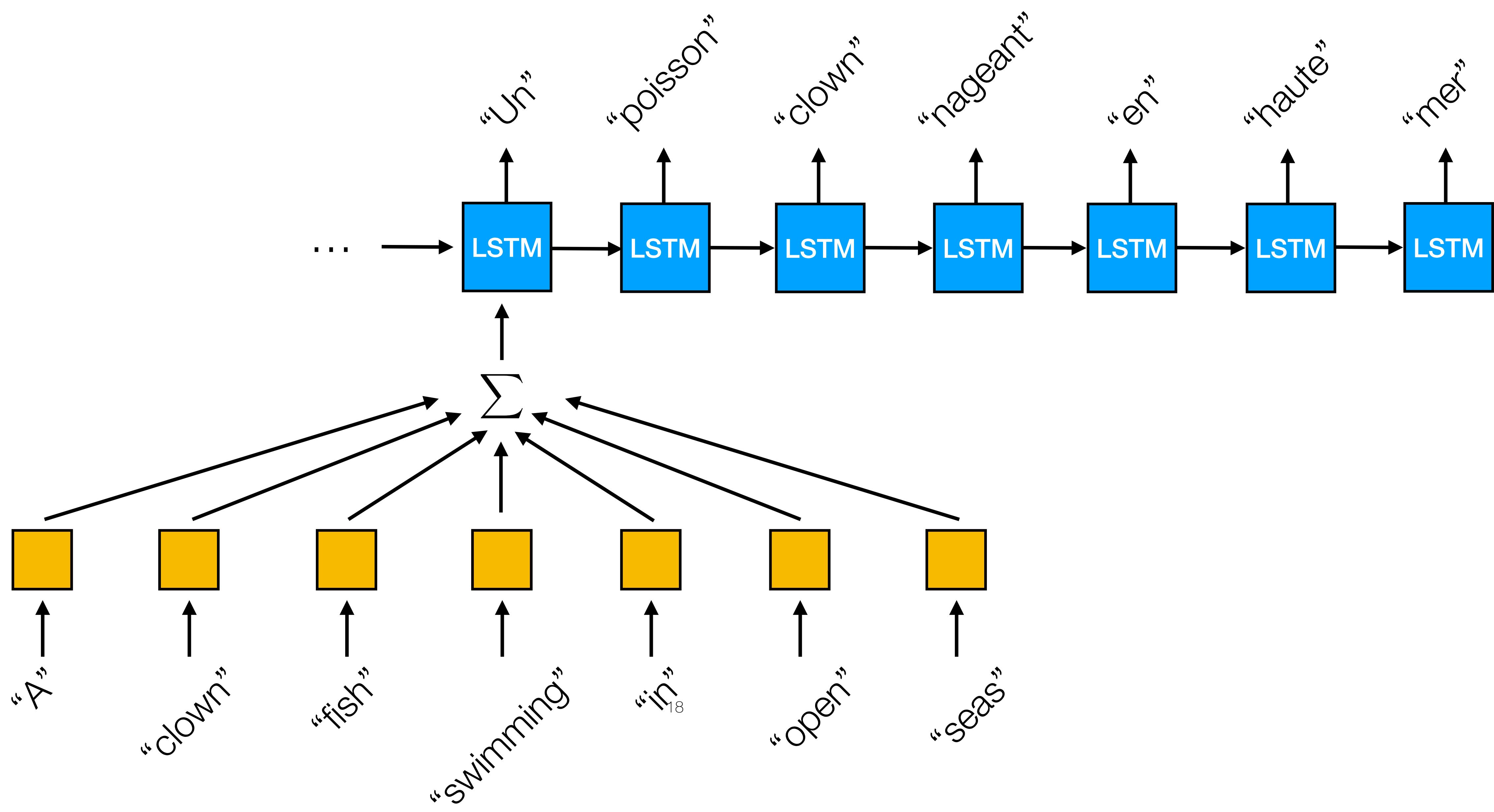


Pooling

Outputs

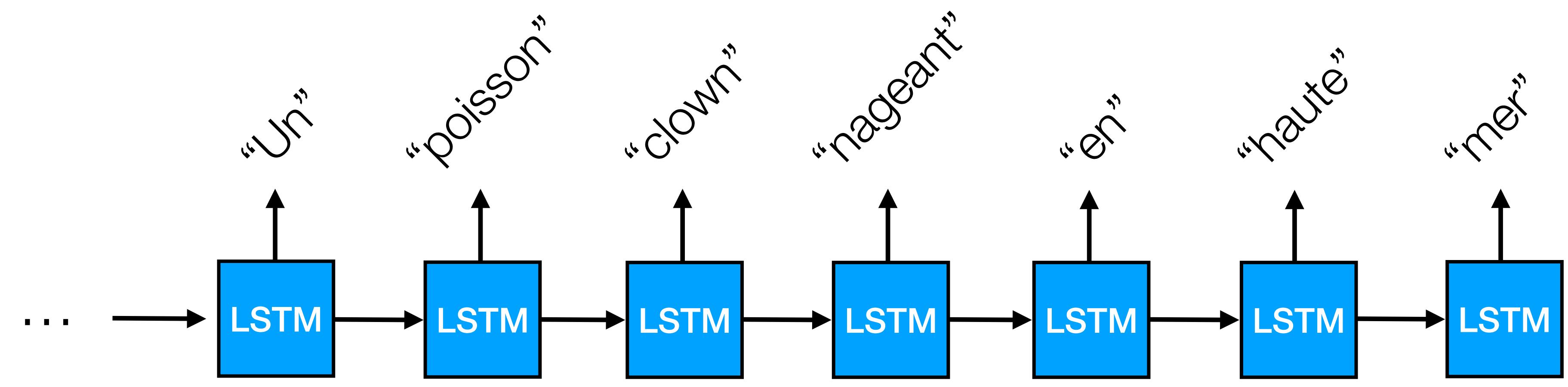
Input

Hidden

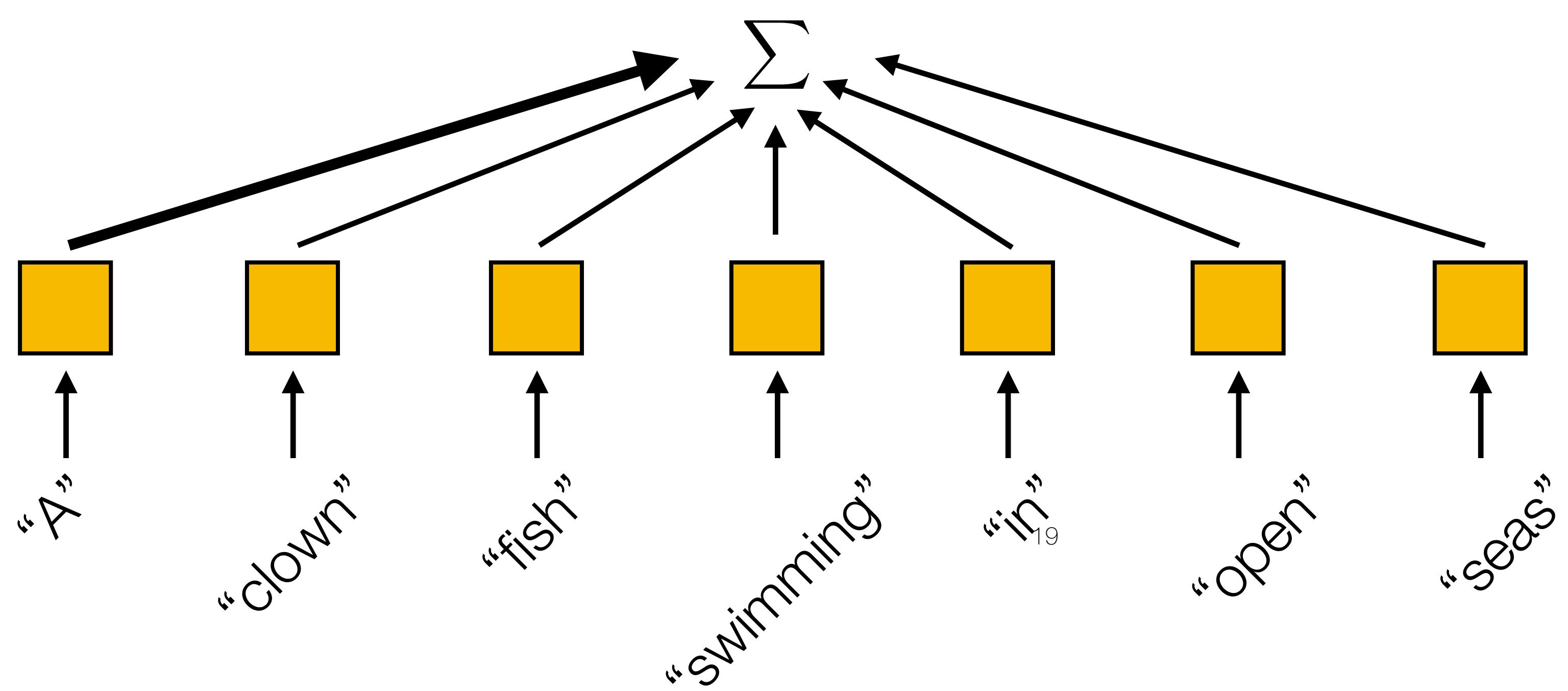


Attention

Outputs

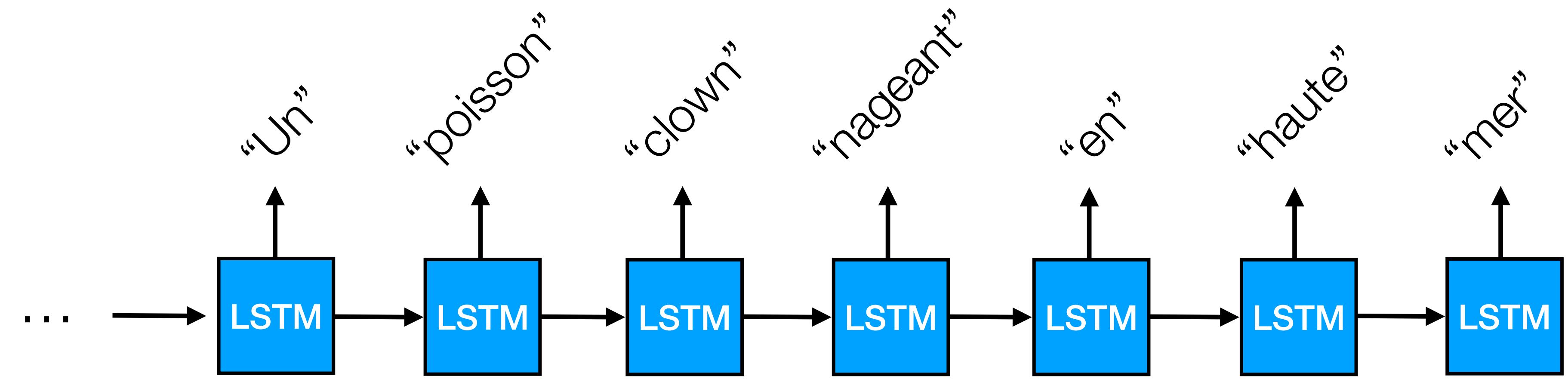


Hidden
Input

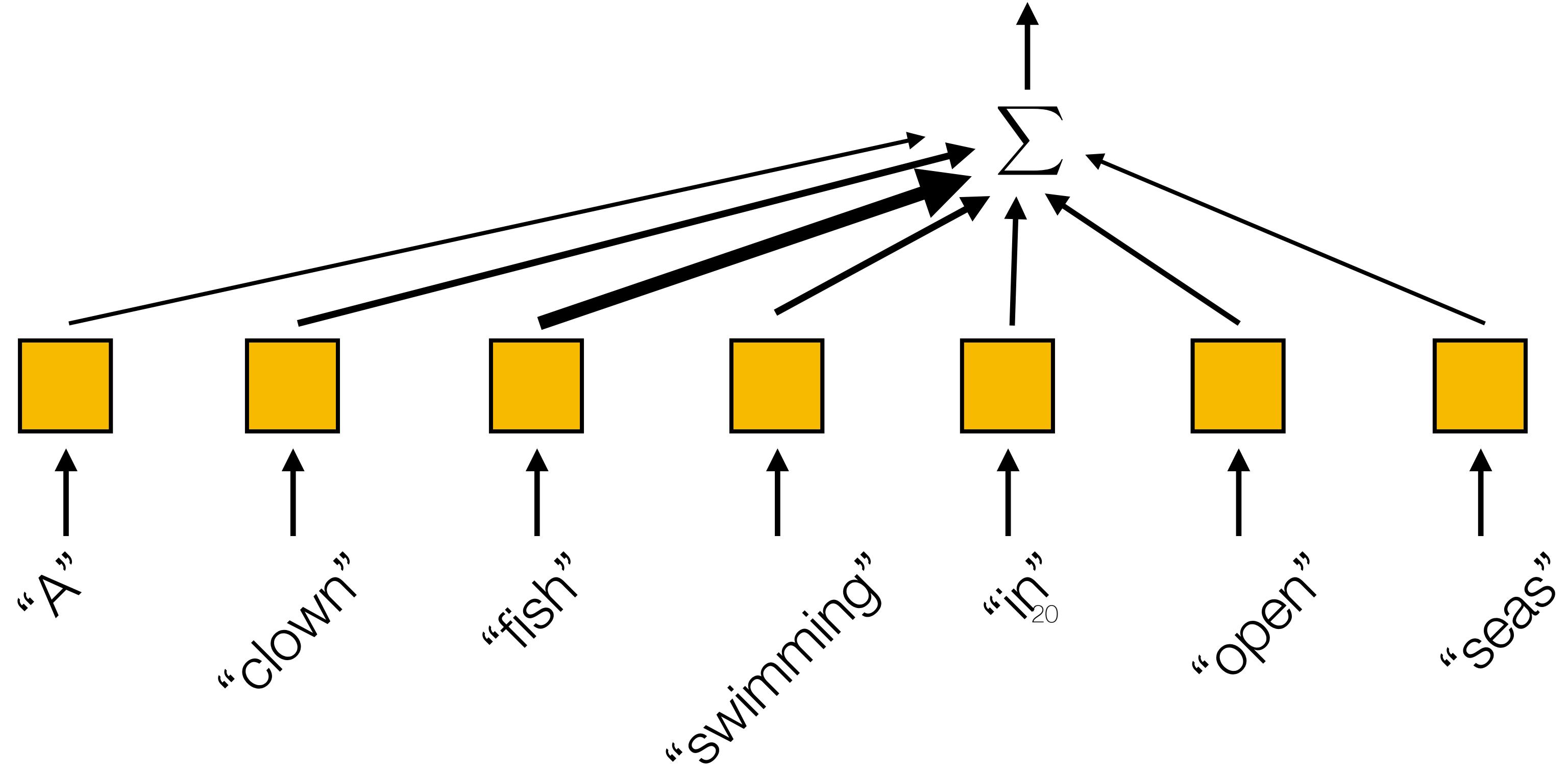


Attention

Outputs

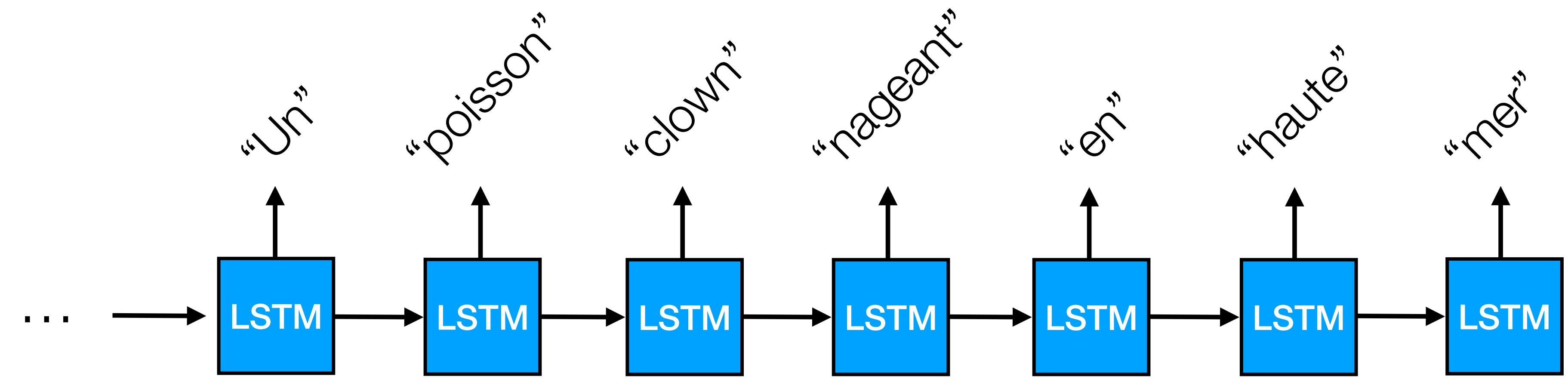


Hidden
Input

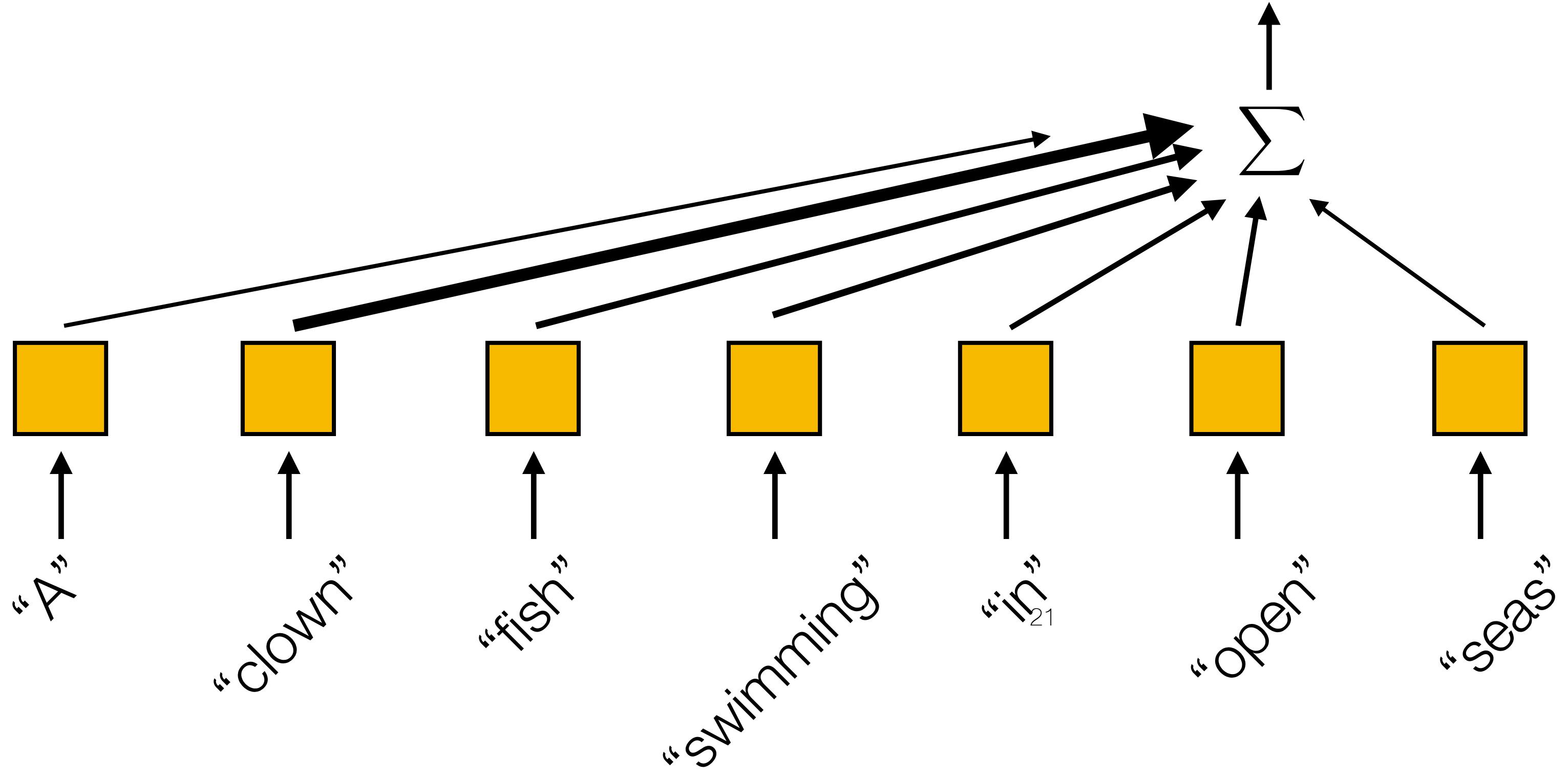


Attention

Outputs

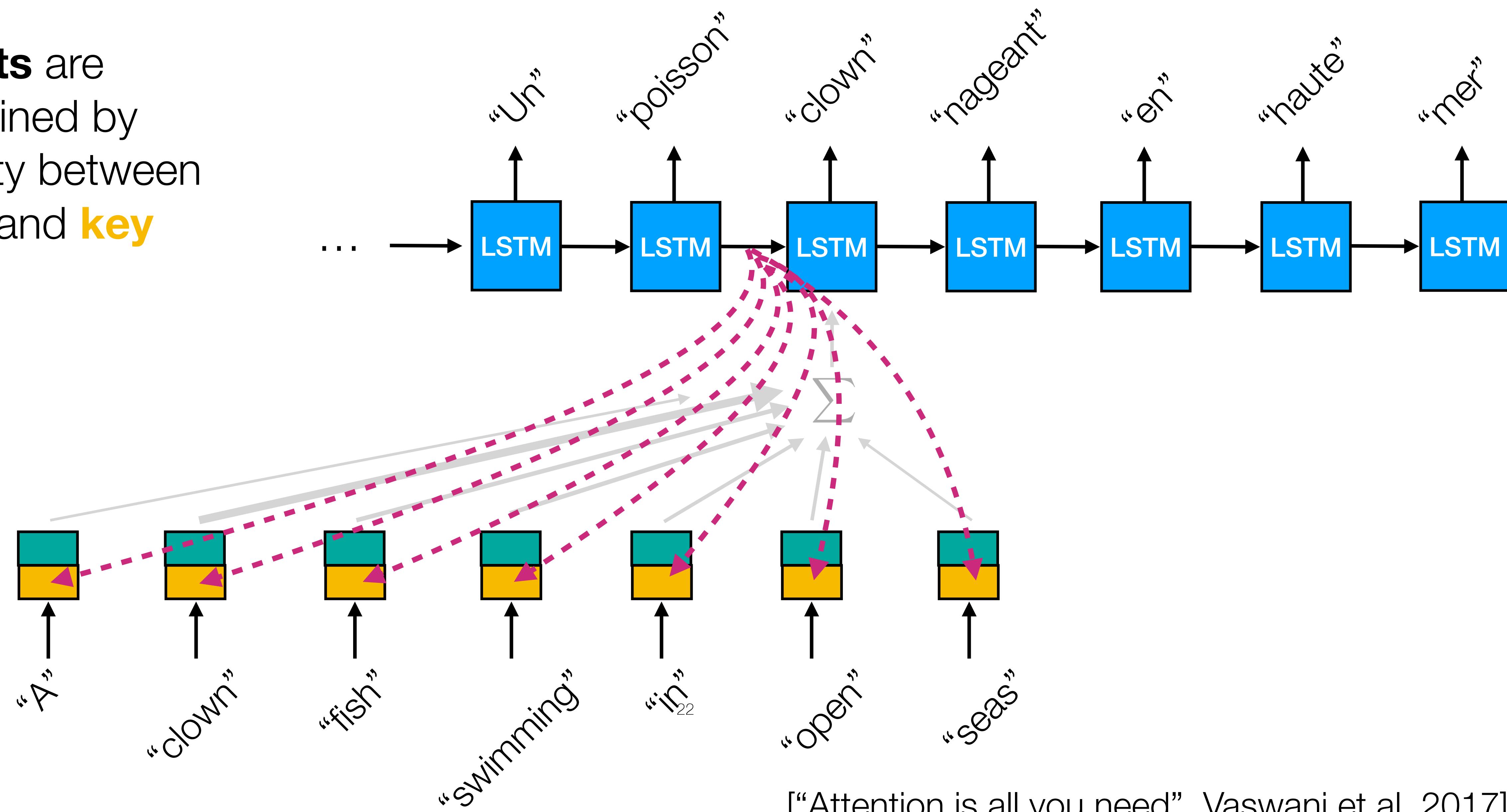


Hidden
Input



Attention

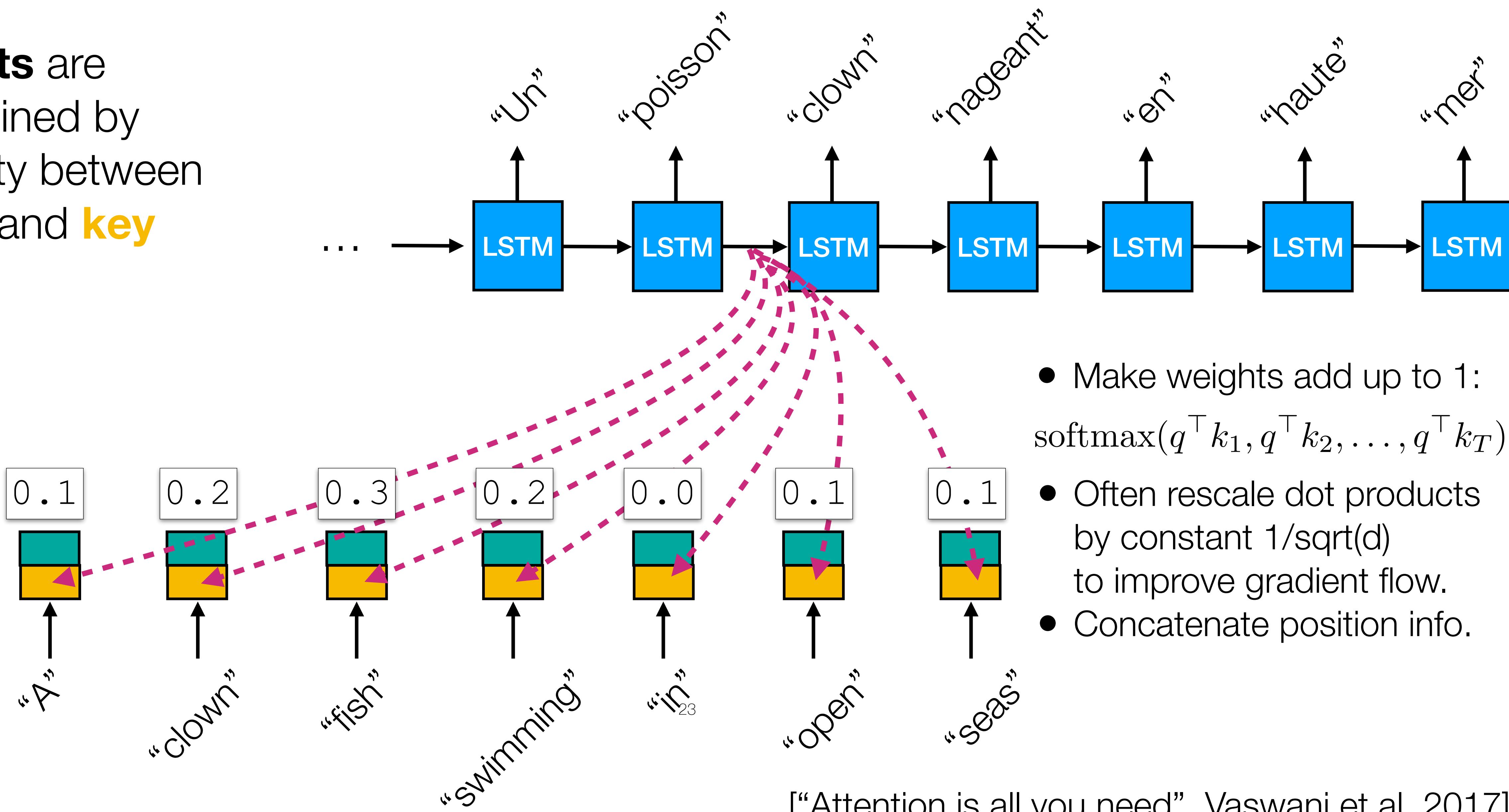
Weights are determined by similarity between **query** and **key**



[“Attention is all you need”, Vaswani et al. 2017]

Attention

Weights are determined by similarity between **query** and **key**



Attention

Weights are determined by similarity between **query** and **key**

summation is over **weights * value**

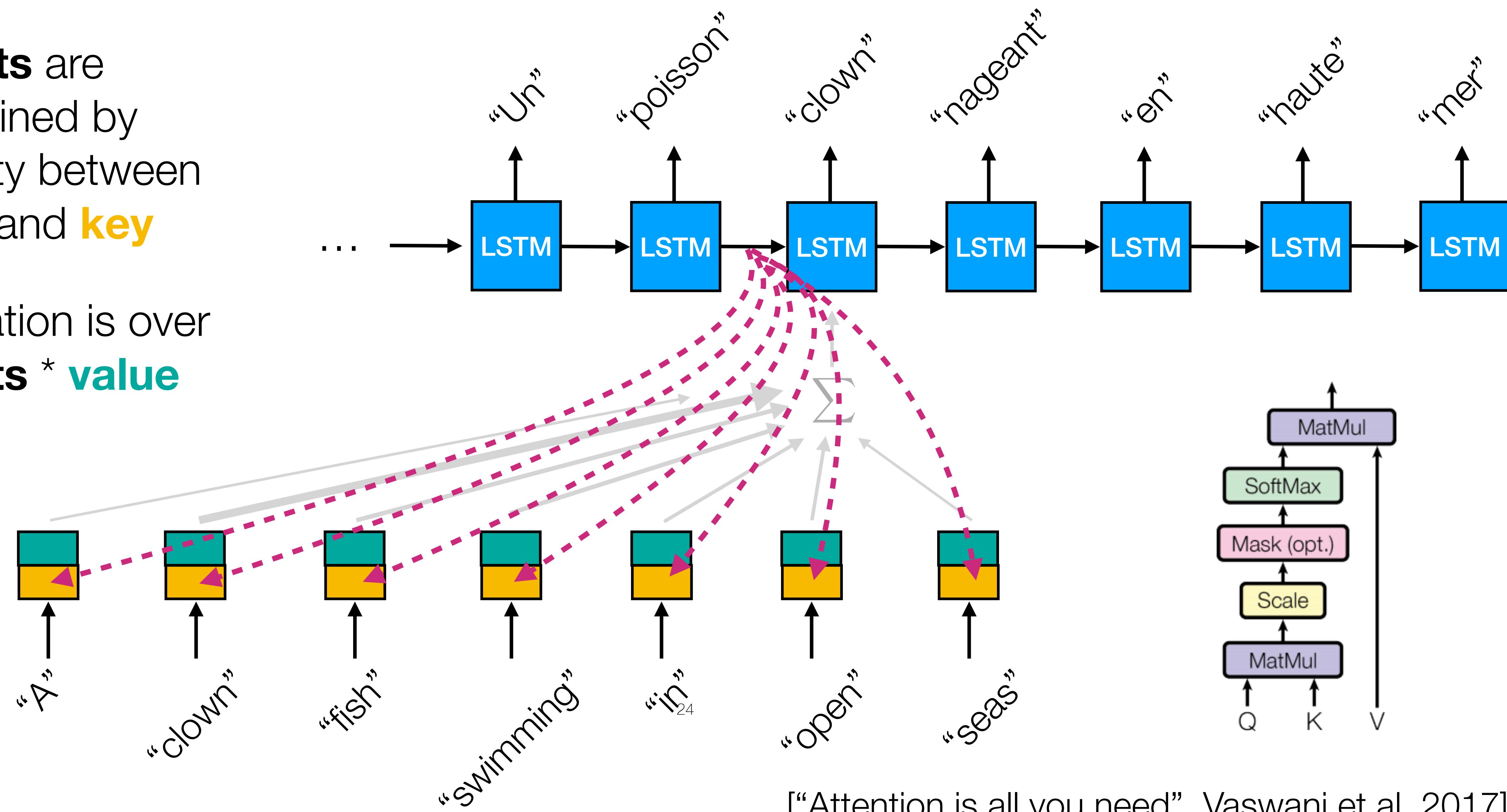


Image captioning with attention

