

IN4325



Neural IR

Arthur Câmara (WIS, TU Delft)

Slides (mostly) by Claudia Hauff

39 Full Papers. (23% acceptance rate)

- 3 contain “neural” in title
- 3 contain “embedding” in title
- 11 others *may* contain neural “stuff”, given title.

~43% of papers with some type of neural IR.

Ian Goodfellow @goodfellow_ian

Usually, at least one of the baselines should be a result published in another paper, where the authors of that other paper had some incentive to get a good result. This way the evaluations are at least incentive-compatible.

Show this thread



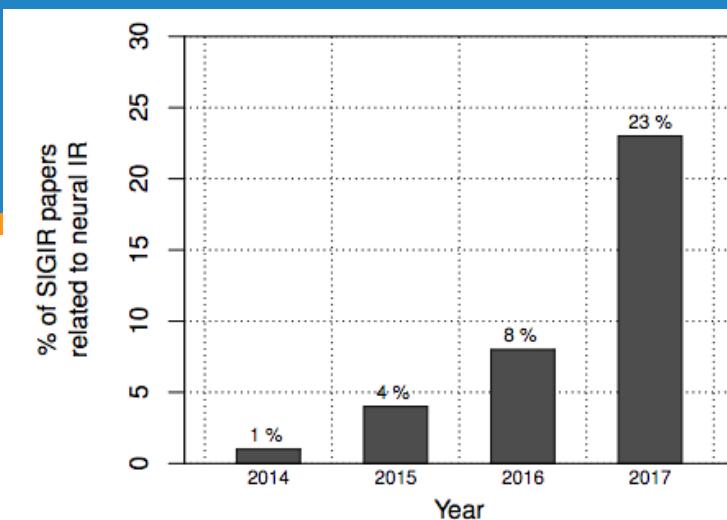
Different fields have different standards ...

A great starting point for Neural
IR - 123 pages of insights!
We follow some of it here.

Neural IR

... is not taking over the IR research field
(yet)

Covered on a high level: this topic by itself could take up all whole quarter of lectures. Unfortunately, we do not have the time.



GoogLeNet
(2014)

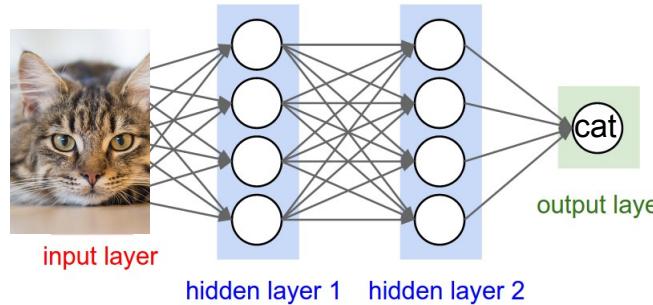
Overview

NeuralIR is the application of **shallow** or **deep** neural networks to IR tasks.

NeuralIR can make use of neural NLP techniques to enhance the retrieval Pipeline.

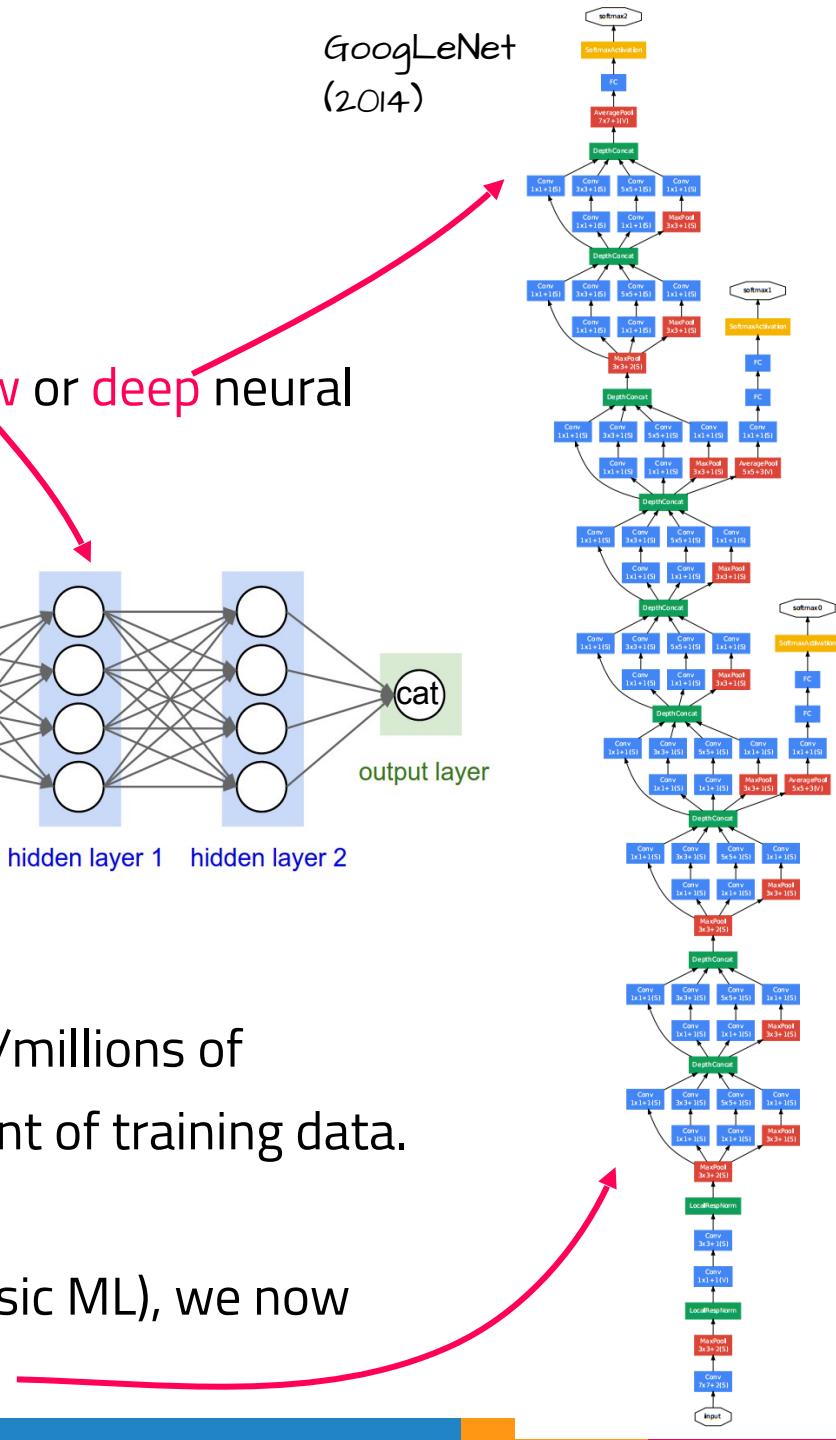


#params
LM, BM25

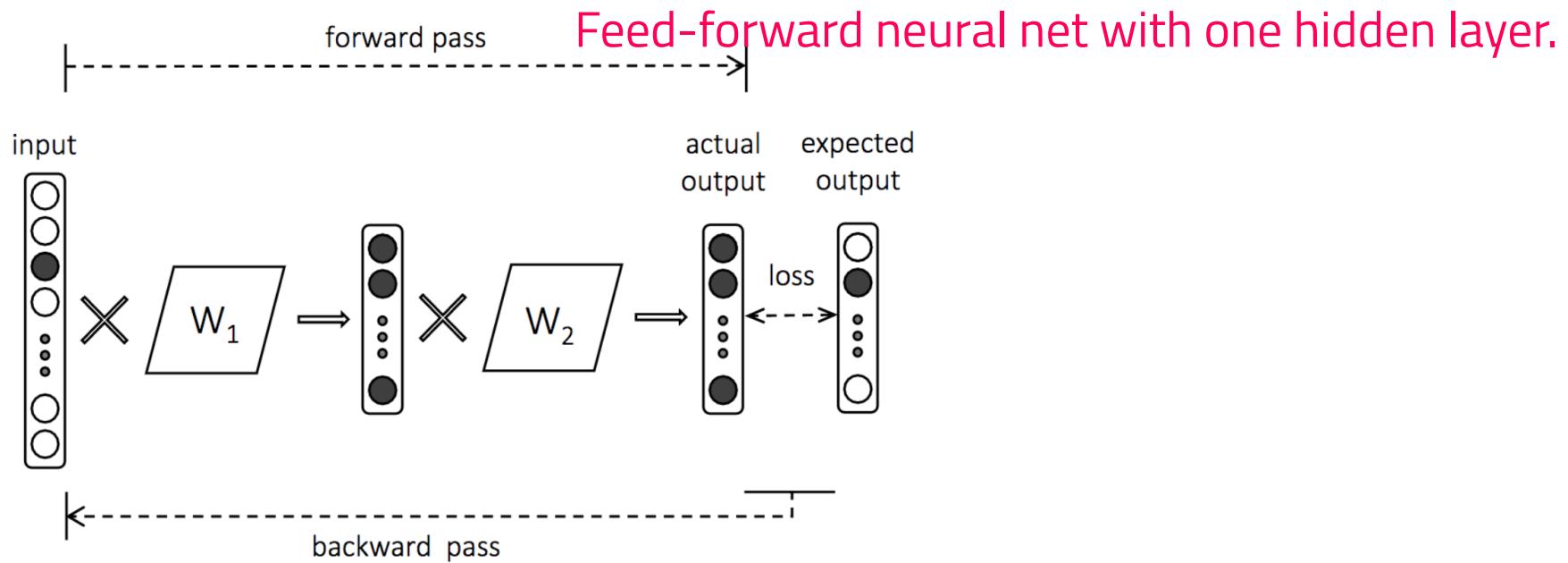


NeuralIR models contain thousands/millions of parameters. They require a large amount of training data.

Instead of hand-crafting **features** (classic ML), we now handcraft NN **architectures**



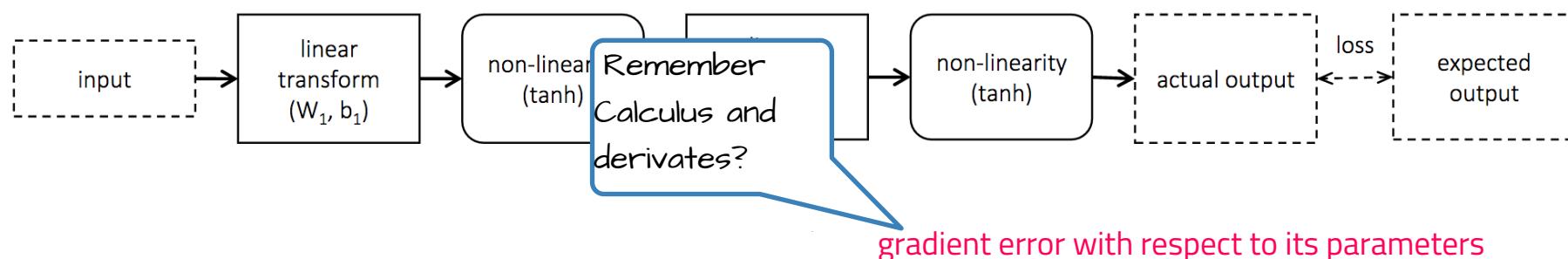
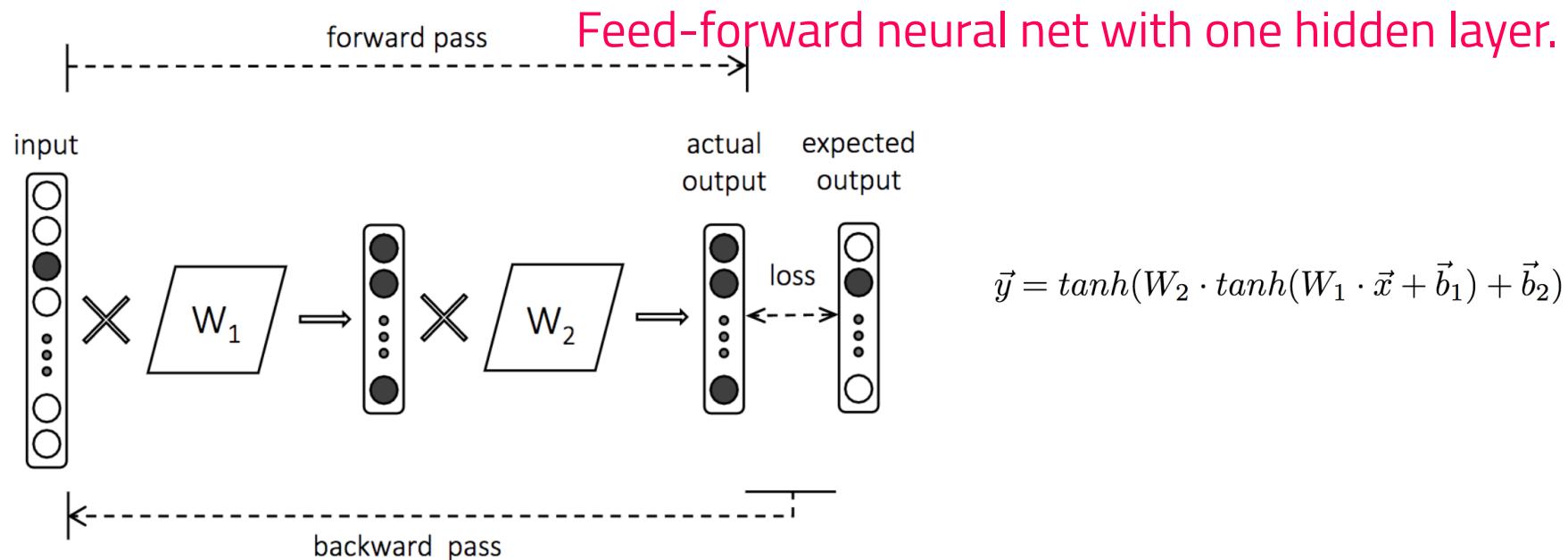
Neural net basics



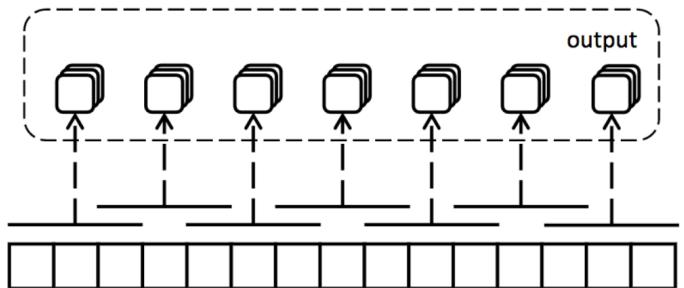
Initially: deep nets save us from expensive and elaborate feature engineering

Now: lets engineer architectures and hardware/software that allows us to efficiently train deep nets.

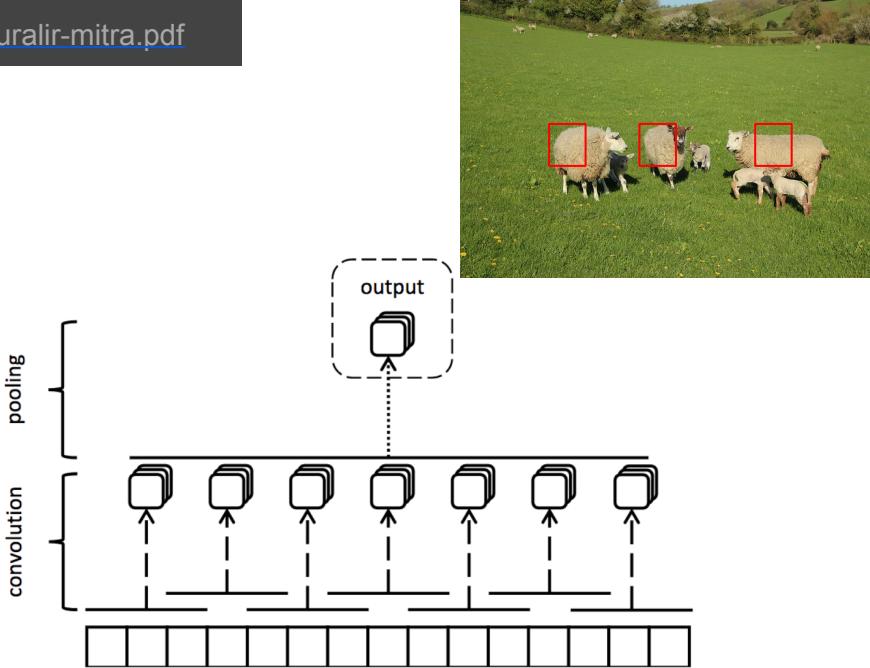
Neural net basics



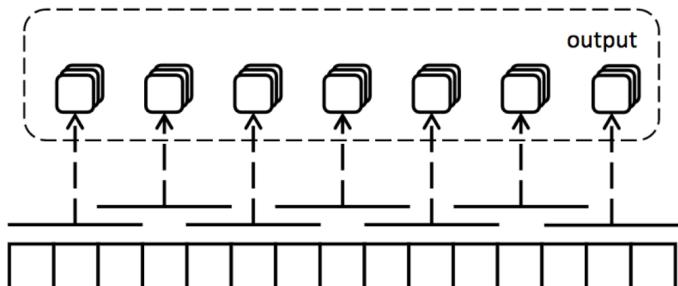
CNNs and RNNs



CNN: convolutional neural network.
Most often found in computer vision setups.

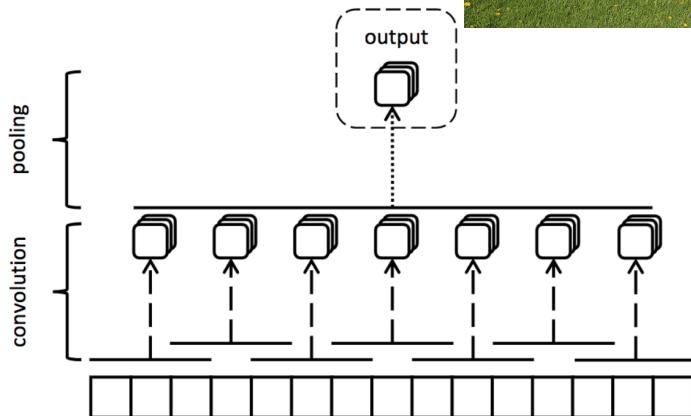


CNNs and RNNs

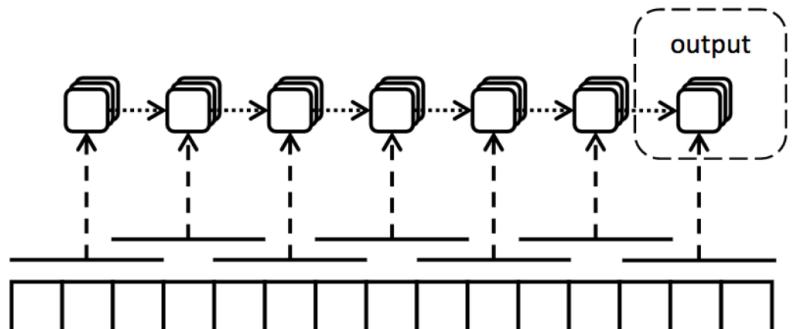


CNN: convolutional neural network.
Most often found in computer vision setups.

* Maybe should be revisited? (ICLR'18): <https://openreview.net/pdf?id=BJEX-H1Pf>

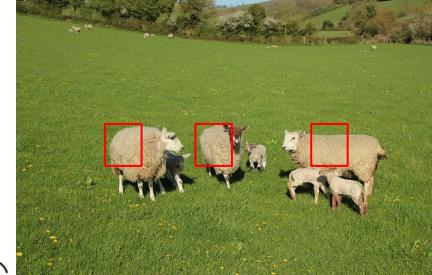


RNN: recurrent neural network.
Most often found in NLP setups.*



All this is not limited to these two applications, including convolutional and recurrent layers move a fixed size window over the input space with fixed stride.

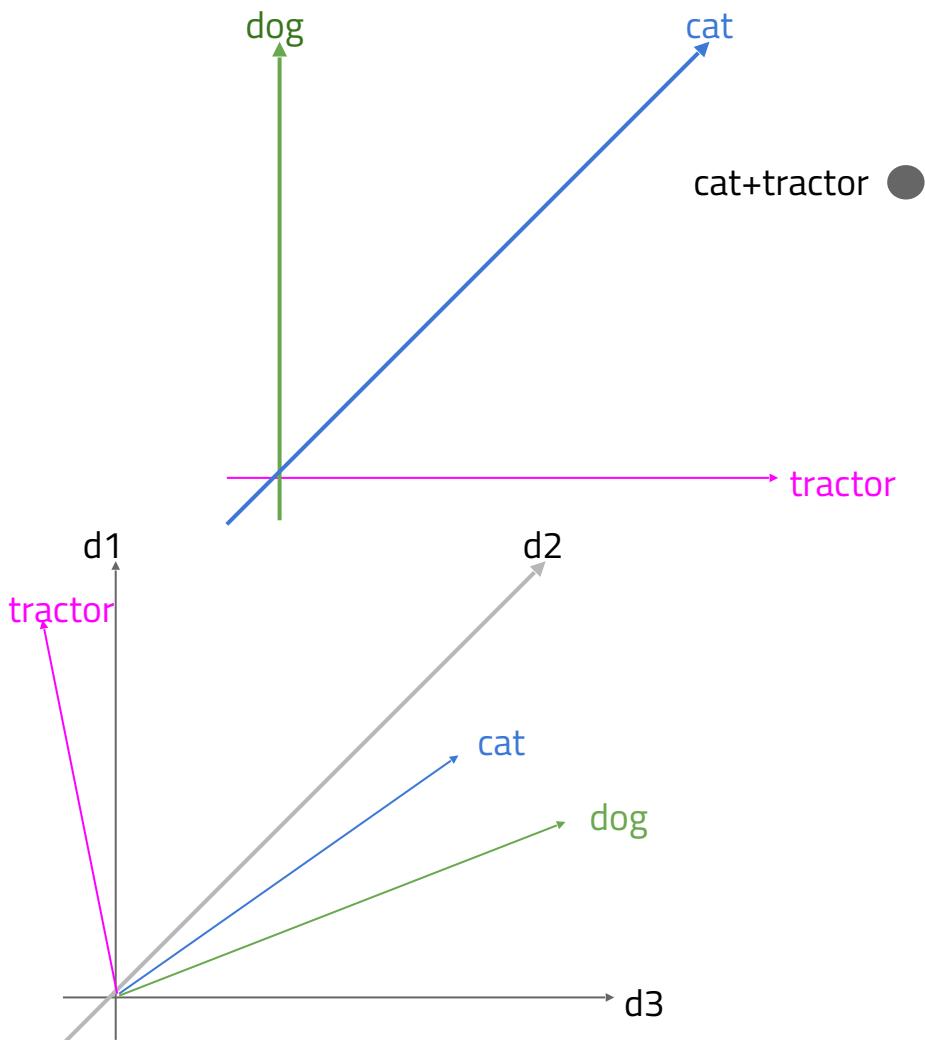
Each window is projected by a parameterized neural operation, often followed by an aggregation step such as (max) pooling.



Text representations

Neural IR

Text representations



In DL, commonly known as "one-hot" encoding/repr.

Local vector representation (sparse, high-dimensional)

$$\text{cat} = (0, 0, 0, 0, 1, 0, 0, 0, 0, 0, \dots, 0)$$

$$\text{tractor} = (0, 0, 0, 1, 0, 0, 0, 0, 0, 0, \dots, 0)$$

$$\text{dog} = (1, 0, 0, 0, 0, 0, 0, 0, 0, 0, \dots, 0)$$

$$\text{feline} = (0, 1, 0, 0, 0, 0, 0, 0, 0, 0, \dots, 0)$$

Semantic gap: the similarity between any of

Individual dimensions are (usually) not interpretable

Distributed vector representation

(dense, often low-dimensional)

$$\text{cat} = (0, 1, 1, 0, 1, 0, 0, 0, 0, 0)$$

$$\text{tractor} = (0, 0, 0, 1, 0, 0, 0, 1, 1, 0)$$

$$\text{dog} = (1, 1, 0, 0, 0, 0, 0, 0, 0, 0)$$

$$\text{feline} = (0, 1, 1, 1, 1, 0, 0, 0, 0, 1)$$

The similarity between some of those terms is no longer zero.

Text representations



Text representations can be learnt in a **supervised** or **unsupervised** fashion.

In IR, supervised approaches use query-document pairs.

In IR, the unsupervised approach uses just queries or just documents.



"Similarity" is not an absolute concept, it depends on the task & context at hand. In IR, it is related to relevance.

In DL, commonly known as "one-hot" encoding/repr.

Local vector representation (sparse, high-dimensional)
 $cat = (0, 0, 0, 0, 1, 0, 0, 0, 0, 0, \dots, 0)$

$tractor = (0, 0, 0, 1, 0, 0, 0, 0, 0, 0, \dots, 0)$

$dog = (1, 0, 0, 0, 0, 0, 0, 0, 0, 0, \dots, 0)$

$feline = (0, 0, 0, 0, 0, 0, 0, 0, 0, 0, \dots, 0)$

Semantic gap: the similarity between any of

Individual dimensions are (usually) not interpretable

Distributed vector representation (dense, often low-dimensional)

$cat = (0, 1, 1, 0, 1, 0, 0, 0, 0, 0)$

$tractor = (0, 0, 0, 1, 0, 0, 0, 1, 1, 0)$

$dog = (1, 1, 0, 0, 0, 0, 0, 0, 0, 0)$

$feline = (0, 1, 1, 1, 1, 0, 0, 0, 0, 1)$

The similarity between some of those terms is

Extreme case: very long document with a single relevant sentence. Still relevant ...



Which representation is best so that:
- "Amsterdam" and "London"
- "Ajax" and "Amsterdam"
are similar? (**Typical vs Topical!**)
in-document features

Feature-based representation examples

$$cat = (0, 0, \underset{doc \ 3}{\downarrow} 1, 0, \underset{doc \ 5}{\downarrow} 1, 0, 0, 0, 1, 0, \dots, \underset{doc \ 1678}{\downarrow} 1, 0)$$

neighbouring-word features

$$cat = (0, 0, \underset{runs}{\downarrow} 1, 0, \underset{eats}{\downarrow} 1, 0, 0, 0, \underset{hurt}{\downarrow} 1, 0, \dots, \underset{is}{\downarrow} 1, 0)$$

neighbouring-word with distances features

$$cat = (0, 0, \underset{runs^{+1}}{\downarrow} 1, 0, \underset{black^{-1}}{\downarrow} 1, 0, 0, 0, \underset{along^{+2}}{\downarrow} 1, 0, \dots, \underset{old^{-1}}{\downarrow} 1, 0)$$

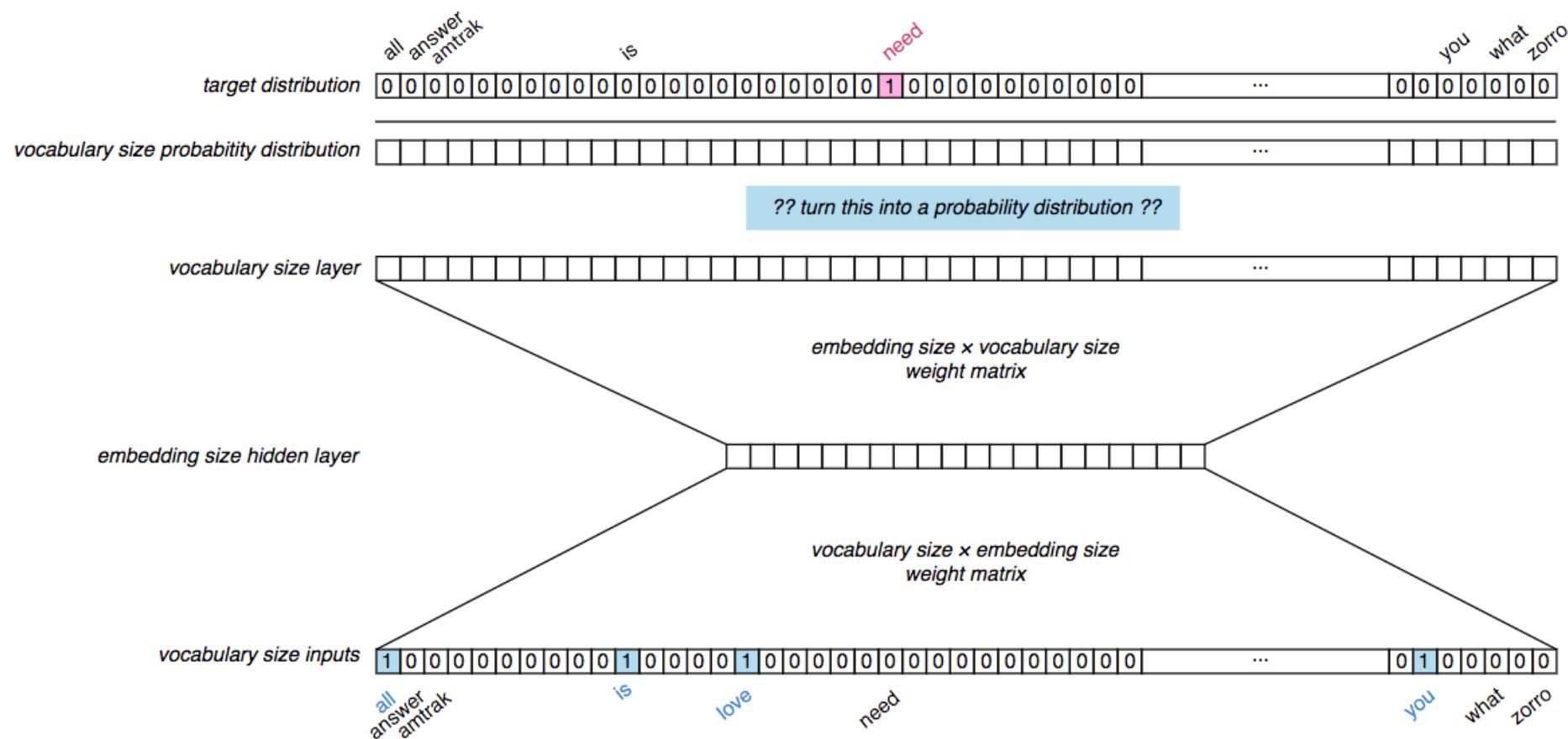
character trigram features

$$kitten = (0, 0, \underset{en\#}{\downarrow} 1, 0, \underset{\#ki}{\downarrow} 1, 0, \underset{tte}{\downarrow} 1, 0, \underset{kit}{\downarrow} 1, 0, \underset{ten}{\downarrow} 1, 0, \dots, \underset{itt}{\downarrow} 1, 0)$$

Embeddings

word2vec/GloVe Will be covered
in more details by Naval!

word2vec/GloVe: the vector of a word should be similar to the vectors of its neighbouring words



... representation in a new space that should preserve the relationships between items of the original representation.

Embeddings

Dense vector representation.

Low-dimensional.

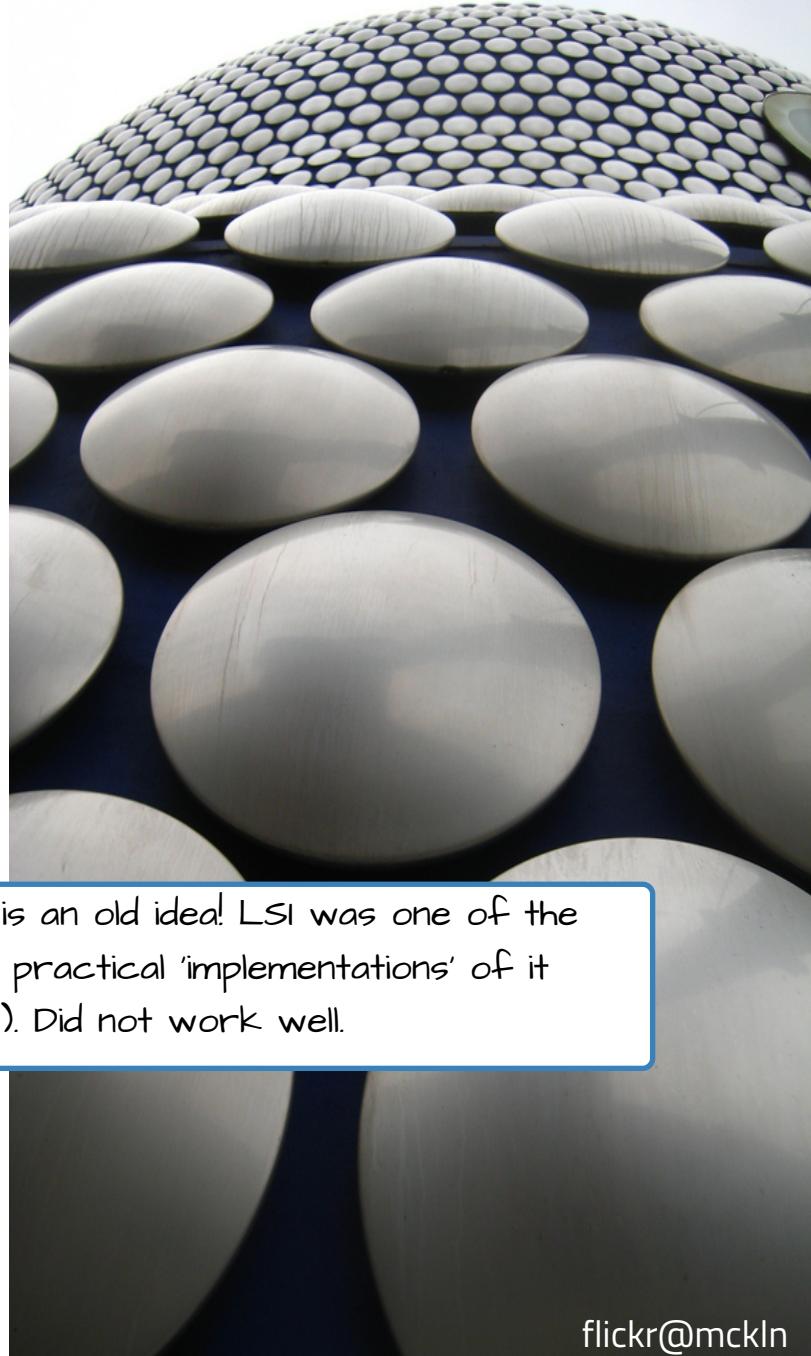
Learnt from data.

Distributed representation most often refers to learnt embeddings.

In today's NLP literature, the default encoding.

Motivated by the **distributional hypothesis** (Harris, 1954): "*a word is characterized by the company it keeps*" (Firth, 1957)

This is an old idea! LSI was one of the first practical 'implementations' of it (1988). Did not work well.



Local and distributed representations are needed

The President of the United States of America (POTUS) is the elected head of state and head of government of the United States. The president leads the executive branch of the federal government and is the commander in chief of the United States Armed Forces. Barack Hussein Obama II (born August 4, 1961) is an American politician who is the 44th and current President of the United States. He is the first African American to hold the office and the first president born outside the continental United States.

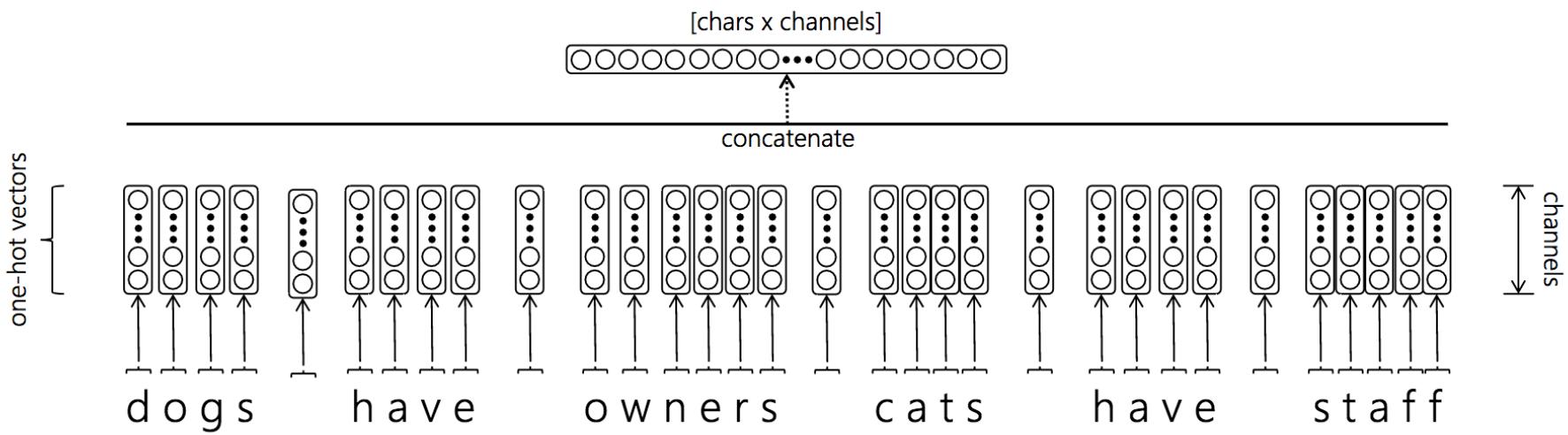
Q="united states president"

Shade of green indicates the drop in retrieval model's document score by individually removing each of the passage terms.

The President of the United States of America (POTUS) is the elected head of state and head of government of the United States. The president leads the executive branch of the federal government and is the commander in chief of the United States Armed Forces. Barack Hussein Obama II (born August 4, 1961) is an American politician who is the 44th and current President of the United States. He is the first African American to hold the office and the first president born outside the continental United States.

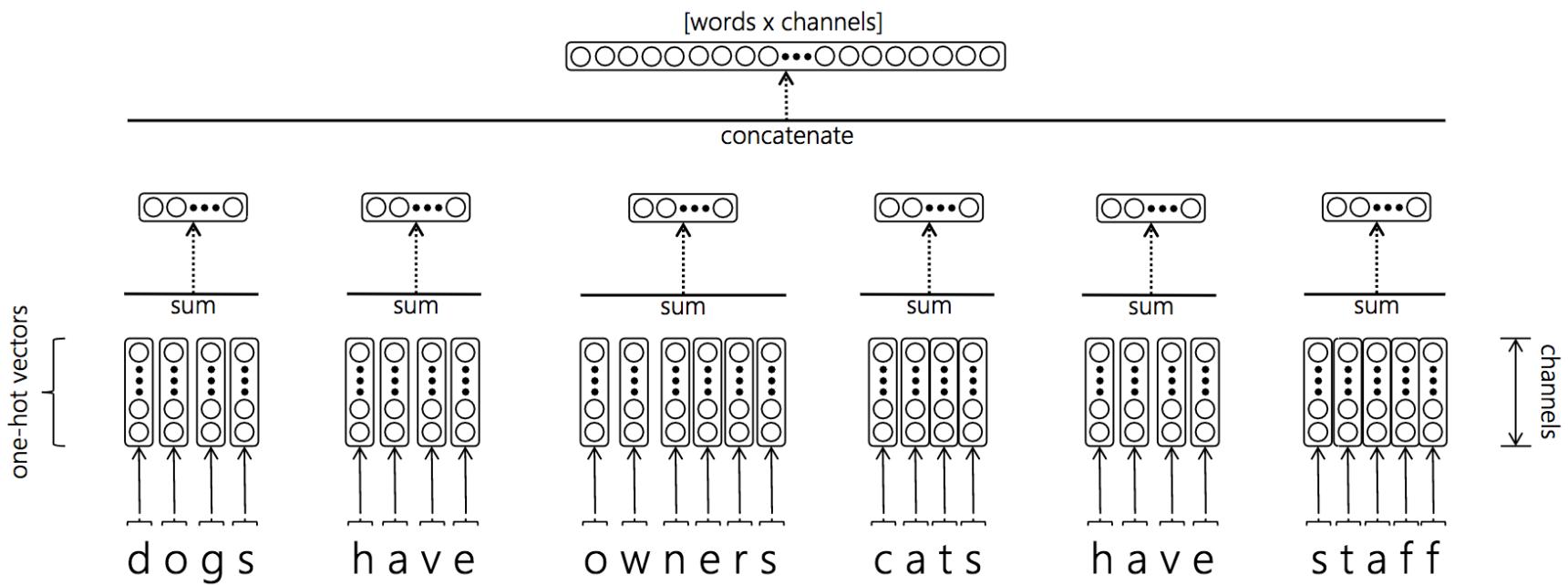
Text input to deep neural nets

character-level input



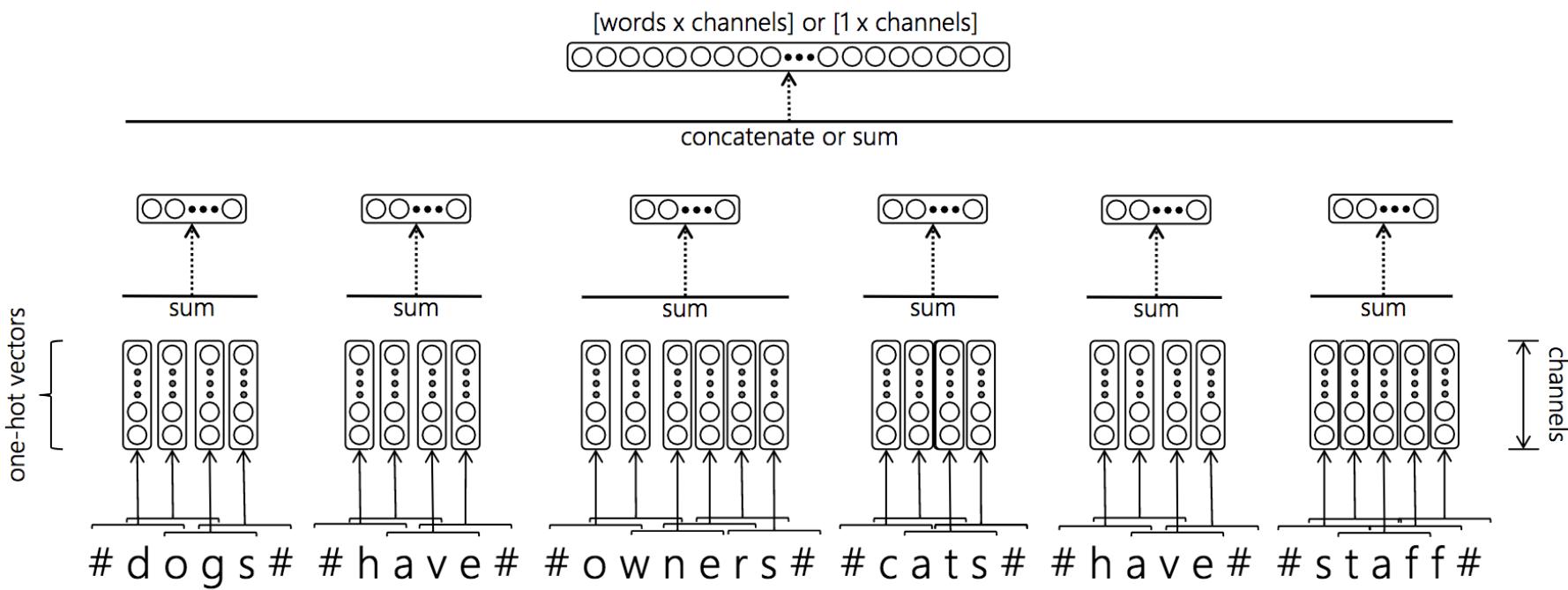
Text input to deep neural nets

term-level input with
bag-of-chars per term



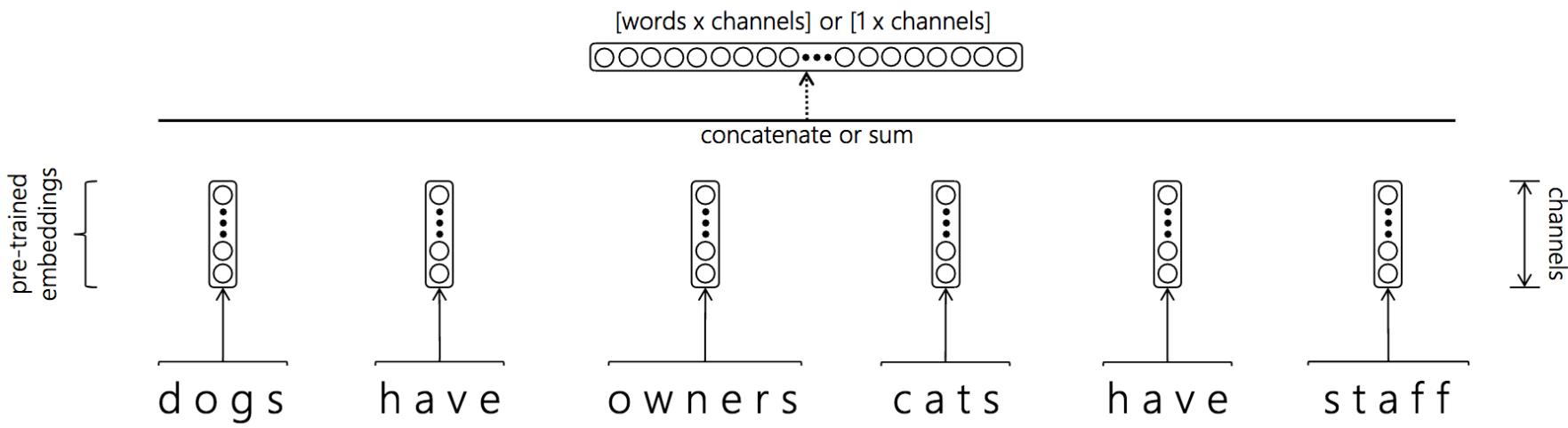
Text input to deep neural nets

term-level input with bag-of-trigrams per term



Text input to deep neural nets

term-level input with pre-trained word embeddings



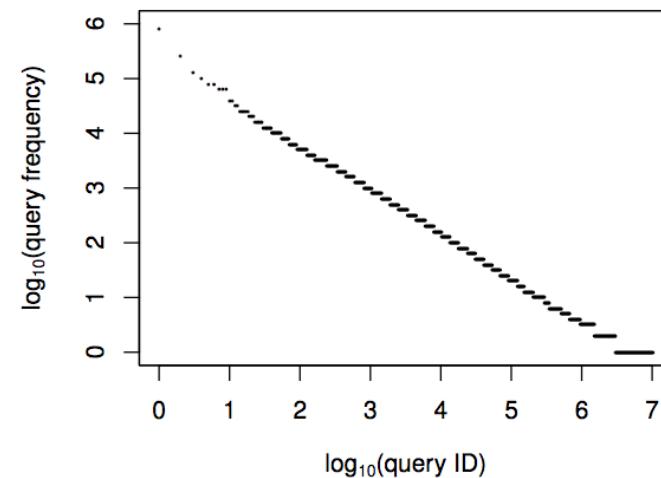
Glove (Nava will explain how this works exactly)

- Pre-trained word vectors. This data is made available under the [Public Domain Dedication and License v1.0](#) whose full text can be found at: <http://www.opendatacommons.org/licenses/pddl/1.0/>.
 - [Wikipedia 2014 + Gigaword 5](#) (6B tokens, 400K vocab, uncased, 50d, 100d, 200d, & 300d vectors, 822 MB download): [glove.6B.zip](#)
 - Common Crawl (42B tokens, 1.9M vocab, uncased, 300d vectors, 1.75 GB download): [glove.42B.300d.zip](#)
 - Common Crawl (840B tokens, 2.2M vocab, cased, 300d vectors, 2.03 GB download): [glove.840B.300d.zip](#)
 - Twitter (2B tweets, 27B tokens, 1.2M vocab, uncased, 25d, 50d, 100d, & 200d vectors, 1.42 GB download): [glove.twitter.27B.zip](#)

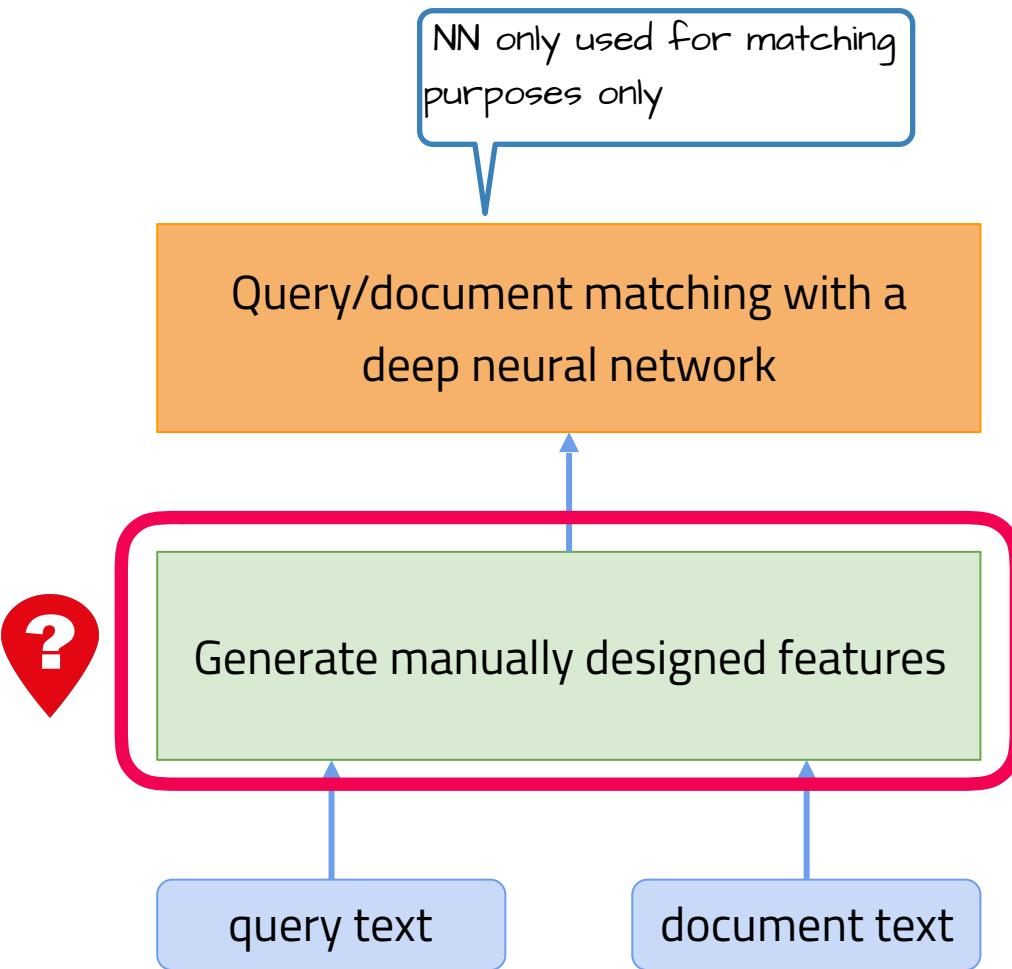
Neural IR architectures

A good retrieval system should ...

1. Have **semantic understanding** and enable **exact term matching** (*cat* and *feline* are similar query terms, but *hot dog* is not similar to *warm puppy; [rare name] gardening*)
2. Be **robust** to rare **inputs**: remember the long tail
3. Be **robust** to **corpus variance**:
BM25 out-of-the-box is a good ranking model for all domains, neural models require a lot of training (training on domain X data and testing on domain Y data can go wrong)
4. Be **robust** to **variable length input**
(classic IR has document length normalization)
5. Be **sensitive** to **context** (e.g. location)



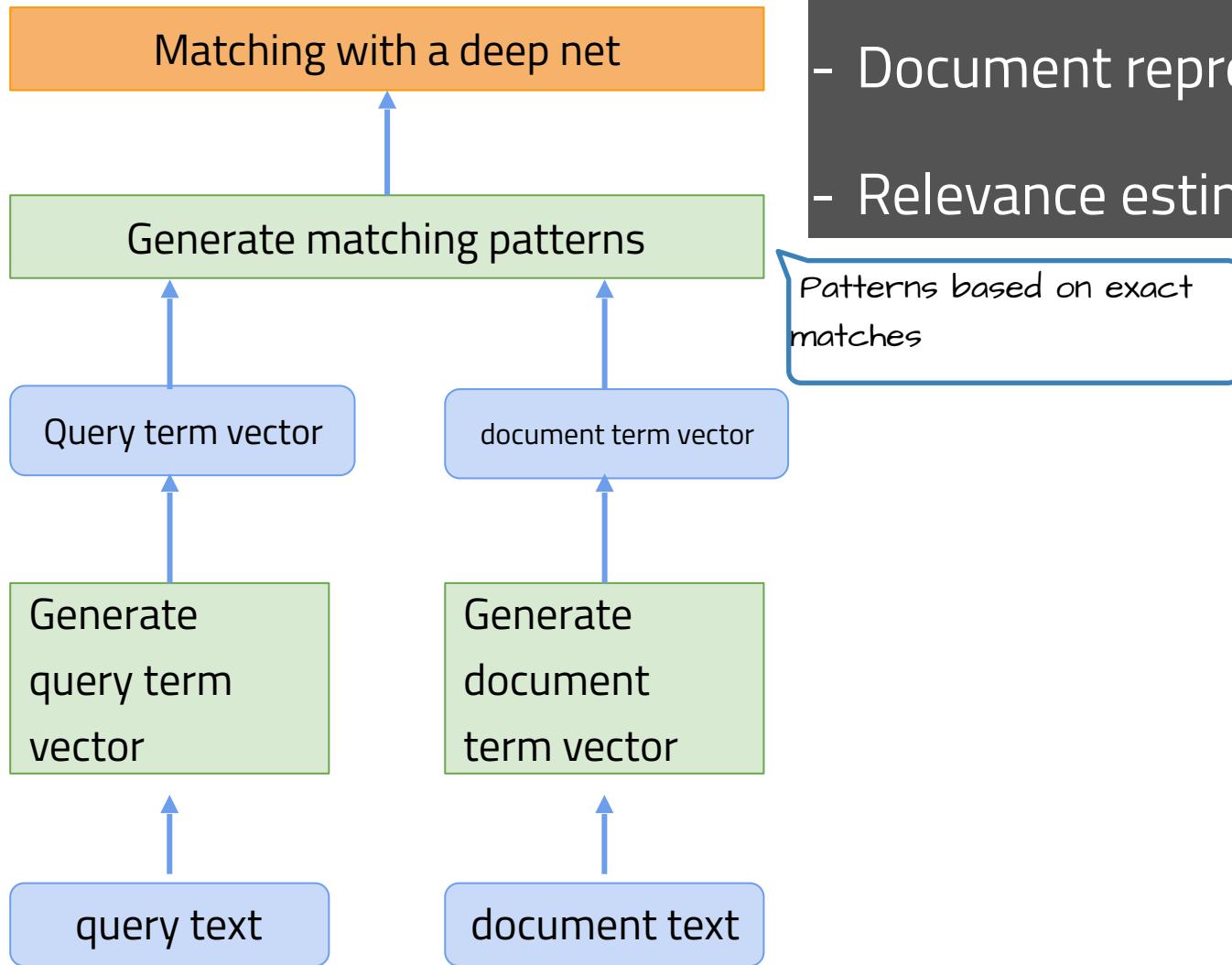
Learning to rank



Categorized according to their:

- Query representation
- Document representation
- Relevance estimation

Interaction focused models

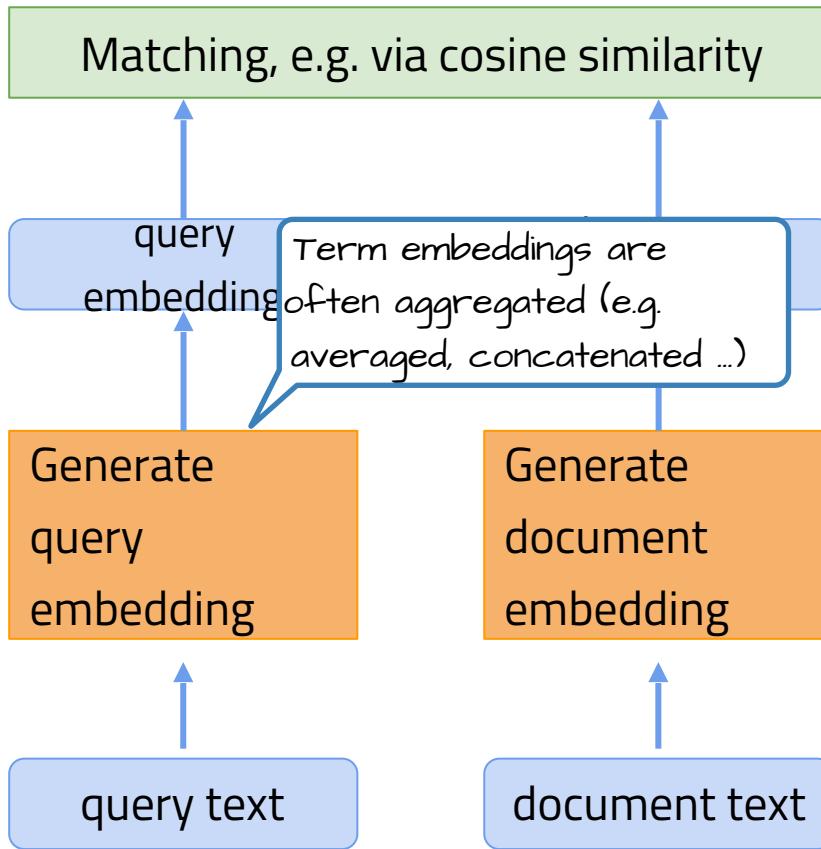


Categorized according to their:

- Query representation
- Document representation
- Relevance estimation

Patterns based on exact matches

Representation focused models

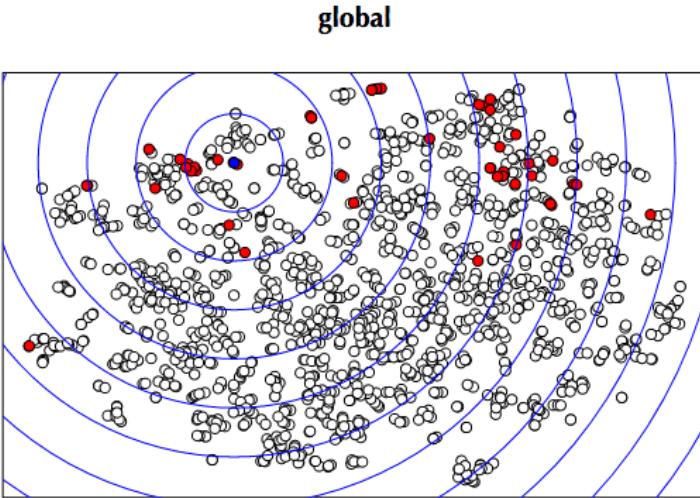


Categorized according to their:

- Query representation
- Document representation
- Relevance estimation

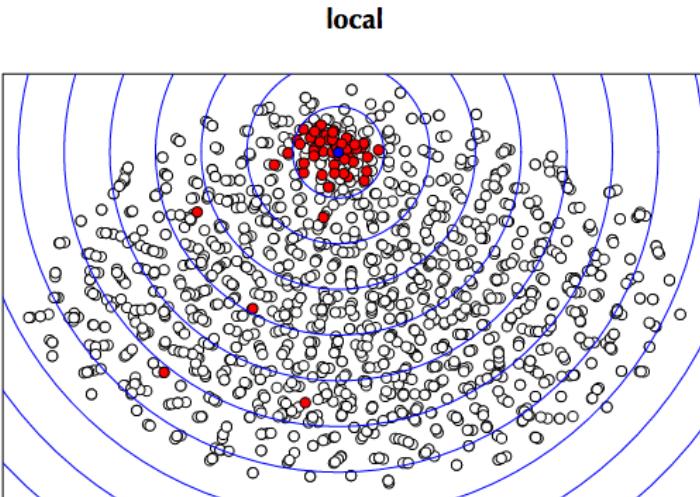
$$sim(q, d) = \cos(\vec{v}_q, \vec{v}_d) = \frac{\vec{v}_q^\top \vec{v}_d}{\|\vec{v}_q\| \|\vec{v}_d\|}$$

Query expansion



Point: term

Blue: query term



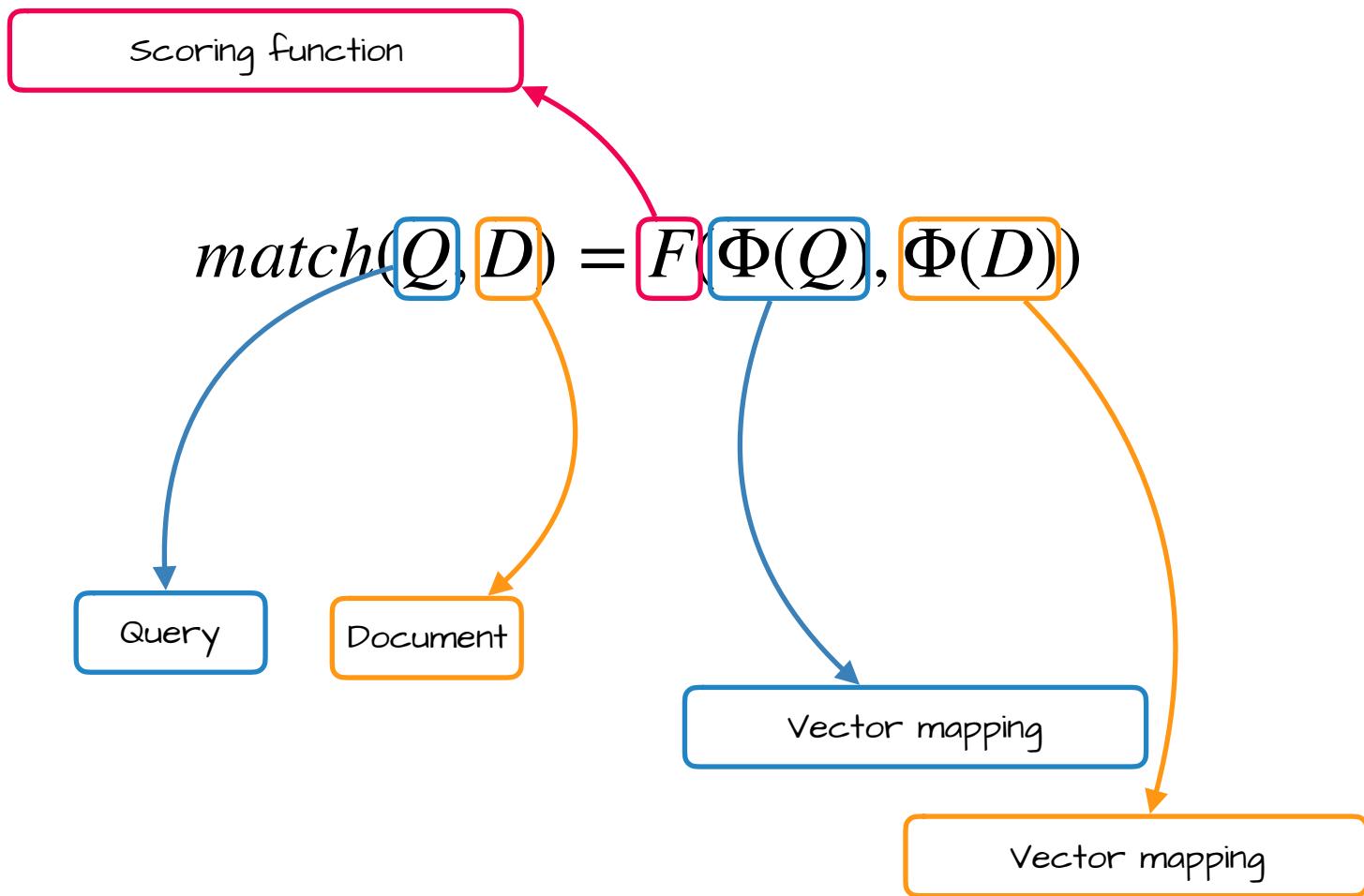
Red: term
occurring in
the relevant
set of
documents

- Categorized according to their:
- Query representation
 - Document representation
 - Relevance estimation

t-SNE projection
(stochastic neighbour
embedding)

A simple idea that works well
for PRF: learn embeddings on
the top retrieved documents
only instead of a large corpus.

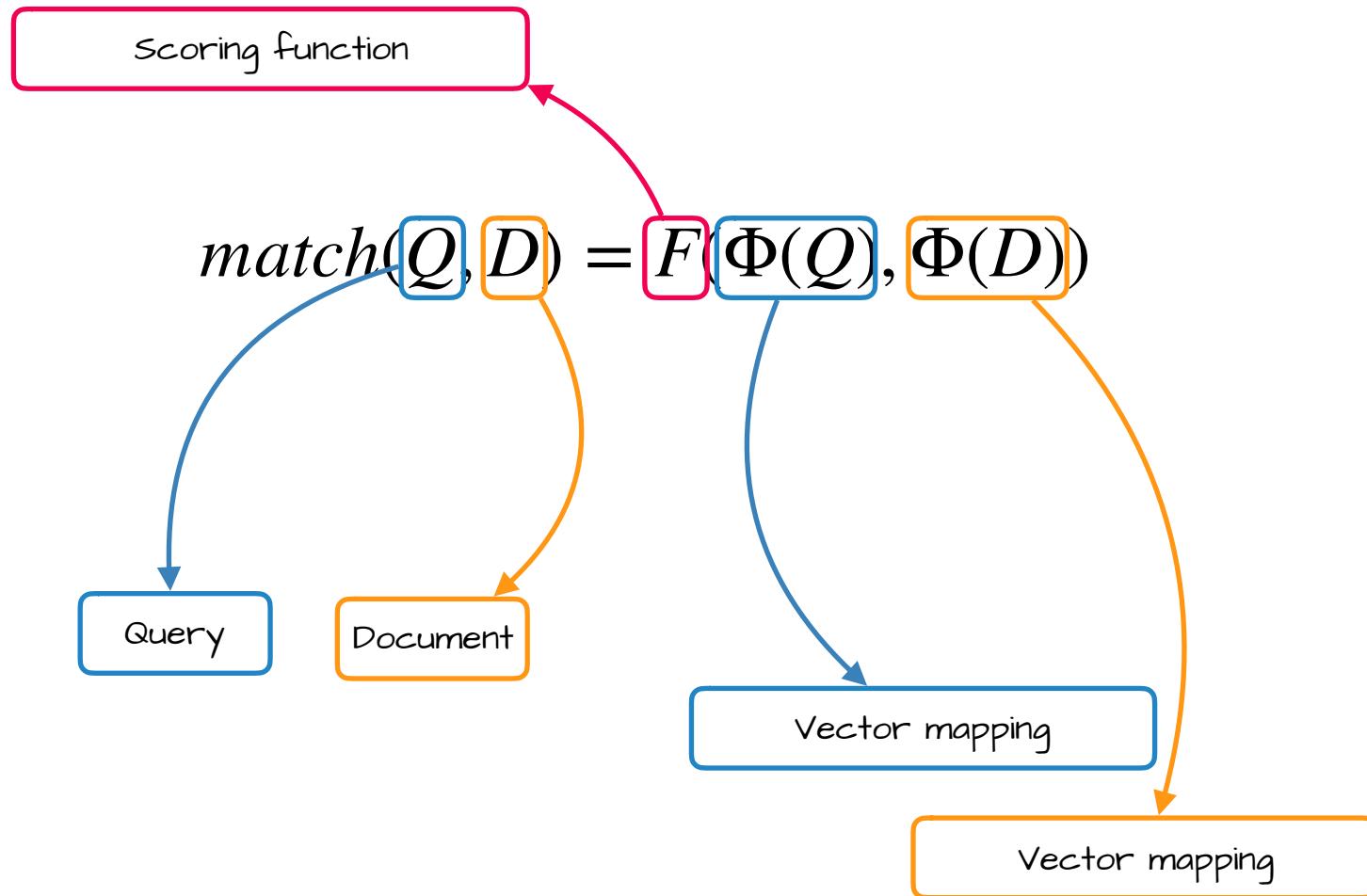
Generally speaking...



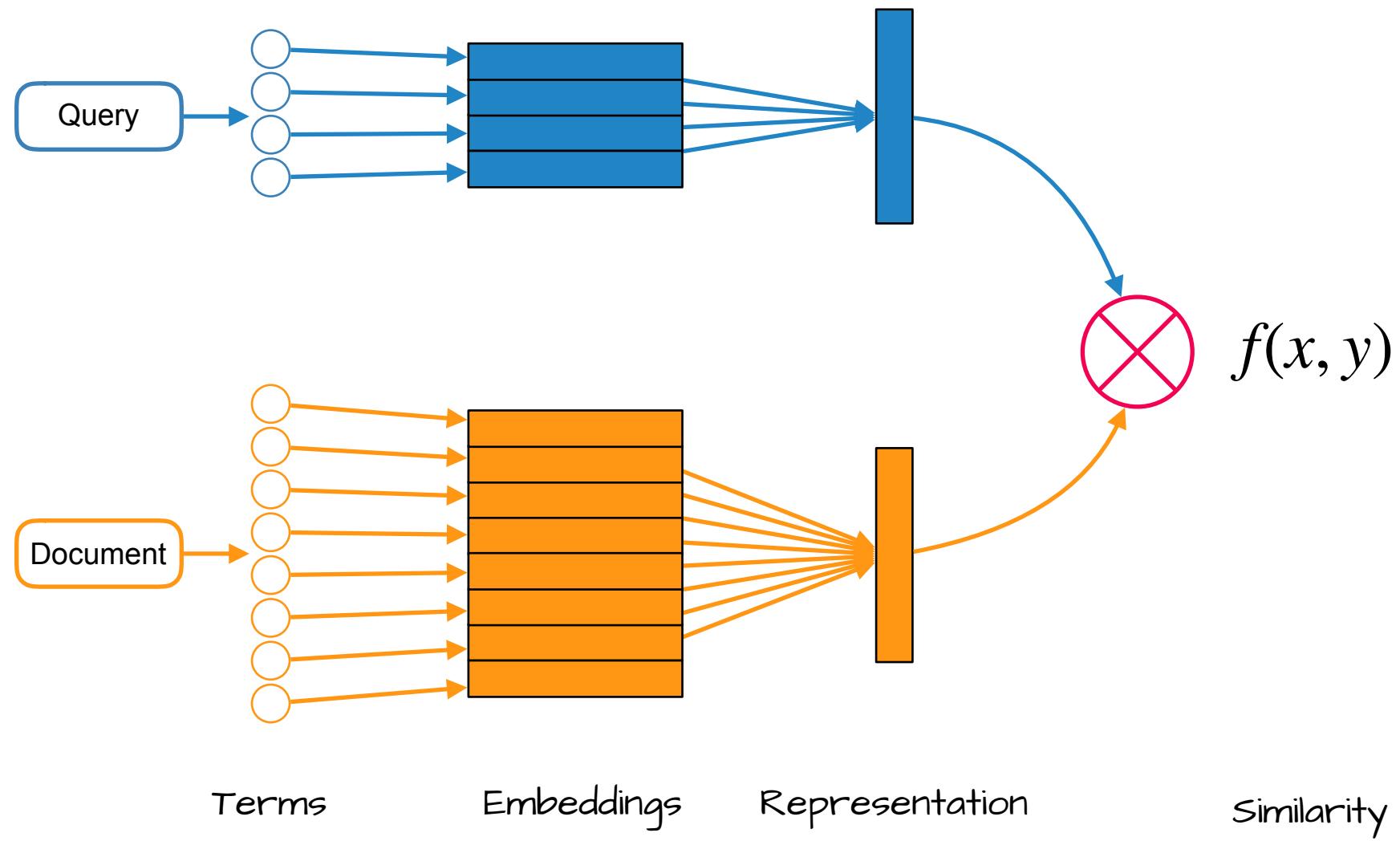
Representation focused models

(Or Siamese Networks)

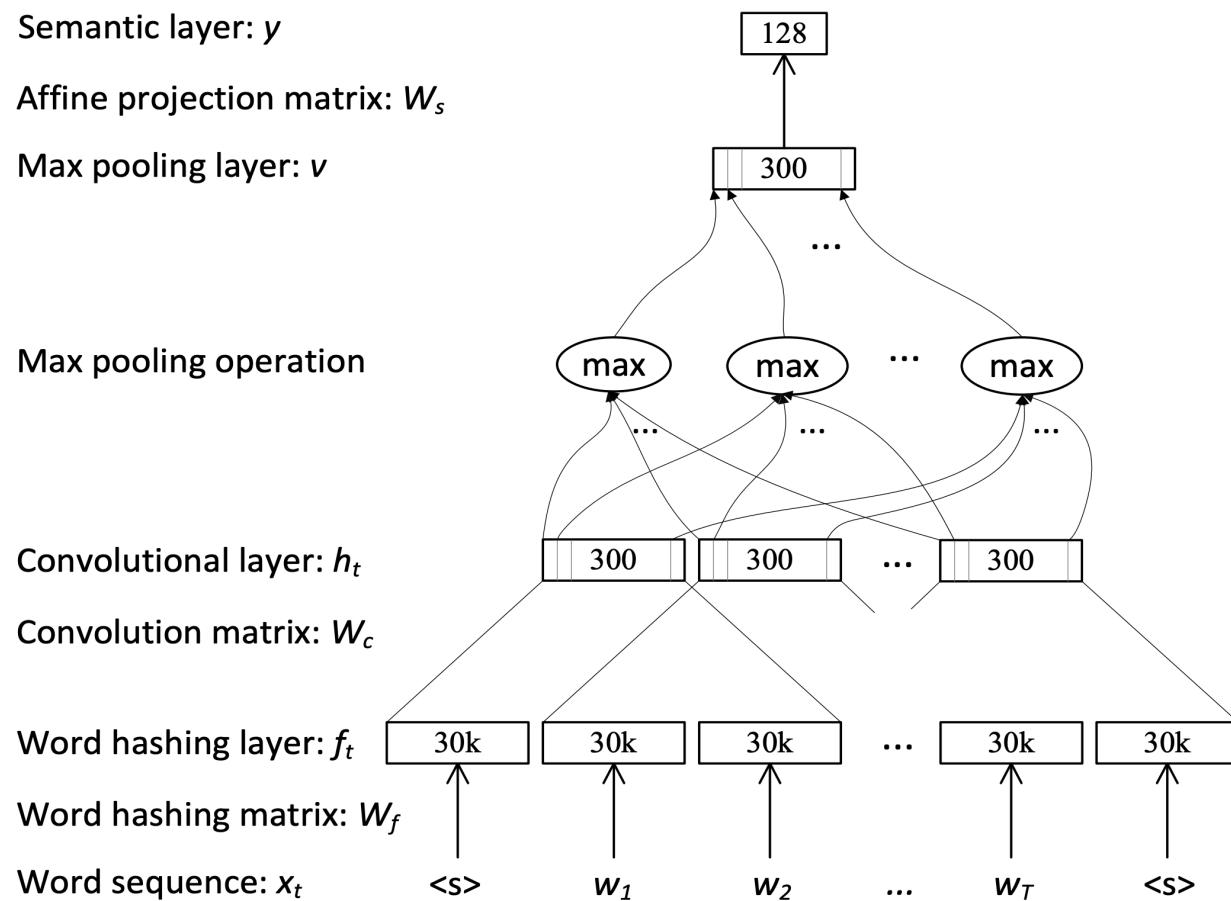
Representation focused models



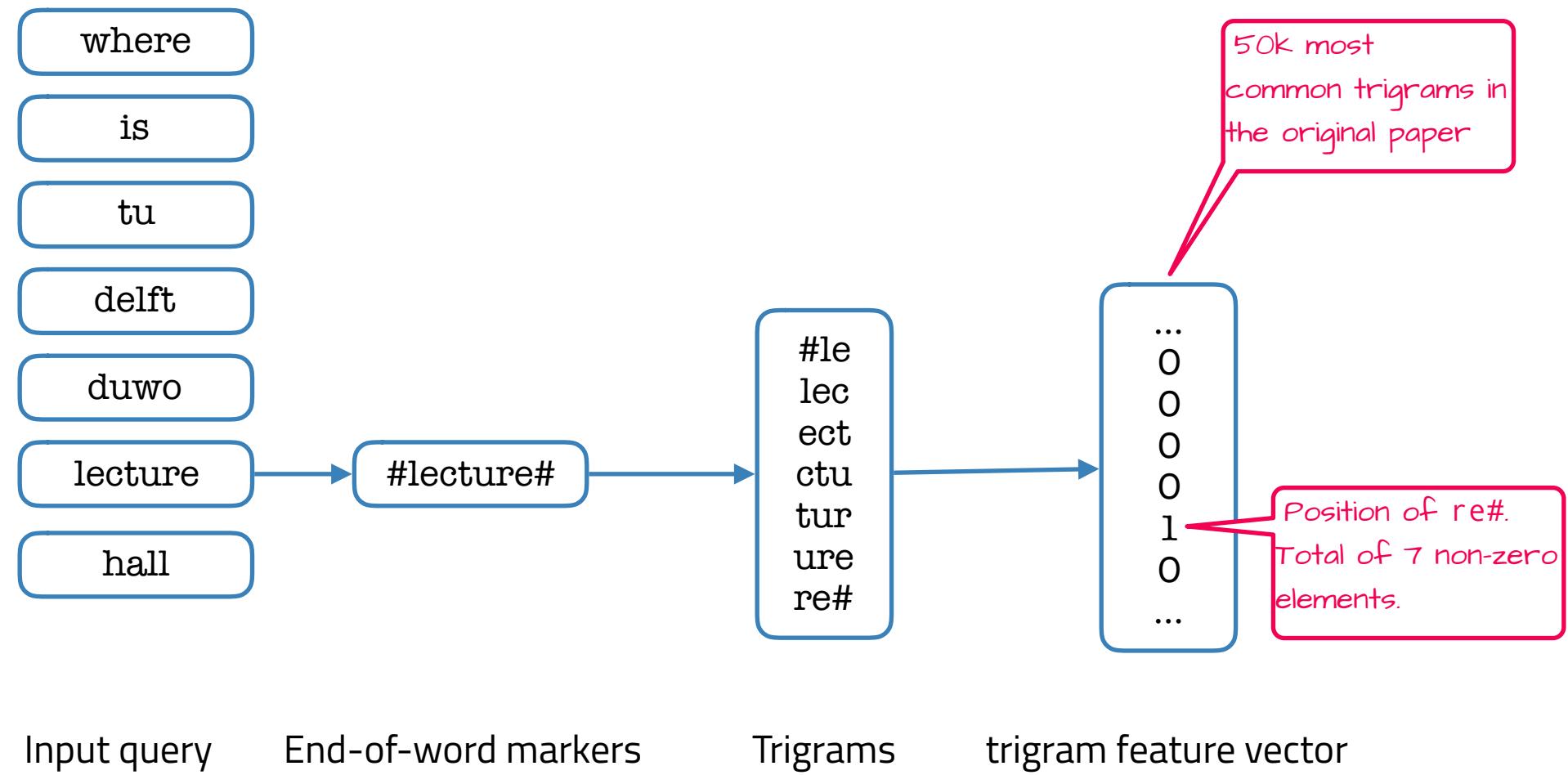
Representation focused models



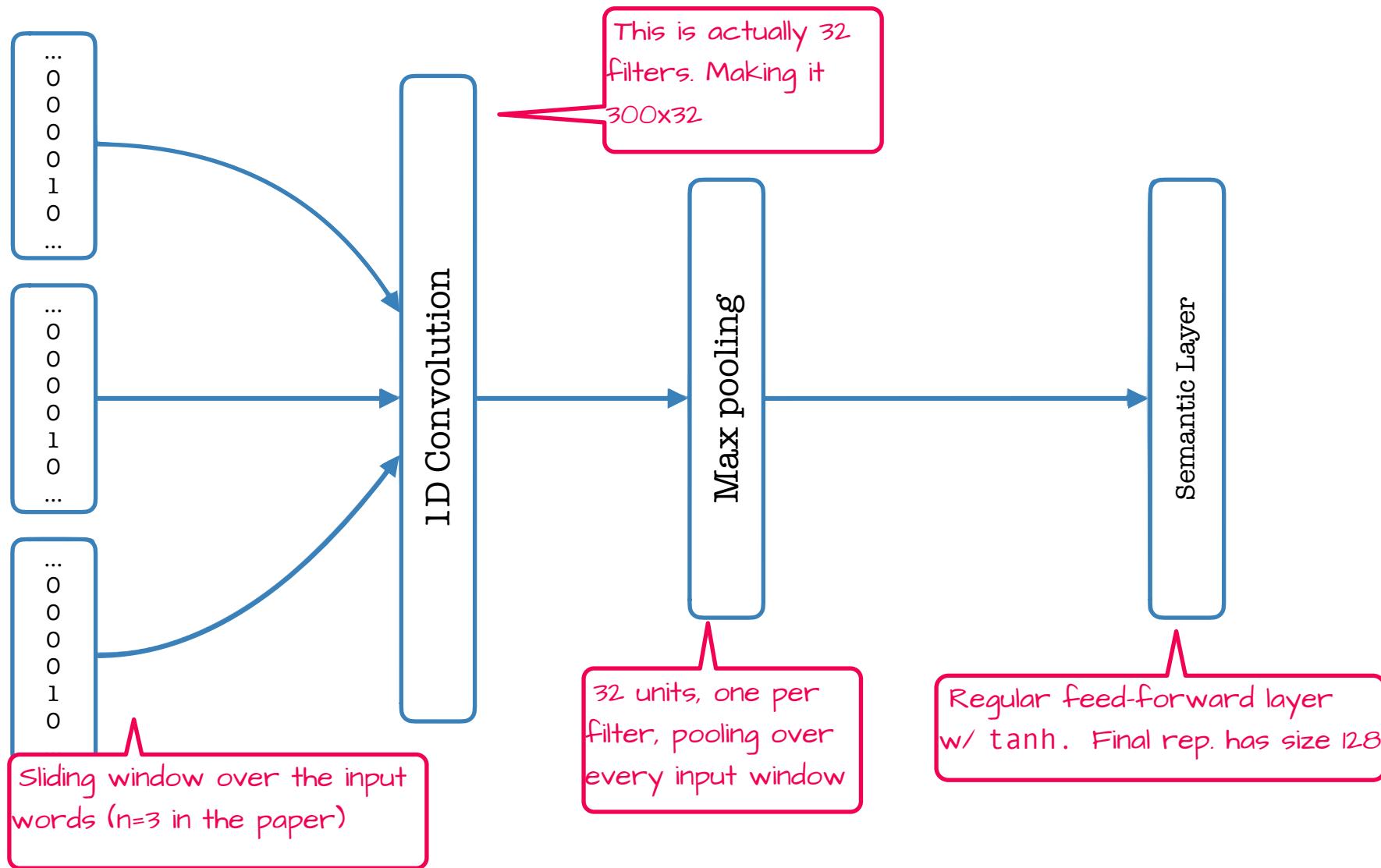
CDSSM - Convolutional Deep Structure Semantic Models



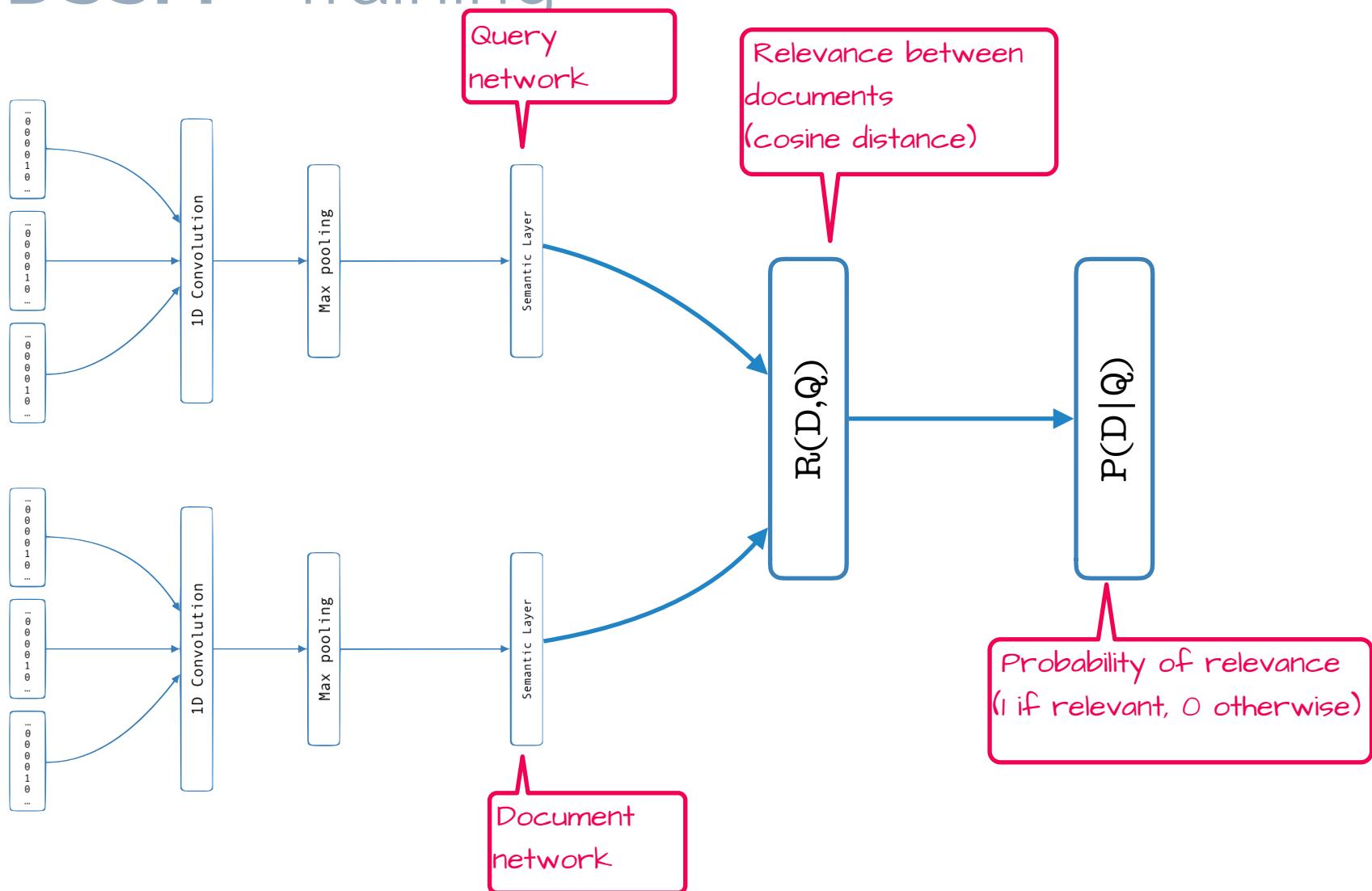
CDSSM - Word Hashing Layer



CDSSM - Convolutional Layer



CDSSM - Training



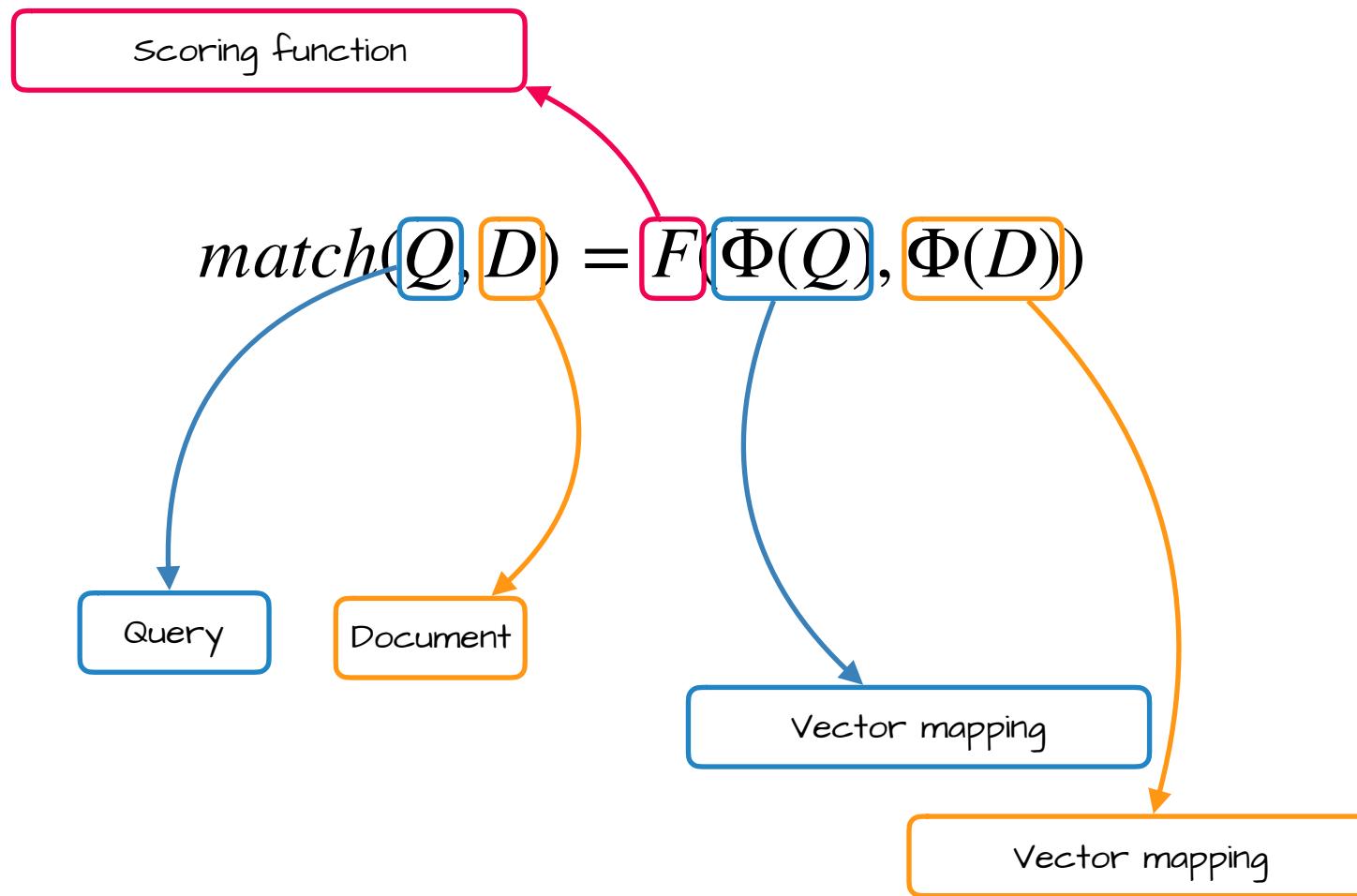
CDSSM - Results

Using bing queries

#	Models	NDCG@1	NDCG@3	NDCG@5
1	BM25	0.305	0.328	0.388
2	ULM	0.304	0.327	0.385
3	WTM	0.315	0.342	0.411
4	PTM (len <=3)	0.319	0.347	0.413
5	DSSM	0.320	0.355	0.431
6	C-DSSM (win=3)	0.342	0.374	0.447

Interaction focused models

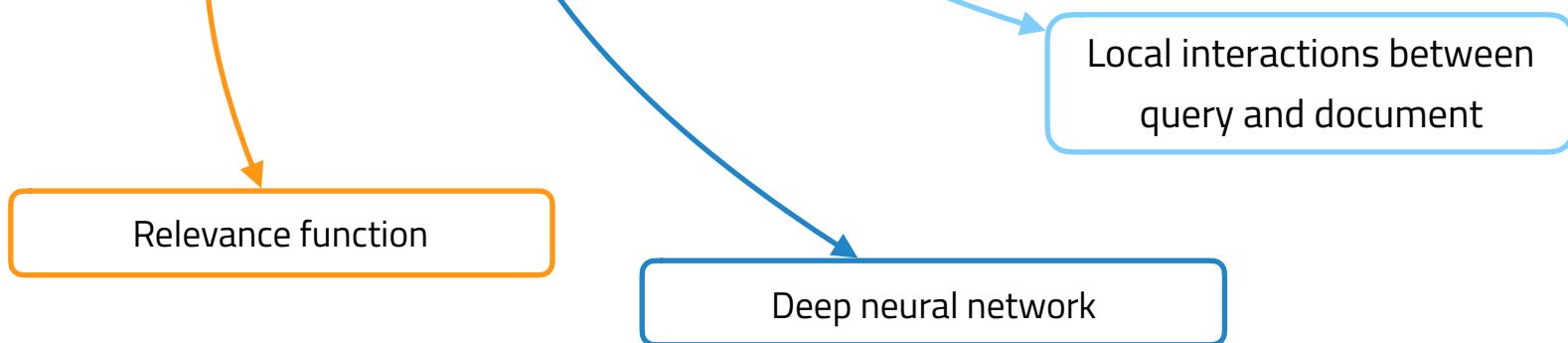
Interaction-focused models



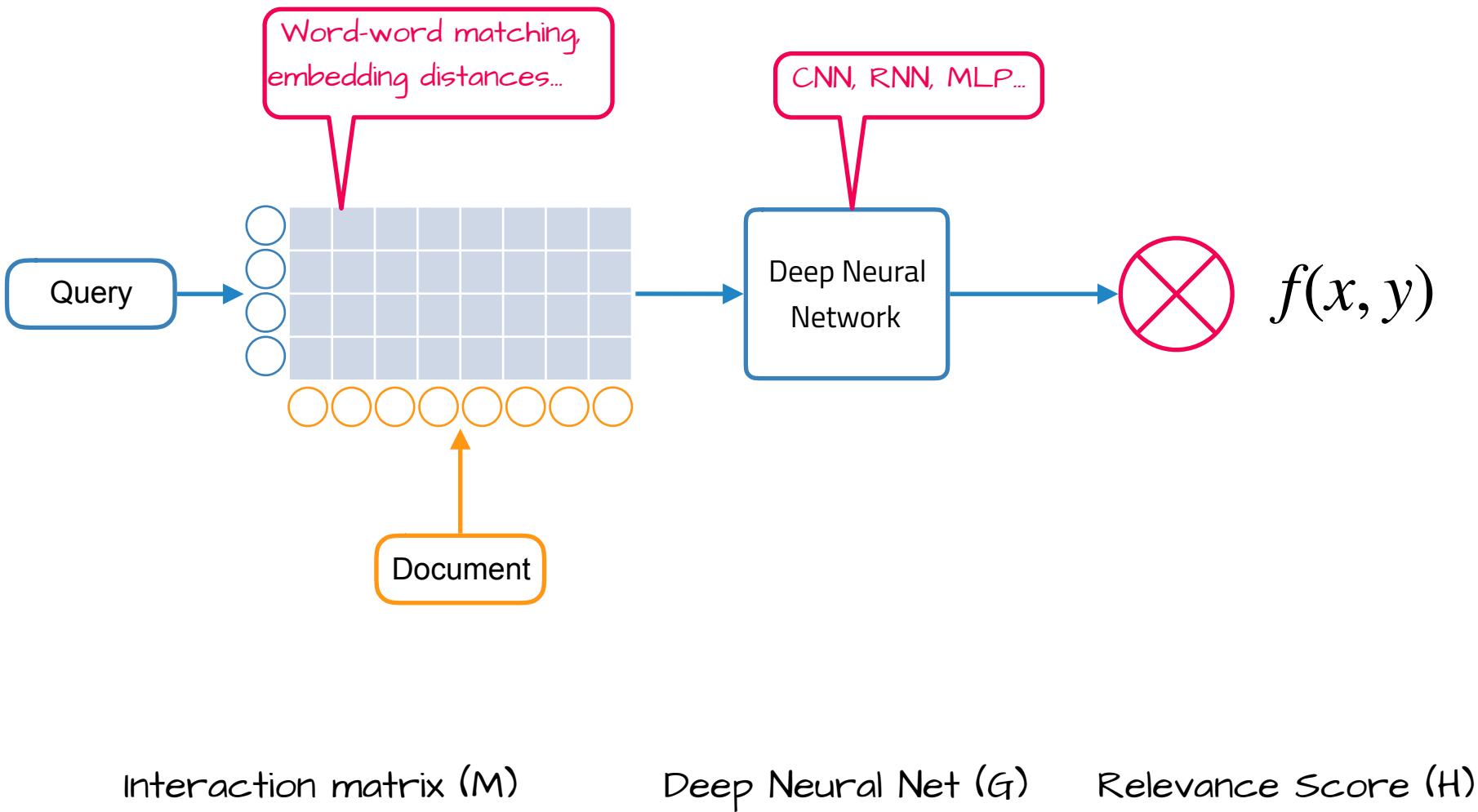
Interaction-focused models

$$match(Q, D) = F(\Phi(Q), \Phi(D))$$

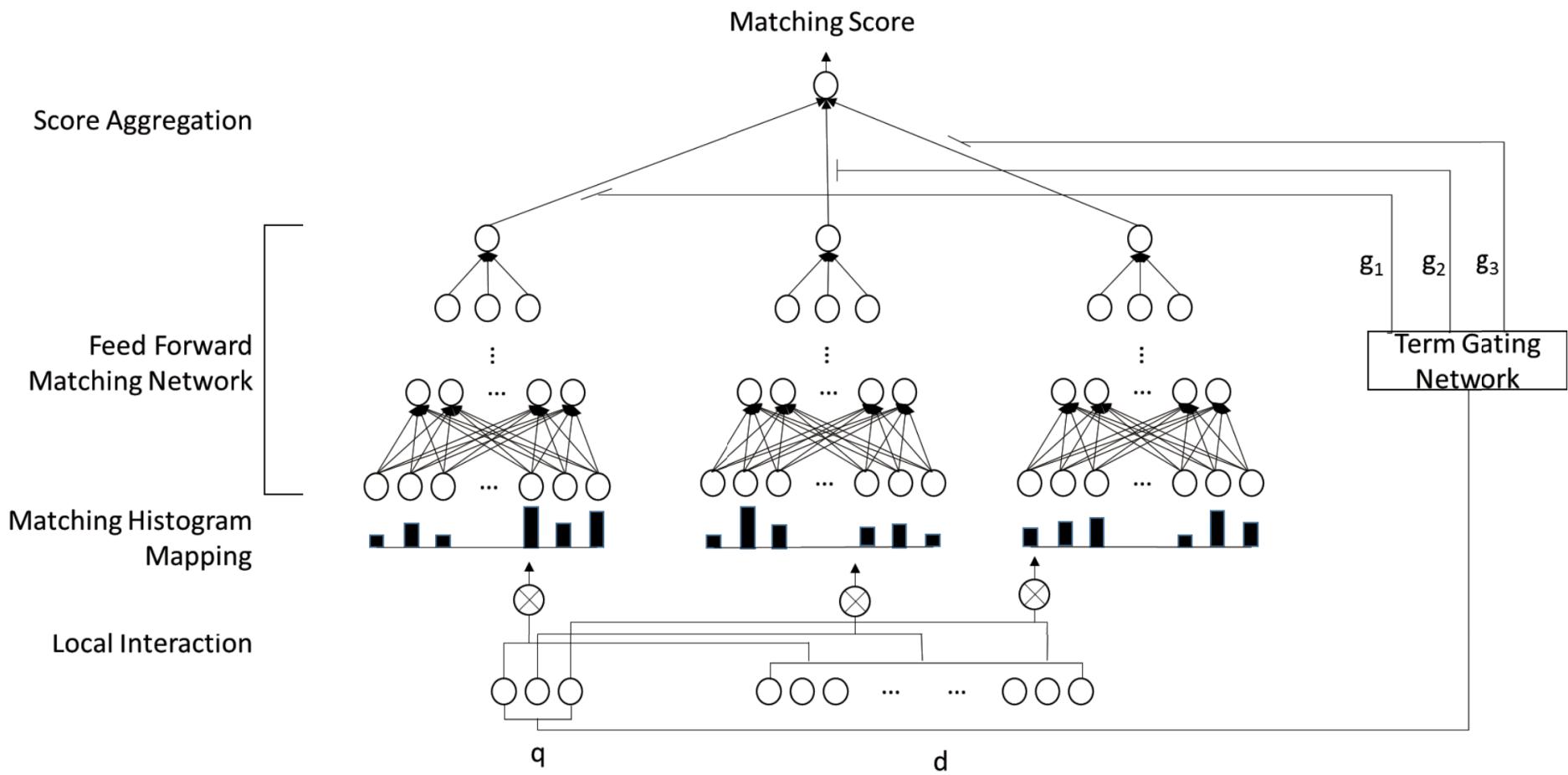
$$H \otimes G \otimes M(\Phi(Q)\Phi(D))$$



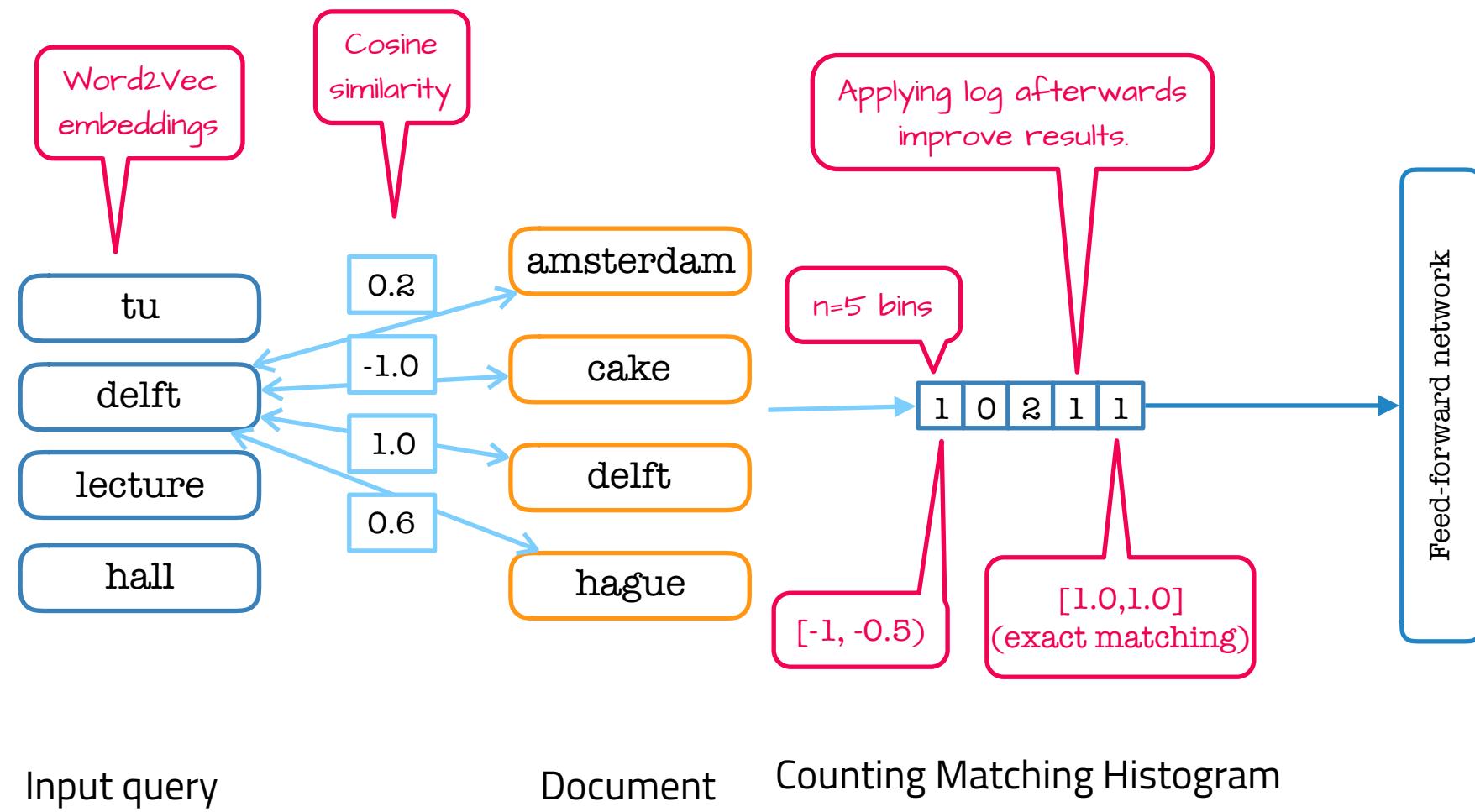
Interaction-focused models



DRMM - Deep Relevance Matching Model



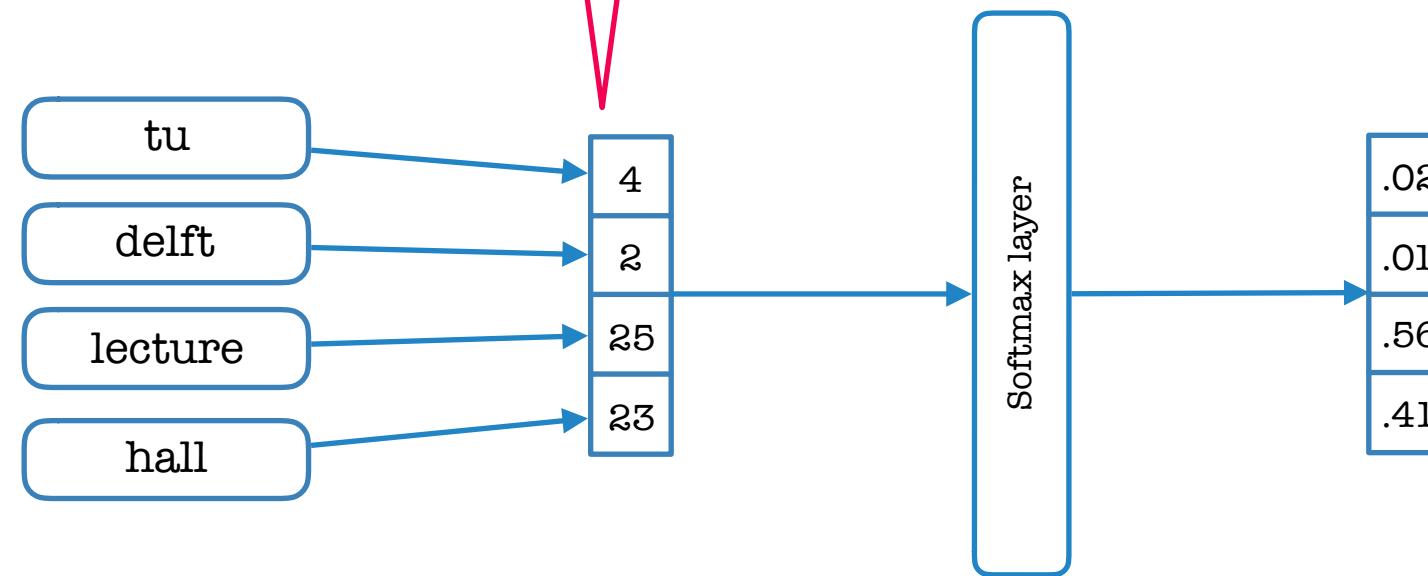
DRMM - Local Interactions and Histogram Mapping



DRMM - Term Gating Network

$$g_i = \frac{\exp(x_i^{(q)})}{\sum_{j=1}^M \exp(x_j^{(q)})}$$

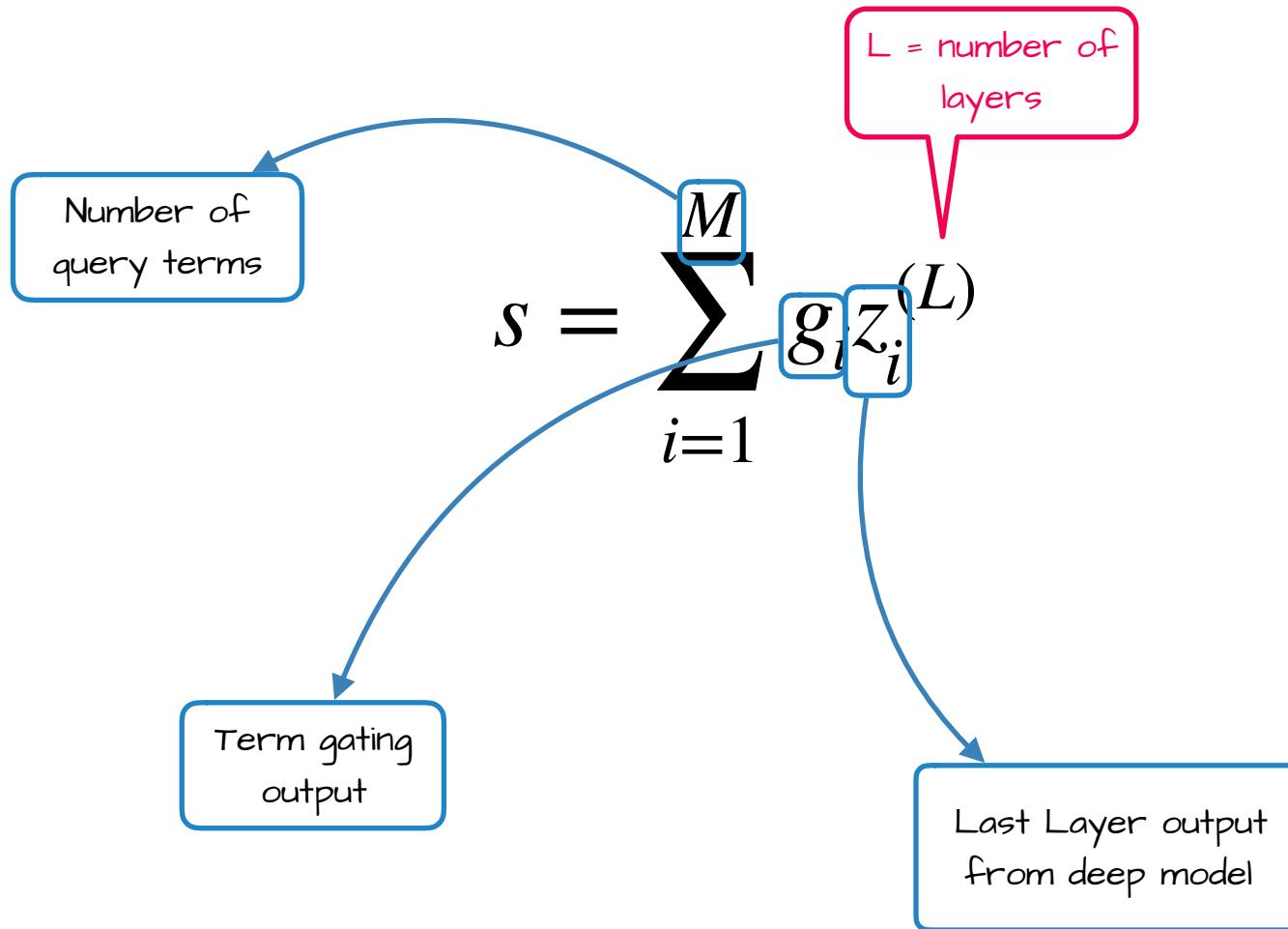
Calculated on
corpus



Input query

IDF Vector

DRMM - Final aggregation



DRMM - Results

Robust04, titles only

#	Models	MAP	NDCG@20	P@20
1	QL	0.253	0.415	0.369
2	BM25	0.255	0.418	0.370
3	CDSSM	0.067	0.146	0.125
4	ARC-II	0.067	0.147	0.128
5	DRRM	0.279	0.431	0.382

Duet: local and distributed representations



We need both exact term matches and semantic matches

Local representation

The President of the United States of America (POTUS) is the elected head of state and head of government of the United States. The president leads the executive branch of the federal government and is the commander in chief of the United States Armed Forces. Barack Hussein Obama II (born August 4, 1961) is an American politician who is the 44th and current President of the United States. He is the first African American to hold the office and the first president born outside the continental United States.

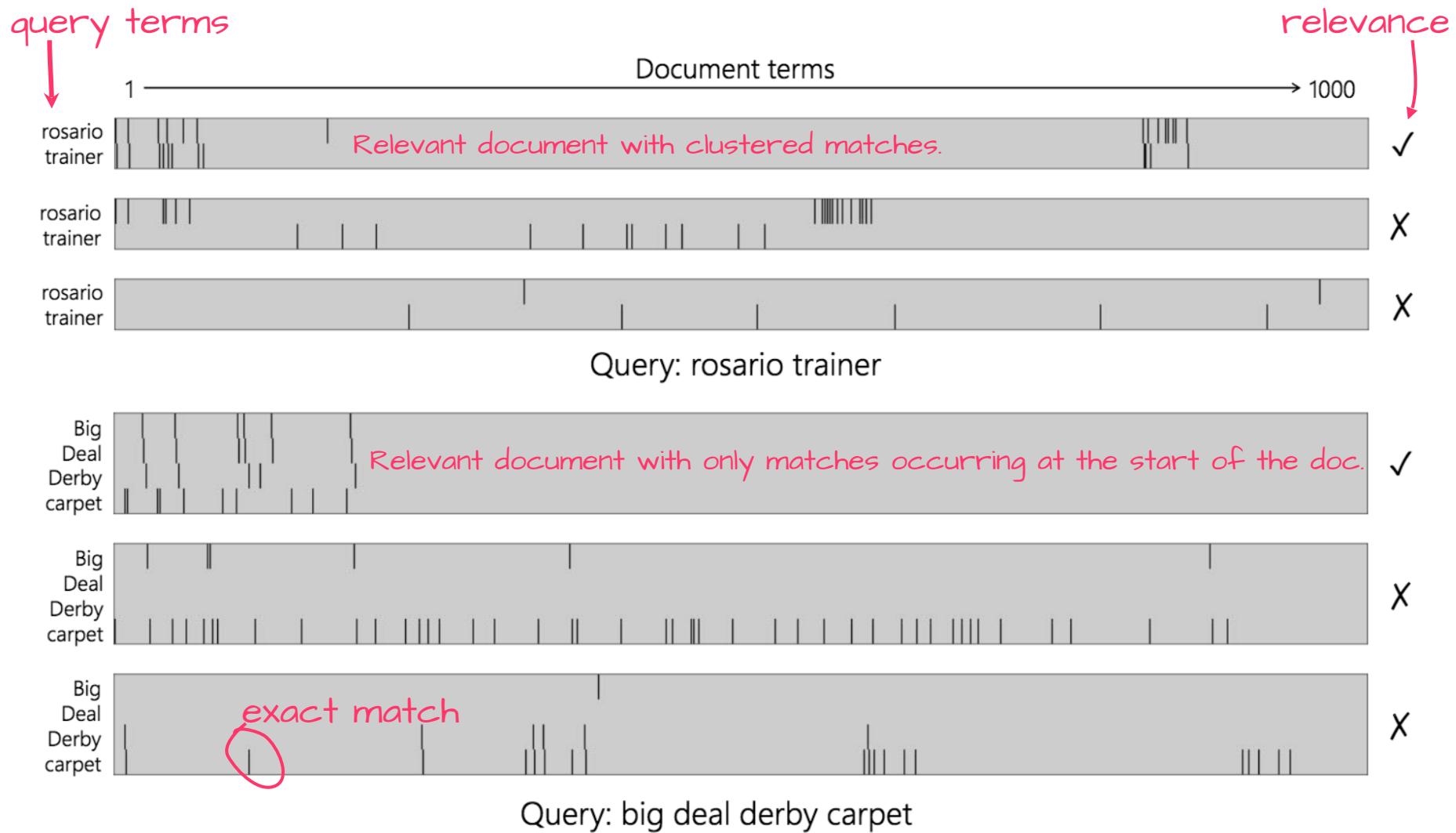
Q="united states president"

Shade of green indicates the drop in retrieval model's document score by individually removing each of the passage terms.

Distributed representation

The President of the United States of America (POTUS) is the elected head of state and head of government of the United States. The president leads the executive branch of the federal government and is the commander in chief of the United States Armed Forces. Barack Hussein Obama II (born August 4, 1961) is an American politician who is the 44th and current President of the United States. He is the first African American to hold the office and the first president born outside the continental United States.

Patterns of query term matches



joint training

Architecture

$$f(\mathbf{Q}, \mathbf{D}) = f_{local}(\mathbf{Q}, \mathbf{D}) + f_{distr}(\mathbf{Q}, \mathbf{D})$$

where $\mathbf{Q} = [q_1, \dots, q_{n_q}]$

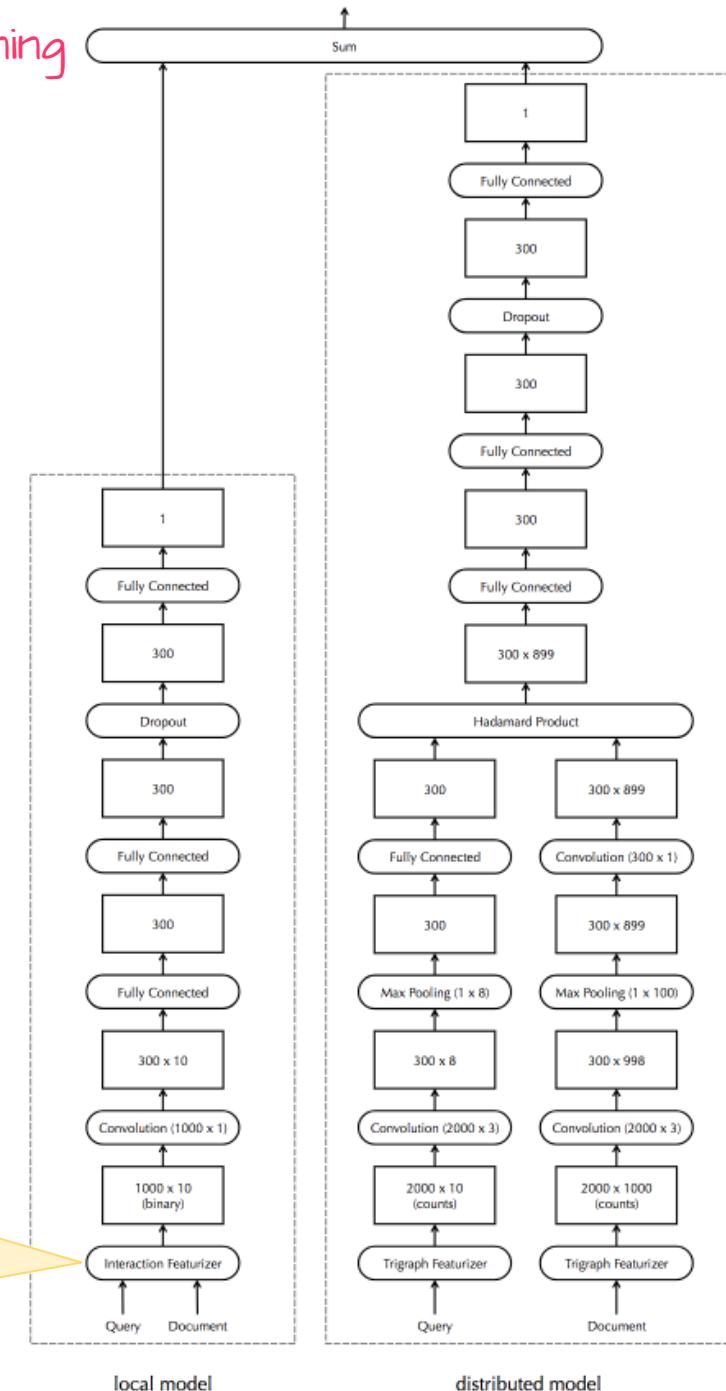
and $\mathbf{D} = [d_1, \dots, d_{n_d}]$;

a term is represented by a $m \times 1$ vector

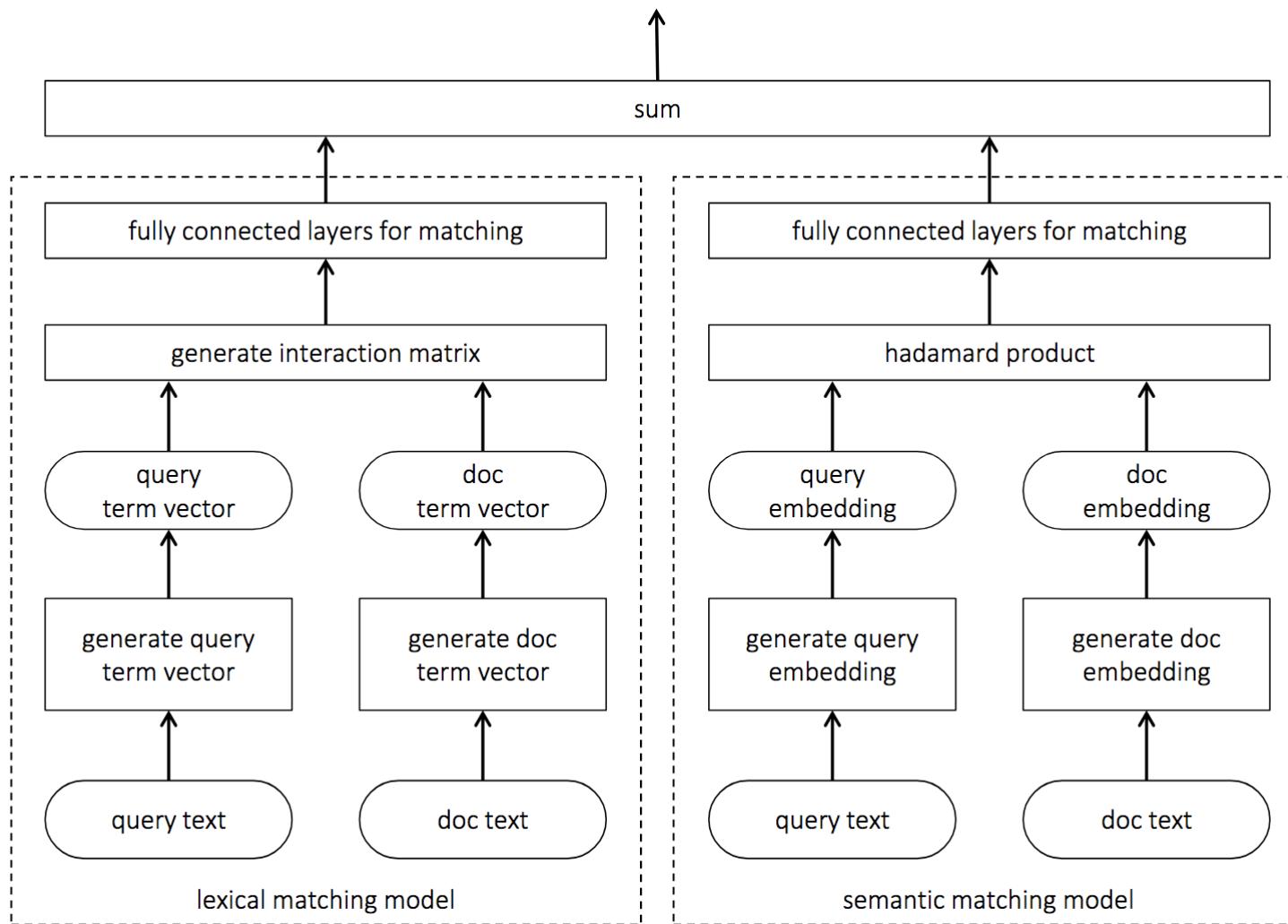
Fixed input length: first 10 terms per query and first 1000 terms per document. Shorter queries/docs are padded with zeros.

Local model takes one-hot encoding as input. Distributed model takes n-graph based representation input
(2000 most frequent uni/bi/tri/4/5-grams).

Patterns of term matches



Architecture (the bigger picture)



Experiments

Domain knowledge is key: "For training, we discarded all documents rated as perfect because a large portion of them fall under the **navigational intent**, which can be better satisfied by historical click based ranking signals"

Training data: **200K queries** sampled from Bing's search logs.

Approximately five documents were judged by human judges as either *perfect, excellent, good, fair, or bad*.

Test data: **8K queries**, with approx. 25 documents per query (**14%** of queries also occurred in the training set)

Four negative documents sampled per relevant document.

	NDCG@1	NDCG@10
Non-neural baselines		
LSA	22.4	44.2
BM25	24.2	45.5
DM	24.7	46.2
QL	24.6	46.3
Neural baselines		
DRMM	24.3	45.2
DSSM	25.8	48.2
CDSSM	27.3	48.2
DESM	25.4	48.3
Our models		
Local model	24.6	45.1
Distributed model	28.6	50.5
Duet model	32.2	53.0

Caveats?

A properly tuned BM25+RM3...

Condition	AP	P20
BM25	0.238	0.354
BM25 + Features	0.250	0.367
Neural _x	0.258	0.372
Neural _y	0.256	0.370
Neural _x + Neural _y	0.373	
A + Neural _y	0.380	
A + Neural _y + M	0.380	
B + Neural _y	0.270	0.383
B + Neural _y + M	0.272	0.386
Anserini: QL	0.2481	0.3517
Anserini: BM25	0.2528	0.3598
Anserini: BM25 + RM3 (independent)	0.2991	0.3901
Anserini: BM25 + RM3 (joint)	0.2956	0.3931

This is a variation
over DRMM...

Maybe that's why DL is not taking over IR yet...

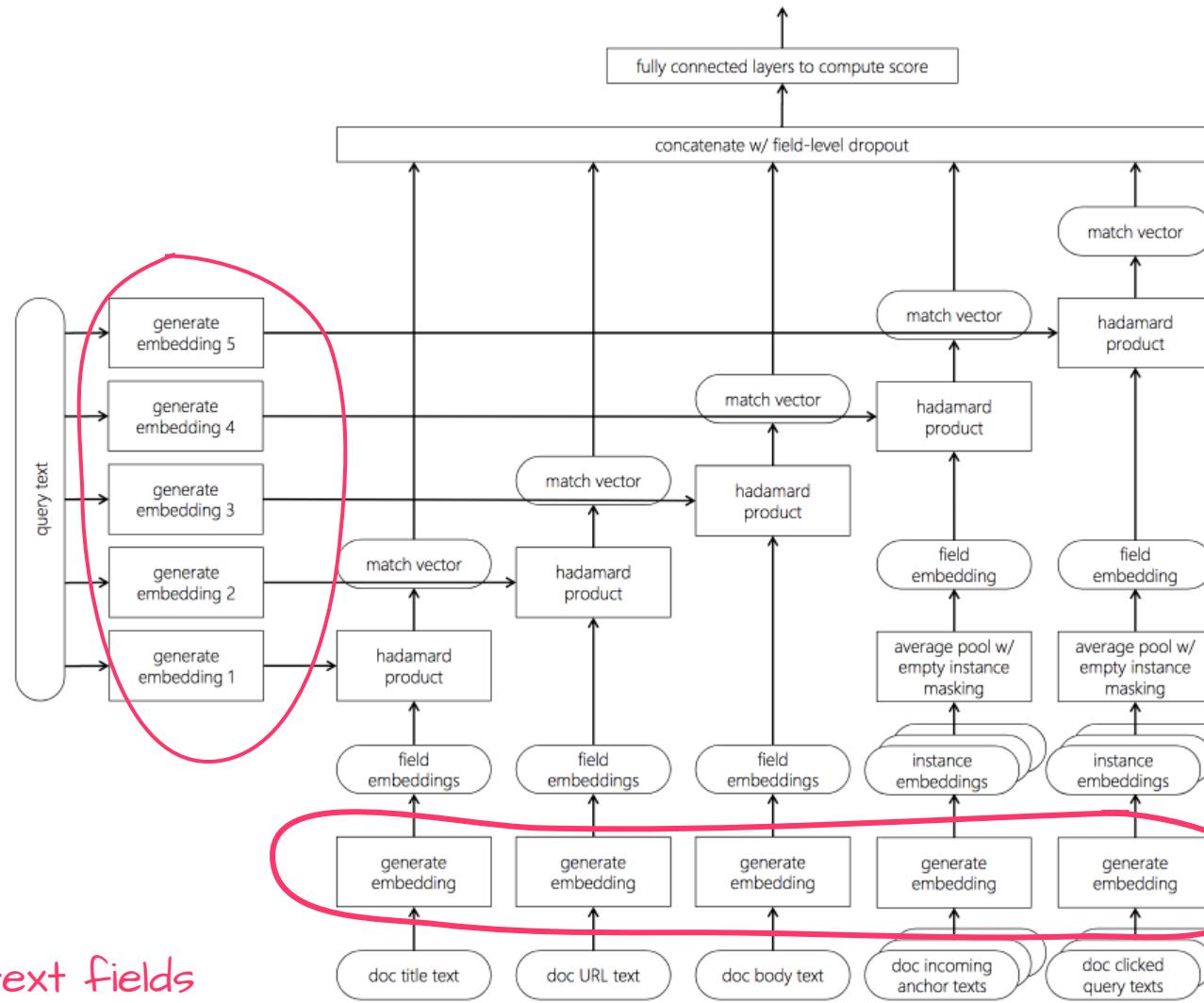
Insights

- Generally, interaction-focused techniques work best (Nie, et al., ICTR '18).
- For short texts, vocabulary mismatch is more serious than for long texts
- MatchZoo has lots of models implemented in Python/TensorFlow (**BEWARE!** **Implementations may be different from original paper!**) [<https://github.com/NTMC-Community/MatchZoo/>].
- Properly turned "classical" IR models are still good (if not better than neural) at most IR tasks.



What now?

Designing more architectures is not a problem ...



text fields

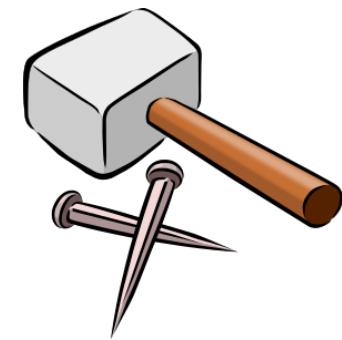
Learning embeddings
separately performs
better than one
common embeddings
space

“

However, given the pace at which the area of deep learning is growing [...] we should be wary of the combinatorial explosion of trying every model on every IR task.

We should not disproportionately focus on maximizing quantitative improvements and in the process neglect theoretical understanding and qualitative insights. [...]

Neural models should not be the hammer that we try on every IR task, or we may risk reducing every IR task to a nail.



Many worthwhile issues to tackle

- Lack of large-scale training data for supervised DNN remains the biggest obstacle (TREC-style evaluations will never scale up, we cannot leave the field to industry); **weak supervision** has attracted a lot of attention recently
- So far, most DNNs have been applied to ad hoc retrieval and Web search, but of course IR is much bigger than that
- Clever tricks from the machine learning community are not yet often used (e.g. **curriculum learning**)
- **Reproducibility** and benchmarking
- **Metrics** (of course!) for new tasks, e.g. answer generation
- Instead of applying DNN building blocks, we should **develop our own** (relevance, IDF, ...)



That's it for today!

Slack: in4325_2019.slack.com

Email: in4325-ewi@tudelft.nl

MSc project ideas? A.BarbosaCamara@tudelft.nl