

University of Waterloo E-Thesis Template for L^AT_EX

by

Zeynep Akkalyoncu Yilmaz

A thesis
presented to the University of Waterloo
in fulfillment of the
thesis requirement for the degree of
Master of Mathematics
in
Computer Science

Waterloo, Ontario, Canada, 2019

© Zeynep Akkalyoncu Yilmaz 2019

Abstract

Standard bag-of-words term-matching techniques in document retrieval fail to exploit rich semantic information embedded in the document texts. One promising recent trend in facilitating context-aware semantic matching has been the development of massively pre-trained language models, culminating in BERT as its most popular example today. In this work, we propose adapting BERT as a neural reranker for document retrieval with large improvements on news articles. Two fundamental issues arise in applying BERT to “ad hoc” document retrieval on newswire collections: relevance judgements in existing test collections are provided only at the document level, and documents often exceed the length that BERT was designed to handle. To overcome these challenges, we compute and aggregate sentence-level relevance scores to rank documents. We solve the problem of lack of appropriate relevance judgements by leveraging sentence-level and passage-level relevance judgements available in collections from other domains to capture cross-domain notions of relevance. We demonstrate that models of relevance can be transferred across domains. By leveraging semantic cues learned across various domains, we propose a model that achieves state-of-the-art results across three standard TREC newswire collections. We explore the effects of cross-domain relevance transfer, and trade-offs between using document and sentence scores for document ranking. We also present an end-to-end document retrieval system that incorporates the open-source Anserini information retrieval toolkit, discussing the related technical challenges and design decisions.

Table of Contents

List of Tables	v
List of Figures	vi
1 Introduction	1
1.1 Contributions	4
1.2 Thesis Organization	4
2 Background and Related Work	6
2.1 Document Retrieval	6
2.2 Pretrained Language Models	7
2.2.1 Feature-based Approaches	7
2.2.2 Fine-tuning Approaches	8
2.3 Machine-Learned Ranking Models	10
2.3.1 Learning to Rank Methods	10
2.3.2 Neural Document Retrieval	11
2.3.3 Comparison of Non-neural and Neural Methods	15
2.4 Evaluation Metrics	16
2.4.1 Mean Average Precision (MAP)	16
2.4.2 Precision at k (P@k)	17
2.4.3 Normalized Discounted Cumulative Gain (NDCG@k)	17

3	Cross-Domain Relevance Transfer with BERT	19
3.1	Datasets	19
3.1.1	TREC Microblog (MB)	20
3.1.2	MicroSoft MACHine Reading Comprehension (MS MARCO)	21
3.1.3	TREC Complex Answer Retrieval (CAR)	22
3.2	Modeling Relevance with BERT	23
3.2.1	Relevance Classifier	23
3.3	Reranking with BERT	25
3.4	Experimental Setup	26
3.4.1	Training and Inference with BERT	26
3.4.2	Evaluation	27
4	Architecture	28
4.1	Anserini	29
4.2	Main Module	30
4.3	Integration	30
4.4	Replicability and Reproducibility	31
5	Experimental Results	33
5.1	Datasets	33
5.2	Results and Analyses	34
5.2.1	Effect of Training Data	36
5.2.2	Number of Sentences	37
5.2.3	Comparison to Other Ranking Models	38
5.2.4	Per-Query Analysis	39
5.2.5	Effect of Length	40
5.2.6	Semantic Matching	42
6	Conclusion and Future Work	45
	References	46

List of Tables

3.1	Statistics about the MB dataset.	20
3.2	Statistics about the MS MARCO dataset.	21
3.3	Statistics about the CAR dataset.	22
5.1	Ranking effectiveness on Robust04.	34
5.2	Ranking effectiveness on Core17.	35
5.3	Ranking effectiveness on Core18.	35
5.4	Average AP with respect to query length on Robust04.	40
5.5	Ranking effectiveness on shortened MS MARCO and CAR evaluated on Robust04.	41
5.6	Retrieval effectiveness on pruned Robust04.	42

List of Figures

1.1	An example of a query-text pair from the TREC Robust04 collection where a relevant piece of text does not contain direct query matches.	1
2.1	The architecture combining BERT an additional output layer for the sentence pair classification task, where E represents the input embedding for each sentence and T_i the contextual representation of token i . Adapted from Devlin et al. [18].	9
2.2	The two types of deep matching architectures: representation-focused (a) and interaction-focused (b). Adapted from Guo et al. [23].	11
2.3	Visualization of best AP scores on Robust04 for 108 papers based on non-neural and neural approaches. Adapted from Yang et al. [66].	15
3.1	A sample relevant query and tweet pair from the MB dataset.	20
3.2	A sample relevant and non-relevant passage pair for a query from the MB dataset	21
3.3	Illustration of BERT input representation adapted from Devlin et al. [18].	24
4.1	Architecture of our system featuring tight integration between Python and the JVM.	28
5.1	Per-query difference in AP between BERT(MS MARCO \rightarrow MB) and the BM25+RM3 baseline on Robust04, Core17 and Core18.	38
5.2	Per-query difference in AP between BERT(CAR) and the BM25+RM3 baseline on Robust04, Core17 and Core18.	39
5.3	Attention visualizations of BERT(MS MARCO \rightarrow MB) for a sentence with a high BERT score for the query <code>international art crime</code>	44

Chapter 1

Introduction

Document retrieval refers to the task of generating a ranking of documents from a large corpus D in response to a query Q . In a typical document retrieval pipeline, an inverted index is constructed in advance from the collection, which often comprises unstructured text documents, for fast access during retrieval. When the user issues a query, the query representation is matched against the index, computing a similarity score for each document. The top most relevant documents based on their closeness to the query are returned to the user in order of relevance. This procedure may be followed by a subsequent re-ranking stage where the candidate documents outputted by the previous step are further re-ranked in a way that maximizes some retrieval metric such as average precision (AP).

Document retrieval systems traditionally rely on term-matching techniques, such as BM25, to judge the relevance of documents in a corpus. More specifically, the more common terms a document shares with the query, the more relevant it is considered. As a result, these systems may fail to detect documents that do not contain exact query terms, but are nonetheless relevant. For example, consider a document that expresses relevant information in a way that cannot be resolved without external semantic analysis. Figure 1 displays

<p>Query: international art crime</p> <p>Text: The thieves demand a ransom of \$2.2 million for the works and return one of them.</p>

Figure 1.1: An example of a query-text pair from the TREC Robust04 collection where a relevant piece of text does not contain direct query matches.

one such query-text pair where words semantically close to the query need to be identified to establish relevance. This “vocabulary mismatch” problem represents a long-standing challenge in information retrieval. To put its significance into context, Zhao et al. [72] show in their paper on term necessity prediction that, statistically, the average query terms do not appear in as many as 30% of relevant documents in TREC 3 to 8 “ad hoc” retrieval datasets.

Clearly, the classic exact matching approach to document retrieval neglects to exploit rich semantic information embedded in the document texts. To overcome this shortcoming, a number of models such as Latent Semantic Analysis [16], which map both queries and documents into high-dimensional vectors, and measure closeness between the two based on vector similarity, has been proposed. This innovation has enabled semantic matching to improve document retrieval by extracting useful semantic signals. With the advent of neural networks, it has become possible to learn better distributed representations of words that capture more fine-grained semantic and syntactic information [40], [47]. More recently, massively unsupervised language models that learn context-specific semantic information from copious amounts of data have changed the tide in NLP research (e.g: ELMo [48], GPT-2 [51]). These models can be applied to various downstream tasks with minimal task-specific fine-tuning, highlighting the power of transfer learning from large pre-trained models. Arguably the most popular example of these deep language representation models is the Bidirectional Encoder Representations from Transformers (BERT) [18]. BERT has achieved state-of-the-art results across a broad range of NLP tasks from question answering to machine translation.

While BERT has enjoyed widespread adoption across the NLP community, its application in information retrieval research has been limited in comparison. Guo et al. [22] suggest that the lackluster success of deep neural networks in information retrieval may be owing to the fact that they often do not properly address crucial characteristics of the “ad hoc” document retrieval task. Specifically, the relevance matching problem in information retrieval and semantic matching problem in natural language processing are fundamentally different in that the former depends heavily on exact matching signals, query term importance and diverse matching requirements. In other words, it is crucial to strike a good balance between exact and semantic matching in document retrieval. For this reason, we employ both document scores based on term-matching and semantic relevance scores to determine the relevance of documents.

In this thesis, we extend the work of Yang et al. [68] by presenting a novel way to apply BERT to “ad hoc” document retrieval on long documents – particularly, newswire articles – with significant improvements. Following Nogueira et al. [43], we adapt BERT for binary relevance classification over text to capture notions of relevance. We then deploy

the BERT-based re-ranker as part of a multi-stage architecture where an initial list of candidate documents is retrieved with a standard bag-of-words term matching technique. The BERT model is used to compute a relevance score for each constituent sentence, and the candidate documents are re-ranked by combining sentence scores with the original document score.

We emphasize that applying BERT to document retrieval on newswire documents is not trivial due to two main challenges: First of all, BERT has a maximum input length of 512 tokens, which is insufficient to accommodate the overall length of most news articles. To put this into perspective, a typical TREC Robust04 document has a median length of 679 tokens, and in fact, 66% of all documents are longer than 512 tokens. Secondly, most collections provide relevance judgements only at the document level. Therefore, we only know what documents are relevant for a given query, but not the specific spans within the document. To further aggravate this issue, a document is considered relevant as long as some part of it is relevant, and most of the document often has nothing to do with the query.

We address the abovementioned challenges by proposing two effective innovations: First, instead of relying solely on document-level relevance judgements, we aggregate sentence-level evidence to rank documents. As mentioned before, since standard newswire collections lack sentence level judgements to facilitate this approach, we instead explore leveraging sentence-level or passage-level judgements already available in collections in other domains, such as tweets and reading comprehension. To this end, we fine-tune BERT models on these out-of-domain collections to learn models of relevance. Surprisingly, we demonstrate that models of relevance can indeed be successfully transferred across domains. It is important to note that the representational power of neural networks come at the cost of challenges in interpretability. For this reason, we dedicate a portion of this thesis to error analysis experiments in an attempt to qualify and better understand the cross-domain transfer effects. We also elaborate on our engineering efforts to ensure reproducibility and replicability, and the technical challenges involved in bridging the worlds of natural language processing and information retrieval from a software engineering perspective.

1.1 Contributions

The main contributions of this thesis can be summarized as follows:

- We present two innovations to successfully apply BERT to *ad hoc* document retrieval with large improvements: integrating sentence-level evidence to address the fact that BERT cannot process long spans posed by newswire documents, and exploiting cross-domain models of relevance for collections without sentence- or passage-level annotations. With the proposed model, we establish state-of-the-art effectiveness on three standard TREC newswire collections at the time of writing. Our results on Robust04 exceed the previous highest known score of 0.3686 [14] with a non-neural method based on ensembles, which has stood unchallenged for ten years.
- We explore through various analyses the effects of cross-domain relevance transfer with BERT as well as the contributions of BM25 and sentence scores to the final document ranking. We investigate the effect of query and document length on retrieval effectiveness with BM25 and BERT, and the reasons behind the substantial improvements introduced with BERT.
- We release an end-to-end pipeline, Birch¹, that applies BERT to document retrieval over large document collections via integration with the open-source Anserini information retrieval toolkit. An accompanying Docker image is also included to ensure that anyone can easily deploy and test our system. We elaborate on the technical challenges in the integration of NLP and IR capabilities, and the rationale behind design decisions.

1.2 Thesis Organization

The remainder of this thesis is organized in the following order: Chapter 2 reviews related work in document retrieval, pretrained language models and document re-ranking models based on machine learning techniques as well as the metrics used in evaluation. Chapter 3 introduces the datasets used for fine-tuning BERT, motivates our approach with detailed information regarding the proposed model, and describes the experimental setup. Chapter 4 outlines an end-to-end pipeline for document retrieval with BERT by elaborating on

¹<https://github.com/castorini/birch>

the design decisions and challenges. Chapter 5 presents our results on three newswire collections – Robust04, Core17 and Core18, and provides a number of further analyses to help interpret the inner workings of BERT. Chapter 6 concludes the thesis by summarizing the main contributions and discussing future work.

Chapter 2

Background and Related Work

2.1 Document Retrieval

Traditional document retrieval techniques have evolved from the simple Boolean model to probabilistic models such as the Binary Independence Model over time to increase matching effectiveness. While these methods perform reasonably well for consistently short text like titles or abstracts, they have fallen short with the development of modern text collections with highly variable lengths. To this end, Okapi BM25 (commonly dubbed BM25) was developed as a bag-of-words ranking function that is sensitive to both term frequency and document length without introducing too many additional parameters [28]. Intuitively, BM25 pays more attention to the rarer terms in a query by increasing their term weight while dampening the matching signal for words that occur too frequently in a document with a term saturation mechanism. Term weights are also normalized with respect to the document length.

In addition to a term weighting scheme, query expansion has also been found to improve retrieval effectiveness by increasing recall. Unlike manual relevance feedback, pseudo relevance feedback allows for automatic local analysis without extended interaction with the user. RM3 is one such pseudo-relevance feedback mechanism where the original query is expanded by adding the top terms that appear in the contents of top k most relevant BM25 documents. While this method still relies on exact matching of query terms, it partly relieves the problem of synonymy. For instance, a query that contains the term “assistance” may be augmented with another high-frequency term “support” in relevant documents, therefore extending the range of matching. One obvious danger, however, is that retrieval may be incorrectly biased towards certain terms that occur frequently in

the most relevant documents, but are not directly relevant to the query. Despite their simplicity, well-tuned BM25+RM3 baselines achieve competitive effectiveness on TREC collections [32].

The most prominent approach to document retrieval today is to employ these techniques as the first step in a multi-stage architecture where an initial list of candidate documents is retrieved with a standard term-matching technique. A more computationally intensive model is then deployed as a re-ranker over the candidate documents to produce a final ranking. These re-rankers facilitate the architecture to detect and incorporate more complex relevance signals in the retrieved documents. A panoply of models have been proposed as re-rankers, including neural models [17] and statistical models based on term occurrence [33] and document similarity [30, 6].

2.2 Pretrained Language Models

Natural language processing tasks have traditionally been addressed with supervised learning on task-specific datasets. Due to the relatively small size of these datasets, training deep neural networks in this manner introduces the risk of overfitting on the training data, and lack of generalization across different datasets. With the increasing availability of large corpora, pretrained deep language models have been rapidly gaining traction among NLP researchers. Language model pretraining has proven extremely effective on many natural language processing tasks ranging from machine translation to reading comprehension. The underlying assumption in applying pretrained language models to downstream NLP tasks is that language modeling inherently captures many facets of language such as resolving long-term dependencies [34] and hierarchical patterns [21]. In general, pretrained language models can be applied to downstream tasks in one of two ways: feature-based and fine-tuning.

2.2.1 Feature-based Approaches

The feature-based approach, such as ELMo [48], employs deep pretrained representations learned with language modeling as additional features in task-specific architectures. This approach has the advantage of being easily incorporated into existing models with significant improvements in performance. ELMo [48] extends traditional word embeddings to learn context-sensitive features with a deep language model. Therefore, instead of taking the final layer of a deep LSTM (Long Short-Term Memory) as a word embedding,

ELMo embeddings are learned as a function of *all* the internal states of a bidirectional deep LSTM language model. This method is motivated by a thread of work in NLP that suggests that the higher levels of a deep LSTM capture context [39] and meaning while the lower levels learn syntactic features well [7]. While traditional pretrained word embeddings like GloVe [47] cannot differentiate between homonyms, ELMo can as it generates different embeddings for them depending on their context. ELMo embeddings are constructed as a shallow concatenation of independently trained left-to-right and right-to-left LSTMs. Peters et al. [48] show that integrating deep contextualized embeddings learned with ELMo into task-specific architectures significantly improves over the original performance in six NLP tasks, including question answering on SQuAD [52] and sentiment analysis on the Stanford Sentiment Treebank (SST-5) [56].

2.2.2 Fine-tuning Approaches

The fine-tuning approach is inspired by the growing trend in transfer learning. These models are first pretrained with respect to a language modeling objective, and then applied to downstream NLP tasks by “freezing” their last layer, and “fine-tuning” on external data for the specific task with minimal task-specific parameters. This approach has been shown to greatly boost the performance of many NLP tasks.

Radford et al. [51] claim that this phenomenon occurs because language models inherently capture many NLP tasks without explicit supervision. Therefore, they propose Generative Pretrained Transformers (GPT-2) to perform zero-shot task transfer on multiple sentence-level tasks from the GLUE benchmark [60] with impressive results. At the core of GPT-2 lies a multi-layer left-to-right transformer [58] decoder, with each layer consisting of a multi-head self-attention mechanism and fully connected feed-forward network [50]. The large capacity of the transformer is exploited by pretraining it on Google BookCorpus dataset [73] (800M words) where long contiguous spans of text allow the transformer to condition on long-range information.

To address the limitation of the unidirectional nature of GPT-2 [51], Bidirectional Encoder Representations from Transformers (BERT) [18] has introduced a novel way to pretrain bidirectional language models, and has since enjoyed widespread popularity across the NLP community. Standard language models cannot be conditioned on bidirectional context as this would cause the model to apply self-attention on the current token in a multi-layered context. However, BERT enables bidirectional language modeling by conditioning on both left and right context in all layers by employing a new pretraining objective called “masked language model” (MLM). Conceptually, MLM randomly masks some of the

input tokens, i.e: 15% of tokens in each sequence, at random with the goal of predicting the masked tokens based only on their left and right context. The final hidden vectors corresponding to the masked tokens are then fed into a softmax layer over the vocabulary as in a standard language model. This objective allows the representation to fuse both left and right context, which is indispensable for token-level tasks such as question answering, according to the authors. Ablation studies confirm that the bidirectional nature of BERT is the single most important factor in BERT’s performance. In addition to the novel language modeling approach, Devlin et al. [18] also propose a “next sentence prediction” task for applications that require an understanding of the relationship between two sentences, such as question answering or language language inference. Essentially, this trains a binary classifier to determine whether or not one sentence follows another sentence.

The underlying model architecture of BERT is a multi-layer bidirectional transformer [58]: The larger BERT model has 24 layers each with 1024 hidden nodes, and 16 self-attention heads in total. It is pretrained on the union of Google BookCorpus [73] (800M words) and English Wikipedia (2,500M words). The input representation for BERT is formed by concatenating the token with segment and position embeddings. Furthermore, the input may contain a single sentence or a sentence pair separated by the meta-token [SEP], i.e: separator. Each sequence is prepended with [CLS], corresponding to the “class” meta-token, whose final hidden state can be used for classification tasks. The words are represented

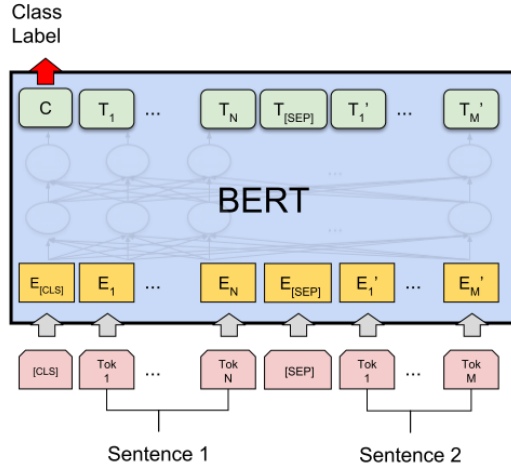


Figure 2.1: The architecture combining BERT an additional output layer for the sentence pair classification task, where E represents the input embedding for each sentence and T_i the contextual representation of token i . Adapted from Devlin et al. [18].

with WordPiece embeddings [61] with a vocabulary of 30,000 tokens. Originally proposed for segmentation problem in Japanese and Korean, the WordPiece model is used to divide words into small sub-word units in order to handle rare or out-of-vocabulary words more effectively. Positional embeddings are learned – not hard-coded – for up to 512, which is the maximum input size allowed by BERT.

To fine-tune BERT for classification tasks, a single-layer neural network is added on top of BERT with the class label as the input, and label probabilities are computed with softmax. The parameters of the additional layer and BERT are fine-tuned jointly to maximize the log-probability of the correct label. For span-level and token-level prediction tasks, the final step needs to be modified to account for multiple tokens. Figure 2.1 visualizes the model for fine-tuning BERT for the “sentence pair classification” task that takes two sentences separated by a [SEP] token as the input.

BERT has been applied to a broad range of NLP tasks from sentence classification to sequence labeling with impressive results. Most relevant to the task of document retrieval, applications of BERT include BERTserini by Yang et al. [67] which integrated BERT with Anserini for question answering over Wikipedia by fine-tuning BERT on SQuAD, and Nogueira et al. [43] who adopted BERT for passage re-ranking over MS MARCO.

2.3 Machine-Learned Ranking Models

2.3.1 Learning to Rank Methods

Learning to rank (LTR) methods have arisen in document retrieval to apply supervised machine learning techniques to automatically learn a ranking model. This trend has started in response to a growing number of relevance signals, particularly in web search, such as anchor texts. In this setup, training samples are extracted from large amounts of search log data in the form of a list of documents and their relevance labels for a number of queries. Unlike traditional ranking models, LTR models require a good deal of feature engineering, which can be time-consuming and difficult to generalize.

Existing algorithms for LTR can be categorized into three categories based on their input representation and loss function: pointwise, pairwise or listwise. [36] Pointwise algorithms such as McRank [31] approximate the LTR problem as a regression problem where a relevance score is predicted for each query-document pair. On the other hand, LTR is framed as a classification problem by pairwise algorithms like RankSVM [27] and LambdaMART [11] to choose the more relevant one given a pair of documents. Listwise

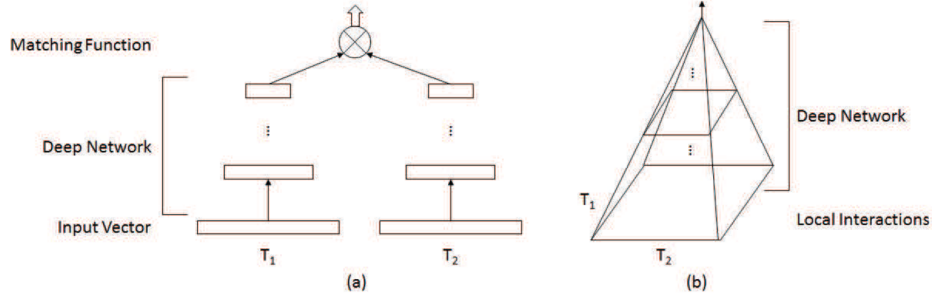


Figure 2.2: The two types of deep matching architectures: representation-focused (a) and interaction-focused (b). Adapted from Guo et al. [23].

algorithms such as ListNet [12] instead optimize the relevance score over the entire list of documents for a query. [36] claim that listwise LTR approaches often produce better rankings in practice.

2.3.2 Neural Document Retrieval

With the impressive results achieved by neural networks in many areas such as computer vision and natural language processing, document retrieval, too, has witnessed a shift from non-neural methods to neural methods over the last few years. Neural models have especially been instrumental in facilitating semantic matching in document retrieval. Neural models developed to address the deep matching problem in document retrieval can be divided into two broad categories based on their underlying architecture: representation-based and interaction-based. The high-level differences between these architectures can be observed in Figure 2.2.

Representation-based Models

Representation-based approaches first construct a representation from the input vectors for the query and document with a deep neural network, and then perform matching between these representations (see Figure 2.2). One set of such models loans the concept of word embeddings from natural language processing to represent query and documents. This paradigm represents words as low-dimensional continuous feature vectors that embody hidden semantic or syntactic dependencies. Some of the most popular pretrained English word embeddings include word2vec [40] trained on Google News, GloVe [47] on Common

Crawl, Wikipedia and Twitter, and fastText [10] on English webcrawl and Wikipedia. These word embeddings can be learned from scratch for a specific corpus or pretrained over large corpora and reused with significant improvements over the former option [57]. There has also been some effort in learning word embeddings to directly capture relevance matching [69, 20] rather than linguistic features as in word2vec [40] or GloVe [47]. Word embeddings are commonly used as input to many representation-based retrieval models.

Other representation-based architectures explore alternative ways to represent text for effective retrieval. DSSM (short for Deep Structured Semantic Models) [25] extends previously dominant latent semantic models to deep semantic matching for web search by projecting query and documents into a common low-dimensional space. In order to accommodate a large vocabulary required by the task, the text sequences are mapped into character-level trigrams with a word hashing layer before computing a similarity matrix through dot product and softmax layers. While shown effective on a private dataset comprised of log files of a commercial search engine, DSSM requires too much training data to be effective. Moreover, DSSM cannot match synonyms because it is based on the specific composition of words and not semantic proximity. C-DSSM [55] was proposed as an extension to DSSM by replacing the multi-layer perceptron with a convolutional layer to devise semantic vectors for search queries and Web document. By performing a max pooling operation to extract local contextual information at the n-gram level, a global vector representation is formed from the local features. Shen et al. [55] demonstrate that both local and global contextual features are necessary for semantic matching for Web search. While C-DSSM improves over DSSM by exploiting the context of each trigram, it still suffers from most of the same issues listed above.

Interaction-based Models

Interaction-based approaches capture local matching signals, and directly compute word-word similarity between query and document representations. In contrast to more shallow representation-based approaches, this setup allows the deep neural network to learn more complex hierarchical matching patterns across multiple layers. Some notable examples of these architectures include DRMM [23], KNRM [62] and DUET [42].

DRMM [23], which stands for Deep Relevance Matching Model, maps variable-length local interactions of query and document into a fixed-length matching histogram. A feed forward matching network is used to learn hierarchical matching patterns from the histogram representation, and a matching score is computed for each term. An overall matching score is obtained by aggregating the scores from each query term with a term gating network. KNRM [62] similarly calculates the word-word similarities between query and document

embeddings, but converts word-level interactions into ranking features with a novel kernel pooling technique. Specifically, a feature vector for each word in the query is constructed from the similarity matrix with k-max pooling. Ranking features are combined to form a final ranking score through a learning-to-rank layer. Unlike DRMM and KNRM, the goal of DUET [42] is to employ both local and distributed representations, therefore leveraging both exact matching and semantic matching signals. DUET is composed of two separate deep neural networks, one to match the query and the document using a one-hot representation, and another using learned distributed representations, which are trained jointly. The former estimates document relevance based on exact matches of the query terms in the document by computing an interaction matrix from one-hot encodings. The latter instead performs semantic matching by computing the element-wise product between the query and document embeddings. Their approach significantly outperform traditional baselines for web search with lots of click-through logs.

Contextualized Language Models

While the models introduced in Section 2.3.2 successfully leverage semantic information to varying degrees, they are limited by the size and variability of available training data. Ideally, these models would be trained on a large number of semantically and syntactically varied labeled query-document pairs; however, it is impractical to automatically gather a sufficient number of such training samples at scale. Instead massively pretrained unsupervised language models hold promises for obtaining better query and document representations, and therefore, achieving unprecedented effectiveness at semantic matching without the need for more relevance information. Section 2.2 outlines some of the most popular unsupervised language models that form the basis of effective retrieval architectures. In general, these language models are deployed as re-rankers over an initial list of candidate documents retrieved with traditional term-matching techniques in Section 2.1.

Modeling relevance requires an understanding of the relationship between two text sequences, e.g: the query and the document. Clearly, traditionally language modeling does not suffice to capture such a relationship. Fortunately, BERT facilitates such relevance classification by pre-training a binary next sentence prediction task based on its masking language model approach as discussed in Section 2.2. However, it is still not trivial to apply BERT to document retrieval because BERT was not designed to handle long spans of text, such as documents, given a maximum input length of 512 tokens. Partly because of this inherent challenge, the majority of work on re-ranking with BERT has focused on passage re-ranking instead of document re-ranking.

Notably, Nogueira et al. [43] proposed to re-rank MS MARCO passages based on a

simple re-implementation of BERT to learn a model of relevance, outperforming the previous state of the art by 27% in MRR@10 and replacing the previous top entry in the leaderboard of the MS MARCO passage retrieval task. Our neural model is inspired by the BERT re-implementation described in their paper. Padigela et al. [45] prioritize studying the reasons behind the gains that come with re-ranking MS MARCO passages with BERT. To put re-ranking with BERT into perspective, they compare their BERT-based re-ranker to feature based learning to rank models such as RankSVM [27] and a number of neural kernel matching models such as KNRM [62] and Conv-KNRM [15], and conclude that fine-tuning BERT is substantially more effective than either neural model. They also test four hypotheses regarding the behavior of matching with BERT compared to BM25; specifically, with respect to term frequency and document length.

To our knowledge, Yang et al. [68] are the first to successfully apply BERT to “ad hoc” document retrieval. They demonstrate that BERT can be fine-tuned to capture relevance matching by following the “next sentence classification” task of BERT on the TREC Microblog Tracks where document length does not pose an issue. They further propose overcoming the challenge of long documents by applying inference on each individual sentence and combining the top scores to compute a final document score. Their approach is motivated by user studies by Zhang et al. [71] which suggest that the most relevant sentence or paragraph in a document provides a good proxy for overall document relevance. The work of Yang et al. [68] pave the way for future work that culminated in this thesis.

More recently, MacAvaney et al. [37] shifted focus from incorporating BERT as a re-ranker to using its representation capabilities to improve existing neural architectures. By computing a relevance matrix between the query and each candidate document at each layer of a contextualized language model – in particular, ELMo or BERT – they report a high score of NDCG@20 0.5381 on Robust04 by combining CEDR (Contextualized Embeddings for Document Ranking) [37] with KNRM [62]. They also propose a joint model that combines the classification mechanism of BERT into existing neural architectures to help benefit from both deep semantic matching with BERT *and* relevance matching with traditional ranking architectures.

A recent arXiv preprint by Qiao et al. [49] also examines the performance and behavior of BERT when used as a re-ranker for passage ranking on MS MARCO and for document ranking on the TREC Web Track. Their findings are consistent with those of Nogueira et al. [43] in that BERT outperforms previous neural models on the passage re-ranking task on MS MARCO. For ad hoc document ranking, they explore using BERT both as representation-based and interaction-based rankers and in combination with KNRM [62] and Conv-KNRM [15]. However, they find that their re-ranking TREC Web Track documents with BERT performs worse than Conv-KNRM and feature-based learning-to-rank

models trained on user clicks in Bing’s search log.

2.3.3 Comparison of Non-neural and Neural Methods

Despite growing interest in neural models for document ranking, researchers have recently voiced concern as to whether or not they have truly contributed to progress [32], at least in the absence of large amounts of behavioral log data only available to search engine companies. Some of the models discussed in this section are designed for the web search task where a variety of other signals are available, such as large amounts of log data and the webgraph. However, this is not the case for “ad hoc” document retrieval where the only available data is the document text, which is the main focus of this thesis. The SIGIR Forum piece by Lin [32] also echoes the general skepticism concerning the empirical rigor and contributions of machine learning applications in Lipton et al. [35] and Sculley et al. [54]. In particular, Lin et al. [32] lament that comparisons to weak baselines sometimes inflate the merits of certain neural information retrieval methods.

To rigorously study the current state of document retrieval literature, Yang et al. [66] recently conducted a thorough meta-analysis of over 100 papers that report results on the TREC Robust 2004 Track. Their findings are illustrated in Figure 2.3 where the empty circles correspond to the baselines and filled circles to the best AP scores of each paper. The solid black line represents the best submitted run at AP 0.333, and the dotted black line the median TREC run at AP 0.258. The other line is a RM3 baseline run with default

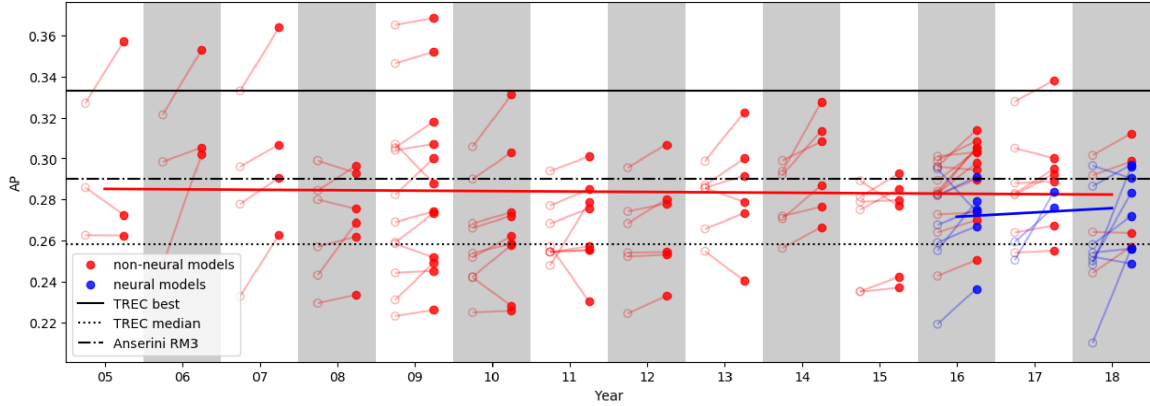


Figure 2.3: Visualization of best AP scores on Robust04 for 108 papers based on non-neural and neural approaches. Adapted from Yang et al. [66].

parameters from the Anserini open-source information retrieval toolkit [65] at AP 0.3903. The untuned RM3 baseline is more effective than 60% of all studied papers, and 20% of them report results below the TREC median. More surprisingly, only six of the papers report AP scores higher than the TREC best, with the highest being by Cormack et al. [14] in 2009 at AP 0.3686. Their approach is based on building an ensemble of multiple TREC runs with reciprocal rank fusion. Among the neural models, the highest encountered score is by Zamani et al. [70] in 2018 at AP 0.2971. Deviating from the dominant approach of deploying neural models as re-rankers, Zamani et al. [70] propose a standalone neural ranking model to learn a latent representation for each query and document, which is sparse enough to construct an inverted index for the entire collection. However, their reported result is still much lower than the TREC best and far below the best reported result of Cormack et al. [14]. Moreover, about half of the neural papers compare their results to a baseline *below* the TREC median, which is consistent with the claims of Lin et al. [32]. Overall, Figure 2.3 exhibits no clear upward trend in terms of AP on Robust04 from 2005 to 2018.

TREC Common Core 2017 (Core17) [1] and 2018 (Core18) [2] are two of the more recent document collections that we evaluate our models on, which are not nearly as well-studied as Robust04 yet. Excluding runs that make use of past labels or require human intervention, the TREC best run on Core17 is `umass baselnrm` at AP 0.275 and on Core18 `uwrmrg` at AP 0.276. To our knowledge, Neural Vector Spaces for Unsupervised Information Retrieval by Van Gysel et al. [24] represents the only major neural model evaluated on Core17. While their approach has the advantage of not requiring supervised relevance judgements, their reported results are quite low. Otherwise, evaluation of neural retrieval methods on both Core17 and Core18 has been limited.

2.4 Evaluation Metrics

Evaluation in information retrieval relies on the distinction between “relevant” and “irrelevant” documents with respect to an information need as expressed by a query. A number of automatic evaluation metrics has been formalized specifically for ranking tasks, some of which are described below.

2.4.1 Mean Average Precision (MAP)

Precision specifies what fraction of a set of retrieved documents is in fact relevant for a given query q . By extension, average precision (AP) expresses the average of the precision

values obtained for the set of top k documents for the query. Suppose that $D = \{d_1, \dots, d_{m_j}\}$ is the set of all relevant documents for a query q_j , then AP can be formulated as:

$$AP = \frac{1}{m_j} \sum_{k=1}^{m_j} P(R_{jk}) \quad (2.1)$$

where R_{jk} represents the set of top k ranked retrieval results.

The respective AP for each query $q_j \in Q$ can be aggregated to obtain mean average precision (MAP) for the overall retrieval effectiveness in the form of a single-figure measure of quality across various recall levels:

$$MAP = \frac{\sum_{j=1}^{|Q|} AP}{|Q|} = \frac{1}{|Q|} \sum_{j=1}^{|Q|} \frac{1}{m_j} \sum_{k=1}^{m_j} P(R_{jk}) \quad (2.2)$$

MAP is known to have especially good discrimination and stability compared to other evaluation metrics, which makes it the ideal choice for large text collections [38]. It is hence one of the standard metrics among the TREC community.

2.4.2 Precision at k (P@ k)

While MAP factors in precision at all recall levels, certain applications may have a distinctly different notion for ranking quality. Particularly in the case of web search, the user often only cares about the results on the first page or two. This restriction essentially requires measuring precision at fixed low levels of retrieved results, i.e: top k documents – hence the name for the metric “precision at k ”. On the one hand, it eliminates the need for any estimate of the size of the set of relevant documents because it is only concerned with the top documents. However, it also produces the least stable results out of all evaluation metrics. Moreover, precision at k does not average well because it is too sensitive to the total number of relevant documents for a query.

2.4.3 Normalized Discounted Cumulative Gain (NDCG@ k)

Cumulative gain (CG) simply computes the sum of relevance labels for all the retrieved documents, treating the search results as an unordered set. However, since a highly relevant document is inherently more useful when it appears higher up in the search results, CG

has been extended to discounted cumulative gain (DCG). DCG estimates the relevance of a document based on its rank among the retrieved documents. The relevance measure is accumulated from top to bottom, discounting the value of documents at lower ranks. NDCG at k measures DCG for the top k documents, normalizing by the highest possible value for a query; therefore, a perfect ranking yields NDCG equals 1.

NDCG is uniquely useful in applications with a non-binary notion of relevance, e.g: a spectrum of relevance. This makes NDCG comparable across different queries: The NDCG values for all queries can be averaged to reliably evaluate the effectiveness of a ranking algorithm for various information needs across a collection. Given a set of queries $q_j \in Q$ and relevance judgements R_{dj} for a document d :

$$NDCG(Q, k) = \frac{1}{|Q|} \sum_{j=1}^{|Q|} Z_{kj} \sum_{m=1}^k \frac{2^{R_{jm}} - 1}{\log_2(1 + m)} \quad (2.3)$$

where Z_{kj} is the normalization factor.

Chapter 3

Cross-Domain Relevance Transfer with BERT

Our proposed model is based on sentence-level relevance modeling and document re-ranking with BERT. By training BERT as a relevance classifier, we aim to extract valuable semantic matching signals which can be leveraged to re-rank a list of candidate documents retrieved with a term-matching technique such as BM25. We also explore applying cross-domain relevance transfer to exploit models of relevance learned on out-of-domain collections, which is crucial in re-ranking documents that are too long for BERT to directly process. This chapter introduces the datasets that we use and, and describes the details of our BERT-based sentence-level relevance classifier and document re-ranker.

3.1 Datasets

In order to model sentence-level relevance with BERT, we need training pairs of queries and short text, annotated with relevance labels. Fortunately, a number of collections fortuitously contain such relevance judgements at the sentence and passage level, which makes them the ideal choice for training our model. We fine-tune BERT on three such sentence- and passage-level datasets individually and in combination: TREC Microblog, MicroSoft MACHine Reading Comprehension and TREC Complex Answer Retrieval datasets. The details of each dataset are provided below.

Type	Training Set	Validation Set
Number of queries	166	59
Number of tweets	133K	44K

Table 3.1: Statistics about the MB dataset.

3.1.1 TREC Microblog (MB)

The TREC Microblog dataset draws from the Microblog Tracks at TREC from 2011 to 2014, with topics and relevance judgments over tweets. Topics associated with tweets are treated as queries, and each of the four datasets contains 50 queries. The nature of this collection differs from newswire documents that we evaluate our models on in distinct ways: First of all, tweets have much fewer tokens than newswire documents. By definition, tweets are limited to 280 characters. Furthermore, because queries and tweets in this dataset are comparable in length, exact matches of query terms occur less frequently in the tweets than they might in longer documents such as news articles. Therefore, semantic matching signals may take precedence in improving retrieval effectiveness on MB. Related to this point, tweets are expressed in a much less formal language than news articles. Tweets may characteristically contain various abbreviations (partly due to the aforementioned length constraint), informal conventions such as hashtags or typos. Such informal language may result in term mismatches in the case of exact matching. It may therefore be helpful to catch other semantic signals with a deep neural network.

We use the MB data prepared by Rao et al. [53]. We extract the queries, tweets and relevance judgements from the dataset, excluding metadata such as query time and URLs of the tweets. Relevance judgements in MB are reported on a three-point scale where (“irrelevant”, “relevant” and “highly relevant”); however, for the purposes of this work we treat both higher degrees of relevance as equal [44]. Both queries and tweets are segmented

<p>Query: bbc world service staff cuts</p> <p>Text: irish times : bbc world service confirms cuts : the bbc world service will shed around 650 jobs or more</p> <p>Relevance: 1 (“relevant”)</p>

Figure 3.1: A sample relevant query and tweet pair from the MB dataset.

into token sequences. We sample 25% of the data for the validation set, and use the rest for fine-tuning BERT. We experiment with different splits as discussed in Chapter 5, and find this split to be ideal.

3.1.2 MicroSoft MACHine Reading Comprehension (MS MARCO)

MS MARCO is a large-scale machine reading comprehension and question answering dataset that is extensively used in the NLP community. MS MARCO [5] features user queries sampled from Bing’s search logs and passages extracted from web documents. The dataset is composed of tuples of a query with relevant and non-relevant passages. On average, each query has one relevant passage. However, some may have no relevant passage at all as the dataset is constructed from the top-10 passages manually annotated by human judges. Therefore, some relevant passages might not have been retrieved with BM25. MS MARCO can be distinguished from similar datasets by its size and real-world nature. Similar to MB, MS MARCO is representative of a natural, and noisy, distribution

Type	Training Set	Validation Set
Number of queries	809K	6.9K
Number of passages	12.M	6.9M

Table 3.2: Statistics about the MS MARCO dataset.

<p>Query: is a little caffeine ok during pregnancy</p> <p>Relevant Passage: We don’t know a lot about the effects of caffeine during pregnancy on you and your baby. So it’s best to limit the amount you get each day. If you’re pregnant, limit caffeine to 200 milligrams each day. This is about the amount in 1.5 8-ounce cups of coffee or one 12-ounce cup of coffee.</p> <p>Non-relevant Passage: It is generally safe for pregnant women to eat chocolate because studies have shown to prove certain benefits of eating chocolate during pregnancy. However, pregnant women should ensure their caffeine intake is below 200 mg per day.</p>

Figure 3.2: A sample relevant and non-relevant passage pair for a query from the MB dataset

Type	Training Set	Validation Set
Number of queries	3M	700K
Number of passages	30M	7M

Table 3.3: Statistics about the CAR dataset.

of information needs, unlike other datasets that often contain high-quality text that may not reflect the use in real life.

Here we focus on the passage-ranking task of MS MARCO. Following the settings in Nogueira et al. [43], we train BERT on approximately 12.8M training samples. The development set is composed of approximately 6.9k queries, each paired with the top 1000 most relevant passages in the MS MARCO dataset as retrieved with BM25. Similarly, the evaluation set contains approximately 6.8 queries and their top 1000 passages, but without the relevance annotations.

3.1.3 TREC Complex Answer Retrieval (CAR)

TREC CAR [19] uses paragraphs extracted from all paragraphs in the English Wikipedia, except the abstracts. Each query is formed by concatenating an article title and a section heading, with all passages under that section considered relevant. The goal of this TREC track is to automatically collect and condense information for a complex query into a single coherent summary. Rather than focusing on document retrieval, the priority is aggregating synthesized information in the form of references, facts and opinions. However, CAR is a synthetic dataset in the sense that queries and documents do not reflect real-world distributions or information needs. The organizers of TREC CAR 2017 only provide manual annotations for the top-5 passages retrieved, meaning some relevant passages may not be annotated if they rank lower. For this reason, we opt to use automatic annotations that provide relevance judgements for all possible query-passage pairs.

The dataset has five predefined folds over the queries. Paragraphs corresponding to the first four folds are used to construct the training set consisting of approximately 3M queries, and the rest the validation set of around 700K queries. The original test set used to evaluate submissions to TREC CAR is used for testing purposes. A subtle detail to note is that the official BERT models are pre-trained on the entire Wikipedia dump; therefore, they have also been trained on documents in the TREC CAR test collection albeit in an unsupervised fashion. In order to avoid the leak of test data into training, we

use the BERT model pre-trained only on the half of Wikipedia present in CAR training samples [43]. The training pairs for CAR are generated by retrieving the top 10 passages from the entire CAR corpus with BM25.

3.2 Modeling Relevance with BERT

We propose modeling sentence-level and passage-level relevance with BERT to capture semantic signals helpful for relevance prediction. This approach is motivated by the application of transfer learning in NLP where a large transformer model trained for language modeling can be used for various downstream tasks. In our implementation, we choose BERT as our base model. BERT is trained on copious amounts of unsupervised data from the Google BookCorpus and English Wikipedia with masked language modeling. Although the training procedure doesn’t involve any explicit objective to extract linguistic features, it has been shown to implicitly recognize such features as subject-verb agreement and conference resolution [26, 13], which allows a number of NLP tasks to greatly benefit from features implicitly encoded in BERT weights.

3.2.1 Relevance Classifier

The core of our model is a BERT-based sentence-level relevance classifier. In other words, we build a model on top of BERT to predict a relevance score s_i for a sentence or passage d_i given a query q . Because the maximum input length that BERT can handle is 512 tokens, we limit our training data to sentence-level and passage-level datasets. In other words, d_i are either tweets drawn from TREC Microblog or passages from MS MARCO or TREC CAR. Following Nogueira et al. [43], we frame relevance modeling as a binary classification task. More specifically, we feed query-text pairs into the BERT model with their respective relevance judgements (i.e: 0 for non-relevant and 1 for relevant). Through this training process BERT learns to assign a relevance score to unseen text for a given query. The details of the input representation to BERT and specifics of fine-tuning BERT for relevance prediction are discussed at length in the remainder of this chapter.

Input Representation

We form the input to BERT by concatenating the query q and a sentence d into the sequence $[[CLS], q, [SEP], d, [SEP]]$. The $[SEP]$ metatoken is used to distinguish between two non-consecutive token sequences in the input, i.e: query and text, and the $[CLS]$ signifies a

special symbol for classification output. Although BERT supports variable length token sequences, the final input length must be consistent across each batch. Therefore, we pad each sequence in a mini-batch to the maximum length in the batch.

The complete input embeddings of BERT is comprised of token, segmentation and position embeddings. The first is constructed by tokenizing the above sequence with the proper metatokens in place with the BERT tokenizer. Since BERT was trained based on WordPiece tokenization, we use the same tokenizer to achieve optimal performance. WordPiece tokenization may break words into multiple subwords in order to more efficiently deal with out-of-vocabulary words and better represent complex words. During training, the subwords derived with WordPiece tokenization are reconstructed based on the training corpus. After tokenization, each token in the input sequence is converted into token IDs corresponding to the index in BERT’s vocabulary. Tokens that do not exist in the vocabulary are represented with a special [UNK] token.

The segment embeddings indicate the start and end of each sequence, whether it be a single sequence or a pair. For relevance classification where we have two texts in the input sequence, i.e: query and sentence, the segment embeddings corresponding to the tokens of the first sequence, i.e: the query, are all 0’s, and those for the second sequence, i.e: the document, are all 1’s. The position embeddings are learned for sequences up to 512 tokens, and help BERT recognize the relative position of each token in the sequence. The input representation for a sample short query-sentence pair is shown in Figure 3.3.

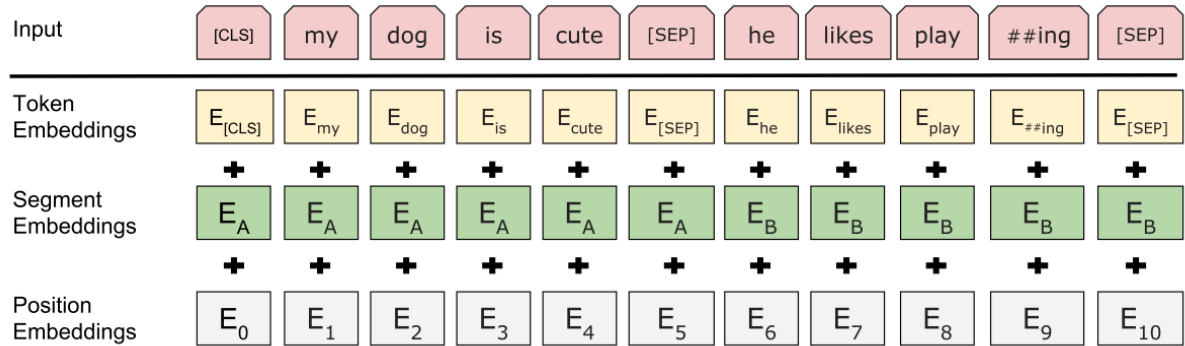


Figure 3.3: Illustration of BERT input representation adapted from Devlin et al. [18].

Fine-Tuning

A variety of useful deep semantic features are already encoded in pretrained BERT weights. It is thus possible to fine-tune BERT for a specific downstream task with less data and time by adding a fully connected layer on top of the network. Intuitively, the lower layers of the network have already been trained to capture features relevant to the task.

To fine-tune BERT for relevance modeling, we add a single layer neural network on top of BERT for classification. This layer consists of $K \times H$ randomly initialized neurons where K is the number of classifier labels and H is the hidden state size. For relevance classification, we have two labels indicating whether the sentence is relevant or non-relevant for the given query ($K = 2$).

The final hidden state corresponding to the first token, i.e: [CLS], provides a H -dimensional aggregate representation of the input sequence that can be used for classification. We feed the final hidden state in the model corresponding to [CLS] into the classification layer. The probability that the sentence d_i is relevant to the query q_i is thus computed with standard softmax:

$$\sigma(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}} \quad (3.1)$$

where $\sigma(y_i)$ maps the arbitrary real value y_i into a probability distribution. Intuitively, $\sigma(y_i)$ represents the relevance score for the sentence d_i . The parameters of BERT and the additional softmax layer are optimized jointly to maximize the log-probability of the correct label with cross-entropy loss.

3.3 Reranking with BERT

Fine-tuning BERT on relevance judgements of query-text pairs allows us to obtain a model of relevance so that we can compute sentence-level relevance scores easily on any collection. However, recall that we train BERT on sentence-level or passage-level datasets so as not to exceed the maximum input size of BERT. These training datasets come from very different distributions than the test collections introduced in Chapter 5. In order to predict relevance on much longer newswire documents, we explore cross-domain relevance transfer by using models trained on MB, MS MARCO and CAR on newswire collections. Our hypothesis is that if a neural network with a large capacity such as BERT can capture relevance in one domain, the model of relevance might successfully transfer to other domains.

To apply cross-domain relevance transfer, we retrieve relevant documents from the collection to depth 1000 with BM25 and split each document into its constituent sentences to match the input size of BERT. We then run inference over the sentences with our models fine-tuned on out-of-domain datasets to compute a score for each sentence. We determine overall document scores by combining exact and semantic matching signals. Based on BM25+RM3 document scores we know a ranking of documents based on exact matches of query terms. Sentence-level scores obtained with BERT reveal other implicit semantic information not evident to BM25. By combining the two sets of relevance matching signals, we establish a more diverse notion of relevance, leading to a better ranking of documents.

Therefore, to determine overall document relevance, we combine the top n scores with the original document score as follows:

$$s_f = a \cdot s_{doc} + (1 - \alpha) \cdot \sum_{i=1}^n w_i \cdot s_i \quad (3.2)$$

where s_{doc} is the original document score and s_i is the i -th top scoring sentence according to BERT. In other words, the relevance score of a document comes from the combination of a document-level term-matching score and relevance evidence contributions from the top sentences in the document as determined by BERT. The parameters α and the w_i 's in Equation 3.2 can be tuned via cross-validation.

3.4 Experimental Setup

3.4.1 Training and Inference with BERT

We fine-tune BERT_{Large} [18] on the datasets introduced earlier in this section: TREC Microblog, MS MARCO AND TREC CAR. In our implementation we adopt BERT_{Large}'s `BertForNextSentencePrediction` interface from the Huggingface `transformers` (previously known as `pytorch-pretrained-bert`) library¹ as our base model. The maximum sequence length, i.e: 512 tokens, is used for BERT in all our experiments.

The fine-tuning procedure introduces few new hyperparameters in addition to those already used in pre-training: batch size, learning rate, and number of training epochs. Due to the large amount of training data in MS MARCO and TREC CAR, BERT is initially trained on Google's TPU's with a batch size of 32 for 400k iterations. We use Adam [29]

¹<https://github.com/huggingface/transformers>

with an initial learning rate of 3×10^{-6} , $\beta_1 = 0.9$ and $\beta_2 = 0.999$ and L2 decay of 0.01. Learning rate warmup is applied over the first 10k steps with linear decay of learning rate. We apply dropout with probability of 0.1 across all layers.

We train all other models using cross-entropy loss for 5 epochs with a batch size of 16. We conduct all our experiments on NVIDIA Tesla P40 GPUs with PyTorch v1.2.0. We use Adam [29] with an initial learning rate of 1×10^{-5} , linear learning rate warmup at a rate of 0.1 and decay of 0.1. We find that applying diminishing learning rates is especially important in fine-tuning BERT in order to preserve the information encoded in the original BERT weights and speed up training.

3.4.2 Evaluation

We retrieve an initial ranking of 1000 documents for each query in Robust04, Core17 and Core18 using the open-source Anserini information retrieval toolkit based on Lucene 8. To ensure fairness across all three collections, we use BM25 with RM3 query expansion with default parameters. Before running inference with BERT to obtain relevance scores, we preprocess the retrieved documents: First, we clean the documents by stripping any HTML/XML tags and split each document into its constituent sentences with NLTK’s Stanford Tokenizer². If the length of a sentence with the meta-tokens still exceeds the maximum input length of BERT, we further segment the spans into fixed sized chunks.

We experiment with the number of top scoring sentences to consider while computing the overall score, and find that using only the top 3 sentences is often enough. In general, considering any more doesn’t yield better results. To tune hyperparameters in Equation 3.2, we apply five-fold cross-validation over TREC topics. For Robust04, we follow the five-fold cross-validation settings in Lin [32] over 250 topics. For Core17 and Core18 we similarly apply five-fold cross validation. We learn parameters α and the w_i on four folds via exhaustive grid search with $w_1 = 1$ and varying $a, w_2, w_3 \in [0, 1]$ with a step size 0.1, selecting the values that yield the highest AP on the remaining fold. We report retrieval effectiveness in terms of AP, P@20, and NDCG@20.

²<https://nlp.stanford.edu/software/tokenizer.shtml>

Chapter 4

Architecture

We apply BERT to document retrieval via integration with the open-source information retrieval toolkit Anserini¹. The proposed architecture of our system follows a two-stage pipeline as shown in Figure 4.1: Anserini is used to retrieve documents from indexed collections with BM25, forming an initial candidate list. Our model introduced earlier in Section 3 is deployed as a re-ranker over the candidate documents to produce a final ranking of documents based on sentence-level relevance. In this section we review each

¹<http://anserini.io/>

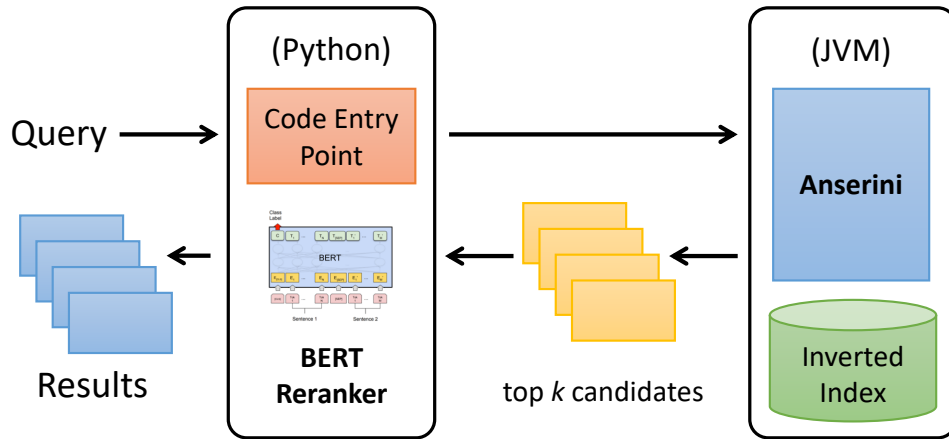


Figure 4.1: Architecture of our system featuring tight integration between Python and the JVM.

component of the architecture and discuss the design choices behind their integration. We also touch upon the issue of reproducibility in information retrieval and our efforts to make our work more reproducible as well as their current limitations.

4.1 Anserini

Technology transfer between the academic and industrial information retrieval communities is at times impeded due to a lack of universal set of tools and infrastructure. Most industry practitioners have adopted Lucene², Solr³ or Elasticsearch⁴ as the de facto platform in the development of search applications with the primary objective of scalability. However, academic systems such as Indri⁵ and Terrier⁶ are far more prevalent among researchers, which prioritize better rankings above all else with little consideration for operational characteristics.

Anserini [63, 64] was developed in response to this disconnect to provide a research-focused information retrieval toolkit on top of the open-source Lucene search library. Anserini facilitates efficient full text indexing and search capabilities over large-scale text collections by providing wrappers and intuitive APIs on top of core Lucene libraries. More importantly, Anserini makes it possible for researchers and industry practitioners alike to systematically evaluate their models over standard test collections in a reproducible and comparable manner.

Related to our work, Anserini can be seamlessly integrated into multi-stage ranking architectures with large improvements in retrieval effectiveness and low latency as demonstrated by Nogueira et al. [43] and Yang et al. [?]. Similar to their work, we initially use Anserini to index our test collections in a multi-threaded manner with Lucene 8.0 (post commit id 75e36f9). For each test collection, we retrieve an initial ranked list of documents 1000 for each query with BM25 using default Anserini parameters.

²<https://lucene.apache.org/>

³<https://lucene.apache.org/solr/>

⁴<https://www.elastic.co/>

⁵<https://www.lemurproject.org/>

⁶<http://terrier.org/>

4.2 Main Module

The main Python module lies at the core of our system, encompassing the preprocessing, training / inference and evaluation components. All the functionalities of our proposed model discussed in Chapter 3 are implemented in this module in Python using the deep learning framework PyTorch⁷.

The preprocessing component of this module consumes the documents retrieved with Anserini, and converts them into a format that can be used by the main component that enables training and inference with BERT. On the one hand, the main component can be used to train BERT as a relevance classifier. This functionality may be used independently of the overall pipeline to fine-tune BERT on new collections. On the other hand, we can run inference over the output of the preprocessing module with previously trained models, producing a list of sentence relevance scores. Finally, this component also serves as a re-ranker where sentence and BM25 document scores are interpolated in order to compute an overall relevance score for each candidate document. Last but not least, the evaluation component integrates directly with Anserini to assess the retrieval effectiveness of our system.

4.3 Integration

Our two-stage pipeline marries NLP and IR capabilities to implement a document retrieval system that successfully leverages semantic cues in documents. For an effective integration, we need to address the technical challenge of connecting the two components that have fundamentally different infrastructural requirements. In this section we discuss the design choices in bridging the worlds of NLP and IR from a software engineering viewpoint.

Anserini, which is responsible for indexing and retrieval in our system, runs on the Java Virtual Machine (JVM) as it is mostly implemented in Java or provides Python wrappers on Java. However, our deep learning framework of choice PyTorch, similar to alternatives such as TensorFlow⁸, are implemented in Python with a C++ backend.

There exist two immediate solutions to connecting Python and the JVM. “Loosely-coupled” integration approaches involve using an intermediary medium between Python and the JVM. For example, we may pass text files between the two in order to facilitate communication without direct interaction. However, this is not an efficient solution

⁷<https://pytorch.org/>

⁸<https://www.tensorflow.org/>

as it requires writing / reading potentially large files to / from disk, not to mention the memory requirements. Furthermore, this approach requires diligent monitoring to ensure that changing file formats and APIs do not break code. Integration via REST APIs is plagued with similar issues as passing intermediate text files. Specifically, it may require frequent HTTP calls, thus introducing significant overhead. Additionally, imperfect solutions for enforcing API contracts risk stability of the system. Ultimately, neither approach is suitable for rapid experimentation in a research environment.

Therefore, we explore ways to achieve “tightly-coupled” integration. One solution is to adopt the Java Virtual Machine (JVM) as the primary code entry point, and connect to PyTorch’s C++ backend via the Java Native Interface (JNI). However, this would result in two separate code paths (JVM to C++ for execution and Python to C++ for model development), leading to maintainability issues similar to those mentioned with regard to REST APIs.

For this reason, we finally chose Python as our primary development environment, integrating Anserini using the Pyjnius Python library⁹ for accessing Java classes. Pyjnius was originally developed to facilitate Android development in Python, and allows Python code to directly manipulate Java classes and objects. Thus, our system supports Python as the main development language (and code entry point, as shown in Figure 4.1), connecting to the JVM to access retrieval capabilities of Anserini.

4.4 Replicability and Reproducibility

Over the last decade, it has become increasingly challenging to verify reported results and compare various performance metrics due to growing number of information retrieval systems both in academia and the industry. Unlike some fields of computer science where it is practical to manually corroborate findings or visually inspect results, the amount and type of data involved in document retrieval deems this approach infeasible. This challenge has prompted one of the largest IR conferences in the world, SIGIR, to issue a task force to determine guidelines to establish repeatability, replicability and reproducibility principles in IR projects.¹⁰

The first dimension of this movement, repeatability, emphasises a researcher’s ability to reliably repeat her own runs. The path to this goal is through rigorous logging, good data management practices and consistent use of virtual environments. We don’t delve further

⁹<https://pyjnius.readthedocs.io/>

¹⁰<http://sigir.org/wp-content/uploads/2018/07/p004.pdf>

into the details of repeatability as the practices we follow are universal to all research endeavors.

The second dimension, replicability, highlights the ability of an independent group to obtain the same results using the researcher’s original artifacts. We strive to make our work replicable by building a Docker image to accompany our system that allows anyone to deploy and test our system on any operating system easily. By adhering to the requirements defined in the SIGIR Open-Source IR Replicability Challenge (OSIRRC), we ensure that our system can seamlessly work with their evaluation infrastructure in the future. The OSIRRC jig¹¹ needs to be set up first to run the commands on Docker hub. The OSIRRC Docker container contract includes three “hooks” for interacting with the system: The **init** hook has to be called first, whose purpose is to run any preparatory steps for the retrieval run including downloading and compiling the source code, downloading pre-built artifacts such as JAR files and other external resources such as pretrained models. In our case, we pull the source code, data and pretrained models from Google Cloud Storage buckets; build Anserini with Maven, and the TREC evaluation tool. Next the **index** hook is called to, as the name indicates, build the necessary indexes. Finally, the **search** hook helps perform multiple ad-hoc retrieval runs in a row. Each of the hook scripts accepts a JSON file that defines the various arguments for the respective script such as path to the relevance judgements file.

Finally, the third dimension, reproducibility, refers to the the ability of an independent group of researchers to implement the author’s proposed artifacts from scratch with the same results. This final goal is indeed the hardest to achieve; as a matter of fact, it may even be impossible in certain cases due to non-determinism. Unfortunately, we found this to be true with some aspects of our work with BERT as well. For example, the fine-tuning and inference processes described in Chapter 3 produces slightly different sentence scores (i.e: third decimal) unless they are performed on the same GPU. To further aggravate this issue, these small differences add up over floating point operations, leading to as much as a 0.5 point difference in AP. We try to relieve this issue by releasing our tuned hyperparameters which help reproduce the same results despite minor differences in sentence scores.

¹¹<https://github.com/osirrc/jig>

Chapter 5

Experimental Results

5.1 Datasets

We conduct end-to-end document ranking experiments on three TREC newswire collections: the Robust Track from 2004 (Robust04) and the Common Core Tracks from 2017 and 2018 (Core17 and Core18).

Robust04

Robust04 draws from the set of documents in TREC Disks 4 and 5, spanning news articles from Financial Times and LA Times, except the Congressional Record. The collection comprises 250 topics, with relevance judgments over 500K documents. The goal of the Robust track is to improve the consistency and robustness of retrieval methods by focusing ad hoc search on poorly performing topics [59]. Notably the distribution of document lengths in Robust04 is highly skewed, with the majority of documents having fewer than 200K tokens and a couple of documents having more than 1M tokens. Existing neural re-rankers such as DRMM [22] and KNRM [62] are known to struggle with documents in the tail of the distribution.

Core17 & Core18

Core17 and Core18 build on the TREC 2017 and 2018 Common Core Tracks respectively. The motivation behind these tracks is to build up-to-date test collections based on more

recently created documents. Core17 contains 1.8M articles from the New York Times Annotated Corpus while Core18 has around 600K articles from the TREC Washington Post Corpus. Core17 and Core18 have only 50 topics each, which are drawn from the Robust Track topics. Given their relatively recent release, literature on these collections is still sparse.

5.2 Results and Analyses

Tables 5.1, 5.2 and 5.3 displays our main results on Robust04, Core17 and Core18. The top row (BM25+RM3) corresponds to the BM25 runs with RM3 query expansion using default Anserini parameters. Although higher scores could be obtained on Robust04 with tuned parameters, we present untuned runs for the sake of fairness as no careful tuning has been performed for Core17 or Core18. The remaining five blocks show the retrieval effectiveness of our models trained as described in Chapter 3. The models are labeled to reflect the fine-tuning procedure where the datasets that BERT was trained on are listed in order inside parentheses. The n S preceding the model name indicates that the top n sentences in each document were interpolate to compute an overall document score. The

Model	MAP	P@20	NDCG@20
BM25+RM3	0.2903	0.3821	0.4407
1S: BERT(MB)	0.3408 [†]	0.4335 [†]	0.4900 [†]
2S: BERT(MB)	0.3435 [†]	0.4386 [†]	0.4964 [†]
3S: BERT(MB)	0.3434 [†]	0.4422 [†]	0.4998 [†]
1S: BERT(CAR)	0.3025 [†]	0.3970 [†]	0.4509
2S: BERT(CAR)	0.3025 [†]	0.3970 [†]	0.4509
3S: BERT(CAR)	0.3025 [†]	0.3970 [†]	0.4509
1S: BERT(MS MARCO)	0.3028 [†]	0.3964 [†]	0.4512
2S: BERT(MS MARCO)	0.3028 [†]	0.3964 [†]	0.4512
3S: BERT(MS MARCO)	0.3028 [†]	0.3964 [†]	0.4512
1S: BERT(CAR → MB)	0.3476 [†]	0.4380 [†]	0.4988 [†]
2S: BERT(CAR → MB)	0.3470 [†]	0.4400 [†]	0.5015 [†]
3S: BERT(CAR → MB)	0.3466 [†]	0.4398 [†]	0.5014 [†]
1S: BERT(MS MARCO → MB)	0.3676 [†]	0.4610 [†]	0.5239 [†]
2S: BERT(MS MARCO → MB)	0.3697[†]	0.4657 [†]	0.5324 [†]
3S: BERT(MS MARCO → MB)	0.3691 [†]	0.4669[†]	0.5325[†]

Table 5.1: Ranking effectiveness on Robust04.

Model	MAP	P@20	NDCG@20
BM25+RM3	0.2823	0.5500	0.4467
1S: BERT(MB)	0.3091 [†]	0.5620	0.4628
2S: BERT(MB)	0.3137 [†]	0.5770	0.4781
3S: BERT(MB) 0.3154 [†]	0.5880	0.4852 [†]	
1S: BERT(CAR)	0.2814 [†]	0.5500	0.4470
2S: BERT(CAR)	0.2814 [†]	0.5500	0.4470
3S: BERT(CAR)	0.2814 [†]	0.5500	0.4470
1S: BERT(MS MARCO)	0.2817 [†]	0.5500	0.4468
2S: BERT(MS MARCO)	0.2817 [†]	0.5500	0.4468
3S: BERT(MS MARCO)	0.2817 [†]	0.5500	0.4468
1S: BERT(CAR → MB)	0.3103 [†]	0.5830	0.4758
2S: BERT(CAR → MB)	0.3140 [†]	0.5830	0.4817 [†]
3S: BERT(CAR → MB)	0.3143 [†]	0.5830	0.4807
1S: BERT(MS MARCO → MB)	0.3292 [†]	0.6080 [†]	0.5061 [†]
2S: BERT(MS MARCO → MB)	0.3323[†]	0.6170 [†]	0.5092[†]
3S: BERT(MS MARCO → MB)	0.3314 [†]	0.6200[†]	0.5070 [†]

Table 5.2: Ranking effectiveness on Core17.

Model	MAP	P@20	NDCG@20
BM25+RM3	0.3135	0.4700	0.4604
1S: BERT(MB)	0.3393 [†]	0.4930	0.4848 [†]
2S: BERT(MB)	0.3421 [†]	0.4910	0.4857 [†]
3S: BERT(MB)	0.3419 [†]	0.4950 [†]	0.4878 [†]
1S: BERT(CAR)	0.3120	0.4680	0.4586
2S: BERT(CAR)	0.3116	0.4670	0.4585
3S: BERT(CAR)	0.3113	0.4670	0.4584
1S: BERT(MS MARCO)	0.3121	0.4670	0.4594
2S: BERT(MS MARCO)	0.3121	0.4670	0.4594
3S: BERT(MS MARCO)	0.3121	0.4670	0.4594
1S: BERT(CAR → MB)	0.3385 [†]	0.4860	0.4785
2S: BERT(CAR → MB)	0.3386 [†]	0.4810	0.4755
3S: BERT(CAR → MB)	0.3382 [†]	0.4830	0.4731
1S: BERT(MS MARCO → MB)	0.3486 [†]	0.4920	0.4953[†]
2S: BERT(MS MARCO → MB)	0.3496 [†]	0.4830	0.4899 [†]
3S: BERT(MS MARCO → MB)	0.3522[†]	0.4850	0.4899 [†]

Table 5.3: Ranking effectiveness on Core18.

main result tables also highlights statistically significant results based on paired t -tests compared to the BM25+RM3 baseline with a \dagger . We report significance at the $p < 0.01$ level, with appropriate Bonferroni corrections for multiple hypothesis testing.

5.2.1 Effect of Training Data

By fine-tuning BERT on three different datasets alone and in combination, we hope to study the effect of nature and amount of training data on the power of our learned relevance matching model. As seen in Tables 5.1, 5.2 and 5.3, the particular source of relevance judgements that we train BERT on substantially influences retrieval effectiveness across all three test collections.

First of all, we find that fine-tuning BERT on MB alone (BERT(MB)) significantly outperforms the BM25+RM3 baseline for all metrics on Robust04. We also observe significant increases in AP on Core17 and Core18, as well as in P@20 and NDCG@20 in some cases. These results confirm that relevance models learned from tweets can be successfully transferred to news articles in spite of the considerable differences in domain. This surprising finding may be attributed to the relevance matching power enabled with deep semantic information learned by BERT.

To see if BERT would be able to extract more useful information from a larger number of tweets, we experiment with different training-dev splits for MB. We first reserve 67% of the entire dataset for training, which leads to AP 1-1.5 points lower than the reported results. While using 82% and 90% of the data as the training set produces better results, we find that effectiveness peaks when we use exactly 75% of all data, showing that more data is not always better to learn an effective relevance model. We report performance of models trained on 75% of all tweets and validated on the remaining 25%.

Contrary to the large gains brought by fine-tuning on MB, fine-tuning on CAR or MS MARCO alone results in marginal gains over the baseline on Robust04. Re-ranking with these models in fact hurts effectiveness on Core17 and Core18 for most metrics. The synthetic nature of CAR data especially does not appear to be useful for relevance modeling on newswire collections. For instance, using BERT(CAR) gives 0.3120 AP on Core18 compared to the 0.3135 AP of the baseline, which means that the model actually fails to score relevant sentences highly. As the 2S and 3S results for the same model show, the effectiveness on Core18 in fact progressively degrades the more sentences are considered in final score aggregation. Intuitively, these results indicate that using these models incorrectly disrupts the better order imposed by the baseline.

Results for BERT(MS MARCO) are more surprising since the web passages in the MS MARCO dataset are “closer” to the news articles in the test collections than MB. Given the proximity of the domains, it would be reasonable to expect relevance transfer between MS MARCO and newswire collections to be more effective. However, our results reveal that this is not necessarily the case, and fine-tuning on MS MARCO alone is far less effective for relevance transfer than fine-tuning on MB alone.

Although fine-tuning on only CAR or MS MARCO does not yield large improvements, we actually obtain considerably higher results by fine-tuning these models further on MB. With BERT(CAR \rightarrow MB) we achieve effectiveness that is slightly better than fine-tuning on MB alone in some cases. We hypothesize that CAR might have a similar effect to language pre-training in that it doesn’t directly apply to the downstream document retrieval task, but provides a better representation that can benefit from fine-tuning on MB. More surprisingly, fine-tuning on MS MARCO first and then on MB represents our best model (BERT(MS MARCO \rightarrow MB)) as shown in the final block of the table. This model is able to exploit data from both MS MARCO and MB, with a score that is higher than fine-tuning on each dataset alone.

5.2.2 Number of Sentences

We consider up to three top scoring sentences in each document to re-rank documents. Our main results suggest that the top scoring sentence by itself is often a good indicator of overall document relevance; in fact we achieve the best AP in about half of the experiments by only considering the most relevant sentence of the document. This finding is consistent with the results of [71] who found through user studies that the most relevant sentence or paragraph in a document provides a reliable proxy for document relevance.

Considering the next most relevant sentence in addition to the top sentence yields a noticeable increase in some of the experiments, such as BERT(MS MARCO \rightarrow MB) on Robust04 and Core18. However, we find that adding a third in fact causes effectiveness to drop in some cases. Preliminary experiments show that in general looking beyond the top three sentence doesn’t help effectiveness. These results suggest that document ranking may be distilled into relevance prediction primarily at the sentence level, and pave the way for new directions in this task.

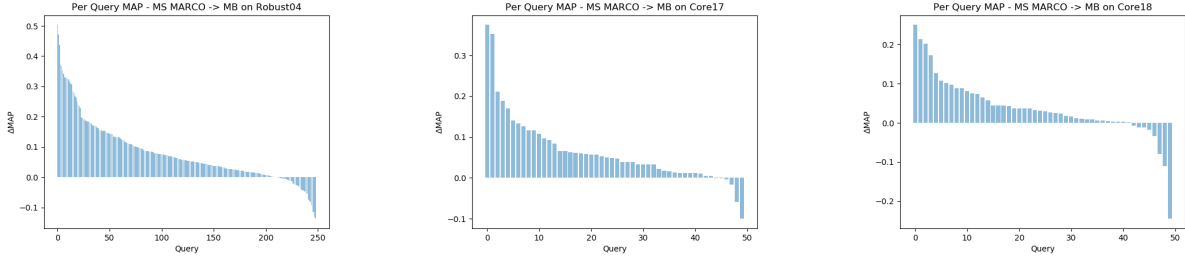


Figure 5.1: Per-query difference in AP between BERT(MS MARCO \rightarrow MB) and the BM25+RM3 baseline on Robust04, Core17 and Core18.

5.2.3 Comparison to Other Ranking Models

In this section we assess our results in the broader context of document re-ranking literature. The meta-analysis¹ of over 100 papers up until 2019 by [66] currently provides the most thorough overview of related work on Robust04. Following the same cross-validation settings as in this thesis to re-rank a strong BM25 baseline, they report the most effective neural model to be DRMM [22] at 0.3152 AP and 0.4718 NDCG@20. In comparison, with our best model BERT(MS MARCO \rightarrow MB) we report the highest AP that we are aware of at 0.3697. Furthermore, our results also exceed the previous highest known score of 0.3686, which is a non-neural method based on ensembles [14].

More recently, [?] reported 0.5381 NDCG@20 on Robust04 by integrating contextualized word embeddings into existing neural ranking models. Our best NDCG@20 on Robust04 at 0.5325 approaches their results even though we optimize for AP instead of NDCG@20. Furthermore, since we are only using Robust04 data for hyperparameter tuning in Eq (3.2), and not for fine-tuning BERT itself, it is less likely that we are overfitting.

Our best model also achieves a higher AP on Core17 than the best TREC submission that does not make use of past labels or human intervention (`umass_baselnrm`, 0.275 AP) [1]. Under similar conditions, we beat every TREC submission in Core18 as well (with the best run being `uwmrgr`, 0.276 AP) [2]. Core17 and Core18 are relatively new and thus have yet to receive much attention from researchers, but to our knowledge, these figures represent the state of the art.

5.2.4 Per-Query Analysis

Our results in Table 5.1 serve as an overview of the effect of training data on cross-domain relevance transfer. However, they do not reveal much with respect to the particular strengths and weaknesses of each model compared to the baseline. To gain further insight into the characteristics of our models, we analyze the per-query retrieval effectiveness of each model compared to the baseline on Robust04, Core17 and Core18.

Figure 5.1 plot the Δ AP per query between our best model, BERT(MS MARCO \rightarrow MB), and the baseline sorted in descending order. BERT(MS MARCO \rightarrow MB) performs better than the baseline for 83%, 88% and 84% of the queries on Robust04, Core17 and Core18 respectively, highlighting the usefulness of relevance signals learned from MS MARCO and MB with BERT. One of the best performing queries across all three collections is **human stampede** while one of the worst is **flavr savr tomato**. Manual inspection of top scoring sentences for **flavr savr tomato** show that our model matches terms semantically close to “tomato” but not relevant to the query phrase itself while the BM25+RM3 baseline looks for direct query matches, therefore ranking documents correctly.

Repeating the same analysis for our worst model BERT(CAR), we present the results in Figure 5.2. This model outperforms the baseline in 57%, 80% and 34% of the Robust04, Core17 and Core18 queries respectively. Not only does this model perform less accurately for much fewer queries, but it is also less stable than BERT(MS MARCO \rightarrow MB). The queries **polygamy polyandry polygyny** and **women driving in Saudi Arabia** appear to be particularly difficult for BERT based models to correctly match.

¹<https://github.com/lintool/robust04-analysis>

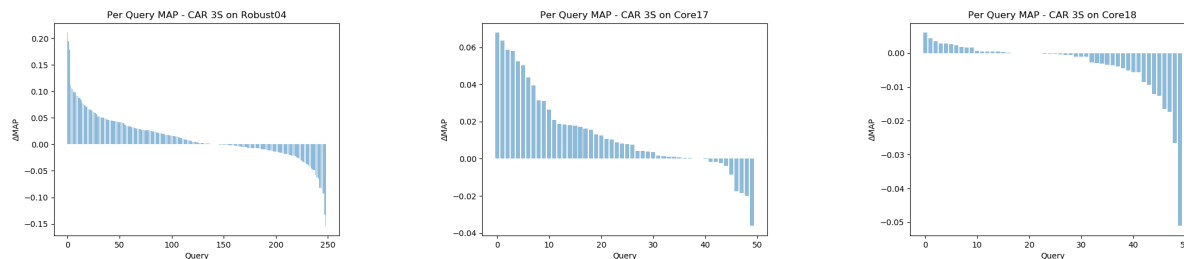


Figure 5.2: Per-query difference in AP between BERT(CAR) and the BM25+RM3 baseline on Robust04, Core17 and Core18.

5.2.5 Effect of Length

Query Length

To investigate the ability of BERT to exploit context-aware representations of the queries, we study the trend of effectiveness across increasing query lengths on Robust04. We conjecture that longer queries would give our models richer context to work with, therefore increasing retrieval effectiveness. To this end, we categorize the 250 queries from the Robust Track by the number of tokens, and calculate the average AP per query length, i.e: from 1 to 5. The majority of the queries (56%) are composed of three tokens, and only a small fraction (5%) have either only a single token or five tokens. This approach is somewhat limited by the narrow range of query lengths, but still gives insight into the effect of query length on each model’s effectiveness.

The average AP per query length for each model is shown in Table 5.4. We only report the best results for each model considering any number of top sentences. The average AP for the BM25+RM3 baseline drops by 24% from a single token query to a query with five tokens, indicating that exact term-matching is more effective for fewer number of terms. Furthermore, the improvements introduced by our models follow a similar pattern to when they are evaluated over all queries as shown in Table ??, with BERT(MS MARCO \rightarrow MB) being the best model across all query lengths. Similar to the baseline, BERT(CAR) and BERT(MS MARCO) experience a gradual decline in AP as the query length increases from 1 to 5, leading to an overall 29% and 28% decrease respectively. Note, however, that the decrease in AP is not continuous at each step but reaches a minimum at query length 3 and a maximum at query length 4.

More interestingly, while increasing query length results in a decrease in AP across

Model	Query Length				
	1	2	3	4	5
BM25+RM3	0.3889	0.2975	0.2760	0.3089	0.2972
BERT(MB)	0.4145	0.3624	0.3230	0.3681	0.4172
BERT(CAR)	0.3978	0.3112	0.2871	0.3288	0.2844
BERT(MS MARCO)	0.3997	0.3117	0.2872	0.3289	0.2871
BERT(CAR \rightarrow MB)	0.4202	0.3603	0.3296	0.3757	0.4926
BERT(MS MARCO \rightarrow MB)	0.4253	0.3885	0.3480	0.4077	0.5627

Table 5.4: Average AP with respect to query length on Robust04.

Model	Robust04		
	AP	P@20	NDCG@20
BERT(CAR[*])	0.3030	0.3980	0.4520
BERT(MS MARCO[*])	0.3300	0.4309	0.4906

Table 5.5: Ranking effectiveness on shortened MS MARCO and CAR evaluated on Robust04.

all the other models, BERT(MB), BERT(CAR \rightarrow MB) and BERT(MS MARCO \rightarrow MB) in fact perform better at query length 5 than at 1. However, due to the fairly small sample size for query length 5 it is not possible to draw generalizable conclusions from this observation. Nevertheless, the effectiveness of the best model BERT(MS MARCO \rightarrow MB) degrades much less where the minimum usually occurs, i.e: query length 3, compared to the baseline, BERT(CAR) and BERT(MS MARCO), showing that they are indeed able to use contextualized query representations most effectively.

Document Length

It is interesting to note that fine-tuning on MS MARCO or CAR alone results in marginal improvements over the baseline, if any, as shown in the second and third blocks of Table ???. Considering any number of top scoring sentences, BERT(CAR) and BERT(MS MARCO) both outperform the BM25+RM3 baseline by 1.2 AP on Robust04 with similar gains in P@20 and NDCG@20. Although still statistically significant, these improvements are much lower than those gained with fine-tuning on MB (5 AP). Moreover, both models in fact perform worse than the baseline on both Core17 and Core18.

We attribute this observation to the length mismatch between the training and evaluation text lengths. After splitting documents into chunks that fit the maximum input that BERT can handle, the average number of tokens in Robust04, Core17 and Core17 sentences is 19, which is fairly close to the the average number of tokens in tweets (15). However, MS MARCO and CAR passages are much longer. It is possible that the difference in document length may be causing issues with relevance transfer across domains.

To validate this hypothesis, we divide MS MARCO and CAR passages into chunks the same size as Robust04 sentences, and train BERT for relevance prediction. The ranking effectiveness of BERT(CAR^{*}) and BERT(MS MARCO^{*}) trained on the shortened data is shown in Table 5.5. For BERT(CAR^{*}), AP only increases by 0.5 upon fine-tuning on

Model	Pruned Robust04		
	AP	P@20	NDCG@20
BM25+RM3	0.2903	0.3821	0.4407
BERT(MB)	0.3031	0.4014	0.4580
BERT(CAR)	0.2959	0.3936	0.4480
BERT(MS MARCO)	0.2962	0.3936	0.4483
BERT(CAR \rightarrow MB)	0.3037	0.3998	0.4527
BERT(MS MARCO \rightarrow MB)	0.3101	0.4102	0.4639

Table 5.6: Retrieval effectiveness on pruned Robust04.

the shortened data, and is still not comparable to BERT(MB). However, the change in document length results in an increase of almost 3 AP for BERT(MS MARCO*), thus approaching BERT(MB) in all metrics. From this finding we infer that while comparable document length is an important consideration for cross-domain relevance transfer, there may be other factors at play that need to be investigated.

5.2.6 Semantic Matching

To isolate the contribution of BERT, we filter sentences in Robust04 that don’t contain any of the query terms. This essentially eliminates the impact of exact matching on the sentence relevance scores, allowing us to verify whether BERT successfully leverages semantic cues in the documents. Table 5.6 displays the retrieval effectiveness of all models on the “pruned” Robust04 dataset. Similar to Table 5.4 we only report the best results for each model.

As expected, filtering Robust04 sentences leads to a decrease in all metrics across all models, which indicates that exact matching signals are still valuable in relevance predictions over long documents. However, notice that the all models still perform better than the baseline; in other words, that they indeed successfully perform semantic matching with notable gains. The improvements over the baseline follow the same trend as the results in Table ?? for all models. While BERT(CAR) and BERT(MS MARCO) yield minor improvements over the baseline when sentences that contain exact query matches are removed, the best performing model BERT(MS MARCO \rightarrow MB) still significantly outperforms the baseline. Furthermore, the drop in AP caused by filtering sentences is also the

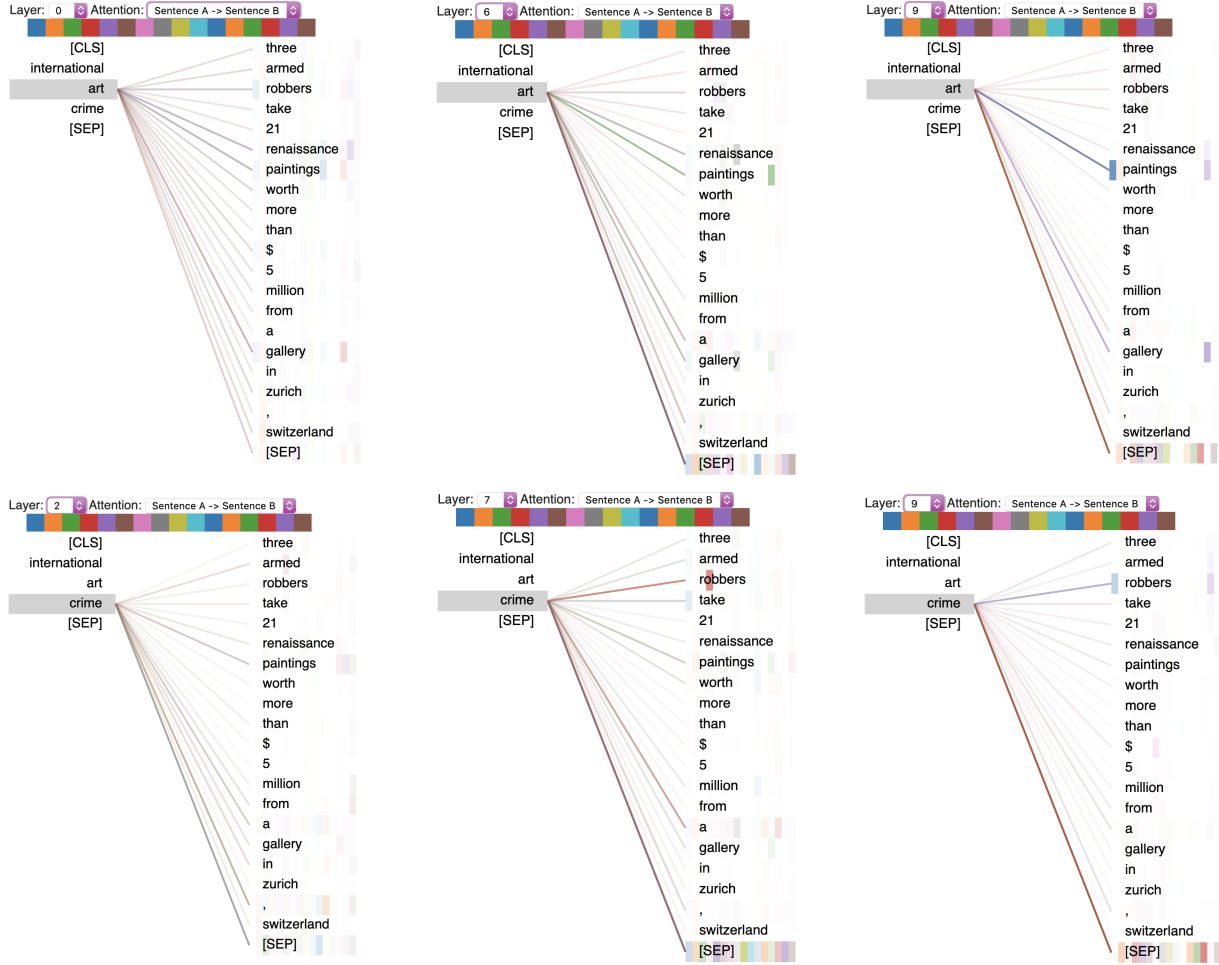


Figure 5.3: Attention visualizations of BERT(MS MARCO \rightarrow MB) for a sentence with a high BERT score for the query international art crime.

highest for the best performing model, indicating that this model is able to exploit both exact and semantic matching signals. The overall effectiveness of this model on the original Robust04 dataset may be owing to the joint relevance matching power demonstrated in this experiment.

To further investigate the semantic matching power of BERT, we visualize the attention of our models at different layers. Having been pretrained on massive amounts of data for language modeling, we expect BERT to capture various semantic relationships helpful for relevance prediction. Figure 5.3 visualizes the attention of our best model BERT(MS MARCO \rightarrow MB) for one of the top scoring sentences for the query 322 **international art crime**. Note that the sentence does not contain any exact matches with the query terms and is nonetheless relevant to the query. The top three images illustrate the attention for the query term “art” at increasingly deeper layers (0, 6 and 9) and the bottom three for the query term “crime”. We see that the model attends to related terms such as “renaissance” and “paintings” for the term “art”, and “robbers” for crime, therefore giving this particular sentence a high relevance score. On the other hand, these terms are not recognized by term-matching techniques like BM25 as relevant. This highlights how the semantic knowledge captured by BERT directly help with relevance modeling on newswire articles.

Chapter 6

Conclusion and Future Work

In this thesis, we propose two innovations to successfully apply BERT to document retrieval with significant improvements on three TREC newswire collections: Robust04, Core17 and Core18. First, to overcome the maximum input length restriction imposed by BERT, we focus on integrating sentence-level evidence to re-rank newswire documents. This approach requires sentence-level relevance labels to train BERT for relevance prediction; however, relevance judgements in most test collections are provided only at the document level. We address this challenge by leveraging sentence-level and passage-level relevance judgements fortuitously available in out-of-domain collections. We fine-tune BERT with the goal of capturing cross-domain notions of relevance, which can be used to re-rank the longer documents in the newswire collections.

We show that relevance models learned with BERT can indeed be transferred across domains in a straightforward manner. Combined with sentence-level relevance modeling, our simple model achieves state-of-the-art results across all three test collections. Furthermore, our results suggest that judging only a small number of most relevant sentences in a document may be sufficient for effective document re-ranking. Through our analyses we investigate the successes and failures of BERT in document re-ranking with respect to training data, query and document length, and semantic matching. Our findings illustrate the relevance matching power of deep semantic information learned by BERT. However, the analyses conducted in this thesis do not come close to fully understanding the effect of BERT in relevance matching. A promising future direction based on our findings includes extended analyses of cross-domain transfer transfer and comparison of exact and semantic matching signals involved in re-ranking documents.

References

- [1] James Allan, Donna Harman, Evangelos Kanoulas, Dan Li, and Christophe Van Gysel. Trec 2017 common core track overview.
- [2] James Allan, Donna Harman, Evangelos Kanoulas, and Ellen Voorhees. TREC 2018 Common Core Track overview. In *Proceedings of the Twenty-Seventh Text REtrieval Conference (TREC 2018)*, Gaithersburg, Maryland, 2018.
- [3] Rie Kubota Ando and Tong Zhang. A framework for learning predictive structures from multiple tasks and unlabeled data. *Journal of Machine Learning Research*, 6(Nov):1817–1853, 2005.
- [4] N. Asadi and J. Lin. Effectiveness/efficiency tradeoffs for candidate generation in multi-stage retrieval architectures. In *SIGIR*, pages 997–1000, 2013.
- [5] Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew McNamara, Bhaskar Mitra, Tri Nguyen, Mir Rosenberg, Xia Song, Alina Stoica, Saurabh Tiwary, and Tong Wang. MS MARCO: A human generated MACHine Reading COMprehension dataset. *arXiv:1611.09268v3*, 2018.
- [6] Jaroslaw Baliński and Czesław Daniłowicz. Re-ranking method based on inter-document distances. *Information processing & management*, 41(4):759–775, 2005.
- [7] Yonatan Belinkov, Nadir Durrani, Fahim Dalvi, Hassan Sajjad, and James Glass. What do neural machine translation models learn about morphology? *arXiv preprint arXiv:1704.03471*, 2017.
- [8] David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.

- [9] John Blitzer, Ryan McDonald, and Fernando Pereira. Domain adaptation with structural correspondence learning. In *Proceedings of the 2006 conference on empirical methods in natural language processing*, pages 120–128. Association for Computational Linguistics, 2006.
- [10] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146, 2017.
- [11] Christopher JC Burges. From ranknet to lambdarank to lambdamart: An overview.
- [12] Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and Hang Li. Learning to rank: from pairwise approach to listwise approach. In *Proceedings of the 24th international conference on Machine learning*, pages 129–136. ACM, 2007.
- [13] Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D Manning. What does bert look at? an analysis of bert’s attention. *arXiv preprint arXiv:1906.04341*, 2019.
- [14] Gordon V. Cormack, Charles L A Clarke, and Stefan Buettcher. Reciprocal rank fusion outperforms condorcet and individual rank learning methods. In *Proceedings of the 32Nd International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR ’09, pages 758–759, New York, NY, USA, 2009. ACM.
- [15] Zhuyun Dai, Chenyan Xiong, Jamie Callan, and Zhiyuan Liu. Convolutional neural networks for soft-matching n-grams in ad-hoc search. In *Proceedings of the eleventh ACM international conference on web search and data mining*, pages 126–134. ACM, 2018.
- [16] Scott Deerwester, Susan T Dumais, George W Furnas, Thomas K Landauer, and Richard Harshman. Indexing by latent semantic analysis. *Journal of the American society for information science*, 41(6):391–407, 1990.
- [17] Mostafa Dehghani, Hamed Zamani, Aliaksei Severyn, Jaap Kamps, and W Bruce Croft. Neural ranking models with weak supervision. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 65–74. ACM, 2017.
- [18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings*

of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota, June 2019.

- [19] Laura Dietz, Manisha Verma, Filip Radlinski, and Nick Craswell. TREC Complex Answer Retrieval overview. In *Proceedings of the Twenty-Sixth Text REtrieval Conference (TREC 2017)*, Gaithersburg, Maryland, 2017.
- [20] Debasis Ganguly, Dwaipayan Roy, Mandar Mitra, and Gareth JF Jones. Word embedding based generalized language model for information retrieval. In *Proceedings of the 38th international ACM SIGIR conference on research and development in information retrieval*, pages 795–798. ACM, 2015.
- [21] Kristina Gulordava, Piotr Bojanowski, Edouard Grave, Tal Linzen, and Marco Baroni. Colorless green recurrent networks dream hierarchically. *CoRR*, abs/1803.11138, 2018.
- [22] Jiafeng Guo, Yixing Fan, Qingyao Ai, and W Bruce Croft. A deep relevance matching model for ad-hoc retrieval. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*, pages 55–64. ACM, 2016.
- [23] Jiafeng Guo, Yixing Fan, Qingyao Ai, and W. Bruce Croft. A deep relevance matching model for ad-hoc retrieval. *CoRR*, abs/1711.08611, 2017.
- [24] Christophe Van Gysel, Maarten de Rijke, and Evangelos Kanoulas. Neural vector spaces for unsupervised information retrieval. *ACM Trans. Inf. Syst.*, 36(4):38:1–38:25, June 2018.
- [25] Po-Sen Huang, Xiaodong He, Jianfeng Gao, Li Deng, Alex Acero, and Larry Heck. Learning deep structured semantic models for web search using clickthrough data. In *Proceedings of the 22nd ACM international conference on Information & Knowledge Management*, pages 2333–2338. ACM, 2013.
- [26] Ganesh Jawahar, Benoît Sagot, Djamé Seddah, et al. What does bert learn about the structure of language?
- [27] Thorsten Joachims. Optimizing search engines using clickthrough data. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 133–142. ACM, 2002.
- [28] K Sparck Jones, Steve Walker, and Stephen E. Robertson. A probabilistic model of information retrieval: development and comparative experiments: Part 2. *Information processing & management*, 36(6):809–840, 2000.

- [29] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv:1412.6980*, 2014.
- [30] Kyung-Soon Lee, Young-Chan Park, and Key-Sun Choi. Re-ranking model based on document clusters. *Information processing & management*, 37(1):1–14, 2001.
- [31] Ping Li, Qiang Wu, and Christopher J Burges. Mcrank: Learning to rank using multiple classification and gradient boosting. In *Advances in neural information processing systems*, pages 897–904, 2008.
- [32] Jimmy Lin. The neural hype and comparisons against weak baselines. In *ACM SIGIR Forum*, volume 52, pages 40–51. ACM, 2019.
- [33] Yang Lingpeng, Ji Donghong, and Tang Li. Document re-ranking based on automatically acquired key terms in chinese information retrieval. In *Proceedings of the 20th international conference on Computational Linguistics*, page 480. Association for Computational Linguistics, 2004.
- [34] Tal Linzen, Emmanuel Dupoux, and Yoav Goldberg. Assessing the ability of lstms to learn syntax-sensitive dependencies. *CoRR*, abs/1611.01368, 2016.
- [35] Zachary C Lipton and Jacob Steinhardt. Troubling trends in machine learning scholarship. *arXiv preprint arXiv:1807.03341*, 2018.
- [36] Tie-Yan Liu. *Learning to rank for information retrieval*. Springer Science & Business Media, 2011.
- [37] Sean MacAvaney, Andrew Yates, Arman Cohan, and Nazli Goharian. CEDR: contextualized embeddings for document ranking. *CoRR*, abs/1904.07094, 2019.
- [38] Christopher Manning, Prabhakar Raghavan, and Hinrich Schütze. Introduction to information retrieval. *Natural Language Engineering*, 16(1):100–103, 2010.
- [39] Oren Melamud, Jacob Goldberger, and Ido Dagan. context2vec: Learning generic context embedding with bidirectional lstm. In *Proceedings of the 20th SIGNLL conference on computational natural language learning*, pages 51–61, 2016.
- [40] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013.

- [41] Bhaskar Mitra and Nick Craswell. Neural models for information retrieval. *arXiv preprint arXiv:1705.01509*, 2017.
- [42] Bhaskar Mitra, Fernando Diaz, and Nick Craswell. Learning to match using local and distributed representations of text for web search. In *Proceedings of the 26th International Conference on World Wide Web*, pages 1291–1299. International World Wide Web Conferences Steering Committee, 2017.
- [43] Rodrigo Nogueira and Kyunghyun Cho. Passage re-ranking with BERT. *arXiv:1901.04085*, 2019.
- [44] Iadh Ounis, Craig Macdonald, Jimmy Lin, and Ian Soboroff. Overview of the trec-2011 microblog track.
- [45] Harshith Padigela, Hamed Zamani, and W. Bruce Croft. Investigating the successes and failures of BERT for passage re-ranking. *arXiv:1905.01758*, 2019.
- [46] Harshith Padigela, Hamed Zamani, and W. Bruce Croft. Investigating the successes and failures of BERT for passage re-ranking. *CoRR*, abs/1905.01758, 2019.
- [47] Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.
- [48] Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. *arXiv preprint arXiv:1802.05365*, 2018.
- [49] Yifan Qiao, Chenyan Xiong, Zhenghao Liu, and Zhiyuan Liu. Understanding the behaviors of BERT in ranking. *arXiv:1904.07531*, 2019.
- [50] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training.
- [51] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners.
- [52] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.

- [53] Jinfeng Rao, Wei Yang, Yuhao Zhang, Ferhan Türe, and Jimmy Lin. Multi-perspective relevance matching with hierarchical convnets for social media search. *CoRR*, abs/1805.08159, 2018.
- [54] D Sculley, Jasper Snoek, Alex Wiltschko, and Ali Rahimi. Winner’s curse? on pace, progress, and empirical rigor. 2018.
- [55] Yelong Shen, Xiaodong He, Jianfeng Gao, Li Deng, and Grégoire Mesnil. Learning semantic representations using convolutional neural networks for web search. In *Proceedings of the 23rd International Conference on World Wide Web*, pages 373–374. ACM, 2014.
- [56] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642, 2013.
- [57] Joseph Turian, Lev Ratinov, and Yoshua Bengio. Word representations: a simple and general method for semi-supervised learning. In *Proceedings of the 48th annual meeting of the association for computational linguistics*, pages 384–394. Association for Computational Linguistics, 2010.
- [58] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [59] Ellen M. Voorhees. Overview of the TREC 2004 Robust Track. In *TREC 2004*, pages 52–69, 2004.
- [60] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*, 2018.
- [61] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*, 2016.
- [62] Chenyan Xiong, Zhuyun Dai, Jamie Callan, Zhiyuan Liu, and Russell Power. End-to-end neural ad-hoc ranking with kernel pooling. *CoRR*, abs/1706.06613, 2017.

- [63] P. Yang, H. Fang, and J. Lin. Anserini: Enabling the use of Lucene for information retrieval research. In *SIGIR 2017*, pages 1253–1256, 2017.
- [64] P. Yang, H. Fang, and J. Lin. Anserini: Reproducible ranking baselines using Lucene. *Journal of Data and Information Quality*, 10(4):Article 16, 2018.
- [65] Peilin Yang, Hui Fang, and Jimmy Lin. Anserini: Reproducible ranking baselines using lucene. *J. Data and Information Quality*, 10(4):16:1–16:20, October 2018.
- [66] Wei Yang, Kuang Lu, Peilin Yang, and Jimmy Lin. Critically examining the “neural hype”: Weak baselines and the additivity of effectiveness gains from neural ranking models. In *Proceedings of the 42nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2019)*, pages 1129–1132, Paris, France, 2019.
- [67] Wei Yang, Yuqing Xie, Aileen Lin, Xingyu Li, Luchen Tan, Kun Xiong, Ming Li, and Jimmy Lin. End-to-end open-domain question answering with bertserini. *arXiv preprint arXiv:1902.01718*, 2019.
- [68] Wei Yang, Haotian Zhang, and Jimmy Lin. Simple applications of bert for ad hoc document retrieval. *arXiv preprint arXiv:1903.10972*, 2019.
- [69] Hamed Zamani and W. Bruce Croft. Relevance-based word embedding. *CoRR*, abs/1705.03556, 2017.
- [70] Hamed Zamani, Mostafa Dehghani, W Bruce Croft, Erik Learned-Miller, and Jaap Kamps. From neural re-ranking to neural ranking: Learning a sparse representation for inverted indexing. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, pages 497–506. ACM, 2018.
- [71] Haotian Zhang, Mustafa Abualsaud, Nimesh Ghelani, Mark D Smucker, Gordon V Cormack, and Maura R Grossman. Effective user interaction for high-recall retrieval: Less is more. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, pages 187–196. ACM, 2018.
- [72] Le Zhao and Jamie Callan. Term necessity prediction. In *Proceedings of the 19th ACM international conference on Information and knowledge management*, pages 259–268. ACM, 2010.
- [73] Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Aligning books and movies: Towards story-like visual

explanations by watching movies and reading books. In *Proceedings of the IEEE international conference on computer vision*, pages 19–27, 2015.