

# Natural Language Processing

## Neural Networks and Neural LMs

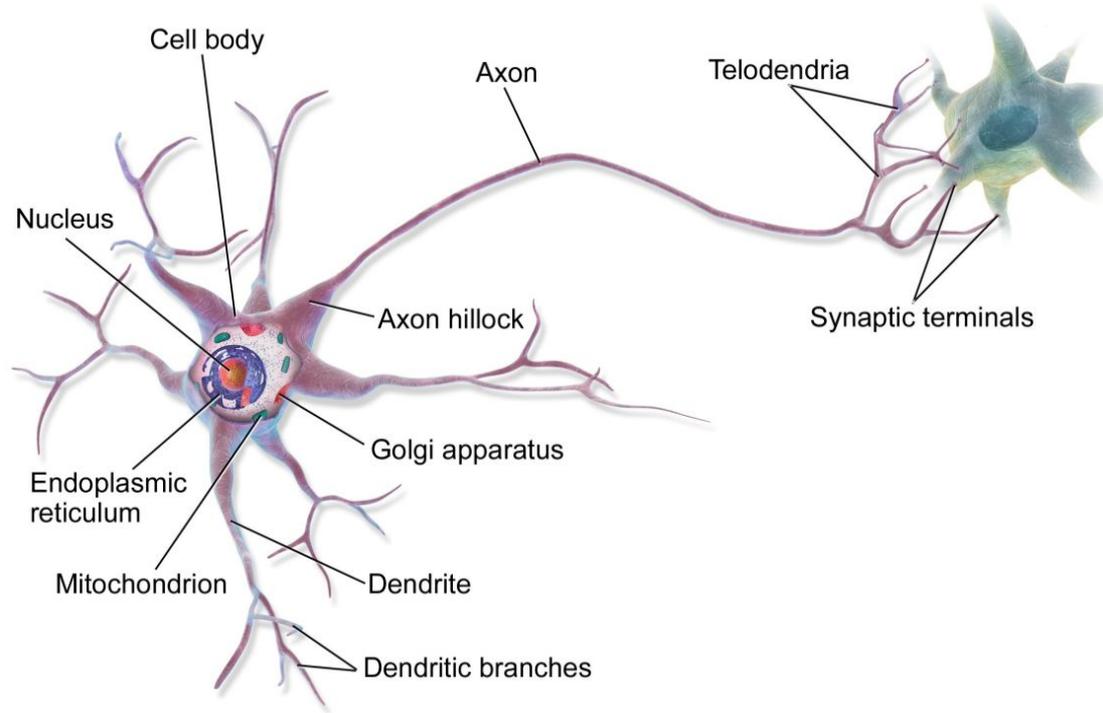
Yulia Tsvetkov

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

# Readings

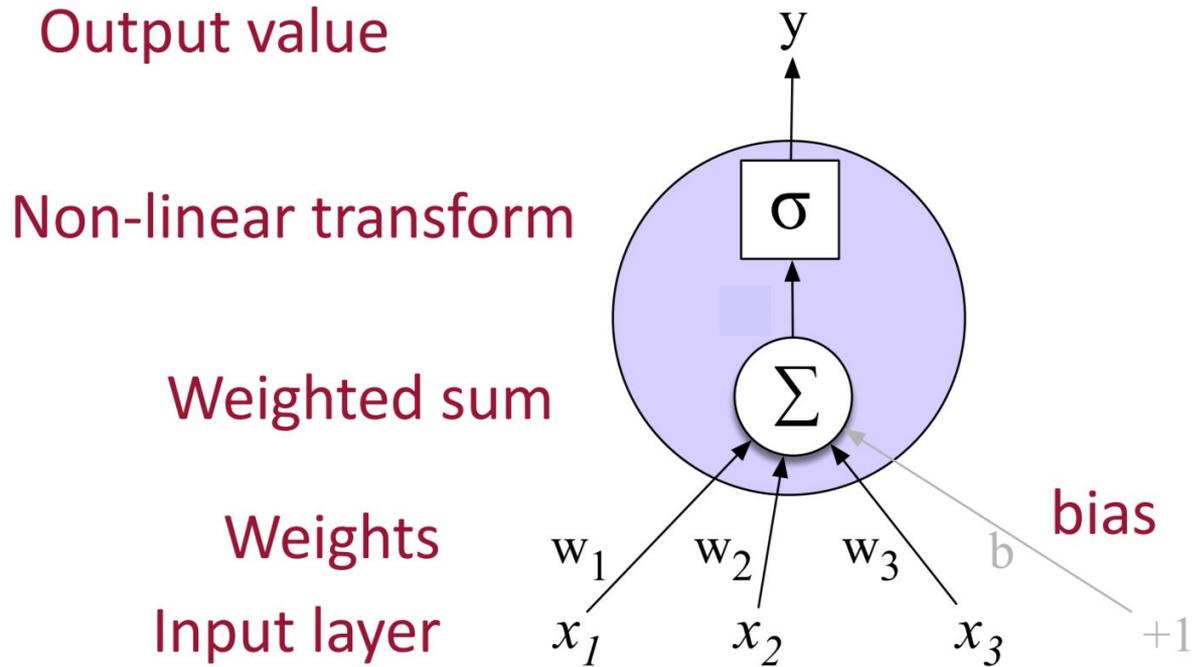
- Neutral networks chapters in J&M 6:
  - <https://web.stanford.edu/~jurafsky/slp3/6.pdf>
- Hundreds of blog posts and tutorials

# 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  $f$ :

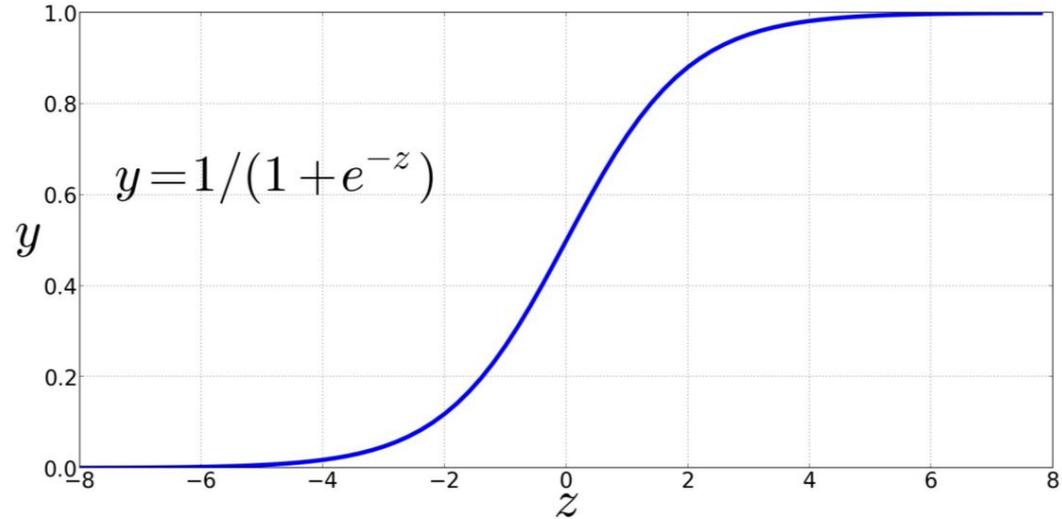
$$y = a = f(z)$$

# Non-Linear Activation Functions

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

**Sigmoid**

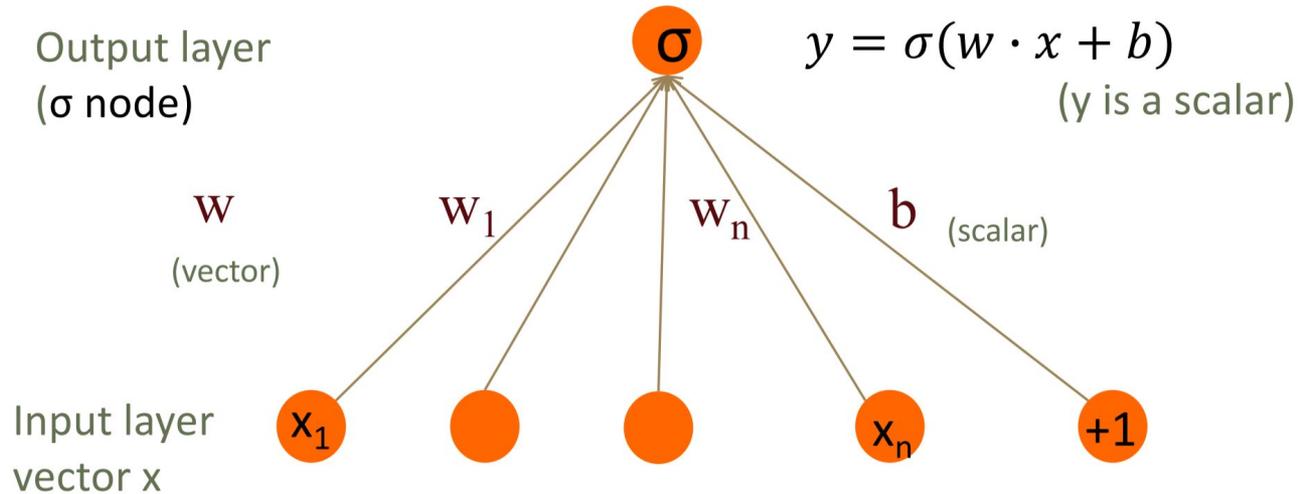
$$y = \sigma(z) = \frac{1}{1 + e^{-z}}$$



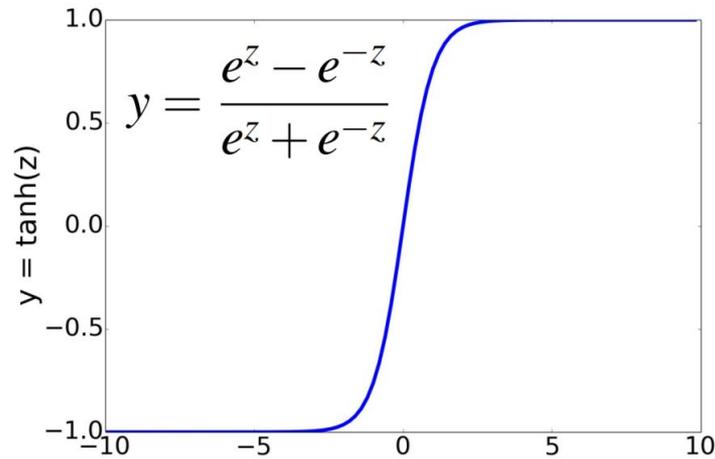
# Final function the unit is computing

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

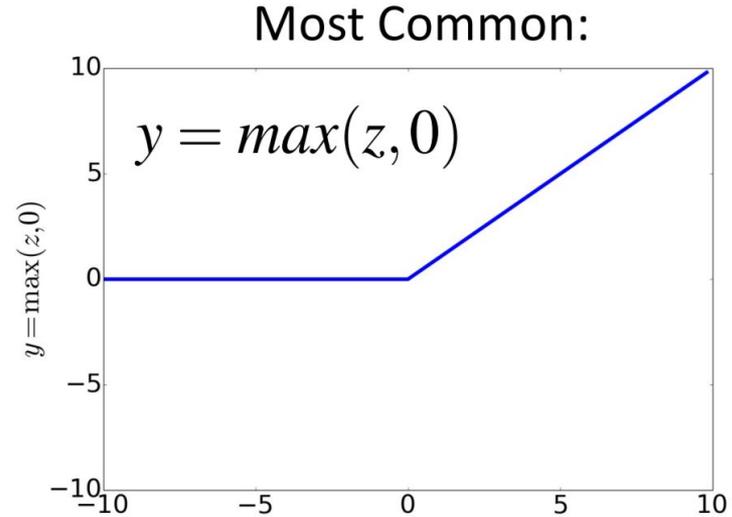
# Binary Logistic Regression as a 1-layer network



# Non-Linear Activation Functions besides sigmoid



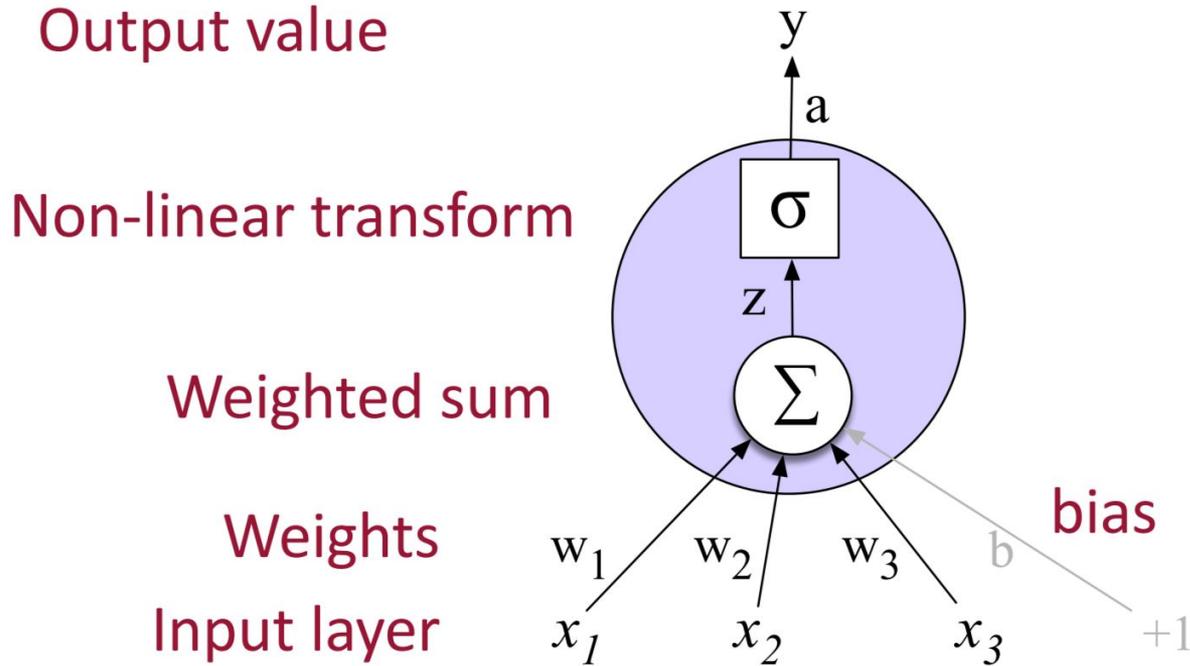
**tanh**



**ReLU**

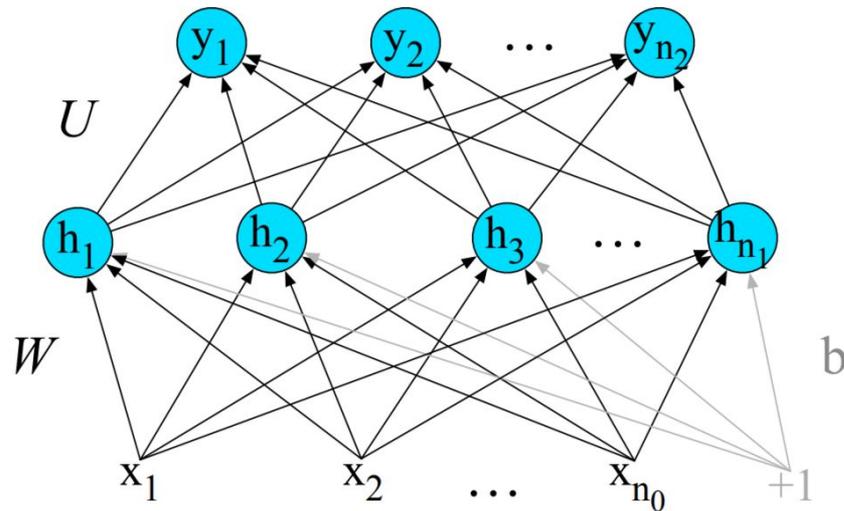
**Rectified Linear Unit**

# 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!)





# softmax: a generalization of sigmoid

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

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

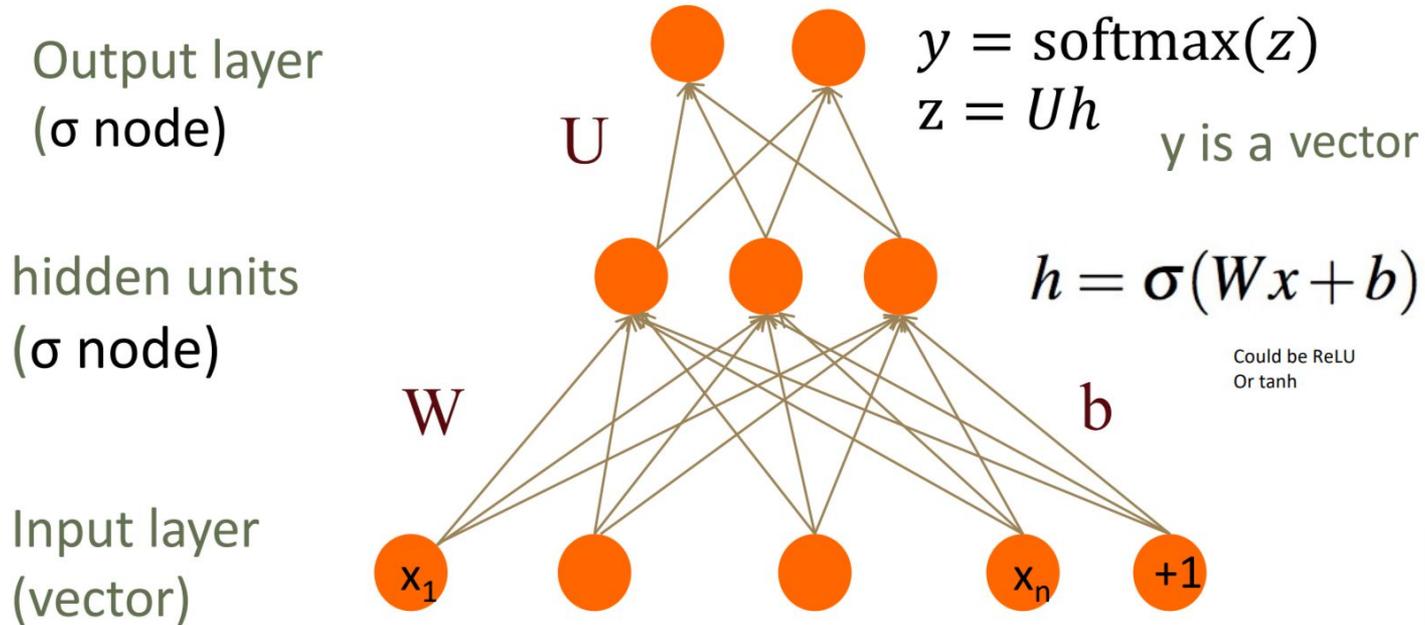
$$\text{softmax}(z_i) = \frac{\exp(z_i)}{\sum_{j=1}^k \exp(z_j)} \quad 1 \leq i \leq k$$

Example:

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

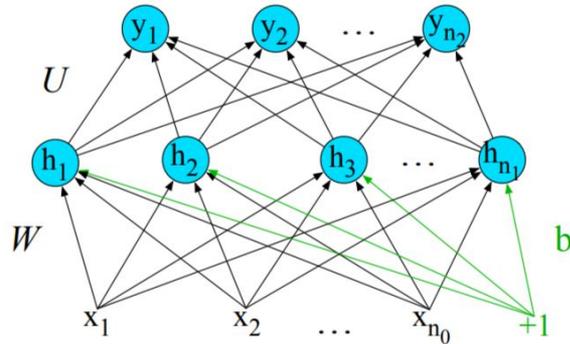
$$\text{softmax}(z) = [0.055, 0.090, 0.006, 0.099, 0.74, 0.010]$$

# Two-Layer Network with softmax output

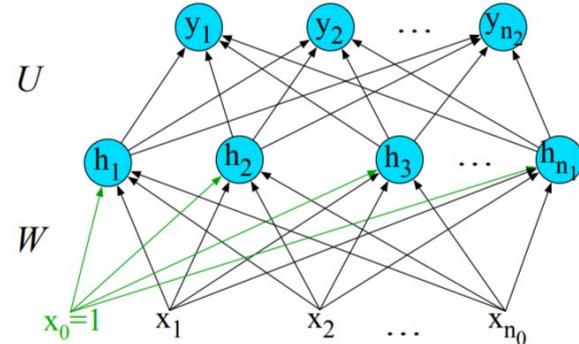


# Replacing the bias unit

Instead of:



We'll do this:



# Learning the weights

- Cross-entropy loss
- Backpropagation algorithm

---

## Algorithm 1 Backpropagation Algorithm

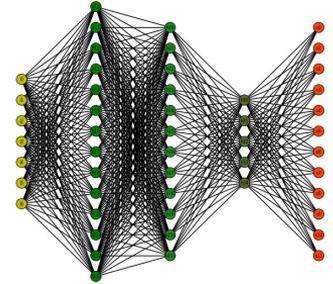
---

```

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!

# Common methods for getting short dense vectors

- In count-based models - Singular Value Decomposition (SVD)
  - A special case of this is called LSA – Latent Semantic Analysis
- “Neural Language Model”-inspired models
  - The weight matrix in the input layer is often used as “word embeddings”
  - Compute one embeddings for each word (word types)
- Alternative to these "static embeddings":
  - Contextual Embeddings (ELMo, BERT)
  - Compute distinct embeddings for a word in its context
  - Separate embeddings for each token of a word

Simple static embeddings you can download!

word2vec (Mikolov et al)

<https://code.google.com/archive/p/word2vec/>

GloVe (Pennington, Socher, Manning)

<http://nlp.stanford.edu/projects/glove/>

# Low-dimensional word representations

- Learning representations by back-propagating errors
  - Rumelhart, Hinton & Williams, 1986
- A neural probabilistic language model
  - Bengio et al., 2003
- Natural Language Processing (almost) from scratch
  - Collobert & Weston, 2008
- Word representations: A simple and general method for semi-supervised learning
  - Turian et al., 2010
- Distributed Representations of Words and Phrases and their Compositionality
  - Word2Vec; Mikolov et al., 2013

# word2vec

- Popular embedding method
- Very fast to train
- Code available on the web
- Idea: predict rather than count



Why GitHub? ▾ Enterprise Explore ▾ Marketplace Pricing ▾

Search

Sign in

Sign up

tmikolov / word2vec

Watch 55

★ Star 840

🔗 Fork 346

Code

Issues 38

Pull requests 4

Projects 0

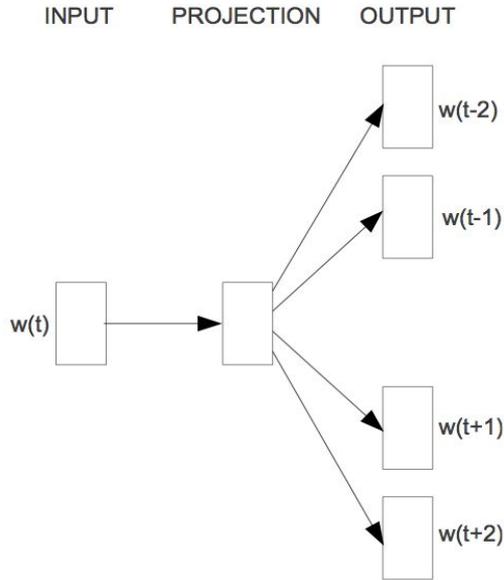
Security

Insights

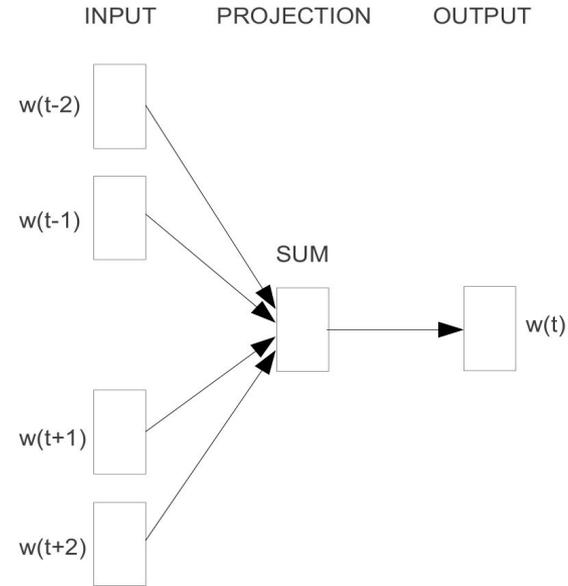
# word2vec

- 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)

# word2Vec



**Skip-gram**



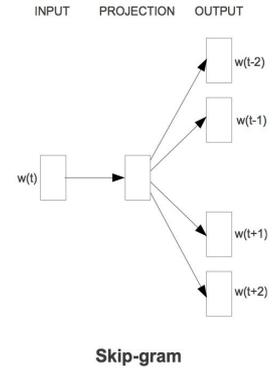
**CBOW**

- [Mikolov et al.' 13]

# Skip-gram Prediction

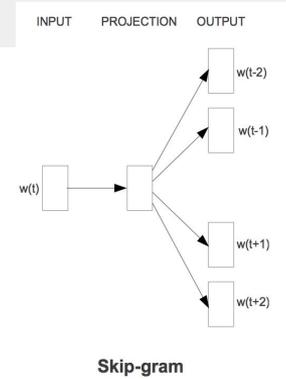
- Predict vs Count

the cat sat on the mat



# Skip-gram Prediction

- Predict vs Count



context size = 2

# Skip-gram Prediction

- Predict vs Count

the cat sat on the mat

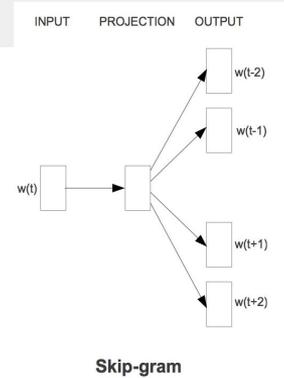
$w_t = \text{cat}$



**CLASSIFIER**



$w_{t-2} = \langle \text{start}_{-1} \rangle$   
 $w_{t-1} = \text{the}$   
 $w_{t+1} = \text{sat}$   
 $w_{t+2} = \text{on}$



context size = 2

# Skip-gram Prediction

- Predict vs Count

the cat sat on the mat

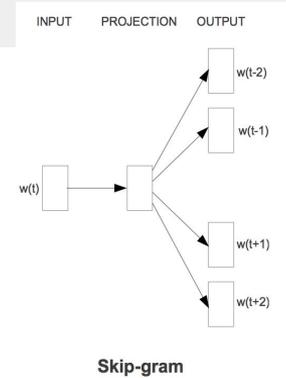
$w_t = \text{sat}$



**CLASSIFIER**



$w_{t-2} = \text{the}$   
 $w_{t-1} = \text{cat}$   
 $w_{t+1} = \text{on}$   
 $w_{t+2} = \text{the}$



context size = 2

# Skip-gram Prediction

- Predict vs Count

the cat sat on the mat

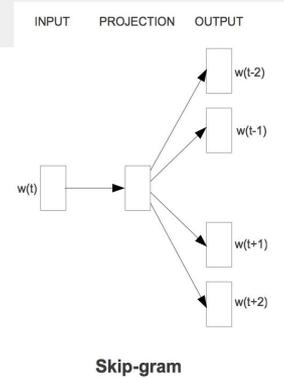
$w_t = \text{on}$



**CLASSIFIER**



$w_{t-2} = \text{cat}$   
 $w_{t-1} = \text{sat}$   
 $w_{t+1} = \text{the}$   
 $w_{t+2} = \text{mat}$



context size = 2

# Skip-gram Prediction

- Predict vs Count

the cat sat on the mat

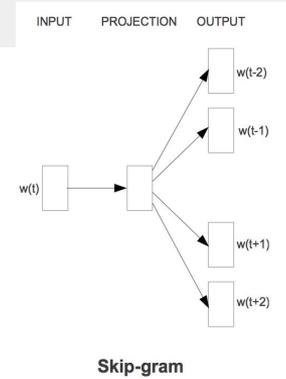
$w_t = \text{the}$



**CLASSIFIER**



$w_{t-2} = \text{sat}$   
 $w_{t-1} = \text{on}$   
 $w_{t+1} = \text{mat}$   
 $w_{t+2} = \langle \text{end}_{+1} \rangle$



context size = 2

# Skip-gram Prediction

- Predict vs Count

the cat sat on the mat

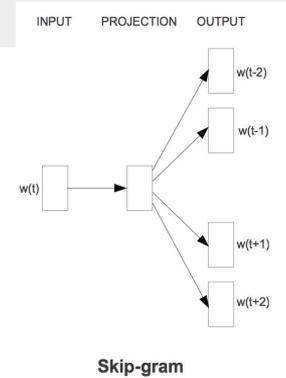
$w_t = \text{mat}$



**CLASSIFIER**



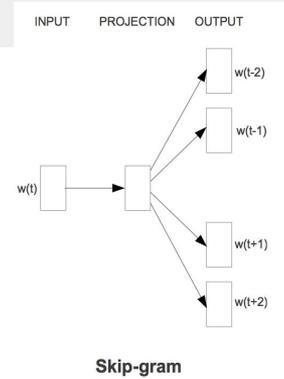
$w_{t-2} = \text{on}$   
 $w_{t-1} = \text{the}$   
 $w_{t+1} = \langle \text{end}_{+1} \rangle$   
 $w_{t+2} = \langle \text{end}_{+2} \rangle$



context size = 2

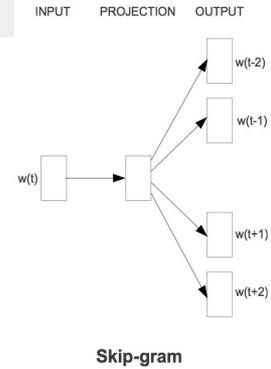
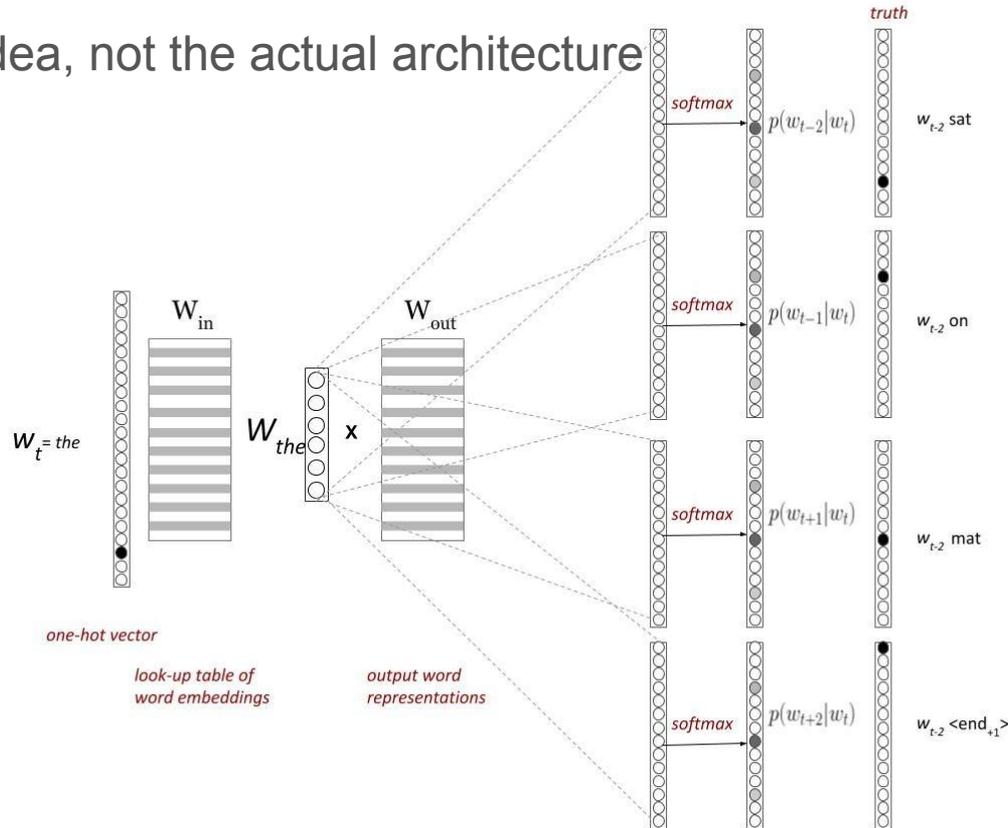
# Skip-gram Prediction

- Predict vs Count

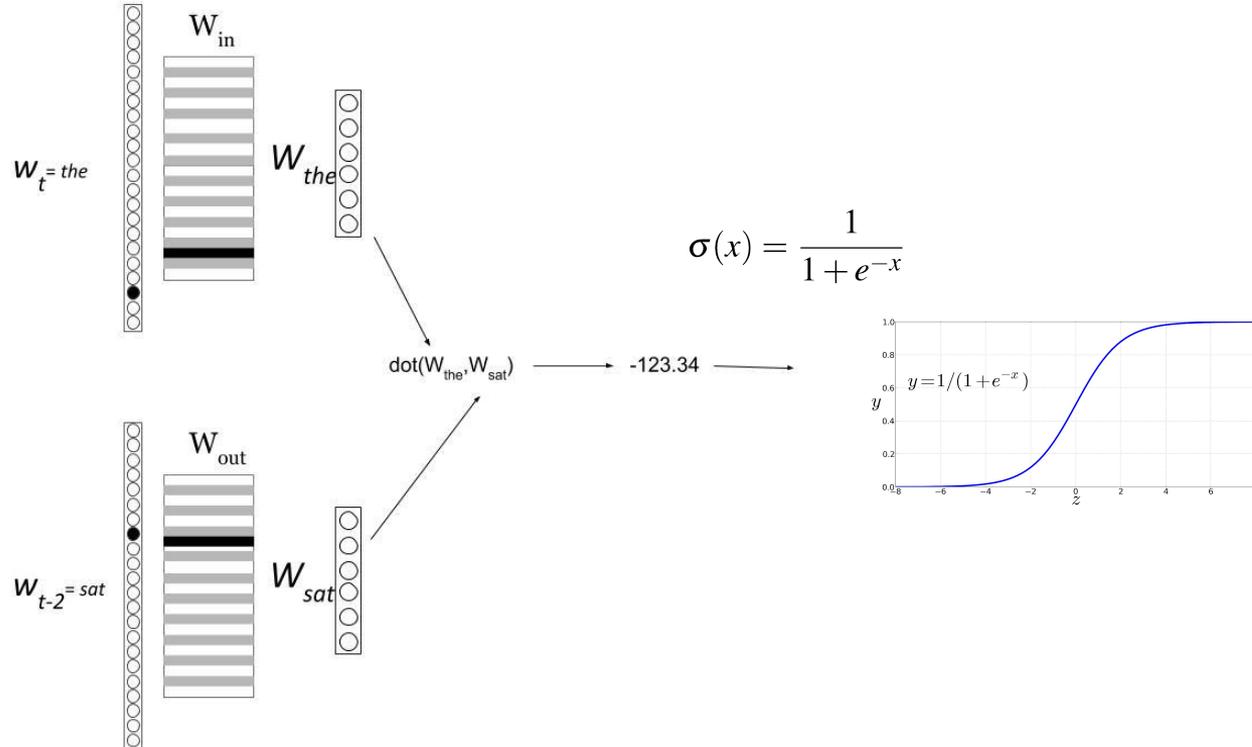


# Skip-gram Prediction

- Conceptual idea, not the actual architecture



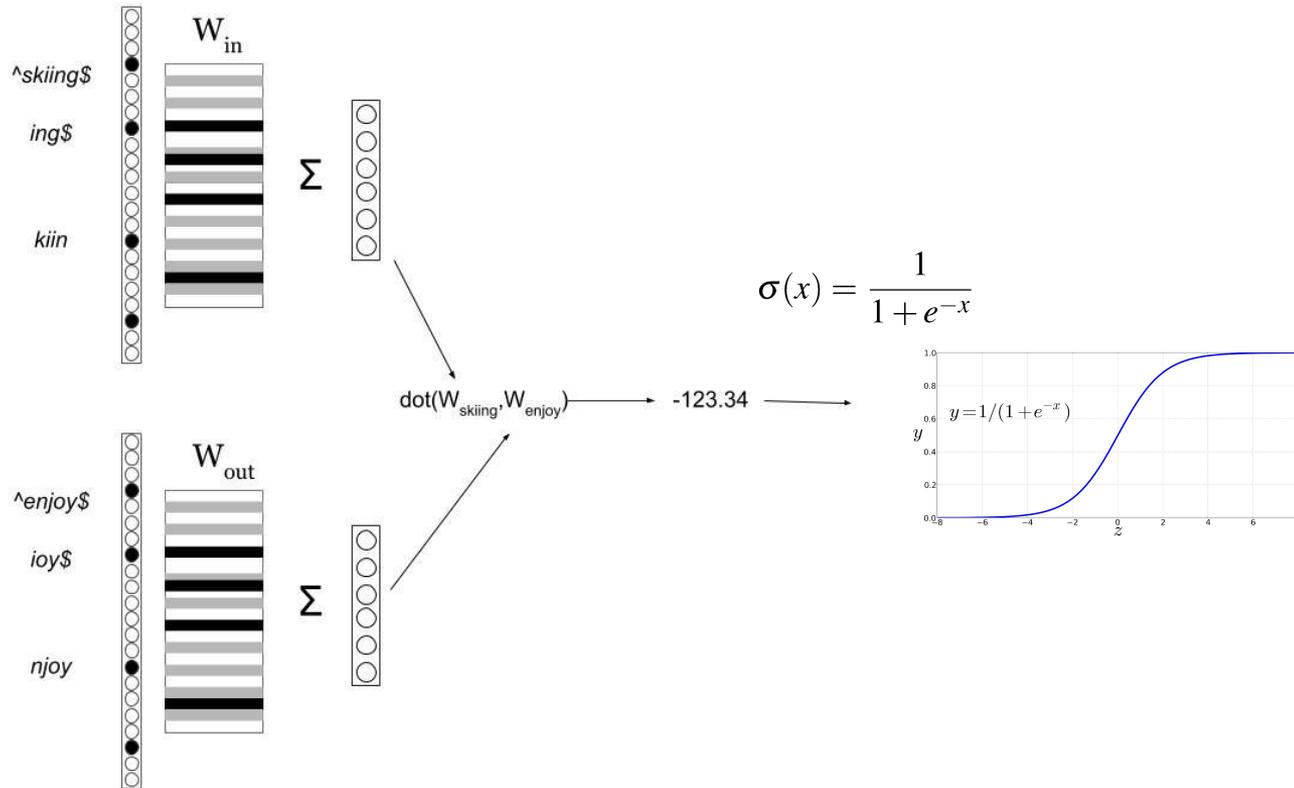
# How to compute $p(+ | t, c)$ ?



# Approach: predict if candidate word $c$ is a “neighbor”

1. Treat the target word  $t$  and a neighboring context word  $c$  as **positive examples**
2. Randomly sample other words in the lexicon to get negative examples
3. Use logistic regression to train a classifier to distinguish those two cases
4. Use the learned weights as the embeddings

# FastText



# SGNS

Given a tuple  $(t,c)$  = target, context

- $(\text{cat}, \text{sat})$
- $(\text{cat}, \text{aardvark})$

Return probability that  $c$  is a real context word:

$$P(+|t,c) = \frac{1}{1 + e^{-t \cdot c}}$$

$$\begin{aligned} P(-|t,c) &= 1 - P(+|t,c) \\ &= \frac{e^{-t \cdot c}}{1 + e^{-t \cdot c}} \end{aligned}$$

# Learning the classifier

- Iterative process
  - We'll start with 0 or random weights
  - Then adjust the word weights to
    - make the positive pairs more likely
    - and the negative pairs less likely
  - over the entire training set:

$$\sum_{(t,c) \in +} \log P(+|t,c) + \sum_{(t,c) \in -} \log P(-|t,c)$$

- Train using gradient descent

# 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>

# 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

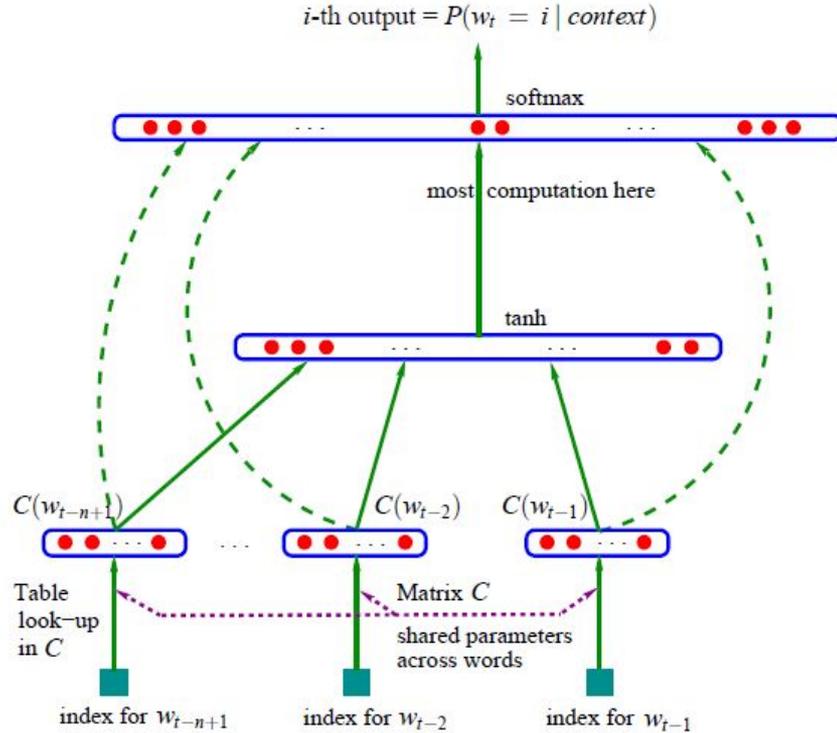


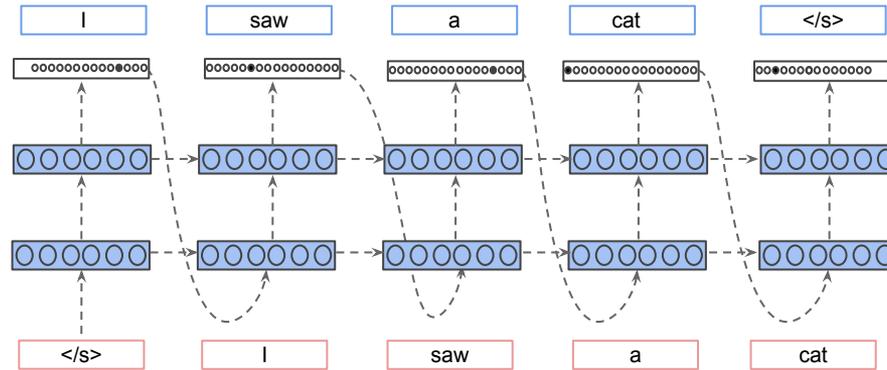
Image: (Bengio et al, 03)

# Neural LMs

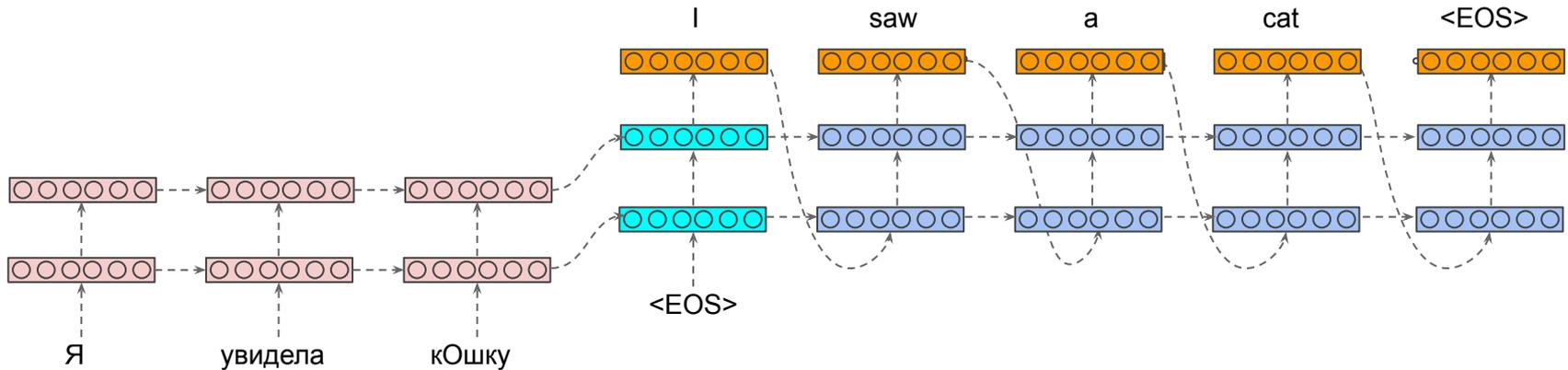
	n	c	h	m	direct	mix	train.	valid.	test.
MLP1	5		50	60	yes	no	182	284	268
MLP2	5		50	60	yes	yes		275	257
MLP3	5		0	60	yes	no	201	327	310
MLP4	5		0	60	yes	yes		286	272
MLP5	5		50	30	yes	no	209	296	279
MLP6	5		50	30	yes	yes		273	259
MLP7	3		50	30	yes	no	210	309	293
MLP8	3		50	30	yes	yes		284	270
MLP9	5		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

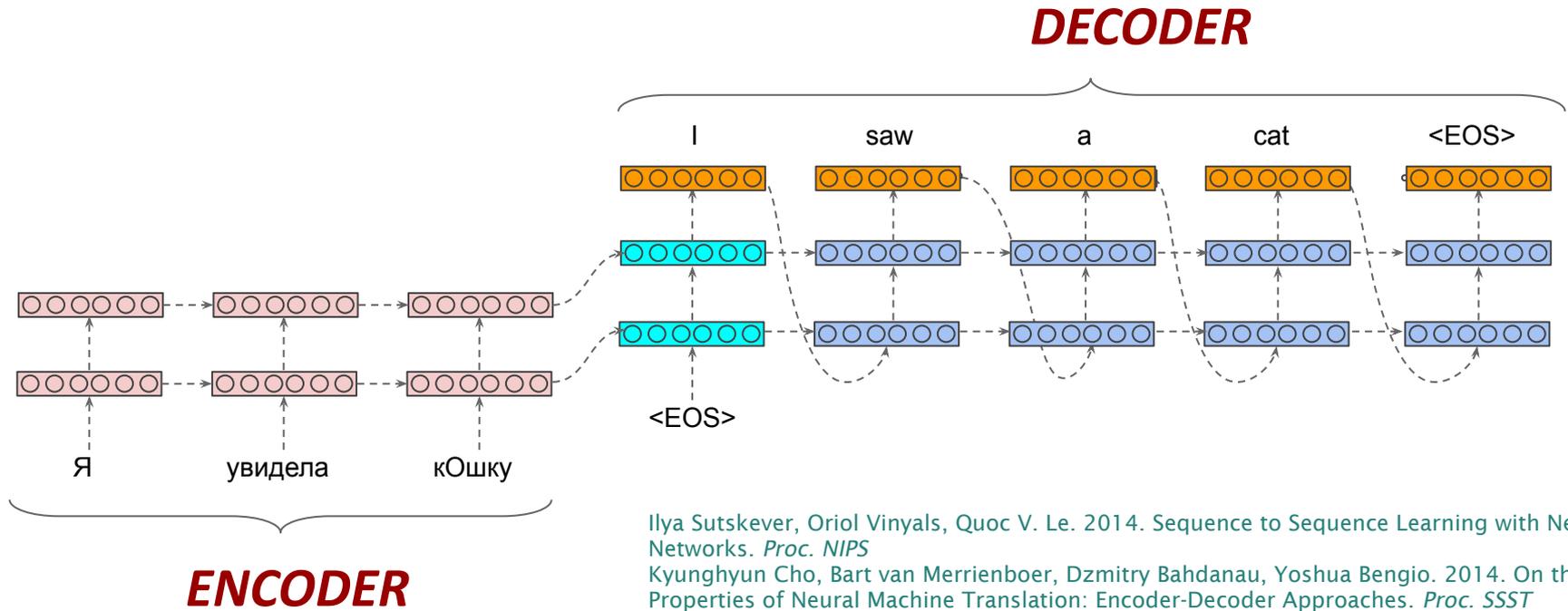


# Sequence-to-Sequence Models



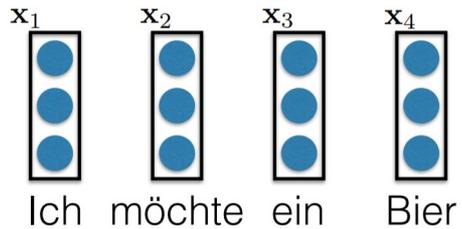
Ilya Sutskever, Oriol Vinyals, Quoc V. Le. 2014. Sequence to Sequence Learning with Neural Networks. Proc. NIPS

# Sequence-to-Sequence Models for Neural Machine Translation

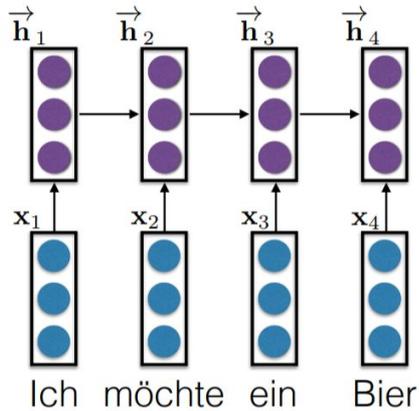




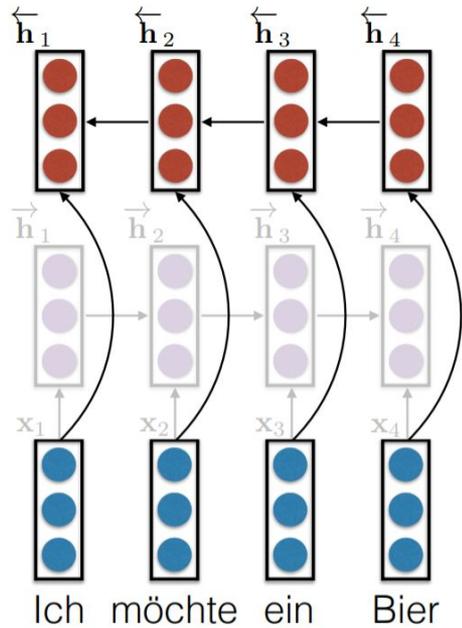
# Encoder: Bidirectional RNN



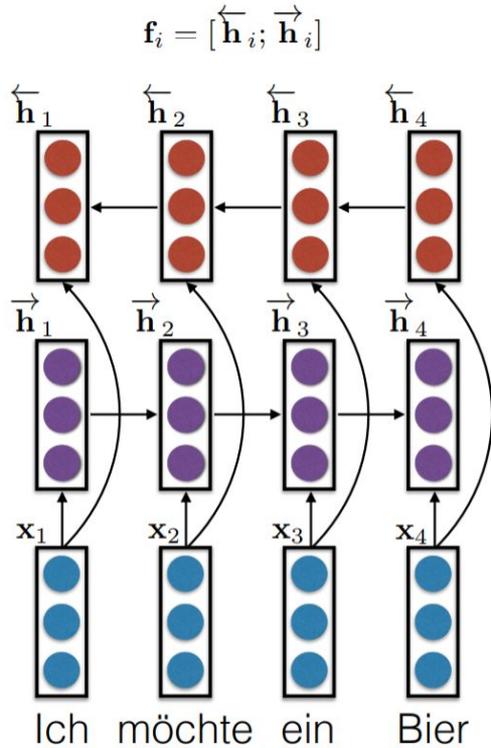
# Encoder: Bidirectional RNN



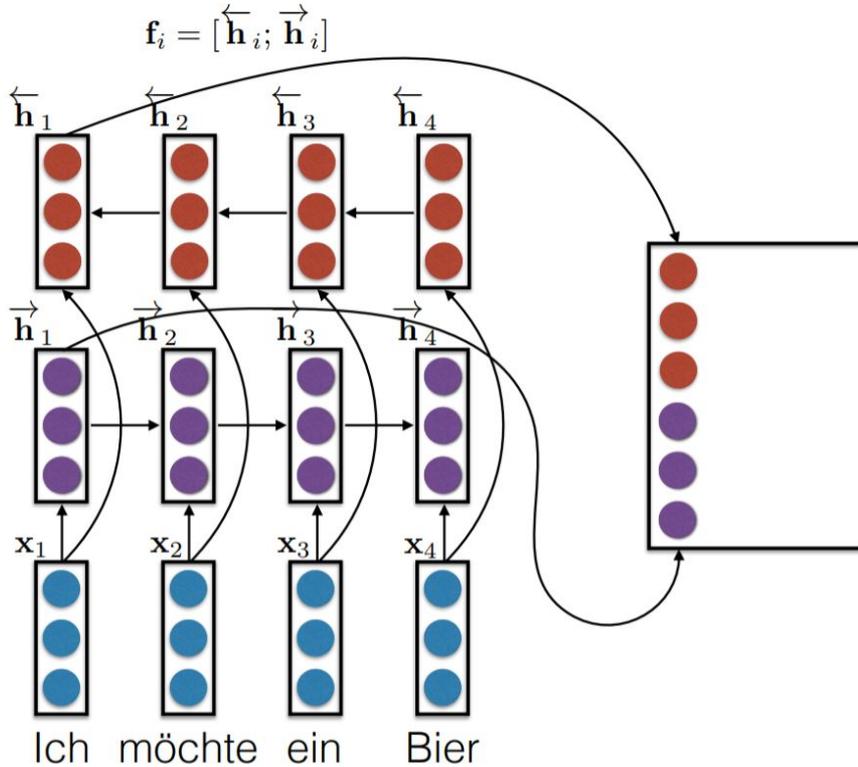
# Encoder: Bidirectional RNN



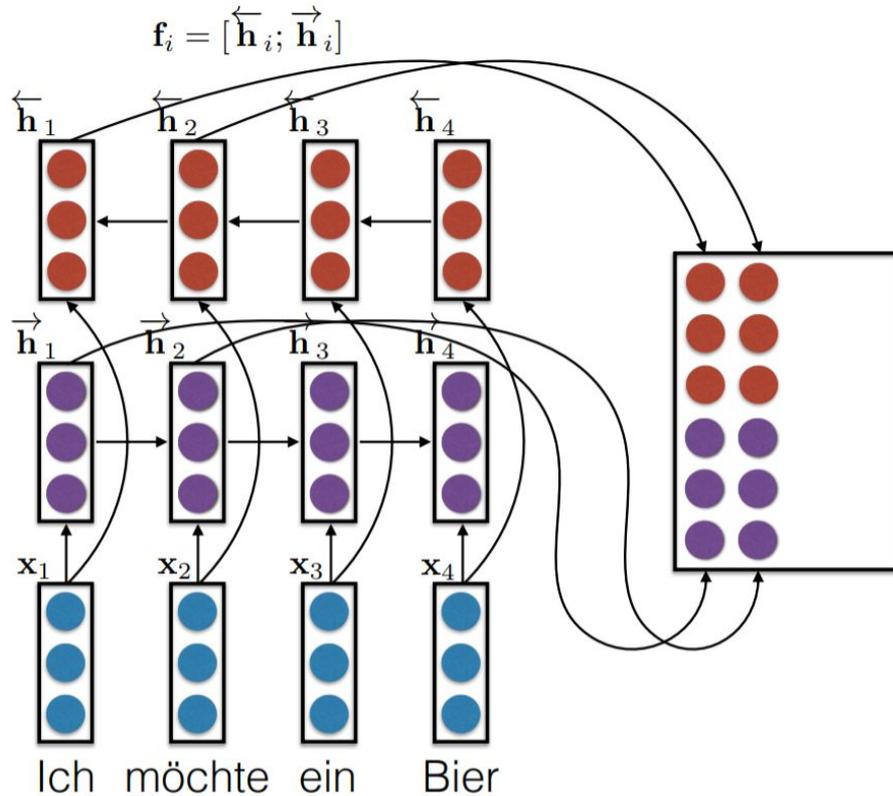
# Encoder: Bidirectional RNN



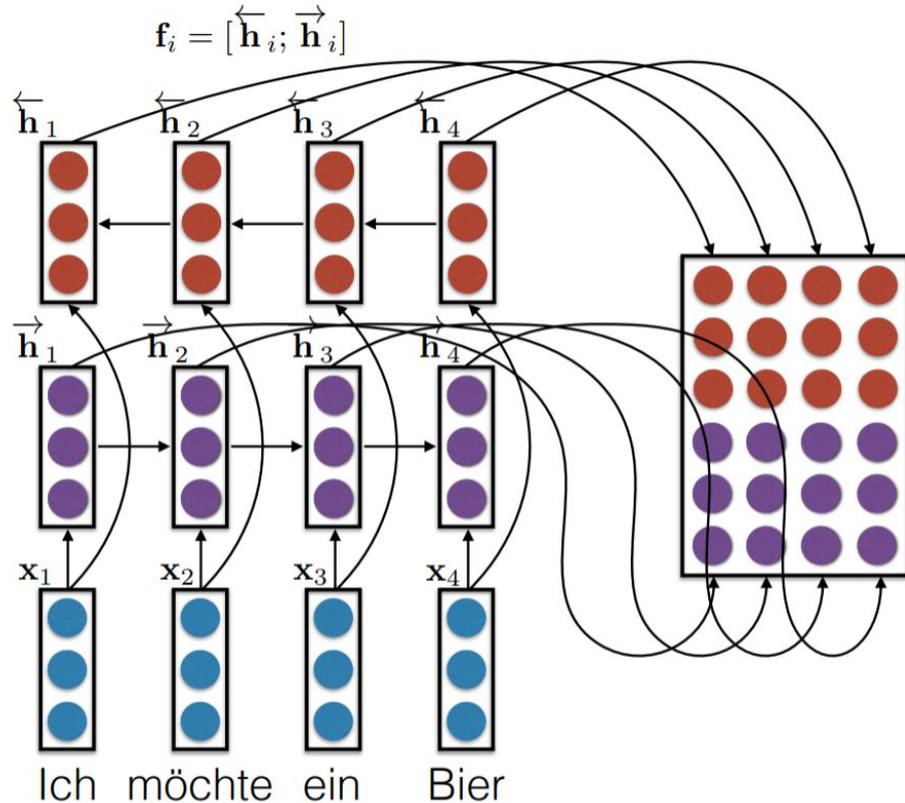
# Encoder: Bidirectional RNN



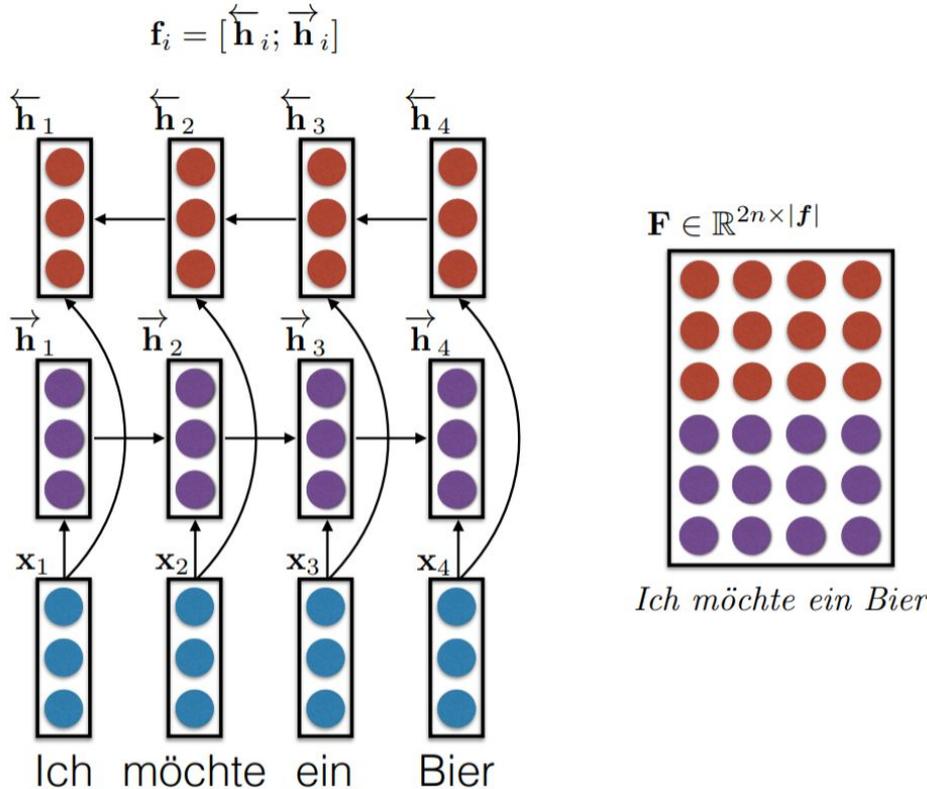
# Encoder: Bidirectional RNN



# Encoder: Bidirectional RNN

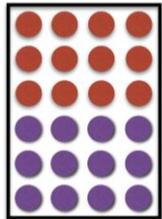
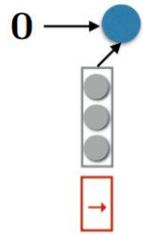


# Matrix Sentence Encoding

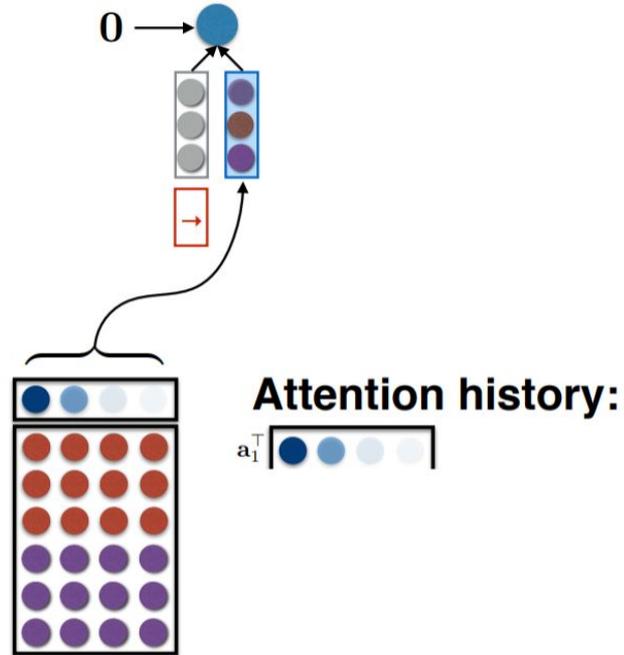


- matrix-encoded sentence

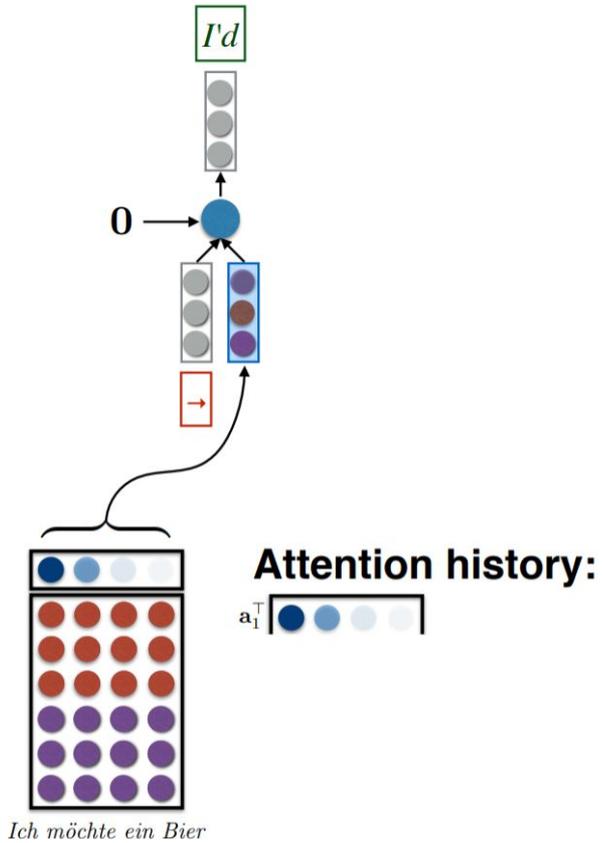
# Decoder: RNN + Attention

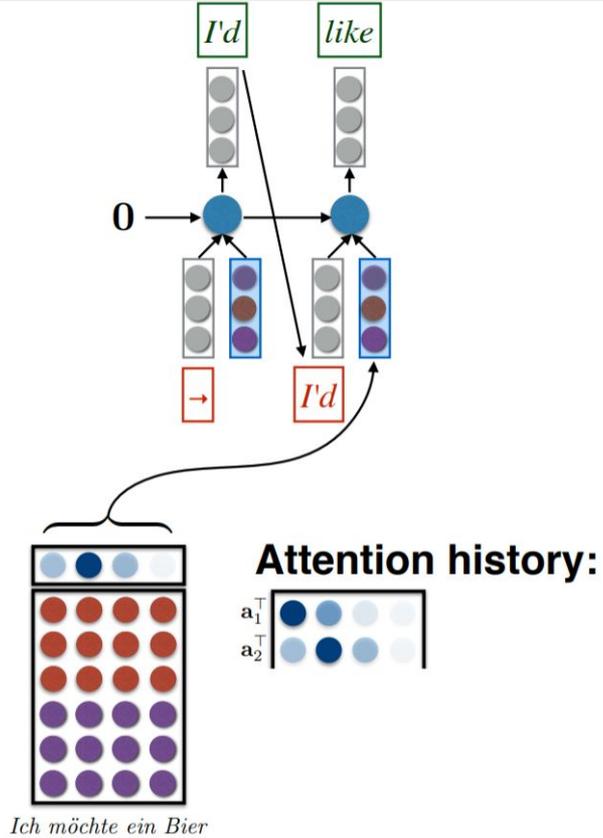


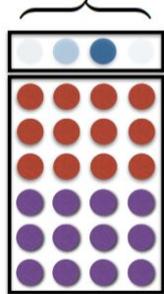
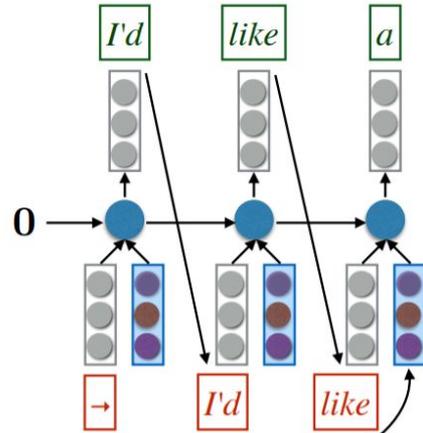
*Ich möchte ein Bier*



Ich möchte ein Bier



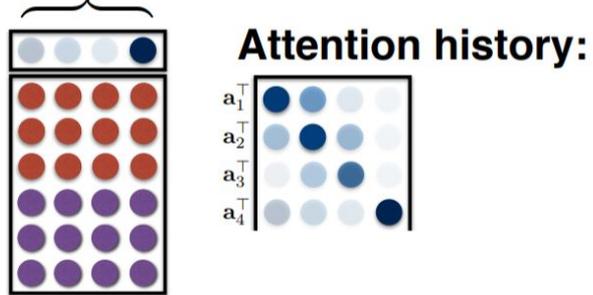
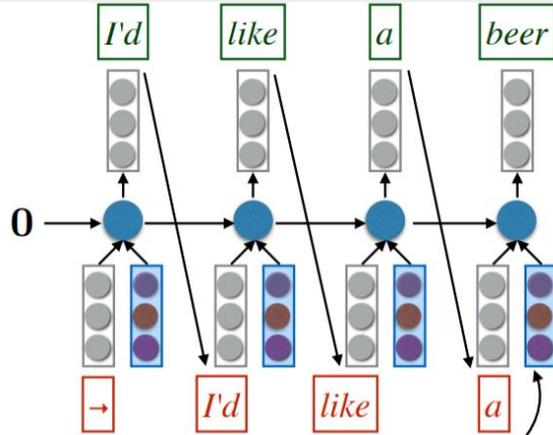




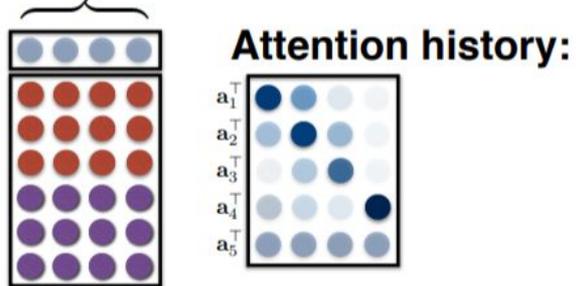
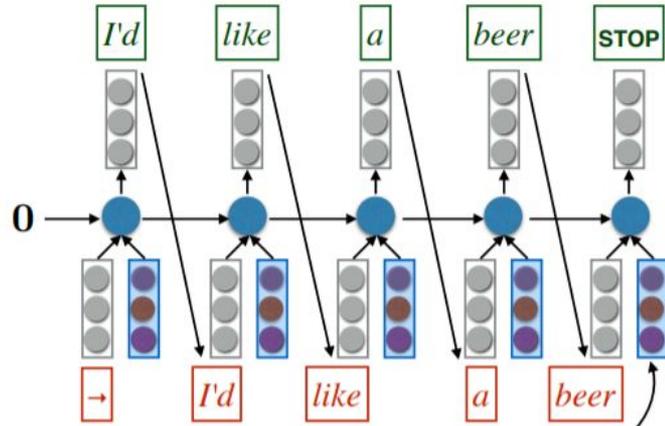
**Attention history:**



*Ich möchte ein Bier*



*Ich möchte ein Bier*



*Ich möchte ein Bier*