

GPL: Generative Pseudo Labeling for Unsupervised Domain Adaptation of Dense Retrieval

Anonymous ACL submission

Abstract

Dense retrieval approaches can overcome the lexical gap and lead to significantly improved search results. However, they require large amounts of training data which is not available for most domains. As shown in previous work (Thakur et al., 2021b), the performance of dense retrievers severely degrades under a domain shift. This limits the usage of dense retrieval approaches to only a few domains with large training datasets.

In this paper, we propose the novel unsupervised domain adaptation method *Generative Pseudo Labeling* (GPL), which combines a query generator with pseudo labeling from a cross-encoder. On six representative domain-specialized datasets, we find the proposed GPL can outperform an out-of-the-box state-of-the-art dense retrieval approach by up to 8.9 points nDCG@10. GPL requires less (unlabeled) data from the target domain and is more robust in its training than previous methods.

We further investigate the role of six recent pre-training methods in the scenario of domain adaptation for retrieval tasks, where only three could yield improved results. The best approach, TSDAE (Wang et al., 2021) can be combined with GPL, yielding another average improvement of 1.0 points nDCG@10 across the six tasks. The code is available.¹

1 Introduction

Information Retrieval (IR) is a central component of many natural language applications. Traditionally, lexical methods (Robertson et al., 1994) have been used to search through text content. However, these methods suffer from the lexical gap (Berger et al., 2000) and are not able to recognize synonyms and distinguish between ambiguous words.

Recently, information retrieval methods based on dense vector spaces have become popular to address these challenges. These dense retrieval

methods map queries and passages² to a shared, dense vector space and retrieve relevant hits by nearest-neighbor search. Significant improvement over traditional approaches has been shown for various tasks (Karpukhin et al., 2020; Xiong et al., 2021). This method is also adapted increasingly by industry to enhance the search functionalities of various applications (Choi et al., 2020; Huang et al., 2020).

However, as shown in Thakur et al. (2021b), dense retrieval methods require large amounts of training data to work well.³ Most importantly, dense retrieval methods are extremely sensitive to domain shifts: Models trained on MS MARCO perform rather poorly for questions for COVID-19 scientific literature (Wang et al., 2020; Voorhees et al., 2021). The MS MARCO dataset was created before COVID-19, hence, it does not include any COVID-19 related topics and models did not learn how to represent this topic well in a vector space.

In this work, we present *Generative Pseudo Labeling* (GPL), an unsupervised domain adaptation technique for dense retrieval models (see Figure 1). For a collection of paragraphs from the desired domain, we use an existing pre-trained T5 encoder-decoder to generate suitable queries. For each generated query, we retrieve the most similar paragraphs using an existing dense retrieval model which will serve as negative passages. Finally, we use an existing cross-encoder to score each (query, passage)-pair and train a dense retrieval model on these generated, pseudo-labeled queries using MarginMSE-Loss (Hofstätter et al., 2020).

We use publicly available models for query generation, negative mining, and the cross-encoder, which have been trained on the MS MARCO

²We use passage to refer to text of any length.

³For reference, the popular MS MARCO dataset (Bajaj et al., 2018) has about 500k training instances; the Natural Questions dataset (Kwiatkowski et al., 2019) has more than 100k training instances.

¹Anonymous link.

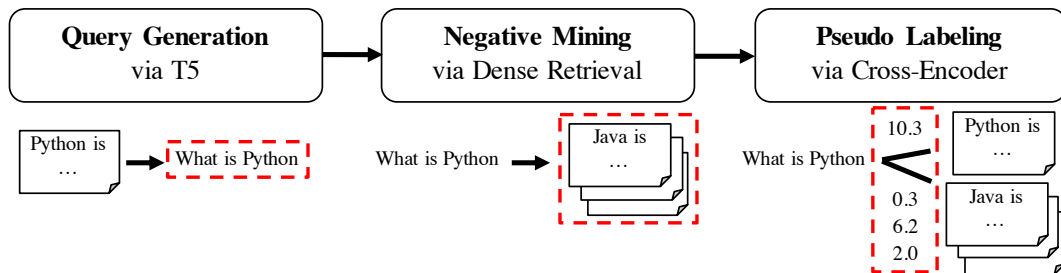


Figure 1: Generative Pseudo Labeling (GPL) for training domain-adapted dense retriever. First, synthetic queries are generated for each passage from the target corpus. Then, the generated queries are used for mining negative passages. Finally, the query-passage pairs are labeled by a cross-encoder and used to train the domain-adapted dense retriever. The output at each step is marked with dashed boxes.

dataset (Nguyen et al., 2016a), a large-scale dataset from Bing search logs combined with relevant passages from diverse web sources. We evaluate GPL on six representative domain-specific datasets from the BeIR benchmark (Thakur et al., 2021b). **GPL improves the performance by up to 8.9 points nDCG@10** compared to state-of-the-art model trained solely on MS MARCO. Compared to the previous state-of-the-art domain-adaption method QGen (Ma et al., 2021; Thakur et al., 2021b), GPL improves the performance by up to 5.2 nDCG@10 points. Training with GPL is easy, fast, and data efficient.

We further analyze the role of six recent pre-training methods in the scenario of domain adaptation for retrieval tasks. The best approach is TSDAE (Wang et al., 2021), that outperforms the second best approach (Masked Language Modeling (Devlin et al., 2019)) on average by 2.5 points nDCG@10. TSDAE can be combined with GPL, yielding another average improvement of 1 point nDCG@10.

2 Related Work

Pre-Training based Domain Adaptation. The most common domain adaption technique for transformer models is domain-adaptive pre-training (Gururangan et al., 2020), which continues pre-training on in-domain data before fine-tuning with labeled data. However, for retrieval it is often difficult to get in-domain labeled data and models are applied in a zero-shot setting on a given corpus. Besides Masked Language Modeling (MLM) (Devlin et al., 2019), different pre-trained strategies specifically for dense retrieval methods have been proposed. Inverse Cloze Task (ICT) (Lee et al., 2019) generates query-passage pair by randomly selecting one sentence from the passage as the query and

the other part as the paired passage. ConDensor (CD) (Gao and Callan, 2021) applies MLM on top of the CLS token embedding from the final layer and the other context embeddings from a previous layer to force the model to learn meaningful CLS representation. SimCSE (Gao et al., 2021a; Liu et al., 2021) passes the same input twice through the network with different dropout masks and minimizes the distance of the resulting embeddings, while Contrastive Tension (CT) (Carlsson et al., 2021) passes the input through two different models. TSDAE (Wang et al., 2021) uses a denoising auto-encoder architecture with bottleneck: Words from the input text are removed and passed through an encoder to generate a fixed-sized embedding. A decoder must reconstruct the original text without noise. As we show in Appendix D, just using these unsupervised techniques is not sufficient and the resulting models perform poorly.

So far, ICT and CD have only been studied on in-domain performance, i.e. a large in-domain labeled dataset is available which is used for subsequent supervised fine-tuning. SimCSE, CT, and TSDAE have been only studied for unsupervised sentence embedding learning. As our results show in Appendix D, they do not work at all for purely unsupervised dense retrieval.

If these pre-training approaches can be used for unsupervised domain adaptation for dense retrieval was so far unclear. In this work, we transfer the setup from Wang et al. (2021) to dense retrieval and first pre-train on the target corpus, followed by supervised training on labeled data from MS MARCO (Nguyen et al., 2016b). Performance is then measured on the target corpus.

Query Generation. Query generation has been used to improve retrieval performances. Doc2query (Nogueira et al., 2019a,b) expands passages with

153 predicted queries, generated by a trained encoder-
 154 decoder model, and uses traditional BM25 lexical
 155 search. This performed well in the zero-shot re-
 156 trieval benchmark BeIR (Thakur et al., 2021b). Ma
 157 et al. (2021) propose QGen, that uses a query gen-
 158 erator trained on general domain data to synthesize
 159 domain-targeted queries for the target corpus, on
 160 which a dense retriever is trained from scratch. Fol-
 161 lowing this idea, Thakur et al. (2021b) views QGen
 162 as a post-training method to adapt powerful MS
 163 MARCO retrievers to the target domains.

164 Despite the success of QGen, previous methods
 165 only consider the cross-entropy loss with in-batch
 166 negatives, which provides coarse-grained relevance
 167 and thus limits the performance. In this work, we
 168 show that extending this approach by using pseudo-
 169 labels from a cross-encoder together with hard neg-
 170 atives can boost the performance by several points
 171 nDCG@10.

172 **Other Methods.** Recently, Xin et al. (2021) pro-
 173 poses MoDIR to use Domain Adversarial Training
 174 (DAT) (Ganin et al., 2016) for unsupervised do-
 175 main adaptation of dense retrievers. MoDIR trains
 176 models by generating domain invariant represen-
 177 tations to attack a domain classifier. However, as
 178 argued in Karouzos et al. (2021), DAT trains mod-
 179 els by minimizing the distance between represen-
 180 tations from different domains and such learning
 181 objective can result in bad embedding space and
 182 unstable performance. For sentiment classification,
 183 Karouzos et al. (2021) proposes UDALM based on
 184 multiple stages of training. UDALM first applies
 185 MLM training on the target domain; and it then ap-
 186 plies multi-task learning on the target domain with
 187 MLM and on the source domain with a supervised
 188 objective. However, as shown in section 5, we find
 189 this method cannot yield improvement for retrieval
 190 tasks.

191 **Pseudo Labeling and Cross-Encoders:** Bi-
 192 Encoders map queries and passage independently
 193 to a shared vector space from which the query-
 194 passage similarity is computed. In contrast, cross-
 195 encoders (Humeau et al., 2020) work on the con-
 196 catenation of the query and passage and predict
 197 a relevance score using cross-attention between
 198 query and passage. This can be used in a re-ranking
 199 setup (Nogueira and Cho, 2019), where the rele-
 200 vancy is predicted for all query-passage-pairs for
 201 a small candidate set. Previous work has shown
 202 that cross-encoders achieve much higher perfor-
 203 mances (Thakur et al., 2021a; Hofstätter et al.,

204 2020; Ren et al., 2021) and are less prone to domain
 205 shifts (Thakur et al., 2021b). But cross-encoders
 206 come with an extremely high computational over-
 207 head, making them less suited for a production set-
 208 ting. Transferring knowledge from cross-encoder
 209 to bi-encoders have been shown previous for sen-
 210 tence embeddings (Thakur et al., 2021a) and for
 211 dense retrieval: Hofstätter et al. (2020) predict
 212 cross-encoder scores for (query, positive)-pairs and
 213 (query, negative)-pairs and learns a bi-encoder to
 214 predict the margin between the two scores. This has
 215 been shown highly effective for in-domain dense
 216 retrieval.

217 3 Method

218 This section describes our proposed *Generative*
 219 *Pseudo Labeling* (GPL) method for the unsuper-
 220 vised domain adaptation of dense retrievers. Fig-
 221 ure 1 illustrates the idea of GPL.

222 For a given target corpus, we generate for each
 223 passage three queries (cf. Table 3) using an T5-
 224 encoder-decoder model (Raffel et al., 2020). For
 225 each of the generated queries, we use an exist-
 226 ing retrieval system to retrieve 50 negative pas-
 227 sages. Dense retrieval with a pre-existing model
 228 was slightly more effective than BM25 lexical re-
 229 trieval (cf. Appendix A). For each (query, posi-
 230 tive, negative)-tuple we compute the margin $\delta =$
 231 $\text{CE}(Q, P^+,) - \text{CE}(Q, P^-)$ with CE the score as
 232 predicted by a cross-encoder, Q the query and
 233 P^+ / P^- the positive / negative passage.

234 We use the synthetic dataset $D_{\text{GPL}} =$
 235 $\{(Q_i, P_i, P_i^-, \delta_i)\}_i$ with the MarginMSE loss (Hof-
 236 stätter et al., 2020) for training a domain-adapted
 237 dense retriever that maps queries and passages into
 238 the shared vector space.

239 Our method requires from the target domain just
 240 an unlabeled collection of passages. Further, we
 241 use pre-existing T5- and cross-encoder models
 242 that have been trained on the MS MARCO passages
 243 dataset.

244 **Query Generation:** To enable supervised train-
 245 ing on the target corpus, synthetic queries can be
 246 generated for the target passages using a query
 247 generator trained on a different, existing dataset
 248 like MS MARCO. Previous work QGen (Ma et al.,
 249 2021) used the simple MultipleNegativesRanking
 250 (MNRL) loss (Henderson et al., 2017; van den
 251 Oord et al., 2018) with in-batch negatives to train
 252 the model:

$$L_{\text{MNRL}}(\theta) = -\frac{1}{M} \sum_{i=0}^{M-1} \log \frac{\exp(\tau \cdot \sigma(f_{\theta}(Q_i), f_{\theta}(P_i)))}{\sum_{j=0}^{M-1} \exp(\tau \cdot \sigma(f_{\theta}(Q_i), f_{\theta}(P_j)))}$$

where P_i is a relevant passage for Q_i ; σ is a certain similarity function for vectors; τ controls the sharpness of the softmax normalization; M is the batch size.

MarginMSE loss: MultipleNegativesRanking loss considers only the coarse relationship between queries and passages, i.e. the matching passage is considered as relevant while all other passages are considered irrelevant. However, the query encoder is not without flaws and might generate queries that are not answerable by the passage. Further, other passages might actually be relevant as well for a given query, which is especially the case if training is done with hard negatives as we do it for GPL.

In contrast, MarginMSE loss (Hofstätter et al., 2020) uses a powerful cross-encoder to soft-label (query, passage) pairs. It then teaches the dense retriever to mimic the score margin between the positive and negative query-passage pairs. Formally,

$$L_{\text{MarginMSE}}(\theta) = -\frac{1}{M} \sum_{i=0}^{M-1} |\hat{\delta}_i - \delta_i|^2 \quad (1)$$

where $\hat{\delta}_i$ is the corresponding score margin of the student dense retriever, i.e. $\hat{\delta}_i = f_{\theta}(Q_i)^T f_{\theta}(P_i) - f_{\theta}(Q_i)^T f_{\theta}(P_i^-)$. Here the dot-product is usually used due to the infinite range of the cross-encoder scores.

This loss is a critical component of GPL, as it solves two major issues from the previous QGen method: A badly generated query for a given passage will get a low score from the cross-encoder, hence, we do not expect the dense retriever to put the query and passage close in the vector space. A false negative will lead to a high score from the cross-encoder, hence, we do not force the dense retriever to assign a large distance between the corresponding embeddings. In section 6.3, we show that GPL is a lot more robust to badly generated queries than the previous QGen method.

4 Experiments

In this section, we describe the experimental setup, the datasets used and the baselines for comparison.

4.1 Experimental Setup

We use the MS MARCO passage ranking dataset (Bajaj et al., 2018) as the data from the source domain. It has 8.8M passages and 532.8K query-passage pairs labeled as relevant in the training set. As Table 1 shows, a state-of-the-art dense retrieval model, achieving an MRR@10 of 33.2 points on the MS MARCO passage ranking dataset, performs poorly on the six selected domain-specific retrieval datasets when compared to simple BM25 lexical search.

We use the DistilBERT (Sanh et al., 2019) for all the experiments. We use a maximum sequence length of 350 with mean pooling and dot-product similarity. For QGen, we use the default setting in Thakur et al. (2021b): 1-epoch training and batch size 75. For GPL, we train the models with 140k training steps and batch size 32. To generate queries for both QGen and GPL, we use the docT5query (Nogueira et al., 2019a) generator trained on MS MARCO and generate ⁴ 3 queries per passage using nucleus sampling with temperature 1.0, $k = 25$ and $p = 0.95$. To retrieve hard negatives for both GPL and the zero-shot setting of MS MARCO training, we use two dense retrievers using cosine-similarity trained on MS MARCO: *msmarco-distilbert-base-v3* and *msmarco-MiniLM-L-6-v3* from Sentence-Transformers⁵. The zero-shot performance of these two dense retrievers are available in Appendix B. We retrieve 50 negatives using each retriever and uniformly sample one negative passage and one positive passage for each training query to form one training example. For pseudo labeling, we use the *ms-marco-MiniLM-L-6-v2*⁶ cross-encoder. For TSDAE pre-training, we train the models with 100K training steps and batch size 8.

4.2 Evaluation

As our methods focus on domain adaptation to specialized domains, we selected six domain-specific text retrieval tasks from the BeIR benchmark (Thakur et al., 2021b): FiQA (financial domain) (Maia et al., 2018), SciFact (scientific papers) (Wadden et al., 2020), BioASQ (biomedical Q&A) (Tsatsaronis et al., 2015), TREC-COVID

⁴We use the script from BeIR at <https://github.com/UKPLab/beir>.

⁵<https://github.com/UKPLab/sentence-transformers>

⁶<https://huggingface.co/cross-encoder/ms-marco-MiniLM-L-6-v2>

(scientific papers on COVID-19) (Roberts et al., 2020), CQADupStack (12 StackExchange subforums) (Hoogeveen et al., 2015) and Robust04 (news articles) (Voorhees, 2005). These selected datasets each contain a corpus with a rather specific language and can thus act as a suitable test bed for domain adaptation.

The detailed information for all the target datasets is available at Appendix C. We make modification on BioASQ and TREC-COVID. For efficient training and evaluation on BioASQ, we randomly remove irrelevant passages to make the final corpus size to 1M. In TREC-COVID, the original corpus has many documents with a missing abstract. The retrieval systems that were used to create the annotation pool for TREC-COVID often ignored such documents, leading to a strong annotation bias for these documents. Hence, we removed all documents with a missing abstract from the corpus. The evaluation results on the original BioASQ and TREC-COVID are available at Appendix C. Evaluation is done using nDCG@10.

4.3 Baselines

Zero-Shot Models: We apply supervised training on MS MARCO or PAQ (Lewis et al., 2021) and evaluate the trained retrievers on the target datasets. (a) **MS MARCO** represents a distilbert-base dense retrieval model trained with MarginMSE on the MS MARCO dataset with batch-size 75 for 70k steps. (b) **PAQ** (Oguz et al., 2021) represents MNRL training on the PAQ dataset. (c) **PAQ + MS MARCO** represents MNRL training on PAQ followed by MarginMSE training on MS MARCO. (d) **TSDAE_{MS MARCO}** represents TSDAE (Wang et al., 2021) pre-training on MS MARCO followed by MarginMSE training on MS MARCO. (e) **BM25** system based on lexical matching from Elasticsearch⁷.

Previous Domain Adaptation Methods: We include two previous unsupervised domain adaptation methods, UDALM (Karouzos et al., 2021) and MoDIR (Xin et al., 2021). For UDALM, we apply MLM training on the target corpus and then apply the multi-task training of MarginMSE training on MS MARCO and MLM training on the target corpus. For MoDIR, it starts from the ANCE checkpoint and apply domain adversarial training on MS MARCO and the target dataset. As of writing, the training code of MoDIR is not public, but domain

adapted models for 5 out of 6 datasets have been released by the authors.

Pre-Training based Domain Adaptation: We follow the setup proposed in Wang et al. (2021) on domain-adapted pre-training: We pre-train the dense retrievers with different methods on the target corpus and then continue to train the models on MS MARCO with MarginMSE loss. The pre-training methods consist of: (a) **CD** (Gao and Callan, 2021) extracts the hidden representations from an intermediate layer and applies MLM on the CLS token representation and these extracted hidden representations⁸. (b) **SimCSE** (Gao et al., 2021b; Liu et al., 2021) simply encode the same text twice with different dropout masks in combination with MNRL loss. (c) **CT** (Carlsson et al., 2021) is similar to SimCSE but it uses two independent encoders to encode a pair of text. (d) **MLM** (Devlin et al., 2019) uses the default setting in original paper, where 15% tokens in a text are sampled to be masked and are needed to be predicted. (e) **ICT** (Lee et al., 2019) uniformly samples one sentence from a passage as the pseudo query to that passage and uses MNRL loss on the synthetic data. We follow the setting in Lee et al. (2019) and masked out the selected sentence 90% of the time. (f) **TSDAE** (Wang et al., 2021) uses a denoising autoencoder to pre-train the dense retrievers with 60% random tokens deleted in the input texts.

Generation-based Domain Adaptation: We use the training script⁹ from Thakur et al. (2021b) to train QGen models with the default setting. For each passage, 3 queries are generated using the same sampling strategy as for GPL. Cosine similarity is used and the models are fine-tuned for 1 epoch with MNRL. The default QGen is trained with in-batch negatives. For a fair comparison, we also test QGen with hard negatives as used in GPL, noted as **QGen (w/ Hard Negatives)**. Further, We we test the combination of TSDAE and QGen (**TSDAE + QGen**).

Re-Ranking with Cross-Encoders: We also include results of the powerful but inefficient re-ranking methods for reference. Three retrievers for the first-phrase retrieval are tested: BM25 from Elasticsearch, the zero-shot MS MARCO retriever and the enhanced GPL retriever by TSDAE pre-training. We use the cross-

⁷<https://www.elastic.co>

⁸CD can only be applied with CLS pooling.

⁹<https://github.com/UKPLab/beir>

Dataset \ Method	FiQA	SciFact	BioASQ	TRECC.	CQADup.	Robust04	Avg.
<i>Zero-Shot Models</i>							
MS MARCO	26.7	57.1	52.9	66.1	29.6	39.0	45.2
PAQ	15.2	53.3	44.0	23.8	24.5	31.9	32.1
PAQ + MS MARCO	26.7	57.6	53.8	63.4	30.6	37.2	44.9
TSDAE _{MS MARCO}	26.7	55.5	51.4	65.6	30.5	36.6	44.4
BM25	23.9	66.1	70.7	60.1	31.5	38.7	48.5
<i>Previous Domain Adaptation Methods</i>							
UDALM	23.3	33.6	33.1	57.1	24.6	26.3	33.0
MoDIR	29.6	50.2	47.9	66.0	29.7	–	–
<i>Pre-Training based Domain Adaptation: Target → MS MARCO</i>							
CT	28.3	55.6	49.9	63.8	30.5	35.9	44.0
CD	27.0	62.7	47.7	65.4	30.6	34.5	44.7
SimCSE	26.7	55.0	53.2	68.3	29.0	37.9	45.0
ICT	27.0	58.3	55.3	69.7	31.3	37.4	46.5
MLM	30.2	60.0	51.3	69.5	30.4	38.8	46.7
TSDAE	29.3	62.8	55.5	76.1	31.8	39.4	49.2
<i>Generation-based Domain Adaptation (Previous State-of-the-Art)</i>							
QGen	28.2	61.7	60.0	72.8	33.6	38.5	49.1
QGen (w/ Hard Negatives)	26.0	59.6	57.7	65.0	33.2	36.5	46.3
TSDAE + QGen (Ours)	30.3	64.7	60.5	73.8	35.1	38.4	50.5
<i>Proposed Method: Generative Pseudo Labeling</i>							
GPL	33.1	65.2	61.6	71.7	34.4	42.1	51.4
TSDAE + GPL	33.3	67.3	62.8	74.0	35.1	42.1	52.4
<i>Re-Ranking with Cross-Encoders (Upper Bound, Inefficient at Inference)</i>							
BM25 + CE	33.1	67.6	72.8	71.2	36.8	46.7	54.7
MS MARCO + CE	33.0	66.9	57.4	65.1	36.9	44.7	50.7
TSDAE + GPL + CE	36.4	68.1	68.0	71.4	38.1	48.3	55.1

Table 1: Evaluation using nDCG@10. The best results of the single-stage dense retrievers are bold. TRECC. and CQADup. are short for TREC-COVID and CQADupStack. Our proposed GPL significantly outperforms other domain adaptation methods. For the first time, we investigate the TSDAE pre-training in domain adaptation for dense retrieval and find it can significantly improve both QGen and GPL.

encoder *ms-marco-MiniLM-L-6-v2* from Sentence-Transformers, which is also for pseudo labeling for GPL.

5 Results

Pre-Training based Domain Adaptation:

The results are shown in Table 1. Compared with the zero-shot MS MARCO model, TSDAE, MLM and ICT can improve the performance if we first pre-train on the target corpus and then perform supervised training on MS MARCO. Among them, TSDAE is the most effective method, outperforming the zero-shot baseline by 4.0 points nDCG@10 on average. CD, CT and SimCSE are not able to adapt to the domains in a pre-training setup and achieve a performance worse than the zero-shot model.

To ensure that TSDAE actually learns domain specific terminology, we include TSDAE_{MS MARCO} in our experiments: Here, we performed TSDAE pre-training on the MS MARCO dataset follow

by supervised learning on MS MARCO. This performs slightly weaker than the zero-shot MS MARCO model.

We also tested the pre-training methods without any supervised training on MS MARCO. We find all of them fail miserably compared as shown in Appendix D.

Previous Domain Adaptation Methods: We test MoDIR on the datasets except Robust04¹⁰. MoDIR performs on-par with our zero-shot MS MARCO model on FiQA, TREC-COVID and CQADupStack, while it performs much weaker on SciFact and BioASQ. An improved training setup with MoDIR could improve the results.

We also test UDALM, which first does MLM pre-training on the target corpus, and then runs multitask learning with MLM objective and supervised training on MS MARCO. The results show that UDALM in this case greatly harms the perfor-

¹⁰The original author did not train the model on Robust04 and the code is also not available.

mance by 12.2 points in average, when compared with the MLM-pre-training approach. We suppose this is because unlike text classification, the dense retrieval models usually do not have an additional task head and the direct MLM training conflicts with the supervised training.

Generation-based Domain Adaptation: The results show that the previous best method, QGen, can successfully adapt the MS MARCO models to the new domains, improving the performance on average by 3.9 points. It performs on par with TSDAE-based domain-adaptive pre-training. Combining TSDAE with QGen can further improve the performance by 1.4 points.

When using QGen with hard negatives instead of random in-batch negatives, the performance decreases by 2.8 points in average. QGen is sensitive to false negatives, i.e. negative passage that are actually relevant for the query. This is a common issue for hard negative mining. GPL solves this issue by using the cross-encoder to determine the distance between the query and a passage. We give more analysis in [Appendix F](#).

Generative Pseudo Labeling (GPL, proposed method): We find GPL is significantly better on almost all the datasets compared to the other tested method, outperforming QGen by up to 4.9 points (on FiQA) and in average by 2.3 points. One exception is TREC-COVID, but as this dataset has just 50 test queries, so this difference can be due to noise.

As a further enhancement, we find that TSDAE-based domain-adaptive pre-training combined with GPL (i.e. TSDAE + GPL) can further improve the performance on all the datasets, achieving the new state-of-the-art result of 52.4 nDCG@10 points in average. It outperforms the out-of-the-box MS MARCO model 7.2 points on average.

Re-ranking with Cross-Encoders: Cross-encoders perform well in a zero-shot setting and outperform dense retrieval approaches significantly ([Thakur et al., 2021b](#)), but they come with a significant computational cost at inference. TSDAE and GPL can narrow but not fully close the performance gap. Due to the much lower computational costs at inference, the TSDAE + GPL model would be preferable in a production setting.

6 Analysis

In this section, we analyze the influence of training steps, corpus size, query generation and choices of

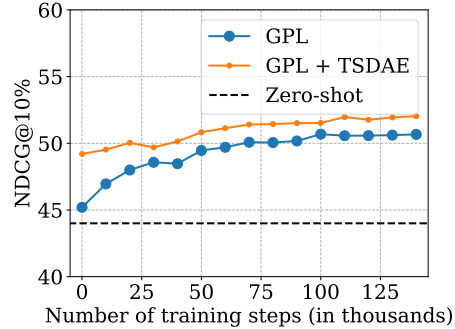


Figure 2: Influence of the number training steps on the averaged performance. The performance of GPL begins to be saturated after 100K steps. TSDAE helps improve the performance during the whole training stage.

Method	Size				
	1K	10K	50K	250K	528K
QGen	35.3	36.9	38.3	37.2	38.5
GPL	37.2	41.3	42.6	42.9	42.1
Zero-shot	39.0				

Table 2: Influence of corpus size on performance on Robust04. The full size is 528K. GPL can achieve the best performance with as little as 50K passages.

starting checkpoints on GPL.

6.1 Influence of Training Steps

We first analyze the influence of the number of training steps on the model performance. We evaluate the models every 10K training steps and end the training after 140K steps. The results for the change of averaged performance on all the datasets are shown in [Figure 2](#). We find the performance of GPL begins to be saturated after around 100K steps. With the TSDAE pre-training, the performance can be improved consistently during the whole training stage. For reference, training a distilbert-base model for 100k steps takes about 9.6 hours on a single V100 GPU.

6.2 Influence of Corpus Size

We next analyze the influence of different corpus sizes. We use Robust04 for this analysis, since it has a relatively large size. We sample 1K, 10K, 50K and 250K passages from the whole corpus independently to form small corpora and train QGen and GPL on the same small corpus. The results are shown in [Table 2](#). We find with more than 10K passages, GPL can already significantly outperform the zero-shot baseline by 2.3 NDCG@10 points; with more than 50K passages, the performance begins to saturate. On the other hand, QGen falls

Method \ QPP	QPP				
	1	2	3	5	10
QGen	57.4	61.6	61.7	62.1	61.3
GPL	60.4	63.0	65.2	64.8	65.6
Zero-shot	57.1				

Table 3: Influence of number of generated Queries Per Passage (QPP) on performance on SciFact. Using a large QPP (e.g. 5 or 10) cannot further improve the performance.

554 behind the zero-shot baseline for each corpus size.

555 6.3 Robustness against Query Generation

556 Next, we study how the query generation influences
 557 the model performance. First, we train QGen and
 558 GPL on SciFact and generate 1 up to 10 queries
 559 per passage. The results are shown in Table 3.
 560 Generating 3 queries per passages appears to be
 561 optimal, generating more queries per passages does
 562 not yield further improvements.

563 The temperature plays an important role in nu-
 564 cleus sampling, higher values make the generated
 565 queries more diverse, but of lower quality. We
 566 train QGen and GPL on FiQA with different tem-
 567 peratures: 0.1, 1, 1.3, 3, 5 and 10. Examples of
 568 generated queries are available in Appendix E. We
 569 generated 3 queries per passage. The results are
 570 shown in Figure 3. We find the performance of
 571 QGen and GPL both peaks at 1.0. With a higher
 572 temperature, the next-token distribution will be flat-
 573 ter and more diverse queries, but of lower quality,
 574 will be generated. With high temperatures, the gen-
 575 erated queries have nearly no relationship to the
 576 passage. QGen will perform poorly in these cases,
 577 worse than the zero-shot model. In contrast, GPL
 578 performs still well even when the generates queries
 579 are of such low quality.

580 6.4 Sensitivity to Starting Checkpoints

581 We also analyze the influence of initialization
 582 on GPL. In the default setting, we start from a
 583 distilbert-model supervised on MS MARCO using
 584 MarginMSE loss. We also evaluate to directly
 585 fine-tune a distilbert-model using QGen, GPL and
 586 TSDAE + GPL. The performance averaged on all
 587 the datasets are shown in Table 4. We find the MS
 588 MARCO training has relatively small effect on the
 589 performance of GPL (with 0.9-point difference in
 590 average), while QGen highly relies on the choice
 591 of the initialization checkpoint (with 3.7-point dif-
 592 ference in average).

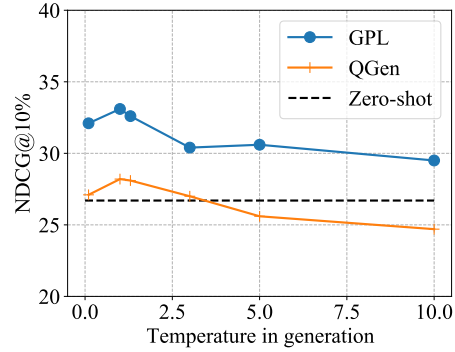


Figure 3: Influence of the temperature in generation on the performance on FiQA. A higher temperature means more diverse queries but of lower quality. GPL can still yield around 3.0-point improvement over the zero-shot baseline with high temperature value of 10.0, where the generated queries have nearly no connection to the passages.

Method \ Init.	Init.	
	Distilbert	MS MARCO
QGen	45.4	49.1
GPL	50.5	51.4
TSDAE + GPL	50.9	52.4
Zero-shot	–	45.2

Table 4: Influence of initialization checkpoint on performance in average. GPL yields similar performance when starting from different checkpoints.

593 7 Conclusion

594 In this work we propose GPL, a novel unsuper-
 595 vised domain adaptation method for dense retrieval
 596 models. It generates queries for a target corpus and
 597 pseudo labels these with a cross-encoders. Pseudo-
 598 labeling overcomes two important short-comings
 599 of previous methods: Not all generated queries are
 600 of high quality and pseudo-labels efficiently detects
 601 those. Further, training with mined hard negatives
 602 is possible as the pseudo labels performs efficient
 603 denoising.

604 We observe GPL performs well on all the
 605 datasets and significantly outperforms other ap-
 606 proaches. As a limitation, GPL requires a relatively
 607 complex training setup and future work can focus
 608 on simplify this training pipeline.

609 In this work, we also evaluated different
 610 pre-training strategies in a domain-adaptive pre-
 611 training setup: We first pre-trained on the target
 612 domain, then performed supervised training on MS
 613 MARCO. ICT and MLM were able to yield a small
 614 improvement, while TSDAE was able to yield a sig-
 615 nificant improvement of 4 points. Other approaches
 616 degraded the performance.

References

- 618 Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng,
619 Jianfeng Gao, Xiaodong Liu, Rangan Majumder, An-
620 drew McNamara, Bhaskar Mitra, Tri Nguyen, Mir
621 Rosenberg, Xia Song, Alina Stoica, Saurabh Tiwary,
622 and Tong Wang. 2018. *Ms marco: A human gener-
623 ated machine reading comprehension dataset*. *arXiv
624 preprint arXiv:1611.09268*.
- 625 Adam L. Berger, Rich Caruana, David Cohn, Dayne
626 Freitag, and Vibhu O. Mittal. 2000. *Bridging the lex-
627 ical chasm: statistical approaches to answer-finding*.
628 In *SIGIR 2000: Proceedings of the 23rd Annual Inter-
629 national ACM SIGIR Conference on Research and
630 Development in Information Retrieval, July 24-28,
631 2000, Athens, Greece*, pages 192–199. ACM.
- 632 Fredrik Carlsson, Amaru Cuba Gyllensten, Evan-
633 gelia Gogoulou, Erik Ylipää Hellqvist, and Magnus
634 Sahlgren. 2021. *Semantic re-tuning with contrastive
635 tension*. In *International Conference on Learning
636 Representations*.
- 637 Jason Ingyu Choi, Surya Kallumadi, Bhaskar Mi-
638 tra, Eugene Agichtein, and Faizan Javed. 2020.
639 *Semantic product search for matching structured
640 product catalogs in e-commerce*. *arXiv preprint
641 arXiv:2008.08180*.
- 642 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
643 Kristina Toutanova. 2019. *BERT: pre-training of
644 deep bidirectional transformers for language under-
645 standing*. In *Proceedings of the 2019 Conference of
646 the North American Chapter of the Association for
647 Computational Linguistics: Human Language Techno-
648 logies, NAACL-HLT 2019, Minneapolis, MN, USA,
649 June 2-7, 2019, Volume 1 (Long and Short Papers)*,
650 pages 4171–4186. Association for Computational
651 Linguistics.
- 652 Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pas-
653 cal Germain, Hugo Larochelle, François Laviolette,
654 Mario Marchand, and Victor S. Lempitsky. 2016.
655 *Domain-adversarial training of neural networks*. *J.
656 Mach. Learn. Res.*, 17:59:1–59:35.
- 657 Luyu Gao and Jamie Callan. 2021. *Condenser: a pre-
658 training architecture for dense retrieval*. In *Proceed-
659 ings of the 2021 Conference on Empirical Methods
660 in Natural Language Processing*, pages 981–993,
661 Online and Punta Cana, Dominican Republic. Asso-
662 ciation for Computational Linguistics.
- 663 Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021a.
664 *SimCSE: Simple contrastive learning of sentence em-
665 beddings*. In *Proceedings of the 2021 Conference
666 on Empirical Methods in Natural Language Process-
667 ing*, pages 6894–6910, Online and Punta Cana, Do-
668 minican Republic. Association for Computational
669 Linguistics.
- 670 Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021b.
671 *SimCSE: Simple contrastive learning of sentence em-
672 beddings*. In *Proceedings of the 2021 Conference
on Empirical Methods in Natural Language Process-
ing*, pages 6894–6910, Online and Punta Cana, Do-
minican Republic. Association for Computational
Linguistics.
- 673 Suchin Gururangan, Ana Marasović, Swabha
674 Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey,
675 and Noah A. Smith. 2020. *Don’t stop pretraining:
676 Adapt language models to domains and tasks*. In
677 *Proceedings of the 58th Annual Meeting of the
678 Association for Computational Linguistics*, pages
679 8342–8360, Online. Association for Computational
680 Linguistics.
- 681 Matthew L. Henderson, Rami Al-Rfou, Brian Strope,
682 Yun-Hsuan Sung, László Lukács, Ruiqi Guo, Sanjiv
683 Kumar, Balint Miklos, and Ray Kurzweil. 2017. *Effi-
684 cient natural language response suggestion for smart
685 reply*. *arXiv preprint arXiv:1705.00652*.
- 686 Sebastian Hofstätter, Sophia Althammer, Michael
687 Schröder, Mete Sertkan, and Allan Hanbury. 2020.
688 *Improving efficient neural ranking models with cross-
689 architecture knowledge distillation*. *arXiv preprint
690 arXiv:2010.02666*.
- 691 Doris Hoogeveen, Karin M Verspoor, and Timothy Bald-
692 win. 2015. *Cquadupstack: A benchmark data set for
693 community question-answering research*. In *Proceed-
694 ings of the 20th Australasian document computing
695 symposium*, pages 1–8.
- 696 Jui-Ting Huang, Ashish Sharma, Shuying Sun, Li Xia,
697 David Zhang, Philip Pronin, Janani Padmanabhan,
698 Giuseppe Ottaviano, and Linjun Yang. 2020.
699 *Embedding-based retrieval in facebook search*. In
700 *KDD ’20: The 26th ACM SIGKDD Conference
701 on Knowledge Discovery and Data Mining, Virtual
702 Event, CA, USA, August 23-27, 2020*, pages 2553–
703 2561. ACM.
- 704 Samuel Humeau, Kurt Shuster, Marie-Anne Lachaux,
705 and Jason Weston. 2020. *Poly-encoders: Architec-
706 tures and pre-training strategies for fast and accurate
707 multi-sentence scoring*. In *International Conference
708 on Learning Representations*.
- 709 Constantinos Karouzos, Georgios Paraskevopoulos, and
710 Alexandros Potamianos. 2021. *UDALM: Unsuper-
711 vised domain adaptation through language modeling*.
712 In *Proceedings of the 2021 Conference of the North
713 American Chapter of the Association for Computa-
714 tional Linguistics: Human Language Technologies*,
715 pages 2579–2590, Online. Association for Computa-
716 tional Linguistics.
- 717 Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick
718 Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and
719 Wen-tau Yih. 2020. *Dense passage retrieval for open-
720 domain question answering*. In *Proceedings of the
721 2020 Conference on Empirical Methods in Natural
722 Language Processing (EMNLP)*, pages 6769–6781,
723 Online. Association for Computational Linguistics.

728	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: A benchmark for question answering research . <i>Transactions of the Association for Computational Linguistics</i> , 7:453–466.	787
729		788
730		789
731		790
732		791
733		
734		792
735		793
736		794
737	Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. 2019. Latent retrieval for weakly supervised open domain question answering . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 6086–6096, Florence, Italy. Association for Computational Linguistics.	795
738		796
739		
740		797
741		798
742		799
743	Patrick S. H. Lewis, Yuxiang Wu, Linqing Liu, Pasquale Minervini, Heinrich Küttler, Aleksandra Piktus, Pontus Stenetorp, and Sebastian Riedel. 2021. PAQ: 65 million probably-asked questions and what you can do with them . <i>arXiv preprint arXiv:2102.07033</i> .	800
744		801
745		802
746		803
747		804
748		805
749	Fangyu Liu, Ivan Vulić, Anna Korhonen, and Nigel Collier. 2021. Fast, effective, and self-supervised: Transforming masked language models into universal lexical and sentence encoders . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 1442–1459, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	806
750		807
751		808
752		809
753		810
754		811
755		
756	Kyle Lo, Lucy Lu Wang, Mark Neumann, Rodney Kinney, and Daniel Weld. 2020. S2ORC: The semantic scholar open research corpus . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4969–4983, Online. Association for Computational Linguistics.	812
757		813
758		814
759		815
760		816
761		817
762		818
763		819
764	Ji Ma, Ivan Korotkov, Yinfei Yang, Keith Hall, and Ryan McDonald. 2021. Zero-shot neural passage retrieval via domain-targeted synthetic question generation . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 1075–1088, Online. Association for Computational Linguistics.	820
765		821
766		822
767		823
768		824
769		825
770		826
771	Macedo Maia, Siegfried Handschuh, André Freitas, Brian Davis, Ross McDermott, Manel Zarrouk, and Alexandra Balahur. 2018. Www’18 open challenge: Financial opinion mining and question answering . In <i>Companion Proceedings of the The Web Conference 2018</i> , WWW ’18, page 1941–1942, Republic and Canton of Geneva, CHE. International World Wide Web Conferences Steering Committee.	827
772		828
773		829
774		830
775		831
776		832
777		833
778		834
779	Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016a. MS MARCO: A human generated machine reading comprehension dataset . In <i>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)</i> , Barcelona, Spain, December 9, 2016, volume 1773 of <i>CEUR Workshop Proceedings</i> . CEUR-WS.org.	835
780		836
781		837
782		838
783		
784		839
785		840
786		841
		842
	Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016b. MS MARCO: A human generated machine reading comprehension dataset . <i>arXiv preprint arXiv:1611.09268</i> .	
	Rodrigo Nogueira and Kyunghyun Cho. 2019. Passage re-ranking with BERT . <i>arXiv preprint arXiv:1901.04085</i> .	
	Rodrigo Nogueira, Jimmy Lin, and AI Epistemic. 2019a. From doc2query to docttttquery . <i>Online preprint</i> .	
	Rodrigo Nogueira, Wei Yang, Jimmy Lin, and Kyunghyun Cho. 2019b. Document expansion by query prediction . <i>arXiv preprint arXiv:1904.08375</i> .	
	Barlas Oguz, Kushal Lakhota, Anchit Gupta, Patrick S. H. Lewis, Vladimir Karpukhin, Aleksandra Piktus, Xilun Chen, Sebastian Riedel, Wen-tau Yih, Sonal Gupta, and Yashar Mehdad. 2021. Domain-matched pre-training tasks for dense retrieval . <i>arXiv preprint arXiv:2107.13602</i> .	
	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer . <i>Journal of Machine Learning Research</i> , 21(140):1–67.	
	Ruiyang Ren, Yingqi Qu, Jing Liu, Wayne Xin Zhao, Qiaoqiao She, Hua Wu, Haifeng Wang, and Ji-Rong Wen. 2021. RocketQAv2: A joint training method for dense passage retrieval and passage re-ranking . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 2825–2835, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	
	Kirk Roberts, Tasmee Alam, Steven Bedrick, Dina Demner-Fushman, Kyle Lo, Ian Soboroff, Ellen Voorhees, Lucy Lu Wang, and William R Hersh. 2020. TREC-COVID: rationale and structure of an information retrieval shared task for COVID-19 . <i>Journal of the American Medical Informatics Association</i> , 27(9):1431–1436.	
	Stephen E. Robertson, Steve Walker, Susan Jones, Micheline Hancock-Beaulieu, and Mike Gatford. 1994. Okapi at TREC-3 . In <i>Proceedings of The Third Text REtrieval Conference, TREC 1994, Gaithersburg, Maryland, USA, November 2-4, 1994</i> , volume 500-225 of <i>NIST Special Publication</i> , pages 109–126. National Institute of Standards and Technology (NIST).	
	Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of BERT: smaller, faster, cheaper and lighter . <i>arXiv preprint arXiv:1910.01108</i> .	
	Nandan Thakur, Nils Reimers, Johannes Daxenberger, and Iryna Gurevych. 2021a. Augmented SBERT: Data augmentation method for improving bi-encoders for pairwise sentence scoring tasks . In	

843	<i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 296–310, Online. Association for Computational Linguistics.	
844		
845		
846		
847		
848	Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. 2021b. Beir: A heterogenous benchmark for zero-shot evaluation of information retrieval models . <i>arXiv preprint arXiv:2104.08663</i> .	
849		
850		
851		
852		
853	George Tsatsaronis, Georgios Balikas, Prodromos Malakasiotis, Ioannis Partalas, Matthias Zschunke, Michael R Alvers, Dirk Weissenborn, Anastasia Krithara, Sergios Petridis, Dimitris Polychronopoulos, et al. 2015. An overview of the bioasq large-scale biomedical semantic indexing and question answering competition . <i>BMC bioinformatics</i> , 16(1):138.	
854		
855		
856		
857		
858		
859		
860	Aäron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding . <i>CoRR</i> , abs/1807.03748.	
861		
862		
863	Ellen Voorhees. 2005. Overview of the trec 2004 robust retrieval track . Special Publication (NIST SP), National Institute of Standards and Technology, Gaithersburg, MD.	
864		
865		
866		
867	Ellen Voorhees, Tasmeer Alam, Steven Bedrick, Dina Demner-Fushman, William R. Hersh, Kyle Lo, Kirk Roberts, Ian Soboroff, and Lucy Lu Wang. 2021. TREC-COVID: Constructing a pandemic information retrieval test collection . <i>SIGIR Forum</i> , 54(1).	
868		
869		
870		
871		
872	David Wadden, Shanchuan Lin, Kyle Lo, Lucy Lu Wang, Madeleine van Zuylen, Arman Cohan, and Hannaneh Hajishirzi. 2020. Fact or fiction: Verifying scientific claims . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 7534–7550, Online. Association for Computational Linguistics.	
873		
874		
875		
876		
877		
878		
879	Kexin Wang, Nils Reimers, and Iryna Gurevych. 2021. TSDAE: Using transformer-based sequential denoising auto-encoder for unsupervised sentence embedding learning . In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 671–688, Punta Cana, Dominican Republic. Association for Computational Linguistics.	
880		
881		
882		
883		
884		
885		
886	Lucy Lu Wang, Kyle Lo, Yoganand Chandrasekhar, Russell Reas, Jiangjiang Yang, Doug Burdick, Darrin Eide, Kathryn Funk, Yannis Katsis, Rodney Michael Kinney, Yunyao Li, Ziyang Liu, William Merrill, Paul Mooney, Dewey A. Murdick, Devvret Rishi, Jerry Sheehan, Zhihong Shen, Brandon Stilson, Alex D. Wade, Kuansan Wang, Nancy Xin Ru Wang, Christopher Wilhelm, Boya Xie, Douglas M. Raymond, Daniel S. Weld, Oren Etzioni, and Sebastian Kohlmeier. 2020. CORD-19: The COVID-19 open research dataset . In <i>Proceedings of the 1st Workshop on NLP for COVID-19 at ACL 2020</i> , Online. Association for Computational Linguistics.	
887		
888		
889		
890		
891		
892		
893		
894		
895		
896		
897		
898		
	Ji Xin, Chenyan Xiong, Ashwin Srinivasan, Ankita Sharma, Damien Jose, and Paul N. Bennett. 2021. Zero-shot dense retrieval with momentum adversarial domain invariant representations . <i>arXiv preprint arXiv:2110.07581</i> .	899
		900
		901
		902
		903
	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 <i>International Conference on Learning Representations</i> .	904
		905
		906
		907
		908
		909

A Performance of Using Different Retrievers for Negative Mining in GPL

The performance of using different retrievers (BM25, dense and BM25 + dense) for mining hard negatives in GPL is shown in Table 5. The results show GPL performs best when using hard negatives mined by dense retrievers.

Method \ Dataset	FiQA	SciFact	BioASQ	TRECC.	CQADup.	Robust04	Avg.
GPL (w/ BM25 + dense)	32.9	64.4	61.1	68.6	33.8	41.3	50.4
GPL (w/ BM25)	31.1	60.9	57.8	67.5	33.5	35.9	47.8
GPL (w/ dense)	33.1	65.2	61.6	71.7	34.4	42.1	51.4
MS MARCO	26.7	57.1	52.9	66.1	29.6	39.0	45.2

Table 5: Performance of using different retrievers for hard-negative mining in GPL. The scores of the baseline MS MARCO and the scores of GPL with dense retrievers are copied from Table 1.

B Performance of the Zero-Shot Retrievers in Hard-Negative Mining

The performance of directly using the zero-shot retrievers for hard-negative mining in GPL is shown in Table 6. Compared with the strong baseline (MS MARCO in Table 6) trained with MarginMSE, *msmarco-distilbert-base-v3* and *msmarco-MiniLM-L-6-v3* are much worse in terms of zero-shot generalization on each dataset. This comparison supports GPL can indeed train powerful domain-adapted dense retrievers with minimum reliance on choices of the retrievers for hard-negative mining.

Method \ Dataset	FiQA	SciFact	BioASQ	TRECC.	CQADup.	Robust04	Avg.
<i>msmarco-distilbert-base-v3</i>	24.0	52.3	45.6	61.1	24.3	30.6	39.7
<i>msmarco-MiniLM-L-6-v3</i>	23.3	48.8	41.9	57.9	24.3	28.5	37.5
MS MARCO	26.7	57.1	52.9	66.1	29.6	39.0	45.2

Table 6: Performance of different zero-shot retrievers. *msmarco-distilbert-base-v3* and *msmarco-MiniLM-L-6-v3* are used in GPL for hard-negative mining. The scores of the baseline MS MARCO are copied from Table 1.

C Target Datasets

FiQA is for the task of opinion question answering over financial data. It contains 648 queries and 5.8K passages from StackExchange posts under the Investment topic in the period between 2009 and 2017. The labels are binary (relevant or irrelevant) and there are 2.6 passages in average labeled as relevant for each query.

SciFact is for the task of verifying scientific claims using evidence from the abstracts of the scientific papers. It contains 300 queries and 5.2K passages built from S2ORC (Lo et al., 2020), a publicly-available corpus of millions of scientific articles. The labels are binary and there are 1.1 passages in average labeled as relevant for each query.

BioASQ is for the task of biomedical question answering. It originally contains 500 queries and 15M articles from PubMed¹¹. The labels are binary and it has 4.7 passages in average labeled as relevant for each query. For efficient training and evaluation, we randomly remove irrelevant passages to make the final corpus size to 1M.

TREC-COVID is an ad-hoc search challenge for scientific articles related to COVID-19 based on the COVID-19 dataset (Wang et al., 2020). It originally contains 50 queries and 171K documents. The original corpus has many documents with only a title and an empty body. We remove such documents and the final corpus size is 129.2K. The labels in TREC-COVID are 3-level (i.e. 0, 1 and 2) and there are 430.8 passages in average labeled as 1 or 2 in the clean-up version.

CQADupStack is a dataset for community question-answering, built from 12 StackExchange subforums: Android, English, Gaming, Gis, Mathematica, Physics, Programmers, Stats, Tex, Unix, Webmasters

¹¹<https://pubmed.ncbi.nlm.nih.gov/>

and WordPress. The task is to retrieve duplicate question posts with both a title and a body text given a post title. It has 13.1K queries and 457.2k passages. The labels are binary and there are 1.4 passages in average labeled as relevant for each query. As in [Thakur et al. \(2021b\)](#), the average score of the 12 sub-tasks is reported.

Robust04 is a dataset for news retrieval focusing on poorly performing topics. It has 249 queries and 528.2K passages. The labels are 3-level and there are in average 69.9 passages labeled as relevant for each query.

The detailed statistics of these target datasets are shown in [Table 7](#).

Dataset	Domain	Title	Relevancy	#Queries	#Passages	PPQ	Query Len.	Passage Len.
FiQA	Financial	✗	Binary	648	57.6K	2.6	10.8	132.2
SciFact	Scientific	✓	Binary	300	5.2K	1.1	12.4	213.6
BioASQ	Bio-Medical	✓	Binary	500	1.0M	4.7	8.1	204.1
BioASQ*	Bio-Medical	✓	Binary	500	14.9M	4.7	8.1	202.6
TREC-COVID	Bio-Medical	✓	3-Level	50	129.2K	430.8	10.6	210.3
TREC-COVID*	Bio-Medical	✓	3-Level	50	171.3K	493.5	10.6	160.8
CQADupStack	Forum	✓	Binary	13,145	457.2K	1.4	8.6	129.1
Robust04	News	✗	3-Level	249	528.2K	69.9	15.3	466.4

Table 7: Statistics of the target datasets used in the experiments. Column **Title** indicates whether there is (✓) a title for each passage or not (✗). Column **PPQ** represents number of Passages Per Query. Query/passage lengths are counted in words. Symbol * marks the original version from the BeIR benchmark ([Thakur et al., 2021b](#))

We also evaluate the models trained in this work on the original version of BioASQ and TREC-COVID datasets from BeIR ([Thakur et al., 2021b](#)). The results are shown in [Table 8](#).

Method	BioASQ*	TRECC.*
GPL	42.5	71.8
TSDAE + GPL	42.6	73.7
QGen	37.8	43.1

Table 8: Performance on the original version of BioASQ and TREC-COVID in BeIR ([Thakur et al., 2021b](#)).

D Performance of Unsupervised Pre-Training

The performance of the unsupervised pre-training methods without access to the MS MARCO data is shown in [Table 9](#). We find ICT is the best method, achieving highest scores on all the datasets. However, all the unsupervised pre-training methods cannot directly yield improvement in performance compared with the zero-shot baseline.

Method	FiQA	SciFact	BioASQ	TRECC.	CQADup.	Robust04	Avg.
CD	6.6	0.6	0.3	9.8	8.1	3.8	4.9
CT	0.2	0.7	0.0	2.5	0.9	0.0	0.7
MLM	5.4	27.8	4.7	16.0	8.5	6.1	11.4
TSDAE	7.8	37.2	6.9	9.4	14.3	10.1	14.3
SimCSE	5.5	25.0	13.1	26.0	14.6	9.8	15.7
ICT	10.2	42.6	39.0	47.5	23.0	16.5	29.8
MS MARCO	26.7	57.1	52.9	66.1	29.6	39.0	45.2

Table 9: Performance of unsupervised pre-training methods with only access to the target corpus as the training data. The scores of the zero-shot baseline MS MARCO are copied from [Table 1](#).

E Examples of Generated Queries under Different Temperatures

The generation temperature controls the sharpness of the next-token distribution. The examples for one passage from FiQA are shown in [Table 10](#) Higher temperature results in longer and less duplicate queries

Item	Text	Pseudo Label
Input Passage	You can never use a health FSA for individual health insurance premiums. Moreover, FSA plan sponsors can limit what they are will to reimburse. While you can't use a health FSA for premiums, you could previously use a 125 cafeteria plan to pay premiums, but it had to be a separate election from the health FSA. However, under N. 2013-54, even using a cafeteria plan to pay for individual premiums is effectively prohibited.	-
Temperature 0.1	can you use a cafeteria plan for premiums	9.1
	can you use a cafeteria plan for premiums	9.1
	can you use a cafeteria plan for premiums	9.1
Temperature 1.0	can i use my fsa to pay for a health plan	9.7
	can i use my health fsa for an individual health plan?	9.9
	can fsa pay premiums	9.2
Temperature 3.0	cafe a number cafe plan is used by	-10.5
	what type of benefits do the health savings accounts cover when applying for medical terms health insurance	-7.2
	why can't an individual file medical premium on their insurance account with an fsa plan instead of healthcare policy.	6.0
Temperature 5.0	which one does not apply after an emergency medical	-11.1
	is medicare cafe used exclusively as plan funds (health savings account	-7.2
	how soon to transfer coffee bean fses to healthcare	-11.0
Temperature 10.0	will employer limit premiums reimbursement on healthcare expenses with caeatla cafetartil and capetarians account on my employer ca. plans and deductible accounts a.f,haaq and asfrhnta,	-2.5
	kfi what is allowed as personal health account or ca	-10.2
	do people put funds back to buy plan plans before claiming an deductible without the provider or insurance cover f/f associator funds of the person you elect? healthfin depto of benefit benefits deduct all oe premiumto payer for individual care	-4.5

Table 10: Examples of generated queries under different temperature value for a passage from FiQA.

Item	Text	GPL	QGen
Query	what is futures contract	-	-
Positive	Futures contracts are a member of a larger class of financial assets called derivatives ...	10.3	1
Negative 1	... Anyway in this one example the s&p 500 futures contract has an "initial margin" of \$19,250, meaning ...	2.0	0
Negative 2	... but the moment you exercise you must have \$5,940 in a margin account to actually use the futures contract ...	0.3	0
Negative 3	... a futures contract is simply a contract that requires party A to buy a given amount of a commodity from party B at a specified price...	8.2	0
Negative 4	... A futures contract commits two parties to a buy/sell of the underlying securities, but ...	6.9	0

Table 11: Examples of the labels assigned to different query-passage pairs in FiQA by GPL and QGen. The key term "futures contract" are marked in bold. QGen uses only 0-1 scores. GPL uses raw logits, which can be any value between positive and negative infinity.

F Case Study: Fine-Grained Labels

GPL uses continuous pseudo labels from a cross-encoder, which can provide more fine-grained information and is more informative than the simple 0-1 labels as in QGen. In this section, we give a more detailed insight into it by a case study.

One example from FiQA is shown in Table 11. The generated query for the positive passage asks for the definition of "futures contract". Negative 1 and 2 only mention futures contract without explaining

the term (with low GPL labels below 2.0), while Negative 3 gives the required definition (which high GPL label 8.2). As an interesting case, Negative 4 gives a partial explanation of the term (with medium GPL label 6.9). GPL assigns suitable fine-grained labels to different negative passages. In contrast, QGen simply labels all of them as 0, i.e. as irrelevant. Such difference explains the advantage of GPL over QGen and why using hard negatives harms the performance of QGen in [Table 1](#).

965
966
967
968
969