ADAM: Dense Retrieval Distillation with Adaptive Dark Examples

Anonymous ACL submission

Abstract

To improve the performance of the dualencoder retriever, one effective approach is knowledge distillation from the cross-encoder ranker. Existing works prepare training instances by pairing each query with one positive and a batch of negatives. However, most hard negatives mined by advanced dense retrieval methods are still too trivial for the teacher to distinguish, preventing the teacher from transferring abundant dark knowledge to the student through its soft label. To alleviate this issue, we propose ADAM, a knowledge distillation framework that can better transfer the dark knowledge held in the teacher with Adaptive Dark exAMples. Different from previous works that only rely on one positive and hard negatives as candidate passages, we create dark examples that all have moderate relevance to the query by strengthening negatives and masking positives in the discrete space. Furthermore, as the quality of knowledge held in different training instances varies as measured by the teacher's confidence score, we propose a self-paced distillation strategy that adaptively concentrates on a subset of high-quality instances to conduct our dark-example-based knowledge distillation to help the student learn better. We conduct experiments on two widely-used benchmarks and verify the effectiveness of our method.

1 Introduction

001

017

024

037

041

Information retrieval (IR) that aims to identify relevant passages for a given query is an important topic for both academic and industrial areas, and has powered many downstream tasks such as opendomain QA (Chen et al., 2017) and knowledgegrounded conversation (Dinan et al., 2018). Typically, IR systems usually follow the retrieve-andre-rank paradigm (Hofstätter et al., 2020; Huang et al., 2020; Zou et al., 2021) where a fast retriever first retrieved a bundle of relevant passages from a large-scale corpus through pre-built indices and



Figure 1: Distributions of the prediction for the crossencoder of R^2 anker (Zhou et al., 2023) over MS-MARCO. POS and NEG mean the distribution of positive and hard negatives respectively. The hard negatives are provided by RocketQAv2 (Ren et al., 2021c).

then a more sophisticated ranker comes to re-rank these candidate passages to further obtain more accurate retrieval results.

Under this paradigm, recent years have witnessed a growing number of works that utilize pre-trained language models (PLMs) (Qu et al., 2021; Gao and Callan, 2021b) as retrievers and rankers to build IR systems. Among these efforts, there are two commonly adopted architectures: cross-encoder (Devlin et al., 2019a) that measure the relevance of a query-passage pair through jointly modeling their deep interactions; dual-encoder (Karpukhin et al., 2020; Qu et al., 2021) that encodes queries and passages separately into dense representations and calculate the similarity. Although dual-encoders are efficient for billions of indices, they suffer from inferior performance compared with cross-encoders since they can't capture the fine-grained semantic relevance between the query and the passage due to the absence of their deep interactions (Luan et al., 2021a). To help dual-encoders achieve better retrieval performance, a common practice is to draw on the powerful but cumbersome cross-encoder through knowledge distillation (Yang et al., 2020; Zhang et al., 2022; Ren et al., 2021c; Zeng et al., 2022; Lin et al., 2023). Along this line of research, various techniques are proposed to improve the knowledge

043

084

095

transfer including data curriculum (Lin et al., 2023; Zeng et al., 2022), on-the-fly distillation (Zhang et al., 2022; Ren et al., 2021c) and new distillation objectives (Lu et al., 2022; Menon et al., 2022).

Though effective, we argue that existing dense retrieval distillation methods may not fully exploit the dark knowledge deeply held by the teacher. In knowledge distillation (Xu et al., 2018; Lin et al., 2023), the student learns not just the highestscored class from the soft labels provided by the teacher, but also the entire probability distribution over classes, as this contains comprehensive finegrained information referred to as "dark knowledge". However, we empirically find that for existing distillation methods, the soft labels (i.e., the probability distributions over one positive and multiple negatives for a query) given by the teacher are too "sharp", despite they already adopted hard negatives (Ren et al., 2021c). As illustrated in Figure 1, we draw the score distributions of the positive and negative pairs using a pre-trained cross-encoder teacher. It can be observed that the scores for most hard negatives are quite low (concentrated in (-7.5, -2.5)) and distributed far from the positives that have high scores. A similar observation is also drawn by Menon et al. (2022). This phenomenon indicates that even the hard negatives mined by the dense retriever are still too trivial for a well-trained cross-encoder teacher to distinguish, losing most of the utile dark knowledge.

To alleviate this issue, we propose ADAM, a 100 knowledge distillation framework that can better 101 exploit dark knowledge deeply held in the teacher 102 by distillation with adaptive dark examples. Our 103 method originated from the intuition that a good 104 soft label for the retriever to learn should be more 105 smooth, which implies that the provided query-106 passage pairs should diversely distribute from highly-relevant pairs to loosely-relevant pairs from 108 the view of the teacher. To fill the gap between highly-relevant pairs and loosely-relevant pairs ex-110 isting in current negative sampling methods, we 111 propose two approaches to construct dark exam-112 ples that all have moderate relevance to the query. 113 The first approach is to make negatives more rel-114 evant to the query by strengthening the negatives 115 with the positive passage. The second approach is 116 117 to make positives less relevant to the query by replacing some randomly selected tokens with mask 118 tokens. Considering that the newly created pas-119 sages have moderate relevance to the query, we believe it is more appropriate to call them dark 121

examples instead of negatives. With these dark examples added, we successfully make the score distribution smoother as shown in Figure 3(b), so that we can transfer more useful dark knowledge from the teacher. Moreover, since the soft label for different query-positive-negatives have different "sharpness" which we consider as an indication of how well the dark knowledge has been exploited, we further propose a self-paced distillation strategy that adaptively selects those examples whose soft labels are sharp to conduct our dark-example-based distillation to better transfer the dark knowledge. 122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

151

153

154

155

156

157

158

159

160

161

163

164

165

166

167

168

169

170

171

We conduct experiments on two benchmarks, including MS-MARCO (Nguyen et al., 2016) and TREC Deep Learning 2019 (Craswell et al., 2020). In both benchmarks, the model is required to select the best response from a candidate pool. Evaluation results indicate that our method is significantly better than existing models on two benchmarks. *We will release all codes for the easy reproduction*. To sum up, our contributions is three-fold:

- Propose to augment dark examples including reinforced negatives and noisy positives for more effective knowledge distillation in IR;
- Propose to adaptively concentrate on highconfidence training instances to better transfer knowledge;
- Empirical verify of the effectiveness of the proposed approach on two public datasets.

2 Related Works

There are two lines of research related to our work: dense retriever and knowledge distillation.

Dense Retriever. To overcome the vocabulary and semantic mismatch problems existing in conventional term-based approaches such as BM25 (Robertson and Zaragoza, 2009), researchers began to build neural retrievers upon pre-trained language models (Devlin et al., 2019b; Liu et al., 2019). In this way, the whole input text can be represented as a dense vector in a lowdimensional space (e.g., 768) and efficient retrieval can be achieved by approximate nearest neighbor search (ANN) algorithms such as FAISS (Johnson et al., 2019). To learn a good dense retriever, various attempts have been made including hard negative mining (Karpukhin et al., 2020; Luan et al., 2021a; Qu et al., 2021; Xiong et al., 2021; Zhan et al., 2021a), retrieval-oriented pre-training (Lee et al., 2019; Gao and Callan, 2021a,b), knowledge distillation (Ren et al., 2021c; Zhang et al., 2022;

221

222

223

229 230 231

232 233

234

235

237

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

258

259

261

262

263

(1)

To fulfill this goal, the retriever is typically trained with supervised contrastive loss:

retriever is responsible for collecting a bubble of

candidate passages and the ranker further re-ranks

them. Considering the trade-off between efficiency

and accuracy, dual-encoders (DE) (Karpukhin et al.,

2020; Luan et al., 2021a; Qu et al., 2021) are of-

ten chosen as the retriever while cross-encoders

(DE) (Devlin et al., 2019b) are usually adopted as

The dual-encoder-based retriever Enc_{de} is re-

sponsible for encoding the given query q_i and each

of the candidate passage p_i into dense vectors

 $Enc_{de}(q_i), Enc_{de}(p_j) \in \mathbb{R}^h$. Then the relevance

score for q_i and p_j is simply calculated as the inner

 $\mathcal{R}_{\mathsf{de}}(q_i, p_i) = \mathsf{Enc}_{\mathsf{de}}(q_i)^\top \cdot \mathsf{Enc}_{\mathsf{de}}(p_i).$

product of their representations:

the ranker.¹

$$\mathcal{L}_{sup} = -\log \frac{\exp^{\mathcal{R}_{de}(q_i, p_i^+)}}{\exp^{\mathcal{R}_{de}(q_i, p_i^+)} + \sum_{p_{i,j}^- \in \mathbb{P}_i^-} \exp^{\mathcal{R}_{de}(q_i, p_{i,j}^-)}}.$$

Cross-Encoders The cross-encoder ranker Enc_{ce} is in charge of calculating the matching score of q_i and p_i more accurately as it can model their finegrained interactions, and re-ranking the retrieved candidate passages provided by the retriever to improve the retrieval results. Concretely, given a query q_i and a passage p_i , the input is formed as the concatenation of q and p with [CLS] in the beginning and [SEP] as their separation and is fed into transformer (Vaswani et al., 2017). The representation of [CLS] in the top layer is used to calculate the relevance score with a projection head $f(\cdot)$:

$$\mathcal{R}_{ce}(q_i, p_j) = f(Enc_{ce}([CLS], q_i, [SEP], p_j)).$$
(2)

Knowledge Distillation in IR As cross-encoders are more capable of measuring the relevance of q_i and p_j than dual-encoders but at a cost of computational inefficiency, it's promising to transfer the knowledge from the strong cross-encoders to the weak dual-encoders through knowledge distillation (Zhang et al., 2022; Ren et al., 2021c; Zeng et al., 2022; Lu et al., 2022; Lin et al., 2023). In dense retrieval distillation, as both the positive passage p_i^+ and the negatives \mathbb{P}_i^- can be treated uniformly, we use $\mathbb{P}_i = \{p_i^+\} \cup \mathbb{P}_i^-$ to denote the whole candidate set of passages. The relevance

172

173

199

201

202

203

210

211

212

213 214

215

216

217

218

Knowledge Distillation. Knowledge distilla-174 tion (Hinton et al., 2015) aims to transfer the knowl-175 edge from a powerful teacher model to a student 176 model to help it learn better. To achieve this goal, 177 the student model is provided with the teacher's out-178 puts as the supervision signal that it is enforced to 179 mimic. There are multiple types of supervision sig-180 nals for the student to learn, including the teacher's 181 output logits (Hinton et al., 2015), intermediate representations (Romero et al., 2014), relations of 183 representations (Park et al., 2019), etc. In the context of dense retrieval distillation, researchers basically adopt the cross-encoder as the teacher and 186 use the teacher's probability distribution over can-187 didate passages as the supervision signal. On this 188 basis, several studies (Ren et al., 2021c; Zhang et al., 2022; Lu et al., 2022) explored on-the-fly distillation to jointly optimize the teacher and the 191 student, Zeng et al. (2022) and Lin et al. (2023) 192 combined knowledge distillation with curriculum 193 strategies to gradually improve the student. Dif-194 ferent from existing work, we focus on the quality 195 of knowledge held in the teacher's soft label and 196 propose to distill with adaptive dark examples to 197 198 better transfer the dark knowledge to the student.

3 Methodology

In this section, we first introduce the preliminaries in dense retrieval distillation, then present our dark example augmentation method and adaptive distillation with dynamic data selection.

3.1 Preliminary

Task Description In this work, we study the learning of the dense retriever following the general setting of dense retrieval in existing work (Qu et al., 2021; Ren et al., 2021c; Zhang et al., 2022). Formally, there is a training set $\mathcal{D} = \{(q_i, \mathbb{P}_i)\}_{i=1}^n$ where q_i is the query and \mathbb{P}_i is the set of candidate passages. Commonly, \mathbb{P}_i consists of a positive passage p_i^+ and m negative passages $\mathbb{P}_i^- = \{p_{i,j}^-\}_{j=1}^m$ constructed by random negative sampling (Henderson et al., 2017; Gillick et al., 2018) or hard negative mining (Xiong et al., 2020; Karpukhin et al., 2020; Qu et al., 2021). Based on \mathcal{D} , we aim to learn a retriever that can select the most relevant passage from the whole candidate pool.

Dual-Encoders A typical text retrieval system adopts the retrieve-and-rank paradigm, where the 220

Lu et al., 2022), etc. We mainly focus on knowledge distillation in this paper.

¹We will use retriever and dual-encoder interchangeably.

294



Figure 2: Illustration of dark examples. The solid rectangle and triangles mean the gold passage and the negative passages respectively. Dotted rectangles and circles denote noisy positives and mixed samples respectively.

score of q_i and each $p_j \in \mathbb{P}_i$ can be calculated using a dual-encoder Enc_{de} and a cross encoder Enc_{ce} using Eq. 1 and Eq. 2. Then, the probability distributions over candidate passages of the dualencoder and the cross-encoder $p_{de,i}, p_{ce,i} \in \mathbb{R}^{|\mathbb{P}_i|}$ are calculated by normalizing the relevance scores over \mathbb{P}_i , where each element is calculated as:

264

265

271

272

273

275

277

278

279

281

285

292

293

$$\hat{\mathcal{R}}_{de,i}^{j} = \frac{\exp^{\mathcal{R}_{de}(q_{i},p_{j})}}{\sum_{p_{k}\in\mathbb{P}_{i}}e^{\mathcal{R}_{de}(q_{i},p_{k})}}$$

$$\hat{\mathcal{R}}_{ce,i}^{j} = \frac{\exp^{\mathcal{R}_{ce}(q_{i},p_{j})}}{\sum_{p_{k}\in\mathbb{P}_{i}}\exp^{\mathcal{R}_{ce}(q_{i},p_{k})}}.$$
(3)

To distill the knowledge from the cross-encoder to the dual-encoder, the distribution of the crossencoder $\hat{\mathcal{R}}_{ce,i}$ is considered as the soft label that guides the learning of the dual-encoder by minimizing the KL-divergence between $\hat{\mathcal{R}}_{ce,i}$ and $\hat{\mathcal{R}}_{de,i}$:

$$\mathcal{L}_{kd} = -\sum_{(q_i, \mathbb{P}_i) \in \mathcal{D}} \text{KL-Div}(\hat{\mathcal{R}}_{ce,i} || \hat{\mathcal{R}}_{de,i}) \quad (4)$$

3.2 Dark Examples Construction

When transferring the knowledge from the crossencoder teacher to the dual-encoder student using Eq. 4, the set of candidate passages \mathbb{P}_i plays a vital role. Previous works in dense retrieval distillation (Zhang et al., 2022; Ren et al., 2021c; Zeng et al., 2022; Lu et al., 2022; Lin et al., 2023) simply follow the supervised learning setting where they utilize $\mathbb{P}_i = \{p_i^+\} \cup \mathbb{P}_i^-$ as the candidate set. However, by empirical analyses on Fig. 1, we have found that the negative set \mathbb{P}_i^- produced by existing hard negative mining approaches (Qu et al., 2021) is too trivial for the cross-encoder teacher, which makes the soft label provided by the crossencoder teacher too sharp at the positive passage and therefore prevents the student from learning utile dark knowledge hidden in the distribution of other passages (i.e., negatives).

We suppose smoother soft labels naturally obtained (instead of scaled by softmax temperature) can be better knowledge carriers that transfer the dark knowledge. Given the teacher and the query, we point out that the natural way to smoothen the soft label is to operate on the set of candidate passages, or more precisely, to replace the original set of candidate passages \mathbb{P}_i that are either too relevant or too irrelevant from the teacher's view with new ones $\tilde{\mathbb{P}}_i$ whose relevance to the query cannot be easily tell apart by the cross-encoder teacher.

To construct the new set of candidate passages that satisfy this desired characteristic, we propose two dual approaches that operate on the original positive passage p_i^+ and the negative set \mathbb{P}_i^- respectively. We name the newly constructed passages in $\tilde{\mathbb{P}}_i$ dark examples to demonstrate that can no longer be simply categorized into positives and negatives as they have moderate relevance to the query. An illustration of dark examples is shown in Figure 2. It should be noticed that it is the specific setting of knowledge distillation where the supervision signal is derived from the teacher's soft label instead of human labels that make it possible to learn from dark examples.

Sampled Negatives. Early works (Henderson et al., 2017; Gillick et al., 2018) randomly choose negative passages by considering the passages of other query-passage pairs within the same mini-batch as the negatives. More recently, researchers use BM25 (Karpukhin et al., 2020) or dual-encoders (Xiong et al., 2020) to select hard negatives globally from the whole candidate passages with the fast retrieval method (Qu et al., 2021; Ren et al., 2021c). We will compare the effective-ness of random negatives (denoted as Rand) and hard negatives (denoted as Hard) with our method (denoted as Dark) in experiments.

Dark Examples with Reinforced Negatives The reasonable way to create dark examples based on \mathbb{P}_i^- is to make hard negatives harder, or in other words, more relevant to the query. To achieve this goal, it is non-trivial to accurately edit the semantics of a negative passage towards increasing its relevance to the query with controllable text generation techniques. Instead, we propose a rather simple yet effective approach that mixes up query-relevant content with negative passages to directedly strength their relevance to the query. Based on

this motivation, we consider mixing up hard negatives with the positive passage². Formally, given a training example $(q_i, p_i^+, \mathbb{P}_i^-)$, we concatenate p_i^+ with each of the negative passage $p_{i,j}^-$ to form the set of dark examples for q_i :

$$\mathcal{N}_{i}^{rein} = \{p_{i}^{+}[\text{SEP}]p_{i,j}^{-}\}_{j=1}^{m}.$$
 (5)

Here, we choose to mix-up passages at the lexical level instead of the embedding space (Guo et al., 2019) because our method can produce valid language inputs and can preserve the relevant cues while introducing some less-relevant content. We also tried mixing-up negatives with the positive in the embedding space but found this kind of mixup resulted in low-quality predictions of the crossencoder teacher since it has never seen samples based on mixed embeddings during training.

354

364

371

372

373

376

377

387

388

Dark Examples with Noisy Positive Different from the above approach that creates dark examples by making hard negatives harder, we also consider the opposite direction: making the positive passage p^+ not that relevant to the query by introducing noise. We achieve this goal by input-masking (Devlin et al., 2019b). Given the positive passage p_i^+ for the query q_i , we randomly sample a subset of tokens from p_i^+ and replace them with the special token [MASK] with the masking ratio m_r :

$$\mathcal{N}_i^{mask} = \{ \mathsf{MASK}_{m_r}(p_i^+) \}_{m_r}.$$
 (6)

To generate noisy positives with more diverse relevant to the query, we use masking with a variety of masking ratios.

3.3 Distillation with Adaptive Dark Examples

We have elaborated our motivation and approach to create dark examples, the remaining question is how to conduct effective knowledge distillation with dark examples. Existing knowledge distillation methods using all the labeled data without distinction, which we argue is sub-optimal. As knowledge distillation relies on the teacher's prediction as the supervision signal, the "quality" of knowledge held in the teacher's soft label naturally varies among different training examples. We assume that those training examples that the teacher is more confident than others are better carriers of knowledge for three reasons: (1) These instances are far from the decision boundaries of the model, and thus the corresponding passages are more likely to be true positives and true negatives, avoiding data noise. (2) Only the knowledge held in the instances that the teacher can cope with well are reliable and worth to be learned by the student. (3) The teacher's soft label for the high-confidence instances is too sharp, which indicates the dark knowledge held in these reliable instances has not been well exploited. 389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

Therefore, we propose to adaptively concentrate on these high-confidence training instances during the training process to conduct our dark-examplebased knowledge distillation. Formally, for a training instance, we can calculate the log-probability of the positive passage p_i^+ against negatives $\mathbb{P}_i^$ with the teacher as the confidence score:

$$\mathcal{C}(q_i) = \log \frac{\exp^{\mathcal{R}_{ce}(q_i, p_i^+)}}{\exp^{\mathcal{R}_{ce}(q_i, p_i^+)} + \sum_{p_{i,j}^- \in \mathbb{P}_i^-} \exp^{\mathcal{R}_{ce}(q_i, p_{i,j}^-)}}.$$
(7)

Suppose the training process consists of T epochs, in each epoch t, we can sort a batch of training instances \mathcal{B}_t in ascending order based on the confidence scores. Then we adaptively select the subset of instances $\hat{\mathcal{B}}_t$ in the batch that have the highest confidence scores with the ratio $(1 - \frac{t}{2*T})$ to construct dark examples:

$$\tilde{\mathcal{B}}_t = \operatorname*{arg\,max}_{q_i \in \mathcal{B}_t, \tilde{\mathcal{B}}_t \subset \mathcal{B}_t, \|\tilde{\mathcal{B}}_t\| = (1 - \frac{t}{2*T}) \times b} \mathcal{C}(q_i).$$
(8)

where b is the batch size for training.

Thereby, we have two sets in each step of the *t*-th training epoch: the original training batch \mathcal{B}_t and the subset with the highest confidence that has both original candidate passages and our created dark examples $\tilde{\mathcal{B}}_t$. We jointly optimize the student with the supervised loss (Eq. 3.1) on \mathcal{B}_t and the knowledge distillation loss (Eq. 4) on $\tilde{\mathcal{B}}_t$:

$$\mathcal{L}_{t} = \lambda \cdot \sum_{\mathcal{B}_{t} \in \mathcal{D}} \sum_{(q_{i}, \mathbb{P}_{i}) \in \mathcal{B}_{t}} \mathcal{L}_{sup} + \sum_{\hat{\mathcal{B}}_{t} \in \mathcal{D}} \sum_{(q_{i}, \tilde{\mathbb{P}}_{i}) \in \tilde{\mathcal{B}}_{t}} \mathcal{L}_{kd}.$$
(9)

where $\tilde{\mathbb{P}}_i = \{\mathbb{P}_i^- \cup \mathcal{N}_i^{mix} \cup \mathcal{N}_i^{mask}\}\$ is the new candidate set for q_i , and λ is a hyper-parameter as a trade-off between the supervised objective and distillation objective with adaptive dark examples.

4 Experiments

We evaluate our method on two public humanannotated real-world benchmarks, namely MS-Marco and TREC Deep Learning 2019.

 $^{^{2}}$ We also tried to make the hard negatives even harder by mixing up hard negatives with the query following Kalantidis et al. (2020), however, we found little change in performance.

			MS-MARCO Dev			TREC DL 19	
Methods	PLM	KD	MRR@10	R@50	R@1000	NDCG@10	R@100
Sparse retrieval							
BM25 (anserini) (Yang et al., 2017a)	-	-	18.7	59.2	85.7	50.6	-
doc2query (Nogueira et al., 2019b)	-	-	21.5	64.4	89.1	-	-
DeepCT (Dai and Callan, 2019b)	BERThase	-	24.3	69.0	91.0	55.1	-
docTTTTTquery (Nogueira et al., 2019a)	-	-	27.7	75.6	94.7	-	-
UHD-BERT (Jang et al., 2021)	BERThase	-	29.6	77.7	96.1	-	-
COIL-full (Gao et al., 2021)	BERThase	-	35.5	-	96.3	70.4	-
UniCOIL (Lin and Ma, 2021)	BERThase	-	35.2	80.7	95.8	-	-
SPLADE-max (Formal et al., 2021)	BERThase	-	34.0	-	96.5	68.4	-
Unifier _{lexicon} (Shen et al., 2023)	coCon _{base}	\checkmark	39.7	-	98.1	73.3	-
Dense retrieval							
DPR-E (Ren et al., $2021c$)	ERNIE	-	32.5	82.2	97.3	-	-
ANCE (single) (Xiong et al., 2020)	RoBERTabase	-	33.0	-	95.9	65.4	44.5
TAS-Balanced (Hofstätter et al., 2021a)	BERThase	\checkmark	34.0	-	-	71.2	_
ME-BERT (Luan et al., 2021b)	BERTlarge	-	34.3	-	-	-	-
ColBERT (Khattab and Zaharia, 2020a)	BERThase	-	36.0	82.9	96.8	67.0	-
ColBERT v2 (Santhanam et al., 2021)	BERThase	\checkmark	39.7	86.8	98.4	72.0	-
ADORE+STAR (Zhan et al., 2021b)	RoBERTabase	-	34.7	-	-	68.3	-
Condenser (Gao and Callan, 2021a)	BERT _{base}	-	36.6	-	97.4	-	-
RocketQA (Qu et al., 2021)	ERNIE _{base}	-	37.0	85.5	97.9	-	-
PAIR (Ren et al., 2021a)	ERNIE _{base}	-	37.9	86.4	98.2	-	-
CoCondenser (Gao and Callan, 2022)	BERT _{base}	-	38.2	-	98.4	-	-
RocketQAV2 (Ren et al., 2021c)	BERT _{base}	\checkmark	38.8	86.2	98.1	-	-
AR2 (Zhang et al., 2022)	BERT _{base}	\checkmark	39.5	-	98.6	-	-
CL-DRD (Zeng et al., 2022)	DistilBERT	\checkmark	38.2	-	-	72.5	45.3
ERNIE-Search (Lu et al., 2022)	BERT _{base}	\checkmark	40.1	87.7	98.2	-	-
RetroMAE (Xiao et al., 2022)	BERT _{base}	\checkmark	39.3	87.0	98.5	-	-
Unifier _{dense} (Shen et al., 2023)	coCon _{base}	\checkmark	38.8	-	97.6	71.1	-
bi-SimLM (Wang et al., 2023)	BERT _{base}	\checkmark	39.1	87.3	98.6	69.8	-
PROD (Lin et al., 2023)	ERNIE-2.0-BASE	\checkmark	39.3	87.1	98.4	73.3	48.4
InDi (Cohen et al., 2024)	coCon _{base}	-	38.8	86.6	98.5	-	-
Rand KD (<i>Teacher</i> = <i>RocketQAV2</i>)	BERT _{base}	\checkmark	38.14	86.92	98.17	-	-
Hard KD (<i>Teacher</i> = <i>RocketQAV2</i>)	BERT _{base}	\checkmark	39.13	87.60	98.51	-	-
ADAM (<i>Teacher</i> = <i>RocketQAV2</i>)	BERT _{base}	\checkmark	39.79	88.07	98.64	72.1	50.3
Rand KD (<i>Teacher</i> = R^2 anker)	BERT _{base}	\checkmark	38.13	85.96	97.87	-	-
Hard KD (<i>Teacher</i> = R^2 anker)	BERT _{base}	\checkmark	39.99	87.62	98.12	-	-
ADAM (<i>Teacher</i> = R^2 anker)	BERT _{base}	\checkmark	<u>41.00</u>	<u>88.54</u>	98.48	<u>73.4</u>	49.8

Table 1: Passage retrieval results on MS-MARCO and TREC DL 19 datasets. PLM is the abbreviation of the pre-trained language Model. KD indicates whether a model is distilled by a ranker. We copy the results from original papers and leave them blank if the original paper does not report the result. The best results are in underlined fonts.

4.1 Datasets and Evaluation Metrics

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

Consisting with previous studies on dense information retrieval (Hofstätter et al., 2021b; Xiong et al., 2021), we use popular passage retrieval datasets, MS-MARCO (Nguyen et al., 2016). The dataset contains 8.8M passages from Web pages gathered from Bing's results to real-world queries. The training set contains about 500k pairs of query and relevant passage, and the dev set consists of 6,980 queries. Based on the queries and passages in the dataset, MS-MARCO passage retrieval and ranking tasks were created. Following previous works (Zeng et al., 2022), we report the performance on MS-MARCO Dev set as well as TREC Deep Learning (DL) 2019 set (Craswell et al., 2020) which includes 43 queries. We report MRR@10 and Recall@50/1K for MS-MARCO, and nDCG@10 and Recall@100 for TREC DL 19.

449

450

451

452

453

454

455

456

457

458

459

460

461

4.2 Baselines

To make a comprehensive comparison, we choose both sparse and dense passage retrievers as baselines. Please refer to Appendix A.1 for the details of baseline methods.

4.3 Implementation Details

Consisting with the setting of RocketQA V2 (Ren et al., 2021c), we choose the learned dual-encoder in the first step of RocketQA (Qu et al., 2021) as the initialization of our dense retriever³. We adopt two advanced cross-encoder rankers as our teacher model: RocketQAV2 (Ren et al., 2021c)

³The retriever can also be replaced with other trained retriever. We observed that using the trained model to initialize the retriever can help achieve slightly better results.

and \mathbb{R}^2 anker (Zhou et al., 2023)⁴. We randomly 462 select m hard negatives provided by Ren et al. 463 (2021c) for each query. For supervised learning, 464 a positive passage and all the selected negatives 465 are used. While for distillation, the candidate passage set for a query consists of m original nega-467 tives, m dark examples in $\mathcal{N}_i^{mix},$ and 5 dark exam-468 ples in \mathcal{N}_i^{mask} with different masking ratios $m_r \in$ 469 $\{0.15, 0.25, 0.35, 0.45, 0.55\}$. We set the number 470 of negatives m to 10 from $\{5, 10, 15, 20, 25, 30\}^5$. 471 We set the maximum lengths for queries and pas-472 sages as 32 and 128. The dropout rate is set 473 to 0.1 on the cross-encoder. In training, we use 474 AdamW (Loshchilov and Hutter, 2017) as the op-475 timizer to train the model. We set the batch size 476 as 128, the peak learning rate as 5e - 5, and the 477 warm-up steps as 100. We set the weight λ for 478 the supervised objective as 0.01 by varying it in 479 $\{0.001, 0.01, 0.05, 0.1, 0.5\}.$ 480

4.4 Overall Performance

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

503

504

505

We report the overall evaluation results on MS-MARCO and TREC Deep Learning 2019 respectively. On both benchmarks, we not only show the performance of our dual-encoder retriever under knowledge distillation from two different crossencoder teachers, but also provide comparisons between different choices of construction of candidate set \mathbb{P}_i . The main results are shown in Table 1. We can draw three main conclusions:

Our created dark examples improve the performance of knowledge distillation over hard negatives and random negatives. With the same crossencoder as the teacher, we analyze the impact of how the candidate set of passages is constructed. It can be observed that using random negatives results in poor performance and the integration of hard negative mining indeed improve the performance. When equipped with our created dark examples which are even harder than existing hard negatives, our model further makes a substantial improvement over that using hard negatives.

Our framework ADAM is compatible with different teachers. To test the generalization ability over different teachers, we conduct experiments using two advanced cross-encoders (\mathbb{R}^2 anker and

Methods	MRR@10
Adam	<u>38.99</u>
w/o. \mathcal{N}^{rein} (Eq.5)	38.82
w/o. \mathcal{N}^{mask} (Eq.6)	38.76
w/o. { \mathcal{N}^{rein} & \mathcal{N}^{mask} }	38.64
w/o. { \mathcal{N}^{rein} & \mathcal{N}^{mask} & ADA }	38.61
w/o. { \mathcal{N}^{rein} & \mathcal{N}^{mask} & ADA & \mathcal{L}_{sup} }	38.36

Table 2: Ablation results on MS-Marco. We report the reranking performance.

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

RocketQAV2) as the teacher. Consistent improvement can be observed when using our proposed dark examples for knowledge distillation with the two different teachers. Moreover, we can compare the effectiveness of the two teachers. When using random negatives, knowledge distillation with the two teachers results in comparable results. But when using hard negatives and dark examples, the model distilled by R^2 anker yields significantly better performance than its counterparts. Therefore, for the remaining ablation studies and analyses, we use R^2 anker as the teacher by default.

With R^2 anker as the teacher, our method (the bottom line) achieves superior performance over most baselines. Our model achieves 41.00 on MRR@10 on the development set of MS-MARCO, outperforming most of the existing methods and is comparable with SimLM (Wang et al., 2023) which is obtained by a time-consuming large-scale pretraining followed with a cumbersome multi-stage supervised fine-tuning.

4.5 Ablation study

We have analyzed the overall performance on two benchmarks and proved the effectiveness of our method. Here, we conduct ablation studies to verify the indispensability of each crucial design. We provide the results of the ablation study in Table 2.

Dark examples. Recall that we propose two types of methods to construct dark examples: (1) strengthening negatives (\mathcal{N}^{rein}) by mixing with the positive to make negatives more relevant to the query, and (2) polluting positives ((\mathcal{N}^{mask}))) to make positives not that relevant. We first test the individual effect of \mathcal{N}^{rein} and \mathcal{N}^{mask} . When removing each of them individually, performance drops can be observed. And when we remove both of them, the model performs worse. This observation indicates that the incorporation of both \mathcal{N}^{rein} and \mathcal{N}^{mask} is beneficial to the overall performance.

Distillation with adaptive dark examples. In addition to dark examples, we also introduce a self-paced distillation algorithm that can better trans-

 $^{^{4}}$ The results of BM25-reranking on MS-MARCO Dev for R²anker (Zhou et al., 2023) and RocketQAV2 (Ren et al., 2021c) are 40.1 and 40.7 respectively.

⁵We found m = 15 to be the optimal parameter. However, considering that our method will expand the number of negatives with the augmented dark examples, we set m=10 in our experiment.



Figure 3: (a) The impact of m; (b) Distributions of model prediction for the R²anker over MS-MARCO.

fer dark knowledge with adaptive dark examples. When this strategy is removed, we create dark examples for all the training instances. It can be seen that distillation adaptively using the subset of instances that the teacher is most confident is better than using the whole training set, which is in accord with our assumption that the instances with higher confidence are a better carrier of knowledge.

Distillation with additional supervised loss. Although the teacher's soft label provides abundant dark knowledge for the student to learn, we also involve the traditional supervised loss. We can observe that although the weight λ for supervised loss is quite small (i.e., 0.01), we find this term indispensable for the overall performance.

4.6 Discussions

The impact of the number of negatives. When constructing the training set, the number of negatives plays a vital role as it also indirectly controls the number of dark examples. To explore the effect of the number of negative samples as well as to find the best choice for m, we conduct experiments on different m^6 . As illustrated in Figure 3(a), when m is small, increasing m brings a positive effect and leads to the best performance when m = 15. But as the curve indicates, incorporating more negatives brings no benefit, which is also in line with existing findings (Karpukhin et al., 2020). The above trend also indicates that too many trivial negatives (m > 15) can not always bring improvement while incorporating our dark examples can still bring improvement to the knowledge transfer. The phenomenon also reveals the importance of distillation data in IR knowledge transfer.

The impact of mask rate and the number of noisy positives. We further investigate how the mask rate and the number of noisy positives influence the performance of ADAM. Due to the limited computation resource, we test several typical set-

Methods	MRR@10
$m_r = \{0.15, \cdots, 0.55\}$ (ADAM)	<u>38.99</u>
$m_r = 5 \times \{0.15\}$	38.73
$m_r = 5 \times \{0.35\}$ $m_r = 5 \times \{0.55\}$	38.82 38.85
$m_r = 5 \times \{0.75\}$ $m_r = \{0.15, \cdots, 0.45\}$	38.77 38.95
$m_r = \{0.15, \cdots, 0.65\}$	38.96

Table 3: The impact of mask ration and the number of noisy positives.

588

589

590

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

tings for comparison, as shown in Table 3. First, for a fixed number of masking positives (a.k.a., 5), the performance increases until the mask ratio reaches a certain value, and then drops when the mask rate keeps increasing. The results are rational since too smaller mask ratio results in too many highly-relevant candidates while a larger mask ratio leads to too many loosely-relevant candidates. Notably, we can observe that masking with a variety of masking ratios is better than masking with a mono masking ratio.

The impact of dark examples on the output distribution of ranker. Finally, we examine the impact of dark examples on the output distribution of the ranker. As illustrated in Figure 3(b), we draw the score distributions of the positive, negative candidates, and negative candidates plus dark examples using a teacher (R²anker) over MS-MARCO. It can be observed that the scores for most original hard negatives are quite low and distributed far from the positives that have high scores. By incorporating these dark examples, we are able to improve the smoothness of the score distribution and prob our teacher model with a wider range of candidates that are more diversely relevant to the query. This enables us to more effectively transfer valuable "dark" knowledge from the teacher model.

5 Conclusion

In this paper, we propose a knowledge distillation framework that can better transfer the dark knowledge in the cross-encoder with adaptive dark examples to help the dual-encoder achieve better performance. We propose two approaches to create dark examples that are much harder for the crossencoder teacher to distinguish than typical hard negatives to transfer more dark knowledge. Further, we propose a self-paced distillation strategy that transfers the knowledge adaptively with highconfidence training instances. Experimental results in two widely-used benchmarks verify the effectiveness of our proposed method.

584

585

586

587

⁶To better analyze the impact of the number of negative samples, we conduct the experiment on the model without adaptive dark examples.

Limitations

629

647

650

651

662

665

667

670

671

672

673

674

675

676

677

678

679

(i) Training computation overheads: although having the same inference complexity as any other 631 dense retrieval models, our approach requires more 632 computation resources during training as it expands 633 the number of negatives with the augmented dark examples. (ii) More analysis on noisy positives: 635 due to the limited computation resource, we only test and compare several typical settings of noisy positives, better strategies for constructing noisy positives (e.g., better masking methods and varying the number of noisy positives) can be explored to further improve the performance.

References

- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. Reading wikipedia to answer opendomain questions. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 -August 4, Volume 1: Long Papers, pages 1870–1879.
- Nachshon Cohen, Hedda Cohen Indelman, Yaron Fairstein, and Guy Kushilevitz. 2024. Indi: Informative and diverse sampling for dense retrieval.
- Nick Craswell, Bhaskar Mitra, Emine Yilmaz, Daniel Campos, and Ellen M. Voorhees. 2020. Overview of the TREC 2019 deep learning track. *CoRR*, abs/2003.07820.
- Zhuyun Dai and Jamie Callan. 2019a. Context-aware sentence/passage term importance estimation for first stage retrieval. *CoRR*, abs/1910.10687.
- Zhuyun Dai and Jamie Callan. 2019b. Deeper text understanding for IR with contextual neural language modeling. In Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2019, Paris, France, July 21-25, 2019, pages 985–988.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019a. 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), pages 4171–4186.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019b. 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),

pages 4171–4186. Association for Computational Linguistics.

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

- Emily Dinan, Stephen Roller, Kurt Shuster, Angela Fan, Michael Auli, and Jason Weston. 2018. Wizard of wikipedia: Knowledge-powered conversational agents. In *International Conference on Learning Representations*.
- Thibault Formal, Carlos Lassance, Benjamin Piwowarski, and Stéphane Clinchant. 2021. SPLADE v2: Sparse lexical and expansion model for information retrieval. *CoRR*, abs/2109.10086.
- Luyu Gao and Jamie Callan. 2021a. 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, pages 981–993. Association for Computational Linguistics.
- Luyu Gao and Jamie Callan. 2021b. Unsupervised corpus aware language model pre-training for dense passage retrieval. *CoRR*, abs/2108.05540.
- Luyu Gao and Jamie Callan. 2022. Unsupervised corpus aware language model pre-training for dense passage retrieval. In *ACL*.
- 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*, pages 3030–3042.
- Daniel Gillick, Alessandro Presta, and Gaurav Singh Tomar. 2018. End-to-end retrieval in continuous space. *arXiv preprint arXiv:1811.08008*.
- Hongyu Guo, Yongyi Mao, and Richong Zhang. 2019. Augmenting data with mixup for sentence classification: An empirical study. *arXiv preprint arXiv:1905.08941*.
- Matthew L. Henderson, Rami Al-Rfou, Brian Strope, Yun-Hsuan Sung, László Lukács, Ruiqi Guo, Sanjiv Kumar, Balint Miklos, and Ray Kurzweil. 2017. Efficient natural language response suggestion for smart reply. *CoRR*, abs/1705.00652.
- Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. 2015. Distilling the knowledge in a neural network. *CoRR*, abs/1503.02531.
- 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.
- Sebastian Hofstätter, Sheng-Chieh Lin, Jheng-Hong Yang, Jimmy Lin, and Allan Hanbury. 2021a. Efficiently teaching an effective dense retriever with balanced topic aware sampling. *CoRR*, abs/2104.06967.

736

- 751 752 754 755 756 757 758 761 763
- 765 768 769 770
- 773 774 775 776 778
- 781

786 787

790

793

- Sebastian Hofstätter, Sheng-Chieh Lin, Jheng-Hong Yang, Jimmy Lin, and Allan Hanbury. 2021b. 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, pages 113-122. ACM.
- 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, pages 2553-2561.
- Kyoungrok Jang, Junmo Kang, Giwon Hong, Sung-Hyon Myaeng, Joohee Park, Taewon Yoon, and Hee-Cheol Seo. 2021. UHD-BERT: bucketed ultra-high dimensional sparse representations for full ranking. CoRR, abs/2104.07198.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with gpus. IEEE Transactions on Big Data, 7(3):535–547.
- Yannis Kalantidis, Mert Bulent Sariyildiz, Noe Pion, Philippe Weinzaepfel, and Diane Larlus. 2020. Hard negative mixing for contrastive learning. Advances in Neural Information Processing Systems, 33:21798-21809.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick S. H. Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020, pages 6769-6781. Association for Computational Linguistics.
- Omar Khattab and Matei Zaharia. 2020a. 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, pages 39-48.
- Omar Khattab and Matei Zaharia. 2020b. 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, pages 39-48. ACM.
- 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, pages 6086-6096. Association for Computational Linguistics.

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.

794

795

796

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

- Zhenghao Lin, Yeyun Gong, Xiao Liu, Hang Zhang, Chen Lin, Anlei Dong, Jian Jiao, Jingwen Lu, Daxin Jiang, Rangan Majumder, et al. 2023. Prod: Progressive distillation for dense retrieval. In Proceedings of the ACM Web Conference 2023, pages 3299–3308.
- 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.
- Ilya Loshchilov and Frank Hutter. 2017. Fixing weight decay regularization in adam. ArXiv. abs/1711.05101.
- Jing Lu, Gustavo Hernández Ábrego, Ji Ma, Jianmo Ni, and Yinfei Yang. 2020. Neural passage retrieval with improved negative contrast. CoRR, abs/2010.12523.
- Shuqi Lu, Chenyan Xiong, Di He, Guolin Ke, Waleed Malik, Zhicheng Dou, Paul Bennett, Tie-Yan Liu, and Arnold Overwijk. 2021. Less is more: Pre-training a strong siamese encoder using a weak decoder. *CoRR*, abs/2102.09206.
- Yuxiang Lu, Yiding Liu, Jiaxiang Liu, Yunsheng Shi, Zhengjie Huang, Shikun Feng Yu Sun, Hao Tian, Hua Wu, Shuaiqiang Wang, Dawei Yin, et al. 2022. Erniesearch: Bridging cross-encoder with dual-encoder via self on-the-fly distillation for dense passage retrieval. arXiv preprint arXiv:2205.09153.
- Yi Luan, Jacob Eisenstein, Kristina Toutanova, and Michael Collins. 2021a. Sparse, dense, and attentional representations for text retrieval. Trans. Assoc. Comput. Linguistics, 9:329–345.
- Yi Luan, Jacob Eisenstein, Kristina Toutanova, and Michael Collins. 2021b. Sparse, dense, and attentional representations for text retrieval. Transactions of the Association for Computational Linguistics, 9:329-345.
- Aditya Menon, Sadeep Jayasumana, Ankit Singh Rawat, Seungyeon Kim, Sashank Reddi, and Sanjiv Kumar. 2022. In defense of dual-encoders for neural ranking. In International Conference on Machine Learning, pages 15376-15400. PMLR.
- Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. MS MARCO: A human generated machine reading comprehension dataset. In Proceedings of the Workshop on Cognitive Computation: Integrating neural and symbolic approaches 2016 co-located with the 30th Annual Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain, December 9, 2016, volume 1773 of CEUR Workshop Proceedings. CEUR-WS.org.

Rodrigo Nogueira, Jimmy Lin, and AI Epistemic. 2019a. From doc2query to doctttttquery. *Online preprint*.

850

853

861

870

871

873

874

875

876

877

878

879

881

887

890

891

895

897

900

901

902 903

- Rodrigo Nogueira, Wei Yang, Jimmy Lin, and Kyunghyun Cho. 2019b. Document expansion by query prediction. *CoRR*, abs/1904.08375.
- Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. 2019. Relational knowledge distillation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 3967–3976.
 - 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 opendomain 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, pages 5835–5847. Association for Computational Linguistics.
- Ruiyang Ren, Shangwen Lv, Yingqi Qu, Jing Liu, Wayne Xin Zhao, QiaoQiao She, Hua Wu, Haifeng Wang, and Ji-Rong Wen. 2021a. PAIR: Leveraging passage-centric similarity relation for improving dense passage retrieval. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP* 2021, pages 2173–2183.
- Ruiyang Ren, Shangwen Lv, Yingqi Qu, Jing Liu, Wayne Xin Zhao, Qiaoqiao She, Hua Wu, Haifeng Wang, and Ji-Rong Wen. 2021b. 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, volume ACL/IJCNLP 2021 of *Findings of ACL*, pages 2173– 2183. Association for Computational Linguistics.
- Ruiyang Ren, Yingqi Qu, Jing Liu, Wayne Xin Zhao, Qiaoqiao She, Hua Wu, Haifeng Wang, and Ji-Rong Wen. 2021c. 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, pages 2825–2835. Association for Computational Linguistics.
- Stephen E. Robertson and Hugo Zaragoza. 2009. The probabilistic relevance framework: BM25 and beyond. *Found. Trends Inf. Retr.*, 3(4):333–389.
- Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. 2014. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*.
- 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.

Tao Shen, Xiubo Geng, Chongyang Tao, Can Xu, Guodong Long, Kai Zhang, and Daxin Jiang. 2023.
Unifier: A unified retriever for large-scale retrieval.
In Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD '23, page 4787–4799, New York, NY, USA.
Association for Computing Machinery. 904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2023. SimLM: Pre-training with representation bottleneck for dense passage retrieval. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2244–2258, Toronto, Canada. Association for Computational Linguistics.
- Shitao Xiao, Zheng Liu, Yingxia Shao, and Zhao Cao. 2022. RetroMAE: Pre-training retrieval-oriented language models via masked auto-encoder. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 538–548, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul Bennett, Junaid Ahmed, and Arnold Overwijk. 2020. Approximate nearest neighbor negative contrastive learning for dense text retrieval. *CoRR*, abs/2007.00808.
- 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.
- Kai Xu, Dae Hoon Park, Chang Yi, and Charles Sutton. 2018. Interpreting deep classifier by visual distillation of dark knowledge. *arXiv preprint arXiv:1803.04042*.
- Peilin Yang, Hui Fang, and Jimmy Lin. 2017a. Anserini: Enabling the use of lucene for information retrieval research. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval, Shinjuku, Tokyo, Japan, August 7-11, 2017*, pages 1253–1256.
- Peilin Yang, Hui Fang, and Jimmy Lin. 2017b. Anserini: Enabling the use of lucene for information retrieval research. In Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval, Shinjuku, Tokyo, Japan, August 7-11, 2017, pages 1253–1256. ACM.
- Yinfei Yang, Ning Jin, Kuo Lin, Mandy Guo, and Daniel Cer. 2020. Neural retrieval for question answering

with cross-attention supervised data augmentation.*CoRR*, abs/2009.13815.

963

964

965 966

967

968

969

970

971

972

973 974

975

976

977 978

979

981

983

984

985

986

987

991

992 993

994

- Hansi Zeng, Hamed Zamani, and Vishwa Vinay. 2022. Curriculum learning for dense retrieval distillation. In Proceedings of the 45th International ACM SI-GIR Conference on Research and Development in Information Retrieval, pages 1979–1983.
- Jingtao Zhan, Jiaxin Mao, Yiqun Liu, Jiafeng Guo, Min Zhang, and Shaoping Ma. 2021a. Optimizing dense retrieval model training with hard negatives. In SI-GIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021, pages 1503–1512. ACM.
 - Jingtao Zhan, Jiaxin Mao, Yiqun Liu, Jiafeng Guo, Min Zhang, and Shaoping Ma. 2021b. Optimizing dense retrieval model training with hard negatives. *CoRR*, abs/2104.08051.
 - 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.
- Yucheng Zhou, Tao Shen, Xiubo Geng, Chongyang Tao, Can Xu, Guodong Long, Binxing Jiao, and Daxin Jiang. 2023. Towards robust ranker for text retrieval. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 5387–5401, Toronto, Canada. Association for Computational Linguistics.
- 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, pages 4014–4022. ACM.

1000

1001

1002

1003

1005

1006

1007

1008

1009

1010

1013

1014

1015

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1030

1031

1032

1033

A Appendix

A.1 Baselines

To make a comprehensive comparison, we choose the following state-of-the-art approaches as baselines. These methods contain both sparse and dense passage retrievers.

The sparse retrieval methods include the traditional retriever BM25 (Yang et al., 2017b) and several representative sparse retrievers, including doc2query (Lu et al., 2020), DeepCT (Dai and Callan, 2019a), docTTTTT-query (Nogueira et al., 2019a), UHD-BERT (Jang et al., 2021), COILfull (Gao et al., 2021), UniCOIL (Lin and Ma, 2021), and SPLADE-max (Formal et al., 2021).

The dense retrieval methods produce continuous neural vectors for each passage and The methods include DPR-E (Qu query. et al., 2021), ANCE (Xiong et al., 2021), TAS-Balanced (Hofstätter et al., 2021b), ME-BERT (Luan et al., 2021a), ColBERT (Khattab and Zaharia, 2020b), ColBERT v2 (Santhanam et al., 2021), NPRINC (Lu et al., 2021), ADORE+STAR (Zhan et al., 2021a), Condenser (Gao and Callan, 2021a), RocketOA (Ou et al., 2021), PAIR (Ren et al., 2021b), CoCondenser (Gao and Callan, 2022), RoketQAV2 (Ren et al., 2021c), AR2 (Zhang et al., 2022), CL-DRD (Zeng et al., 2022), ERNIE-Search (Lu et al., 2022), RetroMAE (Xiao et al., 2022), Unifier (Shen et al., 2023), bi-SimLM (Wang et al., 2023), PROD (Lin et al., 2023) and InDi (Cohen et al., 2024). Some of them are enhanced by knowledge distillation from the ranker. For example, RoketQAV2, AR2, and ERNIE-Search introduce the on-the-fly distillation method. CL-DRD and PROD propose progressive distillation with a data curriculum to gradually improve the student.

A.2 More Discussions

Methods	MRR@10
Adaptive-THC (ADAM)	38.99
Adaptive-SHC	38.68
Adaptive-TLC	38.65
Full	38.71

Table 4: Comparison of different data curriculums. Adaptive-THC and Adaptive-TLC mean selecting highconfidence samples and low-confidence samples given by the ranker respectively based on Equation (7) during training. Full means the model is trained with all samples.

Comparison of different data curriculums. To 1034 demonstrate the effect of our adaptive strategy 1035 (Eq.8), we compare our strategy with several differ-1036 ent strategies, including selecting low-confidence 1037 samples given by the ranker (denoted as Adaptive-TLC), selecting high-confidence samples given by 1039 the student (denoted as Adaptive-SHC), and using 1040 all training samples (denoted as Full). The evalua-1041 tion results are shown in Table 4. First, we can find 1042 that selecting high-confidence samples given by the 1043 teacher lead to better performance than using low-1044 confidence samples given by the teacher. Second, 1045 signals provided by the teacher are better than that 1046 provided by the student, as Adaptive-THC outper-1047 forms Adaptive-SHC. Finally, our ADAMchieves 1048 better results than using the whole instances in training, which indicates the effectiveness of our 1050 proposed method. 1051