

# Natural Language Processing

NN review and Neural LMs

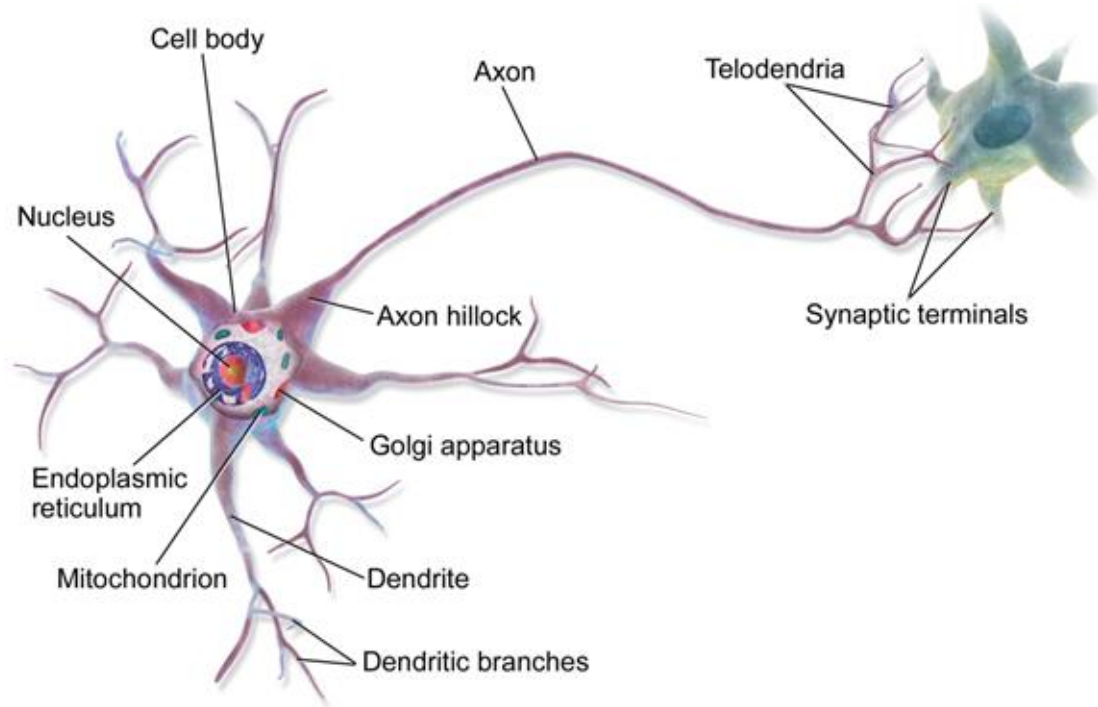
Luke Zettlemoyer

[lsz@cs.washington.edu](mailto:lsz@cs.washington.edu)

# Readings

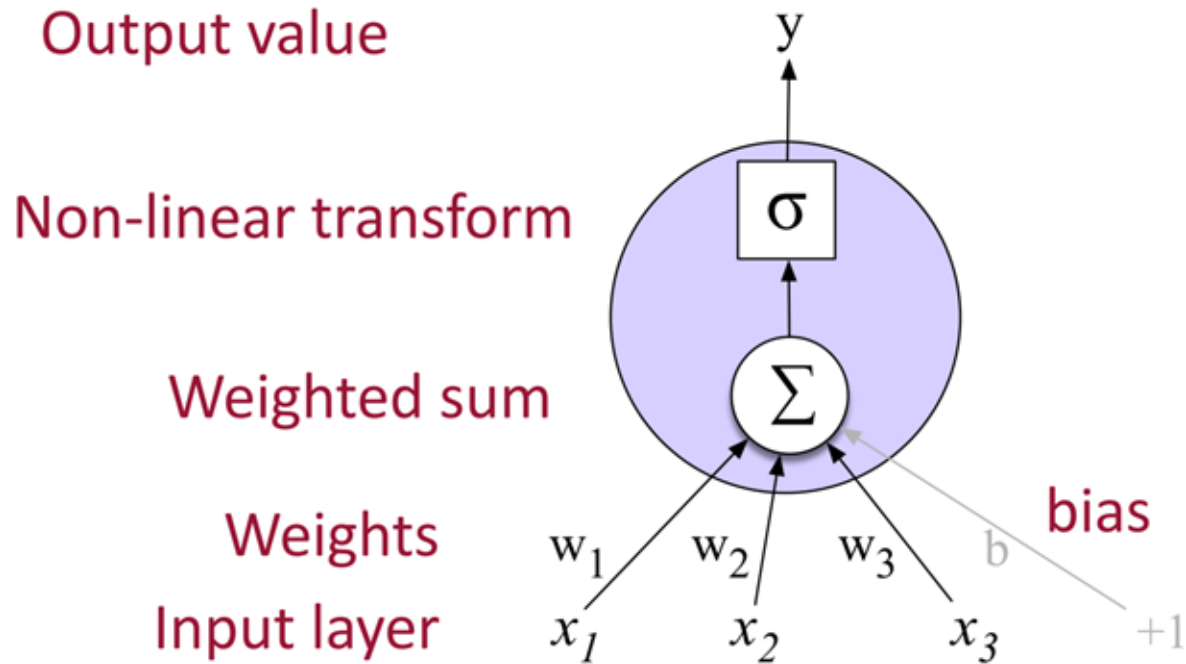
- Neutral networks chapters in J&M 6:
  - <https://web.stanford.edu/~jurafsky/slp3/6.pdf>
- Hundreds of blog posts and tutorials

# This is in your brain



By BruceBlais - Own work, CC BY 3.0,  
<https://commons.wikimedia.org/w/index.php?curid=28761830>

# Neural Network Unit (this is not in your brain)



# Neural unit

- Take weighted sum of inputs, plus a bias

$$z = b + \sum_i w_i x_i$$

$$z = w \cdot x + b$$

- Instead of just using  $z$ , we'll apply a nonlinear activation function  $f$ :

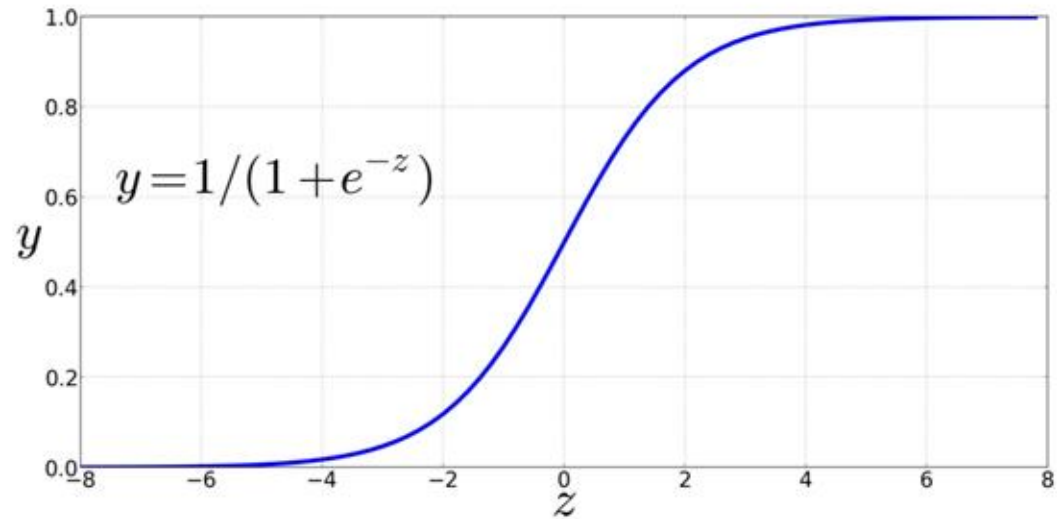
$$y = a = f(z)$$

# Non-Linear Activation Functions

- We've already seen the sigmoid for logistic regression:

Sigmoid

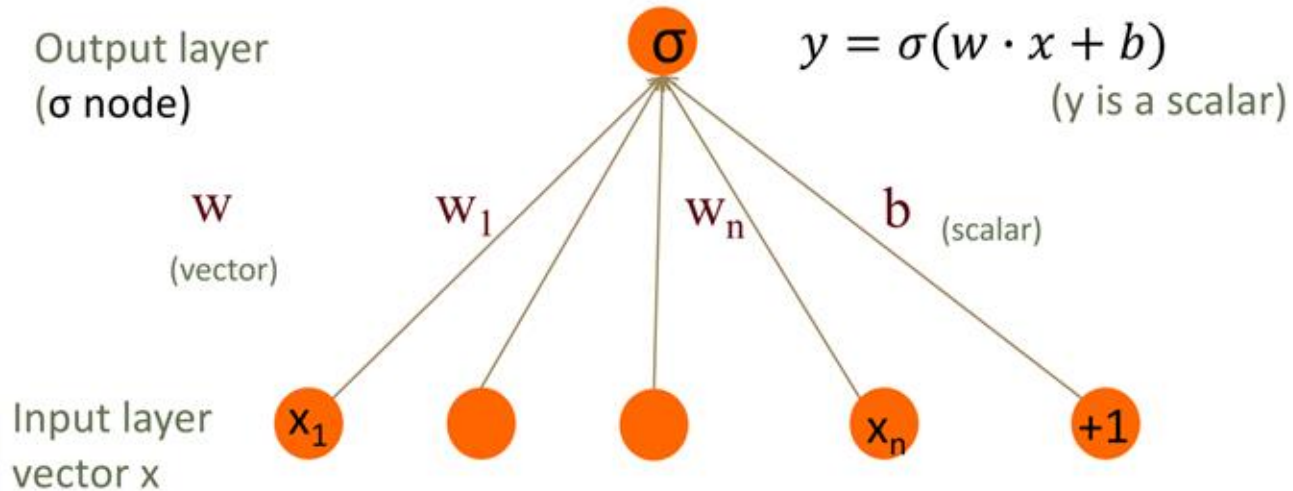
$$y = \sigma(z) = \frac{1}{1 + e^{-z}}$$



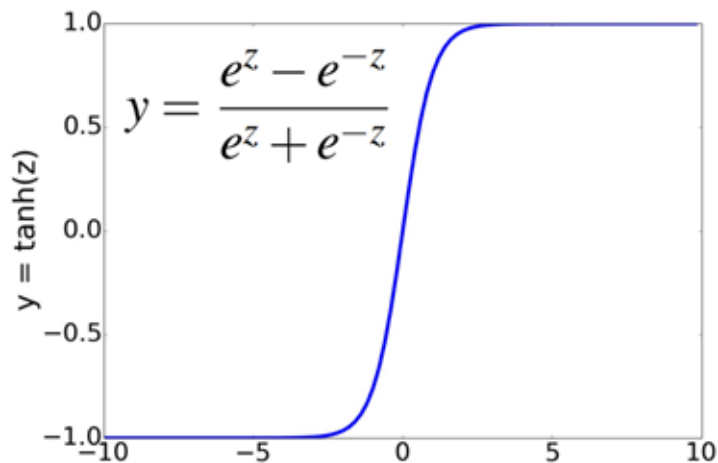
# Final function the unit is computing

$$y = \sigma(w \cdot x + b) = \frac{1}{1 + \exp(-(w \cdot x + b))}$$

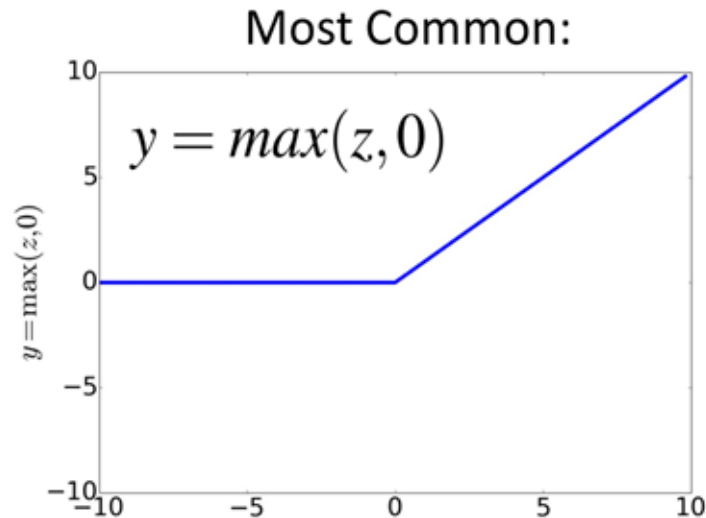
# Binary Logistic Regression as a 1-layer network



# Non-Linear Activation Functions besides sigmoid



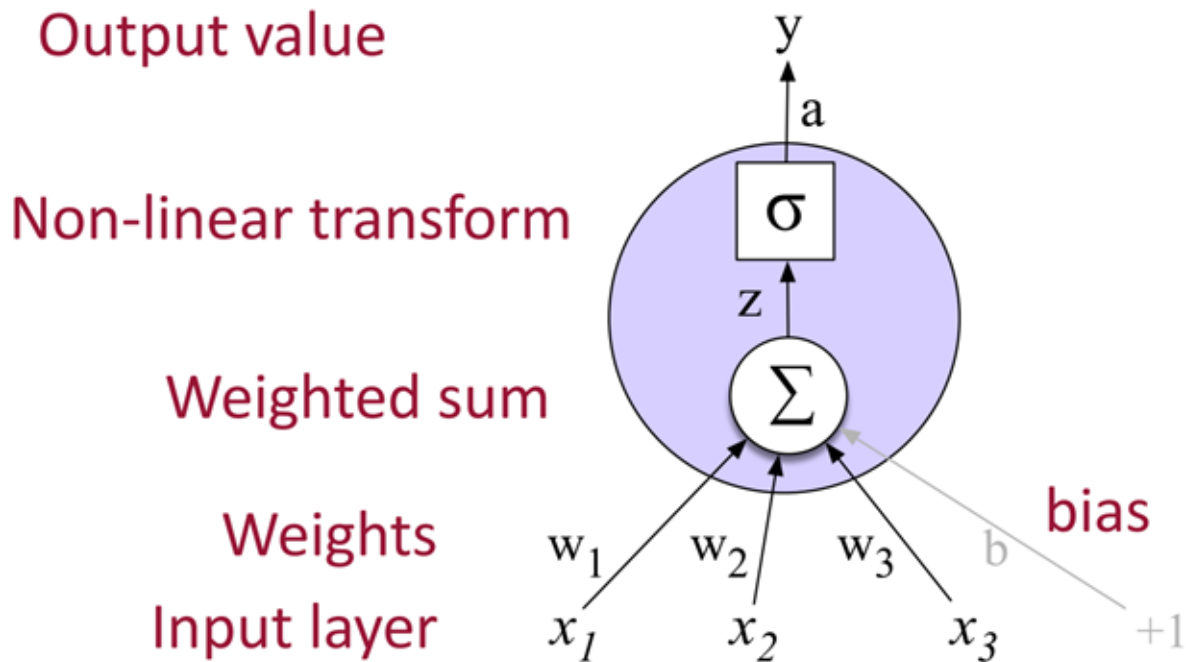
**tanh**



**ReLU**

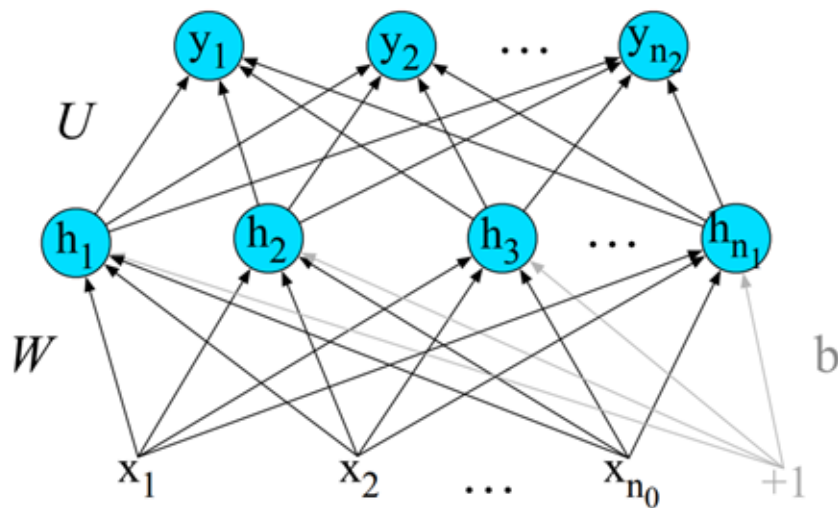
**Rectified Linear Unit**

# Final unit again



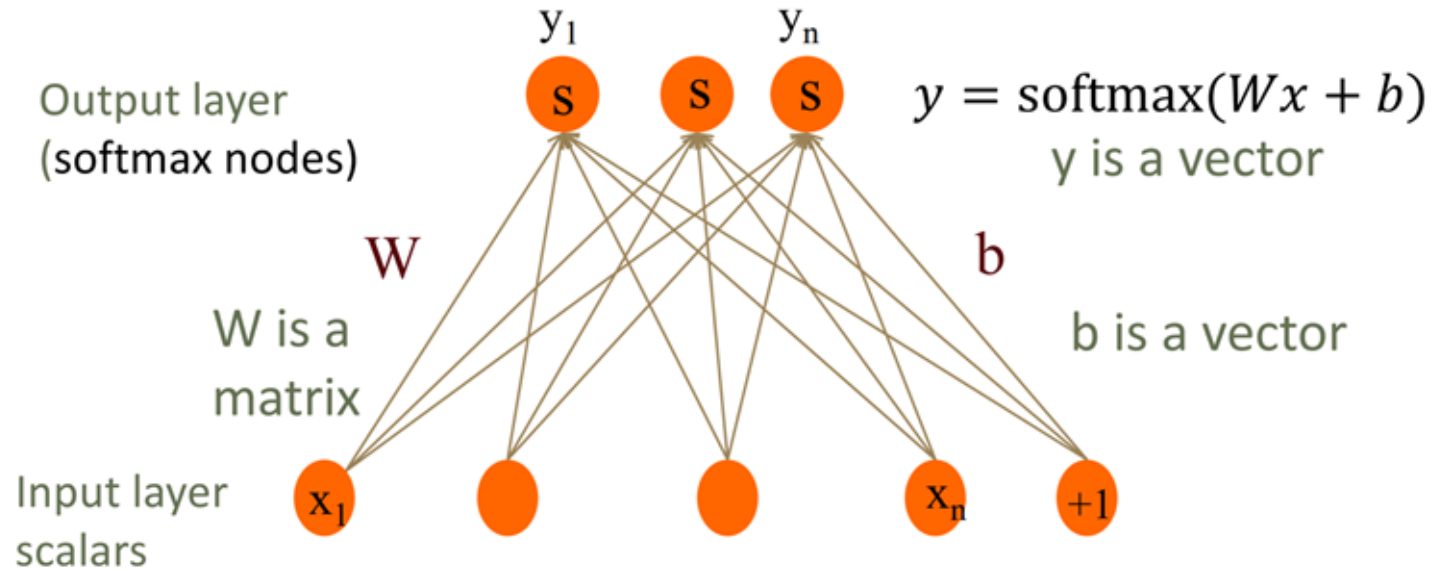
# Feedforward Neural Networks

- Can also be called **multi-layer perceptrons (or MLPs)** for historical reasons
  - (we don't count the input layer in counting layers!)



# Multinomial Logistic Regression as a 1-layer Network

Fully connected single layer network



# softmax: a generalization of sigmoid

- For a vector  $z$  of dimensionality  $k$ , the softmax is:

$$\text{softmax}(z) = \left[ \frac{\exp(z_1)}{\sum_{i=1}^k \exp(z_i)}, \frac{\exp(z_2)}{\sum_{i=1}^k \exp(z_i)}, \dots, \frac{\exp(z_k)}{\sum_{i=1}^k \exp(z_i)} \right]$$

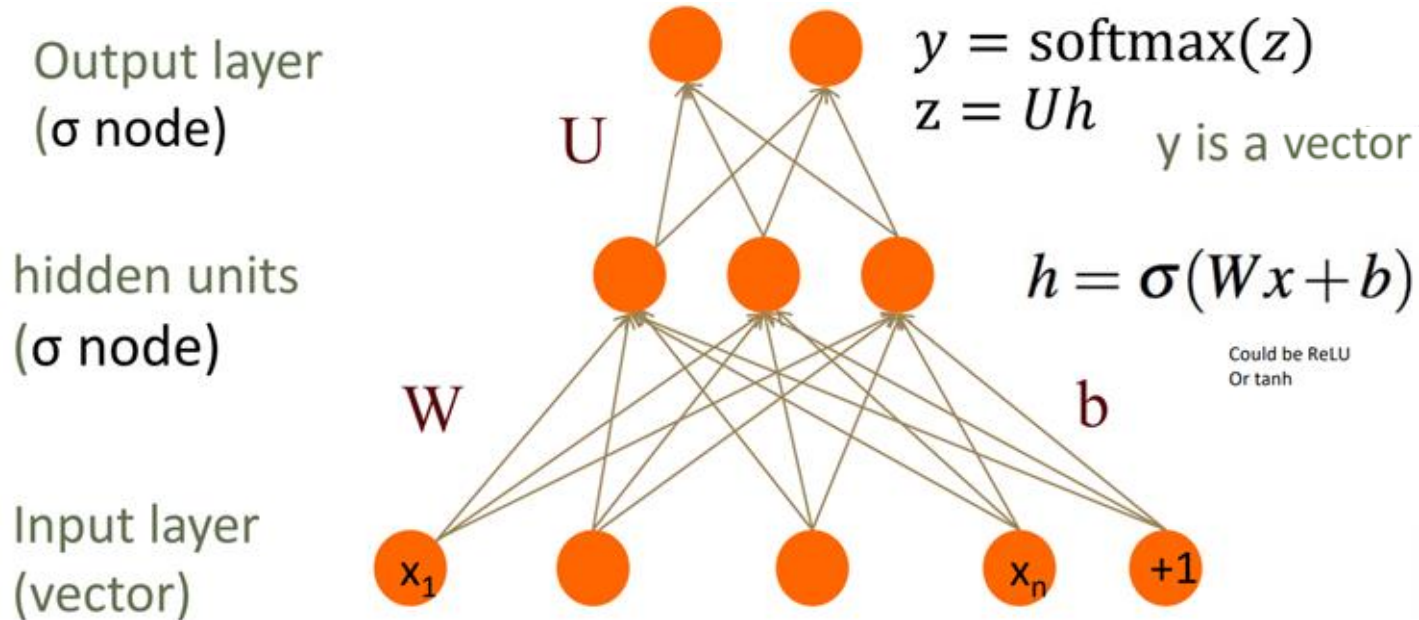
$$\text{softmax}(z_i) = \frac{\exp(z_i)}{\sum_{j=1}^k \exp(z_j)} \quad 1 \leq i \leq k$$

Example:

$$z = [0.6, 1.1, -1.5, 1.2, 3.2, -1.1]$$

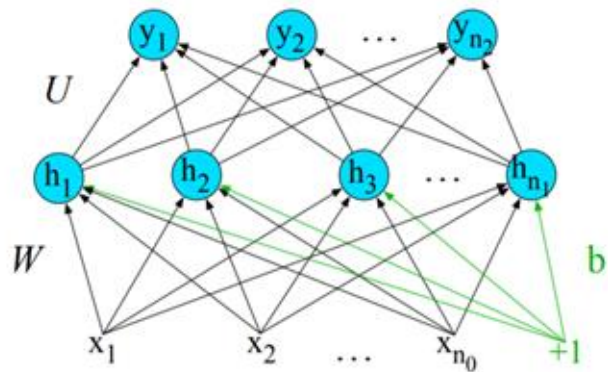
$$\text{softmax}(z) = [0.055, 0.090, 0.006, 0.099, 0.74, 0.010]$$

# Two-Layer Network with softmax output

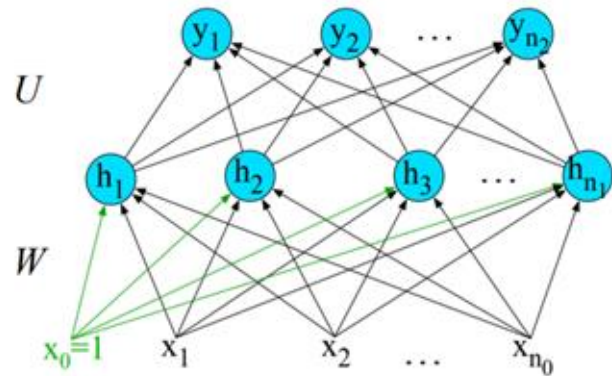


# Replacing the bias unit

Instead of:



We'll do this:



# Learning the weights

- Cross-entropy loss
- Backpropagation algorithm

---

## Algorithm 1 Backpropagation Algorithm

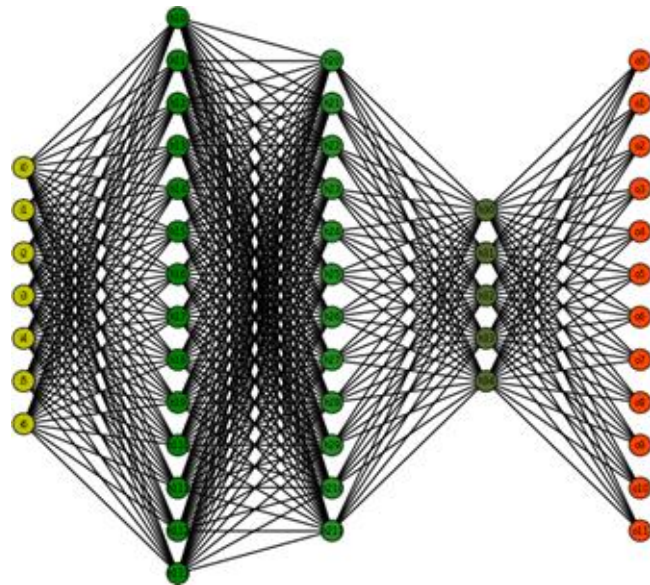
---

```

1: procedure TRAIN
2:    $X \leftarrow$  Training Data Set of size  $m \times n$ 
3:    $y \leftarrow$  Labels for records in  $X$ 
4:    $w \leftarrow$  The weights for respective layers
5:    $l \leftarrow$  The number of layers in the neural network,  $1 \dots L$ 
6:    $D_{ij}^{(l)} \leftarrow$  The error for all  $l, i, j$ 
7:    $t_{ij}^{(l)} \leftarrow 0$ . For all  $l, i, j$ 
8:   For  $i = 1$  to  $m$ 
9:      $a^l \leftarrow \text{feedforward}(x^{(i)}, w)$ 
10:     $d^l \leftarrow a(L) - y(i)$ 
11:     $t_{ij}^{(l)} \leftarrow t_{ij}^{(l)} + a_j^{(l)} \cdot t_i^{l+1}$ 
12:    if  $j \neq 0$  then
13:       $D_{ij}^{(l)} \leftarrow \frac{1}{m} t_{ij}^{(l)} + \lambda w_{ij}^{(l)}$ 
14:    else
15:       $D_{ij}^{(l)} \leftarrow \frac{1}{m} t_{ij}^{(l)}$ 
16:    where  $\frac{\partial}{\partial w_{ij}^{(l)}} J(w) = D_{ij}^{(l)}$ 

```

---



# Applying neural networks to NLP tasks



# Use cases for feedforward networks

- Word representations
- Text classification
- Language modeling

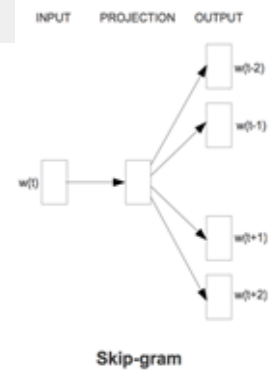
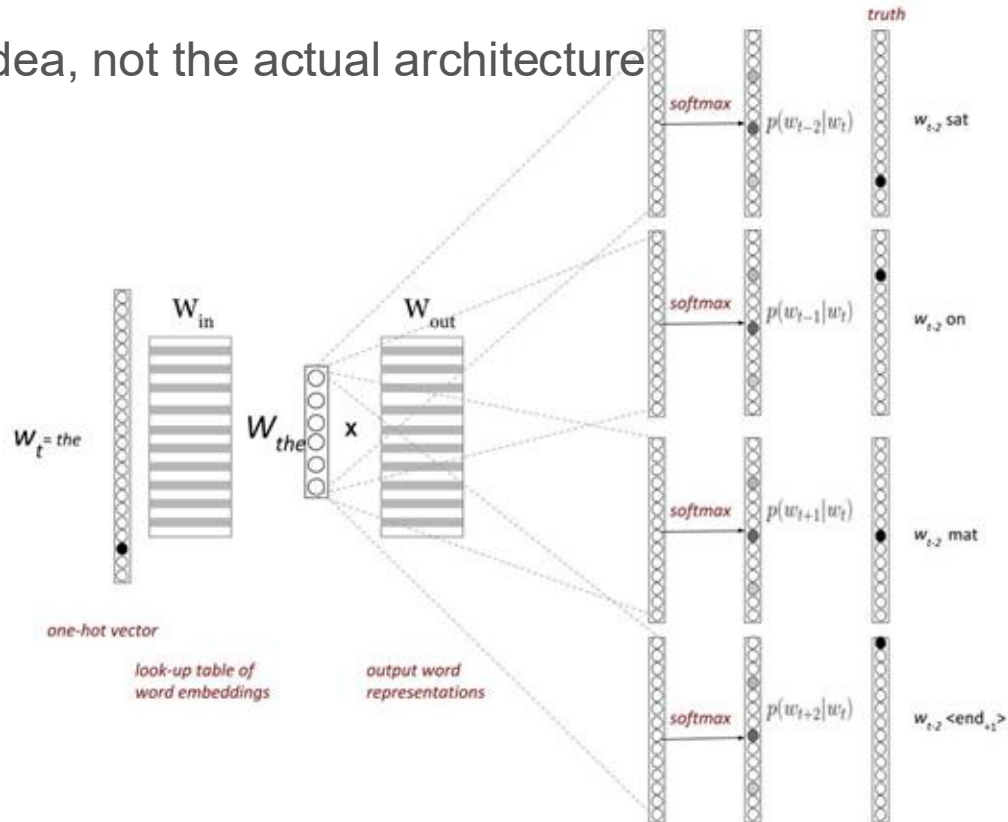
State of the art systems use more powerful neural architectures (we will learn transformers architectures on Friday), but simpler models are useful to consider!

# word2vec

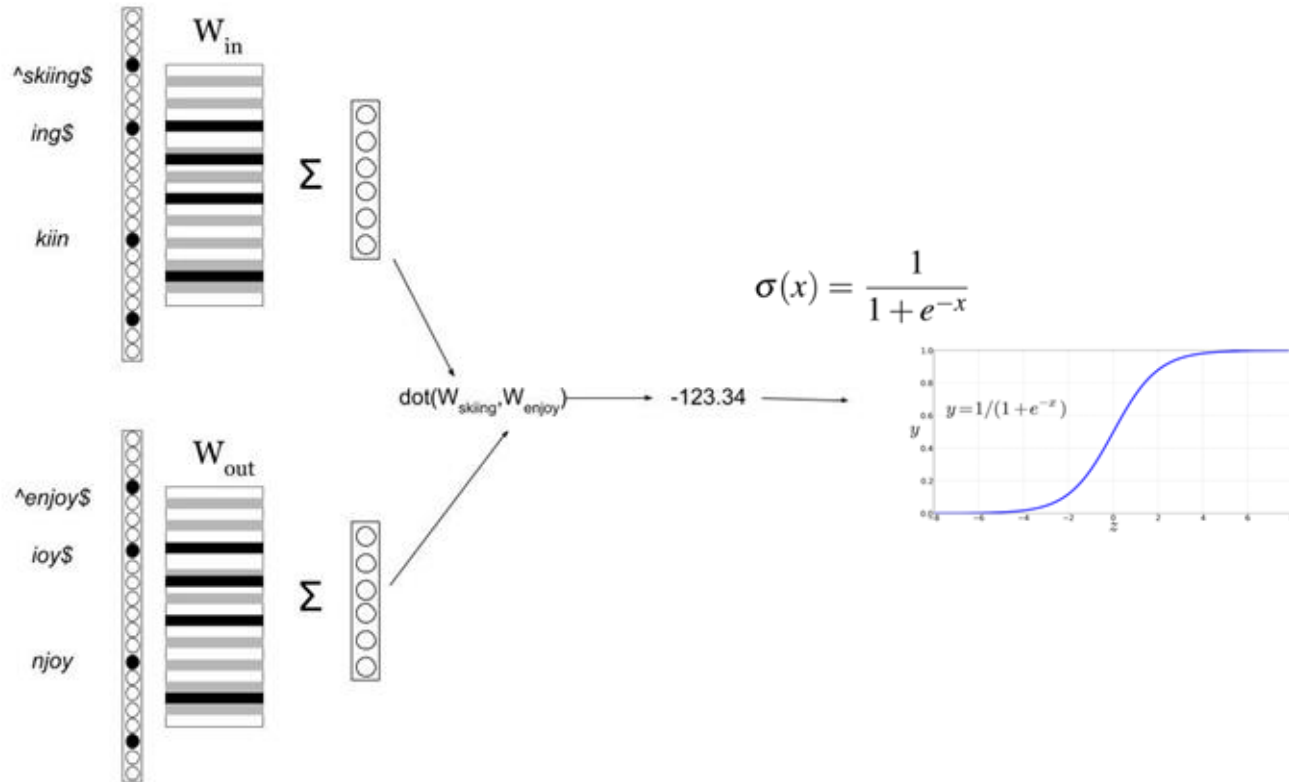
- Instead of counting how often each word  $w$  occurs near "apricot"
  - Train a classifier on a binary **prediction** task:
    - Is  $w$  likely to show up near "apricot"?
- We don't actually care about this task
  - But we'll take the learned classifier weights as the word embeddings
- Big idea: **self-supervision**:
  - A word  $c$  that occurs near apricot in the corpus acts as the gold "correct answer" for supervised learning
  - No need for human labels
  - Bengio et al. (2003); Collobert et al. (2011)

# Skip-gram Prediction

- Conceptual idea, not the actual architecture



# FastText



# BERT

## **BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding**

**Jacob Devlin   Ming-Wei Chang   Kenton Lee   Kristina Toutanova**

Google AI Language

{jacobdevlin, mingweichang, kentonl, kristout}@google.com

<https://ai.googleblog.com/2018/11/open-sourcing-bert-state-of-art-pre.html>

# Preview: BERT as a text classification problem

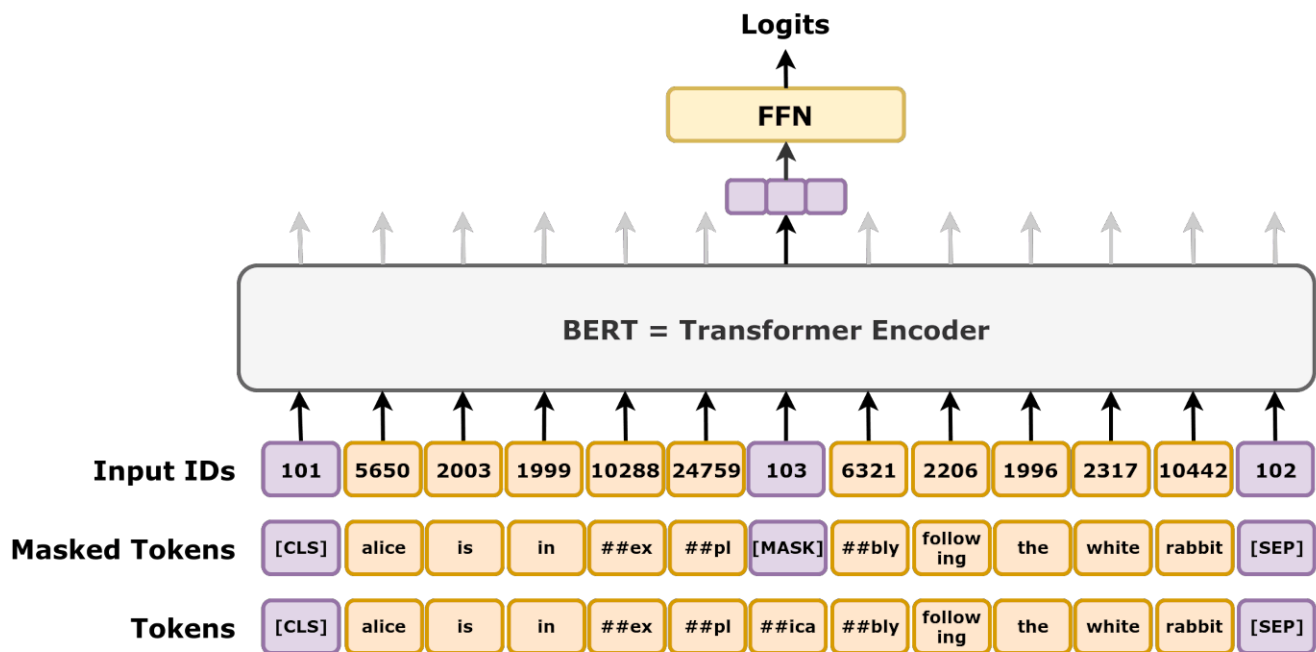


Image: [https://en.wikipedia.org/wiki/BERT\\_\(language\\_model\)](https://en.wikipedia.org/wiki/BERT_(language_model))

# Use cases for feedforward networks

- Word representations
- Text classification
- **Language modeling**

State of the art systems use more powerful neural architectures (we will learn transformers architectures on Monday), but simple models are useful to consider!

# Neural LMs

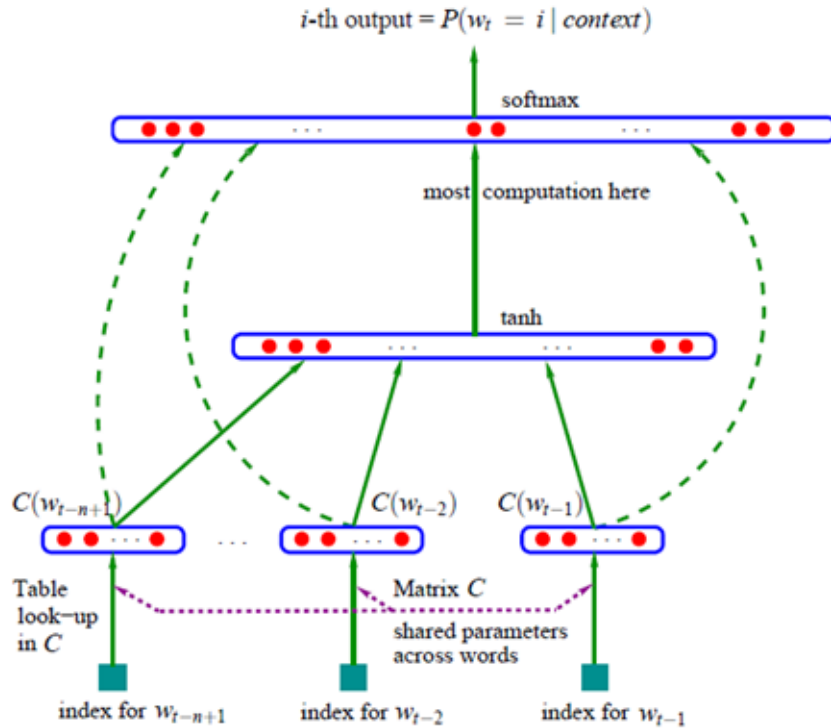


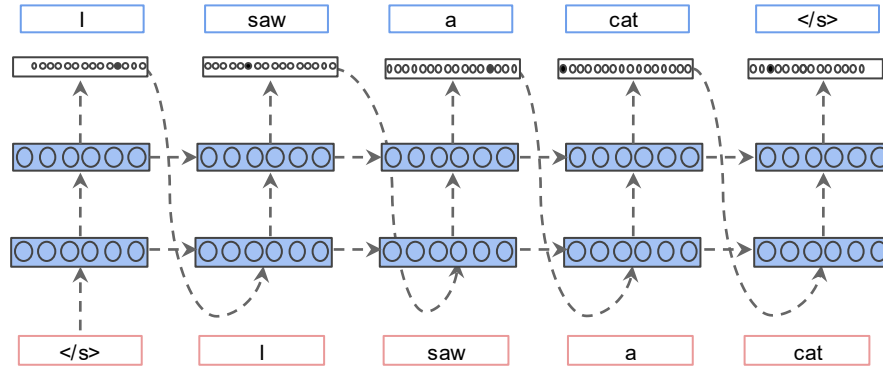
Image: (Bengio et al, 03)

# Neural LMs

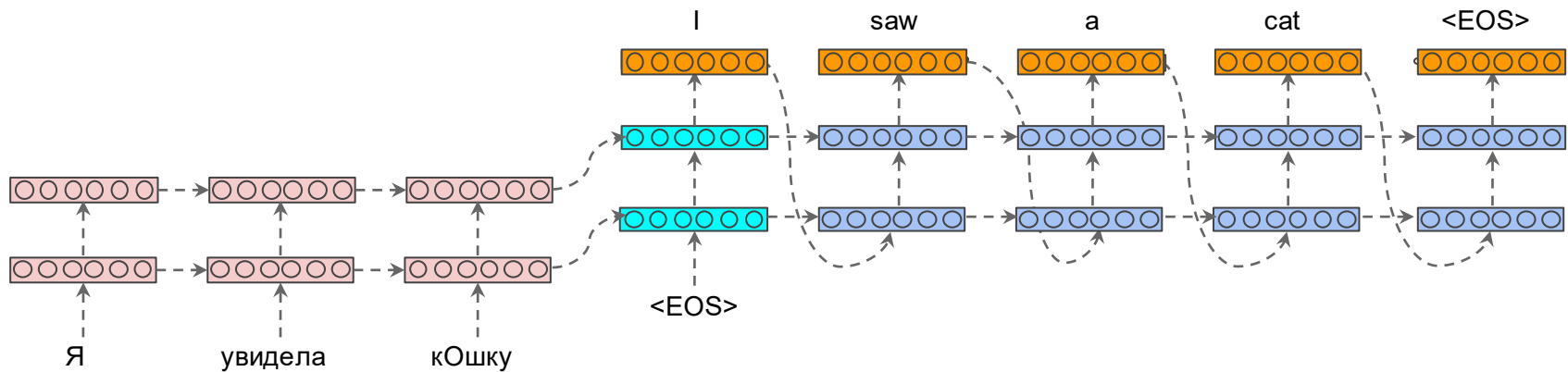
	n	c	h	m	direct	mix	train.	valid.	test.
MLP1	5		50	60	yes	no	182	284	268
MLP2	5		50	60	yes	yes		275	257
MLP3	5		0	60	yes	no	201	327	310
MLP4	5		0	60	yes	yes		286	272
MLP5	5		50	30	yes	no	209	296	279
MLP6	5		50	30	yes	yes		273	259
MLP7	3		50	30	yes	no	210	309	293
MLP8	3		50	30	yes	yes		284	270
MLP9	5		100	30	no	no	175	280	276
MLP10	5		100	30	no	yes		265	<b>252</b>
Kneser-Ney back-off	3							334	323
Kneser-Ney back-off	4							332	321
Kneser-Ney back-off	5							332	321

(Bengio et al, 03)

# Recurrent LMs

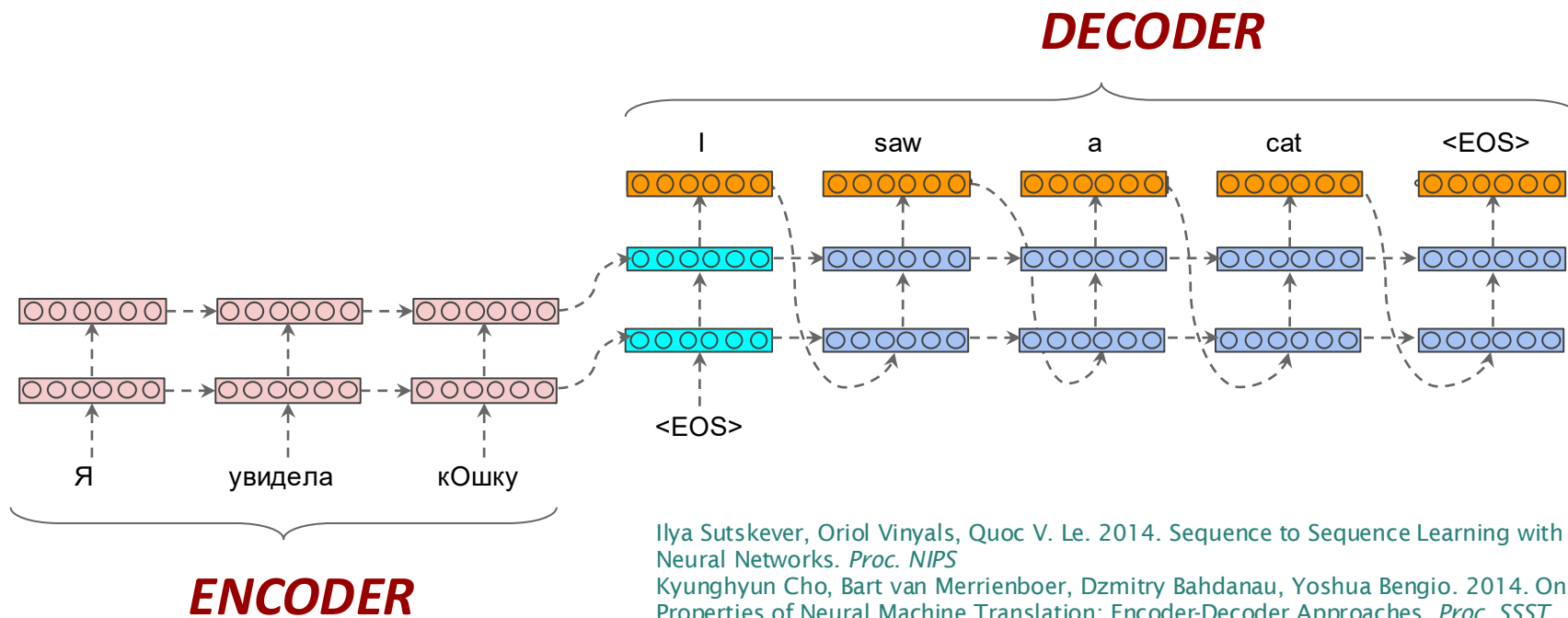


# Sequence-to-Sequence Models

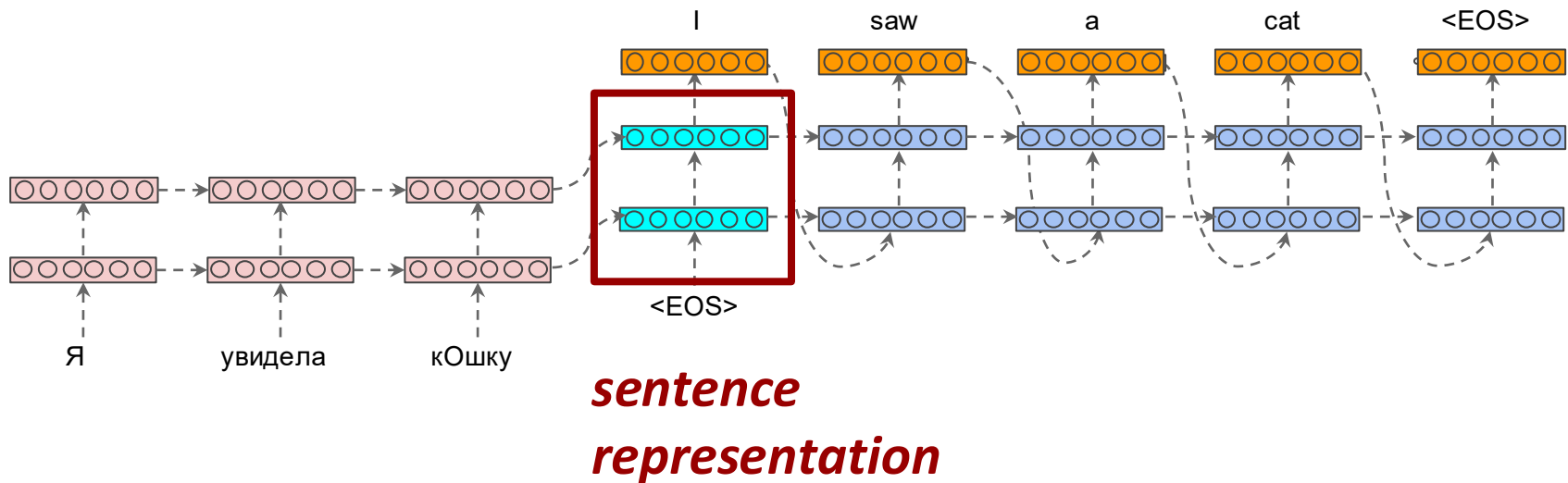


Ilya Sutskever, Oriol Vinyals, Quoc V. Le. 2014. Sequence to Sequence Learning with Neural Networks. Proc. NIPS

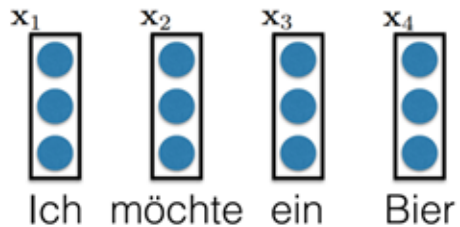
# Sequence-to-Sequence Models for Neural Machine Translation



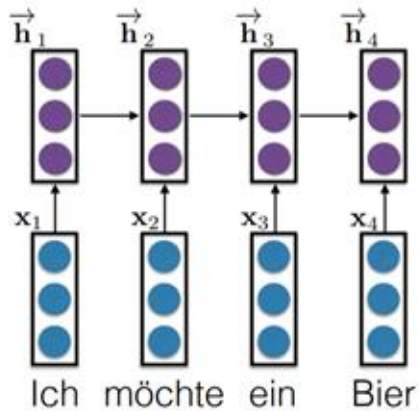
# Sequence-to-Sequence Models for NMT



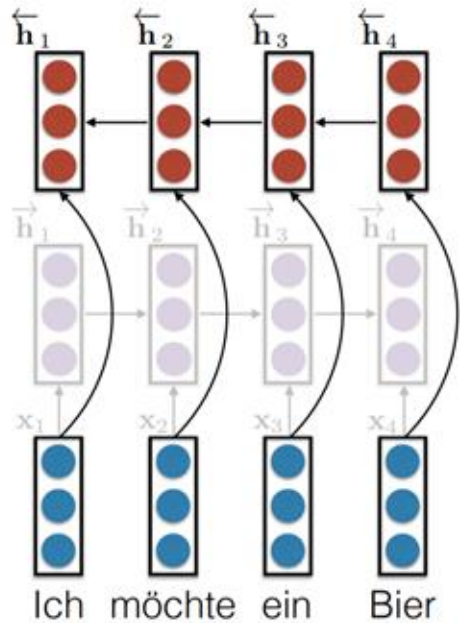
# Encoder: Bidirectional RNN



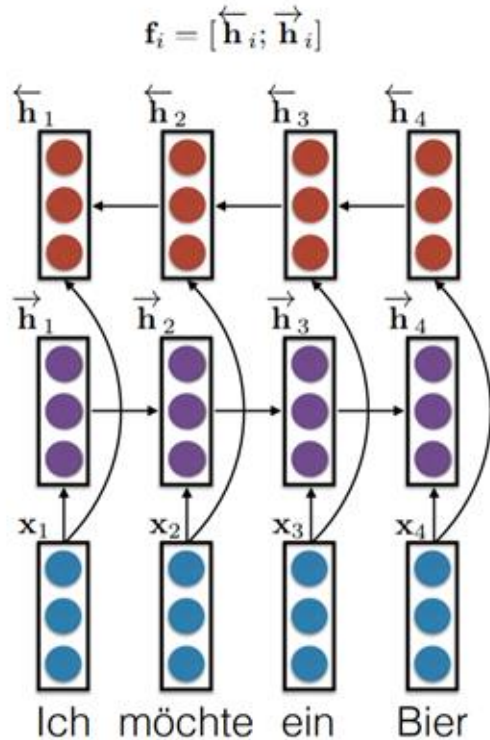
# Encoder: Bidirectional RNN



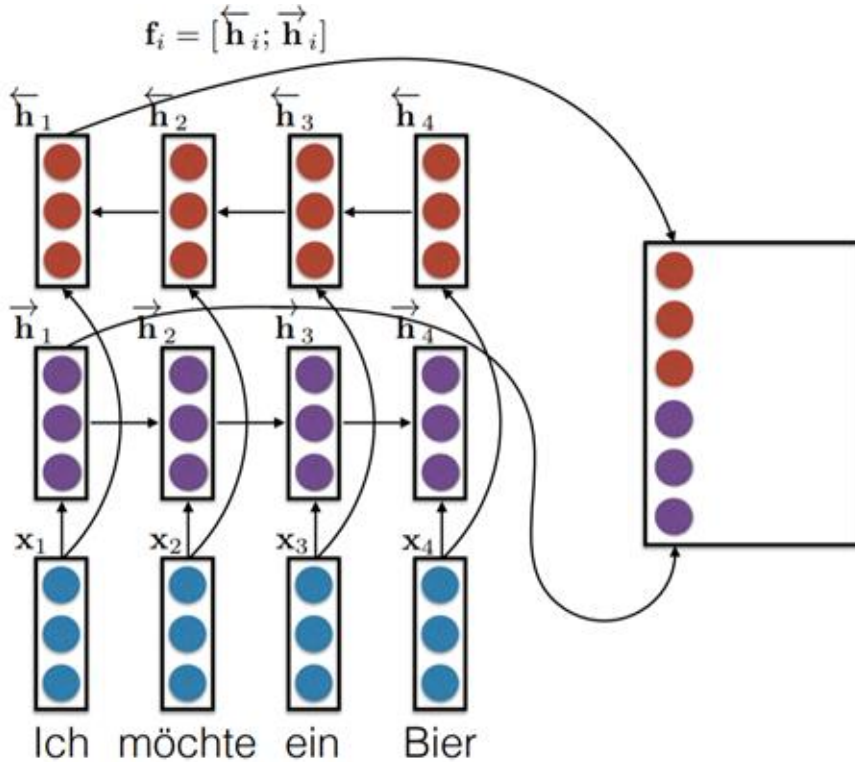
# Encoder: Bidirectional RNN



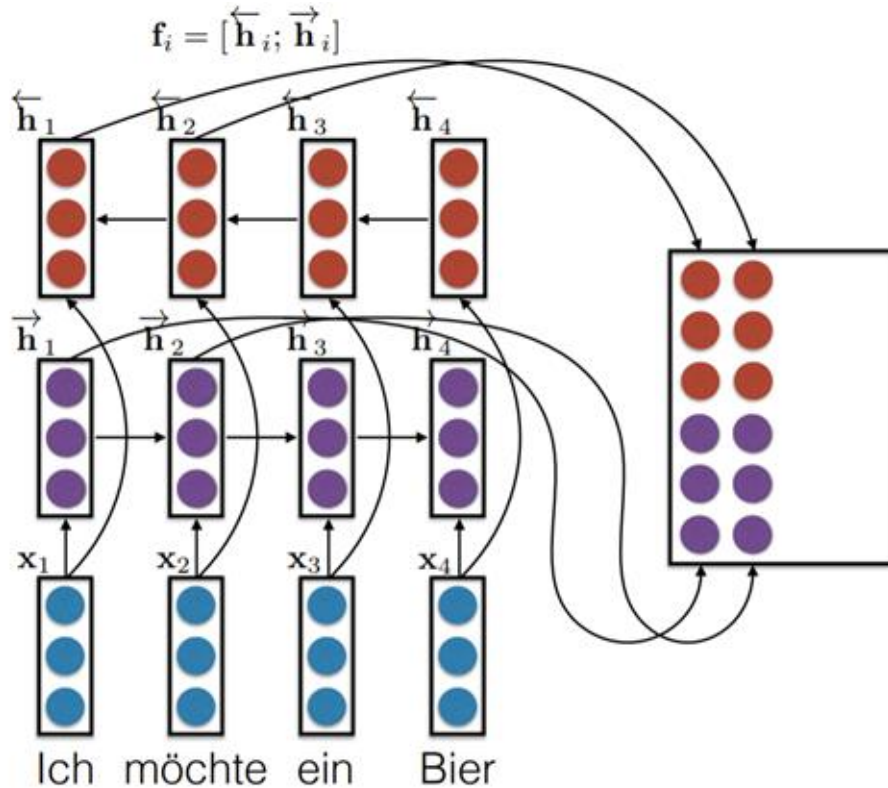
# Encoder: Bidirectional RNN



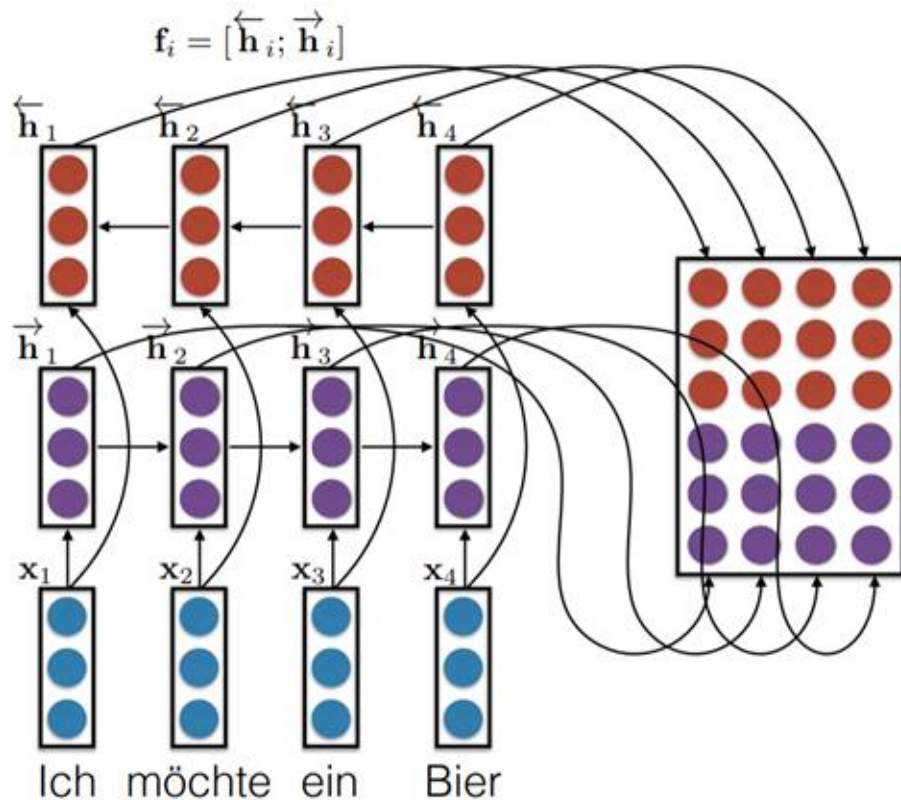
# Encoder: Bidirectional RNN



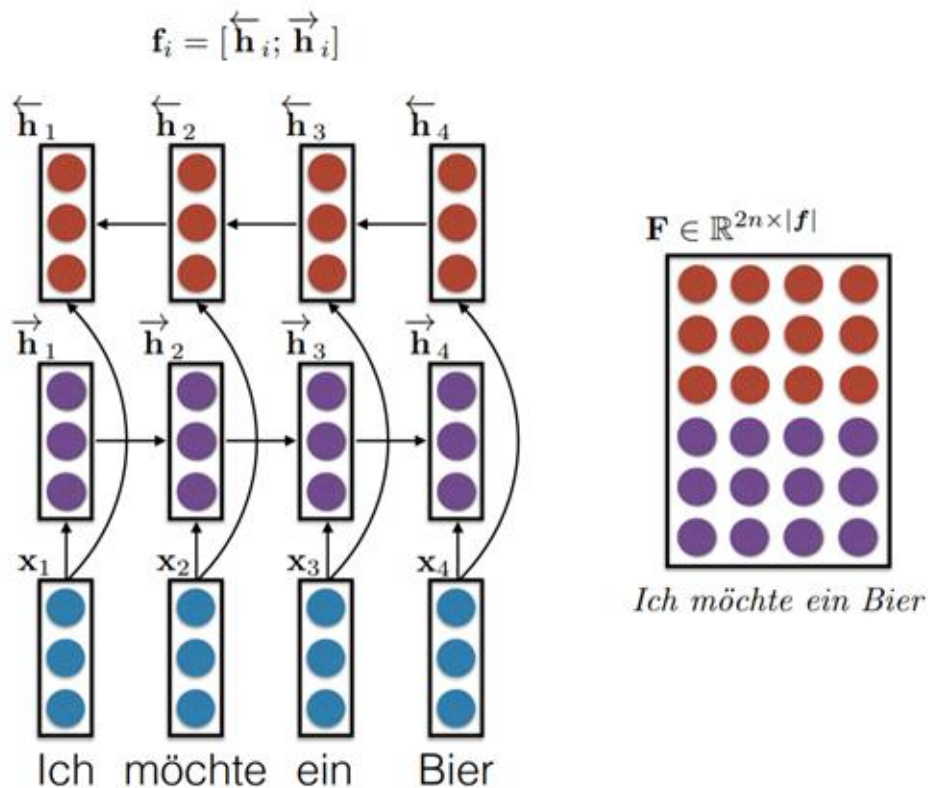
# Encoder: Bidirectional RNN



# Encoder: Bidirectional RNN

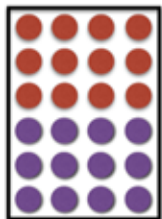
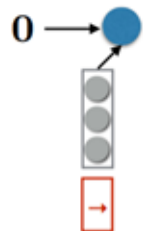


# Matrix Sentence Encoding

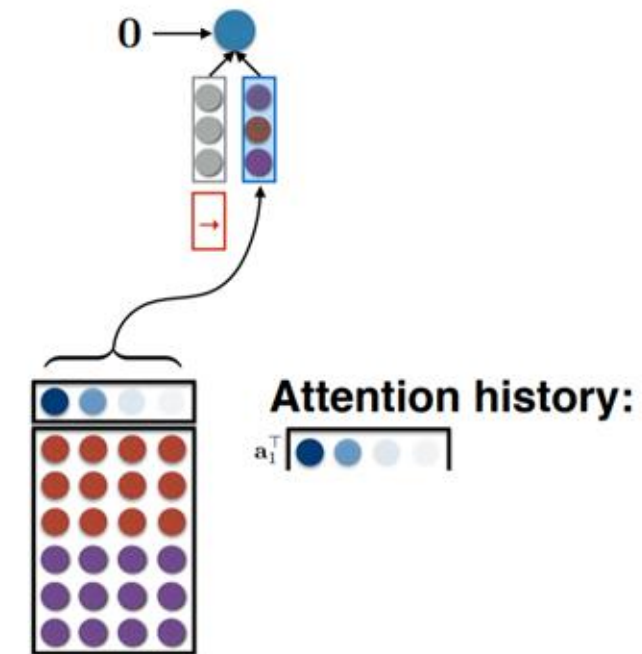


- matrix-encoded sentence

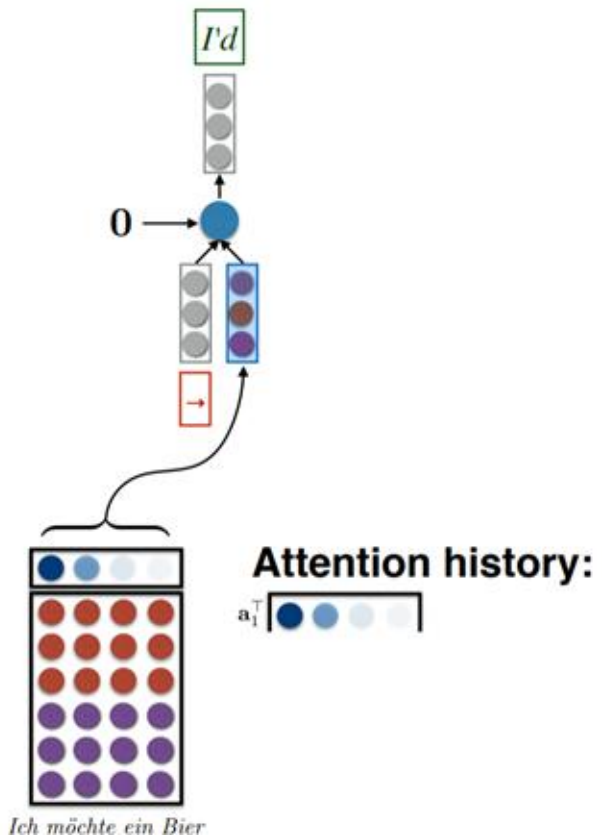
# Decoder: RNN + Attention

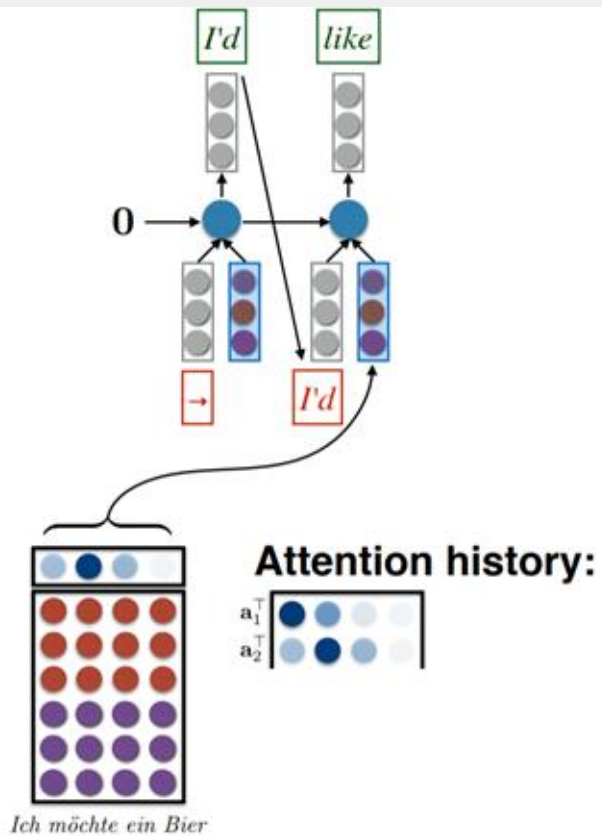


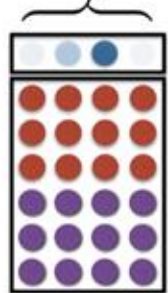
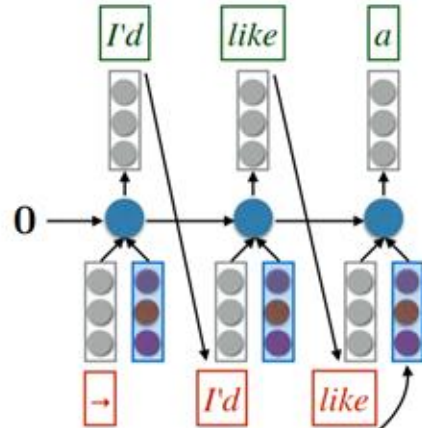
*Ich möchte ein Bier*



Ich möchte ein Bier



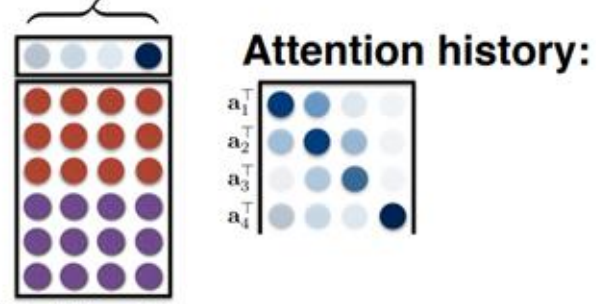
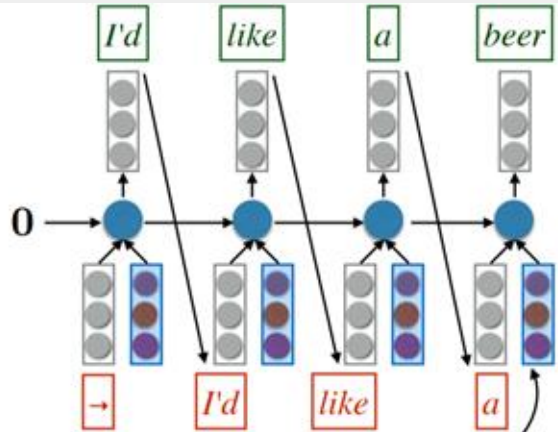




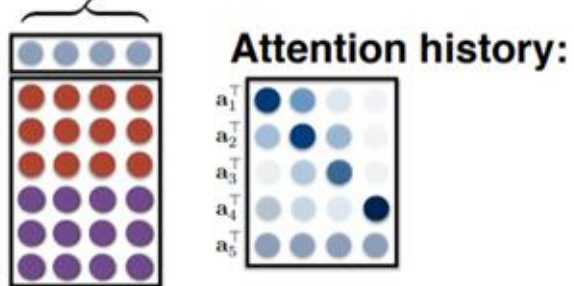
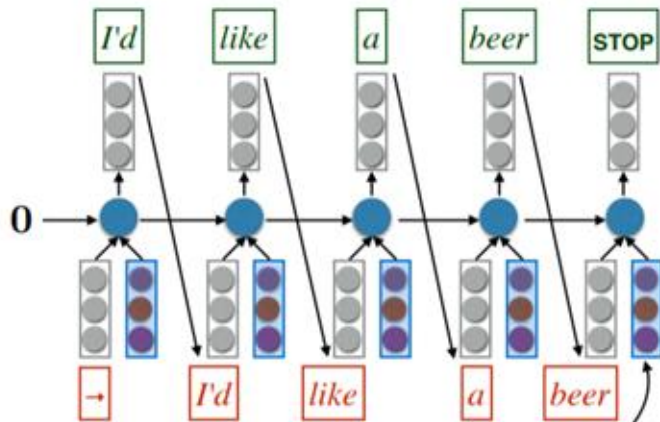
Ich möchte ein Bier

Attention history:





*Ich möchte ein Bier*



Ich möchte ein Bier