

---

# Deep Learning for Semantic Matching in Search

Jun Xu

[junxu@ict.ac.cn](mailto:junxu@ict.ac.cn)

ICT, CAS

---

# Tutorial Talk @ ADL 52 & NLPCC 2014



ADL & NLPCC 2014 Tutorial  
December 6, 2014  
ShenZhen China

## Semantic Matching in Search

Jun Xu

[junxu@ict.ac.cn](mailto:junxu@ict.ac.cn)

Institute of Computing Technology  
Chinese Academy of Sciences

# Discussed Traditional Approaches to Semantic Matching in Search

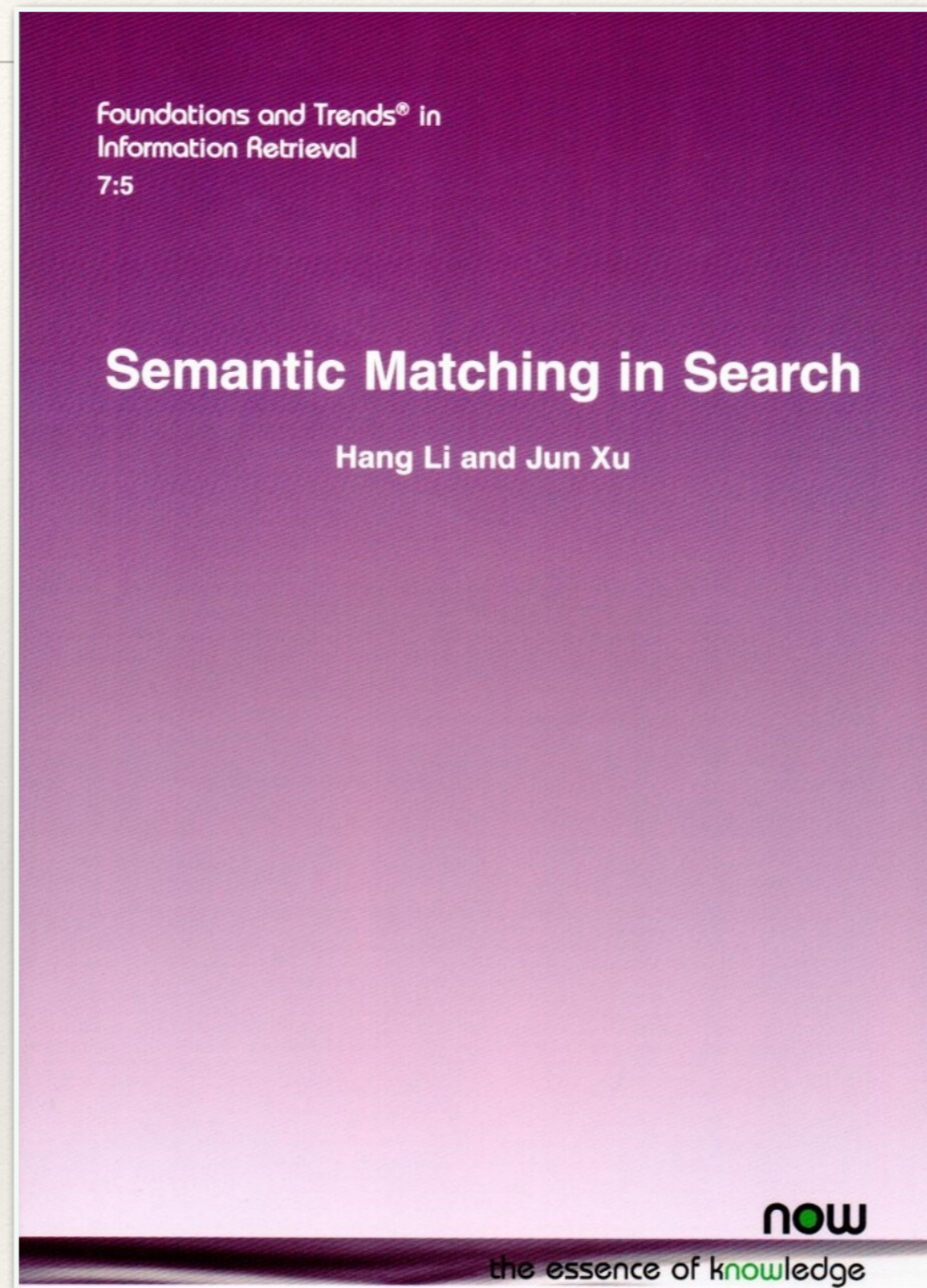
## Outline of Tutorial

- Semantic Matching between Query and Document
- Approaches to Semantic Matching
  1. Matching by Query Reformulation
  2. Matching with Term Dependency Model
  3. Matching with Translation Model
  4. Matching with Topic Model
  5. Matching with Latent Space Model
- Summary

---

# Details Introduced in a Monograph

---



<http://www.nowpublishers.com/articles/foundations-and-trends-in-information-retrieval/INR-035>  
[http://www.hangli-hl.com/uploads/3/1/6/8/3168008/ml\\_for\\_match-step2.pdf](http://www.hangli-hl.com/uploads/3/1/6/8/3168008/ml_for_match-step2.pdf)

---

# Growing Interest in “Deep” IR in the Past Three Years

---

- ❖ Success of deep learning in other fields
  - ❖ Speech recognition, computer vision, and NLP
- ❖ Growing presence of deep learning in IR research
  - ❖ SIGIR 2016 keynote, Tutorial, and Neu-IR workshop
- ❖ Adopted by industry
  - ❖ ACM News: Google Turning its Lucrative Web Search Over to AI Machines (Oct. 26, 2015)
  - ❖ WIRED: AI is Transforming Google Search. The Rest of the Web is Next (April 2, 2016)
- ❖ Chris Manning (Stanford)’s SIGIR keynote:  
“I’m certain that *deep learning will come to dominate SIGIR over the next couple of years* ... just like speech, vision, and NLP before it.”



---

# “Deep” Semantic Matching also Gain a Lot of Attention

---

- ❖ Before 2014, a few studies, e.g.,
  - ❖ Paraphrase detection [Socher et al., 2011]
  - ❖ Ad-hoc retrieval (DSSM)[Huang et al., 2013]
- ❖ 2014 ~ 2017, a lot of studies (as summarized in the tutorial)
  - ❖ Paraphrase identification
  - ❖ Ad-hoc retrieval
  - ❖ Question answering
  - ❖ Dialog
  - ❖ Result diversification
  - .....

This tutorial:  
Update the survey with newly  
developed deep matching methods

---

# Outline

---

- ❖ Semantic matching in search
- ❖ Word-level matching: bridging the semantic gap
- ❖ Sentence-level matching: capturing the proximity
- ❖ Summary and discussion



---

# A Good Web Search Engine

---

- ❖ Must be good at
  - ❖ Relevance
  - ❖ Coverage
  - ❖ Diversity
  - ❖ Freshness
  - ❖ Response time
  - ❖ User interface.....
- ❖ Relevance is particularly important



# Query-Document Mismatch Challenge

**Table 1.1:** Examples of query document mismatch.

query	document	term match	semantic match
seattle best hotel	seattle best hotels	partial	yes
pool schedule	swimming pool schedule	partial	yes
natural logarithm transform	logarithm transform	partial	yes
china kong	china hong kong	partial	no
why are windows so expensive	why are macs so expensive	partial	no

---

# Why Query-Document Mismatch Happens?

---

- ❖ Search is still mainly based on term-level matching signals
- ❖ Some search intent can be represented by different queries (representations)
- ❖ Query document mismatch occurs, when searcher and author use different terms (representations) to describe the same concept

# Same Search Intent, Different Query Representations

**Table 1.2:** Queries about “distance between sun and earth”.

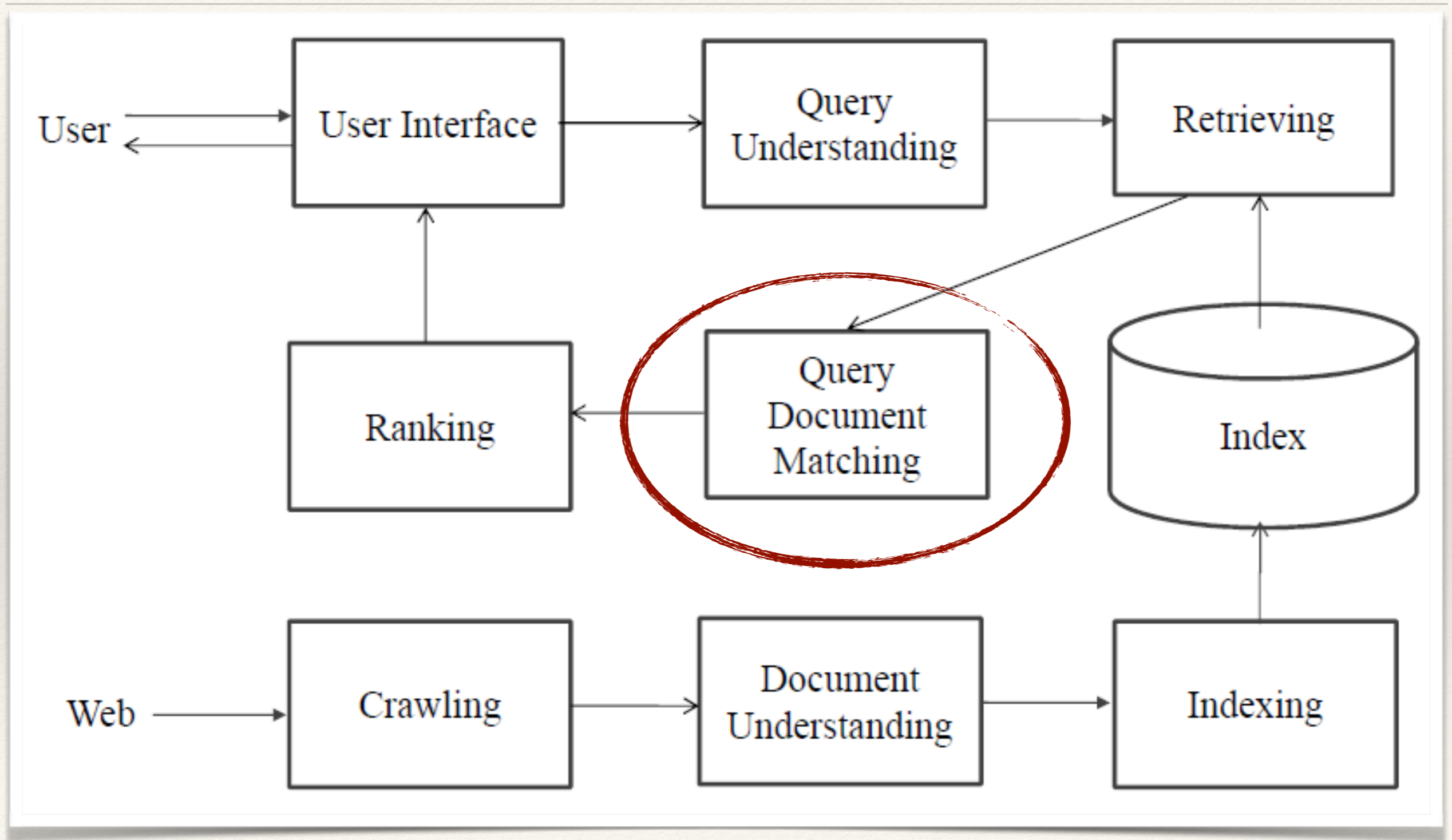
“how far” earth sun	average distance from the earth to the sun
“how far” sun	how far away is the sun from earth
average distance earth sun	average distance from earth to sun
how far from earth to sun	distance from earth to the sun
distance from sun to earth	distance between earth and the sun
distance between earth & sun	distance between earth and sun
how far earth is from the sun	distance from the earth to the sun
distance between earth sun	distance from the sun to the earth
distance of earth from sun	distance from the sun to earth
“how far” sun earth	how far away is the sun from the earth
how far earth from sun	distance between sun and earth
how far from earth is the sun	how far from the earth to the sun
distance from sun to the earth	

# Same Search Intent, Different Query Representations

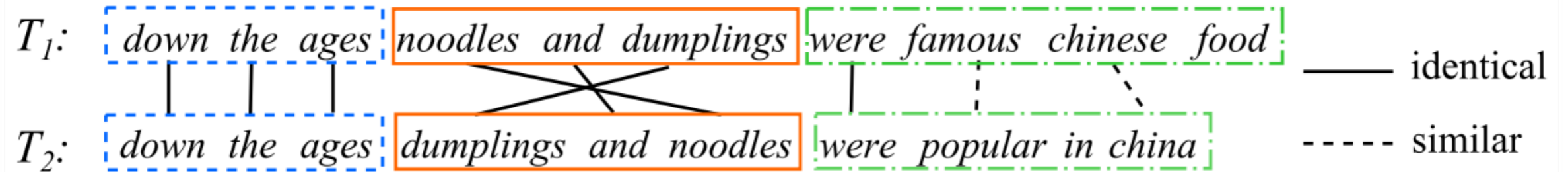
**Table 1.3:** Queries about “Youtube”.

yutube	yuotube	yuo tube
ytube	youtubr	yu tube
youtubo	youtuber	youtubecom
youtube om	youtube music videos	youtube videos
youtube	youtube com	youtube co
youtub com	you tube music videos	yout tube
youtub	you tube com yourtube	your tube
you tube	you tub	you tube video clips
you tube videos	www you tube com	www youtube com
www youtube	www youtube com	www youtube co
yotube	www you tube	www utube com
ww youtube com	www utube	www u tube
utube videos	utube com	utube
u tube com	utub	u tube videos
u tube	my tube	toutube
outube	our tube	toutube

# Semantic Matching in Search



# Challenges



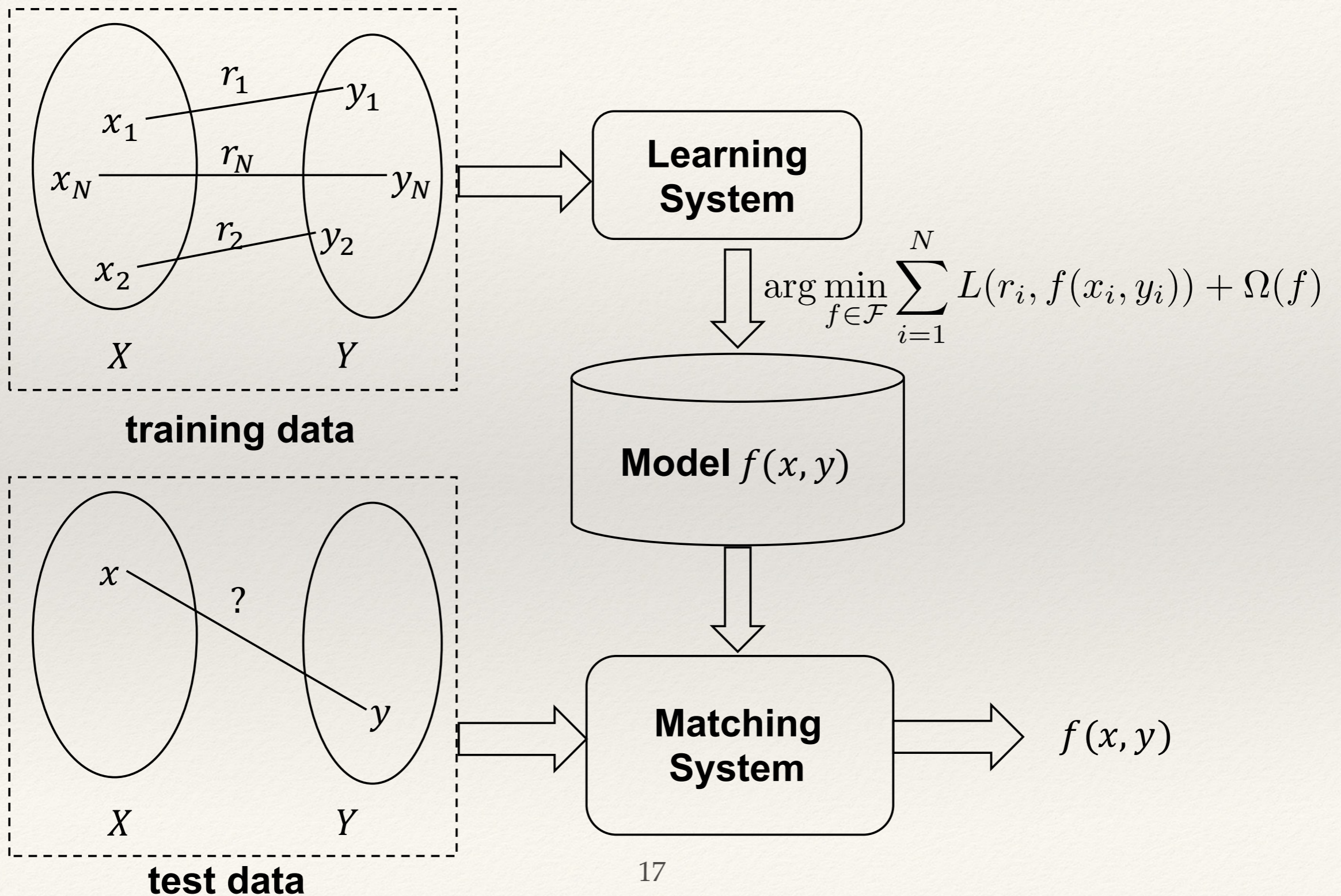
- ❖ Word-level matching: semantic gap between words
  - ❖ Two words has similar meanings
  - ❖ “popular” ~ “famous”; “china” ~ “chinese”
- ❖ Sentence-level: proximity matching between sentences
  - ❖ The matching positions do matter
  - ❖ “noodles and dumplings” ~ “dumplings and noodles”
- ❖ Need to consider them simultaneously

Ideally: Understanding the Natural Language

Current Approaches: Avoid Understanding and  
Conduct Matching



# Learning to Match



---

# Why Deep?

---

- ❖ Representation

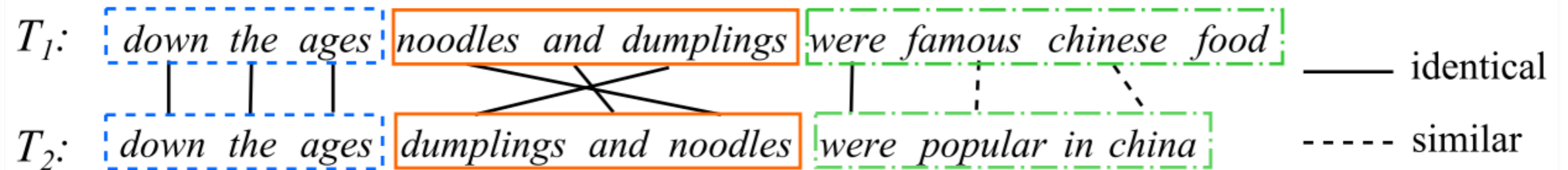
- ❖ Word: one hot distributed
- ❖ Sentence: bag-of-words  $\rightarrow$  distributed representation
- ❖ Better representation ability, better generalization ability

- ❖ Matching function

- ❖ Inputs (features): handcrafted  $\rightarrow$  automatically learned
- ❖ Function: simple functions (e.g., cosine, dot product)  $\rightarrow$  nonlinear neural networks
- ❖ Involving richer matching signals
- ❖ Considering soft matching patterns

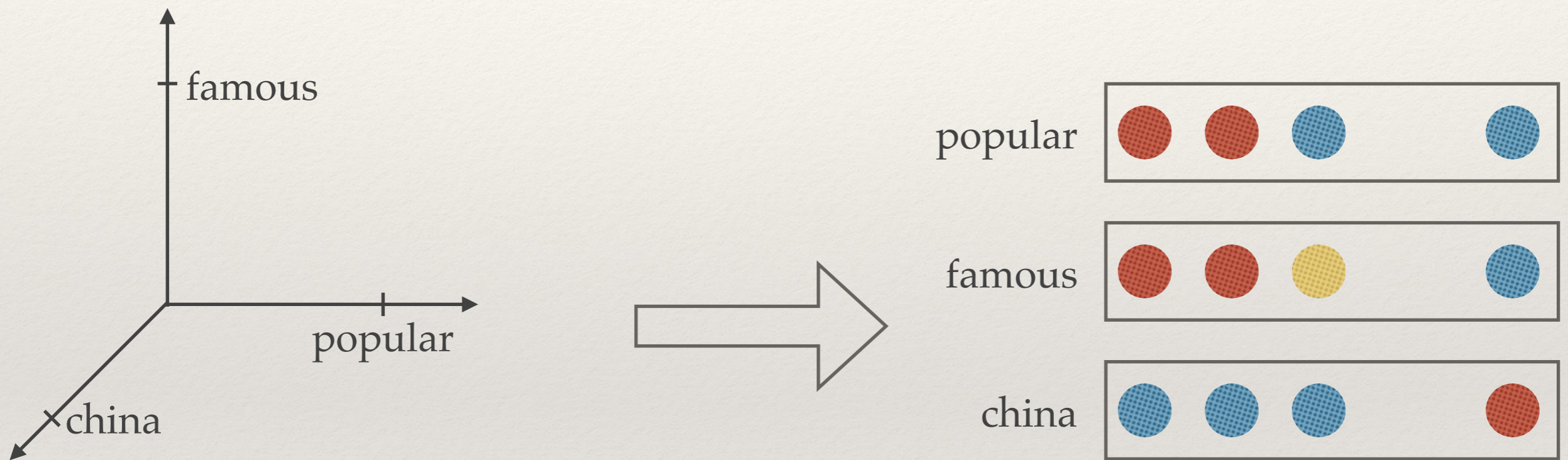
# Outline

- ❖ Semantic matching in search
- ❖ **Word-level matching: bridging the semantic gap**
- ❖ Sentence-level matching: capturing the proximity
- ❖ Summary and discussion



Measure the similarity between “famous” and “popular”

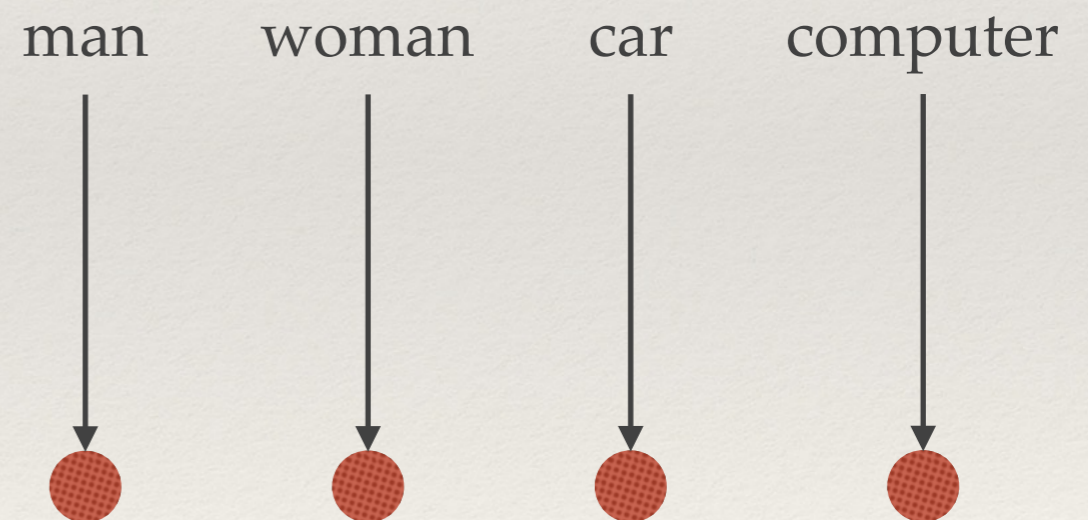
# From Local to Distributed Representations



# Local Representation of Words

- ❖ Words are building blocks of queries / documents
- ❖ Conventional IR models considers words as atomic symbols, also known as “one-hot” or local representations

Local(One-Hot) Representation	
man	[1,0,...,0,0,...,0,0]
woman	[0,1,...,0,0,...,0,0]
car	[0,0,...,1,0,...,0,0]
automobile	[0,0,...,0,0,...,1,0]



Each word is locally represented by a distinct node

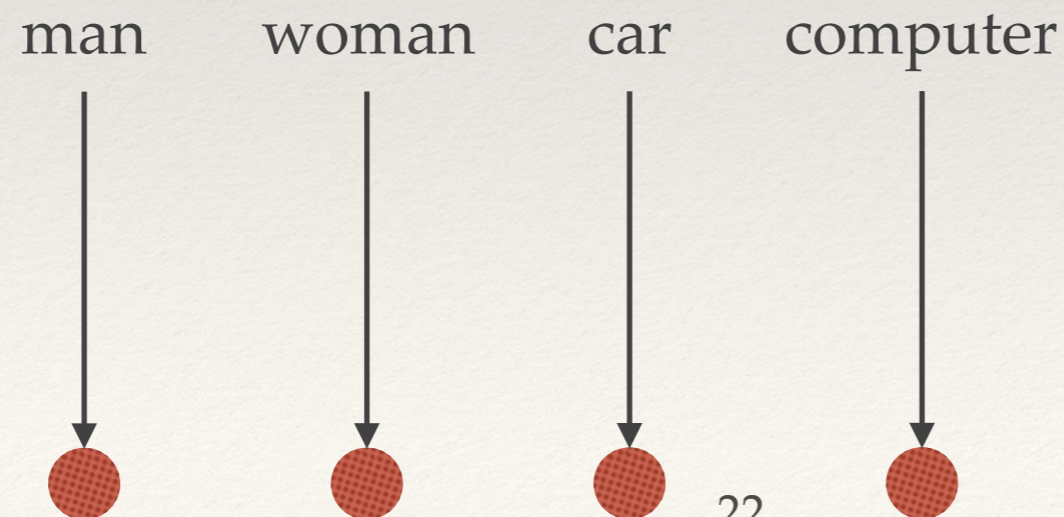
# Limitation of Local Representations

## ❖ Independent assumption

Local (one-hot) representations	
man	[1,0,...,0,0,...,0,0]
woman	[0,1,...,0,0,...,0,0]
car	[0,0,...,1,0,...,0,0]
automobile	[0,0,...,0,0,...,1,0]

$$\cos(\text{man}, \text{woman}) = 0$$
$$\cos(\text{man}, \text{automobile}) = 0$$

## ❖ Inefficient: N dimensions for N words



# Limitation of Local Representations (cont')

- ❖ Poor generalization ability
- ❖ Using language modeling as an example
  - ❖ Cannot generalize to unseen bigram “three groups”

Doc1: There are **three teams** left for the qualification

Doc2: **Four teams** have passed the first round

Doc3: **Four groups** are playing in the field

$$P(\text{teams}|\text{three}) > 0$$

$$P(\text{teams}|\text{four}) > 0$$

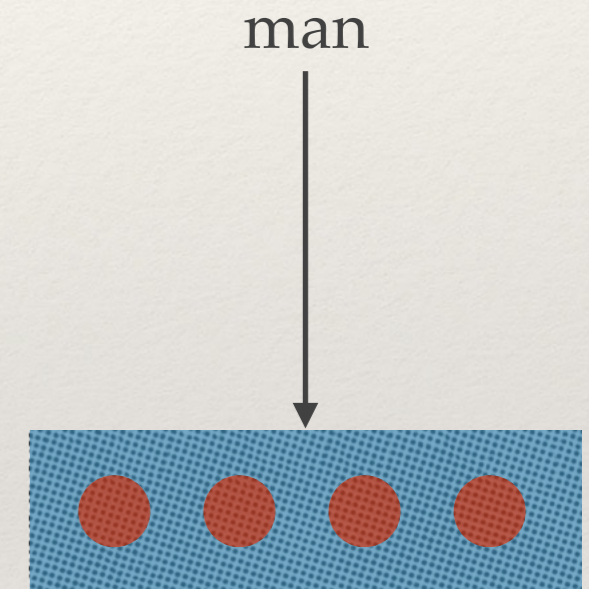
$$P(\text{groups}|\text{four}) > 0$$

$$P(\text{groups}|\text{three}) = 0$$

# Distributional Representation of Words

- ❖ Each word is represented by a low-dimensional dense vector

	Distributed representations
man	[0.326172, ..., 0.00524902, ..., 0.0209961]
woman	[0.243164, ..., -0.205078, ..., -0.0294189]
car	[0.0512695, ..., -0.306641, ..., 0.222656]
automobile	[0.107422, ..., -0.0375977, ..., -0.0620117]



Hinton, G. E., et al. Distributed representations. In Rumelhart, D. E., McClelland, J. L., and PDP Research Group, C., editors, *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, Vol. 1, 1986, pages 77–109. MIT Press, Cambridge, MA, USA.



# Advantages of Distributed Representations

- ❖ Beyond the independent assumption

	Distributed representations
man	[0.326172, ..., 0.00524902, ..., 0.0209961]
woman	[0.243164, ..., -0.205078, ..., -0.0294189]
car	[0.0512695, ..., -0.306641, ..., 0.222656]
automobile	[0.107422, ..., -0.0375977, ..., -0.0620117]

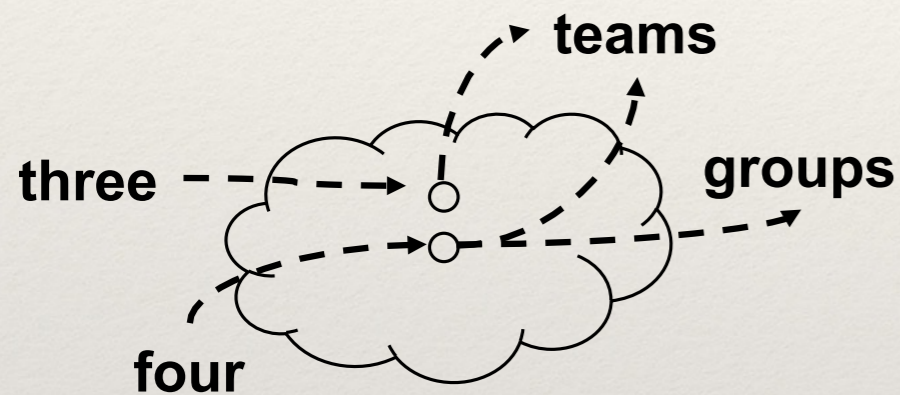
$$\cos(\text{man}, \text{woman}) = 0.77$$

$$\cos(\text{man}, \text{automobile}) = 0.25$$



# Advantages of Distributed Representations (cont')

- ❖ Better generalization ability: semantically similar words are mapped to nearby points



- ❖ Assigning probability to unseen bigram “**three groups**”  
 $P(\text{groups} \mid \text{three}) > 0$

Doc1: There are **three teams** left for the qualification

Doc2: **Four teams** have passed the first round

Doc3: **Four groups** are playing in the field

Language modeling with distributed word representations can assign probabilities to unseen bigrams according to their semantics

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

–*J. R. Firth (1957)*

“Words that occur in the same context tends to have similar meanings.”

–*Zelling Harris (1954)*

---

# What is the Meaning of “bardiwac”?

---

He handed her a glass of **bardiwac**.

Beef dishes are made to complement the **bardiwacs**.

Nigel staggered to his feet, face flushed from too much **bardiwac**.

Malbec, one of the lesser-known **bardiwac grapes**, responds well to Australia’s sunshine.

I dined off bread and cheese and this excellent **bardiwac**.

The **drinks** were delicious: blood-**red bardiwac** as well as light, sweet Rhenish.

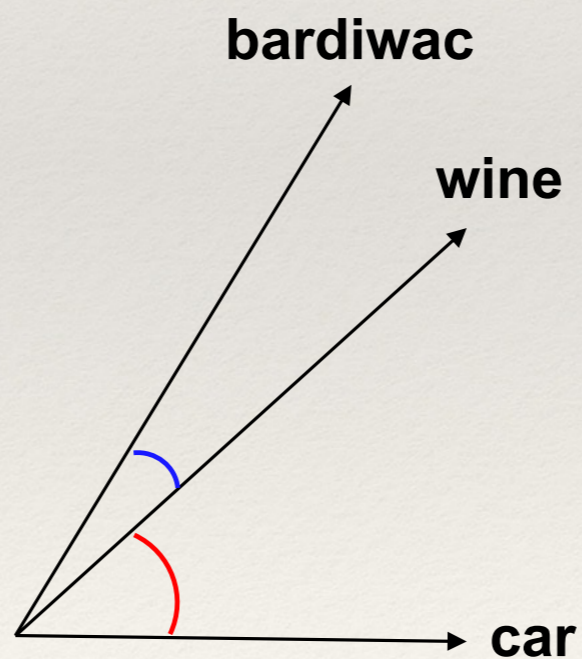
- ❖ A **red** alcoholic **beverage** made from **grapes**

# Surrounding Words

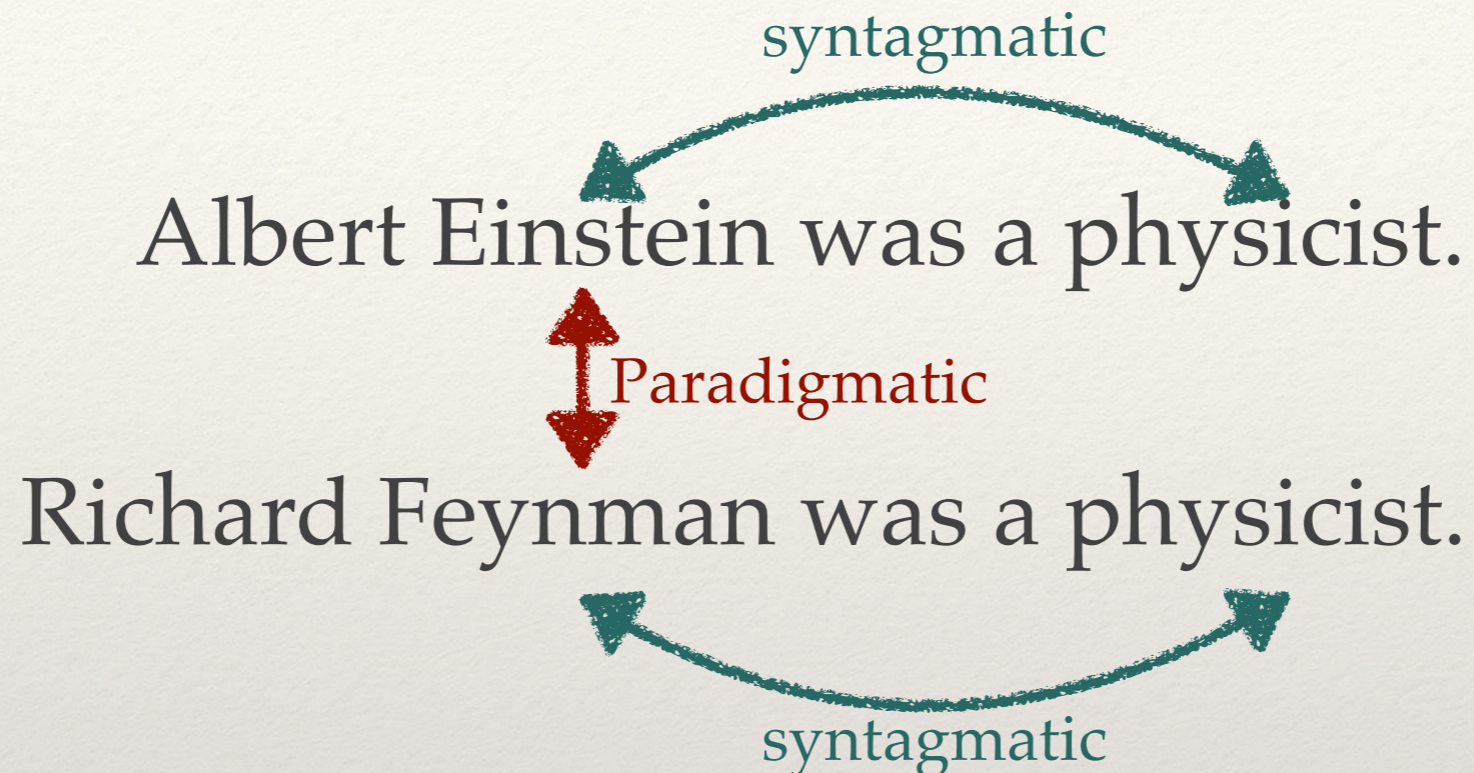
Just checking on the	<b>bardiwac</b>	he boomed as he come back
I hope you'll take to a good French	<b>bardiwac</b>	chimed in Arthur Iverson jovially
our host did slip out to attend to the	<b>bardiwac</b>	that was before the shrimp
Iverson did when he went through to see to	<b>bardiwac</b>	before dinner. Henry rubbed his hands
and drinking red win from France -- sour	<b>bardiwac</b>	, which bad proved hard to sell.
eyes were alight and he was drinking the	<b>bardiwac</b>	down like water. It is like Hallow-fair
quizzically at him and offering him some more	<b>bardiwac</b>	. He shook his head. 'I will sleep
drinks (as Queen Victoria reputedly did with	<b>bardiwac</b>	and malt whisky), but still the result
do we really 'wash down' a good meal with	<b>bardiwac</b>	? Port is immediately suggested by Stilton
completely different: cheap and cheerful	<b>bardiwac</b>	. Two good examples from Victoria Wine are
examples from Victoria Wine are its house	<b>bardiwac</b>	, juicy and touch almondy, a good buy
opened a bottle of rather rust-coloured	<b>bardiwac</b>	. I ate too much and drank nearly three-quarters
elections, it was apparent the SDP of '	<b>bardiwac</b>	and chips' mould-breaking fame at the time
the black hills. Not a night of vintage	<b>bardiwac</b>	. burnley: Pearce, Measham, McGrory
SONS Old School – the Marlborian navy,	<b>bardiwac</b>	and slim-white stripe. Heavy woven silk
white-hot passion, We are like a good bottle of	<b>bardiwac</b>	; we both have sediment in our shoes
few minutes later he was uncorking a fine	<b>bardiwac</b>	in Masha's room, saying he had something
the phone, Surkov silently offered me more	<b>bardiwac</b>	but I indicated a bottle of Perrier

# Word-Word Co-occurrence

	glass	drink	grape	rex	meal
<b>bardiwac</b>	10	22	43	16	29
<b>wine</b>	14	10	4	15	45
<b>car</b>	5	0	0	10	0



# Two Interpretations of Distributed Hypothesis



- ❖ **Syntagmatic:** words co-occur in the same text region
- ❖ **Paradigmatic:** words occur in the same context, may not at the same time

Sahlgren, M. (2008). The distributional hypothesis. *Italian Journal of Linguistics*, 20(1):33–54.

Fei Sun et al. Learning Word Representations by Jointly Modeling Syntagmatic and Paradigmatic Relations. In *Proceedings of ACL*. 2015, 136–145



# Modeling the Syntagmatic Relation

syntagmatic

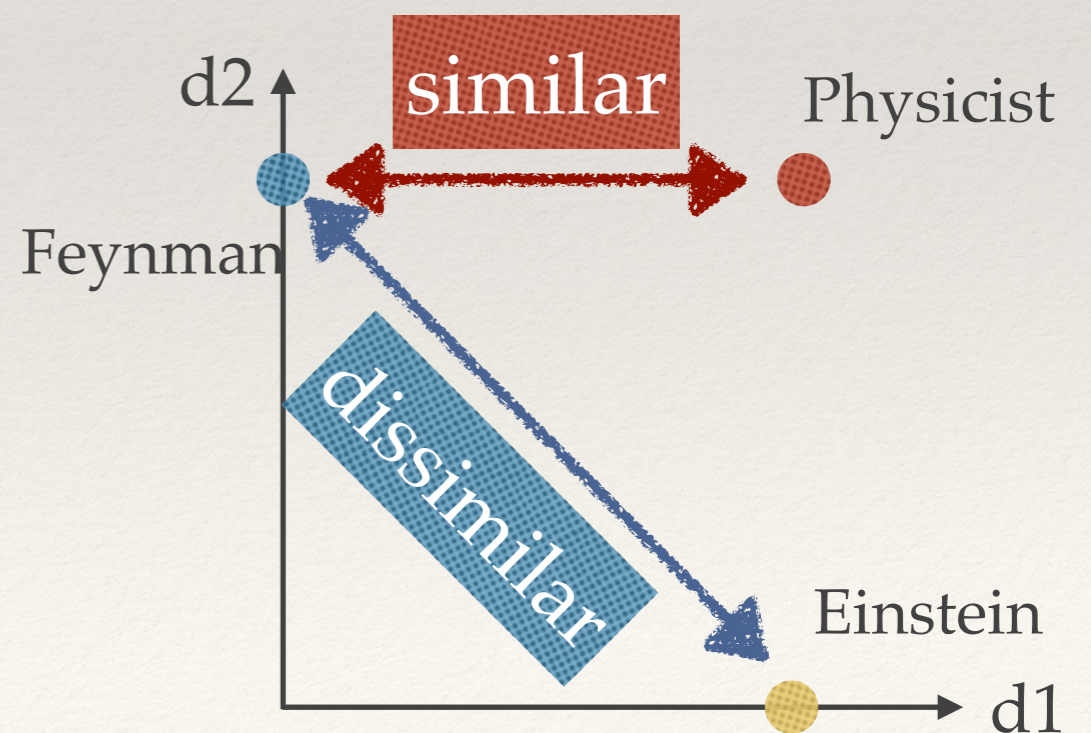
Albert Einstein was a physicist.

Richard Feynman was a physicist.

syntagmatic

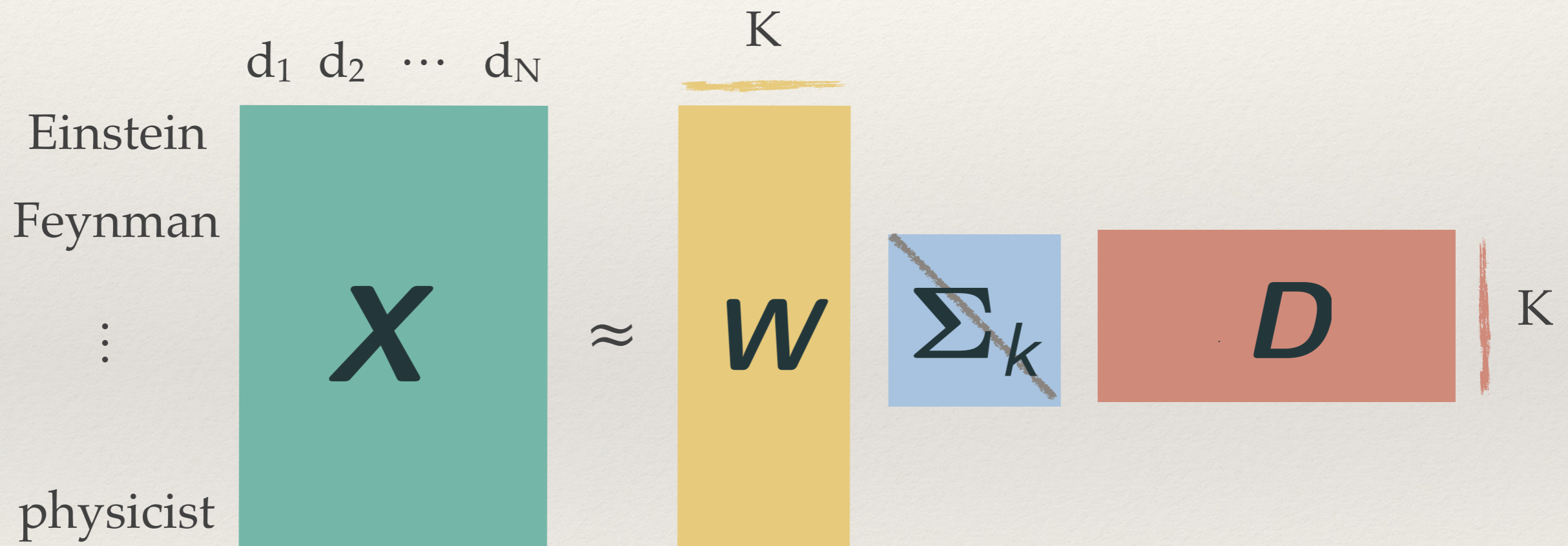
Word-document co-occurrence matrix  
(words represented by documents)

	d1	d2
Einstein	1	0
Feynman	0	1
Physicist	1	1



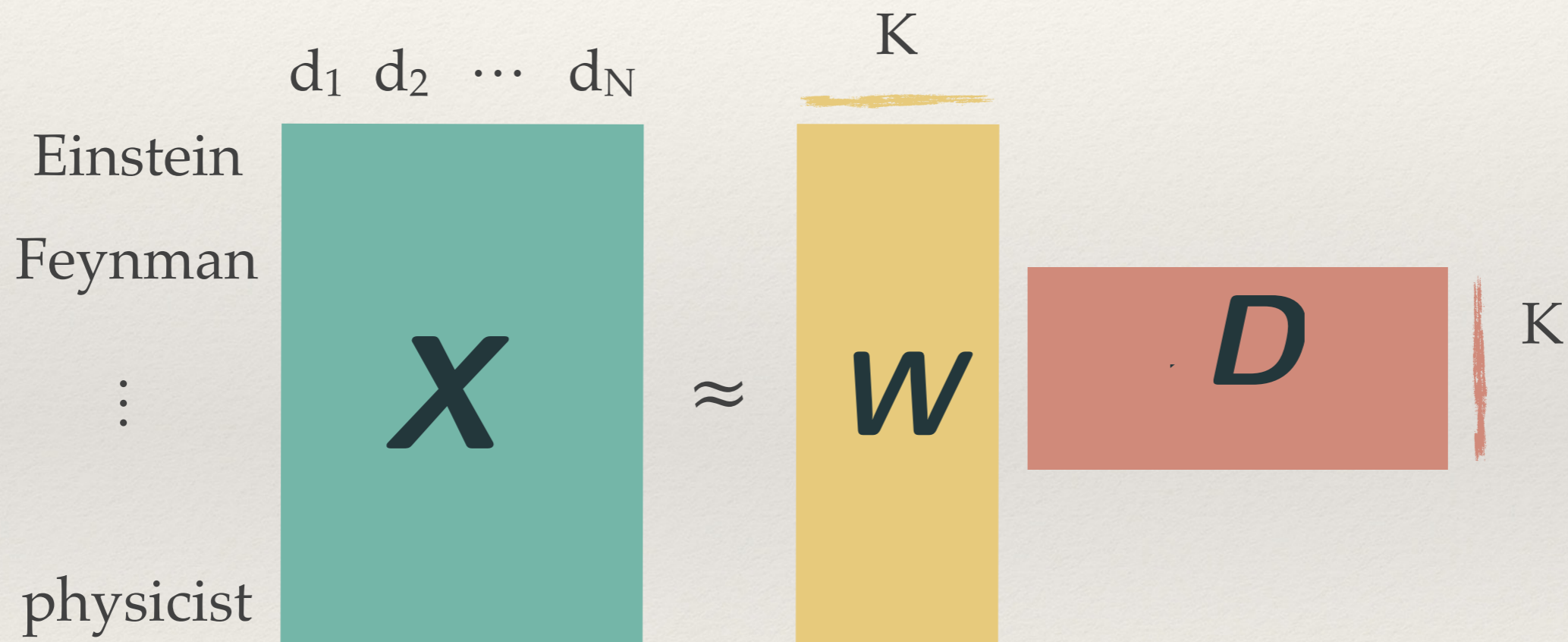
# Modeling Syntagmatic Relation – LSI

- ❖ Rank-reduced SVD of document-word co-occurrence matrix



$$X \approx W \Sigma_k D^T$$

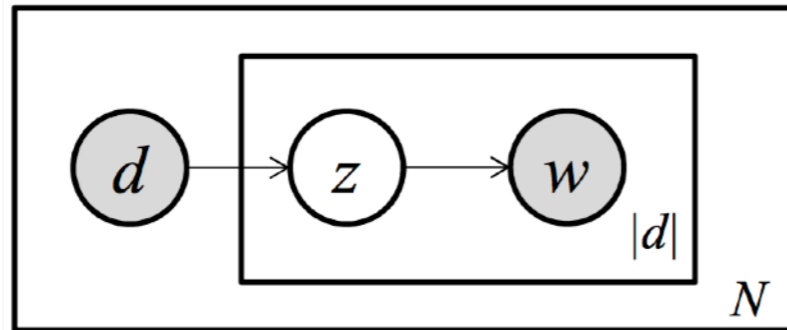
# Modeling Syntagmatic Relation – NMF



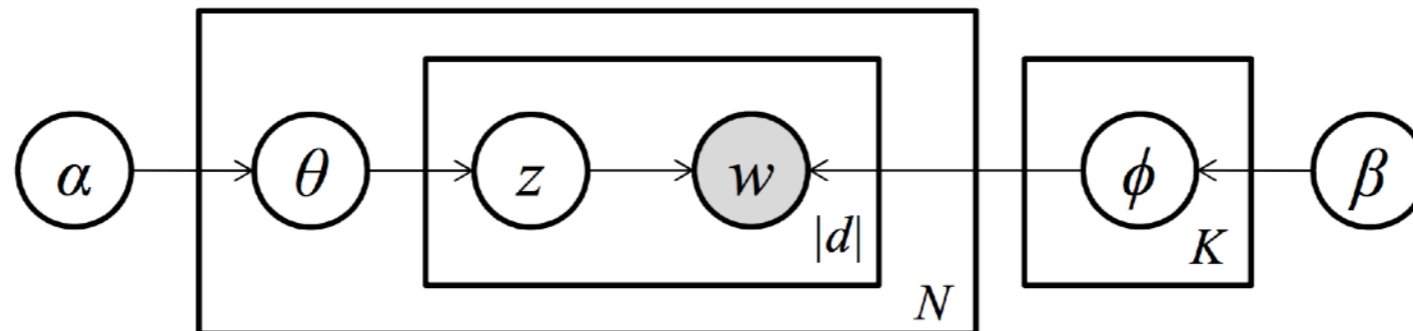
$$\mathbf{X} \approx \mathbf{W} \mathbf{D}^T$$

# Modeling Syntagmatic Relation – PLSA and LDA

PLSA



LDA



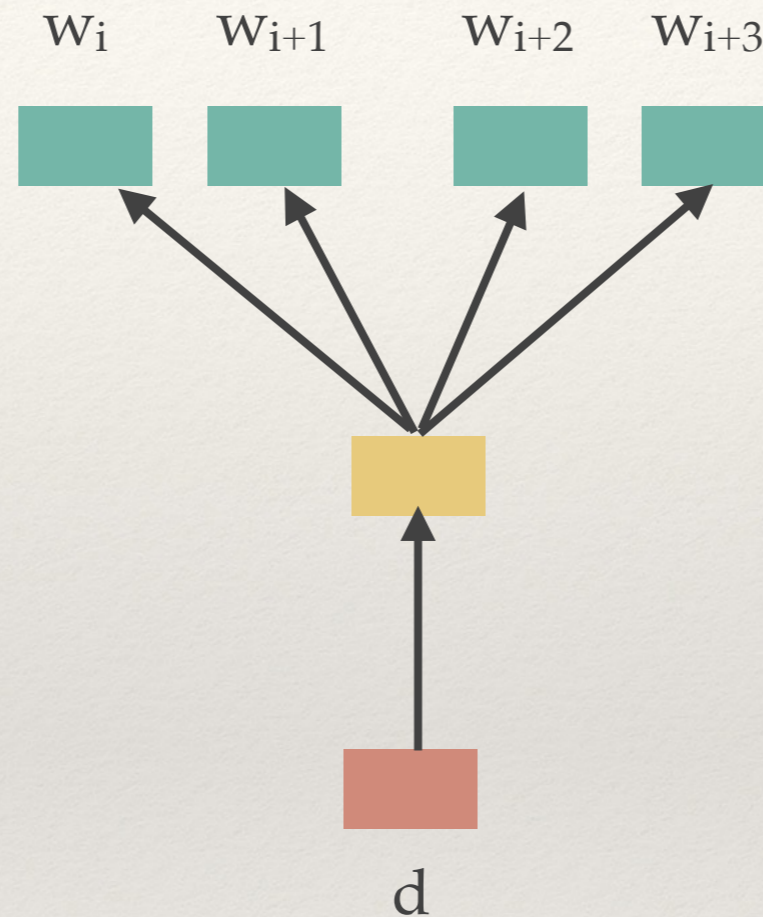
- ❖ Maximum likelihood solution of PLSA is NMF with KL divergence

---

## Modeling Syntagmatic Relation

### – Distributed Bag of Words Version of Paragraph Vector (PV-DBOW)

---



Predict word vector using document vector.

# Modeling the Paradigmatic Relation

Albert Einstein was a physicist.

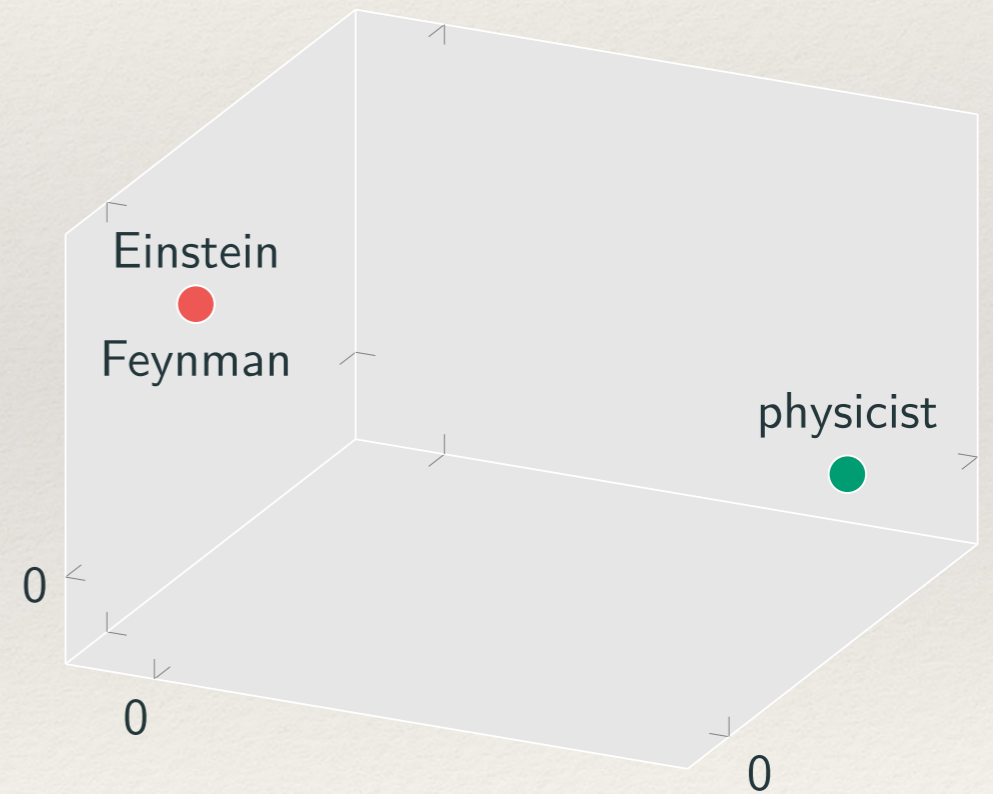


Paradigmatic

Richard Feynman was a physicist.

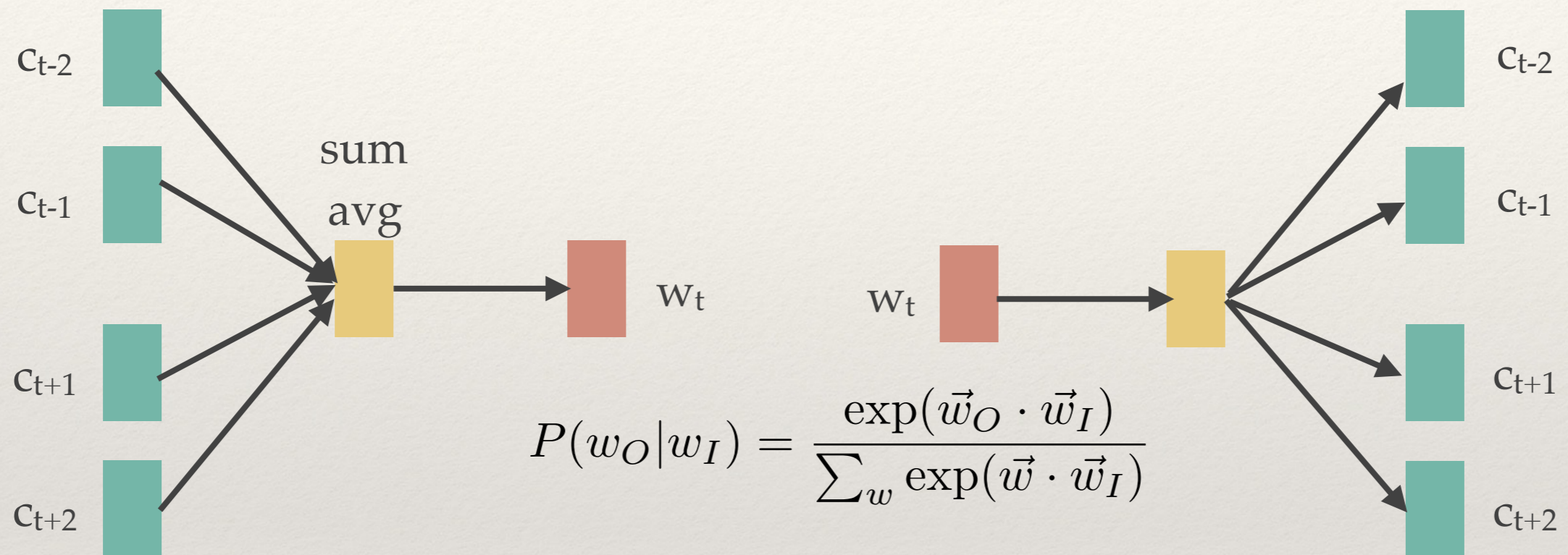
Word-word co-occurrence matrix  
(words represented by other words)

	Einstein	Feynman	Physicist
Einstein	0	0	1
Feynman	0	0	1
Physicist	1	1	0



More suitable for learning the embeddings from short documents.

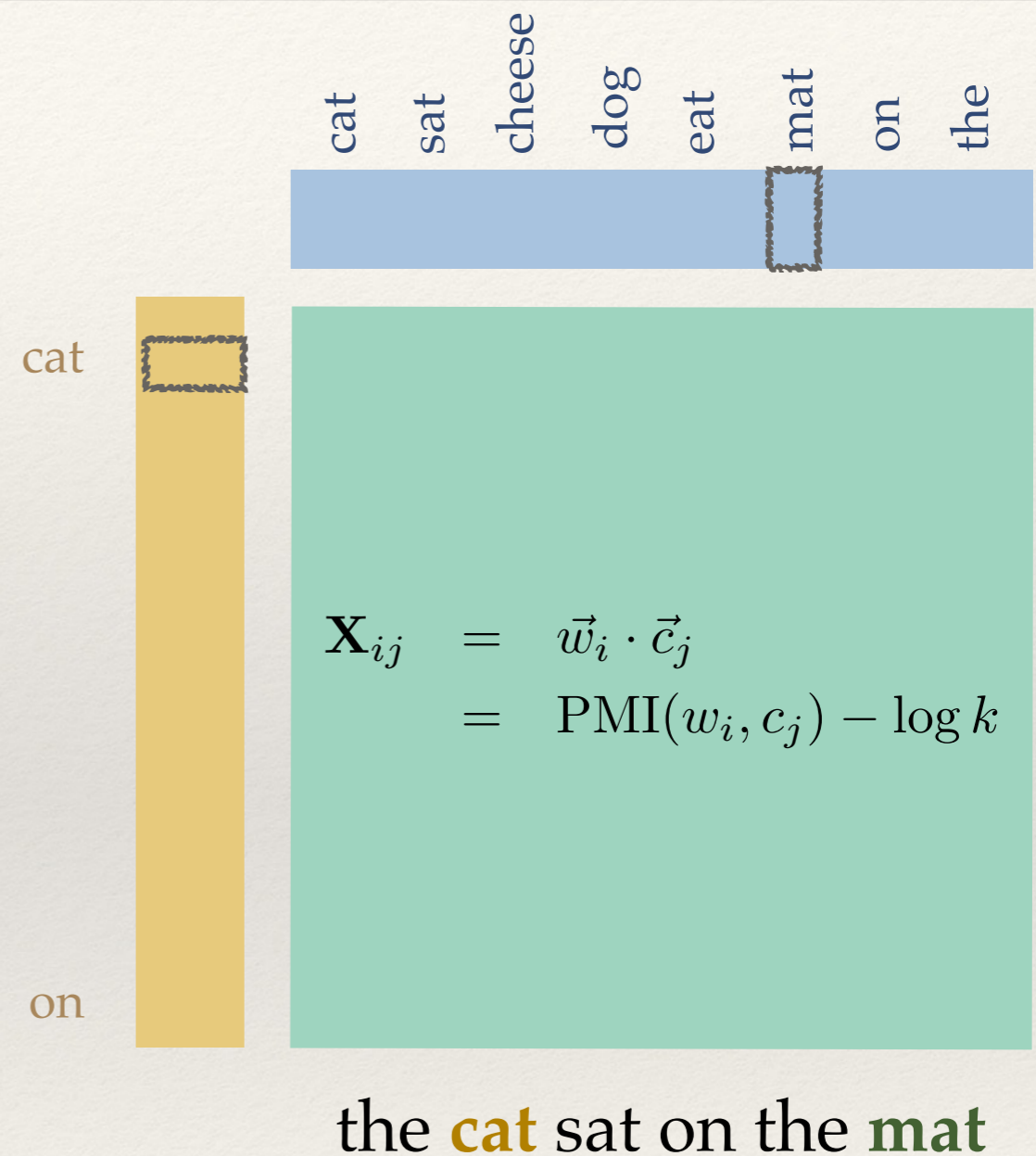
# Modeling the Paradigmatic Relation – Word2Vec



- ❖ Usually optimized with negative sampling

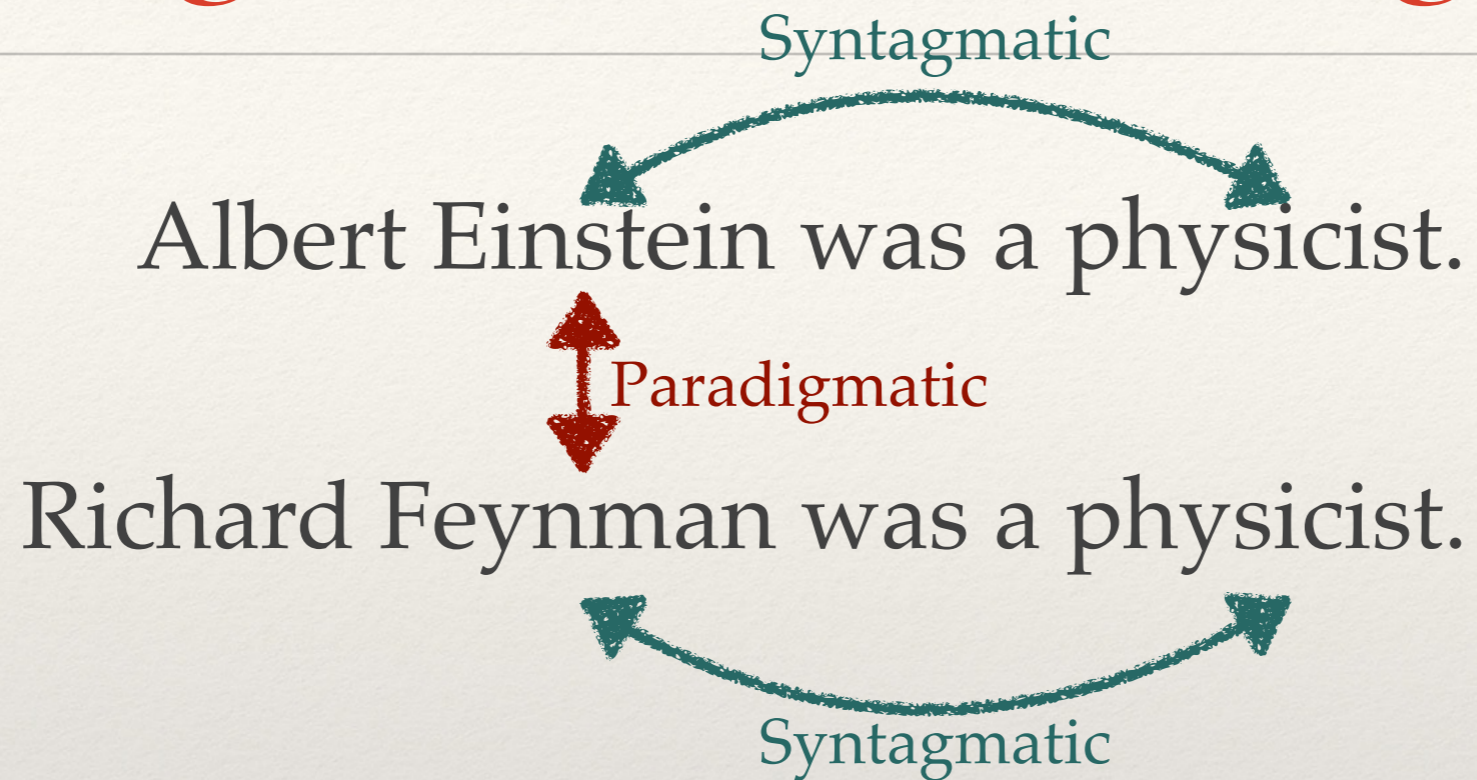
# Word Embedding as Matrix Factorization

- ❖ Skip-gram negative sampling (with sampling values  $k > 1$ ) is factorizing the shifted point wise mutual information (PMI) matrix





# Syntagmatic v.s. Paradigmatic



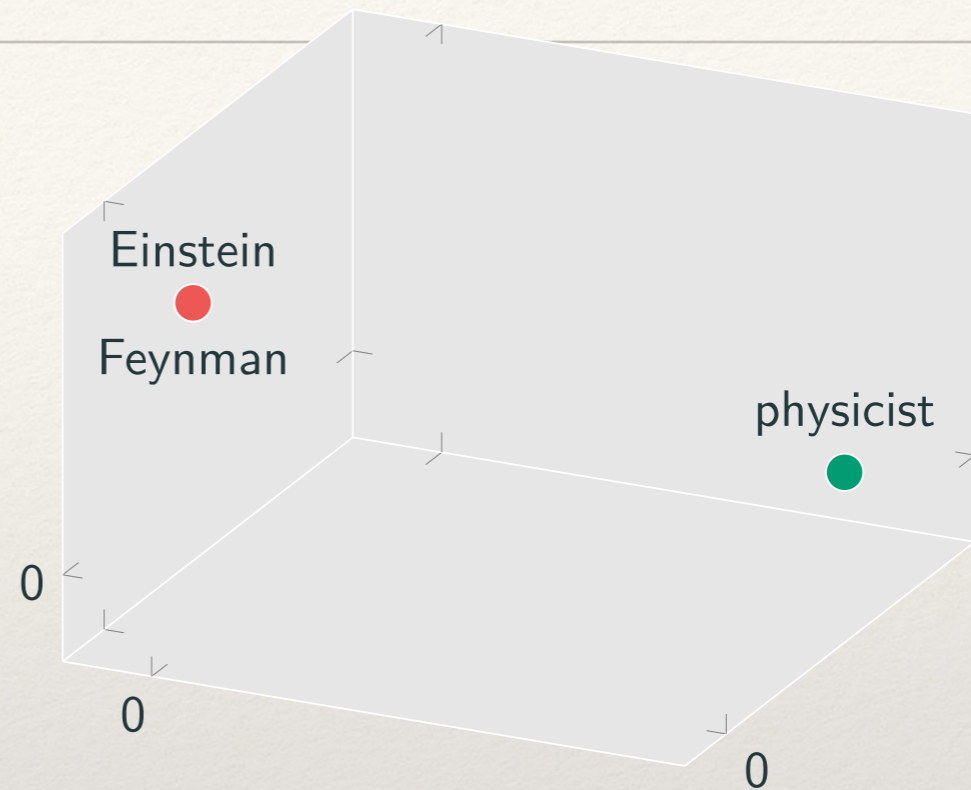
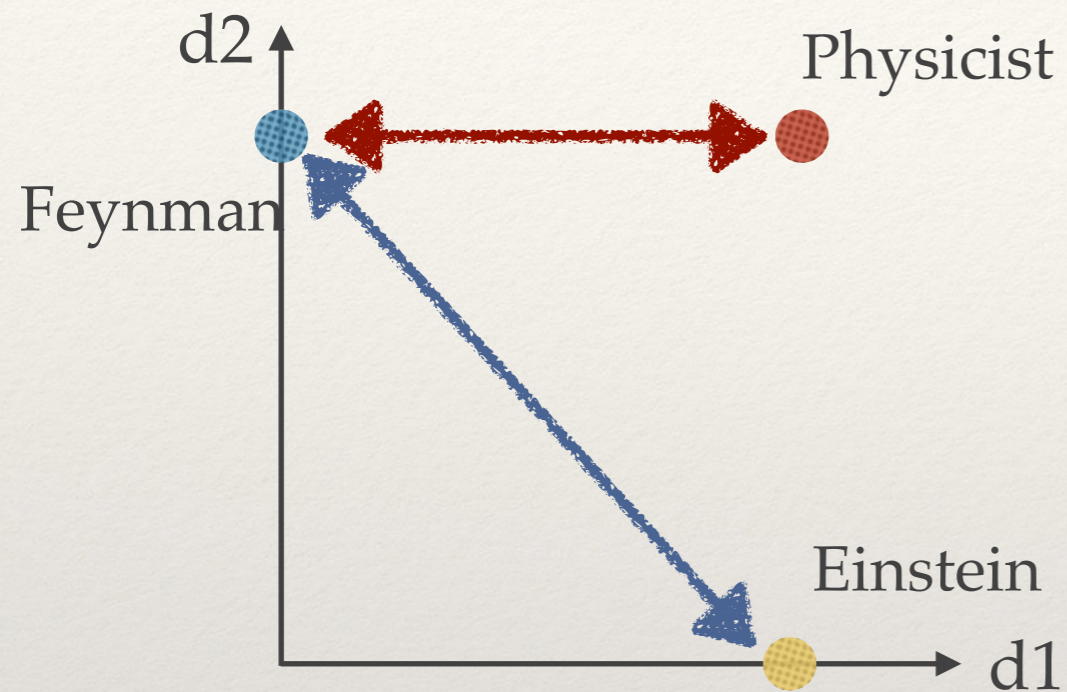
Word-document co-occurrence matrix  
(words represented by documents)

	d1	d2
Einstein	1	0
Feynman	0	1
Physicist	1	1

Word-word co-occurrence matrix  
(words represented by other words)

	Einstein	Feynman	Physicist
Einstein	0	0	1
Feynman	0	0	1
Physicist	1	1	0

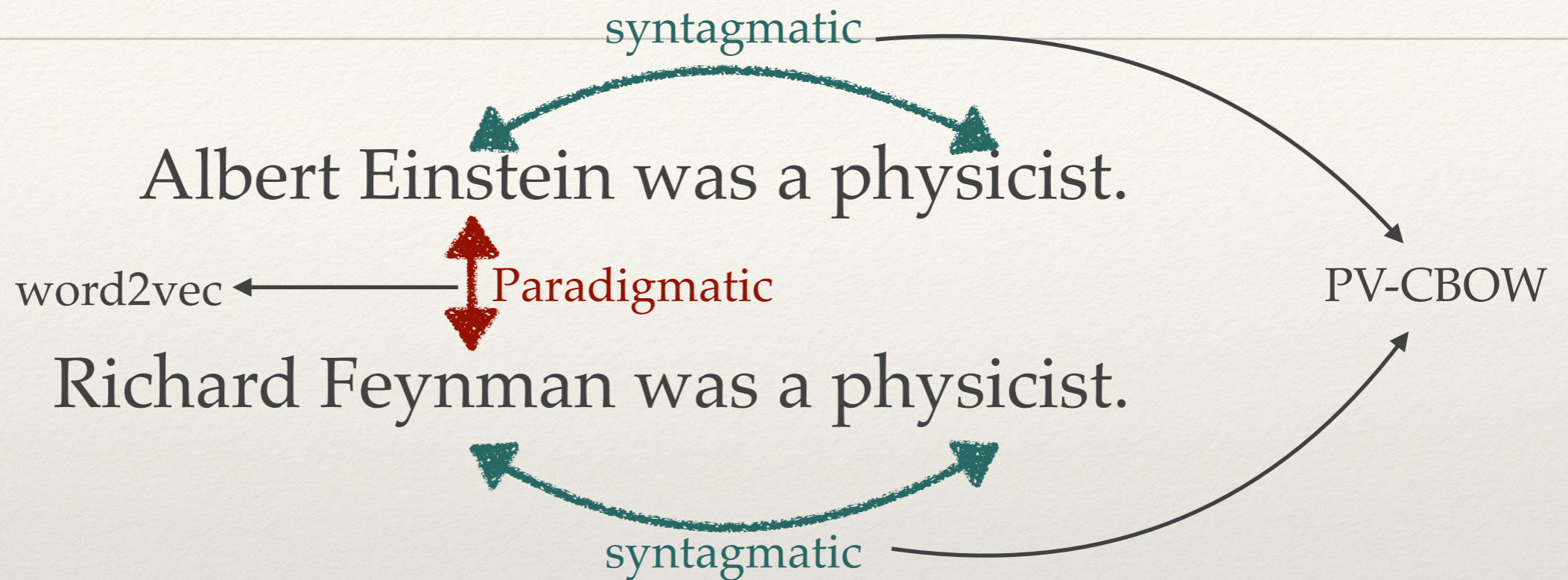
# Syntagmatic v.s. Paradigmatic (cont')



## Similar words to "Feynman"

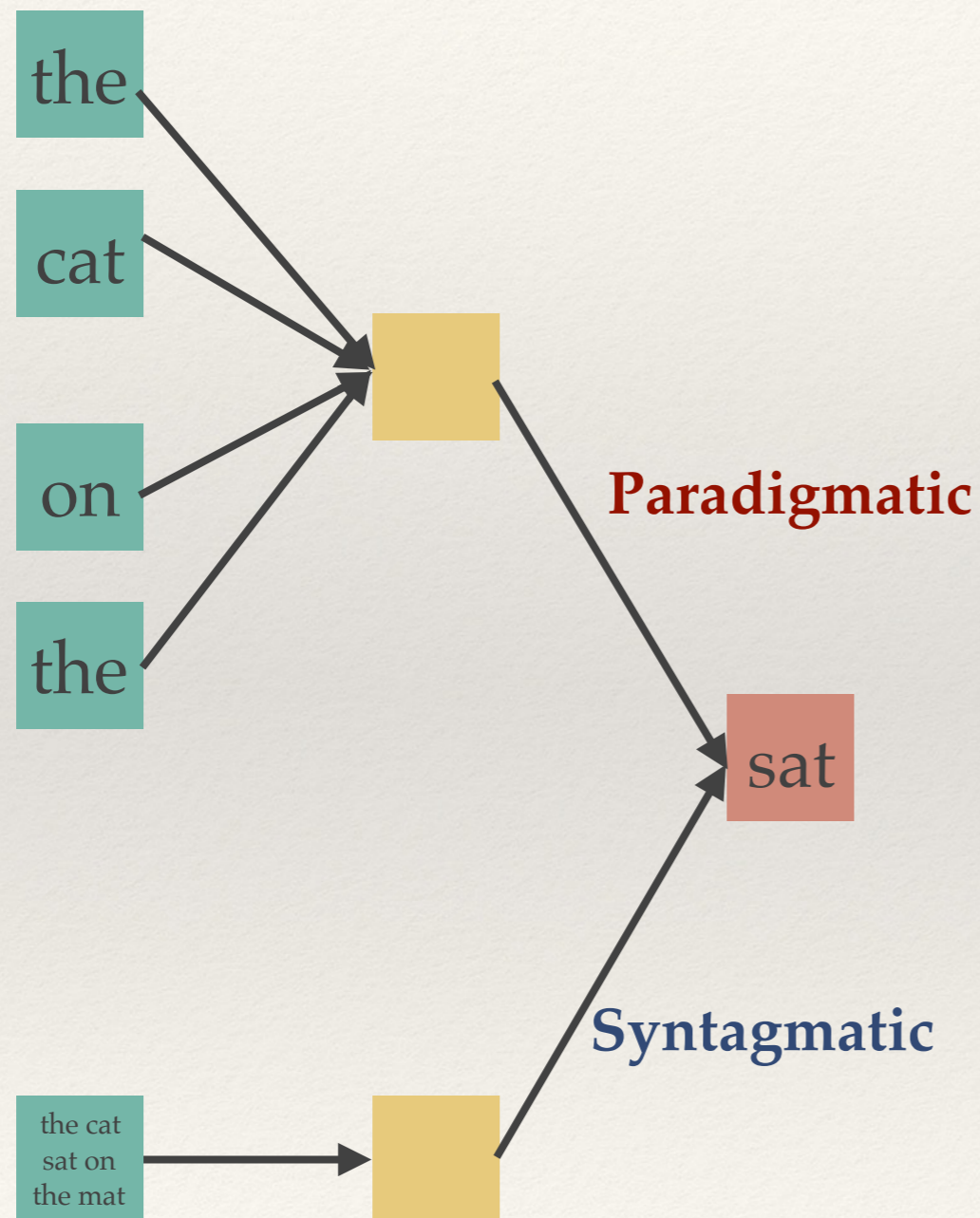
Syntagmatic	Paradigmatic
quantum	Einstein
physicist	Schwinger
electrodynamics	Bethe
relativity	Bohm

# A Natural Extension: Modeling them Jointly



- ❖ Construct the model under word2vec framework
  - ❖ **Paradigmatic**: modeling with word2vec
  - ❖ **Syntagmatic**: modeling with PV-CBOW

# Parallel Document Content (PDC) Model



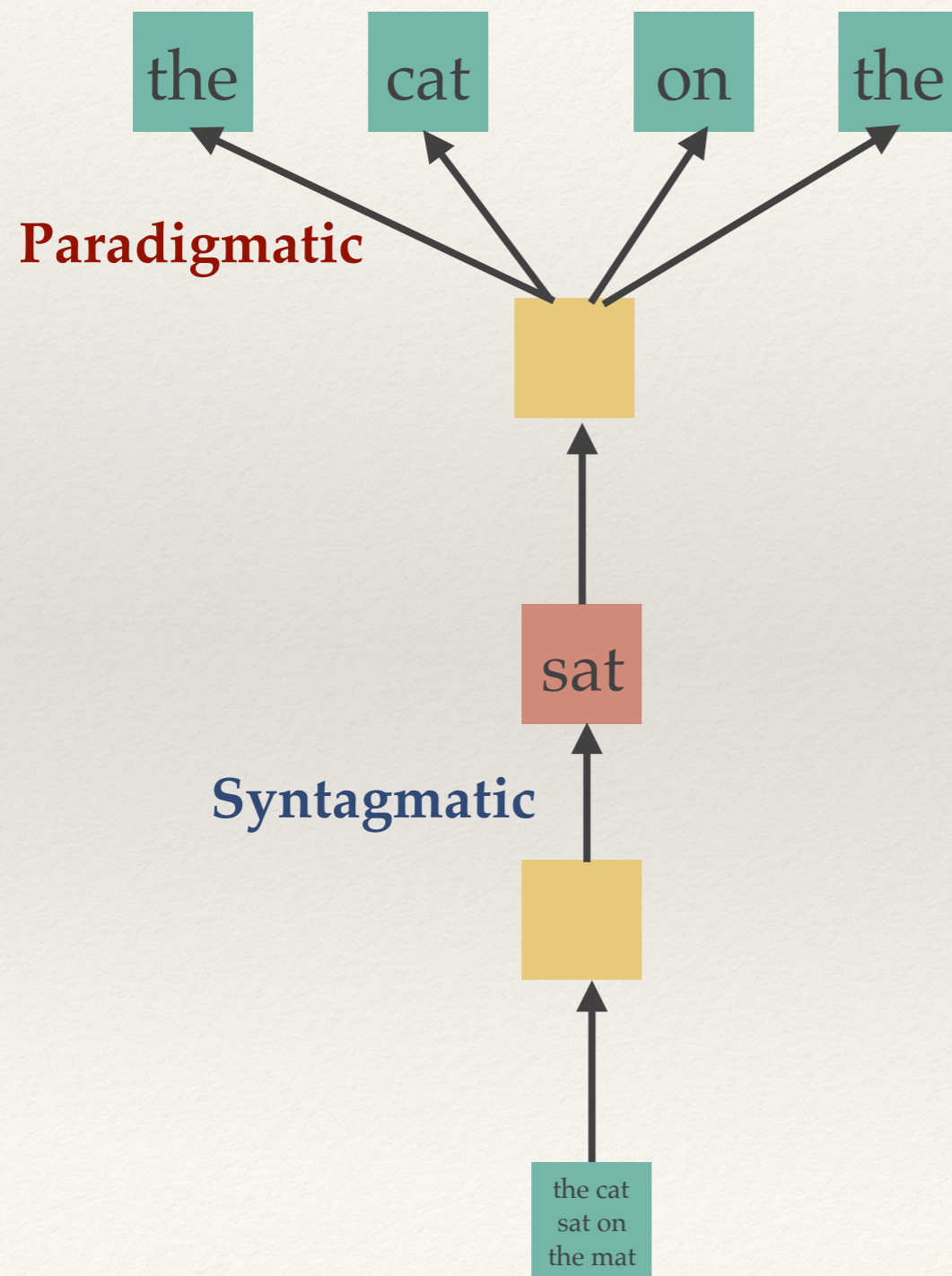
$$\ell = \sum_{n=1}^N \sum_{w_i^n \in d_n} \left( \log p(w_i^n | h_i^n) + \log p(w_i^n | d_n) \right)$$

Negative Sampling

$$\ell = \sum_{n=1}^N \sum_{w_i^n \in d_n} \left( \log \sigma(\vec{w}_i^n \cdot \vec{h}_i^n) + \log \sigma(\vec{w}_i^n \cdot \vec{d}_n) \right. \\ \left. + k \cdot \mathbb{E}_{w' \sim P_{nw}} \log \sigma(-\vec{w}' \cdot \vec{h}_i^n) \right. \\ \left. + k \cdot \mathbb{E}_{w' \sim P_{nw}} \log \sigma(-\vec{w}' \cdot \vec{d}_n) \right)$$

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

# Hierarchical Document Context Model (HDC)



$$\ell = \sum_{n=1}^N \sum_{w_i^n \in d_n} \left( \log p(w_i^n | d_n) + \sum_{\substack{j=i-L \\ j \neq i}}^{i+L} \log p(c_j^n | w_i^n) \right)$$

Negative Sampling

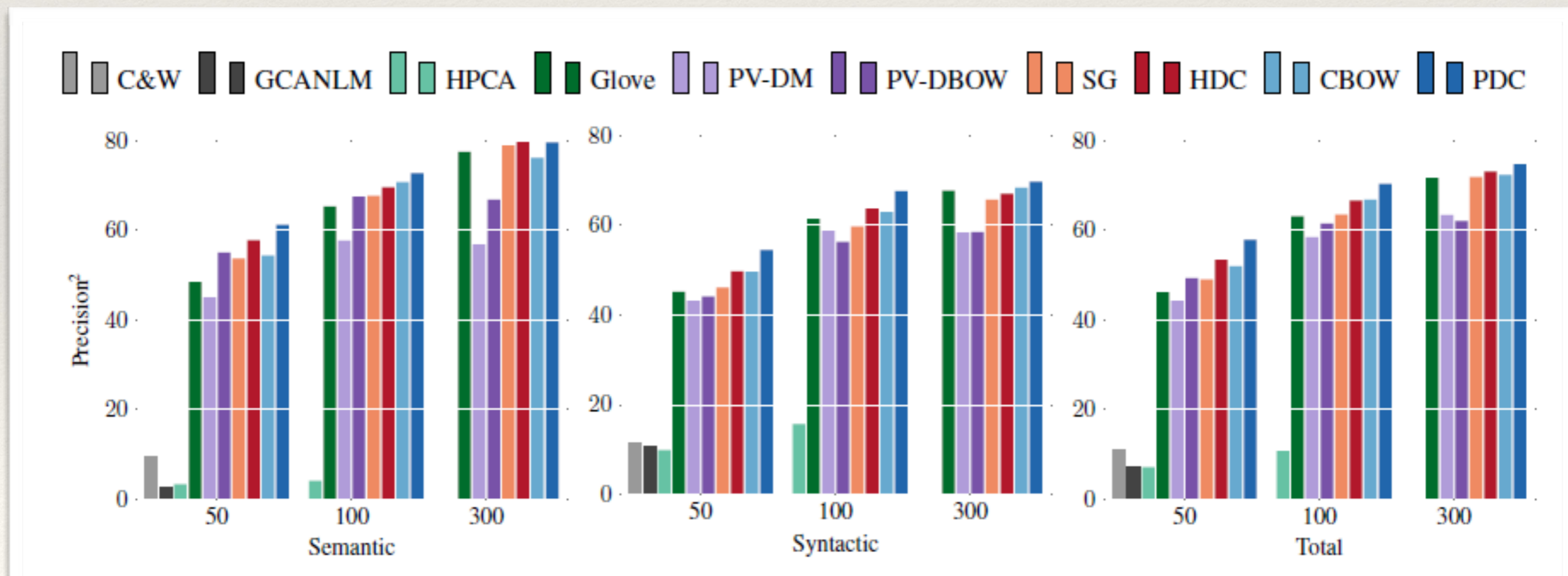
$$\ell = \sum_{n=1}^N \sum_{w_i^n \in d_n} \left( \sum_{\substack{j=i-L \\ j \neq i}}^{i+L} \left( \log \sigma(\vec{c}_j^n \cdot \vec{w}_i^n) \right) \right.$$

$$\left. + k \cdot \mathbb{E}_{c' \sim P_{nc}} \log \sigma(-\vec{c}' \cdot \vec{w}_i^n) \right)$$

$$+ \log \sigma(\vec{w}_i^n \cdot \vec{d}_n) + k \cdot \mathbb{E}_{w' \sim P_{nw}} \log \sigma(-\vec{w}' \cdot \vec{d}_n)$$

# Empirical Evaluation of PDC and HDC

- ❖ Word analogy based on Google test set
  - ❖ Semantic: “Beijing is to China as Paris is to \_\_\_\_”
  - ❖ Syntactic: “big is to bigger as deep is to \_\_\_\_”



# More Diverse Results

Paradigmatic



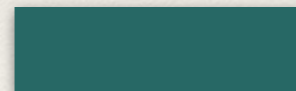
Syntagmatic

Top five similar words to "Feynman"

CBOW	SG	PDC	HDC	PV-DBOW
Einstein	Schwinger	geometroynamics	Schwinger	physicists
Schwinger	quantum	Bethe	electrodynamics	spacetime
Bohm	Bethe	semiclassical	Bethe	geometroynamics
Bethe	Einstein	Schwinger	semiclassical	tachyons
relativity	semiclassical	peturbative	quantum	Einstein

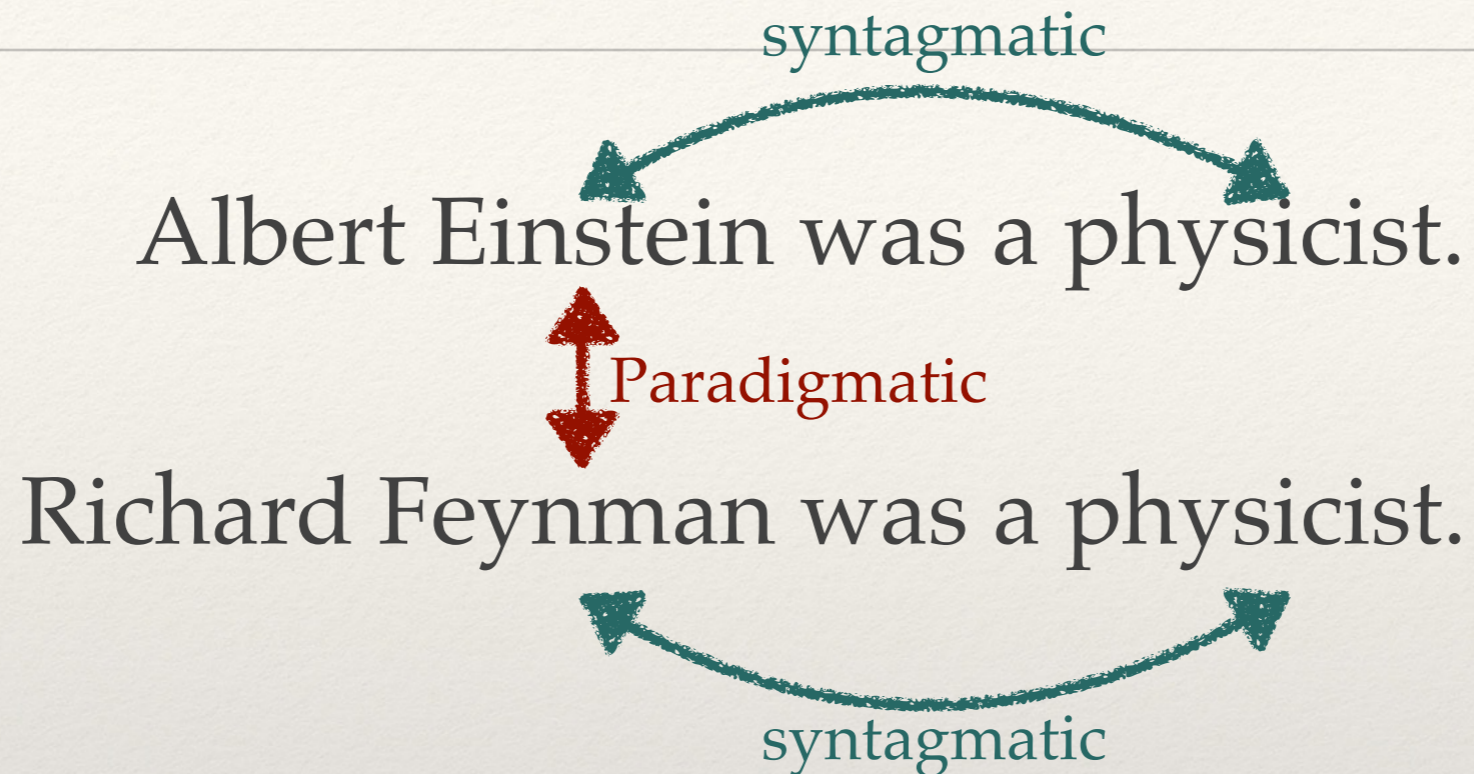


Paradigmatic



Syntagmatic

# Re-examine the Distributed Hypothesis



- ❖ **Syntagmatic**: words co-occur in the same text region
- ❖ **Paradigmatic**: words occur in the same context, may not at the same time
- ❖ Distributed hypothesis considers words as IDs
  - ❖ However, words are constructed by more fine-grained elements, e.g., breakable —> break, able



---

# Beyond Distributed Hypothesis

---

- ❖ Distributed hypothesis: discovering semantics of words from **external** information
- ❖ Beyond distributed hypothesis: discovering semantics of words from both **external** and **internal** information
  - ❖ External: distributed hypothesis
  - ❖ Internal: words are built from morphemes, e.g., breakable → break, able

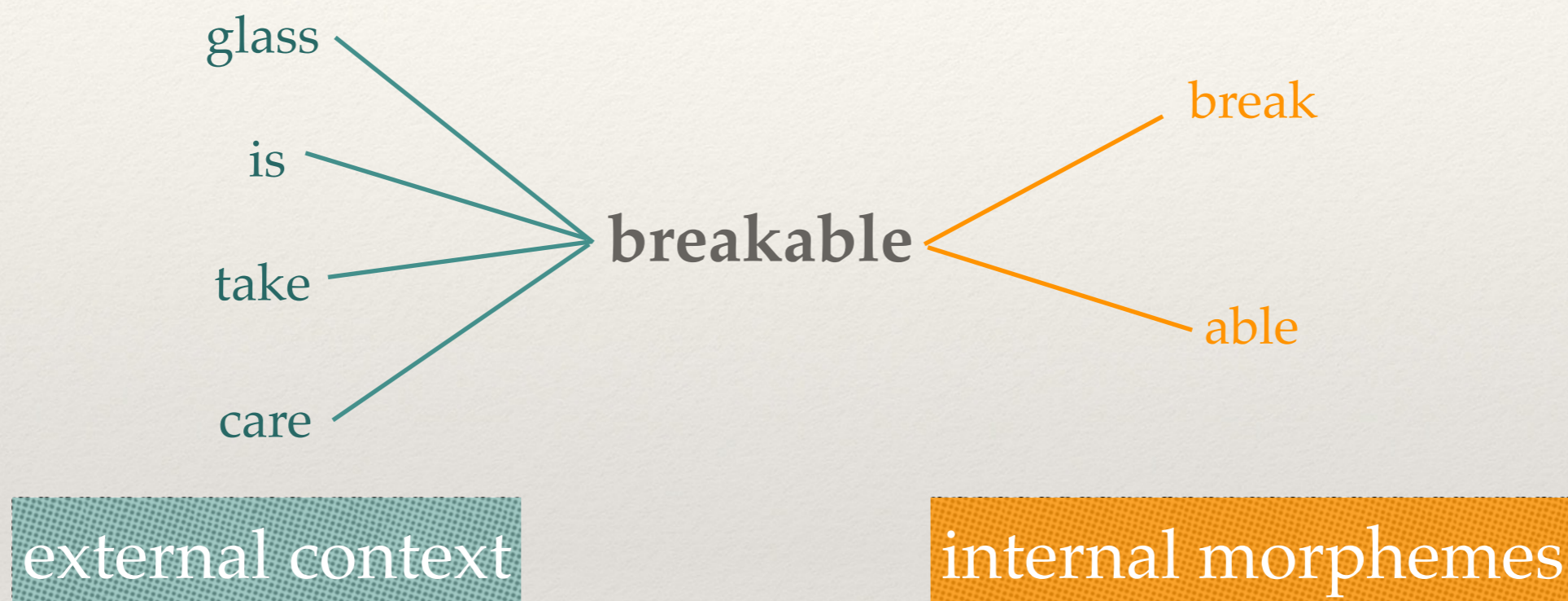
glass is **breakable**, take care

He **breaks** the glass

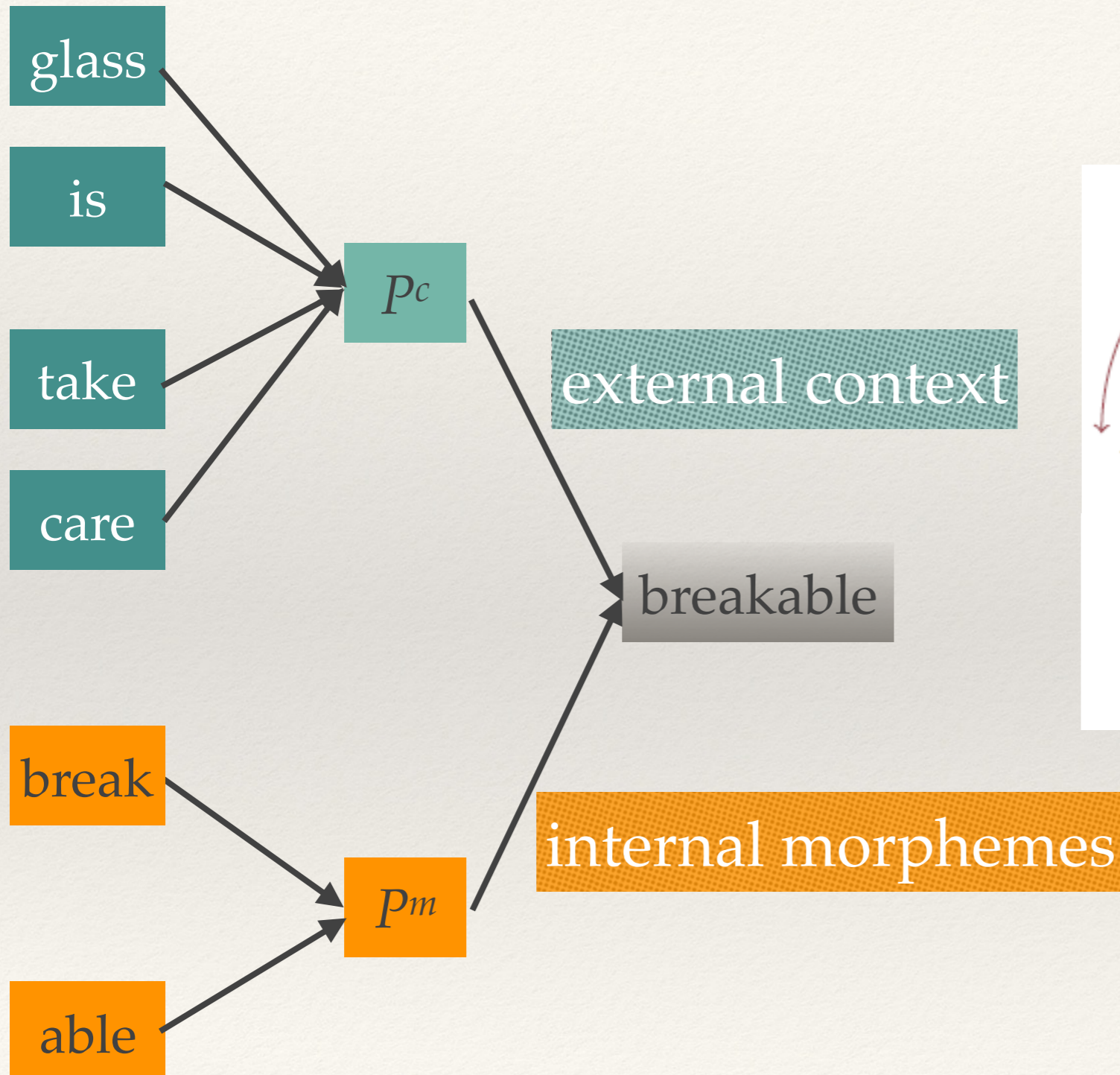
Similar embeddings for “**breakable**” and “**break**”

# Word Embedding with Morphemes

“... glass is **breakable**, take care ...”



# Continuous Bag of External and Internal Gram (BEING)

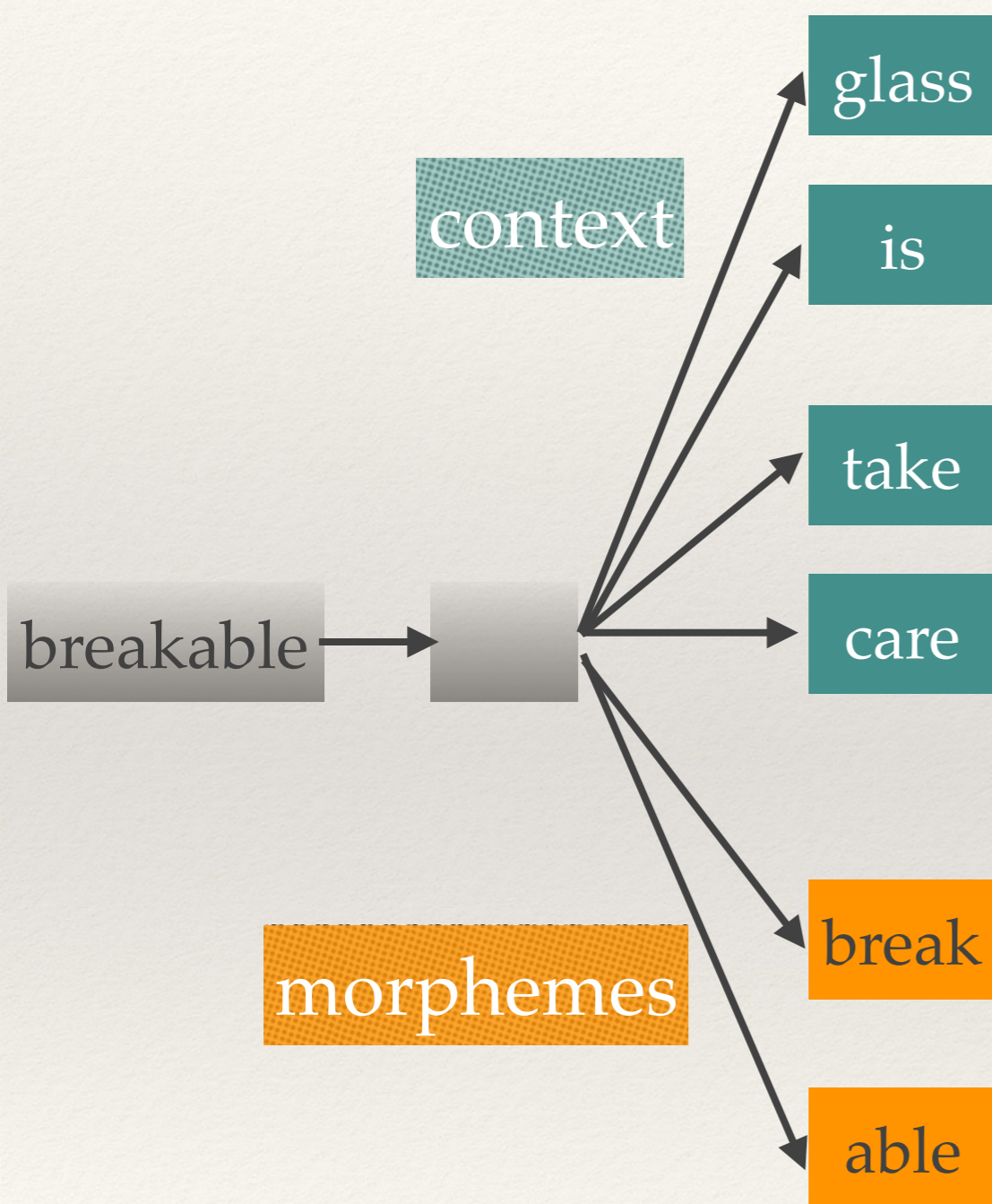


$$\mathcal{L} = \sum_{i=1}^N (\log p(w_i | \mathcal{P}_i^c) + \log p(w_i | \mathcal{P}_i^m))$$

Negative Sampling

$$\mathcal{L} = \sum_{i=1}^N (\log \sigma(\vec{w}_i \cdot \vec{\mathcal{P}}_i^c) + k \cdot \mathbf{E}_{\vec{w} \sim P_{\vec{w}}} \log \sigma(-\vec{w} \cdot \vec{\mathcal{P}}_i^c) + \log \sigma(\vec{w}_i \cdot \vec{\mathcal{P}}_i^m) + k \cdot \mathbf{E}_{\vec{w} \sim P_{\vec{w}}} \log \sigma(-\vec{w} \cdot \vec{\mathcal{P}}_i^m))$$
$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

# Continuous Skip External and Internal Gram (SEING)



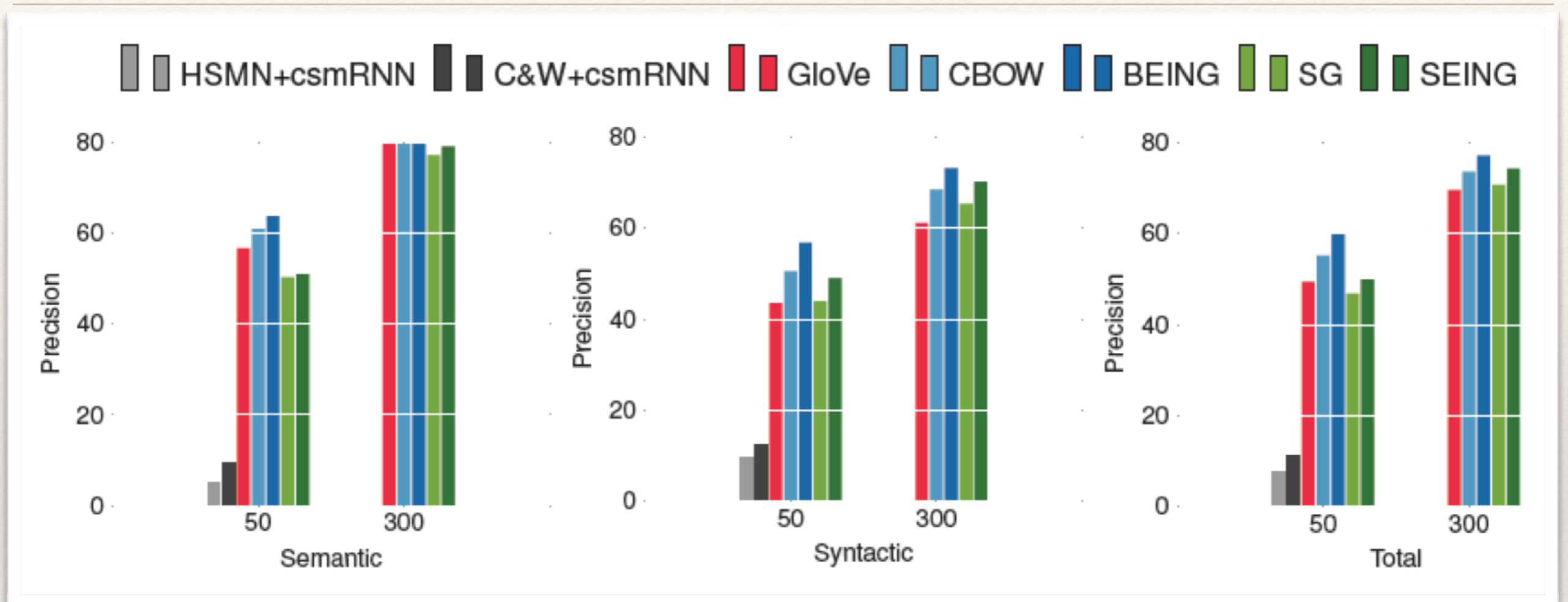
$$\mathcal{L} = \sum_{i=1}^N \left( \sum_{\substack{j=i-1 \\ j \neq i}}^{i+1} \log p(c_j | w_i) + \sum_{z=1}^{s(w_i)} \log p(m_i^{(z)} | w_i) \right)$$

Negative Sampling

$$\mathcal{L} = \sum_{i=1}^N \left( \sum_{\substack{j=i-1 \\ j \neq i}}^{i+1} \left( \log \sigma(\vec{c}_j \cdot \vec{w}_i) + k \cdot \mathbf{E}_{\vec{c} \sim P_C} \log \sigma(-\vec{c} \cdot \vec{w}_i) \right) + \sum_{z=1}^{s(w_i)} \left( \log \sigma(\vec{m}_i^{(z)} \cdot \vec{w}_i) + k \cdot \mathbf{E}_{\vec{m} \sim P_M} \log \sigma(-\vec{m} \cdot \vec{w}_i) \right) \right)$$

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

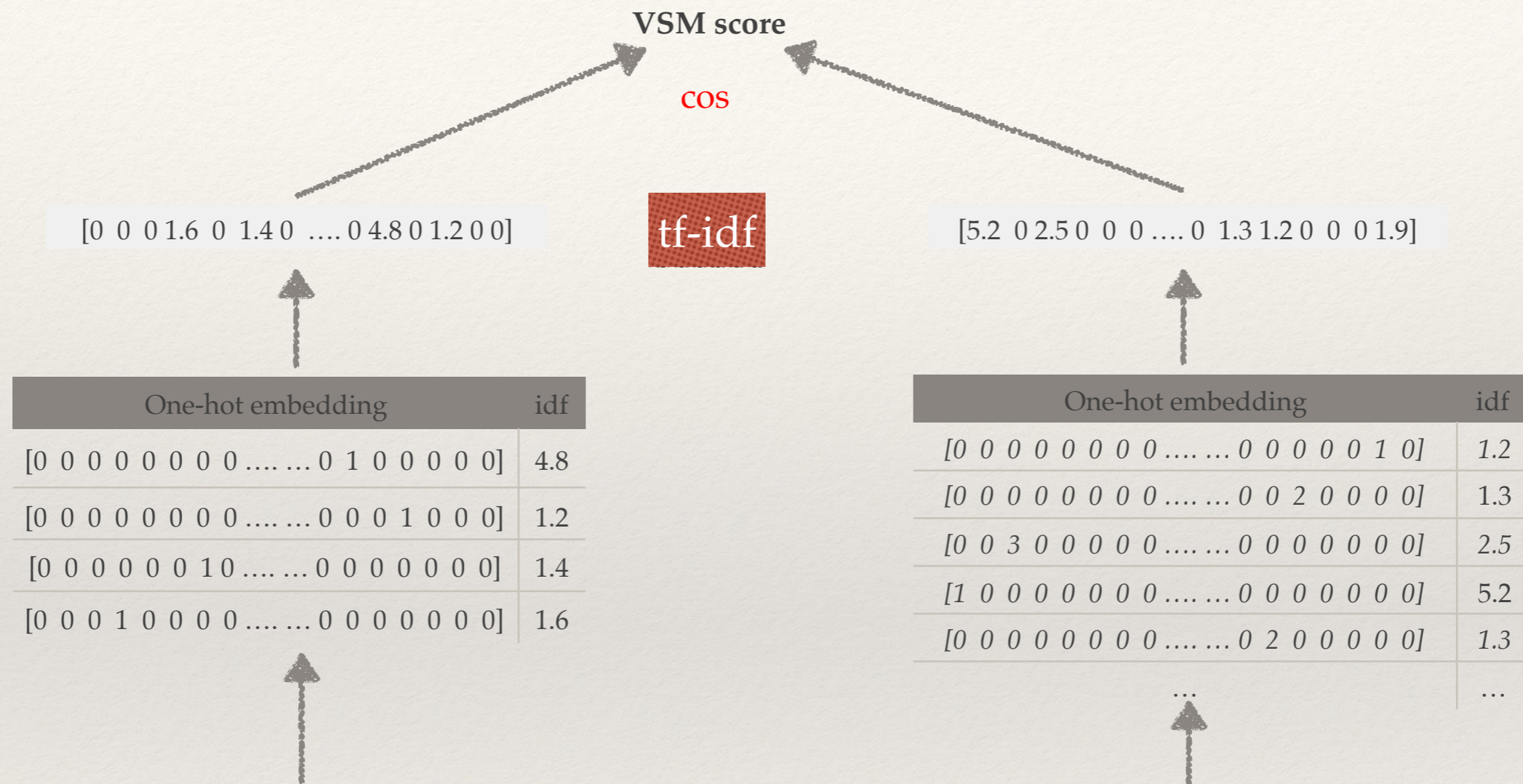
# Empirical Evaluation of BEING and SEING



- ❖ BEING and SEING outperformed C&W and SG, respectively
- ❖ Significant improvements achieved on **syntactic** task

# Direct Matching with Word Embeddings

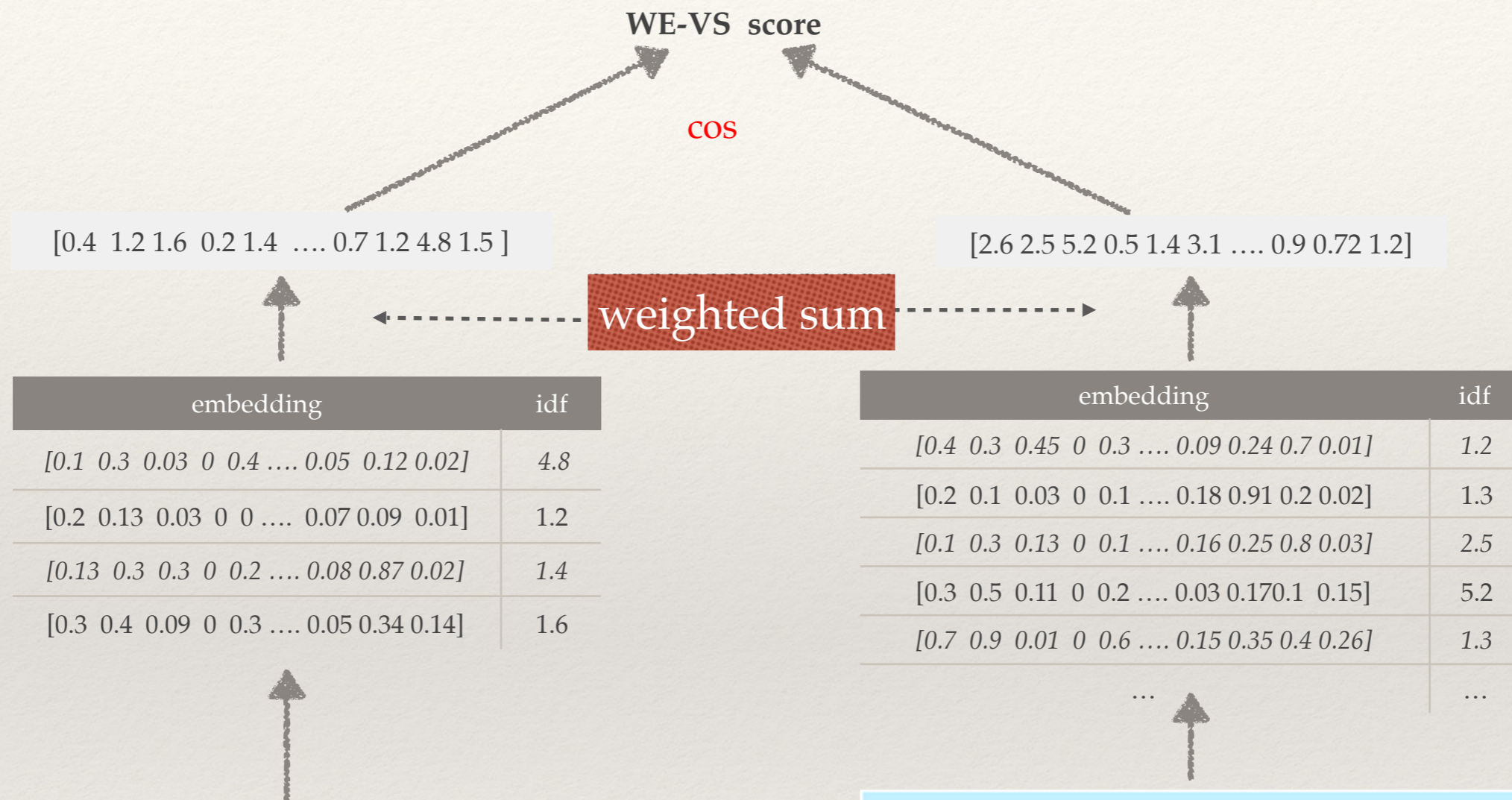
# Query-Document Matching based on Local Representations



peace process in the Middle East

As for the Arabian and Palestinian voices that are against the current negotiations and the so-called peace process, they are not against peace per se, but rather for their well-founded predictions that Israel would NOT give an inch of the West bank (and most probably the same for Golan Heights) back to the Arabs. An 18 months of "negotiations" in Madrid, and Washington proved these predictions. Now many will jump on me saying why are you blaming israelis for no-result negotiations. I would say why would the Arabs stall the negotiations, what do they have to loose ?

# When Embedding Comes ...



embedding	idf
[0.1 0.3 0.03 0 0.4 ... 0.05 0.12 0.02]	4.8
[0.2 0.13 0.03 0 0 ... 0.07 0.09 0.01]	1.2
[0.13 0.3 0.3 0 0.2 ... 0.08 0.87 0.02]	1.4
[0.3 0.4 0.09 0 0.3 ... 0.05 0.34 0.14]	1.6

idf	embedding	idf
	[0.4 0.3 0.45 0 0.3 ... 0.09 0.24 0.7 0.01]	1.2
	[0.2 0.1 0.03 0 0.1 ... 0.18 0.91 0.2 0.02]	1.3
	[0.1 0.3 0.13 0 0.1 ... 0.16 0.25 0.8 0.03]	2.5
	[0.3 0.5 0.11 0 0.2 ... 0.03 0.17 0.1 0.15]	5.2
	[0.7 0.9 0.01 0 0.6 ... 0.15 0.35 0.4 0.26]	1.3
	...	...

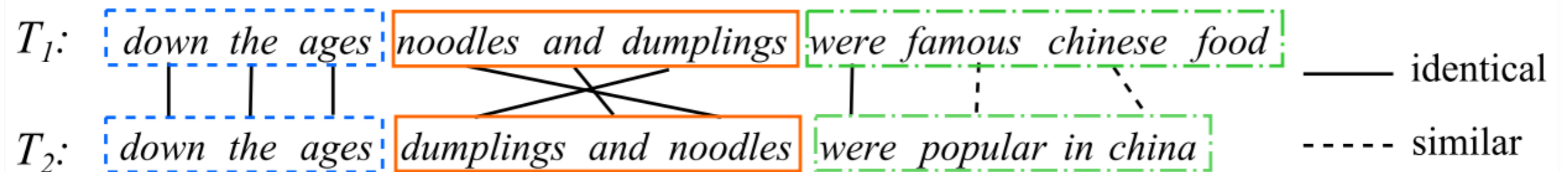
peace process in the Middle East

As for the Arabian and Palestinian voices that are against the current negotiations and the so-called peace process, they are not against peace per se, but rather for their well-founded predictions that Israel would NOT give an inch of the West bank (and most probably the same for Golan Heights) back to the Arabs. An 18 months of "negotiations" in Madrid, and Washington proved these predictions. Now many will jump on me saying why are you blaming israelis for no-result negotiations. I would say why would the Arabs stall the negotiations, what do they have to loose ?



# Outline

- ❖ Semantic matching in search
- ❖ Word-level matching: bridging the semantic gap
- ❖ Sentence-level matching: capturing the proximity
- ❖ Summary and discussion



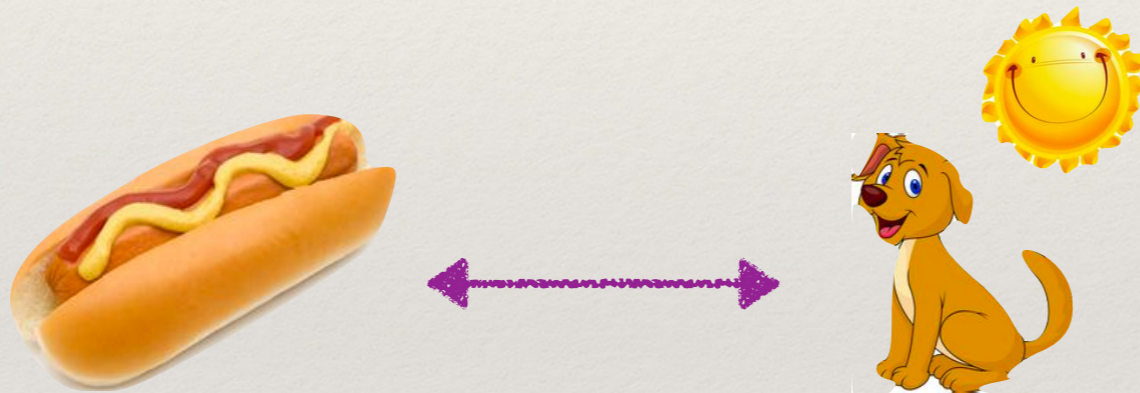
Similarity between “noodles and dumplings”  
and “dumplings and noodles”

---

# Problems with Direct Matching

---

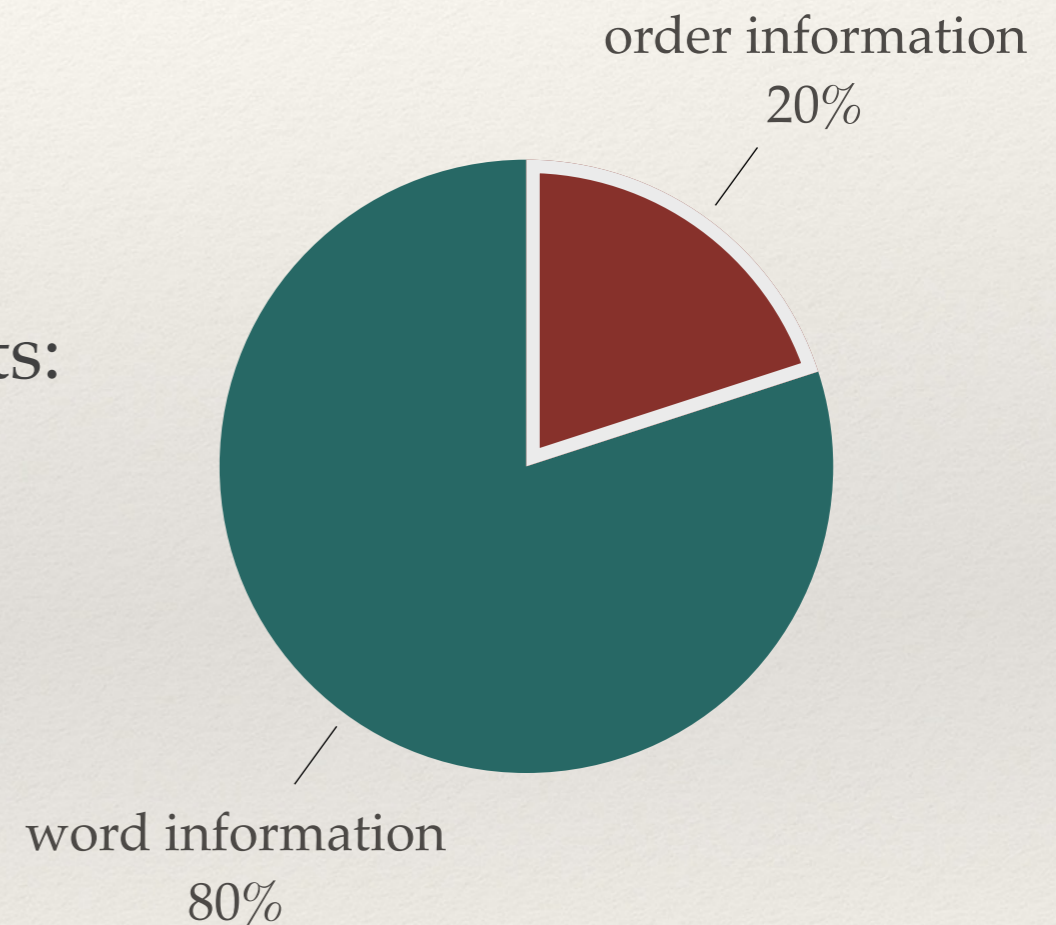
- ❖ *Problem 1: information on the words order is missing*



- ❖ Bag of words: Dog Hot = Hot Dog
- ❖ In real world: Dog Hot  $\neq$  Hot Dog

# The Importance of Words Order

- ❖ Assume:
  - ❖ size of vocabulary = 100,000
  - ❖ average sentences length = 20
- ❖ Rough contributions of information in bits:
  - ❖ From the selection of words:  $\log_2(100000^{20})$
  - ❖ From the order of words:  $\log_2(20!)$
- ❖ Conclusion: over 80% of the potential information in language being in the choice of words without regard to the order in which they appear

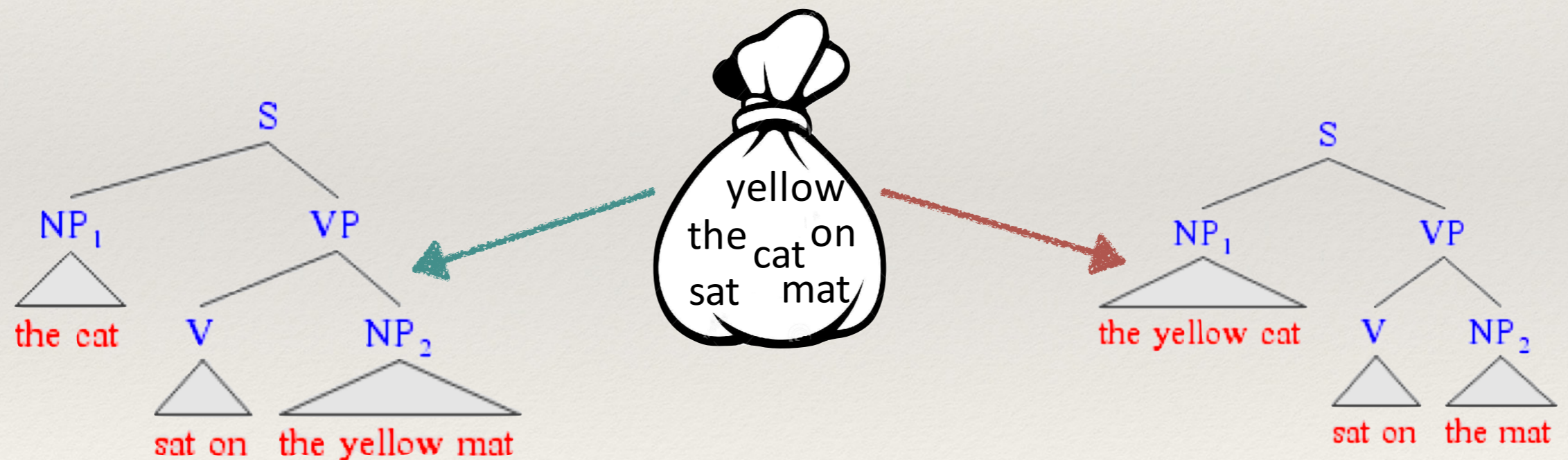


# Problem with Direct Methods

## ❖ *Problem 2: simple sentence representation*

With bag-of-words assumption:

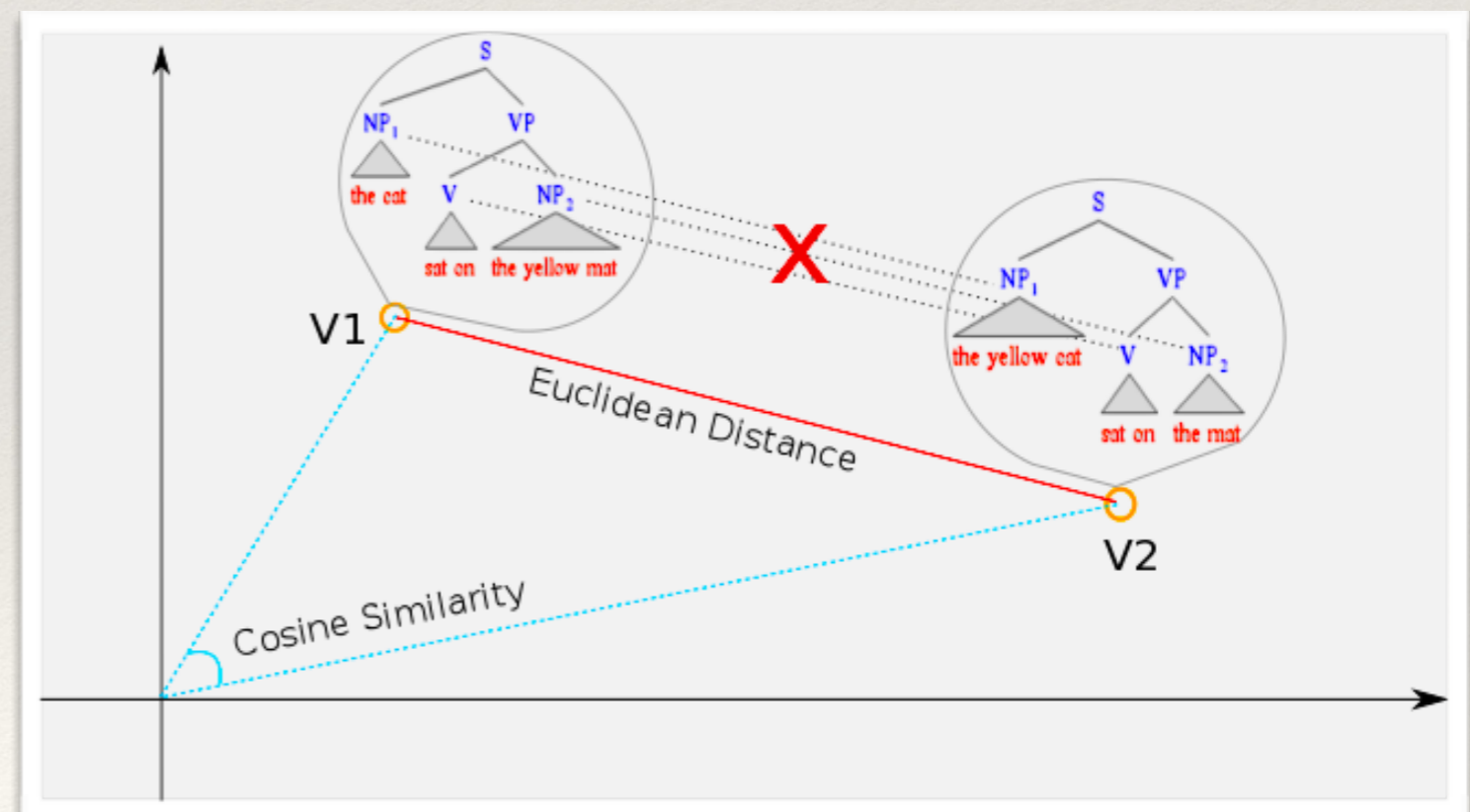
“the **yellow** cat sat on the mat” = “the mat sat on the **yellow** mat”



# Problem with Direct Methods

- ❖ *Problem 3: Heuristic matching function*
  - ❖ A vector for representing the whole sentence
  - ❖ Based on distance measures between two vectors, e.g., Cosine, dot product, Euclidean distance

Limited information is taken into consideration



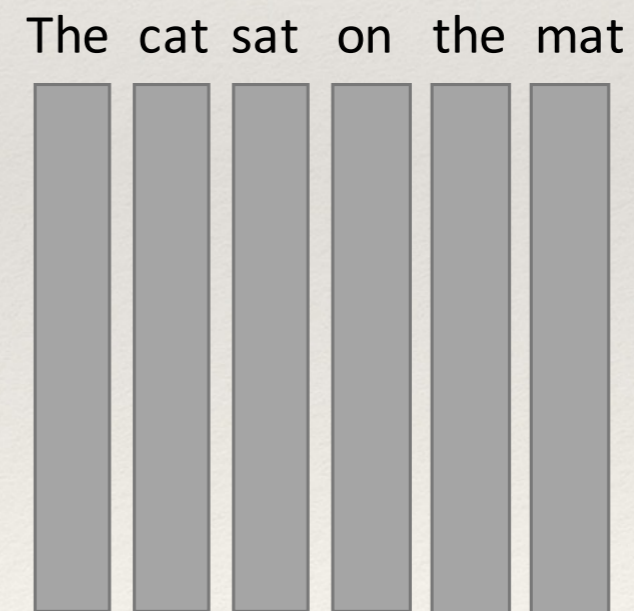
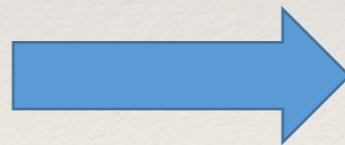
How to design a deep model for  
semantic text matching?

# Keeping Order Information

- ❖ Sequence of word embeddings as the inputs
  - ❖ Convert words to embeddings (e.g. word2vec)
  - ❖ Concatenate embeddings to a sequence



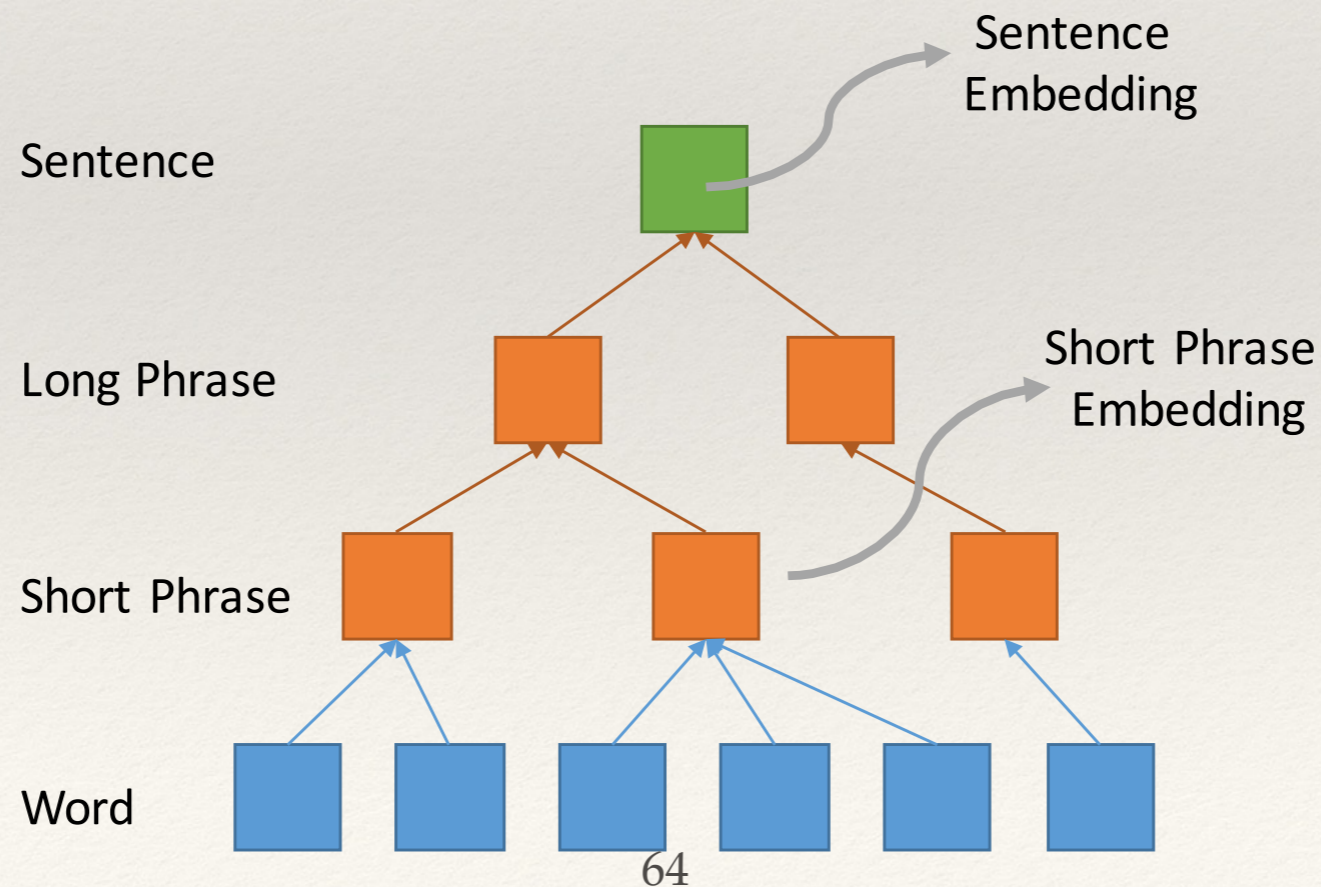
Bag of Word Embeddings



Sequence of Word Embeddings

# Rich Sentence Representation

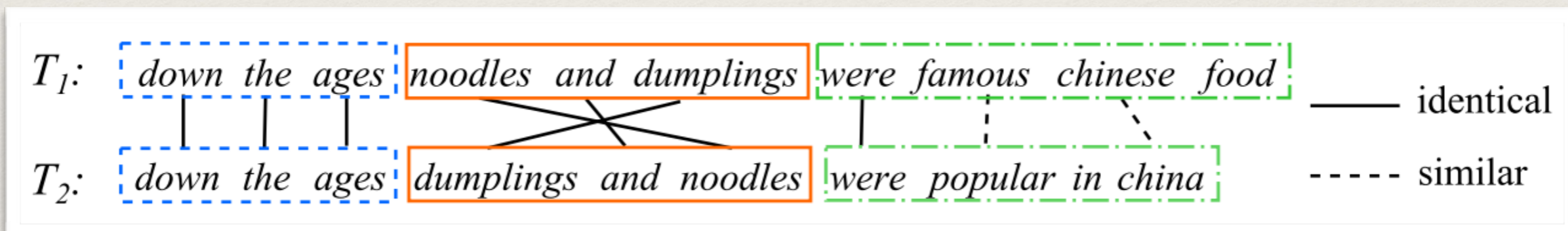
- ❖ Hierarchical structure of sentence representation
  - ❖ Different levels of embeddings
  - ❖ Involving sentence structure





# Powerful Matching Function

- ❖ Considering different levels / types of matching signals



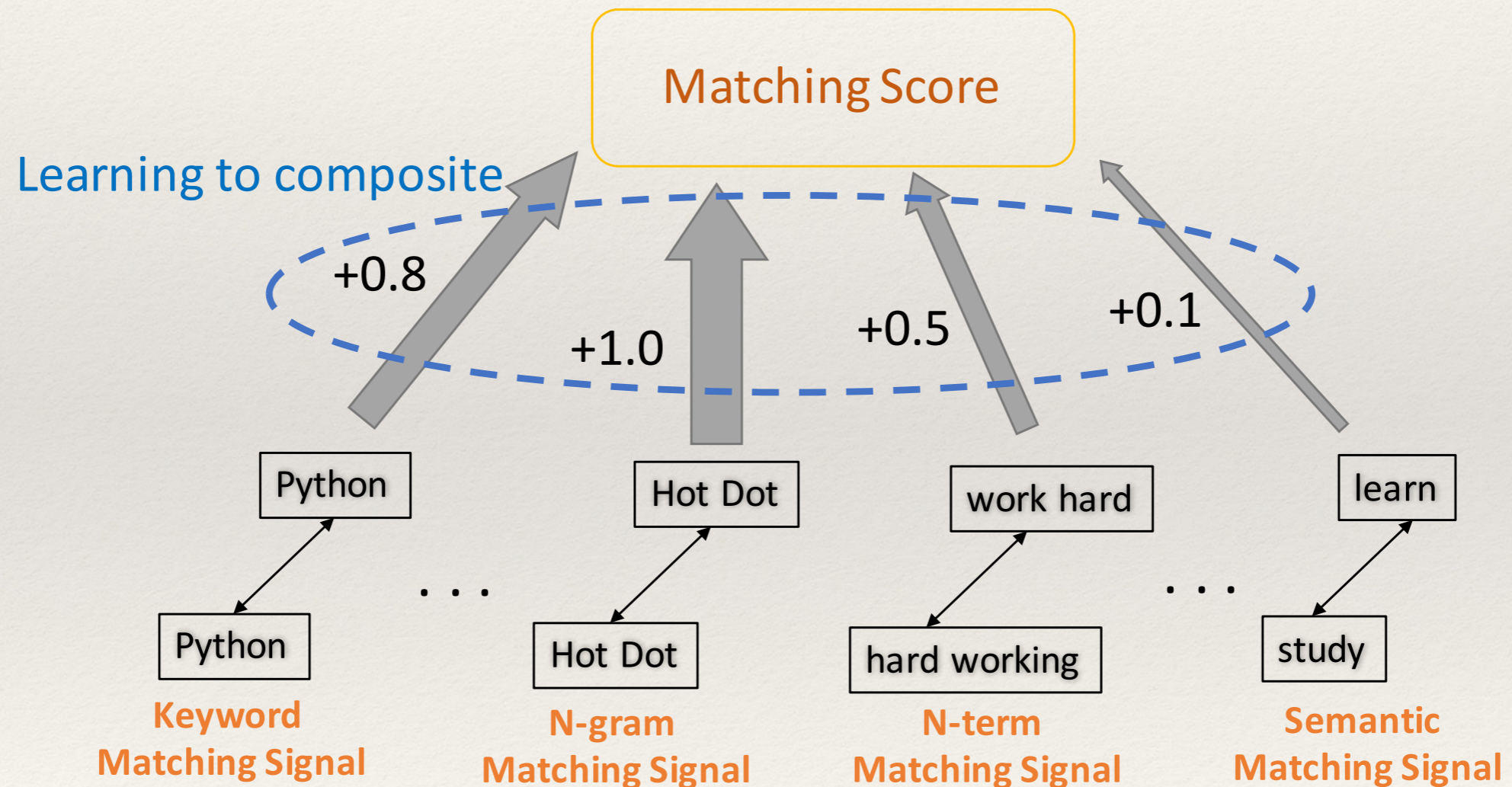
N-gram

N-term

Proximal N-term

# Learning the Matching Function

- ❖ Data-driven approaches for determining the parameters



# Existing Deep Matching Models for Semantic Text Matching

---

# Existing Deep Text Matching Models

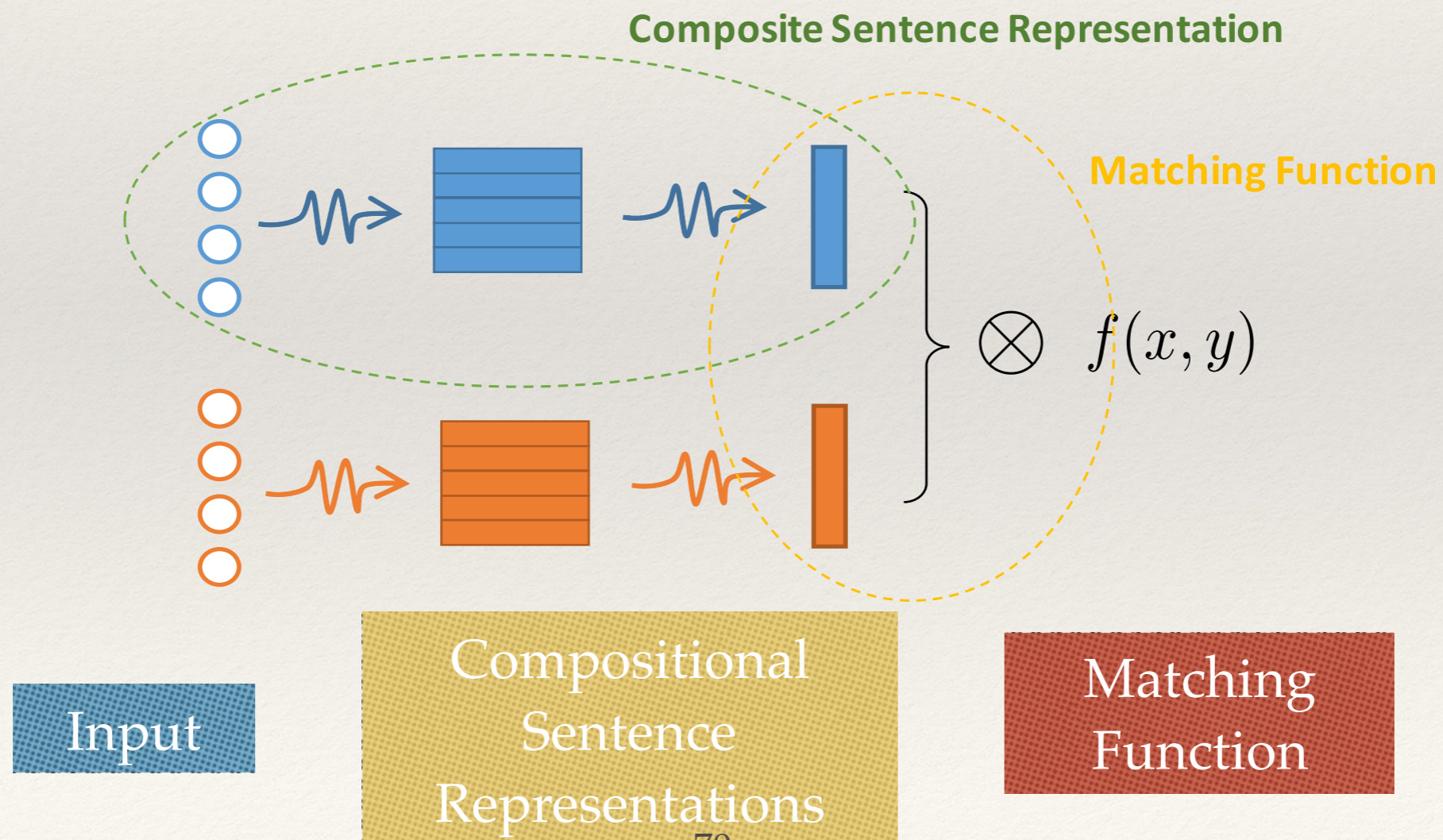
---

- ❖ **Composition Focused Methods [Problem 1] [Problem 2]**
  - ❖ Composite each sentence into one embedding
  - ❖ Measure the similarity between the two embeddings
- ❖ **Interaction Focused Methods [Problem 1] [Problem 3]**
  - ❖ Two sentences meet before their own high-level representations mature
  - ❖ Capture complex matching patterns

# Composition Focused Methods

# Composition Focused Methods

- ❖ Step 1: Compositional sentence representation  $\phi(x)$
- ❖ Step 2: Matching function  $F(\phi(x), \phi(y))$



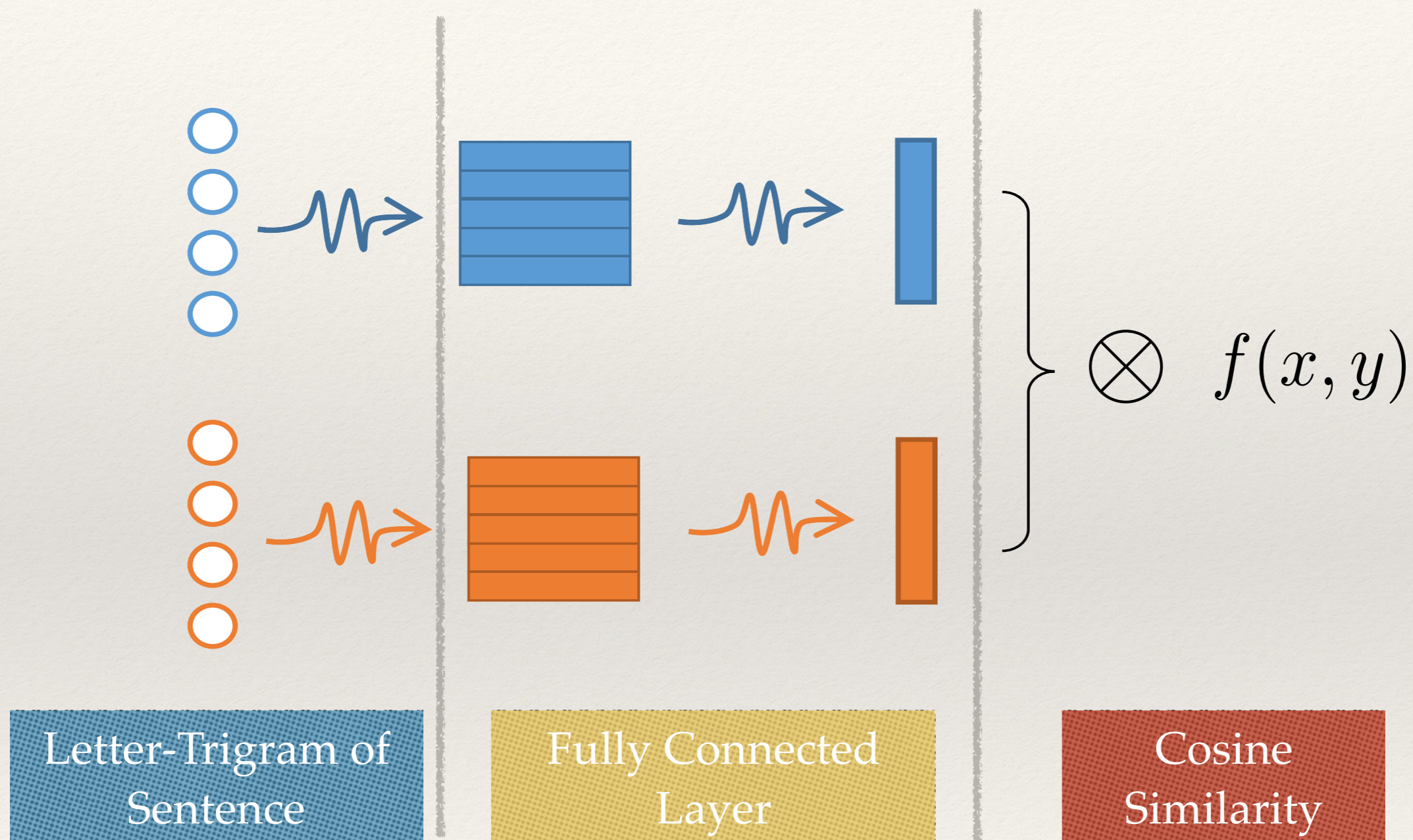
---

# Typical Composition Focused Deep Matching Models

---

- ❖ Based on DNN
  - ❖ **DSSM**: Learning Deep Structured Semantic Models for Web Search using Click-through Data (Huang et al., CIKM'13)
- ❖ Based on CNN
  - ❖ **CDSSM**: A latent semantic model with convolutional-pooling structure for information retrieval (Shen et al. CIKM'14)
  - ❖ **ARC I**: Convolutional Neural Network Architectures for Matching Natural Language Sentences (Hu et al., NIPS'14)
  - ❖ **CNTN**: Convolutional Neural Tensor Network Architecture for Community-Based Question Answering (Qiu and Huang., IJCAI'15)
- ❖ Based on RNN
  - ❖ **LSTM-RNN**: Deep Sentence Embedding Using the Long Short Term Memory Network: Analysis and Application to Information Retrieval (Palangi et al., TASLP'2016)

# Deep Structured Semantic Model (DSSM)





---

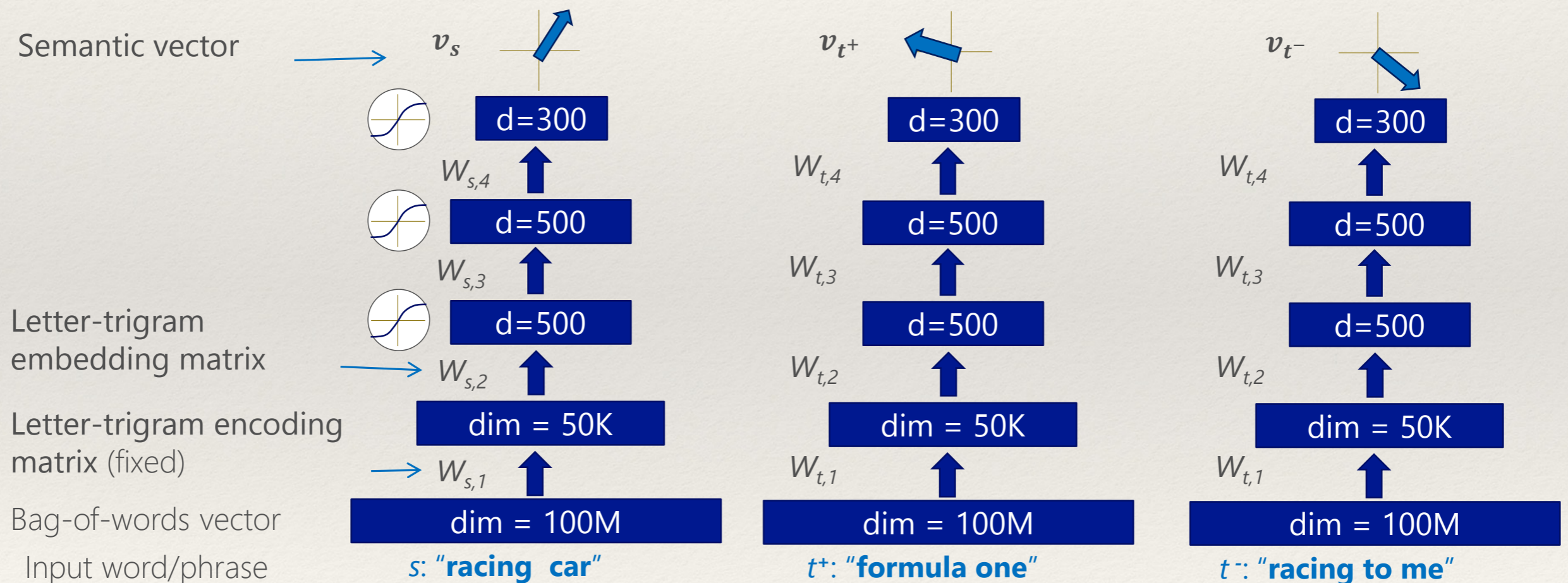
# DSSM Inputs: Letter-trigram

---

- ❖ Bag of words representation
  - ❖ “candy store”: [ 0 0 0 1 0 0 0 1 0 0 0 ... ]
- ❖ Letter-trigram representation
  - ❖ “#candy# #store#”  $\Rightarrow$  #ca | can | and | ndy | dy# | #st | sto | tor | ore | re#
  - ❖ [ 0 0 1 0 0 ... 0 1 0 1 ... 0 0 ... ]
- ❖ Advantages:
  - ❖ Compact representation: # words: 500K  $\Rightarrow$  # letter-trigram: 30K
  - ❖ Generalize to unseen words
  - ❖ Robust to misspelling, inflection, etc.

# DSSM Sentence Representation: DNN

- ❖ Model: DNN (auto-encoder) to capture the compositional sentence representations



---

# DSSM Matching Function

---

- ❖ Cosine similarity between semantic vectors

$$S = \frac{x^T \cdot y}{|x| \cdot |y|}$$

- ❖ Training

- ❖ A query  $q$  and a list of docs  $D = \{d^+, d_1^-, \dots, d_k^-\}$

- ❖  $d^+$  positive doc,  $d_1^-, \dots, d_k^-$  negative docs to query

- ❖ Objective:  $P(d^+ | q) = \frac{\exp(\gamma \cos(q, d^+))}{\sum_{d \in D} \exp(\gamma \cos(q, d))}$

- ❖ Optimizing with SGD

---

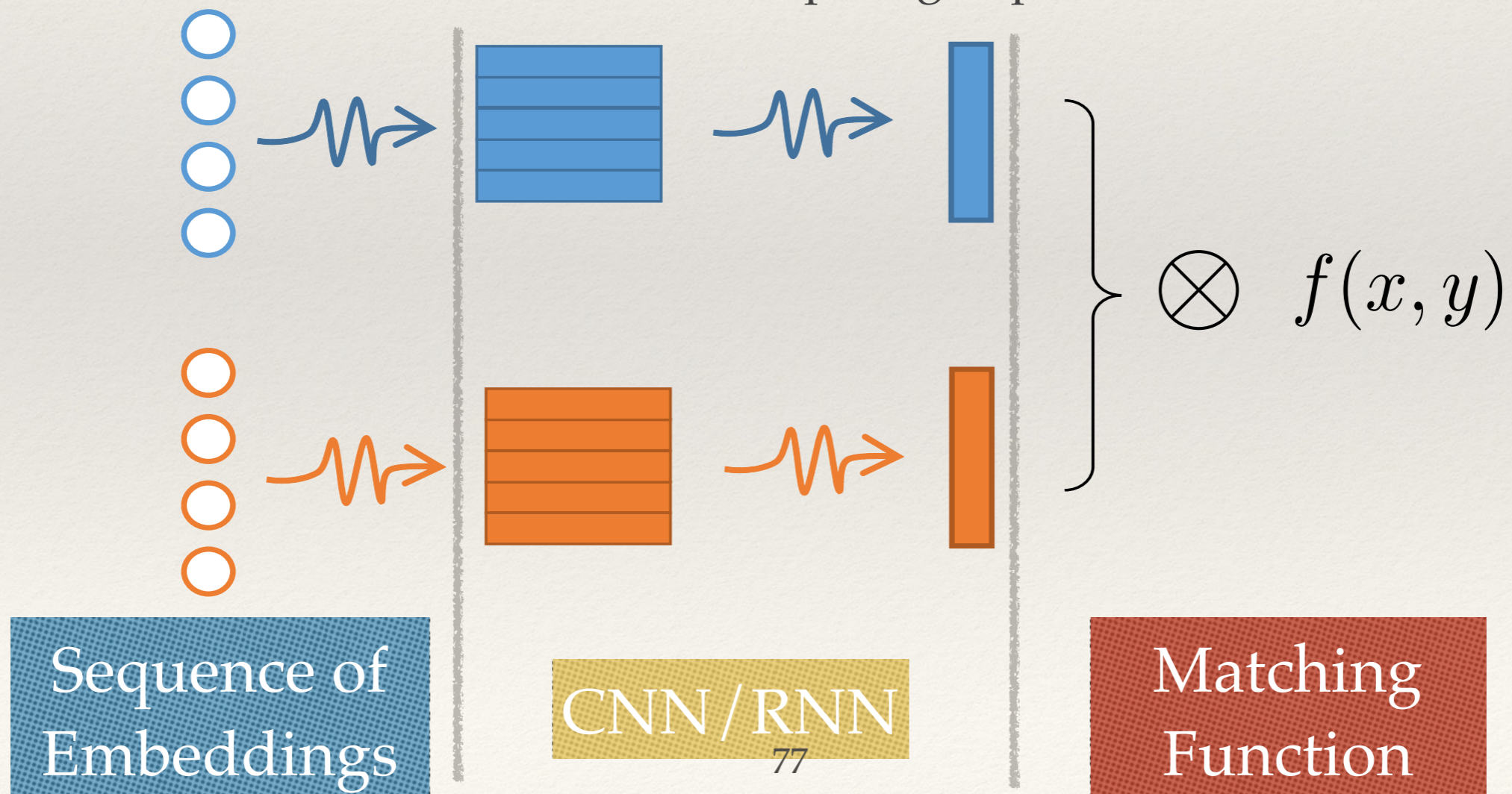
# DSSM: Brief Summary

---

- ❖ Inputs: sub-word units (i.e. letter-trigram) as input for scalability and generalizability
- ❖ Representations: mapping sentences to vectors (i.e. DNN): semantically similar sentences close to each other
- ❖ Matching: cosine similarity as the matching function
- ❖ Problem: bag of letter-trigrams, **the order information of words is missing**

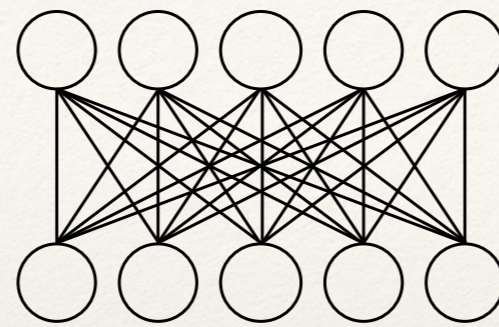
# Capturing Order Information?

- ❖ Input: word sequence instead of bag of letter-trigrams
- ❖ Model:
  - ❖ Convolutional based methods can keep locally order
  - ❖ Recurrent based methods can keep long dependence relations



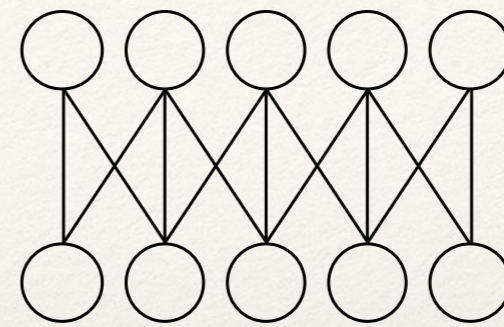
# CNN can Model the Order Information

- ❖ Inspired by the cat's visual cortex [Hubel '68]



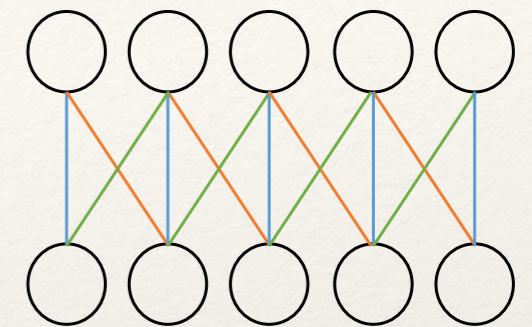
Fully Connected Layer

All different weights



Locally Connected Layer

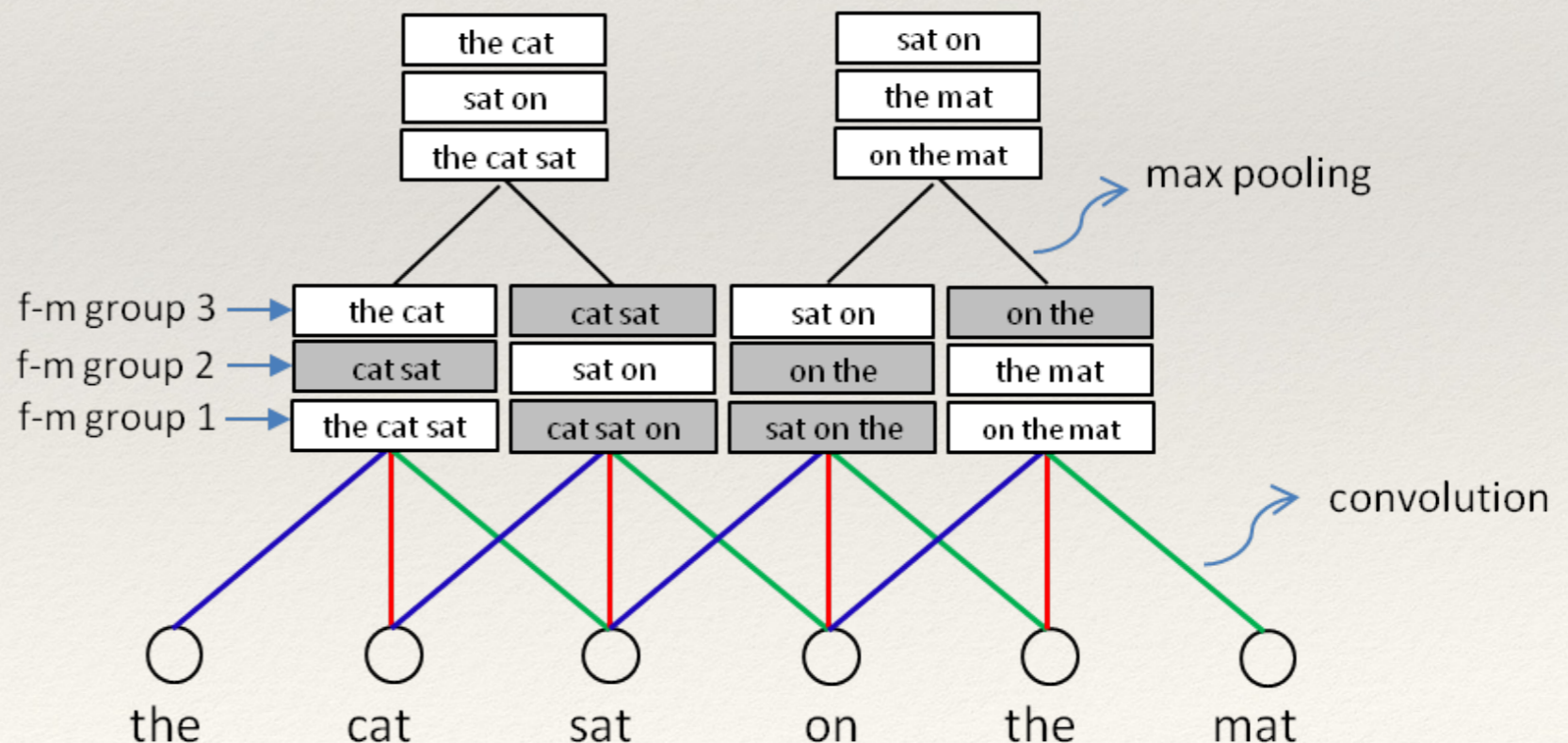
All different weights



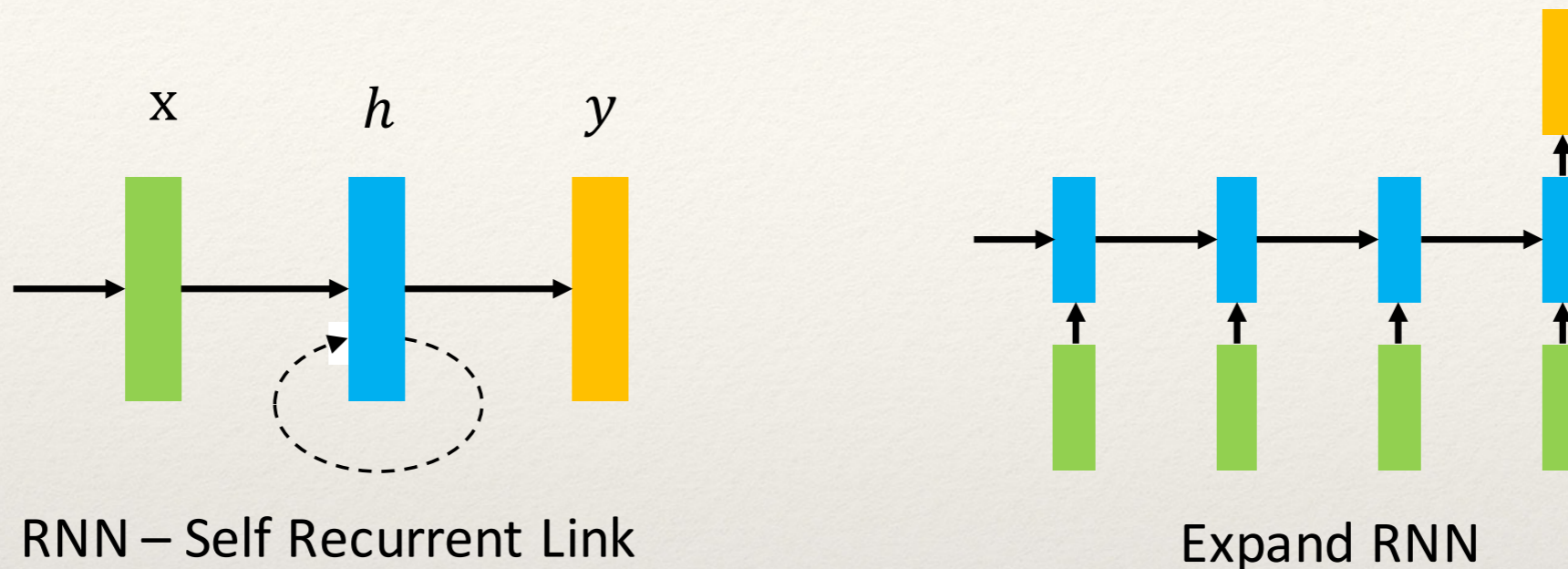
Convolutional Layer

Shared weights

- ❖ Convolution & max pooling operations on text



# RNN can Model the Order Information



- ❖ RNNs implement dynamical systems
- ❖ RNNs can approximate arbitrary dynamical systems with arbitrary precision
- ❖ Training: Back Propagation Through Time

$$s(t) = f(\mathbf{U}w(t) + \mathbf{W}s(t-1) + b)$$

- ❖ Two popular variations: long-short term memory (LSTM) and gated recurrent unit (GRU)

# Using CNN: CDSSM

- ❖ Input: encode **each word** as bag of letter-trigram
- ❖ Model: the convolutional operation in CNN compacts each **sequence of k words**

Semantic layer:  $y$

Affine projection matrix:  $W_s$

Max pooling layer:  $v$

Max pooling operation

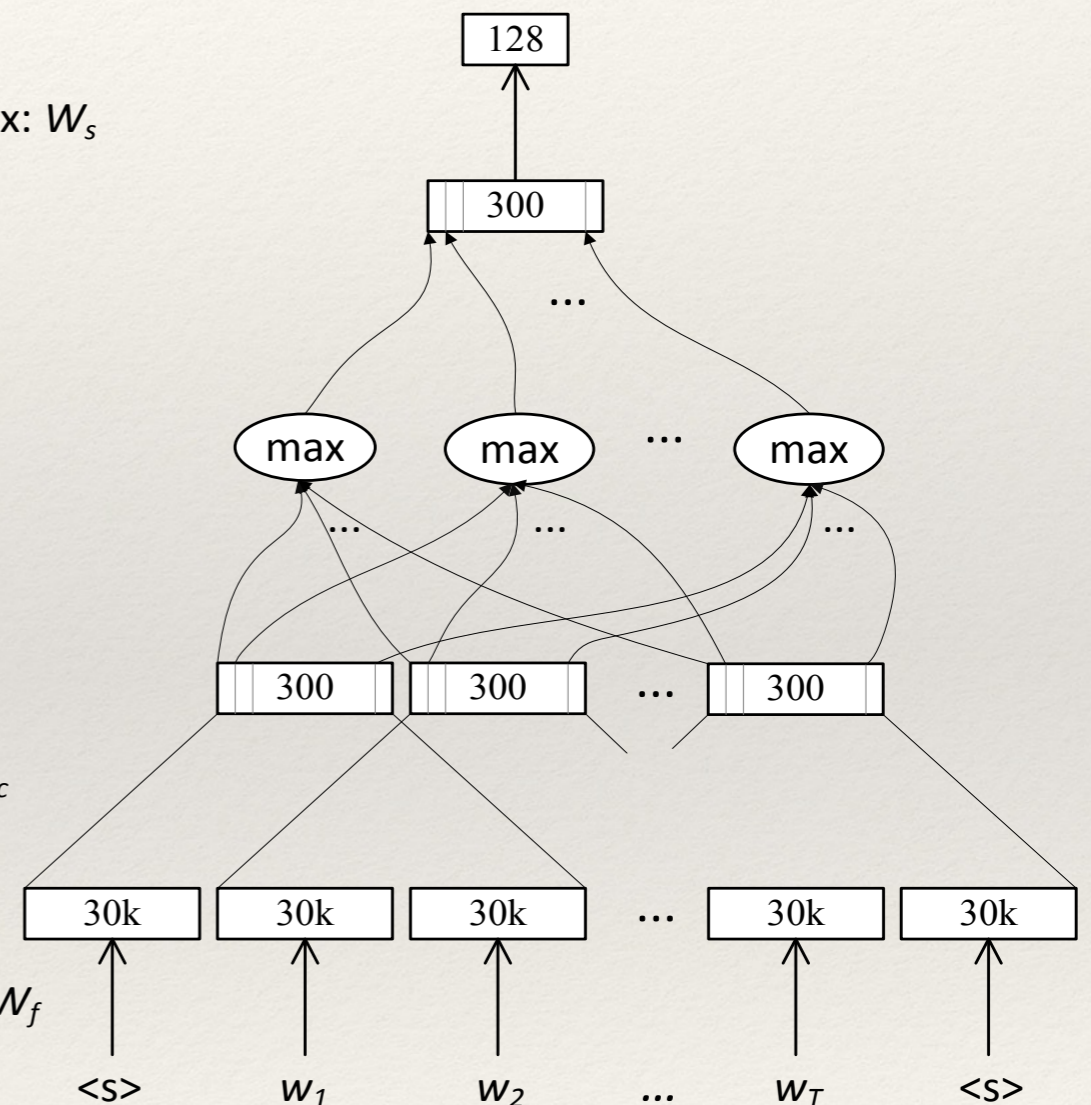
Convolutional layer:  $h_t$

Convolution matrix:  $W_c$

Word hashing layer:  $f_t$

Word hashing matrix:  $W_f$

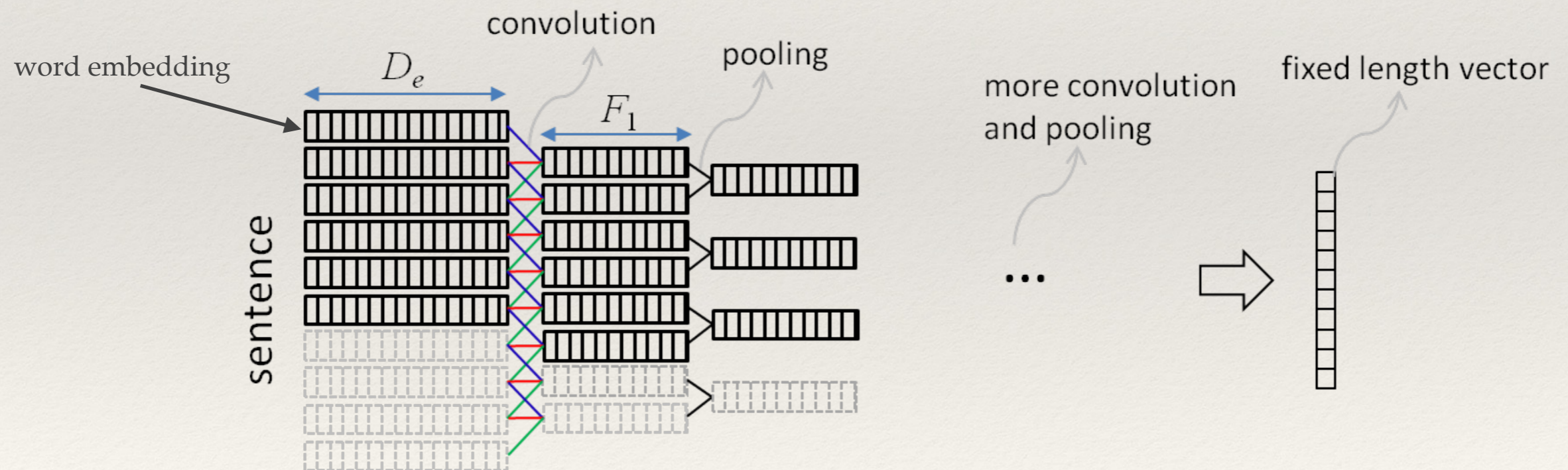
Word sequence:  $x_t$





# Using CNN: ARC-I / CNTN

- ❖ Input: sequence of word embeddings
  - ❖ Word embeddings from word2vec model train on large dataset
- ❖ Model: the convolutional operation in CNN compacts each **sequence of k words**

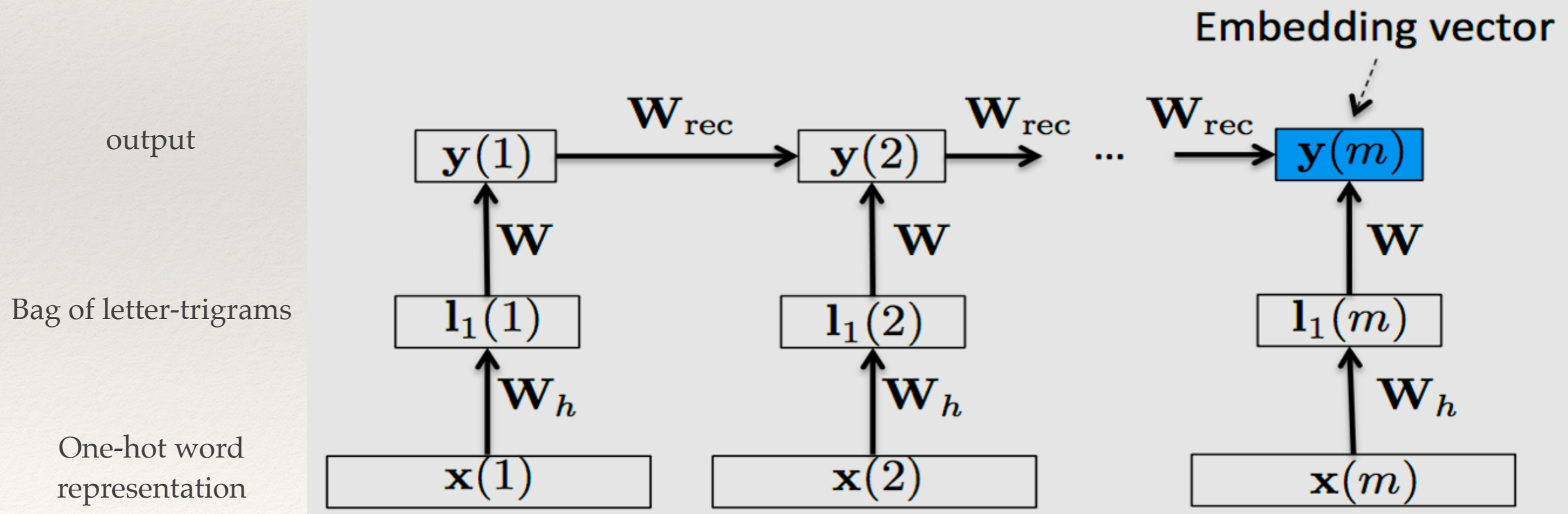


Baotian Hu, Zhengdong Lu, Hang Li, Qingcai Chen. Convolutional Neural Network Architectures for Matching Natural Language Sentences. Proceedings of Advances in Neural Information Processing Systems 27 (NIPS'14), 2042-2050, 2014.

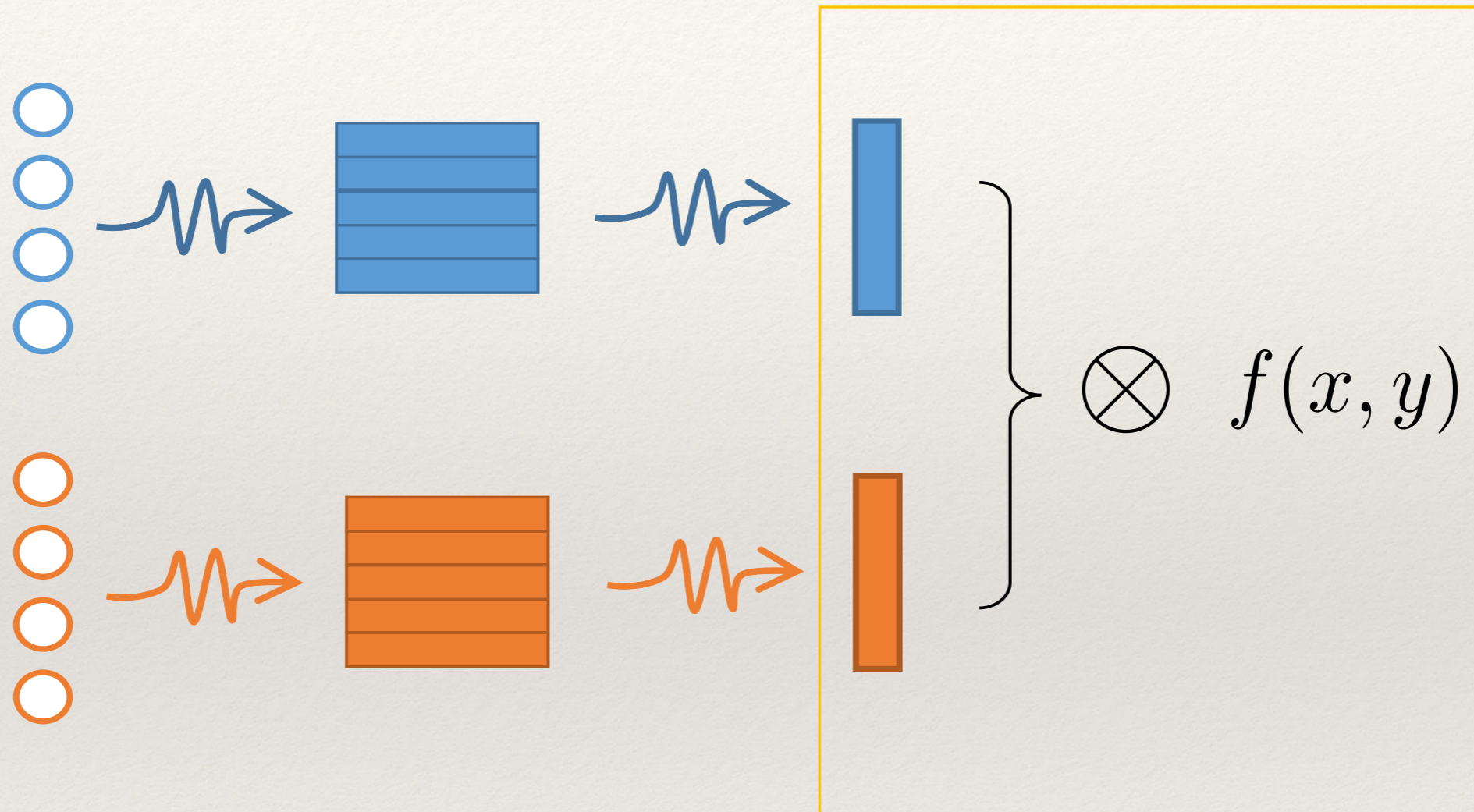
Qiu X, Huang X. Convolutional neural tensor network architecture for community-based question answering//Proceedings of the 24th (IJCAI), Buenos Aires, Argentina, 2015: 1305-1311.

# Using RNN: LSTM-RNN

- ❖ Input: sequence letter trigrams
- ❖ Model: Long-short term memory (LSTM)
  - ❖ The last output as the sentence representation



# Matching Function



Heuristic: cosine, dot product

Learning: MLP, Neural tensor networks

---

# Matching Functions (cont')

---

- ❖ Given the representations of two sentences:  $x$  and  $y$ .
- ❖ Similarity between these two embeddings:

- ❖ Cosine Similarity (DSSM, CDSSM, RNN-LSTM)

$$S = \frac{x^T \cdot y}{|x| \cdot |y|}$$

- ❖ Dot Product

$$S = x^T \cdot y$$

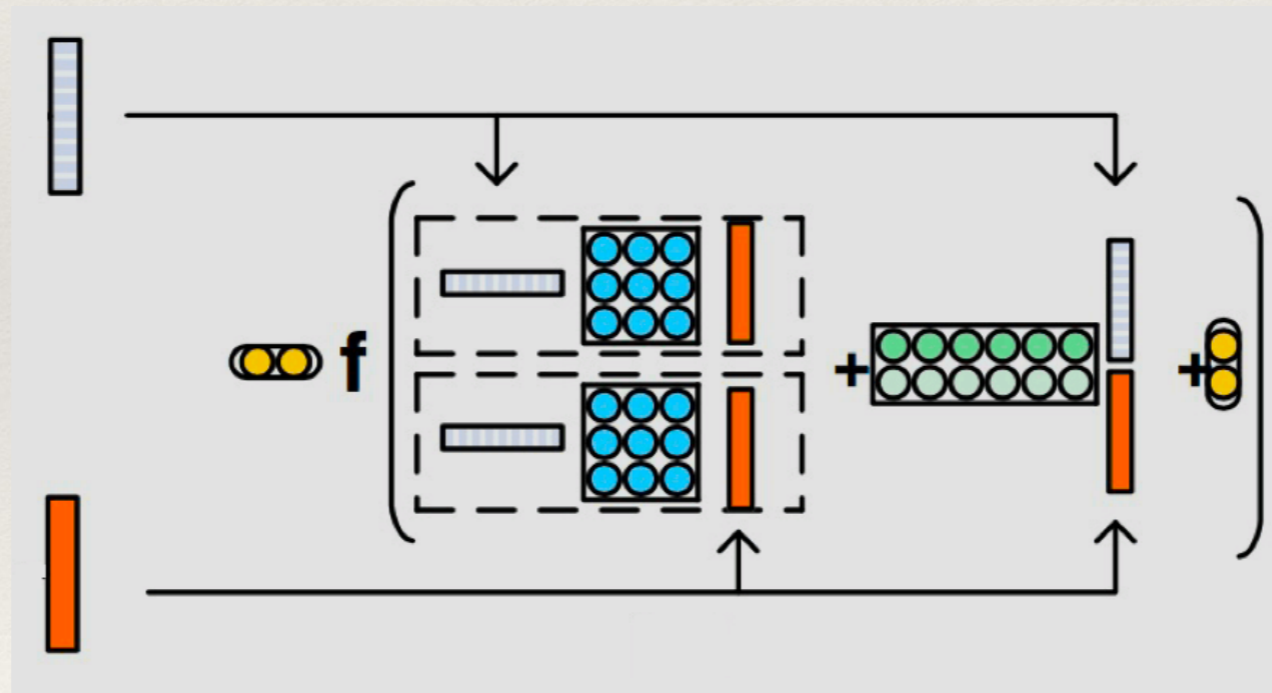
- ❖ Multi-Layer Perception (ARC-I)

$$S = W_2 \cdot \left( W_1 \cdot \begin{bmatrix} x \\ y \end{bmatrix} + b_1 \right) + b_2$$

# Matching Functions (cont')

## ❖ Neural Tensor Networks (CNTN)

$$S = u^T f \left( x^T \mathbf{M}^{[1:r]} y + V \begin{bmatrix} x \\ y \end{bmatrix} + b \right)$$



Qiu X, Huang X. Convolutional neural tensor network architecture for community-based question answering//Proceedings of the 24th (IJCAI), Buenos Aires, Argentina, 2015: 1305-1311.

---

# Performance Evaluation on QA Task

---

- ❖ Dataset: Yahoo! Answers
  - ❖ 60,564 (question, answer) pairs



- ❖ Example:
  - ❖ *Q: How to get rid of memory stick error of my sony cyber shot?*
  - ❖ *A: You might want to try to format the memory stick but what is the error message you are receiving.*

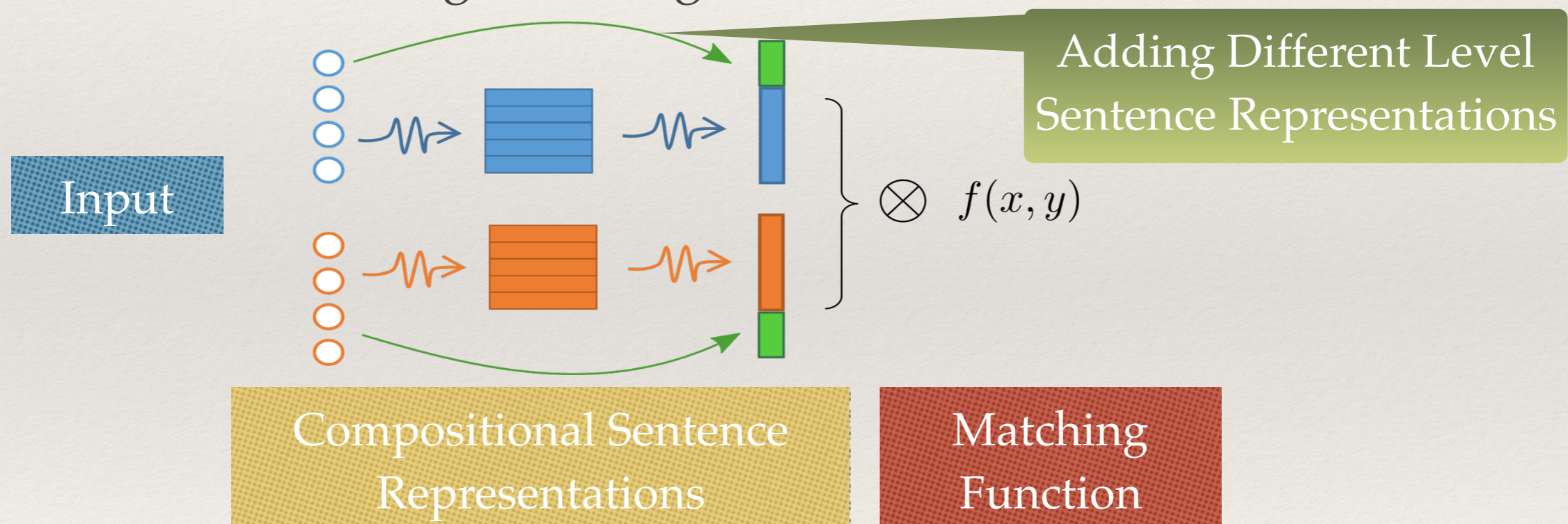
# Experimental Results

	Model	P@1	MRR
Random	Random	0.200	0.457
Traditional	BM25	0.579	0.726
Comosition Focused	ARC-I	0.581	0.756
	CNTN	0.626	0.781
	LSTM-RNN	0.690	0.822

- ❖ Composition focused methods outperformed the baselines
- ❖ Semantic representation is important
- ❖ LSTM-RNN is the best performed method
- ❖ Modeling the order information does help

# Extensions to Composition Focused Methods

- ❖ Problem: sentence representations are too coarse to conduct text match
  - ❖ Experience in IR: **combining topic level** and word level matching signals usually achieve better performances Adding more fine-grained matching signals
- ❖ Solution: add fine-grained signals



- **MultiGranCNN**: An Architecture for General Matching of Text Chunks on Multiple Levels of Granularity. (Yin W, Schütze T, Hinrich. ACL2015)
- **U-RAE**: Dynamic Pooling and Unfolding Recursive Autoencoders for Paraphrase Detection, (Richard Socher, Eric H. Huang, Jeffrey Pennington, Andrew Y. Ng, Christopher D. Manning, NIPS2011)
- **MV-LSTM**: A Deep Architecture for Semantic Matching with Multiple Positional Sentence Representations. (Shengxian Wan, Yanyan Lan, Jiafeng Guo, Jun Xu, and Xueqi Cheng. AAAI 2016)



# Performance Evaluations on QA Task

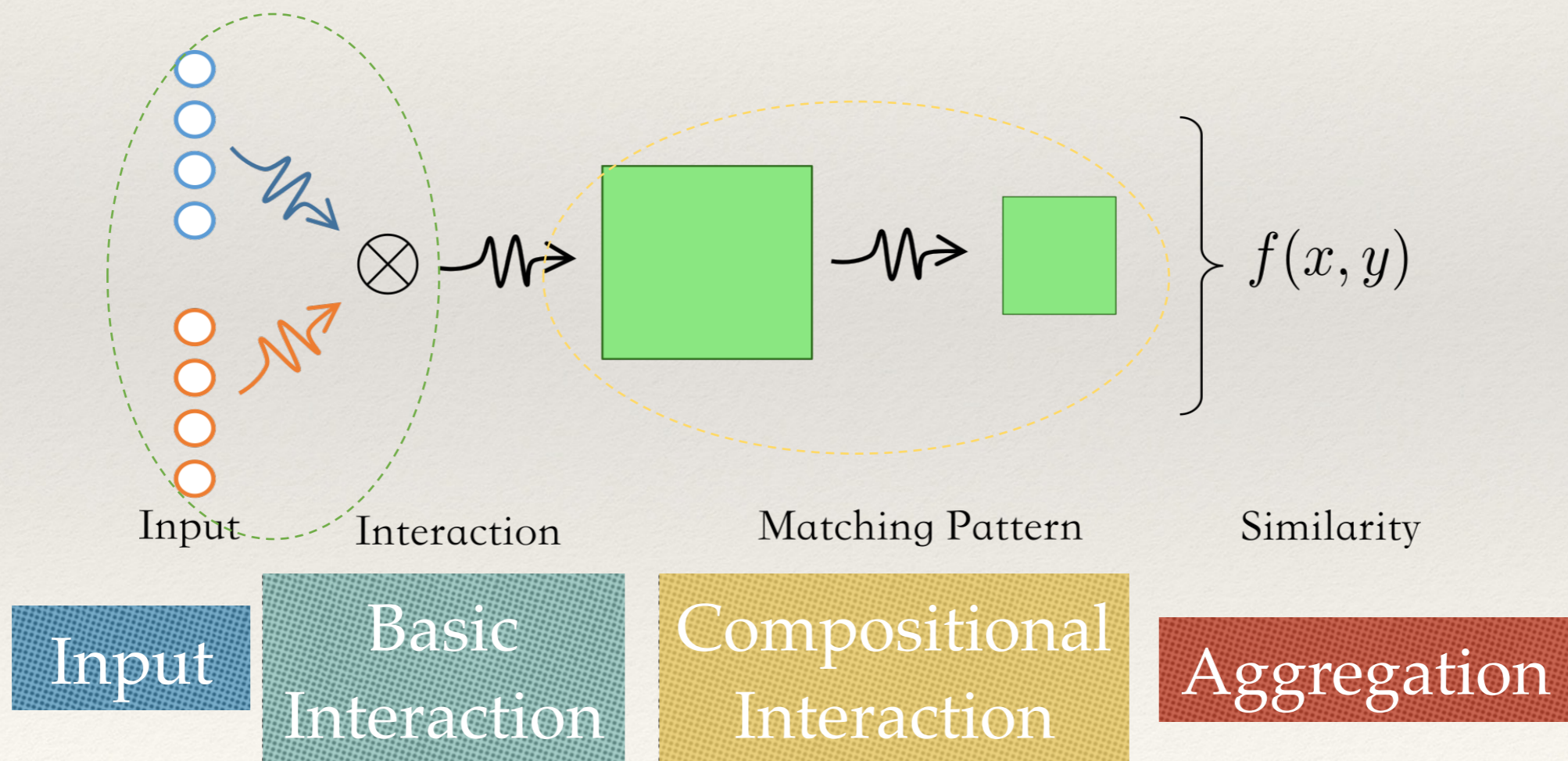
	Model	P@1	MRR
Statistic Traditional	Random	0.200	0.457
	BM25	0.579	0.726
Comosition Focused	ARC-I	0.581	0.756
	CNTN	0.626	0.781
	LSTM-RNN	0.690	0.822
	uRAE	0.398	0.652
	MultiGranCNN	0.725	0.840
	MV-LSTM	0.766	0.869

- ❖ MultiGranCNN and MV-LSTM achieved the best performances
  - ❖ Fine-grained matching signals are useful

# Interaction Focused Methods

# Interaction Focused Methods

- ❖ Step 1: Construct basic low-level interaction signals
- ❖ Step 2: Aggregate matching patterns



---

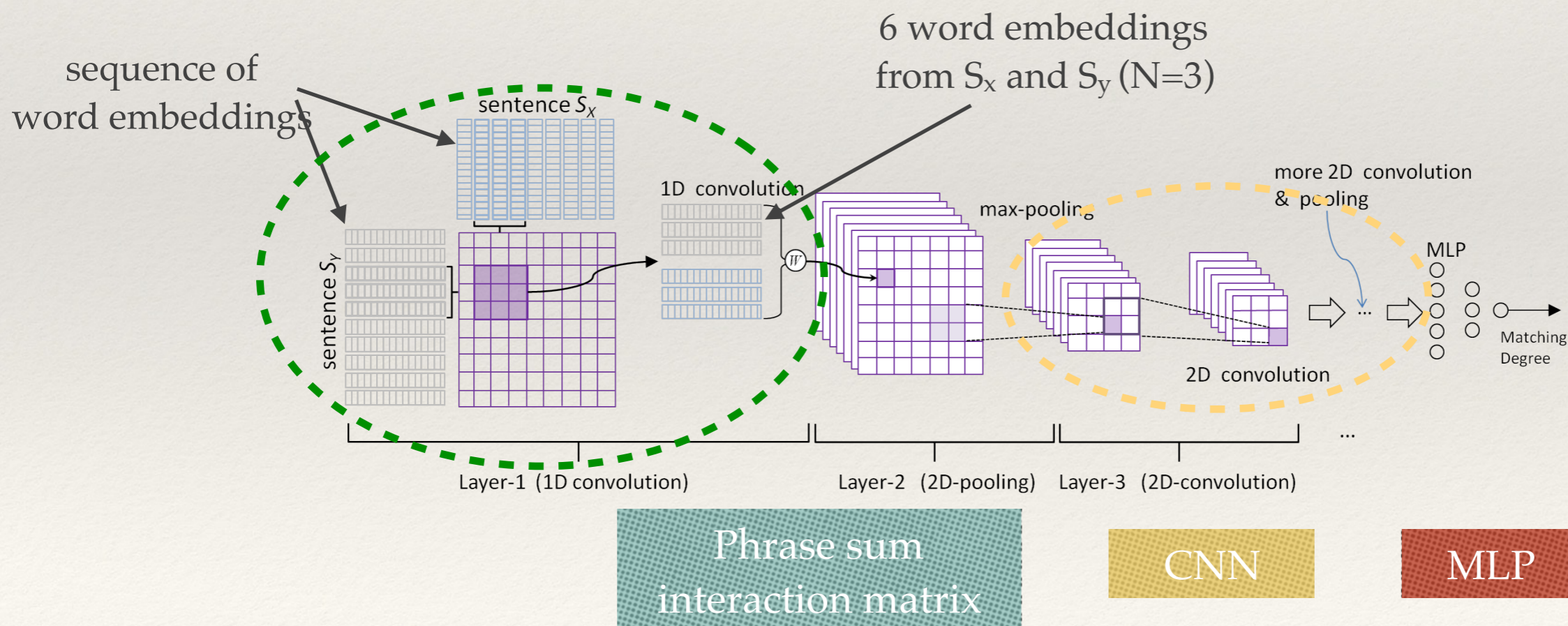
# Typical Interaction Focused Methods

---

- ❖ **ARC II**: Convolutional Neural Network Architectures for Matching Natural Language Sentences (Hu et al., NIPS'14)
- ❖ **MatchPyramid**: Text Matching as Image Recognition. (Pang et al. AAAI'16)
- ❖ **Match-SRNN**: Modeling the Recursive Matching Structure with Spatial RNN. (Wan et al. IJCAI'16)

# ARC-II

- ❖ Let two sentences meet before their own high-level representations mature.
- ❖ Basic interaction: phrase sum interaction matrix
- ❖ Compositional interaction: CNN to capture the local interaction structure
- ❖ Aggregation Function: MLP



# ARC-II (cont')

- ❖ Order Preservation

- ❖ Both the convolution and pooling have order preserving property

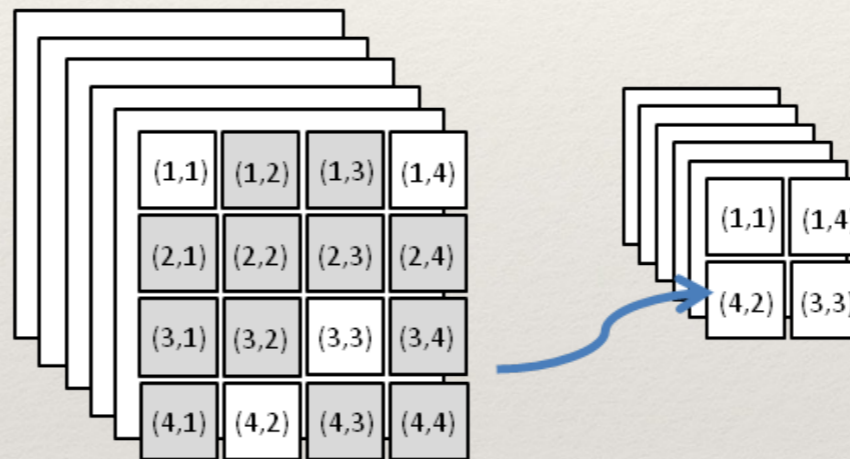


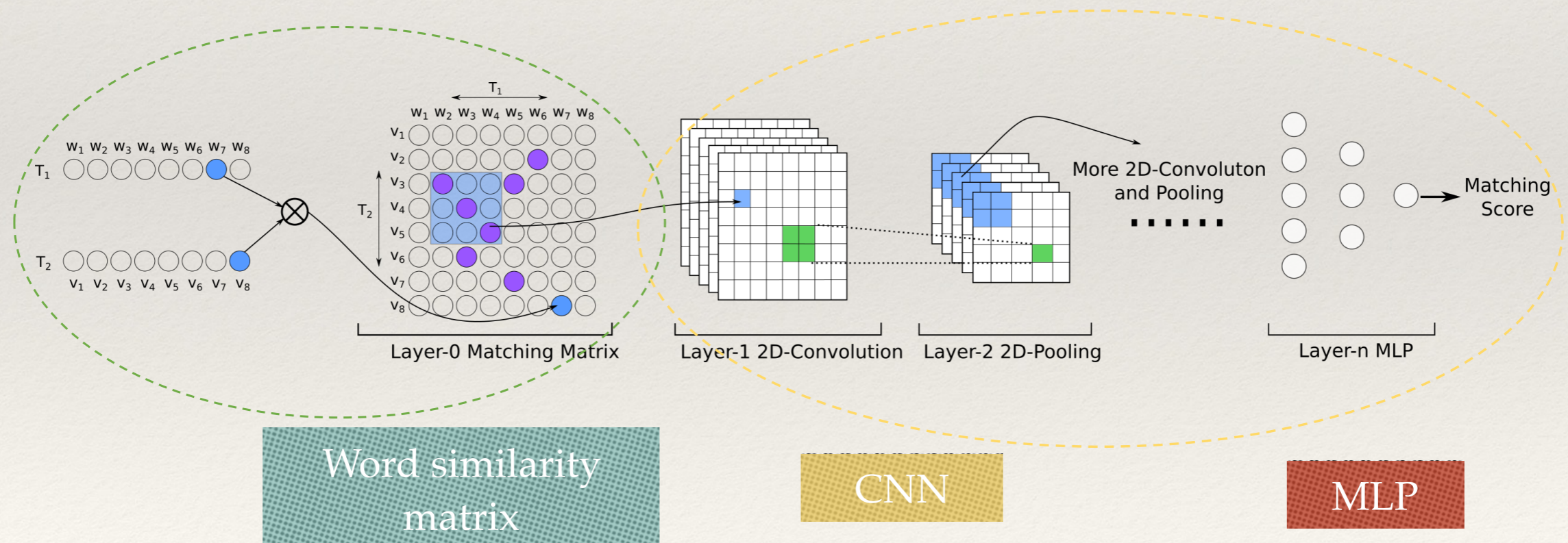
Figure 5: Order preserving in 2D-pooling.

- ❖ However, the **word level matching signals are lost**

- ❖ 2-D matching matrix is construct based on the embedding of the words in two N-grams

# MatchPyramid

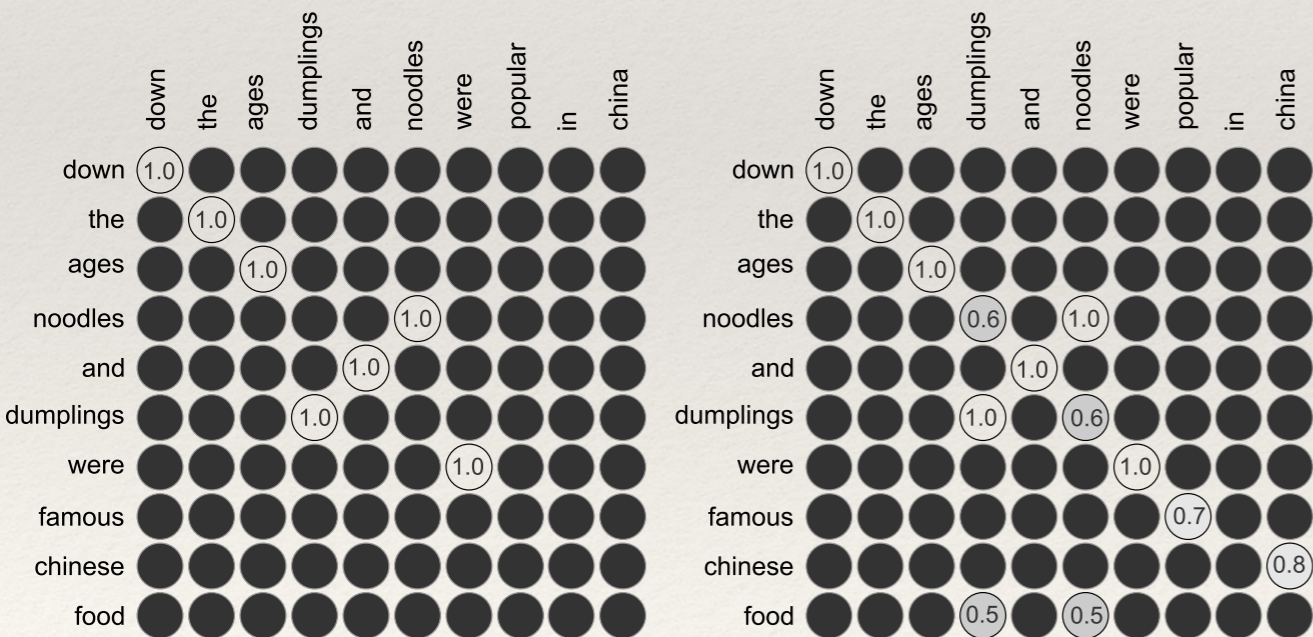
- ❖ Inspired by image recognition task
- ❖ Basic Interaction: word-level matching matrix
- ❖ Compositional interaction: hierarchical convolution
- ❖ Aggregation: MLP



# MatchPyramid: Matching Matrix

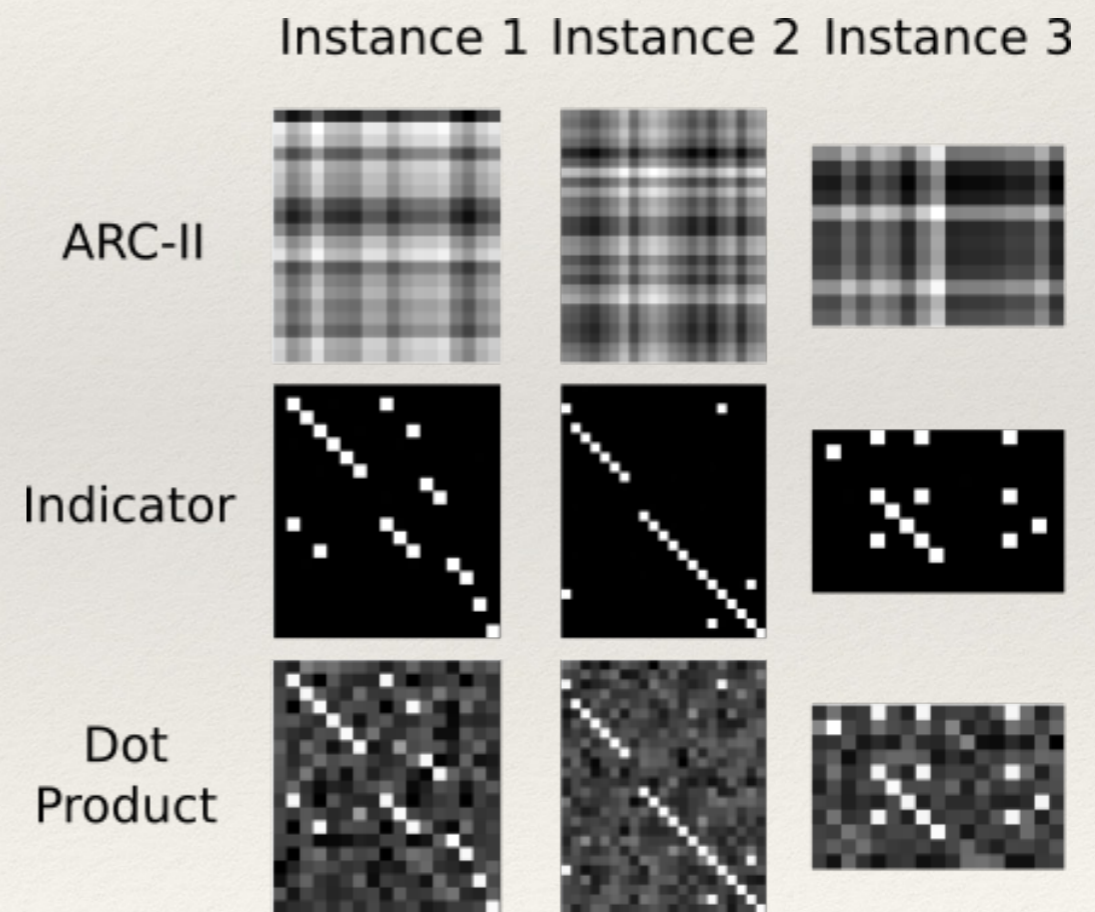
- ❖ Basic Interaction: word similarity matrix
  - ❖ Strength of the word-level matching
  - ❖ Positions of the matching occurs

$$\mathbf{M}_{ij} = w_i \otimes v_j$$



(a) Indicator

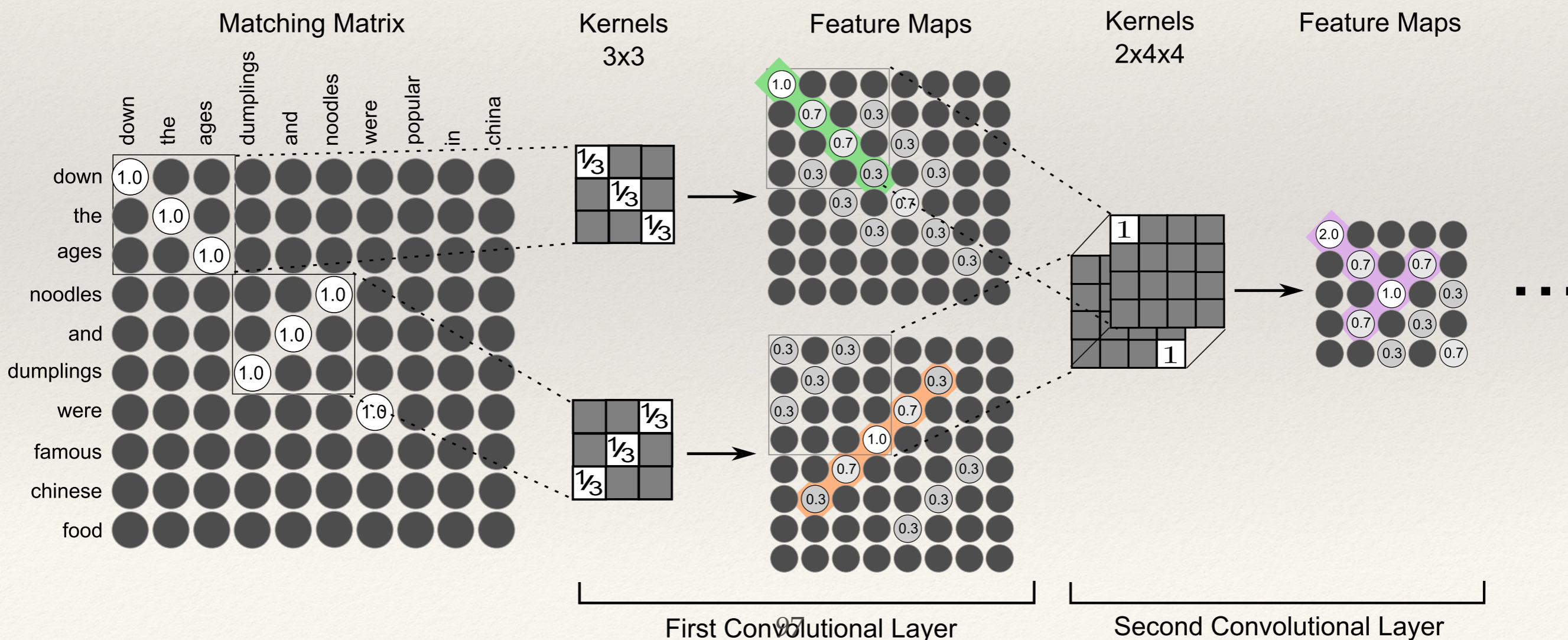
(b) Cosine





# MatchPyramid - Hierarchical Convolution

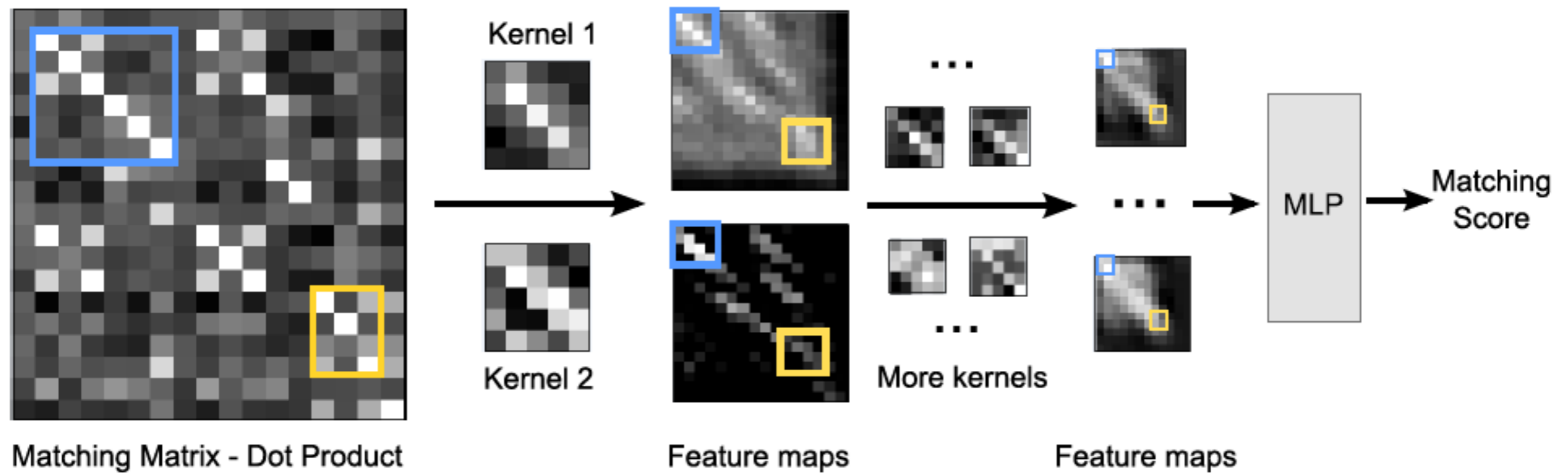
- ❖ Compositional interaction: CNN to capture different levels of matching patterns, based on word-level matching signals



# Matching Patterns Discovered by MathPyramid

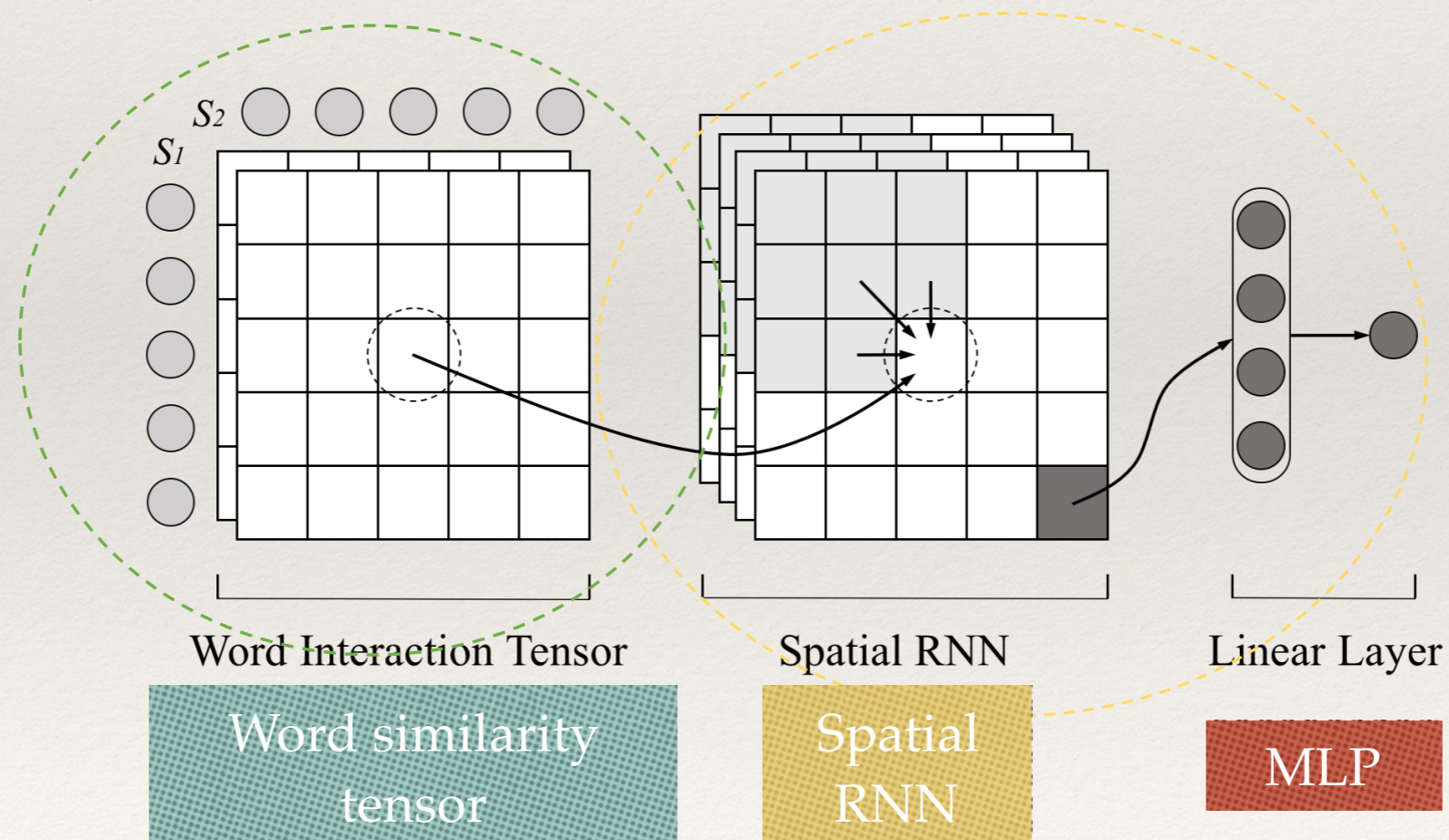
T<sub>1</sub>: PCCW's chief operating officer, Mike Butcher, and Alex Arena, the chief financial officer, will report directly to Mr So.

T<sub>2</sub>: Current Chief Operating Officer Mike Butcher and Group Chief Financial Officer Alex Arena will report to So.

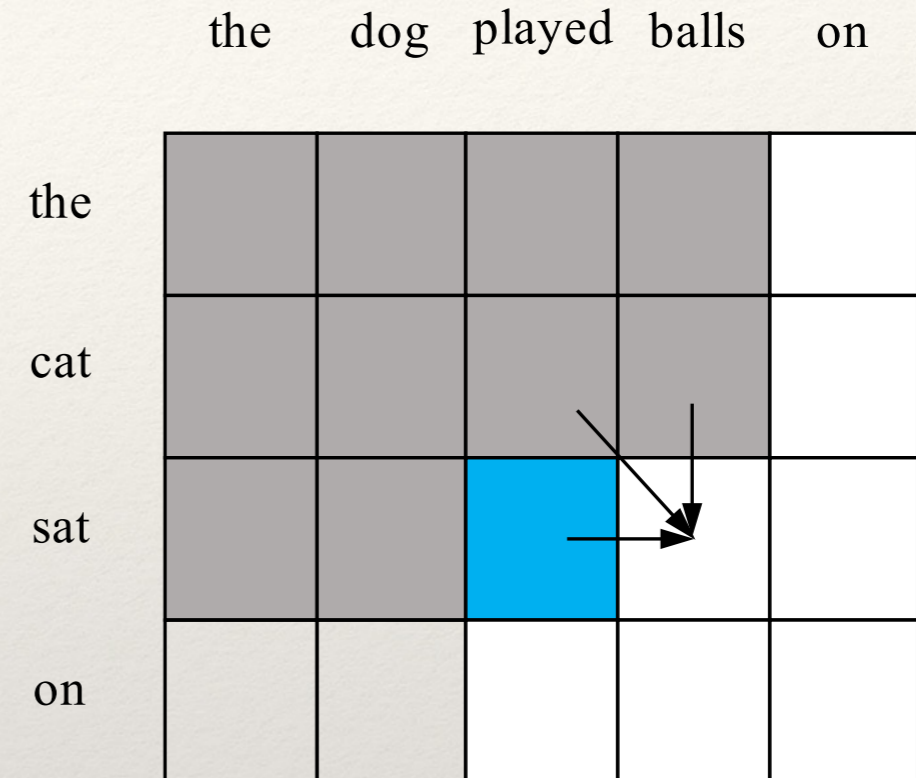
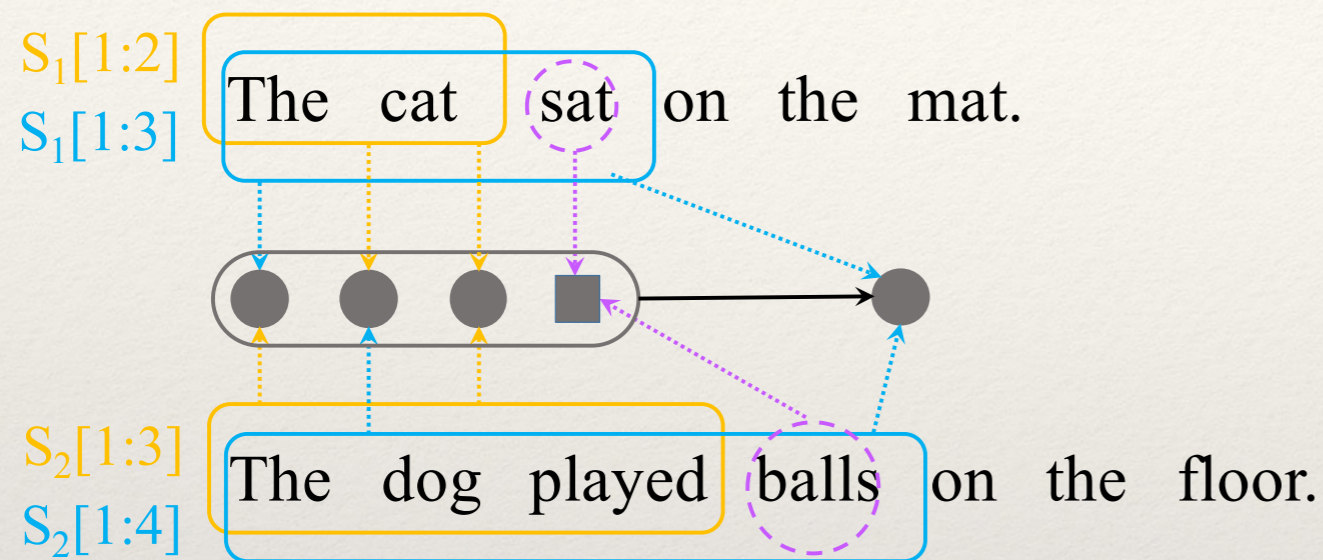


# Match-SRNN

- ❖ Spatial recurrent neural network (SRNN) for text matching
- ❖ Basic interaction: word similarity tensor
- ❖ Compositional interaction: recursive matching
- ❖ Aggregation: MLP



# Match-SRNN: Recursive Matching Structure



- ❖ Matching scores are calculated recursively (from top left to bottom right)
- ❖ We can see all matching between sub sentences have been utilized
  - ❖ sat  $\longleftrightarrow$  balls
  - ❖ The cat  $\longleftrightarrow$  the dog played
  - ❖ The cat  $\longleftrightarrow$  The dog played balls
  - ❖ The cat sat  $\longleftrightarrow$  The dog played

---

# Match-SRNN: Recursive Matching Structure (cont')

---

- ❖ Definition

$$S_1 = w_1, \dots, w_m, S_2 = v_1, \dots, v_n$$

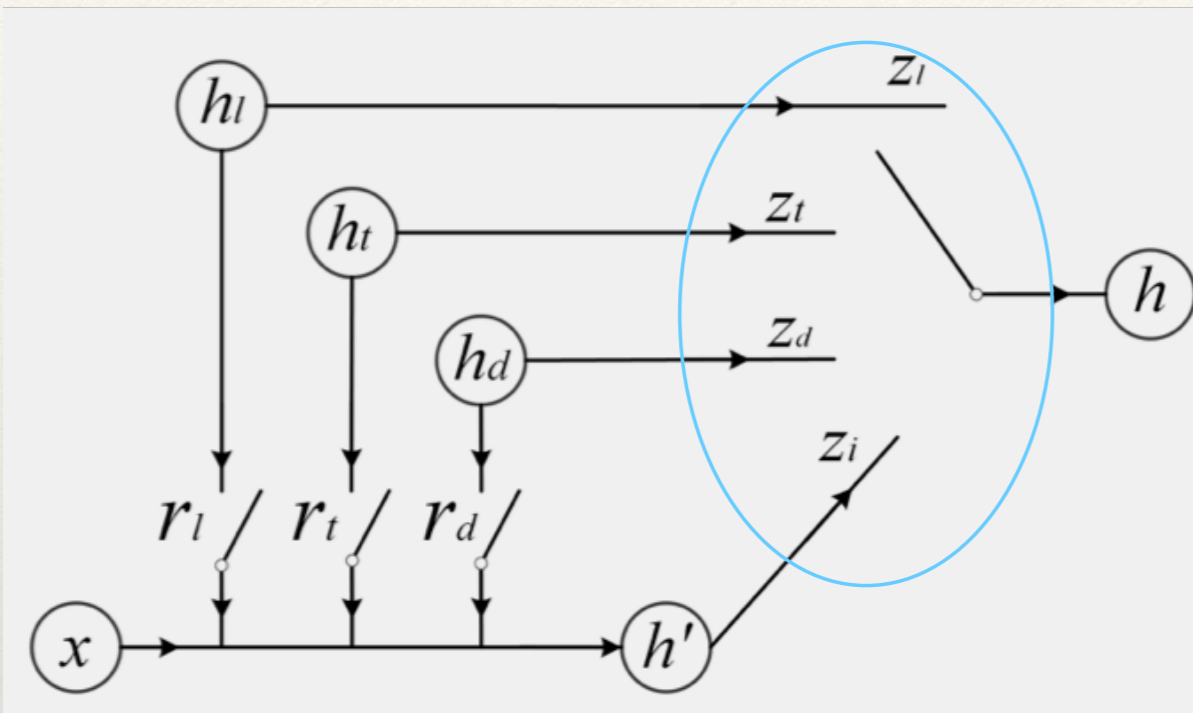
- ❖  $S_1[1 : i]$ : prefix of length  $i$

- ❖  $S_2[1 : j]$ : prefix of length  $j$

- ❖  $h_{ij}$ : match representation between  $S_1[1:i]$  and  $S_2[1:j]$

- ❖ We have  $h_{i,j} = f(h_{i-1,j} h_{i,j-1} h_{i-1,j-1} s(w_i, v_j))$

# Using Spatial GRU (two dimensions)



Softmax function is used to 'soft' choose connections among four choices.

$$q^T = [h_{i-1,j}^T, h_{i,j-1}^T, h_{i-1,j-1}^T, s_{ij}^T]^T,$$

$$r_l = \sigma(W^{(r_l)}q + b^{(r_l)}),$$

$$r_t = \sigma(W^{(r_t)}q + b^{(r_t)}),$$

$$r_d = \sigma(W^{(r_d)}q + b^{(r_d)}),$$

$$r^T = [r_l^T, r_t^T, r_d^T]^T,$$

$$z'_i = W^{(z_i)}q + b^{(z_i)},$$

$$z'_l = W^{(z_l)}q + b^{(z_l)},$$

$$z'_t = W^{(z_t)}q + b^{(z_t)},$$

$$z'_d = W^{(z_d)}q + b^{(z_d)},$$

$$[z_i, z_l, z_t, z_d] = \text{SoftmaxByRow}([z'_i, z'_l, z'_t, z'_d]),$$

$$h'_{i,j} = \phi(Ws_{ij} + U(r \odot [h_{i,j-1}^T, h_{i-1,j}^T, h_{i-1,j-1}^T]^T) + b),$$

$$h_{i,j} = z_l \odot h_{i,j-1} + z_t \odot h_{i-1,j} + z_d \odot h_{i-1,j-1} + z_i \odot h'_{i,j}.$$

# Connection with LCS

- ❖ Longest Common Sub-Sequence

- ❖ S1: A B C D E

- ❖ S2: F A C G D

- ❖ LCS: A C D

	(A)	B	(C)	(D)	E
F	0	0	0	0	0
(A)	1 ← 1		1	1	1
(C)	1	1	2	2	2
G	1	1	2	2	2
(D)	1	1	2	3 ← 3	

- ❖ Solving LCS with dynamic programming

- ❖ Step function:  $c[i, j] = \max(c[i, j - 1], c[i - 1, j], c[i - 1, j - 1] + \mathbb{I}_{x_i=y_j})$

- ❖ Backtrace: depends on the selection of “max” operation

# Connection with LCS

- ❖ Matching-SRNN can be explained with LCS
- ❖ Simplify Match-SRNN
  - ❖ Using exact word level matching signals only
  - ❖ remove the reset gate  $r$  and set hidden dimensions to 1

$$h_{i,j} = z_l \cdot h_{i,j-1} + z_t \cdot h_{i-1,j} + z_d \cdot h_{i-1,j-1} + z_i \cdot h'_{ij}$$

- ❖ Simplified Match-SRNN simulates LCS

$$c[i,j] = \max(c[i,j-1], c[i-1,j], c[i-1,j-1] + \mathbb{I}_{x_i=y_j})$$

- ❖  $z$  is obtained by SOFTMAX
- ❖ Backtrace by the value of  $z$  in simplified Match-SRNN



# Simulation with Simplified Math-SRNN

## ❖ Simulation data

- ❖ random sampled sequence
- ❖ ground truth obtained by DP
- ❖ the label is the length of LCS

	(A)	B	(C)	(D)	E
F	0	0	0	0	0
(A)	1	1	1	1	1
(C)	1	1	2	2	2
G	1	1	2	2	2
(D)	1	1	2	3	3

(a)

	(A)	B	(C)	(D)	E
F	0.0	0.0	0.0	0.0	0.0
(A)	1.0	1.0	1.0	1.0	0.9
(C)	1.0	1.0	2.1	2.1	2.0
G	1.0	1.0	2.1	2.0	2.0
(D)	1.0	1.0	2.0	3.1	3.1

(b)

	(A)	B	(C)	(D)	E
F					
(A)	0.8	0.0	0.0	0.1	
(C)	0.0	0.8	0.8	0.0	
G			0.0	0.9	
(D)			0.1	0.7	0.1
			0.1	0.1	0.9

(c)

$z_l$   $z_t$   $z_d$

Match-SRNN simulates LCS!



# Performance Evaluations on QA Task

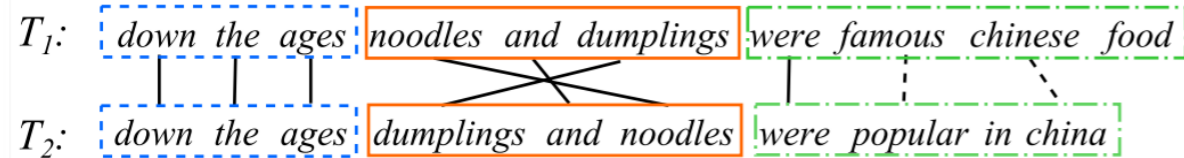
	Model	P@1	MRR
Statistic Traditional	Random	0.200	0.457
	BM25	0.579	0.726
Comosition Focused	ARC-I	0.581	0.756
	CNTN	0.626	0.781
	LSTM-RNN	0.690	0.822
	uRAE	0.398	0.652
	MultiGranCNN	0.725	0.840
Interaction Focused	MV-LSTM	0.766	0.869
	DeepMatch	0.452	0.679
	ARC-II	0.591	0.765
	MatchPyramid	0.764	0.867
	Match-SRNN	0.790	0.882

- ❖ Interaction focused methods outperformed the composition focused ones
  - ❖ Low level interaction (word level) signals are important
- ❖ Match-SRNN performs the best
  - ❖ Powerful recursive matching structure

# Application to Search — Document Level Matching

# Document Level Matching: Aggregating Matching Signals

## Sentence-sentence matching



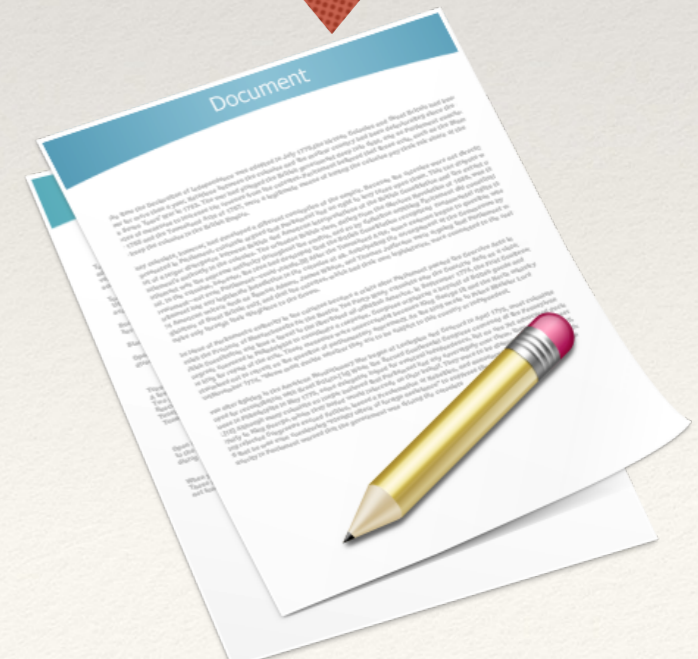
## Query-document matching

query

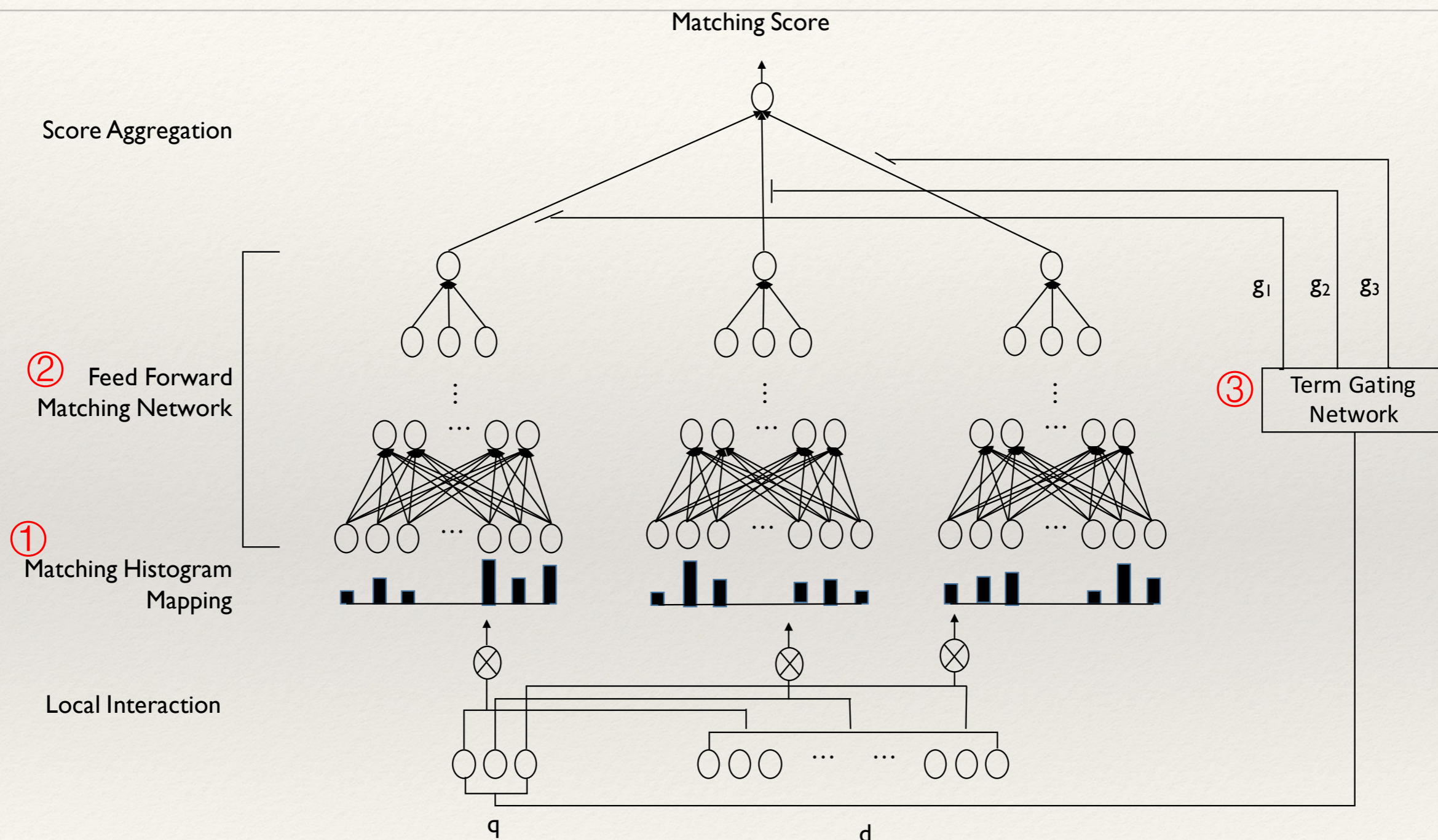
deep semantic matching 



document



# Deep Relevance Matching Model (DRMM)



Jiafeng Guo, Yixing Fan, Qingyao Ai, and W. Bruce Croft. 2016. A Deep Relevance Matching Model for Ad-hoc Retrieval. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management (CIKM '16)*. ACM, New York, NY, USA, 55-64.

# Deep Relevance Matching Model (cont')

- ❖ Learning the parameters
  - ❖ Pairwise loss:  $\ell(\mathbf{q}, \mathbf{d}^+, \mathbf{d}^-) = \max(0, 1 - F(\mathbf{q}, \mathbf{d}^+) + F(\mathbf{q}, \mathbf{d}^-))$
  - ❖ Optimization with stochastic gradient descent
- ❖ Experimental results (Robust-04 collection)

	Using topic titles			Using topic descriptions		
	MAP	nDCG@20	P@20	MAP	nDCG@20	P@20
DSSM	0.095	0.201	0.171	0.078	0.169	0.145
CDSSM	0.067	0.146	0.125	0.050	0.113	0.093
ARC-I	0.041	0.066	0.065	0.030	0.047	0.045
ARC-II	0.067	0.147	0.128	0.042	0.086	0.074
MP-IND	0.169	0.319	0.281	0.067	0.142	0.118
MP-COS	0.189	0.330	0.290	0.094	0.190	0.162
MP-DOT	0.083	0.159	0.155	0.047	0.104	0.092
<b>DRMM</b>	<b>0.279</b>	<b>0.431</b>	<b>0.382</b>	<b>0.275</b>	<b>0.437</b>	<b>0.371</b>

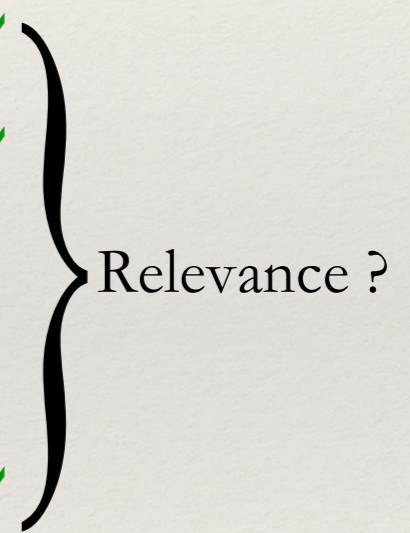
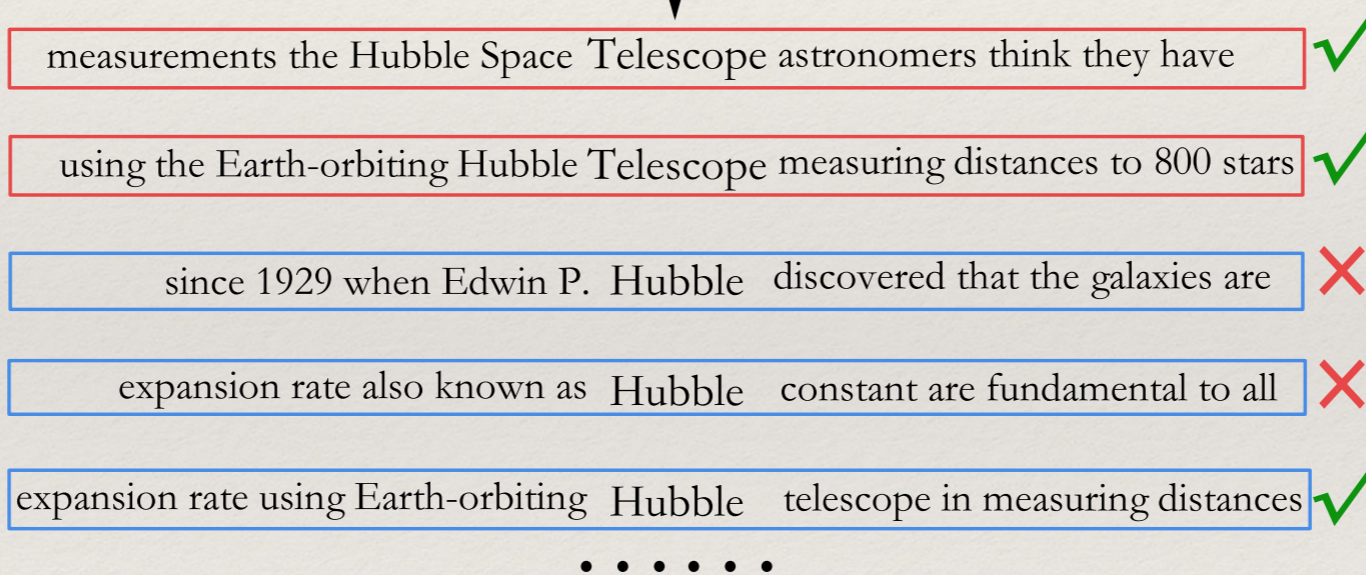
# DeepRank: Semantic Query-Document Matching

## ❖ Motivation: mimicking human-judgment of relevance

Query : Hubble Telescope Achievements

Document

Query-Centric Assumption



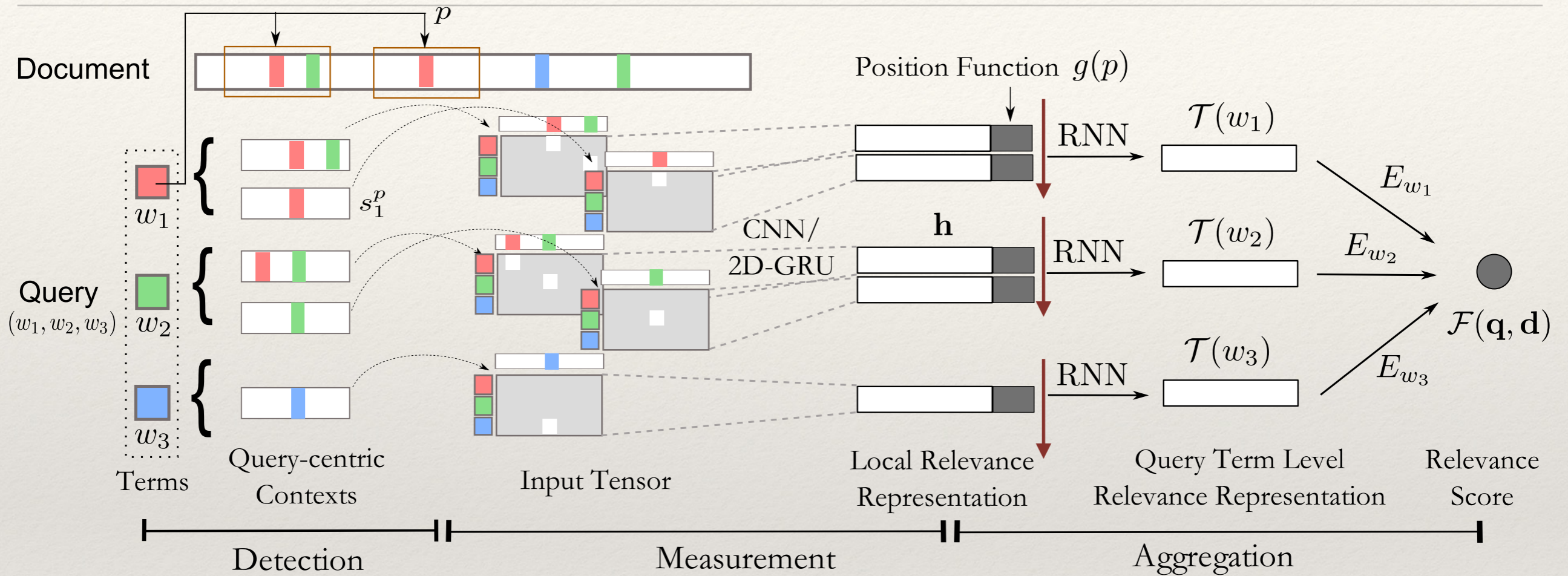
① Detection

② Measurement

③ Aggregation



# DeepRank



Focusing on the location of query terms when scanning the whole document

Determine local relevance – relevance between query and each query term -centric context, using MatchPyramid/MatchSRNN

Query term level aggregation

$$F(\mathbf{q}, \mathbf{d}) = \sum_{w \in \mathbf{q}} (E_w \mathbb{I})^T \cdot \mathcal{T}(w)$$

# Learning and Empirical Evaluation

- ❖ Learning the parameters

- ❖ Pairwise loss:  $\ell(\mathbf{q}, \mathbf{d}^+, \mathbf{d}^-) = \max(0, 1 - F(\mathbf{q}, \mathbf{d}^+) + F(\mathbf{q}, \mathbf{d}^-))$

- ❖ Optimization: stochastic gradient decent

- ❖ Experimental results

MQ2007									
Model	NDCG@1	NDCG@3	NDCG@5	NDCG@10	P@1	P@3	P@5	P@10	MAP
BM25-TITLE	0.358 <sup>-</sup>	0.372 <sup>-</sup>	0.384 <sup>-</sup>	0.414 <sup>-</sup>	0.427 <sup>-</sup>	0.404 <sup>-</sup>	0.388 <sup>-</sup>	0.366 <sup>-</sup>	0.450 <sup>-</sup>
RANKSVM	0.408 <sup>-</sup>	0.405 <sup>-</sup>	0.414 <sup>-</sup>	0.442 <sup>-</sup>	0.472 <sup>-</sup>	0.432 <sup>-</sup>	0.413 <sup>-</sup>	0.381 <sup>-</sup>	0.464 <sup>-</sup>
RANKBOOST	0.401 <sup>-</sup>	0.404 <sup>-</sup>	0.410 <sup>-</sup>	0.436 <sup>-</sup>	0.462 <sup>-</sup>	0.428 <sup>-</sup>	0.405 <sup>-</sup>	0.374 <sup>-</sup>	0.457 <sup>-</sup>
ADARANK	0.400 <sup>-</sup>	0.410 <sup>-</sup>	0.415 <sup>-</sup>	0.439 <sup>-</sup>	0.461 <sup>-</sup>	0.431 <sup>-</sup>	0.408 <sup>-</sup>	0.373 <sup>-</sup>	0.460 <sup>-</sup>
LAMBDA MART	0.412 <sup>-</sup>	0.418 <sup>-</sup>	0.421 <sup>-</sup>	0.446 <sup>-</sup>	0.481 <sup>-</sup>	0.444 <sup>-</sup>	0.418 <sup>-</sup>	0.384 <sup>-</sup>	0.468 <sup>-</sup>
DSSM	0.290 <sup>-</sup>	0.319 <sup>-</sup>	0.335 <sup>-</sup>	0.371 <sup>-</sup>	0.345 <sup>-</sup>	0.359 <sup>-</sup>	0.359 <sup>-</sup>	0.352 <sup>-</sup>	0.409 <sup>-</sup>
CDSSM	0.288 <sup>-</sup>	0.288 <sup>-</sup>	0.297 <sup>-</sup>	0.325 <sup>-</sup>	0.333 <sup>-</sup>	0.309 <sup>-</sup>	0.301 <sup>-</sup>	0.291 <sup>-</sup>	0.364 <sup>-</sup>
ARC-I	0.310 <sup>-</sup>	0.334 <sup>-</sup>	0.348 <sup>-</sup>	0.386 <sup>-</sup>	0.376 <sup>-</sup>	0.377 <sup>-</sup>	0.370 <sup>-</sup>	0.364 <sup>-</sup>	0.417 <sup>-</sup>
SQA-NOFEAT	0.309 <sup>-</sup>	0.333 <sup>-</sup>	0.348 <sup>-</sup>	0.386 <sup>-</sup>	0.375 <sup>-</sup>	0.373 <sup>-</sup>	0.372 <sup>-</sup>	0.364 <sup>-</sup>	0.419 <sup>-</sup>
DRMM	0.380 <sup>-</sup>	0.396 <sup>-</sup>	0.408 <sup>-</sup>	0.440 <sup>-</sup>	0.450 <sup>-</sup>	0.430 <sup>-</sup>	0.417 <sup>-</sup>	0.388 <sup>-</sup>	0.467 <sup>-</sup>
ARC-II	0.317 <sup>-</sup>	0.338 <sup>-</sup>	0.354 <sup>-</sup>	0.390 <sup>-</sup>	0.379 <sup>-</sup>	0.378 <sup>-</sup>	0.377 <sup>-</sup>	0.366 <sup>-</sup>	0.421 <sup>-</sup>
MATCHPYRAMID	0.362 <sup>-</sup>	0.364 <sup>-</sup>	0.379 <sup>-</sup>	0.409 <sup>-</sup>	0.428 <sup>-</sup>	0.404 <sup>-</sup>	0.397 <sup>-</sup>	0.371 <sup>-</sup>	0.434 <sup>-</sup>
MATCH-SRNN	0.392 <sup>-</sup>	0.402 <sup>-</sup>	0.409 <sup>-</sup>	0.435 <sup>-</sup>	0.460 <sup>-</sup>	0.436 <sup>-</sup>	0.413 <sup>-</sup>	0.384 <sup>-</sup>	0.456 <sup>-</sup>
DEEPRANK-2DGRU	0.439	0.439	0.447	0.473	<b>0.513</b>	0.467	0.443	0.405	0.489
DEEPRANK-CNN	<b>0.441</b>	<b>0.447</b>	<b>0.457</b>	<b>0.482</b>	0.508	<b>0.474</b>	<b>0.452</b>	<b>0.412</b>	<b>0.497</b>

---

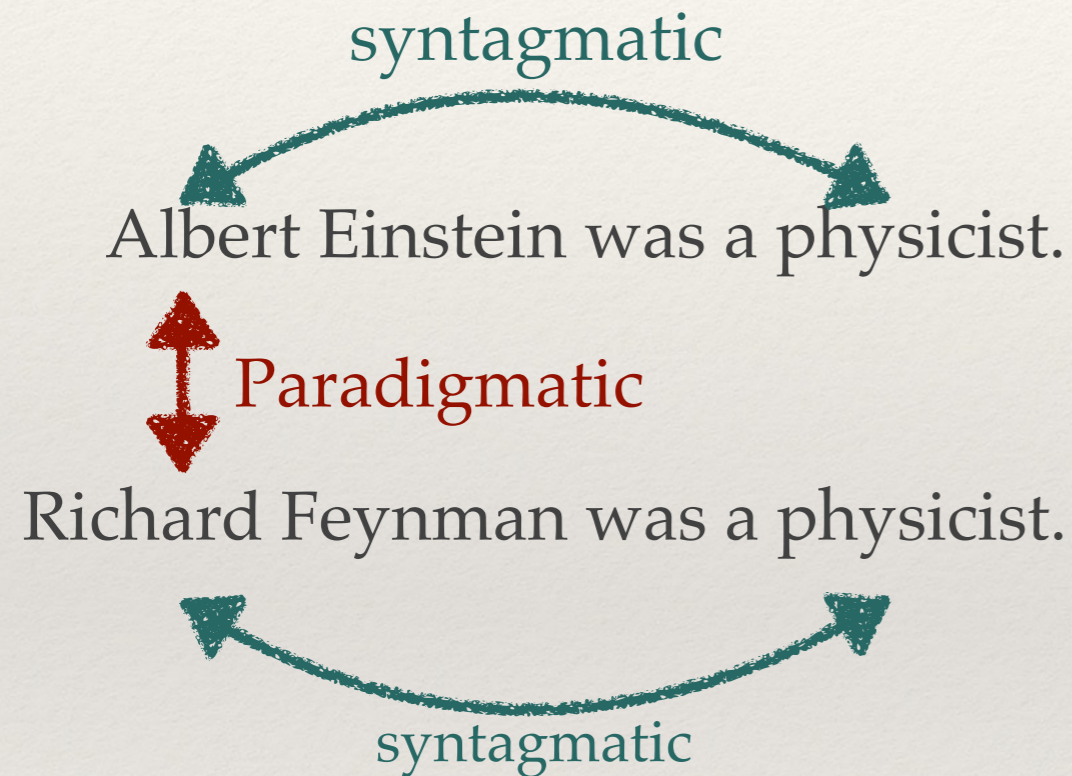
# Outline

---

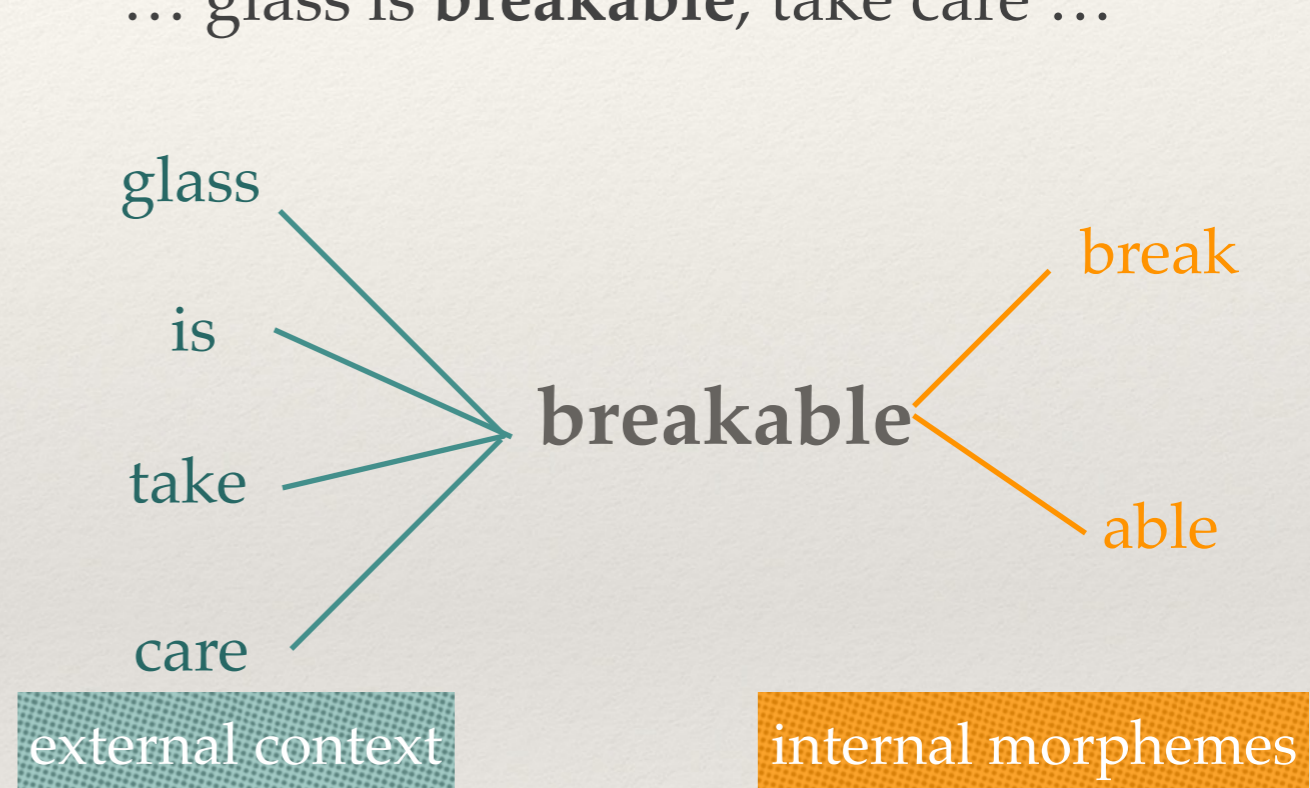
- ❖ Semantic matching in search
- ❖ Word-level matching: bridging the semantic gap
- ❖ Sentence-level matching: capturing the proximity
- ❖ Summary and discussion

# Summary

- ❖ Word level matching: bridging the semantic gap



"... glass is breakable, take care ..."



Two interpretations of distributed hypothesis

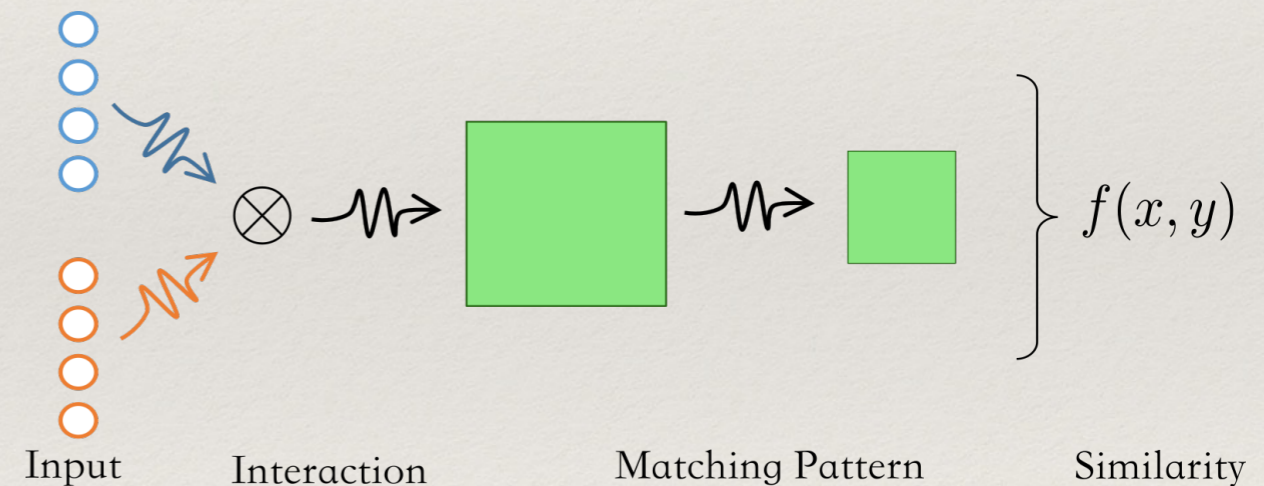
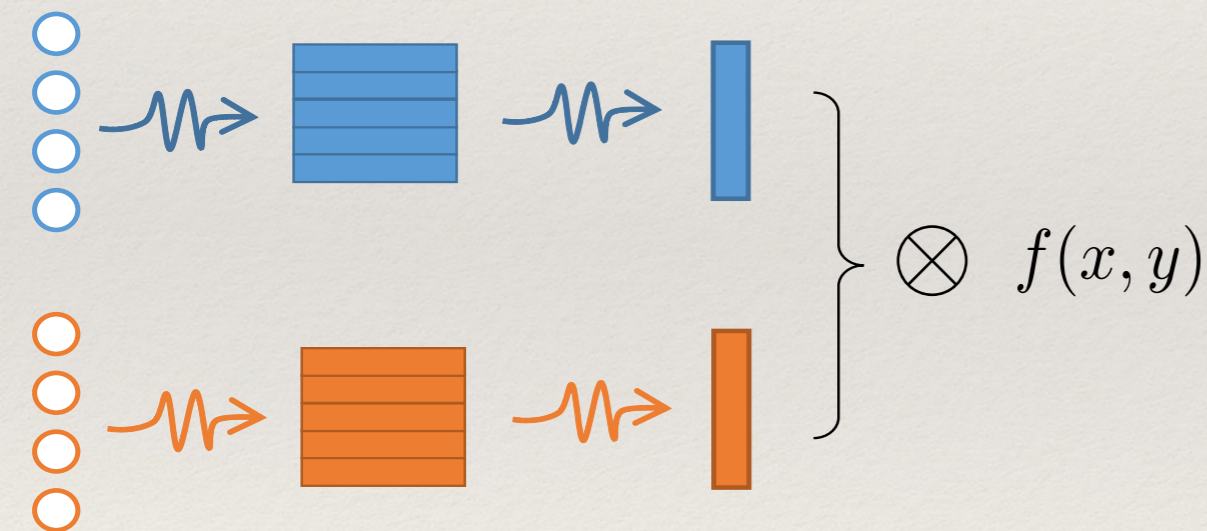
Beyond distributed hypothesis

# Summary

- ❖ Sentence level matching: capturing the proximity

Semantic representation of sentences

Aggregating fine-grained matching signals

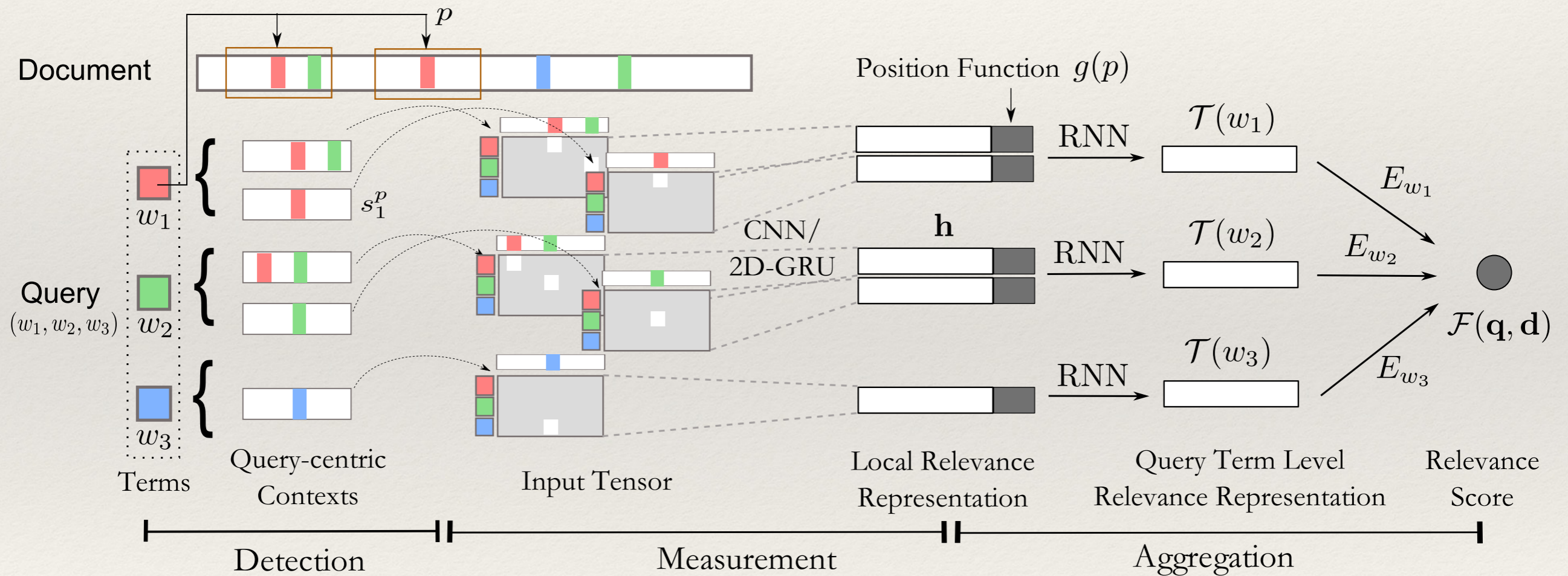


Compositional focused methods:  
representing queries and document in  
semantic space

Interaction focused methods: discovering  
the query-document matching patterns

# Summary

- ❖ Document level matching: aggregating matching signals



---

# Challenges

---

- ❖ Data: building benchmarks
  - ❖ Current: lack of large scale text matching data
  - ❖ Deep learning models has a lot of parameters
- ❖ Model: leveraging human knowledge
  - ❖ Current: most models are purely data-driven
  - ❖ Prior information (e.g., large scale knowledge base) should be helpful
- ❖ Application
  - ❖ Domain specific matching models: different application have different matching goal, e.g., in IR, relevance  $\neq$  similarity

# Easy Machine Learning Github Project

<https://github.com/ICT-BDA/EasyML>

- ❖ Purpose: ease the process of applying machine learning algorithms to real tasks
  - ❖ Machine learning tasks as data-flow DAG
  - ❖ Interactive GUI for creating, running, and managing scalable machine learning tasks
  - ❖ Deployed as web service <http://159.226.40.104:18080/dev/>

The screenshot displays the EasyML web interface. On the left, there is a sidebar with a list of examples, including tasks like Titanic Demo, Twitter Demo, and various machine learning algorithms. The main area shows a data flow DAG (Directed Acyclic Graph) with nodes such as Row\_Normalize, Word\_Segment, Word\_Filter, TFIDF, Feature\_Index, File\_Split, LogisticRegression\_Train, and BinaryClassification\_Evaluate. On the right, there are job specifications and component specifications for the selected task.

**Job Specifications**

- Job Name: 【实例】分布式 移动垃圾短信分类
- Job Owner: bdaict@hotmail.com
- Job ID: 0000139-160606112228201-bda-o
- Job Status: SUCCEEDED
- Start Time: 2017-03-28 09:29:49
- End Time: 2017-03-28 09:40:12
- Use Time: 00:10:23
- Description: 【实例】分布式 移动垃圾短信分类

**Component Specifications**

Name	Type	Value
LogisticRegression_Train		
Description		Spark版本的LogisticRegression Train
Determinacy	Boolean	false
Version		0.7
Create Time		2016-05-17 10:48:56 PM
Owner		fortianyou@hello.net
Deprecated	Boolean	no
是否使用 GraphX	Boolean	false
max_iter	Int	10
reg	Double	0.01
learn_rate	Double	0.0001

The screenshot shows the Github repository page for ICT-BDA / EasyML. The repository is public and has 125 stars, 1,463 forks, and 294 forks. The repository description states: "Easy Machine Learning is a general-purpose dataflow-based system for easing the process of applying machine learning algorithms to real world tasks." The repository has 73 commits, 2 branches, 0 releases, 5 contributors, and is licensed under Apache-2.0. The latest commit is by sinlychen, updating the baidu cloud url to fix mysql container restart script, committed 19 days ago. Other recent commits include removing local repository jar and updating settings, and fixing a 3rd jar build problem.



Thanks!