

Social Media & Text Analysis

lecture 10 - Deep Learning for NLP



CSE 5539-0010 Ohio State University
Instructor: Wei Xu
Website: socialmedia-class.org

Last Class

- December 7 (Review and Q&A)
- No final exam
- Also one-on-one office hour by appointment on December 7 for homework #3 or research projects or Q&A

Word Embeddings

- Skip-gram — predicts surrounding “outside” words given the “center” word

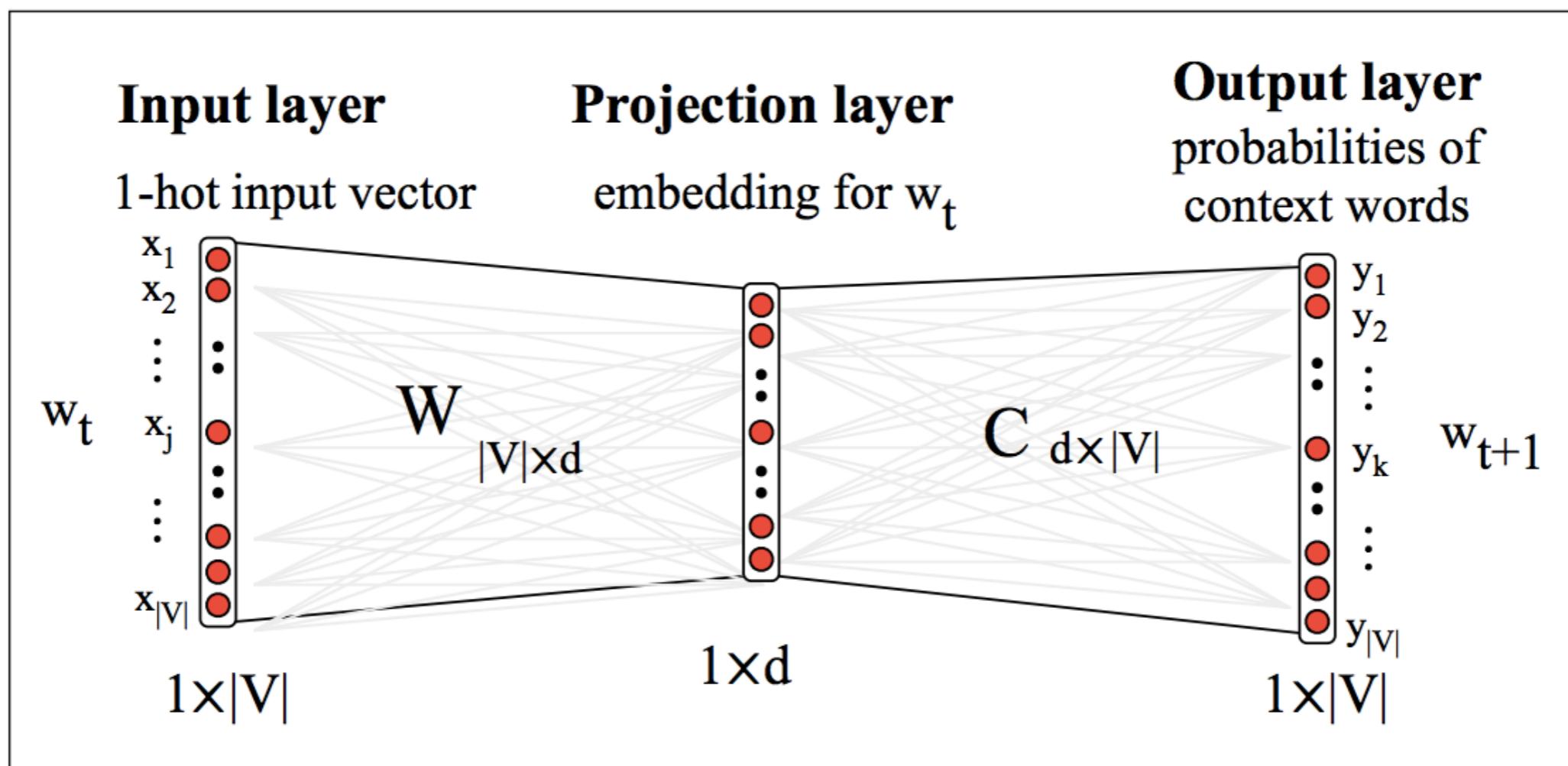
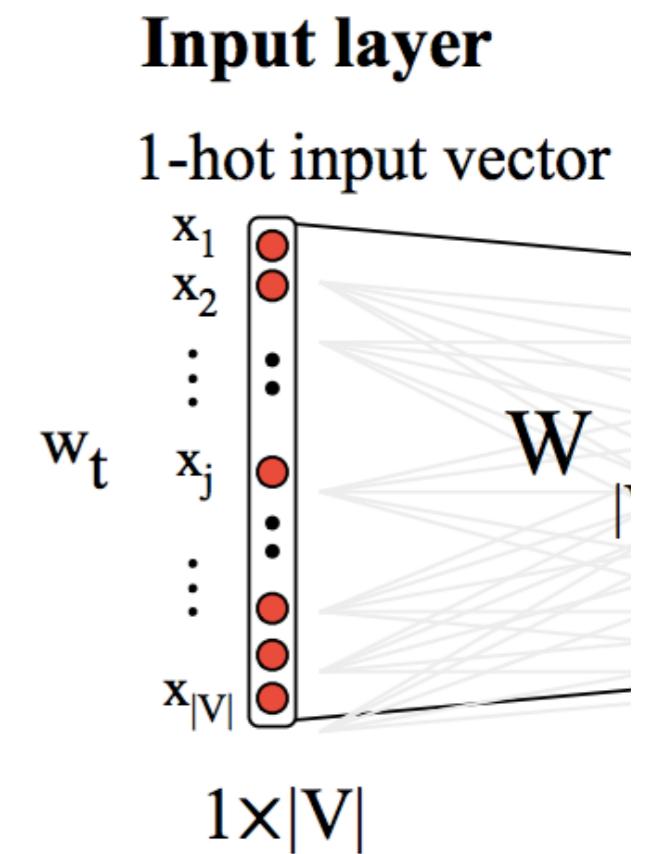


Figure 16.5 The skip-gram model viewed as a network (Mikolov et al. 2013, Mikolov et al. 2013a).

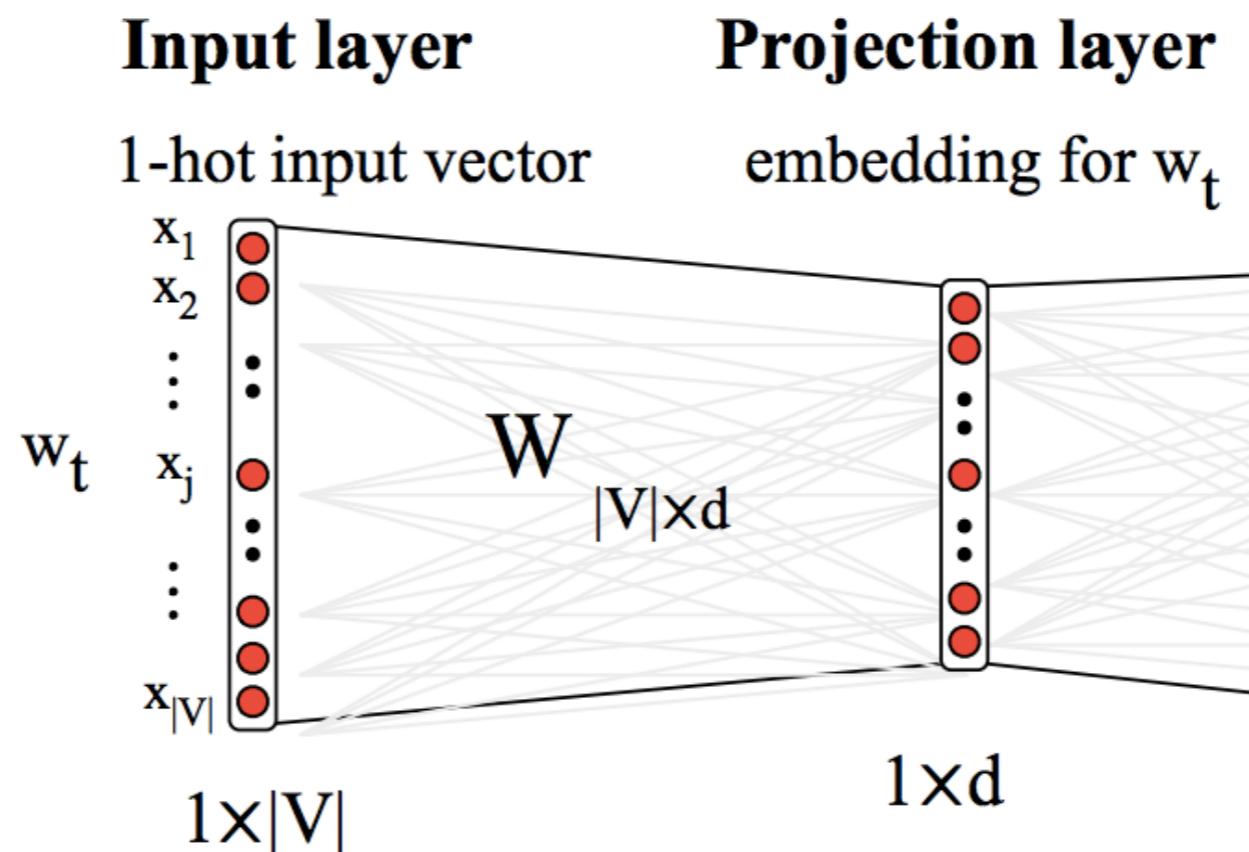
Input Layer

- “one-hot” word vectors
 - a vector of dimension $|V|$ (size of vocabulary)
 - all “0”s except a single “1” in the vector
 - different positions of that “1” represent different words



Hidden (Projection) Layer

- A simple look up — the rows of this weight matrix are actually “input” word vectors



Hidden (Projection) Layer

- A simple look up — the rows of this weight matrix are actually “input” word vectors

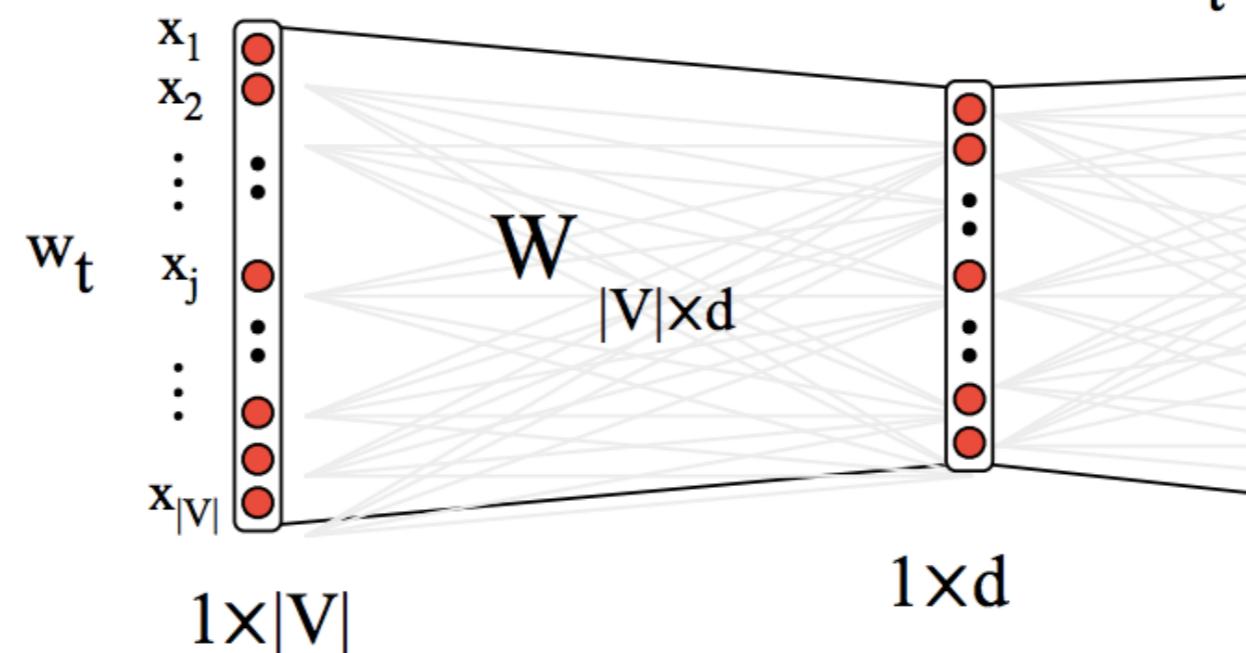
$$\begin{bmatrix} 0 & 0 & 0 & \boxed{1} & 0 \end{bmatrix} \times \begin{bmatrix} 17 & 24 & 1 \\ 23 & 5 & 7 \\ 4 & 6 & 13 \\ \boxed{10} & \boxed{12} & \boxed{19} \\ 11 & 18 & 25 \end{bmatrix} = [10 \quad 12 \quad 19]$$

Input layer

1-hot input vector

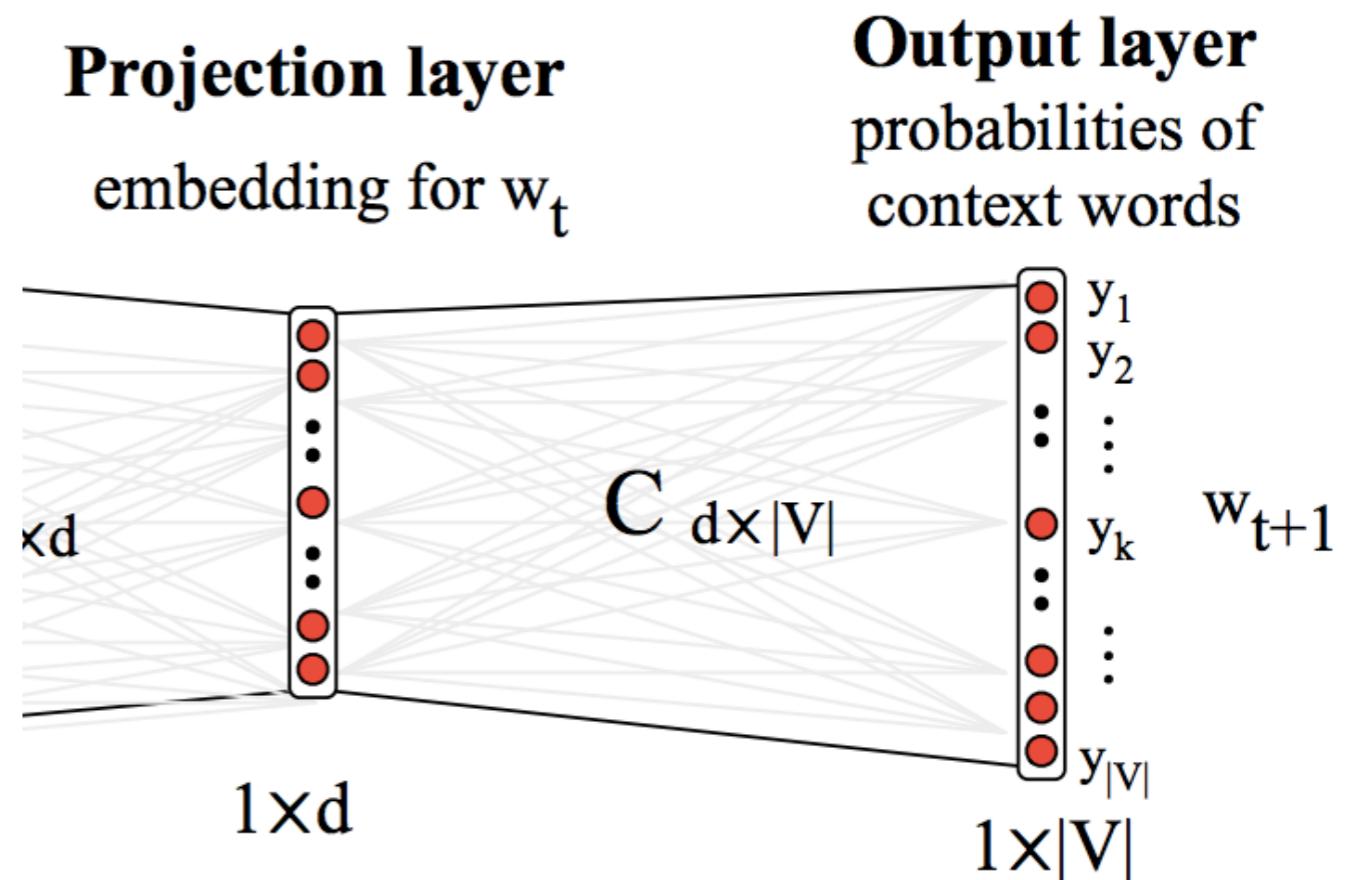
Projection layer

embedding for w_t



Output Layer

- predicts surrounding “outside” (context) words given the “center” word → A classification problem!
- Softmax Regression = Multi-class Logistic Regression



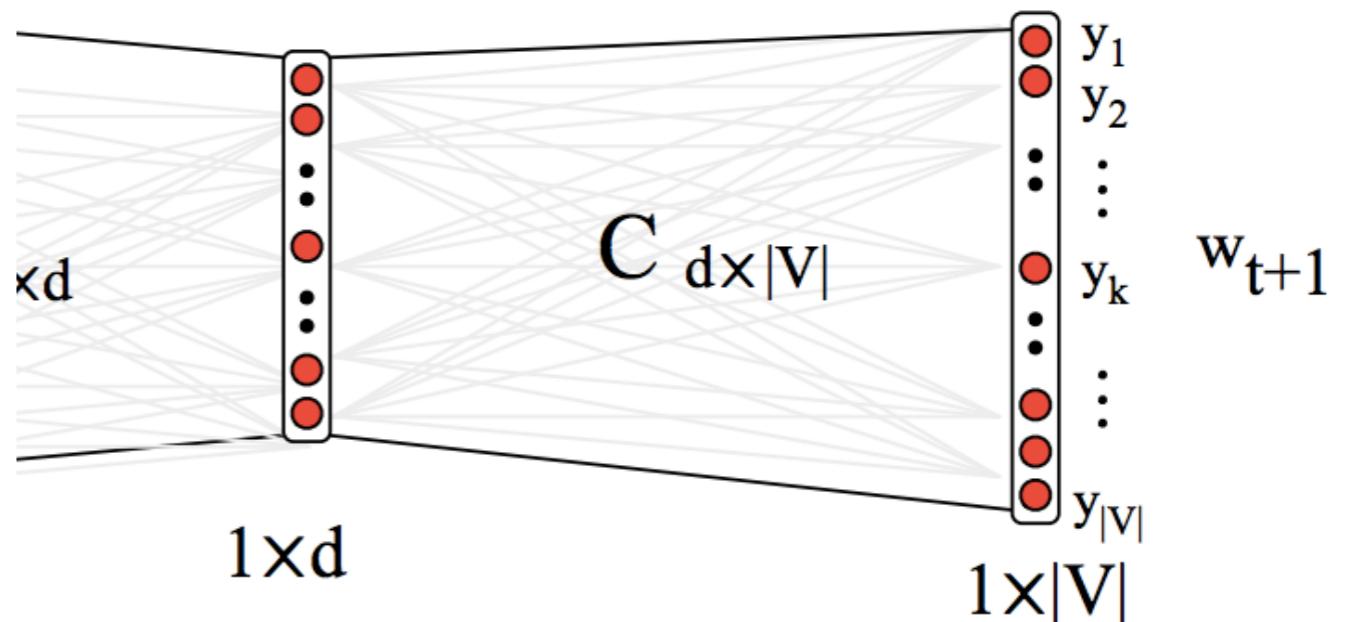
Output Layer

- Objective function: maximize the log probability of any “outside” (context) word given the “center” word

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

Projection layer
embedding for w_t

Output layer
probabilities of
context words

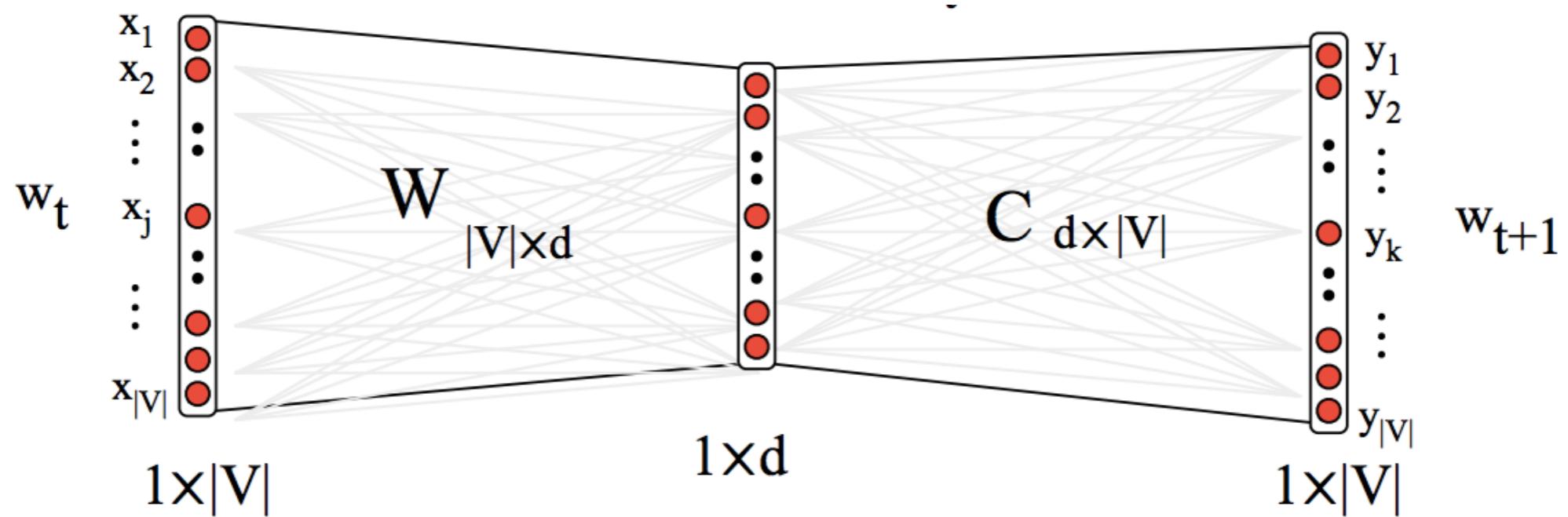


Output Layer

- predicts surrounding “outside” (context) words given the “center” word

$$p(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w=1}^W \exp(u_w^T v_c)}$$

- so, every word has two vectors!



Gradient Descent

- Cost/Objective function:

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

- For a “center” word and an “outside” word:

$$\log p(o|c) = \log \frac{\exp(u_o^T v_c)}{\sum_{w=1}^W \exp(u_w^T v_c)}$$

Gradient Descent

- Basics:

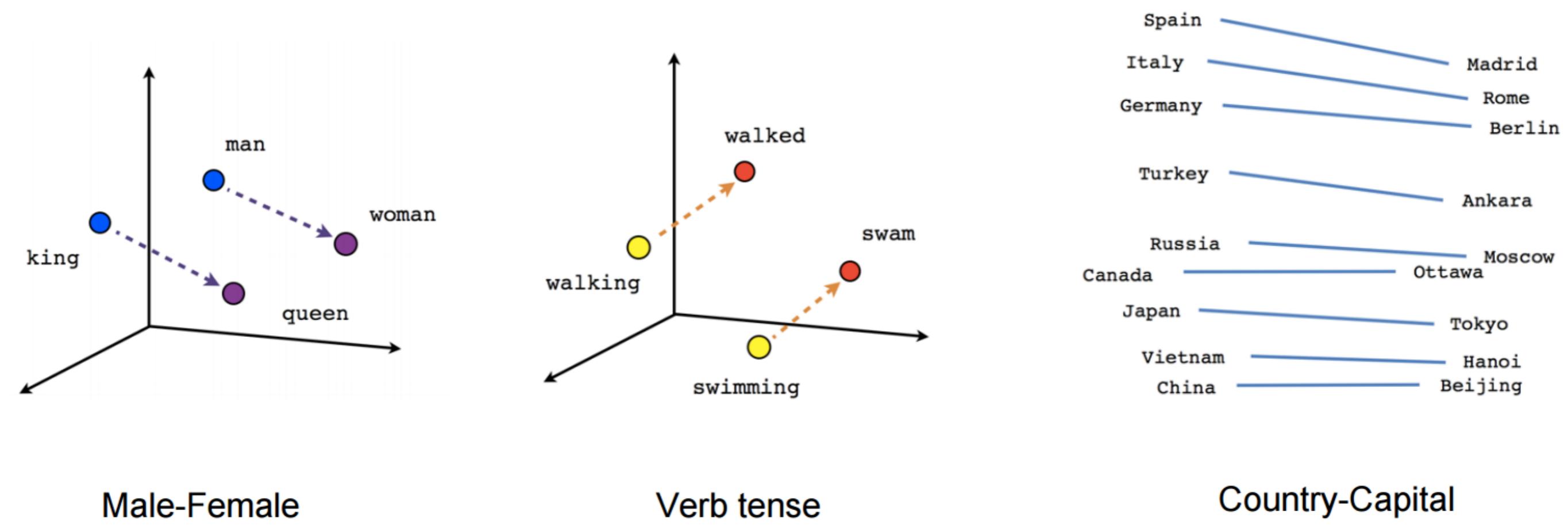
$$\frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a}$$

$$\frac{\partial e^x}{\partial x} = e^x \qquad \qquad \frac{\partial \log x}{\partial x} = \frac{1}{x}$$

- Chain Rule:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial x} = \frac{\partial f(g)}{\partial g} \frac{\partial g(x)}{\partial x}$$

Visualization



Male-Female

Verb tense

Country-Capital

Visualization

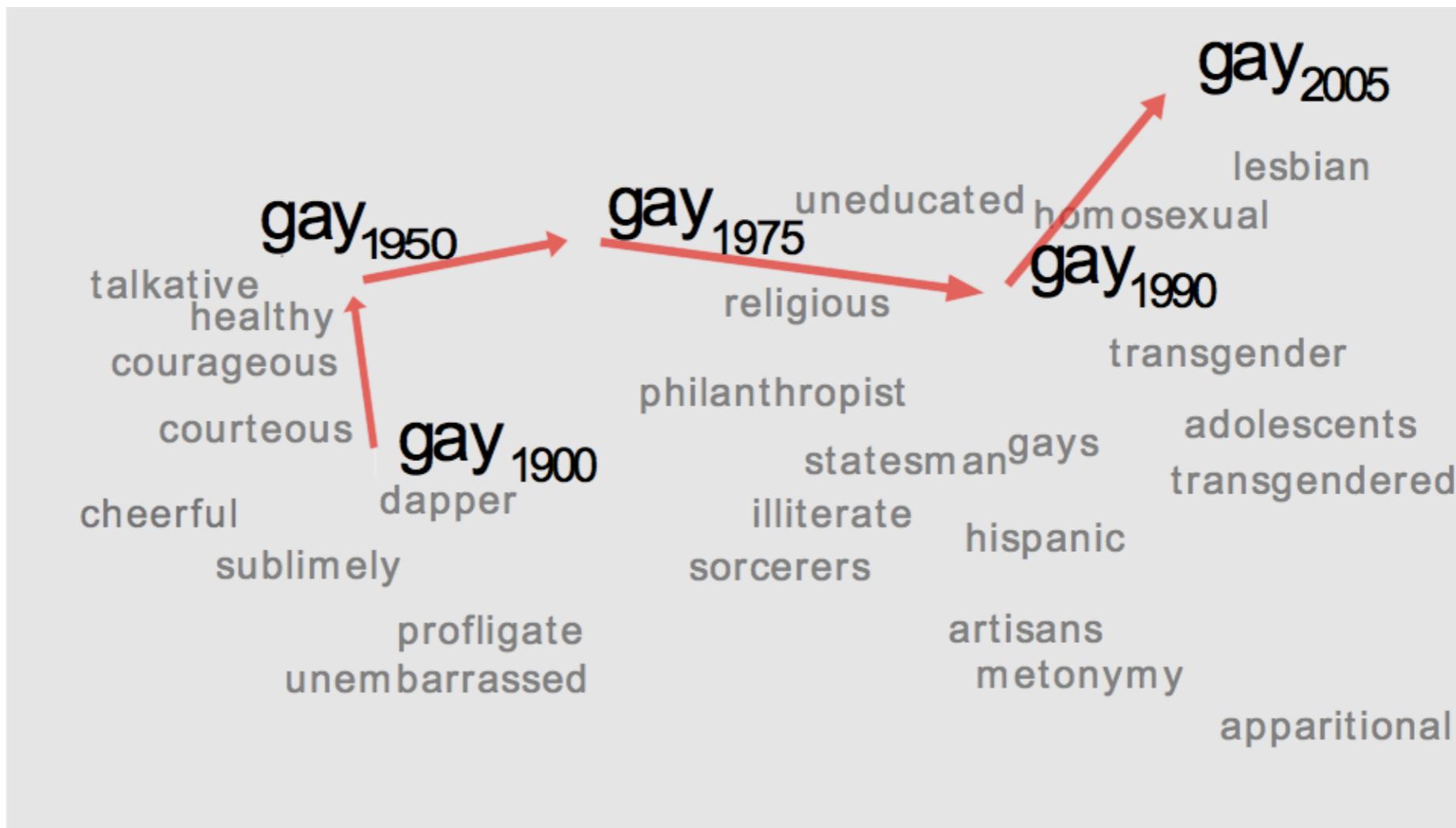


Figure 1: A 2-dimensional projection of the latent semantic space captured by our algorithm. Notice the semantic trajectory of the word **gay** transitioning meaning in the space.

Source: Kulkarni et al. (WWW 2015)
Statistically Significant Detection of Linguistic Change

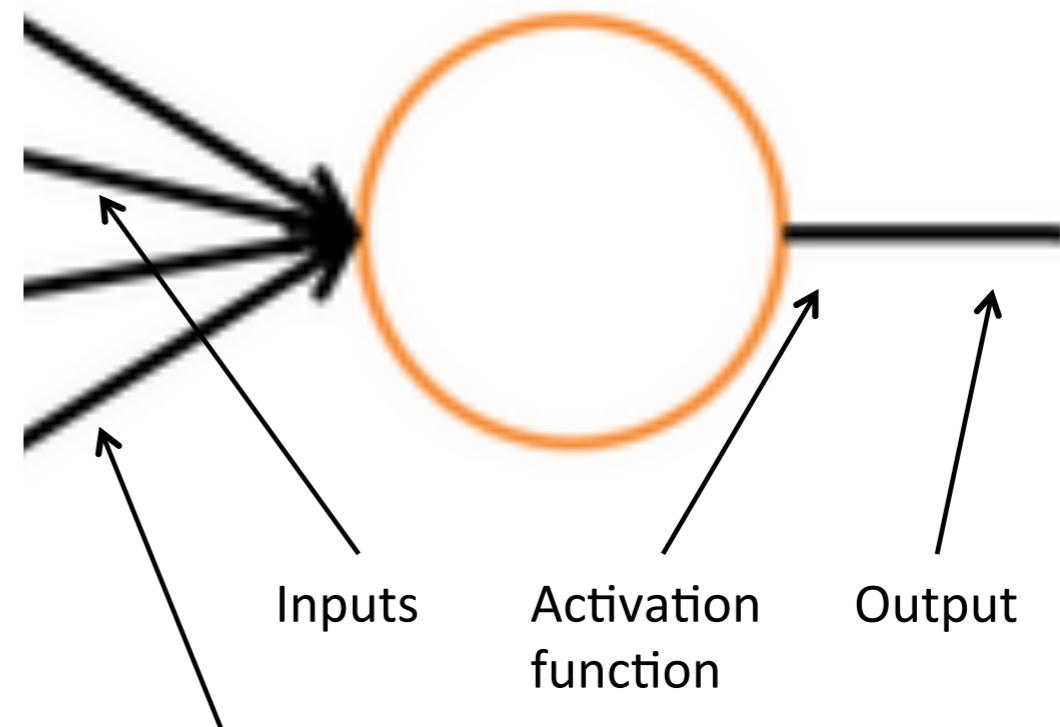
DeepNLP

- Word vector and Recurrent Neural Networks (RNN) are the two most important concepts for deepNLP
- also, Convolutional Neural Networks (CNN), Long Short-Term Memory Networks (LSTM) ...

A Neuron

- If you know Logistic Regression, then you already understand a basic neural network neuron!

A single neuron
A computational unit with n (3) inputs
and 1 output
and parameters W, b



Bias unit corresponds to intercept term

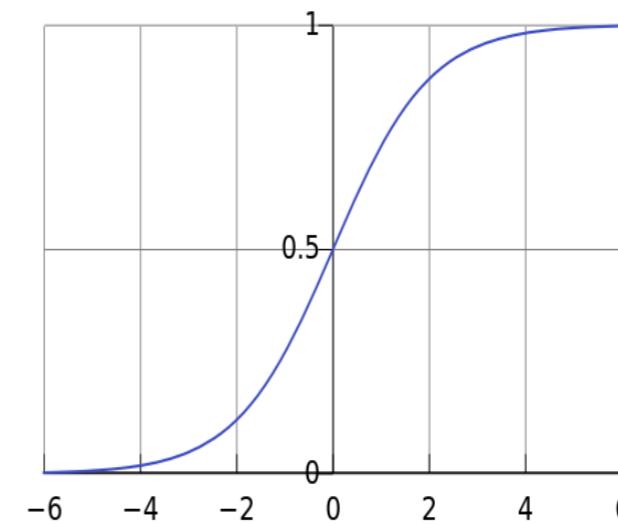
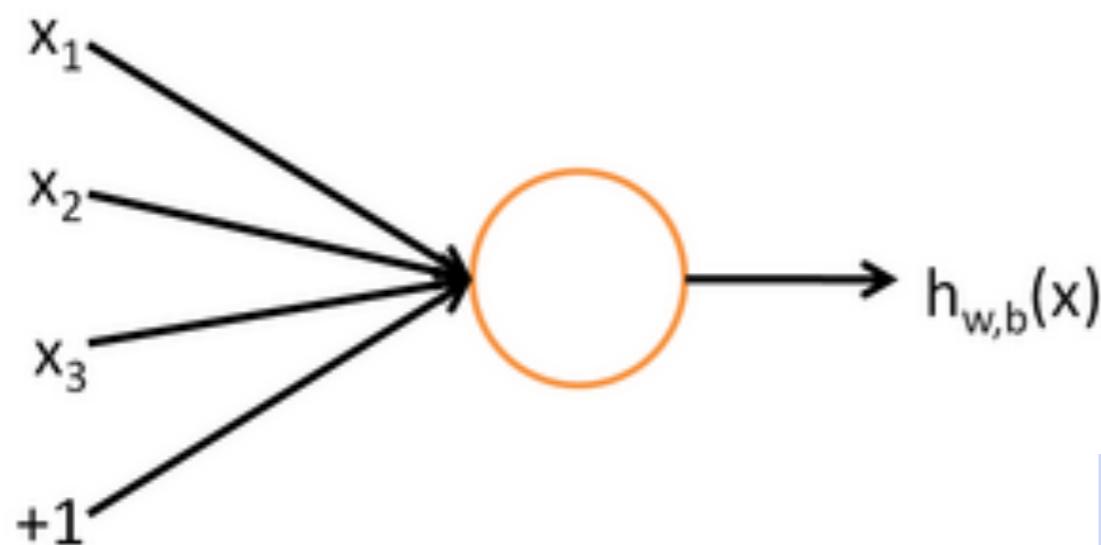
A Neuron

is essentially a binary logistic regression unit

$$h_{w,b}(x) = f(w^\top x + b)$$

b: We can have an “always on” feature, which gives a class prior, or separate it out, as a bias term

$$f(z) = \frac{1}{1 + e^{-z}}$$

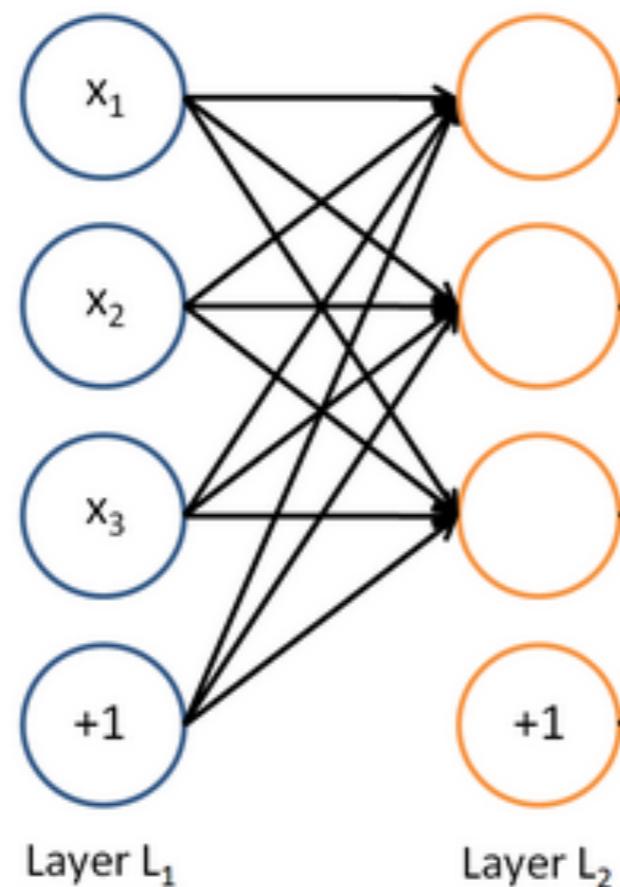


w, *b* are the parameters of this neuron
i.e., this logistic regression model

A Neural Network

= running several logistic regressions at the same time

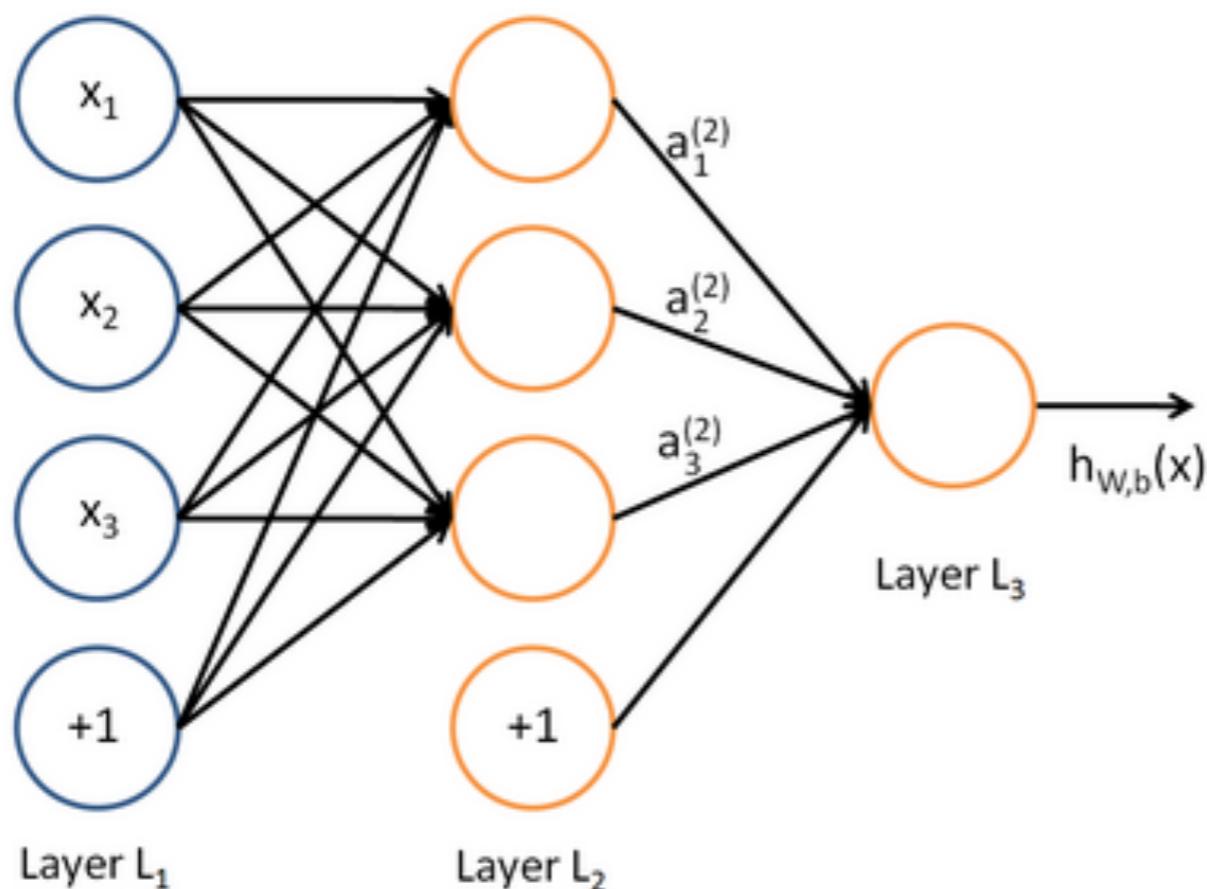
If we feed a vector of inputs through a bunch of logistic regression functions, then we get a vector of outputs ...



A Neural Network

= running several logistic regressions at the same time

... which we can feed into another logistic regression function

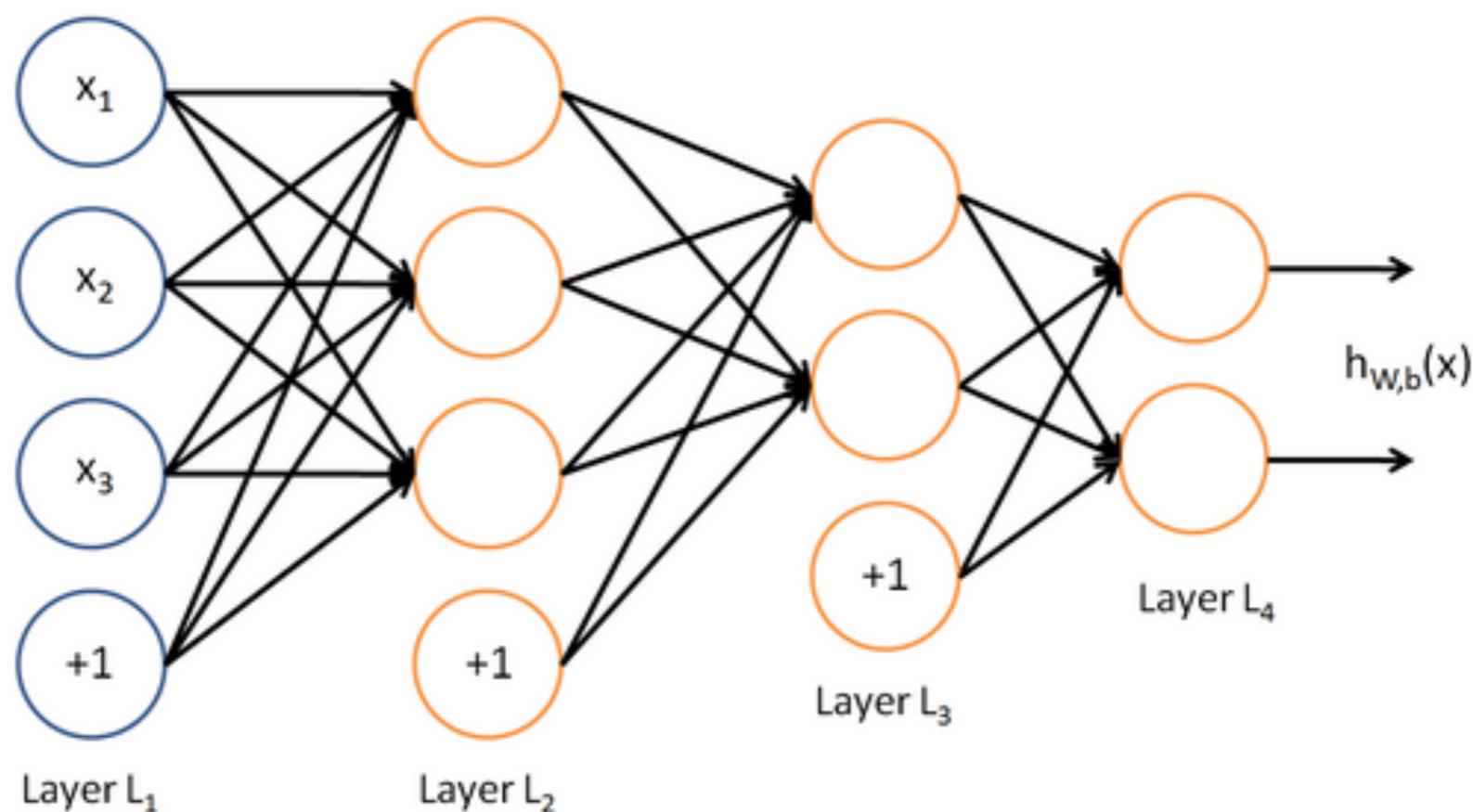


It is the loss function that will direct what the intermediate hidden variables should be, so as to do a good job at predicting the targets for the next layer, etc.

A Neural Network

= running several logistic regressions at the same time

Before we know it, we have a multilayer neural network....



f : Activation Function

We have

$$a_1 = f(W_{11}x_1 + W_{12}x_2 + W_{13}x_3 + b_1)$$

$$a_2 = f(W_{21}x_1 + W_{22}x_2 + W_{23}x_3 + b_2)$$

etc.

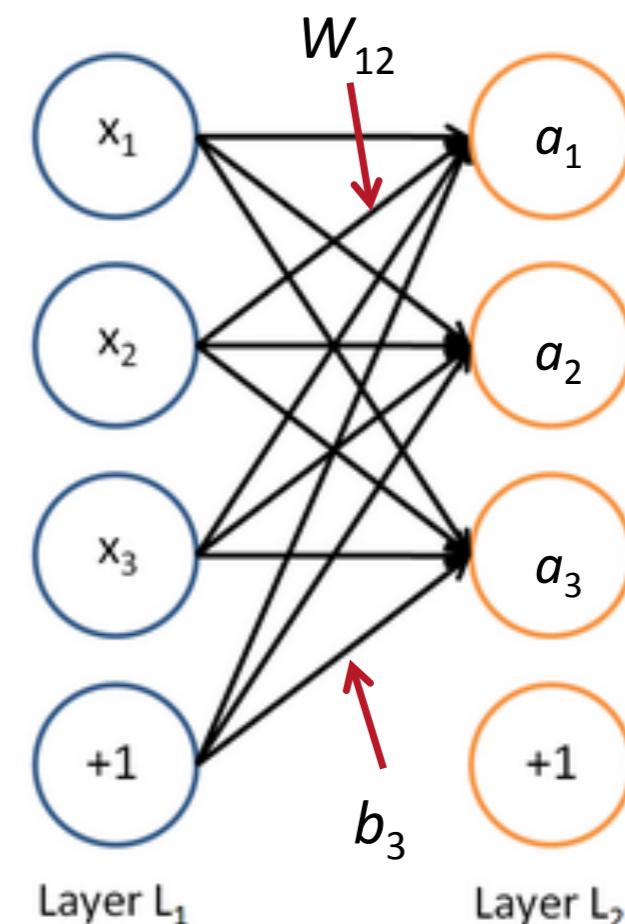
In matrix notation

$$z = Wx + b$$

$$a = f(z)$$

where f is applied element-wise:

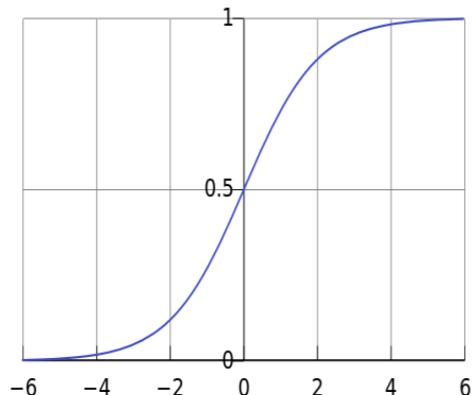
$$f([z_1, z_2, z_3]) = [f(z_1), f(z_2), f(z_3)]$$



Activation Function

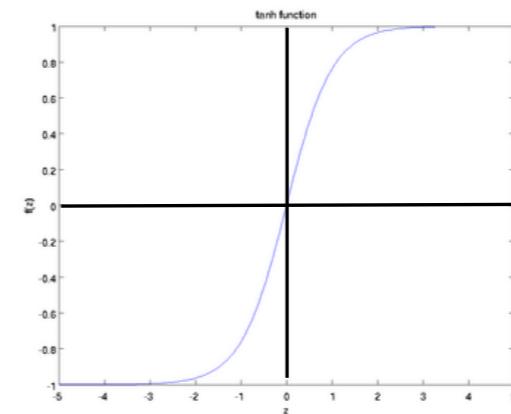
logistic (“sigmoid”)

$$f(z) = \frac{1}{1 + \exp(-z)}.$$



tanh

$$f(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}},$$



$$f'(z) = f(z)(1 - f(z))$$

$$f'(z) = 1 - f(z)^2$$

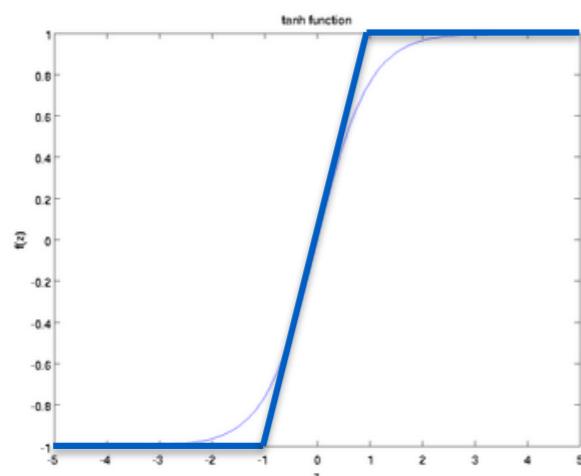
tanh is just a rescaled and shifted sigmoid

$$\tanh(z) = 2\text{logistic}(2z) - 1$$

Activation Function

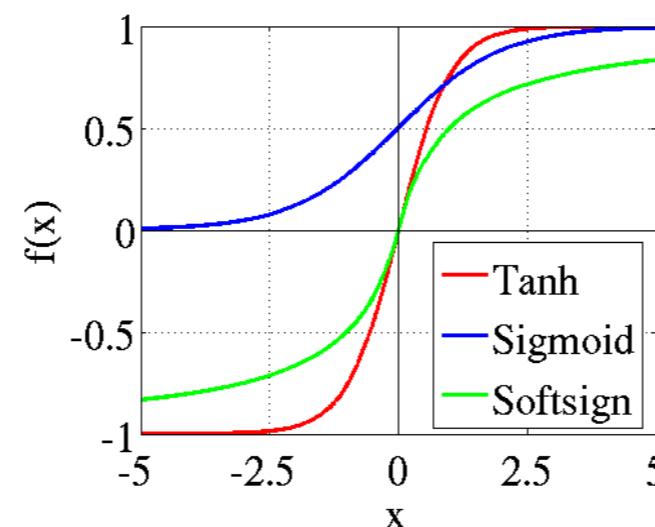
hard tanh

$$\text{HardTanh}(x) = \begin{cases} -1 & \text{if } x < -1 \\ x & \text{if } -1 \leq x \leq 1 \\ 1 & \text{if } x > 1 \end{cases}$$



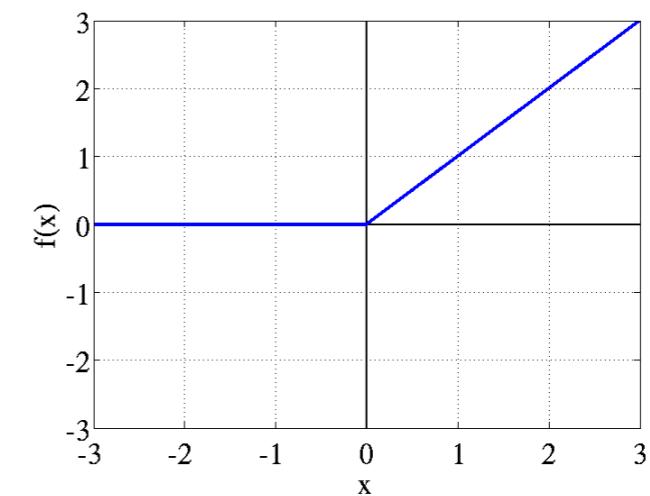
soft sign

$$\text{softsign}(z) = \frac{a}{1+|a|}$$



rectified linear (ReLU)

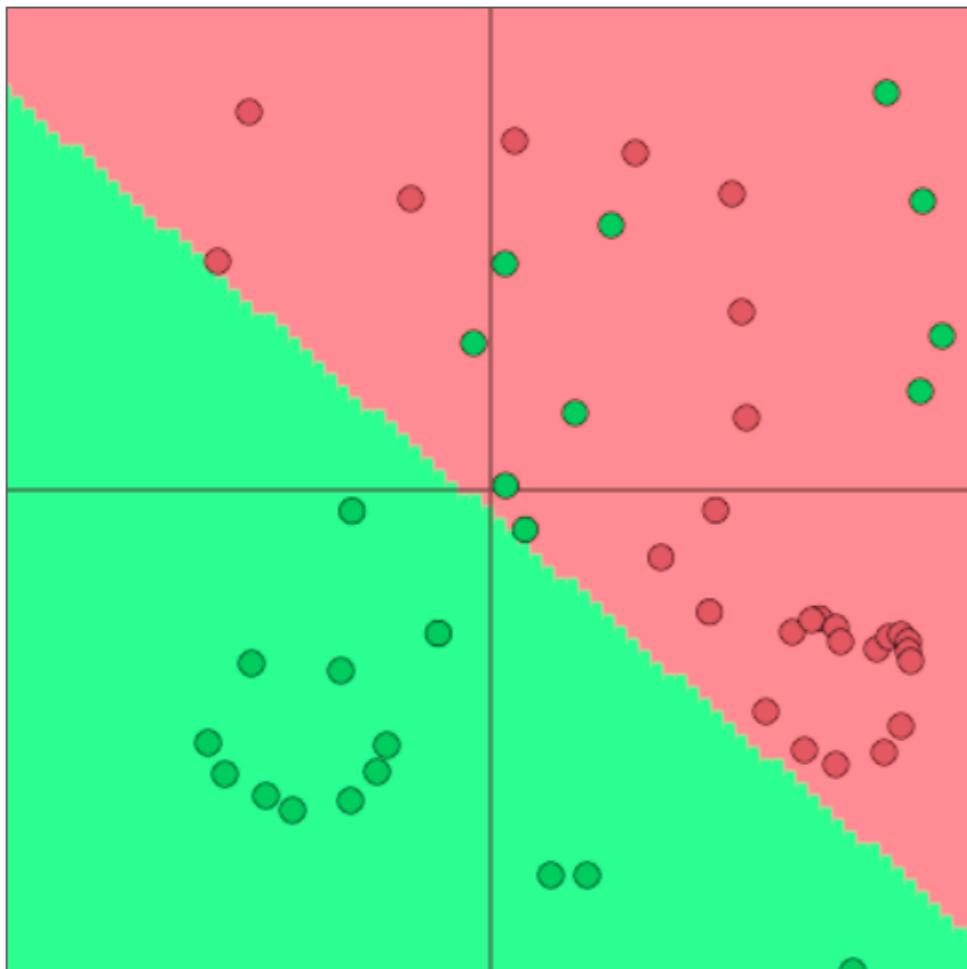
$$\text{rect}(z) = \max(z, 0)$$



- hard tanh similar but computationally cheaper than tanh and saturates hard.
- Glorot and Bengio, *AISTATS 2011* discuss softsign and rectifier

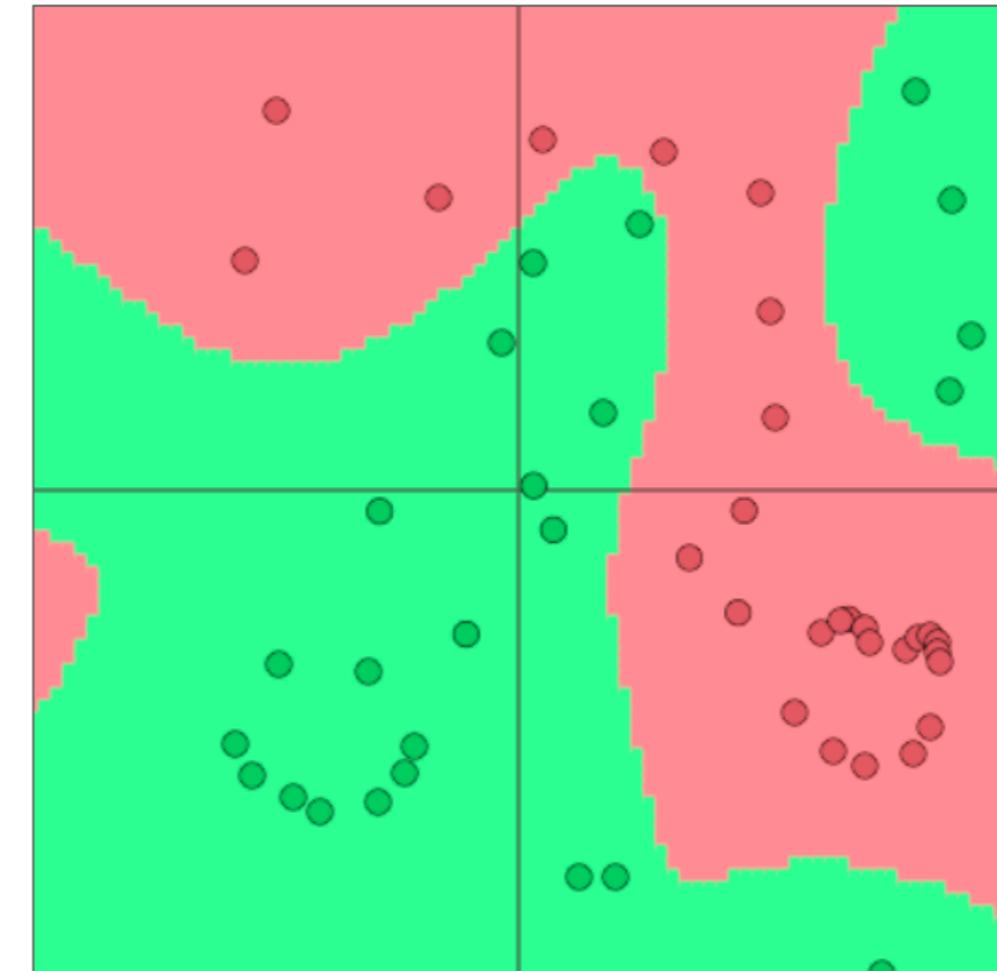
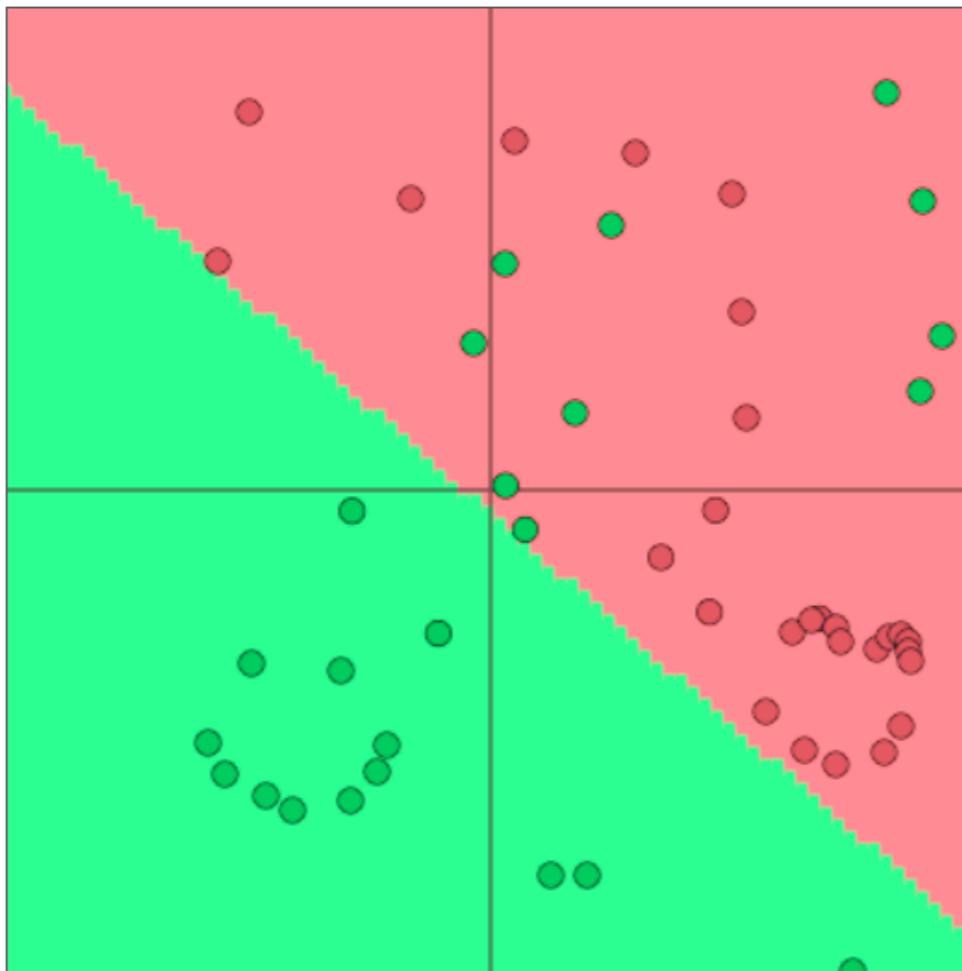
Non-linearity

- Logistic (Softmax) Regression only gives linear decision boundaries

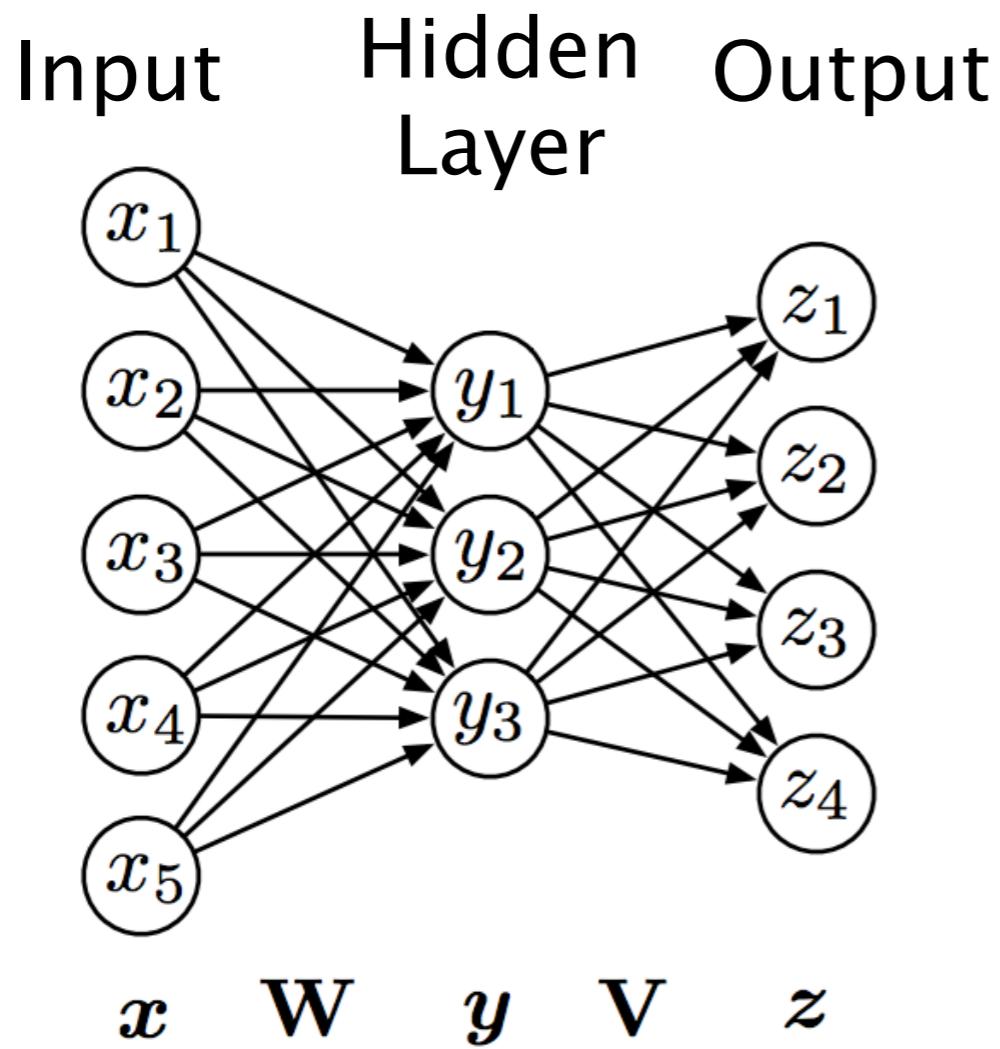


Non-linearity

- Neural networks can learn much more complex functions and nonlinear decision boundaries!



Non-linearity



$$\mathbf{y} = g(\mathbf{W}\mathbf{x} + \mathbf{b})$$

$$\mathbf{z} = g(\mathbf{V}g(\mathbf{W}\mathbf{x} + \mathbf{b}) + \mathbf{c})$$

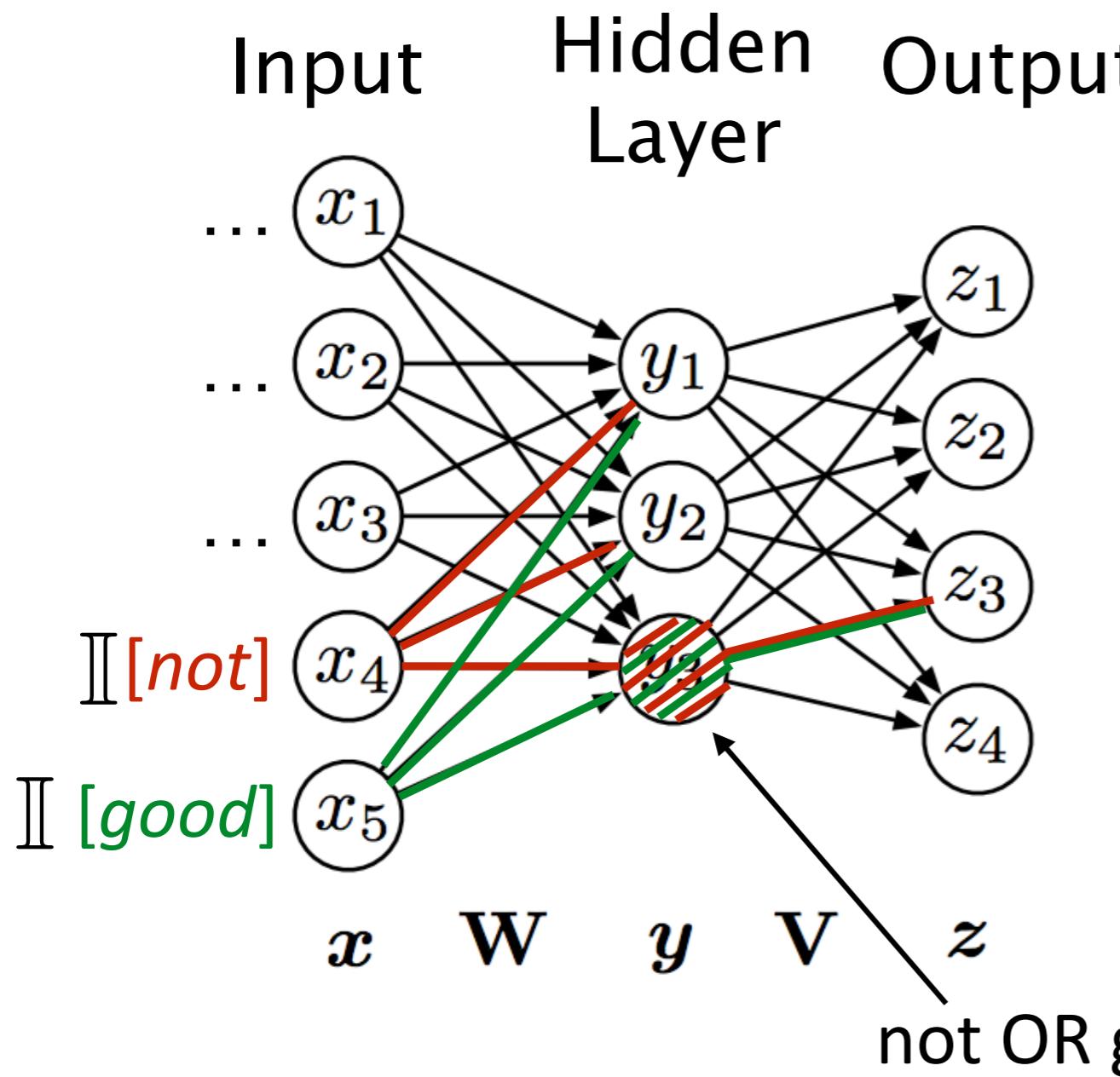
output of first layer

With no nonlinearity:

$$\mathbf{z} = \mathbf{V}\mathbf{W}\mathbf{x} + \mathbf{V}\mathbf{b} + \mathbf{c}$$

Equivalent to $\mathbf{z} = \mathbf{U}\mathbf{x} + \mathbf{d}$

Non-linearity



Nodes in the hidden layer can learn interactions or conjunctions of features

$$y = -2x_1 - x_2 + 2 \tanh(x_1 + x_2)$$

What about Word2vec
(Skip-gram and CBOW)?

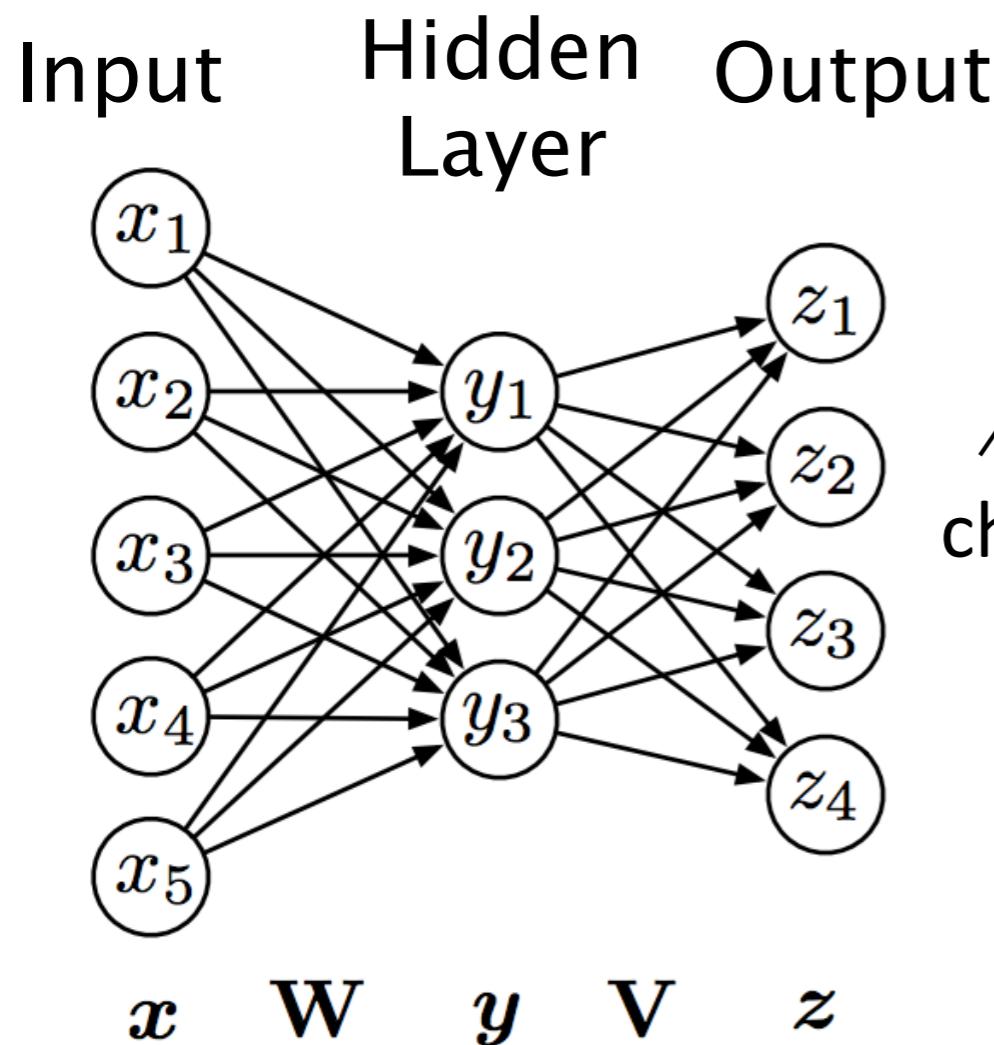
So, what about Word2vec (Skip-gram and CBOW)?

It is not deep learning — but “shallow” neural networks.

It is — in fact — a log-linear model (softmax regression).

So, it is faster over larger dataset yielding better embeddings.

Learning Neural Networks



$$\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$$

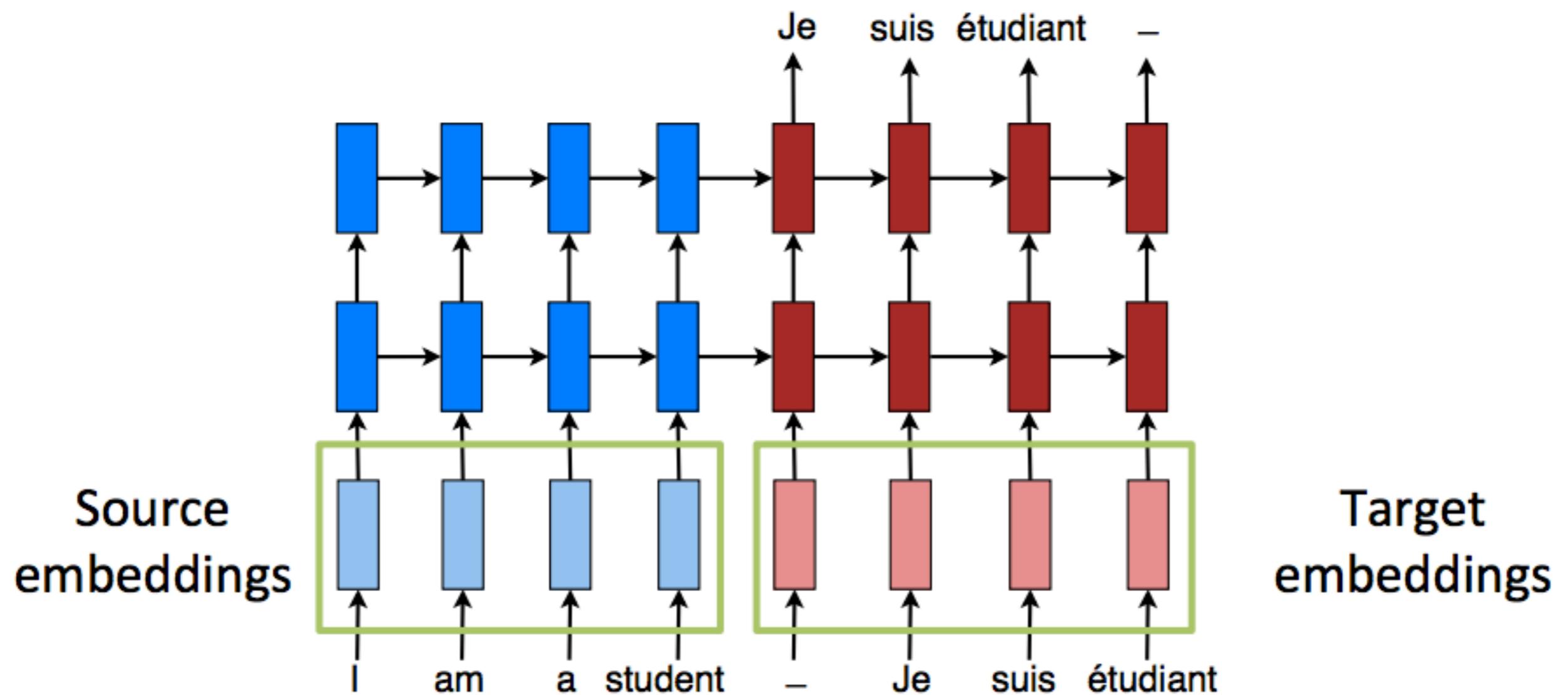
change in output w.r.t. hidden
change in hidden w.r.t. input
change in output w.r.t. input

Computing these looks like running this network in reverse (backpropagation)

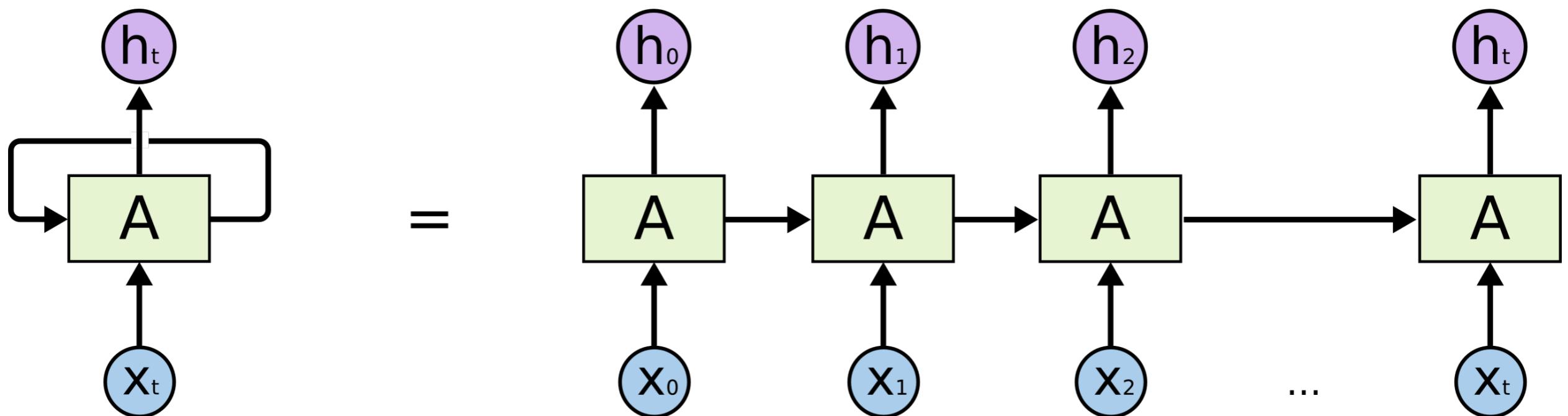
Strategy for Successful NNs

- Select network structure appropriate for problem
 - Structure: Single words, fixed windows, sentence based, document level; bag of words, recursive vs. recurrent, CNN, ...
 - Nonlinearity
- Check for implementation bugs with gradient checks
- Parameter initialization
- Optimization tricks
- Check if the model is powerful enough to overfit
 - If not, change model structure or make model “larger”
 - If you can overfit: regularize

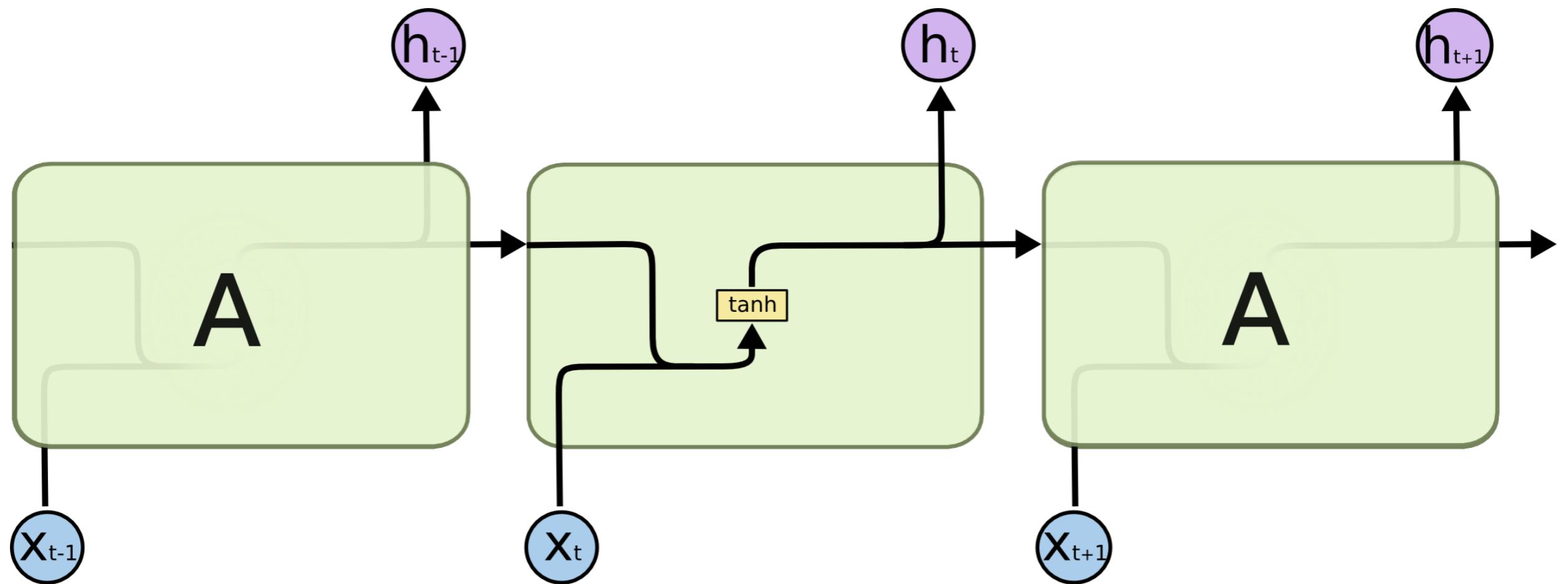
Neural Machine Translation



Recurrent Neural Network (RNN)

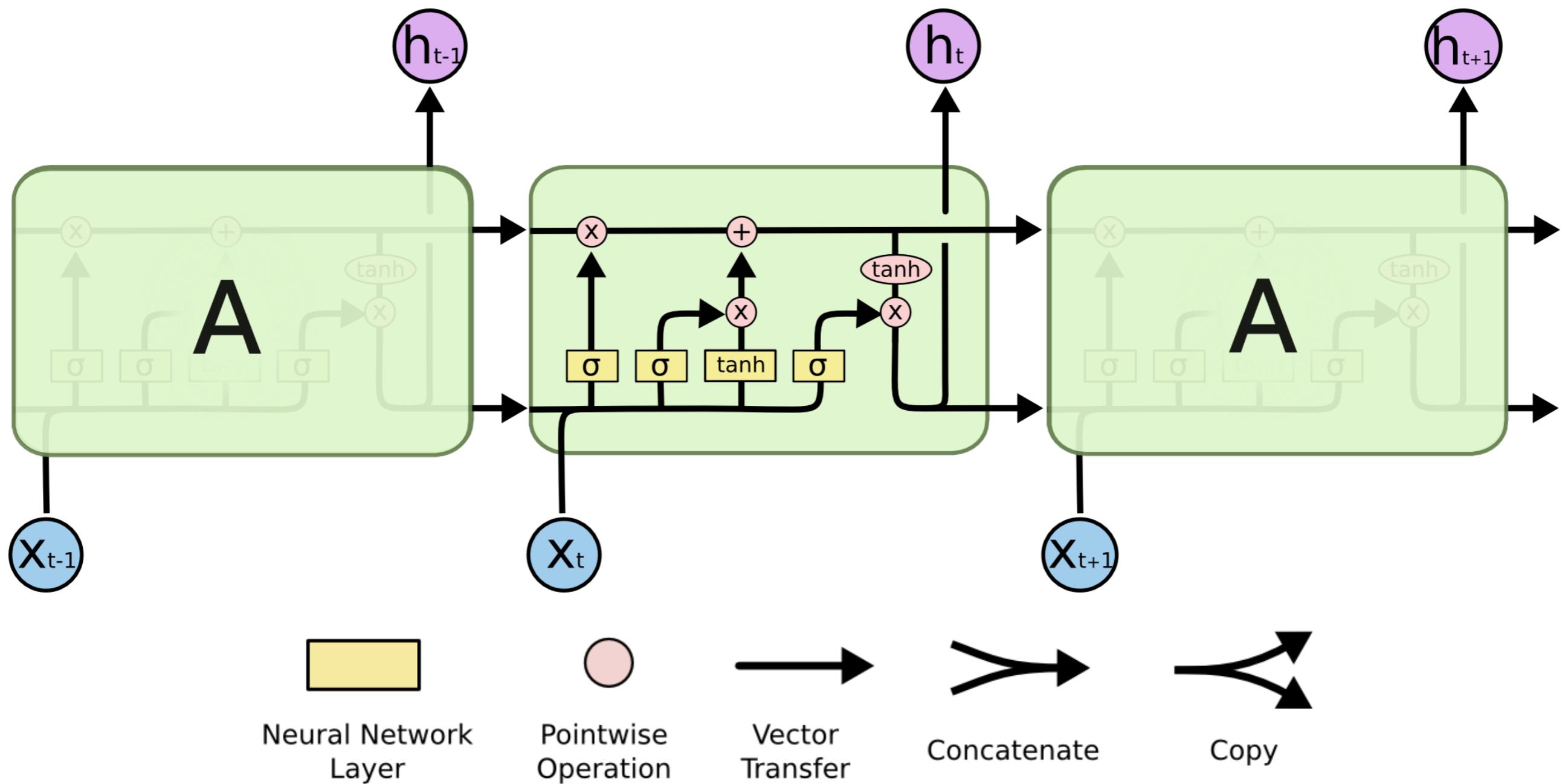


Recurrent Neural Network (RNN)

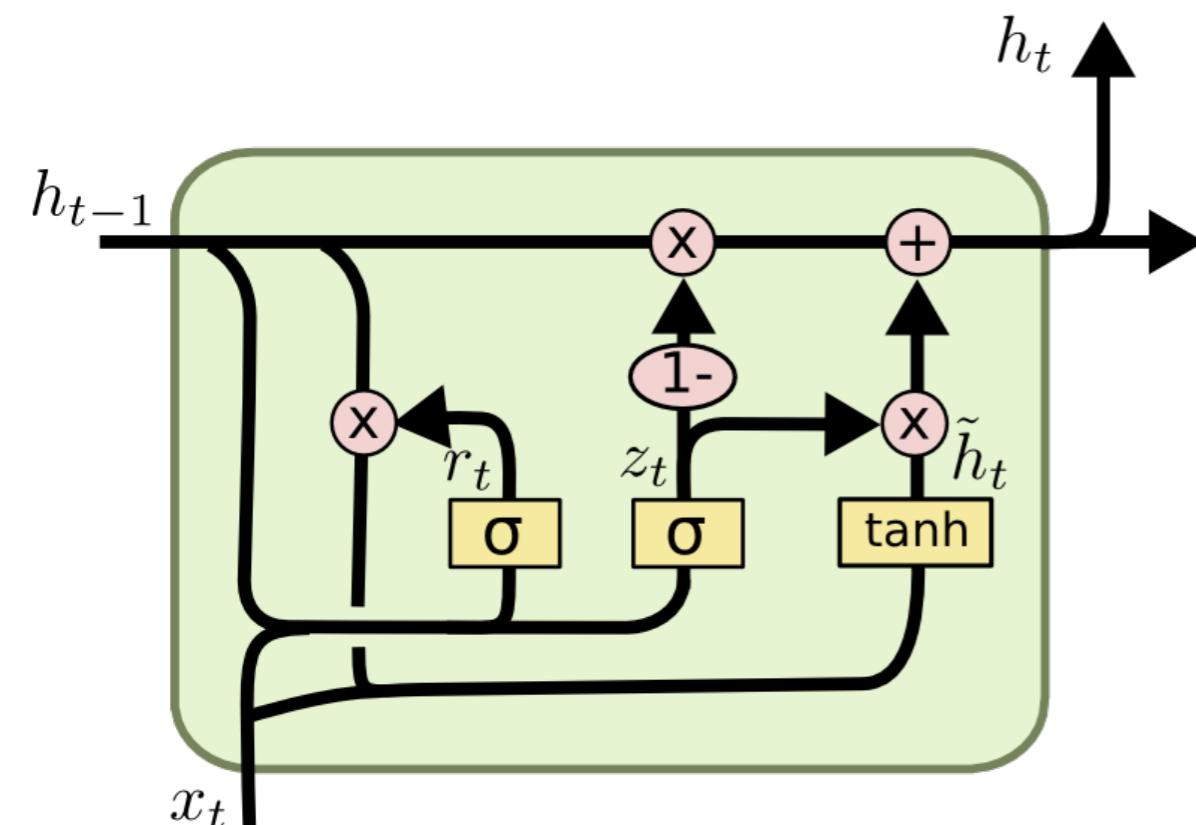


Long Short-Term Memory Networks (LSTM)

(Hochreiter & Schmidhuber, 1997)



Long Short-Term Memory Networks (LSTM)

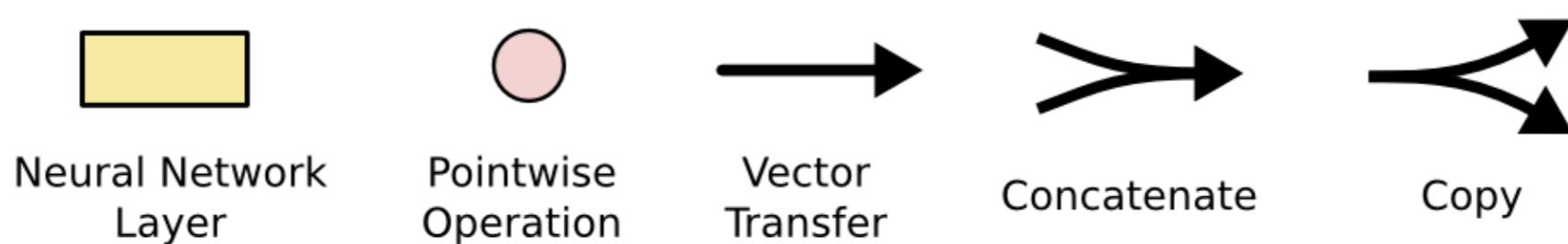


$$z_t = \sigma (W_z \cdot [h_{t-1}, x_t])$$

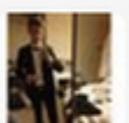
$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t])$$

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

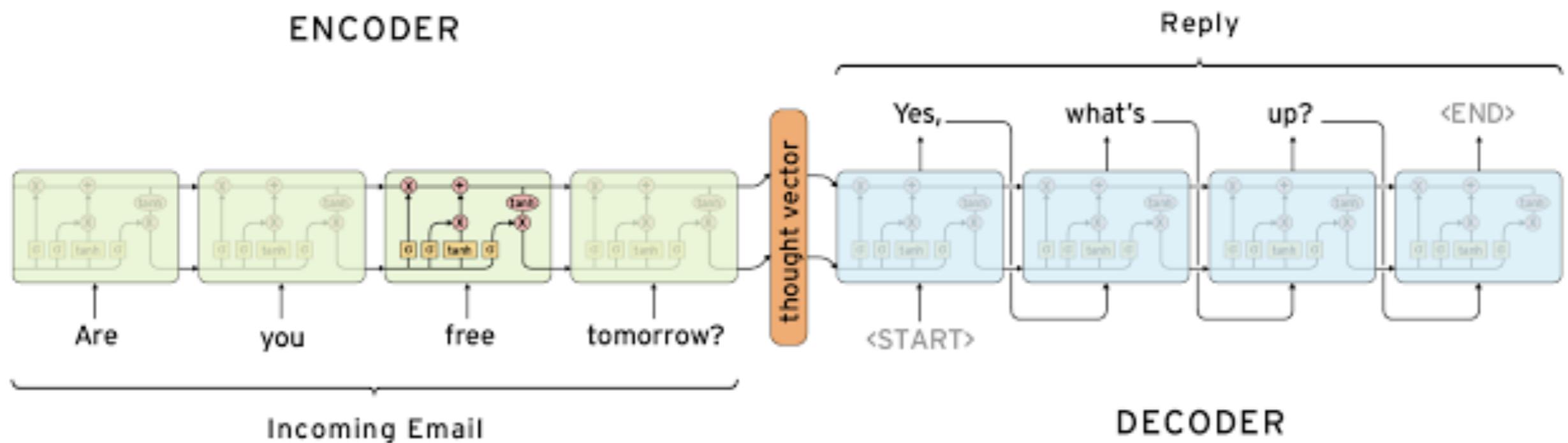
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$



Twitter Conversation Data

	Tesco Mobile @tescomobile	10 Nov
	@RiccardoEspaa7 She's clearly going through a rough time in life.	
	Details	
	Riccardo Esposito @RiccardoEspaa7	10 Nov
	@tescomobile I know either that or shit is going insane! There is nothing wrong with tesco mobile	
	Details	
	Tesco Mobile @tescomobile	10 Nov
	@RiccardoEspaa7 Together Riccardo, we'll make this world a better place.	
	Details	
	Riccardo Esposito @RiccardoEspaa7	10 Nov
	@tescomobile me and you against the world	
	Details	
	Tesco Mobile @tescomobile	10 Nov
	@RiccardoEspaa7 Yeah baby.	
	Details	
	Riccardo Esposito @RiccardoEspaa7	10 Nov
	@tescomobile no one can stop us	
	Details	

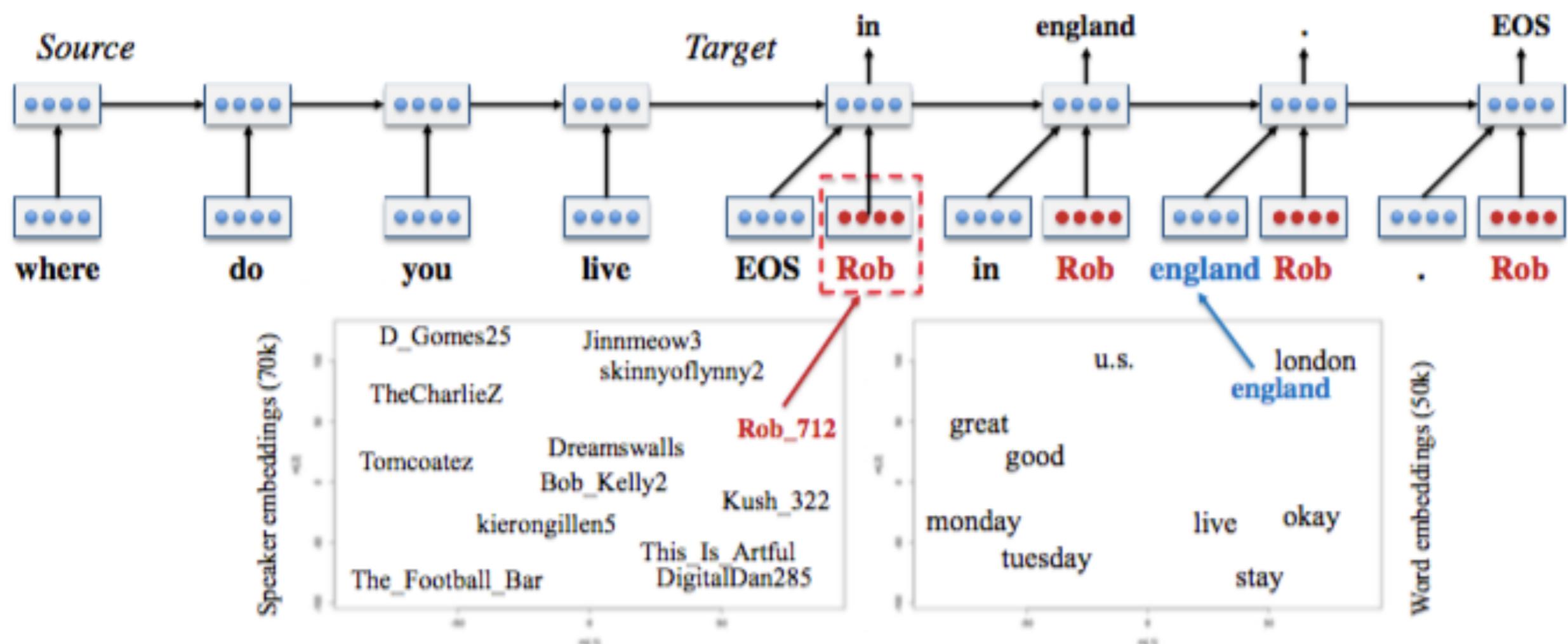
Neural Conversation



Source: Google's blog

Neural Conversation

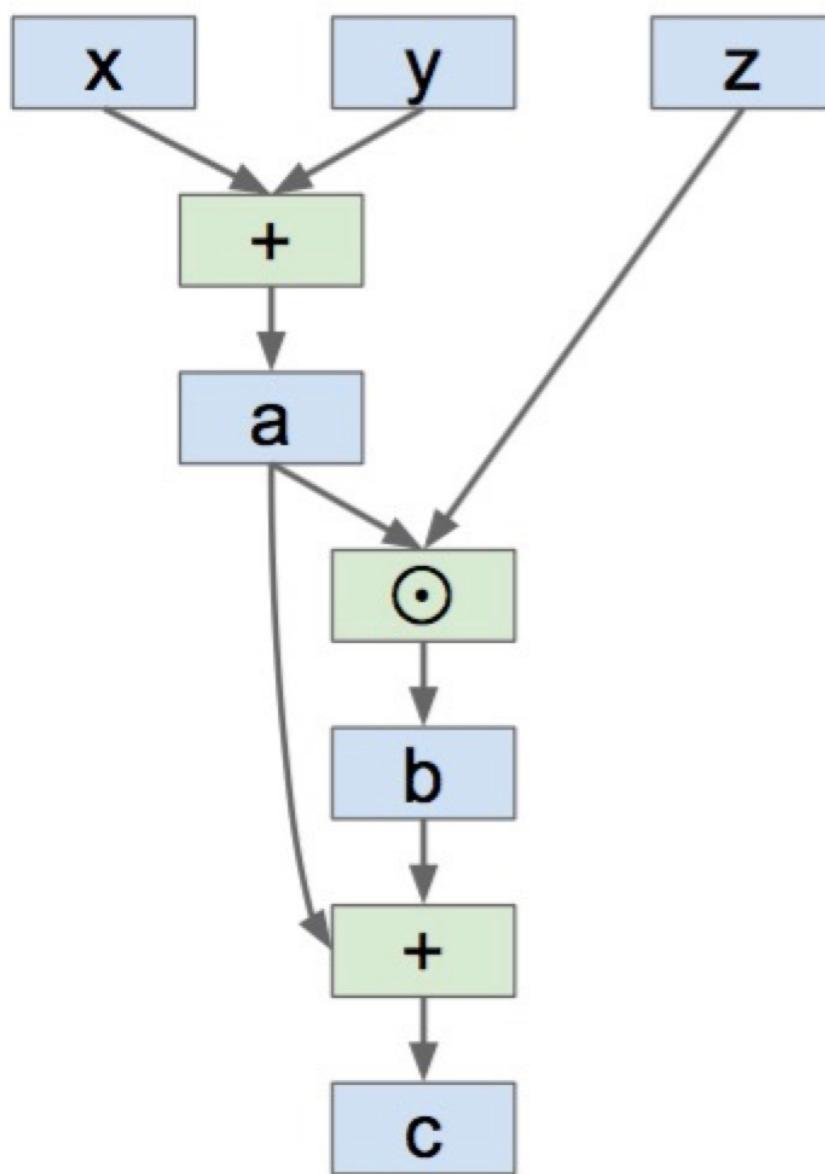
modeling speakers



Neural Network Toolkits

- **Tensorflow**: <https://www.tensorflow.org/>
 - By Google, actively maintained, bindings for many languages
- **Theano**: <http://deeplearning.net/software/theano>
 - University of Montreal, less and less maintained
- **Torch**: <http://torch.ch/>
 - Facebook AI Research, Lua
- **DyNet**: <https://github.com/clab/dynet>
 - CMU and others, dynamic structures that change for every training instance
- **Caffe**: <http://caffe.berkeleyvision.org/>
 - UC Berkeley, for vision

Neural Network Toolkits



```
import theano
import theano.tensor as T

# Define symbolic variables
x = T.matrix('x')
y = T.matrix('y')
z = T.matrix('z')

# Compute some other values symbolically
a = x + y
b = a * z
c = a + b

# Compile a function that computes c
f = theano.function(
    inputs=[x, y, z],
    outputs=c
)

# Evaluate the compiled function
# on some real values
xx = np.random.randn(4, 5)
yy = np.random.randn(4, 5)
zz = np.random.randn(4, 5)
print f(xx, yy, zz)

# Repeat the same computation
# explicitly using numpy ops
aa = xx + yy
bb = aa * zz
cc = aa + bb
```

Compile a function that produces c from x , y , z (generates code)





Instructor: Wei Xu

www.cis.upenn.edu/~xwe/

Course Website: socialmedia-class.org