

# CoRT: Complementary Rankings from Transformers

Marco Wrzalik     Dirk Krechel

RheinMain University of Applied Sciences, Germany

{firstname.lastname}@hs-rm.de

## Abstract

Many recent approaches towards neural information retrieval mitigate their computational costs by using a multi-stage ranking pipeline. In the first stage, a number of potentially relevant candidates are retrieved using an efficient retrieval model such as BM25. Although BM25 has proven decent performance as a first-stage ranker, it tends to miss relevant passages. In this context we propose CoRT, a simple neural first-stage ranking model that leverages contextual representations from pre-trained language models such as BERT to complement term-based ranking functions while causing no significant delay at query time. Using the MS MARCO dataset, we show that CoRT significantly increases the candidate recall by complementing BM25 with missing candidates. Consequently, we find subsequent re-rankers achieve superior results with less candidates. We further demonstrate that passage retrieval using CoRT can be realized with surprisingly low latencies.

## 1 Introduction

The successful development of neural ranking models over the past few years has rapidly advanced state-of-the-art performance in information retrieval (Guo et al., 2020; Craswell et al., 2020). One key aspect of the success is the exploitation of *query-document interactions* based on token representations from self-supervised language models (LMs) (Devlin et al., 2019; Radford et al., 2019; Pennington et al., 2014). Due to high computational effort, however, these *interaction-focused* approaches are limited to re-ranking scenarios and thus they depend on the effectiveness of first-stage ranking models for candidate retrieval. Term-based retrieval models such as BM25 have proven decent performance in this task, but tend to miss relevant passages. In this context, we propose *COmplementary Rankings from Transformers* (CoRT), a simple neural first-stage ranking model that lever-

ages contextual representations from *transformer-based language models* (Vaswani et al., 2017; Devlin et al., 2019) to complement term-based first-stage rankings. CoRT optimizes an underlying text encoder towards representations that reflect relevance through vector similarity. The model is trained to act complementary to term-based retrieval by using passages from BM25 rankings as negative examples. We study the characteristics of CoRT with four types of experiments based on the MS MARCO dataset. First, we measure various ranking metrics and compare the results with first-stage ranking baselines and competitors. In course of this, we demonstrate the portion of relevant candidates that are added by CoRT. Second, we combine the candidates from CoRT and BM25 with a state-of-the-art re-ranker based on BERT (Nogueira and Cho, 2019) and investigate how many candidates are needed to saturate the ranking quality. Third, we train CoRT with various representation sizes and measure its impact on the first-stage ranking quality. Fourth, we measure the retrieval latencies of CoRT with two retrieval modalities: a distributed exhaustive search on four GPUs and an approximate search based on a graph-based nearest-neighbor index with pruning heuristics (Iwasaki and Miyazaki, 2018). Finally, we build an exemplary end-to-end ranking pipeline using our first-stage ranking to demonstrate its efficiency. Our contribution is a first-stage ranking framework with the potential to improve end-to-end ranking pipelines by adding candidates that term-based retrieval models typically miss. With this it is possible to reduce the total number of re-ranking candidates without hurting end-to-end ranking quality. As a secondary contribution, we provide an open-source implementation<sup>1</sup> that enables other researchers to reproduce our results and test CoRT on other datasets.

<sup>1</sup><https://github.com/lavis-nlp/CoRT>

## 2 Background and Related Work

In this section we describe key concepts of neural ranking and refer to related work. We then present neural first-stage ranking approaches that allow direct comparison with our results.

### 2.1 Key Concepts of Neural Ranking

According to Guo et al. (2016), Neural ranking approaches can be categorized into two types of models depending on the architecture. *Representation-focused* approaches (Huang et al., 2013; Shen et al., 2014; Zamani et al., 2018) produce representations for queries and documents to predict relevancy scores using a simple distance or similarity measure. In this context, exploiting *local interactions* between neighboring terms is a commonly used technique (Shen et al., 2014; Zamani et al., 2018). Models of the *interaction-focused* type exploit interactions between query and document terms (Guo et al., 2016; Xiong et al., 2017; Dai et al., 2018; Nogueira and Cho, 2019). Although this leads to superior ranking quality (Guo et al., 2016, 2020; Qiao et al., 2019), it is computationally much more intensive, since a given query has to be processed together with each potentially relevant document. Hence, this type of neural ranking model is only applicable in a ranking pipeline, where limited numbers of documents are given as potentially relevant candidates. These candidates are selected by an efficient retrieval model that is able to retrieve documents directly from the corpus in a reasonable amount of time. The multi-stage ranking technique is also known as *cascade ranking* and optimizing the configuration of such a pipeline towards maximized efficiency and effectiveness has been extensively studied in the past (Wang et al., 2011; Chen et al., 2017). Many *interaction-focused* neural ranking models employ a dedicated layer to explicitly perform a matching between query and document terms (Guo et al., 2016; Xiong et al., 2017; Dai et al., 2018). Another approach is using the *attention mechanism* (Vaswani et al., 2017), or more specifically, a pretrained transformer encoder such as BERT (Devlin et al., 2019) to exploit both local and query-document interactions (Nogueira and Cho, 2019; Qiao et al., 2019). Recently, some hybrid approaches have been proposed that combine typical representation-focused techniques with interaction-focused approaches to reduce computational cost: Gao et al. (2020) propose a model architecture comprising three modules for document

understanding, query understanding and relevance judging respectively. The *understanding modules* produce token-level representations, which can be cached as usual in representation-focused approaches. The *relevance judging module* uses those cached representations to apply query-document interactions more quickly. Each module is a stack of transformer layers (Vaswani et al., 2017), initialized with weights from BERT. In a related approach, MacAvaney et al. (2020) investigate the relationship between different numbers of dedicated layers of BERT for query-document interactions and measure the resulting speedup that is due to token representation caching, as well as its impact on the end-to-end ranking quality. Khattab and Zaharia (2020) propose a related approach, namely *ColBERT*. The model architecture incorporates an inexpensive max-similarity mechanism to perform token-level query-document interactions. The authors propose to store *token* representations in an *Approximate Nearest Neighbor* (ANN) index to quickly retrieve only those documents that have token representations in the proximity to those of the query. Thus, *ColBERT* can be described as an end-to-end ranking approach that brings its own first-stage retrieval mechanism allowing to perform end-to-end ranking in a reasonable amount of time.

### 2.2 Neural First-stage Ranking

Now we discuss neural ranking approaches that can be used to retrieve passages or documents directly from an entire corpus in a reasonable amount of time and thus qualify for first-stage ranking. Many proposed methods make use of existing infrastructure for sparse bag-of-words retrieval or at least inverted indexing (Manning et al., 2008). Zamani et al. (2018) propose *SNRM*, a representation-focused approach with sparse representations that can be used with an inverted index as if each feature dimensions corresponds to a term in a bag-of-words representation. SNRM uses pretrained *GloVe Word Embeddings* (Pennington et al., 2014) to model soft-matched n-grams which are encoded in sparse representations. Nogueira et al. (2019) predict queries for given documents to expand those documents by corresponding query terms. In their first work, known as *doc2query*, they used a sequence-to-sequence (seq2seq) transformer model (Vaswani et al., 2017). In a subsequent work, Nogueira and Lin (2019) reported large effectiveness gains for their follow-up model *docTTTTTquery* by replac-

ing the seq2seq model with *T5* (Raffel et al., 2020). Another approach aims at predicting optimal document term weights as a function of the term’s context. *DeepCT*, proposed by Dai and Callan (2020), utilizes BERT to predict these context-aware weights based on associated queries in the training data.

Inverted indexing is only applicable to sparse representations. Representation-focused models using dense representations can instead employ an ANN index, which heuristically prunes documents that are unlikely to be in the top- $k$  proximity of the query representation to realize low response latencies (Boytsov et al., 2016; Gysel et al., 2018). Karpukhin et al. (2020) recently used this technique in combination with a fine-tuned BERT encoder for open question answering.

### 3 Proposed Approach

We describe a first-stage ranking model that acts as a complementary ranker to existing term-based retrieval models such as BM25. To achieve this, we make use of a transformer-based pretrained language model and its inherent ability to make use of token-level local interactions. Its complementary behavior is further supported by negative sampling from BM25 rankings.

#### 3.1 Architecture

The model architecture of CoRT, illustrated in Figure 1, follows the idea of a *Siamese Neural Network* (Bromley et al., 1993). Passages and queries are encoded using an identical model with shared weights except for one detail: The passage encoder  $\psi_\alpha$  and the query encoder  $\psi_\beta$  use different segment embeddings (Devlin et al., 2019). CoRT computes relevance scores as angular similarity between query and passage representations while training a *pair-wise* ranking objective.

#### 3.2 Encoding

CoRT can incorporate any BERT-like encoder as underlying text encoder. Here, we use a pretrained ALBERT (Lan et al., 2020) encoder for its smaller model size, the tougher sentence coherence pre-training and increased first-stage ranking quality throughout our early-stage experiments compared to BERT. The tokenizer of ALBERT is a Word-Piece tokenizer (Wu et al., 2016) including the special tokens [CLS] and [SEP] known from BERT. From the text encoder we seek a single representa-

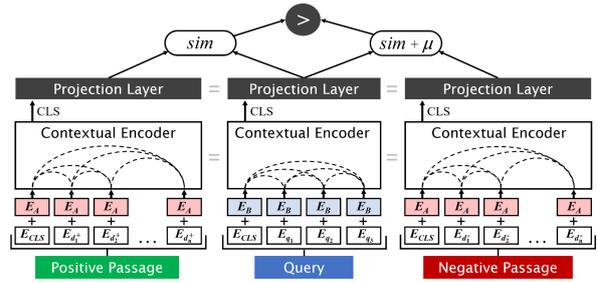


Figure 1: CoRT’s model architecture and pair-wise learning objective (simplified).

tion vector for the whole passage or query, which we call *context representation*. From ALBERT we take the [CLS] embedding of the last layer for this purpose. We denote the context representation obtained from the underlying encoder for an arbitrary string  $s$  with  $\tau(s) \in \mathbb{R}^h$  where  $h$  is the output representation size.

ALBERT’s language modeling approach involves sentence coherence prediction for which segment embeddings are used to signal different input segments. Although we only feed single segments to the encoder, i.e. a query or a passage, we use segment embeddings which allow the model to encode queries differently from passages. The segment embeddings  $E_A$  and  $E_B$  (illustrated in Figure 1) are part of the context encoder functions  $\tau_\alpha$  and  $\tau_\beta$  for passages and queries respectively. The context representation is further projected to the desired representation size  $e$  using a linear layer followed by a tanh activation function. Thus, the complete passage encoder function is  $\psi_\alpha(s) := \tanh(W\tau_\alpha(s) + b)$  where  $W \in \mathbb{R}^{h \times e}$  and  $b \in \mathbb{R}^e$  are parameters of the linear layer. The query encoder  $\psi_\beta$  is defined analogous.

#### 3.3 Training

Training CoRT corresponds to updating the parameters of the encoder  $\psi$  towards representations that reflect relevance between queries and passages through vector similarity. Each training sample is a triple comprising a query  $q$ , a positive passage  $d^+$  and a negative passage  $d^-$ . While positive passages are taken from relevance assessments, negative passages are sampled from term-based rankings (i.e. BM25) to support the complementary property of CoRT. The relevance score for a query-passage pair  $(q, d)$  is calculated using the angular cosine similarity function<sup>2</sup>:

<sup>2</sup>Similar to Cer et al. (2018), we found angular similarity performs better than cosine similarity.

$$\text{sim}(q, d) := 1 - \arccos\left(\frac{\psi_\beta(q) \cdot \psi_\alpha(d)}{\|\psi_\beta(q)\| \|\psi_\alpha(d)\|}\right) / \pi$$

As illustrated in Figure 1, the training objective is to score the positive example  $d^+$  by at least the margin  $\mu$  higher than the negative one  $d^-$ . As part of our loss function, we use the triplet margin objective:

$$l(q, d^+, d^-) := \max(0, \text{sim}(q, d^-) - \text{sim}(q, d^+) + \mu)$$

Inspired by Song et al. (2016), we aim to take full advantage of the whole training batch. For each query, each passage in the batch is used as a negative example, except for the given positive passage. Thus, we define our batch-wise loss function as follows:

$$\mathcal{L} := \sum_{1 \leq i \leq n} \left( \sum_{1 \leq j \leq n} l(q_i, d_i^+, d_j^-) + \sum_{1 \leq k \leq n, k \neq i} l(q_i, d_i^+, d_k^+) \right)$$

$q_i$ ,  $d_i^+$  and  $d_i^-$  denote the triple of the  $i_{th}$  sample in the batch and  $n$  the number of samples per batch. We found this technique to make the training process more robust against exploding gradients. Otherwise we need to employ gradient clipping (Zhang et al., 2020) to stabilize the training process. Also, it positively affects first-stage ranking results<sup>3</sup>.

### 3.4 Indexing and Retrieval

For retrieval with CoRT, each passage must be encoded by the passage encoder  $\psi_\alpha$ . Subsequent normalization of each vector allows us to use the dot product as a proxy score function for  $\text{sim}$ , which is sufficient to form rankings accurately. Given a query  $q$ , we calculate its representation  $\psi_\beta(q)$  and the dot product with each normalized passage vector. From those, the  $k$  highest scores are selected and sorted to form the CoRT ranking. This procedure can be implemented heavily parallelized using GPU matrix operations. Alternatively, the passage representations can be indexed in an ANN index to avoid exhaustive similarity search. In contrast to the first-stage ranking of Khatib and Zaharia (2020) and MacAvaney et al. (2020), we

<sup>3</sup>We achieve 2.0 p.p. higher MRR@10 compared to the plain triplet margin loss on the MS MARCO passage task.

only index one representation per passage rather than one per token. Finally, we combine the resulting ranking of CoRT with the respective BM25 ranking by *interleaving* the positions beginning with CoRT to create a single merged ranking of equal length. During this process, each passage that was already added by the other ranking is omitted. For example, merging two ranking lists beginning with  $[a, b, c, d, \dots]$  and  $[e, c, f, a, \dots]$  would result in  $[a, e, b, c, f, d, \dots]$ . The interleaving procedure stops as soon as the desired ranking size has been reached. The result is a compound ranking of CoRT and BM25, which we denote with  $\text{CoRT}_{\text{BM25}}$ .

## 4 Experiments

We present four experiments studying the ranking quality and recall of CoRT, the connection between the number of candidates and re-ranking effectiveness, the impact of the representation size  $e$ , and CoRT’s retrieval latencies. Finally, we outline a competitive end-to-end ranking setup with CoRT and a BERT-based re-ranker.

### 4.1 Datasets

#### 4.1.1 MS MARCO Passage Retrieval

The *Microsoft Machine Reading Comprehension* (Nguyen et al., 2016) dataset for passage ranking was introduced in 2018. It provides a benchmark for passage retrieval with real-world queries and passages gathered from Microsoft’s *Bing* search. The MS MARCO passage ranking task comprises 8.8M passages sampled from web pages and about 1M queries that are formulated as questions. The objective is to rank those passages high that were labeled as relevant to answer the respective question. The annotations, however, are sparse. There are 530k positive relevance labels distributed over 808k queries in the training set, whereby most queries are associated to one passage. The validation and evaluation sets, `dev` and `eval`, comprise 101k queries each. An official subset of `dev`, called `dev.small` comprises 6980 queries and 7437 relevance labels and is often used for publicly reported evaluations. We follow this convention and use `dev.small` for testing. The creators suggest to use the *mean reciprocal rank* cut at the tenth position (MRR@10) as primary evaluation measure. Additionally, we measure NDCG@20 (Manning et al., 2008) as less punishing ranking quality measure and the recall at various positions

Table 1: First-stage ranking results on the MS MARCO `dev.small` set. The asterisk (\*) denotes merged rankings using an instance of CoRT that was not specifically trained to complement the corresponding term-based ranker.

MS MARCO Passage ( <code>dev.small</code> )	MRR @10	NDCG @20	@50	@100	RECALL @200	@500	@1000
BM25	18.7	25.8	59.2	67.0	73.8	81.2	85.7
doc2query	21.5	-	-	-	-	-	89.3
DeepCT	24.3	32.1	68.5	75.2	81.0	87.3	90.9
docTTTTTquery	27.7	36.5	75.6	81.9	86.9	91.6	94.7
CoRT	27.1	34.0	66.4	73.1	78.6	84.4	88.0
CoRT <sub>BM25</sub>	27.4	35.9	74.3	81.6	87.3	92.5	94.9
CoRT <sub>DeepCT*</sub>	28.3	36.8	75.6	82.3	87.6	92.5	94.9
CoRT <sub>DocTTTTTquery*</sub>	<b>28.8</b>	<b>38.0</b>	<b>78.5</b>	<b>85.3</b>	<b>90.0</b>	<b>94.4</b>	<b>96.5</b>

to indicate how many relevant passages a re-ranker would miss if the number of candidates is reduced.

#### 4.1.2 TREC 2019 DL Passage Retrieval

The passage retrieval section of the *TREC 2019 Deep Learning Track* (Craswell et al., 2020) provides on average 215 manual relevance assessments per query for a set of 43 MS MARCO queries. Each assessment corresponds to a rating on a scale from 0 (not relevant) to 3 (perfectly relevant). We adopt the evaluation metrics MRR (uncut), NDCG@10 and MAP from the official TREC overview. In contrast to the original MS MARCO benchmark, this evaluation set provides dense annotations, but only for few queries.

## 4.2 First-Stage Ranking

We train CoRT as described in Section 3.3 while using a representation size of  $e = 768$ . In this section we discuss the first-stage ranking results of our model using the datasets and their associated metrics described in Section 4.1.

### 4.2.1 MS MARCO Passage Retrieval

The results of CoRT and its baselines on the MS MARCO passage retrieval task (`dev.small`) are reported in Table 1. Next to BM25 as a baseline, we include *DeepCT* (Dai and Callan, 2020), *doc2query* (Nogueira et al., 2019) and its successor *docTTTTTquery* (Nogueira and Lin, 2019). All three are recent first-stage rankers with average retrieval latencies below 100ms per query on the MS MARCO passage corpus. The metrics MRR@10 and NDCG@20 reveal a quite decent ranking quality for the standalone CoRT ranker. Since CoRT’s primary use is candidate retrieval rather than standalone ranking, we pay particular attention to the recall at various cuts. From the perspective of BM25, the absolute increase of recall due to merging with CoRT ranges between

Table 2: First-stage ranking results on the *TREC 2019 DL* passage task.

TREC DL 2019	MRR @1000	NDCG @10	MAP @1000	RECALL @500
BM25	68.5	49.7	29.0	69.4
CoRT	84.3	<b>60.0</b>	29.7	58.3
CoRT <sub>BM25</sub>	<b>86.2</b>	59.7	<b>35.1</b>	<b>76.9</b>

15.1 (RECALL@50) and 9.2 (RECALL@1000), which we consider the complementary portion of CoRT. Greater increases of recall can be noticed for lower cuts, which is particular useful when re-ranking is performed with low numbers of candidates. The top-200 candidates from CoRT<sub>BM25</sub> comprise higher recall than the top-1000 candidates from BM25. The metrics for *DeepCT* and *docTTTTTquery* have been calculated using published top-1000 rankings from the respective authors. Thus, we were able to merge those rankings with CoRT. However, the used instance of CoRT was trained on BM25 and not on the external ranker.

### 4.2.2 TREC 2019 DL Passage Retrieval

Although we consider the relevance assessments from TREC 2019 DL to be dense, we found 112 unlabeled passages among the 43 top-10 CoRT rankings while the assessments for the BM25 rankings are complete. This means there are, on average, 2.6 unlabeled passages in the top-10 CoRT rankings, which might make this evaluation somewhat unfavourable for CoRT and explain the drop in recall. Still, Table 2 shows superior results for CoRT compared to BM25 in terms of ranking quality (MRR, NDCG and MAP). Merging CoRT with BM25 slightly increases MRR and NDCG, while a decent gain in terms of recall can be noticed. We can not report any results for *DeepCT* or *docTTTT-*

$T_{query}$ , since only rankings for the MS MARCO dev .small set are available online from the respective authors.

### 4.3 Candidate Re-ranking

We re-rank candidates from both BM25 and  $\text{CoRT}_{\text{BM25}}$  to study the impact of the candidates on a subsequent *interaction-focused* re-ranking. By varying the numbers of candidates, we investigate at which point adding more candidates becomes ineffective.

#### 4.3.1 Re-ranking Model

Similar to (Nogueira and Cho, 2019), we use a simple binary classifier based on BERT. The query-passage pair  $(q, p)$  is concatenated to one token sequence of two segments. This sequence is processed by the BERT encoder while the [CLS] embedding of the last layer, which we denote with  $\phi(q, p)$ , is projected to a single classification logit. We then apply the sigmoid activation function  $\sigma$  to obtain the relevance confidence for query  $q$  and passage  $p$ . This procedure can be formalized as  $\zeta(q, p) = \sigma(W'\phi(q, p) + b')$  where  $W' \in \mathbb{R}^{h \times 1}$  and  $b' \in \mathbb{R}$  are the parameters of a linear layer with a single output activation. To form a ranking at inference time, we sort the candidates by the model’s confidence. Following (Nogueira and Cho, 2019), this model is trained using a point-wise objective. We sample query-passage pairs, each associated with a binary relevance label  $y \in \{0, 1\}$  and minimize the binary cross-entropy loss:

$$l'(q, p, y) = y \cdot \log \zeta(q, p) + (1 - y) \cdot \log (1 - \zeta(q, p))$$

#### 4.3.2 Re-ranking Results

As illustrated in figure 2, using  $\text{CoRT}_{\text{BM25}}$  as first-stage ranking appears to result in superior end-to-end ranking quality in terms of MRR@10. This is especially true, if low numbers of candidates are used. We also notice earlier saturation<sup>4</sup> of MRR@10 for  $\text{CoRT}_{\text{BM25}}$ , which is illustrated in Figure 2. Only 64 candidates from  $\text{CoRT}_{\text{BM25}}$  are sufficient to achieve top results with this re-ranker. In contrast, 256 candidates from BM25 are needed to reach the point of saturation, which translates in quadrupled re-ranking time.

<sup>4</sup>In our definition, saturation is reached, when doubling the number of candidates results in less than 0.5% increase of the respective metric

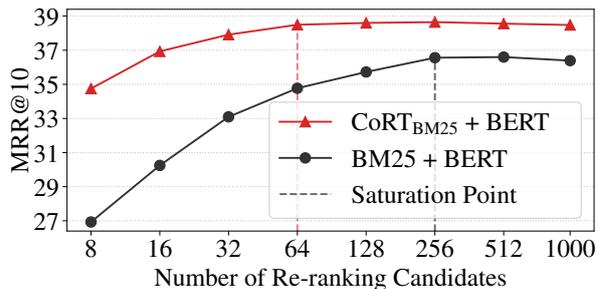


Figure 2: Re-ranking quality by number of candidates for BM25 (black) and  $\text{CoRT}_{\text{BM25}}$  (red) on the MS MARCO passage task. Dashed lines indicate effectiveness saturation (<0.5% increase).

### 4.4 Impact of Representation Size

As described in Section 3.2, CoRT projects the context representation of the underlying encoder  $\tau$  to an arbitrary representation size  $e$ . This size determines the size of the final index and also influences the retrieval latency. The total size of the encoded corpus is easy to calculate. For example, with  $e = 128$  and the MS MARCO corpus, the index size (without overhead) amounts  $8.8M \text{ documents} \times 128 \text{ floats/document} \times 4 \text{ bytes/float} \approx 4.5 \times 10^9 \text{ bytes} \approx 4.2 \text{ GB}$ . Thus,  $e$  is proportional to the total size and reducing  $e$  to 64 would halve the memory footprint. If  $e$  is small, however, it is more difficult to attain the training objective. Thus,  $e$  can be used for a trade-off between ranking quality and computational effort / resource cost. We investigate the relation between the representation size  $e$  and the ranking quality by conducting identical training runs with different numbers for  $e$ . The results in Table 3 show that MRR@10 already saturates at  $e = 128$ . Interestingly, even with an representation size of  $e = 32$ , CoRT outperforms BM25 in terms of MRR@10 and nDCG@20 by a big margin.

Table 3: First-stage ranking results for various representation sizes.

	$e$	MRR @10	nDCG @20	RECALL @200	RECALL @1k
CoRT	32	23.6	29.8	70.6	81.8
	64	25.6	32.3	75.6	85.7
	128	26.8	33.4	77.2	87.1
	256	26.8	33.6	78.2	87.8
	768	27.1	34.0	78.6	88.0
CoRT <sub>BM25</sub>	32	25.2	33.5	85.2	93.7
	64	26.6	34.9	86.2	94.4
	128	27.4	35.7	87.0	94.8
	256	27.2	35.6	87.0	94.9
	768	27.4	36.0	87.3	94.9

## 4.5 Latency Measurements

We propose two methods for the deployment of CoRT. The first exhaustively calculates similarity scores using multiple GPUs while the second incorporates an *Approximate Nearest Neighbor* index (ANN). We measure retrieval latencies of those methods and compare them with BM25 as representative for term-based retrieval models based on inverted indexing. We conduct the latency measurement based on the top-1000 retrieval for the `dev.small` split of the MS MARCO passage corpus. Since some approaches profit from batch computing, we also measure the latency for batches of 32 queries. As representation size, we have chosen  $e = 128$ , since it is the smallest representation size investigated in Section 4.4 that does not hurt the ranking quality of  $\text{CoRT}_{\text{BM25}}$ .

### 4.5.1 Lucene BM25 Baseline

As retrieval latency baseline, we use a *Lucene* index generated by the *Anserini* toolkit (Yang et al., 2017). Please note, this is not perfectly representative for sparse bag-of-words retrieval: Retrieval latency can be reduced due to index pruning without significantly hurting retrieval quality (Mackenzie et al., 2020). The retrieval was performed on a machine with an Intel Core i9-9900KS processor (16 logical cores, 8 physical) and enough memory to hold the whole corpus. Single queries were processed using the single-threaded search function, while batch-wise search has been performed with 16 threads.

### 4.5.2 Retrieval using multiple GPUs

Multiple GPUs can be used to deploy CoRT for fast large-scale ranking. We propose to uniformly distribute the vector representations of the corpus on the available GPUs. Each GPU ranks its own partition of the corpus as described in Section 3.4. Afterwards, the results for each partition are aggregated by selecting the top-k candidates with highest scores.

### 4.5.3 Retrieval using ANN

Since CoRT operates on vector similarities, it can make use of ANN search. We measure the retrieval latency and the loss of recall, which occurs due to the pruning heuristics. We use a graph-based index optimized with the *ONNG* method (Iwasaki and Miyazaki, 2018). An implementation of this method is publicly available as part of the *NGT*

Table 4: Retrieval latencies averaged over 6980 queries of the `dev.small` split.

MS MARCO Passage ( <code>dev.small</code> )	RECALL @200	LATENCY (ms)	
		Single	Batch32
BM25 ( <i>anserini</i> )	73.8	38	290
<b>CoRT</b> ( $e = 128$ )			
<u>Query Encoding</u>			
- Single GPU	-	8	17
<u>Retrieval</u>			
- Single GPU	77.2	68	164
- Quad GPU	77.2	17	35
- ANN $_{\epsilon = 0.01}$	76.6	4	-
- ANN $_{\epsilon = 0.1}$	76.9	17	-
- ANN $_{\epsilon = 0.4}$	77.2	71	-
<b>CoRT<sub>BM25</sub> Total</b>			
- Quad-GPU	87.0	~63	~342
- ANN $_{\epsilon = 0.1}$	86.9	~63	-
- ANN $_{\epsilon = 0.01}$	86.8	~50	-

*Library*<sup>5</sup>. To adjust the trade-off between retrieval latency and accuracy of the index, we alter the search range coefficient  $\epsilon$ . We always retrieve 1000 candidates from the ANN index, even if we use a smaller number of candidates in a ranking pipeline.

### 4.5.4 Latency Measurements

The latency measurements are reported in Table 4. For CoRT the total retrieval latency per query consists of two factors: Query encoding and retrieval. The query encoding has to be performed by the query encoder  $\psi_{\beta}$ , which we highly recommend to run on a GPU. The latency of the retrieval depends on the retrieval methods described above. Employing multiple GPUs appears to reduce retrieval latency on a linear scale: The exhaustive search using four GPUs takes 17ms for a single query, while a single GPU takes four times as long. The total retrieval time per query sums up to  $17 + 8 = 25ms$  for the quad GPU setting, which is below the BM25 baseline. It is worth noting, that batch-wise processing results in a substantial efficiency increase: Retrieval for 32 queries at once only takes about twice as long as a single query. This can be useful if multiple queries queue up while the system is busy. The tested BM25 index (*Anserini*), on the other side, seems to suffer from multiprocessing overhead or some sort of bottleneck. The latencies for the ANN index has been measured with three different values for the search range coefficient  $\epsilon$ . While this significantly affects the retrieval latency, only slight differences regarding the quality of the first-stage ranking are

<sup>5</sup><https://github.com/yahoojapan/NGT>

observed. Latencies for  $\text{CoRT}_{\text{BM25}}$  comprise latencies from BM25, CoRT’s query encoding and the corresponding retrieval method. For simplicity, we assume sequential processing of all three components, although BM25 could be processed in parallel.

#### 4.6 End-to-end Retrieval

Intrigued by the remarkable ratio of retrieval latency and ranking quality of ColBERT’s full-ranking approach (Khattab and Zaharia, 2020), we used our above findings to create a competitive end-to-end ranking setup. We suggest to re-rank the top-64 candidates from  $\text{CoRT}_{\text{BM25}}$  with  $e = 128$ , retrieved by an ANN index ( $\epsilon = 0.1$ ). The end-to-end latency comprises  $8ms$  for query encoding,  $17ms$  for CoRT retrieval based on ONNG,  $38ms$  for BM25 retrieval, and  $192ms$  for re-ranking. Although the BM25 candidates could be retrieved in parallel, we report sequential processing latencies. As shown in Table 5, we outperform ColBERT’s end-to-end ranking performance in terms of  $\text{MRR}@10$  and retrieval latency. It is worth noting that the "RTX 2080 Ti" we used for latency measurements is less powerful than the "Tesla V100" Khattab and Zaharia (2020) used for their measurements. CoRT’s representations for the MS MARCO corpus only weight 4.3GB when  $e$  is set to 128, or 7.0GB when indexed in an ONNG index. The size of the query encoder only amounts about 50MB, which is due to ALBERT’s parameter sharing. To compile the full  $\text{CoRT}_{\text{BM25}}$  candidates, the corresponding BM25 index is needed, which amounts 2.2 GB on disk. Although more memory is needed to deploy and operate both indexes, this is by far less than the 154GB footprint reported by Khattab and Zaharia (2020) for ColBERT’s end-to-end approach.

Table 5: Measured end-to-end ranking quality and latency. \* ColBERT’s latency was measured with different hardware

MS MARCO Passage (dev. small)	MRR @10	LATENCY (ms)
ColBERT <sub>L2</sub> (Khattab and Zaharia, 2020)	36.0	458*
CoRT <sub>BM25</sub> (ANN <sub><math>\epsilon = 0.1</math></sub> , top-64) + BERT Re-ranking	38.4	255

## 5 Conclusion

In this paper, we propose CoRT, a framework and neural first-stage ranking model that leverages contextual representations from transformer-based language models to complement term-based ranking functions. As a result, we observe decently increased recall measures and improved end-to-end ranking quality on the MS MARCO passage task. Also, we are able to decrease the number of candidates for re-ranking without hurting the final performance. Our further experiments reveal sweet spots for CoRT’s representation size and the number of re-ranking candidates. We presented two deployment strategies for CoRT and measured their performances in terms of efficiency and effectiveness. Finally, we demonstrate CoRT can be used with a simple BERT-based re-ranker to create a competitive ranking pipeline.

### Acknowledgments

We would like to thank Felix Hamann and Prof. Dr. Adrian Ulges for helpful discussions and comments on the manuscript, as well as the anonymous reviewers for their valuable feedback. This work was funded by German Federal Ministry of Education and Research (Program FHprofUnt, Project DeepCA (13FH011PX6)).

### References

- Leonid Boytsov, David Novak, Yury Malkov, and Eric Nyberg. 2016. [Off the beaten path: Let’s replace term-based retrieval with k-nn search](#). In *Proceedings of the 25th ACM International Conference on Information and Knowledge Management, CIKM 2016, Indianapolis, IN, USA, October 24-28, 2016*, pages 1099–1108. ACM.
- Jane Bromley, James W. Bentz, Léon Bottou, Isabelle Guyon, Yann LeCun, Cliff Moore, Eduard Säckinger, and Roopak Shah. 1993. [Signature verification using A "siamese" time delay neural network](#). *Int. J. Pattern Recognit. Artif. Intell.*, 7(4):669–688.
- Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St. John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, Brian Strope, and Ray Kurzweil. 2018. [Universal sentence encoder for english](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, EMNLP 2018: System Demonstrations, Brussels, Belgium, October 31 - November 4, 2018*, pages 169–174. Association for Computational Linguistics.

- Ruey-Cheng Chen, Luke Gallagher, Roi Blanco, and J. Shane Culpepper. 2017. [Efficient cost-aware cascade ranking in multi-stage retrieval](#). In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval, Shinjuku, Tokyo, Japan, August 7-11, 2017*, pages 445–454. ACM.
- Nick Craswell, Bhaskar Mitra, Emine Yilmaz, Daniel Campos, and Ellen M. Voorhees. 2020. [Overview of the TREC 2019 deep learning track](#). *CoRR*, abs/2003.07820.
- Zhuyun Dai and Jamie Callan. 2020. [Context-aware term weighting for first stage passage retrieval](#). In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, SIGIR 2020, Virtual Event, China, July 25-30, 2020*, pages 1533–1536. ACM.
- Zhuyun Dai, Chenyan Xiong, Jamie Callan, and Zhiyuan Liu. 2018. [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, WSDM 2018, Marina Del Rey, CA, USA, February 5-9, 2018*, pages 126–134. ACM.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [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. Association for Computational Linguistics.
- Luyu Gao, Zhuyun Dai, and Jamie Callan. 2020. [Modularized transformer-based ranking framework](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 4180–4190. Association for Computational Linguistics.
- Jiafeng Guo, Yixing Fan, Qingyao Ai, and W. Bruce Croft. 2016. [A deep relevance matching model for ad-hoc retrieval](#). In *Proceedings of the 25th ACM International Conference on Information and Knowledge Management, CIKM 2016, Indianapolis, IN, USA, October 24-28, 2016*, pages 55–64. ACM.
- Jiafeng Guo, Yixing Fan, Liang Pang, Liu Yang, Qingyao Ai, Hamed Zamani, Chen Wu, W. Bruce Croft, and Xueqi Cheng. 2020. [A deep look into neural ranking models for information retrieval](#). *Inf. Process. Manag.*, 57(6):102067.
- Christophe Van Gysel, Maarten de Rijke, and Evangelos Kanoulas. 2018. [Neural vector spaces for unsupervised information retrieval](#). *ACM Trans. Inf. Syst.*, 36(4):38:1–38:25.
- Po-Sen Huang, Xiaodong He, Jianfeng Gao, Li Deng, Alex Acero, and Larry P. Heck. 2013. [Learning deep structured semantic models for web search using clickthrough data](#). In *22nd ACM International Conference on Information and Knowledge Management, CIKM'13, San Francisco, CA, USA, October 27 - November 1, 2013*, pages 2333–2338. ACM.
- Masajiro Iwasaki and Daisuke Miyazaki. 2018. [Optimization of indexing based on k-nearest neighbor graph for proximity search in high-dimensional data](#). *CoRR*, abs/1810.07355.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. [Dense passage retrieval for open-domain question answering](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6769–6781, Online. Association for Computational Linguistics.
- Omar Khattab and Matei Zaharia. 2020. [Colbert: Efficient and effective passage search via contextualized late interaction over bert](#). In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '20*, page 39–48, New York, NY, USA. Association for Computing Machinery.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. [ALBERT: A lite BERT for self-supervised learning of language representations](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Ilya Loshchilov and Frank Hutter. 2019. [Decoupled weight decay regularization](#). In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.
- Sean MacAvaney, Franco Maria Nardini, Raffaele Perego, Nicola Tonello, Nazli Goharian, and Ophir Frieder. 2020. [Efficient document re-ranking for transformers by precomputing term representations](#). In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, SIGIR 2020, Virtual Event, China, July 25-30, 2020*, pages 49–58. ACM.
- Joel M. Mackenzie, Zhuyun Dai, Luke Gallagher, and Jamie Callan. 2020. [Efficiency implications of term weighting for passage retrieval](#). In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, SIGIR 2020, Virtual Event, China, July 25-30, 2020*, pages 1821–1824. ACM.
- Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. 2008. *Introduction to information retrieval*. Cambridge University Press.

- Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. [MS MARCO: A human generated machine reading comprehension dataset](#). In *Proceedings of the Workshop on Cognitive Computation: Integrating neural and symbolic approaches 2016 co-located with the 30th Annual Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain, December 9, 2016*, volume 1773 of *CEUR Workshop Proceedings*. CEUR-WS.org.
- Rodrigo Nogueira and Kyunghyun Cho. 2019. [Passage re-ranking with BERT](#). *CoRR*, abs/1901.04085.
- Rodrigo Nogueira and Jimmy Lin. 2019. [From doc2query to docttttquery](#). *Online preprint*.
- Rodrigo Nogueira, Wei Yang, Jimmy Lin, and Kyunghyun Cho. 2019. [Document expansion by query prediction](#). *CoRR*, abs/1904.08375.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. [Pytorch: An imperative style, high-performance deep learning library](#). In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 8024–8035.
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. [Glove: Global vectors for word representation](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1532–1543. ACL.
- Yifan Qiao, Chenyan Xiong, Zhenghao Liu, and Zhiyuan Liu. 2019. [Understanding the behaviors of BERT in ranking](#). *CoRR*, abs/1904.07531.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8):9.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *J. Mach. Learn. Res.*, 21:140:1–140:67.
- Yelong Shen, Xiaodong He, Jianfeng Gao, Li Deng, and Grégoire Mesnil. 2014. [Learning semantic representations using convolutional neural networks for web search](#). In *23rd International World Wide Web Conference, WWW '14, Seoul, Republic of Korea, April 7-11, 2014, Companion Volume*, pages 373–374. ACM.
- Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. 2016. [Deep metric learning via lifted structured feature embedding](#). In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 4004–4012. IEEE Computer Society.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008.
- Lidan Wang, Jimmy J. Lin, and Donald Metzler. 2011. [A cascade ranking model for efficient ranked retrieval](#). In *Proceeding of the 34th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2011, Beijing, China, July 25-29, 2011*, pages 105–114. ACM.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. 2019. [Huggingface’s transformers: State-of-the-art natural language processing](#). *CoRR*, abs/1910.03771.
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. [Google’s neural machine translation system: Bridging the gap between human and machine translation](#). *CoRR*, abs/1609.08144.
- Chenyan Xiong, Zhuyun Dai, Jamie Callan, Zhiyuan Liu, and Russell Power. 2017. [End-to-end neural ad-hoc ranking with kernel pooling](#). In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval, Shinjuku, Tokyo, Japan, August 7-11, 2017*, pages 55–64. ACM.
- Peilin Yang, Hui Fang, and Jimmy Lin. 2017. [Anserini: Enabling the use of lucene for information retrieval research](#). In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval, Shinjuku, Tokyo, Japan, August 7-11, 2017*, pages 1253–1256. ACM.
- Hamed Zamani, Mostafa Dehghani, W. Bruce Croft, Erik G. Learned-Miller, and Jaap Kamps. 2018. [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, CIKM 2018, Torino, Italy, October 22-26, 2018*, pages 497–506. ACM.
- Jingzhao Zhang, Tianxing He, Suvrit Sra, and Ali Jadbabaie. 2020. [Why gradient clipping accelerates training: A theoretical justification for adaptivity](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

## A Appendix: Implementation Details

### A.1 Hardware & Software

We use up to four NVIDIA GTX 2080 TI graphic cards in combination with 128GB DDR4 RAM and an Intel Core i9-9900KS processor. We use *PyTorch* (Paszke et al., 2019) and *HuggingFace’s Transformers* (Wolf et al., 2019) as deep learning libraries. BM25 rankings are generated using the *Anserini* toolkit (Yang et al., 2017).

### A.2 CoRT Training

We train CoRT based on the pretrained ALBERT model "albert-base-v2", which is the lightest available version in *HuggingFace’s* repository<sup>6</sup>. Each model is trained for 10 epochs, where each epoch includes all queries that are associated to at least one relevant document plus one randomly sampled positive and one negative passage. Negative examples are sampled from unlabeled passages of top-100 BM25 rankings. There, we exclude the first 8 ranks to reduce the probability of drawing actual relevant passages and thus give contradictory signals less often. We find this slightly increases CoRT’s ranking quality. As usual for BERT-based models we use the ADAM optimizer with weight decay fix (Loshchilov and Hutter, 2019) and the default parameters  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $eps = 10^{-6}$ , a weight decay rate of  $\lambda = 0.1$  and a linearly decreasing learning rate schedule starting with  $lr = 2 \times 10^{-5}$  after

2.000 warm-up steps. We train mini-batches of size  $n = 6$  samples (triples) while accumulating the gradients of 100 mini-batches before performing one update step. The triplet margin is set to  $\mu = 0.1$ , which has been tuned within the range of  $[0.01, 0.2]$ .

### A.3 Re-ranker Training

Our BERT re-ranking experiment utilizes the pretrained "bert-base-uncased" model, hosted by *HuggingFace*. We use equal optimizer settings than for CoRT except for the learning rate, which we empirically set to  $5 \times 10^{-5}$ . We use a batch-size of 8 and accumulate the gradients of 16 batches.

## B Appendix: Retrieval Examples

Table 6 shows top-1 retrieval examples of CoRT and BM25. The first query exemplifies the advantage of local interactions in the query encoder. We hypothesize, the query could be interpreted as a question about the density of aluminum although the term density was not included. The second query is an example, where BM25 works well due to favorable keywords in the passage. Although CoRT’s top result is not labeled, it clearly is relevant to the query. Since the passage misses the keyword "insane", it is difficult to retrieve for a term-based model. We hypothesize, due to the terms "hallucinations" and "paranoia", CoRT is able to match the contexts in this example.

Table 6: Retrieval Examples with highlighted keywords. Ranks beyond the top-1000 are denoted with "n/a".

Query	Sample Passage	Label	Rank	
			BM25	CoRT
how much does aluminum weigh	Question and answer. how much does a western 14 ft. aluminum boat weigh? what is the weight difference between a 12 ft. and 14 ft. aluminum boat. 12 ft aluminum boat with no gear or anything would probably weigh about 115-150 lbs, 14 ft aluminum boat with no gear or anything would probably weigh about 250-300 lbs.	n/a	1	n/a
	Quick Answer. One cubic inch of aluminum weighs 1.56 ounces. The metal sinks in water, but it is still relatively lightweight. The density of aluminum is 2.7 grams per milliliter. Aluminum is used as a metal foil, a conductor of electricity and in the construction of airplane fuselages.	True	735	1
how many days of no sleep until insane	Although there are many articles stating one can be declared legally insane if... How long must you go without sleep to be declared legally insane Each state would have different laws regarding the requirements of declaring... Is it true that after three days of complete sleep deprivation you are considered legally insane? Not after 3 days but any longer will.	True	1	n/a
	The longest recorded time a human has ever gone without sleep is 18 days, 21 hours, and 40 minutes, which resulted in hallucinations, paranoia, etc. However most people can only last 4-6 days without stimulants, and about 7-10 days before the body will be unable to function and long term damage can be caused.	n/a	n/a	1

<sup>6</sup>[https://huggingface.co/transformers/pretrained\\_models.html](https://huggingface.co/transformers/pretrained_models.html)