

Neural models of language

CSC401/2511 – Natural Language Computing – Spring 2022
Lecture 5

CSC401/2511 – Spring 2022

University of Toronto

Logistics

- Assignment 1: due Feb 11, 2022
- Assignment 2: release Feb 12, 2022
- Lecture delivery:
 - Online (*as is*) until Feb 18
 - Reading week break: Feb 21-25 (*no lectures*)
 - In-person Feb 28th onwards
- Final exam: planned in-person

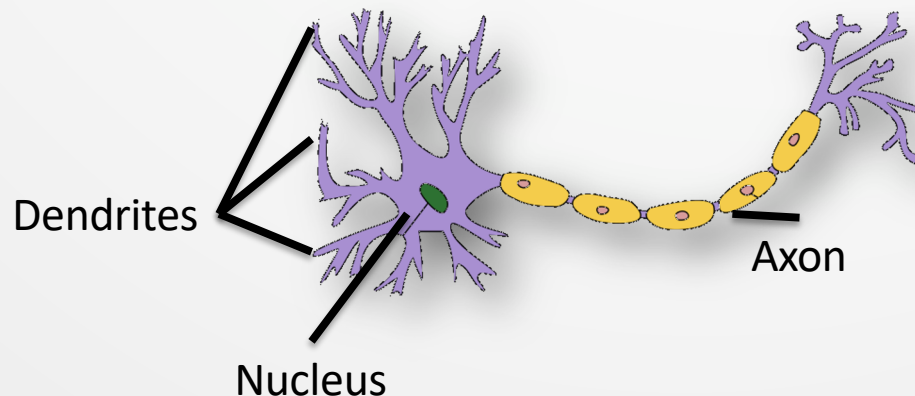
Neural networks

- Introduction
- Word-level representations
- Neural language models
- Recurrent neural networks
- Sequence-to-sequence modelling
- Some recent developments

With material from Phil Blunsom, Piotr Mirowski, Adam Kalai, and James Zou

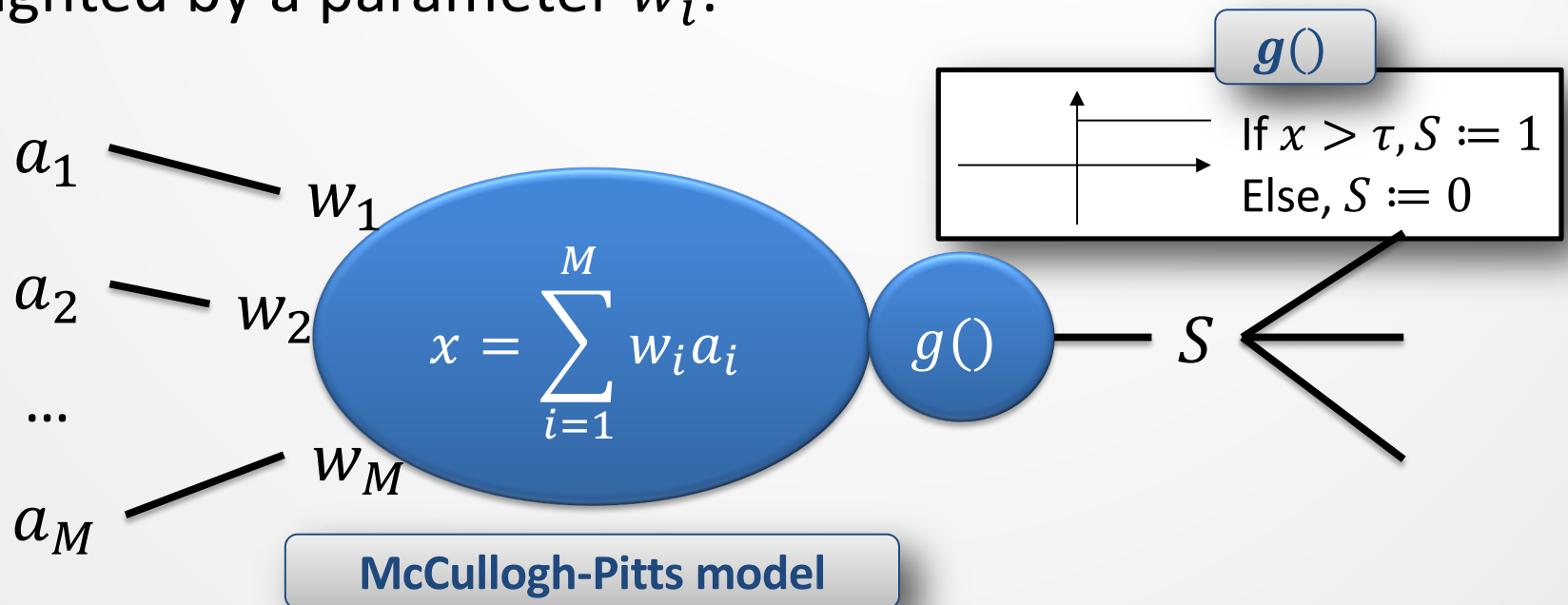
Artificial neural networks

- **Artificial neural networks (ANNs)** were (kind of) inspired from neurobiology (Widrow and Hoff, 1960).
 - Each unit has many inputs (**dendrites**), one output (**axon**).
 - The **nucleus** fires (sends an electric signal along the axon) given input from other neurons.
 - ‘Learning’ occurs at the **synapses** that connect neurons, either by amplifying or attenuating signals.



Perceptron: an artificial neuron

- Each neuron calculates a **weighted sum** of its inputs and compares this to a threshold, τ . If the sum exceeds the threshold, the neuron fires.
- Inputs a_i are activations from adjacent neurons, each weighted by a parameter w_i .

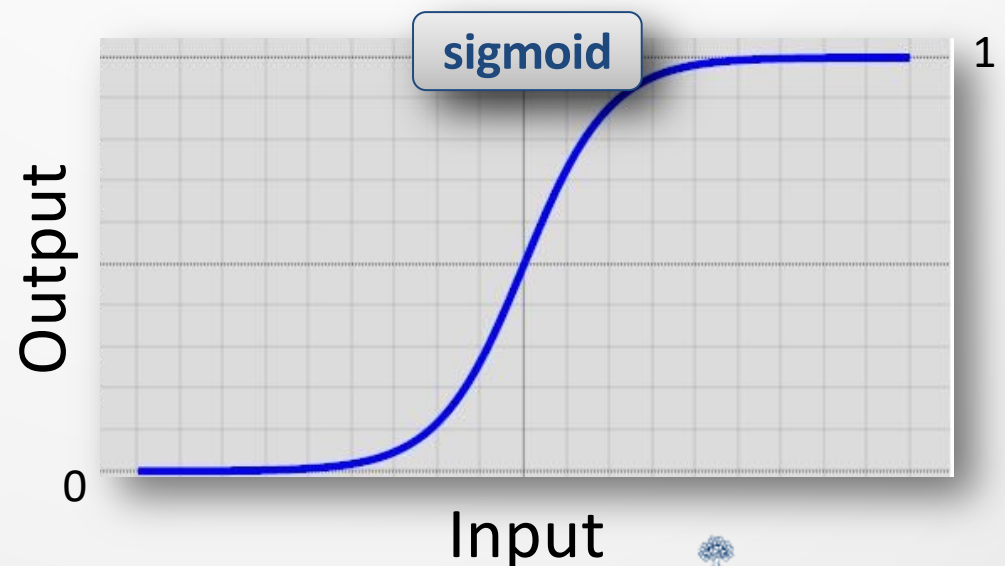
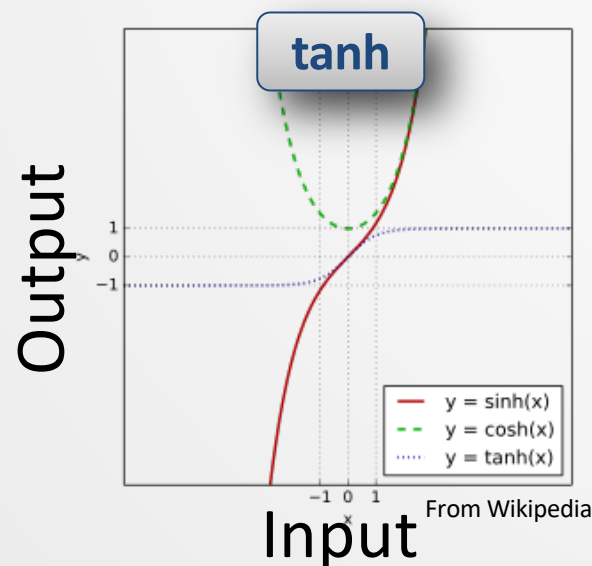


Perceptron output

- Perceptron output is determined by **activation functions**, $g()$, which **can be non-linear functions** of weighted input.
- Popular activation functions include **tanh** and the **sigmoid**:

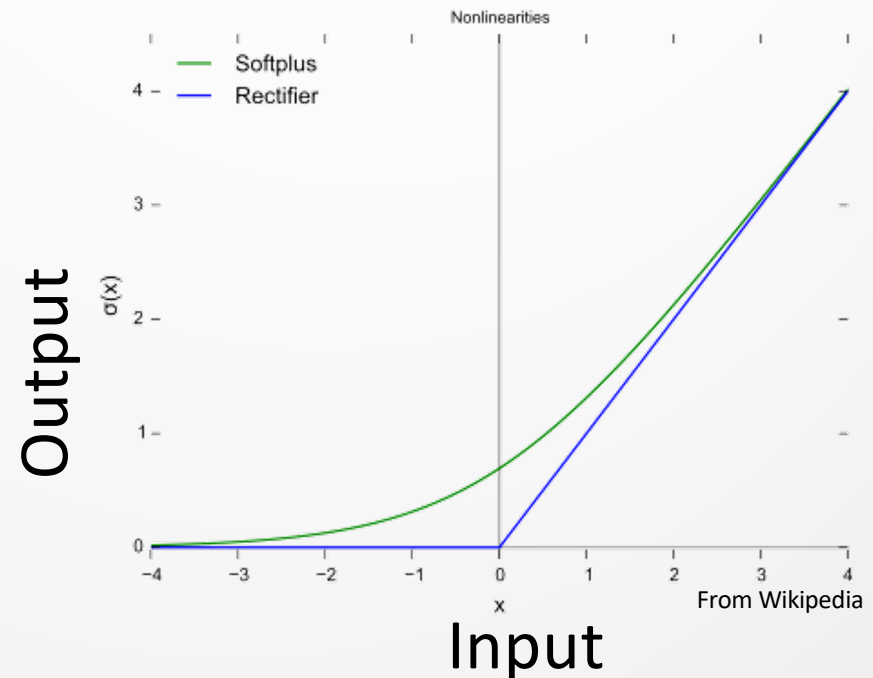
$$g(x) = \sigma(x) = \frac{1}{1 + e^{\rho x}}$$

- 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 more popular.
 - More biologically plausible, sparse activation, limited (vanishing) gradient problems, efficient computation.
- 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?*



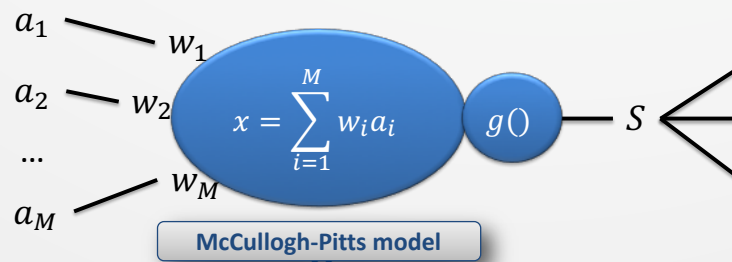
X Glorot, A Bordes, Y Bengio (2011). Deep sparse rectifier neural networks. AISTATS.

Perceptron learning

- 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, α (e.g., 0.05), we can repeatedly adjust each of the weight parameters by

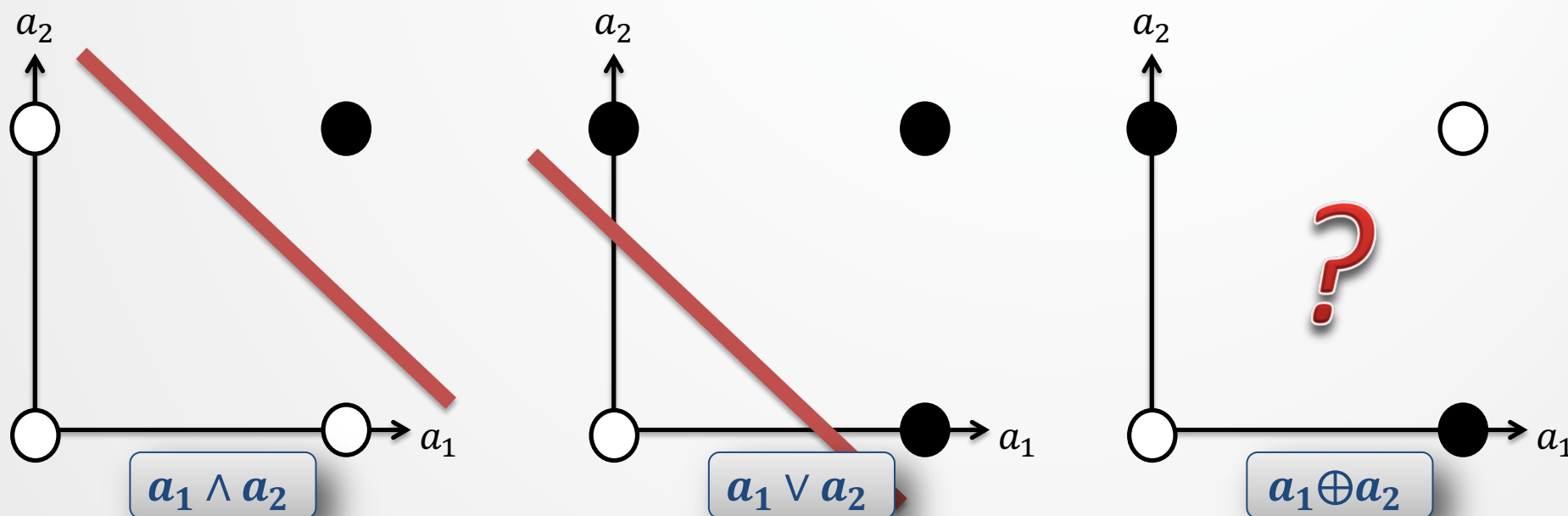
$$w_j := w_j + \alpha \cdot \sum_{i=1}^R \text{Err}_i \cdot g'(x_i) \cdot a_j[i] \quad \left. \vphantom{\sum_{i=1}^R} \right\} \begin{array}{l} \text{Assumes} \\ \text{mean-square} \\ \text{error objective} \end{array}$$

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



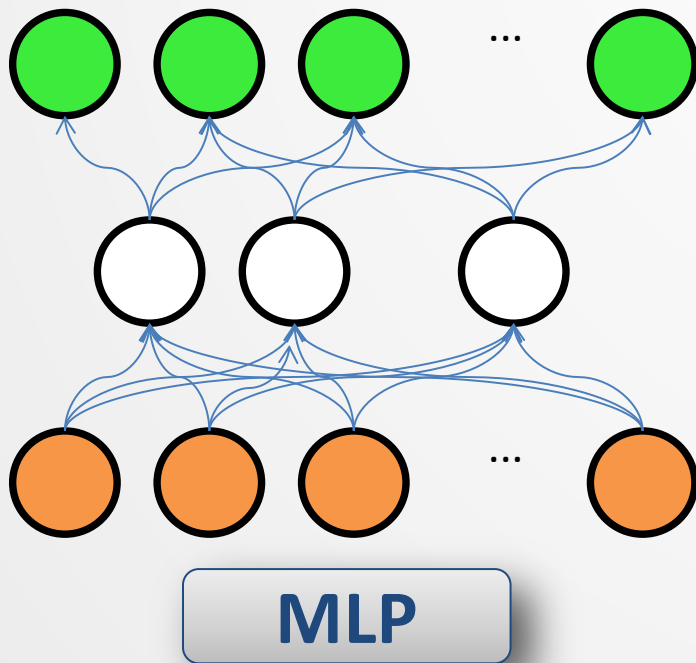
Threshold perceptrons and XOR

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



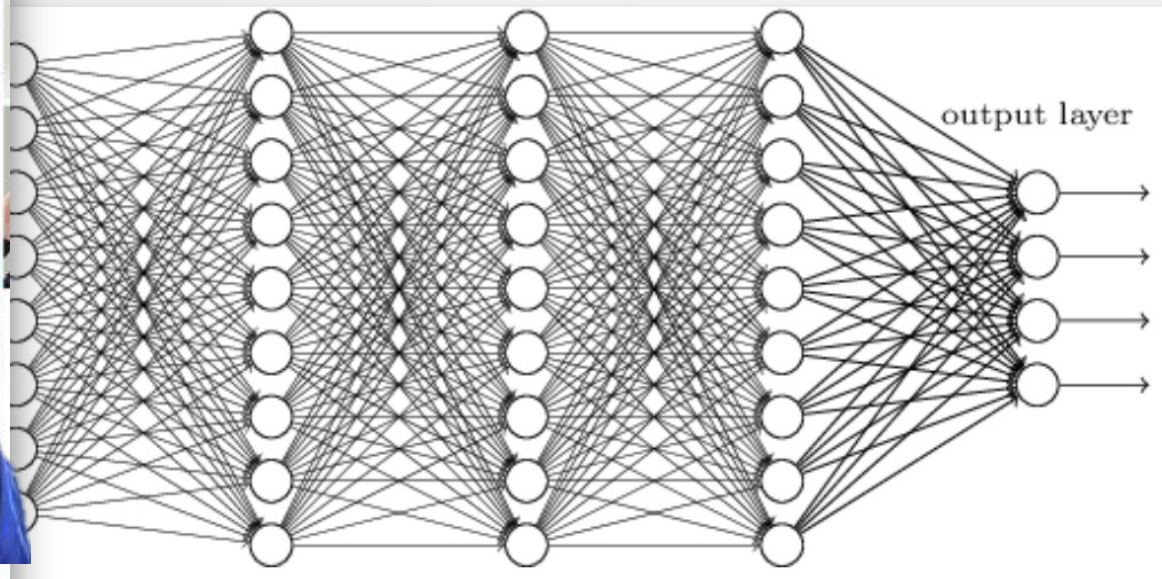
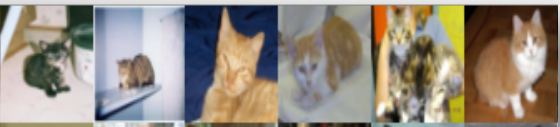
Artificial neural networks

- 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**, and are trained specifically to reduce error.

Deep



Depression.

‘hidden’ representations are learned here
Can we find hidden patterns in words?

Words

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

lugubrious

0	0	0	0	..	0	1	0	...	0
---	---	---	---	----	---	---	---	-----	---

 D

- Classic **word-feature representation** assigns **features** to each index in a much denser vector.
 - E.g., 'VBG', 'negative', 'age-of-acquisition'.

1	0.8	2.5	0.81	...	99
---	-----	-----	------	-----	----

 $d \ll D$

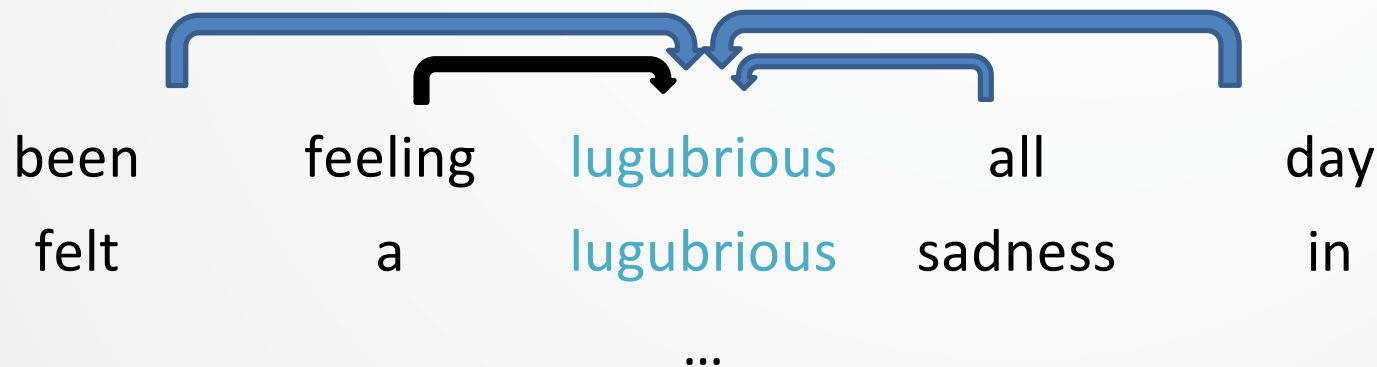
- Can we learn a dense representation? What will it give us?

Learning word semantics

"You shall know a word by the company it keeps."

— J.R. Firth (1957)

$$P(w_t = \text{lugubrious} | w_{t-1} = \text{feeling}, w_{t-2} = \text{been}, \dots)$$

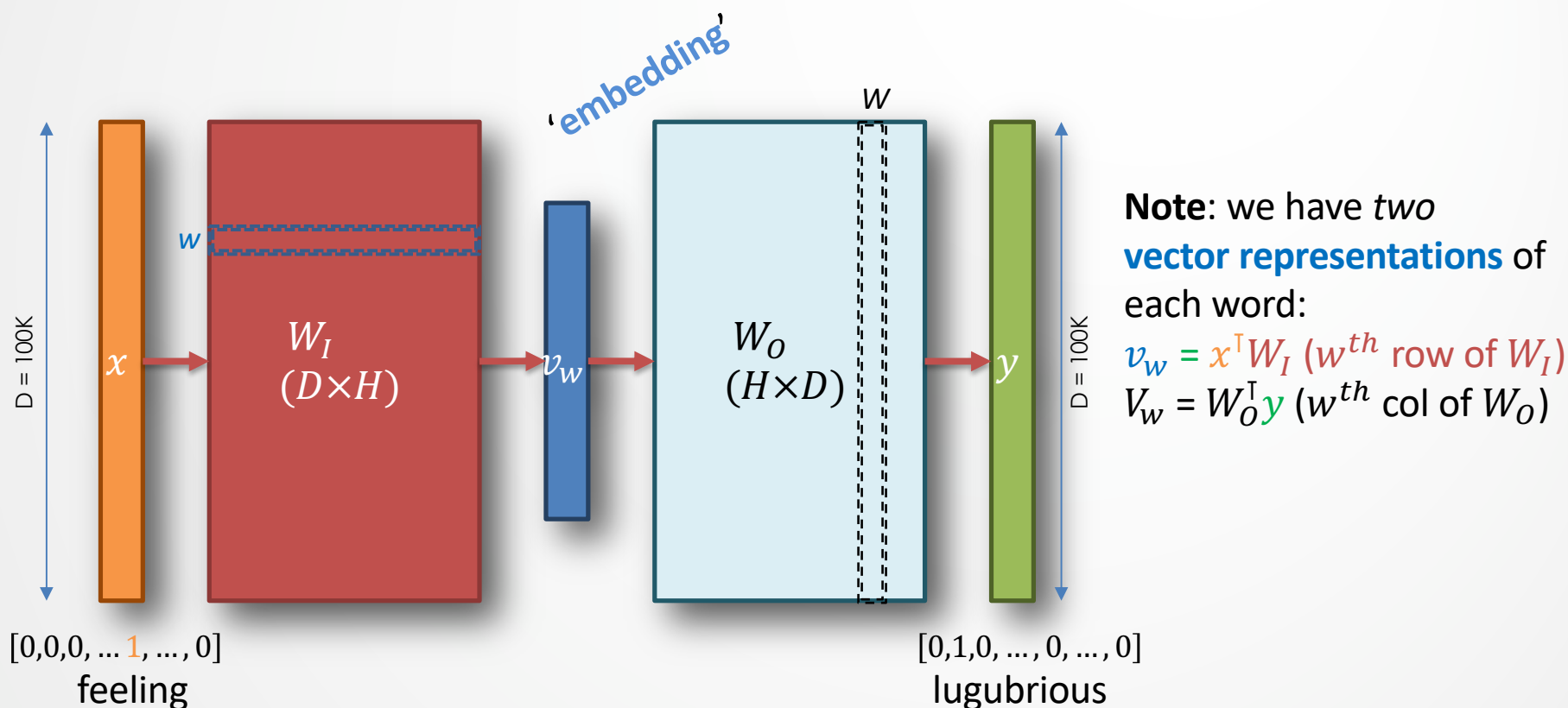


Here, we're predicting the *center* word given the context.
This is called the '**continuous bag of words**' (**CBOW**) model¹.

¹ Mikolov T, Corrado G, Chen K, *et al.* Efficient Estimation of Word Representations in Vector Space. *Proc (ICLR 2013)* 2013;;1–12.

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

Continuous bag of words (1 word context)



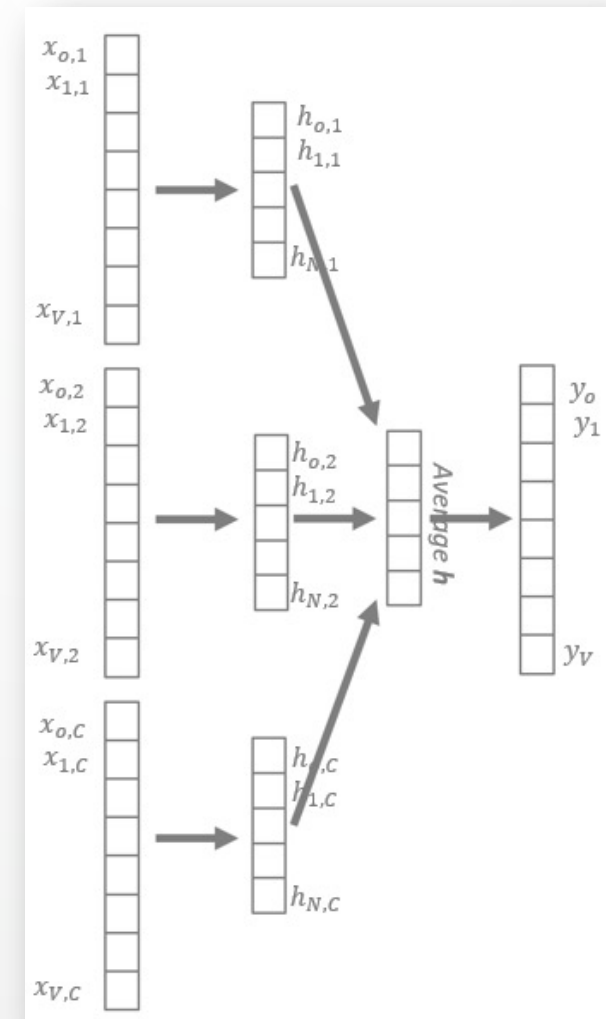
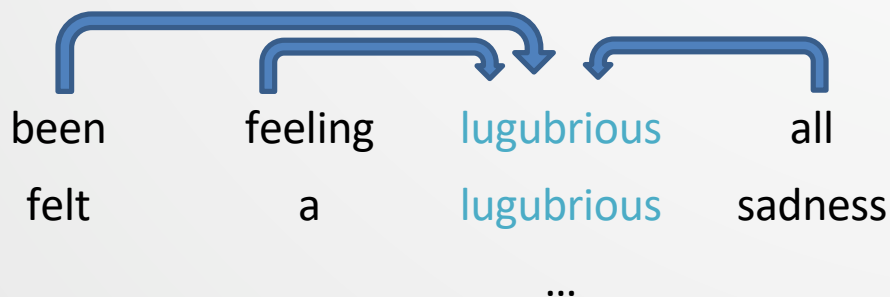
'softmax': $P(w_o | w_i) = \frac{\exp(V_{w_o}^T v_{w_i})}{\sum_{w=1}^W \exp(V_w^T v_{w_i})}$

Where

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

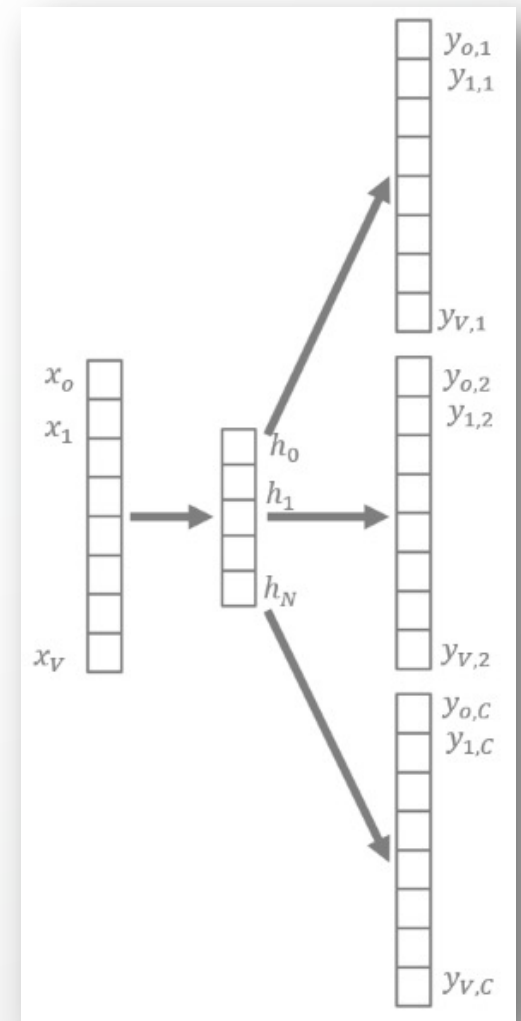
Continuous bag of words (C words context)

- If we want to use **more context**, C , we need to change the network architecture somewhat.
 - Each input word will produce one of C embeddings
 - We just need to add an **intermediate layer**, usually this just averages the embeddings.



Skip-grams

- **Skip-grams** invert the task – we predict context words given the current word.
- According to Mikolov,
Skip-gram: works well with small amounts of training data, represents rare words.
- CBOW**: several times faster to train, slightly better accuracy for frequent words



Mikolov T, Corrado G, Chen K, *et al.* Efficient Estimation of Word Representations in Vector Space. *Proc (ICLR 2013)* 2013;:1–12.

<https://arxiv.org/pdf/1301.3781.pdf>

Actually doing the learning

- Given H -dimensional embeddings, and V word types, our parameters, θ , are:

$$\theta = \begin{bmatrix} v_a \\ v_{aardvark} \\ \vdots \\ v_{zymurgy} \\ V_a \\ V_{aardvark} \\ \vdots \\ V_{zymurgy} \end{bmatrix} \in \mathbb{R}^{2V \times H}$$

Actually doing the learning

We have many options. Gradient descent is popular.
We want to optimize, given T tokens of training data,

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

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

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

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 ,

Actually doing the learning

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

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

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

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

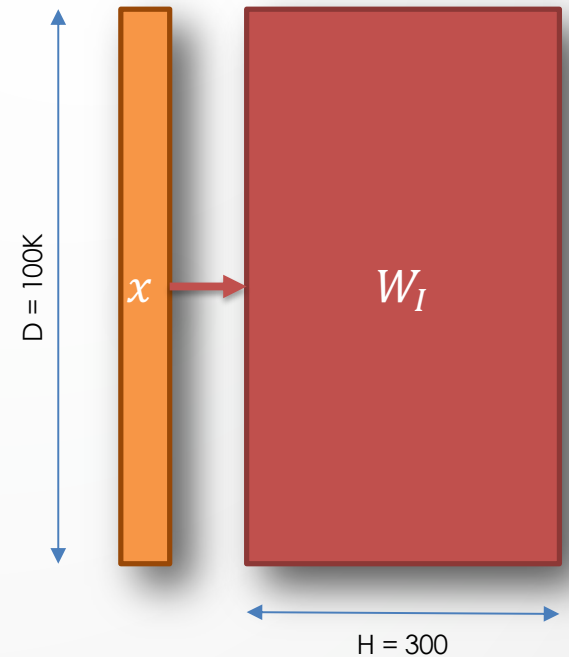
$$v_w = x^T W_I$$

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



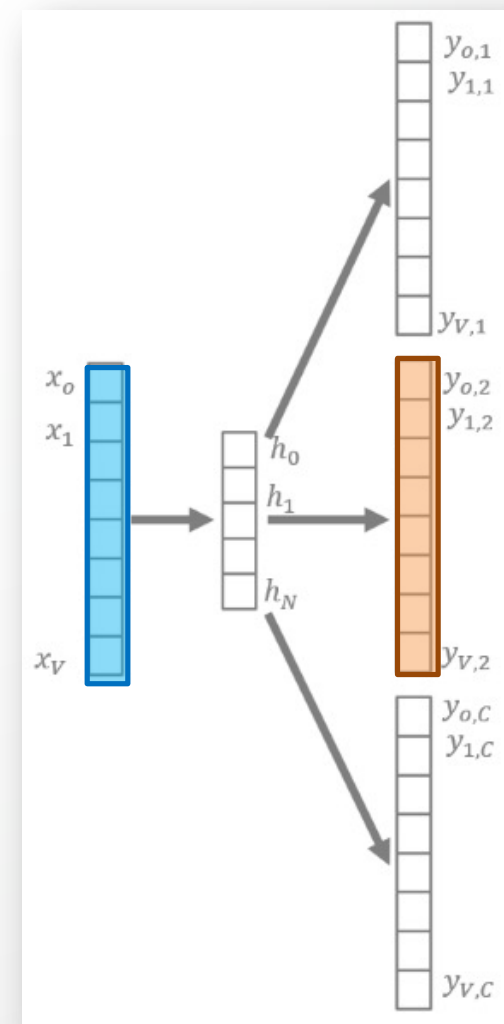
Reminder:

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

Skip-grams with negative sampling

- The default process is inefficient.
 - For one – **what a waste of time!**
We don't want to update $H \times D$ weights!
 - For two – **we want to avoid confusion!**
'Hallucinated' (negative) contexts should be minimized.
- 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:

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



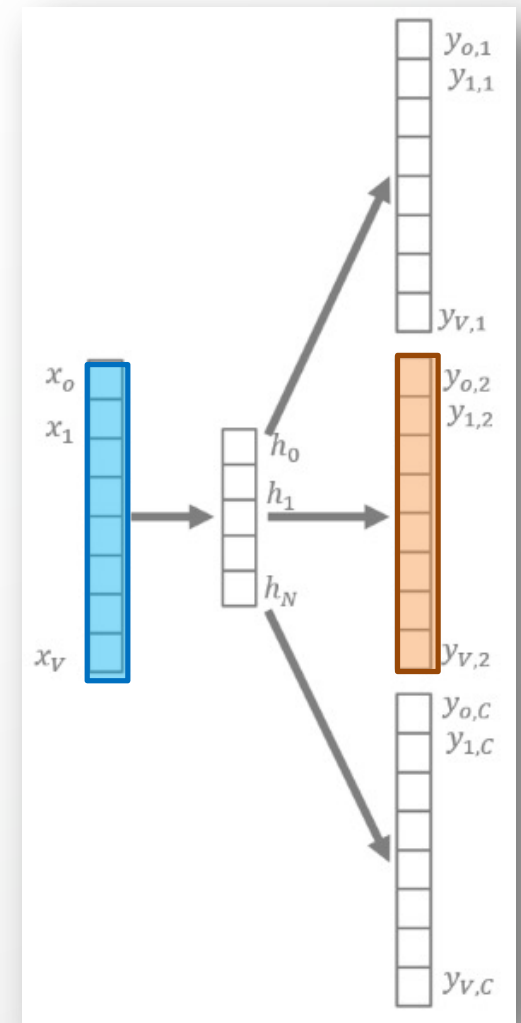
Skip-grams with negative sampling

- 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*



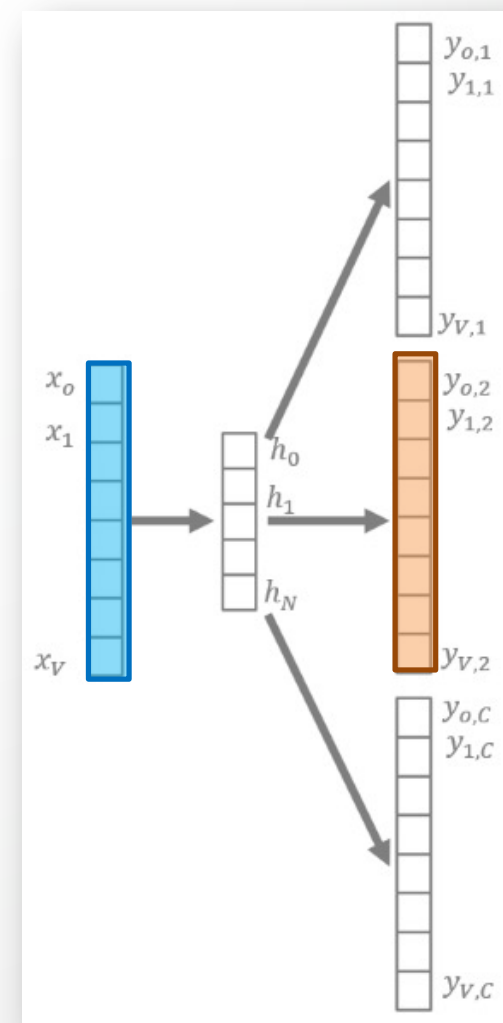
Skip-grams with negative sampling

- Choose a small number k of ‘negative’ words, and just update the weights for the ‘positive’ word plus 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.

$$J(\theta) = \log \sigma(v_o^T v_c) + \sum_{i=1}^k \mathbb{E}_{j \sim P(w)} [\log \sigma(-v_j^T v_c)]$$

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

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



Mimno, D., & Thompson, L. (2017). The strange geometry of skip-gram with negative sampling. *EMNLP 2017*, 2873–2878. [\[link\]](#)

Smell the GloVe

- **GloVe** ('**G**lobal **V**ectors') is an alternative method of obtaining word embeddings.
 - Instead of predicting words at particular positions, look at the **co-occurrence matrix**.

	<i>I</i>	<i>like</i>	<i>enjoy</i>	<i>deep</i>	<i>learning</i>	<i>NLP</i>	<i>flying</i>	<i>.</i>
<i>I</i>	0	2	1	0	0	0	0	0
<i>like</i>	2	0	0	1	0	1	0	0
<i>enjoy</i>	1	0	0	0	0	0	1	0
<i>deep</i>	0	1	0	0	1	0	0	0
<i>learning</i>	0	0	0	1	0	0	0	1
<i>NLP</i>	0	1	0	0	0	0	0	1
<i>flying</i>	0	0	1	0	0	0	0	1
<i>.</i>	0	0	0	0	1	1	1	0

Word w_i occurs
 $X_{i,j} (= X_{j,i})$
times with word w_j ,
within some context
window (e.g., 10 words,
a sentence, ...).

Pennington J, Socher R, Manning CD. (2014) GloVe: Global Vectors for Word Representation.

Proc EMNLP 2014:1532–43. doi:10.3115/v1/D14-1162 <https://nlp.stanford.edu/projects/glove/>

Smell the GloVe

- Populating the co-occurrence matrix requires a complete pass through the corpus, but needs only be done once.
- Let $P_{i,j} = P(w_j|w_i) = X_{i,j}/X_i$,

Table 1: Co-occurrence probabilities for target words *ice* and *steam* with selected context words from a 6 billion token corpus. Only in the ratio does noise from non-discriminative words like *water* and *fashion* cancel out, so that large values (much greater than 1) correlate well with properties specific to ice, and small values (much less than 1) correlate well with properties specific of steam.

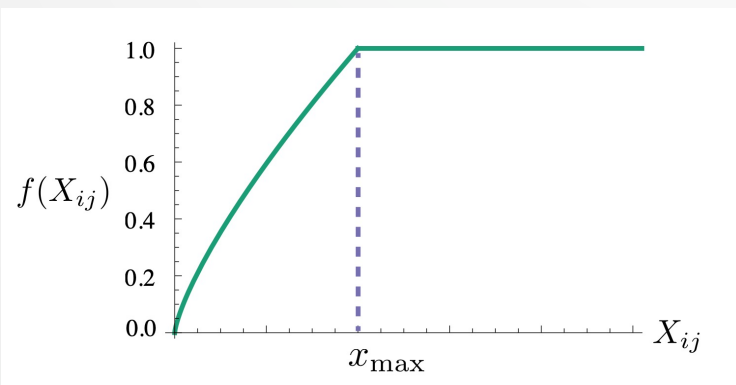
Probability and Ratio	$k = solid$	$k = gas$	$k = water$	$k = fashion$
$P(k ice)$	1.9×10^{-4}	6.6×10^{-5}	3.0×10^{-3}	1.7×10^{-5}
$P(k steam)$	2.2×10^{-5}	7.8×10^{-4}	2.2×10^{-3}	1.8×10^{-5}
$P(k ice)/P(k steam)$	8.9	8.5×10^{-2}	1.36	0.96

Pennington J, Socher R, Manning CD. (2014) GloVe: Global Vectors for Word Representation.
Proc EMNLP 2014:1532–43. doi:10.3115/v1/D14-1162 <https://nlp.stanford.edu/projects/glove/>

Aside – smell the GloVe

- Minimize $J = \sum_{i,j=1}^V f(X_{i,j}) \left(v_{w_i}^T v_{w_j} + b_i + \tilde{b}_j - \log X_{i,j} \right)^2$
where, b_i and \tilde{b}_j are input and output bias terms associated with w_i and w_j , respectively
- Weighting function $f(X_{i,j})$:

$$f(x) = \begin{cases} (x/x_{\max})^\alpha & \text{if } x < x_{\max} \\ 1 & \text{otherwise} \end{cases}.$$



Weighting function f with $\alpha = 3/4$, $x_{\max} = 100$

- $f(0) = 0$. If f is viewed as a continuous function, it should vanish as $x \rightarrow 0$ fast enough that the $\lim_{x \rightarrow 0} f(x) \log^2 x$ is finite.
- $f(x)$ should be non-decreasing so that rare co-occurrences are not overweighted.
- $f(x)$ should be relatively small for large values of x , so that frequent co-occurrences are not overweighted.

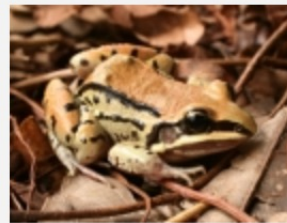
Aside – evaluation

- **Intrinsic evaluation:** popular Redacted method was to cherry-pick a few k -nearest neighbours examples that match expectations.

0. frog
1. frogs
2. toad
3. litoria
4. leptodactylidae
5. rana
6. lizard
7. eleutherodactylus



3. litoria



4. leptodactylidae



5. rana



7. eleutherodactylus

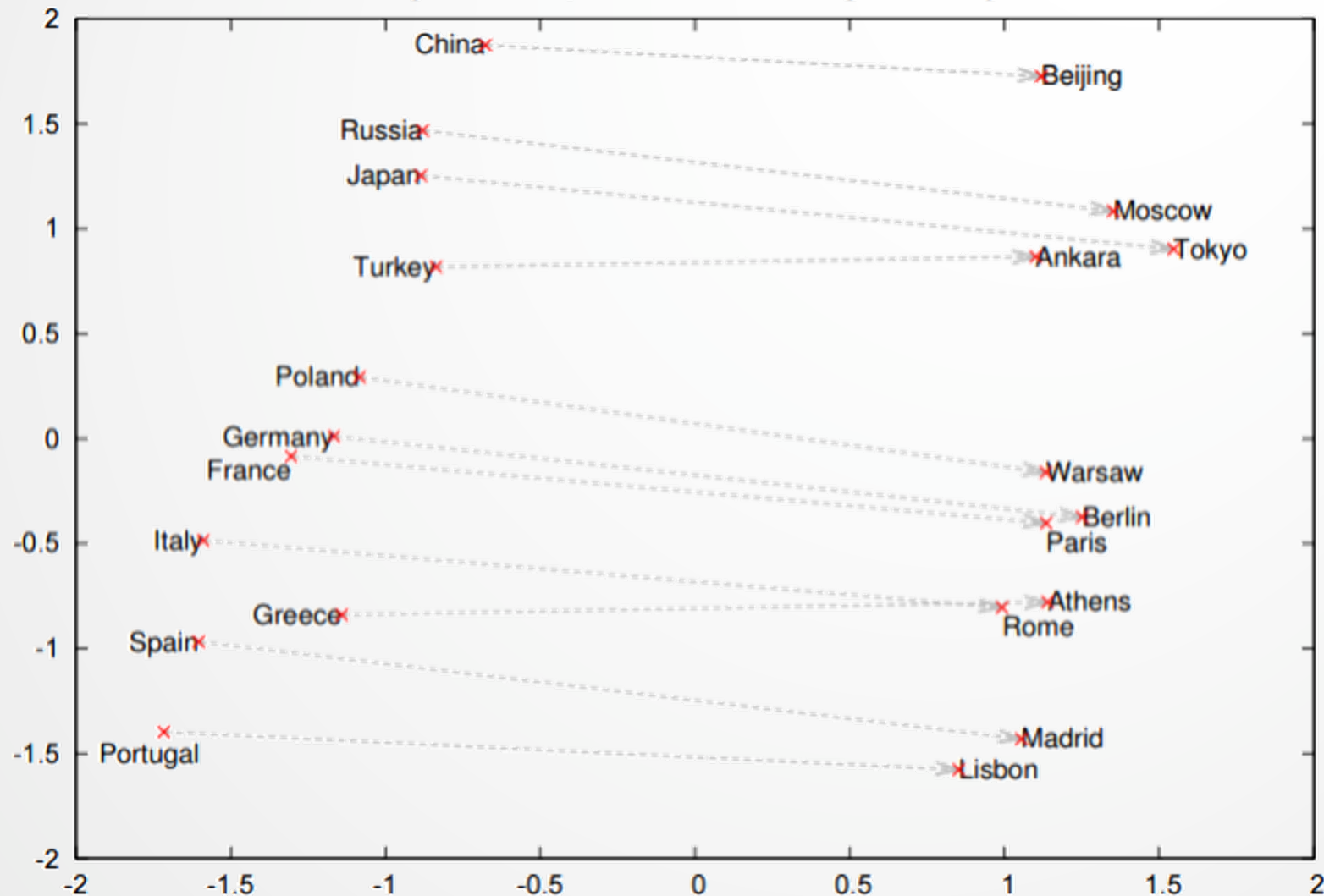
- **Extrinsic evaluation:** embed resulting vectors into a variety of tasks^[1,2].

SuperGLUE		GLUE		Leaderboard Version: 2.0										
Rank	Name	Model	URL	Score	BoolQ	CB	COPA	MultiRC	ReCoRD	RTE	WiC	WSC	AX-b	AX-g
1	Liam Fedus	SS-MoE		91.0	92.3	96.9/98.0	99.2	89.2/65.2	95.0/94.2	93.5	77.4	96.6	72.3	96.1/94.1
2	Microsoft Alexander v-team	Turing NLR v5		90.9	92.0	95.9/97.6	98.2	88.4/63.0	96.4/95.9	94.1	77.1	97.3	67.8	93.3/95.5
3	ERNIE Team - Baidu	ERNIE 3.0		90.6	91.0	98.6/99.2	97.4	88.6/63.2	94.7/94.2	92.6	77.4	97.3	68.6	92.7/94.7

¹ <https://gluebenchmark.com/tasks>

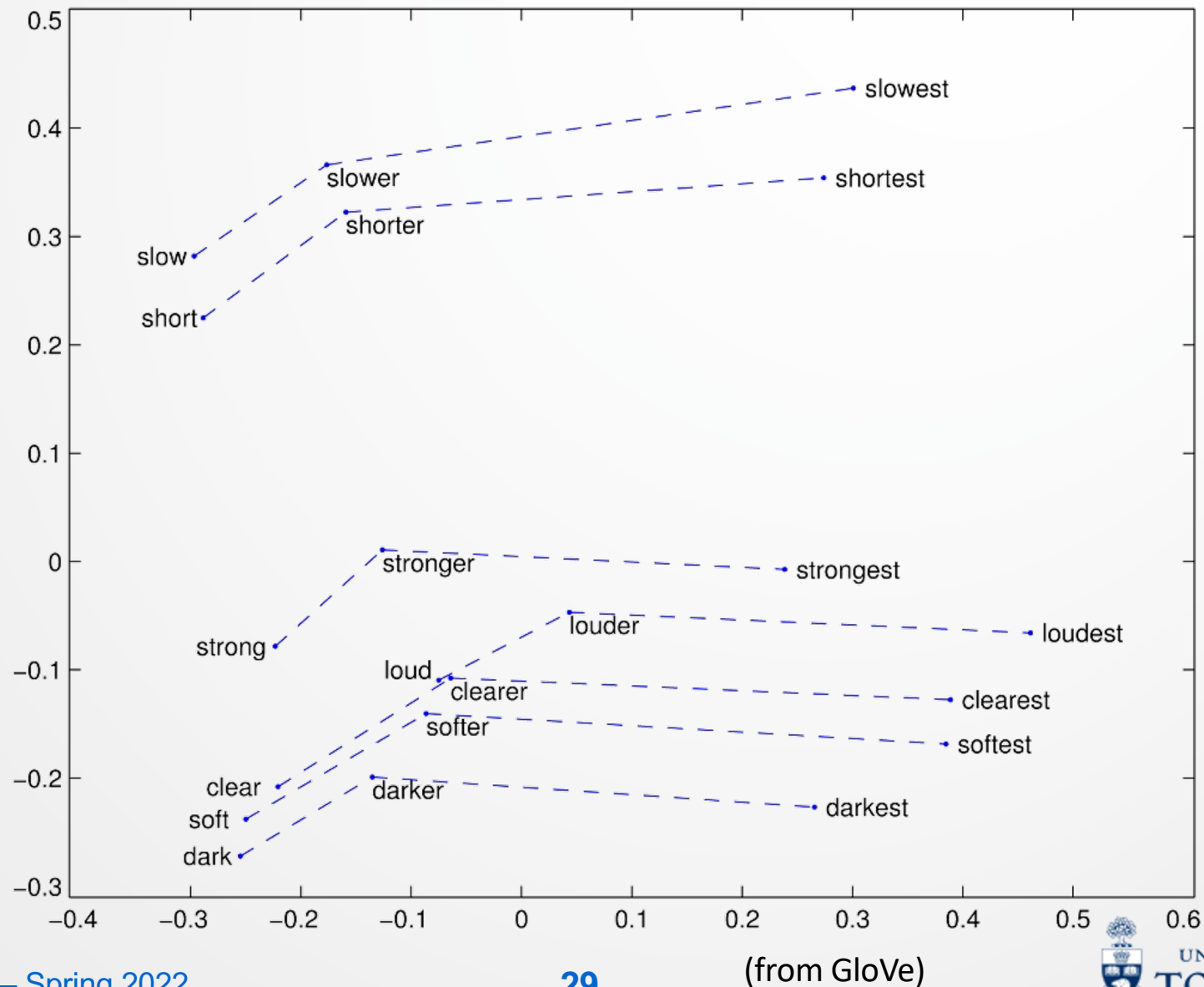
² <https://super.gluebenchmark.com/tasks>

Linguistic regularities in vector space



Trained on the Google news corpus with over 300 billion words.

Linguistic regularities in vector space



Linguistic regularities in vector space

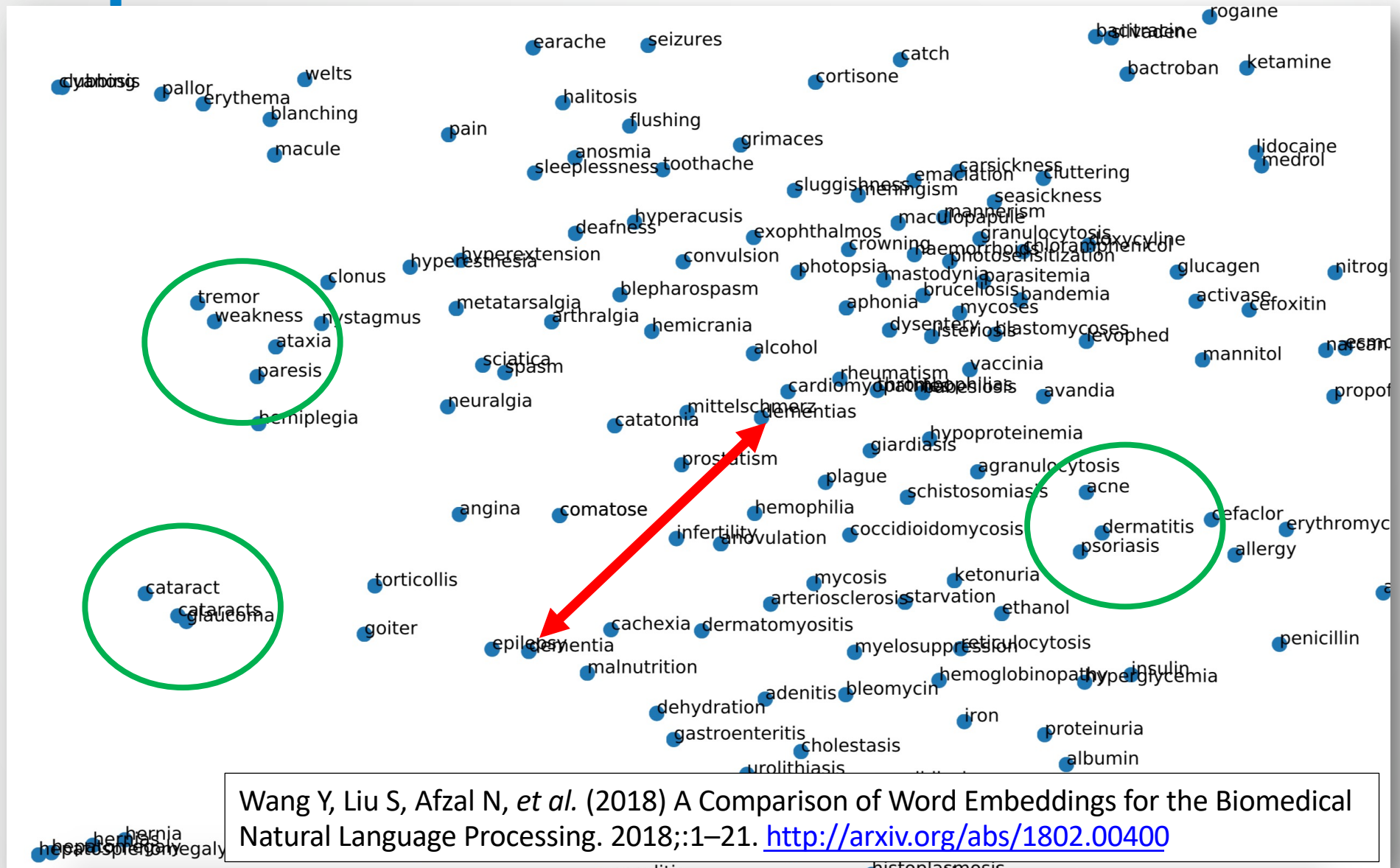
Expression	Nearest token
Paris – France + Italy	Rome
Bigger – big + cold	Colder
Sushi – Japan + Germany	bratwurst
Cu – copper + gold	Au
Windows – Microsoft + Google	Android

Analogies: apple:apples :: octopus:octopodes

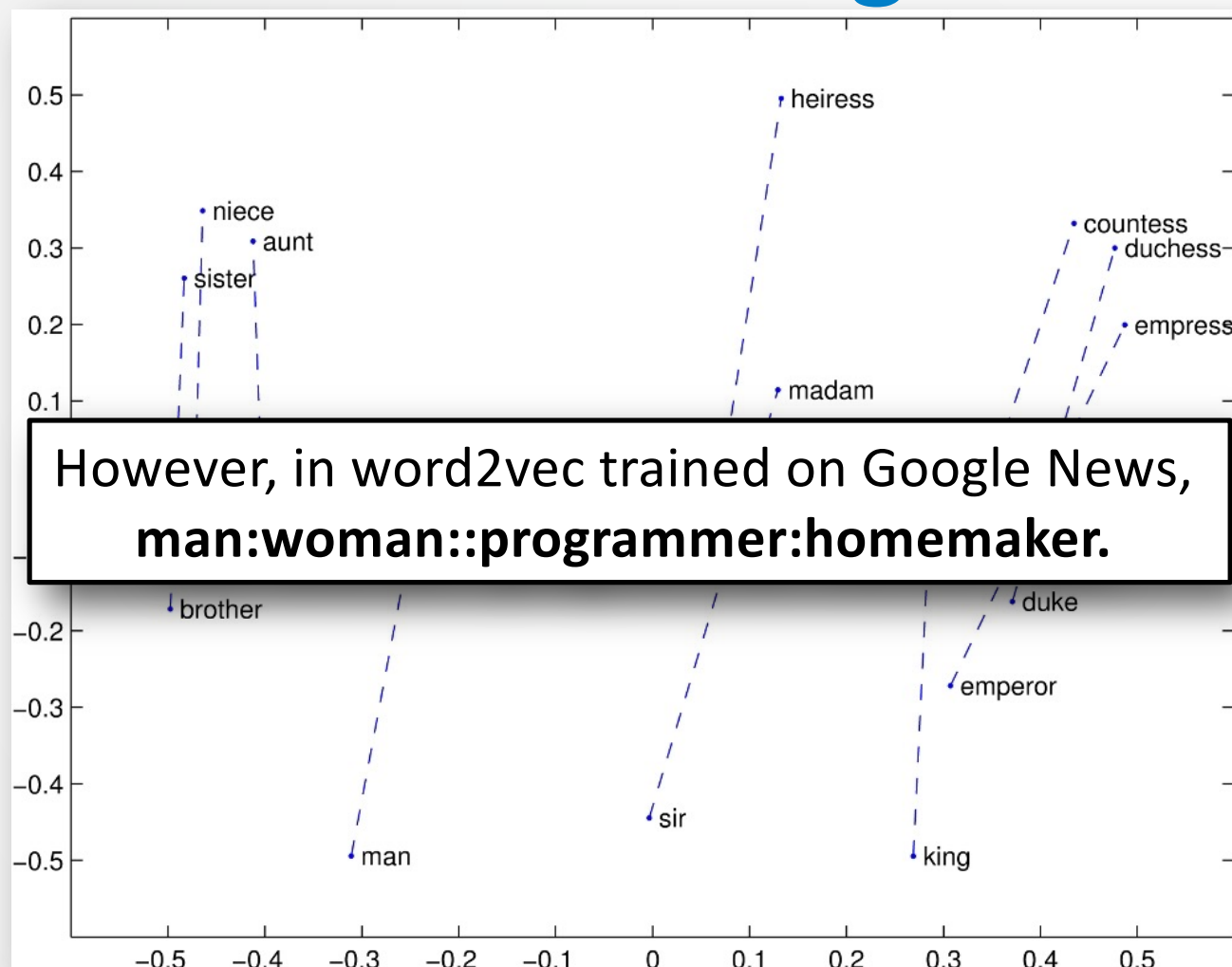
Hypernymy: shirt:clothing :: chair:furniture

Semantic: queen – king \approx woman – man

Importance of in-domain data



Biases: let's talk about gender



Bolukbasi T, Chang K, Zou J, *et al.* Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings. In: *NIPS*. 2016. 1–9.

Biases: let's talk about gender

Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings

Tolga Bolukbasi¹, Kai-Wei Chang², James Zou², Venkatesh Saligrama^{1,2}, Adam Kalai²

¹Boston University, 8 Saint Mary's Street, Boston, MA

²Microsoft Research New England, 1 Memorial Drive, Cambridge, MA

tolgab@bu.edu, kw@kwchang.net, jamesyzou@gmail.com, srv@bu.edu, adam.kalai@microsoft.com

Abstract

The blind application of machine learning runs the risk of amplifying biases present in data. Such a danger is facing us with *word embedding*, a popular framework to represent text data as vectors which has been used in many machine learning and

Extreme *she*

1. homemaker
2. nurse
3. receptionist
4. librarian
5. socialite
6. hairdresser
7. nanny
8. bookkeeper
9. stylist
10. housekeeper

Extreme *he*

1. maestro
2. skipper
3. protege
4. philosopher
5. captain
6. architect
7. financier
8. warrior
9. broadcaster
10. magician

sewing-carpentry
nurse-surgeon
blond-burly
giggle-chuckle
sassy-snappy
volleyball-football

queen-king
waitress-waiter

Gender stereotype *she-he* analogies

registered nurse-physician
interior designer-architect
feminism-conservatism
vocalist-guitarist
diva-superstar
cupcakes-pizzas

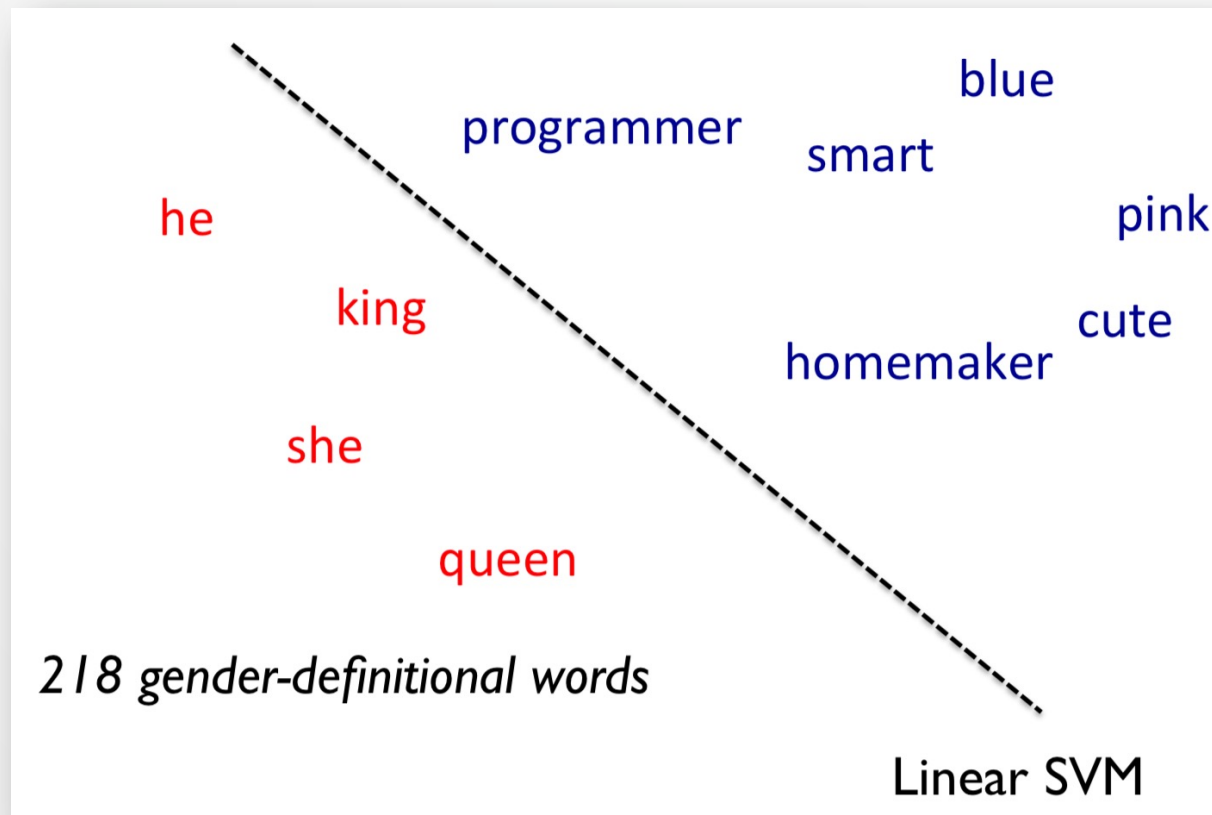
Gender appropriate *she-he* analogies

sister-brother
ovarian cancer-prostate cancer
mother-father
convent-monastery

housewife-shopkeeper
softball-baseball
cosmetics-pharmaceuticals
petite-lanky
charming-affable
lovely-brilliant

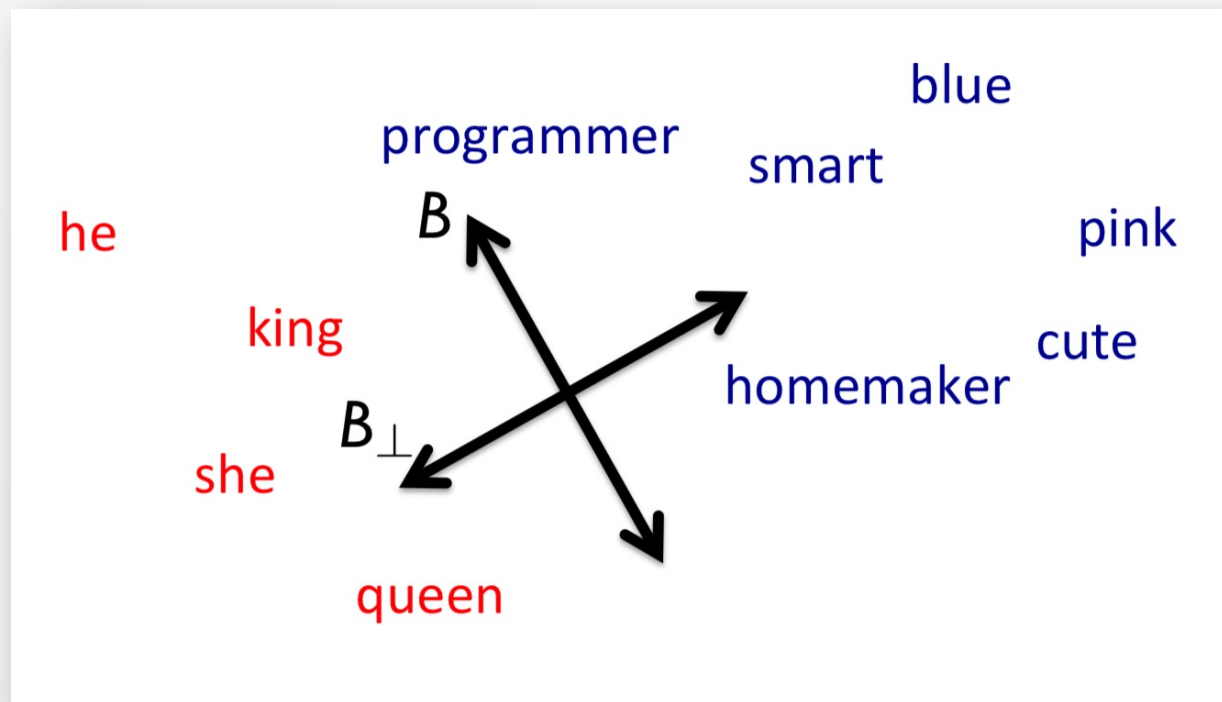
Solution?

1. Hand-pick words S_0 that are 'gender definitional'.
'Neutral' words are the complement, $N = V \setminus S_0$.



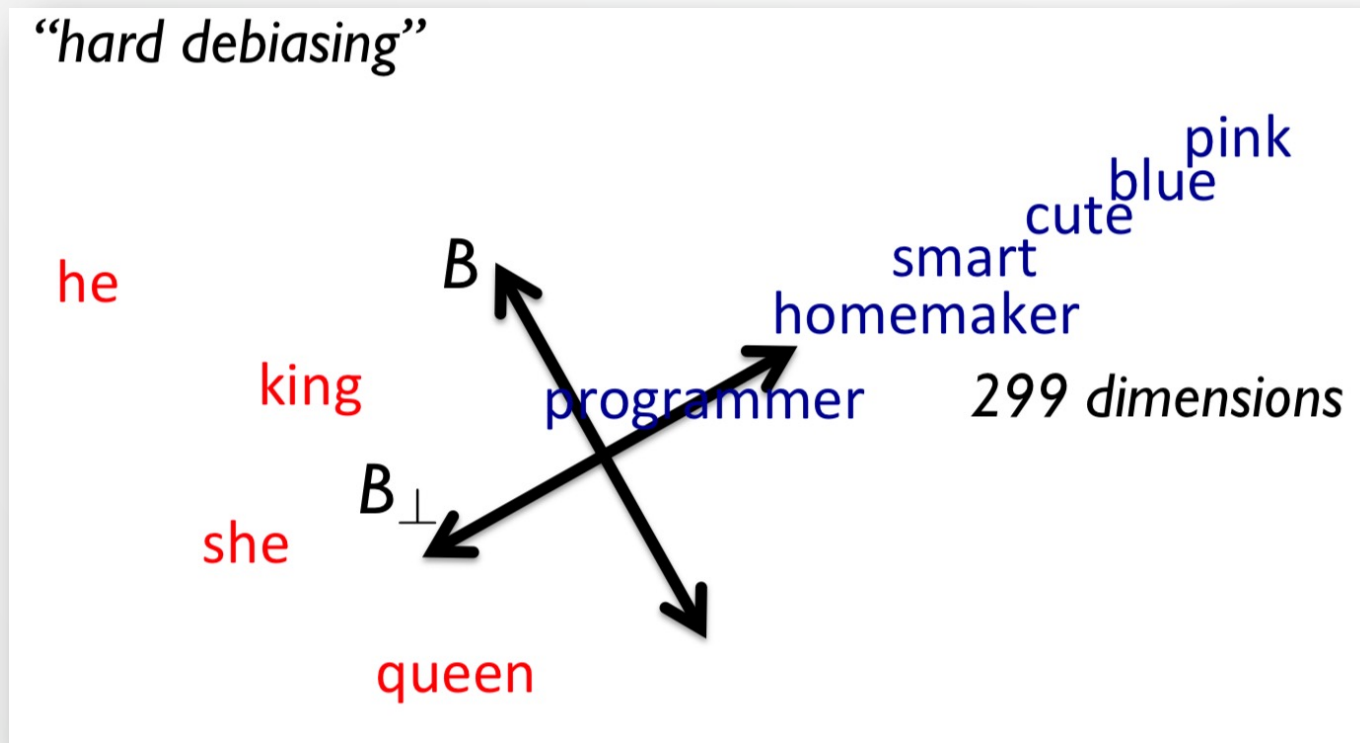
Solution?

2. Project away gender subspace from gender-neutral words, $w := w - w \cdot B$ for $w \in N$, where B is the gender subspace.



Solution?

2. Project away gender subspace from gender-neutral words, $\tilde{w} := w - w \cdot B$ for $w \in N$, where B is the gender subspace.



Results

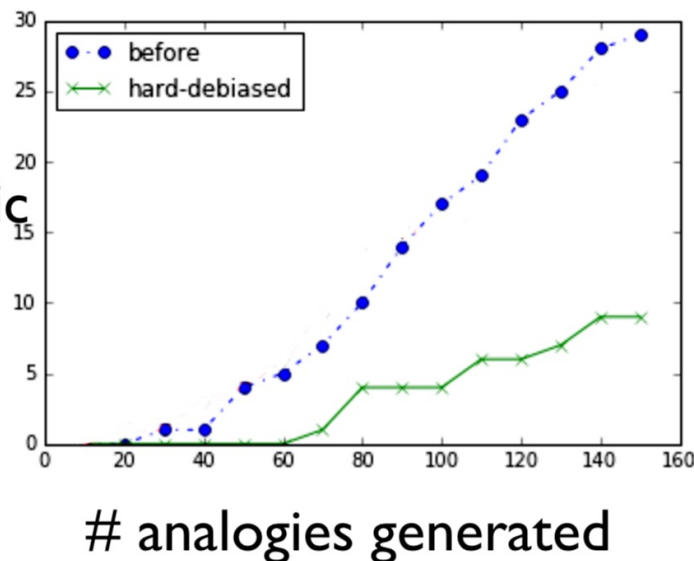
- Generate many analogies, see which ones preserve gender stereotypes.

He:Blue :: She: ?

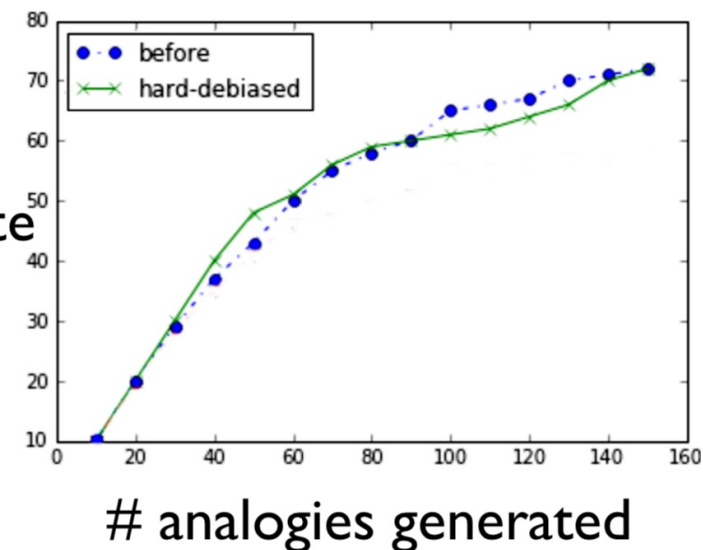
He:Doctor :: She: ?

He:Brother :: She: ?

stereotypic
analogies



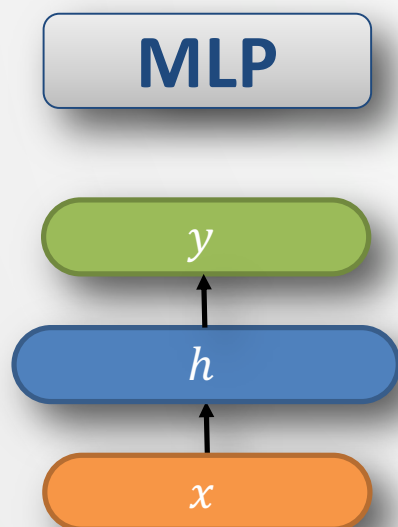
appropriate
analogies



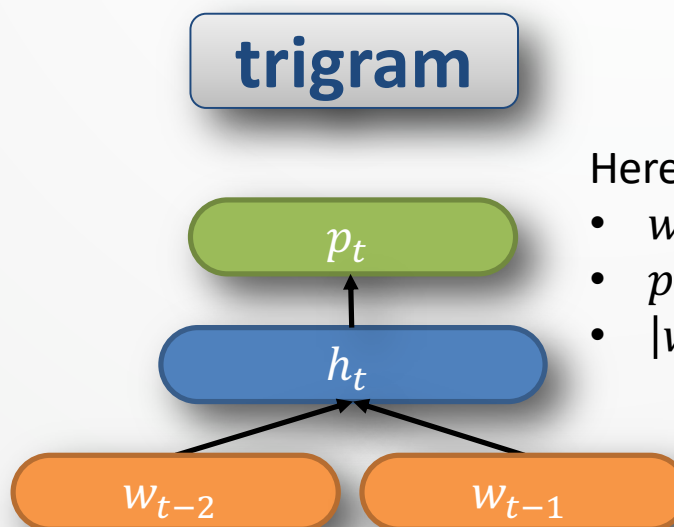
NEURAL LANGUAGE MODELS

Trigram models

- CBOW: prediction of current word w_t given w_{t-1} .
- Let's reconsider predicting w_t given multiple w_{t-j} ?
 - I.e., let's think about **language modelling**.



$$h = g(W_I x + c)$$
$$y = W_O h + b$$



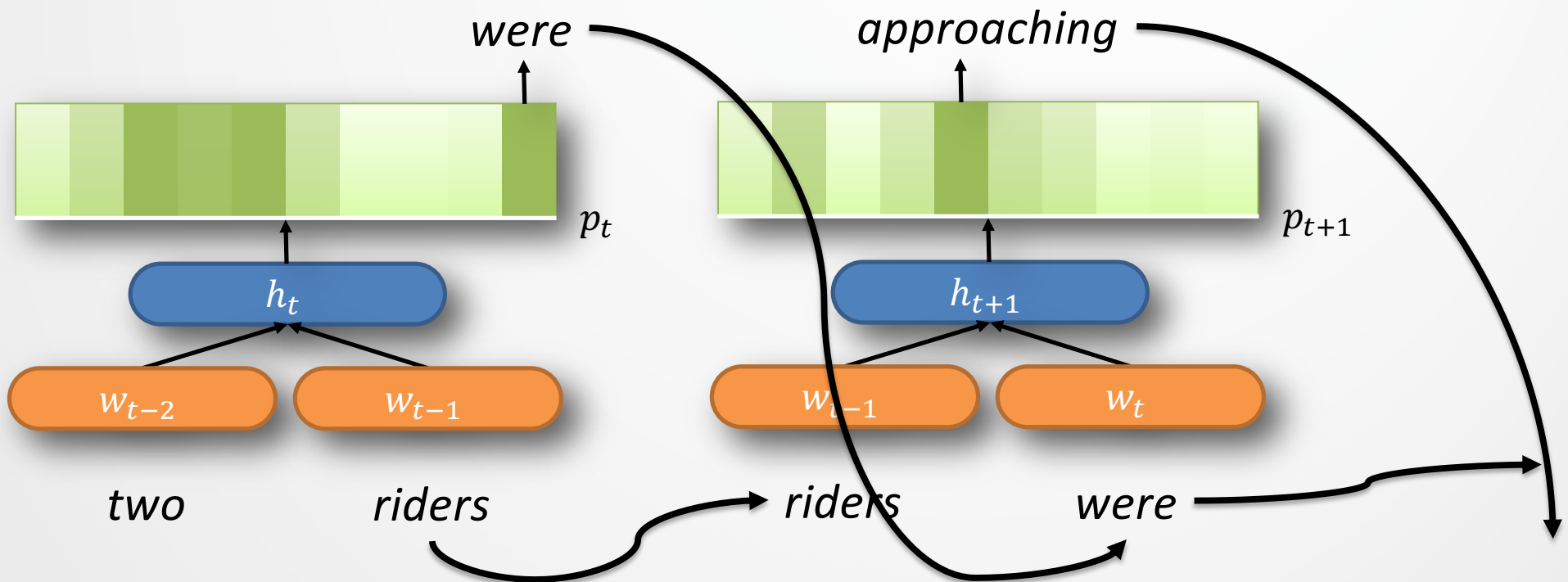
$$h_t = g(W_I [w_{t-2}; w_{t-1}] + c)$$
$$p_t = \text{softmax}(W_O h_t + b)$$

Here:

- w_i is a one-hot vector,
- p_t is a distribution, and
- $|w_i| = |p_t| = |V|$
(i.e., the size of the vocabulary)

Sampling from trigram models

- Since $p_t \sim P(w_t | w_{t-2} w_{t-1})$, we just feed forward and sample from the output vector.



Training trigram models

- Here's one approach:
 1. Randomly choose a batch (e.g., 10K consecutive words)
 2. Propagate words through the current model
 3. Obtain word likelihoods (loss)
 4. Back-propagate loss
 5. Gradient step to update model
 6. Go to (1)

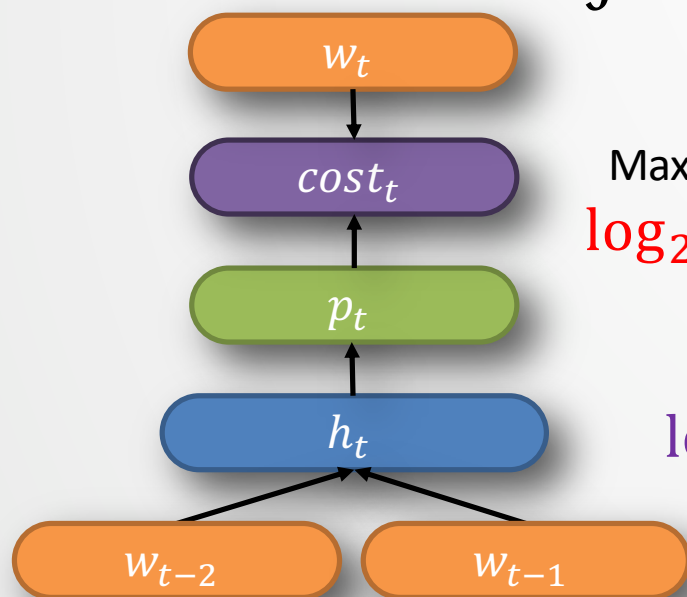
Training trigram models

- The typical training objective is the cross entropy (see Lecture 4) of the corpus C given the model M :

$$\mathcal{F} = H(C; M) = - \frac{\log_2 P_M(C)}{\|C\|}$$

Minimize

$$\text{Maximize } \log_2 P_M(C) = \log_2 \prod_{t=0}^T P(w_t) = \sum_{t=0}^T \log_2 P(w_t)$$



$$\log_2 P(w_t) = \mathbf{w}_t^T \log p_t$$

$$h_t = g(W_I[\mathbf{w}_{t-2}; \mathbf{w}_{t-1}] + c)$$

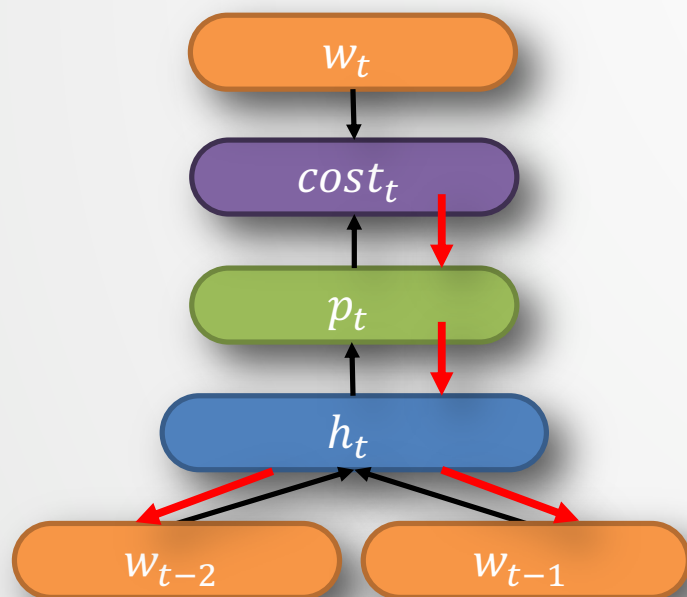
$$p_t = \text{softmax}(W_O \mathbf{h}_t + b)$$

Here:

- w_i is a one-hot vector, and
- p_t is a distribution.

Training trigram models

- Compute our gradients, using $\mathcal{F} = -\frac{\log_2 P_M(C)}{\|C\|}$ and $\log_2 P(w_t) = \mathbf{w}_t^\top \log p_t$ and **back-propagate**.



$$\frac{\delta \mathcal{F}}{\delta W_O} = -\frac{1}{\|C\|} \sum_t \frac{\delta cost_t}{\delta p_t} \frac{\delta p_t}{\delta W_O}$$

$$\frac{\delta \mathcal{F}}{\delta W_I} = -\frac{1}{\|C\|} \sum_t \frac{\delta cost_t}{\delta p_t} \frac{\delta p_t}{\delta h_t} \frac{\delta h_t}{\delta W_I}$$

$$h_t = g(W_I[\mathbf{w}_{t-2}; \mathbf{w}_{t-1}] + c)$$

$$p_t = \text{softmax}(W_O \mathbf{h}_t + b)$$

Here:

- w_i is a one-hot vector, and
- p_t is a distribution.

So what?

- 😊 Neural language models of this type:
 - Can generalize better than MLE LMs to unseen n -grams,
 - Can use *semantic* information as in word2vec.

$$P(\text{the cat sat on the } \mathbf{mat}) \approx P(\text{the cat sat on the } \mathbf{rug})$$

- ☹️ Neural language models of this type:
 - Can take *relatively* long to train. “GPUs kill the Earth.”
 - Number of parameters scale poorly with increasing context.

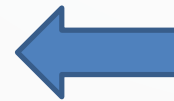
Let's improve both of these issues...

Dealing with that bottleneck

- Traditional datasets for neural language modeling include:
 - AP News (14M tokens, 17K types)
 - HUB-4 (1M tokens, 25K types)
 - Google News (6B tokens, 1M types)
 - Wikipedia (3.2B tokens, 2M types)
- Datasets for **medical/clinical** LM include:
 - EMERALD/ICES (3.5B tokens, 13M types)
- Much of the computational effort is in the initial embedding, and in the softmax.
 - Can we simplify and speed up the process?

Dealing with that bottleneck

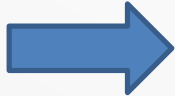
- **Replace** rare words with <out-of-vocabulary> token.
- **Subsample** frequent words.
- Hierarchical softmax.
- Noise-contrastive estimation.
- Negative sampling.



[Morin & Bengio, 2005, Mikolov et al, 2011, 2013b;
Mnih & Teh 2012, Mnih & Kavukcuoglu, 2013]

Hierarchical softmax with grouping

- Group words into distinct classes, c , e.g., by frequency.
 - E.g., c_1 is top 5% of words by frequency, c_2 is the next 5%, ...
- Factorize $p(w_o | w_i) = p(c | w_i) p(w_o | w_i, c)$

'softmax': $P(w_o | w_i) = \frac{\exp(V_{w_o}^T v_{w_i})}{\sum_{w=1}^W \exp(V_w^T v_{w_i})}$  $\frac{\exp(c_j v_{w_i})}{\sum_c \exp(c v_{w_i})} \times \frac{\exp(V_{w_o}^T v_{w_i})}{\sum_{w \in c} \exp(V_w^T v_{w_i})}$

Where

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

[Mikolov et al, 2011, Auli et al, 2013]

RECURRENT NEURAL NETWORKS

Statistical language models

- Probability is conditioned on (window of) n previous words*
- A necessary (but incorrect) Markov assumption: each observation only depends on 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

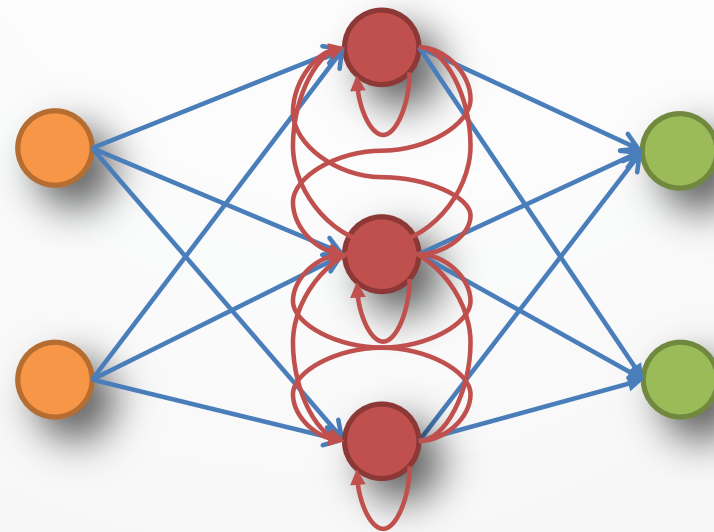
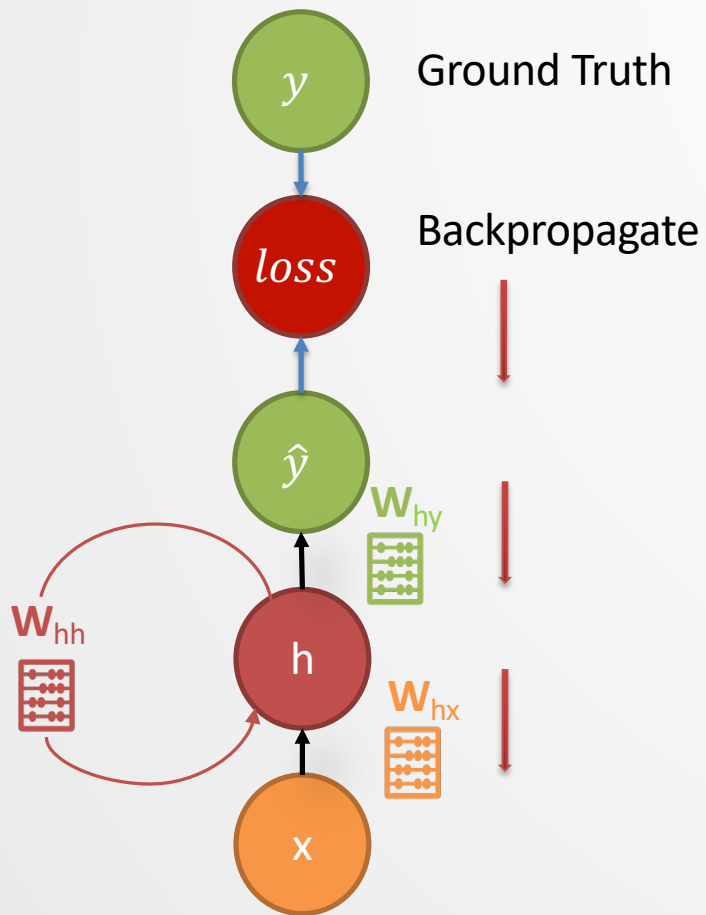
Statistical language models

- Using higher n-gram counts (with smoothing) improves performance*
- *Computational burden*: too many n-grams (combinations)
 - Infeasible RAM requirements
- *RNN intuition*:
 - Use the **same set of weight** parameters for each word (or across all time steps)
 - Condition the neural network on all previous words (or time steps)
 - Memory requirement now scales with number of words

*From Lecture 2

Recurrent neural networks (RNNs)

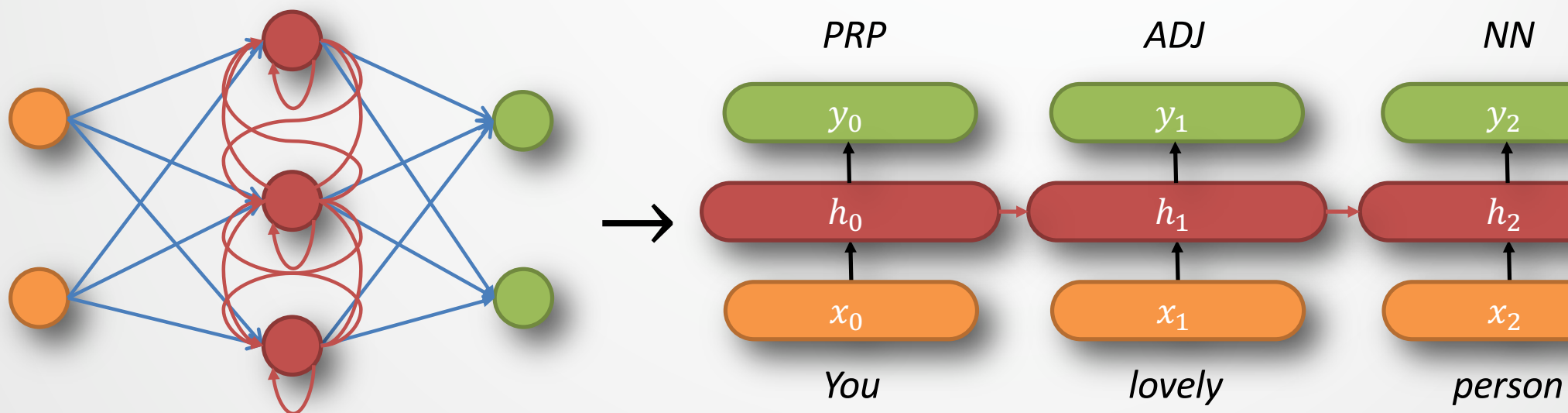
- An RNN has **feedback** connections in its structure so that it *'remembers'* previous states, when reading a sequence.



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.
 - These can be applied to various tasks.



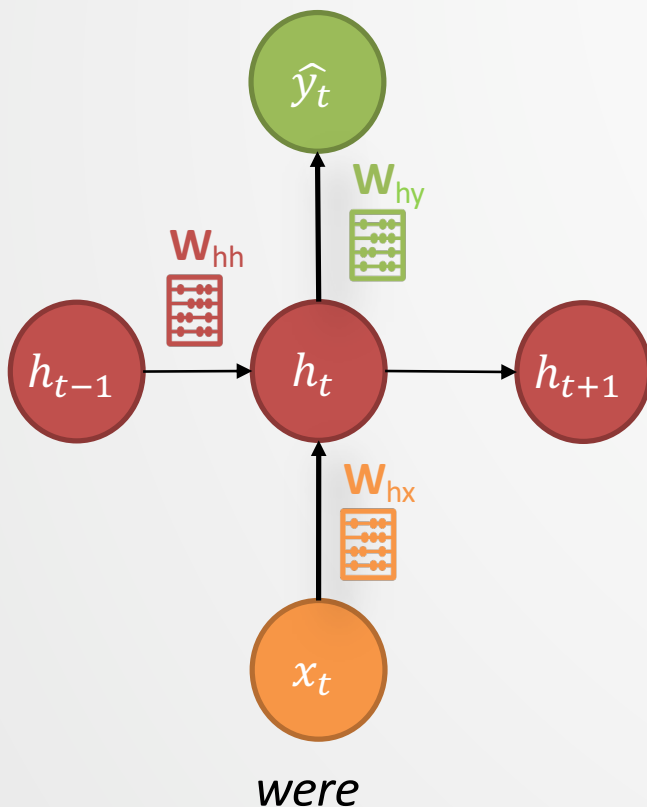
$$h_t = g(W_I[h_{t-1}; \mathbf{x}] + c)$$
$$\mathbf{y}_t = W_O \mathbf{h}_t + b$$

RNNs: One time step snapshot

Two riders .. approaching .. horses.

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

approaching



- 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_hx, x))

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

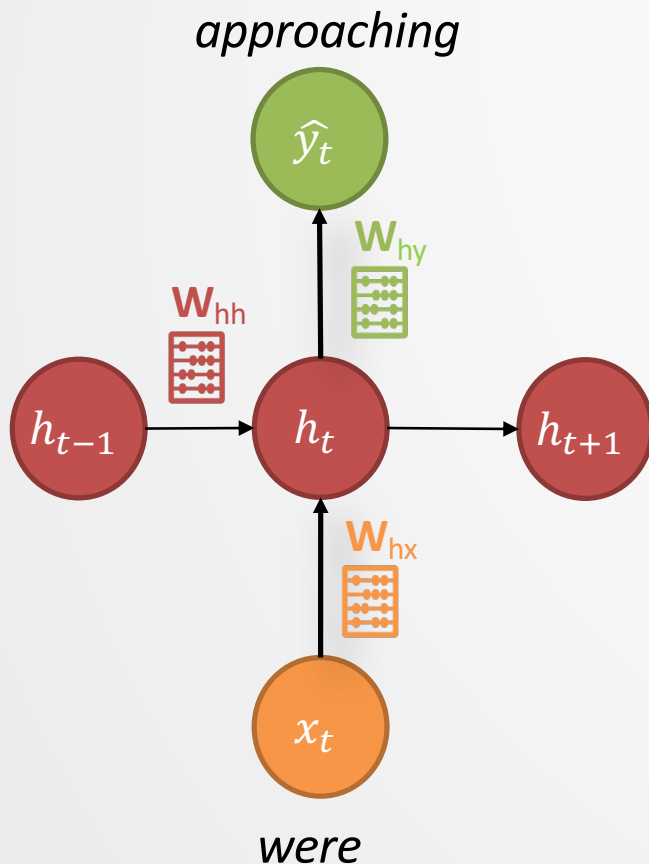
        return softmax(y) if is_normalized else y
```

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

RNNs: Training

Two riders .. approaching .. horses.

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



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

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

Evaluation

- Same **cross-entropy** loss function

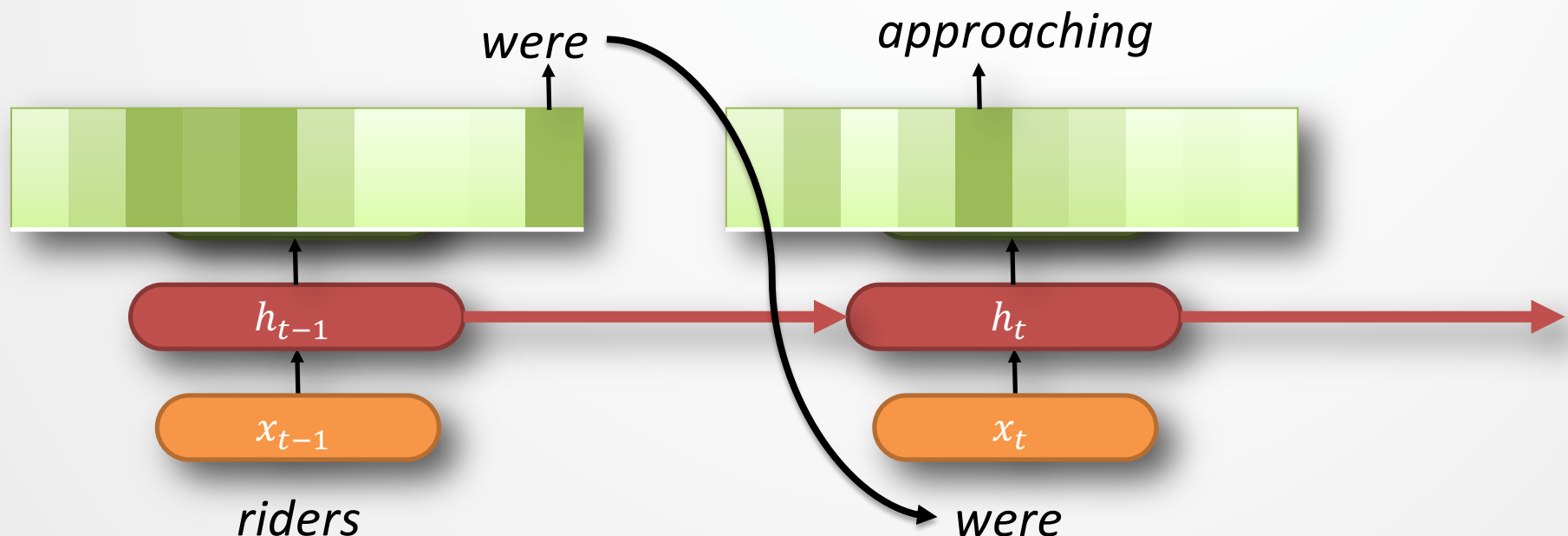
$$J^{(t)}(\theta) = - \sum_{j=1}^{|V|} \overbrace{y_{t,j}}^{\text{Ground truth}} \log \underbrace{\hat{y}_{t,j}}_{\text{prediction}}$$

- 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 from the trigram NN.
 - In 'theory', information is maintained in h_i across arbitrary lengths of time...



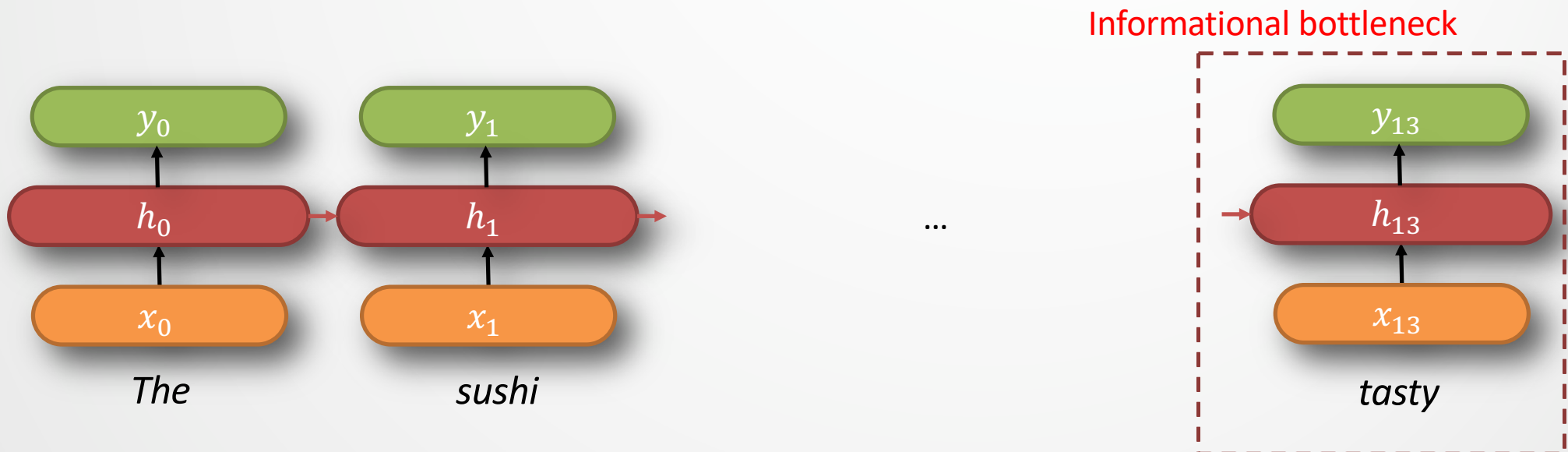
$$h_t = g([W_{hh}h_{t-1}; W_{hx}x_t] + c)$$

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

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

RNNs and retrograde amnesia

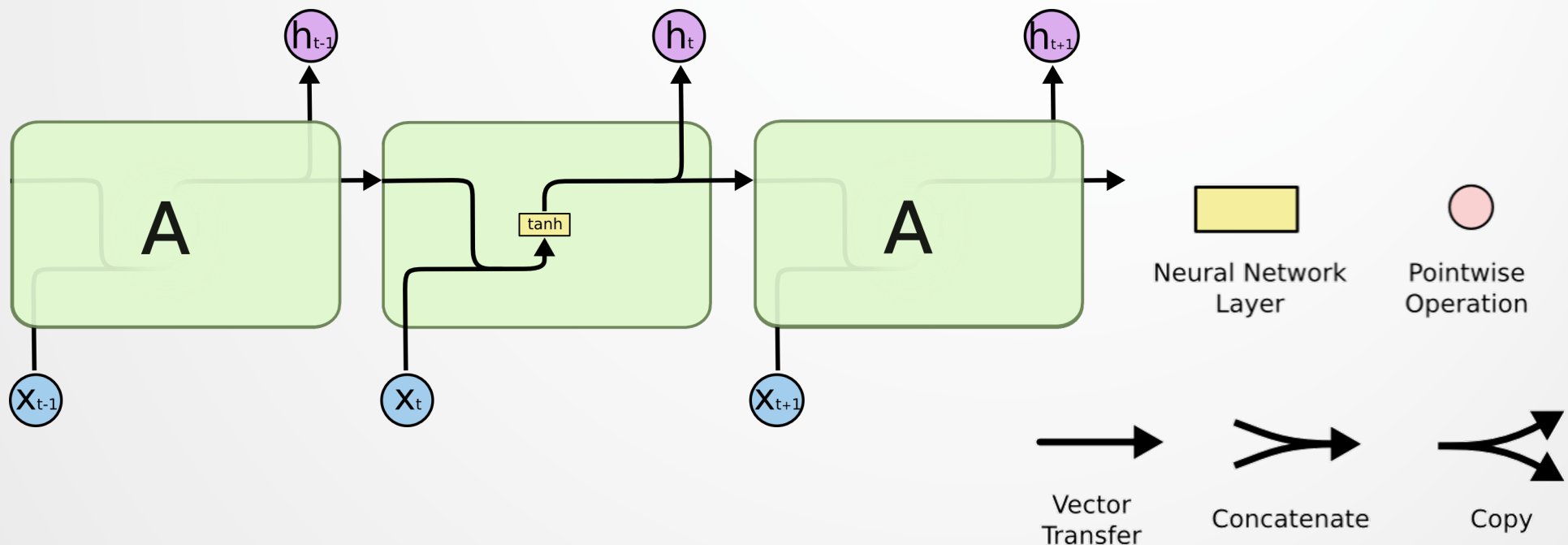
- Unfortunately, **catastrophic forgetting** is common.
 - E.g., the **relevant** context in “*The sushi the sister of your friend’s programming teacher told you about was...*” has likely been **overwritten** by the time h_{13} is produced.



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

RNNs and retrograde amnesia

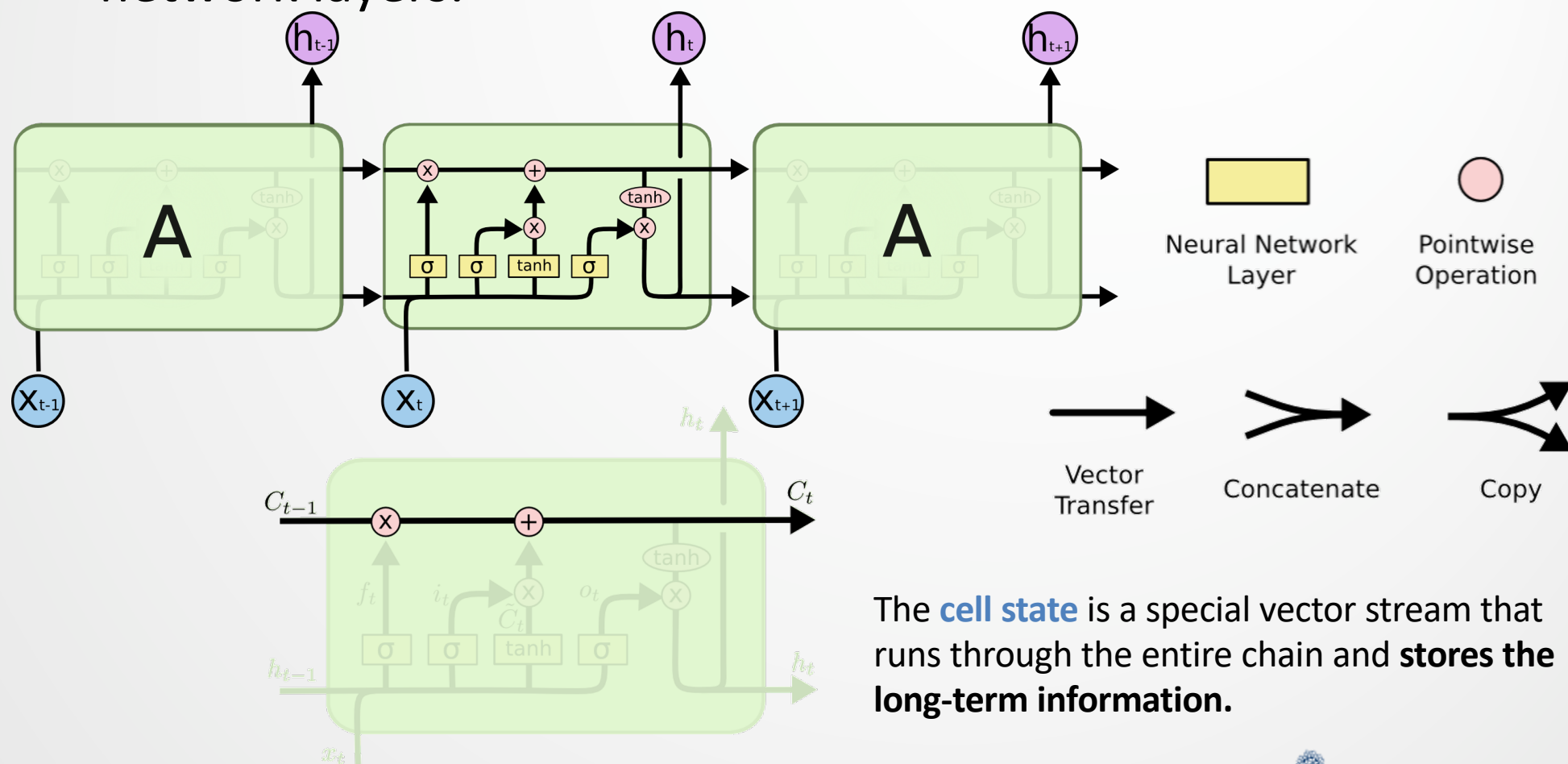
- One challenge with RNNs is that the **gradient** decays quickly as one pushes it back in time. Can we store relevant information?



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

Long short-term memory (LSTM)

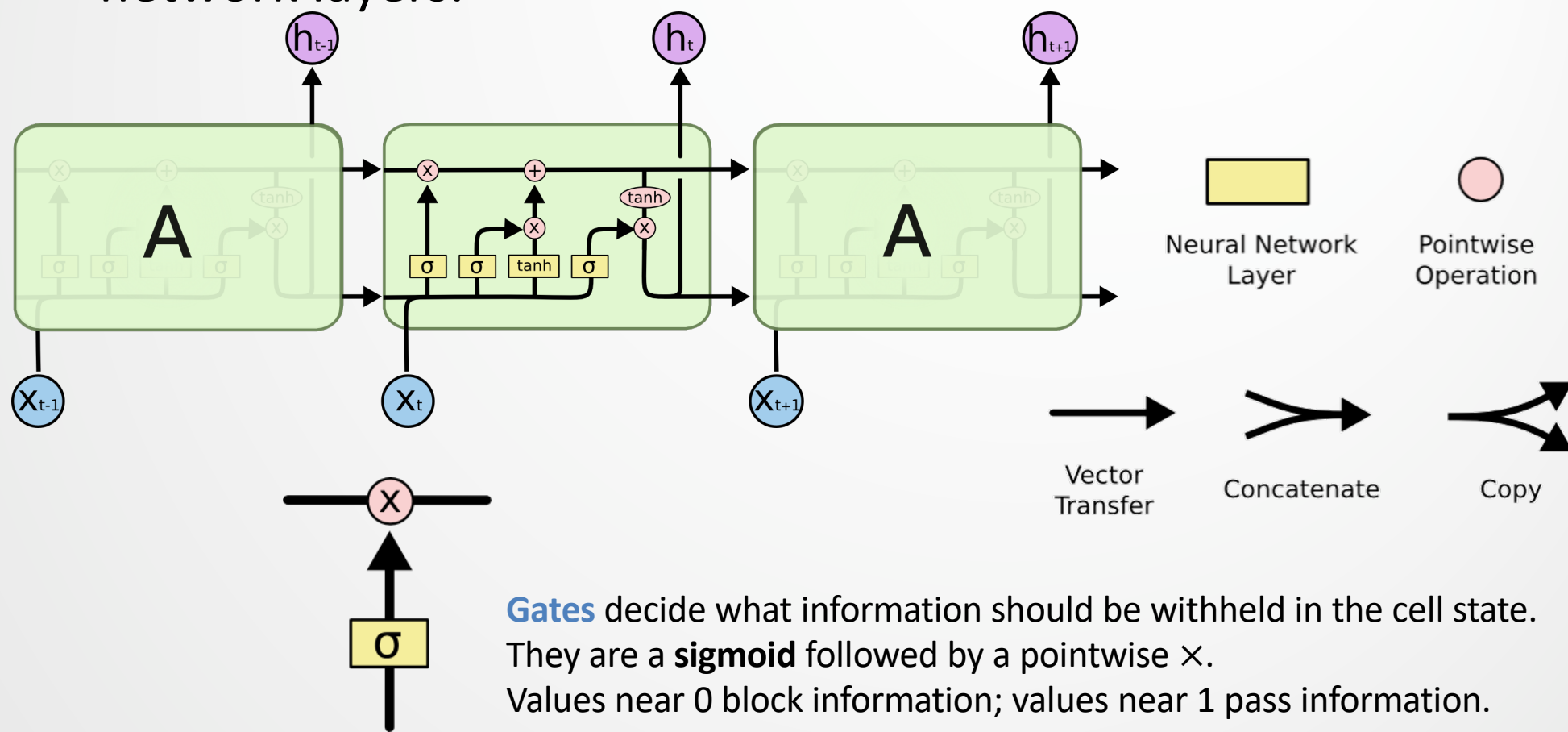
- In each **module**, in an LSTM, there are four interacting neural network layers.



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

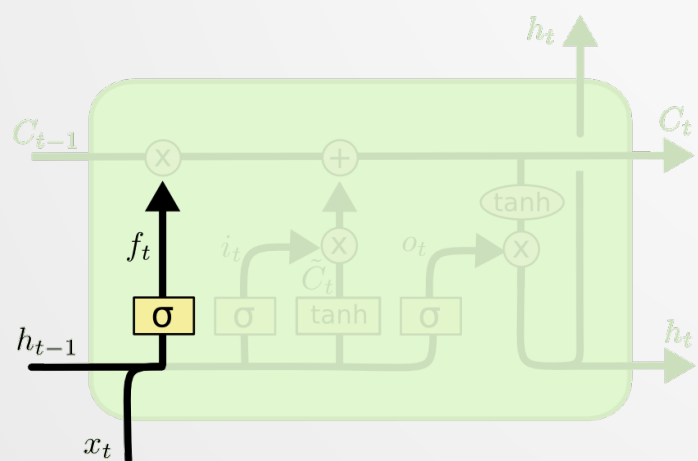
Long short-term memory (LSTM)

- In each **module**, in an LSTM, there are four interacting neural network layers.

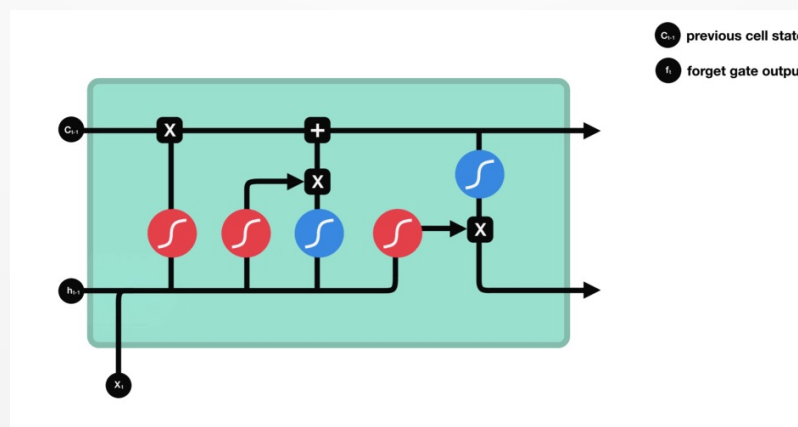


LSTM step 1: decide what to forget

- 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 explicit.)

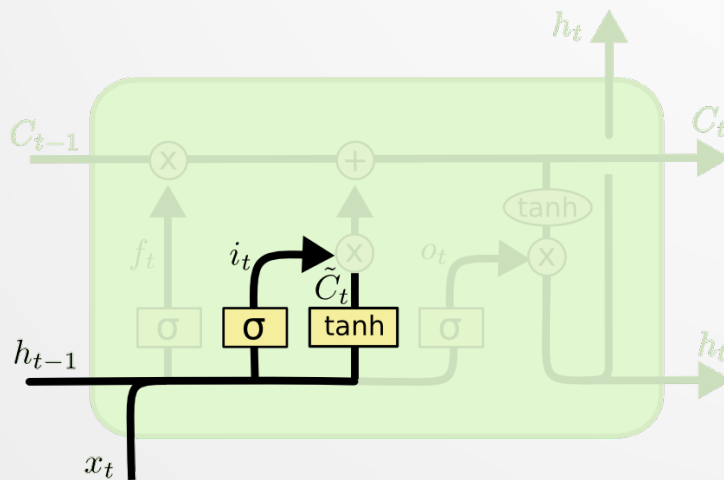


$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$



LSTM step 2: decide what to store

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

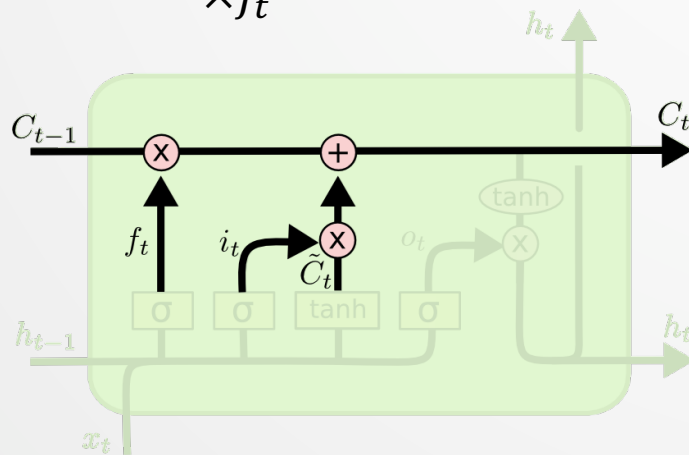
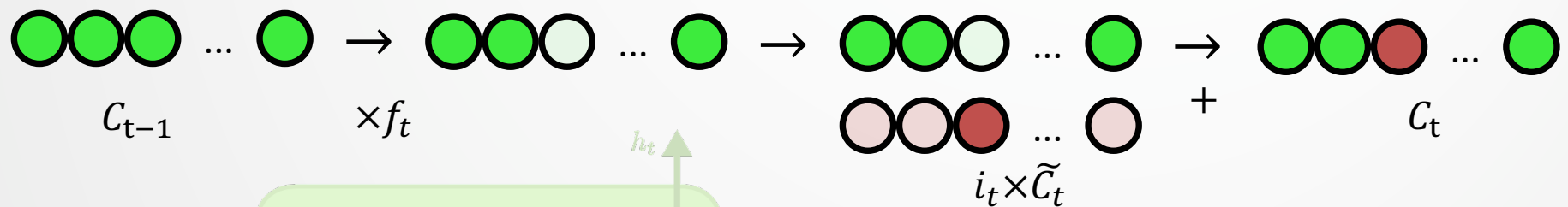


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

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

LSTM step 3: update the cell state

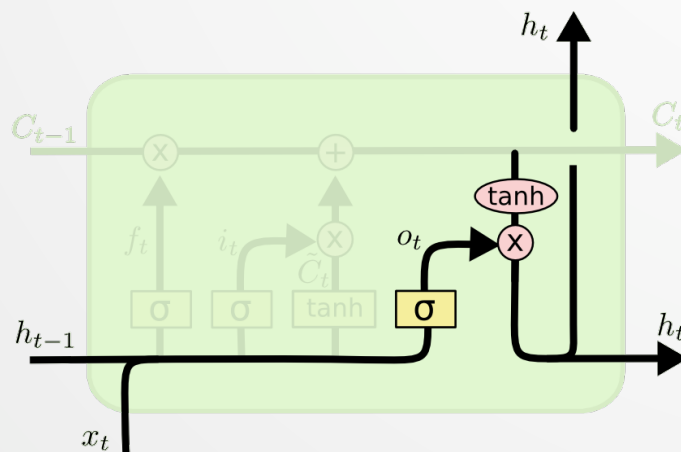
- Update C_{t-1} to C_t .
 - First, forget what we want to forget: multiply C_{t-1} by f_t .
 - Then, create a 'mask vector' of information we want to store, $i_t \times \tilde{C}_t$.
 - Finally, write this information to the new cell state C_t .



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

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).

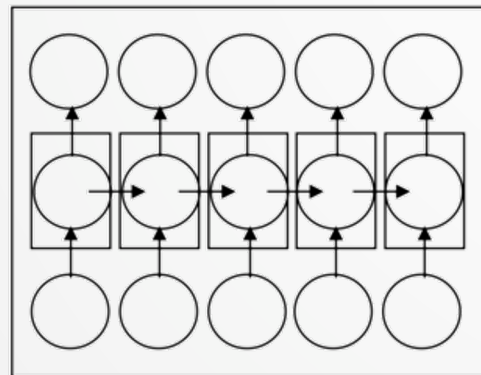


$$o_t = \sigma(W_o[h_{t-1}, x_t] + b_o)$$

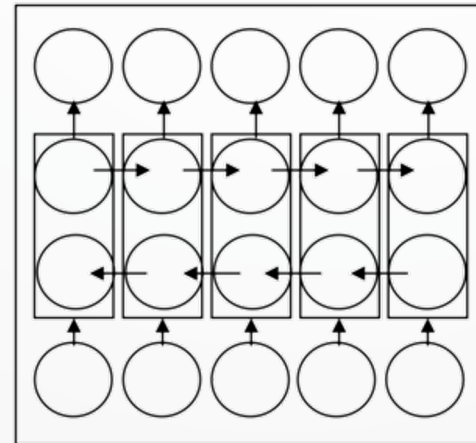
$$h_t = o_t \times \tanh(C_t)$$

Variants of LSTMs

- There are various variations on LSTMs.
 - ‘*Bidirectional LSTMs*’ (and bidirectional RNNs generally), learn



(a)



(b)

Structure overview

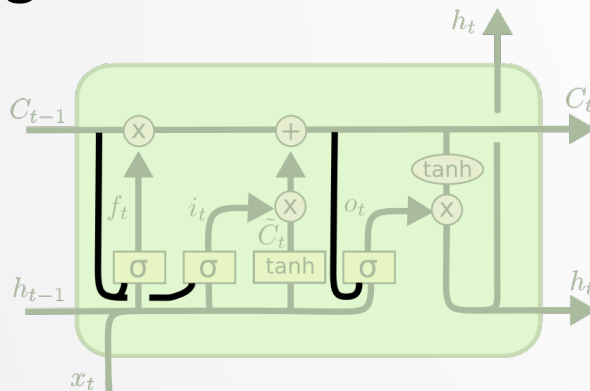
(a) unidirectional RNN

(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.

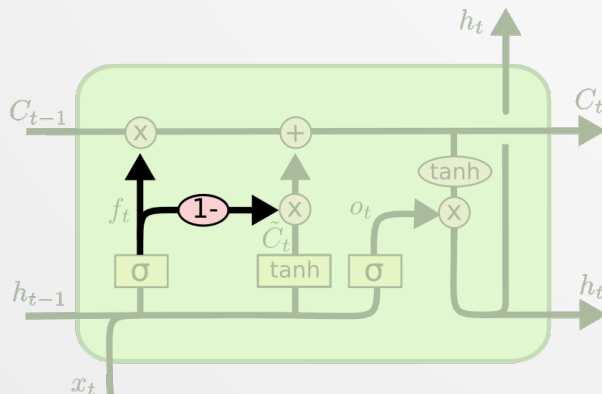


$$f_t = \sigma(W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f)$$

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

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

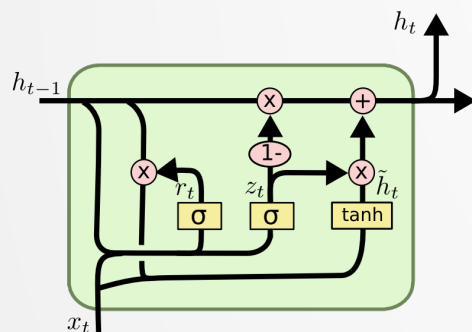
- We can **couple** the '*forget*' and '*input*' gates.
 - Joint decisioning is more efficient.



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

Aside - Variants of LSTMs

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



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

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

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

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

- 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.

RECENT_{-ISH} BREAKTHROUGHS

Deep contextualized representations

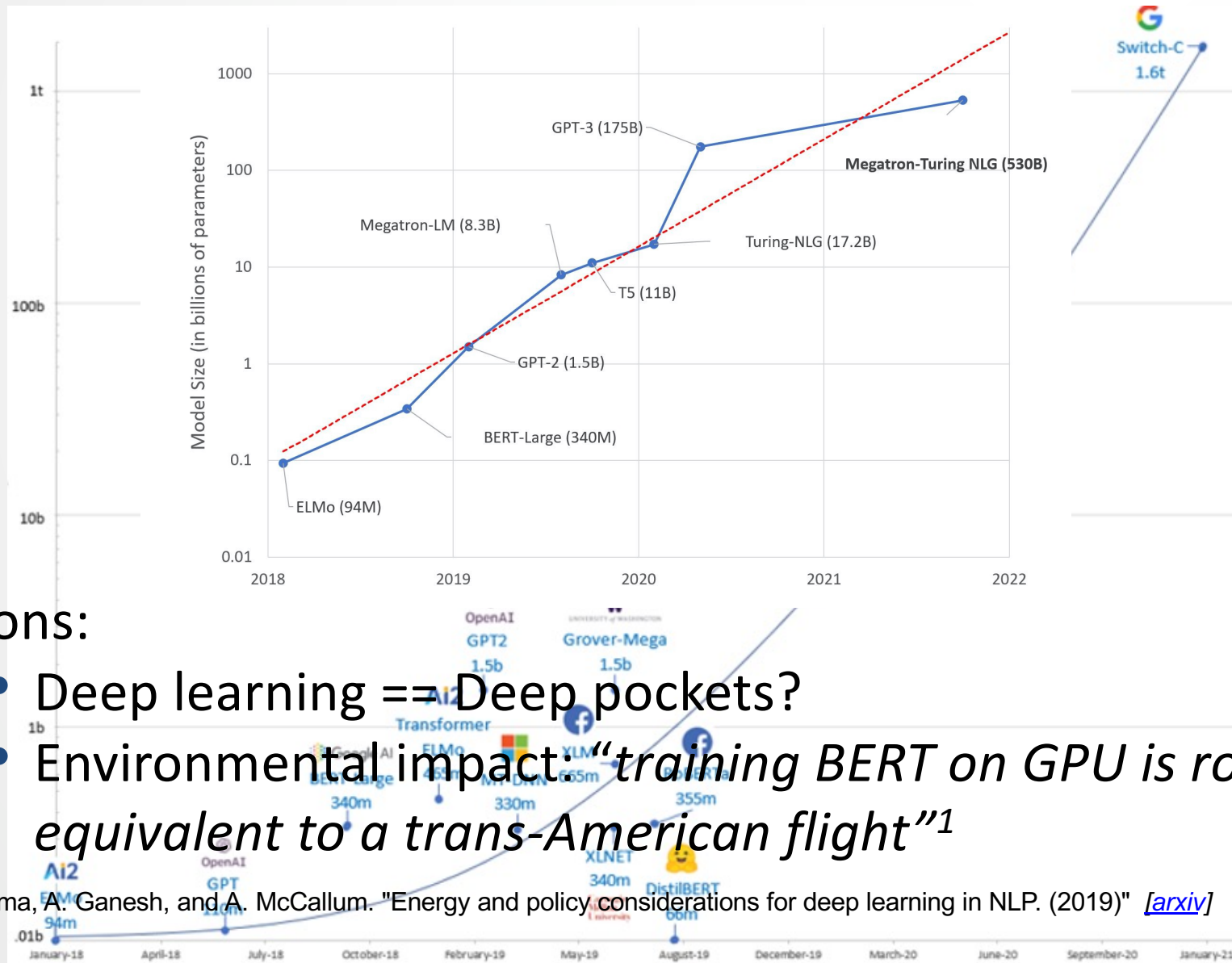
- What does the word *play* mean?



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>

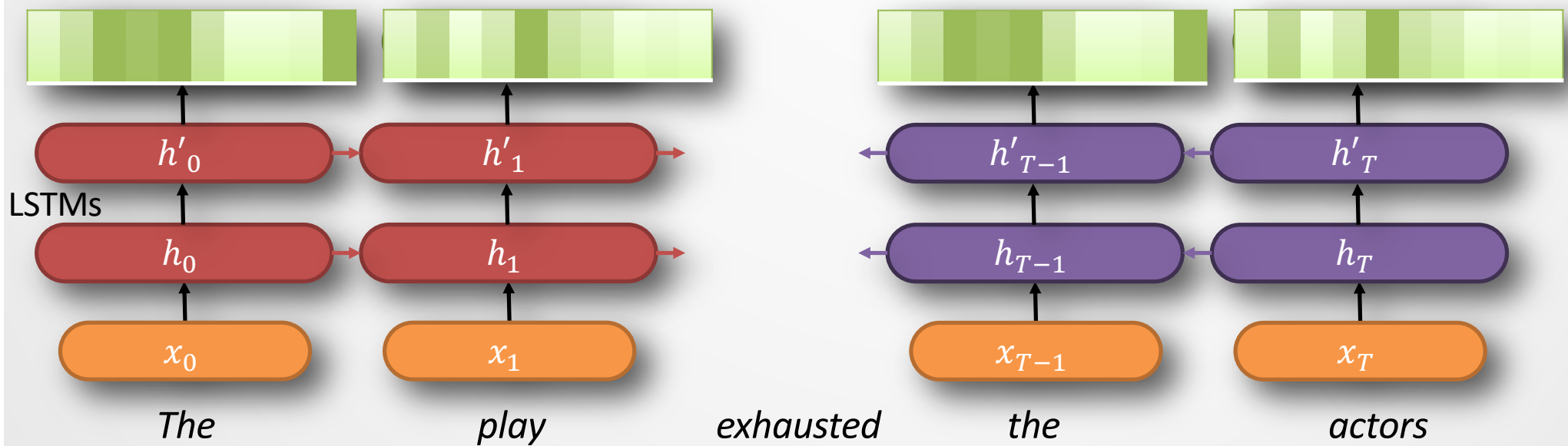
NLM: the bigger is better trend



¹ S. Emma, A. Ganesh, and A. McCallum. "Energy and policy considerations for deep learning in NLP. (2019)" [\[arxiv\]](#)

ELMo: Embeddings from Language Models

- Instead of a fixed embedding for each word **type**, ELMo considers the entire sentence before embedding each **token**.
 - It uses a **bi-directional** LSTM trained on a specific task.
 - Outputs are softmax probabilities on words, as before.



ELMo: Embeddings from Language Models

For each token, a L-layer biLM computes (2L+1) representations:

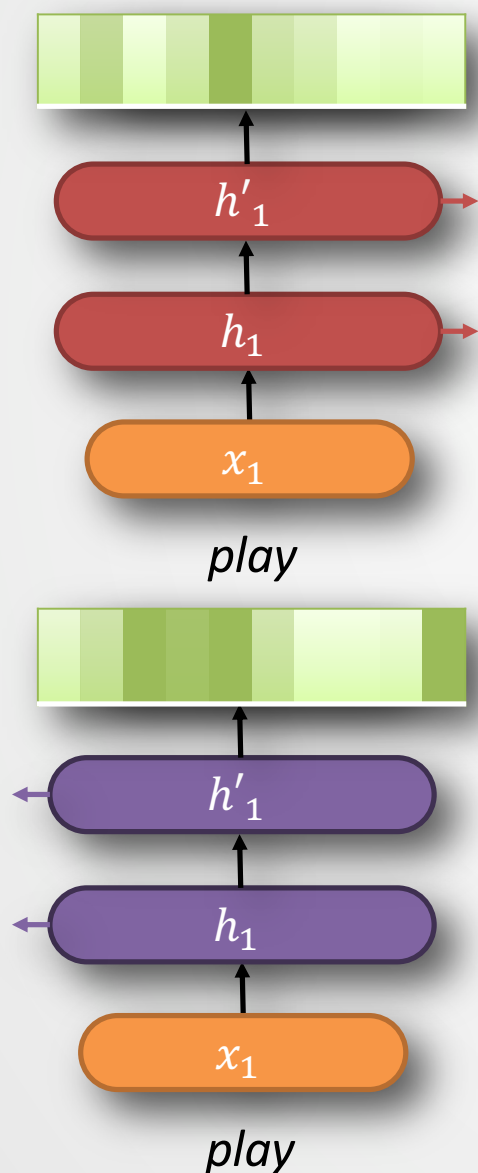
$$\begin{aligned} R_k &= \{\mathbf{x}_k^{LM}, \vec{\mathbf{h}}_{k,j}^{LM}, \overleftarrow{\mathbf{h}}_{k,j}^{LM} \mid j = 1, \dots, L\} \\ &= \{\mathbf{h}_{k,j}^{LM} \mid j = 0, \dots, L\}, \end{aligned}$$

- Task specific weighting produces the final embedding for word token k .

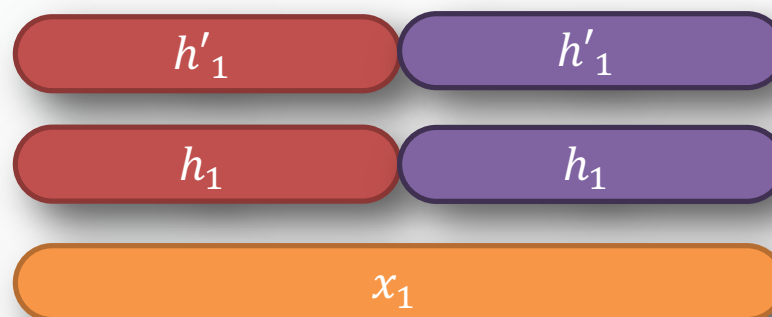
$$\mathbf{ELMo}_k^{task} = E(R_k; \Theta^{task}) = \gamma^{task} \sum_{j=0}^L s_j^{task} \mathbf{h}_{k,j}^{LM}$$

- where R_K is the set of all L hidden layers, $\mathbf{h}_{k,j}$
 s_j^{task} is the task's weight on the layer, and
 γ^{task} is a weight on the entire task

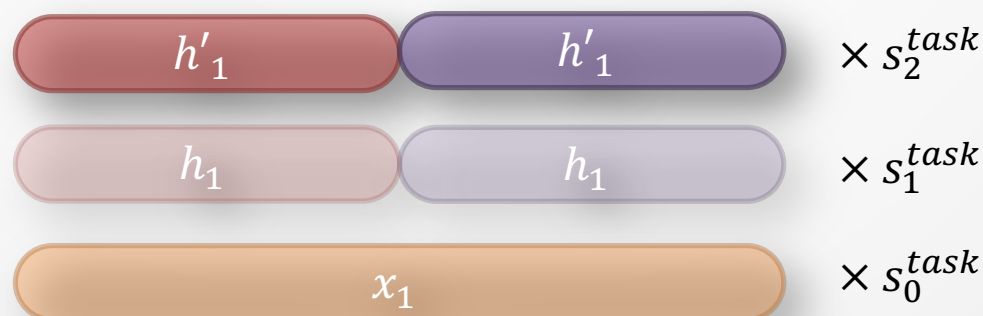
ELMo: Embeddings from Language Models



1. Concatenate



2. Multiply by weight vectors



3. Sum



ELMo: Embeddings from Language Models

- What does the word *play* mean?

Source		Nearest Neighbors
GloVe	play	playing, game, games, played, players, plays, player, Play, football, multiplayer
biLM	Chico Ruiz made a spectacular <u>play</u> on Alusik 's grounder {...}	Kieffer , the only junior in the group , was commended for his ability to hit in the clutch , as well as his all-round excellent <u>play</u> .
	Olivia De Havilland signed to do a Broadway <u>play</u> for Garson {...}	{...} they were actors who had been handed fat roles in a successful <u>play</u> , and had talent enough to fill the roles competently , with nice understatement .

Table 4: Nearest neighbors to “play” using GloVe and the context embeddings from a biLM.

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>

ELMo: Embeddings from Language Models

Q&A
Textual entailment
Semantic role labelling
Coreference resolution
Name entity resolution
Sentiment analysis

TASK	PREVIOUS SOTA		OUR BASELINE	ELMo + BASELINE	INCREASE (ABSOLUTE/ RELATIVE)
SQuAD	Liu et al. (2017)	84.4	81.1	85.8	4.7 / 24.9%
SNLI	Chen et al. (2017)	88.6	88.0	88.7 ± 0.17	0.7 / 5.8%
SRL	He et al. (2017)	81.7	81.4	84.6	3.2 / 17.2%
Coref	Lee et al. (2017)	67.2	67.2	70.4	3.2 / 9.8%
NER	Peters et al. (2017)	91.93 ± 0.19	90.15	92.22 ± 0.10	2.06 / 21%
SST-5	McCann et al. (2017)	53.7	51.4	54.7 ± 0.5	3.3 / 6.8%

Table 1: Test set comparison of ELMo enhanced neural models with state-of-the-art single model baselines across six benchmark NLP tasks. The performance metric varies across tasks – accuracy for SNLI and SST-5; F_1 for SQuAD, SRL and NER; average F_1 for Coref. Due to the small test sizes for NER and SST-5, we report the mean and standard deviation across five runs with different random seeds. The “increase” column lists both the absolute and relative improvements over our baseline.

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: Bidirectional encoder representations from transformers

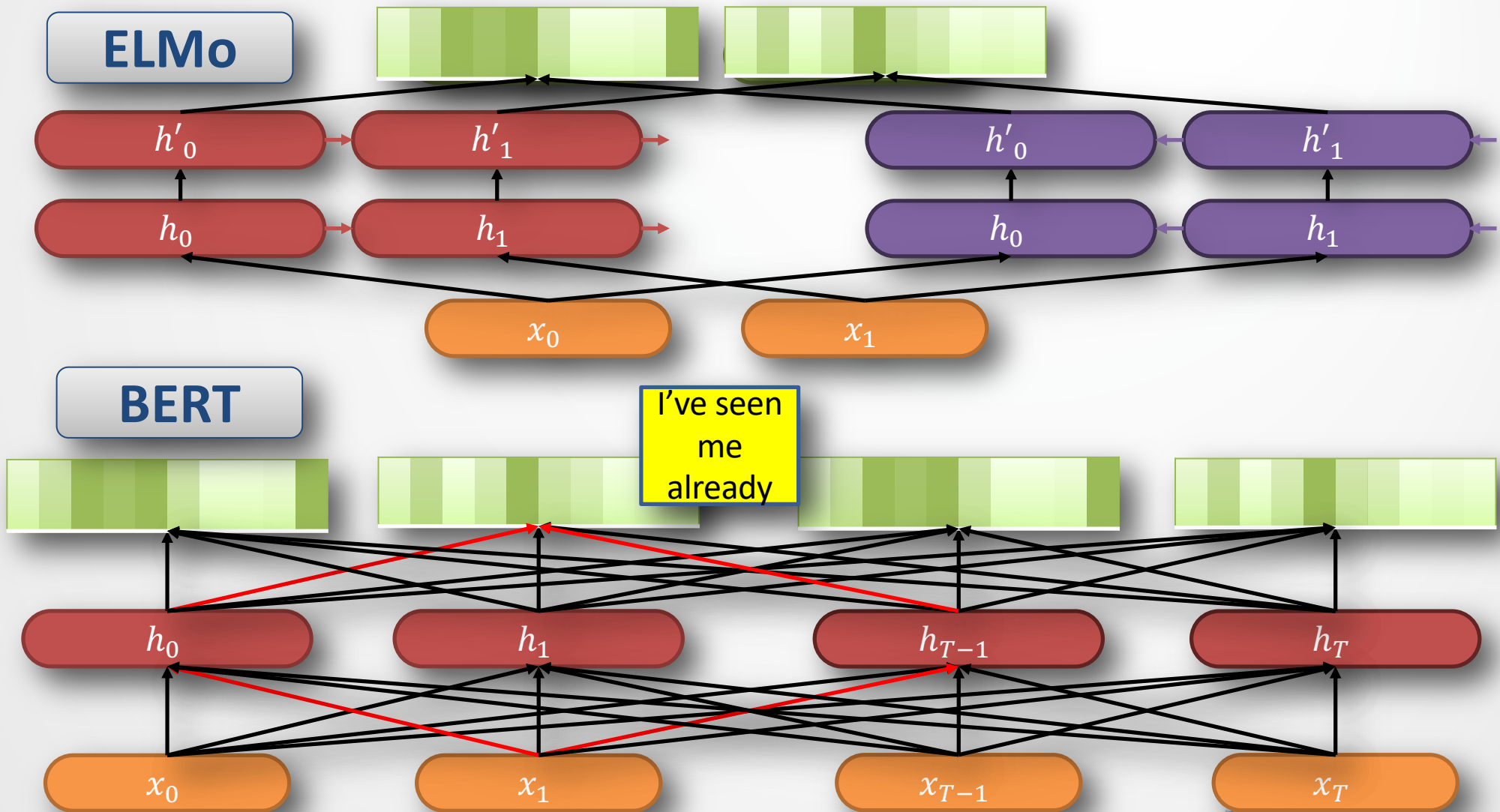
- Unlike ELMo, BERT is **deeply** bidirectional.
 - i.e., every embedding conditions every other in the next layer.
- This is difficult, because when predicting word x_t , you would already have 'seen' that word in modelling its own contexts.



Code and models: <https://github.com/google-research/bert>

Devlin J, Chang M-W, Lee K, *et al.* BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. <http://arxiv.org/abs/1810.04805>

BERT: Bidirectional encoder representations from transformers



BERT: Bidirectional encoder representations from transformers

- This can be solved by **masking** the word being predicted.

```
Input: The man went to the [MASK]1 . He bought a [MASK]2 of milk .  
Labels: [MASK]1 = store; [MASK]2 = gallon
```

- (actually, 80% we use [MASK]. 10% we replace the target word with another actual word; 10% we keep the word as-is, to bias 'towards the observation'.)
- We can also predict other relationships, like whether one sentence follows another.

```
Sentence A = The man went to the store.  
Sentence B = He bought a gallon of milk.  
Label = IsNextSentence
```


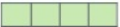






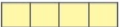


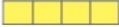

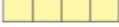

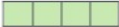


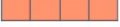


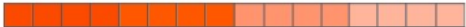
```
Sentence A = The man went to the store.  
Sentence B = Penguins are flightless.  
Label = NotNextSentence
```

- (actually, you can fine-tune on *many* different tasks)

Aroca-Ouellette S, Rudzicz F (2020) [On Losses for Modern Language Models](#), EMNLP.

BERT: Bidirectional encoder representations from transformers

What is the best contextualized embedding for “Help” in that context?
For named-entity recognition task CoNLL-2003 NER

		Dev F1 Score
12 	First Layer Embedding 	91.0
...		
7 	Last Hidden Layer 12 	94.9
6 		
5 	Sum All 12 Layers	95.5
4 	12  + ... + 2  + 1  = 	
3 	Second-to-Last Hidden Layer 11 	95.6
2 		
1 	Sum Last Four Hidden	95.9
	12  + 11  + 10  + 9  = 	
Help	Concat Last Four Hidden	96.1
	9 10 11 12 	

(From <http://jalammar.github.io/illustrated-bert/>)

BERT: Bidirectional encoder representations from transformers

- The age of humans is over?

Rank Name			Model	URL	Score	CoLA	SST-2	MRPC	STS-B	QQP	
1	T5 Team - Google	T5			89.7	70.8	97.1	91.9/89.2	92.5/92.1	74.6/90.4	
2	ALBERT-Team Google Language	ALBERT (Ensemble)			89.4	69.1	97.1	93.4/91.2	92.5/92.0	74.2/90.5	
+	3	王玮	ALICE v2 large ensemble (Alibaba DAMO NLP)			89.0	69.2	97.1	93.6/91.5	92.7/92.3	74.4/90.7
	4	Microsoft D365 AI & UMD	FreeLB-RoBERTa (ensemble)			88.8	68.0	96.8	93.1/90.8	92.4/92.2	74.8/90.3
	5	Facebook AI	RoBERTa			88.5	67.8	96.7	92.3/89.8	92.2/91.9	74.3/90.2
	6	XLNet Team	XLNet-Large (ensemble)			88.4	67.8	96.8	93.0/90.7	91.6/91.1	74.2/90.3
+	7	Microsoft D365 AI & MSR AI	MT-DNN-ensemble			87.6	68.4	96.5	92.7/90.3	91.1/90.7	73.7/89.9
	8	GLUE Human Baselines	GLUE Human Baselines			87.1	66.4	97.8	86.3/80.8	92.7/92.6	59.5/80.4
	9	Stanford Hazy Research	Snorkel MeTaL			83.2	63.8	96.2	91.5/88.5	90.1/89.7	73.1/89.9
	10	XLM Systems	XLM (English only)			83.1	62.9	95.6	90.7/87.1	88.8/88.2	73.2/89.8

BERT

Humans

Aside – ClosedAI

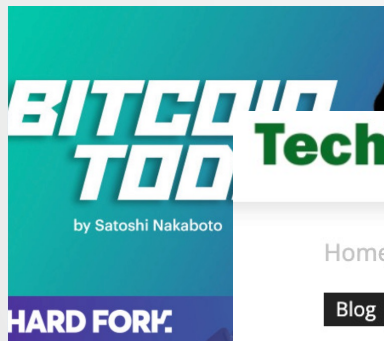
- There are, of course, alternatives.
- **FastText**: Represent each word as a bag of character-grams
Paper: <https://arxiv.org/abs/1607.04606>
Code: <https://fasttext.cc>
- **ULMFit**: Model fine-tuning for classification tasks
Paper: <https://arxiv.org/abs/1801.06146>
Code: [Here](#)
- **GPT-2/3**: Spooky, closed uni-directional model
Paper: [Here](#)
Blog: [Here](#)

A college student posts and ends

He says he wanted to prove the

By Kim Lyons | Aug 16, 2020, 1:55pm EDT

f t SHARE



thenextweb.com

Satoshi Nakaboto: 'biggest thing since'

Our robot colleague Satoshi Nakaboto: Welcome to another edition of what's been going on with Bitcoin

[Read more](#)

Home

Blog

The

the transformation of OpenAI

By Ben Dickson - August 17, 2020

GPT-3 CREATIVE FICTION

Creative writing by OpenAI's GPT-3 model, demonstrating poetry, dialogue, puns, literary parodies, and storytelling.

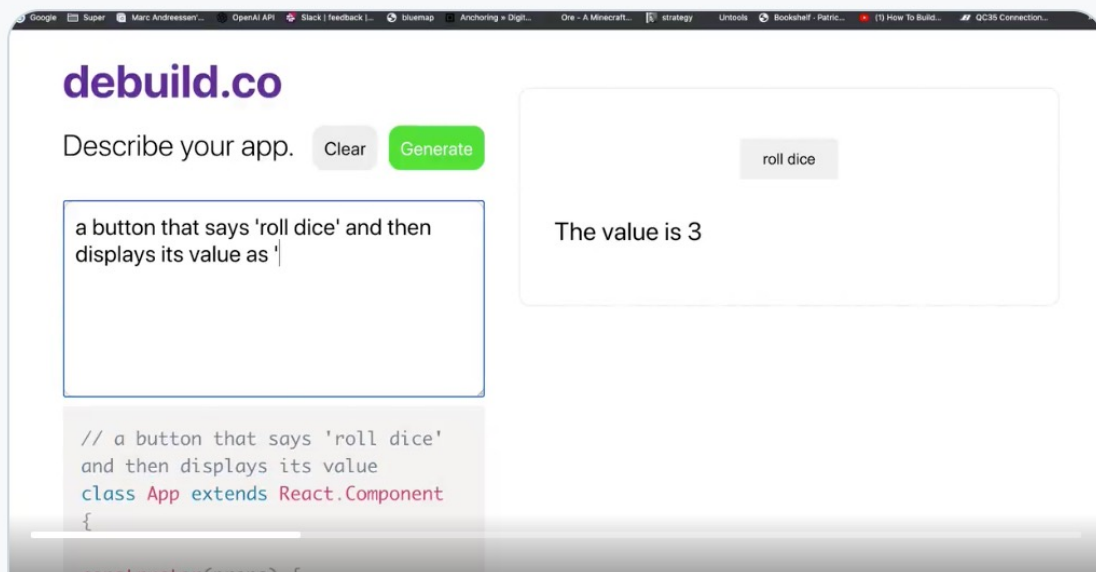


Sharif Shameem @sharifshameem · Jul 19

Wow.

I built a React dice component with GPT-3.

This feels far more fun than writing JSX.

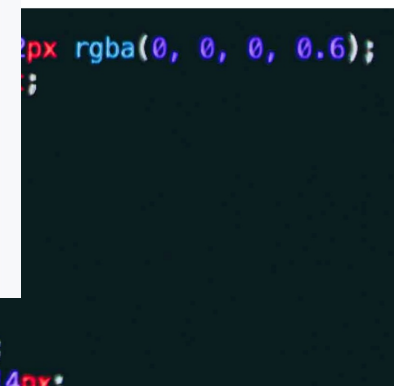


Language
is
—and
less

created and can generate
can't bring us closer to

Can Now Generate The
e For You

20/07/2020



Google GPT 3 + Medium

About 2,030,000 results (0.37 seconds)

medium.com · crazy-gpt-3-use-cases-232c22142044
Crazy GPT-3 Use Cases. Discover how powerful ... - Medium
code (websites, machine learning models, design); completed sentences; layout; simple reasoning. GPT-3 ...
Jul 27, 2020 · Uploaded by Przemek Chojek

onzero.medium.com · gpt-3-is-an-amazing-research-...
GPT-3 is an Amazing Research Tool. But OpenAI Isn't Sharing ...
18 hours ago · At its heart, GPT-3 is an incredibly powerful tool for writing in the English language. ... GPT-3 can be accessed only through an API run by OpenAI, similar ... Create a free Medium account to get Pattern Matching in your inbox.

medium.com · gpt-3-conversational-ai-chatbots-3fb1cf...
GPT-3: Conversational AI & Chatbots | by Cobus ... - Medium
Jul 26, 2020 · There has been much talk and hype regarding GPT-3, the last couple of days. Even though I have not built any prototypes with GPT yet, ...

medium.com · what-gpt-3-means-for-non-techs-...
What GPT-3 Means for Non-Technical Professionals - Medium
Jul 23, 2020 · Earlier this week I was lucky to get early access to a beta version of OpenAI's latest generative pre-trained transformer model (GPT-3) — a new ...

medium.com · openai-unveils-175-billion-parameter-g...
OpenAI Unveils 175 Billion Parameter GPT-3 ... - Medium
May 29, 2020 · The researchers show through GPT-3 training that scaling up language models can greatly improve task-agnostic, few-shot performance. ...

medium.com · crypto-oracle · openai-gpt-3-is-the-fut...
OpenAI's GPT-3 Is The Future We've Been Waiting For - Medium
Jul 19, 2020 · What we do with it, for good or evil, is only limited by our imaginations. Generative Pre-Trained Transformer (GPT) — 3. The original paper on ...

Google GPT3

About 1,340,000 results (0.56 seconds)

Did you mean: **GPT 3**

openai.com · blog · openai-api
OpenAI API
Jun 11, 2020 · Today the API runs models with weights from the GPT-3 family with many speed and throughput improvements. Machine learning is moving ...

www.theverge.com · gpt-3-explainer-openai-examples-...
OpenAI's latest breakthrough is astonishingly powerful, but still ...
Jul 30, 2020 · As the name suggests, GPT-3 is the third in a series of autocomplete tools designed ... I made a fully functioning search engine on top of GPT3.

Google GPT3 + Medium

About 400,000 results (0.50 seconds)

Did you mean: **GPT 3 + Medium**

medium.com · crypto-oracle · openai-gpt-3-is-the-fut...
OpenAI's GPT-3 Is The Future We've Been Waiting For - Medium
Jul 19, 2020 · OpenAI is an artificial intelligence research laboratory started in 2015 by Elon Musk, Sam Altman, Peter Thiel, Reid Hoffman, Marc Benioff and ...

medium.com · crazy-gpt-3-use-cases-232c22142044
Crazy GPT-3 Use Cases. Discover how powerful ... - Medium
Jul 27, 2020 · code (websites, machine learning models, design); completed sentences; layout; simple reasoning. GPT-3 definitely will influence how we ...

medium.com · gpt-3-conversational-ai-chatbots-3fb1cf...
GPT-3: Conversational AI & Chatbots | by Cobus ... - Medium
Jul 26, 2020 · The fact that OpenAI is in the process of releasing an API will have a significant impact on the Conversational AI marketplace. OpenAI states ...

medium.com · openai-unveils-175-billion-parameter-g...
OpenAI Unveils 175 Billion Parameter GPT-3 ... - Medium
May 29, 2020 · When it comes to large language models, it turns out that even 1.5 billion parameters is not large enough. While that was the size of the GPT-2 ...

YouTube GPT 3

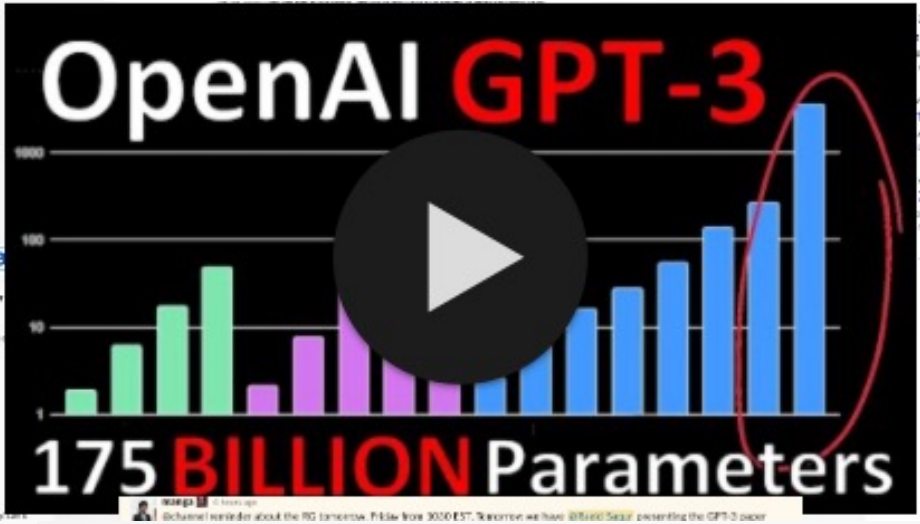
GPT 3 Demo and Explanation - An AI revolution from OpenAI
Half Ideas · Startups and Entrepreneurship · 97K views · 1 month ago
GPT 3 can write poetry, translate text, chat convincingly, and answer abstract questions. It's being used to code, design and much ...

OpenAI GPT-3 - Good At Almost Everything!
Two Minute Papers · 205K views · 1 week ago
Check out Weights & Biases and sign up for a free demo here: <https://www.wandb.com/papers> Their instrumentation of a ...

OpenAI GPT-3
175 BILLION Param · 1:04:30
Yannic Kilcher · 78K views · 2 months ago
gpt3 #openai #gpt3 How far can you go with ONLY language modeling? Can a large enough language model perform NLP task ...

Thoughts about GPT-3 | Dileep George and Lex Fridman
Lex Clips · 10K views · 4 days ago
Full episode with Dileep George (Aug 2020): https://www.youtube.com/watch?v=tg_mLuoRWdM Clips channel (Lex Clips): ...

GPT-3: An Even Bigger Language Model - Computerphile
Computerphile · 148K views · 1 month ago
Basic mathematics from a language model? Rob Miles on GPT3, where it seems like size does matter! More from Rob Miles: ...



manga · 4 hours ago
@channel reminder about the RG tomorrow, Friday from 1030 EST. Tomorrow we have @Raeid Saqur presenting the GPT-3 paper <https://arxiv.org/abs/2005.14165>

arXiv.org
Language Models are Few-Shot Learners
Recent work has demonstrated substantial gains on many NLP tasks and benchmarks by pre-training on a large corpus of text followed by fine-tuning on a specific task. While typically task-agnostic...

Raeid Saqur · 10 minutes ago
sigh ... now I know how it would feel for an underground grunge metal band to cover Nickelback. 😞 It's soo'...' mainstream/pop-culture-esque ... Shia LaBeouf would +1 me easy 😊

Google GPT3 tutorial

About 113,000 results (0.46 seconds)

Did you mean: **GPT 3 tutorial**

Videos

OpenAI GPT-3: Beginners Tutorial
Bhavesh Bhatt · YouTube · Jul 23, 2020

GPT 3 Tutorial | GPT 3 Explained | What is GPT 3(Generative ...
Simplelearn · YouTube · Aug 10, 2020

OpenAI's GPT-3 Model Generates Code!
Bhavesh Bhatt · YouTube · Jul 30, 2020

www.twilio.com · blog · openai-gpt-3-chatbot-python-...
Building a Chatbot with OpenAI's GPT-3 engine, Twilio SMS ...
Aug 3, 2020 · In this tutorial I'm going to show you how easy it is to build a chatbot for Twilio SMS using the OpenAI platform and the Flask framework for ...

mariodian.com · blog · how-to-get-early-access-to-gpt-...
How to get an early access to GPT-3 and how to talk to it
Jul 23, 2020 · In the meantime, follow this short tutorial on how to get around it and start playing with it now. AI Dungeon. AI Dungeon is a free text based game ...

openai.com · blog · openai-api
OpenAI API
Jun 11, 2020 · Today the API runs models with weights from the GPT-3 family with many speed and throughput improvements. Machine learning is moving ...

The Open AI GPT Papers

- The GPT papers:
 - GPT (2018)
 - GPT2 (2019)
 - GPT3 (2020)
- Each builds on the predecessor

Improving Language Understanding by Generative Pre-Training

Alec Radford OpenAI alec@openai.com	Karthik Narasimhan OpenAI karthikn@openai.com	Tim Salimans OpenAI tim@openai.com	Ilya Sutskever OpenAI ilyasu@openai.com
---	---	--	---

Language Models are Unsupervised Multitask Learners

Alec Radford *¹ Jeffrey Wu *¹ Rewon Child¹ David Luan¹ Dario Amodei **¹ Ilya Sutskever **¹

Language Models are Few-Shot Learners

Tom B. Brown*	Benjamin Mann*	Nick Ryder*	Melanie Subbiah*	
Jared Kaplan†	Prafulla Dhariwal	Arvind Neelakantan	Pranav Shyam	Girish Sastry
Amanda Askell	Sandhini Agarwal	Ariel Herbert-Voss	Gretchen Krueger	Tom Henighan
Rewon Child	Aditya Ramesh	Daniel M. Ziegler	Jeffrey Wu	Clemens Winter
Christopher Hesse	Mark Chen	Eric Sigler	Mateusz Litwin	Scott Gray
Benjamin Chess	Jack Clark	Christopher Berner		
Sam McCandlish	Alec Radford	Ilya Sutskever	Dario Amodei	

OpenAI

Approach: Model & Architectures

The three settings we explore for in-context learning

Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 cheese => ..... ← prompt
```

One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 sea otter => loutre de mer ← example
3 cheese => ..... ← prompt
```

Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 sea otter => loutre de mer ← examples
3 peppermint => menthe poivrée
4 plush girafe => girafe peluche
5 cheese => ..... ← prompt
```

Traditional fine-tuning (not used for GPT-3)

Fine-tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.

```
1 sea otter => loutre de mer ← example #1
↓ gradient update
1 peppermint => menthe poivrée ← example #2
↓ gradient update
...
1 plush giraffe => girafe peluche ← example #N
↓ gradient update
1 cheese => ..... ← prompt
```

Figure 2.1: Zero-shot, one-shot and few-shot, contrasted with traditional fine-tuning. The panels above show

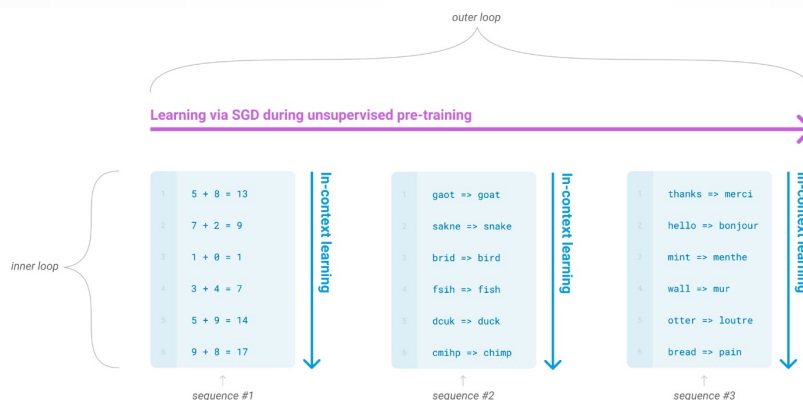


Figure 1.1: Language model meta-learning. During unsupervised pre-training, a language model develops a broad set of skills and pattern recognition abilities. It then uses these abilities at inference time to rapidly adapt to or recognize the desired task. We use the term “in-context learning” to describe the inner loop of this process, which occurs within

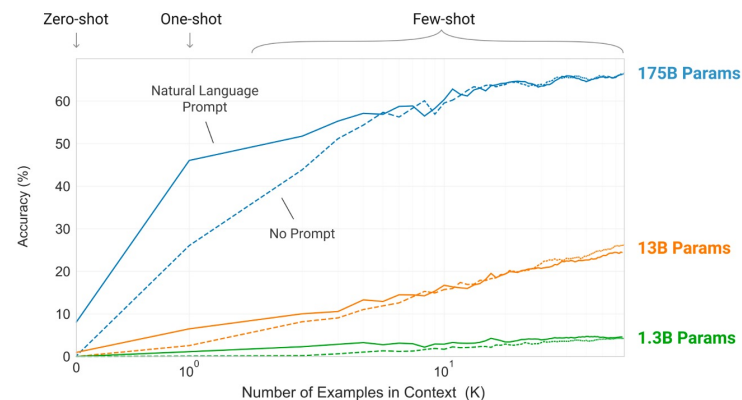


Figure 1.2: Larger models make increasingly efficient use of in-context information. We show in-context learning

Approach: Model & Architectures – GPT

- Architecture evolution: $\text{GPT3} \leftarrow \text{GPT2+mods} \leftarrow \text{GPT+mods}$
- Core architecture is classic ‘language modeling’: $p(x) = \prod_{i=1}^n p(s_i | s_1, \dots, s_{i-1})$
- Learning to perform a task as estimating distribution $P(\text{output} | \text{input})$
- Original GPT¹ trains a standard LM objective to maximize the likelihood:

$$L_1(\mathcal{U}) = \sum_i \log P(u_i | u_{i-k}, \dots, u_{i-1}; \Theta)$$

- Given an unsupervised corpus of tokens $\mu = \{\mu_1, \dots, \mu_n\}$, where k is context window, P is modelled using a neural network with parameters θ
- GPT uses a multi-layer **Transformer** decoder for the language model

[1] Radford, Alec, et al. "Improving language understanding by generative pre-training." (2018): 12.

Aside: GPT Architecture – Transformer

- Also used in (w/ caveats) in current SOTA language modeling and NLP architectures like BERT and BERT-variants (RoBERTa, XLNet, Transformer XL etc.)
- GPT vs. BERT-variants:
 - GPT uses ‘transformer’ blocks as decoders, and BERT as encoders.
 - Underlying (block level) ideology is same
 - GPT (later Transformer XL, XLNet) is an **autoregressive** model, BERT is not
 - At the cost of auto-regression, BERT has bi-directional context awareness
 - GPT, like traditional LMs, outputs one token at a time

[1] Radford, Alec, et al. "Improving language understanding by generative pre-training." (2018): 12.

Neural networks

- Research in neural networks is exciting, expansive, and explorative.
- We have many **hyperparameters** we can tweak (e.g., activation functions, number and size of layers).
- We have many **architectures** we can use (e.g., deep networks, LSTMs, attention mechanisms).
- Given the fevered hype, it's important to retain our scientific skepticism.
 - What are our **biases** and expectations?
 - When are neural networks the **wrong choice**?
 - How are we actually **evaluating** these systems?

