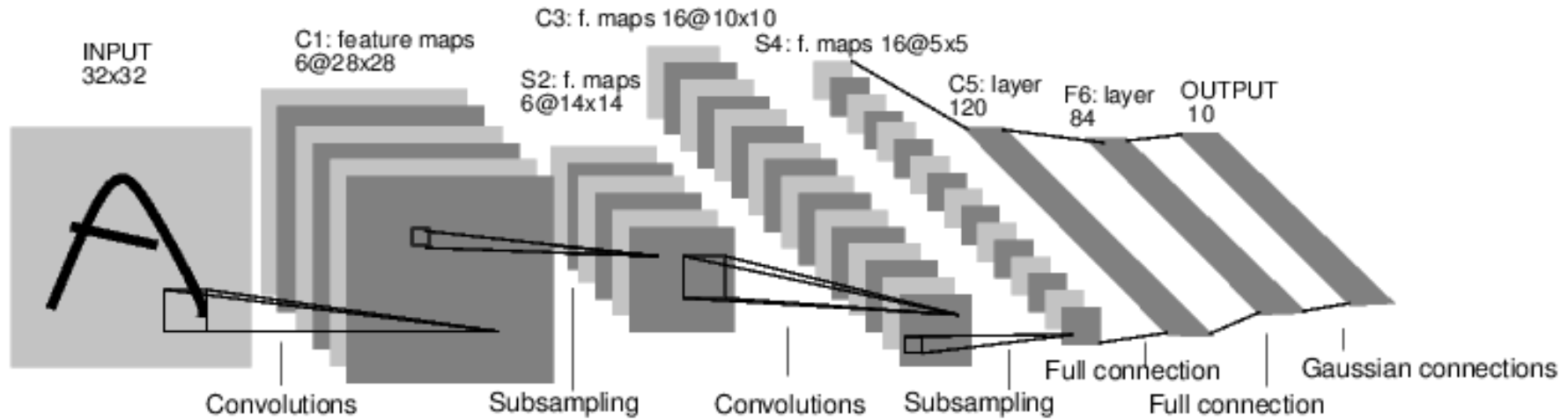


# Deep learning

LeNet-5 (1998)

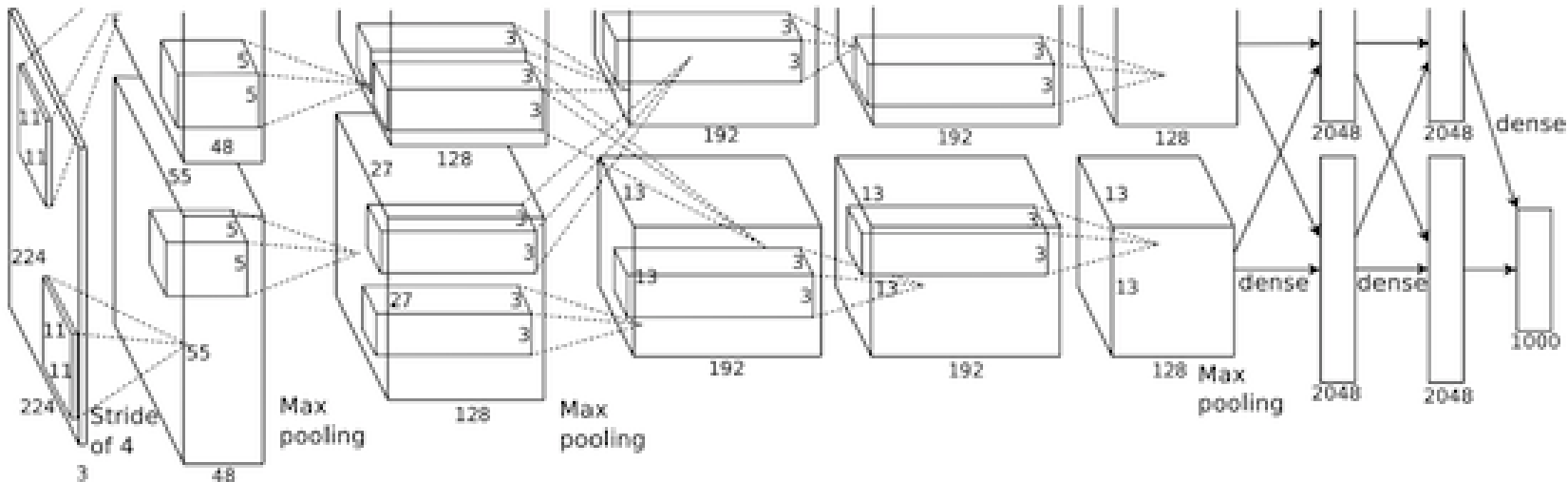


$(5 \times 5) \times 6$

+

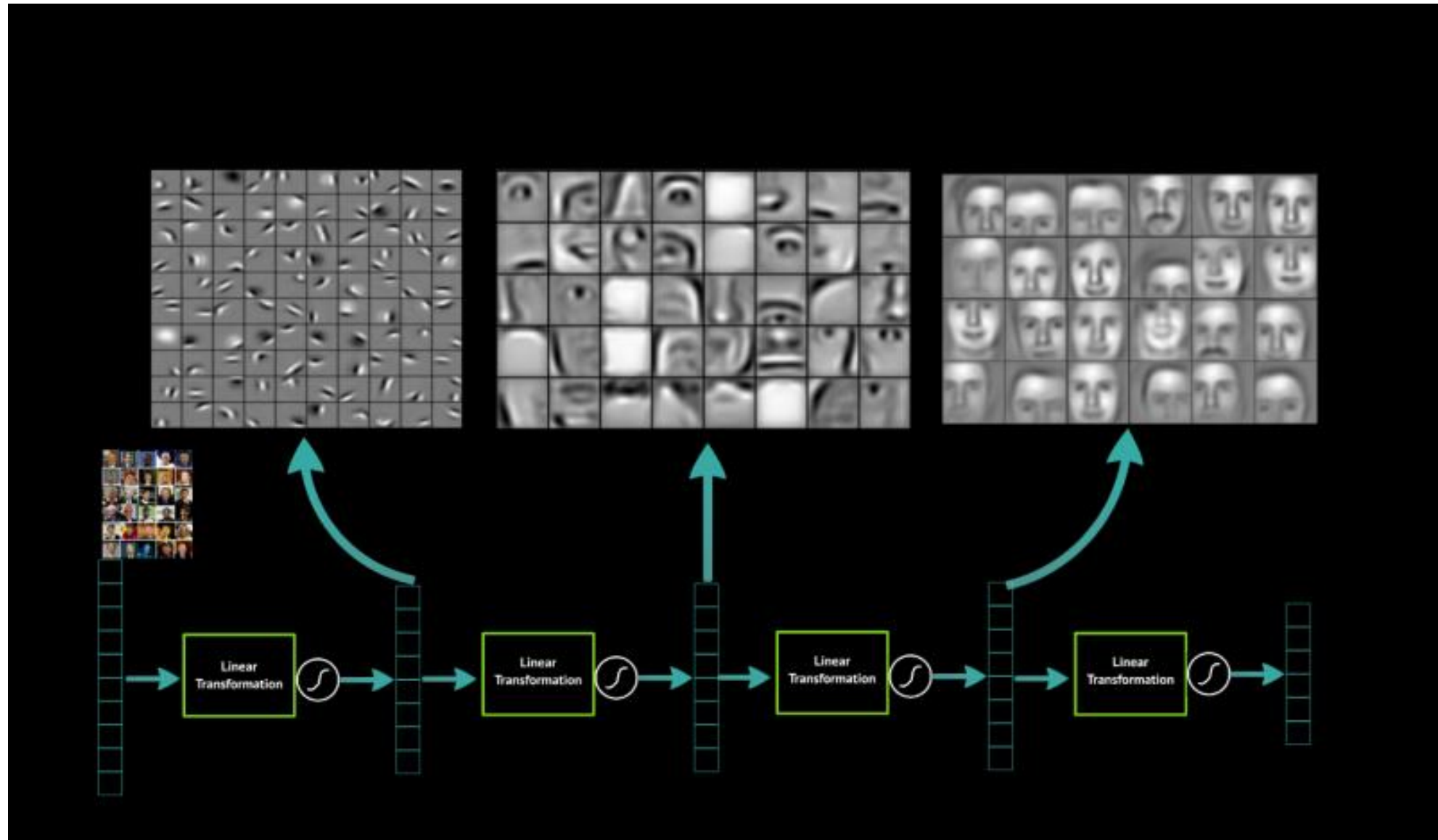
$(5 \times 5 \times 6) \times 16 + (16 \times 5 \times 5) \times 120 + 120 \times 84 + 84 \times 10 = 61470$

AlexNet (2012)



$(11 \times 11 \times 3) \times 48 + (5 \times 5 \times 48) \times 128 + (3 \times 3 \times 128) \times 192 \times 2 + (3 \times 3 \times 192) \times 192 + (3 \times 3 \times 192) \times 128 + (13 \times 13 \times 128) \times 2048 \times 2 + 2048 \times 2048 \times 2 + 2048 \times 1000 = 100\,207\,632$

# Convolution layers

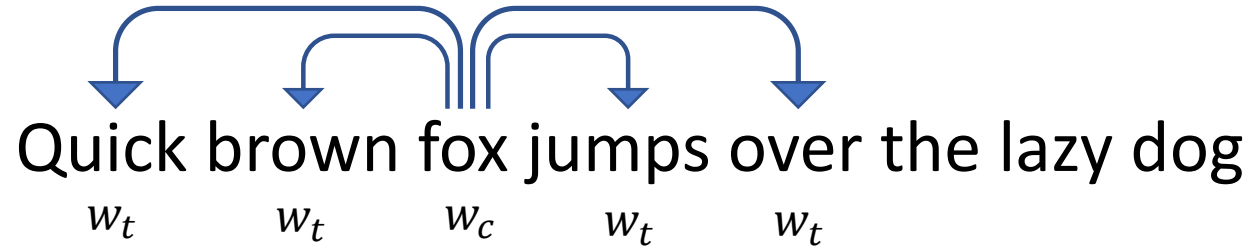


# Deep Learning for NLP

## (Natural Language Processing)

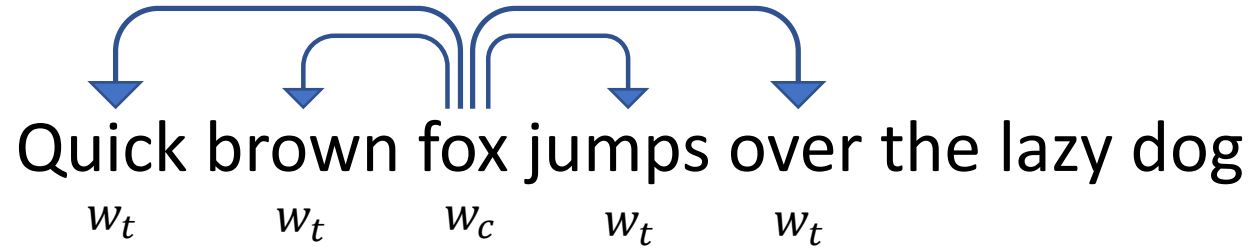
# word2vec (Google 2013)

The Idea : Let's maximize cooccurrence probability for cooccurring words (don't care about order – bag-of-words).



# word2vec (Google 2013)

The Idea : Let's maximize cooccurrence probability for cooccurring words (don't care about order – bag-of-words).

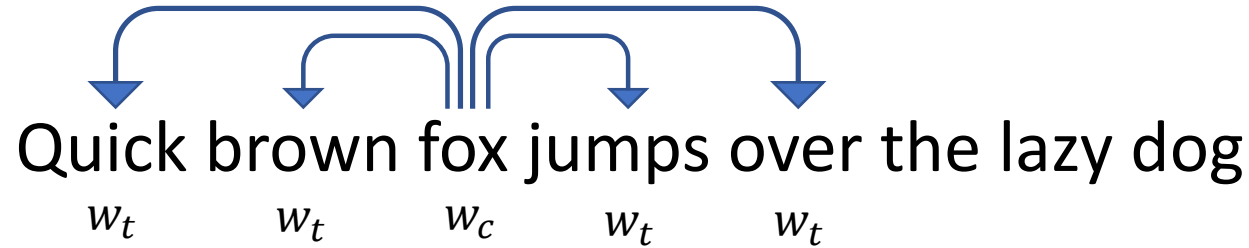


$$P(w_t|w_c) = \frac{\exp(f(w_t, w_c))}{\sum_{w_i \in Dict} \exp(f(w_i, w_c))}$$

$$\text{Loss function } J = \frac{1}{T} \sum J_t, \text{ where } J_t = -\log P(w_t|w_c) = -f(w_t, w_c) + \log \left( \sum_{w_i \in Dict} \exp(f(w_i, w_c)) \right)$$

# word2vec (Google 2013)

The Idea : Let's maximize cooccurrence probability for cooccurring words (don't care about order – bag-of-words).



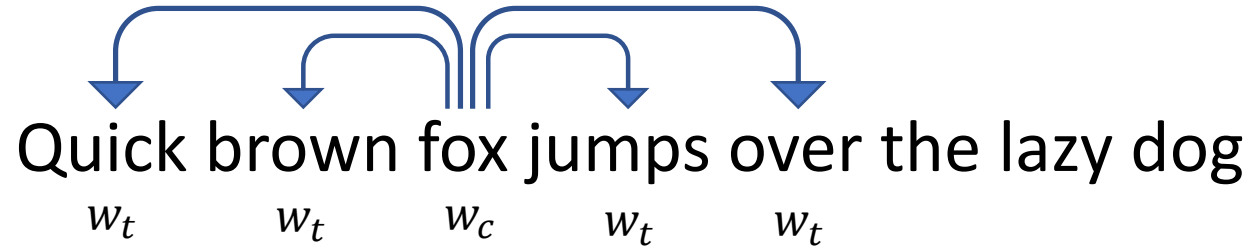
$$P(w_t|w_c) = \frac{\exp(f(w_t, w_c))}{\sum_{w_i \in Dict} \exp(f(w_i, w_c))}$$

$$\text{Loss function } J = \frac{1}{T} \sum J_t, \text{ where } J_t = -\log P(w_t|w_c) = -f(w_t, w_c) + \log \left( \sum_{w_i \in Dict} \exp(f(w_i, w_c)) \right)$$

The simplest (but effective approach):  $f(w_t, w_c) = u_{w_t}^T v_{w_c}$ , where  $u$  and  $v$  are two-vector word embeddings.

# word2vec (Google 2013)

The Idea : Let's maximize cooccurrence probability for cooccurring words (don't care about order – bag-of-words).



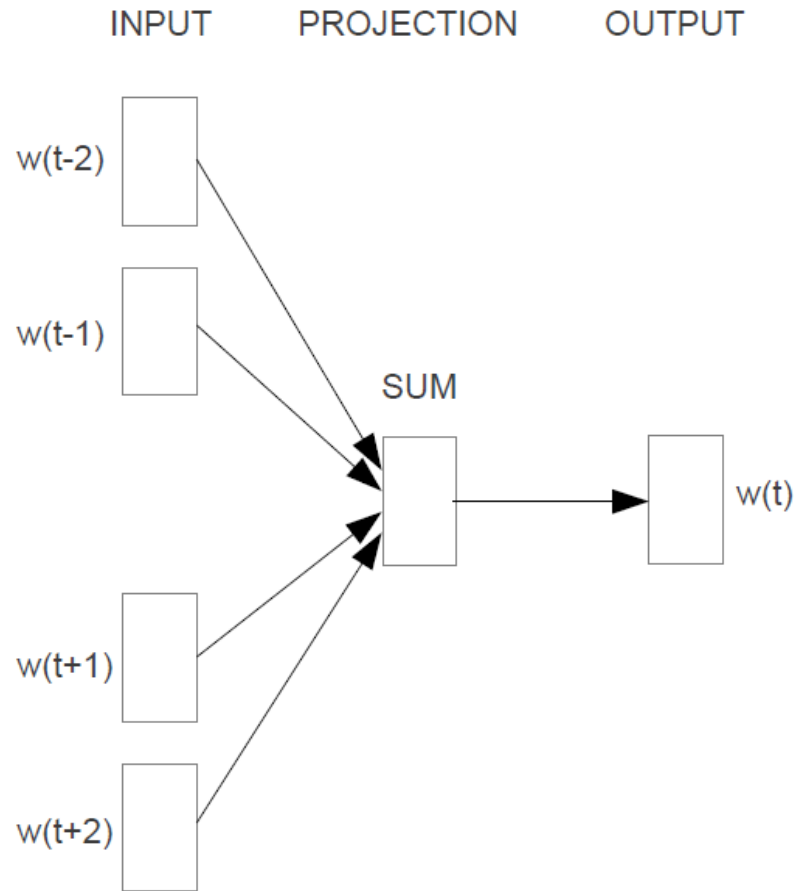
$$P(w_t|w_c) = \frac{\exp(f(w_t, w_c))}{\sum_{w_i \in Dict} \exp(f(w_i, w_c))}$$

$$\text{Loss function } J = \frac{1}{T} \sum J_t, \text{ where } J_t = -\log P(w_t|w_c) = -f(w_t, w_c) + \log \left( \sum_{w_i \in Dict} \exp(f(w_i, w_c)) \right)$$

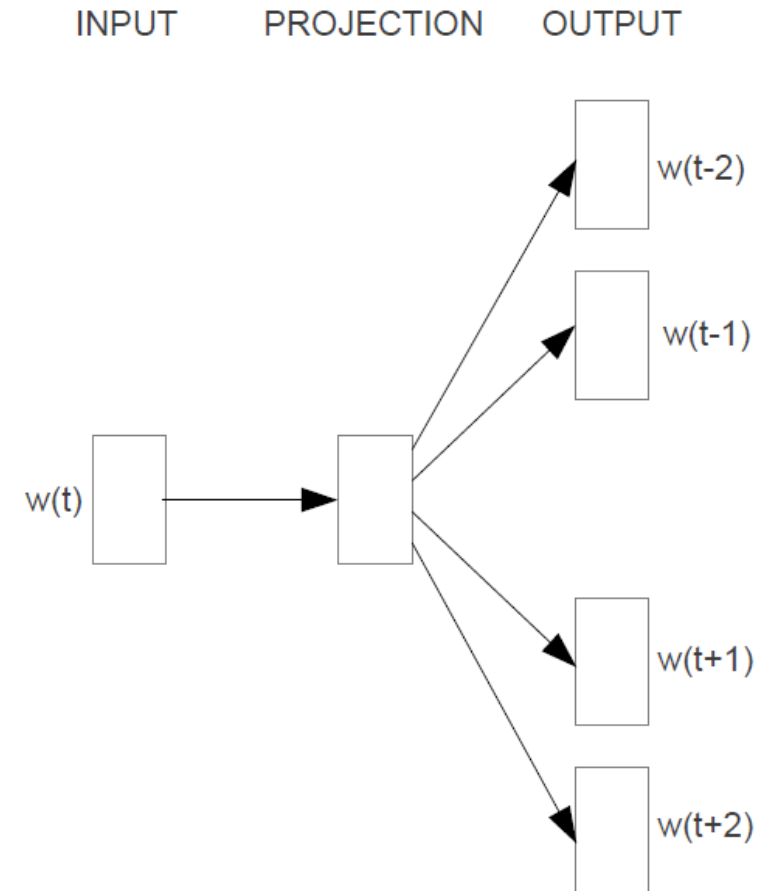
The simplest (but effective approach):  $f(w_t, w_c) = u_{w_t}^T v_{w_c}$ , where  $u$  and  $v$  are two-vector word embeddings.

We can use gradient descent.

# CBOW & Skip-gram



**CBOW**



**Skip-gram**

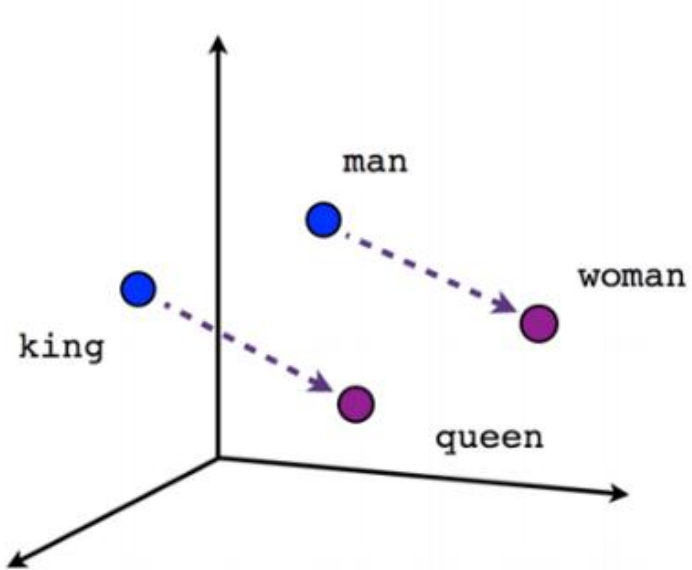


# Fasttext

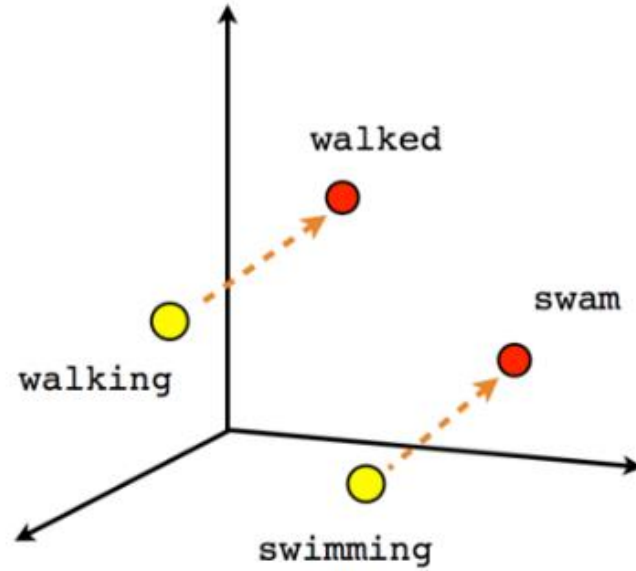
Adding n-grams to the dictionary:

where → <where> + <wh + whe + her + ere + re>

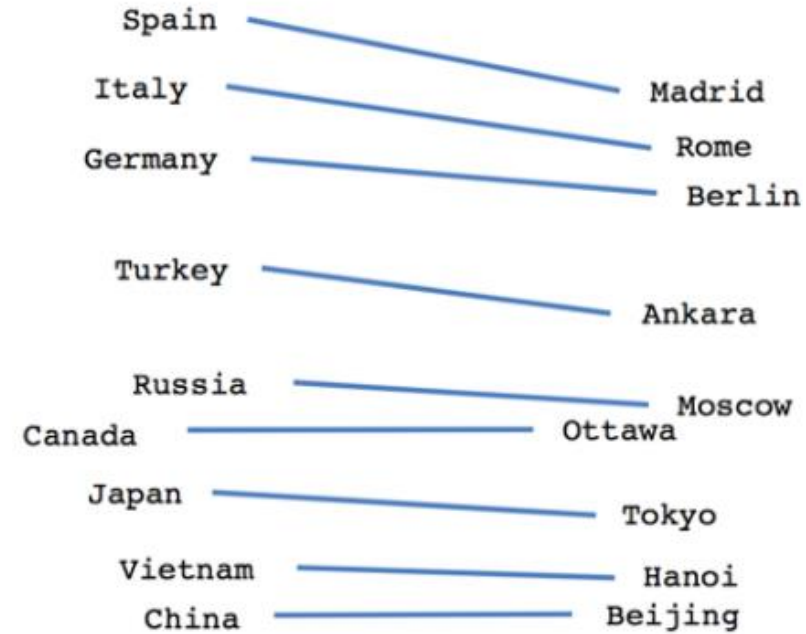
# Word Embeddings



Male-Female



Verb tense



Country-Capital

# Semantic and syntactic relations

Type of relationship	Word Pair 1		Word Pair 2	
Common capital city	Athens	Greece	Oslo	Norway
All capital cities	Astana	Kazakhstan	Harare	Zimbabwe
Currency	Angola	kwanza	Iran	rial
City-in-state	Chicago	Illinois	Stockton	California
Man-Woman	brother	sister	grandson	granddaughter
Adjective to adverb	apparent	apparently	rapid	rapidly
Opposite	possibly	impossibly	ethical	unethical
Comparative	great	greater	tough	tougher
Superlative	easy	easiest	lucky	luckiest
Present Participle	think	thinking	read	reading
Nationality adjective	Switzerland	Swiss	Cambodia	Cambodian
Past tense	walking	walked	swimming	swam
Plural nouns	mouse	mice	dollar	dollars
Plural verbs	work	works	speak	speaks

# skip-gram pair results

Relationship	Example 1	Example 2	Example 3
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - France	Berlusconi: Italy	Merkel: Germany	Koizumi: Japan
copper - Cu	zinc: Zn	gold: Au	uranium: plutonium
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza

# Metrics comparison

Model	Dim.	Dataset	Semantic	Syntactic	Average
PCA	100	1.6B	4.2	16.4	10.8
GloVe	100	1.6B	<b>67.5</b>	<b>54.3</b>	<b>60.3</b>
SG	300	1.6B	61	61	61
CBOW	300	1.6B	16.1	52.6	36.1
GloVe	300	1.6B	<b>80.8</b>	<b>61.5</b>	<b>70.3</b>
SVD	300	6B	6.3	8.1	7.3
CBOW	300	6B	63.6	<b>67.4</b>	65.7
GloVe	300	6B	<b>77.4</b>	67.0	<b>71.7</b>
CBOW	1000	6B	57.3	68.9	63.7
SG	1000	6B	66.1	65.1	65.6
SVD	300	42B	38.4	58.2	49.2
GloVe	300	42B	<b>81.9</b>	<b>69.3</b>	<b>75.0</b>

Model	ACE	MUC7
SVD	77.3	73.7
PCA	81.7	80.7
CBOW	82.2	81.1
GloVe	<b>82.9</b>	<b>82.2</b>

# BioVectors

Original Sequence

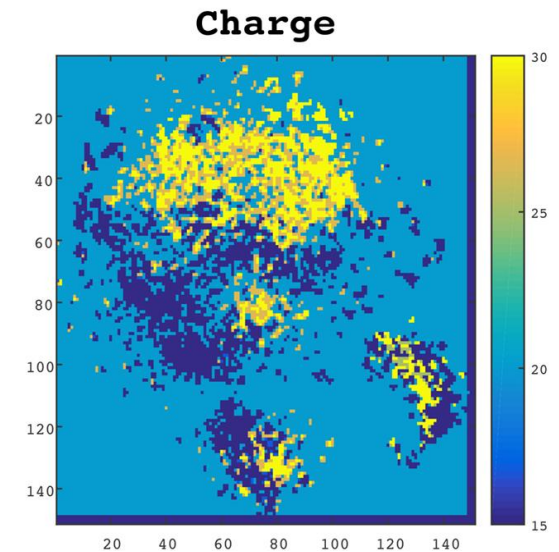
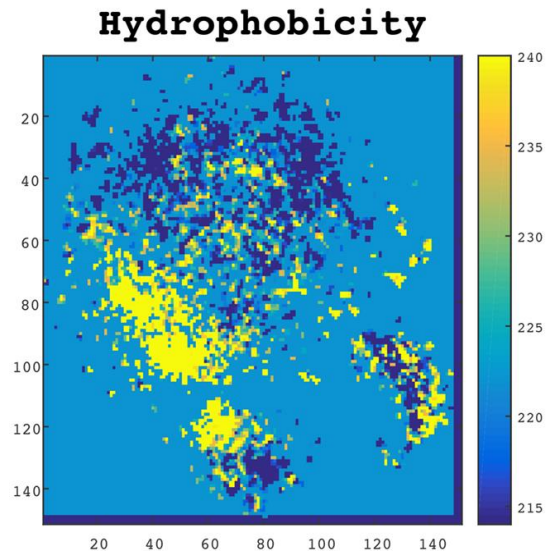
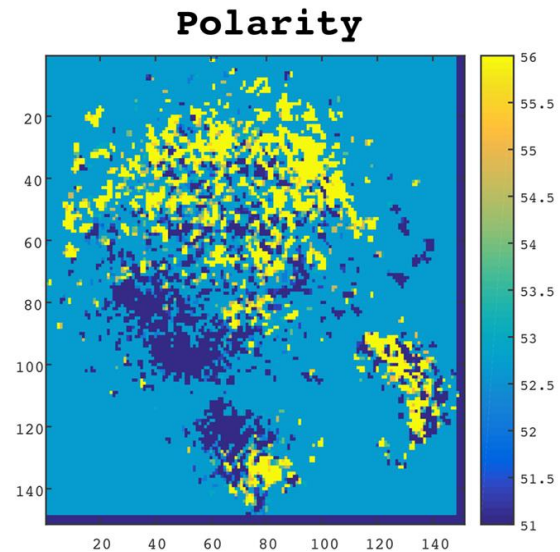
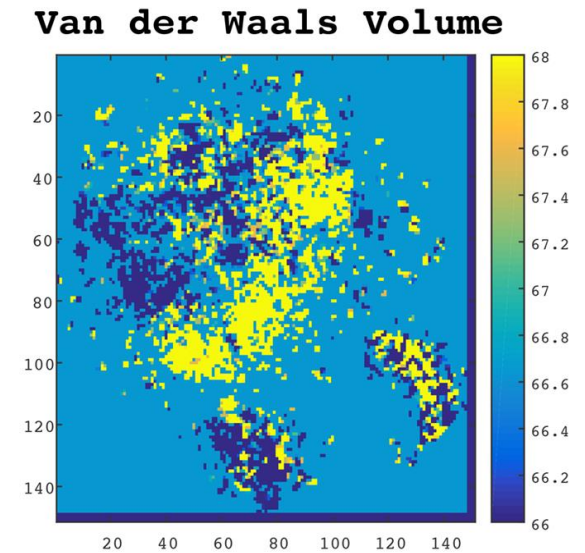
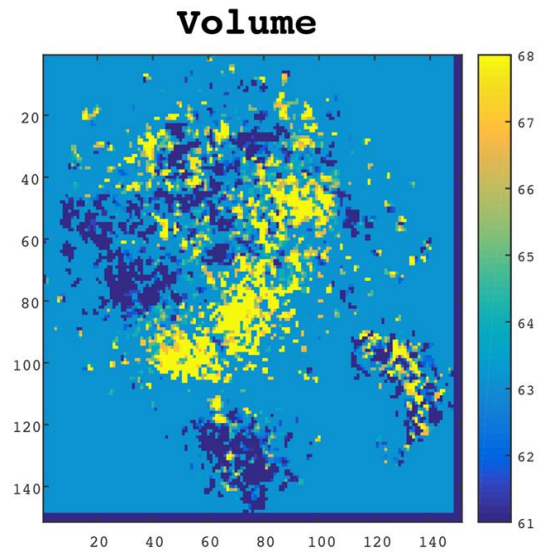
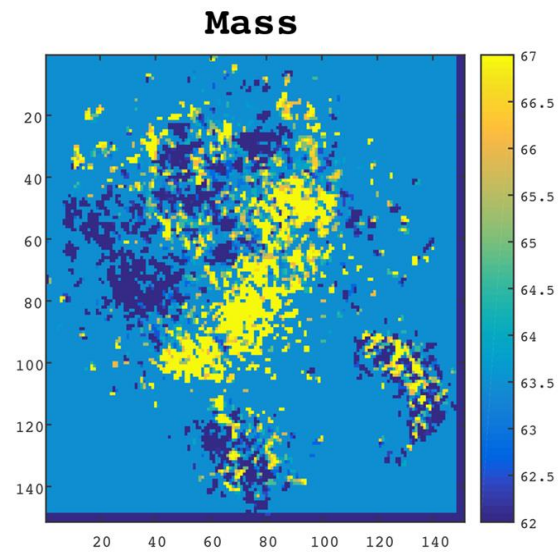
(1)  $\vec{M}$  (2)  $\vec{A}$  (3)  $\vec{F}$  *SAEDV LKEY DRRRR MEAL..*

Splittings

$\left\{ \begin{array}{l} \textcolor{red}{1)} \text{ MAF, SAE, DVL, KEY, DRR, RRM, ..} \\ \textcolor{blue}{2)} \text{ AFS, AED, VLK, EYD, RRR, RME, ..} \\ \textcolor{green}{3)} \text{ FSA ,EDV, LKE, YDR, RRR, MEA, ..} \end{array} \right.$



# BioVecs

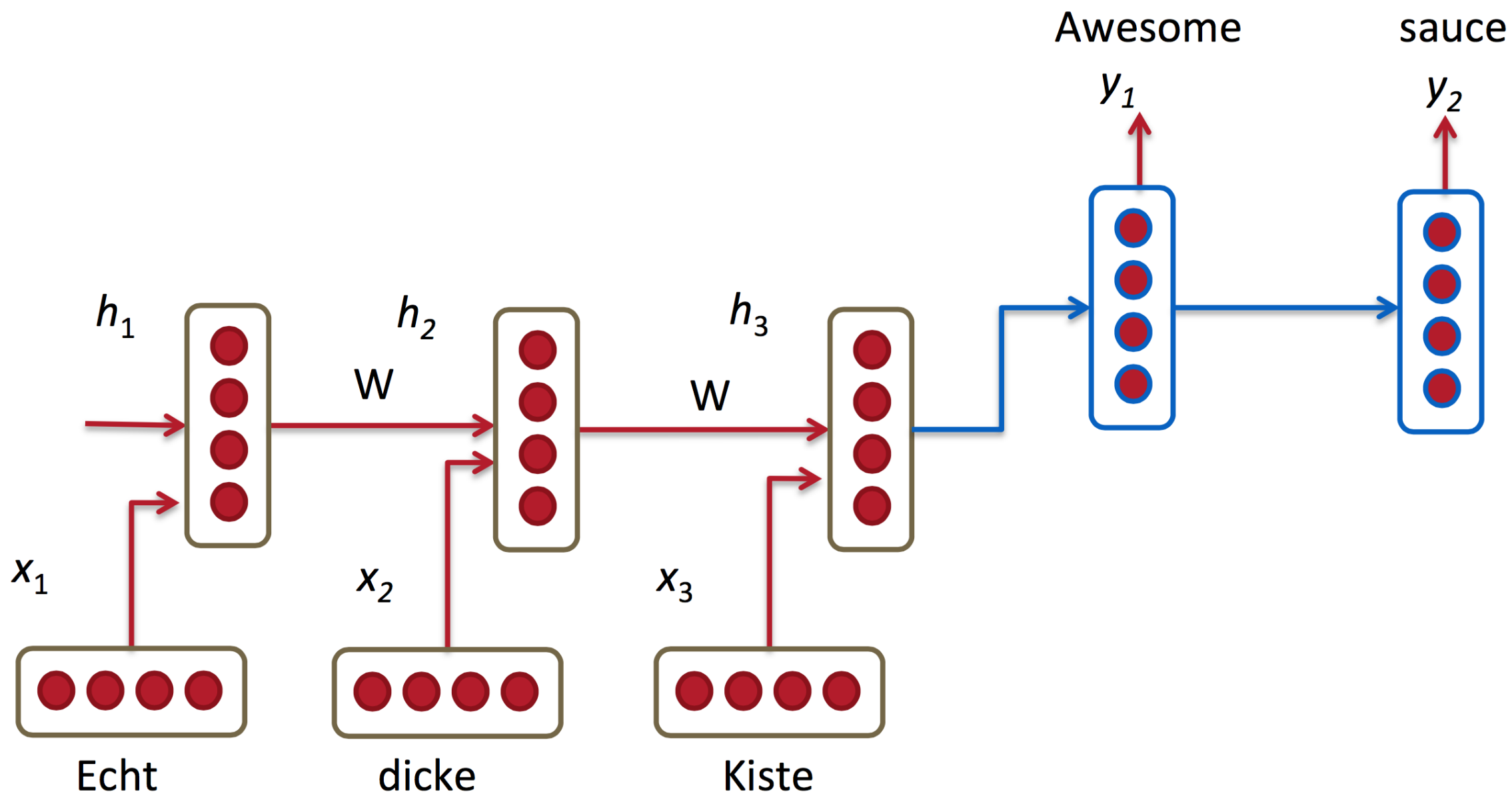


# Protein family classification

Family name	Training instances		Classification Result		
	# of positive sequences	# of negative sequences	Specificity	Sensitivity	Accuracy
50S ribosome-binding GTPase	3,084	3,084	0.95	0.93	0.94
Helicase conserved C-terminal domain	2,518	2,518	0.83	0.80	0.82
ATP synthase alpha-beta family, nucleotide-binding domain	2,387	2,387	0.98	0.97	0.97
7 transmembrane receptor (rhodopsin family)	1,820	1,820	0.95	0.96	0.95
Amino acid kinase family	1,750	1,750	0.91	0.92	0.91
ATPase family associated with various cellular activities (AAA)	1711	1711	0.92	0.90	0.91
tRNA synthetases class I (I, L, M and V)	1,634	1,634	0.97	0.97	0.97
tRNA synthetases class II (D, K and N)	1,419	1,419	0.88	0.83	0.85
Major Facilitator Superfamily	1,303	1,303	0.95	0.97	0.96
Hsp70 protein	1,272	1,272	0.97	0.97	0.97
NADH-Ubiquinone-plastoquinone (complex I), various chains	1,251	1,251	0.97	0.97	0.97
Histidine biosynthesis protein	1,248	1,248	0.96	0.97	0.97
TCP-1-cpn60 chaperonin family	1,246	1,246	0.95	0.96	0.95

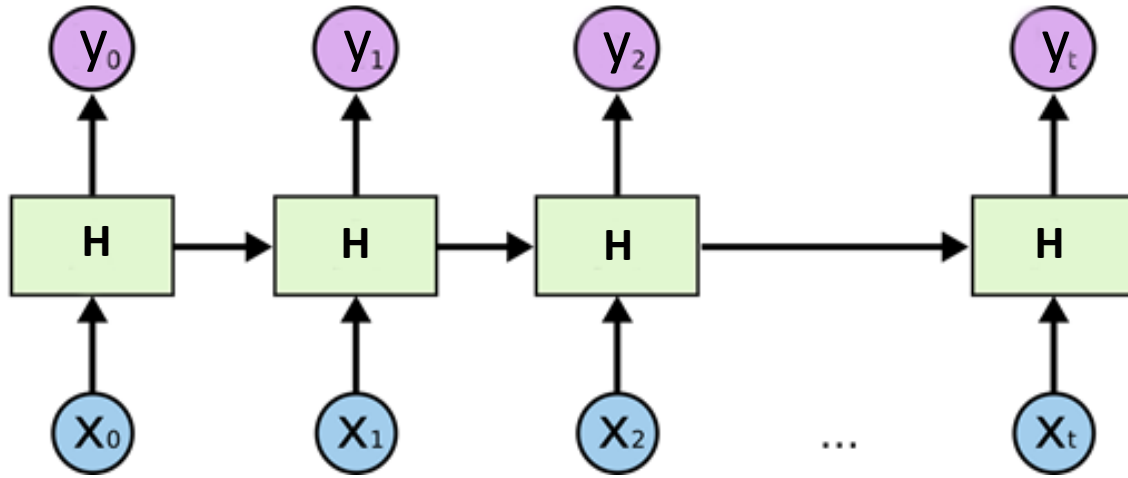


# Deep Learning for Natural Language Processing

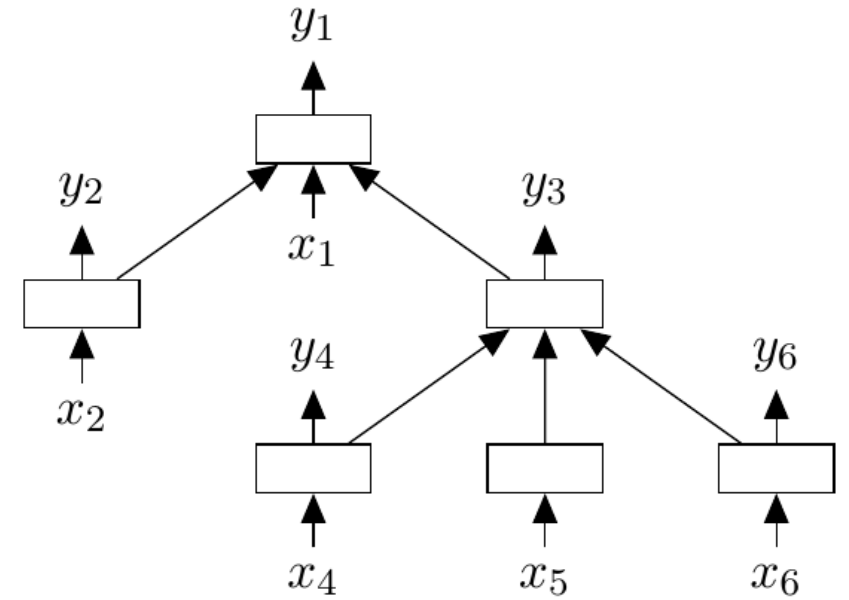


# RNN

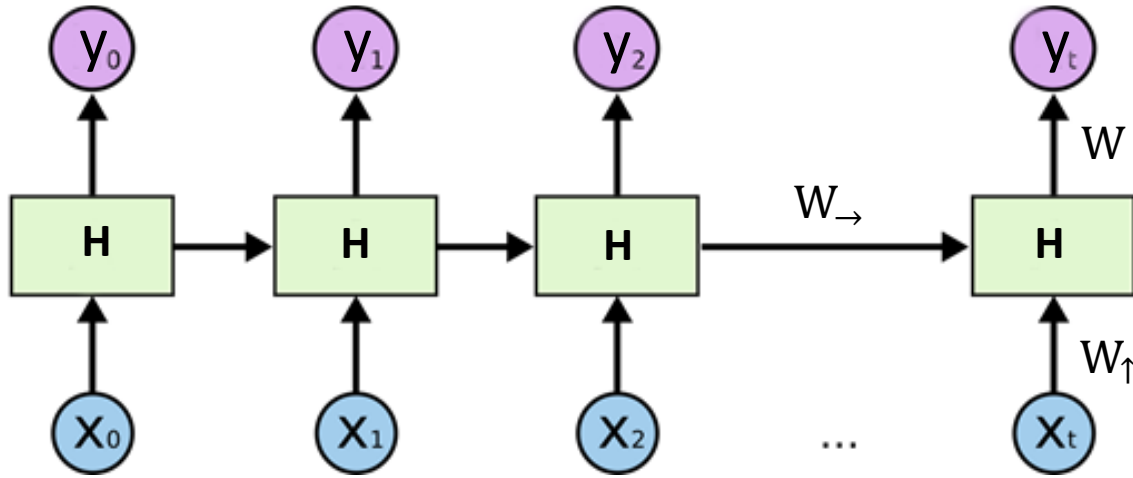
Recurrent



Recursive



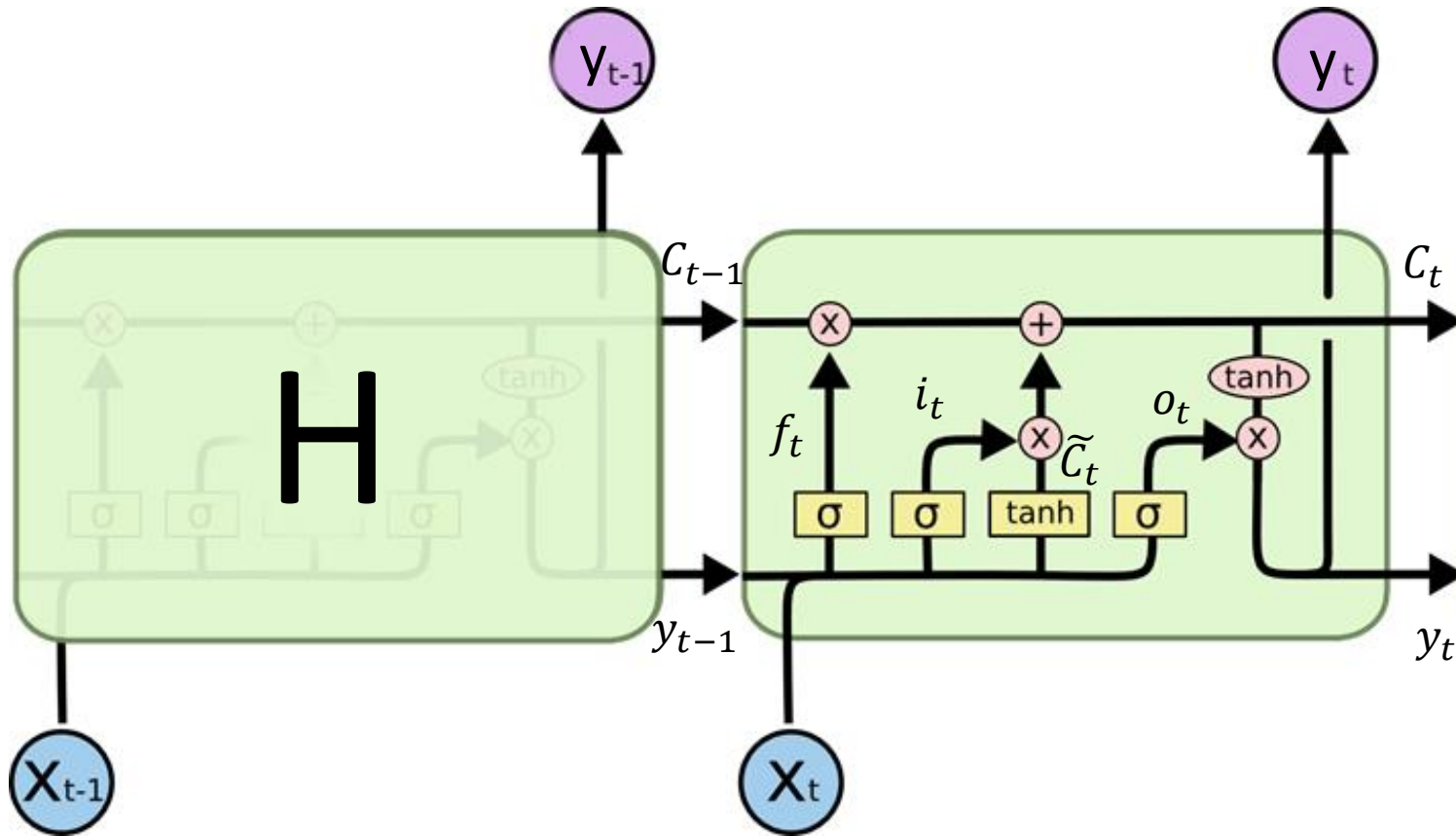
# Recurrent Neural Network



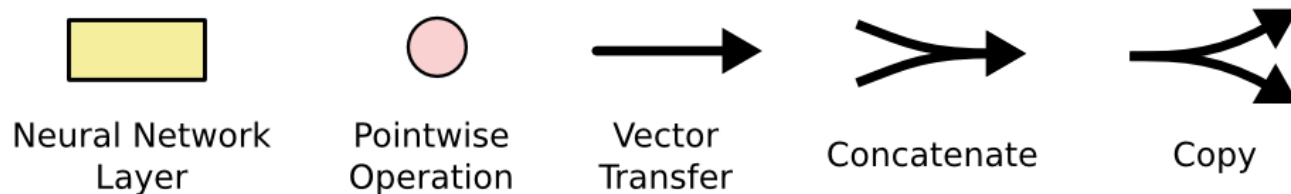
$$H_t = W_{\rightarrow} \sigma(H_{t-1}) + W_{\uparrow}(x_t)$$

$$\hat{y}_t = W \sigma(H_t)$$

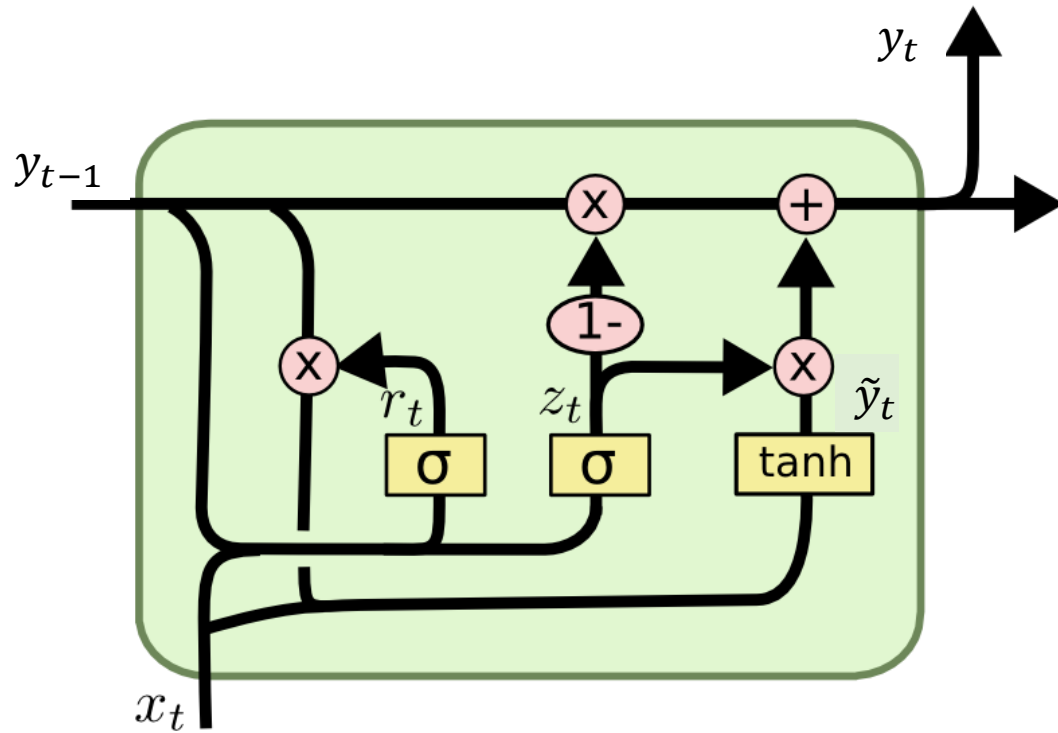
# LSTM (Long short-term memory)



**forget gate** –  $f_t = \sigma(W_f \cdot [y_{t-1}, x_t] + b_f)$   
**input gate** –  $i_t = \sigma(W_i \cdot [y_{t-1}, x_t] + b_i)$   
 $\tilde{c}_t = \tanh(W_c \cdot [y_{t-1}, x_t] + b_c)$   
 $c_t = f_t * c_{t-1} + i_t * \tilde{c}_t$   
**output gate** –  $o_t = \sigma(W_o \cdot [y_{t-1}, x_t] + b_o)$   
 $y_t = o_t * \tanh(c_t)$



# GRU (Gated Recurrent Unit)



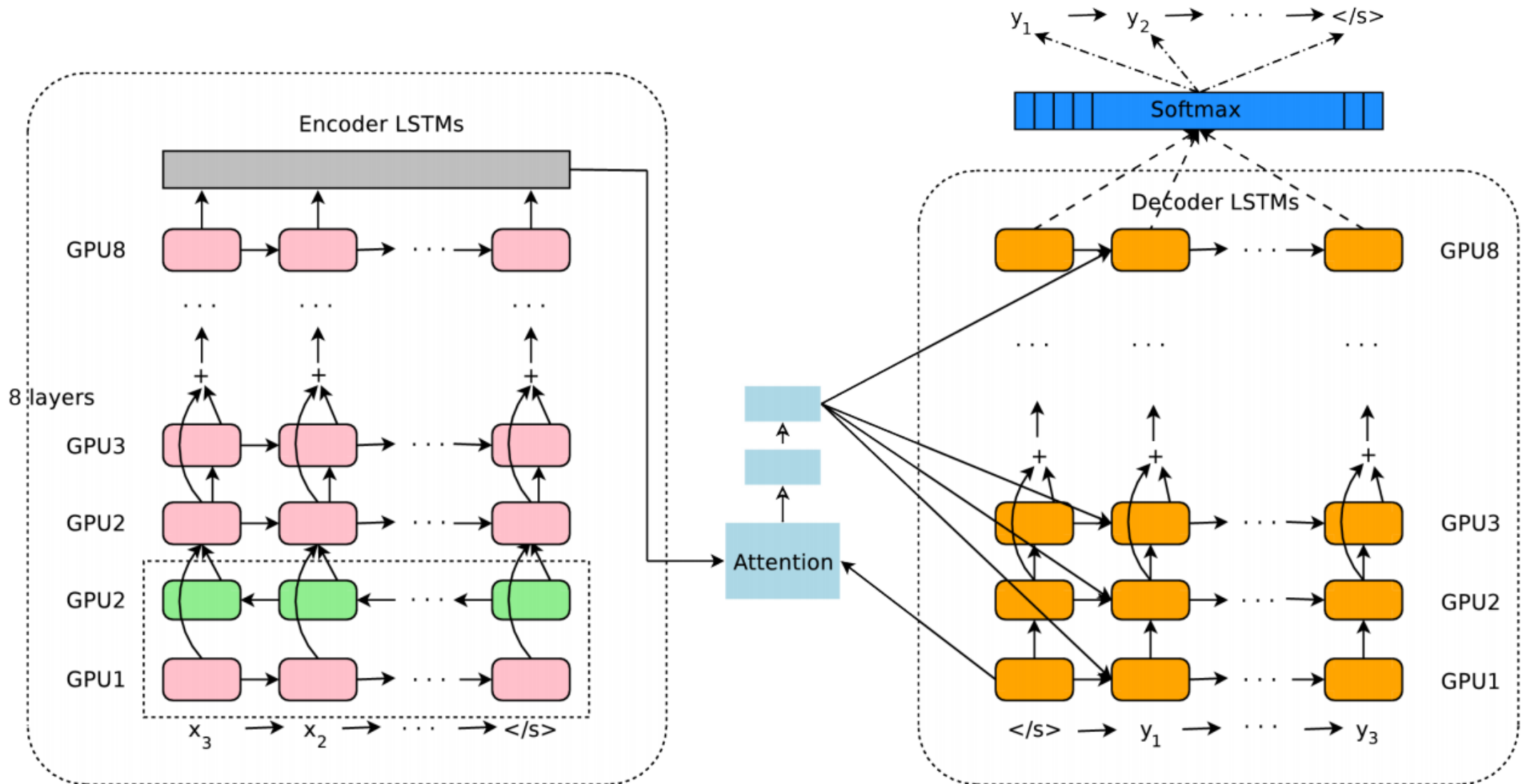
*update gate*  $- z_t = \sigma(W_z \cdot [y_{t-1}, x_t])$

*reset gate*  $- r_t = \sigma(W_r \cdot [y_{t-1}, x_t])$

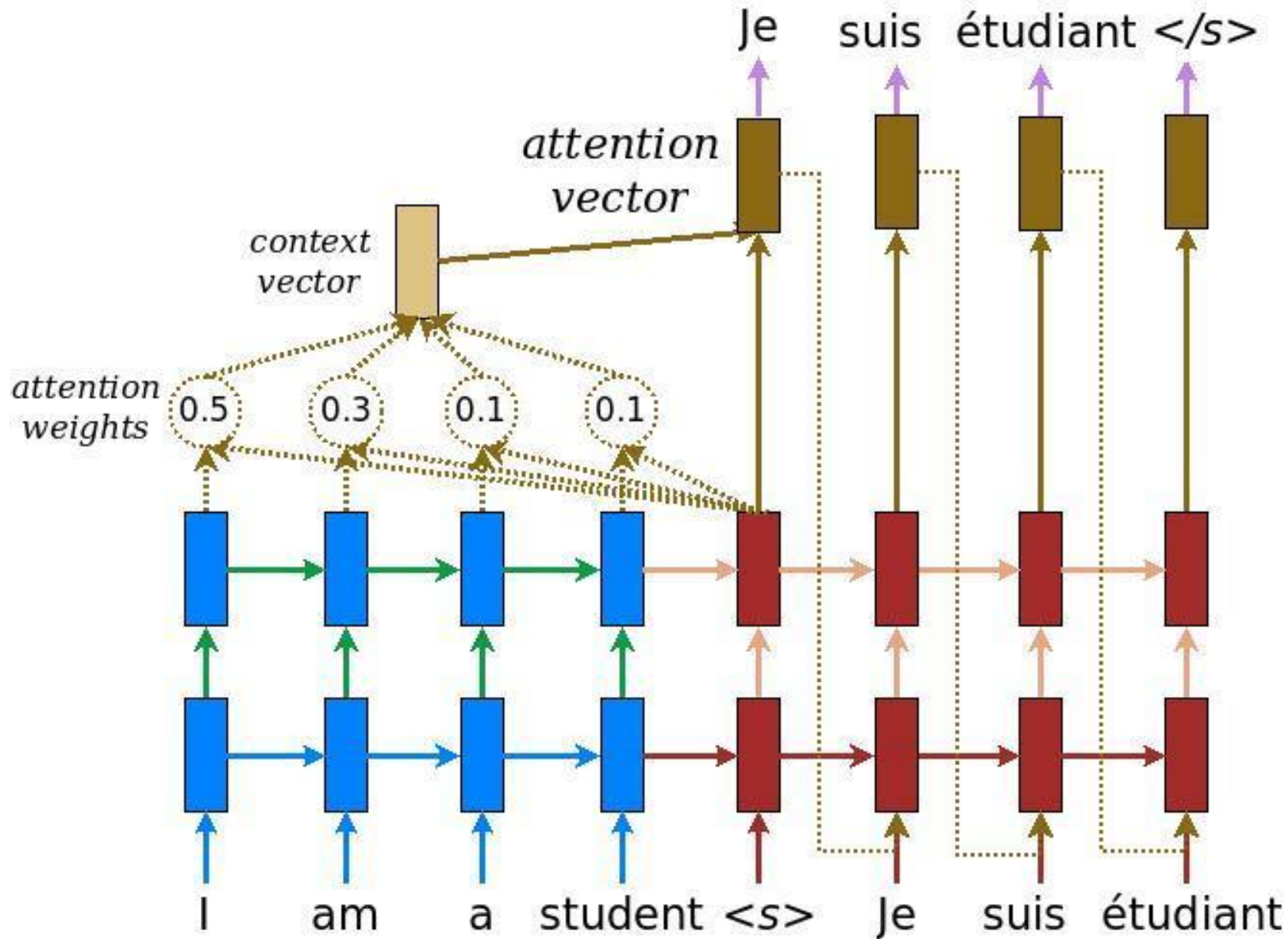
$$\tilde{y}_t = \tanh(W \cdot [r_t * y_{t-1}, x_t])$$

$$y_t = (1 - z_t) * y_{t-1} + z_t * \tilde{y}_t$$

# Google Neural Machine Translation

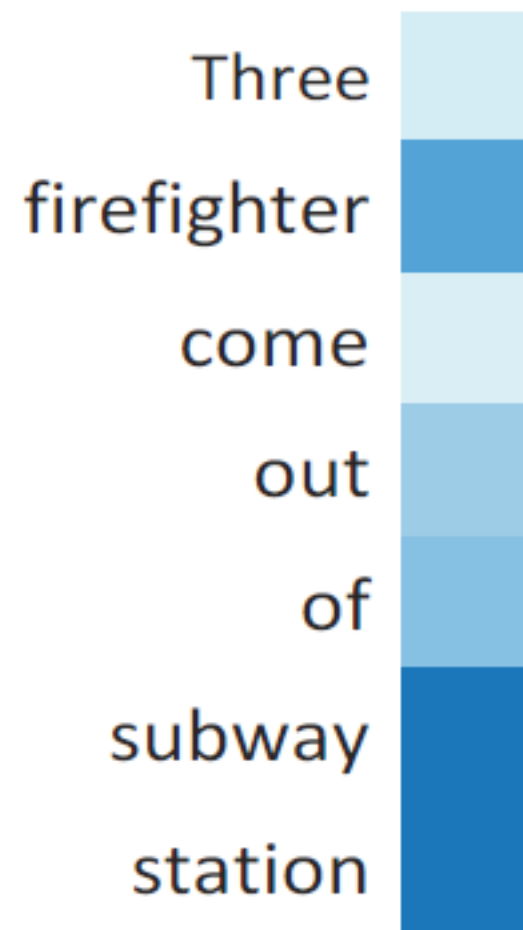
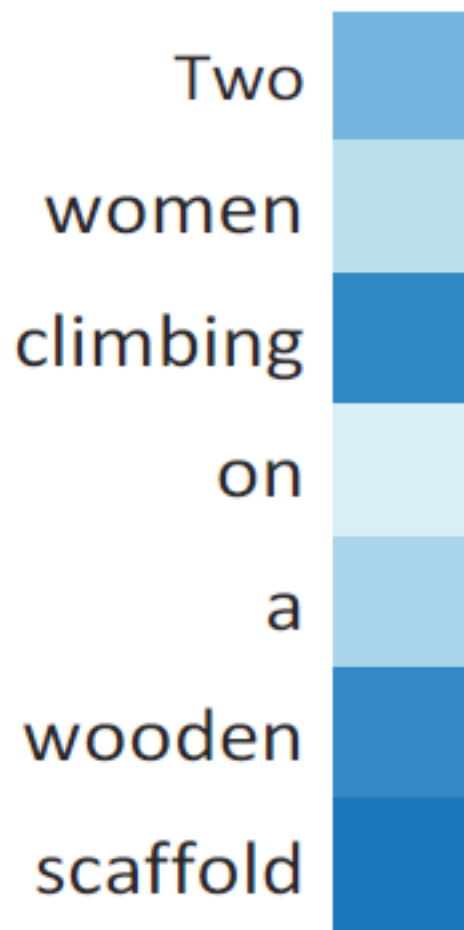
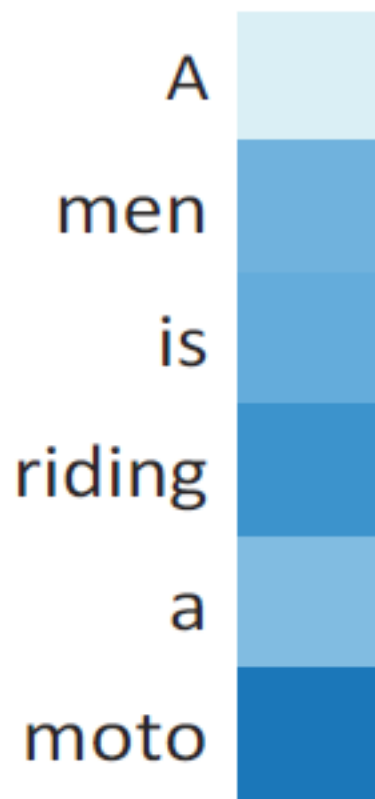


# Attention

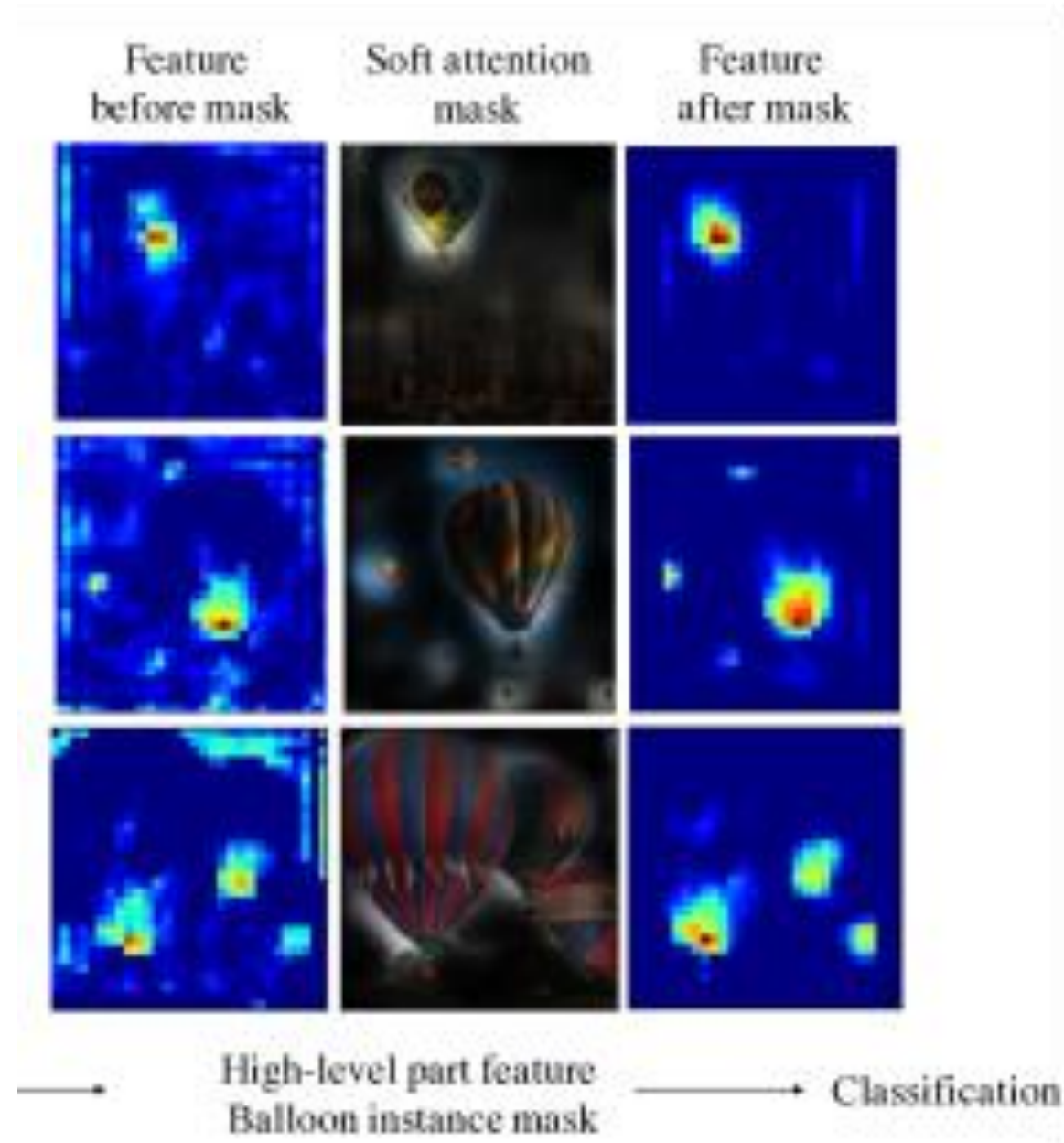




# Attention



# Visual attention



# SQuAD1

Rank	Model	EM	F1
	Human Performance <i>Stanford University</i> (Rajpurkar & Jia et al. '18)	86.831	89.452
1 Sep 13, 2018	nlnet (single model) <i>Microsoft Research Asia</i>	74.238	77.022
2 Sep 17, 2018	Unet (ensemble) <i>Fudan University &amp; Liulishuo Lab</i>	71.553	75.011
2 Aug 15, 2018	Reinforced Mnemonic Reader + Answer Verifier (single model) <i>NUDT</i> <a href="https://arxiv.org/abs/1808.05759">https://arxiv.org/abs/1808.05759</a>	71.699	74.238
2 Aug 28, 2018	SLQA+ (single model) <i>Alibaba DAMO NLP</i> <a href="http://www.aclweb.org/anthology/P18-1158">http://www.aclweb.org/anthology/P18-1158</a>	71.451	74.422
3 Sep 14, 2018	SAN (ensemble model) <i>Microsoft Business Applications Research Group</i> <a href="https://arxiv.org/abs/1712.03556">https://arxiv.org/abs/1712.03556</a>	71.282	73.658

# SQuAD2

Rank	Model	EM	F1
	Human Performance <i>Stanford University</i> (Rajpurkar & Jia et al. '18)	86.831	89.452
1 Jul 22, 2019	XLNet + DAAF + Verifier (ensemble) <i>PINGAN Omni-Sinitic</i>	88.592	90.859
2 Jul 19, 2019	XLNet + SG-Net Verifier (ensemble) <i>Shanghai Jiao Tong University &amp; CloudWalk</i>	88.050	90.645
3 Jul 23, 2019	XLNet + SG-Net Verifier (single model) <i>Shanghai Jiao Tong University &amp; CloudWalk</i>	87.046	89.899
3 Mar 20, 2019	BERT + DAE + AoA (ensemble) <i>Joint Laboratory of HIT and iFLYTEK Research</i>	87.147	89.474
3 Jul 20, 2019	RoBERTa (single model) <i>Facebook AI</i>	86.820	89.795
4 Mar 15, 2019	BERT + ConvLSTM + MTL + Verifier (ensemble) <i>Layer 6 AI</i>	86.730	89.286
5 Mar 05, 2019	BERT + N-Gram Masking + Synthetic Self-Training (ensemble) <i>Google AI Language</i> <a href="https://github.com/google-research/bert">https://github.com/google-research/bert</a>	86.673	89.147