

---

# Neural Models for Information Retrieval

---

**Bhaskar Mitra**  
Microsoft, UCL\*  
Cambridge, UK  
bmitra@microsoft.com

**Nick Craswell**  
Microsoft  
Bellevue, USA  
nickcr@microsoft.com

## Abstract

Neural ranking models for information retrieval (IR) use shallow or deep neural networks to rank search results in response to a query. Traditional learning to rank models employ machine learning techniques over hand-crafted IR features. By contrast, neural models learn representations of language from raw text that can bridge the gap between query and document vocabulary. Unlike classical IR models, these new machine learning based approaches are data-hungry, requiring large scale training data before they can be deployed. This tutorial introduces basic concepts and intuitions behind neural IR models, and places them in the context of traditional retrieval models. We begin by introducing fundamental concepts of IR and different neural and non-neural approaches to learning vector representations of text. We then review shallow neural IR methods that employ pre-trained neural term embeddings without learning the IR task end-to-end. We introduce deep neural networks next, discussing popular deep architectures. Finally, we review the current DNN models for information retrieval. We conclude with a discussion on potential future directions for neural IR.

## 1 Introduction

Since the turn of the decade, there have been dramatic improvements in performance in computer vision, speech recognition, and machine translation tasks, witnessed in research and in real-world applications [112]. These breakthroughs were largely fuelled by recent advances in neural network models, usually with multiple hidden layers, known as deep architectures [8, 49, 81, 103, 112]. Exciting novel applications, such as conversational agents [185, 203], have also emerged, as well as game-playing agents with human-level performance [147, 180]. Work has now begun in the information retrieval (IR) community to apply these neural methods, leading to the possibility of advancing the state of the art or even achieving breakthrough performance as in these other fields.

Retrieval of information can take many forms. Users can express their information need in the form of a text query—by typing on a keyboard, by selecting a query suggestion, or by voice recognition—or the query can be in the form of an image, or in some cases the need can even be implicit. Retrieval can involve ranking existing pieces of content, such as documents or short-text answers, or composing new responses incorporating retrieved information. Both the information need and the retrieved results may use the same modality (e.g., retrieving text documents in response to keyword queries), or different ones (e.g., image search using text queries). Retrieval systems may consider user history, physical location, temporal changes in information, or other context when ranking results. They may also help users formulate their intent (e.g., via query auto-completion or query suggestion) and/or extract succinct summaries of results for easier inspection.

*Neural IR* refers to the application of shallow or deep neural networks to these retrieval tasks. This tutorial serves as an introduction to neural methods for ranking documents in response to a query, an

---

\*The author is a part-time PhD student at University College London.

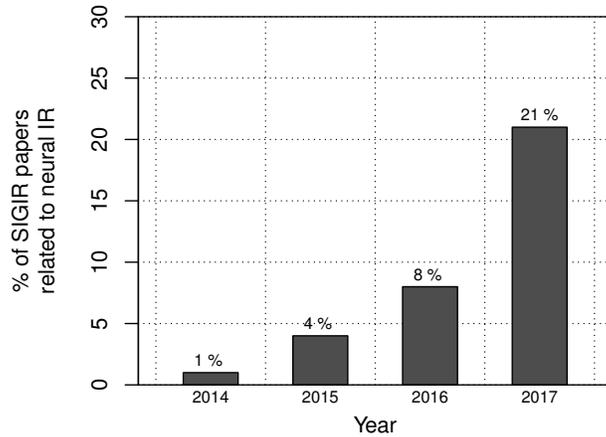


Figure 1: The percentage of neural IR papers at the ACM SIGIR conference—as determined by a manual inspection of the paper titles—shows a clear trend in the growing popularity of the field.

important IR task. A search query may typically contain a few terms, while the document length, depending on the scenario, may range from a few terms to hundreds of sentences or more. Neural models for IR use vector representations of text, and usually contain a large number of parameters that needs to be tuned. ML models with large set of parameters typically require a large quantity of training data [196]. Unlike traditional *learning to rank* (L2R) approaches that train ML models over a set of hand-crafted features, neural models for IR typically accept the raw text of a query and document as input. Learning suitable representations of text also demands large-scale datasets for training [141]. Therefore, unlike classical IR models, these neural approaches tend to be data-hungry, with performance that improves with more training data.

Text representations can be learnt in an unsupervised or supervised fashion. The supervised approach uses IR data such as labeled query-document pairs, to learn a representation that is optimized end-to-end for the task at hand. If sufficient IR labels are not available, the unsupervised approach learns a representation using just the queries and/or documents. In the latter case, different unsupervised learning setups may lead to different vector representations, that differ in the *notion of similarity* that they capture between represented items. When applying such representations, the choice of unsupervised learning setup should be carefully considered, to yield a notion of text similarity that is suitable for the target task. Traditional IR models such as *Latent Semantic Analysis* (LSA) [48] learn dense vector representations of terms and documents. Neural representation learning models share some commonalities with these traditional approaches. Much of our understanding of these traditional approaches from decades of research can be extended to these modern representation learning models.

In other fields, advances in neural networks have been fuelled by specific datasets and application needs. For example, the datasets and successful architectures are quite different in visual object recognition, speech recognition, and game playing agents. While IR shares some common attributes with the field of natural language processing, it also comes with its own set of unique challenges. IR systems must deal with short queries that may contain previously unseen vocabulary, to match against documents that vary in length, to find relevant documents that may also contain large sections of irrelevant text. IR systems should learn patterns in query and document text that indicate relevance, even if query and document use different vocabulary, and even if the patterns are task-specific or context-specific.

The goal of this tutorial is to introduce the fundamentals of neural IR, in context of traditional IR research, with visual examples to illustrate key concepts and a consistent mathematical notation for describing key models. Section 2 presents a survey of IR tasks, challenges, metrics and non-neural models. Section 3 provides a brief overview of neural IR models and a taxonomy for different neural approaches to IR. Section 4 introduces neural and non-neural methods for learning term embeddings, without the use of supervision from IR labels, and with a focus on the notion of similarity. Section 5 surveys some specific approaches for incorporating such embeddings in IR. Section 6 introduces the fundamentals of deep models that are used in IR so far, including popular architectures and toolkits.

Section 7 surveys some specific approaches for incorporating deep neural networks in IR. Section 8 is our discussion, including future work, and conclusion.

**Motivation for this tutorial** Neural IR is an emerging field. Research publication in the area has been increasing (Figure 1), along with relevant workshops [42–44], tutorials [97, 119, 140], and plenary talks [41, 129]. Because this growth in interest is fairly recent, some researchers with IR expertise may be unfamiliar with neural models, and other researchers who have already worked with neural models may be unfamiliar with IR. The purpose of this tutorial is to bridge the gap, by describing the relevant IR concepts and neural methods in the current literature.

## 2 Fundamentals of text retrieval

We focus on text retrieval in IR, where the user enters a text query and the system returns a ranked list of search results. Search results may be passages of text or full text documents. The system’s goal is to rank the user’s preferred search results at the top. This problem is a central one in the IR literature, with well understood challenges and solutions. This section provides an overview of those, such that we can refer to them in subsequent sections.

### 2.1 IR tasks

Text retrieval methods for full text documents and for short text passages have application in ad hoc retrieval systems and question answering systems respectively.

**Ad-hoc retrieval** Ranked document retrieval is a classic problem in information retrieval, as in the main task of the Text Retrieval Conference [205], and performed by popular search engines such as Google, Bing, Baidu, or Yandex. TREC tasks may offer a choice of query length, ranging from a few words to a few sentences, whereas search engine queries tend to be at the shorter end of the range. In an operational search engine, the retrieval system uses specialized index structures to search potentially billions of documents. The results ranking is presented in a search engine results page (SERP), with each result appearing as a summary and a hyperlink. The engine can instrument the SERP, gathering implicit feedback on the quality of search results such as click decisions and dwell times.

A ranking model can take a variety of input features. Some ranking features may depend on the document alone, such as how popular the document is with users, how many incoming links it has, or to what extent document seems problematic according to a Web spam classifier. Other features depend on how the query matches the text content of the document. Still more features match the query against document metadata, such as referred text of incoming hyperlink anchors, or the text of queries from previous users that led to clicks on this document. Because anchors and click queries are a succinct description of the document, they can be a useful source of ranking evidence, but they are not always available. A newly created document would not have much link or click text. Also, not every document is popular enough to have past links and clicks, but it still may be the best search result for a user’s rare or tail query. In such cases, when text metadata is unavailable, it is crucial to estimate the document’s relevance primarily based on its text content.

In the text retrieval community, retrieving documents for short-text queries by considering the long body text of the document is an important challenge. The *ad-hoc* and *Web* tracks<sup>2</sup> at the popular Text REtrieval Conference (TREC) [204] focus specifically on this task. The TREC participants are provided a set of, say fifty, search queries and a document collection containing 500-700K newswire and other documents. Top ranked documents retrieved for each query from the collection by different competing retrieval systems are assessed by human annotators based on their relevance to the query. Given a query, the goal of the IR model is to rank documents with better assessor ratings higher than the rest of the documents in the collection. In Section 2.4, we describe popular IR metrics for quantifying model performance given the ranked documents retrieved by the model and the corresponding assessor judgments for a given query.

**Question-answering** Question-answering tasks may range from choosing between multiple choices (typically entities or binary true-or-false decisions) [78, 80, 165, 212] to ranking spans of text or

<sup>2</sup><http://www10.wwwconference.org/cdrom/papers/317/node2.html>

passages [3, 55, 162, 206, 221], and may even include synthesizing textual responses by gathering evidence from one or more sources [145, 154]. TREC question-answering experiments [206] has participating IR systems retrieve spans of text, rather than documents, in response to questions. IBM’s DeepQA [55] system—behind the Watson project that famously demonstrated human-level performance on the American TV quiz show, "Jeopardy!"—also has a primary search phase, whose goal is to find as many potentially answer-bearing passages of text as possible. With respect to the question-answering task, the scope of this tutorial is limited to ranking answer containing passages in response to natural language questions or short query texts.

Retrieving short spans of text pose different challenges than ranking documents. Unlike the long body text of documents, single sentences or short passages tend to be on point with respect to a single topic. However, answers often tend to use different vocabulary than the one used to frame the question. For example, the span of text that contains the answer to the question "what year was Martin Luther King Jr. born?" may not contain the term "year". However, the phrase "what year" implies that the correct answer text should contain a year—such as '1929' in this case. Therefore, IR systems that focus on the question-answering task need to model the patterns expected in the answer passage based on the intent of the question.

## 2.2 Desiderata of IR models

Before we describe any specific IR model, it is important for us to discuss the attributes that we desire from a good retrieval system. For any IR system, the relevance of the retrieved items to the input query is of foremost importance. But relevance measurements can be nuanced by the properties of *robustness*, *sensitivity* and *efficiency* that we expect the system to demonstrate. These attributes not only guide our model designs but also serve as yard sticks for comparing the different neural and non-neural approaches.

**Semantic understanding** Most traditional approaches for ad-hoc retrieval count repetitions of the query terms in the document text. *Exact term matching* between the query and the document text, while simple, serves as a foundation for many IR systems. Different weighting and normalization schemes over these counts leads to a variety of TF-IDF models, such as BM25 [166].

However, by only inspecting the query terms the IR model ignores all the evidence of *aboutness* from the rest of the document. So, when ranking for the query "Australia", only the occurrences of "Australia" in the document are considered, although the frequency of other words like "Sydney" or "kangaroo" may be highly informative. In the case of the query "what channel are the seahawks on today", the query term "channel" implies that the IR model should pay attention to occurrences of "ESPN" or "Sky Sports" in the document text—none of which appears in the query itself.

Semantic understanding, however, goes beyond mapping query terms to document terms. A good IR model may consider the terms "hot" and "warm" related, as well as the terms "dog" and "puppy"—but must also distinguish that a user who submits the query "hot dog" is not looking for a "warm puppy" [118]. At the more ambitious end of the spectrum, semantic understanding would involve logical reasons by the IR system—so for the query "concerts during SIGIR" it associates a specific edition of the conference (the upcoming one) and considers both its location and dates when recommending concerts nearby during the correct week.

These examples motivate that IR models should have some latent representations of intent as expressed by the query and of the different topics in the document text—so that *inexact matching* can be performed that goes beyond lexical term counting.

**Robustness to rare inputs** Query frequencies in most IR setups follow a Zipfian distribution [216] (see Figure 2). In the publicly available AOL query logs [159], for example, more than 70% of the distinct queries are seen only once in the period of three months from which the queries are sampled. In the same dataset, more than 50% of the distinct documents are clicked only once. A good IR method must be able to retrieve these infrequently searched-for documents, and perform reasonably well on queries containing terms that appear extremely rarely, if ever, in its historical logs.

Many IR models that learn latent representations of text from data often naively assume a fixed size vocabulary. These models perform poorly when the query consists of terms rarely (or never) seen in the training data. Even if the model does not assume a fixed vocabulary, the quality of the latent

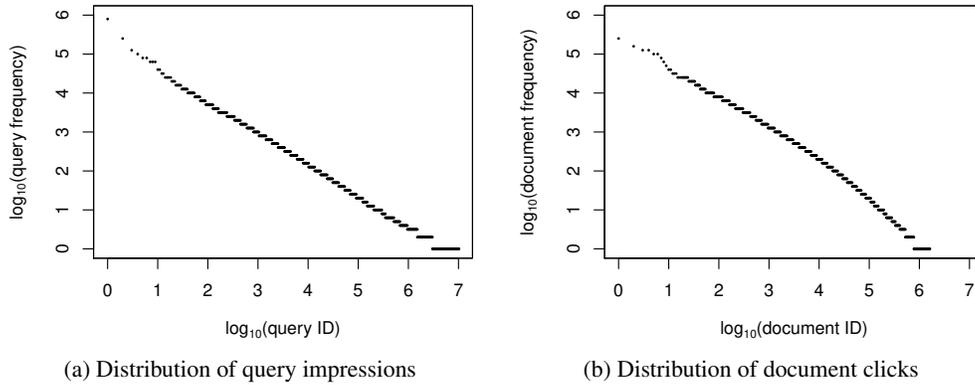


Figure 2: A Log-Log plot of frequency versus rank for query impressions and document clicks in the AOL query logs [159]. The plots highlight that these quantities follow a Zipfian distribution.

representations may depend heavily on how frequently the terms under consideration appear in the training dataset. *Exact matching* models, like BM25 [166], on the other hand can precisely retrieve documents containing rare terms.

Semantic understanding in an IR model cannot come at the cost of poor retrieval performance on queries containing rare terms. When dealing with a query such as “pekarovic land company” the IR model will benefit from considering exact matches of the rare term “pekarovic”. In practice an IR model may need to effectively trade-off exact and inexact matching for a query term. However, the decision of when to perform exact matching can itself be informed by semantic understanding of the context in which the terms appear in addition to the terms themselves.

**Robustness to corpus variance** An interesting consideration for IR models is how well they perform on corpuses whose distributions are different from the data that the model was trained on. Models like BM25 [166] have very few parameters and often demonstrate reasonable performance “out of the box” on new corpuses with little or no additional tuning of parameters. Deep learning models containing millions (or even billions) of parameters, on the other hand, are known to be more sensitive to distributional differences between training and evaluation data, and has been shown to be especially vulnerable to adversarial inputs [194].

Some of the variances in performance of deep models on new corpuses is offset by better retrieval on the test corpus that is distributionally closer to the training data, where the model may have picked up crucial corpus specific patterns. For example, it maybe understandable if a model that learns term representations based on the text of Shakespeare’s Hamlet is effective at retrieving passages relevant to a search query from The Bard’s other works, but performs poorly when the retrieval task involves a corpus of song lyrics by Jay-Z. However, the poor performances on new corpus can also be indicative that the model is overfitting, or suffering from the Clever Hans<sup>3</sup> effect [187]. For example, an IR model trained on recent news corpus may learn to associate “Theresa May” with the query “uk prime minister” and as a consequence may perform poorly on older TREC datasets where the connection to “John Major” may be more appropriate.

ML models that are hyper-sensitive to corpus distributions may be vulnerable when faced with unexpected changes in distributions or “black swans”<sup>4</sup> in the test data. This can be particularly problematic when the test distributions naturally evolve over time due to underlying changes in the user population or behavior. The models, in these cases, may need to be re-trained periodically, or designed to be invariant to such changes.

**Robustness to variable length inputs** A typical text collection contains documents of varied lengths (see Figure 3). For a given query, a good IR system must be able to deal with documents of different lengths without over-retrieving either long or short documents. Relevant documents may

<sup>3</sup>[https://en.wikipedia.org/wiki/Clever\\_Hans](https://en.wikipedia.org/wiki/Clever_Hans)

<sup>4</sup>[https://en.wikipedia.org/wiki/Black\\_swan\\_theory](https://en.wikipedia.org/wiki/Black_swan_theory)

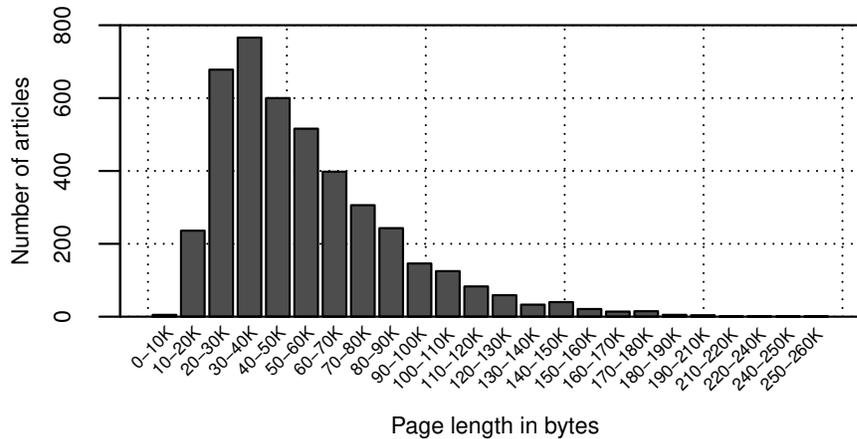


Figure 3: Distribution of Wikipedia featured articles by document length (in bytes) as of June 30, 2014. Source: [https://en.wikipedia.org/wiki/Wikipedia:Featured\\_articles/By\\_length](https://en.wikipedia.org/wiki/Wikipedia:Featured_articles/By_length).

contain irrelevant sections, and the relevant content may either be localized in a single section of the document, or spread over different sections. Document length normalization is well-studied in the context of IR models (e.g., pivoted length normalization [181]), and this existing research should inform the design of any new IR models.

**Robustness to errors in input** No IR system should assume error-free inputs—neither when considering the user query nor when inspecting the documents in the text collection. While traditional IR models have typically involved specific components for error correction—such as automatic spell corrections over queries—new IR models may adopt different strategies towards dealing with such errors by operating at the character-level and/or by learning better representations from noisy texts.

**Sensitivity to context** Retrieval in the wild can leverage many implicit and explicit context information.<sup>5</sup> The query “weather” can refer to the weather in Seattle or in London depending on where the user is located. An IR model may retrieve different results for the query “decorations” depending on the time of the year. The query “giants match highlights” can be better disambiguated if the IR system knows whether the user is a fan of baseball or American football, whether she is located on the East or the West coast of USA, or if the model has knowledge of recent sport fixtures. In a conversational IR system, the correct response to the question “When did she become the prime minister?” would depend on disambiguating the correct entity based on the context of references made in the previous turns of the conversation. Relevance, therefore, in many applications is situated in the user and task context, and is an important consideration in the design of IR systems.

**Efficiency** Efficiency of retrieval is one of the salient points of any retrieval system. A typical commercial Web search engine may deal with tens of thousands of queries per second<sup>6</sup>—retrieving results for each query from an index containing billions of documents. Search engines typically involve large multi-tier architectures and the retrieval process generally consists of multiple stages of pruning the candidate set of documents [131]. The IR model at the bottom of this *telescoping* setup may need to sift through billions of documents—while the model at the top may only need to re-rank between tens of promising documents. The retrieval approaches that are suitable at one level of the stack may be highly impractical at a different step—models at the bottom need to be *fast* but mostly focus on eliminating irrelevant or junk results, while models at the top tend to develop more sophisticated notions of *relevance*, and focus on distinguishing between documents that are much closer on the relevance scale. So far, much of the focus on neural IR approaches have been limited to re-ranking top-*n* documents.

<sup>5</sup>As an extreme example, in the *proactive retrieval* scenario the retrieval can be triggered based solely on implicit context without any explicit query submission from the user.

<sup>6</sup><http://www.internetlivestats.com/one-second/#google-band>

Table 1: Notation used in this tutorial.

Meaning	Notation
Single query	$q$
Single document	$d$
Set of queries	$Q$
Collection of documents	$D$
Term in query $q$	$t_q$
Term in document $d$	$t_d$
Full vocabulary of all terms	$T$
Set of ranked results retrieved for query $q$	$R_q$
Result tuple (document $d$ at rank $i$ )	$\langle i, d \rangle$ , where $\langle i, d \rangle \in R_q$
Ground truth relevance label of document $d$ for query $q$	$rel_q(d)$
$d_i$ is more relevant than $d_j$ for query $q$	$rel_q(d_i) > rel_q(d_j)$ , or succinctly $d_i \succ_q d_j$
Frequency of term $t$ in document $d$	$tf(t, d)$
Number of documents in $D$ that contains term $t$	$df(t)$
Vector representation of text $z$	$\vec{v}_z$
Probability function for an event $\mathcal{E}$	$p(\mathcal{E})$

While this list of desired attributes of an IR model is in no way complete, it serves as a reference for comparing many of the neural and non-neural approaches described in the rest of this tutorial.

### 2.3 Notation

We adopt some common notation for this tutorial shown in Table 1. We use lower-case to denote vectors (e.g.,  $\vec{x}$ ) and upper-case for tensors of higher dimensions (e.g.,  $X$ ). The ground truth  $rel_q(d)$  in Table 1 may be based on either manual relevance annotations or be implicitly derived from user behaviour on SERP (e.g., from clicks).

### 2.4 Metrics

A large number of IR studies [52, 65, 70, 84, 92, 93, 106, 144] have demonstrated that users of retrieval systems tend to pay attention mostly to top-ranked results. IR metrics, therefore, focus on rank-based comparisons of the retrieved result set  $R$  to an ideal ranking of documents, as determined by manual judgments or implicit feedback from user behaviour data. These metrics are typically computed at a rank position, say  $k$ , and then averaged over all queries in the test set. Unless otherwise specified,  $R$  refers to the top- $k$  results retrieved by the model. Next, we describe a few popular metrics used in IR evaluations.

**Precision and recall** Precision and recall both compute the fraction of relevant documents retrieved for a query  $q$ , but with respect to the total number of documents in the retrieved set  $R_q$  and the total number of relevant documents in the collection  $D$ , respectively. Both metrics assume that the relevance labels are binary.

$$Precision_q = \frac{\sum_{\langle i, d \rangle \in R_q} rel_q(d)}{|R_q|} \quad (1)$$

$$Recall_q = \frac{\sum_{\langle i, d \rangle \in R_q} rel_q(d)}{\sum_{d \in D} rel_q(d)} \quad (2)$$

**Mean reciprocal rank (MRR)** Mean reciprocal rank [40] is also computed over binary relevance judgments. It is given as the reciprocal rank of the first relevant document averaged over all queries.

$$RR_q = \max_{\langle i, d \rangle \in R_q} \frac{rel_q(d)}{i} \quad (3)$$

**Mean average precision (MAP)** The average precision [235] for a ranked list of documents  $R$  is given by,

$$AveP_q = \frac{\sum_{\langle i,d \rangle \in R_q} Precision_{q,i} \times rel_q(d)}{\sum_{d \in D} rel_q(d)} \quad (4)$$

where,  $Precision_{q,i}$  is the precision computed at rank  $i$  for the query  $q$ . The average precision metric is generally used when relevance judgments are binary, although variants using graded judgments have also been proposed [167]. The mean of the average precision over all queries gives the MAP score for the whole set.

**Normalized discounted cumulative gain (NDCG)** There are few different variants of the discounted cumulative gain ( $DCG_q$ ) metric [90] which can be used when graded relevance judgments are available for a query  $q$ —say, on a five-point scale between zero to four. A popular incarnation of this metric is as follows.

$$DCG_q = \sum_{\langle i,d \rangle \in R_q} \frac{2^{rel_q(d)} - 1}{\log_2(i + 1)} \quad (5)$$

The ideal DCG ( $IDCG_q$ ) is computed the same way but by assuming an ideal rank order for the documents up to rank  $k$ . The normalized DCG ( $NDCG_q$ ) is then given by,

$$NDCG_q = \frac{DCG_q}{IDCG_q} \quad (6)$$

## 2.5 Traditional IR models

In this section, we introduce a few of the traditionally popular IR approaches. The decades of insights from these IR models not only inform the design of our new neural based approaches, but these models also serve as important baselines for comparison. They also highlight the various desiderata that we expect the neural IR models to incorporate.

**TF-IDF** There is a broad family of statistical functions in IR that consider the number of occurrences of each query term in the document (term-frequency) and the corresponding inverse document frequency of the same terms in the full collection (as an indicator of the informativeness of the term). One theoretical basis for such formulations is the probabilistic model of IR that yielded the popular BM25 [166] ranking function.

$$BM25(q, d) = \sum_{t_q \in q} idf(t_q) \cdot \frac{tf(t_q, d) \cdot (k_1 + 1)}{tf(t_q, d) + k_1 \cdot \left(1 - b + b \cdot \frac{|d|}{avgdl}\right)} \quad (7)$$

where,  $avgdl$  is the average length of documents in the collection  $D$ , and  $k_1$  and  $b$  are parameters that are usually tuned on a validation dataset. In practice,  $k_1$  is sometimes set to some default value in the range [1.2, 2.0] and  $b$  as 0.75. The  $idf(t)$  is popularly computed as,

$$idf(t) = \log \frac{|D| - df(t) + 0.5}{df(t) + 0.5} \quad (8)$$

BM25 aggregates the contributions from individual terms but ignores any phrasal or proximity signals between the occurrences of the different query terms in the document. A variant of BM25 [229] also considers documents as composed of several fields (such as, title, body, and anchor texts).

**Language modelling (LM)** In the language modelling based approach [79, 161, 230], documents are ranked by the posterior probability  $p(d|q)$ .

$$p(d|q) = \frac{p(q|d) \cdot p(d)}{\sum_{\bar{d} \in D} p(q|\bar{d}) \cdot p(\bar{d})} \propto p(q|d) \cdot p(d) \quad (9)$$

$$= p(q|d) \quad , \text{ assuming } p(d) \text{ is uniform} \quad (10)$$

$$= \prod_{t_q \in q} p(t_q|d) \quad (11)$$

$$= \prod_{t_q \in q} \left( \lambda \hat{p}(t_q|d) + (1 - \lambda) \hat{p}(t_q|D) \right) \quad (12)$$

$$= \prod_{t_q \in q} \left( \lambda \frac{tf(t_q, d)}{|d|} + (1 - \lambda) \frac{\sum_{\bar{d} \in D} tf(t_q, \bar{d})}{\sum_{\bar{d} \in D} |\bar{d}|} \right) \quad (13)$$

where,  $\hat{p}(\mathcal{E})$  is the maximum likelihood estimate (MLE) of the probability of event  $\mathcal{E}$ .  $p(q|d)$  indicates the probability of generating query  $q$  by randomly sampling terms from document  $d$ . For smoothing, terms are sampled from both the document  $d$  and the full collection  $D$ —the two events are treated as mutually exclusive, and their probability is given by  $\lambda$  and  $(1 - \lambda)$ , respectively.

Both TF-IDF and language modelling based approaches estimate document relevance based on the count of only the query terms in the document. The position of these occurrences and the relationship with other terms in the document are ignored.

**Translation models** Berger and Lafferty [17] proposed an alternative method to estimate  $p(t_q|d)$  in the language modelling based IR approach (Equation 11), by assuming that the query  $q$  is being generated via a "translation" process from the document  $d$ .

$$p(t_q|d) = \sum_{t_d \in d} p(t_q|t_d) \cdot p(t_d|d) \quad (14)$$

The  $p(t_q|t_d)$  component allows the model to garner evidence of relevance from non-query terms in the document. Berger and Lafferty [17] propose to estimate  $p(t_q|t_d)$  from query-document paired data similar to popular techniques in statistical machine translation [22, 23]—but other approaches for estimation have also been explored [236].

**Dependence model** None of the three IR models described so far consider proximity between query terms. To address this, Metzler and Croft [132] proposed a linear model over proximity-based features.

$$\begin{aligned} DM(q, d) = & (1 - \lambda_{ow} - \lambda_{uw}) \sum_{t_q \in q} \log \left( (1 - \alpha_d) \frac{tf(t_q, d)}{|d|} + \alpha_d \frac{\sum_{\bar{d} \in D} tf(t_q, \bar{d})}{\sum_{\bar{d} \in D} |\bar{d}|} \right) \\ & + \lambda_{ow} \sum_{c_q \in ow(q)} \log \left( (1 - \alpha_d) \frac{tf_{\#1}(c_q, d)}{|d|} + \alpha_d \frac{\sum_{\bar{d} \in D} tf_{\#1}(c_q, \bar{d})}{\sum_{\bar{d} \in D} |\bar{d}|} \right) \\ & + \lambda_{uw} \sum_{c_q \in uw(q)} \log \left( (1 - \alpha_d) \frac{tf_{\#uwN}(c_q, d)}{|d|} + \alpha_d \frac{\sum_{\bar{d} \in D} tf_{\#uwN}(c_q, \bar{d})}{\sum_{\bar{d} \in D} |\bar{d}|} \right) \end{aligned} \quad (15)$$

where,  $ow(q)$  and  $uw(q)$  are the set of all contiguous  $n$ -grams (or phrases) and the set of all bags of terms that can be generated from query  $q$ .  $tf_{\#1}$  and  $tf_{\#uwN}$  are the ordered-window and unordered-window operators from Indri [186]. Finally,  $\lambda_{ow}$  and  $\lambda_{uw}$  are the tunable parameters of the model.

**Pseudo relevance feedback (PRF)** PRF-based methods, such as Relevance Models (RM) [108, 109], typically demonstrate strong performance at the cost of executing an additional round of retrieval. The set of ranked documents  $R_1$  from the first round of retrieval is used to select expansion terms to augment the query for the second round of retrieval. The ranked set  $R_2$  from the second round are presented to the user.

The underlying approach to scoring a document in RM is by computing the KL divergence [105] between the query language model  $\theta_q$  and the document language model  $\theta_d$ .

$$score(q, d) = - \sum_{t \in T} p(t|\theta_q) \log \frac{p(t|\theta_q)}{p(t|\theta_d)} \quad (16)$$

Without PRF,

$$p(t|\theta_q) = \frac{tf(t, q)}{|q|} \quad (17)$$

But under the popular RM3 [2] formulation the new query language model  $\bar{\theta}_q$  is estimated by,

$$p(t|\bar{\theta}_q) = \alpha \frac{tf(t, q)}{|q|} + (1 - \alpha) \sum_{d \in R_1} p(t|\theta_d) p(d) \prod_{\bar{t} \in q} p(\bar{t}|\theta_d) \quad (18)$$

By expanding the query using the results from the first round of retrieval PRF based approaches tend to be more robust to the vocabulary mismatch problem plaguing many other traditional IR models.

## 2.6 Learning to rank (L2R)

In learning to rank, a query-document pair is represented by a vector of numerical features  $\vec{x} \in \mathbb{R}^n$ , and a model  $f : \vec{x} \rightarrow \mathbb{R}$  is trained that maps the feature vector to a real-valued score. The training dataset for the model consists of a set of queries and a set of documents per query. Depending on the flavour of L2R, in addition to the feature vector, each query-document pair in the training data is augmented with some relevance information. Liu [121] categorized the different L2R approaches based on their training objectives.

- In the *pointwise approach*, the relevance information  $rel_q(d)$  is in the form of a numerical value associated with every query-document pair with feature vector  $\vec{x}_{q,d}$ . The numerical relevance label can be derived from binary or graded relevance judgments or from implicit user feedback, such as clickthrough information. A regression model is typically trained on the data to predict the numerical value  $rel_q(d)$  given  $\vec{x}_{q,d}$ .
- In the *pairwise approach*, the relevance information is in the form of preferences between pairs of documents with respect to individual queries (e.g.,  $d_i \succ_q d_j$ ). The ranking problem in this case reduces to binary classification for predicting the more relevant document.
- Finally, the *listwise approach* involves directly optimizing for a rank-based metric—which is difficult because these metrics are often not continuous (and hence not differentiable) with respect to the model parameters.

The input features for L2R models typically belong to one of three categories.

- *Query-independent* or *static* features (e.g., PageRank or spam score of the document)
- *Query-dependent* or *dynamic* features (e.g., BM25)
- *Query-level* features (e.g., number of words in query)

Many machine learning models—including support vector machines, neural networks, and boosted decision trees—have been employed over the years for the learning to rank task, and a correspondingly

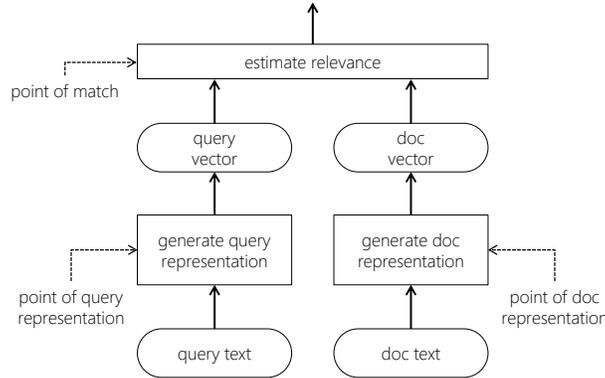


Figure 4: Document ranking typically involves a query and a document representation steps, followed by a matching stage. Neural models can be useful either for generating good representations or in estimating relevance, or both.

large number of different loss functions have been explored. Next, we briefly describe RankNet [26] that has been a popular choice for training neural L2R models and was also—for many years—an industry favourite, such as at the commercial Web search engine Bing.<sup>7</sup>

**RankNet** RankNet [26] is *pairwise* loss function. For a given query  $q$ , a pair of documents  $\langle d_i, d_j \rangle$ , with different relevance labels, such that  $d_i \succ_q d_j$ , and feature vectors  $\langle \vec{x}_i, \vec{x}_j \rangle$ , is chosen. The model  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , typically a neural network but can also be any other machine learning model whose output is differentiable with respect to its parameters, computes the scores  $s_i = f(\vec{x}_i)$  and  $s_j = f(\vec{x}_j)$ , such that ideally  $s_i > s_j$ . Given the output scores  $\langle s_i, s_j \rangle$  from the model corresponding to the two documents, the probability that  $d_i$  would be ranked higher than  $d_j$  is given by,

$$p_{ij} \equiv p(d_i \succ_q d_j) \equiv \frac{1}{1 + e^{-\sigma(s_i - s_j)}} \quad (19)$$

where,  $\sigma$  determines the shape of the sigmoid. Let  $S_{ij} \in \{-1, 0, +1\}$  be the true preference label between  $d_i$  and  $d_j$  for the training sample—denoting  $d_i$  is more, equal, or less relevant than  $d_j$ , respectively. Then the desired probability of ranking  $d_i$  over  $d_j$  is given by  $\bar{p}_{ij} = \frac{1}{2}(1 + S_{ij})$ . The cross-entropy loss  $\mathcal{L}$  between the desired probability  $\bar{p}_{ij}$  and the predicted probability  $p_{ij}$  is given by,

$$\mathcal{L} = -\bar{p}_{ij} \log(p_{ij}) - (1 - \bar{p}_{ij}) \log(1 - p_{ij}) \quad (20)$$

$$= \frac{1}{2} (1 - S_{ij}) \sigma(s_i - s_j) + \log(1 + e^{-\sigma(s_i - s_j)}) \quad (21)$$

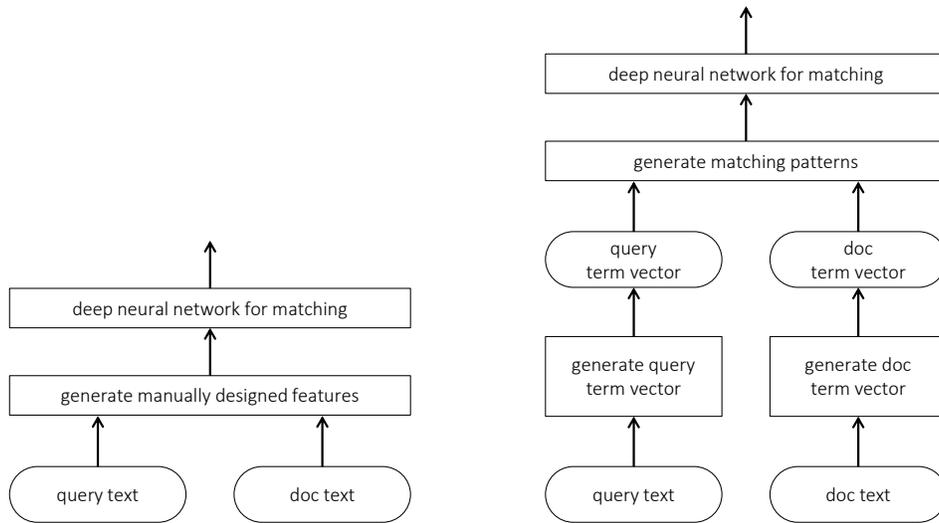
$$= \log(1 + e^{-\sigma(s_i - s_j)}) \quad \text{if, documents are ordered such that } d_i \succ_q d_j (S_{ij} = 1) \quad (22)$$

Note that  $\mathcal{L}$  is differentiable with respect to the model output  $s_i$  and hence the model can be trained using gradient descent. We direct the interested reader to [27] for more detailed derivations for computing the gradients for RankNet and for the evolution to the *listwise* models LambdaRank [28] and LambdaMART [214].

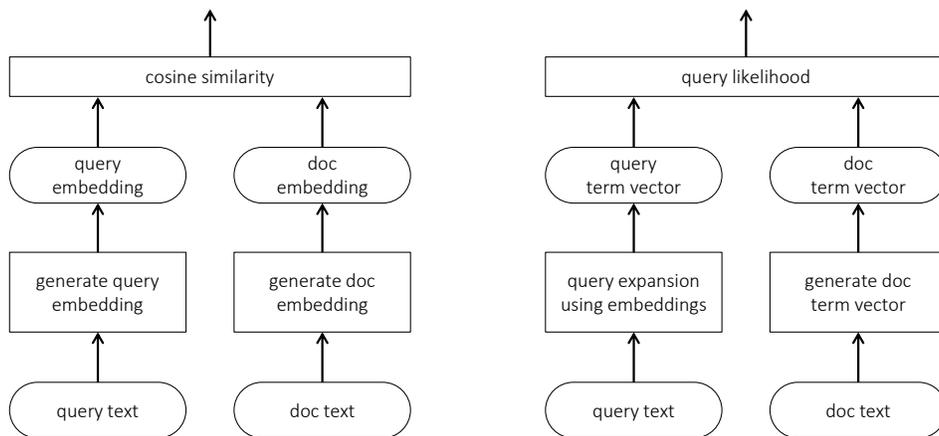
### 3 Anatomy of a neural IR model

At a high level, document ranking comprises of performing three primary steps—generate a representation of the query that specifies the information need, generate a representation of the document

<sup>7</sup><https://www.microsoft.com/en-us/research/blog/ranknet-a-ranking-retrospective/>



(a) Learning to rank using manually designed features (e.g., Liu [121])      (b) Estimating relevance from patterns of exact matches (e.g., [71, 141])



(c) Learning query and document representations for matching (e.g., [88, 143])      (d) Query expansion using neural embeddings (e.g., [51, 170])

Figure 5: Examples of different neural approaches to IR. In (a) and (b) the neural network is only used at the point of matching, whereas in (c) the focus is on learning effective representations of text using neural methods. Neural models can also be used to expand or augment the query before applying traditional IR techniques, as shown in (d).

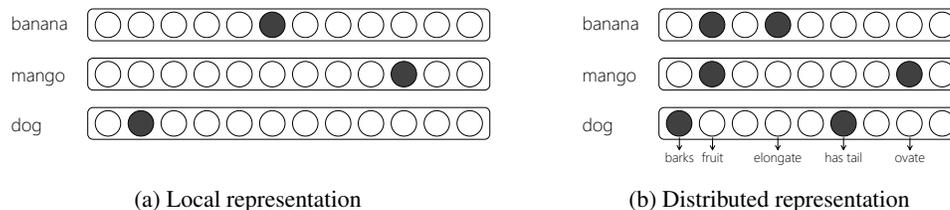


Figure 6: Under local representations the terms “banana”, “mango”, and “dog” are distinct items. But distributed vector representations may recognize that “banana” and “mango” are both fruits, but “dog” is different.

that captures the distribution over the information contained, and match the query and the document representations to estimate their mutual relevance. All existing neural approaches to IR can be broadly categorized based on whether they influence the query representation, the document representation, or in estimating relevance. A neural approach may impact one or more of these stages shown in Figure 4.

Neural networks are popular as learning to rank models discussed in Section 2.6. In these models, a joint representation of the query and the document is generated using manually designed features and the neural network is used only at the *point of match* to estimate relevance, as shown in Figure 5a. In Section 7.4, we will discuss deep neural network models, such as [71, 141], that estimate relevance based on patterns of exact query term matches in the document. Unlike traditional learning to rank models, however, these architectures (shown in Figure 5b) depend less on manual feature engineering and more on automatically detecting regularities in good matching patterns.

In contrast, many (shallow and deep) neural IR models depend on learning good low-dimensional vector representations—or *embeddings*—of query and document text, and using them within traditional IR models or in conjunction with simple similarity metrics (e.g., cosine similarity). These models shown in Figure 5c may learn the embeddings by optimizing directly for the IR task (e.g., [88]), or separately in an unsupervised fashion (e.g., [143]). Finally, Figure 5d shows IR approaches where the neural models are used for query expansion [51, 170].

While the taxonomy of neural approaches described in this section is rather simple, it does provide an intuitive framework for comparing the different neural approaches in IR, and highlights the similarities and distinctions between these different techniques.

## 4 Term representations

### 4.1 A tale of two representations

Vector representations are fundamental to both information retrieval and machine learning. In IR, terms are typically the smallest unit of representation for indexing and retrieval. Therefore, many IR models—both neural and non-neural—focus on learning good vector representations of terms.

Different vector representations exhibit different levels of generalization—some consider every term as distinct entities while others learn to identify common attributes. Different representation schemes derive different notions of similarity between terms from the definition of the corresponding vector spaces. Some representations operate over fixed-size vocabularies, while the design of others obviate such constraints. They also differ on the properties of compositionality that defines how representations for larger units of information, such as passages and documents, can be derived from individual term vectors. These are some of the important considerations for choosing a term representation suitable for a specific task.

**Local representations** Under local (or *one-hot*) representations, every term in a fixed size vocabulary  $T$  is represented by a binary vector  $\vec{v} \in \{0, 1\}^{|T|}$ , where only one of the values in the vector is one and all the others are set to zero. Each position in the vector  $\vec{v}$  corresponds to a term. The term “banana”, under this representation, is given by a vector that has the value one in the position corresponding to “banana” and zero everywhere else. Similarly, the terms “mango” and “dog” are represented by setting different positions in the vector to one.

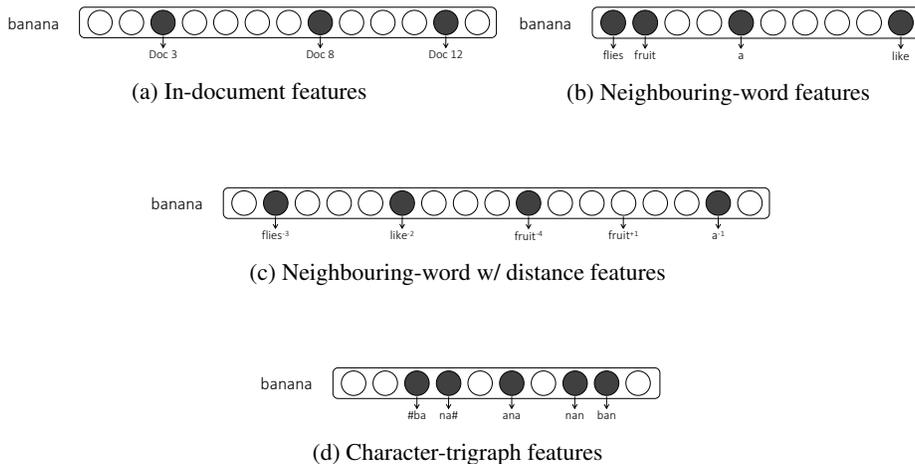


Figure 7: Examples of different feature-based distributed representations of the term “banana”. The representations in (a), (b), and (c) are based on external contexts in which the term frequently occurs, while (d) is based on properties intrinsic to the term. The representation scheme in (a) depends on the documents containing the term, while the scheme shown in (b) and (c) depends on other terms that appears in its neighbourhood. The scheme (b) ignores inter-term distances. Therefore, in the sentence “Time flies like an arrow; fruit flies like a banana”, the feature “fruit” describes both the terms “banana” and “arrow”. However, in the representation scheme of (c) the feature “fruit<sup>-4</sup>” is positive for “banana”, and the feature “fruit<sup>+1</sup>” for “arrow”.

Figure 6a highlights that under this scheme each term is a unique entity, and “banana” is as distinct from “dog” as it is from “mango”. Terms outside of the vocabulary either have no representation, or are denoted by a special “UNK” symbol, under this scheme.

**Distributed representations** Under distributed representations every term is represented by a vector  $\vec{v} \in \mathbb{R}^{k_l}$ .  $\vec{v}$  can be a sparse or a dense vector—a vector of hand-crafted features or a learnt representation in which the individual dimensions are not interpretable in isolation. The key underlying hypothesis for any distributed representation scheme, however, is that by representing a term by its attributes allows for defining some notion of similarity between the different terms based on the chosen properties. For example, in Figure 6b “banana” is more similar to “mango” than “dog” because they are both fruits, but yet different because of other properties that are not shared between the two, such as shape.

A key consideration in any feature based distributed representation is the choice of the features themselves. A popular approach involves representing terms by features that capture their distributional properties. This is motivated by the *distributional hypothesis* [75] that states that terms that are used (or occur) in similar context tend to be semantically similar. Firth [56] famously purported this idea of *distributional semantics*<sup>8</sup> by stating “a word is characterized by the company it keeps”. However, both *distribution* and *semantics* by themselves are not well-defined and under different context may mean very different things. Figure 7 shows three different sparse vector representations of the term “banana” corresponding to different distributional feature spaces—documents containing the term (e.g., LSA [48]), neighbouring words in a window (e.g., HAL [125], COALS [168], and [24]), and neighbouring words with distance (e.g., [117]). Finally, Figure 7d shows a vector representation of “banana” based on the character trigrams in the term itself—instead of external contexts in which the term occurs. In Section 4.2 we will discuss how choosing different distributional features for term representation leads to different nuanced notions of semantic similarity between them.

<sup>8</sup>Readers should take note that while many *distributed* representations take advantage of *distributional* properties, the two concepts are not synonymous. A term can have a distributed representation based on non-distributional features—e.g., parts of speech classification and character trigrams in the term.

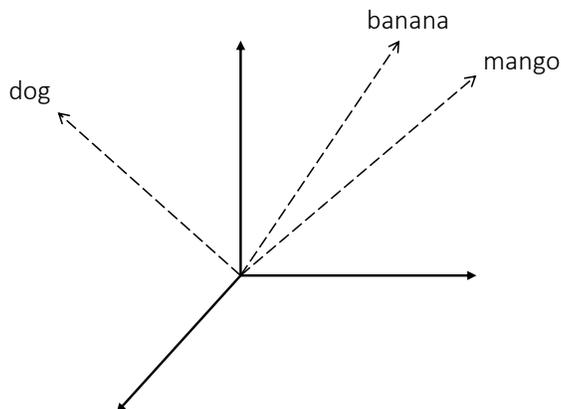


Figure 8: A vector space representation of terms puts “banana” closer to “mango” because they share more common attributes than “banana” and “dog”.

When the vectors are high-dimensional, sparse, and based on distributional feature they are referred to as *explicit* vector representations [117]. On the other hand, when the vectors are dense, small ( $k \ll |T|$ ), and learnt from data then they are commonly referred to as *embeddings*. For both explicit and embedding based representations several distance metrics can be used to define similarity between terms, although cosine similarity is commonly used.

$$\text{sim}(\vec{v}_i, \vec{v}_j) = \cos(\vec{v}_i, \vec{v}_j) = \frac{\vec{v}_i^\top \vec{v}_j}{\|\vec{v}_i\| \|\vec{v}_j\|} \quad (23)$$

Most embeddings are learnt from explicit vector space representations, and hence the discussions in 4.2 about different notions of similarity are also relevant to the embedding models. In Section 4.3 and 4.4 we briefly discuss explicit and embedding based representations.

With respect to *compositionality*, it is important to understand that *distributed* representations of items are often derived from local or distributed representation of its parts. For example, a document can be represented by the sum of the one-hot vectors or embeddings corresponding to the terms in the document. The resultant vector, in both cases, corresponds to a distributed bag-of-word representation. Similarly, the character trigraph representation of terms in Figure 7d is simply an aggregation over the one-hot representations of the constituent trigraphs.

In the context of neural models, distributed representations generally refer to learnt embeddings. The idea of ‘local’ and ‘distributed’ representations has a specific significance in the context of neural network models. Each concept, entity, or term can be represented within a neural network by the activation of a single neuron (local representation) or by the combined pattern of activations of several neurons (distributed representation) [82].

## 4.2 Notions of similarity

Any vector representation inherently defines some notion of relatedness between terms. Is “Seattle” closer to “Sydney” or to “Seahawks”? The answer depends on the type of relationship we are interested in. If we want terms of similar *type* to be closer, then “Sydney” is more similar to “Seattle” because they are both cities. However, if we are interested to find terms that co-occur in the same document or passage, then “Seahawks”—Seattle’s football team—should be closer. The former represents a *Typical*, or type-based notion of similarity while the latter exhibits a more *Topical* sense of relatedness.

If we want to compare “Seattle” with “Sydney” and “Seahawks” based on their respective vector representations, then the underlying feature space needs to align with the notion of similarity that we are interested in. It is, therefore, important for the readers to build an intuition about the choice of features and the notion of similarity they encompass. This can be demonstrated by using a toy corpus, such as the one in Table 2. Figure 9a shows that the “in documents” features naturally lend

Table 2: A toy corpus of short documents that we consider for the discussion on different notions of similarity between terms under different distributed representations. The choice of the feature space that is used for generating the distributed representation determines which terms are closer in the vector space, as shown in Figure 9.

Sample documents			
doc 01	Seattle map	doc 09	Denver map
doc 02	Seattle weather	doc 10	Denver weather
doc 03	Seahawks jerseys	doc 11	Broncos jerseys
doc 04	Seahawks highlights	doc 12	Broncos highlights
doc 05	Seattle Seahawks Wilson	doc 13	Denver Broncos Lynch
doc 06	Seattle Seahawks Sherman	doc 14	Denver Broncos Sanchez
doc 07	Seattle Seahawks Browner	doc 15	Denver Broncos Miller
doc 08	Seattle Seahawks Ifedi	doc 16	Denver Broncos Marshall

to a Topical sense of similarity between the terms, while the “neighbouring terms with distances” features in Figure 9c gives rise to a more Typical notion of relatedness. Using “neighbouring terms” without the inter-term distances as features, however, produces a mixture of Topical and Typical relationships. This is because when the term distances are considered in feature definition then the document “Seattle Seahawks Wilson” produces the bag-of-features  $\{Seahawks^{+1}, Wilson^{+2}\}$  for “Seattle” which is non-overlapping with the bag-of-features  $\{Seattle^{-1}, Wilson^{+1}\}$  for “Seahawks”. However, when the feature definition ignores the term-distances then there is a partial overlap between the bag-of-features  $\{Seahawks, Wilson\}$  and  $\{Seattle, Wilson\}$  corresponding to “Seattle” and “Seahawks”. The overlap increases significantly when we use a larger window-size for identifying neighbouring terms pushing the notion of similarity closer to a Topical definition. This effect of the windows size on the Topicality of the representation space was reported by Levy and Goldberg [115] in the context of learnt embeddings.

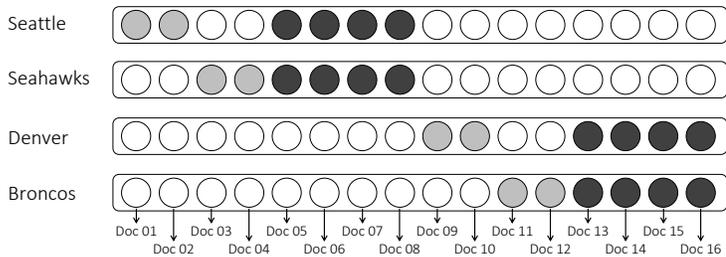
Readers should take note that the set of all inter-term relationships goes far beyond the two notions of Typical and Topical that we discuss in this section. For example, vector representations could cluster terms closer based on linguistic styles—e.g., terms that appear in thriller novels versus in children’s rhymes, or in British versus American English. However, the notions of Typical and Topical similarities popularly come up in discussions in the context of many IR and NLP tasks—sometimes under different names such as *Paradigmatic* and *Syntagmatic* relations<sup>9</sup>—and the idea itself goes back at least as far as Saussure [30, 47, 74, 172].

### 4.3 Explicit vector representations

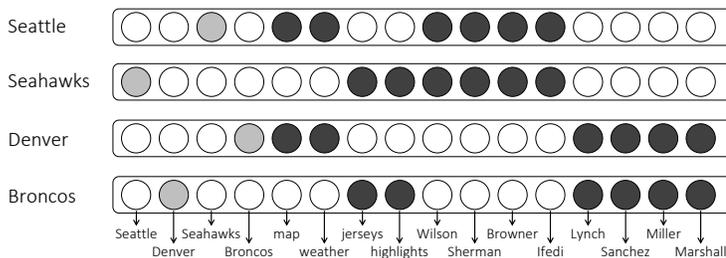
Explicit vector representations can be broadly categorized based on their choice of distributional features (e.g., in documents, neighbouring terms with or without distances, etc.) and different weighting schemes (e.g., TF-IDF, positive pointwise mutual information, etc.) applied over the raw counts. We direct the readers to [12, 199] which are good surveys of many existing explicit vector representation schemes.

Levy et al. [117] demonstrated that explicit vector representations are amenable to the term analogy task using simple vector operations. A term analogy task involves answering questions of the form “*man* is to *woman* as *king* is to \_\_\_\_?”—the correct answer to which in this case happens to be “queen”. In NLP, term analogies are typically performed by simple vector operations of the following form followed by a nearest-neighbour search,

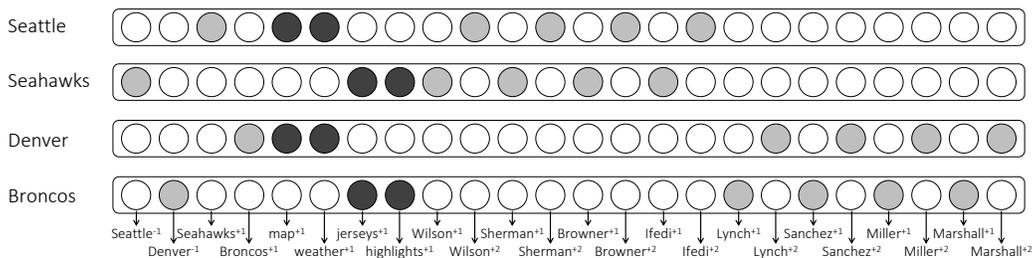
<sup>9</sup>Interestingly, the notion of Paradigmatic (Typical) and Syntagmatic (Topical) relationships show up almost universally—not just in text. In vision, for example, the different images of “noses” bear a Typical similarity to each other, while they share a Topical relationship with images of “eyes” or “ears”. Curiously, Barthes [13] even extended this analogy to garments—where paradigmatic relationships exist between items of the same type (e.g., between hats and between boots) and the proper Syntagmatic juxtaposition of items from these different Paradigms—from hats to boots— forms a fashionable ensemble .



(a) "In-documents" features



(b) "Neighbouring terms" features



(c) "Neighbouring terms w/ distances" features

Figure 9: The figure shows different distributed representations for the four terms—"Seattle", "Seahawks", "Denver", and "Broncos"—based on the toy corpus in Table 2. Shaded circles indicate non-zero values in the vectors—the darker shade highlights the vector dimensions where more than one vector has a non-zero value. When the representation is based on the documents that the terms occur in then "Seattle" is more similar to "Seahawks" than to "Denver". The representation scheme in (a) is, therefore, more aligned with a Topical notion of similarity. In contrast, in (c) each term is represented by a vector of neighbouring terms—where the distances between the terms are taken into consideration—which puts "Seattle" closer to "Denver" demonstrating a Typical, or type-based, similarity. When the inter-term distances are ignored, as in (b), a mix of Typical and Topical similarities is observed. Finally, it is worth noting that neighbouring-terms based vector representations leads to similarities between terms that do not necessarily occur in the same document, and hence the term-term relationships are less sparse than when only in-document features are considered.

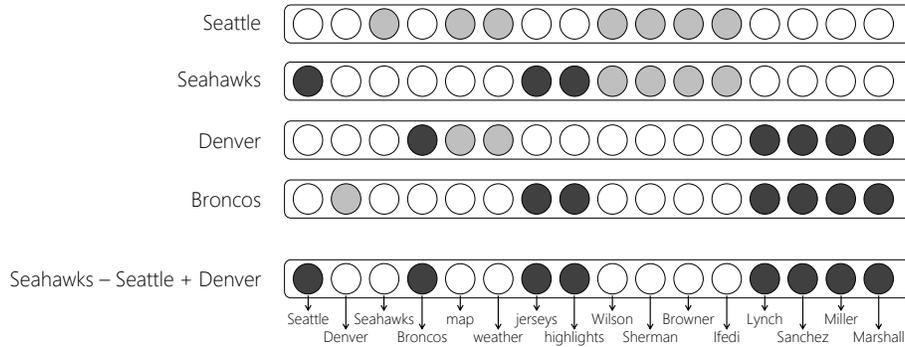


Figure 10: A visual demonstration of term analogies via simple vector algebra. The shaded circles denote non-zero values. Darker shade is used to highlight the non-zero values along the vector dimensions for which the output of  $\vec{v}_{Seahawks} - \vec{v}_{Seattle} + \vec{v}_{Denver}$  is positive. The output vector is closest to  $\vec{v}_{Broncos}$  as shown in this toy example.

$$\vec{v}_{king} - \vec{v}_{man} + \vec{v}_{woman} \approx \vec{v}_{queen} \quad (24)$$

It may be surprising to some readers that the vector obtained by the simple algebraic operations  $\vec{v}_{king} - \vec{v}_{man} + \vec{v}_{woman}$  produces a vector close to the vector  $\vec{v}_{queen}$ . We present a visual intuition of why this works in practice in Figure 10, but we refer the readers to [7, 117] for a more rigorous mathematical explanation.

#### 4.4 Embeddings

While explicit vector representations based on distributional features can capture interesting notions of term-term similarity they have one big drawback—the resultant vector spaces are highly sparse and high-dimensional. The number of dimensions is generally in the same order as the number of documents or the vocabulary size, which is unwieldy for most practical tasks. An alternative is to learn lower dimensional representations of terms from the data that retains similar attributes as the higher dimensional vectors.

An *embedding* is a representation of items in a new space such that the properties of, and the relationships between, the items are preserved. Goodfellow et al. [64] articulate that the goal of an embedding is to generate a *simpler* representation—where simplification may mean a reduction in the number of dimensions, an increase in the sparseness of the representation, disentangling the principle components of the vector space, or a combination of these goals. In the context of term embeddings, the explicit feature vectors—like those we discussed in Section 4.3—constitutes the original representation. An embedding trained from these features assimilate the properties of the terms and the inter-term relationships observable in the original feature space.

The most popular approaches for learning embeddings include either factorizing the term-feature matrix (e.g. LSA [48]) or using gradient descent based methods that try to predict the features given the term (e.g., [15, 134]). Baroni et al. [11] empirically demonstrate that these feature-predicting models that learn lower dimensional representations, in fact, also perform better than explicit counting based models on different tasks—possibly due to better generalization across terms—although some counter evidence the claim of better performances from embedding models have also been reported in the literature [116].

The sparse feature spaces of Section 4.3 are easier to visualize and leads to more intuitive explanations—while their corresponding embeddings are more practically useful. Therefore, it makes sense to *think sparse, but act dense* in many scenarios. In the rest of this section, we will describe some of the popular neural and non-neural embedding models.

**Latent Semantic Analysis (LSA)** LSA [48] involves performing *singular value decomposition* (SVD) [63] on a term-document (or term-passage) matrix  $X$  to obtain its low-rank approximation

[130]. SVD on  $X$  involves finding a solution to  $X = U\Sigma V^T$ , where  $U$  and  $V$  are orthogonal matrices and  $\Sigma$  is a diagonal matrix.<sup>10</sup>

$$\begin{array}{c}
 \begin{array}{c} X \\ (\vec{d}_j) \\ \downarrow \end{array} \\
 \begin{array}{c} \left[ \begin{array}{ccc} x_{1,1} & \dots & x_{1,|D|} \\ \vdots & \ddots & \vdots \\ x_{|T|,1} & \dots & x_{|T|,|D|} \end{array} \right] \\ (\vec{t}_i^T) \rightarrow \end{array}
 \end{array}
 =
 \begin{array}{c}
 U \\
 \begin{array}{c} \left[ \begin{array}{c} \vec{u}_1 \\ \vdots \\ \vec{u}_l \end{array} \right] \dots \left[ \begin{array}{c} \vec{u}_1 \\ \vdots \\ \vec{u}_l \end{array} \right] \end{array} \\
 \begin{array}{c} \left[ \begin{array}{ccc} \sigma_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_l \end{array} \right] \\ \Sigma \end{array} \\
 \begin{array}{c} V^T \\ (\vec{d}_j) \\ \downarrow \\ \left[ \begin{array}{c} \vec{v}_1 \\ \vdots \\ \vec{v}_l \end{array} \right] \end{array}
 \end{array}
 \end{array}
 \quad (25)$$

where,  $\sigma_1, \dots, \sigma_l, \vec{u}_1, \dots, \vec{u}_l$ , and  $\vec{v}_1, \dots, \vec{v}_l$  are the singular values, the left singular vectors, and the right singular vectors, respectively. The  $k$  largest singular values, and corresponding singular vectors from  $U$  and  $V$ , is the rank  $k$  approximation of  $X$  ( $X_k = U_k \Sigma_k V_k^T$ ). The embedding for the  $i_{th}$  term is given by  $\Sigma_k \vec{t}_i$ .

While LSA operate on a term-document matrix, matrix factorization based approaches can also be applied to term-term matrices [25, 111, 168].

*Neural term embedding* models are typically trained by setting up a prediction task. Instead of factorizing the term-feature matrix—as in LSA—neural models are trained to predict the term from its features. Both the term and the features have *one-hot* representations in the input and the output layers, respectively, and the model learns dense low-dimensional representations in the process of minimizing the prediction error. These approaches are based on the *information bottleneck method* [197]—discussed in more details in Section 6.2—with the low-dimensional representations acting as the bottleneck. The training data may contain many instances of the same term-feature pair proportional to their frequency in the corpus (e.g., word2vec [134]), or their counts can be pre-aggregated (e.g., GloVe [160]).

**Word2vec** For word2vec [61, 134, 136, 137, 169], the features for a term are made up of its neighbours within a fixed size window over the text from the training corpus. The *skip-gram* architecture (see Figure 11a) is a simple one hidden layer neural network. Both the input and the output of the model is in the form of one-hot vectors and the loss function is as follows,

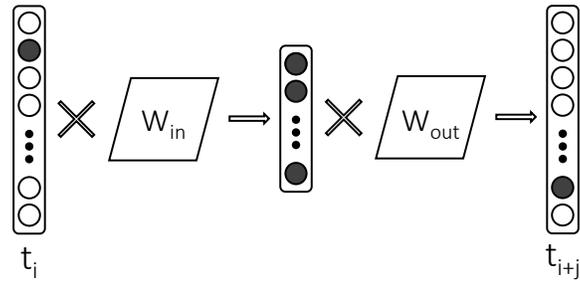
$$\mathcal{L}_{skip-gram} = -\frac{1}{|S|} \sum_{i=1}^{|S|} \sum_{-c \leq j \leq +c, j \neq 0} \log(p(t_{i+j}|t_i)) \quad (26)$$

$$\text{where, } p(t_{i+j}|t_i) = \frac{\exp((W_{out} \vec{v}_{t_{i+j}})^T (W_{in} \vec{v}_{t_i}))}{\sum_{k=1}^{|T|} \exp((W_{out} \vec{v}_{t_k})^T (W_{in} \vec{v}_{t_i}))} \quad (27)$$

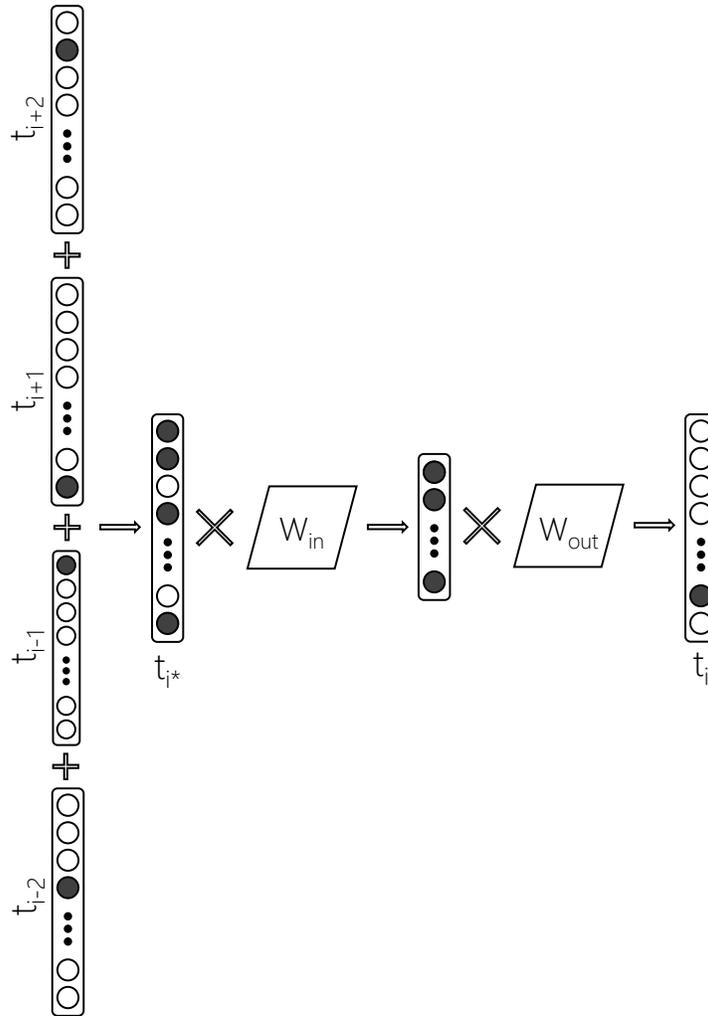
$S$  is the set of all windows over the training text and  $c$  is the number of neighbours we need to predict on either side of the term  $t_i$ . The denominator for the softmax function for computing  $p(t_{i+j}|t_i)$  sums over all the words in the vocabulary. This is prohibitively costly and in practice either hierarchical-softmax [149] or negative sampling is employed. Also, note that the model has two different weight matrices  $W_{in}$  and  $W_{out}$  that are learnable parameters of the models.  $W_{in}$  gives us the IN embeddings corresponding to all the input terms and  $W_{out}$  corresponding to the OUT embeddings for the output terms. Generally, only  $W_{in}$  is used and  $W_{out}$  is discarded after training, but we will discuss an IR application that makes use of both the IN and the OUT embeddings later in Section 5.1.

The *continuous bag-of-words* (CBOW) architecture (see Figure 11b) is similar to the skip-gram model, except that the task is to predict the middle term given the sum of the one-hot vectors of the neighbouring terms in the window. Given a middle term  $t_i$  and the set of its neighbours

<sup>10</sup>The matrix visualization is adapted from [https://en.wikipedia.org/wiki/Latent\\_semantic\\_analysis](https://en.wikipedia.org/wiki/Latent_semantic_analysis).



(a) Skip-gram



(b) Continuous bag-of-words (CBOW)

Figure 11: The (a) skip-gram and the (b) continuous bag-of-words (CBOW) architectures of word2vec. The architecture is a neural network with a single hidden layer whose size is much smaller than that of the input and the output layers. Both models use one-hot representations of terms in the input and the output. The learnable parameters of the model comprise of the two weight matrices  $W_{in}$  and  $W_{out}$  that corresponds to the embeddings the model learns for the input and the output terms, respectively. The skip-gram model trains by minimizing the error in predicting a term given one of its neighbours. The CBOW model, in contrast, predicts a term from a bag of its neighbouring terms.

$\{t_{i-c}, \dots, t_{i-1}, t_{i+1}, \dots, t_{i+c}\}$ , the CBOW model creates a single training sample with the sum of the one-hot vectors of all the neighbouring terms as input and the one-hot vector  $\vec{v}_{t_i}$ , corresponding to the middle term, as the expected output.

$$\mathcal{L}_{CBOW} = -\frac{1}{|S|} \sum_{i=1}^{|S|} \log(p(t_i | \sum_{-c \leq j \leq +c, j \neq 0} t_{i+j})) \quad (28)$$

Contrast this with the skip-gram model that creates  $2 \times c$  samples by individually pairing each of the neighbouring terms with the middle term. During training, given the same number of windows of text, the skip-gram model, therefore, trains orders of magnitude slower than the CBOW model [134] because it creates  $2 \times c$  the number of training samples.

Word2vec gained particular popularity for its ability to perform word analogies using simple vector algebra, similar to what we have already discussed in Section 4.3. For domains where the interpretability of the embeddings may be important, Sun et al. [191] introduced an additional constraint in the loss function to encourage more sparseness in the learnt representation.

$$\mathcal{L}_{sparse-CBOW} = \mathcal{L}_{CBOW} - \lambda \sum_{t \in T} \|\vec{v}_t\|_1 \quad (29)$$

**GloVe** The skip-gram model trains on individual term-neighbour pairs. If we aggregate all the training samples into a matrix  $X$ , such that  $x_{ij}$  is the frequency of the pair  $\langle t_i, t_j \rangle$  in the training data, then the loss function changes to,

$$\mathcal{L}_{skip-gram} = - \sum_{i=1}^{|T|} \sum_{j=1}^{|T|} x_{ij} \log(p(t_i | t_j)) \quad (30)$$

$$= - \sum_{i=1}^{|T|} x_i \sum_{j=1}^{|T|} \frac{x_{ij}}{x_i} \log(p(t_i | t_j)) \quad (31)$$

$$= - \sum_{i=1}^{|T|} x_i \sum_{j=1}^{|T|} \bar{p}(t_i | t_j) \log(p(t_i | t_j)) \quad (32)$$

$$= \sum_{i=1}^{|T|} x_i H(\bar{p}(t_i | t_j) \log(p(t_i | t_j))) \quad (33)$$

$H(\dots)$  is the cross-entropy error between the actual co-occurrence probability  $\bar{p}(t_i | t_j)$  and the one predicted by the model  $p(t_i | t_j)$ . This is similar to the loss function for GloVe [160] if we replace the cross-entropy error with a squared-error and apply a saturation function  $f(\dots)$  over the actual co-occurrence frequencies.

$$\mathcal{L}_{GloVe} = - \sum_{i=1}^{|T|} \sum_{j=1}^{|T|} f(x_{ij}) (\log(x_{ij}) - \vec{w}_i^\top \vec{w}_j)^2 \quad (34)$$

GloVe is trained using AdaGrad [53]. Similar to word2vec, GloVe also generates two different (IN and OUT) embeddings, but unlike word2vec it generally uses the sum of the IN and the OUT vectors as the embedding for each term in the vocabulary.

**Paragraph2vec** Following the popularity of word2vec [134, 136], similar neural architectures [4, 5, 68, 69, 110, 189] have been proposed that trains on term-document co-occurrences. The training typically involves predicting a term given the ID of a document or a passage that contains the term. In some variants, as shown in Figure 12, neighbouring terms are also provided as input.

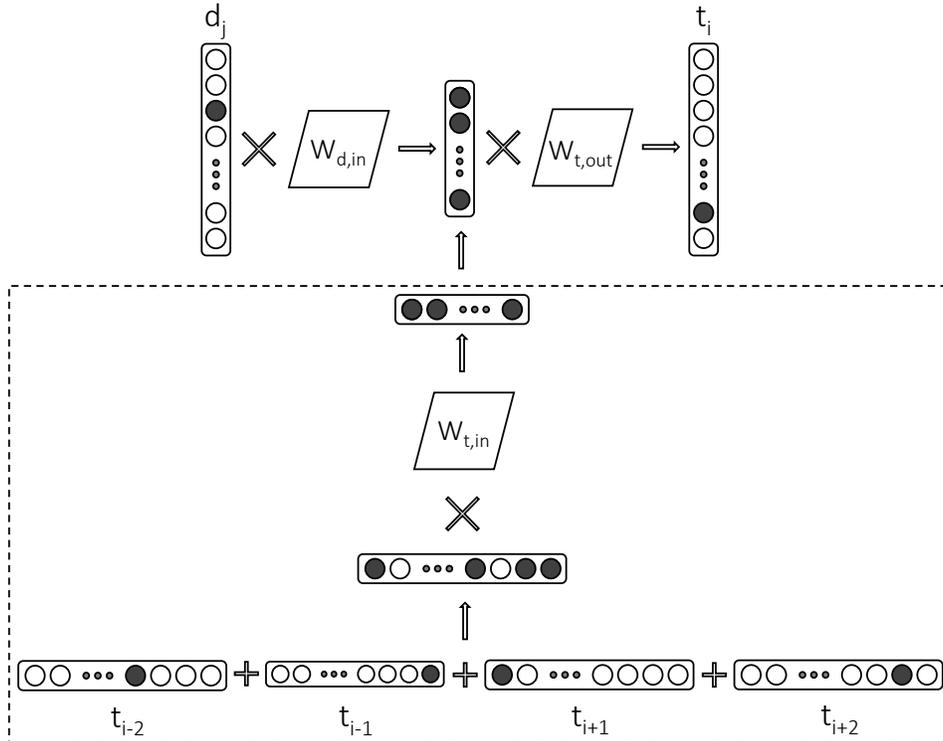


Figure 12: The paragraph2vec architecture as proposed by Le and Mikolov [110] trains by predicting a term given a document (or passage) ID containing the term. By trying to minimize the prediction error, the model learns an embedding for the term as well as for the document. In some variants of the architecture, optionally the neighbouring terms are also provided as input—as shown in the dotted box.

The key motivation for training on term-document pairs is to learn an embedding that is more aligned with a Topical notion of term-term similarity—which is often more appropriate for IR tasks. The term-document relationship, however, tends to be more sparse [219]—including neighbouring term features may compensate for some of that sparsity.

In the context of IR tasks, Ai et al. [4, 5] proposed a number of IR-motivated changes to the original Paragraph2vec [110] model training—including, document frequency based negative sampling and document length based regularization.

## 5 Term embeddings for IR

Traditional IR models use local representations of terms for query-document matching. The most straight-forward use case for term embeddings in IR is to enable *inexact* matching in the embedding space. In Section 2.2, we argued the importance of inspecting non-query terms in the document for garnering evidence of relevance. For example, even from a shallow manual inspection, it is easy to conclude that the passage in Figure 13a is *about* Albuquerque because it contains “metropolitan”, “population”, and “area” among other informative terms. On the other hand, the passage in Figure 13b contains “simulator”, “interpreter”, and “Altair” which seems to suggest that the passage is instead more likely related to computers and technology. In traditional term counting based IR approaches these signals are often ignored.

Albuquerque is the most populous city in the U.S. state of New Mexico. The high-altitude city serves as the county seat of Bernalillo County, and it is situated in the central part of the state, straddling the Rio Grande. The city population is 557,169 as of the July 1, 2014 population estimate from the United States Census Bureau, and ranks as the 32nd-largest city in the U.S. The Albuquerque metropolitan statistical area (or MSA) has a population of 907,301 according to the United States Census Bureau’s most recently available estimate for 2015.

Allen suggested that they could program a BASIC interpreter for the device; after a call from Gates claiming to have a working interpreter, MITS requested a demonstration. Since they didn’t actually have one, Allen worked on a simulator for the Altair while Gates developed the interpreter. Although they developed the interpreter on a simulator and not the actual device, the interpreter worked flawlessly when they demonstrated the interpreter to MITS in Albuquerque, New Mexico in March 1975; MITS agreed to distribute it, marketing it as Altair BASIC.

(a) About Albuquerque

(b) Not about Albuquerque

Figure 13: Two passages both containing exactly a single occurrence of the query term “Albuquerque”. However, the passage in (a) contains other terms such as “population” and “area” that are relevant to a description of the city. In contrast, the terms in passage (b) suggest that it is unlikely to be about the city, and only mentions the city potentially in a different context.

Most existing shallow neural methods for IR focus on inexact matching using term embeddings. These approaches can be broadly categorized as those that compare the query with the document directly in the embedding space; and those that use embeddings to generate suitable query expansion candidates from a global vocabulary and then perform retrieval based on the expanded query. We discuss both these classes of approaches in the remainder of this section.

### 5.1 Query-document matching

A popular strategy for using term embeddings in IR involves deriving a dense vector representation for the query and the document from the embeddings of the individual terms in the corresponding texts. The term embeddings can be aggregated in different ways, although using the *average word (or term) embeddings* (AWE) is quite popular [96, 101, 110, 143, 151, 190, 207]. Non-linear combinations of term vectors—such as using Fisher Kernel Framework [35]—have also been explored, as well as other families of aggregate functions of which AWE has been shown to be a special case [228].

The query and the document embeddings themselves can be compared using a variety of similarity metrics, such as cosine similarity or dot-product. For example,

$$sim(q, d) = cos(\vec{v}_q, \vec{v}_d) = \frac{\vec{v}_q^\top \vec{v}_d}{\|\vec{v}_q\| \|\vec{v}_d\|} \tag{35}$$

$$\text{where, } \vec{v}_q = \frac{1}{|q|} \sum_{t_q \in q} \frac{\vec{v}_{t_q}}{\|\vec{v}_{t_q}\|} \tag{36}$$

$$\vec{v}_d = \frac{1}{|d|} \sum_{t_d \in d} \frac{\vec{v}_{t_d}}{\|\vec{v}_{t_d}\|} \tag{37}$$

An important consideration here is the choice of the term embeddings that is appropriate for the retrieval scenario. While, LSA [48], word2vec [136], and GloVe [160] are popularly used—it is important to understand how the notion of inter-term similarity modelled by a specific vector space may influence its performance on a retrieval task. In the example in Figure 13, we want to rank documents that contains related terms, such as “population” or “area” higher—these terms are Topically similar to the query term “Albuquerque”. Intuitively, a document about “Tucson”—which is Typically similar to “Albuquerque”—is unlikely to satisfy the user intent. The discussion in Section 4.2 on how input features influence the notion of similarity in the learnt vector space is relevant here.

Models, such as LSA [48] and Paragraph2vec [110], that consider term-document pairs generally capture Topical similarities in the learnt vector space. On the other hand, word2vec [136] and GloVe [160] embeddings may incorporate a mixture of Topical and Typical notions of relatedness. These

neural models behave more Typical when trained with short window sizes or on short text, such as on keyword queries [115] (refer to Section 4.2 for more details).

In Section 4.4, we made a note that the word2vec model learns two different embeddings—IN and OUT—corresponding to the input and the output terms. Mitra et al. [143] point out that when using word2vec embeddings for IR it is more appropriate to represent the query terms using the IN embeddings and the document terms using the OUT embeddings of the trained model. In this *Dual Embedding Space Model* (DESM)<sup>11</sup> [143, 151] the word2vec embeddings are trained on search queries, which empirically performs better than training on document body text. Training on short queries, however, makes the inter-term similarity more pronouncedly Typical (where, “Yale” is closer to “Harvard” and “NYU”) when both terms are represented using their IN vectors—better retrieval performance is achieved instead by using the IN-OUT similarity (where, “Yale” is closer to “faculty” and “alumni”) that mirrors more the Topical notions of relatedness.

$$DESM_{in-out}(q, d) = \frac{1}{|q|} \sum_{t_q \in q} \frac{\vec{v}_{t_q, in}^\top \vec{v}_{d, out}}{\|\vec{v}_{t_q, in}\| \|\vec{v}_{d, out}\|} \quad (38)$$

$$\vec{v}_{d, out} = \frac{1}{|d|} \sum_{t_d \in d} \frac{\vec{v}_{t_d, out}}{\|\vec{v}_{t_d, out}\|} \quad (39)$$

An alternative to representing queries and documents as an aggregate of their term embeddings is to incorporate the term representations into existing IR models, such as the ones we discussed in Section 2.5. Zuccon et al. [236] proposed the *Neural Translation Language Model* (NTLM) that uses the similarity between term embeddings as a measure for term-term translation probability  $p(t_q|t_d)$  in Equation 15.

$$p(t_q|t_d) = \frac{\cos(\vec{v}_{t_q}, \vec{v}_{t_d})}{\sum_{t \in T} \cos(\vec{v}_t, \vec{v}_{t_d})} \quad (40)$$

On similar lines, Ganguly et al. [58] proposed the *Generalized Language Model* (GLM) which extends the Language Model based approach in Equation 13 to,

$$p(d|q) = \prod_{t_q \in q} \left( \lambda \frac{tf(t_q, d)}{|d|} + \alpha \frac{\sum_{t_d \in d} (sim(\vec{v}_{t_q}, \vec{v}_{t_d}) \cdot tf(t_d, d))}{\sum_{t_{d_1} \in d} \sum_{t_{d_2} \in d} sim(\vec{v}_{t_{d_1}}, \vec{v}_{t_{d_2}}) \cdot |d|^2} \right. \\ \left. + \beta \frac{\sum_{\bar{t} \in N_t} (sim(\vec{v}_{t_q}, \vec{v}_{\bar{t}}) \cdot \sum_{\bar{d} \in D} tf(\bar{t}, \bar{d}))}{\sum_{t_{d_1} \in N_t} \sum_{t_{d_2} \in N_t} sim(\vec{v}_{t_{d_1}}, \vec{v}_{t_{d_2}}) \cdot \sum_{\bar{d} \in D} |\bar{d}| \cdot |N_t|} + (1 - \alpha - \beta - \lambda) \frac{\sum_{\bar{d} \in D} tf(t_q, \bar{d})}{\sum_{\bar{d} \in D} |\bar{d}|} \right) \quad (41)$$

where,  $N_t$  is the set of nearest-neighbours of term  $t$ . Ai et al. [4] incorporate paragraph vectors [110] into the query-likelihood model [161].

Another approach, based on the Earth Mover’s Distance (EMD) [171], involves estimating similarity between pairs of documents by computing the minimum distance in the embedding space that each term in the first document needs to travel to reach the terms in the second document. This measure, commonly referred to as the *Word Mover’s Distance* (WMD), was originally proposed by Wan et al. [210, 211], but used WordNet and topic categories instead of distributed representations for defining distance between terms. Term embeddings were later incorporated into the model by Kusner et al. [87, 104]. Finally, Guo et al. [72] incorporated similar notion of distance into the *Non-linear Word Transportation* (NWT) model that estimates relevance between a query and a document. The NWT model involves solving the following constrained optimization problem,

<sup>11</sup>The dual term embeddings trained on Bing queries is available for download at <https://www.microsoft.com/en-us/download/details.aspx?id=52597>

$$\max \sum_{t_q \in q} \log \left( \sum_{t_d \in u(d)} f(t_q, t_d) \cdot \max(\cos(\vec{v}_{t_q}, \vec{v}_{t_d}), 0)^{idf(t_q)+b} \right) \quad (42)$$

$$\text{subject to } f(t_q, t_d) \geq 0, \quad \forall t_q \in q, t_d \in d \quad (43)$$

$$\text{and } \sum_{t_q \in q} f(t_q, t_d) = \frac{tf(t_d) + \mu \frac{\sum_{\bar{d} \in D} tf(t_q, \bar{d})}{\sum_{\bar{d} \in D} |\bar{d}|}}{|d| + \mu}, \quad \forall t_d \in d \quad (44)$$

$$\text{where, } idf(t) = \frac{|D| - df(t) + 0.5}{df(t) + 0.5} \quad (45)$$

$u(d)$  is the set of all unique terms in document  $d$ , and  $b$  is a constant.

Another term-alignment based distance metric was proposed by Kenter and de Rijke [98] for computing short-text similarity. The design of the *saliency-weighted semantic network* (SWSN) is motivated by the BM25 [166] formulation.

$$swsn(s_l, s_s) = \sum_{t_i \in s_l} idf(t_i) \cdot \frac{sem(t_i, s_s) \cdot (k_1 + 1)}{sem(t_i, s_s) + k_1 \cdot \left(1 - b + b \cdot \frac{|s_s|}{avgsl}\right)} \quad (46)$$

$$\text{where, } sem(t, s) = \max_{\vec{t} \in s} \cos(\vec{v}_t, \vec{v}_{\vec{t}}) \quad (47)$$

Here  $s_s$  is the shorter of the two sentences to be compared, and  $s_l$  the longer sentence.

**Telescoping evaluation** Figure 14 highlights the distinct strengths and weaknesses of matching using local and distributed representations of terms for retrieval. For the query ‘‘Cambridge’’, a local representation (or exact matching) based model can easily distinguish between the passage on Cambridge (Figure 14a) and the one on Oxford (Figure 14b). However, the model is easily duped by an non-relevant passage that has been artificially injected with the term ‘‘Cambridge’’ (Figure 14d). The distributed representation based matching, on the other hand, can spot that the other terms in the passage provide clear indication that the passage is not *about* a city, but fails to realize that the passage about Oxford (Figure 14b) is inappropriate for the same query.

Embedding based models often perform poorly when the retrieval is performed over the full document collection [143]. However, as seen in the example of Figure 14, the errors made by embedding based models and exact matching models are typically different—and the combination of the two performs better than exact matching models alone [4, 58, 143]. Another popular technique is to use the embedding based model to re-rank only a subset of the documents retrieved by a different—generally an exact matching based—IR model. The chaining of different IR models where each successive model re-ranks a smaller number of candidate documents is called *Telescoping* [131]. Telescoping evaluations are popular in the neural IR literature [71, 88, 141, 143, 177] and the results are representative of performances of these models on re-ranking tasks. However, as Mitra et al. [143] demonstrate, good performances on re-ranking tasks may not be indicative how the model would perform if the retrieval involves larger document collections.

## 5.2 Query expansion

Instead of comparing the query and the document directly in the embedding space, an alternative approach is to use term embeddings to find good expansion candidates from a global vocabulary, and then retrieving documents using the expanded query. Different functions [51, 170, 227] have been proposed for estimating the relevance of candidate terms to the query—all of them involves comparing the candidate term individually to every query term using their vector representations, and then aggregating the scores. For example, [51, 170] estimate the relevance of candidate term  $t_c$  as,

$$score(t_c, q) = \frac{1}{|q|} \sum_{t_q \in q} \cos(\vec{v}_{t_c}, \vec{v}_{t_q}) \quad (48)$$

the city of **cambridge** is a university city and the county town of cambridgeshire , england . it lies in east anglia , on the river cam , about 50 miles ( 80 km ) north of london . according to the united kingdom census 2011 , its population was ( including students ) . this makes **cambridge** the second largest city in cambridgeshire after peterborough , and the 54th largest in the united kingdom . there is archaeological evidence of settlement in the area during the bronze age and roman times ; under viking rule **cambridge** became an important trading centre . the first town charters were granted in the 12th century , although city status was not conferred until 1951 .

(a) Passage about the city of Cambridge

oxford is a city in the south east region of england and the county town of oxfordshire . with a population of . it is the 52nd largest city in the united kingdom , and one of the fastest growing and most ethnically diverse . oxford has a broad economic base . its industries include motor manufacturing , education , publishing and a large number of information technology and businesses , some being academic offshoots . the city is known worldwide as the home of the university of oxford , the oldest university in the world . buildings in oxford demonstrate examples of every english architectural period since the arrival of the saxons , including the radcliffe camera . oxford is known as the city of dreaming spires , a term coined by poet matthew arnold .

(b) Passage about the city of Oxford

the giraffe ( giraffa camelopardalis ) is an african ungulate mammal , the tallest living terrestrial animal and the largest ruminant . its species name refers to its shape and its colouring . its chief distinguishing characteristics are its extremely long neck and legs , its , and its distinctive coat patterns . it is classified under the family , along with its closest extant relative , the okapi . the nine subspecies are distinguished by their coat patterns . the scattered range of giraffes extends from chad in the north to south africa in the south , and from niger in the west to somalia in the east . giraffes usually inhabit savannas , grasslands , and open woodlands .

(c) Passage about giraffes

the **cambridge** ( giraffa camelopardalis ) is an african ungulate mammal , the tallest living terrestrial animal and the largest ruminant . its species name refers to its shape and its colouring . its chief distinguishing characteristics are its extremely long neck and legs , its , and its distinctive coat patterns . it is classified under the family , along with its closest extant relative , the okapi . the nine subspecies are distinguished by their coat patterns . the scattered range of giraffes extends from chad in the north to south africa in the south , and from niger in the west to somalia in the east . giraffes usually inhabit savannas , grasslands , and open woodlands .

(d) Passage about giraffes, but 'giraffe' is replaced by 'Cambridge'

Figure 14: A visualization of IN-OUT similarities between terms in different passages with the query term “Cambridge”. The visualization—adapted from <https://github.com/bmitra-msft/Demos/blob/master/notebooks/DESM.ipynb>—reveal that, besides the term “Cambridge”, many other terms in the passages about both Cambridge and Oxford have high similarity to the query term. The passage (d) is adapted from the passage (c) on giraffes by replacing all the occurrences of the term “giraffe” with “cambridge”. However, none of the other terms in (d) are found to be relevant to the query term. An embedding based approach may be able to determine that passage (d) is non-relevant to the query “Cambridge”, but fail to realize that passage (b) is also non-relevant. A term counting-based model, on the other hand, can easily identify that passage (b) is non-relevant, but may rank passage (d) incorrectly high.

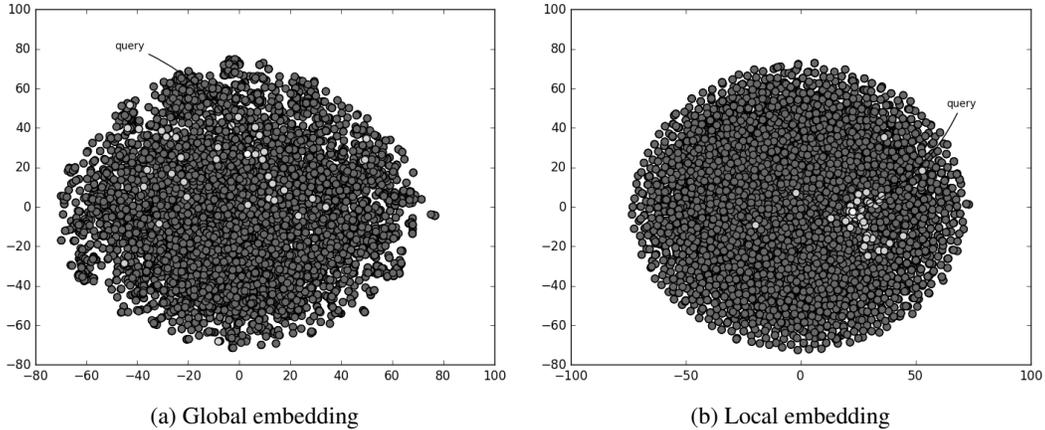


Figure 15: A two-dimensional visualization of term embeddings when the vector space is trained on a (a) global corpus and a (b) query-specific corpus, respectively. The grey circles represent individual terms in the vocabulary. The white circle represents the query “ocean remote sensing” by averaging the embeddings of the individual terms in the query, and the light grey circles correspond to good expansion terms for this query. When the representations are query-specific then the meaning of the terms are better disambiguated, and more likely to result in the selection of good expansion terms.

Term embedding based query expansion on its own performs worse than pseudo-relevance feedback [170]. But like the models in the previous section, shows better performances when used in combination with PRF [227].

Diaz et al. [51] explored the idea of query-specific term embeddings and found that they are much more effective in identifying good expansion terms than a global representation (see Figure 15). The local model proposed by Diaz et al. [51] incorporate relevance feedback in the process of learning the term embeddings—a set of documents is retrieved for the query and a query-specific term embedding model is trained. This *local* embedding model is then employed for identifying expansion candidates for the query for a second round of document retrieval.

Term embeddings have also been explored for re-weighting query terms [233] and finding relevant query re-writes [69], as well as in the context of other IR tasks such as cross-lingual retrieval [207] and entity retrieval [200]. In the next section, we move on to neural network models with deeper architectures and their applications to retrieval.

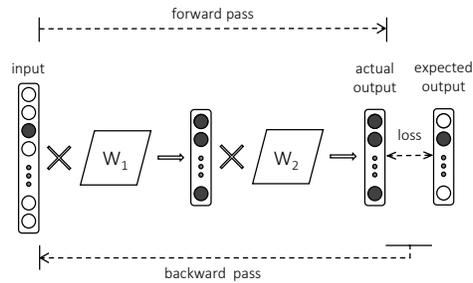
## 6 Deep neural networks

Deep neural network models consist of chains of tensor operations. The tensor operation can range from parameterized linear transformations (e.g., multiplication with a weight matrix, addition of a bias vector) to elementwise application of non-linear functions, such as *tanh* or *rectified linear units* (ReLU) [73, 89, 150]. Figure 16 shows a simple *feed-forward* neural network with *fully-connected* layers. For an input vector  $\vec{x}$ , the model produces the output  $\vec{y}$  as follows,

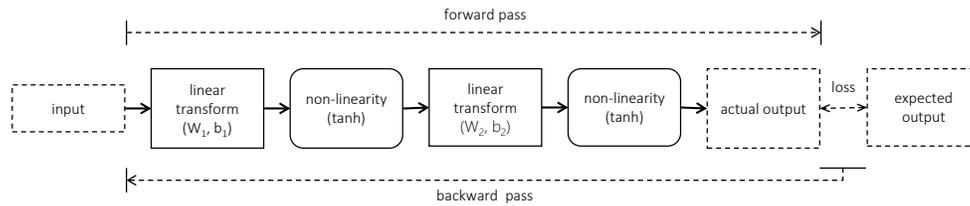
$$\vec{y} = \tanh(W_2 \cdot \tanh(W_1 \cdot \vec{x} + \vec{b}_1) + \vec{b}_2) \quad (49)$$

The model training involves tuning the parameters  $W_1$ ,  $\vec{b}_1$ ,  $W_2$ , and  $\vec{b}_2$  to minimize the loss between the expected output and the actual output of the final layer. The parameters are usually trained discriminatively using backpropagation [14, 77, 175]. During forward-pass each layer generates an output conditioned on its input, and during backward pass each layer computes the error gradient with respect to its parameters and its inputs.

The design of a DNN typically involves many choices of architectures and hyper-parameters. Neural networks with as few a single hidden layer—but with sufficient number of hidden nodes—can



(a) A neural network with a single hidden layer.



(b) The same neural network viewed as a chain of computational steps.

Figure 16: Two different visualizations of a feed-forward neural network with a single hidden layer. In (a), the addition of the bias vector and the non-linearity function is implicit. Figure (b) shows the same network but as a sequence of computational nodes. Most popular neural network toolkits implement a set of standard computational nodes that can be connected to build more sophisticated neural architectures.

theoretically approximate any function [85]. In practice, however, deeper architectures—sometimes with as many as 1000 layers [76]—have been shown to perform significantly better than shallower networks. For readers who are less familiar with neural network models, we present a simple example in Figure 17 to illustrate how hidden layers enable these models to capture non-linear relationships. We direct readers to [148] for further discussions on how additional hidden layers help.

The rest of this section is dedicated to the discussion of input representations and popular architectures for deep neural models.

## 6.1 Input text representations

Neural models that learn representations of text take raw text as input. A key consideration is how the text should be represented at the input layer of the model. Figure 18 shows some of the popular input representations of text.

Some neural models [66, 94, 100, 192] operate at the character-level. In these models, each character is typically represented by a one-hot vector. The vector dimensions—referred to as *channels*—in this case equals the number of allowed characters in the vocabulary. These models incorporate the least amount of prior knowledge about the language in the input representation—for example, these models are often required to learn about tokenization from scratch by treating space as just another character in the vocabulary. The representation of longer texts, such as sentences, can be derived by concatenating or summing the character-level vectors as shown in Figure 18a.

The input text can also be pre-tokenized into terms—where each term is represented by either a sparse vector or using pre-trained term embeddings (Figure 18d). Terms may have a one-hot (or local) representation where each term has a unique ID (Figure 18b), or the term vector can be derived by aggregating one-hot vectors of its constituting characters (or character  $n$ -graphs) as shown in Figure 18c. If pre-trained embeddings are used for term representation, then the embedding vectors can be further tuned during training, or kept fixed.

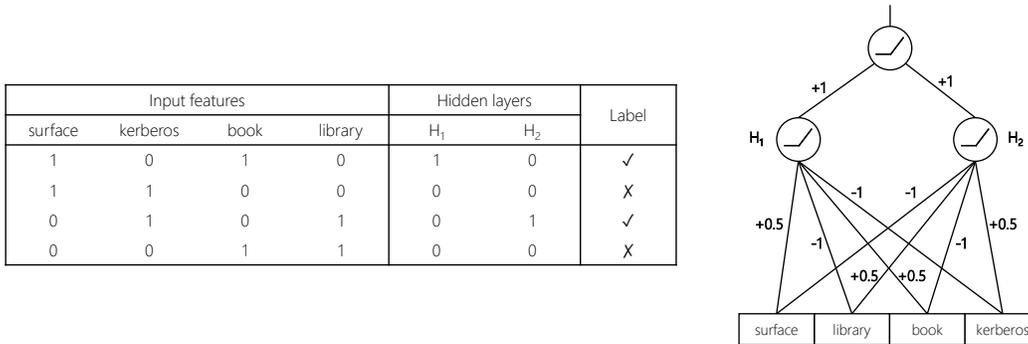


Figure 17: Consider a toy binary classification task on a corpus of four short texts—“surface book”, “kerberos library”, “library book”, and “kerberos surface”—where the model needs to predict if the text is related to computers. The first two texts—“Surface Book” and “kerberos library”—are positive under this classification, and the latter two negative. The input feature space consists of four binary features that indicate whether each of the four terms from the vocabulary is present in the text. The table shows that the specified classes are not linearly separable with respect to the input feature space. However, if we add couple of hidden nodes, as shown in the diagram, then the classes can be linearly separated with respect to the output of the hidden layer.

Similar to character-level models, the term vectors are further aggregated (by concatenation or sum) to obtain the representation of longer chunks of text, such as sentences. While one-hot representations of terms (Figure 18b) are common in many NLP tasks, pre-trained embeddings (e.g., [86, 158]) and character  $n$ -graph based representations (e.g., [88, 141]) are more popularly employed in IR.

## 6.2 Popular architectures

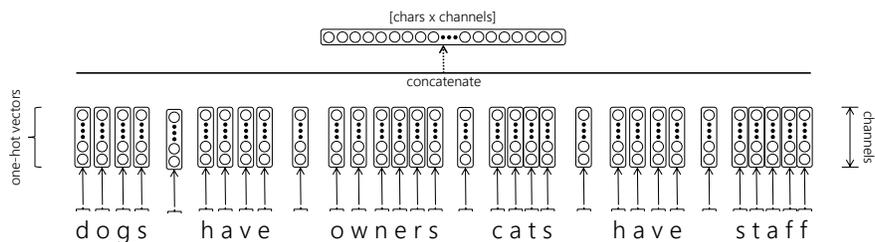
In this section, we describe few neural operations and architectures popular in IR. For broader overview of different neural architectures and design patterns please refer to [64, 112, 175].

**Shift-invariant neural operations** Convolutional [89, 103, 113, 114] and recurrent [67, 83, 135, 173] architectures are commonplace in most deep learning applications. These neural operations are part of a broader family of shift-invariant architectures. The key intuition behind these architectures stem from the natural regularities observable in most inputs. In vision, for example, the task of detecting a face should be invariant to whether the image is shifted, rotated, or scaled. Similarly, the meaning of an English sentence should, in most cases, stay consistent independent of which part of the document it appears in. Therefore, intuitively a neural model for object recognition or text understanding should not learn an independent logic for the same action applied to different parts of the input space.

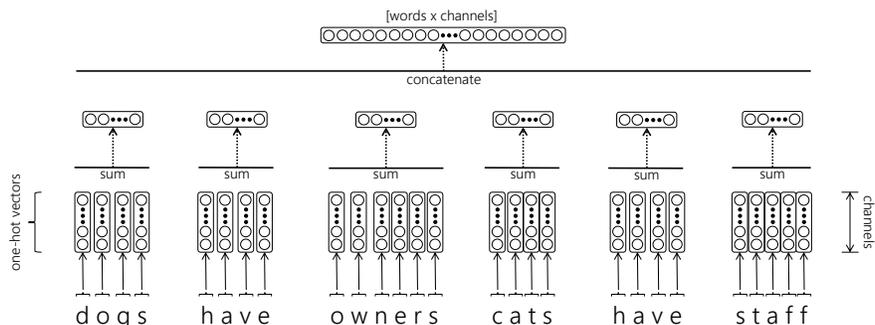
All shift-invariant neural operations fundamentally employ a window-based approach. A fixed size window is moved over the input space with fixed stride in each step. A (typically parameterized) function—referred to as a *kernel*, or a *filter*, or a *cell*—is applied over each instance of the window. The parameters of the cell are shared across all the instances of the input window. The shared parameters not only implies less number of total parameters in the model,, but also more supervision per parameter per training sample due to the repeated application.

Figure 19a shows an example of a cell being applied on a sequence of terms—with a window size of three terms—in each step. A popular cell implementation involves multiplying with a weight matrix—in which case the architecture in Figure 19a is referred as *convolutional*. An example of a cell without any parameters is *pooling*—which consists of aggregating (e.g., by computing the max or the average) over all the terms in the window<sup>12</sup>. Note, that the length of the input sequence can be variable in both cases and the length of the output of a convolutional (or pooling) layer is a function of the input length. Figure 19b shows an example of *global pooling*—where the window spans over the

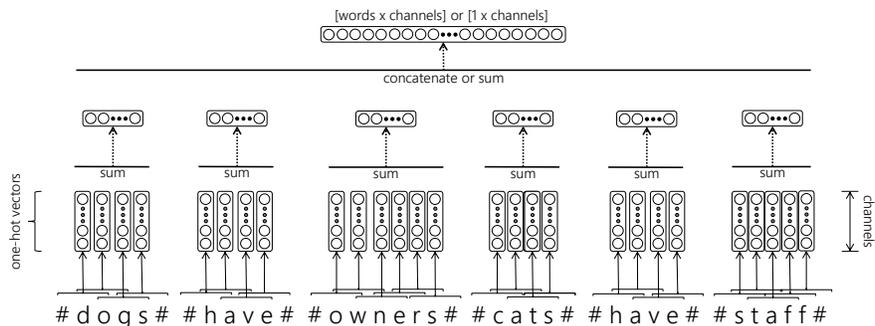
<sup>12</sup>If the input has multiple channels per term then the aggregation is performed per channel.



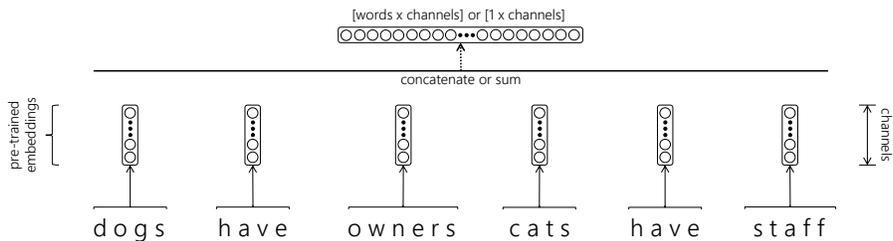
(a) Character-level input



(b) Term-level input w/ bag-of-characters per term



(c) Term-level input w/ bag-of-trigraphs per term



(d) Term-level input w/ pre-trained term embeddings

Figure 18: Examples of different representation strategies for text input to deep neural network models. The smallest granularity of representation can be a character or a term. The vector can be a sparse local representation, or a pre-trained embedding.

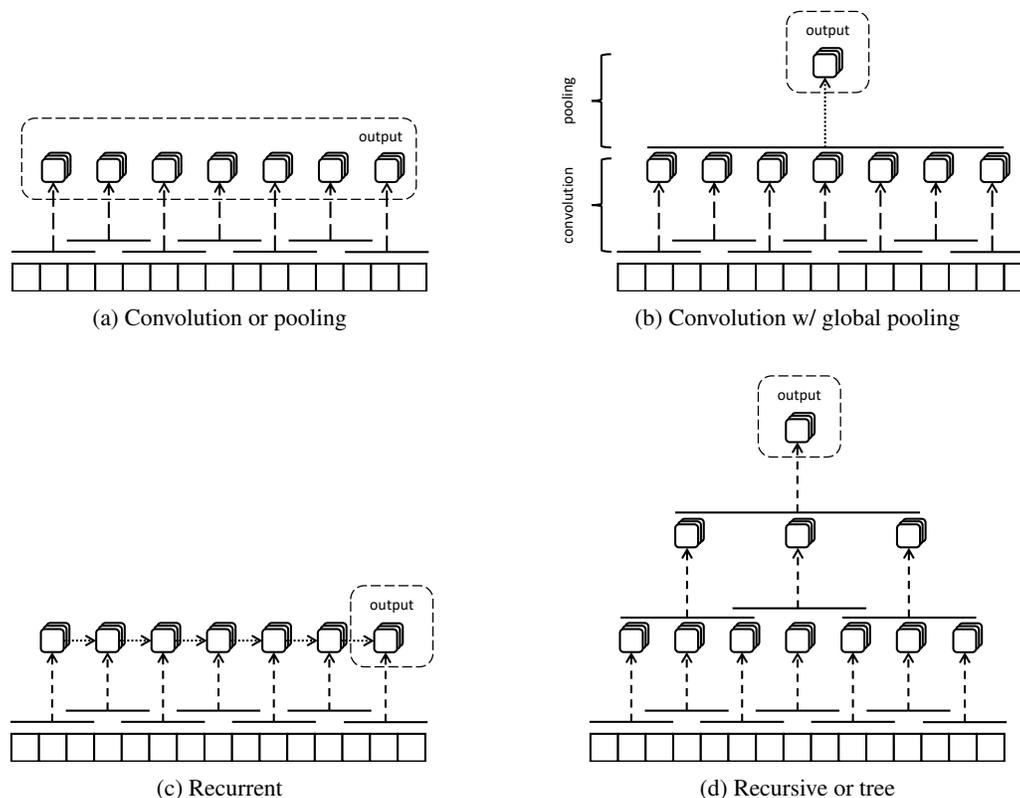


Figure 19: Popular shift-invariant neural architectures including convolutional neural networks (CNN), recurrent neural networks (RNN), pooling layers, and tree-structured neural networks.

whole input—being applied on top of a convolutional layer. The global pooling strategy is common for generating a fixed size output from a variable length input.<sup>13</sup>

In convolution or pooling, each window is applied independently. In contrast, in the *recurrent* architecture of Figure 19c the cell not only considers the input window but also the output of the previous instance of the cell as its input. Many different cell architectures have been explored for recurrent neural networks (RNN)—although Elman network [54], Long Short-Term Memory (LSTM) [83], and Gated Recurrent Unit (GRU) [32, 34] are popular. RNNs are popularly applied to sequences, but can also be useful for two (and higher) dimensional inputs [209].

One consideration when using convolutional or recurrent layers is how the window outputs are aggregated. Convolutional layers are typically followed by pooling or fully-connected layers that perform a global aggregation over all the window instances. While a fully-connected layer is aware of each window position, a global pooling layer is typically agnostic to it. However, unlike a fully-connected layer, a global max-pooling operation can be applied to a variable size input. Where a global aggregation strategy may be less appropriate (e.g., long sequences), recurrent networks with memory [18, 188, 213] and/or attention [33, 78, 126, 146, 217] may be useful.

Finally, Figure 19d shows *tree-structured* (or *recursive*) neural networks [20, 62, 182, 183, 195] where the same cell is applied at multiple levels in a tree-like hierarchical fashion.

<sup>13</sup>It is obvious, but may be still worth pointing out, that a *global convolutional* layer is exactly the same as a fully-connected layer.

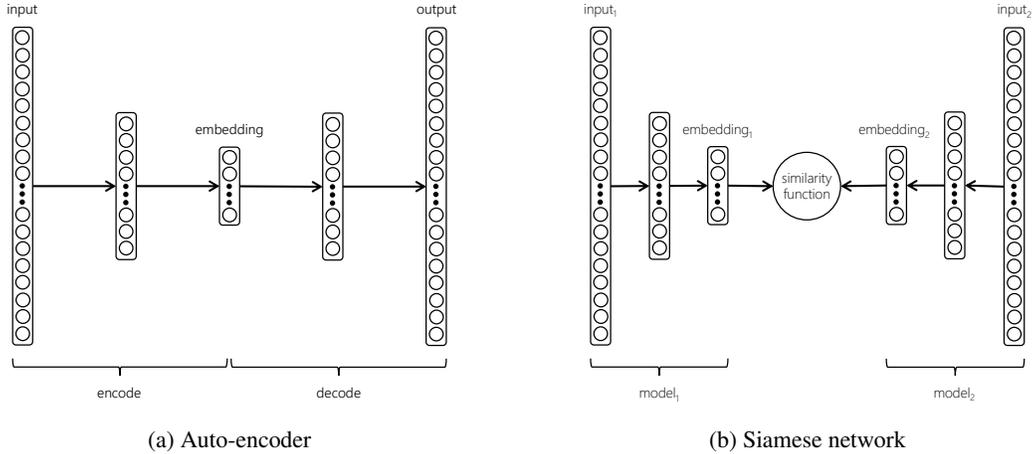


Figure 20: Both (a) the auto-encoder and (b) the Siamese network architectures are designed to learn compressed representations of inputs. In an auto-encoder the embeddings are learnt by minimizing the self-reconstruction error, whereas a Siamese network focuses on retaining the information that is necessary for determining the similarity between a pair of items (say, a query and a document).

**Auto-encoders** The auto-encoder architecture [14, 16, 164] is based on the *information bottleneck method* [197]. The goal is to learn a compressed representation  $\vec{x} \in \mathbb{R}^k$  of items from their higher-dimensional vector representations  $\vec{v} \in \mathbb{R}^K$ , such that  $k \ll K$ . The model has an hour-glass shape as shown in Figure 20a and is trained by feeding in the high-dimensional vector inputs and trying to re-construct the same representation at the output layer. The lower-dimensional middle layer forces the encoder part of the model to extract the *minimal sufficient statistics* of  $\vec{v}$  into  $\vec{x}$ , such that the decoder part of the network can reconstruct the original input back from  $\vec{x}$ . The model is trained by minimizing the reconstruction error between the input  $\vec{v}$  and the actual output of the decoder  $\vec{v}'$ . The squared-loss is popularly employed.

$$\mathcal{L}_{\text{auto-encoder}}(\vec{v}, \vec{v}') = \|\vec{v} - \vec{v}'\|^2 \quad (50)$$

**Siamese networks** Siamese networks were originally proposed for comparing fingerprints [10] and signatures [21]. Yih et al. [222] later adapted the same architecture for comparing short texts.

The siamese network, as seen in Figure 20b, resembles the auto-encoder architecture (if you squint hard enough!)—but unlike the latter is trained on pairs of inputs  $(input_1, input_2)$ . The architecture consists of two models ( $model_1$  and  $model_2$ ) that project  $input_1$  and  $input_2$ , respectively, to  $\vec{v}_1$  and  $\vec{v}_2$  in a common embedding space. A pre-defined metric (e.g., cosine similarity) is used to then compute the similarity between  $\vec{v}_1$  and  $\vec{v}_2$ . The model parameters are optimized such that  $\vec{v}_1$  and  $\vec{v}_2$  are closer when the two inputs are expected to be similar, and further away otherwise.

One possible loss function is the logistic loss. If each training sample consist of a triple  $\langle \vec{v}_q, \vec{v}_{d1}, \vec{v}_{d2} \rangle$ , such that  $sim(\vec{v}_q, \vec{v}_{d1})$  should be greater than  $sim(\vec{v}_q, \vec{v}_{d2})$ , then we minimize,

$$\mathcal{L}_{\text{siamese}}(\vec{v}_q, \vec{v}_{d1}, \vec{v}_{d2}) = \log\left(1 + e^{\gamma(sim(\vec{v}_q, \vec{v}_{d2}) - sim(\vec{v}_q, \vec{v}_{d1}))}\right) \quad (51)$$

where,  $\gamma$  is a constant that is often set to 10. Typically both the models— $model_1$  and  $model_2$ —share identical architectures, but can also choose to share the same parameters.

It is important to note that, unlike the auto-encoder, the minimal sufficient statistics retained by a Siamese network is dictated by which information it deems important for determining the similarity between the paired items.

### 6.3 Neural toolkits

In recent years, the advent of numerous flexible toolkits [1, 6, 31, 38, 91, 152, 198, 225] has had a catalytic influence on the area of neural networks. Most of the toolkits define a set of common neural operations that—like Lego<sup>14</sup> blocks—can be composed to build complex network architectures.<sup>15</sup> Each instance of these neural operations or *computation nodes* can have associated learnable parameters that are updated during training, and these parameters can be shared between different parts of the network if necessary. Every computation node under this framework must implement the appropriate logic for,

- computing the output of the node given the input (forward-pass)
- computing the gradient of the loss with respect to the inputs, given the gradient of the loss with respect to the output (backward-pass)
- computing the gradient of the loss with respect to its parameters, given the gradient of the loss with respect to the output (backward-pass)

A deep neural network, such as the one in Figure 16 or ones with much more complex architectures (e.g., [76, 107, 193]), can then be specified by chaining instances of these available computation nodes, and trained end-to-end on large datasets using backpropagation over GPUs or CPUs. In IR, various application interfaces [142, 201] bind these neural toolkits with existing retrieval/indexing frameworks, such as Indri [186].

Refer to [179] for a comparison of different neural toolkits based on their speed of training using standard performance benchmarks.

## 7 Deep neural models for IR

Traditionally, deep neural network models have much larger number of learnable parameters than their shallower counterparts. A DNN with a large set of parameters can easily overfit to smaller training datasets [231]. Therefore, during model design it is typical to strike a balance between the number of model parameters and the size of the data available for training. Data for ad-hoc retrieval mainly consists of,

- Corpus of search queries
- Corpus of candidate documents
- Ground truth—in the form of either explicit human relevance judgments or implicit labels (e.g., from clicks)—for query-document pairs

While both large scale corpora of search queries [46, 159] and documents [9, 29, 45] are publicly available for IR research, the amount of relevance judgments that can be associated with them are often limited outside of large industrial research labs—mostly due to user privacy concerns. We note that we are interested in datasets where the raw text of the query and the document is available. Therefore, this excludes large scale public labelled datasets for learning-to-rank (e.g., [122]) that don't contain the textual contents.

The proportion of labelled and unlabelled data that is available influences the *level of supervision* that can be employed for training these deep models. Most of the models we covered in Section 5 operate under the data regime where large corpus of documents or queries is available, but limited (or no) labelled data. Under such settings where no direct supervision or relevance judgments is provided, typically an *unsupervised* approach is employed (e.g., [174]). The unlabelled document (or query) corpus is used to learn good text representations, and then these learnt representations are incorporated into an existing retrieval model or a query-document similarity metric. If small amounts of labelled data are available, then that can be leveraged to train a retrieval model with few parameters that in turn uses text representations that is pre-trained on larger unlabelled corpus. Examples of such *semi-supervised* training includes models such as [71, 157, 158]. In contrast, *fully-supervised*

<sup>14</sup><https://en.wikipedia.org/wiki/Lego>

<sup>15</sup><http://www.inference.vc/content/images/2016/01/9k-.jpg>

Table 3: Comparing the nearest neighbours for "seattle" and "taylor swift" in the CDSSM embedding spaces when the model is trained on query-document pairs vs. query prefix-suffix pairs. The former resembles a Topical notion of similarity between terms, while the latter is more Typical in the definition of inter-term similarities.

seattle		taylor swift	
Query-Document	Prefix-Suffix	Query-Document	Prefix-Suffix
weather seattle	chicago	taylor swift.com	lady gaga
seattle weather	san antonio	taylor swift lyrics	meghan trainor
seattle washington	denver	how old is taylor swift	megan trainor
ikea seattle	salt lake city	taylor swift twitter	nicki minaj
west seattle blog	seattle wa	taylor swift new song	anna kendrick

models such as [37, 88, 141, 176], optimize directly for the target task by training on large number of labelled query-document pairs.

It is also useful to distinguish between deep neural models that focus on ranking long documents, from those that rank short texts (e.g., for the question-answering task, or for document ranking where the document representation is based on the title or on clicked queries). The challenges in short text ranking are somewhat distinct from those involved in the ad-hoc retrieval task [36]. When computing similarity between pairs of short-texts, vocabulary mismatches are more likely than when the retrieved items contain long text descriptions [133]. Neural models that perform matching in an embedding space tends to be more robust towards the vocabulary mismatch problem compared to lexical term-based matching models. On the other hand, documents with long body texts may contain mixture of many topics and the query matches may be spread over the whole document. A neural document ranking model (NDRM) must effectively aggregate the relevant matches from different parts of a long document.

In the rest of this section, we discuss different types of NDRM architectures and approaches that have been explored in the literature.

## 7.1 Document auto-encoders

Salakhutdinov and Hinton [174] proposed one of the earliest deep neural models for ad-hoc retrieval. The model is a deep auto-encoder trained on unlabelled document corpus. The model treats each document as a bag-of-terms and uses a one-hot vector for representing the terms themselves—considering only top two thousand most popular terms in the corpus after removing stopwords. Salakhutdinov and Hinton [174] first pre-train the model layer-by-layer, and then train it further end-to-end for additional tuning. The model uses binary hidden units and therefore the learnt vector representations of documents are also binary.

The *Semantic Hashing* model generates a condensed binary vector representation (or a hash) of documents. Given a search query, a corresponding hash is generated and the relevant candidate documents quickly retrieved that match the same hash vector. A standard IR model can then be employed to rank between the selected documents.

Semantic hashing is an example of a document encoder based approach to IR. The vocabulary size of two thousand distinct terms may be too small for most practical IR tasks. A larger vocabulary or a different term representation strategy—such as the character trigram based representation of Figure 18c—may be considered. Another shortcoming of the auto-encoder architecture is that it minimizes the document reconstruction error which may not align exactly with the goal of the target IR task. A better alternative may be to train on query-document paired data where the choice of what constitutes as the minimal sufficient statistics of the document is influenced by what is important for determining relevance of the document to likely search queries. In line with this intuition, we next discuss the Siamese architecture based models.

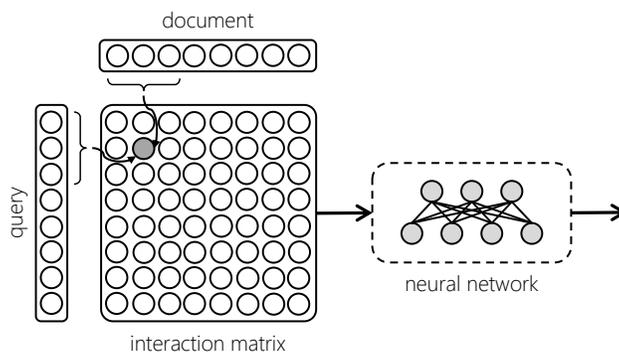


Figure 21: Schematic view of an interaction matrix generated by comparing windows of text from the query and the document. A deep neural network—such as a CNN—operates over the interaction matrix to find patterns of matches that suggest relevance of the document to the query.

## 7.2 Siamese networks

In recent years, several deep neural models based on the Siamese architecture have been explored especially for short text matching. The *Deep Semantic Similarity Model* (DSSM) [88] is one such architecture that trains on query and document title pairs where both the pieces of texts are represented as bags-of-character-trigrams. The DSSM architecture consists of two deep models—for the query and the document—with all fully-connected layers and cosine distance as the choice of similarity function in the middle. Huang et al. [88] proposed to train the model on clickthrough data where each training sample consists of a query  $q$ , a positive document  $d^+$  (a document that was clicked by a user on the SERP for that query), and a set of negative documents  $D^-$  randomly sampled with uniform probability from the full collection. The model is trained by minimizing the cross-entropy loss after taking a softmax over the model outputs for all the candidate documents,

$$\mathcal{L}_{dssm}(q, d^+, D^-) = -\log\left(\frac{e^{\gamma \cdot \cos(\vec{q}, \vec{d}^+)}}{\sum_{d \in D} e^{\gamma \cdot \cos(\vec{q}, \vec{d})}}\right) \quad (52)$$

$$\text{where, } D = \{d^+\} \cup D^- \quad (53)$$

While, DSSM [88] employs deep fully-connected architecture for the query and the document models, more sophisticated architectures involving convolutional layers [59, 86, 177, 178], recurrent layers [155, 156], and tree-structured networks [195] have also been explored. The similarity function can also be parameterized and implemented as additional layers of the neural network as in [176]. Most of these models have been evaluated on the short text matching task, but Mitra et al. [141] recently reported meaningful performances on the long document ranking task from models like DSSM [88] and CDSSM [177]. Mitra et al. [141] also show that sampling the negative documents uniformly from the collection is less effective to using documents that are closer to the query intent but judged as non-relevant by human annotators.

**Notions of similarity** It is important to emphasize that our earlier discussion in Section 4.2 on different notions of similarity between terms that can be learnt by shallow embedding models is also relevant in the context of these deeper architectures. In the case of Siamese networks, such as the convolutional-DSSM (CDSSM) [177], the notion of similarity being modelled depends on the choice of the paired data that the model is trained on. When the CDSSM is trained on query and document title pairs [177] then the notion of similarity is more *Topical* in nature. Mitra and Craswell [139] trained the same CDSSM architecture on query prefix-suffix pairs which, in contrast, captures a more *Typical* notion of similarity, as shown in Table 3. In a related work, Mitra [138] demonstrated that the CDSSM model when trained on session-query pairs is amenable to vector-based text analogies.

The **President** of the **United States** of America (PO-TUS) 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**.

The President of the **United States** of America (PO-TUS) 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**.

(a) Lexical model

(b) Semantic model

Figure 22: Analysis of term importance for estimating the relevance of a passage to the query “United States President” by a lexical and a semantic deep neural network model. The lexical model only considers the matches of the query terms in the document, but gives more emphasis to earlier occurrences. The semantic model is able to extract evidence of relevance from related terms such as “Obama” and “federal”.

$$\vec{v}_{\text{things to do in london}} - \vec{v}_{\text{london}} + \vec{v}_{\text{new york}} \approx \vec{v}_{\text{new york tourist attractions}} \quad (54)$$

$$\vec{v}_{\text{university of washington}} - \vec{v}_{\text{seattle}} + \vec{v}_{\text{denver}} \approx \vec{v}_{\text{university of colorado}} \quad (55)$$

$$\vec{v}_{\text{new york}} + \vec{v}_{\text{newspaper}} \approx \vec{v}_{\text{new york times}} \quad (56)$$

By modelling different notions of similarity these deep neural models tend to be more suitable for other IR tasks, such as query auto-completion [139] or session-based personalization [138].

### 7.3 Interaction-based networks

Siamese networks represent both the query and the document using single embedding vectors. Alternatively, we can individually compare different parts of the query with different parts of the document, and then aggregate these partial evidence of relevance. Especially, when dealing with long documents—that may contain a mixture of many topics—such a strategy may be more effective than trying to represent the full document as a single low-dimensional vector. Typically, in these approaches a sliding window is moved over both the query and the document text and each instance of the window over the query is compared against each instance of the window over the document text (see Figure 21). The terms within each window can be represented in different ways including, one-hot vectors, pre-trained embeddings, or embeddings that are updated during the model training. A neural model (typically convolutional) operates over the generated interaction matrix and aggregates the evidence across all the pairs of windows compared.

The interaction matrix based approach have been explored both for short text matching [86, 124, 158, 208, 220, 223], as well as for ranking long documents [141, 157].

### 7.4 Lexical and semantic matching networks

Much of the explorations in neural IR models have focused on learning good representations of text. However, these representation learning models tend to perform poorly when dealing with rare terms and search intents. In Section 2.2, we highlighted the importance of modelling rare terms in IR. Based on similar motivations, Guo et al. [71] and Mitra et al. [141] have recently emphasized the importance of modelling lexical matches using deep neural networks. Mitra et al. [141] argue that Web search is a “tale of two queries”. For the query “pekarovic land company”, it is easier to estimate relevance based on patterns of exact matches of the rare term “pekarovic”. On the other hand, a neural model focused on matching in the embedding space is unlikely to have a good representation for this rare term. In contrast, for the query “what channel are the seahawks on today”, the target document likely contains “ESPN” or “Sky Sports”—not the term “channel”. A representation learning neural model can associate occurrences of “ESPN” in the document as positive evidence towards the document being relevant to the query. Figure 22 highlights the difference between the terms that influence

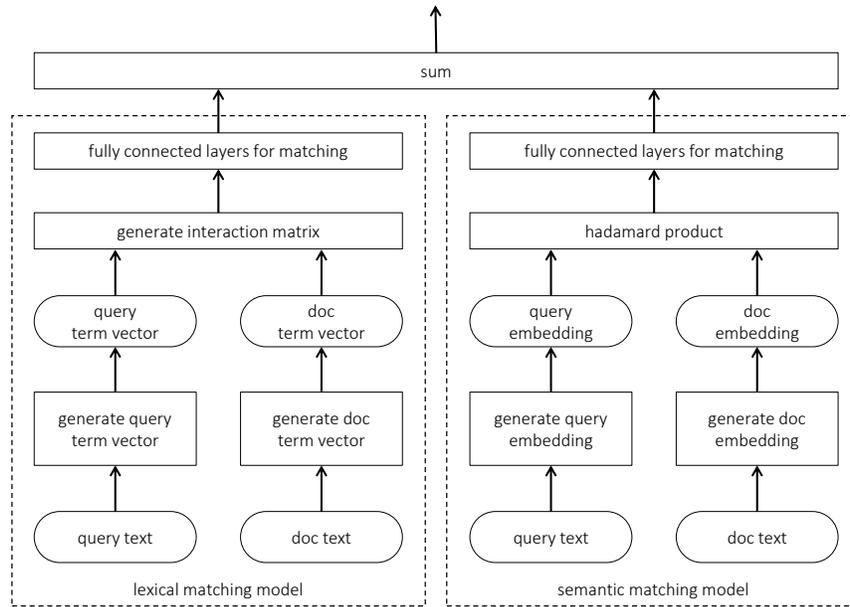


Figure 23: In the Duet architecture [141], the two sub-networks are jointly trained and the final output is a linear combination of the outputs of the lexical and the semantic matching sub-networks. The lexical matching sub-network (left) uses a convolutional model that operates over a binary interaction matrix.<sup>16</sup>The semantic matching sub-network (right) learns representations of query and document text for effective matching in the embedding space. Cross-entropy loss is used to train the network similar to other models in Section 7.2.

the estimation of relevance of the same query-passage pair by a lexical matching and a semantic matching model. A good neural IR model should incorporate both lexical and semantic matching signals [141].

Guo et al. [71] proposed to use histogram-based features in their DNN model to capture lexical notion of relevance. Mitra et al. [141] leverage large scale labelled data from Bing to train a *Duet* architecture (Figure 23) that learns to identify good patterns of both lexical and semantic matches jointly. Neural models that focus on lexical matching typically have fewer parameters, and can be trained under small data regimes—unlike their counterparts that focus on learning representations of text.

Interestingly, a query level analysis seems to indicate that both traditional non-neural IR approaches and more recent neural methods tend to perform well on different segments of queries depending on whether they focus on lexical or semantic matching. Figure 24 plots a few of these models based on their per-query NDCG values on a test set.

## 8 Conclusion

We present a tutorial on neural methods for information retrieval. For machine learning researchers who may be less familiar with IR tasks, we introduced the fundamentals of traditional IR models and metrics. For IR researchers, we summarized key concepts related to representation learning with (shallow or deep) neural networks. Finally, we presented some of the recent neural methods for document ranking and question-answer matching.

<sup>16</sup>It is important to emphasize, that while Mitra et al. [141] and others have used interaction-based representation for modelling lexical matches, the two ideas are distinct. Some interaction-matrix based representations compare texts using their pre-trained embeddings [86, 223]. Similarly, lexical matching can be modelled without employing an interaction matrix based representation [71].

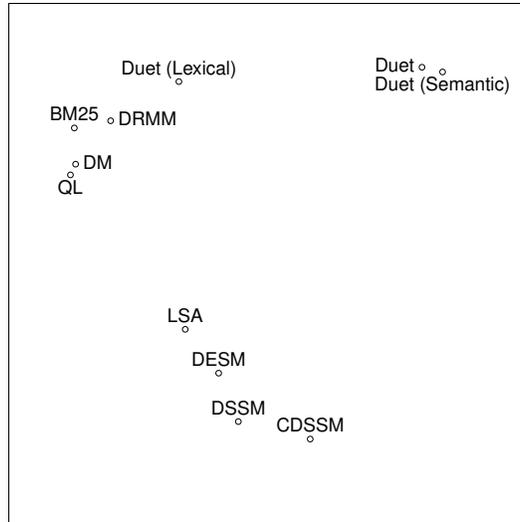


Figure 24: A demonstration that IR models that focus on lexical matching tend to perform well on queries that are distinct from queries on which semantic matching models achieve good relevance. Each model is represented by a vector of NDCG scores achieved on a set of test queries. For visualization, t-SNE [128] is used to plot the points in a two-dimensional space. Lexical matching models (BM25, QL, DM, DRMM, and Duet-Lexical) are seen to form a cluster—as well as the models that focus on representation learning.

We have focused on retrieval of long and short text. In the case of long text, the model must deal with variable length documents, where the relevant sections of a document may be surrounded by irrelevant text. For both long and short text, but particularly for short, IR models should also deal with the query-document vocabulary mismatch problem, by learning how patterns of query words and (different) document words can indicate relevance. Models should also consider lexical matches when the query contains rare terms—such as a person’s name or a product model number—not seen during training, and to avoid retrieving semantically related but irrelevant results.

An ideal model for information retrieval would be able to infer the meaning of a query from context. Given a query about the Prime Minister of UK, for example, it may be obvious from context whether it refers to John Major or Teresa May—perhaps due to the time period of the corpus, or it may need to be disambiguated based on other context such as the other query terms or the user’s short or long-term history. The ideal IR model may need to encode this context, which means that the model is like a *library* that effectively memorizes massive number of connections between entities and contexts. The number of learnable parameters of a ML model, however, is typically fixed, which may imply that there is a limited budget for how much real world knowledge the model can incorporate during training. An ideal model, therefore, may also need to learn to be like a *librarian* with incomplete domain knowledge, but capable of reading documents related to the current query and reasoning about the meaning of the query as part of the retrieval process.

Many of the breakthroughs in deep learning have been motivated by the needs of specific application areas. Convolutional neural networks, for example, are particularly popular with the vision community, whereas recurrent architectures find more applications in speech recognition and NLP. It is likely that the specific needs and challenges of IR tasks may motivate novel neural architectures and methods. Future IR explorations may also be motivated by developments in related areas, such as NLP. For example, neural architectures that have been evaluated on non-IR tasks [39, 50, 95, 99, 232] can be investigated in the retrieval context. Similarly, new methods for training deep models for NLP—e.g., using reinforcement learning [163, 224] and generative adversarial networks (GANs) [226]—may carry over to the IR setup.

However, given the pace at which the area of deep learning is growing, in terms of the number of new architectures and training regimes, 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. It would

be a bad outcome for the field if these explorations do not grow our understanding of the fundamental principles of machine learning and information retrieval. 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.<sup>17</sup> A better metaphor for the neural models may be a mirror that allows IR researchers to gain new insights into the underlying principles of IR. This may imply that we prefer neural models that, if not interpretable, then at least are amenable to analysis and interrogation. We may elicit more insights from simpler models while more sophisticated models may achieve state-of-the-art performances. As a community, we may need to focus on both to achieve results that are both impactful as well as insightful.

The focus of this article has been on ad-hoc retrieval and to a lesser extent on question-answering. However, neural approaches have shown interesting applications to other existing retrieval scenarios, including query auto-completion [139], query recommendation [184], session modelling [138], modelling diversity [215], modelling user click behaviours [19], knowledge-based IR [153], and even optimizing for multiple IR tasks [123]. In addition, recent trends suggest that advancements in deep neural networks methods are also fuelling emerging IR scenarios such as conversational IR [218, 234] and multi-modal retrieval [127]. Neural methods may have an even bigger impact on some of these other IR tasks.

IR also has a role in the context of the ambitions of the machine learning community. Retrieval is key to many one-shot learning approaches [102, 202]. Ghazvininejad et al. [60] proposed to “search” external information sources in the process of solving complex tasks using neural networks. The idea of learning local representations proposed by Diaz et al. [51] may be applicable to non-IR tasks. While we look at applying neural methods to IR, we should also look for opportunities to leverage IR techniques as part of—or in combination with—neural and other machine learning models.

Finally, we must also renew our focus on the fundamentals, including benchmarking and reproducibility. An important prerequisite to enable the “neural IR train” to steam forward is to build shared public resources—e.g., large scale datasets for training and evaluation, and repository of shared model implementations—and to ensure that appropriate bindings exist (e.g., [142, 201]) between popular IR frameworks and popular toolkits from the neural network community. The emergence of new IR tasks also demands rethinking many of our existing metrics. The metrics that may be appropriate for evaluating document ranking systems may be inadequate when the system generates textual answers in response to information seeking questions. In the latter scenario, the metric should distinguish between whether the response differs from the ground truth in the information content or in phrasing of the answer [57, 120, 145]. As multi-turn interactions with retrieval systems become more common, the definition of task success will also need to evolve accordingly. Neural IR should not only focus on novel techniques, but should also encompass all these other aspects.

## References

- [1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, and others. 2016. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467* (2016).
- [2] Nasreen Abdul-Jaleel, James Allan, W Bruce Croft, Fernando Diaz, Leah Larkey, Xiaoyan Li, Mark D Smucker, and Courtney Wade. 2004. UMass at TREC 2004: Novelty and HARD. (2004).
- [3] Eugene Agichtein, David Carmel, Dan Pelleg, Yuval Pinter, and Donna Harman. 2015. Overview of the TREC 2015 LiveQA Track.. In *TREC*.
- [4] Qingyao Ai, Liu Yang, Jiafeng Guo, and W Bruce Croft. 2016. Analysis of the paragraph vector model for information retrieval. In *Proc. ICTIR*. ACM, 133–142.
- [5] Qingyao Ai, Liu Yang, Jiafeng Guo, and W Bruce Croft. 2016. Improving language estimation with the paragraph vector model for ad-hoc retrieval. In *Proc. SIGIR*. ACM, 869–872.
- [6] Rami Al-Rfou, Guillaume Alain, Amjad Almahairi, Christof Angermueller, Dzmitry Bahdanau, Nicolas Ballas, Frédéric Bastien, Justin Bayer, Anatoly Belikov, Alexander Belopolsky, and

---

<sup>17</sup>[https://en.wikipedia.org/wiki/Law\\_of\\_the\\_instrument](https://en.wikipedia.org/wiki/Law_of_the_instrument)

- others. 2016. Theano: A Python framework for fast computation of mathematical expressions. *arXiv preprint arXiv:1605.02688* (2016).
- [7] Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. 2015. Rand-walk: A latent variable model approach to word embeddings. *arXiv preprint arXiv:1502.03520* (2015).
- [8] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473* (2014).
- [9] Peter Bailey, Nick Craswell, and David Hawking. 2003. Engineering a multi-purpose test collection for web retrieval experiments. *Information Processing & Management* 39, 6 (2003), 853–871.
- [10] Pierre Baldi and Yves Chauvin. 1993. Neural networks for fingerprint recognition. *Neural Computation* 5, 3 (1993), 402–418.
- [11] Marco Baroni, Georgiana Dinu, and Germán Kruszewski. 2014. Don't count, predict! a systematic comparison of context-counting vs. context-predicting semantic vectors. In *Proc. ACL*, Vol. 1. 238–247.
- [12] Marco Baroni and Alessandro Lenci. 2010. Distributional memory: A general framework for corpus-based semantics. *Computational Linguistics* 36, 4 (2010), 673–721.
- [13] Roland Barthes. 1977. *Elements of semiology*. Macmillan.
- [14] Yoshua Bengio and others. 2009. Learning deep architectures for AI. *Foundations and trends® in Machine Learning* 2, 1 (2009), 1–127.
- [15] Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. 2003. A neural probabilistic language model. *Journal of machine learning research* 3, Feb (2003), 1137–1155.
- [16] Yoshua Bengio, Pascal Lamblin, Dan Popovici, Hugo Larochelle, and others. 2007. Greedy layer-wise training of deep networks. *Proc. NIPS* 19 (2007), 153.
- [17] Adam Berger and John Lafferty. 1999. Information retrieval as statistical translation. In *Proc. SIGIR*. ACM, 222–229.
- [18] Antoine Bordes, Nicolas Usunier, Sumit Chopra, and Jason Weston. 2015. Large-scale simple question answering with memory networks. *arXiv preprint arXiv:1506.02075* (2015).
- [19] Alexey Borisov, Ilya Markov, Maarten de Rijke, and Pavel Serdyukov. 2016. A neural click model for web search. In *Proc. WWW*. Proc. WWW, 531–541.
- [20] Samuel R Bowman, Jon Gauthier, Abhinav Rastogi, Raghav Gupta, Christopher D Manning, and Christopher Potts. 2016. A fast unified model for parsing and sentence understanding. *arXiv preprint arXiv:1603.06021* (2016).
- [21] Jane Bromley, James W. Bentz, Léon Bottou, Isabelle Guyon, Yann LeCun, Cliff Moore, Eduard Säckinger, and Roopak Shah. 1993. Signature Verification Using A "Siamese" Time Delay Neural Network. *IJPRAI* 7, 4 (1993), 669–688.
- [22] Peter F Brown, John Cocke, Stephen A Della Pietra, Vincent J Della Pietra, Fredrick Jelinek, John D Lafferty, Robert L Mercer, and Paul S Roossin. 1990. A statistical approach to machine translation. *Computational linguistics* 16, 2 (1990), 79–85.
- [23] Peter F Brown, Vincent J Della Pietra, Stephen A Della Pietra, and Robert L Mercer. 1993. The mathematics of statistical machine translation: Parameter estimation. *Computational linguistics* 19, 2 (1993), 263–311.
- [24] John A Bullinaria and Joseph P Levy. 2007. Extracting semantic representations from word co-occurrence statistics: A computational study. *Behavior research methods* 39, 3 (2007), 510–526.

- [25] John A Bullinaria and Joseph P Levy. 2012. Extracting semantic representations from word co-occurrence statistics: stop-lists, stemming, and SVD. *Behavior research methods* 44, 3 (2012), 890–907.
- [26] Chris Burges, Tal Shaked, Erin Renshaw, Ari Lazier, Matt Deeds, Nicole Hamilton, and Greg Hullender. 2005. Learning to rank using gradient descent. In *Proceedings of the 22nd international conference on Machine learning*. ACM, 89–96.
- [27] Christopher JC Burges. 2010. From ranknet to lambdarank to lambdamart: An overview. *Learning* 11, 23-581 (2010), 81.
- [28] Christopher JC Burges, Robert Ragno, and Quoc Viet Le. 2006. Learning to rank with nonsmooth cost functions. In *NIPS*, Vol. 6. 193–200.
- [29] Jamie Callan, Mark Hoy, Changkuk Yoo, and Le Zhao. 2009. Clueweb09 data set. (2009).
- [30] Daniel Chandler. 1994. Semiotics for beginners. (1994).
- [31] Tianqi Chen, Mu Li, Yutian Li, Min Lin, Naiyan Wang, Minjie Wang, Tianjun Xiao, Bing Xu, Chiyuan Zhang, and Zheng Zhang. 2015. Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems. *arXiv preprint arXiv:1512.01274* (2015).
- [32] Kyunghyun Cho, Bart Van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. 2014. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259* (2014).
- [33] Jan K Chorowski, Dzmitry Bahdanau, Dmitriy Serdyuk, Kyunghyun Cho, and Yoshua Bengio. 2015. Attention-based models for speech recognition. In *Proc. NIPS*. 577–585.
- [34] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. 2014. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555* (2014).
- [35] Stéphane Clinchant and Florent Perronnin. 2013. Aggregating continuous word embeddings for information retrieval. In *Proceedings of the Workshop on Continuous Vector Space Models and their Compositionality*. 100–109.
- [36] Daniel Cohen, Qingyao Ai, and W Bruce Croft. 2016. Adaptability of neural networks on varying granularity ir tasks. *arXiv preprint arXiv:1606.07565* (2016).
- [37] Daniel Cohen and W Bruce Croft. 2016. End to End Long Short Term Memory Networks for Non-Factoid Question Answering. In *Proc. ICTIR*. ACM, 143–146.
- [38] Ronan Collobert, Koray Kavukcuoglu, and Clément Farabet. 2011. Torch7: A matlab-like environment for machine learning. In *BigLearn, NIPS Workshop*.
- [39] Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. 2011. Natural language processing (almost) from scratch. *The Journal of Machine Learning Research* 12 (2011), 2493–2537.
- [40] Nick Craswell. 2009. Mean reciprocal rank. In *Encyclopedia of Database Systems*. Springer, 1703–1703.
- [41] Nick Craswell. 2017. Neural Models for Full Text Search. In *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining*. ACM, 251–251.
- [42] Nick Craswell, W Bruce Croft, Maarten de Rijke, Jiafeng Guo, and Bhaskar Mitra. 2017. Neu-IR’17: Neural Information Retrieval. In *Proc. SIGIR*. ACM.
- [43] Nick Craswell, W Bruce Croft, Jiafeng Guo, Bhaskar Mitra, and Maarten de Rijke. 2016. Neu-IR: The SIGIR 2016 Workshop on Neural Information Retrieval. (2016).
- [44] Nick Craswell, W Bruce Croft, Jiafeng Guo, Bhaskar Mitra, and Maarten de Rijke. 2016. Report on the SIGIR 2016 Workshop on Neural Information Retrieval (Neu-IR). *ACM Sigir forum* 50, 2 (2016), 96–103.

- [45] Nick Craswell, David Hawking, Ross Wilkinson, and Mingfang Wu. 2003. Overview of the TREC 2002 Web track. In *TREC*, Vol. 3. 12th.
- [46] Nick Craswell, Rosie Jones, Georges Dupret, and Evelyne Viegas. 2009. *Proceedings of the 2009 workshop on Web Search Click Data*. ACM.
- [47] Ferdinand De Saussure, Wade Baskin, and Perry Meisel. 2011. *Course in general linguistics*. Columbia University Press.
- [48] Scott C. Deerwester, Susan T Dumais, Thomas K. Landauer, George W. Furnas, and Richard A. Harshman. 1990. Indexing by latent semantic analysis. *JASIS* 41, 6 (1990), 391–407.
- [49] Li Deng, Dong Yu, and others. 2014. Deep learning: methods and applications. *Foundations and Trends® in Signal Processing* 7, 3–4 (2014), 197–387.
- [50] Misha Denil, Alban Demiraj, Nal Kalchbrenner, Phil Blunsom, and Nando de Freitas. 2014. Modelling, visualising and summarising documents with a single convolutional neural network. *arXiv preprint arXiv:1406.3830* (2014).
- [51] Fernando Diaz, Bhaskar Mitra, and Nick Craswell. 2016. Query Expansion with Locally-Trained Word Embeddings. In *Proc. ACL*.
- [52] Fernando Diaz, Ryen White, Georg Buscher, and Dan Liebling. 2013. Robust models of mouse movement on dynamic web search results pages. In *Proceedings of the 22nd ACM international conference on Conference on information & knowledge management*. ACM, 1451–1460.
- [53] John Duchi, Elad Hazan, and Yoram Singer. 2011. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research* 12, Jul (2011), 2121–2159.
- [54] Jeffrey L Elman. 1990. Finding structure in time. *Cognitive science* 14, 2 (1990), 179–211.
- [55] David Ferrucci, Eric Brown, Jennifer Chu-Carroll, James Fan, David Gondek, Aditya A Kalyanpur, Adam Lally, J William Murdock, Eric Nyberg, John Prager, and others. 2010. Building Watson: An overview of the DeepQA project. *AI magazine* 31, 3 (2010), 59–79.
- [56] John R Firth. 1957. A synopsis of linguistic theory, 1930-1955. (1957).
- [57] Michel Galley, Chris Brockett, Alessandro Sordoni, Yangfeng Ji, Michael Auli, Chris Quirk, Margaret Mitchell, Jianfeng Gao, and Bill Dolan. 2015. deltaBLEU: A discriminative metric for generation tasks with intrinsically diverse targets. *arXiv preprint arXiv:1506.06863* (2015).
- [58] Debasis Ganguly, Dwaipayan Roy, Mandar Mitra, and Gareth JF Jones. 2015. Word Embedding based Generalized Language Model for Information Retrieval. In *Proc. SIGIR*. ACM, 795–798.
- [59] Jianfeng Gao, Patrick Pantel, Michael Gamon, Xiaodong He, Li Deng, and Yelong Shen. 2014. Modeling interestingness with deep neural networks. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*.
- [60] Marjan Ghazvininejad, Chris Brockett, Ming-Wei Chang, Bill Dolan, Jianfeng Gao, Wen-tau Yih, and Michel Galley. 2017. A Knowledge-Grounded Neural Conversation Model. *arXiv preprint arXiv:1702.01932* (2017).
- [61] Yoav Goldberg and Omer Levy. 2014. word2vec Explained: deriving Mikolov et al.’s negative-sampling word-embedding method. *arXiv preprint arXiv:1402.3722* (2014).
- [62] Christoph Goller and Andreas Kuchler. 1996. Learning task-dependent distributed representations by backpropagation through structure. In *Neural Networks, 1996., IEEE International Conference on*, Vol. 1. IEEE, 347–352.
- [63] Gene H Golub and Christian Reinsch. 1970. Singular value decomposition and least squares solutions. *Numerische mathematik* 14, 5 (1970), 403–420.
- [64] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. 2016. *Deep learning*. MIT Press.

- [65] Laura A Granka, Thorsten Joachims, and Geri Gay. 2004. Eye-tracking analysis of user behavior in WWW search. In *Proceedings of the 27th annual international ACM SIGIR conference on Research and development in information retrieval*. ACM, 478–479.
- [66] Alex Graves. 2013. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850* (2013).
- [67] Alex Graves, Marcus Liwicki, Santiago Fernández, Roman Bertolami, Horst Bunke, and Jürgen Schmidhuber. 2009. A novel connectionist system for unconstrained handwriting recognition. *IEEE transactions on pattern analysis and machine intelligence* 31, 5 (2009), 855–868.
- [68] Mihajlo Grbovic, Nemanja Djuric, Vladan Radosavljevic, and Narayan Bhamidipati. 2015. Search Retargeting using Directed Query Embeddings. In *Proc. WWW*. International World Wide Web Conferences Steering Committee, 37–38.
- [69] Mihajlo Grbovic, Nemanja Djuric, Vladan Radosavljevic, Fabrizio Silvestri, and Narayan Bhamidipati. 2015. Context-and Content-aware Embeddings for Query Rewriting in Sponsored Search. In *Proc. SIGIR*. ACM, 383–392.
- [70] Zhiwei Guan and Edward Cutrell. 2007. An eye tracking study of the effect of target rank on web search. In *Proceedings of the SIGCHI conference on Human factors in computing systems*. ACM, 417–420.
- [71] Jiafeng Guo, Yixing Fan, Qingyao Ai, and W Bruce Croft. 2016. A Deep Relevance Matching Model for Ad-hoc Retrieval. In *Proc. CIKM*. ACM, 55–64.
- [72] Jiafeng Guo, Yixing Fan, Qingyao Ai, and W Bruce Croft. 2016. Semantic Matching by Non-Linear Word Transportation for Information Retrieval. In *Proc. CIKM*. ACM, 701–710.
- [73] Richard HR Hahnloser, Rahul Sarpeshkar, Misha A Mahowald, Rodney J Douglas, and H Sebastian Seung. 2000. Digital selection and analogue amplification coexist in a cortex-inspired silicon circuit. *Nature* 405, 6789 (2000), 947–951.
- [74] Roy Harris. 2001. *Saussure and his Interpreters*. Edinburgh University Press.
- [75] Zellig S Harris. 1954. Distributional structure. *Word* 10, 2-3 (1954), 146–162.
- [76] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 770–778.
- [77] Robert Hecht-Nielsen and others. 1988. Theory of the backpropagation neural network. *Neural Networks* 1, Supplement-1 (1988), 445–448.
- [78] Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. In *Proc. NIPS*. 1693–1701.
- [79] Djoerd Hiemstra. 2001. *Using language models for information retrieval*. Taaluitgeverij Neslia Paniculata.
- [80] Felix Hill, Antoine Bordes, Sumit Chopra, and Jason Weston. 2015. The Goldilocks Principle: Reading Children’s Books with Explicit Memory Representations. *arXiv preprint arXiv:1511.02301* (2015).
- [81] Geoffrey Hinton, Li Deng, Dong Yu, George E Dahl, Abdel-rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara N Sainath, and others. 2012. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *Signal Processing Magazine, IEEE* 29, 6 (2012), 82–97.
- [82] Geoffrey E Hinton. 1984. Distributed representations. (1984).
- [83] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.

- [84] Kajita Hofmann, Bhaskar Mitra, Filip Radlinski, and Milad Shokouhi. 2014. An Eye-tracking Study of User Interactions with Query Auto Completion. In *Proc. CIKM*. ACM, 549–558.
- [85] Kurt Hornik, Maxwell Stinchcombe, and Halbert White. 1989. Multilayer feedforward networks are universal approximators. *Neural networks* 2, 5 (1989), 359–366.
- [86] Baotian Hu, Zhengdong Lu, Hang Li, and Qingcai Chen. 2014. Convolutional neural network architectures for matching natural language sentences. In *Proc. NIPS*. 2042–2050.
- [87] Gao Huang, Chuan Guo, Matt J Kusner, Yu Sun, Fei Sha, and Kilian Q Weinberger. 2016. Supervised Word Mover’s Distance. In *Proc. NIPS*. 4862–4870.
- [88] Po-Sen Huang, Xiaodong He, Jianfeng Gao, Li Deng, Alex Acero, and Larry Heck. 2013. Learning deep structured semantic models for web search using clickthrough data. In *Proc. CIKM*. ACM, 2333–2338.
- [89] Kevin Jarrett, Koray Kavukcuoglu, Yann LeCun, and others. 2009. What is the best multi-stage architecture for object recognition?. In *Computer Vision, 2009 IEEE 12th International Conference on*. IEEE, 2146–2153.
- [90] Kalervo Järvelin and Jaana Kekäläinen. 2002. Cumulated gain-based evaluation of IR techniques. *ACM Transactions on Information Systems (TOIS)* 20, 4 (2002), 422–446.
- [91] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. 2014. Caffe: Convolutional architecture for fast feature embedding. In *Proceedings of the 22nd ACM international conference on Multimedia*. ACM, 675–678.
- [92] Thorsten Joachims, Laura Granka, Bing Pan, Helene Hembrooke, and Geri Gay. 2005. Accurately interpreting clickthrough data as implicit feedback. In *Proceedings of the 28th annual international ACM SIGIR conference on Research and development in information retrieval*. Acm, 154–161.
- [93] Thorsten Joachims, Laura Granka, Bing Pan, Helene Hembrooke, Filip Radlinski, and Geri Gay. 2007. Evaluating the accuracy of implicit feedback from clicks and query reformulations in web search. *ACM Transactions on Information Systems (TOIS)* 25, 2 (2007), 7.
- [94] Rafal Jozefowicz, Oriol Vinyals, Mike Schuster, Noam Shazeer, and Yonghui Wu. 2016. Exploring the limits of language modeling. *arXiv preprint arXiv:1602.02410* (2016).
- [95] Nal Kalchbrenner, Edward Grefenstette, and Phil Blunsom. 2014. A convolutional neural network for modelling sentences. *arXiv preprint arXiv:1404.2188* (2014).
- [96] Tom Kenter, Alexey Borisov, and Maarten de Rijke. 2016. Siamese cbow: Optimizing word embeddings for sentence representations. *arXiv preprint arXiv:1606.04640* (2016).
- [97] Tom Kenter, Alexey Borisov, Christophe Van Gysel, Mostafa Dehghani, Maarten de Rijke, and Bhaskar Mitra. 2017. Neural Networks for Information Retrieval (NN4IR). In *Proc. SIGIR*. ACM.
- [98] Tom Kenter and Maarten de Rijke. Short Text Similarity with Word Embeddings. In *Proc. CIKM*, Vol. 15. 115.
- [99] Yoon Kim. 2014. Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882* (2014).
- [100] Yoon Kim, Yacine Jernite, David Sontag, and Alexander M Rush. 2015. Character-aware neural language models. *arXiv preprint arXiv:1508.06615* (2015).
- [101] Ryan Kiros, Richard Zemel, and Ruslan R Salakhutdinov. 2014. A multiplicative model for learning distributed text-based attribute representations. In *Proc. NIPS*. 2348–2356.
- [102] Gregory Koch. 2015. *Siamese neural networks for one-shot image recognition*. Ph.D. Dissertation. University of Toronto.

- [103] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. Imagenet classification with deep convolutional neural networks. In *Proc. NIPS*. 1097–1105.
- [104] Matt J Kusner, EDU Yu Sun, EDU Nicholas I Kolkin, and WUSTL EDU. From Word Embeddings To Document Distances. (????).
- [105] John Lafferty and Chengxiang Zhai. 2001. Document language models, query models, and risk minimization for information retrieval. In *Proceedings of the 24th annual international ACM SIGIR conference on Research and development in information retrieval*. ACM, 111–119.
- [106] Dmitry Lagun, Chih-Hung Hsieh, Dale Webster, and Vidhya Navalpakkam. 2014. Towards better measurement of attention and satisfaction in mobile search. In *Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval*. ACM, 113–122.
- [107] Gustav Larsson, Michael Maire, and Gregory Shakhnarovich. 2016. Fractalnet: Ultra-deep neural networks without residuals. *arXiv preprint arXiv:1605.07648* (2016).
- [108] Victor Lavrenko. 2008. *A generative theory of relevance*. Vol. 26. Springer Science & Business Media.
- [109] Victor Lavrenko and W Bruce Croft. 2001. Relevance based language models. In *Proceedings of the 24th annual international ACM SIGIR conference on Research and development in information retrieval*. ACM, 120–127.
- [110] Quoc V Le and Tomas Mikolov. 2014. Distributed Representations of Sentences and Documents.. In *ICML*, Vol. 14. 1188–1196.
- [111] Rémi Lebreton and Ronan Collobert. 2013. Word emdeddings through hellinger PCA. *arXiv preprint arXiv:1312.5542* (2013).
- [112] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. 2015. Deep learning. *Nature* 521, 7553 (2015), 436–444.
- [113] Yann LeCun, Fu Jie Huang, and Leon Bottou. 2004. Learning methods for generic object recognition with invariance to pose and lighting. In *Computer Vision and Pattern Recognition, 2004. CVPR 2004. Proceedings of the 2004 IEEE Computer Society Conference on*, Vol. 2. IEEE, II–104.
- [114] Yann LeCun, Koray Kavukcuoglu, and Clément Farabet. 2010. Convolutional networks and applications in vision. In *Circuits and Systems (ISCAS), Proceedings of 2010 IEEE International Symposium on*. IEEE, 253–256.
- [115] Omer Levy and Yoav Goldberg. 2014. Dependencybased word embeddings. In *Proc. ACL*, Vol. 2. 302–308.
- [116] Omer Levy, Yoav Goldberg, and Ido Dagan. 2015. Improving distributional similarity with lessons learned from word embeddings. *Transactions of the Association for Computational Linguistics* 3 (2015), 211–225.
- [117] Omer Levy, Yoav Goldberg, and Israel Ramat-Gan. 2014. Linguistic regularities in sparse and explicit word representations. *CoNLL-2014* (2014), 171.
- [118] Steven Levy. 2011. *In the plex: How Google thinks, works, and shapes our lives*. Simon and Schuster.
- [119] Hang Li and Zhengdong Lu. Deep Learning for Information Retrieval. (????).
- [120] Chia-Wei Liu, Ryan Lowe, Iulian V Serban, Michael Noseworthy, Laurent Charlin, and Joelle Pineau. 2016. How NOT to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation. *arXiv preprint arXiv:1603.08023* (2016).
- [121] Tie-Yan Liu. 2009. Learning to Rank for Information Retrieval. *Foundation and Trends in Information Retrieval* 3, 3 (March 2009), 225–331.

- [122] Tie-Yan Liu, Jun Xu, Tao Qin, Wenyong Xiong, and Hang Li. 2007. Letor: Benchmark dataset for research on learning to rank for information retrieval. In *Proceedings of SIGIR 2007 workshop on learning to rank for information retrieval*. 3–10.
- [123] Xiaodong Liu, Jianfeng Gao, Xiaodong He, Li Deng, Kevin Duh, and Ye-Yi Wang. Representation Learning Using Multi-Task Deep Neural Networks for Semantic Classification and Information Retrieval. *Proc. NAACL, May 2015* (????).
- [124] Zhengdong Lu and Hang Li. 2013. A deep architecture for matching short texts. In *Proc. NIPS*. 1367–1375.
- [125] Kevin Lund and Curt Burgess. 1996. Producing high-dimensional semantic spaces from lexical co-occurrence. *Behavior Research Methods, Instruments, & Computers* 28, 2 (1996), 203–208.
- [126] Minh-Thang Luong, Hieu Pham, and Christopher D Manning. 2015. Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025* (2015).
- [127] Lin Ma, Zhengdong Lu, Lifeng Shang, and Hang Li. 2015. Multimodal convolutional neural networks for matching image and sentence. In *Proceedings of the IEEE International Conference on Computer Vision*. 2623–2631.
- [128] Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of Machine Learning Research* 9, Nov (2008), 2579–2605.
- [129] Christopher Manning. 2016. Understanding Human Language: Can NLP and Deep Learning Help?. In *Proc. SIGIR*. ACM, 1–1.
- [130] Ivan Markovskiy. 2011. *Low rank approximation: algorithms, implementation, applications*. Springer Science & Business Media.
- [131] Irina Matveeva, Chris Burges, Timo Burkard, Andy Laucius, and Leon Wong. 2006. High accuracy retrieval with multiple nested ranker. In *Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval*. ACM, 437–444.
- [132] Donald Metzler and W Bruce Croft. 2005. A Markov random field model for term dependencies. In *Proc. SIGIR*. ACM, 472–479.
- [133] Donald Metzler, Susan Dumais, and Christopher Meek. 2007. Similarity measures for short segments of text. In *European Conference on Information Retrieval*. Springer, 16–27.
- [134] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781* (2013).
- [135] Tomas Mikolov, Martin Karafiát, Lukáš Burget, Jan Cernocký, and Sanjeev Khudanpur. 2010. Recurrent neural network based language model.. In *Interspeech*, Vol. 2. 3.
- [136] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Proc. NIPS*. 3111–3119.
- [137] Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. 2013. Linguistic Regularities in Continuous Space Word Representations.. In *HLT-NAACL*. Citeseer, 746–751.
- [138] Bhaskar Mitra. 2015. Exploring Session Context using Distributed Representations of Queries and Reformulations. In *Proc. SIGIR*. ACM, 3–12.
- [139] Bhaskar Mitra and Nick Craswell. 2015. Query Auto-Completion for Rare Prefixes. In *Proc. CIKM*. ACM.
- [140] Bhaskar Mitra and Nick Craswell. 2017. Neural Text Embeddings for Information Retrieval. In *Proc. WSDM*. ACM, 813–814.
- [141] Bhaskar Mitra, Fernando Diaz, and Nick Craswell. 2017. Learning to Match Using Local and Distributed Representations of Text for Web Search. In *Proc. WWW*. 1291–1299.

- [142] Bhaskar Mitra, Fernando Diaz, and Nick Craswell. 2017. Luandri: a Clean Lua Interface to the Indri Search Engine. In *Proc. SIGIR*. ACM.
- [143] Bhaskar Mitra, Eric Nalisnick, Nick Craswell, and Rich Caruana. 2016. A Dual Embedding Space Model for Document Ranking. *arXiv preprint arXiv:1602.01137* (2016).
- [144] Bhaskar Mitra, Milad Shokouhi, Filip Radlinski, and Katja Hofmann. 2014. On User Interactions with Query Auto-Completion. In *Proc. SIGIR*. 1055–1058.
- [145] Bhaskar Mitra, Grady Simon, Jianfeng Gao, Nick Craswell, and Li Deng. 2016. A Proposal for Evaluating Answer Distillation from Web Data. In *Proceedings of the SIGIR 2016 WebQA Workshop*.
- [146] Volodymyr Mnih, Nicolas Heess, Alex Graves, and others. 2014. Recurrent models of visual attention. In *Proc. NIPS*. 2204–2212.
- [147] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, and others. 2015. Human-level control through deep reinforcement learning. *Nature* 518, 7540 (2015), 529–533.
- [148] Guido F Montufar, Razvan Pascanu, Kyunghyun Cho, and Yoshua Bengio. 2014. On the number of linear regions of deep neural networks. In *Proc. NIPS*. 2924–2932.
- [149] Frederic Morin and Yoshua Bengio. 2005. Hierarchical Probabilistic Neural Network Language Model.. In *Aistats*, Vol. 5. Citeseer, 246–252.
- [150] Vinod Nair and Geoffrey E Hinton. 2010. Rectified linear units improve restricted boltzmann machines. In *Proc. ICML*. 807–814.
- [151] Eric Nalisnick, Bhaskar Mitra, Nick Craswell, and Rich Caruana. 2016. Improving Document Ranking with Dual Word Embeddings. In *Proc. WWW*.
- [152] Graham Neubig, Chris Dyer, Yoav Goldberg, Austin Matthews, Waleed Ammar, Antonios Anastasopoulos, Miguel Ballesteros, David Chiang, Daniel Clothiaux, Trevor Cohn, and others. 2017. DyNet: The Dynamic Neural Network Toolkit. *arXiv preprint arXiv:1701.03980* (2017).
- [153] Gia-Hung Nguyen, Lynda Tamine, Laure Soulier, and Nathalie Bricon-Souf. 2016. Toward a deep neural approach for knowledge-based ir. *arXiv preprint arXiv:1606.07211* (2016).
- [154] Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. MS MARCO: A Human Generated MACHine Reading COMprehension Dataset. *arXiv preprint arXiv:1611.09268* (2016).
- [155] H Palangi, L Deng, Y Shen, J Gao, X He, J Chen, X Song, and R Ward. 2014. Semantic Modelling with Long-Short-Term Memory for Information Retrieval. *arXiv preprint arXiv:1412.6629* (2014).
- [156] Hamid Palangi, Li Deng, Yelong Shen, Jianfeng Gao, Xiaodong He, Jianshu Chen, Xinying Song, and Rabab Ward. 2015. Deep Sentence Embedding Using the Long Short Term Memory Network: Analysis and Application to Information Retrieval. *arXiv preprint arXiv:1502.06922* (2015).
- [157] Liang Pang, Yanyan Lan, Jiafeng Guo, Jun Xu, and Xueqi Cheng. 2016. A study of match-pyramid models on ad-hoc retrieval. *arXiv preprint arXiv:1606.04648* (2016).
- [158] Liang Pang, Yanyan Lan, Jiafeng Guo, Jun Xu, Shengxian Wan, and Xueqi Cheng. 2016. Text Matching as Image Recognition. In *Proc. AAAI*.
- [159] Greg Pass, Abdur Chowdhury, and Cayley Torgeson. 2006. A picture of search. In *Proc. InfoScale*. ACM.
- [160] Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. *Proc. EMNLP* 12 (2014), 1532–1543.

- [161] Jay M Ponte and W Bruce Croft. 1998. A language modeling approach to information retrieval. In *Proc. SIGIR*. ACM, 275–281.
- [162] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250* (2016).
- [163] Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. 2015. Sequence level training with recurrent neural networks. *arXiv preprint arXiv:1511.06732* (2015).
- [164] Marc’Aurelio Ranzato, Christopher Poultney, Sumit Chopra, and Yann LeCun. 2006. Efficient learning of sparse representations with an energy-based model. In *Proceedings of the 19th International Conference on Neural Information Processing Systems*. MIT Press, 1137–1144.
- [165] Matthew Richardson, Christopher JC Burges, and Erin Renshaw. 2013. MCTest: A Challenge Dataset for the Open-Domain Machine Comprehension of Text.. In *EMNLP*, Vol. 3. 4.
- [166] Stephen Robertson, Hugo Zaragoza, and others. 2009. The probabilistic relevance framework: BM25 and beyond. *Foundations and Trends® in Information Retrieval* 3, 4 (2009), 333–389.
- [167] Stephen E Robertson, Evangelos Kanoulas, and Emine Yilmaz. 2010. Extending average precision to graded relevance judgments. In *Proceedings of the 33rd international ACM SIGIR conference on Research and development in information retrieval*. ACM, 603–610.
- [168] Douglas LT Rohde, Laura M Gonnerman, and David C Plaut. 2006. An improved model of semantic similarity based on lexical co-occurrence. *Commun. ACM* 8 (2006), 627–633.
- [169] Xin Rong. 2014. word2vec Parameter Learning Explained. *arXiv preprint arXiv:1411.2738* (2014).
- [170] Dwaipayan Roy, Debjyoti Paul, Mandar Mitra, and Utpal Garain. 2016. Using Word Embeddings for Automatic Query Expansion. *arXiv preprint arXiv:1606.07608* (2016).
- [171] Yossi Rubner, Carlo Tomasi, and Leonidas J Guibas. 1998. A metric for distributions with applications to image databases. In *Computer Vision, 1998. Sixth International Conference on*. IEEE, 59–66.
- [172] Magnus Sahlgren. 2006. *The Word-Space Model: Using distributional analysis to represent syntagmatic and paradigmatic relations between words in high-dimensional vector spaces*. Ph.D. Dissertation. Institutionen för lingvistik.
- [173] Hasim Sak, Andrew W Senior, and Françoise Beaufays. 2014. Long short-term memory recurrent neural network architectures for large scale acoustic modeling.. In *Interspeech*. 338–342.
- [174] Ruslan Salakhutdinov and Geoffrey Hinton. 2009. Semantic hashing. *International Journal of Approximate Reasoning* 50, 7 (2009), 969–978.
- [175] Jürgen Schmidhuber. 2015. Deep learning in neural networks: An overview. *Neural networks* 61 (2015), 85–117.
- [176] Aliaksei Severyn and Alessandro Moschitti. 2015. Learning to rank short text pairs with convolutional deep neural networks. In *Proc. SIGIR*. ACM, 373–382.
- [177] Yelong Shen, Xiaodong He, Jianfeng Gao, Li Deng, and Gregoire Mesnil. 2014. A latent semantic model with convolutional-pooling structure for information retrieval. In *Proc. CIKM*. ACM, 101–110.
- [178] Yelong Shen, Xiaodong He, Jianfeng Gao, Li Deng, and Grégoire Mesnil. 2014. Learning semantic representations using convolutional neural networks for Web search. In *Proc. WWW*. 373–374.
- [179] Shaohuai Shi, Qiang Wang, Pengfei Xu, and Xiaowen Chu. 2016. Benchmarking State-of-the-Art Deep Learning Software Tools. *arXiv preprint arXiv:1608.07249* (2016).

- [180] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, and others. 2016. Mastering the game of Go with deep neural networks and tree search. *Nature* 529, 7587 (2016), 484–489.
- [181] Amit Singhal, Chris Buckley, and Mandar Mitra. 1996. Pivoted document length normalization. In *Proc. SIGIR*. ACM, 21–29.
- [182] Richard Socher, Cliff C Lin, Chris Manning, and Andrew Y Ng. 2011. Parsing natural scenes and natural language with recursive neural networks. In *Proceedings of the 28th international conference on machine learning (ICML-11)*. 129–136.
- [183] Richard Socher, Jeffrey Pennington, Eric H Huang, Andrew Y Ng, and Christopher D Manning. 2011. Semi-supervised recursive autoencoders for predicting sentiment distributions. In *Proceedings of the conference on empirical methods in natural language processing*. Association for Computational Linguistics, 151–161.
- [184] Alessandro Sordani, Yoshua Bengio, Hossein Vahabi, Christina Lioma, Jakob G Simonsen, and Jian-Yun Nie. 2015. A Hierarchical Recurrent Encoder-Decoder for Generative Context-Aware Query Suggestion. *arXiv preprint arXiv:1507.02221* (2015).
- [185] Alessandro Sordani, Michel Galley, Michael Auli, Chris Brockett, Yangfeng Ji, Margaret Mitchell, Jian-Yun Nie, Jianfeng Gao, and Bill Dolan. 2015. A neural network approach to context-sensitive generation of conversational responses. *arXiv preprint arXiv:1506.06714* (2015).
- [186] Trevor Strohman, Donald Metzler, Howard Turtle, and W Bruce Croft. 2005. Indri: A language model-based search engine for complex queries. In *Proceedings of the International Conference on Intelligent Analysis*, Vol. 2. Citeseer, 2–6.
- [187] Bob L Sturm. 2014. A simple method to determine if a music information retrieval system is a “horse”. *IEEE Transactions on Multimedia* 16, 6 (2014), 1636–1644.
- [188] Sainbayar Sukhbaatar, Jason Weston, Rob Fergus, and others. 2015. End-to-end memory networks. In *Proc. NIPS*. 2440–2448.
- [189] Fei Sun, Jiafeng Guo, Yanyan Lan, Jun Xu, and Xueqi Cheng. 2015. Learning word representations by jointly modeling syntagmatic and paradigmatic relations. In *Proc. ACL*.
- [190] Fei Sun, Jiafeng Guo, Yanyan Lan, Jun Xu, and Xueqi Cheng. 2016. Semantic Regularities in Document Representations. *arXiv preprint arXiv:1603.07603* (2016).
- [191] Fei Sun, Jiafeng Guo, Yanyan Lan, Jun Xu, and Xueqi Cheng. 2016. Sparse word embeddings using l1 regularized online learning. In *Proc. IJCAI*. 2915–2921.
- [192] Ilya Sutskever, James Martens, and Geoffrey E Hinton. 2011. Generating text with recurrent neural networks. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*. 1017–1024.
- [193] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. 2015. Going deeper with convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1–9.
- [194] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. 2013. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199* (2013).
- [195] Kai Sheng Tai, Richard Socher, and Christopher D Manning. 2015. Improved semantic representations from tree-structured long short-term memory networks. *arXiv preprint arXiv:1503.00075* (2015).
- [196] Michael Taylor, Hugo Zaragoza, Nick Craswell, Stephen Robertson, and Chris Burges. 2006. Optimisation methods for ranking functions with multiple parameters. In *Proceedings of the 15th ACM international conference on Information and knowledge management*. ACM, 585–593.

- [197] Naftali Tishby, Fernando C Pereira, and William Bialek. 2000. The information bottleneck method. *arXiv preprint physics/0004057* (2000).
- [198] Seiya Tokui, Kenta Oono, Shohei Hido, and Justin Clayton. 2015. Chainer: a next-generation open source framework for deep learning. In *Proceedings of workshop on machine learning systems (LearningSys) in the twenty-ninth annual conference on neural information processing systems (NIPS)*.
- [199] Peter D Turney and Patrick Pantel. 2010. From frequency to meaning: Vector space models of semantics. *Journal of artificial intelligence research* 37 (2010), 141–188.
- [200] Christophe Van Gysel, Maarten de Rijke, and Evangelos Kanoulas. 2016. Learning latent vector spaces for product search. In *Proc. CIKM*. ACM, 165–174.
- [201] Christophe Van Gysel, Evangelos Kanoulas, and Maarten de Rijke. 2017. Pyndri: a Python Interface to the Indri Search Engine. *arXiv preprint arXiv:1701.00749* (2017).
- [202] Oriol Vinyals, Charles Blundell, Tim Lillicrap, Daan Wierstra, and others. 2016. Matching networks for one shot learning. In *Proc. NIPS*. 3630–3638.
- [203] Oriol Vinyals and Quoc Le. 2015. A neural conversational model. *arXiv preprint arXiv:1506.05869* (2015).
- [204] Ellen M Voorhees and Donna Harman. 2000. Overview of the eighth text retrieval conference (TREC-8). (2000), 1–24.
- [205] Ellen M Voorhees, Donna K Harman, and others. 2005. *TREC: Experiment and evaluation in information retrieval*. Vol. 1. MIT press Cambridge.
- [206] Ellen M Voorhees and Dawn M Tice. 2000. Building a question answering test collection. In *Proceedings of the 23rd annual international ACM SIGIR conference on Research and development in information retrieval*. ACM, 200–207.
- [207] Ivan Vulić and Marie-Francine Moens. 2015. Monolingual and cross-lingual information retrieval models based on (bilingual) word embeddings. In *Proc. SIGIR*. ACM, 363–372.
- [208] Shengxian Wan, Yanyan Lan, Jiafeng Guo, Jun Xu, Liang Pang, and Xueqi Cheng. 2015. A deep architecture for semantic matching with multiple positional sentence representations. *arXiv preprint arXiv:1511.08277* (2015).
- [209] Shengxian Wan, Yanyan Lan, Jun Xu, Jiafeng Guo, Liang Pang, and Xueqi Cheng. 2016. Match-srnn: Modeling the recursive matching structure with spatial rnn. *arXiv preprint arXiv:1604.04378* (2016).
- [210] Xiaojun Wan. 2007. A novel document similarity measure based on earth mover’s distance. *Information Sciences* 177, 18 (2007), 3718–3730.
- [211] Xiaojun Wan and Yuxin Peng. 2005. The earth mover’s distance as a semantic measure for document similarity. In *Proceedings of the 14th ACM international conference on Information and knowledge management*. ACM, 301–302.
- [212] Jason Weston, Antoine Bordes, Sumit Chopra, Alexander M Rush, Bart van Merriënboer, Armand Joulin, and Tomas Mikolov. 2015. Towards ai-complete question answering: A set of prerequisite toy tasks. *arXiv preprint arXiv:1502.05698* (2015).
- [213] Jason Weston, Sumit Chopra, and Antoine Bordes. 2014. Memory networks. *arXiv preprint arXiv:1410.3916* (2014).
- [214] Qiang Wu, Christopher JC Burges, Krysta M Svore, and Jianfeng Gao. 2010. Adapting boosting for information retrieval measures. *Information Retrieval* 13, 3 (2010), 254–270.
- [215] Long Xia, Jun Xu, Yanyan Lan, Jiafeng Guo, and Xueqi Cheng. 2016. Modeling Document Novelty with Neural Tensor Network for Search Result Diversification. In *Proc. SIGIR*. ACM, 395–404.

- [216] Yinglian Xie and David O’Hallaron. 2002. Locality in search engine queries and its implications for caching. In *INFOCOM 2002. Twenty-First Annual Joint Conference of the IEEE Computer and Communications Societies. Proceedings. IEEE*, Vol. 3. IEEE, 1238–1247.
- [217] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. 2015. Show, attend and tell: Neural image caption generation with visual attention. In *International Conference on Machine Learning*. 2048–2057.
- [218] Rui Yan, Yiping Song, and Hua Wu. 2016. Learning to respond with deep neural networks for retrieval-based human-computer conversation system. In *Proc. SIGIR*. ACM, 55–64.
- [219] Xiaohui Yan, Jiafeng Guo, Shenghua Liu, Xueqi Cheng, and Yanfeng Wang. 2013. Learning topics in short texts by non-negative matrix factorization on term correlation matrix. In *Proceedings of the SIAM International Conference on Data Mining*.
- [220] Liu Yang, Qingyao Ai, Jiafeng Guo, and W Bruce Croft. 2016. aNMM: Ranking Short Answer Texts with Attention-Based Neural Matching Model. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*. ACM, 287–296.
- [221] Yi Yang, Wen-tau Yih, and Christopher Meek. 2015. WikiQA: A Challenge Dataset for Open-Domain Question Answering.. In *EMNLP*. Citeseer, 2013–2018.
- [222] Wen-tau Yih, Kristina Toutanova, John C Platt, and Christopher Meek. 2011. Learning discriminative projections for text similarity measures. In *Proceedings of the Fifteenth Conference on Computational Natural Language Learning*. Association for Computational Linguistics, 247–256.
- [223] Wenpeng Yin, Hinrich Schütze, Bing Xiang, and Bowen Zhou. 2015. Abcnn: Attention-based convolutional neural network for modeling sentence pairs. *arXiv preprint arXiv:1512.05193* (2015).
- [224] Dani Yogatama, Phil Blunsom, Chris Dyer, Edward Grefenstette, and Wang Ling. 2016. Learning to compose words into sentences with reinforcement learning. *arXiv preprint arXiv:1611.09100* (2016).
- [225] Dong Yu, Adam Eversole, Mike Seltzer, Kaisheng Yao, Zhiheng Huang, Brian Guenter, Oleksii Kuchaiev, Yu Zhang, Frank Seide, Huaming Wang, and others. 2014. *An introduction to computational networks and the computational network toolkit*. Technical Report. Tech. Rep. MSR, Microsoft Research, 2014, <http://codebox/cntk>.
- [226] Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. 2016. Seqgan: sequence generative adversarial nets with policy gradient. *arXiv preprint arXiv:1609.05473* (2016).
- [227] Hamed Zamani and W Bruce Croft. 2016. Embedding-based query language models. In *Proc. ICTIR*. ACM, 147–156.
- [228] Hamed Zamani and W Bruce Croft. 2016. Estimating embedding vectors for queries. In *Proc. ICTIR*. ACM, 123–132.
- [229] Hugo Zaragoza, Nick Craswell, Michael J Taylor, Suchi Saria, and Stephen E Robertson. 2004. Microsoft Cambridge at TREC 13: Web and Hard Tracks.. In *TREC*, Vol. 4. 1–1.
- [230] Chengxiang Zhai and John Lafferty. 2001. A study of smoothing methods for language models applied to ad hoc information retrieval. In *Proc. SIGIR*. ACM, 334–342.
- [231] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. 2016. Understanding deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530* (2016).
- [232] Han Zhao, Zhengdong Lu, and Pascal Poupart. 2015. Self-Adaptive Hierarchical Sentence Model. *arXiv preprint arXiv:1504.05070* (2015).
- [233] Guoqing Zheng and Jamie Callan. 2015. Learning to Reweight Terms with Distributed Representations. In *Proc. SIGIR*. ACM, 575–584.

- [234] Xiangyang Zhou, Daxiang Dong, Hua Wu, Shiqi Zhao, R Yan, D Yu, Xuan Liu, and H Tian. 2016. Multi-view response selection for human-computer conversation. *EMNLP'16* (2016).
- [235] Mu Zhu. 2004. Recall, precision and average precision. *Department of Statistics and Actuarial Science, University of Waterloo, Waterloo 2* (2004), 30.
- [236] Guido Zuccon, Bevan Koopman, Peter Bruza, and Leif Azzopardi. 2015. Integrating and evaluating neural word embeddings in information retrieval. In *Proceedings of the 20th Australasian Document Computing Symposium*. ACM, 12.