

# LED: Lexicon-Enlightened Dense Retriever for Large-Scale Retrieval

Kai Zhang\*  
The Ohio State University  
Columbus, Ohio, USA  
zhang.13253@osu.edu

Chongyang Tao  
Microsoft Corporation  
Beijing, China  
chotao@microsoft.com

Tao Shen  
AAIL, FEIT, University of  
Technology Sydney  
Sydney, Australia  
tao.shen@uts.edu.au

Can Xu  
Microsoft Corporation  
Beijing, China  
caxu@microsoft.com

Xiubo Geng  
Microsoft Corporation  
Beijing, China  
xigeng@microsoft.com

Binxing Jiao  
Microsoft Corporation  
Beijing, China  
binxjia@microsoft.com

Daxin Jiang<sup>†</sup>  
Microsoft Corporation  
Beijing, China  
djjiang@microsoft.com

## ABSTRACT

Retrieval models based on dense representations in semantic space have become an indispensable branch for first-stage retrieval. These retrievers benefit from surging advances in representation learning towards compressive global sequence-level embeddings. However, they are prone to overlook local salient phrases and entity mentions in texts, which usually play pivot roles in first-stage retrieval. To mitigate this weakness, we propose to make a dense retriever align a well-performing lexicon-aware representation model. The alignment is achieved by weakened knowledge distillations to enlighten the retriever via two aspects – 1) a lexicon-augmented contrastive objective to challenge the dense encoder and 2) a pair-wise rank-consistent regularization to make the dense model’s behavior incline to the other. We evaluate our model on three public benchmarks, which shows that with a comparable lexicon-aware retriever as the teacher, our proposed dense one can bring consistent and significant improvements, and even outdo its teacher. In addition, we show our lexicon-aware distillation strategies are compatible with the standard ranker distillation, which can further lift state-of-the-art performance.<sup>1</sup>

## CCS CONCEPTS

• Information systems → Retrieval models and ranking.

## KEYWORDS

Dense retrieval, Lexicon-aware retrieval, Lexicon augmentation, Contrastive learning, Knowledge distillation

\*Work done during the internship at Microsoft.

<sup>†</sup>Corresponding author.

<sup>1</sup>Code is available at <https://github.com/drogozhang/LED>.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

WWW '23, May 1–5, 2023, Austin, TX, USA

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-9416-1/23/04...\$15.00

<https://doi.org/10.1145/3543507.3583294>

## ACM Reference Format:

Kai Zhang, Chongyang Tao, Tao Shen, Can Xu, Xiubo Geng, Binxing Jiao, and Daxin Jiang. 2023. LED: Lexicon-Enlightened Dense Retriever for Large-Scale Retrieval. In *Proceedings of the ACM Web Conference 2023 (WWW '23)*, May 1–5, 2023, Austin, TX, USA. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3543507.3583294>

## 1 INTRODUCTION

Large-scale passage retrieval [2] aims to fetch relevant passages from a million- or billion-scale collection for a given query to meet users’ information needs, serving as an important role in many downstream applications including open domain question answering [20], search engine [56], and recommendation system [54], etc. Recent years have witnessed an upsurge of interest and remarkable performance of dense passage retrievers on first-stage retrieval. Built upon powerful pre-trained language models (PLM) [9, 30, 36], dense retrievers [20, 35, 49] encode queries and passages into a joint low-dimensional semantic space in a Siamese manner (i.e. dual-encoder), so that the passages could be offline pre-indexed and query could be encoded online and searched via approximate nearest neighbor [19], reaching an efficiency-effectiveness trade-off.

Although dense retrieval becomes indispensable in modern systems, a long-term challenge is that the dense representations in a latent semantic space are abstractive and condensed, exposing the systems to a risk that pivot phrases and mentions may be overlooked and thus leading to sub-optimal efficacy. For example, DPR [20] didn’t regard “*Thoros of Myr*” as an entity mention in the query “*Who plays Thoros of Myr in Game of Thrones?*”. Analogously, given the query “*What is an active margin*”, ANCE [49] overlooked the “*active margin*” as an entire local salient phrase and hence retrieved passages related to the financial term “*margin*”. As a remedy, prior works resort to either coupling a dense retriever with the term matching scores (e.g., TF-IDF, BM25) [7, 14, 26, 33] or learning BM25 ranking into a dense model as additional features to complement the original one [4]. But, these approaches are limited by superficial combinations and almost unlearnable BM25 scoring.

To circumvent demerits from the superficial hybrid or learning with inferior lexicon-based representations upon PLM, we propose a brand-new lexicon-enlightened dense (LED) retriever learning framework to inject rich lexicon information into a single

dense encoder, while keeping its sequence-level semantic representation capability. Instead of prevailing BM25 as lexicon-rich sources, we propose to leverage the recently advanced lexicon-centric representation learning model transferred from large-scale masked language modeling (MLM), and attempt to align a dense encoder with two brand-new weakened distilling objectives. On the one hand, we present lexicon-augmented contrastive learning that incorporates the hard negatives provided by lexicon-aware retrievers for contrastive training. Intuitively, the negatives given by the lexicon-aware models could be regarded as adversarial examples to challenge the dense one, so as to transfer lexical knowledge to the dense model. On the other hand, inspired by previous work [1, 8], we propose a pair-wise rank-consistent regularization as a weak supervision to guide dense model’s behavior incline to the lexicon-aware ones. Compared to distribution regularization such as KL-divergence [53] and strict fine-grained distillation like Margin-MSE [16], LED provides weak supervision signals from the lexicon-aware retrievers, leading to desirable partial knowledge injection while maintaining the dense retriever’s own properties.

We evaluate our method on three real-world human-annotated benchmarks. Experimental results show that our methods consistently and significantly improve dense retriever’s performance, even outdoing its teacher. Notably, these significant improvements are brought by the supervision of a performance-comparable lexicon-aware retriever. Besides, a detailed analysis of retrieval results shows that our knowledge distillation strategies indeed equip the dense retriever with lexicon-aware capabilities. Lastly, we show our lexicon-aware distillation strategies are compatible with the standard ranker distillation, achieving further improvement and a new state-of-the-art performance.

Our contributions are three-fold: (1) We consider improving the dense retriever by imitating the retriever based on the lexicon-aware representation model upon PLM; (2) We propose two strategies including lexicon-augmented contrastive training and pair-wise rank-consistent regularization to inject lexical knowledge into the dense retriever; (3) Evaluation results on three benchmarks show that our method brings consistent and significant improvements to the dense retriever with a comparable lexicon-aware retriever as a teacher and a new state-of-the-art performance is achieved.

## 2 RELATED WORK

Current passage retrieval systems are widely deployed as retrieve-then-rank pipelines [18, 56]. The first-stage retriever (i.e., dual-encoder) [27, 29, 35, 38, 47, 49] selects a small number of candidate passages (usually at most thousands) from the entire collection, and the second-stage ranker (i.e., cross-encoder [55]) scores these candidates again to provide a more accurate passages order. In this paper, we focus on enhancing the first-stage retriever.

**Dense Retriever.** Built upon Pre-trained Language Models [9, 30], dense retriever [20, 35] is to capture the semantic meaning of an entire sequence by encoding sequential text as a continuous representation into a low-dimensional space (e.g., 768). In this way, the dense retriever could handle vocabulary and semantic mismatch issues within the traditional term-based techniques like BM25 [40]. To train a better dense retriever, various techniques are proposed for providing hard negatives including reusing in-batch negatives [20,

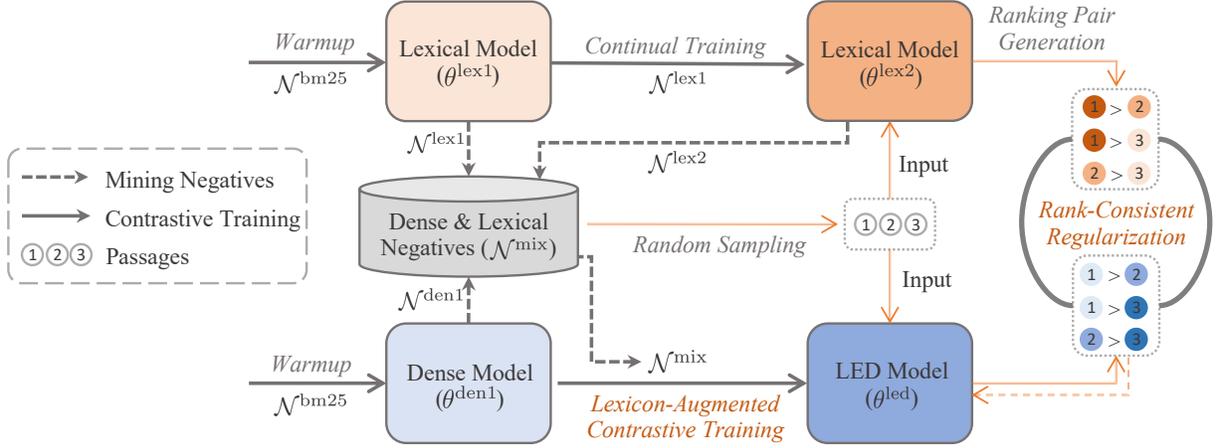
31, 35], iteratively sampling [49], mining by a well-trained model or dynamic sampling [52, 53], and denoising by cross-encoder [35]. To build retrieval-specific pre-trained language models, Lee et al. [23] proposed an unsupervised pre-training task, namely Inverse Cloze Task (ICT), Gao and Callan [12] decoupled model architecture during pre-training and further designed corpus-level contrastive learning [13] for better passage representations.

**Lexicon-Aware Retriever.** Another paradigm of work [11, 14, 26] takes advantage of strong PLMs to build lexicon-aware sparse retrievers by term-importance [7, 26] and top coordinate terms [10, 11]. These models have lexical properties and could be coupled with inverted indexing techniques. Based on contextualized representation generated by PLMs [9], Dai and Callan [7] learned to estimate individual term weights, Mallia et al. [33] further optimized the sum of query terms weights for better term interaction, COIL-series works used token-level interactions on weight vector [14] or scalar [26] to obtain exact word matching scores, and Formal et al. [10, 11] trained a retriever encoding passages as vocabulary-size highly sparse embeddings. Recently, Chen et al. [4] trained a PLM-based retriever from scratch with data generated by BM25. The trained lexicon-aware retriever could encode texts as low-dimensional embeddings and have identical lexical properties and comparable performance with BM25.

**Hybrid Retriever.** Arguably, dense sequence-level retrievers and lexicon-aware retrievers have distinctive pros and are complementary to each other. This fact triggered researchers to investigate how to combine their advantages, such as direct score aggregation [22], weighted sum [25, 26, 28, 31, 46], multiplication [50], or concatenation [4, 32, 43] in an ensemble system. The above hybrid retrievers require two dense and lexical retrievers to first compute individually and then combine their results in feature-level [4, 32, 43] or relevance-score-level [22, 25, 26, 28, 31, 46, 50] to obtain a final result. In contrast, our method only needs one model to achieve both dense and lexicon-aware retrieval behaviors, significantly decreasing the memory footprint and inference speed meanwhile providing a more in-depth fusion of lexicon-aware and dense retrieval views.

**Knowledge Distillation.** Cross-encoder empirically outperforms dual-encoder for it inputs query and passage as a whole, so that attention mechanism will be applied between them, leading to in-depth token-level interactions. Its superior performance motivates many works [16, 53] to enhance dual-encoders by knowledge distillation from cross-encoder. KL-Divergence, which minimizes the distances of distributions between teacher and student, has proven effective in many works [42, 53]. Margin-MSE [16] aims to minimize the difference of margins in two passages, and it’s been applied in later works [10, 17]. Reddi et al. [37] used the teacher’s top passages as positive examples to teach students point-wisely. ListNet [3, 48] ensured the consistency of list-wise ranking order by minimizing the difference in score distributions over passages.

The above methods are designed for the cross-encoder teacher that is much stronger than the dense student. But in our situation, the lexicon-aware teacher can only achieve comparable performance with a dense retriever. In practice (i.e., Tab. 2), we find that existing distillation methods may not be perfect choices in our



**Figure 1: The training framework of our LED retriever. The Lexical teacher is independently trained following a two-stage process. After warming up, LED is trained with negatives mined by self and two lexicon-aware retrievers for lexicon-augmented contrastive learning, during which the Lexical Model ( $\theta^{\text{lex}2}$ ) enhances LED with pair-wise rank-consistent regularization.**

setting. Therefore, we propose two novel strategies, namely lexicon-augmented contrastive training and pair-wise rank-consistent regularization to transfer lexical knowledge.

### 3 METHODOLOGY

We first introduce the task formalization, general training framework, and retriever architectures in Sec. 3.1. Then we present our Lexical- Enlightened Dense (LED) retriever in Sec. 3.2.

#### 3.1 Preliminary

**Task Definition.** In the first-stage retrieval, given a query  $q$ , a retriever is required to fetch top- $k$  relevant passages from a million-even billion-scale passage collection  $C$ . Due to the efficiency requirement, dual-encoder architecture is widely applied in this task for its lightweight metric calculation. Formally, dual-encoder represents text  $x$  (could be query  $q$  or passage  $p$ ) to  $d$ -dimensional embeddings, i.e.,

$$\mathbf{x} = \text{Dual-Enc}(x; \theta) \in \mathbb{R}^d, \quad (1)$$

where  $\theta$  could be dense retriever ( $\theta^{\text{den}}$ ) or lexicon-aware retriever ( $\theta^{\text{lex}}$ ). With separately encoded query  $q$  and passage  $p$ , we could calculate the relevance score via dot product for retrieval, i.e.,

$$\mathcal{R}(q, p; \theta) = \mathbf{q}^T \mathbf{p}. \quad (2)$$

The dual-encoder architecture and lightweight dot product evaluation enable us to encode and index all passages in the collection  $C$  beforehand, so we only need to encode the given query for online retrieval, achieving more efficiency.

**Learning Framework for Retriever.** To train the dual-encoder  $\theta$ , we utilize contrastive learning following previous works [15, 49]. Specifically, with a given query  $q$ , a labeled positive passage  $p^+$ , and negative passages  $\mathcal{N}$ , contrastive loss can be applied to optimize the dual-encoder  $\theta$  by maximizing the relevance of the  $q$  and  $p^+$

while minimizing that of  $q$  and  $p \in \mathcal{N}$ , i.e.,

$$\mathcal{L}_\theta^{\text{cl}} = -\log \frac{\exp(\mathcal{R}(q, p^+; \theta))}{\sum_{p \in \{p^+\} \cup \mathcal{N}} \exp(\mathcal{R}(q, p; \theta))}, \quad (3)$$

where negative passage set  $\mathcal{N}$  can be generated from top-ranked non-answer passages in retrieval results of BM25 model [34] or a trained retrievers [52, 53], i.e.,

$$\mathcal{N} = \{p \mid p \sim P(C \setminus \{p^+\} \mid q; \theta^{\text{samp}})\}, \quad (4)$$

where  $P$  is a probability distribution over  $C$ , which can be defined as non-parametric (e.g.,  $\theta^{\text{samp}} = \emptyset$ ) or parametric (e.g.,  $\theta^{\text{samp}} \neq \emptyset$ ).

**Dense & Lexicon-Aware Retrievers.** Both dense retriever ( $\theta^{\text{den}}$ ) and lexicon-aware retriever ( $\theta^{\text{lex}}$ ) follow dual-encoder architecture and the encoders are built upon PLMs like BERT [9]. Precisely, a PLM ( $\theta^{\text{plm}}$ ) encodes a given text (i.e., query  $q$  or passage  $p$ ),  $x = \{t_1, t_2, \dots, t_n\}$ , to contextualized embeddings, i.e.,

$$\begin{aligned} \mathbf{H}^x &= \text{PLM}(x; \theta^{\text{plm}}) \\ &= \text{PLM}([\text{CLS}], t_1, t_2, \dots, t_n, [\text{SEP}]; \theta^{\text{plm}}), \end{aligned} \quad (5)$$

eventually  $\mathbf{H}^x = [\mathbf{h}_{[\text{CLS}]}^x, \mathbf{h}_1^x, \dots, \mathbf{h}_n^x, \mathbf{h}_{[\text{SEP}]}^x]$ . [CLS] and [SEP] are special tokens designed for sentence representation and separation by recent PLMs [9, 30]. Dense retriever [35, 49] represents text by using the embedding of special token [CLS] (i.e.,  $\mathbf{h}_{[\text{CLS}]}^x$ ) as follows,

$$\mathbf{x}^{\text{den}} = \text{Dual-Enc}(x; \theta^{\text{den}}) = \text{CLS-Pool}(\mathbf{H}^x), \quad (6)$$

where  $\theta^{\text{den}} = \theta^{\text{plm}}$  with no additional parameters.

For lexicon-aware retriever, we adopt SPLADE [10] which learns to predict the weights of terms in PLM vocab for each token in the input  $x$  by the Masked Language Modeling (MLM) layer and sparse regularization, then max-pooling these weights into a discrete text representation after log-saturation. Formally, with  $\mathbf{H}^x$  encoded by the PLM ( $\theta^{\text{plm}}$ ), a MLM layer ( $\theta^{\text{mlm}}$ ) linearly transform it into  $\tilde{\mathbf{H}}^x$ ,

then term weight representation of  $x$  could be obtained as follows,

$$\begin{aligned} \mathbf{x}^{\text{lex}} &= \text{Dual-Enc}(x; \theta^{\text{lex}}) \\ &= \text{MAX-Pool}(\log(1 + \text{ReLU}(\mathbf{W}^e \cdot \tilde{\mathbf{H}}^x))), \end{aligned} \quad (7)$$

where  $\mathbf{W}^e \in \mathbb{R}^{|V| \times e}$  is the transpose of the input embedding matrix in PLM as the MLM head, and the  $\theta^{\text{lex}} = \{\theta^{\text{plm}}, \theta^{\text{mlm}}, \mathbf{W}^e\}$ .

The dense encoder represents texts as global sequence-level embeddings and is good at global semantic matching, while the lexicon-aware encoder represents local term-level embeddings and handles salient phrases and entity mentions well. Both encoders can be optimized with the Eq. 3.

### 3.2 Lexical Enlightened Dense Retriever

Fig. 1 illustrates the training workflow of our LED retriever. Specifically, we follow a two-stage training procedure. In the *Warmup* stage, we independently train the dense and lexicon-aware retrievers by Eq. 3, both with BM25 negatives ( $\mathcal{N}^{\text{bm25}}$ ). This stage ends up with two retrievers, namely the Lexical Warm-up ( $\theta^{\text{lex1}}$ ) and the Dense Warm-up ( $\theta^{\text{den1}}$ ).

Then, we sample negative passages ( $\mathcal{N}^{\text{lex1}}$ ) with the Lexical Warm-up checkpoint ( $\theta^{\text{lex1}}$ ) for the second stage, namely *Continual Training* stage. With the fixed negative passages ( $\mathcal{N}^{\text{lex1}}$ ), we continually train the lexical retriever initialized from the warm-up checkpoint ( $\theta^{\text{lex1}}$ ) by Eq. 3. After the second stage, we could obtain the model named Lexical ( $\theta^{\text{lex2}}$ ), which plays a role of a teacher for later lexical knowledge teaching.

With a well-trained lexicon-aware teacher ( $\theta^{\text{lex2}}$ ) and dense student after warming up ( $\theta^{\text{den1}}$ ), we enlighten the student by transferring knowledge from the teacher. The knowledge transfer is achieved from two perspectives – 1) a lexicon-augmented contrastive objective to challenge the dense encoder and 2) a rank-consistent regularization to make the dense model’s behavior inclined to its lexicon-aware teacher. We will detail the two objectives in the following paragraphs.

**Lexicon-Augmented Contrastive.** Following previous work [52], we use the dense negatives ( $\mathcal{N}^{\text{den1}}$ ) sampled by the Dense Warm-up ( $\theta^{\text{den1}}$ ) to boost dense retrieval. Meanwhile, inspired by Chen et al. [4] who trained a lexical retriever with negatives provided by term-based techniques such as BM25, we try to enhance the dense retriever from the negatives augmentation perspective.

Considering the similar backbone (i.e., both are fine-tuned on PLMs) and the same optimization objectives (i.e., Eq. 3) of both dense and lexicon-aware retrievers, their significant differences in retrieval behaviors may partially stem from training with different negative passages. So intuitively, dense one can use the lexical negatives ( $\mathcal{N}^{\text{lex1}}$ ) to partially imitate the training process of the lexical teacher ( $\theta^{\text{lex2}}$ ), thus learning lexicon-aware ability. One step further, we use the lexical negatives ( $\mathcal{N}^{\text{lex2}}$ ) for learning more and harder lexical knowledge. Meanwhile, compared to the dense negatives ( $\mathcal{N}^{\text{den1}}$ ), these lexical negatives ( $\mathcal{N}^{\text{lex1}}$  and  $\mathcal{N}^{\text{lex2}}$ ) can provide more diverse examples and could be regarded as adversarial examples to challenge the dense retriever for robust retriever training. Formally,

the lexicon-augmented contrastive loss for LED is,

$$\mathcal{L}_{\theta^{\text{led}}}^{\text{cl}} = -\log \frac{\exp(\mathcal{R}(q, p^+; \theta^{\text{led}}))}{\sum_{p \in \{p^+\} \cap \mathcal{N}^{\text{mix}}} \exp(\mathcal{R}(q, p; \theta^{\text{led}}))}, \quad (8)$$

where  $\mathcal{N}^{\text{mix}} = \{\mathcal{N}^{\text{lex1}} \cap \mathcal{N}^{\text{lex2}} \cap \mathcal{N}^{\text{den1}}\}$ .

**Rank-Consistent Regularization.** From the retrieval behavior perspective, for given query-passage pairs, we utilize the lexicon-aware teacher ( $\theta^{\text{lex2}}$ ) to generate ranking pairs to regularize and guide LED’s retrieval behavior.

Specifically, given a query  $q$  and passages from  $\mathcal{D}^q = \{p \in \{p^+\} \cap \mathcal{N}^{\text{mix}}\}$ , the Lexical ( $\theta^{\text{lex2}}$ ) scores each query-passage pair (abbr.  $\mathcal{R}(p; \theta^{\text{lex2}})$ ) with Eq. 2 and generate ranking pairs as follows,

$$\mathcal{K}^q = \{(p_i, p_j) | p_i, p_j \in \mathcal{D}^q, \mathcal{R}(p_i; \theta^{\text{lex2}}) > \mathcal{R}(p_j; \theta^{\text{lex2}})\}. \quad (9)$$

Then, a pair-wise rank-consistent regularization is employed to make the dense model’s behavior incline to the lexicon-aware one by minimizing the following margin-based ranking loss,

$$\mathcal{L}_{\theta^{\text{led}}}^{\text{ll}} = \frac{1}{|\mathcal{K}^q|} \sum_{p_i, p_j \in \mathcal{K}^q} \max[0, \mathcal{R}(p_i; \theta^{\text{led}}) - \mathcal{R}(p_j; \theta^{\text{led}})], \quad (10)$$

where  $\mathcal{R}(p; \theta^{\text{led}})$  is the abbreviation of relevance  $\mathcal{R}(q, p; \theta^{\text{led}})$  calculated by LED ( $\theta^{\text{led}}$ ) with Eq. 2. Compared to logits distillation with a list-wise KL-divergence loss, our training objective provides a weak supervision signal pair-wisely, thus keeping the effects of injecting lexical knowledge on dense properties at a minimum level. Especially since we don’t punish the dense student as long as its ranking of a given pair is the same as the teacher, without strict requirements on the score gap like Margin-MSE [16]. Experiments in Tab. 2 demonstrate the merit of our method.

**Training and Inference.** To incorporate lexicon-aware ability while keeping its sequence-level semantic representation ability for passage retrieval, we combine contrastive loss ( $\mathcal{L}_{\theta^{\text{led}}}^{\text{cl}}$ ) in Eq. 8 and lexical learning loss ( $\mathcal{L}_{\theta^{\text{led}}}^{\text{ll}}$ ) in Eq. 10 to train our LED retriever ( $\theta^{\text{led}}$ ) as follows,

$$\mathcal{L}_{\theta^{\text{led}}} = \mathcal{L}_{\theta^{\text{led}}}^{\text{cl}} + \lambda \mathcal{L}_{\theta^{\text{led}}}^{\text{ll}}, \quad (11)$$

where  $\lambda$  is a hyperparameter to control how intensive the training inclines to transfer lexical-aware knowledge from the lexicon-aware teacher ( $\theta^{\text{lex2}}$ ).

For inference, LED pre-computes the embeddings of all passages in the entire collection  $\mathcal{C}$  and builds indexes with FAISS [19]. Then LED encodes queries online and retrieves top-ranked  $k$  passages based on the relevance score.

**Remark.** Our framework injects lexicon-aware capability into a sequence-level representation model, showing two-fold advantages in comparison to previous methods superficially combining dense and lexicon-aware retrievers: 1) Compared to using two separate PLM-based dense and lexicon-aware retrievers [4, 26, 28, 44], our LED retriever could achieve hybrid retrieval results and comparable performances with only one model. We no longer need to maintain multiple index systems for both retrieval models in an ensemble system or to encode an online query twice with different retrievers, reducing memory footprint and inference time. 2) Compared to fusing PLM-based dense and traditional term-matching retrievers like BM25 [4, 24–26, 28, 31, 32, 43, 46], our single method could

achieve better results since the lexicon-aware capability of LED is learned from a strong retriever (shown in Tab. 3).

## 4 EXPERIMENTS

We evaluate our retriever on three public human-annotated real-world benchmarks, namely MS MARCO [34], TREC Deep Learning 2019 [6], and TREC Deep Learning 2020 [5]. MS-MARCO Dev has 6980 queries, TREC 2019 has 43 queries, and TREC 2020 includes 54 queries. In all three benchmarks, first-stage retrievers are required to fetch relevant passages from an 8-million scale collection. We report MRR@10, Recall@50, and Recall@1000 for MS MARCO Dev, as well as NDCG@10 for both TREC Deep Learning 2019 and TREC Deep Learning 2020. For all three datasets, we use the official TREC evaluation files to conduct the evaluation protocol.<sup>2</sup>

### 4.1 Baselines

We compare with previous state-of-the-art baselines including traditional term-based techniques like BM25 [40], and dense [13, 17, 21, 28, 38, 39, 42, 49, 51–53] as well as lexicon-aware retrievers [4, 7, 10, 14, 26]. More details about the baselines are provided in Appendix A. We report the models used during the two-stage training pipeline for a detailed comparison. For lexicon-aware retrievers, we report the models after the warm-up training, namely LEX (Warm-up), and continual training, namely LEX (Continue). Note that LEX (Continue) is the lexical teacher used for teaching. Similarly, for dense retrievers, we show DEN (Warm-up) and DEN (Continue). Note that the DEN (Warm-up) is the dense student where our LED model starts from and DEN (Continue) is independently trained with the hard negatives provided by DEN (Warm-up) by Eq. 3. We also report the DEN (Continue) further enhanced with a strong ranker distillation (i.e., DEN (w/ RT)).

### 4.2 Implementation Details

All experiments run on 1 NVIDIA Tesla A100 GPU having 80GB memories with a fixed random seed. We train our models with mixed precision to speed up the training and meet the huge memory need. The training time will last about 32 hours. For the lexical teacher, we train a DistilBERT [41] following SPLADE-max [10].<sup>3</sup> Following previous work [55], in the warm-up stage, we train the lexical retriever with batch size 48, 5 negatives for each query randomly sampled from BM25 negatives, and a learning rate  $3e^{-5}$  for three epochs. In the second stage, we remain all hyperparameters unchanged except lower the learning rate to  $2e^{-5}$  and use negative passages randomly sampled from the top 200 self-mining hard negatives. For dense retriever, we train coCondenser [13] checkpoint with batch size 16, 7 negatives per query, and a learning rate of  $1e^{-5}$  for three epochs.<sup>4</sup>

Particularly, in the LED training stage, with other hyperparameters unchanged, we set the learning rate  $5e^{-6}$  and randomly select 32 hard negatives from the mixture of each top 200 negatives mined by the warm-up dense student, warm-up lexical teacher, and final

lexical teacher. The number of negatives per query 32 is selected from {8, 16, 24, 32}. The higher number of negatives per query indicates the more pair-wise ranks constructed by the lexical teacher, leading to more lexical knowledge transfer.

For rank-consistent regularization, we set loss weight  $\lambda = 1.2$  after searching from {1.0, 1.2, 1.5, 1.8, 2.0}.

### 4.3 Main Results

As shown in Tab. 1, we present the evaluation results on the aforementioned three public benchmarks.

Firstly, our LED retriever achieves comparable performance with state-of-the-art methods ColBERTv2 [42] and AR2 [53] on MS MARCO Dev, although both baselines are taught by the powerful ranking model (i.e., cross-encoder). After coupling with a similar ranker distillation, our LED retriever (i.e., LED (w/ RT)) can be further improved and meanwhile outperforms state-of-the-art baselines on all three datasets, showing the compatibility of distillation from lexicon-aware sparse retriever.<sup>5</sup> Note that we neither use heavy ranker teacher in AR2 [53] nor multiple vector representation applied in ColBERTv2 [42].

Secondly, LED (w/ RT) achieves better performance than DEN (w/ RT) on all three datasets, demonstrating that our training method can transfer some complementary lexicon-aware knowledge not covered by the cross-encoder. The weak intensity of the supervision signal makes our lexical-enlightened strategy a promising plug-and-play technique for other dense retrievers.

Thirdly, our LED retriever taught by a smaller lexicon-aware retriever is similarly performant as the dense retriever taught by a strong cross-encoder (i.e., DEN (w/ RT)), showing the effectiveness of injecting lexical knowledge into the dense retriever. The reasons are two-fold: (1) The dual-encoder architecture of the lexicon-aware teacher enables the relevance calculation can be easily integrated into in-batch techniques to scale up the teaching amount. (2) More importantly, lexicon-aware retriever could provide self-mining hard negatives for more direct supervision while cross-encoder can only provide score distribution over given passages.

### 4.4 Further Analysis

**Comparison of Teaching Strategies.** Tab. 2 shows the comparison of our proposed pair-wise rank-consistent regularization with other teaching strategies. Filter means the negatives with high scores (i.e., false negatives) are filtered by LEX (Continue). The other three strategies (e.g., Margin-MSE, ListNet, and KL-Divergence) are borrowed from knowledge distillation in IR domain. From the table, we can find that all strategies can bring performance gain, even in an indirect way like Filter. This observation proves that learning from the lexicon-aware representation model leads to a better dense retriever. Also, our rank-consistent regularization outperforms other baselines on MRR@10 metric by a large margin, showing the superiority of our method. Besides, we can find that the point-wise objective (i.e., Margin-MSE) brings the least gain, followed by the list-wise objectives (i.e., ListNet and KL-Divergence) and our pair-wise rank-consistent regularization brings the most significant

<sup>2</sup>[https://github.com/usnistgov/trec\\_eval](https://github.com/usnistgov/trec_eval)

<sup>3</sup><https://huggingface.co/distilbert-base-uncased>

<sup>4</sup><https://huggingface.co/Luyu/co-condenser-marco>. We chose the coCondenser as the dense retriever due to its retrieval-oriented pre-training and its superior performance compared to vanilla BERT and DistilBERT checkpoints. We tested our proposed strategies on these models as well and saw similar improvements as with coCondenser.

<sup>5</sup>Like AR2 [53] and ColBERTv2 [42], we use KL-divergence to distill the ranker’s scores into the LED model, but we use ERNIE-2.0-base [45] instead of ERNIE-2.0-large in AR2. The KL loss is directly added to Eq. 11 during the training of LED.

**Table 1: Experimental results on MS MARCO, TREC DL 2019 (DL'19), and TREC DL 2020 (DL'20) datasets (%). We mark the best results in bold and the second-best underlined. Numbers marked with “\*” mean that the improvement is statistically significant compared with the baseline (t-test with  $p$ -value  $< 0.05$ ).**

Methods	PLM	Ranker Taught	Multi Vector	MS MARCO Dev			DL'19	DL'20
				MRR@10	R@50	R@1k	NDCG@10	NDCG@10
<i>Lexicon-Aware Retriever</i>								
BM25 [40]	-			18.7	59.2	85.7	50.6	48.0
DeepCT [7]	BERT <sub>base</sub>			24.3	69.0	91.0	55.1	55.6
COIL-full [14]	BERT <sub>base</sub>			35.5	-	96.3	70.4	-
UniCOIL [26]	BERT <sub>base</sub>			35.2	80.7	95.8	-	-
SPLADE-max [10]	DistilBERT			34.0	-	96.5	68.4	-
DistilSPLADE-max [10]	DistilBERT	✓		36.8	-	97.9	<u>72.9</u>	-
UniCOIL $\Lambda$ [4]	BERT <sub>base</sub>			34.1	82.1	97.0	-	-
<i>Dense Retriever</i>								
ANCE [49]	RoBERTa <sub>base</sub>			33.0	-	95.9	64.5	64.6
ADORE [52]	RoBERTa <sub>base</sub>			34.7	-	-	68.3	66.6
TAS-B [17]	DistilBERT	✓		34.7	-	97.8	71.7	68.5
TAS-B + CL-DRD [51]	DistilBERT	✓		38.2	-	-	72.5	68.7
TCT-ColBERT [28]	BERT <sub>base</sub>	✓		35.9	-	97.0	71.9	-
ColBERTv1 [21]	BERT <sub>base</sub>		✓	36.0	82.9	96.8	67.0	66.8
ColBERTv2 [42]	BERT <sub>base</sub>	✓	✓	<u>39.7</u>	<u>86.8</u>	<b>98.4</b>	72.0	62.1
coCondenser [13]	BERT <sub>base</sub>			38.2	-	<b>98.4</b>	-	-
PAIR [38]	ERNIE <sub>base</sub>	✓		37.9	86.4	98.2	-	-
RocketQAv2 [39]	ERNIE <sub>base</sub>	✓		38.8	86.2	98.1	-	-
AR2-G [53]	BERT <sub>base</sub>	✓		39.5	-	-	-	-
<i>Our Models</i>								
LEX (Warm-up)	DistilBERT			36.1	84.2	97.5	67.4	66.4
LEX (Continue)	DistilBERT			38.3	85.9	98.0	72.8	67.7
DEN (Warm-up)	BERT <sub>base</sub>			36.1	83.5	97.7	64.7	65.9
DEN (Continue)	BERT <sub>base</sub>			38.1	86.3	<b>98.4</b>	69.1	67.8
DEN (w/ RT)	BERT <sub>base</sub>	✓		39.6	86.7	<b>98.4</b>	71.8	<u>69.7</u>
LED	BERT <sub>base</sub>			39.6	86.6	<u>98.3</u>	70.5	67.9
LED (w/ RT)	BERT <sub>base</sub>	✓		<b>40.2*</b>	<b>87.6*</b>	<b>98.4</b>	<b>74.4*</b>	<b>70.2*</b>

**Table 2: Evaluation results of different teaching strategies on MS MARCO Dev (%). “\*” refers to statistical significance.**

Methods	MRR@10	R@1k
No Distillation	38.1	<b>98.4</b>
Filter [35]	38.4	<b>98.4</b>
Margin-MSE [16]	38.5	98.3
ListNet [48]	38.7	98.2
KL-Divergence [53]	39.0	<b>98.4</b>
Ours	<b>39.6*</b>	98.3

gain. The phenomenon implies that a soft teaching objective is more functional for transferring knowledge from the lexicon-aware model than strict objectives. In fact, enforcing dense retrievers to be aligned with fine-grained differences between scores of the LEX often leads to training collapse. Concretely, only equipped with

carefully chosen hyperparameters, especially small distillation loss weight, Margin-MSE can enhance the dense retriever.

**Comparison of Ensemble Retrievers.** We are also curious whether our LED can improve the performance of ensemble retrievers. With LEX (Continue) ( $\theta^{\text{lex}2}$ ) and LED ( $\theta^{\text{led}}$ ), we simply use the summation of the normalized relevance scores of two retrievers, and then return a new order of retrieval results. Tab. 3 gives the evaluation results of our systems and other strong baselines reported in previous work [4, 26, 28]. Note that previous work [26, 28] utilized weighted score sum after hyper-parameter searching while we directly sum the normalized scores of two retrievers without any tuning. From the results in Tab. 3, we can observe:

(1) Aligned with results in SPAR [4], the ensemble of two dense retrievers (i.e., DEN (Continue) + LED and DEN (Continue) + DEN (w/ RT)) is not as performant as that of one dense and one lexicon-aware retriever. In particular, the ensemble of two dense retrievers is even less competitive than a single LED or DEN (w/ RT). The results are rational because two base models have similar retrieval

**Table 3: Comparison with Ensemble Systems on MS MARCO Dev (%). The first block results are copied from [4, 26, 28].  $\Lambda$  [4] refers to a dense retriever trained with data generated by lexicon-based methods such as BM25 and UniCOIL. “\*\*” indicates statistical significance compared to their counterparts without our training strategies.**

Ensemble Systems	MRR@10	R@50	R@1k
TCT-ColBERT + BM25 [28]	36.9	-	-
TCT-ColBERT + UniCOL [26]	37.8	-	-
TCT-ColBERT + UniCOL [26]	38.2	-	-
ANCE + BM25 [4]	34.7	81.6	96.9
RocketQA + BM25 [4]	38.1	85.9	98.0
RocketQA + UniCOIL [4]	38.8	86.5	97.3
RocketQA + BM25 $\Lambda$ [4]	37.9	85.7	98.0
RocketQA + UniCOIL $\Lambda$ [4]	38.6	86.3	98.5
<hr/>			
DEN (Continue) + BM25	30.4	87.1	98.6
DEN (Continue) + LED	39.3	86.9	98.5
DEN (Continue) + DEN (w/ RT)	39.4	87.0	98.5
DEN (Continue) + LEX (Continue)	40.4	88.4	<b>98.7</b>
DEN (w/ RT) + LEX (Continue)	40.7	88.4	<b>98.7</b>
LED + LEX (Continue)	<b>40.9*</b>	88.3	98.6
LED (w/ RT) + LEX (Continue)	<b>41.1*</b>	<b>88.5</b>	<b>98.7</b>

behaviors and the strong one will be impeded by the weak one if they have the same weight in the ensemble system. The latter reason could also be used to explain why the ensemble of dense and the traditional term-based technique like BM25 (i.e., DEN (Continue) + BM25) is less good than the single DEN (Continue).

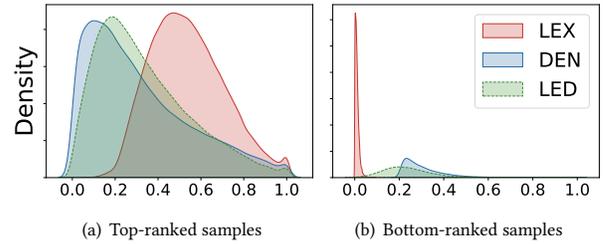
(2) Although coupling with LEX (Continue) will not introduce new knowledge to the hybrid ensemble system where LED is the base retriever, LED + LEX (Continue) can further boost the performance of DEN (Continue) + LEX (Continue). The reason behind this is that the LED scores golden query-passage pairs higher than DEN, so these pairs are ranked higher in the later ensemble process. This behavior could be regarded as an instance-level weighted score aggregation inside the network and it is more feasible to obtain than tuning the weights of retrievers for each query in the ensemble system. This observation could from the side prove that our dense and lexicon-aware abilities fusion inside the network is better than a superficial ensemble.

(3) LED (w/ RT) + LEX (Continue) is slightly better than LED + LEX (Continue) and DEN (w/ RT) + LEX (Continue), once again proving that our lexical rank-consistent regularization is complementary to the ranker distillation.

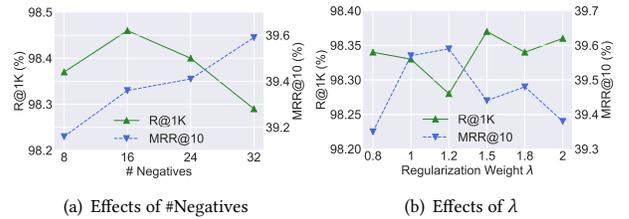
**Impact of Different Components.** We conduct an ablation study to further investigate the impact of lexical hard negatives and rank-consistent regularization method. Tab. 4 reports the results of removing each component. We can observe that pair-wise rank-consistent regularization plays an important role in lexical learning because removing it will bring significant performance degradation on MRR@10 metric. In addition, we can find that both negatives provided by LEX (Warm-up) and LEX (Continue) are both helpful for the contrastive training of the dense retriever, and removing both of them results in a more obvious performance drop.

**Table 4: Ablation Study on MS MARCO Dev (%). Negs is short for negatives. “\*” indicates statistical significance.**

Retrievers	MRR@10	R@1k
LED	<b>39.6*</b>	98.3
w/o Rank Regularization	37.9	<b>98.5</b>
w/o LEX Continue Negs ( $\mathcal{N}^{\text{lex2}}$ )	39.4	98.3
w/o LEX Warm-up Negs ( $\mathcal{N}^{\text{lex1}}$ )	39.4	98.3
w/o LEX Mixed Negs ( $\mathcal{N}^{\text{lex1}} \cap \mathcal{N}^{\text{lex2}}$ )	39.2	98.4



**Figure 2: Distributions of model prediction for DEN (Continue), LEX (Continue), and LED retrievers over MS MARCO Dev. For visual clarity, we use the query-passage pairs which the LEX and DEN predict discrepantly as data samples. The discrepancy is determined by that there is a  $> 0.2$  margin between their predicted scores normalized over passages retrieved for a  $q$ . To ensure diversity, we consider two normalization cases, LEX-biased pairs (i.e., LEX’s top-100) and LEX-unbiased pairs (i.e., LEX’s bottom-100 out of 1000).**



**Figure 3: (a) Effects of the number of negatives per query on MS MARCO Dev. (b) Effects of the regularization weight  $\lambda$  on MS MARCO Dev.**

**Effects on Model Predictions.** To further check the effects of learning lexicon-aware capability on the LED, we illustrate the distribution shift of predictions of dense retrievers before and after lexical enlightenment in Fig. 3. We can make the following observations: (1) In both two sets of query-passage pairs, compared to DEN distributions, the score distributions of LED are clearly shifting to the LEX, showing the success of lexical knowledge learning. (2) LED’s distribution remains more overlaps with DEN instead of LEX, which proves that our rank-consistent regularization method could keep LED’s dense retriever properties, thanks to the weak supervision signal.

**Table 5: The average rank of golden passages by four retrievers on MS MARCO Dev. We bin the dev examples into buckets with the rank predicted by LEX and calculate the average ranking of other retrievers by group.**

Ranges	Count	Average Ranking			
		LEX	DEN	LED	LED (w/ RT)
Top 1	1,787	1.0	2.5	2.3	2.4
(1, 5]	2,242	3.1	6.4	5.2	4.8
(5, 10]	875	7.8	14.7	13.8	12.7
(10, 50]	1,428	23.3	31.4	31.1	28.9
(50, 100]	358	70.5	80.0	75.9	74.0
(100, 500]	445	216.8	156.3	166.7	154.1
(500, 1000]	69	698.0	298.9	334.7	289.1

**Impacts of Hyperparameters.** We conduct extra experiments to explore the impact of hyperparameters on LED retriever training. Fig. 3(a) illustrates the impact of changing negative passages on the LED. We can observe that as the number of negative passages increases, the MRR@10 performance goes up and the R@1k performance reaches the peak when 16 and decreases gradually. The main table shows that Lexical is less performant than Dense at R@1k metric ( $98.0 < 98.4$ ). So the trend of increasing the number of negatives proves that imitating too much the lexical retriever will also be negatively influenced by the weakness of the teacher. These two trends indicate that, with more negatives, the teacher will construct more rank pairs for more lexical knowledge transfer. Fig. 3(b) shows the performance with regard to different regularization weights  $\lambda$ . It is observed that the performances don't fluctuate significantly as the weight  $\lambda$  changes, demonstrating the robustness of lexical enhancement strategies. Interestingly, the increase in MRR@10 comes with the drop in R@1k to some extent, once again showing that a well-enhanced LED also inherits the weakness of Lexical.

**Zoom-in Study of Retrieval Ranking.** Tab. 5 shows how the average rank of golden passages varies across different rank ranges, bucketed by LEX-predicted ranks of the golden passages. We can observe that: (1) More than 50% golden passages are ranked in the top 5 by the LEX, paving the way for good lexical teaching. (2) The average ranking of golden passages by LED is consistently improved until the top 100, which means approximately 90% of answers are ranked higher by the retriever after learning lexical knowledge, proving the effectiveness of our lexical knowledge transfer. Meanwhile, similar even more gain can also be observed in LED (w/ RT), once again proving that our method is complementary to distillation from a cross-encoder. (3) In the queries that the LEX performs unfavorably (i.e., ground truth ranked lower than 100), LED and LED (w/ RT) are negatively impacted by the lexicon-aware teacher's mistakes. Interestingly, their original rankings of these ground truths are not very high, either. So these queries are intractable for both dense and lexicon-aware retrievers, which we leave for future work.

#### 4.5 Case Study

Tab. 6 shows the three case queries with rankings of 4 retrievers. With lexicon-aware capability, LED and LED w/ RT could retrieve

**Table 6: Case study on MS-Marco Dev. 'Passage+' denotes the golden passage of the corresponding query. 'Rank' indicates the ranking of golden passage by retrievers.**

<b>Query</b>	ID: 1090413// state the benefits of internet
<b>Passage+</b>	ID: 7998365// What Are Some Benefits of Using the Internet?<sep>Some of the <b>benefits of the Internet</b> include reduced geographical distance and fast communication. The Internet is also a hub of information where users can simply upload, download and publish ideas...
<b>Rank</b>	LEX: 1; DEN: 3; LED: 1; LED w/ RT: 1
<b>Retrieved</b>	<b>DEN's 1st.</b> ID: 7339157 // -<sep>Advantages of the Internet. The Internet provides opportunities galore, and can be used for a variety of things. Some of the things that you can do via the Internet are: 1 E-mail: E-mail is an online correspondence system. 2 With e-mail you can send and receive instant electronic messages, which works like writing letters.
<b>Query</b>	ID:1033652// what is the purpose of pencil tool
<b>Passage+</b>	ID: 7212314// Pencil<sep>Should I remove Pencil by Evolus Co? Pencil is built for the purpose of providing a free and <b>open-source GUI prototyping tool</b> that people can easily install and use to create mockups in popular desktop platforms.
<b>Rank</b>	LEX: 1; DEN: 11; LED: 1; LED w/ RT: 1
<b>Retrieved</b>	<b>DEN's 1st.</b> ID:313304 // Pencil<sep>This article is about the writing implement. For other uses, see Pencil (disambiguation). A pencil is a writing implement or art medium constructed of a narrow, solid pigment core inside a protective casing which prevents the core from being broken or leaving marks on the users hand during use.

golden passages as the top-1 result like their teacher LEX. In particular, in the first case, DEN mismatches the "benefits" in the query with "advantages" in the passage since they are both positive words. On the contrary, the LEX and LED-series retrievers could exactly match the phrase "benefits of the Internet". In the second query, the "Pencil tool" refers to a specific GUI prototyping tool (as highlighted in the positive passage). DEN misunderstands the mention "Pencil tool" in the query and returns passages about the vanilla pencil, which is non-relevant to the user's intention. The above two cases show that retrievers with lexicon-aware capability (i.e., LEX, LED, LED w/ RT) could well capture the salient phrases and entity mentions, providing more precise retrieval results.

## 5 CONCLUSION

In this paper, we consider developing a lexicon-enlightened dense retriever by transferring knowledge from a lexicon-aware sparse representation model into a dense one. To achieve this end, we propose to enlighten a dense retriever from two aspects, namely the lexicon-augmented contrastive objective and the pair-wise rank-consistent regularization. Experimental results on three real-world retrieval benchmarks show that with a performance-comparable lexicon-aware representation model as the teacher, our strategies can improve a dense retriever consistently and significantly, even outdoing its teacher. Further extensive analysis and discussions demonstrate the effectiveness and compatibility of our training strategies, as well as the interpretability of the LED retriever.

## REFERENCES

- [1] Christopher J. C. Burges, Tal Shaked, Erin Renshaw, Ari Lazier, Matt Deeds, Nicole Hamilton, and Gregory N. Hullender. 2005. Learning to rank using gradient descent. In *Machine Learning, Proceedings of the Twenty-Second International Conference (ICML 2005), Bonn, Germany, August 7-11, 2005 (ACM International Conference Proceeding Series, Vol. 119)*, Luc De Raedt and Stefan Wrobel (Eds.). ACM, 89–96. <https://doi.org/10.1145/1102351.1102363>
- [2] Yinqiong Cai, Yixing Fan, Jiafeng Guo, Fei Sun, Ruqing Zhang, and Xueqi Cheng. 2021. Semantic Models for the First-stage Retrieval: A Comprehensive Review. *CoRR abs/2103.04831* (2021). arXiv:2103.04831 <https://arxiv.org/abs/2103.04831>
- [3] Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and Hang Li. 2007. Learning to rank: from pairwise approach to listwise approach. In *Machine Learning, Proceedings of the Twenty-Fourth International Conference (ICML 2007), Corvallis, Oregon, USA, June 20-24, 2007 (ACM International Conference Proceeding Series, Vol. 227)*, Zoubin Ghahramani (Ed.). ACM, 129–136. <https://doi.org/10.1145/1273496.1273513>
- [4] Kilun Chen, Kushal Lakhotia, Barlas Oguz, Anchit Gupta, Patrick S. H. Lewis, Stan Peshterliev, Yashar Mehdad, Sonal Gupta, and Wen-tau Yih. 2021. Salient Phrase Aware Dense Retrieval: Can a Dense Retriever Imitate a Sparse One? *CoRR abs/2110.06918* (2021). arXiv:2110.06918 <https://arxiv.org/abs/2110.06918>
- [5] Nick Craswell, Bhaskar Mitra, Emine Yilmaz, and Daniel Campos. 2020. Overview of the TREC 2020 Deep Learning Track. In *Proceedings of the Twenty-Ninth Text REtrieval Conference, TREC 2020, Virtual Event [Gaithersburg, Maryland, USA], November 16-20, 2020 (NIST Special Publication, Vol. 1266)*, Ellen M. Voorhees and Angela Ellis (Eds.). National Institute of Standards and Technology (NIST). <https://trec.nist.gov/pubs/trec29/papers/OVERVIEW.DL.pdf>
- [6] 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* (2020). arXiv:2003.07820 <https://arxiv.org/abs/2003.07820>
- [7] Zhuyun Dai and Jamie Callan. 2019. Context-Aware Sentence/Passage Term Importance Estimation For First Stage Retrieval. *CoRR abs/1910.10687* (2019). arXiv:1910.10687 <http://arxiv.org/abs/1910.10687>
- [8] Mostafa Dehghani, Hamed Zamani, Aliaksei Severyn, Jaap Kamps, and W. Bruce Croft. 2017. Neural Ranking Models with Weak Supervision. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval, Shinjuku, Tokyo, Japan, August 7-11, 2017*, Noriko Kando, Tetsuya Sakai, Hideo Joho, Hang Li, Arjen P. de Vries, and Ryan W. White (Eds.). ACM, 65–74. <https://doi.org/10.1145/3077136.3080832>
- [9] 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, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, Jill Burstein, Christy Doran, and Thamar Solorio (Eds.). Association for Computational Linguistics, 4171–4186. <https://doi.org/10.18653/v1/n19-1423>
- [10] Thibault Formal, Carlos Lassance, Benjamin Piwowarski, and Stéphane Clinchant. 2021. SPLADE v2: Sparse Lexical and Expansion Model for Information Retrieval. *CoRR abs/2109.10086* (2021). arXiv:2109.10086 <https://arxiv.org/abs/2109.10086>
- [11] Thibault Formal, Benjamin Piwowarski, and Stéphane Clinchant. 2021. SPLADE: Sparse Lexical and Expansion Model for First Stage Ranking. In *SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021*, Fernando Diaz, Chirag Shah, Torsten Suel, Pablo Castells, Rosie Jones, and Tetsuya Sakai (Eds.). ACM, 2288–2292. <https://doi.org/10.1145/3404835.3463098>
- [12] Luyu Gao and Jamie Callan. 2021. Condenser: a Pre-training Architecture for Dense Retrieval. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (Eds.). Association for Computational Linguistics, 981–993. <https://doi.org/10.18653/v1/2021.emnlp-main.75>
- [13] Luyu Gao and Jamie Callan. 2021. Unsupervised Corpus Aware Language Model Pre-training for Dense Passage Retrieval. *CoRR abs/2108.05540* (2021). arXiv:2108.05540 <https://arxiv.org/abs/2108.05540>
- [14] Luyu Gao, Zhuyun Dai, and Jamie Callan. 2021. COIL: Revisit Exact Lexical Match in Information Retrieval with Contextualized Inverted List. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tür, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou (Eds.). Association for Computational Linguistics, 3030–3042. <https://doi.org/10.18653/v1/2021.naacl-main.241>
- [15] Tianyu Gao, Kingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple Contrastive Learning of Sentence Embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (Eds.). Association for Computational Linguistics, 6894–6910. <https://doi.org/10.18653/v1/2021.emnlp-main.552>
- [16] Sebastian Hofstätter, Sophia Althammer, Michael Schröder, Mete Sertkan, and Allan Hanbury. 2020. Improving Efficient Neural Ranking Models with Cross-Architecture Knowledge Distillation. *CoRR abs/2010.02666* (2020). arXiv:2010.02666 <https://arxiv.org/abs/2010.02666>
- [17] Sebastian Hofstätter, Sheng-Chieh Lin, Jheng-Hong Yang, Jimmy Lin, and Allan Hanbury. 2021. Efficiently Teaching an Effective Dense Retriever with Balanced Topic Aware Sampling. In *SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021*, Fernando Diaz, Chirag Shah, Torsten Suel, Pablo Castells, Rosie Jones, and Tetsuya Sakai (Eds.). ACM, 113–122. <https://doi.org/10.1145/3404835.3462891>
- [18] Jui-Ting Huang, Ashish Sharma, Shuying Sun, Li Xia, David Zhang, Philip Pronin, Janani Padmanabhan, Giuseppe Ottaviano, and Linjun Yang. 2020. Embedding-based Retrieval in Facebook Search. In *KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020*, Rajesh Gupta, Yan Liu, Jiliang Tang, and B. Aditya Prakash (Eds.). ACM, 2553–2561. <https://doi.org/10.1145/3394486.3403305>
- [19] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2021. Billion-Scale Similarity Search with GPUs. *IEEE Trans. Big Data* 7, 3 (2021), 535–547. <https://doi.org/10.1109/TBDATA.2019.2921572>
- [20] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick S. H. Lewis, Edell 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 2020, Online, November 16-20, 2020*, Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (Eds.). Association for Computational Linguistics, 6769–6781. <https://doi.org/10.18653/v1/2020.emnlp-main.550>
- [21] 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 2020, Virtual Event, China, July 25-30, 2020*, Jimmy Huang, Yi Chang, Xueqi Cheng, Jaap Kamps, Vanessa Murdock, Ji-Rong Wen, and Yiqun Liu (Eds.). ACM, 39–48. <https://doi.org/10.1145/3397271.3401075>
- [22] Saar Kuzi, Mingyang Zhang, Cheng Li, Michael Bendersky, and Marc Najork. 2020. Leveraging Semantic and Lexical Matching to Improve the Recall of Document Retrieval Systems: A Hybrid Approach. *CoRR abs/2010.01195* (2020). arXiv:2010.01195 <https://arxiv.org/abs/2010.01195>
- [23] Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. 2019. Latent Retrieval for Weakly Supervised Open Domain Question Answering. In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28–August 2, 2019, Volume 1: Long Papers*, Anna Korhonen, David R. Traum, and Lluís Màrquez (Eds.). Association for Computational Linguistics, 6086–6096. <https://doi.org/10.18653/v1/p19-1612>
- [24] Hang Li, Ahmed Mourad, Shengyao Zhuang, Bevan Koopman, and Guido Zuccon. 2021. Pseudo Relevance Feedback with Deep Language Models and Dense Retrievers: Successes and Pitfalls. *CoRR abs/2108.11044* (2021). arXiv:2108.11044 <https://arxiv.org/abs/2108.11044>
- [25] Hang Li, Shuai Wang, Shengyao Zhuang, Ahmed Mourad, Xueguang Ma, Jimmy Lin, and Guido Zuccon. 2022. To Interpolate or not to Interpolate: PRF, Dense and Sparse Retrievers. In *SIGIR '22: The 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, Madrid, Spain, July 11-15, 2022*, Enrique Amigó, Pablo Castells, Julio Gonzalo, Ben Carterette, J. Shane Culpepper, and Gabriella Kazai (Eds.). ACM, 2495–2500. <https://doi.org/10.1145/3477495.3531884>
- [26] Jimmy Lin and Xueguang Ma. 2021. A Few Brief Notes on DeepImpact, COIL, and a Conceptual Framework for Information Retrieval Techniques. *CoRR abs/2106.14807* (2021). arXiv:2106.14807 <https://arxiv.org/abs/2106.14807>
- [27] Sheng-Chieh Lin, Minghan Li, and Jimmy Lin. 2022. Aggretriever: A Simple Approach to Aggregate Textual Representation for Robust Dense Passage Retrieval. *CoRR abs/2208.00511* (2022). <https://doi.org/10.48550/arXiv.2208.00511>
- [28] Sheng-Chieh Lin, Jheng-Hong Yang, and Jimmy Lin. 2021. In-Batch Negatives for Knowledge Distillation with Tightly-Coupled Teachers for Dense Retrieval. In *Proceedings of the 6th Workshop on Representation Learning for NLP, Repl4NLP@ACL-IJCNLP 2021, Online, August 6, 2021*, Anna Rogers, Iacer Calixto, Ivan Vulic, Naomi Saphra, Nora Kassner, Oana-Maria Camburu, Trapit Bansal, and Vered Schwartz (Eds.). Association for Computational Linguistics, 163–173. <https://doi.org/10.18653/v1/2021.repl4nlp-1.17>
- [29] Zhenghao Lin, Yeyun Gong, Xiao Liu, Hang Zhang, Chen Lin, Anlei Dong, Jian Jiao, Jingwen Lu, Daxin Jiang, Rangan Majumder, and Nan Duan. 2022. PROD: Progressive Distillation for Dense Retrieval. *CoRR abs/2209.13335* (2022). <https://doi.org/10.48550/arXiv.2209.13335>
- [30] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *CoRR abs/1907.11692* (2019). arXiv:1907.11692 <http://arxiv.org/abs/1907.11692>
- [31] Yi Luan, Jacob Eisenstein, Kristina Toutanova, and Michael Collins. 2021. Sparse, Dense, and Attentional Representations for Text Retrieval. *Trans. Assoc. Comput. Linguistics* 9 (2021), 329–345. [https://doi.org/10.1162/tacl\\_a\\_00369](https://doi.org/10.1162/tacl_a_00369)

- [32] Ji Ma, Ivan Korotkov, Yinfei Yang, Keith B. Hall, and Ryan T. McDonald. 2021. Zero-shot Neural Passage Retrieval via Domain-targeted Synthetic Question Generation. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, EAACL 2021, Online, April 19 - 23, 2021*, Paola Merlo, Jörg Tiedemann, and Reut Tsarfaty (Eds.). Association for Computational Linguistics, 1075–1088. <https://doi.org/10.18653/v1/2021.eacl-main.92>
- [33] Antonio Mallia, Omar Khattab, Torsten Suel, and Nicola Tonello. 2021. Learning Passage Impacts for Inverted Indexes. In *SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021*, Fernando Diaz, Chirag Shah, Torsten Suel, Pablo Castells, Rosie Jones, and Tetsuya Sakai (Eds.). ACM, 1723–1727. <https://doi.org/10.1145/3404835.3463030>
- [34] 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 (CEUR Workshop Proceedings, Vol. 1773)*, Tarek Richard Besold, Antoine Bordes, Artur S. d'Avila Garcez, and Greg Wayne (Eds.). CEUR-WS.org. [http://ceur-ws.org/Vol-1773/CoCoNIPS\\_2016\\_paper9.pdf](http://ceur-ws.org/Vol-1773/CoCoNIPS_2016_paper9.pdf)
- [35] Yingqi Qu, Yuchen Ding, Jing Liu, Kai Liu, Ruiyang Ren, Wayne Xin Zhao, Daxiang Dong, Hua Wu, and Haifeng Wang. 2021. RocketQA: An Optimized Training Approach to Dense Passage Retrieval for Open-Domain Question Answering. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tür, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou (Eds.). Association for Computational Linguistics, 5835–5847. <https://doi.org/10.18653/v1/2021.naacl-main.466>
- [36] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training. (2018).
- [37] Sashank J. Reddi, Rama Kumar Pasumarthi, Aditya Krishna Menon, Ankit Singh Rawat, Felix X. Yu, Seungyeon Kim, Andreas Veit, and Sanjiv Kumar. 2021. RankDistil: Knowledge Distillation for Ranking. In *The 24th International Conference on Artificial Intelligence and Statistics, AISTATS 2021, April 13-15, 2021, Virtual Event (Proceedings of Machine Learning Research, Vol. 130)*, Arindam Banerjee and Kenji Fukumizu (Eds.). PMLR, 2368–2376. <http://proceedings.mlr.press/v130/reddi21a.html>
- [38] Ruiyang Ren, Shangwen Lv, Yingqi Qu, Jing Liu, Wayne Xin Zhao, Qiaoqiao She, Hua Wu, Haifeng Wang, and Ji-Rong Wen. 2021. PAIR: Leveraging Passage-Centric Similarity Relation for Improving Dense Passage Retrieval. In *Findings of the Association for Computational Linguistics: ACL/IJCNLP 2021, Online Event, August 1-6, 2021 (Findings of ACL, Vol. ACL/IJCNLP 2021)*, Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (Eds.). Association for Computational Linguistics, 2173–2183. <https://doi.org/10.18653/v1/2021.findings-acl.191>
- [39] Ruiyang Ren, Yingqi Qu, Jing Liu, Wayne Xin Zhao, Qiaoqiao She, Hua Wu, Haifeng Wang, and Ji-Rong Wen. 2021. RocketQAv2: A Joint Training Method for Dense Passage Retrieval and Passage Re-ranking. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (Eds.). Association for Computational Linguistics, 2825–2835. <https://doi.org/10.18653/v1/2021.emnlp-main.224>
- [40] Stephen E. Robertson and Hugo Zaragoza. 2009. The Probabilistic Relevance Framework: BM25 and Beyond. *Found. Trends Inf. Retr.* 3, 4 (2009), 333–389. <https://doi.org/10.1561/15000000019>
- [41] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *CoRR abs/1910.01108* (2019). [arXiv:1910.01108](http://arxiv.org/abs/1910.01108) <http://arxiv.org/abs/1910.01108>
- [42] Keshav Santhanam, Omar Khattab, Jon Saad-Falcon, Christopher Potts, and Matei Zaharia. 2021. ColBERTv2: Effective and Efficient Retrieval via Lightweight Late Interaction. *CoRR abs/2112.01488* (2021). [arXiv:2112.01488](https://arxiv.org/abs/2112.01488) <https://arxiv.org/abs/2112.01488>
- [43] Min Joon Seo, Jinhyuk Lee, Tom Kwiatkowski, Ankur P. Parikh, Ali Farhadi, and Hannaneh Hajishirzi. 2019. Real-Time Open-Domain Question Answering with Dense-Sparse Phrase Index. In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, Anna Korhonen, David R. Traum, and Lluís Màrquez (Eds.). Association for Computational Linguistics, 4430–4441. <https://doi.org/10.18653/v1/p19-1436>
- [44] Tao Shen, Xiubo Geng, Chongyang Tao, Can Xu, Kai Zhang, and Daxin Jiang. 2022. Unifier: A Unified Retriever for Large-Scale Retrieval. *CoRR abs/2205.11194* (2022). <https://doi.org/10.48550/arXiv.2205.11194> [arXiv:2205.11194](https://arxiv.org/abs/2205.11194)
- [45] Yu Sun, Shuohuan Wang, Yu-Kun Li, Shikun Feng, Hao Tian, Hua Wu, and Haifeng Wang. 2020. ERNIE 2.0: A Continual Pre-Training Framework for Language Understanding. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*. AAAI Press, 8968–8975. <https://ojs.aaai.org/index.php/AAAI/article/view/6428>
- [46] Shuai Wang, Shengyao Zhuang, and Guido Zuccon. 2021. BERT-based Dense Retrievers Require Interpolation with BM25 for Effective Passage Retrieval. In *ICTIR '21: The 2021 ACM SIGIR International Conference on the Theory of Information Retrieval, Virtual Event, Canada, July 11, 2021*, Faegheh Hasibi, Yi Fang, and Akiko Aizawa (Eds.). ACM, 317–324. <https://doi.org/10.1145/3471158.3472233>
- [47] Xing Wu, Guangyuan Ma, Meng Lin, Zijia Lin, Zhongyuan Wang, and Songlin Hu. 2022. Contextual mask auto-encoder for dense passage retrieval. *arXiv preprint arXiv:2208.07670* (2022).
- [48] Shitao Xiao, Zheng Liu, Weihao Han, Jianjin Zhang, Defu Lian, Yeyun Gong, Qi Chen, Fan Yang, Hao Sun, Yingxia Shao, Denvy Deng, Qi Zhang, and Xing Xie. 2022. Distill-VQ: Learning Retrieval Oriented Vector Quantization By Distilling Knowledge from Dense Embeddings. *CoRR abs/2204.00185* (2022). <https://doi.org/10.48550/arXiv.2204.00185> [arXiv:2204.00185](https://arxiv.org/abs/2204.00185)
- [49] Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul N. Bennett, Junaid Ahmed, and Arnold Overwijk. 2021. Approximate Nearest Neighbor Negative Contrastive Learning for Dense Text Retrieval. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net. <https://openreview.net/forum?id=zeFrgyZln>
- [50] Canwen Xu, Daya Guo, Nan Duan, and Julian J. McAuley. 2022. LaProDoR: Unsupervised Pretrained Dense Retriever for Zero-Shot Text Retrieval. In *Findings of the Association for Computational Linguistics: ACL 2022, Dublin, Ireland, May 22-27, 2022*, Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (Eds.). Association for Computational Linguistics, 3557–3569. <https://doi.org/10.18653/v1/2022.findings-acl.281>
- [51] Hansi Zeng, Hamed Zamani, and Vishwa Vinay. 2022. Curriculum Learning for Dense Retrieval Distillation. In *SIGIR '22: The 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, Madrid, Spain, July 11 - 15, 2022*, Enrique Amigó, Pablo Castells, Julio Gonzalo, Ben Carterette, J. Shane Culpepper, and Gabriella Kazai (Eds.). ACM, 1979–1983. <https://doi.org/10.1145/3477495.3531791>
- [52] Jingtao Zhan, Jiabin Mao, Yiqun Liu, Jiafeng Guo, Min Zhang, and Shaoping Ma. 2021. Optimizing Dense Retrieval Model Training with Hard Negatives. In *SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021*, Fernando Diaz, Chirag Shah, Torsten Suel, Pablo Castells, Rosie Jones, and Tetsuya Sakai (Eds.). ACM, 1503–1512. <https://doi.org/10.1145/3404835.3462880>
- [53] Hang Zhang, Yeyun Gong, Yelong Shen, Jiancheng Lv, Nan Duan, and Weizhu Chen. 2022. Adversarial Retriever-Ranker for Dense Text Retrieval. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=MR7XubKUFb>
- [54] Shuai Zhang, Lina Yao, Aixin Sun, and Yi Tay. 2019. Deep Learning Based Recommender System: A Survey and New Perspectives. *ACM Comput. Surv.* 52, 1 (2019), 5:1–5:38. <https://doi.org/10.1145/3285029>
- [55] Yucheng Zhou, Tao Shen, Xiubo Geng, Chongyang Tao, Can Xu, Guodong Long, Bingxing Jiao, and Daxin Jiang. 2022. Towards Robust Ranker for Text Retrieval. <https://doi.org/10.48550/ARXIV.2206.08063>
- [56] Lixin Zou, Shengqiang Zhang, Hengyi Cai, Dehong Ma, Suqi Cheng, Shuaiqiang Wang, Daiting Shi, Zhicong Cheng, and Dawei Yin. 2021. Pre-trained Language Model based Ranking in Baidu Search. In *KDD '21: The 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, Singapore, August 14-18, 2021*, Feida Zhu, Beng Chin Ooi, and Chunyan Miao (Eds.). ACM, 4014–4022. <https://doi.org/10.1145/3447548.3467147>

## A BASELINES

We compare with previous state-of-the-art baselines including traditional term-based techniques like BM25 [40], and dense as well as lexicon-aware retrievers. For lexicon-aware retrievers, DeepCT [7] was trained to predict term weights. COIL [14] used contextualized representation for exact term matching and UniCOIL [26] compressed vectors in COIL into scalars. DistilSPLADE-max and SPLADE-max [10] were both trained with Eq. 7 and the latter one was further enhanced by a cross-encoder. The UniCOIL  $\Lambda$  was a lexicon-aware model trained with UniCOIL’s top-ranked passages and negatives [4]. For dense retrievers, ANCE [49] selected hard training negatives from the entire collection. ADORE [52] used self-mining static negatives and then dynamic negatives. TAS-B [17] proposed balanced topic-aware negative sampling strategies for effective teaching. CL-DRD [51] taught the retriever in a curriculum learning fashion, starting from coarse-grained pair examples and progressing to fine-grained ones. ColBERTv1 [21] and ColBERTv2 [42] utilized late-interaction and the latter one further incorporates ranker distillation. TCT-ColBERT [28] utilized ColBERTv1 as the tightly-coupled teacher to enable in-batch distillation. The coCondenser [13] augmented MLM loss with contrastive learning and based a model architecture [12] with a decoupled sentence and token interaction. PAIR [38] introduced passage-centric loss to assist the contrastive loss and combine cross-encoder teaching. RocketQAv2 [39] utilized K-L divergence to align the list-wise distributions between retriever and ranker and proposed hybrid data augmentation. AR2-G [53] used an adversarial framework to train the retriever and ranker simultaneously. Notably, AR2 used a different Recall@N evaluation from the official TREC Recall@N.<sup>6</sup> Therefore, we don’t report their Recall@N performances in Tab. 1.

---

<sup>6</sup><https://github.com/microsoft/AR2/tree/main/AR2>