

Deep Learning for Natural Language Processing

Hien T. Nguyen

Director, Artificial Intelligence Laboratory

Faculty of Information Technology

Ton Duc Thang University, Vietnam

Email: hien@tdt.edu.vn or hiennguyenthanh@gmail.com

Outline

- Introduction
- Review on basic machine learning models
- Deep Learning
- Distributed Representation
- Deep Learning for NLP
 - [Named Entity Recognition](#) and [Sequence Labelling](#) tasks
 - Sentence Modeling
 - Sentiment Analysis
 - Question Answering
 - Machine Translation
 - Syntactic Parsing
 - Visual Recognition
 - Speech Synthesis
- Conclusion

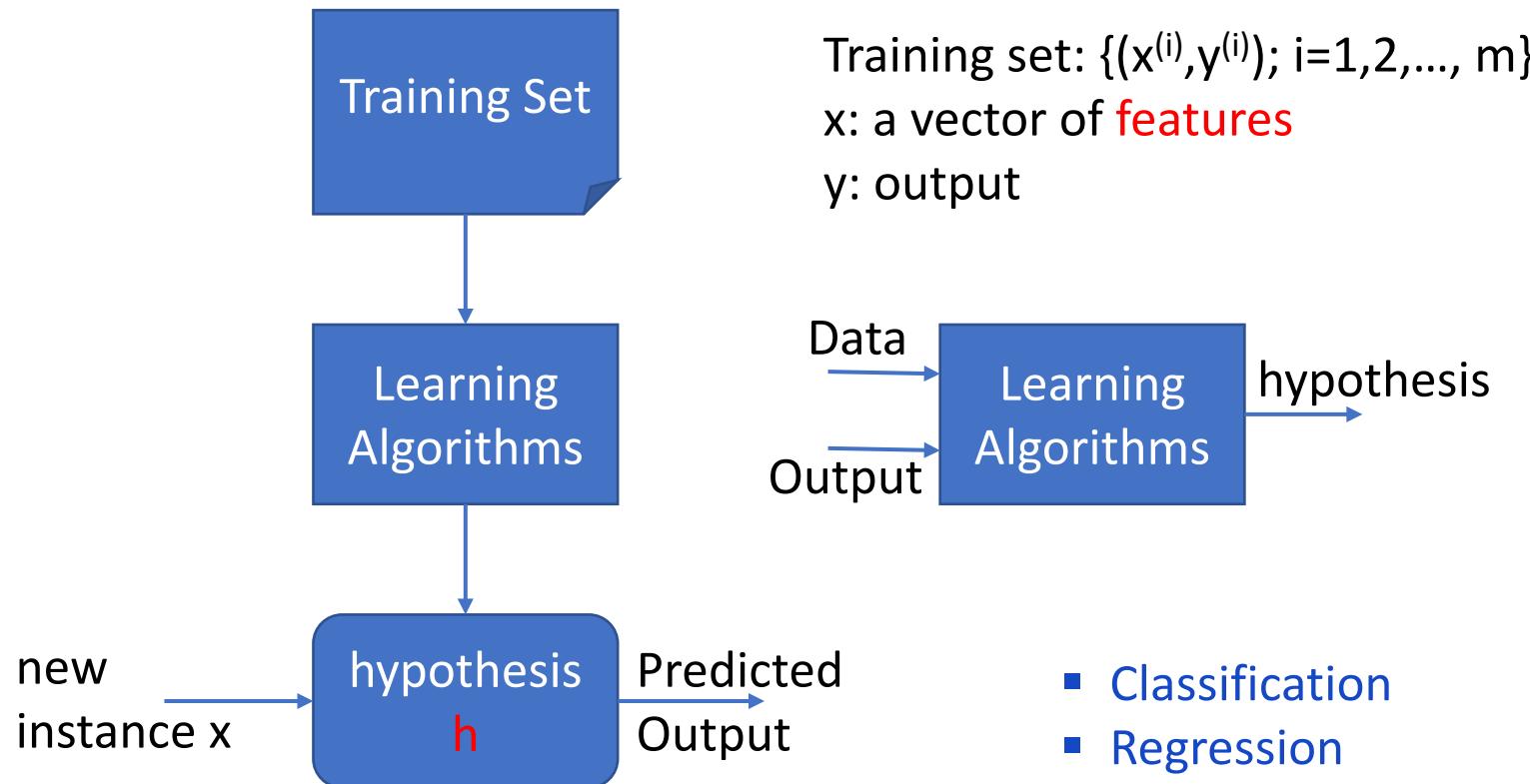
Review

- Learning models
- Logistic regression and Classification
- Perceptron
- Softmax regression
- Basic Neural Networks

Learning models

- Supervised learning
- Unsupervised learning
- Reinforcement learning

Supervised Learning



Supervised Learning

Living area (feet ²)	#bedrooms	Price (1000\$s)
2104	3	400
1600	3	330
2400	3	369
1416	2	232
3000	4	540
⋮	⋮	⋮

Linear Regression

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

$$h(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x, \quad (x_0=1)$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2.$$

(loss function)

Learning using Gradient descent

$$\theta_j := \theta_j - \alpha \boxed{\frac{\partial}{\partial \theta_j} J(\theta)}$$

learning rate

Supervised Learning

- Gradient descent
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$
 - Batch gradient descent
 - Stochastic gradient descent
 - Minibatch
- Loss function
 - Cross-entropy
 - Square error
 - Max-margin (Hinge loss)

Supervised vs. Unsupervised Learning

Supervised Learning

Data: (x, y)

x is data, y is label

Goal: Learn a *function* to map $x \rightarrow y$

Examples: Classification,
regression, object detection,
semantic segmentation, image
captioning, etc.

Unsupervised Learning

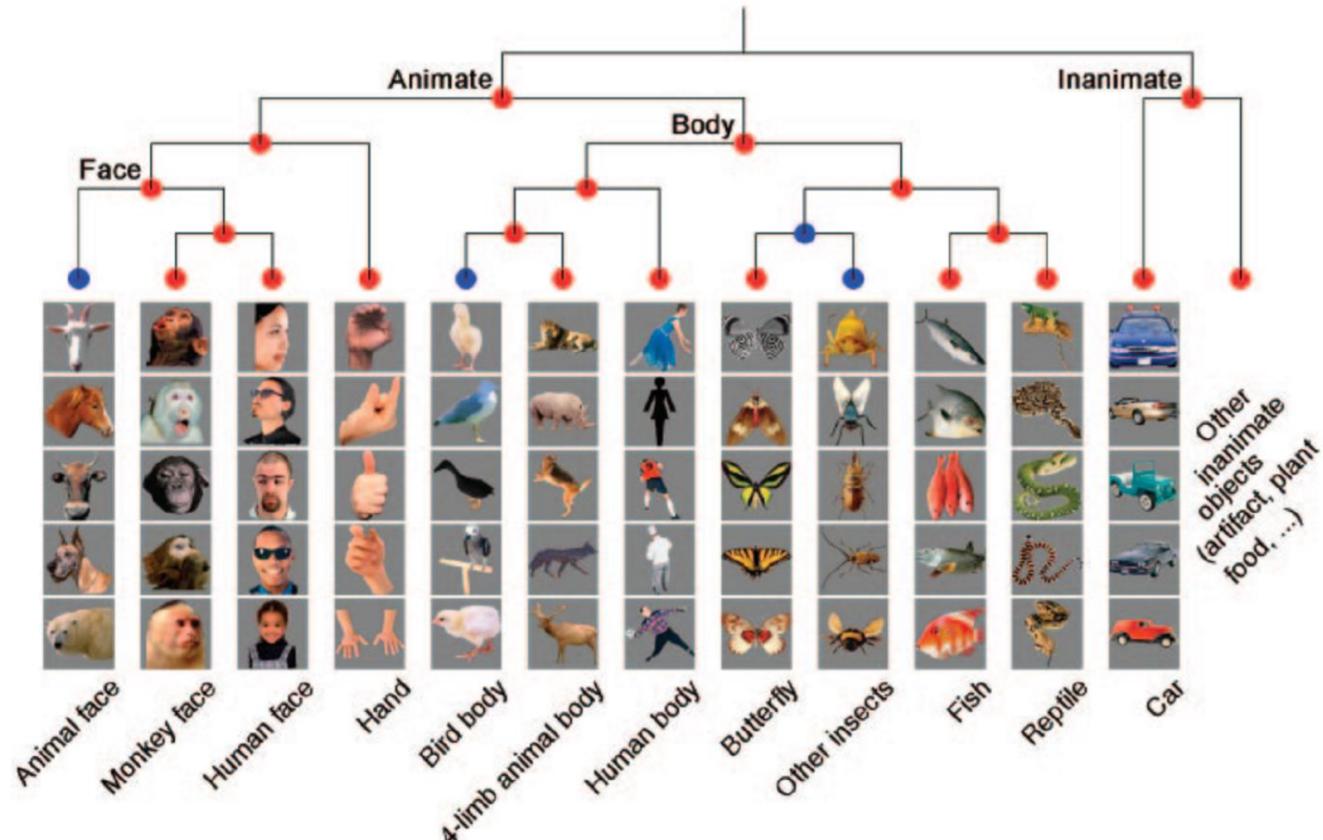
Data: x

Just data, no labels!

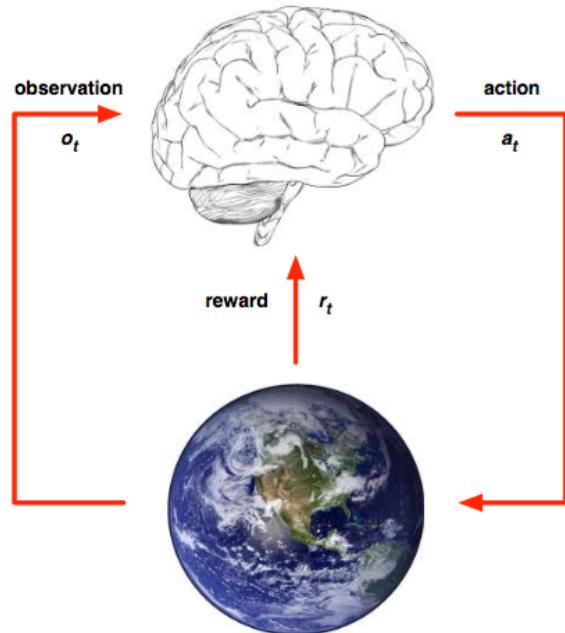
Goal: Learn some underlying
hidden *structure* of the data

Examples: Clustering,
dimensionality reduction, feature
learning, density estimation, etc.

Unsupervised Learning



Reinforcement Learning

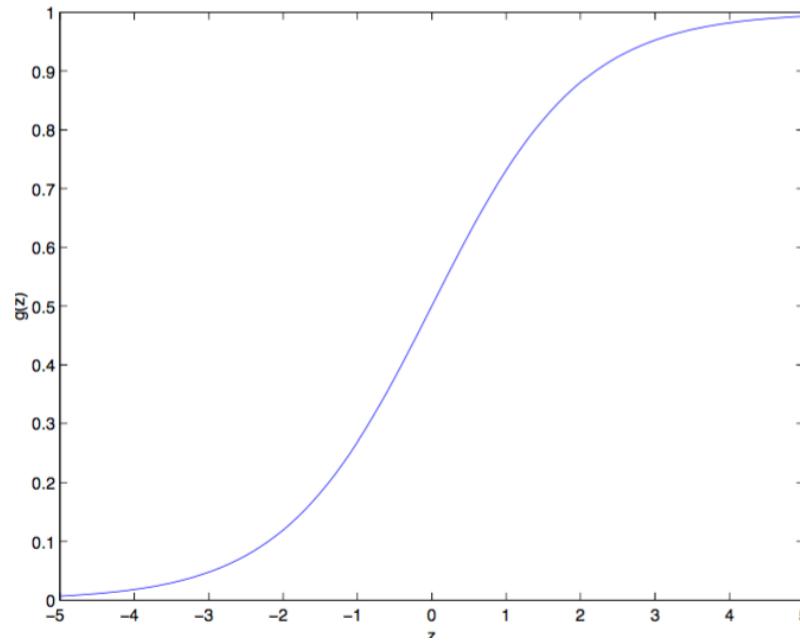


- ▶ At each step t the agent:
 - ▶ Executes action a_t
 - ▶ Receives observation o_t
 - ▶ Receives scalar reward r_t
- ▶ The environment:
 - ▶ Receives action a_t
 - ▶ Emits observation o_{t+1}
 - ▶ Emits scalar reward r_{t+1}

Logistic Regression and Classification

- Logistic Function (sigmoid)

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



Logistic Regression and Classification

- Logistic Function (sigmoid)

$$\begin{aligned}g'(z) &= \frac{d}{dz} \frac{1}{1+e^{-z}} \\&= \frac{1}{(1+e^{-z})^2} (e^{-z}) \\&= \frac{1}{(1+e^{-z})} \cdot \left(1 - \frac{1}{(1+e^{-z})}\right) \\&= g(z)(1-g(z)).\end{aligned}$$

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

- Our hypotheses $h_{\theta}(x)$

Logistic Regression and Classification

Let us assume that

$$\begin{aligned} P(y = 1 \mid x; \theta) &= h_\theta(x) \\ P(y = 0 \mid x; \theta) &= 1 - h_\theta(x) \end{aligned}$$

Note that this can be written more compactly as

$$p(y \mid x; \theta) = (h_\theta(x))^y (1 - h_\theta(x))^{1-y}$$

Assuming that the m training examples were generated independently, we can then write down the likelihood of the parameters as

$$\begin{aligned} L(\theta) &= p(\vec{y} \mid X; \theta) \\ &= \prod_{i=1}^m p(y^{(i)} \mid x^{(i)}; \theta) \\ &= \prod_{i=1}^m (h_\theta(x^{(i)}))^{y^{(i)}} (1 - h_\theta(x^{(i)}))^{1-y^{(i)}} \end{aligned}$$

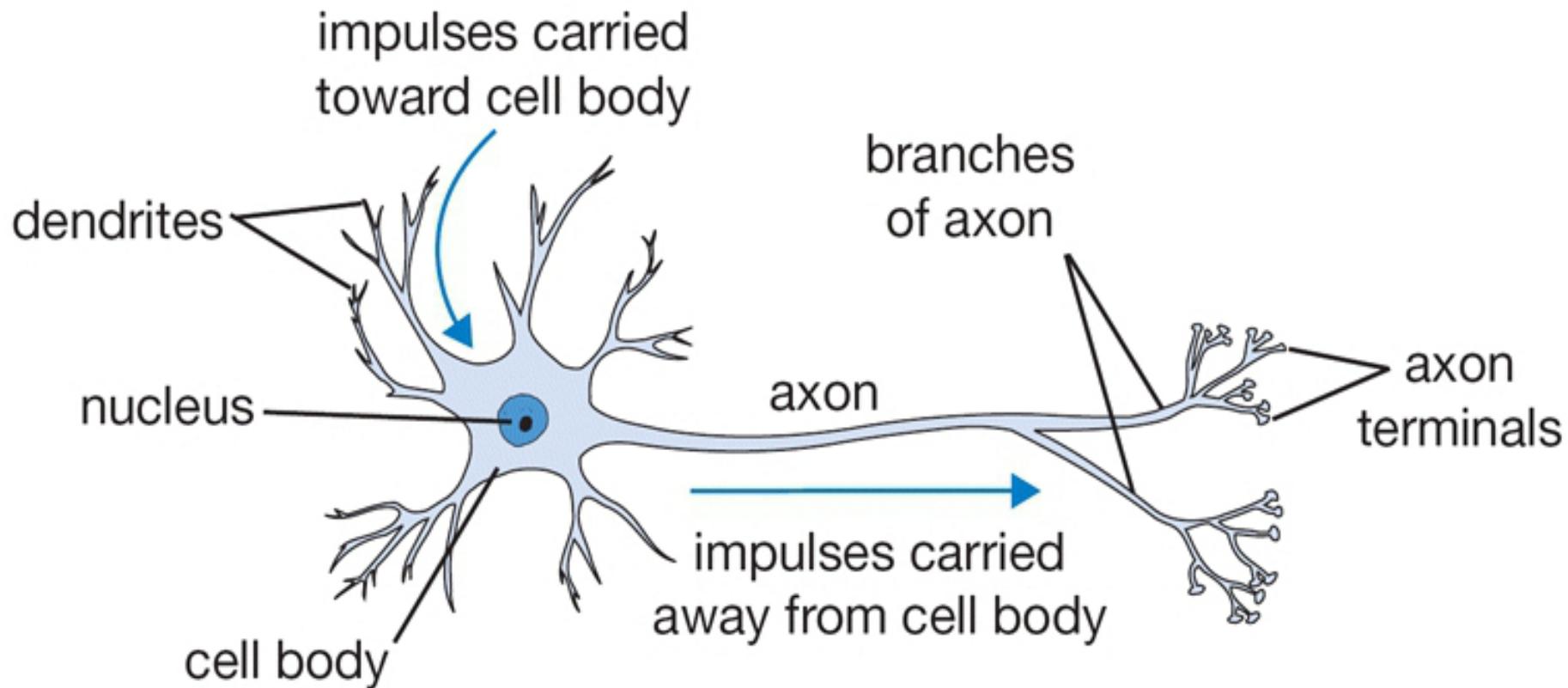
Logistic Regression and Classification

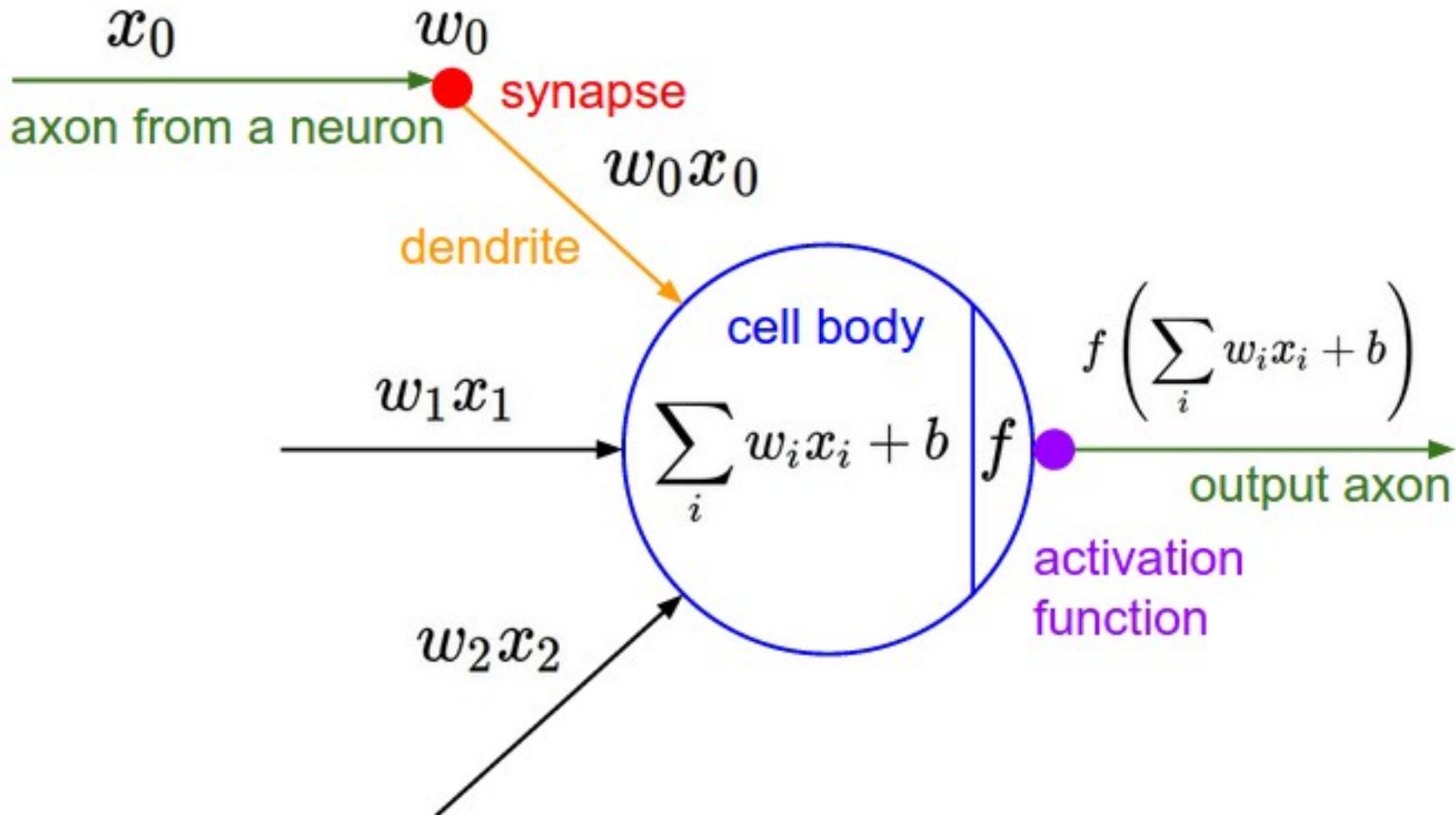
■ Log likelihood

$$\begin{aligned}\ell(\theta) &= \log L(\theta) \\ &= \sum_{i=1}^m y^{(i)} \log h(x^{(i)}) + (1 - y^{(i)}) \log(1 - h(x^{(i)}))\end{aligned}$$

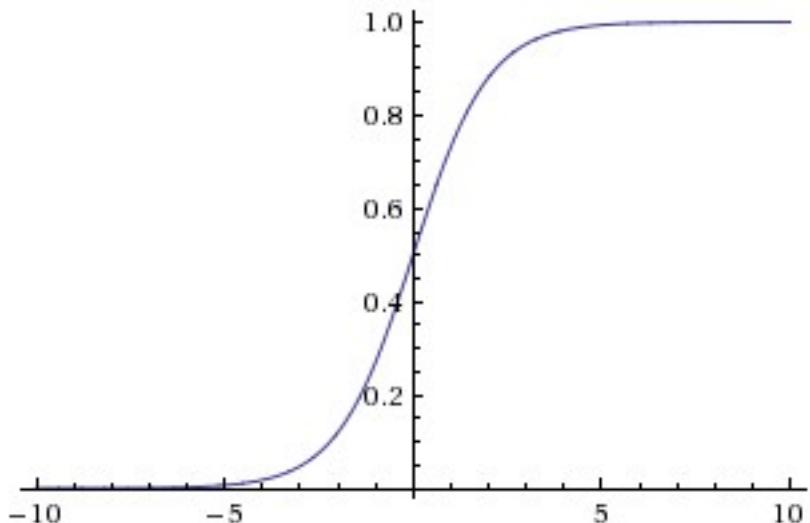
$$\begin{aligned}\frac{\partial}{\partial \theta_j} \ell(\theta) &= \left(y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)} \right) \frac{\partial}{\partial \theta_j} g(\theta^T x) \\ &= \left(y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)} \right) g(\theta^T x) (1 - g(\theta^T x)) \frac{\partial}{\partial \theta_j} \theta^T x \\ &= (y(1 - g(\theta^T x)) - (1 - y)g(\theta^T x)) x_j \\ &= (y - h_\theta(x)) x_j\end{aligned}$$

Perceptron

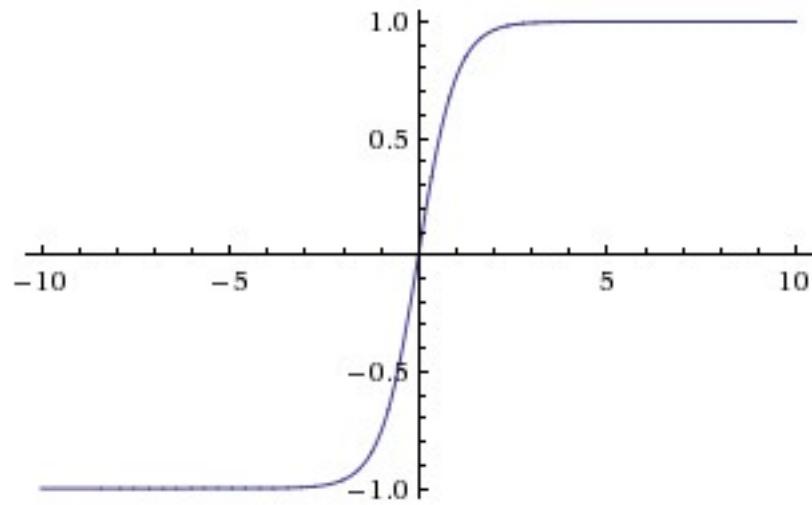




Activation function

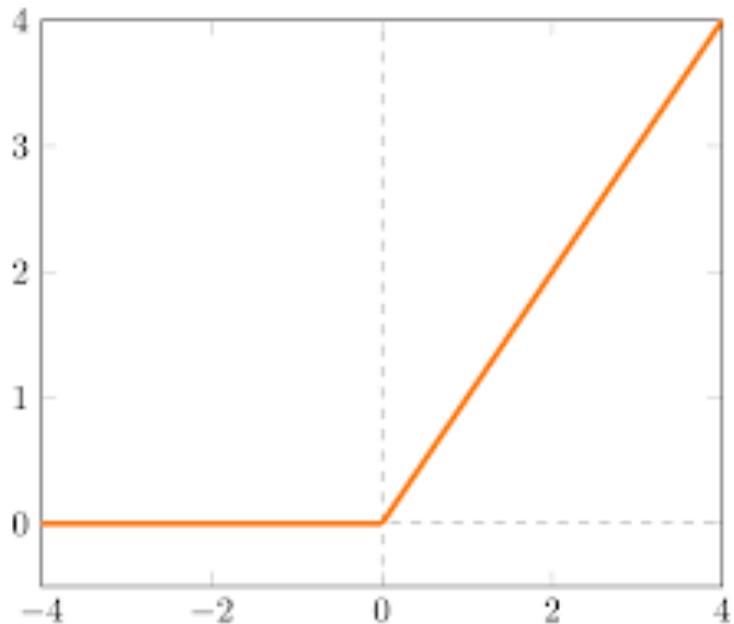


Sigmoid function



$$\text{Tanh}(x) = 2\text{sigmoid}(2x) - 1$$

Activation function



Rectified Linear Unit (ReLU)

$$\text{ReLU}(x) = \max(0, x)$$

Softmax

$$\text{softmax}(z_i) = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

matrix multiply + bias offset

0.01	-0.05	0.1	0.05
0.7	0.2	0.05	0.16
0.0	-0.45	-0.2	0.03

W

-15
22
-44
56

+

0.0
0.2
-0.3

b

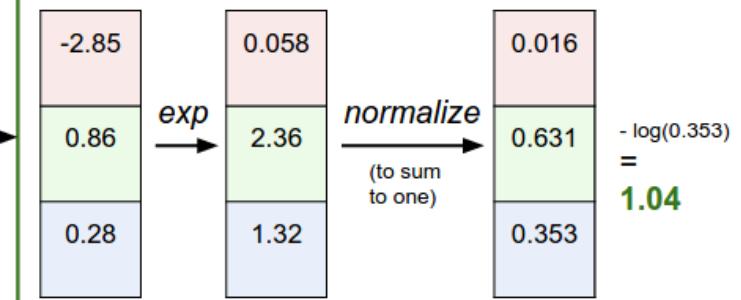
y_i 2

-2.85
0.86
0.28

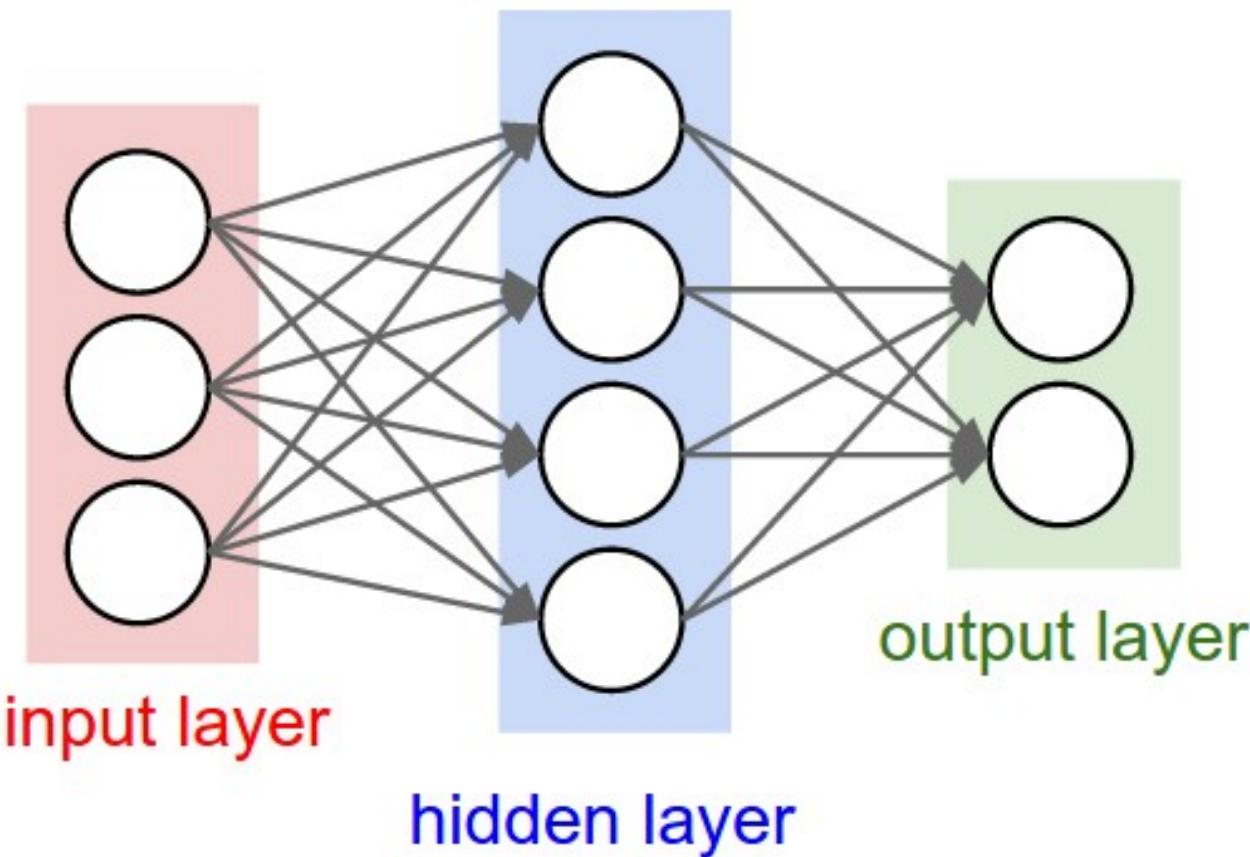
hinge loss (SVM)

$$\begin{aligned} & \max(0, -2.85 - 0.28 + 1) + \\ & \max(0, 0.86 - 0.28 + 1) \\ & = \\ & \mathbf{1.58} \end{aligned}$$

cross-entropy loss (Softmax)



Basic Neural Networks



Derivatives

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x + h) - f(x)}{h}$$

$$f(x + h) = f(x) + h \frac{df(x)}{dx}$$

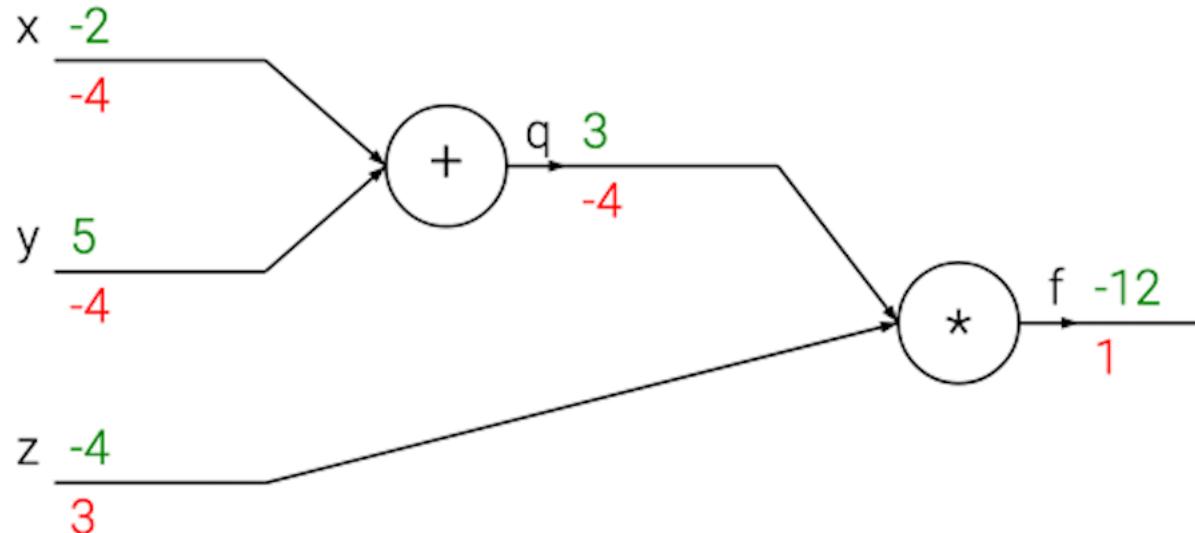
Chain rule

- $f(x,y,z) = (x+y)z$
- $q(x,y) = x + y$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

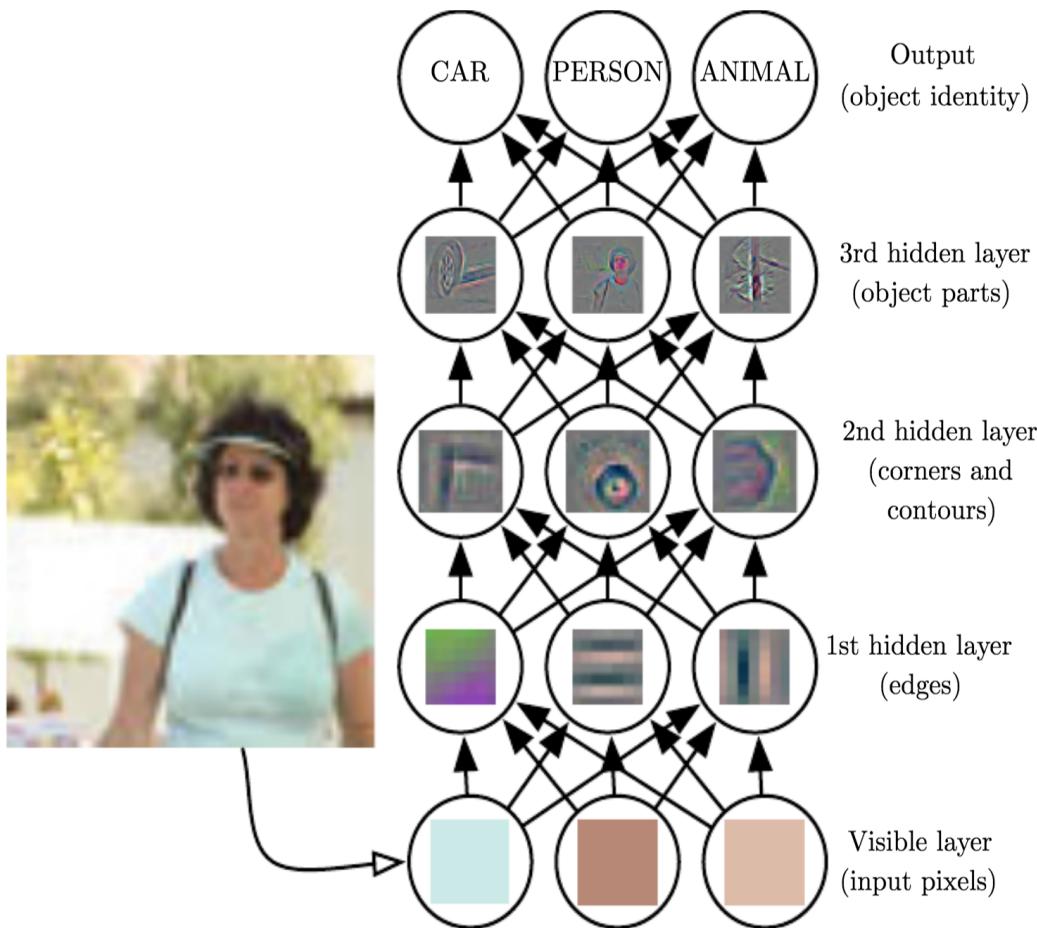
Backpropagation

- Forward pass
- Backward pass

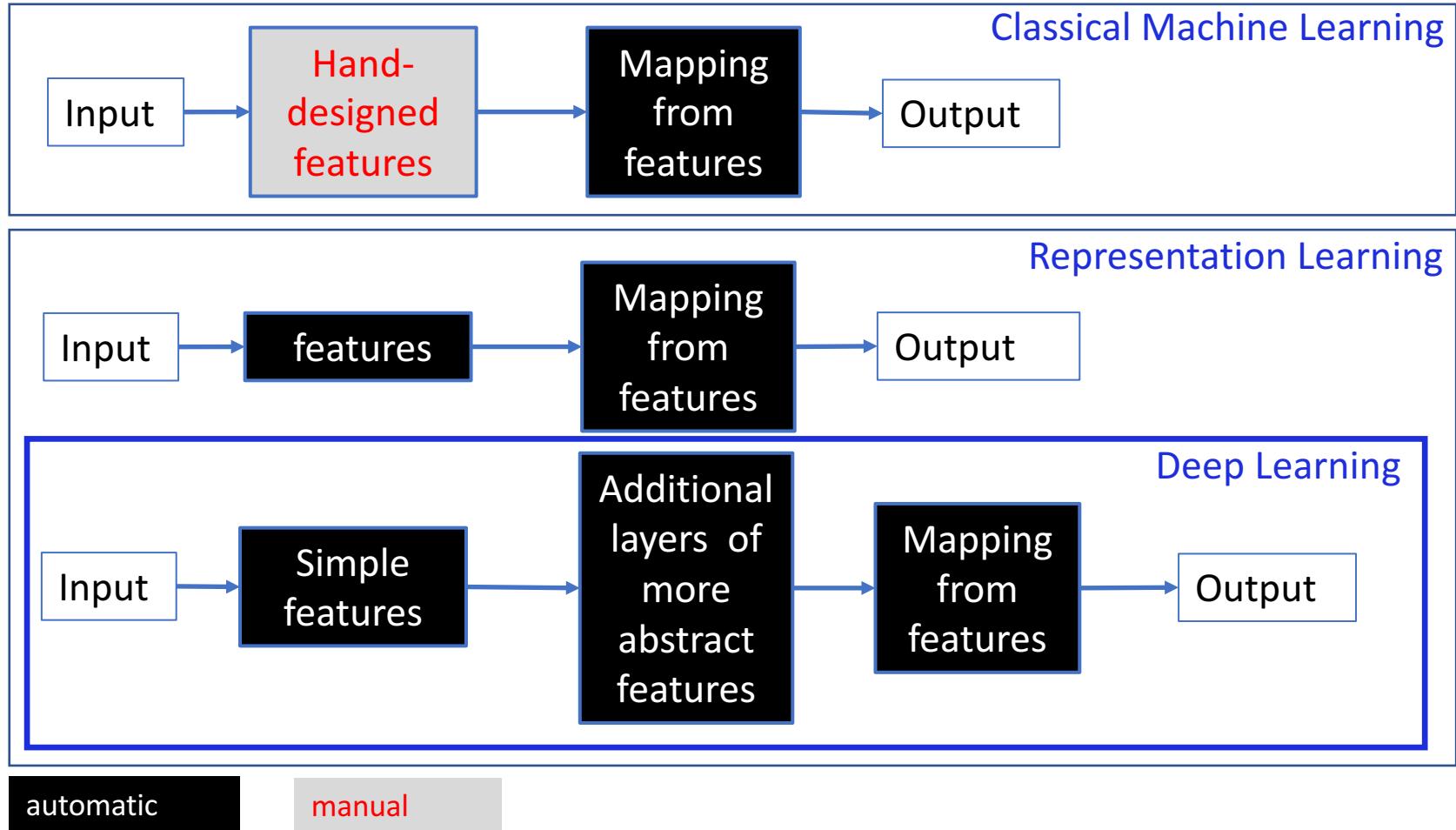


Deep Learning

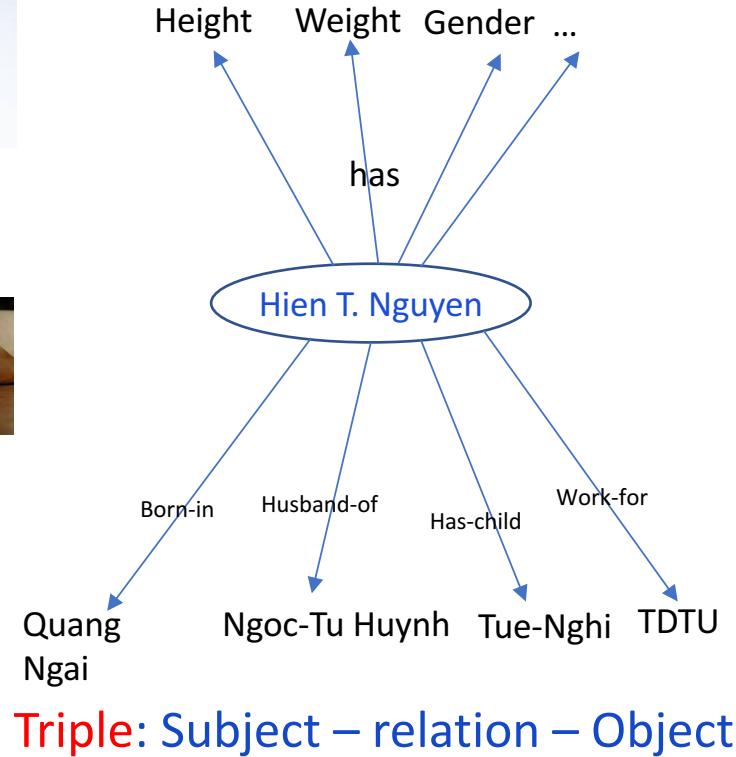
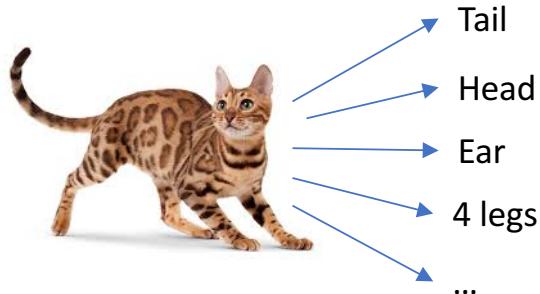
- Automatically learn features or representations at multiple levels of abstraction
- Use computational models composing of multiple processing layers



Deep Learning vs. Classsis ML



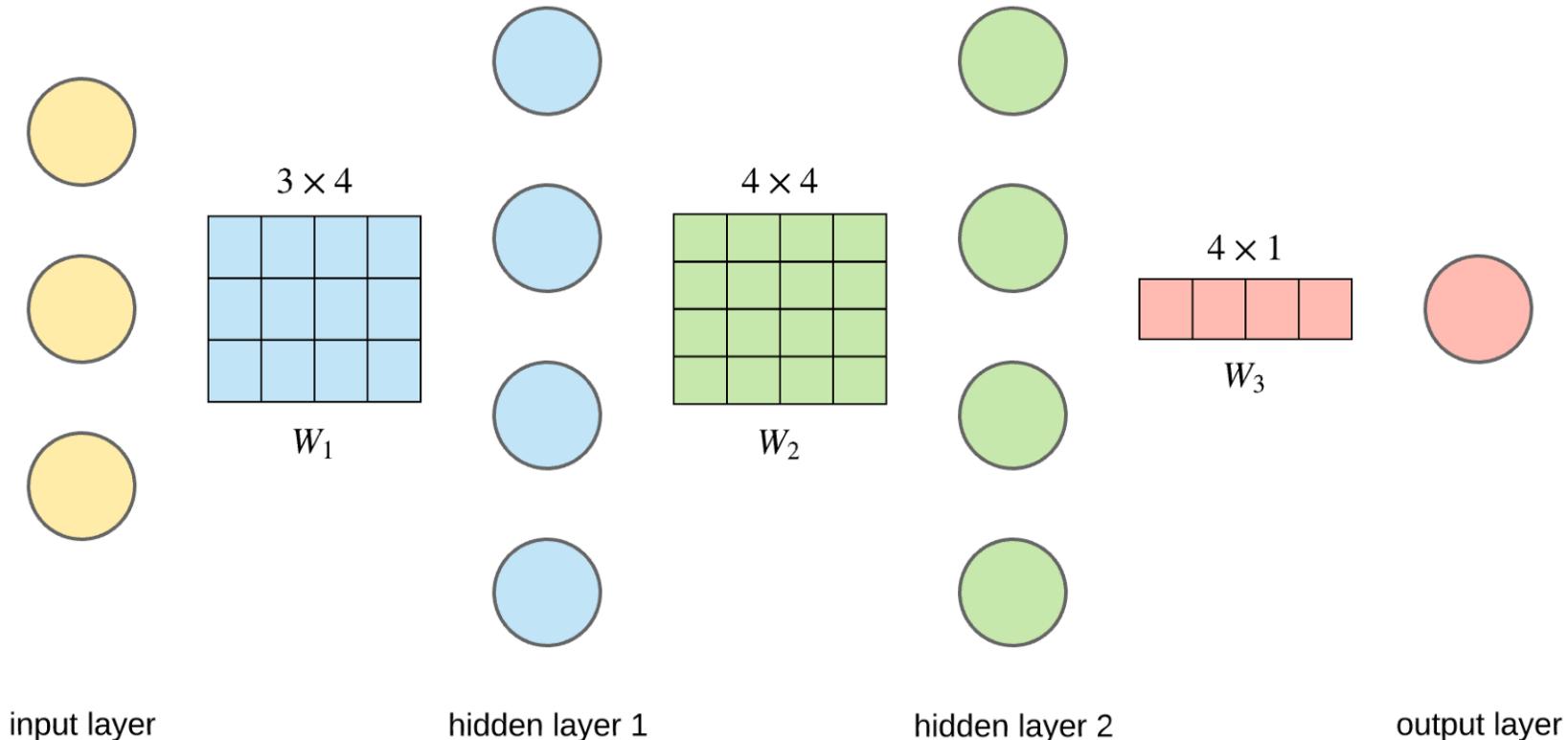
Representation Learning



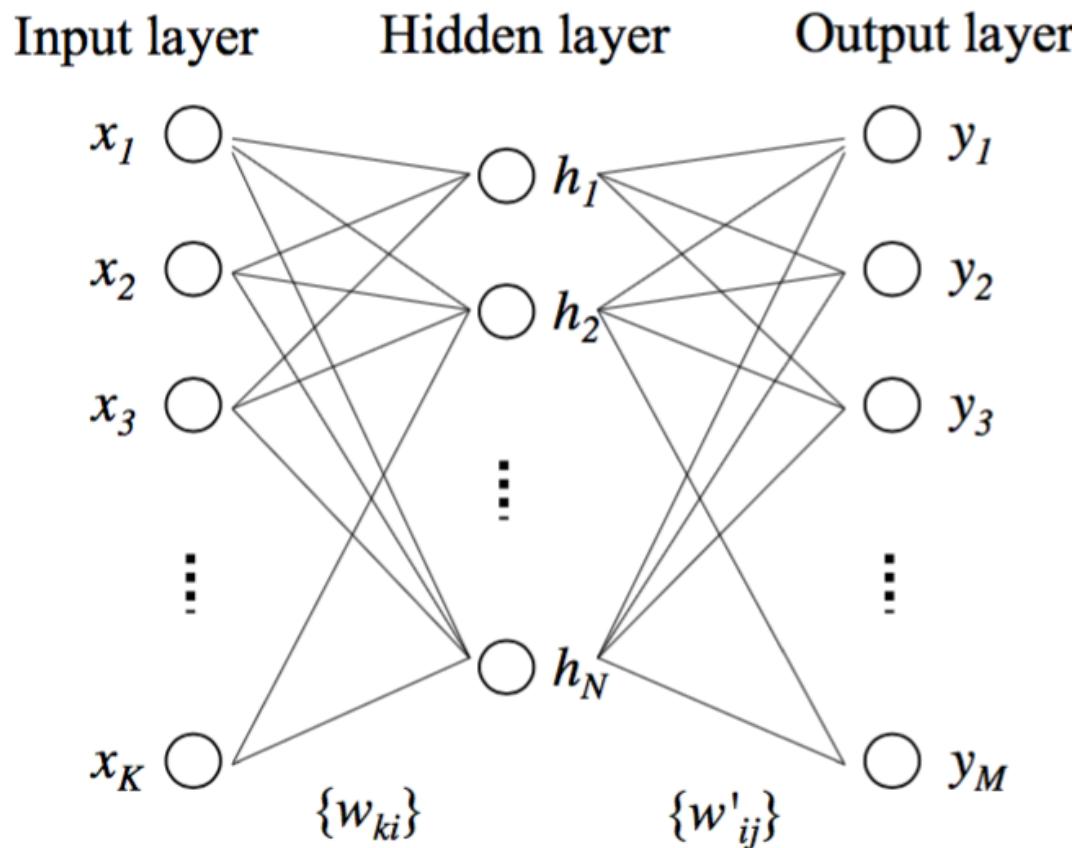
Deep Neural Networks

- Feed-Forward Neural Networks (**FFNN**)
- Convolutional Neural Networks (**CNN**)
- Recurrent Neural Networks (**RNN**)
 - Long-Short Term Memory (**LSTM**)
 - Gated Recurrent Unit (**GRU**)

Multi-Layer Networks



Multi-Layer Networks



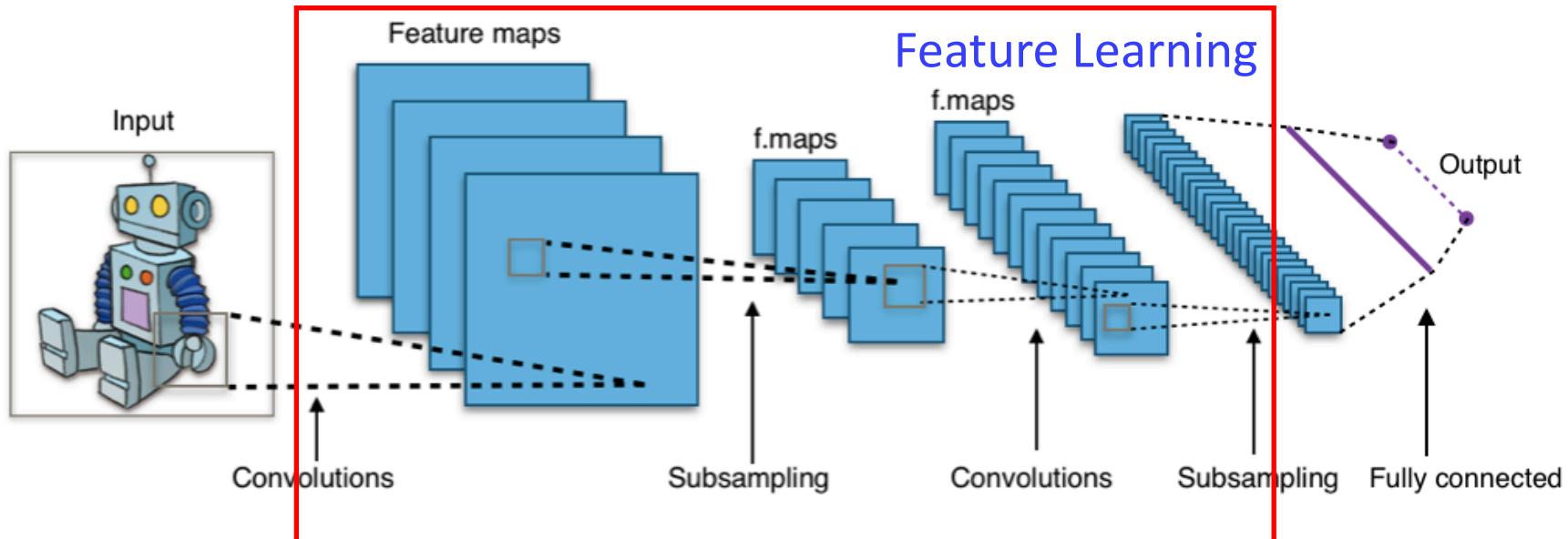
Loss Function

$$E(\mathbf{x}, \mathbf{t}, \mathbf{W}, \mathbf{W}') = \frac{1}{2} \sum_{j=1}^M (y_j - t_j)^2$$

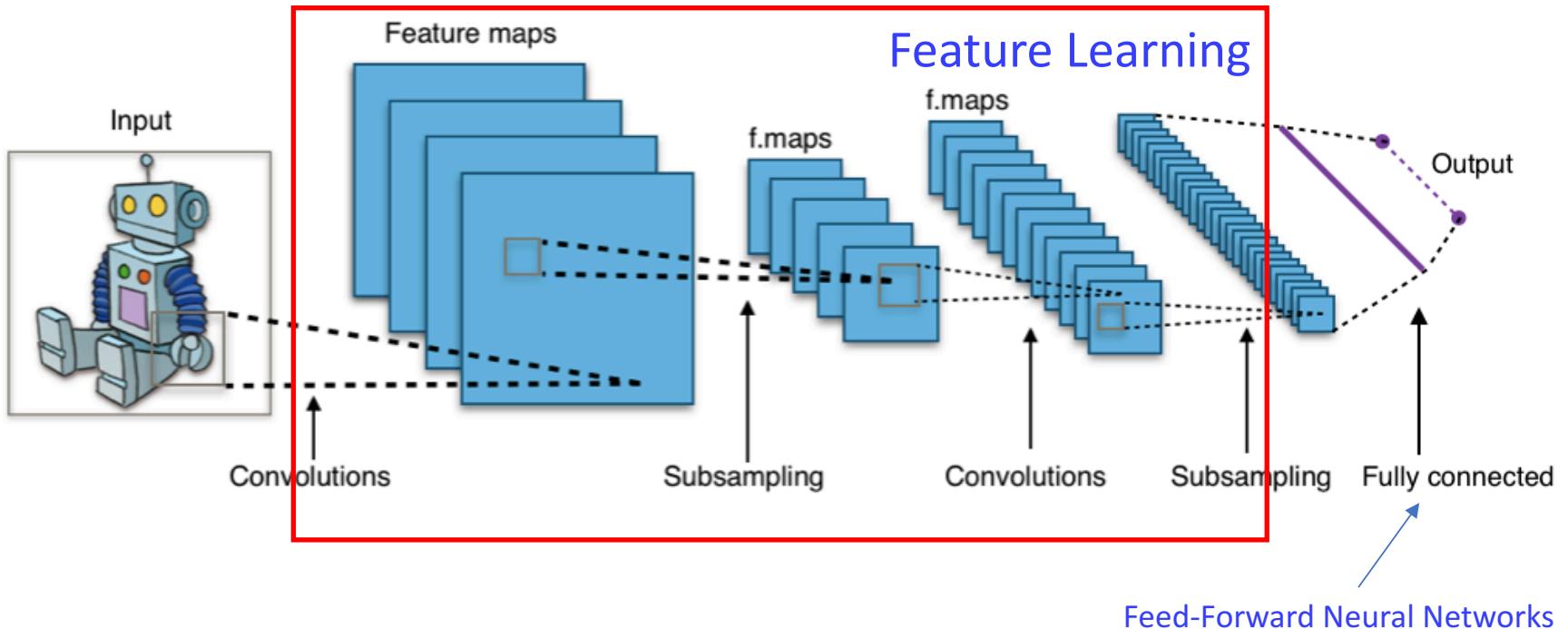
$$h_i = \sigma(u_i) = \sigma \left(\sum_{k=1}^K w_{ki} x_k \right)$$

$$y_j = \sigma(u'_j) = \sigma \left(\sum_{i=1}^N w'_{ij} h_i \right)$$

Convolutional Neural Networks



Convolutional Neural Networks



Convolutional Neural Networks

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Input

1	0	1
0	1	0
1	0	1

Filter / Kernel

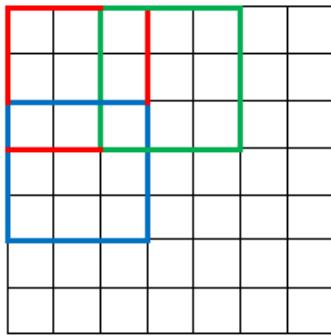
Convolutional Neural Networks

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

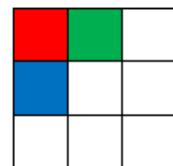
4		

Convolutional Neural Networks

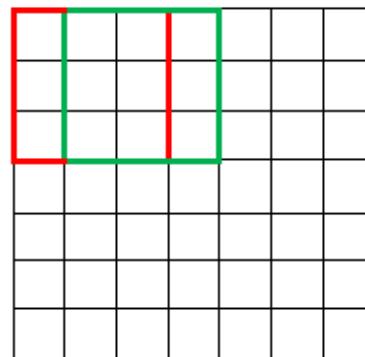
7 x 7 Input Volume



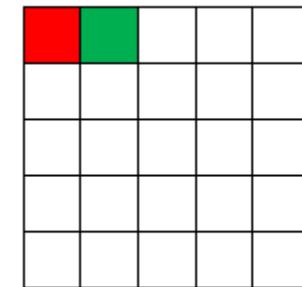
3 x 3 Output Volume



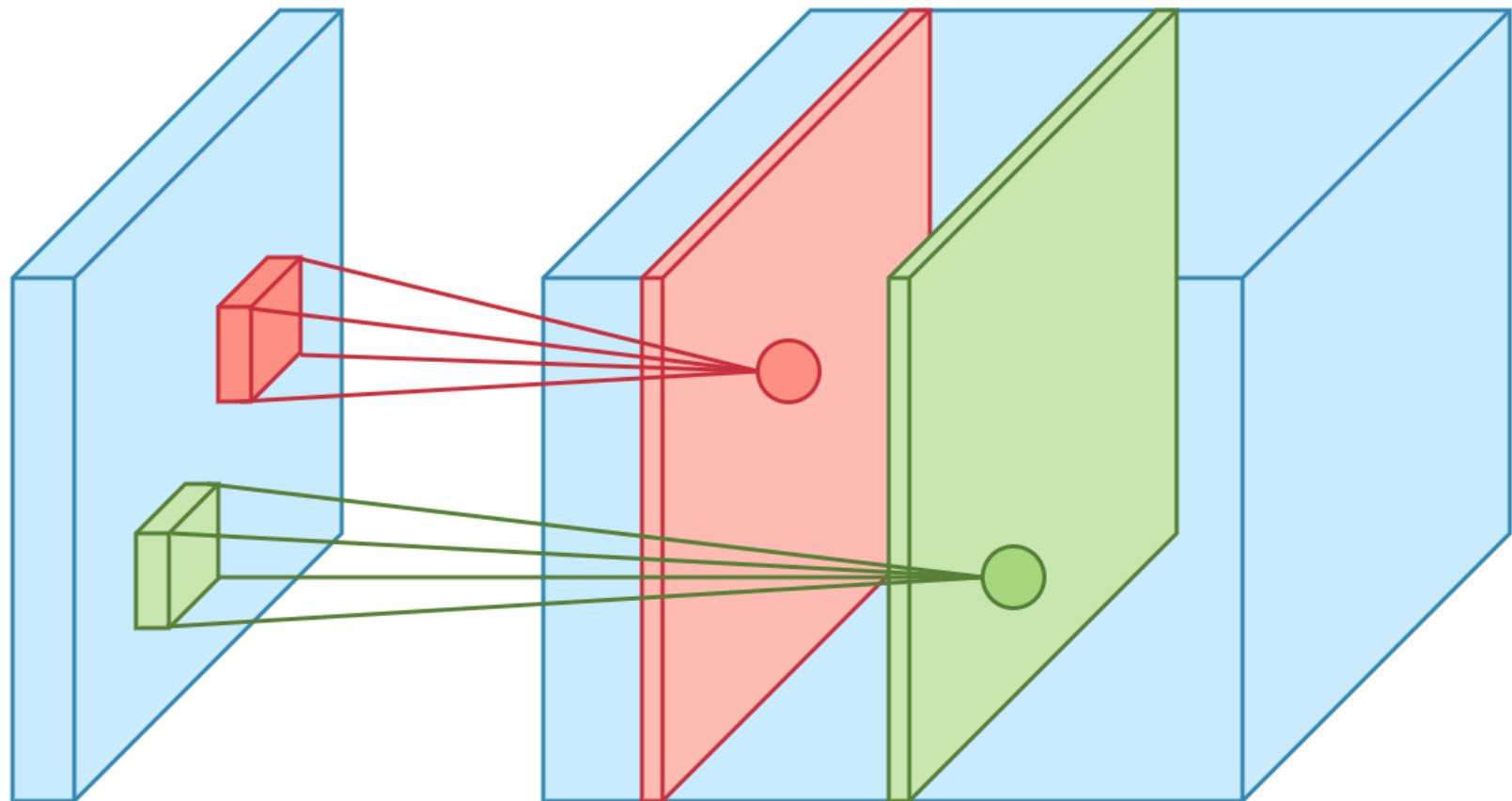
7 x 7 Input Volume



5 x 5 Output Volume



Convolutional Neural Networks



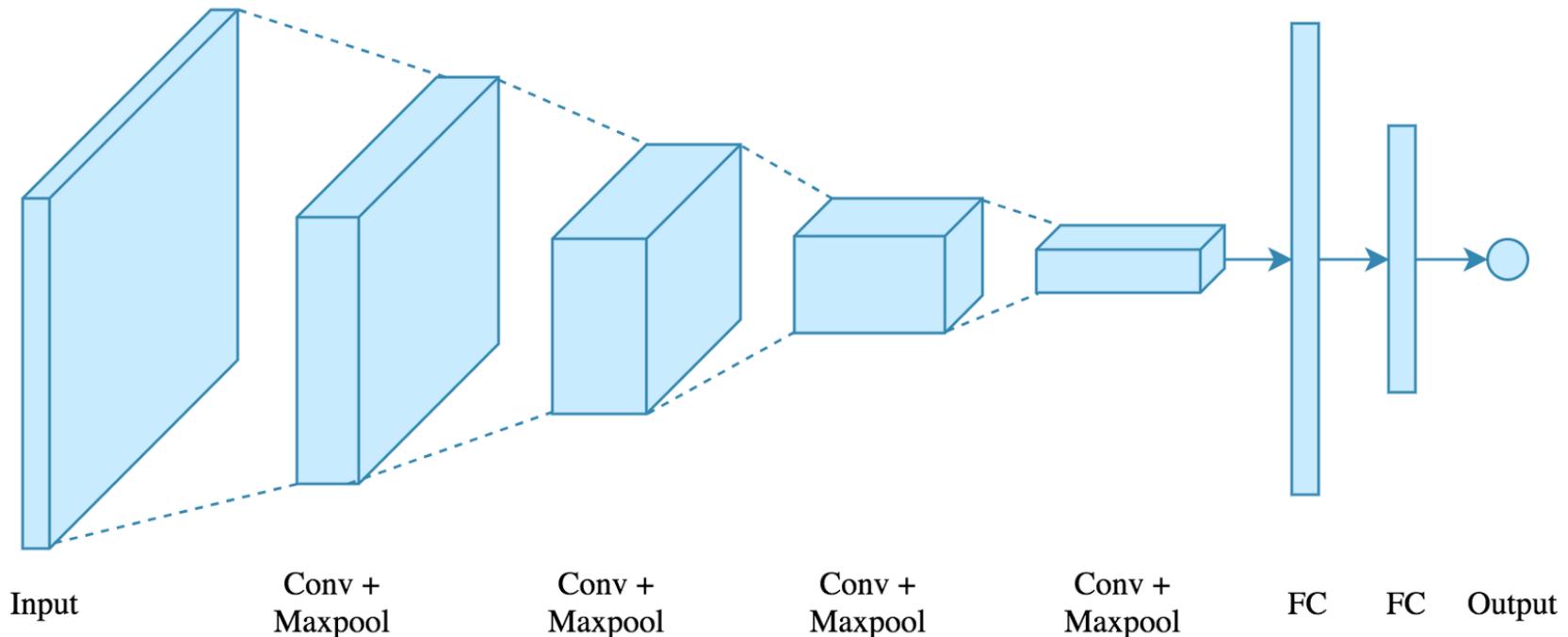
Convolutional Neural Networks

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

max pool with 2x2
window and stride 2

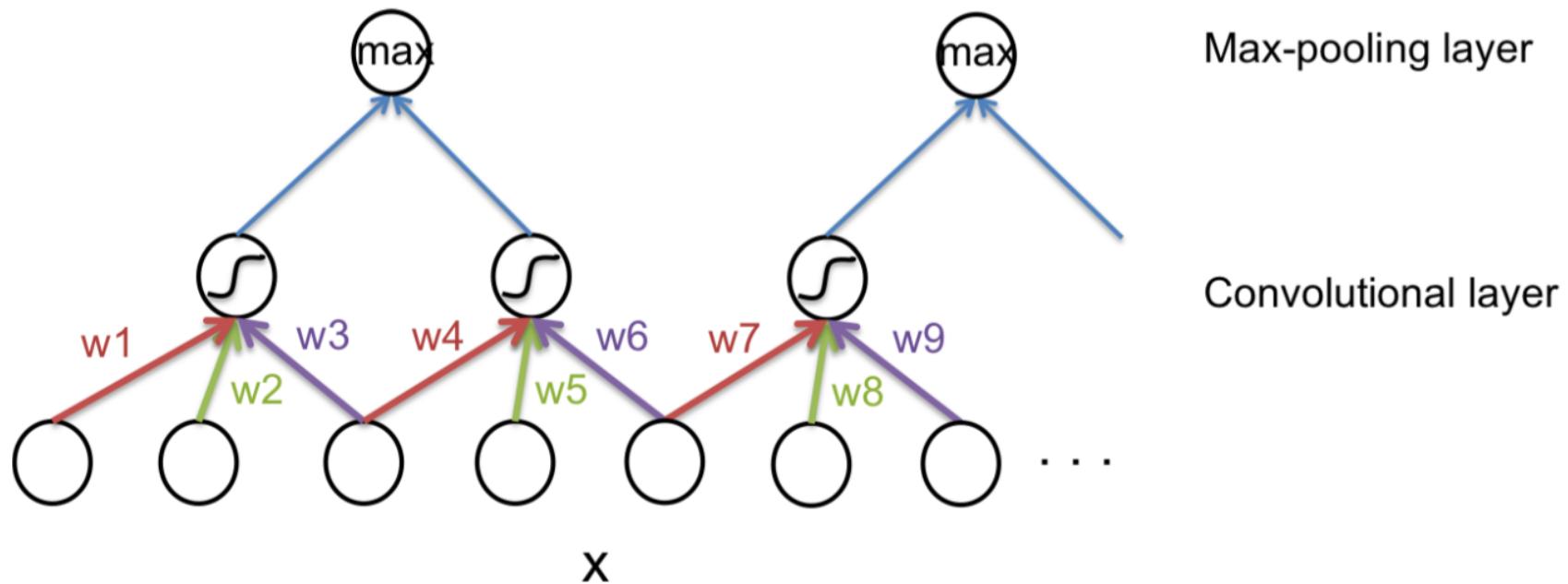
6	8
3	4

Convolutional Neural Networks



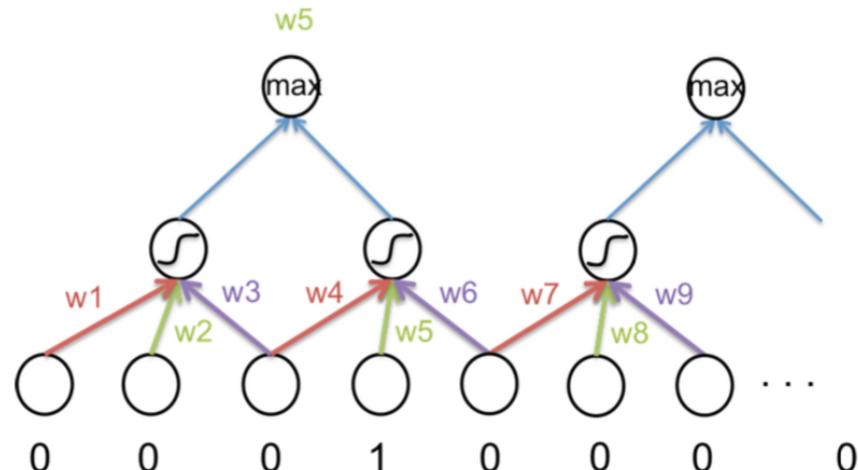
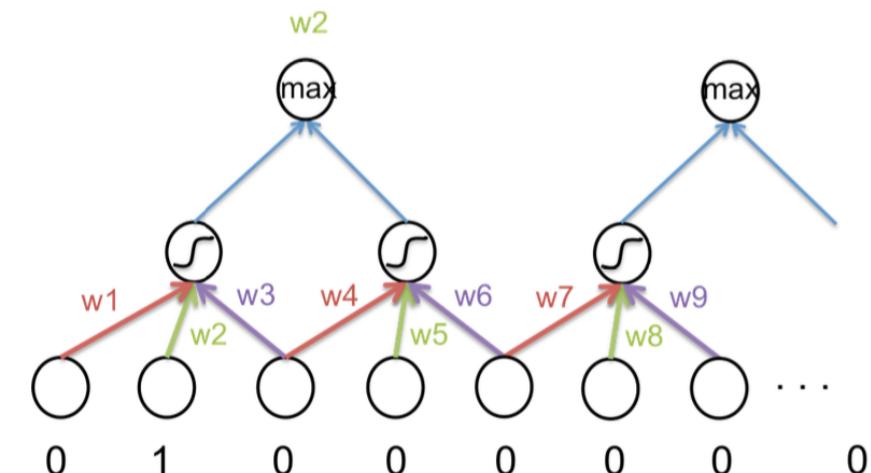
Convolutional Neural Networks

- Weight Sharing

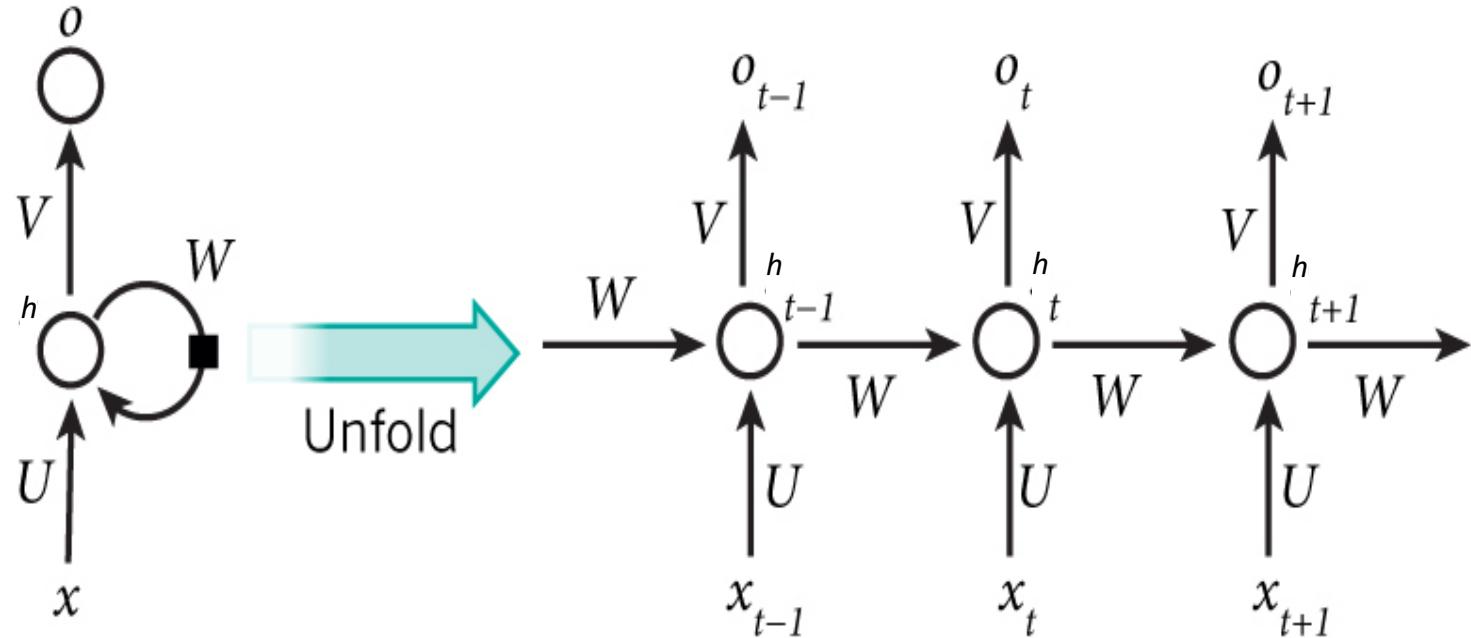


Convolutional Neural Networks

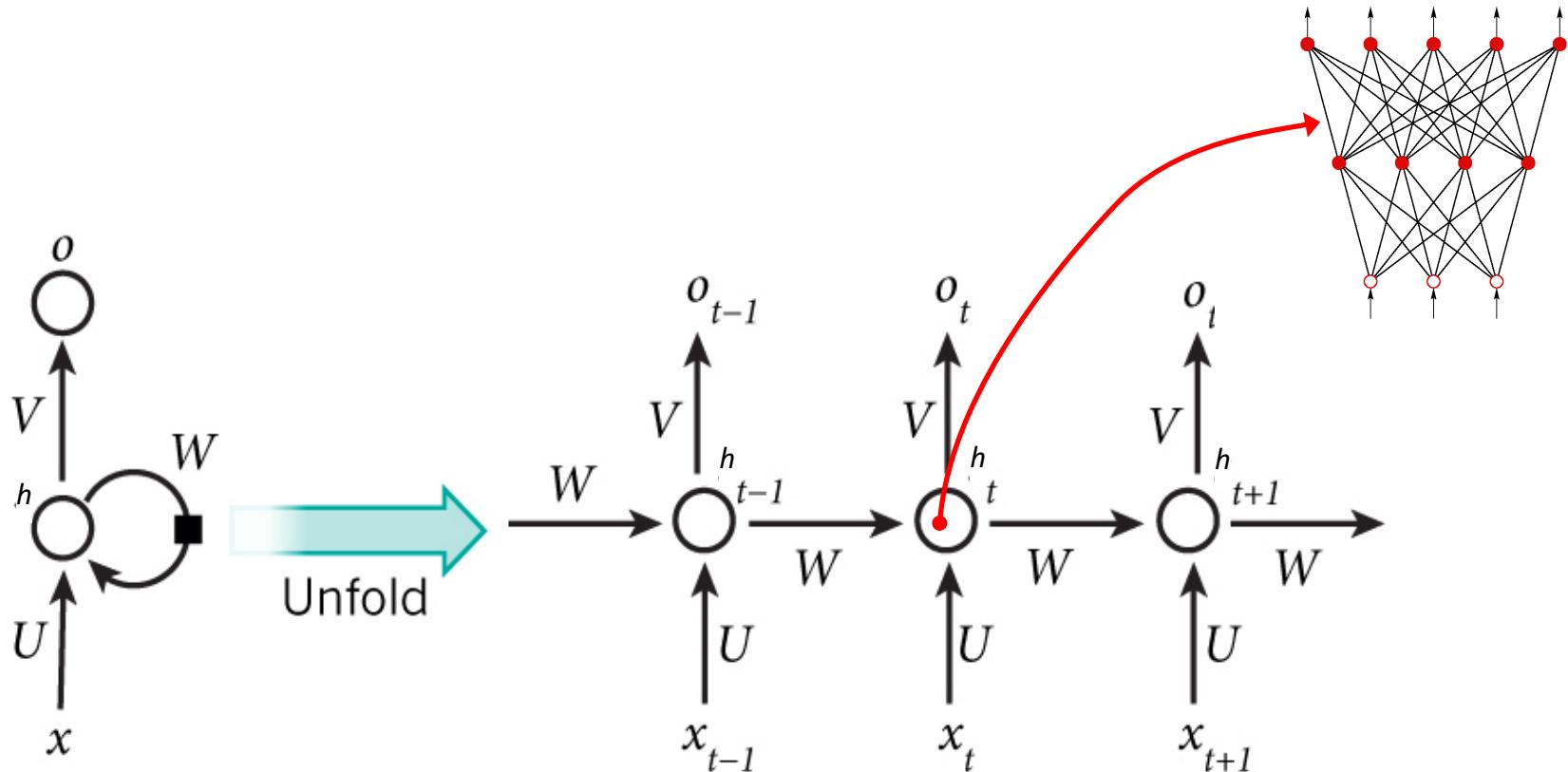
- The max-pooling neurons are invariant to shifts in the inputs



Recurrent Neural Networks

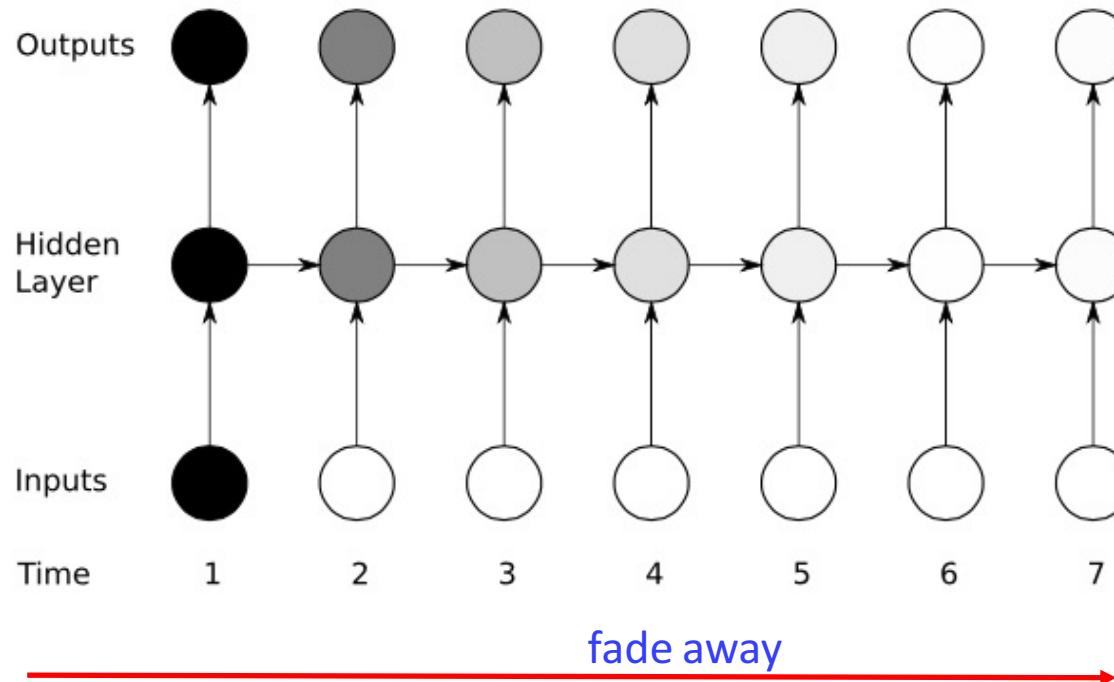


Recurrent Neural Networks



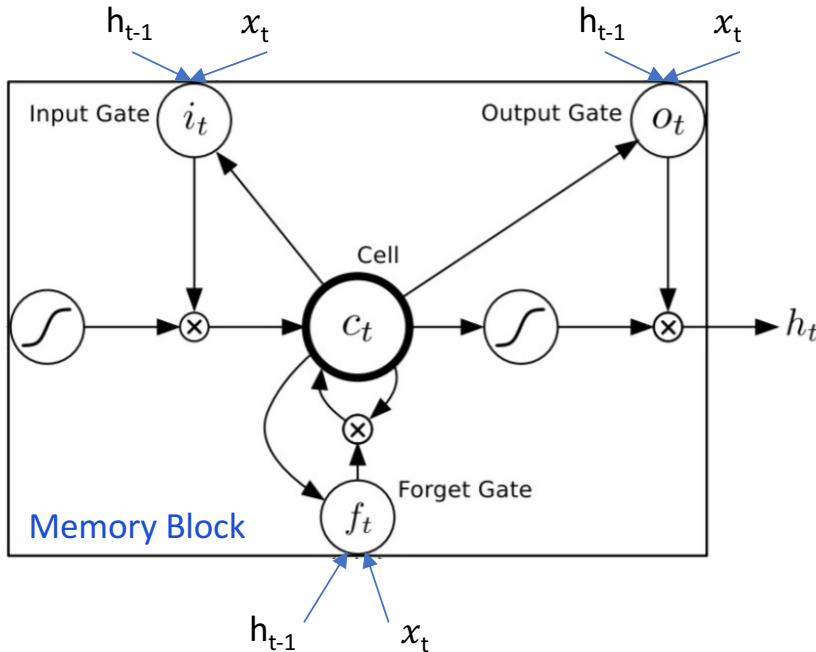
Recurrent Neural Networks

- Limitations of RNN: Vanishing and Exploding



Gated Recurrent Neural Networks

- Long-Short Term Memory (LSTM)

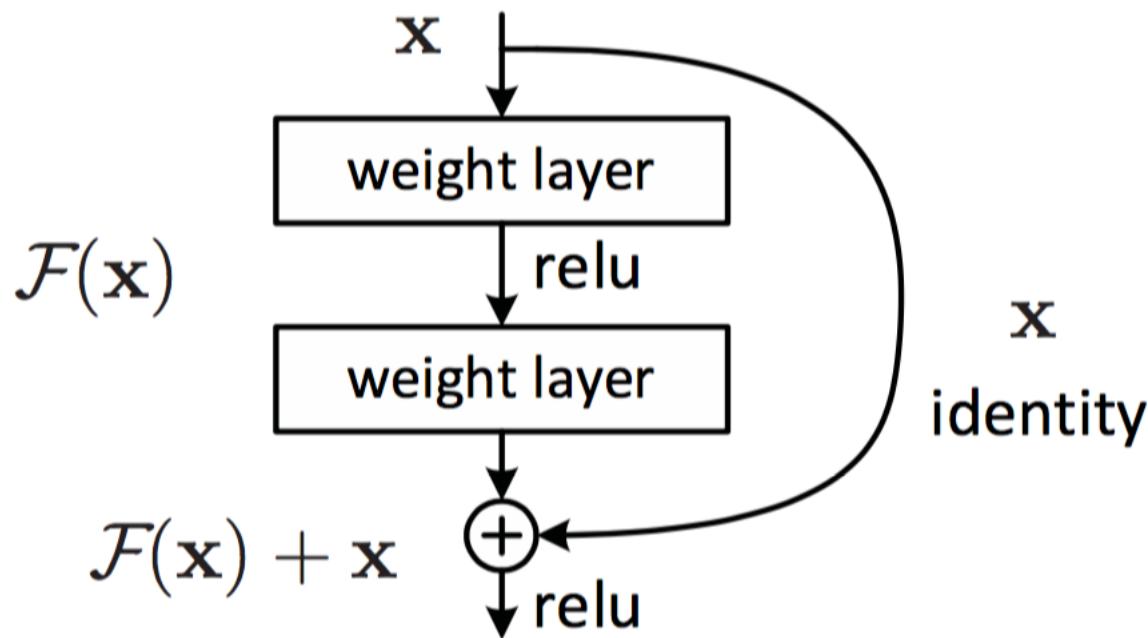


$$\begin{aligned} f_t &= \sigma(W_f[x_t, h_{t-1}] + b_f) \\ i_t &= \sigma(W_i[x_t, h_{t-1}] + b_i) \\ \tilde{c}_t &= \tanh(W_c[x_t, h_{t-1}] + b_c) \\ c_t &= f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \\ o_t &= \sigma(W_o[x_t, h_{t-1}] + b_o) \\ h_t &= o_t \odot \tanh(c_t) \end{aligned}$$

- Gated Recurrent Unit (GRU): a special case of LSTM

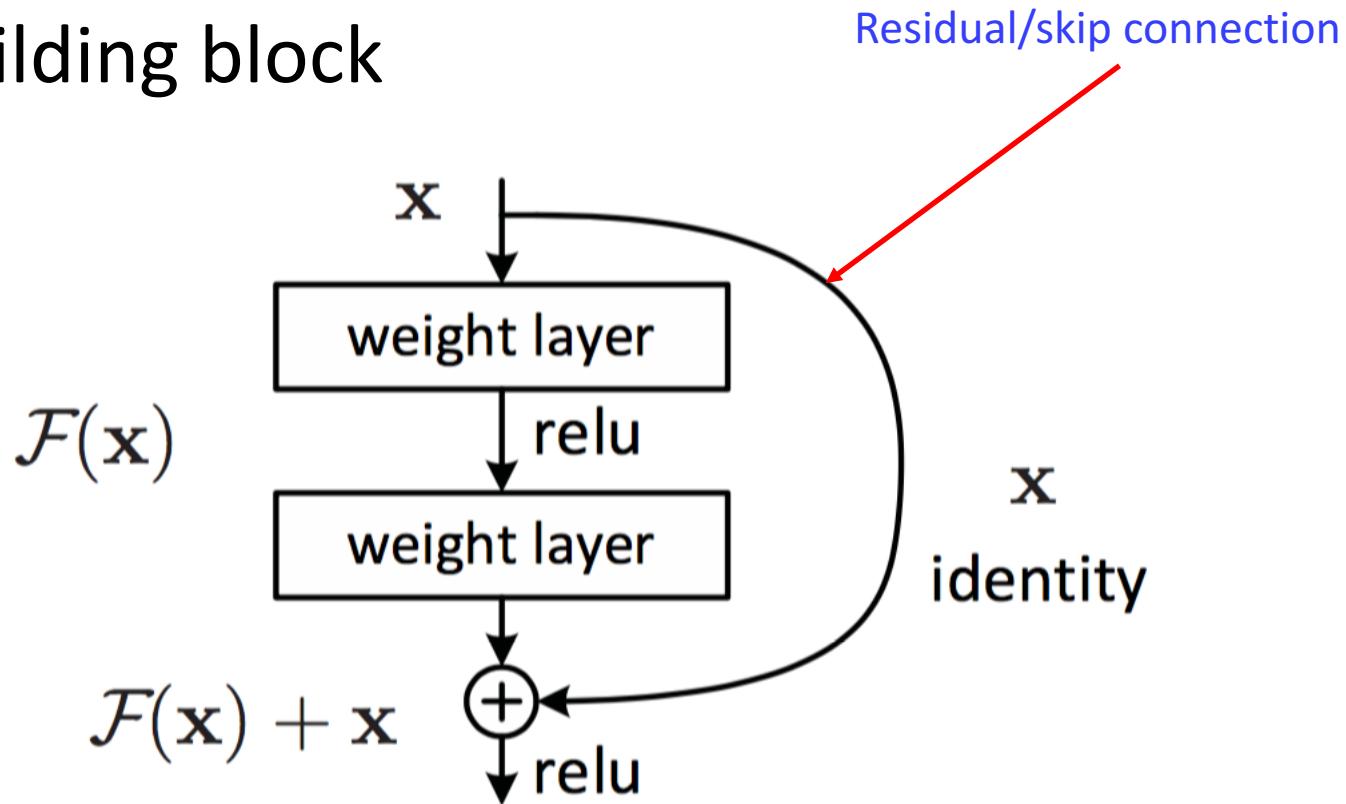
Residual Network

- Building block



Residual Network

- Building block



Residual Network

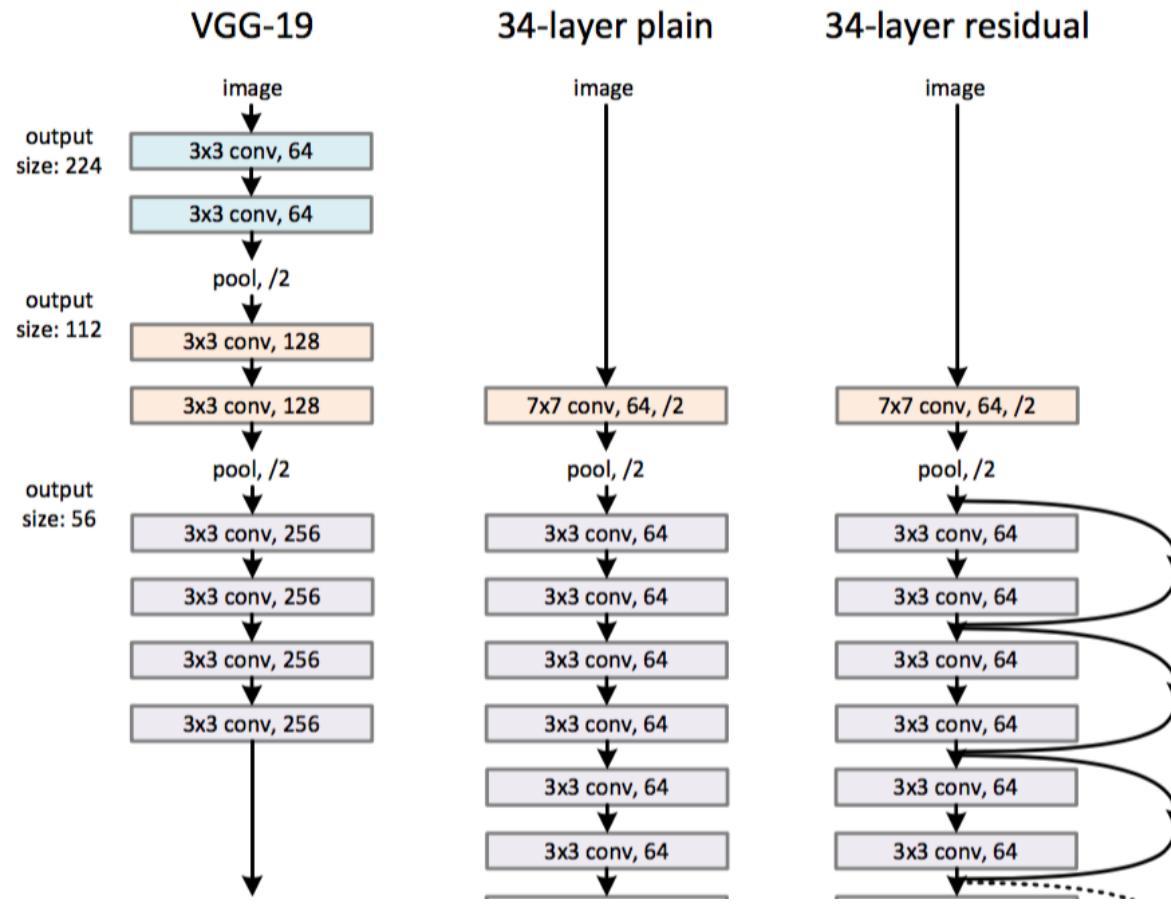


Image taken from He et. al. 2015

Attention Mechanism



A woman is throwing a frisbee in a park.

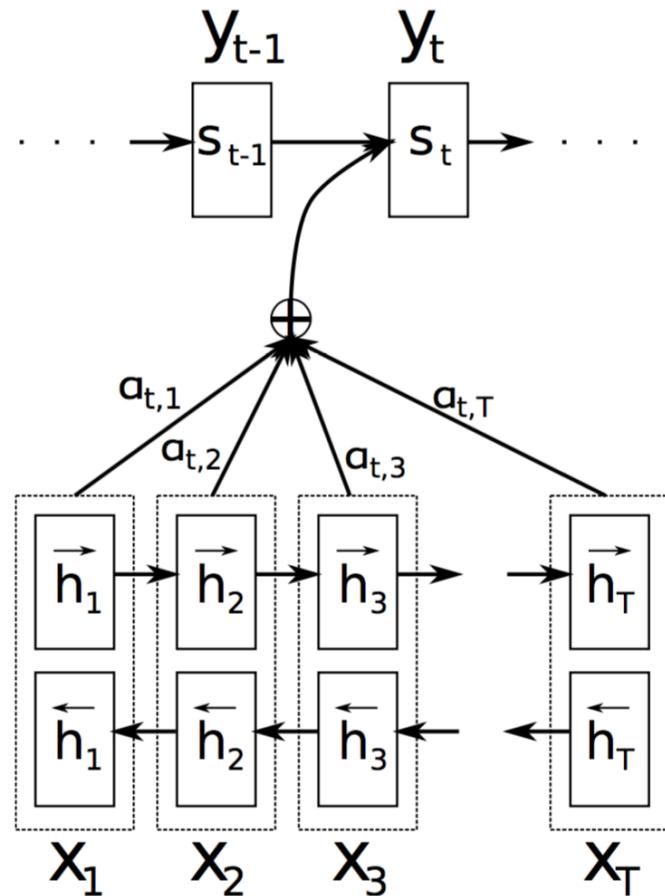


A stop sign is on a road with a mountain in the background.

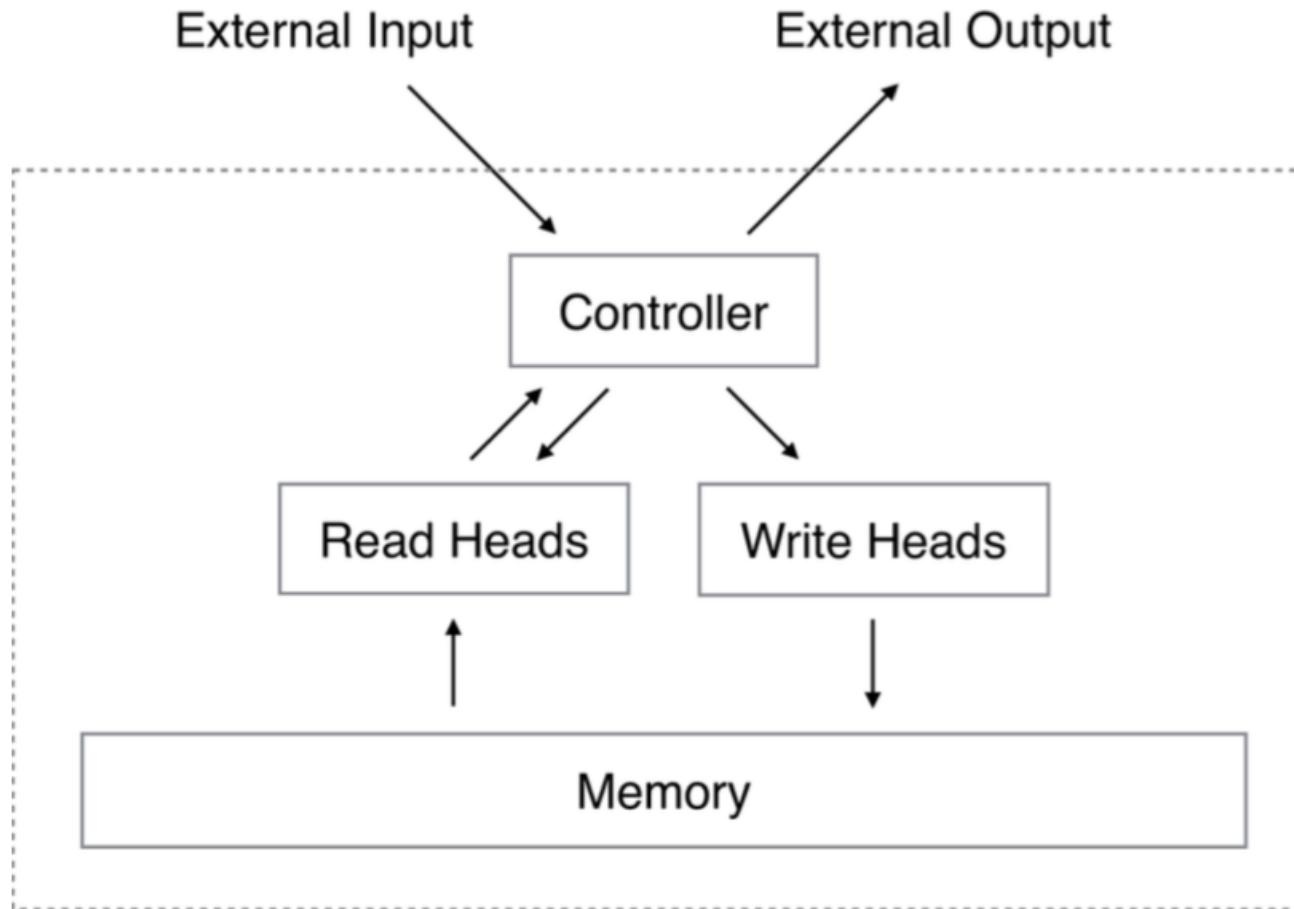


A dog is standing on a hardwood floor.

Attention Mechanism



Neural Turing Machine



Natural Language Processing

- Named Entity Recognition (NER)
 - identifying precisely entity mentions
 - classifying them into the corresponding entity types according to their context
 - E.g., Extract entities of a small set of types: people (PER), organizations (ORG), locations (LOC), and MISC.

“**Germany**’s representative to the **European Union**’s veterinary committee **Werner Zwingman** said on Wednesday consumers should ...”

Natural Language Processing

- IOB scheme for Sequence Labelling
 - NER: B-PER, I-PER, B-ORG, I-ORG, B-LOC, I-LOC, B-MISC, I-MISC, O

“**Germany**’s representative to the **European Union**’s veterinary . . .

B-LOC O O O B-ORG I-ORG O

- Part-of-Speech Tagging
- Chunking
- Semantic Role Labeling
- Etc.

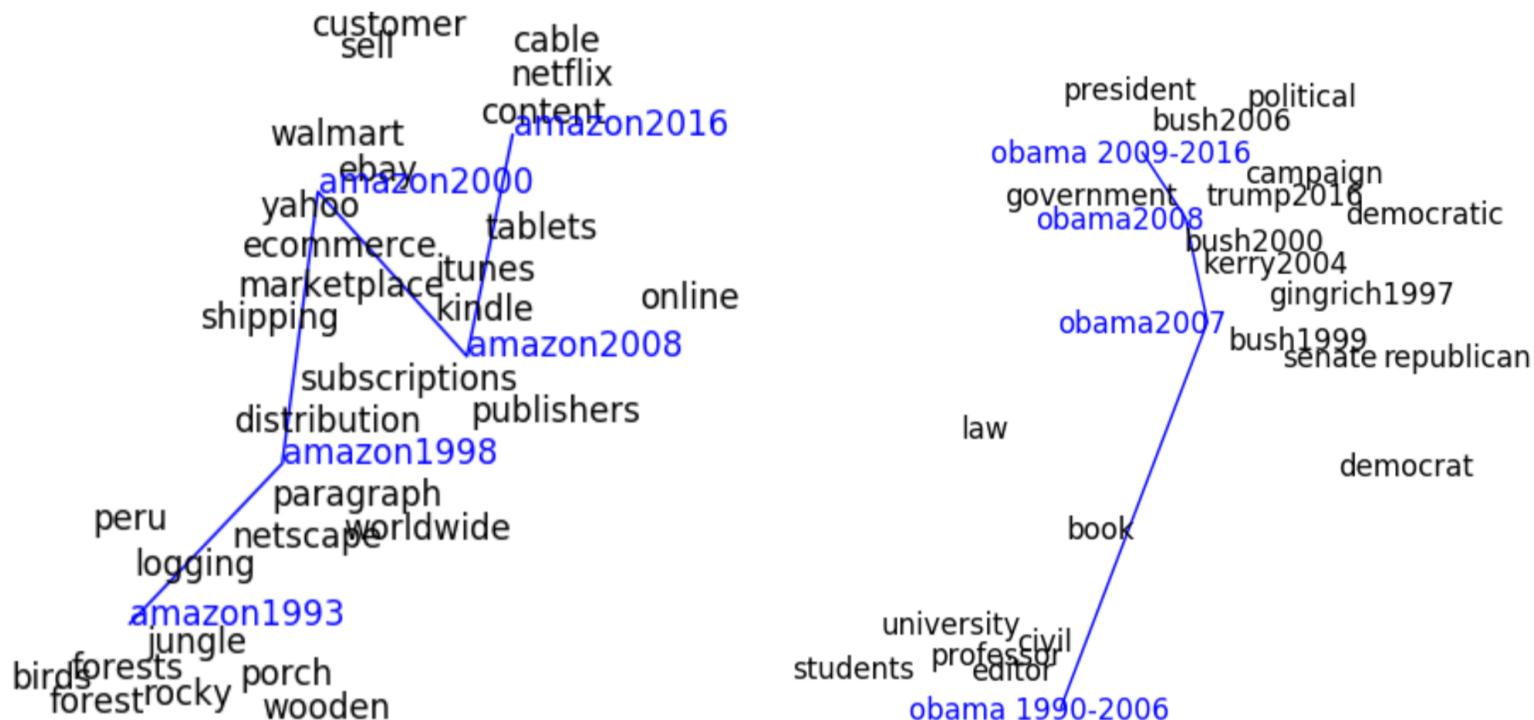
Natural Language Processing

- Sentiment Analysis: negative, neutral, positive
- Question Answering (QA)
 - Model
 - Retrieval-based
 - Generative
 - Domain: open and close
- Machine Translation: translate a sentence in a source language into a target language
- Visual Recognition
 - Image Captioning
 - Visual QA
- Speech Synthesis: generate audio from text

Distributed Representation

- Many-to-many relationship between two types of representations
 - A word is represented as a vector and a vector encodes words as its entries
- Multi-level Representation
 - Word embedding
 - Character-level representation of words
 - Mention and context representation
 - Knowledge representation of entities

Word Embedding



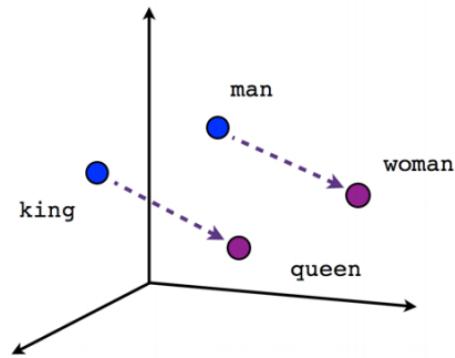
“You shall know a word by the company it keeps”

(J. R. Firth 1957: 11)

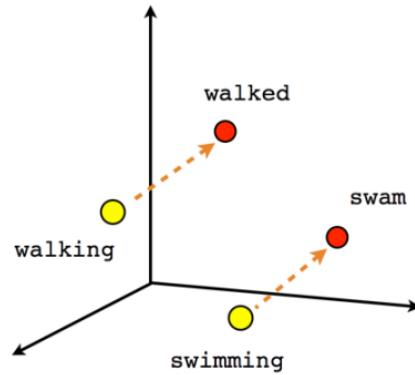
Word Embedding

- SENNA' word embedding
 - <https://ronan.collobert.com/senna/>
- Word2vec
 - <https://github.com/dav/word2vec>
- Glove
 - <https://nlp.stanford.edu/projects/glove/>

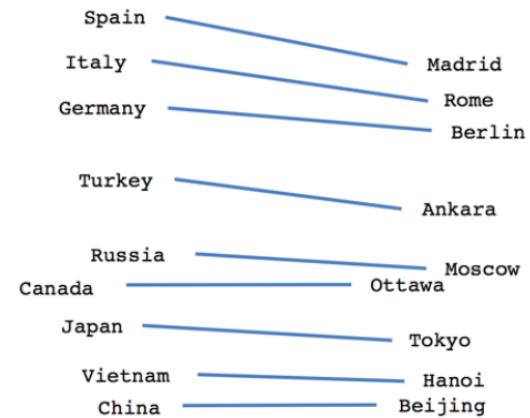
Word Embedding



Male-Female



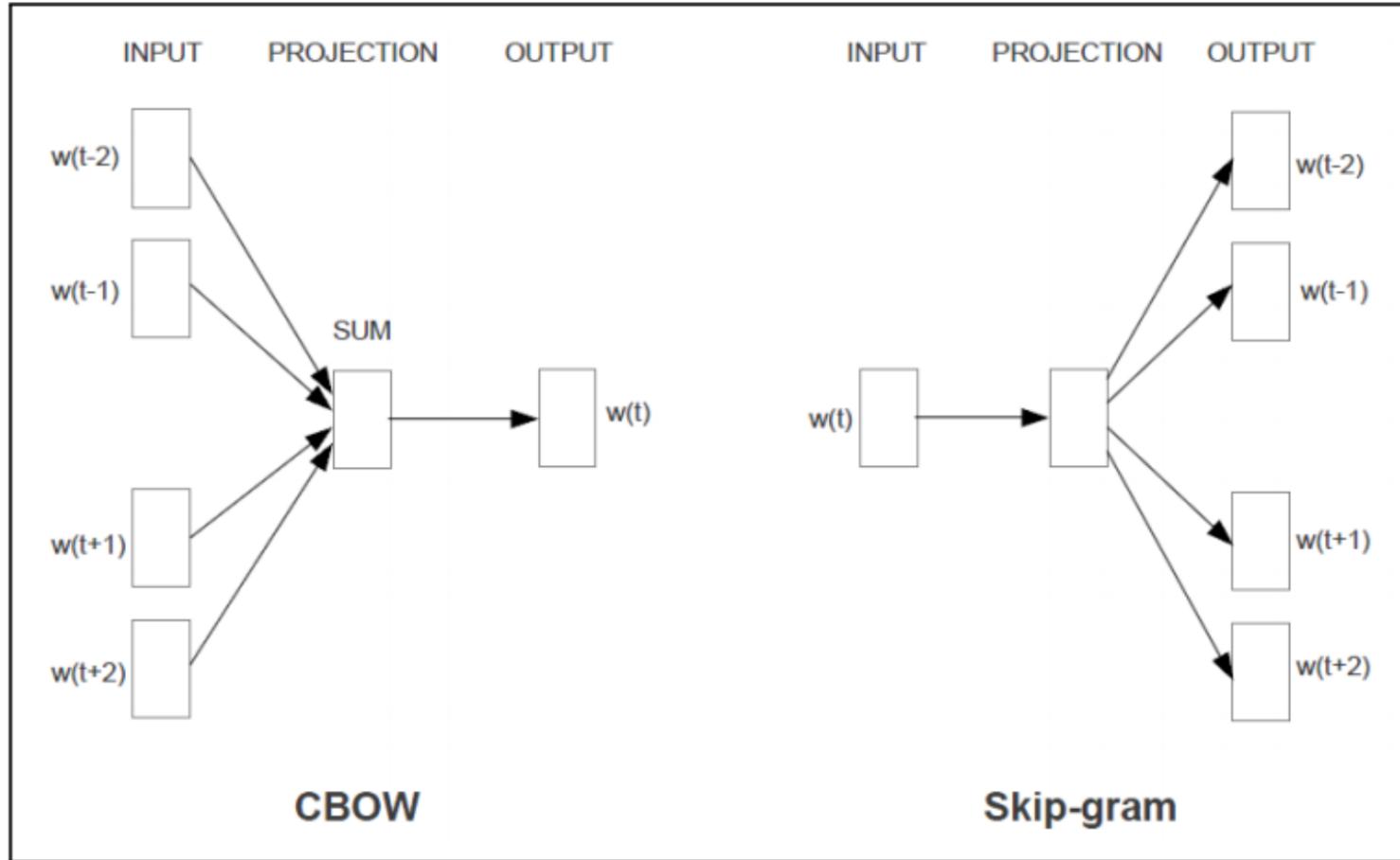
Verb tense



Country-Capital

'king' - 'man' + 'woman' is close to 'queen'

Word2vec

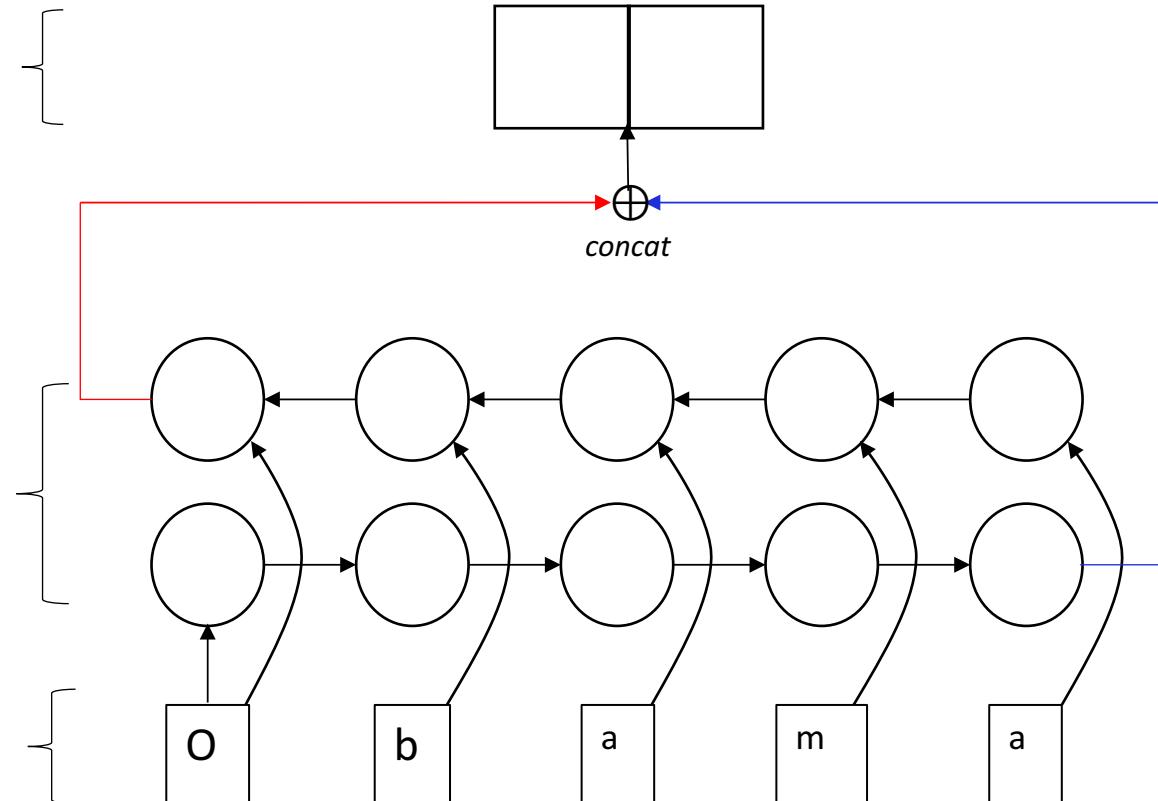


Char-level Representation of words using LSTM

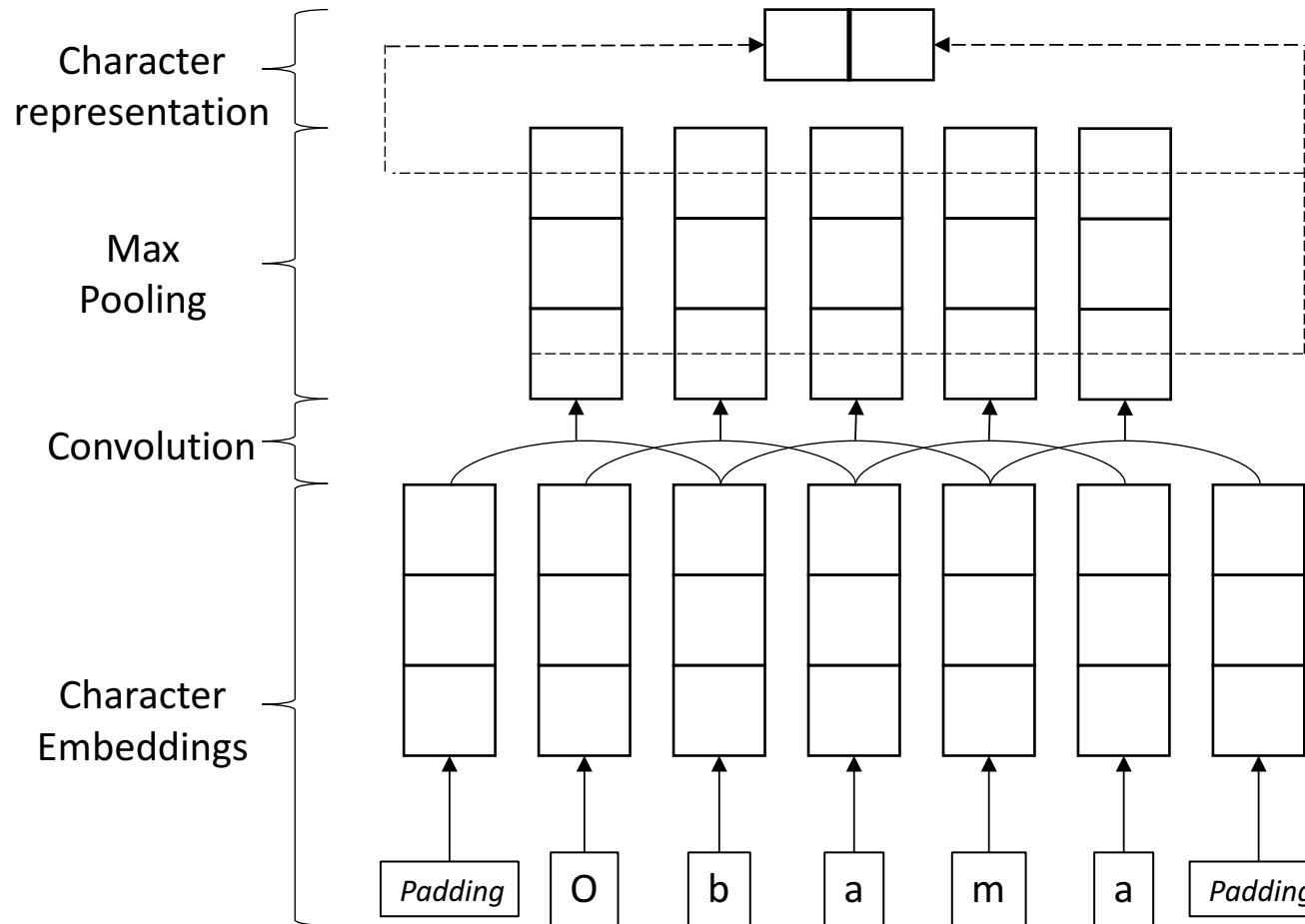
Character representation

Bidirectional RNN

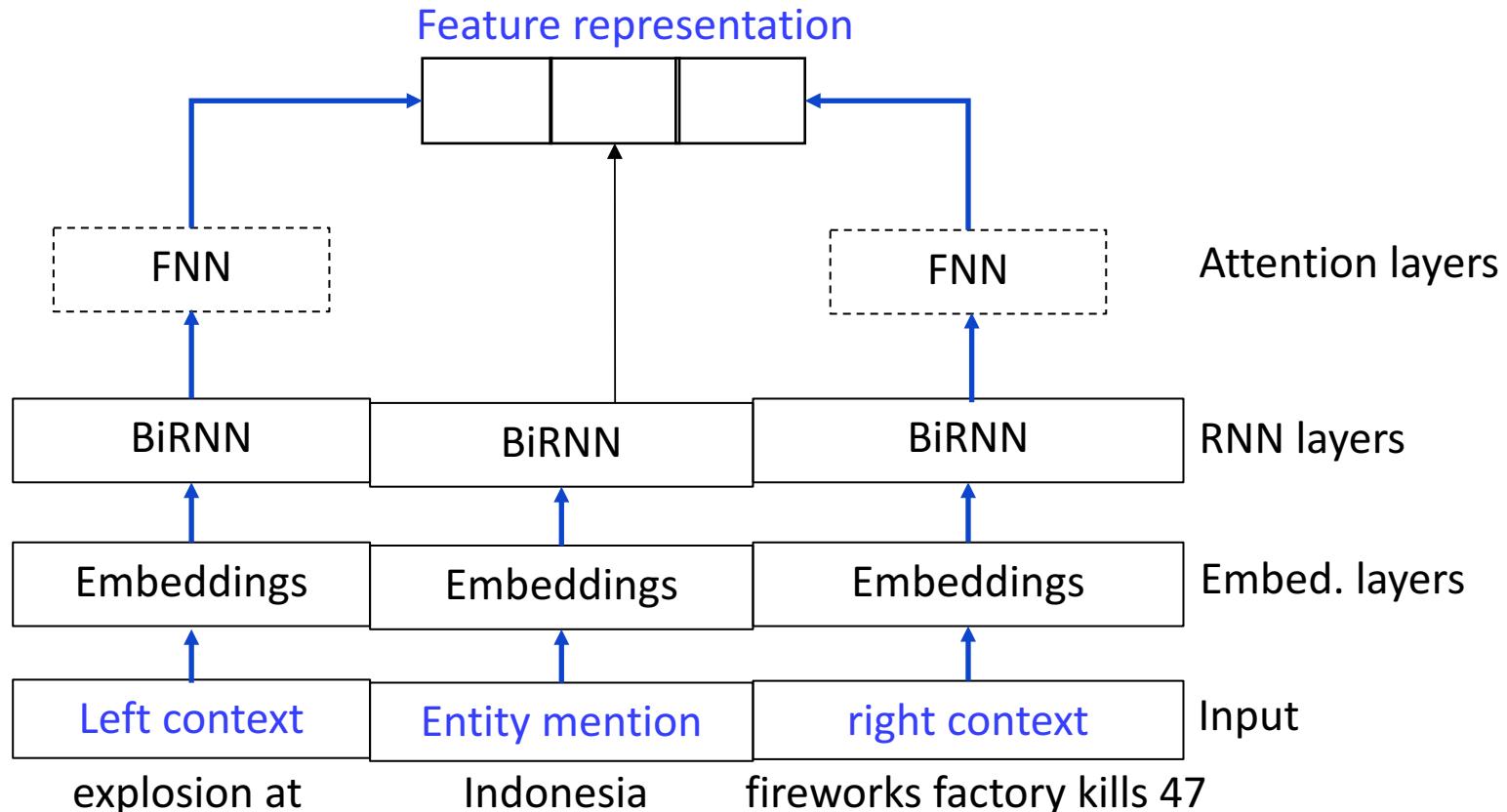
Character embeddings



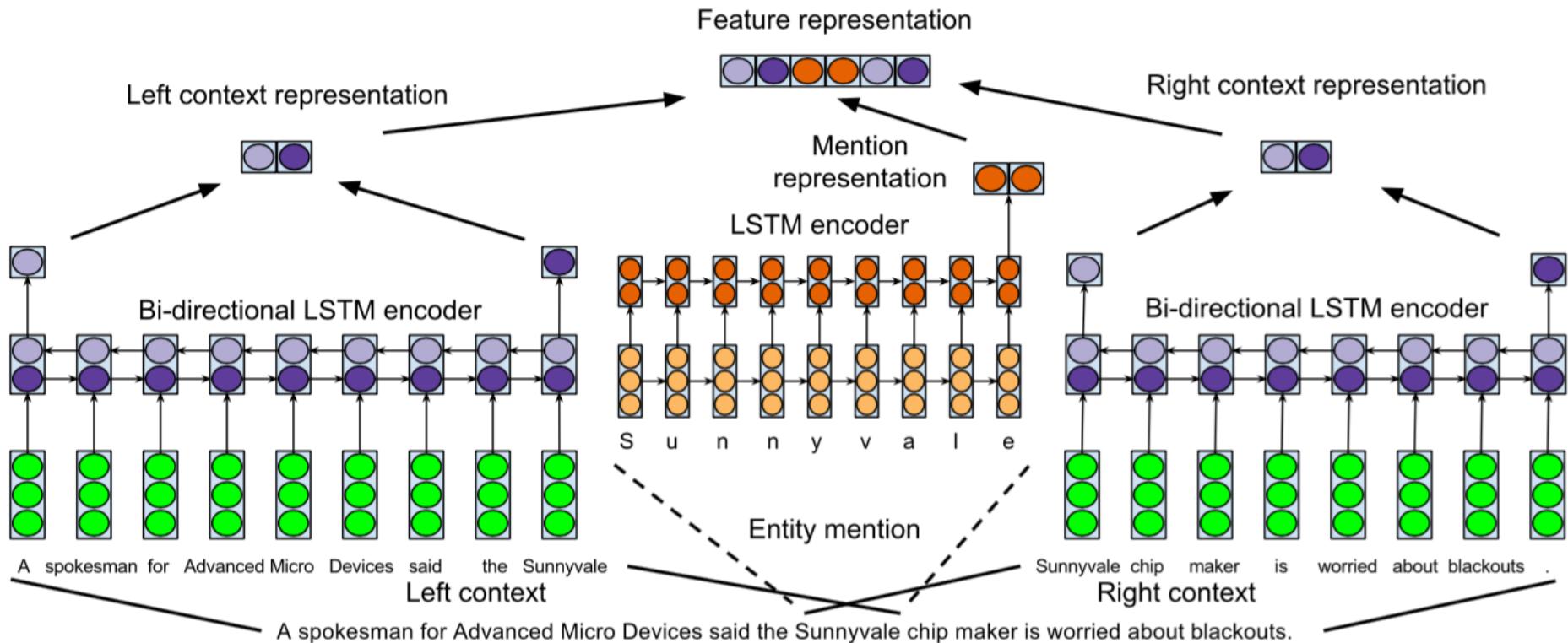
Char-level Representation of words Using CNN



Mention and Context Representation



Mention and Context Representation

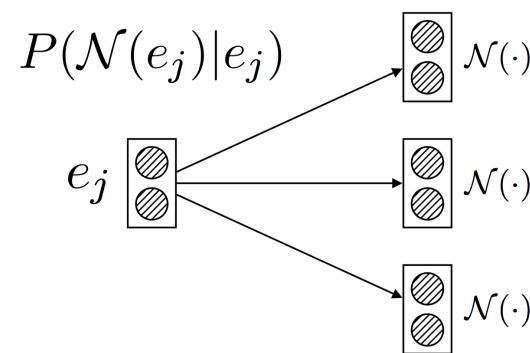


Knowledge representation of entities

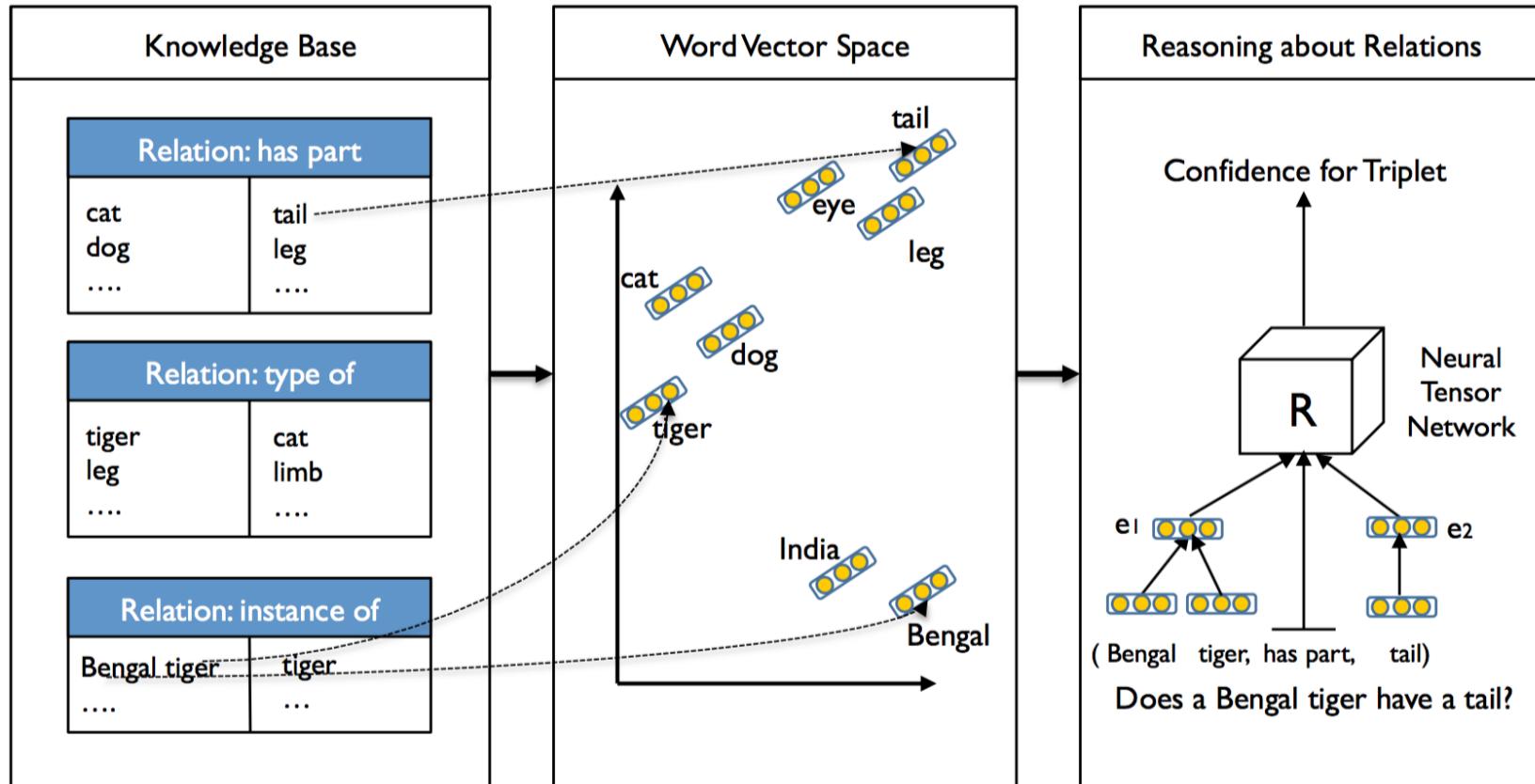
- Entities sharing more common neighbors tend to have similar representations
- Given a \mathcal{KB} , let $\mathcal{E} = \{e_j\}$ be a set of entities
 - maximizing the log probability

$$\mathcal{L}_e = \sum_{e_j \in \mathcal{E}} \log P(\mathcal{N}(e_j) | e_j)$$

Neighbor Entities



Knowledge representation of entities



Sentic LSTM

$$f_i = \sigma(W_f[x_i, h_{i-1}, \mu_i] + b_f)$$

$$I_i = \sigma(W_I[x_i, h_{i-1}, \mu_i] + b_I)$$

$$\tilde{C}_i = \tanh(W_C[x_i, h_{i-1}] + b_C)$$

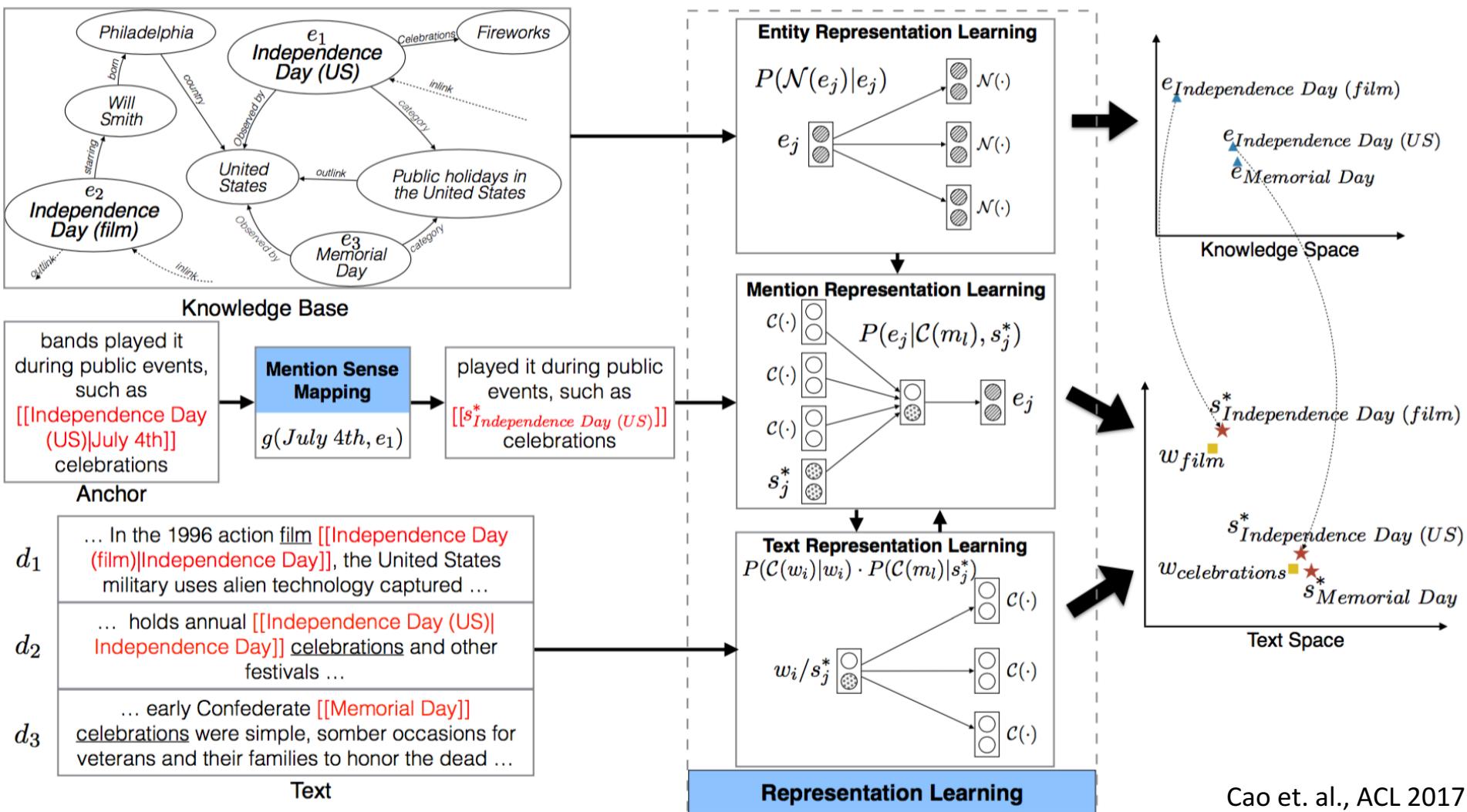
$$C_i = f_i * C_{i-1} + I_i * \tilde{C}_i$$

$$o_i = \sigma(W_o[x_i, h_{i-1}, \mu_i] + b_o)$$

$$o_i^c = \sigma(W_{co}[x_i, h_{i-1}, \mu_i] + b_{co})$$

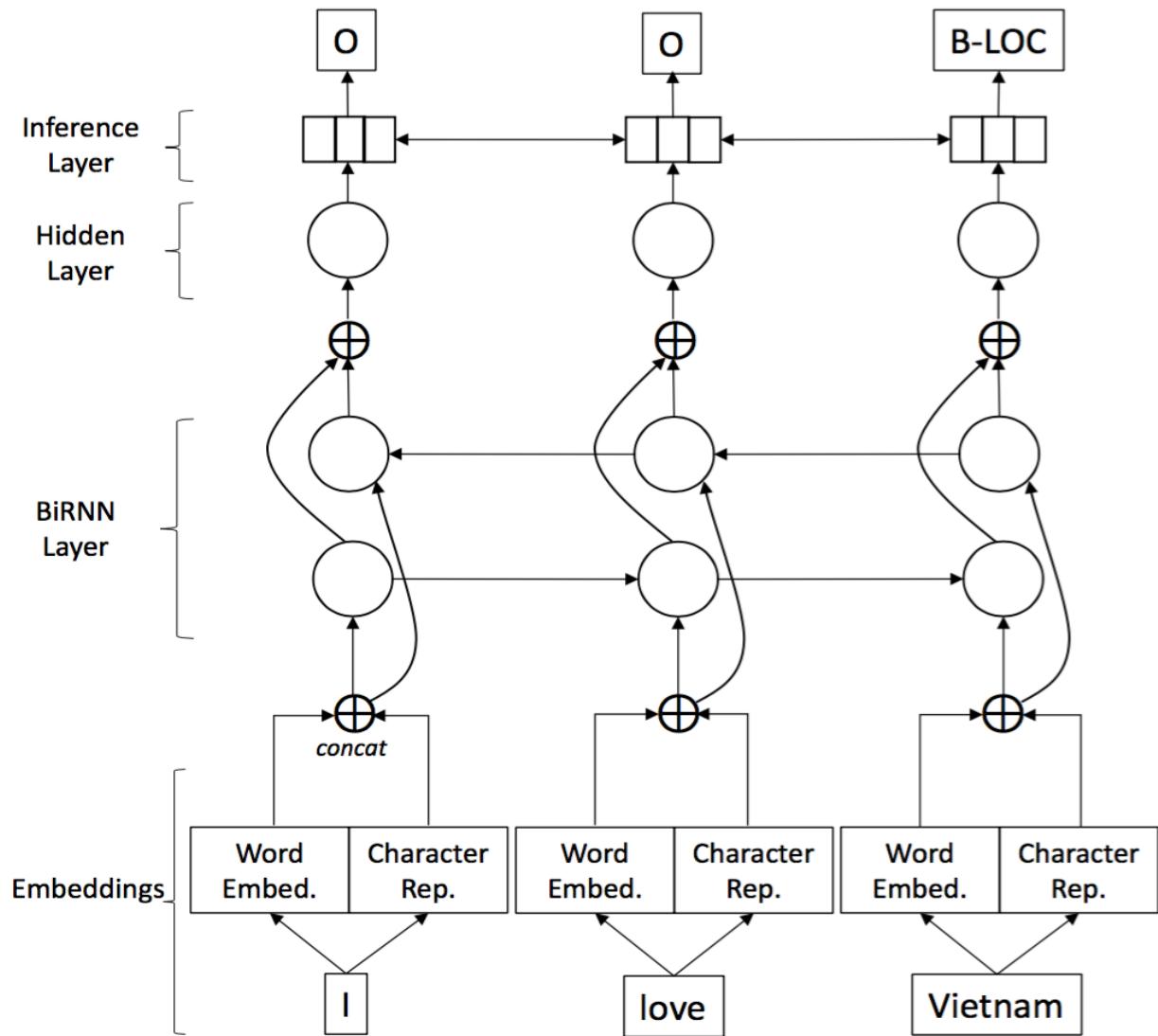
$$h_i = o_i * \tanh(C_i) + o_i^c * \tanh(W_c \mu_i)$$

Unify Text and Knowledge into a semantic space



State-of-the-Art NER Systems

- BiRNN-CRF
with some
variants



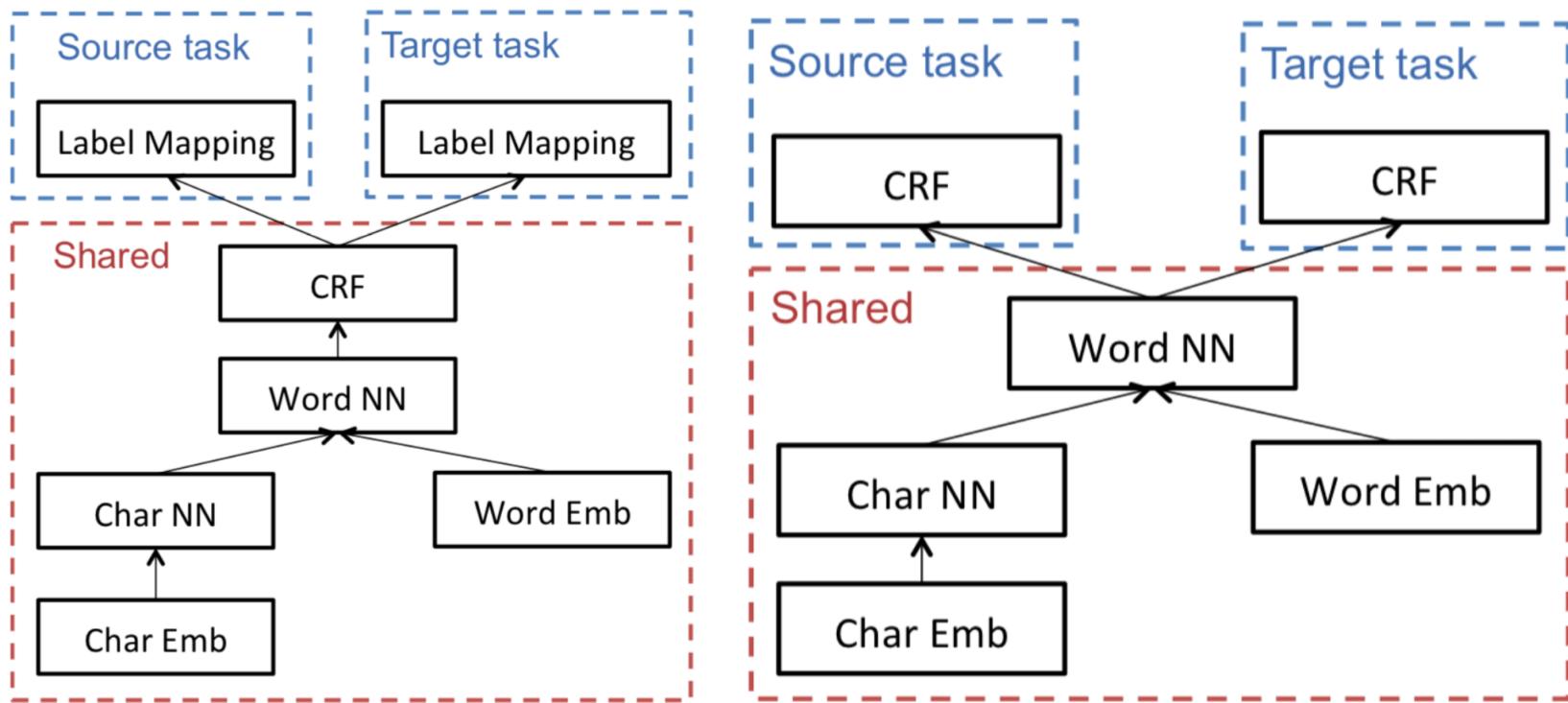
State-of-the-Art NER Systems

- English CoNLL 2003 dataset

	No.	Method	F1-score
BiRNN - CRF	1	BiLSTM-CRF (Huang <i>et. al.</i> , 2015) [31]	90.10%
	2	BiLSTM-CNN-CRF (Ma <i>et. al.</i> , 2016) [23]	91.21%
	3	BiLSTM-CRF, char-BiLSTM (Lample <i>et. al.</i> , 2016) [21]	90.94%
	4	BiLSTM-CRF, char-CNN (Chiu <i>et. al.</i> , 2016) [19]	91.62%
	5	Neural Reranking (Yang <i>et. al.</i> , 2017) [18]	91.62%
BiRNN-CRF+	6	Transfer Learning with RNN (Yang <i>et. al.</i> , 2017) [17]	91.26%
FOFE+FNN	7	FOFE+FNN (Xu <i>et. al.</i> , 2017) [16]	90.85%
Res-RNN	8	3 Res-RNN + foreLM + backLM + bias (Tran <i>et. al.</i> , 2017) [12]	91.69%
ID-CNN	9	ID-CNN (Strubell <i>et. al.</i> , 2017) [10]	90.65%
BiRNN-CRF + LM	10	TagLM (Peters <i>et. al.</i> , 2017) [8]	91.93%
	11	LM-BiLSTM-CRF (Liu <i>et. al.</i> , 2017) [7]	91.71%

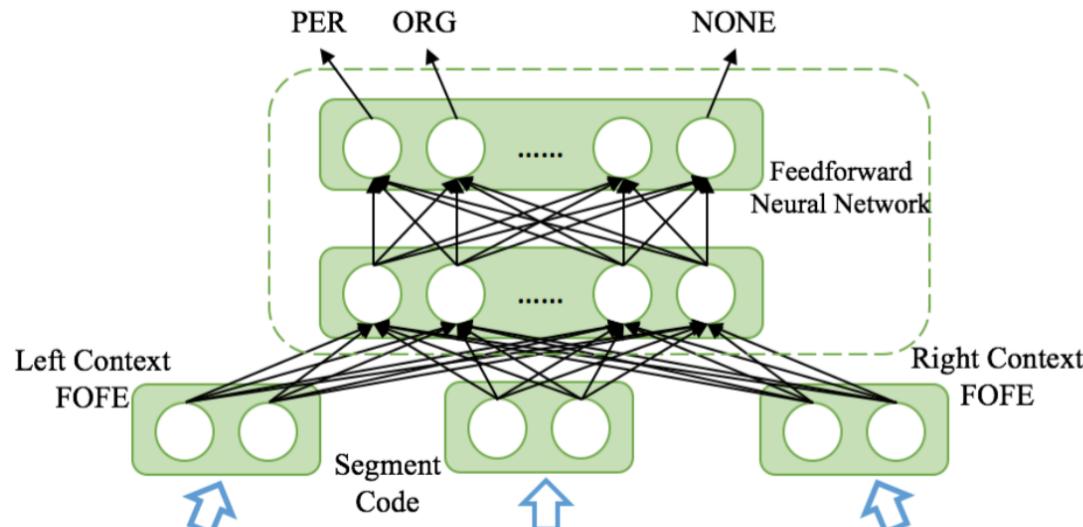
State-of-the-Art CGER Systems

■ Transfer Learning



State-of-the-Art CGER Systems

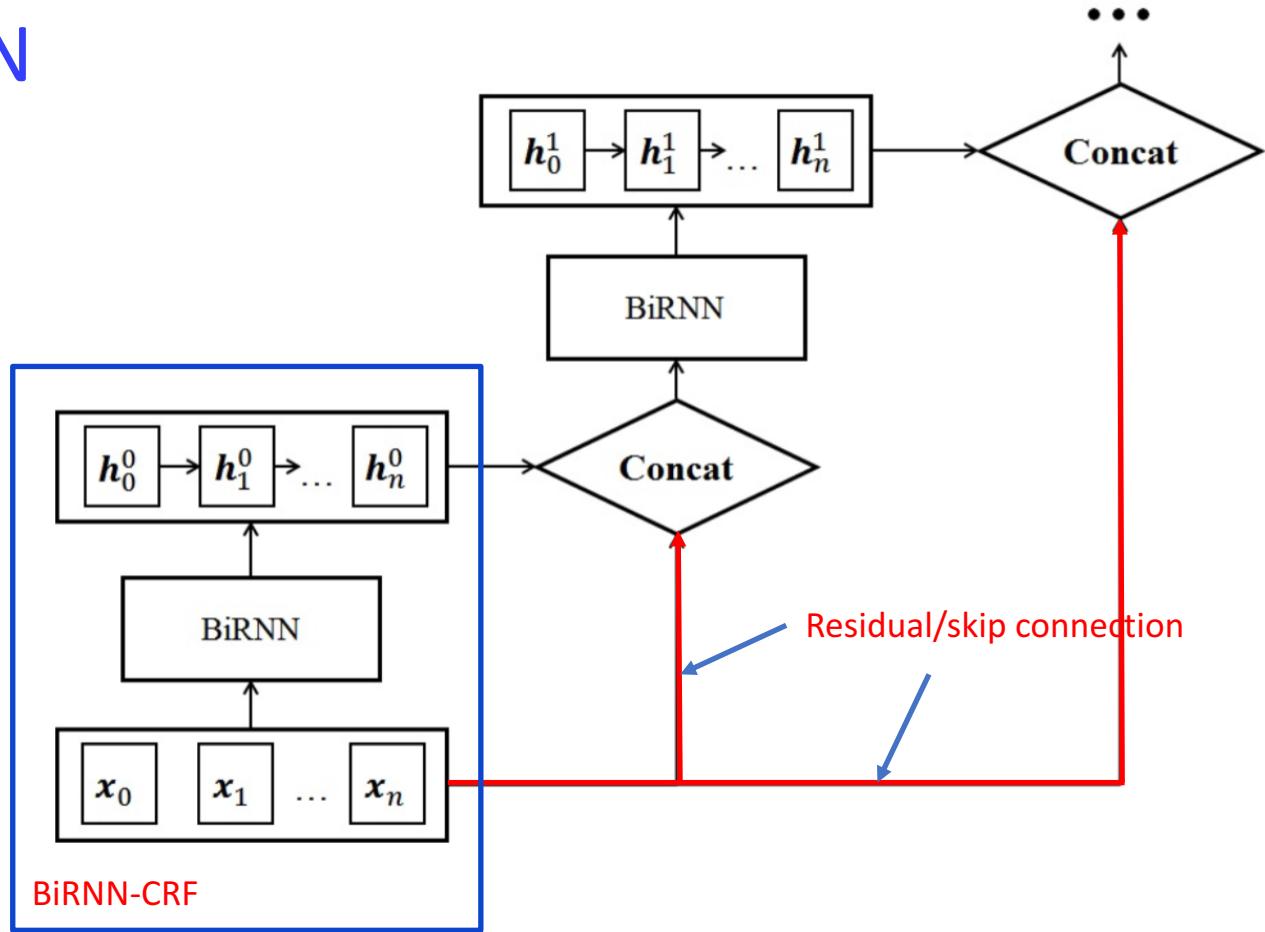
- FOFE+FNN



He dropped a puck from space for the [Toronto Maple Leafs] [home opener against the Buffalo Sabres.]

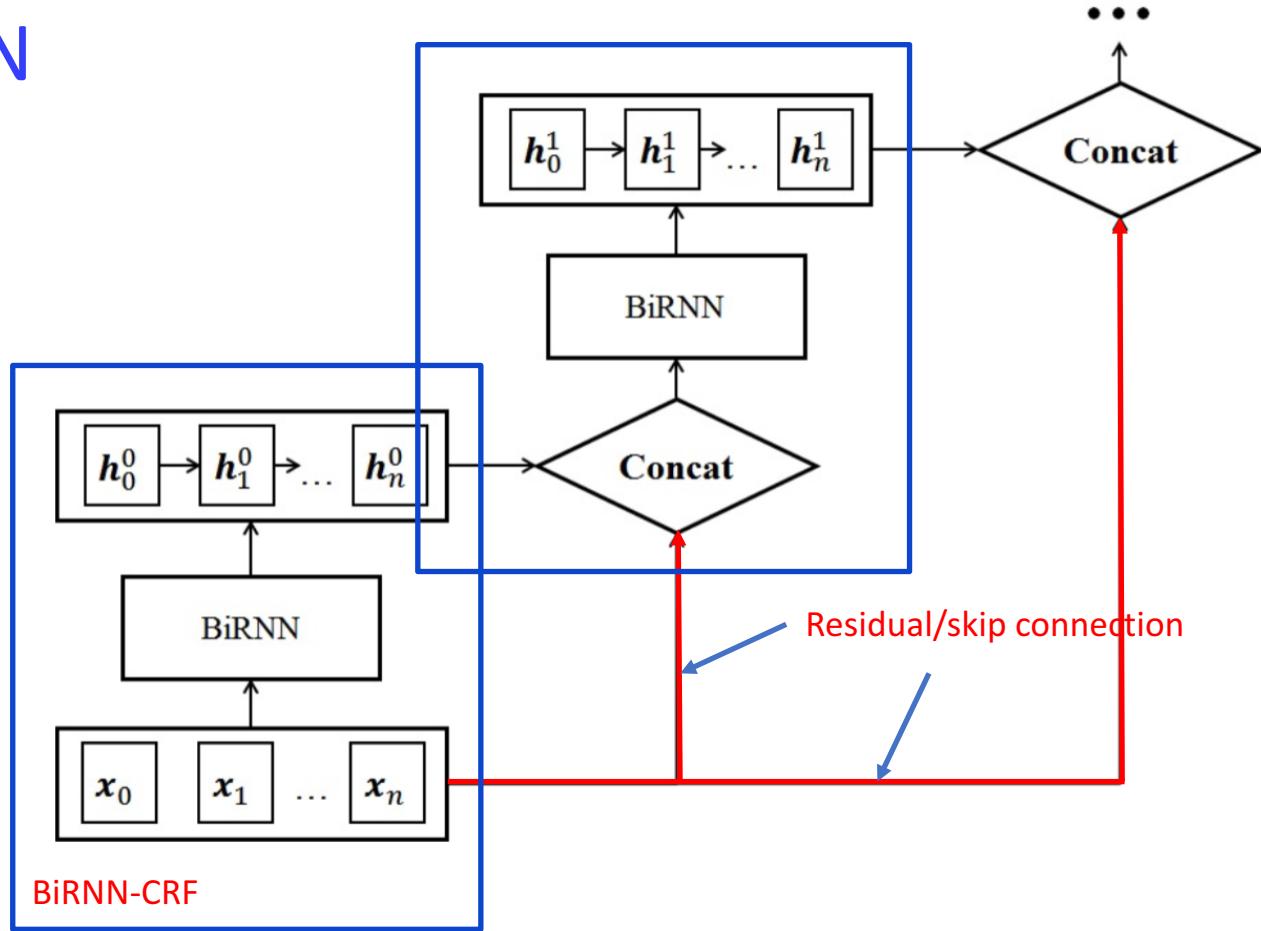
State-of-the-Art CGER Systems

■ Res-RNN



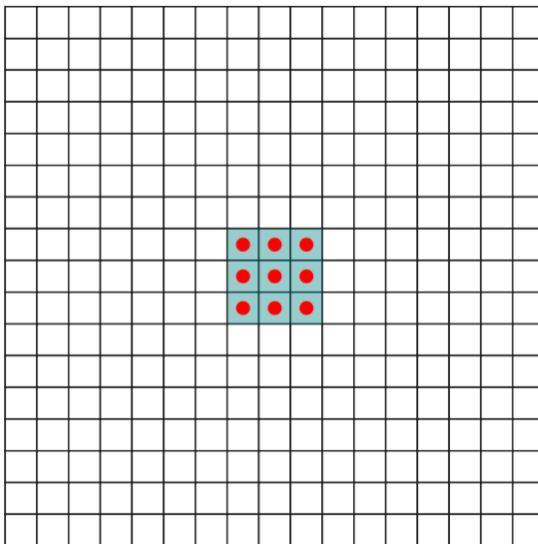
State-of-the-Art CGER Systems

■ Res-RNN

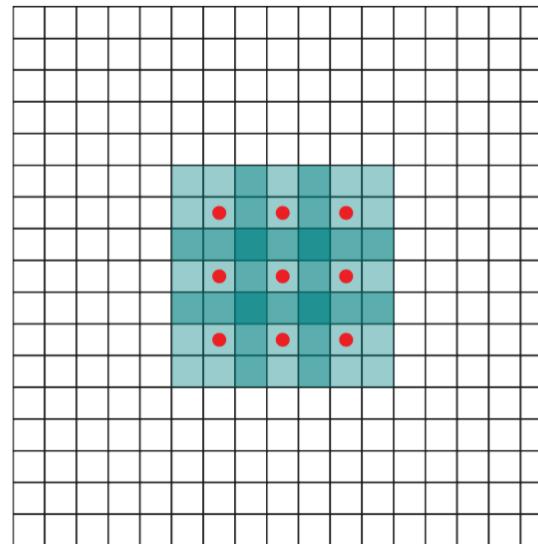


State-of-the-Art CGER Systems

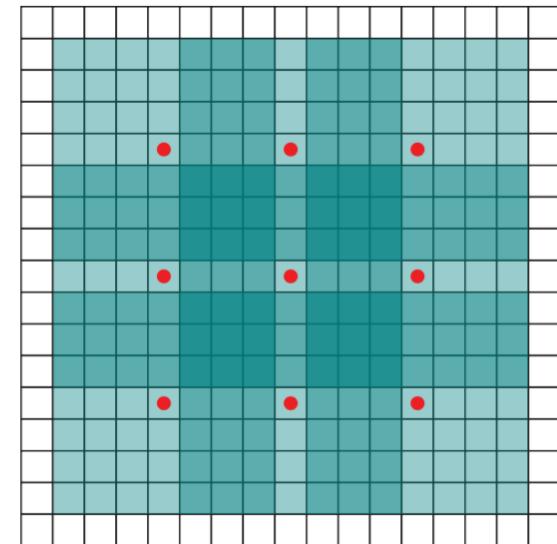
- ID-CNN (Iterated Dilated CNNs)



1-dilated convolution



2-dilated convolution



4-dilated convolution

State-of-the-Art CGER Systems

- ID-CNN (Iterated Dilated CNNs)

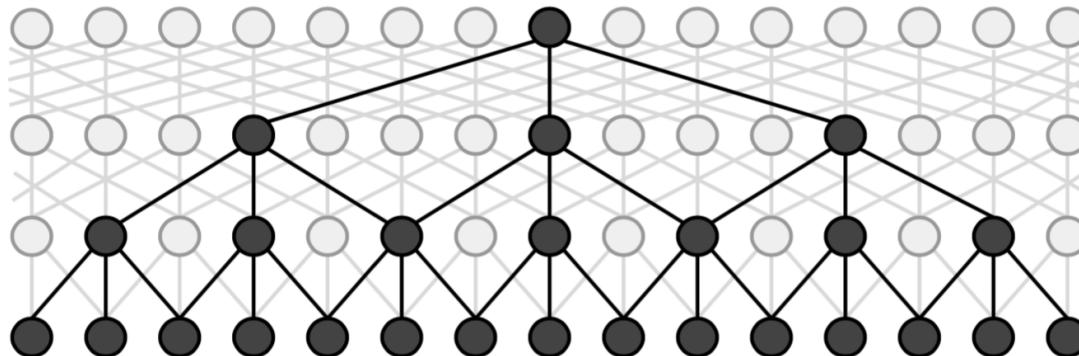


Figure 1: A dilated CNN block with maximum dilation width 4 and filter width 3. Neurons contributing to a single highlighted neuron in the last layer are also highlighted.

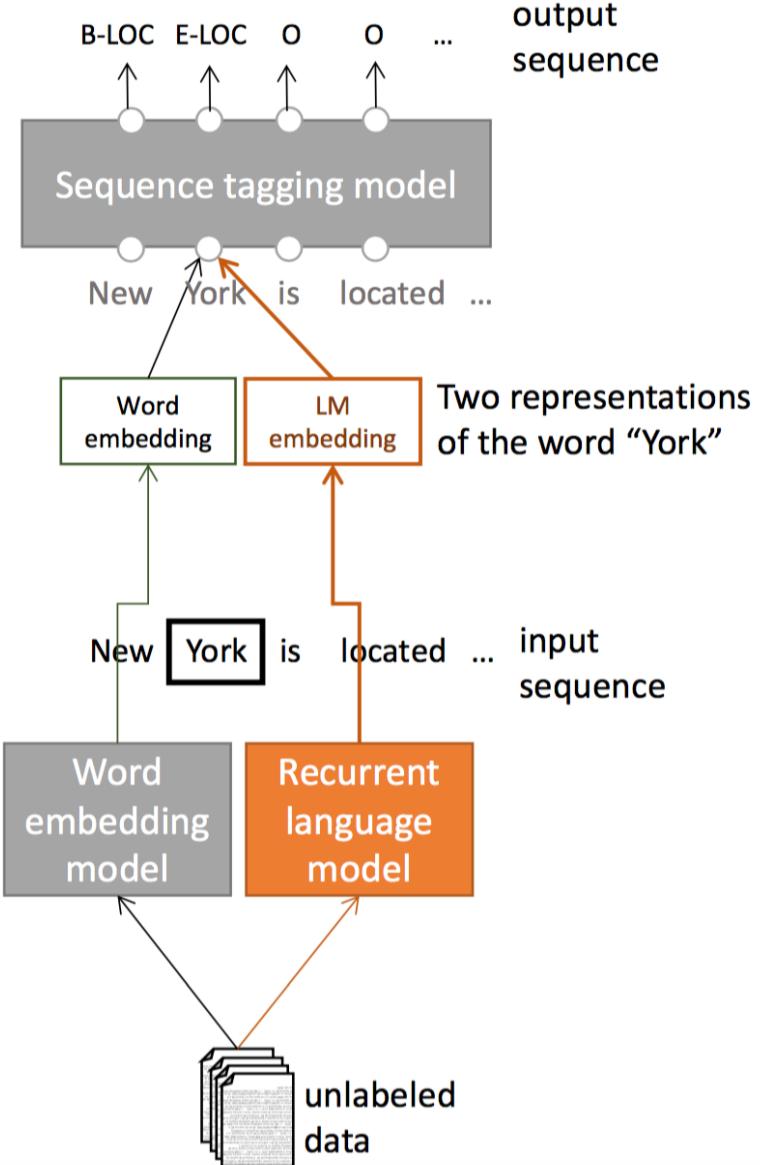
State-of-the-Art CGER Systems

■ BiRNN-CRF + LM

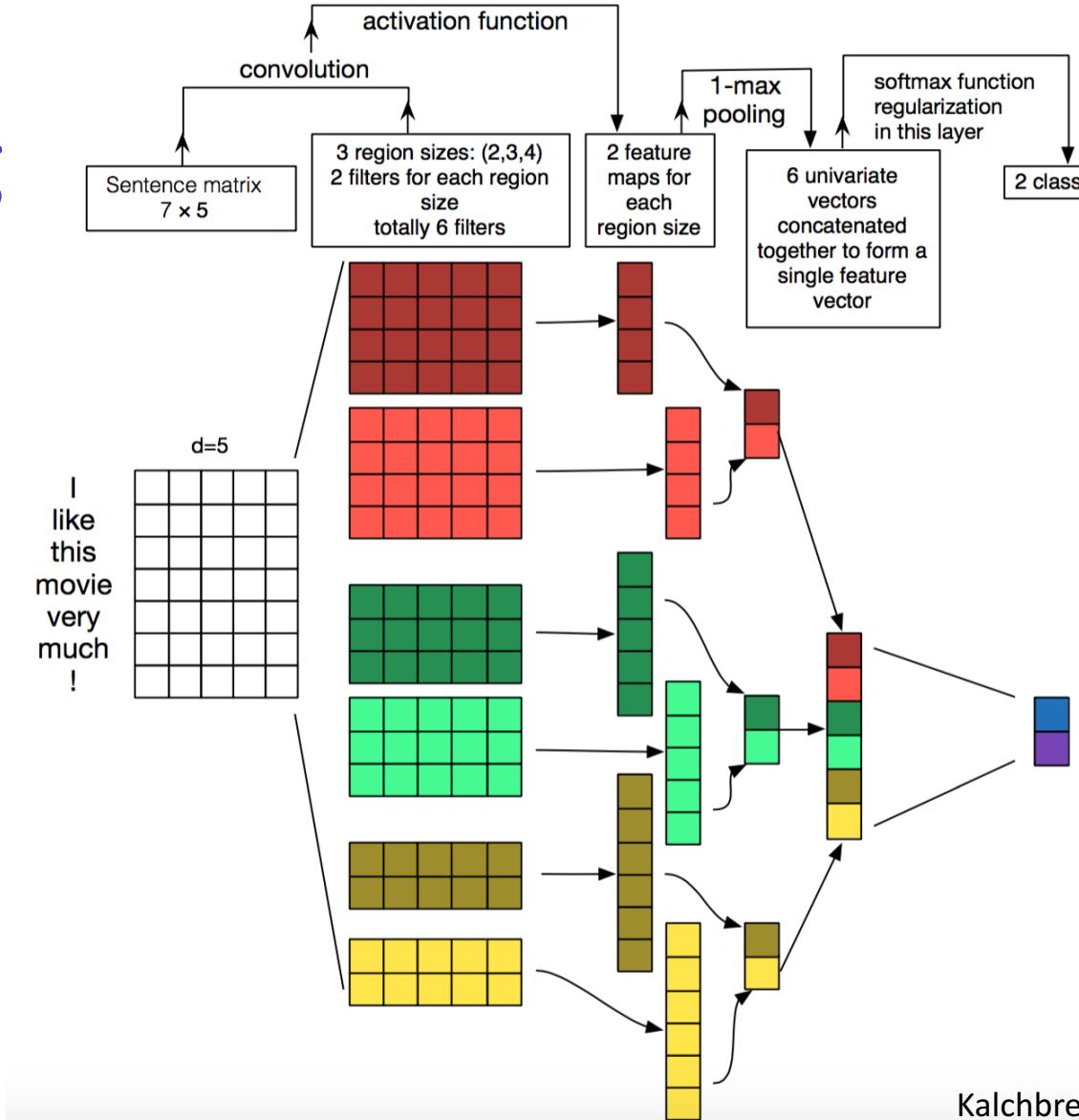
Step 3:
Use both word embeddings and LM embeddings in the sequence tagging model.

Step 2: Prepare word embedding and LM embedding for each token in the input sequence.

Step 1: Pretrain word embeddings and language model.

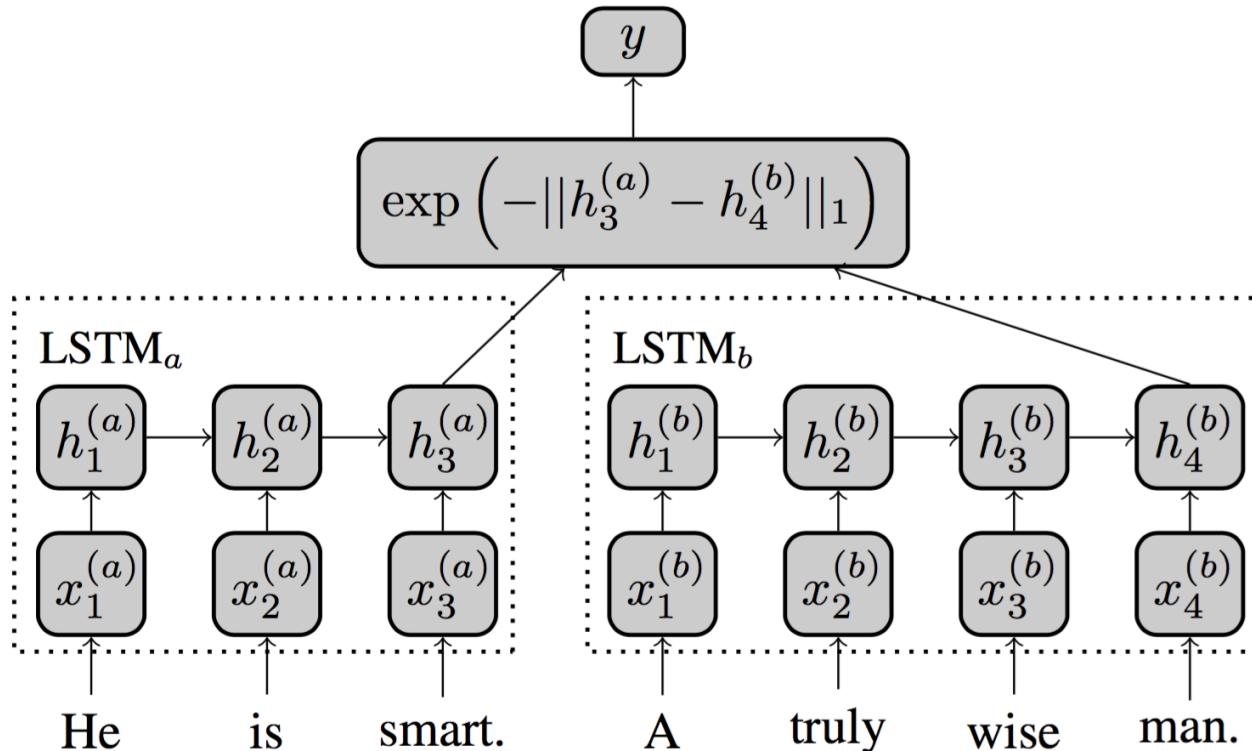


Sentence Modelling



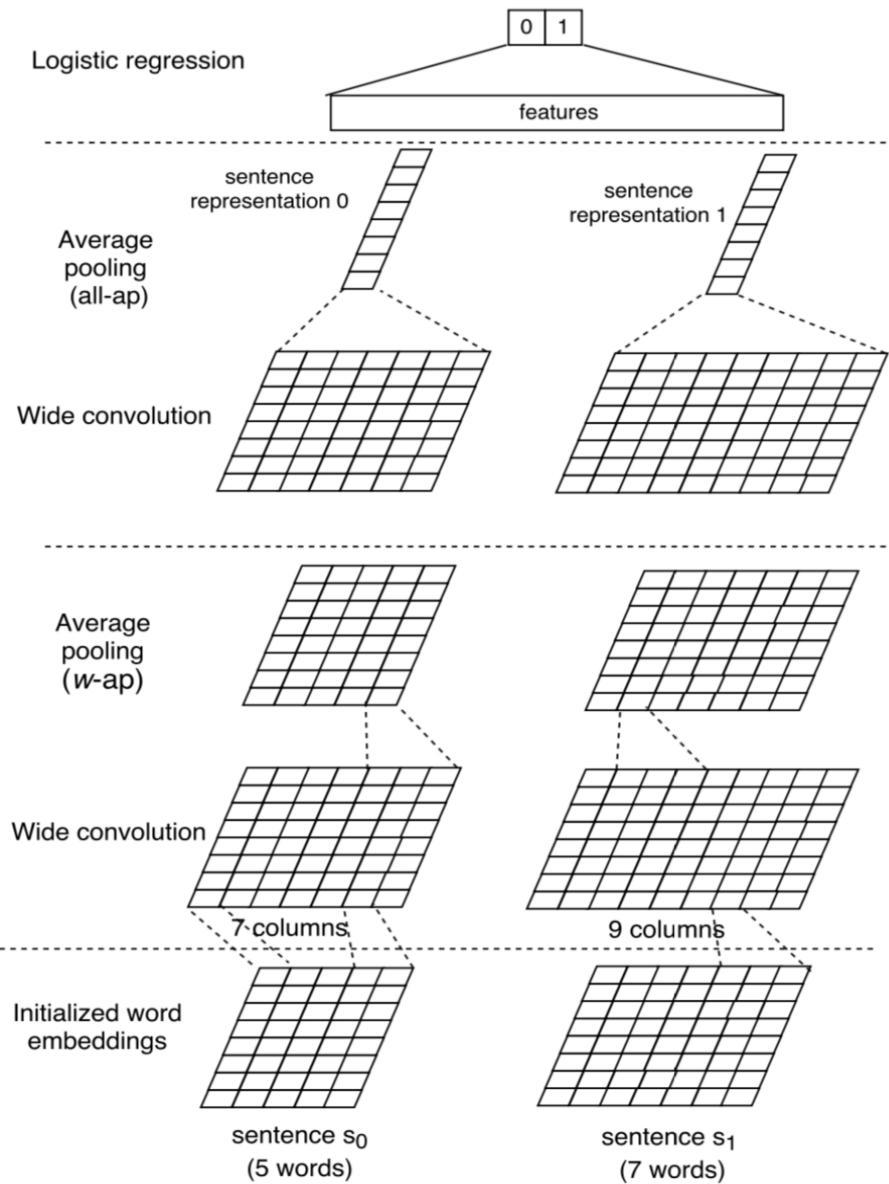
Sentence Modelling

■ Siamese Recurrent Architectures (RNN)

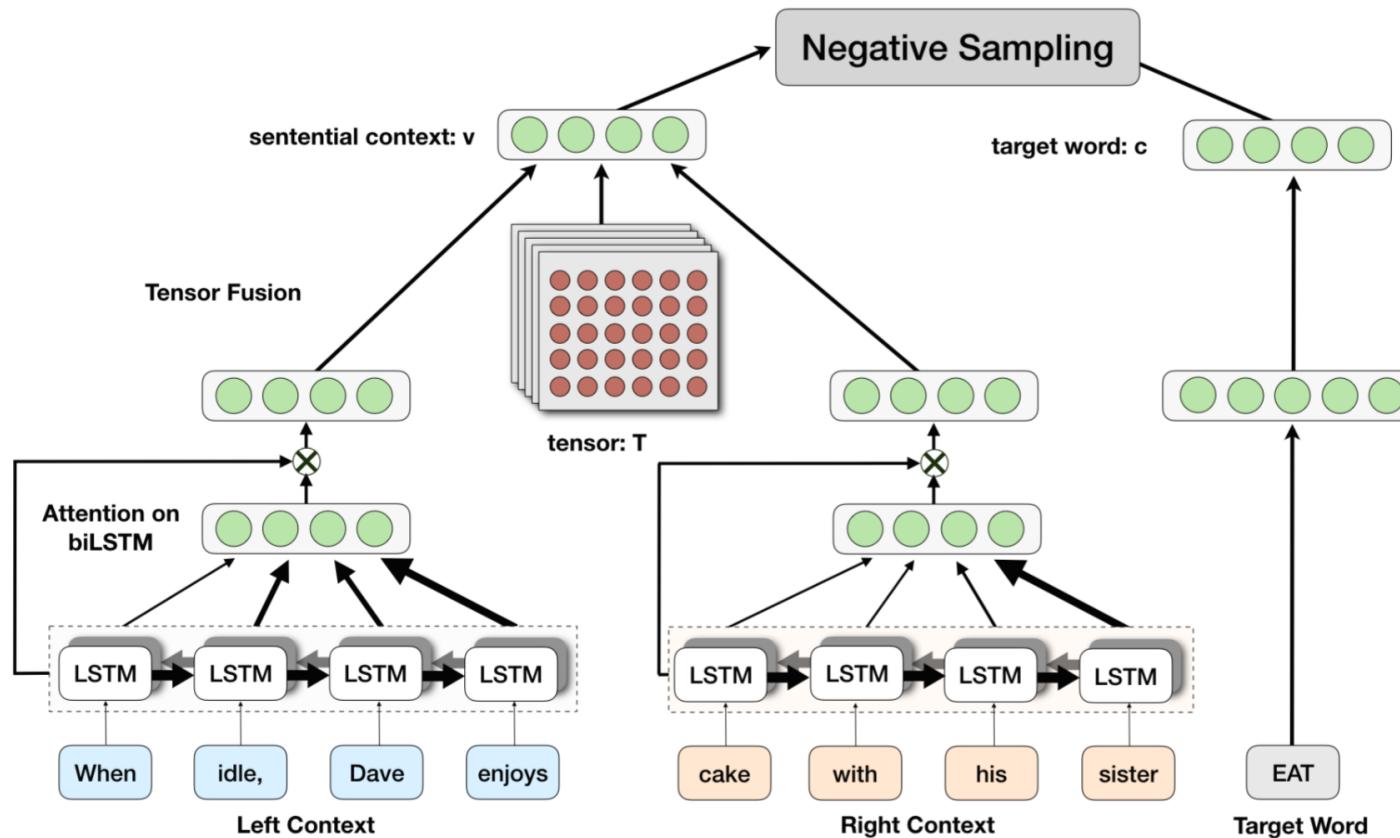


Sentence Modelling

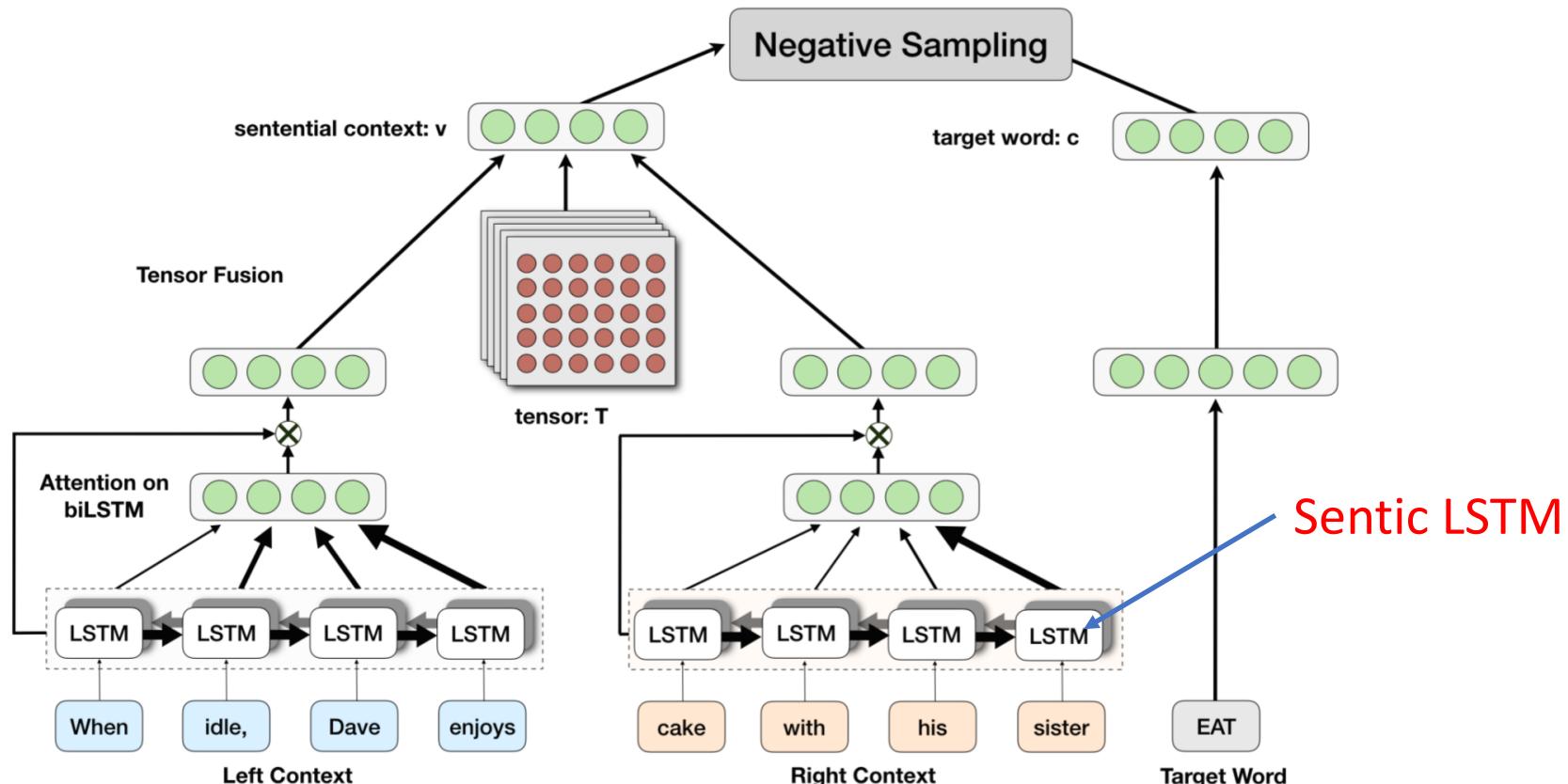
- Siamese Recurrent Architectures (CNN)



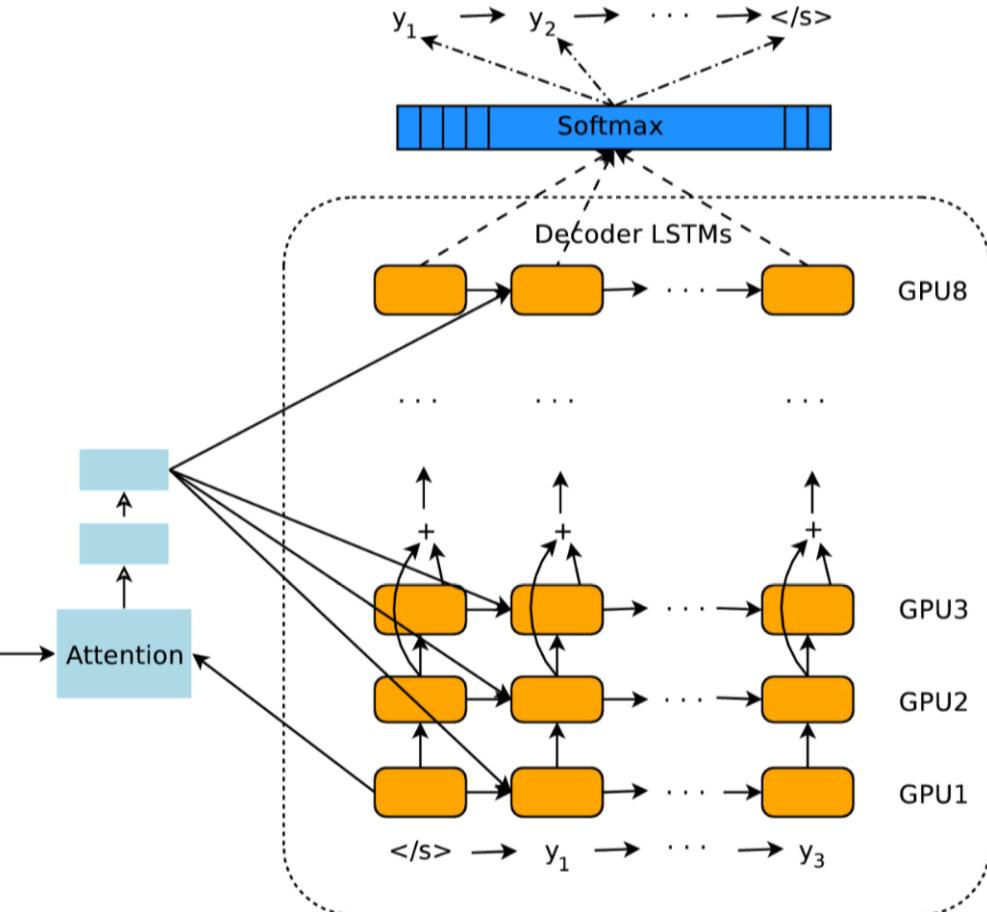
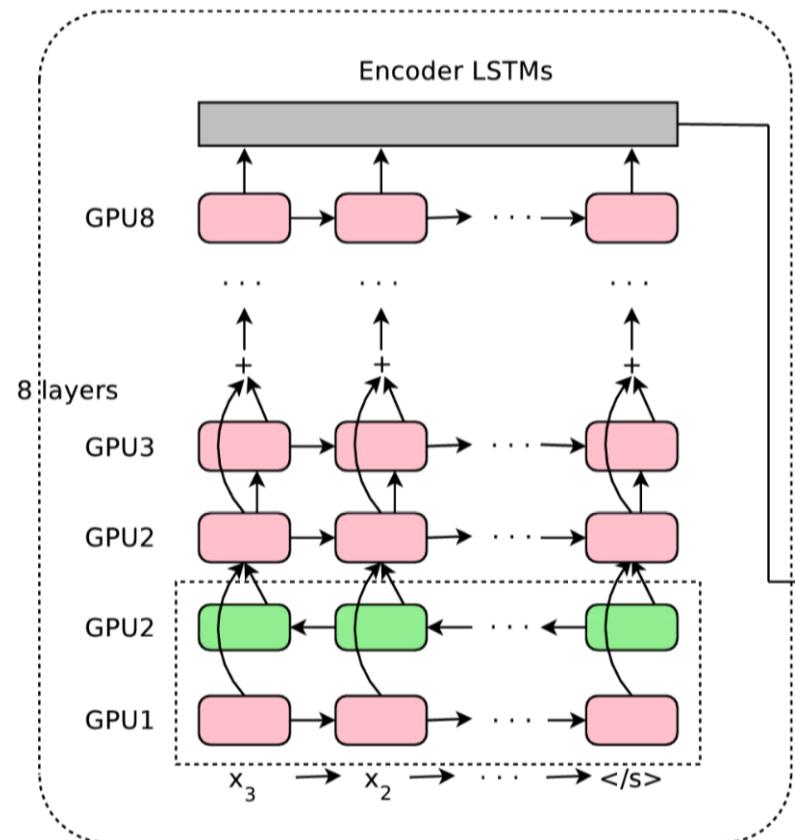
Learning Representation for Sentiment Analysis



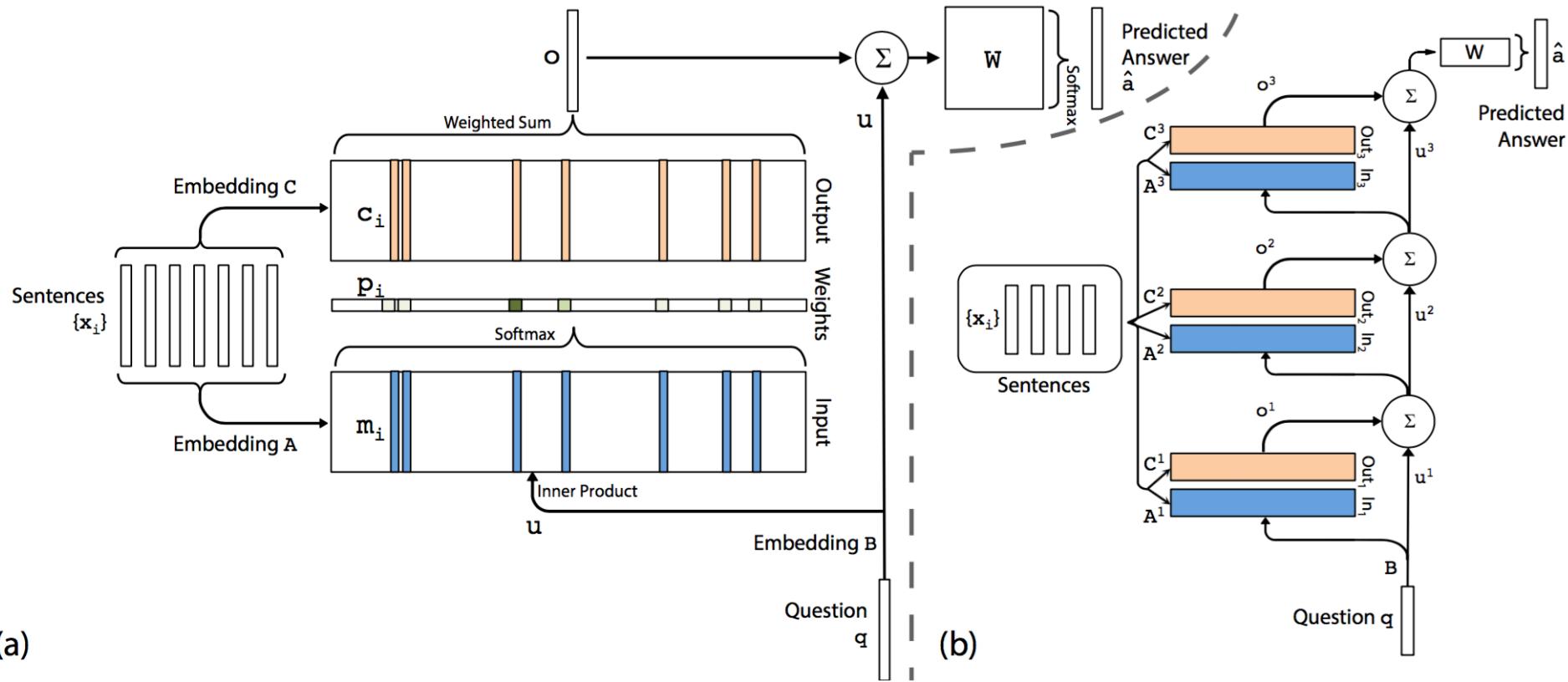
Learning Representation for Sentiment Analysis



Google's Neural Machine Translation



End-To-End Memory Networks



Sequence to Sequence Learning

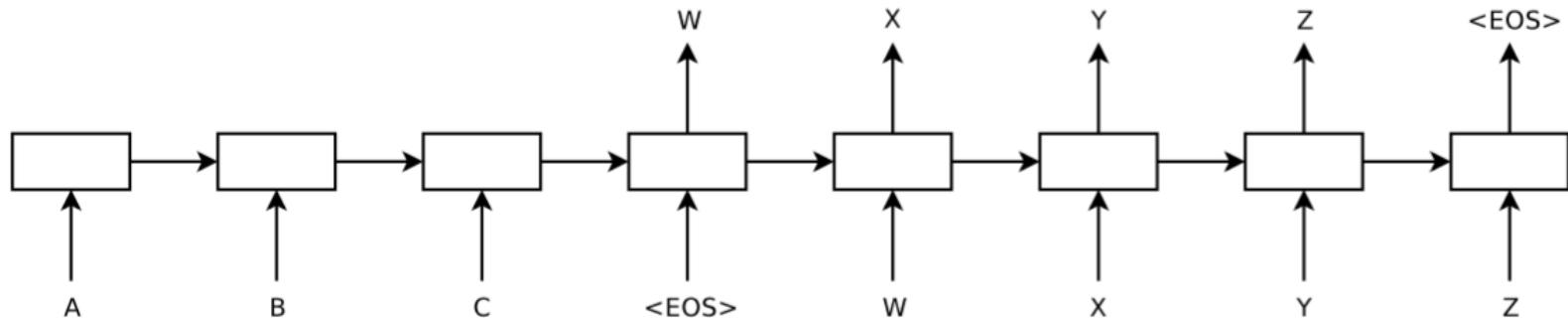
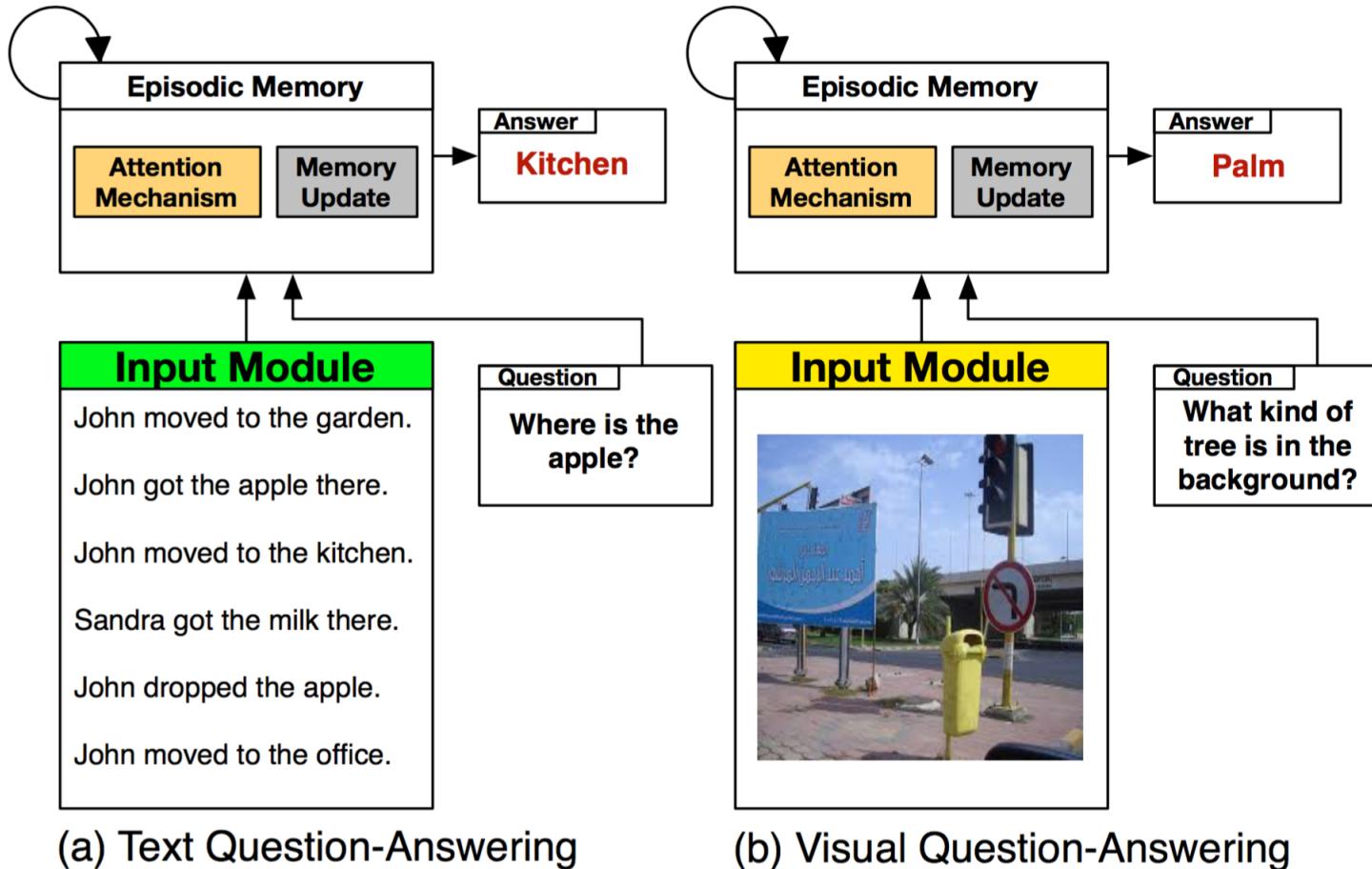


Figure 1: Our model reads an input sentence “ABC” and produces “WXYZ” as the output sentence. The model stops making predictions after outputting the end-of-sentence token. Note that the LSTM reads the input sentence in reverse, because doing so introduces many short term dependencies in the data that make the optimization problem much easier.

Dynamic Memory Networks



(a) Text Question-Answering

(b) Visual Question-Answering

Dynamic Memory Networks

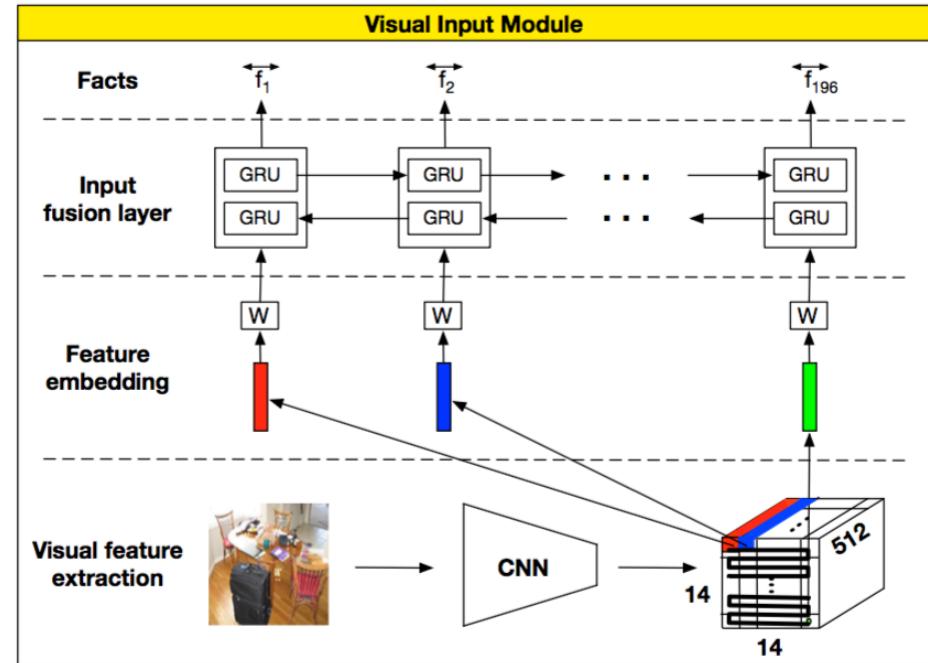
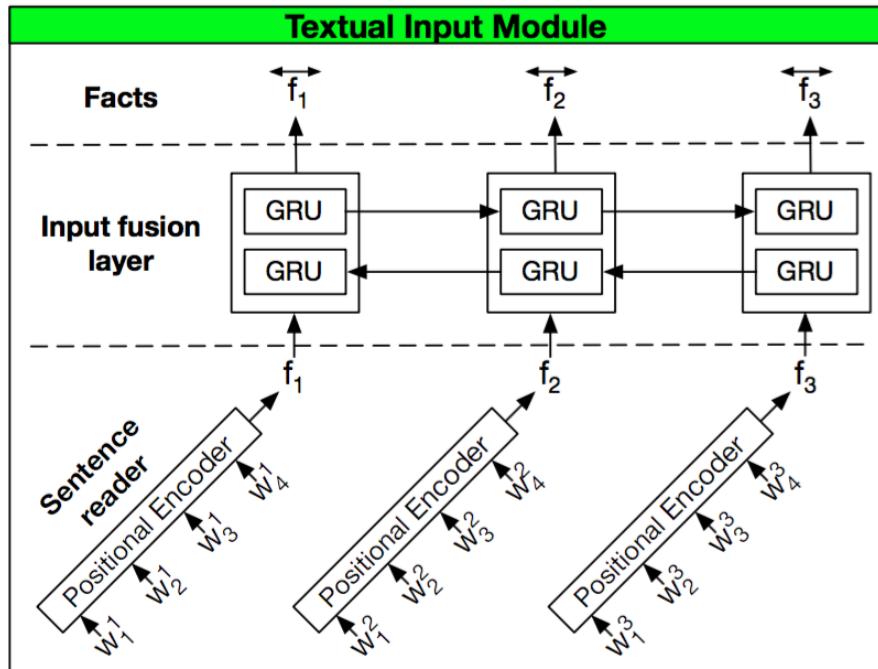


Image Captioning

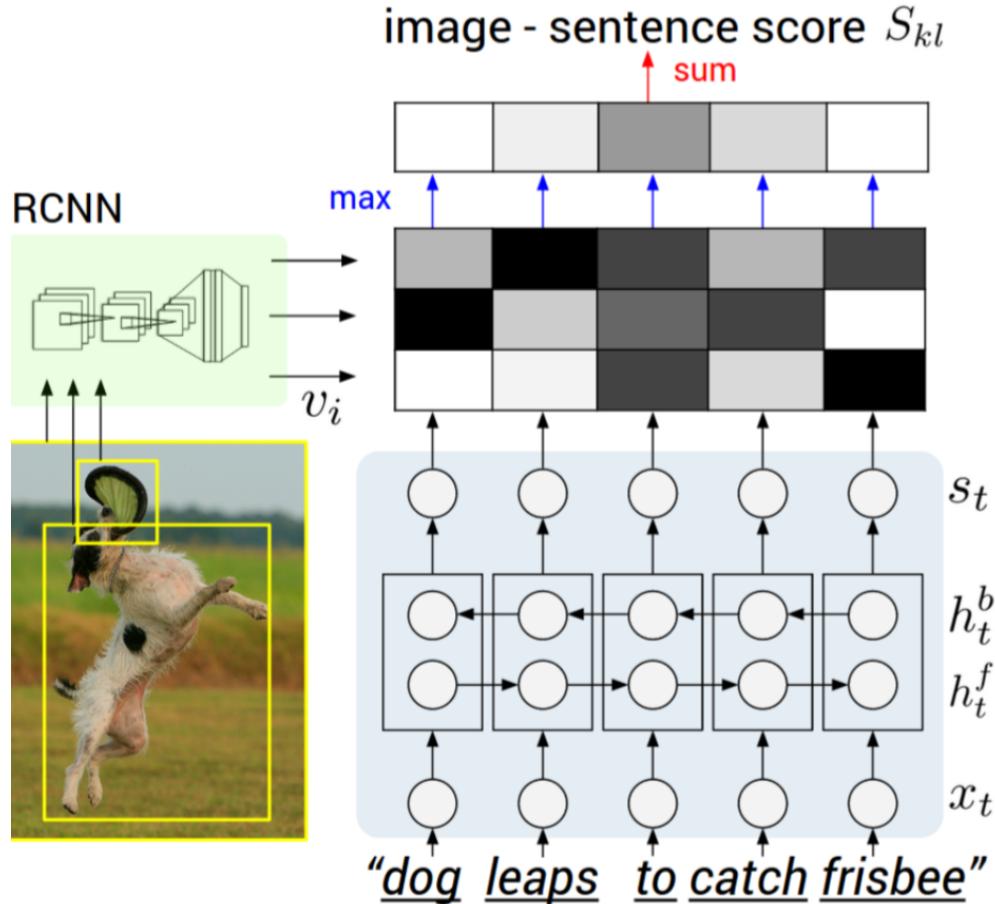
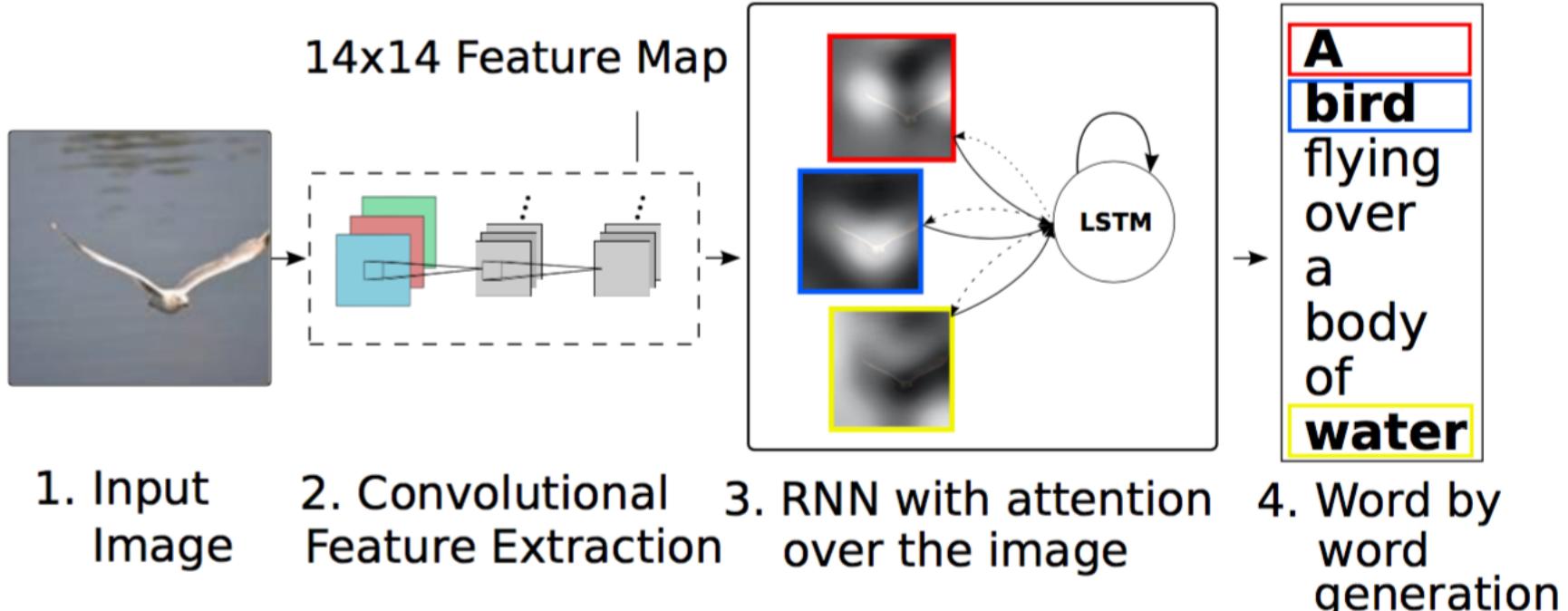
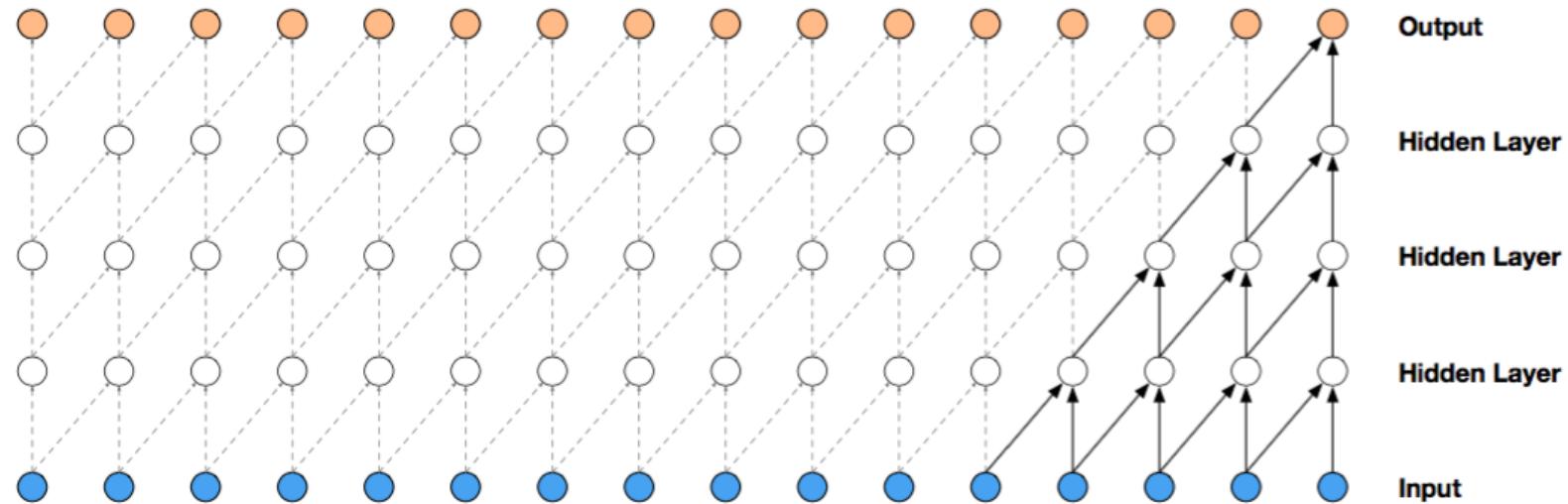


Image Captioning



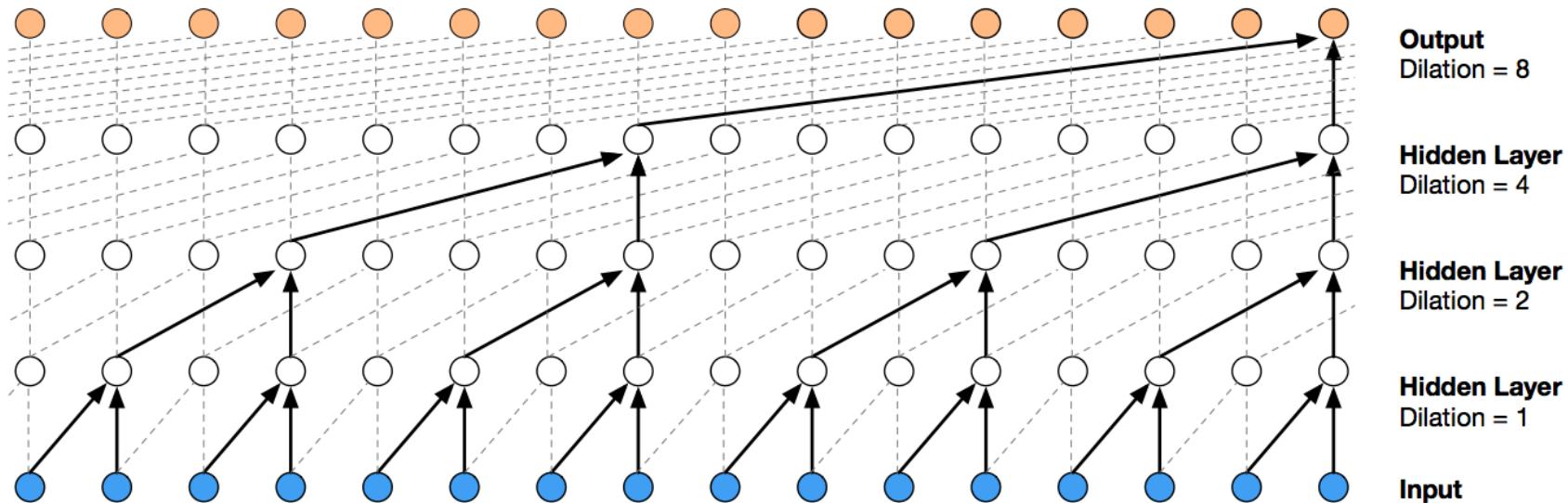
WaveNet

- Stack of causal convolutional layers



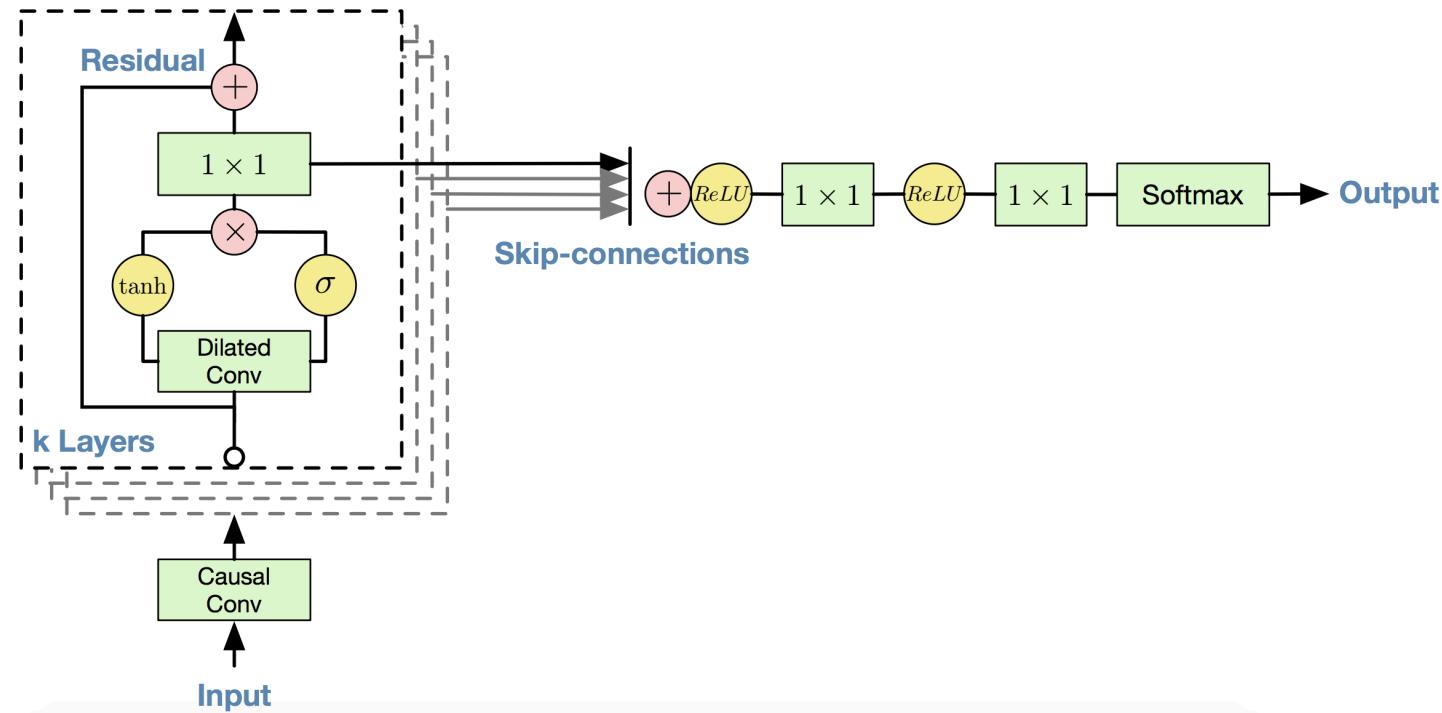
WaveNet

- Stack of **dilated** convolutional layers



WaveNet

- Residual block and overall architecture



Thank you.