

Image: AI text2img generated from lecture's title. Background: Viva Magenta (Pantone 18-1750) 2024 Color of the year!

# Neural models of language



CSC401/2511 – Fall 2024

Copyright © 2024. Raed Saqr, University of Toronto

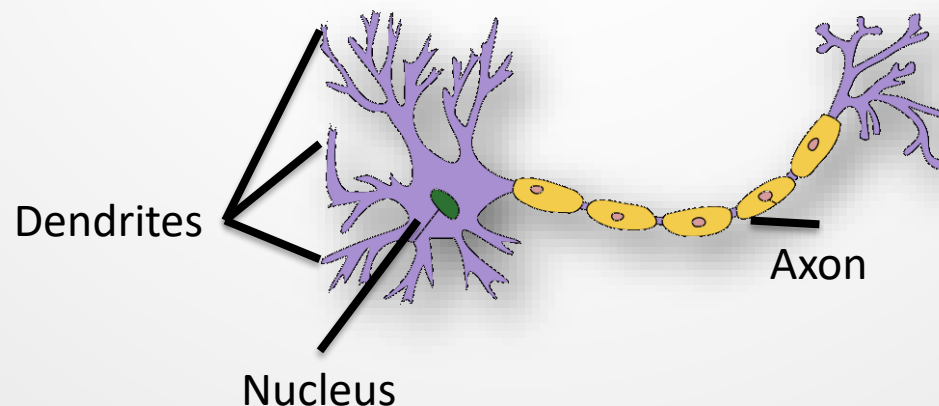
CSC401/2511 – Natural Language Computing – Fall 2024

Lecture 5

University of Toronto

# Artificial neural networks

- **Artificial neural networks (ANNs)** were loosely inspired by networks of cytoplasmic protrusions in the brain.
  - Each unit has many inputs (~**dendrites**), one output (~**axon**).
  - The **nucleus** fires (sending an electric signal along the axon) given input from other neurons.
  - ‘Learning’ was formerly thought to occur at the **synapses** that connect neurons, either by amplifying or attenuating signals.

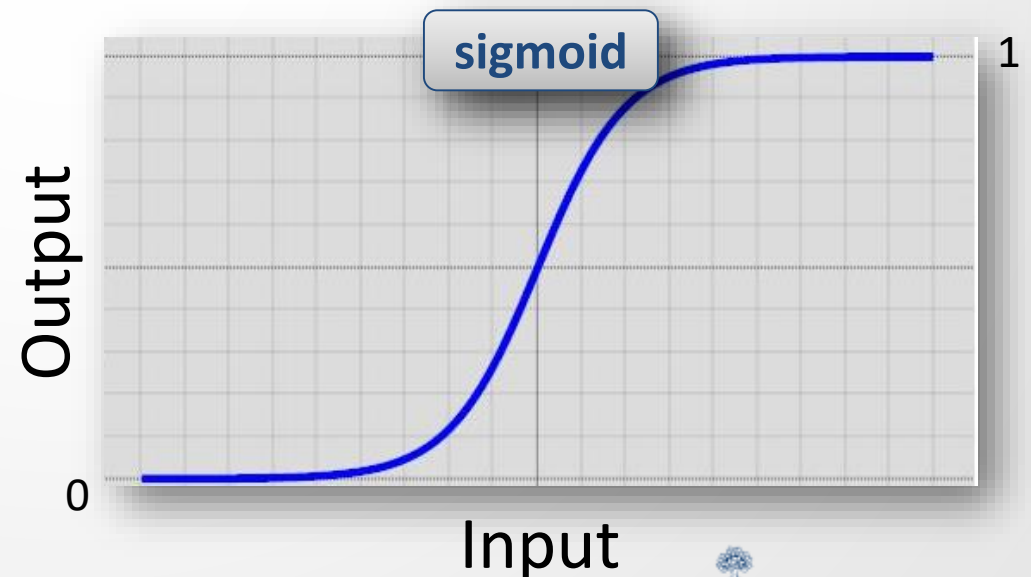
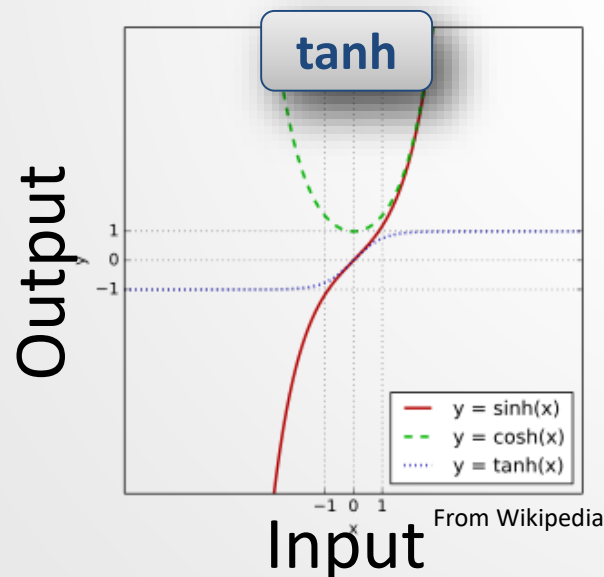


# Feed-forward output

- Output is determined by an **activation function**,  $g()$ , which **can be non-linear** (of weighted input). Activation is empirically determined, but not learned as a parameter.
- Popular activation functions include **tanh** and the **sigmoid**:

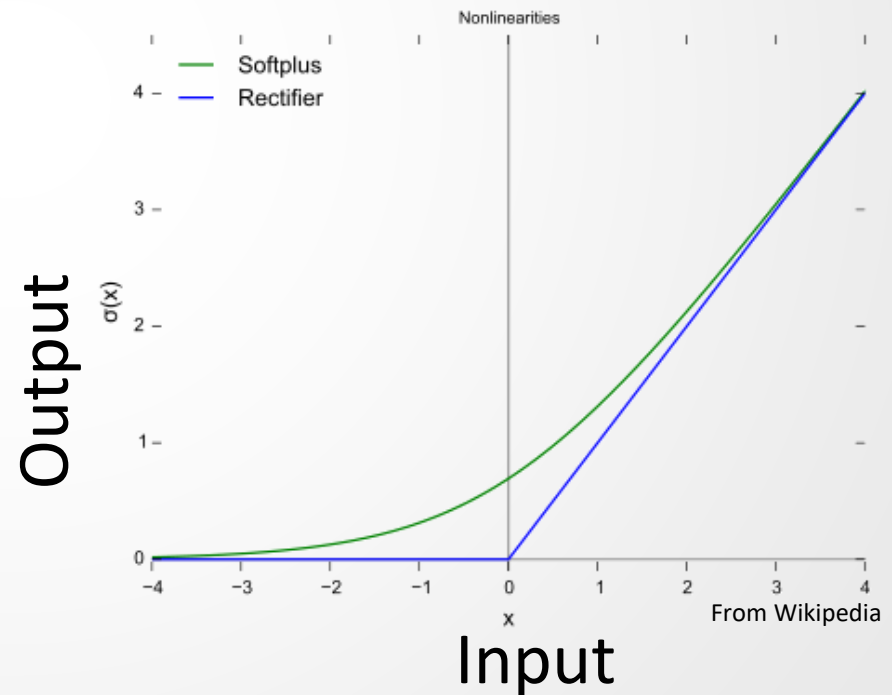
$$g(x) = \sigma(x) = \frac{1}{1 + e^{\rho x}}$$

- The sigmoid's derivative is the easily computable  $\sigma' = \sigma \cdot (1 - \sigma)$



# Rectified Linear Units (ReLUs)

- Since 2011, the **ReLU**  $S = g(x) = \max(0, x)$  has become popular.
  - More appeals to biological plausibility, but sparse activation, and reduced likelihood of vanishing gradients are very practical reasons.
- A smooth approximation is the **softplus**  $\log(1 + e^x)$ , which has a simple derivative  $1/(1 + e^{-x})$
- *Why do we care about the derivatives?*



X Glorot, A Bordes, Y Bengio (2011). Deep sparse rectifier neural networks. AISTATS.

# Parameter estimation

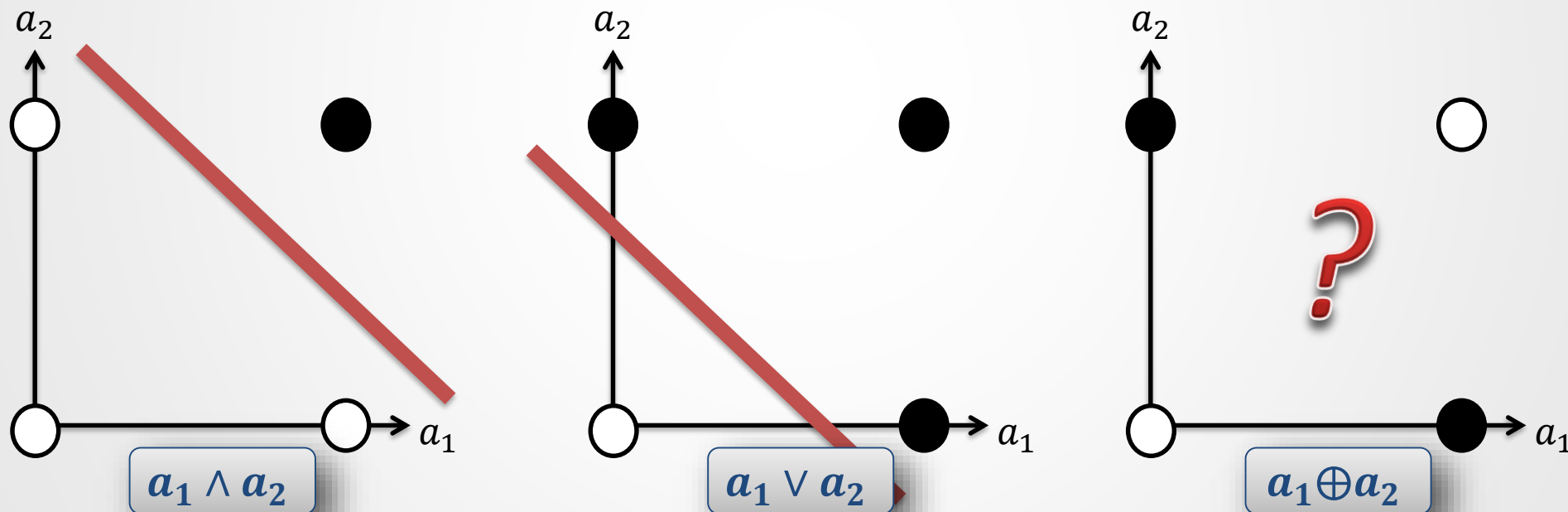
- Weights are adjusted in **proportion to the error** (i.e., the **difference** between the desired,  $y$ , and the actual output,  $S$ ).
- The **derivative**  $g'$  allows us to assign blame proportionally.
- Given a small learning rate,  $\alpha$  (e.g., 0.05), we can repeatedly adjust each of the weight parameters by

$$w_j := w_j + \alpha \cdot \sum_{i=1}^R \text{Err}_i \cdot g'(x_i) \cdot a_j[i] \quad \left. \vphantom{\sum_{i=1}^R} \right\} \begin{array}{l} \text{Assumes} \\ \text{mean-square} \\ \text{error objective} \end{array}$$

where  $\text{Err}_i = (y_i - S_i)$ , among  $R$  training examples.

# Threshold perceptrons and XOR

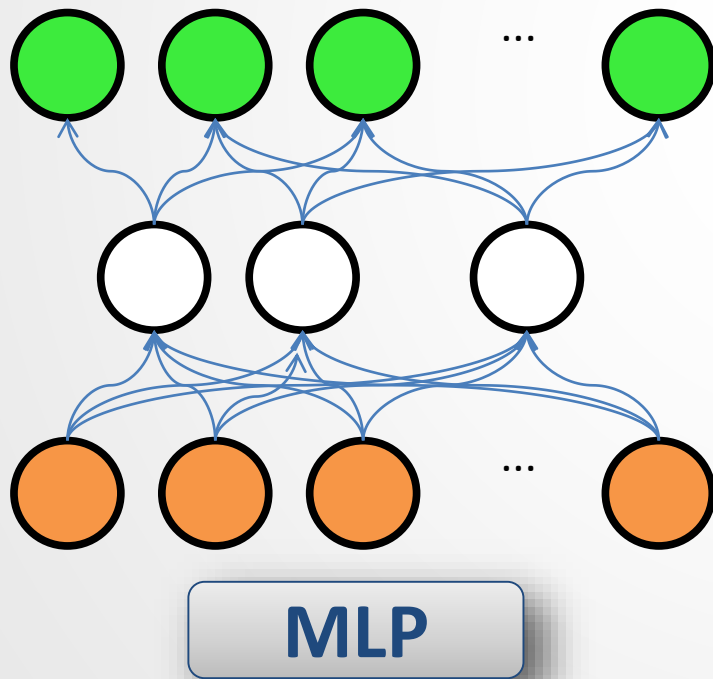
- Some relatively simple logical functions cannot be learned by threshold perceptrons (since they are not linearly separable).





# Multi-layer neural networks

- Complex functions can be represented by layers of perceptron (**multi-layer perceptron, MLPs**).



- Inputs are passed to the **input layer**.
- Activations** are propagated through **hidden layers** to the **output layer**.
- MLPs are quite **robust to noise**. Sometimes, we even add noise.

# Parameter estimation

We have many options. Gradient descent is popular.  
Given  $T$  tokens of training data, optimize objective:

$$J(\theta) = \frac{1}{T} \sum_{t=1}^T \sum_{-c < j < c, j \neq 0} \log P(w_{t+j} | w_t)$$

And we want to update vectors  $V_{w_{t+j}}$  then  $v_{w_t}$  within  $\theta$

$$\theta^{(new)} = \theta^{(old)} - \alpha \nabla_{\theta} J(\theta)$$

So, we'll need to take the derivative of the (log of the) softmax function:

$$P(w_o | w_i) = \frac{\exp(V_{w_o}^T v_{w_i})}{\sum_{w=1}^W \exp(V_w^T v_{w_i})}$$

Where  $v_w$  is the 'input' vector for word  $w$ ,  
and  $V_w$  is the 'output' vector for word  $w$ ,



# Parameter estimation

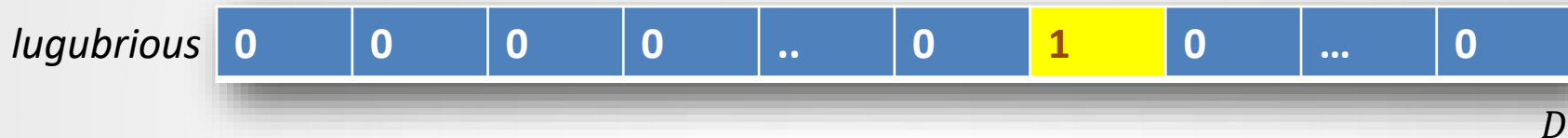
We need the derivative of the (log of the) softmax function:

$$\begin{aligned}\frac{\delta}{\delta v_{w_t}} \log P(w_{t+j}|w_t) &= \frac{\delta}{\delta v_{w_t}} \log \frac{\exp(V_{w_{t+j}}^\top v_{w_t})}{\sum_{w=1}^W \exp(V_w^\top v_{w_t})} \\&= \frac{\delta}{\delta v_{w_t}} \left[ \log \exp(V_{w_{t+j}}^\top v_{w_t}) - \log \sum_{w=1}^W \exp(V_w^\top v_{w_t}) \right] \\&= V_{w_{t+j}} - \left[ \frac{\delta}{\delta v_{w_t}} \log \sum_{w=1}^W \exp(V_w^\top v_{w_t}) \right] \\&\quad \text{[apply the chain rule } \frac{\delta f}{\delta v_{w_t}} = \frac{\delta f}{\delta z} \frac{\delta z}{\delta v_{w_t}} \text{]} \\&= V_{w_{t+j}} - \sum_{w=1}^W p(w|w_t) V_w\end{aligned}$$

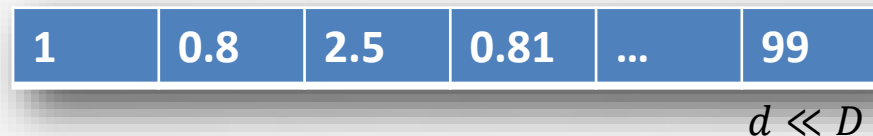
More details: <http://arxiv.org/pdf/1411.2738.pdf>

# Words

- Given a corpus with  $D$  (e.g.,  $= 100K$ ) unique words, the **one-hot approach** uniquely assigns **each word** an index in  $D$ -dimensional vectors ('one-hot' representation).



- In psychology, **word-feature representations** assign **features** to each index in a much denser vector.
  - E.g., concept-based features 'cheerful', 'emotional-tone'.



- Neither of these is learned.

# Using word representations

Without a latent space,

lugubrious =  $[0,0,0, \dots, 0, 1, 0, \dots, 0]$ , &

sad =  $[0,0,0, \dots, 0, 0, 1, \dots, 0]$  so

Similarity =  $\cos(x, y) = 0.0$

EMBEDDING

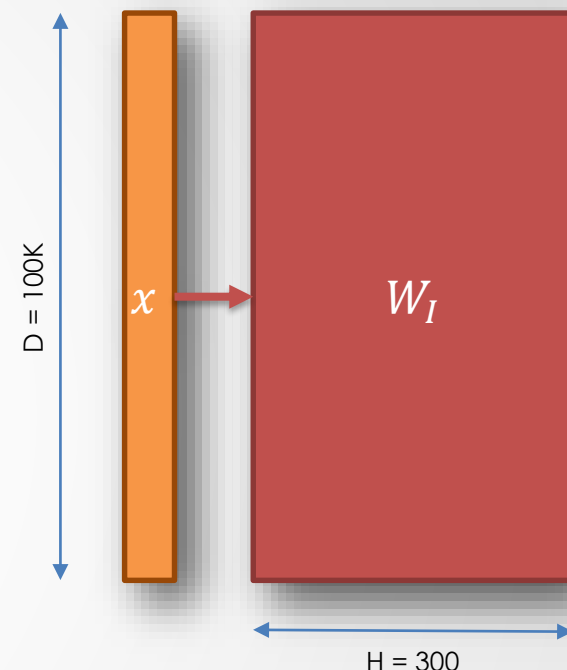
$$v_w = x^T W_I$$

In latent space,

lugubrious =  $[0.8, 0.69, 0.4, \dots, 0.05]_H$ , &

sad =  $[0.9, 0.7, 0.43, \dots, 0.05]_H$  so

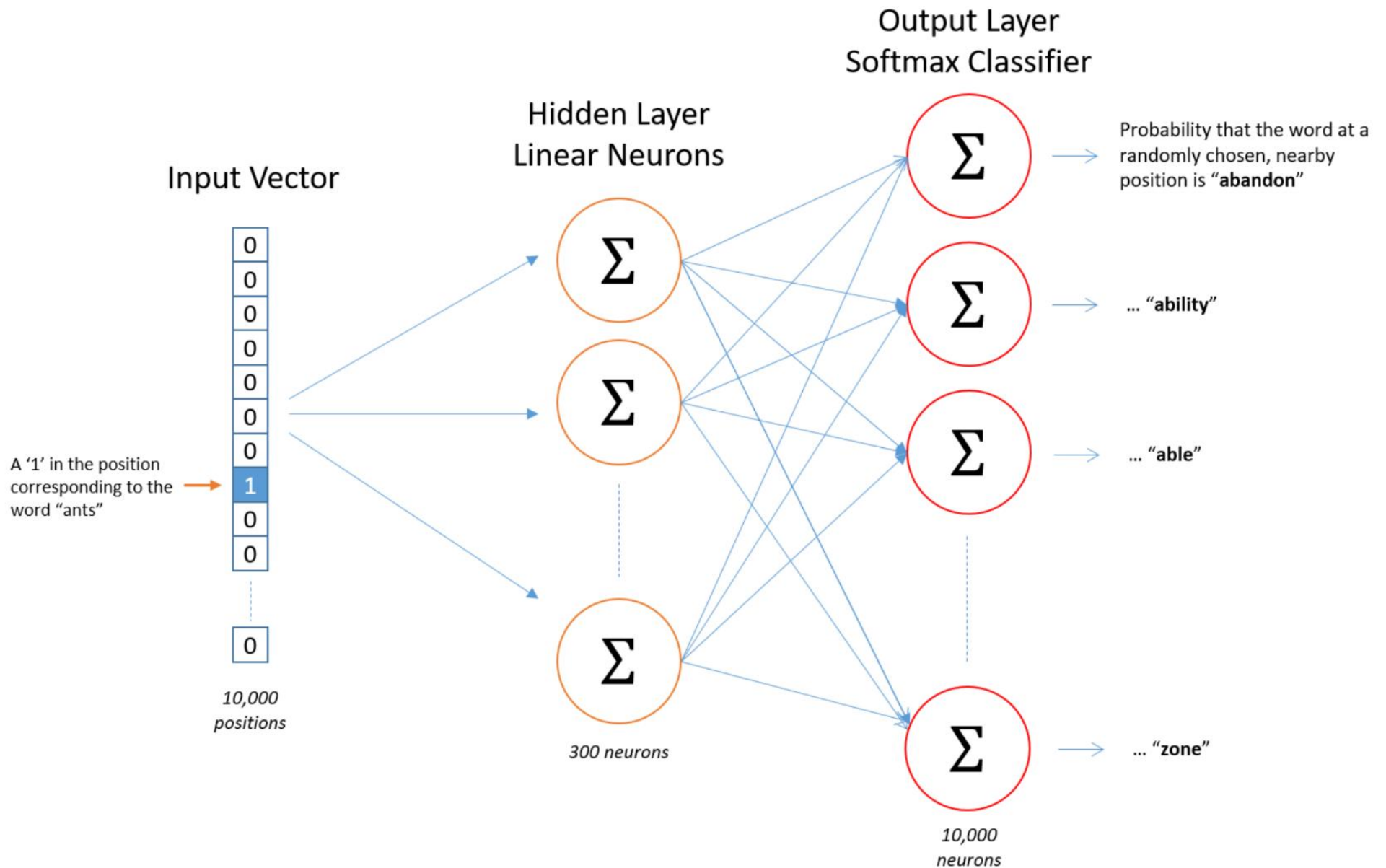
Similarity =  $\cos(x, y) = 0.9$



Reminder:

$$\cos(u, v) = \frac{u \cdot v}{||u|| \times ||v||}$$

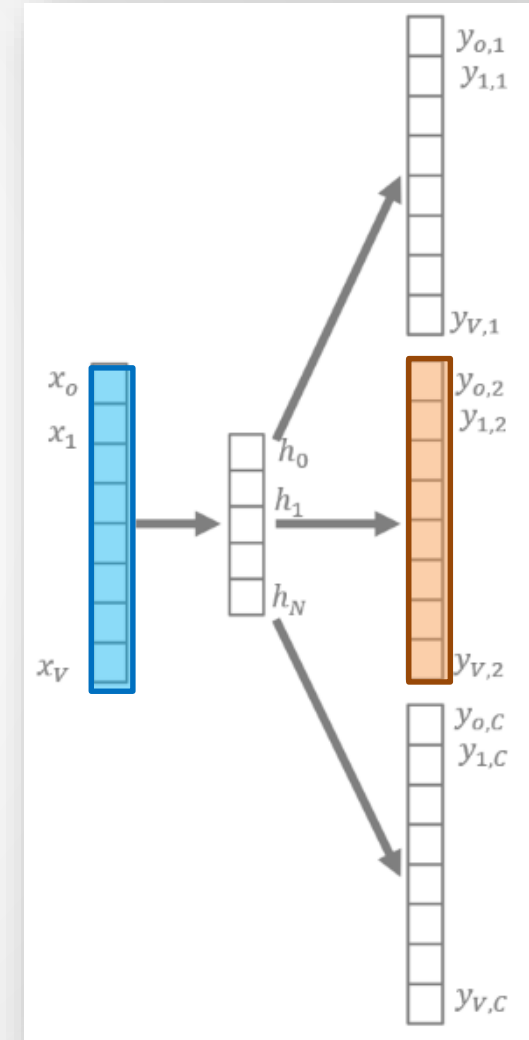
# word2vec training regimen



# Skip-grams with negative sampling

- Most word types do not appear together within a small window. The default process does not know this.
  - Also, not all that efficient – would be nice not to update  $H \times D$  weights
  - **Contrastive learning:** push away from negative examples as well as towards positives.
- For the observed (true) pair (*lugubrious*, *sadness*), only the output neuron for *sadness* should be 1, and all  $D - 1$  others should be 0.
- Mathematical Intuition:

$$P(w_o | w_c) = \frac{\exp(v_o^T v_c)}{\sum_{w=1}^D \exp(v_w^T v_c)} \left. \vphantom{\sum_{w=1}^D} \right\} \text{Computationally infeasible}$$



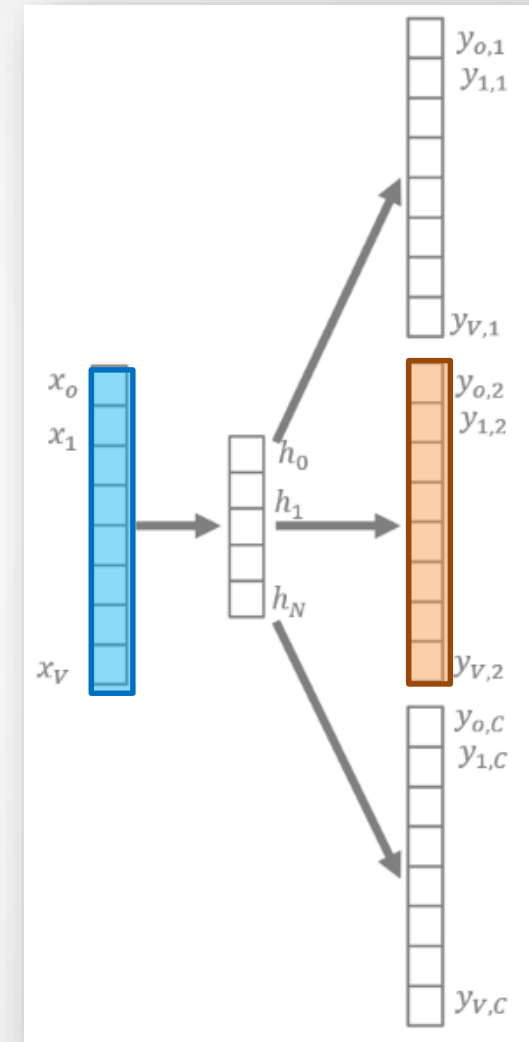
# Skip-grams with negative sampling

- We want to **maximize** the association of **observed** (positive) contexts:

*lugubrious*   *sad*  
*lugubrious*   *feeling*  
*lugubrious*   *tired*

- We want to **minimize** the association of **'hallucinated'** contexts:

*lugubrious*   *happy*  
*lugubrious*   *roof*  
*lugubrious*   *truth*



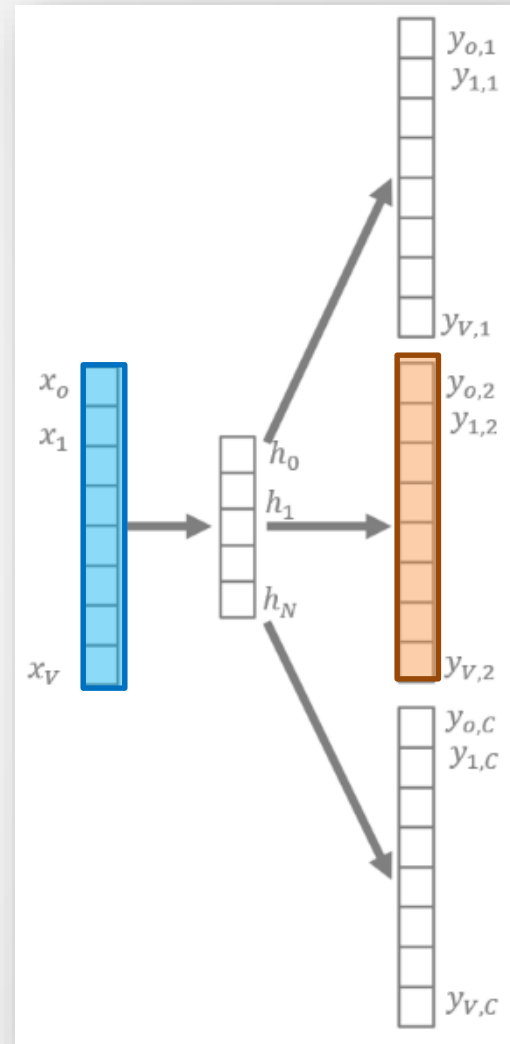
# Skip-grams with negative sampling

- Choose a small number  $k$  of ‘negative’ words, then update the weights only for all the **positive** and the  $k$  **negative** words.
  - $5 \leq k \leq 20$  can work in practice for fewer data.
  - For  $D = 100K$ , we only update 0.006% of the weights in the output layer.

$$J(\theta) = \log \sigma(v_o^T v_c) + \sum_{i=1}^k \mathbb{E}_{i \sim \underbrace{P(w)}_{\text{Unigram dist.}}} [\log \sigma(-v_i^T v_c)]$$

- Mimno and Thompson (2017) choose the top  $k$  words by **modified unigram probability**:

$$\underline{P^*(w_{t+1})} = \frac{C(w_{t+1})^{\frac{3}{4}}}{\sum_w C(w)^{\frac{3}{4}}}$$



Mimno, D., & Thompson, L. (2017). The strange geometry of skip-gram with negative sampling. *EMNLP 2017*. [\[link\]](#)



# RECURRENT NEURAL NETWORKS

# Statistical language models

- Probability is conditioned on (window of)  $n$  previous words\*
- A necessary (but incorrect) Markov assumption: each observation only factors through a **short linear history** of length  $L$ .

$$P(w_n | w_{1:(n-1)}) \approx P(w_n | w_{(n-L+1):(n-1)})$$

- Probabilities are estimated by computing unigrams and bigrams

$$P(s) = \prod_{i=1}^t P(w_i | w_{i-1})$$

bigram

$$P(s) = \prod_{i=2}^t P(w_i | w_{i-2} w_{i-1})$$

trigram

\*From Lecture 2

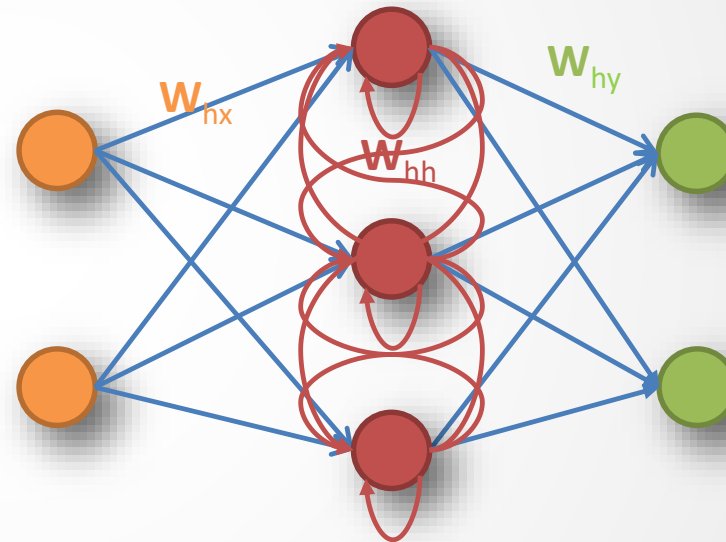
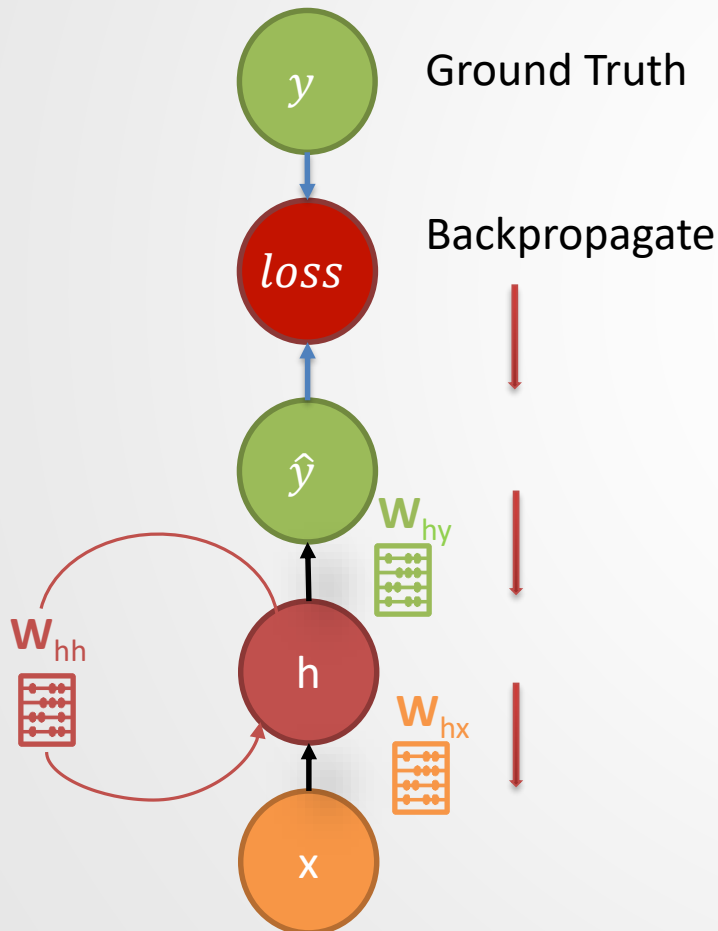
# Statistical language models

- Using higher n-gram counts (with smoothing) improves performance\*
- *RNN intuition:*
  - Use as much history as we need to use
  - Use the **same set of weight** parameters for each word (or across all time steps) to keep the size of the network down
  - Memory requirement now scales with number of words

\*From Lecture 2

# Recurrent neural networks (RNNs)

- An RNN has **feedback** connections in its structure so that it *'remembers'* previous states, when reading a sequence.

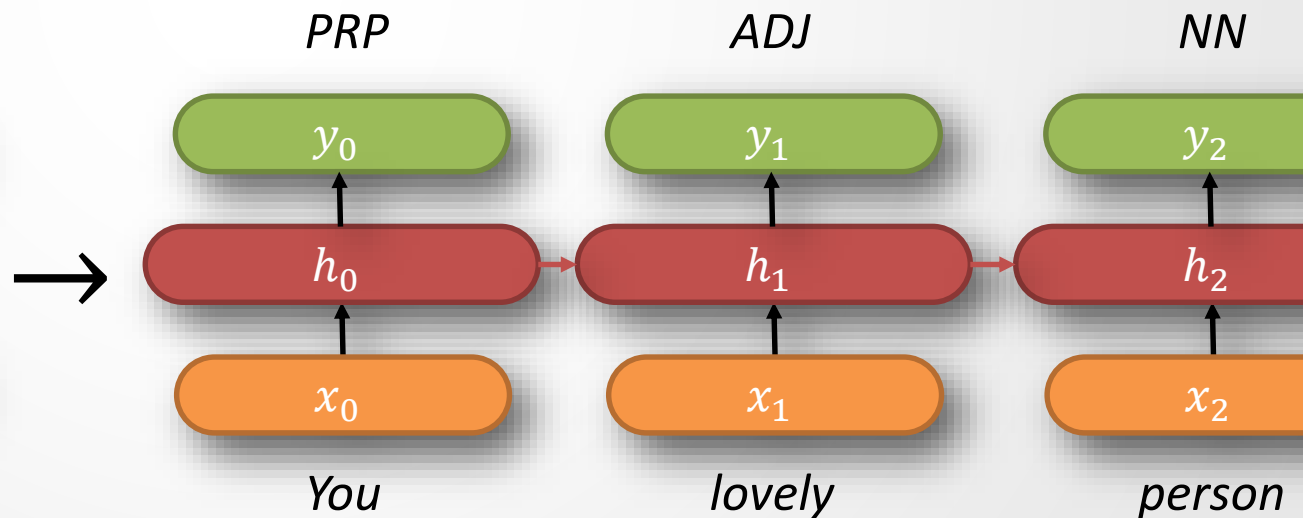
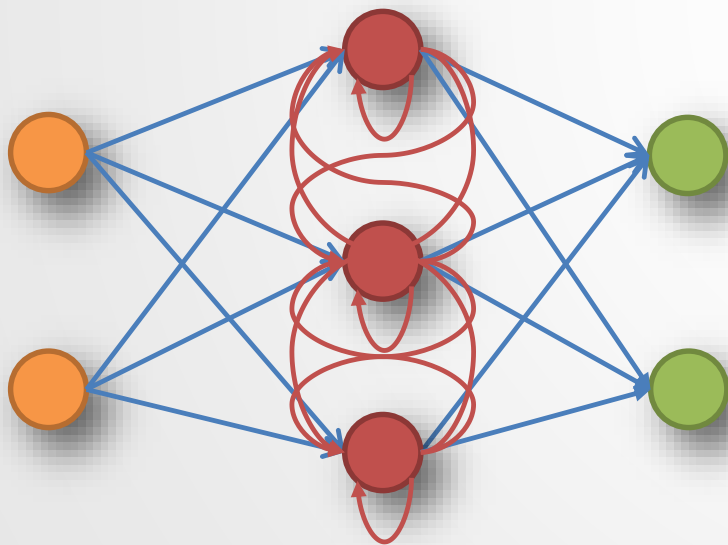


Elman network feed hidden units back

Jordan network (not shown) feed output units back

# RNNs: Unrolling the $h_i$

- Copies of the same network can be applied (i.e., **unrolled**) at each point in a time series.
- Now we can use an approximation: backpropagation through time (BPTT).

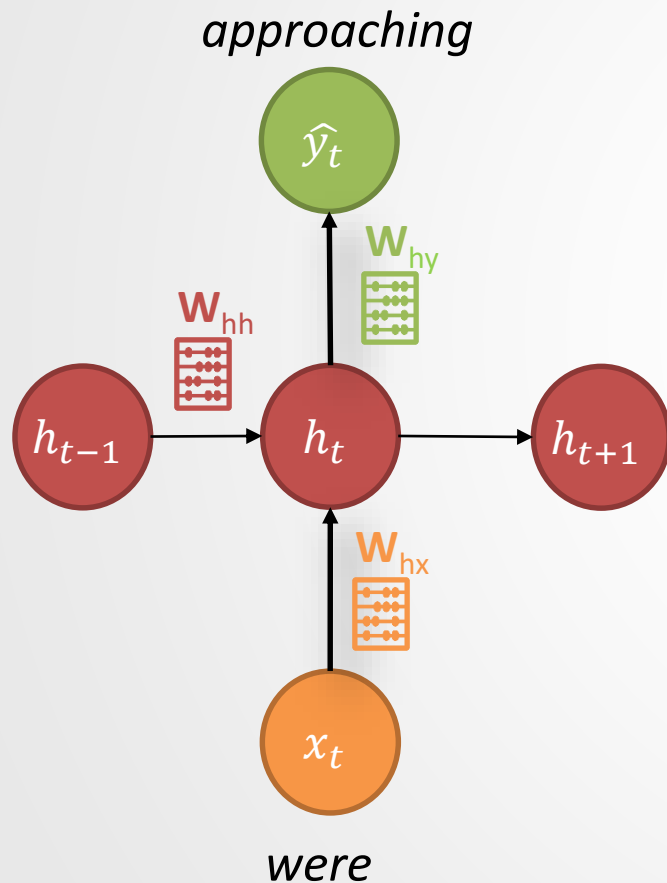


$$h_t = g(W_I[h_{t-1}; x] + c)$$
$$y_t = W_O h_t + b$$

# RNNs: One time step snapshot

*Two riders .. approaching .. horses.*

- Given a list of word vectors  $\mathbf{X}$ :  $x_1, x_2, \dots, x_t, x_{t+1}, \dots, x_T$



- At a single time-step:

$$h_t = g([W_{hh}h_{t-1} + W_{hx}x_t] + c)$$

$$h_t = g(W_I[h_{t-1}; x_t] + c) \quad (\text{equivalent notation})$$

$$\hat{y}_t = \text{softmax}(W_{hy}h_t + b)$$

```
import numpy as np

def softmax(x):
    f_x = np.exp(x) / np.sum(np.exp(x))
    return f_x

class RNN:
    # ...
    def step(self, x, is_normalized=False):
        # update the hidden state
        self.h = np.tanh(np.dot(self.W_hh, self.h) + np.dot(self.W_xh, x))

        # compute the output vector
        y = np.dot(self.W_hy, self.h)

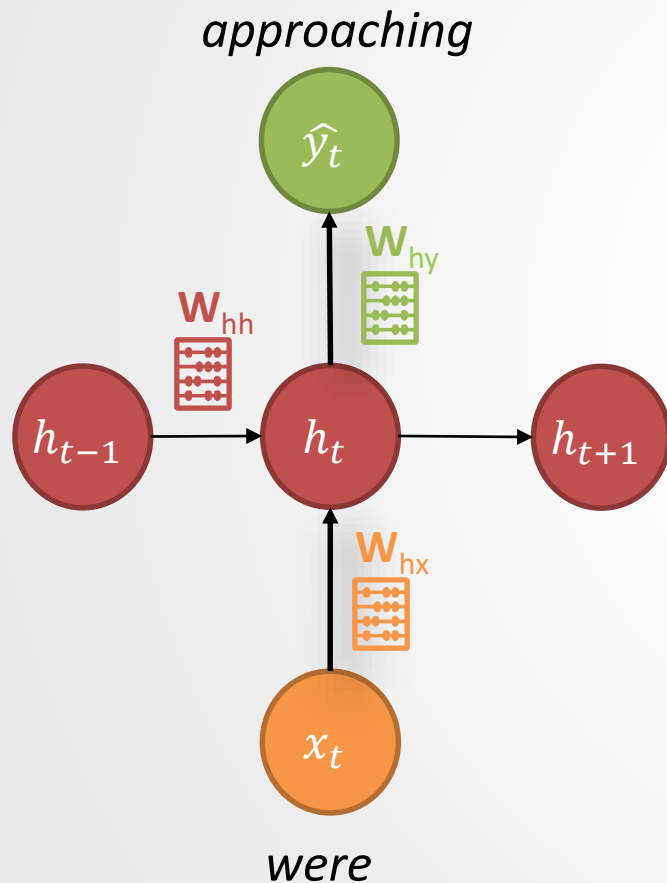
        return softmax(y) if is_normalized else y
```

$$P(x_{t+1} = v_j | x_t, \dots, x_1) = \hat{y}_{t,j}$$

# RNNs: Training

*Two riders .. approaching .. horses.*

- Given a list of word vectors  $\mathbf{X}$ :  $x_1, x_2, \dots, x_t, x_{t+1}, \dots, x_T$



$\hat{\mathbf{y}} \in \mathbb{R}^{|V|}$  is a probability distribution over the vocabulary

The output  $\hat{y}_{t,j}$  is the word (index) prediction of the next word ( $\mathbf{x}_{t+1}$ )

## Evaluation

- Same **cross-entropy** loss function

$$J^{(t)}(\theta) = - \sum_{j=1}^{|V|} \underbrace{y_{t,j}}_{\text{Ground truth}} \log \underbrace{\hat{y}_{t,j}}_{\text{prediction}}$$

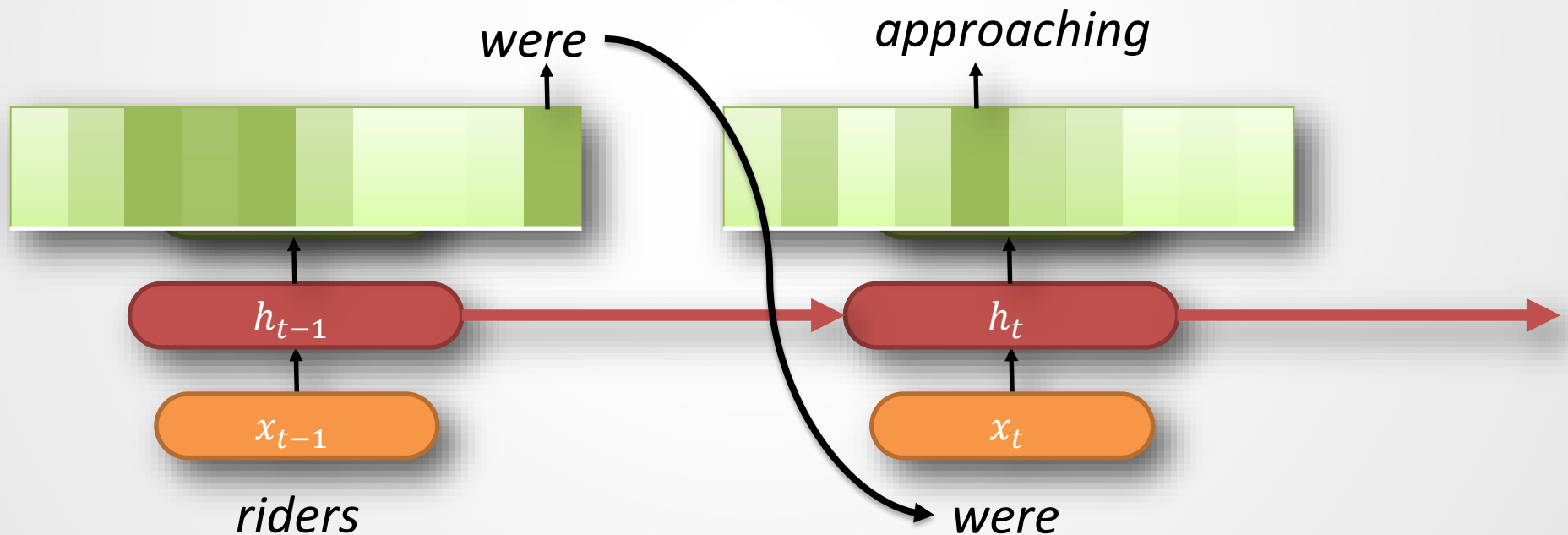
- Perplexity**:  $2^J$  (lower is better)

$$P(x_{t+1} = v_j | x_t, \dots, x_1) = \hat{y}_{t,j}$$



# Sampling from a RNN LM

- If  $|h_i| < |V|$ , we've already reduced the number of parameters relative to trigrams.
- Good news: NN encodings tend to be very compact.

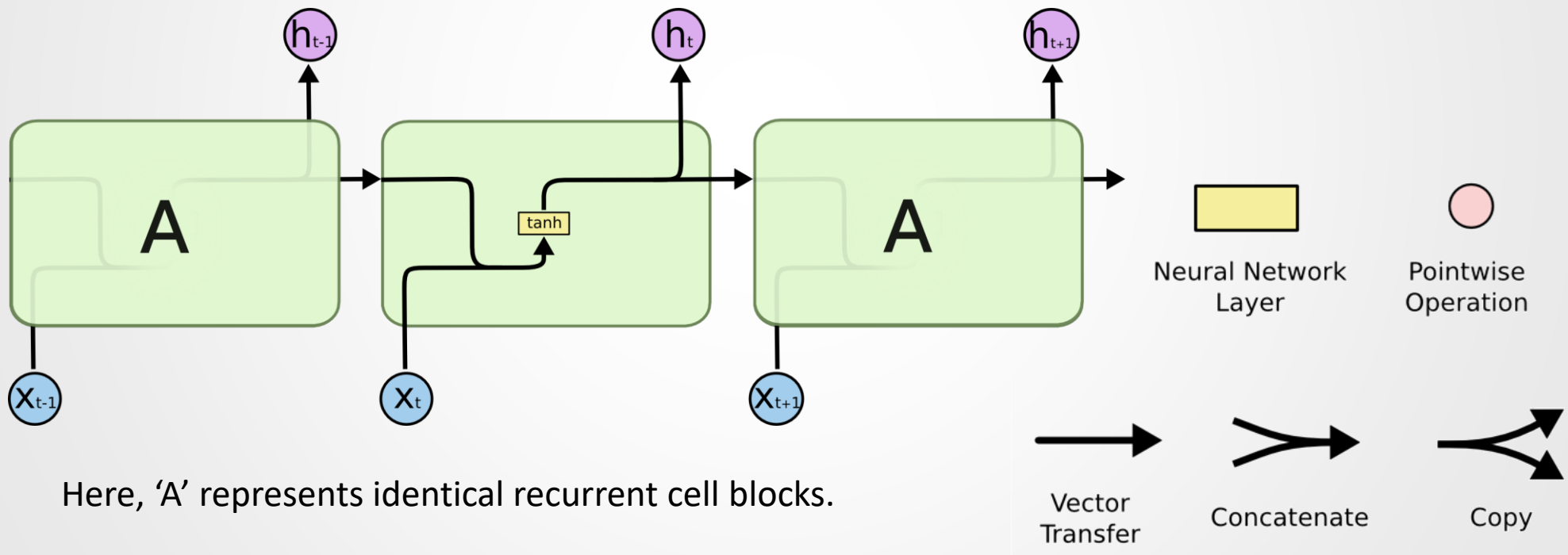


$$h_t = g([W_{hh}h_{t-1}; W_{hx}x_t] + c)$$
$$\hat{y}_t = \text{softmax}(W_{hy}h_t + b)$$

Karpathy (2015),  
[The Unreasonable Effectiveness of Recurrent Neural Networks](#)

# RNNs and retrograde amnesia

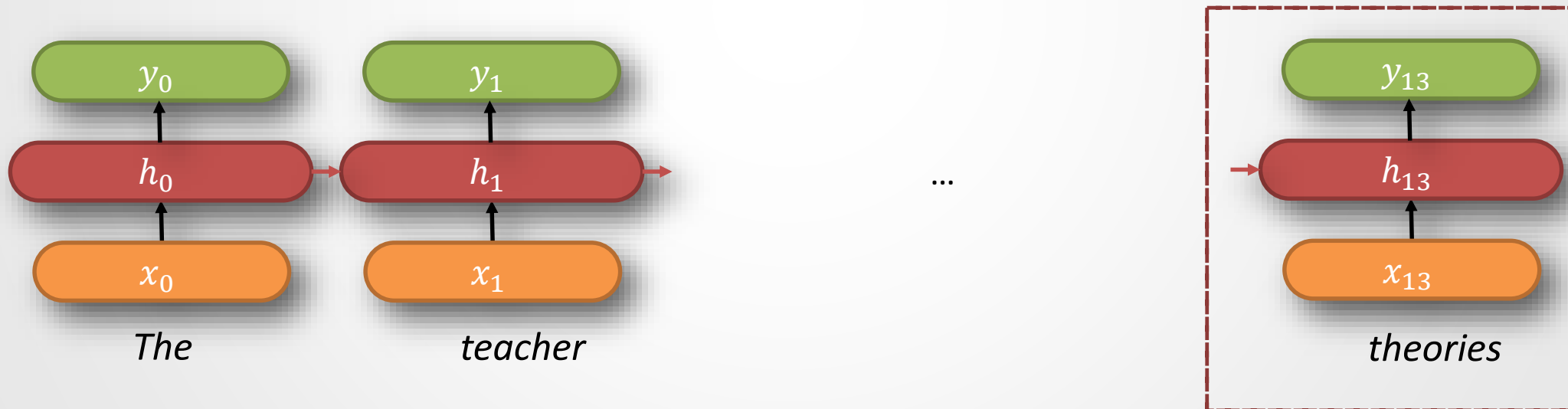
- Bad news: gradients don't multiply out well over long distances (**gradient decay**).
- Can we spend some parameters to store extra information?



Imagery and sequence from <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

# RNNs and retrograde amnesia

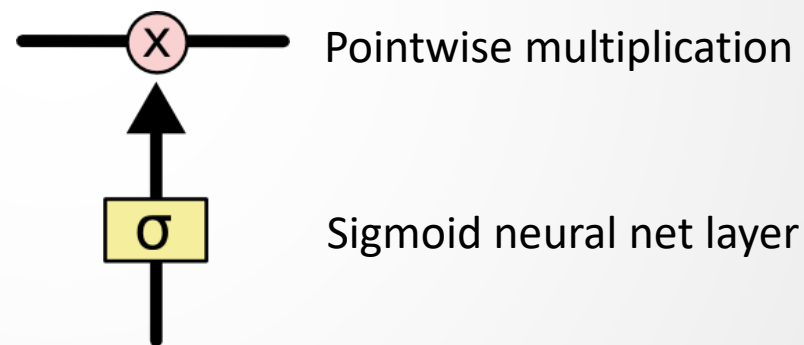
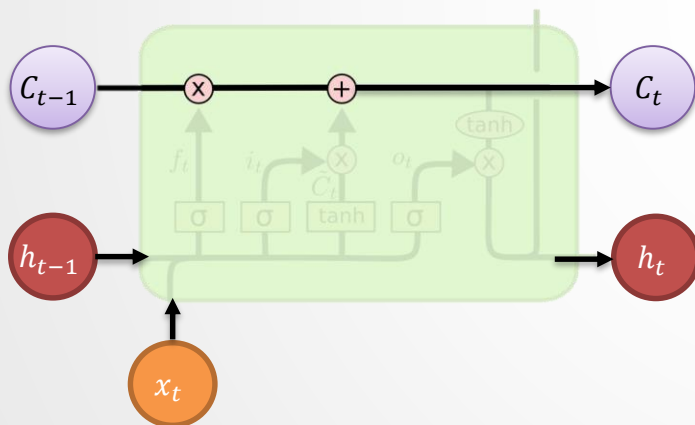
- **Catastrophic forgetting** is common.
  - E.g., the **relevant** context in “*The teacher taught transformers terribly telling tiring, tortuous theories ...*” has likely been **overwritten** by the time  $h_{13}$  is produced.



Bengio Y, Simard P, Frasconi P. (1994) Learning Long-Term Dependencies with Gradient Descent is Difficult. IEEE Trans. Neural Networks.;5:157–66. doi:10.1109/72.279181

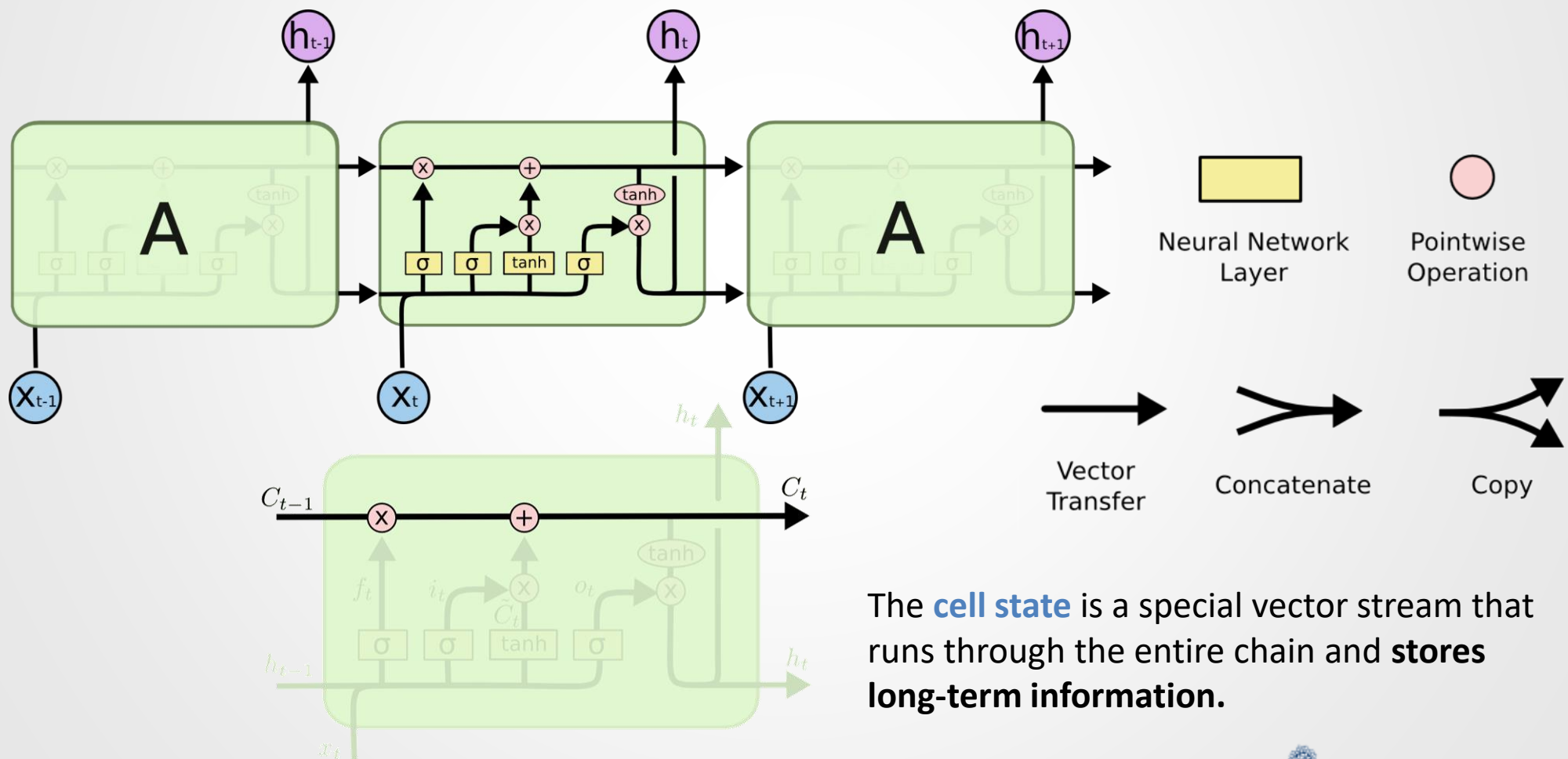
# Long short-term memory (LSTM)

- Within each *recurrent unit or cell*:
  - Self-looping recurrence for **cell state** using vector  $C$
  - Information flow regulating structures called **gates**



# LSTM – core ideas

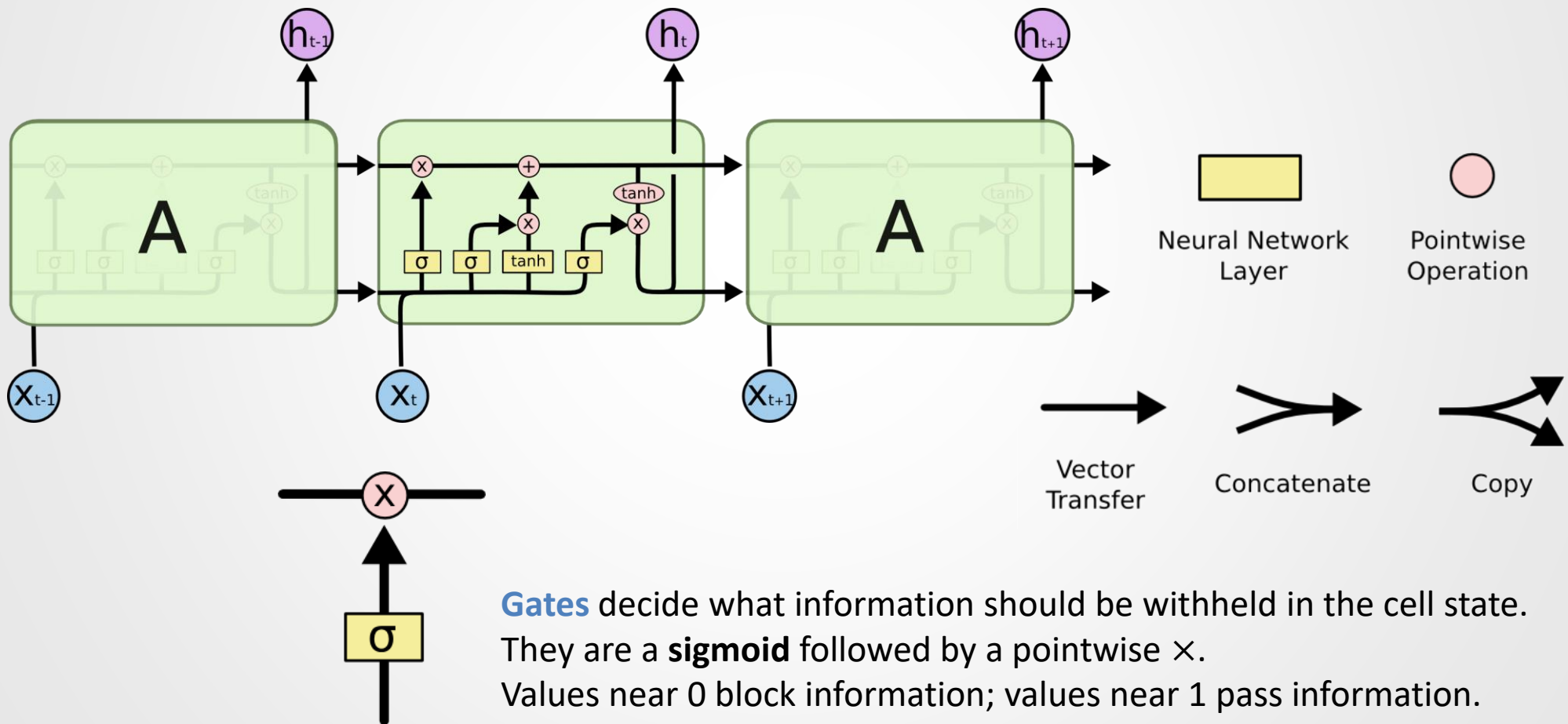
- In each **cell** (i.e. recurrent unit) in an LSTM, there are four interacting neural network layers.



The **cell state** is a special vector stream that runs through the entire chain and **stores long-term information**.

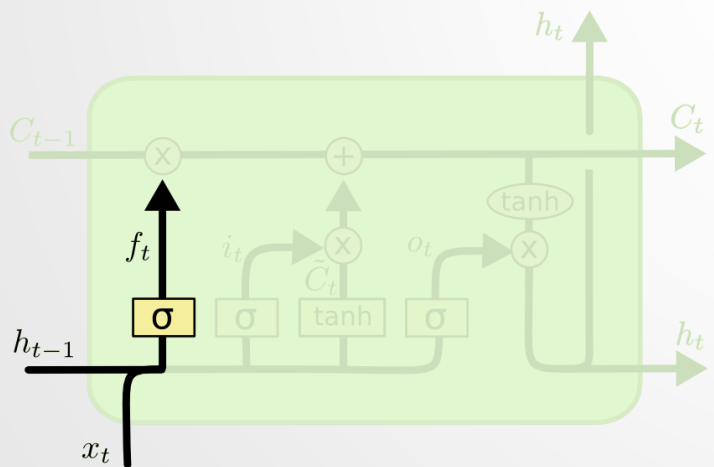
# LSTM – core ideas

- In each **cell** (i.e. recurrent unit) in an LSTM, there are four interacting neural network layers.

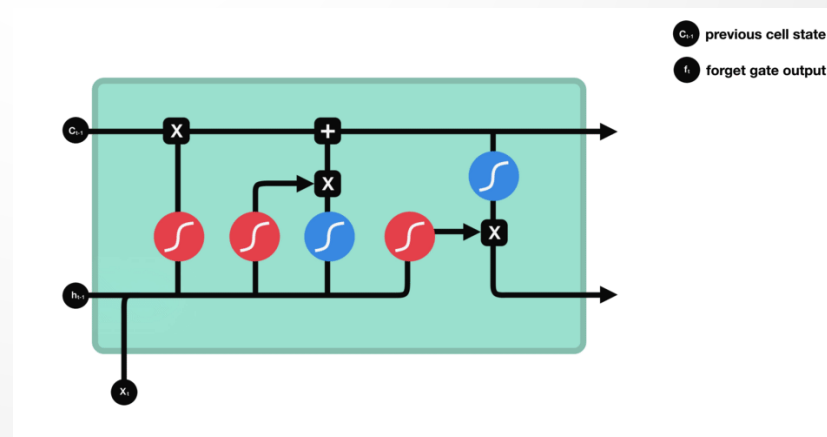


# LSTM step 1: decide what to forget

- The **forget gate layer** compares  $h_{t-1}$  and the current input  $x_t$  to decide which elements in cell state  $C_{t-1}$  to keep and which to turn off.
  - E.g., the cell state might ‘remember’ the number (sing./plural) of the current subject, in order to predict appropriately conjugated verbs, but decide to forget it when a new subject is mentioned at  $x_t$ .
    - (There’s scant evidence that such information is so readily interpretable.)



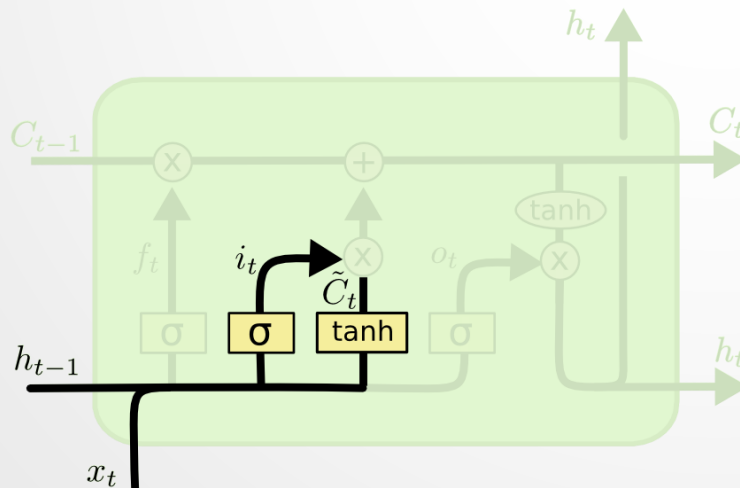
$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$





# LSTM step 2: decide what to store

- The **input gate layer** has two steps.
  - First, a sigmoid layer  $\sigma$  decides which cell units to update.
  - Next, a **tanh** layer creates new candidate values  $\tilde{C}_t$ .
  - E.g., the  $\sigma$  can turn on the 'number' units, and the tanh can push information on the current subject.
  - The  $\sigma$  layer is important – we don't want to push information on units (i.e., latent dimensions) for which we have no information.

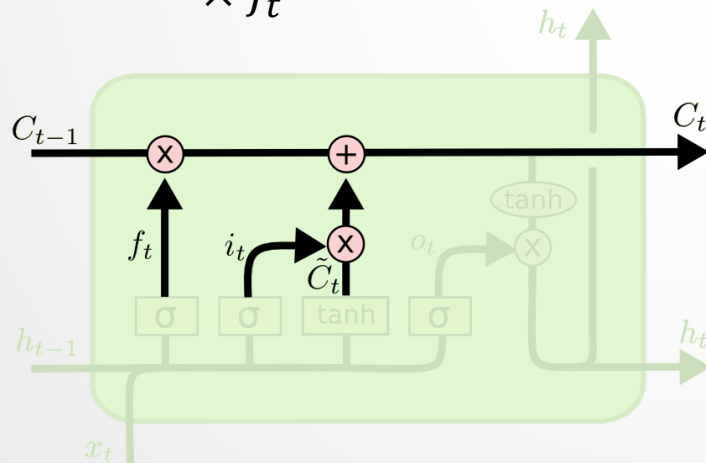
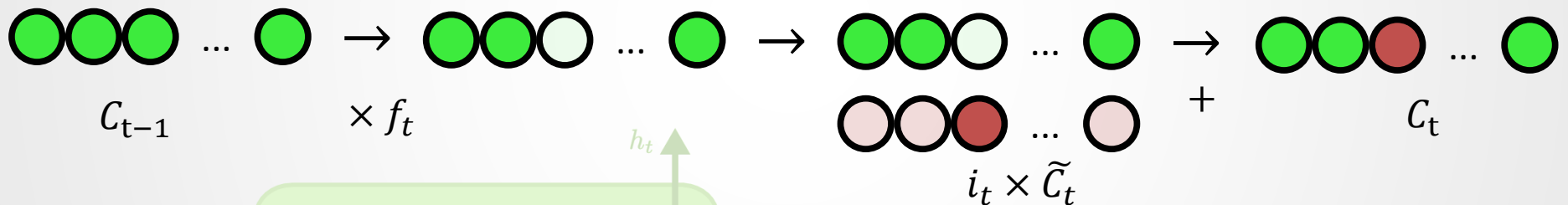


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

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

# LSTM step 3: update the cell state

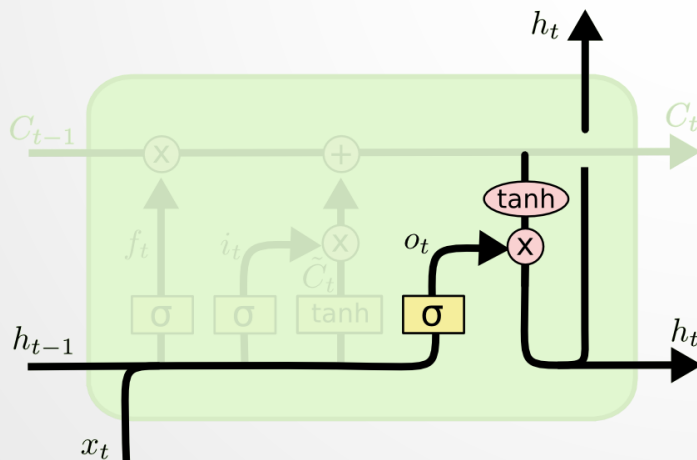
- Update  $C_{t-1}$  to  $C_t$ .
  - First, forget what we want to forget: multiply  $C_{t-1}$  by  $f_t$ .
  - Then, create a 'mask vector' of information we want to store,  $i_t \times \tilde{C}_t$ .
  - Finally, write this information to the new cell state  $C_t$ .



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

# LSTM step 4: output and feedback

- Output something,  $o_t$ , based on the current  $x_t$  and  $h_{t-1}$ .
- Combine the output with the cell to give your  $h_t$ .
  - Normalize cell  $C_t$  on  $[-1,1]$  using tanh and combine with  $o_t$
- In some sense,  $C_t$  is **long-term** memory and  $h_t$  is the **short-term memory** (hence the name).

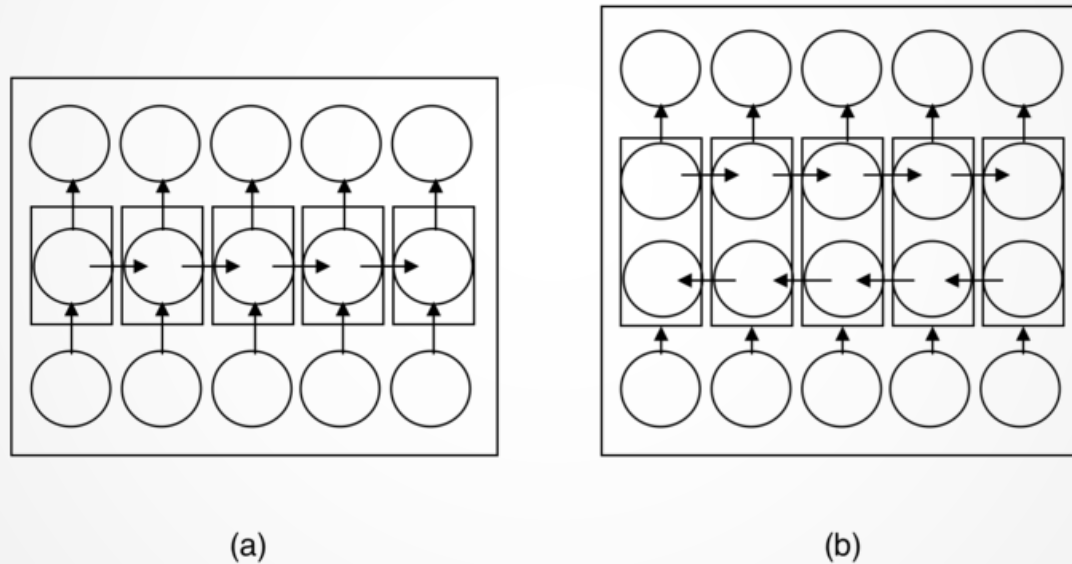


$$o_t = \sigma(W_o[h_{t-1}, x_t] + b_o)$$

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

# Variants of LSTMs

- There are many variations on LSTMs.
  - ‘*Bidirectional LSTMs*’ (and bidirectional RNNs generally), learn. (Similar: *Multi-stack RNNs*)



Structure overview

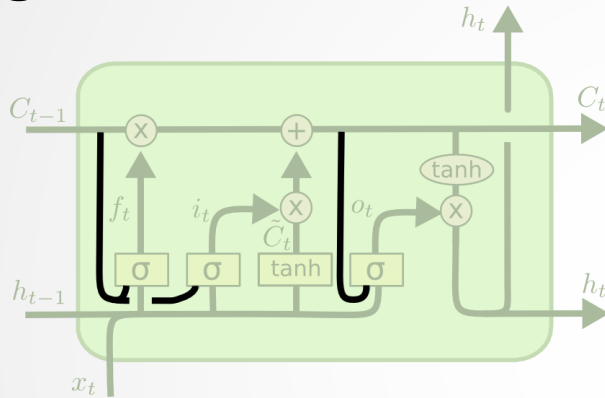
(a) unidirectional RNN

(b) bidirectional RNN

Schuster, Mike, and Kuldip K. Paliwal. (1997) Bidirectional recurrent neural networks. *Signal Processing, IEEE Transactions on* **45**(11) (1997): 2673-2681.2.

# Variants of LSTMs

- Gers & Schmidhuber (2000) add '**peepholes**' that allow all sigmoids to read the cell state.

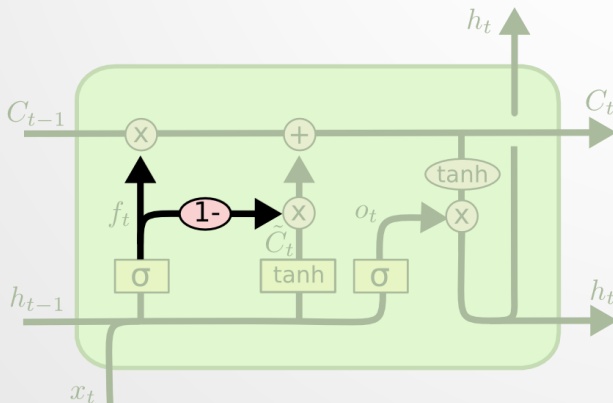


$$f_t = \sigma(W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f)$$

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

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

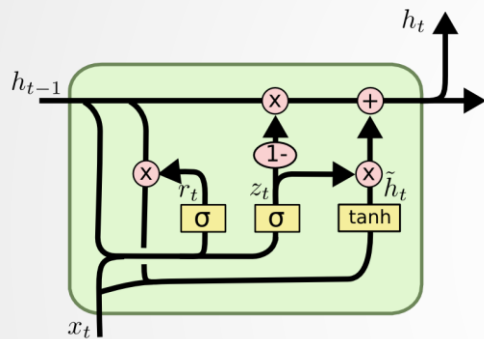
- We can **couple** the '*forget*' and '*input*' gates.
  - Joint decisioning is more efficient.



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

# Aside - Variants of LSTMs

- **Gated Recurrent units** (GRUs; [Cho et al \(2014\)](#)) go a step further and also merge the cell and hidden states.



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t]) \text{ Update gate}$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t]) \text{ Reset gate (0: replace units in } h_{t-1} \text{ with those in } x_t)$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

- Which of these variants is best? Do the differences matter?
  - [Greff, et al. \(2015\)](#) do a nice comparison of popular variants, finding that **they're all about the same**
  - [Jozefowicz, et al. \(2015\)](#) tested more than ten thousand RNN architectures, finding some that worked better than LSTMs on certain tasks.

# CONTEXTUAL WORD EMBEDDINGS



# Deep contextualized representations

- What does the word *play* mean?

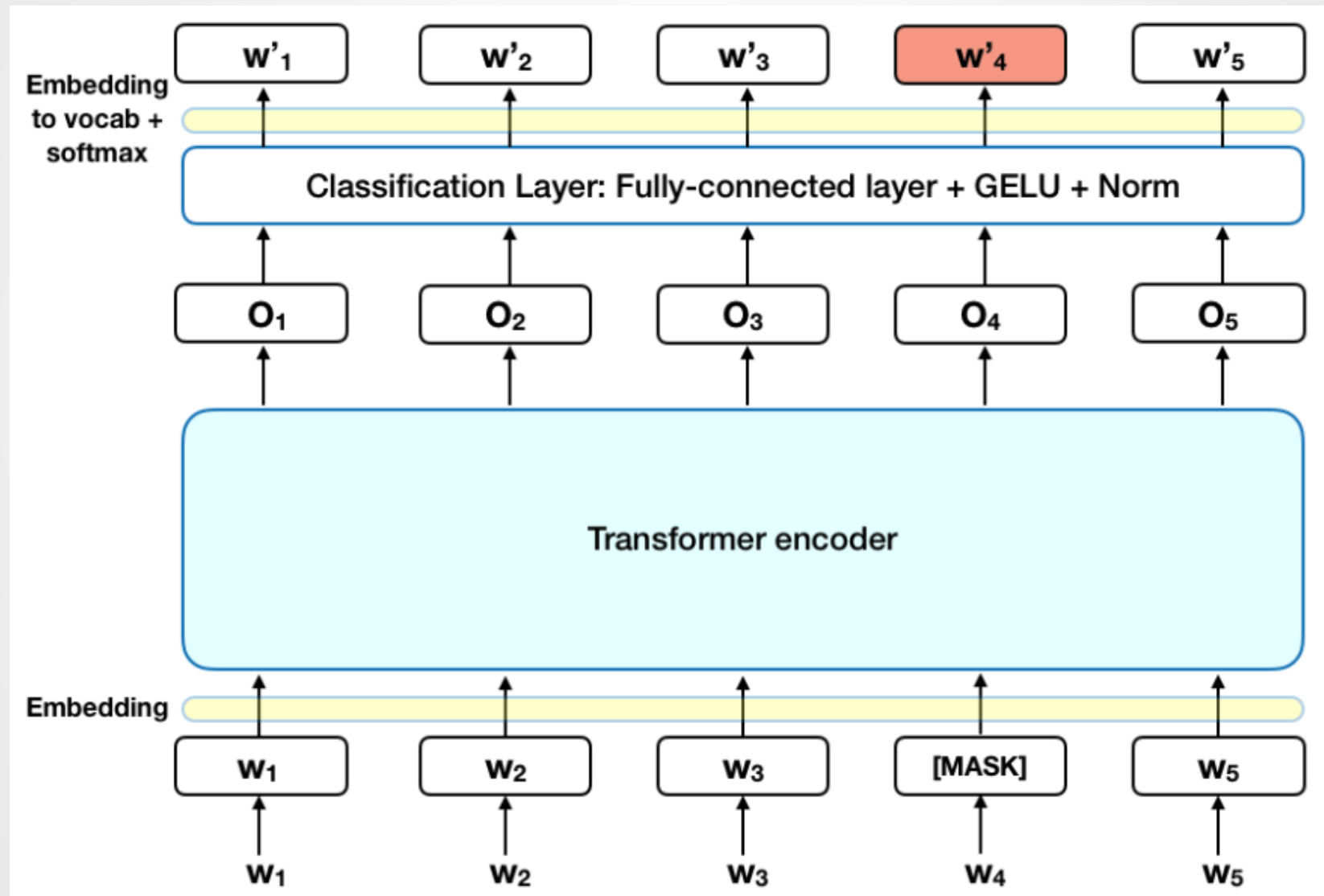


AllenNLP

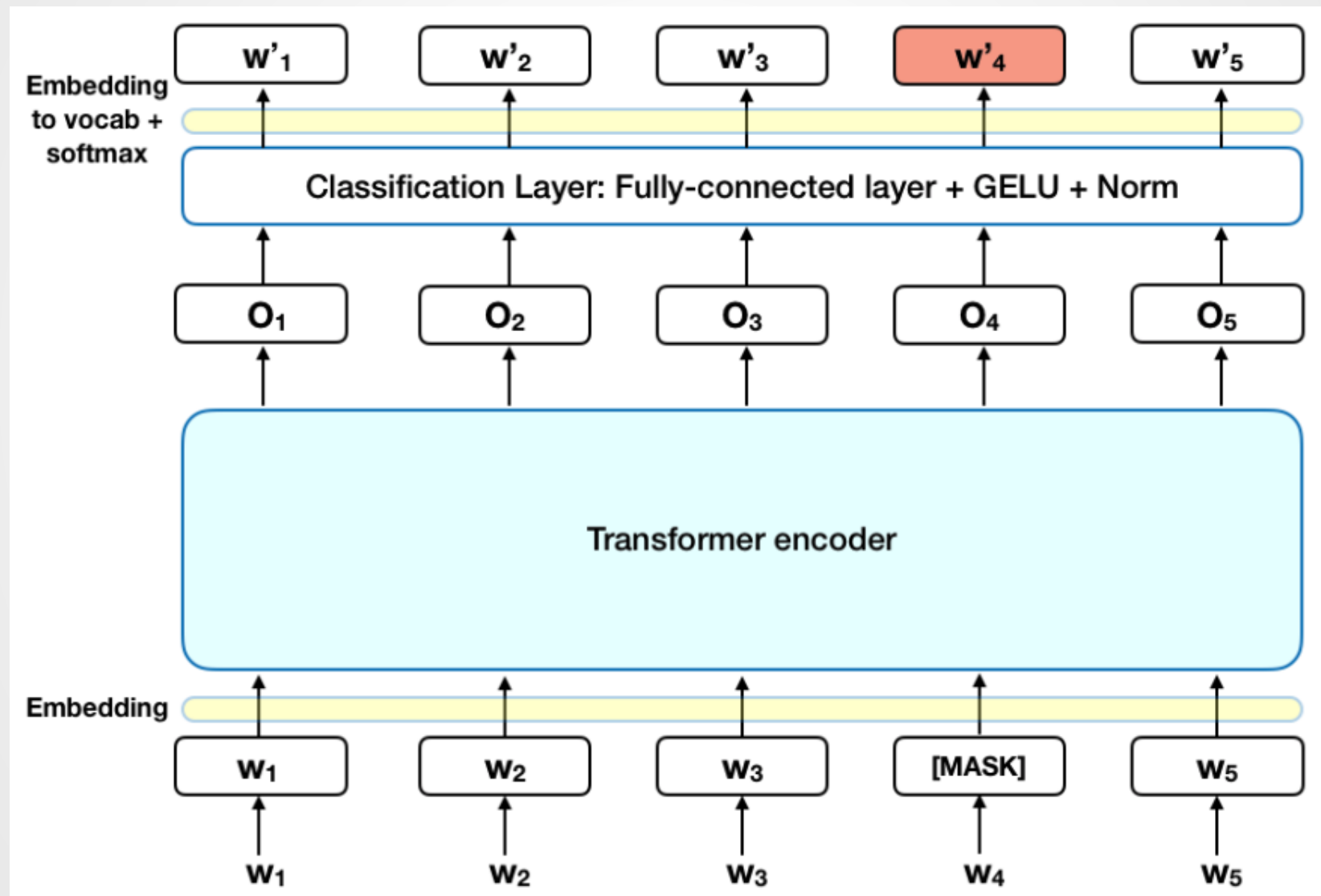
Peters ME, Neumann M, Iyyer M, *et al.* (2018) Deep contextualized word representations.

Published Online First: 2018. doi:10.18653/v1/N18-1202; <http://arxiv.org/abs/1802.05365>

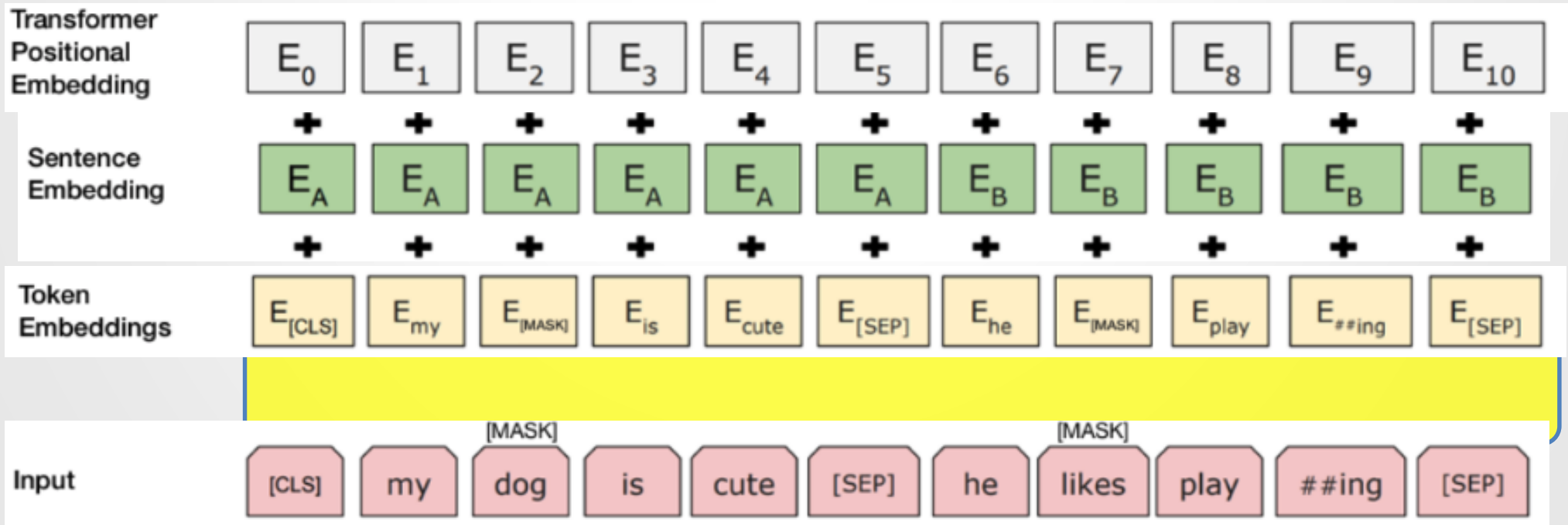
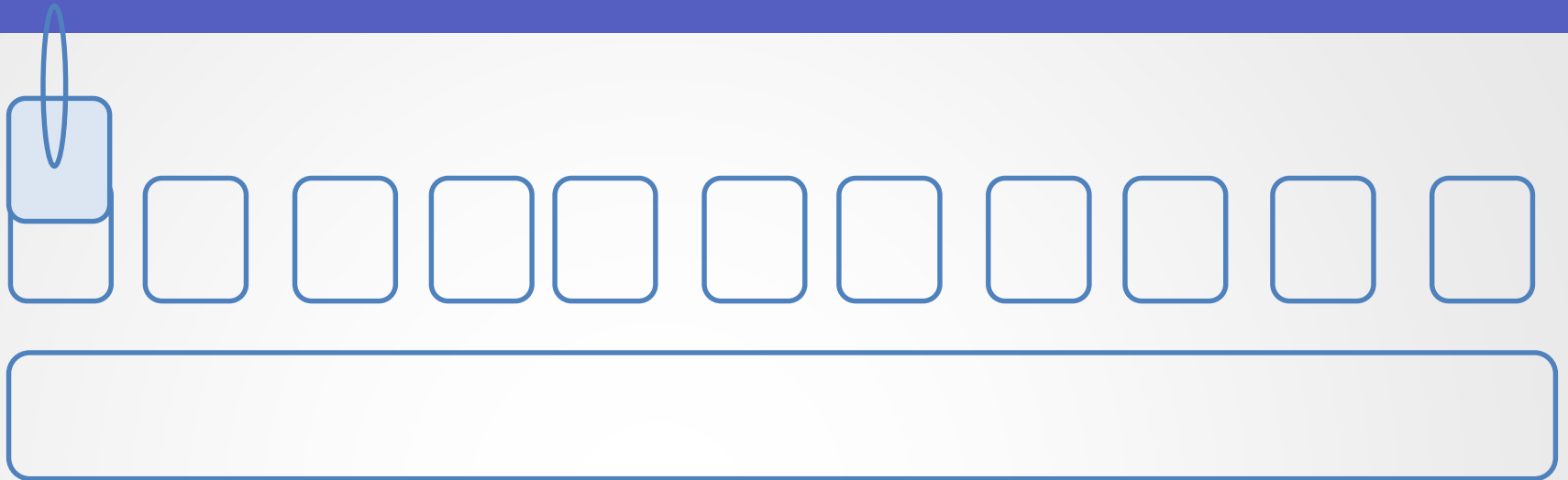
# • BERT



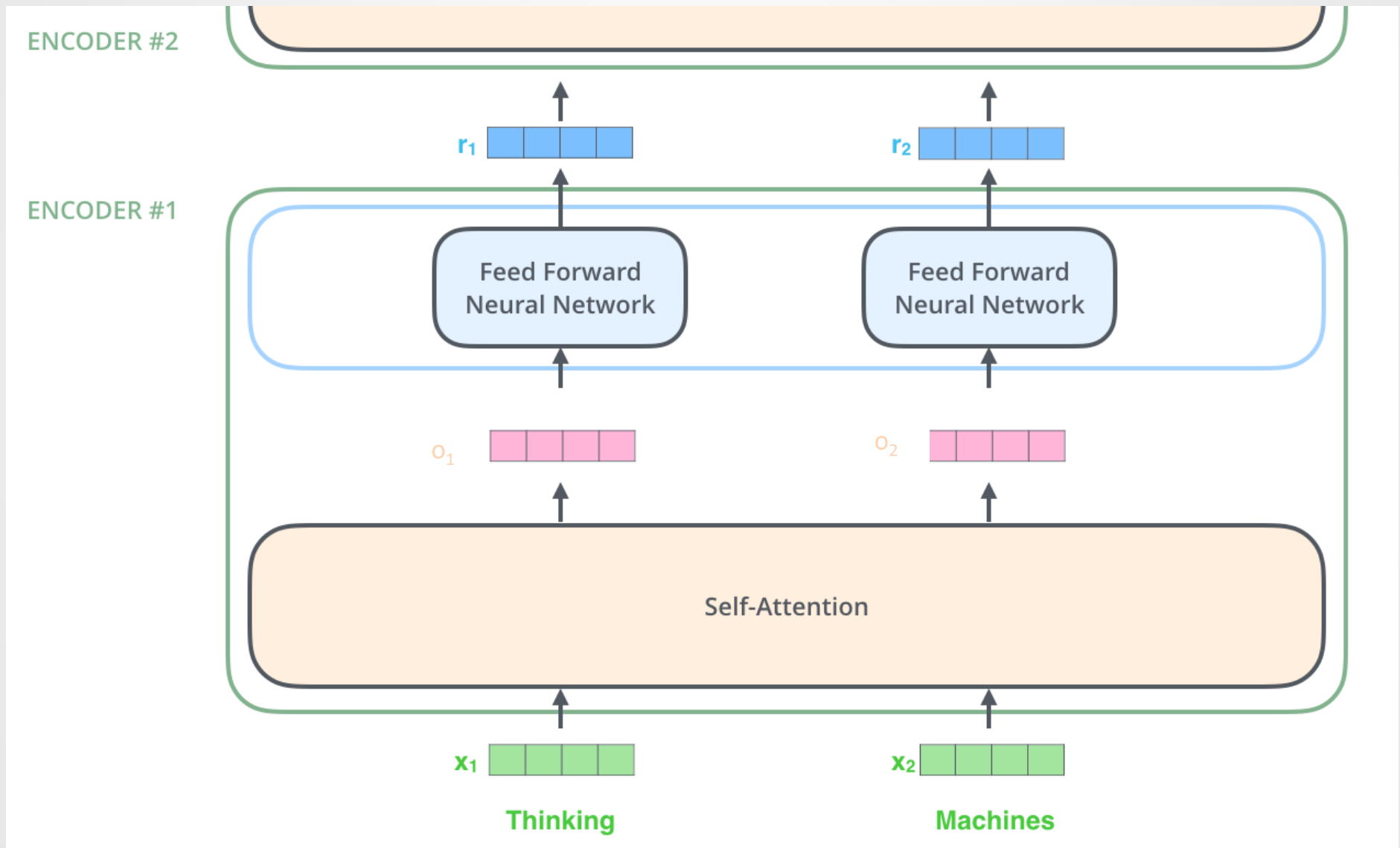
# Training task 1: Masking



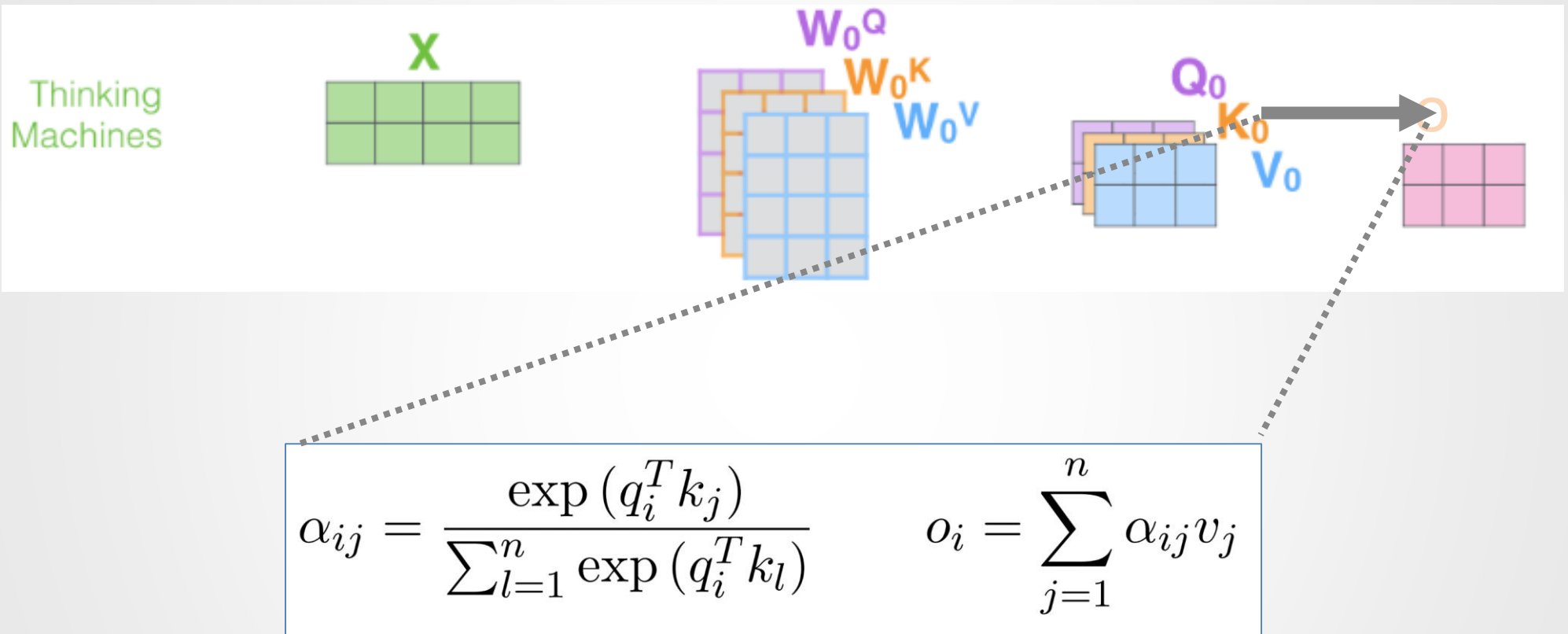
# Training task 2: Next Sent.



# .Transformers



# .Self-attention



# .Multiheaded Self attention

1) This is our input sentence\*

Thinking  
Machines

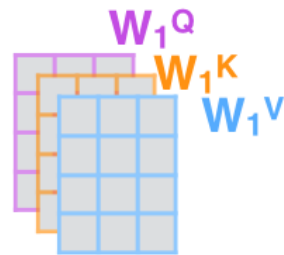
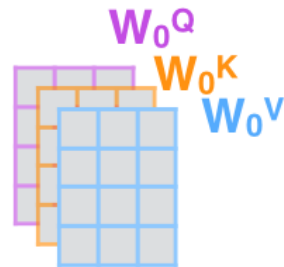
2) We embed each word\*



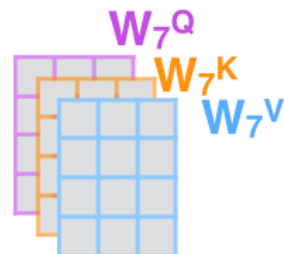
\* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



3) Split into 8 heads. We multiply  $X$  or  $R$  with weight matrices



...



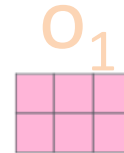
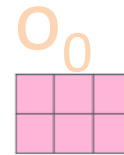
4) Calculate attention using the resulting  $Q/K/V$  matrices



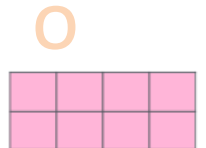
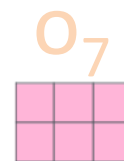
...



5) Concatenate the resulting  $Z$  matrices, then multiply with weight matrix  $W^O$  to produce the output of the layer



...



# .Positional encodings

$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \sin(\omega_2 \cdot t) \\ \cos(\omega_2 \cdot t) \\ \vdots \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_{d \times 1}$$

where

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases}$$

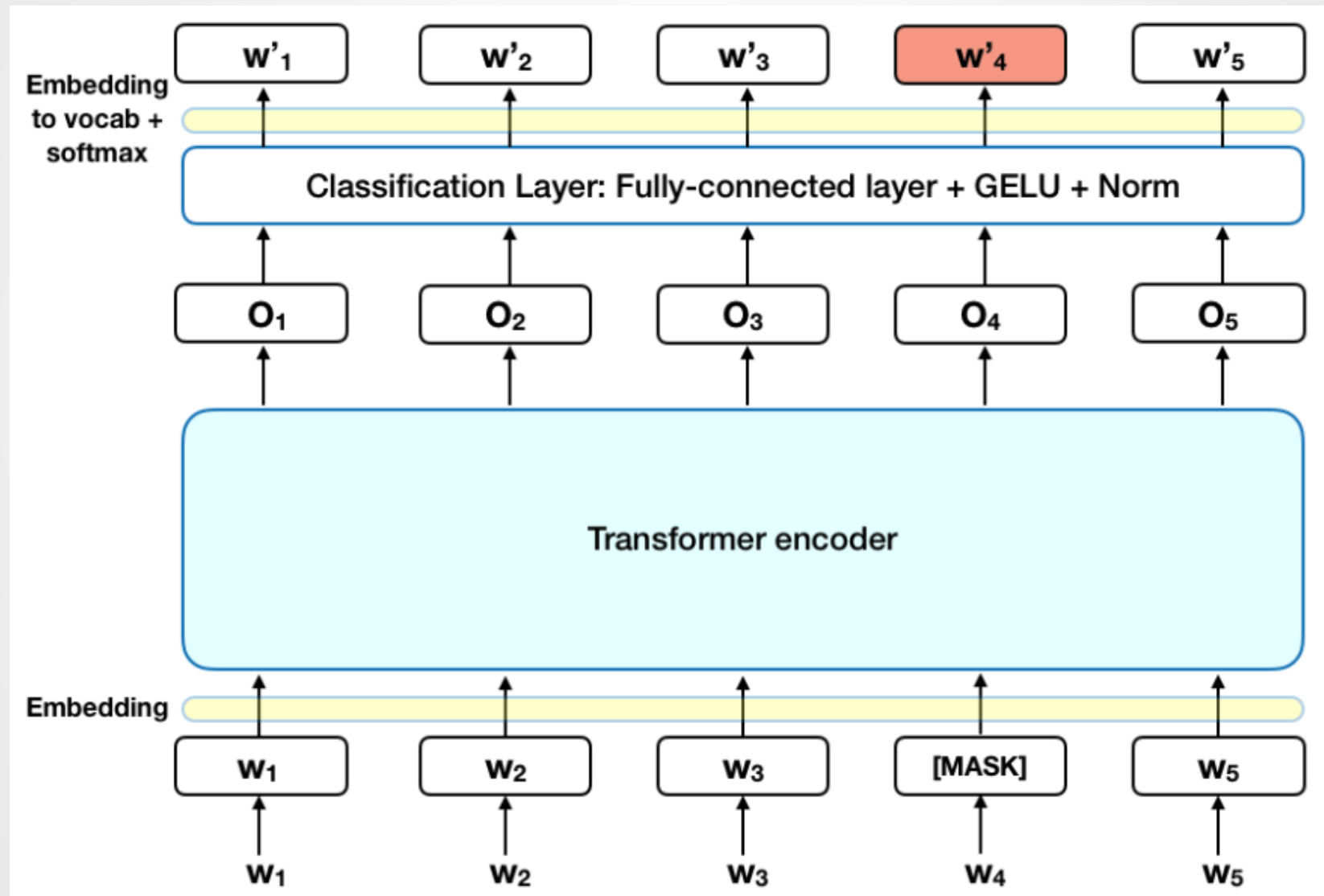
$$\omega_k = \frac{1}{10000^{2k/d}}$$



# .Huh?

- Encodings of any two distinct positions are distinct
- Each position maps to only one encoding
- Test sentences may be longer than training
- Distance between two positions should be constant across sentences (of varying lengths).

# Training task 1: Masking



# .The truth about masking

- Real easy to do well on MASKed position and nothing else
- Real easy to learn to copy the context-independent embedding
- So...
  - 80% of the time: MASK
  - 10% of the time: correct word
  - 10% of the time: another random word