

# Lecture 9: Natural Language Processing

Data Science, Fall 2018

Hong-Han Shuai

Thanks to Prof. Hung-yi Lee from NTU for the slides.

# Language Technology

## spam detection



(<http://spam-filter-review.toptenreviews.com/>)

## Part-of-speech Tagging

John saw the saw.  
↓ ↓ ↓ ↓  
PN V D N

## Sentiment Analysis

這部電影太糟了

Negative (負雷)

## Retrieval



## Name Entity Recognition

這位是明 金城武  
Name of People

## Translation

“Machine learning .....



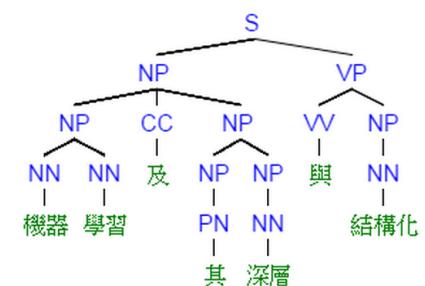
“機器學習 .....

## Speech Recognition

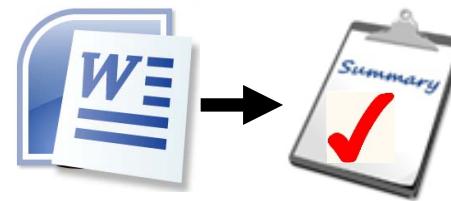


大家好.....

## Syntactic Analysis



## Summarization



document summary

# Do machine really understand human language?

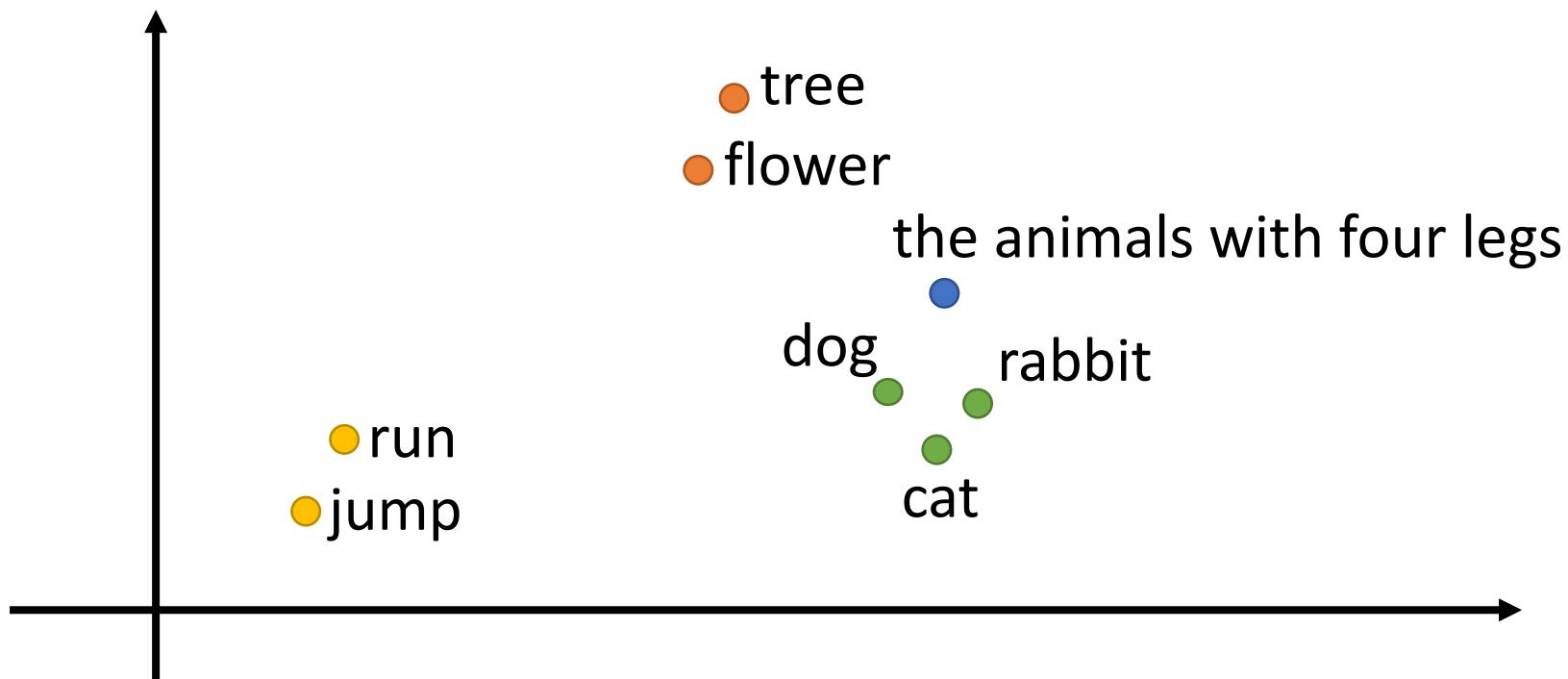


*“A word is known by the  
company it keeps”*

John Rupert Firth

# Meaning Representation

Do machine know the meaning of a word or word sequence?



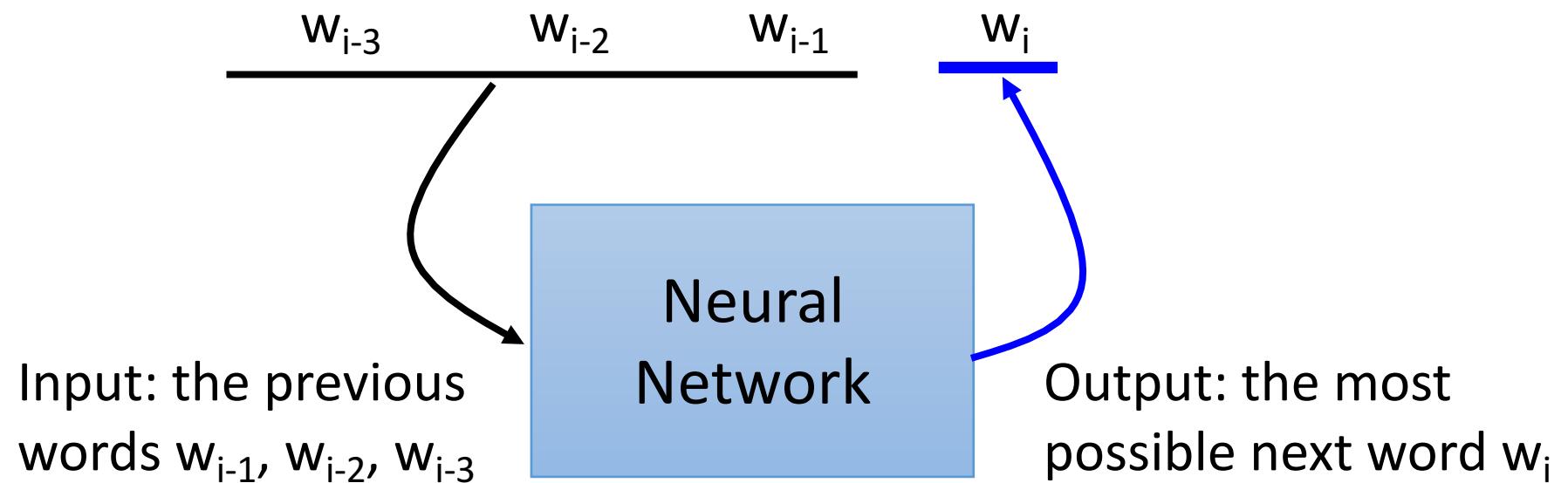
# Meaning of Word

# Predicting the next word

噓	hsiaoyoshye:	麻煩這系列的請到政黑或其他地方討論好嗎?這裡是八	04/27 00:40
推	lrfnc:	仙	04/27 00:51
推	headiron:	樂	04/27 00:52
推	victorshu:	園	04/27 00:57

# Fill in the Blank

..... 哈密瓜 有 一種 \_\_\_\_\_ 味



Each word should be represented as a feature vector.

# Fill in the Blank

## 1-of-N Encoding

lexicon = {apple, bag, cat, dog, elephant}

apple = [ 1 0 0 0 0] The vector is lexicon size.

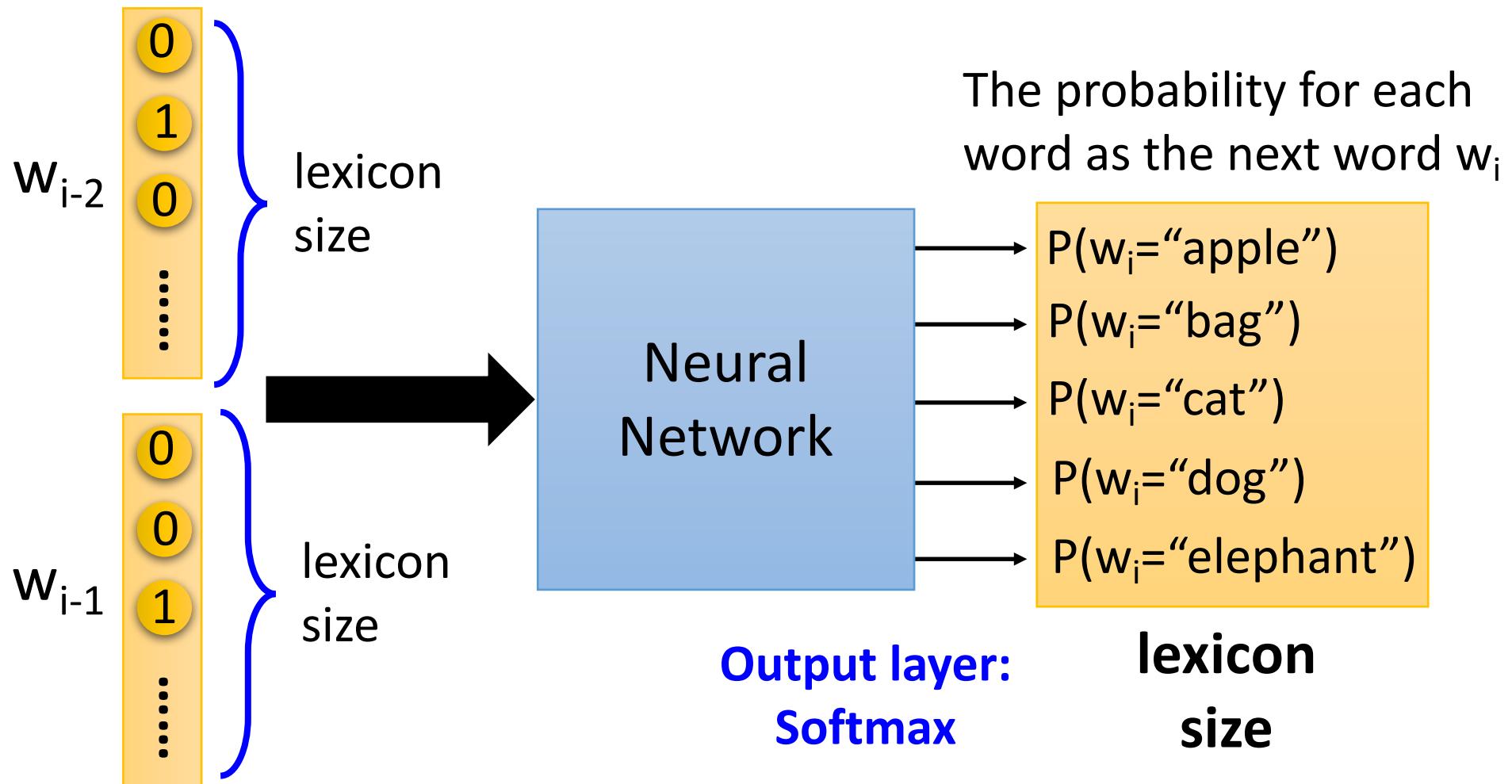
bag = [ 0 1 0 0 0] Each dimension corresponds

cat = [ 0 0 1 0 0] to a word in the lexicon

dog = [ 0 0 0 1 0] The dimension for the word

elephant = [ 0 0 0 0 1] is 1, and others are 0

# Fill in the Blank



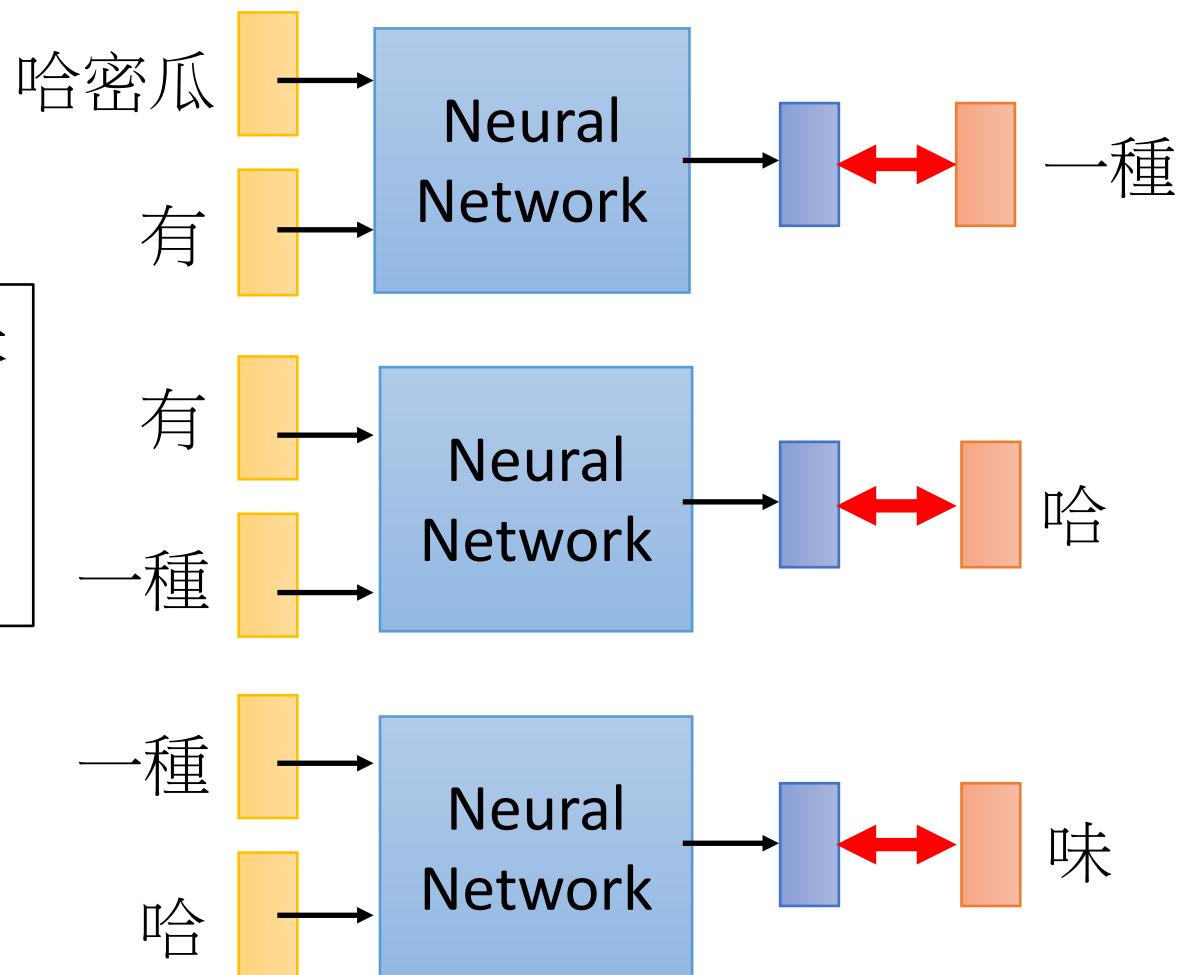
# Fill in the Blank

- Training:

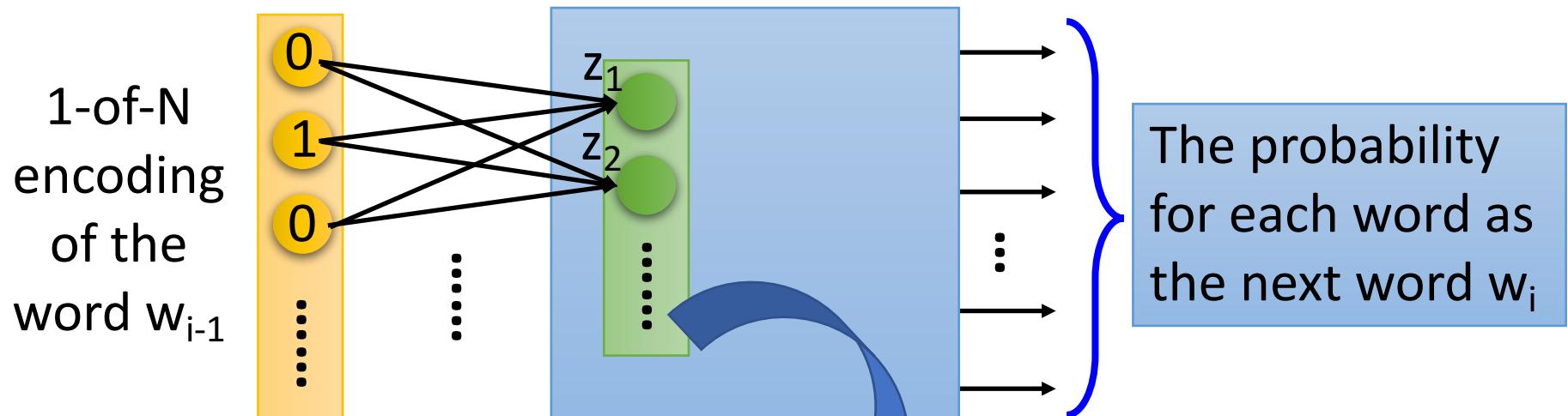
Collect data:

哈密瓜 有一種 哈 味  
不爽 不要 買  
公道價 八萬 一  
.....

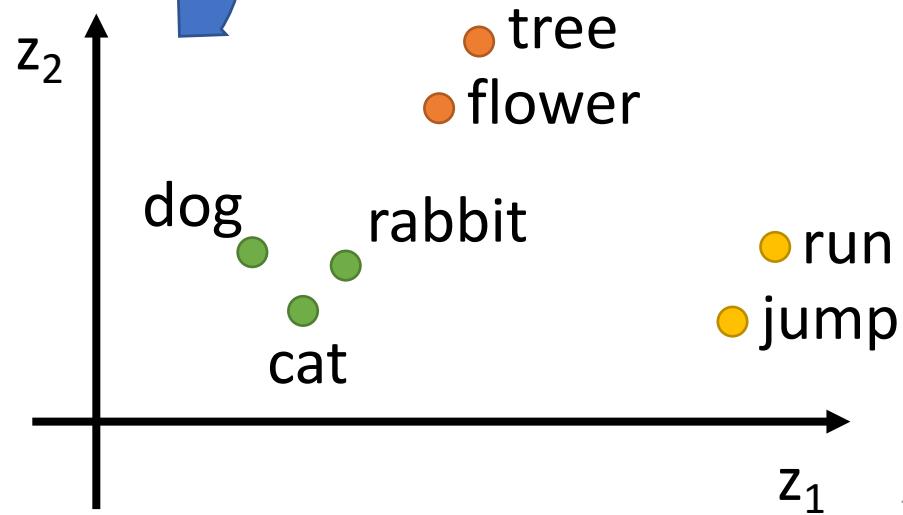
**Minimizing  
cross entropy**



# Word Vector

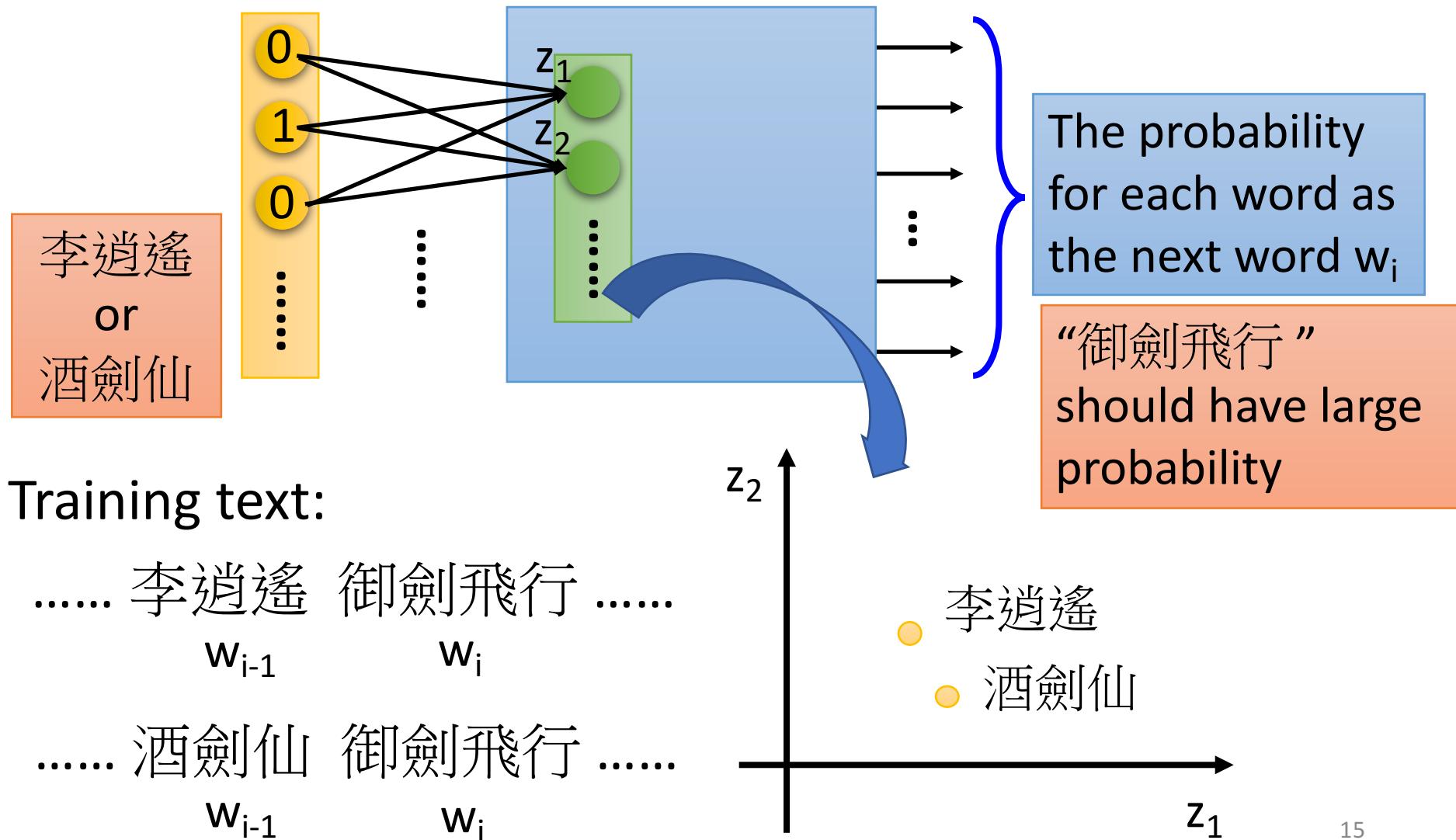


- Take out the input of the neurons in the first layer
- Use it to represent a word  $w$
- Word vector, word embedding feature:  $V(w)$

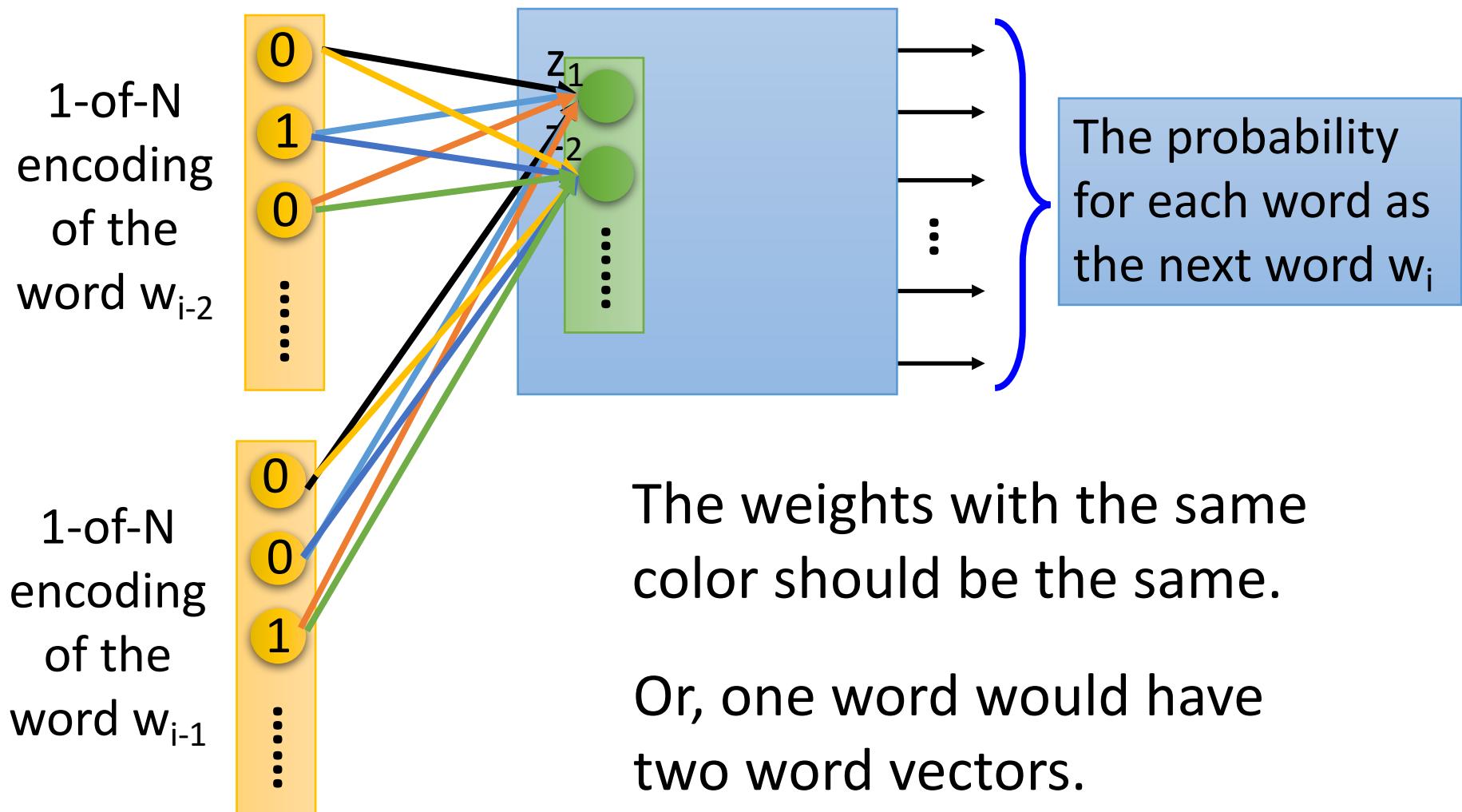


# Word Vector

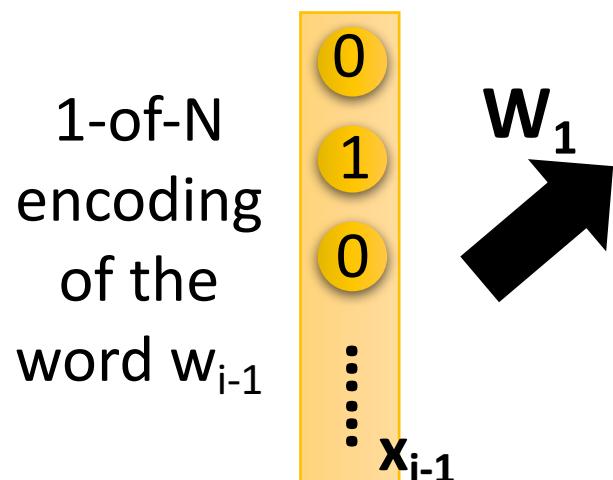
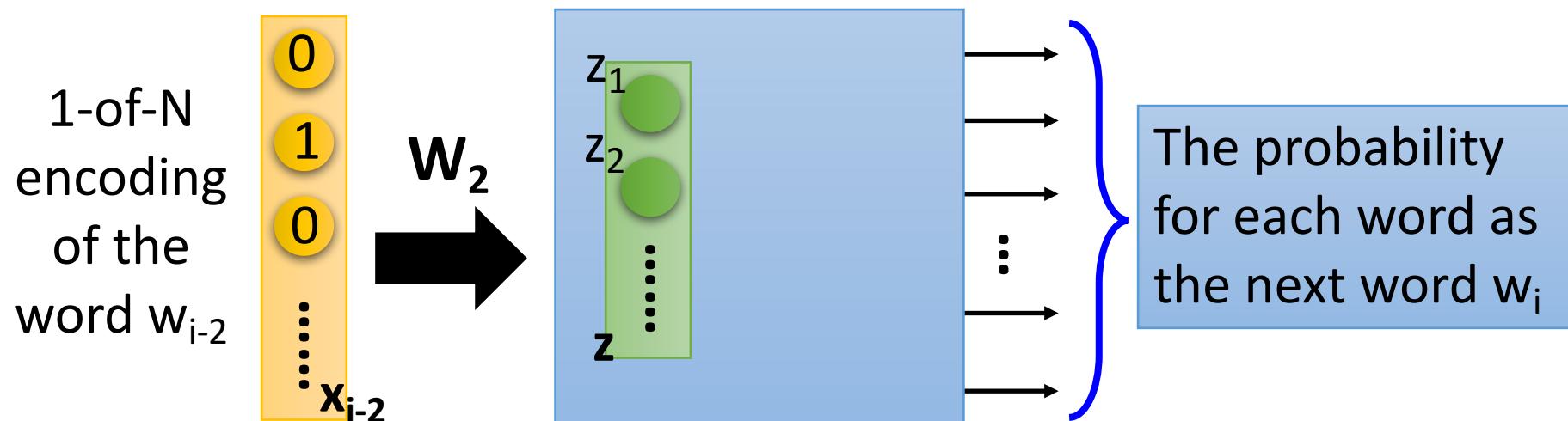
You shall know a word  
by the company it keeps



# Word Vector – Sharing Parameters

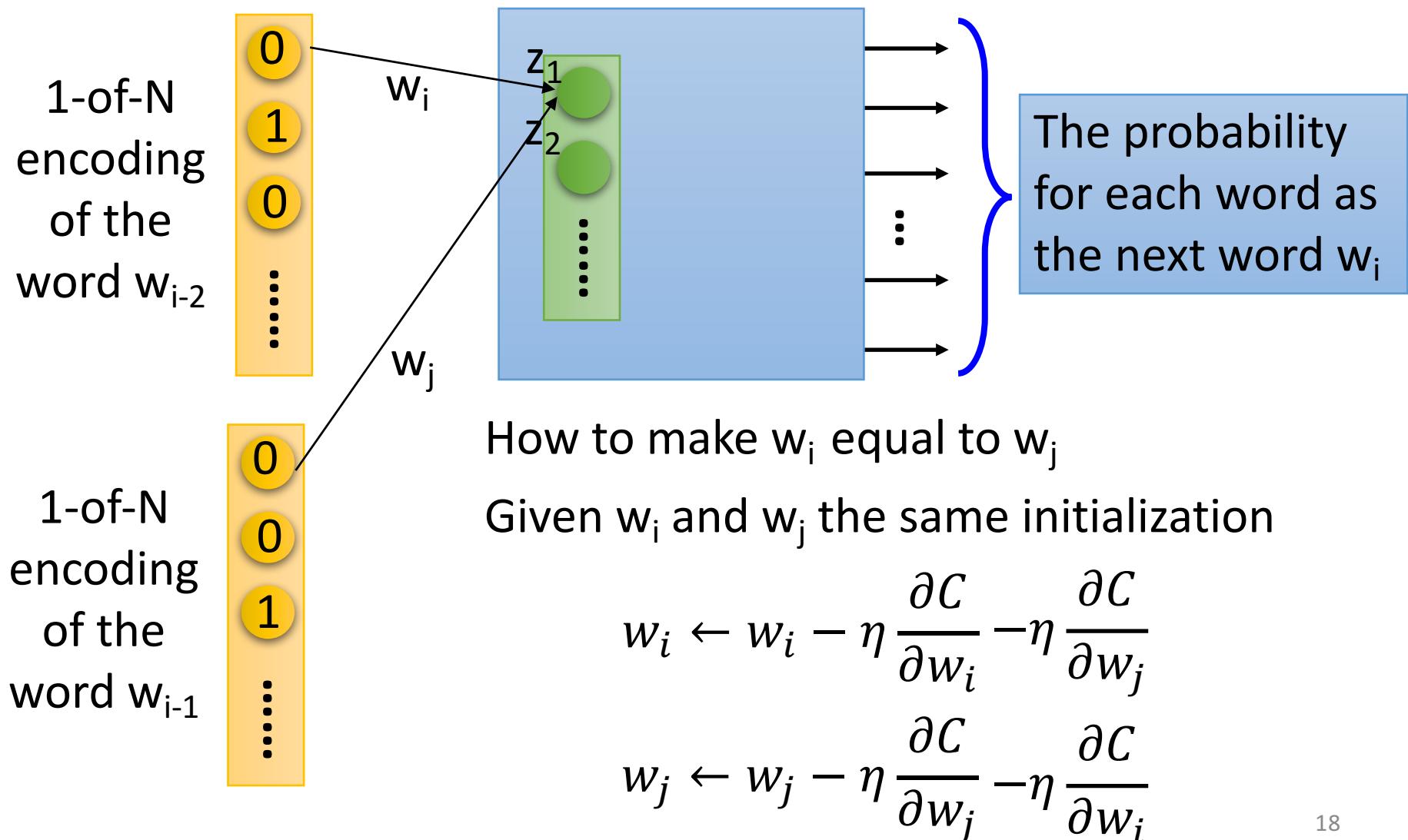


# Word Vector – Sharing Parameters



The length of  $x_{i-1}$  and  $x_{i-2}$  are both  $|V|$ .  
The length of  $z$  is  $|Z|$ .  
$$z = W_1 x_{i-1} + W_2 x_{i-2}$$
  
The weight matrix  $W_1$  and  $W_2$  are both  $|Z| \times |V|$  matrices.  
$$W_1 = W_2 = W \rightarrow z = W (x_{i-2} + x_{i-1})$$

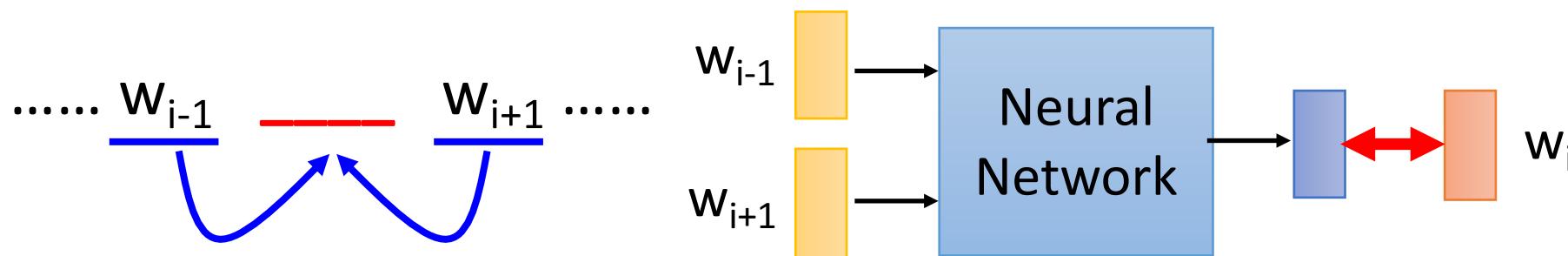
# Word Vector – Sharing Parameters



# Word Vector

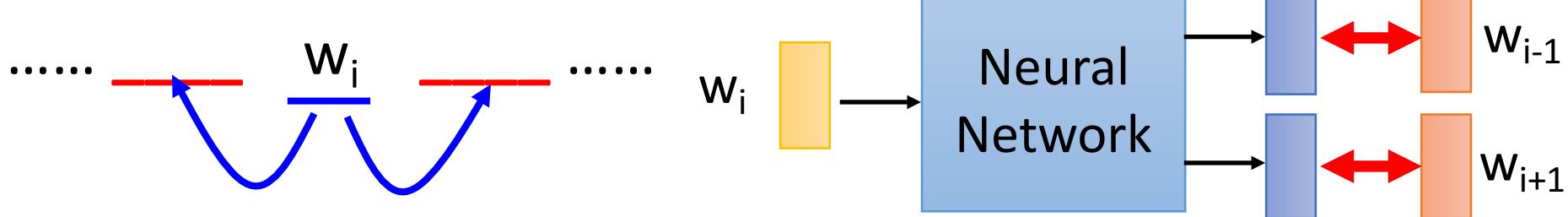
## – Various Architectures

- Continuous bag of word (CBOW) model



*predicting the word given its context*

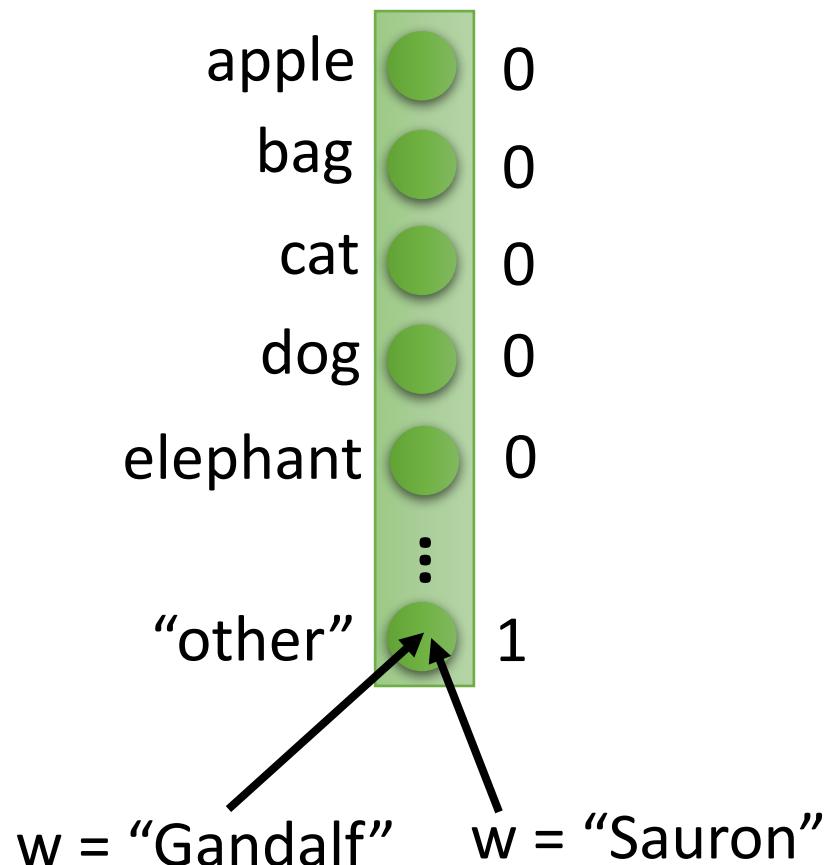
- Skip-gram



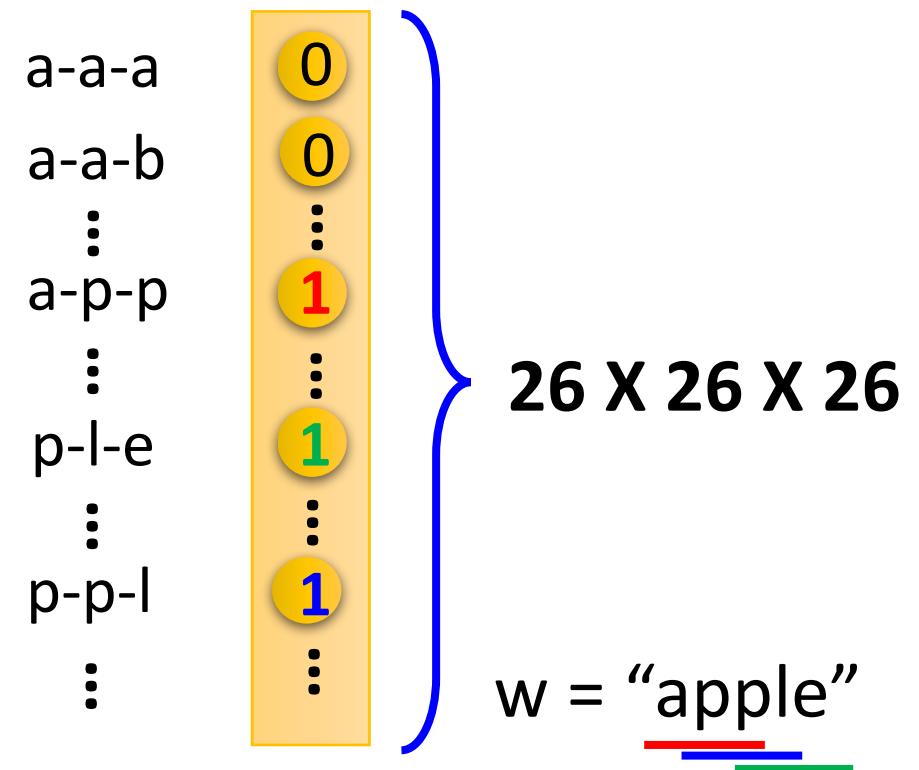
*predicting the context given a word*

# Beyond 1-of-N encoding

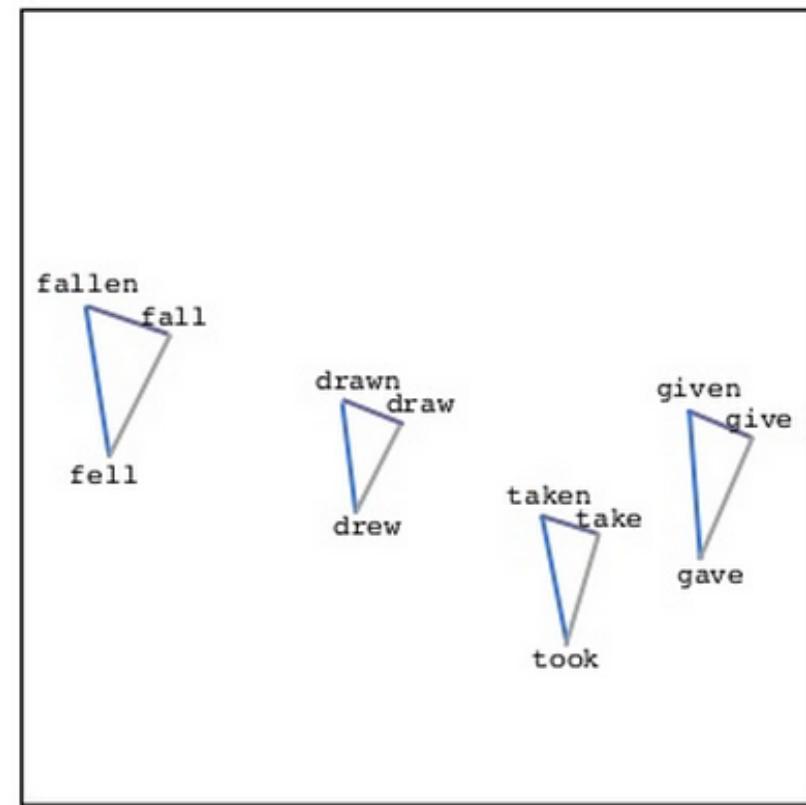
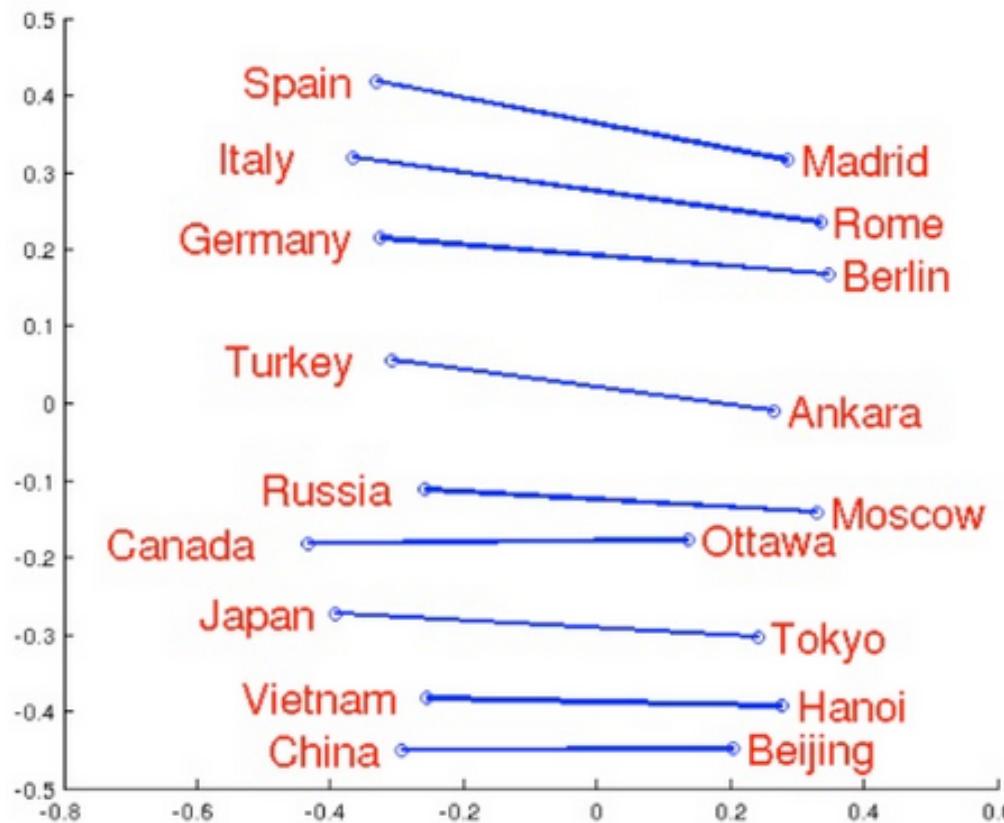
## *Dimension for “Other”*



## *Word hashing*

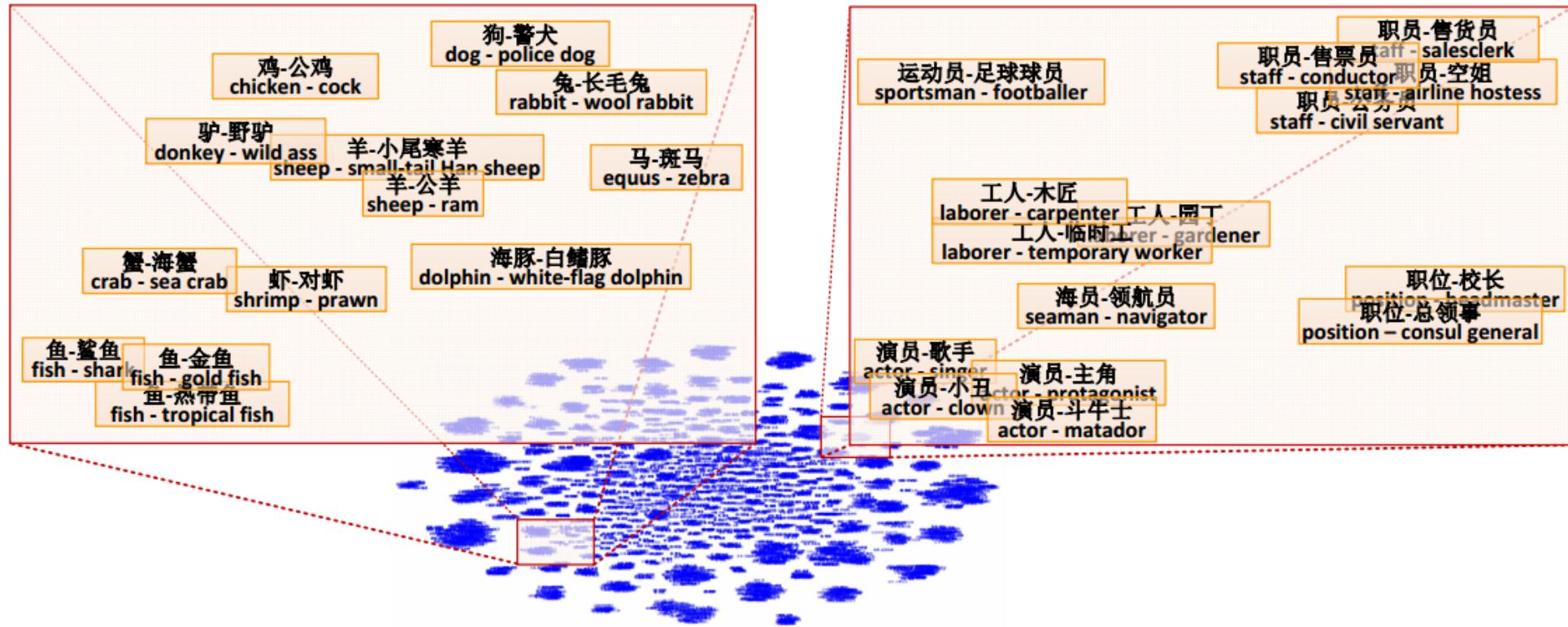


# Word Vector



Source: <http://www.slideshare.net/hustwj/cikm-keynotenov2014>

# Word Vector



Fu, Ruiji, et al. "Learning semantic hierarchies via word embeddings." *Proceedings of the 52th Annual Meeting of the Association for Computational Linguistics: Long Papers*. Vol. 1. 2014.

Word Vector      
$$\approx V(Berlin) - V(Rome) + V(Italy)$$

- Characteristics

$$V(hotter) - V(hot) \approx V(bigger) - V(big)$$

$$V(Rome) - V(Italy) \approx V(Berlin) - V(Germany)$$

$$V(king) - V(queen) \approx V(uncle) - V(aunt)$$

- Solving analogies

Rome : Italy = Berlin : ?

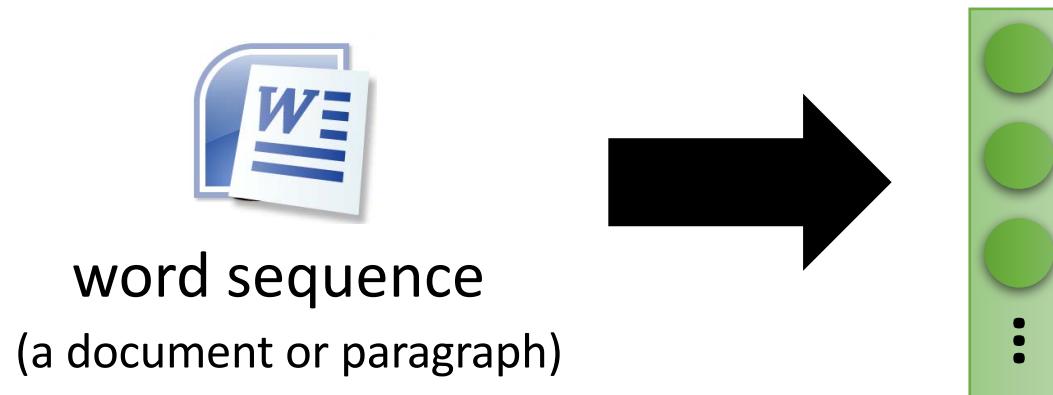
Compute  $V(Berlin) - V(Rome) + V(Italy)$

Find the word w with the closest  $V(w)$

# Meaning of Word Sequence

# Meaning of Word Sequence

- word sequences with different lengths → the vector with the same length
  - The vector representing the meaning of the word sequence
  - A word sequence can be a document or a paragraph



# Outline

Deep Structured  
Semantic Model  
(DSSM)

- Application: Information Retrieval (IR)

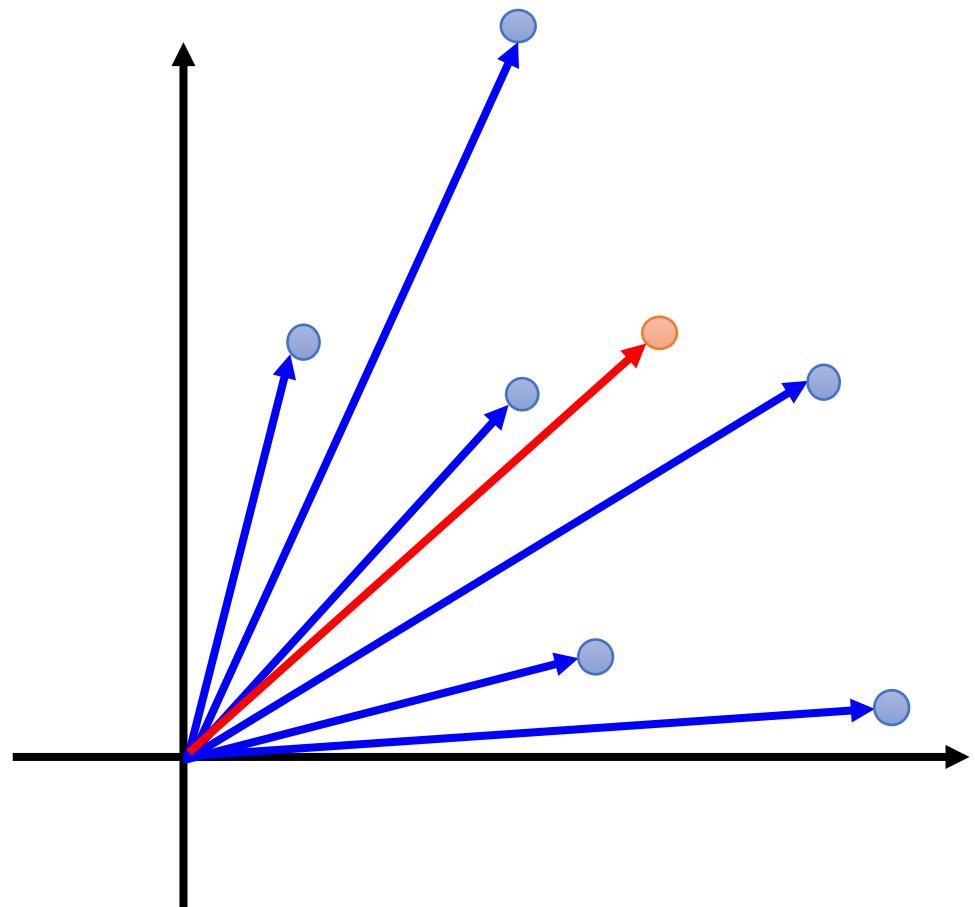
Recursive Neural  
Network

- Application: Sentiment Analysis,  
Sentence Relatedness

Unsupervised

- Paragraph Vector
- Sequence-to-sequence auto-encoder

# Information Retrieval (IR)



## **Vector Space Model**

The documents are vectors in the space.

The query is also a vector.

How to use a vector to represent word sequences

# Information Retrieval (IR)

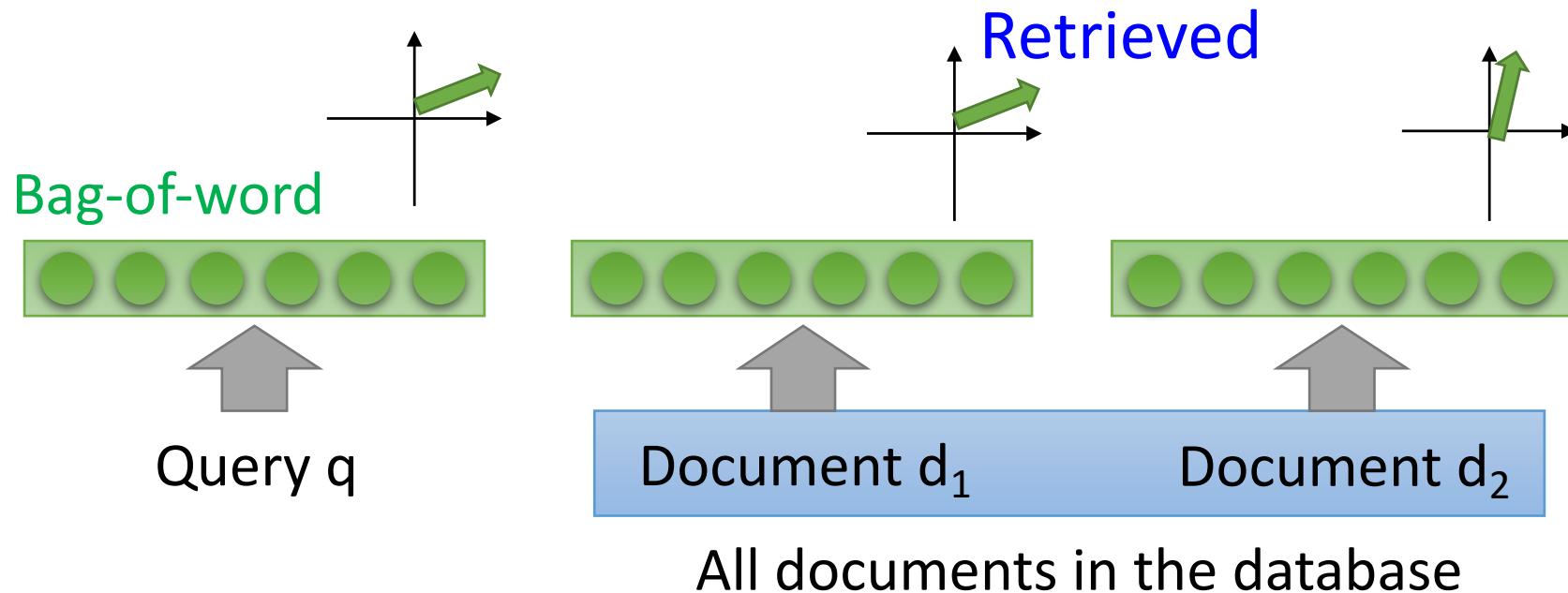
## ***Bag-of-word***

	this		1		this		1
	is		1		is		1
word string s1:	a		0	word string s2:	a		1
“This is an apple”	an		1	“This is a pen”	an		0
	apple		1		apple		0
	pen		0		pen		1
		:				:	

Weighted by IDF  
28

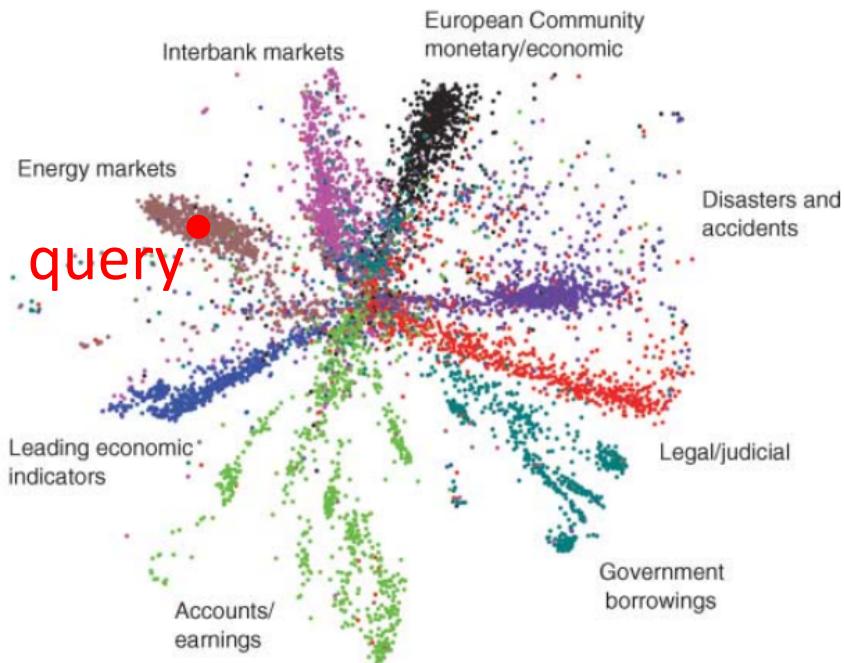
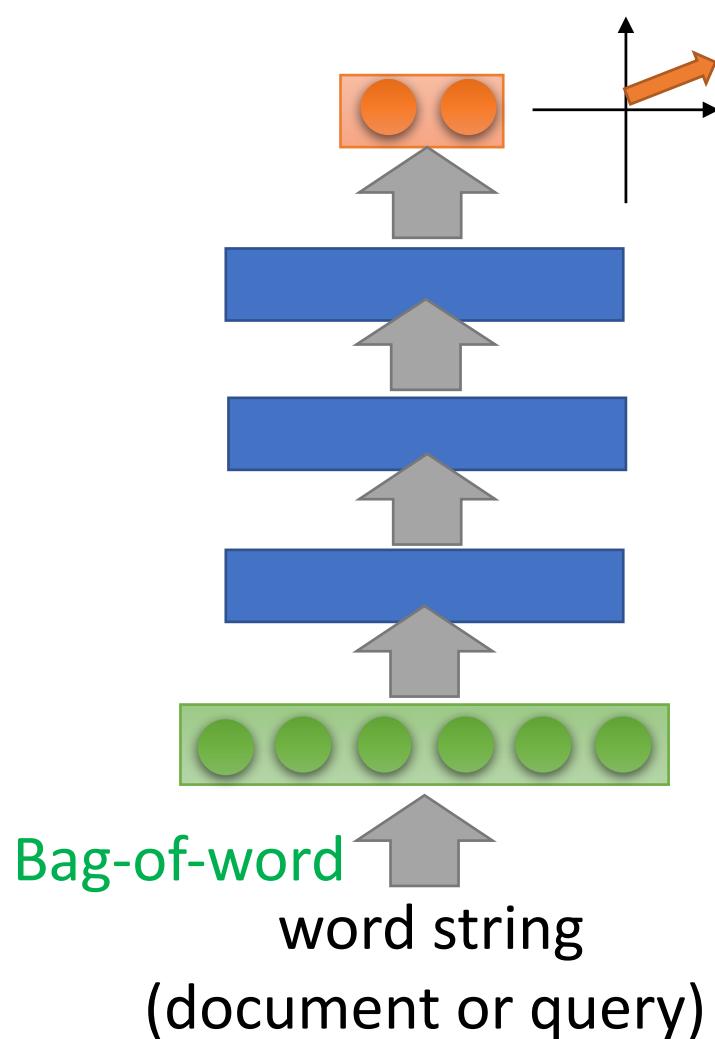
# Information Retrieval (IR)

## ***Vector Space Model + Bag-of-word***



- All the words are treated as discrete tokens.
- Never considered: Different words can have the same meaning, and the same word can have different meanings.

# IR - Semantic Embedding



Reference: Hinton, Geoffrey E., and Ruslan R. Salakhutdinov. "Reducing the dimensionality of data with neural networks." *Science* 313.5786 (2006): 504-507

How to achieve that? (No target .....)

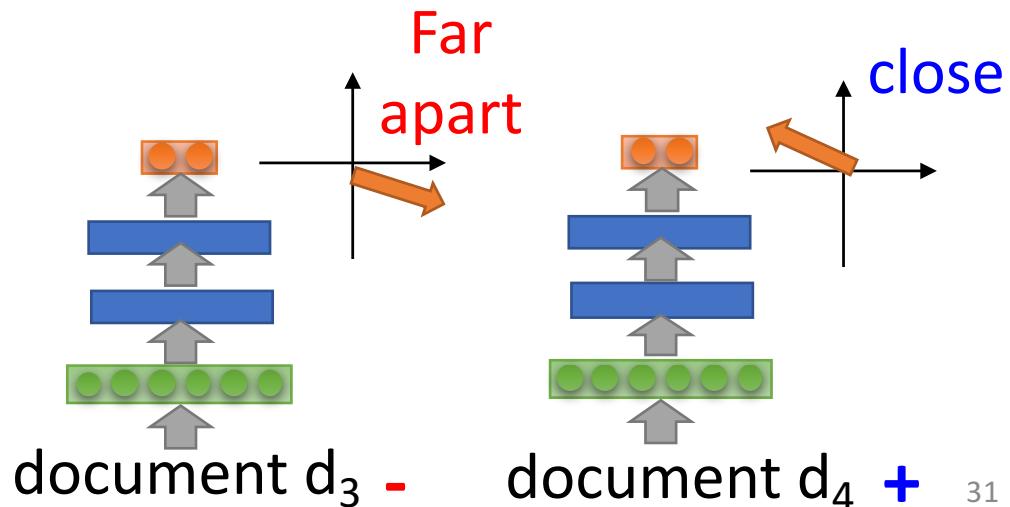
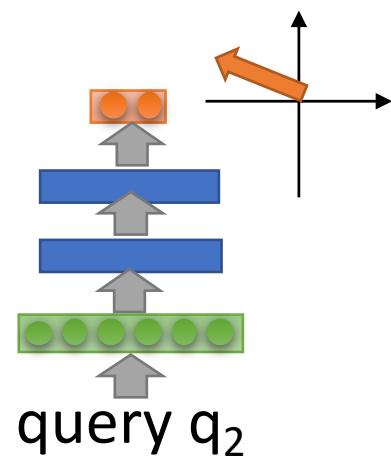
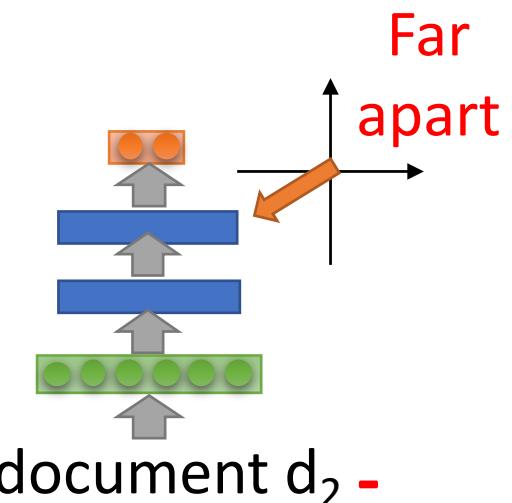
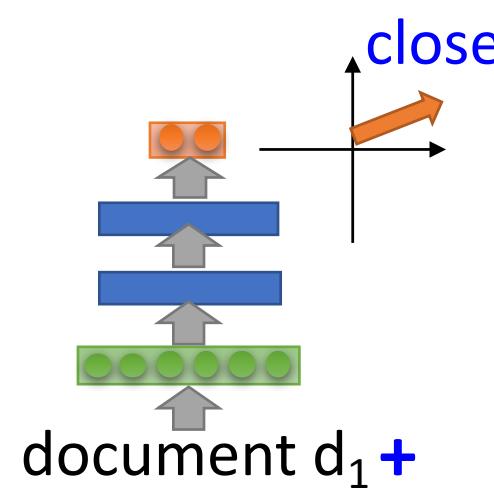
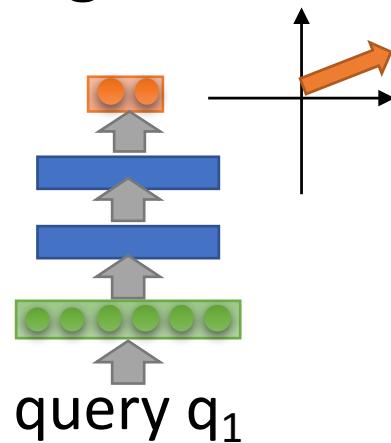
# DSSM

Click-through data:  $q_1 \rightarrow d_1 : + \quad d_2 : -$   
 $q_2 \rightarrow d_3 : - \quad d_4 : +$



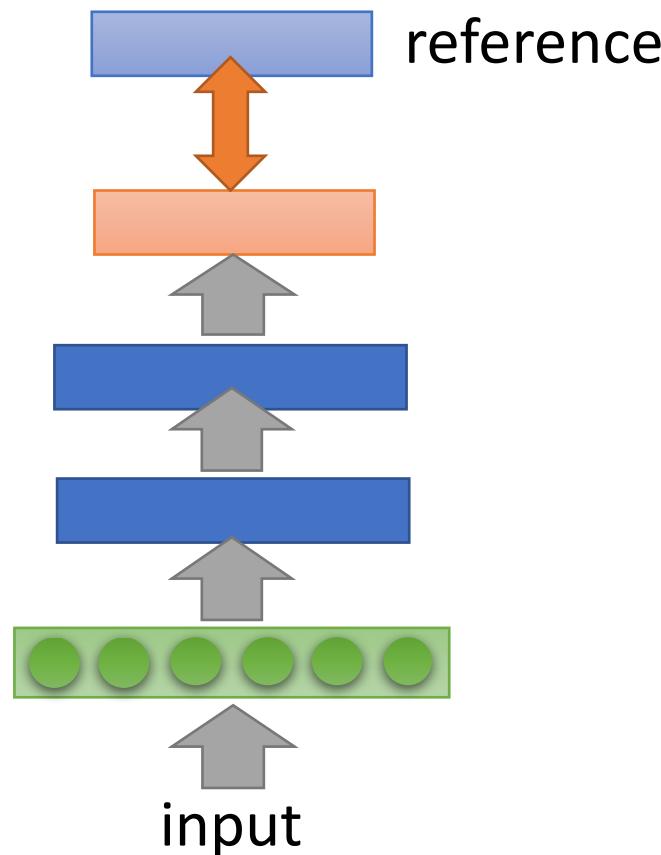
.....

Training:

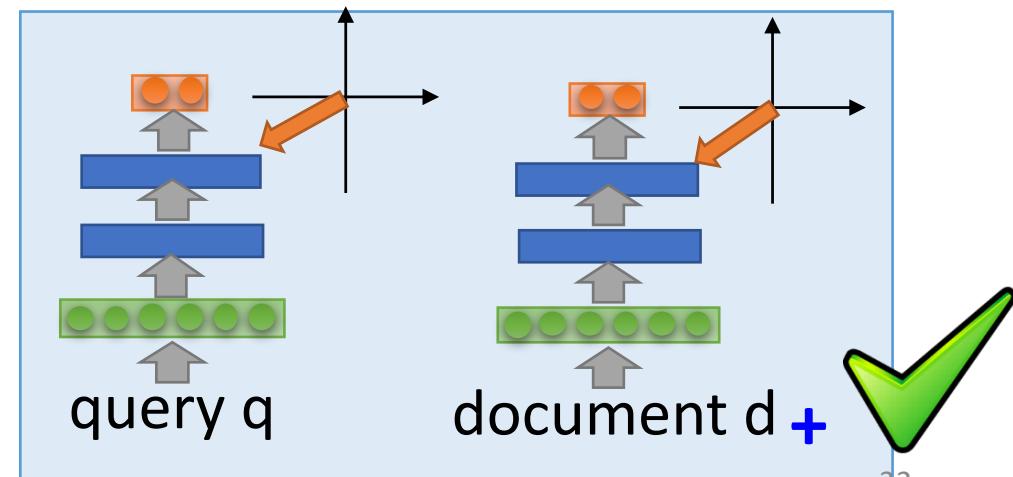
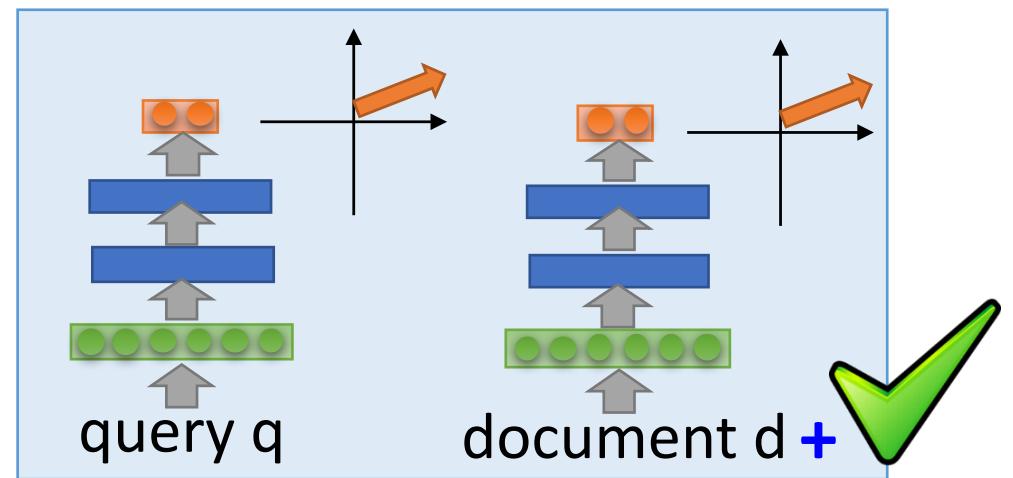


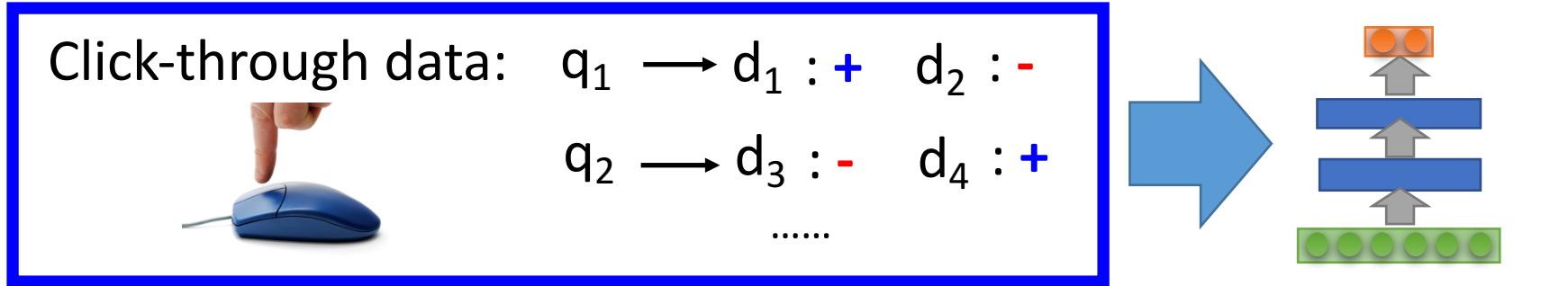
# DSSM v.s. Typical DNN

## Typical DNN

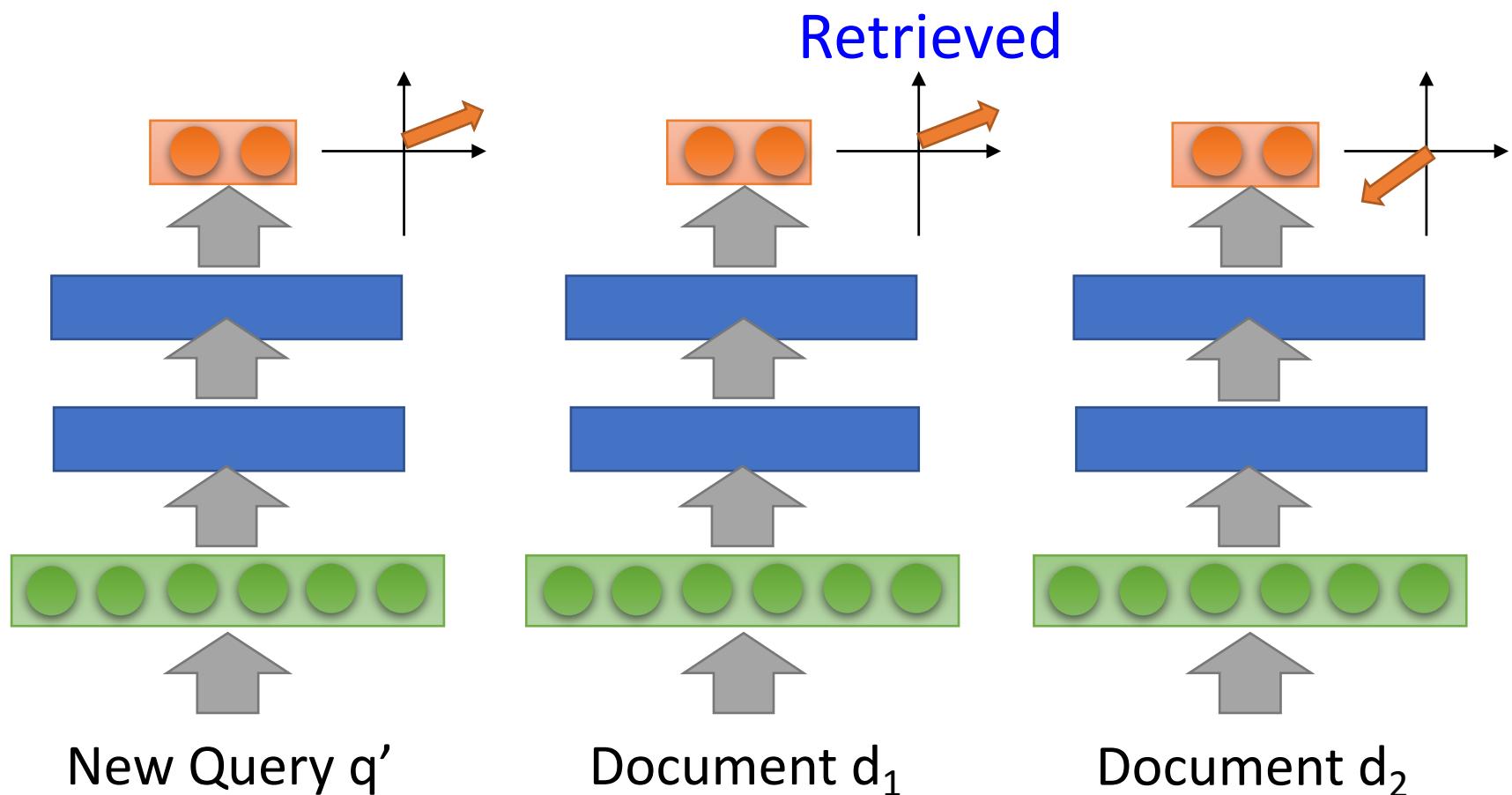


## DSSM



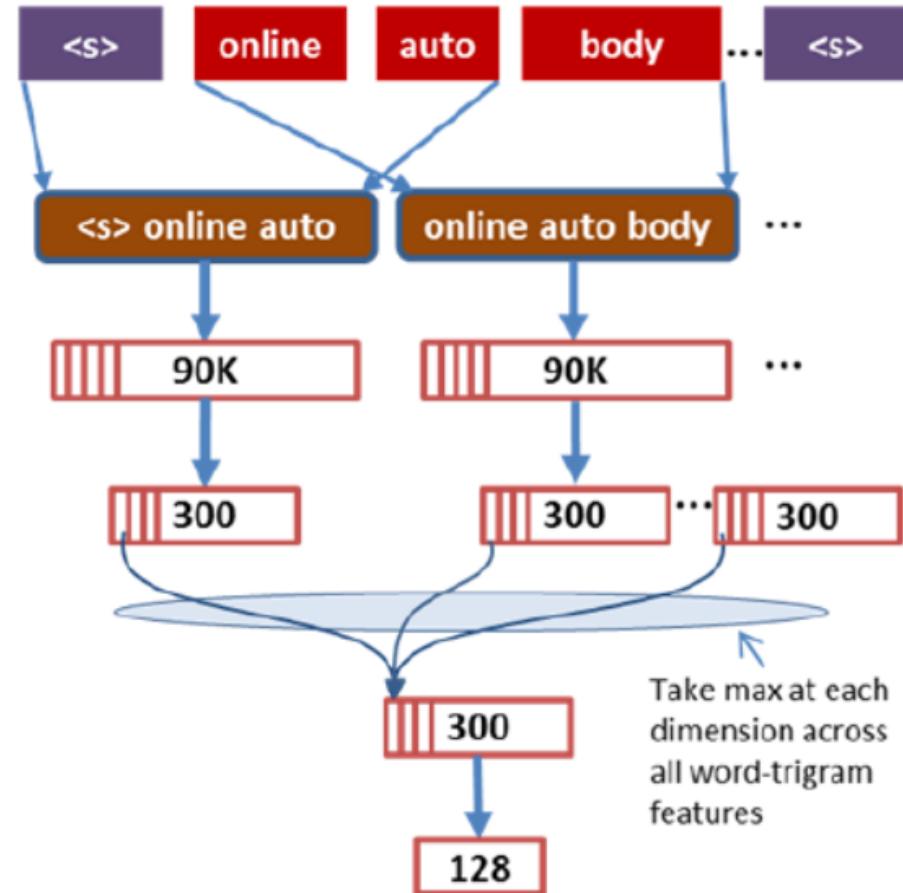


- How to do retrieval?



# Reference

- Huang, Po-Sen, et al. "Learning deep structured semantic models for web search using clickthrough data." ACM, 2013.
- Shen, Yelong, et al. "A latent semantic model with convolutional-pooling structure for information retrieval." ACM, 2014.



# Outline

Deep Structured  
Semantic Model  
(DSSM)

- Application: Information Retrieval (IR)

Recursive  
Neural Network

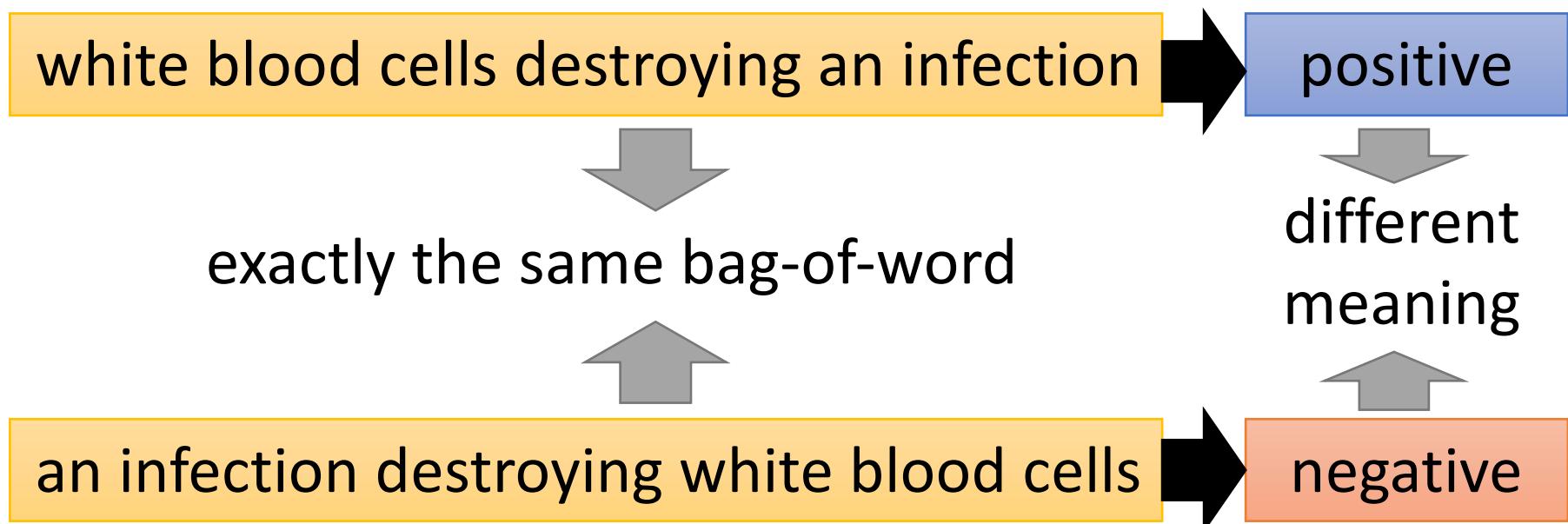
- Application: Sentiment Analysis,  
Sentence Relatedness

Unsupervised

- Paragraph Vector
- Sequence-to-sequence auto-encoder

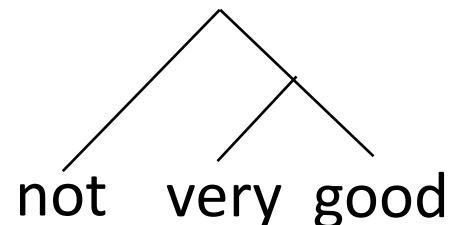
# Recursive Deep Model

- To understand the meaning of a word sequence, the order of the words can not be ignored.



# Recursive Deep Model

syntactic structure



How to do it is out  
of the scope

word sequence:

not

very

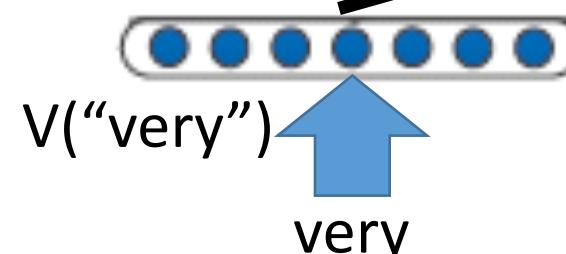
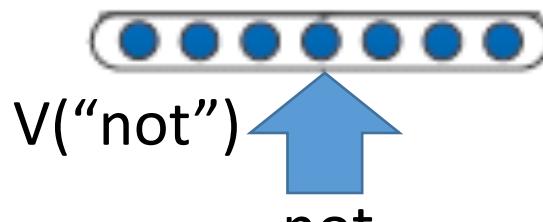
good

# Recursive Deep Model

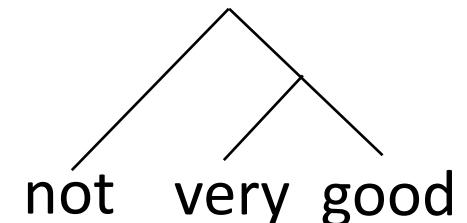
By composing the two meaning, what should the meaning be.

Dimension of word vector =  $|Z|$

Input:  $2 \times |Z|$ , output:  $|Z|$



syntactic structure



Meaning of "very good"



NN

# Recursive Deep Model

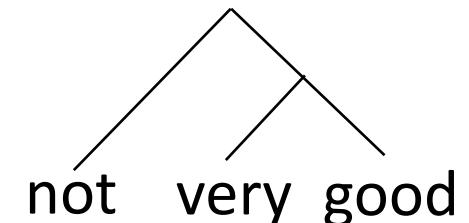
$$V(w_A w_B) \neq V(w_A) + V(w_B)$$

“not”: neutral

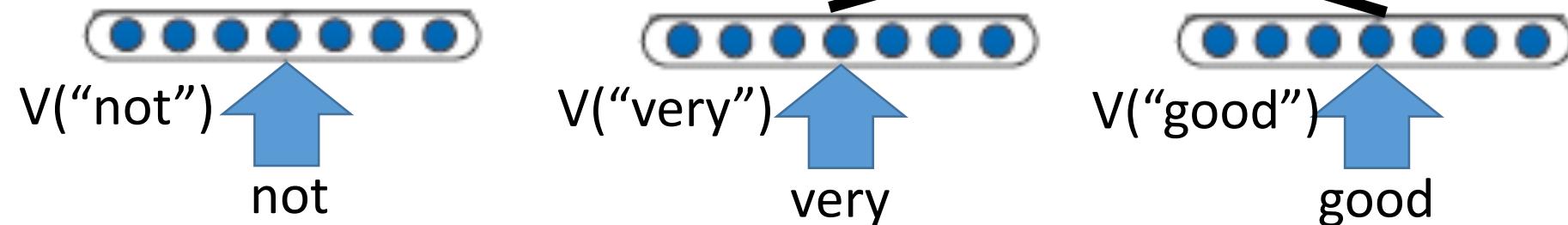
“good”: positive

“not good”: negative

syntactic structure



Meaning of “very good”



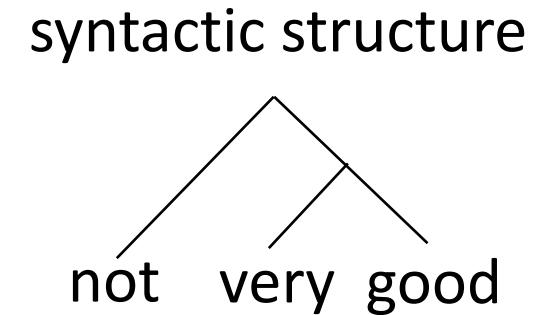
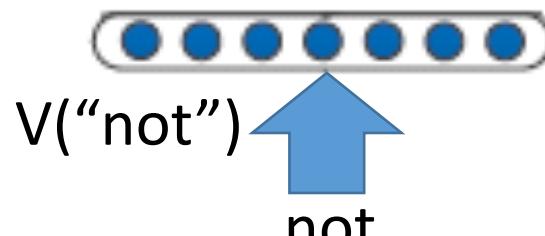
# Recursive Deep Model

$$V(w_A w_B) \neq V(w_A) + V(w_B)$$

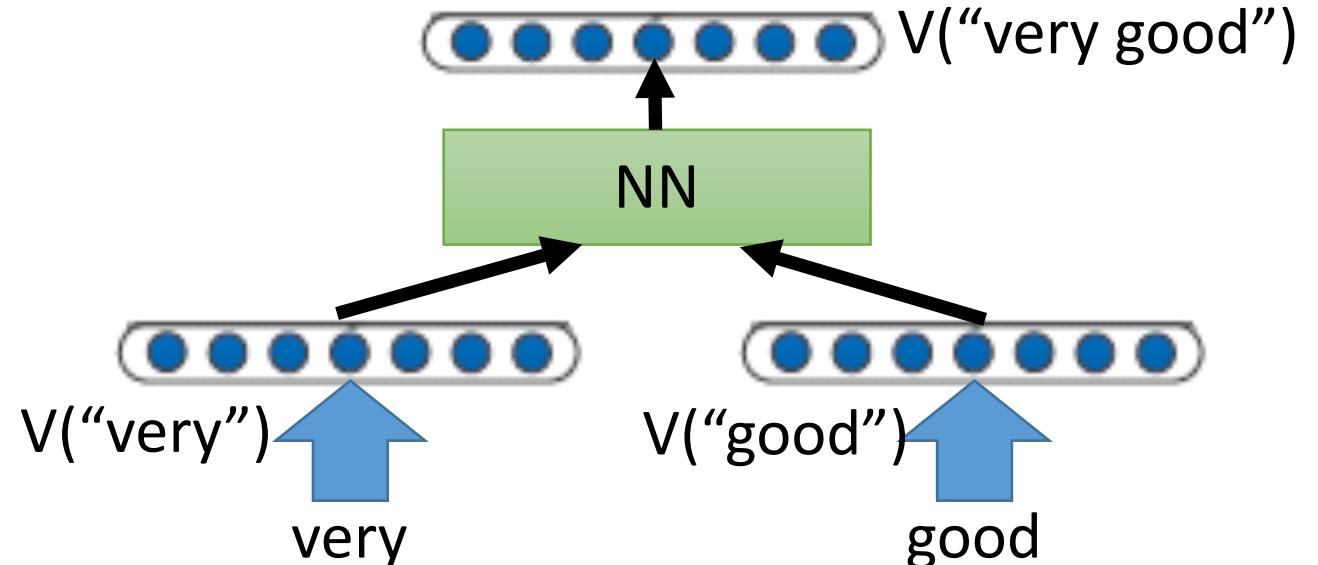
“棒”: positive

“好棒”: positive

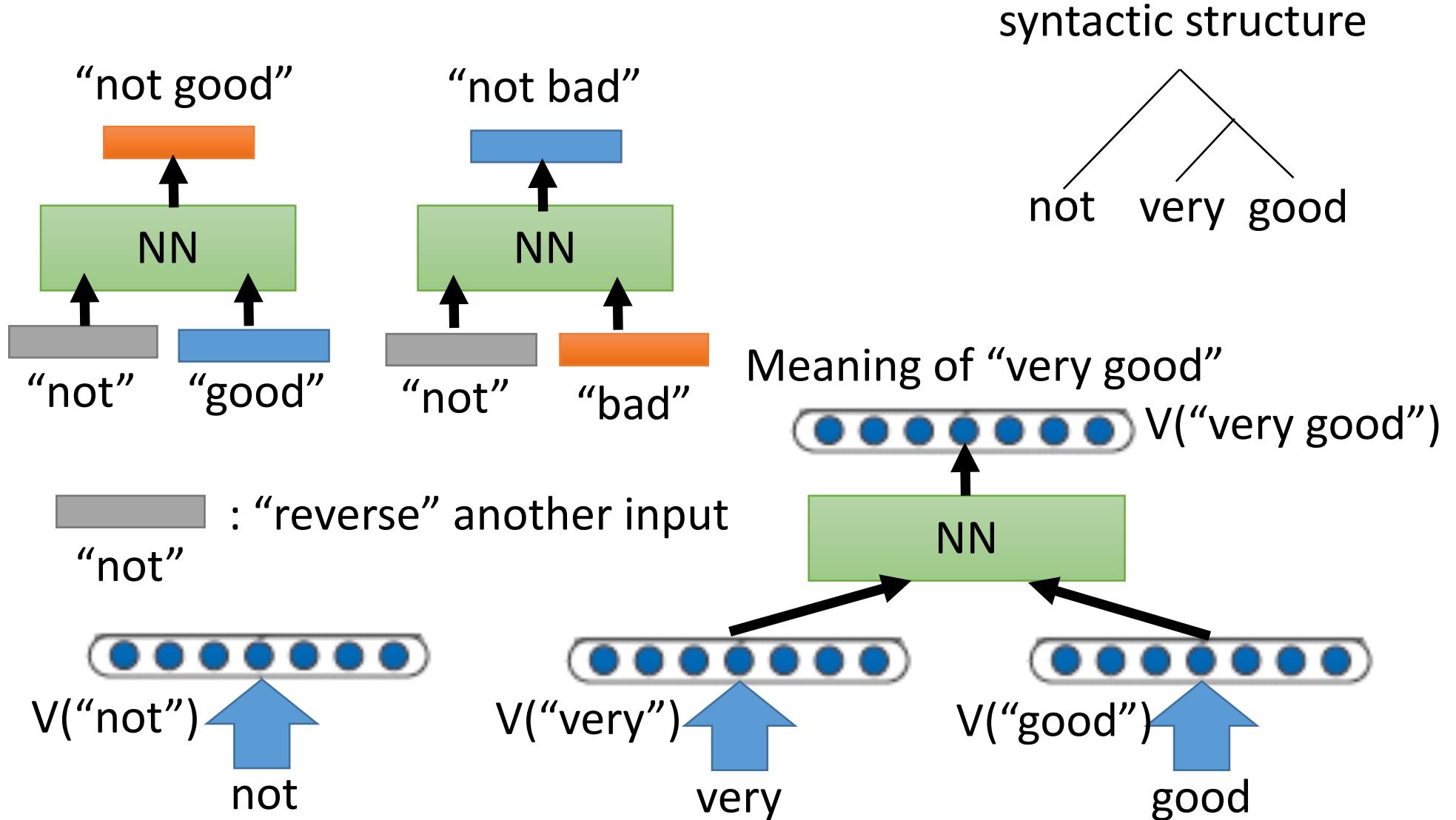
“好棒棒”: negative



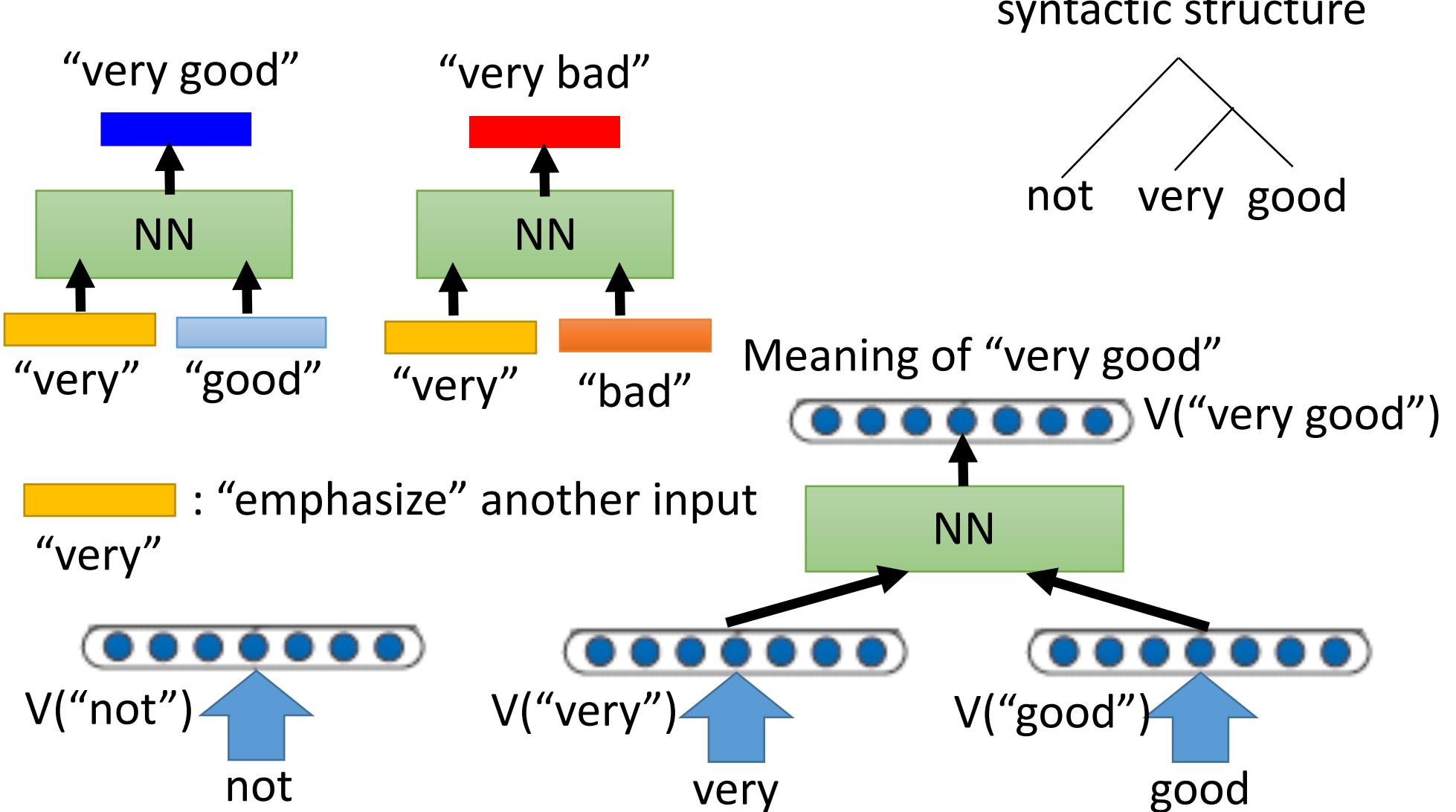
Meaning of “very good”



# Recursive Deep Model

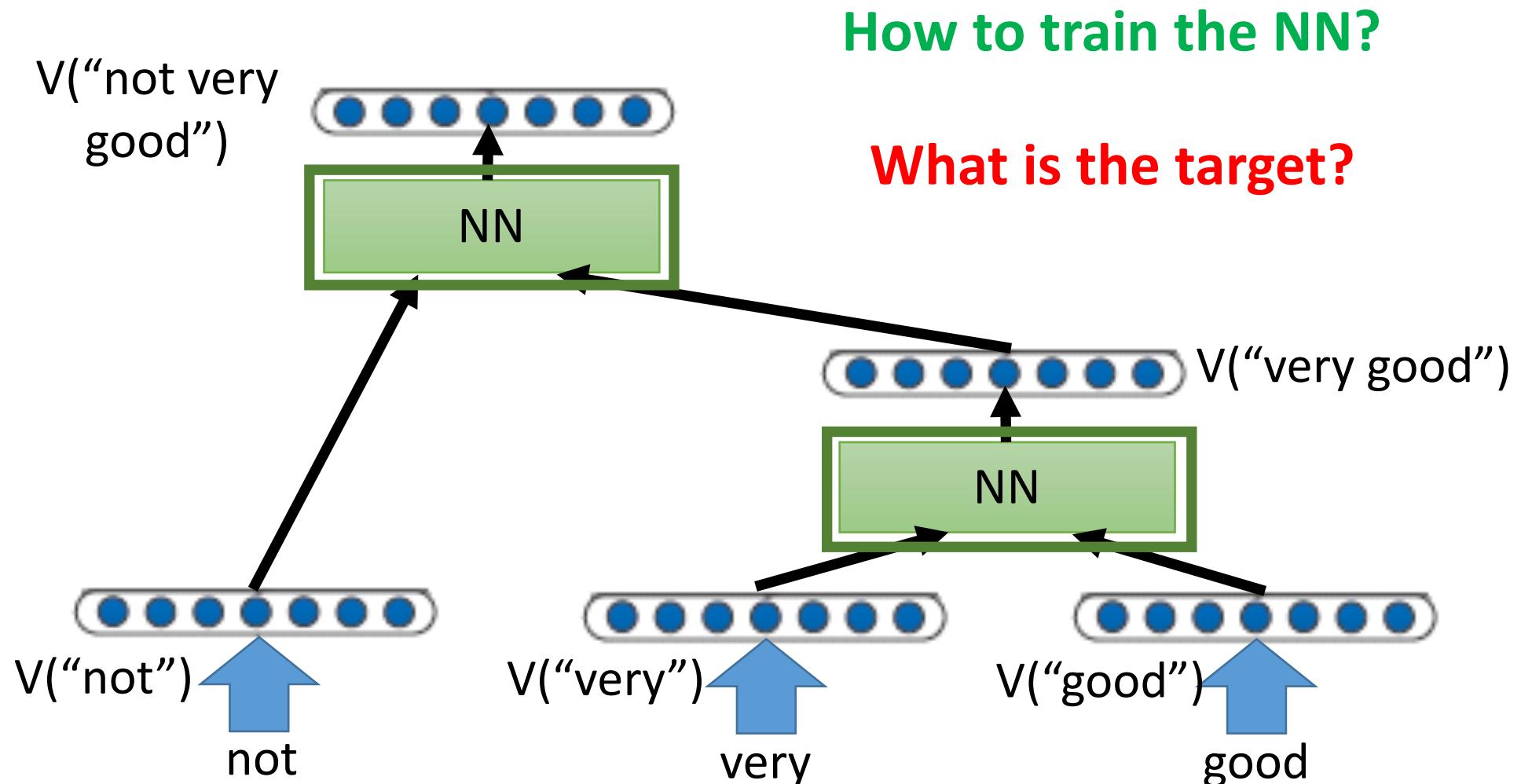


# Recursive Deep Model



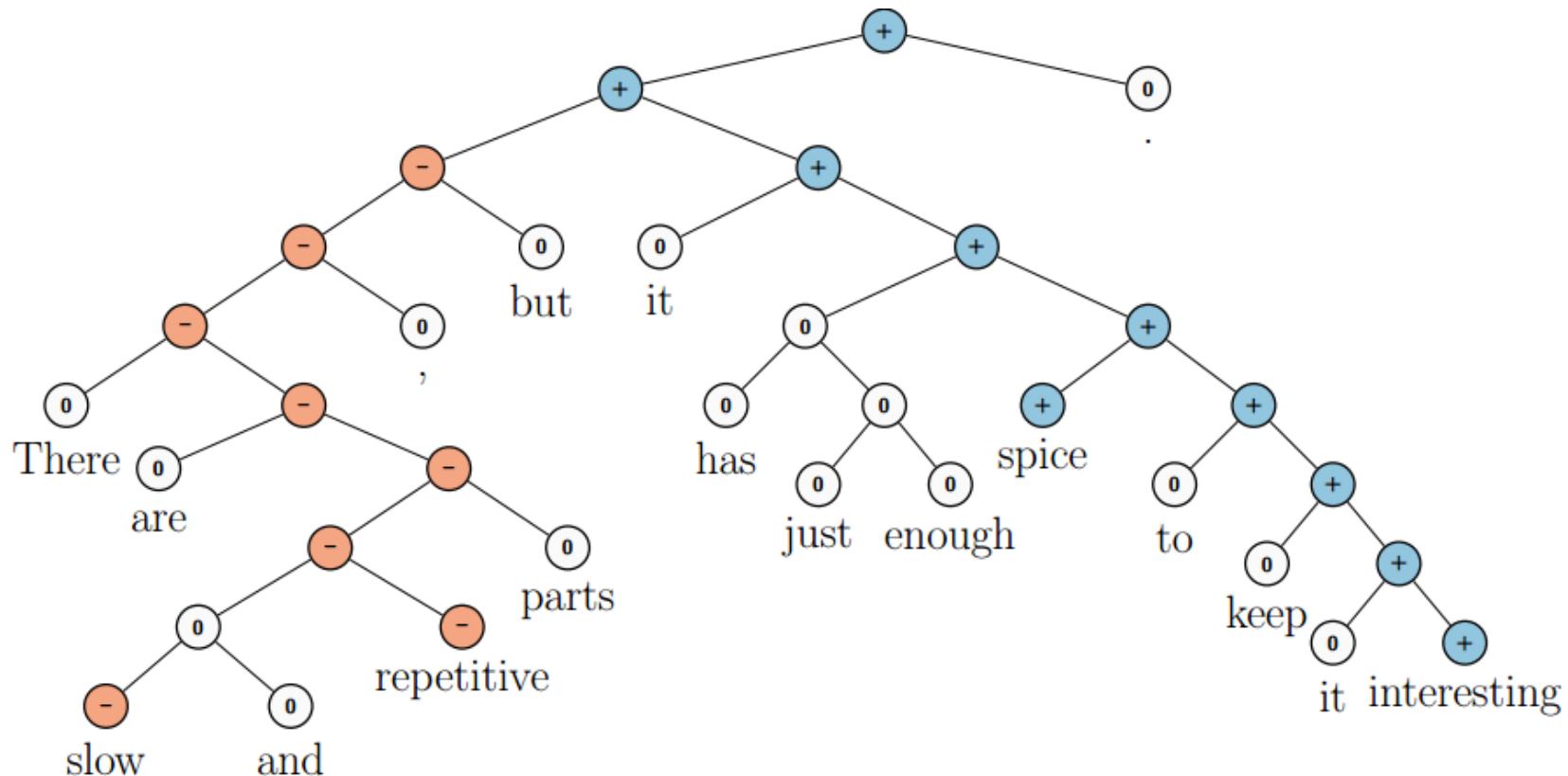
The word order is considered.

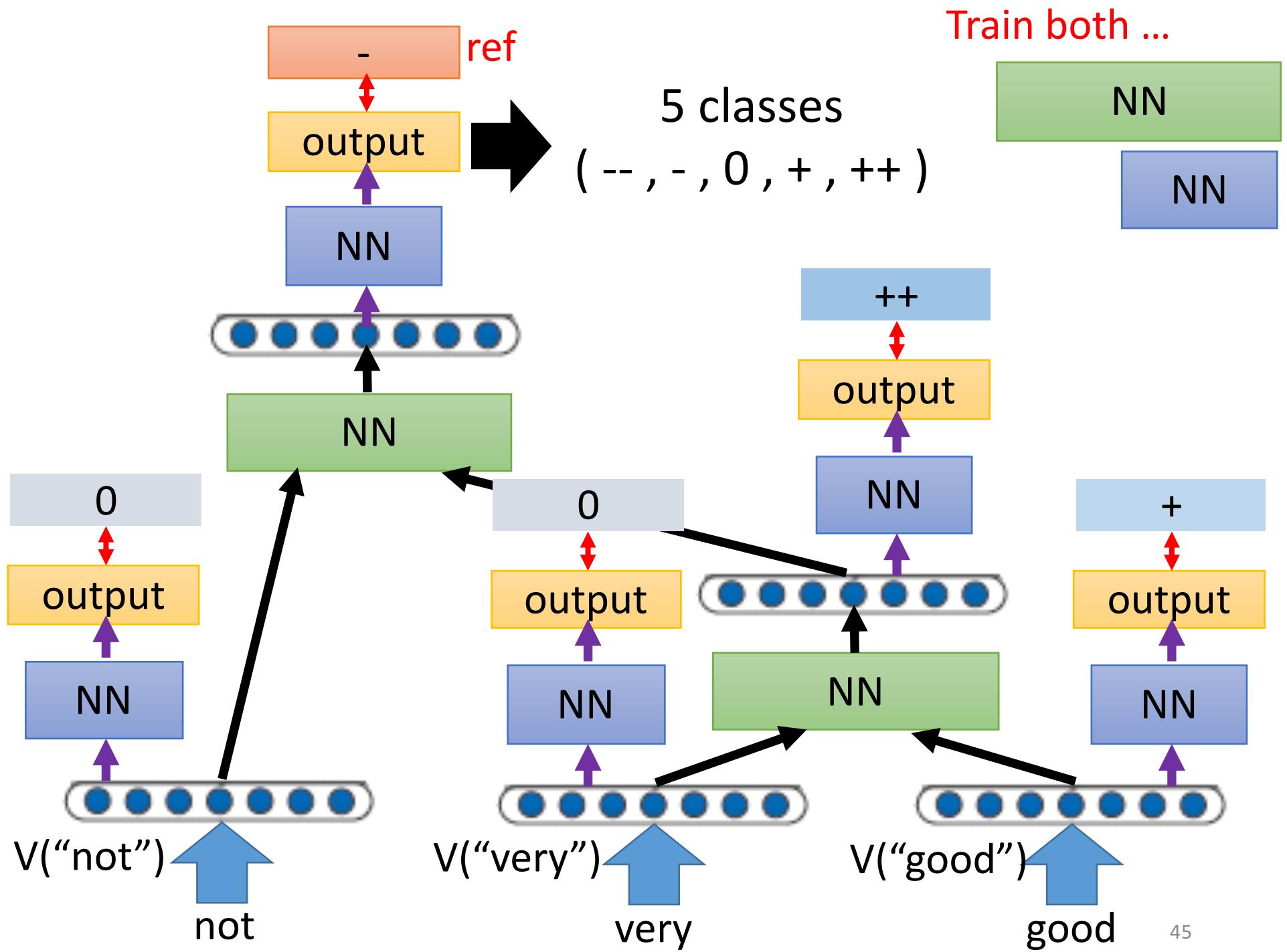
The representation of the sequence will change if the order of the words are changed



# Need a Training Target .....

# 5-class sentiment classification ( -- , - , 0 , + , ++ )





# Outline

Deep Structured  
Semantic Model  
(DSSM)

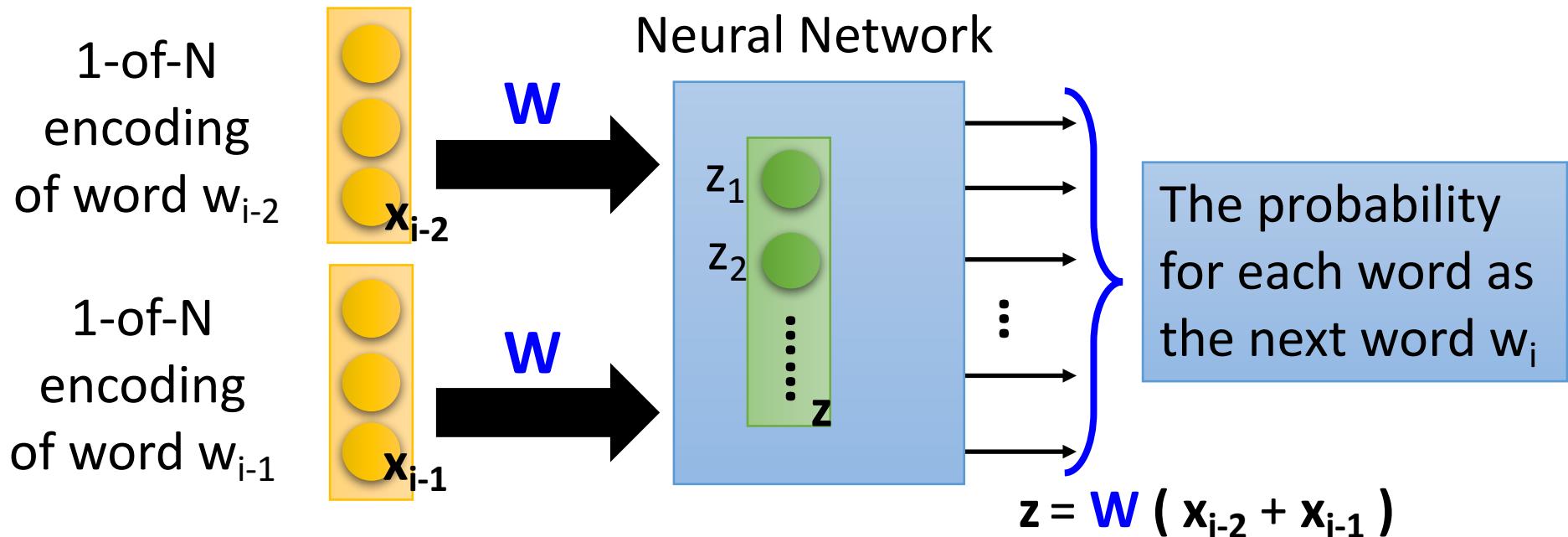
- Application: Information Retrieval (IR)

Recursive  
Neural Network

- Application: Sentiment Analysis,  
Sentence Relatedness

Unsupervised

- Paragraph Vector
- Sequence-to-sequence auto-encoder



Paragraph  $d_1$ : (The paragraph is from "The lord of the ring")

..... 魔君 名叫 索倫 (Sauron) .....

$w_{i-2}$        $w_{i-1}$        $w_i$

$$z = W(x_{i-2} + x_{i-1})$$

Paragraph  $d_2$ : (The paragraph is from "仙五")

..... 魔君 名叫 姜世離 .....

$w_{i-2}$        $w_{i-1}$        $w_i$

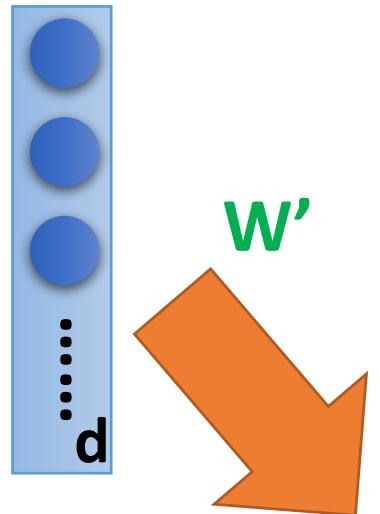
$$z = W(x_{i-2} + x_{i-1})$$

the same → Same output

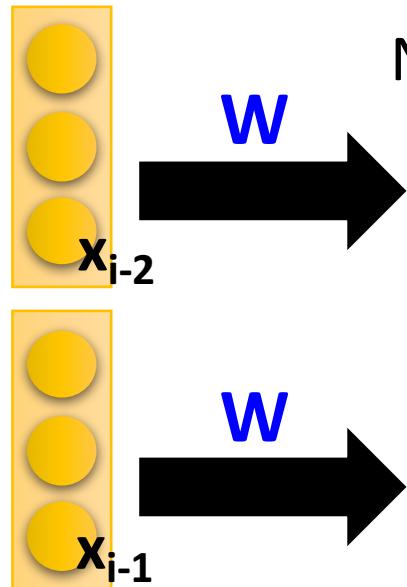
# Paragraph Vector

Le, Quoc, and Tomas Mikolov. "Distributed Representations of Sentences and Documents." ICML, 2014

1-of-N  
encoding  
of paragraph d



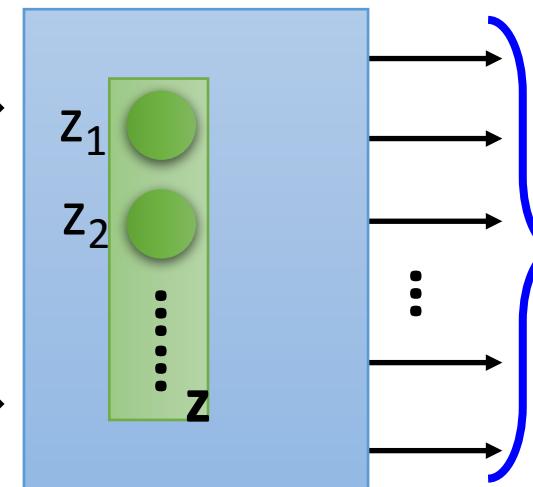
1-of-N  
encoding  
of word  $w_{i-2}$



1-of-N  
encoding  
of word  $w_{i-1}$

Original word vector:  $z = W (x_{i-2} + x_{i-1})$   
Paragraph vector:  
$$z = W (x_{i-2} + x_{i-1}) + W' d$$

Neural Network



The probability  
for each word  
as the next word  $w_i$

# Paragraph Vector

Le, Quoc, and Tomas Mikolov. "Distributed Representations of Sentences and Documents." ICML, 2014

Original word vector:

$$z = \mathbf{W} (x_{i-2} + x_{i-1})$$

Paragraph vector:

$$z = \mathbf{W} (x_{i-2} + x_{i-1}) + \mathbf{W}' d$$

Then error of the prediction can be explained by the meaning of the paragraphs.

Paragraph  $d_1$ : (The paragraph is related to  
“The lord of the ring”)

$$\dots \dots \text{魔君} \quad \text{名叫} \quad \underline{\text{索倫 (Sauron)}} \quad \dots \dots$$
$$w_{i-2} \quad w_{i-1} \quad \underline{w_i}$$
$$z = \mathbf{W} (x_{i-2} + x_{i-1}) + \mathbf{W}' d_1$$

Paragraph  $d_2$ : (The document is related to  
“仙五”)

$$\dots \dots \text{魔君} \quad \text{名叫} \quad \underline{\text{姜世離}} \quad \dots \dots$$
$$w_{i-2} \quad w_{i-1} \quad \underline{w_i}$$
$$z = \mathbf{W} (x_{i-2} + x_{i-1}) + \mathbf{W}' d_2$$

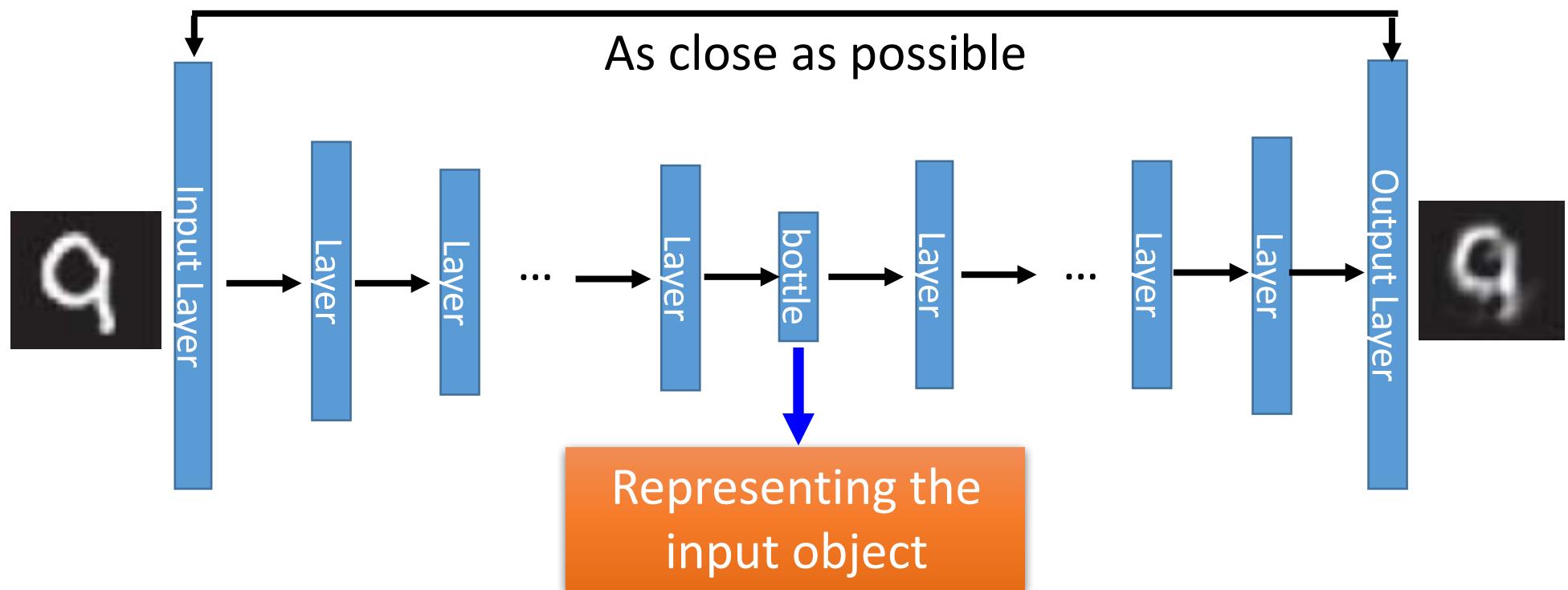
different

Paragraph vector of  $d: V(d) = \mathbf{W}' d$

Meaning of the paragraph

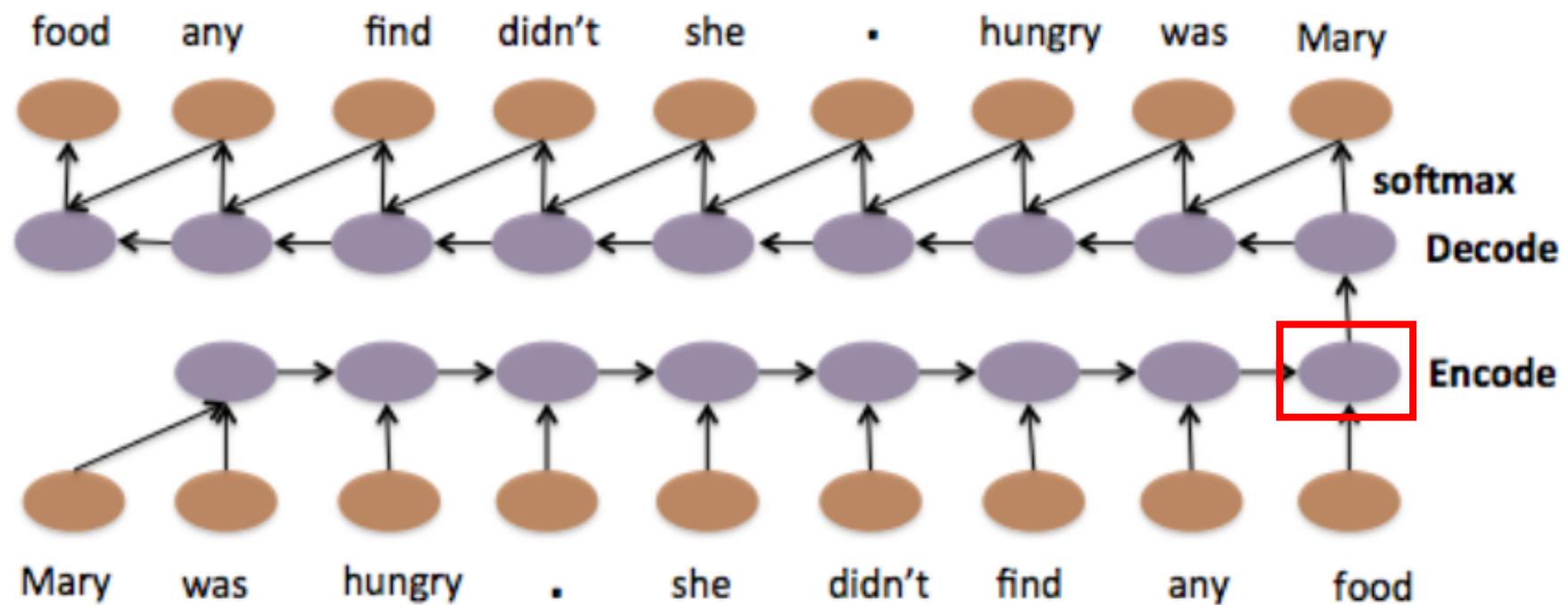
# Sequence-to-sequence Auto-encoder

- Original Auto-encoder



Reference: Hinton, Geoffrey E., and Ruslan R. Salakhutdinov. "Reducing the dimensionality of data with neural networks." *Science* 313.5786 (2006): 504-507

# Sequence-to-sequence Auto-encoder



Li, Jiwei, Minh-Thang Luong, and Dan Jurafsky. "A hierarchical neural autoencoder for paragraphs and documents." *arXiv preprint arXiv:1506.01057*(2015).

# Summary

Deep Structured  
Semantic Model  
(DSSM)

- Application: Information Retrieval (IR)

Recursive  
Neural Network

- Application: Sentiment Analysis,  
Sentence Relatedness

Unsupervised

- Paragraph Vector
- Sequence-to-sequence auto-encoder

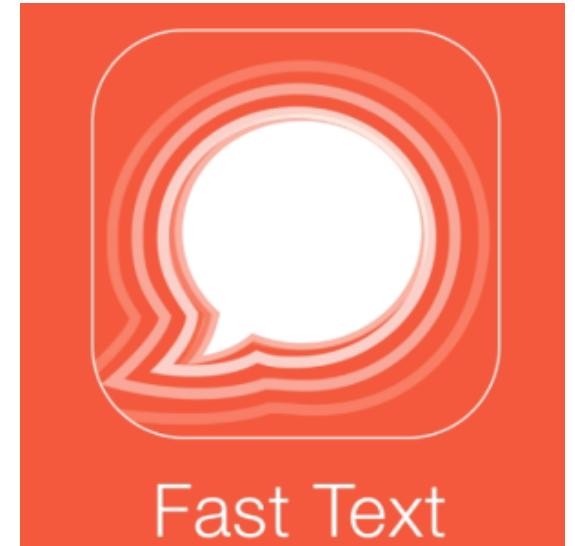
# Tools

- 中文斷詞
  - Jieba
- 英文stopwords
  - NLTK stopwords
- Easiest way to use word2vec is via the Gensim library for Python (tends to be slowish, even though it tries to use C optimizations like Cython, NumPy)

<https://radimrehurek.com/gensim/models/word2vec.html>

# FastText

- Library for fast text representation and classification
    - Courtesy of Facebook Research
  - Recent state-of-the-art English word vectors.
  - Word vectors for 157 languages trained on Wikipedia and Crawl.
  - Models for language identification and various supervised tasks.
- 
- [https://github.com/facebookresearch/fastText?utm\\_source=mybridge&utm\\_medium=blog&utm\\_campaign=read\\_more](https://github.com/facebookresearch/fastText?utm_source=mybridge&utm_medium=blog&utm_campaign=read_more)



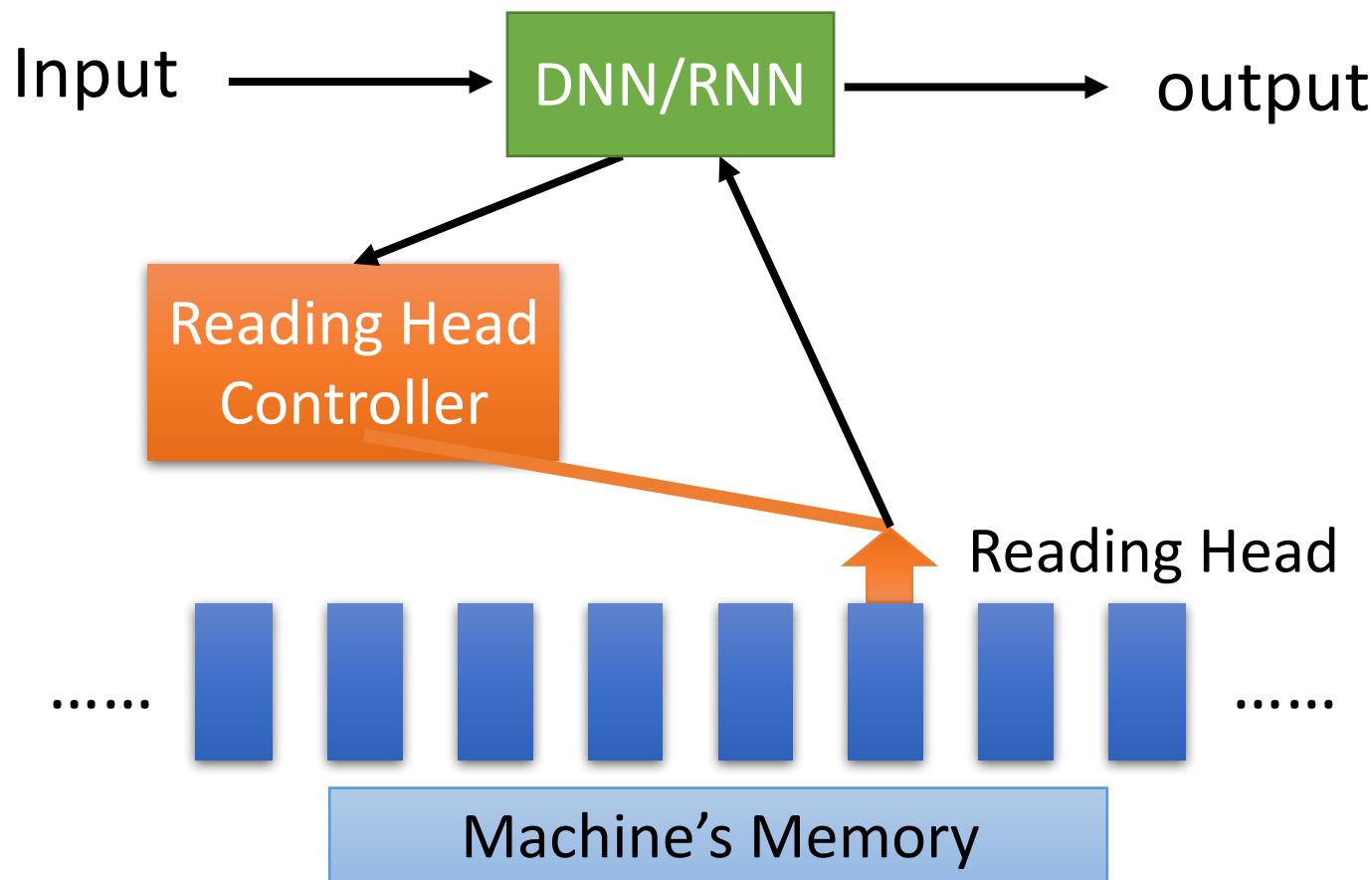
# Fairseq

- sequence-to-sequence learning toolkit for Torch from Facebook AI Research
- English to French, English to German and English to Romanian translation

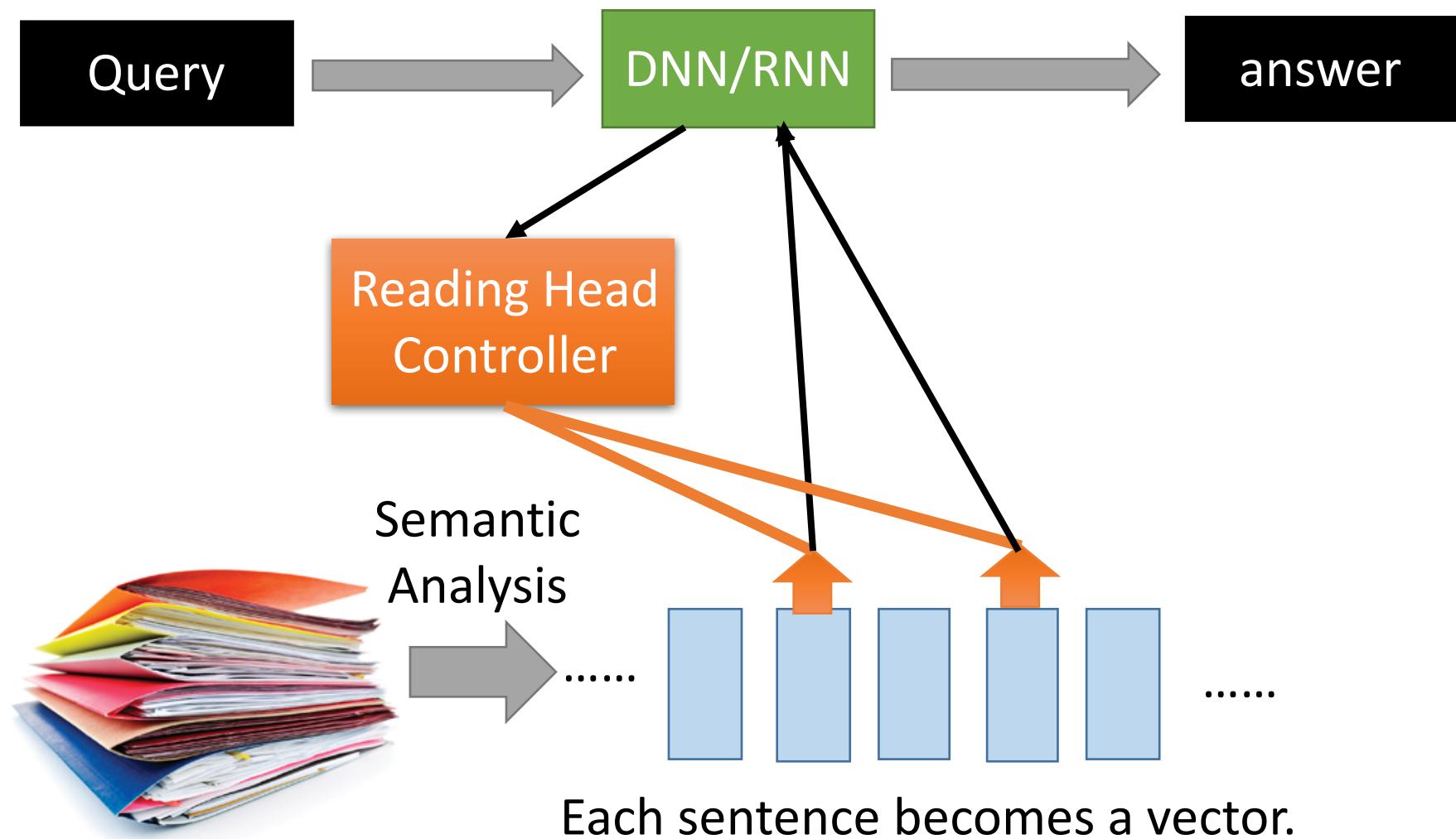
.. la maison de Léa <end> ..

# Attention-based Model

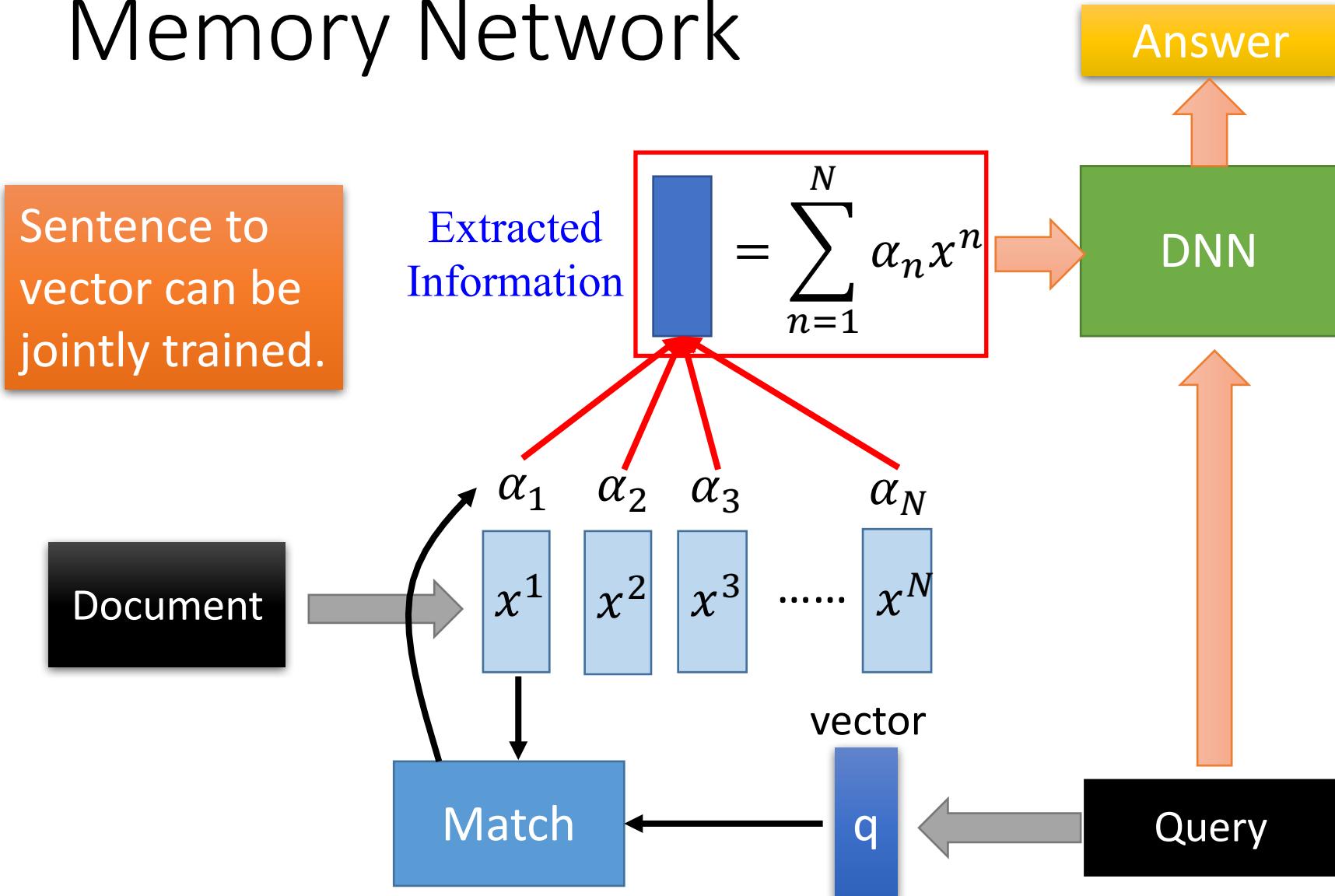
# External Memory



# Reading Comprehension

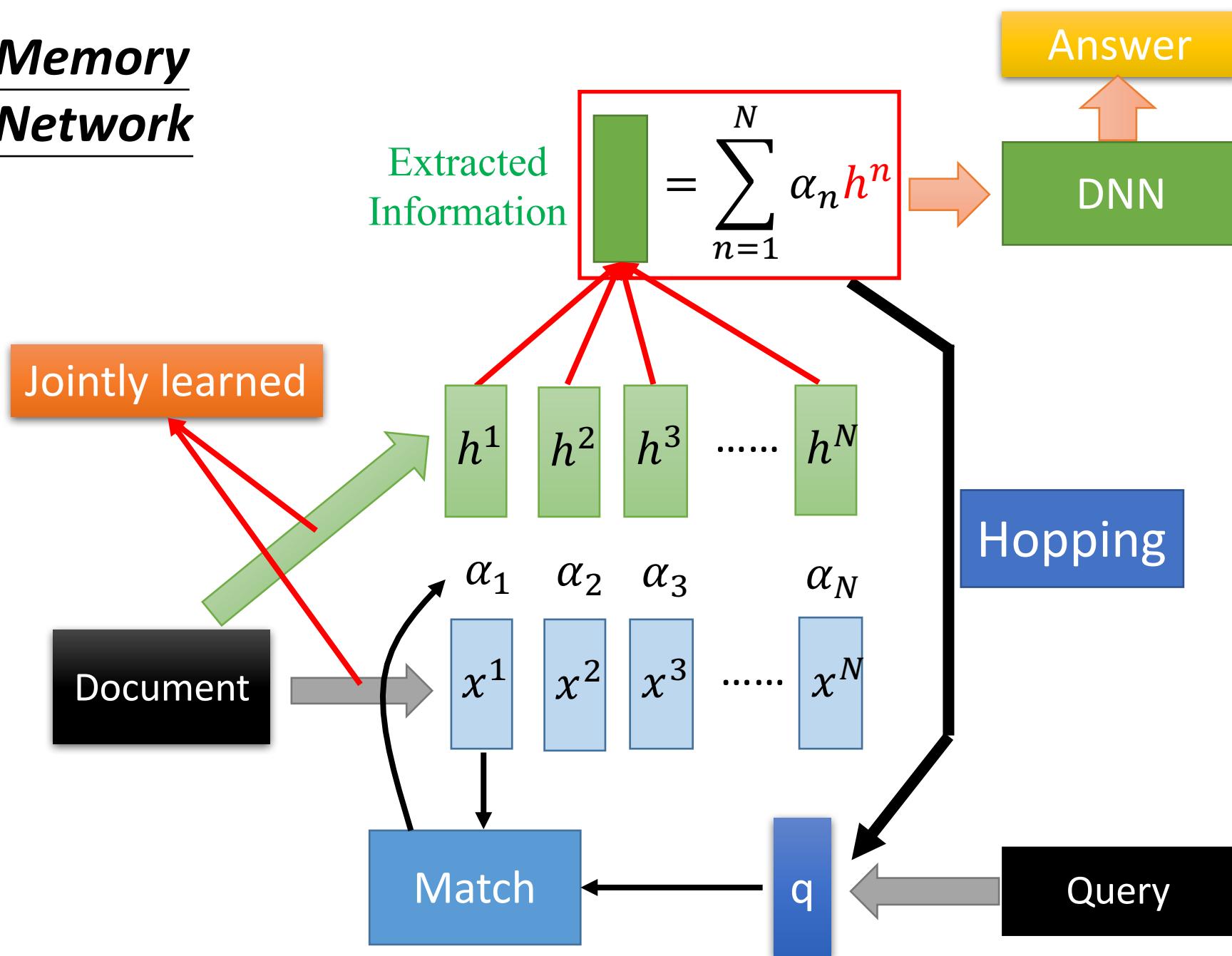


# Memory Network

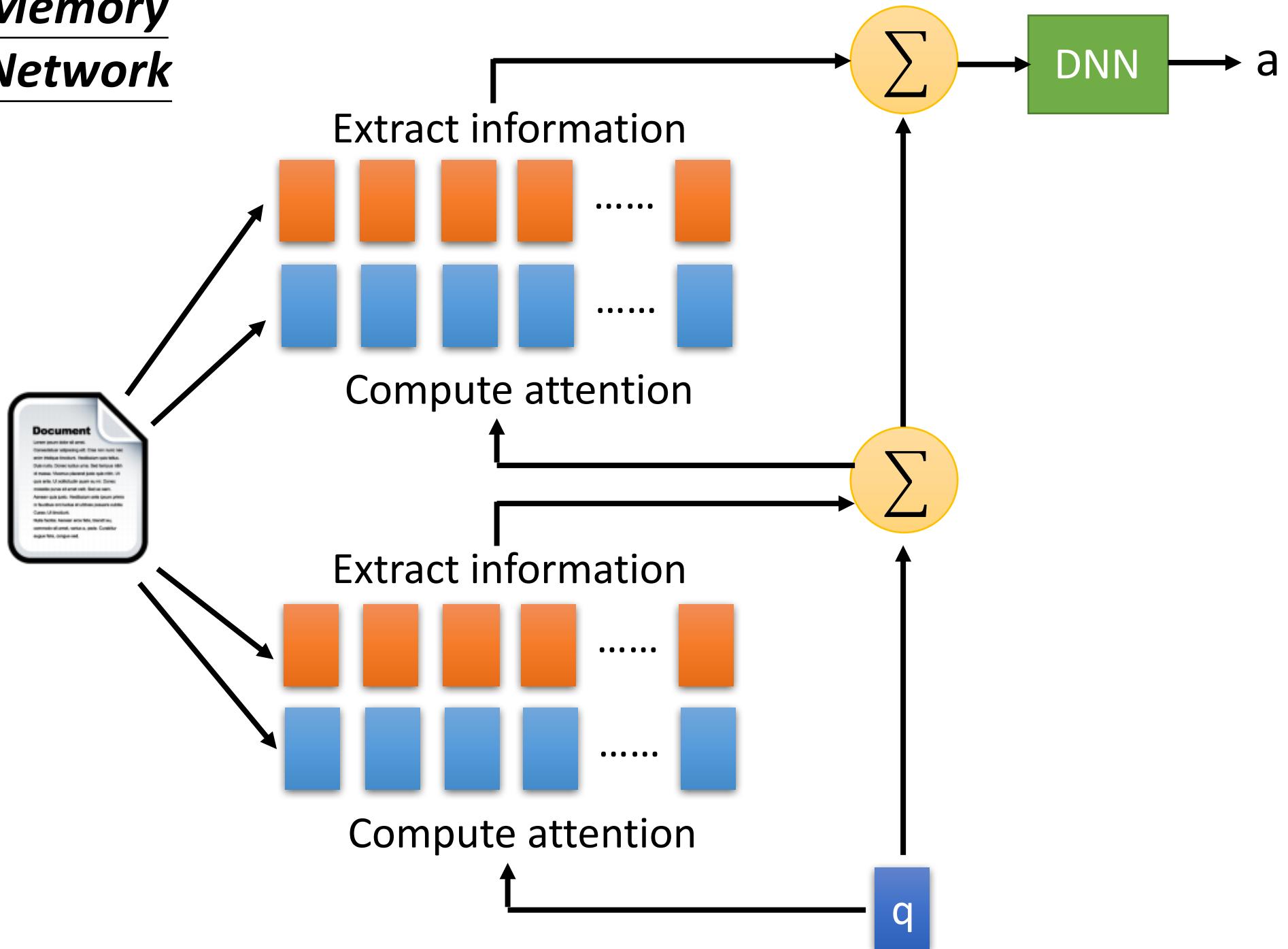


Sainbayar Sukhbaatar, Arthur Szlam, Jason Weston, Rob Fergus, "End-To-End Memory Networks", NIPS, 2015

# Memory Network



# Memory Network



# Multiple-hop

- End-To-End Memory Networks. S. Sukhbaatar, A. Szlam, J. Weston, R. Fergus. NIPS, 2015.

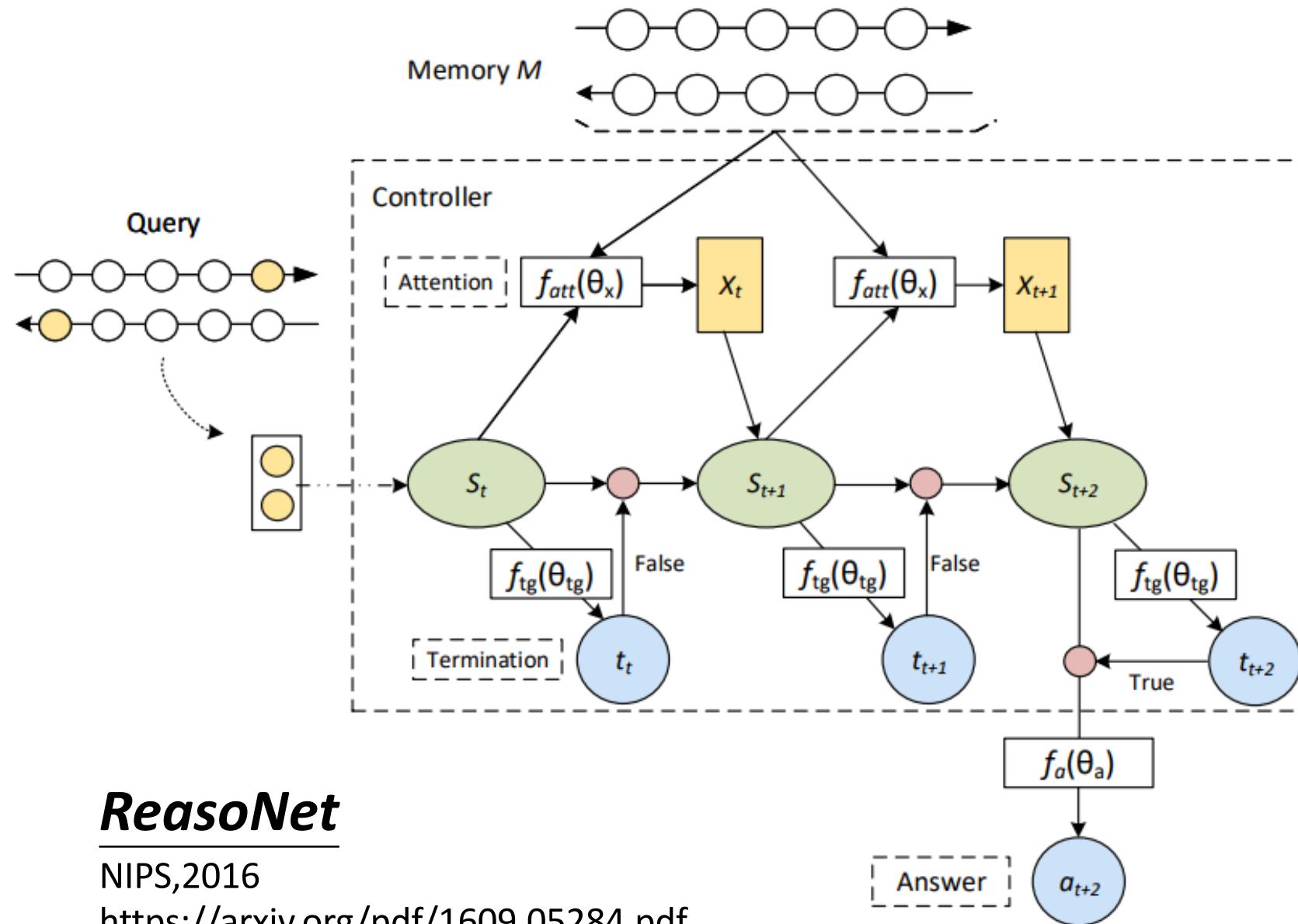
The position of reading head:

Story (16: basic induction)	Support	Hop 1	Hop 2	Hop 3
Brian is a frog.	yes	0.00	0.98	0.00
Lily is gray.		0.07	0.00	0.00
Brian is yellow.	yes	0.07	0.00	1.00
Julius is green.		0.06	0.00	0.00
Greg is a frog.	yes	0.76	0.02	0.00
<b>What color is Greg? Answer: yellow Prediction: yellow</b>				

Keras has example:

[https://github.com/fchollet/keras/blob/master/examples/babi\\_memnn.py](https://github.com/fchollet/keras/blob/master/examples/babi_memnn.py)

# Multiple-hop

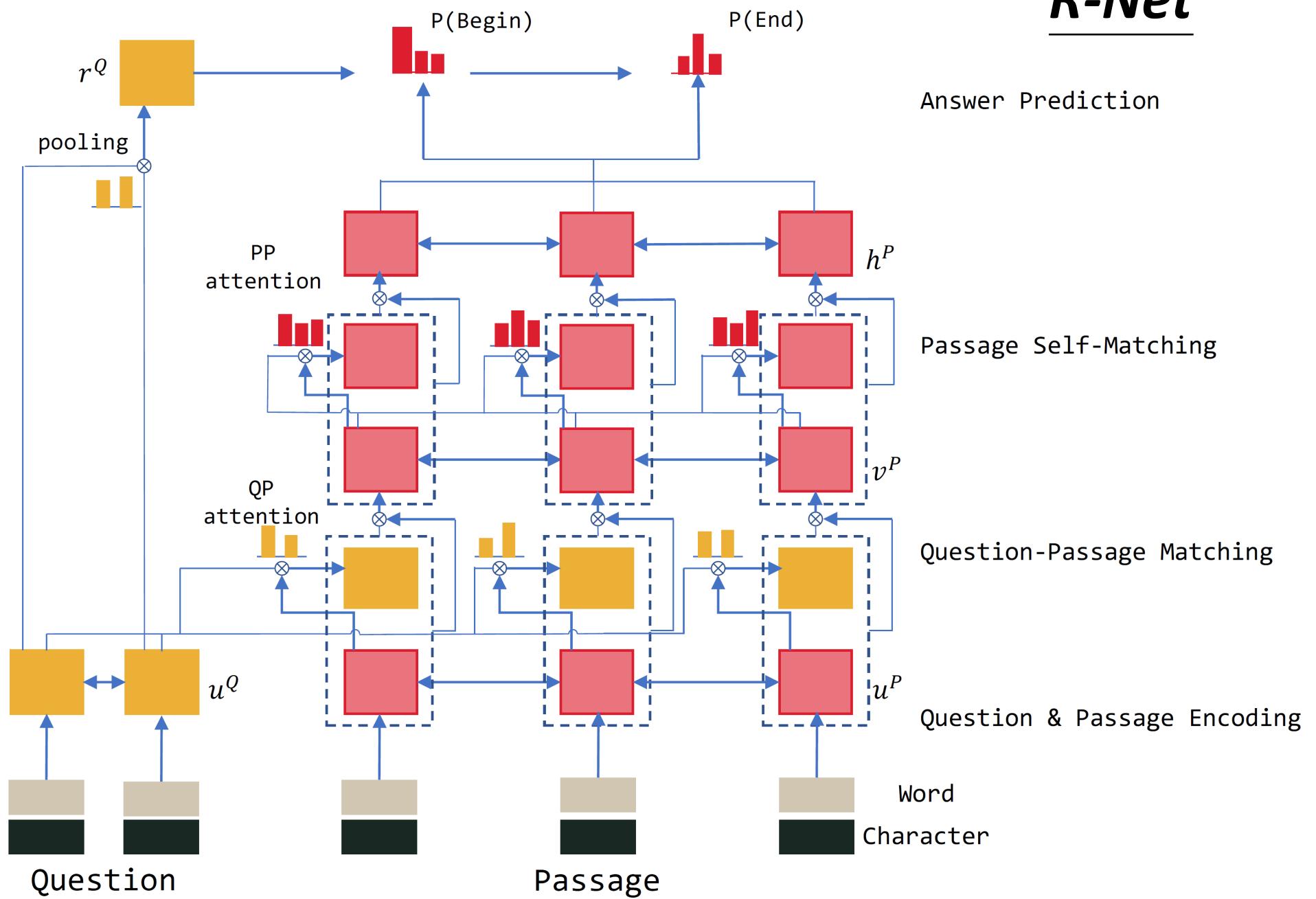


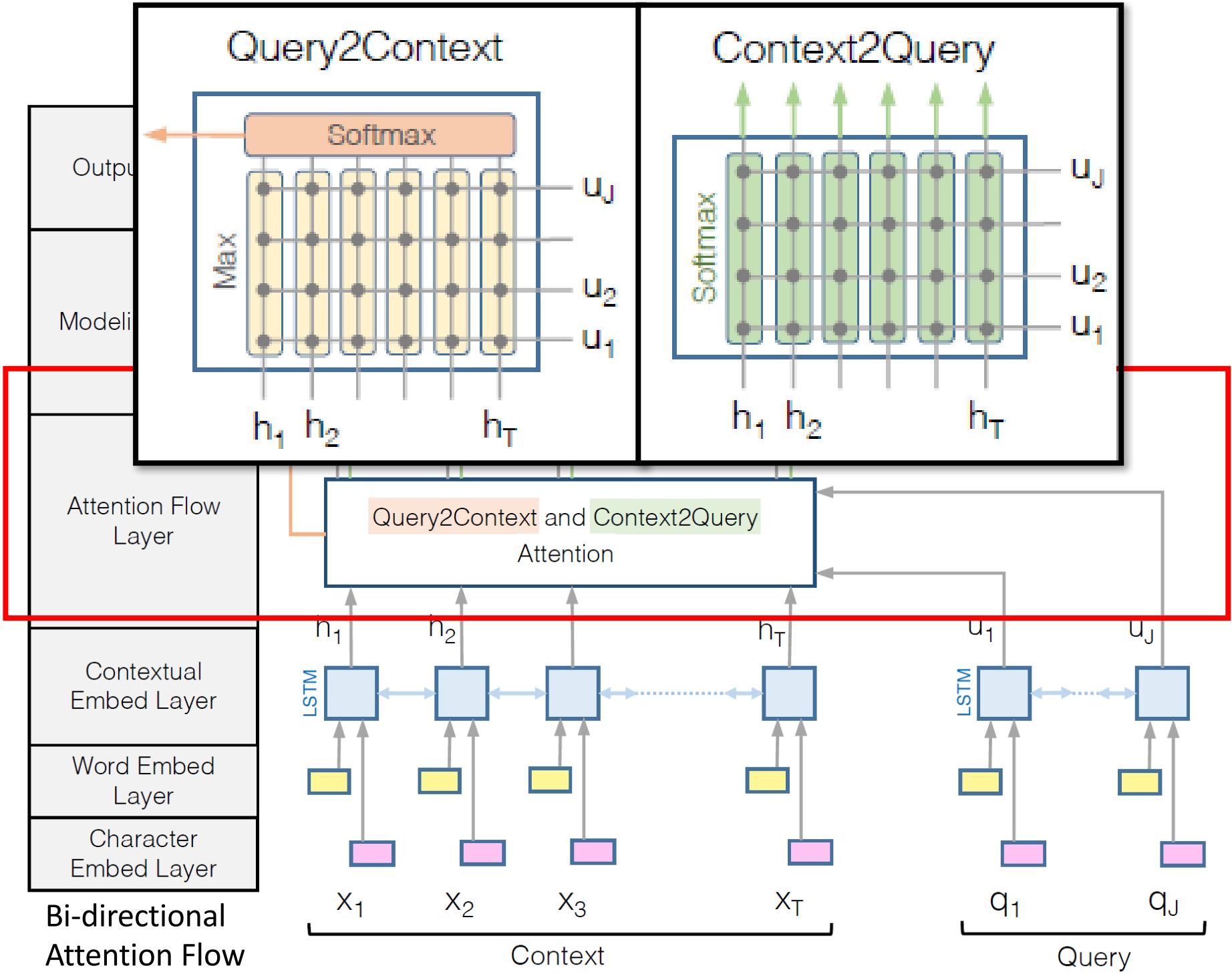
**ReasoNet**

NIPS, 2016

<https://arxiv.org/pdf/1609.05284.pdf>

# R-Net





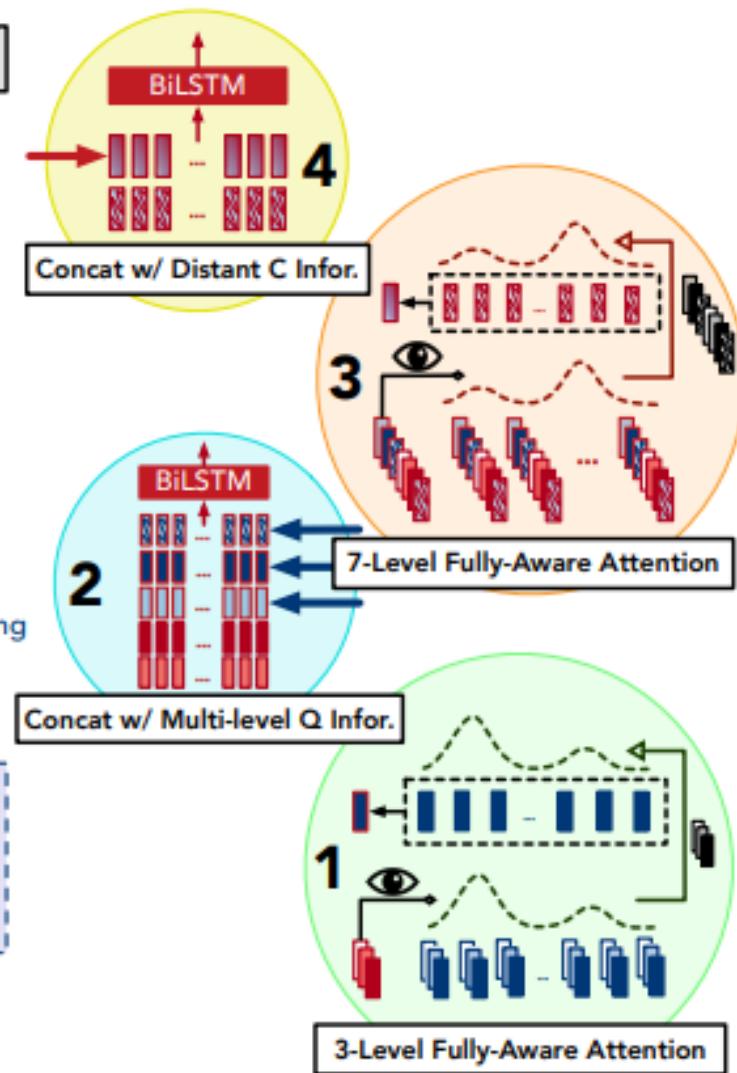
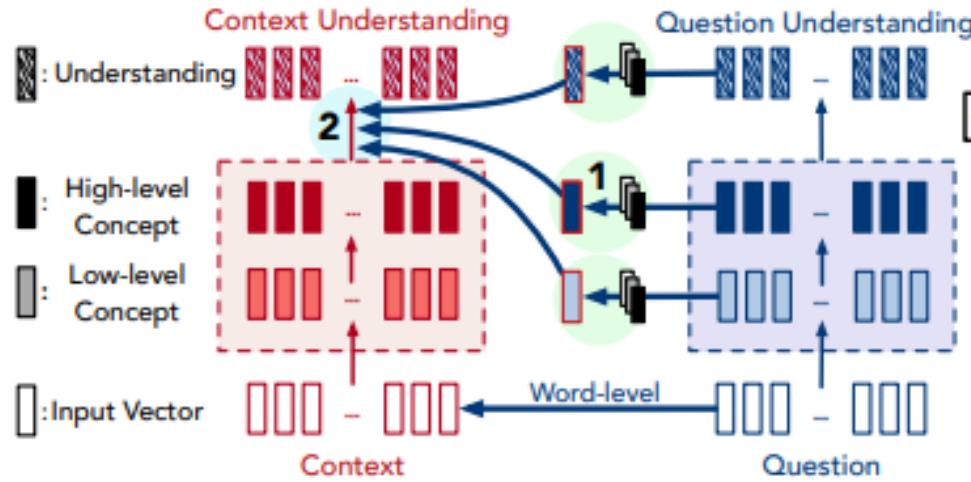
## Fully-Aware Fusion Network

### Fully-Aware Self-Boosted Fusion

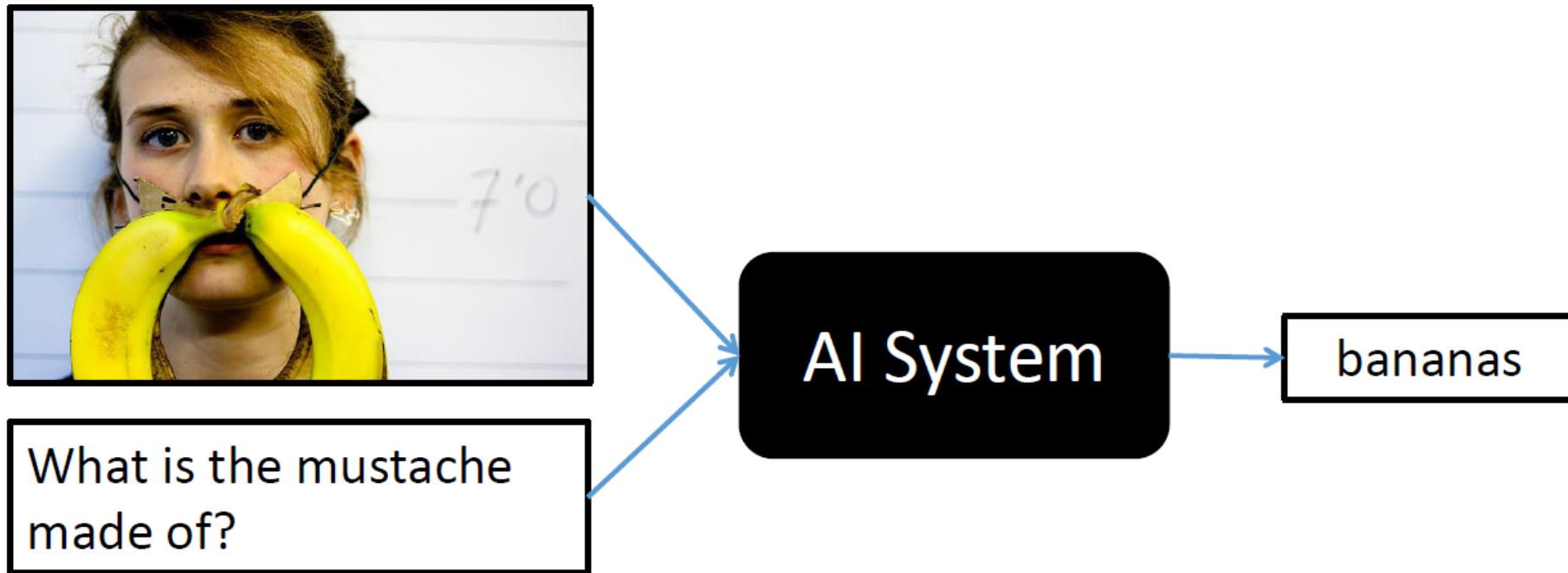


The above can be used to capture long range info.

### Fully-Aware Multi-level Fusion

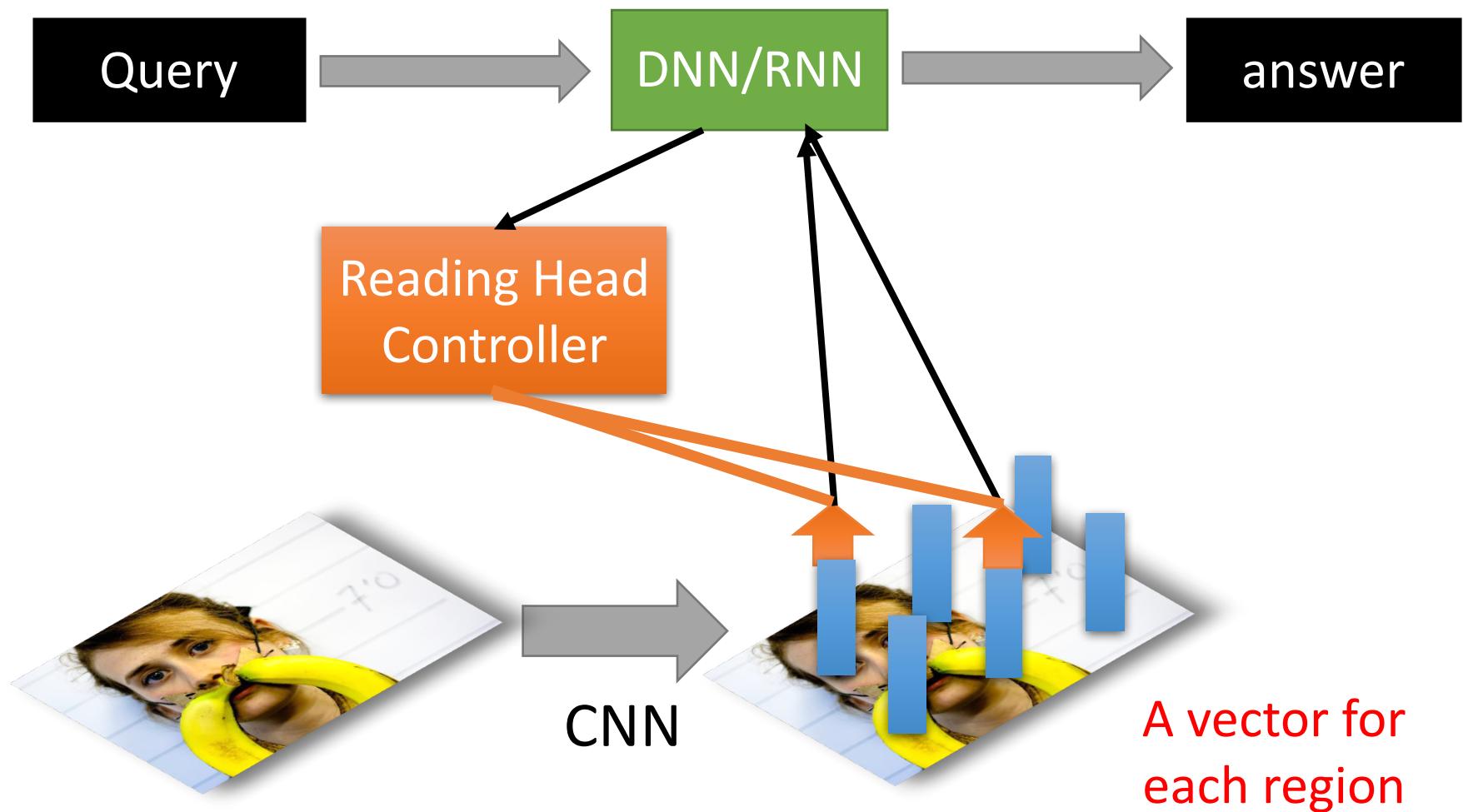


# Visual Question Answering



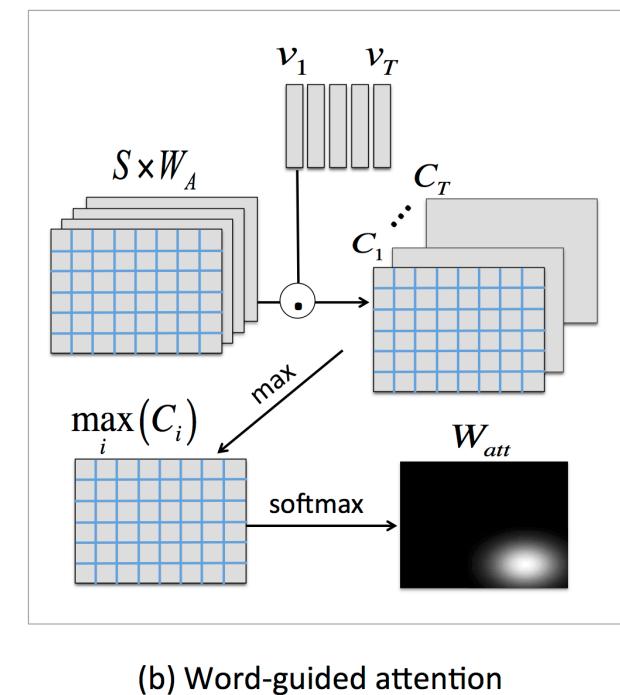
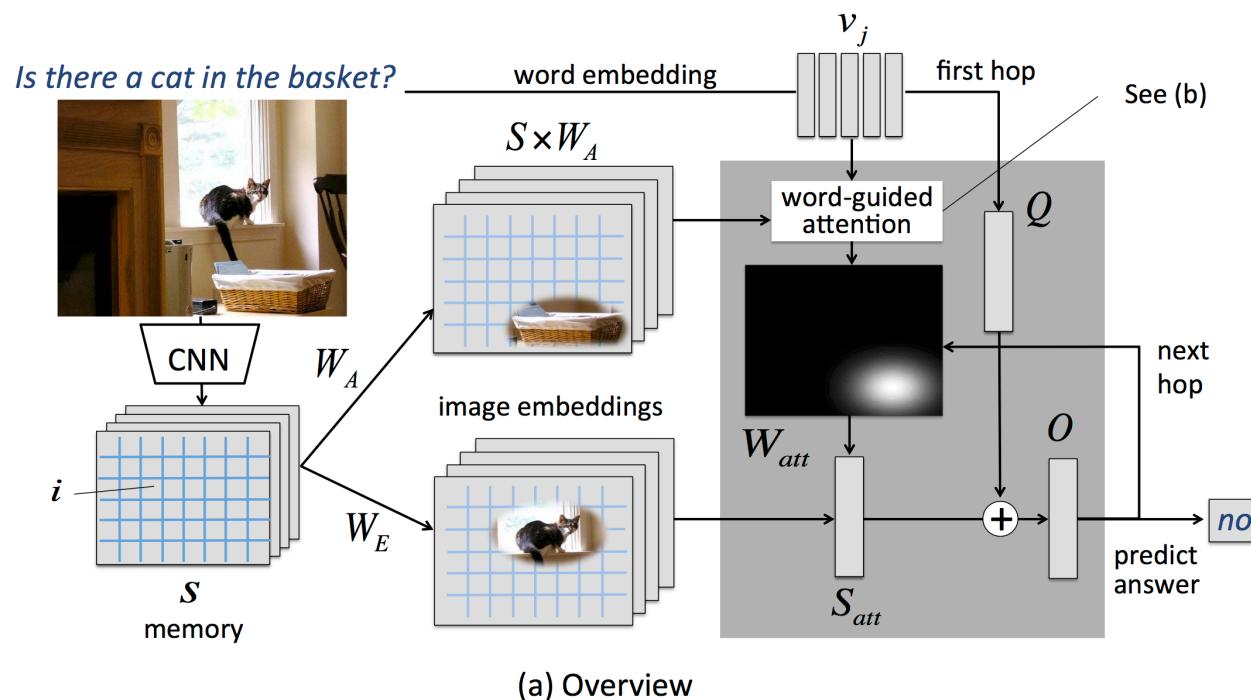
source: <http://visualqa.org/>

# Visual Question Answering



# Visual Question Answering

- Huijuan Xu, Kate Saenko. Ask, Attend and Answer: Exploring Question-Guided Spatial Attention for Visual Question Answering. arXiv Pre-Print, 2015



# Visual Question Answering

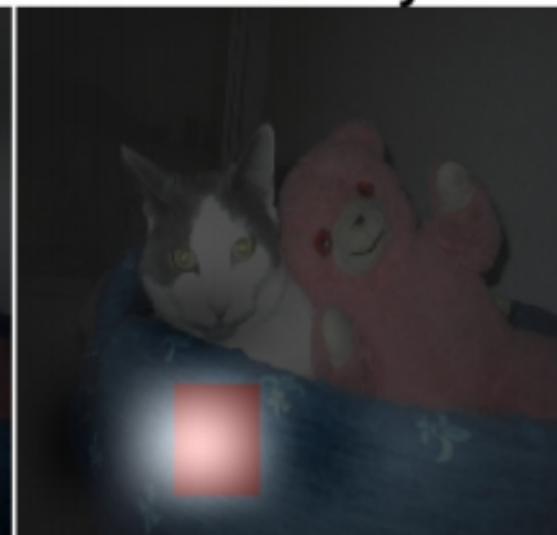
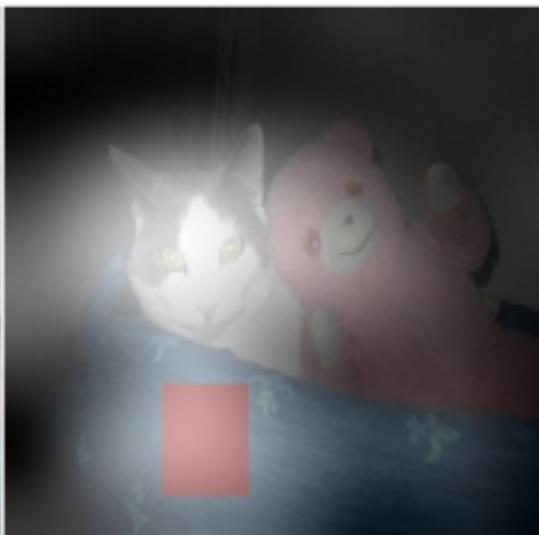
- Huijuan Xu, Kate Saenko. Ask, Attend and Answer: Exploring Question-Guided Spatial Attention for Visual Question Answering. ECCV, 2016.

**Is there a red square on the bottom of the cat?**

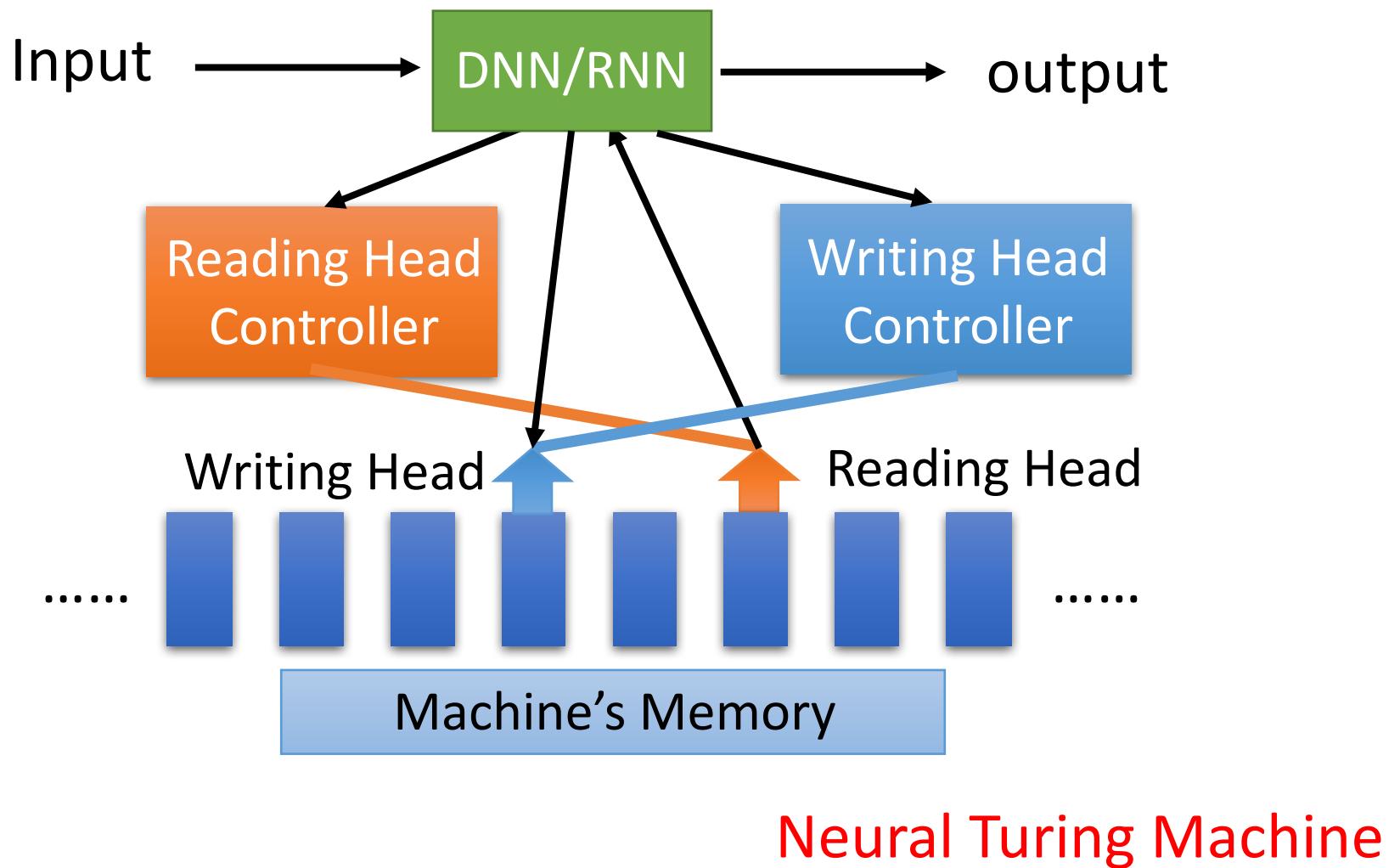
**GT: yes**



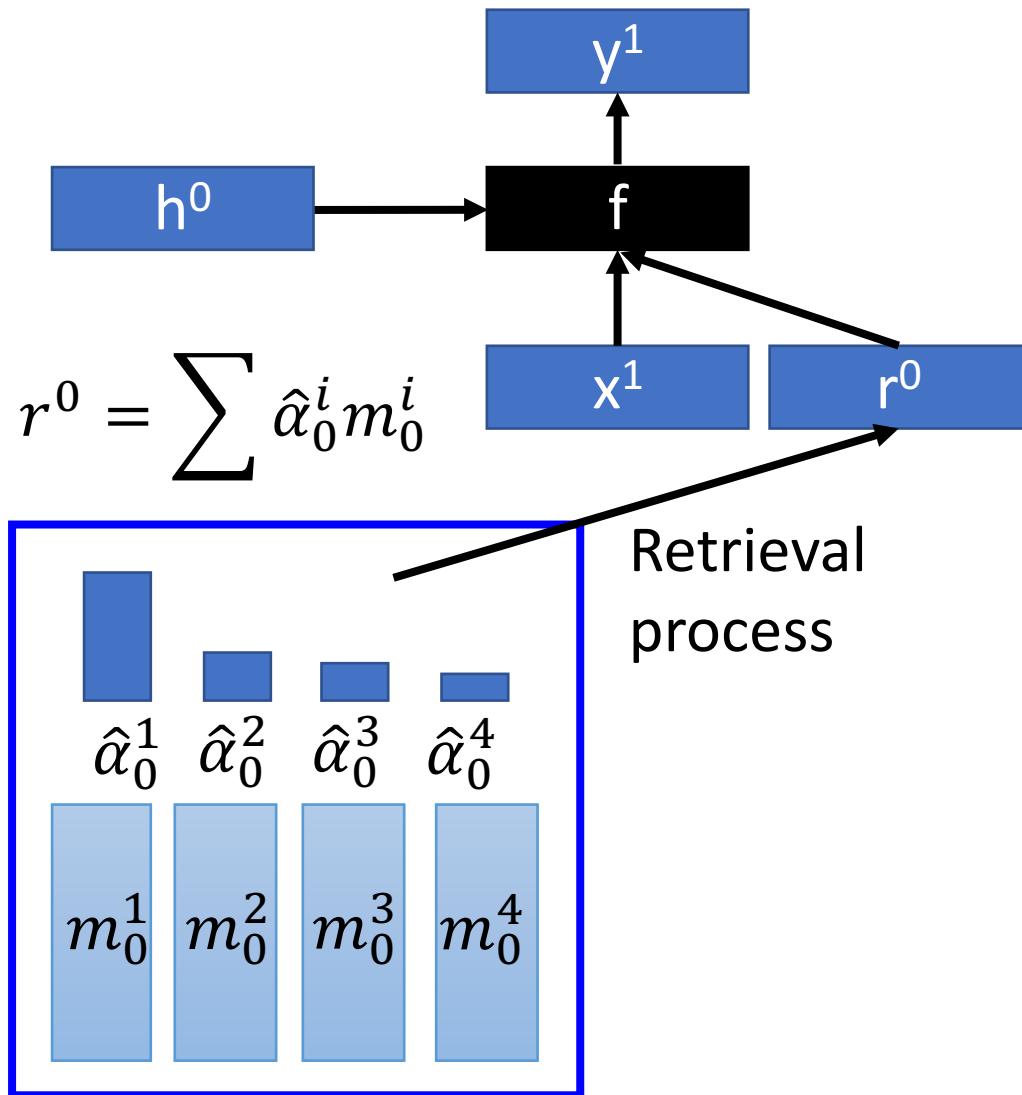
**Prediction: yes**



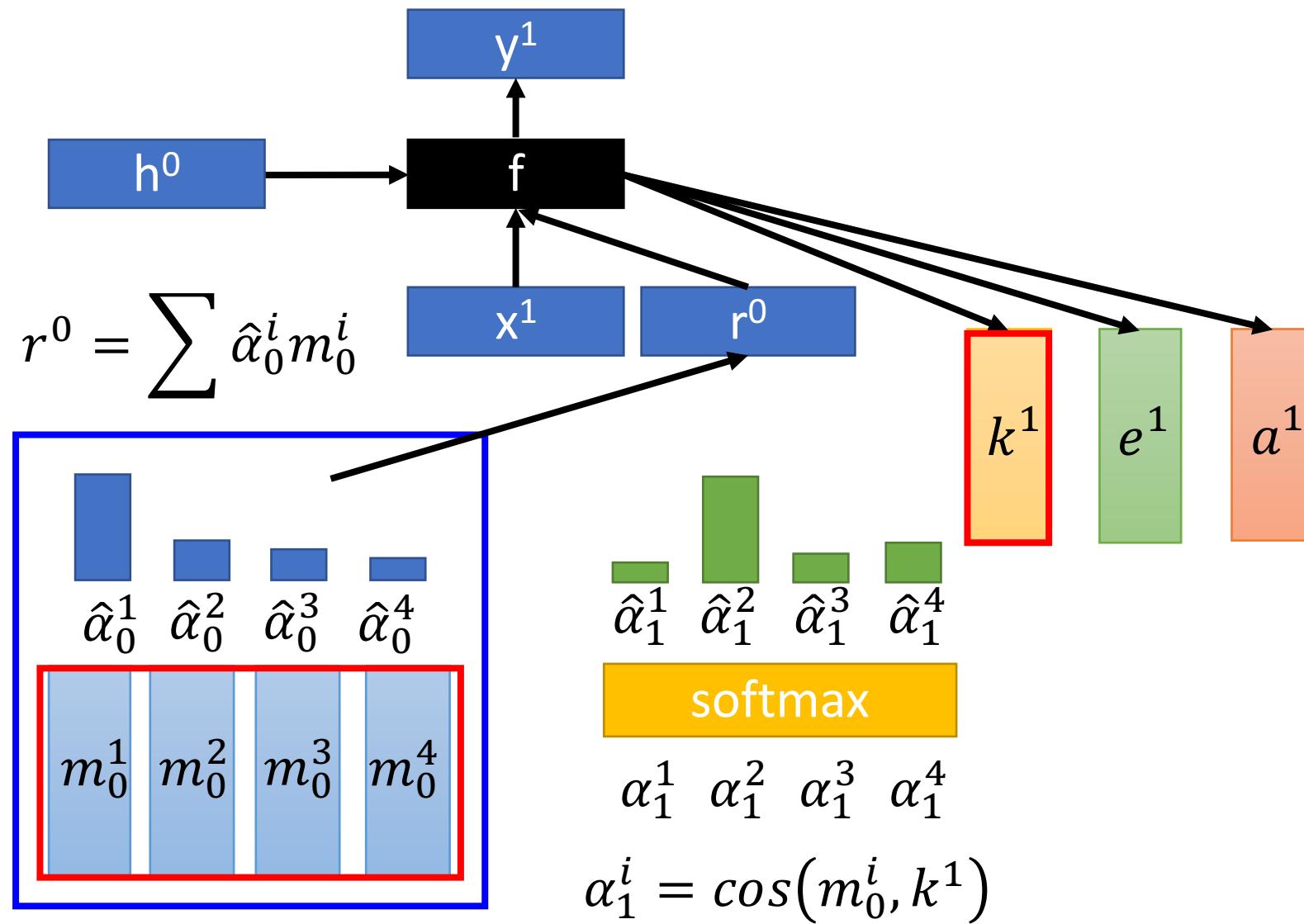
# External Memory v2



# Neural Turing Machine



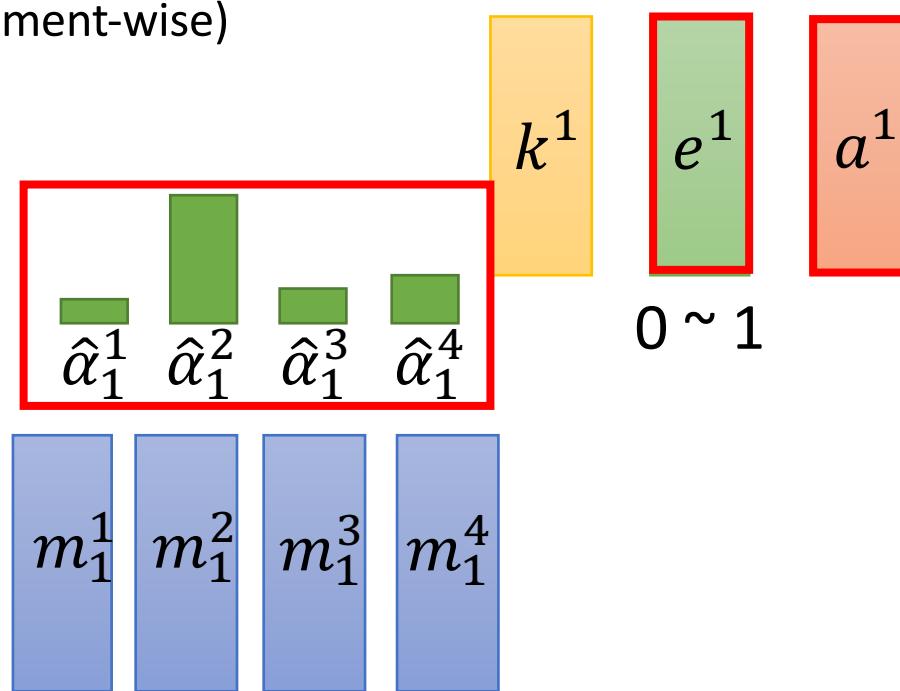
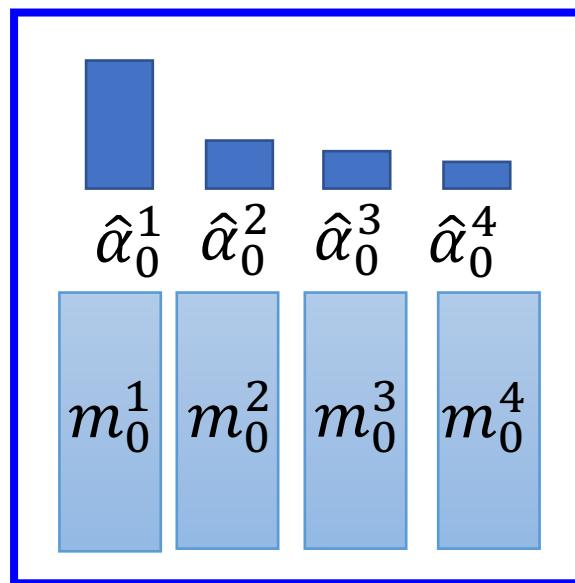
# Neural Turing Machine



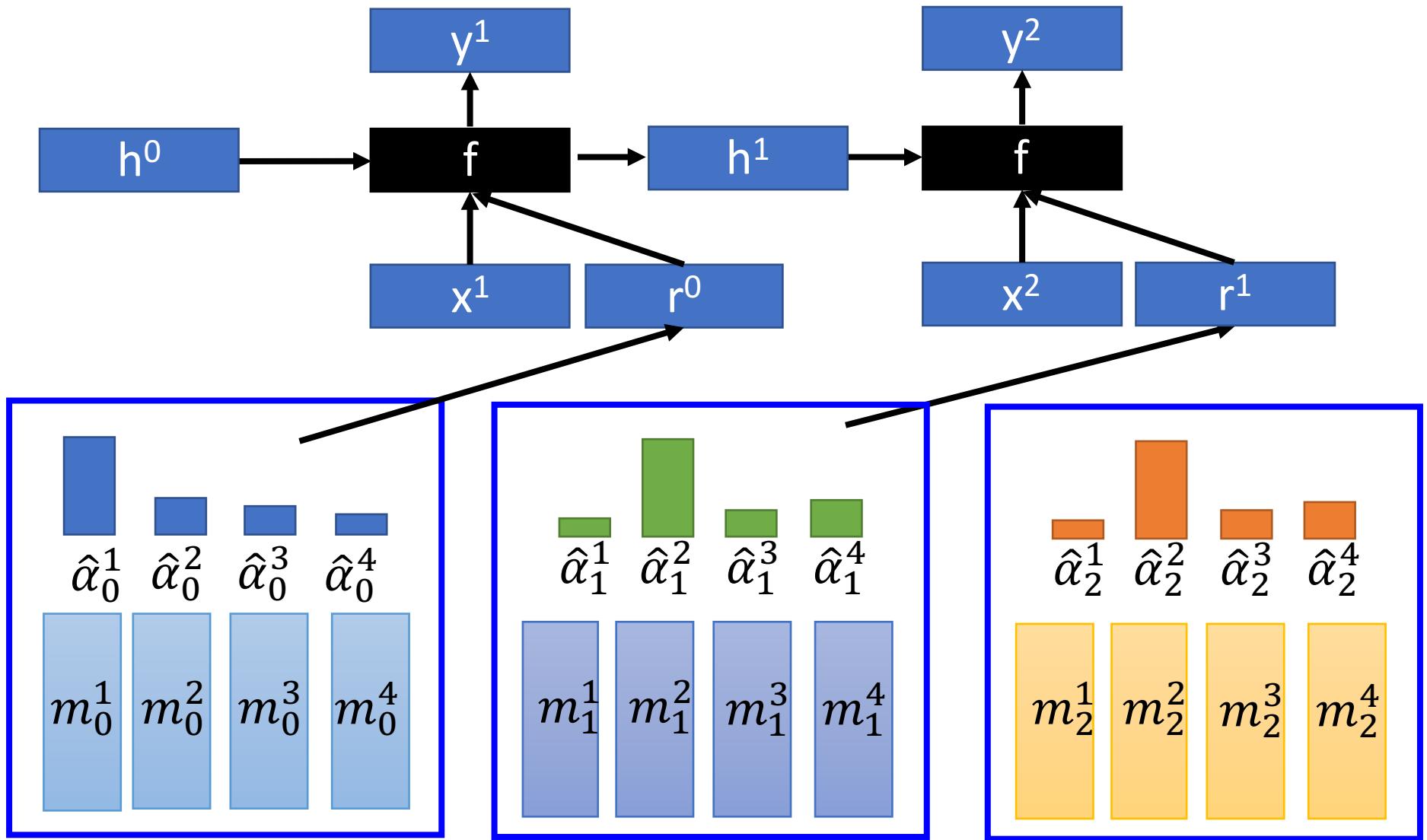
# Neural Turing Machine

$$m_1^i = m_0^i - \hat{\alpha}_1^i e^1 \odot m_0^i + \hat{\alpha}_1^i a^1$$

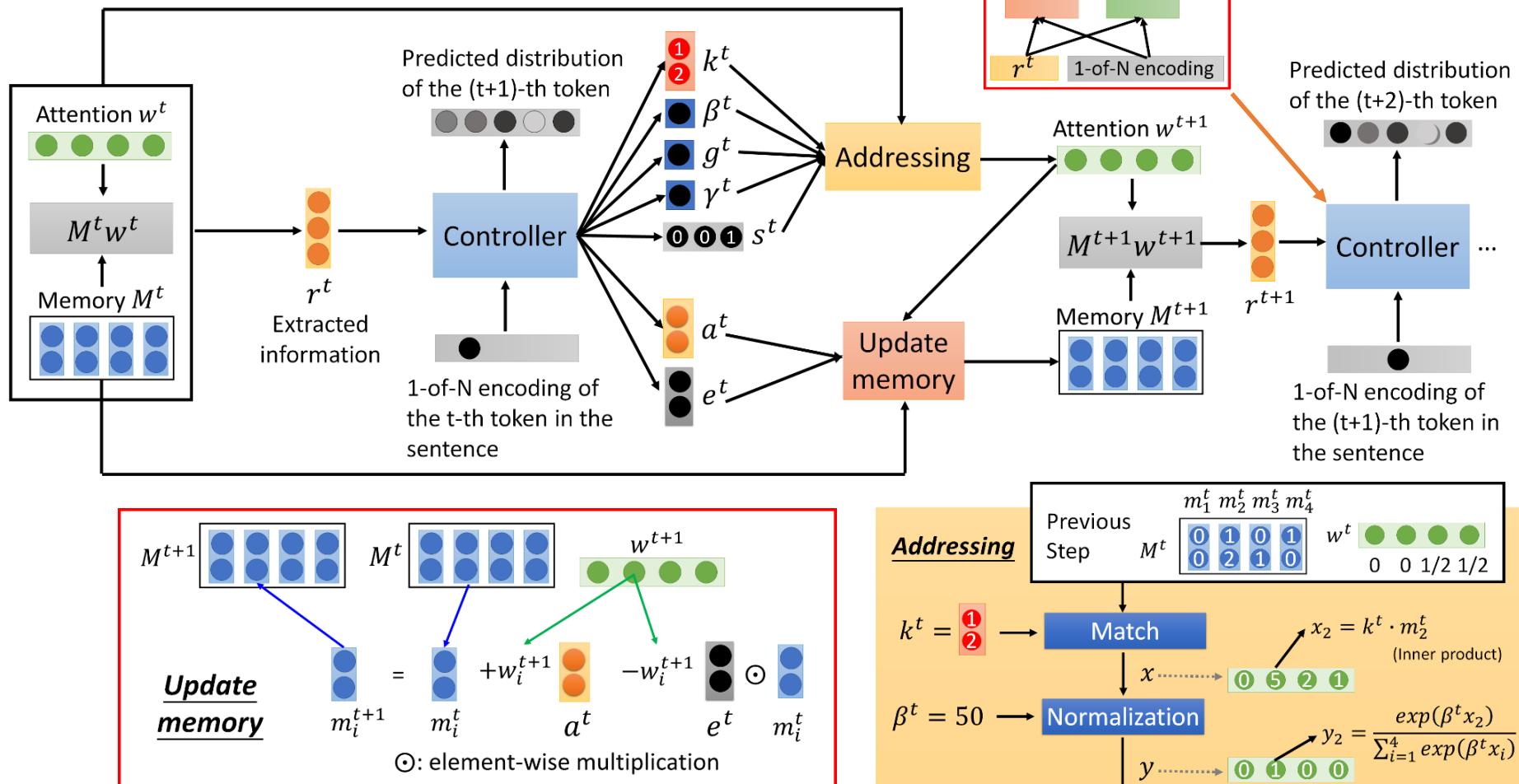
(element-wise)



# Neural Turing Machine

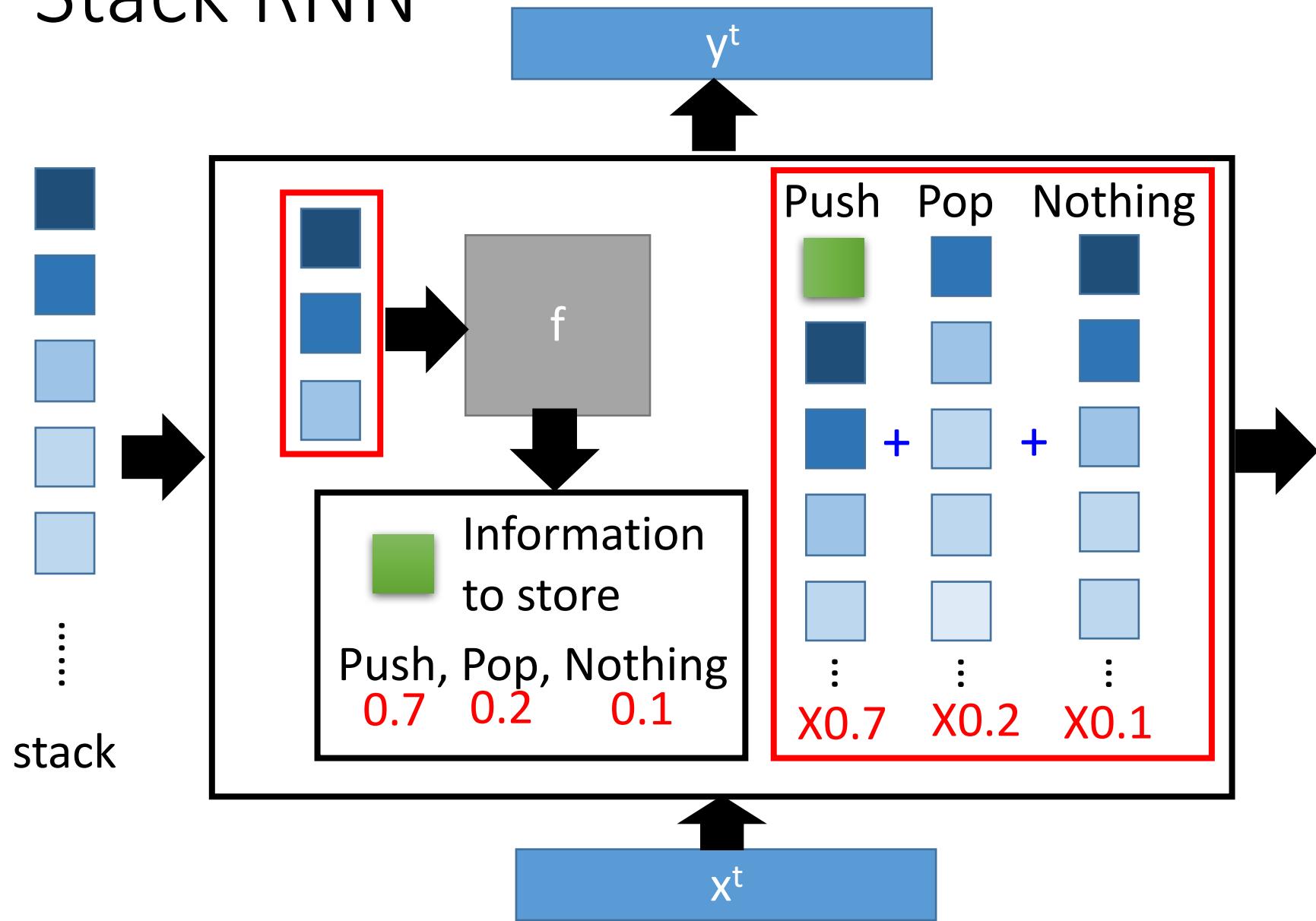


# Neural Turing Machine for LM



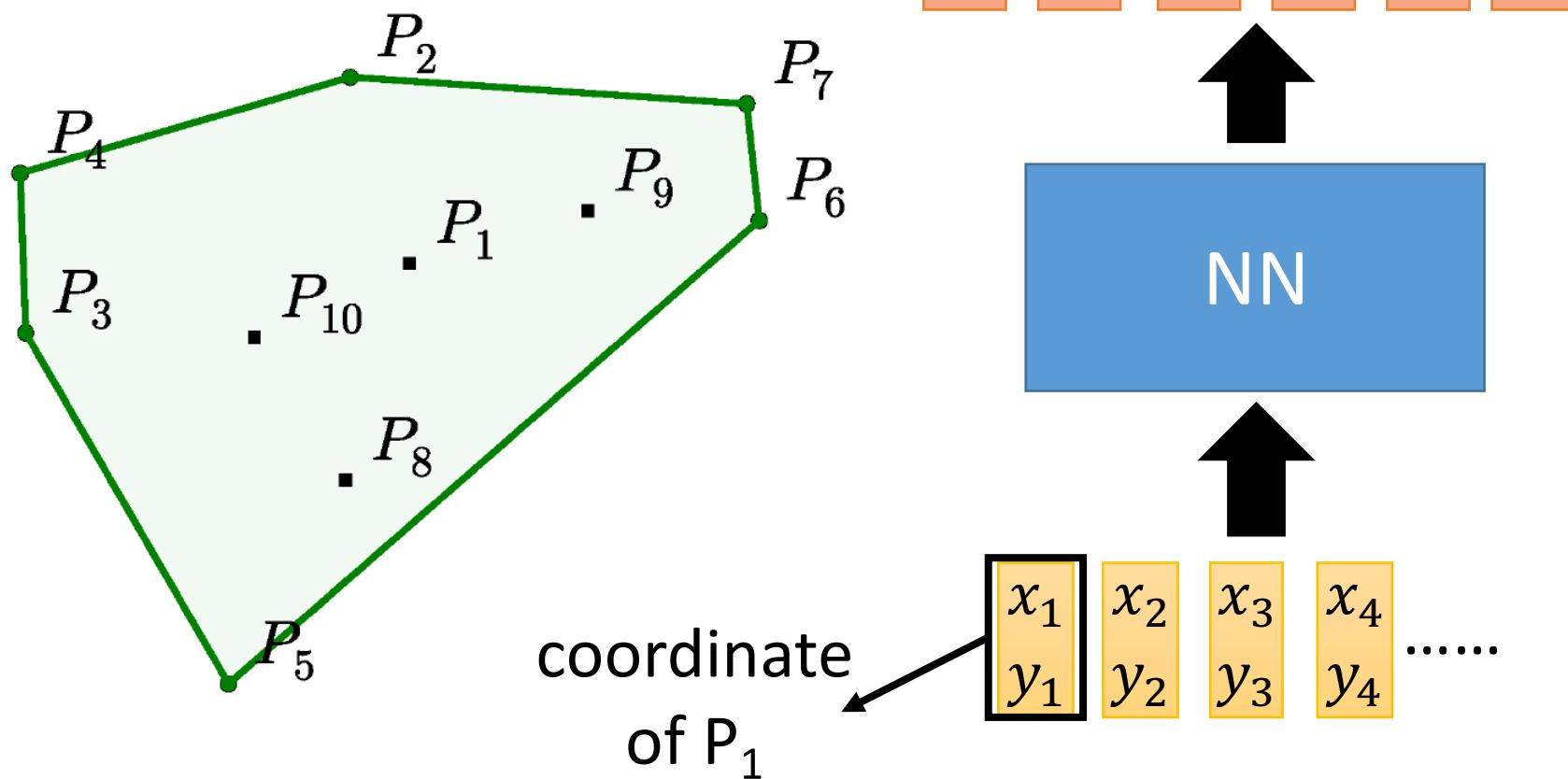
Wei-Jen Ko, Bo-Hsiang Tseng, Hung-yi Lee,  
 “Recurrent Neural Network based Language  
 Modeling with Controllable External Memory”,  
 ICASSP, 2017

# Stack RNN



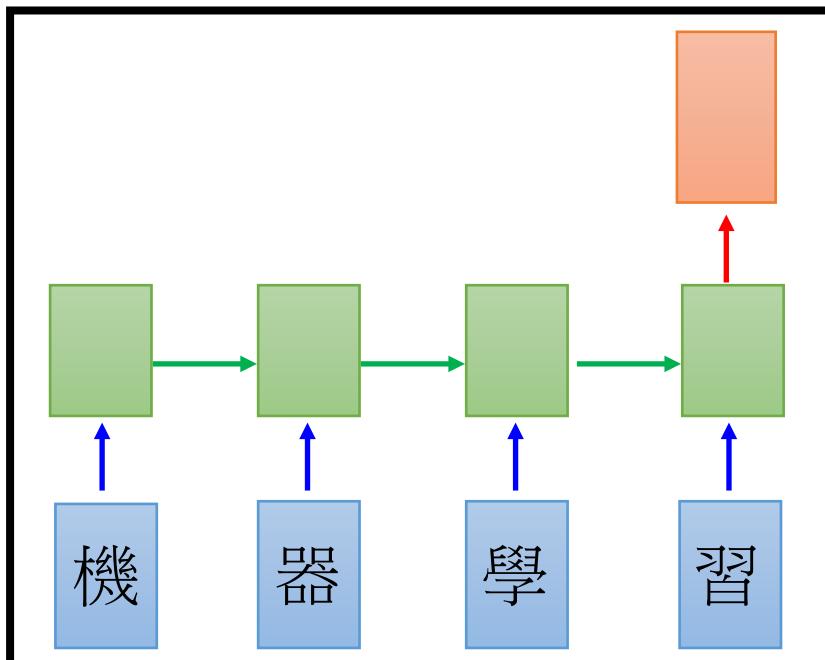
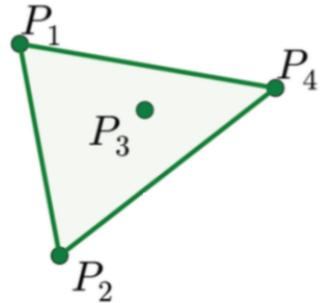
# Pointer Network

# Pointer Network

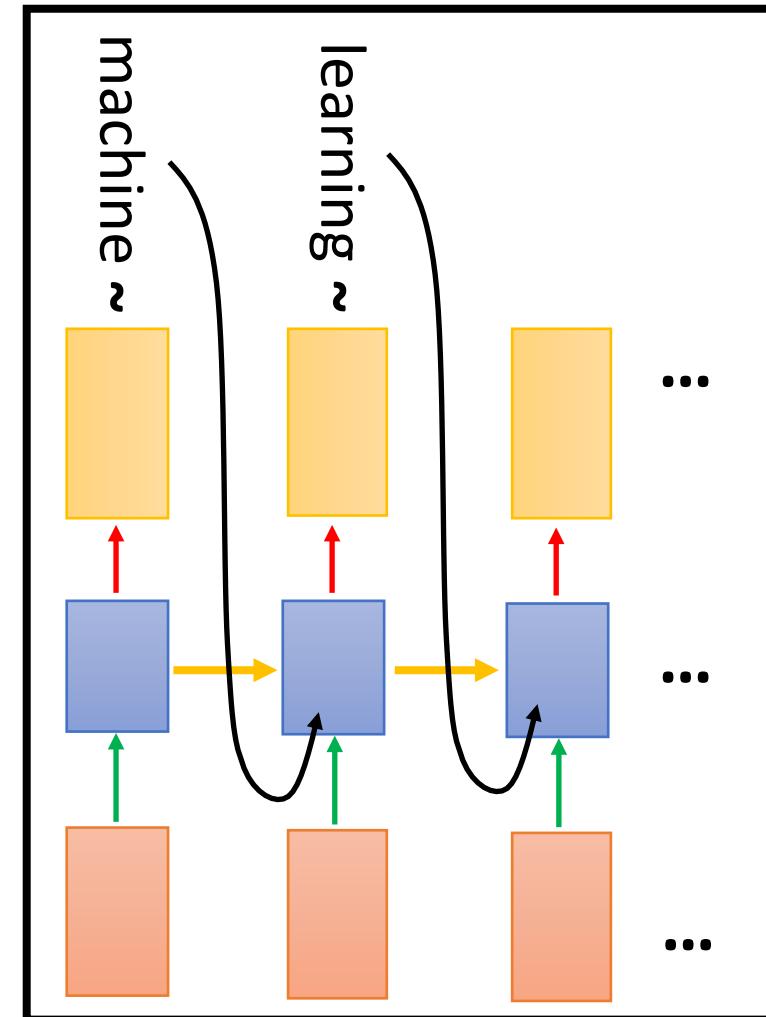


planar convex hulls, computing Delaunay triangulations, and the planar Travelling Salesman Problem

# Sequence-to-sequence?



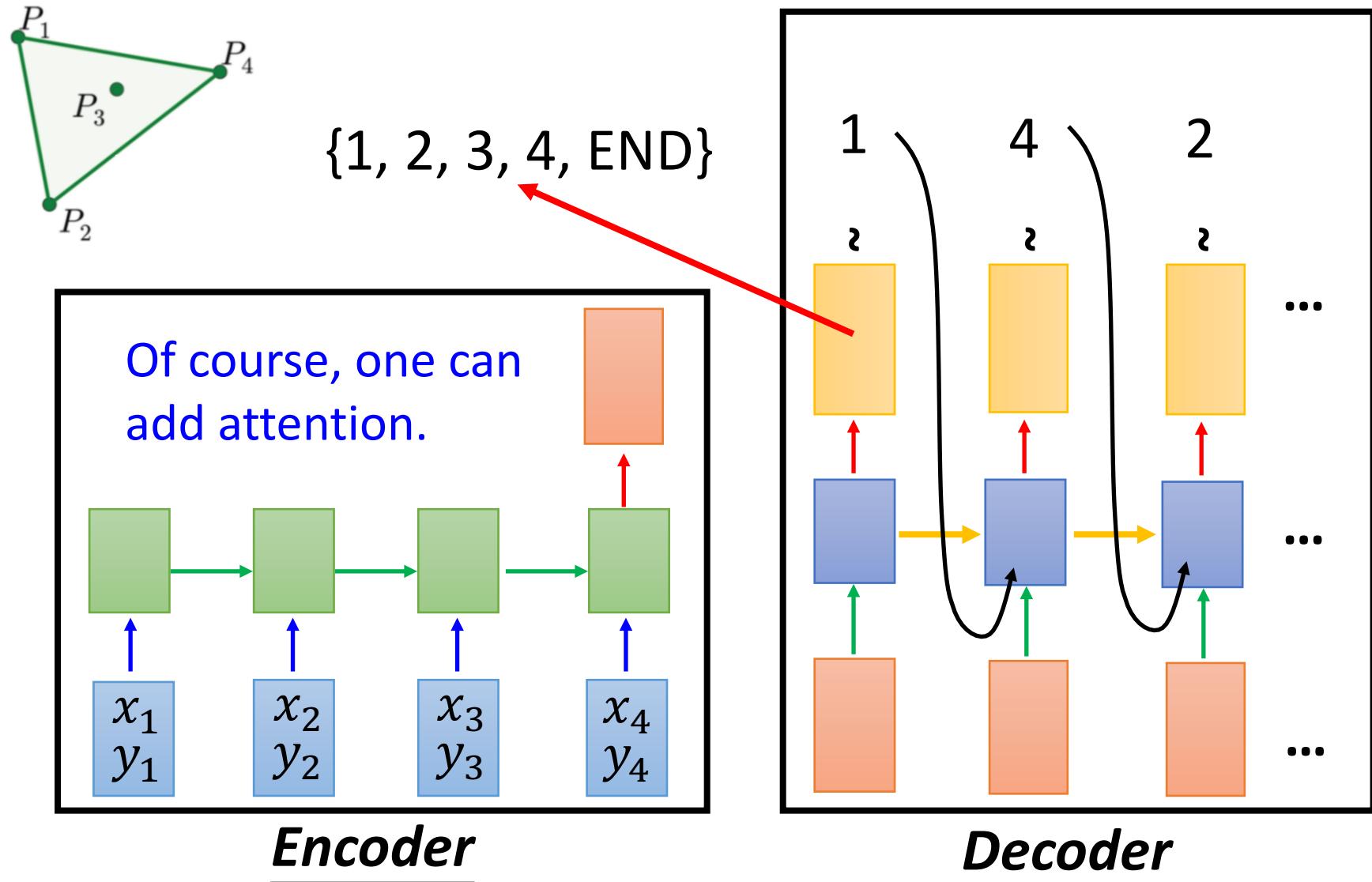
*Encoder*



*Decoder*

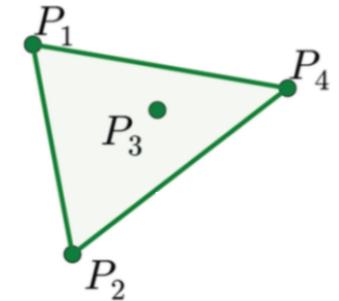
Problem?

# Sequence-to-sequence?

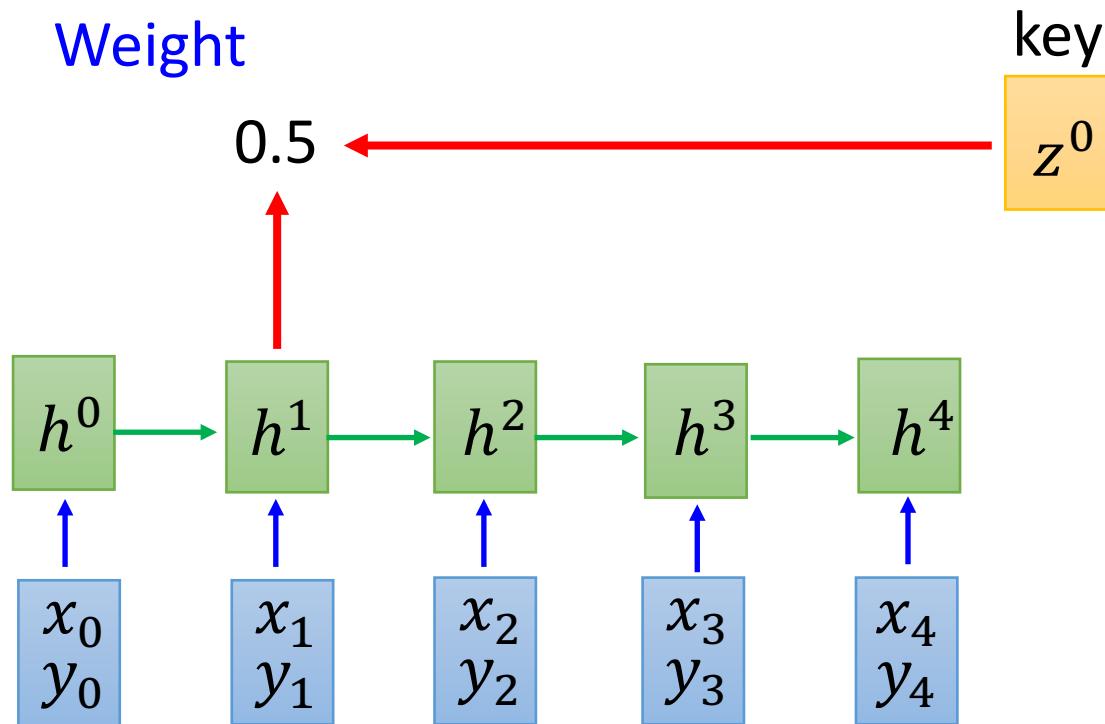


# Pointer Network

$x_0$   
 $y_0$  : END

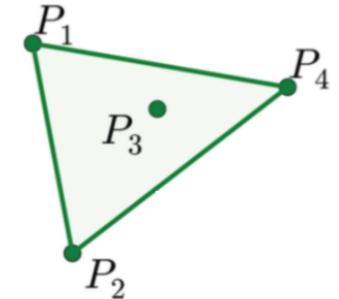


Attention  
Weight



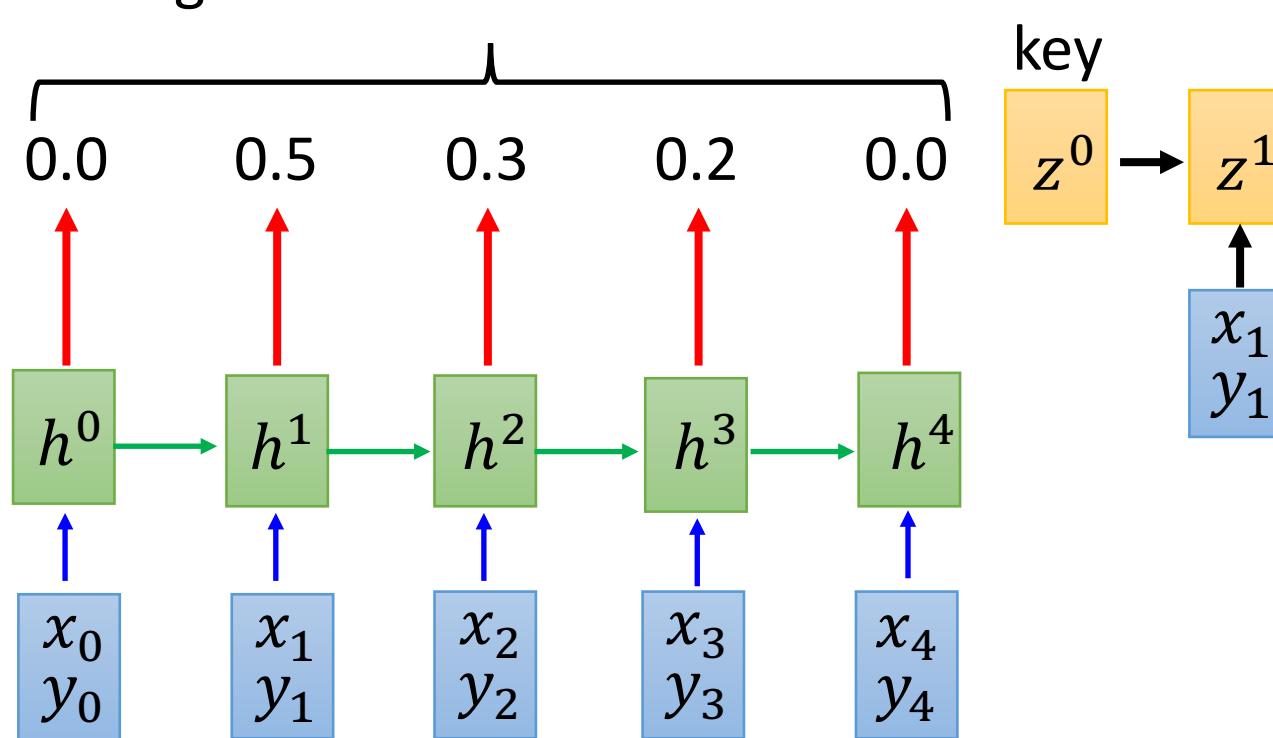
# Pointer Network

$x_0$   
 $y_0$  : END



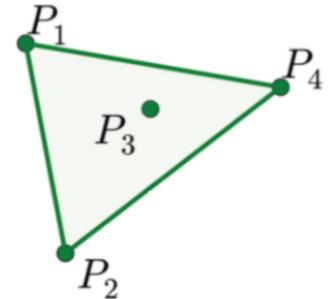
Output: 1  
? argmax from this distribution

What decoder can output depends on the input.



# Pointer Network

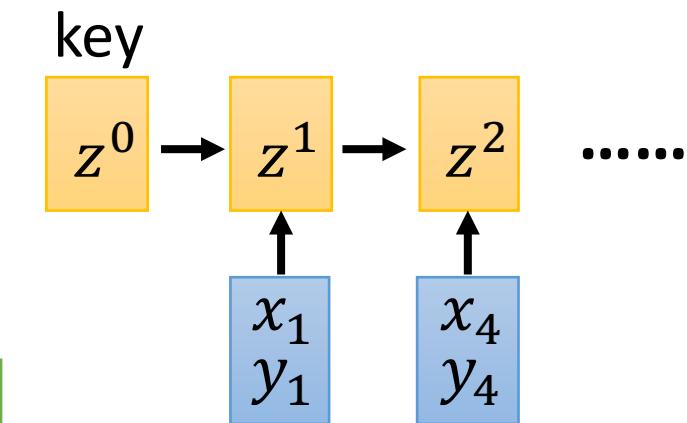
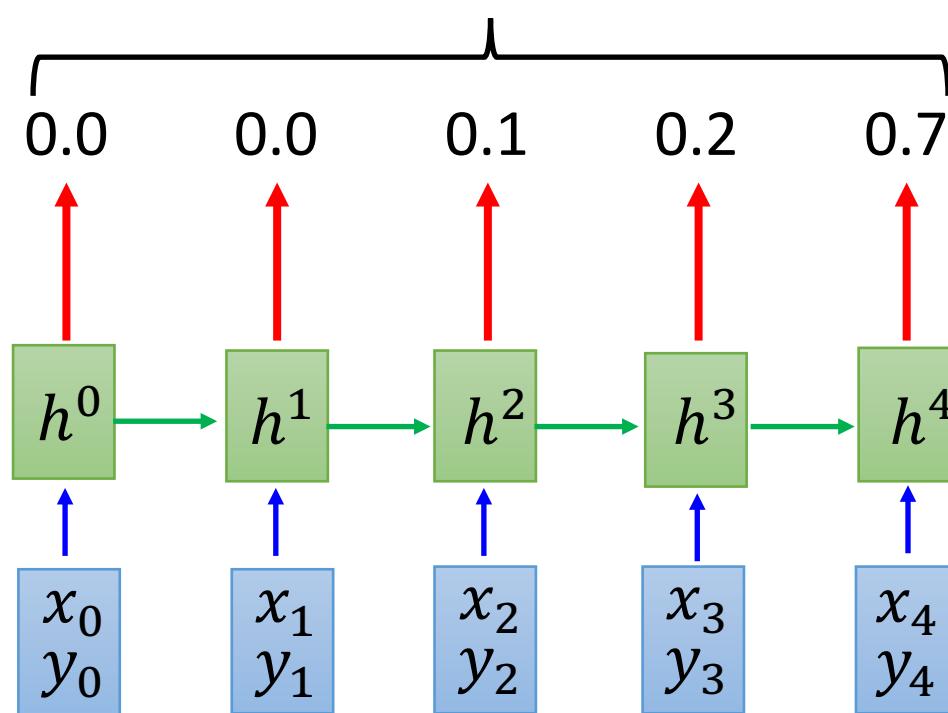
$x_0$   
 $y_0$  : END



Output: 4

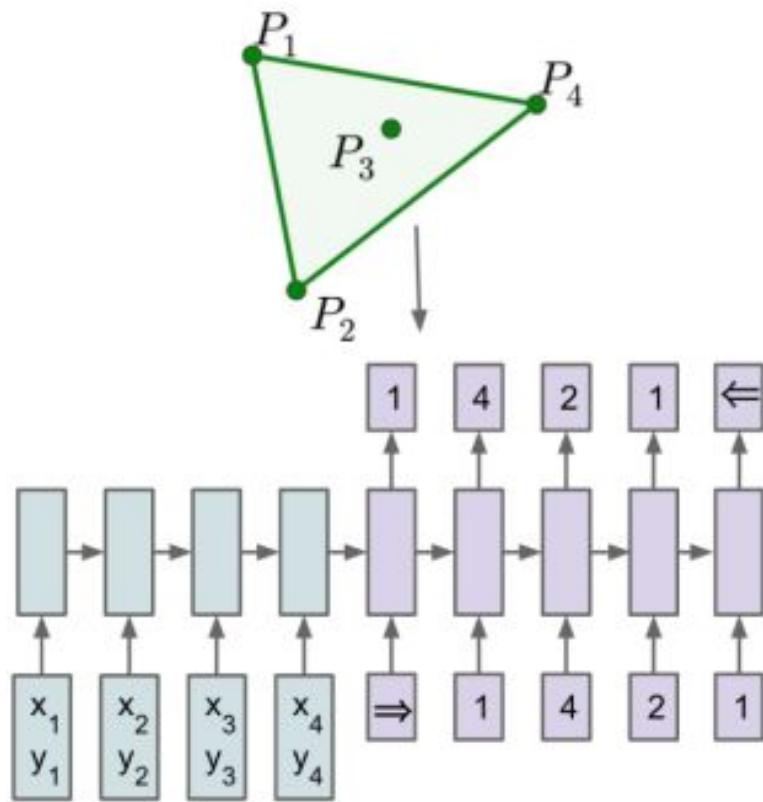
?

argmax from this distribution

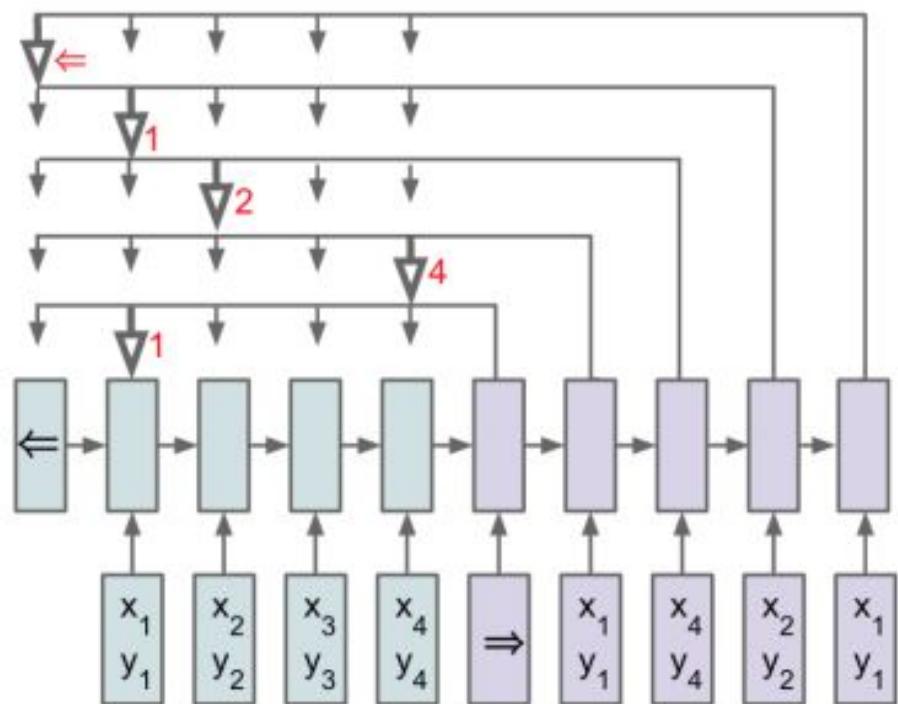


The process stops when  
“END” has the largest  
attention weights.

# Pointer Network

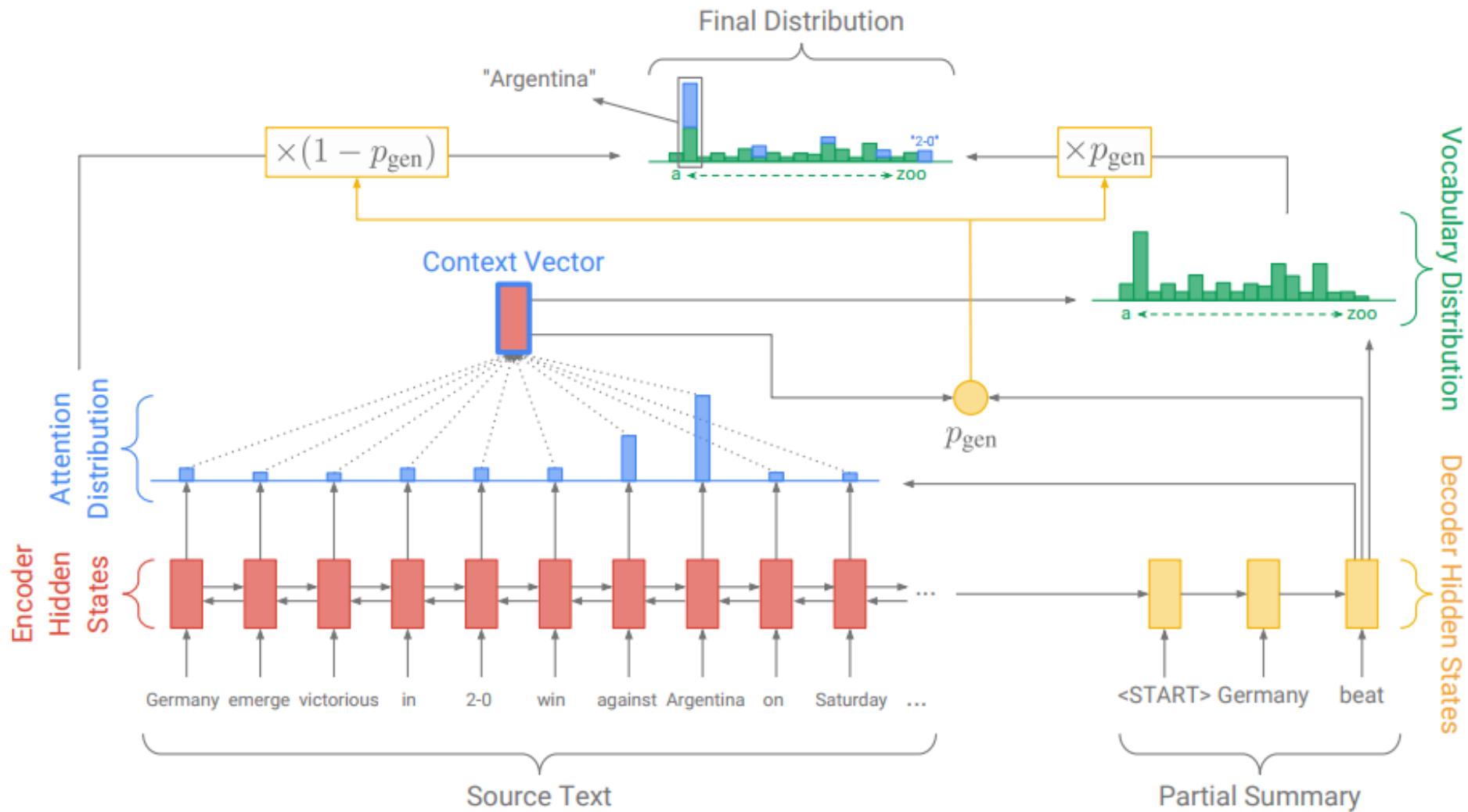


(a) Sequence-to-Sequence



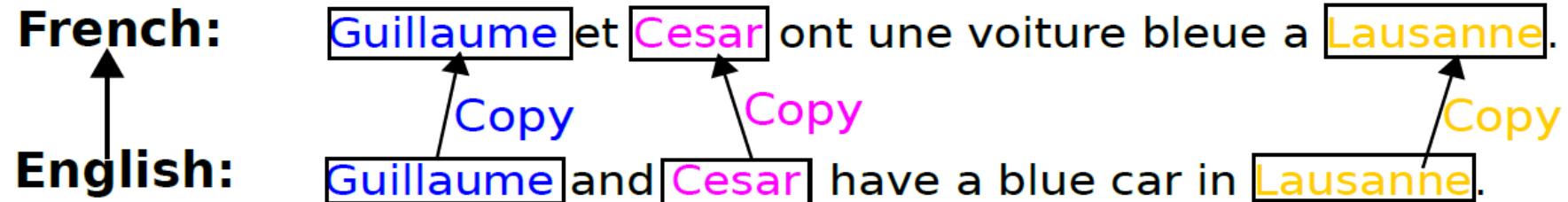
(b) Ptr-Net

# Applications - Summarization



# More Applications

## *Machine Translation*



## *Chat-bot*

User: X寶你好，我是庫洛洛

Machine: 庫洛洛你好，很高興認識你

# More Applications

- Article summarization - 《Neural Summarization by Extracting Sentences and Words》
- Information retrieval - 《End-to-End Information Extraction without Token-Level Supervision》
- Sentence ordering = 《End-to-End Neural Sentence Ordering Using Pointer Network》
- Auto programming for cards of board games  
《Latent predictor networks for code generation》