

WWW 2018 Tutorial  
April 23, 2018  
Lyon, France

# Deep Learning for Matching in Search and Recommendation

Jun Xu

Chinese Academy  
of Sciences

Xiangnan He

National University of  
Singapore

Hang Li

Toutiao AI Lab

# Outline of Tutorial

- Unified view of matching in search and recommendation
- Part 1: Traditional Approaches to Matching
- Part 2: Deep Learning Approaches to Matching
- Summary

# Overview of Web Search Engine

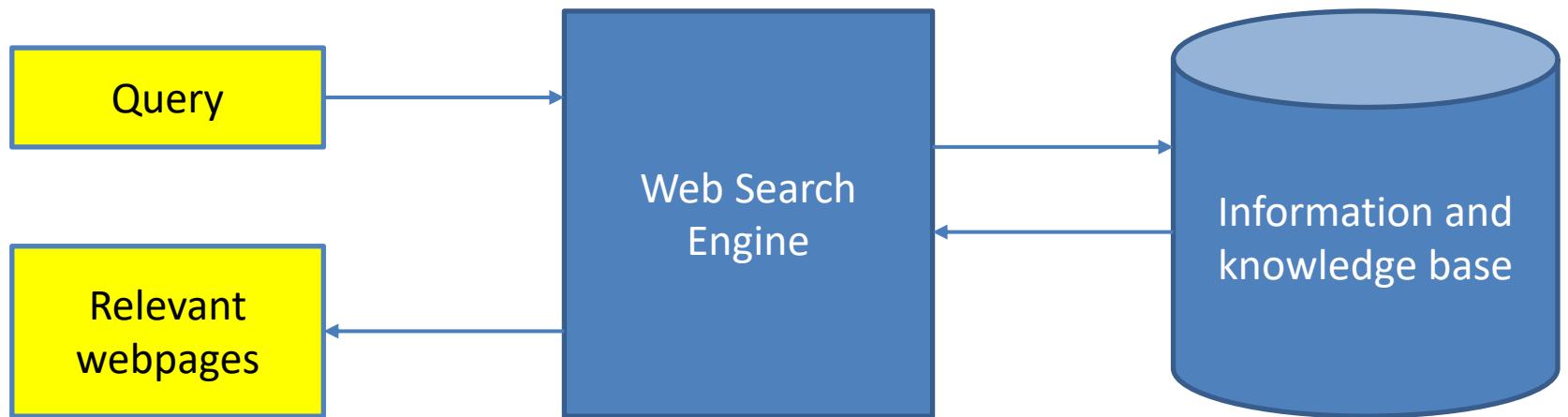
Information pull: a user pulls information by making a specific request

**User intent** is explicitly reflected in query:

- Keywords, Question

**Content** is in

- Webpages, images, ...



**Key challenge:** query-document semantic gap

# Example of Query-Document Mismatch

Query	Document	Term matching	Semantic matching
seattle best hotel	seattle best hotels	no	yes
pool schedule	swimmingpool schedule	no	yes
natural logarithm transformation	logarithm transformation	partial	yes
china kong	china hong kong	partial	no
why are windows so expensive	why are macs so expensive	partial	no

# Same Search Intent Different Query Representations

## Example: “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

## Example: “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	wwwwww youtube com
www youtube	www youtube com	wwwwww 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

---

# Overview of Recommendation Engine

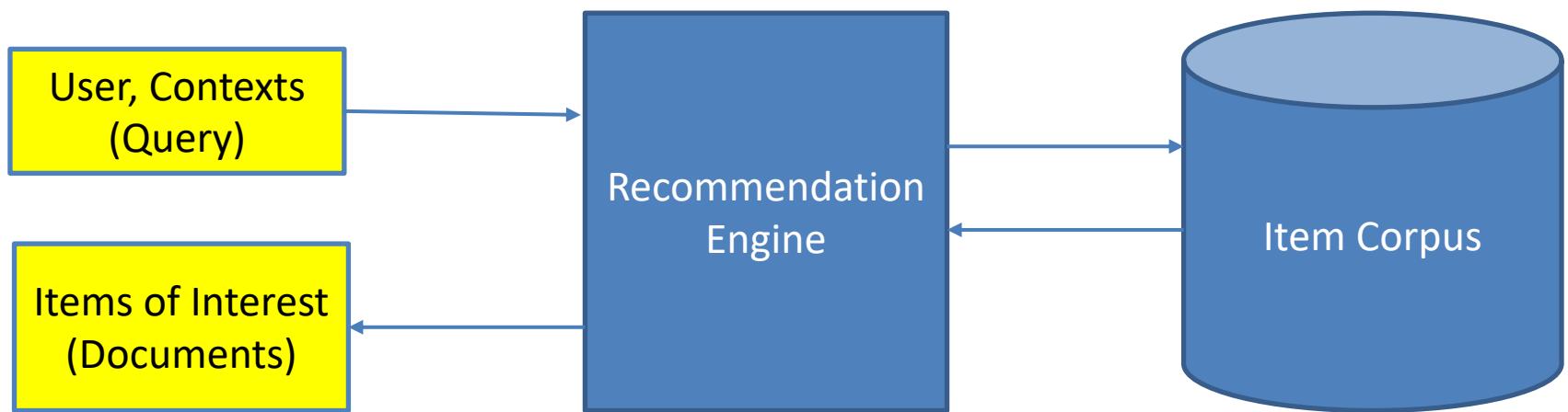
Information push: the system pushes information to a user by guessing the user interest

**User Interest** is implicitly reflected in:

- Interaction history
- Demographics
- Contexts

**Items** can be:

- Products, News, Movies, Videos, Friends ...

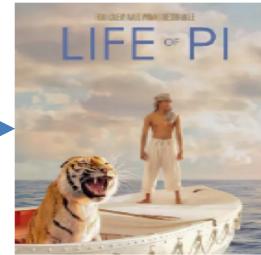


**Key challenge:** user-item semantic gap

- Even severe than search, since user and item are two **different types of entities** and are represented by different features

# Example of User-Item Semantic Gap

Movie Recommendation



## User Profile (query):

- User ID
- Rating history
- Age, Gender
- Income level
- Time of the day

.....

## Item Profile (document):

- Item ID
- Description
- Category
- Price
- Image

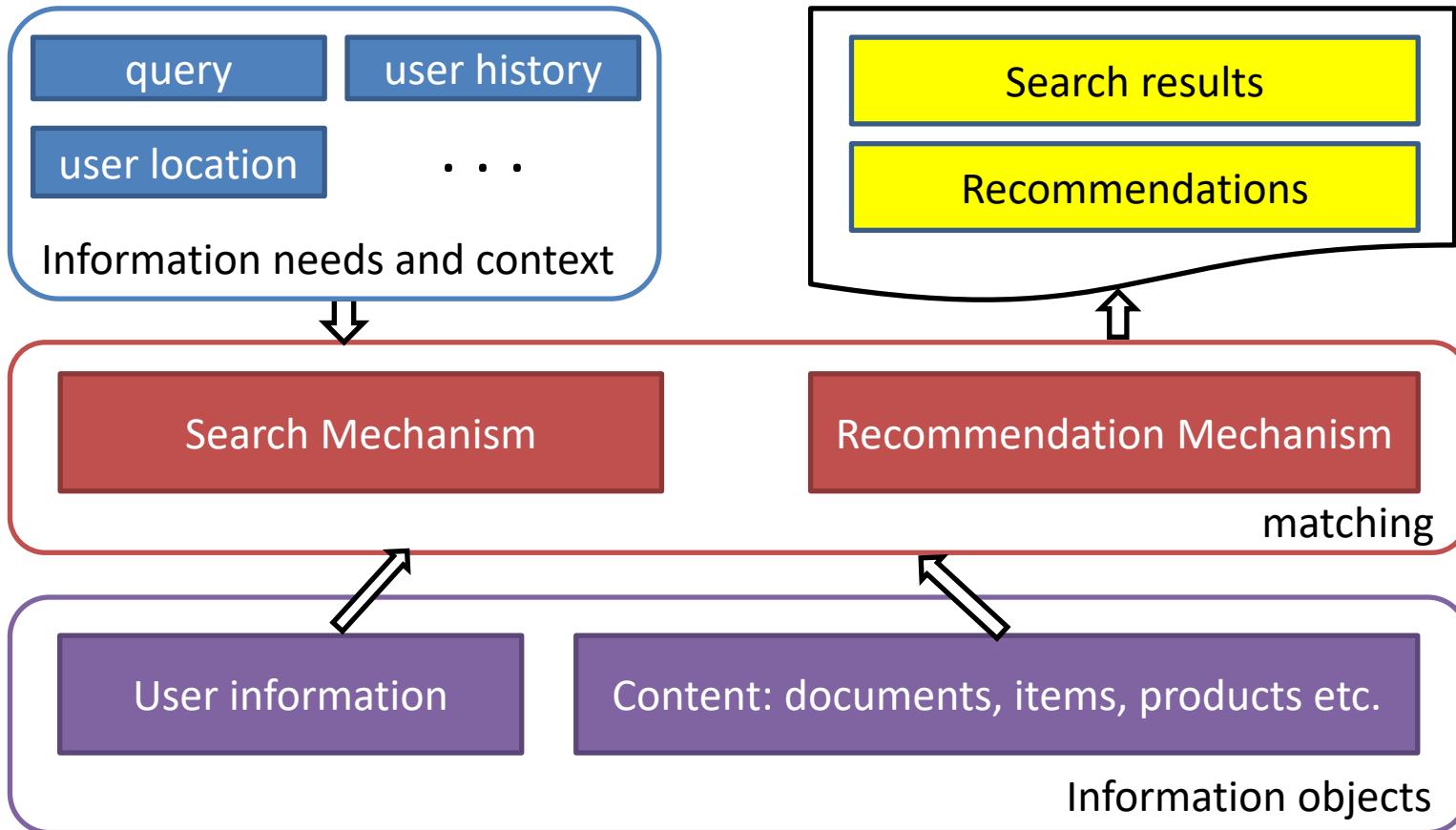
.....

There may be **no overlap** between user features and item features  
Matching cannot be done on the superficial feature level!

# Information Providing Mechanisms of Search and Recommendation (Hector et al., CACM' 11)

	Search	Recommendation
Delivery model	Pull	Push or pull
Beneficiary	User	User and provider
Unexpected good?	No	Yes
Collective knowledge	Maybe	Maybe
Query available	Yes	Maybe
Context dependent	Maybe	Maybe

# Unified View on Matching in Search and Recommendation (Hector et al, CACM'11)



Difference for search and recommendation: **features** used for matching!

# Semantic Gap is Biggest Challenge in both Search and Recommendation

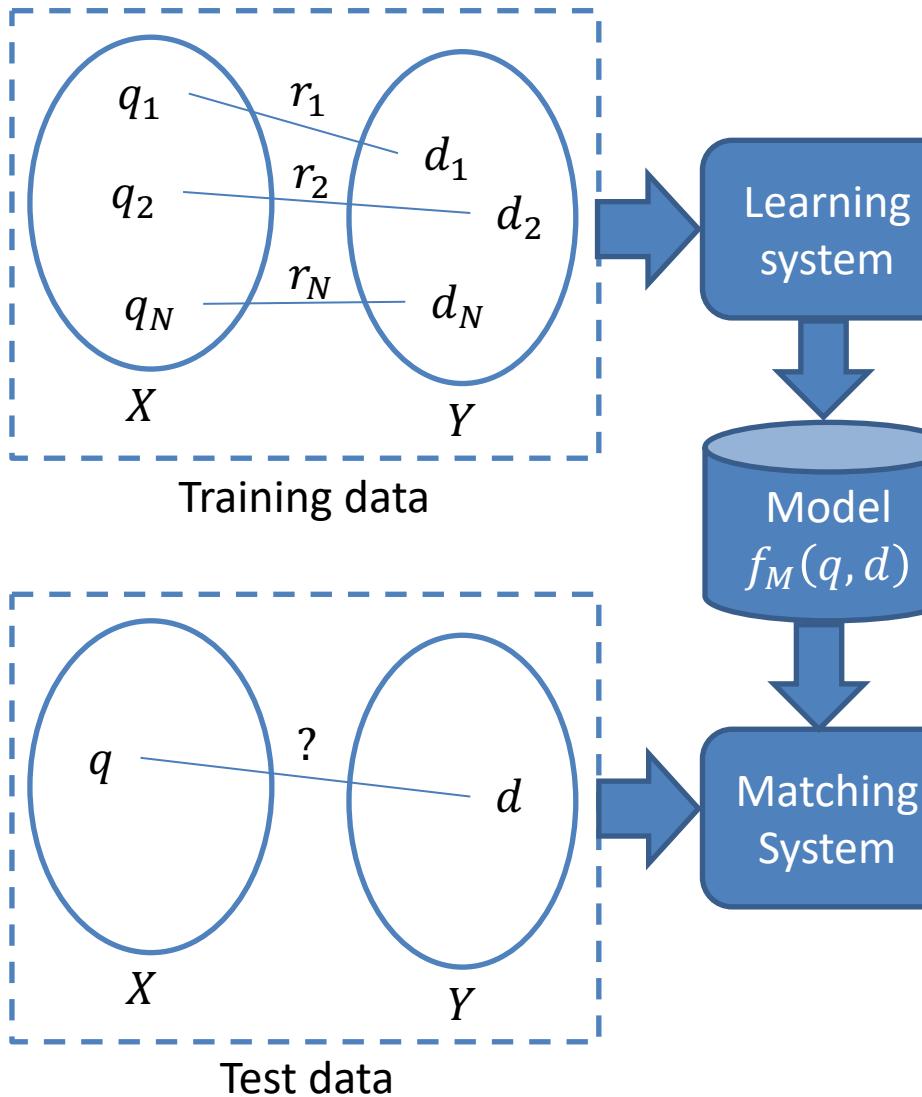
## Query-document Mismatch

- Same intent can be represented by different queries (representations)
- Search is still mainly based on term level matching
- Query document mismatch occurs, when searcher and author use different representations

## User-item Semantic Gap

- Features are used to represent a user and an item may be totally different (e.g., ID feature)
- Even when they partially overlap in features, it is insensible to conduct direct matching

# Machine Learning for Matching



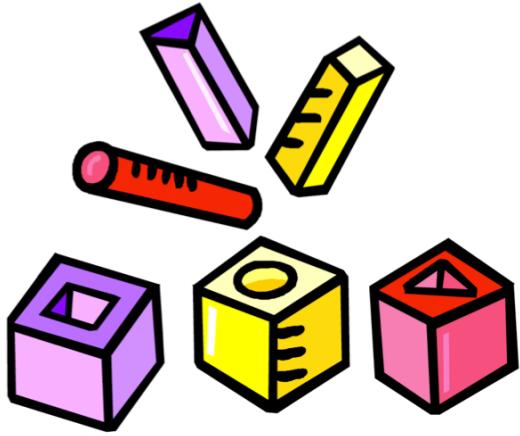
- Using relations in data for learning the matching function  $f_M(q, d)$  or  $P(r|q, d)$
- Training data  $\{(q_i, d_i, r_i)\}_{i=1}^N$ 
  - Queries and documents (users and items) represented with feature vectors or ID's
  - Target can be binary or numerical values

# Organization of the Tutorial

- Unified view of matching in search and recommendation (Jun Xu)
- Part 1: Traditional Approaches to Matching
  - Traditional matching models for search (Jun Xu)
  - Traditional matching models for recommendation (Xiangnan He)
- Part 2: Deep Learning Approaches to Matching
  - Overview (Jun Xu)
  - Deep matching models for search (Jun Xu)
  - Deep matching models for recommendation (Xiangnan He)
- Summary (Xiangnan He)

# Outline of Tutorial

- Unified view of matching in search and recommendation
- **Part 1: Traditional Approaches to Matching**
  - Traditional matching models for search
  - Traditional matching models for recommendation
- Part 2: Deep Learning Approaches to Matching
- Summary



# QUERY-DOCUMENT MATCHING

# Key Factors for Query-Document Matching

## Query:

Down the ages noodles and dumplings were famous Chinese food

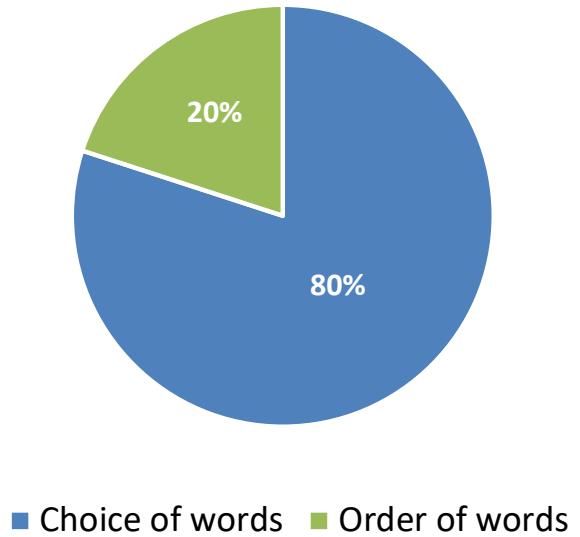
## Document:

... down the ages dumplings and noodles were popular in China ...

- Bridging the semantic gap between words
    - Semantically similar words: famous ~ popular, Chinese ~ China
  - Capturing order of words
    - N-grams: “down the ages” ~ “down the ages”
    - N-terms: “noodles and dumplings” ~ “dumplings and noodles”
- .....

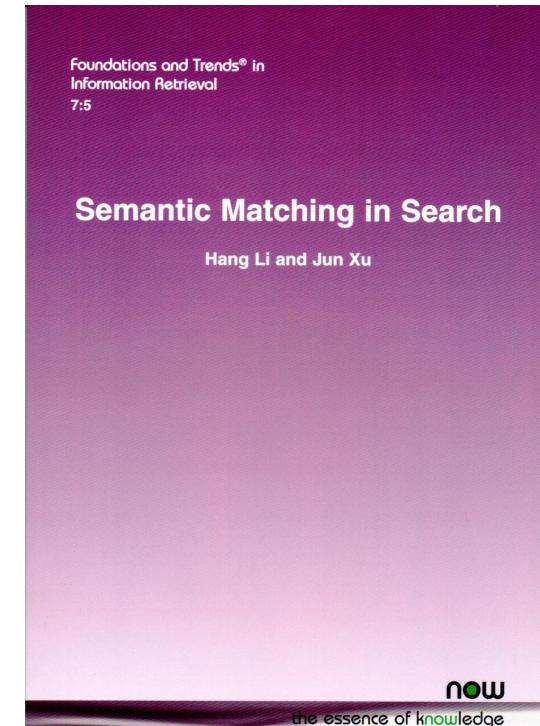
# Information from Choice of Words and Order of Words (Ross, '02)

- Assume:
  - Size of vocabulary = 10,000
  - Average sentence length = 20
- Rough contributions of information in bits
  - From the selection of words:  $\log_2(10000^2)$
  - From the order of words:  $\log_2(20!)$
- “Over 80% of the potential information in language being in the choice of words without regard to the order in which they appear”
  - 80%: choice of words
  - 20%: order of words



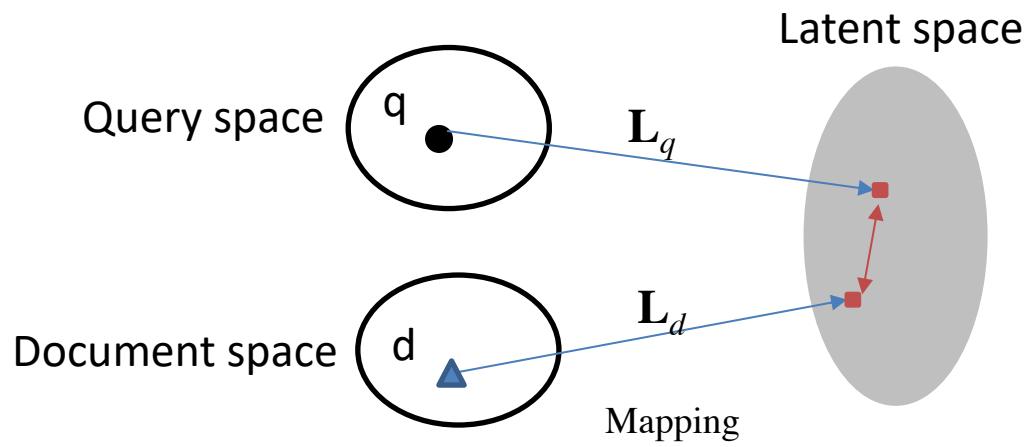
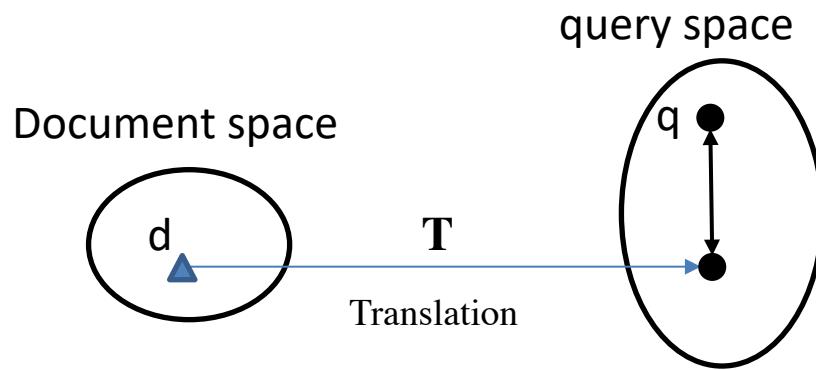
# Traditional Approaches to Query-Document Semantic Matching

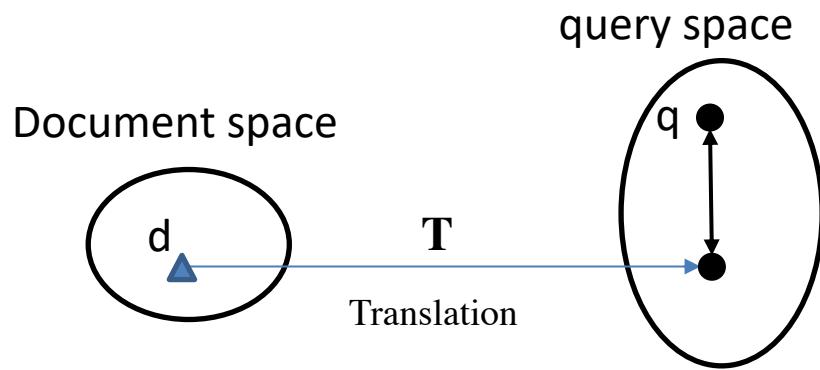
- Matching by query formulation
- Matching with term dependency
- Matching with topic model
- Matching with translation model
- Matching in latent space model



# Traditional Matching Models for Search

- **Matching with machine translation:** mapping document to query space
- **Matching in latent space:** mapping query and document into a latent space





# MATCHING WITH TRANSLATION MODEL

# Statistical Machine Translation (SMT)

- Given a sentence  $C$  in source language, translates it into sentence  $E$  in target language

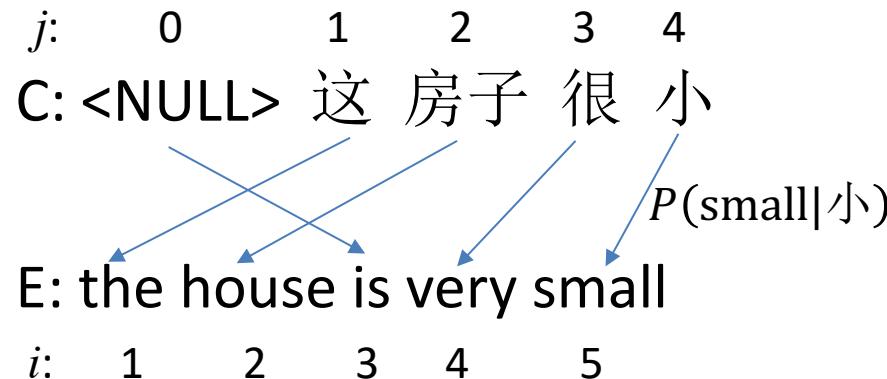
$$E^* = \operatorname{argmax}_E P(E|C)$$

- Linear combination of features

$$P(E|C) = \frac{1}{Z(C,E)} \exp \sum_i \lambda_i h(C, E)$$

$$E^* = \operatorname{argmax}_E \sum_i \lambda_i h(C, E)$$

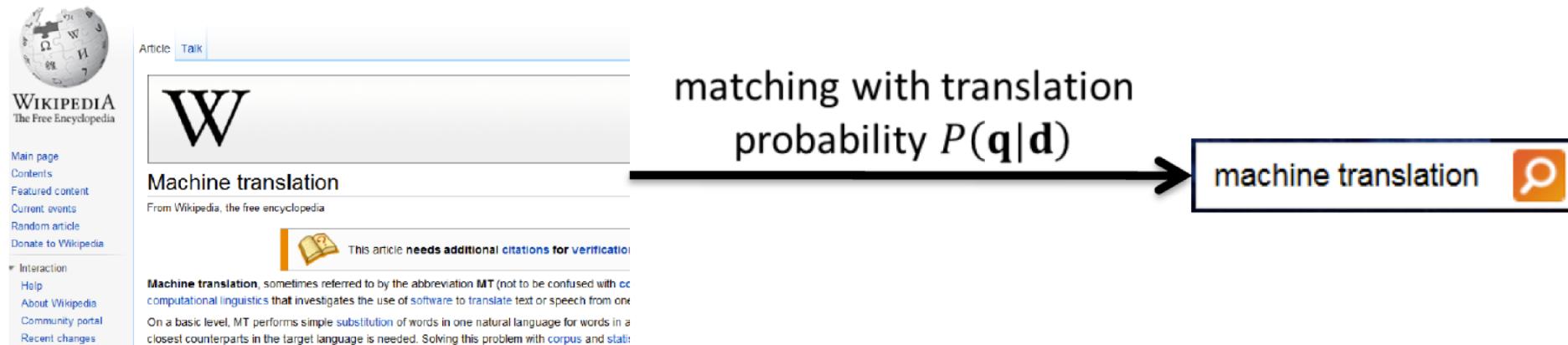
# Word-based Model: IBM Model One (Brown et al., 1993)



- Generating target sentence
  - Choose length of target sentence  $I$
  - For each position  $i (i = 1, 2, \dots, I)$ 
    - Choose position  $j$  in source sentence  $C$
    - Generate target word  $e_i$  according to  $P(e_i | c_j)$

$$P(E|C) = \frac{\epsilon}{(J+1)^I} \prod_{i=1}^I \sum_{j=0}^J P(e_i | c_j)$$

# Statistical Machine Translation for Query-Document Matching



- Translating document  $\mathbf{d}$  to query  $\mathbf{q}$
- Matching degree: translation probability  $P(\mathbf{q}|\mathbf{d})$
- Key difference from conventional translation model
  - Translation within the same language (need to handle self-translation)

# Matching with Word-based Translation Models

- Basic model

$$P(\mathbf{q}|\mathbf{d}) = \prod_{q \in \mathbf{q}} P(q|\mathbf{d}) = \prod_{q \in \mathbf{q}} \sum_{w \in \mathbf{d}} P(q|w)P(w|d)$$

translation probability

Document language model

- Smoothing to avoid zero translation probability (Berger & Lafferty '99)

$$P(\mathbf{q}|\mathbf{d}) = \prod_{q \in \mathbf{q}} \left( \alpha \textcolor{red}{P}(q|C) + (1 - \alpha) \left( \sum_{w \in \mathbf{d}} P(q|w)P(w|d) \right) \right)$$

background language model

Document language model

- Adding self-translation (Gao et al., '10)

$$P(\mathbf{q}|\mathbf{d}) = \prod_{q \in \mathbf{q}} \left( \alpha \textcolor{red}{P}(q|C) + (1 - \alpha) \left( \beta \textcolor{red}{P}(q|\mathbf{d}) - (1 - \beta) \left( \sum_{w \in \mathbf{d}} P(q|w)P(w|d) \right) \right) \right)$$

unsmoothed document language model

Smoothed document language model

# Bridging the Semantic Gap between Words

- Translation matrix can bridge the semantic gap between query words and document words

$q$	$t(q   w)$
solzhenitsyn	0.319
citizenship	0.049
exile	0.044
archipelago	0.030
alexander	0.025
soviet	0.023
union	0.018
komsomolskaya	0.017
treason	0.015
vishnevskaya	0.015

$w = \text{solzhenitsyn}$

$q$	$t(q   w)$
carcinogen	0.667
cancer	0.032
scientific	0.024
science	0.014
environment	0.013
chemical	0.012
exposure	0.012
pesticide	0.010
agent	0.009
protect	0.008

$w = \text{carcinogen}$

$q$	$t(q   w)$
zubin_mehta	0.248
zubin	0.139
mehta	0.134
philharmonic	0.103
orchestra	0.046
music	0.036
bernstein	0.029
york	0.026
end	0.018
sir	0.016

$w = \text{zubin}$

$q$	$t(q   w)$
pontiff	0.502
pope	0.169
paul	0.065
john	0.035
vatican	0.033
ii	0.028
visit	0.017
papal	0.010
church	0.005
flight	0.004

$w = \text{pontiff}$

$q$	$t(q   w)$
everest	0.439
climb	0.057
climber	0.045
whittaker	0.039
expedition	0.036
float	0.024
mountain	0.024
summit	0.021
highest	0.018
reach	0.015

$w = \text{everest}$

$q$	$t(q   w)$
wildlife	0.705
fish	0.038
acre	0.012
species	0.010
forest	0.010
environment	0.009
habitat	0.008
endangered	0.007
protected	0.007
bird	0.007

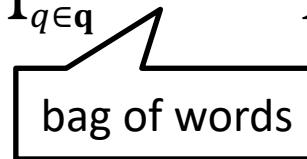
$w = \text{wildlife}$

# Capturing the Proximity

- Word-based models cannot capture the proximity

$$P(\mathbf{q}|\mathbf{d}) = \prod_{q \in \mathbf{q}} P(q|\mathbf{d}) = \prod_{q \in \mathbf{q}} \sum_{w \in \mathbf{d}} P(q|w)P(w|d)$$

bag of words



- Phrase-based translation models can capture the proximity (Gao et al., '10)

$\mathbf{d}$ :	... cold home remedies ...	<i>title</i>
$S$ :	[“cold”, “home remedies”]	<i>segmentation</i>
$T$ :	[“stuffy nose”, “home remedy”]	<i>translation</i>
$M$ :	( $1 \rightarrow 2, 2 \rightarrow 1$ )	<i>permutation</i>
$\mathbf{q}$ :	“home remedy stuffy nose”	<i>query</i>

# Learned Phrase Translation Probabilities

$\mathbf{q}$	$P(\mathbf{q} \mathbf{w})$	$\mathbf{q}$	$P(\mathbf{q} \mathbf{w})$
titanic	0.43195	sierra vista	0.61717
rms titanic	0.03793	sv	0.02260
titanic sank	0.02114	vista	0.01678
titanic sinking	0.01695	sierra	0.01581
titanic survivors	0.01537	az	0.00417
titanic ship	0.01112	bella vista	0.00320
titanic sunk	0.00960	arizona	0.00223
titanic pictures	0.00593	dominoes sierra	0.00221
		vista	
titanic exhibit	0.00540	dominos sierra vista	0.00221
ship titanic	0.00383	meadows	0.00029

$\mathbf{w} = \text{rms titanic}$

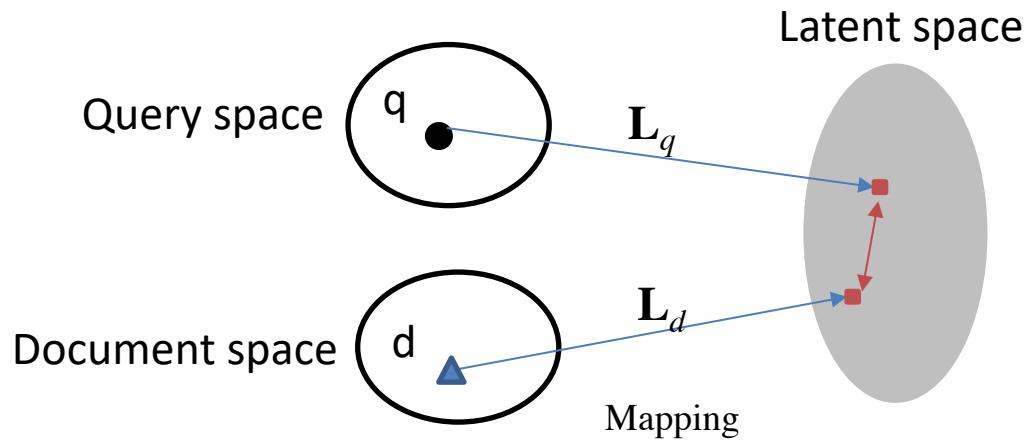
$\mathbf{w} = \text{sierra vista}$

# Experimental Results

Models	NDCG@1	NDCG@3	NDCG@10
BM25 (baseline)	0.3181	0.3413	0.4045
WTM (without self-translation)	0.3210	0.3512	0.4211
WTM (with self-translation)	0.3310	0.3566	0.4232
PTM	0.3355	0.3605	0.4254

Based on a large scale real world data set containing 12,071 English queries sampled from one-year query log files of a commercial search engine (Gao et al., 2010)

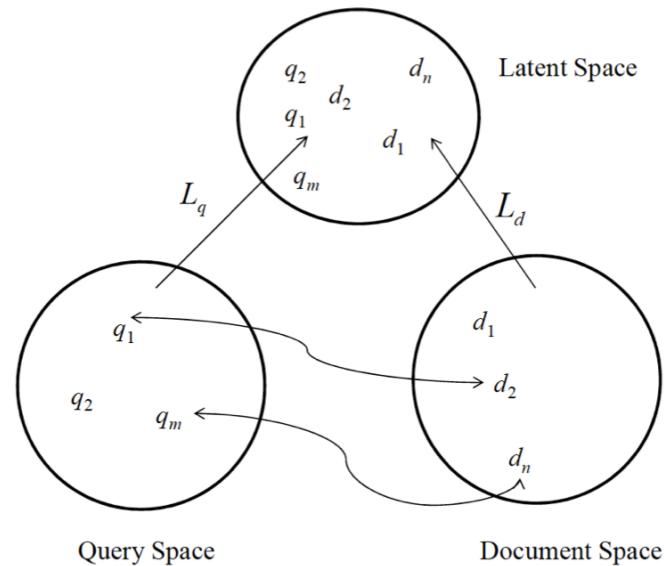
- Word-based translation model (WTM) outperformed the baseline
  - Translation probabilities bridge the semantic gap between query words and document words
  - Self-translation is effective
- Phrase-based translation model (PTM) further improved the performances though capturing the proximity information



# MATCHING IN LATENT SPACE

# Matching in Latent Space

- Assumption
  - Queries/documents have similarities
  - Click-through data represent “similarity” relations between queries and documents
- Approach
  - Project queries and documents to latent space
  - With some regularization or constraints



# Partial Least Square (PLS)

- Two spaces:  $\mathcal{X} \subset \mathbb{R}^m$  and  $\mathcal{Y} \subset \mathbb{R}^n$
- Training data:  $\{(x_i, y_i, r_i)\}_{i=1}^N$ ,  $r_i \in \{+1, -1\}$  or  $r_i \in \mathbb{R}$
- Model
  - Dot product as similarity:  $f(x, y) = \langle L_X^T x, L_Y^T y \rangle = x^T L_X L_Y^T y$
  - $L_X$  and  $L_Y$  are two linear (and orthonormal) transformations
- Objective function

$$\begin{aligned} & \operatorname{argmax}_{L_X, L_Y} \sum_{r_i=+1} x_i^T L_X L_Y^T y_i - \sum_{r_i=-1} x_i^T L_X L_Y^T y_i \\ \text{s.t. } & L_X^T L_X = I_{K \times K}, L_Y^T L_Y = I_{K \times K} \end{aligned}$$

# Regularized Mapping to Latent Space (RMLS)

- Two spaces:  $\mathcal{X} \subset \mathbb{R}^m$  and  $\mathcal{Y} \subset \mathbb{R}^n$
- Training data:  $\{(x_i, y_i, r_i)\}_{i=1}^N$ ,  $r_i \in \{+1, -1\}$  or  $r_i \in \mathbb{R}$
- Model
  - Dot product as similarity:  $f(x, y) = \langle L_X^T x, L_Y^T y \rangle = x^T L_X L_Y^T y$
  - $L_X$  and  $L_Y$  are two linear (and orthonormal) transformations with  $\ell_1$  and  $\ell_2$  regularizations (sparse transformations)
- Objective function

$$\operatorname{argmax}_{L_X, L_Y} \sum_{r_i=+1} x_i^T L_X L_Y^T y_i - \sum_{r_i=-1} x_i^T L_X L_Y^T y_i$$

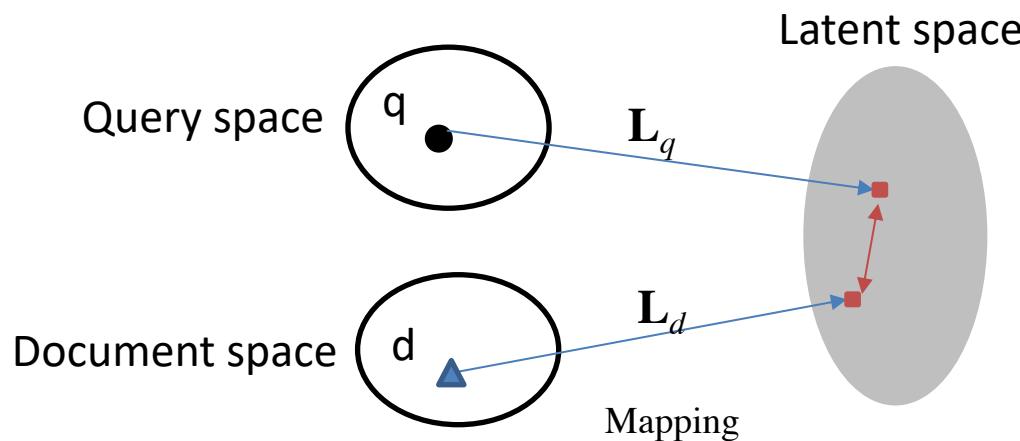
s. t.  $|L_X| \leq \lambda_X, |L_Y| \leq \lambda_Y, \|L_X\| \leq \vartheta_X, \|L_Y\| \leq \vartheta_Y$

# PLS v.s. RMLS

	PLS	RMLS
Transformation Assumption	orthonormal	and regularization
Optimization Method	singular value decomposition	coordinate ascent
Optimality	global optimum	local optimum
Efficiency	low	high
Scalability	low	high

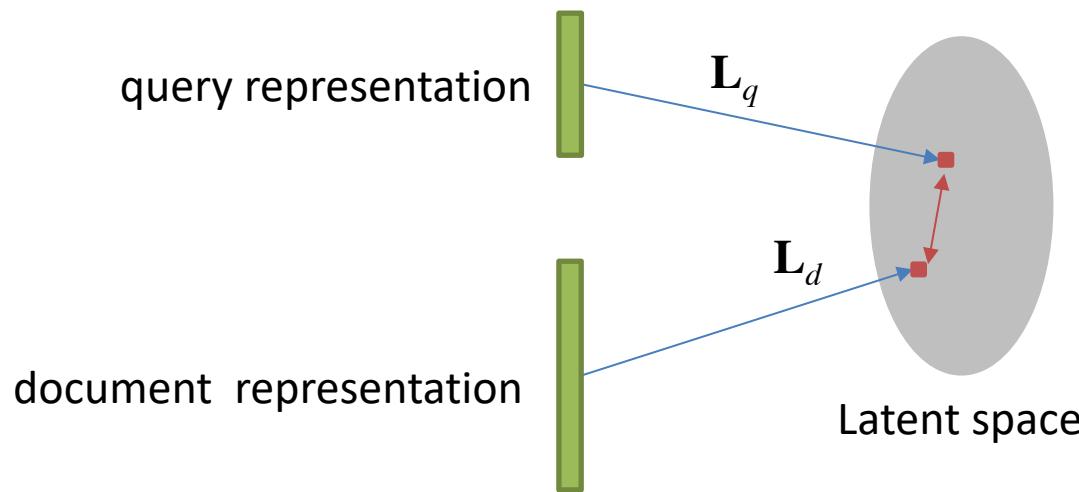
# Bridging the Semantic Gap

- Latent space models bridge semantic gap between words through
  - Reducing dimensionality of latent space (from term level matching to semantic matching)
  - Correlating semantically similar terms (matrices are not diagonal)
- Automatically learning mapping functions from data



# Capturing the Order Information

- Depends on the features for representing the queries and documents
  - Bag-of-words representations: order of words missed
  - Bag of phrases or other proximity features: capture the order of words



# Experimental Results

	NDCG@1	NDCG@3	NDCG@5
BM25	0.637	0.690	0.690
SSI	0.538	0.621	0.629
BLTM	0.657	0.702	0.701
PLS	0.676	0.728	<b>0.736</b>
RMLS	<b>0.686</b>	<b>0.732</b>	0.729

Based on a web search data set containing 94,022 queries and 111,631 documents.  
Click through associated with the queries and documents at a search engine is used.

- Latent space models work better than baseline (BM25)
- RMLS works equally well as PLS, with higher learning efficiency and scalability

# References

- Adam Berger, Rich Caruana, David Cohn, Dayne Freitag, and Vibhu Mittal. Bridging the Lexical Chasm: Statistical Approaches to Answer-Finding. In SIGIR 2000.
- Jianfeng Gao, Xiaodong He, and JianYun Nie. Click-through-based Translation Models for Web Search: from Word Models to Phrase Models. In CIKM 2010.
- Jianfeng Gao, Kristina Toutanova, and Wen-tau Yih. Clickthrough-based latent semantic models for web search. In SIGIR 2011.
- Jianfeng Gao : Statistical Translation and Web Search Ranking. <http://research.microsoft.com/en-us/um/people/jfgao/paper/SMT4IR.res.pptx>
- Dustin Hillard, Stefan Schroedl, and Eren Manavoglu, Hema Raghavan, and Chris Leggetter. Imrpoved Ad Relevance in Sponsored Search. In WSDM 2010.
- Jian Huang, Jianfeng Gao, Jiangbo Miao, Xiaolong Li, Kuansan Wang, Fritz Behr, and C. Lee Giles. Exploring web scale language models for search query processing. In WWW 2010.
- Ea-Ee Jan, Shih-Hsiang Lin, and Berlin Chen. Translation Retrieval Model for Cross Lingual Information Retrieval. In AIRS 2010.
- Rong Jin, Alex G. Hauptmann, and Chengxiang Zhai. Title Language Model for Information Retrieval. In SIGIR 2002.
- Maryan Karimzadehgan and Chengxiang Zhai. Estimation of Statistical Translation Models based on Mutual Information for Ad Hoc Information Retrieval. In SIGIR 2010.
- David Mimno , Hanna M. Wallach , Jason Naradowsky , David A. Smith, Andrew McCallum. Polylingual topic models. In EMNLP 2009.

# References

- Adam Berger and John Lafferty. Information Retrieval as Statistical Translation. In SIGIR 1999.
- Jae-Hyun Park, W. Bruce Croft, and David A. Smith. Qusi-Synchronous Dependence Model for Information Retrieval. In CIKM 2011.
- Stefan Riezler and Yi Liu. Query Rewriting Using Monolingual Statistical Machine Translation. In ACL 2010.
- Dolf Trieschnigg, Djoerd Hiemstra, Franciska de Jong, and Wessel Kraaij. A cross-lingual Framework for Monolingual Biomedical Information Retrieval. In CIKM 2010.
- Elisabeth Wolf, Delphine Bernhard, and Iryna Gurevych. Combining Probabilistic and Translation-based Models for Information Retrieval based on Word Sense Annotations. In CLEF Workshop 2009.
- D.R. Hardoon, S. Szedmak, and J. Shawe-Taylor. Canonical correlation analysis: An overview with application to learning methods. Neural Computation, 2004.
- Jianfeng Gao, Kristina Toutanova and Wen-tau Yih. Clickthrough-based latent semantic models for web search. In Proc. of SIGIR, 2011.
- R. Rosipal and N. Krämer. Overview and recent advances in partial least squares. Subspace, Latent Structure and Feature Selection, 2006.
- Wei Wu, Hang Li, and Jun Xu. Learning Query and Document Similarities from Click-through Bipartite Graph with Metadata. Microsoft Research Technical Report, 2011.
- Jun Xu, Hang Li, Chaoliang Zhong, Relevance Ranking Using Kernels, In Proceedings of the 6th Asian Information Retrieval Societies Symposium (AIRS'10), 1-12, 2010.
- Hector Garcia-Molina, Georgia Koutrika, Aditya Parameswaran, Information Seeking: Convergence of Search, Recommendations, and Advertising Communications of the ACM, Vol. 54 No. 11, Pages 121-130.
- Brian H. Ross. Psychology of Learning and Motivation: Advances in Research and Theory. Elsevier. 2002.

# Outline of Tutorial

- Unified view of matching in search and recommendation
- **Part 1: Traditional Approaches to Matching**
  - Traditional matching models for search
  - Traditional matching models for recommendation
    - Collaborative Filtering Models
    - Generic Feature-based Models
- Part 2: Deep Learning Approaches to Matching
- Summary

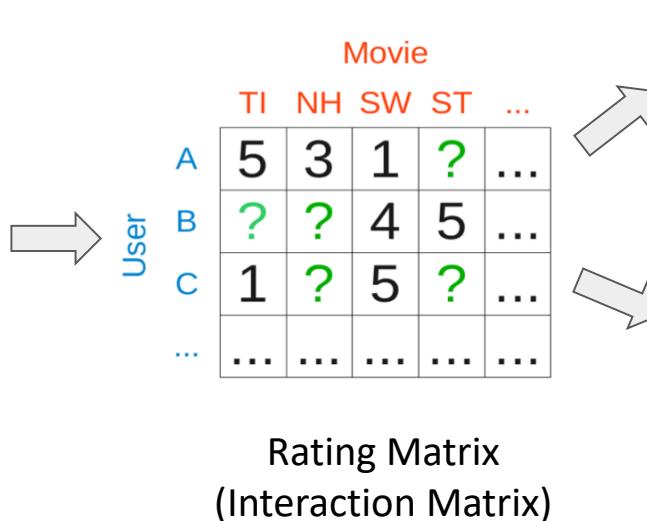
# Collaborative Filtering

- Collaborative Filtering (CF) is the most well-known technique for recommendation.

*“CF makes predictions (**filtering**) about a user’s interest by collecting preferences information from many users (**collaborating**)” ---Wikipedia*

User	Movie	Rating
Alice	Titanic	5
Alice	Notting Hill	3
Alice	Star Wars	1
Bob	Star Wars	4
Bob	Star Trek	5
Charlie	Titanic	1
Charlie	Star Wars	5
...	...	...

Input Tabular data



1. Memory-based CF:  
Predict by **memorizing** similar users’ (or items’) ratings

2. Model-based CF:  
Predict by **inferring** from an underlying model.

# Memory-based CF

Problem: predict user  $u$ 's rating on item  $i$ .

- User-based CF leverages the ratings of  $u$ 's **similar users** on the target item  $i$ .

$$\hat{y}_{ui} = \sum_{u' \in S_u(u)} sim(u, u') \cdot [y_{u'i}]$$

Similar users of  $u$

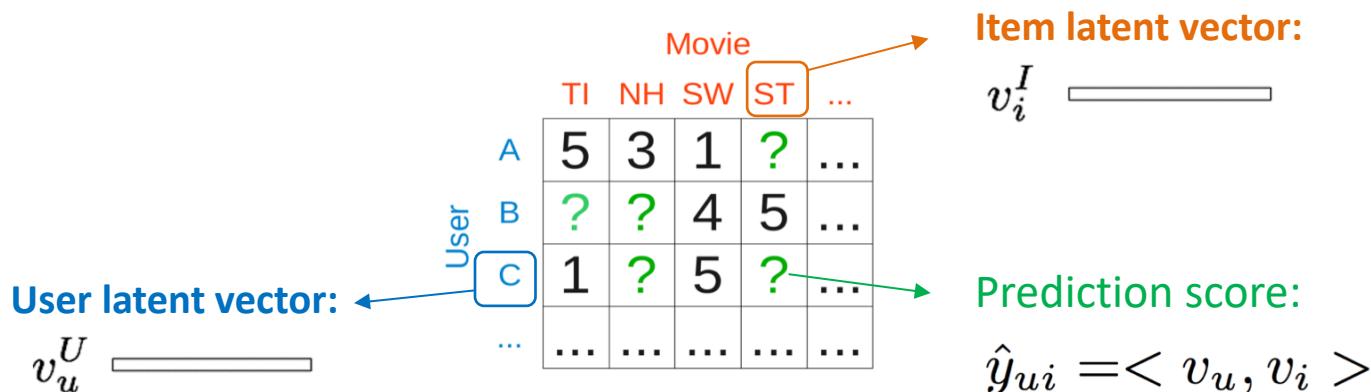
Rating of a similar user on  $i$

	Movie				
	TI	NH	SW	ST	...
A	5	3	1	?	...
B	?	?	4	5	...
C	1	?	5	?	...
...	...	...	...	...	...

- Item-based CF leverages the ratings of  $u$  on other **similar items** of  $i$ .
- $$\hat{y}_{ui} = \sum_{i' \in S_i(i)} sim(i, i') \cdot [y_{ui'}]$$
- Similar items of  $i$
- Rating of  $u$  on a similar item
- Many similarity measures can be used, e.g., Jaccard, Cosine, Pearson Correlation. Recent advance learns the similarity from data.

# Model-based CF

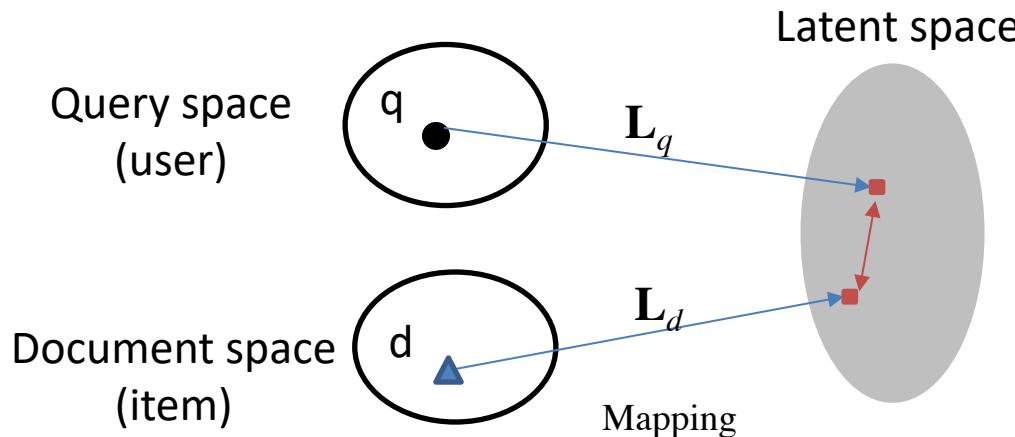
- Matrix Factorization (MF) is the most popular and effective model-based CF method.
- It represents a user and an item as a vector of latent factors.
- The score is estimated as the **inner product** of user latent vector and item latent vector.



- Optimizing a loss to minimize the prediction error on training data can get the latent vectors.

# Convergence of Recommendation and Search Methods

- MF is similar to “Matching in Latent Space” methods in Search!



1. Using one-hot encoding on the ID feature of user and item

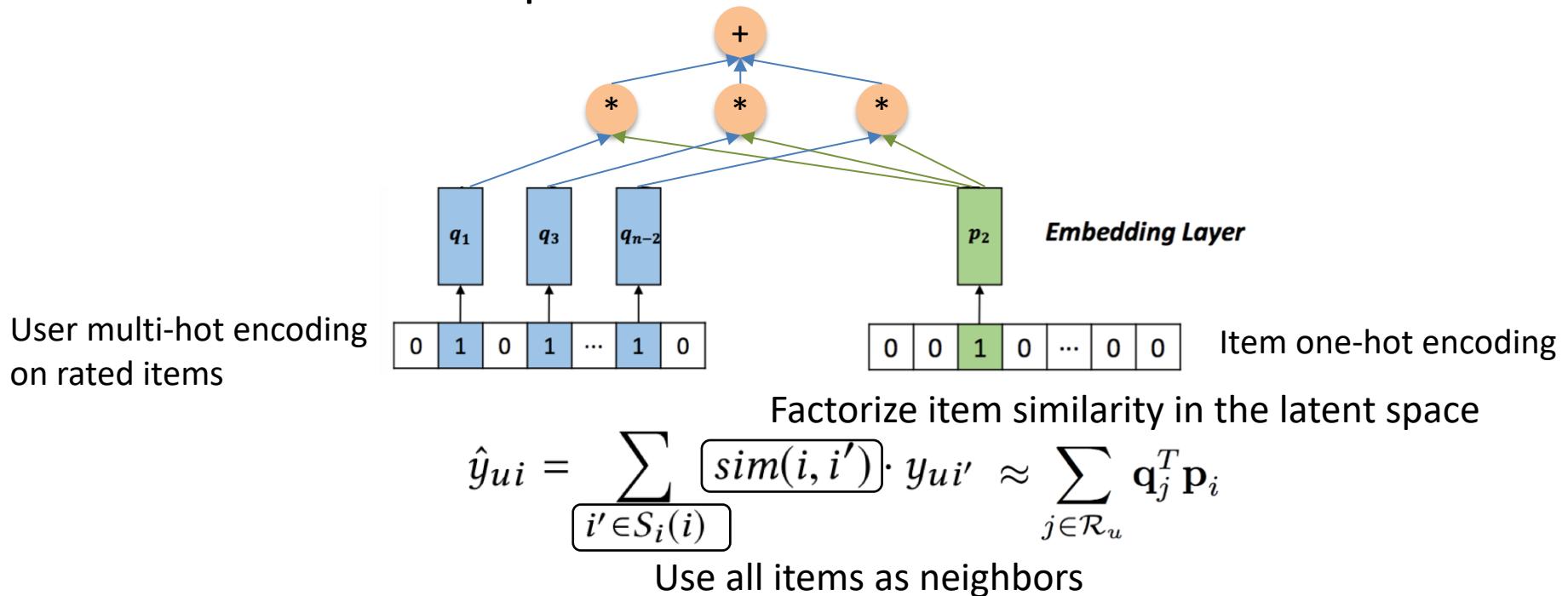


2. Using a linear mapping function, i.e.,  $\mathbf{v}_u = \mathbf{L}_q u$ ,  $\mathbf{v}_i = \mathbf{L}_d i$

3. Using inner product as the matching function in the latent space.

# Item-based CF in Latent Space (Kabbur et al., KDD'14)

- Instead of only using an ID to encode a user, we can make the encoding more meaningful by using the user's rated items.
- This can be interpreted as an item-based CF model.



- Known as the **Factored Item Similarity Model (FISM)** (Kabbur et al, KDD'14)

# Fusing User-based and Item-based CF in Latent Space (Koren, KDD'08)

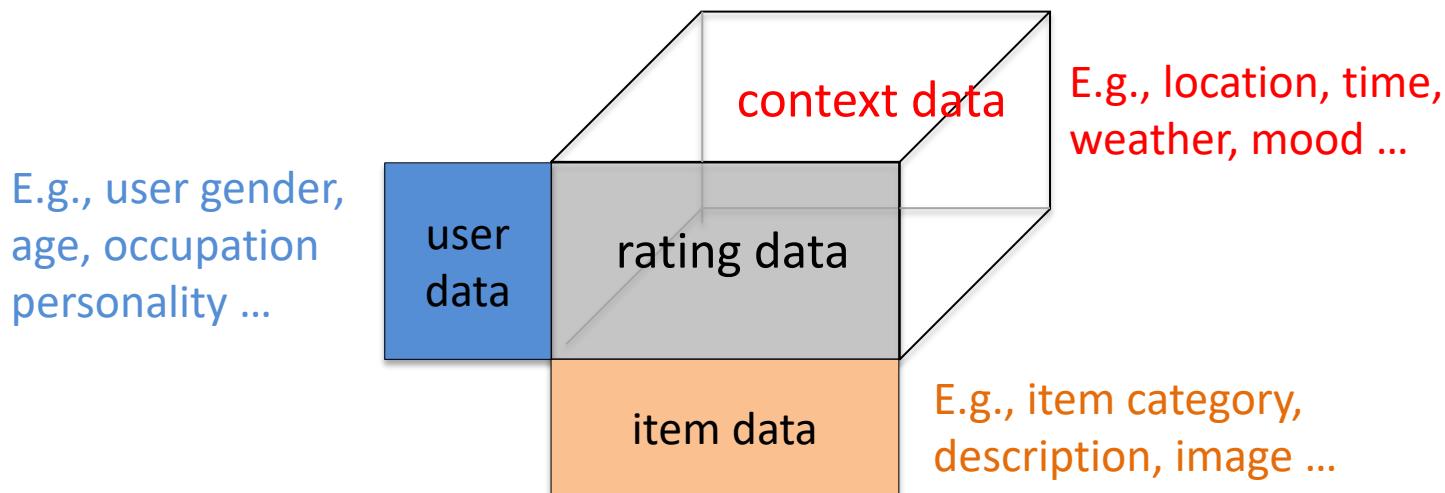
- MF (user-based CF) represents a user as her ID.
  - Directly projecting the ID into latent space
- FISM (item-based CF) represents a user as her interacted items.
  - Projecting interacted items into latent space
- SVD++ fuses the two types of models in the latent space:

$$\hat{y}_{ui} = (\mathbf{v}_u + \underbrace{\sum_{j \in \mathcal{R}_u} \mathbf{q}_j}_{\text{User representation in latent space}})^T \mathbf{p}_i$$

- This is the best single model for rating prediction in the Netflix challenge.

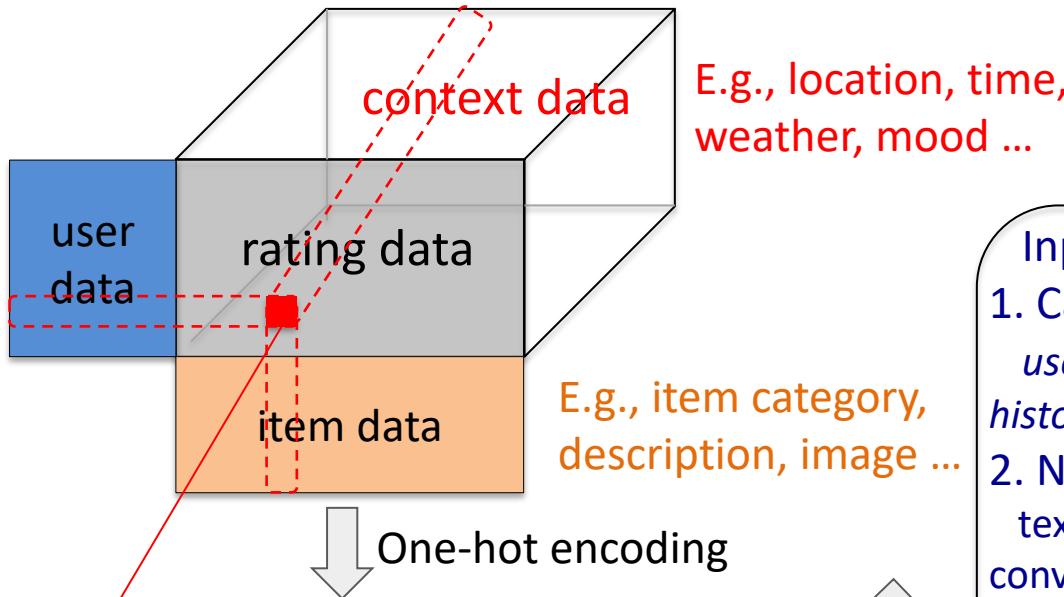
# Feature-based Recommendation

- CF utilizes only the interaction matrix only to build the predictive model.
- How about other information like user/item attributes and contexts?
- Example data used for building a RecSys:



# Feature-based Recommendation

E.g., user gender,  
age, occupation  
personality ...



Each row encodes all info for a rating

	Feature vector $\mathbf{x}$				Target $y$			
	$x^{(1)}$	$x^{(2)}$	$x^{(3)}$	$x^{(4)}$	$x^{(5)}$	$x^{(6)}$	$x^{(7)}$	
	1 0 0 ...	1 0 0 0 ...	0.3 0.3 0.3 0 ...	5 $y^{(1)}$	3 $y^{(2)}$	1 $y^{(3)}$	4 $y^{(4)}$	5 $y^{(5)}$
$x^{(1)}$	1 0 0 ...	0 1 0 0 ...	0.3 0.3 0.3 0 ...	3 $y^{(6)}$	1 $y^{(7)}$			
$x^{(2)}$	1 0 0 ...	0 0 1 0 ...	0.3 0.3 0.3 0 ...					
$x^{(3)}$	1 0 0 ...	0 0 1 0 ...	0.3 0.3 0.3 0 ...					
$x^{(4)}$	0 1 0 ...	0 0 1 0 ...	0 0 0.5 0.5 ...					
$x^{(5)}$	0 1 0 ...	0 0 0 1 ...	0 0 0.5 0.5 ...					
$x^{(6)}$	0 0 1 ...	1 0 0 0 ...	0.5 0 0.5 0 ...					
$x^{(7)}$	0 0 1 ...	0 0 1 0 ...	0.5 0 0.5 0 ...					
	A B C ...	TI NH SW ST ...	TI NH SW ST ...					
	User	Movie	Other Movies rated					

Input Features:

1. Categorical features:  
*user/item ID, bag-of-words, historical features...*
2. Numerical features:  
*textual/visual embeddings, converted features (e.g. TFIDF, GBDT)...*

Sparse Predictive Model

# FM: Factorization Machine (Rendle, ICDM'10)

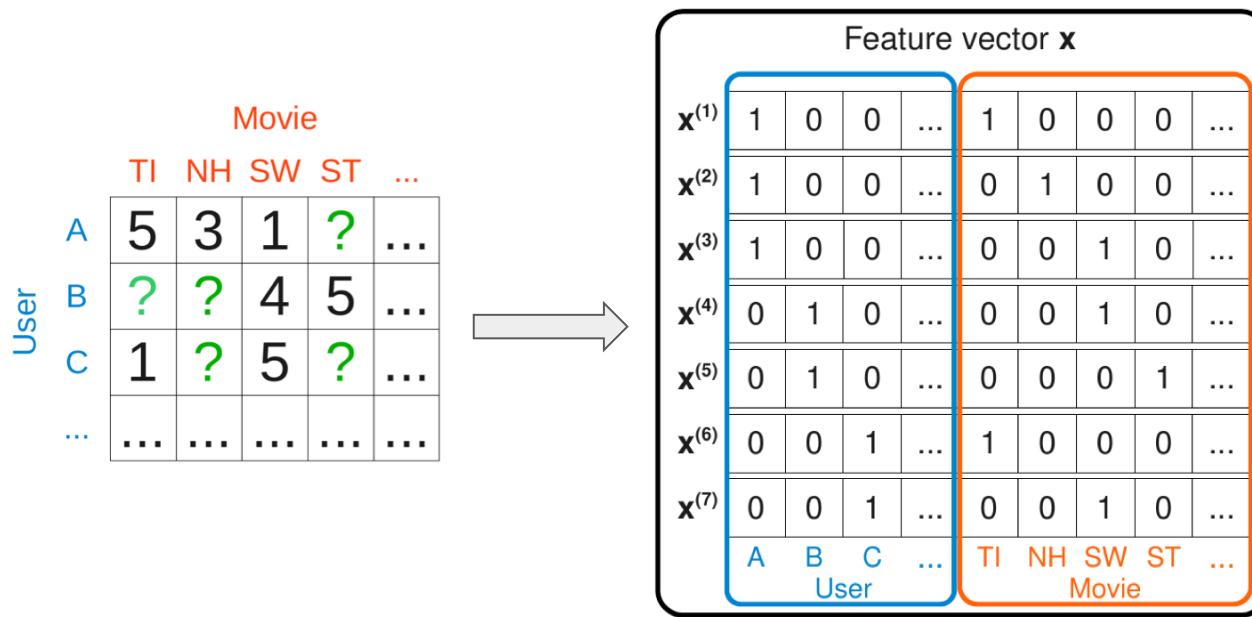
- FM is inspired from previous factorization models
- It represents each feature as a latent vector (embedding), and models the second-order feature interactions:

$$\hat{y}(\mathbf{x}) = w_0 + \underbrace{\sum_{i=1}^p w_i x_i}_{\text{First-order: Linear Regression}} + \underbrace{\sum_{i=1}^p \sum_{j>i}^p \langle \mathbf{v}_i, \mathbf{v}_j \rangle x_i x_j}_{\text{Second-order: pair-wise interactions between features}}$$

- Note: self-interaction is not included:  $\langle \mathbf{v}_i, \mathbf{v}_i \rangle$ .
- FM allows easy feature engineering for recommendation, and can mimic many existing models (that are designed for a specific task) by inputting different features.
  - E.g., MF, SVD++, timeSVD (Koren, KDD'09), PIFT (Rendle, WSDM'10) etc.

# Matrix Factorization with FM

- Input: 2 variables <user (ID), item (ID)>.

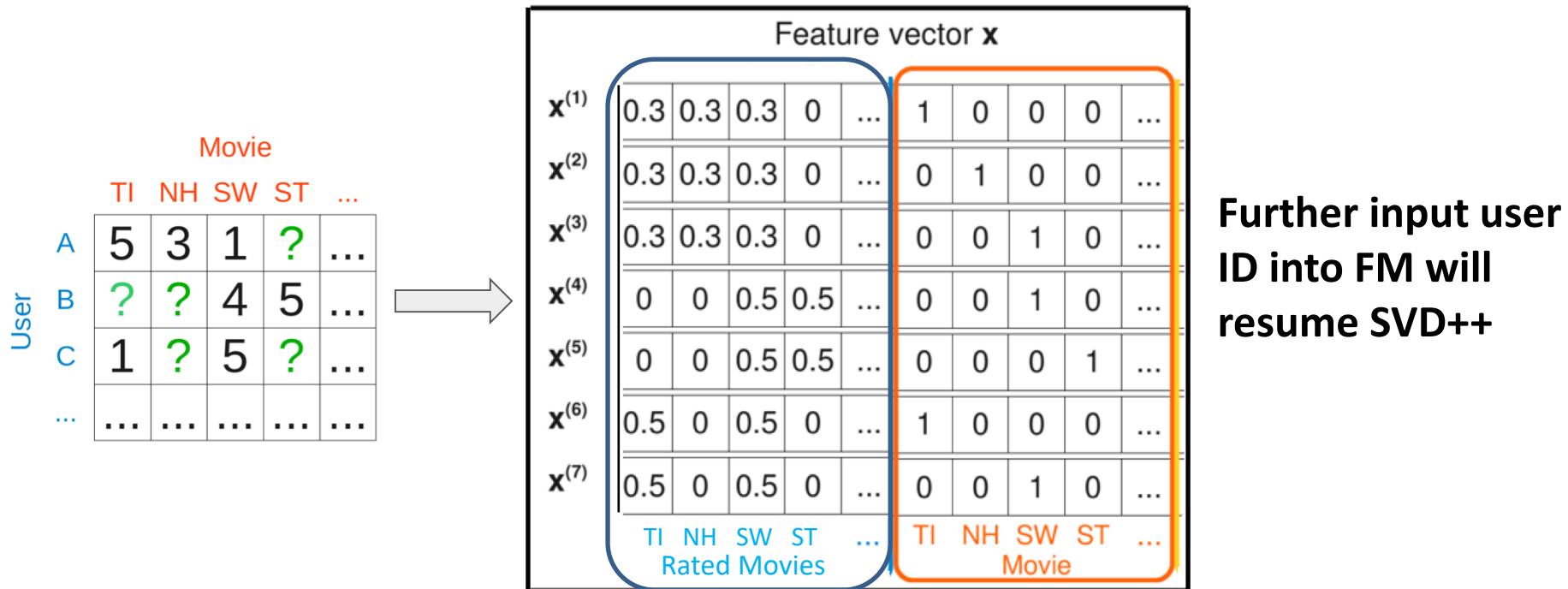


With this input, FM is identical to MF with bias:

$$\hat{y}(\mathbf{x}) = w_0 + w_u + w_i + \underbrace{\langle \mathbf{v}_u, \mathbf{v}_i \rangle}_{\text{MF}}$$

# Factored Item Similarity Model with FM

- Input: 2 variables <user (historical items ID), item (ID)>.



With this input, FM subsumes FISM with additional terms:

$$\hat{y}(\mathbf{x}) = bias + \underbrace{\sum_{j \in \mathcal{R}_u} \langle \mathbf{v}_j, \mathbf{v}_i \rangle}_{\text{FISM}} + \sum_{j \in \mathcal{R}_u, j' > j} \langle \mathbf{v}_j, \mathbf{v}_{j'} \rangle$$

FISM

# Learning Recommender Models

- For ranking novel items for a user (i.e., top-K recommendation), it is crucial to account for **the missing data** (negative signal)

3 common loss functions (for a user  $u$ ):

» 1. Pointwise Regression Loss (explicit & implicit data):

$$L_u = \sum_{i \in \mathcal{R}_u} (y_{ui} - \hat{y}_{ui})^2 + w_0 \sum_{j \in \mathcal{R}_u^-} (0 - \hat{y}_{uj})^2 + \lambda_\Theta \|\Theta\|^2 \quad (\text{Bayer et al, WWW'17})$$

» 2. Pointwise Classification Loss (implicit data):

$$L_u = \sum_{i \in \mathcal{R}_u} \log \sigma(\hat{y}_{ui}) + w_0 \sum_{j \in \mathcal{R}_u^-} \log(1 - \sigma(\hat{y}_{uj})) + \lambda_\Theta \|\Theta\|^2 \quad (\text{He et al, WWW'17})$$

» 3. Pairwise Classification loss (implicit data):

$$L_u = - \sum_{i \in \mathcal{R}_u} \sum_{j \in \mathcal{R}_u^-} \log(\sigma(\hat{y}_{ui} - \hat{y}_{uj})) + \lambda_\Theta \|\Theta\|^2 \quad (\text{Rendle et al., UAI'09})$$

**$L_2$  regularizer must be tuned to prevent overfitting.**

User	Movie				
	TI	NH	SW	ST	...
A	5	3	1	?	...
B	?	?	4	5	...
C	1	?	5	?	...
...	...	...	...	...	...

# References

- [https://en.wikipedia.org/wiki/Collaborative\\_filtering](https://en.wikipedia.org/wiki/Collaborative_filtering)
- Xiangnan He, Hanwang Zhang, Min-Yen Kan, and Tat-Seng Chua. Fast matrix factorization for online recommendation with implicit feedback. In SIGIR 2016.
- Yehuda Koren, and Robert Bell. Advances in collaborative filtering. Recommender systems handbook. Springer, Boston, MA, 2015. 77-118.
- Santosh Kabbur, Xia Ning, and George Karypis. Fism: factored item similarity models for top-n recommender systems. In KDD 2013.
- Yehuda Koren. Factorization meets the neighborhood: a multifaceted collaborative filtering model. In KDD 2018.
- Steffen Rendle. Factorization machines. In ICDM 2010.
- Yehuda Koren. Collaborative filtering with temporal dynamics. Communications of the ACM 53, no. 4 (2010): 89-97.
- Steffen Rendle, and Lars Schmidt-Thieme. Pairwise interaction tensor factorization for personalized tag recommendation. In WSDM 2010.
- Immanuel Bayer, Xiangnan He, Bhargav Kanagal, and Steffen Rendle. A generic coordinate descent framework for learning from implicit feedback. In WWW 2017.
- Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. Neural collaborative filtering. In WWW 2017.
- Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. BPR: Bayesian personalized ranking from implicit feedback. In UAI 2009.

# Outline of Tutorial

- Unified view of matching in search and recommendation
- Part 1: Traditional Approaches to Matching
- Part 2: Deep Learning Approaches to Matching
  - Overview
  - Deep matching models for search
  - Deep matching models for recommendation
- Summary

# Growing Interests in “Deep Matching”

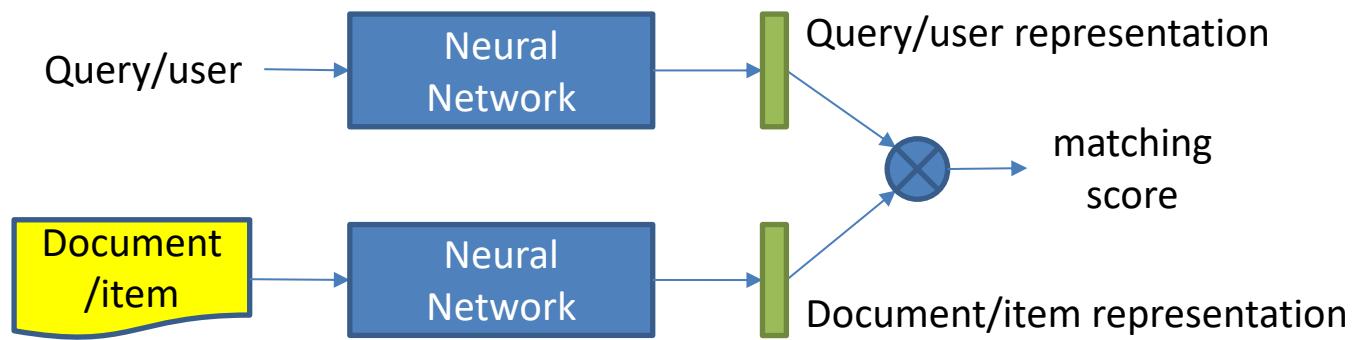
- Success of deep learning in other fields
  - Speech recognition, computer vision, and natural language processing
- Growing presence of deep learning in IR research
  - SIGIR 2016 keynote, Tutorial, and Neu-IR workshop
- Adopted by industry
  - ACM News: *Google Turning its Lucrative Web Search Over 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

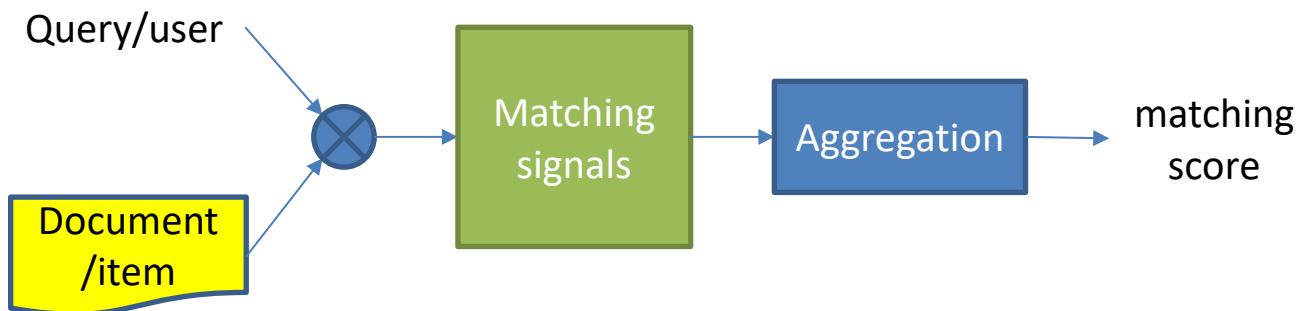
- Representation
  - Word: one hot —> distributed
  - Sentence: bag-of-words —> distributed representation
  - Better representation ability, better generalization ability
- Matching function
  - Inputs (features): handcrafted —> automatically learned
  - Function: simple functions (e.g., cosine, dot product) —> neural networks (e.g., MLP, neural tensor networks)
  - Involving richer matching signals
  - Considering soft matching patterns

# Deep Learning Paradigms for Matching

- Methods of representation learning

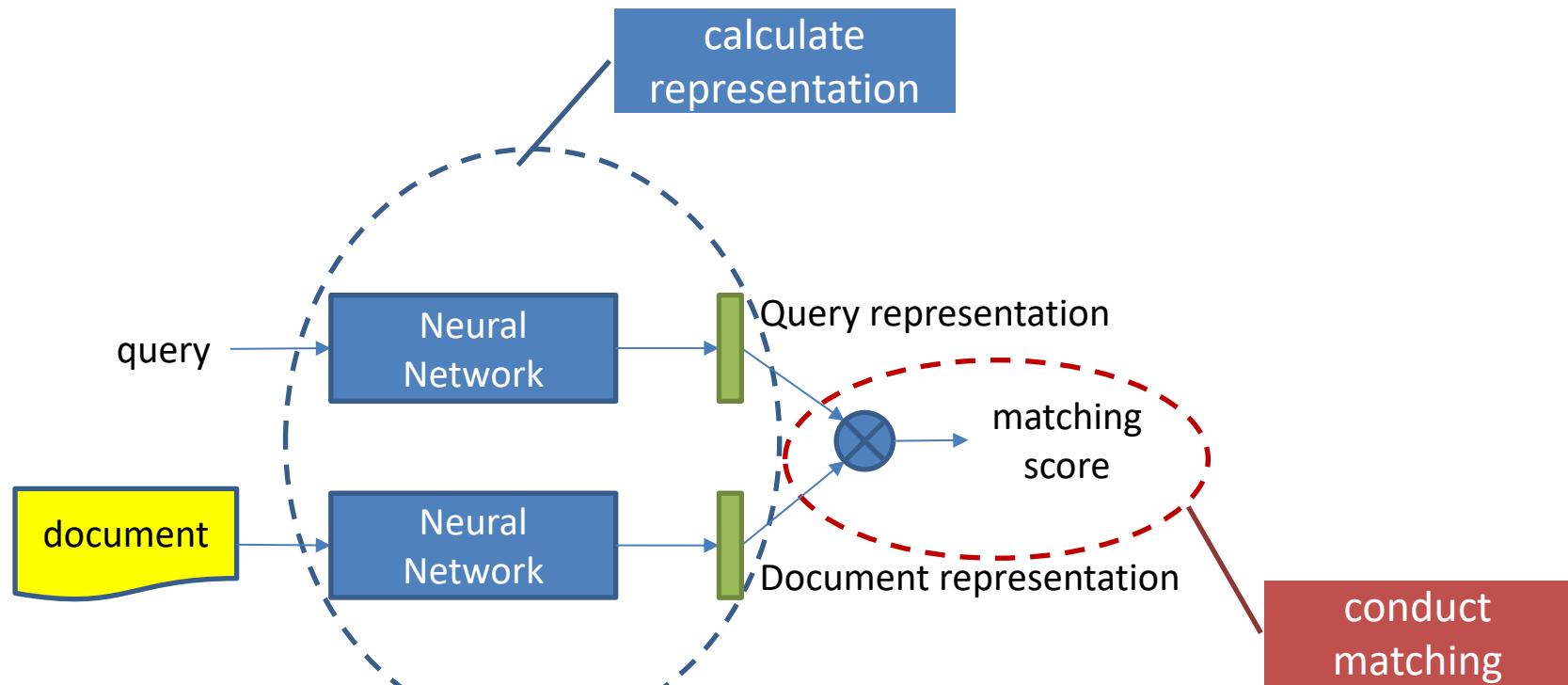


- Methods of matching function learning



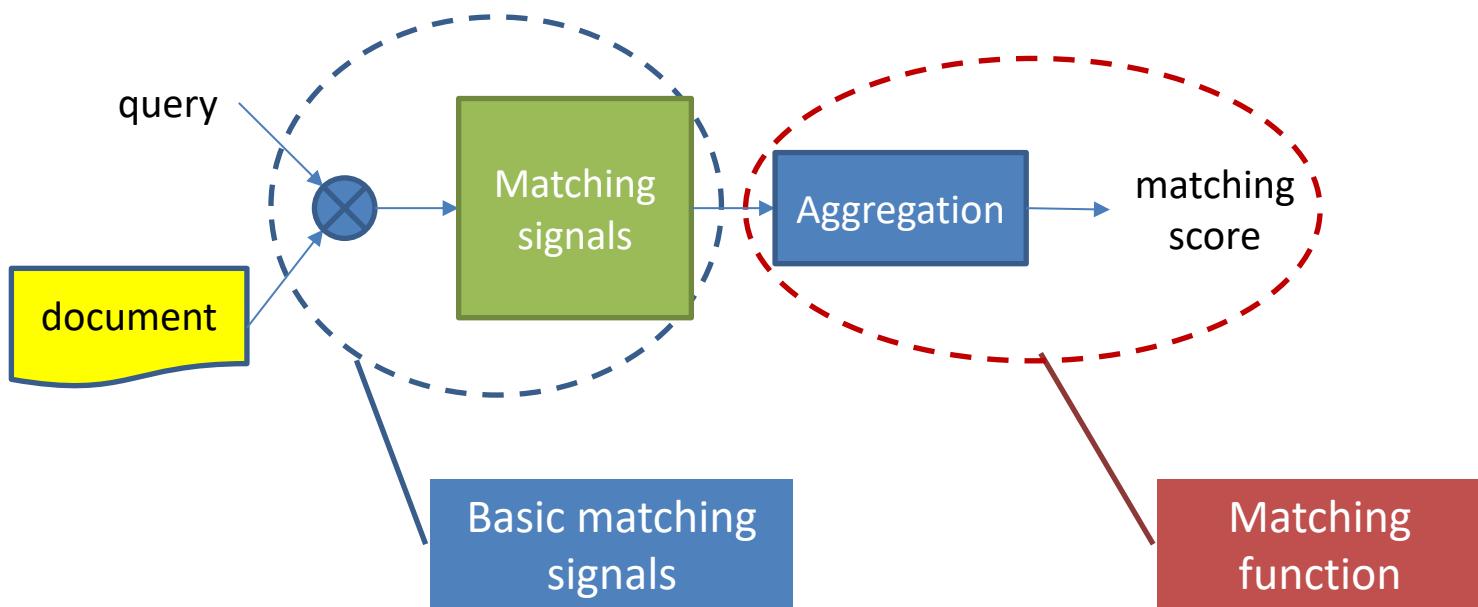
# Methods of Representation Learning

- Step 1: calculate representation  $\phi(x)$
- Step 2: conduct matching  $F(\phi(x), \phi(y))$



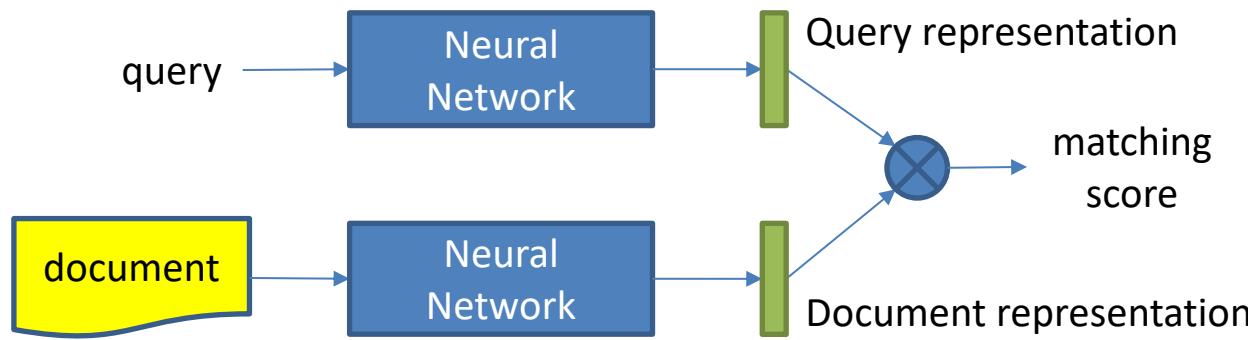
# Methods of Matching Function Learning

- Step 1: construct basic low-level matching signals
- Step 2: aggregate matching patterns



# Outline of Tutorial

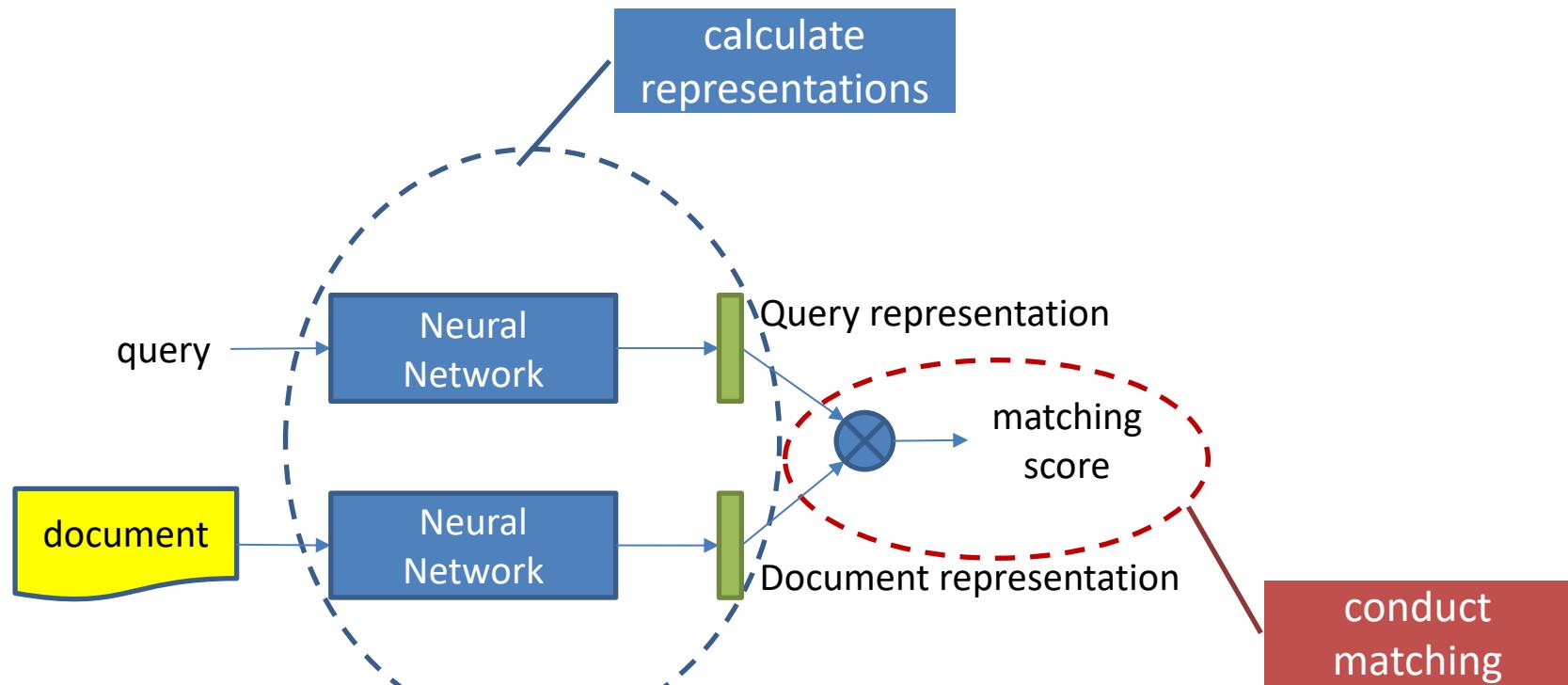
- Unified view of matching in search and recommendation
- Part 1: Traditional Approaches to Matching
- Part 2: Deep Learning Approaches to Matching
  - Overview
  - Deep matching models for search
  - Deep matching models for recommendation
- Summary



# METHODS OF REPRESENTATION LEARNING

# Representation Learning for Query-Document Matching

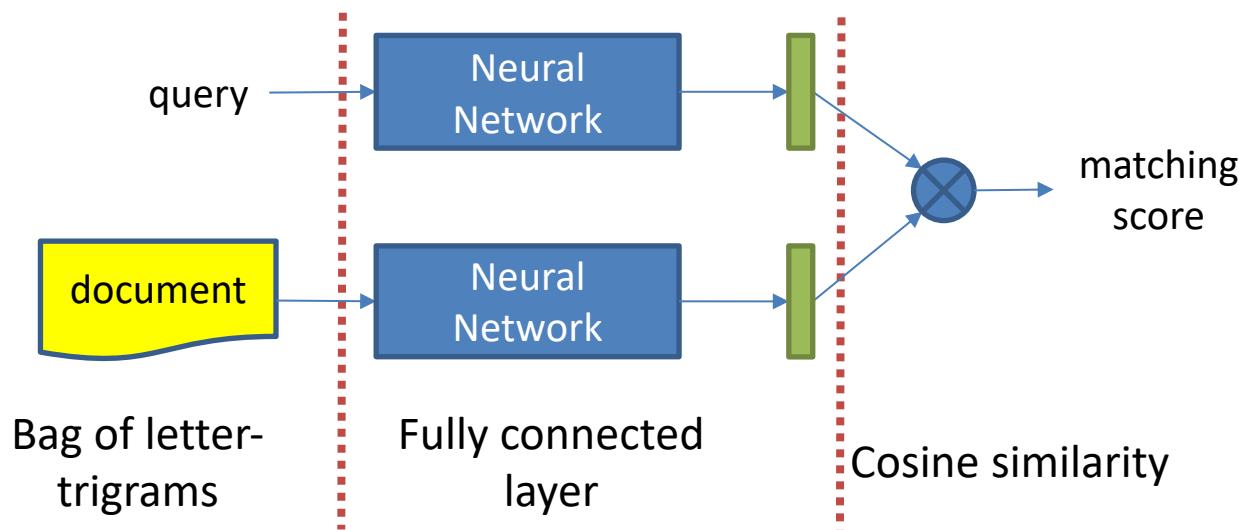
- Step 1: calculate query and document representation
- Step 2: conduct query-document matching



# Typical Methods of Representation Learning for Matching

- 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)



- Bag-of-words representation
  - “candy store”: [0, 0, 1, 0, ..., 1, 0, 0]
- Bag of letter-trigrams representation
  - “#candy# #store#” --> #ca can and ndy dy# #st sto tor ore re#
  - Representation: [0, 1, 0, 0, 1, 1, 0, ..., 1]
- Advantages of using bag of letter-trigrams
  - Reduce vocabulary: #words 50K → # letter-trigram: 30K
  - Generalize to unseen words
  - Robust to misspelling, inflection etc.

# DSSM Matching Function

- Cosine similarity between semantic vectors

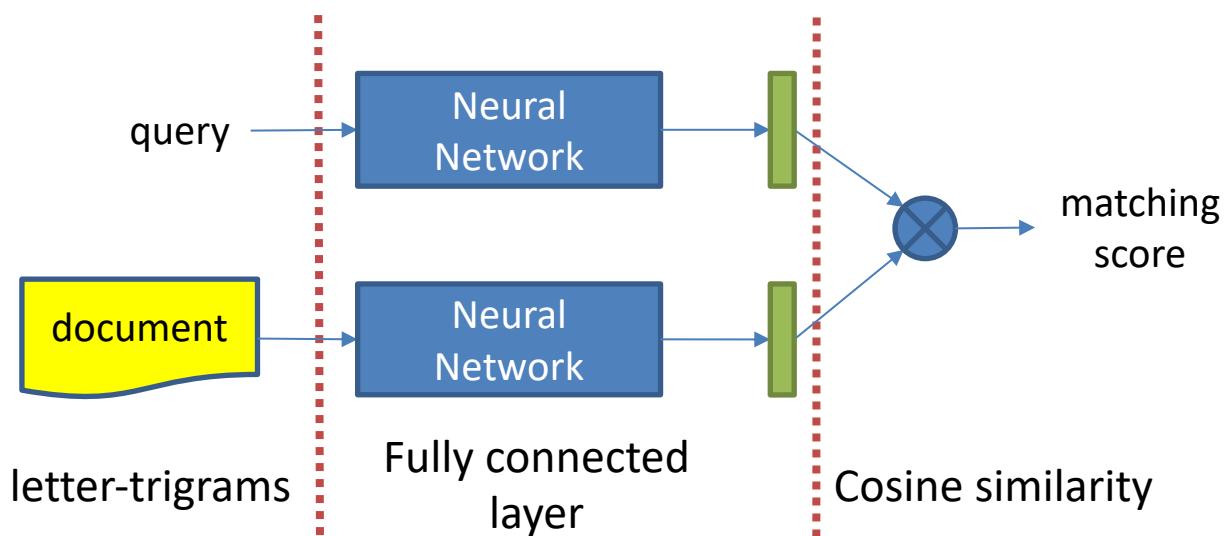
$$S = \frac{x^T \cdot y}{|x| \cdot |y|}$$

- 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:

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

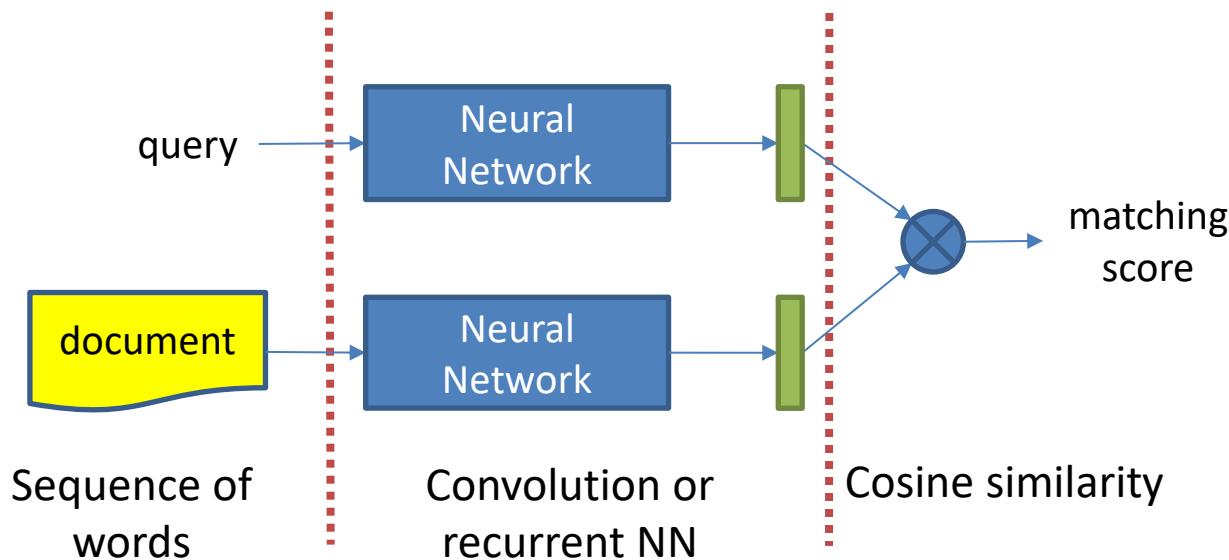
# DSSM: Brief Summary

- **Inputs:** Bag of letter-trigrams as input for improving the scalability and generalizability
- **Representations:** mapping sentences to vectors with DNN: semantically similar sentences are close to each other
- **Matching:** cosine similarity as the matching function
- **Problem:** *the order information of words is missing* (bag of letter-trigrams cannot keep the word order information)



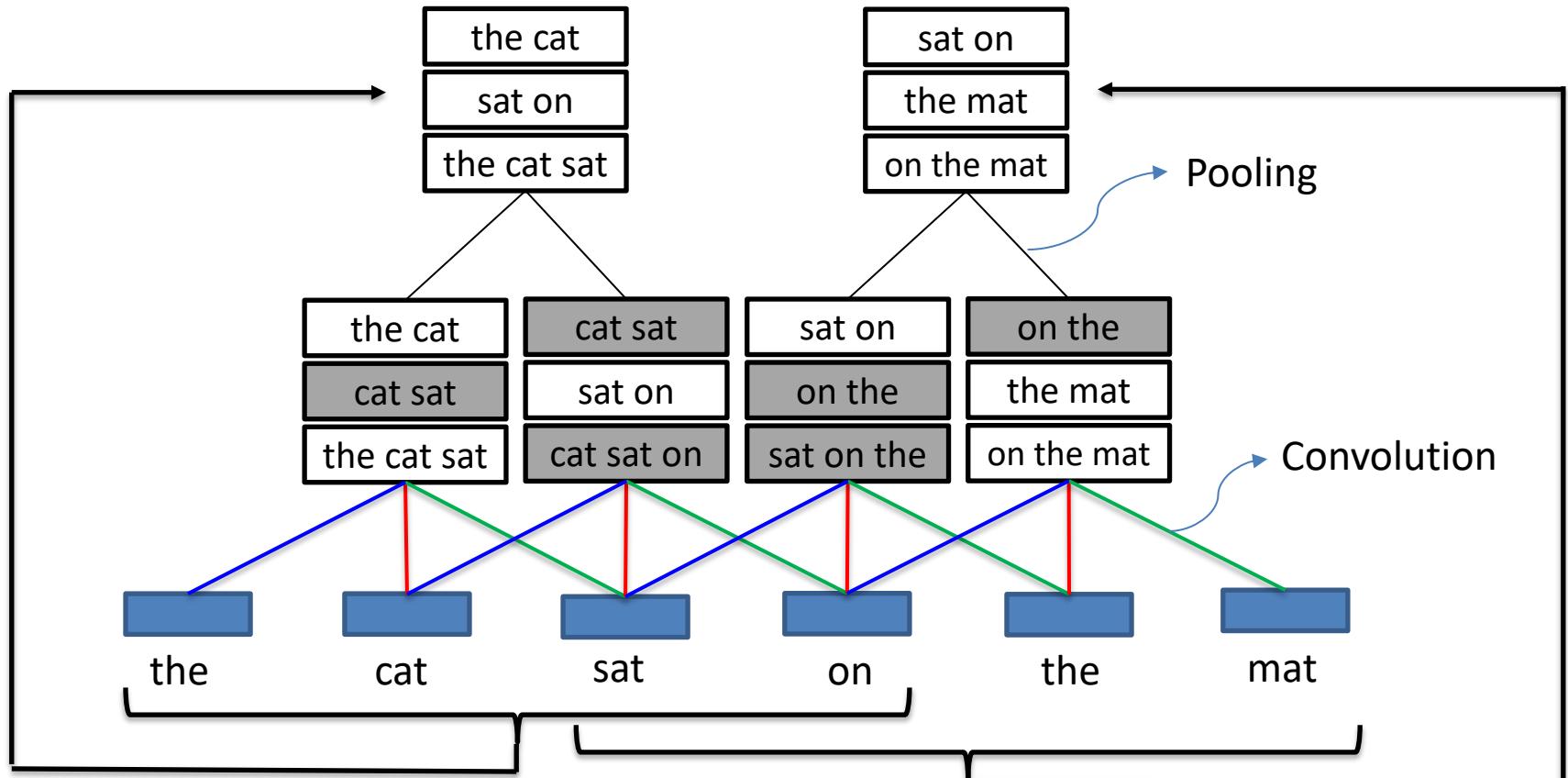
# How to Capture Order Information?

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



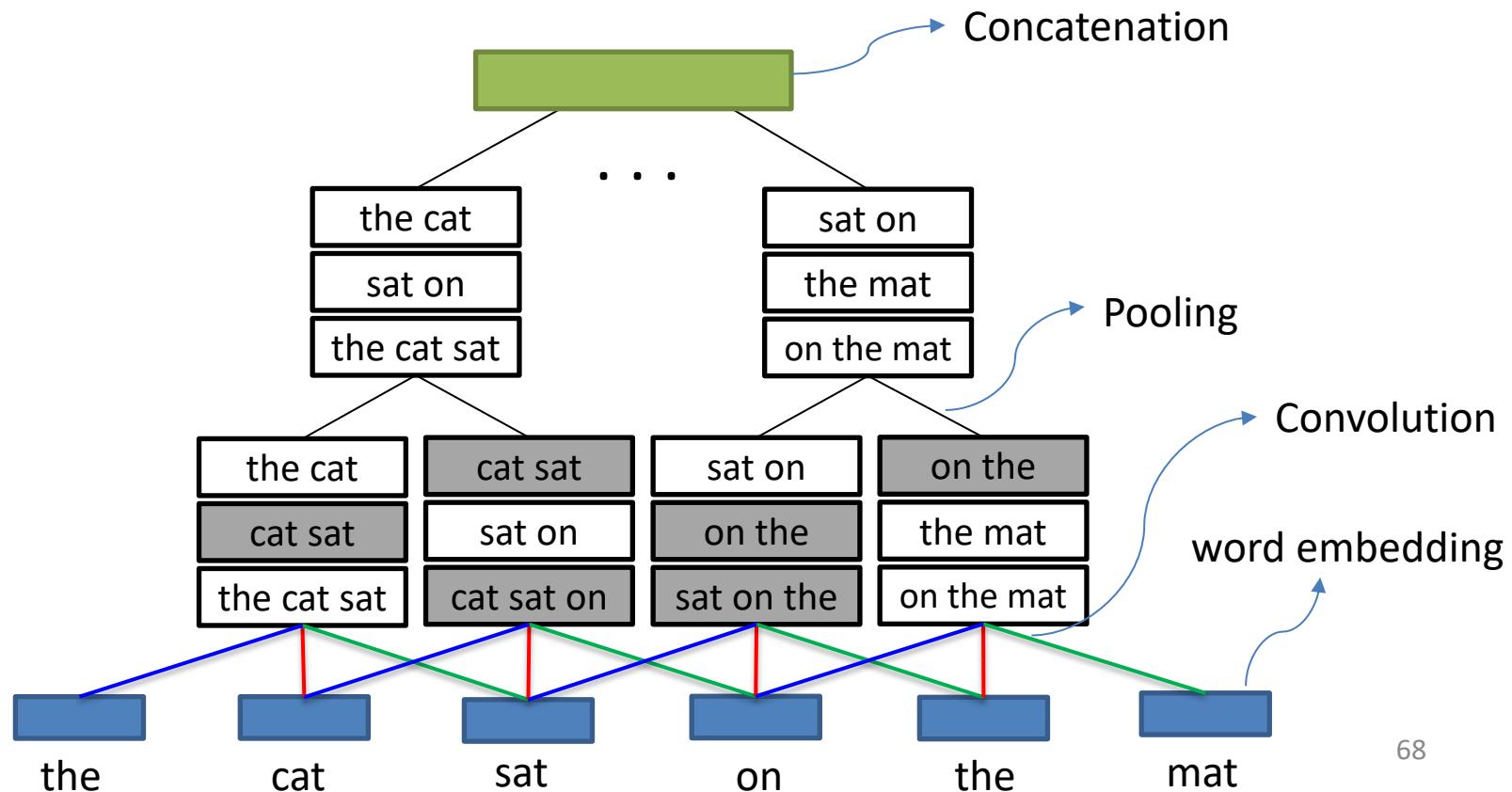
# CNN can Keep the Order Information

1-D convolution and pooling operations can keep the word order information



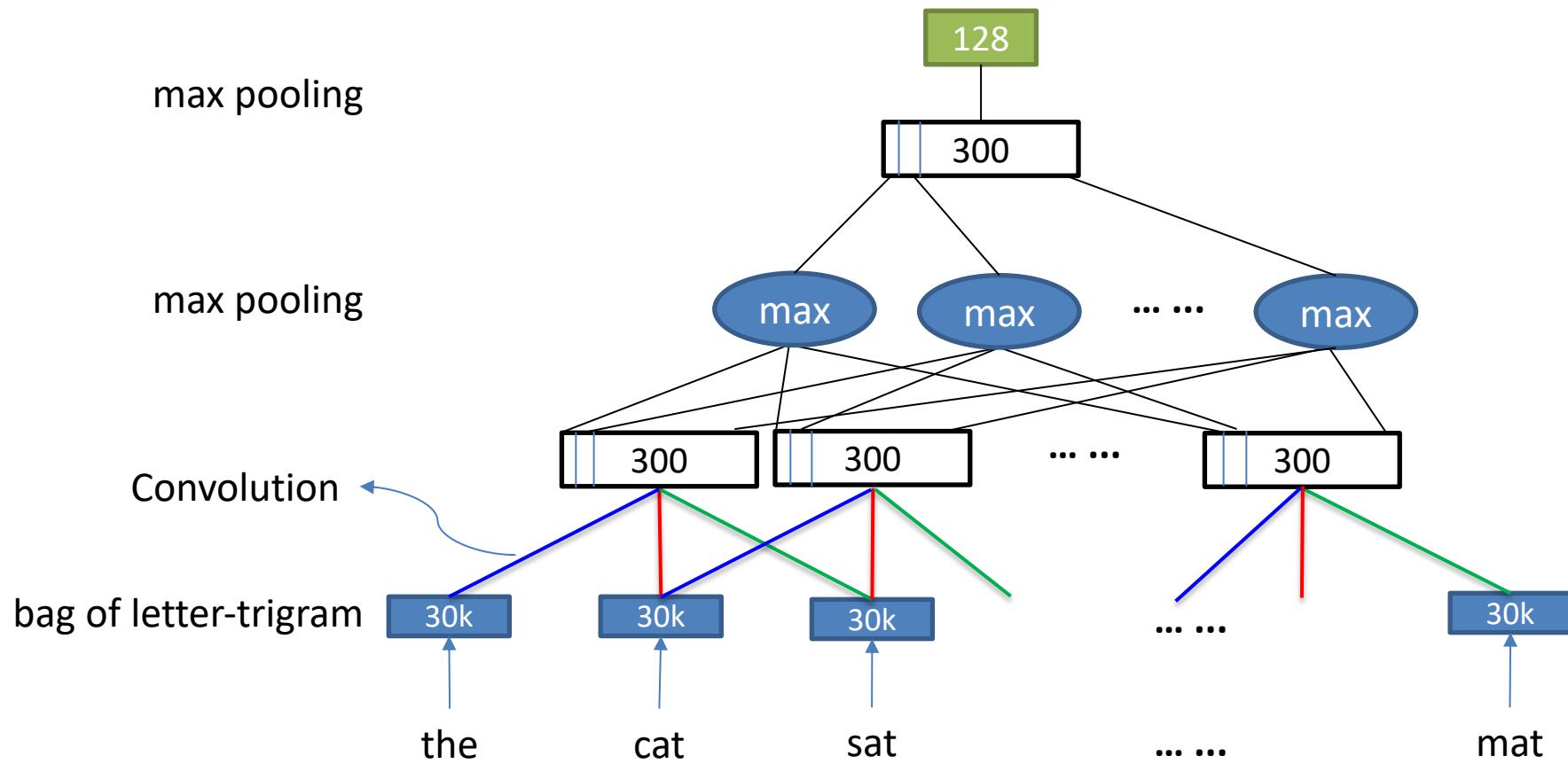
# Using CNN: ARC-I (Hu et al., 2014) and CNTN (Qiu et al., 2015)

- Input: sequence of word embeddings trained on a large dataset
- Model: the convolutional operation in CNN compacts each **sequence of  $k$  words**

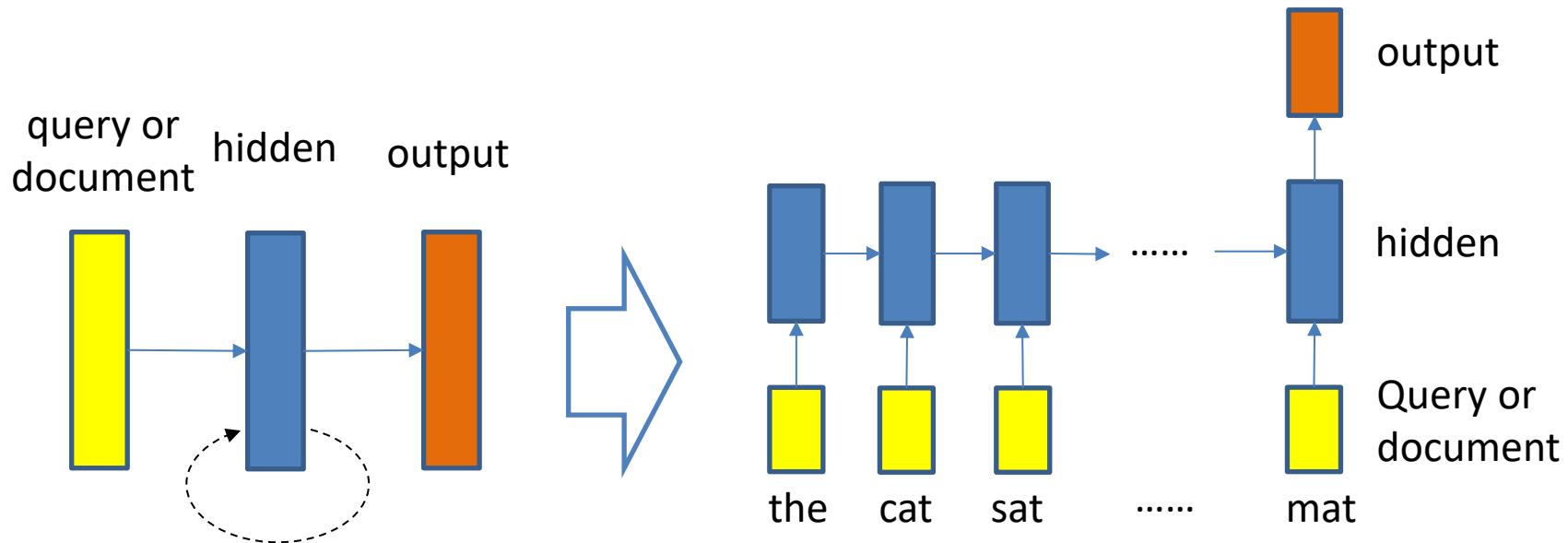


# Using CNN: CDSSM (Shen et al., '14)

The convolutional operation in CNN compacts **each sequence of  $k$  words**



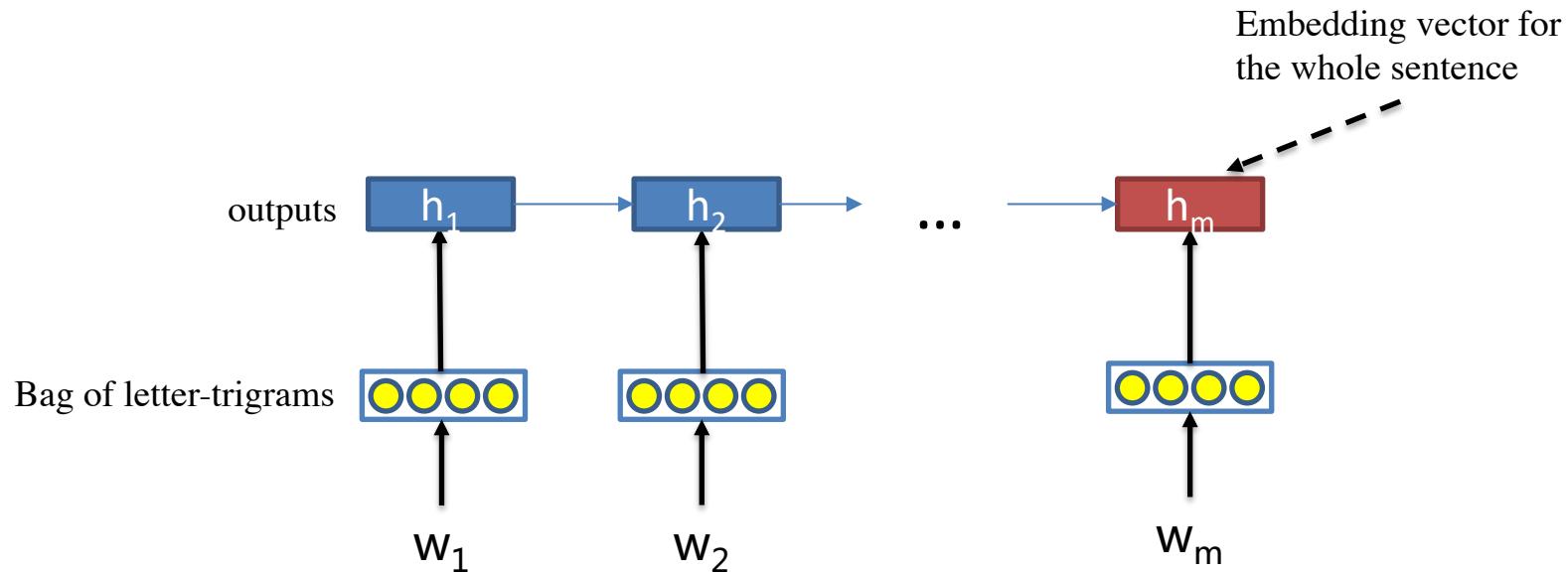
# RNN can Keep the Order Information



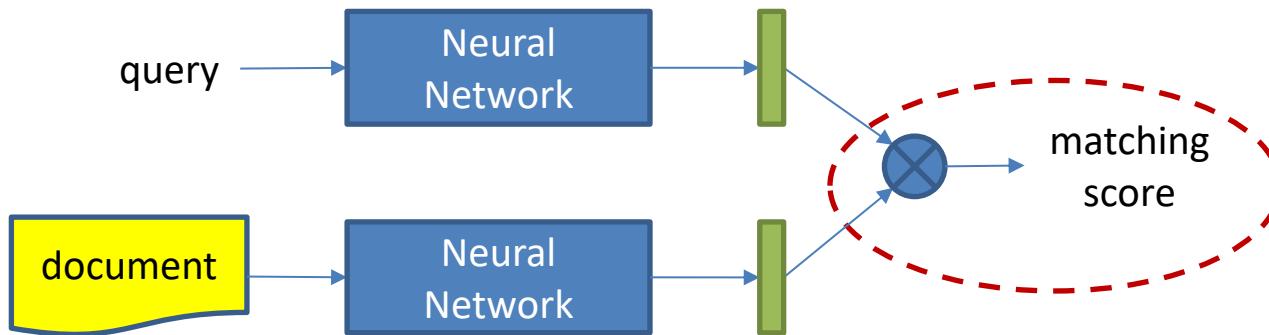
- RNNs implement dynamical systems
- RNNs can approximate arbitrary dynamical systems with arbitrary precision
- Two popular variations: long-short term memory (LSTM) and gated recurrent unit (GRU)

# Using RNN: LSTM-RNN (Palangi et al., '16)

- 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 representations of query and document :  $q$  and  $d$
- Similarity between these two representations:
  - Cosine Similarity (DSSM, CDSSM, RNN-LSTM)
  - Dot Product

$$s = \frac{q^T \cdot d}{|q| \cdot |d|}$$

$$s = q^T \cdot d$$

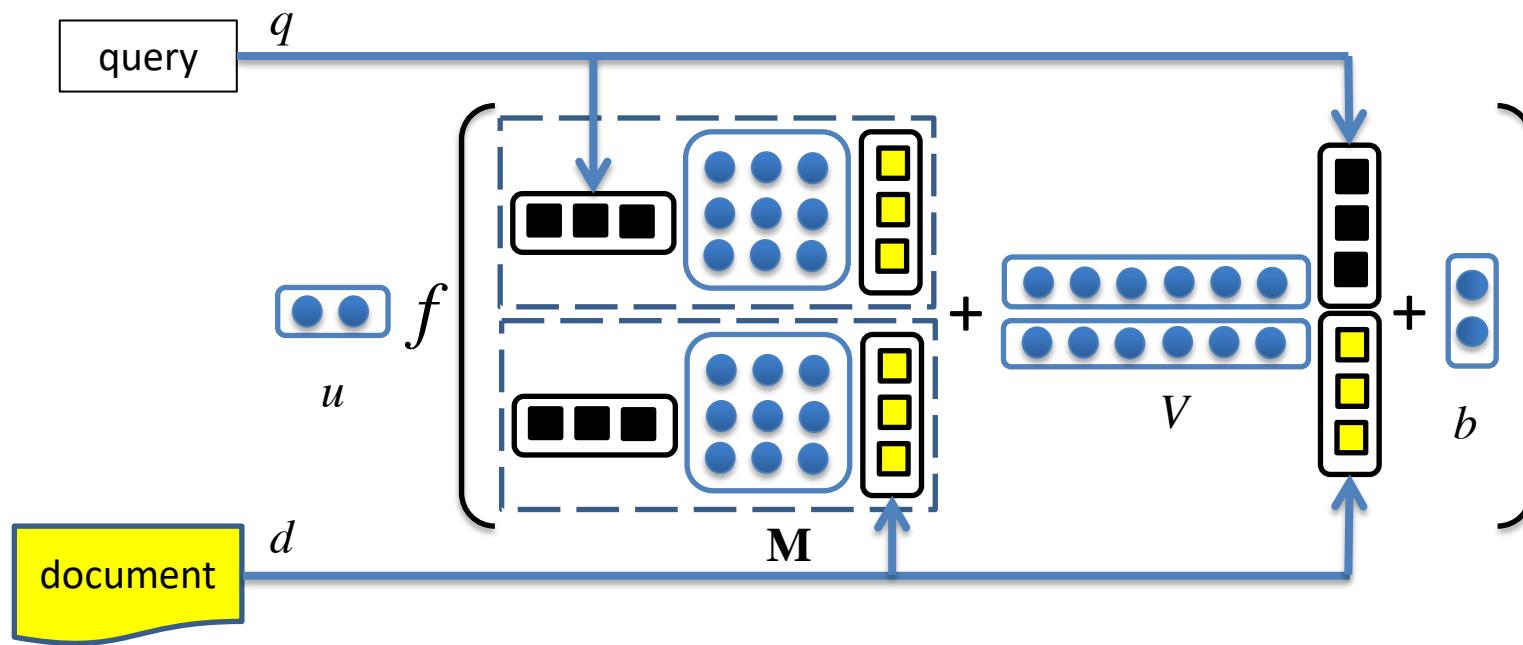
- Multi-Layer Perception (ARC-I)

$$s = W_2 \cdot \sigma \left( W_1 \cdot \begin{bmatrix} q \\ d \end{bmatrix} + b_1 \right) + b_2$$

# Matching Functions (cont')

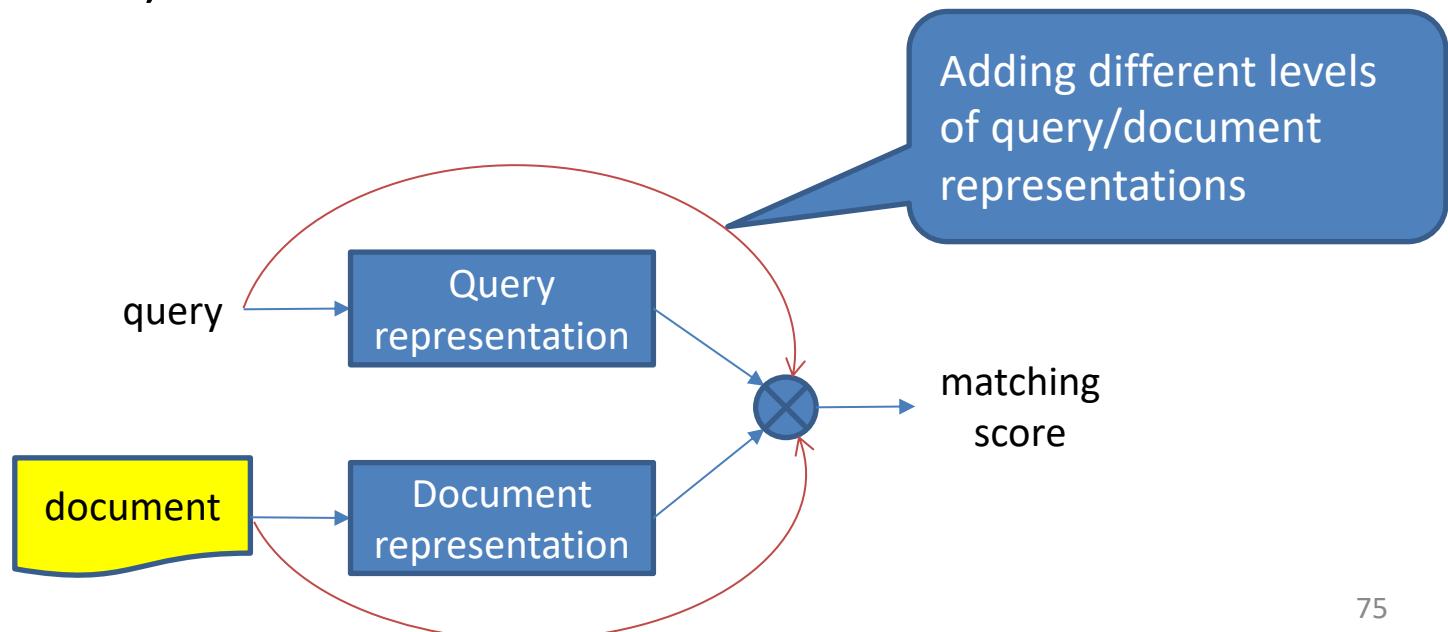
- Neural Tensor Networks (CNTN) (Qiu et al., '15)

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



# Extensions to Representation Learning Methods

- Problem: representations are too coarse to conduct text match
  - Experience in IR: combining topic-level and word-level matching signals usually achieve better performances
- Solution: add fine-grained signals, include MultGranCNN(Yin et al., ACL 2015), U-RAE (Socher et al., NIPS 2011), MV-LSTM (Wan et al., AAAI 2016)

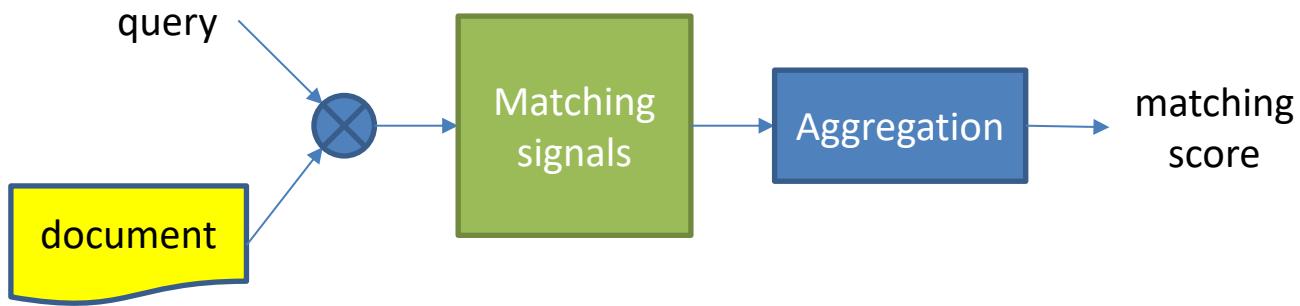


# Experimental Results

	Model	P@1	MRR
Traditional methods	BM25	0.579	0.726
	ARC-I	0.581	0.756
	CNTN	0.626	0.781
Representation learning for matching	LSTM-RNN	0.690	0.822
	uRAE	0.398	0.652
	MultiGranCNN	0.725	0.840
	MV-LSTM	0.766	0.869

Based on Yahoo! Answers dataset (60,564 question-answer pairs)

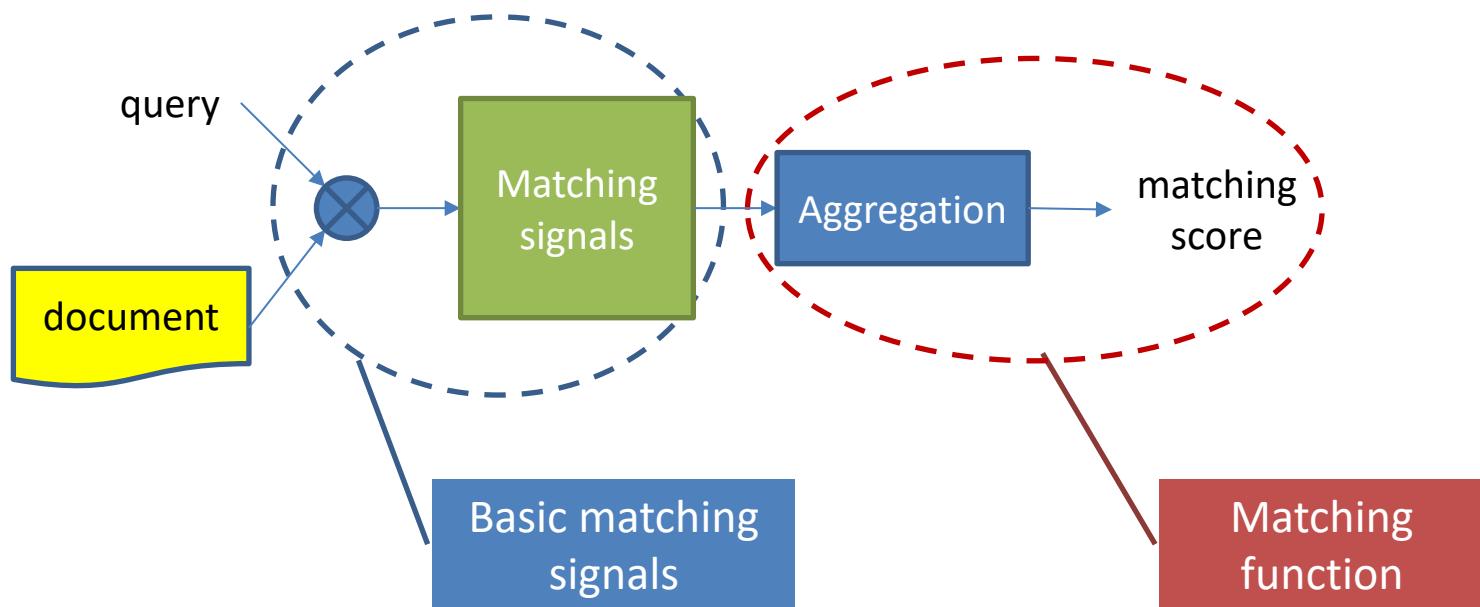
- Representation learning methods outperformed baselines
  - Semantic representation is important
- LSTM-RNN performed better than ARC-I and CNTN
  - Modeling the order information does help
- MultiGranCNN and MV-LSTM are the best performing methods
  - Fine-grained matching signals are useful



# METHODS OF MATCHING FUNCTION LEARNING

# Matching Function Learning

- Step 1: construct basic low-level matching signals
- Step 2: aggregate matching patterns

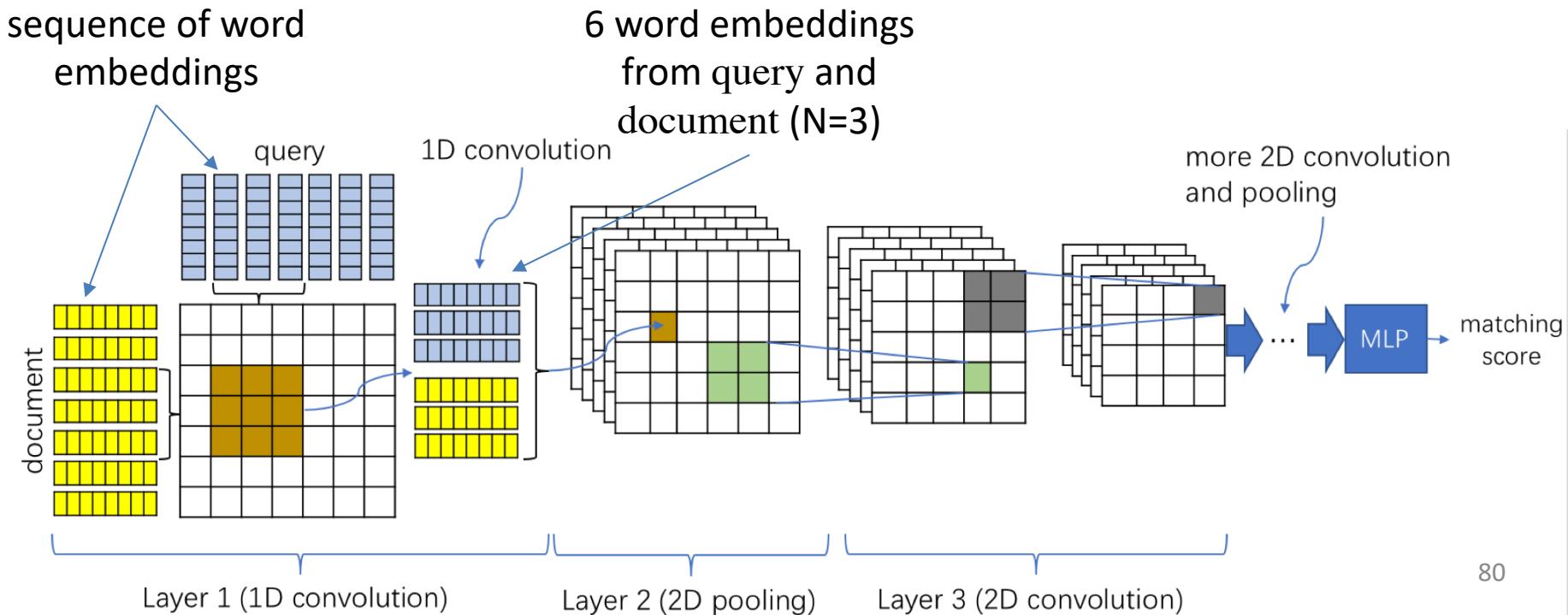


# Typical Matching Function Learning Methods

- Matching with query-document matching matrix :
  - ARC II (Hu et al., NIPS'14)
  - MatchPyramid (Pang et al. AAAI'16)
  - Match-SRNN (Wan et al. IJCAI'16)
  - K-NRM (Xiong et al., SIGIR 2017)
  - Conv-KNRM (Dai et al., WSDM 2018)
- Matching with attention model (Parikh et al., EMNLP 2016)

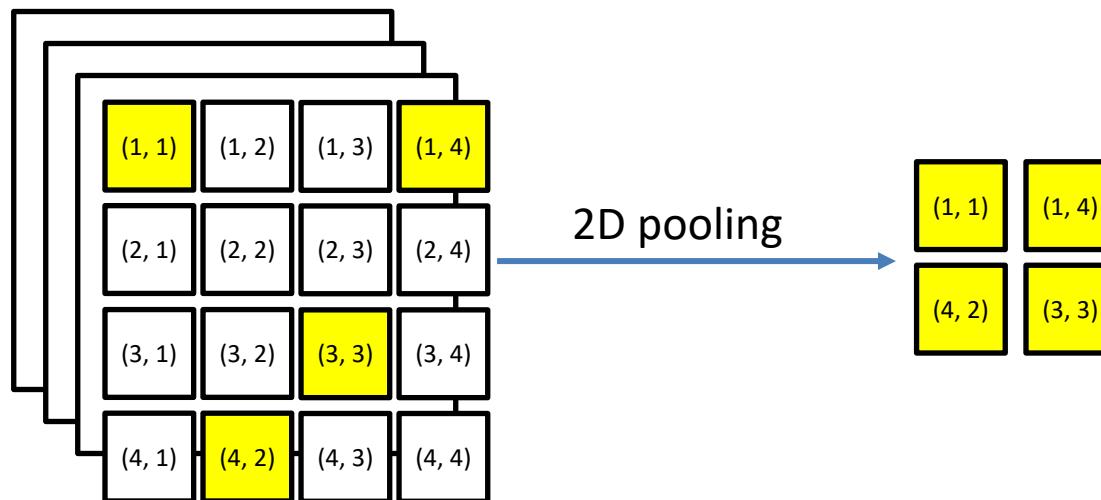
# ARC-II

- Let two sentences meet **before** their own high-level representations mature
- Basic matching signals: phrase sum interaction matrix
- Interaction: CNN to capture the local interaction structure
- Aggregation Function: MLP



# ARC-II (cont')

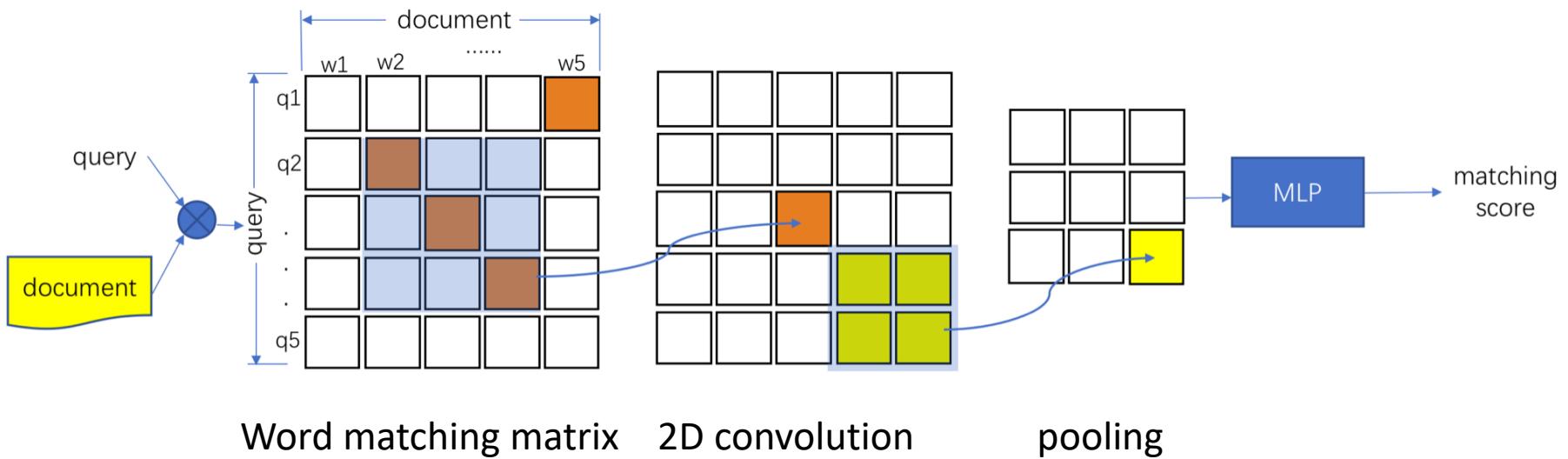
- Keeping word order information
  - Both the convolution and pooling are order preserving



- However, word level exact matching signals are lost
  - 2-D matching matrix is constructed based on the embedding of the words in two N-grams

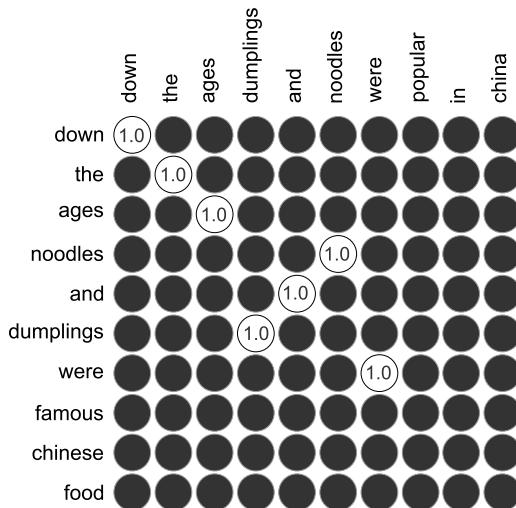
# MatchPyramid

- Inspired by image recognition
- Basic matching signals: word-level matching matrix
- Matching function: 2D convolution + MDP

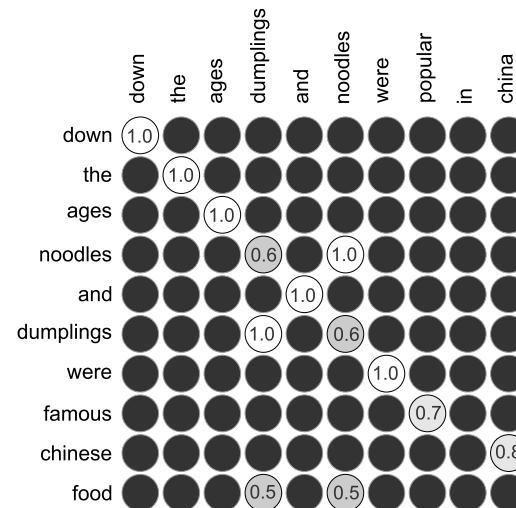


# Matching Matrix: Basic Matching Signals

- Each entry calculated based on
  - Word-level exact matching (0 or 1)
  - Semantic similarity based on embeddings of words



Exact match

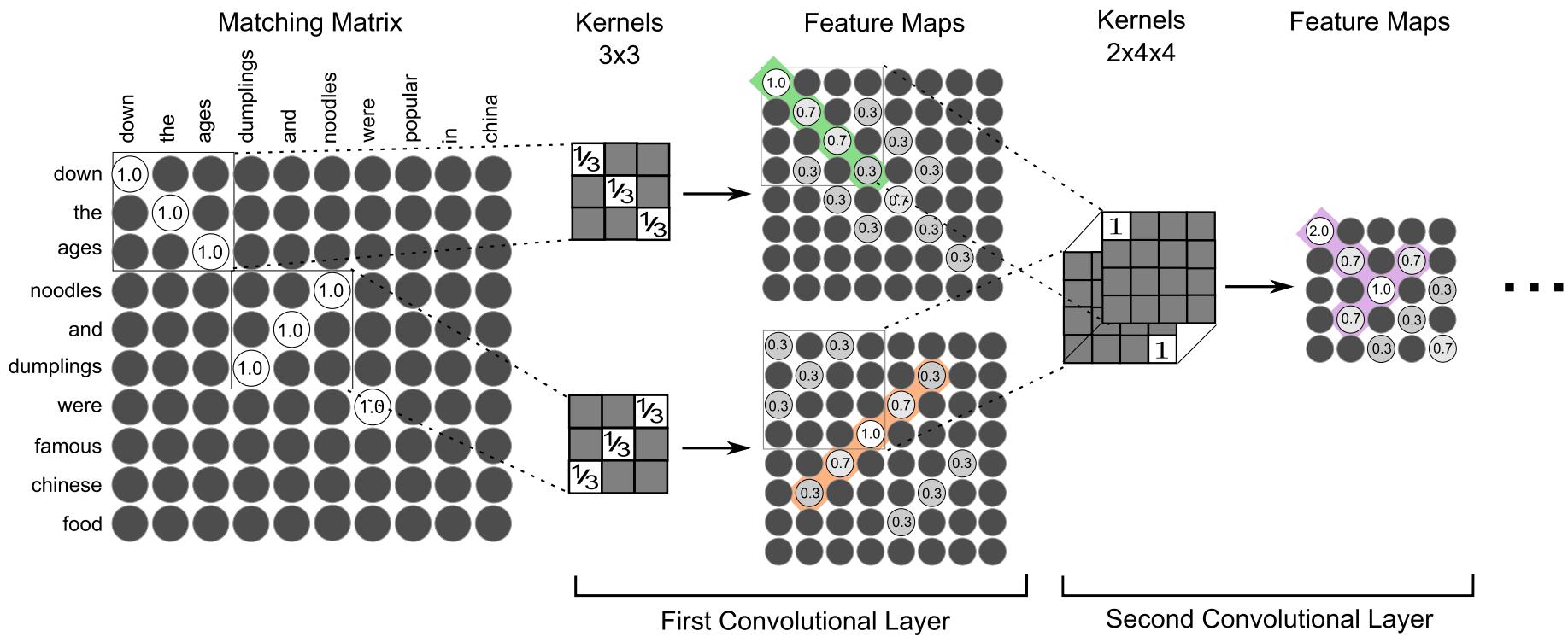


Cosine similarity

- Positions information of words is kept

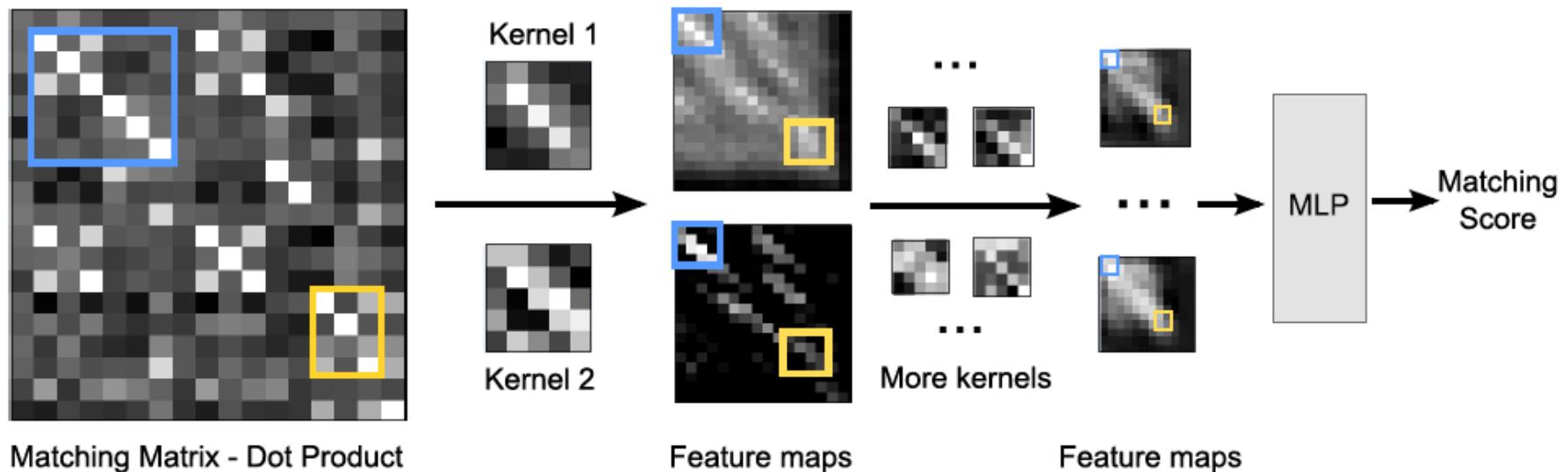
# Matching Function: 2D Convolution

- Discovering the matching patterns with CNN, stored in the kernels



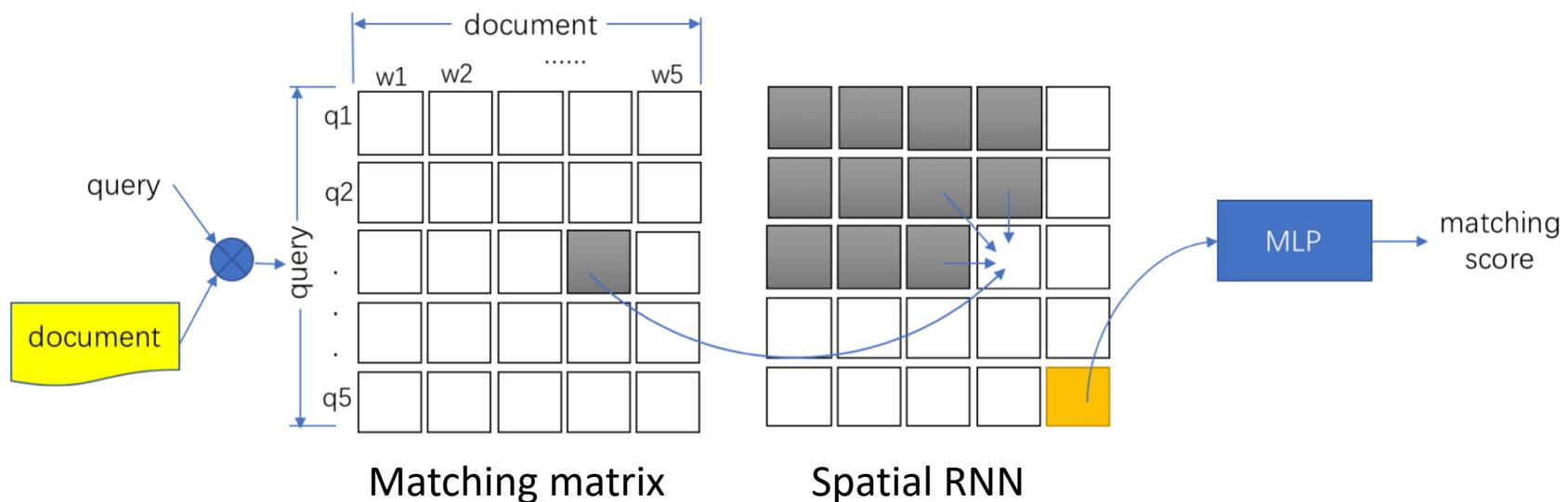
# Discovered Matching Patterns

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.

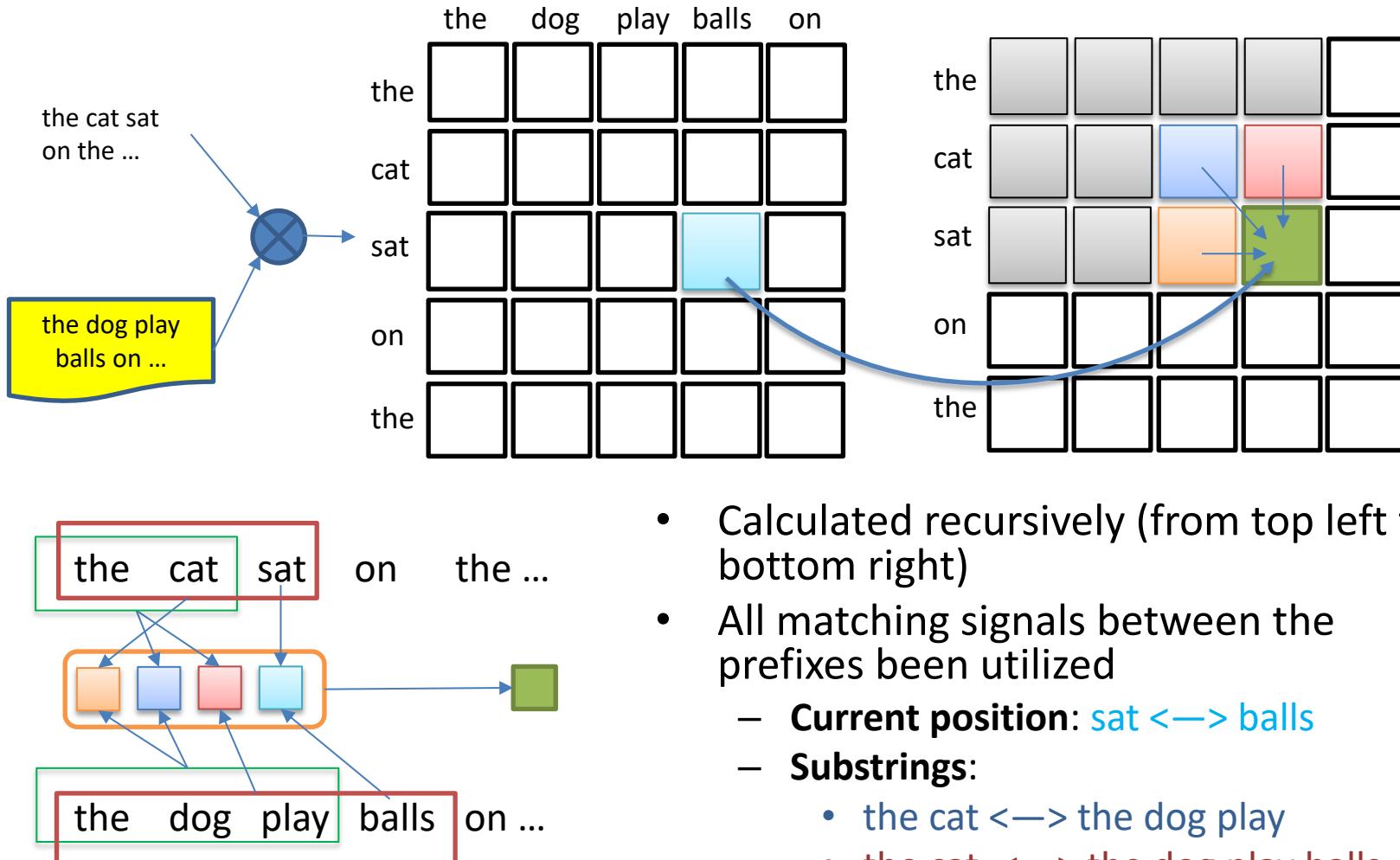


# Match-SRNN (Wan et al., 16)

- Based on spatial recurrent neural network (SRNN)
- Basic matching signals: word-level matching matrix
- Matching function: Spatial RNN + MLP

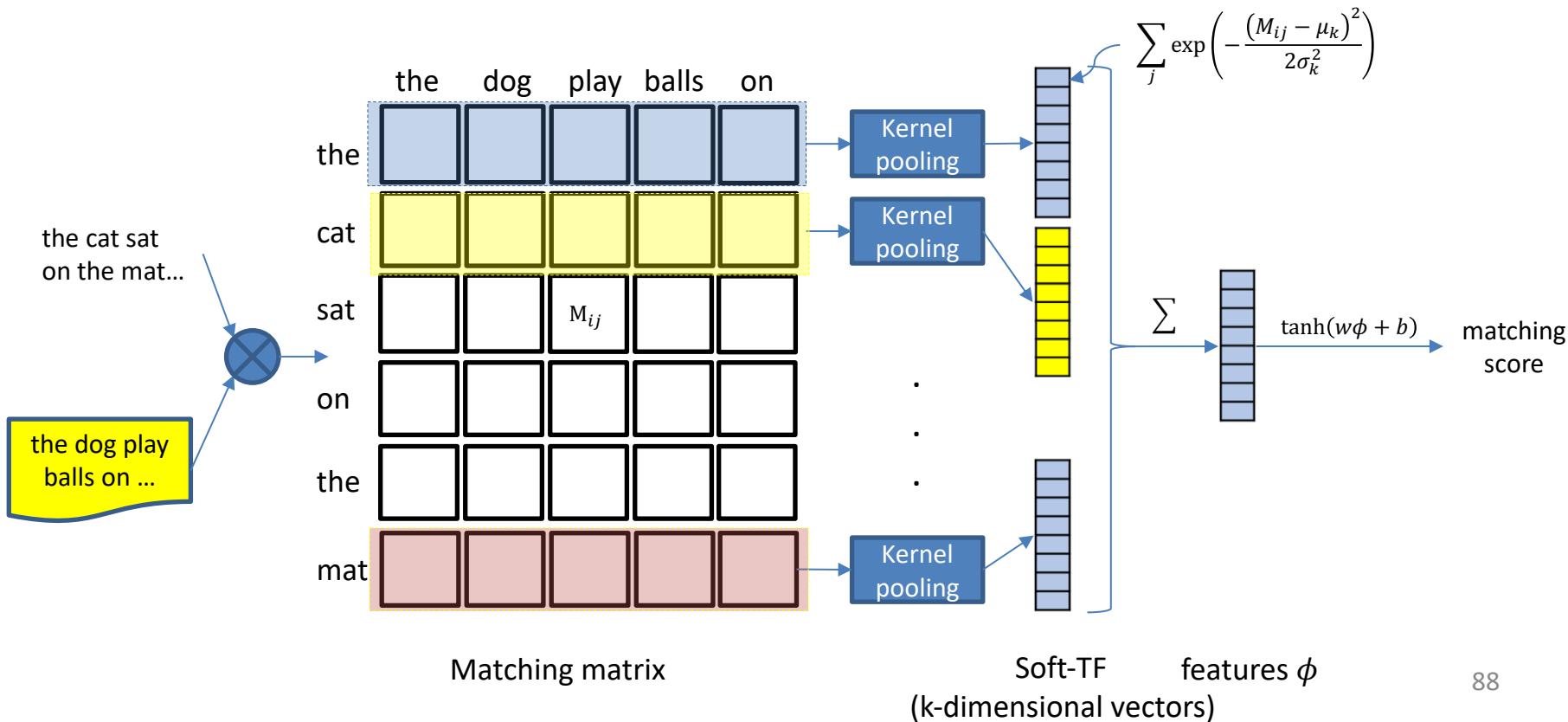


# Match-SRNN: Recursive Matching Structure



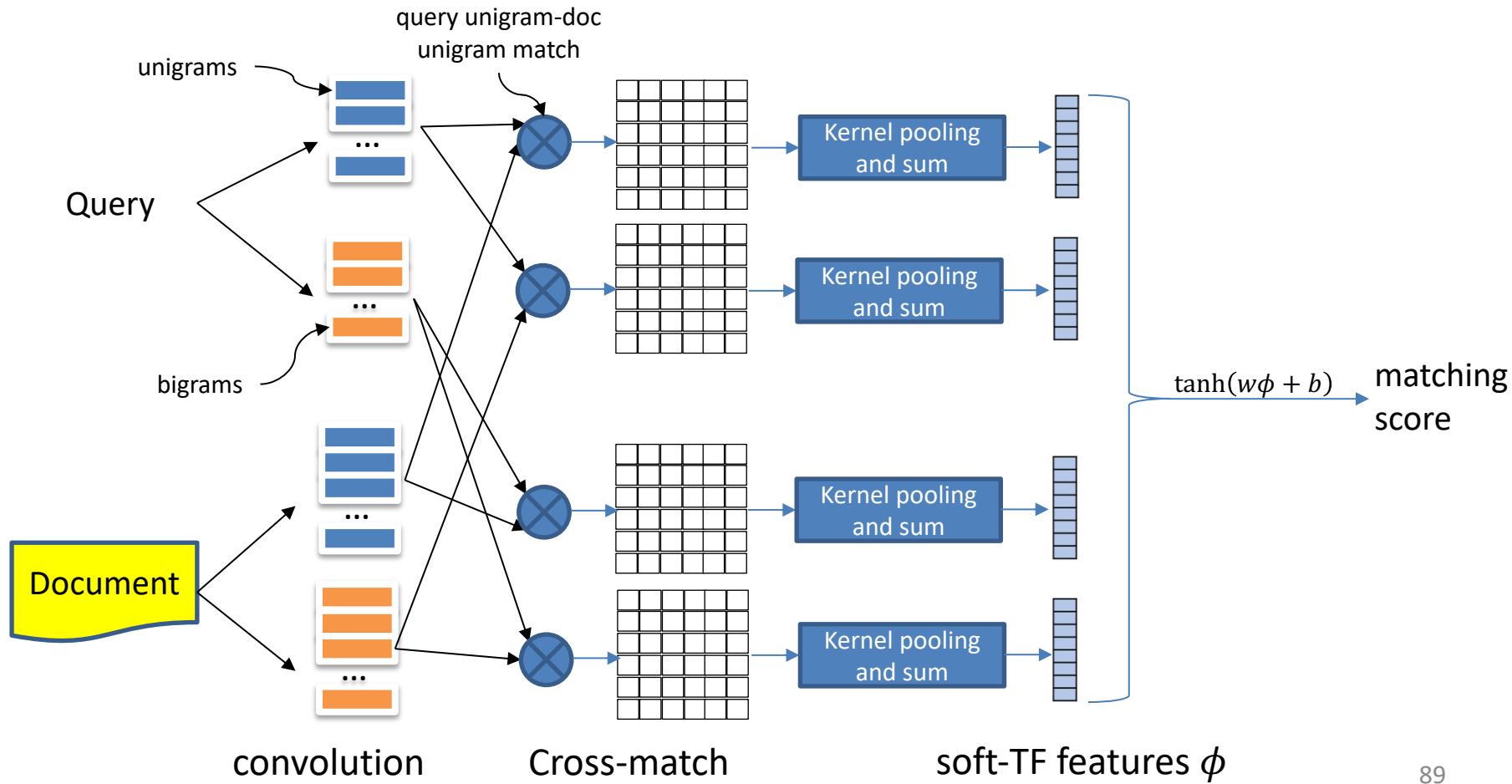
# K-NRM: Kernel Pooling as Matching Function (Xiong et al., SIGIR 2017)

- Basic matching signals: cosine similarity of word embeddings
  - Ranking function: kernel pooling + nonlinear feature combination
  - Semantic gap: embedding and soft-TF bridge the semantic gap
  - Proximity: kernel pooling and sum operations **lost word order** information



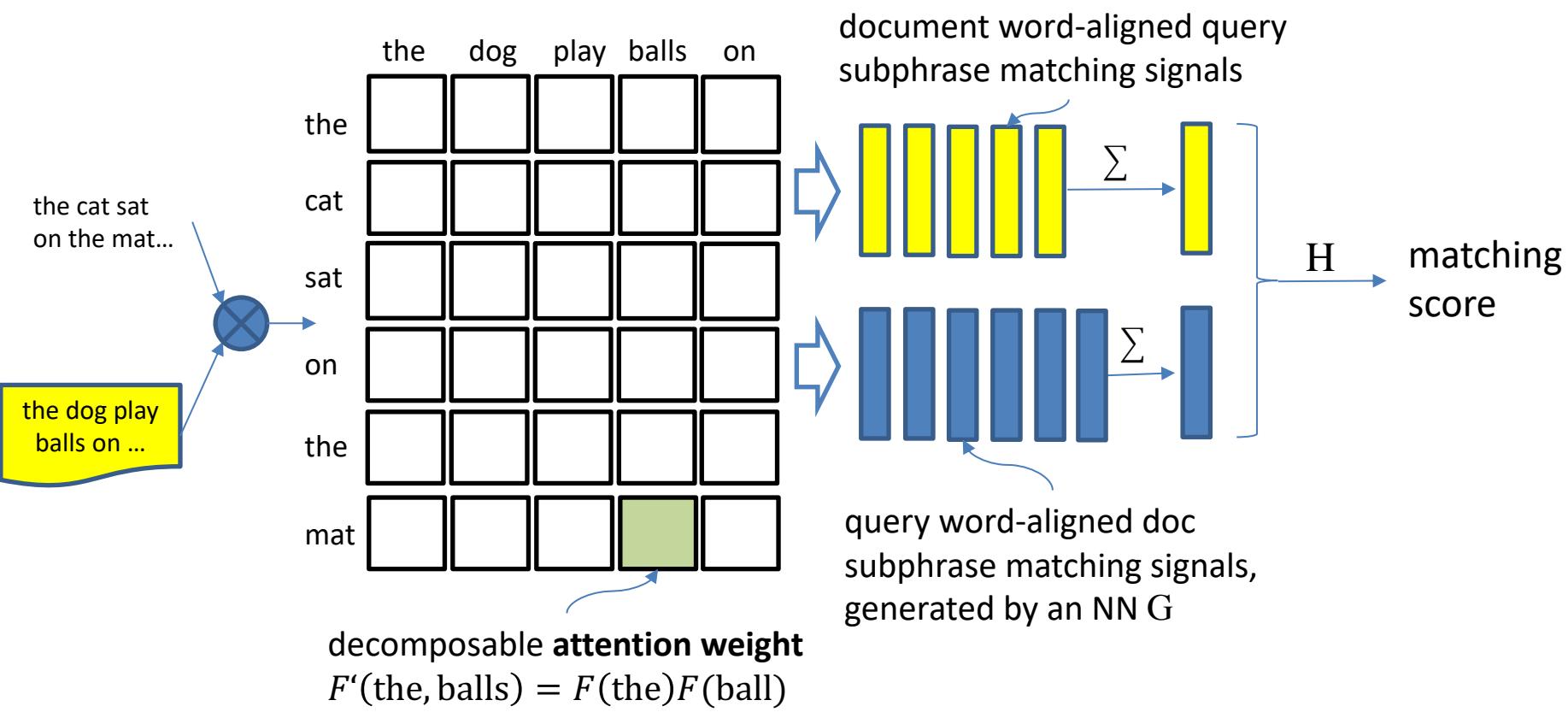
# Conv-KNRM (Dai et al., WSDM 2018)

- Based on KNRM
- N-gram cross-matching to capture the word order information



# Decomposable Attention Model for Matching (Parikh et al., EMNLP 2016)

- Based on decomposable attention model
  - Three steps: attend-compare-aggregate
    - **Attend**: soft-align words of query and document
    - **Compare**: separately compare word-aligned subphrase, get matching signals
    - **Aggregate**: aggregate the matching signals for produce final matching score



# Experimental Evaluation

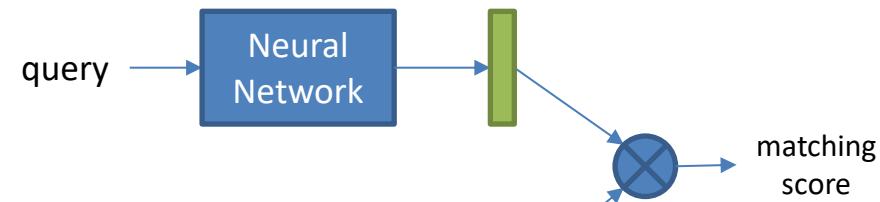
	Method	P@1	MRR
Traditional IR	BM25	0.579	0.457
Representation Learning methods	ARC-I	0.581	0.756
	CNTN	0.626	0.781
	LSTM-RNN	0.690	0.822
	uRAE	0.398	0.652
	MultiGranCNN	0.725	0.840
	MV-LSTM	0.766	0.869
Matching Function Learning	ARC-II	0.591	0.765
	MatchPyramid	0.764	0.867
	Match-SRNN	0.790	0.882

Based on Yahoo! Answers dataset (60,564 question-answer pairs)

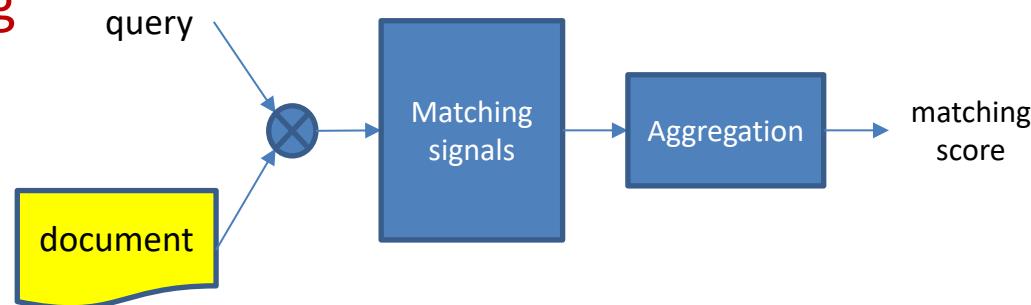
- Matching function learning based methods outperformed the representation learning ones

# Summary of Deep Matching Models in Search

- Representation learning:  
**representing queries and document in semantic space**



- Matching function learning:  
**discovering and aggregating the query-document matching patterns**



# References

- Clark J. Google turning its lucrative web search over to ai machines[J]. Bloomberg Technology. Publicado em, 2015, 26.
- Metz C. AI is transforming Google search[J]. The rest of the web is next. WIRED Magazine, 2016.
- Huang P S, He X, Gao J, et al. Learning deep structured semantic models for web search using clickthrough data[C]//Proceedings of the 22nd ACM international conference on Conference on information & knowledge management. ACM, 2013: 2333-2338.
- Hu B, LuShen Y, He X, Gao J, et al. A latent semantic model with convolutional-pooling structure for information retrieval[C]//Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management. ACM, 2014: 101-110.
- Z, Li H, et al. Convolutional neural network architectures for matching natural language sentences[C]//Advances in neural information processing systems. 2014: 2042-2050.
- Qiu X, Huang X. Convolutional Neural Tensor Network Architecture for Community-Based Question Answering[C]//IJCAI. 2015: 1305-1311.
- Palangi H, Deng L, Shen Y, et al. Deep sentence embedding using long short-term memory networks: Analysis and application to information retrieval[J]. IEEE/ACM Transactions on Audio, Speech and Language Processing (TASLP), 2016, 24(4): 694-707.
- Yin W, Schütze H. Multigrancnn: An architecture for general matching of text chunks on multiple levels of granularity[C]//Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers). 2015, 1: 63-73.

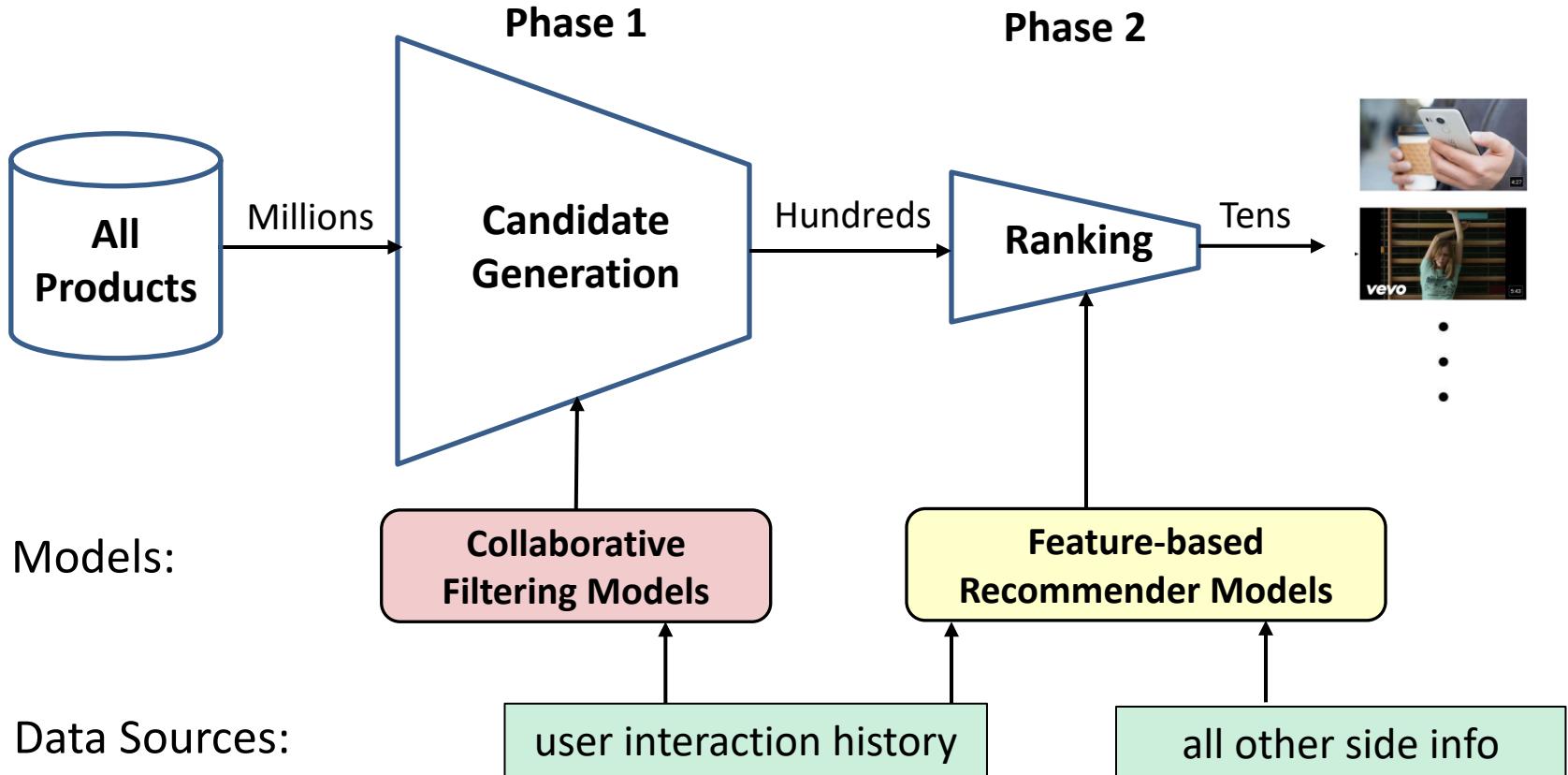
# References

- Socher R, Huang E H, Pennin J, et al. Dynamic pooling and unfolding recursive autoencoders for paraphrase detection[C]//Advances in neural information processing systems. 2011: 801-809.
- Wan S, Lan Y, Guo J, et al. A Deep Architecture for Semantic Matching with Multiple Positional Sentence Representations[C]//AAAI. 2016, 16: 2835-2841.
- Pang L, Lan Y, Guo J, et al. Text Matching as Image Recognition[C]//AAAI. 2016: 2793-2799.
- Shengxian Wan, Yanyan Lan, Jun Xu, Jiafeng Guo, Liang Pang, and Xueqi Cheng. 2016. Match-SRNN: modeling the recursive matching structure with spatial RNN. In Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence (IJCAI'16), 2922-2928.
- Ankur P. Parikh, Oscar Tackstrom, Dipanjan Das, and Jakob Uszkoreit. A Decomposable Attention Model for Natural Language Inference. In Proceedings of EMNLP, 2016.
- Zhuyun Dai, Chenyan Xiong, Jamie Callan, and Zhiyuan Liu. Convolutional Neural Networks for Soft-Matching N-Grams in Ad-hoc Search. In Proceedings of WSDM 2018.
- Chenyan Xiong, Zhuyun Dai, Jamie Callan, Zhiyuan Liu, Russell Power. End-to-End Neural Ad-hoc Ranking with Kernel Pooling. In Proceedings of SIGIR 2017.

# Outline of Tutorial

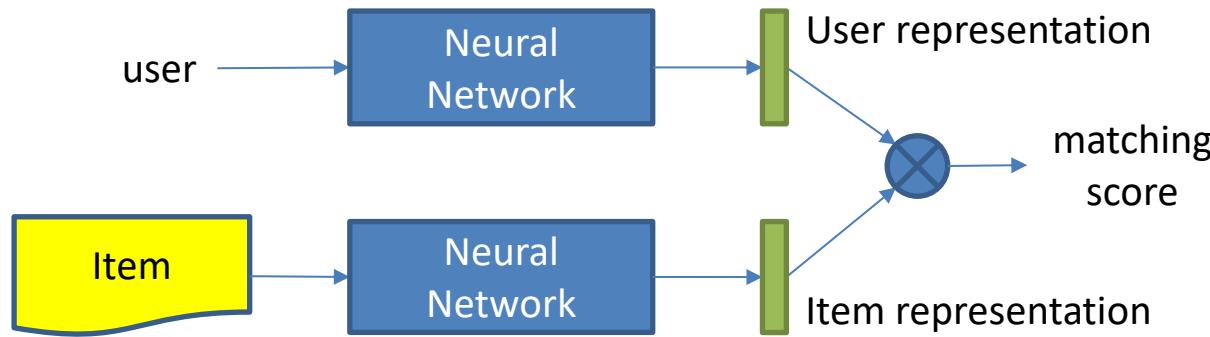
- Unified view of matching in search and recommendation
- Part 1: Traditional Approaches to Matching
- Part 2: Deep Learning Approaches to Matching
  - Deep matching models for search
  - Deep matching models for recommendation
- Summary

# Modern RecSys Architecture (Covington et al, Recsys'16)

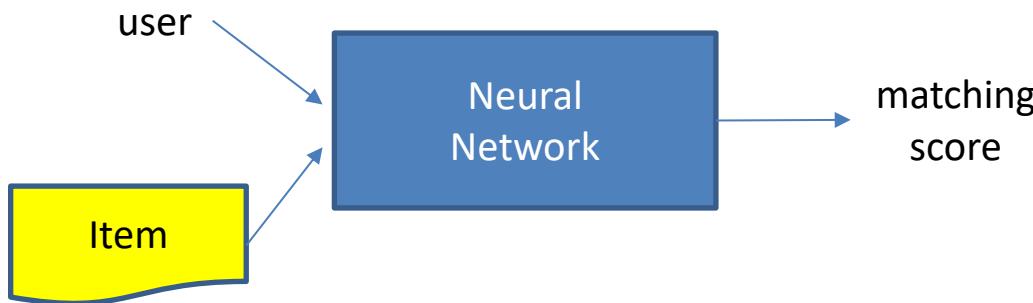


# Deep Matching Models for Recommendation

- Methods of representation learning



- Methods of matching function learning



# Methods of Representation Learning

## 1. Pure CF models:

Only ID or interaction history is used as input.

- **DeepMF**: Deep Matrix Factorization (Xue et al, IJCAI'17)
- **AutoRec**: Autoencoders Meeting CF (Sedhain et al, WWW'15)
- **CDAE**: Collaborative Denoising Autoencoder (Wu et al, WSDM'16)

## 2. CF with side information:

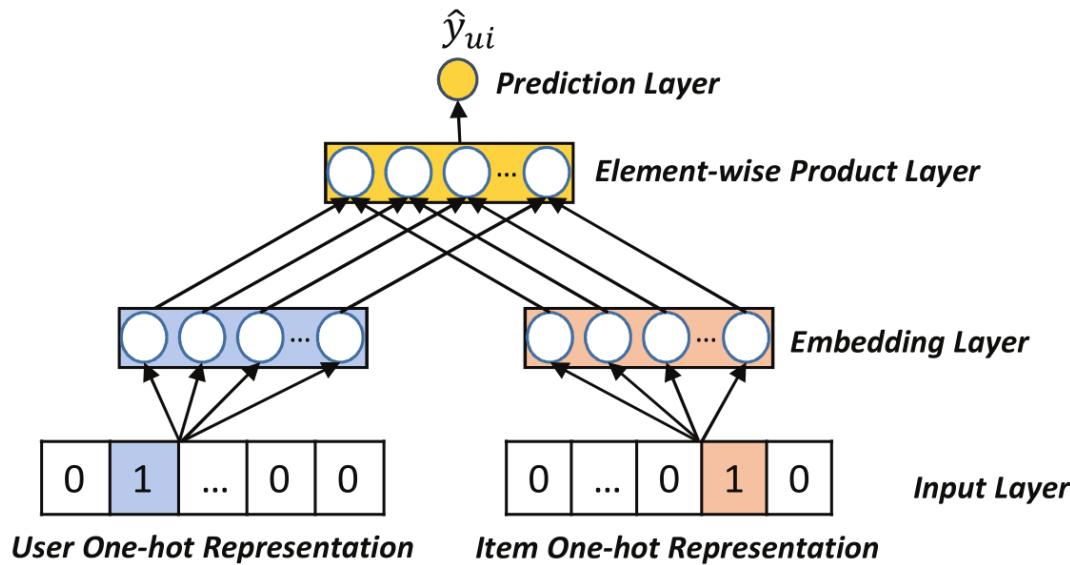
Any available data can be used as input.

- **DCF**: Deep Collaborative Filtering via Marginalized DAE (Li et al, CIKM'15)
- **DUIF**: Deep User-Image Feature (Geng et al, ICCV'15)
- **ACF**: Attentive Collaborative Filtering (Chen et al, SIGIR'17)
- **CKB**: Collaborative Knowledge Base Embeddings (Zhang et al, KDD'16)

# Matrix Factorization as a Neural Network (Wang et al, SIGIR'17)

- **Input:** user  $\rightarrow$  ID (one-hot), item  $\rightarrow$  ID (one-hot).
- **Representation Function:** linear embedding layer.
- **Matching Function:** inner product.

$$f_{MF}(u, i | \mathbf{p}_u, \mathbf{q}_i) = \mathbf{p}_u^\top \mathbf{q}_i = \sum_{k=1}^K p_{uk} q_{ik},$$



# Deep Matrix Factorization (Xue et al, IJCAI'17)

- **Input:**

user -> historically rated items (multi-hot), i.e., row vector of  $\mathbf{Y}$

indicates the user's global preference

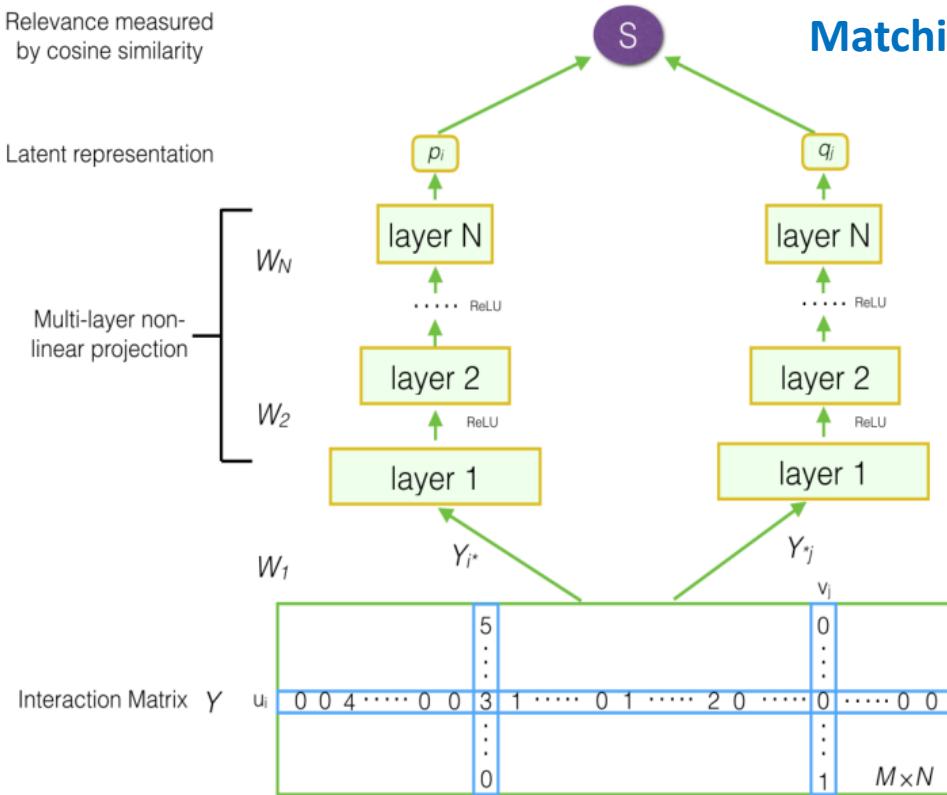
item -> users who have rated it (multi-hot), i.e., column vector of  $\mathbf{Y}$

indicates the item's rating profile.

Interaction Matrix $\mathbf{Y}$		$u_i$	$0 \ 0 \ 4 \ \dots \ 0 \ 0$	$3 \ 1 \ \dots \ 0 \ 1 \ \dots \ 2 \ 0 \ \dots \ 0 \ \dots \ 0 \ 0$	$M \times N$
5	$\vdots$	0	0	0	1

# Deep Matrix Factorization (Xue et al, IJCAI'17)

- **Representation Function:**
  - Multi-Layer Perceptron (same as DSSM).



**Matching Function: cosine similarity**

$$\text{cosine}(p_i, q_j) = \frac{p_i^T q_j}{\|p_i\| \|q_j\|}$$

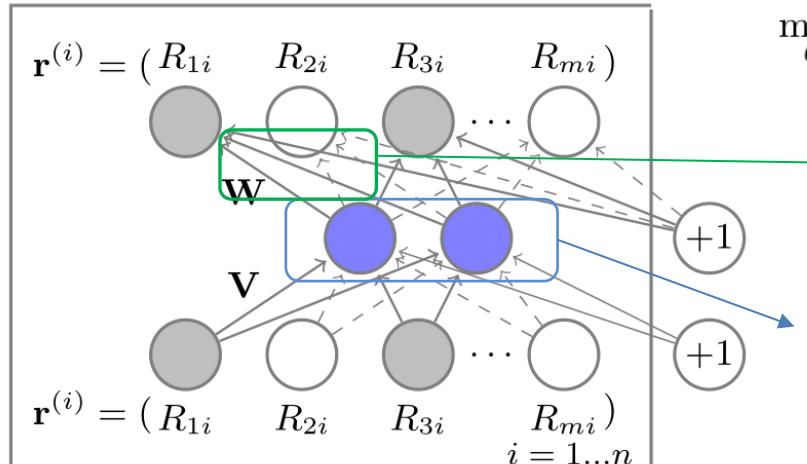
$$\begin{aligned} l_1 &= W_1 x \\ l_i &= f(W_{i-1} l_{i-1} + b_i), i = 2, \dots, N-1 \\ h &= f(W_N l_{N-1} + b_N) \end{aligned}$$

$$p_i = f_{\theta_N^U}(\dots f_{\theta_3^U}(W_{U2} f_{\theta_2^U}(Y_{i*} W_{U1}))\dots)$$

$$q_j = f_{\theta_N^I}(\dots f_{\theta_3^I}(W_{V2} f_{\theta_2^I}(Y_{*j}^T W_{V1}))\dots)$$

# AutoRec (Sedhain et al, WWW'15)

- **Input: (similar to DeepMF)**
  - user -> historically rated items -> **user-based autoencoder**.
  - item -> users who have rated it -> **item-based autoencoder**.
- **Representation Function:** Multi-Layer Perceptron
- **Matching Function:** inner product



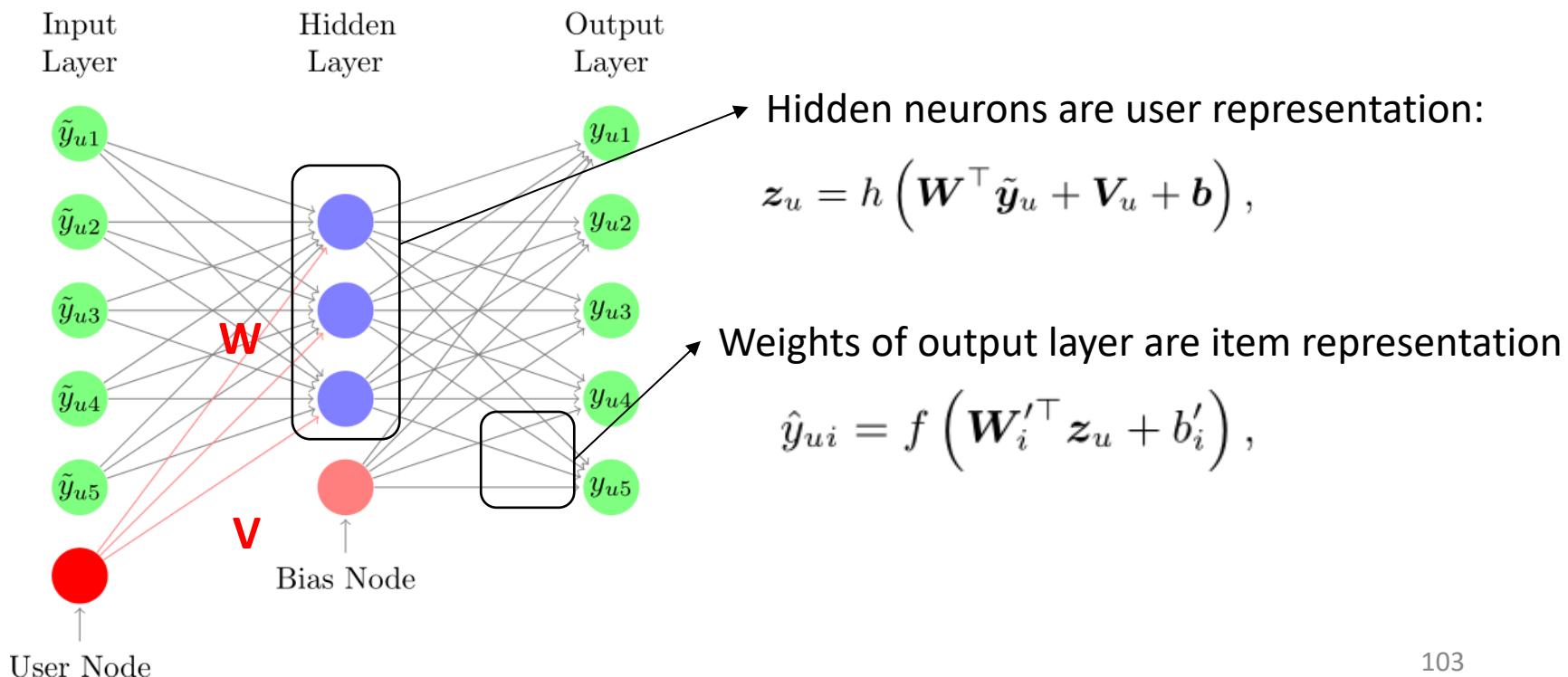
$$\text{Input reconstruction: } h(\mathbf{r}; \theta) = f(\mathbf{W} \cdot g(\mathbf{V}\mathbf{r} + \boldsymbol{\mu}) + \mathbf{b})$$
$$\min_{\theta} \sum_{i=1}^n \|\mathbf{r}^{(i)} - h(\mathbf{r}^{(i)}; \theta)\|_{\mathcal{O}}^2 + \frac{\lambda}{2} \cdot (\|\mathbf{W}\|_F^2 + \|\mathbf{V}\|_F^2),$$

→ Output weights denote item representation

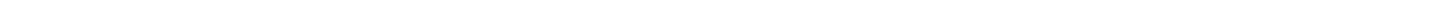
Hidden neurons denote user representation

# Collaborative Denoising Auto-Encoder (Wu et al, WSDM'16)

- **Input:**  
user  $\rightarrow$  ID & historically rated items (similar to SVD++)  
item  $\rightarrow$  ID
- **Representation:** Multi-Layer Perceptron



# Short Summary

- Either ID or history is used as the profile of user/item
  - Models with history as input are more expressive, but are also more expensive to train.
  - The Auto-Encoder architecture is essentially identical to  
  
MLP (representation learning) + MF (matching function).  

- Nonlinear                                    Linear

# Methods of Representation Learning

## 1. Pure CF models:

Only ID or interaction history is used as input.

- **DeepMF**: Deep Matrix Factorization (Xue et al, IJCAI'17)
- **AutoRec**: Autoencoders Meeting CF (Sedhain et al, WWW'15)
- **CDAE**: Collaborative Denoising Autoencoder (Wu et al, WSDM'16)

## 2. CF with side information:

Any available data can be used as input.

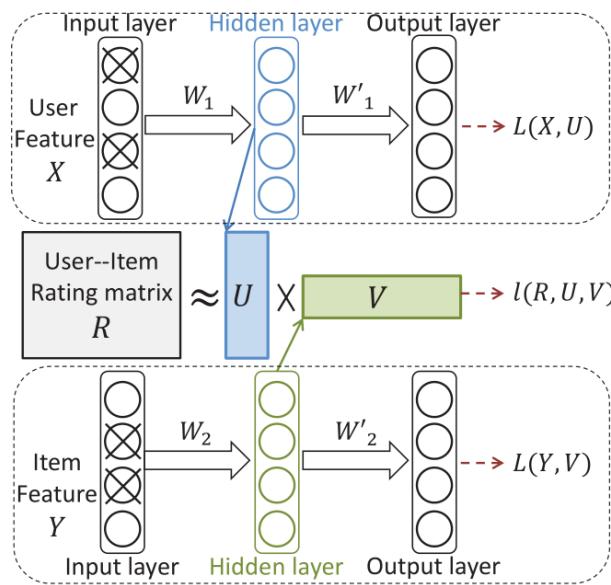
- **DCF**: Deep Collaborative Filtering via Marginalized DAE (Li et al, CIKM'15)
- **DUIF**: Deep User-Image Feature (Geng et al, ICCV'15)
- **ACF**: Attentive Collaborative Filtering (Chen et al, SIGIR'17)
- **CKB**: Collaborative Knowledge Base Embeddings (Zhang et al, KDD'16)

# Deep Collaborative Filtering via Marginalized DAE (Li et al, CIKM'15)

- Denoising Auto-Encoder is used to learn features (hidden layers) of user and item from side information.
- The predictive model is MF.

age, gender,  
city, occupation,  
locations ...

genres,  
title, texts



$$\arg \min_{U, V, W_1, W_2, P_1, P_2}$$

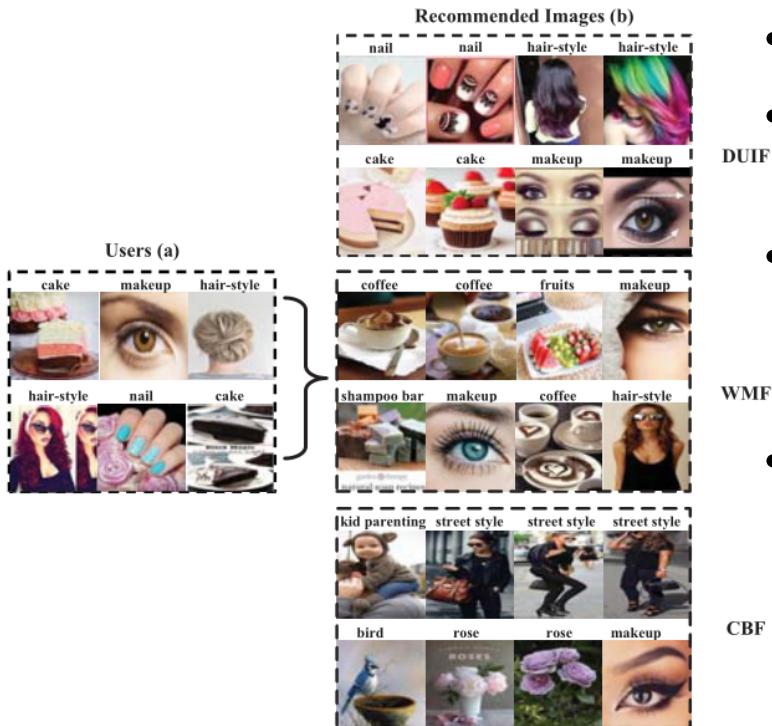
User features reconstruction      Item features reconstruction

$$\mathcal{L}_U(W_1, P_1, U) + \mathcal{L}_V(W_2, P_2, V) +$$

$$\alpha \|A \odot (R - UV^\top)\|_F^2 + \beta(\|U\|_F^2 + \|V\|_F^2)$$

**Matrix Factorization Kernel**

# DUIF: Deep User and Image Feature Learning (Geng et al, ICCV'15)



- Task: collaborative image recommendation
  - Deep CNN (AlexNet) is used to extract features for images
  - The deep image features (dim=4096) are projected to user latent space (dim=300) by using linear projection.
  - The predictive model is MF:
- $$\hat{y}_{ui} = \langle \mathbf{p}_u, \mathbf{W}^T \text{CNN}(\mathbf{f}_i) \rangle,$$
- Linear Projection      Image raw feature
- The overall model (MF+W+CNN) is trained end-to-end.

# ACF: Attentive Collaborative Filtering (Chen et al, SIGIR'17)

- **Input:**

user -> ID & historical interacted items.

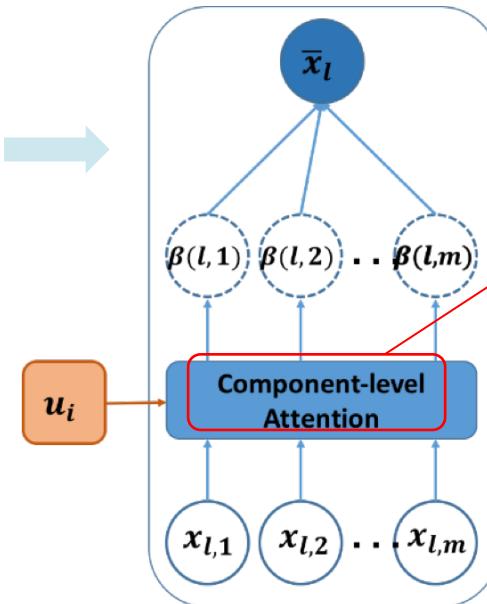
Item -> ID & visual features.

- **Item Representation:**

**Component-level attention** -> not all components contribute equally to an item's representation



A user's preference on different **components** of the item is unknown & **not equal!**



$$\bar{x}_l = \sum_{m=1}^{|\{x_{l*}\}|} \beta(i, l, m) \cdot x_{lm}$$

$$\beta(i, l, m) = \frac{\exp(b(i, l, m))}{\sum_{n=1}^{|\{x_{l*}\}|} \exp(b(i, l, n))}$$

$$b(i, l, m) = \mathbf{w}_2^T \phi(\mathbf{W}_{2u} \mathbf{u}_i + \mathbf{W}_{2x} \mathbf{x}_{lm} + \mathbf{b}_2) + c_2$$

# ACF: Attentive Collaborative Filtering (Chen et al, SIGIR'17)

- **Input:**
  - user -> ID & historical interacted items.
  - item -> ID & visual features.
- **User Presentation:**
  - **Item-level attention** -> not all historical items contribute equally to a user's representation



$$\hat{R}_{ij} = \left( \mathbf{u}_i + \sum_{l \in \mathcal{R}(i)} \alpha(i, l) \mathbf{p}_l \right)^T \mathbf{v}_j$$

Attention weights learned by a neural net  
↔ Attentive SVD++ model.

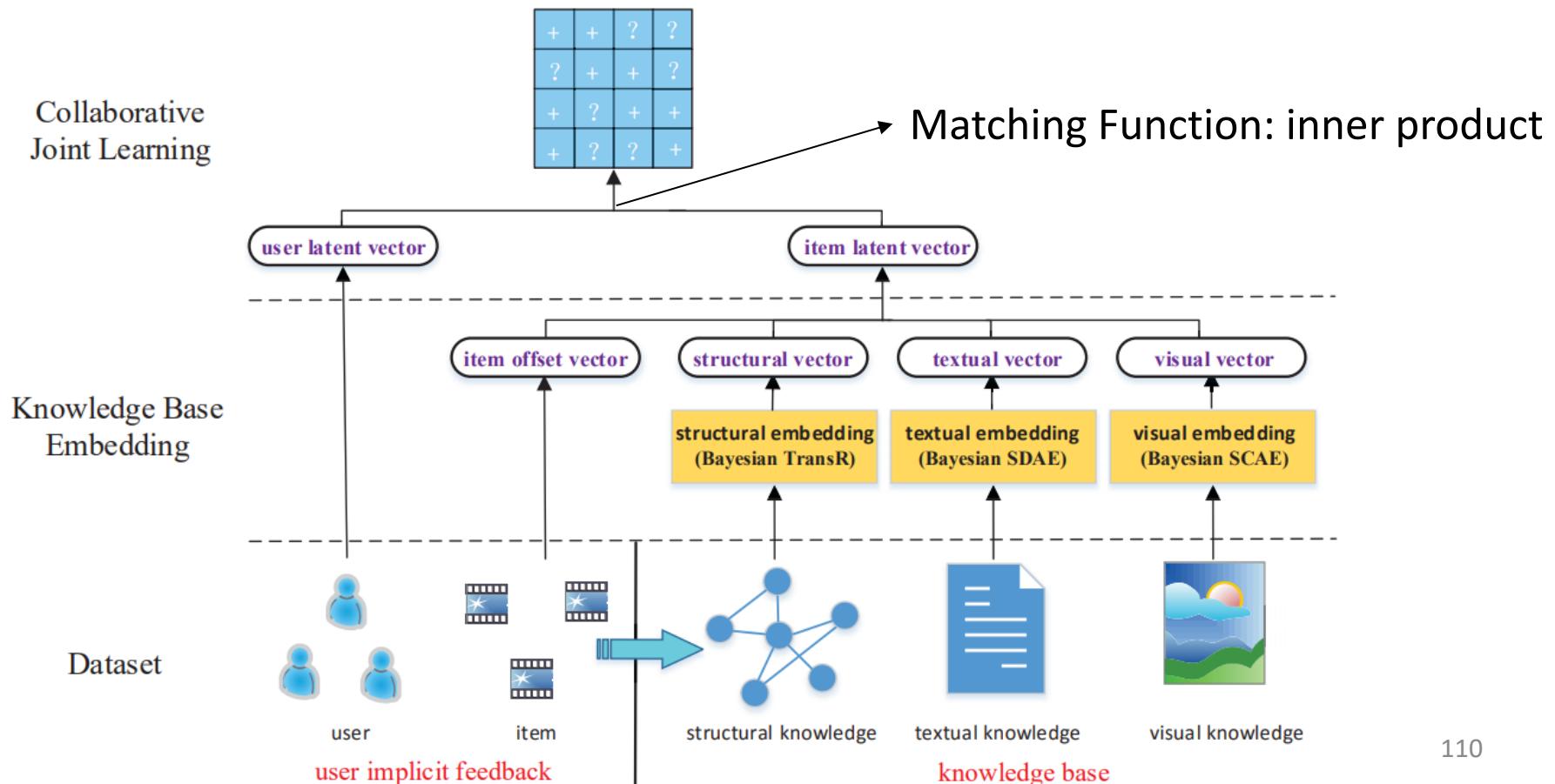
A user's preference on  
**historical items** is unknown &  
not equal!

# CKE: Collaborative Knowledge Base Embedding (Zhang et al, KDD'16)

- **Input:**

user  $\rightarrow$  ID

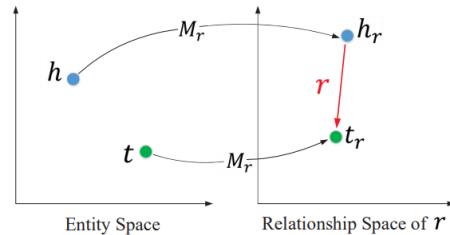
item  $\rightarrow$  ID + Information in KB (structural, textual, visual)



# CKE: Collaborative Knowledge Base Embedding (Zhang et al, KDD'16)

- **Representation:**

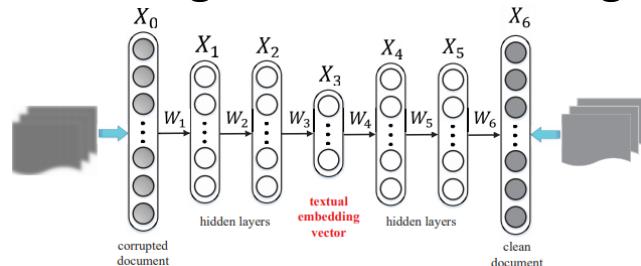
- Structural embedding: TransR, TransE, ...



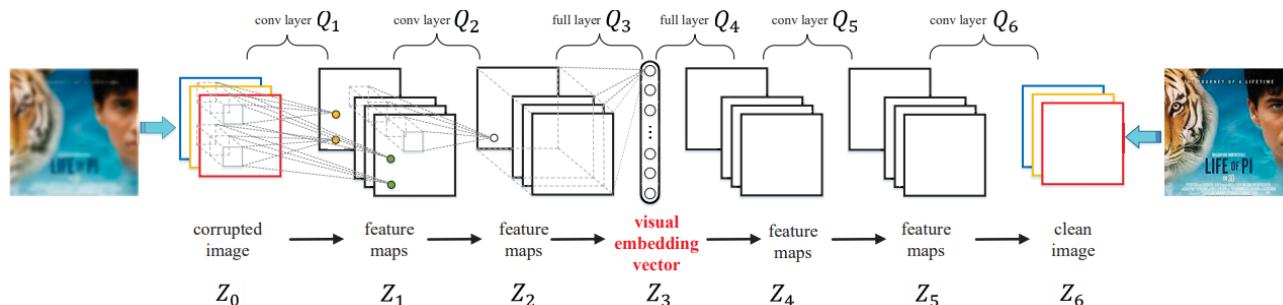
$$\mathbf{v}_h^r = \mathbf{v}_h \mathbf{M}_r, \quad \mathbf{v}_t^r = \mathbf{v}_t \mathbf{M}_r.$$

$$f_r(v_h, v_t) = \|\mathbf{v}_h^r + \mathbf{r} - \mathbf{v}_t^r\|_2^2.$$

- Textual embedding: stacked denoising auto-encoders (S-DAE)

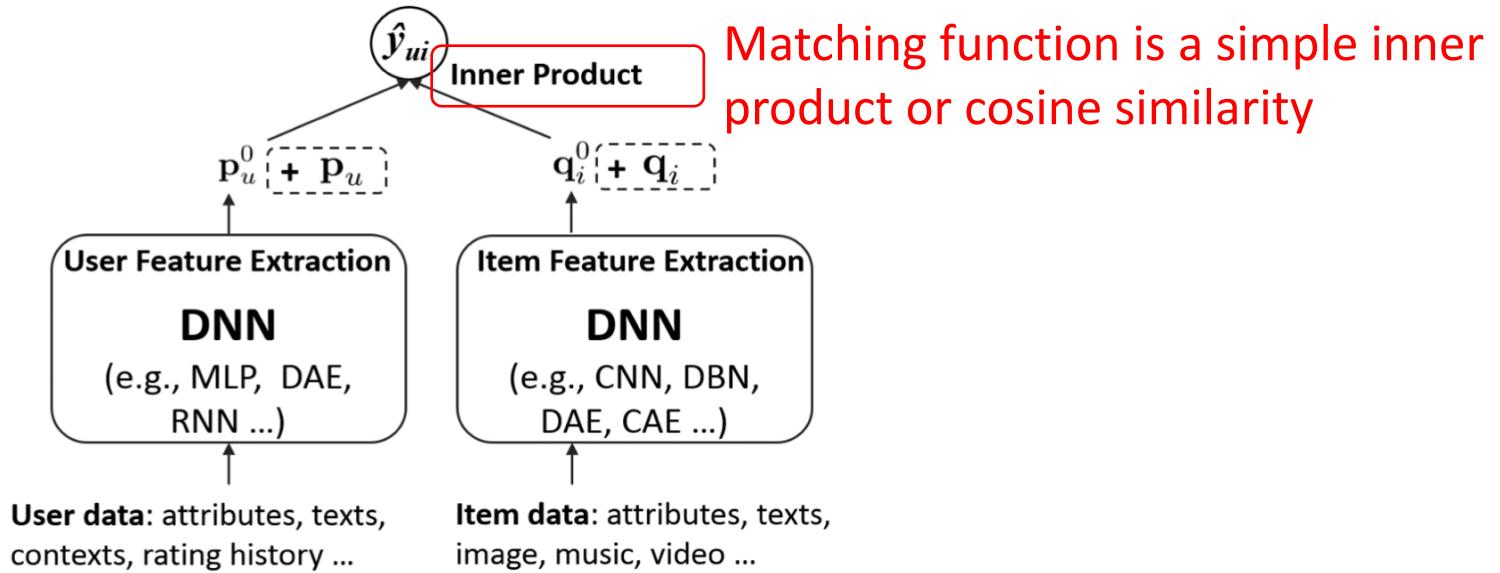


- Visual embedding: stacked convolutional auto-encoders (SCAE)



# Short Summary

- A General framework to summarize the above works:



- Depending on the available data to describe a user/item, we can choose appropriate DNN to learn representation.  
E.g., Textual Attributes -> AutoRec, Image -> CNN, Video -> RNN etc.

# Next: Methods of Matching Function Learning

## 1. Pure CF models:

- Based on Neural Collaborative Filtering (NCF) framework:  
**NeuMF**: Neural Matrix Factorization (He et al, WWW'17)  
**NNCF**: Neighbor-based NCF (Bai et al, CIKM'17)  
**ConvNCF**: Outer Product-based NCF (He et al, IJCAI'18)
- Based on Translation framework:  
**TransRec**: Translation-based Recommendation (He et al, Recsys'17)  
**LRML**: Latent Relational Metric Learning (Tay et al, WWW'18)

## 2. Feature-based models:

- Based on Multi-Layer Perceptron:  
**Wide&Deep** (Cheng et al, DLRS'16), **Deep Crossing** (Shan et al, KDD'16)
- Based on Factorization Machines (FM):  
**Neural FM** (He and Chua, SIGIR'17), **Attentional FM** (Xiao et al, IJCAI'17),
- Based on Trees:  
**GB-CENT**: Categorical Embedding and Numerical Trees (Zhao et al, WWW'18)  
**DEF**: Deep Embedding Forest (Zhu et al, KDD'17)  
**TEM**: Tree-enhanced Embedding Model (Wang et al, WWW'18)

# Why Using Neural Networks to Learn the Matching Function?

- The simple choice of **inner product** can limit the **expressiveness** of an embedding-based matching model.

$$\hat{y}_{ui} = \mathbf{U}_i^T \mathbf{V}_j \simeq \cos(\mathbf{U}_i, \mathbf{V}_j) \quad (\text{E.g., assuming a unit length})$$

- Example:

	$i_1$	$i_2$	$i_3$	$i_4$	$i_5$
$u_1$	1	1	1	0	1
$u_2$	0	1	1	0	0
$u_3$	0	1	1	1	0
$u_4$	1	0	1	1	1

$$\text{sim}(u_1, u_2) = 0.5$$

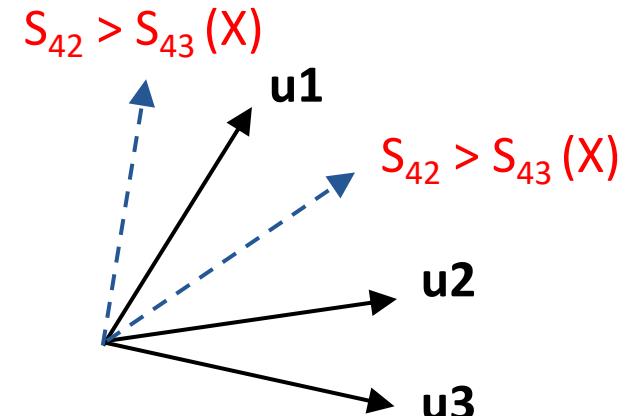
$$\text{sim}(u_3, u_1) = 0.4$$

$$\text{sim}(u_3, u_2) = 0.66$$

$$\text{sim}(u_4, u_1) = 0.6 \quad *****$$

$$\text{sim}(u_4, u_2) = 0.2 \quad *$$

$$\text{sim}(u_4, u_3) = 0.4 \quad ***$$



# Neural Collaborative Filtering Framework (He et al, WWW'17)

- NCF is a general framework that replaces the inner product with a neural network to learn the matching function.  $\hat{y}_{ui} = f(\mathbf{p}_u, \mathbf{q}_i)$

Matching function based on NN

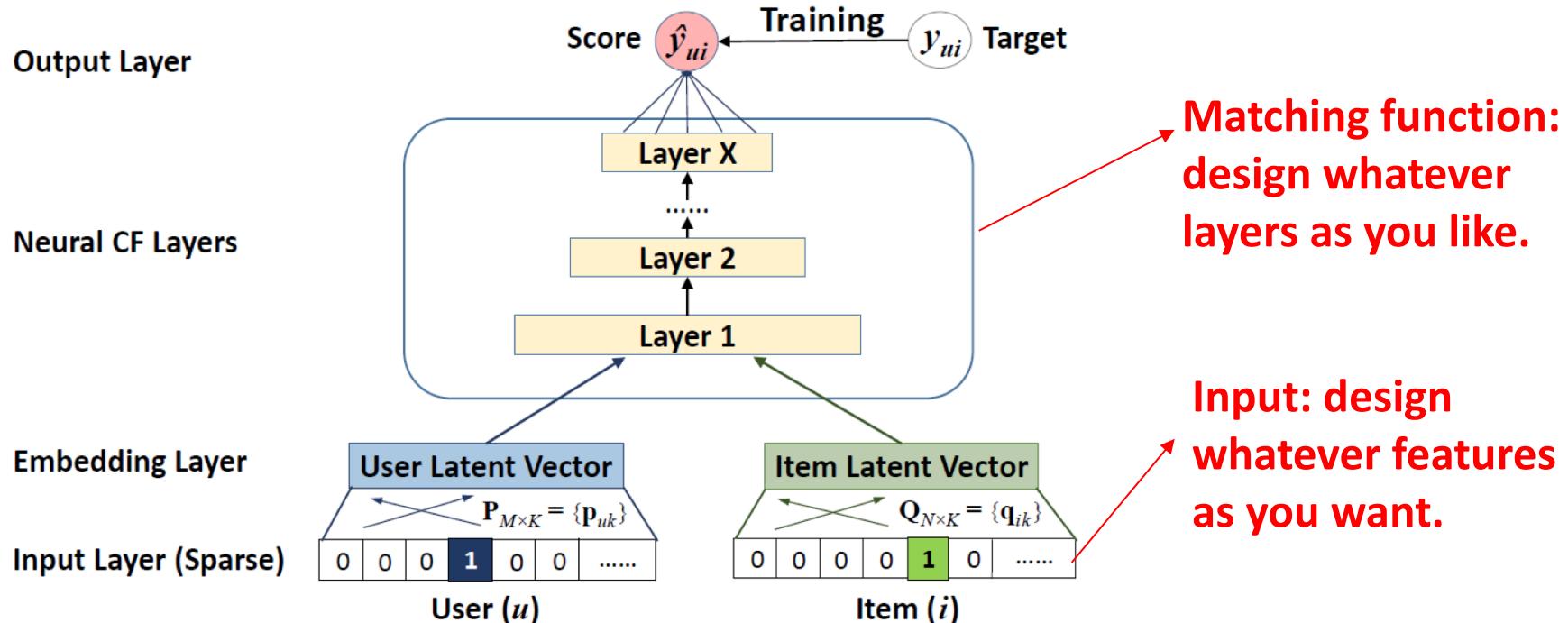
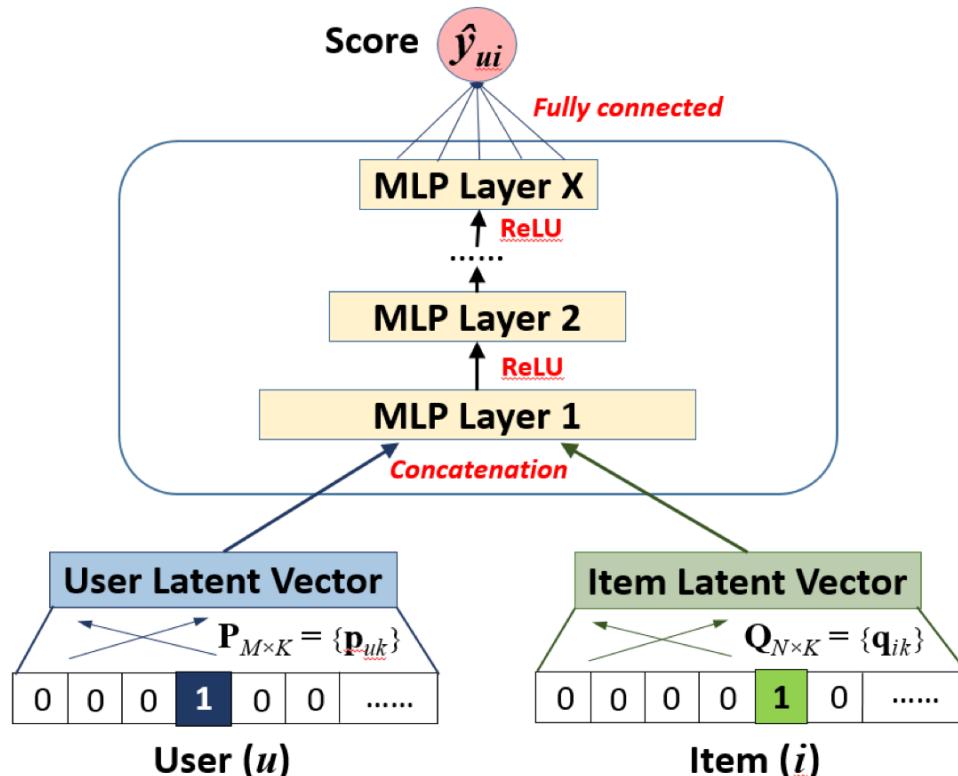


Figure 2: Neural collaborative filtering framework

# Multi-Layer Perceptron for CF

- The most intuitive idea is to use a Multi-Layer Perceptron as the matching function.



Unfortunately, MLP doesn't perform well and underperforms MF.

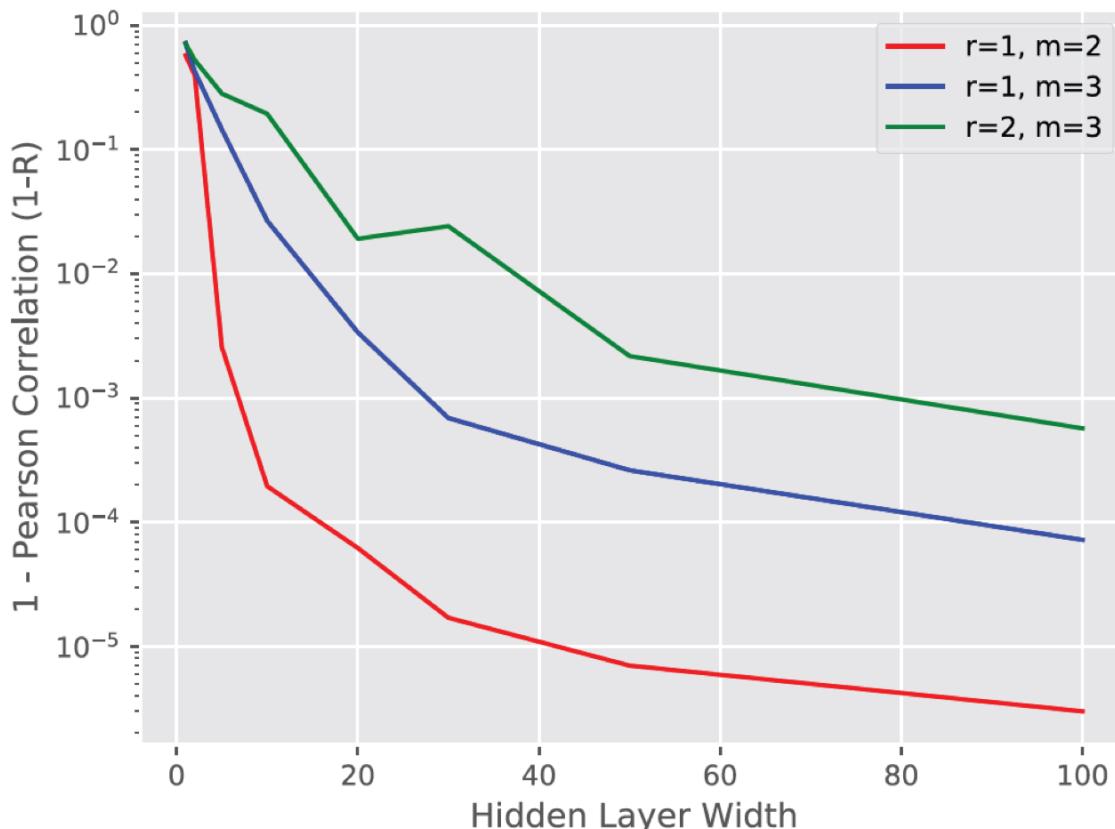
Why?

# MLP is Weak in Capturing Low-Rank Relation (Beutel et al, WSDM'18)

Setting: Generating low-rank data, and using one-layer MLP to fit it.

$r$ : rank size;

$m$ : data dimension (2 -> matrix; 3 -> 3D tensor).

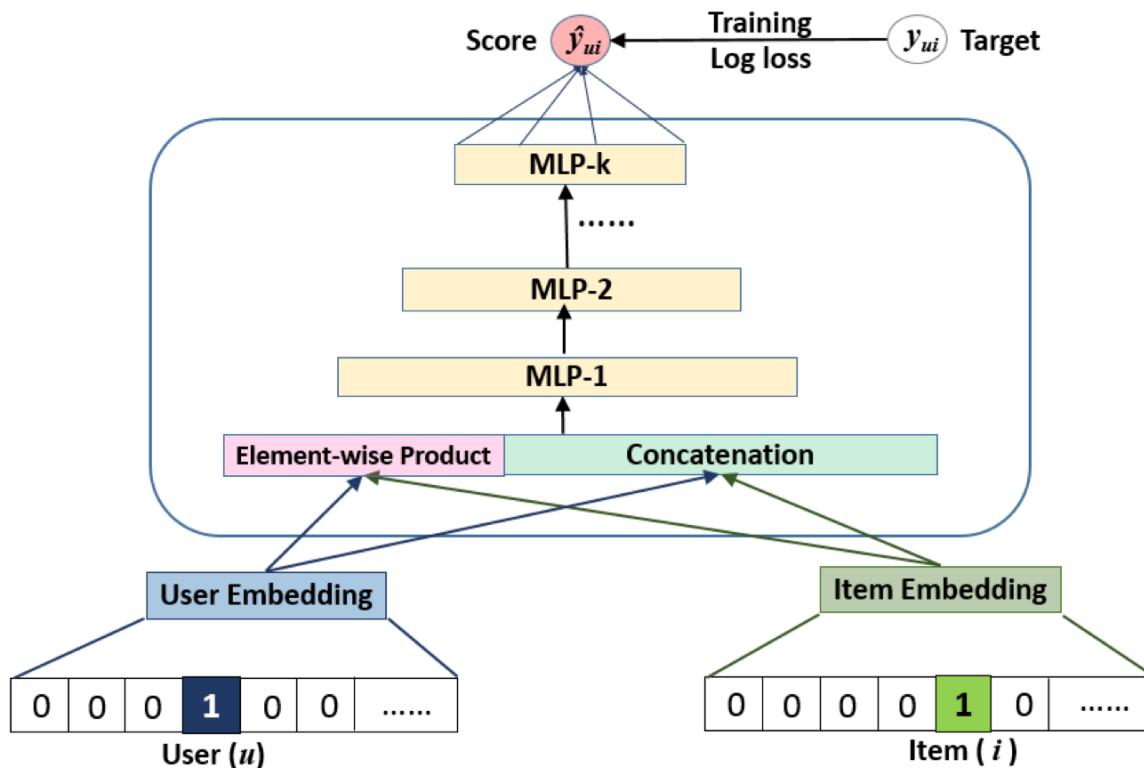


We have to design  
more effective models  
to make DNN work for  
CF!

**MLP can learn to approximate the low-rank relation, but is inefficient in doing so.**

# NeuMF: Neural Matrix Factorization (He et al, WWW'17)

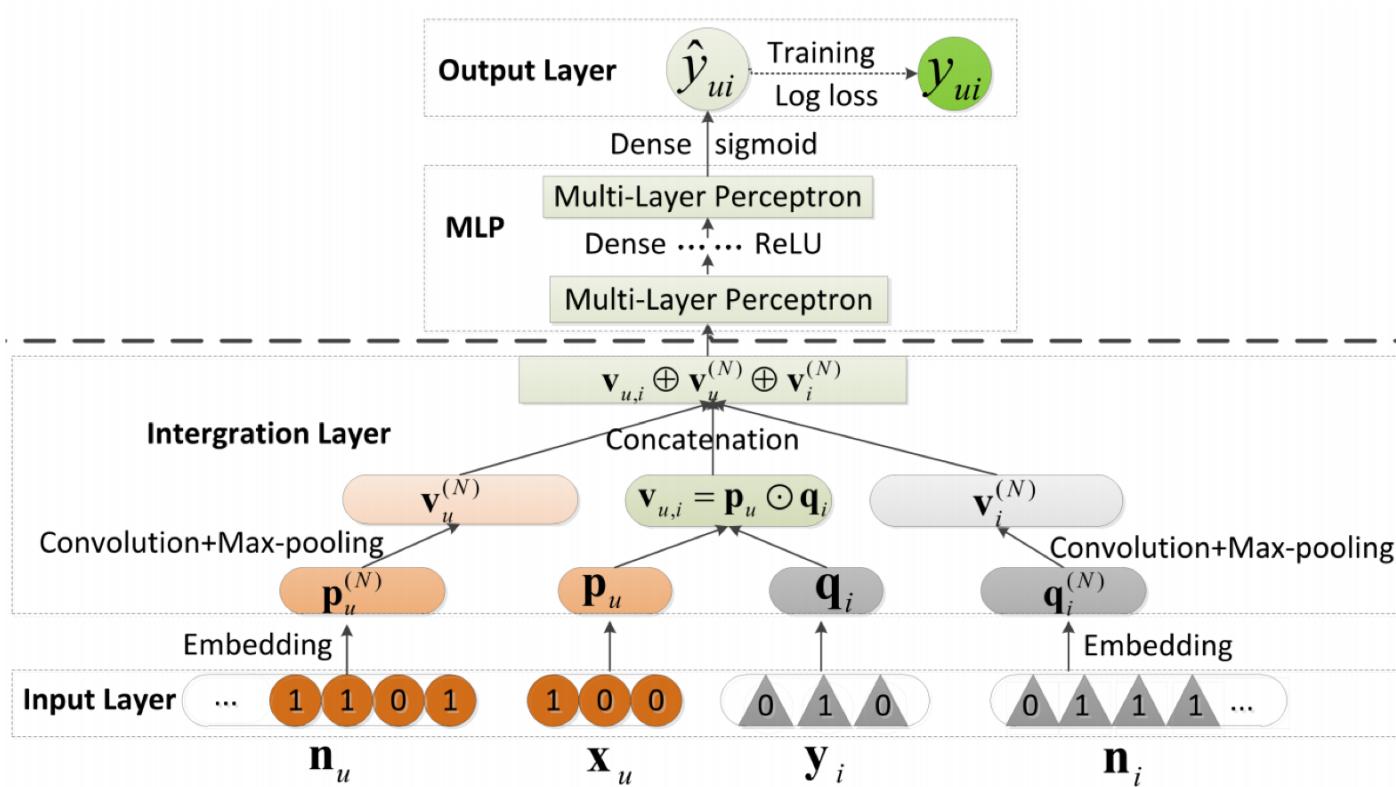
- NeuMF unifies the strengths of MF and MLP in learning the matching function:
  - MF uses inner product to capture the **low-rank** relation
  - MLP is more **flexible in using DNN** to learn the matching function.



# NNCF: Neighbor-based NCF

(Bai et al, CIKM'17)

- Feeding user and item neighbors into the NCF framework
  - Direct neighbors or community neighbors are considered.



# Experiment Evidence

Datasets	#Interaction	# Users	#Items	Sparsity
Delicious	437,593	1,867	69,223	99.66%
MovieLens	1,000,209	3,706	6,040	95.53%

*Performance Comparison on Item Recommendation (%)*

Datasets	Delicious		MovieLens	
Models	HR@5	NDCG@5	HR@5	NDCG@5
ItemPop	5.41	3.22	31.49	20.18
ItemKNN	59.69	55.90	45.01	30.14
MF-BPR	73.77	74.11	51.03	36.21
NeuMF	85.53	80.68	56.55	38.30
NNCF	87.31	84.58	62.00	42.21

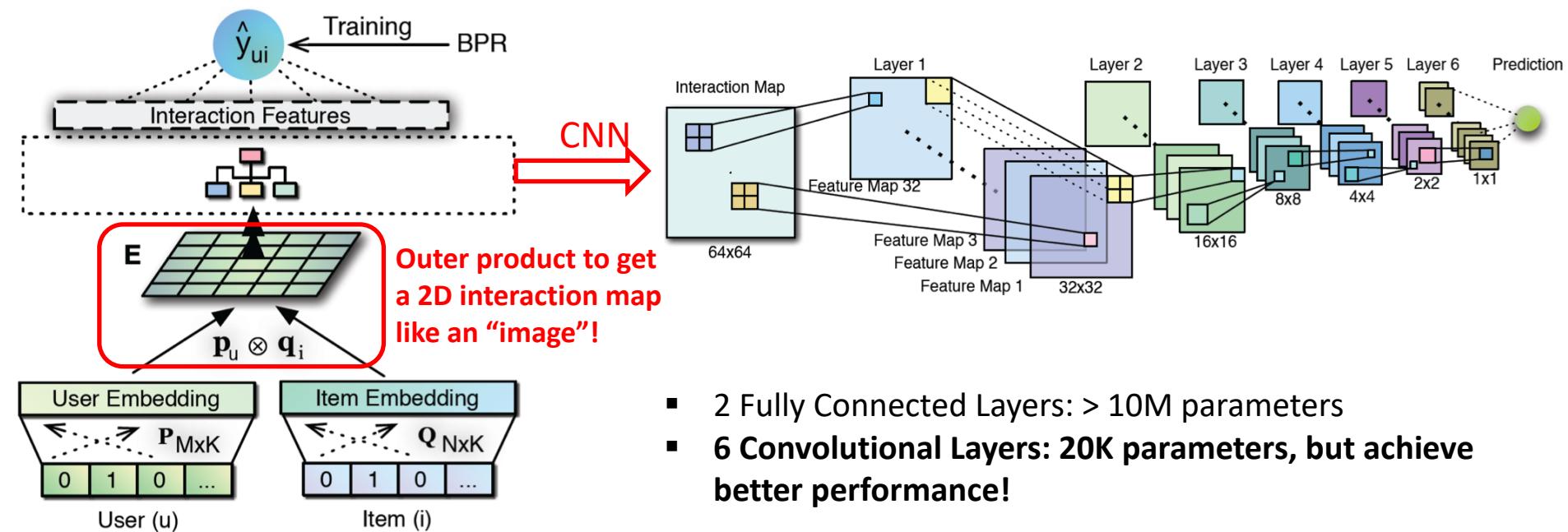
CF method is better than non-personalized method

Model-based CF is better than memory-based CF

Deep NCF models are better than shallow MF models by a large margin.

# Convolutional NCF (He et al, IJCAI'18)

- Although fully connected layers are popular in learning the matching function, they have too many parameters.
- Recently, we propose to use the **locally connected CNN** to build deeper NCF models.



- 2 Fully Connected Layers: > 10M parameters
- 6 Convolutional Layers: 20K parameters, but achieve better performance!

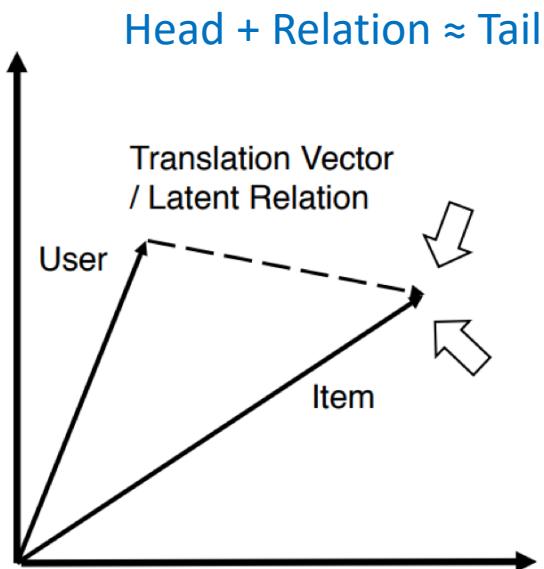
# Experiment Evidence

Datasets	#Interactions	#Users	#Items	Sparsity
Yelp	730,791	25,815	25,677	99.89%
Gowalla	1,249,703	54,156	52,400	99.95%

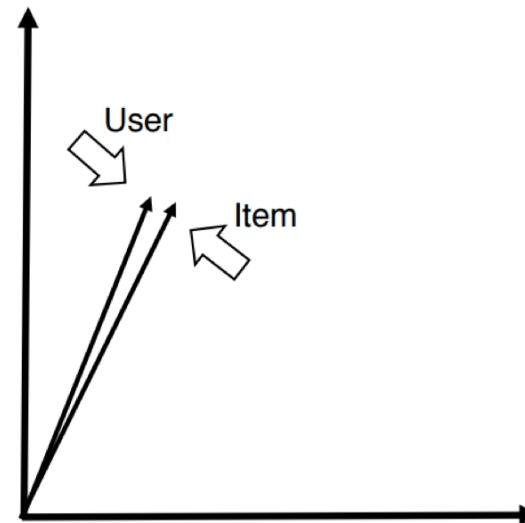
Datasets	Gowalla		Yelp		
	Models	HR@5	NDCG@5	HR@5	NDCG@5
ItemPop		20.03	10.99	7.10	3.65
MF-BPR		62.84	48.25	17.52	11.04
MLP		63.59	48.02	17.66	11.03
IRGAN		63.89	49.58	18.61	11.98
NeuMF		67.44	53.19	18.81	11.89
ConvNCF		69.14	54.94	19.06	12.09

ConvNCF are better than NeuMF and MLP with much fewer parameters.

# Overview of Translation-based Recommendation (Tay et al, WWW'18)



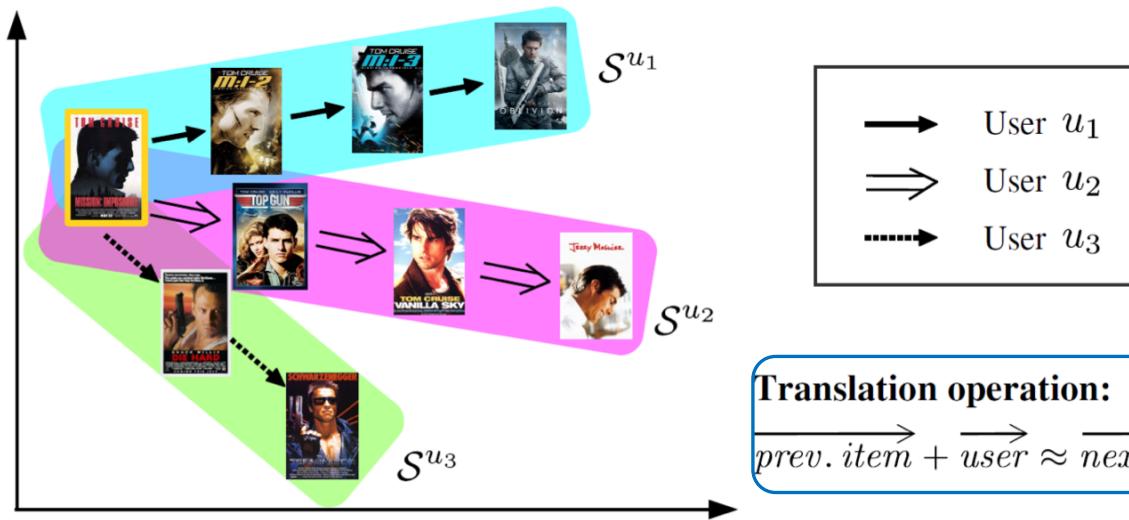
(a) Translation-based



(b) Matrix Factorization-based

# TransRec (He et al, Recsys'17)

- Focused on next-item recommendation
  - Third-order relationship between <user, current item, next item>
  - Current item is the “Relation”: **Head + Relation  $\approx$  Tail**



**Translation operation:**  
 $\xrightarrow{\text{prev. item}} + \xrightarrow{\text{user}} \approx \xrightarrow{\text{next item}}$

Predictive Model:  
 $\text{Prob}(j | u, i) \propto \beta_j - d(\vec{\gamma}_i + \vec{T}_u, \vec{\gamma}_j),$

Item bias      Translation distance

**Figure 1: TransRec as a sequential model:** Items (movies) are embedded into a ‘transition space’ where each user is modeled by a *translation* vector. The transition of a user from one item to another is captured by a user-specific *translation operation*.

# Latent Relational Metric Learning

## (Tay et al, WWW'18)

- Distance-based predictive model:

$$s(p, q) = \| p + r - q \|_2^2$$

where  $r$  is the **latent relation vector**, formed by an attentive sum over **memory vectors**:

$$r = \sum_i a_i m_i$$

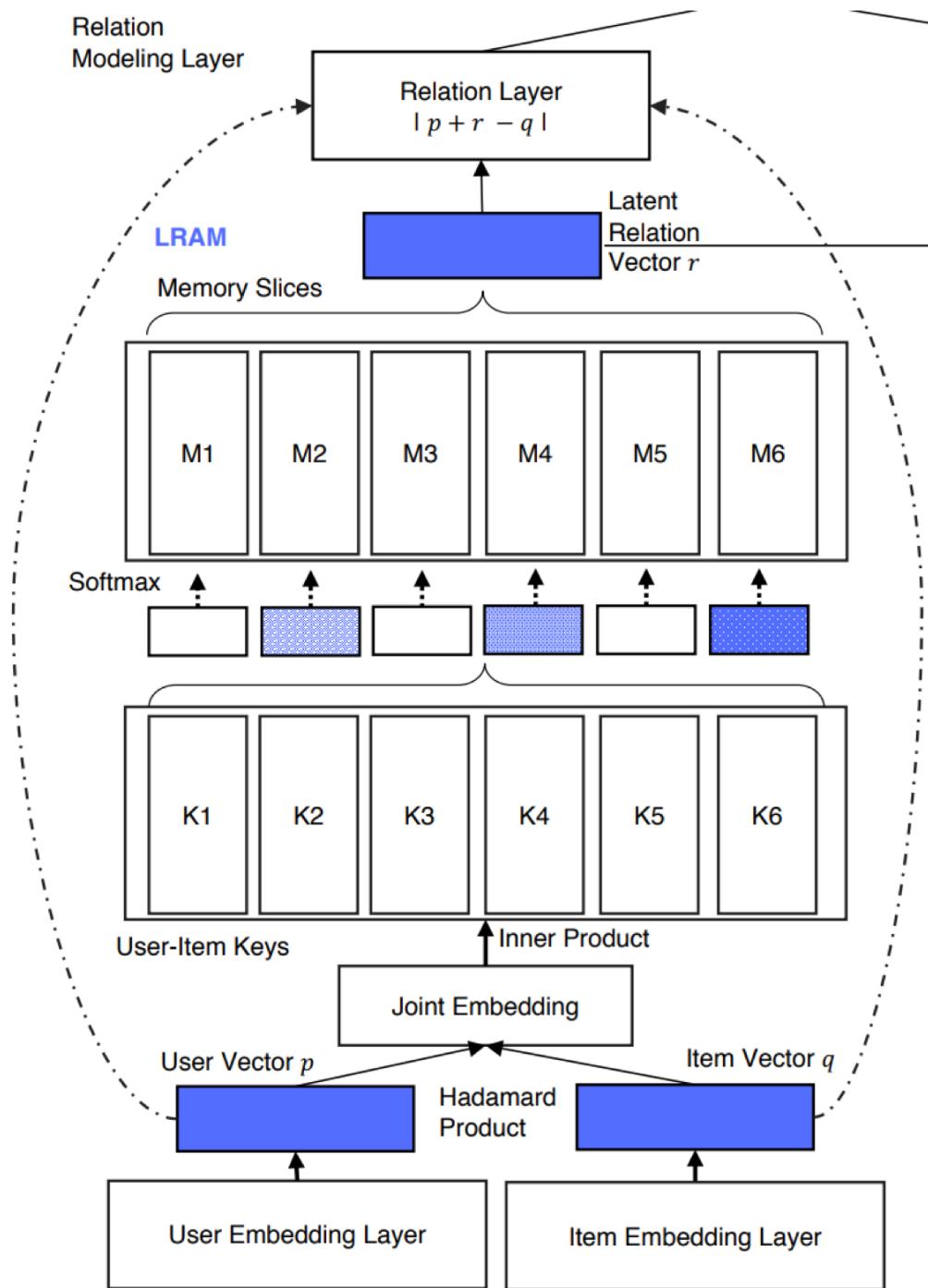
Attention weight, with inner product

$$s = p \odot q$$

as input.

Memory vector, which can encode user attributes/interest.

## Overview of the LRML's predictive model:



# Methods of Matching Function Learning

## 1. Pure CF models:

- Based on Neural Collaborative Filtering (NCF) framework:  
**NeuMF**: Neural Matrix Factorization (He et al, WWW'17)  
**NNCF**: Neighbor-based NCF (Bai et al, CIKM'17)  
**ConvNCF**: Outer Product-based NCF (He et al, IJCAI'18)
- Based on Translation framework:  
**TransRec**: Translation-based Recommendation (He et al, Recsys'17)  
**LRML**: Latent Relational Metric Learning (Tay et al, WWW'18)

## 2. Feature-based models:

- Based on Multi-Layer Perceptron:  
**Wide&Deep** (Cheng et al, DLRS'16), **Deep Crossing** (Shan et al, KDD'16)
- Based on Factorization Machines (FM):  
**Neural FM** (He and Chua, SIGIR'17), **Attentional FM** (Xiao et al, IJCAI'17),
- Based on Trees:  
**GB-CENT**: Categorical Embedding and Numerical Trees (Zhao et al, WWW'18)  
**DEF**: Deep Embedding Forest (Zhu et al, KDD'17)  
**TEM**: Tree-enhanced Embedding Model (Wang et al, WWW'18)

# Recall: Input to Feature-based Models

Feature vector $\mathbf{x}$												Target $y$						
$\mathbf{x}^{(1)}$	1	0	0	...	$\mathbf{x}^{(2)}$	1	0	0	0	...	$\mathbf{x}^{(3)}$	0.3	0.3	0.3	0	...	$y^{(1)}$	5
$\mathbf{x}^{(4)}$	1	0	0	...	$\mathbf{x}^{(5)}$	0	1	0	0	...	$\mathbf{x}^{(6)}$	0.3	0.3	0.3	0	...	$y^{(2)}$	3
$\mathbf{x}^{(7)}$	1	0	0	...	$\mathbf{x}^{(8)}$	0	0	1	0	...	$\mathbf{x}^{(9)}$	0.3	0.3	0.3	0	...	$y^{(3)}$	1
$\mathbf{x}^{(10)}$	0	1	0	...	$\mathbf{x}^{(11)}$	0	0	1	0	...	$\mathbf{x}^{(12)}$	0	0	0.5	0.5	...	$y^{(4)}$	4
$\mathbf{x}^{(13)}$	0	1	0	...	$\mathbf{x}^{(14)}$	0	0	0	1	...	$\mathbf{x}^{(15)}$	0	0	0.5	0.5	...	$y^{(5)}$	5
$\mathbf{x}^{(16)}$	0	0	1	...	$\mathbf{x}^{(17)}$	1	0	0	0	...	$\mathbf{x}^{(18)}$	0.5	0	0.5	0	...	$y^{(6)}$	1
$\mathbf{x}^{(19)}$	0	0	1	...	$\mathbf{x}^{(20)}$	0	0	1	0	...	$\mathbf{x}^{(21)}$	0.5	0	0.5	0	...	$y^{(7)}$	5
User				Movie				Other Movies rated										
A	B	C	...	TI	NH	SW	ST	...	TI	NH	SW	ST	...	Other Movies rated				

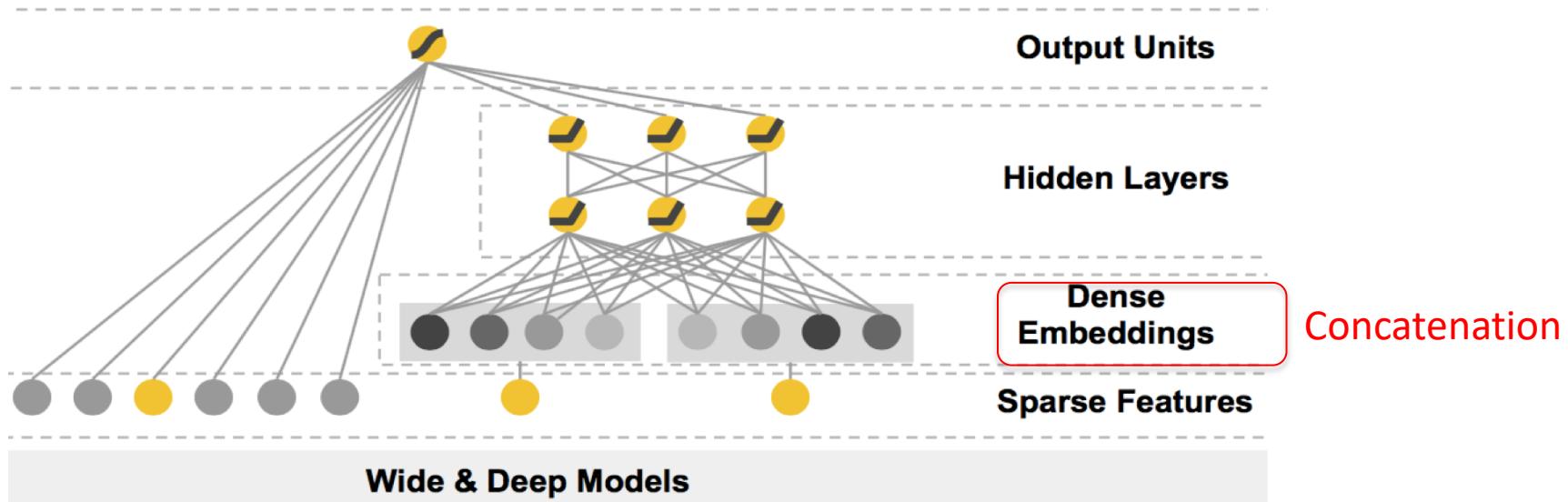
Raw features:

1. Categorical features  
One-hot encoding on ID features
2. Continuous features  
E.g., time, frequency.  
Need feature normalization

Transformed features:

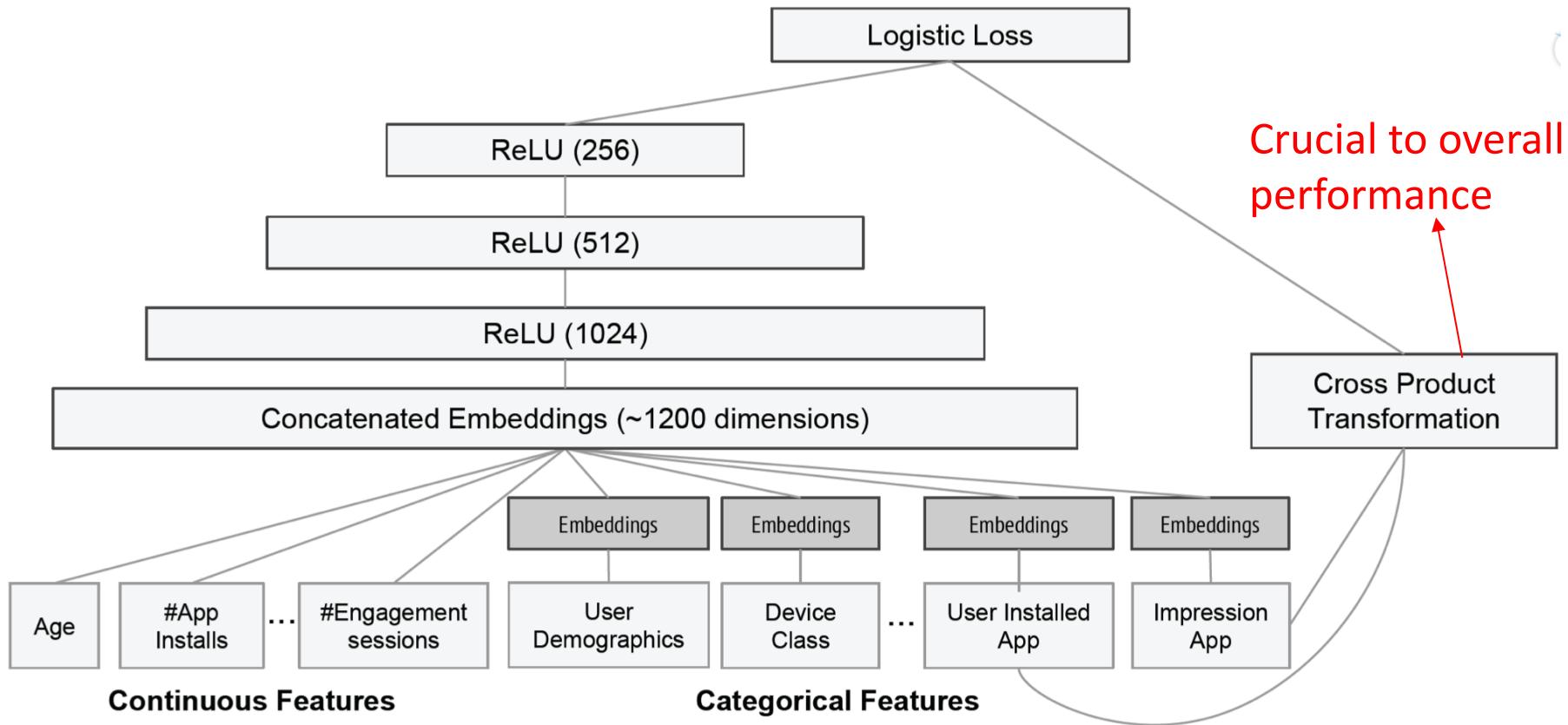
1. Categorical features  
Cross features are important
2. Continuous features  
E.g., outputs of other models like visual embeddings.

# Wide&Deep (Cheng et al, Recsys'16)



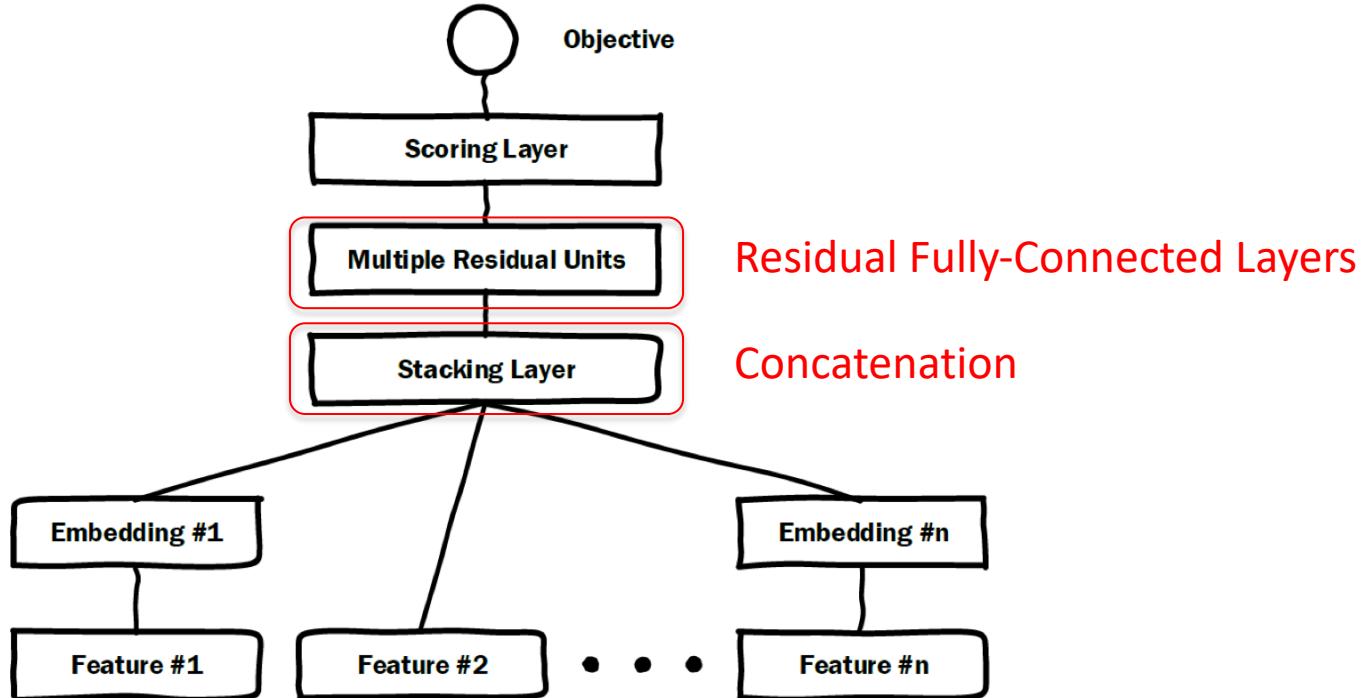
- The wide part is linear regression for memorizing seen feature interactions, which requires careful engineering on cross features.  
E.g.,  $AND(gender=female, language=en)$  is 1 iff both single features are 1
- The deep part is for generalizing to unseen feature interactions.

# Wide&Deep for App Recommendation (Cheng et al, Recsys'16)



# Deep Crossing (Shan et al, KDD'16)

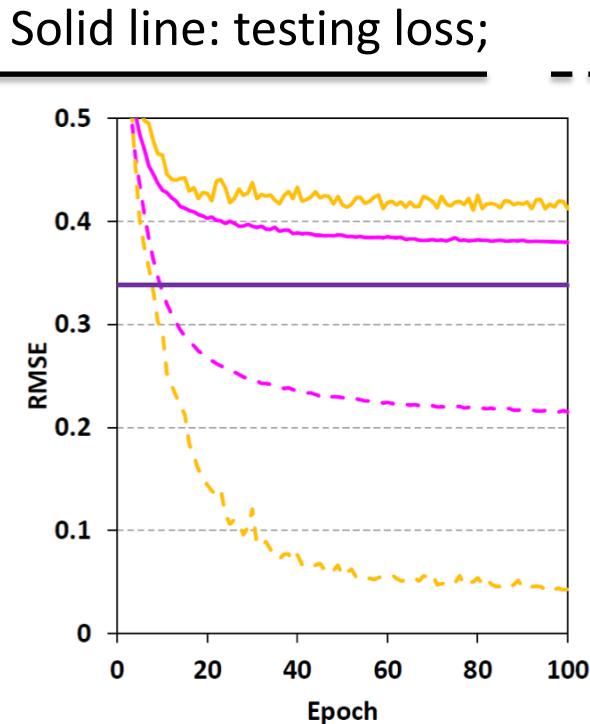
Microsoft's Sponsor Search Solution in 2016:



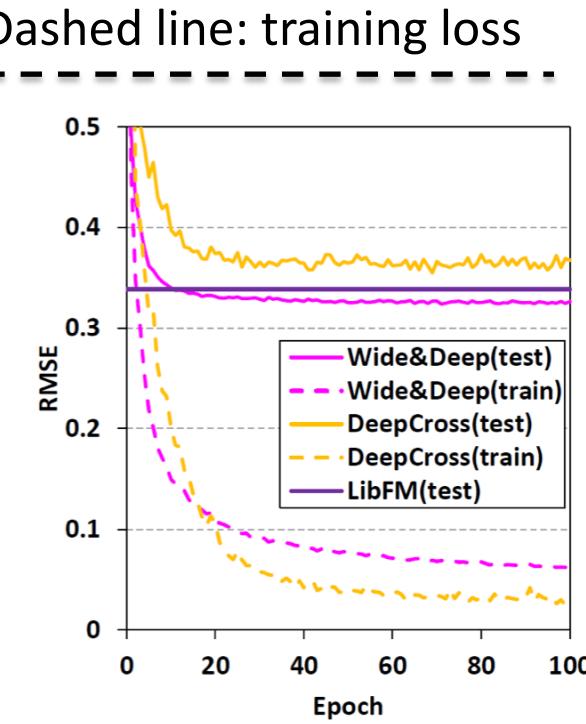
- The deep part can learn feature interactions in an implicit way.
- The use of residual layers makes the network be deep.

# Empirical Evidence

- However, when only raw categorical features are used, both DL models underperform the shallow FM in learning unseen feature interactions.



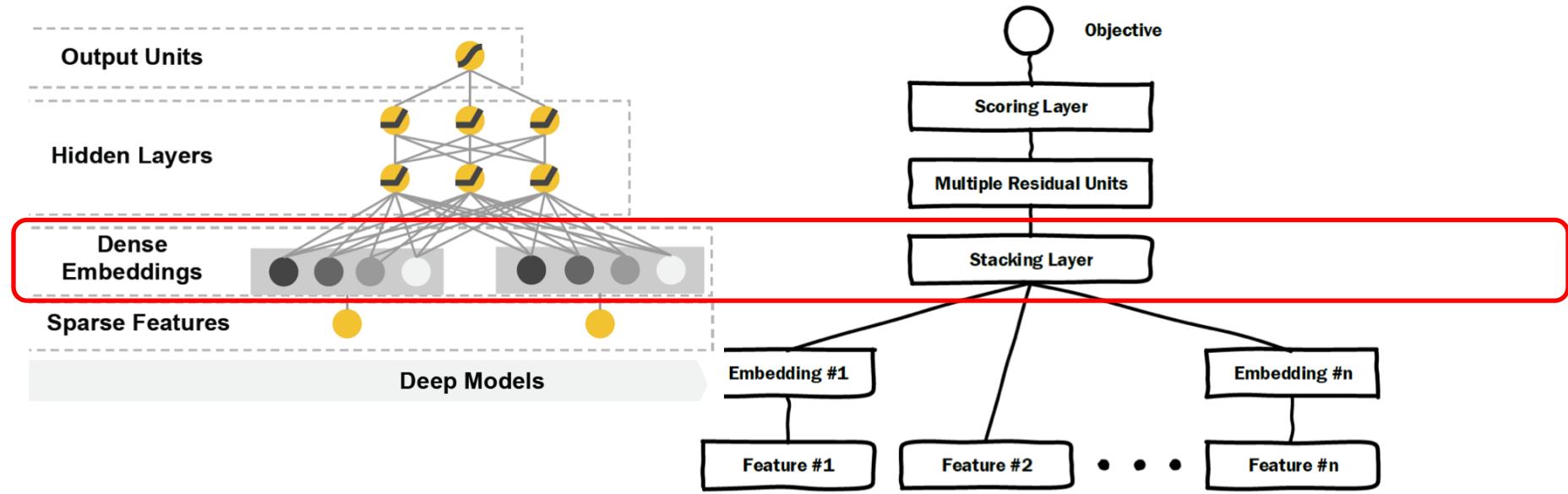
(a) Random initialization



(b) FM as pre-training

# Why MLP is Ineffective?

Besides optimization difficulties, one reason might be the model design:



1. Embedding concatenation carries little information about feature interactions in the low level!
2. The structure of Concat+MLP is ineffective to learn the multiplicative relation (Beutel et al, WSDM'18).

# NFM: Neural Factorization Machine (He and Chua, SIGIR'18)

- Inspired by FM, NFM models pairwise interactions between feature embeddings with multiplication.

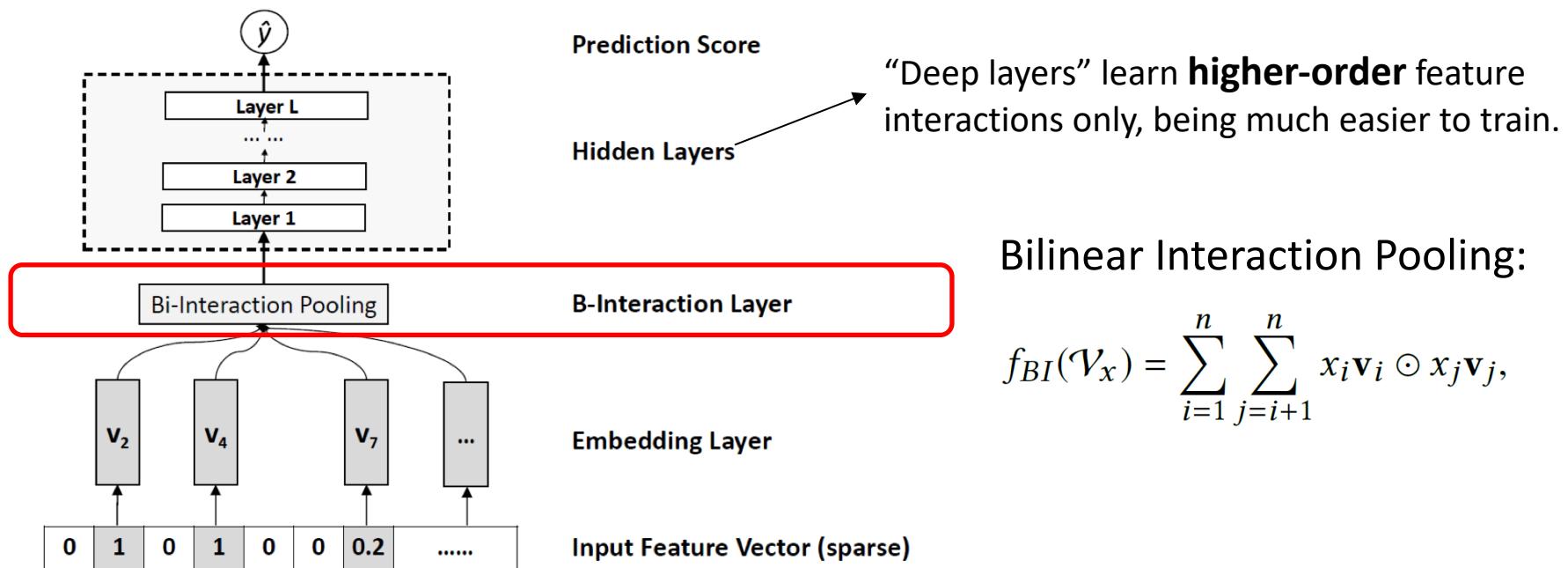


Figure 2: Neural Factorization Machines model (the first-order linear regression part is not shown for clarity).

# Experiment Evidence

Task #1: Context-aware App Usage Prediction  
- Frappe data: instance #: 288,609, feature #: 5,382

Task #2: Personalized Tag Recommendation  
- MovieLens data: Inst #: 2,006,859, Feat #: 90,445

**Table: Parameter # and testing RMSE at embedding size 128**

Method	Frappe		MovieLens	
	Param#	RMSE	Param#	RMSE
Logistic Regression	5.38K	0.5835	0.09M	0.5991
FM	1.38M	0.3385	23.24M	0.4735
High-order FM	2.76M	0.3331	46.40M	0.4636
Wide&Deep (3 layers)	4.66M	0.3246	24.69M	0.4512
DeepCross (10 layers)	8.93M	0.3548	25.42M	0.5130
<b>Neural FM (1 layer)</b>	<b>1.45M</b>	<b>0.3095</b>	<b>23.31M</b>	<b>0.4443</b>

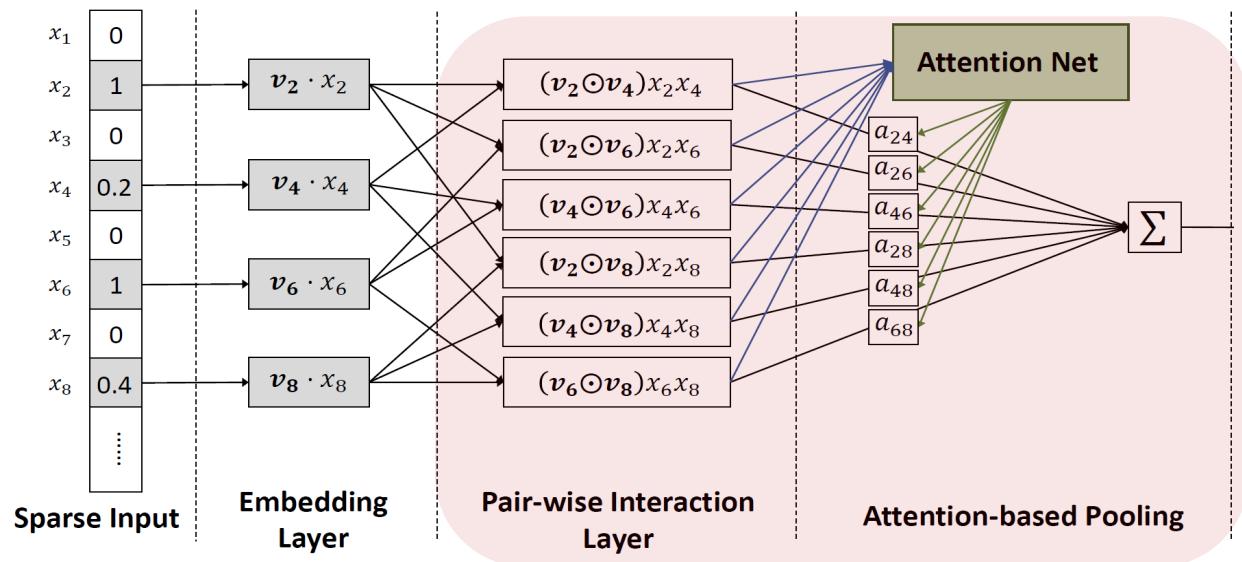
Codes: [github.com/hexiangnan/neural\\_factorization\\_machine](https://github.com/hexiangnan/neural_factorization_machine)

1. Shallow embedding methods learn interactions, better than simple linear models
2. Deep embedding methods:  
**Wide&Deep = Concat+3 layers**  
**DeepCross = Concat+10 layers**
3. Our methods:  
**Neural FM = BI pooling + 1 layer**  
Shallower but outperforming existing deeper methods with less parameters.

# AFM: Attentional Factorization Machine

(Xiao et al, IJCAI'18)

- Neural FM treats all second-order feature interactions as contributing equally.
- Attentional FM uses an attention network to learn the weight of a feature interaction.



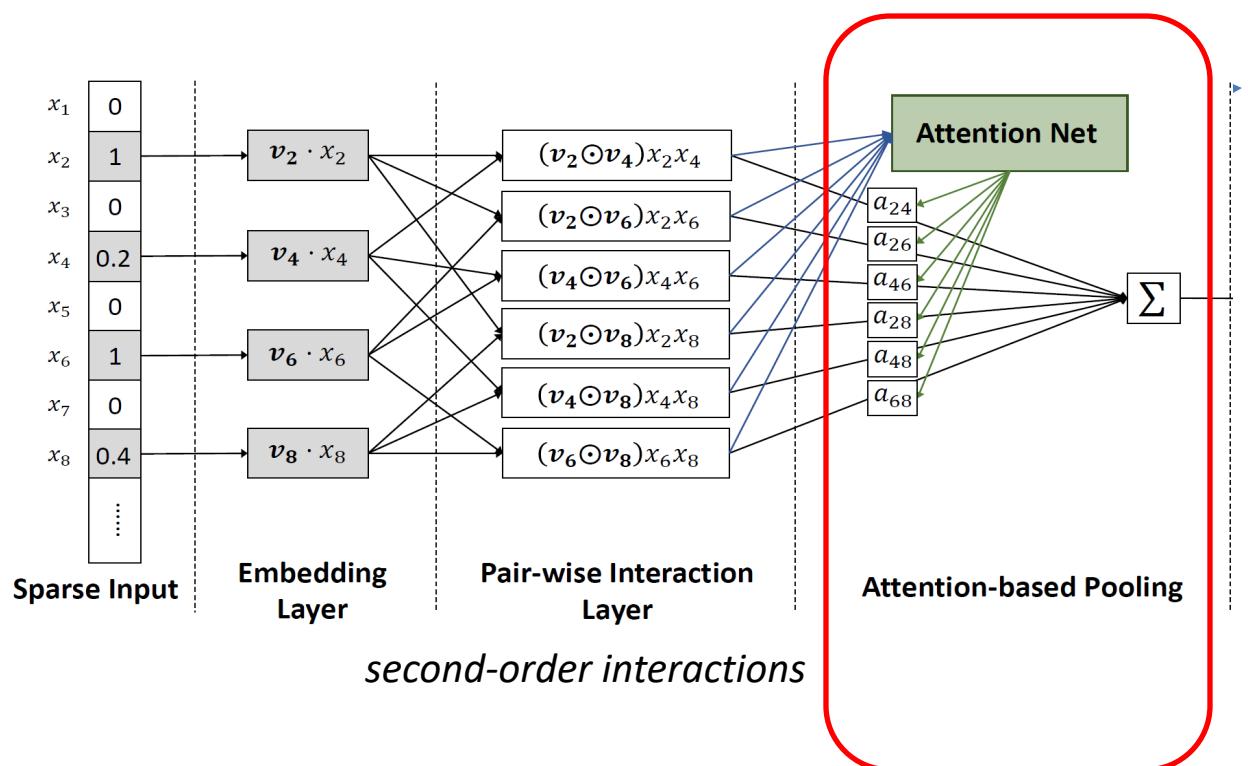
$$f_{ABI}(\mathcal{V}_x) = \sum_{i=1}^n \sum_{j=i+1}^n (x_i \mathbf{v}_i \odot x_j \mathbf{v}_j) a_{ij}$$

$$a'_{ij} = \mathbf{h}^T \text{ReLU}(\mathbf{W}(\mathbf{v}_i \odot \mathbf{v}_j)x_i x_j + \mathbf{b}),$$

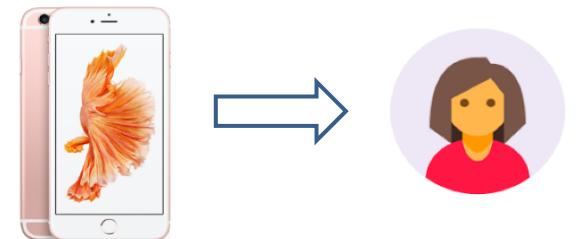
$$a_{ij} = \frac{\exp(a'_{ij})}{\sum_{(i,j) \in \mathcal{R}_x} \exp(a'_{ij})},$$

# Explaining Recommendation with AFM

The attention scores can be used to select the most predictive second-order feature interactions as explanations.



Example: explainable recommendation with second-order cross features:



<Female, Age 20>  
<Age 20, iPhone>  
<Female, Color Pink>  
.....

# Experiment Evidence

Task #1: Context-aware App Usage Prediction

- Frappe data: instance #: 288,609, feature #: 5,382

Task #2: Personalized Tag Recommendation

- MovieLens data: Inst #: 2,006,859, Feat #: 90,445

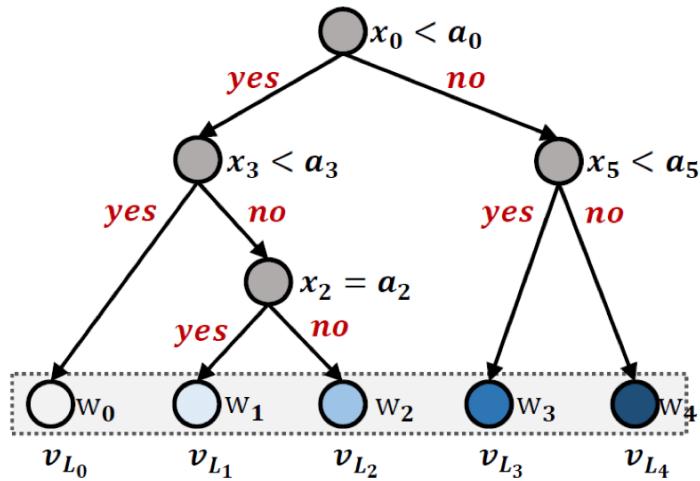
**Table: Parameter # and testing RMSE at embedding size 128**

Method	Frappe		MovieLens	
	Param#	RMSE	Param#	RMSE
Logistic Regression	5.38K	0.5835	0.09M	0.5991
FM	1.38M	0.3385	23.24M	0.4735
High-order FM	2.76M	0.3331	46.40M	0.4636
Wide&Deep (3 layers)	4.66M	0.3246	24.69M	0.4512
DeepCross (10 layers)	8.93M	0.3548	25.42M	0.5130
<b>Neural FM (1 layer)</b>	<b>1.45M</b>	<b>0.3095</b>	<b>23.31M</b>	<b>0.4443</b>
<b>Attentional FM (0 layer)</b>	<b>1.45M</b>	<b>0.3102</b>	<b>23.26M</b>	<b>0.4325</b>

AFM without hidden layers can be better than NFM with 1 hidden layer.

Codes: [github.com/hexiangnan/attentional\\_factorization\\_machine](https://github.com/hexiangnan/attentional_factorization_machine)

# Tree-based Model



Working mechanism of tree-based models:

- Each node **splits** a feature into two decision edges according to a value.
- Given a feature vector, there exists a **path** from the root to a leaf, which forms a decision rule (like a cross feature).
- The leaf node corresponds to the **prediction value**.

E.g., meaning of a path:  $v_{L_1}$ : [Age < 18] & [Country ≠ France] & [Restaurant Tag = French].

- Since a single tree may not be expressive enough, a typical way is to build a forest, i.e., an ensemble of multiple trees:

$$\hat{y}_{GBDT}(\mathbf{x}) = \sum_{s=1}^S \hat{y}_{DT_s}(\mathbf{x}),$$

S → # of trees  
s → Prediction of the s-th tree

# Tree-based vs. Embedding-based Model

Tree-based Model (e.g., GBDT)	Embedding-based Model (e.g., DNN, FM)
+ Strong at continuous features	+ Strong at categorical features
+ Explainable	- Blackbox
+ Low serving cost	- High serving cost
- Weak generalization ability to unseen feature combinations.	+ Strong generalization ability to unseen feature combinations.

**Why not combining the strengths of the two types of models?**

*In the next:*

- Gradient Boosted Categorical Embedding and Numerical Trees (Zhao et al, WWW'17)
- Deep Embedding Forest (Zhu et al, KDD'17)
- Tree-enhanced Embedding Model (Wang et al, WWW'18)

# GB-CENT: Gradient Boosted Categorical Embedding and Numerical Trees (Zhao et al, WWW'17)

- GB-CENT unifies the strengths of embeddings in **categorical feature learning** and trees in **continuous feature learning**.
  - SVDFeature is applied on categorical features.
  - GBDT is applied on continuous features.

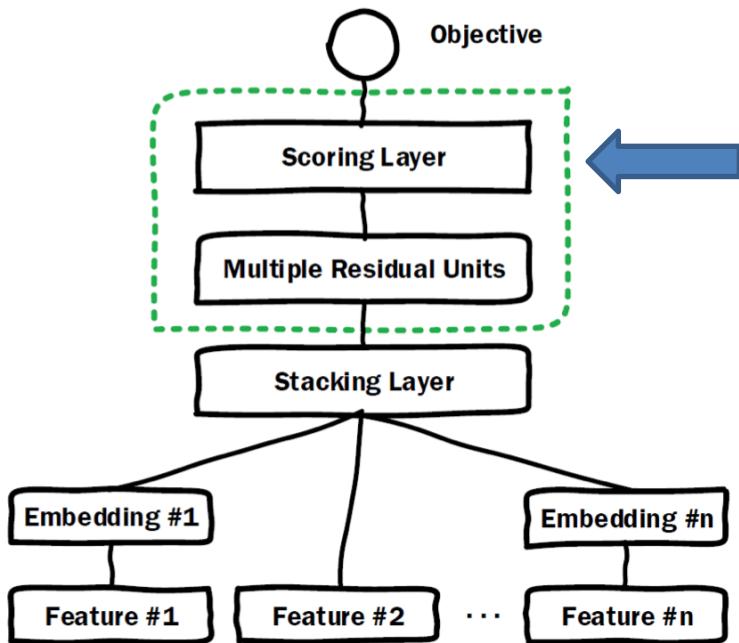
$$\hat{y}(x) = \underbrace{\sum_{i=0}^k w_{a_i}}_{\text{bias}} + \underbrace{\left( \sum_{a_i \in U(a)} Q_{a_i} \right)^T \left( \sum_{a_i \in I(a)} Q_{a_i} \right)}_{\text{SVDFeature}} + \underbrace{\sum_{i=0}^k T_{a_i}(b)}_{\text{GBDT}}$$

- Each categorical feature corresponds to a tree  
(i.e., # of categorical features = # of trees)

# Deep Embedding Forest

## (Zhao et al, KDD'17)

- DEF uses **forest** (e.g., LightGBM or XGBoost) as the hidden layers to reduce the online serving time of embedding-based models.



Deep Crossing (Shan et al, KDD'16)

Using Forest Layer instead.

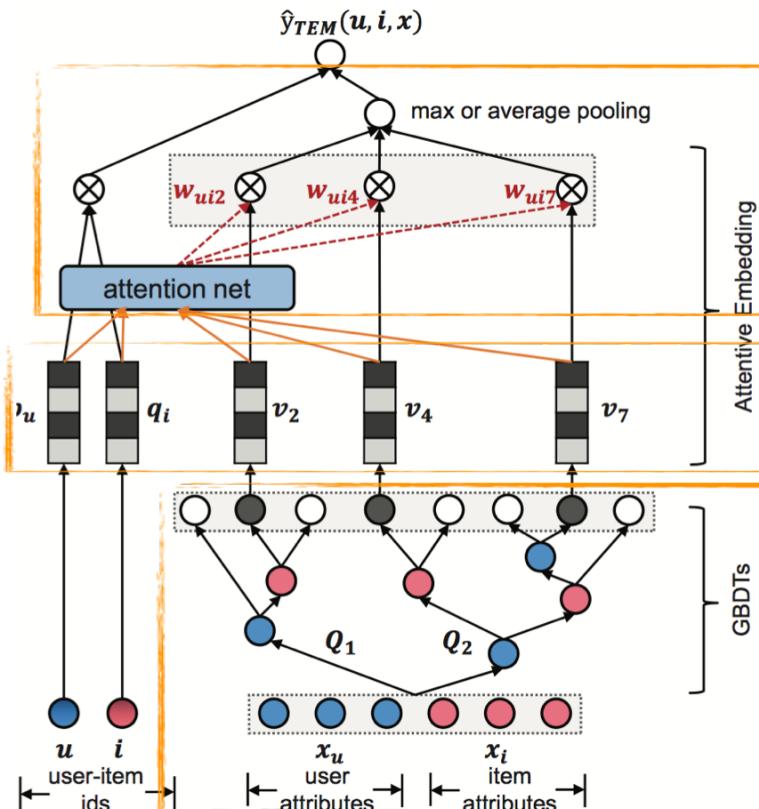
- Two step training
- Initialize DEF using Deep Crossing

Experiment evidence

Methods	Relative Log Loss	Time(ms)
Deep Crossing	100	2.272
DEF (XGBoost)	99.96	0.168
DEF (LightGBM)	99.94	0.204

# Tree-enhanced Embedding Model (Wang et al, WWW'18)

- TEM **explicitly** learns which **cross features** are more important for a **<user, item>** prediction.



## Attention layer

- Based on the embeddings, it learns the personalized weights of cross features

## Embedding layer:

- User ID, item ID, and cross features are embedded.

## Forest layer for cross feature learning:

- User attributes and item attributes are fed into a GBDT.
- Activated leaf nodes denote prominent cross features for an input feature vector.

# Experiment Evidence

- Datasets: TripAdvisor in London, New York City, & Singapore:
  - Profile of Users
  - Public Description of POIs/Restaurants
  - User-POI/Restaurant Interactions



- Age: Female
- Gender: 25 - 34
- Country: China
- City: Beijing
- Traveler Styles: Foodies, History Buff, Urban Explorer, Art & Architecture Lover, Like a Local



National Gallery, LON

- Attributes: Art Museums, Museums
- Tags: Trafalgar Square, Van Gogh, Sainsbury Wing, Da Vinci, Famous Paintings, Special Exhibitions, Beautiful Buildings
- Rating: 4.5

Table: Logloss of predictive models (the lower, the better)

Dataset	LON-A			NYC-R		
Method	logloss	ndcg@5	p-value	logloss	ndcg@5	p-value
XGBoost	0.1251	0.6785	8e-5	0.1916	0.3943	4e-5
GBDT+LR	0.1139	0.6790	2e-4	0.1914	0.3997	4e-4
GB-CENT	0.1246	0.6784	6e-5	0.1918	0.3995	4e-5
FM	0.0939	0.6809	1e-2	0.1517	0.4018	5e-5
NFM	0.0892	0.6812	2e-2	0.1471	0.4020	8e-4
<b>TEM-avg</b>	0.0818	0.6821	-	0.1235	0.4019	-
<b>TEM-max</b>	<b>0.0791</b>	<b>0.6828</b>	-	<b>0.1192</b>	<b>0.4038</b>	-

TEM outperforms pure embedding-based methods & generates personalized reasons (cross features) for a recommendation.

# Short Summary

- Feature interaction learning (i.e., cross feature effect) is crucial for matching function learning in recommendation.
- Many models have been explored, e.g., DNN, FM, Tree-based, Attention Net etc.
- One insight is that doing early cross on raw features (or feature embeddings) is important to performance. E.g.,
  - Wide&Deep do manual cross on raw features
  - FM-based methods do second-order cross on feature embeddings
  - Tree-based methods do trainable cross on raw features.
- It remains challenging to build **explainable** matching function with strong generalization ability.
  - I.e., explainable high-order interaction learning.

# References

- Paul Covington, Jay Adams, and Emre Sargin. Deep neural networks for youtube recommendations. In Recsys 2016.
- Xiang Wang, Xiangnan He, Liqiang Nie, and Tat-Seng Chua. Item silk road: Recommending items from information domains to social users. In SIGIR 2017.
- Hong-Jian Xue, Xin-Yu Dai, Jianbing Zhang, Shujian Huang, and Jiajun Chen. Deep matrix factorization models for recommender systems. IJCAI 2017.
- Suvash Sedhain, Aditya Krishna Menon, Scott Sanner, and Lexing Xie. Autorec: Autoencoders meet collaborative filtering. In WWW 2015.
- Yao Wu, Christopher DuBois, Alice X. Zheng, and Martin Ester. Collaborative denoising auto-encoders for top-n recommender systems. In WSDM 2016.
- Sheng Li, Jaya Kawale, and Yun Fu. Deep collaborative filtering via marginalized denoising auto-encoder. In CIKM 2015.
- Xue Geng, Hanwang Zhang, Jingwen Bian, and Tat-Seng Chua. Learning image and user features for recommendation in social networks. In ICCV 2015.
- Jingyuan Chen, Hanwang Zhang, Xiangnan He, Liqiang Nie, Wei Liu, and Tat-Seng Chua. Attentive collaborative filtering: Multimedia recommendation with item-and component-level attention. In SIGIR 2017.
- Fuzheng, Zhang, Nicholas Jing Yuan, Defu Lian, Xing Xie, and Wei-Ying Ma. Collaborative knowledge base embedding for recommender systems. In KDD 2016.
- Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. Neural collaborative filtering. In WWW 2017.
- Ting Bai, Ji-Rong Wen, Jun Zhang, and Wayne Xin Zhao. A Neural Collaborative Filtering Model with Interaction-based Neighborhood. CIKM 2017.

# References

- Xiangnan He, Xiaoyu Du, Xiang Wang, Feng Tian, Jinhui Tang, and Tat-Seng Chua. Out Product-based Neural Collaborative Filtering. In IJCAI 2018.
- Ruining He, Wang-Cheng Kang, and Julian McAuley. Translation-based Recommendation. In Recsys 2017.
- Yi Tay, Luu Anh Tuan, and Siu Cheung Hui. Latent Relational Metric Learning via Memory-based Attention for Collaborative Ranking. In WWW 2018.
- Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishi Aradhye, Glen Anderson et al. Wide & deep learning for recommender systems. In DLRS 2016.
- Ying Shan, T. Ryan Hoens, Jian Jiao, Haijing Wang, Dong Yu, and J. C. Mao. Deep crossing: Web-scale modeling without manually crafted combinatorial features. In KDD 2016.
- Xiangnan He, and Tat-Seng Chua. Neural factorization machines for sparse predictive analytics. In SIGIR 2017.
- Jun Xiao, Hao Ye, Xiangnan He, Hanwang Zhang, Fei Wu, and Tat-Seng Chua. Attentional factorization machines: Learning the weight of feature interactions via attention networks. IJCAI 2017.
- Qian Zhao, Yue Shi, and Liangjie Hong. GB-CENT: Gradient Boosted Categorical Embedding and Numerical Trees. In WWW 2017.
- Jie Zhu, Ying Shan, J. C. Mao, Dong Yu, Holakou Rahamanian, and Yi Zhang. Deep embedding forest: Forest-based serving with deep embedding features. In KDD 2017.
- Xiang Wang, Xiangnan He, Fuli Feng, Liqiang Nie, and Tat-Seng Chua. TEM: Tree-enhanced Embedding Model for Explainable Recommendation. WWW 2018.

# Outline of Tutorial

- Unified view of matching in search and recommendation
- Part 1: Traditional Approaches to Matching
- Part 2: Deep Learning Approaches to Matching
- Summary

# Summary

- Search and Recommendation are two sides of the same coin
  - Search -> *Information Pull* with *explicit* info request (query)
  - Recommendation -> *Information Push* with *implicit* info request (user profile, contexts)
- Technically, they can be unified under the same matching view
  - Though they are studied by different communities: SIGIR vs. RecSys
- Deep learning-based matching methods
  - Representation learning-focused
  - Matching function learning-focused
- Matching is a generic problem for a wide range of applications
  - E.g., online advertising, question answering, image annotation, drug design

# Challenges

- Data: building better **benchmarks**
  - Large-scale text matching data
  - Large-scale user-item matching data with rich attributes.
- Model: data-driven + **knowledge-driven**
  - Most current methods are purely data-driven
  - Prior information (e.g., domain knowledge, large-scale knowledge based) is helpful and should be integrated into data-driven learning in a principled way.
- Task: **multiple criteria**
  - Existing work have primarily focused on similarity
  - Different application scenarios should have different matching goals
  - Other criteria such as novelty, diversity, and explainability should be taken into consideration

# Thanks!