
EmbedDistill: A Geometric Knowledge Distillation for Information Retrieval

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Large neural models (such as Transformers) achieve state-of-the-art performance
2 for information retrieval (IR). In this paper, we aim to improve distillation methods
3 that pave the way for the resource-efficient deployment of such models in practice.
4 Inspired by our theoretical analysis of the teacher-student generalization gap for
5 IR models, we propose a novel distillation approach that leverages the relative
6 geometry among queries and documents learned by the large teacher model. Unlike
7 existing teacher score-based distillation methods, our proposed approach employs
8 embedding matching tasks to provide a stronger signal to align the representations
9 of the teacher and student models. In addition, it utilizes query generation to
10 explore the data manifold to reduce the discrepancies between the student and the
11 teacher where training data is sparse. Furthermore, our analysis also motivates
12 novel asymmetric architectures for student models which realizes better embedding
13 alignment without increasing online inference cost. On standard benchmarks like
14 MSMARCO, we show that our approach successfully distills from both dual-
15 encoder (DE) and cross-encoder (CE) teacher models to 1/10th size asymmetric
16 students that can retain 95-97% of the teacher performance.

17 1 Introduction

18 Neural models for information retrieval (IR) are increasingly used to model the true ranking function
19 in various applications, including web search [38], recommendation [65], and question-answering
20 (QA) [6]. Notably, the recent success of Transformers [59]-based pre-trained language models [11,
21 30, 49] on a wide range of natural language understanding tasks has also prompted their utilization in
22 IR to capture query-document relevance [see, e.g., 10, 34, 43, 26, 20].

23 A typical IR system comprises two stages: (1) A *retriever* first selects a small subset of potentially
24 relevant candidate documents (out of a large collection) for a given query; and (2) A *re-ranker* then
25 identifies a precise ranking among the candidates provided by the retriever. *Dual-encoder* (DE)
26 models are the de-facto architecture for retrievers [26, 20]. Such models independently embed queries
27 and documents into a common space, and capture their relevance by simple operations on these
28 embeddings such as the inner product. This enables offline creation of a document index and supports
29 fast retrieval during inference via efficient maximum inner product search implementations [12, 19],
30 with *online* query embedding generation primarily dictating the inference latency. *Cross-encoder* (CE)
31 models, on the other hand, are preferred as re-rankers, owing to their excellent performance [43, 9, 62].
32 A CE model jointly encodes a query-document pair while enabling early interaction among query
33 and document features. Employing a CE model for retrieval is often infeasible, as it would require
34 processing a given query with *every* document in the collection at inference time. In fact, even in
35 the re-ranking stage, the inference cost of CE models is high enough [22] to warrant exploration of
36 efficient alternatives [14, 22, 37]. Across both architectures, scaling to larger models brings improved
37 performance at increased computational cost [41, 39].

38 *Knowledge distillation* [5, 13] provides a general strategy to address the prohibitive inference cost
 39 associated with high-quality large neural models. In the IR literature, most existing distillation
 40 methods only rely on the teacher’s query-document relevance scores [see, e.g., 31, 14, 8, 51, 56] or
 41 their proxies [16]. However, given that neural IR models are inherently embedding-based, it is natural
 42 to ask: *Is it useful to go beyond matching of the teacher and student models’ scores, and directly aim*
 43 *to align their embedding spaces?*

44 With this in mind, we propose a novel distillation method for IR models that utilizes an *embedding*
 45 *matching* task to train student models. The proposed method is inspired by our rigorous treatment
 46 of the generalization gap between the teacher and student models in IR settings. Our theoretical
 47 analysis of the *teacher-student generalization gap* further suggests novel design choices involving
 48 *asymmetric configurations* for student DE models, intending to further reduce the gap by better
 49 aligning teacher and student embedding spaces. Notably, our proposed distillation method supports
 50 *cross-architecture distillation* and improves upon existing (score-based) distillation methods for both
 51 retriever and re-ranker models. When distilling a large teacher DE model into a smaller student DE
 52 model, for a given query (document), one can minimize the distance between the query (document)
 53 embeddings of the teacher and student (after compatible projection layers to account for dimension
 54 mismatch, if any). In contrast, a teacher CE model doesn’t directly provide document and query
 55 embeddings, and so to effectively employ embedding matching-based distillation requires modifying
 56 the scoring layer with *dual-pooling* [61] and adding various regularizers. Both of these changes
 57 improve geometry of teacher embeddings and facilitate effective knowledge transfer to the student
 58 DE model via embedding matching-based distillation.

59 Our key contributions toward improving IR models via distillation are:

- 60 • We provide the first rigorous analysis of the teacher-student generalization gap for IR settings
 61 which captures the role of alignment of embedding spaces of the teacher and student towards
 62 reducing the gap (Sec. 3).
- 63 • Inspired by our analysis, we propose a novel distillation approach for neural IR models, namely
 64 EmbedDistill, that goes beyond score matching and aligns the embedding spaces of the teacher and
 65 student models (Sec. 4). We also show that EmbedDistill can leverage synthetic data to improve a
 66 student by further aligning the embedding spaces of the teacher and student (Sec. 4.3).
- 67 • Our analysis motivates novel distillation setups. Specifically, we consider a student DE model with
 68 an *asymmetric* configuration, consisting of a small query encoder and a *frozen* document encoder
 69 inherited from the teacher. This significantly reduces inference latency of query embedding
 70 generation, while leveraging the teachers’ high-quality document index (Sec. 4.1).
- 71 • We provide a *comprehensive* empirical evaluation of EmbedDistill (Sec. 5) on two standard IR
 72 benchmarks – Natural Questions [23] and MSMARCO [40]. We also evaluate EmbedDistill on
 73 BEIR benchmark [57] which is used to measure the *zero-shot* performance of an IR model.

74 Note that prior works have utilized embedding alignment during distillation for *non-IR* setting [see,
 75 e.g., 52, 55, 18, 1, 64, 7]. However, to the best of our knowledge, our work is the first to study
 76 embedding matching-based distillation method for IR settings which requires addressing multiple
 77 IR-specific challenges such as cross-architecture distillation, partial representation alignment, and en-
 78 abling novel asymmetric student configurations. Furthermore, unlike these prior works, our proposed
 79 method is theoretically justified to reduce the teacher-student performance gap.

80 2 Background

81 Let \mathcal{Q} and \mathcal{D} denote the query and document spaces, respectively. An IR model is equivalent to
 82 a scorer $s : \mathcal{Q} \times \mathcal{D} \rightarrow \mathbb{R}$, i.e., it assigns a (relevance) score $s(q, d)$ for a query-document pair
 83 $(q, d) \in \mathcal{Q} \times \mathcal{D}$. Ideally, we want to learn a scorer such that $s(q, d) > s(q, d')$ iff the document d is
 84 more relevant to the query q than document d' . We assume access to n labeled training examples
 85 $\mathcal{S}_n = \{(q_i, \mathbf{d}_i, \mathbf{y}_i)\}_{i \in [n]}$. Here, $\mathbf{d}_i = (d_{i,1}, \dots, d_{i,L}) \in \mathcal{D}^L$, $\forall i \in [n]$, denotes a list of L documents
 86 and $\mathbf{y}_i = (y_{i,1}, \dots, y_{i,L}) \in \{0, 1\}^L$ denotes the corresponding labels such that $y_{i,j} = 1$ iff the
 87 document $d_{i,j}$ is relevant to the query q_i . Given \mathcal{S}_n , we learn an IR model by minimizing

$$R(s; \mathcal{S}_n) := \frac{1}{n} \sum_{i \in [n]} \ell(s_{q_i, \mathbf{d}_i}, \mathbf{y}_i), \quad (1)$$

88 where $s_{q_i, \mathbf{d}_i} := (s(q_i, d_{1,i}), \dots, s(q_i, d_{L,i}))$ and $\ell(s_{q_i, \mathbf{d}_i}, \mathbf{y}_i)$ denotes the loss s incurs on $(q_i, \mathbf{d}_i, \mathbf{y}_i)$.
 89 Due to space constraint, we defer concrete choices for the loss function ℓ to Appendix A.

90 While this learning framework is general enough to work with any IR models, next, we formally
 91 introduce two families of Transformer-based IR models that are prevalent in the recent literature.

92 2.1 Transformer-based IR models: Cross-encoders and Dual-encoders

93 Let query $q = (q^1, \dots, q^{m_1})$ and document $d = (d^1, \dots, d^{m_2})$ consist of m_1 and m_2 tokens, respec-
 94 tively. We now discuss how Transformers-based CE and DE models process the (q, d) pair.

95 **Cross-encoder model.** Let $p = [q; d]$ be the sequence obtained by concatenating q and d . Further,
 96 let \tilde{p} be the sequence obtained by adding special tokens such [CLS] and [SEP] to p . Given an
 97 encoder-only Transformer model Enc , the relevance score for the (q, d) pair is

$$s(q, d) = \langle w, \text{pool}(\text{Enc}(\tilde{p})) \rangle = \langle w, \text{emb}_{q,d} \rangle, \quad (2)$$

98 where w is a d -dimensional classification vector, and $\text{pool}(\cdot)$ denotes a pooling operation that
 99 transforms the contextualized token embeddings $\text{Enc}(\tilde{p})$ to a joint embedding vector $\text{emb}_{q,d}$. [CLS]-
 100 pooling is a common operation that simply outputs the embedding of the [CLS] token as $\text{emb}_{q,d}$.

101 **Dual-encoder model.** Let \tilde{q} and \tilde{d} be the sequences obtained by adding appropriate special tokens
 102 to q and d , respectively. A DE model comprises two (encoder-only) Transformers Enc_Q and Enc_D ,
 103 which we call query and document encoders, respectively.¹ Let $\text{emb}_q = \text{pool}(\text{Enc}_Q(\tilde{q}))$ and emb_d
 104 $= \text{pool}(\text{Enc}_D(\tilde{d}))$ denote the query and document embeddings, respectively. Now, one can define
 105 $s(q, d) = \langle \text{emb}_q, \text{emb}_d \rangle$ to be the relevance score assigned to the (q, d) pair by the DE model.

106 2.2 Score-based distillation for IR models

107 Most distillation schemes for IR [e.g., 31, 14, 8] rely on teacher relevance scores. Given a training set
 108 \mathcal{S}_n and a teacher with scorer s^t , one learns a student with scorer s^s by minimizing

$$R(s^s, s^t; \mathcal{S}_n) = \frac{1}{n} \sum_{i \in [n]} \ell_d(s_{q_i}^s, s_{q_i}^t), \quad (3)$$

109 where ℓ_d captures the discrepancy between s^s and s^t . See Appendix A for common choices for ℓ_d .

110 3 Teacher-student generalization gap: Inspiration for embedding alignment

111 Our main objective is to devise novel distillation methods to realize high-performing student DE
 112 models. As a first step in this direction, we rigorously study the teacher-student generalization
 113 gap as realized by standard (score-based) distillation in IR settings. Informed by our analysis, we
 114 subsequently identify novel ways to improve the student model’s performance. In particular, our
 115 analysis suggests two natural directions to reduce the teacher-student generalization gap: 1) enforcing
 116 tighter alignment between embedding spaces of teacher and student models; and 2) exploring novel
 117 asymmetric configuration for student DE model.

118 Let $R(s) = \mathbb{E} [\ell(s_{q,d}, \mathbf{y})]$ be the population version of the empirical risk in Eq. 1, which measures
 119 the test time performance of the IR model defined by the scorer s . Thus, $R(s^s) - R(s^t)$ denotes the
 120 *teacher-student generalization gap*. In the following result, we bound this quantity (see Appendix C.1
 121 for a formal statement and proof). We focus on distilling a teacher DE model to a student DE model
 122 and $L = 1$ (cf. Sec. 2) as it leads to easier exposition without changing the main takeaways. Our
 123 analysis can be extended to $L > 1$ or CE to DE distillation with more complex notation.

124 **Theorem 3.1** (Teacher-student generalization gap (informal)). *Let \mathcal{F} and \mathcal{G} denote the function*
 125 *classes for the query and document encoders for the student model, respectively. Suppose that the*
 126 *score-based distillation loss ℓ_d in Eq. 3 is based on binary cross entropy loss (Eq. 12 in Appendix A).*
 127 *Let one-hot (label-dependent) loss ℓ in Eq. 1 be the binary cross entropy loss (Eq. 10 in Appendix A).*
 128 *Further, assume that all encoders have the same output dimension and embeddings have their ℓ_2 -norm*
 129 *bounded by K . Then, we have*

$$\begin{aligned} R(s^s) - R(s^t) &\leq \mathcal{E}_n(\mathcal{F}, \mathcal{G}) + 2K R_{\text{Emb},Q}(t, s; \mathcal{S}_n) + 2K R_{\text{Emb},D}(t, s; \mathcal{S}_n) \\ &\quad + \Delta(s^t; \mathcal{S}_n) + K^2 (\mathbb{E} [|\sigma(s_{q,d}^t) - y|] + \frac{1}{n} \sum_{i \in [n]} |\sigma(s_{q_i, d_i}^t) - y_i|), \end{aligned} \quad (4)$$

¹It is common to employ dual-encoder models where query and document encoders are shared.

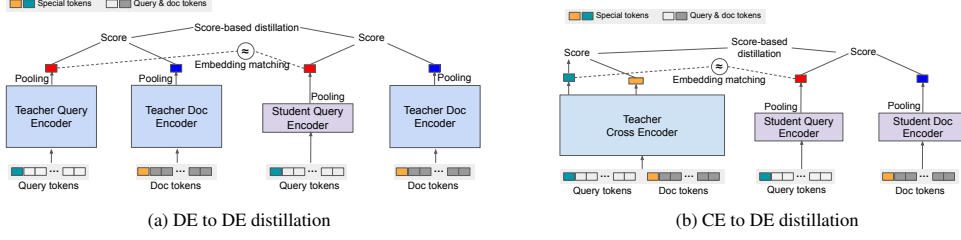


Figure 1: Proposed distillation method with query embedding matching. **Left:** The setting where student employs an asymmetric DE configuration with a small query encoder and a large (non-trainable) document encoder inherited from the teacher DE model. The smaller query encoder ensures small latency for encoding query during inference, and large document encoder leads to a good quality document index. **Right:** Similarly the setting of CE to DE distillation using EmbedDistill, with teacher CE model employing dual pooling.

130 where $\mathcal{E}_n(\mathcal{F}, \mathcal{G}) := \sup_{s^s \in \mathcal{F} \times \mathcal{G}} |R(s^s, s^t; \mathcal{S}_n) - \mathbb{E} \ell_d(s_{q,d}^s, s_{q,d}^t)|$; σ denotes the sigmoid function;
 131 and $\Delta(s^t; \mathcal{S}_n)$ denotes the deviation between the empirical risk (on \mathcal{S}_n) and population risk of the
 132 teacher s^t . Here, $R_{\text{Emb},Q}(t, s; \mathcal{S}_n)$ and $R_{\text{Emb},D}(t, s; \mathcal{S}_n)$ measure misalignment between teacher and
 133 student embeddings by focusing on queries and documents, respectively (cf. Eq. 7 & 8 in Sec. 4.1).

134 The last three quantities in the bound in Thm. 3.1, namely $\Delta(s^t; \mathcal{S}_n)$, $\mathbb{E}[|\sigma(s_{q,d}^t) - y|]$, and
 135 $\frac{1}{n} \sum_{i \in [n]} |\sigma(s_{q_i, d_i}^t) - y_i|$, are independent of the underlying student model. These terms solely
 136 depend on the quality of the underlying teacher model s^t . That said, the teacher-student gap can be
 137 made small by reducing the following three terms: 1) uniform deviation of the student’s empirical
 138 distillation risk from its population version $\mathcal{E}_n(\mathcal{F}, \mathcal{G})$; 2) misalignment between teacher student query
 139 embeddings $R_{\text{Emb},Q}(t, s; \mathcal{S}_n)$; and 3) misalignment between teacher student document embeddings
 140 $R_{\text{Emb},D}(t, s; \mathcal{S}_n)$.

141 The last two terms motivate us to propose an *embedding matching*-based distillation that explicitly
 142 aims to minimize these terms during student training. Even more interestingly, these terms also
 143 inspire an *asymmetric DE configuration* for the student which strikes a balance between the goals of
 144 reducing the misalignment between the embeddings of teacher and student (by inheriting teacher’s
 145 document encoder) and ensuring serving efficiency (small inference latency) by employing a small
 146 query encoder. Before discussing these proposals in detail in Sec. 4 and Fig. 1, we explore the first
 147 term $\mathcal{E}_n(\mathcal{F}, \mathcal{G})$ and highlight how our proposals also have implications for reducing this term. Towards
 148 this, the following result bounds $\mathcal{E}_n(\mathcal{F}, \mathcal{G})$. Due to space constraints, we present an informal statement
 149 of the result (see Appendix C.2 for a more precise statement and proof).

150 **Proposition 3.2.** Let ℓ_d be a distillation loss which is L_{ℓ_d} -Lipschitz in its first argument. Let \mathcal{F} and \mathcal{G}
 151 denote the function classes for the query and document encoders, respectively. Further assume that,
 152 for each query and document encoder in our function class, the query and document embeddings
 153 have their ℓ_2 -norm bounded by K . Then,

$$\mathcal{E}_n(\mathcal{F}, \mathcal{G}) \leq \mathbb{E}_{\mathcal{S}_n} \frac{48KL_{\ell_d}}{\sqrt{n}} \int_0^\infty \sqrt{\log(N(u, \mathcal{F})N(u, \mathcal{G}))} du. \quad (5)$$

154 Furthermore, with a fixed document encoder, i.e., $\mathcal{G} = \{g^*\}$,

$$\mathcal{E}_n(\mathcal{F}, \{g^*\}) \leq \mathbb{E}_{\mathcal{S}_n} \frac{48KL_{\ell_d}}{\sqrt{n}} \int_0^\infty \sqrt{\log N(u, \mathcal{F})} du. \quad (6)$$

155 Here, $N(u, \cdot)$ is the u -covering number of a function class.

156 Note that Eq. 5 and Eq. 6 correspond to uniform deviation when we train *without* and *with* a frozen
 157 document encoder, respectively. It is clear that the bound in Eq. 6 is less than or equal to that in
 158 Eq. 5 (because $N(u, \mathcal{G}) \geq 1$ for any u), which alludes to desirable impact of employing a frozen
 159 document encoder as one of our proposal seeks to do via *inheriting teacher’s document encoder* (for
 160 instance in an asymmetric DE configuration). Furthermore, our proposal of employing an embedding-
 161 matching task will regularize the function class of query encoders; effectively reducing it to \mathcal{F}' with
 162 $|\mathcal{F}'| \leq |\mathcal{F}|$. The same holds true for document encoder function class when document encoder is
 163 trainable (as in Eq. 5), leading to an effective function class \mathcal{G}' with $|\mathcal{G}'| \leq |\mathcal{G}|$. Since we would have
 164 $N(u, \mathcal{F}') \leq N(u, \mathcal{F})$ and $N(u, \mathcal{G}') \leq N(u, \mathcal{G})$, this suggests desirable implications of embedding
 165 matching for reducing the uniform deviation bound.

166 **4 Embedding-matching based distillation**

167 Informed by our analysis of teacher-student generalization gap in Sec. 3, we propose EmbedDistill – a
 168 novel distillation method that explicitly focuses on aligning the embedding spaces of the teacher and
 169 student. Our proposal goes beyond existing distillation methods in the IR literature that only use the
 170 teacher scores. Next, we introduce EmbedDistill for two prevalent settings: (1) distilling a large DE
 171 model to a smaller DE model;² and (2) distilling a CE model to a DE model.

172 **4.1 DE to DE distillation**

173 Given a (q, d) pair, let emb_q^t and emb_d^t be the query and document embeddings produced by the
 174 query encoder Enc_Q^t and document encoder Enc_D^t of the teacher DE model, respectively. Similarly,
 175 let emb_q^s and emb_d^s denote the query and document embeddings produced by a student DE model
 176 with $(\text{Enc}_Q^s, \text{Enc}_D^s)$ as its query and document encoders. Now, EmbedDistill optimizes the following
 177 embedding alignment losses in addition to the score-matching loss from Sec. 2.2 to align query and
 178 document embeddings of the teacher and student:

$$R_{\text{Emb},Q}(t, s; \mathcal{S}_n) = \frac{1}{n} \sum_{q \in \mathcal{S}_n} \|\text{emb}_q^t - \text{proj}(\text{emb}_q^s)\|; \quad (7)$$

$$R_{\text{Emb},D}(t, s; \mathcal{S}_n) = \frac{1}{n} \sum_{d \in \mathcal{S}_n} \|\text{emb}_d^t - \text{proj}(\text{emb}_d^s)\|. \quad (8)$$

179 **Asymmetric DE.** We also propose a novel student DE configuration where the student employs the
 180 teacher’s document encoder (i.e., $\text{Enc}_D^s = \text{Enc}_D^t$) and only train its query encoder, which is much
 181 smaller compared to the teacher’s query encoder. For such a setting, it is natural to only employ the
 182 embedding matching loss in Eq. 7 as the document embeddings are aligned by design (cf. Fig. 1a).

183 Note that this asymmetric student DE does not incur an increase in latency despite the use of a
 184 large teacher document encoder. This is because the large document encoder is only needed to
 185 create a good quality document index offline, and only the query encoder is evaluated at inference
 186 time. Also, the similarity search cost is not increased as the projection layer ensures the same small
 187 embedding dimension as in the symmetric DE student. Thus, for DE to DE distillation, we prescribe
 188 the asymmetric DE configuration universally. Our theoretical analysis (cf. Sec. 3) and experimental
 189 results (cf. Sec. 5) suggest that the ability to inherit the document tower from the teacher DE model
 190 can drastically improve the final performance, especially when combined with query embedding
 191 matching task (cf. Eq. 7).

192 **4.2 CE to DE distillation**

193 Given that CE models jointly encode query-document pairs, individual query and document embed-
 194 dings are not readily available to implement embedding matching losses as per Eq. 7 and 8. This
 195 makes it challenging to employ EmbedDistill for CE to DE distillation.

196 As a naïve solution, for a (q, d) pair, one can simply match a joint transformation of the student’s query
 197 embedding emb_q^s and document embedding emb_d^s to the teacher’s joint embedding $\text{emb}_{q,d}^t$, produced
 198 by (single) teacher encoder Enc^t . However, we observed that including such an embedding matching
 199 task often leads to severe over-fitting, and results in a poor student. Since $s^t(q, d) = \langle w, \text{emb}_{q,d}^t \rangle$,
 200 during CE model training, the joint embeddings $\text{emb}_{q,d}^t$ for relevant and irrelevant (q, d) pairs are
 201 encouraged to be aligned with w and $-w$, respectively. This produces degenerate embeddings that
 202 do not capture semantic query-to-document relationships. We notice that even the final query and
 203 document token embeddings lose such semantic structure (cf. Appendix G.2). Thus, a teacher CE
 204 model with $s^t(q, d) = \langle w, \text{emb}_{q,d}^t \rangle$ does not add value for distillation beyond score-matching; in
 205 fact, it *hurts* to include naïve embedding matching. Next, we propose a modified CE model training
 206 strategy that facilitates EmbedDistill.

207 **CE models with dual pooling.** A *dual pooling* scheme is employed in the scoring layer to produce
 208 two embeddings $\text{emb}_{q \leftarrow (q,d)}^t$ and $\text{emb}_{d \leftarrow (q,d)}^t$ from a CE model that serve as the *proxy* query and
 209 document embeddings, respectively. Accordingly, we define the relevance score as $s^t(q, d) =$
 210 $\langle \text{emb}_{q \leftarrow (q,d)}^t, \text{emb}_{d \leftarrow (q,d)}^t \rangle$. We explore two variants of dual pooling: (1) special token-based pooling
 211 that pools from [CLS] and [SEP]; and (2) segment-based weighted mean pooling that separately

²CE to CE distillation is a special case of this with classification vector w (cf. Eq. 2) as trivial second encoder.

Table 1: Full recall performance of various student DE models on NQ dev set, including symmetric DE student model (67.5M or 11.3M transformer for both encoders), and asymmetric DE student model (67.5M or 11.3M transformer as query encoder and document embeddings inherited from the teacher). All distilled students used the same teacher (110.1M parameter BERT-base models as both encoders), with the full Recall@5 = 72.3, Recall@20 = 86.1, and Recall@100 = 93.6.

Method	6-Layer (67.5M)			4-Layer (11.3M)		
	R@5	R@20	R@100	R@5	R@20	R@100
Train student directly	36.2	59.7	80.0	24.8	44.7	67.5
+ Distill from teacher	65.3	81.6	91.2	44.3	64.9	81.0
+ Inherit doc embeddings	69.9	83.9	92.3	56.3	70.9	82.5
+ Query embedding matching	72.7	86.5	93.9	61.2	75.2	85.1
+ Query generation	73.4	86.3	93.8	64.3	77.8	87.9
Train student using only embedding matching and inherit doc embeddings	71.4	84.9	92.6	64.6	50.2	76.8
+ Query generation	71.8	85.0	93.0	54.2	68.9	80.8

Table 2: Performance of EmbedDistill for DE to DE distillation on NQ test set. While prior works listed in the table rely on techniques such as negative mining and multi-stage training, we explore the orthogonal direction of embedding-matching that improves *single-stage* distillation, which can be combined with them.

Method	#Layers	R@20	R@100
DPR [20]	12	78.4	85.4
DPR + PAQ [47]	12	84.0	89.2
DPR + PAQ [47]	24	84.7	89.2
ACNE [60]	12	81.9	87.5
RocketQA [48]	12	82.7	88.5
MSS-DPR [53]	12	84.0	89.2
MSS-DPR [53]	24	84.8	89.8
Our teacher [63]	12 (220.2M)	85.4	90.0
EmbedDistill	6 (67.5M)	85.1	89.8
EmbedDistill	4 (11.3M)	81.2	87.4

212 performs weighted averaging on the query and document segments of the final token embeddings.
 213 See Appendix B for details.

214 In addition to dual pooling, we also utilize a reconstruction loss during the CE training, which
 215 measures the likelihood of predicting each token of the original input from the final token embed-
 216 dings. This loss encourages reconstruction of query and document tokens based on the final token
 217 embeddings and prevents the degeneration of the token embeddings during training. Given proxy
 218 embeddings from the teacher CE, we can perform EmbedDistill with the embedding matching loss
 219 defined in Eq. 7 and Eq. 8 (cf. Fig. 1b).

220 4.3 Task-specific online data generation

221 Data augmentation as a general technique has been previously considered in the IR literature [see, e.g.,
 222 45, 47, 17], especially in data-limited, out-of-domain, or zero-shot settings. As EmbedDistill aims
 223 to align the embeddings spaces of the teacher and student, the ability to generate similar queries or
 224 documents can naturally help enforce such an alignment globally on the task-specific manifold. Given
 225 a set of unlabeled task-specific query and document pairs \mathcal{U}_m , we can further add the embedding
 226 matching losses $R_{\text{Emb,Q}}(t, s; \mathcal{U}_m)$ or $R_{\text{Emb,D}}(t, s; \mathcal{U}_m)$ to our training objective. Interestingly, for
 227 DE to DE distillation setting, our approach can even benefit from a large collection of task-specific
 228 queries \mathcal{Q}' or documents \mathcal{D}' . Here, we can independently employ embedding matching losses
 229 $R_{\text{Emb,Q}}(t, s; \mathcal{Q}')$ or $R_{\text{Emb,D}}(t, s; \mathcal{D}')$ that focus on queries and documents, respectively. Please refer
 230 to Appendix E describing how the task-specific data were generated.

231 5 Experiments

232 We now conduct a comprehensive evaluation of the proposed distillation approach. Specifically, we
 233 highlight the utility of the approach for both DE to DE and CE to DE distillation. We also showcase
 234 the benefits of combining our distillation approach with query generation methods.

235 5.1 Setup

236 **Benchmarks and evaluation metrics.** We consider two popular IR benchmarks — Natural Questions
 237 (NQ) [24] and MSMARCO [40], which focus on finding the most relevant passage/document given
 238 a question and a search query, respectively. NQ provides both standard test and dev sets, whereas
 239 MSMARCO provides only the dev set that are widely used for common benchmarks. In what
 240 follows, we use the terms query (document) and question (passages) interchangeably. For NQ, we
 241 use the standard full recall (*strict*) as well as the *relaxed* recall metric [20] to evaluate the retrieval
 242 performance. For MSMARCO, we focus on the standard metrics *Mean Reciprocal Rank* (MRR)@10,
 243 and *normalized Discounted Cumulative Gain* (nDCG)@10 to evaluate both re-ranking and retrieval
 244 performance. For the re-ranking, we restrict to re-ranking only the top 1000 candidate document
 245 provided as part of the dataset to be fair, while some works use stronger methods to find better
 246 top 1000 candidates for re-ranking (resulting in higher evaluation numbers) See Appendix D for a
 247 detailed discussion on these evaluation metrics. Finally, we also evaluate EmbedDistill on the BEIR
 248 benchmark [57] in terms of nDCG@10 and recall@100 metrics.

249 **Model architectures.** We follow the standard Transformers-based IR model architectures similar
 250 to Karpukhin et al. [20], Qu et al. [48], Oğuz et al. [47]. We utilized various sizes of DE models based
 251 on BERT-base [11] (12-layer, 768 dim, 110M parameters), DistilBERT [55] (6-layer, 768 dim, 67.5M
 252 parameters – $\sim 2/3$ of base), or BERT-mini [58] (4-layer, 256 dim, 11.3M parameters – $\sim 1/10$ of
 253 base). For query generation (cf. Sec. 4.3), we employ BART-base [27], an encoder-decoder model, to
 254 generate similar questions from each training example’s input question (query). We randomly mask
 255 10% of tokens and inject zero mean Gaussian noise with $\sigma = \{0.1, 0.2\}$ between the encoder and
 256 decoder. See Appendix E for more details on query generation and Appendix F.1 for hyperparameters.

257 5.2 DE to DE distillation

258 We employ AR2 [63]³ and SentenceBERT-
 259 v5 [50]⁴ as teacher DE models for NQ
 260 and MSMARCO. Note that both models
 261 are based on BERT-base. For DE to DE
 262 distillation, we consider two kinds of con-
 263 figurations for the student DE model: (1)
 264 *Symmetric*: We use identical question and
 265 document encoders. We evaluate Distil-
 266 BERT and BERT-mini on both datasets. (2)
 267 *Asymmetric*: The student inherits document
 268 embeddings from the teacher DE model
 269 and *are not* trained during the distillation.
 270 For query encoder, we use DistilBERT or
 271 BERT-mini which are smaller than docu-
 272 ment encoder.

273 **Student DE model training.** We train stu-
 274 dent DE models using a combination of
 275 (i) one-hot loss (cf. Eq. 9 in Appendix A)
 276 on training data; (ii) distillation loss in
 277 (cf. Eq. 11 in Appendix A); and (iii) em-
 278 bedding matching loss in Eq. 7. We used [CLS]-pooling for all student encoders. Unlike DPR [20]
 279 or AR2, we do not use hard negatives from BM25 or other models, which greatly simplifies our
 280 distillation procedure.

281 **Results and discussion.** To understand the impact of various proposed configurations and losses, we
 282 train models by sequentially adding components and evaluate their retrieval performance on NQ and
 283 MSMARCO dev set as shown in Table 1 and Table 3 respectively. (See Table 6 in Appendix F.2 for
 284 performance on NQ in terms of the relaxed recall and Table 7 in Appendix F.3 for MSMARCO in
 285 terms of nDCG@10.)

286 We begin by training a symmetric DE without distillation. As expected, moving to distillation brings
 287 in considerable gains. Next, we swap the student document encoder with document embeddings
 288 from the teacher (non-trainable), which leads to a good jump in the performance. Now we can
 289 introduce EmbedDistill with Eq. 7 for aligning query representations between student and teacher.
 290 The two losses are combined with weight of 1.0 (except for BERT-mini models in the presence of
 291 query generation with 5.0). This improves performance significantly, e.g., it provides ~ 3 and ~ 5
 292 points increase in recall@5 on NQ with students based on DistilBERT and BERT-mini, respectively
 293 (Table 1). We further explore the utility of EmbedDistill in aligning the teacher and student embedding
 294 spaces in Appendix G.1.

295 On top of the two losses (standard distillation and embedding matching), we also use $R_{\text{Emb}, Q}(t, s; Q')$
 296 from Sec. 4.3 on 2 additional questions (per input question) generated from BART. We also try a
 297 variant where we eliminate the standard distillation loss and only employ the embedding matching
 298 loss in Eq. 7 along with inheriting teacher’s document embeddings. This configuration without the
 299 standard distillation loss leads to excellent performance (with query generation again providing
 300 additional gains in most cases.)

Table 3: Performance of various DE models on MSMARCO dev set for both *re-ranking* and *retrieval* tasks (full corpus). The teacher model (110.1M parameter BERT-base models as both encoders) for re-ranking achieves MRR@10 of 36.8 and that for retrieval get MRR@10 of 37.2. The table shows performance (in MRR@10) of the symmetric DE student model (67.5M or 11.3M transformer as both encoders), and asymmetric DE student model (67.5M or 11.3M transformer as query encoder and document embeddings inherited from the teacher).

Method	Re-ranking		Retrieval	
	67.5M	11.3M	67.5M	11.3M
Train student directly	27.0	23.0	22.6	18.6
+ Distill from teacher	34.6	30.4	35.0	28.6
+ Inherit doc embeddings	35.2	32.1	35.7	30.3
+ Query embedding matching	36.2	35.0	35.4	40.8
+ Query generation	36.2	34.4	37.2	34.8
Train student using only embedding matching and inherit doc embeddings	36.5	33.5	36.6	31.4
+ Query generation	36.4	34.1	36.7	32.8

³<https://github.com/microsoft/AR2/tree/main/AR2>

⁴<https://huggingface.co/sentence-transformers/msmarco-bert-base-dot-v5>

301 It is worth highlighting that DE models trained with
 302 the proposed methods (e.g., asymmetric DE with em-
 303 bedding matching and generation) achieve 99% of
 304 the performance in both NQ/MSMARCO tasks with
 305 a query encoder that is 2/3rd the size of that of the
 306 teacher. Furthermore, even with 1/10th size of the
 307 query encoder, our proposal can achieve 95-97% of
 308 the performance. This is particularly useful for la-
 309 tency critical applications with minimal impact on
 310 the final performance.

311 Finally, we take our best student models, i.e., one
 312 trained using with additional embedding matching
 313 loss and using data augmentation from query gen-
 314 eration, and evaluate on test sets. We compare with
 315 various prior work and note that most prior work used
 316 considerably bigger models in terms of parameters,
 317 depth (12 or 24 layers), or width (upto 1024 dims). For NQ test set results are reported in Table 2, but
 318 as MSMARCO does not have any public test set, we instead present results for the BEIR benchmark
 319 in Table 4. Note we also provide evaluation of our SentenceBERT teacher achieving very high
 320 performance on the benchmark which can be of independent interest (please refer to Appendix F.4
 321 for details). For both NQ and BEIR, our approach obtains competitive student model with fewer than
 322 50% of the parameters: even with 6 layers, our student model is very close (98-99%) to its teacher.

323 5.3 CE to DE distillation

324 We consider two CE teachers for MSMARCO re-
 325 ranking task⁵: a standard [CLS]-pooled CE teacher,
 326 and the Dual-pooled CE teacher (cf. Sec. 4.2). Both
 327 teachers are based on RoBERTa-base and trained on
 328 triples in the training set for 300K steps with cross-
 329 entropy loss.

330 **Student DE model training.** We considered the fol-
 331 lowing distillation variants: standard score-based di-
 332 stillation from the [CLS]-pooled teacher, and our novel
 333 Dual-pooled CE teacher (with and without embed-
 334 ding matching loss). For each variant, we initialize en-
 335 coders of the student DE model with two RoBERTa-
 336 base models and train for 500K steps on the train-
 337 ing triples. We performed the naïve joint embedding
 338 matching for the [CLS]-pooled teacher (cf. Sec. 4.2) and employed the query embedding matching
 339 (cf. Eq. 7) for the Dual-pooled CE teacher. In either case, embedding-matching loss is added on top of
 340 the standard cross entropy loss with the weight of 1.0 (when used).

341 **Results and discussion.** Table 5 evaluates the effectiveness of the dual pooling and the embedding
 342 matching for CE to DE distillation. As described in Sec. 4.2, the traditional [CLS]-pooled teacher did
 343 not provide any useful embedding for the embedding matching (see Appendix G.2 for the further
 344 analysis of the resulting embedding space). However, with the Dual-pooled teacher, embedding
 345 matching does boost student’s performance.

346 6 Related work

347 Here, we position our EmbedDistill work with respect to prior work on distillation and data augmenta-
 348 tion for Transformers-based IR models. We also cover prior efforts on aligning representations during
 349 distillation for *non-IR* settings. Unlike our problem setting where the DE student is factorized, these
 350 works mainly consider distilling a single large Transformer into a smaller one.

351 **Distillation for IR.** Traditional distillation techniques have been widely applied in the IR literature,
 352 often to distill a teacher CE model to a student DE model [28, 8]. Recently, distillation from a DE

Table 4: Average BEIR performance of our DE teacher and EmbedDistill student models and their numbers of trainable parameters. Both models are trained on MSMARCO and evaluated on 14 other datasets (the average does not include MSMARCO). The full table is at Appendix F.4. With EmbedDistill, student materializes most of the performance of the teacher on the unforeseen datasets.

Method	#Layers	nDCG@10	R@100
DPR [21]	12	22.5	47.7
ANCE [60]	12	40.5	60.0
TAS-B [15]	6	42.8	64.8
GenQ [57]	6	42.5	64.2
Our teacher [50]	12 (220.2M)	45.7	65.1
EmbedDistill	6 (67.5M)	44.0	63.5

Table 5: Performance of DE models distilled from [CLS]-pooled and Dual-pooled CE models on MSMARCO re-ranking task (original top1000 dev). While both teacher models perform similarly, embedding matching-based distillation only works with the Dual-pooled teacher. See Appendix F for nDCG@10 metric.

Method	MRR@10
[CLS]-pooled teacher	37.1
Dual-pooled teacher	37.0
Standard distillation from [CLS]-pooled teacher	33.0
+Joint matching	32.4
Standard distillation from Dual-pooled teacher	33.3
+Query matching	33.7

⁵Note: Full retrieval is prohibitively expensive with CE models.

353 model (with complex late interaction) to another DE model (with inner-product scoring) has also been
354 considered [29, 15]. As for distilling across different model architectures, Lu et al. [31], Izacard and
355 Grave [16] consider distillation from a teacher CE model to a student DE model. Hofstätter et al. [14]
356 conduct an extensive study of knowledge distillation across a wide-range of model architectures. Most
357 existing distillation schemes for IR rely on only teacher scores; by contrast, we propose a geometric
358 approach that also utilizes the teacher *embeddings*. Many recent efforts [48, 51, 56] show that iterative
359 multi-stage (self-)distillation improves upon single-stage distillation [48, 51, 56]. These approaches
360 use a model from the previous stage to obtain labels [56] as well as mine harder-negatives [60]. We
361 only focus on the single-stage distillation in this paper. Multi-stage procedures are complementary to
362 our work, as one can employ our proposed embedding-matching approach in various stages of such a
363 procedure. Interestingly, we demonstrate in Sec. 5 that our proposed EmbedDistill can successfully
364 benefit from high quality models trained with such complex procedures [50, 63]. In particular, our
365 single-stage distillation method can transfer almost all of their performance gains to even smaller
366 models. Also to showcase that our method brings gain orthogonal to how teacher was trained, we
367 conduct experiments with single-stage trained teacher in Appendix F.5.

368 **Distillation with representation alignments.** Outside of the IR context, a few prior works proposed
369 to utilize alignment between hidden layers during distillation [52, 55, 18, 1, 64]. Chen et al. [7] utilize
370 the representation alignment to re-use teacher’s classification layer for image classification. Unlike
371 these works, our work is grounded in a rigorous theoretical understanding of the teacher-student
372 (generalization) gap for IR models. Further, our work differs from these as it needs to address multiple
373 challenges presented by an IR setting: 1) cross-architecture distillation such as CE to DE distillation;
374 2) partial representation alignment of query or document representations as opposed to aligning for
375 the entire input, i.e., a query-documents pair; and 3) catering representation alignment approach to
376 novel IR setups such as asymmetric DE configuration. To the best of our knowledge, our work is first
377 in the IR literature that goes beyond simply matching scores (or its proxies) for distillation.

378 **Semi-supervised learning for IR.** Data augmentation or semi-supervised learning has been previ-
379 ously used to ensure data efficiency in IR [see, e.g., 35, 66]. More interestingly, data augmentation
380 have enabled performance improvements as well. Doc2query [45, 44] performs document expansion
381 by generating queries that are relevant to the document and appending those queries to the
382 document. Query expansion has also been considered, e.g., for document re-ranking [67]. Notably,
383 generating synthetic (query, passage, answer) triples from a text corpus to augment existing training
384 data for QA systems also leads to significant gains [2, 47]. Furthermore, even zero-shot approaches,
385 where no labeled query-document pairs are used, can also perform competitively to supervised
386 methods [26, 17, 33, 54]. Unlike these works, we utilize query-generation capability to ensure tighter
387 alignment between the embedding spaces of the teacher and student.

388 **Richer transformers-based architectures for IR.** Besides DE and CE models (cf. Sec. 2), interme-
389 diate configurations [36, 22, 42, 32] have been proposed. Such models independently encode query
390 and document before applying a more complex *late interaction* between the two. Nogueira et al.
391 [46] explore *generative* encoder-decoder style model for re-ranking. In this paper, we focus on basic
392 DE/CE models to showcase the benefits of our proposed geometric distillation approach. Exploring
393 embedding matching for aforementioned architectures is an interesting avenue for future work.

394 7 Conclusion

395 We propose EmbedDistill — a novel distillation method for IR that goes beyond simple score matching.
396 En route, we provide a theoretical understanding of the teacher-student generalization gap in an IR
397 setting which not only motivated EmbedDistill but also inspired new design choices for the student DE
398 models: (a) reusing the teacher’s document encoder in the student and (b) aligning query embeddings
399 of the teacher and student. This simple approach delivers consistent quality and computational gains
400 in practical deployments and we demonstrate them on MSMARCO, NQ, and BEIR benchmarks.
401 Finally, we found EmbedDistill retains 95-97% of the teacher performance to with 1/10th size students.

402 **Limitations.** As discussed in Sec. 4.2 and 5.3, EmbedDistill requires modifications in the CE scoring
403 function to be effective. In terms of underlying IR model architectures, we only explore Transformer-
404 based models in our experiments; primarily due to their widespread utilization. That said, we expect
405 our results to extend to non-Transformer architectures such as MLPs. Finally, we note that our
406 experiments only consider NLP domains, and exploring other modalities (e.g., vision) or multi-modal
407 settings (e.g., image-to-text search) is left as an interesting avenue for future work.

408 **References**

- 409 [1] Gustavo Aguilar, Yuan Ling, Yu Zhang, Benjamin Yao, Xing Fan, and Chenlei Guo. Knowledge
410 distillation from internal representations. In *Proceedings of the AAAI Conference on Artificial*
411 *Intelligence*, volume 34, pages 7350–7357, 2020.
- 412 [2] Chris Alberti, Daniel Andor, Emily Pitler, Jacob Devlin, and Michael Collins. Synthetic
413 QA corpora generation with roundtrip consistency. In *Proceedings of the 57th Annual*
414 *Meeting of the Association for Computational Linguistics*, pages 6168–6173, Florence, Italy,
415 July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1620. URL
416 <https://aclanthology.org/P19-1620>.
- 417 [3] Yoshua Bengio and Jean-SÉbastien Senecal. Adaptive importance sampling to accelerate
418 training of a neural probabilistic language model. *IEEE Transactions on Neural Networks*, 19
419 (4):713–722, 2008. doi: 10.1109/TNN.2007.912312.
- 420 [4] Olivier Bousquet, Stéphane Boucheron, and Gábor Lugosi. *Introduction to Statistical Learning*
421 *Theory*, pages 169–207. Springer Berlin Heidelberg, Berlin, Heidelberg, 2004. ISBN
422 978-3-540-28650-9. doi: 10.1007/978-3-540-28650-9_8. URL [https://doi.org/10.1007/](https://doi.org/10.1007/978-3-540-28650-9_8)
423 [978-3-540-28650-9_8](https://doi.org/10.1007/978-3-540-28650-9_8).
- 424 [5] Cristian Bucilă, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In
425 *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and*
426 *Data Mining*, KDD '06, pages 535–541, New York, NY, USA, 2006. ACM.
- 427 [6] Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. Reading Wikipedia to answer
428 open-domain questions. In *Proceedings of the 55th Annual Meeting of the Association for*
429 *Computational Linguistics (Volume 1: Long Papers)*, pages 1870–1879, Vancouver, Canada,
430 July 2017. Association for Computational Linguistics. doi: 10.18653/v1/P17-1171. URL
431 <https://aclanthology.org/P17-1171>.
- 432 [7] Defang Chen, Jian-Ping Mei, Hailin Zhang, Can Wang, Yan Feng, and Chun Chen. Knowledge
433 distillation with the reused teacher classifier. In *Proceedings of the IEEE/CVF Conference on*
434 *Computer Vision and Pattern Recognition*, pages 11933–11942, 2022.
- 435 [8] Xuanang Chen, Ben He, Kai Hui, Le Sun, and Yingfei Sun. Simplified tinybert: Knowledge
436 distillation for document retrieval. In Djoerd Hiemstra, Marie-Francine Moens, Josiane Mothe,
437 Raffaele Perego, Martin Potthast, and Fabrizio Sebastiani, editors, *Advances in Information*
438 *Retrieval*, pages 241–248, Cham, 2021. Springer International Publishing. ISBN 978-3-030-
439 72240-1.
- 440 [9] Zhuyun Dai and Jamie Callan. Deeper text understanding for IR with contextual neural language
441 modeling. In Benjamin Piwowarski, Max Chevalier, Éric Gaussier, Yoelle Maarek, Jian-Yun
442 Nie, and Falk Scholer, editors, *Proceedings of the 42nd International ACM SIGIR Conference*
443 *on Research and Development in Information Retrieval, SIGIR 2019, Paris, France, July 21-25,*
444 *2019*, pages 985–988. ACM, 2019.
- 445 [10] Zhuyun Dai and Jamie Callan. Context-aware sentence/passage term importance estimation for
446 first stage retrieval. *arXiv preprint arXiv:1910.10687*, 2019.
- 447 [11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of
448 deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and
449 Tamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter*
450 *of the Association for Computational Linguistics: Human Language Technologies, NAACL-*
451 *HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages
452 4171–4186. Association for Computational Linguistics, 2019.
- 453 [12] Ruiqi Guo, Philip Sun, Erik Lindgren, Quan Geng, David Simcha, Felix Chern, and Sanjiv
454 Kumar. Accelerating large-scale inference with anisotropic vector quantization. In *International*
455 *Conference on Machine Learning*, 2020. URL <https://arxiv.org/abs/1908.10396>.
- 456 [13] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network,
457 2015.

- 458 [14] Sebastian Hofstätter, Sophia Althammer, Michael Schröder, Mete Sertkan, and Allan Hanbury.
459 Improving efficient neural ranking models with cross-architecture knowledge distillation. *CoRR*,
460 abs/2010.02666, 2020. URL <https://arxiv.org/abs/2010.02666>.
- 461 [15] Sebastian Hofstätter, Sheng-Chieh Lin, Jheng-Hong Yang, Jimmy Lin, and Allan Hanbury.
462 Efficiently teaching an effective dense retriever with balanced topic aware sampling. In *Pro-*
463 *ceedings of the 44th International ACM SIGIR Conference on Research and Development in*
464 *Information Retrieval*, SIGIR '21, page 113–122, New York, NY, USA, 2021. Association
465 for Computing Machinery. ISBN 9781450380379. doi: 10.1145/3404835.3462891. URL
466 <https://doi.org/10.1145/3404835.3462891>.
- 467 [16] Gautier Izacard and Edouard Grave. Distilling knowledge from reader to retriever for question
468 answering. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=NTEz-6wysdb>.
- 470 [17] Gautier Izacard, Mathild Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand
471 Joulin, and Edouard Grave. Unsupervised dense information retrieval with contrastive learning.
472 *arXiv preprint arXiv:2112.09118*, 2021.
- 473 [18] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and
474 Qun Liu. TinyBERT: Distilling BERT for natural language understanding. In *Findings of the*
475 *Association for Computational Linguistics: EMNLP 2020*, pages 4163–4174, Online, November
476 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.372.
477 URL <https://aclanthology.org/2020.findings-emnlp.372>.
- 478 [19] Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with gpus. *IEEE*
479 *Transactions on Big Data*, 7(3):535–547, 2021. doi: 10.1109/TBDATA.2019.2921572.
- 480 [20] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov,
481 Danqi Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering.
482 In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Process-*
483 *ing (EMNLP)*, pages 6769–6781, Online, November 2020. Association for Computational
484 Linguistics.
- 485 [21] Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov,
486 Danqi Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering.
487 *arXiv preprint arXiv:2004.04906*, 2020.
- 488 [22] Omar Khattab and Matei Zaharia. *ColBERT: Efficient and Effective Passage Search via*
489 *Contextualized Late Interaction over BERT*, page 39–48. Association for Computing Machinery,
490 New York, NY, USA, 2020. ISBN 9781450380164.
- 491 [23] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris
492 Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion
493 Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav
494 Petrov. Natural questions: A benchmark for question answering research. *Transactions of the*
495 *Association for Computational Linguistics*, 7:452–466, 2019. doi: 10.1162/tacl_a_00276. URL
496 <https://aclanthology.org/Q19-1026>.
- 497 [24] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris
498 Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. Natural questions: a
499 benchmark for question answering research. *Transactions of the Association for Computational*
500 *Linguistics*, 7:453–466, 2019.
- 501 [25] Michel Ledoux and Michel Talagrand. *Probability in Banach spaces*. Springer-Verlag, 1991.
- 502 [26] Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. Latent retrieval for weakly supervised
503 open domain question answering. In Anna Korhonen, David R. Traum, and Lluís Màrquez,
504 editors, *Proceedings of the 57th Conference of the Association for Computational Linguistics,*
505 *ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 6086–6096.
506 Association for Computational Linguistics, 2019.

- 507 [27] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer
508 Levy, Veselin Stoyanov, and Luke Zettlemoyer. BART: Denoising sequence-to-sequence pre-
509 training for natural language generation, translation, and comprehension. In *Proceedings of*
510 *the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880,
511 Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.
512 703. URL <https://aclanthology.org/2020.acl-main.703>.
- 513 [28] Canjia Li, Andrew Yates, Sean MacAvaney, Ben He, and Yingfei Sun. Parade: Passage repre-
514 sentation aggregation for document reranking. *arXiv preprint arXiv:2008.09093*, 2020.
- 515 [29] Sheng-Chieh Lin, Jheng-Hong Yang, and Jimmy Lin. In-batch negatives for knowledge dis-
516 tillation with tightly-coupled teachers for dense retrieval. In *Proceedings of the 6th Work-*
517 *shop on Representation Learning for NLP (Repl4NLP-2021)*, pages 163–173, Online, August
518 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.repl4nlp-1.17. URL
519 <https://aclanthology.org/2021.repl4nlp-1.17>.
- 520 [30] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike
521 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining
522 approach. *arXiv preprint arXiv:1907.11692*, 2019.
- 523 [31] Wenhao Lu, Jian Jiao, and Ruofei Zhang. Twinbert: Distilling knowledge to twin-structured
524 compressed bert models for large-scale retrieval. In *Proceedings of the 29th ACM International*
525 *Conference on Information & Knowledge Management, CIKM '20*, page 2645–2652, New
526 York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450368599. doi:
527 10.1145/3340531.3412747. URL <https://doi.org/10.1145/3340531.3412747>.
- 528 [32] Yi Luan, Jacob Eisenstein, Kristina Toutanova, and Michael Collins. Sparse, dense, and
529 attentional representations for text retrieval. *Transactions of the Association for Computational*
530 *Linguistics*, 9:329–345, 2021. doi: 10.1162/tacl_a_00369. URL [https://aclanthology.org/](https://aclanthology.org/2021.tacl-1.20)
531 [2021.tacl-1.20](https://aclanthology.org/2021.tacl-1.20).
- 532 [33] Ji Ma, Ivan Korotkov, Yinfei Yang, Keith Hall, and Ryan McDonald. Zero-shot neural pas-
533 sage retrieval via domain-targeted synthetic question generation. In *Proceedings of the 16th*
534 *Conference of the European Chapter of the Association for Computational Linguistics: Main*
535 *Volume*, pages 1075–1088, Online, April 2021. Association for Computational Linguistics. doi:
536 10.18653/v1/2021.eacl-main.92. URL <https://aclanthology.org/2021.eacl-main.92>.
- 537 [34] Sean MacAvaney, Andrew Yates, Arman Cohan, and Nazli Goharian. CEDR: Contextualized
538 embeddings for document ranking. In *Proceedings of the 42nd International ACM SIGIR*
539 *Conference on Research and Development in Information Retrieval, SIGIR'19*, page 1101–1104,
540 New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450361729. doi:
541 10.1145/3331184.3331317. URL <https://doi.org/10.1145/3331184.3331317>.
- 542 [35] Sean MacAvaney, Andrew Yates, Kai Hui, and Ophir Frieder. Content-based weak supervision
543 for ad-hoc re-ranking. In *Proceedings of the 42nd International ACM SIGIR Conference on*
544 *Research and Development in Information Retrieval, SIGIR'19*, page 993–996, New York, NY,
545 USA, 2019. Association for Computing Machinery. ISBN 9781450361729. doi: 10.1145/
546 3331184.3331316. URL <https://doi.org/10.1145/3331184.3331316>.
- 547 [36] Sean MacAvaney, Franco Maria Nardini, Raffaele Perego, Nicola Tonellotto, Nazli Goharian,
548 and Ophir Frieder. *Efficient Document Re-Ranking for Transformers by Precomputing Term*
549 *Representations*, page 49–58. Association for Computing Machinery, New York, NY, USA,
550 2020. ISBN 9781450380164.
- 551 [37] Aditya Menon, Sadeep Jayasumana, Ankit Singh Rawat, Seungyeon Kim, Sashank Reddi,
552 and Sanjiv Kumar. In defense of dual-encoders for neural ranking. In Kamalika Chaudhuri,
553 Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings*
554 *of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of*
555 *Machine Learning Research*, pages 15376–15400. PMLR, 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/menon22a.html>.
556

- 557 [38] Bhaskar Mitra and Nick Craswell. An introduction to neural information retrieval. *Foundations*
558 *and Trends® in Information Retrieval*, 13(1):1–126, 2018. ISSN 1554-0669. doi: 10.1561/
559 1500000061. URL <http://dx.doi.org/10.1561/1500000061>.
- 560 [39] Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Tworek,
561 Qiming Yuan, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, et al. Text and code embeddings
562 by contrastive pre-training. *arXiv preprint arXiv:2201.10005*, 2022.
- 563 [40] Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder,
564 and Li Deng. MS MARCO: A human generated machine reading comprehension dataset. In
565 Tarek Richard Besold, Antoine Bordes, Artur S. d’Avila Garcez, and Greg Wayne, editors,
566 *Proceedings of the Workshop on Cognitive Computation: Integrating neural and symbolic*
567 *approaches 2016*, volume 1773 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2016.
- 568 [41] Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernandez Abrego, Ji Ma, Vincent Zhao,
569 Yi Luan, Keith Hall, Ming-Wei Chang, and Yinfei Yang. Large dual encoders are generalizable
570 retrievers. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language*
571 *Processing*, pages 9844–9855, Abu Dhabi, United Arab Emirates, December 2022. Association
572 for Computational Linguistics. URL <https://aclanthology.org/2022.emnlp-main.669>.
- 573 [42] Ping Nie, Yuyu Zhang, Xiubo Geng, Arun Ramamurthy, Le Song, and Daxin Jiang. DC-BERT:
574 decoupling question and document for efficient contextual encoding. In Jimmy Huang, Yi Chang,
575 Xueqi Cheng, Jaap Kamps, Vanessa Murdock, Ji-Rong Wen, and Yiqun Liu, editors, *Proceedings*
576 *of the 43rd International ACM SIGIR conference on research and development in Information*
577 *Retrieval, SIGIR 2020, Virtual Event, China, July 25-30, 2020*, pages 1829–1832. ACM, 2020.
578 doi: 10.1145/3397271.3401271. URL <https://doi.org/10.1145/3397271.3401271>.
- 579 [43] Rodrigo Nogueira and Kyunghyun Cho. Passage re-ranking with BERT. *CoRR*, abs/1901.04085,
580 2019. URL <http://arxiv.org/abs/1901.04085>.
- 581 [44] Rodrigo Nogueira, Jimmy Lin, and AI Epistemic. From doc2query to docttttquery. *Online*
582 *preprint*, 6, 2019.
- 583 [45] Rodrigo Nogueira, Wei Yang, Jimmy Lin, and Kyunghyun Cho. Document expansion by query
584 prediction. *arXiv preprint arXiv:1904.08375*, 2019.
- 585 [46] Rodrigo Nogueira, Zhiying Jiang, Ronak Pradeep, and Jimmy Lin. Document ranking
586 with a pretrained sequence-to-sequence model. In *Findings of the Association for Com-*
587 *putational Linguistics: EMNLP 2020*, pages 708–718, Online, November 2020. Associa-
588 tion for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.63. URL <https://aclanthology.org/2020.findings-emnlp.63>.
- 590 [47] Barlas Oğuz, Kushal Lakhota, Anchit Gupta, Patrick Lewis, Vladimir Karpukhin, Aleksandra
591 Piktus, Xilun Chen, Sebastian Riedel, Wen-tau Yih, Sonal Gupta, et al. Domain-matched
592 pre-training tasks for dense retrieval. *arXiv preprint arXiv:2107.13602*, 2021.
- 593 [48] Yingqi Qu, Yuchen Ding, Jing Liu, Kai Liu, Ruiyang Ren, Wayne Xin Zhao, Daxiang Dong, Hua
594 Wu, and Haifeng Wang. RocketQA: An optimized training approach to dense passage retrieval
595 for open-domain question answering. In Kristina Toutanova, Anna Rumshisky, Luke Zettle-
596 moyer, Dilek Hakkani-Tür, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty,
597 and Yichao Zhou, editors, *Proceedings of the 2021 Conference of the North American Chapter*
598 *of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT*
599 *2021, Online, June 6-11, 2021*, pages 5835–5847. Association for Computational Linguistics,
600 2021.
- 601 [49] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,
602 Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified
603 text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020. URL
604 <http://jmlr.org/papers/v21/20-074.html>.
- 605 [50] Nils Reimers, Iryna Gurevych, and Iryna Gurevych. Sentence-BERT: Sentence embeddings
606 using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods*
607 *in Natural Language Processing*. Association for Computational Linguistics, 11 2019. URL
608 <http://arxiv.org/abs/1908.10084>.

- 609 [51] Ruiyang Ren, Yingqi Qu, Jing Liu, Wayne Xin Zhao, Qiaoqiao She, Hua Wu, Haifeng Wang,
610 and Ji-Rong Wen. Rocketqav2: A joint training method for dense passage retrieval and passage
611 re-ranking. In *Proceedings of EMNLP*, 2021.
- 612 [52] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and
613 Yoshua Bengio. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*, 2014.
- 614 [53] Devendra Sachan, Mostofa Patwary, Mohammad Shoeybi, Neel Kant, Wei Ping, William L.
615 Hamilton, and Bryan Catanzaro. End-to-end training of neural retrievers for open-domain
616 question answering. In *Proceedings of the 59th Annual Meeting of the Association for
617 Computational Linguistics and the 11th International Joint Conference on Natural Lan-
618 guage Processing (Volume 1: Long Papers)*, pages 6648–6662, Online, August 2021. As-
619 sociation for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.519. URL <https://aclanthology.org/2021.acl-long.519>.
- 621 [54] Devendra Singh Sachan, Mike Lewis, Mandar Joshi, Armen Aghajanyan, Wen-tau Yih, Joelle
622 Pineau, and Luke Zettlemoyer. Improving passage retrieval with zero-shot question generation.
623 *arXiv preprint arXiv:2204.07496*, 2022.
- 624 [55] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version
625 of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- 626 [56] Keshav Santhanam, Omar Khattab, Jon Saad-Falcon, Christopher Potts, and Matei Zaharia. Col-
627 bertv2: Effective and efficient retrieval via lightweight late interaction. *CoRR*, abs/2112.01488,
628 2021.
- 629 [57] Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych.
630 BEIR: A heterogeneous benchmark for zero-shot evaluation of information retrieval models. In
631 *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks
632 Track (Round 2)*, 2021. URL <https://openreview.net/forum?id=wCu6T5xFjeJ>.
- 633 [58] Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Well-read students learn
634 better: On the importance of pre-training compact models. *arXiv preprint arXiv:1908.08962*,
635 2019.
- 636 [59] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
637 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proceedings of the 31st
638 International Conference on Neural Information Processing Systems, NIPS’17*, page 6000–6010,
639 Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.
- 640 [60] Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul N. Bennett, Junaid
641 Ahmed, and Arnold Overwijk. Approximate nearest neighbor negative contrastive learning
642 for dense text retrieval. In *International Conference on Learning Representations*, 2021. URL
643 <https://openreview.net/forum?id=zeFrfgyzln>.
- 644 [61] Nishant Yadav, Nicholas Monath, Rico Angell, Manzil Zaheer, and Andrew McCallum. Efficient
645 nearest neighbor search for cross-encoder models using matrix factorization. In *Proceedings of
646 the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2171–2194,
647 Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics.
648 URL <https://aclanthology.org/2022.emnlp-main.140>.
- 649 [62] Zeynep Akkalyoncu Yilmaz, Wei Yang, Haotian Zhang, and Jimmy Lin. Cross-domain modeling
650 of sentence-level evidence for document retrieval. In *Proceedings of the 2019 Conference on
651 Empirical Methods in Natural Language Processing and the 9th International Joint Conference
652 on Natural Language Processing (EMNLP-IJCNLP)*, pages 3490–3496, Hong Kong, China,
653 November 2019. Association for Computational Linguistics.
- 654 [63] Hang Zhang, Yeyun Gong, Yelong Shen, Jiancheng Lv, Nan Duan, and Weizhu Chen. Ad-
655 versarial retriever-ranker for dense text retrieval. In *International Conference on Learning
656 Representations*, 2022. URL <https://openreview.net/forum?id=MR7XubKUFB>.
- 657 [64] Linfeng Zhang and Kaisheng Ma. Improve object detection with feature-based knowledge
658 distillation: Towards accurate and efficient detectors. In *International Conference on Learning
659 Representations*, 2020.

- 660 [65] Shuai Zhang, Lina Yao, Aixin Sun, and Yi Tay. Deep learning based recommender system:
661 A survey and new perspectives. *ACM Comput. Surv.*, 52(1), feb 2019. ISSN 0360-0300. doi:
662 10.1145/3285029. URL <https://doi.org/10.1145/3285029>.
- 663 [66] Chen Zhao, Chenyan Xiong, Jordan Boyd-Graber, and Hal Daumé III. Distantly-supervised
664 dense retrieval enables open-domain question answering without evidence annotation. In
665 *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*,
666 pages 9612–9622, Online and Punta Cana, Dominican Republic, November 2021. As-
667 sociation for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.756. URL
668 <https://aclanthology.org/2021.emnlp-main.756>.
- 669 [67] Zhi Zheng, Kai Hui, Ben He, Xianpei Han, Le Sun, and Andrew Yates. BERT-QE: Con-
670 textualized Query Expansion for Document Re-ranking. In *Findings of the Association for*
671 *Computational Linguistics: EMNLP 2020*, pages 4718–4728, Online, November 2020. As-
672 sociation for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.424. URL
673 <https://aclanthology.org/2020.findings-emnlp.424>.